



HAL
open science

Formalisme pragmatiste pour le développement de schèmes cognitifs en robotique autonome

Cédric Coussinet

► **To cite this version:**

Cédric Coussinet. Formalisme pragmatiste pour le développement de schèmes cognitifs en robotique autonome. Autre [cs.OH]. Université Paris Sud - Paris XI, 2007. Français. NNT : . tel-00319702v2

HAL Id: tel-00319702

<https://theses.hal.science/tel-00319702v2>

Submitted on 16 Feb 2009 (v2), last revised 23 Feb 2009 (v3)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



N° d'ordre :

UNIVERSITÉ PARIS XI – PARIS-SUD

U. F. R. Scientifique d'Orsay

École doctorale n°427 : Informatique de Paris-Sud

THÈSE

présentée en vue de l'obtention du grade de

Docteur en Sciences de l'UNIVERSITÉ PARIS XI

Discipline : informatique

par

Cédric COUSSINET

préparée au sein du laboratoire LIMSI-CNRS (UPR3251)

FORMALISME PRAGMATISTE POUR LE DÉVELOPPEMENT DE SCHÈMES COGNITIFS EN ROBOTIQUE AUTONOME

Soutenue publiquement le 25 septembre 2007 devant le jury composé de :

M. Jérôme DOKIC (Rapporteur)

M. Philippe GAUSSIER (Rapporteur)

M^{me} Hélène PAUGAM-MOISY (Rapporteur)

M. Jean-Paul SANSONNET (Président)

M^{me} Michelle SEBAG (Examineur)

M. Philippe TARROUX (Directeur)

« Savoir que sera mauvaise l'œuvre que nous ne réaliserons jamais. Plus mauvaise encore, malgré tout, serait celle que nous ne réaliserions jamais. Celle que nous réalisons a au moins le mérite d'exister. Elle ne vaut pas grand-chose, mais elle existe, comme la plante rabougrie du seul et unique pot de ma voisine infirme. Cette plante fait sa joie, et parfois la mienne aussi. Ce que j'écris et qui est mauvais, je le sais bien, peut néanmoins apporter à son tour quelques instants de distraction, qui le détournent de quelque chose de pire, à tel ou tel esprit triste ou malheureux. Cela me suffit ou ne me suffit pas, mais cela a toujours son utilité, et il en est ainsi de la vie tout entière. »

Le livre de l'intranquillité p. 50 de Fernando Pessoa (1999).

« Je reste toujours ébahi quand j'achève quelque chose. Ébahi et navré. Mon instinct de perfection devrait m'interdire d'achever ; il devrait même m'interdire de commencer. Mais voilà : je pêche par distraction, et j'agis. Et ce que j'obtiens est le résultat, en moi, non pas d'un acte de ma volonté, mais bien d'une défaillance de sa part. Je commence parce que je n'ai pas la force de penser ; je termine parce que je n'ai pas le courage de m'interrompre. Ce livre est celui de ma lâcheté. »

Le livre de l'intranquillité p. 174 de Fernando Pessoa (1999).

REMERCIEMENTS

*« Je suis Arnaut qui amasse le vent
Chasse le lièvre avec le bœuf
Et nage contre le courant »*

Arnaut Daniel (vers 1190)

Cette thèse a été initiée en 2001 dans le cadre d'une Bourse Ministérielle de l'Éducation Nationale de la Recherche et de la Technologie et d'un contrat de monitorat sur trois années rattaché à l'Université Paris-Sud. Elle s'est ensuite poursuivie pendant une année avec un demi-poste d'Attaché Temporaire d'Enseignement et de Recherche à l'Université Paris-Sud.

Cependant, le travail présenté dans ce mémoire de thèse n'aurait pu être entrepris ni mené à bien sans le soutien de mon directeur de thèse M. Philippe Tarroux, Professeur à l'École Normale Supérieure d'Ulm et directeur adjoint du Laboratoire d'Informatique pour la Mécanique et les Sciences de l'Ingénieur (UPR3251). Il a été le garant d'une grande liberté intellectuelle toujours guidée par l'exigence et l'esprit critique. Tout en respectant mes choix scientifiques, il a eu la générosité, avec M. Patrick Le Quéré, Directeur de Recherche au CNRS et directeur du LIMSI, de pourvoir à d'excellentes conditions matérielles y compris pendant les deux dernières années où les financements m'ont fait défaut.

Malgré un projet de recherche s'établissant dans la durée, l'encadrement administratif fut assuré grâce à la compréhension toujours reconduite de M^{me} Christine Paulin, Professeur à l'Université Paris-Sud et Directrice de l'École Doctorale N°427 d'informatique, ainsi que celle de M. Yannis Manoussakis, Professeur à l'Université Paris-Sud et délégué aux thèses.

De même, la soutenance n'aurait jamais pu avoir lieu sans la patience et la confiance des rapporteurs, M. Jérôme Dokic, Directeur d'Études à l'École des Hautes Études en Sciences Sociales, M. Philippe Gaussier, Professeur à l'Université de Cergy Pontoise et M^{me} Hélène Paugam-Moisy, Professeur à l'Université de Lyon, et également sans le dévouement du président du jury, M. Jean-Paul Sansonnet, Directeur de Recherche au CNRS.

Par ailleurs, la qualité de ce mémoire, notamment concernant le premier chapitre, doit beaucoup aux relectures de : M^{me} Amandine Afonso, M^{me} Malika Auvray, M. Jean-Baptiste Berthelin, M. Jean-Philippe Lebœuf, M. Gérard Ligozat, M. Patrick Paroubek, M. Philippe Tarroux et pour sa qualité littéraire dans son ensemble, à ma première lectrice, M^{me} Sandrine Enault.

Au delà des livres et des articles, cette thèse ne peut se soustraire à la dette intellectuelle qu'elle a envers M. David Philipona pour le questionnement sur la notion d'espace, M^{me} Malika Auvray au sujet des théories Piagétienne et M. Philippe Tarroux à propos des sciences cognitives en général.

Enfin, la réalisation des expériences qui est le véritable résultat de ce travail doit amplement aux matériels trouvés, confectionnés ou adaptés par les techniciens du laboratoire M. Pierre Durand et M. Daniel Lerin.

Pour tout ce qui précède et ce qui n'a pas été exprimé, je remercie tous ceux qui ont aidé ou du moins autorisé cette ambition aux figures incertaines, et leur demande pardon de ne pas avoir su faire autrement.

RÉSUMÉ

Le progrès actuel de la robotique repose sur le développement de méthodes liées à la spécification des tâches et des techniques associées. Toutefois, cette démarche occulte toutes les situations où l'environnement se révèle imprévisible à cause de la méconnaissance du milieu ou de la présence d'agents autonomes comme les humains par exemple. Dans ces situations, le robot doit être autonome c'est-à-dire être en mesure de spécifier lui-même ses objectifs ainsi que les moyens d'y parvenir. Afin de réaliser un tel robot, un grand nombre d'approches inspirées des sciences cognitives existent. Cependant, aucune d'entre elle n'a permis d'atteindre cet objectif au cours des cinquante dernières années.

Afin de déterminer les raisons de cette impasse générale, une grille d'analyse philosophique et épistémologique a été établie pour mener une étude transversale des différentes approches en sciences cognitives. Le point commun de ces approches se révèle être alors l'utilisation directement ou indirectement d'au moins une hypothèse ontologique qui définit a priori la notion d'objet et la notion de vérité. L'examen des différentes critiques sur les principales conceptions de la vérité révèle alors que l'emploi d'hypothèse(s) ontologique(s) entraîne des paradoxes qui se répercutent dans les théories de la connaissance sur lesquelles se fonde les sciences cognitives.

Une manière de dépasser cette impasse consiste alors à définir la notion de vérité fondée sur une philosophie de la cognition et non sur une philosophie du monde. Dans ce cadre, le pragmatisme de James (1905) se propose de définir la vérité comme l'acceptation d'une croyance en vertu des avantages dont l'évaluation se traduit également par des croyances. Cette circularité amène à considérer la constitution des croyances comme une dynamique sémiotique (Peirce, 1934) toujours dépendante des interactions avec l'environnement. En utilisant la notion d'autopoïèse de Maturana et Varela (1989), il a été alors proposé de définir la cognition comme une autopoïèse sémiotique. A partir de cette conception et des remarques éthologiques de Lorenz (1937), une généalogie en quatre stades de la cognition au cours de l'évolution des espèces a été avancée.

L'analyse du concept d'autopoïèse sémiotique montre que le projet d'artificialisation de la cognition demeure envisageable sur des robots munis de senseurs et d'effecteurs. Par ailleurs, la formalisation de la sémiologie a permis de définir une architecture cognitive et de mettre en évidence la différence entre un système logique décrivant des relations logiques et une autopoïèse sémiotique réalisant des relations effectives. En s'inspirant des systèmes de classeurs (Holland, 1976) considérés comme des structures dissipatives subsymboliques, une architecture cognitive a été spécifiée et implémentée afin de produire une autopoïèse sémiotique minimale.

Considérer la cognition comme une autopoïèse sémiotique implique que la sémiologie débute à partir d'un ensemble de croyances précédent toute expérience cognitive : les *proto-croyances*. Par ailleurs, la nécessité des proto-croyances exprime l'impossibilité d'éliminer le normatif (Putnam, 1983). L'évolution des capacités cognitives d'un individu dépend, en plus de l'environnement, de la complexité de ces schèmes cognitifs initiaux. Pour l'implémentation robotique, les schèmes cognitifs initiaux ont représenté une boucle

sensorimotrice secondaire simple et un détecteur de régularité par prédiction. L'utilisation de ces deux schèmes au cours de l'étude expérimentale a montré d'une part, que l'architecture cognitive proposée possède des auto-organisations permettant de classifier son environnement sans a priori, et d'autre part, que cette auto-organisation native peut être elle-même orientée par des croyances.

Enfin, le concept d'autopoièse sémiotique offre un cadre explicatif pour des domaines divers et vastes, traditionnellement isolés dans leur tentative d'artificialisation telle que l'émergence du couplage sensorimoteur, de la notion d'objet, des comportements d'apprentissage ou du langage. Cette transversalité constitue la clé pour résoudre les difficultés à concevoir un système cognitif artificiel qui intègre nécessairement les propriétés recherchées par les quatre axes de recherche

TABLE DES MATIÈRES

Introduction.....	1
Chapitre I Études philosophiques des approches en robotique cognitive	9
1. Introduction	9
2. Présentation de la grille d'analyse employée.....	10
2.1. Les principales hypothèses philosophiques	11
2.1.1. Les hypothèses ontologiques.....	11
2.1.2. Les hypothèses cognitives.....	16
2.2. Les notions liées à la description.....	23
2.2.1. Les trois types d'investigation	24
2.2.2. Les trois types de description	26
3. Analyse critique de la cognition artificielle	27
3.1. La description de l'individu en soi	29
3.1.1. Le mentalisme.....	30
3.1.1.1. Le symbolisme.....	31
3.1.1.2. Le subsymbolisme.....	37
3.1.2. L'éliminativisme.....	43
3.1.2.1. Le neuromimétisme	44
3.1.2.2. Le connexionnisme.....	53
3.2. La description de la construction d'un individu.....	70
3.2.1. Le fonctionnalisme écologique.....	72
3.2.2. L'évolutionnisme.....	94
3.2.3. L'interactionnisme.....	105
4. Conclusion.....	128
Chapitre II Analyse de l'impasse et de son dépassement par le pragmatisme	137
1. Introduction	137
2. Critique du concept de vérité.....	138
2.1. La vérité comme une correspondance	139
2.2. La vérité comme une cohérence	147
2.3. La vérité comme une vérification.....	150
2.4. La vérité comme un consensus	152
3. Le pragmatisme comme philosophie cognitive.....	154
3.1. Le pragmatisme et l'utilité du vrai.....	155
3.2. Une histoire phylogénétique de la cognition.....	159

3.3. La valeur morale de ce pragmatisme	165
4. Conclusion.....	167
Chapitre III Proposition d'une architecture cognitive générale	171
1. Introduction	171
2. Analyse de l'autopoièse sémiotique	172
2.1. Autopoièse sémiotique et autopoièse physique	173
2.2. L'architecture cognitive	177
2.3. La sémiologie au sein d'une architecture cognitive	182
3. Formalisation de l'architecture cognitive	193
3.1. Systèmes formels et première formalisation de l'architecture cognitive ...	195
3.2. L'architecture cognitive comme une logique non monotone	197
3.3. L'analyse de la contradiction et de son intégration.....	202
3.4. Gestion de l'incertitude des inférences	205
3.5. Les deux niveaux de typage.....	208
3.6. La notion de variable et la logique d'ordre n	211
4. Les systèmes de classeurs	213
4.1. Remarques sur l'architecture générale des systèmes de classeurs.....	214
4.2. L'architecture originelle, le CS1	215
4.3. Les architectures simplifiées de type ZCS et XCS.....	222
4.4. Les architectures avec anticipation comme le ACS.....	224
4.5. Les systèmes de classeurs hiérarchiques ou motivationnels.....	225
5. Spécification de l'architecture cognitive.....	227
5.1. La structure générale	228
5.1.1. Les éléments constituants	228
5.1.2. Les deux mémoires	231
5.1.3. Les schèmes cognitifs fondamentaux.....	234
5.2. Le processus de sélection d'une règle.....	243
5.3. Les bases d'une autopoièse sémiotique.....	245
5.4. L'implémentation	251
6. Discussion sur l'expressivité de l'architecture cognitive générale proposée	252
7. Conclusion.....	256
Chapitre IV Étude comportementale de l'implémentation robotique	259
1. Introduction	259
2. Présentation du robot utilisé.....	261

2.1. Les généralités	261
2.2. La caractérisation des capteurs infrarouges	262
2.3. Les différents types de contrôles moteur.....	266
3. Élaboration des règles initiales	268
3.1. Construction de la boucle sensorimotrice	268
3.2. Caractérisation statistique de l'environnement	271
3.2.1. L'influence de l'environnement sur les données brutes	272
3.2.2. Étude des composantes principales des capteurs.....	277
3.2.3. Construction de la primitive sensorielle.....	284
3.2.4. Construction des environnements imposés	289
3.3. Analyse des prémisses des règles initiales	290
3.3.1. Les règles initiales pour l'environnement réel.....	290
3.3.2. Les règles initiales pour l'environnement imposé.....	292
4. Étude du système dans des environnements imposés	294
4.1. Caractérisation des paramètres fondamentaux	295
4.1.1. Les taux de taxe et d'enchère.....	295
4.1.2. Le taux de remboursement.....	308
4.1.3. Le nombre de règles.....	309
4.1.4. Le seuil d'élimination	311
4.1.5. Les incertitudes initiales des prémisses	314
4.2. La périodicité des états sensoriels imposés.....	317
4.3. L'apport de la création de règles	320
4.4. Introduction de la prédiction.....	325
4.4.1. Vérification de l'ordonnancement	325
4.4.2. Répartition des règles prédictives.....	327
5. Étude du système dans un environnement réel.....	329
5.1. Étude de la stabilité en environnement simple	330
5.1.1. Les états sensoriels observés.....	330
5.1.2. La conservation du comportement	334
5.1.3. L'évolution des populations de règles.....	338
5.1.4. Les dynamiques	340
5.1.5. Le positionnement des règles finales.....	342
5.1.6. Les structures temporelles	347
5.2. Influence de la complexité de l'environnement.....	348

5.2.1. Les états sensoriels observés.....	348
5.2.2. La conservation du comportement	351
5.2.3. L'évolution des populations de règles.....	352
5.2.4. Les dynamiques	355
5.2.5. Le positionnement des règles finales.....	357
5.3. Extension de la base de règles par création spontanée.....	361
5.4. Étude d'un mécanisme prédictif	367
6. Conclusion.....	371
Conclusions et perspectives	375
Bibliographie	383

INTRODUCTION

« L'ÉLÈVE. — *Quatre moins trois... Quatre moins trois... Quatre moins trois ?... ça ne fait tout de même pas dix?*

LE PROFESSEUR. — *Oh, certainement pas, Mademoiselle. Mais il ne s'agit pas de deviner, il faut raisonner. Tâchons de le déduire ensemble. Voulez-vous compter ?*

L'ÉLÈVE. — *Oui, Monsieur. Un..., deux... euh.*

LE PROFESSEUR. — *Vous savez bien compter ? Jusqu'à combien savez-vous compter ?*

L'ÉLÈVE. — *Je puis compter... à l'infini.*

LE PROFESSEUR. — *Cela n'est pas possible, Mademoiselle.*

L'ÉLÈVE. — *Alors, mettons jusqu'à seize.*

LE PROFESSEUR. — *Cela suffit. Il faut savoir se limiter. Comptez donc, s'il vous plaît, je vous en prie.*

L'ÉLÈVE. — *Un, deux..., et puis après deux, il y a trois... quatre...*

LE PROFESSEUR. — *Arrêtez-vous, Mademoiselle. Quel nombre est plus grand ? Trois ou quatre ?*

L'ÉLÈVE. — *Euh... trois ou quatre ? Quel est le plus grand ? Le plus grand de trois ou quatre ? Dans quel sens le plus grand ? »*

La Leçon d'Eugène Ionesco p. 65-66 (1950).

A - Le développement de la robotique industrielle

Les concepts, les méthodes ainsi que les techniques de la robotique constituent des enjeux scientifiques, économiques et sociaux majeurs. Selon la Japan Robotics Association, le marché global de la robotique, estimé à 11 milliards de dollars en 2005, pourrait passer à 24,9 milliards de dollars en 2010, puis décoller en 2025 en atteignant 66,4 milliards de dollars. Les fantasmes de la science fiction ont joué et jouent encore un rôle dans la robotisation de l'ensemble des domaines d'activité. La corrélation entre une forte culture de la robotique de science fiction et l'effort de robotisation de la société s'illustre par la situation particulière du Japon où l'industrie manufacturière en 2004 utilisait 322 robots pour 10000 personnes employées, alors qu'en Europe le nombre de robots pour 10000 personnes employées atteignait seulement 93.

Cependant, les raisons fondamentales liées au développement de la robotique se révèlent plus pratiques, notamment avec les avantages liés à la sécurité des individus, au contrôle de la qualité et à la vitesse de production. La naissance de la robotique industrielle remonte au début des années 60. Un robot industriel se définit succinctement comme une machine possédant plusieurs degrés de liberté et effectuant automatiquement ou semi-automatiquement des opérations de fabrication. Il se conçoit dans un cadre complètement déterminé tant au niveau des objectifs qu'au niveau des contraintes environnementales. La complexité de sa réalisation prend sa source dans la gestion des nombreux axes de liberté et dans leurs mises en œuvre en vue de répondre aux objectifs. En plus des avantages cités précédemment, la robotisation des chaînes de production offre une meilleure flexibilité, ce qui est devenu un critère important en économie. Entre 1990 et 2003, l'indice des prix des robots industriels chute de 100 à 16 en tenant compte à la fois de l'amélioration des performances et de l'augmentation de l'indice de la main d'œuvre dans les entreprises françaises. Cette diminution des coûts ainsi que l'amélioration des techniques entretient en moyenne une croissance régulière de 13,5% par an. En 2004, le parc de la robotique industrielle mondial regroupait environ 850 000 unités.

Il y a une dizaine d'années, le monde de la robotique se composait essentiellement de robots industriels. En 2005, la robotique de services représentait 42% du marché total, en 2025 elle devrait représenter les trois quarts. Les robots de services professionnels se scindent en deux catégories : d'une part les robots assistants qui participent à la réalisation d'une tâche en collaboration étroite avec un opérateur humain, et d'autre part les robots autonomes qui réalisent, indépendamment d'un opérateur humain, une tâche n'ayant pas pour finalité la participation à la fabrication.

La robotique d'assistance s'imisce dans des activités de plus en plus variées. La robotique de services ne se limite plus à la conquête spatiale, bien que, dans ce domaine, elle demeure toujours incontournable. Avec 5400 unités, les robots sous-marins participaient à 21% du nombre total de robots de services pour les professionnels en 2004. L'exploitation des ressources de la mer et la maintenance des centrales nucléaires (Kansai) mobilisent la majorité de ces robots. Les services de nettoyage, pour les avions par exemple (Skywash) et les services de manutention ou de démolition se répandent également, représentant respectivement 14% et 13% du parc des robots de services à destination des professionnels.

Bien que la robotique d'assistance à usage militaire soit devenue une réalité, il reste difficile de parvenir à une estimation chiffrée. Néanmoins, la société iRobot affirme que 150 de leurs robots participent au conflit irakien. Par ailleurs, dans le domaine de la santé, le robot permet aujourd'hui d'améliorer considérablement les pratiques de haute technicité pour des opérations bien définies et toujours sous le contrôle total du praticien (MKM, Zeus, CASPAR). Leur nombre dépassait déjà les 3000 unités en 2004. La robotisation du monde médical commence à s'étendre à l'assistance de personnes souffrant de cécité, en les accompagnant ainsi dans la vie quotidienne tout en veillant à leur sécurité et à leur santé. Ainsi apparaissent des robots prenant la forme d'un bras manipulateur sur fauteuil roulant (MANUS) ou de petits robots mobiles pouvant aider à transporter des objets (Care-O-Bot). La conception de ces derniers doit prendre en compte, pour la première fois, la notion de convivialité qui passe par une réflexion sur l'ergonomie et sur la communication homme-machine. Toutefois, l'interaction avec un robot de relation publique dans un cadre normal reste très limitée ; cela se traduit par un nombre restreint de mises en service de ce type de robot, 20 dans le monde en 2004.

Ce rapide aperçu des activités susceptibles d'utiliser la robotique d'assistance ne se trouve pas réservé aux domaines précédemment évoqués. Les besoins émergent de façon imprévisible, comme le montre l'insolite robotisation des jockeys pour la course de dromadaires en Arabie Saoudite. Tous ces robots d'assistance restent sous la supervision d'un opérateur humain. Les robots mobiles complètement autonomes se limitent à des tâches de nettoyage, des navettes de transport ou de surveillance (CyberGuard) dans des environnements clos et connus. L'autonomie se comprend comme la non intervention de l'homme pour indiquer ou pour aider la réalisation par le robot de la tâche définie lors de sa conception.

Malgré les difficultés dans la communication homme-machine, la diversité des systèmes robotiques ainsi que celle de leur mise en œuvre parfois en milieu hostile démontrent la robustesse des techniques employées et contribuent à encourager leur diffusion. Selon le World Robotics 2004, le nombre de robots de service pour professionnels aura triplé en 2008 pour atteindre 70 000 unités. En ce qui concerne la robotique de services pour particuliers, le marché explose. Introduit en 2001, le nombre de robots aspirateurs autonomes s'élève à 1,2 millions en 2004 dont 500 000 vendus cette

même année. Le marché des tondeuses à gazon autonomes prend la même envolée. Les prévisions pour 2008 estiment à 7 millions le nombre de robots ménagers. Un autre marché connaît un véritable essor, celui de la robotique ludique. Aibo, le précurseur et le plus renommé d'entre eux, s'apparente à un chien possédant plusieurs fonctionnalités dont certaines peuvent être adaptées à l'utilisateur comme la reconnaissance vocale. L'annonce de l'arrêt de la production de ce robot chien n'a pas découragé la concurrence qui propose déjà des modèles similaires. Le Robosapien, jouet en forme de robot humanoïde partiellement programmable, commercialisé par Wow-Wee, a été l'un des succès de 2004, avec plus de 1,5 millions d'exemplaires vendus. Les prix des robots domestiques et des robots jouets baissent continuellement. Ainsi l'offre s'ajuste de plus en plus à la demande. Toutefois, le nombre de robots vendus restent encore faible par rapport à la population susceptible d'acquiescer ce nouveau type de matériel. Les enjeux économiques se trouvent donc principalement liés à la robotique d'assistance et domestique.

B - Deux approches pour la robotique contemporaine

La recherche et le développement de la robotique contemporaine se scindent alors schématiquement en deux directions. La première, qui représente le courant majoritaire, collecte, améliore et intègre les techniques résolvant des problèmes élémentaires mais non-triviaux, comme l'extraction d'informations, la gestion de l'espace, la prise de décision, ou encore l'optimisation des coordinations motrices. Les techniques viennent de différents horizons scientifiques comme la modélisation d'inspiration biologique ou la physique statistique. L'approche cartésienne qui consiste à décomposer un problème en sous-problèmes reste privilégiée. Un effort important est consacré à la spécification des problèmes. De véritables plates-formes de développement pourraient être constituées afin d'accélérer le développement et l'évolution des produits, d'augmenter les performances et les capacités d'adaptation des systèmes face à la diversité des tâches et de leur environnement. L'INRIA a depuis longtemps orienté ses efforts sur le concept de modularité. Microsoft a annoncé, en 2004, la mise au point d'un kit de développement pour robots. Des tentatives de standardisation des problèmes existent, tel le réseau thématique CLAWAR ou la spécification JAUS du département de la défense américaine. Pour l'instant, chaque système robotique conserve ses spécificités matérielles et logicielles adaptées à chaque situation et à des projets précis.

Lors de la conférence RoboBusiness, en 2005, Breazeal, chercheuse au Medialab du MIT résume assez bien les objectifs de la seconde direction de la recherche robotique : « Un robot n'est pas un outil, mais un partenaire. Les robots doivent être dotés d'une personnalité, et être capables de comprendre nos intentions et ce que nous sommes, pour devenir de véritables partenaires, utiles, capables de collaborer avec nous et de s'intégrer socialement à notre environnement ». Dans les faits, cela se comprend comme l'analyse des différents traits d'émotions perçus chez l'homme dans le but de créer une base communicationnelle à réutiliser dans les systèmes robotiques. La compréhension des intentions doit ici s'entendre comme la reconnaissance comportementale de problèmes humains types auxquels le robot aura associé un comportement plus ou moins appris. L'intégration à son environnement se veut être l'ensemble des réactions assurant le bon accomplissement de ses tâches et l'intégrité du système robotique. La décomposition fonctionnelle de la première direction de recherche se retrouve également dans la seconde mais elle reste centrée sur les aspects cognitifs de l'homme inhérents à une bonne interaction avec les machines, notamment la capacité de chaque être humain à projeter des intentions voire une « humanité » sur des choses animées qui en sont dépourvues.

Ces enjeux scientifiques inscrivent le développement de la robotique contemporaine dans une vision de « convergence intégrative » entre les sciences physiques, les sciences de l'ingénieur associées aux technologies de l'information, les sciences cognitives et les sciences sociales. Un défi majeur concerne leur intégration au sein d'un même système robotique.

C - Limitations de la robotique autonome

Néanmoins, ces approches avec leurs réalisations de haute technologie semblent être en contradiction ou du moins en dissonance avec les promesses médiatiques sur leur totale autonomie et leur intelligence. En effet, il existe une différence entre l'imitation d'un comportement qualifié d'intelligent et un comportement produit par un système intelligent, la différence résidant dans la capacité d'adaptation à un changement brusque et imprévu de l'environnement qui nécessite, a posteriori, une réévaluation des enjeux et une redéfinition des acteurs.

L'illustration de ce propos peut se faire en imaginant la situation suivante. Un robot autonome explore un conduit d'aération. À intervalles réguliers, il trace de petits ronds à l'aide d'un bras muni d'un feutre sur la surface plane la plus proche. Ces traces permettront ensuite au système de retrouver son chemin. Le bras peut toucher toutes les parois, mais le sol et les côtés se situent le plus souvent à proximité. Maintenant se présente une situation imprévue, le nouveau conduit à explorer se trouve être également le logement d'un rat qui apprécie le goût du feutre. Le rat efface systématiquement toutes les marques en les léchant. Un système se limitant à la stricte réalisation de la tâche ne perçoit pas le problème, et sans autre solution de secours, il sera condamné à rester bloqué dans le conduit. En robotique classique, la difficulté rencontrée par ce robot constitue une occasion pour mettre à jour les spécifications du problème de l'exploration dans les conduits d'aération, voire de reformuler toute la problématique. Les concepteurs parviendraient sans difficulté à obtenir un robot accomplissant sa fonction malgré la présence de rats. S'il y a adaptation ou apprentissage par essai-erreur, elle se trouve au niveau des concepteurs. Un robot doté de capacités cognitives serait capable d'identifier la présence de quelque chose, de repérer que celle-ci abîme les précieuses traces et de constater que cette même chose ne semble pas pouvoir atteindre le plafond ; puis il serait capable d'en conclure qu'il faut modifier son comportement en écrivant seulement au plafond. Ce type de scénario n'exige pas des capacités cognitives d'ordre supérieur, un corbeau dans une situation analogue y arriverait sans difficulté comme le montre les expériences éthologiques de Lorenz (1950) mais aucun robot pour l'instant n'en est capable dans le principe. Le problème n'est plus celui de l'exploration du conduit d'aération mais celui de la construction d'une problématique. Par ailleurs, les limites de la conception par dessins renvoient aux questions fondamentales posées par l'évolution biologique.

Mais la robotique a-t-elle réellement besoin de ces capacités cognitives si difficiles à définir ? La méthodologie de la robotique classique permet de résoudre des problèmes divers et variés. En 2004, le DARPA a organisé une compétition de robotique autonome dont l'épreuve consistait à parcourir 200 km dans un désert rocailleux. La meilleure équipe, sur plus d'une vingtaine participantes, ne dépassa pas 13 km. L'année suivante, toutes les équipes atteignirent l'objectif à la même épreuve. En un an, il n'y a eu aucune révolution des outils ou des algorithmes en robotique autonome mais une réévaluation précise des caractéristiques environnementales, des besoins et des fins. Pour 2008, l'épreuve sera de traverser une grande ville. La puissance de la robotique moderne s'appuie sur une démarche d'ingénierie itérative, incrémentale. Il est légitime de penser que les 20 prochaines

années se nourriront de ces technologies. Le marché de la robotique s'ouvre et la technologie commence à être à la hauteur des premiers besoins. Les robots grand public seront peu onéreux, simples et surtout, ils assureront des tâches bien identifiées au sein du foyer. Cependant, la variété des problèmes (et leurs imbrications) pour un robot de service aux tâches complexes, reposant sur une interaction, est considérable ; non seulement à cause de la variabilité des personnes et de leurs intérêts, mais également à cause de leurs situations environnementales et sociales. Les robots classiques multifonctions et interactifs obligent à une spécification perpétuelle, au gré de l'activité du robot et des événements. Chaque robot devra être suivi par une équipe de roboticiens. Cet artisanat robotique restreint et restreindra toujours la robotique classique dans sa démocratisation et dans sa diffusion au sein des activités humaines.

L'espoir d'un avenir à long terme pour la robotique passe alors obligatoirement par une troisième voie de recherche : l'étude des propriétés cognitives et leurs possibles intégrations dans des systèmes robotiques en vue de leur conférer une autonomie cognitive. Cette autonomie permettrait au système robotique de spécifier lui-même ses propres problèmes. Dans ce cadre, la communication entre l'homme et une entité cognitive robotique s'apparenterait davantage aux interactions que l'homme peut avoir avec un animal doté de capacités cognitives, tel qu'un chien. Les représentations et les dispositions à agir étant intimement liées au corps de l'entité cognitive, il semble impossible qu'un humain puisse considérer un robot comme son homologue. En effet, de nombreux travaux en psychologie mettent en évidence l'importance des contingences sensori-motrices (Dewey, 1896 ; Piaget, 1987) dans l'élaboration des capacités cognitives. Dans la perspective plus large de concevoir un robot comme un animal, celui-ci ne serait plus programmé mais dressé pour une tâche. La phase de dressage ne doit pas être ici considérée comme un handicap mais comme une souplesse à la fois pour la réutilisation des systèmes robotiques et pour la prise en main par les utilisateurs. Le résultat d'un dressage pourrait être mémorisé et réutilisé pour un robot morphologiquement identique.

Toutefois, est-il concevable de créer une entité artificielle dotée de capacités cognitives ? Si oui, comment ? Si non, pourquoi ? Est-il possible de définir ces capacités cognitives ? Ces interrogations dépassent largement les enjeux évoqués précédemment. En effet, elles questionnent la source de notre propre cognition, notre faculté de penser. Naturellement, les questions relatives à la connaissance sont depuis longtemps traitées par les philosophes, et particulièrement par ceux du XX^e siècle. À ces philosophes se sont rajoutés les psychologues, les pédagogues, les épistémologues, les neurologues, les biologistes et les mathématiciens. Grâce à l'apparition de l'informatique, la robotique offre des moyens expérimentaux pour la simulation des modèles de processus ou principes cognitifs, ce qui représente un véritable tournant pour ces différentes investigations. Pour les positions qui admettent la faisabilité d'un tel projet, la robotique cognitive associée à l'intelligence artificielle représente l'incarnation d'un problème philosophique.

D - Objectifs de la thèse

Les cinquante dernières années de recherche rattachées aux sciences cognitives furent l'objet de nombreuses controverses sur le choix du paradigme. Bien qu'elles aient permis de dégager un grand nombre de concepts favorisant la discussion, à l'heure actuelle, aucun paradigme ne domine entièrement les sciences cognitives. Même s'il existe un continuum entre eux, quatre pôles participant à l'exploration des processus cognitifs apparaissent dans le paysage de la robotique : le fonctionnalisme, l'évolutionnisme, la modélisation neurophysiologique, et l'interactionnisme. Le premier courant revendique une approche

similaire à celle de la robotique classique en décomposant la cognition en fonctions élémentaires mais s'inspirant des situations et des concepts issus de l'éthologie et de la psychologie. L'évolutionnisme cherche les propriétés élémentaires de systèmes primaires dans leurs processus de transformation qui évolueront naturellement vers des systèmes cognitifs. La modélisation neurophysiologique considère que la compréhension de la cognition passe par l'analyse du substrat des propriétés cognitives, le système neuronal. Le dernier courant, l'interactionnisme, considère la cognition comme étroitement liée à deux principaux aspects : le couplage continu par le biais des activités sensori-motrices ainsi que l'ensemble d'activités endogènes auto-organisées en accord avec le contexte écologique. Les algorithmes dégagés au cours de ces années de recherche continuent à être exploités en robotique classique, cependant ils ne parviennent pas à atteindre les prémices d'une autonomie cognitive. Selon Brooks (2001), l'un des pères du renouveau de la recherche en robotique, le manque de résultats suggère une ou plusieurs impasses conceptuelles qui demeurent floues au regard de la communauté scientifique. Ici, ce mémoire souhaite contribuer au déblocage de ce verrou conceptuel.

L'identification de ce verrou conceptuel passe par une analyse fine des différents paradigmes présents en robotique cognitive. Avec cet objectif en vue, le premier chapitre présentera les quatre pôles des sciences cognitives sous un nouvel angle, c'est-à-dire définir une arborescence des approches à partir de la manière dont celles-ci abordent le problème de la description de l'individu autonome. La première bifurcation sépare les approches souhaitant décrire l'individu en soi de celles souhaitant décrire la construction de l'individu par ses capacités d'agir. Pour chacun de ces niveaux de description de l'individu, une synthèse des principales approches sera proposée. Au préalable, afin de les comparer, une grille d'analyse sera établie. Celle-ci, se voulant la plus large possible, reposera sur certains concepts philosophiques et épistémologiques qui seront rappelés. Cette grille d'analyse permettra de mettre en avant malgré leur grande variabilité, le point commun de toutes ces approches. Ce point commun situé au cœur des paradigmes révélera la nécessité d'analyser le concept de vérité afin de savoir si le projet d'une robotique cognitive est sensé, et si oui, à quelles conditions.

L'analyse du concept de vérité ouvrira le deuxième chapitre expliquant les quatre types de vérités implicitement liées aux paradigmes des sciences cognitives fondant la recherche en robotique. Les critiques épistémologiques et logiques de chacune d'entre elles conduiront à considérer la vérité comme une notion indéfinissable, expliquant ainsi l'origine du verrou conceptuel en robotique cognitive. Cependant, la présentation d'une définition pragmatiste de la vérité offrira l'occasion de sortir de cette impasse. L'explicitation des conséquences épistémologiques permettra d'avancer une généalogie de la cognition.

Dans ce cadre, le troisième chapitre définira précisément la cognition ainsi que les concepts nécessaires à sa compréhension. De cette définition, seront extraites les notions d'architecture cognitive et schèmes cognitifs qui offriront des éléments de réponses concernant la faisabilité et la pertinence de l'artificialisation de la cognition, soit de la robotique cognitive. À partir des algorithmes issus des systèmes de classeurs qui seront analysés, seront présentées successivement une spécification d'une architecture cognitive puis son implémentation.

Le dernier chapitre exposera l'évaluation de l'architecture cognitive implémentée avec un schème cognitif élémentaire afin de vérifier que celle-ci possède bien les propriétés nécessaires au développement de la cognition. Dans un premier temps, l'évaluation du

système robotique portera sur le matériel utilisé et la méthodologie employée afin d'apprécier les capacités d'adaptation. Dans un second temps, l'évaluation du système robotique portera d'une part sur les expériences réalisées en environnement simulé et d'autre part sur celles réalisées en conditions réelles. L'analyse des résultats montrera que l'architecture cognitive développée possède les propriétés que renferme ce modèle élémentaire et par extension le potentiel du formalisme dans lequel il s'inscrit.

La conclusion de ce mémoire rappellera l'enchaînement des principales étapes conduisant à la formalisation d'une philosophie pragmatiste de la cognition. Cette dernière se montera suffisamment explicite pour imaginer et implémenter des modèles cognitifs dans un robot mobile immergé dans un environnement réel. Les résultats obtenus confirmeront la pertinence de cette approche tout en ouvrant de nombreuses perspectives de recherche en robotique cognitive. Cette conclusion se terminera par l'exposition détaillée des pistes de recherche les plus prometteuses.

CHAPITRE I ÉTUDES PHILOSOPHIQUES DES APPROCHES EN ROBOTIQUE COGNITIVE

*« Qu'y a-t-il de plus réel, par exemple, dans notre univers, qu'une vie d'homme, et comment espérer la faire mieux revivre que dans un film réaliste ? Mais à quelles conditions un tel film sera-t-il possible ? A des conditions purement imaginaires. Il faudra en effet supposer une caméra idéale fixée, nuit et jour, sur cet homme et enregistrant sans arrêt ses moindres mouvements. Le résultat serait un film dont la projection elle-même durerait une vie d'homme [...]. Même à ces conditions, ce film inimaginable ne serait pas réaliste. Pour cette raison simple que la réalité d'une vie d'homme ne se trouve pas seulement là où il se tient. Elle se trouve dans d'autres vies qui donnent forme à la sienne [...]. Il n'y a donc qu'un seul film réaliste possible, celui-là même qui sans cesse est projeté devant nous par un appareil invisible sur l'écran du monde. »
Discours de Suède p. 44-45 d'Albert Camus (1957).*

1. Introduction

L'autonomie ne se réduit pas à réaliser des ambitions indépendamment d'autrui, elle implique également l'auto-détermination de ses fins. Étymologiquement, la notion d'autonomie renvoie à une personne ou à une collectivité qui détermine elle-même les lois auxquelles elle se soumet. Cette définition succincte sépare d'une part le choix des règles pouvant se traduire en termes d'objectifs et d'autre part leurs applications, c'est-à-dire les moyens permettant de les respecter ou de les atteindre. Ces deux aspects se trouvent au cœur des sciences de la cognition qui ont pour objet l'acquisition et l'utilisation de la connaissance. L'observation d'un individu autonome permet d'entrevoir tout d'abord le second aspect puis d'en déduire les buts issus du premier aspect. Toutefois le mécanisme de leur détermination se trouve hors de portée. Cette dissymétrie d'observation entre la détermination des fins et l'utilisation des moyens pour les atteindre engendre la confusion entre un individu autonome et un individu semi-autonome. L'individu semi-autonome possède des lois qui lui sont propres mais à l'élaboration desquelles il n'a pas participé, et il agit de façon à les respecter. Les robots munis de systèmes de régulation ou d'asservissement se trouvent dans cette catégorie, ceux qui ne possèdent pas ce minimum d'autonomie s'appellent des automates.

L'impossibilité d'une observation directe des mécanismes et des principes à l'origine de l'autonomie pénalise leur l'étude. Afin de pallier cette difficulté, les sciences cognitives mobilisent les trois types de discours utilisés dans le cadre de la description d'un système : la description procédurale, la description mathématisée et la description componentielle qui seront définies plus en détail. Chacun de ces types de discours possède ses avantages et ses inconvénients. Par ailleurs, la portée de ces discours varie selon les interprétations métaphysiques. Le nombre des approches possibles en sciences cognitives provient de la combinaison de ces facteurs. Cette pluralité se reflète dans le paysage de la robotique autonome. Les systèmes experts, l'inférence bayésienne, les réseaux de neurones

constituent autant d'outils mis en relation avec des positions précises en sciences cognitives. Toutefois, la nécessité pratique de réaliser des robots opérationnels favorise l'apparition d'architectures hybrides. Les implications philosophiques dont héritent les algorithmes employés deviennent alors plus floues. Les techniques peuvent être assujetties à des approches parfois contradictoires avec les idées qui peuvent être à leur origine. La simple énumération des techniques utilisées en robotique ne semble donc pas le meilleur examen pour identifier les impasses théoriques. La compréhension de celles-ci, même en souhaitant se restreindre à la robotique autonome, nécessite alors une étude précise de l'imbrication des concepts et des approches en sciences cognitives.

Afin de mener à bien cette étude indispensable pour pouvoir identifier l'impasse conceptuelle de la robotique mobile, la première partie de ce chapitre s'attachera à distinguer les principales hypothèses philosophiques liées à la cognition ainsi que les trois formes de description des systèmes. L'ensemble de ces notions formera la grille d'analyse qui sera utilisée pour analyser et comparer les approches en sciences cognitives dans la seconde partie. Cette dernière se structurera autour du problème de la description de l'individu en distinguant deux mouvements en sciences cognitives. Le premier mouvement abordant la description de *l'individu en soi* ouvrira le débat sur de dualité corps et esprit. De ce débat se dégagera quatre principaux courants possédant chacun différentes approches dont les limites seront chaque fois illustrées avec des réalisations en intelligence artificielle ou en robotique. L'étude détaillée de l'ensemble de ces courants appartenant à de ce premier mouvement montera que, tant que l'objet de la description reste l'individu en tant que tel, la question sur l'origine de ses connaissances demeurera. Le deuxième mouvement quant à lui oriente la description de l'individu dans la perspective de son évolution permanente de part ses interactions avec l'environnement, c'est dire que la description porte sur *la construction de l'individu*. L'étude des trois principaux courants avec leurs approches révélera qu'en définitive ce changement de perspective déplace le verrou conceptuel qui se trouvait sur l'origine de la connaissance vers l'origine de la téléonomie. La conclusion de ce chapitre proposera une synthèse des courants et de leurs approches selon la grille d'analyse utilisée au cours de leur présentation. Cette synthèse cernera alors le facteur commun à l'ensemble des approches qui seraient à l'origine de l'impasse conceptuelle générale et dont l'analyse dans le chapitre suivant établira un cadre conceptuel suffisamment solide pour proposer une cognition artificielle.

2. Présentation de la grille d'analyse employée

La science, au sens large, se manifeste comme un ensemble d'énoncés à propos du monde. Écrits dans un système de signes conventionnels partagé par une communauté, ces énoncés se modifient suivant une dynamique sociale basée sur l'échange et la révision. Un énoncé provient de la transcription de la connaissance conceptuelle d'un sujet sur le monde. Une définition minimale de ce monde peut-être *ce qui s'impose à la volonté*. Cette définition rejette de fait le solipsisme radical qui stipule qu'il n'est rien d'autre que soi. La discussion sur les rapports entre les concepts médiatisés par les énoncés et les choses perçues s'appuie sur l'acceptation ou le refus de deux types d'hypothèses : les hypothèses ontologiques et les hypothèses cognitives. La connaissance des diverses positions vis-à-vis de ces hypothèses se justifie par le fait qu'elles conditionnent l'interprétation des énoncés. Toutefois, les arguments et les critiques qui accompagnent les différentes positions ne seront pas exposés. De même, les exemples n'ont uniquement vocation qu'à illustrer la diversité des interprétations des hypothèses et non à expliquer précisément le système philosophique sous-tendant la position exemplifiée. L'objectif consiste à se munir des outils

conceptuels qui permettront de dégager le point commun entre toutes les positions concourant actuellement à l'étude des systèmes cognitifs.

L'ensemble de ces hypothèses et de ces méthodes constituera la grille d'analyse utilisée pour expliquer les différentes approches en robotique cognitive au cours de la seconde partie de ce chapitre, mais également afin de les comparer au travers d'un tableau récapitulatif de leurs positions vis-à-vis de ces hypothèses et de ces méthodes. La présentation de ces dernières concernera dans un premier temps les hypothèses ontologiques et les hypothèses cognitives portant sur la perception et la conception. Les premières permettent une interprétation suivant les hypothèses admises de la valeur de la connaissance dans l'absolu. Les secondes renseignent sur les moyens et la possibilité d'atteindre cette connaissance. Dans un second temps, les trois différents types de discours employés pour la description d'un système seront définis, ainsi que les trois méthodes scientifiques employées pour développer des connaissances.

2.1. Les principales hypothèses philosophiques

2.1.1. Les hypothèses ontologiques

L'ontologie se définit comme la science de l'être en tant qu'être, c'est-à-dire qu'elle qualifie et informe sur l'être en tant que tel. Les hypothèses ontologiques sont des hypothèses métaphysiques, elles ne s'appuient sur aucune expérience concrète. Toutefois, il existe un grand nombre d'arguments visant à attaquer ou à défendre telle ou telle position. Les deux hypothèses présentées, l'une (A) sur la physique et l'autre (B) sur l'idéal, ne sont pas forcément celles directement formulées dans les doctrines qu'elles souhaitent inclure. Leurs formulations visent une abstraction compatible avec le plus grand nombre de doctrines. (C) Quatre familles philosophiques seront définies selon l'adhésion à ces hypothèses. L'adhésion à ces familles philosophiques des différents courants en sciences cognitives et leurs approches seront régulièrement rappelés afin de révéler les points de convergences et de divergences sur leur fondement philosophique.

A - L'hypothèse ontologique d'une physique

La première hypothèse ontologique porte sur le monde sensible, le monde physique, ce qui s'impose à la volonté. Le monde sensible n'est pas une hypothèse, mais un constat. En revanche, tout jugement allant au-delà de ce constat est considéré comme une hypothèse. L'hypothèse, qui sera considérée ici comme fondamentale, réside dans la décomposition en essences. En d'autres termes, l'hypothèse ontologique d'une physique (HOP) : *le monde extralinguistique qui impressionne nos sens est constitué d'objets dotés d'une essence propre indépendamment des jugements portés sur eux*. L'ensemble de ces objets forme le monde physique. Cette décomposition du monde peut se traduire de manières très différentes. Par exemple, Paul peut comprendre l'énoncé E1 : « trois pommes se trouvent sur la table » comme un énoncé renvoyant à quatre objets : une table et trois pommes. Chacun de ces objets possède une réalité propre, c'est-à-dire qu'ils ont en eux-mêmes et par eux-mêmes toutes les propriétés de cette table et de ces pommes. Cette conception se rapproche par exemple de celle d'Aristote. Cependant, Jean peut considérer que cette table et ces pommes ne sont pas des objets par eux-mêmes. Pour lui, ces mots ne renvoient qu'à une description pratique ne touchant pas l'essence des choses en elle-même. Néanmoins, Jean peut considérer que les particules élémentaires constituant la table et les pommes dénotent les véritables objets essentiels. Cette dernière interprétation se retrouve par exemple chez

Démocrite. Dans les deux cas, il trouve légitime de penser que le monde est constitué d'objets et que cette décomposition est indépendante des représentations du sujet qui les observe. La connaissance physique résulte de la découverte des objets appartenant au monde.

B - L'hypothèse ontologique d'un idéal

La seconde hypothèse ontologique porte sur les concepts permettant de discuter du monde physique. Le maniement des concepts par l'esprit révèle pour certains un monde différent de celui du monde physique. Dans ce cas, les concepts, ou ce qui relève de l'intellect, trouvent leur origine dans un monde idéal. L'existence de ce monde ainsi que sa nature se traduisent en une hypothèse ontologique d'un idéal (HOI) : *Il existe des concepts objectifs, autrement dit, indépendamment du sujet qui les manipule.* Ces concepts permettent, entre autre, de décrire la réalité physique en termes d'objets avec leurs caractéristiques et leurs propriétés. L'ensemble de ces concepts forme un monde idéal. Comme pour l'hypothèse ontologique d'une physique, il peut y avoir diverses interprétations. Dans l'énoncé E1 précédent, les mots « trois pommes » renvoient à trois objets distincts qui sont tous les trois des instances du concept pomme. Ce concept peut être lui-même précisé en utilisant d'autres concepts « fruit à chair pâle et dure ». Pierre peut croire que les pommes en tant qu'objets possèdent l'ensemble des propriétés qui font d'elles des pommes et le récapitulatif de ces propriétés précède les pommes elles-mêmes. Alors Pierre croit que de toute éternité le concept de pomme existe comme également le penserait par exemple Platon. Mais ces propos peuvent être jugés excessifs par Jules. Pour celui-ci, l'ensemble des qualités d'une pomme n'est qu'une moyenne de l'ensemble des qualités des pommes observées. Toutefois, le refus de croire que le concept pomme soit un concept objectif n'empêche pas Jules de conserver l'hypothèse d'un monde idéal. En effet, dans l'énoncé E1, le concept du nombre trois ne fait pas directement référence à un objet physique. Chacun peut revendiquer un concept de pomme différent suivant son vécu mais tout le monde comprend le concept trois de la même façon. Les concepts logicomathématiques peuvent alors être considérés pour certains, par exemple Russell, comme des concepts objectifs et purs, c'est-à-dire sans aucune mise en relation avec les sens. La question des concepts objectifs peut aussi se poser en termes de concepts éthique ou esthétique comme la justice ou la beauté. Dans tous les cas, pour ceux qui acceptent HOI, la connaissance spirituelle est le résultat de la *découverte* des concepts objectifs. Les exemples illustrant chacune des deux hypothèses montrent un aperçu de la diversité et de l'incompatibilité de positions philosophiques reposant pourtant sur des hypothèses ontologiques communes.

C - Les quatre familles philosophiques

Une catégorisation plus fine des doctrines est autorisée par l'indépendance entre HOP et HOI. De la combinaison des hypothèses ontologiques se dégagent *quatre familles philosophiques*. Chacune de ces familles interprète la connaissance de manière différente, à la fois sur ce qu'elle vise et sur sa valeur. L'étude de ces familles repose sur la distinction entre deux types d'énoncés : les énoncés analytiques et les énoncés synthétiques. La signification des énoncés analytiques dépend uniquement de l'implication des concepts utilisés : l'énoncé E2 « les trois pommes sont et ne sont pas sur la table » est logiquement faux, de même l'énoncé E3 « les trois pommes sont ou ne sont pas sur la table » est nécessairement vraie. La signification de E2 et E3 provient d'un jugement formel. L'énoncé E1 est un énoncé synthétique, sa signification dépend de son rapport (plus ou moins direct) avec le monde physique, c'est un jugement factuel. La distinction entre ces deux types d'énoncés ressemble à celle de Kant (1781) entre les jugements synthétiques reposant sur l'expérience

physique et les jugements analytiques reposant sur le contenu des concepts mais sans y associer une valeur puisque celle-ci se fonderait sur l'adhésion à telles ou telles hypothèses ontologiques et hypothèses cognitives. Sans ordre d'importance, les quatre familles seront présentées successivement : (i) la première famille, moniste, la famille du réalisme conceptuel (FC), qui regroupe toutes les philosophies s'articulant uniquement sur l'HOI, (ii) la deuxième, dualiste, la famille du réalisme conceptuel et physique (FCP), qui revendique HOI et HOP, (iii) la troisième, moniste, la famille du réalisme physique (FP) qui adhère uniquement à HOP, et la dernière, la famille de l'agnosticisme (FA) qui refuse les hypothèses ontologiques.

	HOP : les objets existent indépendamment du sujet	HOP : les objets n'existent pas indépendamment du sujet
HOI : les concepts existent indépendamment du sujet	FCP	FC
HOI : les concepts n'existent pas indépendamment du sujet	FP	FA

Tableau I-1 : Les quatre familles philosophiques selon leur position sur les hypothèses ontologiques.

i - La famille de l'idéalisme moniste (FC)

Pour la famille philosophique FC, en vertu de HOI, l'expérience sensible ne permet pas d'acquiescer toutes les connaissances, en particulier celles liées à la logique et aux mathématiques. Le monde idéal n'est composé que de concepts objectifs purs. Toutefois, ces derniers peuvent servir à construire des concepts servant à la description du monde sensible. Ainsi, la raison est nécessaire pour connaître le monde idéal et décrire le monde physique. Cependant, en fonction des hypothèses cognitives adoptées, la raison peut ne pas être suffisante. Dans un tel cas, la nécessaire prise en compte de l'expérience sensible ne signifie pas l'acceptation de HOP. Rien n'existe en dehors de la raison, c'est-à-dire que tout respecte l'harmonie définie par le monde idéal. Par conséquent, quelle que soit la nature du monde physique, le monde idéal doit nécessairement répondre à une description rationnelle puisqu'il est rationnel : « ce qui est rationnel est réel, ce qui est réel est rationnel » (Hegel, 1821). S'appuyer sur l'expérience pour découvrir des concepts rationnels peut ne pas être une hérésie. Toutefois, l'assujettissement du monde physique au monde idéal donne la primauté à la raison dans l'établissement de la vérité. La science vise ici une description globale des lois régissant le monde physique, chose qui sera possible seulement si tous les concepts nécessaires ont été découverts. Les concepts de pomme ou de quark sont des outils conceptuels pour la construction de l'édifice mais ne renvoient pas à des objets essentiels. La science invente des objets physiques, elle ne les découvre pas. La recherche de la cohérence globale des énoncés devient alors plus importante que la recherche d'une description fidèle du monde physique se fondant sur la qualité des prédictions. Autrement dit, la vérité issue d'un jugement formel et la vérité issue d'un jugement factuel ne possèdent pas la même valeur. La vérité formelle est absolue, de toute éternité, et ne dépend que de l'énoncé lui-même, alors que ce n'est pas le cas pour la vérité factuelle d'un énoncé synthétique. En effet, ici « pomme » et « table » ne sont que des concepts conventionnels pratiques pour désigner un ensemble de propriétés définies de manière consensuelle. Au-delà d'un problème de langage, un martien qui n'a jamais vu de pomme ou entendu parler de ce fruit ne pourra pas saisir l'énoncé et donc ne pourra pas statuer sur sa valeur de vérité. La vérité factuelle se trouve relative au cadre dans lequel est formulé l'énoncé.

ii - La famille du réalisme et de l'idéalisme (FCP)

La famille philosophique FCP accepte l'existence de concepts objectifs et d'objets objectifs. Les concepts objectifs sont soit purs, soit en relation avec le monde sensible. Les premiers justifient leur essence par eux-mêmes, comme le concept d'identité ($A=A$). Les seconds justifient leur essence par l'existence d'instances de ce qu'ils qualifient. Une propriété sur le plan idéal *est* parce qu'il existe un objet physique qui l'incarne et un objet physique *est* parce qu'il actualise des propriétés existant sur le plan idéal. L'objet objectif et les concepts objectifs qui le caractérisent sont indissociables. Le concept de couleur existe parce que, par exemple, les pommes sont colorées, et les pommes ne peuvent pas ontologiquement se soustraire à la couleur. Au sein de la famille FCP, le rôle des concepts objectifs caractérisant l'objet visé reste libre. Une manière d'illustrer la diversité des interprétations des conséquences des deux hypothèses ontologiques aurait été de comparer le système philosophique de Platon et celui de Descartes (1641) par exemple. Mais plus simplement, Pierre peut considérer que la pomme est la réunion des concepts rouge, fruit etc. en accord avec la substance qui porte son existence. En opposition, Jules, lui, considère que l'impression de rouge provient de l'interaction entre la lumière et la matière de la pomme. Il en est de même pour toutes les autres qualités dont Pierre s'est servi pour saisir le concept de pomme. La pomme de Pierre peut être le seul objet au monde. Or, Jules pense que les grains de matière ou d'énergie interagissent selon des lois physiques. L'électron est l'électron de par les lois régissant ses interactions avec les autres objets. L'existence de l'électron n'est pas concevable sans cette interaction. Mais l'électron de Jules peut être considéré à son tour comme la réunion d'un ensemble de lois en accord avec la substance qui les sous-tend. En définitive, Jules appartient à la même famille que Pierre.

Dans tous les cas, les philosophies appartenant à la FCP visent la connaissance à la fois des concepts purs et des concepts objectifs en même temps que les objets du monde sensible auquel ils sont reliés. Quel que soit le choix sur les hypothèses cognitives, en principe, la signification d'un énoncé synthétique dispose de la même puissance qu'une signification formelle. L'objet étant défini de manière objective par des concepts eux-mêmes objectifs, la correspondance existe et la valeur de vérité est absolue. Dans ce cas, un martien d'un niveau de technologie égal pourra statuer sur tout énoncé synthétique puisqu'il aura obligatoirement découvert des concepts objectifs et les objets objectifs.

iii - La famille du réalisme moniste (FP)

Pour la famille philosophique FP qui admet uniquement l'HOP, les concepts ne sont pas objectifs par nature, les concepts purs non plus. Ces derniers résultent d'une certaine inspiration du monde physique et surtout d'une convention entre les hommes. La logique n'est qu'un ensemble de définitions conventionnelles de règles dans le maniement de symboles. La signification des énoncés analytiques perd ici de son importance puisqu'elle ne traduit plus l'absolue harmonie d'un idéal mais le respect de règles inventées. La vérité formelle devient relative au système logique choisi. Les concepts liés au monde physique deviennent des métaphores des objets qui le composent. Ces concepts peuvent s'appuyer sur des concepts purement conventionnels. Mais la convention des concepts physiques est guidée par ce qu'ils visent : les objets objectifs du monde, à la grande différence avec les concepts purement conventionnels qui ne renvoient à rien. Le concept d'électron n'est pas une convention dans l'idée qu'il dénote réellement un objet authentique du monde cependant, les concepts décrivant son existence sont purement conventionnels. FC prône un conventionnalisme physique alors que FP prône un conventionnalisme formel. La signification des énoncés synthétiques perd aussi de sa force, elle ne fait que valider la

bonne métaphore d'une théorie par ses capacités prédictives mais dans une moindre mesure que FC. En effet, un concept d'objet qui ne vise pas un réel objet devient obligatoirement faux. Par ailleurs, quelle que soit la théorie perceptive, les doctrines de FP espèrent converger vers un ensemble de métaphores saisissant selon nos conventions tous les objets du monde physique. La vérité absolue n'est qu'une limite ontologiquement inaccessible vers une métaphore globale parfaite suivant les conventions employées pour tous les objets composants le monde physique. La FP regroupe notamment les naturalistes comme Dennett (1991) par exemple.

iv - La famille de l'agnosticisme (FA)

Enfin, les doctrines de la famille philosophique FA, qui refusent les deux hypothèses ontologiques, adhèrent à la fois au conventionnalisme formel et au conventionnalisme physique. Toutes les constructions de concepts purs et de concepts relatifs au monde sont a priori permises, seul le rapport au monde informe de ce qui est impossible. Ces doctrines ne se revendiquent pas sceptiques (si par ce terme est entendu l'impossibilité de statuer sur la vérité des choses) puisque cette position dépend des hypothèses cognitives. En revanche, elles sont nécessairement agnostiques (si par ce terme est entendu l'impossibilité de définir la vérité) puisque la vérité ne renvoie à aucune réalité en dehors de l'existence de quelque chose. La vérité ne peut prendre aucune des trois formes ontologiques suivantes qui seront détaillées dans le chapitre : l'harmonie conceptuelle, la mise en correspondance des choses et des concepts, la prédictibilité des événements. La vérité n'ambitionne plus d'être absolue, elle résulte d'un consensus sur les principes de sa détermination. La critique de la vérité qu'implique cette famille sera également approfondie dans le chapitre suivant. La diversité des doctrines à l'intérieur de ce courant provient essentiellement des différentes prises de positions vis-à-vis des hypothèses cognitives. Par exemple, cette famille philosophique recouvre la philosophie de Hume (1740) ou celle de James (1907).

*

Les quatre familles philosophiques ainsi définies peuvent s'analyser comme le croisement de lignes d'oppositions (Figure I-1) : celle opposant FC et FP et celle opposant FCP et FA. Mais en fait, la synthèse des deux qui transparaîtra dans la conclusion de ce chapitre verra la première, la ligne d'opposition (FC - FP) se fondre dans la seconde (FCP - FA) pour ne laisser apparaître qu'une seule ligne d'opposition : FC, FCP, FP – FA. En d'autres termes, la principale ligne d'opposition se trouvant au sein des quatre familles philosophiques réside dans la séparation des doctrines possédant une métaphysique réaliste, c'est-à-dire revendiquant HOP ou HOI, et les doctrines avec une métaphysique antiréaliste en refusant les hypothèses ontologiques. En définitive, la Figure I-1 représente la « rose des vents » de tout le paysage de la philosophie.

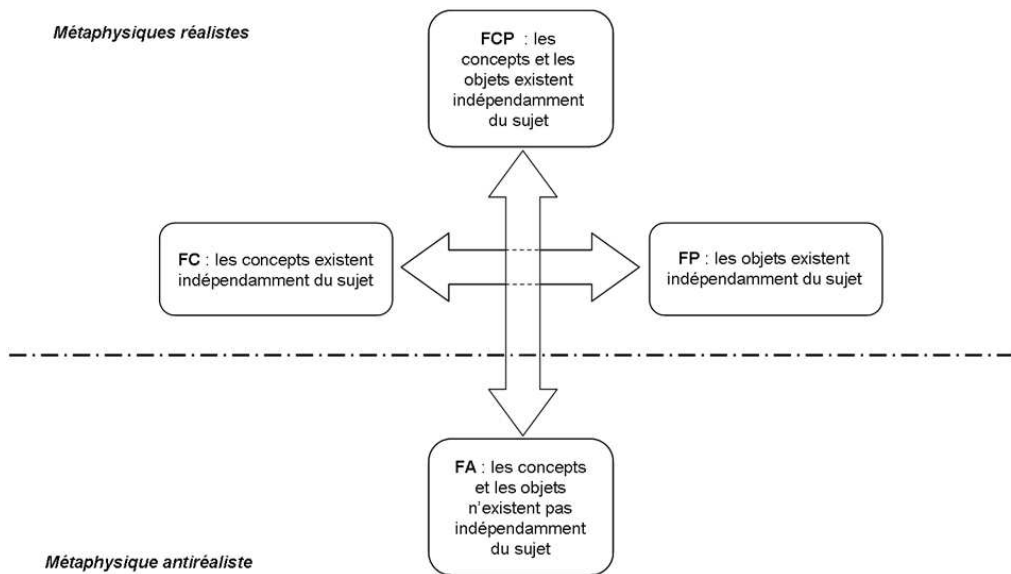


Figure I-1 : Schéma des oppositions entre les quatre familles philosophiques.

2.1.2. Les hypothèses cognitives

Les hypothèses ontologiques définissent en principe la connaissance ou ce qu'elle devrait viser. Les hypothèses cognitives conditionnent les moyens de l'acquérir ainsi que les limites d'une telle entreprise. En fonction des hypothèses ontologiques admises, le terme d'acquisition de la connaissance est synonyme d'accession à la connaissance ou de construction de la connaissance. Toutefois, l'indépendance des hypothèses cognitives et des hypothèses ontologiques se fonde sur la distinction entre accession/construction de la connaissance et la valeur ou la signification qui lui est attribuée. Cette distinction sera détaillée plus attentivement dans le chapitre suivant.

Les hypothèses cognitives portent d'une part sur (A) la perception et d'autre part sur (B) la conception. Les six hypothèses cognitives dégagées se révéleront toutes indépendantes des unes et des autres. Leur présentation a pour objet uniquement de déceler des concepts communs principaux ou des problématiques communes dans les sciences cognitives et ceci, sans prise de position. Les philosophies de la perception et de la conception seront exposées plus précisément avec leurs critiques lors de la présentation des différents courants composant les sciences cognitives.

A - La perception

La perception se définit comme l'action de connaître par l'intermédiaire des sens. Elle est dans l'instant à la fois univoque, indubitable et intuitive. Les deux premiers termes sous-tendent la notion de jugement perceptif et le dernier la notion de jugement existentiel. Le jugement perceptif associe l'impression sensible à ce qui doit la générer. Un jugement perceptif peut prendre la forme d'un énoncé d'observation comme « Pierre a l'impression de voir une pomme ». L'univocité de cet énoncé vient du fait que Pierre ne peut avoir l'impression de voir à la fois une pomme et une poire. Si ses sensations sont confuses, à cause de la présence d'une lumière éblouissante par exemple, il pourra osciller entre les deux perceptions mais il verra soit une pomme, soit une poire. De même, le dessin de la Figure I-2 évoque à chaque instant soit un lapin, soit un canard.

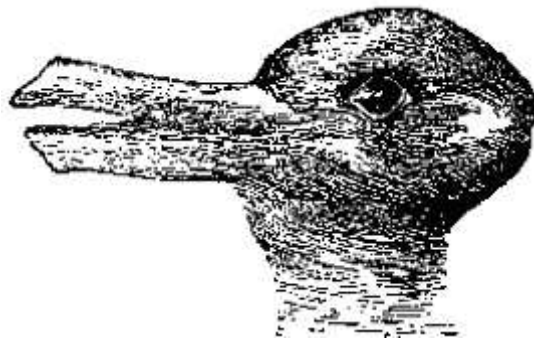


Figure I-2 : Dessin de Jastrow (1901).

Pierre peut douter de ce qui est à l'origine de son impression de pomme ou de poire mais indubitablement il a réellement l'impression de voir l'un des deux fruits. Le jugement existentiel porte sur la réalité effective de ce qui doit générer les impressions sensibles. Intuitivement, dans de bonnes conditions de visibilité, Pierre tient pour vrai « il y a une pomme ». En pratique, ces deux types de jugements deviennent transparents avec l'énoncé « Pierre voit une pomme ». La nécessité de poser des hypothèses perceptives provient du fait que nos jugements perceptifs et nos jugements existentiels changent sans que cela vienne a priori d'un changement du monde extérieur. Deux types de situations révèlent cette tension : l'illusion et l'hallucination. L'illusion se produit lorsque le jugement existentiel n'accrédite pas le jugement perceptif. Pierre voit une pomme, puis s'en approche et se rend compte que c'est une imitation en pâte à papier. Il se replace comme au début, et il a toujours l'impression de voir une pomme bien qu'il sache que cela n'en est pas une. Il peut avoir aussi des illusions dues à un biais sensoriel. L'illusion d'optique donnant l'impression dans la Figure I-3 que le cercle A est plus grand que le cercle B offre une situation analogue. L'hallucination se produit lorsque consécutivement deux jugements perceptifs différents pour une même scène. Pierre passe rapidement devant un pommier, il a l'impression de voir une pomme, il considère voir une pomme. En revenant plus tranquillement sur ses pas, il se rend compte qu'il s'agissait d'une balle retenue par les branches.

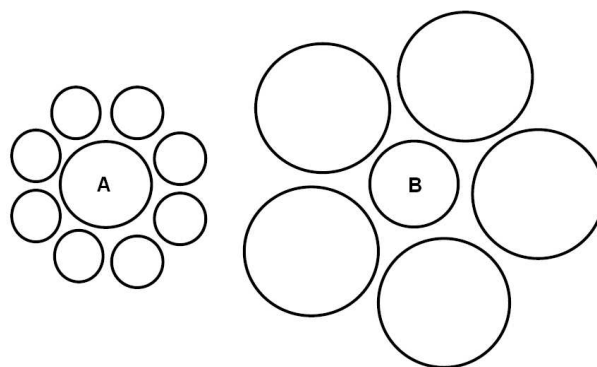


Figure I-3 : Illusion de Titchener (1898).

Deux points se dégagent de ces exemples. Premièrement, avant que l'illusion ou l'hallucination ne soit révélée, rien ne diffère avec les autres perceptions. Suivant les hypothèses ontologiques admises, ce premier point pose le problème de la possibilité de l'accessibilité d'une signification, dans l'absolu, d'un énoncé synthétique. Par ailleurs, la

signification d'E4 « Pierre voit une pomme et il y a réellement une pomme », ne change pas le fait que Pierre croit voir une pomme. L'énoncé E4 peut être vrai comme faux car la pomme est une très bonne réplique en pâte à papier, ou la pomme provient des effets de psychotropes ingérés par Pierre sans le savoir, ou l'univers y compris la pomme résulte d'une gigantesque manipulation de son esprit par quelque chose d'autre. Mais cela ne change pas le fait que Pierre a l'impression de voir une pomme et qu'il peut croire ou non à son existence. La perception n'est pas influencée par une signification extérieure au couple sujet/monde quelles que soient les hypothèses ontologiques. Par conséquent, les hypothèses perceptives doivent être centrées sur le sujet, et leurs constructions doivent être indépendantes des hypothèses ontologiques.

Deuxièmement, la révélation de l'illusion ou de l'hallucination repose obligatoirement sur deux capacités du sujet : l'exploration et la mémorisation (l'action). Ces capacités se justifient réciproquement d'une part par la nécessité de générer une différence de perception s'il y a lieu et d'autre part, par la nécessité de comparer les jugements passés et présents pour s'apercevoir de cette différence. Ces deux capacités impliquent que le monde soit suffisamment stable pour établir des relations entre les différentes perceptions dans l'espace et le temps.

La perception dans ce cadre se réduit à la simultanéité de deux jugements : perceptif et existentiel, indépendamment d'une vérité extérieure. Toutefois, les processus et les propriétés de ces jugements restent indéfinis bien qu'ils se caractérisent par trois principales interrogations. Ces interrogations, indépendantes les unes des autres, seront successivement examinées : (i) La perception est-elle totalement conceptuelle ? (ii) La perception est-elle totalement consciente ? (iii) La perception est-elle totalement indépendante des états mentaux ?

i - La perception est-elle totalement conceptuelle ?

La première interrogation invite à poser la première hypothèse sur la perception (HP1) : *le jugement perceptif révèle des qualités sensibles ou des entités qui renvoient immédiatement à des concepts*. Le contenu de la perception devient exclusivement conceptuel et cela explique nos capacités d'inférence sur les choses perçues. Autrement dit, les concepts illustrés par des images apparaissent à la conscience qui ensuite les manipule. Plusieurs interprétations de cette hypothèse existent tout comme pour les hypothèses ontologiques ; toutefois, trois remarques s'imposent. La première est que le sujet ne perçoit pas ce qu'il ne conçoit pas (Kant, 1781). Cela ne veut pas dire que Pierre devant un Martien ne voit rien, mais il percevra un ensemble de qualités sensibles telles que la couleur, l'odeur ou le son. En revanche, il ne regroupera pas cet ensemble d'impressions sensorielles renvoyant à des concepts précis au concept de Martien. Il lui faut pour cela passer par une phase de conceptualisation. La seconde remarque naît de la première, la circonscription d'impressions sensibles qui ne sont pas nommées nécessite l'introduction de concepts déictiques comme « ceci » ou « celle-là ». Ces concepts déictiques étant fortement liés à leur contexte d'utilisation, leurs interprétations sont source de débat (Dokic, 2004). La troisième remarque porte sur la définition du terme concept. Un concept peut se comprendre comme une représentation mentale abstraite et générale. En acceptant HP1, cette généralité ne se porte plus uniquement sur les choses (le concept pomme englobe toutes les pommes) mais aussi sur des impressions sensibles dénotant une chose singulière (le concept de Pierre englobe toutes les sensations produites par la présence de Pierre). HP1 oblige à distinguer concept général et concept singulier (Frege, 1892). De nombreuses critiques existent à l'encontre de HP1 (Dokic, 2004). Pierre a faim, voit une pomme, tend le bras, la saisit puis

la mange. Cette action a-t-elle, nécessite-telle une conceptualisation aussi fine de la part de Pierre ? Qu'en est-il du singe qui réalise le même type d'action avec une banane ? Sans rentrer dans les différents débats, le refus de HP1 n'implique pas de nier que des concepts puissent être mis en correspondance avec les impressions sensibles et qu'ils puissent être utilisables lors de raisonnement. Mais de manière primitive, la conscience utilise des jugements perceptifs et existentiels sans les mettre en correspondance avec un concept associé à un symbole et ses implications dans sa manipulation avec d'autres symboles. Ici, le concept fait partie d'un langage. En somme, le résultat de la perception consciente peut prendre une forme non-conceptuelle qui est étroitement liée à l'évolution des impressions sensibles et aux actions associées. Le déclenchement de l'action provient davantage de l'anticipation des perceptions que d'un raisonnement s'appuyant sur des contraintes logiques ou conceptuelles.

ii - La perception est-elle totalement consciente ?

La réponse par l'affirmative à la seconde interrogation traduit la seconde hypothèse sur la perception (HP2) : *tout ce qui est perçu vient à la conscience, autrement dit tout le traitement de l'information perceptive converge vers un centre d'analyse global*. Cette hypothèse implique deux choses, d'une part que toutes les informations mnésiques ont transité par le champ de la conscience, et d'autre part que toutes les actions conscientes s'appuient exclusivement sur des perceptions conscientes. Ces points peuvent être attaqués. Par exemple, le premier point contredit les situations dans lesquelles un sujet se souvient après coup des détails d'une scène passée qu'il avait négligés jusqu'alors. Ce genre de situation montre qu'une perception non consciente peut non seulement être présente en mémoire mais aussi être accessible a posteriori à la conscience. Refuser HP2 revient à affirmer qu'il existe plusieurs traitements de l'information fonctionnant en parallèle tout en conservant des moyens d'interagir. Il est possible de considérer, par exemple, que Wittgenstein (1921) acceptait cette hypothèse contrairement à Gibson (1979).

iii - La perception est-elle totalement indépendante des états mentaux ?

La dernière interrogation introduit la troisième hypothèse sur la perception (HP3) : *le processus de perception repose uniquement sur les données sensorielles*. L'acceptation de cette hypothèse conduit à considérer la perception comme étant passive et neutre. Néanmoins, elle peut être active dans le sens où l'attention peut être dirigée. Le sujet dirige son champ d'attention comme il le désire et récolte les perceptions se trouvant dans le champ. Mais ce qui est récolté ne dépend pas seulement des données sensorielles produites par la scène. La perception est un traitement d'informations sensorielles ascendant, des sens à l'esprit, dirigé par l'attention qui, elle, peut-être à la fois descendante et ascendante.

La Figure I-2, selon notre stratégie exploratoire, présente le dessin soit d'un lapin soit d'un canard. Un sujet averti oriente consciemment son attention (qui inclut la stratégie pour acquérir l'information) afin de percevoir l'une des deux interprétations du dessin. La perception induite par la Figure I-2 ne remet pas en cause HP3. En revanche, d'autres situations perceptives suggèrent des failles. Par exemple, Pierre et Paul vont dans un bar un peu sombre, Pierre commande un thé et Paul un café. Pierre croit que le sachet de thé se trouve déjà dans la théière et remplit sa tasse seulement d'eau chaude. Pierre boit le contenu de sa tasse sans surprise. Paul a vu toute la scène et demande à Pierre si le thé est bon. Ce dernier répond oui. Pierre a eu une hallucination, il ne s'est pas rendu compte qu'il buvait de l'eau chaude. Il est pourtant capable de différencier le thé de l'eau. Ceux qui refusent HP3 concluent que la perception utilise des informations qui ne sont pas dans la

scène présente. Par exemple, l'empirisme de Locke (1689) et le positivisme logique de Mach (1900) peuvent se comprendre comme souscrivant à HP3 contrairement à la phénoménologie de Husserl (1913).

*

Ces hypothèses et leurs négations peuvent se combiner et ouvrir de nouvelles interrogations. Mais cette présentation visait simplement à cerner les concepts élémentaires liés à la perception et utilisés par les courants composant les sciences cognitives. De la même manière, la section suivante abordera la présentation des hypothèses sur la conception.

B - La conception

La conception se comprend comme l'acte de concevoir et de manier des idées. Afin de prévenir les malentendus, il est nécessaire de distinguer idée et concept. L'idée est une représentation mentale abstraite, elle est générée par le sujet et possède donc un caractère subjectif quelle que soit la position ontologique adoptée. Pour certaines positions philosophiques, la notion d'idée peut recouvrir la notion d'image mentale ou toute autre évocation sensorielle par mémoire ou par imagination. Les énoncés suivants, compris par un sujet, renvoient à des idées : « un chat » et « le chat chasse les souris ». Une pensée, pour reprendre la terminologie de Frege (1892), devient une idée de forme propositionnelle comme « le chat chasse les souris » où une signification est possible. Concernant le concept, la nature de celui-ci varie selon l'acceptation ou le refus de l'existence d'un monde conceptuel : soit les concepts existent par eux-mêmes et sont objectifs, soit les concepts n'existent pas par eux-mêmes et sont intersubjectifs. Dans le premier cas, l'idée se veut d'être une approximation subjective du concept objectif. Ainsi une idée peut être identique pour deux sujets différents parce que leur idée reflète le même concept. Selon une position qui accepte HOI et qui suppose une capacité extrasensorielle portant sur l'idéal, Pierre voit une pomme. De même que les impressions de pomme sont un aperçu de la pomme, l'idée de ce qu'est une pomme est un aperçu du concept de pomme. Dans le second cas, les idées de deux sujets sont considérées comme identiques parce qu'elles aboutissent aux mêmes conclusions pratiques ou théoriques. Pierre demande à Paul de couper la pomme en deux, Paul s'exécute. Pour Pierre, Paul possède les mêmes idées à la fois pour la pomme et pour l'action de couper en deux.

Le manque de sensation ne permet pas de réutiliser les notions de jugements perceptif et existentiel. Pourtant, la tension révélée par les illusions et les hallucinations ressemble aux erreurs de raisonnement réalisées de bonne foi : les paralogismes. Ces erreurs peuvent être commises au quotidien, mais un exemple issu des mathématiques est plus neutre vis-à-vis du monde physique. Le problème provient du travail de Cantor sur la théorie des ensembles (Volken, 2003). Dans la Figure I-4, y a-t-il plus de points sur un segment de longueur 1 ou dans un carré de côté 1 ? Le segment étant inclus dans le carré, un sujet naïf donne intuitivement le carré pour bonne réponse. Or, en posant rigoureusement les implications des concepts mis en jeux, une construction permet d'aboutir à la conclusion qu'il y en a autant, c'est-à-dire qu'il existe une bijection entre les points constituant le segment et ceux constituant la surface. Malgré la démonstration, la première intuition demeure. Cette situation rappelle celle des illusions sensorielles où les cercles de la Figure I-3 semblent de tailles différentes même si l'illusion est connue.

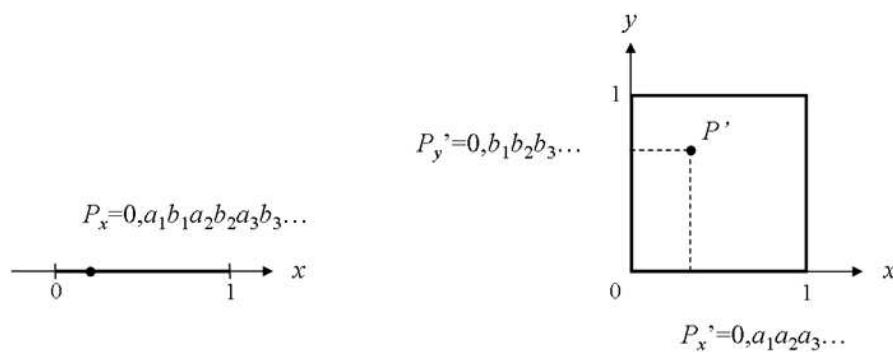


Figure I-4 Le problème de Cantor : il existe une bijection entre les points formant le segment de gauche et les points formant la surface du carré de droite comme le suggère la construction d'une correspondance univoque entre leurs coordonnées.

Pour expliquer cette situation, deux types de jugements analogues à ceux de la perception peuvent être dégagés : le jugement propositionnel et le jugement de vérité. Le jugement propositionnel est le résultat du processus qui formule une idée. Le jugement de vérité est la signification du jugement d'assertabilité. Tout comme pour la perception, ces deux jugements sont transparents. Pierre propose une tarte à ses six invités, deux déclinent la proposition, combien Pierre doit-il couper de parts sachant qu'il souhaite en manger ? Pierre se dira certainement que la tarte doit être coupée en cinq. Cette pensée vient d'un jugement d'assertabilité et implicitement Pierre en admettra la véracité : le jugement de vérité. Deux types de paralogismes apparaissent. En premier lieu l'illusion, auquel cas le jugement d'assertabilité parvient à se maintenir malgré un jugement de vérité qui lui est défavorable (paralogisme illustré par le problème de Cantor). En second lieu l'hallucination, dans ce cas le jugement d'assertabilité se révèle finalement sans fondement, autrement dit avec l'impossibilité de retrouver les prémisses de sa construction. Remarquer ce genre d'erreurs nécessite également des capacités de mémorisation et des capacités d'exploration. Ces capacités d'exploration ne reflètent pas ici une action dans un espace mais une évaluation des conséquences, des implications lors de la manipulation des idées.

Les ressemblances de structure dans la présentation de la conception et de la perception ne suffisent pas à affirmer une symétrie sur leur nature. Cette dernière signifierait une symétrie entre HOP et HOI alors que ce n'est pas le cas. Le monde physique est un postulat qui autorise la formulation de HOP tandis que HOI est l'hypothèse d'un monde idéal. Les illusions et les hallucinations montrent seulement l'existence d'une mémoire et d'une capacité à agir soit physiquement soit mentalement. La raison du postulat du monde physique est quelque chose qui s'impose à soi. Or, rien ne s'impose à l'action mentale, la rigidité d'un raisonnement vient seulement des règles d'inférence admises par le penseur. Admettre HOI, et même supposer une sensibilité extrasensorielle qui permette juste de reconnaître la vérité parmi les idées produites, ne remet pas en cause cette liberté. À souligner également que les limites de l'imagination ne sont pas celles de la raison. Un espace non euclidien peine à être imaginé bien que sa conception en termes mathématiques soit possible.

La symétrie de structure se poursuit par l'établissement de trois interrogations sur la conception : (i) est-elle seulement un processus d'inférence ? (ii) est-elle totalement consciente ? (iii) est-elle indépendante des émotions ? La réponse par l'affirmative à ces interrogations amène trois hypothèses, chacune indépendante les unes des autres ainsi que des hypothèses ontologiques.

i - La conception est-elle seulement un processus d'inférence ?

Répondre oui à la première question revient à tenir la première hypothèse sur la conception (HC1) : *les idées se produisent par l'induction et la déduction à partir d'autres idées dont la perception est à l'origine*. La première de ces opérations mentales consiste à passer du particulier au général. Pierre coupe les mangues apportées par Paul et voit que leur chair est colorée d'un jaune orangé, il en déduira qu'il en va de même pour toutes les mangues. La seconde consiste à suivre les implications associées aux idées manipulées. Pierre pense que la chair de toutes les mangues possède une teinte jaune orangée, Paul revient avec une nouvelle mangue, Pierre en déduit la couleur de la chair de celle-ci. Nier HC1 revient à introduire un principe de création d'idées qui ne relève ni de l'induction ni de la déduction à partir des propriétés des éléments du problème. Peirce (1934) propose le terme d'abduction : cette opération mentale génère une idée compatible avec d'autres sans pour autant en être déductible ou inductible, c'est une idée ad hoc. Pierre n'était pas là lorsque Paul a posé la nouvelle mangue sur la table. Aucun indice n'indique que c'est Paul qui l'a apporté mais Pierre le suppose quand même.

ii - La conception est-elle totalement consciente ?

La conception est-elle totalement consciente ? Cette interrogation est analogue à la seconde interrogation concernant la perception. Se construit alors sur le même schéma que HP2 la seconde hypothèse sur la conception (HC2) : *toutes les idées conçues viennent entièrement et immédiatement à la conscience, et le raisonnement s'applique sur elles seules, autrement dit, tout le processus d'inférence est centralisé et clos*. Le raisonnement s'appuie exclusivement sur des idées conscientes, le flot d'idées se générant par lui-même et par les idées provenant de la perception. En somme, HC2 considère que le jugement d'assertabilité est entièrement conscient. Malgré la forte impression qu'a un sujet de contrôler son discours, trois points remettent en question HP2. Le premier point est qu'un sujet n'a pas conscience de tout ce qu'il sait. Le second est que la justification d'une assertion peut s'appuyer sur des inférences implicites. L'énoncé E10 « Pierre ne doit pas conduire, il a trop bu » est supporté par un grand nombre d'hypothèses qui n'ont pas été formulées mais qui sont implicitement liées et dont la présence peut être dévoilée à la conscience au fil du discours ; voici quelques exemples d'hypothèses implicites : Pierre sait conduire, Pierre a bu de l'alcool, l'alcool diminue les performances, Pierre est en France où la loi condamne le conducteur ivre. Ainsi, toute connaissance est liée à d'autres connaissances en un réseau (Quine, 1960). Le dernier point apparaît quand la solution à un problème resté sans solution vient à l'esprit bien que l'attention se porte sur autre chose. Cette situation suggère une activité mentale en dehors du champ de la conscience, en parallèle ou pendant le sommeil.

iii - La conception est-elle indépendante des émotions ?

La dernière interrogation mérite d'être développée puisqu'en effet l'influence de l'émotion sur la raison n'est plus à démontrer : des propos avancés sous la colère peuvent paraître ridicules rétrospectivement. Cette influence, en revanche, se réduit seulement à une forte perturbation dans l'enchaînement des idées. Le terme « indépendance » de cette troisième interrogation se comprend comme la non participation de façon positive (constructive) des émotions aux raisonnements. Cela revient à poser la troisième hypothèse sur la conception (HC3) : *le raisonnement se construit exclusivement sur des idées issues soit de la perception sensorielle, soit de la combinaison d'autres idées*. Autrement dit, la perception d'émotions ne débouche pas sur de réelles idées et freine la conception. HC3 considère alors la décision comme le résultat des traitements intellectuels s'appuyant sur des idées abstraites

associées à des idées d'origines sensorielles présentes, passées ou anticipées. En refusant HC3, l'émotion prend part au processus de décision comme le soutient Damasio (1997). Cependant, la nature des émotions reste à déterminer ainsi que les mécanismes sous-tendant leur influence. Deux pôles dominent le débat : soit les émotions jouent le rôle de médiateur entre des classes de perceptions et des classes de comportements et de raisonnements : Marie voit son chat mourir écrasé, elle devient triste donc elle pleure ; soit l'émotion traduit un mode de fonctionnement concernant aussi bien la perception que la conception. Ainsi l'idée d'une émotion s'identifie a posteriori par le résultat de ce mode : Marie voit son chat mourir écrasé, elle pleure donc elle est triste. L'émotion dans le second cas devient une modalité de fonctionnement sur lequel s'effectue une perception des émotions. Il reste encore à définir les caractéristiques de ces modalités ainsi que leur influence et les mécanismes de transition entre ces modalités.

**

L'ensemble des hypothèses formulées offre une grille de lecture des positions philosophiques employées suffisamment large pour contenir celles qui participent de près ou de loin au courant de la robotique cognitive. L'indépendance des hypothèses présentées permet de définir 4 familles métaphysiques, 64 positions concernant la perception et autant pour la conception. Ainsi, en considérant ces thèmes indépendamment, le nombre de combinaisons atteint 256. Bien qu'il faille prendre en considération le fait que certaines hypothèses font l'objet d'un fort consensus, ce qui réduit le nombre de combinaisons, ce chiffre élevé donne une explication sur la difficulté d'établir un consensus sur un paradigme unique d'autant plus que les discours emploient souvent ces hypothèses de façon implicite. Par ailleurs, chaque hypothèse ou combinaison d'hypothèses amène des problématiques uniquement compréhensibles en leur sein. Mais l'étude et la classification des différentes positions philosophiques existantes à partir de cette grille d'analyse constituent un travail philosophique à part entière. Ici, l'application de cette grille d'analyse vise l'étude des idées composantes les sciences cognitives et plus particulièrement leur conséquence pour les sciences de l'artificiel. Dans cette perspective, la grille d'analyse doit être complétée par les types de discours et les méthodes d'investigation liées à toute activité scientifique.

2.2. Les notions liées à la description

Le terme « système » traduit ici un ensemble de phénomènes considéré comme un tout en raison de leur organisation et de leur forte interrelation. Cette définition se veut vaste afin de recouvrir tous les champs scientifiques. Un système résulte d'un découpage perceptif qui le distingue du reste des perceptions : l'arrière plan. La description d'un système revient à dégager les observations et les concepts associés afin de le détailler et de mieux le cerner. La question concernant la justification de ce découpage du monde sera abordée dans le chapitre suivant. L'ensemble des descriptions d'un système forme un modèle. Une réflexion sur le monde peut alors s'appuyer sur ce modèle comme une représentation conceptuelle suffisante d'un aspect de la réalité extralinguistique. Si un modèle ne vise plus un système particulier mais un ensemble de systèmes, alors les caractéristiques générales constituent une théorie. Mais cette théorie peut être interprétée comme un modèle du monde compris comme un méta-système. L'activité scientifique s'attache à constituer collectivement des modèles et des théories. Le jugement concernant la valeur ou la finalité d'une telle entreprise varie selon les hypothèses ontologiques adoptées.

Les trois types de description employés par les modèles ne conditionnent pas le type d'argument utilisé pour constituer et légitimer un modèle. En revanche, les propriétés liées aux systèmes modélisés conditionnent le choix du type d'argument. Les trois principaux types d'arguments résultent de la compétition qui caractérise l'activité scientifique. Cette compétition provient de l'obligation de convaincre d'autres personnes pour obtenir des ressources intellectuelles et matérielles nécessaires à la poursuite des investigations. En ce qui concerne l'activité scientifique contemporaine, Latour (2001) analyse le capitalisme scientifique comme étant la résultante des cycles de la crédibilité, Figure I-5 :

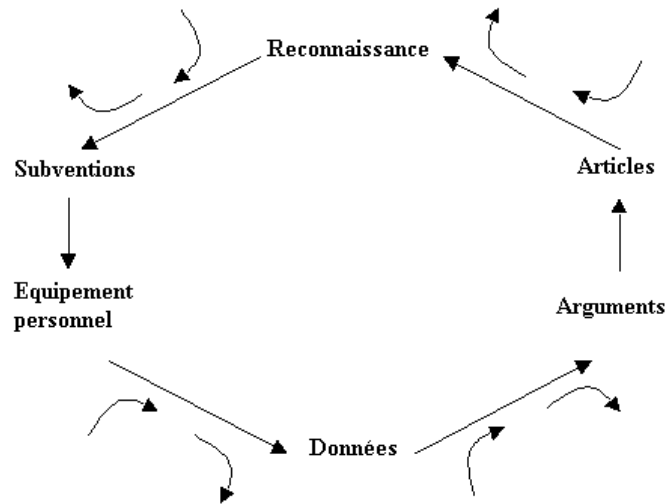


Figure I-5 : Le schéma des cycles de crédibilité selon Latour (1995).

Chaque partie du cycle peut être fatale aux recherches. Les arguments présentés ici sont aussi des méthodes d'investigations : (A) méthode hypothético-déductive, (B) méthode par la cohérence, (C) méthode artefactuelle. La connaissance de ces méthodes offrira un éclairage particulier sur les avantages de l'approche robotique concernant l'étude de l'autonomie.

2.2.1. Les trois types d'investigation

A - La méthode hypothético-déductive

La première méthode, la méthode hypothético-déductive (MHD), consiste à poser un ensemble d'hypothèses sur les relations ou la constitution d'un système considéré, puis d'en déduire certaines conclusions que l'observation peut corroborer ou réfuter. Popper (1934) souligne que l'interprétation dans les deux cas est asymétrique. Un résultat en accord avec les prédictions ne suffit pas à prouver que la description capture parfaitement le réel, seulement que les hypothèses ont été compatibles avec la réalité considérée à un moment donné. En revanche, l'inverse signifie que la théorie est fautive. Toutefois, si un modèle est corroboré plusieurs fois, alors un certain crédit peut lui être accordé.

Cette méthode est particulièrement adaptée à la construction de concepts qui ne renvoient pas directement à l'observation. Cette création devient nécessaire dans la description lorsque, malgré une apparence identique, le comportement du système diffère. Une règle de plexiglas attire ou non les bouts de papier sans changer d'apparence pour les

sens. Une solution est de rajouter une caractéristique supplémentaire, un aspect caché mais révélé indirectement. Cependant, l'introduction d'un concept supplémentaire, dans le cadre de la MHD, doit respecter deux contraintes : celle de précision et celle de portée explicative. La première découle du fait qu'une explication vague confère à la prédiction de trop grandes incertitudes qui n'autorisent alors pas à conclure sur le résultat. La seconde contrainte vient du fait qu'un concept expliquant tout, quel que soit le résultat de l'expérience, n'autorise plus non plus à conclure. Par exemple, l'énoncé « Dieu décide quand le plexiglas attire ou non les bouts de papier » ne permet pas d'imaginer une expérience qui pourrait l'invalider.

La reproduction d'une expérience sur le système se révèle être une contrainte forte. Une expérience est une observation provoquée selon un protocole. Ce dernier doit être suffisamment précis pour constater des régularités dans les observations successives. L'expérience doit être appliquée au même système ou à un système identique en tout point.

Les champs d'études autorisant l'expérimentation forment les sciences empiriques, en opposition avec les sciences historiques. Ces dernières doivent alors s'appuyer sur d'autres types de méthodes.

B - La méthode par cohérence

La deuxième méthode, la méthode par cohérence, se trouve utilisée par toutes les sciences historiques. Celles-ci portent sur tous les systèmes dont l'expérimentation n'est pas de l'ordre de l'humain ou dont le comportement dépend de son vécu. Par exemple, les fossiles ainsi que les données issues de la génétique offrent des informations pour théoriser l'évolution de l'espèce considérée dans son ensemble comme un système. Des simulations informatiques de certains mécanismes théoriques de l'évolution, qui feront l'objet d'une analyse critique dans la troisième partie de ce chapitre, peuvent participer à l'investigation ; en revanche l'expérimentation concerne uniquement la simulation informatique d'un modèle. L'expérience qui consisterait à revenir à l'état initial de la terre pour observer l'influence des contingences historiques appartient à l'imaginaire. Par ailleurs, le système humain avec ses capacités d'apprentissage peut ne plus jamais se comporter de la même façon après certaines mises en situations. Ces deux cas montrent une impossibilité à reproduire certaines expériences.

La modélisation de ce genre de systèmes s'appuie alors exclusivement sur l'observation et sur la cohérence des modèles. Contrairement à la MHD qui se focalise sur une observation bien précise et contrôlée du système étudié, la méthode par la cohérence (MC) cherche à récolter l'ensemble des observations sur le système. L'ensemble de ces observations forme alors un réseau de contraintes, sur lequel s'élabore un modèle cohérent du point de vue des indices et du point de vue de la logique interne. Le modèle résultant peut alors conduire à la déduction de certaines observations qui le corroboreront ou l'invalideront. De manière identique à la MHD, les concepts invoqués pour maintenir la cohérence d'un modèle ou d'une théorie doivent fournir des attentes sur des observations suffisamment précises pour conclure, qu'elles soient réalisables ou non.

C - La méthode artefactuelle

La méthode artefactuelle (MA) consiste à reproduire matériellement l'ensemble des propriétés et des relations décrites par le modèle. La structure ainsi réalisée doit mimer les caractéristiques considérées du système soit dans les faits, ce qui correspond à une simulation effective (par exemple l'utilisation de mannequins pour les crash-tests) soit

métaphoriquement médiatisée par les mesures, ce qui correspond à une simulation logique (par exemple le calcul de trajectoire de trois astres). La simulation est limitée par l'expressivité du modèle sur lequel elle s'appuie. Elle permet de vérifier que le niveau d'explication du modèle correspond bien aux attentes. Le résultat peut soit corroborer le modèle, soit l'invalider, comme pour la MHD. Toutefois, à la corroboration s'ajoute au modèle un gage d'objectivité dans le sens où la mécanique des relations relevées par le modèle se déroule sans un interprétant. Le terme « objectif » se rapporte à la constatation d'une similitude mais non à une existence conceptuelle objective du modèle dans un monde idéal. Dans le cadre de la compétition entre les modèles ou théories, qui pousse à ce que les énoncés deviennent plus neutres vis-à-vis des opinions de chacun afin de convaincre le plus grand nombre, l'implémentation d'un modèle est l'aboutissement de cette logique. Par sa réalisation, le modèle affirme son indépendance vis-à-vis des diverses interprétations. Les instruments de mesure sont également des réalisations concrètes de théories.

2.2.2. Les trois types de description

A - La description procédurale

La description procédurale (DP) porte sur les relations entre les événements observés ou supposés. Ce mode de description mobilise uniquement un enchaînement de règles pour caractériser le système étudié. Les règles tendent à prendre la forme de propositions logiques. Elles associent des prémisses à des conclusions qui peuvent à leur tour devenir prémisses : R1 « Si Paul a mangé trop de cerises, alors il va être malade. », R2 « Si Paul est malade alors il va désirer se reposer ». En restant dans le cadre de la description, une prémisses annonce une conclusion mais sans revendiquer être la cause réelle de ce qui est conclu. L'ensemble des descriptions se définit extensionnellement, c'est-à-dire que le rajout de règles s'effectue au gré des situations nouvelles. S'appuyant sur le langage naturel, cette forme de description est souple mais peut aussi se révéler très imprécise : d'un point de vue quantitatif que signifie « trop de cerises » ?

A noter que, si la description procédurale d'un même système peut se traduire par un ensemble de symboles et d'implications qui ne souffre pas de contradictions, alors il est possible de simuler le comportement du système. En effet, Turing (1937) a prouvé que si une fonction est calculable, alors il existe une machine abstraite pour la calculer. Celle-ci possède deux mémoires : la première possède un contenu modifiable, la seconde un contenu figé pendant le déroulement du calcul. La première mémoire contient les données initiales et finales ainsi que toutes les valeurs intermédiaires. La seconde mémoire contient toutes les instructions pour modifier la première en fonction de la valeur de l'emplacement mémoire désigné dans celle-ci. Cette architecture correspond à celle des ordinateurs. Ainsi, toute description procédurale pouvant se mettre sous forme logique peut être simulée (isomorphisme de Curry-Howard).

B - La description mathématisée

La description mathématisée (DM) essaye également de capturer les relations entre les événements observés ou supposés mais en remplaçant le langage naturel par un langage mathématique. Ce langage mathématique offre une définition précise des concepts employés ainsi que de leurs utilisations. Ces concepts tels que variable, opérateur, fonction, etc. permettent de formuler des équations qui décrivent de manière concise un

comportement du système. La manipulation de leurs symboles étant codifiée et mécanisée, la conclusion des manipulations ne souffre pas d'une interprétation subjective. Surtout, ces équations peuvent offrir une description d'un système sur l'ensemble des possibles, la description est intensionnelle. La tension U aux bornes d'une résistance électrique R suivant l'intensité du courant I qui la traverse se traduit par l'équation $U=R*I$, quels que soient R et I , cette équation prédit U sans pour autant avoir testé toutes les valeurs.

C - La description componentielle

La dernière description, la description componentielle (DC), porte sur l'organisation des choses observées ou supposées. Ce mode descriptif consiste soit à dégager des règles de classification pour ordonner et répertorier les éléments appartenant au système, soit à identifier les composants du système ainsi que leur place dans celui-ci. Dans le premier cas, le choix d'une règle de classification ne doit pas souffrir d'exception et doit procéder par analogie sur un ensemble de données. La taxinomie repose sur cette méthode. Dans le second cas, la description vise à circonscrire des sous-systèmes et à établir leurs connexions. De cette manière, si la description procédurale ou la description mathématisée de ces sous-systèmes existe, alors leurs combinaisons donnent une description procédurale du système les englobant. Cette démarche est tributaire de la schématisation des sous-systèmes.

3. Analyse critique de la cognition artificielle

L'ensemble des concepts évoqués jusqu'à présent constitue une grille d'analyse permettant d'étudier les nombreuses positions abordant la question de la description d'individus autonomes et de dégager l'origine de l'impasse commune afin de la dépasser. En reprenant la définition donnée au début de ce chapitre, un agent autonome est un système qui tente de suivre les lois qu'il s'est construites au fil de son activité, ce qui rend l'agent autonome historique et unique. Ces caractéristiques conditionnent directement l'utilisation des trois méthodes d'investigation pour la description du comportement et par conséquent conditionne tout projet d'artificialisation de la cognition notamment en robotique.

Le behaviorisme radical incarne l'utilisation de la première méthode d'investigation, la MHD. Pour décrire le comportement de l'individu, le behaviorisme impose deux hypothèses fortes qui vont se révéler intenable. La première hypothèse consiste à réduire l'autonomie à la capacité d'apprentissage de règles comportementales sous forme de couples stimulus externe puis réponse. Cet apprentissage s'élaborerait mécaniquement par la corrélation entre les stimuli, les actions et les renforcements négatifs ou positifs. L'agent est bien historique, c'est-à-dire que son actuel comportement dépend de son vécu, mais il se trouve entièrement déterminé par l'extérieur. Cependant, dans certaines conditions, la mise en évidence de variables internes influençant l'apprentissage ne permet pas de conserver cette hypothèse. Un rat apprendra beaucoup plus vite un labyrinthe avec récompense à la clé s'il est privé de nourriture que s'il est rassasié (Blodgett, 1929). La seconde hypothèse porte aussi sur l'apprentissage : celui-ci est guidé par des conditions externes identifiables parce que cet apprentissage est lié directement à la survie ou à la reproduction de l'agent. Selon cette hypothèse, la connaissance des intérêts physiologiques de l'animal permet de cerner les facteurs guidant l'apprentissage. Cependant, les travaux sur le rat de Tolman (1948) montrent l'existence d'un apprentissage latent, c'est-à-dire un apprentissage sans renforcement positif ou négatif. Ces deux exemples remettent en cause

l'approche MHD pour caractériser globalement un agent autonome. Plus précisément, ils illustrent deux raisons fondamentales. La première raison vient du fait que la caractérisation d'une variable interne n'est possible que dans un cadre expérimental précis. Or, ce genre de situation ne correspond pas à la variabilité naturelle des situations que l'agent rencontre. Dans ces conditions, il semble impossible de discerner le nombre, le rôle et l'interdépendance des variables cachées, d'autant plus qu'à ces difficultés se rajoutent la vicariance et la variabilité interindividuelle. La seconde raison vient de l'impossibilité de savoir ce que perçoit et ce que considère un agent autonome comme l'illustre Nagel (1974) montrant l'incapacité pour un humain de connaître l'effet que cela fait d'être une chauve-souris. De ce fait, comment imaginer un modèle décrivant le comportement global d'un agent si les éléments et les intérêts sur lesquels il se fonde sont inconnus ? La MHD se trouve adaptée pour la caractérisation de certaines capacités cognitives comme la mémoire ou le conditionnement, elle permet de rendre compte de certaines facettes du système mais pas de les unifier pour définir un modèle global qui puisse, dans de nouvelles situations, créer de nouveaux comportements.

Pourtant, c'est ce que semble faire l'humain tous les jours avec ses congénères et autres animaux. Pour cela, il s'appuie sur des concepts relatifs aux croyances, aux intentions, aux émotions (Dennett, 1991) et utilise une méthode d'investigation plus proche de la méthode par cohérence (MC). L'observateur projette sur l'agent des états mentaux et des contenus intentionnels à partir desquels il prédira le comportement futur. Cette interprétation est le résultat soit d'une simulation interne à partir de son propre vécu, autrement dit « se mettre à la place de », l'empathie fusionnelle ; soit d'un raisonnement développé à partir de règles psychologiques issues de son vécu, de sa culture ou de principes innés, c'est-à-dire une empathie distale ; soit enfin, un mélange des deux. Naturellement, la révision des croyances ou des règles psychologiques s'effectue en fonction de la qualité des prédictions. L'assignation de croyances et d'intentions aux gens, aux animaux et même aux choses se révèle être un acte familier. Cependant, l'animisme ne possède plus sa place dans la démarche scientifique, puisque les multiples interprétations sur l'attitude des esprits et autres génies de la nature ne permettent pas l'accord du plus grand nombre. Toutefois, sans utiliser les concepts liés aux états mentaux et aux contenus intentionnels, la MDH n'offre pas de théorie pratique pour les sciences humaines et certains aspects de l'éthologie. Cette incapacité continue à justifier l'emploi de ces concepts. Les sciences interprétatives s'opposent aux sciences explicatives. Les premières utilisent des raisons pour expliquer le déroulement de certains événements, les secondes recourent à des causes directement liées aux propriétés physico-chimiques (Soler, 2000). Les sciences humaines s'appuient principalement sur la MC en affirmant que les notions de croyance et d'intention se révèlent communes aux humains et de ce fait qu'il est possible de les prendre en tant que données.

La troisième méthode, la MA, propose une autre perspective pour dépasser ce blocage méthodologique concernant l'étude de la cognition. La MA ne souhaite pas seulement produire une expérience analogue à une autre, mais également concevoir au final un agent autonome avec toutes ses implications. En effet, la volonté de concevoir un agent introduit de nouvelles contraintes liées à sa réalisation, contraintes sur lesquelles les théories peuvent se développer. Leurs implémentations montrent à la fois une validation de principe et la puissance de la théorie. Cependant, l'évaluation du système réalisé passe par un protocole expérimental. La conception de l'agent se trouve fortement orientée vers cet objectif, ce qui amène une contradiction puisque, cela revient finalement à admettre que la méthode MHD est suffisante pour décrire un agent autonome ou pour juger de son

autonomie. Chaque approche explique cette contradiction comme une étape nécessaire vers la conception d'un agent autonome qui aura des réactions imprévisibles et pertinentes lors de situations inconnues. Les conséquences de la notion d'évaluation dans la conception constitueront l'un des traits majeurs de la présentation des principales approches.

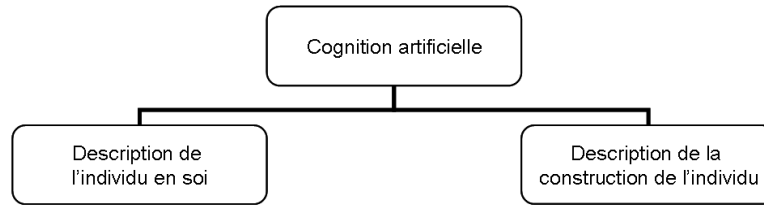


Figure I-6 : Premier embranchement de l'arborescence des approches constituant la cognition artificielle proposé dans ce chapitre.

L'histoire des idées en sciences cognitives se révèle d'autant plus complexe qu'elle couvre des domaines très divers et des sujets pouvant être abordés de multiples façons. Cependant, l'effort de synthèse pousse à sacrifier une herméneutique ou une généalogie des sciences cognitives. En conservant la perspective de la cognition artificielle, la constitution d'une arborescence des courants commence par une première interrogation (Figure I-6) : la description de l'individu en soi est-elle suffisante pour la compréhension du comportement ou faut-il se focaliser davantage sur les mécanismes construisant l'individu ? A partir de cet embranchement, l'ensemble des approches en sciences cognitives sera présenté et analysé grâce à la grille de lecture dégagée. Le détail de l'arborescence apparaîtra au cours des analyses afin de mieux situer les nœuds conceptuels. Par ailleurs, un tableau récapitulatif sur la position des approches vis-à-vis des hypothèses et des méthodes constituant la grille d'analyse conclura sur chaque courant des sciences cognitives. La synthèse de l'ensemble de ces tableaux permettra de déceler les divergences et les convergences. Malgré les différentes définitions de la cognition, l'objectif consiste à comprendre les difficultés à homogénéiser les paradigmes de la cognition artificielle, afin de déterminer l'origine de l'impasse conceptuelle de celle-ci.

3.1. La description de l'individu en soi

Comprendre l'individu en tant que tel implique qu'il existe une distinction forte voire ontologique entre celui-ci et son environnement, autrement dit l'individu s'assimile à un système traitant les informations venant de l'extérieur. Dans ce cadre, la description d'un individu se trouve directement confrontée à la problématique de la valeur de l'explication causale entre un énoncé issu des sciences interprétatives et des sciences explicatives concernant le comportement d'un individu. Comment comprendre les deux énoncés suivant : « Pierre regarde dans le placard parce qu'il cherche un tire-bouchon » et « Pierre regarde dans le placard parce que la configuration de l'ensemble de son système nerveux l'y a conduit » ? Le premier énoncé fait référence aux croyances et désirs de Pierre et le second fait référence aux propriétés physico-chimiques des éléments qui le composent. Le dualisme de Descartes (1637) propose que ces deux énoncés renvoient à deux réalités de nature distincte, et que chacun des énoncés décrive des causes réelles. Toutefois, afin de conserver la primauté de l'esprit pour agir, il faut accepter le fait qu'il puisse y avoir certaines interactions entre ces deux plans de réalités. Or, ce dernier rajout de la position dualiste la rend contradictoire. En effet, deux réalités distinctes ne peuvent interagir sinon elles forment une même réalité. En admettant cette contrainte et en n'acceptant qu'une

seule réalité causale, la physique, il reste néanmoins difficile de nier qu'un sujet agit en fonction de ses désirs. La philosophie de l'esprit a beaucoup contribué à éclaircir le débat et à proposer de nombreuses solutions. Les présentations des positions qui suivent s'appuieront sur les deux principales familles (Figure I-7) : le mentalisme et l'éliminativisme.

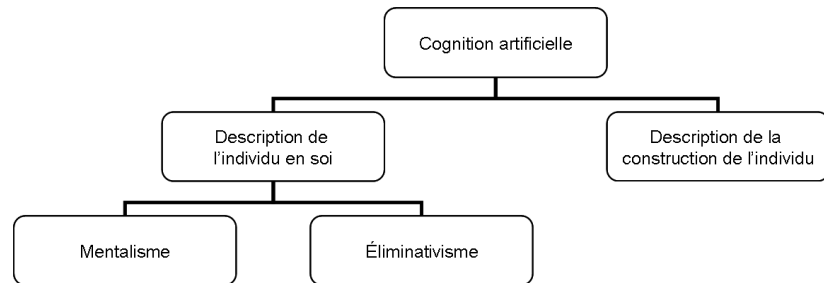


Figure I-7 : Le mentalisme et l'éliminativisme au sein de l'arborescence de la cognition artificielle.

3.1.1. Le mentalisme

Le mentalisme accepte que les seules causes efficientes soient dues à la physique, toutefois, il suggère que la compréhension des raisonnements donne un niveau d'explication plus pertinent dans le sens où l'organisation de la matière semble juste permettre l'incarnation des concepts, c'est-à-dire un dualisme de propriété (Searle, 1984) et non de substance. L'esprit se réduit à un ensemble de calculs portant sur des représentations symboliques dont le résultat aboutit à l'action si les conditions physiques l'autorisent. Dans ce cadre, le cerveau se compare à l'ordinateur et l'esprit au programme. L'avantage de l'étude de l'esprit sur le physique se justifie ici par le fait qu'un même algorithme puisse s'effectuer avec des supports matériels différents. Parmi les courants revendiquant cette position, deux formes se dégagent (Figure I-8) : le symbolisme, et le subsymbolisme.

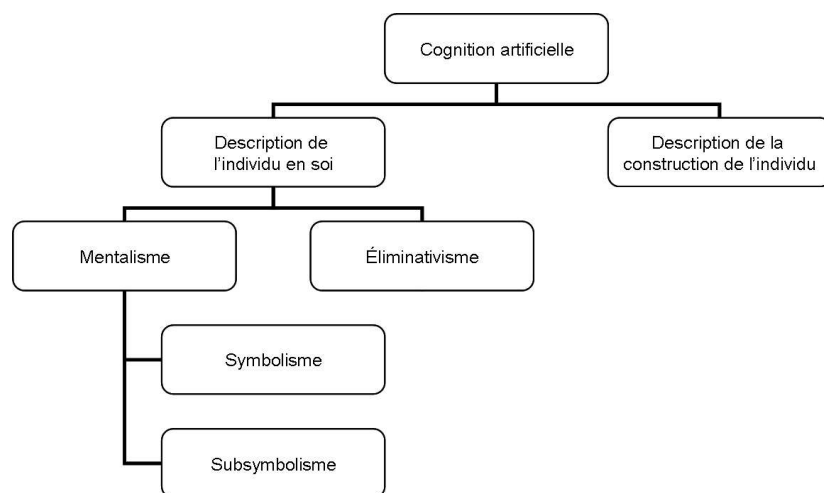


Figure I-8 : Le courant du symbolisme et celui du subsymbolisme au sein de l'arborescence de la cognition artificielle.

Bien que le mentalisme revendique le rôle fondamental des états mentaux au sens large, il s'inspire principalement des concepts informatiques et mathématiques. Ces derniers

vont profondément influencer la manière d’aborder la psychologie cognitive et plus particulièrement en ce qui concerne la notion de représentation mentale. Cette section sur le mentalisme abordera donc uniquement le problème de la représentation mentale en amont des théories psychologiques. En revanche, l’étude de ces théories psychologiques sur les représentations mentales se révélera indispensable pour la compréhension des approches interactionnistes abordées dans la partie consacrée aux courants souhaitant décrire la construction de l’individu.

3.1.1.1. Le symbolisme

Le symbolisme appartient à la famille philosophique FCP, c’est-à-dire la famille pour laquelle le langage reflète ontologiquement des objets du monde. Autrement dit, comprendre le langage revient à comprendre le monde. Chaque proposition ou concept vrai trouve une correspondance dans le monde. À ce niveau, deux voies s’ouvrent à la réflexion : (A) soit le langage naturel possède les propriétés nécessaires et suffisantes, et dans ce cas il faut les découvrir ; (B) soit le langage naturel est imparfait et dans ce cas il faut trouver un meilleur langage. Sans s’ignorer, les partisans de la première voie vont s’orienter vers la linguistique (Saussure, 1916), ceux de la seconde vers la logique (Frege, 1892). Bien que l’intelligence artificielle ne soit pas leur principale ambition, ces courants de pensées ont influencé certaines de ces conceptions. Le langage, comme échange de symboles dans ces deux positions, privilégie la description procédurale (DP) pour la construction des modèles cognitifs. Cette section souhaite montrer les conséquences de ces deux influences majeures et leurs apports à la réflexion sur l’autonomie robotique (Figure I-9).

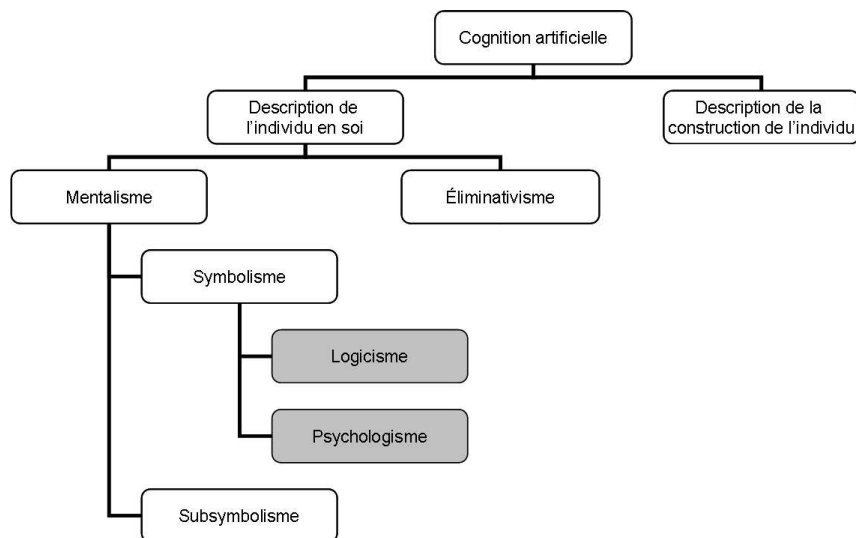


Figure I-9 : Les approches issues du logicisme et celles issues du psychologisme au sein de l’arborescence de la cognition artificielle.

A - L’intelligence artificielle et le langage naturel

L’influence de la première voie se révèle dans la conception de l’intelligence artificielle de Turing (1950). Celui-ci a prouvé que n’importe quel langage avec un nombre fini de symboles et avec un nombre fini de règles d’inférence peut se dérouler mécaniquement à l’aide d’un simple dispositif : un ruban infini sur lequel se juxtaposent en longueur des cases régulières et dans lesquelles peut être inscrit un seul symbole, une tête

de lecture qui ne lit qu'une case à la fois ; une liste de symboles, une liste d'états internes, des actions élémentaires et une table mettant en correspondance un symbole, un état interne et une action. Dans l'optique selon laquelle la pensée n'est que manipulation de symboles, lorsque les règles du langage naturel seront répertoriées, il restera à fournir à la machine les concepts élémentaires pour que son apprentissage soit similaire à celui d'un enfant. Turing propose que l'évaluation de l'intelligence de la machine se fasse de la manière suivante : Un être humain A discute par le biais d'un ordinateur avec deux interlocuteurs, l'un est un être humain B, l'autre un programme informatique C. Si A ne peut distinguer la différence de nature entre C et B, alors ils sont équivalents. Tout n'étant que verbe, les hypothèses concernant la perception deviennent superflues. En revanche, toutes les hypothèses concernant la conception doivent être acceptées.

Les études sur les grammaires deviennent nécessaires dans cette démarche ainsi que l'élaboration d'ontologies afin de mettre en relation les concepts entre eux. Toutefois, l'objectif étant de donner l'illusion d'une réelle interaction avec un être humain, quelques règles peuvent suffire pour certains types de dialogues. Eliza, le premier programme réalisé par Weizenbaum (1966), destiné à communiquer avec un être humain, reproduit la rhétorique d'un psychothérapeute rogéien. Cette dernière s'appuie sur un répertoire de schémas conversationnels et sur la reformulation des informations provenant de l'interlocuteur. La réutilisation des mots par ce procédé confère à Eliza une faculté d'adaptation, donnant ainsi l'illusion d'un savoir étendu et commun. Cependant, en dehors d'un contexte spécifique avec des échanges standardisés comme pour la réservation de billet de train, les programmes de ce type ne fournissent pas de réponses pertinentes. De plus, chaque situation nécessite des règles spécifiques en fonction de la finalité du dialogue ou de l'ontologie employée, rendant la généralisation ou la réutilisation de programmes quasiment impossible.

L'intelligence artificielle basée uniquement sur le dialogue subit trois principales critiques : les deux premières expriment des difficultés, la dernière exprime une impossibilité. La première critique porte sur la difficulté de construire une ontologie exhaustive afin de couvrir l'ensemble des situations quotidiennes ou de saisir les glissements sémantiques. La seconde affirme que l'étude de la grammaire d'une langue ne révèle qu'une structure de surface ; en effet, tout le monde communique sans pour autant maîtriser parfaitement la grammaire. Une structure du langage existe mais elle n'est pas directement détectable dans la syntaxe. Ce second point a particulièrement été défendu par Chomsky (1965) qui a proposé une grammaire générative qui serait liée à des processus cognitifs innés. La dernière critique pose le problème de l'acquisition du langage et de sa nature ; celui-ci est illustré par l'exemple de la chambre chinoise de Searle (1980) : enfermé, un homme ignorant la langue chinoise doit communiquer par courrier en mandarin. Pour cela, il utilise un manuel d'instructions permettant de fournir une réponse en fonction des symboles reçus. Sans aucune mise en correspondance avec le monde, l'homme n'apprendra jamais le mandarin. Un programme n'effectuant rien de plus que l'homme se trouvant dans la chambre chinoise, il est légitime de penser qu'il ne comprend pas les symboles manipulés. Autrement dit, la sémantique d'un contenu mental ne se réduit pas seulement à la manipulation symbolique. Ici, la critique va au-delà d'un problème de mise en correspondance, il manque une visée intentionnelle. Une visée intentionnelle implique qu'une assertion génère une attente de la même façon que n'importe quel acte intentionnel. Ce sont les conséquences d'une assertion sur autrui qui deviennent source de sens. Par exemple, Pierre et Paul se trouvent dans une salle de travail avec une fenêtre ouverte. Pierre dit à Paul qui se trouve près de la fenêtre : « j'ai un peu froid ». Cet énoncé est vrai en cet

instant et à cet endroit, Pierre a une sensation de froid et il s'attend à ce que Paul ferme la fenêtre. Pierre aurait également pu mentir pour camoufler les véritables raisons pour lesquelles il souhaitait fermer la fenêtre. Il convient alors de distinguer la vérité d'une assertion et sa justification par ses croyances et ses attentes. L'argument de Searle (1980) rend caduc à la fois le test de Turing (1950) et la faisabilité d'une machine intelligente sans motivation et sans perception. Toutefois, des recherches sur les systèmes multi-agents, autrement dit plusieurs programmes communiquant entre eux, essaient de formaliser la notion d'acte de langage comme une couche de communication supplémentaire. De manière plus formelle, les systèmes multi-agents se définissent (Ferber, 2005) comme un ensemble B d'entités plongées dans un environnement E caractérisé par l'ensemble des états de l'environnement et un ensemble A d'agents avec $A \subseteq B$ auquel est associé un système d'opérations permettant à des agents d'agir dans E . De plus, un système de communication entre agents définit leurs interactions sur la base d'une organisation structurant l'ensemble des agents et définissant les fonctions remplies par les agents (notion de rôle et éventuellement de groupes). Éventuellement, une relation à des utilisateurs U agissant dans ce système multi-agent est spécifiée via des agents interfaces $U \subseteq A$

B - L'intelligence artificielle et la logique

La seconde voie pour qui le langage naturel est imparfait réduit également la pensée à la manipulation d'idées, reflets des concepts appartenant à un monde idéal. Cependant, les rouages psychologiques de la pensée humaine n'atteignent qu'approximativement la véritable intelligence de ces manipulations. La compréhension ne passe pas par l'étude de ces rouages mais par l'étude de ce qu'ils visent, c'est-à-dire les concepts et leurs règles logiques. La principale influence de l'intelligence artificielle provient de l'atomisme logique (Russel, 1918). Ce courant stipule qu'il existe des propositions élémentaires qui renvoient à des faits simples et irréductibles de la réalité extralinguistique comme « la pomme est sur la table ». L'agrégation de ces atomes logiques par des mots de connections tels « ou » ou « et » forme des propositions complexes comme « la pomme est sur la table et le ciel est bleu ». La signification des propositions composées dépend de la valeur d'un jugement synthétique sur les atomes qui la composent. Frege (1892) proposa un système formel pour rendre compte de cette réalité conceptuelle, la logique propositionnelle. Un système formel est constitué d'une convention d'écriture de formules et d'un ensemble de règles d'inférence transformant une formule en une autre. Un axiome représente une formule donnée a priori. Un théorème correspond à une formule qui se déduit mécaniquement par règles d'inférence à partir des axiomes. Une démonstration se résume par liste exhaustive des étapes ayant permis l'obtention d'un théorème. Par exemple, le système « MUI » (Hofstadter, 1979) possède un alphabet de trois lettres M, U et I. Le mot MI représente le seul axiome. À cela s'ajoute quatre règles d'inférence qui permettent de générer des mots :

1. Si un mot se termine par I, un U peut y être concaténé : MI donne MIU.
2. Un mot de type MX, où X représente un mot quelconque, peut se transformer en MXX : MI donne MII.
3. Si un mot contient III, cette chaîne peut être remplacée par U : MIII donne MIU ou MUI.
4. Si un mot contient UU, cette chaîne peut être supprimée : MIUU donne MI.

Dans ce système « MIU » se révèle être un théorème puisque l'application de la règle 1 sur l'axiome MI donne MIU à partir duquel le mot MIII se déduit par deux applications

successives de la règle 2. Ce dernier mot se transforme en MIU par la règle 3. L'arbre de la Figure I-10 représente le début de la production de l'ensemble des théorèmes de ce système.

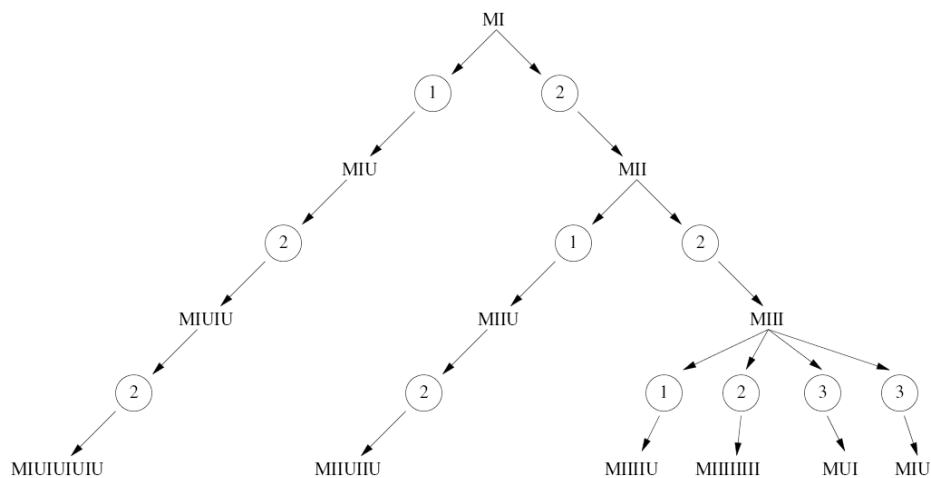


Figure I-10 : Arbre de production d'un théorème, les numéros indiquent la règle de transformation utilisée.

La notion de système formel permet de retrouver la distinction entre vérité analytique et vérité synthétique. Si une formule se révèle démontrable dans un système, alors cette formule est admise comme vraie par rapport à celui-ci, ce qui correspond à une vérité analytique. Une proposition renvoyant à un fait du monde est alors considérée comme un axiome traduisant la vérité synthétique. Une étude plus précise des propriétés des systèmes formels sera entreprise dans le chapitre suivant, néanmoins trois critiques doivent être évoquées afin de comprendre l'influence de ce courant en intelligence artificielle.

La première critique concerne l'établissement de la vérité synthétique. En effet, la mise en correspondance entre un concept et un objet du monde suppose qu'il existe soit une identification, soit des conditions de vérité, autrement dit des techniques validant la mise en correspondance. La première condition ne semble pas raisonnable car l'identification impose que la comparaison se fasse entre deux entités de même nature, ce qui n'est pas le cas au vu des hypothèses HOP et HOI. La seconde rend le problème récursif, en effet, les techniques de validation doivent être elles-mêmes validées et ainsi de suite à l'infini. Contrairement à ce qui a été évoqué précédemment, cette dernière remarque ne permet pas d'accepter les faits perçus comme des axiomes d'un système formel objectif. La vérité synthétique devient subjective ainsi que faillible et ne prétend plus à l'absolu. La vérité synthétique dépend de la méthode employée pour la mise en correspondance de la même manière que la vérité analytique dépend des axiomes et des règles d'inférence choisies. L'impossibilité d'une mise en correspondance absolue n'élimine cependant pas l'idée d'une vérité synthétique absolue : un partisan peut argumenter que le monde idéal et le monde physique se reflètent mutuellement l'un dans l'autre sans qu'il soit possible d'établir une correspondance exacte à cause de notre subjectivité. Autrement dit, la vérité synthétique est inaccessible et n'a de sens que d'un point de vue qui se trouve en dehors de l'existant.

En intelligence artificielle, les partisans du symbolisme pour dépasser ce problème vont transposer le modèle cognitif dans un monde virtuel, autorisant ainsi l'identification entre concepts et percepts. La conception d'un monde virtuel se fonde soit sur

l'imagination des chercheurs comme le monde des blocs de Winograd (1972), soit sur une représentation de la réalité. Cette dernière provient d'une construction géométrique à partir d'informations issues de différents capteurs. Cette construction nécessite un ensemble hiérarchisé de primitives ad hoc, c'est-à-dire un modèle du monde *a priori*, et des connaissances sur le dispositif des capteurs. Le concepteur s'autorise à prendre un point de vue situé en dehors de l'existant. La justification de cette position s'explique par le fait que la représentation géométrique de l'espace ou des objets physiques des concepteurs est considérée ici comme un objectif, indépendamment de leur perception. Dans ce contexte, les hypothèses HP1 et HP2 sont acceptées, HP3 en revanche sera refusée puisque la reconstruction dépend des objets recherchés. La perception descendante extrait les objets à percevoir. Hormis les difficultés techniques pour effectuer une reconstruction en trois dimensions d'une scène, la reconnaissance d'objets pose le problème de la catégorisation. Par exemple, le concept chaise doit être défini exhaustivement par un ensemble de propriétés permettant de désigner de manière intensionnelle l'ensemble des chaises. Mais l'établissement d'un ensemble de ces propriétés est impossible car la définition d'une chaise se rapporte également à sa fonction. L'autre solution, la constitution d'un ensemble en extension du concept de chaise, est également impossible dans un environnement ouvert c'est-à-dire avec des objets aux formes inconnues à l'avance. Le subsymbolisme proposera des solutions afin de contourner ces impasses.

La seconde critique porte sur l'indépendance des propositions élémentaires relativement à leurs valeurs de vérité. La valeur de vérité d'une proposition conditionne la valeur de vérité d'une autre ; juger vraie la proposition : « La tour Eiffel se trouve à Paris » conduit à juger fausse : « La tour Eiffel ne se trouve pas à Paris », et de manière intensionnelle toutes les propositions négatives relatives à la localisation de la tour Eiffel en dehors de Paris. La dépendance entre propositions ne limite pas seulement à la négation. Deux cas de figures traduisent l'impossibilité d'évaluer une proposition complexe par la somme des évaluations individuelles des propositions qui la composent. Le premier cas de figure intervient dans les situations où la proposition fait référence à la croyance d'un sujet : « Paul croit que la tour Eiffel se trouve à Paris ». La valeur de vérité de cette dernière proposition ne dépend pas de la vérité de la proposition « La tour Eiffel se trouve à Paris » mais de la croyance de Paul ; ici, la décomposition en propositions élémentaires n'apporte rien. Les logiques épistémiques modélisent ce type de situation que la logique propositionnelle ne prend pas en compte. Le second cas de figure englobe les logiques non monotones. Il s'agit de toutes les propositions dont la valeur de vérité est statuée a priori mais pouvant être révisée selon l'acceptation d'autres propositions dont une définition en intension ou en extension se révèle impossible.

Par exemple, « Si une pièce d'un euro est introduite dans la machine à café alors un gobelet de café sera servi » est vrai, mais si le réservoir de café le permet, si un extraterrestre ne la sabote pas, etc. Afin de caractériser ce nombre de possibilités ou l'apparition d'exceptions, Kripke (1959) a conçu la logique modale de laquelle le chapitre suivant traitera plus en profondeur. Cette dépendance entre propositions détruit l'ambition de l'atomisme logique. En intelligence artificielle, elle contraint à ne plus concevoir la connaissance comme un stockage de propositions successives. Chaque nouvelle proposition doit être reliée à d'autres, la construction de base de connaissances appelée aussi « système expert » relève d'une architecture logique. La maintenance de tels systèmes passe par l'intégration de nouvelles informations en fonction du reste de la base et de son objectif. Ces bases de connaissances permettent de prendre en compte l'évolution d'une situation seulement si elles possèdent les éléments conceptuels qui la prévoient. Ces

systèmes logiques ne simulent pas le « bon sens » qui permet de comprendre les implications entre propositions même si elles sont inconnues au préalable. En définitive, concevoir un système organisant les données de façon à ce que toutes les possibilités soient toujours visibles revient à accepter HC2. Dans ce cadre, HC3 est également acceptée ou du moins le traitement des émotions se réduit à prendre en compte des variables supplémentaires.

La dernière critique remet en question HC1. Cette hypothèse affirme que toute connaissance issue d'un système peut être construite par déduction. La résolution d'un problème se réduit à trouver un théorème satisfaisant les contraintes initiales, soit la bonne combinaison des règles d'inférences séparant les formules initiales de la formule finale souhaitée. Une méthode radicale s'appuie sur l'examen de toutes les combinaisons possibles. Dans la plupart des cas, l'explosion combinatoire interdit l'emploi de cette méthode, et la recherche en résolution de problème se concentre sur les différentes heuristiques possibles. Le programme de Newell et Simon (1960) : General Problem Solver, propose de décomposer le problème en sous-problèmes jusqu'à le rendre trivial. Ce programme complexe ne permet pas de résoudre tous les problèmes, seulement une classe de problèmes restreints par les connaissances injectées dans le système. Mais cette limitation peut être comprise également comme une limitation théorique de certains systèmes formels.

En 1931, Gödel démontra deux théorèmes. Le premier stipule que pour tout système formel suffisamment complexe (par exemple un système formel contenant l'arithmétique), il se peut qu'un théorème soit démontré bien qu'il existe des contre-exemples, autrement dit le système formel pour découvrir mécaniquement des vérités devient trop puissant dans le sens où elle peut démontrer certaines formules malgré l'existence de contre-exemples, cela s'appelle l'inconsistance. Le second théorème de Gödel stipule qu'il existe des formules dont l'application soit toujours correcte bien qu'aucune démonstration ne puisse les prouver, cela s'appelle l'incomplétude. Par ailleurs, Turing (1936) démontra que l'indécidabilité d'une formule ne peut être prouvée dans le système qui doit l'évaluer.

Par exemple, dans le système MUI présenté au début de cette section, il est impossible de prouver que la formule MU se révèle valide ou invalide à partir des règles de transformations et des axiomes du système. Cependant en utilisant un métalangage qui permet d'analyser ce système formel comme l'arithmétique, il devient alors possible de conclure que MU n'appartient pas au système MUI. Pour générer MU il faut que le mot précédent soit de la forme MXU où X représente une chaîne contenant un nombre de I multiple de 3 afin de pouvoir utiliser la règle 3, or les règles 1 et 4 n'interviennent pas dans la production de I et la règle 2 génère un multiple de trois seulement si la chaîne à dupliquer l'est déjà. Puisque l'axiome ne contient pas de multiple de trois, alors le mot MU est indémontrable au sein du système MUI. La construction de cette solution a nécessité un métalangage, autrement dit, un moyen de pouvoir sortir du cadre imposé par le système. Toutefois, ce métalangage n'échappe pas lui-même aux limitations révélées par les théorèmes de Gödel et de Turing. Ce problème réflexif constitue souvent des paradoxes, comme le paradoxe du menteur « je mens ».

Cette limite empêchant de pouvoir rendre compte du monde de manière cohérente au sein d'un unique système pose un problème majeur pour le symbolisme. Un système formel se restreint alors à une problématisation d'un ensemble de situations dont certaines éventuellement nécessiteraient pour être traitées un nouveau système plus complet correspondant à un métalangage. Il est illusoire néanmoins d'imaginer un métalangage

incluant tous les sous-systèmes formels de manière cohérente. Le problème pour le symbolisme porte alors d'une part sur le choix de la problématique du système formel initial et d'autre part, sur les mécanismes permettant de rajouter de nouveaux principes à ce système qui n'étaient pas directement déductibles à partir des données du problème. Cette capacité à créer des relations ou des objets supplémentaires sans utiliser la mécanique symbolique du système de référence correspond à la notion d'abduction de Peirce (1934) évoquée précédemment. Celle-ci aborde à la fois le problème de la créativité et de la gestion de l'incertain qui lui est liée. Selon le symbolisme, la formalisation de l'abduction suit essentiellement deux directions.

Les travaux de Hempel (1966) représentent la première ; il caractérise l'abduction comme l'utilisation de manière inverse d'une règle déductive. Étant donné un fait B et sachant que A implique B, l'abduction permet de supposer l'existence du fait A. Cependant, cette interprétation de l'abduction se révèle être également une théorie déductive, par conséquent, avec des capacités de créativité limitée puisqu'elle se limite à rechercher la meilleure cause possible dans une liste exhaustive. La seconde direction peut-être illustrée par les travaux de Gärdenfors (1988) où l'abduction devient un procédé pour la révision de la base de connaissances : une théorie T et un fait B qui ne peut être déduit à partir de T oblige à l'introduction d'une hypothèse A dans T qui permette la déduction de B. Le choix de A pour la révision de T doit toutefois respecter des critères de parcimonie et de cohérence. Ces critères renvoient à l'idée d'une dialectique entre la théorie et le monde qui tende asymptotiquement vers la théorie fidèle au monde. La seconde direction concernant la formalisation de l'abduction privilégie la vérité par cohérence à défaut de pouvoir avoir accès à une vérité par correspondance. La position de la vérité par cohérence autorise que ce second formalisme de l'abduction soit également soutenu par les partisans de la famille FCP. De par la précision et la facilité d'utilisation, les concepteurs en intelligence artificielle privilégient le premier formalisme de l'abduction. Étant donné que celui-ci se réduit à une forme de théorie déductive, cela revient à accepter HC1.

**

En robotique, l'approche symbolique oblige le concepteur à définir les problèmes, les concepts de bases pour les résoudre et la représentation de l'environnement. Or, ces contraintes s'opposent à la définition de l'autonomie. Afin de pallier cette lacune, l'approche subsymbolique souhaite conserver la puissance de la formalisation tout en reconsidérant le rôle de la perception.

3.1.1.2. Le subsymbolisme

Cette section se propose (A) de dégager les caractéristiques du courant subsymbolique en comparant les interprétations sur une expérience de pensée avec le symbolisme. Les principes fondamentaux du subsymbolisme précisés, (B) leur adéquation avec les données neurophysiologiques sera ensuite abordée. Enfin, (C) la position du subsymbolisme vis-à-vis des hypothèses cognitives sera étudiée.

A - La différence entre le symbolisme et le subsymbolisme

Contrairement au symbolisme qui se concentre seulement sur la manipulation des idées, le subsymbolisme réintègre la description physique dans l'étude de la cognition. Plus précisément, le subsymbolisme affirme que les descriptions neurophysiologiques et psychologiques renvoient à une même réalité, c'est-à-dire qu'il existe une identification

entre une configuration neuronale et un état mental. Ce changement de perspective permet de modifier le problème de la vérité synthétique dans l'approche symbolique. L'exemple de la « terre jumelle » de Putnam (1975) illustre ce changement : Une planète P ressemble en tout point à la terre excepté l'eau qui se trouve remplacée par un liquide L ayant la même apparence, bien que sa composition chimique soit différente. Le même vocable est utilisé pour les deux liquides. Un terrien ou un habitant de P aura alors un état mental identique (donc selon ce courant un état cérébral identique également) en désignant l'un des deux liquides par le vocable eau, puisque l'un comme l'autre est incapable de les distinguer.

Cependant, le symbolisme considère que l'intentionnalité est différente, puisqu'en effet, une pensée vise toujours quelque chose et d'une certaine manière (Husserl, 1929). Lorsque le terrien énonce « cette eau est bonne » alors qu'il est sur P, il souhaite exprimer l'idée de l'eau qui dénote l'eau H₂O, le terrien se trompe. Autrement dit, les conditions de vérité de la proposition « cette eau est bonne » se révèlent être différentes selon le locuteur et selon le lieu où elle est prononcée. Malgré la similitude neurophysiologique supposée, les pensées de nature transcendante se révèlent alors différentes. Néanmoins, une autre expérience de pensée similaire à celle de Putnam (1975) peut illustrer une seconde interprétation. La terre et la planète P sont identiques en tout point, c'est-à-dire que Paul-T habitant la terre vit de la même manière que Paul-P sur la planète P. Un démon s'amuse à échanger depuis leur naissance Paul-T et Paul-P. Pour chacun, sans le savoir, l'apprentissage du concept de l'eau s'est forgé à la fois sur H₂O et XYZ. Si tous les deux prononcent « cette eau est bonne », les conditions de vérité deviennent identiques bien que les pensées puissent renvoyer à deux objets différents. La position d'un observateur omniscient ne permet pas de distinguer les pensées, donc elles peuvent être réductibles aux états cérébraux qui les sous-tendent. En définitive, la première interprétation se trouve biaisée par l'introduction d'un observateur qui possède des facultés que les habitants ne possèdent pas et qui projette des propriétés supplémentaires au concept de chacun des habitants. Un partisan du symbolisme peut arguer que les concepts possèdent une hiérarchie et de ce fait les deux hommes ont en commun un hyperonyme du concept eau H₂O ou XYZ. Cet argument permet de conserver l'idée du symbolisme qu'une pensée est transcendante tout en acceptant qu'elle puisse être équivalente à un état cérébral. Mais le subsymbolisme va plus loin en affirmant que cette équivalence révèle le caractère subjectif de la vérité d'une proposition et rejette l'idée d'une correspondance ontologique. Les concepts relatifs au monde perdent alors leur statut de concepts objectifs.

Les concepts relatifs au monde s'identifient avec les traitements des entrées sensorielles qui définissent de manière intensionnelle ce qui doit se placer sous un même symbole. Le sens, autrement dit ce qui relie l'extérieur et l'intérieur, n'est plus capturé par le symbole en lui-même mais par ce qui l'amène, le subsymbolique. La manipulation symbolique continue d'exister mais elle représente la manipulation du contenu étroit du concept en opposition avec le contenu large du concept qui se trouve être l'application du traitement incarnant le concept avec l'environnement. La définition intensionnelle d'un concept s'identifie à la notion de fonction qui relie des entrées à une sortie. La cognition devient un ensemble interconnecté de fonctions dont la sortie de l'une peut être aiguillée vers l'une des entrées d'une autre. Toutes les fonctions calculables étant réductibles à une machine de Turing, si les concepts sont compris comme des fonctions calculables, alors il existe un langage formel qui spécifie ces fonctions. Le but des sciences cognitives serait de retrouver ces spécifications, Krivine (2004) propose le lambda calcul comme le langage de référence.

B - Le subsymbolisme et les neurosciences

Le subsymbolisme conserve les hypothèses HOP et HOI. Le monde reste constitué d'objets élémentaires dont le comportement peut être étudié collectivement et tendre ainsi vers une description intersubjective sans prétendre à l'existence de concepts miroirs. Toutefois, la logique et les mathématiques qui restent pures peuvent décrire n'importe quel élément du monde. Ces éléments, puisqu'ils suivent ontologiquement des propriétés logiques, peuvent s'organiser de telle manière qu'apparaissent des propriétés structurelles réductibles à l'architecture d'une machine de Turing. Schématiquement, le subsymbolique se fonde alors sur deux conjectures : (i) l'identification des états mentaux avec des états cérébraux et (ii) la réduction du cerveau à une machine de Turing. Réciproquement, ces deux conjectures qui seront présentées trouvent des justifications dans deux types de données issues de l'étude du cerveau.

i - Les états mentaux et les états cérébraux

La première conjecture s'appuie sur la neuropsychologie. La corrélation entre des troubles cognitifs et des lésions cérébrales montre que certaines parties du cerveau participent essentiellement à des tâches mentales spécifiques telles que le langage (Broca, 1861). À l'inverse, des stimulations électriques précises peuvent déclencher des sensations, des souvenirs ou des schèmes comportementaux (Penfield, 1937). Par ailleurs, l'imagerie cérébrale fonctionnelle qui regroupe un ensemble de techniques permettant de suivre l'activité cérébrale selon différents points de vue conforte l'idée que le cerveau est organisé en régions spécialisées. Elle met également en évidence que ces régions fortement interconnectées travaillent ensemble. La faculté du langage, par exemple, ne se trouve pas localisée en un centre mais résulte de la coopération de plusieurs centres. La fonction cognitive du langage émerge d'une activité collective. Ce concept d'émergence traduit le fait qu'une fonction résulte de l'interaction d'autres fonctions sans qu'aucune d'elles ne possède les propriétés de la nouvelle fonction. La description componentielle devient alors particulièrement bien adaptée pour capturer ce phénomène. Pour ce qui concerne les états mentaux, même s'ils ne peuvent pas être liés directement à l'activité d'une région cérébrale, le concept d'émergence permet de les mettre en correspondance avec des ensembles d'activité cérébrale. Les sciences interprétatives se concentreraient sur les méta-fonctions cognitives alors que les sciences explicatives se concentreraient sur les fonctions physiologiques qui les sous-tendent.

Dans cette vision fonctionnelle des processus mentaux, deux aspects contribuent à développer la notion d'autonomie ainsi qu'à la limiter. Le premier aspect concerne la notion d'activité cérébrale ; celle-ci ne doit pas être comprise comme un état statique précis, mais comme une dynamique entretenue par une assemblée de neurones. La notion dynamique renvoie à la notion de rétroaction. Les systèmes rétroactifs caractérisés par Wiener (1948) reçoivent en entrée des informations provenant des effets de leur sortie. Ce mécanisme permet au système de tendre vers un état d'équilibre sans forcément l'atteindre. Cet état d'équilibre peut se comprendre comme un attracteur ou comme une consigne qui, invisible directement, traduirait la téléonomie du système. Selon les subsymbolistes, les états mentaux viseraient les consignes, les modes du système cognitif global à un moment donné. Ainsi par ce niveau d'abstraction supplémentaire, les états mentaux restent pertinents. Ce premier aspect introduit la notion de régulation voire d'autogestion, condition nécessaire à l'autonomie mais insuffisante. En effet, le choix des consignes dépend soit de l'extérieur, soit d'une détermination innée. En d'autres termes, la rétroaction autorise tout au plus la conception de robots semi-autonomes. Le second aspect se focalise

sur la nature de ces fonctions et sur leurs rôles qui restent difficiles à définir. Le subsymbolisme substitue le rôle du programmeur au mécanisme de l'évolution, la phylogénèse. Toutefois, le nombre de l'ensemble des connexions synaptiques d'un cerveau s'élevant à un milliard de milliards, leurs organisations ne peuvent pas être codées précisément dans l'ADN. Les fonctions cognitives nécessitent ainsi une phase d'initialisation, l'ontogénèse. Cependant, cette conception interdit au système de construire ses propres fonctions en termes de finalité alors qu'il s'agit d'un principe fondamental de l'autonomie.

ii - Les neurones et la machine de Turing

La seconde conjecture du subsymbolisme consiste à apparenter le cerveau à l'ordinateur et à apparenter la plasticité des fonctions cognitives aux diverses applications d'un ordinateur. Cette position implique la distinction et l'indépendance entre matériel et logiciel. Le logiciel se traduit et opère par des modifications physiques mineures (réversible, conservant l'intégrité du système) permettant toutefois de changer les propriétés du système. Pour l'ordinateur, ces modifications mineures se traduisent par le changement d'état électrique de composants électroniques comme les transistors, alors que pour le système nerveux, composé de cellules neuronales connectées les unes aux autres (Cajal, 1893), le subsymbolique réduit le neurone à un composant possédant également plusieurs états électriques. Schématiquement, le neurone possède plusieurs entrées, les dendrites, et une unique sortie, l'axone. Les dendrites pondèrent les entrées en fonction des poids de connexion. La Figure I-11 représente un neurone formel, modélisation classique d'un neurone.

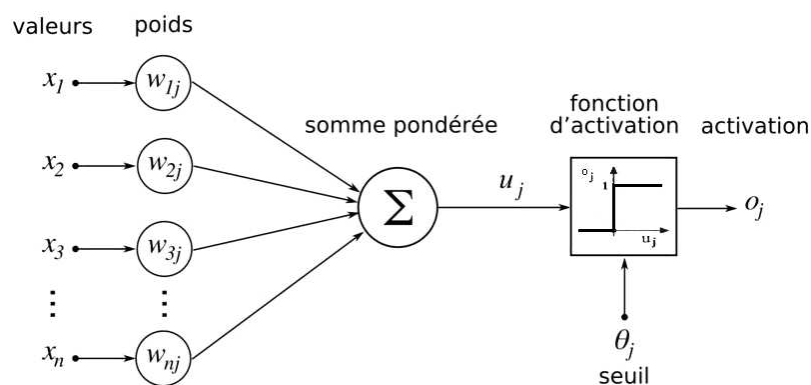


Figure I-11 : Modèle classique de neurone (McCulloch et Pitts, 1943).

Cette modélisation permet au subsymbolisme de considérer le neurone comme une unité de calcul élémentaire bien que cette modélisation ne prenne pas en compte un nombre important de facteurs comme le rôle des cellules gliales ou la complexité et la spécificité des échanges électrochimiques. Il est à noter que ces unités de calcul peuvent travailler en parallèle alors qu'un ordinateur classique centralise la puissance de calcul. Le neurone se conçoit comme un petit ordinateur fonctionnant avec un programme défini mais avec des valeurs de fonctionnement modifiables, les poids de connexion. La description de l'organisation matérielle devient également une description componentielle. Cette organisation se traduit dans la modélisation de réseaux de neurones par leurs différentes interconnexions, comme l'illustre la Figure I-12 :

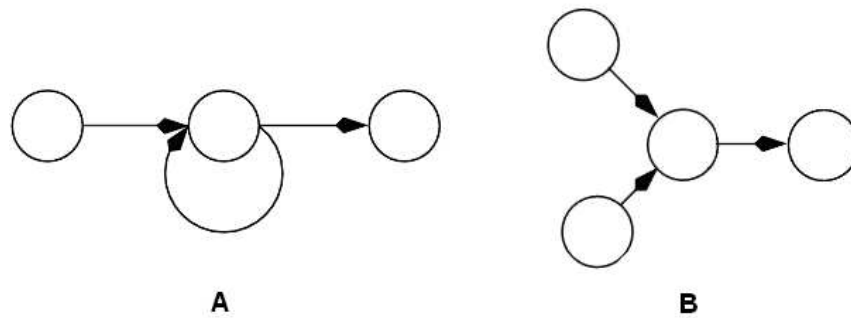


Figure I-12 : Deux exemples d'interconnexions, les schémas A et B représentent respectivement un réseau récurrent, c'est-à-dire comprenant une rétroaction et un réseau de type perceptron (Rosenblatt, 1958).

En acceptant un ensemble de machines de Turing se réduit à une machine de Turing, le subsymbolisme comprendra cette décomposition comme une preuve de la similitude entre cerveau et ordinateur. Par ailleurs, la programmation du cerveau peut se concevoir comme la combinaison de ces unités de calcul et l'apprentissage comme le moyen de modifier les poids de connexion. Dans ce cadre, trois types d'apprentissage existent : l'apprentissage supervisé, l'apprentissage par auto-organisation et l'apprentissage par renforcement. Tous trois peuvent s'intégrer dans un projet robotique du point de vue du subsymbolisme mais les développements des deux derniers se réalisent généralement dans un cadre différent. Tous deux seront abordés dans la section portant sur le connexionnisme et sur le fonctionnalisme écologique. L'apprentissage supervisé sera expliqué lors de l'interprétation des hypothèses d'accessibilité.

*

En résumé, le subsymbolique propose deux voies d'investigation pour la cognition qui se révèlent très ardues (Figure I-13) : le désassemblage et le fonctionnalisme. Le désassemblage correspond à l'opération inverse à celle de la traduction d'un programme écrit dans un langage formel de haut niveau en un langage machine qui décrit l'activité des composants. Dans ce cas, le désassemblage correspond à décrypter les comportements inscrits dans le cerveau comme un programme écrit en λ -calcul pour tenter d'obtenir une spécification des fonctions cognitives (Krivine, 2004). Cependant, le désassemblage pose aussi bien des problèmes d'ordre qualitatif (quelles informations choisir ? quelle méthodologie appliquer ?) que des problèmes d'ordre quantitatif (comment suivre simultanément des milliards de neurones ?). Le fonctionnalisme qui correspond à la décomposition fonctionnelle des processus mentaux se trouve compliqué par le phénomène d'émergence qui masque les fonctions sous-jacentes. Il s'agit pourtant de l'unique description valable puisque la description procédurale (DP) n'arrive pas à saisir les variables cachées en toutes circonstances. Plus concrètement, la première investigation n'a pas d'impact sur la réalisation robotique, contrairement à la seconde qui va orienter quasiment toutes les positions sur les hypothèses cognitives.

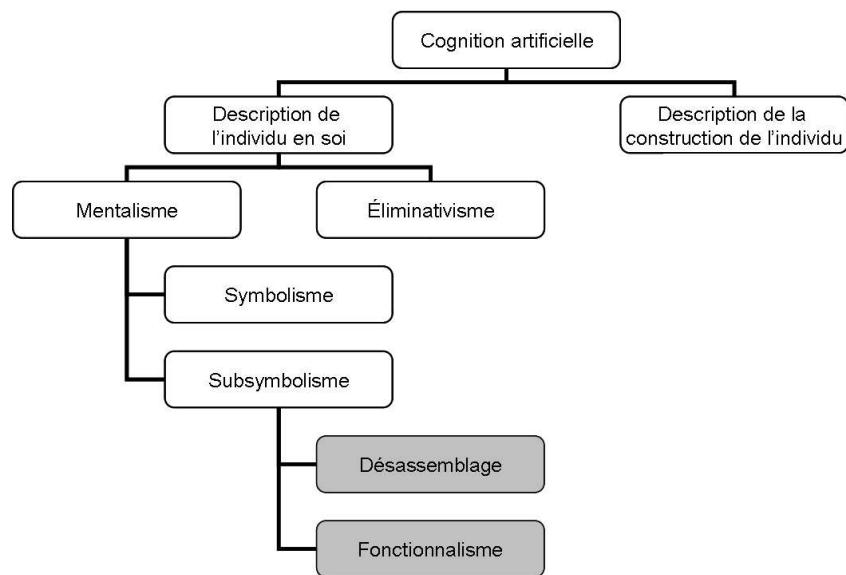


Figure I-13 : Les approches issues du désassemblage et celles issues du fonctionnalisme au sein de l'arborescence de la cognition artificielle.

C - Les conséquences sur les hypothèses cognitives

Concernant les hypothèses perceptives, le subsymbolique, comme le symbolique, adopte HP1. Dans la logique du subsymbolisme le capteur correspond à une segmentation du monde élémentaire et de ce fait constitue un concept en tant que fonction élémentaire. À chaque traitement, l'intégration devient de plus en plus abstraite, jusqu'au système de décision. Cette structure pyramidale ne permet pas d'exclure HP2 ; certes les traitements élémentaires, autrement dit les concepts de bases, ne participent jamais directement au mécanisme décisionnel, mais seuls les concepts de haut niveau, comme pour le symbolisme, participent au système décisionnel central. HP3 est nécessairement acceptée puisque les concepts se réduisent à des fonctions dont les entrées dépendent seulement de l'extérieur (Marr, 1980). À l'inverse du symbolisme, la perception est ascendante. La construction des concepts peut s'effectuer par un apprentissage supervisé. Il consiste à trouver mécaniquement une relation entre l'état des capteurs et le symbole associé, grâce à la rétropropagation du gradient d'erreur. La perception se réduit à la reconnaissance, mais le système ne crée pas lui-même les spécifications des fonctions qu'il doit produire.

Cette rigidité se révèle d'autant plus indispensable pour toute activité de communication entre deux entités robotiques. La subjectivité du concept provient de l'imperfection de la perception et non de l'algorithme qui l'amène. Toutefois, deux robots ramasseurs de pommes avec des capteurs différents (sonar et caméra par exemple) auront des algorithmes de reconnaissance de formes différents, même si la communication impose un symbole identique pour la pomme. Dans cette situation, l'accord sur le concept de pomme vient du fait que les fonctions de reconnaissance des deux systèmes se ressemblent d'un point de vue extérieur. Le subsymbolique, limitant l'apprentissage à l'initialisation de fonctions, oblige à discerner l'acquis de l'inné : la capacité de lier des images sensorielles à un symbole relève de l'inné et l'association entre des images d'une pomme et le symbole « pomme » relève de l'acquis. Mais seuls les concepteurs contrôlent ce qui doit être associé.

La position sur les hypothèses relatives à la conception reste inchangée dans le principe par rapport au symbolisme. Cependant, le subsymbolique ne recherche plus des raisonnements bâtis sur des systèmes formels figés puisque la correspondance n'est plus

assurée. La vérité par correspondance se transforme donc en une croyance dépendante de diverses justifications. La notion d'incertitude associée se représente soit par des probabilités soit par de la logique floue. L'évaluation des possibles permet de trier les inférences abductrices comprises comme des déductions inversées. Toutefois, le problème concernant la création d'inférence reste ouvert de la même manière que pour le symbolisme. Le maintien de l'hypothèse HC2 s'explique par l'organisation fonctionnelle imaginée par le subsymbolique qui se structure de manière très hiérarchisée. En effet, les inférences des étages inférieurs ne servant qu'à dégager les concepts pour les niveaux d'abstraction supérieurs, seul le dernier réseau d'inférence contenant tous les concepts de hauts niveaux décide ou raisonne. Tout comme le symbolisme, le subsymbolique n'intègre pas les émotions directement dans la cognition et conserve HC3.

**

Les fonctions cognitives, en ce qu'elles capturent et en ce qu'elles visent, ne pouvant être définies précisément à partir des études neuropsychologiques, les réalisations robotiques subsymboliques créent des organisations fonctionnelles orientées vers un ou plusieurs buts. Mais, contrairement au symbolisme, l'apprentissage supervisé permet d'adapter la reconnaissance des objets ; le système cherche la manière d'effectuer au mieux la mise en correspondance au lieu de suivre des règles prédéfinies. En résumé, la difficulté des approches subsymboliques se trouvent dans la tentative de réunir le bas et le haut niveau, soit, la physique et la psychologie. L'éliminativiste souhaite pallier cette difficulté en éliminant le second terme. Mais avant de l'aborder, un premier tableau (Tableau I-2) récapitulatif des diverses positions vis-à-vis des hypothèses et méthodes constituant la grille d'analyse peut être avancé concernant les approches issues du mentalisme.

		Types de famille métaphysique				Hypothèses sur la perception			Hypothèses sur la conception			Types de méthode d'investigation			Types de description		
		FC	FCP	FP	FA	HP1	HP2	HP3	HC1	HC2	HC3	MHD	MC	MA	DP	DM	DCO
Symbolisme	Logicisme	Red	Green	Red	Red	Green	Green	Red	Green	Green	Yellow	Yellow	Green	Yellow	Green	Red	Red
	Psychologisme	Red	Green	Red	Red	Green	Green	Green	Green	Green	Yellow	Green	Yellow	Green	Red	Red	Red
Subsymbolique	Désassemblage	Yellow	Green	Red	Red	Green	Green	Green	Green	Green	Green	Yellow	Green	Yellow	Green	Red	Red
	Fonctionnalisme	Red	Green	Red	Red	Green	Green	Green	Green	Green	Green	Yellow	Green	Yellow	Yellow	Red	Green

Tableau I-2 : Récapitulation de diverses positions relatives à la grille d'analyse des approches provenant du symbolisme et du subsymbolisme. Les cases vertes correspondent à un avis favorable, les cases jaunes à un avis mitigé et les cases rouges à un avis défavorable.

3.1.2. L'éliminativisme

L'éliminativisme prône une seule description, celle issue des sciences explicatives. Les croyances et les intentions sont considérées comme des descriptions naïves qui ont survécu de par leur aspect pratique et normatif. Ainsi, ces concepts renvoient à des propriétés et à des entités mentales qui n'ont aucune réalité (Churchland, 1981). Les descriptions psychologiques restent obligatoirement singulières comparativement aux lois de la physique qui visent l'universel. Seules les entités ou les lois physiques possèdent un pouvoir causal, autrement dit la conscience de soi devient un épiphénomène, un spectateur impuissant. Cette position met entre parenthèses les hypothèses portant sur la perception et la

conception. Dans le cadre de l'intelligence artificielle, l'éliminativisme peut se réduire à deux approches (Figure I-14) : le neuromimétisme et le connexionnisme.

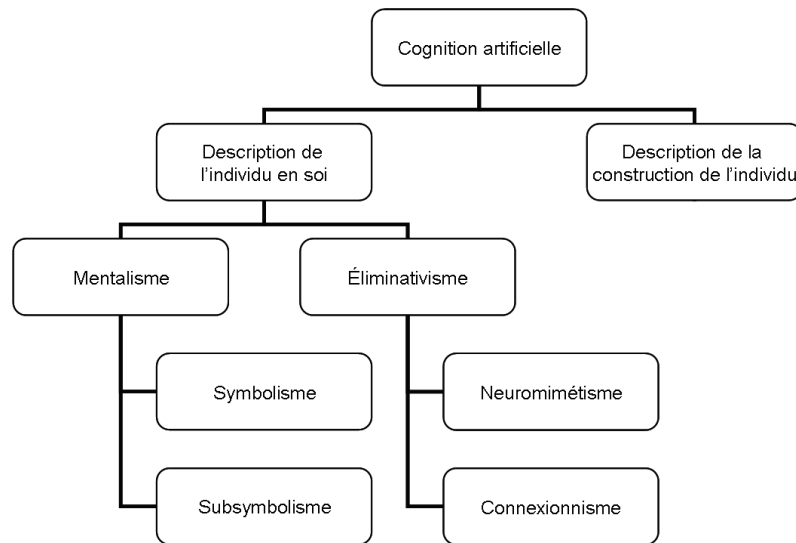


Figure I-14 : Le courant du neuromimétisme et celui du connexionnisme au sein de l'arborescence de la cognition artificielle.

De par la place qu'occupent les sciences physiques dans l'éliminativisme, les diverses interprétations sur le rôle et la nature des sciences physiques seront détaillées au sein de chacune de ces deux approches. Ces diverses interprétations physicalistes permettront de mieux comprendre comment l'éliminativisme dissocie l'épistémologie et la cognition, distinction qui sera remise en cause dans la partie portant sur la description de la construction de l'individu. Cette section s'attachera surtout à montrer comment cette distinction et ses interprétations conditionnent le développement de la robotique cognitive dans les approches éliminativistes. En effet, il existe un véritable continuum entre ces positions mais, pour des raisons de lisibilité, seules les principales positions seront abordées.

3.1.2.1. Le neuromimétisme

Le neuromimétisme, dont la démarche consiste à mimer l'activité cérébrale, considère que seule l'étude du substrat permet d'expliquer le comportement des individus cognitifs (Churchland, 1986). Cette position peut être défendue par des partisans de la famille métaphysique FCP ou FP. Cette section abordera l'interprétation des sciences physiques selon ces deux familles, c'est-à-dire selon deux positions physicalistes : (A) le matérialisme radical et (B) le matérialisme rationnel. Ces deux positions conditionneront la valeur de chacune des trois méthodes d'investigations et celle de chacune des trois types de description. (C) Cette analyse permettra de mieux comprendre le problème de la cognition perçu par le matérialisme et éclairera les raisons de l'utilisation de systèmes robotiques. (D) Les différences avec le subsymbolisme, qui appartient également à la famille FCP, seront soulignées en particulier sur la notion d'apprentissage. Enfin, il sera démontré que l'approche neuromimétique, malgré sa puissance de description, possède des limites intrinsèques liées à ce que signifie la compréhension d'un phénomène.

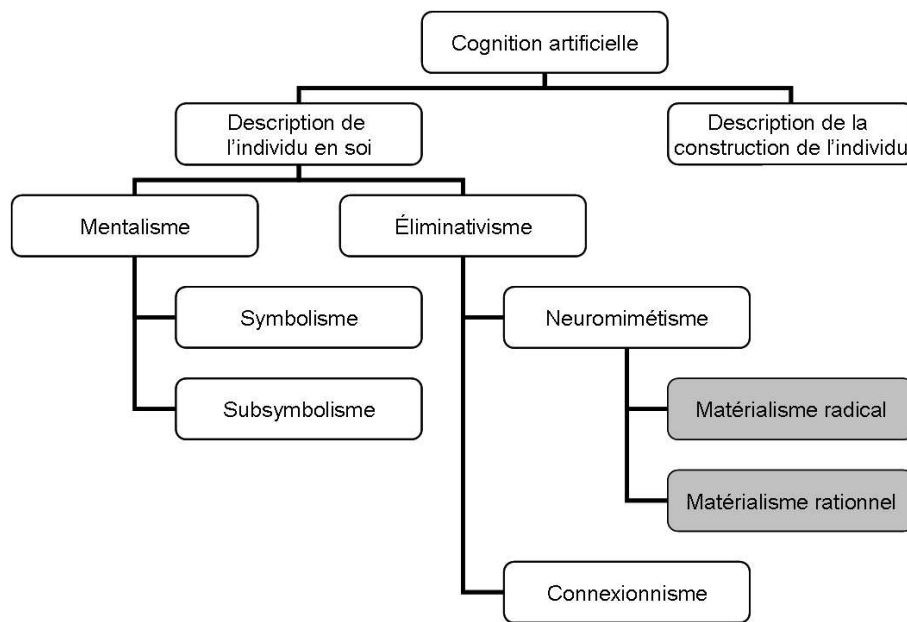


Figure I-15 : Les approches issues du matérialisme radical et celles issues du matérialisme rationnel au sein de l'arborescence de la cognition artificielle.

A - Le matérialisme radical

La première position physicaliste, le matérialisme radical, repose seulement sur l'HOP, il appartient ainsi à la famille FP. Ce mouvement considère que l'univers se compose de particules élémentaires, irréductibles, dont l'agrégation permet d'obtenir des formes et des comportements de plus en plus complexes. Seules les particules élémentaires peuvent prétendre être des objets objectifs. En effet, les objets constitués d'une agrégation de particules se révèlent être un regroupement perceptif arbitraire ou pratique. Le concept de pomme se comprend comme une invention alors que le concept caractérisant une particule élémentaire résulte d'une découverte. Dans cette perspective, l'homme construit tous les concepts, même ceux issus de la logique et des mathématiques. Ces derniers permettent de représenter les propriétés des objets physiques sans être des concepts purs possédant une objectivité transcendante.

Cette dissymétrie ontologique entre le physique et l'idéal se répercute sur la valeur des méthodes d'investigation et d'argumentation. La méthode par la cohérence (MC) qui repose sur la juste imbrication conceptuelle ne suffit plus à convaincre puisque ces concepts ne sont que des outils conceptuels pratiques destinés à rendre compte de certaines formes du réel et non le réel. La manipulation conceptuelle sans connexion avec le sensible perd de son sens. La seule méthode valable réside alors dans l'approche hypothético-déductive (MHD). Toutes constructions conceptuelles ou hypothèses se confrontent à l'expérience qui sanctionne et oriente l'investigation. Plus un énoncé prédit l'évolution d'un phénomène, plus il s'approche de la vérité ; cependant, il n'y a pas de correspondance ontologique entre un énoncé et le phénomène visé. La dernière méthode, la méthode artefactuelle (MA) permet de pallier la difficulté de réaliser certaines expériences. Ainsi, elle participe à l'argumentation d'un modèle physique mais ne participe pas à son investigation ou sa conception proprement dite, puisque celle-ci s'appuie uniquement sur l'observation de l'objet d'étude initial.

Concernant l'appréciation des trois types de description, elle diffère selon que l'objet visé correspond à une particule élémentaire ou à un agrégat de celles-ci. Obligatoirement, pour le premier cas, seule la description procédurale ou la description mathématisée peut prétendre à décrire les particules élémentaires. En cela, ces deux types de description forment la base de la connaissance, toutefois, pour des raisons d'efficacité, la physique contemporaine utilise exclusivement la description mathématisée. Mais la valeur de cette description se trouve dans la certitude d'avoir localisé, à la dichotomie, l'une des particules élémentaires. Sans tenir compte de la position prônant une dichotomie à l'infini, deux critères principaux revendiquent la légitimité d'être le signe que la dichotomie a atteint sa limite. Chacun de ces deux critères conditionne une définition du déterminisme.

Le premier critère est positif : si la précision des prédictions scientifiques peut s'effectuer de manière absolue, c'est-à-dire sans incertitude alors cela signifiera que les théories arrivent à capturer le comportement des particules élémentaires. Cette position introduit une contrainte supplémentaire qui n'est pas incluse dans l'HOP : le mécanisme. Le mécanisme stipule que l'évolution d'une particule élémentaire dépend uniquement des interactions passées ou présentes avec les autres particules élémentaires. L'introduction de cette hypothèse peut être interprétée comme une acceptation de l'HOI dans le sens où l'idée du mécanisme participe à l'interprétation des résultats scientifiques sans être soumise à l'expérience. Dans tous les cas, ce critère aboutit à un déterminisme historique : si les conditions initiales de l'univers étaient parfaitement connues, alors le monde serait à tout moment (passé, présent et futur) prévisible. Cette idée a été notamment professée par Laplace (1814).

Le second critère est négatif : si les appareils de mesure ne permettent plus de déceler les raisons de l'imprécision des résultats et que la théorie tout en intégrant ce facteur conserve un pouvoir prédictif, alors la dichotomie est achevée. Par exemple, la théorie de la mécanique quantique propose que les incertitudes sont une propriété des objets étudiés (Gell-Mann, 1995). Autrement dit, il y aurait à la fois une causalité extrinsèque provenant de l'interaction avec les autres particules et une causalité intrinsèque dépendant de l'essence même de la particule. Une particule élémentaire A possède une position dans l'espace suivant une distribution de probabilité, la notion de particule devient une commodité de langage, elle ne correspond plus à la métaphore intuitive de la sphère dure qui autorisait un modèle planétaire de l'atome (Rutherford, 1911). Ainsi, chaque particule aurait un champ des possibles en interaction avec les autres particules, ce champ des possibles se réduirait pour donner le présent, tout en s'ouvrant pour potentialiser le futur. Le déterminisme ne peut plus être historique, absolument prédictible, puisque la mesure devient ontologiquement incertaine au niveau subatomique. Cependant, le monde reste déterminé par les particules élémentaires qui le composent ; alors, en prenant en compte leurs propriétés, le déterminisme devient un déterminisme probabiliste. Dans le champ des possibles, l'univers est imprévisible, néanmoins l'univers ne peut être autrement.

En somme, concernant l'étude des objets élémentaires, quel que soit le critère défendu, la DP et la DM restent équivalentes alors que la DC devient inappropriée. En revanche, lorsque la description porte sur un objet composé, la DC acquiert une valeur particulière au détriment de la DP et de la DM. Ces deux dernières tentent de décrire un phénomène compris comme un tout alors qu'il est composé de plusieurs phénomènes interagissant avec d'autres. D'une part, la description se révèle alors ponctuée d'approximations. Or, ces approximations de mesures ou de facteurs peuvent engendrer de grandes erreurs de prédiction. En effet, Poincaré (1908) démontre que certains modèles dynamiques présentent une grande sensibilité aux conditions initiales, ce sont les

phénomènes chaotiques. Dans la perspective du second critère, la distinction entre l'incertitude des systèmes chaotiques et l'incertitude des systèmes quantiques se situe sur leur origine, la première provient de l'observateur (incapacité à relever précisément toutes les conditions initiales), alors que la seconde provient de l'objet étudié. D'autre part, la différence entre les types de description vient aussi du fait qu'elles ne prétendent pas à la même valeur du point de vue de la causalité ; les modèles de la mécanique classique peuvent être compris comme des raccourcis dans la chaîne causale réelle qui serait la prise en compte de toutes les particules élémentaires avec leurs propriétés quantiques. Mais à ce jour, la physique théorique ne semble pas être capable d'expliquer les relations entre les propriétés relevées à un niveau macroscopique et à un niveau subatomique. Ainsi, en définitive, pour un même objet, une DC sera toujours plus proche de la réalité (de la vérité) qu'une DP ou une DM. Le matérialisme radical aboutit à la méthodologie cartésienne, si le modèle possède trop d'imprécisions, alors, cela signifie que l'objet considéré en tant que tel doit être décomposé afin d'obtenir un grain plus fin de l'image de la réalité.

B - Le matérialisme rationnel

La seconde position physicaliste, le matérialisme rationnel, qui appartient à la famille FCP, considère également que les objets élémentaires sont les seuls à être objectifs, mais accepte en plus que les concepts logicomathématiques soient des concepts objectifs et purs. La logique définit les principes fondamentaux de toute existence et les mathématiques permettent de définir les lois particulières qui régissent les choses. Chaque découverte mathématique offre de nouveaux termes pour décrire les lois de la nature. Par rapport au matérialisme radical, l'adhésion à l'PHOI a des répercussions sur l'appréciation de la MC et de la DM. La MHD conserve le statut de méthode principale et la MA reste un palliatif, un exercice d'imitation. En revanche, puisque les concepts logiques traduisent ontologiquement la nécessité de ce qui est, la MC reprend le statut de méthode scientifique sous couvert d'une étroite relation avec les données recueillies. Cette position amène à compléter la vérité par l'observation des prédictions avec la vérité par la cohérence.

Commune à toute position matérialiste, l'analyse du critère d'arrêt de la dichotomie s'applique également ici. Toutefois, l'appréciation des trois types de description diffère quelque peu. En effet, la place privilégiée des mathématiques restreint, tout au plus, la DP aux travaux préliminaires pour cerner le phénomène à étudier. Mais la modification majeure est que les descriptions de type mathématique (DM) et componentiel (DC) peuvent pareillement prétendre à saisir la réalité. Les particules élémentaires suivent des lois qui ne semblent pas compatibles avec les lois macroscopiques. Le matérialiste radical résout le paradoxe en enlevant toute pertinence ontologique de la visée des secondes. Pour le matérialisme rationnel, bien qu'il conserve que les seuls objets physiques objectifs soient les particules élémentaires, les lois macroscopiques ne se réduisent pas à de commodes approximations. Ces lois visent également la structure de l'univers. Les particules élémentaires possèdent des lois et leur combinaison forme de nouvelles lois et ainsi de suite. Les lois macroscopiques émergent des lois des phénomènes plus petits et participent à l'émergence de lois des phénomènes qui les englobent. Ces conceptions s'appuient sur la notion d'échelle, plus particulièrement en mathématique sur la notion d'objet fractal (Mandelbrot, 1975).

La métaphore des poupées russes permet de mieux comprendre une telle conception. Chaque poupée contient une poupée qui elle-même en contient une autre et ainsi de suite jusqu'à la dernière. Chaque forme s'emboîte et conditionne ainsi à la fois la précédente et la suivante. La forme de la première poupée s'identifie aux propriétés de la matière. Ainsi,

pour le matérialisme rationnel, chacune de ces formes traduit une réalité structurelle dont l'étude permet de découvrir indirectement les autres. En revanche, le matérialisme radical refuse l'emboîtement, la forme de la poupée conditionne la taille mais pas la forme de ce qu'elle contient. Néanmoins, pour les deux mouvements, la référence dans la constitution du savoir scientifique reste les objets élémentaires, la matière. Ils privilégieront l'étude des éléments constituant les phénomènes plutôt que l'inverse.

C - Le physicalisme et la cognition artificielle

En définitive, ces deux mouvements matérialistes examinent la cognition d'un individu comme n'importe quel autre phénomène physique. La psychologie naïve échoue dans l'étude de la cognition car elle donne un modèle ad hoc très évasif, bâti sur des intérêts autres que scientifiques. Elle ne peut pas prétendre approcher la réalité. Toutefois, même avec une approche scientifique, de sérieuses difficultés apparaissent pour proposer des modèles satisfaisants. Cela signifie ici que le grain d'analyse, l'individu, ne convient pas. Plus précisément, selon le matérialisme radical, il est nécessaire de descendre au niveau des cellules voire des macromolécules qui les constituent. Selon le matérialisme rationnel, la conception de structures imbriquées offre deux alternatives. La première propose que le phénomène sera plus facilement compréhensible en étudiant la structure du réel qui régit les phénomènes constituant l'individu. Dans ce cas, l'étude du substrat reste primordiale mais elle doit être complétée par la recherche des différents niveaux d'analyse permettant d'expliquer les dynamiques engendrées par la forte connectivité des neurones. La seconde alternative suggère que la cognition sera plus compréhensible avec les théories d'un niveau supérieur, c'est-à-dire du point de vue de l'évolution des espèces. Cette seconde voie sortant du cadre de la description de l'individu, elle sera reprise plus en détail dans la deuxième section de la partie sur la description de la construction de l'individu.

Ces deux approches de la cognition individuelle se fondent sur les différents procédés d'exploration du système nerveux qui forment les neurosciences. Elles se scindent principalement en deux disciplines, toujours en interaction, la neurobiologie et la neuropsychologie. La première se concentre exclusivement sur les mécanismes mis en œuvre dans le cerveau indépendamment du comportement de l'individu. La seconde, sans prétendre à l'identification entre intention et substrat biologique, désire dépasser le problème du behaviorisme radical en complétant les observations environnementales et comportementales avec les données neurologiques. De façon caricaturale, en comparant le cerveau avec une automobile, la première démarche souhaite relever toutes les pièces ainsi que leurs degrés de liberté avec leurs voisines et la seconde souhaite regarder l'évolution des pièces pendant le fonctionnement de la voiture, plus particulièrement observer des dommages sur certaines pièces afin de mieux saisir leur rôle. Mais, le nombre de pièces et de niveaux d'analyse d'une automobile n'ont pas de commune mesure avec celles d'un cerveau. Les deux disciplines conservent en leur sein une grande variété d'objets d'études et d'échelles comme le montre le large éventail d'appareils de mesures et de protocoles expérimentaux.

Par ailleurs, les études peuvent porter également sur le système nerveux d'autres animaux. L'intérêt de ces études est double : d'une part, éviter des expériences traumatisantes pour les sujets humains. En effet, des ressemblances structurales entre le cerveau d'une souris (ou d'un singe) et celui d'un homme autorisent sous certaines conditions à formuler des hypothèses sur le rôle de certaines zones cérébrales. D'autre part, si une méthodologie s'avère efficace pour la compréhension du fonctionnement d'un système nerveux, même simple, alors il est probable que celle-ci puisse être réutilisée pour

l'étude d'un autre organisme. À noter que ce type de travaux provoque un autre désaccord entre le neuromimétisme et le symbolisme. Ce dernier considérant le langage comme le support et la raison d'être de la cognition, les organismes ne possédant pas de langage n'offrent aucun intérêt pour les sciences cognitives. La question sur la nécessité du langage sera abordée dans la partie sur la description de la construction de l'individu. Avant de continuer l'étude des différences avec le mentalisme et plus précisément avec le subsymbolisme, les deux grandes tendances de la modélisation et de la robotique dans le cadre du neuromimétisme doivent être présentées. Ces deux tendances illustrent respectivement le matérialisme radical et le matérialisme rationnel.

La première tendance, le matérialisme radical, refuse la modélisation numérique. Autrement dit, seule la modélisation analogique qui ne s'appuie pas sur des calculs effectués par un ordinateur se révèle fiable. Pour les tenants de la famille FP, cette position se justifie par le fait que les concepts et leur calcul ne possèdent aucune légitimité pour simuler le réel puisque ce ne sont que des inventions intellectuelles. Par ailleurs, la discrétisation qu'impose la numérisation ne correspond pas à la continuité temporelle et spatiale de la physique macroscopique. Or, la simulation d'un phénomène cognitif est suffisamment complexe pour ne pas y rajouter une distorsion induite par la digitalisation.

La démarche et les travaux de Franceschini (2002) sur le système nerveux de la drosophile peuvent être interprétés dans ce cadre. Ces recherches se focalisent plus particulièrement sur le système visuel. Cet insecte possède environ un million de neurones et quelques milliers d'entre eux sont dédiés aux traitements visuels. En plus de l'avantage de posséder un petit nombre de neurones, la configuration neurale se révèle suffisamment stable d'un individu à un autre pour répertorier les neurones et leurs connectivités. Ces études dégagent un certain nombre de mécanismes supposés être à l'origine de certains comportements tel que le suivi ou la stabilisation. Afin de conforter ou non ces hypothèses, l'implémentation robotique de ces seuls mécanismes et leurs mises en œuvre deviennent obligatoires. Concernant la vision de la mouche, Franceschini (2002) a montré la pertinence de certains mécanismes, mais pour obtenir le reste du comportement de la mouche il faut continuer à compter chaque neurone et à comprendre leur rôle.

À l'inverse, la seconde tendance, plus répandue, accepte la modélisation numérique. Principalement, deux types d'étude utilisent la MA : celles portant sur l'activité d'un neurone et celles portant sur l'activité d'une assemblée de neurones. En effet, le premier souhaite décrire fidèlement le comportement d'un neurone pour mieux évaluer l'influence et l'importance des différents facteurs. L'estimation de ces derniers se révèle parfois difficile, la modélisation offre alors la possibilité d'ajuster les valeurs et de vérifier que les équations imaginées se rapprochent des observations. Mais seul le second type d'étude conduit à la réalisation robotique. Deux objectifs majeurs incitent à la modélisation d'un réseau de neurones. Le premier souhaite profiter de l'étude comparative des dynamiques issues des réseaux de neurones biologiques et des réseaux de neurones artificiels pour évaluer l'incidence des approximations des neurones artificiels, et ainsi optimiser les modèles de neurones dans l'étude de leur dynamique. En effet, les modèles réalistes des neurones prennent trop de ressources de calcul pour simuler dans des temps raisonnables la dynamique des réseaux. Le second objectif veille à dégager des lois d'organisation qui aboutissent dans la pratique à certains types de capacité.

Cette démarche ressemble à celle évoquée précédemment avec le système visuel de la mouche, cependant elle diffère par le point suivant : le matérialiste radical cherche simplement un niveau d'approximation satisfaisant alors que le matérialiste rationnel

cherche un niveau d'approximation satisfaisant qui de plus renvoie à l'existence d'une loi de la nature. Mais au final, dans les deux cas, la vérification de l'intérêt pratique des mécanismes dégagés reste la raison d'être de la robotique, comme le montre les travaux sur le rôle de l'hippocampe. Ceux-ci s'appuient sur des études *in vivo* sur le rat qui ont mis en évidence une relation topologique entre l'activité neuronale de l'hippocampe et l'emplacement de son proche environnement. En synthétisant les données relatives à l'organisation neurale de cette région cérébrale, Gaussier (1994) a reproduit les capacités mnésiques suffisantes pour qu'un robot puisse se construire une représentation d'un environnement simple et contrasté puis s'orienter en fonction de la récompense octroyée par le programmeur. Toutefois, comme avec le système visuel de la mouche, il faut poursuivre l'examen neurophysiologique pour aller plus loin que des comportements élémentaires et circonscrits.

D - Les différences avec le subsymbolisme

La modélisation neuromimétique vise à imiter les observations, mais les retombées de cette recherche participe aussi à améliorer les neurones formels et à élargir leurs applications. Les travaux de Thorpe (2001) illustrent ce transfert qui se situe à mi-chemin entre le neuromimétisme et le subsymbolisme. Ils s'appuient sur le fait que les neurones ne déchargent pas simultanément et que cette latence entre les potentiels d'actions (impulsion électrique délivrée par un neurone stimulé) ne s'assimile pas seulement à du bruit mais qu'elle contient également de l'information pertinente pour le système. Ainsi, la prise en compte de la dimension temporelle dans un modèle de neurones, même simplifié, offre la possibilité d'utiliser un codage temporel pouvant accélérer certains calculs. Cette approche a permis de développer de nouveaux algorithmes de reconnaissance de visages tout en conservant des techniques classiques de l'apprentissage supervisé. Ces réalisations conjuguent les outils bio-inspirés et les objectifs d'ingénierie. Cette position hybride se trouve facilitée par la proximité des hypothèses métaphysiques entre le subsymbolisme et le neuromimétisme fondé sur un matérialisme rationnel, mais en déduire que seul le statut de l'intentionnalité fait la différence entre les deux courants revient à négliger trois autres points de désaccord, qui seront successivement décrits.

i - Premier point de désaccord avec le subsymbolisme

Le premier point de désaccord se porte sur l'assimilation du neurone à une unité de calcul élémentaire réductible à une machine de Turing. Le neuromimétisme ne se résout pas à réduire le neurone biologique à une simple unité de calcul pour deux raisons. La première raison, qui invoque la complexité des mécanismes physiques parallèles, ne permet pas de comparer le fonctionnement d'un neurone à la machine séquentielle tel que le propose le subsymbolisme. La Figure I-16 qui schématise la structure type d'un neurone illustre la complexité de ce type de cellule. Les traits les plus saillants sont : la topologie des connexions qui joue un rôle important dans l'intégration des potentiels d'action et de ce fait la réactivité du neurone, et les boucles de régulation avec l'extérieur, ou exclusivement internes, qui transforme le neurone en un système actif. La compréhension de ces mécanismes passe par la compréhension des phénomènes physiques qui régissent les macromolécules constituant la cellule. Par ailleurs, il existe un large éventail de types de neurones de par leur forme ou leurs neuromédiateurs. La seconde raison de ne pas réduire le neurone à une unité de calcul porte sur la nature des mécanismes : les machines de Turing fonctionnent à partir d'un calcul discret alors que le neurone possède à la fois une capacité de discrétisation mais aussi d'intégration du continu. La seule manière de rendre compatible les deux approches, le matérialisme rationnel et le subsymbolisme, serait de

considérer les particules élémentaires comme des machines de Turing. L'identification se situant au niveau des particules élémentaires, il n'y a plus de différence entre matière et information, c'est la physique digitale (Zuse, 1967). Dans ce cas, l'univers pourrait être interprété comme une machine de Turing. Le chapitre suivant reviendra sur les difficultés d'un tel point de vue.

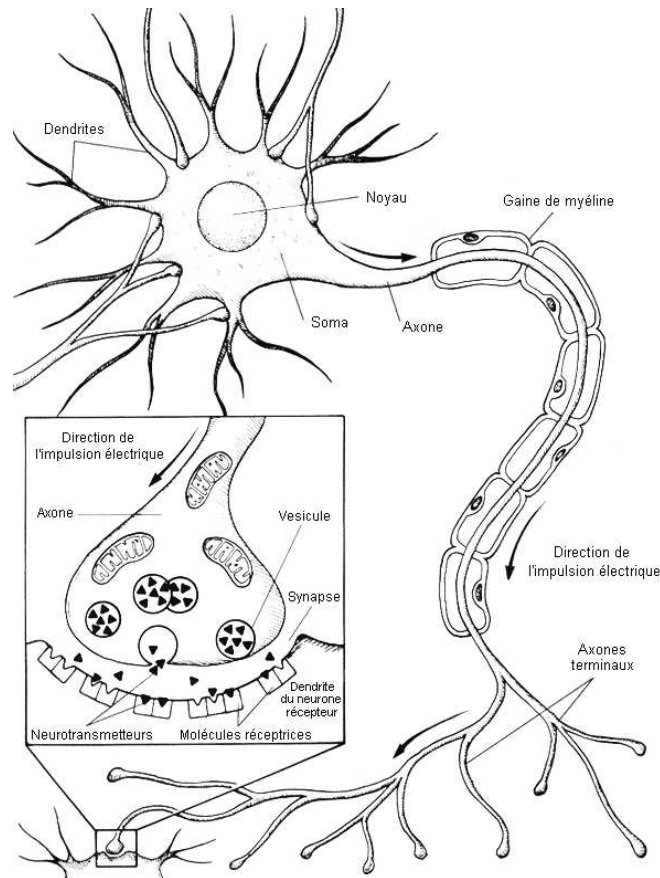


Figure I-16 : Structure type d'un neurone et de sa connexion synaptique.

ii - Deuxième point de désaccord avec le subsymbolisme

Le deuxième point de désaccord concerne l'apprentissage. Le subsymbolisme réduit l'émergence à l'encapsulation de modules où chaque module possède une tâche bien définie. L'apprentissage s'assimile à une optimisation de l'accomplissement de la tâche. Un module peut être constitué d'autres modules, l'apprentissage aura alors pour objet la coordination et la pondération des modules mais non leur modification. Un module est une fonction, soit une relation qui à chaque élément de son ensemble de départ associe au plus un autre élément de son ensemble d'arrivée. Dans cette perspective, lorsque le module se substitue à un réseau de neurones artificiels, l'apprentissage s'applique sur tout le réseau en vue de l'objectif, c'est-à-dire que la modification du réseau s'effectue selon un algorithme extérieur à celui-ci. Or, les neurosciences remettent en cause deux aspects : l'hémogénéité de l'apprentissage supervisé et le cloisonnement de l'apprentissage. Concernant le premier aspect, Schulz (1990) a montré effectivement l'existence des diffusions de neurotransmetteurs sur des zones liées à des processus d'apprentissage en fonction de la réalisation ou non d'une attente ; cependant la précision des connectivités indique, en général, que les mécanismes d'apprentissage s'effectuent localement de neurone à neurone.

À l'aide du conditionnement pavlovien de l'aplysie, Kandel (1988) a mis en évidence que l'apprentissage résulte d'une adaptation des synapses en réaction avec la fréquence du stimulus. Dans cette expérience, le conditionnement résulte de la réorganisation en cascade des neurones par leurs interactions mutuelles. Ce principe renvoie à celui d'auto-organisation qui sera traité dans la section suivante sur le connexionnisme. Le second aspect ne correspond pas aux multiples mécanismes liés à la plasticité neuronale. Les mécanismes cytoplasmiques peuvent jouer sur la sensibilité post-synaptique ou pré-synaptique ou somatique et peuvent également créer de nouvelles connexions ou multiplier les synapses à certains endroits. Enfin, des modifications génétiques de la cellule ne sont pas exclues. La diversité de ces mécanismes montre que de nombreuses pressions agissent sur le neurone ou l'assemblée de neurones, contrairement à la vision du subsymbolisme qui ne considère qu'une pression à la fois. Ces divergences sur l'apprentissage soulignent en définitive le troisième point de désaccord entre le neuromimétisme et le subsymbolisme.

iii - Troisième point de désaccord avec le subsymbolisme

Ce troisième point de désaccord se trouve dans la différence d'appréciation du terme « fonction ». Le subsymbolique présuppose une organisation fonctionnelle de telle sorte que chaque bloc possède une finalité, indépendamment des autres a priori, dans sa définition mais non dans son application. Fin et moyen se dissocient. Ainsi la connaissance de toutes les fins suffit à comprendre le comportement global. Or, à cause de la forte connexité du tissu cérébral, les neurosciences ne présupposent ni cette décomposition ni cette dissociation. Pour le neuromimétisme, le terme *fonction* se comprend plutôt au sens de *rôle*, autrement dit de comportement d'un élément dans un ensemble. Ainsi les modules se définissent pour et par les autres, soit la fonction se définit toujours dans un contexte. Cependant, ici, le contexte se réduit souvent à l'activité cérébrale et n'inclut pas l'environnement contrairement à ce que supposerait une démarche écologique (abordée dans la partie examinant l'individu comme un processus qui se construit). Ce troisième point explique que le subsymbolique s'autorise à attribuer une téléonomie au système robotique alors que les réalisations robotiques de type neuromimétique se restreignent à l'imitation de l'observé.

En conclusion, le neuromimétisme, en assimilant la cognition à un phénomène physique comme un autre, permet de dégager des propriétés pertinentes telles que la notion de rôle dans le sens de contexte de fonctionnement et d'influence sur les autres éléments, ainsi que la multiplicité des interactions participant à la plasticité d'un « module ». Mais cette démarche n'aboutit cependant pas à la conception d'un système autonome. Cette situation s'interprète de deux façons : soit la capacité des appareils de mesures et de calculs se trouve insuffisante pour appréhender et simuler le cerveau humain, soit les théories physiques négligent des aspects fondamentaux de la mécanique quantique qui pourraient expliquer la liberté individuelle comme le suggère Penrose (1994). Mais dans les deux cas, même si l'avenir permettrait de combler ces manques, le neuromimétisme n'expliquera pas la cognition, mais seulement le fonctionnement de son substrat. Certes, les connaissances acquises permettront de soigner des lésions, des dégénérescences ou autres déséquilibres, mais elles ne permettront pas de construire un système autonome pour un environnement spécifique, car la connaissance de la cognition s'arrêtera à la réplication. Selon ce critère, la description du substrat de la cognition ne permet pas d'accéder à sa compréhension. Face à cette impasse, le connexionnisme se détache de la matière comme point de départ ou de référence, en se concentrant sur la notion d'organisation.

3.1.2.2. Le connexionnisme

Le terme connexionnisme regroupe habituellement toutes les démarches utilisant des réseaux de neurones artificiels. Ici, sera retenue une autre définition, se fondant davantage sur l'objet d'étude que sur la technique employée : le connexionnisme étudie la coordination spontanée d'un ensemble d'éléments dont l'activité de chacun dépend de celle des autres suivant leurs relations, en un mot leur connexion, mais non de la finalité de la coordination. Par rapport à la première définition, ce changement aura pour conséquence d'évincer certaines approches utilisant des réseaux de neurones artificiels et d'amener d'autres qui n'en font pas usage. Avant de présenter ces approches, les enjeux du connexionnisme doivent être élucidés afin de comprendre la diversité des courants qui le composent. Tout d'abord, (A) il sera montré à l'aide d'une analyse critique du positivisme logique dans quelles mesures les positions épistémologiques des deux familles philosophiques soutenant le connexionnisme diffèrent des positions matérialistes évoquées précédemment, et quelle est l'influence sur le rôle des méthodes d'investigation, sur l'appréciation des trois types de description et surtout sur la conception de la cognition. De cette conception, (B) deux tendances seront identifiées dans le connexionnisme : l'étude des propriétés nécessaires à l'émergence de la cognition et la formalisation du rôle de la cognition. Ensuite (C), de nouveaux points de discordances, en plus de ceux soulevés par le neuromimétisme, viendront compléter la distinction entre l'éliminativisme et le subsymbolisme. Enfin, la synthèse de ces différentes approches permettra de dégager les limites du connexionnisme.

A - Analyse critique du positivisme logique

Le connexionnisme peut naître soit dans la famille FC, soit dans la famille FP. Le connexionnisme issu de la première, qui ne tolère que l'HOI, perçoit l'univers comme un mouvement éternel dont les seuls êtres permanents sont les lois qui le régissent. La position qu'assigne à la science physique le connexionnisme issu de la seconde famille (FP), bien qu'il accepte l'existence d'objets objectifs permanents, ressemble fortement à celle que lui assigne la première famille (FC). Cependant, les raisons de cette similitude doivent être examinées en détail car elles expliquent également la différence de conception de la cognition entre le neuromimétisme et le connexionnisme. Afin de mieux les saisir, quelques commentaires doivent être apportés. Plusieurs conceptions de la science physique peuvent coexister au sein de la famille FP. Plus particulièrement, le point de litige entre la conception de la physique du neuromimétisme et celle du connexionnisme se situe sur la valeur de la dichotomie. Pour le matérialisme rationnel, cette dernière se justifie par le fait que les théories se construisent à partir de faits avérés permettant de découvrir de nouveaux faits et ainsi de suite. L'idée qu'il existe une base de faits sur laquelle la science puisse se fonder exclusivement constitue le fer de lance du positivisme logique. Or, des objections fortes peuvent être levées contre une telle position épistémologique, mettant ainsi à mal la légitimité de la dichotomie. La compréhension de ces arguments passe par une explication plus fine du positivisme logique.

Le positivisme logique dont Carnap (1928) fut l'un des fondateurs souhaite constituer l'unique et véritable base de connaissance sur le monde, soit la connaissance scientifique, à l'aide d'une méthodologie irréprochable qui servirait également de critère pour se démarquer des connaissances non scientifiques. Cette méthodologie irréprochable s'appuie sur une nouvelle classification des énoncés relatifs au monde en fonction de leur vérification : les énoncés d'observation, les énoncés théoriques et les énoncés métaphysiques. Le premier type d'énoncé présente l'avantage d'être directement vérifiable

par l'expérience immédiate puisqu'il ne décrit qu'un état ou événement du monde ; « sur la table se trouvent trois pommes » peut être vérifié par la constatation que sur la table se trouvent trois pommes. Le deuxième type d'énoncé, l'énoncé théorique, ne se vérifie qu'indirectement, c'est-à-dire qu'il se vérifie seulement par le contrôle des énoncés d'observation auquel il peut aboutir par déduction. Enfin, l'énoncé métaphysique introduit des dogmes empêchant toute déduction d'un énoncé d'observation qui permettrait une vérification indirecte. Le positivisme logique déduit de cette classification que l'édifice scientifique doit se former à partir des deux premiers types d'énoncés uniquement. Plus précisément, ce sont les énoncés d'observation qui doivent constituer la base empirique sur laquelle se fonde cet édifice. Cette base empirique doit être incontestable et irrévocable ; toutefois la valeur de certains faits peut éventuellement être ajustée suivant l'amélioration des appareils de mesures, mais cela ne remet pas en cause le rôle de la base empirique.

En définitive, les énoncés d'observation reflètent la réalité extralinguistique, autrement dit, le monde est un langage matériel dont une traduction en un langage formel existe potentiellement. Cette conception rappelle celle du symbolisme et par conséquent des problèmes liés à l'établissement d'une vérité synthétique. En réaction, le positivisme logique peut radicaliser son discours en distinguant les concepts portant sur les choses qui peuvent toujours être sujets à discussion comme « chaise » et les concepts portant sur les données sensorielles comme « rouge », distinction ignorée par le symbolisme. Ainsi, le positivisme logique pense restreindre le problème de la correspondance en ne considérant que les concepts sensoriels stricts. Malgré tout, la démarche conserve en elle-même deux hypothèses fondamentales qui n'échappent pas à la critique : (i) l'hypothèse qu'un langage d'observation neutre existe et (ii) celle que la perception suffit à justifier un énoncé d'observation.

i - Critique sur l'existence d'un langage d'observation neutre

La ligne de critique au sujet de la première hypothèse se décompose en trois points qui peuvent se comprendre soit individuellement soit ensemble. *Le premier point de la ligne de critique* souligne que les mots utilisés dans un énoncé d'observation renvoient à des concepts donc à des théories qui les sous-tendent. Par exemple, « cheval » désigne un ensemble d'individus présentant certaines ressemblances perceptives, en revanche « Ourasi » désigne un célèbre cheval de course présentant certaines dissemblances perceptives avec tous les autres chevaux. En étendant le principe à tous les concepts descriptifs, apparaît alors un réseau hiérarchisé de dissemblances et de ressemblances. Par ce tissage, le langage impose un découpage conceptuel de la réalité. Toutefois, ce découpage ne se révèle pas universel, de nombreuses classifications existent s'appuyant sur des jeux de concepts différents et, plus grave, celles-ci peuvent se montrer incompatibles. Autrement dit, tout énoncé se trouvant chargé de théories, il n'existe donc pas de langage neutre (Popper, 1934). Le positivisme logique peut arguer qu'il faut se limiter strictement aux sensations (Mach, 1900), mais cette position souffre de deux objections, dont la dernière est irrévocable. La première objection rappelle que le nombre de sens communément recensé ou du moins naïvement considéré (l'ouïe, l'odorat, le toucher, le goût, la vue et la proprioception) se trouve bien inférieur au nombre de types de capteurs sensoriels révélés par l'anatomie. Dans ce cas, si la définition des sensations n'est pas donnée spontanément et sans ambiguïté par le simple ressenti, il ne semble pas pertinent de réduire la sensation à la notion de capteur puisque la correspondance n'est pas aussi évidente (Grice, 1962). La deuxième objection part du principe qu'il est effectivement possible de réduire la perception à la valeur d'un capteur et que celle-ci s'identifie à la

brique de base d'un énoncé d'observation. Il n'en demeure pas moins, de par la nature même du capteur, que cette information reflète un point de vue, un aspect particulier de la réalité et qu'elle aurait pu être autre avec un autre type de capteur. La physique du capteur incarne donc également une théorie. Ainsi, il n'existe pas de segmentation de la réalité qui soit neutre, et refuser de découper le monde revient à interdire tout énoncé donc toute science.

Le deuxième point de la première ligne de critique met en évidence deux types d'énoncé d'observation : ceux qui se cantonnent aux impressions sensorielles et ceux qui invoquent les entités qui les produisent. Le second type introduit des croyances dont la révision peut modifier la signification d'un énoncé d'observation. Par exemple, en s'inspirant de celui de Frege (1892), la planète Vénus apparaît dès le coucher du soleil et avant son lever. Comme il ne fut pas établi au début qu'il s'agissait du même astre, plusieurs noms lui furent attribués : dans le premier cas l'étoile du matin ou Eosphorus, et dans le second cas l'étoile du soir ou Hesperus. En regardant le ciel à l'aurore, Pierre déclare, E6 « je vois Hesperus ». Si Jean, qui l'accompagne, est un astronome, il conviendra qu'il s'agit de Vénus, quel que soit le nom invoqué, et que l'E6 est vrai. Maintenant, si Jean considère que Hesperus et Eosphorus ne dénotent pas le même astre, alors il jugera E6 faux. La signification d'E6 ne dépend donc pas de la perception mais des croyances associées aux entités qui doivent la provoquer. Selon Feyerabend (1975), cette sensibilité aux croyances peut mener à ce que deux théories soient totalement incommensurables bien qu'elles travaillent sur des phénomènes semblables. Comparativement, un énoncé d'observation du premier type qui repose uniquement sur les sensations est plus stable : « au matin, je vois un point brillant dans le ciel proche de l'horizon ». Cette différence qualitative, également soulignée par Frege (1892), incite le sensualisme à restreindre les véritables énoncés d'observation au premier type puisque ceux du second se révèlent facilement révocables, en contradiction avec le principe d'une base des faits immuables.

Néanmoins, se limiter à ce genre d'énoncé conduit à reconsidérer la notion de causalité. En effet, dans le cas où les deux types d'énoncé sont acceptés, les expériences permettent de dégager les objets du monde ou du moins de tendre vers une définition fiable, autrement dit de dépasser les données sensorielles. En reprenant le vocabulaire de Kant (1787), se distingue alors le phénomène « a », l'impression sur les sens « un point brillant », qui est produit par le noumène « A », l'entité essentielle « Vénus », la chose en soi (Figure I-17). Cette transcendance permet ensuite de raisonner directement sur les objets et ainsi de prédire des faits : invariablement, nécessairement « A » implique « B ». La science découvre et collectionne ces objets (l'électron, la gravité, etc.) puis les manipule pour prédire et façonner le monde des hommes. Or, dans le cas où le second type d'énoncé est désavoué, la causalité se réduit à une succession observations : à un moment donné « a » précède « b ». Pour Hume (1748), les entités physiques étant toujours construites sur des croyances, l'homme ne peut pas déceler une nécessité ontologique entre la succession de deux événements. La science ne se réduit pas pour autant à une description journalière du monde. Des régularités et des corrélations peuvent être relevées, offrant la possibilité d'effectuer des calculs statistiques afin d'évaluer des probabilités d'apparition. Ces règles inductives se comprennent comme un moyen de condenser l'information issue d'un cahier d'observations. Par exemple, jour 1, « a » puis « b » ; jour 2, « a » puis « b » ; ... ; jour n, « a » puis « b » est remplacé par la règle : « a » précède généralement « b ». Selon cette perspective, les entités physiques deviennent des subterfuges destinés à effectuer cette compression d'informations au même titre qu'une recette de cuisine. Cependant, les anticipations déduites de ces règles ne possèdent qu'une légitimité pratique, puisque

l'induction du passé vers l'avenir ne trouve pas de justification rationnelle et qu'aucune connaissance ontologique n'est accessible. Mais comment une physique se retréignant à la description des régularités peut-elle formuler des hypothèses devant la faire avancer ? Comment peut-elle donner une sémantique aux règles qu'elle produit puisqu'il n'y a pas de différence entre : « il fait chaud dans la chambre parce que le thermomètre est élevé » et « le thermomètre est élevé parce qu'il fait chaud dans la chambre » ?

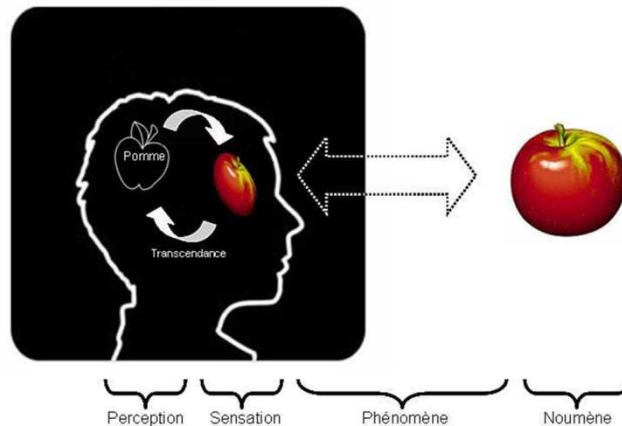


Figure I-17 : Schéma illustrant une position Kantienne de la perception, le noumène est la pomme en soi, le phénomène est la manière dont la pomme se révèle à l'observateur, la sensation est l'impression sur les sens du phénomène, et la perception transcende les sensations pour approximer le noumène. Le terme « transcender » renvoie à l'idée que la perception non seulement introduit des choses qui sont au-delà de l'expérience sensible présente mais également aboutit à la conscience d'un objet qui dépasse ses apparences.

Le dernier point de la ligne critique sur l'existence d'un langage d'observation neutre concerne la notion de traductibilité d'un langage vers un autre. L'importance de cet aspect s'impose parce qu'un langage neutre est universel et doit se traduire par conséquent dans toutes les langues. Mais l'analyse de cette opération de traduction met en évidence le fait que le langage souffre toujours d'une indétermination de ce à quoi il renvoie. Engel (1994) résume ainsi la conclusion des travaux de Quine (1969) : « Selon Quine, pour un ensemble D de données comportementales et physiques, s'il existe un manuel de traduction T1 traduisant les phrases d'une langue L dans une langue L', il est toujours possible de forger au moins un autre manuel T2 incompatible avec T1 (i.e. établissant des corrélations entre L et L' distinctes de celles qu'établit T1), mais néanmoins compatible avec D. Ceci signifie que la traduction, et par conséquent l'attribution de signification au langage des locuteurs, est toujours indéterminée ». Cet argument met en défaut le projet du positivisme logique qui reposait justement sur la non ambiguïté de la base de faits dont résultait le langage scientifique universel.

ii - Critique de l'hypothèse stipulant que la perception suffit à justifier un énoncé d'observation

La seconde ligne de critique concernant la perception comme justification d'un énoncé d'observation ne comporte qu'un seul argument majeur : les illusions et les hallucinations ne se distinguant pas des autres sensations perceptives avant qu'elles ne soient discernées, cela conduit à considérer la perception comme toujours révisable. Autrement dit, un énoncé d'observation se trouve motivé par une perception mais pas justifié par celle-ci (Popper, 1934). Par ailleurs, comme il a été montré précédemment, les

illusions et les hallucinations altèrent la solidité à la fois du jugement perceptif et du jugement existentiel. Par conséquent, les deux types d'énoncés d'observation tombent sous la critique de la faillibilité de la perception. En somme, un énoncé d'observation, en plus d'être chargé de théorie, devient une hypothèse à part entière. Par exemple, la perception motive à tenir pour vrai l'énoncé E1 « sur la table se trouvent trois pommes », hypothèse qui autorise ensuite à réfléchir sur la décision d'en saisir une. L'action de manger la pomme ne reviendra pas à porter une preuve de sa réalité mais simplement à renforcer les relations conceptuelles entre les choses perçues, soit la cohérence.

L'objection à cet argument perceptif ne nie pas qu'un individu ne peut pas prétendre à l'objectivité puisque sa perception résulte d'un point de vue. En revanche, elle affirme que la mise en commun des points de vue, des subjectivités, permet un consensus qui élimine les affabulations parasites, soit une intersubjectivité qui vise l'objectivité (Nagel, 1986). Autrement dit, le travail collectif permet de dépasser la perception individuelle par la confrontation des points de vue. Plus précisément, pour le positivisme logique non sensualiste, l'activité sociale complète qualitativement la perception (transcendance individuelle) qui est une tentative de reconstruction du noumène par l'observation de son phénomène (Figure I-18).

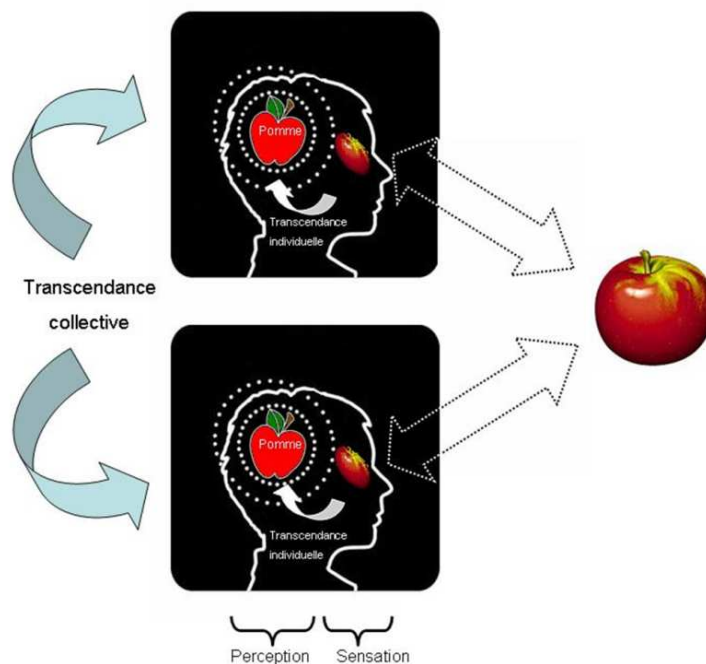


Figure I-18 : Illustration de la transcendance collective selon le positivisme logique qui complète et harmonise les transcendencies individuelles par rapport à la figure précédente.

Pour le sensualisme, le travail collectif accède à la neutralisation des observations, en ce sens où la neutralisation transcende la subjectivité de la sensation. Toutefois, cette objection s'expose à deux critiques. La première critique porte sur la supposée différence qualitative entre la perception individuelle et la perception collective. En raisonnant par l'absurde, dans le cadre du positivisme logique, un groupe d'individus de même culture possède des percepts (formes perceptives) communs ainsi qu'une intelligence reposant sur une rationalité également commune. Les illusions et les hallucinations proviennent d'un mauvais point de vue. Pour compenser cela, il convient de multiplier les points de vue en

continuant l'exploration. Statistiquement, plus les données relevées viendront de points de vue différents, moins l'influence des erreurs perceptives sera grande. Compris ainsi, le travail collectif n'est qu'un moyen de paralléliser l'activité exploratoire et la transcendance sociale se révèle être une illusion produite par l'accélération de la théorisation de l'objet. Selon les principes du positivisme, la différence entre la perception individuelle et la perception collective se révèle quantitative, et non qualitative, or l'objection reste valable uniquement si la transcendance sociale a lieu.

La seconde critique de l'objection à l'argument perceptif distingue la notion de normalisation et celle de neutralisation. Selon l'objection émise par le sensualisme, les énoncés se cumulent indépendamment les uns des autres, permettant de moyenniser, de neutraliser le discours et ainsi d'obtenir un énoncé d'un point de vue de nulle part, la neutralisation se confond avec la normalisation. Pourtant, la psychologie sociale montre qu'il y a là une confusion en deux points. Le premier point souligne que la neutralisation doit améliorer la théorie, la mesure, en tant qu'énoncés d'observations collectifs, mais en aucun cas modifier la perception individuelle. Or, les travaux de Shérif (1935) sur la construction d'une norme par un groupe dément cette affirmation par l'expérience suivante : la tâche des participants consiste à évaluer la distance parcourue par un point lumineux qui en définitive est fixe car le mouvement résulte d'une illusion d'optique, l'autocinétique. Dans ce contexte ambigu, chaque participant définit individuellement sa propre valeur, sa propre norme. Lorsqu'une seule valeur est demandée à un groupe de participants, la norme ne résulte pas forcément de la moyenne des valeurs estimées, les meneurs spontanés possèdent une crédibilité plus élevée et plusieurs groupes peuvent aboutir à des normes différentes. Ensuite, de nouveau individuellement, les participants accomplissent encore une fois la tâche et les réponses données avoisinent la norme décidée par le groupe, ce qui signifie que les participants ont intériorisé la norme perceptive.

Dans ce contexte, les normes peuvent être interprétées comme des positions d'équilibre de la dynamique sociale (Lewin, 1959). Ces positions d'équilibre se comparent avec la notion d'attracteur dans le sens où les personnes tendent à résister à des propositions éloignées de la norme, autrement dit, le système reste stable malgré de petites divergences. Pour le positivisme logique sensualiste, les notions d'intériorité et d'équilibre peuvent s'intégrer au projet de neutralisation. L'intériorité se comprend comme une intégration légitime puisque le groupe a apporté des connaissances supplémentaires et que l'existence de plusieurs positions d'équilibre est normale dans la zone d'incertitude liée à l'observation de l'objet. Si l'apport de nouvelles connaissances réduit la zone d'incertitude, alors elle réduira le nombre de positions d'équilibre.

Mais ces considérations incluent deux hypothèses : que la majorité des individus donne toujours des informations justes et que la dynamique de la normalisation se nourrit uniquement de la réalité des perceptions, l'influence sociale jouant un rôle uniquement dans les situations ambiguës. Toutefois, le second point souligné par la psychologie sociale porte sur la notion du conformisme qui discrédite les deux suppositions précédentes. Le conformisme traduit le changement d'opinion sous l'influence d'un groupe qui possède déjà sa norme. L'expérience Asch (1951) montre que ce phénomène ne se cantonne pas à la zone d'incertitude liée à la qualité des sensations mais qu'il apparaît même dans des situations supposées indiscutables : 7 participants croyant participer à des tests perceptifs devaient désigner parmi 3 segments celui qui correspondait au segment de référence (Figure I-19). Cet exercice fut présenté 18 fois avec des arrangements différents, mais 6 des 7 participants étaient des complices de l'expérimentateur et donnaient une fausse réponse 12 fois sur les 18 arrangements.



Figure I-19 : Exemple de dessin présenté au cours de l'expérience d'Asch (1951), à droite la ligne de référence et à gauche les lignes de comparaison.

Dans ces conditions, d'une part un tiers des réponses fournies par les participants était conforme à l'opinion de la majorité, au lieu de 1% en situation normale, et d'autre part 76% des participants se sont conformés au moins une fois aux autres lorsque les compères donnaient une mauvaise réponse. La majorité des participants naïfs se rallie aux groupes sans renier leurs croyances mais certains se persuadent malgré tout d'avoir bien répondu. Donc, contrairement à la seconde supposition, l'influence sociale agit également sur la construction d'un énoncé d'observation commun et sur la perception individuelle bien que la scène observée soit non ambiguë. Toutefois, la diversité des opinions montre que le conformisme est un mécanisme social parmi d'autres, comme son opposé l'influence des minorités (Faucheux et Moscovici, 1971), chacun ayant ses propriétés et ses facteurs. Mais, lister exhaustivement les mécanismes et leurs facteurs soulève également des difficultés liées à l'étude d'un individu cognitif.

En définitive, accepter l'une de ces deux lignes de critique suffit pour ébranler le positivisme logique. Si tel est le cas, une doctrine issue de la famille FP se transforme en l'une des trois positions suivantes, dont seule la dernière sera propice au connexionnisme. La première position consiste à considérer que ces critiques révèlent des facteurs d'égarement et des limites à l'intelligence humaine mais qu'il est possible de minimiser les facteurs et d'approcher asymptotiquement ces limites suivant une méthode scientifique. La minimisation des facteurs peut s'interpréter comme une réactualisation de la minimisation de l'influence des idoles de Bacon (1620) et l'approche asymptotique peut être illustrée par le progrès de la science. Cette première position conduit à se retourner vers le neuromimétisme. Par rapport au matérialisme rationnel évoqué précédemment, en plus de l'incertitude liée à l'objet d'étude et aux théories, se rajoutent à la perception et à la dynamique sociale les incertitudes liées au langage. Toutefois, en pratique, il n'y a pas de différence entre le neuromimétisme issu d'un matérialisme rationnel et celui issu d'un matérialisme rationnel averti. Seules leurs visions épistémologiques divergent, visions qui seront abordées dans la conclusion de ce chapitre.

La deuxième position revient à un changement radical sur l'ambition de la science et de son objet. En effet, les critiques du positivisme logique montrent que les bases d'une science physique sont toujours liées à nos capacités cognitives, ce qui signifie que la science doit davantage se pencher sur la psychologie humaine pour savoir comment la connaissance se construit. Cette attitude est similaire à celle de Hume (1748) suite à la remise en question de la causalité. Cette deuxième position conduit nécessairement à sortir de l'éliminativisme et à trouver sa place dans une approche qui privilégie la psychologie cognitive. La dernière position intègre l'impossibilité de connaître objectivement les objets du monde, toutefois, celle-ci considère que la connaissance relative aux lois qui sont à

L'origine de leurs organisations et de leurs structures est accessible. En effet, des lois de structures apparaissent indépendamment des objets. La science au sens large dégage des régularités de structure et de dynamique sur des sujets très variés de l'économie, de la physique, de la sociologie, de la biologie, etc., ce qui laisse à penser que la connaissance stable et absolue soit constituée par ces régularités et non par les objets eux-mêmes qui dépendent du découpage du monde. Cette position revient au dualisme de propriétés évoqué précédemment, mais pour la physique. Ce mouvement scientifique peut prendre le nom de science des systèmes ou cybernétique, et plus particulièrement pour la cognition, de connexionnisme. Dans ces conditions, hormis la conviction de principe que le monde se compose d'objets, rien dans la démarche ou dans l'objectif ne différencie le connexionnisme issu de la famille FC du connexionnisme issu de la famille FP. Dans ce contexte épistémologique, la description mathématique (DM) se voit privilégiée par rapport à la description componentielle, elle-même privilégiée par rapport à la description procédurale. L'abstraction de l'objet scientifique (l'organisation prenant la place de la matière) rend théoriquement équivalentes les diverses méthodes d'investigation, mais la méthode artefactuelle, qui offre la possibilité d'éprouver directement les théories imaginées, est favorisée.

B - Le connexionnisme, héritier du positivisme logique

Maintenant que les principaux concepts ont été dégagés, une définition connexionniste de la cognition peut être envisagée : la cognition résulte d'un processus d'auto-organisation et elle vise à dégager les régularités du monde et ses structures, en somme à compresser les informations sensorielles. Concrètement, cette définition de la cognition implique que les processus perceptifs et les processus neurophysiologiques sont isomorphes. Cette compréhension de la cognition trouve écho aussi bien dans des domaines tels que la psychologie de la forme, les neurosciences, ou les mathématiques. La psychologie de la forme (Gestalt) s'intéresse plus précisément à la perception comme la capacité à saisir directement la structure globale de la scène visuelle qui ne peut être réductible à la somme des stimuli (Koffka, 1935), c'est-à-dire que l'organisation prédomine sur les éléments.

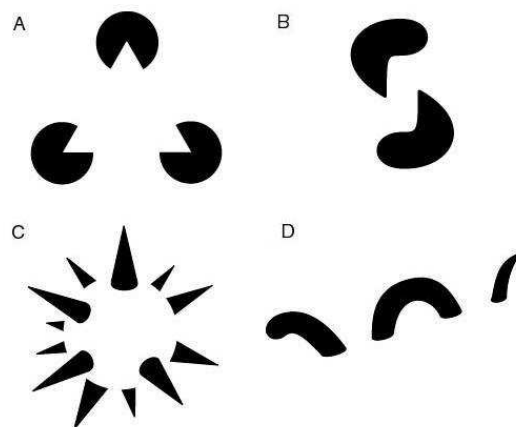


Figure I-20 : Exemple de réification où le sujet humain verra un triangle dans l'image A, bien qu'aucun triangle n'ait été réellement dessiné. Dans les images B et D, il identifiera des formes disparates comme « appartenant » à une forme simple, dans C une forme tridimensionnelle complète est vue (Lehar, 2003).

Par exemple, la Figure I-20 illustre la loi de réification qui est l'aspect constructif ou génératif de la perception par lequel le percept expérimenté contient une information spatiale plus explicite que le stimulus sensoriel sur lequel il est basé. L'influence culturelle dans l'élaboration de ces lois n'est pas forcément ignorée en avançant, dans ce cas, que les lois perceptives sont issues de principes plus généraux et du développement.

Parallèlement, les neurosciences observent que certaines zones cérébrales (même au niveau neurone) fonctionnent comme des détecteurs de coïncidences, autrement dit que le système nerveux possède des mécanismes relevant l'organisation d'événements temporels. Ces mécanismes, résultant uniquement des interactions entre neurones sans un coordinateur général, amènent à penser que les mécanismes d'extraction d'une organisation du monde se trouvent en amont de toute sémantique. Cette faculté à trouver des organisations ou des formes semble toutefois inévitable, comme le démontre le mathématicien Ramsey (1930). Un rapide résumé du théorème de Ramsey serait le suivant : « pour tout entier k (taille des parties), tout entier m (nombre de classes), et tout entier r , il existe un entier $N(k, m, r)$ tel que pour toute partition en au plus m classes de $P_k(E)$ d'un ensemble E d'au moins N éléments, il existe un sous-ensemble F de E d'au moins r éléments tel que $P_k(F)$ est entièrement contenu dans une des classes de la partition » (Sakarovitch, 2003). Ainsi, avec un minimum de descripteur et de choses à décrire, il apparaît des structures. Le désordre complet est impossible ou une structuration minimale est inévitable. Ces diverses structures représentent par exemple le type de connaissances visé par le connexionnisme. Les configurations ne dépendent ni de la nature de la relation ni des éléments mais de la combinatoire, il y a bien d'une part l'organisation et d'autre part les objets qui l'instancient.

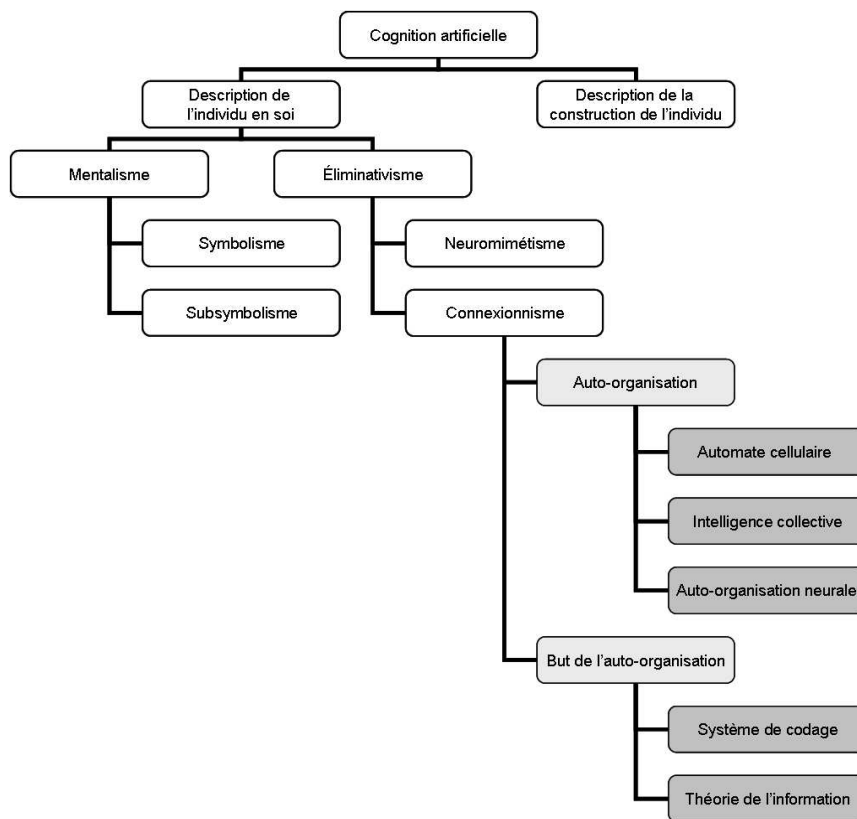


Figure I-21 : Les cinq types d'approche issus des deux facettes du connexionnisme au sein de l'arborescence de la cognition artificielle.

Dans les sciences de l'artificiel, le connexionnisme explore les deux facettes de sa conception de la cognition (Figure I-21) : l'auto-organisation dont la cognition émerge et la capacité à créer une représentation du monde qui constitue sa finalité. Plus précisément, la première facette repose sur l'idée que la cognition résulte de systèmes s'auto-organisant, dont les éléments eux-mêmes résultent d'une auto-organisation, et ainsi de suite. Cette hiérarchisation ressemble à celle évoquée dans la section sur le neuromimétisme, mais elle en diffère par le fait que cette hiérarchisation découle d'une construction. Autrement dit, elle repose sur l'évolution des organismes vivants en complexité. Sans s'intéresser directement à cette genèse, la première facette du connexionnisme se décline en trois types d'approche : (i) l'étude des principes d'auto-organisation élémentaire, (ii) l'apparition de comportements intelligents par l'activité parallèle de plusieurs éléments et enfin (iii) l'étude des propriétés organisationnelles des assemblées de neurones.

i - L'étude des principes d'auto-organisation élémentaire

Le premier type d'approche de la première facette du connexionnisme considère que la cognition est un phénomène complexe issu de multiples niveaux d'organisation et que sa compréhension passe tout d'abord par celle du vivant : pas de cognition sans le vivant. En évitant une définition du vivant, le connexionnisme souligne que les organismes dotés de capacités cognitives se composent de cellules considérées comme des organismes vivants irréductibles. Le schéma selon lequel la cognition résulte d'une complexification croissante est respecté, la cellule devient la brique élémentaire. Dans ce contexte, celle-ci se réduit à une entité avec des règles élémentaires de comportement qui réagissent en fonction de son proche environnement. La faculté de se répliquer se trouve intégrée dans le système de règles pour simuler la propension du vivant à se développer tant que l'environnement le permet. Cette faculté offre un moyen de discriminer les systèmes d'auto-organisation associés au vivant et ceux qui ne le sont pas, comme les structures dissipatives (Prigogine, 1986) illustrées par la formation de petits tourbillons, sortes de structures locales qui ne se maintiennent que par la pression des conditions environnementales. Dans la logique de complexification progressive des propriétés des systèmes, l'auto-organisation doit être complétée par la notion d'autoréplication (Von Neumann, 1966). Généralement, les automates cellulaires qui modélisent un regroupement de cellules identiques se résument en un tableau de n dimensions où chaque case représente un espace discret contenant une cellule avec un état interne, et à chaque pas de temps les cellules sont réactualisées suivant des règles comportementales fixées. Selon les conditions initiales, l'évolution des états des cases de l'ensemble du tableau forme des motifs dont la dynamique peut être étudiée. L'émergence de propriétés se comprend ici comme l'observation de motifs présentant des régularités et des comportements identifiables, alors que rien dans les règles cellulaires ne le présageait. L'automate cellulaire « le jeu de la vie » de Conway (1970) illustre cette émergence, avec une grille à deux dimensions dans laquelle chaque cellule possède deux possibles états : 0 ou 1 (« morte » ou « vivante »). À chaque pas de temps, son état et celui de ses huit voisines déterminent son état futur, selon les deux règles suivantes :

1. Une cellule à 0 passe à 1 si exactement trois voisines sont à 1.
2. Une cellule à 1 passe à 0 si le nombre de voisines à 1 n'est pas deux ou trois.

A partir de ces deux règles, cet automate cellulaire présente un très large panel de motifs cycliques : motifs oscillants (Figure I-22, ligne A), motifs oscillants en se déplaçant dans la grille (Figure I-22, ligne B), motifs oscillants produisant des motifs oscillants, etc. Certains des motifs peuvent également être stables et irréversibles, c'est-à-dire des structures ayant plusieurs antécédents possibles.

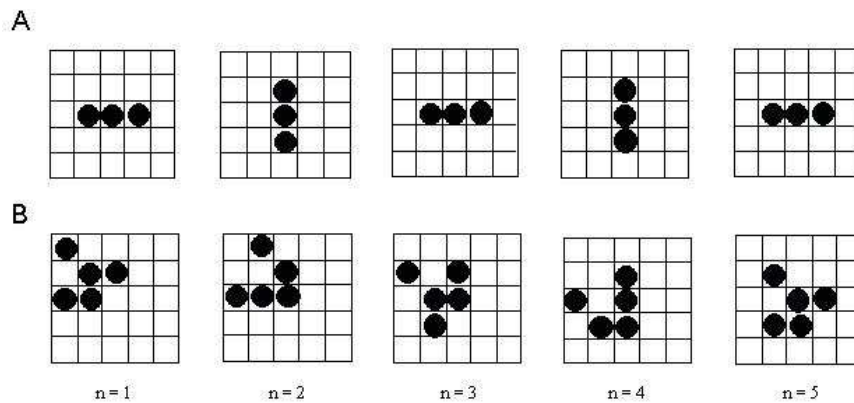


Figure I-22 : Exemple de motifs cycliques, les ronds noirs indiquent que les cellules se trouvent à l'état 1, les autres étant à l'état 0. La ligne A représente une séquence d'un motif oscillant sur lui-même, la ligne B représente la séquence d'un motif oscillant en se déplaçant

Le comportement de ces automates cellulaires révèle principalement deux points. Le premier point porte sur les conditions d'apparition de ces motifs. En effet, tous les jeux de règles ne permettent pas la génération de motifs cycliques complexes. Certaines règles de transition conduisent, quelles que soient les conditions initiales, à un état stable (ordonné) et à l'inverse d'autres conduisent, quelles que soient les conditions initiales, à un état erratique (désordonné). Entre les deux se trouvent les conditions optimales pour l'apparition des motifs différents et de comportements élaborés (Wolfram, 1983 ; Heudin, 1998). Le second point met l'accent sur les capacités de certains automates cellulaires à implémenter une machine de Turing universelle, autrement dit, sur le fait qu'ils peuvent calculer n'importe quelle fonction calculable. Le connexionnisme interprète cette propriété comme une confirmation de ce que la cognition émerge de l'activation d'unités de calculs parallèles. Malgré ces deux points, la critique majeure de cette démarche, pour accéder à la compréhension de la cognition, réside dans le fait que l'émergence ne se produit que dans un univers clos et homogène. Des réalisations robotiques explorent les capacités de systèmes basés sur la théorie des automates cellulaires (Zykov, 2005) mais ceux-ci restent confinés à l'intérieur d'un environnement contrôlé et restent en totale autarcie. Cependant, si la cognition doit servir à maîtriser les échanges entre l'intérieur et l'extérieur d'un organisme afin de maintenir son autonomie, alors les automates cellulaires ne semblent pas propices à l'émergence de la cognition en dehors de leur monde virtuel. Une façon radicale de dépasser cette impasse consiste à identifier les atomes aux cellules de ces automates. Des modélisations de phénomènes physiques comme l'évolution des molécules d'un gaz peuvent être avancées à l'appui mais il s'agit toujours de phénomènes s'exprimant dans un milieu suffisamment homogène et clos.

ii - L'apparition de comportements intelligents par l'activité parallèle de plusieurs éléments

Sans qu'il se justifie en ces termes, le second type d'approche de la première facette du connexionnisme propose une autre solution pour dépasser le confinement des automates cellulaires. Les cellules deviennent des agents évoluant dans un monde ouvert. Le comportement de l'agent est déterminé par sa configuration sensorielle riche, c'est-à-dire qui n'est pas exclusivement centrée sur leur homologue, et ce comportement agit sur l'environnement. À la différence d'une « communication » directe entre cellules, la communication entre agents se trouve médiatisée par l'environnement : c'est la stigmergie (Grassé, 1959). Par ailleurs, la modification de l'environnement influe sur le déclenchement

des règles comportementales et cette rétroaction par l'environnement permet de retrouver les propriétés des systèmes dynamiques mis en avant par la cybernétique. Ainsi, l'organisation des agents se façonne sur et par le monde, contrairement aux automates cellulaires qui sont coupés du monde. Ce type d'organisation s'inspire des travaux d'éthologie sur les animaux sociaux. Dans la plupart des cas, chaque individu réagit selon des règles réactives instinctives et aucun individu ne coordonne l'évolution du groupe, le système évolue de manière décentralisée.

Par exemple, pour la fourmi, la recherche de la nourriture la plus proche de la fourmilière s'appuie sur des mécanismes qui peuvent se décrire en deux règles : (1) une fourmi laisse des traces olfactives (des phéromones) après son passage (celle-ci s'évaporent après un certain laps de temps) et (2) la présence de phéromones affecte la probabilité d'orientation d'une fourmi au cours de son déplacement. Arrivant sur une bifurcation vierge de phéromone, une colonne de fourmis pouvant aller de manière équiprobable soit à droite soit à gauche se divisera en deux dans un premier temps, les fourmis parcourant le chemin le plus court pour atteindre une source de nourriture reviendront plus vite et laisseront ainsi plus de phéromones que celles ayant pris le chemin le plus long. Au final, la majorité de la colonne prendra un seul chemin. Ce mécanisme ressemble aux automates cellulaires sur deux points. Le premier point vient de l'émergence de motifs sous forme de coordination du travail, l'idée de rechercher le chemin le plus court ne se trouve pas dans les règles initiales. Le second point souligne l'importance des paramètres des règles, ici, la quantité de phéromones qui est liée à sa persistance et l'influence des phéromones dans le choix de la direction. En effet, une influence maximale imposera l'ordre absolu chez les fourmis mais elle ne permet plus la découverte d'autres sources de nourriture et une influence nulle introduira un désordre tel qu'il ne permettra pas d'optimiser les trajets. La différence avec les automates cellulaires est que ces motifs révèlent un avantage sélectif pour la survie du groupe dans son environnement.

Mais pour autant, la fourmilière peut-elle être reconnue comme un individu cognitif distribué ? La cognition humaine est-elle le résultat de l'auto-organisation d'agents enfermés dans le crâne, avec une vue plus ou moins commune sur le monde associée à une certaine capacité d'action ? Il faut distinguer comportement émergent et finalité attendue du comportement émergent. La recherche sur l'auto-organisation des animaux sociaux explore, à l'aide de la simulation informatique (Drogoul, 1993) ou de la simulation robotique (Mataric, 1992), les différents facteurs participant à l'émergence de comportements collectifs (spécialisation des agents, capacité de communication, etc.). Néanmoins, l'appréciation du comportement, c'est-à-dire de l'intérêt de la finalité de ces systèmes coopératifs reste à l'appréciation de l'observateur extérieur (le chercheur). L'autonomie consisterait à ce que le système reconnaisse l'intérêt de sa propre dynamique, or aucun de ces systèmes ne génère un recul sur sa propre activité.

iii - L'étude des propriétés organisationnelles des assemblées de neurones

Le troisième type d'approche peut également s'interpréter comme une proposition pour concilier les automates cellulaires et le monde en assimilant les réseaux de neurones à des automates cellulaires et les données sensorielles aux états initiaux du réseau. Cette analogie s'appuie sur le fait que la dynamique de certains automates cellulaires aboutit à des motifs stables identiques pour des ensembles d'états initiaux sensiblement différents. Ces motifs stables peuvent alors être considérés comme des attracteurs. Cette caractéristique se retrouve en physique statistique avec le modèle d'Ising (1924) qui s'identifie à un automate cellulaire. Celui-ci explique la stabilisation due à la minimisation de l'énergie d'un ensemble

de particules dont le moment magnétique, soit M^- soit M^+ , dépend du moment de leurs voisins à l'état initial. A partir de ce modèle, Hopfield (1982) propose de modéliser la mémoire comme un ensemble d'attracteurs supporté par un système dynamique, en insistant sur la ressemblance de structure entre l'ensemble des particules ferromagnétiques qui s'influencent mutuellement en fonction de leur proximité et le réseau de neurones fortement connexe dont les poids synaptiques modulent les interactions (Figure I-23). La mémoire ainsi définie est adressable par le contenu, c'est-à-dire qu'une stimulation sensorielle (état initial du réseau) partielle ou dégradée amène à une sensation prototypique connue (l'état stable final). L'apprentissage de ces sensations prototypiques dans le réseau de neurones formels totalement interconnectés s'appuie sur la règle de Hebb (1949) qui place l'apprentissage au niveau du renforcement des synapses des neurones dont l'activité est corrélée.

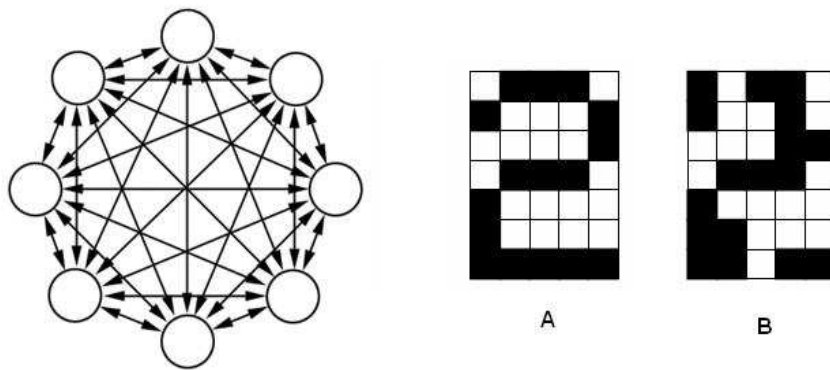


Figure I-23 : Modèle mnémotique de Hopfield (1982) : à gauche, exemple d'un réseau totalement interconnecté, à droite, les deux grilles représentent un état d'un réseau de 35 neurones (image 5x7) : la grille A est le prototype, l'état stable ; la grille B est un état initial (image du '2' bruité) qui convergera vers A.

Ce type de réseau montre que des lois d'auto-organisation et que l'association ou la corrélation entre les événements locaux comme unique source de connaissance permettent de simuler des capacités mnésiques telles que la reconnaissance de caractères, de photos ou même d'informations hétérogènes (nom et visage par exemple). Ce mode de stockage possède la fiabilité et la robustesse d'une représentation distribuée. Toutefois, cette voie de recherche pose deux catégories de problèmes. La première contient les difficultés théoriques concernant l'amélioration de l'apprentissage hebbien, la présence d'attracteurs parasites, les techniques de désapprentissage qu'ils nécessitent, la capacité de stockage liée au nombre de neurones, de leur connectivité, ainsi que la dissemblance des motifs prototypiques, et la signification de la vitesse de stabilisation. La seconde catégorie regroupe les problèmes liés à l'approche en elle-même. Bien que le principe de l'apprentissage de la mémoire et son exploitation découlent de l'auto-organisation de celle-ci, le motif appris reste fourni par l'utilisateur. Le réseau de type Hopfield diffère dans son principe du perceptron (Figure I-11) évoqué dans la section sur le subsymbolisme, mais en définitive, tous deux effectuent un apprentissage complètement supervisé. Les prototypes étant définis à l'avance, le réseau de type Hopfield propose une sorte de catégorisation alors que la cognition classe, autrement dit construit, ses catégories au fil des expériences. Ce problème prend ici plus d'importance car la cognition doit elle-même être un processus issu de l'auto-organisation et si la mémoire fonctionne effectivement comme un réseau

dynamique possédant des attracteurs, rien ne permet d'expliquer comment fonctionne le superviseur.

La seconde facette du connexionnisme recherche les mécanismes permettant d'extraire les objets et leurs liens dans le foisonnement d'informations sensorielles sans superviseur. Cet axe de recherche demeure encore dans le cadre épistémologique évoqué précédemment qui nie la possibilité d'accéder à une connaissance absolue des objets physiques. Cependant, la cause de l'inaccessibilité de l'objet réel se trouve ici liée à l'existence inévitable de médiateurs entre le monde et le sujet : les capteurs ou les sensations. En revanche, les principes de construction conceptuelle d'un objet restent soumis à des lois d'organisation dont la connaissance est accessible puisque indépendante des capteurs. Par conséquent, l'auto-organisation doit nécessairement dégager des associations, des classes de motifs et pas simplement retrouver des motifs. En incorporant l'idée d'un emboîtement de l'auto-organisation, la classification se stratifie, c'est-à-dire qu'une première classification ou une première reconnaissance de motifs se dégage des entrées sensorielles dont le résultat servira de base pour une nouvelle classification : c'est l'abstraction.

Cette idée de hiérarchisation ressemble à celle proposée par le subsymbolisme mais les motivations et les justifications diffèrent. Ces dernières proviennent des données neurophysiologiques qui confortent cette position avec l'observation de la sélectivité de certains neurones à des motifs spatiotemporels tels que la direction et l'axe d'un mouvement (Barlow et Levick, 1965). Le neurone sélectif capture par l'intermédiaire de son champ récepteur (poids synaptiques sur les entrées sensorielles) une association significative entre les différents états sensoriels. Le processus devant se répéter, toute l'attention se concentre sur les premiers traitements des données sensorielles qui se comprennent comme des systèmes de codage, un codage neuronal (Figure I-24). Plus précisément, la présentation de cette seconde facette du connexionnisme comporte deux types d'approches : le premier (iv) porte sur le contexte et les enjeux des systèmes de codage, le deuxième (v) porte sur la formalisation mathématique de ces enjeux.

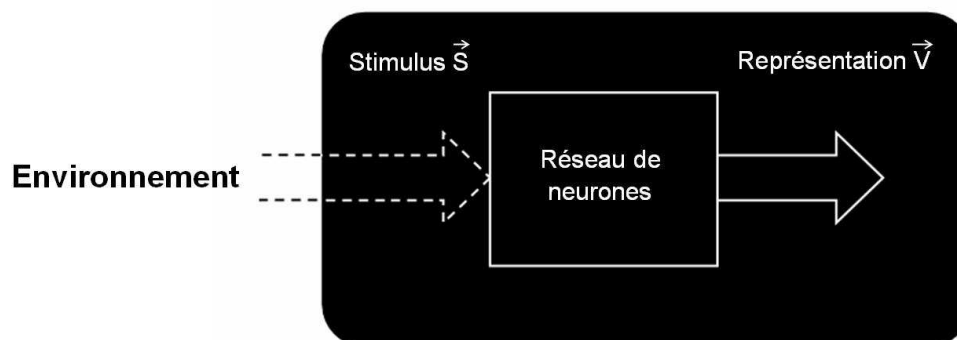


Figure I-24 : L'environnement génère sur l'organisme un stimulus qui est ensuite traité par le réseau de neurones afin d'obtenir une représentation plus adéquate pour les traitements ultérieurs.

iv - La cognition et le codage

Le premier type d'approche de la seconde facette du connexionnisme part du principe que les statistiques des images de scènes naturelles conditionnent directement les descripteurs ou champs récepteurs. Concernant la vision, les travaux dans ce domaine

montrent d'une part que les images naturelles appartiennent à un sous-ensemble de l'ensemble des images possibles (Field, 1987), et dont l'une des particularités est l'invariance par rapport à l'échelle de leurs propriétés statistiques (Simoncelli et Olshausen, 2001), et, d'autre part, que la distribution des niveaux de gris des images révèle que la valeur d'un pixel dépend de celle des autres (Ruderman, 1994). Ces résultats confortent l'hypothèse de Barlow (1961) stipulant que l'environnement perçu contient beaucoup de redondances et que l'économie des représentations peut prétendre à un avantage sélectif au cours de la phylogénèse des systèmes biologiques. Plus précisément, la minimisation de la redondance implique l'existence d'un codage optimal de type factoriel où chaque neurone doit coder une caractéristique du signal statiquement indépendante de celles codées par les autres neurones.

Dans ce sens, Kohonen (1982) propose un algorithme permettant l'auto-adaptation des champs récepteurs, en insistant sur l'organisation tonotopique de cellules corticales de bas niveau et sur leur forte connectivité excitatrice ou inhibitrice. Selon cette architecture de réseaux de neurones (Figure I-25) les entrées sensorielles prennent la forme de connexions synaptiques. Schématiquement, l'algorithme d'apprentissage se déroule en 5 étapes :

1. Initialisation : affecter aléatoirement les poids synaptiques des neurones formels du réseau qui couvrent l'ensemble des champs récepteurs, soit toutes les entrées.
2. Stimulation : présenter un motif sensoriel tiré aléatoirement d'une base d'apprentissage.
3. Élection : rechercher le neurone le plus fortement excité par les entrées.
4. Modification : ajuster des poids synaptiques sensoriels afin d'améliorer la reconnaissance et ajuster les poids synaptiques des neurones appartenant au voisinage afin d'assurer le monopole de cette sélectivité.
5. Continuation : retourner à l'étape 2 si l'ajustement reste significatif.

Au final, l'activation des neurones forme le relief d'une carte pouvant être interprétée comme une traduction, une représentation de la scène observée ou le résultat d'une classification.

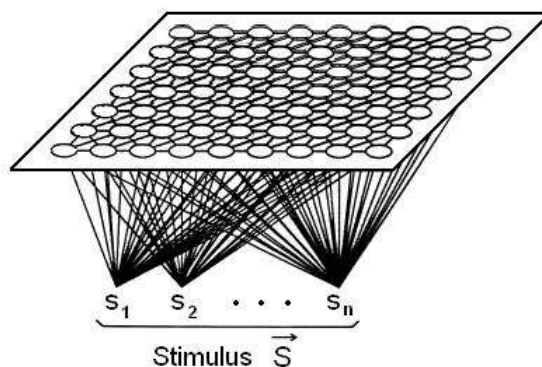


Figure I-25 : Architecture type d'un réseau de neurone auto-associatif

v - La cognition et la théorie de l'information

Le second point de la présentation de cette seconde facette du connexionnisme s'appuie sur la théorie de l'information. Celle-ci offre un formalisme mathématique pour le problème du codage qui a prouvé son efficacité dans les télécommunications. Plus précisément, la théorie se fonde sur trois concepts. Le premier concept, la quantité d'information d'un mot, traduit la pertinence d'un mot selon sa probabilité d'apparition et non pas selon sa sémantique. Par exemple, le mot « desquamer » aura plus de pertinence d'un point de vue statistique que le mot « le ». Le deuxième concept, l'entropie, se comprend comme la moyenne de la quantité d'information sur l'ensemble du vocabulaire dans un corpus donné, autrement dit, l'entropie évalue le désordre ou plutôt mesure l'incertitude du langage employé. Lorsque tous les mots sont équiprobables, toutes les prédictions se valent, l'entropie est alors maximale. Le dernier concept, l'information mutuelle, s'applique lors de la comparaison entre deux sources d'informations S1 et S2. L'information mutuelle mesure la quantité d'information apportée en moyenne par une réalisation de S1 sur les probabilités de réalisation de S2. La maximisation de l'information mutuelle traduit ainsi l'optimisation de la transmission de l'information.

Dans ce cadre, Linsker (1988) a montré que ce critère peut permettre d'auto-organiser un réseau de neurones où chaque neurone représente une source. Par ailleurs, Nadal et Parga (1994) ont démontré que la maximisation de l'information mutuelle conduit à une analyse en composantes indépendantes montrant l'équivalence formelle entre la minimisation de la redondance et la maximisation de l'information mutuelle. Ainsi, il devient possible d'interpréter la théorie de l'information comme un moyen pour expliquer les critères d'auto-organisation et leur raison d'être.

Rétrospectivement, de nombreuses études confortent l'intérêt et la portée de ce paradigme pour l'étude du système cérébral. Parmi celles-ci deux exemples complètent particulièrement bien le discours précédent. Le premier concerne la modélisation de la rétine de Atick (1992) qui en s'appuyant sur l'idée de la maximisation de la transmission de l'information tout en minimisant la consommation d'énergie liée à celle-ci trouve des résultats compatibles avec les courbes de sensibilité observées par des biologistes (Enroth-Cugell et Robson, 1966). Le second exemple, qui est plus une remarque, concerne également la vision. Il souligne la similitude entre la sensibilité des champs récepteurs du cortex visuel primaire et les filtres calculés à partir d'une analyse en composantes indépendantes effectuée à partir d'images naturelles (Bell et Sejnowski, 1997). Par ailleurs, la modélisation de ces principes dans le cadre de systèmes informatiques montre leur efficacité pour la reconnaissance de formes multi-échelles (Machrouh et Tarroux, 2001), la classification automatique d'images (Denquive, 2002) et permet même d'aboutir à des modèles attentionnels (Itti et Koch, 2001) à mi-chemin avec le neuromimétisme.

Malgré ce socle théorique, cette voie de recherche rencontre deux types de difficultés. L'une provient de la nécessité d'estimer la distribution statistique de l'environnement et des bruits intervenants à tous les niveaux ainsi que de choisir des expressions mathématiques définissant la quantité d'information et l'entropie qui influencent l'interprétation des propriétés statistiques. Dans tous les cas, ces considérations interviennent dans l'écriture des équations sur lesquelles se fonde la conception d'un algorithme. L'autre difficulté est que cette thématique du connexionnisme est confrontée à la différence entre traitement du signal et interprétation. Le premier s'effectue sur l'ensemble des capteurs et le second vise à regrouper, à classifier, à localiser une sous-partie du flux d'information traité. Le défi lancé aux théories de l'information consiste à trouver de nouveaux critères objectifs, c'est-à-dire

se justifiant mathématiquement, pour dégager les invariants structurels et temporels qui semblent être l'étape suivante en neurosciences afin d'aboutir à une analyse de scène (Poggio, 1995) et qui semblent également être indiqués par les travaux issus de la psychologie de la forme. À cela plusieurs questions s'ajoutent : le rôle de l'attention ne serait-il pas alors un mécanisme de sélection destiné à créer un sous-ensemble statistique ? Peut-on considérer la statistique de l'ensemble des capteurs comme le contexte général orientant le choix d'une estimation de distributions plus spécifique pour la recherche de motifs associés ? Répondre par l'affirmative équivaut-il à refuser l'HP3 ? Par ailleurs, les données neurophysiologiques confirment la minimisation de la redondance mais en même temps, elles montrent qu'il existe une grande diversité de modalités : spécialisation corticale pour la couleur, l'intensité lumineuse, le son, mais également la spécialisation corticale multimodale (Berthoz, 2000). Autrement dit, le processus de la minimisation de la redondance semble être lié à une structure qui multiplie les points de vue. Pour autant le problème de la recherche d'invariance se réduit-il à la fusion de données hétérogènes et à leurs influences mutuelles ?

L'étude de ces questions reste emprisonnée dans un cadre qui considère que seuls les capteurs sensoriels (sans prendre en compte les hypothèses sur les distributions de probabilité) offrent les informations nécessaires à l'élaboration de la connaissance. Le traitement statistique, ici assimilé au traitement cognitif, est passif. La légitimité de ce paradigme semble se justifier pour le bas niveau mais qu'en est-il pour l'appréhension de l'espace ou des objets ? La temporalité de l'information se trouve intimement liée à l'activité du système ; or l'épistémologie qui façonne la définition du connexionnisme considère que la subjectivité de la connaissance prend uniquement sa source dans la notion de mesure, indépendamment de l'activité. Pour le connexionnisme, l'individu optimise ses estimations internes et ses classifications uniquement sur la base de critères absolus et fixes. L'intégration de l'activité, avec la proprioception, introduit une notion d'évolution dynamique des critères de convergence qui participe aux soubassements de l'approche interactionniste développée dans la partie traitant les descriptions de la construction d'un individu.

vi - La différence entre le connexionnisme et le subsymbolisme

Avant de conclure sur le connexionnisme, il faut souligner sa forte ressemblance avec l'architecture cognitive très hiérarchisée et modulaire proposée par le subsymbolisme. Cette ressemblance structurelle induit une grande perméabilité entre les deux approches. Néanmoins, les différences relevées dans la section concernant le neuromimétisme conservent leur pertinence, et plus encore, la notion d'auto-organisation proposée par le connexionnisme met en évidence la confusion que fait le symbolisme entre émergence et synergie. En effet, le subsymbolisme affirme que le concept cognitif précède le processus ou la fonction, alors que c'est l'inverse pour le connexionnisme où ce sont les lois d'organisation, d'optimisation, qui déterminent le processus ou la fonction puis le concept cognitif qui en est issu. La notion d'émergence possède la notion d'inattendu, et la notion de synergie possède la notion de coordination, autrement dit, la finalité se trouve pour l'un dans l'optimisation du critère de convergence de l'organisation, et pour l'autre dans l'approximation du concept visé.

**

En somme, pour la robotique, le connexionnisme, en se focalisant sur des théories basées sur des critères de convergence pour expliquer certains phénomènes, offre des

méthodes pour trouver des solutions provenant d'une certaine classe de problèmes qui restent à être inventés dans la pratique, c'est-à-dire définir un environnement propice (respectant les hypothèses statistiques) et une tâche associée (réalisable). Il est à noter que la problématique se retrouve inversée : les solutions préexistent au problème, sans remarquer la pertinence de ce principe dans la cognition elle-même qui sera approfondie dans le second chapitre. Par ailleurs, l'autonomie est ici identifiée à l'auto-organisation, ce qui revient à confondre apprentissage et optimisation. Les traitements sur les données des capteurs adaptés selon la statistique de l'environnement optimisent la relation sensation/effecteur de tous les animaux réactifs comme la grenouille, mais n'offrent pas de piste solide pour l'abstraction.

Avant d'aborder le second mouvement de la robotique cognitive, la description de la construction d'un individu, un tableau récapitulatif des diverses positions issues des approches éliminativistes vis-à-vis des hypothèses et des méthodes constituant la grille d'analyse peut être avancé (Tableau I-3).

		Types de famille métaphysique				Hypothèses sur la perception			Hypothèses sur la conception			Types de méthode d'investigation			Types de description		
		FC	FCP	FP	FA	HP1	HP2	HP3	HC1	HC2	HC3	MHD	MC	MA	DP	DM	DCO
Neuromimétisme	Matérialisme radical	Red	Red	Green	Red	Grey	Grey	Grey	Grey	Grey	Grey	Green	Red	Yellow	Green	Green	Yellow
	Matérialisme rationnel	Red	Green	Red	Red	Grey	Grey	Grey	Grey	Grey	Grey	Green	Yellow	Yellow	Yellow	Green	Green
Connexionnisme	Automate cellulaire	Green	Red	Yellow	Red	Grey	Grey	Grey	Grey	Grey	Grey	Green	Yellow	Green	Yellow	Green	Green
	Intelligence collective	Yellow	Red	Green	Red	Grey	Grey	Grey	Grey	Grey	Grey	Green	Yellow	Green	Yellow	Green	Green
	Auto-organisation neuronale	Yellow	Red	Green	Red	Grey	Grey	Grey	Grey	Grey	Grey	Green	Yellow	Green	Yellow	Green	Green
	Système de codage	Green	Red	Yellow	Red	Grey	Grey	Grey	Grey	Grey	Grey	Green	Green	Green	Red	Green	Red
	Théorie de l'information	Green	Red	Red	Red	Grey	Grey	Grey	Grey	Grey	Grey	Green	Green	Green	Red	Green	Red

Tableau I-3 : Récapitulatif de diverses positions relatives à la grille d'analyse des approches provenant du neuromimétisme et du connexionnisme. Les cases vertes correspondent à un avis favorable, les cases jaunes à un avis mitigé et les cases rouges à un avis défavorable. Les cases grises signifient qu'une prise de position se trouve hors de propos dans le cadre du paradigme éliminativiste.

3.2. La description de la construction d'un individu

Se concentrer sur la description de l'individu en soi autorise à dissocier la connaissance du système cognitif (que sait le système ?) et la connaissance sur le système cognitif (comment fonctionne le système ?). L'exploration des approches appartenant à ce cadre de pensée a contribué à l'amélioration ou à la construction de nombreuses notions telles que la mécanisation des raisonnements déductifs inductifs et pseudo abductifs ainsi que celle d'auto-organisation. Mais la distinction ontologique entre l'individu et le monde conduisant à la distinction entre connaissance sur le monde et connaissance de l'individu introduit également la disjonction entre facultés innées et facultés acquises, ce qui génère une tension pour les concepteurs souhaitant réaliser un robot autonome. En effet, l'acquisition du *bon* réglage ou le choix de la *bonne* clause logique s'effectue selon des critères ou des mécanismes préétablis et déterminés par les objectifs du concepteur ; or, les objectifs d'un système autonome ne doivent pas provenir directement de celui-ci.

Afin de réduire cette tension, la perspective écologique propose de se concentrer davantage sur la notion d'interaction entre l'individu et l'environnement. Simon (1969)

illustre ce changement de perspective en reprenant l'exemple de la tortue mécanique à trois roues du psychologue behavioriste Walter (1950). Chacun des deux capteurs photosensibles de la tortue mécanique, situés à l'avant du véhicule, commande directement une des deux roues motrices arrière. Selon la configuration du câblage choisie entre capteur et moteur, le robot aura un comportement photophobe ou photophile. Vu de l'extérieur, le déplacement du robot peut paraître complexe mais il ne dépend en fait que de la configuration de l'environnement et du câblage. Toutefois, l'examen de la tortue en dehors de toute connaissance sur le monde ne peut prévoir la trajectoire de celle-ci et la seule connaissance de l'environnement ne le permet pas non plus. La complexité naît de l'interaction entre l'intérieur et l'extérieur du système, frontière difficile à cerner et qui dépend essentiellement de l'observateur. Si la cognition correspond à une réponse adaptative pour mieux développer et gérer à son avantage ses interactions avec l'environnement, alors une conception directe d'une machine autonome est impossible du fait que la frontière interactive perçue sera celle du concepteur et non celle du système. Ici, la subjectivité n'est plus seulement un point de vue du sujet dépendant de la nature de ses capteurs sensoriels mais un vécu indissociable de son être au monde ; autrement dit, le sujet-dans-le-monde (l'individu se distingue du monde et le regarde) devient le sujet-par-le-monde (l'individu se façonne par et avec le monde). Moins qu'une métamorphose conceptuelle, la notion de subjectivité se complète inévitablement par la notion de relativité : elles forment les faces d'une même pièce.

Ainsi, la cognition comprise comme un processus gérant, participant et développant une interaction permanente nécessite un corps immergé dans le monde, faisant partie du monde de manière immanente et indéfectible. Cela oblige, contrairement à d'autres approches, à l'utilisation de la robotique dans la recherche de la cognition artificielle (Brook, 1990). L'interaction qui cimente l'individu au monde se traduit par un échange permanent justifié par les besoins dynamiques que réclame le maintien de ce corps. Cet ancrage dans le réel de la pensée a pour incidence de rajouter une nouvelle source de limitation à celle-ci pour appréhender ce réel. Les tentatives de décrire un individu ont mis en évidence deux limitations intrinsèques à la rationalité du monde par un individu. La première limitation concerne à la fois l'aspect pratique (la puissance de calcul pour représenter le monde dans un système logique) et l'aspect théorique (la capacité à fonder un système logique complexe cohérent). La deuxième limitation est celle liée à la subjectivité issue des capteurs et des concepts de description. Le changement de perspective introduit une troisième limitation : les interactions dont les intérêts associés biaisent la perception du monde. Ainsi, la valeur d'une information ne dépend plus seulement des caractéristiques environnementales. En prenant en compte les deux premières limitations, l'individu immergé dans le monde se trouve dépendant d'intérêts très divers (ainsi que des activités associées) liés à la subsistance et au développement, mais ceux-ci peuvent se présenter sous plusieurs formes, de telle façon qu'aucune conceptualisation ou modélisation globale qui aurait autorisé à imaginer une solution optimale ne soit possible. Prix Nobel d'économie, Simon (1969) illustre cette situation en expliquant qu'en fonction des représentations (ou idées) du phénomène économique, ni les problématiques abordées, ni les méthodes de résolution ne sont identiques. Toutefois, certaines problématiques peuvent se recouvrir : sortir d'une crise économique par exemple. Dans ce cas, les solutions issues des différentes visions économiques peuvent être contradictoires, laissant dans l'impossibilité de prendre une décision de manière rationnelle. Il n'y a pas un système unique pouvant éclairer l'ensemble des problématiques économiques, et tenter de se limiter à un seul système serait une grave erreur. Deux systèmes ne sont comparables que sur des problématiques communes, cependant cela ne suffit pas à réfuter entièrement une vision économique par

rapport à une autre. Une théorie ne subsiste que grâce à son efficacité sur ses propres problématiques.

La description de l'individu débouche obligatoirement sur la question insoluble du dualisme, cristallisée par le choix entre « Pierre regarde dans le placard parce qu'il cherche un tire-bouchon » et « Pierre regarde dans le placard parce que la configuration de l'ensemble de son système nerveux l'y a conduit », ainsi que sur la difficulté à concevoir un individu autonome. En se focalisant sur la notion d'interaction comme véritable expressivité de l'individu et source de son développement, le choix entre les deux énoncés perd son sens et les difficultés du concepteur s'évanouissent ou plutôt se décalent. Principalement, trois courants proposent un cadre théorique pour orienter ce décalage (Figure I-26) : le fonctionnalisme écologique, l'évolutionnisme et l'interactionnisme.

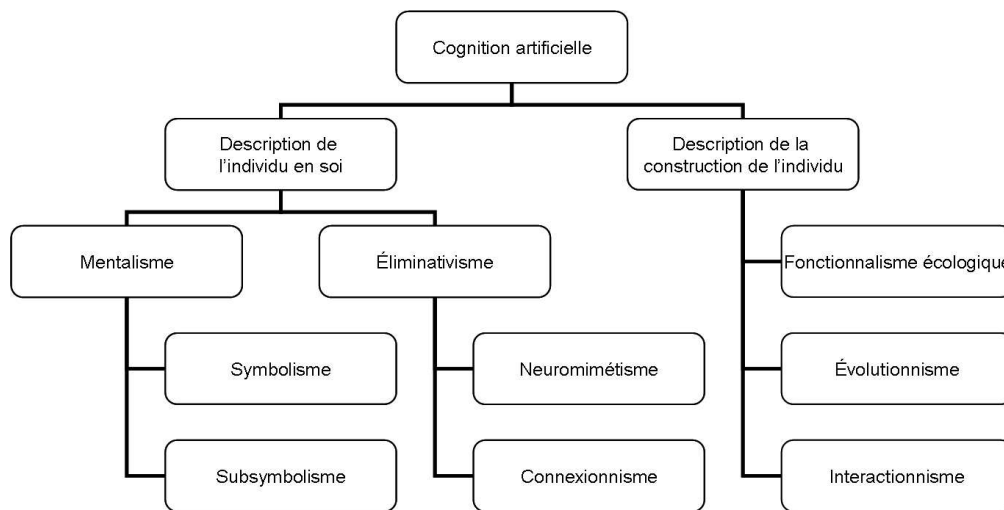


Figure I-26 : Les courants issus du mouvement de la description de la construction au sein de l'arborescence de la cognition artificielle : le fonctionnalisme écologique, l'évolutionnisme et l'interactionnisme.

3.2.1. Le fonctionnalisme écologique

Le fonctionnalisme écologique exploite plus particulièrement la notion de besoin, jointe à celles de plaisir, de douleur, de désir et de répulsion. Ce choix se justifie parce que ce point permet à lui seul de dépasser l'opposition entre mentalisme et éliminativisme. Cette solution n'impose pas pour autant une doctrine unique, elle invite ces approches à une nouvelle coopération, entraînant une pluralité de positions et de techniques associées ou tolérées. Cette section évaluera une tendance moyenne en esquissant dans un premier temps (A) les raisons du dépassement selon chacune des approches, puis dans un second temps en proposant une nouvelle définition de la cognition. Cette dernière amènera à considérer l'incidence d'une position écologique sur la notion de croyance en trois points : (B) l'interprétation de sa crédibilité ou probabilité, (C) son rôle entretenu entre besoin, désir et action puis (D) sa valeur avec le principe de Ramsey (1978). La déduction des conséquences de ces trois points sur les diverses hypothèses et méthodes servira de fil directeur à ce chapitre. Ensuite, les deux principaux axes de recherche des sciences de l'artificiel qui s'inscrivent dans le fonctionnalisme écologique seront présentés : (E) architecture pour la décision et (F) apprentissage par renforcement. Enfin, les difficultés intrinsèques de cette entreprise seront dégagées.

A - La notion de besoin contre le mentalisme et l'éliminativisme

Le premier courant de la description de l'individu en soi, le mentalisme, considère la cognition comme la capacité à manipuler des symboles représentant l'enchaînement logique des concepts. Ces derniers se justifient soit par leur correspondance avec les objets du monde, comme l'énoncé E1 : « sur la table se trouvent trois pommes », soit par eux-mêmes, autrement dit par leur nécessité à exister de fait, ontologiquement. Cela concerne les concepts purs comme la justice, les mathématiques ou le concept d'intentionnalité. Cette autojustification pose le problème de la nature du sens qui entraîne la fracture entre la raison intentionnelle (« Pierre regarde dans le placard parce qu'il cherche un tire-bouchon ») et l'explication matérielle (« Pierre regarde dans le placard parce que la configuration de l'ensemble de son système nerveux l'y a conduit »). L'introduction de l'interaction comme nécessité à subsister, donc de besoin, propose une autre justification à cette seconde catégorie de concepts, une utilité. Cette position rejoint la doctrine utilitariste de Bentham (1811) ou de Mill (1861) qui assimile le comportement humain à la tentative de maximiser les plaisirs et de minimiser les souffrances. Cette maximisation et cette minimisation portent sur le cumul tout le long de la vie du sujet, il y a donc intérêt à discerner les plaisirs brefs qui peuvent ensuite entraîner de plus grandes souffrances et les brèves souffrances pouvant déboucher sur une satisfaction durable. Dans cette logique, l'abstraction sert à éviter ces pièges et à chercher des solutions pérennes. Ainsi, aucun concept ne se trouve détaché du monde, exceptés peut-être les concepts mathématiques. Le concept de justice n'est plus un concept pur préexistant à toute activité humaine qui trouve écho dans la société des hommes comme l'aurait imaginé Platon. Le concept de justice se comprend comme un concept élaboré par les sociétés humaines pour maximiser le bonheur (ou minimiser les malheurs produit par les conflits) de la société, qui aura un impact généralement positif sur chaque individu. Il apparaît alors possible de hiérarchiser les besoins, ceux qui sont directs et instantanés pour l'individu et ceux qui se révèlent indirects et liés à la condition globale de l'individu.

Le second courant de la description de l'individu en soi, l'éliminativisme, fonde l'explication de la cognition d'un individu sur l'étude des lois de la matière (le neuromimétisme) ou sur l'étude des lois de l'organisation (le connexionnisme) qui régissent les éléments le constituant. Ainsi, le neuromimétisme dévoile la manière dont les mécanismes cérébraux participent à la cognition mais n'explique pas le principe de la cognition, son rôle. Le fonctionnalisme écologique affirme que ces mécanismes sont orientés et adaptés pour un environnement précis, une niche écologique où leur survie se trouve favorisée et que cette dimension est indispensable pour comprendre leur intérêt et par conséquent leur sens. Le connexionnisme, en s'appuyant sur la minimisation de la redondance comme seul principe directeur à la cognition, se restreint aux traitements de bas niveau pour développer une catégorisation automatique des objets du monde. Ici, la notion de besoin, incluse dans le fonctionnalisme écologique, apporte un nouveau critère pour orienter la construction de catégories et de contextes.

Dans ce contexte, une définition unificatrice de la cognition serait la suivante : la cognition correspond à un processus permettant de se constituer des croyances utilisées dans l'analyse des besoins pour la définition des buts et dans le choix de l'action par rapport à ces derniers. Deux remarques viennent compléter cette définition. D'une part, le terme croyance rentre directement et pleinement en conflit avec l'éliminativisme. Toutefois, le fonctionnalisme écologique propose ici un accord en réduisant la croyance à trois critères : (1) la croyance est un élément participant à la justification d'une action (Peirce, 1878), (2) la crédibilité d'une croyance est évaluable par rapport à une autre et (3) la

crédibilité aussi bien que la croyance elle-même peuvent évoluer. Ainsi définie, la croyance devient suffisamment précise et souple pour considérer comme telle une assemblée de neurones codant la distribution de probabilité d'un motif appris. D'autre part, l'action plus que le besoin introduit une composante normative supplémentaire. En effet, en mettant de côté le principe de la mise en correspondance, un système de connaissances se construirait passivement sur la prédiction des observations et sur la cohérence des concepts utilisés ; ici, la réussite des actions participe également à la construction des connaissances et les uniformise si tous les agents agissent de façon similaire.

Ces deux remarques invitent à approfondir trois questions qui seront successivement abordées : Quelle est la signification de cette crédibilité comprise comme probabilité, et en quoi diffère-t-elle des significations implicitement employées dans les sections précédentes ? Pourquoi une croyance semble-t-elle nécessaire à l'action et pourquoi doit-elle être comprise différemment d'une représentation explicite du monde comme le préconise le symbolisme ? Comment évaluer la vérité d'une croyance, la réussite de l'action est-elle suffisante ?

B - Quelle est la signification de cette crédibilité comprise comme probabilité et en quoi diffère-t-elle des significations implicitement employées dans les sections précédentes ?

La réponse à la première interrogation passe par l'analyse des interprétations existantes au sujet de la probabilité. La notion de probabilité traduit le degré de crédibilité ou de vraisemblance concernant soit une connaissance soit une prédiction. Cette distinction entre les deux types d'objets de la probabilité sera davantage discutée dans la réponse à la troisième interrogation. Selon le mode d'évaluation, la probabilité peut être soit objective soit subjective. La probabilité objective résulte d'une méthodologie qui ne fait pas intervenir l'opinion de l'observateur sur le résultat. Deux méthodologies existent : la probabilité fréquentielle ou empirique et la probabilité formelle ou logique. La probabilité empirique se fonde sur la fréquence d'apparition d'un événement soit effectif (lancer une pièce de monnaie), soit informatif (sondage d'opinion). Plus précisément, la probabilité d'apparition d'un événement revient au quotient du nombre d'occurrences de l'événement par le nombre d'épreuves lorsque ce dernier tend vers l'infini, ou du moins devient suffisamment grand pour que la fréquence devienne stable dans l'intervalle de précision considéré. La seconde méthode, la probabilité logique, présuppose un modèle du monde, d'un univers qui représente l'ensemble de tous les résultats possibles de l'expérience. Le rapport entre le nombre d'occurrences d'un résultat ou d'un ensemble de résultats et le nombre total donne la probabilité de l'apparition de ce résultat ou de cet ensemble de résultats. Par exemple, la probabilité de tirer une boule rouge dans une urne qui en contient dix dont deux rouges est de $2/10$.

La probabilité subjective correspond à l'impossibilité d'employer strictement l'une des deux méthodes précédemment décrites. Par exemple, s'il est demandé d'évaluer la probabilité d'obtenir trois avec un jet de dés, la première réponse est a priori $1/6$. Après avoir précisé que celui-ci était pipé, la meilleure réponse reste $1/6$ mais cette fois la réponse ne prétend plus être exacte. Elle approxime l'erreur minimale de la probabilité objective inconnue. De même, avant d'avoir effectué un grand nombre d'essais, la probabilité fréquentielle reste subjective à son échantillonnage, que le dé soit pipé ou non. Les probabilités subjectives sont des probabilités imprécises. Les probabilités objectives modélisent l'incertitude sur le monde mais sans douter de leur propre modélisation ; autrement dit, elles ne tiennent pas compte de leur propre ignorance. Il existe de multiples sources d'ignorance et des taxinomies de celles-ci sont proposées (Smithson, 1989 ; Ha-

Duong, 2005). La reconnaissance de l'ignorance devient une source de connaissance utile pour la compenser, à défaut de la combler. Du point de vue de la théorie de la probabilité logique, la compensation revient à introduire des distributions de probabilité a priori. Toutefois, plusieurs théories sortent du cadre trop restreint des probabilités classiques en proposant d'autres concepts afin d'explorer certaines facettes de l'ignorance : la possibilité, la plausibilité, le flou, etc. Néanmoins, cette section conservera autant que possible le cadre théorique probabiliste tout en acceptant l'introduction de distributions à priori. En effet, les autres théories évoquées peinent à trouver une cohérence bien que des efforts notables aillent dans ce sens (Ha-Duong, 2005). Par ailleurs, la théorie classique des probabilités permet de formuler sobrement la relation entre croyance ou hypothèse, probabilité et connaissance : $P(h,c)$ où P représente une fonction de probabilité, h une hypothèse et c l'ensemble des informations établies. Les quatre interprétations associées à cette expression proposées ici s'inspirent de l'analyse de Schilpp (1963).

i - $P(h,c)$ comme une approximation d'une probabilité existante

La première interprétation, d'ordre statique, comprend $P(h,c)$ comme une approximation d'une probabilité existante, autrement dit objective, en dehors de toute sémantique liée à la causalité de h et de c . Plus précisément, cette interprétation se décline de manière différente selon qu'elle traduit une probabilité sur un événement ou un état du monde, chacune d'elle mettant en exergue un point faible de cette interprétation, justifiant une présentation plus approfondie. Selon l'interprétation statistique d'un événement, la prédiction s'appuie sur une probabilité subjective $P(h)$ avant toute expérience. Les expériences successives ne modifient pas la fonction de probabilité primitive en $P'(h)$ mais elles mettent à jour continuellement son argument par h sachant c qui représente les résultats des expériences précédentes (De Finetti, 1937). Cette interprétation assure la cohérence et la continuité dans l'élaboration d'un système de connaissance. Pour le jeu de pile ou face, les séquences d'événements de longueurs équivalentes ayant la même proportion de tirage pile et face possèdent une probabilité identique. Alors, la probabilité d'obtenir face F_{n+1} à la suite de n lancer sachant c soit « il y a eu r tirages face précédemment », en faisant tendre n vers l'infini devient :

$$P(F_{n+1} / c) \approx \frac{r}{n} .$$

Dans le cas où la probabilité qu'il y aura r faces dans n épreuves est la même quelle que soit la valeur de r , la règle de succession s'utilise directement :

$$P(F_{n+1} / c) = \frac{r+1}{n+2} .$$

Toutefois, ce type de raisonnement inductif, prôné par le positivisme logique évoqué dans la section sur le connexionnisme ne permet pas d'intégrer de situation nouvelle ou plus précisément des ruptures. Popper (1934) illustre cette lacune en calculant la probabilité que le soleil se lève demain. En prétendant que l'histoire à 5000 ans et qu'aucune trace n'indique l'absence d'un levé de soleil alors la probabilité que celui-ci se lève demain est de 0,999 999 4 (1 représentant la certitude). Si le soleil venait à ne pas se lever demain, la probabilité qu'il se lève le surlendemain restera de l'ordre de 0,999 998 9.

L'interprétation statistique d'un état du monde s'avère conflictuelle avec une approche strictement fréquentielle lorsque le recensement exhaustif se trouve exclu. En effet, toute déduction ou inférence s'appuyant sur un sondage contient une part de

subjectivité. Par exemple, soit h l'évaluation du nombre de partisans pour le « oui » avant un référendum sur la base d'un sondage c dépendant selon le théorème de Bayes des probabilités a priori $P(h)$ et $P(c)$:

$$P(h/c) = \frac{P(h)P(c/h)}{P(c)}$$

En appliquant le théorème des probabilités totales, la formule devient avec l'indice n qui indique le nombre de partisans au « oui » et N qui représente le nombre d'individus composant la population soumis au référendum :

$$P(h_n/c) = \frac{P(h_n)P(c/h_n)}{\sum_{k=1}^N P(c/h_k)P(h_k)}$$

Le terme $P(c/h_n)$ appelé également fonction de vraisemblance de h_n étant calculable, l'estimation de $P(h_n/c)$ dépend alors de la distribution a priori de $P(h)$ sur n qui peut cependant toujours être réajustée. Les deux sous-interprétations restent fidèles à l'interprétation commune de $P(h/c)$ qui demeure l'estimation de la fréquence relative d'une propriété X dans une population donnée formée d'éléments non inclus dans c . Cette interprétation apparaît au sein des familles FC et FCP.

ii - $P(h/c)$ est comprise comme une mesure de confirmation

La seconde interprétation de l'expression $P(h/c)$ la comprend comme une mesure de confirmation selon laquelle la croyance de h sur la base de c , ou comme le degré de corrélation entre h et c . En assimilant ces corrélations à des liens causaux, le théorème de Bayes (1763) permet à partir de l'observation de remonter aux « causes » ou du moins d'estimer leur probabilité. Par exemple, un tirage s'effectue de façon équiprobable entre deux urnes qui ne contiennent que des boules noires ou blanches. La répartition de la première se traduit par une probabilité de 75% de tirer une boule blanche. La seconde urne, pour le même événement, possède une probabilité de 50%. La probabilité de l'hypothèse h_1 que la boule provienne de la première urne égale la probabilité de celle de h_2 que la boule provienne de la seconde, soit $P(h_1) = P(h_2) = 50\%$. Toutefois, l'information c du premier tirage qui est « une boule blanche » permet de préciser la probabilité concernant son origine, $P(h_1/c) = 60\%$ puisque :

$$P(h_1/c) = \frac{P(h_1)P(c/h_1)}{P(h_1)P(c/h_1) + P(h_2)P(c/h_2)}$$

En considérant les états comme des nœuds et les corrélations comme des liens se tisse alors un réseau d'inférence, appelé aussi réseau bayésien. Cette interprétation de l'expression $P(h/c)$ trouve particulièrement écho dans deux approches. Tout d'abord, le subsymbolique qui considère la possibilité de réaliser n'importe quelle inférence au sein de ce réseau comme le schéma d'un raisonnement cognitif. Une mesure entre h et c reflète la pertinence entre l'explanandum (ce qu'il y a à expliquer) et l'explanans (ce qui explique). Plus l'explanans s'applique de façon décisive à des situations variées, plus il aura tendance à être élevé au rang de loi universelle. L'autre approche, le connexionnisme, utilise également cette interprétation bien qu'elle se fonde sur la théorie de l'information. En effet, la minimisation de la redondance oblige à trouver la meilleure représentation des données et conduit à estimer indirectement leurs distributions. Or, les réseaux bayésiens visent

explicitement à trouver les estimateurs des probabilités a priori correspondant le mieux aux données. Ces deux techniques sont donc équivalentes (Nadal, 1997). La définition de l'information mutuelle s'identifie à l'interprétation de $P(h/c)$ comme une mesure de corrélation entre h et c . Mais au-delà du délicat problème de trouver le bon estimateur, il demeure le problème de la construction de la topologie du réseau. Dans le cas du connexionnisme, les capteurs s'identifient directement à des nœuds, et cette identification résout le problème de la détermination des nœuds. Mais les réseaux dont les nœuds représentent des objets ou des qualités doivent tous être définis par avance. Dans ce cas l'évolution « cognitive » se limite à la modulation de liens entre les nœuds déjà existants. Le problème se complique également lorsque la temporalité est prise en compte : Quelle granularité choisir ? Quelle profondeur mnémonique retenir ? Par ailleurs, la prise en compte de tous les liens potentiels nécessite un grand nombre de calculs et introduit du bruit ; il est donc conseillé de supprimer les liens entre les nœuds qui sont a priori indépendants.

iii - $P(h/c)$ comme le quotient d'un pari juste

La troisième interprétation, à la différence des deux précédentes, rentre directement dans le cadre du fonctionnalisme écologique du fait que l'expression $P(h/c)$ devient le quotient d'un pari juste par rapport à h sachant c et par conséquent introduit la notion de valeur d'échange qui peut-être utile pour une interaction avec l'environnement. Le principe du pari s'applique traditionnellement dans la situation où Paul promet à Jean une somme S_1 si h ne se réalise pas. Dans le cas inverse, Jean promet de donner une somme S_2 à Paul. Le quotient du pari relatif à h se traduit par S_1 divisé par la somme des mises. Le quotient devient juste lorsque les deux parties profitent équitablement du pari :

$$P(h/c) = \frac{S_1}{S_2 + S_1}$$

Cette interprétation de l'expression $P(h/c)$ exprime la notion de risque à entreprendre une certaine action. Par exemple, Paul connaît plusieurs baies de couleur rouge comestibles mais il sait aussi que certaines provoquent des douleurs à l'estomac ; Paul va-t-il prendre le risque de goûter une baie qu'il ne connaît pas ? Il met en jeu un mal de ventre contre le plaisir gustatif. Dans cette situation, ce sont les conséquences du pari qui vont moduler les estimations a priori du théorème de Bayes.

iv - $P(h/c)$ comme la crédibilité de h sachant c en fonction de l'utilité

La quatrième et dernière interprétation se trouve en totale adéquation avec l'approche du fonctionnalisme écologique, au sens où $P(h/c)$ exprime la crédibilité de h sachant c en fonction de l'utilité que peut représenter h . Par rapport au pari qui exprime le risque d'une seule action, ici la notion d'utilité permet d'exprimer le risque à agir d'une manière parmi d'autres et en fonction de ses propres besoins. Cette interprétation très présente au sein des théories économiques peut être illustrée par le problème du marchand de glace (Ha-Duong, 2005) qui montre l'impossibilité de définir un agent rationnel idéal ou autrement dit l'impossibilité de définir un agent en dehors d'une perspective écologique. Un marchand de glace doit choisir entre quatre emplacements $E \in \{e_1, e_2, e_3, e_4\}$. Le profit prévisionnel G varie selon l'emplacement et selon la météo $M \in \{\text{Ensoleillé, Pluvieux}\}$ (Tableau I-4).

		Emplacement			
		e_1	e_2	e_3	e_4
Météo	Ensoleillé	10	6	11	8
	Pluvieux	2	4	0	3
Espérance		6	5	5,5	5,5

Tableau I-4 : Profit prévisionnel $G(E,M)$ en k€ et espérance associée lorsque les annonces météorologiques sont équiprobables (Duong, 2004).

Un marchand de glace optimiste pensera que la place offrant un maximum de gain est la plus opportune, alors qu'un marchand pessimiste choisira au contraire la place maximisant ses gains même lorsque le temps s'assombrit. Un marchand de glace ordinaire ira, dans l'ignorance des pronostics météorologiques, à l'emplacement maximisant l'espérance mathématique des gains. Cependant, plusieurs stratégies de décision rationnelle se révèlent possibles lors d'une situation incertaine. Le choix de la stratégie ne semble pas être relié uniquement au gain monétaire. En effet, un marchand de glace devant rembourser impérativement une dette à court terme sous peine de lourds désagréments s'arrangera pour s'installer à l'emplacement e_2 . La théorie de l'utilité propose un cadre unificateur décrivant ces diverses stratégies en introduisant une fonction d'utilité U qui se comprend comme une évaluation de la nécessité des gains par rapport au besoin ou une évaluation du choix minimisant le risque d'obtention des gains par rapport à l'importance des besoins. L'emploi de la fonction d'utilité entraîne l'interprétation hors du cadre de la théorie des probabilités, $P(h,c)$ devenant une mesure de pertinence liée au choix h en sachant c et en fonction de U . Avec C l'ensemble des degrés de croyance dans la réalisation des différents événements et h_i représentant l'hypothèse de choisir e_i , l'expression des pronostics météorologiques devient ici :

$$P(h_1, C) = \sum_{n=1}^2 U(G(h_1, c_n)) * P(c_n)$$

La décision de l'action s'effectue en prenant l'hypothèse dont la valeur de probabilité subjective est maximale :

$$Action = \arg \max_{i \in N} P(h_i, C)$$

Pour le marchand de glace, quatre fonctions recouvrent les comportements évoqués : $U=G$, $U=\ln(G)$, $U=-1/G$, $U=G^2$. Les résultats montrent des choix très disparates lorsque l'incertitude est maximale puis ils tendent à s'uniformiser lorsque la certitude augmente, c'est-à-dire lorsque la notion de risque s'efface et que la notion de maximisation des gains devient le seul critère (Tableau I-5). Si le marchand a la certitude que le temps sera ensoleillé, il choisira l'emplacement e_3 quelle que soit la fonction d'utilité adoptée et a contrario ce sera l'emplacement e_2 qui sera désigné.

		Emplacement			
Annonce Météo	Fonction d'utilité	e1	e2	e3	e4
P(Ensoleillé)=0% P(Pluvieux)=100%	G	2	4	0	3
	ln(G)	0,69	1,38	-∞	1,09
	-1/G	-0,50	-0,25	-∞	-0,33
	G ²	4	16	0	9
P(Ensoleillé)=25% P(Pluvieux)=75%	G	4	4,5	2,75	4,25
	ln(G)	1,09	1,48	-∞	1,34
	-1/G	-0,40	-0,23	-∞	-0,28
	G ²	28	21	30,25	22,75
P(Ensoleillé)=50% P(Pluvieux)=50%	G	6	5	5,5	5,5
	ln(G)	1,49	1,58	-∞	1,59
	-1/G	-0,30	-0,20	-∞	-0,23
	G ²	52	26	60,5	36,5
P(Ensoleillé)=75% P(Pluvieux)=25%	G	8	5,5	8,25	6,75
	ln(G)	1,9	1,69	-∞	1,83
	-1/G	-0,20	-0,18	-∞	-0,17
	G ²	76	31	90,75	50,25
P(Ensoleillé)=100% P(Pluvieux)=0%	G	10	6	11	8
	ln(G)	2,30	1,79	2,39	2,07
	-1/G	-0,1	-0,16	-0,09	-0,12
	G ²	100	36	121	64

Tableau I-5 : Pertinence des choix selon les pronostics météorologiques et selon la fonction d'utilité adoptée. Les cases surlignées représentent l'emplacement le plus pertinent suivant la stratégie employée.

Un système reposant sur la théorie de l'utilité se trouve en perpétuelle évolution concernant les estimations des probabilités des événements répertoriés, la modification de la fonction d'utilité suivant l'état des besoins et l'ajustement des profits attendus, tout cela orienté par la pertinence effective ou non de l'action entreprise. Ce dernier point est particulièrement important car il explique le caractère historique et situé des représentations ou croyances que possède l'individu. Pour finir, cette quatrième interprétation propose une explication de la rationalisation limitée de l'agent mais également de la rationalité limitée qu'a celui-ci sur le monde. La théorie des jeux montre que l'équilibre de certaines situations entre plusieurs agents ne se révèle pas être le maximum global de satisfaction (Nash, 1951). Cette situation s'illustre traditionnellement avec le dilemme du prisonnier. Pris avec un comparse, il doit choisir entre le dénoncer ou nier les faits. Les conséquences des choix sur les condamnations, qui dépendent également de l'attitude du comparse, se résument par le tableau suivant :

	P1 nie les faits	P1 dénonce P2
P2 nie les faits	(6 mois, 6 mois)	(relaxé, 10 ans)
P2 dénonce P1	(10 ans, relaxé)	(5 ans, 5 ans)

Tableau I-6 : Peines des prisonniers selon leur déposition respective : (P1, P2).

Dans le doute, la stratégie la plus sûre s'avère être la dénonciation mutuelle, ce qui ne correspond pas au maximum global. Pour un écosystème, le principe d'équilibre de Nash se comprend comme une situation d'interaction stable et viable parmi un ensemble d'agents qui n'ont pas intérêt à changer unilatéralement leurs stratégies.

En conclusion, à la première question (quelle est la signification de cette crédibilité comprise comme probabilité et en quoi diffère-t-elle des significations implicitement employées dans les sections précédentes ?), le développement précédent répond que le fonctionnalisme écologique s'approprie la troisième et plus particulièrement la quatrième interprétation, et que celles-ci diffèrent des deux premières par le fait que ce n'est pas tant la prédiction sur l'état du monde qui guide l'estimation des probabilités que la récompense suite aux actions décidées à partir des probabilités antérieures.

C - Pourquoi une croyance semble-t-elle nécessaire à l'action et pourquoi doit-elle être comprise différemment d'une représentation explicite du monde comme le préconise le symbolisme ?

Répondre à la deuxième interrogation oblige à définir les termes besoin, plaisir, douleur, désir et répulsion. Le besoin d'un agent résulte d'un manque de quelque chose qu'il doit s'approprier pour compléter son être. Le plaisir s'assimile à une rétribution positive signalant la satisfaction d'un besoin ou signalant que la situation est potentiellement intéressante pour satisfaire un besoin qui est pour l'instant comblé. La douleur s'assimile à la rétribution négative signalant que la situation génère anormalement des besoins ou que les besoins atteignent un seuil critique. Le désir représente la propension à aller vers un objet considéré comme bénéfique c'est-à-dire source de plaisir et indirectement capable de satisfaire un besoin. À l'inverse, la répulsion traduit la crainte d'un objet potentiellement dangereux, autrement dit, pouvant soit être douloureux soit indirectement générer des besoins en dehors de la dynamique normale du système. Le plaisir et la douleur restent difficiles à définir sans recourir à l'utilisation du vocabulaire relatif aux sensations propres, toutefois, un sens supplémentaire est introduit dans le cadre du fonctionnalisme écologique.

En admettant que le plaisir et la douleur incarnent des directions opposées sur une même dimension sensitive, la troisième interprétation de l'expression $P(h,c)$ suffit à harmoniser un système et son environnement. Le pari est que l'objet désiré procure du plaisir. L'objet désiré se traduit plus généralement par un stimulus désiré. En effet, un système extrait seulement la quantité d'information suffisante pour assurer la bonne action. Par exemple, pour capturer les mouches, le mécanisme visuel de la grenouille associe un déplacement rapide d'un point à une certaine distance avec l'action consistant à déployer sa langue vers celui-ci, le « concept » de mouche s'arrête là. Concernant le plaisir, celui-ci permet un mécanisme d'apprentissage. Le renforcement par épreuves successives oriente l'apprentissage de la forme de ce qui doit être désiré. Ce mécanisme assure l'adéquation entre le monde et l'individu ainsi qu'entre l'objet du désir issu de l'environnement et le besoin qui résulte de la physiologie de l'individu. Toutefois, beaucoup de systèmes se dispensent d'un mécanisme de renforcement, l'adéquation entre l'action et le stimulus résultant alors de la phylogénèse, comme pour la grenouille. Mais le mécanisme par renforcement possède l'avantage de découvrir de nouvelles sources de satisfaction. Jusqu'ici, le besoin existe sur le plan physiologique sans apparaître sur le plan cognitif, c'est-à-dire dans une boucle décisionnelle, la récompense guidant implicitement sa satisfaction. Ainsi dans le cadre de l'apprentissage par renforcement, les probabilités et leurs estimations (les croyances) visent des profits prévisionnels sans avoir nécessité de spécifier besoin et état du monde.

Cependant, comme le montre la quatrième interprétation de l'expression $P(h,c)$, des situations plus complexes poussent à coupler la relation de pari avec l'utilité des gains. Afin d'identifier l'origine de cette complexification, le chemin d'analyse sera renversé en reprenant l'exemple du marchand de glace : quel paramètre de l'énoncé doit-il être modifié pour revenir à une simple situation de pari ? Deux solutions se révèlent envisageables. La première correspond à la situation où il y a une totale certitude sur le pronostic météorologique. Le pari porte alors sur l'estimation du profit maximum. La seconde solution propose de réduire le nombre d'emplacements à un. Dans ce cas, le pari reprend du sens si un risque existe que la recette ne compense pas les coûts d'une ouverture. Dans aucune de ces deux situations, le besoin se révèle indispensable, ce qui signifie que la notion de besoin intervient lorsque le champ des possibles augmente en même temps que l'incertitude sur les états des mondes. Il ne faut pas pour autant comprendre que seules la liberté d'action et l'ignorance se trouvent à l'origine de la cognition. Cette situation n'a de raison d'être seulement parce que parallèlement de multiples besoins s'entremêlent dans un jeu de compensation complexe où la satisfaction des uns oblige à entamer celle des autres et *vice versa*. Par ailleurs, les besoins ne se régulent généralement pas linéairement : chasser dépense de l'énergie mais permet d'en reprendre lorsque la proie se trouve prise. En somme, la diversité des actions et l'incertitude révèlent des choix, qui n'apparaissent pas auparavant, relatifs à des besoins déjà existants.

Comme pour la mise en place d'un mécanisme de pari, la notion de croyance n'est pas indispensable pour ce genre de mécanisme décisionnel. Chez la grenouille, deux besoins ne se trouvant pas directement liés au maintien de l'homéostasie renvoient à des comportements antagonistes. Le premier besoin concerne la tendance à fuir toute masse se dirigeant vers l'animal, mécanisme de défense procurant un avantage sélectif. Le second besoin concerne l'accouplement ; or, si le mécanisme de défense reste activé, l'approche du ou de la partenaire provoquera une réaction de fuite. Un mécanisme de régulation hormonale résout ce dilemme en synchronisant sur une courte période une inhibition du comportement de fuite, favorisant ainsi les rencontres entre les partenaires. Le taux d'hormone peut-être considéré comme la valeur d'une fonction d'utilité. D'une certaine manière, l'organisme incarne la croyance de l'espèce mais il n'a pas de croyance au sens défini précédemment, c'est-à-dire que le couple croyance et crédibilité n'est pas révisable dans son ensemble au cours de la vie de l'organisme.

L'imbrication des besoins, le nombre de degrés de liberté et l'incertitude sur l'avenir, ou de ce qui est perçu, obligent le système à décider sur la base de valeurs endogènes et exogènes, mais ne justifie pas à eux seuls une construction de croyances propres relatives au monde. Le fonctionnalisme écologique va alors argumenter en trois points la nécessité pour des systèmes complexes de se servir de croyances, c'est-à-dire d'un modèle interne du monde révisable. Le premier argument développé notamment par Sigaud (2001) vient de la constatation que certaines situations, indépendamment de la qualité sensorielle, apparaissent identiques, bien qu'elles ne le soient pas. Par exemple, le robot muni de capteurs de proximité ne percevra pas de différence entre la position A et la position B du labyrinthe (Figure I-27). Pour distinguer ces deux positions, il convient de posséder des informations supplémentaires ; ces informations mémorisées, constituées, ne correspondent plus à des variables internes relatives à un besoin mais bien à un modèle relatif au monde.

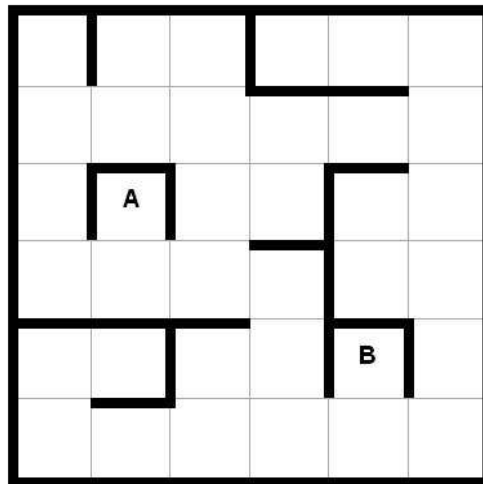


Figure I-27 : Un robot muni de capteur de proximité ayant une portée d'une case ne peut pas distinguer les emplacements A et B à moins de posséder des informations complémentaires.

Le second argument repose sur l'idée que des besoins spécifiques dans un milieu incertain n'autorisent pas le processus phylogénique à constituer un ensemble de réflexes pour les satisfaire. Un individu limité à la contingence d'un grain spatiotemporel peut trouver des solutions pertinentes uniquement au sein de ce grain. Autrement dit, la mise en correspondance des objets désirés, des actions et des besoins doit se réaliser en partie au cours de la vie de l'organisme. Cet argument amorce également l'idée que la planification découle directement de ce genre de mécanisme. Le troisième argument revient sur le principe brièvement abordé précédemment de stratégie optimum au sein d'une population d'agents. Nash (1951) montre que, dans certaines situations, si seule la compétition oriente le comportement d'une population, le point d'équilibre de satisfaction moyen peut être sous-optimal. Dans ce cas, la coopération devient un avantage sélectif mais elle semble obliger à utiliser un modèle d'autrui. Dans un deuxième temps, le faire-semblant possède aussi certainement une efficacité. Le cas des insectes sociaux qui coopèrent par le phénomène de la stigmergie sans l'aide de croyance représente une objection directe contre cet argument. Mais cette objection reste attaquable parce que le besoin ou l'avantage sélectif de coopérer a pu apparaître à un moment différent au cours de la phylogénèse d'autres espèces pour lesquelles les agents sont plus individualisés au départ, et par conséquent le développement phylogénétique sous cette pression a pu s'exprimer autrement. Deux autres intérêts peuvent être avancés pour la prise en compte des croyances sur les intentions d'autrui. Le premier souligne que l'apprentissage par imitation permet de factoriser l'apprentissage et surtout d'identifier quels comportements peuvent être fatals sans se mettre en danger, et d'avertir la communauté. Le second intérêt revendique l'importance du langage dans le développement des concepts ; toutefois il ne peut pas constituer un argument à part entière, le langage pouvant nécessiter la notion de croyance, mais pas l'inverse, bien qu'il puisse jouer le rôle de catalyseur.

Synthétiquement, la réponse à la deuxième interrogation (pourquoi une croyance semble-t-elle nécessaire à l'action, et pourquoi doit-elle être comprise différemment d'une représentation explicite du monde comme le préconise le symbolisme ?) montre que les arguments n'identifient pas une cause précise à l'apparition de la croyance. Ils proposent néanmoins que sa nécessité apparaisse dans des conjonctures particulières regroupant à la fois des besoins, des capacités d'agir, et un environnement ambigu pour le système

concerné. Comparativement au mentalisme, les arguments du fonctionnalisme écologique en faveur des croyances ne préjugent pas de leur nature : iconique, sensorimotrice ou autre. Par ailleurs, le fonctionnalisme écologique considère les croyances comme des éléments complétant les informations sensorielles afin de guider au mieux l'action, et comme des éléments appartenant au monde, contrairement au mentalisme qui considère les croyances comme un monde parallèle, obligeant à avoir des représentations totalisantes iconiques ou géométriques qui peuvent servir éventuellement à simuler une action.

D - Comment évaluer la vérité d'une croyance, la réussite de l'action est-elle suffisante ?

La réponse à la dernière des trois interrogations sur la croyance se construit à partir de l'analyse de Dokic (1999) sur le principe de Ramsey (1978) et de l'action située. Ce principe se fonde sur le raisonnement suivant : une croyance fautive peut nuire à la réussite d'une action, contrairement à une croyance vraie. Autrement dit, la vérité des croyances garantit la réussite des actions. Mais cette dernière reformulation bien qu'elle semble triviale n'aboutit pas à une équivalence logique avec le principe de Ramsey. Selon une première lecture, la croyance E7 : « pour obtenir l'objet désiré D dans un contexte C il faut agir selon A » se trouve validée par la réussite, auquel cas elle est vraie. Cela entraîne l'identification entre réussite et vérité. Cependant, cette définition de la vérité ne porte plus alors uniquement sur les croyances. En effet, en admettant que la réussite se trouve signalée par le plaisir délivré par la satisfaction d'un besoin, la rétribution produite par celle-ci porte à la fois sur l'adéquation entre l'objet désiré avec ce qu'il procure et sur l'action qui a permis de l'acquérir. Il existe alors une possibilité qu'une réussite résulte d'une croyance juste, et d'une action manquée ou inversement. Cette distinction encourage le fonctionnalisme écologique à différencier la notion de vérité, qui concerne uniquement la croyance, et la notion de réussite, qui porte sur la croyance et l'action. La réussite devient une condition de vérité et non plus une définition de la vérité.

Plus précisément, la croyance E7 avec pour indice de réussite la satisfaction d'un besoin se traduit par une croyance implicite supplémentaire qui permet alors de décomposer la situation en deux croyances : une croyance d'orientation E8 « l'objet désiré D procure satisfaction du besoin B » et une croyance instrumentale E7. La réussite, comme la satisfaction d'un besoin, forme une boucle de rétribution qui ne distingue ni E7 ni E8. Cette confusion ne porte pas préjudice dans une logique de construction de schèmes comportementaux basés sur l'habitude. Il importe peu qu'une croyance soit fautive ou que l'action soit manquée (au sens de ne pas saisir l'objet désiré) pourvu que cette dernière génère une satisfaction, une utilité. La survie d'un individu dépend de son opportunisme. Mais lorsque l'individu échoue, la révision des croyances doit-elle être totale ? Un principe d'économie ne serait-il pas de découvrir laquelle d'E7 ou d'E8 est la croyance source d'erreur ? Un tel mécanisme nécessite donc une prédiction sur les conséquences sensorielles de l'action et non plus sur ses conséquences relatives à la modulation des besoins. Ainsi, il devient possible de vérifier que l'objet désiré (identifié comme tel) soit acquis avant d'évaluer son efficacité. Deux boucles de rétribution imbriquées apparaissent comme le montre la Figure I-28. Le besoin B donne une propension à désirer l'objet D car il existe une croyance que D doit satisfaire B ; cette situation génère une attente de satisfaction S' si le sujet repère D. La seconde croyance tient pour vrai que dans le contexte C, l'action A semble la meilleure pour obtenir D, ce qui génère une seconde attente D'. L'attente S' sachant la réalisation de l'attente D' permet d'évaluer la croyance unissant D avec B et l'attente D' permet d'évaluer la croyance portant sur l'efficacité de A.

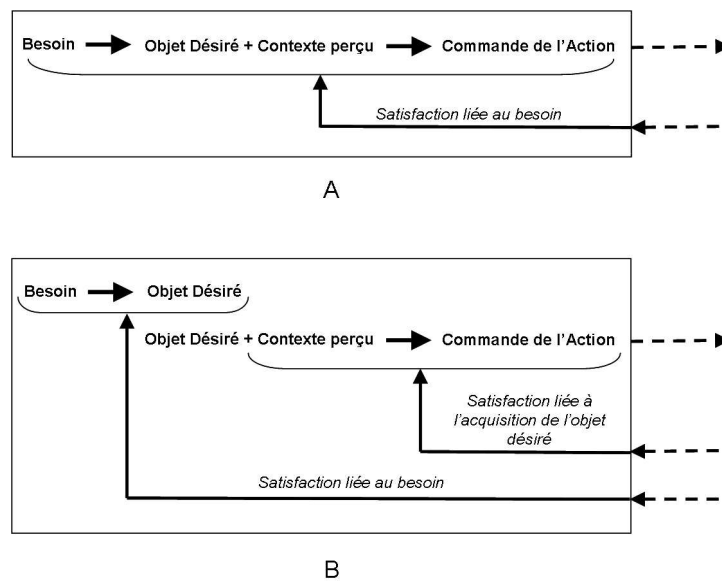


Figure I-28 : Le schéma A représente un système possédant une seule boucle de rétribution, le schéma B introduit une seconde boucle de rétribution permettant de cibler l'évaluation a posteriori des prémisses.

Cette décomposition autorise alors une seconde lecture du principe de Ramsey pour chacune des boucles. L'évaluation de la première, de la croyance E8, s'entreprennent seulement lorsque la réussite de l'action se trouve assurée, autrement dit, ce n'est pas E8 qui est évalué mais c'est E8 sachant E7. Cette condition de la réussite de l'action rend le principe trivial : « *D* procure satisfaction du besoin *B* » est vrai parce que *D* procure satisfaction du besoin *B*. Cela révèle une conception de la vérité comme redondance, position dont le prochain chapitre discutera. Néanmoins, cette apparente trivialité vaut uniquement si la satisfaction se trouve réellement liée au besoin *B*. Si une incertitude subsiste sur l'origine de la satisfaction, E8 peut se révéler faux malgré cette dernière. Pour éviter ce risque, E8 doit se fractionner afin de donner : « *D* procure satisfaction ». La sémantique liant *B* et *D* se réduit alors à une corrélation qui permet d'orienter le choix des désirs. Cette incertitude incite le fonctionnalisme écologique à continuer de considérer le principe de Ramsey comme une condition, et non comme une définition de vérité. À noter que cette incertitude rejoint celle évoquée dans le mentalisme, où deux systèmes d'hypothèses *ad hoc* différents pouvaient aboutir à des prédictions identiques.

La seconde interprétation d'E7 ne s'appuie plus sur un critère de réussite lié à la satisfaction d'un besoin, mais sur un critère mesurant la qualité des prédictions sensorielles. Bien que ce critère de réussite soit plus précis, il continue à porter à la fois sur une croyance et sur le bon déroulement de l'action. En effet, trois causes peuvent expliquer l'échec d'une prédiction associée à une action :

1. Erreur concernant la reconnaissance du contexte.
2. Erreur dans le choix de l'action à mener en fonction du contexte reconnu.
3. Intervention extérieure imprévisible pour l'agent.

La dernière cause n'impute pas l'échec à une croyance fautive mais à un facteur inconnu. Pour que les croyances assurent entièrement le bon déroulement de l'action, il faudrait un sujet : soit omniscient, autrement dit qui connaisse toutes les clauses logiques pouvant intervenir dans le déroulement de l'action, ce qui n'est pas raisonnable comme il a

été montré dans la section sur le symbolisme concernant les logiques non monotones, soit omnipotent ce qui reviendrait à adhérer au solipsisme, or cette position a été écartée au début de ce chapitre. Le principe de Ramsey dépendant de la réalisation de l'action n'est donc pas valide en ces termes.

Toutefois, dans la troisième source d'erreur se distinguent deux types d'interventions extérieures perturbatrices, le premier concernant la réalisation de l'action, et le second concernant les événements imprévus sur *D* ou la nature réelle de *D*. Pour le premier type d'interventions, la proprioception peut aider à orienter un mécanisme de révision des croyances en s'assurant de l'application des commandes motrices. Une nouvelle décomposition devient alors envisageable en posant « l'action *A* entraîne l'information proprioceptive *P* » (Figure I-29). Mais cet artifice, pour se retrouver dans la situation confortable d'évaluer *C* et *A* sachant que *A* a bien été réalisé ne produit qu'un décalage dans l'analyse, puisque la commande motrice une fois envoyée reste soumise à la contingence (par exemple une paralysie temporaire) et que la finalité de l'action dépend de l'absence d'une intervention du second type. À noter que la forme de la commande motrice (pas à pas, schème moteur ou autres) n'a pas de raison d'intervenir dans cette analyse.

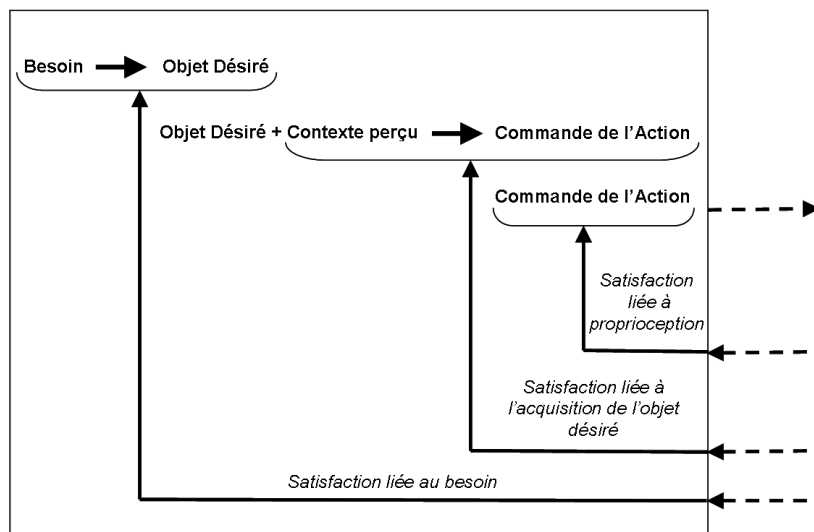


Figure I-29 : Schéma d'un système possédant trois boucles de rétribution.

Les croyances, en définitive, restent dépendantes de la réalisation des actions qu'elles engendrent. L'éventualité d'une compartimentation de clauses logiques liées à cette réalisation n'empêche pas que le soubassement de l'édifice des croyances demeure inextricable de l'action. Par conséquent, les représentations des croyances ne peuvent pas être neutres vis-à-vis de l'action. Cet aspect sera plus particulièrement développé dans la section concernant l'interactionnisme.

Néanmoins, la réponse à la seconde partie de la troisième question (Comment évaluer la vérité d'une croyance, la réussite de l'action est-elle suffisante ?) ne se veut pas totalement négative. Pour les croyances qui évacuent la réalisation de l'action, comme E8 sachant E7, le principe de Ramsey devient un pré-requis c'est-à-dire que les croyances participant à une satisfaction peuvent être vraies ou fausses. En revanche, un échec signifie la présence d'une croyance fautive. A ce stade, cette incertitude pour le fonctionnalisme écologique prend sa source dans le critère de réussite lui-même qui oriente les croyances vers la satisfaction, leur faisant perdre ainsi leur neutralité vis-à-vis de ce qu'elles disent du

monde, en plus du fait que les croyances soient de toute manière situées, autrement dit subjectives de par les sens et les actions qui les sous-tendent. Toutefois, le fonctionnalisme écologique prétend retrouver cette neutralité en multipliant l'intervention d'une même croyance dans des raisonnements justes, conduits par des besoins divers, et ainsi obtenir un noyau de connaissances neutre approchant l'idée de vérité.

Afin d'appliquer une analyse similaire sur le second type de croyances qui n'évacue pas la réalisation de l'action comme E7, la distinction entre croyances épistémiques et croyances événementielles devient essentielle. En effet, en la considérant universelle, E7 devient une croyance épistémique et sa crédibilité dérive des expériences passées, soit $P(E7/c)$, mais lorsque la croyance E7 se trouve employée, elle devient événementielle et la crédibilité porte sur la réussite de l'action présente, soit $P(R/E7,c,i)$ avec la réussite R et les informations présentes i . Cela signifie qu'une règle peut être vraie tout en acceptant des exceptions, ces dernières ne prouvant pas que la règle soit fautive, mais qu'elle existe dans un contexte particulier. La distinction de ces deux utilisations de la croyance E7 autorise à imaginer une gestion de l'ignorance à la fois sur la croyance et sur le monde, autrement dit à introduire directement la notion de probabilité imprécise discutée précédemment. Toute la difficulté de cette position réside dans la transition entre une loi en train de se construire, qui corrige et intègre les nouveaux événements, et une loi établie, inébranlable, qui toutefois collectionne la juxtaposition des exceptions.

Dans tous les cas, le principe de Ramsey se trouve contraint de prendre une forme plus faible : la vérité des croyances augmente la probabilité de leur réussite à défaut de la garantir. Sous réserve de ces réponses, le fonctionnalisme continue à prendre racine à la fois au sein de la famille FP et au sein de la famille FCP pour qui la notion de besoin offre implicitement une mise en correspondance entre le monde idéal et le monde réel. Les hypothèses perceptives vont toutes être rejetées. L'HP1, considérant que le jugement perceptif révèle des qualités sensibles ou des entités qui renvoient immédiatement à des concepts, se trouve en contradiction avec la notion d'affordance développée par Gibson (1979), précurseur du fonctionnalisme écologique.

L'affordance se comprend comme un stimulus évoquant intuitivement et directement la finalité d'une action. L'intellectualisation de celle-ci ne pouvant s'effectuer uniquement par une simulation sensorimotrice interne, échappant ainsi à la possibilité d'une manipulation symbolique. L'affordance s'appuie alors sur un raccourci entre le stimulus et le besoin ou une délégation de la gestion de l'action au bas niveau. Sur le plan cognitif, l'affordance se traduit par une génération automatique d'invitations à agir-pour. Plus exactement, le concept d'affordance en lui-même ne récuse pas HP1, mais c'est la notion de schème sensorimoteur qui le sous-tend. HP2 ne correspond pas à l'idée d'un système dont la gestion des besoins peut être modulaire et locale. Enfin, HP3 est incompatible avec l'idée d'une perception modulée en fonction de l'utilité. Le fonctionnalisme écologique écarte également les trois hypothèses concernant la conception. Toutefois, le refus de HC1 comprend l'abduction au sens faible c'est-à-dire comme une déduction inversée et non comme une véritable construction conceptuelle. HC2 se trouve rejetée pour les mêmes raisons que HP2. La notion de courbe d'utilité offre des perspectives d'explications nouvelles concernant le rôle des émotions, motivant ainsi le refus de HC3.

En ce qui concerne les méthodes d'investigation, aucune ne prend véritablement l'ascendant bien que la MHD doive impérativement prendre en compte la dimension écologique dans le protocole expérimental, c'est-à-dire comprendre l'harmonie qu'a créée

l'individu avec son environnement. De manière identique au connexionnisme, la MA permet d'explorer différents types d'architectures difficilement évaluable sans leur implémentation. Le fonctionnalisme a la particularité d'exploiter tous les types de descriptions, en hybridant la description procédurale et la description mathématisée par l'utilisation de réseaux bayésiens et en utilisant la description componentielle pour exprimer la hiérarchisation des besoins et les tâches associées. La présentation des deux principales approches (Figure I-30) qui se dégagent du fonctionnalisme écologique correspondent chacun à un axe de recherche particulier mais qui, en définitive, se complètent dans la réalisation pratique de systèmes robotiques.

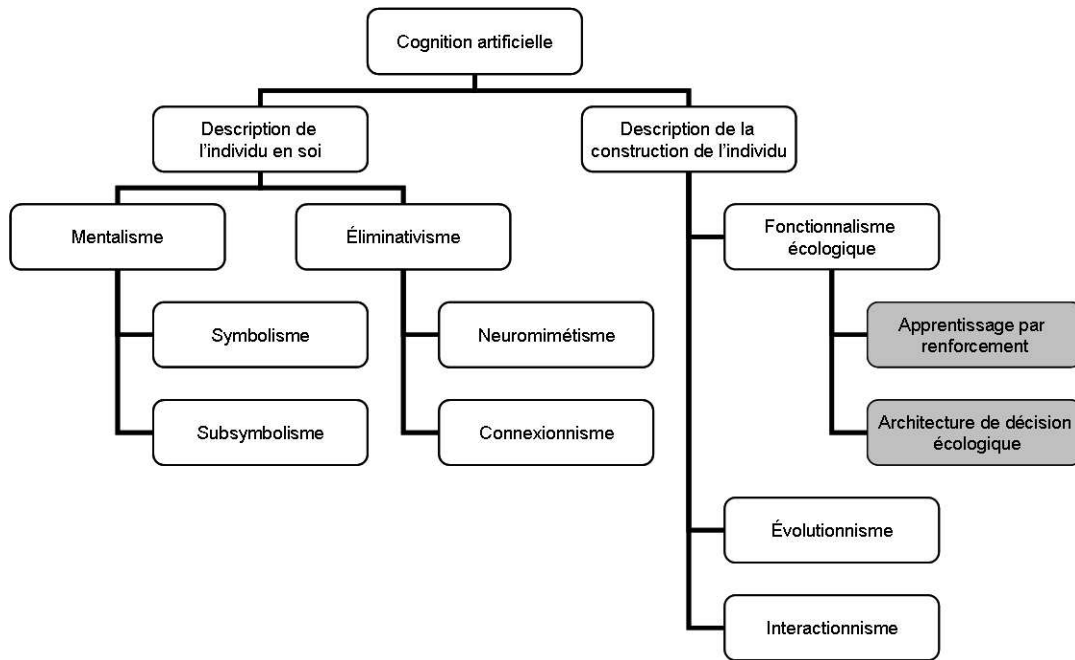


Figure I-30 : Les approches issues de l'apprentissage par renforcement radical et celles issues de la recherche des architectures de décisions écologiques au sein de l'arborescence de la cognition artificielle.

E - Architecture fonctionnelle pour la décision

Le premier axe de recherche se place sur le problème de la hiérarchisation des fonctions. Contrairement au fonctionnalisme subsymbolique dont l'objectif des fonctions consiste à participer à une construction conceptuelle, le fonctionnalisme écologique oriente les fonctions directement vers la subsistance de l'individu et s'inscrit de fait dans le processus évolutif de l'espèce, assurant ainsi une symbiose totale entre l'individu et son environnement. Plus particulièrement, ce dernier point impose une contrainte pratique sur la construction des fonctions. En considérant que les fonctions sont apparues les unes après les autres tout en s'adaptant, il devient nécessaire d'imaginer que les fonctions se greffent au-dessus des fonctions existantes au lieu de s'immiscer dans leur chaîne de traitement. La survie d'un organisme prévalant, les boucles fonctionnelles existantes et fiables ne doivent pas être mises en péril. Par ailleurs, les données paléontologiques indiquent que l'évolution du cerveau se structure par l'apparition de couches successives sans scinder la boucle fonctionnelle inférieure, celle-ci pouvant toujours s'adapter ou s'atrophier sans toutefois se scinder. Le schéma de la figure I-21 montre les deux types

d'architectures les plus représentatifs de l'opposition entre le fonctionnalisme écologique et le fonctionnalisme subsymbolique.

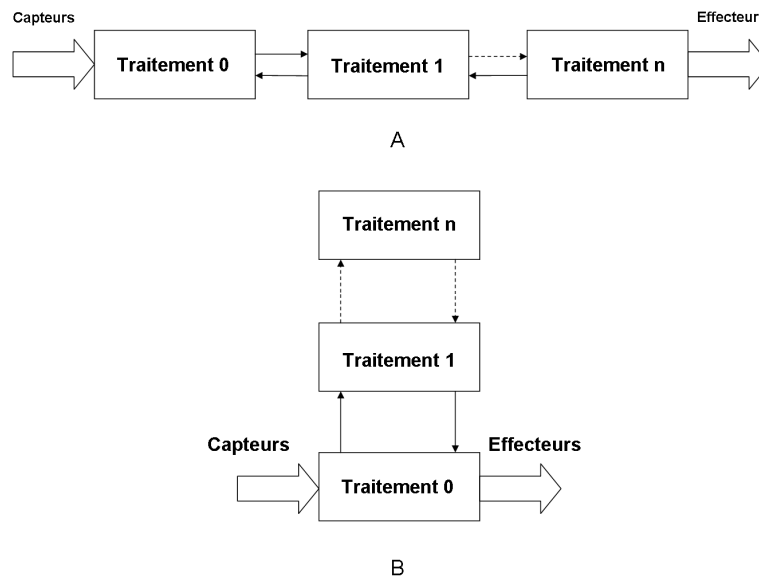


Figure I-31 : Architectures typiques les plus représentatives de l'opposition entre le fonctionnalisme subsymbolique et le fonctionnalisme écologique, respectivement : (A) architecture linéaire et (B) architecture avec subsomption.

L'architecture linéaire repose sur le traitement séquentiel de plus en plus haut niveau pour aboutir à la commande suivante : identifier, modéliser, planifier, réévaluer des tâches puis envoyer des commandes, alors que l'architecture avec subsomption proposée par Brooks (1986) se munit de fonctions élémentaires dont le fonctionnement ne dépend pas en général d'ordres supérieurs. Ainsi, une organisation élémentaire est constituée de fonctions en parallèle comme l'évitement d'obstacle, l'exploration, la cartographie, l'identification des objectifs, etc. Deux avantages significatifs peuvent être avancés pour expliquer l'intérêt de ce développement stratifié. Le premier avantage, en s'appuyant sur ce qui a été avancé précédemment, met l'accent sur le fait qu'à chaque niveau supérieur le problème traité se situe dans un grain spatiotemporel plus important que celui des niveaux inférieurs. L'évitement d'obstacle, par exemple, se focalise sur le présent, alors que l'exploration oblige à utiliser une mémoire et un processus décisionnel allant au-delà de l'état actuel des capteurs. L'avantage de la gestion par échelle de temps différents est de maximiser les satisfactions à court, moyen et long terme. Le second avantage d'une architecture avec subsomption se trouve dans la gestion des exceptions. En effet, un réflexe peut être justifié statiquement mais au lieu de complexifier le processus pour qu'il puisse le prendre en compte, il est parfois plus facile de détecter la défaillance afin de renvoyer le problème vers une fonction plus adaptée. D'autres avantages peuvent être cités : la simplicité des modules qui permet une certaine robustesse, de même que le principe consistant à employer des représentations minimales et locales. Néanmoins, ces recommandations constituent davantage une méthodologie pour la conception de systèmes robotiques destinés à des ingénieurs plutôt qu'un algorithme construisant l'architecture d'un système capable d'autonomie. Les ingénieurs réaliseront des robots en plus grande harmonie avec leur environnement et d'une meilleure fiabilité mais resteront cantonnés aux besoins définis par chacune des fonctions, bien que chacune d'elles puisse faire l'objet d'optimisations améliorant la performance du système.

F - Apprentissage par renforcement

Le second axe de recherche se concentre sur les mécanismes de renforcement qui présentent l'avantage de découvrir les actions menant à la satisfaction ou aux désagréments lorsqu'ils sont inconnus a priori. La présence de ce mécanisme au cœur de différents types de conditionnement animal qui seront détaillés ultérieurement au second chapitre montre son importance dans une démarche écologique. L'apprentissage s'effectue au cours d'une suite d'essais-erreurs où le critère d'apprentissage n'est pas une distance avec ce que le système doit apprendre mais une rétribution le guidant vers l'objectif. En cela l'apprentissage par renforcement se trouve qualifié de semi-supervisé. Toutefois, le renforcement peut être compris également comme un moyen de sélection d'une sous-base d'apprentissage car une analyse statistique sur une grande base de données peut écraser un sous-ensemble pertinent pour le système. Néanmoins, la notion de rétribution contribuant fortement au formalisme proposé dans ce mémoire, l'analyse proposée ici se focalise uniquement sur le renforcement compris comme maximisation de la récompense.

L'apprentissage par renforcement se comprend comme l'exploitation de la quatrième interprétation de l'expression $P(h,c)$, mais la compréhension de cette exploitation nécessite une présentation succincte des concepts liés aux processus stochastiques. Ces derniers représentent une évolution d'une variable aléatoire comme le lancer sans cesse renouvelé d'une pièce de monnaie dont la valeur du résultat serait la variable aléatoire. Cet exemple porte sur une évolution temporelle mais elle aurait pu être spatiale en prenant pour processus stochastique une image dont la variable aléatoire aurait été la valeur d'un pixel. Souvent employés en l'intelligence artificielle, les processus Markoviens constituent un sous-ensemble des processus stochastiques se déclinant en deux types : les automates de Markov et les automates de Markov à états cachés. Le premier type de processus markovien décrit la probabilité de transition entre les états d'un automate en précisant que cette probabilité dépend uniquement des k états précédents, ici $k = 1$; soit la distribution de probabilité conditionnelle d'une variable aléatoire X :

$$P(X_{n+1} / X_0, X_1, \dots, X_n) = P(X_{n+1} / X_n)$$

En ces termes, un automate markovien appelé également chaîne de Markov représente un sous-ensemble des réseaux bayésiens. Pour se conformer à ce cadre, l'exemple des urnes se modifie de telle sorte que le choix de l'urne dépende du tirage : lorsqu'une boule blanche est tirée, il y a changement d'urne. En nommant, la première urne A et la seconde B, le tout représente un processus stochastique ayant pour variable aléatoire X l'urne employée (Figure I-32). Ce formalisme permet d'évaluer la probabilité d'apparition des états, ici des urnes.

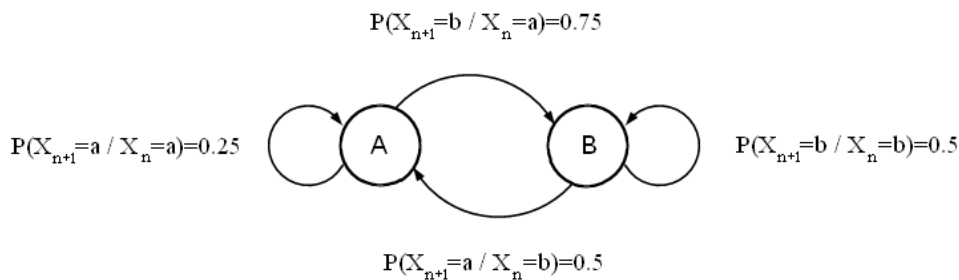


Figure I-32 : Exemple d'un automate markovien à deux états.

Le second type de processus Markovien, les automates de Markov à états cachés, décrit avec le même principe la probabilité de transition entre les états d'un automate, mais le changement d'état ne peut être directement observé. En associant à chacune des urnes de l'exemple précédent une urne contenant dans des proportions différentes des boules vertes et rouges, deux tirages simultanés peuvent être imaginés comme suit : un premier tirage caché s'effectue dans l'une des urnes initiales, la boule blanche indiquant toujours le changement, mais le résultat reste secret ; en parallèle, le résultat du tirage de l'urne associée est divulgué sans préciser l'urne d'origine (Figure I-33). Cela traduit l'enchevêtrement de deux variables aléatoires : X l'urne employée et Y le résultat du tirage de l'urne associée, mais seule la seconde se trouve visible. La popularité de ce paradigme vient du fait que les automates markoviens à états cachés modélisent le problème de la perception du phénomène avec son incertitude et ce qui le génère.

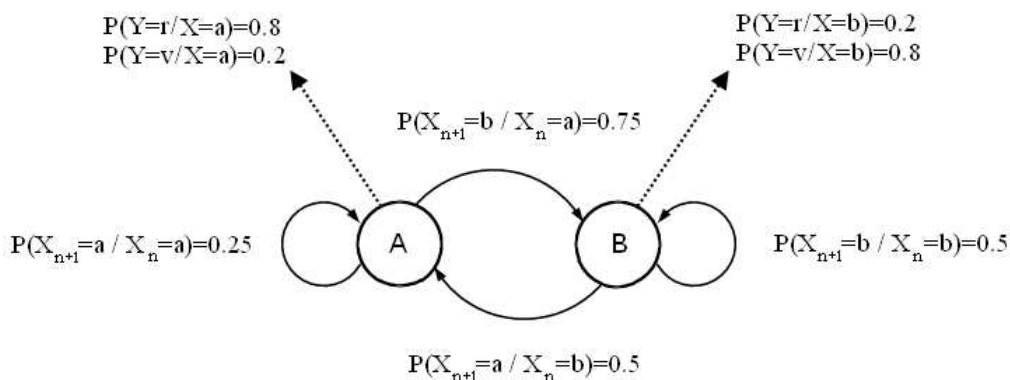


Figure I-33 : Exemple d'un automate markovien avec les variables aléatoires Y correspondant aux observables et X correspondant aux états cachés.

Dans ce formalisme, l'algorithme de Viterbi (1967) permet de calculer d'une part la probabilité d'une séquence d'observations et d'autre part la séquence d'états la plus probable pour une séquence d'observations donnée (Tableau I-7). Cependant, l'application de cet algorithme nécessite que les états de l'automate et leurs distributions soient connus ou estimés, ce qui peut se révéler très difficile a priori. Par ailleurs, cette estimation du nombre d'états a priori, comme le souligne Sigaud (2001), représente pour un système possédant un processus markovien à états cachés un modèle du monde alors qu'un processus markovien classique décrit seulement un comportement réactif.

Séquence de la variable aléatoire X correspondant aux états cachés	a b b a b a a b b b a b a b b
Séquence de la variable aléatoire Y correspondant aux observations	r v r v r r r v v r v r r r v

Tableau I-7 : Exemples de séquences issues de 15 tirages successifs.

En robotique mobile, l'emploi de ces représentations oblige à poser comme hypothèse forte la relation entre un agent et son environnement comme un processus markovien réglé en fonction de récompenses externes exprimant la composante décisionnelle. Autrement dit, l'agent possède à chaque pas l'information suffisante pour décider au mieux puisque le principe de Markov stipule que la probabilité de transition d'un

état à un autre dépend uniquement des k états précédents mémorisés, généralement $k=1$. Les deux types de processus markovien se transforment alors réciproquement en processus de décision markovien (PMD) et processus de décision markovien partiellement observable (POMDP).

Cette section souhaite uniquement présenter rapidement les différents algorithmes d'apprentissage concernant les PMD afin de montrer à la fois le parallèle avec la quatrième interprétation de P(h,c) et des éléments de comparaisons lors de la présentation du formalisme proposé dans ce mémoire. Les techniques utilisées pour les POMDP dérivent des mêmes principes concernant le renforcement mais elles se complexifient du fait qu'à l'incertitude sur l'effet de l'action se rajoute celle de l'état courant. L'exemple illustré par la Figure I-27 correspond à un environnement type pour les POMDP : indépendamment de l'incertitude des capteurs, il existe une incertitude entre les états A et B . L'ambiguïté vient du fait que les états sont partiellement observables à partir des seules informations sensorielles.

La problématique de l'apprentissage d'un PMD se formalise à l'aide des deux fonctions suivantes : une fonction de transition T qui évalue la probabilité d'obtenir un état s_k suite à l'action a_j (ou plus justement la commande de l'action a_j) dans la situation s_i :

$$T(s_i, a_j, s_k) = P(S_{n+1} = s_k / S_n = s_i, A_n = a_j)$$

La fonction de récompense R qui donne l'espérance de la récompense lorsque s_i devient s_k en effectuant a_j :

$$R(s_i, a_j, s_k) = E(r_{n+1} / S_n = s_i, A_n = a_j, S_{n+1} = s_k)$$

Parallèlement à cette fonction qui détermine la récompense moyenne à chaque pas afin d'intégrer la possibilité de sacrifier une récompense à court terme pour en obtenir une à long terme plus conséquente, une seconde fonction β_n cumule les récompenses attendues à long terme. L'expression β_n diffère selon l'importance accordée aux récompenses futures, toutefois en voici un exemple typique :

$$\beta_n = \sum_{k=0}^{\infty} \gamma^k r_{n+k+1}$$

En ces termes, l'objectif de l'apprentissage consiste à déterminer une stratégie π qui traduit la probabilité d'effectuer l'action a pour chaque état s . La détermination de cette stratégie s'appuie traditionnellement sur une fonction de valeur V_π qui traduit l'espérance de β_n pour une stratégie π :

$$V_\pi(s) = E(\beta_n^\pi / S_n = s)$$

Cette expression de la fonction de valeur utilise uniquement l'état sensoriel et une autre expression permet de saisir le couple état sensoriel et action :

$$Q_\pi(s, a) = E(\beta_n^\pi / S_n = s, A_n = a)$$

La mise en équation de ces fonctions correspond à des équations de Bellman. Ces dernières autorisent à choisir, dans un environnement markovien, la stratégie optimum en recherchant les valeurs maximales de la fonction de valeur. À ce stade, une mise en correspondance avec la quatrième interprétation se réalise en identifiant la fonction de récompense R avec la fonction du gain provisionnel G et la fonction de valeur V ou Q avec la fonction d'utilité U . La fonction de transition n'apparaît pas explicitement dans l'exemple

du marchand de glace, parce que l'action consistant à se placer à l'endroit choisi est certaine et que la météo demeure indépendante du marchand de glace quelles que soient ses actions.

Trois classes d'algorithmes (Cornuéjol et Miclet, 2002) proposent des solutions pour approximer les fonctions V ou Q , et π : les algorithmes de programmation dynamique, les méthodes de Monte-Carlo et les méthodes par différence temporelle. Les algorithmes de la première classe supposant les fonctions T et R se composent de deux phases qui se répètent jusqu'à convergence : l'estimation de la fonction de valeur optimale par résolution des équations de Bellman pour une stratégie a priori puis l'amélioration de la stratégie par rapport à cette nouvelle fonction de valeur. Cette première classe offre la possibilité d'un apprentissage incrémental, c'est-à-dire que l'apprentissage s'effectue en même temps que les expériences de l'agent. Cependant, la nécessité de connaître les fonctions T et R pose les mêmes problèmes évoqués à propos du symbolisme concernant la difficulté de créer un modèle a priori exact du monde. La deuxième classe d'algorithmes explore plus ou moins aléatoirement dans un premier temps l'environnement pour évaluer la triade état-action-récompense, puis calcule dans un second temps la stratégie optimale. Aucun modèle du monde n'est ici directement employé, mais l'algorithme d'apprentissage n'est plus incrémental.

La troisième classe d'algorithmes se sert de la récompense pour évaluer le couple état-action mais également la fonction de valeur, en relevant la différence entre la récompense prédite et la récompense réelle. Cette double utilisation de la récompense permet une évaluation itérative de la fonction de valeur et de la stratégie optimale sans nécessiter un modèle a priori du monde : la troisième classe d'algorithme surpasse en cela les deux premières. Barto (1995) dégage ce double processus en proposant une architecture fondée sur l'erreur de prédiction avec un module « acteur » et un module « critique », le premier régissant et le second évaluant la fonction de valeur. Parallèlement, des données électrophysiologiques sur le singe suggèrent que les ganglions de la base avec les décharges dopaminergiques pour signal de renforcement correspondraient à une architecture de type acteur-critique (Suri et Shultz, 1998). Ici, le fonctionnalisme écologique et le neuromimétisme se rejoignent dans l'approche animat (Wilson, 1990) qui souhaite s'inspirer des données neurobiologiques à la limite de la modélisation, tout en se prévalant du caractère situé avec les spécificités d'un corps robotique. Cependant, SARSA et le Q-Learning (Watkins, 1989) restent les algorithmes les plus populaires, bien que la portée temporelle de leur prédiction n'excède pas un pas de temps.

Dans l'ensemble, malgré divers problèmes liés à leur implémentation, comme la discrétisation des états et des capteurs, les méthodes de renforcement demeurent performantes pour un environnement ayant toujours la même fonction R . Une solution (Fox et al., 2005) se trouve dans la constitution d'un ensemble de cartes de navigation dont la sélection s'effectue par un algorithme d'Estimation et Maximisation (Dempster et al., 1977) qui sera détaillée au troisième chapitre. Néanmoins, la multiplication des cartes permet seulement de posséder plusieurs comportements pour divers environnements spécifiques, mais n'offre pas de solution pour dépasser le stade du renforcement qui semble intrinsèquement limité par la définition explicite de l'objet de la rétribution par le concepteur, ce qui constitue la limitation principale du fonctionnalisme écologique.

**

En conclusion, le fonctionnalisme écologique repose sur l'apport de deux notions : le besoin et l'action. La notion de besoin autorise le dépassement des approches de la

description de l'individu en soi, dans le sens où le besoin guide la construction de schèmes comportementaux et justifie de surcroît l'introduction d'une représentation du monde (ou de croyances), quelle que soit leur forme, grâce à la multiplication des besoins matériels sous couvert de certaines conditions environnementales. En contrepartie, cette notion de besoin se trouve à l'origine d'un tropisme supplémentaire que subissent les représentations. Seule la satisfaction des besoins encourage l'association entre croyance et action. L'aspect écologique souligne que l'agent recueille seulement les indices suffisant à conduire une action réussie. Cet opportunisme oblige les croyances à toujours être « utiles à quelque chose ». La carte n'est pas le territoire (Korzybski, 1933) : aussi fidèle qu'elle soit, elle existe comme guide pour voyager. Néanmoins pour les tenants de la famille FCP, sa participation dans des cas variés suggère qu'elle peut tendre à la limite vers l'homologue du territoire situé dans le monde idéal. L'agent cognitif dans ce paradigme se définit par son interaction avec son environnement en fonction de ses besoins et de ses capacités, autrement dit par sa téléonomie. Une contradiction apparaît précisément avec le caractère situé des agents puisque la définition des besoins se fait à la troisième personne. La notion de l'action rappelle que l'agent subit toujours la contingence du monde extérieur et qu'il faut ainsi distinguer la satisfaction globale d'une croyance et la probabilité de sa bonne utilisation pour une situation spécifique. L'anticipation des conséquences se révèle être le seul mécanisme permettant de dégager ces deux types de croyances et de cibler au mieux la révision des croyances. L'importance de ce mécanisme se traduit par son efficacité dans les algorithmes l'utilisant comme les méthodes à différences temporelles, mais également dans d'autres méthodes comme les systèmes de classeurs (Sigaud, 2001) qui seront détaillées dans le troisième chapitre. Mais son importance se retrouve surtout chez les animaux cognitivement évolués (Berthoz, 2000).

En somme, l'approche écologique, tout en refusant les hypothèses liées à la perception et à la conception (Tableau I-8), propose que les représentations se construisent grâce à l'interaction environnementale. Ces constructions ne peuvent pas rendre compte fidèlement du monde à cause de l'action, cependant la vérité absolue demeure en limite. Par ailleurs, l'anticipation joue un rôle déterminant pour diminuer l'impact sur le développement des croyances à un niveau supérieur de l'incertitude liée à l'action. Malgré tout, la définition a priori des besoins ne correspond pas aux exigences qu'impose la notion d'autonomie, et ce paradigme qui repose quasiment uniquement sur celle-ci n'offre pas d'alternatives. Le courant évolutionniste souhaite répondre à ces difficultés en décalant l'analyse au niveau de l'espèce pour aller au-delà des besoins individuels.

		Types de famille métaphysique				Hypothèses sur la perception			Hypothèses sur la conception			Types de méthode d'investigation			Types de description		
		FC	FC	FP	FA	HP1	HP2	HP3	HC1	HC2	HC3	MHD	MC	MA	DP	DM	DC
Fonctionnalisme écologique	Apprentissage par renforcement	Red	Green	Green	Red	Red	Red	Red	Yellow	Red	Red	Yellow	Green	Green	Green	Yellow	Yellow
	Architecture décisionnelle	Red	Green	Green	Red	Red	Red	Red	Yellow	Red	Red	Yellow	Green	Green	Green	Yellow	Green

Tableau I-8 : Récapitulation de diverses positions relatives à la grille d'analyse des deux approches provenant du fonctionnalisme écologique. Les cases vertes correspondent à un avis favorable, les cases jaunes à un avis mitigé et les cases rouges à un avis défavorable.

3.2.2. L'évolutionnisme

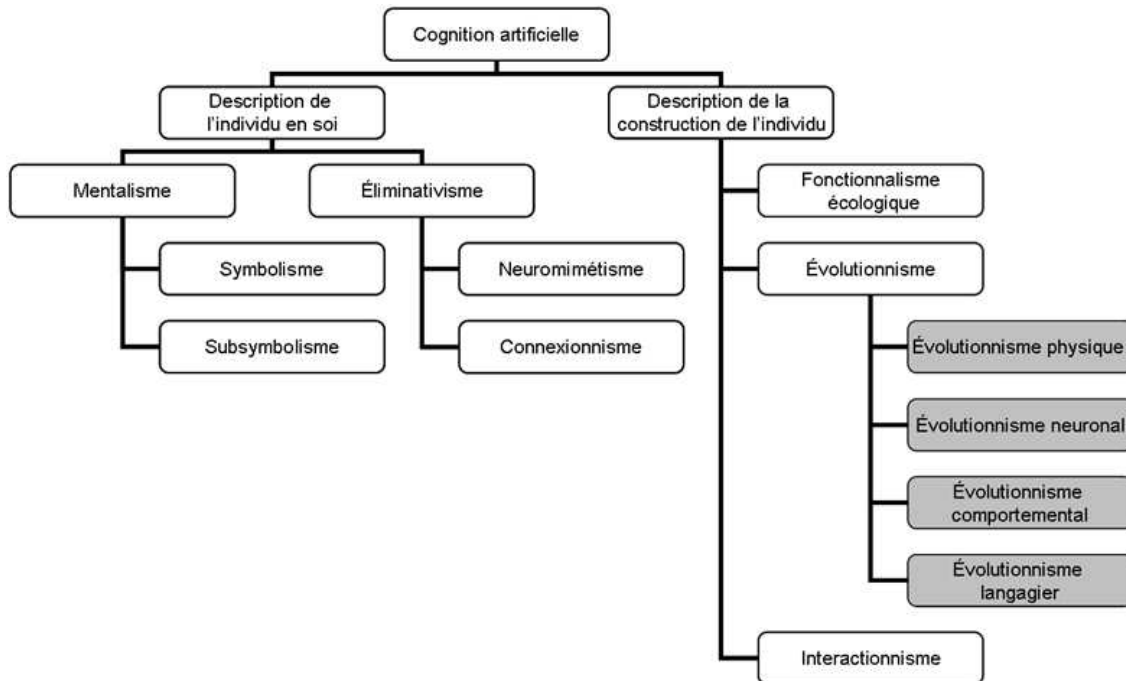


Figure I-34 : Les quatre approches provenant de l'évolutionnisme au sein de l'arborescence de la cognition artificielle.

Le mouvement évolutionniste aspire à créer une vie artificielle, qui dans un deuxième temps, pourra développer des capacités cognitives. Cette approche s'appuie sur le néodarwinisme qui se résume sommairement à la conjonction de deux thèses : la sélection naturelle, la survie du plus apte, guide l'évolution des espèces, et la génétique représente le mécanisme permettant la reproduction des caractères de l'individu sélectionné. Ces deux principes ont généré un grand nombre de positions proposant un cadre interprétatif pour l'étude des conséquences et les raisons de la sélection, ainsi que sur ce que doit recouvrir le concept de gène. Synthétiquement, deux courants se dégagent. (A) Le premier sert de paradigme à quasiment toutes approches évolutionnistes en sciences de l'artificiel dont les quatre principales seront donc détaillées juste après la présentation du premier courant issue du néodarwinisme. (B) L'exposé du second courant s'attachera à montrer les faiblesses du premier courant et par ricochet celles des approches en vie artificielle associées, en expliquant le concept d'autopoïèse et ses conséquences.

A - La première interprétation du néodarwinisme

En s'appuyant sur l'étude de Stewart (2004), le néodarwinisme résulte de la synthèse des travaux emblématiques de Darwin (1859) et de Mendel (1863). La théorie du premier repose sur deux propositions : tous les êtres vivants descendent d'un ancêtre commun d'organisation minimale et toutes les espèces dérivent les unes des autres dessinant un arbre généalogique, la question sur la forme de ce dernier restant néanmoins ouverte. La théorie de Mendel repose sur l'existence d'entités discrètes, les gènes, conditionnant l'apparition de certains traits de caractères permettant de les différencier d'autres individus, la reproduction mélangeant les gènes provenant des parents. Les implications de ces théories poussent à

dessiner un schéma (Figure I-35) qui sépare très distinctement le soma, le corps de l'individu, et le plasma germinatif, l'ensemble des gènes (Weismann, 1883).

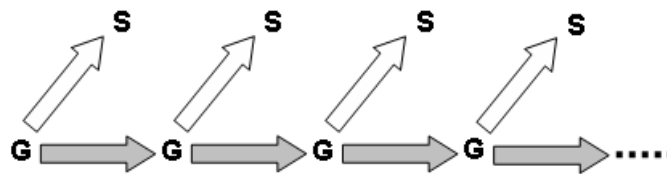


Figure I-35 : Le schéma Weismannien avec le plasma germinatif G qui se transmet de génération en génération sans subir les influences de l'environnement par le biais du soma S, le corps de l'organisme.

À ce schéma se rajoute la spécificité du premier courant néodarwiniste qui s'entend à étendre le concept de gène mendélien portant sur des caractères différentiels à tous les caractères de l'organisme. Par exemple, pour un humain, la couleur des yeux correspond à un gène mendélien, contrairement au nombre de bras. Cela conduit à l'idée que les gènes forment un code, un programme dont l'exécution s'identifie à l'ontogenèse. Plus précisément, l'inné concerne l'expression des caractères génétiques, et l'acquis renvoie aux modifications de l'organisme de par son interaction avec son environnement, mais sans qu'il puisse influencer sur le patrimoine génétique : le plasma germinatif, l'inné. Cette forte séparation entre inné et acquis entraîne une dépréciation de l'organisme lui-même pour le processus évolutif des gènes. L'organisme devient subordonné aux gènes puisque ce sont eux les invariants, ce qui est transmis (Dawkins, 1976).

Dans cette logique, l'histoire du vivant commence par un mécanisme d'autoreproduction de macromolécules qui, suite à des mutations (des erreurs de recopie), se complexifient sans empêcher l'accomplissement de l'autoreproduction. Ces modifications successives ont eu pour conséquence de complexifier le processus au point de confectionner une membrane au sein de laquelle l'unité autorépliatrice a pu se « formaliser » elle-même à travers la molécule d'ADN qui code tout le déroulement du développement de l'individu. Les individus multicellulaires résulteraient alors d'une complexification réussie de leurs ancêtres sous la pression de mutations guidées par la sélection naturelle. Dans ce schéma, l'individu ne possède pas de téléonomie propre, et l'objectif de survivre jusqu'à la reproduction provient indirectement du processus originel.

La girafe n'a pas un long cou pour manger les feuilles mais c'est parce qu'elle a un long cou qu'elle les mange. L'intentionnalité ou le vécu de l'individu n'interfèrent en rien dans la construction par reproduction d'un individu ultérieur. Mais au niveau de l'espèce, ceux qui se reproduisent le plus correspondent à ceux qui se sont le mieux adaptés. Si un individu possède des caractéristiques qui constituent un avantage dans son environnement, alors le principe de la sélection propagera son patrimoine génétique. L'évolution, bien qu'aveugle, s'oriente vers l'optimisation des individus. En somme, tout en dénonçant le fixisme qui prône une création *ex nihilo* de toutes les espèces dans un passé lointain sans autre modification, le néodarwinisme suggère une téléologie du vivant : soit une évolution vers une complexité croissante avec l'argument existentialiste invoquant la nécessité pour l'univers d'accéder à la conscience d'être, soit une évolution vers une espèce optimale. Plusieurs arguments allant à l'encontre de ces positions seront exposés lors de la présentation du second courant.

Le néodarwinisme ainsi compris se fonde au choix sur l'une des trois premières familles philosophiques avec des convictions épistémologiques similaires à celles

rencontrées dans la partie consacrée à l'éliminativisme. Les courants issus de la famille FC favoriseront les arguments téléologiques basés sur la complexité et la nécessité de logique de la conscience d'être, alors que les familles FCP et FP se confondent quasiment en réduisant les particules élémentaires à des mots matériels intégrant leurs propres grammaires avec ou non leur reflet logique dans le monde idéal. Les hypothèses portant sur la perception et sur la conception se révèlent subsidiaires, puisque la véritable source de la cognition et sa compréhension se trouve dans la phylogénie. Ici, l'évolution du vivant se réduit à l'évolution de l'ADN assimilé à un langage de programmation dont les mots parfois au-delà de leur combinaison logique peuvent subir des modifications aléatoires. Les sujets de recherche sur l'impact des mutations, sur la dérive d'une population, etc. montrent que la dynamique de l'évolution est plurifactorielle et complexe. Toutefois, ces aspects ne remettent pas directement en cause le schéma du néodarwinisme représenté par l'algorithme génétique de Holland (1975).

Parmi les quatre approches évoquées en introduction de cette section, trois d'entre elles héritent directement de la structure de l'algorithme génétique. Ce dernier se décompose en trois étapes (Figure I-36).

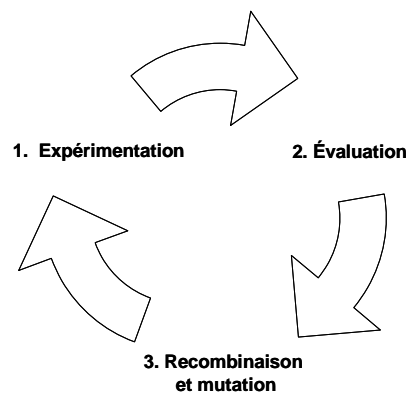


Figure I-36 : Cycle d'un algorithme génétique type.

i - L'évolution des formes d'organismes

La première approche se focalise sur l'évolution des formes d'organismes. En cela, cette approche se revendique comme étant la plus proche du principe de l'évolution naturelle. Une population initiale générée aléatoirement se trouve confrontée à un environnement. Après une période définie, un score est calculé pour chaque individu, selon des critères généralement liés à la motricité, et les gènes des individus ayant obtenu les meilleurs résultats se recombinent pour produire un nouvel individu. Chaque gène renvoie directement à une caractéristique de l'individu. La puissance heuristique de ce type d'algorithmes fut particulièrement mise en avant par Sims (1994) qui a montré qu'ils permettent de générer des individus d'une grande variété de formes, certaines familières, d'autres surprenantes dans un monde physique virtuel avec des critères de sélection rudimentaires.

ii - L'évolution des schèmes comportementaux

La deuxième approche exploite les capacités heuristiques des algorithmes génétiques non plus pour trouver des caractéristiques physiques adaptées, mais pour déterminer un comportement adapté. Le code génétique représente alors les règles comportementales qui

peuvent prendre la forme de poids synaptiques d'un réseau de neurones. De génération en génération, l'attitude évolue en fonction des critères de sélection. Ces critères s'interprètent comme une forme de renforcement mais visible uniquement au niveau de l'espèce. La conception de la cognition que développe cette approche ressemble à celle du fonctionnalisme écologique tout en ignorant toute la problématique concernant les probabilités. Par ailleurs, l'utilisation d'algorithmes génétiques ne facilite pas une construction par couche ; en effet, pour éviter un démantèlement des solutions précédentes, la détermination des modifications à réaliser s'effectue uniquement par le concepteur, soit à la troisième personne. Cette approche a également permis de modéliser des dynamiques de coévolution de stratégie entre proie et prédateur (Floreano et al., 1998).

iii - L'évolution appliquée à une population de neurones

La troisième approche propose aussi d'employer les concepts issus du néodarwinisme pour expliquer le fonctionnement de la cognition, mais directement, sur le substrat de la cognition. Autrement dit, les neurones prennent la place des individus dans l'algorithme génétique. La cognition résulte de la sélection naturelle entre neurones (Changeux, 1983 ; Edelman, 1987). Le critère de sélection peut être par exemple lié à l'activité du neurone avec celle des autres. Cette vision des réseaux de neurones possède des liens avec le connexionnisme, plus particulièrement avec les approches qui ne préjugent pas de la finalité de l'auto-organisation.

iv - L'évolution appliquée à la constitution d'un langage

Sans employer exactement les algorithmes génétiques, la quatrième approche applique l'idée d'évolution dynamique par sélection et modifications successives à la construction d'un langage, la glossogénèse. Dans cette approche, les travaux (Kaplan, 1999) se scindent en deux catégories : ceux qui circonscrivent le problème à un jeu de langage, et ceux qui considèrent que le jeu de langage s'appuie nécessairement sur un vécu. La première catégorie correspond en fait à une facette d'une approche symbolique, avec les défauts et qualités déjà évoqués. La seconde catégorie souligne l'importance de l'ancrage des symboles dans la construction d'un lexique et ressemble à une démarche subsymbolique du point de vue de l'individu. L'avantage par rapport à la première catégorie est que la conjonction du subsymbolique et de l'évolutionniste révèle que l'interaction ne se limite pas à l'environnement physique mais inclut également l'environnement social. À l'aide de caméras communiquant entre elles des symboles en fonction des éléments (figures géométriques colorées sur un tableau blanc) se trouvant dans leur champ de vision, Steels (1999) montre qu'un lexique et une hiérarchisation de concepts descriptifs peuvent se constituer grâce à un jeu de pointage. Le jeu consiste à jouer alternativement le rôle de locuteur et d'interlocuteur. Le locuteur identifie un élément ou des éléments et y associe un symbole qu'il communique ensuite à son interlocuteur. Ce dernier doit interpréter ce symbole puis désigner le ou les éléments qu'il dénote. En cas d'erreur, le locuteur pointe la zone à laquelle il faisait allusion, permettant à l'interlocuteur de se corriger la fois suivante.

Ces résultats montrent qu'une dynamique de communication permet d'aboutir à des opérations de généralisation ou d'induction sans que les individus soient pourvus de telles capacités. Malgré les apparences, cette propriété sociétale ne réfute pas les conclusions concernant l'apport transcendantal des interactions sociales dans la perception qui ont été évoquées lors de l'analyse épistémologique dans la section portant sur le connexionnisme. En effet, les éléments dans le champ de vision étant définis par un vecteur de catégorie portant aussi bien sur la couleur que sur la forme ou la position, l'ensemble des classes

associatives possibles est identique pour tous les agents et l'interaction sociale telle qu'elle est imaginée ici n'apporte qu'une catégorie perceptive de base. Elle permet de construire un lexique et une hiérarchie conceptuelle descriptive consensuelle avec les autres agents mais a priori équivalente à d'autres. Un agent qui possède déjà un ensemble de classeurs fixes et définis *catégorise* les informations perceptives selon ces catégories prédéfinies alors que l'agent qui possède un ensemble extensible de classeurs modifiables *classifie* les informations perceptives en ajustant ou créant des classes s'il y a lieu. Si l'agent possède un algorithme de catégorisation et non de classification, alors l'interaction sociale transcendant la cognition individuelle devient indispensable pour effectuer la classification ; sinon l'interaction sociale joue le rôle de catalyseur néanmoins crucial dans le processus de dénomination et de classification. Dans le premier cas la boucle de rétroaction se trouve en dehors de l'agent et dans le second cas une boucle de rétroaction se situe dans l'agent tout en acceptant une autre boucle de rétroaction facultative et externe (Figure I-37). Mais dans tous les cas, ici, l'agent ne modifie pas les percepts de base donnés a priori permettant de construire par association les catégories.

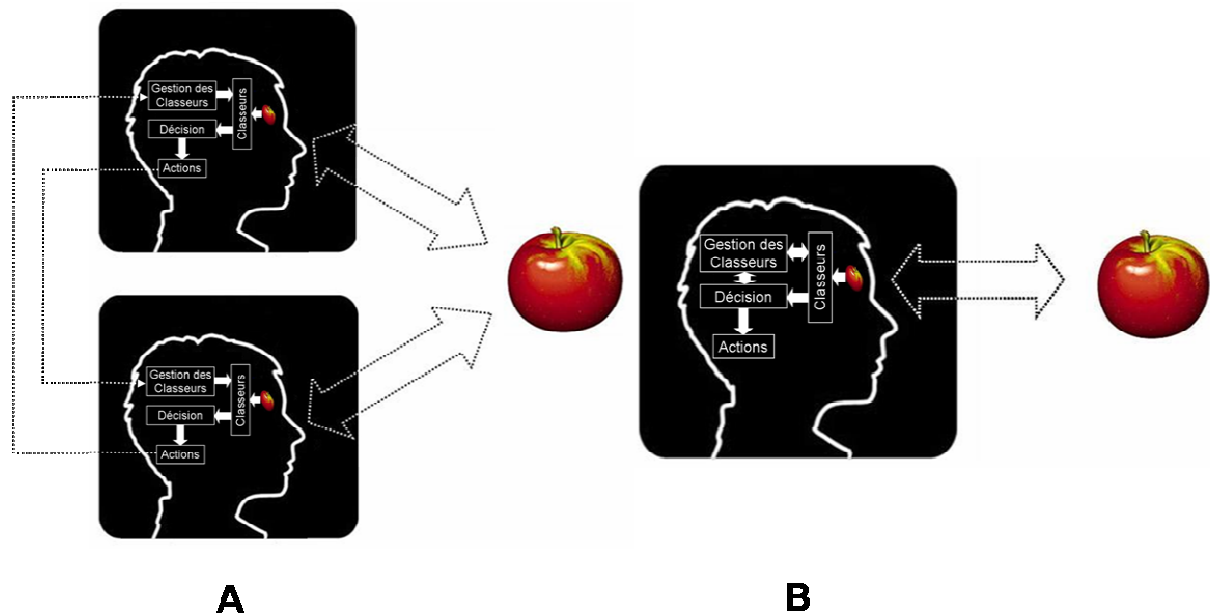


Figure I-37 : Le schéma A illustre deux agents dont le fonctionnement de l'architecture cognitive dépend cruciallement de l'interaction sociale puisqu'elle constitue l'unique boucle de rétroaction et le schéma B illustre un agent qui possède une boucle de rétroaction intrinsèquement dans l'architecture cognitive pouvant ainsi se dispenser d'une interaction sociale pour définir ses classeurs

L'association entre les éléments et le symbole s'ajuste et se propage à travers les agents, ce qui entraîne une compétition lexicale pour un même concept, par exemple le concept descriptif [EN-HAUT, A-GAUCHE, NIVEAU-ELEVE-DE-ROUGE] peut se nommer de plusieurs manières. L'analogie avec le principe de l'évolution s'accroît mais ne s'applique pas aux agents. Les mots constituent alors la population soumise à l'évolution. Ils se transmettent d'agent en agent et rentrent en compétition entre eux tout en ajustant les concepts auxquels ils renvoient. Dawkins (1976) a plus particulièrement développé la transposition entre une interprétation du néodarwinisme radical (le gène se trouve à l'origine de la construction du corps afin de mieux se multiplier) et l'évolution des idées. Dans cette optique, la cognition des personnes constitue le milieu écologique des idées

appelées « mèmes » ou « entités répliquatives d'information ». L'imitation comme moyen de reproduction devient la définition et la raison d'être de la cognition. De la même façon que les gènes subordonnent la biologie du corps, les mèmes subordonnent la psychologie de l'individu. Cette position se défend en montrant l'existence de transmissions culturelles de génération en génération aussi bien chez les hommes (des chansons enfantines aux principes religieux) que chez le singe (le nettoyage de pomme de terre en utilisant de l'eau) (Watson, 1979).

B - La seconde interprétation du néodarwinisme

La seconde interprétation rejette la place centrale conférée au gène par le néodarwinisme qui conduit à l'hégémonie de la génétique comme science explicative du vivant puisque les organismes résultent uniquement du bon déroulement d'un programme codé par les gènes, le programme génétique. La notion de programme génétique subit principalement deux critiques, l'une quantitative et l'autre de principe. La première souligne que le génome d'un mammifère (dont l'homme) ne possède que 30000 gènes et que ce chiffre paraît dérisoire au vue de la complexité du corps humain et de son système nerveux. Ce sentiment perdure même en comprenant la compression de l'information comme un moyen de production de ce qu'il représente (Chaitin, 2002). Ainsi, par exemple, l'algorithme générant le nombre Pi représente en puissance le nombre Pi. La seconde critique, plus fondamentale et longuement développée par Stewart (2004), porte sur la capacité des gènes à décrire ce qui appartient aux caractéristiques communes à une espèce. Plus précisément, les gènes mendéliens représentent des entités qui se transmettent lors de la reproduction. Ils s'expriment par une différence par rapport aux individus qui ne possèdent pas ces entités. Les gènes mendéliens se définissent uniquement dans ce cadre, et étendre la définition du gène à toutes les caractéristiques d'un individu dépasse ce cadre. La première interprétation du néodarwinisme estime que les algorithmes illustrent bien la possibilité d'étendre le concept de gène mendélien en associant à chaque génotype un phénotype et en affirmant que l'ensemble des phénotypes forme un individu. Néanmoins, cet argument se révèle fallacieux étant donné que l'encodage du génome et de l'algorithme l'interprétant qui indique la manière de combiner le matériel génétique sont autant d'éléments établis a priori par le concepteur et qui tombent en dehors de la modélisation d'une évolution génétique.

Cette seconde critique traduit en fait l'attitude paradoxale consistant à utiliser la métaphore de l'ordinateur sans ordinateur. Un programme se constitue d'une suite de commandes qu'un ordinateur exécute mécaniquement ; or le programme génétique n'est pas interprété par un ordinateur, il doit alors en lui-même posséder le code et les moyens de son expression. Cette autoréférence rend paradoxale la métaphore. La notion d'auto-organisation lève en partie le paradoxe, en offrant une explication à l'émergence d'une organisation matérielle sans coordinateur et sans interpréteur, mais la forme de l'organisme n'est alors plus codée. Les flocons de glace, par exemple, possèdent tous six branches quasiment identiques, mais chaque flocon de glace possède une structure de branche différente. Le processus de cristallisation se révèle extrêmement sensible aux conditions atmosphériques. Il n'existe dans ce cas aucun programme déterminant la forme d'un flocon, son ontogenèse, car celle-ci dépend uniquement des propriétés auto-organisatrices de la matière et des conditions environnementales. La génétique se comprend alors comme une internalisation de facteurs influant sur l'ontogenèse. Autrement dit, l'ontogenèse est soumise à des facteurs endogènes, les gènes, et à des facteurs exogènes, les conditions environnementales, et elle devient donc un processus historique. Les gènes ne sont plus des

constructeurs de caractères précis pour un organisme, mais des aiguilleurs participant ensemble au paysage épigénétique de celui-ci. Autrement dit, ils forment l'arbre des équilibres (des développements) possibles du processus chaotique qu'est l'ontogenèse toujours en fonction des conditions environnementales. Par conséquent, un gène X ne codant plus une caractéristique possible, son phénotype diffère selon qu'il est présent au sein de l'espèce Y ou Z puisque son expression dépend de l'ontogenèse de l'organisme, contrairement à la première interprétation du néodarwinisme.

Le principe d'auto-organisation permet effectivement de résoudre le paradoxe, mais en éliminant la notion de programme génétique, il retire également aux gènes et à la reproduction leurs rôles principaux dans la construction des êtres vivants, et par conséquent dans la définition du vivant. Pour autant, la notion d'auto-organisation à elle seule n'autorise pas une définition du vivant. Le flocon n'est pas considéré comme vivant, ni les auto-organisations dynamiques telles que les cyclones qui correspondent à des structures dissipatives évoquées dans la section sur le connexionnisme. Afin de différencier les êtres vivants des autres phénomènes d'auto-organisation, Maturana et Varela (1989) proposent le concept d'autopoïèse (du grec *autos* soi et *poiein* produire) : « Un système autopoïétique est organisé comme un réseau de processus de production de composants qui (a) régénèrent continuellement par leurs transformations et leurs interactions le réseau qui les a produits, et qui (b) constituent le système en tant qu'unité concrète dans l'espace où il existe, en spécifiant le domaine topologique où il se réalise comme réseau ».

Bourguin et Stewart (2004) proposent comme modèle minimal de l'autopoïèse un automate de tessellation pouvant être décrit de la manière suivante (Figure I-38) : une cellule dont la membrane se compose d'éléments C baigne dans une solution d'éléments A de concentration fixe k_0 . Selon un taux de désintégration par unité de surface k_1 , les éléments C se transforment spontanément en D qui s'échappent vers l'environnement extracellulaire laissant un trou dans la membrane. Catalysée selon un paramètre k_2 par la surface intérieure de la membrane, la réaction entre deux éléments A donne un élément B . L'élimination des A de cette manière entraîne une différence de concentration des deux cotés de la membrane et génère une pression osmotique amenant les A de l'extérieur vers l'intérieur. Les éléments B peuvent sortir de l'environnement intracellulaire par les trous membranaires uniquement. Cependant, si un élément B rencontre un C qui n'a plus qu'un seul voisin C alors il se transforme en C . Ainsi, la subsistance de la cellule est permise uniquement lorsque les valeurs de k_0 et k_2 permettent de compenser suffisamment la valeur k_1 pour maintenir une régénération perpétuelle. Cet exemple de système autopoïétique permet de préciser quatre points.

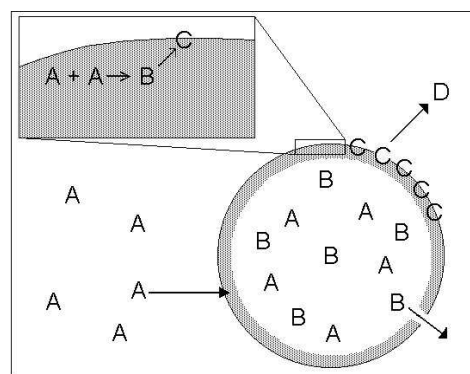


Figure I-38 : Illustration de l'automate de tessellation représentant un modèle minimal de l'autopoïèse (Bourguin et Stewart, 2004).

i - Système de régulation et système autopoïétique

Le premier point porte sur la différence entre un système de régulation et un système autopoïétique. Un système de régulation compense l'influence extérieure sur une variable par rapport à une référence définie, imposée. Puisque la référence est fixe, l'histoire du système se trouve entièrement déterminée par l'extérieur comme les systèmes dépourvus de boucle de retour et entièrement soumis aux conditions environnementales, comme pour les cyclones. En revanche, aucune référence ne guide le comportement des systèmes autopoïétiques et les contraintes extérieures ne déterminent pas totalement leur comportement. L'équilibre d'un système autopoïétique résulte de la rencontre entre deux réactions de nature différente, possédant en commun des composants ou des chaînes causales sur ceux-ci. L'interdépendance de réactions vis-à-vis de ces composants produit une dynamique qui les couple. Dans l'exemple précédant, la première dynamique correspond à la désintégration des C et la seconde correspond à la chaîne de réactions permettant la production de C , bien que celle-ci ne puisse s'effectuer que dans un milieu confiné par des éléments C . Ainsi, la membrane devient une condition nécessaire pour catalyser le métabolisme ($A+A \rightarrow B$) et retenir les métabolites (B) et ce même métabolisme assure la réparation de la membrane. Ici, l'objectif n'est pas le maintien de la membrane ; elle se régénère parce qu'elle se trouve à l'origine du couplage entre deux réactions et par conséquent constitutive de leur l'équilibre. Contrairement à un système de régulation, la rupture de l'équilibre entraîne la dislocation de tout le système. Une machine autopoïétique « accomplit ce processus de remplacement de ses composants parce qu'elle est continuellement forcée de compenser ces perturbations » (Varela, 1989). En ce sens, l'organisation d'un système autopoïétique traduit une situation d'équilibre dynamique vis-à-vis du couplage.

ii - L'identité d'un système

Le deuxième point souligne qu'un système autopoïétique possède une identité propre. Cette notion d'identité s'étaye sur une relation entre structure et organisation qui ressemble à la dualité de propriété évoquée précédemment, en reprenant les termes de Varela (1989) : « L'ensemble des relations qui définissent une machine comme une unité constitue l'organisation » et « l'ensemble des relations effectives entre les composants présents sur une machine concrète dans un espace donné constitue sa structure ». Plusieurs structures de natures différentes peuvent alors correspondre à une même organisation. Plus exactement, se dégagent des systèmes dont l'organisation conserve ses composants (comme une voiture) et des systèmes dont la dynamique d'organisation autorise un renouvellement total de la structure sans changer forcément la nature des composants (comme un tourbillon). Les systèmes autopoïétiques, qui se régénèrent continuellement sous la pression extérieure, appartiennent à cette seconde classe de systèmes. La présence éphémère et l'interchangeabilité avec des éléments extérieurs de l'ensemble des composants pendant la dynamique du système les rendent anonymes vis-à-vis de l'organisation. Autrement dit, l'identification des composants structurels ne suffit pas à identifier l'organisation du système puisqu'ils sont transitoires. Par conséquent, la description componentielle d'un système possédant une identité revient à modéliser tout l'univers susceptible de participer à son maintien, comme pour les motifs des automates cellulaires. En revanche, la description mathématisée centrée sur le système reste possible, car la notion de variables capture cet anonymat, et l'exemple de la cellule se formalise avec des équations aux dérivées partielles spatialisées.

iii - La valeur de l'organisation

Le troisième point rappelle que la notion d'organisation se comprend indépendamment de HOI. En effet, comme le remarque Varela (1989), la notion d'organisation ne prétend nullement à une valeur explicative en soi. Elle traduit un point de vue pour la description de l'invariance d'un phénomène. Ainsi, la définition de l'autopoïèse demeure indépendante des hypothèses métaphysiques.

iv - Le couplage structurel

Pour terminer, le quatrième point revient sur la conservation de l'identité d'un système autopoïétique. Dans le cas de l'exemple de la cellule, le couplage entre réaction interne et réaction externe produit une organisation qui se caractérise par des points d'équilibre selon les conditions environnementales traduisant la fluctuation de la structure. En d'autres termes, la structure détermine l'état du système et le domaine des perturbations admises. Mais l'univers clos de l'exemple ne permet pas d'exprimer tout le potentiel transformationnel d'un système. En effet, les propriétés de compensation induisant la régénération confèrent au système une grande plasticité. Une perturbation externe (introduction de nouveaux éléments) ou interne (une nouvelle réaction apparaissant entre les cellules B à partir d'une certaine concentration) vis-à-vis de la réaction en chaîne habituelle peut entraîner une modification de celle-ci sans impliquer son arrêt. Autrement dit, en dehors des perturbations admises pour la conservation d'une organisation, le couplage entre les réactions internes et les réactions externes entraîne des modifications structurelles qui redéfinissent l'organisation. La conservation de l'identité d'un système autopoïétique passe par le couplage structurel. Ce couplage structurel mis en avant par Maturana (1978) correspond au processus d'ontogenèse et montre bien que les systèmes autopoïétiques sont nécessairement historiques puisqu'ils dépendent à la fois de modifications internes et externes contingentes : « une machine autopoïétique engendre et spécifie continuellement sa propre organisation » (Varela, 1989). Deux remarques s'ensuivent : la première indique que la contrainte de l'ininteruption du couplage pour l'évolution du système concorde avec celle qui motive une architecture de subsomption. La seconde montre que l'organisation ne représente qu'un point stable dans l'histoire du couplage structurel du système et que l'arrêt du couplage traduit la mort du système.

*

En considérant les organismes comme des systèmes autopoïétiques, l'histoire de l'apparition du vivant s'en trouve modifiée. Elle ne coïncide plus avec l'idée que les organismes complexes se sont développés à partir de macromolécules répliquatives puisque, de par le principe d'organisation dynamique, il n'y a pas d'état intermédiaire entre un système classique et un système autopoïétique. Cela signifie donc que seules des conditions environnementales spécifiques (la réunion de composants appropriés et la concaténation de leurs interactions) ont permis la construction de systèmes autopoïétiques moléculaires. Dans cette logique, parmi tous ces systèmes autopoïétiques, des transformations autopoïétiques, des transformations engendrées par un couplage structurel ont abouti à une phase de reproduction dans le cycle autopoïétique. Jusqu'à présent, le discours s'est appuyé sur les organismes unicellulaires, mais il se transpose aisément aux animaux multicellulaires, en les considérant comme des systèmes autopoïétiques dont les composants le sont également. La prochaine section ajoutera toutefois la notion de clôture opérationnelle afin de saisir les conséquences des multiples interrelations dynamiques des couplages entre les composants autopoïétiques.

Cette notion d'autopoièse se trouve radicalement opposée à la première interprétation du néodarwinisme, puisque la reproduction ne devient plus la propriété fondatrice de la vie, bien qu'elle soit primordiale pour sa propagation. Associée à la critique de l'idée du programme génétique, la première forme du néodarwinisme ne peut se retrancher malgré tout derrière les deux positions téléologiques évoquées précédemment. Mais toutes deux tiennent difficilement face à la critique.

La première position, qui imagine que la complexification guide l'évolution, se révèle en fait due à une illusion perceptive. Car, en considérant qu'aucune téléologie ne guide la phylogénèse, la probabilité d'engendrer à chaque reproduction un organisme plus ou moins complexe est de 50%. Néanmoins, la symétrie se trouve brisée par le seuil de la complexité minimale pouvant assurer l'autopoièse. « Ainsi, un processus qui commence juste au-dessus du seuil minimal, et qui à chaque pas de temps donne lieu à deux formes, l'une plus complexe et l'autre moins complexe, produira au cours du temps un ensemble de formes dont la majorité sera relativement peu complexe, mais dont le degré de complexité des plus complexes augmentera progressivement. » (Stewart, 2004).

La seconde position affirme que l'évolution vise à rechercher l'espèce optimale en absolu, pour l'instant l'homme, et que le processus s'effectue de manière linéaire. Deux arguments principaux contraignent à abandonner cette position. Le premier argument porte sur l'idée de pouvoir comparer les espèces entre elles pour décréter laquelle domine, et donc définir une espèce optimale. Les critères de comparaison ne doivent pas concerner les capacités individuelles (intelligence, durée de vie) puisqu'elles ne sont que des moyens pour la victoire, et non le but. La suprématie d'une espèce doit donc se révéler par des critères globaux tels que le nombre d'individus, ou l'étendue de leur présence dans des milieux différents. Cependant, avec ces critères, l'espèce humaine perd de très loin la première place. Un critère plus macabre, la capacité à exterminer les autres espèces en compétition, peut être proposé, mais ce serait oublier d'une part que l'homme restera toujours menacé par des pandémies puisque les bactéries sont indissociables des organismes multicellulaires, et d'autre part que si l'espèce humaine ou une autre espèce vient à détruire toutes les espèces de son écosystème, elle périra également. Le second argument attaque plus particulièrement la linéarité de l'évolution et l'exposer nécessite de revenir plus précisément sur l'histoire du vivant. L'apparition de la vie unicellulaire remonte à 3,5 milliards d'années. Brusquement, il y a 600MA, de nombreuses espèces d'animaux multicellulaires apparaissent. Cette transition s'explique par l'augmentation de la concentration de l'oxygène dans l'eau jusqu'à un certain seuil afin de pouvoir diffuser l'oxygène au centre d'un organisme multicellulaire. Par rapport aux animaux unicellulaires, les organismes multicellulaires possèdent un espace de formes possibles beaucoup plus vaste. Environ une trentaine d'architectures fondamentales appelées « phyla » peut-être imaginée en tenant compte des propriétés topologiques d'un espace tridimensionnel. La découverte de fossiles sur le site de Burgess atteste l'existence d'au moins 13 phyla dont 7 subsistent encore de nos jours. Depuis aucune apparition de nouveaux phyla n'a été observée.

Gould (1991) interprète ces données paléontologiques de la façon suivante : à l'apparition des organismes multicellulaires, l'espace et les ressources disponibles rendaient la compétition quasiment inexistante, permettant ainsi la survie de formes très variées, même si toutes n'ont peut-être pas été explorées pour des raisons de hasard ou des raisons physico-chimiques qui auraient pu limiter certains types d'auto-organisation. Parallèlement, ces communautés, de par leur nombre d'individus devaient être très sensibles aux contingences environnementales. Autrement dit, il existait un facteur aléatoire à la sélection

des phyla. Une fois les principaux phyla constitués et l'exploration de niches écologiques terminée, la compétition pour les ressources favorise la coévolution et peut même aboutir à la suppression de certaines d'entre elles. Une fois cette compétition instituée, aucune nouvelle espèce ne peut apparaître dans cet écosystème puisque toutes les autres sont déjà adaptées pour accaparer les ressources. Ce schéma se retrouve également avec la rapide disparition des dinosaures, probablement due indirectement à la chute d'une météorite il y a 70MA. Les mammifères existant depuis 100MA, tous de petits rongeurs, se retrouvèrent devant un nouvel espace de possibilités. Bien que confinée à l'intérieur de l'architecture des mammifères, une explosion (sur l'échelle de temps de l'évolution) de formes s'ensuivit, de la baleine à la chauve-souris. L'évolution se révèle alors être un équilibre entre espèces, ponctué de grands bouleversements qui rouvrent l'espace des formes possibles soumises à la contingence. Autrement dit, il ne peut plus y avoir une finalité de forme dans le processus évolutif, « le progrès est une fiction collective » (Gould, 1991). Plus précisément, Stewart (2004) explique que la permanence d'une espèce est due au fait que l'organisme adulte possède une organisation autopoïétique viable dans le domaine des perturbations environnementales et que l'ontogenèse de sa progéniture est identique à la sienne, donc robuste, de génération en génération face aux mêmes types de perturbations environnementales. Si l'environnement vient radicalement à changer, ni la survie de l'organisme ni la re-production de l'ontogenèse ne sont assurées, la sélection et l'ontogenèse généreront alors une ouverture des formes possibles.

En somme, l'opposition entre les deux courants néodarwinistes peut être schématisée en deux principales oppositions. La première confronte le programme génétique et l'auto-organisation comme processus de développement. En caricaturant le problème du néodarwinisme par la question : qui apparut le premier, la poule ou l'œuf ? La première interprétation du néodarwinisme penche pour l'œuf alors que la seconde interprétation penche pour la poule. La seconde opposition porte sur la téléonomie. Le premier courant néodarwiniste considère que l'évolution possède une finalité, l'évolution devient alors une progression vers un objectif. Dans cette logique, les individus n'ont pas une finalité propre, hormis la reproduction. À l'inverse, le second courant néodarwiniste ne prête aucune finalité à l'évolution en soi, en revanche, l'individu possède une téléonomie qui est la continuité de son autopoïèse.

Cependant, les arguments présentés autorisent à trancher en faveur du second courant néodarwiniste. Le processus de l'évolution se complète par le processus d'autopoïèse. Les gènes deviennent des facteurs transmis qui vont orienter l'ontogenèse ancestrale si les conditions environnementales sont similaires. Cette conception offre une plasticité théorique du matériel génétique qui permet d'expliquer l'assimilation génétique mise en évidence par les expériences de Waddington (1950) qui, tout en conservant le schéma weismannien, montrent sur la drosophile qu'un même caractère peut-être indifféremment « acquis » ou « inné » et que la transition entre les deux états est relativement aisée (Stewart, 2004).

**

En conclusion, en désavouant le premier courant du néodarwinisme, toutes les approches (Tableau I-9) qui se fondent uniquement sur l'évolution génétique pour concevoir un système cognitif ne couvrent qu'une partie du phénomène, principalement le problème du consensus et de l'heuristique. En effet, avec des exemples en éthologie, Lorenz (1950) insiste bien sur l'intérêt de travailler sur la phylogenèse des comportements et sur leur coévolution avec celle de la morphologie, ainsi que sur les signaux

conventionnels entre congénères, mais il affirme tout autant qu'il faut la différencier de l'étude de la cognition et *vice versa*. L'évolution ne répond donc pas à elle seule au problème de la conception d'un système autonome, néanmoins, cette étude a permis de mettre en avant le concept d'autopoïèse qui offre une alternative au système de régulation qu'induit la notion de besoin du fonctionnalisme écologique. L'interactionnisme revendique cette alternative pour la description de la construction d'un individu.

		Types de famille métaphysique				Hypothèses sur la perception			Hypothèses sur la conception			Types de méthode d'investigation			Types de description		
		FC	FCP	FP	FA	HP1	HP2	HP3	HC1	HC2	HC3	MHD	MC	MA	DP	DM	DC
Evolutionnisme	Évolutionnisme physique	Red	Green	Green	Red	Grey	Grey	Grey	Grey	Grey	Grey	Yellow	Green	Green	Green	Yellow	Green
	Évolutionnisme neuronal	Red	Green	Green	Red	Grey	Grey	Grey	Grey	Grey	Grey	Yellow	Green	Green	Green	Yellow	Green
	Évolutionnisme comportemental	Green	Green	Red	Red	Grey	Grey	Grey	Grey	Grey	Grey	Yellow	Green	Green	Green	Yellow	Green
	Évolutionnisme du langage	Green	Red	Green	Red	Grey	Grey	Grey	Grey	Grey	Grey	Yellow	Green	Green	Green	Yellow	Green

Tableau I-9 : Récapitulation de diverses positions relatives à la grille d'analyse des quatre principales approches provenant du courant évolutionniste. Les cases vertes correspondent à un avis favorable, les cases jaunes à un avis mitigé et les cases rouges à un avis défavorable. Les cases grises signifient qu'une prise de position se trouve hors de propos dans le cadre du paradigme évolutionniste.

3.2.3. L'interactionnisme

L'interactionnisme étudie la coordination des actions et des sensations ainsi que tous les mécanismes y participant. Contrairement au fonctionnalisme écologique qui nécessite la notion de besoin pour découvrir le monde, l'interactionnisme prétend que seule l'activité du sujet dans le monde suffit. Dans ce cas, le couplage sensorimoteur devient le premier mode d'exploration et la notion de besoin devient un rapport au monde particulier se greffant sur la coordination globale des activités du sujet. Sommairement, la cognition se présente alors comme l'ensemble des mécanismes qui dégagent ou font émerger, selon les positions internes de ce mouvement, cette coordination.

Deux approches se proposent d'identifier et de définir ces mécanismes (Figure I-39) : la phénoménologie et le constructivisme. La présentation de la première approche montrera tout l'enjeu et la richesse de l'interactionnisme par rapport aux philosophies évoquées jusqu'ici. Toutefois, l'analyse de ce courant dévoilera une impasse méthodologique qui empêche toute exploitation directe par les sciences de l'artificiel. La seconde approche sera davantage détaillée, en distinguant deux courants qui considèrent la coordination comme résultant, pour l'un, d'un couplage entre les « entrées » et les « sorties » et pour l'autre, d'un couplage par « clôture opérationnelle » notion qui sera définie auparavant. Chacun de ces deux courants conduira à une conception particulière de la notion de représentation qui sera également explicitée. Enfin, la conclusion rappellera les conséquences et les limites sur la cognition artificielle en fonction des différents courants interactionnistes, ainsi que sur les diverses hypothèses qui constituent la grille de lecture de ce chapitre.

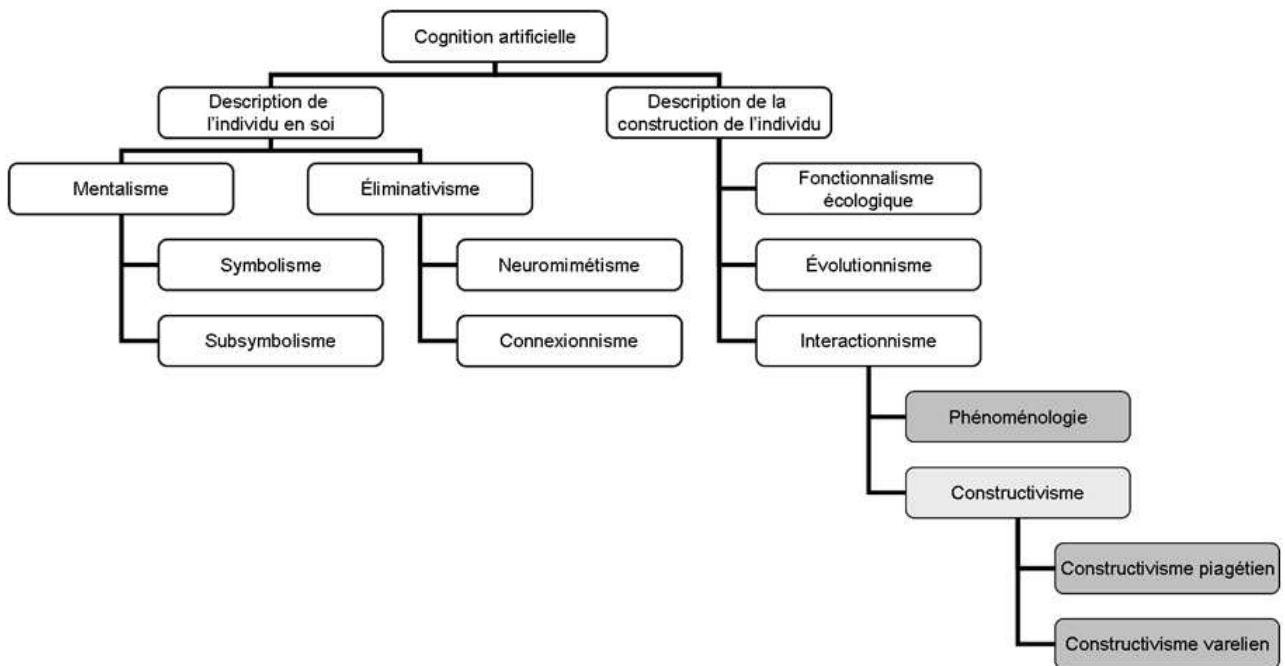


Figure I-39 : Les trois approches provenant de l'interactionnisme au sein de l'arborescence de la cognition artificielle : la phénoménologie, le constructivisme piagétien et le constructivisme varelrien.

A - La phénoménologie

En simplifiant à outrance, les philosophies de la perception évoquées en filigrane au cours de ce chapitre se développent selon l'un des trois schémas types décrivant la relation sujet/objet ou plus exactement sujet/monde (figure I-26). Le schéma A regroupe toutes les philosophies qui attribuent directement ou indirectement au sujet une intuition permettant de saisir ou d'approximer le réel, les objets en soi, les noumènes. Ces philosophies peuvent se rattacher aux familles FCP ou FP. Le jugement existentiel et le jugement perceptif possèdent une signification ontologique. Le symbolisme, par exemple, rentre dans ce schéma. Les philosophies de la famille FC adhérant au schéma B considèrent que les états sensoriels sont suffisamment neutres pour que leurs relations donnent une représentation du monde. Le jugement existentiel devient un truchement pour faciliter cette construction. Le connexionnisme illustre cette position. Le dernier schéma, C, correspond à la phénoménologie, voie médiane entre les deux précédentes.

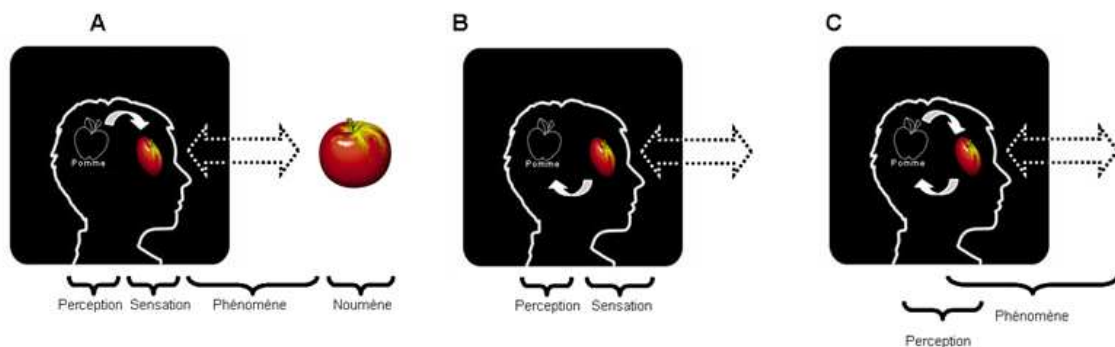


Figure I-40 : Les trois principaux paradigmes types de la perception : le schéma A s'apparente aux philosophies issues de FCP ou de FP, le schéma B s'apparente à celles issues de FC, et le schéma C s'apparente à la phénoménologie.

Certaines philosophies, sous la pression de différentes contraintes, viennent à superposer à divers degrés les schémas A et B, comme le fait le fonctionnalisme écologique. Mais ce procédé n'aboutit pas à une équivalence avec le schéma C, car le sujet et le monde restent bien distincts, et leurs définitions ontologiques demeurent indépendantes. Or, sur ce point crucial, le schéma C, c'est-à-dire celui de la phénoménologie, diffère justement des deux autres. En effet, la conscience saisit immédiatement l'objet mais ce n'est pas une considération distante, auxiliaire. La conscience étant condamnée à toujours être consciente de quelque chose, ce saisissement constitue le sentiment d'existence, le sujet. Le sentiment d'être-là (Dasein) vient nécessairement de cette rencontre. La conscience du sujet se justifie par l'objet qui lui-même n'est ou naît uniquement de la perception suscitée par le monde en fonction de l'intentionnalité du sujet. De ce fait, le phénomène résulte de l'interaction indissociable entre le sujet et l'objet, et c'est de cette indissociabilité que la phénoménologie souhaite partir pour explorer le réel. Ce n'est plus ce qui est perçu qui se trouve au centre de l'étude mais les mécanismes permettant cette perception qui se révèle être, en définitive, la seule réalité immanente.

Le jugement existentiel et le jugement perceptif doivent alors être disséqués et l'introspection constitue l'unique point de vue pour le faire puisque le sentiment d'être, de percevoir, appartient à la sphère privée. Ici se trouve le point de divergence entre la phénoménologie et le constructivisme qui se réclame également du troisième schéma. Avant d'expliquer en détail les raisons de cette divergence, la phénoménologie doit être davantage explicitée dans ses grands principes. La phénoménologie se veut être la science du phénomène ou de l'apparaître qui souhaite déterminer la structure du phénomène ou les conditions générales de l'apparaître. En simplifiant encore, il existe principalement trois formes de phénoménologie qui peuvent aussi se comprendre comme différents stades d'analyse.

La première forme s'attache uniquement à la description ou à l'analyse du donné phénoménal. Elle « se donne pour tâche, non pas d'expliquer le monde ou d'en découvrir "les conditions de possibilité", mais de formuler une expérience du monde, un contact avec le monde qui précède toute pensée *sur* le monde » (Merleau-Ponty, 1948). En effet, certaines hallucinations perceptives forcent à abandonner HP3. Autrement dit, le jugement existentiel influence fortement la perception. Pour dégager le phénomène, l'être véritable, le jugement doit impérativement être suspendu, c'est l'épochè. L'objectif consiste à saisir uniquement le phénomène et d'étudier le jugement perceptif. Cette technique méditative de mise entre parenthèses a permis de mettre en avant l'importance du corps dans la perception. La kinesthésie se révèle alors fondamentale dans l'appréhension des objets et de l'espace et par ailleurs l'indicibilité des expériences sensorimotrices amène à refuser HP1. En ce sens, la phénoménologie sert de préambule à des études psychologiques et même neurobiologiques (Roll, 2003). Malgré tout, les limites de l'épochè montrent que la conscience est toujours consciente de quelque chose et toujours d'une certaine manière. L'épochè se révèle également impuissante face à l'évocation mnémonique incontrôlée, ce qui entraîne le rejet de HP2 et ainsi confirme l'impossibilité du solipsisme. Enfin, cette forme de la phénoménologie appartient à la famille FA, car elle ne nécessite aucune hypothèse métaphysique HOP ou HOI puisqu'elle se veut seulement descriptive et sans concept de vérité associé à ces descriptions subjectives.

La deuxième forme souhaite compléter l'étude descriptive de la première par une étude explicative. En particulier, Sartre (1943) s'appuie sur l'idée que la conscience d'être est une « expérience métaphysique » permettant une dialectique entre l'être et le néant. Ce

faisant, il introduit HOI puisque ladite « expérience métaphysique » lui permet d'acquérir une vérité absolue de la logique, même si c'est la seule.

La troisième forme, la phénoménologie transcendantale défendue par Husserl (1901), veut développer une méthode méditative plus complexe par des mises entre parenthèses successives pour que le phénomène se réduise, s'oriente vers l'interaction entre le sujet avec lui-même, c'est-à-dire qu'elle veut pousser l'introspection jusqu'à ce que le sujet soit face à sa propre conscience, une conscience pure puisqu'indépendante des données empiriques. Cette situation permettrait de dégager les structures essentielles et par conséquent les propriétés essentielles de tout ce qu'il peut connaître. Cette ambition repose également sur l'HOI. Cette phénoménologie diffère du psychologisme qui considère que les concepts, les jugements et les théories sont des événements psychiques, et du logicisme qui considère la logique comme une construction indépendante de la psychologie. Par ailleurs, ces deux courants se reflètent dans les deux courants du symbolisme. Ici, le psychique se différencie de la logique ; néanmoins celle-ci se comprend comme le résultat d'un développement dont les rouages sont à l'origine logiques, et ce sont dans ces rouages que la phénoménologie transcendantale souhaite trouver une vérité apodictique pour bâtir une science pure.

B - Le constructivisme Piagétien

Bien que la phénoménologie soit un tournant dans la façon d'appréhender l'étude du rapport au monde, deux critiques vont montrer les limites méthodologiques de celle-ci. Les critiques exposées se retrouvent notamment dans l'œuvre de Piaget (1965). La première critique, qui porte essentiellement sur l'introspection, se décompose en deux points. Le premier point rappelle qu'un ensemble d'énoncés constitué collectivement sur le monde ou sur le réel est le fruit d'un processus social dont l'objectif est d'obtenir le plus grand nombre d'adhésions pour chaque énoncé. Comme cela a été évoqué dans la première partie de ce chapitre, plusieurs méthodes aident traditionnellement au consensus. Ces méthodes permettent de comparer, et en cela elles sont objectives, des résultats, des faits issus de différents énoncés. Cette comparaison implique la possibilité de quantifier l'observable. Ainsi, la labellisation des énoncés ne dépend plus uniquement de la pression sociale par l'effet du conformisme évoqué dans la section portant sur le connexionnisme. L'extériorisation du discours par ces méthodes autorise le débat scientifique qui produit la norme. Or, la phénoménologie se restreignant à la sphère privée peine à communiquer sur les expériences vécues et arrive encore moins à les comparer. Ce qui est l'objet du discours, le ressenti, appartient toujours à l'intime, le fait se confond avec la norme. Le débat scientifique devient impossible et la construction collective d'un énoncé cohérent résulte uniquement d'un processus psychosocial. Une phénoménologie ayant pour thème les relations sociales peut-être entreprise, elle pourra encore être source d'inspiration pour la psychologie ou psychologie sociale. Néanmoins, cette seule phénoménologie, centrée sur elle-même, ne générera pas davantage un débat scientifique. Le deuxième point de la première critique souligne que l'étude à la première personne portera forcément sur des mécanismes déjà existants du jugement perceptif. Autrement dit, une certaine caractérisation de la genèse de la perception peut être envisagée mais non la genèse des mécanismes qui la produisent puisque le sujet est déjà adulte.

La seconde critique attaque plus particulièrement les deux dernières formes de la phénoménologie. Le sujet transcendantal, c'est-à-dire l'interaction sujet/sujet, résulte d'une évocation du monde pour atteindre l'intuition pure. Mais le sujet transcendantal reste un sujet et le sujet-objet n'est qu'une projection ; aussi, l'intuition correspond à l'activité du sujet dans cette interaction. Se fonder sur cette intuition revient à une forme de psychologisme

(Piaget, 1965). Par exemple, dans l'énoncé : « Pourquoi je pense à ce que je pense ? » le premier « je » englobe le sujet au-delà de sa conscience alors que le second « je » désigne uniquement le sujet conscient. Le problème de l'autoréférence sera exposé plus en détail dans le prochain chapitre, mais proposer sa résolution par l'existence de l'intuition pure et en dégager les principes logiques absolus n'est pas raisonnable : « ou bien la logique est suspendue à l'intuition d'un sujet transcendantal, et elle n'est plus absolue (ce qu'on désirait qu'elle soit), ou bien elle est absolue et il n'est plus besoin d'une intuition transcendantale. » (Cavaillès cité par Piaget, 1965). Cette critique a pour conséquence de rejeter HC2 mais elle n'empêche pas HOI ; les objets idéels sont simplement inaccessibles en soi.

En somme, le discours phénoménologique porte sur la perception phénoménale qui constitue le réel, autrement dit, la seule réalité devient la perception phénoménale et non plus le monde perçu. L'introspection sur cette interaction entre soi et le monde, en évacuant les idées sur celui-ci, offre des sources d'inspiration pour les sciences cognitives. Cependant, les discours à la première personne sont toujours imprégnés du vécu du sujet et par conséquent aucun modèle de cognition artificielle ne peut être imaginé. Le constructivisme va tenter de comprendre les zones d'ombre, mises en évidence par la phénoménologie, grâce à un discours à la troisième personne qui autorisera l'emploi de la méthode artefactuelle.

Le constructivisme s'oppose toujours au schéma A et B. Le sujet ne connaît l'objet qu'en agissant sur lui et aucune affirmation ne peut être formulée avant cette action, de sorte que le sujet le découpe en objets particuliers uniquement au cours de ses actions et par l'interaction entre l'organisme et le milieu. Ce constat est commun avec celui de la phénoménologie, le jugement existentiel est relatif au sujet et à son histoire. Mais là où la phénoménologie va repousser le jugement existentiel et donc s'enfermer dans l'introspection qui l'empêche d'étudier la nature des mécanismes sous-tendant la perception, le constructivisme va admettre que le jugement existentiel est fallacieux, tout en acceptant qu'il participe au choix de l'action et donc de l'acquisition de la connaissance.

Ainsi, imaginé un autre comme moi-même interagissant avec le monde devient un moyen de connaissance sur soi, bien que cette connaissance soit toujours dépendante de soi. Cette problématisation de la relation sujet/objet qui devient autrui/objet permet d'imaginer un point de vue sur l'interaction entre sujet/objet et d'explorer ainsi la zone d'ombre d'un point de vue à la première personne, portant uniquement sur l'impression laissée par l'interaction, le phénomène. Dans cette nouvelle situation, le vécu à la première personne de l'autre reste inaccessible mais qu'importe, puisque le sujet étudie les mécanismes à l'origine de ce vécu et non le vécu lui-même. Le changement de perspective a aussi pour conséquence d'introduire l'espace et le temps comme outils d'expertise, contrairement à l'introspection qui observe la sensation d'espace et de temps tout en étant toujours ici et maintenant. Cependant, la valeur des concepts d'espace et de temps reste de nature pratique.

Dans l'acceptation qu'une analyse à la troisième personne permette de compléter la connaissance de soi, Piaget (1965) définit l'épistémologie génétique comme « une recherche essentiellement interdisciplinaire, qui se propose d'étudier la signification des connaissances, des structures opératoires ou des notions ». En considérant que les autres humains sont des semblables, un autre comme soi-même ou soi-même comme un autre, qui diffèrent toutefois par mille et une choses mais qui se ressemblent par le sentiment supposé d'un vécu, une étude empirique sur ces structures opératoires devient alors possible. La psychologie expérimentale se révèle cruciale et plus particulièrement la

psychologie génétique qui ne se résume pas à repérer et à répertorier les états de développement de la raison, mais qui vise à expliquer expérimentalement le processus de construction des structures logiques fondamentales.

Les principaux travaux de Piaget (1965) en psychologie génétique portent sur les opérations logicomathématiques. Il a mis en évidence par exemple qu'un enfant de 4-5 ans considère qu'il y a moins de jetons dans un paquet de 10 que dans deux paquets de 4 et de 6 résultant d'une séparation, la réversibilité opératoire qui permettrait de percevoir l'égalité n'étant pas comprise. A partir de ces résultats il conclut : « L'irréversibilité est liée à la conscience du sujet individuel qui, centrant tout sur l'action propre et les impressions subjectives qui l'accompagnent, est entraîné par le flux des événements internes et externes et dominé par les configurations apparentes ; au contraire, la découverte de la réversibilité opératoire marque la constitution du sujet épistémique qui se libère de l'action propre au profit des coordinations générales de l'action, c'est-à-dire de ces « formes » permanentes de réunion, d'emboîtement, d'ordination, de correspondance, etc., qui relient les actions les unes aux autres et constituent ainsi leur substructure nécessaire. » Ce faisant, Piaget (1965) montre que, tout en appartenant à la même mouvance de pensée qu'Husserl (1901) (schéma C), une méthodologie à la troisième personne peut fournir des énoncés objectifs sur la genèse de la logique, le terme « objectif » étant compris comme la possibilité d'extérioriser un énoncé par une expérience ou une simulation qui puissent être comparées. Mais surtout, de la notion de réversibilité vont dériver successivement les trois thèses centrales du constructivisme piagétien, qui seront successivement détaillées

i - Les trois thèses centrales du constructivisme piagétien

La première thèse propose que les opérations cognitives se construisent selon l'alternance entre deux phases : l'assimilation et l'accommodation. Ce processus de construction ressemble au dilemme rencontré dans le fonctionnalisme concernant la généralisation d'une loi et l'incorporation d'exceptions. En répétant les expériences, les opérations se généralisent sur certains types d'expériences et marquent ainsi une nouvelle discrimination avec d'autres expériences et d'autres opérations. L'assimilation et l'accommodation font référence indirectement à ces deux principes de généralisation et de discrimination. Dans ce schéma, l'action joue un rôle primordial car l'exploration alimente ce processus. De même l'anticipation (ou la prédiction ou plus simplement l'activité mnésique autorisant la comparaison) a une place importante puisqu'elle seule permet la détection d'invariants qui est un premier pas vers la compréhension de la réversibilité. L'anticipation n'est plus un auxiliaire ou une amélioration pour le processus cognitif comme dans le fonctionnalisme écologique. En définitive, une opération résulte de l'équilibre entre l'assimilation et l'accommodation qui dépend de l'histoire du sujet puisque chaque action du sujet apporte de nouvelles expériences, tant sur l'objet que sur l'action, qui infléchiront les choix des actions futures. Par ailleurs, une opération interagit, se coordonne avec les autres schèmes d'action, et cela ne se réduit pas à la juxtaposition de stratégies gagnantes. L'ensemble de cette coordination propose une nouvelle perception, c'est-à-dire de nouveaux objets pour la pensée et l'action. Ce nouvel horizon issu de l'équilibration des opérations introduit la seconde thèse.

Une structure, comprise comme un ensemble coordonné d'opérations, offre un nouvel horizon d'exploration et de combinaison qui peut conduire à générer une nouvelle structure opératoire en fonction des contraintes exogènes et des propriétés endogènes induites par les structures précédentes. *La seconde thèse* spécifie alors que structure et genèse forment le cycle du développement (Piaget, 1964). Une structure nécessitant et intégrant

obligatoirement les sous-structures, le développement cognitif se décrit inévitablement par des stades ordonnés. Toutefois, il faut souligner que le développement cognitif se trouve limité par les capacités d'agir et par la plasticité du substrat biologique. À noter, par ailleurs, que la notion de subsomption formulée par le fonctionnalisme écologique est compatible avec cette position.

La dernière des trois thèses découle directement de la précédente : la nouvelle structure dépendant toujours de la structure antérieurement constituée, le développement d'un individu est relatif à son histoire. Selon cette position où la connaissance commence par l'action, la première structure opératoire doit concerner les capacités sensorimotrices. Dans l'idée de Piaget, le nouveau-né possède une activité motrice principale endogène qui permet d'explorer le champ des possibles, cette activité initiale s'appuyant sur des boucles sensorimotrices. Le développement de structures opératoires se dégage au fur et à mesure de l'assimilation des régularités des relations sensorimotrices et de l'accommodation de leurs variétés. L'expérience menée par Auvray et al. (2005) avec un dispositif de substitution sensorielle montre bien que ces relations sensorimotrices ne se limitent pas à relier l'activité motrice aux sensations tactiles mais également à toutes les autres modalités sensorielles, et que cette phase demeure surtout un incontournable pour appréhender le monde comme composé d'objets.

Plus précisément, la compréhension des implications de ce paradigme concernant l'appréhension du monde passe par une étude succincte de la notion de représentation : son rôle, sa nature et sa manipulation. Le terme « représentation » se définit soit comme une évocation mentale d'un objet absent, soit comme une reconstruction mentale de l'objet qui se compare avec celui qui est présent. La légitimité de la première définition n'est pas à démontrer. En revanche, la seconde définition, qui correspond à une position symboliste, se révèle intenable face aux expériences psychophysiques. Plus particulièrement, trois types d'expériences psychophysiques méritent un développement afin de compléter les critiques théoriques antérieures sur le symbolisme qui considère la perception comme une reconstruction du monde à l'intérieur de l'esprit où les décisions ou les manipulations s'effectueraient à partir de cette représentation.

ii - Trois types d'expériences s'opposant à la perception comme une reconstruction interne du monde

Le premier des trois types d'expériences perceptives exposés ici provient des travaux d'Yarbus (1967). Ceux-ci montrent d'une part que le balayage oculaire est permanent et d'autre part que la stratégie de celui-ci dépend de l'intentionnalité de l'observateur (Figure I-41). Le fait que l'activité dépende de la tâche que se fixe l'observateur contredit les approches strictement ascendantes du symbolisme ou subsymbolisme pour qui le processus de reconstruction est automatique. Mais cette expérience met à mal également le connexionnisme radical. Celui-ci considère en effet que les mécanismes de la perception, les saccades oculaires incluses, se fondent exclusivement sur les propriétés statistiques des images extraites du monde. Ainsi, selon cette position, la stratégie d'observation dépendrait uniquement de l'image et non de la tâche d'observation, ce qui n'est pas le cas. Par ailleurs, le caractère cyclique du balayage souligne le caractère incrémental de la vision. Cette expérience conforte à la fois la position écologique qui s'appuie sur les données concrètes du monde plutôt que sur une reconstruction mentale et à la fois une position phénoménologique pour qui la perception dépend de l'intentionnalité du sujet donc de la façon d'aborder le monde.

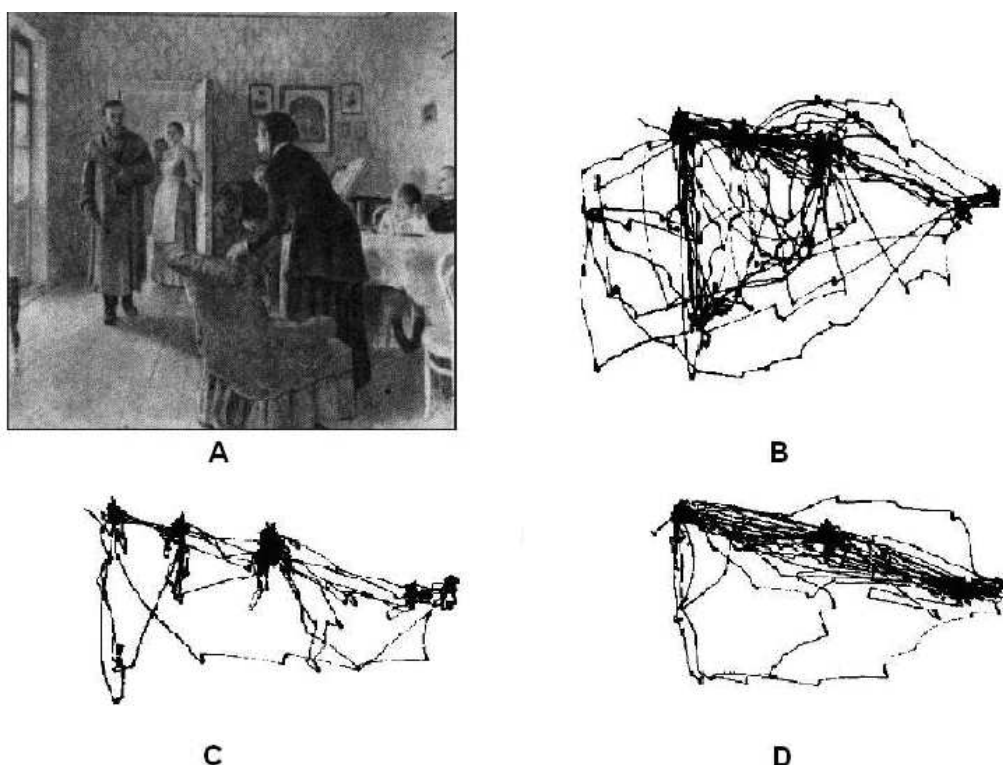


Figure I-41 : Balayage du regard d'un sujet regardant l'image A selon la tâche assignée : examiner librement (B), évaluer l'âge des protagonistes (C) et estimer la durée d'absence de la personne qui arrive (D) (Yarbus, 1967)).

Le second type d'expériences perceptives est d'ordre anatomique et psychophysiologique. Sur la rétine, à cause de l'intersection entre le nerf optique et le globe oculaire, se trouve une zone dépourvue de cellules réceptrices (Mariotte, 1660). Pourtant, même en vision monoculaire, cette zone aveugle n'apparaît pas comme une tache noire dans le champ de vision. Une première raison serait d'avancer que le cerveau compense grâce à une reconstruction interne du monde mais une autre raison serait qu'un sujet ne voit pas ce qui n'est pas situé dans son champ de vision. En d'autres mots, le sujet ne perçoit pas sa rétine puisque c'est précisément elle qui lui permet de voir. Cette seconde raison plus simple et cognitivement la plus économique remet en cause la perception comme une reconstruction du monde.

Par ailleurs, il est possible d'observer son effet en regardant, par exemple, uniquement avec l'œil gauche le rond de la Figure I-42 à environ 25 cm et de s'apercevoir de la disparition du carré. Si la perception était une reconstruction d'un monde à partir des sens, la tache aveugle aurait dû être compensée. Mais aucune compensation n'a lieu parce qu'il n'y a rien à compenser, la tache aveugle ne participe pas à la perception et n'interfère pas dans cette dernière. La problématique de la tache aveugle illustre parfaitement la nécessité du passage de la première à la troisième personne pour décrire certains mécanismes de la perception. À la première personne, si ce phénomène est remarqué, celui-ci paraît étrange par rapport aux autres tout en étant régulier puisqu'il ne s'inscrit dans aucune intentionnalité particulière. Le phénomène induit par la tache aveugle devient incompréhension, mystère. Mais en problématisant ce phénomène dans le cadre de la relation autrui/objet, récepteur/émetteur, il trouve une place et une explication cohérente.



Figure I-42 : Dispositif visuel qui permet d'observer l'effet de la tache aveugle en fixant uniquement avec l'œil gauche le cercle à 25 cm environ et de s'apercevoir de la disparition du carré gauche.

Le troisième type de données expérimentales d'importance dans la critique de la représentation comme une reconstruction interne provient des travaux sur la cécité au changement (O'Regan, 1992 ; Simons, 1996 ; Auvray et O'Regan, 2001). Un observateur ne remarque pas les changements entre deux images successives entrecoupées par un distracteur, bien que les différences soient flagrantes lors de la visualisation simultanée des images (Figure I-43). Le participant ne mémorise donc pas l'ensemble des informations avec exactitude puisqu'il est incapable de comparer avec précision deux vues successives d'une scène visuelle. Pourquoi la perception semble-t-elle globale et parfaite ? Parce que le monde est disponible, toujours là et que, par conséquent, le sujet doit seulement toujours posséder un index (ou pointeur) vers les informations accessibles ainsi que le moyen d'y accéder, « le monde comme une mémoire externe » (O'Regan, 1992). La représentation n'est pas une reconstruction du monde mais une construction perpétuelle du monde à l'aide d'index mentaux qui proposent toujours soit de déchiffrer les réponses en accord avec leur questionnement dans les sensations, soit de trouver les réponses à certains questionnements. Ces index mentaux peuvent s'interpréter comme un réseau de références déictiques qui renvoie à une manière particulière et directe de faire référence à un objet (Ballard et al., 1995).

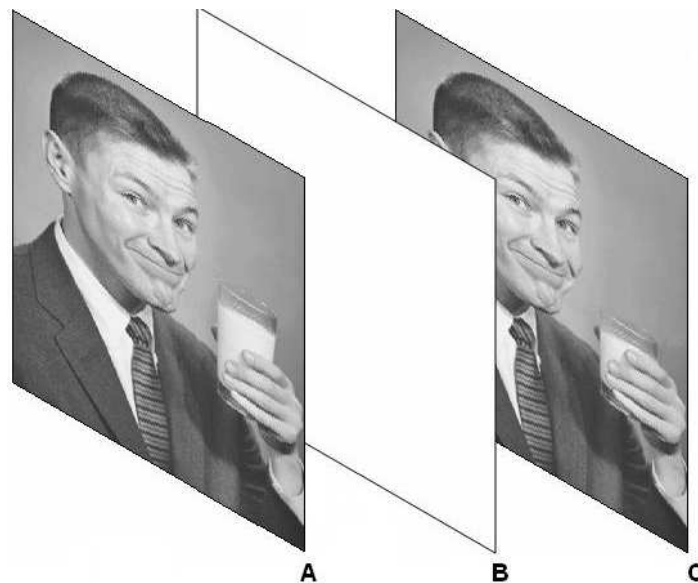


Figure I-43 : Exemple d'un dispositif mettant en évidence la cécité au changement, l'image A apparaît puis l'écran blanc B et enfin l'image C. La différence de hauteur du verre de lait n'est pas remarquée, bien que celle-ci soit évidente lorsque les deux images se trouvent côte à côte (Rensink et al., 1997).

Cet index constitue une disposition mentale pour construire un objet grâce à une unité mnémotique multimodale (sensorielle, motrice) à résolution variable qui est toujours reliée à d'autres index, ce qui autorise ainsi une exploration à la fois incrémentale, d'index en index, et à la fois toujours à propos, guidée. L'ensemble de ces particules d'information forme un espace avec certaines applications. L'attention, par exemple, peut se comprendre

comme l'ensemble des particules participant à la perception de la scène, et l'arrière plan attentionnel inconscient, le pré-attentionnel, correspond à l'ensemble des particules d'informations disponibles, à l'affût, c'est-à-dire des particules avoisinantes à celles impliquées dans la perception en cours. Conformément à la position piagétienne, l'intentionnalité se définit alors comme la conscience de la direction de l'acte et de ses attentes. Selon ce paradigme, les hallucinations deviennent des constructions d'objets dont les informations a priori ont fait la majeure partie du travail mais qui s'effondrent sous la pression d'un réseau d'objets davantage fondé sur les sens.

iii - Débat sur la nature des représentations d'objet absent

Ces trois types de données discréditent le rôle de la représentation comme reconstruction dans la perception du monde, mais la nature des représentations reste ouverte lorsque l'objet est absent. Deux grandes familles sur la nature des représentations cohabitent. La première famille s'occupe des représentations comme une évocation sensorielle de l'objet absent. Les données expérimentales indiquent que les manipulations de ces objets mentaux possèdent des caractéristiques similaires à celles effectuées sur des objets perçus. Par exemple, l'expérience de Finke et Pinker (1982) s'appuie sur le dispositif suivant (Figure I-44) : les participants regardent une configuration de quatre points aléatoires pendant cinq secondes. Après une seconde, une flèche apparaît sur un fond blanc. Les participants doivent ensuite signaler si la flèche pointait vers l'un des points. Sans qu'il soit demandé explicitement d'utiliser une image mentale, les résultats montrent que le temps de réponse est proportionnel à la distance entre la flèche et le point. Des expériences sur d'autres aspects comme les opérations selon la résolution d'images (Cornoldi et al., 1989), la rotation d'image mentale (Shepard, 1975), ou la comparaison de distance (Denis, 1997) font également état de propriétés intrinsèques spatiales des images mentales.

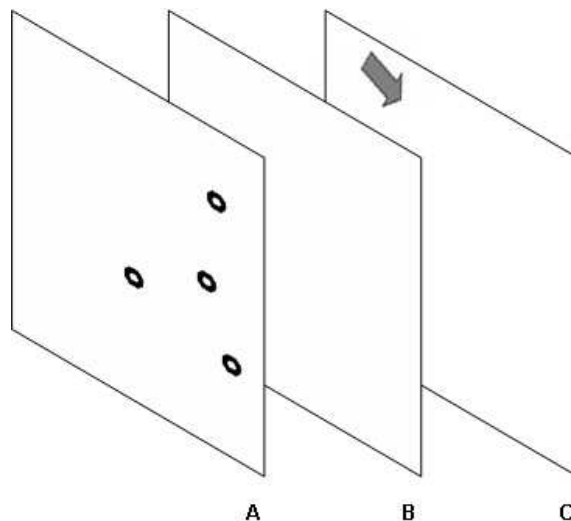


Figure I-44 : Exemple de dispositif type de l'expérience de Finke et Pinker (1982). Après la présentation successive des trois images le sujet doit évaluer si la flèche pointe vers l'un des ronds précédemment vus. L'expérience montre une corrélation entre les temps de réponse et la distance entre la flèche et le rond pointé.

Pour interpréter ces résultats, deux théories s'affrontent dans le cadre de la première, la famille de représentation reposant sur la notion d'image mentale. La première théorie,

défendue par Kosslyn (1980), suggère que les images mentales sont des projections dans une mémoire à court terme, qu'il appelle le *buffer* visuel, d'images élaborées à partir d'informations compressées sous forme propositionnelle se situant dans une mémoire à long terme. Autrement dit, une représentation s'identifie à une reconstruction de l'image du monde et par conséquent l'exploration de cette représentation utiliserait les mêmes opérateurs que ceux utilisés dans le processus de la perception du monde réel, ce qui expliquerait les propriétés communes entre perception et manipulation d'images mentales. En d'autres termes, les représentations correspondent à une reconstruction des sensations déjà vécues ou, autrement dit, à des images sensorielles remémorées. Ainsi, les traitements perceptifs s'appliquent sur ces représentations comme s'il s'agissait d'une image sensorielle réellement vécue, bien que ces représentations puissent être moins précises et moins vives. Ce sont des représentations iconiques qui restent compatibles avec le fonctionnalisme écologique dans le sens où elles peuvent très bien dépendre du mode d'acquisition, de l'action, etc. En revanche, cette position se heurte aux conclusions des expériences désavantageant la thèse de la perception comme reconstruction interne et plus particulièrement par rapport à la notion d'espace qui est la toile de fond du paradigme piagétien. La seconde théorie se résume par la conception de l'espace subjectif de Poincaré (1902) pour qui « l'espace représentatif, sous sa triple forme visuelle, tactile et motrice, est essentiellement différent de l'espace géométrique. [...] Nous ne nous représentons donc pas les corps dans l'espace géométrique, mais nous raisonnons sur ces corps comme s'ils étaient dans l'espace géométrique. ». Autrement dit, les représentations iconiques ne sont pas placées dans un espace géométrique avec ses opérateurs ce qui reviendrait à séparer l'espace et les objets pour reconstruire mentalement un univers pour le percevoir de nouveau. L'argument confortant ce point ressemble à celui avancé pour l'illusion d'une perception globale et parfaite. La représentation ressemble à une reconstruction interne totale parce que l'attention sait toujours quel chemin prendre pour accéder à l'information. Cette potentialisation donne l'illusion de la cohérence et de la totalité de la représentation ou de la perception. Mais surtout, en considérant la représentation comme une perception « à vide », la forme des représentations est nécessairement fondée sur les lois sensorimotrices. Les représentations possèdent toute la potentialité de simulation offerte par leurs lois sensorimotrices dégagées par l'interaction entre le sujet et l'objet.

La seconde famille de représentations mentales, les modèles mentaux, ne lie plus les représentations à des évocations sensorielles se rapportant à un objet mais à l'ensemble des concepts associés de près ou de loin à ce dernier, avec toutes les propositions qui peuvent y être appliquées. Cet ensemble possède une structure hiérarchique permettant des raisonnements sur plusieurs niveaux. Johnson-Laird (1983), le principal défenseur de ces modèles mentaux, considère que les concepts dépendent du vécu et de l'incarnation du sujet mais également de son environnement socioculturel. Ces représentations, sous forme verbale, autorisent des raisonnements sur la configuration spatiale sur des cas simples. Tant que ces raisonnements restent analytiques, aucune similitude avec des propriétés spatiales n'est observée dans les mécanismes de la pensée. Mais lorsqu'une description verbale génère une image, les manipulations effectuées possèdent alors les mêmes propriétés que celles des images mentales produites par remémoration. Ainsi, il existe une traduction dans les deux sens entre ces deux types de représentation.

Au sein de cette famille, peut se développer une branche plus radicale, une position subsymbolique (Pylyshyn, 1984). Elle considère la cognition comme une machine résolvant des problèmes de manière analytique et de ce fait les représentations élémentaires sont amodales, c'est-à-dire que les représentations modales deviennent au mieux des

intermédiaires, au pire des leurres. En somme, rien n'est construit ou reconstruit, tout est calculé. Sans revenir sur les critiques déjà évoquées à l'encontre du subsymbolique, ni sur les expériences qui montrent indubitablement des caractéristiques spatiales des représentations, une critique supplémentaire peut être faite sur l'hypothèse que tout puisse se résoudre analytiquement. En effet, la position de trois astres ne peut se calculer analytiquement au bout d'un temps t (Poincaré, 1902). Seul le calcul pas à pas, en déroulant les équations différentielles qui s'imbriquent, permet de simuler la trajectoire des trois astres jusqu'au temps t . L'idée d'une cognition analytique fondée sur des informations amodales se trouve restreinte à évoluer dans un monde ne prenant pas en compte la dynamique des choses, or ce n'est pas le cas pour le sujet humain. En rejetant l'idée d'un espace géométrique a priori, la simulation sensorimotrice se révèle être le seul moyen pour raisonner intuitivement sur la cinétique des objets.

iv - L'acquisition de la connaissance selon le constructivisme piagétien

Mais cette catégorisation des représentations, entre images mentales et modèles mentaux, reste une grande simplification face aux divers types de mémoire répertoriés par les psychologues : mémoire épisodique, mémoire sémantique, mémoire phonologique, mémoire de travail par exemple. Par ailleurs, ces réseaux mnémoniques fonctionnent tous en interactions mutuelles, ce qui complexifie l'étude de chacun individuellement. Toutefois, toutes ces informations mémorielles se rattachent d'une manière ou d'une autre, ne serait-ce que par leur acquisition, au vécu et par conséquent au sensorimoteur. Apprendre que « Kiev est la capitale de l'Ukraine » ne repose pas a priori sur des aspects sensorimoteurs hormis pour l'encodage de la proposition, mais en fait, savoir que « Kiev est la capitale de l'Ukraine » dépend de tout un réseau sémantique dont les fondations reposent sur des informations sensorimotrices. Par conséquent, le constructivisme rattache très indirectement les informations abstraites telle que « Kiev est la capitale de l'Ukraine » au sensorimoteur.

En effet, selon Piaget (1964), l'extérieur et l'intérieur n'ont pas d'autre signification, à l'origine, pour le nourrisson qu'un mélange de sensations affectives, cénesthésiques et kinesthésiques, autrement dit des événements globaux liés à des proprioceptions diffuses. Dans cet univers primitif, l'intentionnalité (l'objectif des actions et leurs réalisations) ne diffère pas du soi, c'est un solipsisme pratique. La construction de l'objet débute avec la corrélation de certaines données sensorielles avec l'activité immédiate et subjective du sujet. La coordination sensorimotrice isole alors l'objet, perçu comme quelque chose d'indépendant de soi et autorisant certaines actions. Cette relation entre les objets et les actes disponibles au sujet forme la notion d'espace. Dans un premier temps, le sujet construit un espace pratique constitué par l'ensemble des propriétés coordonnant l'action vers l'objet, un référentiel d'activités égocentrées. Dans un second temps, l'espace devient une propriété des choses, définissant ainsi un univers dans lequel tous les déplacements se situent, y compris ceux du sujet : le référentiel devient allocentrique.

Cette construction de l'espace selon un référentiel allocentrique se révèle déterminante pour l'activité du sujet puisqu'il va pouvoir alors sortir de sa perspective propre (à la première personne), c'est-à-dire qu'il pourra se projeter comme un objet parmi d'autres et ainsi se comprendre dans l'espace en mettant en relation ses propres déplacements avec l'ensemble des autres déplacements. Les relations spatiales entre les objets deviennent indépendantes du sujet. Ce bref résumé sur le développement ne décrit que le tout début de la dynamique engendrée : genèse puis structure ouvrant sur une nouvelle genèse et ainsi de suite. Les dimensions sociales et culturelles doivent également

être prises en considération dans ce processus, bien que les proportions entre les facteurs endogènes et exogènes restent l'objet d'un grand débat. Néanmoins, l'étude des premières étapes de la construction de la notion d'espace se révèle incontournable pour une conception robotique de l'interactionnisme.

Par ailleurs, une autre manière de raconter les étapes successives de la construction de l'espace peut être entreprise, de sorte à faciliter l'analogie avec un système robotique. Le sujet découvre un espace sensorimoteur, le terme espace étant compris comme un ensemble sur lequel sont définies des opérations. En plus de la principale distinction entre senseur et effecteur, les dimensions de cet espace ne sont pas homogènes de par la nature des capteurs et de par leur localisation. La dépendance entre les dimensions se trouve conditionnée par l'environnement et le corps, mais également par les boucles sensorimotrices considérées comme des automatismes ou des réflexes innés qui constituent et orientent l'exploration de l'espace sensorimoteur. Par exemple, les sensations motrices doivent être nécessairement plus étroitement corrélées ou liées avec les sensations tactiles qu'avec les sensations olfactives, sans être pour autant inexistantes. Le sujet, dans un premier temps, met à l'épreuve ses boucles sensorimotrices afin de reconnaître les opérations stables, reproductibles et efficaces grâce à la vérification de prédictions internes ou à la compensation des modifications extérieures ou intérieures. Les opérations ainsi dégagées reflètent les liens entre les dimensions sensorimotrices. Ces schémas d'action fortement liés aux situations sensorielles constituent les premières affordances permettant un rapport égocentrique avec le monde. S'il y a simulation interne des actions, cela s'effectue obligatoirement à la première personne.

Dans un second temps, le sujet intègre certaines opérations qui se révèlent réversibles. Cette phase semble être obligatoire pour accéder à un référentiel allocentrique. Piaget (1965) inspiré par les travaux de Poincaré (1902) l'interprète comme une sorte de structuration d'un groupe sensorimoteur au sens mathématique. En effet, la réversibilité d'une opération se trouve au cœur de la notion de groupe, qui est un ensemble G muni d'une opération $*$ qui satisfait les axiomes suivants :

1. *la composition interne* : quels que soient x et y , éléments de G , $x*y$ appartient à G ;
2. *l'associativité* : quels que soient x , y et z , éléments de G , $x*y*z=x*(y*z)=(x*y)*z$;
3. *l'identité ou l'existence d'un élément neutre* : il existe un élément noté e de G tel que, quel que soit x , élément de G , $x*e=e*x=x$;
4. *la réversibilité* : pour tout élément de x de G , il existe un élément x' de G tel que $x*x'=x'*x=e$.

Ce type de structure permettrait alors de caractériser l'espace ou des sous-espaces sensorimoteurs offrant un moyen d'imaginer l'action et de concevoir de nouvelles opérations au lieu de simplement améliorer celles générées par les boucles sensorimotrices. Progressivement, la constitution de ces groupes sur les actions puis sur les opérations et ainsi de suite aboutirait à la construction d'un sous-espace sensorimoteur allocentrique permettant des projections internes (ou des simulations) à la troisième personne. La prédiction se transforme en anticipation puisque le sujet et ses objectifs se situent dans le même espace géométrique imaginé. Les travaux d'Auvray (2004) montrent l'importance de l'activité sensorimotrice sur la constitution de la notion d'objet grâce à la substitution sensorielle, deux autres types de travaux qui seront présentés successivement apportent un éclairage tout particulier sur ce paradigme.

v - Modélisation et études expérimentales avec paradigme sensorimoteur

Concernant le premier type de travaux, Philipona (2003) propose une formalisation de cette théorie sensorimotrice dans un cadre élémentaire, c'est-à-dire, en posant des hypothèses fortes qui ne permettent pas de la transposer directement dans le réel ; elles incluent entre autres l'hypothèse que l'univers possède uniquement des transformations rigides et qu'il existe une relation fonctionnelle instantanée entre les entrées sensorielles et les sorties motrices. Dans ce cadre, Philipona (2005) a développé une méthode pour dégager une partie du groupe de la loi sensorimotrice qui représente ici une fonction reliant les entrées sensorielles aux sorties motrices et à la configuration de l'environnement. Le groupe se constitue à partir des transformations univoques de l'ensemble des couples formés par les sorties motrices et les configurations environnementales qui laissent la loi sensorimotrice invariante. Ainsi posé, ce groupe se trouve indépendant du codage sensoriel et conserve sa structure quel que soit le code moteur. Toutefois, les caractéristiques du groupe restent suffisamment riches pour identifier une loi sensorimotrice à un code sensorimoteur près. En d'autres termes, « à isomorphisme près, les propriétés de la loi sensorimotrice fonctionnelle indépendantes du code sensorimoteur sont précisément les propriétés de ce groupe » (Philipona, 2005). En se fondant sur ces résultats théoriques, la simulation dans un univers virtuel de l'interaction d'une tête de rat (disposant de capteurs visuels, auditifs et tactiles) avec son environnement permet, entre autres, de distinguer les générateurs de translation de ceux de rotations pures ainsi que d'établir une commensurabilité entre eux.

Le second type de travaux, menés par Afonso (2006), porte sur la « plasticité » comportementale de l'individu dès lors qu'il est privé dès la naissance de l'un de ses sens, en l'occurrence, la vue. Plus spécifiquement, Afonso étudie la façon dont cet individu peut se représenter mentalement l'espace dans lequel il évolue, et la façon dont il peut bénéficier des informations provenant de ses autres sens. Il s'agit de comprendre si des étapes indispensables dans la construction des représentations mentales opéraient dès le plus jeune âge, et dans quelle mesure la plasticité comportementale / cérébrale peut pallier cette privation visuelle précoce ; ou si malgré cela, la privation totale d'un sens affecterait définitivement la qualité des représentations mentales. L'originalité du dispositif expérimental repose sur la transposition des expériences classiques de manipulation de représentations mentales avec un apprentissage passif de la configuration initiale (verbale, tactile ou visuelle) dans un environnement avec des sources sonores virtuelles spatialisées qui englobent totalement le sujet et autorisent, par ce fait, une exploration active de la configuration initiale à mémoriser. Les résultats d'Afonso (2006) permettent, entre autres, de répondre à trois questions.

Première question, la privation visuelle précoce affecte-t-elle de manière définitive la qualité des représentations mentales ? Non, l'isomorphisme structural des représentations apparaît également avec les aveugles de naissance. Effectivement, dans le cadre d'une théorie sensorimotrice, les lois sensorimotrices ne dépendent pas du codage sensoriel, les propriétés dégagées traduisent les régularités et les invariants de transformation. Ainsi, des sens suffisamment riches, bien que de nature différente, peuvent avoir des propriétés communes à un isomorphisme près. Sur ce point, Philipona (2005) souligne que cette ressemblance permettrait « au cerveau de réinterpréter la stimulation tactile pour l'extérioriser [...] et retrouver une perception plus proche du visuel que du toucher », en faisant référence aux expériences de substitution sensorielle de Bach-y-Rita (1969) qui transposait des informations visuelles en stimulations tactiles.

Deuxième question, dans quelle mesure la plasticité comportementale et cérébrale peut-elle pallier cette privation visuelle précoce ? Pour un apprentissage de la configuration spatiale d'un espace de déplacement, aucune incidence spécifique de la privation visuelle précoce n'a été relevée. En revanche, pour l'apprentissage de la configuration spatiale d'un petit espace, les aveugles de naissance recourent à des stratégies plus coûteuses (temps de réponse plus longs). La première question montre que les modalités sensorielles possèdent des propriétés sensorimotrices communes à un isomorphisme près. Cependant, il ne s'agit que d'un sous-ensemble des propriétés sensorimotrices, le reste représentant les relations sensorimotrices spécifiques aux sous-espaces sensorimoteurs. Ainsi, pour les petits espaces, la vision doit participer davantage à l'établissement des propriétés sensorimotrices, autrement dit la vision doit compléter ou enrichir les propriétés sensorimotrices générales, et en faciliter l'imagerie.

Ces deux réponses permettent de répondre à la troisième question concluant cette démarche : existe-t-il des limites à la « substitution » sensorielle d'un sens par les autres sens ? Il existe effectivement une limite à la substitution du sens visuel par les autres sens concernant l'expertise de l'individu dans le traitement de la métrique fine.

En conclusion, le constructivisme piagétien montre qu'un point de vue à la troisième personne permet une certaine compréhension des relations qui se façonnent mutuellement entre un sujet et un objet. Toutefois, une ambiguïté demeure : « les structures et les opérations logiques sont-elles dans le répertoire d'activité de l'enfant ou constituent-elles plutôt des caractéristiques formelles de la théorie piagétienne de l'esprit de l'enfant ? » (Bruner, 2000). Autrement dit, soit la cognition, à terme, s'identifie à un groupe d'opérations logiques, soit l'activité cognitive se trouve décrite par le théoricien constructiviste comme un groupe d'opérations logiques. Plus précisément, le premier cas nécessite l'acceptation de HOI puisque le système est guidé par les contingences sensorimotrices, c'est-à-dire les principes mathématiques définissant la théorie sensorimotrice dans l'absolu. Mais cette position trahit le raisonnement originel d'un constructivisme à la troisième personne. En effet, l'identification des structures logicomathématiques avec le sujet observé en soi revient à oublier que la relation autrui/objet forme et reste un objet pour le sujet observant et théorisant : « les structures logicomathématiques n'existent pas en tant que telles *dans* l'esprit de l'enfant mais elles sont *dans* le point de vue – daté historiquement et situé culturellement – du théoricien du développement. » (Troadek et Martinot, 2003). En conséquence, les travaux mathématiques offrent un cadre descriptif puissant qui permettrait de dégager un couplage entre un système et son environnement, entre les entrées sensorielles et les sorties motrices, mais ces mêmes travaux restent muets face à la notion d'autonomie, de la même manière que le fonctionnalisme écologique.

Néanmoins, dans le second cas, en considérant la description mathématique comme un point de vue, l'idée de pouvoir construire artificiellement un sujet se révèle délicate puisque son algorithme ne doit pas contenir en lui-même la description que l'observateur a du phénomène. Le répertoire d'actions doit se constituer d'une manière rudimentaire. Un exemple rentrant dans le cadre du constructivisme présenté ici est le modèle de réseau de neurones de Ziemke (2001) qui se fonde sur la récurrence et sur la prédiction, cette dernière correspondant au critère d'apprentissage. Le système acquiert un ensemble d'actions grâce à l'interaction, et sans être dirigé par des principes mathématiques a priori. Mais, en définitive, une fois cet ensemble extrait, correspondant certainement à un sous-ensemble des mouvements décrits par les générateurs qui aurait pu être calculé, aucun principe n'est proposé pour développer sur cette base la notion d'autonomie. En observant

les interactions sujet/objet, l'autonomie d'autrui peut être décrite à travers cette relation, sans pour autant capturer l'essence de celle-ci. Varela, en s'appuyant sur la notion d'autopoïèse, propose de dépasser le couplage entrée/sortie qui seul caractérise l'étude de la relation sujet/objet dans le constructivisme piagétien.

Une tentative de distinction entre les deux types de constructivisme se présente comme suit : pour le constructivisme piagétien, l'observateur étudie la relation autrui/objet, en supposant que c'est un point de vue indirect sur son propre vécu qui se refuse à l'introspection de l'observateur, alors que pour le constructivisme varelien, l'attention de l'observateur ne doit pas se porter exclusivement sur la relation émergente de l'interaction sujet/objet mais également sur la dynamique interne du sujet qui amène au vécu et produit cette relation. Autrement dit, l'observateur doit se concentrer sur l'élan du sujet vers l'objet qui les entraîne dans une danse existentielle, d'un côté, comment l'action construit le rapport sujet/monde, et de l'autre, comment l'action se génère dans le rapport sujet/monde.

C - Le constructivisme varelien

La différence entre les deux mouvements du constructivisme prend tout son sens avec une vision systémique. Tous deux se fondent sur la distinction entre autrui et le monde, de manière plus générique entre l'unité et le fond pour reprendre les termes de Varela (1989). Si aucune unité ne se dégage du fond, alors aucune description n'est envisageable, aucun ordre ne peut être extrait. Par ailleurs, l'idée d'une unité en dehors de toute chose s'auto-réfute par l'impossibilité de l'existence d'un observateur. Dans l'imaginaire, elle correspondrait à une description absolue et totale où le désordre n'a pas sa place. Entre ces deux pôles se trouvent toutes les descriptions systémiques possibles et imaginables. Plus particulièrement, dans l'étude de l'interaction initiée par une réflexion phénoménologique, les situations intéressantes, instructives, apparaissent lorsqu'il y a une juste indépendance relative entre les événements propres à l'unité et les événements qui appartiennent au milieu. Cette ambivalence fait écho à la richesse maximale que produit l'auto-organisation qui se trouve entre deux pôles, entre ordre et désordre.

Les influences mutuelles entre le milieu et l'unité, situées sur leurs interdépendances, forment une surface de couplage mais celle-ci ne représente pas toute l'unité, « elle ne constitue qu'une ou que quelques-unes de ses dimensions » (Varela, 1989). Ce couplage ponctuel traduit alors la transformation dynamique des états d'un système en fonction de ses entrées, qui se comprend comme le couplage entre les entrées et les sorties. Le formalisme du constructivisme piagétien appartient à ce type de couplage, de même que le béhaviorisme ou le subsymbolisme. Cependant, la signification de ce couplage ainsi que la manière de dégager les fonctions de transition les différencient. Ce type de couplage peut prendre en compte des variables ou des contraintes internes du système pour décrire le couplage ponctuel mais il ne parviendra pas, malgré tout, à prendre en compte les transformations internes caractérisant les systèmes autonomes qui sont justement le point de blocage du constructivisme piagétien.

i - Le couplage par clôture

Afin d'intégrer ces transformations internes, la notion de perturbation vient remplacer la notion d'entrée (*input*) : « Un *input* spécifie la seule façon dont une transformation d'état donnée peut avoir lieu. Une perturbation ne spécifie pas l'agent, elle ne prend en compte que son effet sur la structure de l'unité. En d'autres mots, un *input* fait partie intégrante de la définition d'une unité. Une perturbation peut être couplée à une

unité, elle ne fait pas partie de sa définition. Les diverses façons dont une perturbation donnée peut avoir lieu sont en nombre indéfini. Un *input* donné ne peut avoir lieu que de façon spécifique (c'est la raison pour laquelle on l'a défini) » (Varela, 1989). Cette notion de perturbation s'intègre au concept de système autopoïétique abordé dans la section précédente. En effet, le couplage structurel d'un système autopoïétique ne se comprend qu'avec l'intervention extérieure inattendue obligeant la modification de l'organisation interne pour compenser et ainsi maintenir sa dynamique autopoïétique. Celle-ci constitue l'élan du sujet qui pousse aux échanges et qui entraîne, en même temps, en fonction des perturbations, les transformations de l'organisation interne toujours en devenir.

Jusqu'à présent la notion d'autopoïèse a principalement été illustrée par une cellule élémentaire, mais les organismes multicellulaires conservent cette même propriété. Plus précisément, deux systèmes autopoïétiques ayant une histoire commune arrivent, à terme, à deux situations possibles : soit l'une inclut l'autre, c'est une symbiose endogène ; soit l'une et l'autre se juxtaposent pour créer un espace commun ou une forme, c'est une symbiose exogène. Dans les deux cas, les systèmes autopoïétiques de ces organismes contribuent à un couplage structurel conjoint entraînant une spécification mutuelle qui les lie fermement. Concernant les symbioses exogènes, l'espace commun ou la forme de l'organisme devient un espace privilégié qu'il faut maintenir, et qui, par la suite, constitue un système autopoïétique. L'organisme multicellulaire peut aussi se comprendre comme une clôture opérationnelle ; en effet selon Varela (1989) le système autonome « est opérationnellement clos si son organisation est caractérisée par des processus :

- a) dépendant récursivement les uns des autres pour la génération et la réalisation des processus eux-mêmes, et
- b) constituant le système comme une unité reconnaissable dans l'espace (le domaine) où les processus existent. »

Le couplage par clôture correspond au maintien de la cohérence des processus internes de l'organisme face aux perturbations. L'organisme multicellulaire étant un système autopoïétique, l'histoire de son couplage structurel fondé sur ce couplage par clôture correspond à son identité, de la naissance à la mort. Dans cette logique, à un moment donné, les points d'équilibre de l'organisation résultant du couplage structurel de l'organisme correspondent au couplage ponctuel évoqué précédemment. La cognition émerge du couplage structurel et se définit comme une compréhension incarnée liée à l'action, tout en participant à la clôture opérationnelle. La cognition n'est pas figée, elle se couple constamment avec l'environnement à travers tous les types d'interactions : physiques, culturelles, etc. Par exemple, la personnalité de Pierre correspond à son organisation cognitive, considérée comme un système autopoïétique. Les comportements de Pierre au quotidien reflètent les points d'équilibre de l'interaction entre son organisation cognitive et son environnement. Si Pierre subit un grave traumatisme, son organisation cognitive peut-être modifiée afin de maintenir une cohérence psychique interne. Bien que de nombreux comportements (ou points d'équilibre) demeurent identiques, d'autres disparaissent et d'autres se créent ; la personnalité se modifie, l'organisation cognitive est alors différente. Dans tous les cas, la personnalité n'est pas l'identité de la personne, l'identité représente l'histoire de celle-ci. Le vécu (ou l'interaction) résultant de l'organisation présente façonne l'organisation future.

La notion de processus dynamique résistant aux perturbations externes se retrouve à tous les niveaux. Par exemple, le système immunitaire où le soi n'est pas une simple étiquette définie à la naissance mais un processus complexe d'auto-étiquetage (Varela,

1989). La perception devient également un processus de compensation où les objets représentent en fait les points fixes des perturbations, du bruit que subit la clôture opérationnelle, autrement dit, la dynamique du système nerveux. Les représentations n'existent pas en tant que telles, les réminiscences ne sont que des points de fonctionnement que la dynamique globale du système retrouve bien que les situations sensorielles diffèrent de celles qui ont forcé à la transformation de la structure jusqu'à obtenir ce nouveau point de fonctionnement. Dans ce sens, la représentation comme réactivation d'un point de fonctionnement ne dit rien sur le monde en tant que tel mais reflète uniquement la dynamique propre du système nerveux.

Ainsi, ce paradigme englobe toutes les approches robotiques considérant que les comportements du robot reflètent les attracteurs d'un système chaotique, comme les travaux de Schönner et Dose (1992). Le robot évolue en fonction des perturbations extérieures et les « représentations de l'environnement » correspondent à des points de fonctionnement du système. Cependant, la structure du réseau est figée ou du moins le cadre de son développement est strictement défini. En d'autres termes, le rapport au monde modifie, oriente les points de fonctionnement, révèle la plasticité de l'organisation mais ne présente pas de transformations internes qui généreraient des comportements plus évolués. Les dynamiques du système sont anticipées ainsi que leur domaine de fonctionnement. Il y a toujours un cadre. Le couplage structurel ne peut pas se réaliser. De ce fait, le couplage par clôture porte uniquement sur l'ajustement d'une organisation prédéfinie, ce qui revient à modéliser la génération d'un couplage ponctuel, un couplage sensorimoteur qui peut également s'interpréter dans le cadre du constructivisme piagétien. Sans transformation interne, le système ne spécifie pas sa propre organisation et donc ses propres modes de fonctionnement, ce qui interdit toute autonomie cognitive. L'étude de la dynamique (Gerstner, 2002 ; Gaussier, 2002) des composants de la clôture opérationnelle, le substrat neuronal, devient alors primordiale pour comprendre comment peuvent se réaliser ces transformations internes, le couplage structurel. Autrement dit, l'inspiration neurobiologique se révèle indispensable pour comprendre les transformations internes et leurs conséquences aux niveaux supérieurs. Ainsi, il existe un lien fort avec le neuromimétisme, sauf que dans le cadre du constructivisme varelien, l'étude des mécanismes neurobiologiques ne se réduit pas à une fin en soi.

ii - Le constructivisme varelien et les observations éthologiques de Lorenz

Une autre manière d'aborder le problème revient à considérer que la notion d'apprentissage sous-tendue par l'auto-organisation d'un processus autopoïétique nécessite des concepts supplémentaires que l'étude du substrat peïnera à révéler. Anachroniquement, cette idée se retrouve dans la réflexion de Lorenz sur la formation du concept d'instinct (1937). Lorenz montre que les animaux possèdent leur propre dynamique et que leurs répertoires d'actions représentent les points de fonctionnement de cette dynamique, tout en mettant en garde de ne pas confondre les comportements issus d'un processus de maturation ou d'ajustement et ceux qui dépendent d'une expérience individuelle particulière. Cette analyse se décompose, ici, en trois points.

Le premier point signale que le développement des actes instinctifs ne dépend pas de l'expérience individuelle particulière et qu'il est analogue aux développements des organes, soulignant ainsi l'importance de la phylogenèse. En d'autres termes, avec le vocabulaire du constructivisme varelien, l'organe forme un système autopoïétique dont le couplage structurel dépend de la clôture opérationnelle des systèmes qui le composent et des échanges avec son environnement, les organes limitrophes. Dans ce cas, le développement

de l'organe s'effectue indépendamment d'une expérience particulière en rapport avec son rôle dans l'activité de l'individu, sa clôture opérationnelle adulte. Par exemple, le vol d'un jeune oiseau grandissant en captivité ne s'est jamais développé autrement que celui d'un oiseau vivant en liberté. De même, aucune coordination de mouvements ne s'est développée en s'adaptant aux conditions spatiales imposées. Par ailleurs, la réussite progressive de certains comportements complexes ne signifie pas l'existence d'un apprentissage, comme le montre Lorenz (1937) avec l'exemple de la confection de nids d'oiseaux. La première année, sur trois couples de bouvreuils pivoine, deux couples échouèrent dans la construction d'un nid et le dernier couple séparé des deux autres ne fit aucune tentative. L'année suivante, les trois couples réalisèrent chacun un nid de qualité égale.

Schématiquement, deux situations peuvent induire à tort l'idée d'un apprentissage. La première situation provient du fait que le développement des organes étant principalement lié à leur couplage structurel et non au couplage structurel de l'individu, la maturité des organes peut arriver de manière asynchrone alors que leurs activités pour l'individu sont liées. Cela conduit à deux configurations : soit le système nerveux ou plus précisément les structures cérébrales liées à la coordination sensorimotrice se trouve à terme alors que les organes moteurs ne le sont pas, soit à l'inverse, la coordination sensorimotrice arrive à maturité après le développement des organes moteurs. La seconde situation concerne les organes dont le développement ontogénétique est fortement dépendant des stimulations sensorielles. Il existe de nombreux exemples : la mémorisation des premiers êtres vus par les canetons, la formation du système visuel du chat, etc. Toutefois, tous présentent les mêmes caractéristiques : spécificité des stimuli attendus, durée déterminée de la phase de sensibilisation et inaltérabilité de ce qui a été imprégné. Chacune de ces trois caractéristiques s'interprète dans le cadre d'un couplage structurel (ou de l'ontogenèse) de l'organe. La première caractéristique signifie que, à un stade de développement, le couplage structurel se trouve particulièrement sensible à un certain type de perturbations. La deuxième caractéristique montre bien que le couplage structurel est un processus continu avec sa dynamique propre, c'est-à-dire que, pendant la phase transitoire où la sensibilité à certaines perturbations augmente, l'organe s'accommode s'il y a lieu, mais continuera de toute façon son couplage et aboutira à une structure plus stable. La troisième caractéristique s'explique d'une part par la stabilité de la nouvelle structure, et d'autre part par l'irréversibilité du couplage structurel. Mais il s'agit toujours du couplage structurel de l'organe, et non de celui de l'individu, même s'il existe nécessairement une interdépendance. En revanche, le résultat du couplage structurel de l'organe va ensuite fortement conditionner la clôture opérationnelle de l'individu.

Le second point insiste sur le fait que la plasticité régulatrice ne doit pas être assimilée à l'apprentissage et à l'expérience. Les actes instinctifs tels que la coordination des mouvements de la marche sont susceptibles de régulation mais cette dernière ne provient pas d'une modification adaptative de l'acte instinctif par l'expérience. En effet, les expériences de Bethe (1899) relatées par Lorenz (1937) vont à l'encontre d'une telle supposition : « un chien dont les nerfs sciatiques avaient été sectionnés et recousus de façon croisée, réussit, dès le rétablissement de leur fonctionnement, à coordonner parfaitement et tout à fait normalement les mouvements de la marche. En revanche, en ce qui concerne la sensibilité, il n'y a eut aucune régulation, l'animal réagissant constamment à des excitations douloureuses émises sur une patte de derrière, avec l'autre patte. L'apparition d'une régulation dans le domaine moteur, liée à son absence dans le domaine sensible, est la preuve la plus claire du fait que l'expérience ne joue aucun rôle dans la

formation d'une régularisation motrice » (Lorenz, 1937). En d'autres termes, la régulation correspond à la capacité de l'organisation à s'équilibrer vers un autre point de fonctionnement, sans remise en cause de l'organisation, donc sans couplage structurel.

L'existence de cette régulation ne contredit pas la théorie piagétienne de la construction d'un espace sensorimoteur. En effet, celle-ci s'appuie également sur l'existence de boucles sensorimotrices initiales nécessaires à la découverte du monde. Ici, Lorenz souligne que ces boucles sensorimotrices peuvent être associées à des comportements relativement complexes comme la marche, et que les confondre avec des comportements construits ou appris à partir de l'établissement d'un espace sensorimoteur serait une erreur. Le chapitre suivant reviendra sur les raisons susceptibles d'expliquer alors l'intérêt de construire un espace sensorimoteur.

Le troisième point affirme que le répertoire des actes primitifs est fixe et que la simple modulation de l'excitabilité des actes instinctifs permet d'expliquer généralement l'adéquation entre l'action et son utilité biologique pour les animaux « supérieurs ». Plus précisément, cette modulation peut se décomposer en deux mécanismes : d'une part l'auto-modulation des actes instinctifs sans finalité, et d'autre part la notion d'appétence qui oriente la modulation des actes instinctifs vers une finalité, autrement dit une alternance entre instinct et dressage. L'antériorité du premier mécanisme sur le second souligne bien la différence entre le constructivisme et le fonctionnalisme écologique qui oriente tous les actes instinctifs vers une finalité. L'existence du premier mécanisme se fonde sur trois constats.

Le premier constat est qu'un animal se satisfait aussi bien d'une suite d'actes inachevés sans aucun rôle biologique que du déroulement complet d'actes atteignant un but. L'acte ne vise pas un but, il se déclenche avec une intensité variable selon la situation excitatrice. Le deuxième constat porte sur l'indépendance entre l'acte instinctif et les excitations susceptibles de le déclencher. Deux cas montrent ce phénomène. Le premier cas correspond à la désensibilisation (ou parfois à la sensibilisation) de l'acte instinctif par accoutumance : « L'animal, pendant qu'il s'accoutume à l'excitation, se comporte exactement comme si c'était l'intensité de l'excitation qui allait en décroissant. C'est ainsi que la même excitation peut déclencher des réactions différentes liées de fait à des excitations de force différente. Il peut arriver qu'après le recul progressif de l'intensité d'une réaction surgisse brutalement une autre réaction » (Lorenz, 1937). Le deuxième cas montre que si un acte instinctif n'a pas été déclenché pendant un certain temps, le seuil à partir duquel agissent les excitations nécessaires à son déclenchement s'abaisse considérablement, au point que l'acte instinctif se produit sans excitation apparente. Enfin, le troisième constat révèle qu'un ensemble d'actes instinctifs indépendants peut former une unité fonctionnelle homogène vis-à-vis d'un objet réunissant tous les déclencheurs. Par conséquent, il n'y a pas besoin de notion d'objet pour expliquer certains comportements comme ceux destinés à s'occuper des petits. Par exemple, une cane de Barbarie prend la défense d'un petit canard col-vert comme elle l'aurait fait pour ses propres petits, pour ensuite le traiter avec hostilité. Cette réaction s'explique par le fait que les cris d'alarme se ressemblent fortement alors que les motifs sur la tête, reliés aux comportements de soins, diffèrent beaucoup.

Ce premier mécanisme révèle alors que le schème des actes instinctifs dépend de l'ontogenèse des organes qui les supportent et non de leur finalité par rapport à l'individu. En revanche, la clôture opérationnelle de l'individu se couplant avec l'environnement organise les déclencheurs des actes instinctifs qui se juxtaposent et qui se définissent les uns

par rapport aux autres. Inévitablement, la dynamique du couplage par clôture se fonde alors sur l'enchaînement des réactions déclenchées par le simple fait qu'ils puissent modifier les conditions environnementales.

Le second mécanisme de ce troisième point concernant l'efficacité des comportements instinctifs, qu'il ne faut pas assimiler à des comportements intelligents, coïncide avec certaines idées du fonctionnalisme écologique. Ce second mécanisme permet d'établir un comportement finalisé, c'est-à-dire un comportement susceptible d'accepter une modification adaptative tout en conservant une finalité identique. Autrement dit, un comportement finalisé constitue un enchaînement d'actes instinctifs guidés par une finalité. De manière plus précise, ce mécanisme utilise deux notions, celle d'appétence et celle de satisfaction, ce qui se traduit en terme de fonctionnalisme écologique par « désir » et « plaisir ». D'une part, l'appétence prédispose l'individu à l'excitation de certains actes instinctifs, et d'autre part, la satisfaction permet de reconnaître la situation excitatrice. Ce processus de dressage n'implique pas que l'individu soit conscient de la finalité ou du besoin, c'est simplement la satisfaction qui structure la chaîne comportementale. Ainsi, le déroulement de l'acte instinctif constitue le but et la finalité de tout comportement animal (Lorenz, 1937). L'alternance instinct-dressage aboutit à une série d'actes homogènes d'un point de vue fonctionnel où se succèdent sans transition des fractions de comportements instinctifs et de comportements finalisés.

Par exemple, concernant la confection du nid, les corbeaux possèdent, avant la période de construction, l'acte instinctif de saisir avec leur bec n'importe quel objet et de l'emporter en volant sur de longues distances. Au départ, ils ne manifestent aucune préférence particulière pour le matériau, et parfois certains se sensibilisent pour des objets morphologiquement très éloignés des objets de type brindille, comme des morceaux de tuiles. La préférence pour les brindilles survient après l'apparition, sous-tendue par une appétence, d'un nouvel acte instinctif qui consiste à coordonner ses mouvements de façon à produire des secousses latérales. Ce mouvement aboutit à deux issues possibles lors de la construction du nid : soit l'objet ne s'insère pas et tombe, soit l'objet s'accroche et oppose une résistance. Suite à un certain nombre d'expériences, le type de matériau favorisant la deuxième situation se trouve très rapidement privilégié.

Ce mécanisme d'alternance instinct-dressage formant des comportements finalisés adaptés représente un premier pas vers les comportements intelligents mais cela ne signifie pas qu'il puisse résoudre tous les problèmes et surtout qu'il puisse permettre une certaine abstraction de la finalisation de l'acte ou de la signification. En effet, Lorenz (1937) illustre ce point avec le comportement d'une femelle canari qui avait pris pour habitude au cours de la construction de son nid de maintenir les tiges pour les enchevêtrer, et comme ces tiges étaient du fourrage, elle utilisait la même technique pour s'en nourrir. Cette pratique facilitait grandement la consommation, mais après la période biologique propice à la construction du nid, la femelle canari n'utilisait plus cette technique pour maintenir sa nourriture. Il y a donc une différence considérable entre l'application d'une aptitude et l'identification de la finalité de celle-ci. Par ailleurs, les actes instinctifs se soumettent à l'usage de différentes appétences sans pour autant être modifiés, de même que le bec ne se modifie pas selon les différents usages comme la quête de nourriture, le combat, etc.

En somme, le constructivisme varelien offre un cadre explicatif cohérent avec toutes ces remarques éthologiques qui pourraient se résumer par cette suite de questions-réponses :

1. Pourquoi les actes instinctifs se développent-ils indépendamment de l'expérience individuelle ? Parce que l'ontogenèse de l'organe se distingue du couplage de l'individu avec son environnement.
2. Pourquoi les animaux peuvent-ils être imprégnés à vie par certains stimuli seulement dans leurs premiers stades de développement ? Parce que des événements peuvent marquer le couplage structurel, lequel relève d'un processus dynamique interne, continu, historique donc irréversible, et aboutissant à une structure stable. Mais surtout le couplage structurel ne doit pas être confondu avec le couplage par clôture qui, lui, offre une plasticité organisationnelle.
3. Pourquoi les déclencheurs semblent-ils indépendants des actes instinctifs ? Parce que l'acte instinctif reflète un point de fonctionnement parmi d'autres de la dynamique interne de l'organisme face à une perturbation extérieure. Mais ce déclenchement entraîne également des modifications internes portant sur la réorganisation des points de fonctionnement, les uns par rapport aux autres, ce qui correspond à un couplage par clôture.
4. Pourquoi l'excitabilité des actes instinctifs varie-t-elle a priori sans modifications de l'environnement ? Parce que l'organisme possède sa propre dynamique (son propre cycle) conduisant à des modifications dans l'espace des points de fonctionnement.

Deux conclusions se dessinent. La première souligne que l'importance de l'ontogenèse des organes dans l'acquisition des actes instinctifs suggère une phylogenèse des comportements instinctifs qui pourrait même compléter la taxonomie généralement restreinte à l'anatomie comparée (Lorenz, 1937). La seconde conclusion, la plus importante pour cette analyse du constructivisme varelien, soutient que cette étude sur le comportement animal circonscrit le domaine des comportements appris. En effet, Lorenz dégage un ensemble de comportements complexes qui pourrait être, à tort, qualifié d'intelligent. Il ne s'agit en fait que de régulation ou d'ajustement même s'ils peuvent être complexifiés par le mécanisme d'auto-dressage, l'essentiel des comportements se trouve déjà dans le système, y compris la notion d'empreinte qui n'est qu'une phase de sensibilité particulière dans un processus en cours. L'explication de tous ces mécanismes exploite au mieux tous les concepts de ce constructivisme.

Cependant, ces mécanismes ne proposent aucune piste pour accéder aux comportements intelligents, c'est-à-dire pour expliquer un comportement adopté par le robot imaginé dans le conduit d'aération évoqué dans l'introduction. Sans aborder l'humain, le début de comportement intelligent incluant un apprentissage se trouve pourtant déjà dans certaines espèces citées précédemment. Par exemple, une expérience consiste à placer un bout de papier près d'un morceau de viande appartenant à un corbeau, les corvidés ayant généralement tendance à dissimuler leur nourriture non consommée. Se dessinent alors deux situations : le corbeau a joué auparavant avec ce bout de papier en essayant toutes sortes de manipulations, et il recouvrira la viande directement avec le papier ; dans le cas contraire, ce bout de papier lui est inconnu et le corbeau cessera d'essayer de cacher son butin pour jouer avec ce nouvel objet (Lorenz, 1950). Les seuls mécanismes définis précédemment ne parviennent pas à expliquer ce type de comportement. L'apprentissage possède des liens étroits avec les notions d'objet et d'intentionnalité qui permettent une réflexion sur la situation ou une reconnaissance d'une solution à un problème. Ici, l'apprentissage s'inscrit dans l'action et sa finalité est fondée sur les propriétés dynamiques du comportement de l'individu. Afin d'éviter la confusion avec d'autres conceptions, Varela (1991) propose le terme d'énaction. L'impossibilité d'aller au-

delà de la circonscription du comportement intelligent ne remet pas en cause la puissance du paradigme de Varela puisque celui-ci autorise des descriptions à des niveaux très différents de la cellule à l'humain, ainsi qu'une esquisse de compréhension à un certain niveau du fonctionnement cognitif. Néanmoins, l'incapacité à franchir cette frontière montre que le constructivisme varelien se révèle au mieux incomplet pour expliquer la construction d'un système cognitif.

Avant de conclure sur l'interactionnisme, il convient de préciser que le constructivisme varelien, comme le constructivisme piagétien, reste menacé par l'oubli de la circularité du discours à la troisième personne. Cet écueil conduirait, par exemple, à esquisser rapidement une correspondance entre le vocabulaire des deux dernières formes de phénoménologie et celui du constructivisme varelien. Le Dasein – être-là ou être-*ceci* – correspondrait à l'actuelle frontière toujours mouvante entre le système autopoïétique et l'environnement. « L'essence du Dasein est dans son existence » de Heidegger (1927) ou « l'existence précède l'essence » de Sartre (1946) feraient écho à l'organisation (l'essence) qui se spécifie et se maintient par rapport à l'interaction permanente avec le monde, ainsi qu'au point de fonctionnement ou le couplage ponctuel (l'existence) résultant de cette interaction. Le couplage par clôture ressemblerait alors à l'*in-der-Welt-sein* – l'être-dans-le-monde – (Heidegger, 1927). Mais ces correspondances travestiraient l'idée de ces auteurs qui aspirent à définir des propriétés ontologiques du phénomène alors que le constructivisme accepte de ne pas pouvoir sortir de la subjectivité de la relation sujet/objet même en se penchant sur la relation *autrui/objet* (en tant qu'objet).

		Types de famille métaphysique				Hypothèses sur la perception			Hypothèses sur la conception			Types de méthode d'investigation			Types de description		
		FC	FCP	FP	FA	HP1	HP2	HP3	HC1	HC2	HC3	MHD	MC	MA	DP	DM	DC
Interactionnisme	Phénoménologie	Jaune	Rouge	Rouge	Vert	Rouge	Rouge	Rouge	Jaune	Jaune	Rouge	Grise	Grise	Grise	Grise	Grise	Grise
	Constructivisme piagétien	Jaune	Rouge	Rouge	Vert	Rouge	Rouge	Rouge	Rouge	Rouge	Rouge	Vert	Vert	Jaune	Rouge	Vert	Rouge
	Constructivisme varelien	Jaune	Rouge	Jaune	Vert	Rouge	Rouge	Rouge	Rouge	Rouge	Rouge	Vert	Vert	Jaune	Rouge	Vert	Vert

Tableau I-10 : Récapitulation de diverses positions relatives à la grille d'analyse des trois principales approches provenant du courant interactionniste. Les cases vertes correspondent à un avis favorable, les cases jaunes à un avis mitigé et les cases rouges à un avis défavorable. Les cases grises signifient qu'une prise de position se trouve hors de propos.

En définitive, l'interactionnisme présente une pluralité de positions et de méthodes (Tableau I-10). Cependant, l'interactionnisme diffère des autres mouvements par sa focalisation sur la stricte interaction sujet/objet qui se répercute sur le choix des hypothèses ontologiques. En effet, seul l'interactionnisme se trouve promu par la famille FA. Toutefois, les HOP et HOI peuvent apparaître dans certains courants phénoménologiques avec toutefois une interprétation différente de celle évoquée dans les autres mouvements. En revanche, le constructivisme refuse en théorie ces hypothèses, mais certains confondent le refus des hypothèses ontologiques avec l'impossibilité d'accéder aux propriétés ontologiques visées. Cette distinction sera reprise lors de la conclusion de ce chapitre. Concernant les hypothèses portant sur la perception et la conception, la phénoménologie et le constructivisme auront en général tendance à les rejeter en s'appuyant principalement sur les phénomènes d'hallucination et d'illusion traduisant les limites de l'introspection. Plus spécifiquement, pour le constructivisme, les méthodes d'investigation possèdent une valeur équivalente entre elles, contrairement aux trois types de description. Le constructivisme piagétien privilégie obligatoirement la DM pour décrire la relation entre *autrui/objet* en

terme d'entrée/sortie alors que le constructivisme varelilien opte à la fois pour la DM afin de décrire le domaine de fonctionnement d'une organisation, et pour la DC, afin de décrire le couplage structurel puisque l'organisation est en devenir. Malgré tout, les réalisations robotiques se revendiquant de cette mouvance spécifient entièrement l'environnement et l'unité, soit par les lois régissant la relation sujet/objet, soit par la dynamique du système interagissant c'est-à-dire son domaine de fonctionnement. Les comportements des robots se retrouvent complètement soumis aux contingences environnementales et aux règles prédéfinies interdisant toute autonomie cognitive. Par ailleurs, cette spécification globale et indirecte introduit, involontairement ou non, l'idée d'une maîtrise totale de ce qui est conçu, ce qui revient à réintroduire les hypothèses ontologiques. La conclusion de ce chapitre abordera plus en détail les conséquences d'un tel glissement.

4. Conclusion

Au terme de ce chapitre, une première conclusion un peu rapide serait de considérer que toutes les approches robotiques menées dans l'idée d'étudier la cognition n'aboutissent qu'à des systèmes hétéronomes intégrant des mécanismes de régulation plus ou moins complexes mais restant totalement restreints à leurs domaines de fonctionnement, autrement dit dépourvus d'autonomie cognitive. Cette conclusion invaliderait l'idée, défendue dans l'introduction de ce mémoire, que l'avenir de la robotique passe par les sciences cognitives. Ainsi, la robotique autonome deviendrait finalement une sous-discipline des sciences et techniques de l'automatisme qui pourrait se servir des sciences cognitives comme d'un moyen ludique pour illustrer les théories. Afin d'éviter cet écueil, dans un premier temps, (A) un récapitulatif des diverses conceptions de la cognition ainsi que les approches associées sera entrepris. Dans un second temps, (B) une analyse globale, destinée à déterminer les points communs entre ces courants, sera menée sur les hypothèses et méthodes employées. Enfin, une seconde conclusion, alternative à la première, sera alors proposée.

A - Récapitulatif des courants

Au total, sept courants principaux se dégagent de ce chapitre (Figure I-45). En reprenant le plan de ce dernier, le premier des deux mouvements des sciences cognitives identifiés se concentre sur la description de l'individu en soi, c'est-à-dire que la structure cognitive (la capacité d'apprendre, la personnalité) ne dépend pas de l'environnement. À partir de ce point, deux positions se présentent : soit les concepts mentaux (la psychologie au sens large) reflètent l'activité cognitive, hypothèse défendue par le mentalisme, soit seule l'étude du substrat permet de l'expliquer, position de l'éliminativisme. Le premier cas comporte principalement deux courants.

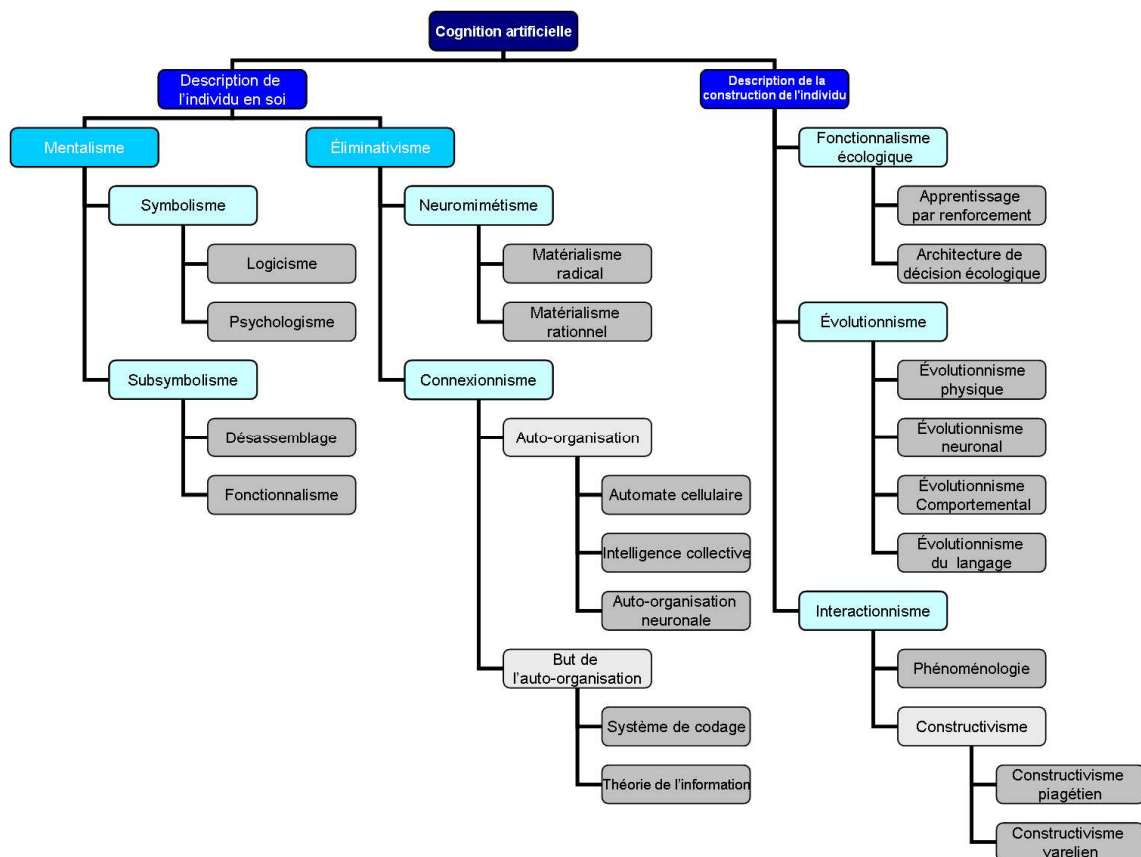


Figure I-45 : Arborescence des approches en cognition artificielle.

i - Le symbolisme

Le premier courant, *le symbolisme*, assimile la pensée à la manipulation d'idées pouvant se formaliser par des symboles qui correspondent à un état du monde. Dans ce cadre, deux approches s'affrontent : d'une part l'étude du langage naturel comme moyen d'accéder aux propriétés cognitives et d'autre part l'étude des systèmes logiques censés décrire les véritables propriétés de l'intelligence que le langage naturel approxime maladroitement. Ces deux approches reflètent respectivement le psychologisme et le logicisme évoqués précédemment dans l'interactionnisme. Concernant les sciences de l'artificiel, ces domaines d'investigation aboutissent d'une part au dialogue homme-machine ou multi-agents et à la reconnaissance de la parole, et d'autre part à des méthodes de résolution de problèmes logiques dans des systèmes formels et à la reconstruction tridimensionnelle de scènes visuelles. Les difficultés rencontrées par ces approches mettent en relief les limites des systèmes formels sur lesquelles le deuxième chapitre reviendra mais surtout l'impossibilité de la mise en correspondance entre le langage et le monde extralinguistique.

ii - Le subsymbolisme

Le subsymbolisme souhaite dépasser le problème de la mise en correspondance en introduisant la notion d'intentionnalité et donc de point de vue dans la perception. *La cognition devient la manipulation de symboles extraits du monde selon un traitement assimilé au concept signifiant mais ce traitement de bas niveau suit les mêmes principes que la manipulation symbolique reposant sur une mise en correspondance.* Les traitements de bas niveau correspondent alors au calcul élémentaire que pourrait effectuer un neurone, ce qui autorise l'identification entre

les états mentaux et l'activité cérébrale. Le bas niveau et le haut niveau se distinguent l'un de l'autre ; l'un est figé alors que l'autre est capable de s'adapter suivant un critère. La cognition résulte ainsi de la synergie de toutes les fonctions cognitives hiérarchisées. Deux approches apparaissent : l'une cherche à répertorier et à définir ces traitements ou fonctions, ce qui correspond à une décomposition fonctionnelle ; l'autre souhaite découvrir le langage élémentaire de programmation de celles-ci, et ainsi effectuer en terme informatique un désassemblage. Cependant, ces deux voies demeurent périlleuses car l'identification des fonctions ou du « langage machine » d'un individu cognitif semble impossible, puisque les fonctions toutes liées ne peuvent s'étudier indépendamment les unes des autres. Par ailleurs, en robotique, l'analyse fonctionnelle permet de définir les objectifs de chaque fonction ainsi que leur critère d'adaptabilité. Cependant, en agissant ainsi, le domaine de fonctionnement se trouve restreint et interdit toute autonomie.

iii - Le neuromimétisme

La deuxième position du premier mouvement décrivant l'individu en soi, l'éliminativisme, se scinde lui-même en deux courants : le neuromimétisme et le connexionnisme. *Le premier courant réduit la cognition à l'activité cérébrale et sa compréhension passe alors par sa description.* Suivant le type de matérialisme adopté, radical ou rationnel, l'étude du substrat a tendance à se tourner soit vers la description anatomique du substrat et de ses mécanismes, sans se soucier du comportement global de l'individu, la neurobiologie, soit vers la description en situation de l'activité cérébrale, la neuropsychologie. De ces diverses études, trois points ressortent. Le premier montre que la complexité et la diversité des neurones interdisent d'identifier le neurone à un calculateur élémentaire, comme le suggérait le subsymbolisme. Le deuxième point révèle la forte connectivité des neurones et le grand nombre de phénomènes entrant en jeu à différentes échelles dans l'apprentissage. Le troisième point insiste sur la différence entre rôle et fonction. Autrement dit, les modules se définissent pour et par rapport aux autres et non par rapport à leur fonction comme le suppose le subsymbolisme. Le type de matérialisme influence également le choix des approches artificielles. La première approche, plus rare, implémente analogiquement des modèles neurobiologiques alors que la seconde approche utilise des implémentations numériques. Néanmoins, les deux approches se bornent à réaliser des imitations plus ou moins fidèles des systèmes nerveux d'organismes existants. Ici se situe le principal problème de ce courant : la seule description du substrat ne parvient pas à accéder à un niveau de compréhension suffisant pour généraliser les concepts et réaliser de nouveau un individu cognitif autonome.

iv - Le connexionnisme

Le deuxième courant de l'éliminativisme, le connexionnisme, étudie la coordination spontanée d'un ensemble d'éléments dont l'activité de chacun dépend de celles des autres suivant leurs relations, en un mot de leur connexion, mais non de la finalité de la coordination. Dans ce cadre, *la cognition se comprend d'une part comme une résultante d'un processus d'auto-organisation et d'autre part comme un moyen de dégager les régularités et les structures du monde.* Le connexionnisme explore ces deux facettes. La première facette comporte trois approches. La première approche s'appuie sur l'étude des automates cellulaires. L'un des résultats montre que la richesse maximale des motifs se situe entre les règles cellulaires conduisant à l'ordre et celles du désordre, mais les automates cellulaires évoluent nécessairement dans un univers clos et homogène. La deuxième approche souhaite dépasser cette difficulté en étudiant les comportements sociaux provoqués par stigmergie. La coordination d'unités dans un monde ouvert permet de résoudre des tâches qu'aucune unité ne peut concevoir et

en ce sens, il y a un comportement émergent. Toutefois, sans introduire une certaine réflexion sur les actes et sur leurs finalités, une unité ou une colonie ne deviendra jamais cognitivement autonome. La troisième approche étudie l'émergence de propriétés cognitives dans les réseaux de neurones formels. Cette approche montre que l'apprentissage peut être totalement distribué et sans superviseur, l'émergence, contrairement au subsymbolisme qui propose un apprentissage également distribué sur plusieurs unités de calcul mais supervisé : la synergie. La deuxième facette du connexionnisme portant sur le moyen de dégager les régularités du monde se décompose en deux approches : l'une se focalise sur l'émergence d'un système de codage du monde issu de réseaux associatifs, l'autre formalise au sein de la théorie de l'information le problème de l'optimisation d'un tel codage. La principale difficulté réside dans la définition du critère d'optimisation. Ce dernier dépend de la définition de la quantité d'information ici reliée à la fréquence d'occurrence. Les systèmes connexionnistes rendent compte passivement du monde, le structurent mais ne parviennent pas à fournir une sémantique.

v - Le fonctionnalisme écologique

Les limitations des quatre courants qui entreprennent la description de l'individu en soi proviennent de la distinction ontologique entre le sujet et le monde. Or, l'attribution du sens peut se faire uniquement par rapport au monde. La compréhension du sujet cognitif passe nécessairement par la compréhension des interactions avec son milieu. Le deuxième mouvement des sciences cognitives se porte alors sur la construction de l'individu au fil de ces échanges. Ce mouvement se décompose en trois courants. *Le premier courant, le fonctionnalisme écologique définit la cognition comme un processus permettant de se constituer des croyances utilisées dans l'analyse des besoins pour la définition des buts et dans le choix de l'action en fonction de ces derniers.* La notion de besoin permet de définir l'origine de la sémantique tout en offrant une grande liberté sur la construction des croyances et par conséquent sur leur ancrage dans le monde. En effet, la vérité d'une croyance porte sur son assertion et sur les conséquences des actions associées. Malgré tout, cette définition de la vérité possède intrinsèquement une limitation due à l'incertitude sur l'action et sur le monde, que seule l'anticipation peut minimiser sans toutefois s'y soustraire. En sciences de l'artificiel, deux approches abordent ce paradigme : l'une s'attache à comprendre les mécanismes d'apprentissage guidés par une rétribution induite par la satisfaction d'un besoin, et l'autre souhaite dégager l'architecture la plus adaptée pour décider de l'action en fonction de l'utilité de celle-ci. Mais en spécifiant les besoins, ce paradigme redéfinit indirectement une description de l'individu en soi avec des capacités d'optimisation adaptées a priori à son milieu.

vi - L'évolutionnisme

Le second courant proposant une description de la construction de l'individu, l'évolutionnisme, affirme que *la cognition résulte de l'évolution de l'espèce qui vise l'optimum, soit de complexité en absolu, soit d'efficacité à survivre.* L'utilisation d'arguments téléologiques sur l'évolution des espèces permet d'aborder la cognition indirectement en recherchant les mécanismes de l'évolution et les contraintes environnementales l'infléchant. Ce courant s'appuie principalement sur l'idée de l'évolution du code génétique et de l'évaluation de son phénotype. Cette idée maîtresse se décline alors en quatre approches. La première souhaite modéliser un monde virtuel pour recréer le processus de l'évolution. La deuxième approche transpose ce principe sur le plan comportemental. La troisième approche exploite le principe d'évolution appliqué à une population de neurones. La dernière approche reproduit le mécanisme de l'évolution sur le langage, les individus devenant secondaires. Ces approches montrent dans l'ensemble que les notions d'évolution et de sélection offrent

une heuristique efficace par rapport à des critères choisis. Mais ces approches se fondent sur une certaine interprétation du néodarwinisme qui repose uniquement sur la notion de codage génétique prédominant sur la notion d'individu. À l'inverse, une seconde interprétation, avec la notion d'autopoïèse, met l'individu au cœur du processus de l'évolution. Ce changement de perspective interdit l'utilisation d'arguments téléologiques pour définir la cognition.

vii - L'interactionnisme

Le dernier courant qui s'inscrit dans le deuxième mouvement des sciences cognitives considère que l'interaction entre le sujet et l'objet prend racine au delà de la notion de besoin qui devient alors secondaire. *La cognition correspond au processus qui dégage ou fait émerger la coordination globale des activités du sujet dans le monde.* Deux stratégies apparaissent pour étudier ce processus. La première stratégie, la phénoménologie, souhaite étudier la relation sujet/objet uniquement à la première personne et par suite entrevoir la possibilité d'étudier les opérations cognitives du sujet en se concentrant sur la relation sujet/sujet. Mais en définitive, la circularité de l'introspection aboutit à l'étude de la relation sujet/sujet-objet et de ce fait le sujet se trouve dans l'impossibilité de se réfléchir totalement. Afin de dépasser cette limitation, la seconde stratégie observe la relation sujet/objet à la troisième personne, c'est-à-dire la relation autrui/objet. Deux approches existent pour décrire cette dernière. La première se concentre sur la détermination des lois régissant le couplage sujet/objet, ce qui oblige à percevoir le sujet comme un système possédant des entrées (senseurs) et des sorties (effecteurs). Le couplage entre les entrées et les sorties traduit des propriétés d'un espace défini à la fois par le monde et par l'activité du sujet. La cognition résulte alors des couplages successifs, ce qui signifie qu'elle devient historique. Toutefois, cette approche ne parvient pas à formaliser la cognition au-delà du couplage sensorimoteur ; en effet, l'établissement de lois de couplage ne correspond pas à la notion d'autonomie définie dans l'introduction de ce chapitre. La seconde approche s'attache davantage à décrire les propriétés systémiques qu'autrui doit présenter pour générer la relation autrui/objet. Autrui possède une dynamique propre qui repose sur les échanges avec le monde. Cette situation génère un couplage entre la dynamique interne du système et les perturbations du monde extérieur. Cette théorie se trouve en adéquation avec un certain nombre de données éthologiques sur l'équilibre des comportements. Cependant, décrire la cognition comme un système autopoïétique ne suffit pas à expliquer en détail l'apprentissage ou la mise en problématique.

*

Les sept courants dégagés ici (Figure I-45) ne représentent que des pôles entre lesquels un véritable continuum de positions existe. Les définitions proposées pour la cognition se révèlent toutes très différentes. Elles semblent parfois complémentaires néanmoins, même en ne tenant pas compte des contradictions internes d'un point de vue philosophique. Aucune réalisation robotique issue de la juxtaposition de toutes ces techniques ne parvient à l'autonomie. L'échec de ce type de démarche, mixte et opportuniste, montre que l'impasse conceptuelle s'avère commune à l'ensemble des courants. Une autre manière d'intégrer la situation consiste à analyser les hypothèses et les méthodes qui sous-tendent consciemment ou non leur conception de la cognition.

B - Analyse globale des hypothèses et des méthodes employées

Les courants de la cognition artificielle s'affrontent sur la définition de la cognition. Toutefois, l'analyse entreprise dans ce chapitre montre que dans certains cas des hypothèses ou des convictions fondatrices de ces courants se révèlent communes bien qu'interprétées différemment. Plus précisément, cette analyse a porté sur les hypothèses concernant à la fois la métaphysique, la perception, la conception ainsi que sur les trois méthodes d'investigation et les trois types de description. L'objectif de cette grille d'analyse est de découvrir le point commun à tous les courants qui les mène à l'impasse. Le Tableau I-11 récapitule les diverses positions des courants vis-à-vis des différentes hypothèses et méthodes.

			Types de famille métaphysique				Hypothèses sur la perception			Hypothèses sur la conception			Types de méthode d'investigation			Types de description		
			FC	FCP	FP	FA	HP1	HP2	HP3	HC1	HC2	HC3	MHD	MC	MA	DP	DM	DC
Description de l'individu en soi	Mentalisme	Symbolisme	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	
		Subsymbolique	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	
	Éliminativisme	Neuromimétisme	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	
		Connexionnisme	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	
Description de la construction de l'individu	Fonctionnalisme écologique	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■		
	Évolutionnisme	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■		
	Interactionnisme	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■		

Tableau I-11 : Récapitulation des positions des différents courants de la cognition artificielle vis-à-vis des hypothèses et méthodes constituant la grille de lecture de ce chapitre.

Le Tableau I-11 tableau montre qu'aucun type de description ne domine sur l'ensemble des courants. De la même manière, l'utilisation des méthodes d'investigation se trouve également répartie. Ces deux points suggèrent que l'impasse générale n'est pas d'ordre méthodologique. Parmi les hypothèses sur la perception et sur la conception, certaines font quasiment l'unanimité, d'autres, à l'inverse, sont unanimement désavouées. Bien que ces hypothèses soient indépendantes entre elles, cette répartition montre que les hypothèses sont en fait rattachées à des conceptions globales de la cognition. Il faut aussi préciser que la variabilité aurait été plus grande en affichant les hypothèses appliquées à l'épistémologie. En effet, le neuromimétisme et le connexionnisme dissocient à tort le problème de la cognition de celui de l'épistémologie qui est considérée en général comme source de neutralité. Ainsi, ils peuvent adopter une hypothèse perceptive dans le cadre de l'accumulation de connaissances scientifiques mais ne pas la transposer comme une propriété fondamentale d'un système cognitif. Cependant, afin de conserver la cohérence de l'analyse, seule la valeur des hypothèses dans le cadre de la construction d'un système cognitif est prise en compte. Quoi qu'il en soit, la disparité des valeurs sur les hypothèses concernant la perception ou la conception disqualifie celles-ci pour correspondre au point commun entre les courants.

Enfin, les quatre familles métaphysiques identifiées possèdent chacune au moins un courant les représentant. Par ailleurs, plusieurs familles peuvent cohabiter au sein d'un même courant et une même famille peut soutenir différentes conceptions de la cognition comme le montre plus exactement le Tableau I-12. Cette situation crée un réseau de dissemblances et de ressemblances entre les différentes approches qui ne facilite ni la fédération des approches ni leur l'examen. Mais surtout, aucune des hypothèses ontologiques définissant les familles métaphysiques n'est commune à l'ensemble des courants, ce qui signifie qu'il n'existerait pas de point commun.

		Types de famille métaphysique				Hypothèses sur la perception			Hypothèses sur la conception			Types de méthode d'investigation			Types de description		
		FC	FCP	FP	FA	HP1	HP2	HP3	HC1	HC2	HC3	MHD	MC	MA	DP	DM	DC
Symbolisme	Logicisme	Red	Green	Red	Red	Green	Green	Red	Green	Green	Yellow	Yellow	Green	Yellow	Green	Red	Red
	Psychologisme	Red	Green	Red	Red	Green	Green	Green	Green	Green	Yellow	Green	Yellow	Green	Red	Red	Red
Subsymbolique	Désassemblage	Yellow	Green	Red	Red	Green	Green	Green	Green	Green	Green	Yellow	Green	Yellow	Red	Red	
	fonctionnalisme	Red	Green	Red	Red	Green	Green	Green	Green	Green	Green	Yellow	Green	Yellow	Red	Green	
Neuromimétisme	Matérialisme radical	Red	Red	Green	Red	Grey	Grey	Grey	Grey	Grey	Grey	Green	Red	Yellow	Green	Green	
	Matérialisme rationnel	Red	Green	Red	Red	Grey	Grey	Grey	Grey	Grey	Grey	Green	Yellow	Yellow	Green	Green	
Connexionnisme	Automate cellulaire	Green	Red	Yellow	Red	Grey	Grey	Grey	Grey	Grey	Green	Yellow	Yellow	Green	Green	Green	
	Intelligence collective	Yellow	Red	Green	Red	Grey	Grey	Grey	Grey	Grey	Green	Yellow	Yellow	Green	Green	Green	
	Auto-organisation neuronale	Yellow	Red	Green	Red	Grey	Grey	Grey	Grey	Grey	Green	Yellow	Yellow	Green	Green	Green	
	Système de codage	Green	Red	Yellow	Red	Grey	Grey	Grey	Grey	Grey	Green	Green	Green	Red	Green	Red	
	Théorie de l'information	Green	Red	Red	Red	Grey	Grey	Grey	Grey	Grey	Green	Green	Green	Green	Green	Green	
Fonctionnalisme écologique	Apprentissage par renforcement	Red	Green	Green	Red	Red	Red	Red	Yellow	Red	Red	Yellow	Green	Green	Green	Yellow	
	Architecture décisionnelle	Red	Green	Green	Red	Red	Red	Red	Yellow	Red	Red	Yellow	Green	Green	Green	Green	
Évolutionnisme	Évolutionnisme physique	Red	Green	Green	Red	Grey	Grey	Grey	Grey	Grey	Grey	Yellow	Green	Green	Green	Yellow	
	Évolutionnisme neuronal	Red	Green	Green	Red	Grey	Grey	Grey	Grey	Grey	Grey	Yellow	Green	Green	Green	Green	
	Évolutionnisme comportemental	Green	Green	Red	Red	Grey	Grey	Grey	Grey	Grey	Grey	Yellow	Green	Green	Green	Green	
Interactionnisme	Évolutionnisme du langage	Green	Red	Green	Red	Grey	Grey	Grey	Grey	Grey	Grey	Yellow	Green	Green	Green	Yellow	
	Phénoménologie	Yellow	Red	Red	Green	Red	Red	Red	Yellow	Yellow	Red	Grey	Grey	Grey	Grey	Grey	
	Constructivisme piagétien	Yellow	Red	Red	Green	Red	Red	Red	Red	Red	Red	Green	Green	Yellow	Red	Green	
	Constructivisme varelilien	Yellow	Red	Yellow	Green	Red	Red	Red	Red	Red	Green	Green	Yellow	Red	Green	Green	

Tableau I-12 : Récapitulation des positions de différentes approches des courants de la cognition artificielle vis-à-vis des hypothèses et méthodes constituant la grille de lecture de ce chapitre.

Toutefois, il faut remarquer que seul l'interactionnisme se revendique de la famille FA. Plus précisément, au sein de l'interactionnisme, ceux qui se réclament d'une hypothèse ontologique rentrent en conflit avec la possibilité de tendre vers une connaissance sur le monde physique ou idéal. Au-delà de l'inaccessibilité de ces connaissances, l'interactionnisme n'indique aucun moyen pour tendre vers ces connaissances. Autrement dit, les autres familles métaphysiques soutenant l'interactionnisme se trouvent soumises à une tension intenable. Mais paradoxalement, le constructivisme issu de la famille FA peine à concevoir un modèle précis autorisant une implémentation. Au mieux, les réalisations dans le cadre de l'interactionnisme se transforment involontairement en une modélisation déterminant le domaine de fonctionnement du système, ce qui autorise une interprétation de l'interactionnisme à l'aide de la famille FC ou FP. Dans ce cas, le point commun à toutes les tentatives de réalisation d'un robot cognitif serait que celles-ci s'appuient sur l'une des trois premières familles métaphysiques : FC, FCP ou FP, c'est-à-dire qu'elles adhèrent à au moins une des deux hypothèses ontologiques. Par conséquent, l'interactionnisme offrirait un point de départ privilégié à la cognition artificielle, d'autant plus que ce courant avance une argumentation convaincante pour rejeter définitivement les hypothèses concernant la

perception et la conception tout en proposant une théorie suffisamment riche pour expliquer les données éthologiques et psychologiques.

**

L'identification de ce point commun aux courants de la cognition artificielle permet d'examiner dans quelles mesures il génère une impasse conceptuelle pour la robotique cognitive et ainsi compléter éventuellement l'interactionnisme de telle sorte à pouvoir proposer un formalisme destiné à concevoir des systèmes autonomes. Plus précisément, le prochain chapitre montrera que les hypothèses ontologiques conditionnent le projet épistémologique, c'est-à-dire l'unification de la connaissance scientifique qui renvoie directement à la question de la vérité.

CHAPITRE II ANALYSE DE L'IMPASSE ET DE SON DÉPASSEMENT PAR LE PRAGMATISME

« HAMM. — [...] Père ! (*Un temps. Plus fort.*) Père ! (*Un temps.*)
Va voir s'il a entendu.

Clov va à la poubelle de Nagg, soulève le couvercle, se penche dessus. Mots confus. Clov se redresse.

CLOV. — Oui.

HAMM. — Les deux fois ?

Clov se penche. Mots confus. Clov se redresse.

CLOV. — Une seule.

HAMM. — La première ou la seconde ? »

Fin de Partie p. 88-89 de Samuel Beckett (1957).

« [...], c'est tout ce que je sais, je dis je en sachant que ce n'est pas moi, moi je suis loin, [...] »

L'innommable p. 197 de Samuel Beckett (1953).

« Il faut être deux pour inventer. L'un forme des combinaisons, l'autre choisit, reconnaît ce qu'il désire, et ce qui lui importe dans l'ensemble des produits du premier. »

De la simulation p. 620 de Paul Valéry (1927) cité par F. Varela (1989).

« L'absurde n'a de sens que dans la mesure où l'on n'y consent pas. »

Le mythe de Sisyphe p. 50 d'Albert Camus (1942).

1. Introduction

Afin d'identifier l'impasse de la robotique cognitive, le premier chapitre a voulu présenter les différentes approches et relever ainsi leurs ressemblances et leurs dissemblances. Toutefois, la grande disparité entre les approches en sciences cognitives a nécessité au préalable d'établir une grille d'analyse commune suffisamment vaste et fondamentale pour couvrir tout le spectre conceptuel des théories cognitives. Les hypothèses constituant cette grille abordent l'ontologie, la perception, la conception, les méthodes d'investigation et les types de description. Ces hypothèses se veulent indépendantes les unes des autres bien qu'elles se rattachent à trois problématiques classiques de la philosophie qui, elles, se chevauchent : les théories de la vérité, les théories de la connaissance et les théories épistémologiques. Brièvement, les théories de la vérité se penchent dans un premier temps, sur la signification du qualificatif « vrai » pour un énoncé. Par exemple, que signifie l'affirmation : l'énoncé E1 « la neige est blanche » est vrai ? Est ce que le terme vrai renvoie à une vérité absolue ou à une simple convention ? Dans un second temps, elles s'intéressent aux conditions justifiant l'emploi de ce qualificatif, autrement dit, les critères en soi de la vérité : Comment s'assurer de la relation entre E1 et le monde ? Comment certifier que l'énoncé E1 « la neige est blanche » soit vrai ? Concernant ces deux points, il a été montré au début du premier chapitre que l'interprétation des hypothèses ontologiques est déterminante. Les théories de la connaissance, qui forment la gnoséologie, s'articulent également autour de deux thèmes. Le premier thème porte sur la définition de la connaissance qui peut se

comprendre comme une réflexion sur la nature de l'activité mentale touchant aux croyances, aux représentations ou aux imaginations. Par rapport, à l'exemple précédent, le sujet d'étude devient l'idée de neige ou de blanc et leurs manipulations. Parallèlement, le second thème s'occupe de la justification pour soi et pour les autres de ces connaissances, ce qui amène, par exemple, à s'interroger sur la possibilité d'asserter l'énoncé E1, ce qui est différent de vouloir s'assurer de la justesse de l'assertion. Ces deux thèmes regroupent des problématiques très hétérogènes mais elles incluent généralement des réflexions liées aux hypothèses concernant la perception et la conception. Enfin, les théories épistémologiques sont très proches de l'objectif des théories de la connaissance mais elles privilégient l'étude de la constitution du savoir scientifique, autrement dit, comment collecter et développer un ensemble d'énoncés traduisant la connaissance sur le monde. Par conséquent, les théories épistémologiques se trouvent étroitement liées à l'évaluation des méthodes d'investigation et des types de descriptions. Toutefois, la valeur de ce projet épistémologique et sa forme dépendent de la conception de vérité admise. En reprenant l'exemple de l'énoncé E1, la perspective épistémologique orienterait l'interrogation vers : comment intégrer E1 à un ensemble d'énoncés déjà constitués ? Et comment convaincre la communauté du bien-fondé de cette intégration ?

Néanmoins, la circonscription et la distinction de ces problématiques au sein des théories cognitives ne s'effectuent pas spontanément, de la même manière que la décomposition de la perception en un jugement existentiel et en un jugement perceptif n'apparaît pas immédiatement parce que, en pratique, ils fusionnent constamment. Aussi, afin de clarifier le discours, le premier chapitre s'est concentré exclusivement sur les hypothèses identifiées, mettant ainsi ces diverses théories en arrière plan. Cependant, la conclusion de cette étude philosophique des approches en cognition artificielle incite à examiner plus attentivement l'une de ces trois problématiques. En effet, le point commun entre toutes les tentatives de réalisation d'un système cognitif réside dans l'acceptation ou la reconnaissance d'au moins une hypothèse ontologique, contrairement aux autres hypothèses et méthodes qui sont alternativement rejetées ou acceptées. En conséquence, l'impasse de la robotique cognitive se situe sur la question de la vérité. Explorer plus minutieusement les théories de la vérité devient indispensable soit pour compléter les hypothèses ontologiques afin de palier un défaut conceptuel, soit pour comprendre en quoi la conception ontologique de la vérité est viciée afin de mieux s'en prémunir dans la construction d'une théorie cognitive.

Répondre à cette alternative sera l'enjeu de ce chapitre qui se constituera de deux parties. La première abordera de manière critique les définitions classiques de la vérité, selon un regard philosophique et logique. La synthèse de ces positions révélera que toutes les théories de la vérité possèdent en leur sein des contradictions empêchant une définition unique de la vérité, ce qui conduirait à renoncer à une quelconque naturalisation de la raison, autrement dit à la fabrication d'une cognition artificielle. Toutefois, en seconde partie, une perspective pragmatiste offrira une définition de la vérité et de l'absolu qui se voudra suffisamment flexible pour concilier les contradictions précédentes. Ces définitions serviront alors à comprendre le rôle de la cognition et ainsi à imaginer une généalogie de la cognition étayée par des observations éthologiques.

2. Critique du concept de vérité

La question de la signification de la vérité appartient aux plus vieilles interrogations de la philosophie. Sans vouloir retranscrire une histoire du concept de vérité, cette partie offrira simplement une synthèse des différentes positions résultantes des débats sur le concept de vérité et non une synthèse de ces derniers. Le travail de synthèse avancé s'inspire essentiellement de celle effectuée par Engel (1998) et par celles de Putnam (1971 ; 1983). Toutefois, afin de conserver une cohérence conceptuelle, le vocabulaire employé au cours de cette partie sera le même que celui du premier chapitre.

Parmi la pluralité de position existante, quatre définitions canoniques de la vérité ont été dégagées. Respectivement, chacune de ces définitions assimile la vérité à : une correspondance, une cohérence, une efficacité ou un consensus. Les trois premières définitions reposent sur l'acceptation d'au moins une hypothèse ontologique. Plus précisément, chacune de ces trois définitions se trouve principalement sous-tendue par l'une des trois familles philosophiques revendiquant un réalisme métaphysique. La critique de ces trois définitions constituera alors l'analyse proprement dite de l'impasse de la cognition artificielle.

La quatrième position concernant la vérité, soutenue par la seule famille de métaphysique anti-réaliste (FA), avance une définition déflationniste de la vérité pour solutionner ces paradoxes. En contrepartie, elle enlèvera également le sens commun généralement attribué au mot « vrai », pénalisant l'établissement d'une théorie cognitive formelle et par conséquent de l'artificialisation de la cognition.

L'ensemble des critiques sur ces différentes définitions de la vérité indiquera alors la direction vers laquelle la recherche d'une alternative en cognition artificielle doit s'orienter.

2.1. La vérité comme une correspondance

La première théorie de la vérité comme correspondance apparaît comme la conception la plus intuitive et la plus communément usitée. Paradoxalement, sa présentation se révèle complexe et nécessite de décomposer l'exposé en trois points. (A) Le premier point portera sur les fondements de la théorie de la vérité comme correspondance et sur la philosophie de la logique induite. (B) Le deuxième point s'attachera à expliquer le projet épistémologique et le rôle de la cognition dans cette perspective. (C) Enfin les différentes critiques à l'encontre de cette définition de la vérité seront avancées.

A - Implications philosophiques et logiques de la vérité comme une correspondance

La théorie de la vérité comme correspondance comprend l'énoncé « *'sur la table se trouvent trois pommes'* est vrai » comme une confirmation qu'il existe réellement trois pommes situés sur la table au moment où l'énoncé est prononcé. Ce dernier n'ajoute aucune propriété nouvelle à la scène décrite par *'sur la table se trouvent trois pommes'*, le prédicat vrai qualifie de manière transcendante la relation entre la proposition *'sur la table se trouvent trois pommes'* et le fait qu'il existe dans le monde extralinguistique trois pommes sur la table.

Cette interprétation revient également à comprendre l'un des principes fondateurs de la logique *P1* « *x est identique à x* » comme « l'affirmation de la réflexivité de la relation d'identité : toute chose soutient avec elle-même cette relation » (Putnam, 1971). En d'autres termes par cette relation toute chose *se reconnaît* dans son reflet, dans l'idée de sa propre existence. Une pomme-objet réelle se reflète dans l'ensemble des idées décrivant toutes ses propriétés ou ses caractéristiques. Aucune propriété de cette pomme n'échappe au monde idéal puisque c'est celui-ci qui confirme son existence. Cette conception de la vérité s'inscrit obligatoirement dans la famille FCP puisqu'elle nécessite une définition essentialiste à la fois pour les objets physiques et pour les concepts au sens large, expliquant l'indépendance des jugements ou des opinions portés sur eux.

Cette conception de vérité peut être qualifiée de synthétique dans la mesure où elle ne dépend pas uniquement du monde idéal. Cependant, la vérité synthétique ainsi définie se réduit à la confirmation existentielle des objets (ou des situations) réels passés, présents ou futurs. Or, un grand nombre d'énoncés est qualifié de vrai sans pour autant porter directement sur l'existence d'objets physiques. Tout en restant dans le cadre d'une théorie de la vérité comme correspondance, une seconde forme de vérité apparaît lorsque seul le monde idéal est en jeu.

La vérité comme correspondance possède ainsi plusieurs facettes dont l'étude passe par : (i) la définition des propositions simples qui justifie une certaine interprétation des fondements de la logique et (ii) celles des propositions composées et des schémas logiques qui servira de base pour concevoir une épistémologie.

i - Définition des propositions simples selon la famille FCP

Précédemment la vérité se centrait sur l'*être*, le concept particulier (Jean recouvre l'objet Jean, « Jean est »). Toutefois, certaines propositions produisent un glissement de la notion de vérité de l'*être* vers l'*avoir*, bien que le verbe « être » puisse être conservé : « Jean est grand ». En effet, la notion de vérité porte non plus sur l'existence des objets physiques mais sur l'attribution de propriété ou sur la relation avec d'autres objets. Ainsi, une proposition simple se définit comme un propos exprimant une relation entre un sujet et un prédicat. Par exemple, « cette pomme est rouge » traduit une mise en relation entre le sujet « cette pomme » et le prédicat « rouge ». La proposition minimale « cette pomme est » revient à la notion stricte de la vérité synthétique développée précédemment. L'ajout du prédicat révèle l'existence d'une relation qui étend le concept de vérité synthétique. En effet, l'application de la vérité synthétique suggérée par l'énoncé « cette pomme est rouge » porte sur l'existence de la relation entre un objet et un concept générale et non plus sur l'existence de ceux-ci. Utilisé par commodité, le verbe « être » exprime ici « a la propriété d'être ». Cette reformulation révèle une dissymétrie entre le sujet et le prédicat, autrement dit, la relation est orientée. En reprenant l'exemple de la pomme, la phrase « ce rouge est pomme » ne possède plus du tout le même type de sémantique que l'énoncé « cette pomme est rouge », si sémantique peut encore y être associée. L'objet formant ici une unité sur laquelle se reflète des propriétés, la dissymétrie est naturelle.

Toutefois, la vérité d'une proposition simple ne se réduit pas à la reconnaissance d'une propriété appartenant à l'ensemble des propriétés attribuées à un objet réel. En effet, la forme des propositions exprimées jusqu'ici appartient à un *premier type de proposition* parmi les trois qui peuvent être dégagés. Le *deuxième type de proposition* provient du fait que la structure du monde physique se reflète également dans le monde idéal, conférant un sens aux relations entre concepts. Comparativement, le premier type de proposition porte sur la relation entre un concept et un objet (comme dans « tous les hommes sont mortels ») et le second type de proposition porte sur la relation entre des concepts (comme dans « l'homme est mortel » traduisant le fait que le concept d'homme implique le concept de mortalité). Le *troisième type de proposition* contient toutes les propositions statuant sur des objets imaginaires ou de fiction, c'est-à-dire des objets construits à partir de la combinaison de concepts appartenant à d'autres objets réels, formant ainsi un univers conceptuel clos sans instance dans le monde physique bien qu'inspiré par celui-ci. Par exemple, « Pégase est un cheval volant » est vrai par rapport à l'imaginaire collectif bien que Pégase n'ait jamais existé.

En somme, fonder sur la correspondance, la vérité d'une proposition simple se résume à la reconnaissance :

- soit d'une propriété appartenant à l'ensemble des propriétés attribuées à l'objet réel,
- soit de la validité de la relation entre deux propriétés ou concepts issue du reflet du monde,
- soit de la validité de la relation entre deux propriétés ou concepts résultant de la combinaison de concepts encré dans le monde physique.

Cette définition de la vérité pour la logique prédicative implique que la vérité logique ne porte pas sur la formulation des propositions mais sur leur signification. En d'autres termes, seuls

les énoncés ayant une signification sont des propositions et le monde idéal se compose d'une infinité de proposition ayant soit une valeur vraie soit une valeur fausse. Par exemple, l'énoncé fantaisiste de Putnam (1971) « tous les boojums sont des snarks » ne constitue pas une proposition. Un énoncé écrit dans un certain langage permet d'exprimer une idée mais ne correspond pas à l'idée : il y a une distinction entre le signifiant (l'énoncé) et le signifié (ce qu'il vise). Le souci de la signification vaut également pour des phrases correctement construites avec un vocabulaire courant, ainsi, l'énoncé de Russell entendu au 20^{ème} siècle hors d'un cadre fictionnel : « L'actuel roi de France est chauve » n'a pas de signification car le sujet ne reflète aucune réalité, donc la vérité de relation proposée ne peut-être statuée et cet énoncé ainsi formulé n'est pas une proposition.

La vérité comme une correspondance peut servir également à la justification métaphysique de deux autres principes fondamentaux de la logique classique :

- Le principe de bivalence stipule qu'une proposition p est soit vraie soit fausse, $P2$ « p ou ($non\ p$) ».
- Le principe de contradiction réproouve l'inconsistance de toute proposition p , $P3$ « p et ($non\ p$) ». En d'autres termes, une proposition et sa contradiction ne peuvent être toutes les deux vraies.

L'acceptation de ces deux principes permet d'introduire la notion de tiers exclus ou le raisonnement par l'absurde où montrer que la proposition $non\ p$ est fausse implique obligatoirement que p est vrai. L'interprétation et la justification de ces principes logiques seront contestées par les autres définitions de la vérité mais ici elles conditionnent l'interprétation de la logique dans son ensemble.

ii - Les propositions composé et les schémas logiques

Les propositions simples ne représentent qu'un sous-ensemble des propositions existantes. Les autres propositions représentent des propositions composés c'est-à-dire un ensemble de propositions reliées par des connecteurs logiques : « sur la table se trouvent trois pommes *et* il pleut ». La particularité des propositions simples est que leur valeur de vérité réside uniquement dans l'adéquation entre un fait physique ou conceptuel et un énoncé, contrairement aux propositions composées dont la valeur de vérité dépend de celles des propositions qui la compose et des connecteurs qui les relis. La logique Booléen correspond dans une certaine mesure à un moyen de *révéler* mécaniquement la valeur de vérité de propositions composées (avec les connecteurs de conjonction et disjonction) en connaissant la valeur de vérité des propositions.

Une autre manière de *découvrir* de nouvelles propositions vraies consiste à en déduire à partir d'autres à l'aide d'une règle d'inférence représentant un schéma logique. Par exemple, le syllogisme, un autre principe fondamental de la logique ($P4$) se traduit comme suit : « tous les A sont B », « tous les B sont C » donc « tous les A sont C ». Les termes A , B et C désignent ici respectivement des éléments au sein d'une classe A , d'une classe B et d'une classe C . L'inférence se comprend alors comme la transitivité de la relation « être une sous-classe de » : « la classe A est une sous-classe de la classe B », « la classe B est une sous-classe de la classe C » donc « la classe A est une sous-classe de la classe C ».

Dans ce cadre, les systèmes formels se comprennent comme des outils conceptuels permettant de découvrir mécaniquement des vérités logiques. Toutefois, comprendre finement l'interprétation de ces outils conceptuels est un préalable pour distinguer la théorie de la vérité par correspondance et par cohérence. Un système formel décrit les conventions d'écriture et de manipulations de formules, soit décrit un langage selon le quadruplet suivant : un alphabet, des

règles de formation des mots, une théorie formelle correspondant à un ensemble d'axiomes (des mots donnés a priori) et enfin des règles d'inférence. Un théorème correspond à une formule qui se déduit mécaniquement par règles d'inférence à partir des axiomes. La démonstration d'un théorème représente la liste exhaustive des règles d'inférences utilisées à partir des axiomes afin de l'obtenir.

Une proposition démontrée à partir d'axiomes considérés pour vrais est également considérée pour vraie. Cependant, Gödel (1931) a démontré l'inconsistance des systèmes formels qui intègrent l'arithmétique, c'est-à-dire que ces derniers peuvent démontrer certaines formules malgré l'existence de contre-exemples. De plus, Gödel (1931) démontra qu'il existe des formules dont l'application soit toujours correcte bien qu'aucune démonstration ne puisse les prouver, cela s'appelle l'incomplétude. Le théorème de Turing (1936) précise que l'indécidabilité d'une formule ne peut être prouvée dans le système qui doit l'évaluer (cf. exemple du système MUI).

Les systèmes formels introduisent la notion de démontrabilité et, la notion de vérité n'est qu'une interprétation puisque les axiomes ou les valeurs n'ont pas de valeur de vérité proprement dite. La question de la relation entre ces deux notions constitue le cœur du débat entre réalisme et nominalisme dont l'approfondissement permettra de mieux cerner les implications épistémologiques de la vérité par correspondance.

La relation vérité/démontrabilité devient d'autant plus forte lorsque la notion vérité se trouve modélisée au sein de systèmes formels. En effet, les implications de la valeur de vérité des propositions ne sont pas prises en compte dans les systèmes formels évoqués précédemment. Afin d'étudier ces relations logiques, le prédicat de vérité doit apparaître directement dans le système formel représentant alors un modèle d'un monde considéré. Un système formel décrivant un langage, le prédicat de vérité, les règles de formations ou d'inférences du système formel appartiennent au métalangage qui décrit le système formel. Autrement dit, l'évaluation du prédicat de vérité ne porte pas sur les règles d'inférences. Par exemple, « si certains *A* sont *B*, alors tous les *A* sont *B* » représente une règle d'inférence possible bien qu'elle soit non valide puisque des contre-exemples existent. Étudier la validité de ces règles d'inférences peut s'effectuer par la formalisation du système formel qui les emploie par un méta-métalangage permettant de qualifier la valider par un prédicat de vérité au sein de ce nouveau langage.

Dans ce cas, la question de l'équivalence entre la démonstration d'une proposition évaluée vrai et la vérité se pose : soit la notion vérité peut se réduire à l'évaluation du prédicat de vérité au sein d'un langage (position nominaliste), soit la notion vérité demeure distinct de sa modélisation au sein des systèmes formels (position réaliste). Comprendre la vérité comme une correspondance conduit à la seconde position. Plus précisément, Putnam (1971) défend la position réaliste à l'aide de multiples arguments dont deux principaux. Le résumé des deux arguments majeurs permettra de mieux comprendre la valeur d'un schéma logique et la nature de l'épistémologie de la logique et des mathématiques dans le cadre de la vérité par correspondance.

Le *premier argument de la position réaliste* revendique la différence entre l'ensemble des propositions et l'ensemble des formules et par conséquent la différence entre la notion de vérité et la notion de prédicat de vérité. En effet, la position nominaliste considère que les classes *A*, *B*, *C* de l'inférence « tous les *A* sont *B* », « tous les *B* sont *C* » donc « tous les *A* sont *C* » représente des ensembles de mots alors que la position réaliste considère que ces classes regroupes des choses. En reprenant les exemples de Putnam (1971), l'énoncé « Si tous les corbeaux sont noirs et si toutes les choses noires absorbent la lumière, alors tous les corbeaux absorbent la lumière » et l'énoncé fantaisiste « Si tous les boojums sont des snarks et si tous les snarks sont des eggelumphs, alors tous les boojums sont des eggelumphs » sont vraies *logiquement* selon la position nominaliste puisque toutes les instances de substitution des schémas logiques ne sont que des mots, la syntaxe supplante la signification. Or les schémas logiques (ou règles d'inférences) reposent sur des propositions simples qui doivent par définition avoir une signification, une

correspondance. Par conséquent, « Si tous les boojums sont des snarks et si tous les snarks sont des eggelumphs, alors tous les boojums sont des eggelumphs » n'est pas une proposition, la notion de démontrabilité peut-être appliquée mais pas celle de vérité. La vérité ne découle pas uniquement de la ressemblance syntaxique avec un schéma logique valide.

Les schémas logiques révèlent ici la structure logique des propositions et se définissent uniquement par rapport à l'ensemble des propositions existantes donc signifiantes. De la même façon qu'une proposition simple révèle une structure logique de la relation entre le sujet et le prédicat. La validité d'un schéma logique correspond alors à une correspondance d'ordre supérieur. Le développement d'une logique d'ordre supérieur qui statuerait sur la vérité des schémas logiques des propositions se comprend comme une modélisation de ces structures. Contrairement au nominalisme où l'abstraction s'applique au prédicat de vérité, le réalisme implique une abstraction de la mise en correspondance.

Le *second argument de la position réaliste* souligne que la position nominaliste conduit à ce que la vérité soit toujours relative à un langage puisque la validité d'un schéma logique S revient à affirmer que « toutes les instances de substitutions de S selon le langage L sont vraies ». Cette formulation subit le théorème de la théorie des ensembles stipulant « qu'aucun langage L ne peut contenir de noms pour toutes les collections d'objets susceptibles d'être formées, tout au moins dans le cas où le nombre de ces objets est infini » Putnam (1971). Autrement dit, cette définition génère alors une série infinie de notions de vérité : validité selon L_1 , validité selon L_2, \dots . Pour Putnam (1971), cela signifie que le prédicat de vérité d'un système formel ne parviendra jamais à l'absolu puisqu'il dépend toujours d'un langage particulier et limité, contrairement à la vérité intuitive de la position réaliste qui permet d'entrevoir la validité universelle de « tous les A sont B , tous les B sont C donc tous les A sont C ».

*

La conception de la logique issue d'une théorie de la vérité comme correspondance perçoit le formalisme comme un moyen, un instrument, destiné à faciliter le dégagement des propositions et des schémas logiques dans le paysage idéal ainsi que la discussion sur leur valeur de vérité. Ces outils sont nécessaires pour éviter et débusquer les illusions ou les hallucinations de la conception, de la même manière que le physicien utilise des instruments pour étudier les objets physiques et identifier les illusions ou les hallucinations de la perception. Les questions de fondements de la logique ne se posent pas, les vérités logiques sont découvertes. La preuve logique s'identifie à une expérience mentale assurant les enchaînements où la retranscription symbolique sert de mémoire externe. Par ailleurs, une autre manière d'assurer la maîtrise d'une entité conceptuelle revient à multiplier les raisonnements à son égard, de la même façon que la multiplication des actions sur un objet assure de sa réalité et permet de le connaître. Les méthodes des sciences logiques et mathématiques possèdent alors de forte ressemblance avec les méthodes des sciences empiriques comme l'illustre le premier tiers de la discussion imaginaire au sein d'une classe d'école de Lakatos (1984) sur la définition d'un polyèdre.

Le formalisme radical se comprend ici comme une position qui confondrait fin et moyen. De sorte qu'il devient légitime d'avancer qu'il existe, par exemple, une notion imprédictive d'ensemble (Putnam, 1971) dépassant les restrictions imposées par la théorie des ensembles, autrement dit la signification l'emporte sur la formalisation.

De nombreuses théories de vérité comme correspondance peuvent être élaborées. Toutefois, la notion de correspondance implique nécessairement l'acceptation des hypothèses ontologiques HOI et HOP. La notion de correspondance conduit également à la conception de deux types de vérités : vérité synthétique et vérité analytique. Par ailleurs, cette dernière renvoie à plusieurs, voire une infinité, de types de correspondance qui ne dépend pas de la formalisation du discours logique mais de l'abstraction des objets idéaux quelle vise. Cependant, tout cet édifice

conceptuel repose sur la correspondance fondamentale entre le monde physique et le monde l'idéal. Or, c'est justement sur ce point précis que se situeront la plupart des critiques.

B - Le projet épistémologique et le rôle de la cognition pour la vérité comme correspondance

L'adoption intuitive du réalisme naïf consistant à croire que les choses sont telles qu'elles se présentent et à leur attribuer toutes sortes de qualités favorise l'acceptation implicite de HOP puis de HOI, une fois que ces choses et ces qualités se trouvent représentées au sein d'un langage. Ainsi, la théorie de la vérité comme une correspondance se trouve souvent acceptée tacitement par la communauté scientifique qui maintient alors une épistémologie évolutionniste au sens de Putnam (1983).

Celle-ci se comprend comme un programme de recherche visant à *découvrir* les objets et les concepts qui possèdent une existence propre, c'est-à-dire indépendamment du sujet qui les observe. La méthode hypothético-déductive (MHD) se trouve privilégiée ainsi que la description componentielle (DC) du monde physique. Dans le cadre de la vérité comme une correspondance, principalement deux familles épistémologiques apparaissent selon les différentes opinions concernant les théories de la connaissance. La *première famille* regroupe les épistémologies conservant le réalisme naïf considère que la validation d'un énoncé scientifique vaut pour l'éternité et permet d'accéder à un niveau supérieur, un nouvel espace de découverte. Le projet épistémologique devient la découverte et la description pas à pas du monde physique et du monde idéal. L'édifice scientifique progresse linéairement par essai-erreur bien que perturbé par les aléas historiques. La *seconde famille* rassemble les épistémologies qui considèrent que les objets et concepts objectifs existent mais que l'esprit humain ne parviendra jamais à en rendre compte totalement. Les hypothèses concernant le monde représentent uniquement des approximations qui n'auront de cesse de s'ajuster, de s'améliorer à l'infini tendant ainsi vers la vérité. Cette conception conserve également l'idée d'un progrès scientifique tendant vers la vérité mais que les énoncés validés restent susceptibles d'être modifiés par des hypothèses plus fortes et que le « choix » dans le champ des possibles des hypothèses approximatives résulte des pressions socioéconomiques.

La linéarité et l'idée de progrès se trouvent profondément ancrés dans la communauté scientifique bien que cela soit remis en cause par les discours contraires issus de l'histoire de science (Khun, 1970 ; Feyerabend, 1975) ou de la sociologie des sciences (Latour, 1991). Cette situation s'explique par l'écart entre l'enseignement de la *science qui est faite* (idéalisé voire mystifié) et la *science en train de se faire* avec ses mécanismes socio-économique (Latour, 1991). La culture de leur formation ne permet pas aux scientifiques de percevoir les problématiques épistémologiques mais cela n'empêche pas que concrètement la production scientifique continue indépendamment de la connaissance de ces problématiques épistémologiques. Cependant, l'indépendance de la production scientifique avec les problématiques épistémologiques devient caduque lorsque le sujet d'étude se trouve être la cognition.

En effet, affirmer que l'activité scientifique tend vers une description unique approchant de la vérité implique que l'homme soit prédisposé à découvrir la vérité, autrement dit la cognition est une disposition permettant de découvrir la vérité. Sans tenir compte des arguments sur l'avantage prétendument sélectif de cette disposition discrédité lors de la présentation de l'évolutionnisme au cours du premier chapitre, l'acceptation d'un tel projet épistémologique induit par la vérité comme une correspondance signifie que toutes les approches appartenant à la famille FCP considèrent explicitement ou implicitement la cognition comme une disposition à découvrir la vérité. La critique de la définition de la vérité comme une correspondance conditionne alors toutes ces approches quels que soient leur courant d'origine (symbolisme, subsymbolique, neuromimétisme, fonctionnalisme écologique, ou évolutionnisme).

C - Les différentes critiques à l'encontre d'une vérité comme une correspondance

Les objections à l'encontre de la théorie de la vérité comme une correspondance présentées ici ont déjà été abordées indirectement à des degrés divers lors de la critique des courants animant les sciences cognitives adhérant à la famille FCP. L'analyse de la notion vérité par correspondance permet de mieux comprendre l'origine de ces critiques et leurs portées. Au total six critiques portant cette conception de la vérité seront avancés dans cette section sans revenir pour autant sur les critiques des théories de la connaissance ou sur les théories cognitives induites par celles-ci et évoquées dans le premier chapitre. Enfin, les deux dernières des six critiques seront considérées comme suffisamment fortes pour discréditer la théorie de la vérité comme une correspondance.

i - Le problème de la négation

La première critique présentée ici provient de Russel (1969) : la notion de correspondance a-t-elle un sens avec les faits négatifs comme « Mon verre n'est pas vide » ? La correspondance s'applique t-elle avec l'absence du fait positif ou avec un fait négatif spécifique ? La difficulté apparaît particulièrement dans le cadre d'une approche symbolique. Par exemple, pour définir le déclencheur d'un fait négatif d'un automate, deux options existent : soit connaître tous les autres faits négatifs impliquant l'omniscience ; soit comparer constamment l'état des capteurs avec le fait redouté. Dans les deux cas, aucune détection directe du fait-négatif n'est réalisée.

ii - Le problème de la dépendance

La deuxième critique possède une certaine similarité avec la première. Selon la théorie de la vérité comme une correspondance, la valeur d'une proposition dépendant de son adéquation avec le fait quelle sous-tend. Ainsi la valeur d'une proposition doit être indépendante de la valeur des autres propositions. Or, le fait que « ceci est rouge » soit vrai implique « ceci n'est pas vert » de même, « Naples se trouve au sud de Rome » implique que « Rome se trouve au nord de Naples. Cette seconde critique porte davantage sur la légitimité de l'atomisme logique au sein de la vérité comme correspondance que cette dernière. En effet, la notion de structure ou de schéma logique intègre le fait que le monde physique soit structuré et finalement la compréhension du monde s'effectue plus par la reconnaissance de ces structures que la collection de correspondance entre fait et proposition (Wittgenstein, 1921). Dans les théories cognitives, cette différence s'illustre par l'opposition entre le symbolisme qui se focalise sur cette correspondance et le subsymbolique qui s'intéresse davantage aux structures qui lient les concepts.

iii - Le problème de la signification et dénotation

La troisième critique porte sur la difficulté à distinguer la correspondance entre un objet idéal et un objet physique et la correspondance entre l'idée d'une personne et l'objet quel vise. La question peut également être présentée comme la distinction entre signification et dénotation (Frege, 1892) qui s'illustre par deux remarques.

La première remarque, déjà formulée dans le chapitre précédent, souligne les différentes interprétations de l'énoncé « je vois Hesperus » selon les croyances de celui qui le prononce. La subjectivité devient une condition de vérité et de ce fait la signification d'une proposition ne se réduit plus à la dénotation. Une manière de diminuer l'impact de cette remarque consiste à considérer que le langage naturel, forgé sur les relations humaines avec leurs psychologies et non sur la description de la réalité, biaise toute description du monde. Par conséquent, il devient légitime de restreindre les propositions justes au fait sensoriel évacuant ainsi le problème des

croiances subjectives dans l'interprétation de la signification d'une proposition ce qui conduit au positivisme logique.

La seconde remarque vient du fait que le problème de la distinction entre signification et dénotation dépasse le problème de la condition de vérité. « Etoile du matin est l'étoile du soir » et « Etoile du matin est l'étoile du matin » sont deux énoncés vrais. Leurs dénotations se trouvent équivalente ainsi que leur valeur de vérité mais pourtant l'une est informative et l'autre est tautologique. Cette remarque fait écho à celle déjà rencontrée dans le symbolisme pour indiquer qu'une proposition est avant tout une assertion c'est-à-dire un acte de langage avec ses conséquences. La définition de la signification est alors étroitement liée à l'idée de réaction suscitée.

La prise en compte de ces deux remarques conduit la théorie de la vérité comme une correspondance à rejeter le langage naturel afin de se maintenir. Ce changement s'illustre par l'abandon du symbolisme et par le repli des partisans de la famille FCP dans les théories cognitives tels que le subsymbolique, le neuromimétisme, le fonctionnalisme écologique, etc. Ces mouvements comprennent la sensation comme la mise en correspondance élémentaire avec le monde sur la quelle les réelles propositions pourront être découvertes et ainsi évacuer la psychologie de la description du monde, distinguer la signification (psychologique, subjectif et toujours lié à l'activité humaine) de la dénotation (objective et neutre).

iv - Le problème de l'incertitude du langage

La quatrième critique porte sur l'incertitude du langage qui a été développée lors de présentation des différents types de physicalisme dans le connexionnisme. En résumé, « il existe une multiplicité de manière distinctes de projeter nos propositions et leurs éléments sur la réalité, qui s'accordent toutes avec cette réalité, mais qui sont néanmoins incompatible » (Engel, 1998).

v - Le problème de la correspondance

L'avant dernière critique porte sur la légitimité à imaginer une mise en correspondance. « Un accord ne peut être total que si les choses en accord coïncident, donc ne sont pas de nature différente » (Frege, 1892) Or, la vérité par correspondance soutient le contraire, en dehors du problème relatif à théorie de la connaissance à savoir si l'idée d'un sujet correspond à l'idée objective du monde de idéal ou à une approximation de celle-ci. La correspondance devient ontologiquement impossible puisqu'elle « présuppose l'existence d'une relation entre deux choses différentes, alors que même cette relation semble être l'identité » (Engel, 1998).

vi - Le problème de la légitimité de la correspondance

La dernière objection porte sur la légitimité de n'importe quel jugement. Le jugement synthétique statue sur la correspondance entre une proposition p et les faits. Autrement dit, « p correspond aux faits » est équivalent « p est vrai », soit il est vrai que p si et seulement si p . Mais ce jugement dans sa réalisation constitue également un fait, qu'il faut également jugé (la proposition p « p correspond aux faits » est-elle vrai ? Soit, est-ce que p correspond aux faits ?), et ainsi de ce suite avec ce nouveau jugement. La mise correspondances comprise comme un répertoire de conditions de vérité (Comment cela correspond ?) tombe dans un problème récursif qui la rend intenable pour une définition objective et immuable de la vérité (Kant).

*

Les diverses réponses à ces critiques contribuent à la diversité des interprétations de la théorie de la vérité comme une correspondance et expliquent partiellement les différences entre les approches en cognition artificielle soutenue par la famille philosophique FCP. Surtout,

l'ensemble des critiques et plus particulièrement les deux dernières rendent l'adoption de la théorie de la vérité comme correspondante intenable. Par ailleurs, l'acceptation de ces critiques sur la correspondance désamorce également l'argumentation de Putnam (1981) concernant une logique réaliste puisque la signification d'une proposition ne se définit plus par sa correspondance avec un fait.

**

Ainsi, en acceptant une position métaphysique issue de la FCP, celle-ci conduit à concevoir une théorie de vérité comme une mise en correspondance. Cette conception de la vérité amène à percevoir la science comme un édifice d'énoncés tendant (plus ou moins linéairement) vers un *unique* idéal de savoir. Dans ce cadre, la cognition représente l'ensemble des moyens permettant de construire individuellement et collectivement cet édifice. Les différentes approches en sciences cognitives sous-tendus par la famille FCP (symbolisme, subsymbolique, neuromimétisme, fonctionnalisme écologique, ou évolutionnisme) vont se différencier alors sur la manière d'étudier la cognition mais non sur sa finalité qui est de retrouver ou d'approximer la réalité physique ou idéale. Néanmoins, contrairement aux autres domaines abordés par la science où la position épistémologique ou philosophique n'offre qu'un axe d'interprétation à la production scientifique ou une source de motivation, la conception de la cognition et la manière de mener l'étude se trouve profondément affectées et conditionnées par la définition de vérité adoptée. Or la vérité par correspondance étant très insatisfaisante, par conséquent, ces diverses investigations sont toutes vouées à l'impasse.

2.2. La vérité comme une cohérence

La deuxième définition de la vérité s'appuie sur la notion de cohérence. La présentation de cette définition s'effectuera en trois points. (A) Le premier point évoquera les raisons qui peuvent amener à considérer la cohérence comme un principe métaphysique au delà de sa simple utilité dans le processus d'acquisition de la connaissance. (B) Les conséquences de la vérité comme une cohérence seront alors explicitées notamment sous l'angle de la philosophie de la logique dans le deuxième point. (C) Le dernier point sera consacré aux projets épistémologiques et aux conséquences sur une définition de la cognition dans le cadre d'une position cohérentiste.

A - La cohérence comme un principe métaphysique

La cohérence d'un ensemble de propositions, de jugements ou de croyances se traduit par le fait qu'aucun élément ne remet en cause un autre élément. La première difficulté à élever ce principe au rang de critère de vérité réside dans la possibilité de toujours devoir introduire une nouvelle proposition, croyance ou jugement, pour rendre compatible deux éléments contradictoires. Une manière de répondre à cette difficulté consiste à considérer que l'ensemble des connaissances de base doit se développer à partir de faits avérés. Autrement dit, la vérité-correspondance devient un pré-requis à l'emploi la notion de cohérence. Toutefois, la notion de cohérence ainsi posée reste au niveau de la théorie de connaissance qui vise à comprendre les moyens d'accéder à la vérité et non à la définir. Ainsi, dans le cadre d'une position correspondantialiste radicale, la cohérence n'est qu'un outil conceptuel parmi d'autres.

Cependant, l'importance du rôle joué par la notion de cohérence au sein de la théorie de la connaissance et de l'épistémologie peut amener à penser que la cohérence n'est pas qu'un outil conceptuel mais une propriété du réel. La cohérence devient alors un critère métaphysique supplémentaire pour la vérité-correspondance qui permet par ailleurs d'atténuer les quatre premières critiques à son encontre.

Cependant, les deux dernières critiques demeurent intactes. De plus, l'holisme de Quine (1951) montre les limites de l'emploi de la cohérence pour établir des énoncés vrais dans le cadre de la vérité-correspondance. Une façon de surmonter ces critiques consiste à comprendre que ces tensions proviennent uniquement de l'hypothèse de la correspondance. Dans ce cas, la notion de cohérence reste le seul principe métaphysique pour définir la vérité. Ce choix conduit alors à renoncer à l'HOP et à conserver HOI puisque la notion de cohérence est un principe abstrait considéré comme indépendant du sujet. La vérité-cohérence soutenue par la famille FC peut se traduire dans les termes suivants : « *une proposition, un jugement ou une croyance p est vrai si et seulement si p appartient à un ensemble cohérent de propositions, jugements ou croyances* » (Engel, 1998)

B - Les conséquences de la vérité comme une cohérence

Dans le cas où seul HOI est adopté, seuls les concepts purs existent indépendamment du sujet, les concepts relatifs au monde physique sont toujours subjectifs ou intersubjectifs. Les concepts relatifs au monde physique sont construits collectivement à partir de l'expérience sensible. Ainsi, la vérité-cohérence signifie que la vérité d'une croyance d'un sujet dépend des autres croyances de celui-ci. Cette position conduit à un conventionnalisme physique compatible avec le caractère holistique de la connaissance montré par Quine (1951).

Concernant la découverte des concepts purs, celle-ci s'effectue progressivement par la construction mentale d'un édifice cohérent. La construction de cet édifice est positive, c'est-à-dire que le tiers exclu n'est plus considéré comme un principe logique fondamental. La formalisation joue un rôle prépondérant pour exprimer ces constructions mentales et éviter les erreurs de l'intuition, autrement dit, séparer dans la conception le jugement d'assertabilité et le jugement de vérité. Cependant, les limites du langage formel démontrées par Gödel suggèrent une différence entre la vérité « mathématique » qui porte sur l'idéal dont il est possible d'avoir une intuition et la vérité « logique » qui se réduit à la démonstration au sein d'un système formel. Il est possible d'interpréter cette difficulté comme une limitation intrinsèque de l'esprit humain qui lui interdit d'embrasser l'ensemble des vérités mathématiques.

Toutefois, quelles que soient les modalités pour intégrer de nouvelle connaissance (via l'intuition ou via un formalisme), il reste à déterminer les croyances initiales, ici le(s) concept(s) pur(s) initiaux, fondamentaux sur lesquels réside tout l'édifice de la vérité-cohérence. Dans cette démarche, Husserl a proposé une phénoménologie transcendantale, mais les critiques de Piaget à l'encontre de ce projet, qui ont été évoquées lors de la présentation de l'interactionnisme dans le premier chapitre, montrent son impossibilité. Par ailleurs, la cohérence oblige une connaissance initiale abstraite à partir de laquelle toutes les autres en découleraient ou se grefferaient, et cette connaissance initiale abstraite entraîne nécessairement une expérience métaphysique, or celle-ci ne peut se comprendre comme une mise en correspondance entre les idées du sujet et les concepts purs fondamentaux. Cette réintroduction inévitable de la notion de correspondance amène à une impasse puisque la notion de correspondance a été désavouée précédemment.

C - Projets épistémologiques et définitions de la cognition

Tout en évitant le débat sur l'accessibilité des idées mathématiques (l'intuition comme accessibilité directe des idées ou comme processus de découverte de celles-ci) qui appartient aux problématiques liées à la théorie de la connaissance, la simple idée de réduire la vérité à la cohérence amène à considérer tout système de connaissance en sursis tant que de nouvelles expériences physiques ou de pensées seront susceptibles de le remettre en cause en dehors de son principe fondateur : la cohérence. La révision étant inévitable du fait des limites intrinsèques de l'induction, la dialectique (MHD) se trouve être la méthode privilégiée pour étendre un système de savoir.

Dans ce contexte, le projet épistémologique ne vise pas la description unique du monde, son reflet, mais une simple description totalement cohérente du monde sans supposer qu'elle soit nécessairement unique. Selon Quine, retranscrit par Putnam (1983), trois critères conditionnent le fait qu'un tel système du monde soit correct « :

- 1- il lui faut conjecturer un certain nombre d'énoncés d'observation vrais par stimuli ;
- 2- il doit être axiomatisé de manière finie ;
- 3- il ne doit rien contenir qui ne soit nécessaire à la conjecture d'énoncés d'observation vrais par stimulus, ou d'énoncés conditionnels.»

Ici, la différence entre plusieurs systèmes du monde provient de la différence du langage et du choix de leur formulation, autrement dit toute dialectique se trouve située socio-culturellement. Toutefois, cette pluralité est admise sous couvert qu'en limite les systèmes de connaissances soient tous équivalents.

Ce projet épistémologique se confronte à la question de l'incommensurabilité entre deux théories portant sur un même phénomène, c'est-à-dire que le phénomène se trouve interprété selon deux paradigmes dont la traduction de l'un dans l'autre est impossible. La notion de paradigme se comprend, ici, au sens de Khun (1962) c'est-à-dire un ensemble d'exemple de solutions concrètes dans un cadre problématique défini. Dans ce cas, comment unifier de manière cohérente deux théories incommensurable et satisfaisante pour décrire le monde ?

De plus, cette pluralité de paradigmes se retrouve au sein de chaque individu. Par exemple, Bachelard (1940) explique comment un physicien change sa perception sur le phénomène de masse au cours de son quotidien. Pour Bachelard, cette juxtaposition traduit l'existence de cinq profils épistémologiques adaptés à certain type de problématique, allant de la pratique à l'abstrait.

À chaque niveau, une dialectique s'opère si deux théories dans le cadre du même profil épistémologique se révèlent incommensurables soit il faut monter en abstraction la problématique au profil supérieur, soit la dialectique n'a pas été assez poussée. L'utilisation des profils épistémologiques permet de conserver l'idée d'une dialectique avec le réel pour converger vers un système du monde cohérent sauf que celui-ci possède cinq facettes qu'implique (ontologiquement ?) « la psychologie de l'esprit scientifique », soit une poly-dialectique.

Toutefois, cette façon de concevoir le projet épistémologique oblige à introduire des hypothèses sur la théorie de la connaissance. Conserver un projet épistémologique cohérentiste indépendant de la théorie de la connaissance nécessite d'abstraire la notion de système du monde. Quine selon Putnam (1983) propose un système regroupant l'ensemble des « systèmes du monde » idéaux. Dans ce cas, « il est tentant de caractériser les énoncés émis dans l'une des « formulations théoriques » idéales de Quine comme des *vérités* (en fonction de ce langage et du choix de la formulation [...]) , et en même temps comme constituant toutes les vérités (en fonction toujours des mêmes choix de langage et de formulation) [...] ; seulement, une telle démarche entrerait en conflit avec la *bivalence*, c'est-à-dire avec le principe en vertu duquel *tout* énoncé, dans le langage scientifique idéal que Quine envisage, est vrai ou faux » (Putnam 1983).

En somme, le projet épistémologique cohérentiste vise à concevoir un système de connaissances composé de sous-systèmes qui tend dans sa globalité à rendre compte du monde tout en admettant qu'il possède inévitablement des énoncés ni prouvables ni réfutables sur la base du système. Autrement dit, la position cohérentiste se trouve contrainte soit d'affaiblir le concept de vérité pour tout système de connaissances se traduisant par un système axiomatique fini renfermant les mathématiques classiques soit d'introduire une vérité transcendante dont les fondements restent à définir.

Néanmoins, un tel projet épistémologique suggère que la cognition puisse façonner un ou des systèmes de connaissance suffisamment cohérents qui couvrent même partiellement le monde perçu. La construction de ces systèmes, le processus cognitif, repose sur l'organisation ou plutôt l'auto-organisation des expériences vécues par l'individu en vertu de principes mathématiques supérieurs. En définitive, ce qui différencie les diverses approches de la cognition de la famille FC (Connexionnisme, évolutionnisme, interactionnisme, etc.), c'est la définition des principes organisateurs qui peut se traduire schématiquement par un débat entre internalisme et externalisme.

**

En résumé, en considérant uniquement l'HOC, la notion de vérité n'est pas totalement équivalente entre un énoncé portant sur le monde idéal et le monde physique protéiforme. Un énoncé vrai portant sur le monde idéal signifie qu'il existe dans le monde idéal ce qui implique ontologiquement qu'il soit cohérent avec les principes fondamentaux « évidents » de la logique. L'énoncé portant sur le monde physique protéiforme doit sa vérité à sa cohérence au sein d'un système axiomatique défini dont une partie des axiomes correspondent à une segmentation arbitraire sans prétention ontologique. Le cohérentisme distingue alors une vérité transcendante sous-tendant les principes logico-mathématiques et une vérité mécanisable grâce aux outils façonnés par les premiers.

Mais le cohérentisme se trouve confronté à deux critiques majeures. La première concerne le problème de la mise en correspondance qu'il n'arrive pas à évacuer pour établir les principes logico-mathématiques élémentaires avec l'impasse conceptuelle que cela implique. La seconde critique rappelle que les règles de méthodologie appliquées au travail scientifique s'appliquent également aux mathématiques, c'est-à-dire que les vérités mathématiques doivent être elles aussi garanties en tant que telles en montrant qu'elles sont des théorèmes dans un système. Il apparaît alors une mise en abyme d'une approche systémique de la connaissance qui met à défaut la vérité comme une cohérence, de la même manière que la mise en abyme du critère de la mise en correspondance met en défaut la vérité comme une correspondance.

Les sciences physiques qui utilisent les outils issus des principes logico-mathématiques peuvent se développer malgré ces critiques, contrairement à l'étude de la cognition qui, comprise comme un processus auto-organisé extrayant des organisations, se trouve contrainte de comprendre l'émergence des principes logico-mathématiques qui régissent en même temps le processus cognitif. Par conséquent, si les deux critiques majeures à l'encontre de la vérité comme une cohérence suffisent à montrer que celle-ci mène à une impasse, alors elles impliquent également que les théories cognitives issues de la famille FC se trouvent condamnées à aller vers une impasse.

2.3. La vérité comme une vérification

À l'origine, la notion de vérification correspond à une justification de la valeur de vérité d'un énoncé mais elle peut également devenir un critère de vérité. Cette position se retrouve dans certains courants du pragmatisme (Engel, 1998) toutefois ces liens entre pragmatisme et vérificationnisme seront abordés ultérieurement lors de l'étude sur le pragmatisme. Plus particulièrement, la notion de vérité comme une vérification sera abordée dans le cadre de la FP qui conduit inévitablement à cette définition de la vérité. En effet, la FP considère que le monde physique se compose d'objets indépendamment des croyances de l'observateur (HOP) et qu'en revanche un monde idéal n'existe pas. Dans ce cas, la vérité d'un énoncé ne se comprend pas comme la vérification d'une correspondance transcendante mais comme la vérification d'un énoncé à rendre compte de l'objet du monde par l'expérience. La vérité d'un système de

connaissance correspond à un ensemble d'énoncés cernant au mieux les objets du monde élaborés dans des conditions parfaites, autrement dit un ensemble d'assertions garanties à la limite de l'enquête scientifique.

Comme pour le cohérentisme la correspondance entre une description et le monde physique est déchargée métaphysiquement. Cependant, contrairement au cohérentisme qui doit cette déflation à un conventionnalisme physique, pour le vérificationnisme elle provient d'un conventionnalisme mathématique. La cohérence dans la construction d'énoncés devient alors moins importante que la justification par l'expérience, le principe de la cohérence provient de la structure componentielle du monde physique et non d'une propriété idéelle. Par ailleurs, le vérificationnisme pourra avancer que les mathématiques ne sont qu'une résurgence des lois physiques constituant le substrat de la cognition et que de plus il existe une équivalence entre les définitions des percepts mathématiques et ceux de l'informatique (correspondance Curry-Howard).

Le projet épistémique pour le vérificationnisme correspond à tendre vers un ensemble *unique* d'énoncés décrivant au mieux les objets composant le monde et par conséquent prédisant au mieux l'évolution du monde dans les conditions optimales de l'enquête scientifique. Le cours de l'enquête correspond à une dialectique susceptible de provoquer une révision des croyances. La nécessité de l'unicité de l'asymptote vient du fait qu'elle vise des objets physiques uniques conformément à HOP. Dans ce cadre, la cognition doit posséder les moyens pour mener une enquête sur le monde, c'est-à-dire dégager des assertions garanties. La cognition s'organise grâce à des capacités de détection de régularité et d'anticipations afin d'élaborer des modèles des objets du monde. Contrairement au cohérentisme qui comprenait cette organisation guidée par des principes mathématiques transcendants, le vérificationnisme l'expliquera par les principes physiques ontologiques qui animent son substrat. Les différentes approches en sciences cognitives sous-tendues par la famille FP vont alors soit s'attacher à comprendre la construction d'une représentation du monde à partir de la capacité à prédire et à anticiper en fonction des moyens d'enquête dont dispose l'individu (comme le fonctionnalisme écologique), soit s'attacher à définir le meilleur niveau de description pour dégager les principes organisateurs faisant émerger la cognition (comme l'évolutionnisme ou le neuromimétisme).

Mais le vérificationnisme se heurte à de nombreuses objections, notamment celle de Quine qui attaque la capacité de l'expérience à pouvoir vérifier isolément un énoncé et de ce fait qui discrédite le vérificationnisme. En effet, Quine défend une conception holiste des théories scientifiques c'est-à-dire qu'une théorie ne se confronte jamais seule à l'expérience, mais en bloc, de sorte que la révision de croyances ne peut être isolée. La convergence de proche en proche vers un système du monde unique n'est donc plus assuré puisque deux théories rivales et incompatibles peuvent prédire les mêmes conséquences et que toute théorie moyennant un nombre approprié d'hypothèses *ad hoc* peut être compatible avec l'expérience. Par conséquent, les théories deviennent descriptives et non plus explicatives. Ce point n'entrave pas la production d'énoncé par le physicien vérificationniste mais contrarie juste son idée sur la finalité idéalisée de son travail. En revanche pour le chercheur en sciences cognitives adoptant le vérificationnisme, ce discrédit remet en cause sa démarche scientifique.

Le projet épistémique vérificationniste impliquait que la cognition corresponde ou tende vers l'idéale enquête sur le monde via la vérification. Or le projet épistémique étant discrédité, la légitimité à concevoir la cognition idéale comme un processus construisant une représentation du monde via des modèles anticipatifs n'est également plus tenable sur cette seule base. L'étude indirecte via le substrat se trouve également en difficulté puisque les théories sur le substrat ne deviennent qu'une description locale. Celle-ci ne peut prétendre alors qu'à décrire d'une certaine manière le fonctionnement de la cognition uniquement au sein de ce substrat particulier. En somme, le vérificationnisme induit par la FP mène à une impasse méthodologique pour une compréhension de la cognition en général.

2.4. La vérité comme un consensus

L'appellation de vérité comme un consensus regroupe ici à la fois les trois principales conceptions, non exclusives, du prédicat vrai : (i) comme une redondance, (ii) comme une décitation, (iii) comme une valeur sémantique. Ces conceptions peuvent être avancées sans aucune des hypothèses ontologiques, autrement dit ces conceptions peuvent être utilisées par la FA. La vérité n'a plus de prétention métaphysique puisqu'elle ne vise plus des objets définis ontologiquement. Le « vrai » n'est plus un prédicat ontologique mais une certaine fonction linguistique ou logique, de sorte que la notion de vérité se réduit à un consensus établi au cours des expériences subjectives ou intersubjectives associées à des normes subjectives ou intersubjectives. Les critiques de ces conceptions conduiront à établir les difficultés à définir un projet épistémique et par la suite à définir la cognition qui devrait la sous-tendre. Pour sortir de cette vacuité, le pragmatisme devra apparaître alors comme indispensable.

i - Le prédicat « vrai » comme une redondance

Selon la première conception, « le 'vrai' ne dénote pas une propriété ou une relation *substantielle* qu'auraient nos énoncés, comme la correspondance ou la cohérence, mais un trait superficiel : dire que '*p*' est vrai, c'est simplement *asserter* que *p* » (Engel, 1998). Le prédicat de vérité devient alors redondant. Toutefois, l'équivalence entre '*p*' est vrai et *p* n'est admise seulement si '*p*' est une proposition, soit le schéma d'équivalence (E) suivant (Engel, 1998) :

(E) *La proposition que p est vrai si et seulement si p*

Il faut que celui qui asserte *p* comprenne ce que signifie *p*, autrement dit, il faut que la phrase fasse référence à un vécu ou du moins à un vécu potentiel selon la cosmologie du locuteur.

ii - Le prédicat « vrai » comme une décitation

La deuxième conception reprend l'idée que le prédicat de vérité est superfétatoire mais propose de l'évincer en supprimant simplement les guillemets : « Si 'vrai' n'est qu'un prédicat d'assertion, alors il semble pouvoir s'appliquer à n'importe quelle *phrase* '*p*', que nous en connaissions le sens ou non, et que nous sachions ou non pour quelles raisons elle est assertée » (Engel, 1998). Par exemple, « Les snarks sont boojums' est vrai » devient trivialement « Les snarks sont boojums ». Le *schéma décitationnel* (I) prend alors la forme suivante :

(I) '*p*' est vrai si et seulement si *p*.

Cette position peut sembler délicate puisqu'elle risque d'attribuer le prédicat vrai à une phrase sans savoir ce qu'elle évoque. Cependant, deux lignes de défense se dégagent principalement : d'une part celle qui va considérer que la phrase appartient obligatoirement au langage du locuteur et par conséquent son assertion à un sens au moins pour lui, et d'autre part celle qui va considérer que cette conception n'est valable que pour des propositions hors contexte produites par le raisonnement du locuteur et par conséquent vraies du moins à son point de vue.

iii - Le prédicat « vrai » comme une valeur sémantique

La troisième conception du prédicat de vérité peut être considérée comme une formalisation et une généralisation de la deuxième conception. La définition sémantique de Tarski (1930) souhaite dépasser le problème de la signification du vrai en le réduisant à un résultat logique au sein d'un langage formalisé. Ici, la démarche de Tarski diffère de la démarche cohérentiste dans le sens où pour cette dernière, la formalisation n'est qu'une analyse

conceptuelle intermédiaire pour atteindre la cohérence comme critère ultime de la vérité alors que la définition sémantique élimine les prétentions métaphysiques du prédicat vrai et réduit la définition de la vérité à la résolution analytique. La vérité devient ainsi relative à un langage. De plus selon cette conception, Tarski montre l'impossibilité qu'il existe une vérité pour tous les langages.

Contrairement aux deux précédentes conceptions du prédicat « vrai », le langage naturel est ici considéré mal formé puisqu'il ne permet pas d'identifier les prédicats de vérité et leur système logique associé. Le langage naturel est polymorphe et nécessiterait par conséquent une étude approfondie à part entière. Néanmoins, cette tentative ne permettrait pas de comprendre comment le locuteur reconnaît comme vraie une phrase relative au monde phénoménal.

*

Ces trois conceptions se tiennent sans l'aide d'une seule hypothèse métaphysique, cependant, de nombreuses critiques peuvent être avancées, parfois proches de celles avancées dans les autres conceptions. Par exemple, la vérité-redondance conduit à une décomposition atomique des propositions, or cette position a été disqualifiée précédemment. La vérité-décitation introduit un solipsisme méthodologique dans le sens où l'individu construit son propre système de connaissances à partir de ses expériences et hypothèses *ad hoc*, autrement dit chaque individu devient comparable à une chambre chinoise de Searle avec un sous-ensemble de symboles compatibles mais pas forcément avec le même dictionnaire. Mais la critique majeure de ces deux conceptions est qu'elles assimilent l'assertion à la vérité. Or en faisant cette assimilation, elles oublient les problèmes liés à l'intentionnalité, aux actes de langage, et aux croyances qui ont déjà été évoqués lors de l'étude philosophique des approches en robotique cognitive. Il ne suffit pas de dire avec conviction et sincérité que « la terre est plate » pour qu'elle le soit. Enfin, la définition sémantique de la vérité conduit à une position relativiste en affirmant que toute vérité se trouve relative à un langage mais il apparaît alors une tension intenable lorsque ce principe est appliqué au relativisme lui-même, conduisant à relativiser la vérité du relativisme.

A ce stade, le problème de la position déflationniste du vrai dans le cadre de notre analyse de l'implication du concept de vérité dans la définition de la cognition ne provient pas de ces critiques, mais du fait qu'aucune de ces définitions n'offre un projet épistémologique. Au mieux l'enjeu épistémologique réside dans le fait de se mettre d'accord sur une convention sur les règles de langages et sur l'attribution du prédicat de vérité au sein de celle-ci. La vérité s'identifie en somme à un consensus. Imaginer une auto-organisation sociale reste possible mais cela reviendrait à accepter un réalisme social correspondant à une forme de HOC. Sans projet épistémologique clair, il devient impossible d'en déduire la cognition qui serait capable de réaliser ce projet comme pour les autres définitions de la vérité. En définitive, à partir de cette vacuité et de manière caricaturale, la cognition devient une machine à faire des inférences sans pouvoir expliciter ni leur provenances ni à quoi elles servent, ce qui n'offre aucune piste concrète pour la construction d'un paradigme unique en des sciences cognitives.

La conclusion sur l'étude des quatre types de définition de la vérité se scinde en deux constats. *Le premier constat* concerne les trois premières définitions de la vérité. Chacune de ces trois définitions sous-tendues par une ou des hypothèses métaphysiques propose un projet épistémologique qui circonscrit les théories de la connaissance, autrement dit les définitions de la cognition. Ces positions modifient le cadre et la finalité de l'activité scientifique mais n'introduisent pas d'hypothèse sur les objets d'étude physiques ou mathématiques. Autrement dit, la charge métaphysique de la vérité influence les motivations des scientifiques et leurs choix méthodologiques mais pas l'objet de leur travail qui sera destiné à être évalué, lissé et incorporé par la dynamique sociale de la science. Toutefois, lorsque la cognition devient l'objet d'étude,

L'acceptation de ces définitions de la vérité n'influence plus seulement les motivations mais introduit des hypothèses dans l'objet d'étude, hypothèses qui ne sont pas produites sur la volonté de réaliser un modèle correspondant à un phénomène mais issues d'une réflexion métaphysique. Par conséquent, en adoptant (consciemment ou non) l'un de ces trois types de vérité, la légitimité des approches en sciences cognitives se retrouve dépendante de la légitimité de celle-ci. Or, toutes ces définitions ontologiques conduisent à des apories.

L'impasse conceptuelle de la robotique cognitive s'explique par le fait que le point commun entre quasiment toutes les approches dans ce domaine adoptent au moins une hypothèse ontologique et que celle-ci implique une notion de vérité chargée métaphysiquement, conditionnant la définition de la cognition alors que cette notion de vérité est intenable.

Le second constat porte sur la difficulté à définir la cognition avec une définition de la vérité qui ne s'appuie pas sur une hypothèse ontologique mais simplement sur les propriétés du langage. En effet, il devient difficile de définir la cognition comme le moyen d'accomplir un projet épistémique vide puisque le prédicat de vérité se réduit à une convention relative à un jeu de langage.

Devant cette impasse générale, de nombreux philosophes comme Hume (1748) ou Quine (Putnam, 1983) considèrent alors que seule la psychologie permet de comprendre ce que nous sommes, de « naturaliser la raison ». Cependant, cette démarche n'échappe pas à la circularité entre sujet et objet. Dans ce contexte, Putnam (1983) affirme qu'il est impossible de « naturaliser la raison » parce qu'elle s'appuie sur des normes inextricables. En effet, il rappelle que la raison est à la fois immanente (elle ne peut se trouver hors de l'intuition et d'un jeu de langages corrects) et transcendante (elle régule et critique toutes pratiques et toutes institutions).

Toutefois, cette position se focalise davantage sur l'impossibilité de définir ses propres normes que sur l'impossibilité de définir un processus évoluant avec d'autres normes. En effet, Putnam considère implicitement qu'il n'existe pas d'autres normes cognitives que celles des humains. Le changement de perspective proposé vise ici à l'identification des principes de la dynamique cognitive. Ces principes ne permettraient peut être pas d'en apprendre d'avantage sur les normes humaines inextricables mais permettraient de « naturaliser la dynamique d'une raison » exprimée par une raison humaine.

Mais ce changement de perspective passe par la résolution du problème de la définition de la cognition sans hypothèse ontologique et sans projet épistémologique. Autrement dit, comment motiver la cognition ? A quoi sert la cognition ? Ainsi formulée, les notions d'intérêt, d'utilité et de pratique apparaissent plus fondamentales et incitent à trouver dans le pragmatisme les ressorts nécessaires pour sortir de l'impasse conceptuelle, ainsi qu'à concilier les mots et les choses sans avoir recours à des hypothèses ontologiques. La philosophie première n'est plus une philosophie métaphysique mais une philosophie cognitive.

3. Le pragmatisme comme philosophie cognitive

La position pragmatiste proposée ici s'inspire principalement de celle de James (1907). La présentation de cette philosophie cognitive s'effectuera en deux temps. Dans premier temps, le pragmatisme et l'utilité du vrai seront explicités ainsi que leurs conséquences épistémologiques. Après avoir établi à quoi sert la cognition, ce qu'elle fait, il sera proposé dans un second temps, une histoire de la cognition au travers d'exemples issues de l'éthologie afin de comprendre l'utilité de la cognition selon une perspective évolutionniste. Enfin, les raisons morales qu'implique le paradigme proposé seront présentées et expliquées synthétiquement. En effet, la souplesse de la définition de la vérité pourrait laisser croire à une dérégulation morale et éthique. Cette erreur

amènerait à une vision nihiliste dont les conséquences mortifères sont trop importantes pour être ignorée et point combattues.

3.1. Le pragmatisme et l'utilité du vrai

L'étude précédente sur la vérité rappelle que l'idée de déceler des principes ontologiques régissant la cognition est vaine. Néanmoins, une introspection élémentaire conduit à constater que la cognition repose inextricablement sur un processus de segmentation du monde et de génération de liens entre ces segments. Ce point de départ n'est pas un absolu ontologique mais un absolu pratique ; de même que l'impossibilité de soustraire les normes intrinsèques à nos croyances.

A partir de ces constats, le pragmatisme se concentre sur la manière dont se forment et se sédimentent les croyances. Cette entreprise passe par une définition d'une part (A) des termes de croyance et de norme et d'autre part (B) des capacités fondamentales pour ensuite (C) revenir sur la notion de vérité.

A - Croyances et normes

La définition adoptée du terme croyance sera celle de Pierce (1878) déjà évoquée lors de l'étude du fonctionnalisme : une croyance est un élément participant à la justification d'une action. Les croyances les plus abstraites sont incluses dans cette définition comme le montre James dans son analyse du sentiment religieux (James, 1907). Toute croyance conduit à un comportement. La notion de norme renvoie aux mécanismes *ad hoc* qui sous tendent le déroulement des croyances jusqu'à l'action. Imperceptibles et inextricables à la première personne, ces normes correspondent à des *primitives* (cognitives ou sensorimotrices) sur lesquelles les croyances se développent.

B - Les trois capacités fondamentales de la pensée

À partir d'un constat phénoménologique élémentaire, trois capacités fondamentales liées au déroulement de la pensée peuvent être dégagées :

1. La capacité à évoquer une croyance dans une situation.
2. La capacité à lier ou non deux croyances (créer le lien ou actualiser le lien).
3. La capacité à évaluer une croyance.

Cette triade forme le pendant aux trois conceptions d'une vérité ontologique : l'évocation avec la correspondance, la relation avec la cohérence, l'évaluation avec la vérification. La différence réside dans le fait que ces trois capacités ne reposent pas sur des règles ontologiques mais sur des croyances *et* des normes. En effet, concernant la capacité à lier, la manière de créer des liens ou de réviser des croyances représente uniquement des croyances ou des normes. De même concernant la capacité à évaluer, le critère d'évaluation (désir, besoin, prédiction, etc.) représente uniquement une norme ou une croyance et la crédibilité issue de l'évaluation, une croyance.

Ainsi, toute croyance dépend d'autres croyances dans leurs relations, dans leurs évocations, dans leurs évaluations. Les critères d'évaluations pouvant être multiples et à différent niveau, ils tissent un réseau récurrent de croyances. Cette conception répond au caractère holistique des systèmes de connaissances évoqués précédemment tout en offrant un mécanisme de révision ayant une satisfaction locale. Mais dans le cadre du pragmatisme développé ici, une révision de satisfaction locale des croyances vise une modification soit du désir, soit des moyens, soit des

perceptions ; autrement dit toutes les croyances ayant participé à l'évaluation et le cas échéant du critère d'évaluation lui-même. Le choix de l'endroit où doit se porter la critique repose également sur des croyances ou des normes. En considérant les moyens d'évaluation comme des croyances, le pragmatisme contient le concept de sérépendipité ou d'effet de bord. Les effets de bords correspondent aux implications d'une croyance ou d'une action qui n'ont pas été prises en compte dans la décision de leur déclenchement. Ici, ces implications, si elles sont perceptibles, peuvent générer de nouvelles croyances ou modifier des anciennes.

Mais en définitive, ce sont les critères d'évaluation, à un niveau ou à un autre, qui orientent l'évolution des croyances toujours en fonction du rapport au monde auquel elles conduisent. L'évolution des croyances repose sur une dynamique interne de croyances toujours en contact avec le monde. Cette évolution dynamique associée à la multiplicité des critères d'évaluations conduit à concevoir le réseau de croyances résultant comme un réseau profondément hétérogène ; c'est-à-dire un réseau où des croyances antagonistes ou incommensurables coexistent car leur application reste confinée au sein de perspectives d'usages qui ne se chevauchent jamais. Subjectivement, la cohérence globale est conservée. Dans ce sens, le pragmatisme rejoint le pluralisme ou la rationalité limitée de Simon (1969).

C - La signification du vrai

De la même façon que dans l'analyse précédente de la vérité comme consensus, le prédicat « vrai » se définit par la différence entre deux niveaux d'interprétation. Mais contrairement à la vérité sémantique qui cantonne ces interprétations à des jeux de langage, le pragmatisme les relie à des systèmes d'évaluation de croyances qui renvoient aux motivations et aux compétences de l'individu. Plus simplement, cette différence de niveau d'interprétation peut se comprendre comme la différence entre *savoir-faire* et *savoir*.

Le savoir-faire correspond à la reconnaissance de la solution en fonction du contexte pour un objectif donné. Le savoir-faire correspond à la reconnaissance de la problématique dans laquelle s'inscrivent l'objectif, le contexte et la méthode. Au niveau du savoir-faire, il n'y a pas d'autre explication à la réussite que l'action entreprise ; en ce sens, l'application du prédicat de vérité sur cette action est redondante à la réussite de l'action (qui peut être une assertion). Le savoir-faire provient de l'adaptation des croyances amenant à la réussite. Le savoir provient de la possibilité de résoudre certains conflits entre croyance, expérience et motivation qui ne se résolvent pas avec l'adaptation, soit de la conceptualisation de la situation. Ce besoin de conceptualiser peut provenir de l'identification d'effets de bords qui ne peuvent être pris en compte facilement, ou de la difficulté à réviser des croyances participant à des enchaînements de croyances dont la finalité diffère. La conceptualisation (ou la continuation du processus de segmentation) peut se concevoir comme une manière de réagir au déséquilibre de la dynamique de croyances, provoquée par la nouveauté d'un ensemble d'événements néanmoins liés à cette dynamique. L'enjeu de la conceptualisation consiste alors à identifier les objectifs en les mettant dans la perspective d'autres usages ainsi qu'à identifier les moyens et les contextes afin d'anticiper leurs implications. Au sein de cette représentation ou paradigme évoluant au gré des expériences, la crédibilité d'une croyance devient une croyance jouant le rôle de prédicat de vérité. Mais surtout, la constitution de cette représentation ou paradigme va permettre de problématiser des situations (objectif/contexte), c'est-à-dire ne plus seulement s'adapter au monde mais également adapter le monde, agir de manière proactive.

La dynamique entre *savoir* et *savoir-faire* ressemble beaucoup au constructivisme piagétien (Piaget, 1965) pour qui le système de connaissances correspond à des structures imbriquées se développant et alternant les phases d'assimilation et d'accommodation. En effet, même en considérant le savoir comme la résultante de la théorisation du pratique, la manipulation de cette théorie devient une pratique elle-même théorisable et ainsi de suite, toujours dans la perspective

d'un usage. Par ailleurs, d'un point de vue épistémologique, cette conception de la croissance du savoir rejoint la position de Feyerabend (1975), lorsque celui-ci écrit pour le développement de la connaissance que « tout est bon » puisqu'un comportement peut révéler des effets de bords déstabilisant le système connaissance/croyance l'ayant généré et forçant de ce fait à sa révision. Toutefois, si *a priori* « tout est bon » pour étendre le champ de la connaissance, au final tout ne se vaut pas.

Au delà de la sédimentation des croyances participant au savoir-faire, la notion de vérité attribuée à une croyance se trouve toujours liée à un édifice de croyances façonné par l'expérience selon une analyse critique subjective prenant en compte le coût des révisions. La notion de vérité peut se comprendre alors comme « un nom collectif résumant des processus de vérification, absolument comme "santé, richesse, force" sont des noms désignant d'autres processus relatifs à la vie, d'autres processus qui paient, eux aussi. La vérité est une chose qui se fait... Le vrai consiste simplement dans ce qui est avantageux pour notre pensée, de même que le juste consiste simplement dans ce qui est avantageux pour notre conduite. » (James, 1907).

Toutefois, deux critiques obligent à préciser davantage cette définition de la vérité. La première critique porte sur l'idée que le pragmatisme de James considère qu'une croyance est vraie parce qu'elle fonctionne, alors que d'autres, comme le pragmatisme de Ramsey, considèrent qu'une croyance fonctionne parce qu'elle est vraie. Mais la formulation de cette critique introduit une ambiguïté sur l'interprétation du mot vrai. Pour James, l'usage d'une croyance conduit à lui attribuer le prédicat vrai sans que celui-ci vise au-delà d'un correspondantalisme pratique et non ontologique alors que le pragmatisme de Ramsey présuppose une définition ontologique de la vérité en l'occurrence la vérité comme une vérification. Par ailleurs, avancer que pour James, le prédicat « vrai » s'identifie à ce qui fonctionne serait une erreur d'interprétation sur le terme « avantageux pour la pensée » bien que celle-ci soit toujours en rapport avec le monde.

Cependant, la seconde critique repose également sur un raccourci rapide du pragmatisme de James en comprenant qu'il assimile la vérité à l'utilité. Selon Russell (1910), le pragmatisme déforme le concept de connaissance : savoir que p , c'est savoir que p est vrai ; mais si savoir que p est vrai c'est savoir qu'il est utile de croire que p , ce qui conduit à des absurdités comme par exemple : chercher à savoir si « la neige est blanche » est vrai revient à chercher à savoir s'il est utile de croire en cette proposition. Le problème est que ce raisonnement ne prend pas en compte les différents niveaux de la notion d'utilité ; il existe une différence entre l'utilité de croire que « la neige est blanche » et l'utilité de croire qu'une proposition est vraie. Dans le premier cas, la croyance naît de l'usage, la notion de vérité n'apparaît pas ou au mieux la notion de vérité est redondante. Dans le second cas, la croyance porte sur une autre construisant ainsi une représentation (quelque soit sa nature) qui permet la mise en problématique de situations vécues ou imaginaires conduisant, en fonction de leur utilité subjective, à une prise de décision. Le pragmatisme défendu ici ne réduit pas le vrai à l'utile mais en revanche répond à la question de l'utilité du prédicat de vérité.

Néanmoins, il reste à expliquer pourquoi le pragmatisme ne réduit pas la vérité à la simple assertion. En effet, il ne suffit pas de croire que la terre est plate pour qu'elle le soit. Croire que la terre est plate peut s'avérer avantageux selon certaines cosmologies et répondre à des intérêts psychologiques ou sociaux mais cela reste une opinion fautive. Cette croyance peut être commune à tous les hommes, mais l'intersubjectivité n'empêchera pas qu'elle soit fautive. Dans ce cas, cela signifie-t-il une réintroduction d'une notion d'absolue ontologique, l'idée d'avoir des objets indépendants de nos croyances ?

Non, le pragmatisme évite cet écueil en rappelant le *prima* de l'expérience. Le prédicat de vérité s'applique à une croyance au sein d'un paradigme. La justification d'une croyance peut dépasser le cadre paradigmatique, toutefois tant que cette croyance est maintenue pour les avantages qu'elle procure par ailleurs elle est conservée et demeure au stade de la foi. Par la suite,

la réussite au cours des expériences dans lesquelles elle joue un rôle central augmente sa crédibilité. Cependant, la réussite d'une expérience ne suffit pas pour convaincre qu'une croyance soit une assertion garantie au sein du paradigme, puisque la réussite d'une expérience dépendant toujours du monde indescriptible en totalité empêche d'affirmer la véracité d'un système de connaissances, contrairement à l'échec qui affirme irrémédiablement une erreur. Une vérité au sein d'un paradigme se trouve toujours en sursis tant que toutes les expériences imaginables au sein de celui-ci n'ont pas été réalisées, c'est-à-dire au bout de l'enquête scientifique à supposer qu'elle ait une fin. Mais même dans le cas où toutes les expériences conduisent à considérer un énoncé comme une assertion garantie, elle reste à jamais relative au paradigme qui l'a forgée et qui est dédié à la réalisation de certains usages ; paradigme trouvant sa raison d'être dans une perspective d'usage qui le transcende. Une assertion est garantie seulement si l'usage dans lequel elle s'inscrit est compatible avec le paradigme employé. L'échec de l'expérience ne peut être imputé alors qu'au choix du paradigme employé ou à la façon dont il a été employé.

Dans ce cadre, le pragmatisme contient naturellement l'empirisme ; de plus, la prise en compte de l'asymétrie des conséquences entre la réussite et l'échec le rattache à la démarche épistémologique de Popper, de Lakatos ou de Bachelard, bien que les finalités diffèrent. Cependant le pragmatisme développé ici adhère davantage à la théorie épistémologique de Khun (1962) concernant les structures des révolutions scientifiques. Une communauté scientifique constitue, au gré des expériences, un paradigme ou un ensemble de problématiques qui définissent leurs usages. Une fois constituée, l'évolution du paradigme possède deux phases. La première phase provient de l'ajustement des théories et des usages face aux expériences et à leurs utilités ; les petits déséquilibres sont compensés. La seconde phase correspond à l'abandon du paradigme pour un autre. Cette situation arrive lorsque des usages ou des événements nouveaux apparaissent, rendant intenable la majeure partie du paradigme censé les contenir ; le déséquilibre trop important conduit alors à une révolution scientifique. Néanmoins, si ces paradigmes déchus possédaient en leur cœur des vérités (non mises en défaut au sein du champ expérimental défini par le paradigme), celles-ci perdurent. Autrement dit : « Nous ne disposons pas d'un point d'Archimède ; le langage que nous parlons est toujours celui d'une époque et d'une région ; mais la justesse et la fausseté de ce que nous disons ne vaut pas *simplement* pour une époque et une région. » Putnam(1983).

**

Le pragmatisme, en considérant tous les mécanismes de la pensée comme des normes ou des croyances, propose de rendre compte de l'extraordinaire plasticité de la cognition à travers les différentes orientations critiques offertes à la suite d'un événement insolite, notamment concernant les critères d'évaluation ou les objectifs de l'action. Ce sont ces caractéristiques qui permettent la mise en problématique du monde et du soi compris dans le monde. Les idées permettent des mises en problématique qui vont générer des actions dont les conséquences identifiées seront prises en comptes, pouvant, s'il y a lieu, modifier les idées en formulant ou reformulant des problématiques. Cette construction et la dynamique induite des croyances toujours en interaction avec le monde possèdent des analogies fortes avec la notion d'autopoièse (Varela, 1989) qui renferme à la fois les concepts d'auto-organisation et d'auto-finalisation. Plus précisément, placée au niveau des croyances et de leurs significations, cette autopoièse est comprise comme une autopoièse sémiotique. Cette notion d'autopoièse sémiotique est particulièrement importante car, elle permet, à elle seule, de caractériser toutes les particularités de la cognition et sera de ce fait au cœur du chapitre suivant.

Pour finir, la vérité d'une croyance se trouve relative au paradigme qui la formule ainsi qu'aux champs expérimentaux idéalisés qu'elle évoque au sein de ce paradigme. L'idée de vérité n'étant pas dissociable de la cognition, la prétention de l'ontologie est vaine par conséquent. Il n'y a pas de hors contexte, un point de vue de nulle part. La vérité ontologique n'existe pas. Cette

proposition n'est pas une vérité ontologique, c'est une vérité pratique, une proposition dont l'usage permet de rendre la cognition cohérente à défaut de rendre le monde cohérent. Le pragmatisme développé ici défend radicalement une position agnostique concernant l'ontologie tout en acceptant l'existence d'une notion d'absolu pratique intrinsèque à la cognition issue d'un constat phénoménologique.

Néanmoins, cette présentation sur la cognition et la vérité dans le cadre d'une philosophie pragmatiste suscite deux remarques. La première concerne la possibilité de communiquer entre deux individus ayant construit sa propre vision du monde. Dans le cas de la communication intra espèce, il peut être avancé que la cognition s'appuie sur des croyances et des normes et que ces dernières constituent une base commune de sensations, de besoins, de mises en relation et par suite d'expériences. Par ailleurs, bien que l'individu apprenne seul les contextes l'amenant à apprendre, ce qu'il apprend se trouve façonné par son environnement culturel, sa norme culturelle : « Il n'y a qu'une façon active de connaître, c'est de découvrir et corrélativement, il n'y a qu'une façon d'enseigner, c'est de faire découvrir. » (Bachelard, 1940). Mais de manière plus large, certaines problématiques peuvent être transposées à d'autres animaux ; en effet, il est facile d'imaginer qu'un chien ait soif alors que la transposition est difficile à effectuer sur un dauphin. La seconde remarque insiste sur le fait que seule la pérennité de l'individu viendra conforter ou non les actions, les pensées effectuées. Le plaisir est un moyen de guider nos actions, tout comme l'idée subjective de ce qui est avantageux ; seul l'usage de ces mécanismes a permis leur maintien. La prise en compte de la pérennité de l'individu revient finalement à sortir du système cognitif auto-évaluatif et à prendre une nouvelle perspective pour étudier la cognition.

3.2. Une histoire phylogénétique de la cognition

La proposition d'une histoire de la phylogénèse de la cognition prend sa source essentiellement dans trois essais sur le comportement animal et humain de Lorenz (1937, 1950, 1954), le dernier dégage notamment trois conditions à l'émergence d'espèces douées d'une certaine intelligence, l'homme en particulier. (A) La présentation de ces trois conditions constituera le premier point, permettant ainsi d'introduire le second point qui proposera d'identifier (B) quatre paliers principaux de l'évolution des capacités cognitives en spéculant sur leurs avantages sélectifs.

A - Trois conditions à l'émergence d'espèces douées d'une certaine intelligence

Afin de classifier, d'analyser, puis de comprendre l'enchaînement des comportements, Lorenz propose de construire une phylogénique de la psychologie à partir de la psychologie animale comparée. Les justifications épistémologiques d'une telle proposition seront abordées dans la conclusion sur la présentation du pragmatisme qui éclaircira les rapports existant entre une philosophie cognitive et une phylogénique de la psychologie. Avant de détailler les trois conditions aux comportements intelligents, il faut souligner que la caractéristique commune aux espèces tendant vers ces comportements par rapport aux autres est de présenter moins d'aptitudes spécifiques à un espace vital déterminé et à un mode de vie déterminé. Par exemple, le surmulot n'égale pas le castor à la nage, l'écureuil à l'escalade, le campagnol dans la réalisation de galeries et encore moins la gerboise des steppes à la course, néanmoins le surmulot surpasse chacune des quatre familles apparentées précitées dans les trois aptitudes qui ne sont pas leur « spécialité ». Cela ne signifie pas que le surmulot soit dépourvu d'une spécialisation mais plutôt qu'il s'est spécialisé dans la non-spécialisation. Dans cette perspective, l'homme au vu de ses capacités athlétiques très diverses peut être considéré comme le champion dans la spécialisation de la non-spécialisation. La non-spécialisation élargit le champ des possibles dont l'exploitation se trouve au cœur des trois conditions d'apparition de comportements intelligents.

i - La localisation

La première condition concerne la capacité à localiser des points d'intérêt dans un espace de représentation centrale, quel que soit sa nature. La difficulté à identifier cette capacité réside dans le fait qu'un ensemble de comportements instinctifs peut donner l'illusion de l'existence d'une représentation spatiale centrale. Par exemple, un poisson osseux supérieur percevant une proie derrière une plante aquatique, à travers laquelle il peut voir, mais qui l'empêche de passer. Dans ce cas, le poisson contourne l'obstacle et happe sa proie. La trajectoire résulte à la fois d'une réaction « tigmotactique négative » à la plante et d'une réaction « télotactique positive » à la proie (Lorenz, 1954). De même, des comportements spécialisés peuvent donner l'illusion d'une représentation centrale de l'espace comme le montre le cas de la perdrix qui n'est pas capable de prendre en considération un obstacle perpendiculaire en courant alors qu'elle le peut en volant (Lorenz, 1954). Une manière de discriminer les animaux ayant simplement un comportement réactif et ceux ayant un comportement d'orientation et de scrutation consiste à comparer les conditions de vie. En effet, la richesse comportementale d'une caille huppée de Californie vivant dans la forêt est supérieure à celle de la perdrix des steppes sans différence neuro-anatomique notable (Lorenz, 1954).

La qualité d'une représentation spatiale centrale se trouve corrélée à la complexité de la structure spatiale en fonction des besoins et des moyens d'agir de l'animal. Les habitants des arbres font partie sans nul doute des animaux contraints de maîtriser les structures spatiales les plus complexes ; plus particulièrement encore, les animaux possédant des mains préhensiles contraignant à une certaine analyse avant le saut, contrairement aux animaux possédant des griffes ou des ventouses. Lorenz (1954) établit alors l'existence d'une corrélation entre la nature physiologique de la perception optique de l'espace et la représentation centrale des données spatiales. Toutefois, en s'appuyant sur de nombreux exemples, il situe cette corrélation entre une assez grande précision dans la saisie de l'espace et le fait de fixer les objets du monde environnant plus loin dans la série phylogénétique que la vision binoculaire.

Mais surtout, Lorenz (1954) considère que « l'ensemble de la pensée de l'homme a tiré son origine de ces opérations détachées de l'activité motrice et situées dans l'espace 'représenté', et même que cette fonction originelle constitue la base irremplaçable de nos actes de pensée les plus élevés et les plus complexes » se référant notamment à la présence du vocabulaire spatial dans les rapports abstraits, et ce dans toutes les langues.

ii - La curiosité

La deuxième condition est le comportement de curiosité et d'exploration active que possède les êtres non fixés par une différenciation très spécialisée des organes et des types de comportement instinctif. Cette condition s'inscrit dans la première qui offre le cadre représentationnel pour décomposer le monde en choses ou pour le composer de choses. La curiosité provient d'une sorte de dialogue entre l'objet et l'activité propre de l'animal. Une raison de la valeur conservatoire de l'espèce est que le comportement exploratoire systématique revient à considérer que *tout objet* a potentiellement une importance biologique. Les espèces non spécialisées possèdent l'avantage d'avoir un éventail de comportements plus grand ainsi que des mécanismes de déclenchements moins précis et moins sélectifs. Mais surtout, le processus d'apprentissage exploratoire induit par la curiosité est *indépendant* du besoin de l'instant ou en d'autres termes du motif d'appétence dans lequel le savoir sera utilisé.

Par exemple, même en cas de faim modérée mais reconnaissable, le jeune corbeau va préférer explorer un nouvel objet plutôt que de prendre la friandise présentée. Pour Lorenz, « tout cela signifie, en termes humains : l'animal ne veut pas manger, mais il veut savoir si tel objet précis est mangeable 'théoriquement' ! » Par ailleurs, l'heuristique toujours identique du

corbeau semble conduire à classer les choses selon quatre grandes catégories : ennemi, proie, nourriture, objet sans importance, les objets de cette dernière classe pouvant toujours être recherchés puis utilisés pour cacher de la nourriture. Cette spécialisation de la non-spécialisation des corbeaux leur confère des capacités de survie dans les espaces vitaux les plus divers, du désert à l'Europe centrale. Il est à noter toutefois que la curiosité coïncide avec la jeunesse dans le règne animal, ce qui ne signifie pas que les corbeaux âgés ne puissent plus apprendre mais qu'ils apprennent sous la contrainte de l'environnement. Il n'y a plus d'élan vers l'inconnu.

iii - La néoténie et la diminution des instincts

La troisième et dernière condition avancée par Lorenz (1954) concerne notamment la prolongation de la phase du développement contenant la curiosité. Cette néoténie s'étend jusqu'à la fin de la vieillesse chez l'homme alors qu'elle ne représente qu'une courte phase chez les autres animaux. De plus, la troisième condition contient également la propension à l'indépendance à l'égard d'un instinct rigide, créant ainsi de nouveaux degrés de liberté pour l'action. L'augmentation de cette indépendance favorise par conséquent l'exploration, le dialogue avec le monde conduisant à chosifier le monde puis à conceptualiser (s'apparentant à la chosification) ses sensations et ses actions. D'un point de vue pragmatiste, cette mise entre parenthèses des instincts est comprise comme une libération du sens critique de l'ensemble des facettes de ses activités, qui n'a de sens bien sûr seulement que dans la mesure où potentiellement des réponses à ces critiques existent. En effet, Lorenz (1954) insiste sur l'importance de l'association étroite entre la *praxis* et la *gnosis* pour la direction et l'ajustement des actions contrôlées par le succès. Lorenz (1950) propose enfin d'expliquer l'accentuation de ces caractéristiques chez l'homme par un processus d'auto-domestication qui ne sera pas développé ici afin de se concentrer sur les avantages sélectifs des capacités cognitives et les principes fondamentaux de la cognition, fut-elle rudimentaire.

B - Quatre paliers principaux de l'évolution des capacités cognitives

Partant de l'idée que la cognition est issue d'un processus évolutif, il est intéressant d'entrevoir les facteurs autorisant un tel processus. Au cours de la présentation de l'évolutionnisme dans le premier chapitre, il a été montré que l'évolution ne se trouve pas guidée par une finalité, ce qui n'empêche pas toutefois de concevoir l'évolution des espèces comme un système dynamique qui tend à se coupler avec son environnement via la capacité d'un individu de l'espèce à pouvoir se reproduire.

Dans le cas des comportements réactifs rigides durant la vie des organismes, l'évolution de ces comportements repose entièrement sur l'évolution de l'espèce avec un grain temporel correspondant au cycle de reproduction et un grain spatial correspondant au nombre de situations communes aux individus. En d'autres termes, la sélection naturelle filtre toute variabilité environnementale vécue par l'individu qui n'est pas suffisamment commune par rapport à un seuil critique de la population. Le changement de comportement étant totalement assujéti à la dynamique de l'évolution de l'espèce, les organismes vont converger vers un mode de vie sténoécétique (un espace vital étroit et précis). Les niches écologiques représentent des dynamiques stables mais des changements environnementaux brusques et radicaux à leurs échelles.

Une alternative à cette dynamique de niches apparaît lorsque des capacités d'adaptations viennent répondre à la variabilité environnementale vécue et non critique pour l'espèce. Dès lors, il est raisonnable de penser que les individus pouvant réagir à l'imprévu (même sur une petite marge) augmenteront leurs chances de survie et par conséquent de reproduction. Ce transfert de gestion de l'évolution du comportement de la phylogénétique à l'ontogénèse entrepris, il n'y a aucune raison pour qu'il ne s'amplifie pas puisque la gestion de la variabilité environnementale au

niveau individuel autorise un mode de vie euryoécétique (un espace vital très étendu et varié) qui présente l'avantage de la robustesse. De plus, il faut souligner que ce transfert n'est pas neutre par rapport à la dynamique de l'évolution car, contrairement à cette dernière qui ne possède pas de finalité, l'autopoièse de l'organisme en présente une : le maintien de sa propre dynamique organisationnelle. Ce contexte autorise alors à imaginer une histoire de la cognition décomposée schématiquement en quatre paliers successifs.

i - L'auto-organisation des associations stimuli/comportement

Le stade zéro du comportement correspond à un arc reflexe associant un stimulus à une réponse motrice. Toutefois, Lorentz a montré (1937) que le seuil de déclenchement de certains mouvements instinctifs pouvait varier en fonction de leur utilisation. Par exemple, une période inhabituellement longue de non activation d'un mouvement instinctif conduira à un abaissement du seuil, au point que sa mise œuvre se produise sans excitation extérieure habituellement liée à celui-ci. Mais également, il peut y avoir un glissement d'un comportement vers un autre en fonction de la fréquence d'apparition du stimulus. Le simple schéma excitation-réaction ne suffit plus, il faut introduire un facteur ou un mécanisme endogène. Pour expliquer ce fait, Lorenz emploie la notion d'appétence qu'il apparente à système d'accumulation progressive d'énergie mais il est également possible de proposer un système d'auto-organisation de l'espace des stimuli projetés sur l'espace des actions conduisant à ce que chaque comportement se trouve en compétition avec d'autres. Cette dernière proposition qui sera détaillée dans le chapitre suivant correspond à l'idée d'un couplage élémentaire entre l'organisme et son environnement de sorte à toujours adapter les déclenchements du répertoire des actions sur le champ des possibles offert par les sens, autrement dit l'organisme doit être capable de réagir dans toute situation.

À cette dynamique endogène peut venir se greffer des mécanismes d'inhibition entre mouvement instinctif, par exemple un intrus pourra être considéré comme une proie ou comme un prédateur ; dans le doute, les comportements de fuite ou d'attaque s'inhibent laissant l'occasion à d'autres mouvements instinctifs d'apparaître bien que leur déclenchement puisse être partiel et non efficace. Dans le cas du coq de basse-cour, ce sont les mêmes mouvements effectués quand il picore de la nourriture (Lorentz, 1950). Ce sont ces comportements non finalisés qui au cours de la phylogénèse pourront devenir de véritables signaux pour les congénères.

Par ailleurs, les facteurs endogènes du premier stade cognitif ne se limite pas à l'auto-organisation et aux relations inhibitrices mais englobe également les motivations contextuelles permettant la création de cycles comportementaux, c'est-à-dire différents répertoires d'actions en fonction d'état interne.

Ainsi le premier stade de la cognition, qui peut d'ores et déjà produire des comportements hiérarchiques et cycliques complexes, correspond à la remise en cause des liens donnés entre excitation/réaction par un processus d'auto-organisation, sur une trame donnée, des associations des stimuli, des facteurs endogènes et de réaction. Ce premier stade peut être compris ici comme une sémiose élémentaire ou du moins son prélude et comme l'introduction de la première capacité cognitive avancée précédemment : la capacité évoquée d'une croyance (une association révisable) dans une situation.

ii - Le renforcement

Le second palier correspond essentiellement à la notion de renforcement interne c'est-à-dire à la capacité d'auto-dressage qui a déjà été décrite dans le premier chapitre au cours de l'étude de l'interactionnisme varélien. Cette capacité a pour effet de guider l'association stimulus/action via un mécanisme d'évaluation spécifique et non plus via un processus d'auto-organisation sans

finalité particulière. Ce deuxième palier correspond à l'introduction de la troisième capacité cognitive, la capacité d'évaluer une croyance qui, ici, s'appuie exclusivement sur des normes et dont le résultat est implicite. L'organisme ne réagit plus simplement à un stimulus a priori adéquat mais s'oriente vers un stimulus a posteriori adéquat par rapport à un critère a priori, conférant ainsi à l'espèce une plus grande indépendance vis-à-vis des variations environnementales.

Ces deux stades ne permettent pas de rendre compte de la notion de chose, en revanche ils permettent l'existence de faisceaux de mouvements instinctifs qu'un observateur peut regrouper dans une unité fonctionnelle. De plus, si l'activité d'un organisme ne dépassant pas ces deux premiers stades nécessite une représentation spatiale minimale, une représentation topologique de l'espace (assimilable à un graphe orienté) s'avère suffisante quels que soient les principes dont elle tire son essence. La différence entre les deux premiers stades se situe dans la capacité du deuxième à orienter son apprentissage vers un objectif implicite, contrairement au premier stade où tout apprentissage est latent.

iii - La mise en problématique

Le troisième stade, qui mériterait certainement d'être décomposé en plusieurs, repose essentiellement sur l'explicitation de l'évaluation. Cette explicitation correspond en fait à un besoin, c'est-à-dire à une identification d'une envie, d'un manque contrairement à la notion d'appétence qui se fonde sur des stimuli et à la notion de motivation contextuelle qui se fonde sur des stimuli et des prédispositions à les reconnaître. Le concept de besoin implique la notion de désir. *Dans un premier temps*, le désir va correspondre à la recherche du stimulus désiré dont l'action associée permet d'assouvir le besoin. Tous les stimuli désirables peuvent ne pas être connus a priori par l'organisme, néanmoins ils peuvent être reconnus a posteriori suite à une satisfaction impromptue d'un besoin latent. En d'autres termes, l'organisme peut remarquer un effet de bord lors d'un mouvement instinctif dont le déclenchement ne se trouve pas lié au besoin. Il paraît alors raisonnable que les comportements exploratoires présentent un intérêt conservatoire et que ceux-ci soient favorisés dans la limite des impératifs de la sécurité. La curiosité correspond alors ici à l'application de l'ensemble du répertoire d'actions selon une heuristique sécurisante pour découvrir ces effets de bord. *Dans un second temps*, la complexification des besoins combinés au nombre de stimulus désirables conduit à transformer la somme des stimuli et des actions possibles en une chose localisée dans un espace. La chosification s'appuie ici d'une part sur une représentation métrique de l'espace, quels que soient les principes dont elle tire son essence, et d'autre part sur l'attention active correspondant à la scrutation d'une localisation et non plus passive issue d'une sensibilisation à la présentation de certains stimuli comme cela pouvait être le cas dans les deux premiers stades.

L'avantage sélectif par rapport aux deux premiers stades peut s'illustrer avec l'exemple suivant : Le grèbe huppé évolue au sein d'un environnement vital bien défini de sorte que tous les mouvements instinctifs se trouvent fixés dès le départ (Lorenz, 1937). En effet, la capture de la proie et la nutrition se déclenchent à partir du mouvement du poisson. Cela signifie qu'il est incapable de manger des poissons morts même frais. Il n'est en aucune manière possible de lui apprendre cette solution. Toutefois, il possède certaines capacités d'adaptation qui va l'orienter vers des lieux ou des situations favorables à son activité. Ainsi, cet oiseau répond aux caractéristiques des deux premiers stades. Néanmoins, l'avantage à reconnaître son besoin et à découvrir ce qui pourrait le satisfaire apparaît évident et se distingue du mécanisme de renforcement instinctif qui ne conduit en définitive qu'à une optimisation d'un comportement. Par rapport à tous les exemples évoqués concernant le corbeau, celui-ci peut être considéré comme ayant atteint le troisième stade.

L'explicitation du critère d'évaluation par le besoin conduit à l'explicitation de la satisfaction ou du plaisir qui représente l'évaluation du comportement ou des liens amenant à

celui-ci. Parallèlement, la chosification repose également sur des liens et certainement sur le besoin d'une cohérence entre eux. Ces deux points suggèrent qu'à ce stade, la troisième capacité cognitive identifiée joue un rôle particulièrement important. Ces liens autorisent alors plus que la prédiction qui pouvait faire l'objet d'une évaluation en vue d'une optimisation au cours du deuxième stade, l'anticipation des actions dans la perspective d'un besoin. En d'autres termes, comparés aux stades précédents, la cognition servirait simplement à reconnaître au mieux des solutions a priori et explicites. Quant au troisième stade, il reconnaît un problème a priori et explicite auquel il faut trouver une solution. Le troisième stade par l'explicitation des critères d'évaluations des croyances permet un début de mise en problématique même succincte, et ainsi un approfondissement de l'activité sémiotique, la capacité à donner du sens.

iv - L'imagination

Le quatrième stade pourrait également être davantage décomposé. Toutefois il correspond, une fois les trois capacités cognitives apparues, à la possibilité de compléter, de supplanter les normes rigides qui les sous-tendaient par des croyances. Plus précisément, deux axes évoluent en parallèle. Le premier axe concerne l'augmentation du pouvoir critique sur l'évocation et la mise en relation, afin de formuler ses propres catégories d'objets avec la possibilité de monter en abstraction, contrairement au corbeau qui restera à ses quatre ou cinq catégories innées (proie, nourriture, ennemi, etc.). Par ailleurs, le processus de chosification et d'abstraction ne structure pas seulement le monde extérieur mais également le ressenti provoqué par l'explicitation des évaluations. Le second axe consiste justement à créer et à critiquer de nouvelles évaluations, c'est-à-dire à inventer des motivations dont l'intérêt sera mis à l'épreuve. Dans l'idée de faciliter cette prise de recul critique tout en maximisant la sécurité de l'individu, la capacité d'imagination prend tout son sens, puisque la mise en scène imaginaire, une re-présentation, permet de comparer des scénarii en obtenant ainsi une véritable mise en problématique. En définitive, cette prise de recul conduit à intégrer la notion d'utilité dans la prise de décision.

Pour illustrer ce trait, deux expériences peuvent être comparées (Lorenz, 1950). La première consiste à placer de la nourriture à une hauteur inaccessible pour un chimpanzé, avec toutefois divers objets à sa disposition. Le chimpanzé après une série d'essais et d'erreurs arrive à une solution. La seconde expérience repose sur les mêmes principes mais avec un orang-outan. Ce dernier aura une stratégie toute autre, il trépigne, observe, s'agace et finalement trouve une solution du premier coup. Ces expériences invitent à considérer que l'un pense en même temps qu'il agit et l'autre pense en même temps qu'il imagine, quelle que soit la nature de cette imagination.

La question se pose de savoir ce qui différencie les singes anthropoïdes des hommes pour expliquer la différence entre leurs capacités cognitives. Lorenz (1950) avance que les limites cognitives de ces singes proviennent d'une part du manque de la relation entre « praxis et gnosis » permettant de mener une action avec persévérance constamment réglée par le succès (le sens de l'action) qu'il relie anatomiquement au *gyrus supramarginalis* que seule l'espèce humaine semble posséder ; et d'autre part de la courte phase de curiosité spontanée qui s'efface après l'adolescence alors qu'elle dure quasiment tout au long de la vie d'un humain. Le premier point souligne surtout qu'il existe une différence entre les problématiques pouvant se concevoir comme un emboîtement de problèmes successif et les problématiques dont la résolution des sous-problèmes doit être coordonnée forçant à anticiper l'ensemble ou autrement dit à établir une dialectique entre savoir et savoir-faire. Cette relation se trouve au cœur du développement de la notion d'outil mais également du langage compris comme un moyen d'agir complexe.

Dans le cadre du pragmatisme, il peut être considéré que ce stade est l'aboutissement d'un accroissement de la capacité à se construire des croyances sur les fins, les moyens ou les critiques à tous les niveaux et en même temps de la capacité à les critiquer, les éprouver à travers le monde

créant pleinement une dynamique cognitive au sein de laquelle peut émerger la question « qu'est ce que je fais ? ».

*

En appliquant le même raisonnement sur le début du développement de la capacité cognitive, c'est-à-dire que la complexité environnementale favorise le développement de la cognition, le développement de ce dernier stade a pu être favorisé par la complexité de l'environnement social. Dans ce cas, la rapidité du développement des capacités cognitives devient davantage compréhensible en considérant que la complexité de l'environnement social augmente en même temps que les capacités cognitives des générations successives. En effet, l'individu apprend toujours par lui-même. C'est son entourage qui le guide vers ce qu'il doit apprendre et cela commence bien sûr par l'imitation. Cependant, la pédagogie peut s'enrichir par l'acquisition d'un savoir sur le savoir-faire, procurant de nouveaux moyens pour initier l'apprenti. Ces pratiques, toujours soumises à la critique et à l'évolution de génération en génération, constituent la base culturelle de l'individu cognitif évolué.

Afin de comprendre en quoi l'évolution de la culture, la culturogénèse, peut être comprise comme le pendant de la phylogénèse, il faut revenir sur les mécanismes supposés à l'origine de la vie. Dans le premier chapitre, au cours de l'étude sur l'évolutionnisme, il a été défendu que le vivant prend son origine dans un système organique autopoïétique qui par la suite a intégré dans son ontogénèse un cycle de reproduction amorçant ainsi une phylogénèse. La phylogénèse ajuste la structure dynamique des organismes et le comportement de ceux-ci dépend complètement de cette structure. L'organisme autopoïétique possède sa propre capacité d'évolution irréversible, *hardware*, de même la phylogénèse représente une évolution *hardware* mais au niveau générationnel. Cette évolution *hardware* générationnelle, simplement réactive et inductive, ne sauve pas les organismes présents mais oriente la structure des organismes futurs. L'apparition de couplage structurel sur le plan comportemental et plus uniquement sur le plan physiologique a permis d'introduire une évolution réversible, *software*, individuelle, proactive et par la suite abductive : l'ontogénèse sémiotique. Mais à défaut de se reproduire en même temps que la physiologie de l'individu, l'autopoïèse sémiotique va, en s'exprimant par l'activité de l'individu, conditionner l'environnement cognitif d'autres autopoïèses sémiotiques, transmettre un savoir-faire, amorçant ainsi une culturogénèse, une évolution *software* générationnelle. De part leur co-occurrence, ces quatre types d'évolution s'influencent par un couplage mutuel à divers degrés au point que la culture devient un environnement participant activement au développement des capacités cognitives des individus. Dans cette perspective, la cognition évoluée se construit au sein d'un espace de liberté constitué et défini par des normes physiologiques, des normes comportementales instinctives et des normes culturelles, elles aussi inextricables.

3.3. La valeur morale de ce pragmatisme

En considérant que la cognition évoluée repose sur l'augmentation des degrés de liberté comportementale et corrélativement sur un abaissement des instincts, il apparaît le risque d'éliminer des comportements instinctifs importants pour la conservation de l'espèce. Lorenz (1950, 1954) évoque notamment le réflexe de pitié lors de la confrontation entre congénères qui permet d'éviter la mort. En effet, avec l'introduction des armes (de la pierre à la bombe atomique), le protagoniste dominé risque de mourir bien avant qu'il puisse émettre les signaux d'abdications produisant un sentiment de pitié. Cet instinct court-circuité doit être alors compensé par une responsabilité morale.

Mais sur quelle base élaborer cette responsabilité morale ? La critique de la notion de vérité ontologique a montré que toute hypothèse ontologique conduisait à une impasse, que cette

hypothèse vise le monde physique ou un monde idéal. Ce qui signifie que les philosophies morales portant sur l'idée de Dieu(x) ou autre concept idéalisé comme l'Humanité ne peuvent prétendre à une quelconque légitimité sur le plan rationnel. Par ailleurs, la position impérialiste, qui consisterait à choisir une position dogmatique qui semblerait la moins néfaste à un moment donné et à l'appliquer une bonne fois pour toute, contrevient à l'évolution inévitable de la culture. La dynamique cognitive en exploration constante conduit en effet à reconstruire sans cesse le monde et ses valeurs. Allant à l'encontre de la nature cognitive de l'homme, la rigidité d'un système dogmatique moral conduira alors à moyen ou à long terme à des tensions inéluctables surtout si elle est remise en cause par des événements extérieurs d'envergure.

Face à ce constat, le pragmatisme invite à rechercher des principes moraux qui seraient des absolus pratiques. Dans cette perspective, l'œuvre de Camus (1942, 1953) représente une voie majeure en considérant deux interrogations fondamentales qui ne révèlent pas de la métaphysique mais de la pratique : (i) Pourquoi ne devrions-nous pas nous suicider ? Et (ii) pourquoi ne devrions-nous pas tuer les autres ? Par ailleurs, une troisième question pourrait être rajoutée : (iii) Pourquoi devrions-nous enfanter ou enseigner ? Ici, les réponses avancées doivent être considérées au mieux comme un palliatif intellectuel à un déficit conatif momentané du à l'aseptisation et la facilité de certaines situations offertes par la technique. Ces réponses essaieront d'épuiser le sentiment de dérégulation afin de mieux distinguer les fondamentaux d'une philosophie morale pragmatiste.

i - Pourquoi ne devrions-nous pas nous suicider ?

A priori la mort semble inévitable tôt ou tard avec une fin inévitablement dramatique. Il n'y a aucune lueur d'espoir métaphysique, pas même l'espoir existentialiste à travers le sentiment d'être qui ne vaut pas plus qu'une autre croyance. Alors pourquoi endurer encore plus de souffrance ? Une réponse consiste à rappeler qu'il y a plus de raison à espérer dans la vie incertaine qu'à espérer dans la mort certaine. Cet espoir vient de la capacité cognitive à modifier notre perception du monde, à ajuster ses objectifs à ses moyens et ses moyens à ses objectifs. La métaphysique ne peut répondre au sens de l'existence, cependant la réflexion menée jusqu'ici tend à montrer que le propre de la cognition consiste à problématiser son environnement physique ou mental et à se donner des moyens d'agir sur lui, procurant ainsi une autosatisfaction ; le risque étant d'avoir soit des problèmes impossibles soit une absence de problème. De ce fait, la notion de projet devient cruciale puisqu'elle inclue un retour positif à chacune des trois étapes : conceptualisation, réalisation et contemplation.

ii - Pourquoi ne devrions-nous pas tuer les autres ?

Cette question mériterait comme pour la première un examen plus approfondi, toutefois, une ébauche de réponse peut être avancée. Le pragmatisme définit l'individu cognitif comme une autopoïèse sémiotique, c'est-à-dire comme une dynamique cognitive toujours en interaction avec le monde pour se maintenir. Le développement du psychisme individuel cognitif repose alors sur les interactions sociales. Autrement dit l'autopoïèse sémiotique qui le constitue dépend des interactions avec d'autres individus. De la même manière, plusieurs systèmes autopoïétiques se couplant forment une écologie dont l'équilibre est le garant de la pérennité maximale des entités la constituant ainsi que d'une coévolution viable.

D'un point certain point de vue, tuer tous les autres reviendrait à saborder son équilibre psychique, sa propre autopoïèse sémiotique, autrement dit : soi. Ainsi, l'élimination d'autrui se comprend comme une forme de suicide, ce qui ramène à la première interrogation. Les rapports et les projets avec les autres constituent le soi de sorte que cette dépendance des uns et des autres dans son évolution permanente représente l'essence de la notion de responsabilité toujours partagée de toute activité que l'individu cognitif tisse avec ses semblables.

iii - Pourquoi devrions-nous enfanter ou enseigner ?

De la même manière que la réponse précédente, une piste de réponse pour cette dernière question repose sur la nécessité de créer des liens affectifs avec autrui pour maintenir l'autopoièse du soi. L'enfantement ou transmission au plus jeune permet d'essayer entretenir et de tenter d'améliorer ces rapports jusqu'à la fin, l'amélioration de génération en génération représentant également une forme de projet individuelle et collective.

*

Les trois réponses ne définissent pas en soi un système mais indiquent les fondamentaux que doit suivre une philosophie morale, cette dernière devant toujours être critique face aux raisons et aux conséquences de ses principes en visant une satisfaction globale dans le temps et dans l'espace. Toute la difficulté provient de la diversité des critères de satisfaction et d'insatisfaction dans le fond et dans la forme associés à la dynamique des usages qui doit servir à la critique.

En définitive, tout système entraînant fatalement des effets de bord positif ou négatif, il convient d'agir au mieux lorsqu'un fléau apparaît afin de modérer le plus possible ses effets, de le prévenir ensuite dans la limite du connaissable et de se préparer pour la prochaine calamité. Les philosophies morales défendant ce type de pensée existent depuis longtemps notamment dans l'œuvre de Mill (1863), de Camus (1952) ou plus récemment celle Varela (1996) néanmoins étudier l'apport d'un cadre pragmatisme à ces philosophies dépasse l'ambition de cette thèse.

4. Conclusion

L'étude de la notion de vérité semble suggérer qu'elle ne pouvait reposer sur des hypothèses ontologiques. Cependant, cela ne signifie pas qu'il n'existe pas d'absolu pratique puisque l'expérience cognitive révèle des caractéristiques dont il n'est pas possible de s'abstraire. De même, l'impossibilité de définir ontologiquement la vérité ne signifie pas qu'aucune vérité objective puisse être définie dans le sens où des critères de vérité ne dépendent pas des croyances sur quoi ces critères sont appliqués. Néanmoins, ces critères restent des croyances et leur valeur demeure à la discrétion de ceux qui les appliquent. Dans ce cadre, le pragmatisme considère toute activité mentale comme des croyances évaluables et révisables, y compris celles évaluant et celles révisant, de sorte que la cognition est comprise comme une élaboration dynamique de croyances interagissant toujours avec le monde, soit, de manière synthétique, comme une autopoièse sémiotique.

Paradoxalement la notion de vérité devient annexe dans cette définition de la cognition, néanmoins elle garde une importance décisive dans la compréhension de l'activité cognitive de haut niveau. En effet, encadrée par un paradigme regroupant les outils interprétatifs et modélisateurs de fins et de moyens considérés en une unité pertinente à un moment donné, la notion de vérité permet la problématisation de situations complexes vécues ou imaginées. Dans la perspective de la proposition d'une histoire de la cognition, la cognition représente un moyen de reprendre l'initiative sur le monde, de ne plus le subir mais l'adapter, l'utiliser à ses fins. L'élaboration d'un prédicat de vérité constitue une abstraction supplémentaire dans ce sens. L'individu cognitif est toujours en projet dans le monde physique ou dans l'imaginaire.

Reste maintenant à savoir si la position pragmatiste défendue ici peut prétendre à naturaliser la raison et si oui de quelle manière ? En effet, l'introspection possède intrinsèquement ses limites, la raison étant à la fois immanente et transcendante. Si l'étude de la raison ne peut s'appuyer sur une relation sujet/sujet, deux voies se proposent

immédiatement l'une, matérialiste, considère la relation sujet/substrat, l'autre, psychologue, se concentre sur la relation sujet/autre. Cependant, dans les deux cas, la façon d'aborder cette relation dépend de normes appartenant au sujet d'étude. Lorenz propose de dépasser cette mise en abîme en prenant d'avantage de recul en considérant la raison comme la résultante de la phylogenèse comportementale, de sorte que l'étude successive des points fixes et des transitions puisse révéler le rôle des normes dans la dynamique cognitive. Ainsi, bien que la raison individuelle, ou plus largement la raison humaine, ne puisse pas être naturalisée, cela ne signifie pas qu'il soit impossible de comprendre les principes généraux qui l'ont produite.

En conclusion, le pragmatisme, en définissant une vérité agnostique, revendique le *prima* de la définition des capacités cognitives sur celle du monde, en cela le pragmatisme est une philosophie cognitive bien qu'elle revendique ensuite dans la construction cognitive le *prima* de la pratique. Le pragmatisme étant dynamique, pluraliste et holiste incite toujours à la critique de son activité ce qui autorise toujours la possibilité de refaire ses choix. Par ailleurs, pour tenter d'expliquer les principes de la cognition, le pragmatisme défendu ici invite à prendre la voie de la science ; ce faisant, la philosophie pragmatiste est l'une des rares philosophies proposant une esquisse d'évaluation dans le projet de réaliser une cognition artificielle avec sa propre raison, ineffable.

CHAPITRE III PROPOSITION D'UNE ARCHITECTURE COGNITIVE GÉNÉRALE

« Voilà ce qui est caractéristique. Je ne connais pas de questions et il m'en sort à chaque instant de la bouche. Je crois savoir ce que c'est. C'est pour que le discours ne s'arrête pas, ce discours inutile qui ne m'est pas compté, qui ne me rapproche pas du silence d'une syllabe. »

L'innommable p. 35 de Samuel Beckett (1953).

« Pourquoi le cheval photographié à l'instant où il ne touche pas le sol, en plein mouvement donc, ses jambes presque repliées sous lui, a-t-il l'air de sauter sur place ? Et pourquoi par contre les chevaux de Géricault courent-ils sur la toile dans une posture pourtant qu'aucun cheval au galop n'a jamais prise ? [...] Rodin a ici un mot profond : « C'est l'artiste qui est véridique et c'est la photo qui est menteuse, car, dans la réalité, le temps ne s'arrête pas ». La photographie maintient ouverts les instants que la poussée du temps referme aussitôt, elle détruit le dépassement, l'empiètement, la « métamorphose » du temps, que la peinture rend visibles au contraire, parce que les chevaux ont en eux le « quitter ici, aller là », parce qu'ils ont un pied dans chaque instant. »

L'œil et l'esprit p. 80-81 de Maurice Merleau-Ponty (1964).

1. Introduction

En résumé, le chapitre précédent montre que toute connaissance sur le monde se trouve constitutionnellement relative et subjective, de sorte que l'évolution du savoir devient quantitative à défaut d'être ontologiquement qualitative. Cette restriction provient de l'impossibilité d'éliminer le normatif et ce indépendamment du monde (Putnam, 1983). Bien que sa formulation reste obligatoirement liée à la culture du locuteur, le contenu de cette restriction se trouve absolu puisque tout type de discours se fonde sur une segmentation a priori du monde introduisant le normatif. Autrement dit, cette limitation se révèle intrinsèque à tout système effectuant (quelle que soit la manière) des inférences et des segmentations, soit à tous systèmes cognitifs. Mais cette vérité particulière qui ne porte pas directement sur le monde, demeure confinée au sein d'une réflexion cognitive. Par conséquent, cette restriction représente un absolu pratique et non un absolu ontologique. Toutefois, cela suffit pour rechercher un discours cohérent sur la cognition qui expliquerait l'incohérence de tout discours global sur le monde. Autrement dit, l'étude de la cognition en général est a priori possible ainsi qu'une réflexion sur les conditions de son artificialisation.

Dans cette perspective, l'analyse du pragmatisme dans le chapitre précédent, qui établit les relations entre la philosophie de la vérité et la philosophie cognitive, a permis de définir la cognition comme un système autopoïétique de croyances s'appuyant sur des signes. Seule la notion d'autopoïèse sémiotique explique l'ancrage des symboles qui demeure le problème majeur de toutes les autres approches évoquées dans le premier

chapitre. En d'autres termes, la cognition correspond à la production dynamique des signes et des moyens de leur sémantisation ou de leur évaluation ainsi que leur mise en relation. L'auto-évaluation permet de reconnaître ce qui semble a priori avantageux et d'orienter ainsi la dynamique cognitive dans telle ou telle direction, quitte à redéfinir plus tard le critère désignant ce qui devrait être avantageux. Comme pour tout système autopoïétique (Varela, 1989), seule la pérennité de la dynamique constitue l'évaluation globale d'un système cognitif sans aucune autre contrainte téléologique. Cette définition de la cognition contient alors en elle-même la notion d'autonomie puisque la construction du sens par le système et pour le système correspond à l'auto-finalisation.

Bien que la nécessité de ces propriétés cognitives générales pour tout système cognitif ait été justifiée longuement au cours du chapitre précédent et qu'une généalogie de la cognition ait été proposée, rien ne prouve que ces propriétés avancées soient strictement suffisantes. Seule la tentative d'élaborer artificiellement un système cognitif selon les principes dégagés permettrait de s'assurer de leur validité. Toutefois, il ne s'agit pas de modéliser la cognition d'un individu puisque les états cognitifs demeurent subjectifs (à la première personne) mais de spécifier les conditions de son émergence. Ainsi, l'évaluation des capacités cognitives d'un système artificiel se confrontera aux mêmes difficultés que celles de n'importe quel animal cognitif, soit une évaluation globale et subjective de la part de l'observateur. Par ailleurs, les états informatiques d'un système cognitif artificiel évolué se révéleront certainement aussi complexes et aussi peu explicites que les données neurobiologiques d'un cerveau en activité.

Dans l'objectif de spécifier un système propre à l'émergence de la cognition, la première des trois parties de ce chapitre se consacrera à une analyse d'une autopoïèse sémiotique plus précise que dans le chapitre précédent, en dégageant le concept d'architecture cognitive et celui de schème cognitif. La détermination de l'architecture cognitive se révélera alors comme la première étape dans l'artificialisation de la cognition. Afin de dégager les principes logiques élémentaires qui serviront à formaliser l'architecture cognitive, la deuxième partie comparera les caractéristiques de l'architecture cognitive avec celles des systèmes formels. La ressemblance structurelle entre l'architecture cognitive proposée et les systèmes de classeurs incitera à consacrer la troisième partie à l'étude de ces derniers afin de mieux cerner les convergences et les divergences. À partir de cette étude, la quatrième partie présentera ensuite la spécification d'une architecture cognitive pour la robotique autonome. Dans le cadre de cette spécification, le potentiel et les limites de la conception des schèmes cognitifs seront abordés dans la cinquième partie. Enfin, la conclusion de ce chapitre appliquera la grille d'analyse du premier chapitre à l'approche proposée ici afin de pouvoir la comparer avec les autres approches en cognition artificielle.

2. Analyse de l'autopoïèse sémiotique

Au cours de la présentation du pragmatisme dans le chapitre précédent, un parallèle a été fait entre la production de croyances pour maintenir et développer ce qui semble avantageux avec le concept d'autopoïèse mais appliqué à l'espace des croyances. En considérant que les croyances se fondent sur la production et l'interprétation des signes, ce parallèle a été suffisant pour décrire la cognition comme une autopoïèse sémiotique et pour montrer qu'une théorie de la cognition doit davantage porter sur les propriétés nécessaires au développement de la cognition chez un individu en général que sur la cognition d'un individu en particulier puisque celle-ci dépend de la phylogenèse et de l'ontogenèse.

Cependant, l'artificialisation de la cognition amène à préciser la définition d'autopoïèse sémiotique qui se décomposera en trois sections. La première section examinera l'extension du concept originel d'autopoïèse à celui d'autopoïèse sémiotique. Ce nouvel examen conduira d'une part à définir le concept d'architecture cognitive ainsi que celui de schème cognitif et d'autre part à mieux déterminer les enjeux de l'artificialisation de la cognition. Les deux dernières sections se consacreront alors à l'étude de l'architecture cognitive et des schèmes cognitifs. Plus précisément, la troisième section dégagera les principales caractéristiques de l'architecture cognitive en s'appuyant sur l'étude de la perception et de la conception du premier chapitre ce qui permettra d'aborder la question de l'interface entre l'autopoïèse sémiotique et le corps. La troisième section s'attachera ensuite à étudier la dynamique de l'autopoïèse sémiotique, la sémiose, en s'inspirant de la triade sémiotique proposée par Peirce qui offrira par ailleurs l'occasion de souligner l'originalité de l'approche proposée ici dans son rapport avec la notion d'objet. La mise en perspective de cette étude avec la généalogie de la cognition proposée dans le chapitre précédent permettra alors d'établir plus précisément les propriétés des schèmes cognitifs selon le stade de l'évolution de la cognition. A la fin de cette partie, la définition des concepts participants à la description de la cognition sera suffisante pour entreprendre la formalisation d'un système cognitif.

2.1. Autopoïèse sémiotique et autopoïèse physique

Avant d'aborder les spécificités d'une autopoïèse sémiotique, (A) la définition de l'autopoïèse de Maturana et de Varela (1989) sera rappelée. Ce rappel insistera sur le fait que la notion originale d'autopoïèse de Maturana et de Varela (1989) s'applique uniquement à des systèmes physiques. Toutefois, cette section montrera qu'il est possible de l'étendre au domaine des signes en introduisant le concept d'architecture cognitive et celui de schème cognitif. Cependant, (B) la différence conceptuelle de la cognition, selon que celle-ci se fonde sur la notion d'autopoïèse physique ou sur la notion d'autopoïèse sémiotique, conduira à compléter la critique du constructivisme varelien effectuée à la fin du premier chapitre. Cette dernière critique permettra d'apporter un argument supplémentaire à la faisabilité d'une cognition artificielle. Mais surtout, à partir de cette analyse générale de l'autopoïèse sémiotique, (C) une stratégie pour le développement d'une cognition artificielle sera avancée ainsi que ses difficultés.

A - Rappel de la notion d'autopoïèse suivi de son extension au domaine de la sémiotique

Comme il a été évoqué dans le premier chapitre, Maturana et Varela (1989) considèrent qu'« *un système autopoïétique est organisé comme un réseau de processus de production de composants qui (a) régénèrent continuellement par leurs transformations et leurs interactions le réseau qui les a produits, et qui (b) constituent le système en tant qu'unité concrète dans l'espace où il existe, en spécifiant le domaine topologique où il se réalise comme réseau* ». Les constituants d'un système autopoïétique étant continuellement remplacés, l'invariant du système réside dans son organisation. Face aux perpétuelles perturbations environnementales, la dynamique de régénération s'ajuste de sorte à les compenser. Ce point de compensation correspond à un point d'équilibre dynamique (ou un point de fonctionnement) de l'organisation du système. Dans le cas où les perturbations externes ou internes sont telles qu'aucun point de fonctionnement au sein de l'organisation ne permette de les compenser, deux types d'évolution existent pour une autopoïèse : soit la dégénération de l'organisation conduit à une nouvelle organisation offrant une dynamique adaptée, soit la dégénération de l'organisation conduit inévitablement à la dislocation de tous les constituants mettant un terme à l'autopoïèse.

Ainsi, le couplage ponctuel correspond à l'ajustement du point de fonctionnement du système dans le cadre de son organisation et le couplage structurel correspond à l'ajustement de l'organisation elle-même. En d'autres termes, le couplage ponctuel représente l'ensemble des comportements du système dont les transitions sont réversibles dans le cadre de l'organisation définie et le couplage structurel représente les modifications structurelles amenant de manière irréversible à une nouvelle organisation. En définitive, le propre d'un système autopoïétique est de pouvoir entretenir une dynamique interne (homéostasie) quitte à modifier son organisation selon la contingence du monde (ontogenèse).

Par analogie, une autopoïèse sémiotique correspond à un réseau de processus de production de croyances et de signes qui régénèrent continuellement par leurs ajustements et leurs relations le réseau qui les a produit. Cette régénération compense les perturbations extérieures que sont les expériences significatives, entraînant par exemple la vérification des croyances, la reconnaissance de signes ou l'intégration de nouvelles informations. Dans ce cas, le couplage ponctuel correspond alors à l'utilisation de la croyance la plus cohérente dans un contexte donné par rapport à l'organisation du réseau de croyances, c'est-à-dire à l'interprétation, et le couplage structurel correspond à la révision des croyances et de leurs relations, ce qui modifie l'organisation du réseau de croyances. Cette révision s'opère par apprentissage. Ces analogies seront précisées au cours de l'analyse de l'autopoïèse sémiotique en particulier celle concernant par exemple les facteurs de dégénération qui poussent tout système autopoïétique à s'engendrer perpétuellement. Cependant, toutes ces analogies ne concernent que le point (a) de la définition de Maturana et de Varela (1989).

Le point (b) stipule qu'un système autopoïétique est une unité concrète. La concrétude se comprend dans le cadre de la définition de Maturana et de Varela comme la réalisation effective d'une structure dont la disposition dans l'espace des constituants définit les relations topologiques ainsi que le réseau d'interactions qui les unit. Ce point implique a priori que tout système autopoïétique se comprend dans le cadre d'une physique, autrement dit, que tous les systèmes autopoïétiques représentent une autopoïèse physique. L'exemple de l'automate de tessellation (Bourgine et Stewart, 2004) décrit lors de l'étude sur l'évolutionnisme dans le premier chapitre et illustrant une autopoïèse minimale demeure une autopoïèse physique bien que celle-ci ne soit réalisable uniquement dans un monde virtuel simulé sur ordinateur. En effet, cet automate est conçu dans un espace géométrique avec des éléments définis ainsi que leurs réactions. La construction virtuelle de la structure de ce système autopoïétique passe par une modélisation de l'univers virtuel. Les éléments du monde virtuel affectent concrètement la structure virtuelle de l'automate. Ainsi, l'autopoïèse se comprend encore dans le cadre d'une physique bien qu'elle soit virtuelle. L'importance de la physique dans la définition de l'autopoïèse sera au cœur du complément de la critique du constructivisme, mais avant de l'aborder, la validité conceptuelle de la notion d'autopoïèse sémiotique dans le cadre d'une philosophie pragmatiste doit être démontrée.

Dans le cadre d'un ensemble de croyances et de signes, la démonstration de la validité conceptuelle de l'autopoïèse sémiotique repose sur l'interprétation des deux notions constituant le point (b) : la notion d'unité et celle de concrétude. Ne pas considérer que la notion d'unité soit applicable à un ensemble de croyances et de signes, signifierait que cet ensemble puisse se comprendre comme une juxtaposition de croyances et de signes qui seraient indépendants les uns des autres. Or, à plusieurs reprises dans les deux premiers chapitres, il a été évoqué que d'une part les croyances ne peuvent être testées indépendamment les unes des autres (Quine, 1951) et d'autre part, les concepts descriptifs

forment un réseau hiérarchisé de dissemblances et de ressemblances, de sorte qu'ils ne se comprennent qu'au sein de ce réseau (Soler, 2000). Par ailleurs, il sera montré plus précisément que cette interaction entre les concepts descriptifs se retrouve naturellement avec la notion de signe. Ainsi, un ensemble de croyances et de signes forme un réseau dont l'unité est garantie par l'imbrication de leurs relations.

Afin de mieux saisir ce qu'implique la notion de concrétude, il faut rappeler, au sein de la méthode artefactuelle (MA), la distinction entre la simulation logique qui déroule les calculs d'une description mathématique (DM) ou d'une description procédurale (DP) et la simulation effective qui laisse interagir des éléments imitant ceux décrits par une description componentielle (DC). Par exemple, pour l'automate de tessellation (Bourguine et Stewart, 2004) qui se conçoit dans un univers simple et complètement défini, la simulation logique correspond au déroulement des équations aux dérivées partielles spatialisées modélisant l'organisation de ce système autopoïétique. La simulation effective correspond à la réalisation de la structure de l'automate de tessellation dans le cadre de la modélisation d'un univers virtuel.

Selon le critère de concrétude, la première simulation n'est pas une autopoïèse physique contrairement à la seconde. En effet, la simulation logique rend uniquement compte du couplage ponctuel du système dans le cadre d'une organisation définie. Par exemple, l'introduction d'un nouveau composé avec les réactions associées oblige à imaginer les implications de celui-ci puis à concevoir un nouvel ensemble d'équations sensé décrire l'organisation résultante. En revanche, dans le cas de la simulation effective, cette introduction se résume à la modélisation de nouveaux éléments virtuels rajoutés à ceux déjà définis et ce sont les nouvelles interactions générées qui vont éprouver directement la structure de l'automate de tessellation réalisé dans ce monde virtuel. Ainsi, contrairement à la simulation logique, la simulation effective peut rendre compte du couplage structurel d'un système autopoïétique. En définitive, la concrétude se traduit par la mise à l'épreuve perpétuelle de la structure d'un système autopoïétique face à la contingence du monde qu'une DP ou une DM ne peut rendre compte. Cette définition n'est alors plus strictement associée à la physique. Il reste à présent à transposer cette capacité « à éprouver » à un réseau de croyances et de signes.

Dans le cadre de la philosophie pragmatisme développée au cours du chapitre précédent, la croyance a été définie comme un élément participant à la justification d'une action ou d'une autre croyance de sorte que c'est l'usage qui les forme et les développe. La mise à l'épreuve des croyances correspond aux conséquences de leur application effective pour le réseau de croyances et pour le système physique sur lequel il repose. Par ailleurs, ces conséquences peuvent être évaluées en fonction des attentes justifiées par des croyances portant sur ce qui est avantageux. Autrement dit, le réseau de croyances s'éprouve lui-même. Continuellement sollicité et perturbé par les expériences signifiantes traduisant les expériences physiques du monde, le réseau de croyances se comprend alors bien comme un système autopoïétique.

L'autopoïèse sémiotique peut être décrite par les relations logiques de la même manière que l'automate de tessellation par les équations aux dérivées partielles, mais cette description ne constitue pas une autopoïèse. Un système autopoïétique sémiotique se constitue de relations effectives, ce qui ne signifie pas qu'une autopoïèse sémiotique ne puisse pas utiliser des principes logiques mais que ceux-ci, s'ils sont employés, résultent d'un couplage structurel considéré comme avantageux et se traduisent aux travers de relations effectives. La formalisation de l'autopoïèse sémiotique se consacrera à la

description logique des relations alors que la spécification se concentrera sur la manière de rendre ces relations effectives.

Mais avant d'aborder la formalisation et la spécification qui seront l'objet des deux parties suivantes de ce chapitre, le concept d'architecture cognitive et celui de schème cognitif se dégagent de cette analyse de l'autopoièse sémiotique. En effet, la capacité à éprouver l'effectivité des relations et la traduction des expériences physiques en expériences signifiantes suggèrent l'existence de mécanismes qui sous-tendent la dynamique sémiotique. L'ensemble de ces mécanismes représente les mécanismes subcognitifs (ou parasémiotiques) formant l'architecture cognitive. En d'autres termes, cette dernière répond aux propriétés nécessaires à l'émergence d'une autopoièse sémiotique. La structure de cette autopoièse correspond alors à l'ensemble des schèmes cognitifs. Par analogie avec l'autopoièse physique, l'architecture cognitive coïncide avec la définition des composants du monde et de leurs interactions dans lequel se réalise une autopoièse physique, et les schèmes cognitifs coïncident avec la topologie des constituants de l'autopoièse physique.

B - Complément à la critique du constructivisme varelien

La distinction entre l'autopoièse sémiotique et physique permet de mieux différencier l'approche pragmatiste de la cognition de celle du constructivisme varelien. En effet, les deux approches, en plus de rejeter HOP et HOI, s'accordent sur le fait que le sujet se façonne en fonction des interactions qu'il entretient avec son environnement et que cette manière de se construire suggère que le sujet possède une dynamique propre amenant sans cesse à l'interaction. La différence porte sur les moyens conceptuels pour comprendre les propriétés cognitives de cette dynamique interne. Pour le constructivisme varelien, il a été vu dans le premier chapitre qu'un ensemble de systèmes autopoiétiques peut former une clôture opérationnelle de telle sorte que cet ensemble puisse être considéré lui-même comme une autopoièse. Dans ce cadre, la cognition se comprend comme les propriétés résultantes de la clôture opérationnelle du système nerveux. Ainsi, l'étude de la cognition passe par la compréhension de l'organisation des autopoièses du substrat. Sans revenir sur les conclusions du premier chapitre, cette approche se distingue du neuromimétisme dans la mesure où l'étude de l'organisation prévaut à celle de la structure et se distingue du connexionnisme du fait que les principes de l'organisation ne reposent que sur sa subsistance et non sur des principes mathématiques.

Mais après avoir montré que le concept d'autopoièse originel sous-entend l'existence d'une physique, le constructivisme varelien, en considérant que l'étude de la cognition passe par le concept d'autopoièse originel, se trouve face à la contradiction suivante : HOI et HOP rejetées, la valeur de la physique sous-tendue n'a de sens qu'au travers de la philosophie cognitive devant justement être définie à l'aide du concept d'autopoièse. En définitive, ce paradoxe explique que la distinction avec le neuromimétisme ou le connexionnisme ne soit que de principe lorsque le constructivisme varelien doit concrètement concevoir des modèles cognitifs.

En considérant que le terme « concret » fait référence uniquement à la capacité qu'à un système à être éprouvé de part les relations effectives de ses constituants, la notion de physique s'évanouit, permettant à l'approche pragmatiste de proposer le concept d'autopoièse sémiotique. Mais surtout, cette nouvelle interprétation du terme « concret » entraîne le fait que la cognition ne dérive pas nécessairement d'une autopoièse physique, autrement dit, qu'il devient possible d'imaginer une cognition artificielle dont le support puisse être un robot traditionnel et non un robot dont la structure matérielle résulterait d'une autopoièse physique.

C - Stratégie pour le développement d'une cognition artificielle

Les précisions apportées au concept d'autopoïèse sémiotique suggèrent quatre remarques pour l'artificialisation de la cognition. La première remarque indique que la première étape dans l'artificialisation de la cognition passe par la définition de l'architecture cognitive et que de cette définition résultera le formalisme dans lequel doivent s'exprimer les schèmes cognitifs. La deuxième remarque concerne l'origine de l'autopoïèse sémiotique. En effet, un système ne devient pas progressivement autopoïétique, c'est la réunion simultanée des éléments qui forme une autopoïèse minimale à partir de laquelle le système se développe. Ainsi, la cognition n'apparaît pas progressivement dans une l'architecture cognitive vide de schèmes cognitifs. Autrement dit, la cognition ne se construit pas à partir d'une *tabula rasa* comme le prône l'empirisme radical (Locke, 1689), elle se développe à partir d'une autopoïèse minimale. La deuxième étape de l'artificialisation de la cognition consistera alors à concevoir une autopoïèse sémiotique minimale. La troisième remarque souligne que la difficulté de cette entreprise ne provient pas tant de la conception d'une architecture cognitive et des schèmes cognitifs permettant une autopoïèse minimale que de faire en sorte que cette dernière soit suffisamment riche pour offrir un couplage structurel correspondant à l'apprentissage. En effet, le couplage structurel est toujours lié aux propriétés structurelles du système et aux perturbations environnementales. La richesse d'un système autopoïétique se comprend comme l'ensemble des évolutions structurelles possibles en fonction de la contingence environnementale.

En reprenant l'exemple de la réalisation de l'automate de tessélation dans un monde virtuel, la seule évolution structurelle envisageable de celui-ci correspond à la dégénération progressive menant à la dislocation. L'apparition éventuelle de couplages structurels successifs aboutissant à un système plus complexe voire multicellulaire nécessite d'introduire une nouvelle chimie suffisamment riche et compatible avec l'organisme. Mais la confection d'un tel univers virtuel, de part la complexité nécessaire et de part l'explosion combinatoire qu'elle entraîne, demeure utopique.

Toutefois, contrairement aux autopoïèses physiques dont les constituants proviennent des éléments déjà présents dans l'environnement, l'autopoïèse sémiotique doit générer perpétuellement de nouvelles croyances qui seront intégrées et éprouvées. La génération de croyances provient soit de l'architecture cognitive soit des schèmes cognitifs. L'autopoïèse sémiotique se comprend alors comme un système dont la dynamique interne repose sur l'autostimulation. La difficulté de l'artificialisation de la cognition ne correspond pas à la définition d'un ensemble de croyances préalables mais à la définition des modalités conduisant à la génération de croyances qui autorisent un couplage structurel. Par ailleurs, le terme croyance est défini ici au plus bas niveau comme une relation entre un signe et une action mentale ou physique révisable. Des pistes destinées à surmonter cette difficulté seront avancées dans la section suivante grâce à une mise en perspective avec la généalogie de la cognition.

2.2. L'architecture cognitive

Dans le cadre de l'artificialisation de la cognition, la définition de l'architecture cognitive dont émerge la sémiose demeure prioritaire. Toutefois, avant d'étudier cette sémiose qui sera l'objet de la section suivante, l'architecture cognitive doit répondre à certaines conditions générales relatives à un individu cognitif qui (A) définissent les caractéristiques structurelles de celle-ci, ainsi que (B) son rapport avec le corps qui la sous-tend. Le premier point s'appuiera sur les discussions concernant la perception et la

conception rencontrées essentiellement lors de l'élaboration de la grille d'analyse au début du premier chapitre. Le second point insistera sur la nécessité de l'autopoïèse sémiotique à être rattachée à un corps et établira les caractéristiques de l'interface entre celui-ci et l'architecture cognitive.

A - Les caractéristiques structurelles d'une architecture cognitive

Les caractéristiques structurelles de l'architecture cognitive doivent offrir les conditions élémentaires au développement de la cognition. Une manière d'identifier ces conditions consiste à déterminer les conditions élémentaires à la perception et à la conception qui sous-tendent la cognition. Or, dans le premier chapitre au cours de l'élaboration des hypothèses cognitives, il a été établi que la reconnaissance d'une illusion ou d'une hallucination s'appuie obligatoirement sur des capacités mnémoniques et des capacités d'exploration (physiques ou mentales). Les caractéristiques structurelles d'une architecture cognitive reposent alors sur les capacités mnémoniques et sur leurs manipulations élémentaires. La présentation de ces capacités s'effectuera en deux temps.

Dans un premier temps, la désambiguïsation des jugements (liés à la perception ou à la conception) par la prise en compte de situations vécues passées suggère déjà l'existence de deux mémoires : d'une part une mémoire événementielle contenant des situations passées sur un certain empan sous forme de signes et d'autre part une mémoire épistémique contenant l'ensemble des croyances qui conduiront selon l'interprétation à des jugements. La reconnaissance dans la mémoire événementielle des éléments nécessaires pour effectuer des jugements par la mémoire épistémique correspond alors à une manipulation mnémonique élémentaire.

Dans un second temps, l'analyse de la reconnaissance d'une illusion ou d'une hallucination suggère que le fonctionnement de la mémoire événementielle n'offre pas seulement la simple confrontation d'une situation passée mémorisée avec la situation présente mais que la mémoire événementielle doit offrir également la confrontation de la prédiction effectuée à partir d'informations mémorisées et des actions entreprises avec la situation présente. La mémoire événementielle incorpore les informations liées aux actions entreprises et aux prédictions qui proviennent de jugements prédictifs effectués à partir de la mémoire épistémique. Ici, la prise en compte du temps dans la reconnaissance d'une prédiction conduit à considérer une nouvelle manipulation mnémonique élémentaire. La révélation d'une hallucination ou d'une illusion survient lorsque la prédiction sur un jugement (perceptif, existentiel, de vérité ou d'assertabilité) se révèle erronée. La mémoire événementielle possède alors une structure duale de la gestion de l'écoulement du temps, d'un côté les signes dont la distance avec le présent augmente, et de l'autre les signes dont la distance avec le présent diminue (Figure III-1).

Ainsi, se distinguent deux catégories de signes au sein de la mémoire événementielle : les signes représentant un constat et ceux représentant une intention ou une prévision. Cette distinction offre l'opportunité de préciser quatre points de vocabulaire :

- Une **intention** (ou une prévision) se comprend ici comme la programmation d'effectuer quelque chose (mentalement ou physiquement).
- Une **prédiction** correspond à l'attente d'une certaine observation de manière plus ou moins précise d'un constat ou d'un jugement.

- Une **anticipation** correspond à la simulation interne d'intentions et d'actes à partir de prédictions qui conduisent au final à des intentions et à des actes effectifs.
- Une **prédisposition** correspond aux constats ou aux jugements préalables ou facilitateurs pour en effectuer de nouveaux.

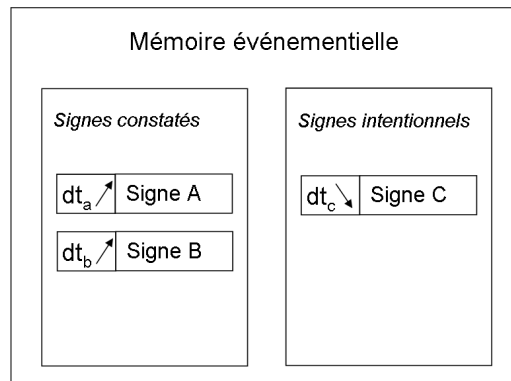


Figure III-1 : Schéma de la mémoire événementielle avec les deux catégories de signes : les signes constatés dont l'écart temporel avec le présent augmente et les signes intentionnels dont l'écart temporel avec le présent diminue

En définitive, les approches précédentes ainsi que les systèmes de classeurs qui seront présentés nomment anticipation ce qui est appelé ici prédiction. L'anticipation devient une notion complexe qui repose sur des schèmes cognitifs lui préexistant, comme la prédiction.

En somme, une architecture cognitive repose au minimum sur deux types de mémoires : la mémoire événementielle et la mémoire épistémique, auxquelles s'associent deux types d'opérateurs : la reconnaissance spontanée et la reconnaissance attendue. En fonction de la mémoire événementielle, la mémoire épistémique émet un jugement dans la mémoire événementielle qui pourra alors participer à la réalisation d'un nouveau jugement. Ainsi, l'interaction entre les deux mémoires offre une réflexivité interne minimale (Figure III-2).

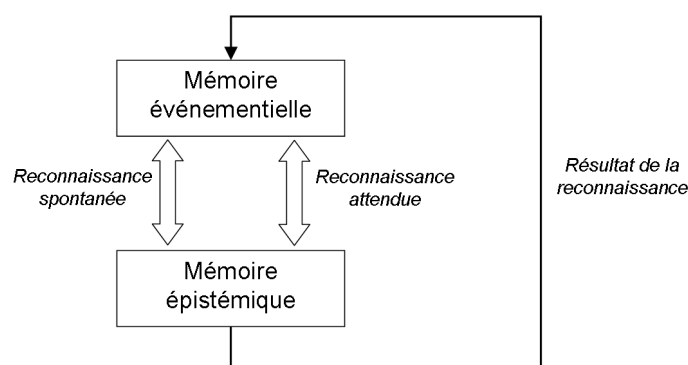


Figure III-2 : Schéma des éléments et des opérateurs de base d'une architecture cognitive

La décomposition en différents types de mémoire (mémoire à court terme et mémoire à long terme ou mémoire sensorielle et mémoire procédurale) a été la base de nombreuses architectures cognitives comme ACT-R (Anderson, 1974) ou SOAR (Newell, 1980). Mais celles-ci, traditionnellement symboliques ou subsymboliques, ne manipulent que des symboles rattachés a priori à des concepts ou percepts figés (ou semi-figés si une

optimisation existe) et définis indépendamment des autres symboles, alors que l'architecture cognitive proposée repose sur la manipulation de signes dépourvus de tout concept et perçut a priori puisque la sémantique doit dépendre de l'auto-organisation des signes en perpétuel ajustement au cours du vécu. Ces architectures cognitives tendent vers l'imitation ou la simulation de la cognition mais non vers son artificialisation.

Par rapport aux concepts de probabilité abordés lors de la présentation du fonctionnalisme écologique dans le premier chapitre, la mémoire événementielle et la mémoire épistémique renvoient respectivement à la notion de probabilité événementielle et à celle de probabilité épistémique. La probabilité événementielle porte sur l'apparition d'un événement dans la mémoire événementielle et la probabilité épistémique porte sur la pertinence globale d'un jugement suivant l'évaluation subjective de ses conséquences identifiées comme telles. Par ailleurs, comme il a été souligné dans le premier chapitre au cours de la présentation du fonctionnalisme, le rôle fondamental de l'opérateur de reconnaissance attendue conforte l'importance de la notion de prédiction, qui conditionne celle d'anticipation dans l'élaboration d'un système cognitif. Mais, en plus d'un opérateur dédié à la prédiction, cette capacité implique que la mémoire événementielle doit être en mesure de gérer l'écoulement du temps pour les signes passés et pour des signes attendus qui peuvent être liés à la prédiction ou au séquençage d'une activité, c'est-à-dire de programmer la transition d'un signe considéré comme attendu (dont la distance avec le présent diminue) à un signe réalisé (dont la distance avec le présent augmente). Ainsi, la mémoire événementielle gère en même temps les événements passés, présents et futurs.

B - L'interface entre l'architecture cognitive et le corps

L'autopoïèse sémiotique ne peut se développer dans un monde clos puisque c'est l'interaction avec le monde qui permet d'éprouver les croyances. Autrement dit, l'interface entre l'architecture cognitive et le corps ne relève pas de la contingence mais de la nécessité. L'architecture cognitive offre un espace dans lequel les croyances s'organisent en fonction des signes apparaissant dans la mémoire événementielle. Ces derniers représentent des événements mentaux qui proviennent de trois sources (Figure III-3) : la sensation, la proprioception et l'activité de la mémoire épistémique réagissant au contenu de la mémoire événementielle.

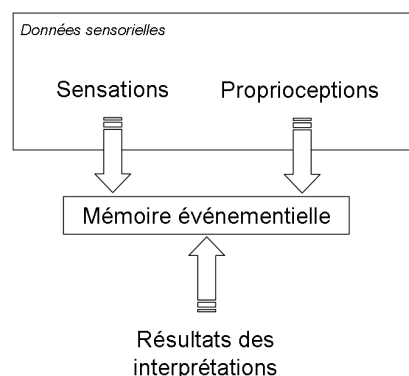


Figure III-3 : Les trois sources alimentant la mémoire événementielle.

Les données sensorielles qui regroupent sensation et proprioception alimentent la mémoire événementielle indépendamment de celle-ci. Les signes sensoriels représentent les moins abstraits de la sémiotique puisqu'ils ne résultent pas d'une interprétation de la sémiotique, soit d'une inférence. L'ensemble des processus participant à la production de ces signes

sensoriels constitue les primitives sensorielles. Concernant la capacité d'agir, l'interprétation d'une croyance conduisant à une action peut prendre deux formes : soit celle d'une action directement appliquée, soit celle d'un signe intentionnel qui explicite et programme l'action à effectuer. Les primitives motrices correspondent à tous les processus situés après l'envoi d'une commande motrice soit directement par l'interprétation, soit indirectement par un signe. En définissant ainsi les primitives (Figure III-4), l'interface entre la sémiose et le corps concerne uniquement les activités corporelles qui synthétisent le signe et l'effet.

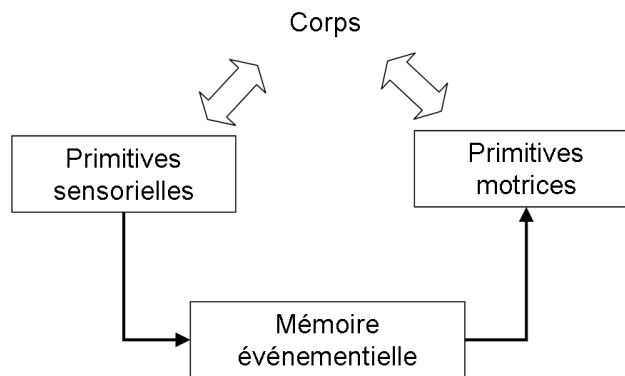


Figure III-4 : Interface entre la mémoire événementielle et le corps.

Quelle que soit la complexité de ces primitives, il suffit qu'elles soient régulières dans leur fonctionnement pour que le système sémiotique puisse façonner des règles cohérentes. Par exemple, une primitive sensorielle peut inclure des réflexes élémentaires comme la contraction de l'iris pour diminuer la luminosité. La définition du réflexe devient alors une relation sensorimotrice se trouvant en dehors du système sémiotique. Même si des signes décrivent le réflexe, ils ne se l'approprient pas : une image produit un signe et parallèlement la luminosité, qui peut transparaître dans le signe, contracte l'iris produisant un signe proprioceptif. Ce scénario offre la possibilité d'observer ses propres réflexes et d'en tenir compte sans toutefois pouvoir les modifier. De même, une primitive motrice peut également se révéler complexe en se basant par exemple sur des systèmes de régulation ou des contrôles dynamiques du mouvement reposant sur des données sensorimotrices. Par ailleurs, rien n'empêche que des mécanismes de corrélation de bas niveau correspondent à un conditionnement pavlovien.

Ainsi, la frontière entre l'activité corporelle et la sémiose n'est pas déterminée a priori par l'introduction des relations sensorimotrices mais par le fait que d'un côté les relations sensorimotrices se révèlent implicites et figées, et de l'autre explicites et modifiables. En reprenant les termes du constructivisme piagétien, les contingences sensorimotrices représentent les primitives et la perception des objets dans l'espace résulte de l'autopoïèse sémiotique. Mais la frontière entre ces deux pôles ne peut pas être précisément identifiée à ce stade de la réflexion. Les boucles sensorimotrices se comprennent soit comme des relations sensorimotrices exprimées par les primitives, soit comme des relations exprimées en termes de croyances liées à l'action. La première forme représente les boucles sensorimotrices primaires et la seconde forme les boucles sensorimotrices secondaires. Ces dernières constituent le premier niveau d'abstraction et se fondent sur les boucles sensorimotrices primaires puisque ces dernières se situent entre les commandes motrices et les effecteurs.

**

L'architecture cognitive ainsi définie présente deux propriétés fondamentales pour réaliser une autopoïèse sémiotique. La première réside dans la boucle de rétroaction entre la mémoire événementielle et la mémoire épistémique. En effet, le résultat de l'interprétation du contenu de la mémoire événementielle vient modifier celle-ci ce qui produit dans le temps une dynamique interprétative. Par ailleurs, la seconde propriété provient de l'utilisation des primitives motrices et sensorielles à travers la mémoire événementielle. Cette dernière devient alors à la fois le support de la dynamique interprétative propre et des perturbations qui conduisent à des transformations d'état ou de dynamique. Autrement dit, l'architecture cognitive permet de réaliser un couplage par clôture (Varela, 1989) tel qu'il a été présenté dans le premier chapitre et qui est indispensable à tout autopoïèse. Maintenant que les principaux traits de l'architecture cognitive ont été établis, il reste à définir la dynamique sémiotique qu'elle doit accueillir et les caractéristiques supplémentaires à l'architecture cognitive qu'elle induit.

2.3. La sémiose au sein d'une architecture cognitive

Afin de transposer au mieux la notion de sémiose à l'activité d'une architecture cognitive, une étude préalable sur (A) la sémiose dans le cadre du pragmatisme sera présentée. Celle-ci se fondera sur la sémiose peircienne élaborée à partir d'une triade conceptuelle. Cependant, cette présentation révélera que la sémiose peircienne repose sur la capacité du sujet à concevoir et à percevoir des objets, des capacités dépassant l'autopoïèse sémiotique élémentaire imaginée. Afin de proposer une alternative à la notion d'objet, (B) une étude sur la correspondance des caractéristiques de l'architecture cognitive et des caractéristiques de la sémiose peircienne sera entreprise, dont il en résultera une nouvelle triade sémiotique. (C) Une fois la définition de l'architecture cognitive complétée, il sera proposé une mise en perspective avec la généalogie de la cognition afin de dégager les schèmes cognitifs d'une autopoïèse sémiotique minimale et ceux pour les stades supérieurs.

A - Définition de la sémiose

Dans le chapitre précédent, l'abandon des hypothèses ontologiques a conduit à replacer la problématique de la cognition dans la perspective de la sémiotique peircienne qui privilégie l'inférence à la référence : « Un signe est quelque chose qui tient lieu pour quelqu'un de quelque chose sous quelque rapport ou à quelque titre » (Peirce, 1897). Cette conception conduit Peirce à percevoir la sémiose comme l'interaction d'un representamen, d'un objet (de pensée) et d'un interprétant. En reprenant ces termes, le signe s'assimile au representamen qui renvoie pour son interprétant à la croyance d'un objet (Figure III-5). Contrairement à la référence qui souhaite relier une entité indépendante du sujet et le signe perçu, l'inférence relie le signe à l'objet pensé par le sujet. En se limitant à l'inférence et à la croyance, la sémiotique se sépare des hypothèses ontologiques. Toutefois, Peirce conserve une position réaliste et considère la démarche scientifique comme une sémiose particulière permettant de faire converger l'idée de l'objet vers l'objet réel. Cependant, cette conception de vérité-correspondance en limite s'est révélée intenable et la conception de la vérité de James (1907) développée dans le second chapitre reste maintenue.

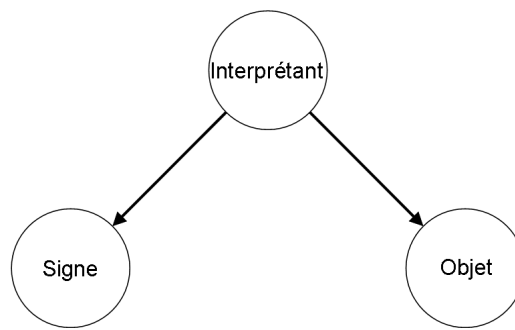


Figure III-5 : Triade sémiotique de Peirce (1904)

La sémiotique imaginée par Peirce propose pour chacune des notions une classification composant la triade sémiotique : le representamen (sur sa nature : qualisigne, sinsigne et légisigne ; sur sa signification : icône, indice, et symbole), l'objet (objet immédiat, objet dynamique) et l'interprétant (interprétant immédiat, interprétant dynamique et interprétant final). Une présentation détaillée de cette ramification serait nécessaire pour s'assurer que des schèmes cognitifs dans le cadre de l'architecture cognitive proposée permettent d'illustrer toutes ces distinctions. Cependant, la conception de l'architecture cognitive se voulant élémentaire, seule la triade sémiotique peircienne sera utilisée dans le cadre de cette thèse.

Selon la sémiotique pragmatiste, le signe en tant que tel ne comporte pas en lui-même les raisons qui le lient à un objet puisqu'en définitive c'est l'interprétant qui les relie. Cette liberté d'association entre le signe et l'objet autorise le fait que plusieurs signes puissent évoquer le même objet et que la relation unissant les signes et les objets puissent être de nature arbitraire, soit un code. Ces principes corroborent les propos sur la sémiotique dans le chapitre précédent qui se résumait de la manière suivante : « le signifié cesse d'être une entité psychique, ontologique ou sociologique : c'est un phénomène culturel, descriptible grâce à un système de relations que le code nous montre comme reçues par un groupe donné à un moment donné. » (Eco, 1987) Mais cette triade sémiotique conserve la notion « d'objet appartenant au monde pour le sujet » qui selon la généalogie de la cognition, vient à la suite de la constitution d'une autopoïèse sémiotique minimale. Par conséquent, la notion d'objet de la triade sémiotique doit être remise en question.

Cette remise en question va s'appuyer sur les deux jugements constituant la perception : le jugement perceptif et le jugement existentiel. En effet, le premier jugement univoque et indubitable ne nécessite pas la notion d'objet, bien qu'il puisse s'en servir, alors que le second jugement intervient directement dans la prise de conscience de l'objet devant se trouver dans le monde. Par ailleurs, le jugement d'assertion et le jugement de vérité offrent une situation analogue : le jugement d'assertion qui est également univoque et indubitable ne présuppose pas un objet conceptuel contrairement au jugement de vérité. En considérant que les premiers stades de l'évolution cognitive reposent uniquement sur le jugement perceptif et que la conception ne vient qu'avec la production de nouvelles interprétations donc après la définition de l'architecture cognitive, la nécessité de la notion d'objet pour réaliser une autopoïèse sémiotique minimale s'évanouit. Ainsi, une tension apparaît entre la généalogie de la cognition et la sémiotique peircienne qui ne peut être dépassée qu'en comprenant la manière de l'ajuster à l'architecture cognitive sans la dénaturer.

B - Architecture cognitive et sémiologie peircienne

Afin de répondre à cette question, la notion d'objet doit être examinée au regard de la sémiologie et des caractéristiques de l'architecture cognitive dégagées précédemment. L'objet se comprend comme le résultat de l'interprétation d'un signe, soit une évocation à la conscience d'une entité dans le monde. De manière générale, le terme « objet de pensée » renvoie à l'existence d'un signe objectal muni d'un ensemble d'opérateurs et d'implications qui explicitent les enjeux et les interactions possibles selon le contexte. La notion d'objet de pensée englobe ici la notion de concept général et celle de concept particulier ainsi que toutes abstractions possibles sur les objets ou les concepts pensés. Dans ce cadre, un objet peut être assimilé à un nœud d'un réseau sémiotique à partir duquel commence le cheminement de la reconnaissance ou de l'anticipation des interactions et des modifications entre l'entité et le monde. La notion d'objet associée à celle d'objectif qui correspond à l'explicitation pour le sujet d'un besoin donné autorise la mise en problématique d'une situation pour évaluer les enjeux ou pour anticiper les conséquences à la troisième personne de divers scénarii. L'évocation de l'objet devient alors elle-même un signe puisque l'attitude du sujet dépendra de son interprétation en fonction du contexte (Figure III-6).

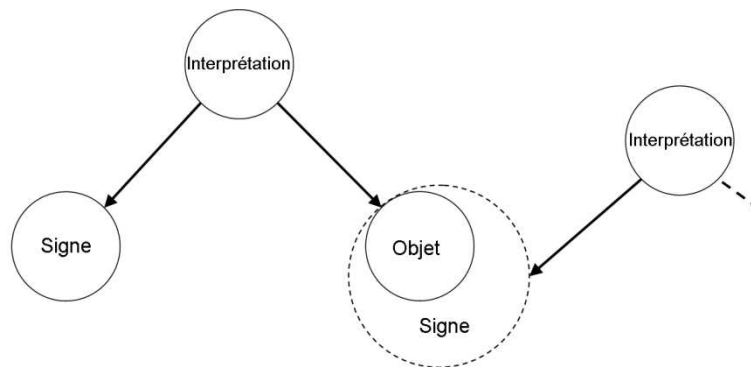


Figure III-6 : Enchaînement de triade sémiotique

Le fait que l'objet résultat de l'interprétation soit considéré comme un signe pour une autre interprétation correspond à la réflexivité de l'architecture cognitive produite par interaction entre une mémoire événementielle et une mémoire épistémique (Figure III-2). La dernière source, la mémoire épistémique, correspond au troisième pôle de la triade sémiotique : l'interprétant. Toutefois, rien ne suggère dans la réflexivité entre la mémoire événementielle et la mémoire épistémique de cette architecture cognitive que toutes les évocations provenant de la mémoire épistémique soient des signes objectaux ou des symboles, c'est-à-dire des objets de pensée comme le suggère la Figure III-5 ou la Figure III-6.

En effet, de manière minimaliste, l'interprétation peut conduire directement à l'action sans que l'évocation d'un objet passe dans la mémoire événementielle ou à des conclusions internes sans dimension conceptuelle. Cette version minimaliste correspondrait par exemple à la situation de la cane de Barbarie qui sauve un caneton col-vert mais le chasse juste après. La cane reçoit comme un signe le cri du caneton qu'elle interprète selon le contexte global comme un signal d'alarme qui déclenche alors un comportement de protection des petits, puis, le calme revenu, le signe produit par le motif au sommet du crâne du petit devient prioritaire et provoque un comportement agressif. Cet exemple montre le total syncrétisme de la cane de Barbarie qui ne permet pas un comportement cohérent dans cette situation. En d'autres termes, « dans le comportement animal, les

signes restent toujours des signaux et ne deviennent jamais des symboles » (Merleau-Ponty, 1942). Un comportement plus cohérent résulterait de l'ajout soit d'une règle d'interprétation supplémentaire liée à un nouveau signe, soit d'une théorie homogène liant ces signaux, un objet.

Une autre manière de présenter ces deux voies consiste à considérer que les principes sous-tendant le jugement perceptif et le jugement existentiel (de même que le jugement d'assertabilité et le jugement de vérité) font partie de l'architecture cognitive mais qu'ils nécessitent des schèmes cognitifs pour exploiter leur potentiel. La signification et l'intérêt d'un jugement existentiel et d'un jugement de vérité sans la conception d'objet physique ou conceptuel seront avancés lors de la formalisation de l'architecture cognitive proposée. Cependant, dans le cadre de ce travail de thèse qui porte exclusivement sur la spécification d'une architecture cognitive, la conception et l'élaboration de schèmes cognitifs complexes tels que les jugements liés à la notion d'objets appartiennent aux perspectives de recherche.

Cette réflexion sur le fait de reculer la définition de la notion d'objet dans l'artificialisation de la cognition donne l'occasion de souligner l'une des particularités de l'approche défendue ici. En effet, les approches classiques revendiquant au moins une hypothèse ontologique définissent de suite ce que doit être un objet en soi et par suite pour le sujet (pour soi). Ainsi, ces approches abordent directement la notion d'objet qui fonde, pour elles, la cognition. Dans ces conditions, la différence entre signaux et symboles devient une différence hiérarchique entre les niveaux de traitements de l'information. Un symbole représente alors un nœud reliant différents signaux et ce nœud peut lui-même être considéré comme un signal à un autre niveau de la hiérarchie, ce qui signifie en définitive que ces approches n'introduisent pas de différence qualitative en signaux et symboles. Les symboles dans ces approches ne correspondent pas aux objets de pensées complexes dont l'interprétation dynamique évolue toujours en fonction des autres. Le concept de nœud n'est qu'une simplification du concept d'objet.

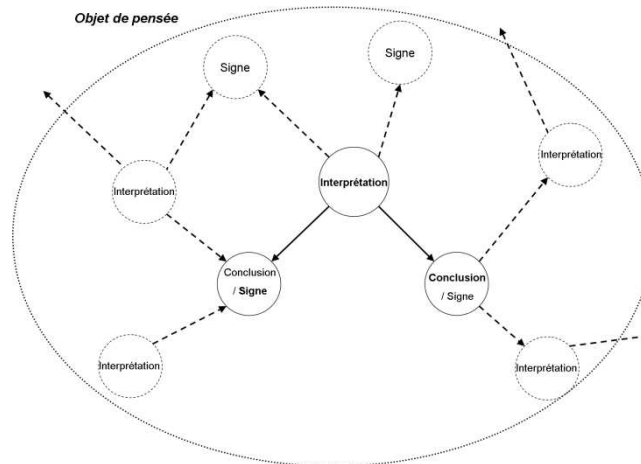


Figure III-7 : Triade sémiotique constituée par le signe, l'interprétation et la conclusion au sein d'un réseau sémiotique

La critiques sur les fondements des approches en sciences cognitives dans le premier chapitre et du concept de vérité qui leur est associé dans le deuxième chapitre a permis de révéler que la notion d'objet représente un saut qualitatif qui arrive en réalité à la fin du développement de la cognition. Par conséquent, la notion d'objet ne peut pas servir de base à la conception d'un système cognitif, contrairement à ce que font implicitement ou

explicitement toutes les approches qui ne se fondent pas sur l'émergence. Avant de comprendre la notion d'objet, il est nécessaire de dégager tous les principes d'une cognition primitive, une sémiologie élémentaire. Au mieux, le concept d'objet renvoie à un sous-réseau de croyances fortement connecté mais dont l'activité ne se réduit pas à la production d'un signe commun (Figure III-7).

La triade sémiotique peircienne devient alors le signe, l'interprétant et la conclusion. La conclusion correspond au résultat de l'interprétation c'est-à-dire, en dehors de toute activité non perceptible par le système sémiotique, quelque chose qui fasse signe ou tout au plus un signe ou des signes associés à un ensemble complexe de relations permettant l'évocation d'un objet. L'introduction de la notion de conclusion se répercute dans l'interprétation des hypothèses concernant la conception et la perception abordées à la fin de ce chapitre. Mais elle n'entraîne pas la fin de la réinterprétation de la triade sémiotique peircienne par rapport à la notion d'architecture cognitive. En effet, une interrogation reste en suspens : quelles contraintes impliquent le jugement perceptif et la notion d'autopoïèse sur la triade sémiotique ?

La réponse à cette interrogation repose sur l'analyse des deux propriétés attribuées au jugement perceptif dans le premier chapitre, qui donneront les éléments pour transposer les principes de l'autopoïèse matérielle à l'autopoïèse sémiotique. La première des deux propriétés correspond à l'indubitabilité de la perception, soit du jugement perceptif. Cette notion a été expliquée dans le premier chapitre avec l'exemple de Pierre pouvant douter de ce qui est à l'origine de son impression de pomme mais qu'indubitablement il a l'impression de voir. Mais cet exemple introduit la notion d'objet et de jugement existentiel. En évacuant ces deux concepts, une autre manière de présenter le caractère indubitable de la perception consiste à dire que le processus de la perception s'appuie sur des signes tout en en produisant et que l'existence de ces signes ne peut être remise en question, seule la provenance peut être sujette à de nouvelles interprétations. Ainsi, en étendant ce principe à toute interprétation, la mémoire événementielle devient indubitable, ce qui implique une totale intégrité, c'est-à-dire qu'en dehors du fonctionnement propre à la mémoire événementielle, les signes mémorisés ne peuvent pas être éliminés ou modifiés par le jeu de nouvelles interprétations. En définitive, cette propriété correspond à la nécessité qu'une segmentation du monde soit stable durant l'élaboration ou l'application d'une théorie.

La seconde propriété du jugement perceptif, l'univocité, correspond au fait qu'à un instant donné, la perception constitue une et une seule interprétation du monde. Dans le premier chapitre, cette impossibilité à superposer des perceptions a été illustrée par le dessin de Jastrow (1901) qui s'interprète traditionnellement comme l'illustration d'une tête soit de canard, soit de lapin. En considérant la perception comme un processus d'interprétation de signes et de production de signes, soit comme une sémiologie, l'univocité implique que pour un champ perceptif donné une seule interprétation de la mémoire événementielle peut s'effectuer. En même temps, il faut souligner que tout champ perceptif produit une interprétation, autrement dit, un sujet éveillé perçoit toujours quelque chose. Exprimé de manière générale, chaque champ spécifique (perceptif, proprioceptif ou conceptuel) possède à un instant donné une seule interprétation. Cette contrainte apparaît indispensable pour la prise de décision puisqu'au final un effecteur ne peut recevoir qu'un seul ordre moteur.

L'univocité porte également sur les signes exprimant une attente dans le cadre de certains champs perceptifs, conceptuels ou autres, autrement dit, deux prédictions au même instant ne peuvent être prises en compte. L'univocité implique alors l'élimination des

signes exprimant une attente lors de l'arrivée de nouveaux signes sur le même champ spécifique. La Figure III-8 représente les champs spécifiques devant être considérés au sein de l'architecture cognitive afin d'assurer l'univocité des conclusions et l'intégrité de la mémoire événementielle. Ce mécanisme appartenant au fonctionnement interne de la mémoire événementielle, l'intégrité de la mémoire reste conservée. Par ailleurs, bien que le terme de champ sera précisé lors de l'explicitation de la structure interprétative élémentaire, la notion d'univocité présentée comme résultant d'une compétition entre diverses interprétations se révèle suffisante pour expliquer la relation non inférentielle qui définit les signes et leurs interprétations entre eux.

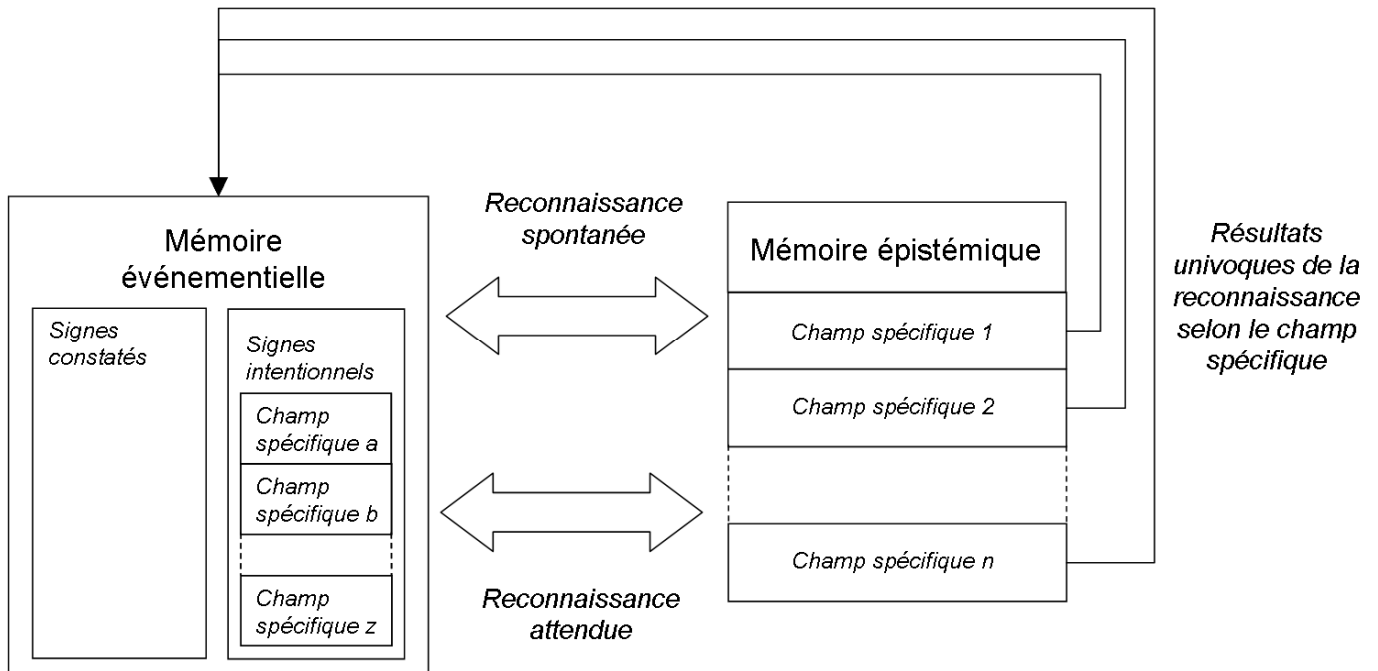


Figure III-8 : Introduction des champs spécifiques nécessaires pour assurer l'univocité au sein de l'architecture cognitive.

La notion de compétition entre les diverses interprétations permet d'écartier définitivement le risque de réintroduire l'HOI. En effet, l'architecture cognitive repose sur une mémoire épistémique et sur une mémoire événementielle, la première interprétant le contenu de la seconde. Cela ne signifie pas qu'une interprétation se définisse par la reconnaissance d'un concept précis dans les signes fournis par la mémoire événementielle car ce type de reconnaissance impliquerait que les concepts correspondent à des objets conceptuels dans un monde idéal. L'interprétation est une reconnaissance d'une situation par rapport à d'autres et le résultat de cette compétition dépend du vécu du sujet, de l'interprétant. L'existence de cette compétition se retrouve alors dans les interactions sociales et culturelles. Cette compétition au niveau socioculturel s'illustre par l'étude comparative de Hjelmslev (1957) des champs sémantiques de plusieurs langues pour un même thème (Tableau III-1), puis repris par Eco (1987) avec le Tableau III-1 : « Un tableau de ce genre ne nous met pas en face « d'idées » [compris comme objet du monde idéal], mais de valeurs émanant du système. Ces valeurs correspondent à ce que l'on peut nommer des concepts, mais ne naissent et ne peuvent être appréhendées que comme pures différences : elles ne se définissent pas par leur contenu, mais par la manière dont elles s'opposent aux autres éléments du système. »

Français	Allemand	Danois	Italien	Anglais
arbre	baum	trae	alberto	tree
bois	holz		legno	timber
	forêt	wald	skov	bosco
				foresta

Tableau III-1 : Comparaison des champs sémantiques entre le français, l'allemand, le danois, l'italien et l'anglais pour décrire un paysage contenant un ou plusieurs arbres (Eco, 1987).

L'univocité impliquant la notion de compétition offre également le moyen de mieux saisir la dynamique de l'autopoièse sémiotique. Une autopoièse exprime la lutte contre sa propre dégénérescence. Pour une autopoièse matérielle, cette dégénérescence se traduit par la perte de ses constituants par des réactions chimiques d'origine interne ou externe. Afin que l'autopoièse sémiotique se comprenne indépendamment de l'autopoièse matérielle, il convient que les principes de la dégénérescence des constituants de la sémiologie ne soient pas totalement liés à la dégénérescence de son substrat. Plus exactement, la dégénérescence d'une autopoièse sémiotique doit porter à la fois sur les signes se trouvant dans la mémoire événementielle et sur les croyances dans la mémoire épistémique. Concernant la mémoire événementielle, l'oubli devient une source de dégénérescence naturelle puisque la mémoire événementielle intègre une notion temporelle au signe par rapport à leur survenance. La modification au cours du temps d'un indice temporel intrinsèque au signe conditionne l'effacement du signe, mais ce processus étant propre au fonctionnement de la mémoire événementielle, l'intégrité de celle-ci n'est pas compromise.

La dégénérescence par l'oubli se trouve compensée par l'afflux permanent de signes sensoriels, de signes proprioceptifs ou de signes résultant de la conclusion d'interprétations. L'oubli représente également une source de dégénérescence pour les croyances de la mémoire épistémique. Cependant, la signification de l'oubli diffère de celle employée pour la mémoire événementielle dans le sens où les traces temporelles liées aux croyances ne participent pas à l'interprétation de la mémoire événementielle mais simplement à l'élimination des croyances qui ne se réalisent plus. L'enjeu de la compétition entre les différentes interprétations prend alors tout son sens : afin de ne pas être oubliée, une croyance doit être utilisée au dépend des autres. Les opérateurs de reconnaissance se complètent d'une évaluation générale des interprétations candidates qui ne peut se concevoir qu'à partir de la définition de la structure interprétative considérée. Par ailleurs, cette compétition doit renforcer la dégénérescence en plus du phénomène d'oubli. Ainsi, même en admettant que l'oubli puisse être lié à la dégénérescence de l'autopoièse matérielle, la dégénérescence due à la compétition ne découle pas directement de l'autopoièse matérielle qui sous-tend l'autopoièse sémiotique mais provient de la sémiologie c'est-à-dire de l'activité propre à l'autopoièse sémiotique. Ainsi, cette notion de compétition entre croyances au cœur de l'architecture cognitive ne va pas a priori à l'encontre d'une artificialisation de la cognition.

Toutefois, à ce stade, l'idée d'une compétition permanente pour résister à l'oubli n'explique pas que certaines situations fassent sens bien qu'elles aient été apprises longtemps auparavant. Ce fait signifie que certaines croyances demeurent au sein de la mémoire épistémique malgré qu'elles ne soient plus utilisées. Une manière de relier ce phénomène avec celui de la compétition consiste à considérer qu'il existe, en plus d'évaluations pour une situation donnée, une évaluation globale de l'activité d'une croyance

dont dépendra l'oubli. Ces deux types d'évaluations correspondent respectivement à une sorte d'évaluation d'une probabilité événementielle et d'une probabilité épistémique. Cette double évaluation autorise une compétition perpétuelle mais qui n'affecterait principalement que les croyances n'ayant pas fait leur preuve.

Par rapport à la définition de l'autopoïèse sémiotique, l'analyse de la sémiose au sein de l'architecture cognitive permet de retrouver la notion de capacité à éprouver et la notion de relation effective qui sous-tend le caractère concret de l'autopoïèse sémiotique. En effet, la capacité à éprouver la structure du système sémiotique correspond à la lutte des croyances contre l'oubli mais également au fait que celles-ci s'ajustent en fonction de leur activité et du contexte. De même, une interprétation représente une relation effective puisque le résultat de l'interprétation se traduit concrètement par l'introduction d'un signe dans la mémoire événementielle et/ou par une activité subcognitive. Par ailleurs, l'interprétation résultant de la compétition d'un ensemble de croyances, cet ensemble forme toujours une unité. Ainsi, les propriétés principales liées à la notion d'autopoïèse (dynamique interne, couplage par clôture, unité concrète) ayant été traduites dans l'architecture cognitive, une réflexion sur la définition des schèmes cognitifs produisant une autopoïèse sémiotique minimale peut maintenant être entreprise.

C - Étude sur l'architecture cognitive et les schèmes cognitifs au travers de la généalogie de la cognition

L'analogie avec l'autopoïèse physique a dégagé le rôle d'une architecture cognitive et celui des schèmes cognitifs. Toutefois, il reste à déterminer quelles sont les propriétés relatives à l'autopoïèse sémiotique minimale (un couplage ponctuel initial), et celles relatives à son épanouissement par le couplage structurel. L'étude de l'ontogenèse d'un individu permet difficilement de déterminer ces propriétés. En effet, pour les individus cognitifs naturels, l'évolution des capacités cognitives repose à la fois sur une autopoïèse physique et sur une autopoïèse sémiotique. De manière caricaturale, dans un premier temps, l'autopoïèse physique prend le pas sur l'autopoïèse sémiotique dans les raisons du développement des capacités cognitives, dans un deuxième temps la prédominance est inversée, enfin dans un troisième temps, la dégénérescence du substrat ne permet plus le développement de l'autopoïèse sémiotique. L'autopoïèse sémiotique de l'individu se déroule dans la continuité bien que le développement de certaines capacités cognitives soit lié au développement de l'organe. Afin de se soustraire de l'influence du développement du substrat et de se limiter à une espèce, l'identification des propriétés liées à l'établissement d'autopoïèse sémiotique minimale et à son développement se fondera alors sur l'étude de la généalogie de la cognition avancée au cours du chapitre précédent. Celle-ci a déterminé principalement quatre stades d'évolution situés après celui des comportements réflexes primitifs. La présentation de ces stades se concentrait sur leur intérêt vis-à-vis du dépassement de l'incomplétude de l'ensemble des schèmes comportementaux. Ici, l'évocation à la généalogie cognitive servira à identifier pour chaque stade les propriétés dynamiques d'une autopoïèse sémiotique nécessaire à l'émergence des capacités cognitives correspondantes.

i - Le premier stade de la généalogie de la cognition

Contrairement au stade zéro dans lequel les arcs réflexes ou autres mécanismes de régulations ne dépendent pas du vécu de l'organisme, au premier stade, le déclenchement des comportements d'un individu s'ajuste en fonction de son environnement sans que cela entraîne une modification de son répertoire d'action et sans que cet ajustement soit orienté

vers une finalité. Cet ajustement traduit une modulation de la segmentation du monde, soit une modulation des couples stimulus-réponse. Par rapport à la discussion sur la sémiose, ce stade se conçoit comme la mise en place de l'interdépendance des signes et de leurs interprétations qui les définissent.

En fait, un couple stimulus-réponse se traduit par une classe d'équivalence de stimuli associée à une réponse particulière, et la modulation d'un couple stimulus-réponse correspond à la modification de cette classe d'équivalence (Figure III-9). La modulation de la croyance ayant servi à l'interprétation influencera la compétition de l'interprétation suivante. Par ailleurs, comme la modulation se produit en fonction des expériences, l'ajustement des croyances peut se comprendre comme un processus d'auto-organisation. Cependant, cette modulation repose uniquement sur une auto-organisation de la segmentation du monde. En effet, le répertoire des actions reste fixe.

Un couple stimulus-réponse représente alors une croyance de la mémoire épistémique. Mais la réponse peut ne pas être explicite à ce stade, c'est-à-dire que l'interprétation (ou l'application de la croyance) conduit directement à une commande motrice sans qu'un signe la représente dans la mémoire événementielle. Par ailleurs, l'auto-organisation des classes d'équivalences constituant les croyances ne provient pas d'un schème cognitif particulier mais de mécanismes subcognitifs appartenant à l'architecture cognitive. Ainsi, une autopoïèse minimale correspond à une architecture cognitive avec un ensemble de croyances dont la forme respecte le schème cognitif minimal que représente le couple sensorimoteur.

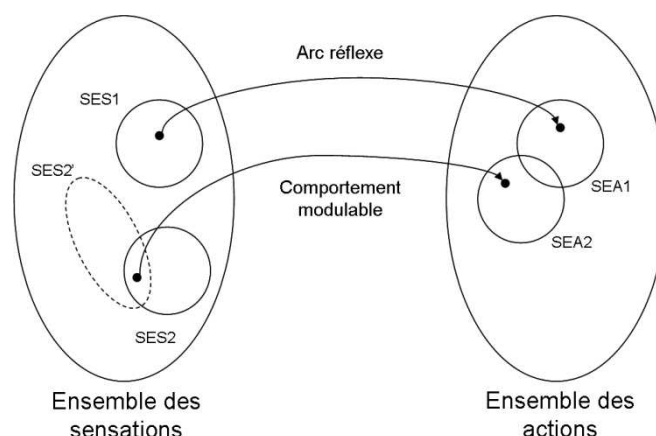


Figure III-9 : Schéma illustrant la différence entre l'arc réflexe et le comportement modulable. Le premier type de comportement relie chaque élément appartenant à un certain sous-ensemble de sensations SES1 à une action appartenant à un sous-ensemble SEA1. Cette relation est figée, elle n'évolue pas en fonction des expériences du sujet. En revanche, le second type de comportement peut voir son ensemble de départ SES2 se modifier en SES2' selon les expériences du sujet.

ii - Le deuxième stade de la généalogie de la cognition

Au deuxième stade, la notion d'acte finalisé apparaît, le résultat d'une action peut inférer sur l'ajustement du comportement ou inciter à la création d'un nouveau couple stimulus-réponse. En d'autres termes, le comportement de l'individu alterne les phases d'instinct avec celles de dressage (Lorenz, 1937).

Dans un premier temps, la réalisation de ce dressage implique la capacité à rétribuer une croyance afin de la favoriser ou défavoriser. Une rétribution positive augmente l'évaluation de la pertinence générale de la croyance et la diminue dans le cas inverse.

L'évaluation de la pertinence globale de la croyance conditionne à terme son élimination. L'acte de rétribuer correspond à un mécanisme subcognitif qui peut être déclenché par l'interprétation d'une situation. À ce stade, la rétribution peut être explicite ou implicite. Par ce mécanisme, le dressage correspond au déclenchement d'une rétribution via une croyance afin de favoriser les croyances précédentes susceptibles d'avoir participé à la situation actuelle. Par ailleurs, comme il a été évoqué lors de la présentation du fonctionnalisme écologique dans le premier chapitre, une manière d'évaluer la stabilité de ces nouvelles croyances consiste à en prédire les conséquences. Ainsi, dès le deuxième stade, l'opérateur de la reconnaissance attendue participe à l'autopoïèse sémiotique.

Dans un second temps, le dressage permet dès lors l'exploration de l'environnement, ce qui implique la génération de nouvelles croyances. Afin d'être viable, ces nouvelles croyances se construisent à partir de l'expérience nécessitant alors que les actions deviennent explicites, c'est-à-dire médiatisées par des signes. L'introduction de nouveaux couples sensorimoteurs complète la segmentation du monde mais elle entraîne également sa réorganisation.

Bien que la génération de nouvelles croyances puisse être déclenchée à partir de l'expérience et utiliser celle-ci, les modalités de production peuvent être définies entièrement dans le cadre de l'architecture cognitive puisque, à ce deuxième stade, les modalités de production ne sont pas vouées à évoluer au cours de l'activité de l'individu. Le deuxième stade de la cognition introduit uniquement de nouveaux schèmes cognitifs liés à la rétribution interne et à la prédiction.

iii - Le troisième stade de la généalogie de la cognition

Au troisième stade, l'individu est capable de décider, c'est-à-dire de choisir une action parmi d'autres en fonction des besoins. Comme il a été vu dans le premier chapitre au cours de la présentation du fonctionnalisme écologique, la décision oblige l'évaluation des gains. L'identification des besoins et des gains passe par une représentation de l'agréable ou du désagréable. En effet, contrairement au deuxième stade, dans lequel certaines situations déclenchent un mécanisme qui va renforcer ou pénaliser certaines croyances, au troisième stade, ce renforcement doit devenir explicite. Cette explicitation introduit une certaine réflexivité puisque le système se rétribue selon ses propres croyances et interprète ses rétributions : de l'agréable et du désagréable dérive la notion de nuisible et celle de profitable.

Par ailleurs, cette réflexivité se retrouve dans la production de croyances. Plus particulièrement, au troisième stade, un individu est capable d'identifier une finalité donc de produire des croyances sur l'agréabilité mais pour éviter les dérives, cette production doit être éprouvée. En d'autres termes, les modalités pour générer une croyance donnant une rétribution doivent être exprimées sous forme de croyances, ces dernières décrivent alors un schème cognitif. De la même manière que pour le deuxième stade, la capacité à éprouver est associée à celle de produire ce qui doit être éprouvé soit, ici, la capacité à générer de nouvelles croyances définissant les modalités de génération d'autres croyances. À ce stade, les modalités sur la génération de ces schèmes cognitifs reposent sur les mécanismes subcognitifs de l'architecture cognitive. L'explicitation de la rétribution associée à la capacité à générer des générateurs de croyance offre des pistes de réflexion pour exprimer le concept d'objet qui n'est pas tant lié à un ensemble de croyances statiques qu'à un ensemble de croyances produisant continuellement de nouvelles croyances dans un certain contexte. De même, l'anticipation devient concevable dès le troisième stade de l'évolution de la cognition.

iv - Le quatrième stade de la généalogie de la cognition

Au dernier stade de l'évolution cognitive apparaît la capacité d'abstraire et d'analyser ses propres motivations. En effet, bien que l'apparition de la conscience de soi puisse commencer au troisième stade comme une simple projection d'un objet particulier dans une simulation interne, elle n'aboutit pleinement qu'à ce quatrième stade avec la réflexion critique du soi. Ce stade suppose alors une réflexivité totale de l'autopoièse sémiotique, de sorte que les croyances finalisées ou non, soient en définitive toujours révisables. En d'autres termes, les modalités liées à l'explicitation de schèmes cognitifs par l'intermédiaire de croyances doivent être elles-mêmes explicitées par des croyances de sorte que le développement des croyances permette de réaliser des récurrences cognitives propices à l'abstraction et que toute production de croyance puisse être éprouvée.

*

Ces quatre stades ne se veulent pas exhaustifs dans la liste des propriétés ou étapes aboutissant à une cognition de niveau équivalent à celle des humains, notamment concernant les deux derniers stades. L'attention, par exemple, n'a pas été évoquée alors qu'elle se trouve indissociable à la réalisation d'une segmentation dynamique du monde permettant de dégager une forme d'un fond, indispensable à la notion d'objet. De même, la problématique du langage ne doit pas être ignorée pour le développement de la sémiologie. Mais seuls les traits les plus fondamentaux ont été privilégiés pour la définition de l'architecture cognitive. En effet, l'attention peut être comprise comme une problématique liée à l'interprétation des signes suivant la dynamique sémiotique du système et de ses conséquences sensorimotrices, ce qui correspond alors à des schèmes cognitifs sans propriétés particulières ou du moins fondamentales pour l'architecture cognitive.

**

En conclusion de cette partie, les propriétés de l'architecture cognitive reposent essentiellement sur la notion d'autoréférence de part l'explicitation des interprétations dans la mémoire événementielle comme les actions, les rétributions et de part également l'explicitation des modalités de production de croyances en termes de croyances. L'autopoièse sémiotique peut se concevoir indépendamment d'une autopoièse physique mais son développement dépendra toujours des caractéristiques du corps qui l'incarne. La définition d'une architecture cognitive représente le pré-requis à l'établissement d'une autopoièse sémiotique minimale sur laquelle pourront se greffer divers schèmes cognitifs dont certains permettraient un couplage structurel auto-finalisé grâce à l'explicitation de la rétribution interne. Néanmoins, les limites cognitives d'un robot pourvu d'une telle architecture résulteraient uniquement des capacités d'interactions avec l'environnement et avec lui-même, et surtout des schèmes cognitifs initiaux. En effet, les stades cognitifs dégagés offrent une réelle graduation dans les capacités cognitives, de la même façon qu'il existe une telle graduation chez les animaux pourvus de capacités cognitives.

La Figure III-10 propose le schéma structurel d'un individu cognitif dans son ensemble. Dans ce schéma, les fréquences associées aux flux d'informations illustrent davantage le fait que les processus les générant ou les utilisant fonctionnent de manière asynchrone et non que ces flux soient nécessairement périodiques. En revanche, la notion de flux cadencé est incompatible avec le flux de commande destiné aux primitives motrices, puisque celui-ci dépend de l'interprétation. Enfin, f_4 représente la fréquence des interprétations de la mémoire épistémique sur la mémoire événementielle. Les flèches à double sens reliant les primitives et le coupe senseur/effecteur symbolise le fait que les

primitives puissent reposer sur un couplage sensorimoteur comme le réflexe de contraction de l'iris ou des processus de régulation pour atteindre des consignes imposées par des commandes. Maintenant que les propriétés et les caractéristiques générales de l'architecture cognitive ont été dégagées, l'identification des principes logiques sous-tendant l'architecture cognitive proposée devient envisageable.

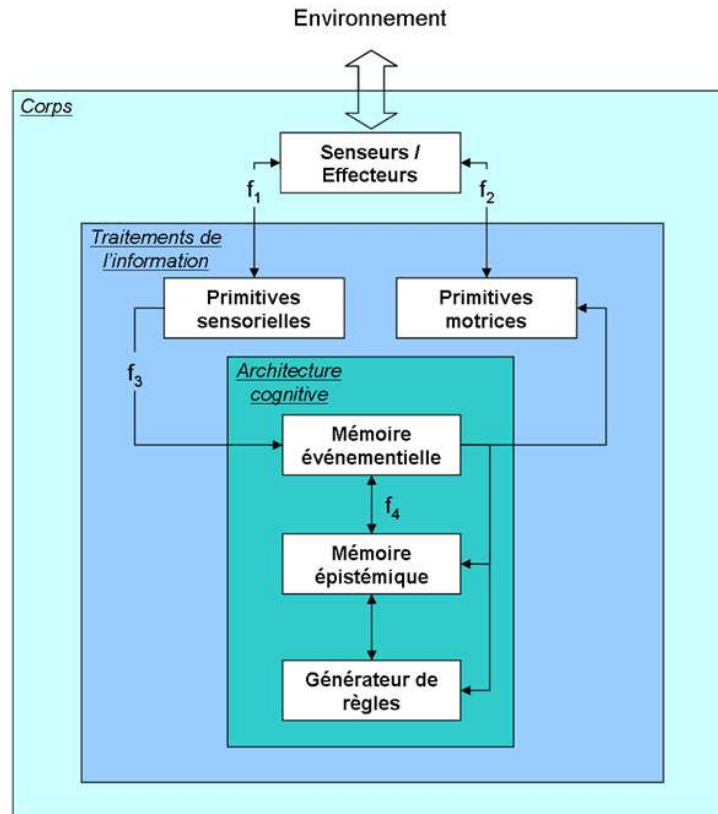


Figure III-10 : schéma d'un système muni d'une architecture cognitive avec les fréquences possibles pour les flux d'informations

3. Formalisation de l'architecture cognitive

La mémoire épistémique représente l'ensemble des procédés permettant d'aboutir à une interprétation de la mémoire événementielle. Indépendamment de la complexité des procédés mis en œuvre, la mémoire épistémique peut être assimilée à une mémoire associative et peut se modéliser par un ensemble de règles d'associations. Dans ce cas, une croyance devient une règle qui se compose d'une prémisse et d'une conclusion. La prémisse contient une représentation des signes escomptés, soit une classe d'équivalence, et la conclusion représente un signe et/ou un lien avec un mécanisme subcognitif (rétribution interne, commande motrice ou génération de règles). L'interprétation représente l'application d'une croyance en fonction de l'état de la mémoire événementielle et des autres croyances. L'architecture cognitive constitue alors un système d'inférence subsymbolique dans lequel la relation entre signe et croyance est effective dans la mesure où d'une part l'application de la croyance se traduit effectivement par l'introduction d'un signe dans la mémoire événementielle ou par une action subcognitive, et d'autre part l'adéquation entre les signes et la croyance n'est jamais automatique puisqu'elle dépend de

la compétition entre les règles toujours perturbée par l'ajustement, l'élimination ou la génération de règles.

Toutefois, la compréhension des systèmes autopoïétiques passe par une description logique qui servira de base pour la spécification de l'architecture cognitive. La formalisation de l'architecture cognitive souhaite établir des relations conceptuelles entre l'architecture cognitive proposée et les systèmes formels utilisés pour modéliser des raisonnements. L'architecture cognitive, de part ses caractéristiques, possède de nombreux points communs avec un certain nombre de logiques formelles très diverses. L'objectif, ici, n'est pas de comparer ces systèmes logiques mais de reconnaître les ressemblances afin de formaliser au mieux l'architecture cognitive proposée et de dégager ses propriétés formelles. L'identification de ces ressemblances s'effectuera par la complexification successive de la formalisation de l'architecture au gré des remarques logiques. Le degré de formalisation s'arrêtera à un niveau suffisant pour situer l'architecture cognitive proposée par rapport aux systèmes logiques existants et pour pouvoir comparer l'architecture cognitive proposée et les systèmes de classeurs. La formalisation complète sera présentée en même temps que la spécification de l'architecture cognitive proposée (la troisième partie de ce chapitre) afin de pouvoir illustrer plus facilement la formalisation détaillée de certains mécanismes qui ne trouvent pas de comparaison au sein des systèmes formels ou au sein des systèmes de classeurs, et qui ne remettent pas en cause les remarques logiques sur l'ensemble de l'architecture.

La discussion sur la formalisation de l'architecture cognitive dans cette section se déroulera en six points :

1. Le premier point rappellera les principes des systèmes formels et proposera une formalisation simplifiée de l'architecture cognitive en ne considérant qu'un seul champ spécifique et qu'une seule règle d'inférence.
2. Le deuxième point proposera un parallèle sur les propriétés de la logique non monotone et les contraintes qu'impliquent les caractéristiques du jugement perceptif sur l'architecture proposée. Ce parallèle conduira à discriminer deux types de négation, l'une portant sur la présence et l'autre sur la relation. La première s'appuiera sur le cadre de la logique non monotone.
3. Le troisième point portera sur l'intégration de la seconde au sein de l'architecture cognitive, ce qui nécessitera d'examiner les différentes manières d'aborder la contradiction.
4. Le quatrième point précisera la signification des évaluations globales et ponctuelles des croyances au cours de la sémiose de l'architecture cognitive proposée tout en soulignant certaines analogies avec la logique floue ou les probabilités subjectives.
5. L'ensemble de ces précisions permettra au cinquième point de compléter en partie le schéma simplifié de l'architecture. Ce nouveau schéma révélera un premier lien avec la notion de type qui se rapportera à la gestion de la multiplicité des champs spécifiques. Le second lien avec la notion de type correspondra à l'élaboration d'une hiérarchie des relations entre les signes au cours de la sémiose. Ces analyses permettront notamment de définir formellement la notion de schème cognitif.

6. Enfin, le sixième point s'intéressera à comprendre comment la génération de croyances peut participer à la production d'abstractions c'est-à-dire réaliser des logiques d'ordre n . Cette analyse passera par l'introduction de la notion de variable au sein de l'architecture cognitive.

La conclusion sur ces six points visera à situer les propriétés logiques essentielles auxquelles se réfère l'architecture cognitive proposée, bien que celle-ci diffère des systèmes logiques évoqués de part son effectivité. Cependant, la dernière partie de ce chapitre montrera que, sous certaines conditions, ces aspects peuvent disparaître et permettent ainsi une équivalence entre l'architecture cognitive et un interpréteur d'un système logique.

3.1. Systèmes formels et première formalisation de l'architecture cognitive

Sommairement, un système formel décrit les conventions d'écriture et de manipulations de formules, soit décrit un langage selon le quadruplet suivant : un alphabet, des règles de formation des mots, une théorie formelle correspondant à un ensemble d'axiomes (des mots donnés a priori) et enfin des règles d'inférence. Le langage utilisé pour décrire les règles de formation ou d'inférence du système formel correspond au métalangage. Une formule est démontrable au sein de ce système si elle est déductible à partir des axiomes. Une formule vraie représente un théorème.

Par exemple, un métalangage élémentaire comporte les symboles suivants :

- a) « \forall » se lit « quelque soit la formule »
- b) « , » se lit « et »
- c) « \vdash » se lit « si...alors »

Avec ce métalangage, un système formel d'ordre zéro ne possédant qu'une règle d'inférence telle que le *modus ponens* peut alors se décrire de la manière suivante, en considérant l'implication notée \supset et un fait i notée X_i comme les symboles qui constituent l'alphabet du langage décrit :

$$(\forall X_1, X_2) X_1 \supset X_2, X_1 \vdash X_2$$

Ainsi, en posant le fait A et l'implication A et $A \supset C$ comme des axiomes, l'application de la règle d'inférence permet de déduire C et que ce dernier est, par conséquent, un théorème.

Les remarques logiques générales sur l'architecture cognitive ne nécessitent pas une formalisation exhaustive. Par conséquent, en première approximation, la formalisation de l'architecture cognitive ne tient compte ni de la provenance des signes ni de l'oubli des signes dans la mémoire événementielle. De même, cette première formalisation ne tiendra compte ni de l'élimination de certaines croyances, ni de la génération de nouvelles croyances dans la mémoire épistémique. De plus, il sera considéré également que cette dernière ne possède qu'un seul champ.

Concernant les règles d'inférence de l'architecture cognitive, la reconnaissance spontanée et la reconnaissance attendue produisent des déductions effectives et non des déductions logiques. Ainsi, la reconnaissance spontanée se comprend comme une règle

d'inférence effective du *modus ponens* ; pour marquer la différence, elle s'écrira de la manière suivante :

- d) « \mapsto » se lit « si...en fonction de l'ensemble des prémisses des autres implications appartenant à la mémoire épistémique alors...est mis dans la mémoire événementielle »

La seconde règle d'inférence effective, la reconnaissance attendue, pouvant se comprendre comme un opérateur modal d'une logique temporelle ne sera pas abordé lors de la première formalisation, elle sera toutefois présentée en détails au cours de la spécification. Par ailleurs, pour simplifier le discours, les signes, les conditions constituant les prémisses et les conclusions seront considérés comme des messages dont le rôle diffère selon leur position dans les formules. Les mécanismes de compétition et d'ajustement des interprétations constituent des mécanismes subcognitifs puisqu'ils ne participent pas directement à la manipulation des messages. Ces mécanismes n'apparaîtront donc pas dans cette première formalisation.

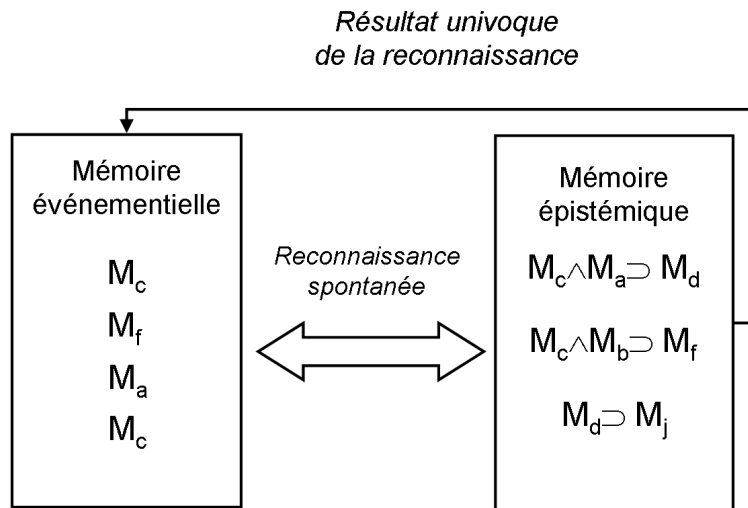


Figure III-11 : Architecture cognitive simplifiée.

En considérant que l'architecture cognitive peut être décrite comme une sorte de système logique avec les schèmes cognitifs représentant la théorie de ce système, le langage des schèmes cognitifs s'appuie alors sur les symboles suivants :

1. « M_a » se lit « message M_a avec le contenu a et une durée temporelle θ_a »
2. « \wedge » représente la conjonction
3. « \supset » représente l'implication assimilée à une croyance de la mémoire épistémique, soit une règle

En précisant que la mémoire événementielle contient uniquement des messages et que la mémoire épistémique contient uniquement des implications, le langage défini suffit à une première formulation de la reconnaissance spontanée :

$$(\forall \{M_i\}, M_j) \{ \wedge_i M_i \} \supset M_j, \{ M_i \} \mapsto M_j$$

Cette première formalisation de l'architecture cognitive, bien que simplifiée, offre déjà un certain nombre de propriétés qui peuvent s'exprimer dans le cadre de la logique non monotone.

3.2. L'architecture cognitive comme une logique non monotone

La définition de l'architecture cognitive proposée s'est fondée sur la capacité élémentaire à s'apercevoir a posteriori d'une illusion ou d'une hallucination relative soit à la conception soit à la perception. Autrement dit, elle est basée sur la capacité à effectuer un raisonnement non monotone dans le sens où les conclusions des jugements passés ont été rétractées par l'ajout de nouvelles informations ; ainsi le nombre de conclusions considérées comme vraies ne s'accroît pas de manière monotone. À l'inverse, un raisonnement monotone correspond à l'invariabilité de la véracité des conclusions malgré l'ajout de nouvelles informations. Dans le premier chapitre, la présentation du symbolisme a montré les limites des systèmes monotones.

Un système logique monotone implique qu'un théorème, dans le cadre d'une théorie formelle définie, reste un théorème après l'ajout de nouveaux axiomes à la théorie formelle originale ; l'accroissement du nombre de théorème est monotone. Par exemple, la logique classique représente un système formel monotone. En notant A une formule, b et c des faits, la monotonie pour un système de logique classique s'exprime de la manière suivante :

$$A \vdash b \Rightarrow A \cup c \vdash b$$

Dans un système logique non monotone, une formule peut perdre son statut de théorème après l'ajout de nouvelles informations, ce qui signifie que l'ensemble des conclusions du système formel n'augmente pas de manière monotone avec la croissance du nombre des axiomes. La logique non monotone permet des inférences à partir de théories formelles incomplètes mais, de fait, cela entraîne une incertitude sur l'inférence. Cette incertitude de l'inférence peut être utilisée dans le cadre d'une logique floue, ce qui sera l'objet du quatrième point concernant les remarques logiques.

Les systèmes logiques non monotones (Brewka *et al.*, 1997) prennent en compte en général cette incertitude en adoptant l'hypothèse (A) d'un monde clos et (B) d'une hiérarchisation par défaut des règles. La présentation successive de ces deux techniques montera qu'en définitive elle répond à des contraintes similaires à celles évoquées lors de la présentation des hypothèses cognitives sur la perception et la conception dans le premier chapitre, l'indubitabilité et l'univocité. (C) La traduction de ces techniques dans l'architecture cognitive introduit la notion de négation qui obligera à revenir sur la notion de vérité afin de bien cerner les implications philosophiques de la forme de cette négation. (D) Cette clarification permettra alors de formaliser cette négation et ainsi de compléter l'architecture cognitive.

A - L'hypothèse d'un monde clos

Le premier problème repose sur l'interprétation de l'ignorance, c'est-à-dire sur la manière dont est interprétée la situation où un fait B ne peut pas être déduit à partir des axiomes considérés. La méthode la plus répandue pour répondre à cette question consiste à supposer que le monde à partir duquel s'effectuent les inférences est clos (Reiter, 1978). Autrement dit, il est considéré par défaut qu'une inférence s'effectue à partir d'un ensemble d'axiomes complets (en intelligence artificielle classique, un axiome est synonyme de fait

avéré). L'hypothèse du monde clos permet alors la relation d'inférence suivante, si B ne peut pas être inféré à partir de A , alors par défaut A infère non B :

$$\neg(A \vdash B) \models (A \vdash \neg B)$$

La négation équivaut ici par défaut à l'absence d'une formule dans la base de connaissance. Par exemple, pour vérifier que Paul n'est pas de l'entreprise A , il n'est pas nécessaire que la base client de cette entreprise contienne l'information « Paul n'est pas un client ». Le simple fait de ne pas le trouver dans la base signifie que Paul n'est pas client. La négation est par défaut puisqu'elle n'est pas explicitée dans la base ; par ailleurs, la conclusion est révisable puisque Paul peut devenir un jour client de l'entreprise. Cependant, si la base de connaissance contient une formule disjonctive comme $A \vee B$, la relation d'inférence conduit à l'inconsistance, ce qui oblige à restreindre les formules de cette forme dans la base de connaissances.

Concernant l'architecture cognitive, celle-ci possède les caractéristiques d'une logique supposant un monde clos, d'une part du fait de l'existence d'une mémoire épistémique qui infère uniquement avec les signes contenus dans la mémoire événementielle et d'autre part du fait que la mémoire ne contient que des signes. Mais la justification de ces caractéristiques diffère de celle des logiques non monotones adoptant l'hypothèse du monde clos. La réflexion sur la perception a conduit à différencier l'incertitude d'un signe et celle de sa justification. Plus exactement, l'indubitabilité des signes signifie que l'incertitude de la présence d'un signe au sein de la mémoire événementielle est nulle, ce qui conduit en définitive à la supposition d'un monde clos. La supposition d'un monde clos permet un raisonnement sur l'ignorance mais qui conduit à rejeter la disjonction dans la base de connaissances sous peine d'inconsistance, ce qui revient à considérer les informations comme connues et leurs présences certaines, soit à l'indubitabilité de la base de la connaissance.

Cependant, l'architecture cognitive n'intègre pas la relation d'inférence décrite par les systèmes logiques non monotones supposant un monde clos. En effet, celle-ci introduit une ambiguïté dans la définition de la négation dans le sens où il n'y a pas de différence entre l'absence de Paul dans la base client de l'entreprise et la présence dans cette même base de l'explication que Paul ne soit pas client. Pour approfondir la question de la négation, la présentation du second problème doit être entreprise.

B - La hiérarchisation par défaut des règles

Le second problème provient du fait qu'à partir d'une même base d'informations des implications contradictoires peuvent s'appliquer, en notant \vdash l'inférence d'une logique non monotone et \neg la négation, le problème de l'inférence s'écrit de la manière suivante :

$$(a \supset \neg c) \cup (a \wedge b \supset c) \cup \{a \cup b\} \vdash ? \quad (1)$$

Par exemple, « le jeudi, Pierre va à la piscine » et « le jeudi quand il pleut, Pierre reste chez lui », en acceptant que l'ubiquité ne soit pas possible, un jeudi pluvieux, Pierre peut logiquement soit être chez lui soit être à la piscine et un mécanisme décisionnel doit être rajouté pour lever l'incertitude.

Dans ce cadre, Reiter (1980) propose la logique non monotone des défauts qui est une logique de premier ordre. Cette logique s'appuie sur un type de règles supplémentaires : les règles par défaut s'appliquent à la plupart des individus mais pas à tous. En d'autres termes, une proposition P est considérée comme vraie tant qu'une information

supplémentaire ne prouve pas le contraire. Ce raisonnement repose sur l'hypothèse d'un monde clos puisque la justification d'une règle par défaut revient à vérifier qu'aucune exception ne peut être dérivée afin d'affirmer qu'a priori la règle par défaut est applicable. Ainsi, les inférences s'effectuent en appliquant autant de règles par défaut que possible.

En définitive, la logique non monotone des défauts utilise la négation par défaut qu'offre un monde clos dans les prémisses des règles afin de hiérarchiser l'application de celles-ci. Mais pour éviter la confusion entre une forme négative d'une proposition et son absence, les conditions d'application des règles se trouvent à côté des prémisses. Sans rentrer dans le détail du formalisme de cette logique, les deux exemples suivant permettent de saisir le fonctionnement de la logique des défauts :

« Si l'individu x est un oiseau alors il vole, s'il est capable de voler »

$$(\forall x) \text{oiseau}(x) \supset \text{vole}(x) : (\text{vole})$$

« Si l'individu x est un oiseau alors il vole, si ce n'est ni un pingouin et ni une autruche »

$$(\forall x) \text{oiseau}(x) \supset \text{vole}(x) : (\neg \text{pingouin}(x) \wedge \neg \text{autruche}(x))$$

Une autre manière de gérer l'ambiguïté entre la forme négative d'une proposition et son absence communément employée réside dans l'utilisation d'une logique non monotone de circonscription (McCarthy, 1986). Cette logique du premier ou deuxième ordre selon les variantes (Brewka *et al.*, 1997) applique la supposition d'un monde clos uniquement à quelques prédicats. Par exemple, en posant le prédicat d'anormalité an circonscrit de façon à assumer $\neg an(x)$ sauf si $an(x)$ est vrai :

$$(\forall x) (\text{oiseau}(x) \wedge \neg an(x) \supset \text{vole}(x)) \cup \text{oiseau}(x) \vdash \text{vole}(x)$$

Concernant l'architecture cognitive simplifiée, celle-ci répond partiellement au second problème. L'univocité de la perception ou de la conception a conduit à la notion de champ, ici, elle s'assimile à la conservation de la consistance d'un ensemble de règles. Cependant, dans le cadre de l'architecture cognitive, la vérification de la consistance ne repose pas sur un enchaînement d'inférences mais sur une compétition directe entre le contenu des règles dont le mécanisme sera présenté au cours de la spécification. Par ailleurs, les règles au sein de la mémoire épistémique ne doivent pas être structurées de façon à posséder une hiérarchisation a priori destinée à effectuer un choix par défaut car cela signifierait qu'il existe une structuration a priori du monde. Les schèmes cognitifs peuvent présenter ce genre d'organisation contrairement à l'architecture cognitive qui doit rester neutre. Ainsi, le choix entre les implications (1) doit résulter de l'auto-organisation produite par la compétition entre les règles.

C - L'architecture cognitive et la notion de négation

Toutefois, sans autre élément, l'architecture cognitive simplifiée formalisée ne peut pas gérer les clauses restrictives comme la logique des défauts ou la logique de circonscription. Pour cela, la notion de négation doit être introduite. La difficulté des deux précédentes logiques à intégrer le concept de négation montre que cette notion est complexe. Mais cette complexité ne provient pas de l'hypothèse d'un monde clos, elle ne fait que révéler la dualité de la notion de négation logique qui reflète en définitive la dualité de la notion philosophique de la vérité abordée dans le chapitre précédent. La notion classique de vérité peut se comprendre soit comme la vérité existentielle (ou vérité synthétique) : est vrai ce qui est, soit comme la vérité relationnelle (ou vérité analytique). Toutes les deux ont été étudiées dans le second chapitre au cours de la présentation de la

vérité comme une mise en correspondance. Ces deux conceptions renvoient traditionnellement à l'adoption de HOI.

Dans le cadre des relations logiques, celles-ci étant avant tout descriptives ne réclament pas de débat philosophique sur la vérité mais, dans le cadre des relations effectives, celles-ci doivent construire avec les propriétés logiques quitte à décrire a posteriori ces relations avec les logiques appropriées. Par conséquent il est nécessaire de s'assurer qu'aucune hypothèse ontologique ne s'insinue dans les principes de bases constituant l'autopoïèse sémiotique. Ainsi, dans un premier temps, la conception de vérité existentielle sera dépourvue de sa charge métaphysique afin de l'intégrer dans l'architecture cognitive et de ce fait servir à compléter la non monotonie de l'architecture cognitive. Dans un second temps, la notion de contradiction se situant au cœur des systèmes logiques, l'analyse de la négation de la vérité relationnelle sera l'objet du troisième point des remarques logiques.

Afin de s'assurer que la notion de vérité existentielle se décharge de sa connotation métaphysique dans le cadre de la définition de l'architecture cognitive, les principes sous-tendant l'architecture cognitive proposée doivent être rappelés. La cognition se fonde sur la segmentation a priori du monde, cependant, la nécessité de segmenter le monde pour initier la cognition n'implique pas que le monde soit ontologiquement segmenté, autrement dit n'implique pas l'adoption de HOI ou HOP. De plus, l'utilisation de cette segmentation dans le cadre d'une activité cognitive repose nécessairement sur l'interaction de deux mémoires : la mémoire événementielle contenant les segments considérés (les signes) et la mémoire épistémique traitant ces segments (les interprétations). Les signes proviennent soit des interprétations, soit des primitives sensorielles et les interprétations de la mémoire épistémique se fondent uniquement sur les signes contenus dans la mémoire événementielle. En d'autres termes, la mémoire événementielle peut être considérée comme un monde clos. La notion de vérité existentielle pour un signe se traduit par la notion de la présence effective de celui-ci dans la mémoire événementielle. L'expression de la négation de la présence effective est uniquement l'absence du signe, c'est-à-dire que l'utilisation d'un opérateur de négation \neg devient non valide pour stipuler, par exemple, que *B* n'est pas présent. Appliquer ainsi la notion de présence effective n'introduit pas d'hypothèse ontologique puisqu'elle repose sur la définition complète et parfaite d'un monde construit, la mémoire événementielle.

Par ailleurs, il ne faut pas confondre la présence effective d'un signe qui porte sur sa présence dans la mémoire événementielle et le jugement existentiel de la perception. En définitive, la présence effective d'un signe correspond à un jugement perceptif d'un signe et le jugement existentiel correspond à la certitude de l'inférence qui a conduit à injecter ce signe dans la mémoire événementielle. Le jugement perceptif d'un objet traduit l'activité d'un réseau de signes et d'interprétations et le jugement existentiel représente la légitimité de leur activation. Le quatrième point traitera plus particulièrement de la traduction d'un jugement existentiel d'un signe.

D - La négation comme une inhibition

Concrètement, pour l'architecture cognitive proposée, cela signifie qu'il n'y a pas d'opérateur exprimant la négation existentielle et par conséquent, seuls les signes présents affectent la mémoire événementielle. La gestion de la restriction dans la non monotonie de l'architecture cognitive se trouve alors dans la manière d'interpréter les signes présents, c'est-à-dire dans la structure des règles afin de spécifier que la présence de certains signes inhibe la règle. Plus précisément, les règles sont des implications possédant une prémisse

excitatrice et une prémisses inhibitrice. La prémisses inhibitrice n'est pas une condition sur la règle mais fait partie de son expression. D'un point de vue existentiel, les clauses sont toujours positives par rapport à l'exemple de la logique des défauts, « si l'individu x est un oiseau alors il vole, si ce n'est ni un pingouin et ni une autruche » devient « si l'individu x est un oiseau alors il vole, sauf si c'est un pingouin ou une autruche ».

Cette démarche est proche de celle de la logique de circonscription sauf qu'ici la logique est d'ordre zéro et que la vérité existentielle n'est pas circonscrite à un prédicat mais intégrée en tant que prémisses inhibitrice dans la règle. Le symbole \neg est traditionnellement employé pour exprimer la négation relative soit à l'existence, soit à la relation. Cependant, ce symbole n'a pas été retenu pour désigner la prémisses inhibitrice afin d'éviter toute ambiguïté. Le langage utilisé pour formaliser l'architecture se complète alors avec deux nouveaux symboles :

4. « $\langle \rangle_I$ » représente la prémisses inhibitrice contenant les faits inhibiteurs disjonctifs
5. « $\langle \rangle_E$ » représente la prémisses excitatrice contenant les faits excitateurs conjonctifs

Avant d'expliquer l'intérêt de limiter la prémisses inhibitrice aux formules disjonctives et de limiter la prémisses excitatrice aux formules conjonctives, le symbole de la disjonction doit être incorporé au langage des schèmes cognitifs :

6. « \vee » représente la disjonction

Une règle de la mémoire épistémique possède la structure suivante :

$$\langle \vee_k M_k \rangle_I \wedge \langle \wedge_l M_l \rangle_E \supset M_n$$

Si M_l est présent, il implique M_n sauf si M_k est également présent. La limitation sur l'emploi de la disjonction et de la conjonction ne nuit pas à l'expressivité des règles puisque, dans le cadre de la logique de l'architecture cognitive, il y a équivalence entre :

$$\left. \begin{array}{l} \langle \rangle_I \wedge \langle A \rangle_E \supset C \\ \langle \rangle_I \wedge \langle B \rangle_E \supset C \end{array} \right\} \equiv \langle \rangle_I \wedge \langle A \vee B \rangle_E \supset C$$

$$\left. \begin{array}{l} \langle B \rangle_I \wedge \langle A \rangle_E \supset C \\ \langle D \rangle_I \wedge \langle A \rangle_E \supset C \end{array} \right\} \equiv \langle B \wedge D \rangle_I \wedge \langle A \rangle_E \supset C$$

La limitation sur l'emploi de la disjonction et de la conjonction est justifiée par le fait que les règles doivent représenter une croyance minimale afin qu'une règle complexe se décompose en plusieurs et que chacune d'entre elles puisse être évaluée séparément dans le cadre de la compétition entre les croyances. De plus, l'atomisation des interprétations facilite la génération progressive d'un ensemble cohérent de règles.

**

En somme, l'architecture cognitive peut être comprise comme un système de logique non monotone qui intègre la négation existentielle comme un mécanisme d'inhibition et qui impose une représentation minimale des implications. Toutefois, il reste à examiner comment l'architecture cognitive intègre la négation relationnelle, soit la contradiction.

3.3. L'analyse de la contradiction et de son intégration

Il existe différentes manières d'aborder la contradiction et le choix dépend en définitive de la notion de vérité adoptée. Afin de bien comprendre la nature de cette relation entre le choix d'un système formel et la conception de la vérité, (A) un résumé d'une partie des discussions du second chapitre doit être présenté. Ensuite, (B) les trois principales interprétations de la contradiction dans les systèmes formels seront successivement présentées afin de choisir celle offrant une charge métaphysique nulle et des pistes pour l'intégration de la notion de contradiction dans l'architecture cognitive.

A - Relation entre le choix d'un système formel et la conception de la vérité

Les systèmes formels présentent l'avantage de rejeter le problème de la signification de la vérité en la définissant parfaitement au sein des systèmes formels qui représentent des mondes construits. En d'autres termes, les systèmes formels permettent de séparer la sémantique et la logique. Cette séparation est pratique pour l'étude de la logique mais peut produire des dilemmes lorsque cette séparation est remise en question comme le fait le symbolisme en souhaitant modéliser la cognition en réunissant la sémantique et les systèmes formels. En effet, cette tentative révèle des tensions intenable comme cela a été montré dans le premier chapitre.

Les systèmes formels définissent l'attribution de la valeur de vérité d'une proposition. De ce fait, une formule qualifiée de vraie n'a pas de sens au-delà de son interprétation au sein d'un système formel. En reprenant l'exemple précédent avec les axiomes $A \supset C$ et A et la règle d'inférence définie (le *modus ponens*), C est vrai uniquement dans ce cadre. Toutefois, quelles valeurs de vérité attribuer aux règles d'inférences qui définissent la vérité des propositions que ces règles d'inférences manipulent ? En effet, des règles d'inférences peuvent être définies dans un système formel sans être valides. Par exemple, en rajoutant les symboles suivants, il devient possible de décrire des systèmes formels du premier ordre :

- a) « $P(x)$ » se lit « x est en relation avec P »
- b) « $\exists x$ » se lit « il existe un x tel que » et indique certaines entités x (au moins une) satisfont une condition ; ainsi « $\exists x P(x)$ » signifie « il existe une entité x qui soit P ».

La règle d'inférence « Si tous les A sont B et tous les B sont C alors tous les A sont C » est valide et s'écrit de la façon suivante :

$$(\forall x (A(x) \supset B(x)), \forall x (B(x) \supset C(x)) \vdash \forall x (A(x) \supset C(x)))$$

En revanche, la règle d'inférence « si certains A sont B , alors tous les A sont B » n'est pas valide dans la mesure où des contre exemples existent :

$$(\exists x A(x) \supset B(x)) \vdash (\forall x (A(x) \supset B(x)))$$

La question de leur validité peut être comprise comme une analyse d'un ordre supérieur, autrement dit, le métalangage utilisé pour décrire le système formel devient le langage décrit par un système formel d'ordre supérieur. La théorie de modèles, en utilisant la logique modale (Kripke, 1959), étudie le domaine de validité des règles d'inférence d'un système. Ainsi, avec une logique modale de type théorie des mondes possibles, la règle d'inférence de *modus ponens* est vraie quel que soit le monde défini au sein de la logique

modale. Cependant, cette vérité demeure relative au système de la logique modale et au langage qu'elle définit. Ici, le problème majeur vient de l'interprétation de ce résultat. Sans revenir sur l'étude des différents arguments effectuée au second chapitre, deux conceptions (avec leurs variantes) s'opposent traditionnellement : le nominalisme qui refuse HOI c'est-à-dire appartenant soit à FA soit à FP, et le réalisme qui accepte HOI c'est-à-dire appartenant soit à FCP soit à FC : soit la logique modale a permis de retrouver une vérité absolue qui correspond à l'intuition qu'effectivement « Si tous les A sont B et tous les B sont C alors tous les A sont C », soit la puissance de cette règle d'inférence est intéressante pour une classe de système formel mais demeure relative au langage qui la décrit comme toute chose en logique.

Deux domaines d'études distincts se dégagent alors : la philosophie de la logique qui s'interroge sur la signification de la vérité des systèmes formels et la logique mathématique qui étudie les propriétés des systèmes formels à attribuer des valeurs de vérité. Toutefois, la volonté de décrire le monde avec les systèmes formels réintroduit le problème de la signification de la vérité. La recherche en logique mathématique peut s'effectuer indépendamment de toute considération philosophique sur la signification de la vérité, mais, dès que la recherche en intelligence artificielle souhaite se servir de ces travaux sur la logique pour décrire le raisonnement et le monde, le problème apparaît et les convictions de chacun influencent le choix du système logique employé. Selon la famille philosophique adoptée, certaines propriétés vont être directement mises au niveau du système formel alors que d'autres estiment que ces propriétés sont des hypothèses trop fortes ou inutiles pour que le langage décrit par le système formel puisse décrire le monde, bien que des équivalences sur le plan formel existent.

B - Les trois interprétations de la négation

La gestion de la négation directement liée à la notion de vérité est caractéristique de la dépendance des convictions philosophiques au choix du système logique utilisé pour décrire le monde. Les différentes interprétations de la négation peuvent être présentées aux travers principalement de trois logiques : (i) la logique classique, (ii) la logique intuitionniste et (iii) la logique minimale. Bien que l'architecture cognitive proposée se situe dans le cadre de la FA, afin de comprendre toutes les implications de la logique adoptée qui servira de modèle pour la définition de l'architecture cognitive, les trois types de logique seront présentés successivement. L'emploi de la première, la logique classique, a tendance à être défendu par les approches issues de la FCP. La seconde, la logique intuitionniste, est promue principalement par la FC. La dernière, la logique minimale, se révèle compatible à la fois avec FP et avec FA. Dans ces trois logiques, la contradiction, notée \perp , se définit comme la réunion d'une proposition A et de sa négation notée $\neg A$:

$$(\neg A \wedge A) \vdash \perp$$

De même, si une proposition A conduit à une contradiction alors $\neg A$, en d'autres termes :

$$A \supset \perp \vdash \neg A$$

i - La contradiction dans la logique classique

La logique classique, en se fondant sur le principe du tiers exclus (A ou non- A) coïncide avec une conception de la vérité comme une mise en correspondance soutenue par la FCP : A correspond avec un fait « a » ou A ne correspond pas avec un fait « a ».

Ainsi, démontrer que non- A est faux revient à démontrer que A est vrai puisque la définition ontologique du vrai est symétrique à celle du faux. En d'autres termes, le raisonnement par l'absurde est valide :

$$\neg A \supset \perp \vdash A$$

Soit,

$$\neg\neg A \supset A$$

Mais, comme il a été souligné dans le second chapitre, de nombreuses critiques existent à l'égard de la vérité par correspondance et par conséquent à l'emploi de la logique classique pour décrire le monde. Russel (Engel, 1998), par exemple, souligne la difficulté de concevoir la notion de correspondance avec les faits négatifs comme « mon verre n'est pas vide », qui concrètement pose problème à la réalisation d'un système en robotique symbolique puisque la définition de tous les faits est impossible dans un monde ouvert.

ii - La contradiction dans la logique intuitionniste

La FC rejette par définition la vérité comme mise en correspondance et soutient que la vérité correspond à la cohérence ultime d'un unique édifice de savoir qui ne peut se construire que progressivement. Dans ce cadre, l'intuitionnisme (Brouwer, 1927) considère que les êtres mathématiques coïncident avec la possibilité de leur construction mentale. De ce fait, la logique intuitionniste refuse la justification ontologique du tiers exclus qui définit indirectement l'existence d'un objet mathématique. En logique intuitionniste, la formule suivante est valide :

$$A \supset \neg\neg A$$

Alors que la formule réciproque ne l'est pas, contrairement à la logique classique qui accepte les deux :

$$\neg\neg A \supset A$$

En revanche, sans définir A , la formule suivante reste valide en logique intuitionniste :

$$\neg\neg\neg A \supset \neg A$$

Par ailleurs, la vérité provenant uniquement de la cohérence du système, la cohérence domine la sémantique. Ainsi, l'intuitionnisme accepte qu'à partir d'une contradiction puisse se déduire n'importe quelle proposition (*ex falso sequitur quodlibet*) : $\perp \supset B$. En d'autres termes, une conditionnelle contrefactuelle est forcément vraie puisque la prémisse est fautive. Par exemple, il est vrai que « si les poules ont des dents alors la tour Eiffel se trouve à Berlin ». Dans le cadre de la logique intuitionniste, cette proposition a autant de valeur que « La tour Eiffel se trouve à Paris ».

iii - La contradiction dans la logique minimale

La logique classique et la logique intuitionniste encouragées par une position ontologique de la vérité logique introduisent directement dans leur système une interprétation fondamentale de la contradiction. Or, l'approche pragmatiste dans laquelle s'inscrit l'architecture cognitive proposée refuse HOI et défend par conséquent un conventionnalisme formel. Le conventionnalisme formel considère que la description du monde ne nécessite a priori qu'un système logique offrant les moyens d'effectuer des

relations entre des faits et que la contradiction vient de la combinaison de ces relations en dehors de toute considération métaphysique. Dans ce cadre, la logique minimale se révèle particulièrement adaptée puisqu'elle ne déduit aucune conséquence de la contraction, comme le tiers exclus pour la logique classique ou *l'ex falso sequitur quodlibet* pour la logique intuitionniste. La contradiction est une convention s'apparentant à n'importe quelle autre formule, en posant la formule C et en définissant $\sim A$ comme synonyme de $A \supset C$ alors :

$$(A \cup \sim A) \vdash C$$

$$A \supset C \vdash \sim A$$

Ainsi, aucun rôle particulier n'est conféré à la contradiction. Les symboles \neg et \perp deviennent des commodités d'écriture. Pour ces raisons, en s'inspirant de la logique minimale, l'architecture cognitive ne gère pas directement les contradictions, leur gestion est reléguée aux schèmes cognitifs, c'est-à-dire à la théorie formelle du système logique. L'architecture cognitive ne nécessite donc pas d'opérateur de négation. La négation relationnelle doit apparaître lors de la construction des règles et de ce fait est soumise à la mise à l'épreuve de sorte qu'aucune relation ne se trouve privilégiée par les principes logiques sous-tendant l'architecture cognitive.

3.4. Gestion de l'incertitude des inférences

En somme, la négation existentielle ou la négation relationnelle ne s'expriment pas par l'application d'un opérateur sur les signes. Cette conclusion ne clôt pas pour autant la discussion sur la vérité sur le plan logique. En effet, toute inférence en logique non monotone est incertaine puisque révisable par de nouvelles, autrement dit toute conclusion possède une certaine légitimité à être. Le terme légitimité a été préféré à celui de crédibilité car il exprime davantage l'idée de bienfondé. L'incertitude de l'inférence porte sur la légitimité du signe et non sur sa présence. Cette préoccupation n'apparaît pas, par exemple, en logique classique puisque toute déduction est certaine. Une manière d'intégrer cette incertitude consiste à attribuer à la conclusion une valeur de légitimité comprise entre 0 et 1. Le détail de l'évaluation de cette légitimité sera présenté lors de la spécification de l'architecture cognitive. Cependant, le formalisme peut déjà être complété en rajoutant un nouvel attribut aux signes :

1. « M_a » se lit « message A avec le contenu a associé à une durée temporelle θ_a et à une valeur de légitimité L_a »

L'apparition de la valeur de la légitimité dans les signes permet alors d'entrevoir l'utilisation de celle-ci afin de combiner et de suivre les conclusions successives. La polyvalence et la combinaison des conclusions en fonction de celles-ci renvoient aux principales particularités des modèles de probabilités imprécises telles que la logique floue (Zadeh, 1978) ou les probabilités subjectives (Savage, 1954). Une comparaison plus approfondie avec ces approches nécessite cependant une définition plus précise des signes et des règles, ce qui sera l'objet de la spécification.

En reprenant les termes utilisés pour l'étude de la perception et de la conception dans le cadre de l'architecture cognitive proposée, le jugement perceptif ou le jugement d'assertabilité d'un signe s'identifie à la vérité existentielle d'un signe qui se résume simplement à la présence de celui-ci dans la mémoire épistémique. La différence entre le jugement perceptif et le jugement d'assertabilité provient de la nature du signe qui sera

développée dans le cinquième point des remarques logiques. Ainsi, le jugement existentiel ou le jugement de vérité d'un signe traduit la légitimité du signe à être présent au sein de la mémoire événementielle. De même, la différence entre le jugement existentiel et le jugement de vérité d'un signe provient de la nature du signe concerné.

Bien que les jugements issus de la conception et de la perception semblent parfaitement symétriques, ils gardent chacun leur spécificité, même en restreignant la réflexion au niveau du signe et non au niveau de l'objet (conceptuel ou physique). En effet, le jugement existentiel d'un signe provient directement d'une primitive sensorielle dont la légitimité est évidente puisqu'elle ne provient pas d'une inférence. Par conséquent, la légitimité d'un signe provenant d'une primitive sensorielle est toujours maximale. À l'inverse, le jugement de vérité d'un signe dépend toujours de l'inférence qui l'a produit. Plus exactement, la légitimité de la conclusion se calcule en fonction de l'appariement entre la prémisse et la mémoire événementielle et de la compétition avec les autres croyances. Une conclusion issue d'une interprétation sans concurrence aura une légitimité élevée (proche de 1), à l'inverse une conclusion provenant d'une interprétation disputée obtiendra une légitimité faible (proche de 0).

Avant d'étudier l'intérêt de la définition de la légitimité d'un signe dans la mémoire événementielle, la réflexion sur la légitimité d'un signe peut être étendue à la légitimité d'une règle. Le problème de la négation des règles ne se pose pas puisqu'elles n'apparaissent pas en tant que tel dans les prémisses. Par ailleurs, de la même manière que pour les signes, la vérité existentielle d'une règle est indubitable et correspond à sa présence au sein de la mémoire épistémique. Mais la notion de légitimité d'un signe est-elle transposable à une règle ? Précédemment, il a été montré que chaque règle se trouve rattachée à une évaluation ponctuelle et une évaluation globale. L'évaluation ponctuelle représente l'adéquation entre la prémisse de la règle et la mémoire événementielle en fonction des prémisses des autres règles, cette évaluation est transmise à la conclusion qui devient la légitimité du signe. L'évaluation globale d'une règle doit s'effectuer au cours de la sémiose en fonction de la manière dont elle s'organise avec les autres et son apport aux rétributions internes. Pour nommer cette évaluation, le terme pertinence correspond davantage à légitimité. Bien que lors de l'arrivée d'une nouvelle règle, l'évaluation globale puisse correspondre dans un premier temps à la légitimité du déclenchement de la règle génératrice, l'évolution de l'évaluation globale dépend de l'apport qu'elle peut avoir au cours du vécu de l'individu auquel elle appartient. Ici, la notion de pertinence d'une règle diffère de beaucoup de la notion de vérité d'une règle dans le sens où elle correspondrait à la probabilité que la règle reflète une description juste du monde, propre aux approches issues des familles FC, FCP ou FP. En effet, le pragmatisme qui appartient à la famille FA conduit à considérer uniquement la rétribution interne pour construire une représentation du monde tout en se fondant sur l'exploration des interactions. La notion de pertinence traduit cette relation de dépendance entre l'activité de la sémiose et les interactions avec le monde qu'elle produit et de ce fait, la notion de pertinence exprime le caractère situé physiquement et cognitivement de l'individu. Dans le cadre de la formalisation de l'architecture cognitive, la prise en compte de la pertinence d'une règle comme son évaluation globale revient à préciser le point 3 du langage des schèmes cognitifs comme suit :

3. « \triangleright_i » représente l'implication assimilée à une règle (ou une croyance) de la mémoire épistémique à laquelle est associée une évaluation globale traduisant la pertinence de la règle $i : P_i$

La définition de l'évaluation globale comme la pertinence conforte l'utilisation de ce critère pour l'élimination des règles (soit l'oubli des croyances) comme il a été évoqué précédemment (2.2 B). Maintenant, il reste à comprendre l'intérêt de la notion de légitimité dans la sémiologie. Pour cela, la notion de légitimité doit être mise en perspective dans les trois premiers stades de la généalogie de la cognition. Au premier stade, les schèmes cognitifs correspondent à des règles stimulus-réponse dont les prémisses se modulent en fonction du vécu de l'organisme. D'un point de vue logique, l'architecture cognitive ne fait que déduire son comportement. La question de l'utilité de la légitimité se pose uniquement pour les comportements choisis puisque les stimuli sont de fait tous légitimes, leur légitimité est identique. Il s'offre alors deux possibilités : la première serait de considérer que la légitimité du comportement n'importe pas dans la réalisation de celle-ci, la seconde possibilité serait que la légitimité module le comportement choisi. Cependant, la première solution n'explique pas pourquoi dans les exemples de Lorenz évoqués, dans le premier chapitre, lors de la présentation du constructivisme varélien, il existerait des variations dans l'intensité des comportements. En revanche, la seconde solution offre des pistes qui seront davantage détaillées lors de la spécification de l'architecture cognitive pour modéliser ce genre de réaction. En définitive, à ce premier stade de la généalogie de la cognition, le rôle de la légitimité et l'idée de la modulation des comportements relèvent davantage des hypothèses sur les primitives motrices que sur l'architecture cognitive.

Au deuxième stade, l'alternance entre instinct et dressage apparaît. Le dressage se comprend ici comme un mécanisme d'auto-rétribution permettant l'évaluation des comportements. Ce mécanisme permet de générer de nouvelles relations, c'est-à-dire de nouvelles règles comportementales. Une manière simple et sûre de les générer repose sur le principe de l'induction. Au sein de l'architecture cognitive, l'induction se comprend comme un schème cognitif générateur en sélectionnant des signes contenus dans la mémoire événementielle. La signification de schème cognitif sera approfondie dans le sixième point de la formalisation et son mécanisme sera détaillé dans la quatrième partie de ce chapitre. Néanmoins, cette génération de règles au deuxième stade suggère que la légitimité des signes peut jouer un rôle dans le mécanisme de sélection de ces derniers lors de la confection de nouvelles règles.

Au troisième stade, le problème de la prise de décision apparaît. En effet, l'auto-rétribution devient un schème cognitif, c'est-à-dire que le processus de rétribution devient apparent pour la sémiologie. Autrement dit, l'individu peut être amené à choisir son comportement en fonction de plusieurs rétributions envisagées. La multiplicité des solutions renvoie au problème évoqué dans le premier chapitre sur la fonction d'utilité. Le choix nécessite l'évaluation des possibles, soit une simulation interne évaluant les conséquences des actions. L'évaluation des possibles doit intégrer la légitimité des signes dans l'évaluation des inférences afin de propager la légitimité qui, dans le cadre d'une simulation interne, devient la crédibilité (ou probabilité) d'un signe envisagé. La propagation de la crédibilité représente la propagation de la valeur de vérité d'un fait en logique floue. Ce type de schèmes se révèle complexe mais il ne nécessite pas la notion d'objet bien qu'elle puisse être intégrée.

En définitive, la légitimité peut être considérée facultative dans une certaine mesure pour les deux premiers stades. Mais, au troisième stade, elle apparaît indispensable pour la prise de décision, prenant en compte l'incertain en utilisant la légitimité comme une valeur de crédibilité se propageant dans le cas de la simulation. Par conséquent, il est préférable d'intégrer directement la notion de la légitimité aux signes dans la définition de

l'architecture cognitive, tout en laissant l'exploitation de cet attribut au niveau des schèmes cognitifs et des primitives motrices.

3.5. Les deux niveaux de typage

Le cinquième point sur les implications logiques de l'architecture cognitive aborde la notion de type. L'introduction de la notion de type est motivée d'une part par (A) des raisons de bas niveau liées à l'architecture cognitive et d'autre part par (B) des raisons de haut niveau liées aux capacités cognitives attribuées aux individus appartenant au troisième ou quatrième stade de la généalogie de la cognition. Afin de comprendre dans quelles mesures le typage de haut niveau produit par la sémiologie repose sur le typage de bas niveau indispensable pour refléter les propriétés de l'architecture cognitive, ces deux niveaux de typage seront successivement présentés.

A - Le typage au niveau de l'architecture cognitive

La nécessité d'un typage de bas niveau provient de la notion de champ dans le formalisme de l'architecture cognitive. Un champ représente l'ensemble des signes qui ne peuvent être présents deux à deux simultanément dans la mémoire événementielle ; par exemple, deux signes traduisant chacun une action sur un même effecteur ne peuvent être introduits en même temps dans la mémoire événementielle car cela remettrait en cause son intégrité, soit son indubitabilité.

Une manière de formaliser la notion de champ spécifique consiste à définir un type de message pour chaque primitive motrice et chaque schème cognitif spécifique. Un type représente un ensemble de données (de structures) et éventuellement de fonctions qu'il est possible de manipuler de manière uniforme à travers un certain nombre d'opérations spécifiques à cet ensemble. La compétition pour l'interprétation s'effectue alors au sein de règles dont le type de la conclusion se trouve identique. Ainsi, un champ spécifique n représentant l'ensemble des implications dont la conclusion est un message de type n , la formalisation de l'inférence pour l'architecture cognitive proposée devient la suivante :

- d. « \vdash^n » se lit « si...en fonction de l'ensemble des prémisses des autres règles impliquant une conclusion de type n appartenant à la mémoire épistémique alors...est mis dans la mémoire événementielle avec une légitimité L »

Il apparaît alors une ressemblance conceptuelle entre l'ensemble des règles ayant une conclusion de type identique à une fonction¹. En considérant l'état de la mémoire événementielle comme un élément de l'ensemble contenant tous les états possibles, la sélection dans la mémoire épistémique d'une conclusion dans l'ensemble des règles d'un type donné correspond au calcul d'une fonction. Dans ce cas, la mémoire épistémique

¹Une fonction au sens mathématique est une correspondance d'un ensemble E vers un ensemble F, qui à tout élément de E associe au plus un élément de F. Une application mathématique représente une opération consistant à mettre en correspondance tout élément d'un ensemble A avec un élément d'un ensemble B et un seul.

deviendrait la somme des images extensionnelles de chacune des fonctions associées à un type de conclusion. En notant M la mémoire événementielle et E^n l'ensemble de règles avec une conclusion de type n , une interprétation s'écrit alors :

$$E^n(M(t)) = C^n \Rightarrow M(t+1) = M(t) + C^n$$

Le découpage de la mémoire épistémique en termes de fonction se trouve définie ici par le typage des conclusions mais ce découpage peut encore s'affiner en considérant qu'une fonction représente l'ensemble des règles de prémisses et de conclusions de nature identique. Une règle de la mémoire épistémique représente alors une relation orientée et restrictive. Par ailleurs, cette décomposition fonctionnelle implique alors le typage également des messages issus des primitives sensorielles.

Dans ce cas, la définition formelle d'un schème cognitif devient la suivante : un schème cognitif constitue un ensemble de règles possédant les mêmes types de conditions et de conclusions. Un schème cognitif K regroupant l'ensemble des règles avec des prémisses représentant un ensemble d'ensembles de messages classés selon leur type et avec une conclusion de type m s'écrit alors comme suit :

$$\left\{ \left\langle \right\rangle_I \wedge \left\langle \left\{ \left\{ \wedge_i M_i \right\}^n \right\} \right\rangle_E \supset M_j^m \right\}_K$$

En définitive, la notion de fonction avec entrée et sortie typées coïncide avec la notion de schème cognitif. Cependant, cette vision fonctionnaliste doit être pondérée par le fait que les règles, et par suite les schèmes cognitifs, sont continuellement modifiés, créés et éliminés. Toutefois, rien n'empêche de concevoir qu'un ensemble de schèmes cognitifs suffisamment stables puisse être dédié à des procédures calculatoires.

B - Le typage au niveau des capacités cognitives

Le typage de haut niveau peut être compris comme une segmentation hiérarchique du monde à partir de la segmentation des primitives sensorielles, soit des classes d'équivalence de haut niveau. Cette segmentation est indispensable à la notion d'objet mais insuffisante puisque la notion d'objet recouvre également les moyens permettant de structurer cette segmentation et de l'accommoder au cours des usages. L'objectif, ici, sera uniquement de présenter la possibilité de réaliser une hiérarchie de types par l'intermédiaire de règles sans modifier l'architecture cognitive.

La théorie des types utilisée pour construire cette hiérarchie est celle des types simples proposés par Ramsey (1978) et Chwistek (1921). Dans ce cadre, un élément de type n peut être représenté comme un ensemble d'éléments de type $n-1$. Par exemple, le type 2 comportera tous les ensembles dont les éléments sont des individus et le type 3 comptera tous les ensembles dont les éléments sont des ensembles d'ensembles d'individus. Plus exactement, un type se définit comme une fonction d'apparence de telle sorte que la fonction d'appartenance à type $n-1$ est définie par le type n .

Ainsi, l'autoréférence n'est pas permise ce qui évite tout paradoxe issu de la théorie des ensembles. De plus, la construction des types s'effectuant successivement, cette théorie des types autorise un développement incrémental de la segmentation. Des pistes dans l'élaboration de schèmes cognitifs réalisant ce genre de hiérarchie de types seront détaillées

dans la dernière partie de ce chapitre car elle repose sur des mécanismes subcognitifs définis lors de la spécification.

En s'appuyant sur la définition formelle des schèmes cognitifs précédente, une fonction d'appartenance peut se comprendre comme un schème cognitif. La traduction de la théorie de types simples nécessite deux sortes de schèmes cognitifs, la première sorte décrit les schèmes cognitifs décrivant les éléments de type 0 et la seconde sorte de schèmes cognitifs décrit ceux décrivant les éléments de type supérieur à 0. Cette traduction nécessite également l'introduction d'un nouveau type de message t qui servira à contenir les étiquettes. Ainsi posés, les schèmes cognitifs constituant une hiérarchie de types simples au niveau de la sémiose s'écrivent de la manière suivante :

$$\left\{ \left\langle \right\rangle_I \wedge \left\langle \left\{ \left\{ \wedge_i M_i \right\}^n \right\} \right\rangle_E \supset M_j^t \right\}_{S_0}$$

$$\left\{ \left\langle \right\rangle_I \wedge \left\langle M_{j \in S_{k-1}}^t \right\rangle_E \supset M_k^t \right\}_{S_k}$$

Concrètement le typage de haut niveau correspond à la production de signes qui conditionnent les futures interprétations. En effet, la coïncidence temporelle des messages liée au principe d'univocité permet d'associer un ensemble de messages de différents types. Il devient alors possible d'exprimer la notion de relation $P(x)$ aperçue lors de la présentation de la logique du premier ordre :

$$\forall x P(x) \supset Q(y)$$

L'introduction dans l'architecture cognitive de l'opérateur \forall semble inévitable mais en réalité cette notion se trouve déjà présente dans l'évaluation de la prémisse qui sera abordée plus en détail au cours de la spécification. La promiscuité temporelle et l'univocité des messages n'autorisent pas d'ambiguïté sur leur association avec les prémisses des règles. Ainsi, en gardant le symbole \forall , la formule précédente est traduite en terme de règles appartenant à la mémoire épistémique par :

$$\left\langle \right\rangle_I \wedge \left\langle M_p^t \wedge M_{\forall i}^n \right\rangle_E \supset M_y^n$$

$$\left\langle \right\rangle_I \wedge \left\langle M_p^t \wedge M_{\forall i}^n \right\rangle_E \supset M_q^t$$

**

La différence entre le typage de bas niveau avec celui de haut niveau réside principalement dans le fait que le premier se trouve complètement défini au sein de l'architecture de sorte qu'aucun autre type de bas niveau supplémentaire ne peut être créée au cours de la sémiose, et le second typage est dynamique, il se construit au cours de la sémiose. Le typage de haut niveau reposant uniquement sur la promiscuité temporelle, il n'existe qu'au travers de l'utilisation des règles. La souplesse d'utilisation des messages étiquettes offre ainsi une totale liberté pour élaborer une segmentation complexe de haut niveau.

3.6. La notion de variable et la logique d'ordre n

La hiérarchisation des types au niveau de la sémiologie permet de représenter des prédicats d'ordre n . Toutefois, cette hiérarchisation n'autorise pas à elle seule de représenter une logique du premier ordre ou d'ordre supérieur. En effet, la formalisation définie jusqu'à présent n'autorise pas la traduction de la règle d'inférence de premier ordre évoquée précédemment :

$$(\forall x (A(x) \supset B(x)) , \forall x (B(x) \supset C(x))) \vdash \forall x (A(x) \supset C(x)) \quad (2)$$

Cette règle d'inférence (2) illustre deux problèmes pour simuler une logique d'ordre supérieur, d'une part (A) celui de l'intégration de la notion de variable et d'autre part (B) celui de la traduction d'une inférence effective sur les implications de l'architecture cognitive de la formule. Ces problèmes seront successivement abordés, toutefois seules des pistes de solutions seront proposées. En effet, ces problèmes complexes se trouvent associés à des propriétés qui n'apparaissent nécessaires qu'au troisième et quatrième stade de la cognition, par conséquent, ils ne sont pas primordiaux dans l'objectif de réaliser une autopoïèse sémiotique minimale.

A - Intégration de la notion de variable

Le premier problème, l'introduction de la notion variable, provient de l'impossibilité à traduire en un schème cognitif la formule (3) qui compose la règle d'inférence (2) :

$$\forall x A(x) \supset B(x) \quad (3)$$

Une manière de représenter cette formule serait de constituer un schème cognitif regroupant l'ensemble des implications pour tout x . Mais cette méthode n'est pas réalisable puisque le système doit acquérir par lui-même ses règles. Dans ce cas, la seule solution consiste à introduire la notion de variable dans l'architecture cognitive de sorte qu'un message puisse se transmettre au cours de la sémiologie. Cette transmission se trouve réalisée par un mécanisme subcognitif qui recopie le message ayant participé au déclenchement de la règle dans la conclusion. La désambiguïsation de la sélection s'effectue grâce à la correspondance entre le signe copié et la conclusion de leur type et de leur indice temporel. Afin de souligner le fait qu'une implication dépende de ce mécanisme subcognitif, cette dernière s'écrira comme suit :

6. « \supset_i » représente l'implication assimilée à une règle (ou une croyance) de la mémoire épistémique à laquelle est associée une évaluation globale traduisant la pertinence de la règle $i : P_i$, et dont le contenu de la conclusion correspond au message de même type et de même indice temporel s'appariant avec la prémisse.

La transition dans le temps d'un message permet alors de représenter la formule (3) par le couple de règles suivant :

$$\begin{aligned} \langle \rangle_I \wedge \langle M_a^t \wedge M_{\forall x}^n \rangle_E &\supset M_{\forall x}^n \\ \langle \rangle_I \wedge \langle M_a^t \wedge M_{\forall x}^n \rangle_E &\supset M_b^t \end{aligned}$$

Ce mécanisme de recopie peut être considéré comme une mémoire dynamique (ou de travail) dans la mesure où un signe se trouve par ce mécanisme reproduit donc réactualisé sans pour autant éliminer le précédent dans la mémoire événementielle.

B - Traduction des inférences portant sur des implications

Le second problème illustre bien la différence entre logique et effectif. Jusqu'à présent, l'inférence effective traduisait l'inférence logique en appareillant les signes de la mémoire événementielle aux prémisses des implications. Cependant, la règle d'inférence logique (2) manipule directement des implications pour en déduire de nouvelles, ce qui ne peut se traduire avec la règle d'inférence effective proposée. Pour répondre à cette difficulté, deux solutions se dessinent.

La première solution consiste à ajouter une nouvelle mémoire épistémique d'ordre supérieur qui permettrait de poser l'inférence effective suivante :

$$(\forall A, B, C) (\forall x (A(x) \supset B(x)), \forall x (B(x) \supset C(x)), ((\forall x (A(x) \supset B(x)) \wedge (\forall x (B(x) \supset C(x))) \supset \forall x A(x) \supset C(x)) \mapsto \forall x (A(x) \supset C(x)) \quad (4)$$

Cette solution entraîne néanmoins de nouvelles interrogations sur l'origine de ces règles et sur leur rôle dans l'autopoïèse sémiotique. Par ailleurs, concevoir l'abstraction comme un emboîtement conduit au problème de l'homoncule.

La seconde solution, celle qui sera privilégiée, propose d'ajouter un nouveau mécanisme subcognitif associé à une implication afin de générer de nouvelles règles. La prémisses d'une telle règle peut se concevoir en deux parties : d'un côté les messages avec un \forall devant leur contenu qui servent à filtrer les signes de la mémoire événementielle nécessaires pour formuler une nouvelle règle, et de l'autre, les messages attendus qui servent de déclencheur mais également de paramètre pour déterminer le rôle des messages filtrés dans la règle créée. Ces mécanismes subcognitifs suggèrent l'existence d'une convention minimale pour la création de règle qui reste à inventer ainsi que les schèmes cognitifs liés à la gestion et à la rétribution de telles règles. Cette convention permettrait par exemple de différencier les règles générant des règles normales (règles d'ordre zéro) et celles générant des règles transmettant des variables (règles d'ordre n en fonction du typage de haut niveau utilisé).

Cette solution suggère que la création de règles transmettant des variables à partir d'une succession d'implications qui transparait grâce à l'ordonnancement temporel des signes dans la mémoire suffit à modéliser la règle d'inférence logique (2). Ce point est très important car sur un plan logique cela signifie que toute création de règle découle d'un processus inductif puisque une succession de signes ne prouve pas l'existence d'une implication les reliant directement. Ce procédé permet ainsi de construire de nouvelles relations à partir de l'expérience cognitive. Néanmoins, rien n'empêche la conception de schèmes cognitifs de telle sorte qu'une partie de la sémiotique puisse être décrite parfaitement par des déductions logiques.

Par ailleurs, la modélisation de la règle d'inférence logique (2) implique qu'il soit possible de déterminer, par exemple en explicitant l'indice temporel, si le contenu du message $M_x(t)$ associé à $M_a(t)$ vaut celui du message $M_x(t-\Delta t)$ associé à $M_b(t-\Delta t)$. L'égalité entre ces deux messages peut se traduire par la une règle construite à partir de l'un des deux messages en ajustant l'indice temporel sur l'autre et dans ce cas la valeur de la légitimité de la conclusion représente le degré d'égalité. La traduction de la règle d'inférence logique (2) requiert alors plusieurs schèmes cognitifs imbriqués se déroulant dans le temps.

Enfin, dans le cadre de la création de règles par recopie du contenu de la mémoire événementielle, l'adduction comprise comme une déduction renversée peut se traduire par une règle ayant parmi ses conditions d'activation la conclusion d'une autre et comme conclusion une partie de sa prémisse. Ce genre de règle suggère l'existence de règle qui en construisent d'autres en inversant l'ordre des indices temporels.

**

En somme, l'établissement d'une logique d'ordre supérieur oblige à un examen plus approfondi concernant la création de règles. Toutefois, le passage à un ordre supérieur de l'autopoïèse sémiotique repose essentiellement sur la notion de recopie qui permet à la fois de représenter la notion variable par la recopie du contenu d'un signe dans une conclusion, et d'induire des règles par recopie partielle de la mémoire événementielle. La gestion des indices temporels se révèle également indispensable pour la désambiguïsation des situations. En d'autres termes, la pensée se déroule dans le temps et elle se reconnaît et se développe dans son déroulement.

L'étude de l'architecture cognitive par le biais des systèmes formels a permis de compléter et de préciser certaines de ses propriétés (la légitimité d'un message, la pertinence d'une règle ou l'inhibition) afin que de la sémiologie puisse émerger des raisonnements similaires à ceux décrits par ces logiques. Maintenant que les propriétés et les caractéristiques générales de l'architecture cognitive ont été dégagées, l'élaboration d'une spécification d'une architecture cognitive devient envisageable. Parmi toutes les familles d'algorithmes utilisés en intelligence artificielle, la famille des systèmes de classeurs se rapproche structurellement le plus de l'architecture cognitive proposée. En effet, ces systèmes sont des architectures de contrôle à base de règles qui se déclenchent en fonction de l'interaction entre deux mémoires : une mémoire à court terme et une mémoire à long terme qui peuvent s'assimiler au premier abord à une mémoire événementielle et à une mémoire épistémique. La spécification sera construite à partir des compatibilités et des incompatibilités conceptuelles avec ce genre de système.

4. Les systèmes de classeurs

Les systèmes de classeurs constituent des architectures de contrôle à base de règles comportementales simples et subsymboliques <condition> : <action> qui permettent l'application d'une méthode évolutionniste. Dans ce cadre, l'apprentissage par renforcement sera utilisé pour parvenir à une base de règles viables. Le domaine d'application de cette famille d'algorithmes est très varié : le diagnostic médical (Bonelli *et al.*, 1991), l'économétrie (Schulenburg *et al.*, 2001), les jeux vidéo (Robert *et al.*, 2002) (Sanza, 2001), l'aide pédagogique (Buche *et al.*, 2004) et la robotique (Dorigo, 1995).

La diversité des paradigmes et des applications reflète celle des systèmes de classeurs existants. Toutefois, dans un premier temps, (A) un ensemble de remarques sur l'architecture générale des systèmes de classeurs comparativement à l'architecture cognitive proposée sera formulé. Dans un second temps, une étude plus fine sur les ressemblances et les incompatibilités sera entreprise pour chacune des quatre architectures types regroupant la majorité des systèmes de classeurs : (B) l'architecture originelle, le CS1 (Holland, 1976), (C) les architectures simplifiées de type ZCS (Wilson, 1994) et XCS (Wilson, 1995), (D) les architectures avec anticipation comme le ACS (Stolzman, 1998) et (E) les architectures hétérogènes introduisant des hiérarchies ou des motivations. Enfin, une analyse générale

des limitations des systèmes de classeurs sera proposée avant d'aborder la spécification d'une architecture cognitive qui tente de les dépasser et ses choix de conception.

4.1. Remarques sur l'architecture générale des systèmes de classeurs

Les systèmes de classeurs appartiennent à la famille algorithmique des systèmes à apprentissage semi-supervisé (c'est-à-dire par renforcement). La Figure III-12 représente le schéma dans lequel s'inscrivent tous ces systèmes. Comparé au schéma de l'architecture cognitive (Figure III-10), celui des systèmes de classeurs possède également trois couches : le corps de l'agent avec ses senseurs et effecteurs, la gestion des interfaces d'entrée et de sortie et l'architecture de contrôle. Les interfaces peuvent être assimilées à des primitives sensorimotrices mais leur conception s'effectue généralement dans un autre cadre de pensée de sorte que la fréquence d'acquisition des capteurs correspond communément à la fréquence du flux d'information traité arrivant à l'architecture de contrôle. Autrement dit, les interfaces sont des processus passifs qui se limitent à un module de codage et de décodage entre les senseurs, les messages et les effecteurs.

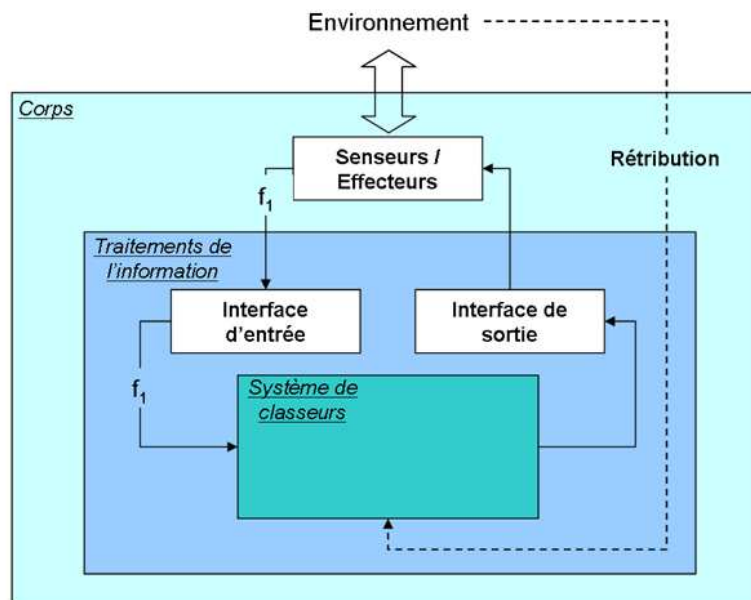


Figure III-12 : Schéma élémentaire d'une architecture de contrôle se fondant sur l'apprentissage par renforcement dans lequel s'inscrivent les systèmes de classeur.

À ce niveau de conception, la différence entre les deux architectures réside dans l'existence, pour les systèmes de classeurs, d'une rétribution extérieure venant directement influencer l'architecture de contrôle. La gestion de cette rétribution sera cependant une source d'inspiration pour imaginer d'une part la dynamique sémiotique et d'autre part les mécanismes d'une rétribution interne. En effet, le point de divergence ne se trouve pas dans le concept de renforcement mais dans la provenance de la rétribution, ainsi que dans la méthode de redistribution qu'elle induit au sein du système. Cette critique qui dépasse le cadre des algorithmes des systèmes de classeurs vise l'apprentissage par renforcement en général, étudié lors de la présentation du fonctionnalisme écologique dans le premier chapitre.

4.2. L'architecture originelle, le CS1

À l'origine, l'architecture des systèmes de classeurs CS1 (*Classifier System One*) conçue par Holland (1976) repose sur l'idée que la mémoire à long terme contient les règles comportementales qui se déclenchent en fonction de l'état d'une mémoire à court terme, alimentée par l'activité du système. La mémoire à court terme correspond alors à une liste de messages provenant soit de l'interface d'entrée (les messages sensoriels), soit de la conclusion d'une règle (les messages d'action et les messages internes). À chaque pas de temps, l'interface de sortie interprète les messages d'action se trouvant dans la mémoire à court terme. Les règles appelées « classeurs » peuvent être déclenchées soit par un message sensoriel, soit par un message interne. Un message, une fois apparié à une règle ou à l'interface de sortie, disparaît de la mémoire à court terme.

Plus formellement, un message se traduit par une chaîne de valeurs discrètes de taille l appartenant à un alphabet, traditionnellement $\{0,1\}$. La condition s'écrit avec une chaîne de valeurs discrètes de longueur l avec le même alphabet mais augmenté du caractère #. Lors de l'appariement, le dièse s'interprète comme un joker ou une valeur indéterminée. Par exemple, la condition $[0,1, 1, \#]$ s'apparie avec les messages $[0, 1, 1, 0]$ ou $[0, 1, 1, 1]$ mais pas avec le message $[0, 0, 1, 1]$. L'introduction du caractère dièse permet la généralisation des règles, ce qui est compris comme une forme d'abstraction.

Dans ce cadre, le problème de l'apprentissage devient la recherche de l'ensemble de règles qui offre à la fois un comportement viable selon un certain critère et possède un nombre réduit de règles, donc les plus générales possibles. En d'autres termes, cela revient à poser deux questions : (i) comment identifier les règles adaptées et éliminer les autres ? Et (ii) comment générer de nouvelles règles pour essayer d'en trouver de meilleures ou de plus générales ?

i - Comment identifier les règles adaptées et éliminer les autres ?

En fait, le problème de l'identification des règles adaptées porte à la fois sur le court terme (la règle la plus adaptée à la situation actuelle) et sur le long terme (celles ayant apporté le plus de satisfaction au cours de l'activité de l'agent). Afin d'opérer ces deux types de sélection, chaque règle possède une *force* correspondant à une évaluation globale de la règle à l'instant t qui résulte notamment des rétributions positives ou négatives reçues antérieurement.

La sélection à court terme se déroule en deux phases. La première phase détermine l'ensemble des règles pouvant s'apparier avec le même message. La seconde phase décide quel message d'action sera envoyé parmi ceux proposés par les règles issues de la première phase. Principalement, trois méthodes existent pour réaliser ce choix : la première méthode, la plus simple, consiste à prendre le message d'action de la règle ayant la plus grande force. Afin de favoriser l'exploration de l'environnement, la seconde méthode, la « roulette de la fortune », tire une règle aléatoirement parmi les règles r présélectionnées avec une probabilité pondérée par leur *force* F :

$$P(r_i) = \frac{F_{r_i}}{\sum_{j=1}^N F_{r_j}}$$

La troisième méthode, variante de la précédente, repose sur l'idée d'une décision collégiale de l'action dans le sens où la probabilité de l'action choisie dépend de la *force* de

l'ensemble des règles R_a désignant cette action A par rapport à la *force* F de l'ensemble de toutes les règles R_i s'appariant avec le même message :

$$P(A_k) = \frac{\sum_{r_i \in R_a} F_{r_i}}{\sum_{r_j \in R_i} F_{r_j}}$$

Toutefois, dans le cadre des systèmes de classeurs proches du CS1, la première et la seconde méthode demeurent les plus utilisées.

La sélection à long terme dépend de la manière de concevoir la rétribution soit comme l'évaluation de chaque action, soit comme la récompense à la réalisation d'un objectif. Dans le cas où la rétribution survient à chaque cycle de fonctionnement, un système de classeurs ressemble à la première des trois catégories d'algorithme pour l'apprentissage par renforcement évoqué lors de la présentation du fonctionnalisme écologique du premier chapitre, c'est-à-dire une sorte de programmation dynamique qui évalue l'action puis l'améliore en fonction de la rétribution. Mais l'inconvénient de cette approche vient de la nécessité de modéliser le monde a priori. Par exemple, dans le cadre d'une rétribution à chaque cycle, un agent situé dans un labyrinthe (Figure III-13) devant apprendre à aller du point A au point B reçoit une rétribution proportionnelle à sa distance au point B , ce qui nécessite à un modèle implicite du monde.

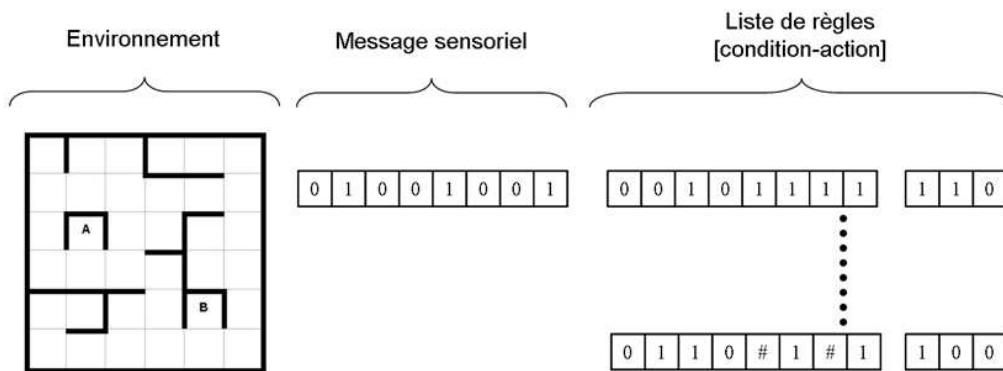


Figure III-13: Exemple du problème du labyrinthe pour aller de A à B avec un système de classeurs : les messages sensoriels et conditionnels traduisent l'accessibilité aux huit cases entourant l'agent et les messages d'actions de celui codés sur 3 bits.

Dans le cas où la rétribution survient de manière irrégulière ou uniquement lors de la réalisation de l'objectif ou, le problème du partage de la rétribution apparaît puisque celle-ci ne résulte pas de la dernière action mais de l'enchaînement de toutes les règles précédentes dont certaines peuvent avoir été plus déterminantes que d'autres. Les systèmes de classeurs de type CS1 souhaitent apporter une solution au second cas, tout en prenant en compte que ce mécanisme de redistribution de la rétribution permet de départager des règles proches et d'évincer les règles jamais appariées, donc jamais évaluées.

Dans cette perspective, Holland (1980) s'appuie pour l'évaluation des classeurs sur une métaphore financière : une règle devient un agent financier dont le capital correspond à sa force, les messages de la mémoire représentent les biens dans lesquels il faut investir. La conclusion d'une règle se comprend alors comme la transformation du bien initial qui est remis sur le marché. Une règle recherche un certain type de biens dont elle sait que la transformation qu'elle opère lui permettra de faire une plus-value lors de la vente.

L'algorithme résultant, le *Bucket brigade*, possède de nombreuses variantes mais son principe reste le même, chaque règle se trouve redevable de celles ayant permis son déclenchement.

Plus précisément, lors de la première phase de la sélection d'une action, toutes les règles font une *enchère* sur le message avec lequel elles s'apparient. L'*enchère* E dépend à la fois d'un coefficient d'enchère C_e et de la force de la règle :

$$E = C_e * F$$

D'autres versions introduisent la notion de *spécificité* S qui prend en compte le degré de généralisation d'une règle telle que :

$$E = C_e * F * S$$

Par ailleurs, un bruit peut être également ajouté au calcul de l'*enchère* pour inciter à l'exploration :

$$E = C_e * F + \text{bruit}$$

Ce dernier rajout peut se comprendre comme un compromis avec la méthode de sélection « roulette de la fortune ». Au final, la règle proposant l'enchère la plus importante l'emporte et se déclenche. Si l'objet de l'enchère est un message interne alors le système redistribue la somme prélevée à la règle gagnante vers toutes les règles dont la conclusion correspond à ce message interne. Si l'objet de l'enchère se révèle être un message sensoriel, la redistribution touche toutes les règles dont la conclusion correspond au dernier message d'action. Ainsi, hormis la rétribution externe, une règle k peut recevoir le *paiement* P d'une règle j dont le montant dépend de l'*enchère* avec le nombre de règles N ayant une conclusion identique à l'objet de l'enchère :

$$P_k = \frac{E_j}{N}$$

Certaines versions du CS1 sont plus restrictives sur la redistribution et ne prennent en compte que la règle qui est effectivement à l'origine du message. En général, ce mécanisme de redistribution permet d'identifier les règles inutiles, c'est-à-dire celles dont le déclenchement n'en entraîne pas d'autres ou n'amène pas de rétribution extérieure, et de favoriser les enchaînements de règles. Ainsi, une chaîne n'aboutissant pas à une rétribution finale disparaîtra progressivement. Par ailleurs, le mécanisme d'enchère permet également de départager deux règles avec des conditions proches et ainsi de diminuer la redondance, sauf cas de rétribution extérieure particulière. Afin d'éliminer les « passagers clandestins », c'est-à-dire les règles qui ne s'apparient jamais et ne paient jamais d'enchère, toutes les règles subissent une *taxe* T proportionnelle à leur *force* F à chaque cycle de fonctionnement :

$$T = C_t * F$$

Ainsi, l'équation de l'évolution dynamique de la force d'une règle prend la forme suivante avec la *rétribution* R :

$$F(n+1) = F(n) - T(n) - E(n) + P(n) + R(n)$$

La dynamique induite par cette équation sépare les règles participant à l'obtention de la rétribution des autres en augmentant la *force* des premières et en diminuant celle des secondes qui, arrivées sous un certain seuil, pourront être éliminées. Dans certaines versions de système de classeurs, les paramètres conditionnant l'évolution de la dynamique des règles peuvent être modifiés au cours de l'activité de l'agent afin de retrouver les phases

d'exploration et d'exploitation traditionnellement utilisées dans les algorithmes d'apprentissage. Par exemple, une première manière d'introduire ces phases consiste à modifier au cours du temps le coefficient d'enchère C_e . Une seconde manière peut influencer la dynamique en amont en utilisant au début une sélection de règles à partir de la méthode « roulette de la fortune » pour ensuite basculer après un certain temps sur la méthode déterministe qui sélectionne la règle à la force maximale.

ii - Comment générer de nouvelles règles ?

Cette évolution dynamique de la force des règles permet de « reconnaître » les règles intéressantes mais ne permet pas de répondre à une situation inconnue ou d'optimiser ces règles. Ces deux problèmes se résolvent par un module de création de règles contenant deux procédés. Le premier procédé, le *covering*, génère une règle lorsqu'aucune règle ne correspond à l'état de la mémoire à court terme : la partie <condition> de la nouvelle règle prend la valeur du message incompatible se trouvant dans la mémoire à court terme avec des caractères # rajoutés aléatoirement et la partie <action> résulte d'un tirage aléatoire des règles selon une probabilité pondérée par leur *force*. En introduisant un tirage aléatoire de l'action lorsque le nombre de règles se trouve insuffisant, le *covering* permet a priori à l'agent de commencer à évoluer dans le monde sans aucune règle.

Le second procédé, l'algorithme génétique, recombine périodiquement les règles existantes en fonction de leur *force*. L'algorithme permet à la fois d'orienter l'exploration de l'espace des règles et de généraliser de manière optimale les règles. Plus précisément, deux types d'implémentation ont été proposés, d'une part la méthode *Pittsburgh* (Smith, 1980) qui considère un sous-ensemble de règles comme un individu et d'autre part la méthode *Michigan* (Holland *et al.*, 1978) qui considère la règle comme un individu. Les deux méthodes possèdent des particularités propres, toutefois, la méthode *Michigan* offre souvent davantage de stabilité et une plus grande rapidité de convergence (Buche, 2006). Les trois types de classeurs présentés dans cette section appartiennent à la famille de systèmes de classeurs *Michigan*.

*

La sélection, l'évaluation et la génération des règles ayant été précisées, le schéma de l'architecture des systèmes de classeurs de type CS1 peut être avancé (Figure III-14) et comparé avec l'architecture cognitive proposée.

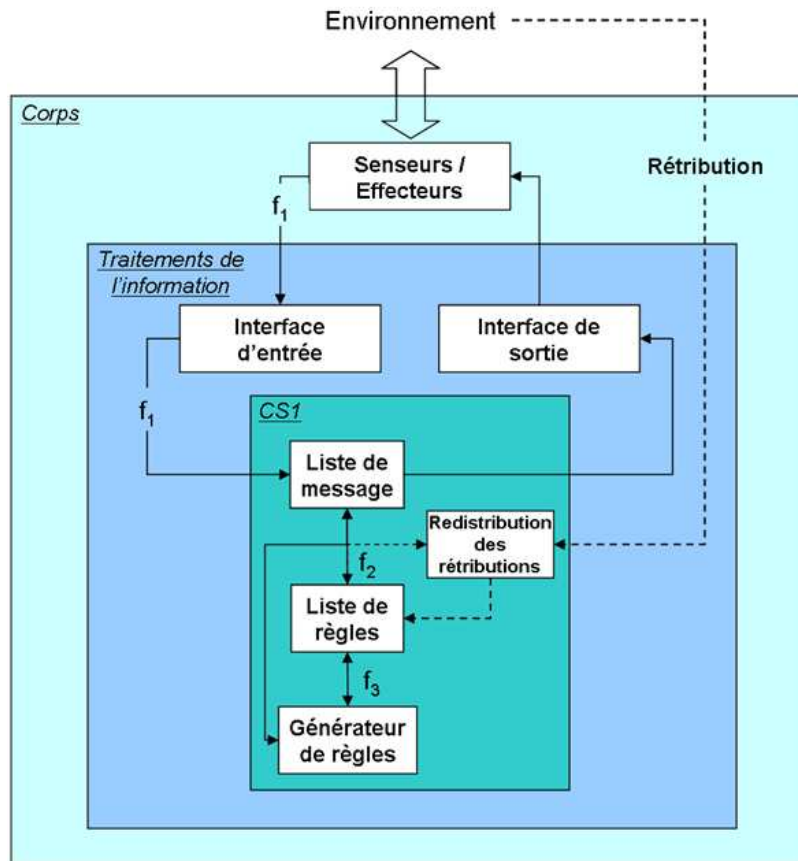


Figure III-14 : Schéma d'un système de classeurs de type CS1

Au premier abord, les deux architectures se ressemblent. En effet, elles se fondent toutes les deux sur l'existence de deux types de mémoire qui interagissent, l'une servant de récepteur et d'émetteur et l'autre de base de règles. Par ailleurs, le mécanisme de messages internes du CS1 traduit bien l'idée d'une réflexivité entre deux mémoires qui permet de modifier son comportement par rapport à un modèle interne c'est-à-dire un contexte interprétatif. Les messages d'action, bien qu'inutilisés par les prémisses des règles, suggèrent également une traduction sémiotique de l'activité de l'agent. Le générateur de règles du CS1 est un principe qui se retrouve également dans l'architecture cognitive mais avec une gestion différente de la création de règle. La ressemblance se poursuit encore sur les propriétés des architectures. L'évolution des règles offre des capacités de généralisation, ce qui autorise qu'un même message puisse être produit par des règles de différentes origines et provoquer le déclenchement de différentes règles comme pour le signe avec sa provenance et son interprétation. De plus, le CS1 parvient à distinguer la sélection ponctuelle d'une règle de son évaluation globale grâce à une dynamique interne de rétribution produite par la compétition entre les règles. En effet, dans ce cas la force correspond à la l'évaluation globale d'une croyance pour l'architecture cognitive et le degré d'appariement entre un message et une règle à une évaluation ponctuelle d'une croyance. Ce point a été dégagé comme un critère primordial pour la pérennité de la sémiose.

Toutefois, le CS1 peut-il être considéré comme une autopoïèse sémiotique ? La mémoire à court terme reçoit des messages qui déclenchent des règles renvoyant des messages et ainsi de suite. Une dynamique interne perturbée par les messages sensoriels peut être imaginée. Le système se régulera en fonction des différents paramètres :

coefficient d'enchère, coefficient de taxe, fréquence de cycle d'interprétation, fréquence de l'arrivée des messages sensoriels et fréquence de la régénération de la population de règles par algorithme génétique. Autrement dit, ce type de système de classeurs peut être considéré comme un système s'autoréglant sans consigne explicite. Par ailleurs, la création de règles par *covering* constitue également un mécanisme de régulation ressemblant à celui du renouvellement des éléments qui constituent le système cellulaire autopoïétique élémentaire expliqué au cours du premier chapitre dans la section traitant de l'évolutionnisme.

Néanmoins, deux critères fondamentaux manquent aux systèmes de classeurs pour les qualifier d'autopoïèse sémiotique. Le premier critère provient du point (a) de la définition d'un système autopoïétique donné par Varela (1989) « Un système autopoïétique est organisé comme un réseau de processus de production de composants qui (a) régénèrent continuellement par leurs transformations et leurs interactions le réseau qui les a produits, [...] ». Or ici, le *covering* ou l'algorithme génétique se déclenche indépendamment du contenu des règles et des messages, autrement dit ce n'est pas l'organisation particulière des règles et des messages qui compense la suppression des règles résultant de la dynamique. Cela justifie l'emploi de mécanismes de création inférentiels dans l'architecture cognitive proposée. Par ailleurs, la rétribution interne due au reversement de l'enchère correspond à la notion de régénération par interaction entre les éléments dans le cadre d'une autopoïèse sémiotique. Cependant, cette rétribution interne native se révèle insuffisante puisqu'en définitive le maintien des règles dépend d'une rétribution extérieure en bout de chaîne. Cela va à l'encontre de la définition du premier stade de la généalogie cognitive qui met en avant l'importance de l'indépendance du système de règles par rapport à toute rétribution, en dehors de la rétribution native résultant de la seule auto-organisation des règles comportementales.

Le deuxième critère fondamental faisant défaut aux systèmes de classeurs réside dans la notion d'interprétation du signe développé précédemment (2.3-B). En effet, un concept ne se définit pas par son contenu mais par son opposition avec d'autres concepts qui, par conséquent, délimitent son signifié. Dans le CS1, il existe bien une compétition entre les règles, mais l'auto-organisation induite correspond à une répartition des messages possibles entre un certain nombre de règles. Une fois cette répartition établie, le domaine d'application de la règle se trouve entièrement déterminé par sa condition. Autrement dit, l'auto-organisation produite par la compétition et aidée de l'algorithme génétique ne propose pas une auto-organisation fondée directement sur l'opposition entre les différents éléments du système autorisant un ajustement perpétuel des conditions. Une première critique sur le problème de la définition du contenu des messages serait d'incriminer l'utilisation de valeurs discrètes qui empêchent la notion de flou nécessaire pour représenter le chevauchement entre deux conditions (Valenzuela-Rendon, 1991). Mais en réalité, l'introduction de la notion de flou ne représente qu'une manière plus fine d'introduire la généralisation dans les prémisses des règles. Le problème est que l'évaluation des règles s'effectue indépendamment les unes des autres. L'évaluation de l'appariement correspond à un droit de passage en fonction d'un seuil vers une deuxième sélection basée cette fois-ci sur la force des règles. Il faut souligner que la force représente l'évaluation globale de la règle et par conséquent n'est pas liée avec l'évaluation de l'appariement présent. Le contenu d'une règle, c'est-à-dire l'ensemble des situations permettant de passer la première sélection, se trouve entièrement défini par sa prémisses indépendamment de celles des autres et la seconde sélection, qui correspond à la compétition, se fonde sur la force. Les contenus ne se trouvent ainsi pas au cœur de la compétition entre les règles, ils s'ajustent indirectement

par la maximisation des récompenses. Par conséquent, les systèmes de classeur manipulent nécessairement des subsymboles et non des signes.

L'absence de ces deux critères, la régénération des règles par la sémiose et l'ajustement des signes, montre que les systèmes de classeurs de type CS1 ne possèdent pas les qualités nécessaires pour être considérés comme des autopoïèses sémiotiques. Cependant cette analyse permet de définir ces systèmes comme des structures dissipatives subsymboliques dans le sens où ils possèdent une dynamique propre basée sur une auto-organisation subsymbolique, la force globale représentant l'énergie. Néanmoins, les systèmes autopoïétiques représentant un sous-ensemble des systèmes dissipatifs (Stewart, 2004), l'analogie entre la dynamique des structures dissipatives et celle des systèmes de classeurs renforce l'idée que ces derniers représentent un solide point de départ pour concevoir une architecture cognitive. Mais au-delà de ces concepts généraux, d'autres points de divergence subsistent entre l'architecture CS1 et l'architecture cognitive proposée, parmi eux quatre points méritent notamment d'être présentés.

Le premier point d'incompatibilité concerne la notion d'intégrité qui empêche toute équivalence entre la mémoire à court terme et la mémoire événementielle. En effet, l'élimination des messages interprétés contrevient à l'idée imaginée dans l'architecture cognitive que la disparition des messages dépend uniquement de la durée de leur présence, indépendamment des interprétations. Par ailleurs, la suppression des messages interprétés empêche l'attribution d'indices temporels concernant la survenue des messages. En second lieu, la mémoire événementielle possède deux registres avec une gestion antagoniste de l'écoulement du temps permettant de prédire et de prévoir des séquences d'actions alors que la mémoire à court terme ne possède qu'un seul registre. En troisième lieu, la dynamique induite par le CS1 oblige qu'une règle satisfaisante soit régulièrement activée pour ne pas succomber à cause de la taxe. Le *Bucket brigade* ne confère pas aux règles bien adaptées un statut qui permettrait de les conserver même si les dernières situations n'ont pas permis de les exploiter, contrairement à ce que doit faire une architecture cognitive. Le dernier point de divergence souligne que le CS1 repose uniquement sur une logique positive qui ne permet pas de modéliser des logiques non monotones, condition indispensable pour un agent dans un monde ouvert. Mais surtout, ce point de divergence pose la question du maintien de la variabilité qui est intrinsèque aux systèmes biologiques qui conduit à des architectures exemptes de téléonomie. Un système biologique n'a pas été construit par dessein de satisfaire un but particulier. Il est apte à s'adapter à la résolution de buts multiples et de ce fait doit inclure une réserve de variabilité qui rend nécessaire de stabiliser à long terme des éléments n'ayant pas d'utilité immédiate.

Cependant, ce ne sont pas ces remarques qui ont incité à modifier ce type d'architecture. Concrètement, la multiplicité des messages des systèmes de classeurs CS1 rend difficile l'analyse en dehors des cas simples. De plus, l'introduction des messages internes offre une richesse intéressante mais en même temps elle agrandit énormément l'espace de recherche par l'algorithme génétique. Par ailleurs, ce dernier peut produire des règles qui s'entretiennent mutuellement, parasitant ainsi le système. Dans la perspective d'une architecture cognitive, le défaut majeur provient de l'algorithme génétique puisque la production de règles doit résulter d'un schème cognitif canalisant l'heuristique vers un objectif dont l'intérêt peut être évalué, ce qui remet en cause la notion de renforcement externe. Dans le cadre de l'architecture cognitive, une recombinaison aléatoire de règles par le biais d'un mécanisme subcognitif reste autorisée mais les modalités de cette recombinaison et de sa réalisation doivent dépendre de la sémiose. Dans une démarche ne

remettant pas en question la notion de rétribution externe, le principal défaut de l'architecture CS1 provient de son trop grand nombre de degrés de liberté.

4.3. Les architectures simplifiées de type ZCS et XCS

Afin de réduire ces degrés de liberté, Wilson (1994) propose de simplifier l'architecture CS1 en supprimant la mémoire à court terme (Figure III-15) pour obtenir un système de classeur minimal, le ZCS (*Zeroth level Classifier System*). En d'autres termes, Wilson enlève la possibilité de générer des messages internes. Cette décision transforme alors les systèmes de classeurs en des systèmes purement réactifs impliquant la synchronisation de flux entrants et sortants des interfaces. Hormis la suppression de la mémoire à court terme, le ZCS conserve les principes du CS1. Pour la sélection de l'action après appariement, le ZCS emploie la méthode de la « roulette de la fortune » pour la phase d'exploration, puis la méthode de la force maximale pour celle d'exploitation. La dynamique des règles repose sur le *Bucket brigade* en prenant en compte le degré de spécialisation des règles. Enfin, la génération des règles par algorithme génétique utilise l'approche *Michigan*.

En limitant les systèmes de classeurs à un comportement purement réactif, les capacités d'apprentissage par renforcement se réduisent aux environnements markoviens définis dans le premier chapitre. Par exemple, le labyrinthe de la Figure III-13 présente un environnement non markovien du fait que A et B se traduisent par un message sensoriel identique (*perceptual aliasing*). En dehors du cas où elles correspondent respectivement au point de départ et à l'objectif, ces positions introduisent des erreurs de décision par rapport à la recherche d'une source de récompense alors que, dans les CS1, ces ambiguïtés peuvent être en théorie levées grâce à des messages pouvant servir de repère interne.

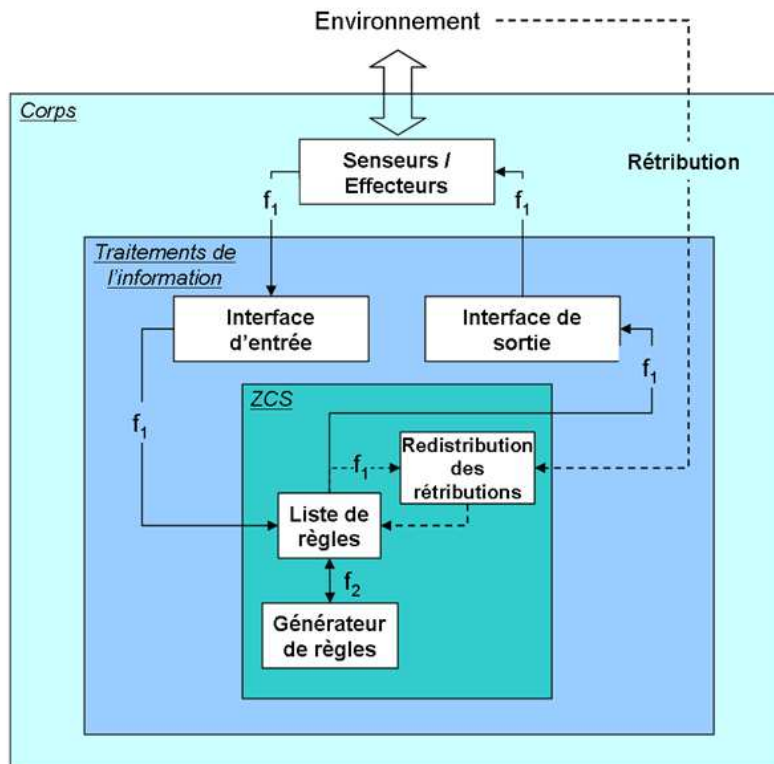


Figure III-15 : Schéma d'un système de classeurs de type ZCS

Parallèlement, la simplification des systèmes de classeurs a permis à Dorigo *et al.* (1994) de montrer l'équivalence de principe entre la dynamique d'apprentissage induite par la *Bucket brigade* et celle induite par le *Q-Learning*. Cet équivalent souligne le problème des systèmes de classeurs à rétribuer les règles éloignées dans la chaîne de règles conduisant à la récompense. En effet, les règles éloignées de la récompense obtiennent des rétributions moindres et sont au final créditées d'une force relativement faible même si la règle se révèle optimale. À l'inverse, une règle sous-optimale se situant à proximité de la récompense sera créditée d'une force supérieure à une règle éloignée optimale. Cette situation remet en question la pertinence d'une sélection sur le long terme avec la force comme critère.

Afin de répondre à ce problème, Wilson (1995) modifie le *Bucket brigade* en s'inspirant du *Q-Learning* qui repose à la fois sur la récompense espérée à chaque couple condition/action et sur l'erreur de cette prédiction. Avec cette modification, le ZCS devient le XCS, la force d'une règle correspond au paiement escompté à laquelle se trouve associé un nouveau paramètre : l'erreur sur la prédiction. En reprenant les termes du premier chapitre concernant l'étude sur l'apprentissage par renforcement, la force représente la fonction d'utilité $Q(s,a)$. Ainsi, le mécanisme de rétribution correspond à la résolution des équations de Bellman. En intégrant le *Q-Learning*, les systèmes de classeurs de type XCS rentrent formellement dans la troisième catégorie d'apprentissage par renforcement représentant les méthodes par différence temporelle. Le système ne cherche plus à trouver des règles adaptées au problème mais à approximer la fonction d'utilité. Par conséquent, le système résout le problème de la maintenance des longues chaînes d'actions puisqu'il ne vise pas directement la maximisation de la récompense attendue. Cela traduit bien l'idée du schéma *Actor-critic* de Barto (1995) qui dissocie le processus d'apprentissage et celui de décision.

Par ailleurs, la méthode de sélection de l'action après appariement des règles utilise un tirage aléatoire pondéré par la force comme la prédiction de paiement et l'erreur sur cette prédiction :

$$P(A_k) = \frac{\sum_{r_i \in R_a} (1 - erreur_{r_i}) * F_{r_i}}{\sum_{r_j \in R_t} (1 - erreur_{r_j}) * (F)_{r_j}}$$

Les synthèses de Gérard (2002) et de Métivier (2004) présentent plus en détail les caractéristiques qui différencient les systèmes ZCS et le XCS. Sans changer l'architecture de base ZCS, le XCS parvient à améliorer significativement les résultats des systèmes de classeurs. La synthèse réalisée par Buche *et al.* (2006) montre une grande diversité dans la conception de systèmes de classeurs développés à partir du ZCS et plus particulièrement du XCS. Mais toutes ces versions se révèlent en définitive très spécifiques au problème que souhaitent résoudre leurs concepteurs.

Les architectures XCS et ZCS s'éloignent de l'architecture proposée, par la suppression de la mémoire à court terme. Cependant, l'amélioration de la performance des systèmes de classeurs, suite à l'intégration d'une dynamique inspirée du *Q-Learning*, incite à étudier sa compatibilité avec l'architecture cognitive proposée. Pour cela, il faut revenir sur l'objectif des algorithmes de *Bucket brigade* et de *Q-Learning* qui consiste à maintenir des enchaînements de règles en vue d'une rétribution finale. Le ZCS se concentre uniquement sur la rétribution finale alors que le XCS s'attache également à la transition entre deux maillons de cette chaîne, ce qui rend ce dernier plus efficace. Dans le cadre de l'architecture cognitive proposé, un mécanisme de rétribution interne et natif se révèle indispensable

pour instaurer une compétition entre les règles par défaut. La rétribution interne et native doit alors provenir uniquement du déclenchement des règles. Cette conception de rétribution native correspond a priori aux rétributions par *paiement* entre les règles instaurées par le *Bucket brigade* et le *Q-Learning*. Ces deux mécanismes sont-ils cependant compatibles avec la conception générale de l'architecture cognitive ? Les deux premiers stades de la généalogie de la cognition n'exigent pas l'établissement de chaîne d'action. Le « bon » enchaînement d'actions provient essentiellement du « bon » enchaînement de conditions environnementales comme il a été évoqué lors de la présentation de l'interactionnisme varelrien dans le premier chapitre, avec les exemples éthologiques de Lorenz. La capacité à détecter des enchaînements d'action avantageux apparaît au troisième stade de la généalogie de la cognition. L'identification d'une chaîne n'a de sens que dans le cadre d'un objectif final, or la détermination de celui-ci contrevient a priori à l'idée que la cognition doit justement créer ces objectifs ou du moins qu'il soit explicité par la sémiose.

Ainsi, dans la perspective de la généalogie de la cognition avancée dans le second chapitre, les mécanismes de dressage assurant l'enchaînement d'actions appartiennent au développement des schèmes cognitifs et non à la dynamique de base instaurée par l'architecture cognitive. Par conséquent, le *Q-Learning* et de manière plus large les méthodes de type différence temporelle, restent intéressantes voire primordiales pour dégager des chaînes d'actions, mais leur intégration dans un agent cognitif apparaît au niveau des schèmes cognitifs et éventuellement au niveau des primitives sensorielles de manière rudimentaire et figée. En somme, la rétribution interne native ne doit pas se bâtir sur la détection de chaîne d'actions. L'alternative consiste à concevoir l'auto-organisation des règles comme la source de la rétribution interne native. Les principes d'une telle auto-organisation seront avancés dans le cadre de la spécification de l'architecture cognitive proposée.

4.4. Les architectures avec anticipation comme le ACS

Les systèmes de classeurs avec anticipation de type ACS (*Anticipatory Classifier System*) (Stolzman, 1998) souhaitent modéliser l'apprentissage latent selon les théories de Hoffman (1993). L'apprentissage latent est un apprentissage qui se réalise en dehors de toute incitation explicite à apprendre et qui demeure invisible si l'environnement n'amène pas l'organisme à émettre un comportement observable issu de cet apprentissage. La prise en compte de l'apprentissage latent montrant les limites de conception de l'apprentissage nécessairement guidé par un renforcement extérieur oblige à dissocier l'apprentissage de l'environnement et celui des actions liées à une récompense. En général, le modèle d'un environnement correspond à une liste des probabilités de passage d'un état à un autre, suite à une certaine action. Dans le cadre des systèmes de classeurs, cela se traduit par des règles comprenant trois parties : <condition> : <action> : <effet>. À chaque déclenchement, la règle se trouve évaluée en fonction du degré de correspondance entre l'effet prévu et l'état réel (Figure III-16). Des caractères # peuvent être introduits dans la partie <effet> afin de dégager des régularités environnementales. Cette capacité de généralisation des règles permet d'établir une liste non exhaustive de toutes les transitions.

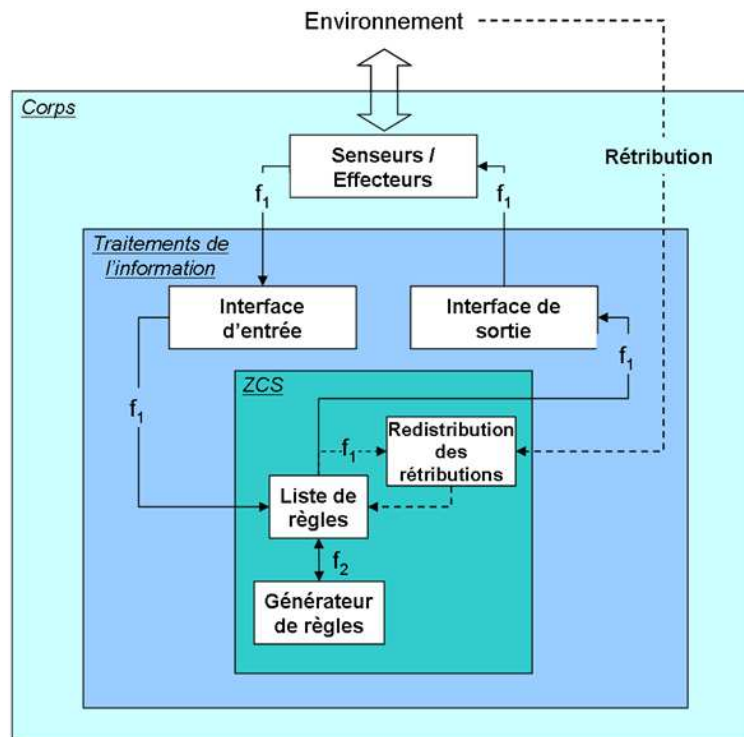


Figure III-16 : Schéma d'un système de classeurs de type ACS

L'introduction de la prédiction dans l'évaluation interne des règles améliore la vitesse de convergence des systèmes de classeurs et peut s'interpréter comme une confirmation de l'importance de la notion d'anticipation dans les processus cognitifs (Sigaud, 2001). Par ailleurs, le but explicite de vouloir anticiper permet de proposer des heuristiques adaptées à ce problème (Gérard, 2002) et ainsi d'améliorer a priori la qualité des règles générées par rapport à celles issues d'algorithmes génétiques classiques. Toutefois, l'efficacité de ces méthodes semble restreinte à certains types d'environnement (Buche, 2006). La conception d'un générateur de règles adaptées au besoin correspond à l'idée de schèmes générateurs de règles orientées selon un objectif interne. Mais ici, le principe d'anticipation et les méthodes de génération de règles font partie de la structure du système alors que, dans le cadre d'une architecture proposée, ces capacités relèvent des schèmes cognitifs qui peuvent être évalués ou modifiés. Les mécanismes d'anticipation sensorielle apparaissent alors à travers les schèmes cognitifs ou à travers les primitives sensorimotrices de manière rudimentaire et figée. Enfin, sans mémoire à court terme ou étiquette temporelle, ces systèmes de classeurs se limitent toujours à la prédiction au pas de temps suivant, ce qui peut être jugé limitatif pour l'anticipation comprise comme imagination d'une chaîne d'actions à pas de temps variables.

4.5. Les systèmes de classeurs hiérarchiques ou motivationnels

La dernière catégorie de systèmes de classeurs peut se comprendre comme l'ensemble des classeurs abordant des situations dont la traduction logique est non monotone. Ces situations peuvent être dues à des traitements à des échelles de temps différents ou à des données hétérogènes dépendantes entre elles. Une façon d'aborder ces problèmes complexes consiste à identifier les buts et à les hiérarchiser afin de proposer des solutions adaptées. Dans ce cadre, chaque sous-système de classeurs correspond à l'une des

solutions adaptées. En définitive, un système de classeurs hiérarchiques se compose d'un module de gestion ou de paramètres internes jouant le rôle d'un multiplexeur pour l'acquisition d'un sous-système de classeurs (Figure III-17).

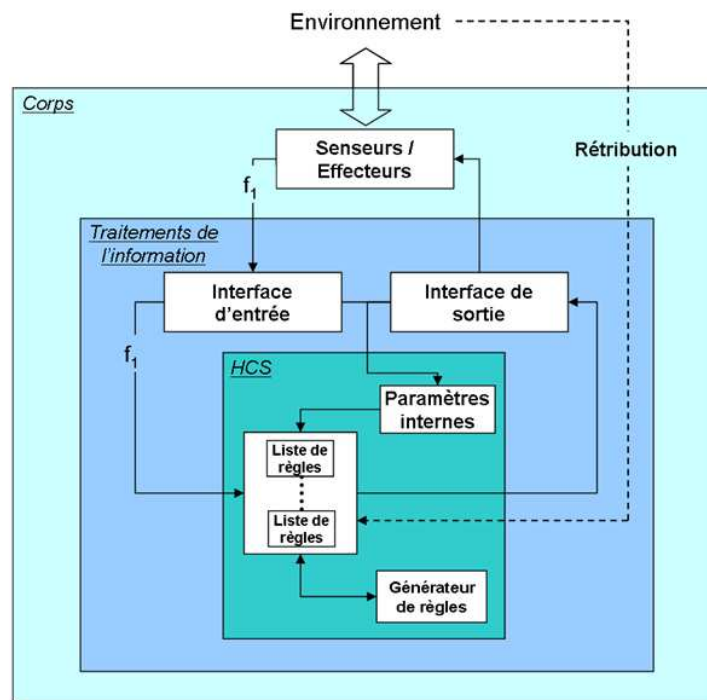


Figure III-17 : Schéma d'un système de classeurs de type HCS

La gestion des paramètres internes peut être également un système de classeurs comme dans le modèle ALECSYS (*A Learning Classifier SYStem*) (Dorigo, 1995) où un système de classeurs directeur choisit l'activation des autres sous-systèmes de classeurs qui apprennent le comportement de base pour une souris artificielle : suivre une source lumineuse mobile et retourner dans son terrier à l'apparition d'un prédateur. Une grande variété de systèmes existe (Buche, 2006) avec leur propre particularité, par exemple les HGCS (*Heterogeneous Genetic Classifier System*) (Sanchez, 2004) qui proposent de remplacer la gestion de paramètres par une classification des données en types distincts ou le MHICS (*Modular Hierarchy Classifier system*) (Robert, 2005) qui assemble des modules distribués sur plusieurs niveaux afin de gérer les motivations et les priorités. Mais les systèmes de classeurs hiérarchiques ne proposent pas de solution générique pour les problèmes nécessitant une logique de description non monotone. Le module de gestion se trouve en dehors de la dynamique de la mémoire à long terme alors qu'elle en résulte pour l'architecture cognitive proposée.

**

En conclusion, malgré leur hétérogénéité, les systèmes de classeurs conservent tous l'idée d'une récompense extérieure définissant indirectement la problématique à laquelle se confronte l'agent. En résumé, la récompense correspond à une donnée objective, c'est-à-dire à une référence réelle. Parallèlement, les systèmes de classeurs proposent une dynamique de sélection et une heuristique permettant la découverte de relations. Toute la difficulté se trouve dans la canalisation de l'heuristique afin de ne pas se perdre dans l'espace des possibles et s'orienter vers une solution viable.

Un système réellement autonome devrait orienter lui-même son heuristique mais cela impliquerait une représentation de la rétribution et du but. Or, la représentation d'une rétribution ou du but soumis à la dynamique d'un système de classeurs, donc modifiable, n'est pas raisonnable selon une conception réaliste puisque la représentation de la rétribution signifierait rajouter de l'erreur inutilement puisque cela reviendrait à estimer une information qui pourrait directement être obtenue. Autrement dit, la rétribution objective étant connue, il serait contre-productif de la rendre subjective. L'emploi de la méthode scientifique MHD renforce cette position en souhaitant dérouler au fur et à mesure les hypothèses et les difficultés, c'est-à-dire réfléchir dans un premier temps à partir d'un cas simple dans lequel le bruit se limite au capteur puis dans un second temps complexifier en rajoutant du bruit sur la récompense.

La conséquence de cette méthode est que la première étape ne peut se faire qu'en extériorisant la problématique de la dynamique des classeurs, soit la canalisation de l'heuristique. En d'autres termes, chaque problème ou type de problèmes nécessite une modification de l'architecture du système des classeurs. Les systèmes de classeurs développés pour des situations simples laissent croire que la simplicité de la problématique révèle des fondamentaux de l'apprentissage. Mais en fait, ces fondamentaux portent davantage sur le type de situation à apprendre que sur l'apprentissage lui-même. Par ailleurs, l'apprentissage d'une nouvelle variante de la situation nécessitera des modifications de l'architecture adaptée pour la situation archétype. Tout cela explique la grande diversité des systèmes de classeurs et leurs bonnes performances dans leurs domaines d'applications. De plus, rien ne présage que le passage à la seconde étape abordera des problèmes d'apprentissages plus généraux puisque le concepteur s'arrangera pour que le système approxime la représentation de la récompense par rapport à sa référence « objective ».

Toutefois, cette impasse reste invisible tant que la définition de l'apprentissage repose sur une conception réaliste. L'introduction de l'anticipation dans les systèmes de classeurs illustre bien l'invisibilité de l'impasse. Le critère de rétribution à partir de l'anticipation permet effectivement d'internaliser la provenance de la récompense au sein de l'agent. Cependant, les processus d'évaluation et de rétribution demeurent à l'extérieur de la dynamique des classeurs, même si cette dynamique y contribue, et adaptés à certains types de situations. De même, la rétribution interne native dépendant uniquement du déclenchement des règles sans autre critère d'évaluation reste orientée vers la problématique extérieure consistant à dégager des chaînes d'actions.

L'architecture cognitive souhaite sortir de cette impasse en éliminant les spécifications qui reposent sur des problématiques faussement constitutives comme l'anticipation ou le chaînage et en rajoutant les mécanismes nécessaires à leur émergence au sein de la sémiose sous la forme de schèmes cognitifs.

5. Spécification de l'architecture cognitive

Avant d'aborder en détail la spécification de l'architecture cognitive proposée, la première grande différence avec les systèmes précédents réside, d'un point de vue général, dans le fait que les interfaces d'entrées et de sorties correspondent maintenant à des primitives sensorielles et motrices qui sont des processus à part entière. De même, la mémoire événementielle et la mémoire épistémique ne représentent plus de simples procédures mais des tâches. En d'autres termes, les composants des systèmes de classeurs classiques sont synchrones alors que l'architecture cognitive et les primitives sont

totallement asynchrones. Plus exactement, en termes informatiques, l'architecture cognitive peut alors être considérée comme un logiciel servant d'intermédiaire entre des applications complexes distribuées (les primitives), soit un intergiciel (*middleware*). Les primitives seront implémentées sous forme de tâches informatiques ainsi que tous les composants constituant l'architecture cognitive comme les mémoires ou les processus subcognitifs annexes (mécanisme de création de règles par exemple). Cependant, dans cette section, la spécification présentée porte uniquement sur l'architecture cognitive, un exemple de primitives sensorielle et motrice sera avancé dans le quatrième chapitre pour l'évaluation de l'architecture cognitive.

La présentation de la spécification s'effectuera en trois temps. Dans un premier temps, les éléments de base de l'architecture seront définis : les messages et les règles. La gestion du temps de ceux-ci sera expliquée lors de la présentation des deux types de mémoires. Dans un deuxième temps, les mécanismes de sélection de règles qui définissent et façonnent leur contenu seront abordés à la fois pour la reconnaissance spontanée et pour la reconnaissance attendue. Dans un troisième temps, la dynamique de la rétribution de règles sera expliquée. Enfin, quelques remarques sur l'implémentation de cette spécification seront évoquées.

5.1. La structure générale

5.1.1. Les éléments constitutants

Toute l'architecture se fonde principalement sur deux hiérarchies d'objets à savoir (A) les messages et (B) les règles. L'étendue de l'arbre de dérivation pour chacune de ces hiérarchies dépend à la fois du nombre de primitives et de la nature des schèmes cognitifs.

A - Les messages

Le message correspond à la brique élémentaire du système cognitif, il est considéré comme la base d'un signe, d'une condition d'une prémisse ou d'une conclusion bien que leur rôle diffère. En plus du typage de bas niveau évoqué dans le cadre de la formalisation, celui des messages repose sur la distinction entre deux catégories de signes : les signes qui représentent un constat et ceux qui représentent une intention. Les signes de la première catégorie proviennent soit des primitives sensorielles, soit des interprétations de type jugement, et ceux de la seconde catégorie proviennent des interprétations de type intentionnel.

La prédiction se traduit par un schème cognitif qui s'appuie sur un opérateur de reconnaissance attendue, et ne constitue donc pas une catégorie de messages. La prédiction, schème cognitif fondamental, sera présentée sommairement en fin de partie 3.1.3 puis sera détaillée au cours de la présentation des opérateurs de reconnaissance et des mécanismes de rétribution. L'anticipation, schème cognitif indispensable pour la cognition d'ordre supérieur, ne sera pas totalement spécifiée mais quelques pistes seront avancées dans la quatrième partie de ce chapitre.

Les deux catégories de messages retenues sont : la catégorie des messages constatés (CMC) et la catégorie des messages intentionnels (CMI). La distinction entre CMC et CMI permet d'introduire la réflexivité du système sur sa propre activité, critère fondamental pour la sémiologie. Dans la mémoire événementielle, un message qui représente une intention doit se distinguer d'un message qui représente une intention réalisée, ce qui permet de tenir

compte de ce qui va se faire et de ce qui a été fait. Il n'existe pas de type message CMI sans son homologue en CMC. La gestion de ces couples de types qui sera explicitée lors de la présentation de la mémoire événementielle s'appuie sur l'étiquette temporelle θ des messages.

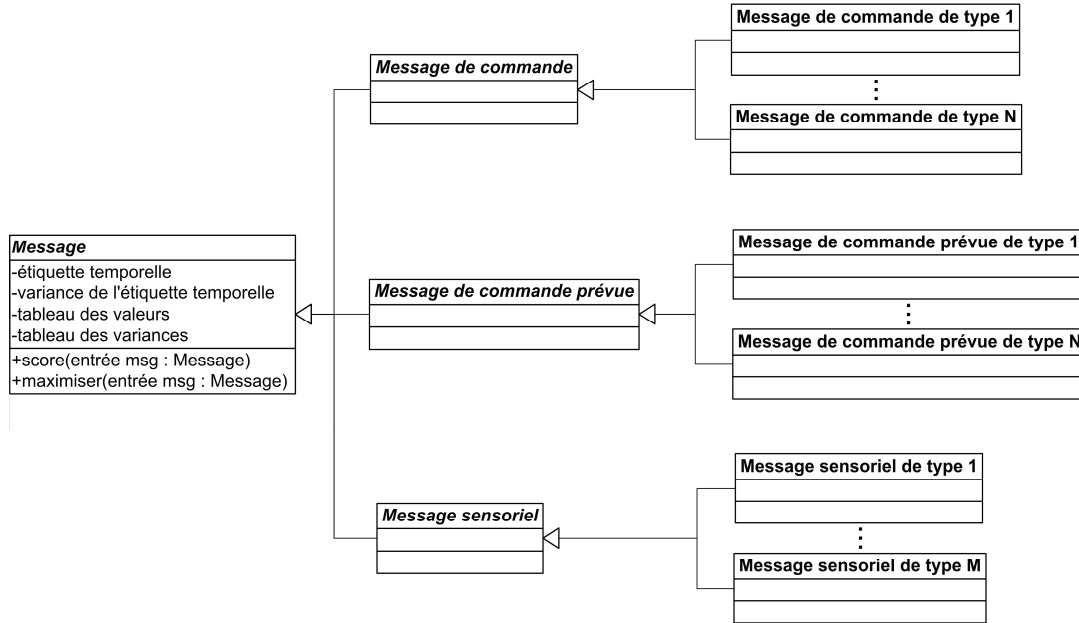


Figure III-18 : Schéma UML simplifié de la hiérarchie minimale de l'objet « message ».

Dans ce contexte, à chaque primitive sensorielle ainsi qu'à chaque sorte de rétribution correspond un type de message appartenant à la catégorie des signes comme un constat (CMC). À chaque primitive motrice est associé à la fois un type de message de catégorie CMC ainsi qu'un type de message de catégorie CMI. La définition des types de messages internes dépend des schèmes cognitifs développés. Par exemple, en prenant en compte uniquement les primitives sensorielles et motrices, la Figure III-18 représente un schéma de spécification UML simplifié de la hiérarchie minimale de l'objet message. L'utilité des fonctions « score » et « maximiser » sera abordée lors de l'explication du mécanisme de sélection d'une règle.

Concrètement, un message représente un tableau de valeurs dont la taille est définie par son type et dont la première valeur correspond à l'étiquette temporelle $\theta(t)$ et la seconde à la légitimité $L(t)$ qui représente les notions hyponymes issues de la définition de l'architecture cognitive et de sa formalisation. Pour des raisons de normalisation, les valeurs des tableaux sont des réels entre 0 et 1 (Figure III-19). L'utilisation de valeurs continues évacue le problème de la granularité et de l'adaptabilité de grain. En effet, ce sont les propriétés d'auto-organisation de l'architecture cognitive qui doivent permettre au système d'aboutir à une segmentation spatio-temporelle du monde, contrairement aux systèmes de classeurs classiques.

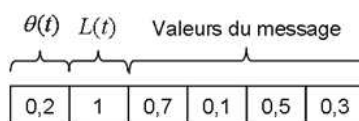


Figure III-19 : Exemple de contenu d'un message dont la longueur dépend de son type.

L'étude sur l'interprétation d'un signe a montré que la condition d'une prémisse peut se comprendre comme une classe d'équivalence qui évolue en fonction des autres. Par ailleurs, la sémiose ne devant pas contraindre la segmentation des primitives, la valeur des messages est considérée comme continue. Dans ce cadre, la représentation choisie pour modéliser ces classes d'équivalence est celle des noyaux gaussiens. Ainsi, toutes les valeurs des messages constituant les prémisses reçoivent une incertitude σ (Figure III-20). La valeur d'incertitude σ correspond à la variance d'une loi gaussienne. Par rapport aux systèmes de classeurs traditionnels, l'incertitude proposée renvoie à la notion de joker avec le caractère #, sauf que la généralisation et la spécialisation induites ici par l'incertitude sous forme de gaussienne permet une plus grande précision dans l'évaluation de l'appariement. Par ailleurs, l'incertitude rend compte à la fois de la distribution des signes et du bruit des processus qui les a générés.

1.10^{-4}	2.10^{-2}	7.10^{-3}	1.10^{-2}	9.10^{-2}	8.10^{-3}	} Valeurs d'incertitudes } Contenu du message
0,2	1	0,7	0,1	0,5	0,3	

Figure III-20 : Exemple de message appartenant à une prémisse.

Par rapport à la notion de probabilité imprécise abordée lors de la formalisation, la solution adoptée appartient aux probabilités subjectives. Une prémisse représente un ensemble de densités de probabilités. Le fait que la légitimité possède également une incertitude permet d'introduire une souplesse sur la propagation des conclusions. En considérant que la légitimité représente une probabilité, le noyau gaussien de la légitimité au sein d'une condition représente une densité de probabilité d'une probabilité.

B - Les règles

Il existe deux sortes de règles selon que leur déclenchement s'effectue par l'opérateur de reconnaissance spontanée ou par celui de reconnaissance attendue. Cependant, ces deux sortes de règles possèdent une structure identique, elles se composent d'une prémisse générale et d'une conclusion. La prémisse générale renferme une prémisse excitatrice et une prémisse inhibitrice. Chacune d'elles peut contenir un nombre quelconque de messages et de n'importe quel type. La conclusion correspond à un message CMC ou à une liste de messages CMI d'un seul type avec des étiquettes temporelles différentes :

$$\left\langle \left\{ \vee_i M_i \right\} \right\rangle_I \wedge \left\langle \left\{ \vee_j M_j \right\} \right\rangle_E \supset \left\{ M_k \right\}^n$$

Les messages d'un type relatif à une primitive sensorielle ne peuvent pas être des conclusions. Le type des conclusions peut appartenir à la catégorie CMC uniquement si aucun type homologue CMI n'existe. En d'autres termes, toute conclusion représente soit un jugement, soit une intention. Le type d'une règle est défini par le type de message de la conclusion. Par conséquent, le nombre de types de règles correspond au nombre de types de messages de catégorie CMC en dehors de ceux relatifs aux primitives sensorielles. Ce mécanisme offre la possibilité d'avoir une réflexivité sur ses actes pour ensuite pouvoir prédire, voire anticiper, et surtout de répondre à la notion de champ abordée lors de la question sur l'univocité de la perception, ainsi que d'assurer l'intégrité de la mémoire événementielle (aucune règle ne peut conclure un message sensoriel).

Chaque règle possède quatre attributs : la *force*, l'*évaluation*, l'*alpha* et l'*étiquette temporelle*. La *force* reflète sa pertinence globale au cours de l'activité du système, traduisant d'une certaine manière l'espérance de rétribution une fois déclenchée. L'*évaluation* correspond à l'adéquation entre la prémisse de la règle et les signes courants en fonction des autres règles. Le calcul de cette adéquation ainsi que l'attribut *alpha* seront détaillés dans la partie 3.2. L'*étiquette temporelle* donne une information temporelle sur le dernier déclenchement de la règle qui permet de définir la portée temporelle d'une rétribution.

La Figure III-21 représente le schéma de spécification UML simplifié des règles construit à partir de celui des messages de la Figure III-18.

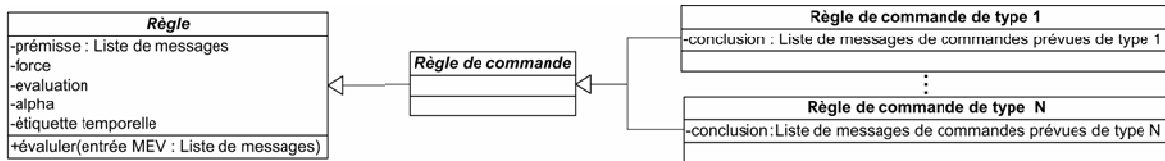


Figure III-21 : Schéma UML simplifié d'une hiérarchie minimale de l'objet « règle ».

5.1.2. Les deux mémoires

L'architecture cognitive repose sur deux mémoires : (A) la mémoire événementielle et (B) la mémoire épistémique.

A - La mémoire événementielle

La mémoire événementielle est exclusivement constituée de messages provenant soit des primitives sensorielles, soit de sa propre gestion des messages, soit des conclusions de règles. Les types des messages provenant des deux premières sources, qui appartiennent obligatoirement à la catégorie CMC contrairement à ceux provenant de la troisième source, appartiennent à la catégorie CMC ou CMI.

Une fois dans la mémoire événementielle, un message représente un signe, quel que soit sa catégorie. Sans introduire explicitement une notion de temps en utilisant une représentation linéaire de celui-ci avec une date t à partir d'une référence extérieure qui viendrait à l'encontre d'une démarche constructiviste, l'étiquette temporelle de ce signe doit toutefois correspondre à une représentation qui relie implicitement le temps physique et le temps représenté. Par ailleurs, l'analyse sur la perception des illusions ou hallucinations montre la nécessité de posséder le moyen d'évaluer une distance temporelle entre le présent et le passé ou le futur.

Afin que cette représentation soit bornée entre 0 et 1, le modèle empirique d'écoulement du temps choisi consiste à diminuer l'étiquette temporelle θ des messages se trouvant dans la mémoire événementielle selon la relation avec λ un réel positif fixé au début de la sémiiose et commun à tous les messages, en notant t_0 le temps absolu d'arrivée du message :

$$\theta(t) = e^{-\lambda(t-t_0)}$$

Lorsque le type du message arrivant appartient à la catégorie CMC :

$$\theta(t = t_0) = \theta_0 = 1$$

Le λ étant positif, l'étiquette temporelle décroît au cours du temps et, arrivé à un certain seuil, le message sera éliminé afin de ne pas surcharger la mémoire événementielle. Toutefois, rien n'empêche d'imaginer un schème cognitif de la mémoire épistémique qui, selon ces critères, recopierait certains signes dans une mémoire à long terme où ils continueraient ainsi leur décroissance exponentielle. Cependant, cette considération reste une perspective annexe puisque l'architecture cognitive minimale repose fondamentalement sur l'interaction entre la mémoire événementielle et la mémoire épistémique.

Pour les messages d'un type de catégorie CMI, l'étiquette temporelle représente une prévision à $t_p > t_0$ qui se traduit par :

$$\theta(t) = e^{-\lambda((t_p - t_0) - (t - t_0))} = e^{-\lambda(t_p - t)}$$

À l'arrivée de ces messages, l'étiquette temporelle devient :

$$\theta(t = t_0) = \theta_0^p = e^{-\lambda(t_p - t_0)}$$

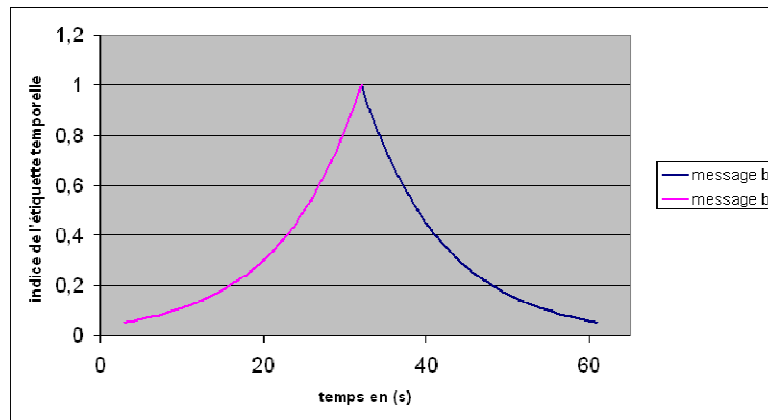


Figure III-22 : Le graphique représente l'évolution de l'étiquette temporelle d'un message b' de catégorie SI programmant sa transformation 30 s après son arrivée en message b selon le type associé, donc de catégorie CMC, et la décroissance de l'étiquette temporelle du message b jusqu'à son élimination.

Dans ce cas, l'étiquette temporelle des messages de catégorie CMI augmente progressivement vers 1 et, arrivée à cette valeur, ils sont convertis en messages de type de catégorie CMC qui leur est associé. Cette conversion correspond à la deuxième source évoquée en introduction de la mémoire événementielle. Les messages convertis étant devenus des messages d'un type de catégorie CMC, leur étiquette temporelle qui était à 1 décroît au cours du temps (Figure III-22).

Cette représentation du temps permet d'une part une symétrie entre la représentation temporelle des messages reliés au passé et ceux reliés au futur et d'autre part elle permet de dater potentiellement tous les événements passés et futurs dans un intervalle entre 0 et 1 puisque la limite à l'infini positif de t de la décroissance exponentielle est 0.

À noter que cette représentation du temps implique que la précision temporelle entre deux événements soit relative à leur éloignement au présent, autrement dit, la différence entre deux étiquettes temporelles paraîtra plus grande si celles-ci sont proches du présent que si elles font référence à des événements lointains dans le passé ou à des prévisions éloignées dans le futur (Figure III-23) avec $t_1 > t_2 > 0$:

$$\Delta\theta = e^{-\lambda(t_1)} - e^{-\lambda(t_2)}$$

$$\Delta\theta = e^{-\lambda(t_2)} \left[e^{-\lambda(t_1-t_2)} - 1 \right]$$

$$\Delta\theta = e^{-\lambda(t_2)} \left[e^{-\lambda(\Delta t)} - 1 \right]$$

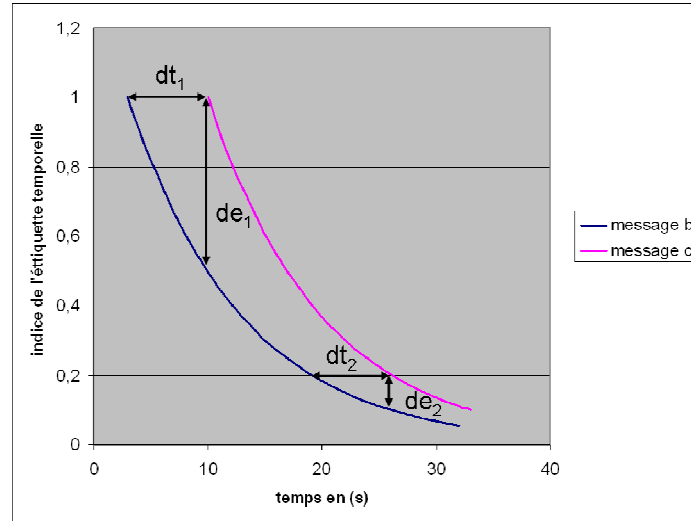


Figure III-23 : Exemple de deux messages arrivant successivement dans la mémoire événementielle, la distance temporelle selon leur étiquette temporelle diminue au cours du temps ($de_1 > de_2$), contrairement à la distance temporelle objective en s qui reste constante ($dt_1 = dt_2$).

Lors de la transformation d'un message d'un type appartenant à la catégorie CMI en son type associé de catégorie CMC, si le type du message est rattaché à une primitive motrice, alors le contenu du message est envoyé à cette primitive. Si le type du message est rattaché à une rétribution, alors son contenu est envoyé à la mémoire épistémique, si le type du message est rattaché à un processus subcognitif lié à un schème, alors son contenu est renvoyé vers ce processus, enfin, si le type de message n'est relié à aucune activité de ce genre, la transformation se déroule normalement.

Lorsqu'un nouveau message ou une liste de nouveaux messages d'un type appartenant à la catégorie CMI arrive, tous les messages appartenant à ce type sont éliminés. Ainsi, l'intégrité de la mémoire événementielle est assurée, et de manière générale aucun message de même type ne peut recevoir la même étiquette temporelle.

B - La mémoire épistémique

La mémoire épistémique contient l'ensemble des règles du système. Chaque règle possède une étiquette temporelle qui décroît selon le même principe que les messages :

$$\theta(t) = e^{-\lambda(t-t_0)}$$

Mais ici, t_0 correspond au dernier déclenchement de la règle. Initialement, t_0 est considéré comme très éloigné dans le passé, c'est-à-dire que toutes les étiquettes temporelles sont proches de zéro.

Contrairement à la mémoire événementielle dont l'activité est rythmée par l'arrivée de messages, la mémoire épistémique possède un rythme propre. À chaque cycle de fonctionnement, la mémoire épistémique interprète la mémoire événementielle c'est-à-dire sélectionne pour chaque type la règle ayant la meilleure évaluation. Si toutes les règles d'un même type ont une évaluation nulle, alors aucune règle de ce type n'est sélectionnée. La sélection d'une règle entraîne l'envoi d'une copie de sa conclusion à la mémoire événementielle et l'ajustement de sa prémisse (valeurs et incertitudes) avec les signes qui lui ont permis d'être sélectionnée. L'ajustement permet de faire en sorte que l'interprétation des signes se façonne dynamiquement par rapport au vécu. La sélection d'une règle a pour effet également de réévaluer sa force représentant l'évaluation globale de la règle et le paramètre alpha dont la signification sera donnée ultérieurement. À la fin du cycle, les rétributions sont délivrées, s'il y a lieu, et les règles dont la force reste inférieure à un certain seuil sont éliminées. Parallèlement à ce cycle, de nouvelles règles provenant d'un processus subcognitif initié ou non par l'activité de certaines règles peuvent enrichir la mémoire épistémique.

5.1.3. Les schèmes cognitifs fondamentaux

Sans détailler le mode de sélection d'une règle et des mécanismes associés, une présentation générale (A) des schèmes sensorimoteurs liés à la gestion de l'action et (B) des schèmes prédictifs liés à la gestion de la prédiction peut être avancée et servir d'exemple pour mieux intégrer les notions précédentes.

A - La gestion de l'action

Le schème sensorimoteur représente l'interface entre les primitives sensorielles et les primitives motrices. Les primitives pouvant posséder leur propre boucle sensorimotrice, il convient de distinguer les boucles sensorimotrices primaires invariables issues des primitives et celles d'ordre supérieur modifiables, issues des schèmes sensorimoteurs. Plus précisément, il sera successivement présenté (i) les schèmes secondaires simples, (ii) les schèmes secondaires complexes et (iii) les schèmes de troisième ordre et plus. Les schèmes sensorimoteurs secondaires peuvent être soit simples, soit complexes, mais dans les deux cas, les *règles de commande* possèdent des *messages sensoriels* dans la prémisse. Les schèmes sensorimoteurs d'ordre supérieur renvoient à des relations sensorimotrices médiatisées par des messages internes intermédiaires, autrement dit nécessitant plusieurs inférences pour relier la sensation à une action.

i - Les schèmes secondaires simples

Un schème sensorimoteur secondaire simple correspond à un ensemble de *règles de commande* dont la prémisse n'est constituée que de *messages sensoriels* (S) associés à des types de primitive sensorielle n différents et la conclusion est une liste de *messages de commande prévue* (CP) associée à un type de primitive motrice m . En prenant pour notation les accolades pour représenter un groupe de messages formant un ensemble de conditions soit conjonctives lorsque les conditions se trouvent dans la prémisse excitatrice, soit disjonctives lorsqu'elles se trouvent dans la prémisse inhibitrice, un schème sensorimoteur simple prend la forme suivante, avec entre parenthèses, l'étiquette temporelle :

$$\left\langle \left\{ S_i^n(\theta_i) \right\}_{\forall n \in N} \right\rangle_I \wedge \left\langle \left\{ S_j^n(\theta_j) \right\}_{\forall n \in N} \right\rangle_E \supset \left\{ CP_k(\theta_k) \right\}_{m \in M}$$

Cette configuration est très proche des classeurs classiques mais pour que les mécanismes soient équivalents, il faut que le seuil de la mémoire événementielle soit paramétré de façon à ne contenir qu'un seul message par type, que la cadence d'envoi de messages de la primitive sensorielle soit identique à la fréquence de sélection de la mémoire épistémique. Concernant la forme de règle, la prémisse inhibitrice doit être vide et la prémisse excitatrice avec la conclusion doit être mono message avec un retard temporel nul. Autrement dit, le message CP découle directement de la sensation présente et se réalise dès qu'il est introduit dans la mémoire événementielle (Tableau III-2) :

$$\left\langle S_j^n(\theta_0) \right\rangle_E \supset CP_k^m(1)$$

t	messages sensoriels	messages de commande prévue	messages de commande
1	S ₁ (1)		
2	S ₂ (1)	CP _a (1) →	C _a (1)
3	S ₃ (1)	CP _b (1) →	C _b (1)

Tableau III-2 : Exemple de trois états successifs de la mémoire événementielle dont le seuil d'oubli se trouve juste en dessous de 1, ce qui signifie que la mémoire événementielle contient au plus un message de chaque type. A côté des messages est indiquée entre parenthèses la valeur de leur étiquette temporelle. Ici la primitive sensorielle est synchrone avec la mémoire épistémique. Le message S₁(1) est à l'origine du déclenchement de la règle A qui se traduit par l'introduction au pas suivant du message de commande prévue CP_a avec un retard nul, ce qui implique l'envoi immédiat à la primitive motrice correspondante de la commande et la conversion du message en message de commande.

Cependant, cette gestion simple ne permet pas de faire la différence entre une action finie (comme prendre un verre d'eau) et une action avec une fin indéterminée (comme marcher jusqu'à trouver de l'eau) puisque l'action est systématiquement remise en question par l'état de la mémoire événementielle, aucune planification n'est possible. Par ailleurs, cette gestion simple n'exploite ni l'utilisation des étiquettes temporelles ni la possibilité de prévoir une chaîne de commandes, autrement dit, cette gestion de l'action n'est efficace que dans un environnement markovien. En baissant le seuil de la mémoire événementielle, une prémisse peut contenir une chaîne de messages sensoriels, ainsi, chaque action sélectionnée repose sur un certain nombre d'états passés (Tableau III-3). La forme des règles devient alors :

$$\left\langle \left\{ S_j^n(\theta_j) \right\} \right\rangle_E \supset CP_k^m(1)$$

Mais, il est impossible a priori de connaître le nombre de messages S à mémoriser pour se retrouver dans une situation markovienne.

t	messages sensoriels	messages de commande prévue	messages de commande
1	S ₁ (1)		
2	S ₂ (1) S ₁ (0,95)		
3	S ₃ (1) S ₂ (0,95)	CP _a (1) →	C _a (1)
4	S ₄ (1) S ₃ (0,95)	CP _b (1) →	C _b (1) C _a (0,95)

Tableau III-3 : Exemple de quatre états successifs de la mémoire événementielle dont le seuil d'oubli se trouve juste en dessous de 0,95, ce qui signifie que la mémoire événementielle contient au plus deux messages de chaque type. A côté des messages est indiquée entre parenthèses la valeur de leur étiquette temporelle. Ici, la primitive sensorielle est synchrone avec la mémoire épistémique. Les messages S1(0,95) et S2(1) sont à l'origine du déclenchement de la règle A qui se traduit par l'introduction au pas suivant du message de commande prévue CP_a avec un retard nul, ce qui implique l'envoi immédiat à la primitive motrice correspondante de la commande et la conversion du message en message de commande.

ii - Les schèmes sensorimoteurs secondaires complexes

Ces limitations peuvent être dépassées par les schèmes sensorimoteurs secondaires complexes qui utilisent des messages internes pour autoriser ou non une nouvelle sélection de l'action. Par exemple, un schème sensorimoteur secondaire complexe élémentaire fait intervenir deux nouveaux types de messages, l'un appartenant à la catégorie CMC et l'autre, son homologue, appartenant à la catégorie CMI : les types de message GC et les types de messages (GPC). Une nouvelle sorte de type de règles apparaît alors : les types de règles GC dont la conclusion ne possédera qu'un seul message GPC.

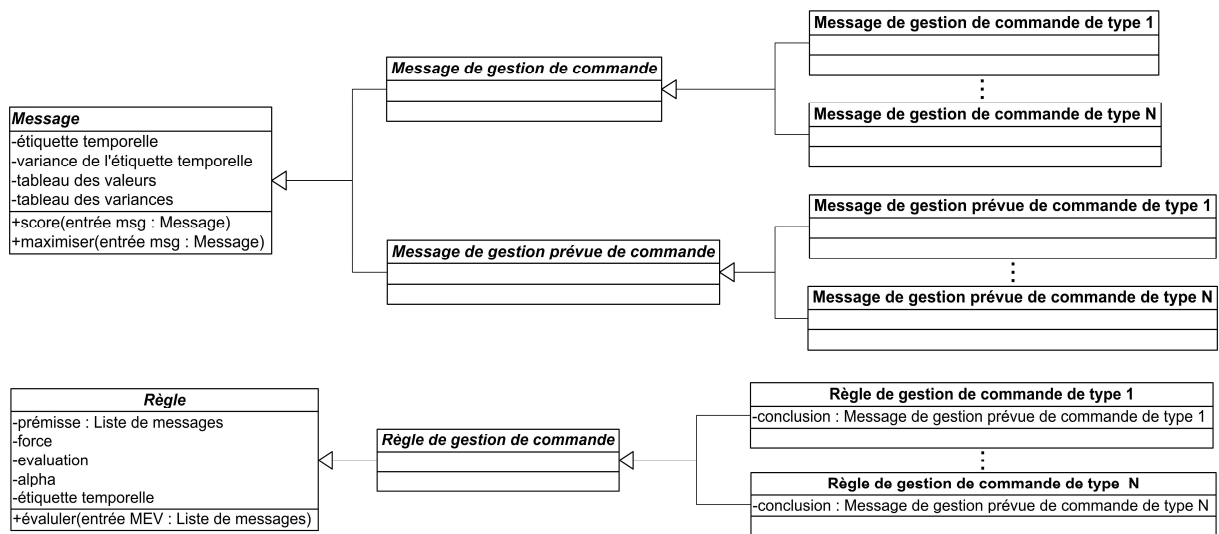


Figure III-24 : Schéma UML simplifié des objets « règle » et « message » liés à la gestion de l'action des N primitives motrices.

Naturellement, les types de messages GC, les types de messages GPC et les types de règles GC sont relatifs aux primitives motrices (Figure III-24). Mais pour alléger le discours, de manière générique, il ne sera évoqué que le type spécifique des messages et des règles en considérant qu'ils sont implicitement homogènes dans le cadre d'un schème sensorimoteur lié à une primitive motrice.

L'introduction de ces nouveaux types de messages et de ce nouveau type de règle permet, pour un schème donné, de coupler les règles sensorimotrices de la façon suivante :

$$\left[\begin{array}{l} \langle \forall GPC^m \rangle_I \wedge \left\langle \left\{ S_j^n(\theta_j) \right\}^{\forall n \in N} \right\rangle_E \supset \left\{ CP_k(\theta_k) \right\}^{m \in M} \\ \langle \forall GPC^m \rangle_I \wedge \left\langle \left\{ S_j^n(\theta_j) \right\}^{\forall n \in N} \right\rangle_E \supset GPC_k^m(\theta_r) \end{array} \right]_{x \in \text{schème}^S}$$

La traduction du symbole \forall (quel que soit) s'effectue par les incertitudes liées aux prémisses et sera expliquée plus en détails lors de la présentation du mécanisme de l'évaluation des règles.

Dans le cadre du couplage des règles, un message GPC est envoyé en même temps que le message CP. Un message de gestion prévue de commande peut avoir un contenu et par conséquent potentiellement avoir une certaine signification pour certaines règles, excepté pour les règles de commandes et celles de gestion de commande constituant le schème sensorimoteur. En effet, pour que le message GPC joue le rôle d'inhibiteur, les prémisses inhibitrices des règles constituant le schème sensorimoteur doivent avoir une incertitude maximale pour le message GPC. Ainsi, tout message GPC dans la mémoire événementielle inhibe toutes les règles du schème sensorimoteur en provoquant une évaluation nulle. Ce mécanisme laisse le temps aux messages CP de se dérouler jusqu'à ce que l'étiquette temporelle du message GPC arrive à 1. Un exemple d'un tel mécanisme peut être présenté avec un schème sensorimoteur (reliant une primitive sensorielle à une primitive motrice) formé de deux couples de règles de la manière suivante (pour alléger l'écriture, la référence au type des primitives n'est pas indiquée) :

$$\left[\begin{array}{l} \langle \forall GPC \rangle_I \wedge \langle S_1(1) \rangle_E \supset \{CP_{a1}(1); CP_{a2}(0,95)\} \\ \langle \forall GPC \rangle_I \wedge \langle S_1(1) \rangle_E \supset GPC_a(0,95) \end{array} \right]_a$$

$$\left[\begin{array}{l} \langle \forall GPC \rangle_I \wedge \langle S_2(1) \rangle_E \supset \{CP_b(1)\} \\ \langle \forall GPC \rangle_I \wedge \langle S_2(1) \rangle_E \supset GPC_b(0,95) \end{array} \right]_b$$

Le Tableau III-4 montre un exemple de l'incidence de ces règles sur la mémoire événementielle dans un contexte sensoriel oscillant entre S_1 et S_2 .

Cette latence entre la sélection de deux règles peut correspondre au fait que certaines actions déclenchées se déroulent sans boucle cognitive de retour (Keele, 1968). Toutefois, une intervention extérieure notable peut nécessiter de changer d'action en cours. Ces cas se traduisent alors par des règles de GC dont la conclusion possède un retard nul et dont la prémisse inhibitrice ne contient pas de message GPC. Ainsi, le déclenchement d'une telle règle n'est pas conditionné par la gestion de la commande en cours, tout en pouvant l'annuler en envoyant un message GPC qui se réalise aussitôt afin de laisser s'accomplir la sélection d'une règle sensorimotrice.

Sur le même principe, une action indéfinie se traduit alors par un message de GPC avec un retard infini qui est annulé uniquement par la réception d'un nouveau message GPC issu d'une règle de GC non couplée avec une règle de CP.

t	messages sensoriels	messages de commande prévue	messages de commande	messages de gestion prévue de commande	messages de gestion de commande
1	$S_1(1)$				
2	$S_2(1)$ $S_1(0,95)$	$CP_{a1}(1)$ → $CP_{a2}(0,95)$	$C_{a1}(1)$	$CPG_a(0,95)$	
3	$S_1(1)$ $S_2(0,95)$	$CP_{a2}(1)$ →	$C_{a2}(1)$ $C_{a1}(0,95)$	$CPG_a(1)$ →	$CG_a(1)$
4	$S_2(1)$ $S_1(0,95)$		$C_{a2}(0,95)$		$CG_a(0,95)$
5	$S_1(1)$ $S_2(0,95)$	$CP_{b1}(1)$ →	$C_{b1}(1)$	$CPG_b(0,95)$	
6	$S_2(1)$ $S_1(0,95)$		$C_{b1}(0,95)$	$CPG_b(1)$ →	$CG_b(1)$

Tableau III-4 : Exemple de six états successifs d'une mémoire événementielle avec un schème sensorimoteur secondaire complexe. À $t=1$, S_1 provoque le déclenchement du couple de règles a , les deux temps suivant le message CPG_a empêche tout déclenchement de règle. À $t=4$, le couple de règles b est sélectionné et inhibera la sélection d'une nouvelle règle jusqu'au temps 6 inclus.

iii - Les schèmes cognitifs de troisième ordre et plus

Les schèmes de troisième ordre correspondent aux schèmes sensorimoteurs dont les règles de commandes ne possèdent pas de messages sensoriels, ce qui signifie que leur déclenchement provient d'un jugement ou d'une hiérarchie intentionnelle d'action. En effet, la perception du temps suggère une hiérarchie de représentations temporelles, les comparaisons s'effectuant à l'intérieur de chaque niveau hiérarchique. À cette hiérarchie temporelle correspond une hiérarchie d'actions, les actions d'un niveau étant des combinaisons d'actions du niveau sous-jacent. Il est possible de représenter cette hiérarchie d'actions comme un complexe simplicial dont les nœuds élémentaires correspondent aux commandes élémentaires. Ainsi, les actions de premier niveau correspondent à une suite de commandes et les actions de niveau n correspondent à une suite d'actions de niveau $n-1$. À cette hiérarchie d'action doit être associée celle des règles permettant leur gestion. En ne considérant que les règles de commandes, chaque niveau posséderait son type de message de gestion de commande hiérarchique (GCH) et sa règle hiérarchique de commande telle que par exemple :

$$\langle GCH_i^m(\theta_0) \rangle_E \supset \left\{ CP_k(\theta_k) \right\}^m$$

$$\langle GCH_j^m(\theta_0) \rangle_E \supset \left\{ GCH_k^m(\theta_k) \right\}_{k \notin j}^m$$

B - La gestion de la prédiction

Dans la prédiction, il convient de distinguer deux aspects (i) l'acte de prédire et (ii) l'utilisation ou la vérification de la prédiction.

i - L'acte de prédire

Les schèmes prédictifs proposent uniquement de programmer la reconnaissance d'un message d'un type appartenant à la catégorie CMC n'ayant pas d'homologue CMI. Cette programmation repose sur l'utilisation de deux types de messages intentionnels, comme le montre par exemple la Figure III-25. L'utilisation de ce message programmé ou de cette prédiction s'effectuera par d'autres schèmes cognitifs.

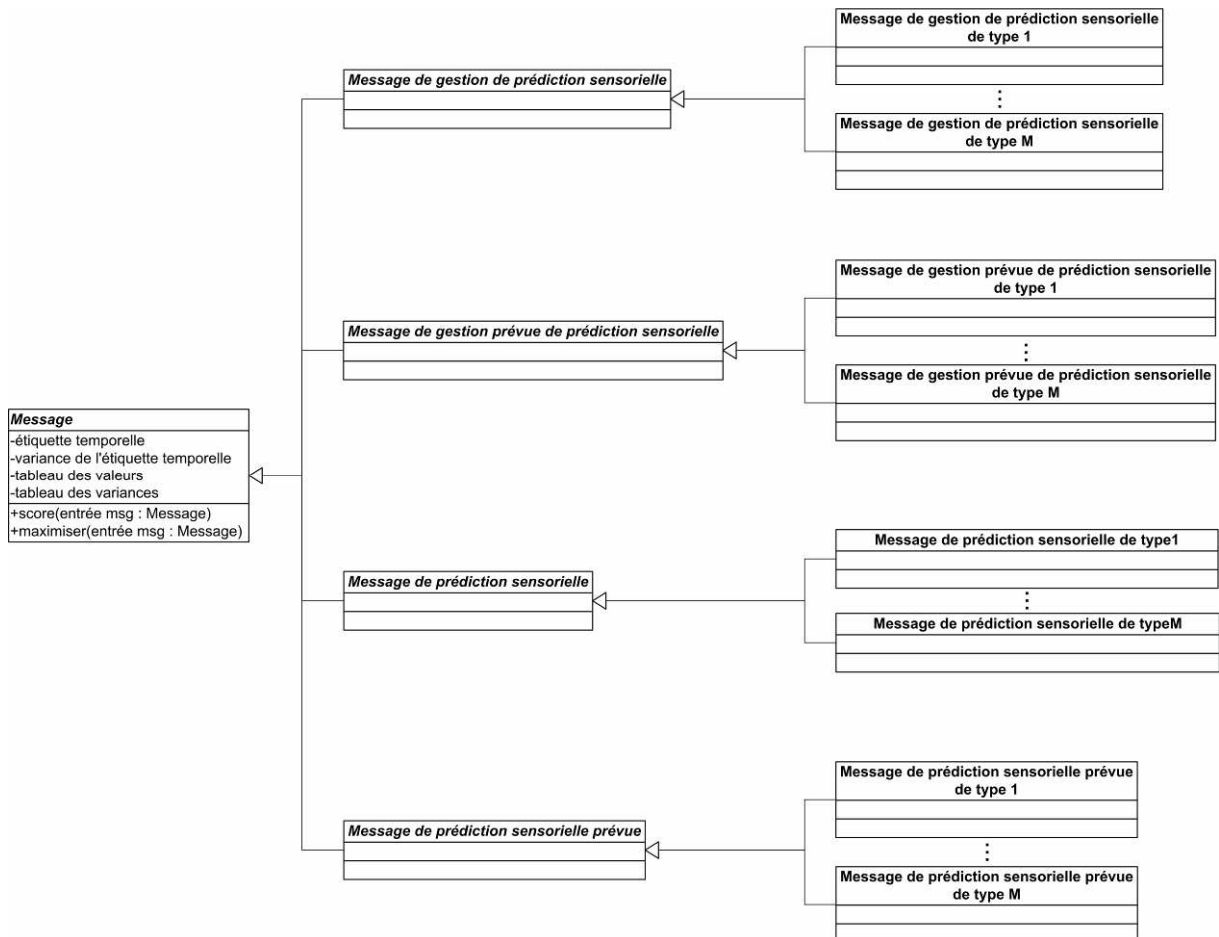


Figure III-25 : Schéma UML simplifié des objets messages sensoriels liés à la gestion de la prédiction

Le premier type de message intentionnel (*message de prédiction sensorielle prévue* (PSP) par exemple) représente la prédiction, c'est-à-dire une réplique du contenu du message prédit mais dans un type différent. Cette reproduction du contenu pourra alors être utilisée par des schèmes d'anticipation. Située en conclusion d'une règle, l'étiquette temporelle de ce type de message représente la portée temporelle de la prévision. Un même contenu d'un message pouvant être prédit avec une portée différente, il est nécessaire d'enlever cette ambiguïté par l'utilisation d'un second type de message intentionnel (*message de gestion prévue de prédiction sensorielle* (GPPS) par exemple) dont le contenu identifie les deux règles prédictives produisant la prédiction. L'étiquette temporelle de ce message correspond à la durée de la prise en compte de cette prédiction. Plus précisément, cette étiquette temporelle représente l'incertitude liée à la précision temporelle de la prédiction. Par convention, le retard supplémentaire de cette étiquette temporelle θ_2 se calcule en fonction de l'incertitude de la première :

$$\theta_2 = \theta_1 - \Delta\theta$$

$$\Delta\theta = 2 * \text{variance}_{\theta_1}$$

Par exemple, un schème prédictif sensoriel correspond à un ensemble de couples de règles (Figure III-26) : une règle de prédiction sensorielle et une règle de gestion de la prédiction sensorielle. La conclusion de la première, un message de prédiction sensorielle prévue, correspond au premier type de message évoqué précédemment. Son contenu représente le contenu du message sensoriel prédit. La conclusion de la seconde règle représente un message de gestion prévue de la prédiction sensorielle et celui-ci permet de déclencher des règles liées à l'utilisation de cette prédiction, ainsi que d'inhiber les autres prédictions, le temps de la prédiction ou de son annulation par une règle en dehors du schème. Ce mécanisme empêche qu'un même schème prédictif puisse effectuer plus d'une prédiction à la fois et ainsi respecter la règle de l'univocité.

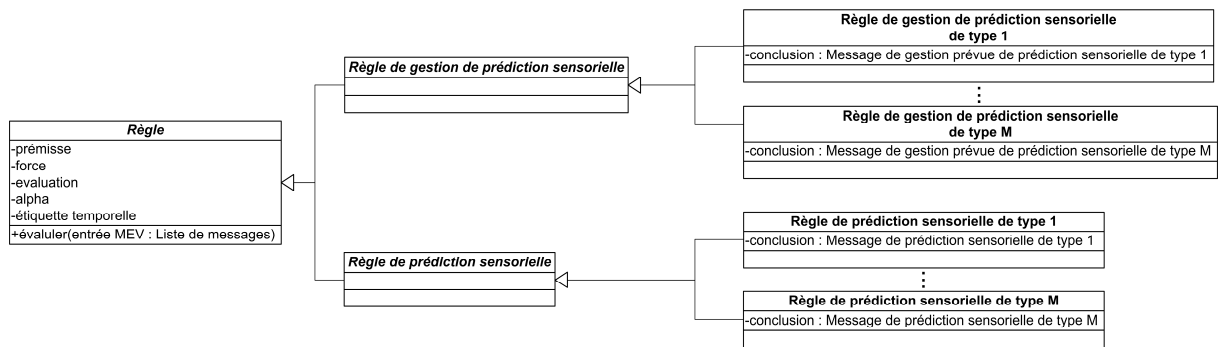


Figure III-26 : Schéma UML simplifié des objets règles sensoriels liés à la gestion de la prédiction

Ainsi, un schème prédictif portant sur le premier aspect de la prédiction, c'est-à-dire la programmation de prédispositions au déclenchement d'une règle, se traduit par un ensemble de couples de règles tel que par exemple :

$$\left[\begin{array}{l} \left\langle \left\langle \forall GPPS^m \right\rangle_I \wedge \left\langle \left\{ S_i^n(\theta_i) \right\}^{\forall n \in N} \wedge \left\{ C_j^m(\theta_j) \right\}^{\forall m \in M} \right\rangle_E \right\rangle \supset PSP_k^n(\theta_k) \\ \left\langle \left\langle \forall GPPS^m \right\rangle_I \wedge \left\langle \left\{ S_i^n(\theta_i) \right\}^{\forall n \in N} \wedge \left\{ C_j^m(\theta_j) \right\}^{\forall m \in M} \right\rangle_E \right\rangle \supset GPPS_k^m(\theta_r) \end{array} \right]_{x \in \text{schème}^P}$$

Cette gestion de la prédiction permet de distinguer θ_k la distance temporelle entre le moment de la prédiction et le moment a priori de sa concrétisation, et θ_r la durée de l'attente autorisée pour vérifier cette prédiction. Avec le jeu des incertitudes sur les étiquettes temporelles et le contenu des messages, toutes les situations peuvent être décrites comme le montera la présentation du schème de rétribution associé à une rétribution. Mais pour l'instant, dans l'exemple d'un schème de prédiction sensorielle, les étiquettes temporelles des conclusions seront identiques.

Ainsi, en reprenant l'exemple précédent sur le schème sensorimoteur secondaire complexe composé de deux couples de règles, un schème de prédiction sensorielle peut être constitué à partir d'un couple de règles comme suit :

$$\left[\begin{array}{l} \langle \forall GPPS \rangle_I \wedge \langle S_1(1) \wedge S_2(0,95) \wedge C_{a1}(1) \wedge C_{a2}(0,95) \rangle_E \supset PSP_1(0,95) \\ \langle \forall GPPS \rangle_I \wedge \langle S_1(1) \wedge S_2(0,95) \wedge C_{a1}(1) \wedge C_{a2}(0,95) \rangle_E \supset GPPS_c(0,95) \end{array} \right]_c$$

Ici, pour simplifier l'exemple, la variance de PSP1 est considérée comme nulle. Ainsi, le Tableau III-5 montre que les cheminements des messages dans la mémoire événementielle issus de la conclusion de ces deux règles sont identiques puisque leurs étiquettes temporelles initiales sont égales.

t	messages sensoriels	messages de commande prévue	messages de commande	messages de gestion prévue de commande	messages de gestion de commande	messages de prédiction sensorielle prévue	messages de prédiction sensorielle	messages de gestion prévue de prédiction sensorielle	messages de gestion de prédiction sensorielle
1	$S_1(1)$								
2	$S_2(1)$ $S_1(0,95)$	$CP_{a1}(1) \rightarrow$ $CP_{a2}(0,95)$	$C_{a1}(1)$	$CPG_a(0,95)$					
3	$S_1(1)$ $S_2(0,95)$	$CP_{a2}(1) \rightarrow$	$C_{a2}(1)$ $C_{a1}(0,95)$	$CPG_a(1) \rightarrow$	$CG_a(1)$				
4	$S_2(1)$ $S_1(0,95)$		$C_{a2}(0,95)$		$CG_a(0,95)$	$PSP_1(0,95)$		$GPPS_c(0,95)$	
5	$S_1(1)$ $S_2(0,95)$	$CP_{b1}(1) \rightarrow$	$C_{b1}(1)$	$CPG_b(0,95)$		$PSP_1(1) \rightarrow$	$PS_1(1)$	$GPPS_c(1) \rightarrow$	$GPS_c(1)$
6	$S_2(1)$ $S_1(0,95)$		$C_{b1}(0,95)$	$CPG_b(1) \rightarrow$	$CG_b(1)$		$PS_1(0,95)$		$GPS_c(0,95)$

Tableau III-5 : Exemple de six états successifs d'une mémoire événementielle avec un schème sensorimoteur secondaire complexe associé à un schème de prédiction sensorielle.

ii - L'utilisation ou la vérification de la prédiction

Les schèmes utilisant la prédiction représentent soit des schèmes liés à la simulation ou à l'évaluation de conséquences, soit des schèmes liés à la vérification. Les règles de ces schèmes sont des règles à déclenchement attendu. Le premier type de schème qui utilise la prédiction pour la simulation constitue une perspective de recherche pour l'anticipation, mais dans le cadre de la spécification de l'architecture cognitive proposée, seul le second type de schème sera développé. Le deuxième type de schème qui vérifie la prédiction correspond à un mécanisme de rétribution interne en fonction de la vérification de la

prédiction. Ce mécanisme repose sur l'emploi de règles à déclenchement attendu et d'un message de rétribution spécifique. En reprenant le schème de la prédiction sensorielle, la Figure III-27 représente les types de règles et de messages nécessaires à la constitution d'un schème de vérification sensorielle.

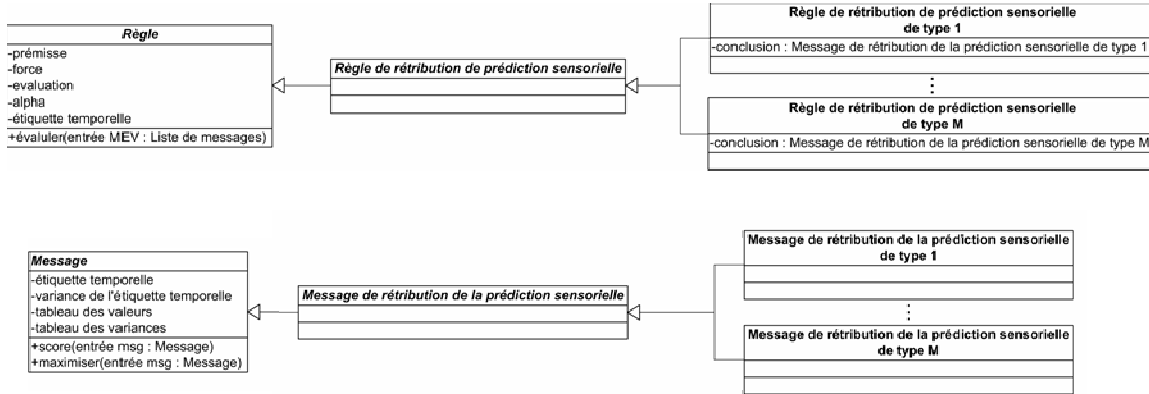


Figure III-27 : Schéma UML simplifié des objets « règle » et « message » de rétribution relatifs à la prédiction sensorielle

Ainsi, le schème de vérification de la prédiction sensorielle regroupe l'ensemble des règles avec la forme suivante, avec R pour le message de rétribution :

$$\langle GPPS_k^m (\forall \theta_r) \wedge S_i^n (\theta_0) \wedge PSP_i^n (\theta_0) \rangle_E \supset R(1)$$

Avec une incertitude minimale sur le contenu du message GPPS mais avec une incertitude maximale pour son étiquette temporelle, le message GPPS permet de sélectionner la règle de rétribution de prédiction sensorielle correspondant précisément à la prédiction que le message GPPS gère sans affecter l'évaluation de la prédiction. Le message PSP sert à vérifier que le signe prédit arrive bien au moment prévu et le message S sert à vérifier son contenu.

Le contenu de la conclusion représentant la rétribution, celui-ci n'est pas prédéfini comme pour les règles spontanées traditionnelles, c'est une spécificité des règles d'attente. Le contenu de la conclusion est calculé à partir de l'évaluation de règle RPS qui sera détaillée ultérieurement, la portée de la prédiction θ_i et d'un coefficient de rétribution interne C_{ri} :

$$rétribution = évaluation(RPS) * (1 - \theta_i) * C_{ri}$$

Le premier terme de ce calcul permet de récompenser les règles les plus précises et le deuxième de favoriser les règles ayant des prédictions les plus éloignées dans le temps. De plus, les règles prédisant tout avec une incertitude maximale obtiendront une rétribution nulle. La rétribution est alors délivrée à toutes les règles s'étant déclenchées entre le moment de la prédiction et la fin de son évaluation. Cette redistribution aveugle repose sur l'idée qu'il n'est pas possible de déterminer les règles déclenchées ayant réellement un rôle dans la rétribution ; c'est le principe de la détection des effets de bords évoqué dans le deuxième chapitre, lors de la présentation du pragmatisme.

La rétribution se représente néanmoins par un message comme un autre et pourra faire l'objet de schème visant sa prédiction. Ainsi, l'activité rétribuée et l'activité de rétribution se trouvent au même niveau de représentation, permettant ainsi d'imaginer des mécanismes de récurrences qui font défaut aux autres systèmes de classeurs.

5.2. Le processus de sélection d'une règle

La partie précédente a montré que l'interprétation d'un signe doit résulter d'une mise en confrontation directe avec d'autres. Pour respecter cette idée, la sélection d'une règle s'effectue en deux étapes : la détermination du score et l'évaluation. La détermination du score est identique à toutes les règles, en revanche, l'évaluation diffère selon que la règle soit à déclenchement spontané ou à déclenchement attendu. En effet, l'évaluation d'une règle à déclenchement spontané se fonde sur un instant donné de la mémoire événementielle, contrairement à l'évaluation d'une règle à déclenchement attendu qui porte sur l'évolution de la mémoire événementielle sur une certaine durée. Dans les deux cas, à la suite de la sélection d'une règle, les messages des prémisses relatives à un type de primitive sensorielle sont ajustés en fonction de son évaluation. Ce procédé permet au contenu des règles en rapport direct avec l'environnement de se façonner en fonction des expériences et permet ainsi l'auto-organisation de la segmentation du monde. Cet ajustement se basant sur l'évaluation, il ne s'applique pas de la même manière entre les règles à déclenchement spontané et celles à déclenchement attendu.

La présentation de ces divers mécanismes prendra l'ordre suivant : (A) le calcul du score, (B) l'évaluation et l'ajustement pour une règle à déclenchement spontané et enfin (C) l'évaluation et l'ajustement pour une règle à déclenchement attendu.

A - Le calcul du score

Le score d'une règle correspond à la distance minimale entre les messages de la prémisses et ceux de la mémoire événementielle. Si la mémoire événementielle ne contient pas assez de messages d'un certain type pour l'appariement, le score de la règle est nul. Afin que le score soit borné entre 0 et 1, le calcul de ce dernier pour la prémisses excitatrice et inhibitrice s'exprime sous la forme d'un produit d'exponentielles, avec K messages P_k de la prémisses ayant N_k valeurs p_{kj} associées chacune à une incertitude σ_{kj} , et K messages X_k de la mémoire événementielle :

$$score_{excitateur}(P_{excitatrice}, X) = \prod_k e^{-\sum_0^{N_k} (p_{kj} - x_{kj})^2 / 2\sigma_{kj}^2}$$

$$score_{inhibiteur}(P_{inhibitrice}, X) = \sum_k e^{-\sum_0^{N_k} (p_{kj} - x_{kj})^2 / 2\sigma_{kj}^2}$$

$$score_{final}(P, X) = score_{excitateur}(P_{excitatrice}, X) - score_{inhibiteur}(P_{inhibitrice}, X)$$

Plus précisément, le score de la prémisses excitatrice correspond à la combinaison par arrangement des distances messages de même type produisant une somme de distance minimale. S'il n'existe pas suffisamment de signes dans la mémoire événementielle pour l'apparier au moins une fois, alors le score de cette prémisses excitatrice est automatiquement nul. L'absence d'un signe est pénalisante pour la prémisses excitatrice. Par opposition, le score de la prémisses excitatrice correspond à la combinaison par arrangement des distances messages de même type produisant une somme de distance maximale. En revanche, l'absence de l'appariement d'un message n'affecte pas le calcul du score ; si aucun message n'est apparié alors le score de la prémisses inhibitrice sera équivalent à 1. Cette dissymétrie illustre le fait que les messages de la prémisses excitatrice sont conjonctifs alors que ceux de la prémisses inhibitrice sont disjonctifs.

B - L'évaluation et l'ajustement pour une règle à déclenchement spontané

En considérant que chaque règle représente un modèle local du monde parmi les règles du même type, alors leur distribution par rapport à l'environnement peut être prise en compte par un facteur *alpha* α borné entre 0 et 1, attribué à chaque règle et initialisé à 1. L'évaluation ponctuelle d'une règle représente sa légitimité à être déclenchée par rapport aux autres règles de même type. L'évaluation ponctuelle se calcule alors de la manière suivante :

$$\text{évaluation}(P_n, X_m)(i) = \frac{\text{score}(P_n, X_m)(i) * \alpha_n}{\sum_{j=1}^N \alpha_j * \text{score}(P_j, X_m)(i)}$$

Cette évaluation correspond, en fait, à la phase d'estimation de l'algorithme *Expectation-maximisation* (EM) de Dempster et al. (1977) qui est un algorithme itératif utilisé pour rechercher le maximum de vraisemblance en fonction de certains paramètres. La phase de maximisation servira ici à actualiser *alpha* et à ajuster les valeurs des messages relatifs aux primitives sensorielles ainsi que leurs incertitudes associées à ces valeurs. Comme les règles sont des modèles à noyau gaussien, en faisant l'hypothèse que les composantes d'un message sensoriel sont décollés, soit statiquement indépendantes, les calculs de la phase de maximisation qui ajustent les valeurs μ_{nk} et les incertitudes σ_{nk} des messages de la prémisse P_n avec l'état X_m de la mémoire événementielle deviennent :

$$S_n(i+1) = \sum_{m=1}^M \text{évaluation}(P_n, X_m)(i)$$

$$\alpha_n(i+1) = \frac{1}{M} * S_n(i+1)$$

$$\mu_{nk}(i+1) = \frac{1}{S_n(i+1)} * \sum_{m=1}^M \text{évaluation}(P_n, X_m)(i) * x_{mk}$$

$$\sigma_{nk}^2(i+1) = \frac{1}{S_n(i+1)} * \sum_{m=1}^M \text{évaluation}(P_n, X_m)(i) * (x_{mk} - \mu_{nk}(i+1))^2$$

Supposer que les primitives sensorielles produisent des messages sensoriels decorrelés constitue une hypothèse forte, mais les travaux d'Atick (1992) évoqués dans le premier chapitre concernant le connexionnisme suggèrent que la maximisation de la transmission de l'information peut se situer à très bas niveau.

C - L'évaluation et l'ajustement pour une règle à déclenchement attendu

Lors de la gestion de la prédiction, la règle de gestion de prédiction envoie un message dont le contenu identifie la prédiction. Ainsi, la sélection de la règle d'attente associée peut être tout ou rien en mettant l'incertitude sur ce type de message au maximum. La compétition ne s'effectue pas entre les règles d'attente puisqu'elles sont tout ou rien (leur *alpha* n'évolue pas et demeure à 1). L'évaluation des règles d'attente s'effectue sur l'ensemble des scores obtenus par la règle durant l'attente fixée par la présence du message GPPS dans la mémoire événementielle :

$$\text{évaluation}(P_n, X_m)(i) = \frac{\text{score}(P_n, X_m)(i)}{\sum_{j=1}^M \text{score}(P_n, X_j)(i)}$$

C'est uniquement à la fin de cette évaluation que les règles d'attente se déclenchent, entraînant l'envoi de la conclusion et l'ajustement des prémisses des règles d'attente. Cette dernière, s'inspirant de l'algorithme EM, s'effectue de la manière suivante :

$$\mu_{nk}(i+1) = \frac{1}{S_n(i+1)} * \sum_{m=1}^M \text{évaluation}(P_n, X_m)(i) * x_{mk}$$

$$\sigma_{nk}^2(i+1) = \frac{1}{S_n(i+1)} * \sum_{m=1}^M \text{évaluation}(P_n, X_m)(i) * (x_{mk} - \mu_{nk}(i+1))^2$$

Les mécanismes subcognitifs liés au déclenchement des règles d'attente sont plus complexes que ceux des règles spontanées qui se traduisent simplement par l'envoi de messages de conclusion dans la mémoire événementielle. En effet, en plus de l'évaluation qui s'effectue sur une certaine durée, l'ajustement doit se reporter sur plusieurs règles puisque la vérification doit permettre d'ajuster les règles à l'origine de la prédiction. Le contenu du message GPPS permet d'identifier ces règles et ainsi de reporter l'ajustement du contenu sur SPS et de l'ajustement temporel sur SPS et GPPS.

5.3. Les bases d'une autopoïèse sémiotique

L'autopoïèse sémiotique repose sur une dynamique interne qui, dans un premier temps, n'est basée que sur la seule activité des règles système. Dans un second temps, elle est associée à une rétribution endogène afin d'orienter cette dynamique. Avant de dégager l'origine du maintien et de l'élimination d'une règle, (A) une étude sur la dynamique basée sur la rétribution positive et la rétribution négative a été entreprise. Cette étude permettra (B) d'adapter la *Bucket brigade* de Holland pour l'évolution de la force des règles puisque la rétribution native ne repose pas sur l'enchaînement des règles mais sur leur auto-organisation.

A - Étude d'une dynamique fondée sur la rétribution

L'évolution dynamique de la *force* des règles repose sur des rétributions positives et négatives successives. Cette dynamique doit discriminer les règles pertinentes des règles redondantes et des règles inactives, tout en prenant en compte qu'une règle ayant été suffisamment pertinente puisse être conservée bien qu'elle ne se soit pas déclenchée depuis un certain temps ou qu'elle se montre inappropriée pour une situation donnée. La solution choisie consiste à rétribuer positivement ou négativement en fonction de la force des règles, ce qui a pour conséquence de rendre les règles faibles sensibles à la rétribution positive et les règles fortes insensibles à la rétribution négative. L'étude de cette dynamique se limitera dans un premier temps (i) à la dynamique induite par une rétribution positive, puis dans un deuxième temps (ii) à la dynamique induite par la rétribution négative, et dans un troisième temps (iii) à la combinaison de ces deux dynamiques.

i - La dynamique induite par une rétribution positive

En ne considérant que la rétribution, l'évolution de la force F entre deux rétributions successives avec une densité ρ^+ à un intervalle Δt se formule de la manière suivante :

$$F(t + \Delta t) = F(t) + \rho^+ (1 - F(t))$$

Lorsque Δt tend vers 0, la rétribution positive devient continue selon un taux λ^+ , cette formule se réduit alors à une équation différentielle ayant la forme suivante :

$$\frac{dy(t)}{dt} = \lambda^+ (1 - y(t))$$

Mais en raisonnant avec des intervalles discrets, la formule de la $F(t+\Delta t)$ devient une suite récurrente qui converge vers 1 avec $\forall F(0) \in]0..1]$ et $\rho^+ > 0$:

$$F(n+1) = F(n) + \rho^+ (1 - F(n))$$

Dans ce cas, l'évolution de la force correspond à une croissance exponentielle limitée. De cette manière, la règle est d'autant plus sensible à une rétribution que sa force est faible (Figure III-28).

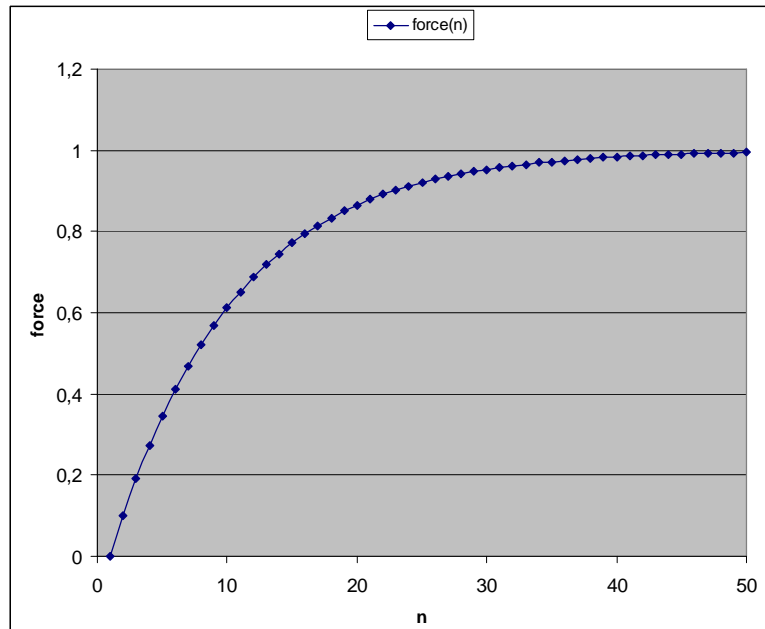


Figure III-28 : Exemple de l'évolution de la force avec une rétribution positive à intervalle régulier avec $F(1)=0,001$ et $\rho^+=0,1$.

ii - La dynamique induite par la rétribution négative

La rétribution négative pourrait correspondre à une décroissance exponentielle afin de conserver une symétrie. Mais dans ce cas, plus les règles possèdent une force élevée, plus elles deviennent sensibles aux rétributions négatives, ce qui n'est pas souhaitable. Une manière de conserver une sensibilité moindre, à la fois dans le cas d'une force élevée et dans le cas d'une force faible, consiste à adopter une décroissance sigmoïde dont l'expression, dans le cadre d'une rétribution négative continue avec un taux λ^- , est :

$$y(t) = \left(1 + e^{\lambda^- (t-\theta)}\right)^{-1}$$

Cette fonction est solution de l'équation différentielle :

$$\frac{dy(t)}{dt} = \lambda^- y(t)(1 - y(t))$$

Cette écriture suggère l'équation aux différences suivantes pour le calcul de la force en cas de rétribution irrégulière :

$$F(t + \Delta t) = F(t) + \rho^- F(t)(1 - F(t))$$

De cette expression, la suite récurrente déduite converge vers 0 avec $\forall F(0) \in [0..1[$ et $\rho^- < 0$, comme l'illustre la Figure III-29 :

$$F(n + 1) = F(n) + \rho^- F(n)(1 - F(n))$$

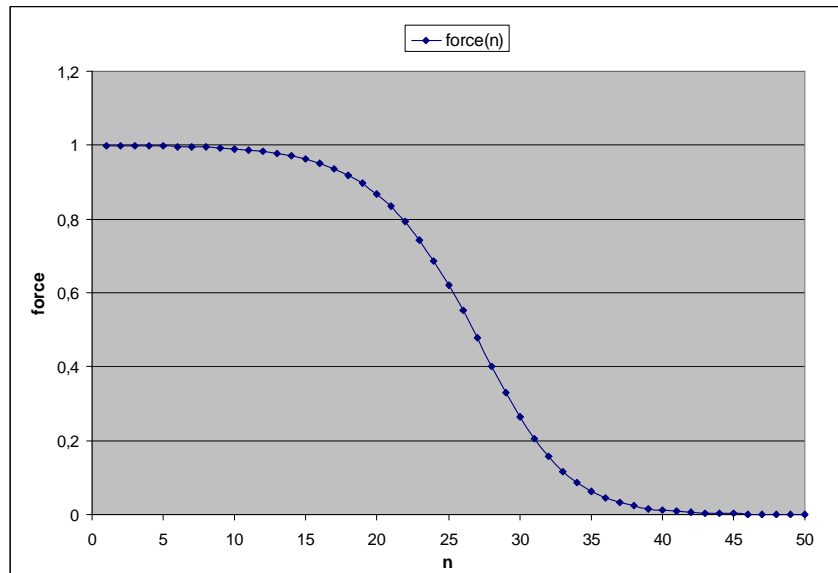


Figure III-29 Exemple de l'évolution de la force avec une rétribution négative à intervalle régulier avec $F(1)=0,999$ et $\rho=0,3$.

iii - La combinaison des dynamiques induites par les rétributions positives et négatives

La combinaison des deux dynamiques doit offrir un système dynamique avec de nombreux points d'équilibre afin qu'un ensemble de règles puisse se maintenir. Dans le cadre des rétributions continues avec des taux constants, l'équation différentielle totale devient la suivante :

$$\frac{dy(t)}{dt} = \lambda^- y(t)(1 - y(t)) + \lambda^+ (1 - y(t))$$

Soit :

$$\frac{dy(t)}{dt} = [\lambda^- y(t) + \lambda^+](1 - y(t))$$

À l'infini, la solution de cette équation présente le point d'équilibre suivant :

$$y(\infty) = -\frac{\lambda^+}{\lambda^-}$$

Cet équilibre survient lorsque la rétribution positive compense la rétribution négative. Plus précisément, cet équilibre correspond au point de croisement entre les courbes d'évolution des quantités de rétributions. Par exemple, la Figure III-30 illustre le cas où le point d'équilibre de la force se situe à 0,5, obtenu avec les taux $\lambda^+ = 1$ et $\lambda^- = -2$.

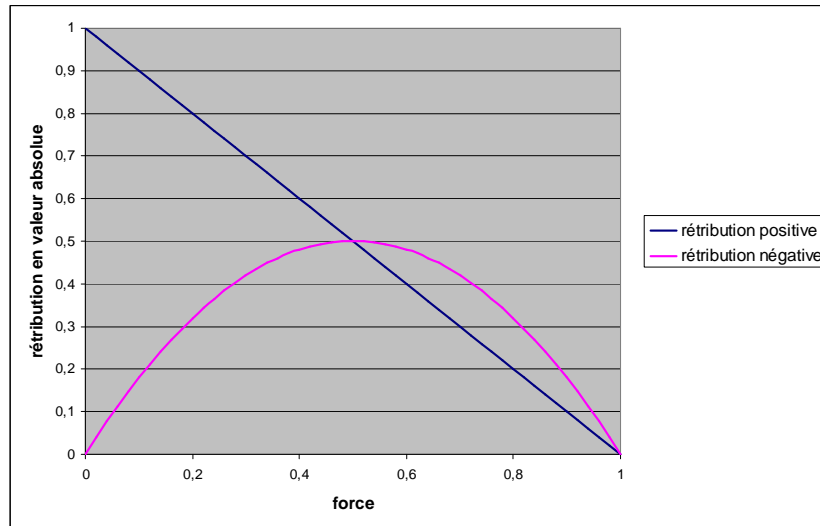


Figure III-30 : Exemple d'évolution des quantités de rétribution avec les taux $\lambda^+=1$ et $\lambda^-=-2$.

La Figure III-31 présente les différentes trajectoires produites par l'attracteur défini par les taux $\lambda^+ = 1$ et $\lambda^- = -2$.

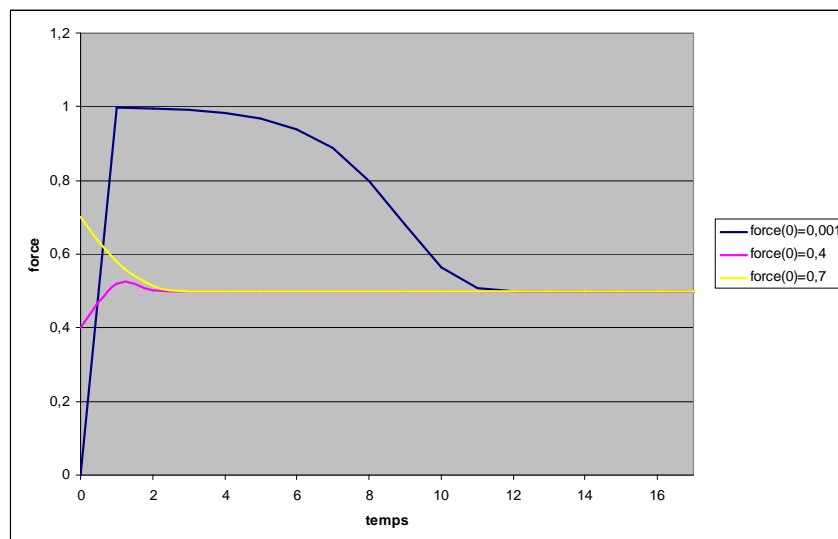


Figure III-31 : Exemple d'évolution de la force avec les taux $\lambda^+=1$ et $\lambda^-=-2$ selon que la force initiale soit à 0,001, 0,4 ou 0,7.

Cet exemple ne reflète pas réellement les conditions de la dynamique de la rétribution des règles dans la mesure où les rétributions ne sont ni continues ni régulières ni

synchrones. Par ailleurs, ρ^+ et ρ^- peuvent ne pas être constants puisque, par exemple, la quantité de rétribution d'une prédiction est liée à son évaluation. Une manière de compenser les faiblesses de cette étude sur la dynamique de l'évolution de la force des règles consiste à effectuer une étude empirique. Toutefois, avant l'implémentation du système, une simulation simplificatrice peut explorer le problème des rétributions non continues, à des fréquences différentes.

Par exemple, en considérant que la fréquence de rétribution négative est supérieure à celle de la rétribution positive et que les taux de rétribution sont constants, le système ne possède que deux points d'équilibre stables possibles : lorsque la force tend avec n à l'infini vers 1 et quand la force tend avec n à l'infini vers 0. Les autres points d'équilibre sont instables, c'est-à-dire que la force oscille entre deux bornes. Les points stables résultent de la compensation en une rétribution positive de toutes les rétributions négatives accumulées, avec une rétribution négative à chaque pas et une rétribution positive tous les x pas :

$$\rho^+(1 - F(n)) + \rho^-(1 - F(n))F(n) = F(n - x) - F(n)$$

$$\rho^+(1 - F(n)) = -\rho^- \sum_{k=n-x}^n (1 - F(k))F(k)$$

$$-\frac{\rho^+}{\rho^-} = \frac{\sum_{k=n-x}^n (1 - F(k))F(k)}{(1 - F(n))}$$

Dans ce cas, l'équilibre de la dynamique dépend des taux de rétribution mais également du nombre de périodes de rétributions négatives au cours d'une période de rétribution positive. Les graphiques des Figure III-32 et Figure III-33 montrent deux équilibres stables issus de deux dynamiques avec des paramètres semblables, hormis pour la fréquence de rétribution positive.

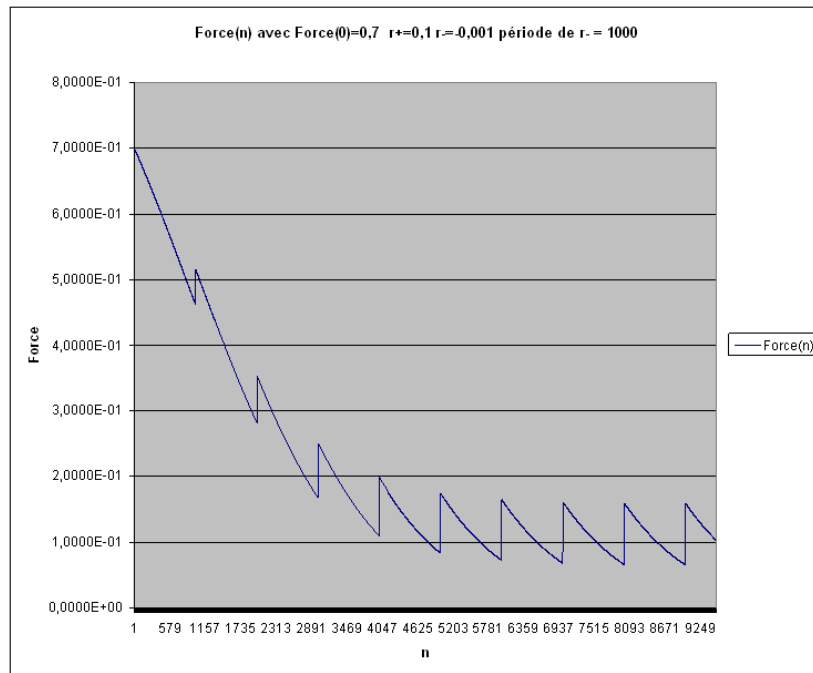


Figure III-32 : Exemple d'un équilibre stable produit par l'apport périodique de rétribution positive et l'apport continu de rétribution négative.

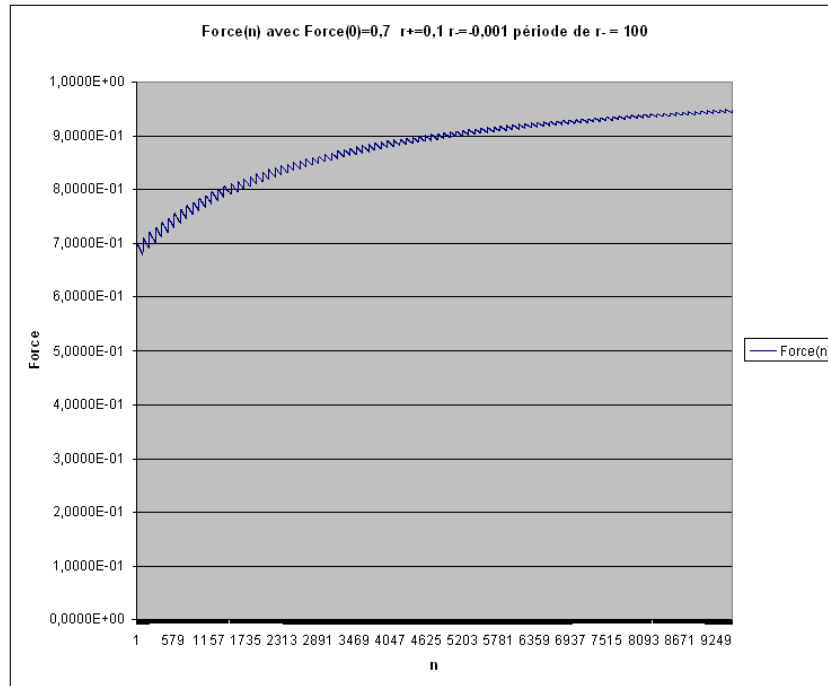


Figure III-33 : Exemple d'un autre équilibre stable mais avec une fréquence de rétribution positive supérieur au premier.

Bien que la dynamique induite soit complexe, elle offre l'avantage de trouver des points d'équilibre stables dès que les rétributions deviennent périodiques. Les règles recevant de manière répétitive une rétribution positive forte atteignent une force suffisante pour supporter les rétributions négatives, et les règles recevant de manière sporadique des rétributions peuvent se maintenir dans une zone de force moyenne ou faible si ces dernières sont toutefois suffisamment régulières, autrement dit dans une zone à plus forte sensibilité aux rétributions, mettant ainsi la règle toujours à l'épreuve.

B - Utilisation de la métaphore financière

Dans les systèmes de classeurs traditionnels, la rétribution est toujours liée aux conséquences produites par les règles (enchaînement, prédiction). La métaphore financière se fonde alors sur l'idée du retour sur investissement. Mais ces critères introduisent toujours un biais téléonomique de bas niveau qui empêche au final le développement de l'autonomie à haut niveau. Pour que l'autopoïèse sémiotique puisse avoir a priori une dynamique propre sans finalité, l'auto-organisation des règles doit résulter de leur propre activité et faciliter celle-ci. Le retour sur investissement devient alors la sélection de la règle elle-même.

La métaphore financière s'interprète de la manière suivante : toutes les règles paient une taxe T de façon à ce qu'une règle qui ne génère pas de rétribution positive soit pénalisée, et toutes les règles paient une enchère E de façon à ce que les règles trop redondantes soient éliminées, avec le taux d'enchère C_e et le taux de la taxe C_t :

$$E(n) = C_e * \text{évaluation}(n) * F(n)(1 - F(n))$$

$$T(n) = C_t * F(n)(1 - F(n))$$

L'enchère est proportionnelle à l'évaluation de la règle de sorte que les règles ayant une évaluation nulle ne paient que la taxe. La taxe et l'enchère représentent des rétributions négatives conduisant à la décroissance de la force. À chaque pas de temps, la force décroît selon l'équation suivante, avec le taux de taxe C_t et le taux enchère C_e réels négatifs :

$$F(n+1) = F(n) + T(n) + E(n)$$

$$F(n+1) = F(n) + (C_t + C_e * \text{évaluation}(n)) * F(n)(1 - F(n))$$

Ces deux mécanismes de rétribution négative constituent des mécanismes subcognitifs, c'est-à-dire que leur activité n'est pas définie par des règles et n'apparaît pas dans la mémoire événementielle. La taxe étant une rétribution négative fréquente, si une règle considérée pertinente doit tendre vers 1, il faut que la taxe soit faible par rapport à la rétribution positive. Le retour sur investissement pour la règle sélectionnée se traduit par un remboursement de son enchère avec intérêt sous forme de rétribution positive, soit avec le taux de remboursement positif C_r :

$$\text{remboursement}(n) = C_r * \text{évaluation}(n) * (1 - F(n))$$

L'évolution dynamique de la force s'écrit alors de la manière suivante, avec un C_r nul si la règle n'est pas sélectionnée :

$$F(n+1) = F(n) + T(n) + E(n) + \text{remboursement}(n)$$

$$F(n+1) = F(n) + ((C_t + C_e * \text{évaluation}(n)) * F(n) + C_r * \text{évaluation}(n)) * (1 - F(n))$$

$$F(n+1) = F(n) + (C_t * F(n) + (C_e * F(n) + C_r) * \text{évaluation}(n)) * (1 - F(n))$$

Ou encore,

$$F(n+1) = F(n) + (C_t * F(n) + \text{évaluation}(n) * |C_e| * \left(\frac{C_r}{|C_e|} - F(n)\right))(1 - F(n))$$

La dynamique induite par cette équation permet à un ensemble de règles dont la pertinence n'a pas été confirmée d'être maintenu vers un équilibre instable en fonction de la fréquence. Pour que les règles sortent de cette zone d'équilibre sensible, une rétribution interne, délivrée par les règles elles-mêmes, doit infléchir positivement ou négativement la trajectoire de leur dynamique. La dynamique propre du système permet de conserver un ensemble de règles a posteriori viable mais a priori sans finalité. C'est la rétribution interne avec un taux C_{ri} qui va privilégier telle ou telle règle comportementale en fonction d'intérêts exprimés également sous la forme d'une règle. Dans le cas d'une rétribution interne positive liée à la vérification de la prédiction, le calcul de la force, avec θ la portée temporelle de la prédiction et avec un C_r nul si la règle n'est pas sélectionnée et un C_{ri} nul si la règle n'a pas été récompensée, devient :

$$F(n+1) = F(n) + T(n) + E(n) + \text{remboursement}(n) + \text{récompense}(n)$$

$$F(n+1) = F(n) + (C_t * F(n) + (C_e * F(n) + C_r + C_{ri} * (1 - \theta)) * \text{évaluation}(n)) * (1 - F(n))$$

5.4. L'implémentation

Toutes les spécifications avancées pour une architecture cognitive pragmatiste ont été développées. Un module de création de règle a été rajouté afin d'étudier le comportement du système face à l'introduction de nouvelles règles. Ce processus est complètement subcognitif dans la mesure où il possède son propre cycle de fonctionnement. Pour

restreindre le champ des possibles, la construction de nouvelles règles s'effectue dans la perspective de compléter un schème cognitif défini comme une boucle sensorimotrice secondaire simple. Concrètement, à intervalle fixe, un message sensoriel et un message de commande se trouvant dans la mémoire événementielle sont copiés afin de former une nouvelle règle qui sera introduite dans la mémoire épistémique. Le développement a été effectué en Ada et représente plus de 14000 lignes de code effectif.

Plus précisément, l'implémentation de l'architecture cognitive proposée a été voulue indépendante du système robotique. La conception abstraite des messages et des règles permet de concevoir n'importe quel schème cognitif nécessitant des prémisses multi-message et multi-type, inhibitrice ou excitatrice. Le nombre de types de messages et de règles gérés par le système se trouve déterminé à l'initialisation par la lecture d'un fichier XML décrivant les paramètres et les règles initiales. La dynamique des règles décrite précédemment a également été implémentée ainsi que les mécanismes de gestion des règles d'attente avec leur rétribution. La prise en compte des événements temporels indépendamment de la fréquence de sélection des règles nécessite une programmation multitâche, autrement dit, la mémoire événementielle et la mémoire épistémique sont représentées par des tâches. L'interfaçage avec les primitives sensorielles et les primitives motrices s'effectue par l'intermédiaire d'une tâche réceptrice et d'une tâche émettrice. Une bibliothèque générique pour les règles et les messages intégrant leur contrainte et leurs relations permet de faciliter cet interfaçage. Cette gestion considère bien les primitives comme des processus concourants et de ce fait représente une solution pour les architectures robotiques s'appuyant sur de tels processus.

Comme l'illustrera le chapitre suivant, le développement respecte les contraintes temps réel qu'impose la robotique. Par ailleurs, l'architecture cognitive permet une gestion asynchrone : la tâche réceptrice recueille les informations sensorielles des primitives sensorielles à une fréquence f_a puis les envoie à la tâche mémoire événementielle. La mémoire épistémique sélectionne une règle avec une fréquence f_b puis renvoie la conclusion à la mémoire événementielle. Les temps de traitement peuvent évoluer en fonction du nombre de règles et de signes dans la mémoire événementielle. La fréquence f_b n'est donc pas assurée. Afin d'éviter ces imprécisions, à chaque sélection, les étiquettes temporelles sont recalculées par rapport à l'horloge interne. Par ailleurs, le mécanisme d'optimisation des prémisses des règles ajuste naturellement les incertitudes temporelles.

Afin de mesurer objectivement certains aspects du comportement du robot, l'ensemble des caractéristiques des systèmes sont enregistrés à chaque sélection de règle sous un format XML. Le traitement de ces données permet d'explorer soit l'évolution des caractéristiques des règles soit l'évolution des caractéristiques concernant l'ensemble des règles. Le développement et l'utilisation d'un logiciel de visualisation permettent de multiplier les points de vue et d'isoler certains phénomènes de sorte que l'exploration des données peut couvrir tous les aspects de la dynamique de l'architecture cognitive.

6. Discussion sur l'expressivité de l'architecture cognitive générale proposée

En définissant la notion de schème cognitif par un ensemble de règles qui peuvent s'enchaîner dans leur déclenchement, la détermination des schèmes initiaux, des protocroyances, peut s'apparenter à une sorte de programmation. Dans le seul objectif d'avoir un aperçu de ce à quoi pourrait aboutir l'approfondissement d'une programmation cognitive,

deux aspects seront très succinctement abordés (A) l'introduction de structures algorithmiques élémentaires et (B) la modélisation de conditionnement ou d'optimisation comportementale, ces deux aspects pouvant se combiner et interagir au cours du couplage.

A - Les capacités algorithmiques

Les deux principales structures algorithmiques qui peuvent être introduites sont la comparaison et la boucle.

La première structure algorithmique, la comparaison, peut se traduire par la compétition entre deux règles. Par exemple, deux règles A et B, avec le même alpha, possèdent en guise de prémisse un message interne de même type dont l'incertitude sur le contenu est totale mais dont les valeurs concernant la légitimité et l'incertitude sont spécifiques. La compétition entre ces deux règles revient alors à comparer la légitimité du message par rapport à un seuil. La Figure III-34 illustre la suprématie d'une règle par rapport à une autre suivant la valeur d'entrée considérée sur laquelle se porte la comparaison.

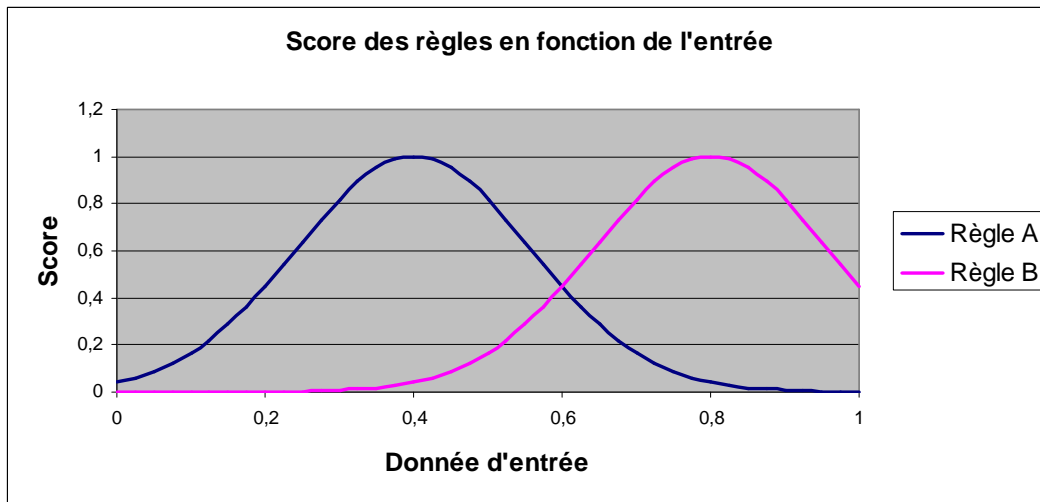


Figure III-34 : Exemple d'une opération de comparaison avec deux règles ayant le même alpha.

La seconde structure algorithmique, la boucle « tant que » s'exprime par trois règles. La première règle sert de déclencheur en introduisant un message interne excitateur (MIE) spécifique dans la mémoire événementielle qui déclenchera une règle produisant ce même message interne excitateur. Ainsi, cette règle s'auto-déclenchera continuellement *tant qu'*une troisième règle n'aura pas produit un message interne inhibiteur (MII) pour cette règle. la boucle en introduisant ce schème cognitif peut s'exprimer alors de la manière suivante :

$$\left\langle \left\{ S_i^n(\theta_i) \right\}_{\forall n \in N} \right\rangle_E \supset MIE_k^m(\theta_k)$$

$$\left\langle MII_k^m(\theta_k) \right\rangle_I \wedge \left\langle MIE_k^m(\theta_k) \right\rangle_E \supset MIE_k^m(\theta_k)$$

$$\left\langle \left\{ S_j^n(\theta_j) \right\}_{\forall n \in N} \right\rangle_E \supset MII_k^m(\theta_k)$$

Ce schème introduit non seulement une boucle mais également une structure de contrôle concurrente. Par ailleurs, ces structures peuvent être figées à l'initialisation comme n'importe quelle règle en fixant la force à 1 et en fixant le nombre d'ajustements au maximum.

B - Les types de dressages modélisables

Le mécanisme de rétribution étant défini via des règles, il est tout à fait possible d'imaginer que dans une situation incertaine, une production de règles comportementales puisse être entreprise jusqu'à ce que l'environnement déclenche une règle de rétribution favorisant le maintien des règles précédentes ou favorisant leur élimination ou leur non activation dans cette situation. Cet apprentissage par renforcement correspond en fait à un conditionnement opérant qui réalise l'association entre une action de l'animal et un stimulus inconditionné (SI).

Ce type d'apprentissage déjà évoqué pose de nombreux problèmes concernant la détermination des actions à rétribuer et incite à créer un type de règles qui étiquetterait les règles susceptibles d'être rétribuées dans un laps de temps déterminé. Cela reste une perspective par rapport au travail effectué. Toutefois, la visualisation de la rétribution sous forme de signe dans la mémoire événementielle permet d'entrevoir la possibilité de modélisations plus complexes des conditionnements de type opérant.

Le conditionnement classique (ou conditionnement pavlovien) est le second type de conditionnement qui serait modélisable à l'aide de l'architecture cognitive proposée. Celui-ci consiste à associer un stimulus neutre à un stimulus inconditionnel (SI). Le stimulus neutre n'entraîne pas de réponse inconditionnée (RI). Au plus, il produit une mise en alerte, autrement dit dans tous les cas une réponse neutre (RN).

SN -> RN

En revanche, la présentation d'un stimulus inconditionné conduit à un comportement de nature "réflexe", un réflexe inconditionné (RI) :

SI -> RI

L'association entre SI et SN s'effectue lors de la présentation régulière de la succession du SN puis du SI. En effet, progressivement, le stimulus qui était neutre va déclencher le réflexe de sorte que celui devient conditionné.

SC -> RC

Le conditionnement classique opère donc uniquement sur des réflexes et par la coïncidence ou la contingence de stimuli dans le temps et l'espace. Une modélisation de ce type de conditionnement serait alors la création de règles en fonction de la coïncidence des événements indiqués par les étiquettes temporelles des signes se trouvant dans la mémoire événementielle. La pertinence de ces nouvelles règles est assurée par l'association d'une règle de prédiction.

Enfin, le dernier type de comportement modélisable serait celui décrit dans le chapitre précédent sur l'indécision de deux solutions radicales comme la fuite ou l'attaque se traduisant par une activité non significative par rapport au contexte, comme l'acte de picorer pour le coq. Cela peut se comprendre comme la modulation de l'alpha et de la variance des règles conditionnant leur déclenchement.

Par exemple, les règles A et B déclenchent respectivement la fuite et l'attaque en fonction d'une valeur perceptive et la règle C sensible à la même variable perceptive déclenche l'action de picorer. Afin de discriminer au mieux les comportements d'attaque et de fuite pour favoriser la radicalisation des actions, il peut être opportun d'introduire un message inhibiteur dans les règles A et B qui aurait pour valeur celle du déclenchement de la règle antagoniste :

$$\begin{aligned} \langle S_i^m(\theta_k) \rangle_I \wedge \langle S_j^m(\theta_k) \rangle_E &\supset C_a^m(\theta_k) \\ \langle S_j^m(\theta_k) \rangle_I \wedge \langle S_i^m(\theta_k) \rangle_E &\supset C_b^m(\theta_k) \\ \langle S_g^m(\theta_k) \rangle_E &\supset C_c^m(\theta_k) \end{aligned}$$

La Figure III-35 illustre les quatre différentes phases comportementales de l'exemple avec trois règles et montre également le pouvoir de discrimination entre les règles A et B qu'introduisent les messages inhibiteurs.

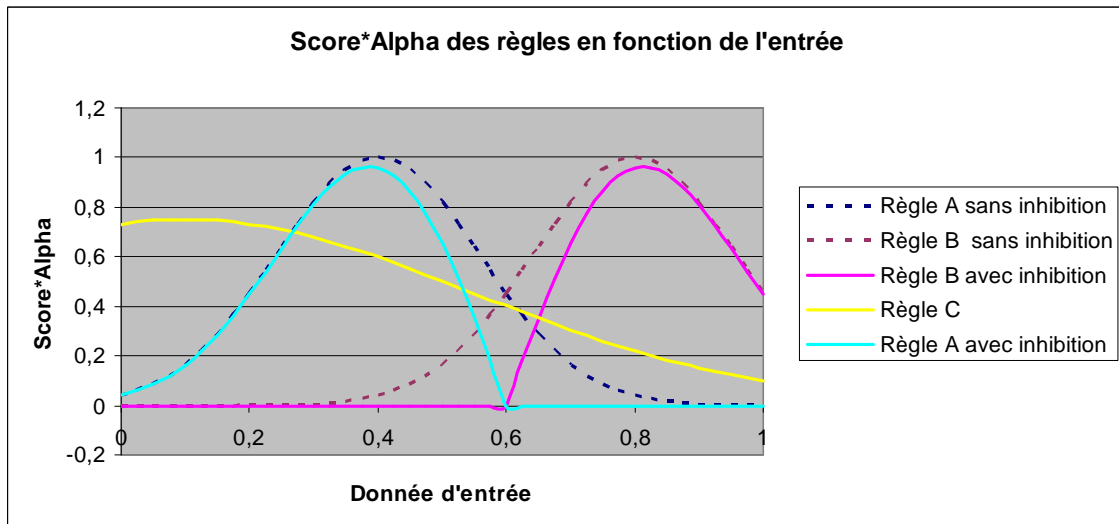


Figure III-35 : Variation des scores entre les règles A, B et C en fonction de la valeur du message d'entrée et avec un Alpha de 1 pour les deux premières règles et de 0,8 pour la dernière.

**

Ces quelques exemples ou quelques pistes de structures algorithmiques ou comportementales illustrent le potentiel d'expressivité de l'architecture cognitive proposée. La conception des schèmes conditionne à la fois la manière de créer des règles et de les renforcer, autrement dit la conception des règles et des « métarègles » s'effectue avec le même formalisme. De plus, la maîtrise de l'évolution des règles, en manipulant la force initiale et le nombre d'ajustements, offre l'idée que la conception de schèmes cognitifs initiaux, les proto-croyances, peut être envisagée comme une sorte de programmation cognitive.

7. Conclusion

La philosophie pragmatiste a amené à considérer la cognition comme un processus dynamique se développant et se couplant en permanence avec le monde extérieur afin de se maintenir. Cette propriété dynamique faisant écho à la notion d'autopoièse, il a été proposé de définir la cognition comme une autopoièse. Plus précisément, la formalisation a conduit à distinguer deux sortes d'autopoièse, l'autopoièse physique correspondant à la notion classique et l'autopoièse sémiotique.

Cette distinction a permis d'une part de mieux comprendre l'enfermement conceptuel de l'interactionnisme varélien et d'autre part de révéler l'importance d'une définition précise de la sémiose. La sémiose minimale comprise comme une interaction entre deux mémoires l'une événementielle et l'autre épistémique permet de définir une architecture cognitive de base. Les propriétés fondamentales de cette architecture ont été élaborées à partir du constat phénoménologique de la perception issue de l'illusion ou de l'hallucination. Aidée par la perspective évolutionniste de la cognition, une distinction des rôles entre les mécanismes subcognitifs et les schèmes cognitifs a été proposée.

Une comparaison avec les systèmes logiques traditionnels a permis de vérifier que la formalisation proposée pouvait a priori les contenir et que la notion de création de règles et d'auto-rétribution permettait d'espérer la construction de règles et de métarègles dans un même formalisme offrant ainsi des perspectives pour une évolution vers l'abstraction.

Le paradigme de la cognition défini, la recherche d'une modélisation de l'architecture cognitive a mis en avant la famille d'algorithmes de type système de classeurs dont la structure interne possède une certaine analogie avec celle proposée. Mais l'architecture proposée se distingue des autres principalement par une rétribution native non orientée, une notion de rétribution explicite, une gestion temporelle des événements, un typage des messages et l'existence de prémisses inhibitrices. Bien que ne définissant pas les mécanismes de création de règle, la spécification de ce premier modèle est suffisamment précise pour entreprendre son implémentation qui permettra de valider les propriétés d'auto-organisation et de couplage ainsi que de vérifier l'importance de l'auto-rétribution pour ces deux propriétés.

CHAPITRE IV ÉTUDE COMPORTEMENTALE DE L'IMPLÉMENTATION ROBOTIQUE

« HAMM. — *Salopard ! Pourquoi m'as-tu fait ?*

NAGG. — *Je ne pouvais pas savoir.*

HAMM. — *Quoi ? Qu'est-ce que tu ne pouvais pas savoir ?*

NAGG. — *Que ce serait toi. »*

Fin de Partie p. 91 de Samuel Beckett (1957).

« *Je ne choisis pas d'être, mais je suis. Une absurdité responsable d'elle-même, voilà ce que je suis »*

Le sang des autres p. 101 Simone de Beauvoir (1945).

« LE RENARD. — *Les hommes ont oublié cette vérité, dit le renard. Mais tu ne dois pas l'oublier. Tu deviens responsable pour toujours de ce que tu as apprivoisé. Tu es responsable de ta rose...*

LE PETIT PRINCE — *Je suis responsable de ma rose... répéta le petit prince, afin de se souvenir. »*

Le petit Prince p. 72 d'Antoine de Saint-Exupéry (1943).

1. Introduction

La généalogie de la cognition et l'analyse de l'autopoièse sémiotique montrent dans les chapitres précédents que la cognition artificielle comme la cognition naturelle se construit graduellement. Dans cette démarche, l'étude empirique de l'autopoièse sémiotique portera ici uniquement sur la réalisation d'une autopoièse sémiotique minimale, c'est-à-dire la réalisation d'un système robotique possédant des capacités cognitives relatives aux deux premiers stades de la cognition. Le premier stade consiste à établir une segmentation propre du monde en dehors de tout but grâce à l'auto-organisation des règles par ajustement de leur prémisses. À ce stade, seul l'opérateur de reconnaissance spontanée existe. Le deuxième stade alterne les phases d'instinct et de dressage se traduisant par la création de règles et un mécanisme de rétribution. À ce stade, la rétribution peut être associée à la capacité de prédire qui repose sur l'opérateur de reconnaissance attendue. La validation de l'architecture consistera en l'observation des propriétés émergentes à ces deux stades, au travers d'une réalisation robotique. Comparée aux autres approches en robotique cognitive, la validation ne correspond pas à la réalisation d'une tâche particulière mais à l'étude de l'interaction entre l'environnement et la dynamique sémiotique du système.

Les principes de la cognition dégagés dans les chapitres précédents se comprennent indépendamment du sujet cognitif et de son environnement puisqu'ils proviennent uniquement de l'idée de la génération d'une interaction entre les deux. Dans ce cadre, l'architecture proposée se veut également indépendante des propriétés physiques du robot et de celles de son environnement. Cependant, toute modélisation d'un système cognitif repose obligatoirement sur des primitives motrices et sur des primitives sensorielles ainsi que sur des comportements réactifs de base (ou schèmes cognitifs initiaux) dont le développement résulte de l'ontogenèse et dont la sélection provient de la phylogenèse.

Les schèmes cognitifs initiaux doivent être élémentaires afin de maîtriser au mieux l'analyse des données de ce système à multiples dimensions : évolution de la population des règles, évolution des règles et de leurs prémisses, sans oublier l'évolution du comportement du robot. Le schème cognitif élémentaire choisi représente un comportement réactif de base qui peut s'assimiler à une boucle sensorimotrice de second ordre, soit un ensemble de règles sensorimotrices homogènes. La boucle sensorimotrice secondaire se fondera sur un seul type de senseur et d'effecteur. Par conséquent, la détermination des primitives se limite à une primitive motrice et à une primitive sensorielle. Concrètement, ces primitives représentent les traitements ou les mécanismes reliant les capteurs à la valeur d'une prémisses sensorielle et reliant la valeur d'une conclusion motrice à l'effet.

L'étude expérimentale du système cognitif développé se focalisant sur l'évolution dynamique des règles et non sur la réalisation d'une tâche, l'étude préliminaire sur les caractéristiques du robot, de son environnement et de ses règles initiales devient indispensable pour l'interprétation des données. Cette étude préliminaire sera l'objet des deux premières parties de ce chapitre. La première partie se consacrera à la présentation générale du robot et à la caractérisation fine des capteurs dans les conditions expérimentales. Les différents types de contrôles moteurs seront également abordés afin de déterminer une primitive motrice. La seconde partie portera sur l'élaboration de la boucle sensorimotrice assurant un comportement exploratoire ainsi que sur les caractéristiques statistiques de l'environnement induites par celui-ci. De cette dernière étude sera alors déterminée une primitive sensorielle.

Le système cognitif développé à partir de l'architecture cognitive représente un système dynamique avec une boucle de rétroaction. Néanmoins, le système cognitif ne constitue pas un système de régulation puisqu'il ne possède pas de consigne. La rétroaction se réalise par le biais des changements sensoriels engendrés par l'accomplissement des commandes motrices issues des règles sélectionnées. Classiquement, l'analyse du comportement du système rétroactif s'effectue dans un premier temps avec la boucle de rétroaction ouverte pour observer le mode propre du système puis dans un second temps avec la boucle de rétroaction fermée. Les deux autres parties de ce chapitre présenteront respectivement ces deux phases dans l'étude empirique du système cognitif.

La troisième partie de ce chapitre correspondant à la première phase de cette étude examinera alors la dynamique du système indépendamment des conséquences induites par les conclusions motrices en appliquant un scénario sensoriel. Le contrôle total sur les entrées du système permettra alors de vérifier le rôle des paramètres fondamentaux et l'influence de la fréquence de la stimulation sur le maintien des règles. Par ailleurs, la capacité à classifier l'environnement par auto-organisation en dehors du retour moteur sera également évaluée ainsi que l'apport de la création de règles dans la classification. Enfin, l'ordonnancement des règles constituant la boucle prédictive décrite dans le chapitre précédent sera vérifié. La quatrième partie de chapitre qui représente l'étude en boucle fermée se consacrera à la stabilité du système robotique en environnement réel simple ou complexe ainsi qu'à l'influence du mécanisme de rétribution par prédiction. La conclusion s'efforcera de synthétiser ces résultats et de proposer une discussion sur leurs implications.

2. Présentation du robot utilisé

2.1. Les généralités

Le robot choisi appartient à la famille des robots Khepera de la première génération. Développé au Laboratoire de Micro Informatique (LAMI) de l'École Polytechnique Fédérale de Lausanne EPFL en 1994, ce robot a connu un large succès au sein de la communauté de la cognition artificielle aussi bien pour l'enseignement que pour la recherche. En effet, ce robot présente l'avantage d'être de petite taille (55mm de diamètre pour une hauteur de 30 mm et une masse de 70 g) et d'avoir une conception modulaire matérialisée par des étages ou tourelles qui s'enfichent les une sur les autres. Ainsi, ce type de robot permet de tester facilement en conditions réelles des algorithmes développés tels que, par exemple, la planification de trajectoire, l'évitement d'obstacles, le prétraitement des données des capteurs ainsi que l'évolution du comportement individuel ou collectif. La diversité des thèmes de recherche utilisant ce type de robot reflète la diversité des approches en sciences cognitives.

Le robot Khepera utilisé possède un processeur Motorola 68331 avec 256 ko de RAM et 18 Ko de ROM. Cependant cette capacité de calcul s'avère insuffisante pour la gestion du modèle cognitif envisagé. Par conséquent, le contrôle du Khepera a toujours été effectué via une communication port série RS232, obligeant à maintenir un fil entre le robot et l'ordinateur maître. Pour des raisons de connectiques, une tourelle avec caméra a été rajoutée (Figure IV-1), bien que cette dernière n'ait jamais été utilisée. Par ailleurs, malgré l'emploi d'un raccord adapté et d'un système pour maintenir le fil, ce dernier provoque parfois une gêne lors du déplacement du robot compte tenu de la faible puissance des moteurs.

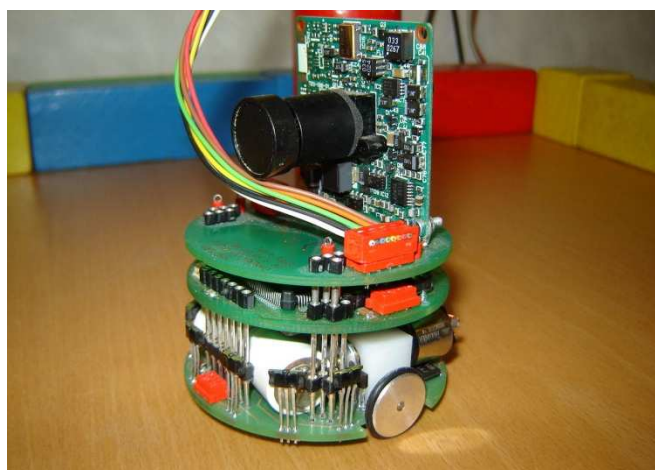


Figure IV-1 : Le robot Khepera.

Le robot khepera dispose de deux roues, soit une configuration de type « char » classique en robotique mobile. Chacune des roues possède un moteur à courant continu muni d'un codeur permettant une mesure de la vitesse et de l'odométrie. Le robot khepera peut rouler au maximum à 1 m/s et au minimum à 0,02 m/s. Les capacités sensorielles se composent de huit capteurs infrarouges passifs qui mesurent la luminosité ambiante et de huit capteurs infrarouges actifs qui mesurent la réflectance des obstacles. Dans l'étude

proposée, seuls les capteurs actifs ont été utilisés afin de faciliter l'analyse comportementale. La caractérisation des capteurs se porte alors exclusivement sur ces derniers.

2.2. La caractérisation des capteurs infrarouges

Les capteurs infrarouges possèdent beaucoup de facteurs d'incertitude, d'une part la verticalité des capteurs par rapport au corps du Khepera ne peut être assurée, ce qui a une incidence sur la mesure, et d'autre part les caractéristiques intrinsèques des capteurs se révèlent très variables. (A) Afin d'apprécier cette variabilité, une étude sur les capacités de détection des capteurs sera entreprise dans un premier temps. Puis, dans un second temps, (B) une évaluation des bruits perturbant la mesure des capteurs sera présentée.

A - Étude de la sensibilité des capteurs infrarouges

L'étude concernant la sensibilité des capteurs infrarouges a été menée avec des obstacles deux types de matériaux : le bois vernis qui constitue le matériau de l'enceinte expérimentale et le polystyrène, le matériau d'une première enceinte prototype.

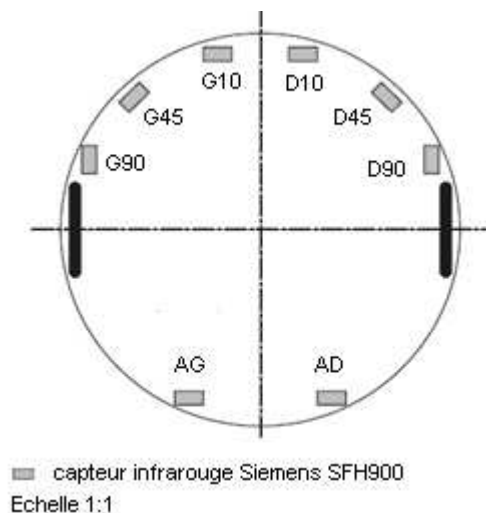


Figure IV-2 : Position des 8 capteurs infrarouges.

La première expérience de caractérisation s'effectue avec un obstacle placé successivement et perpendiculairement à 50cm devant chacun des capteurs (Figure IV-2), sachant que la valeur traduisant la quantité d'infrarouge réfléchi varie entre 0 et 1023. Le Tableau IV-1 montre que l'absorption des infrarouges dépend du matériau de l'obstacle et souligne également la sensibilité des différents capteurs.

	G90	G45	L10	D10	D45	D90	AD	AG
Moyenne avec le bois	68	19	3	83	211	103	35	197
Écart type avec le bois	51	27	7	55	75	92	36	54
Moyenne avec le polystyrène	210	118	163	312	368	311	49	259
Écart type avec le polystyrène	14	9	10	7	3	4	7	4

Tableau IV-1 : Comparaison de la réflexion infrarouge entre le bois vernis et le polystyrène pour chaque capteur.

La variabilité des capteurs s'illustre plus particulièrement avec la Figure IV-3 qui représente les valeurs des capteurs suivant l'éloignement à intervalles réguliers d'un obstacle en polystyrène.

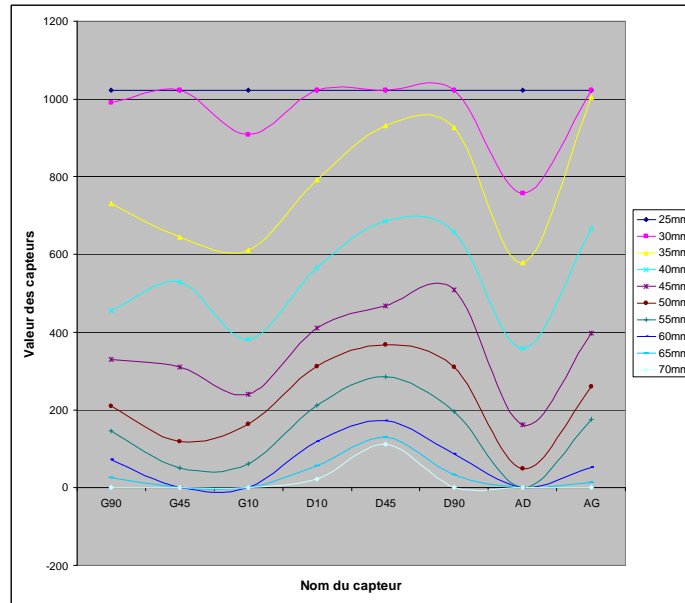


Figure IV-3 : Disparité de la sensibilité des capteurs (sur l'exemple de l'évaluation de la distance sur un objet en polystyrène)

Une autre manière de visualiser cette disparité consiste à tracer la courbe (Figure IV-4) des valeurs moyennes de chaque capteur en fonction de la distance de l'obstacle, soit une représentation plus intuitive de la portée des capteurs infrarouges.

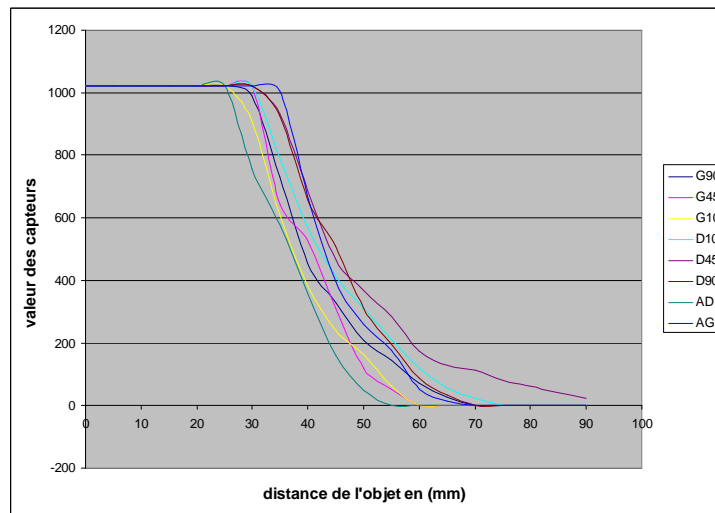


Figure IV-4 : Réponse des capteurs IR en fonction de la distance à un obstacle en polystyrène.

Toutefois, malgré la variabilité des capteurs, la nature d'un obstacle peut être déduite en prenant la moyenne des capteurs comme le montre la Figure IV-5. Dans la zone

pertinente pour la détection fine d'obstacles situés entre 35 mm et 55 mm, la différence entre les valeurs moyennes des capteurs atteint environ 250 pour une même distance. Ainsi approximativement, 8mm séparent deux valeurs moyennes identiques mais issues de matériaux différents.

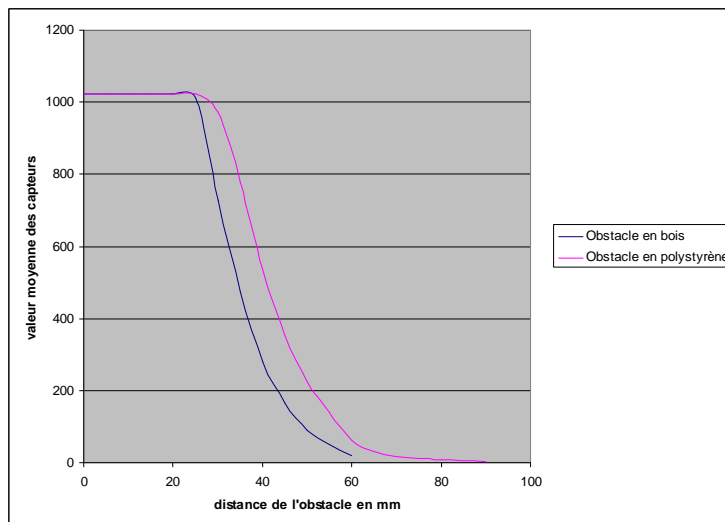


Figure IV-5 : Courbe moyenne des capteurs IR en fonction de la nature de l'obstacle.

B - Évaluation des bruits sur la mesure

En plus de la variabilité de la sensibilité des capteurs, l'interprétation de la mesure se complexifie à cause du bruit. Ce dernier provient de plusieurs sources, le bruit intrinsèque du capteur et le bruit extrinsèque provenant d'une radiation infrarouge. Le bruit intrinsèque varie selon la courbe de réponse du capteur, c'est-à-dire que l'écart type a tendance à augmenter en même temps que la valeur moyenne, comme le montre par exemple la Figure IV-6 sur le capteur D45.

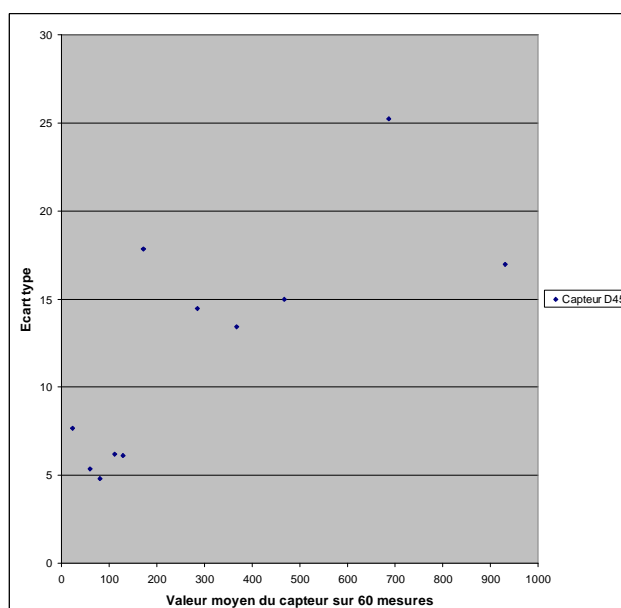


Figure IV-6 : Écart type en fonction de la valeur moyenne du capteur D45.

Le bruit extrinsèque peut prendre des formes très différentes par exemple, des moteurs se trouvant à proximité chauffent les capteurs, modifiant ainsi leurs courbes de réponse, mais ce bruit demeure négligeable. En revanche, les lampes à incandescence ou à néon génèrent un rayonnement infrarouge produisant une forte perturbation. Plus précisément, le bruit de la lampe à incandescence utilisée prend la forme d'un signal oscillant d'une amplitude d'environ 200 et d'une fréquence 1,5 Hz pour une fréquence d'échantillonnage de 10 Hz qui a été employée pour toutes les expériences (Figure IV-7).

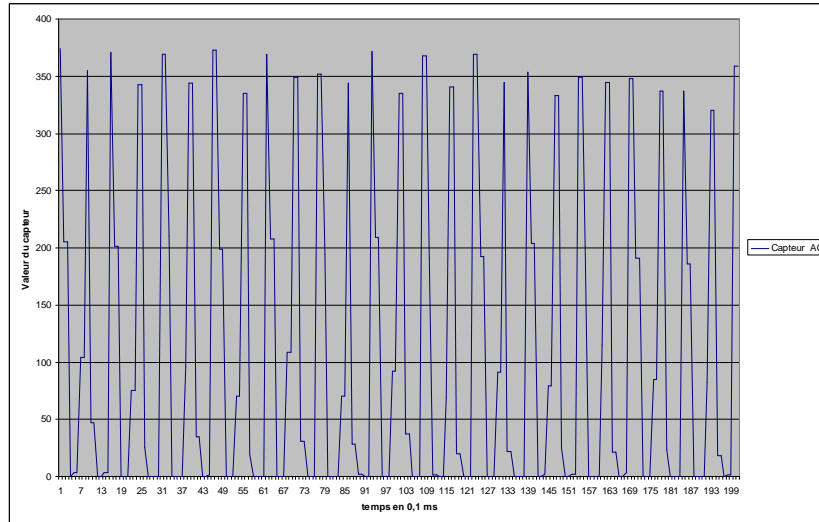


Figure IV-7 : Mesure du capteur AG avec une lampe et un obstacle à proximité.

Afin de mieux percevoir l'incidence d'un tel bruit, la Figure IV-8 et la Figure IV-9 montrent réciproquement deux séries de mesures acquises avec et sans lampe à proximité lors de l'évolution du robot dans l'enceinte selon un algorithme d'évitement d'obstacle de type Braitenberg. L'exploration de l'environnement résulte ici d'un comportement d'évitement d'obstacle.

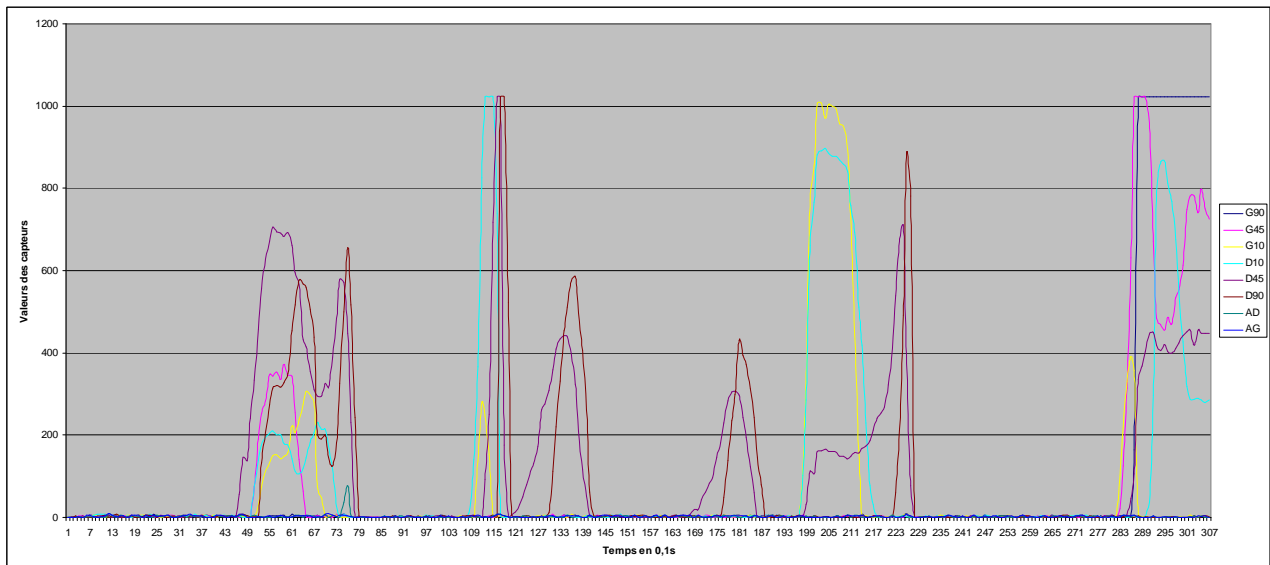


Figure IV-8 : Exploration de l'environnement sans lampe à proximité.

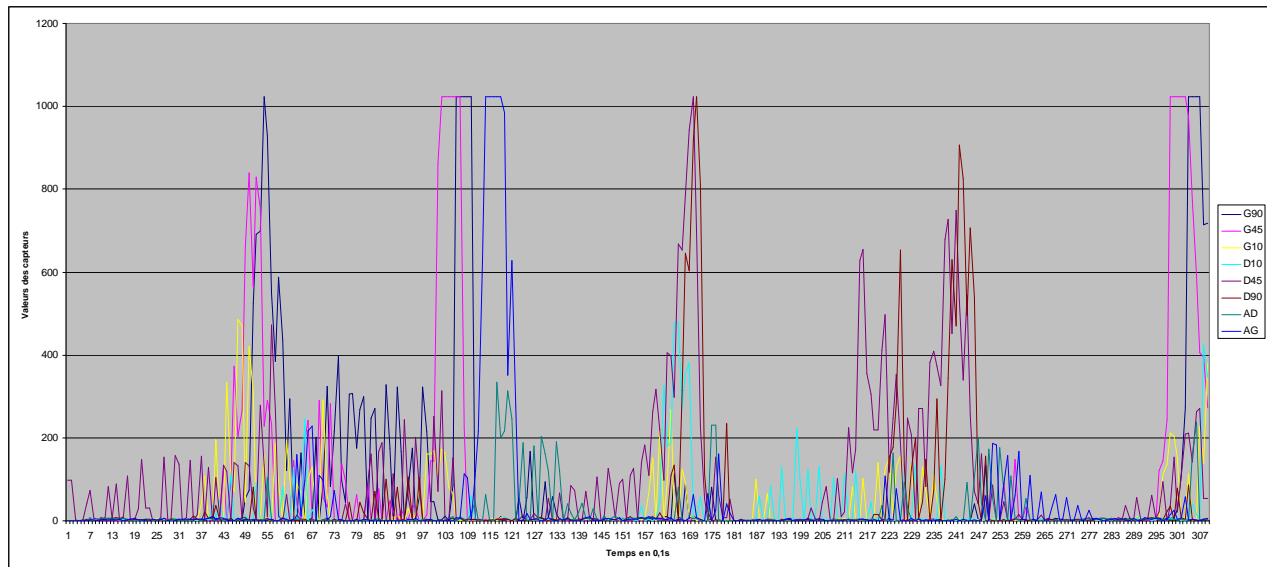


Figure IV-9 : Exploration de l'environnement avec lampe à proximité.

Le dernier facteur influençant la mesure réside dans la luminosité ambiante qui ne provoque pas de changement radical dans la durée des expériences mais peut biaiser la comparaison des résultats entre celles-ci. Le déroulement des expériences a donc été effectué dans la pénombre afin d'avoir une luminosité constante et sans éclairage direct évitant ainsi au maximum les radiations infrarouges émises. Ces conditions expérimentales ont permis de minimiser les bruits de telle sorte que l'écart type d'une valeur moyenne d'un capteur vis-à-vis d'un obstacle situé dans la zone de pertinence représente environ 2% de cette valeur, au lieu de 5 à 8% sans précaution.

2.3. Les différents types de contrôles moteur

Le moteur du Khepera offre trois modes de contrôle (Figure IV-10) : le contrôle direct de la puissance motrice, l'asservissement de position et l'asservissement de vitesse.

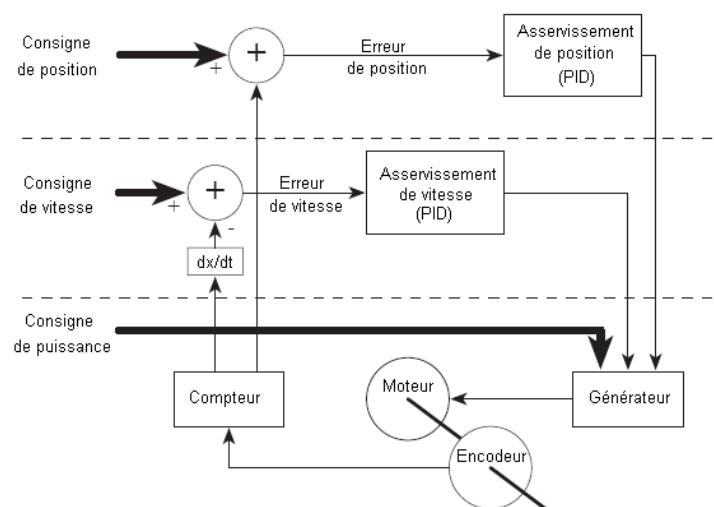


Figure IV-10 : Schéma représentant les trois types de contrôle moteur.

Le premier mode de contrôle consiste à envoyer directement la puissance voulue aux deux moteurs. Cependant, la relation entre la puissance fournie et la vitesse motrice n'est pas linéaire puisque celle-ci dépend notamment de l'inertie du moteur. Cette relation non linéaire empêche de déterminer la sortie motrice avec la simple connaissance de la puissance délivrée. Le contrôle direct de la puissance motrice doit alors s'effectuer par le biais d'un calcul d'asservissement qui prend en compte la dynamique du moteur. Par ailleurs, pour une même puissance délivrée, la courbe de réponse motrice diffère entre les deux moteurs. Cette variabilité implique, par exemple, que le robot tournera au lieu d'aller tout droit lors d'un départ arrêté avec une consigne de puissance identique sur chaque moteur.

Les deux autres modes de contrôle sont des asservissements mis à la disposition de l'utilisateur directement à partir de la ROM. Ces outils de gestion de bas niveau permettent une certaine abstraction dans la gestion du déplacement du robot. Autrement dit, à partir d'une consigne de vitesse ou de position, la puissance délivrée au moteur se calcule automatiquement pour atteindre cette consigne dans les meilleures conditions en termes de stabilité et de rapidité. Toutefois, la vitesse et la position mesurées à partir des encodeurs de chaque roue ne permet pas d'inférer la vitesse et la position exacte du robot. En effet, les déplacements du robot sont non holonomes, en d'autres termes, les glissements ou les frottements interdisent toute relation entre les actions motrices et la position précise du robot dans l'espace. Pour ces raisons, le contrôle direct de la puissance ne sera pas choisi comme primitive motrice.

Le type d'asservissement utilisé est un Proportionnel-Intégral-Dérivé (PID). Plus précisément, la commande appliquée au moteur provient de la combinaison pondérée de trois fonctions appliquées à l'erreur, soit l'écart entre la consigne et la position ou la vitesse mesurée. L'asservissement permet de gommer les différences entre les moteurs mais pas de les effacer. En reprenant l'exemple du départ arrêté, mais cette fois-ci avec une consigne de vitesse, le robot tournera également au début car les courbes de réponses restent différentes mais il redressera sa trajectoire pour terminer en ligne droite.

Bien que ce comportement soit quasiment identique pour les deux types d'asservissement (de vitesse ou de position), seul l'un des deux correspond à l'idée d'une primitive motrice. En effet, l'asservissement de position implique une définition de l'action puisque le début et la fin du mouvement sont définis. Or, le chapitre précédent insiste sur le fait que l'action, en dehors des arcs réflexes, ne doit pas être confondue avec les primitives motrices. Celles-ci ne doivent pas être finalisées afin de pouvoir construire des actes finalisés. Par conséquent, parmi les trois modes de contrôle disponibles, seul l'asservissement vitesse correspond à une primitive motrice.

L'inconvénient est que l'algorithme d'asservissement se trouve intégré à bas niveau dans le système robotique, n'offrant pas la possibilité de recueillir la valeur de la puissance finalement délivrée. Or, cette dernière peut être considérée comme un indice de l'effort nécessaire pour atteindre la consigne. En effet, à consigne égale, la puissance nécessaire est plus importante dans le cas où le robot se trouve dans une montée que sur une surface plane. La différence de puissance peut également traduire la présence d'un obstacle mobile ou fixe. Autrement dit, en plus de la vitesse motrice, une information proprioceptive sur la puissance peut se révéler très intéressante pour l'étude de la genèse de la kinesthésie. Cette voie n'a pas fait l'objet d'études expérimentales, néanmoins, elle sera reprise dans les perspectives présentées à la fin de ce mémoire. Toutefois, par anticipation, une bibliothèque d'asservissement a été réalisée permettant ainsi d'en maîtriser toute la chaîne.

Le programme développé s'appuie sur cette bibliothèque et dans les expériences qui seront présentées, les commandes motrices prendront plus précisément la forme de consignes de vitesse. Indépendamment de l'arrivée de nouvelles consignes, le réajustement de la puissance délivrée s'effectue à une fréquence de 10 Hz.

La prise en compte de la proprioception, soit au niveau de l'asservissement, soit au niveau du système à base de règles, permet d'illustrer la différence évoquée dans le chapitre précédent entre la notion de boucle sensorimotrice primaire et celle de boucle sensorimotrice secondaire. La première se trouve définie par une primitive motrice prenant en compte un retour sensoriel (comme l'asservissement par exemple). La boucle sensorimotrice primaire se révèle immuable dans le sens où les caractéristiques régissant la primitive motrice (comme les paramètres du PID) ne peuvent être modifiées par un processus cognitif. En revanche, la boucle sensorimotrice secondaire s'exprime par un ensemble de règles sensorimotrices simples qui se couple avec l'environnement. Dans ce chapitre, par commodité de langage, l'expression boucle sensorimotrice renverra uniquement à la notion de boucle sensorimotrice secondaire.

3. Élaboration des règles initiales

3.1. Construction de la boucle sensorimotrice

La construction des règles sensorimotrices se déroule en deux phases : (A) la détermination du comportement de base et (B) la modélisation du comportement de base par un ensemble de règles.

A - Le comportement réactif

Le choix du comportement de base repose principalement sur deux critères : la capacité d'explorer son environnement sans se mettre en péril et la simplicité dans le déclenchement du comportement. En d'autres termes, le dernier point signifie que la prémisse d'une règle ne doit, dans un premier temps, posséder qu'un seul message sensoriel. Ainsi, le comportement de base doit dépendre uniquement de l'état actuel des capteurs, un comportement purement réactif.

Dans ce contexte, le comportement d'évitement d'obstacle de type Braitenberg (1984) semble le plus adapté à ces exigences. Le principe d'algorithme consiste à établir une relation linéaire entre les entrées sensorielles et la consigne motrice d'un moteur. Pour le Khepera, cela se traduit par le produit matriciel suivant :

$$\begin{aligned} \text{Consigne}_G &= [(G90, G45, G10, D10, D45, D90, AD, AG) * (0,0,0,-12,-10,0,5,3)]/1100 + 5 \\ \text{Consigne}_D &= [(G90, G45, G10, D10, D45, D90, AD, AG) * (0,-12,-10,0,0,3,5)]/1100 + 5 \end{aligned}$$

L'unité de la consigne est de 8 mm/s. L'ajout de la composante continue (+5) permet au Khepera d'avancer par défaut (c'est-à-dire si tous ses capteurs sont à 0) à une vitesse de 40 mm/s. Cette configuration confère au Khepera un comportement d'évitement d'obstacle ne prenant pas en compte les capteurs situés sur le côté et, par conséquent, le Khepera aura tendance à suivre les murs. Par ailleurs, une relation plus riche entre les senseurs et l'effecteur pourra davantage mettre en valeur les capacités d'auto-organisation entre les règles. Cet enrichissement de la relation entre senseur et effecteur se traduit par le produit matriciel suivant :

$$\text{Consigne}_G = [(G90, G45, G10, D10, D45, D90, AD, AG) * (4, 4, 6, -18, -15, -5, 5, 3)] / 1100 + 5$$

$$\text{Consigne}_D = [(G90, G45, G10, D10, D45, D90, AD, AG) * (-5, -15, -18, 6, 4, 4, 3, 5)] / 1100 + 5$$

D'un point de vue mathématique, les fonctions Consigne_G et Consigne_D sont surjectives, c'est-à-dire que plusieurs configurations sensorielles peuvent donner une même consigne. Par ailleurs, les actions peuvent être interprétées comme le résultat de la projection des capteurs, soit un vecteur de 8 dimensions dans un espace à 2 dimensions.

B - Extraction des règles de commandes initiales

Les règles de commandes initiales seront construites à partir de l'enregistrement des consignes motrices et des senseurs lors du déroulement du programme d'évitement d'obstacle de type Braitenberg. La construction des règles de commandes doit donc tenir compte de l'environnement utilisé pour les expériences. L'environnement utilisé est une enceinte de forme carrée de 30 cm de côté (Figure IV-11), constituée d'un assemblage de blocs de bois vernis de différentes couleurs. La couleur n'a pas d'incidence sur la détection de la présence des blocs par les capteurs infrarouges. Concernant la détermination de la taille de l'enceinte, les mesures de détectabilité avec les blocs de bois montrent que le Khepera ne capte pas la présence de l'obstacle au-delà de 6 cm. Cela reviendrait, pour un homme dépourvu de vision dont les bras tendus serviraient de capteurs de proximité, à évoluer dans une pièce carrée de 5 m de coté.

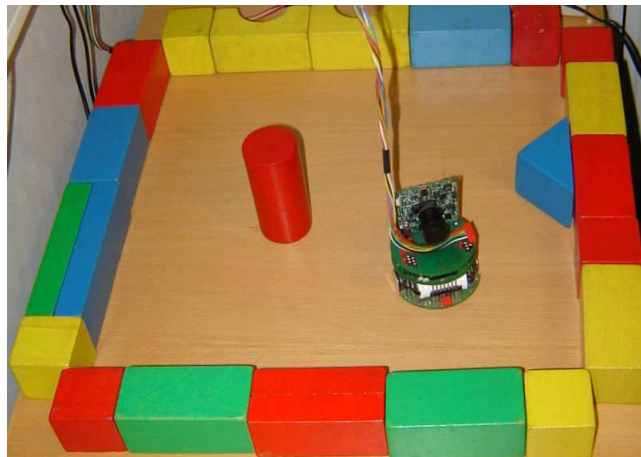


Figure IV-11 : Le robot Khepera dans l'enceinte d'expérimentation.

L'enceinte, illustrée par la Figure IV-11, est qualifiée de complexe dans le cadre des expériences. La complexité pour le robot vient de la présence du prisme bleu accolé au bord droit de l'enceinte et de la présence du cylindre rouge situé au centre gauche. La richesse induite par la présence de ces blocs est visible dans la comparaison entre deux séries de consignes enregistrées lors du déroulement du programme d'évitement d'obstacle, Figure IV-12 : l'une se déroule avec une acquisition à 10 Hz durant 30 min lorsque l'enceinte est vide, l'autre série est réalisée dans les mêmes conditions mais en présence du cylindre et du prisme. L'environnement riche stimule davantage le système bien que toutes les consignes générées par les situations nouvelles restent proches des valeurs produites avec un environnement simple. Toutefois, cette légère différence incitera à extraire la base de règles initiales dans les enregistrements issus de l'environnement complexe afin de maximiser la diversité.

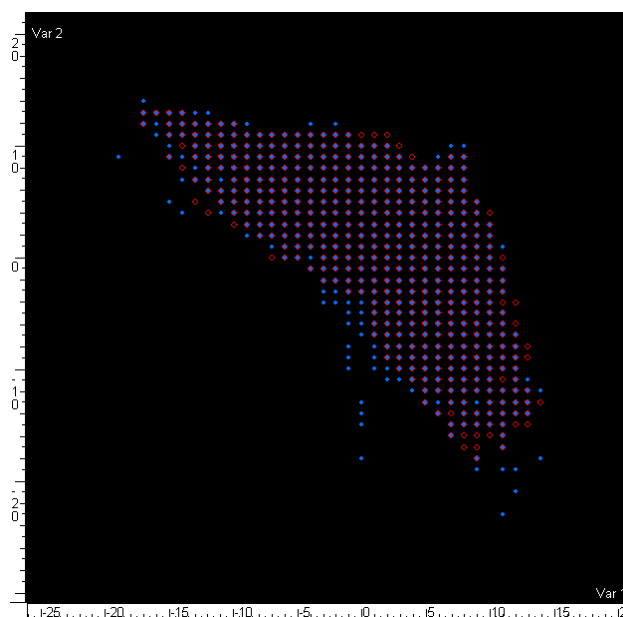


Figure IV-12 : Représentation des consignes motrices en 8mm/s ; Var2 : moteur gauche, Var1 : moteur droit. Les points bleus sont issus d'un environnement simple et les cercles rouges proviennent d'un environnement complexe.

Un tirage aléatoire ne permet pas de s'assurer de la représentativité minimale d'un comportement d'obstacle. En effet, le robot se trouve le plus souvent dans la situation où aucun obstacle ne barre son chemin. La solution choisie consiste à recueillir les consignes et les données des capteurs dans un laps de temps évalué à 30 s avec une fréquence d'acquisition de 20 Hz dans un environnement complexe.

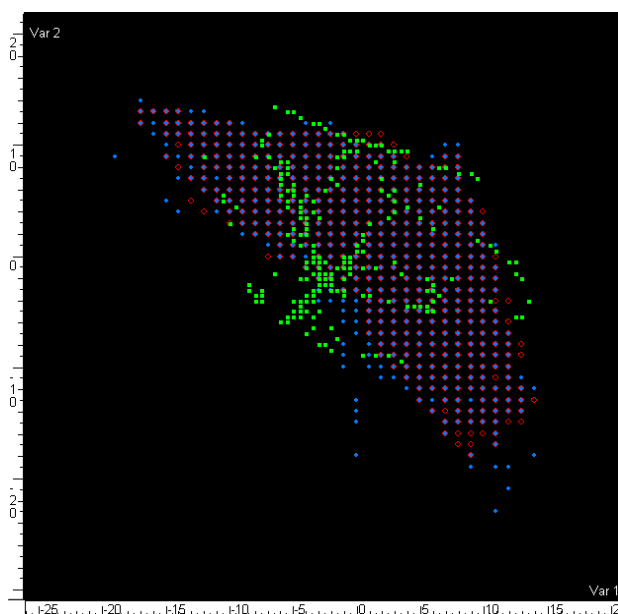


Figure IV-13 : Représentation des consignes motrices en 8mm/s. ; Var2 : moteur gauche, Var1 : moteur droit. Les points verts proviennent du premier essai (avec l'environnement complexe) ; les points bleus représentent le deuxième essai (avec l'environnement simple) et les cercles rouges représentent le troisième essai (avec l'un environnement complexe et avec une position initiale au bord de l'enceinte).

Ce procédé introduit dans la construction de l'ensemble de la base de règles initiales une hétérogénéité intéressante pour observer les capacités d'adaptation. Plus précisément, les consignes des règles sélectionnées ne couvrent pas de façon homogène l'espace des consignes, ainsi, une majorité de règles amènent à tourner plutôt à droite, les carrés verts sur la Figure IV-13. Cette dissymétrie entre la diversité des consignes amenant à tourner d'un côté ou de l'autre constitue une opportunité pour étudier la rééquilibration des règles sensorimotrices dans les expériences présentées dans les parties suivantes.

Par ailleurs, en comparaison avec les deux premiers essais (les cercles rouges et les ronds bleus) sur la Figure IV-13, le troisième présente davantage les consignes simultanément négatives (les points verts). Cela est dû à la position initiale du robot dans l'enceinte. Dans les deux premières expériences préliminaires, le robot a été positionné au centre de l'enceinte alors que dans la dernière, il a été positionné dos à un bloc. Les motivations de ce choix seront abordées au cours de l'analyse des prémisses des règles initiales.

Enfin, le choix sur le nombre de règles initiales nécessite également réflexion. En effet, après 30 min d'expériences avec une fréquence d'acquisition à 10 Hz, le nombre de couples capteurs/consignes s'élève à 18000. Un trop grand nombre de règles qui recouvrait l'ensemble de l'espace sensoriel ne permettrait pas la mise en valeur des capacités d'adaptation attendues du système. À l'inverse, un trop petit nombre de règles favoriserait le déclenchement de règles dans des situations très éloignées de son champ d'application d'origine, le risque étant d'effectuer une action inappropriée pour l'évitement d'obstacle. Après plusieurs essais, le choix de fixer le nombre de règles à 600 en environnement réel apparaît comme un bon compromis.

3.2. Caractérisation statistique de l'environnement

L'ensemble des règles associant une stimulation sensorielle et la consigne vitesse maintenant établies, il reste à déterminer la primitive sensorielle. Celle-ci vise à optimiser l'adéquation entre la prémisses d'une règle et la stimulation, en d'autres termes, à rendre la reconnaissance de configurations sensorielles moins sensible au bruit. La détermination des traitements incarnant cette primitive sensorielle repose sur des principes issus du traitement du signal qui nécessitent une étude statistique de l'environnement. Cette dernière se fondera sur l'analyse en composantes principales (ACP) qui fut introduite en 1901 par Pearson et développée par Hotelling en 1933. L'analyse en composantes principales constitue une méthode de réduction de l'espace des données. Dans un tableau de p variables numériques (en colonnes) décrivant n individus, l'ACP permet une représentation de ces n individus dans un sous-espace de l'espace P défini par les p variables. Autrement dit, elle définit k nouvelles variables, combinaison des p variables de l'espace initial, qui feraient perdre le moins d'information possible. Ces k variables s'intitulent des « composantes principales » et les axes qu'elles déterminent s'appellent les « axes principaux ». L'enregistrement des valeurs des capteurs au cours de l'exécution forme une matrice rectangulaire (8 colonnes, le nombre de capteurs et n lignes, le nombre d'enregistrements).

L'interprétation de cette analyse nécessite toutefois une étude préalable. En effet, dans la partie précédente, la différence entre les consignes produites avec et sans blocs supplémentaires au sein de l'enceinte montre indirectement que l'environnement a une incidence sur la richesse des données sensorielles. La compréhension de la nature de cette richesse sensorielle permettra d'interpréter les différences entre les ACP. Les deux

premières sections de la caractérisation statistique de l'environnement se consacreront respectivement à l'influence de l'environnement sur les données brutes et à l'analyse en composantes principales de celle-ci. La dernière section proposera la primitive sensorielle adoptée à partir de ces résultats.

3.2.1. L'influence de l'environnement sur les données brutes

L'étude préliminaire sur la distribution des données s'est effectuée à partir de trois sources. Les deux premières proviennent des expériences précédentes, celles avec environnement simple et celles avec environnement complexe. La troisième source correspond à une expérience réalisée dans un environnement simple sous l'éclairage d'une lampe afin d'observer si la complexité de l'environnement pouvait être équivalente à l'introduction d'un bruit. Soucieux d'assurer la significativité statistique, chaque colonne représentant un capteur des tableaux de données possède 18000 valeurs.

Traditionnellement, deux critères caractérisent la distribution d'un vaste ensemble de points. Le premier, le Kurtosis, mesure le degré d'écrasement de la distribution d'une variable aléatoire réelle par rapport à une distribution normale. Le second, le Skewness, mesure le degré d'asymétrie de la distribution d'une variable aléatoire réelle par rapport à une distribution normale. Concernant les trois types de données, les résultats élevés des critères (Figure IV-14) suggèrent que les données sont peu gaussiennes bien que celles issues de l'environnement complexe se rapproche globalement d'avantage de la gaussiannité que les deux autres.

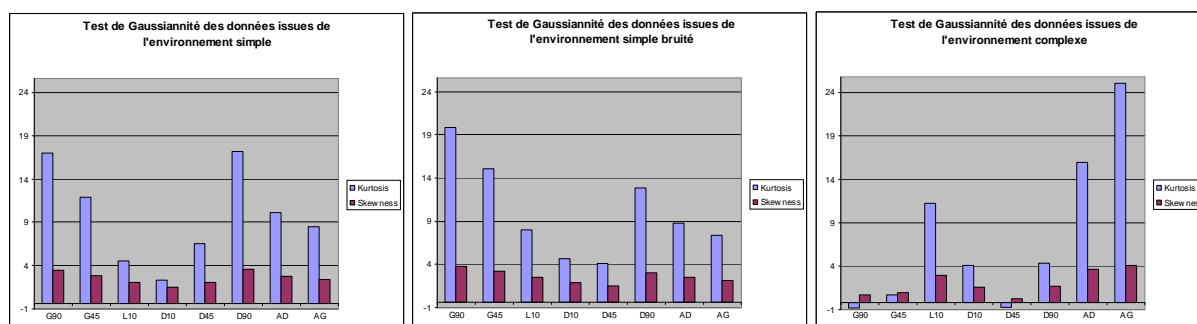


Figure IV-14 : Évaluation pour chaque capteur des critères de gaussiannité.

Toutefois, ces deux critères restent sensibles au point aberrant et la visualisation des données devient indispensable pour vérifier ce point. Afin de contourner la difficulté à fournir une représentation d'un espace à huit dimensions ; Figure IV-15, Figure IV-16 et Figure IV-17 proposent une visualisation deux à deux des variables. Cela se traduit par un tableau de graphiques où la diagonale est l'histogramme des données de la variable indiquée. La construction de cet histogramme s'appuie sur la méthode « Averaged Shifted Histograms » proposée par Scott (1992) afin de pallier le problème lié à la représentation de variables continues ou à forte résolution. Cette méthode calcule plusieurs histogrammes en utilisant la même largeur d'intervalle, mais avec des origines différentes. Les résultats ramenés à une moyenne sont tracés. Cet algorithme possède principalement deux paramètres : le nombre d'intervalles et le nombre d'histogrammes à calculer. Le logiciel de Visualisation Ggobot fixe le premier à deux et le second, paramétrable par l'utilisateur, a été

mis à 20. Ainsi, les histogrammes visualisés résultent de la moyenne de 20 histogrammes décalés.

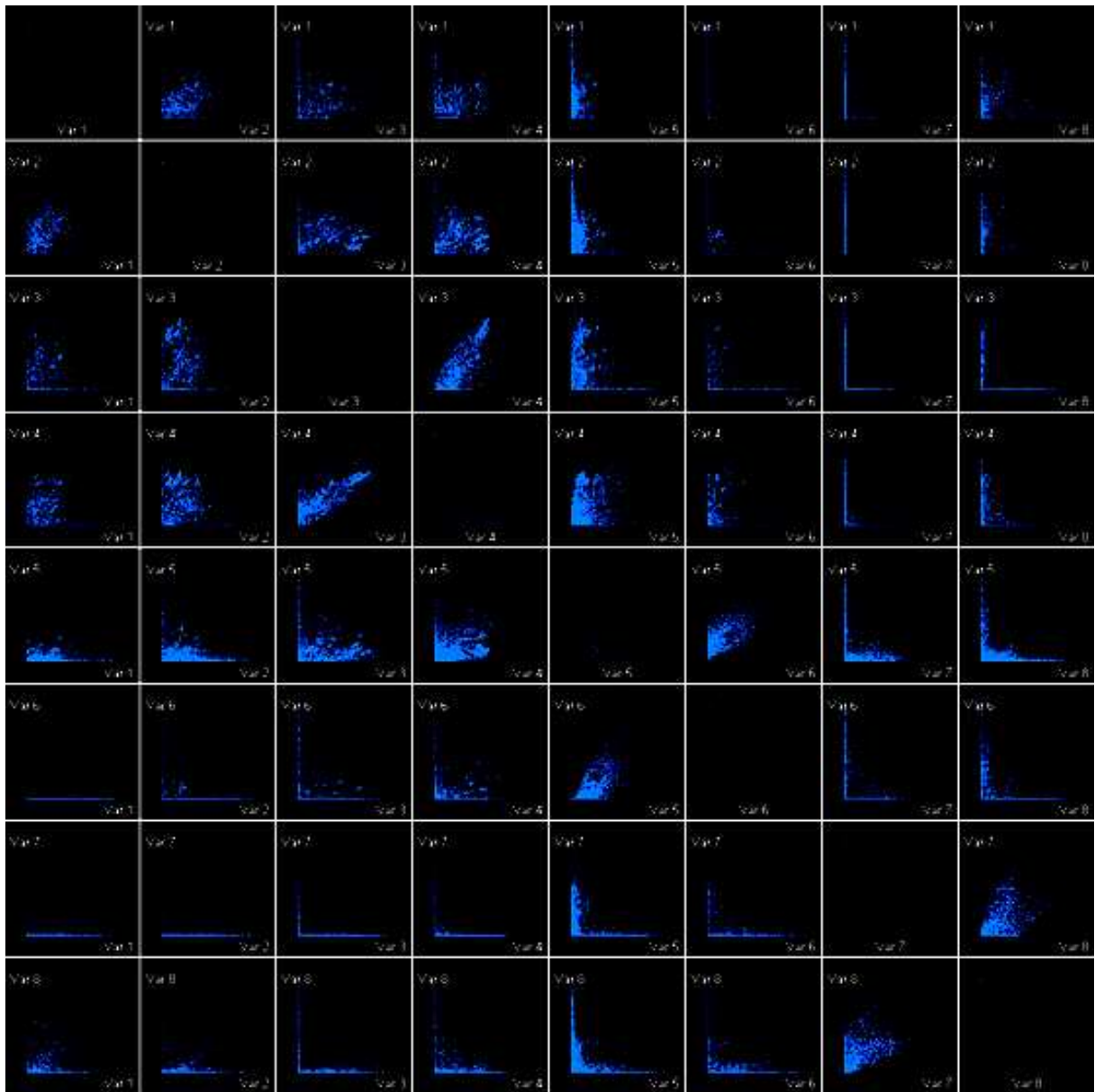


Figure IV-15 : Représentation des mesures dans un environnement simple, la numérotation des variables suit l'ordre : G90, G45, G10, D10, D45, D90, AD, AG.

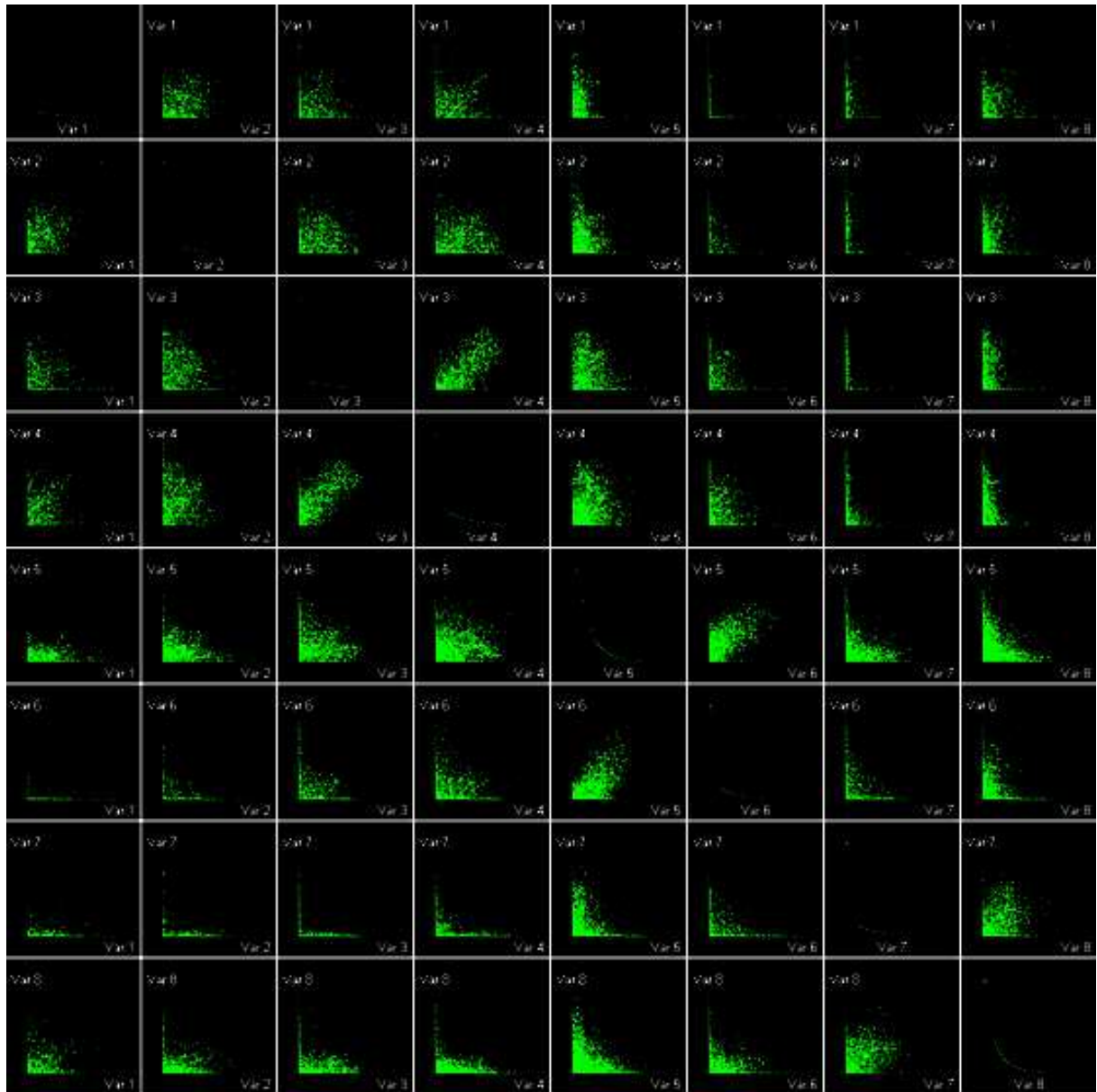


Figure IV-16 : Représentation des mesures dans un environnement simple bruité, la numérotation des variables suit l'ordre : G90, G45, G10, D10, D45, D90, AD, AG.

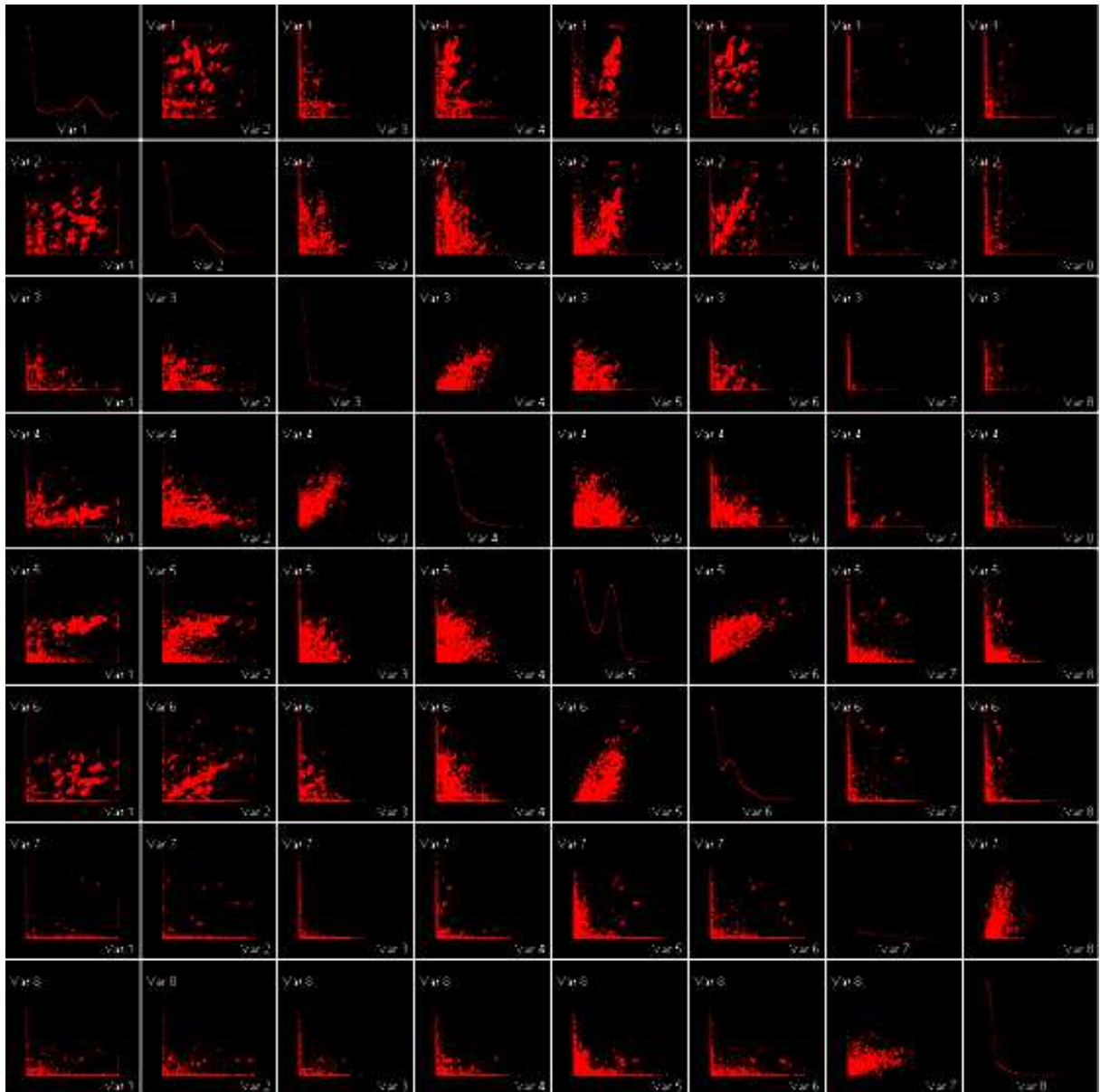


Figure IV-17 : Représentation des mesures dans un environnement complexe, la numérotation des variables suit l'ordre : G90, G45, G10, D10, D45, D90, AD, AG.

La visualisation (Figure IV-15, Figure IV-16 et Figure IV-17) des trois types de données confirme que la répartition des données ne correspond pas à une distribution gaussienne. Malgré tout, deux observations communes aux trois types de données peuvent être avancées : d'une part le centrage des données sur l'origine des axes confirme que le robot se trouve souvent éloigné des obstacles et d'autre part la présence de nombreux points sur les axes correspondant à une variable nulle montre que le robot ne se trouve jamais totalement entouré d'obstacle. Il y a donc toujours un ou des capteurs avec une valeur proche de 0. Ces deux remarques confortent l'idée que l'enceinte est assez grande pour une exploration spatiale de son environnement. Cette visualisation des données révèlent également des différences entre les trois types de données, qui peuvent se comprendre comme l'opposition entre l'environnement simple (Figure IV-15) et

l'environnement simple bruité (Figure IV-16) et l'opposition entre ces deux environnements et l'environnement complexe (Figure IV-17).

La première opposition s'appuie principalement sur trois différences. La première se situe sur la dispersion des points, plus grande dans l'espace des capteurs dans l'environnement simple bruité. La seconde porte sur la différence du nombre de valeurs nulles avec un facteur de 2,5 pour l'environnement simple bruité. La Figure IV-18 représente, par exemple, les histogrammes des trois types de données sur les 20 premières valeurs du capteur G45, soit la Var2.

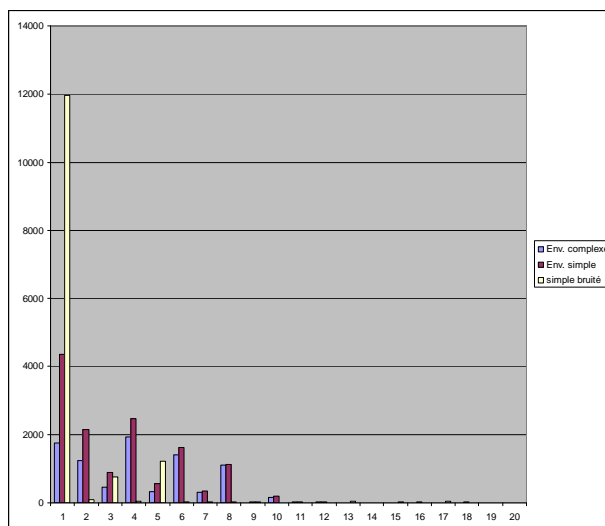


Figure IV-18 : Histogramme du capteur G45 sur les 20 premiers intervalles de valeurs.

Enfin, la troisième différence réside dans la disparition ou l'atténuation d'amas de points visibles dans les données issues de l'environnement simple, se traduisant par l'aplatissement de l'histogramme de la quatrième variable (Var4), par exemple, dans la Figure IV-16. Ces deux dernières différences s'expliquent par le fait que le bruit généré par la lampe augmente les valeurs des capteurs, ce qui a pour incidence d'augmenter la réactivité dans l'évitement d'obstacle. Le robot change ainsi de direction aussitôt qu'un obstacle apparaît. Cette réactivité explique le nombre important de valeurs quasi nulles ainsi que la disparition des amas observés.

La seconde opposition entre les deux premiers types de données et celui issu de l'environnement complexe s'exprime principalement par la quantité de valeurs supérieures à 500. Elle représente, pour tous capteurs confondus, 11,4% des valeurs issues de l'environnement simple, 12,4% des valeurs issues de l'environnement simple bruité et 48,8% des valeurs issues de l'environnement complexe. Cela signifie que l'espace des capteurs a été davantage exploité dans l'environnement complexe. En effet, la disposition des blocs dans l'environnement complexe n'offre pas toujours une échappatoire immédiate : par exemple lorsque le robot se trouve entre l'enceinte et le cylindre, il effectue plusieurs manœuvres avant de se retrouver dans une situation dans laquelle aucun obstacle n'est détecté. Cela explique aussi que l'apparition d'amas soit plus prononcée dans les données issues de l'environnement complexe. La différence de sensibilité des capteurs n'étant pas identique, le robot tourne préférentiellement d'un côté. De même, les moteurs ne possèdent pas une courbe de réponse motrice identique, ce qui implique que le robot tourne plus lentement d'un côté. L'expression de ces dissymétries se retrouve plus

particulièrement dans l'asymétrie entre les variables Var1, Var2 Var3 et les variables Var4, Var5, Var6 qui correspondent respectivement aux capteurs G90, G45, G10 et D10, D45, D90.

3.2.2. Étude des composantes principales des capteurs

Plusieurs manières de présenter l'ACP existent, elles seront abordées ici de façon à favoriser l'interprétation robotique des vecteurs propres résultant de l'analyse. Dans ce cas, l'analyse en composantes principales se comprend comme la recherche des axes ayant un pouvoir de description maximale tout en étant décorrélés entre eux. L'évaluation du pouvoir de description correspond au calcul du moment d'inertie. Le moment d'inertie d'un nuage de points par rapport au barycentre B se définit ainsi, avec la distance euclidienne d :

$$M_B = \frac{1}{n} \sum_{i=1}^n d^2(B, X_i)$$

$$M_B = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^p (x_{ij} - b_j)^2$$

Où n représente le nombre de points X_i de dimension p . Le moment d'inertie revient à la somme des variances empiriques de chaque variable. Il mesure la dispersion du nuage de points autour de son barycentre. L'inertie d'un nuage de points par rapport à un axe Δ passant par B s'écrit avec $h_{\Delta i}$ la projection orthogonale de X_i sur l'axe Δ :

$$M_{\Delta} = \frac{1}{n} \sum_{i=1}^n d_B^2(h_{\Delta i}, X_i)$$

Ce moment d'inertie M_{Δ} mesure la proximité à l'axe Δ du nuage de points. Le moment d'inertie par rapport à un axe est un cas particulier du moment d'inertie par rapport à un sous-espace vectoriel V à une dimension. En posant $h_{V i}$ comme la projection orthogonale de X_i sur le sous-espace vectoriel V , le moment d'inertie s'écrit :

$$M_V = \frac{1}{n} \sum_{i=1}^n d_B^2(h_{V i}, X_i)$$

La relation entre le moment d'inertie par rapport au barycentre et celui par rapport à un sous-espace provient du théorème de Pythagore. Celui-ci décompose les distances comme suit :

$$d^2(B, X_i) = d_B^2(h_{V i}, X_i) + d_B^2(h_{V^c i}, X_i)$$

$$d^2(B, X_i) = d^2(h_{V i}, B) + d^2(h_{V^c i}, B)$$

Où V^c est le complémentaire orthogonal de V dans \mathfrak{R}^p et $h_{V^c i}$ la projection orthogonale de X_i sur V^c . Cette reformulation permet de retrouver le théorème de Huygens :

$$M_B = M_V + M_{V^c}$$

Dans le cas particulier où le sous-espace se réduit à une seule dimension, le moment d'inertie M_{V^c} correspond à la mesure de l'allongement du nuage de points selon cet axe. Autrement dit, cette mesure, appelée aussi inertie expliquée, évalue le pouvoir de description d'un axe. De plus, en décomposant l'espace \mathfrak{R}^p comme la somme de sous-espaces à une dimension et orthogonaux entre eux, l'expression du moment d'inertie devient :

$$M_B = M_{V_1^c} + M_{V_2^c} + \dots + M_{V_p^c}$$

Dans l'ACP, les inerties expliquées portées par les axes correspondent aux valeurs propres de la matrice de corrélation. Le moment d'inertie s'écrit alors :

$$M_B = \lambda_1 + \lambda_2 + \dots + \lambda_p$$

La contribution absolue de l'axe Δ_k à l'inertie totale du nuage de points est égale à λ_k et le pourcentage d'inertie expliquée par Δ_k s'exprime par sa contribution relative :

$$Cr_k = \frac{\lambda_k}{\lambda_1 + \lambda_2 + \dots + \lambda_p}$$

Les graphiques de la Figure IV-19 montrent les valeurs du cumule du pourcentage de l'inertie expliquée et les valeurs propres de la matrice de corrélation selon les trois types d'environnements à l'origine des données. La décroissance du pouvoir de description permet d'utiliser l'analyse en composantes principales comme un moyen de compression de l'information en conservant uniquement les premières composantes pour décrire les données. Ce point est particulièrement intéressant pour alléger le temps de calcul dans le cadre de l'évaluation des règles. Par ailleurs, ce sont généralement les dernières composantes qui représentent le bruit.

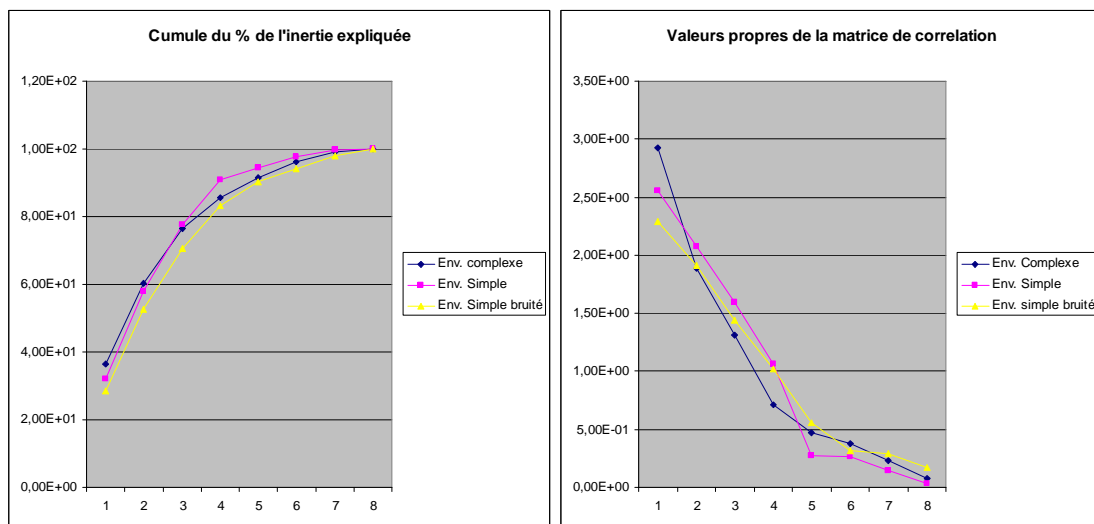


Figure IV-19 : Graphique du cumul de l'inertie expliquée et des valeurs propres des composantes principales.

La détermination du nombre d'axes à retenir s'appuie sur la règle empirique de Kaiser et sur l'interprétation des vecteurs propres. Le critère de Kaiser considère comme significatives toutes les valeurs propres supérieures à 1. En effet, la somme des valeurs

propres demeure toujours égale au nombre de dimensions de l'espace des données dans toute ACP normalisée, par conséquent, la moyenne des valeurs propres valant 1 peut servir de seuil.

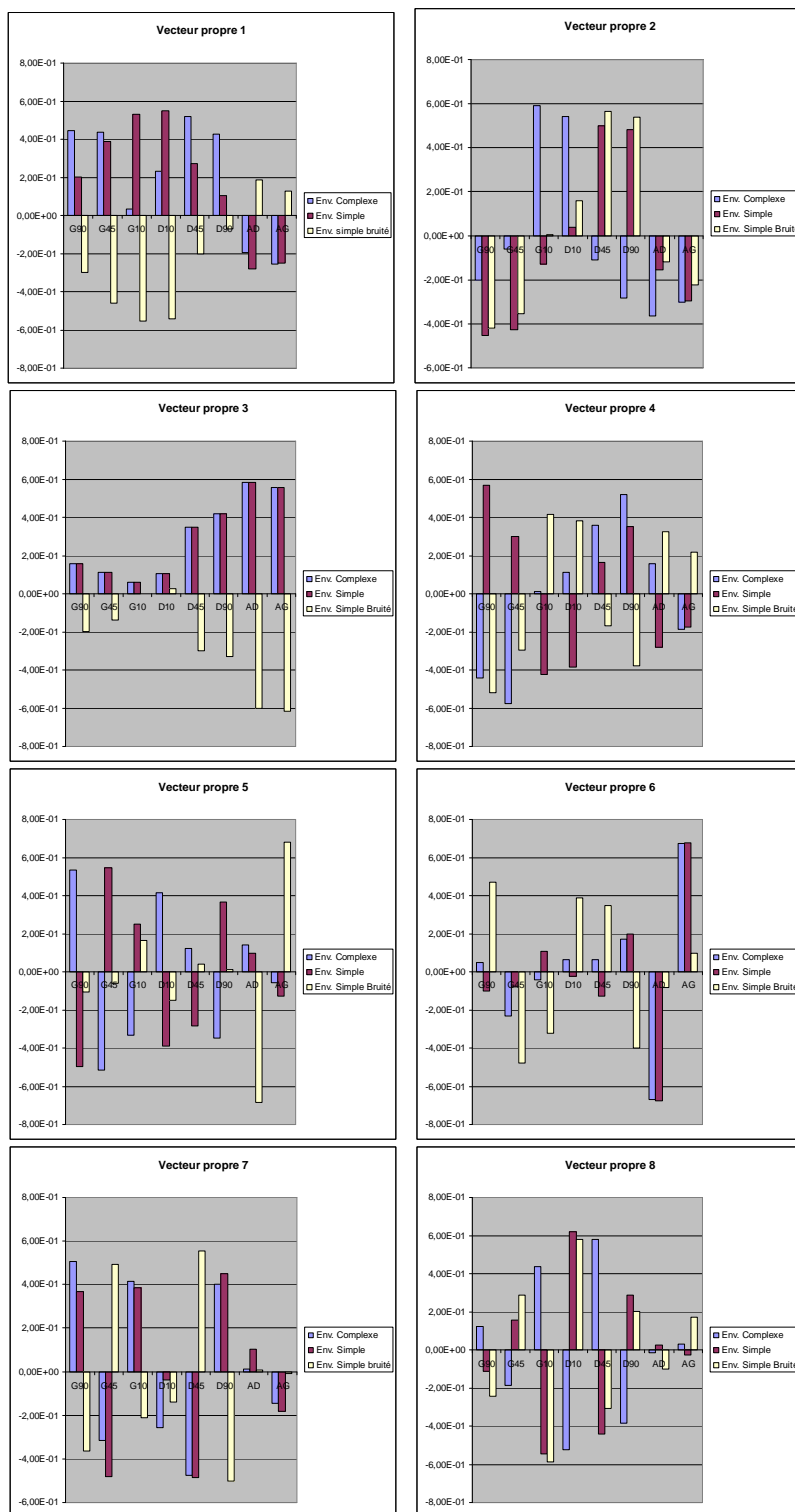


Figure IV-20 : Les huit vecteurs propres issus de l'ACP selon les trois types d'environnements.

Les courbes se ressemblent selon l'origine des données, particulièrement entre les courbes ayant pour origine l'environnement simple et l'environnement simple bruité. Selon ce critère, seules les trois premières composantes suffisent à la description des données pour l'environnement complexe ; en revanche, quatre composantes semblent nécessaires pour les deux autres types d'environnements. Par ailleurs, au regard du cumul du pourcentage de l'inertie expliquée, ces trois composantes portent moins de 80% de l'information quel que soit l'environnement. La prise en compte de la quatrième composante permet de couvrir plus de 85% de l'information. Toutefois, le critère de Kaiser reste seulement indicatif et le choix définitif du nombre d'axes repose sur l'interprétation des vecteurs propres (Figure IV-20).

L'interprétation transversale des vecteurs proposés repose sur l'analyse des valeurs des vecteurs propres dont chaque composante se trouve associée à un capteur. Plus la valeur d'une composante est élevée, plus l'état du capteur se révèle important. Indirectement, chaque vecteur illustre un type de configuration sensorielle :

- Le *vecteur propre n°1* discrimine l'avant et l'arrière pour les trois types d'environnement, bien que cela soit moins net pour l'environnement « complexe ».
- Le *vecteur propre n°2*, pour l'environnement complexe, discrimine les deux capteurs avant des autres capteurs. Pour les environnements « simple » et « simple bruité », ce vecteur discrimine les trois capteurs situés sur un même côté à l'avant. À noter que cette dissymétrie peut s'expliquer par le fait que les deux environnements n'offraient pas la possibilité au Khepera de changer facilement de sens de rotation dans l'enceinte.
- Le *vecteur propre n°3*, dans les trois cas, discrimine la présence ou non d'un objet quel que soit sa position avec une préférence toutefois pour les capteurs situés à l'arrière et sur le côté.
- Le *vecteur propre n°4*, pour l'environnement complexe, différencie la gauche de la droite. En effet, l'environnement faisait que le Khepera devait tourner à gauche, parfois à droite pour éviter au mieux les obstacles. Comparativement, l'environnement « simple » offrait beaucoup moins de situations symétriques puisque le khepera avait tendance à toujours éviter les obstacles de la même manière, autrement dit à tourner dans le même sens. Cela explique que pour les deux environnements « simple » et « simple bruité », ce vecteur distingue l'avant de l'arrière et les informations venant des côtés.
- Le *vecteur propre n°5* : Pour l'environnement « simple » et « complexe », il semblerait que ce vecteur propre discrimine une situation qui ressemble à l'entrée d'un couloir ou d'un coin. Pour l'environnement « simple » et « simple bruité », il y a discrimination principalement entre les deux capteurs situés à l'arrière.
- Le *vecteur propre n°6*, pour l'environnement simple et l'environnement complexe, discrimine principalement les deux capteurs arrière. Plus précisément, pour l'environnement simple et l'environnement bruité, il semblerait que ce vecteur propre discrimine également une situation qui ressemble à l'entrée d'un couloir ou, plus certainement compte, tenu de la topologie de l'enceinte, d'un coin.
- Le *vecteur propre n°7* discrimine la situation lorsque le robot Khepera a détecté un obstacle sur ses deux capteurs avant et sur les côtés.

- Le vecteur propre n°8 ressemble au vecteur propre 7 mais les capteurs avant s'opposent.

Pour les vecteurs 7 et 8, il est difficile de trouver une sémantique claire. Les vecteurs 5, 6, 7 et 8 semblent davantage liés aux motifs topologiques particuliers que le Khepera a pu rencontrer. En revanche, les quatre premiers vecteurs semblent davantage liés, malgré l'incidence de l'environnement, aux caractéristiques physiques (avant/arrière, gauche/droite). Cette remarque convient plus particulièrement aux quatre premiers vecteurs propres issus de l'environnement complexe. Cette « abstraction » ou plutôt cette généralisation s'explique par le fait que l'environnement complexe offre des situations plus variées et moins répétitives.

Cette interprétation conduit à retenir uniquement les quatre premières composantes résultant de l'ACP et plus précisément, celles issues de l'environnement complexe. Ces quatre vecteurs propres constituent le soubassement des primitives sensorielles. Une autre manière de s'assurer qualitativement de ce choix se trouve dans la visualisation des données projetées dans leur espace calculé par l'ACP avec Figure IV-21, Figure IV-22 et Figure IV-23.

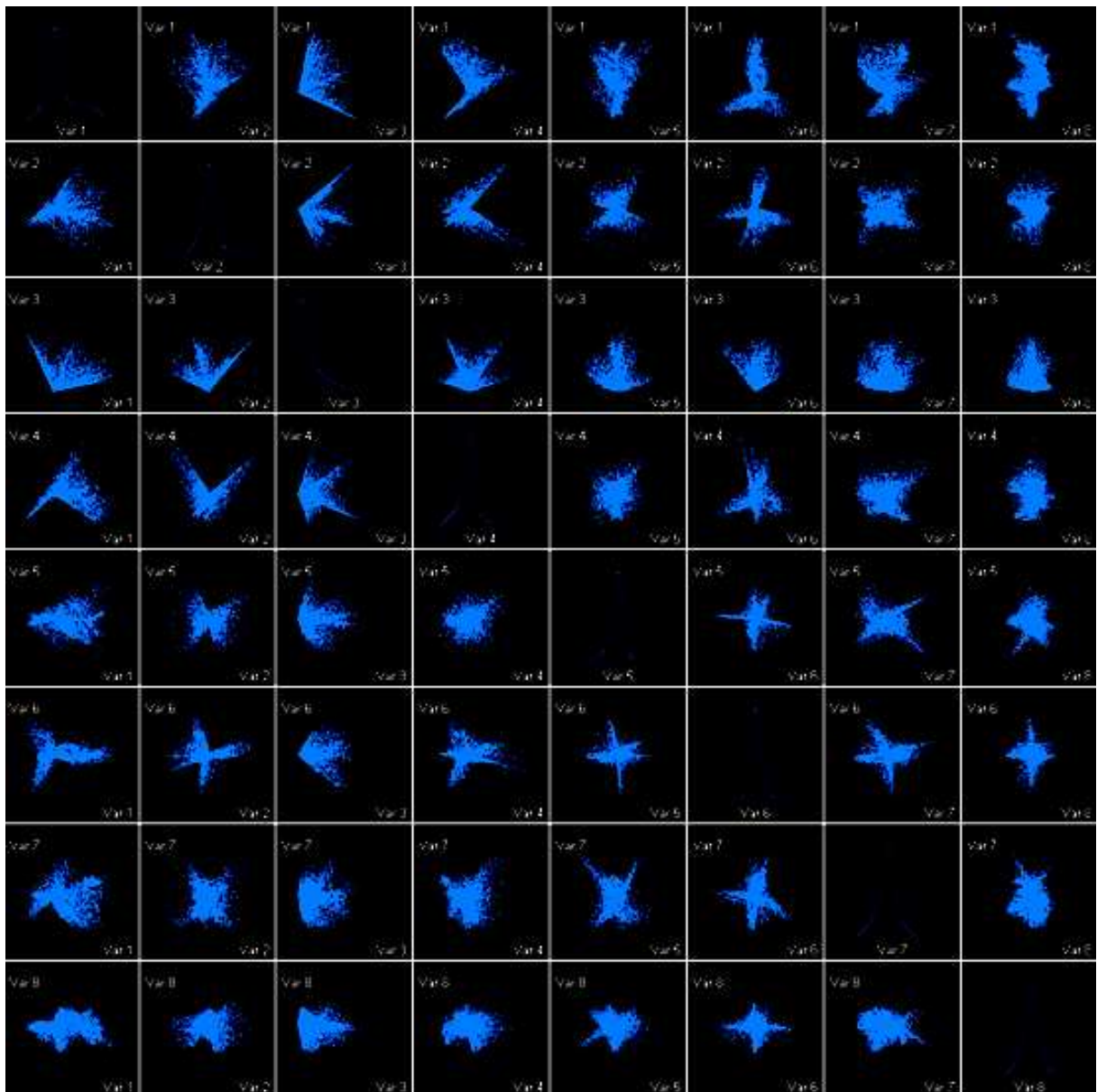


Figure IV-21 : Projection des données issues de l'environnement simple dans l'espace calculé par une ACP.

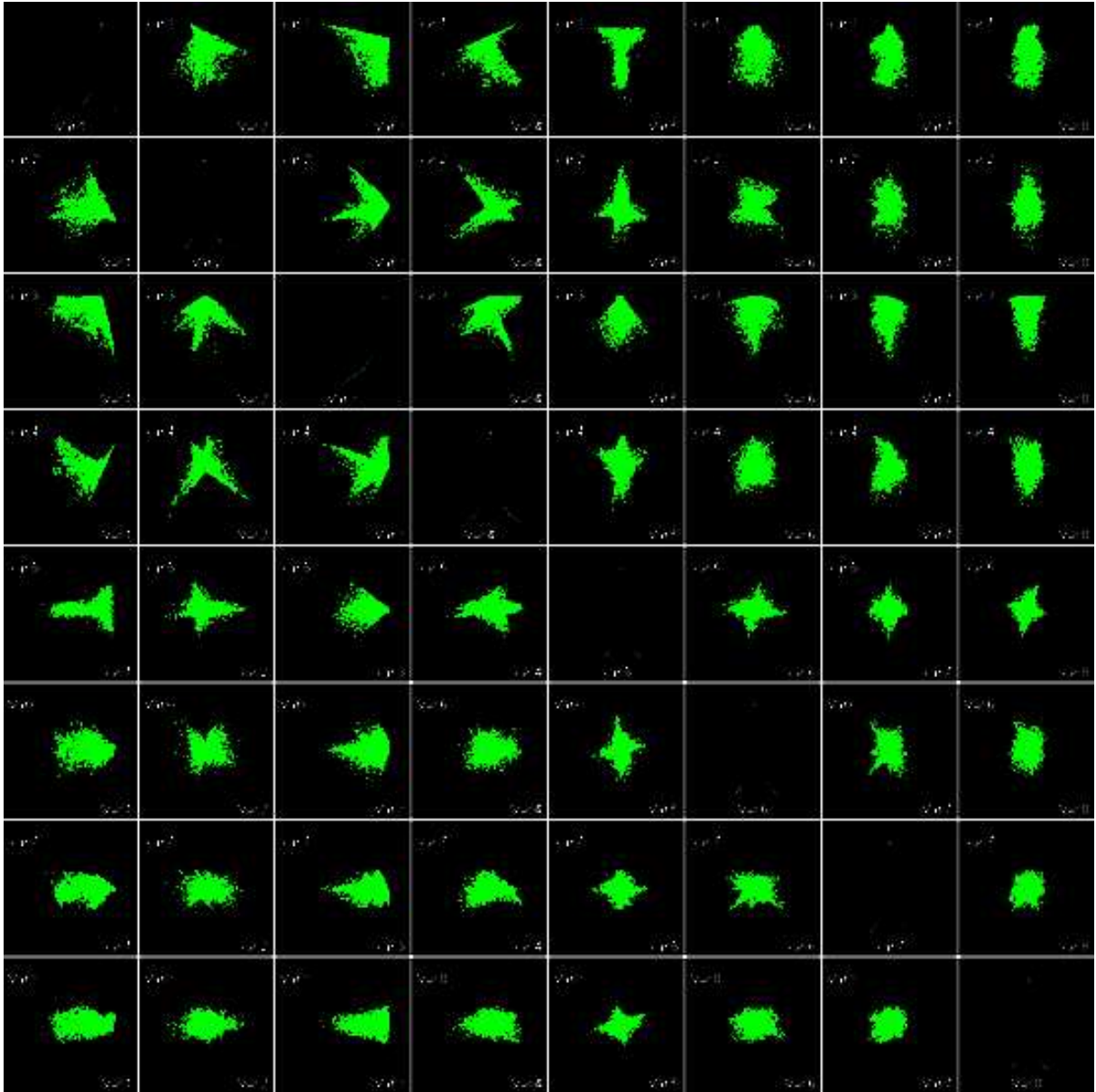


Figure IV-22 : Projection des données issues de l'environnement simple bruité dans l'espace calculé par une ACP.

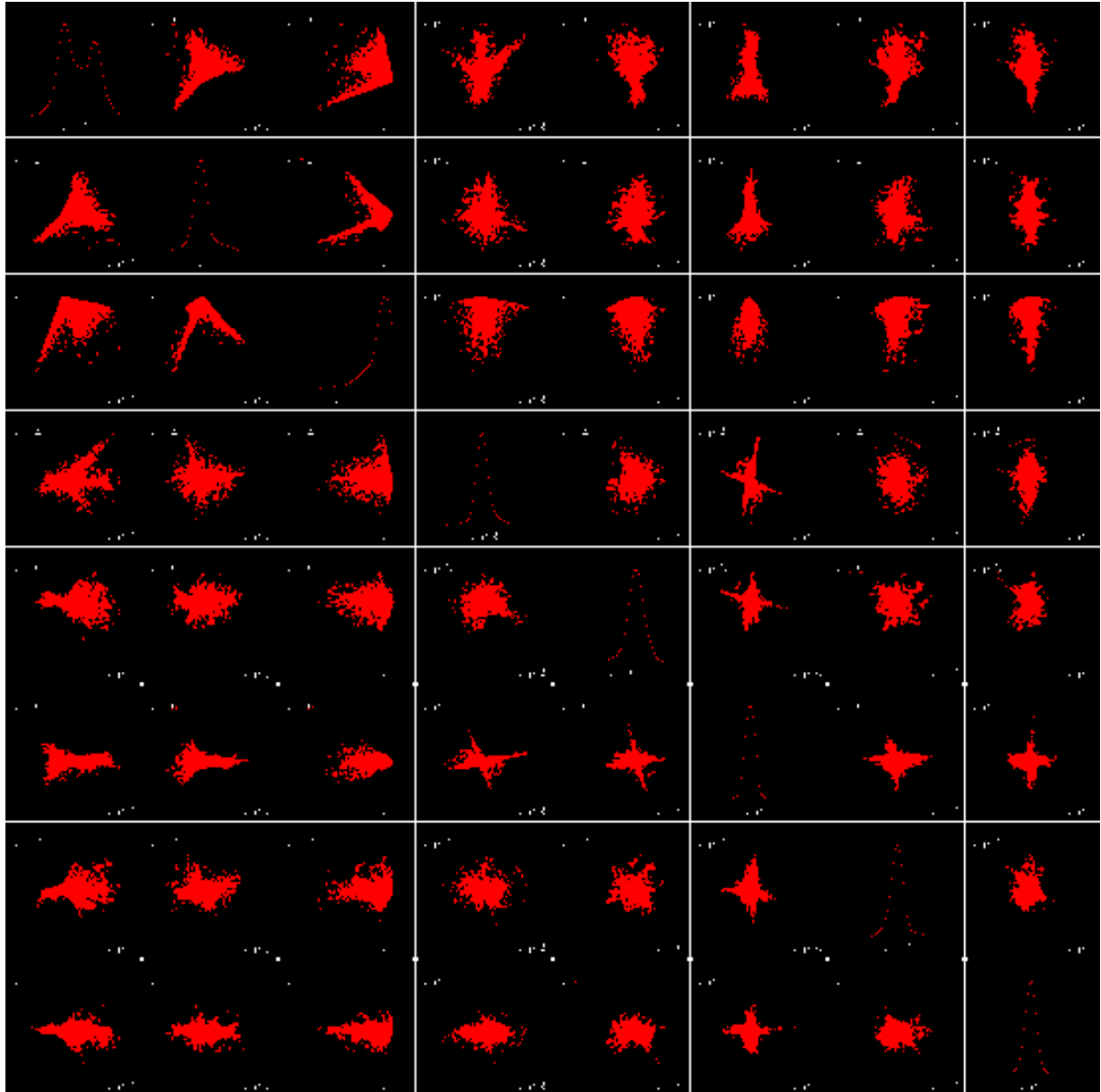


Figure IV-23 : Projection des données issues de l'environnement complexe dans l'espace calculé par une ACP.

3.2.3. Construction de la primitive sensorielle

Les messages sensoriels résultent d'un traitement effectué sur les valeurs issues des capteurs qui constituent la primitive sensorielle. Dans le cadre expérimental proposé, cette primitive sensorielle possède l'avantage de diminuer le bruit et de réduire la taille du vecteur d'information sensorielle. En effet, les messages sensoriels représentent les quatre premières composantes du vecteur sensoriel (les valeurs issues des 8 capteurs) projeté dans l'espace défini par l'ACP, effectuée sur les données issues de l'exploration du robot avec un comportement d'évitement d'obstacles dans l'environnement complexe.

Ainsi, les messages sensoriels se constituent de 5 valeurs, la première correspond à l'étiquette temporelle et les 4 autres sont dédiées à la description des capteurs. Toutes les

valeurs des messages sont comprises entre 0 et 1. Or, la projection des données résultant du produit entre la matrice des vecteurs propres et les vecteurs composés par la valeur des 8 capteurs produit des valeurs supérieures à 1. Une normalisation devient alors inévitable. Par ailleurs, les valeurs des messages sensoriels se trouvent associées à une valeur d'incertitude. L'initialisation de cette incertitude dépend de la variance empirique du bruit et du calcul de normalisation.

La détermination du calcul de normalisation et de l'incertitude initiale va s'appuyer sur l'acquisition de 600 mesures pour chaque capteur, avec le robot Khepera immobile au sein d'une enceinte rectangulaire suffisamment petite pour stimuler tous les capteurs. Les valeurs des huit capteurs forment le vecteur sensoriel. Il est à noter que la forme de l'enceinte explique la non-homogénéité des valeurs. Le Tableau IV-2 affiche la moyenne des valeurs brutes des capteurs et leur variance.

	$C_1 = G90$	$C_2 = G45$	$C_3 = G10$	$C_4 = D10$	$C_5 = D45$	$C_6 = D45$	$C_7 = AD$	$C_8 = AG$
moyenne	497	101	895	1023	190	284	186	480
variance	239	5	212	1	51	60	61	102

Tableau IV-2 : Moyenne et variance des composantes des 600 vecteurs sensoriels, soit les valeurs des capteurs obtenues avec le robot immobile dans une petite enceinte rectangulaire.

Dans le cadre classique de l'utilisation d'une ACP, les valeurs des capteurs sont blanchies sur la base des données ayant servi à l'ACP, puis projetées dans le sous-espace calculé. La projection s'écrit avec C_p , la valeur du capteur p correspondant au numéro d'ordre de la liste employée jusqu'à présent, et avec C_p^* correspondant à une composante du résultat de la projection, la matrice v représentant les vecteurs propres, Moy_p représentant la moyenne des valeurs du capteur p et Var_p en représentant la variance :

$$\begin{pmatrix} C_1^* \\ C_p^* \end{pmatrix} = \begin{pmatrix} v_{11} & v_{1p} \\ v_{p1} & v_{pp} \end{pmatrix} * \begin{pmatrix} (C_1 - Moy_1)/Var_1 \\ (C_p - Moy_p)/Var_p \end{pmatrix}$$

	C_1^*	C_2^*	C_3^*	C_4^*	C_5^*	C_6^*	C_7^*	C_8^*
moyenne	0,391	4,69	-9,81	0,775	1,130	3,41	3,80	-1,83
variance	0,0007	0,00382	0,0236	0,0022	0,0072	0,0199	0,0189	0,0613

Tableau IV-3 : Moyenne et Variance des composantes des 600 vecteurs sensoriels blanchies puis projetées dans l'espace défini par l'ACP des données issues de l'environnement complexe.

Mais, les valeurs moyennes des vecteurs blanchis et projetés affichées dans le Tableau IV-3 se trouvent en dehors des bornes fixées. Une première solution consiste à changer de repère et d'échelle, c'est-à-dire à diviser par la valeur maximale que peut atteindre la première composante puis effectuer un décalage de l'origine. En effet, la première composante possède par définition la plus grande amplitude de valeurs, ainsi la transformation assure le confinement des valeurs entre 0 et 1 pour toutes les composantes. Cette transformation ne touche pas l'intégrité de l'ACP puisque cette transformation est linéaire. En notant, Max_1 deux fois le maximum de la composante C_1^* en valeur absolue de

la première composante du résultat de la projection sans normalisation et d_1 le décalage nécessaire à recentrer les valeurs 0 et 1 de la première composante, la transformation s'écrit :

$$\begin{vmatrix} C_1^* \\ C_p^* \end{vmatrix} = \begin{vmatrix} v_{11} & v_{1p} \\ v_{p1} & v_{pp} \end{vmatrix} * \begin{vmatrix} C_1 \\ C_p \end{vmatrix} * \frac{1}{Max_1} + \begin{vmatrix} d_1 \\ d_p \end{vmatrix}$$

Le calcul du maximum se trouve autorisé par le fait que les valeurs des capteurs sont bornées entre 0 et 1023.

	C_1^*	C_2^*	C_3^*	C_4^*	C_5^*	C_6^*	C_7^*	C_8^*
moyenne	0,614	0,627	0,268	0,500	0,551	0,556	0,555	0,482
variance	4,70E-6	2,49E-6	8,30E-6	2,85E-6	2,74E-6	1,56E-6	3,91E-6	2,17E-6

Tableau IV-4 : Moyenne et Variance des composantes des 600 vecteurs sensoriels normalisés en fonction de la première composante et projetés dans l'espace défini par l'ACP des données issues de l'environnement complexe.

Toutefois, certaines composantes restent confinées dans des bornes étroites comprises entre 0 et 0,5, comme pour la troisième composante par exemple (Tableau IV-4). Afin de pallier ce problème, la transformation précédente doit être adaptée à chaque axe. Ainsi, chaque composante exploite au mieux les bornes, ce qui en facilite la lecture (Tableau IV-5). L'équation reste similaire à la précédente :

$$\begin{vmatrix} C_1^* \\ C_p^* \end{vmatrix} = \begin{vmatrix} v_{11}/max_1 & v_{1p}/max_1 \\ v_{p1}/max_p & v_{pp}/max_p \end{vmatrix} * \begin{vmatrix} C_1 \\ C_p \end{vmatrix} + \begin{vmatrix} d_1 \\ d_p \end{vmatrix}$$

Toutefois, il faut préciser que le calcul du maximum diffère de la précédente, afin d'exploiter au mieux l'intervalle autorisé. Le max_p ne correspond plus au maximum en valeur absolue fois deux mais à l'intervalle entre le minimum et le maximum de la composante C_p^* . Cette nouvelle transformation présente l'inconvénient d'interdire toute comparaison entre les axes, mais cela ne constitue pas un problème pour le fonctionnement du système puisque chaque composante est traitée indépendamment des autres.

	C_1^*	C_2^*	C_3^*	C_4^*	C_5^*	C_6^*	C_7^*	C_8^*
moyenne	0,405	0,802	0,484	0,510	0,612	0,623	0,582	0,447
variance	1,88E-5	1,07E-5	4,11E-5	1,33E-5	1,16E-5	1,04E-5	1,59E-5	1,08E-5

Tableau IV-5 : Moyenne et Variance des composantes des 600 vecteurs sensoriels normalisés en fonction de chaque composante et projetés dans l'espace défini par l'ACP des données issues de l'environnement complexe.

La différence entre les deux types de normalisation est particulièrement visible sur un jeu de données plus riche comme celui recueilli lors de l'expérience avec l'environnement complexe. Le changement de procédé pour déterminer le max_p offre effectivement une meilleure exploitation de l'espace (Figure IV-24 et Figure IV-25).

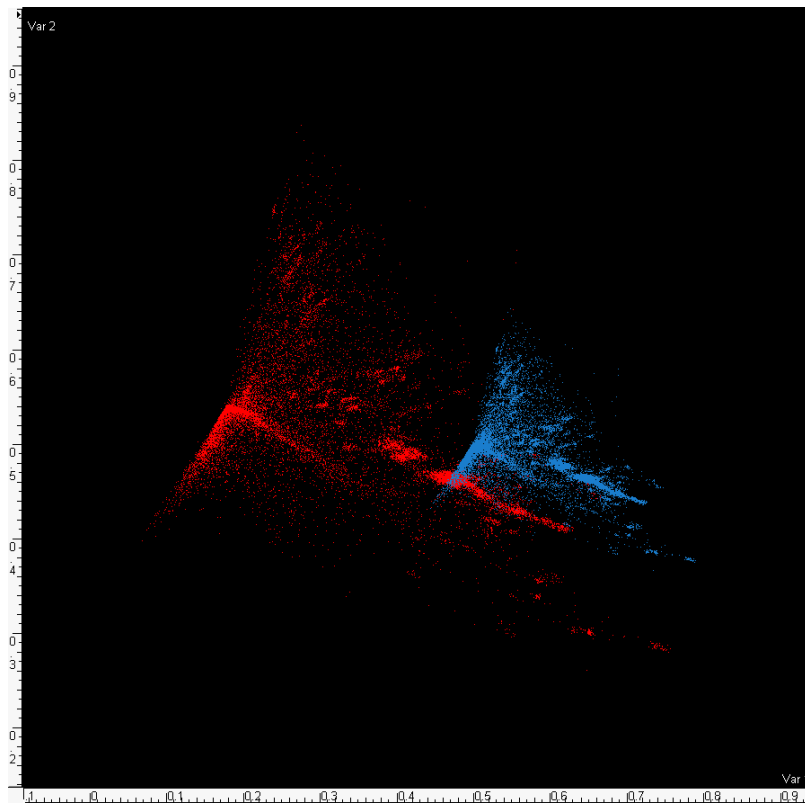


Figure IV-24 : Visualisation des deux premières composantes avec des données identiques mais une normalisation différente: les points bleus correspondent à la normalisation simple, les points rouges à la normalisation adaptée à chaque composante.

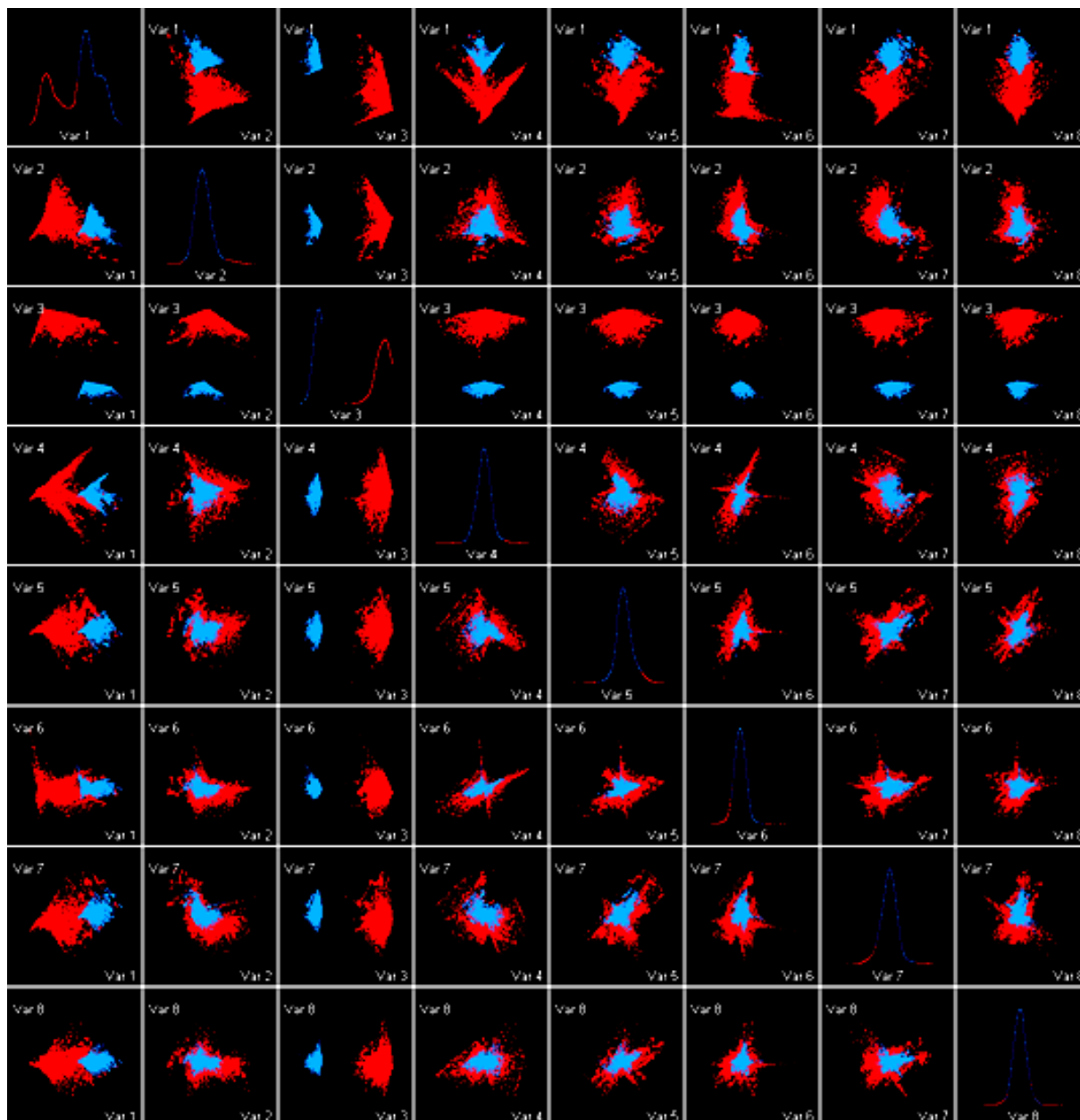


Figure IV-25 : Visualisation des données projetées mais avec une normalisation différente: les points bleus correspondent à la normalisation simple, les points rouges à la normalisation adaptée à chaque composante.

La normalisation adaptée à chaque composante offre une meilleure lisibilité et permet d'initialiser l'incertitude à la même valeur pour toutes les composantes. Tous ces avantages incitent à intégrer la normalisation adaptée dans la primitive sensorielle. Par ailleurs, les variances affichées dans le Tableau IV-5 conduit à choisir 0,001 pour la valeur d'initialisation des variances. Cette incertitude dispose d'un facteur 100 par rapport à celle observée mais elle assure la possibilité d'élire une règle dans n'importe quelle situation au démarrage du système. L'ajustement de cette incertitude se fera ensuite lors de la compétition entre les règles. La variance de l'étiquette temporelle contenue dans le message sera estimée dans un premier temps à 0,0001. Cette différence avec l'incertitude des valeurs sensorielles devrait permettre d'équilibrer l'influence de l'espace sensoriel et le temps dans

le score des règles, bien que les expériences en environnement imposé révéleront que ces chiffres nécessitent une réévaluation.

3.2.4. Construction des environnements imposés

Les environnements imposés s'inscrivent également dans le cadre de la caractérisation de l'environnement puisque leur construction s'appuie sur les observations précédentes qui serviront de matériel de base pour la première étude expérimentale du système. Un environnement imposé correspond à une présentation cyclique d'un certain nombre de stimulations sensorielles définies et ordonnées indépendamment des conclusions motrices et il diffère en cela d'un environnement virtuel qui simule le monde. Ce procédé se trouve motivé par le caractère chaotique de l'architecture cognitive qui rend le système très sensible aux conditions initiales et aux bruits. Ce caractère ne signifie pas que le système soit instable mais que l'histoire du robot se trouve toujours différente et que la base de règles ne se stabilise jamais exactement de la même manière. Cette variabilité induit un manque de rigueur pour l'étude de l'influence de certains paramètres tels que les taux de taxe ou d'enchère. Ainsi, dans ces expériences, l'action n'est pas prise en compte et seules la dynamique d'élection et la capacité d'adaptation des prémisses ont été observées.

L'environnement imposé assure que la série temporelle des données sensorielles soit identique entre deux expériences et autorise ainsi leur comparaison. L'environnement imposé se construit à partir de l'enregistrement d'une courte séquence lors du déplacement du robot. Plus précisément, deux séquences ont été enregistrées, chacune durant 5 secondes soit 50 états sensoriels. Une combinaison de séquences forme un scénario. Afin d'étudier l'influence de l'introduction d'états sensoriels nouveaux sur la dynamique de sélection des règles, trois scénarii ont été conçus :

1. Le scénario simple (S1) répète une même séquence A (AAAAA...).
2. Le scénario complexe (S2) répète un motif composé de deux séquences A suivies d'une séquence B (ABAABA...).
3. Le scénario complexe (S3) répète un motif composé de deux séquences B suivies d'une séquence A (BBABBAB...).

La séquence A représente la trace sensorielle du robot Khepera lorsqu'il arrive sur un mur possédant un angle d'environ 45°. La séquence B représente la trace sensorielle du robot Khepera lors de son démarrage dans un coin : un mur à proximité des capteurs arrière et un autre situé à proximité des capteurs droits. La Figure IV-26 assure que ces deux séquences possèdent un ensemble d'états sensoriels commun mais aussi des états sensoriels très éloignés. La diagonale du tableau des graphiques représente l'histogramme de l'ensemble des points des deux séquences confondues.

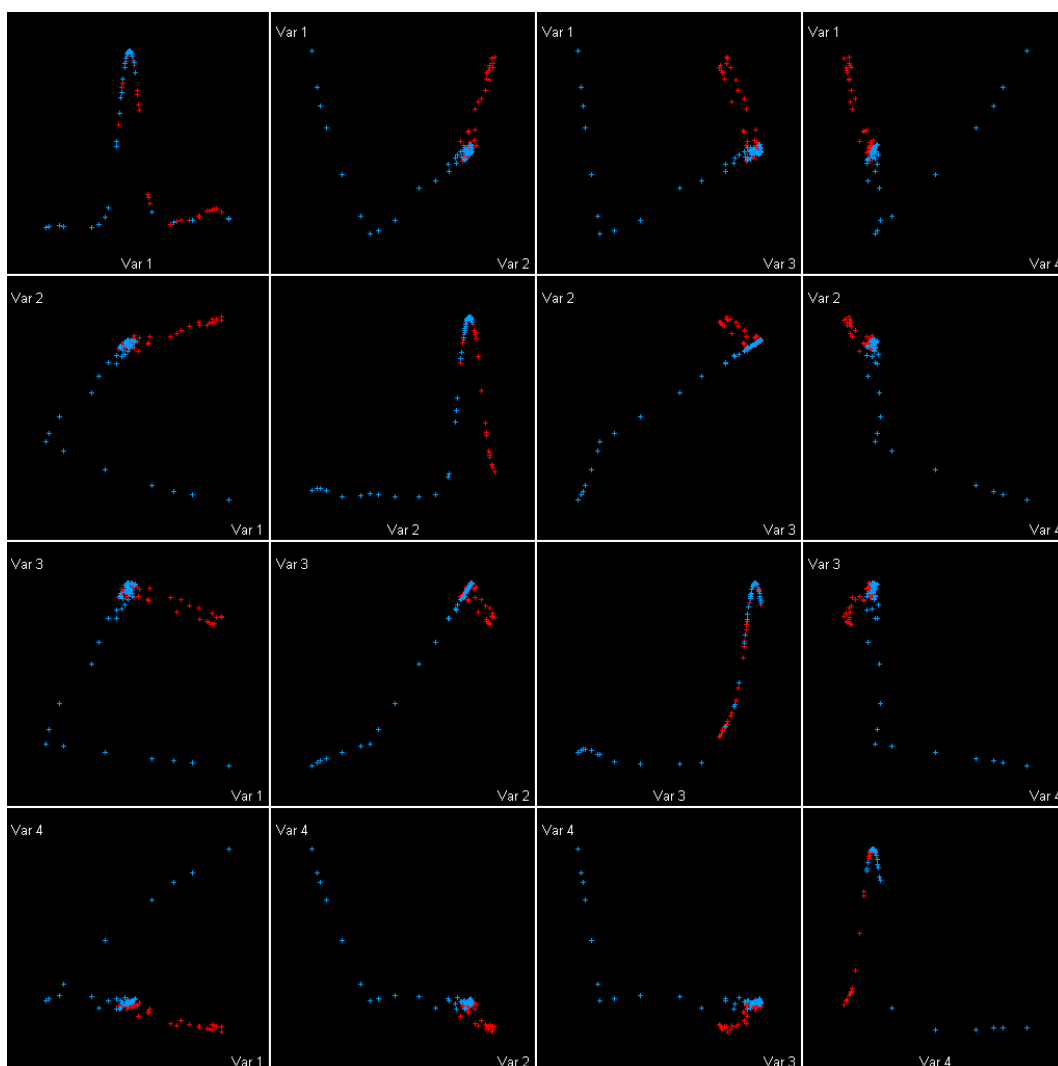


Figure IV-26 : Représentations des messages sensoriels induits par les deux séquences. En rouge, les points provenant de la séquence A. En bleu, les points provenant de la séquence B.

Quantitativement, lors des expériences durant 1800 s, la séquence A est jouée 360 fois avec le scénario S1, 240 fois avec le scénario S2 et 120 avec le scénario S3. De la même manière, la séquence B est jouée 120 fois avec le scénario S2 et 240 fois avec le scénario S3.

3.3. Analyse des prémisses des règles initiales

3.3.1. Les règles initiales pour l'environnement réel

La définition de la primitive sensorielle autorise désormais une caractérisation des règles initiales. Les 600 règles initiales proviennent de l'enregistrement d'une séquence de 30 s échantillonnée à 20 Hz des couples sensorimoteurs pendant que le robot évitait les obstacles dans un environnement complexe grâce à un algorithme de type Braitenberg. La consigne de vitesse pour chacun des deux moteurs représente la conclusion motrice d'une règle sensorimotrice et les quatre premières composantes de la projection normalisée de l'état des capteurs représentent sa prémisse sensorielle. La projection normalisée s'effectue sur la base d'une ACP appliquée aux données issues de l'enregistrement d'un essai de

30 min échantillonné à 10 Hz lors de l'évitement d'obstacle dans l'environnement complexe.

Le choix de prendre pour l'extraction de règles un enregistrement différent de celui utilisé pour l'ACP a été motivé par le fait que certaines valeurs ne sont jamais atteintes comme le montrent les histogrammes. En effet, pour les enregistrements de 30 min, au démarrage la position centrale ne permettait à aucun capteur de détecter la présence d'un obstacle, ensuite la sensibilité de la détection rendait impossible certains cas de figure. Or, un large éventail de prémisses et d'actions possibles est nécessaire afin de vérifier à la fois la faculté d'adaptation du système et sa robustesse, autrement dit les conséquences comportementales de cette adaptation. Pour l'enregistrement en cours servant à l'extraction des règles, la position initiale du robot était dos à un obstacle, soit une situation insolite dans la mesure où le robot ne pouvait pas se retrouver dans une telle situation dans le déroulement normal avec une position initiale neutre. Par ailleurs, l'élargissement du domaine des prémisses des règles sensorimotrices n'implique pas une remise en cause des données utilisées par l'ACP puisque pour celles effectuées sur les données issues de l'environnement complexe, les quatre composantes se révèlent davantage liées à la topologie des capteurs qu'à une configuration spatiale précise de l'environnement.

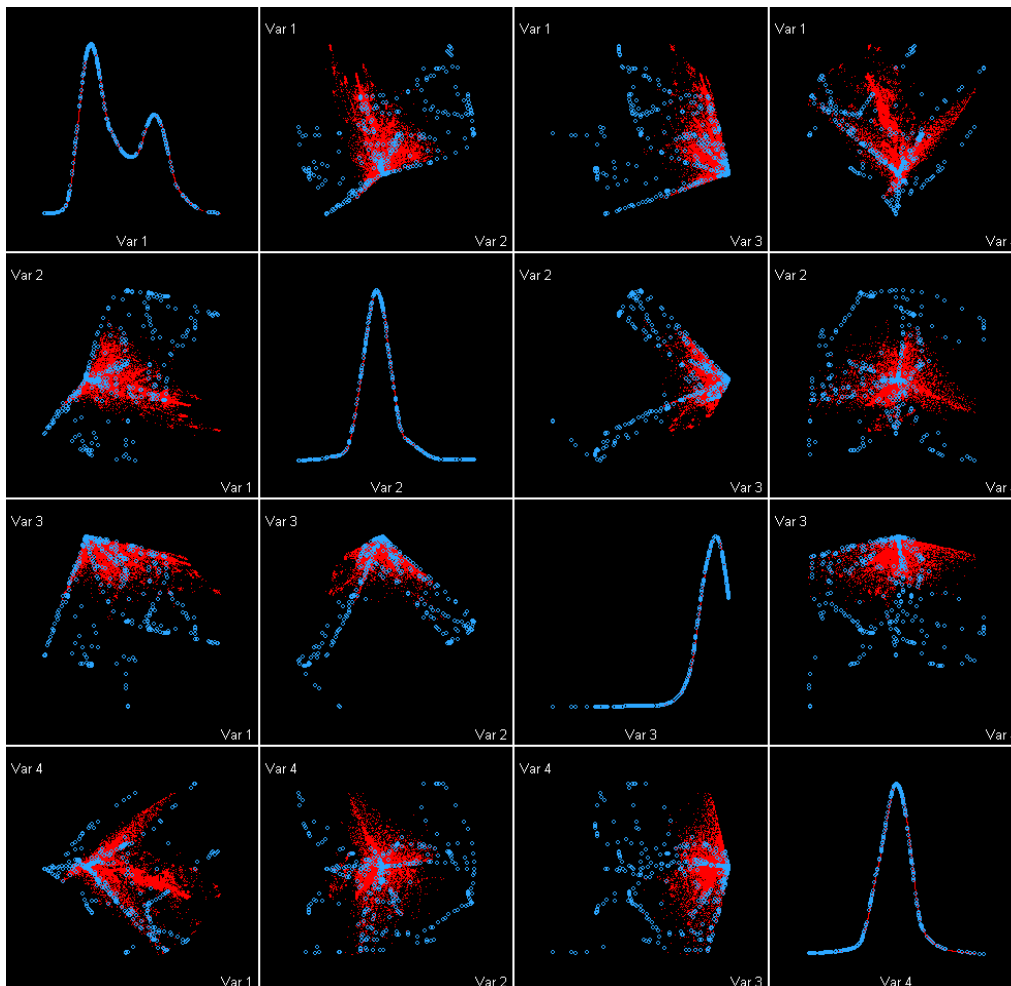


Figure IV-27 Représentation de la projection normalisée des quatre premières composantes des 18000 données issues de l'environnement complexe (points rouges) et des prémisses des 600 règles initiales (cercles bleus).

La Figure IV-27 compare les données issues des types d'enregistrements avec l'environnement complexe, celui de 30 s et celui de 30 min. La saturation de certains capteurs se traduit par les dessins de droites. Ces droites illustrent la projection des arêtes de l'hypercube formé par les extremums des capteurs. L'analogie peut être faite avec l'apparition de losanges dans le dessin de la projection d'un cube sur un plan qui ne soit ni parallèle ni perpendiculaire à une face.

3.3.2. Les règles initiales pour l'environnement imposé

L'environnement imposé étant restreint, les règles initiales ont été choisies afin qu'elles ne couvrent pas directement les états sensoriels de l'environnement imposé et elles ont donc été extraites de l'enregistrement dans un environnement simple pour les expériences en environnement imposé. Cette situation permettra d'évaluer les capacités d'adaptation attendues. Par ailleurs, les expériences sur environnement imposé offre la possibilité d'évaluer l'influence du nombre de règles dans la stabilisation du système. Quatre ensembles de règles ont été constitués sur la base des règles initiales décrites précédemment (Tableau IV-6).

	E1-bis	E1	E2	E3
Nombre de règles de gestion	200	-	-	-
Nombre de règles sensorimotrices	200	200	400	600

Tableau IV-6 : Nombre de règles contenues dans les quatre ensembles de règles utilisés pour les expériences en milieu imposé.

L'ensemble E1-bis introduit les règles de gestion. Celles-ci contrôlent la durée de l'application des ordres moteurs selon le mécanisme expliqué dans le chapitre précédent. Toutes les prémisses inhibitrices des règles (sensorimotrices ou de gestion) possèdent alors un message de gestion. Lorsque celui-ci se trouve présent dans la mémoire événementielle, aucune règle ne peut être élue. Chaque règle de gestion se trouve associée à une règle sensorimotrice et cette association se traduit par des prémisses sensorielles de valeurs identiques. La conclusion d'une règle de gestion correspond à un message interne de gestion. Les prémisses des deux ensembles de règles étant égales, les règles associées se trouvent élues conjointement. Pour l'ensemble E1-bis, l'étiquette temporelle du message de gestion dans la prémisse inhibitrice vaut 1 et l'incertitude sur la valeur du message qui représente l'identifiant de la règle sensorimotrice associée est au minimum.

En débutant avec un seul message sensoriel dans la mémoire événementielle, un cycle de déclenchement se déroule alors comme suit :

1. Élection d'une règle sensorimotrice et de la règle de gestion associée.
2. Transmission du message moteur et du message de gestion à la mémoire événementielle à laquelle s'ajoute un nouveau message sensoriel.
3. Inhibition des règles empêchant une nouvelle élection
4. Introduction d'un nouveau message sensoriel, et mise à jour des étiquettes temporelles pour les autres messages dont le message de gestion autorise ainsi une nouvelle élection.

En somme, l'élection d'une règle sensorimotrice ne peut s'effectuer qu'une fois sur deux, ce qui revient à une fréquence de réactualisation des ordres moteurs de 20 Hz.

L'ensemble E1 représente un sous-ensemble localisé de l'ensemble E3 (Figure IV-28) car il correspond au début de l'enregistrement.

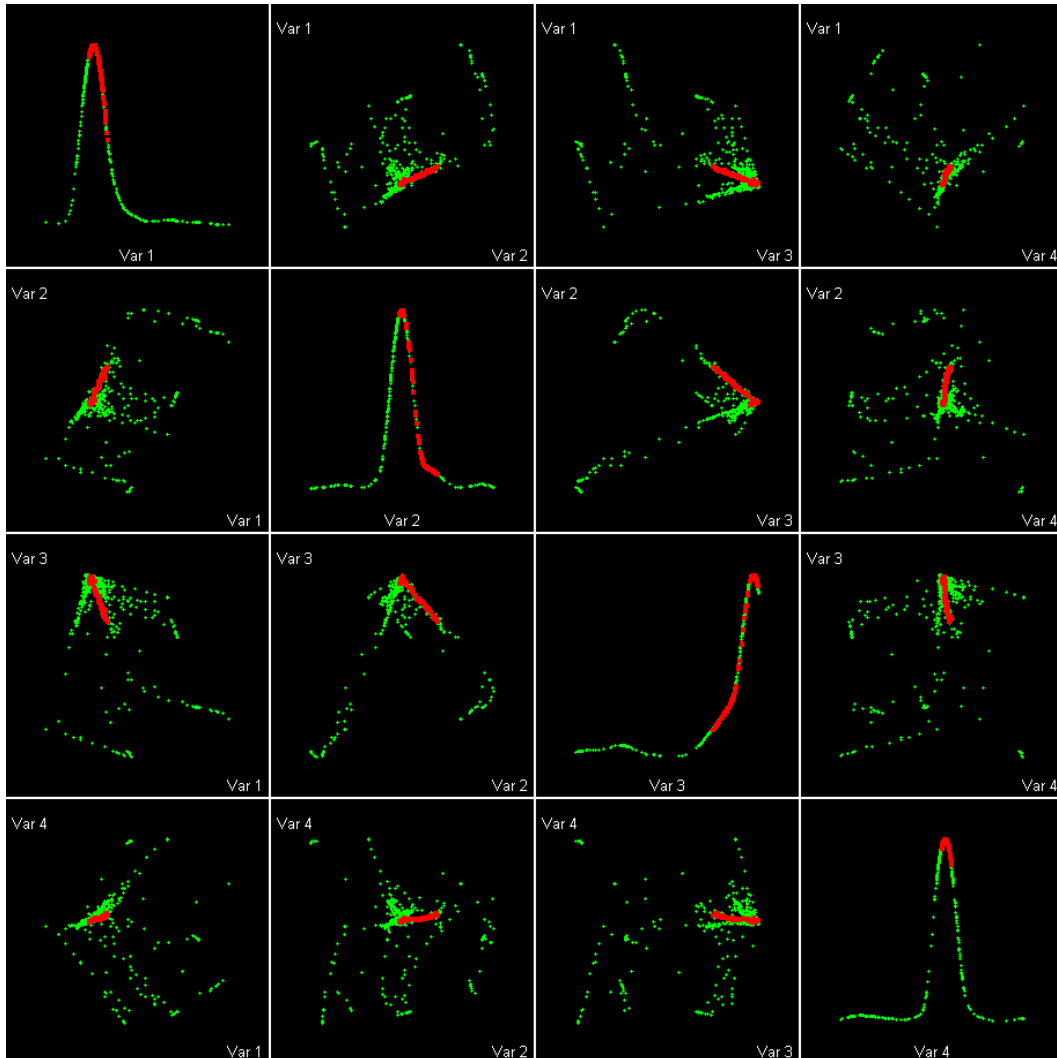


Figure IV-28 Comparaison des prémisses de règles entre E1 (carrés rouges) et E3 (croix vertes).

La comparaison entre les règles initiales et les états sensoriels produits par le scénario S1 illustre les deux problèmes auxquels le système est confronté (Figure IV-29) et auxquels il répondra :

1. Comment une compétition peut s'établir lorsque les prémisses initiales et les états de l'environnement se trouvent confinés au sein d'un petit espace ?
2. Comment des prémisses éloignées et éparées par rapport à la majorité peuvent-elles rejoindre des états sensoriels également éparées et éloignés ?

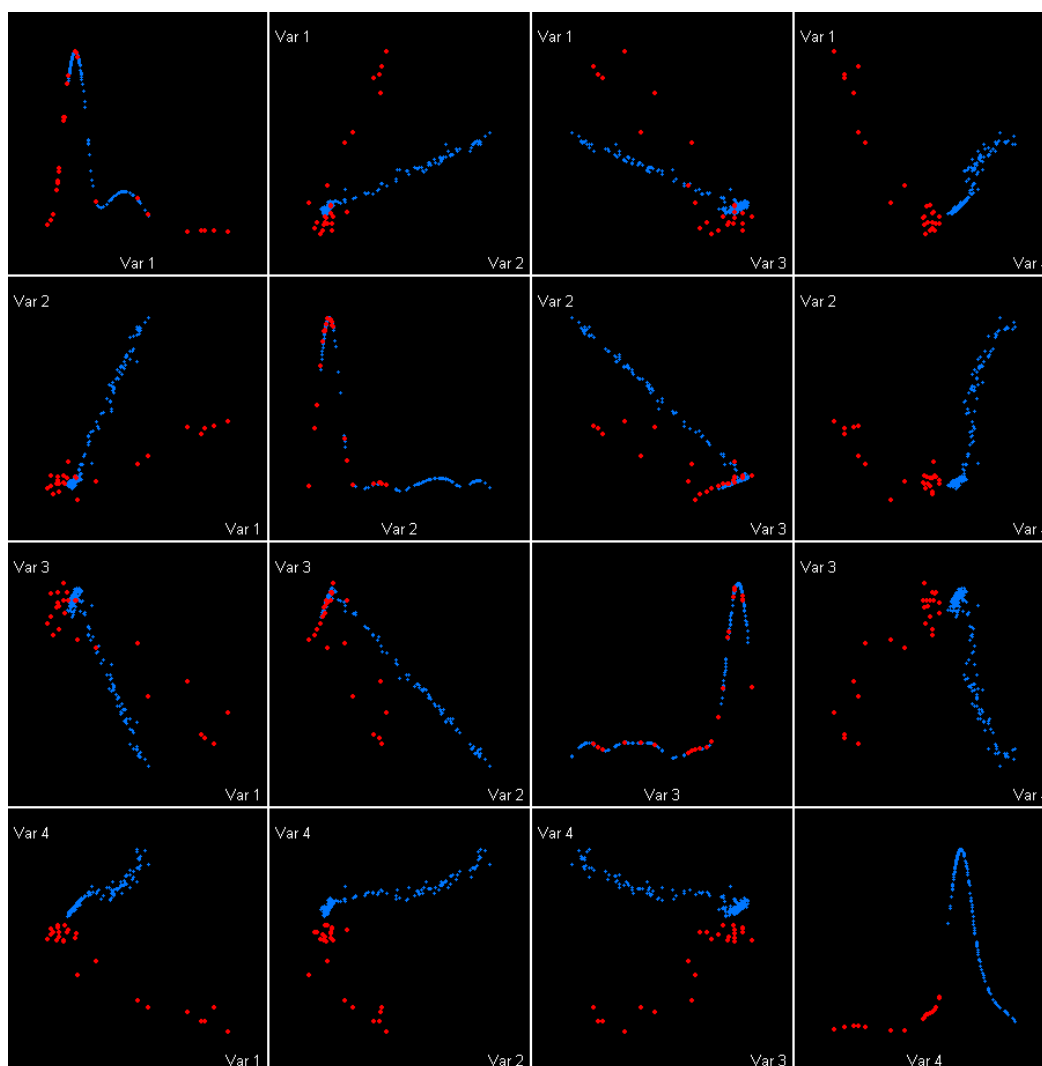


Figure IV-29 Comparaison entre les états sensoriels de l'environnement imposé S1 (ronds rouges) et les prémisses des règles initiales (croix bleues)

4. Étude du système dans des environnements imposés

L'architecture pragmatiste proposée conduit à des systèmes cognitifs possédant une dynamique chaotique, c'est-à-dire avec une forte sensibilité aux conditions initiales. Toutefois, la reproductibilité et la comparabilité des expériences peuvent être assurées en imposant les valeurs sensorielles indépendamment des conclusions motrices. Cette technique occulte l'aspect comportemental du robot mais offre un cadre fiable pour évaluer les principaux facteurs influençant la dynamique du système. Plus exactement, trois types de facteurs ont été étudiés. Les deux premiers, les paramètres fondamentaux du système, tels que l'enchère ou la taxe, etc. et la périodicité des stimulations sensorielles, concernent les capacités d'auto-organisation du système. Le troisième, la manière d'incérer de nouvelles règles, porte sur la stabilité du système face aux perturbations externes ou internes. Ce dispositif expérimental offre également la possibilité de vérifier le déroulement du mécanisme d'anticipation proposé.

Sauf spécification, toutes les expériences durent 30 min avec une fréquence de 10 Hz, ce qui revient à 18000 élections avec les règles E1, E2 ou E3 et à 9000 élections avec les règles E1-bis. L'alpha des règles initiales vaut 1 et la force vaut 0,9 afin que les règles ne soient pas éliminées trop vite par la taxe ou les enchères. Par ailleurs, les graphiques représentent les caractéristiques de la règle nouvellement élue à chaque pas de temps. Le nombre relativement petit de stimulations sensorielles différentes ainsi que leur périodicité entraînent l'élection régulière de certaines règles. Les élections rapprochées de ces règles dessinent sur les graphiques des courbes qui offrent ainsi une indication sur l'évolution de la dynamique.

4.1. Caractérisation des paramètres fondamentaux

4.1.1. Les taux de taxe et d'enchère

Les coefficients d'enchère et de taxe régissent la décroissance de la force des règles. Toutefois, chacun possède un rôle distinct : la taxe doit permettre l'élimination des règles inutiles et l'enchère doit permettre celle des règles redondantes. L'expérience suivante ambitionne d'une part de déterminer leurs influences mutuelles sur la dynamique du système et d'autre part de déterminer les valeurs limites de viabilité. Les 18 expériences ont été réalisées avec le scénario S1 et avec les règles initiales E1-bis, soit une fréquence de sélection aboutie de 5 Hz. De tous ces essais, un cas de référence sera dégagé sur lequel s'appuieront les études sur l'influence des autres paramètres. Le mécanisme d'enchère se trouve étroitement lié au paramètre de remboursement. Afin de préserver l'équilibre de ce mécanisme, le paramètre de remboursement est égal à celui de l'enchère. Le Tableau IV-7 récapitule les expériences réalisées en fonction des paramètres, parmi lesquelles trois groupes se dégagent : (A) les expériences avec une taxe dominante, (B) celles avec une enchère dominante et (C) celles qui ménagent ces deux paramètres. En plus, de ces trois groupes d'expériences, (D) une observation commune à quasiment toutes les expériences sera relevée.

		Remboursement=Enchère						
		0,001	0,01	0,02	0,03	0,04	0,05	0,06
Taxe	0,0001	1	4	6	8	10	13	16
	0,001	2	5			11	14	17
	0,002	3		7		12	15	18
	0,003				9			

Tableau IV-7 : Récapitulatif des expériences réalisées en fonction du taux de la taxe et de celui de l'enchère.

A - Forte domination de la taxe (Exp1, Exp2 et Exp3)

Dans l'expérience 1, l'enchère et la taxe étant très faibles, aucune règle n'est éliminée. Seul l'algorithme d'optimisation des prémisses influence la dynamique. Quelques règles se déclenchent fréquemment. Par exemple, la règle numéro 1148 s'est déclenchée 1802 fois avec une fréquence d'environ 1Hz. Cependant, la majorité (~140) se déclenche plus ou

moins régulièrement avec une fréquence bien inférieure à celle de la séquence imposée (Figure IV-30). Au total, 150 règles sur les 200 ont été déclenchées au moins une fois.

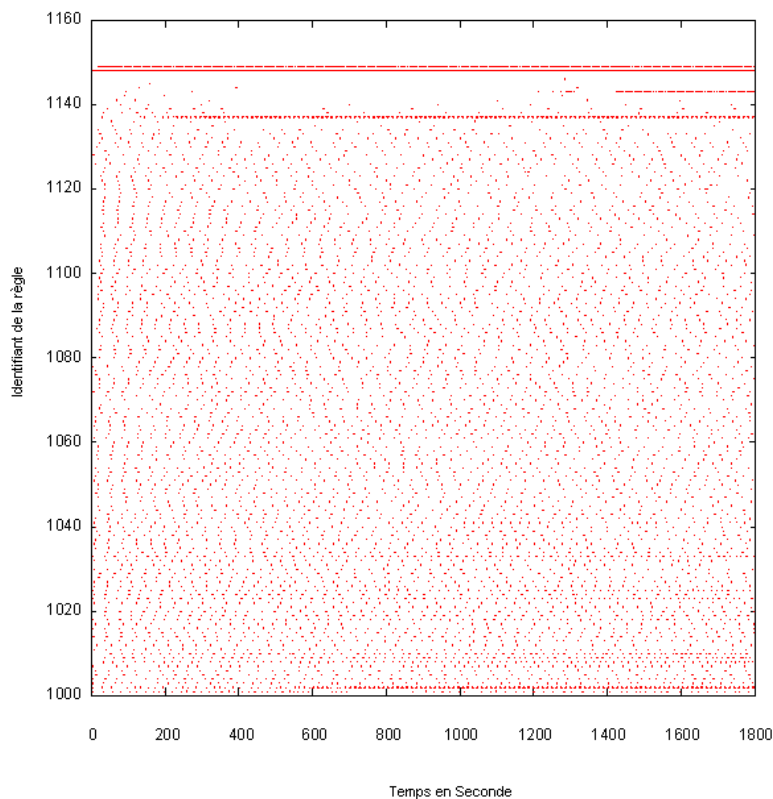


Figure IV-30 : Graphique indiquant le numéro de la règle déclenchée à chaque sélection durant l'expérience 1.

Dans la Figure IV-30, certaines règles se déclenchent suffisamment régulièrement pour dessiner une droite, comme la règle numéro 1148, alors que d'autres se déclenchent sur des intervalles de plusieurs cycles de scénario. Cette situation suggère la présence de règles relativement isolées par rapport aux autres et bien adaptées dès le début. L'isolement se confirme par la rapide augmentation de l'alpha de ces règles (Figure IV-31) puisque l'alpha représente le coefficient de mélange dans l'algorithme Estimation-Maximisation décrit lors de la présentation du processus de sélection d'une règle dans le chapitre précédent.

Concernant l'évolution du score, deux situations se présentent (Figure IV-31) : soit la prémisse sensorielle capture un seul motif et se spécialise sur celui-ci, augmentant ainsi le score, soit la prémisse sensorielle tend à capturer plusieurs motifs en se plaçant sur un motif médian impliquant une diminution du score. Le deuxième cas de figure se réalise seulement si aucune autre règle ne se trouve à proximité. Par ailleurs, l'initialisation des variances étant relativement large afin de couvrir au mieux l'espace sensoriel, toutes les règles se spécialisent, ce qui se traduit par une diminution d'un facteur 2 environ de leurs variances.

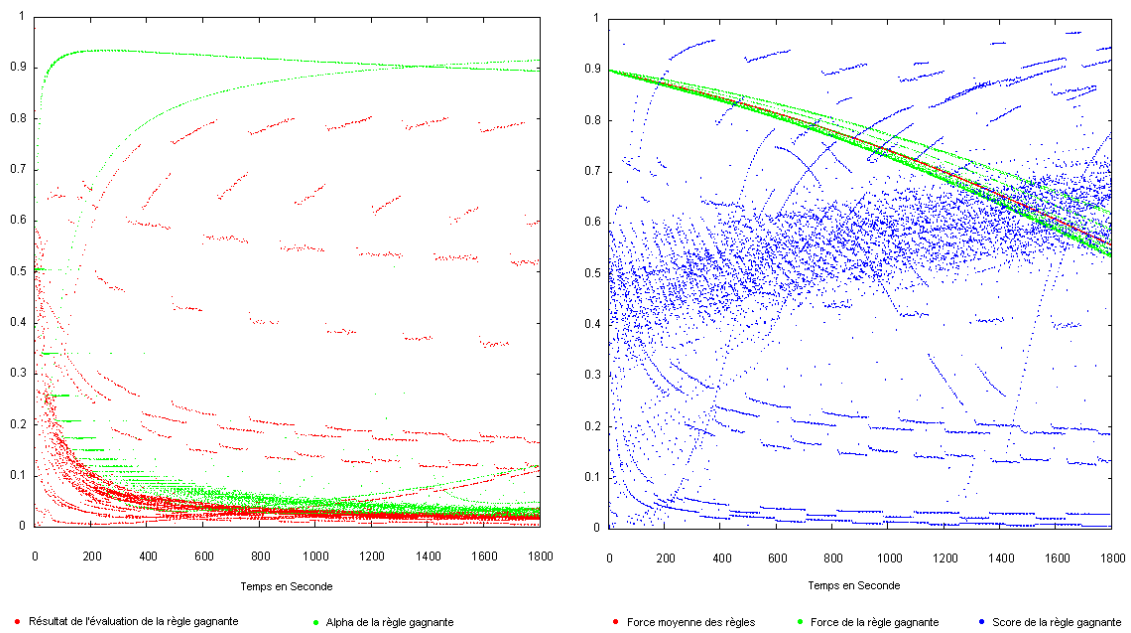


Figure IV-31 : Représentation graphique de l'alpha et du résultat de l'évaluation de la règle déclenchée à chaque sélection durant l'expérience 1, à droite et à gauche de la force et du score de la règle déclenchée ainsi que la force moyenne des règles.

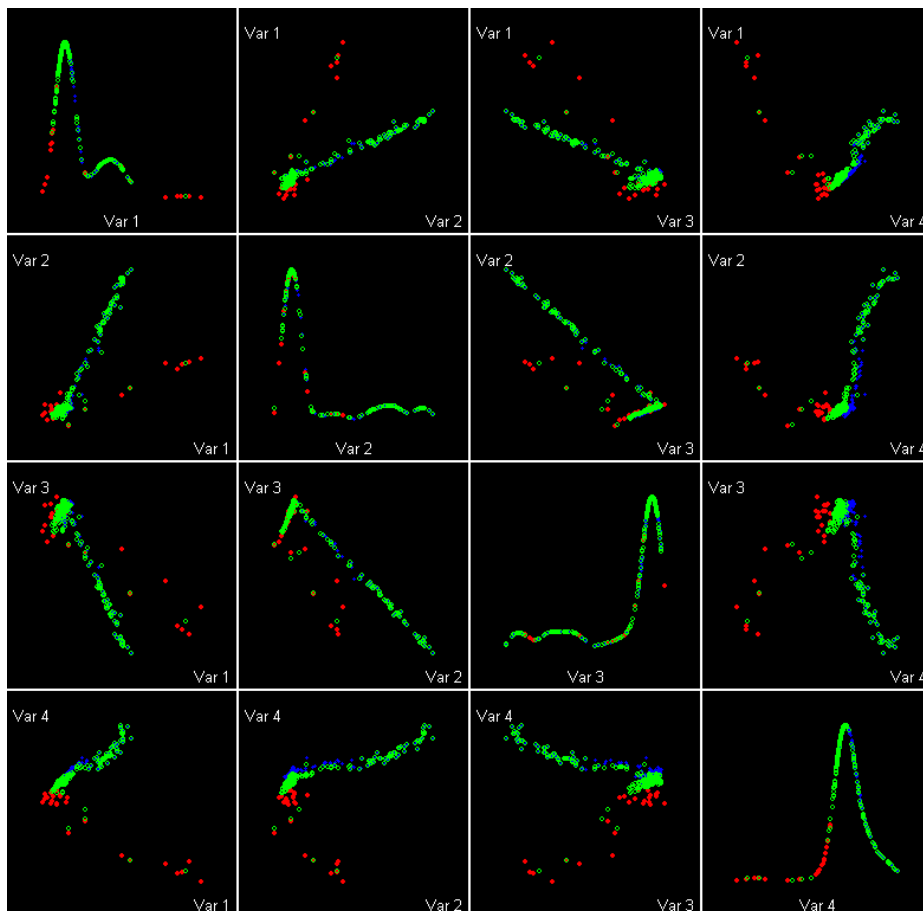


Figure IV-32 : Représentation de l'environnement (ronds rouges), des prémisses appartenant aux règles initiales (croix bleues) et des prémisses appartenant aux règles finales (cercles verts).

Lors de leur déclenchement, certaines règles présentent (Figure IV-31) à la fois un alpha faible et un score élevé. Cette situation révèle une sur-spécialisation, c'est-à-dire qu'un ensemble de règles parvient à se spécialiser sur le même motif, ce que confirment également la faible valeur d'évaluation de ces règles et l'augmentation globale de leurs scores. Ainsi, la compétition ne parvient pas à dégager une règle unique, comme l'illustre l'augmentation globale des scores. Par ailleurs, cette situation est exacerbée par le nombre élevé de règles similaires et par la pauvre variabilité de l'environnement imposé du scénario S1.

La Figure IV-32 représente les prémisses sensorielles initiales et finales par rapport aux 25 stimulations sensorielles effectives de l'environnement imposé. Seuls les motifs sensoriels impairs du scénario S1 participent à la sélection de règles puisque les règles de gestion empêchent la sélection une fois sur deux et que le nombre total de motifs sensoriels de la séquence utilisée par le scénario S1 est pair. Cette figure confirme les propos précédents sur l'existence d'un groupe de prémisses situées dans une zone très concurrentielle, sur l'immigration de règles vers des motifs multiples sans concurrence et enfin se devinent une ou deux règles situées dès le début à proximité de motifs sensoriels spécifiques.

L'ensemble des cercles verts avec une croix bleue au centre traduit les 50 règles qui ne se sont jamais déclenchées à cause de leur éloignement par rapport aux stimulations sensorielles. Le nuage de cercles se situe à mi-chemin entre l'amas de croix bleues et celui de ronds rouges, ce qui montre l'adaptation de cet ensemble de règles. Quelques cercles verts se trouvent à proximité ou au même emplacement que les ronds rouges, ils représentent les règles qui se sont déclenchées le plus fréquemment.

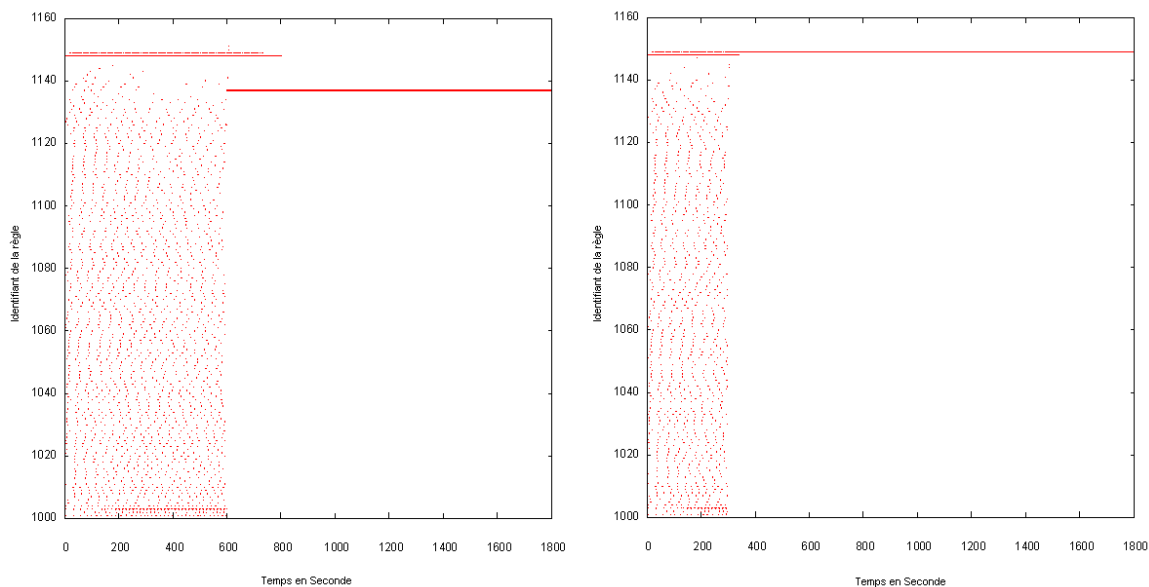


Figure IV-33 : Graphiques indiquant le numéro de la règle gagnante à chaque élection, à gauche concernant l'expérience 2 et à droite concernant l'expérience 3.

Le faible taux d'enchère explique le peu de différence entre les règles concernant la valeur de leur force. Ce trait se trouve particulièrement visible dans les expériences 2 et 3 (Figure IV-33), où l'élimination des règles s'effectue en même temps pour chacune des deux expériences à 600 s et 250 s. Une ou deux exceptions demeurent pour les règles

bénéficiant d'un remboursement élevé. Après l'élimination brutale des autres règles, le remboursement devient maximal à chaque sélection.

En définitive, dans ces trois premières expériences, la valeur de la taxe détermine le moment de l'élimination des règles en dehors du phénomène de la compétition puisque l'enchère est faible.

B - Forte domination de l'enchère (Exp4, Exp6, Exp8, Exp10, Exp13 et Exp16)

Les expériences 4, 6, 8, 10, 13 et 16 montrent l'évolution de la dynamique du système avec une très faible influence de la taxe, particulièrement pour les cinq dernières. Dans l'expérience 4, toutes les règles ont été conservées. Pour les cinq expériences suivantes, l'évolution de la démographie des règles (Figure IV-34) prend la forme suivante : un plateau puis une chute pendant 500 secondes puis une nouvelle stabilisation de la population. Seul le moment de la transition les différencie. La conservation de la forme des courbes s'explique par le maintien de l'équilibre entre le remboursement et l'enchère. Ces cinq expériences possèdent au final le même nombre de règles encore présentes et avec une répartition similaire dans l'espace sensoriel.

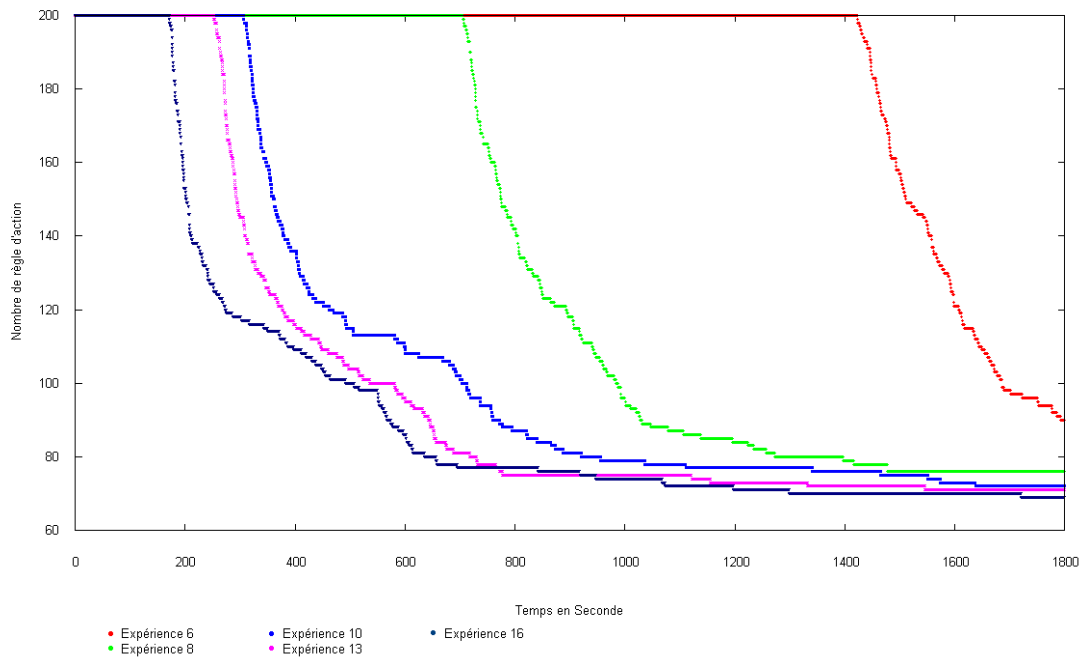


Figure IV-34 : Évolution du nombre de règles sensorimotrices dans les expériences 6, 8, 10, 13 et 16.

Les trois phases décrites se retrouvent également dans les graphiques affichant les caractéristiques de la règle déclenchée. Par exemple, la Figure IV-35 présentant les résultats de l'expérience 13 offre une nouvelle lecture de ces trois phases. Pendant la première phase, beaucoup de règles se déclenchent successivement et parallèlement leur force diminue rapidement à cause de l'enchère élevée due à une forte redondance. Ce phénomène se traduit aussi par la diminution de l'alpha. La deuxième phase commence dès l'élimination des premières règles avec une force inférieure au seuil. La concurrence devient moins sévère, entraînant la diminution des enchères et l'augmentation de la force. Malgré tout, la force et l'alpha des règles diminuent tant qu'il existe un certain niveau de concurrence. Le score peut augmenter ou diminuer, indiquant respectivement la spécialisation de certaines

règles et la généralisation d'autres. La troisième phase apparaît lorsque le remboursement compense la taxe et l'enchère, de sorte que la compétition entre règles voisines ne peut plus entraîner d'élimination. À la fin des expériences 6, 8, 10, 13 et 16, il reste en moyenne 68 règles, dont une cinquantaine n'a jamais été déclenchée.

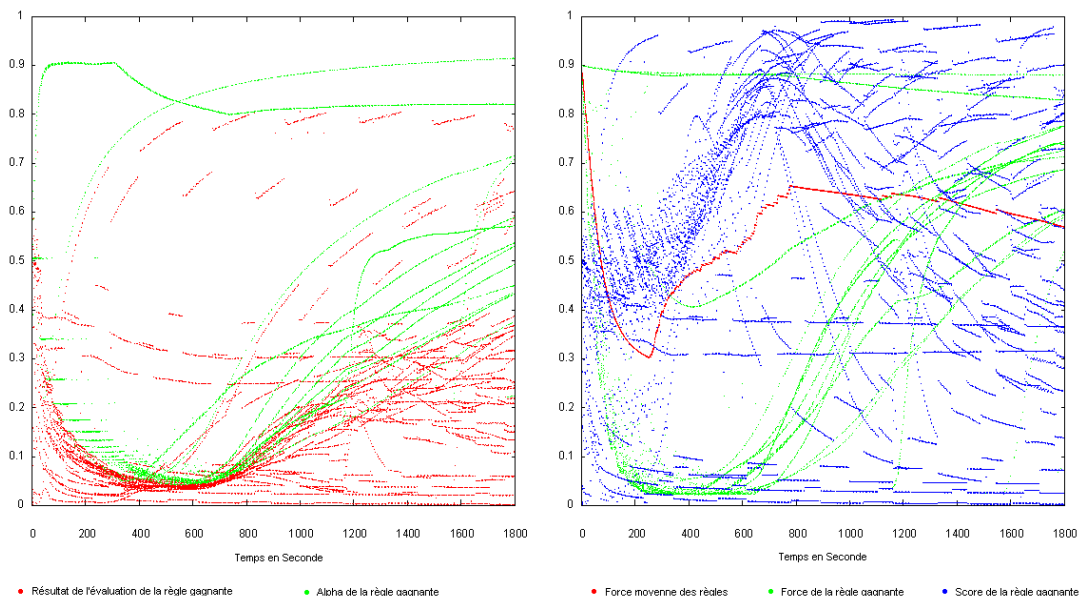


Figure IV-35 : A droite, représentation graphique de l'alpha et du résultat de l'évaluation de la règle déclenchée à chaque sélection durant l'expérience 13. À gauche, représentation de la force et du score de la règle déclenchée ainsi que la force moyenne des règles.

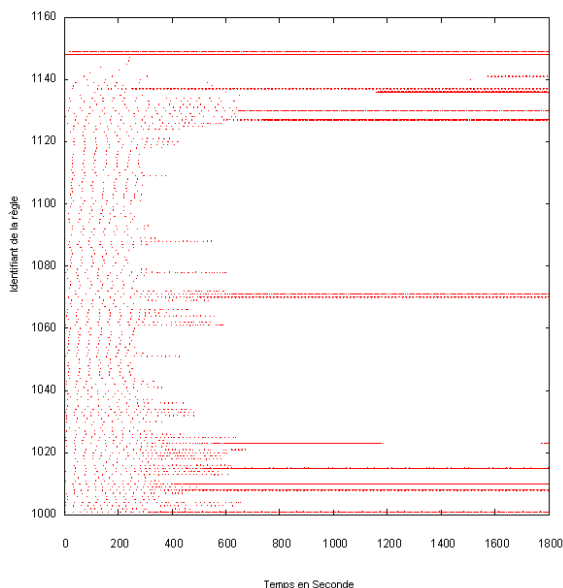


Figure IV-36 : Graphique indiquant le numéro de la règle déclenchée à chaque sélection au cours de l'expérience 13.

De même, ces trois phases transparaissent dans le graphique affichant la règle déclenchée à chaque sélection, la Figure IV-36. La Figure IV-37 affiche en fausse 3D la répartition des prémisses initiales et finales des expériences 6 et 10, ainsi que les états

sensoriels de l'environnement imposé par le scénario S1. Cette visualisation confirme que les prémisses finales des différentes expériences couvrent les mêmes zones. Par ailleurs, les règles s'étant peu déclenchées, et par conséquent étant moins sujettes au réajustement des prémisses, elles sont restées proche de leur zone d'origine. Par exemple, dans la Figure IV-37, une quinzaine de cercles oranges, représentant les prémisses finales des règles qui se sont déclenchées au moins une fois lors de l'expérience 6, se situent à proximité de l'axe dessiné par les croix bleues, qui désignent les prémisses initiales.

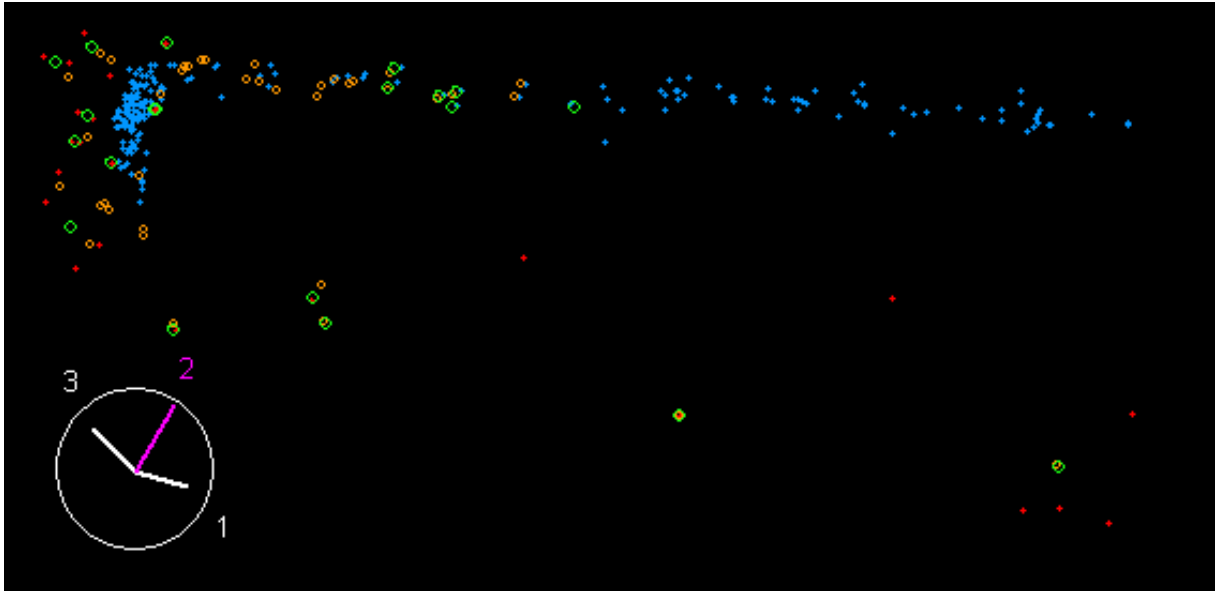


Figure IV-37 : Représentation des trois premiers composants des prémisses initiales (croix bleus), des prémisses de règles restantes s'étant déclenchées au moins une fois (celles de l'expérience 6 : cercles oranges, celles de l'expérience 10 : cercles verts) et les valeurs de l'environnement imposé (croix rouges).

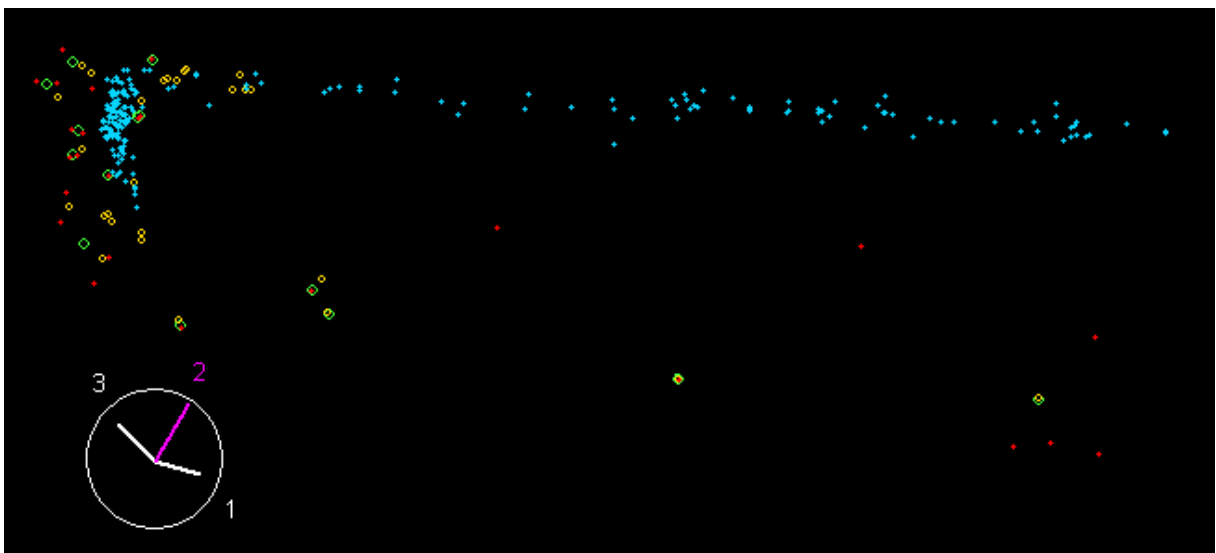


Figure IV-38 : Représentation des trois premiers composants des prémisses initiales (croix bleus), des prémisses de règles restantes s'étant déclenchées plus de 20 fois (celles de l'expérience 6 : cercles oranges, celles de l'expérience 10 : cercles verts) et les valeurs de l'environnement imposé (croix rouges).

Dans l'ensemble des expériences 6, 8, 10, 13 et 16, les prémisses les plus réajustées se concentrent sur les mêmes points ou sur les mêmes zones de points. La Figure IV-38 illustre cette concentration en affichant seulement les règles dont le nombre de déclenchements dépasse 20.

Au final, dans cette série d'expériences, l'enchère détermine le début de l'élimination de règles due à la complétion. Cette série d'expériences suggère également des propriétés d'auto-organisation et de classification puisque les règles se réajustent par rapport aux motifs imposés.

C - L'enchère et la taxe (Exp5, Exp7, Exp9, Exp11, Exp12, Exp14, Exp15, Exp17 et Exp18)

L'exploration de l'espace des paramètres par les expériences 5, 7 et 9 a permis d'identifier l'existence d'une troisième forme d'évolution de la démographie de la base de règles (Figure IV-39). Les expériences 5 et 7 conduisent à l'élimination totale des règles, la transition brusque suggère que la valeur de l'enchère reste trop faible pour générer une compétition.

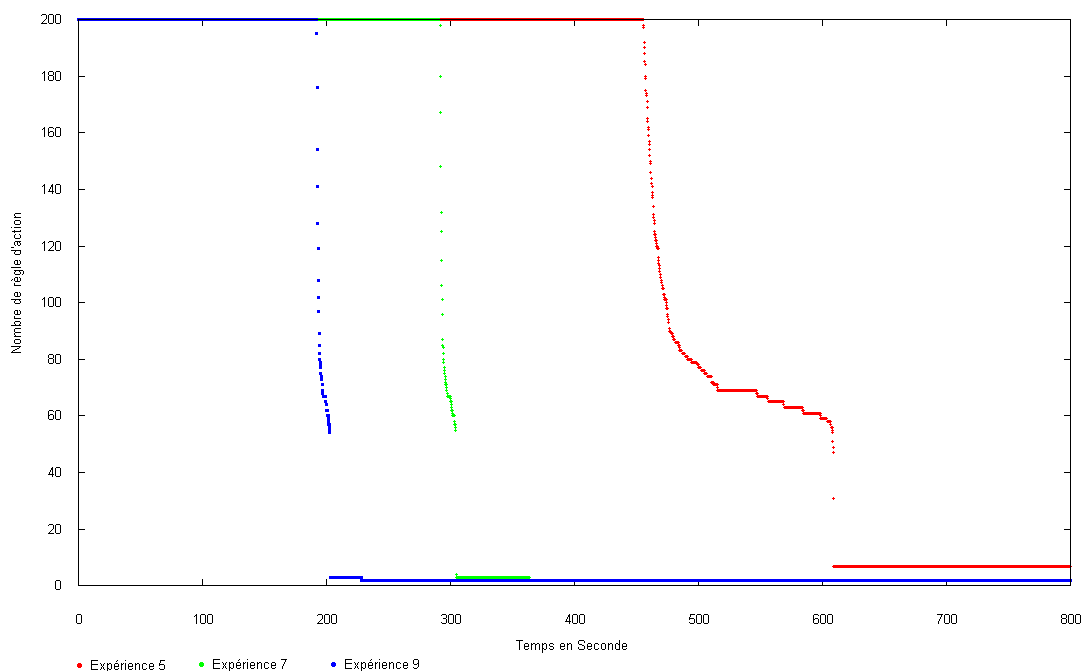


Figure IV-39 : Évolution du nombre de règles sensorimotrices dans les expériences 5, 7, et 9.

La présence de cinq phases particulièrement visibles se dégagent sur les courbes de population des expériences 11, 12, 14, 15, 17 et 18, Figure IV-40 : (1) un plateau suivi (2) d'une transition abrupte au début puis s'inclinant jusqu'à (3) un plateau suivi (4) d'une brusque transition débouchant sur (5) un plateau. Concernant la première transition, le moment d'apparition diffère selon chaque expérience. En revanche, pour la seconde, elle survient à 600 s dans le cas des expériences 11, 14, et 17 alors que dans le cas des expériences 12, 15, et 18 elle survient à 300 s. La première observation montre que les paramètres d'enchère et de taxe déterminent le début des premières éliminations des règles qui n'ont pas réussi à s'adapter suffisamment par rapport à d'autres. La taxe influence plus précisément le début de cette première transition et la durée de la stabilisation puisqu'un

taux élevé revient à atteindre plus rapidement le seuil d'élimination. Quant à l'enchère couplée au remboursement, celle-ci influence légèrement le moment de la transition ainsi que l'inclinaison de la transition puisque plus le taux d'enchère est élevé, plus la compétition devient dangereuse pour les règles ayant une faible fréquence de déclenchement. La seconde observation montre que seule la taxe joue un rôle dans l'apparition de la seconde transition. En effet, la taxe sanctionne les règles qui se sont peu ou pas déclenchées et pour lesquelles le taux d'enchère n'intervient pas. Par exemple, dans les expériences 10, 13 et 18, la taxe est très faible et les règles ne subissant pas ou peu la compétition n'atteignent pas le seuil d'élimination au bout de 30 min.

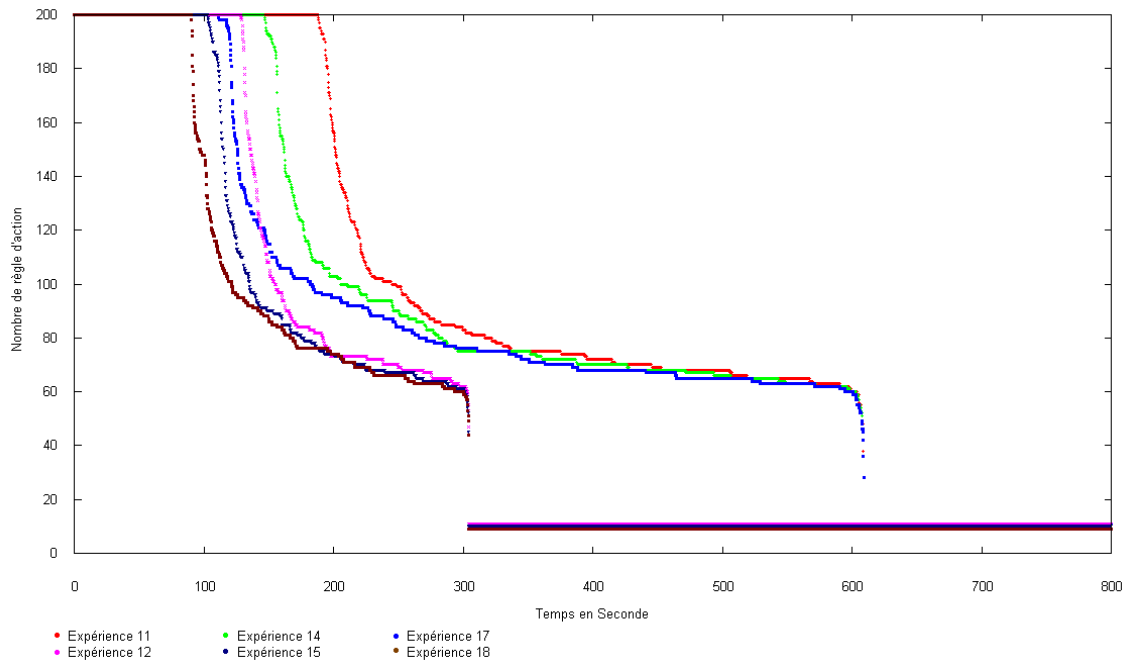


Figure IV-40 : Évolution du nombre de règles sensorimotrices dans les expériences 11, 14, 17, 12, 15 et 18.

Par ailleurs, dans les expériences 10, 13 et 16, les règles à l'équilibre sont au nombre d'environ 70. Dans les expériences 11, 12, 15, 17 et 18, le nombre de règles à l'équilibre atteint environ une quinzaine. Malgré l'augmentation du taux de la taxe, la stabilité de ce nombre confirme l'existence de deux types de règles après la seconde transition : les règles inutilisées et les règles utilisées.

Les cinq phases identifiées transparaissent également sur les courbes de l'évolution du score, de la force, de l'alpha, de l'évaluation, et plus particulièrement avec la courbe de la force moyenne des règles (Figure IV-41). Selon cette dernière, les cinq phases se traduisent par : (1) la décroissance générale de la force des règles à cause du taux d'enchère puis (2) à 180 s la brusque élimination des règles avec une force en dessous du seuil rehausse la moyenne. Les éliminations deviennent plus progressives car le remboursement compense la taxe et l'enchère selon la fréquence de déclenchement de la règle. (3) Il ne reste alors que les règles un peu adaptées et celles qui ne participent pas. (4) La seconde hausse de la force moyenne reflète l'élimination des règles non utilisées, enfin (5) la valeur moyenne demeure stable. Ainsi, une manière de comparer les expériences 11, 12, 15, 17 et 18 consiste à apposer leurs courbes de la force moyenne (Figure IV-42).

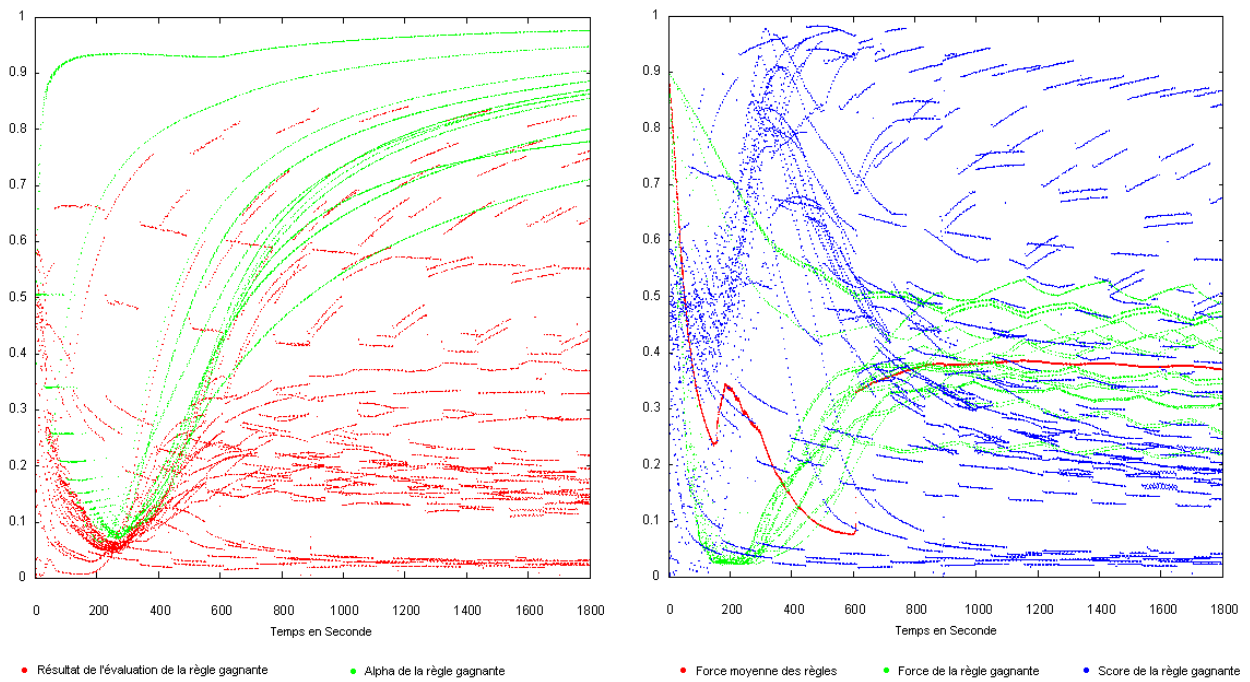


Figure IV-41 : A droite, représentation graphique de l'alpha et du résultat de l'évaluation de la règle déclenchée à chaque sélection durant l'expérience 14 et à gauche, représentation de la force et du score de la règle déclenchée ainsi que la force moyenne des règles.

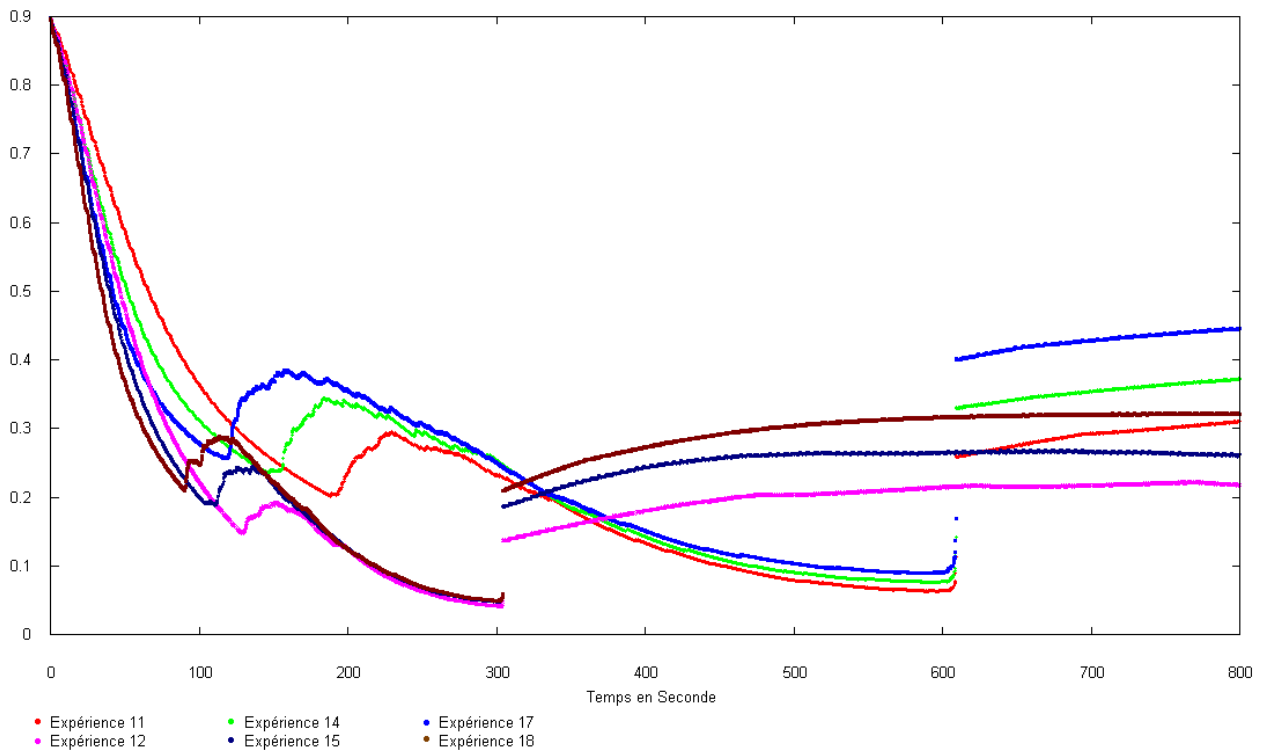


Figure IV-42 : Forces moyennes des règles sensorimotrices des expériences 11, 12, 14, 15, 17 et 18.

L'étude de ces phases indique, comparativement aux expériences 10, 13 et 16 par exemple, que les règles finales doivent correspondre davantage à l'environnement imposé, soit les états impairs de la séquence A composant le scénario S1.

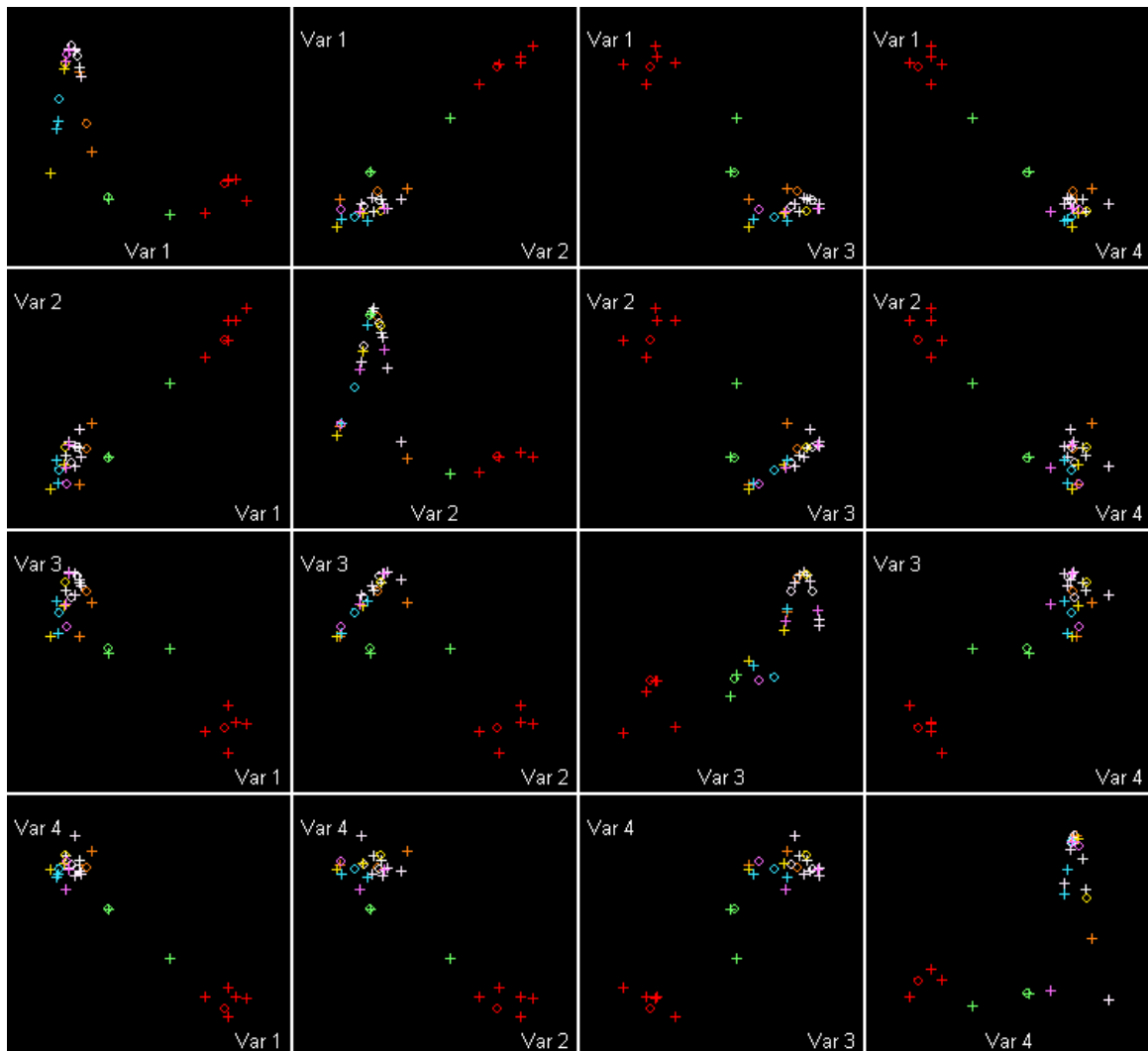


Figure IV-43 : Représentation graphique des prémisses finales qui ont été déclenchées plus d'une fois lors de la dernière séquence (cercles) et, avec la couleur associée, les états sensoriels de l'environnement imposé les déclenchant lors de la dernière séquence du scénario S1 (croix) de l'expérience 17.

A partir de l'expérience 17, Figure IV-43 représente les prémisses sensorielles appartenant aux règles déclenchées par plusieurs états sensoriels durant la dernière séquence, deux prémisses spécialisées sur un seul état sensoriel ont notamment été retirées. Le regroupement des états sensoriels autour de leurs prémisses représentatives se trouve compliqué par leur proximité et par le rapport entre le nombre d'états sensoriels et le nombre de prémisses. Toutefois, la prémisses rouge se centre au sein d'un nuage de points rouges selon les quatre dimensions. Cette prémisses a dû s'ajuster dans un contexte de faible concurrence.

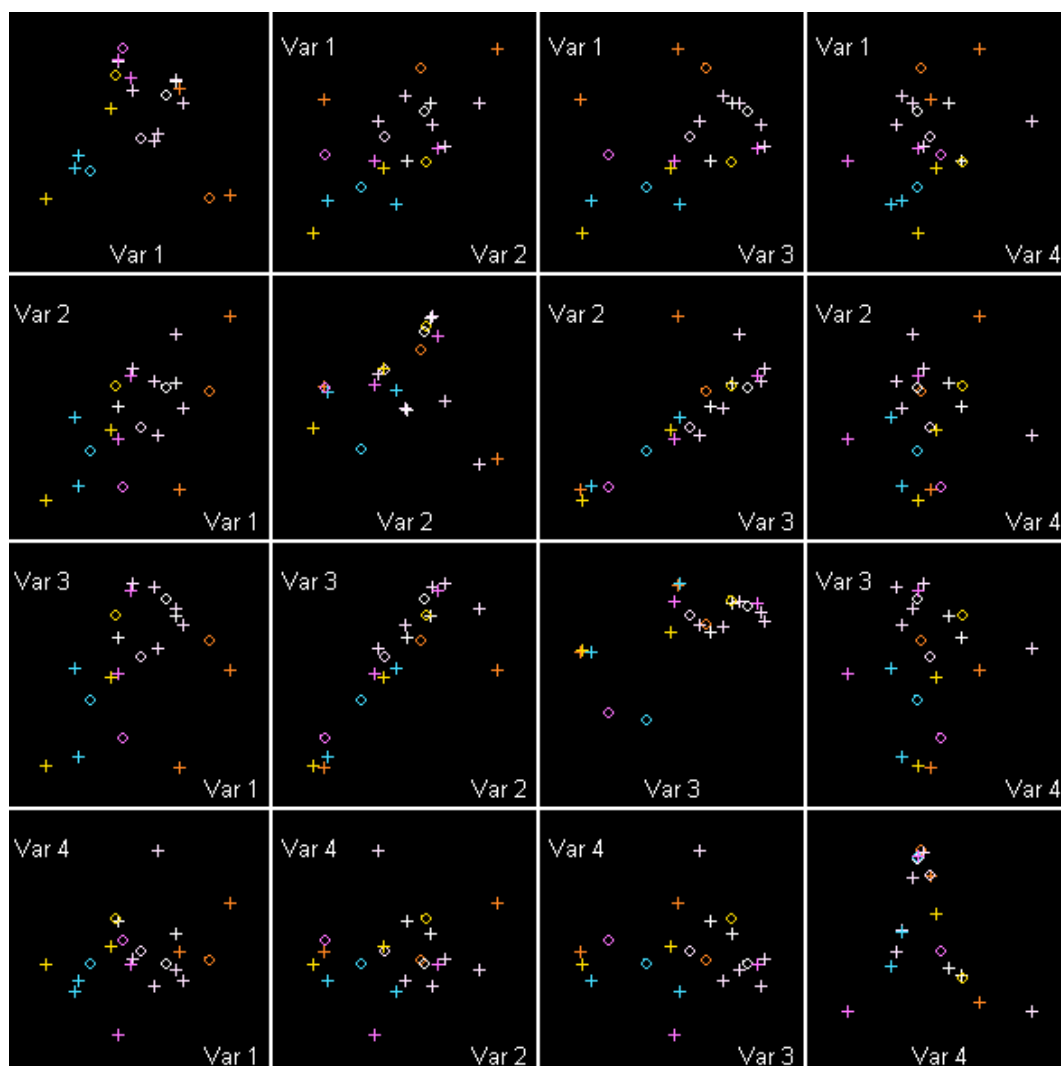


Figure IV-44 : Représentation graphique identique à la figure IV-43 mais sans les règles et les états sensoriels coloriés en vert et en rouge.

Un grossissement (Figure IV-44) sur les prémisses et les états sensoriels proches confirme l'imbrication des champs de recouvrement des prémisses, bien que l'observation en quatre dimensions ne soit pas aisée. Toutefois, la comparaison avec les prémisses finales des autres expériences (Figure IV-45) suggère une certaine régularité dans la configuration spatiale des prémisses finales. Cette régularité signifie que les prémisses sensorielles finales sont relativement indépendantes de l'enchère pour des variations de l'ordre de 0,02 et de l'ordre de 0,002 pour la taxe.

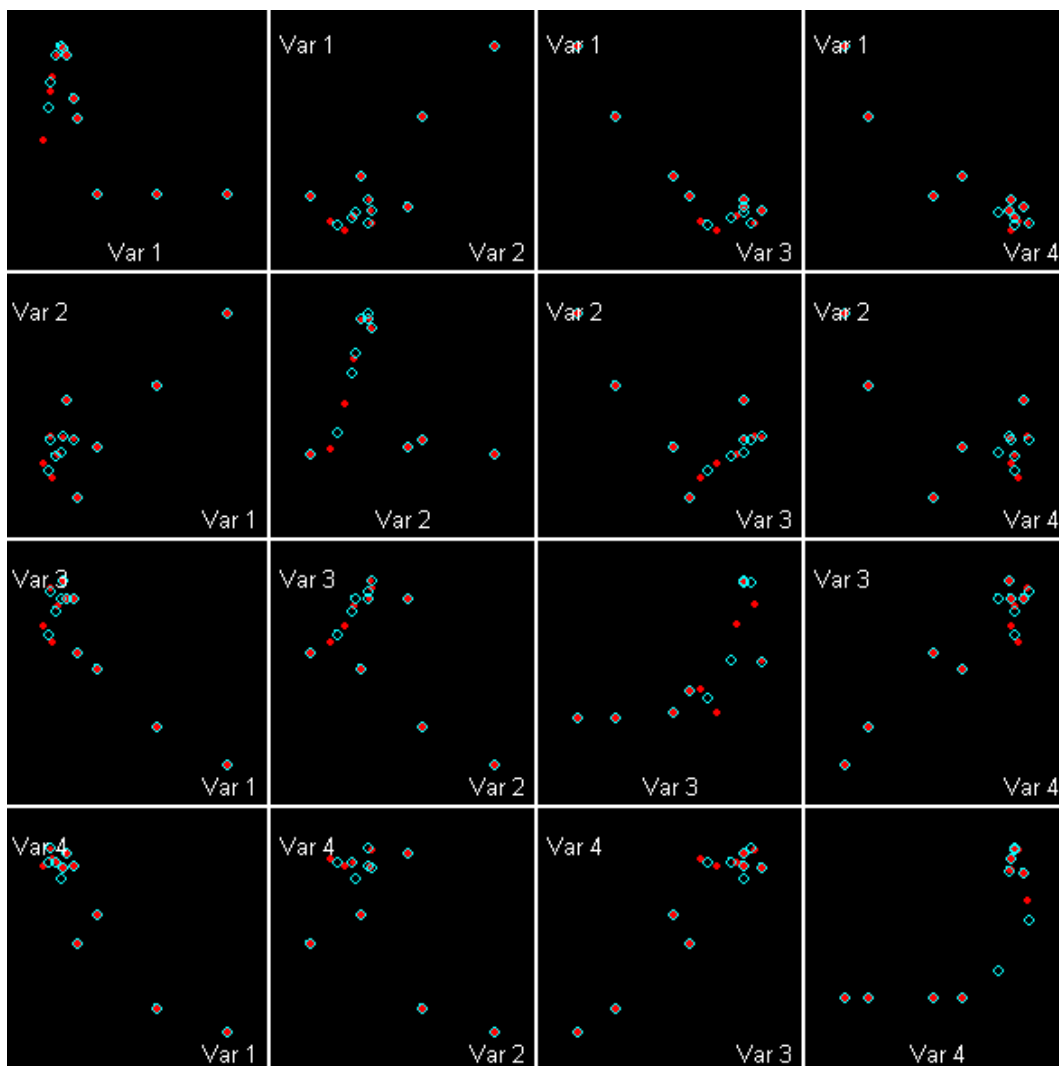


Figure IV-45 : Comparaison des prémisses finales entre l'expérience 11 (cercles bleus) et l'expérience 17 (ronds rouges).

Parmi tous ces essais, l'expérience 14 possède distinctement 5 phases dans son évolution démographique (Figure IV-40) et une stabilisation de la force qui n'est ni trop haute ni trop basse pour posséder une sensibilité plus grande en cas de modification de la dynamique (Figure IV-42). Le taux de la taxe et celui de l'enchère pour les expériences sur l'influence des autres paramètres vaudront ceux de l'essai de référence, l'expérience 14.

D - Observations à part : les oscillations

Sur les graphiques illustrant les évolutions des caractéristiques de la règle déclenchée des expériences 1, 3, 6, 7, 8, 11, 12, 13, 14, 15, 16 et 17, de brusques changements apparaissent périodiquement (Figure IV-31, Figure IV-35 et Figure IV-41). La période de ces oscillations varie entre 800 s et 300 s, celle-ci peuvent être dues à l'alternance entre deux règles ou entre des valeurs de la prémisses d'une règle qui, capturant deux sous-ensembles, oscille entre deux points d'équilibre. Les courbes de l'évolution de l'alpha étant toujours lisses, la seconde solution semble la plus probable puisqu'une alternance de règles conduirait à une alternance dans les valeurs de l'alpha. L'évolution de la force et du score oscille également, ce qui peut s'expliquer par la transition d'un pôle à un autre.

4.1.2. Le taux de remboursement

Dans les expériences précédentes, les taux de remboursement et d'enchère étaient toujours identiques afin de séparer le mécanisme de compétition et le mécanisme visant à supprimer les règles passives. Le Tableau IV-8 récapitule les 5 expériences réalisées, en plus de l'expérience 14, en fonction du taux de remboursement appliqué.

	Remboursement					
	0,03	0,04	0,05	0,06	0,07	0,08
Taxe=0,001 ; Enchère=0,05	19	20	14	21	22	23

Tableau IV-8 : Récapitulatif des expériences réalisées en fonction du taux de la taxe et du taux de remboursement.

À la fin de chaque expérience, le nombre de règles est identique, cette situation se trouve certainement favorisée par le nombre restreint d'états sensoriels (25) et leur faible variabilité. Cependant, l'influence du taux de remboursement sur l'évolution de la démographie apparaît dans la Figure IV-46. La variation du taux de remboursement se remarque également en comparant l'évolution des forces moyennes sur la Figure IV-47.

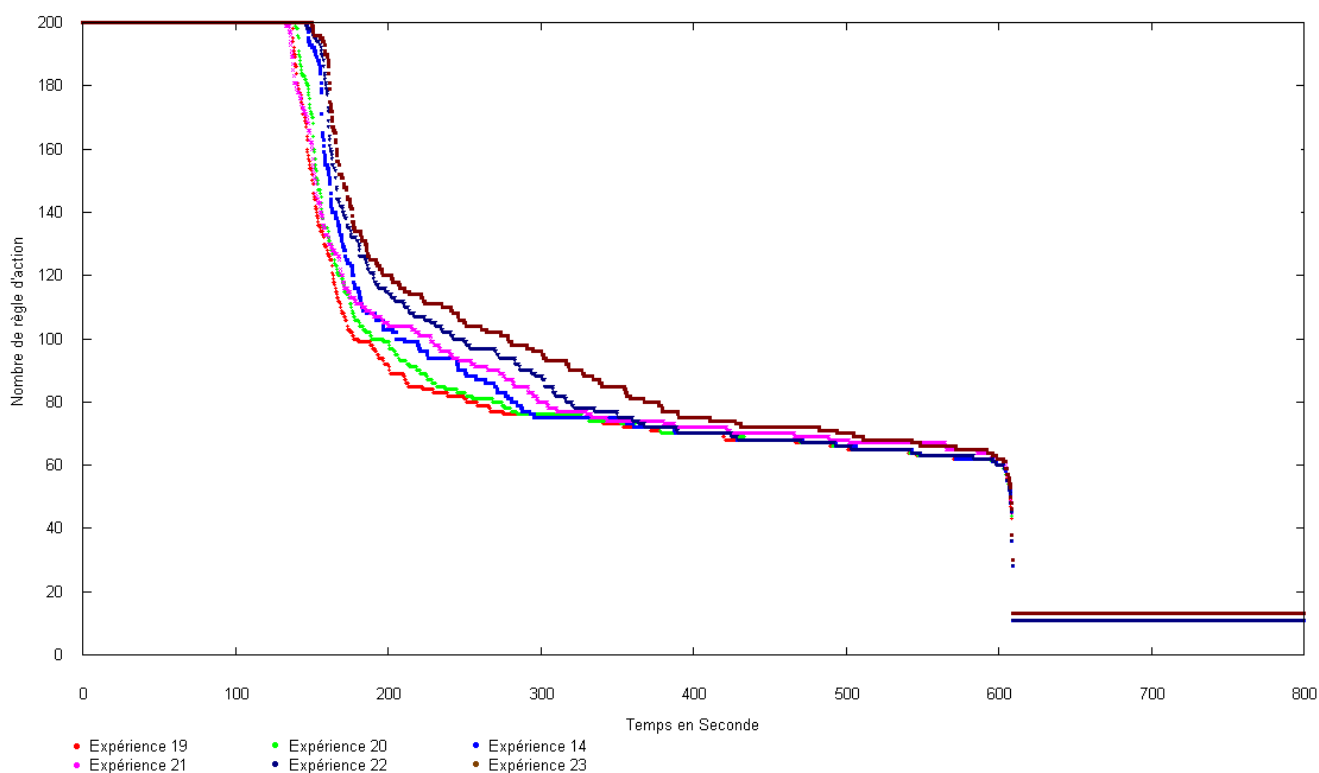


Figure IV-46 : Évolution du nombre de règles sensorimotrices dans les expériences 19, 20, 14, 21, 22 et 23.

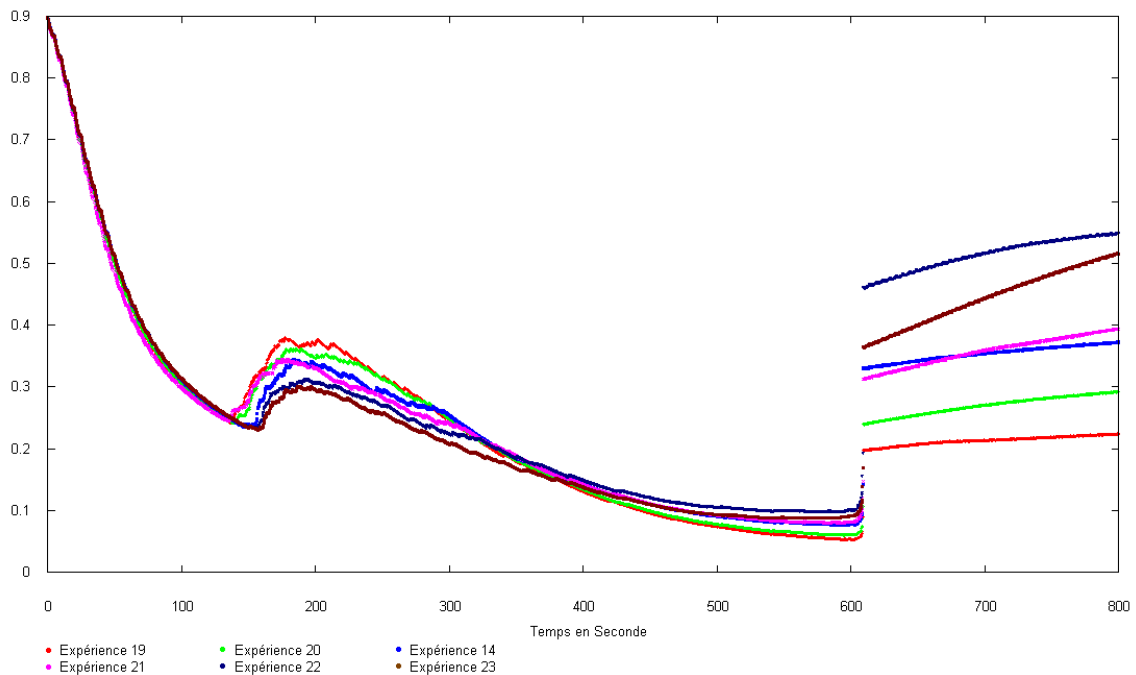


Figure IV-47 : Forces moyennes des règles au cours des expériences 19, 20, 14, 21, 22 et 23.

L'influence du taux de remboursement se situe essentiellement sur l'allongement de la période de compétition après la première brusque transition. Dans les limites de variations proposées, l'augmentation du taux de remboursement ne suffit pas à compenser l'enchère et la taxe de règle peu compétitive. L'augmentation du taux de remboursement freine la décroissance de leur force. En revanche, pour les règles possédant déjà une dynamique positive, cette augmentation accélère leur croissance et élève ainsi leur point d'équilibre. Par commodité, l'abscisse de la Figure IV-47 ne va pas au-delà de 800 s, toutefois, il faut préciser que la valeur moyenne de l'expérience 23 dépasse au final celle de l'expérience 22, de la même manière que la force moyenne de l'expérience 21 l'emporte sur celle de l'expérience 14.

4.1.3. Le nombre de règles

L'étude sur l'influence du nombre de règles initiales a été effectuée avec des taux de taxe, d'enchère et de remboursement identiques à ceux employés pour l'expérience 14. Le Tableau IV-9 récapitule la valeur des paramètres pour les 3 expériences réalisées en complément de l'expérience 14.

	E1-bis	E1 (200)	E2 (300)	E3 (600)
Taxe=0,001 ; Enchère=Remboursement=0,05	14	24	25	26

Tableau IV-9 : Récapitulatif des expériences réalisées en fonction des règles initiales utilisées.

Les règles initiales E1, E2 et E3 n'incluent pas de règles de gestion. Par conséquent, la fréquence de sélection effective s'élève à 10 Hz au lieu de 5 Hz pour les expériences précédentes initialisées avec les règles de E1-bis. L'augmentation de la fréquence induite dans les expériences 24, 25 et 26 entraîne une augmentation de la fréquence de

l'acquiescement de l'enchère. Cela se traduit par l'avancement de la première transition d'environ 80 s dans l'expérience 24 par rapport à l'expérience 14, Figure IV-48. Dans l'ensemble, les expériences convergent à l'équilibre vers une force moyenne similaire à celle de l'expérience 14. Les règles payant plus d'enchère ont une diminution de leur force plus rapide, jusqu'à atteindre le seuil de l'élimination. En revanche, les expériences 14 et 24 ne diffèrent pas aussi bien sur le nombre de règles à l'équilibre, qui est de 12, que sur la répartition des prémisses finales dans l'espace sensoriel. L'augmentation de la fréquence de fonctionnement par l'absence de règles de gestion n'entraîne pas un biais significatif.

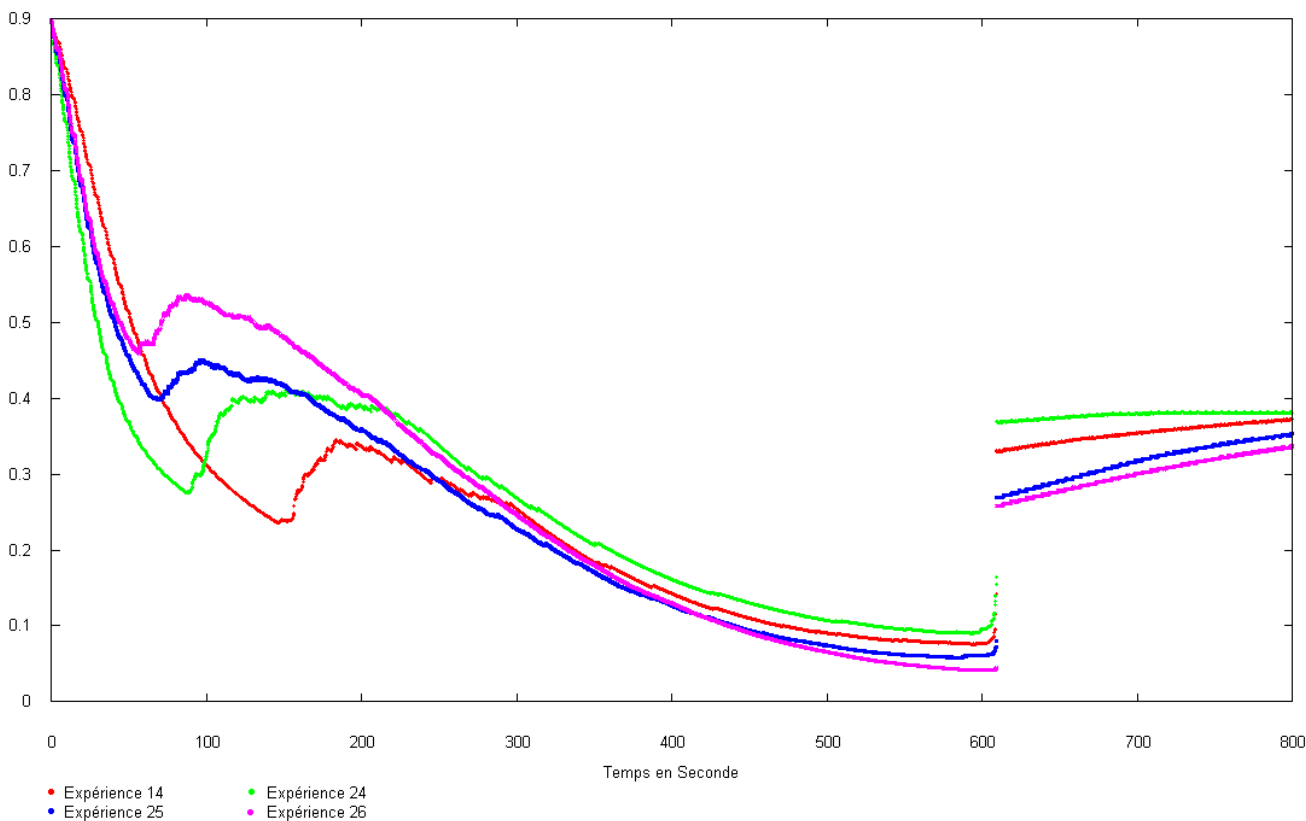


Figure IV-48 : Forces moyennes des règles au cours des expériences 14, 24, 25 et 26.

Par ailleurs, l'augmentation du nombre de règles au départ augmente les chances d'obtenir une règle adaptée à une situation et de s'y maintenir, favorisant la sur-spécialisation. Les expériences 14 et 24 conservent respectivement au final 11 et 12 règles, alors que les expériences 25 et 26 obtiennent respectivement 19 et 18 règles. La répartition des prémisses finales dans l'espace sensoriel se concentre sur la zone de forte densité des états sensoriels de l'environnement imposé (Figure IV-49).

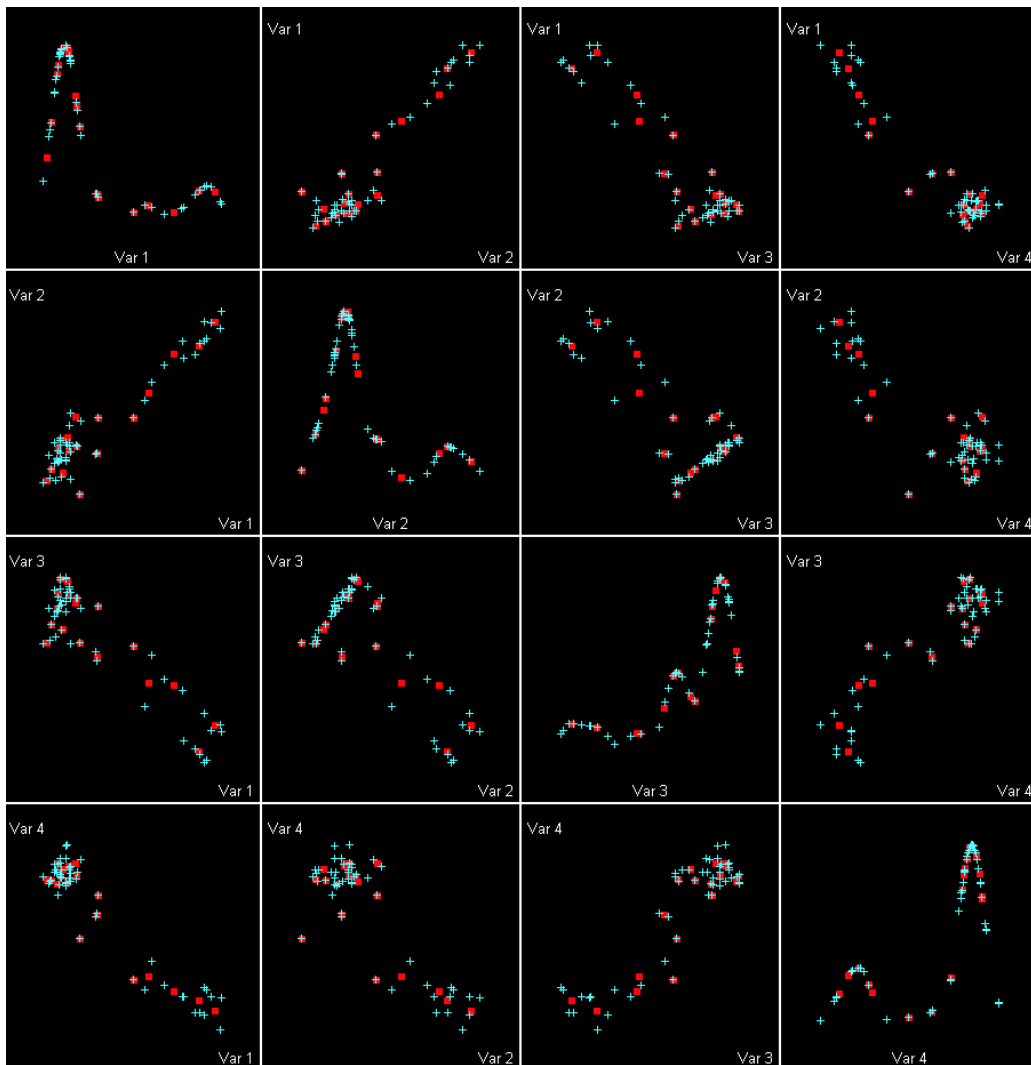


Figure IV-49 Comparaison entre les prémisses finales de l'expérience 26 (carrés rouges) avec tous les états sensoriels de l'environnement imposé (croix bleus).

4.1.4. Le seuil d'élimination

Dans les expériences précédentes, toute règle ayant une force inférieure à 0,02 était éliminée. La modification du seuil d'élimination doit alors influencer la survenue des phases transitoires observées. Le Tableau IV-10 récapitule la valeur des paramètres pour les 3 expériences réalisées en plus de l'expérience 14, dans le cadre de l'étude portant sur l'influence du seuil d'élimination.

	Seuil d'élimination				
	0,001	0,01	0,02	0,05	0,1
Taxe=0,001 ; Enchère=Remboursement=0,05	27	28	14	29	30

Tableau IV-10 : Récapitulatif des expériences réalisées en fonction du seuil d'élimination.

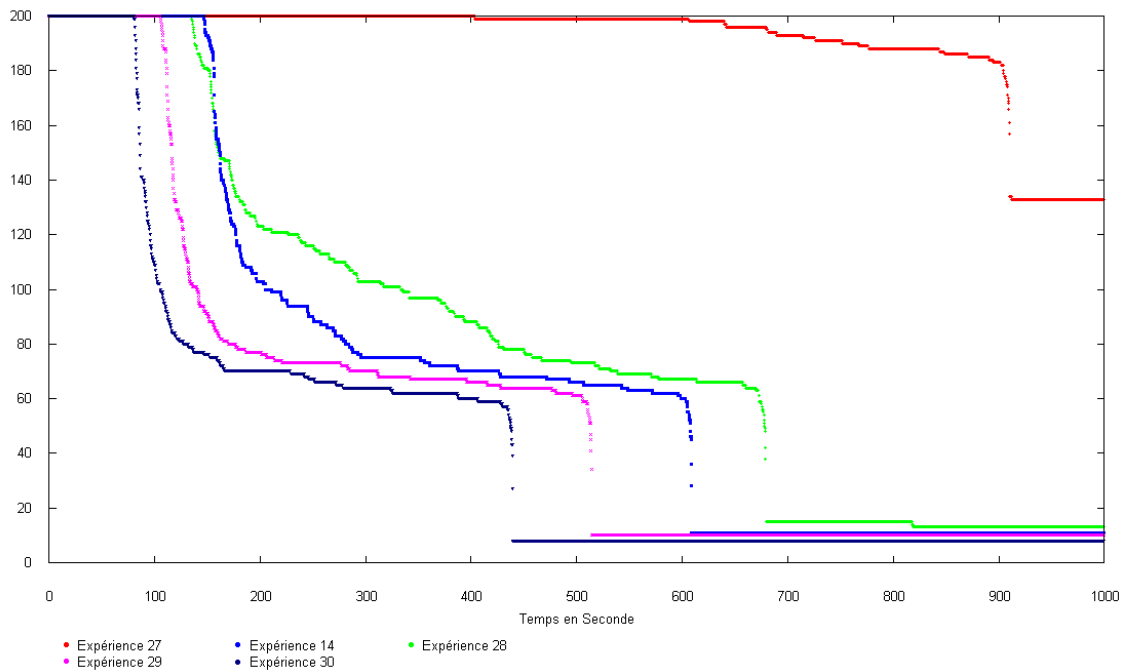


Figure IV-50 : Évolution du nombre de règles d'action dans les expériences 27, 28, 14, 29 et 30

La comparaison des évolutions démographiques (Figure IV-50) montre l'influence du seuil de l'élimination sur l'apparition de la première élimination observée dans les expériences précédentes. La deuxième transition de l'expérience 27 qui n'apparaît pas sur la Figure IV-50 subit également un décalage en fonction de la valeur du seuil, mais en moindre proportion.

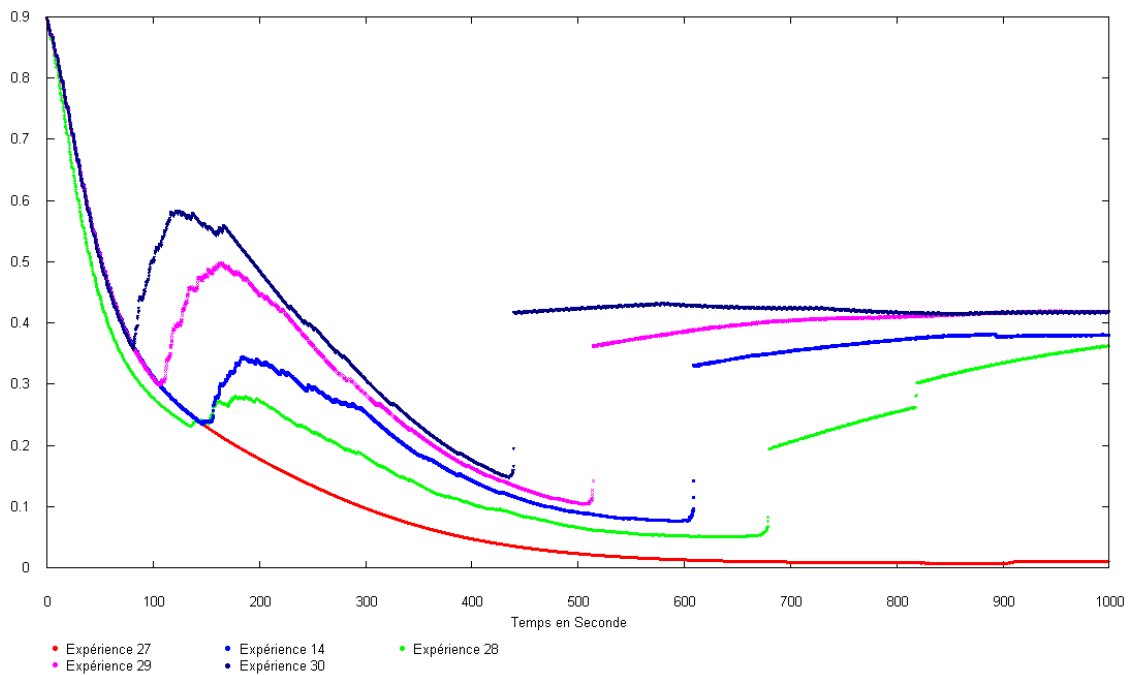


Figure IV-51 : Évolution de la force moyenne des règles sensorimotrices dans les expériences 27, 28, 14, 29 et 30

Un seuil d'élimination élevé écourte la compétition, de sorte que la trajectoire de la force des règles adaptées s'oriente plus rapidement vers l'équilibre. La comparaison des évolutions de la force moyenne montre pour les expériences 28, 14, 29 et 30 que les forces moyennes valent environ $0,396 \pm 0,014$ à l'équilibre. L'allure des courbes d'évolution des paramètres, tels que l'alpha ou le score, est similaire à l'expérience 14.

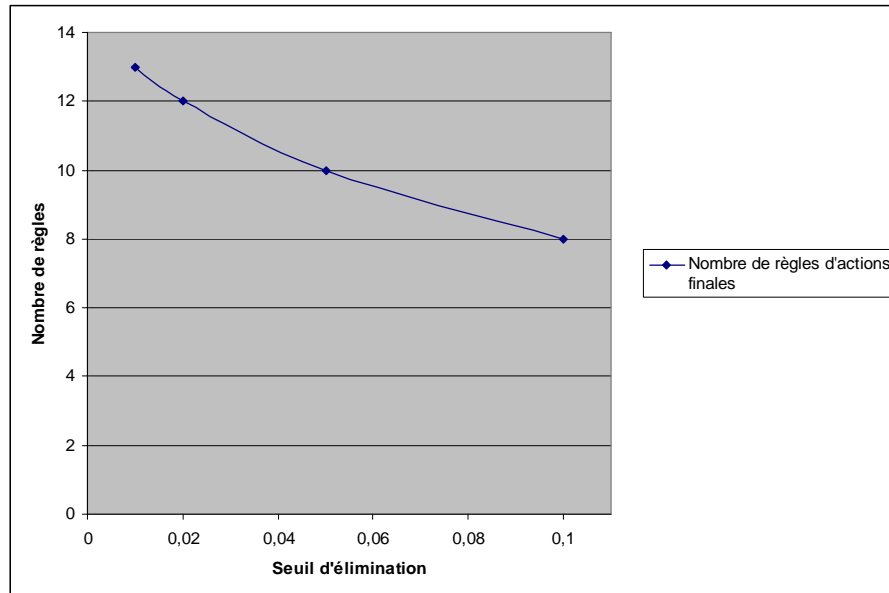


Figure IV-52 : Nombre de règles à l'équilibre en fonction du seuil d'élimination, avec les expériences 28, 14, 29 et 30

L'augmentation du seuil d'élimination est accompagnée d'une diminution du nombre de règles. La démographie de l'expérience 27 se trouve décalée de 900 s et maintient ensuite une population de 130 individus jusqu'à la 1800^{ème} s. Mais cette situation correspond à la troisième phase c'est-à-dire au plateau suivant la première transition ; par conséquent, l'expérience 27 ne peut être comparée avec les autres sur ce point (Figure IV-52). La reprise de l'expérience 27 jusqu'à 2500 s a permis de dépasser la seconde transition. Cette dernière s'étend sur 500 s et aboutit à un début d'équilibre à une population de 100 règles.

Les règles supplémentaires de l'expérience 28, comparativement à l'expérience 30, se situent préférentiellement dans la zone dense des états sensoriels (Figure IV-53). Les prémisses couvrant les données environnementales moins denses sont quasiment identiques.

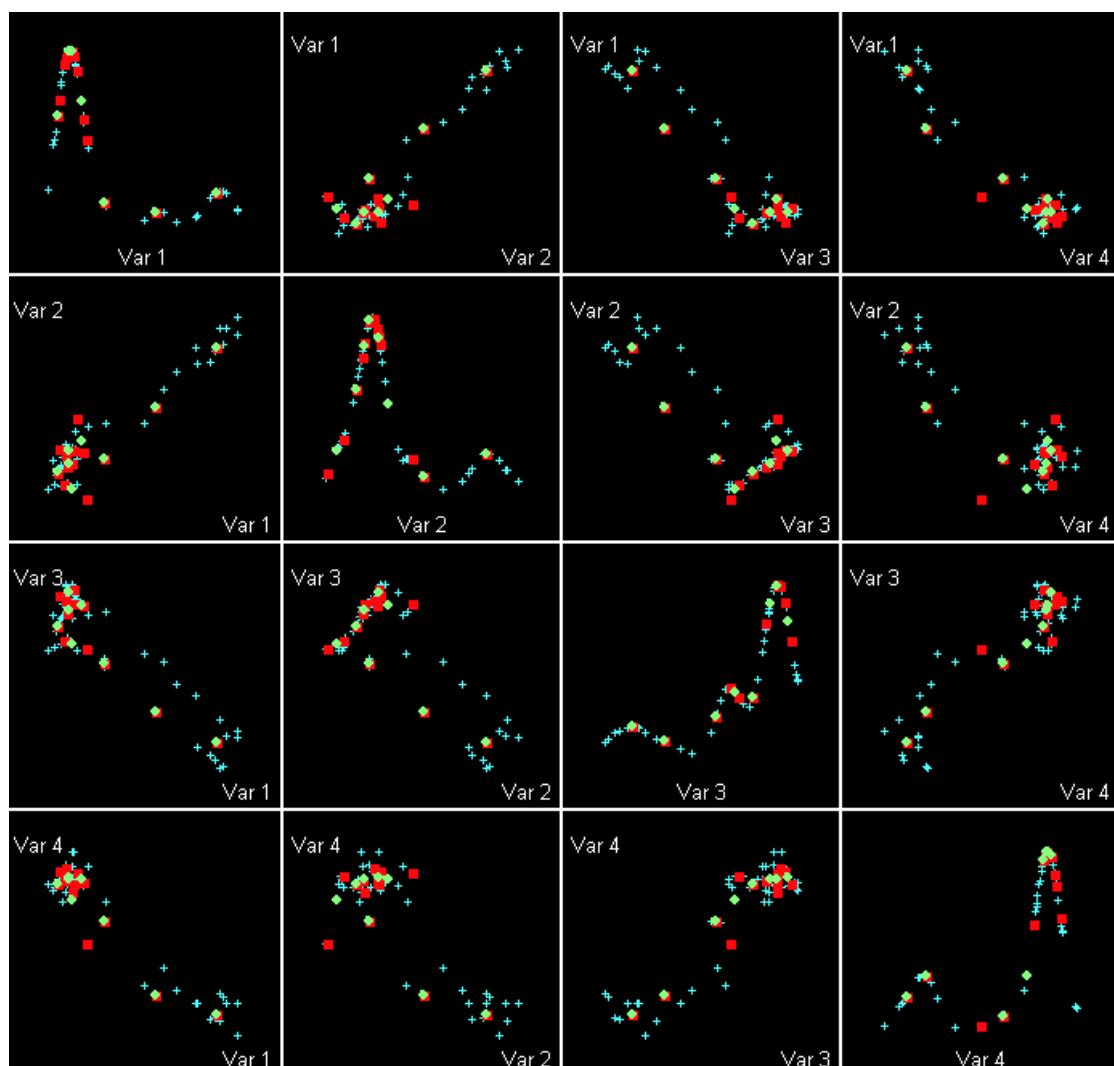


Figure IV-53 : Comparaison entre les prémisses finales de l'expérience 28 (carrés rouges), celles de l'expérience 30 (ronds verts) et de tous les états sensoriels de l'environnement imposé (croix bleus).

4.1.5. Les incertitudes initiales des prémisses

Dans toutes les expériences précédentes, les incertitudes des valeurs sensorielles des prémisses initialisées à 0,01 atteignent au final des valeurs de l'ordre de 10^{-3} et 10^{-4} . La différence entre les valeurs des incertitudes au sein d'un même message peut donc être d'un facteur 10. En revanche, l'incertitude de l'étiquette temporelle est restée de l'ordre de 10^{-4} alors qu'une amélioration est également attendue. En effet, le seuil d'oubli de la mémoire était fixé à 0,925, ce qui correspondait à la mémorisation des 8 derniers messages sensoriels et l'incertitude sur l'étiquette temporelle devait permettre de toujours sélectionner le message le plus récent et devait donc se réduire rapidement puisque la valeur du dernier message était toujours proche de 1. Mais l'incertitude ne devait pas être suffisamment restrictive pour que le poids de l'étiquette temporelle, par rapport au poids du contenu sensoriel, sensibilise les prémisses uniquement sur les messages récents.

Cette déconvenue révèle alors l'existence d'un biais dans les expériences précédentes, dans la mesure où la classification était sensori-temporelle au lieu d'être uniquement

sensorielle. Ce biais expliquerait la concentration de prémisses sur certaines régions de l'espace sensoriel, Figure IV-44. Remédier à ce problème revient soit à diminuer l'incertitude de l'étiquette temporelle, soit à augmenter l'incertitude sur la valeur sensorielle. Seule la seconde solution a été retenue puisque les règles initiales doivent tendre au départ vers le général plutôt que vers le particulier afin de couvrir potentiellement toutes les situations.

L'expérience 31 reprend tous les paramètres de l'expérience 14 en dehors de l'incertitude des valeurs sensorielles des règles initiales qui vaut dorénavant 10^{-2} , soit un facteur cent avec l'incertitude de l'étiquette temporelle. L'amélioration des résultats (Figure IV-54) confirme l'identification du biais. Plus précisément, seules 5 règles subsistent au lieu de 11 pour l'expérience 14 et leur répartition se trouve plus homogène et cohérente avec les états sensoriels imposés. Les incertitudes des prémisses sur les valeurs sensorielles vont de 10^{-3} à 10^{-6} et l'incertitude sur l'étiquette temporelle atteint en moyenne $1,6 \cdot 10^{-6}$.

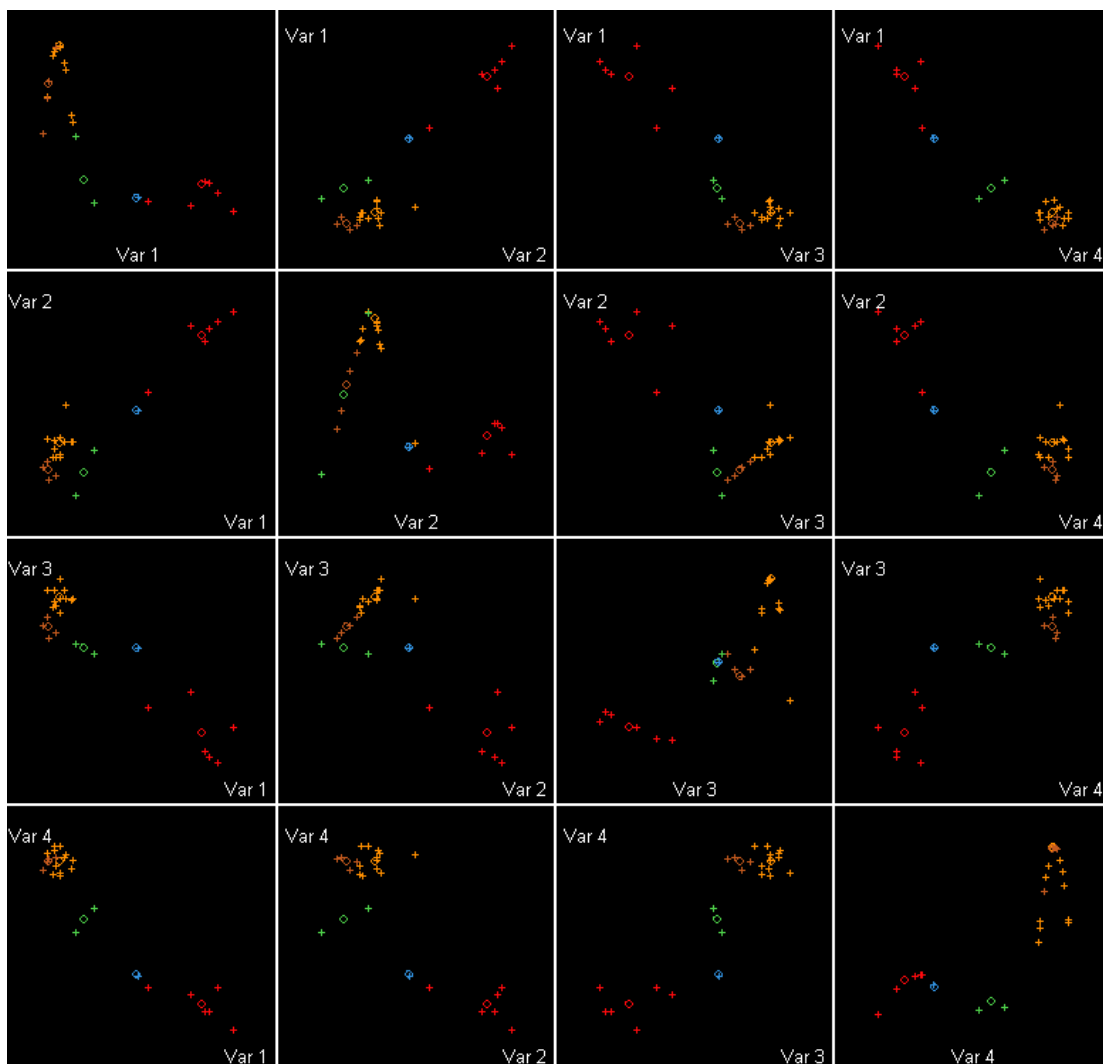


Figure IV-54 : Représentation graphique des prémisses finales qui ont été déclenchées plus d'une fois lors de la dernière séquence (cercles) et des états sensoriels de l'environnement imposé à l'origine d'une des ces prémisses lors de la dernière séquence (croix) de l'expérience 31.

L'évolution démographique de l'expérience 31 présente également les 5 phases décelées dans les expériences 11, 12, 14, 15, 17 et 18 mais avec un allongement du premier plateau de 300 s et un raccourcissement du plateau intermédiaire par rapport à l'expérience 14, Figure IV-55 enfin le dernier plateau se stabilise plus lentement.

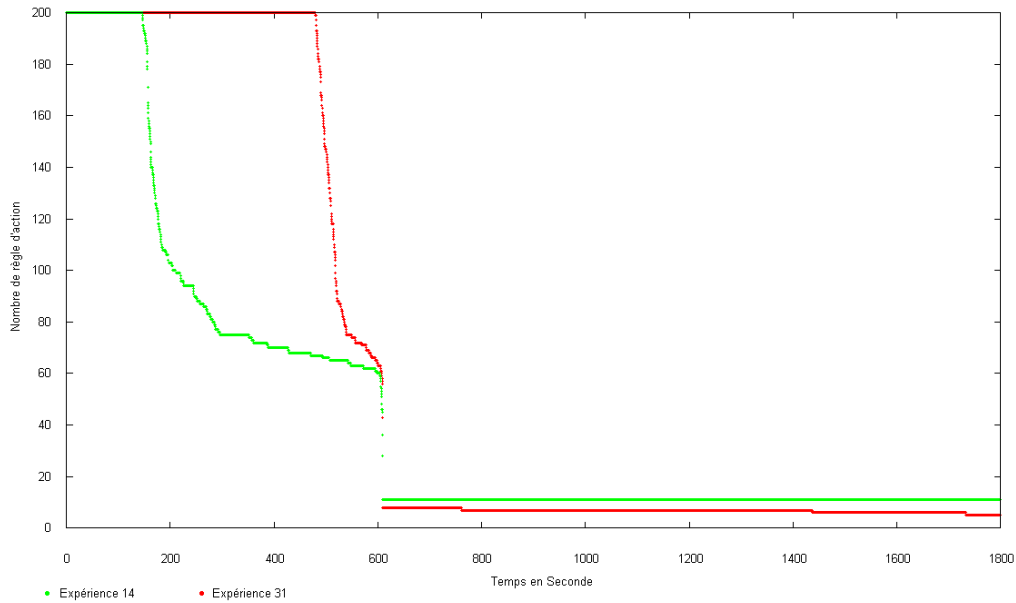


Figure IV-55 : Évolution du nombre de règles d'action dans les expériences 14 et 31.

Le nombre de règles différentes se déclenchant au début (Figure IV-56) se trouve moindre que dans les autres expériences telles que l'expérience 13 (Figure IV-36). L'importance de la dimension temporelle empêche la sélection de règles avec des états sensoriels dépassés, c'est-à-dire que l'objet de la compétition se trouve réduit tout comme l'ensemble des concurrents sérieux.

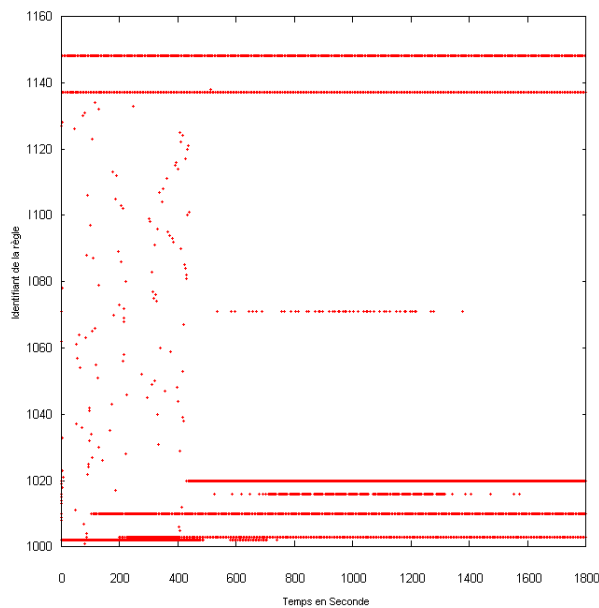


Figure IV-56 : Graphique indiquant le numéro de la règle gagnante à chaque élection de l'expérience 31.

La force moyenne à la fin de l'expérience suit quasiment à l'identique celle de l'expérience 14, avec respectivement un équilibre à 0,38 et 0,4. Cependant, les intervalles de valeur de la force des règles finales ne correspondent pas : les forces valent entre 0,2 et 0,6 pour l'expérience 31 alors qu'elles valent entre 0,25 et 0,5 pour l'expérience 14. Cet élargissement de l'intervalle n'est pas aussi significatif pour les valeurs de l'alpha : entre 0,71 et 0,98 pour l'expérience 31 et entre 0,78 et 0,98 pour l'expérience 14. Malgré l'ajustement des incertitudes qui permet de mieux orienter la compétition, les oscillations des caractéristiques de certaines règles restent présentes comme pour les autres expériences (Figure IV-57).

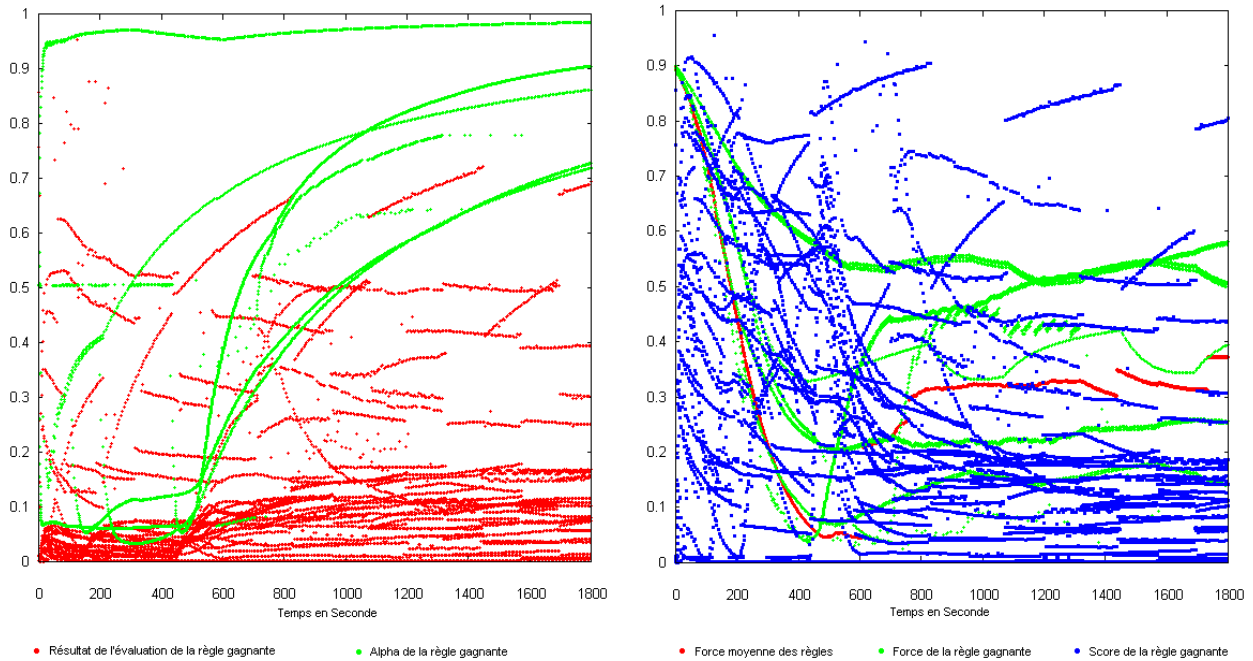


Figure IV-57 : À droite, représentation de l'alpha et du résultat de l'évaluation de la règle déclenchée à chaque sélection durant l'expérience 31, et à gauche représentation de la force et du score de la règle déclenchée ainsi que la force moyenne des règles.

4.2. La périodicité des états sensoriels imposés

Les expériences montrant l'influence de la périodicité des messages sensoriels s'appuient sur des conditions expérimentales identiques à celles de l'expérience 31. Le scénario S2 a été utilisé dans l'expérience 32 et le scénario S3 dans l'expérience 33. Leur courbe d'évolution démographique se décrit également en 5 phases (Figure IV-58). La courbe de l'expérience 31 se trouve entre celle des expériences 32 et 33. L'avancement de la première transition de la courbe démographique de l'expérience 33 s'explique par le fait que la séquence B, la plus fréquente, présente des états sensoriels relativement éloignés par rapport aux prémisses initiales. Face à ces points éloignés, toutes les règles concurrentes se valent, durcissant la compétition. Dans l'expérience 32, le décalage vers la droite de la première transition traduit une compétition plus diversifiée que celle de l'expérience 31 et en même temps plus significative que celle de l'expérience 33.

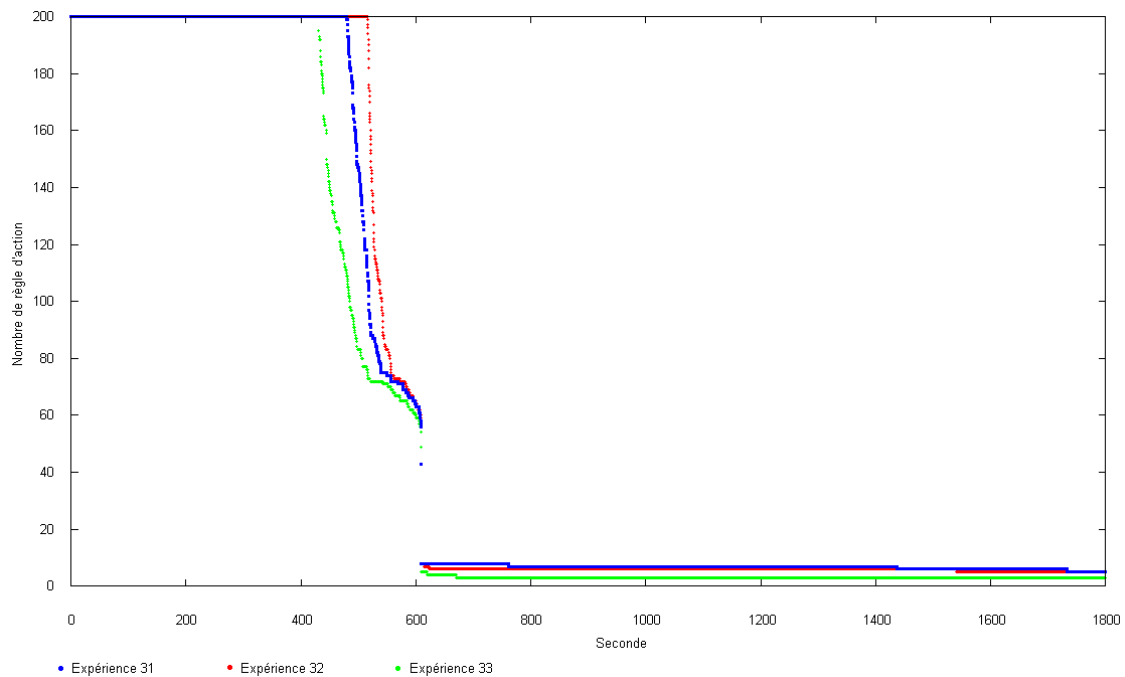


Figure IV-58 : Évolution du nombre de règles sensorimotrices dans les expériences 31, 32 et 33.

La dynamique de la force, du score et de l'alpha reste similaire entre les expériences 32 et 33. La périodicité des séquences transparait sur les courbes représentant la force et l'alpha de la règle déclenchée. Les oscillations observées ont une période de 15 s soit 0,06 Hz correspondant à la durée de trois séquences qui forment le motif répété des scénarios S2 et S3. Cela signifie que certaines règles se sont spécialisées sur un ou plusieurs messages sensoriels se trouvant uniquement sur l'une des deux séquences.

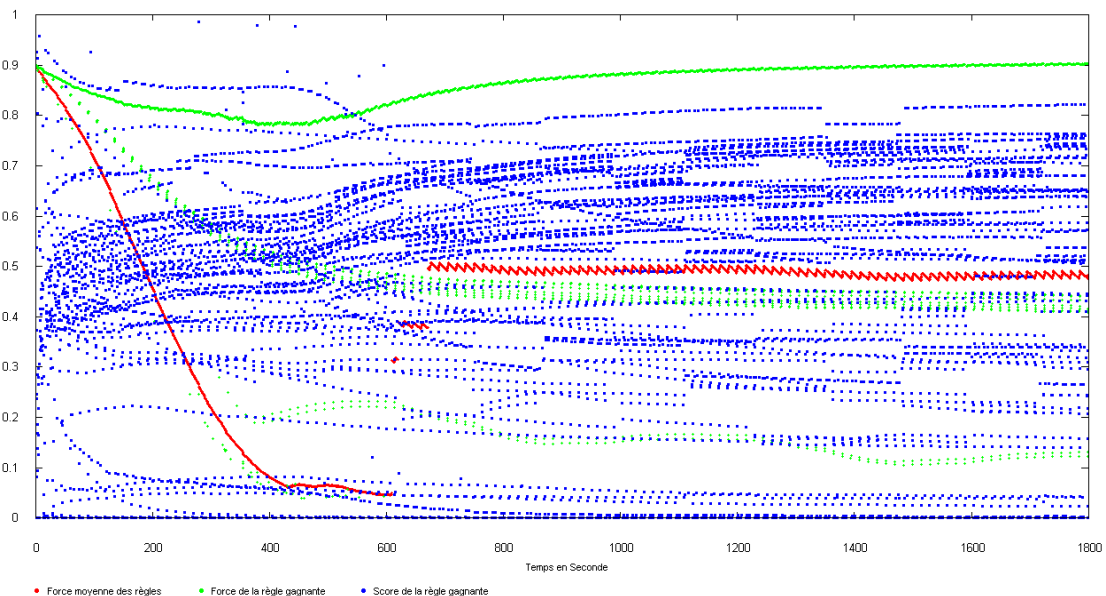


Figure IV-59 : Représentation graphique de la force et du score de la règle gagnante ainsi que la force moyenne des règles au cours du temps de l'expérience 33.

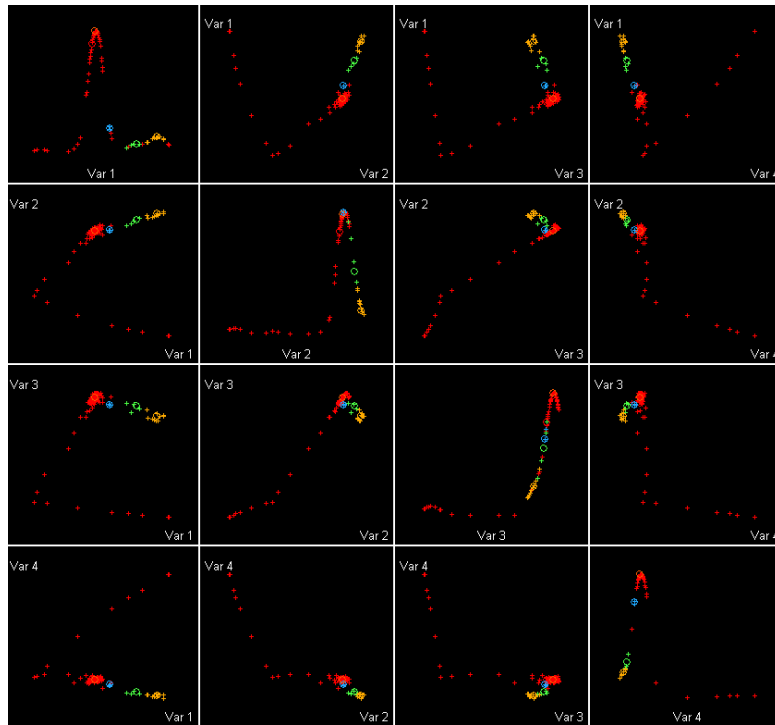


Figure IV-60 : Représentation graphique des prémisses finales qui ont été déclenchées plus d'une fois lors de la dernière séquence (cercles) et des états sensoriels de l'environnement imposé à l'origine de l'une des ces prémisses lors de la dernière séquence (croix) de l'expérience 32.

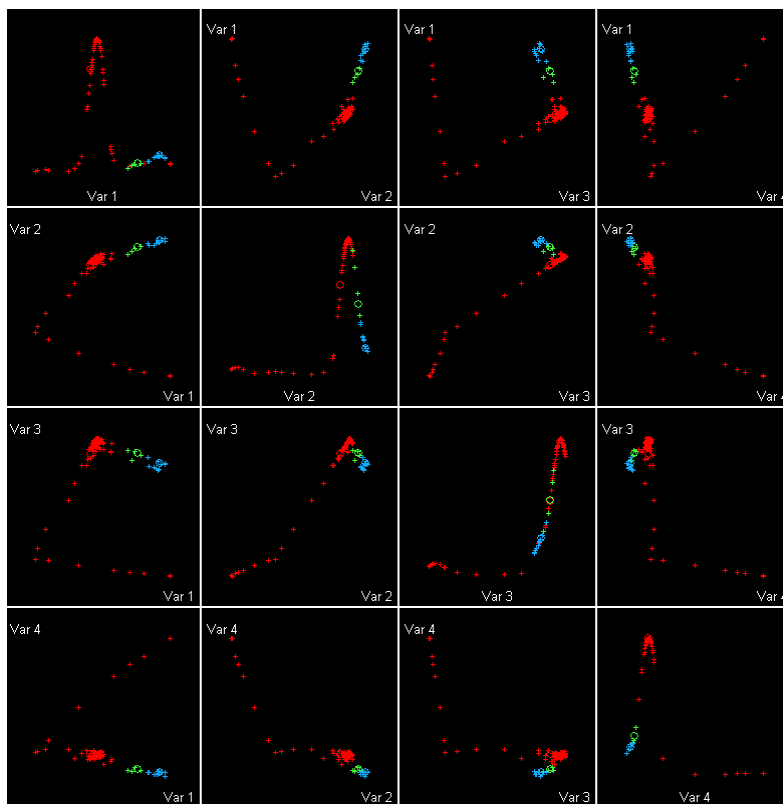


Figure IV-61 : Représentation graphique des prémisses finales qui ont été déclenchées plus d'une fois lors de la dernière séquence (cercles) et des états sensoriels de l'environnement imposé à l'origine de l'une des ces prémisses lors de la dernière séquence (croix) de l'expérience 33.

L'expérience 31 conserve 5 règles alors que les expériences 32 et 33 conservent respectivement 4 et 3 règles sur 200. Ce phénomène s'explique par la diminution de la fréquence de déclenchement des règles spécialisées sur des états sensoriels d'une seule séquence, qui entraîne la diminution de la fréquence de remboursement conduisant à une dynamique négative de la force. La légère différence entre les expériences 33 et 32 sur le nombre de règles finales (Figure IV-60 et Figure IV-61) provient de la fréquence deux fois plus élevée de la séquence A qui se trouve plus proche des règles initiales dans l'espace sensoriel. Dans les deux cas, la généralisation d'une seule règle a permis de couvrir entièrement la séquence B. La même règle se trouve à l'origine de ces règles finales. Les règles finales capturant l'extrémité de la séquence A proviennent également d'une même règle.

4.3. L'apport de la création de règles

La construction de règles au cours des expériences s'appuie sur les données se trouvant dans la mémoire événementielle (MEV). Cette dernière, réglée par un seuil, peut contenir au maximum 8 messages du même type. La MEV reçoit à la fois les messages sensoriels à une fréquence de 10 Hz et les messages moteurs provenant de la conclusion de la règle déclenchée. La construction de la prémisse sensorielle des nouvelles règles s'effectue par recopie d'un nombre aléatoire de messages sensoriels successifs de la MEV et la construction de la conclusion par la recopie des messages moteurs dont l'étiquette temporelle est postérieure à celle du dernier message sensoriel de la prémisse. Les règles de gestion associées sont alors créées à partir de ce modèle. L'introduction de règles sensorimotrices ayant en conclusion une série de messages, l'utilisation des règles de gestion devient indispensable, par conséquent l'ensemble E1-bis correspond aux règles initiales dans les trois expériences qui suivent.

Afin d'empêcher une domination immédiate des nouvelles règles sur les anciennes dont la force résulte d'une première phase très compétitive et dont l'alpha s'est ajusté, la force et l'alpha valent tous les deux initialement 0,5. Afin que la création de règles s'effectue sur l'ensemble d'états sensoriels existant, la fréquence de la création de règles ne doit pas être un modulo de la fréquence d'apparition d'une séquence. En suivant ce critère, la fréquence de création d'une nouvelle règle choisie vaut 0,13 Hz, soit environ toutes les 7,5 s. Les taux de taxe et d'enchère correspondent à ceux de l'expérience 14 soit respectivement à 0,001 et à 0,05 ; en revanche le remboursement a été augmenté à 0,08 afin d'aider les règles nouvellement créées dans une population déjà compétitive. Le Tableau IV-11 récapitule la valeur des paramètres pour les 3 expériences réalisées avec les mêmes conditions que celles de l'expérience 31.

	Scénarii		
	S1	S2	S3
Expériences	34	35	36

Tableau IV-11 : Les scénarii utilisé en fonction des expériences

Les trois évolutions démographiques (figure IV-62) conservent les cinq phases, sauf que le premier plateau devient une croissance linéaire due à la création des règles dès le démarrage du programme. Malgré la création régulière de règles, le nombre de règles se stabilise. Le taux d'élimination de règles compense celui de création. Le nombre de règles à

l'équilibre des expériences 35 et 36 est dix fois supérieur à celui de l'expérience 31. Ce surplus de règles correspond à la présence de règles créées non-pertinentes (non déclenchées) jusqu'à leur élimination. En effet, un nombre de déclenchement inférieur à 5 à la fin du programme correspond à 80% des règles. L'expérience 34 possède un nombre de règles à l'équilibre deux fois supérieur à celui de l'expérience 31. En effet, les règles créées étant semblables et essayant de s'adapter sur les motifs, la compétition via le mécanisme d'enchère pénalise ces nouvelles règles.

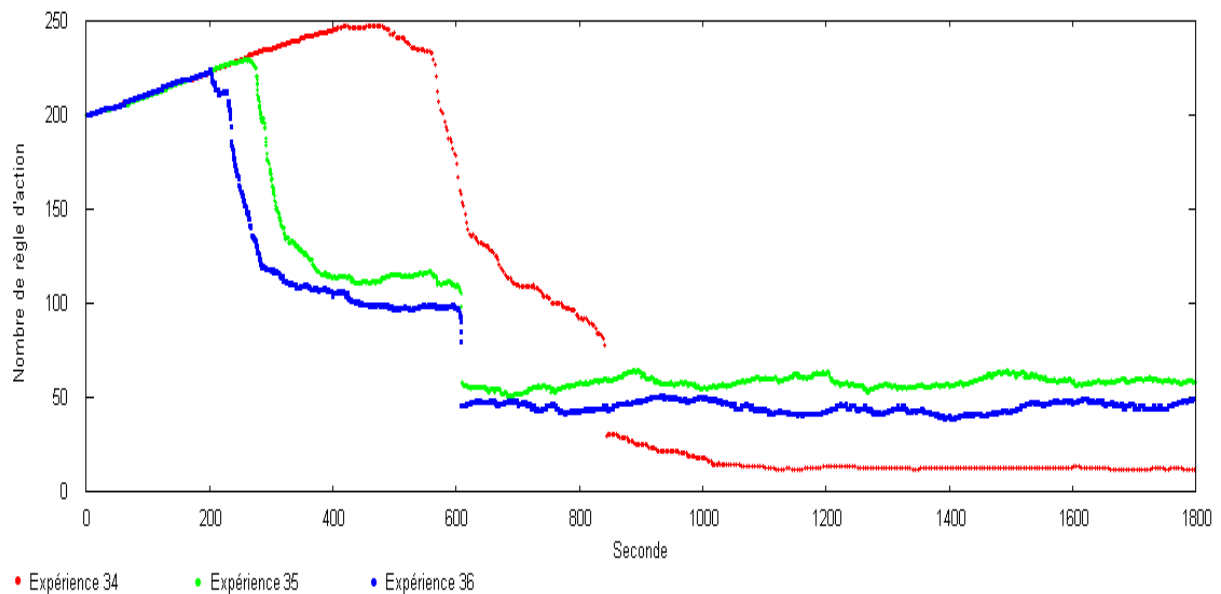


Figure IV-62 : Évolution du nombre de règles sensorimotrices dans les expériences 34, 35 et 36.

A - Résultat avec le scénario simple S1

Les règles créées parviennent difficilement à s'imposer vis-à-vis des règles initiales ou des règles créées qui se sont adaptées dès le début de l'expérience (Figure IV-63).

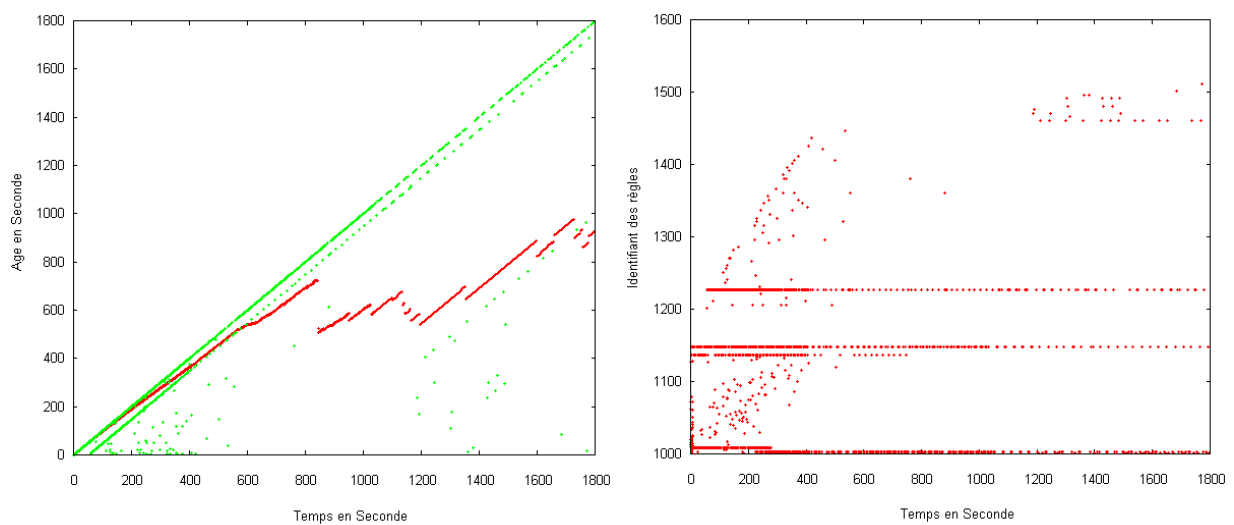


Figure IV-63 : Graphiques représentant à droite l'âge de la règle gagnante en vert et l'âge moyen des règles en rouge ; à gauche le numéro de la règle gagnante de l'expérience 34.

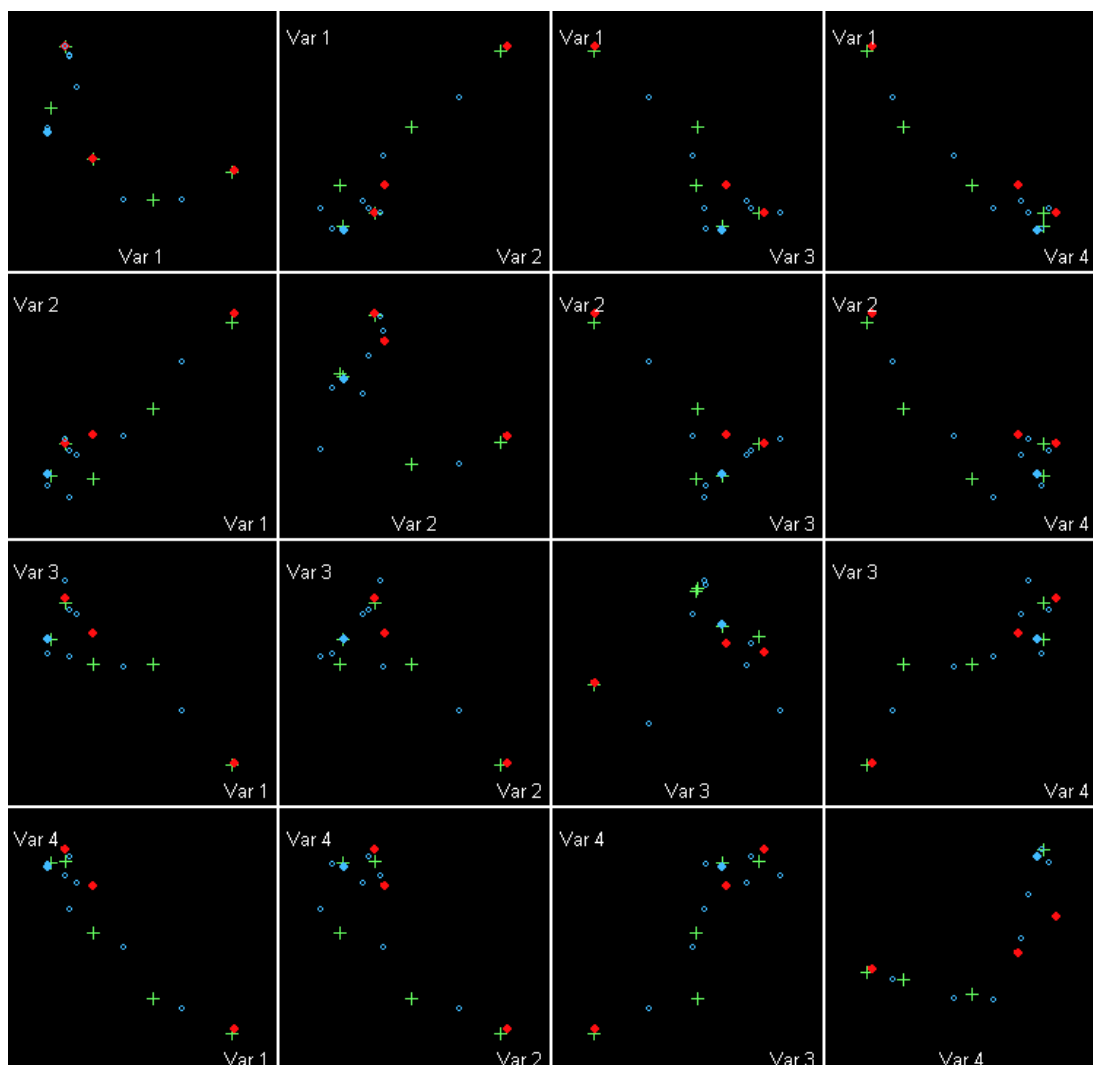


Figure IV-64 : Représentation graphique des prémisses finales de l'expérience 31 en vert et celles des prémisses finales de l'expérience 33 où les ronds rouges illustrent les prémisses de règles issues de la population initiale, les ronds bleus illustrent les prémisses des règles créées déclenchées plus de 20 fois et enfin les cercles bleus illustrent les prémisses des règles créées déclenchées moins de 20 fois.

Les prémisses des règles restantes contiennent majoritairement un ou deux messages, toutefois une prémisses de quatre messages à été observée avec un nombre de déclenchements proche de la centaine. Trois points sur cinq correspondent aux prémisses finales de l'expérience 31. Toutefois, dans la Figure IV-64, seul le premier message de la prémisses sensorielle représente les règles finales en supposant que celui-ci soit suffisamment représentatif des autres. La création intensive de règles n'a pas réussie à déloger celles qui étaient parvenues à s'adapter rapidement.

B - Résultat avec le scénario complexe S2

Le scénario S2 de l'expérience 34 offre un espace sensoriel éloigné des règles initiales. Les règles créées à partir de ces états sensoriels ne se trouvent pas en concurrence avec les règles initiales. Cela se traduit par un nombre plus important de règles créées à s'être déclenchées comparativement à l'expérience 33. Par ailleurs, la première phase de la compétition entre ces nouvelles règles dure également plus longtemps puisqu'elles se

déclenchent au début de leur apparition entre 1 à 15 fois, même pour les règles qui ne parviennent pas à se maintenir (Figure IV-65).

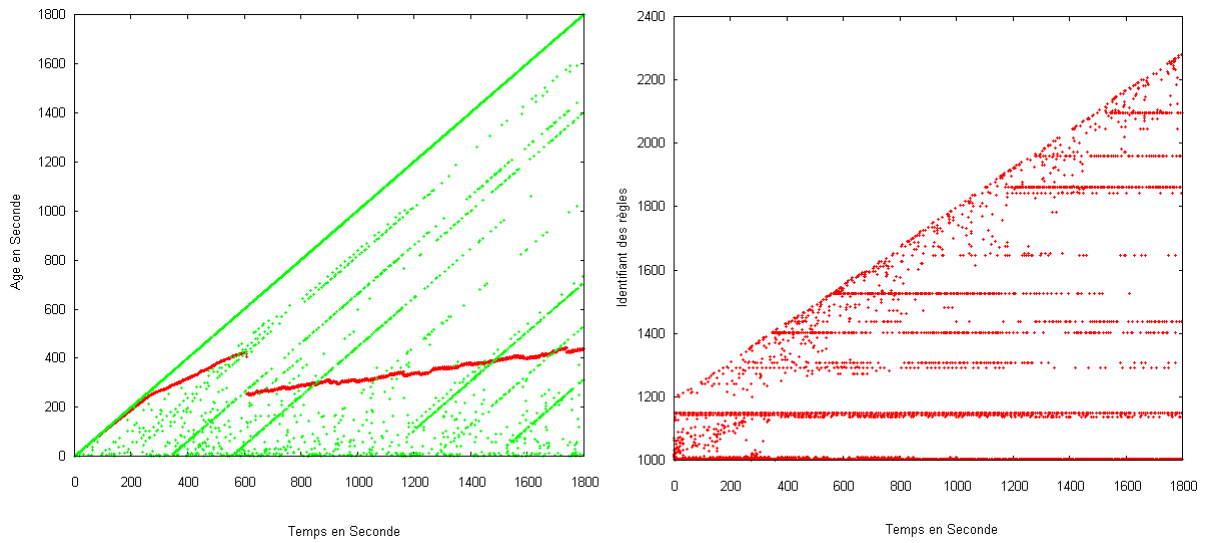


Figure IV-65 : Graphiques représentant à droite à chaque pas de temps l'âge de la règle gagnante en vert et l'âge moyen des règles en rouge, à gauche le numéro de la règle gagnante de l'expérience 34.

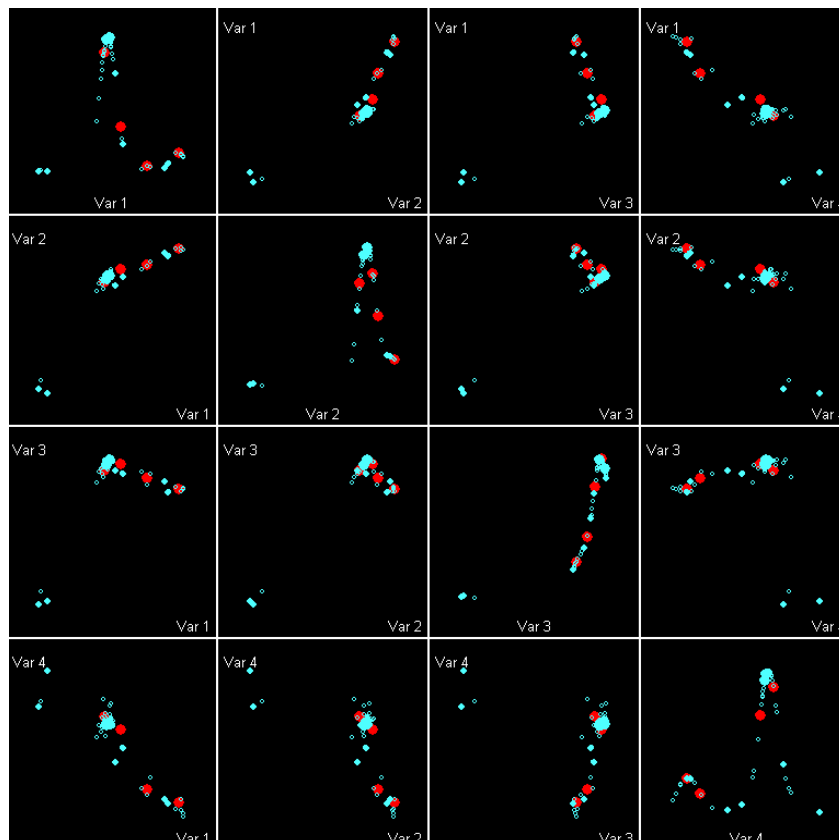


Figure IV-66 : Représentation graphique des prémisses finales de l'expérience 34, où les ronds rouges illustrent les prémisses de règles issues de la population initiales, les ronds bleus illustrent les prémisses des règles créées déclenchées plus de 20 fois et enfin les cercles bleus illustrent les prémisses des règles créées déclenchées moins de 20 fois. Les ronds les plus grands désignent les règles déclenchées plus de 100 fois.

Parmi les règles finales, 85% d'entre elles se sont déclenchées moins de deux fois (Figure IV-66). La séquence A étant deux fois plus présente, le nombre de règles créées se concentre sur celle-ci, sans toutefois s'imposer. Cependant, les prémisses des règles initiales dominent l'espace sensoriel contenant la séquence A. En revanche, les règles créées se sont imposées sur l'espace sensoriel contenant la séquence B.

C - Résultat avec le scénario complexe S3

Pour les mêmes raisons que celles de l'expérience 34, les données de l'expérience 35 (Figure IV-67) révèlent une augmentation du nombre de déclenchements de règles nouvellement créées par rapport à l'expérience 33.

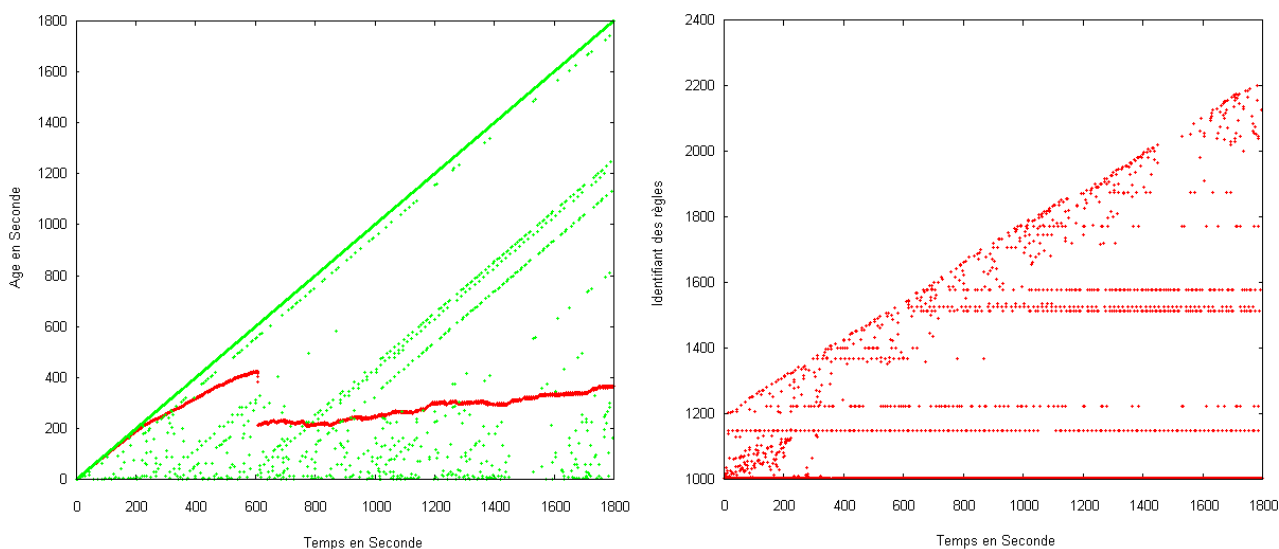


Figure IV-67 : Graphiques représentant à droite l'âge de la règle gagnante en vert et l'âge moyen des règles en rouge ; à gauche l'identifiant de la règle gagnante au cours de l'expérience 35.

La comparaison des prémisses finales entre les expériences 34 et 35 illustre l'importance de la répétition d'une séquence par rapport à une autre dans la stabilisation des prémisses. La Figure IV-68 montre une concentration de nouvelles règles établies (c'est-à-dire dont le nombre de déclenchement dépasse la centaine) deux fois supérieure à l'expérience 34 vers l'espace sensoriel contenant la séquence B. Réciproquement, les règles initiales moins sollicitées ont laissé plus facilement de la place aux nouvelles règles.

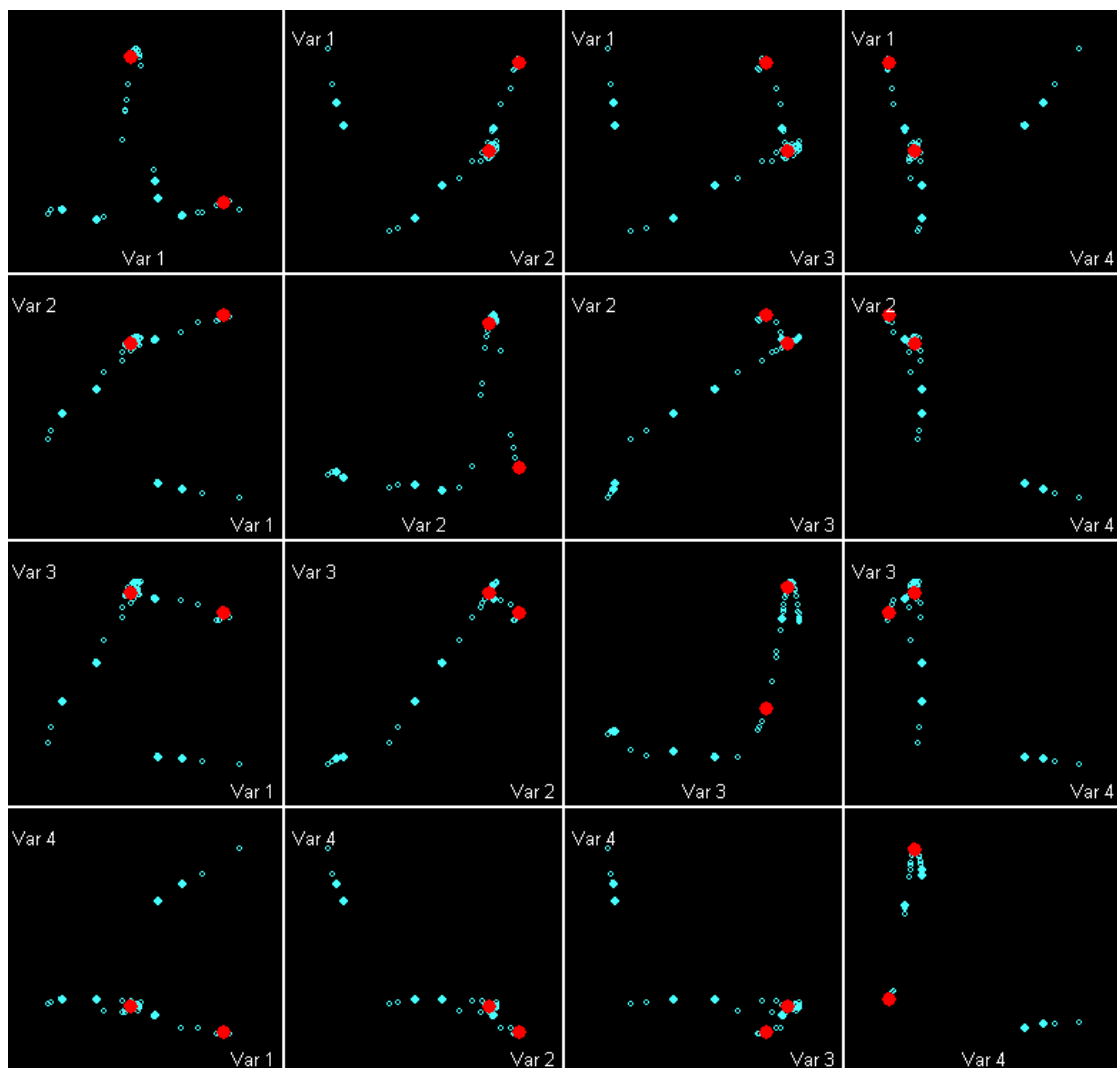


Figure IV-68 : Représentation graphique des prémisses finales de l'expérience 35, où les ronds rouges illustrent les prémisses de règles issues de la population initiale, les ronds bleus illustrent les prémisses des règles créées déclenchées plus de 20 fois et enfin les cercles bleus illustrent les prémisses des règles créées déclenchées moins de 20 fois. Les ronds les plus grands désignent les règles déclenchées plus de 100 fois.

4.4. Introduction de la prédiction

L'étude du mécanisme d'auto-rétribution s'effectue au travers de l'étude d'une boucle de récompense en fonction de la capacité à prédire des états sensoriels futurs en fonction de l'état de la mémoire événementielle. Cette boucle prédictive reprend le modèle prédictif proposé dans le chapitre précédent. Les taux de taxe, d'enchère et de remboursement valent respectivement 1.10^{-3} , 5.10^{-2} et 8.10^{-2} . Les valeurs initiales de la force et de l'alpha des règles créées, tous types confondus, valent 0,5.

4.4.1. Vérification de l'ordonnement

La modélisation de la boucle prédictive se constitue à partir de trois règles : la règle de rétribution (R), la règle de Prédiction Sensorielle Prévues (PSP) et la règle de Gestion

Prévue de Prédiction Sensorielle (GPPS). La première règle contrôle l'adéquation entre la prédiction et les messages sensoriels, la seconde envoie un message prédictif suite à la présence d'un certain motif, enfin la dernière, en même temps que la seconde, envoie un message de gestion de prédiction qui autorisera le déclenchement de la première et interdira celui des autres règles prédictives et de gestion de la prédiction. La création de ce triplet de règles se trouve associée à la création des règles sensorimotrices au cours d'une expérience. Ainsi, l'état sensoriel succédant les messages moteurs sélectionnés constitue le message prédictif.

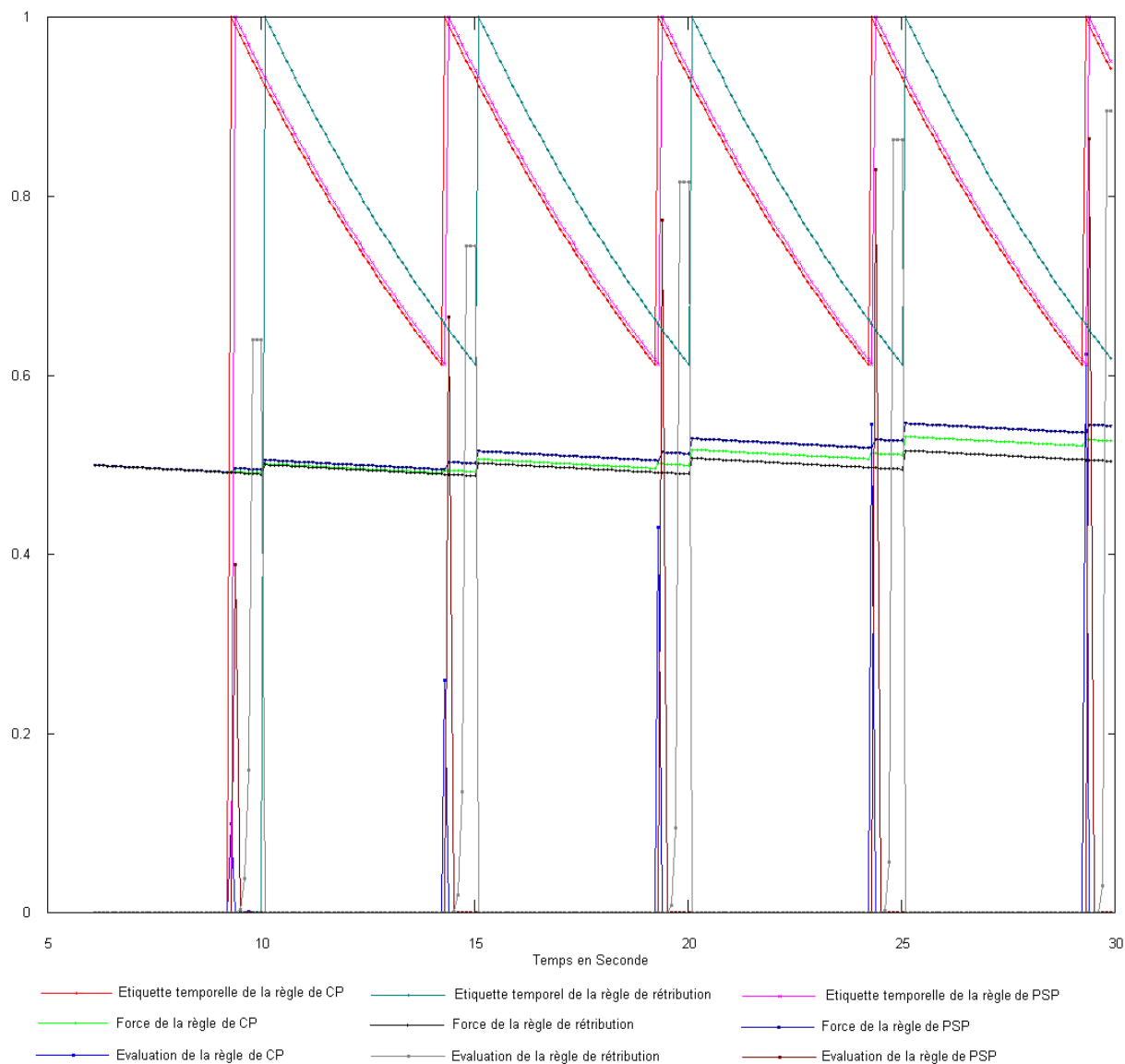


Figure IV-69 : Représentation graphique de l'évolution des caractéristiques des règles d'action, d'anticipation et de rétribution.

L'ordonnancement des règles est identique à celui décrit dans le chapitre précédent, la règle sensorimotrice se déclenche, ce qui a pour effet de mettre son étiquette temporelle à 1 et de mettre son message moteur prévu dans la mémoire événementielle. La règle de gestion se déclenche également en même temps, empêchant le déclenchement d'autres règles sensorimotrices ou de gestion. Le message moteur avec les messages sensoriels

précédents provoque le déclenchement au pas de temps suivant de la règle prédictive et de sa règle de gestion. Le message de gestion de la prédiction entre dans la MEV, la décroissance de son étiquette temporelle temporisera le signalement de la fin de l'évaluation du message de prédiction par la règle de rétribution. En effet, lorsque le message de gestion de la prédiction disparaît de la mémoire événementielle, l'évaluation finale de la règle de rétribution s'opère. Le mécanisme d'évaluation d'une attente fonctionne tant que le message GPPS se trouve dans la MEV. L'évaluation finale entraîne le déclenchement de la règle et du message de rétribution dont la valeur du contenu se calcule en fonction de l'évaluation finale. Les messages de la MEV permettent alors l'évaluation de la règle de rétribution sans les mécanismes d'enchère et de remboursement.

Afin de vérifier l'enchaînement de ces règles, un premier essai a été effectué avec le scénario S1 au cours duquel la création de règles se limite à une seule règle sensorimotrice avec sa règle de gestion associée et son triplet de règles lié à la prédiction des conséquences sensorielles suite à son déclenchement. La Figure IV-69 illustre uniquement l'ordonnancement de règles sensorimotrices, de la règle prédictive associée et de la règle de rétribution, les règles de gestion se déclenchant respectivement en même temps que la règle sensorimotrice et la règle prédictive.

La courbe grise représente l'évaluation maximale calculée depuis le début de l'anticipation. La valeur maximale étant atteinte un peu avant la fin de l'évaluation finale, la courbe grise dessine un plateau jusqu'à celle-ci. L'évolution de la force entre les trois règles se révèle différente, traduisant la dissymétrie de la compétition entre les deux types de règles due au surplus de règles sensorimotrices provenant des règles initiales. Pour la règle de rétribution, il n'y a pas de compétition donc pas de remboursement susceptible d'augmenter sa force. Le cas présenté reste simple mais de multiples cas de figures ont pu être observés tels que l'oscillation d'une règle sensorimotrice autour de deux états sensoriels différents déstabilisant la règle prédictive.

4.4.2. Répartition des règles prédictives

Réalisée dans les mêmes conditions que l'expérience 35, l'expérience 36 crée de nouvelles règles sensorimotrices avec les règles prédictives. Avec un taux de rétribution fixé à 1, la rétribution entraîne une rapide augmentation de la force des règles adaptées. Plus particulièrement, 3 règles atteignent la valeur maximale, soit une force de 1. Comparativement avec l'expérience 35, les règles se maintiennent plus facilement (Figure IV-70), ce qui signifie que la rétribution permet de favoriser la conservation de certaines règles selon des critères décrits par les règles elles-mêmes.

A l'équilibre, le nombre de règles sensorimotrices s'étant déclenchées plus de 20 fois est supérieur à celui observé dans l'expérience 35. Les prémisses finales dans les deux expériences sont quasiment identiques comme le confirme la Figure IV-71. Les règles supplémentaires apparaissent dans les zones où la concentration des états sensoriels est la plus forte. Dans cet essai, le mécanisme de prédiction et de rétribution aide au maintien de certaines règles mais n'influence pas directement la valeur des prémisses finales.

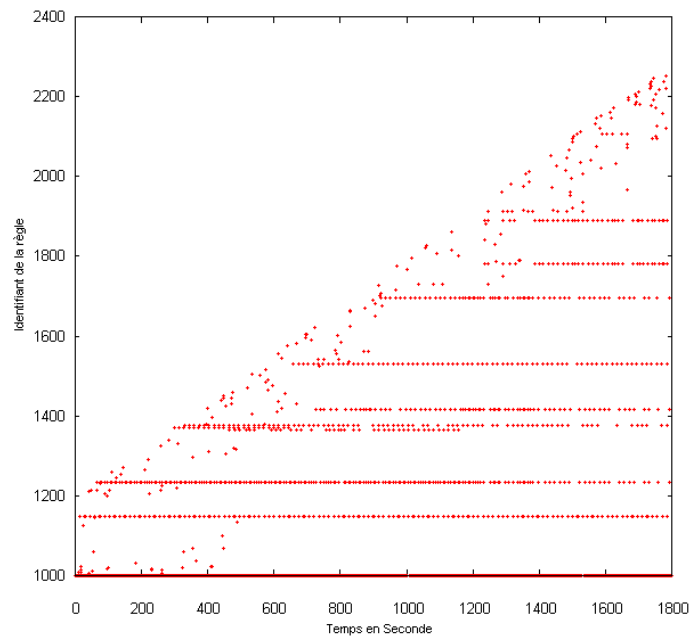


Figure IV-70 : Graphique indiquant le numéro de la règle gagnante à chaque élection de l'expérience 36.

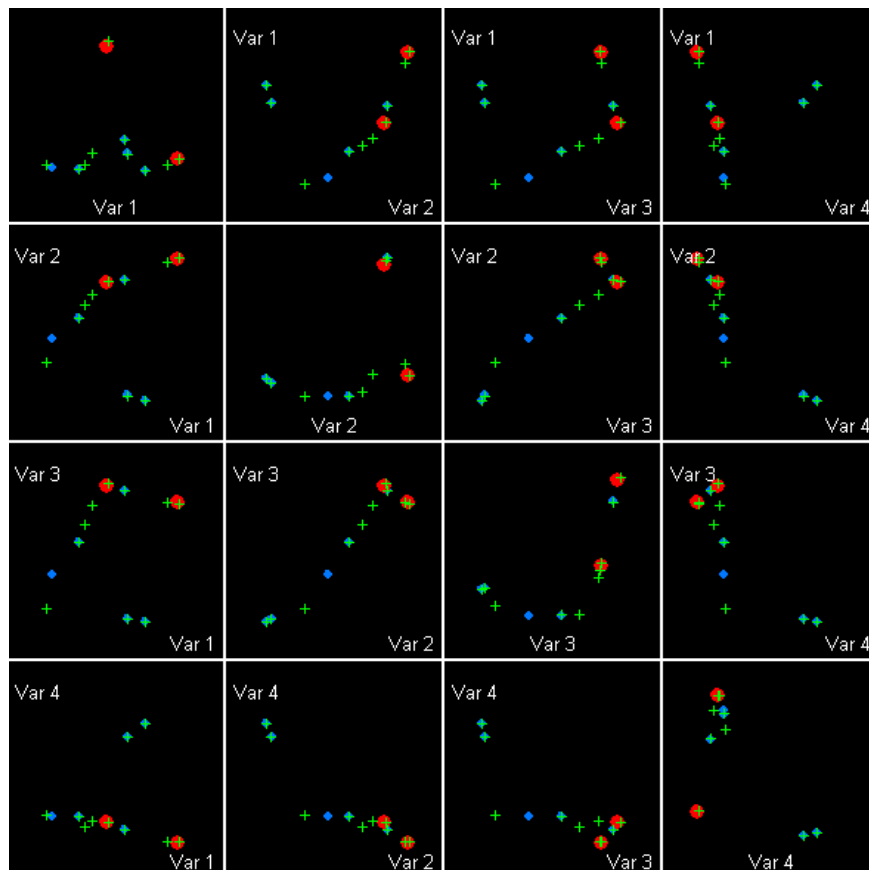


Figure IV-71 : Représentation graphique des prémisses finales de l'expérience 35 (où les ronds rouges illustrent les prémisses de règles issues de la population initiale, les ronds bleus illustrent les prémisses des règles créées déclenchées plus de 20 fois) et de l'expérience 36 (les croix vertes illustrent les prémisses des règles déclenchées plus de 20 fois).

Le nombre de déclenchements pour chaque règle est enregistré afin de pouvoir vérifier que chaque règle sensorimotrice possède un nombre de déclenchements identique à sa règle de gestion ou que chaque règle prédictive a un nombre de déclenchements identique à sa règle de gestion de la prédiction et à sa règle de rétribution associée. En revanche, les règles prédictives obtiennent un nombre de déclenchements inférieur ou égal à celui des règles sensorimotrices créées au cours de l'exécution. Cette situation provient du fait qu'une règle sensorimotrice peut se déclencher avant la fin de l'évaluation d'une prédiction et ainsi sans déclencher les règles liées à la prédiction de ses conséquences.

Mais la différence de la force des règles ne réside pas uniquement dans cette différence du nombre de déclenchements. Elle provient également du mode de calcul de l'évaluation qui prend en compte le score des autres règles. En effet, l'évaluation des règles sensorimotrices créées se fonde aussi sur les règles sensorimotrices initiales, contrairement aux règles prédictives qui dépendent uniquement des règles créées au fur et à mesure de l'essai. Plus particulièrement, la différence d'évolution de la force entre une règle prédictive et sa règle de rétribution réside dans la différence du mode de calcul de l'évaluation. L'évaluation de la seconde règle porte sur une moyenne de plusieurs évaluations sans en subir l'enchère et sans bénéficier du remboursement. Ainsi avec les paramètres choisis, la règle de rétribution se trouve défavorisée par rapport à la règle prédictive associée.

L'évolution de la force des règles apparaît différente, ce qui conduit à ce qu'une règle prédictive ou règle de rétribution puisse être éliminée de manière asynchrone. Par ailleurs, la différence entre la valeur d'évaluation des règles ayant des prémisses sensorielles identiques entraîne qu'elles ne le soient plus au bout d'un certain nombre d'ajustement des prémisses. Cela peut se rectifier par l'utilisation de messages intermédiaires pour déclencher l'action et le motif sensoriel. Ainsi, l'évaluation des règles d'action, d'anticipation et de rétribution serait binaire, assurant ainsi une dynamique plus proche et surtout faisant référence au même état sensoriel. Cependant, la différence d'évolution de la force des règles devient problématique uniquement si les dynamiques font que la règle de rétribution s'élimine avant la règle prédictive, elle-même s'éliminant avant la règle sensorimotrice. Cet ordonnancement correspond à la pire des situations. En effet, une règle sensorimotrice créée, si elle demeure, doit être accompagnée par des règles tentant de prédire ses conséquences ; en revanche, si cette règle sensorimotrice venait à disparaître les règles liées à la prédiction ne pourront plus se déclencher et s'élimineront à cause de la taxe.

5. Étude du système dans un environnement réel

Dans les expériences précédentes, les états sensoriels étaient imposés et la conclusion de la règle sensorimotrice était ignorée et n'entraînait alors aucune incidence sur la succession des états sensoriels. En revanche, en situation réelle, les observations dépendent de l'histoire du robot, de ses actions passées. Le couplage sensorimoteur devient alors une construction dynamique. Les règles initiales correspondent à l'ensemble E3 (600 règles sensorimotrices) auquel il a été ajouté leurs règles de gestion, soit au total 1200 règles constituent le système comportemental minimal du robot. Après plusieurs essais, les paramètres choisis permettent de retrouver une dynamique similaire à celles déjà rencontrées dans la partie consacrée aux essais avec simulation. Les taux de taxe, d'enchère et de remboursement sont respectivement 1.10^{-4} , 8.10^{-2} et $1,6.10^{-1}$. Le seuil d'élimination de la force reste de 2.10^{-2} . La fréquence de sélection a été doublée par rapport à la fréquence d'acquisition des états sensoriels, soit 20 Hz. Ainsi, chaque nouvelle observation peut déclencher une règle sensorimotrice malgré la présence des règles de gestion dont le temps

de latence de leur message soit au minimum. Pour les prémisses des règles initiales, l'incertitude de leurs messages sensoriels et de leur étiquette temporelle reste inchangée par rapport aux dernières expériences en environnement imposé, soit respectivement 1.10^{-2} et 2.10^{-4} . Toutes les expériences en environnement réel ont débuté avec le robot Khepera placé au centre de l'enceinte décrite précédemment (Figure IV-11) et durent 30 min soit 18000 déclenchements effectifs de règles sensorimotrices.

Les expériences en environnement réel reviennent sur quatre thèmes abordés dans l'étude expérimentale en environnement imposé. Ainsi, la première série d'expériences porte sur la capacité d'auto-organisation de la base de règles initiales dans un environnement simple, puis sur la capacité d'adaptation face à une modification dans l'environnement. Plus précisément, l'expérience 37 représente l'évolution du robot dans l'enceinte vide (sans le cylindre ni le prisme). Afin de s'assurer leur convergence, l'expérience 38 consiste en la continuation de l'expérience 37 où, plus exactement, les règles initiales correspondent aux règles finales de l'expérience 37. L'expérience 39 débute avec les règles finales de l'expérience 38 en présence d'un cylindre et d'un prisme (Figure IV-11).

La deuxième série d'expériences souhaite déceler l'incidence de l'environnement sur l'auto-organisation des règles. Afin de comparer avec l'expérience 38, l'expérience 40 représente l'évolution du robot dans l'enceinte en présence d'un cylindre et d'un prisme avec les règles initiales. L'expérience 41 reconduit l'expérience 40, c'est-à-dire que les règles finales de l'une sont les règles initiales de l'autre. L'expérience 42 correspond à la continuation de l'expérience 41. Enfin, un troisième type d'environnement réel participe à l'étude de l'incidence de l'environnement : l'expérience 43 représente l'évolution du robot dans un environnement complexe et bruité, en débutant avec les règles initiales originelles.

La troisième série d'expériences se penche sur la création de règles au cours de l'évolution du robot. Les règles créées possèdent un seul message sensoriel dans leur prémisses et un seul message de commande dans leur conclusion. Les expériences s'inscrivent dans la continuité de l'expérience 38. L'expérience 44, qui se déroule dans un environnement complexe, reprend pour règles initiales les règles finales de l'expérience 38. Afin d'observer la perturbation dans le maintien de ces règles face à un bruit, l'expérience 45 s'effectue avec une enceinte vide mais sous l'éclairage d'une lampe à incandescence. Les règles initiales de cette expérience correspondent aux règles finales de l'expérience 38.

La quatrième et dernière série d'expérience vise à étudier l'effet des règles prédictives concernant la conservation des règles ainsi que l'identification de la structure temporelle. Contrairement aux expériences précédentes en environnement imposé, la taille des prémisses et des conclusions des règles créées a été fixé à quatre. L'anticipation porte alors sur environ une seconde. L'expérience 46 correspondante a été réalisée dans un environnement complexe avec pour règles initiales celles résultantes de l'expérience 42.

5.1. Étude de la stabilité en environnement simple

5.1.1. Les états sensoriels observés

Le premier des six points de vue pour analyser les résultats de la première série d'expériences consiste à comparer les environnements vécus du robot selon qu'il soit contrôlé par l'algorithme Braitenberg ou par le système à base de règles. La Figure IV-72 montre que le système de règles explore davantage l'espace sensoriel par rapport aux données recueillies avec l'algorithme Braitenberg dans un environnement simple. Cette

différence s'explique par le fait que les règles initiales proviennent d'un enregistrement dans lequel la position initiale se trouvait accolée à un obstacle. L'introduction de ces règles a eu pour incidence qu'en situation réelle le robot s'approche davantage d'un obstacle avant de réagir. Le choix d'incorporer ce type de règles se justifie par la volonté de couvrir un nombre plus important de situations et d'observer ensuite sa réduction ou son utilisation au cours de l'évolution du système. Le déclenchement de certaines règles uniquement à une intensité sensorielle maximale se traduit, dans la Figure IV-72, par les points bleus dessinant les arêtes projetées de l'hypercube.

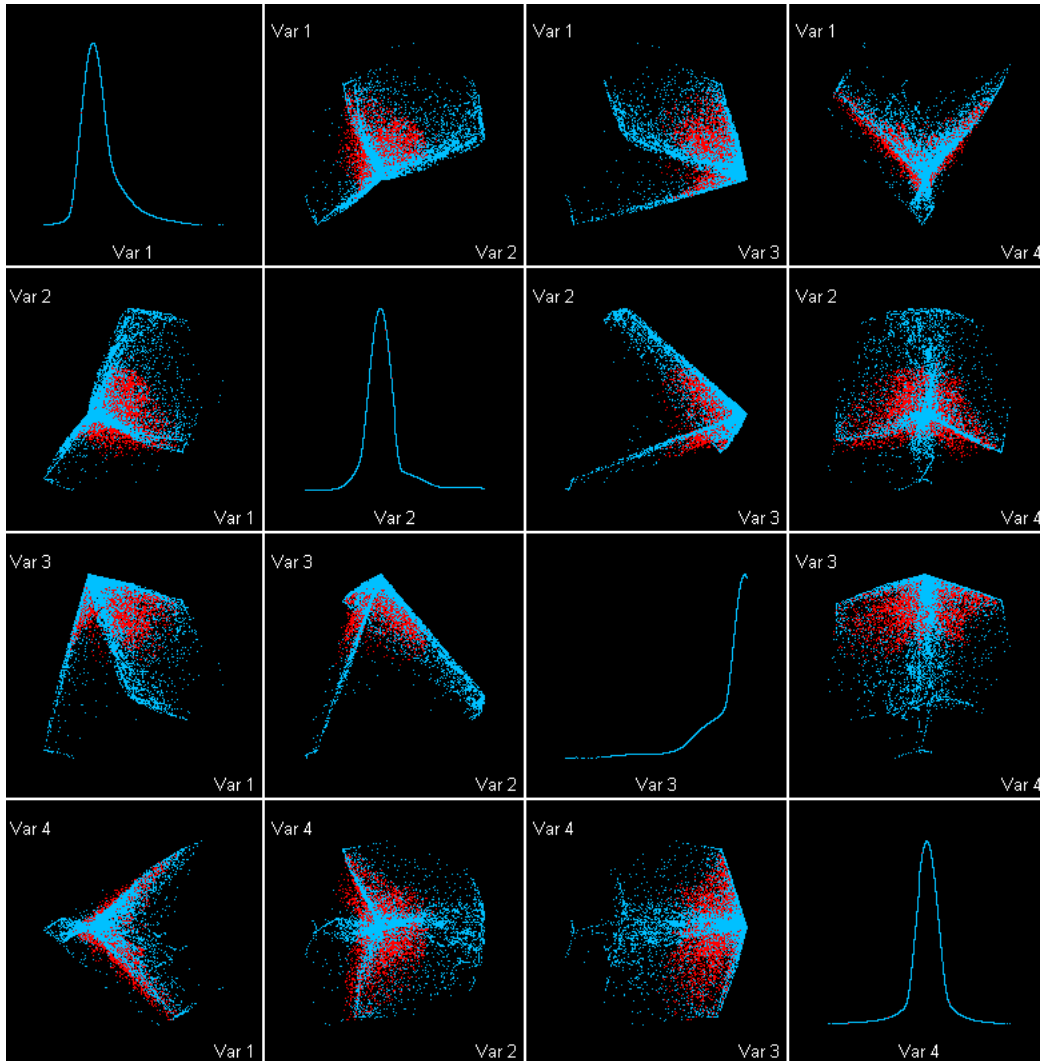


Figure IV-72 : Représentation graphique des états sensoriels observés au cours de l'expérience 37 (points bleus) et les états sensoriels provenant de l'échantillonnage avec l'évitement d'obstacle dans l'environnement simple (points rouges).

Les observations des expériences 37 et 38 couvrent l'espace des états sensoriels de façon similaire (Figure IV-43). Les points provenant de l'expérience 39, réalisée avec un environnement complexe, accèdent légèrement davantage aux extrêmes de l'espace sensoriel (Figure IV-74). L'environnement complexe ainsi exploré ne semble pas offrir de configurations sensorielles caractéristiques différentes selon les quatre premières composantes extraites de l'ACP.

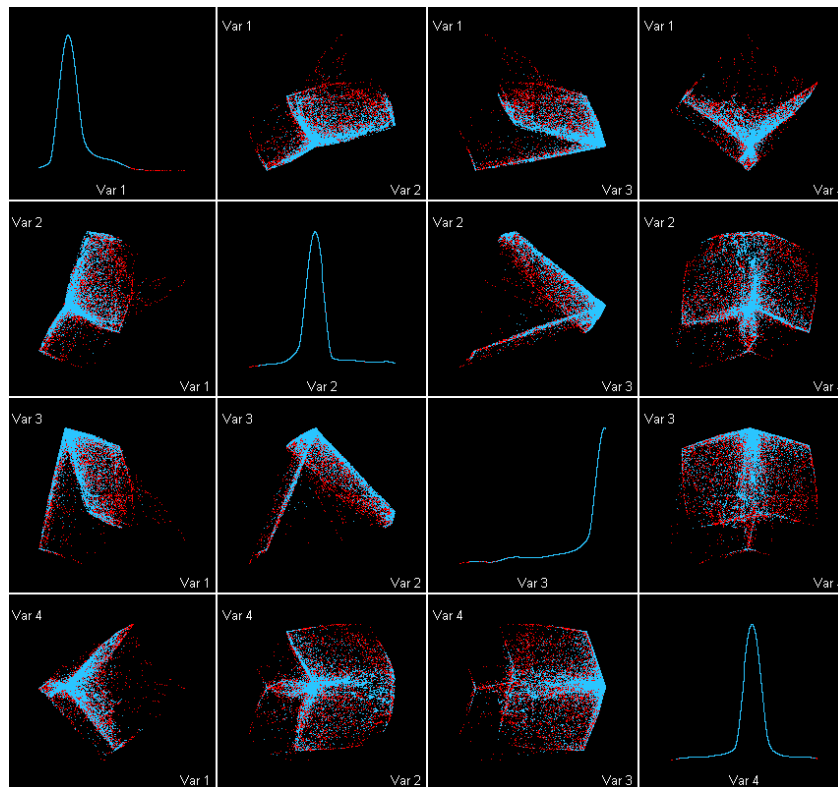


Figure IV-73 : Représentation graphique des états sensoriels observés au cours de l'expérience 37 (points rouges) et ceux observés lors de l'expérience 38 (points bleus).

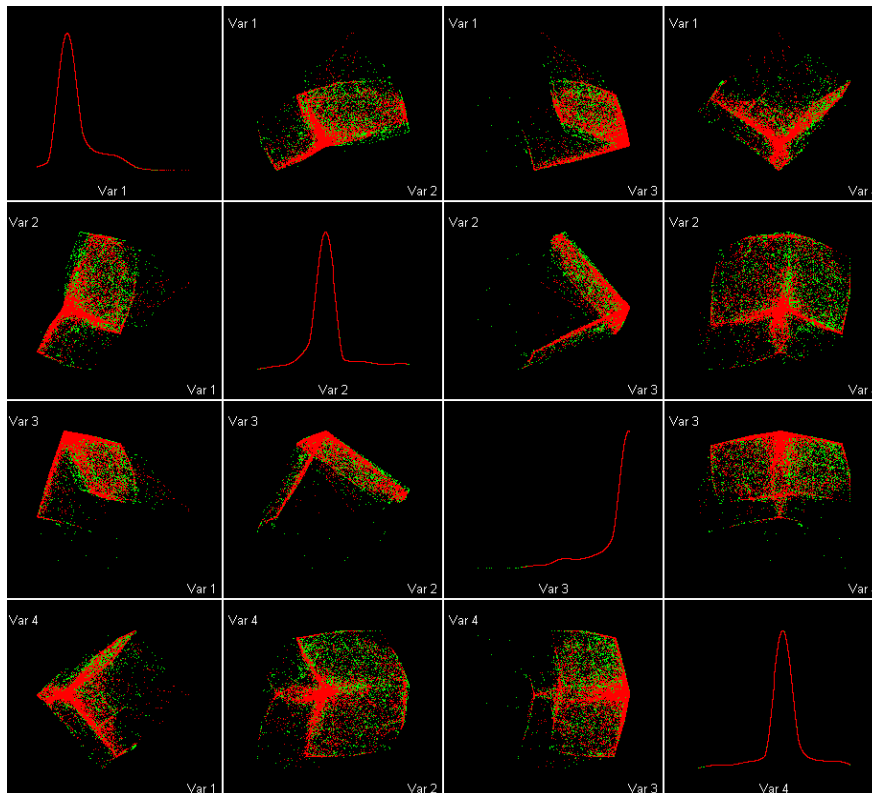


Figure IV-74 : Représentation graphique des états sensoriels observés au cours de l'expérience 37 (points rouges) et ceux observés lors de l'expérience 39 (points vert).

Toutefois, une distinction apparaît dans l'homogénéité de la répartition des états sensoriels entre les expériences 37 et 38 et l'expérience 39, des amas d'états sensoriels se dégagent (Figure IV-75) dans les deux premières. Ces amas signifient que le robot s'est retrouvé plusieurs fois dans la même configuration, cependant, ces configurations ne peuvent être interprétées comme des signatures globales de l'environnement car les positions de ces amas sont différentes entre les deux expériences. L'existence de ces amas pourrait traduire les phases pendant lesquelles le comportement du robot devient cyclique, mais en fait, la majorité des points de ces amas se suit temporellement, indiquant qu'ils résultent des situations de quasi immobilité du robot. La configuration sensorielle change alors lentement, sauf en la présence d'un bruit infrarouge. La différence de positions des amas entre les deux expériences peut venir de la variabilité de la position de départ. Toutefois, ces amas se situent dans une même zone conique de l'espace sensoriel.

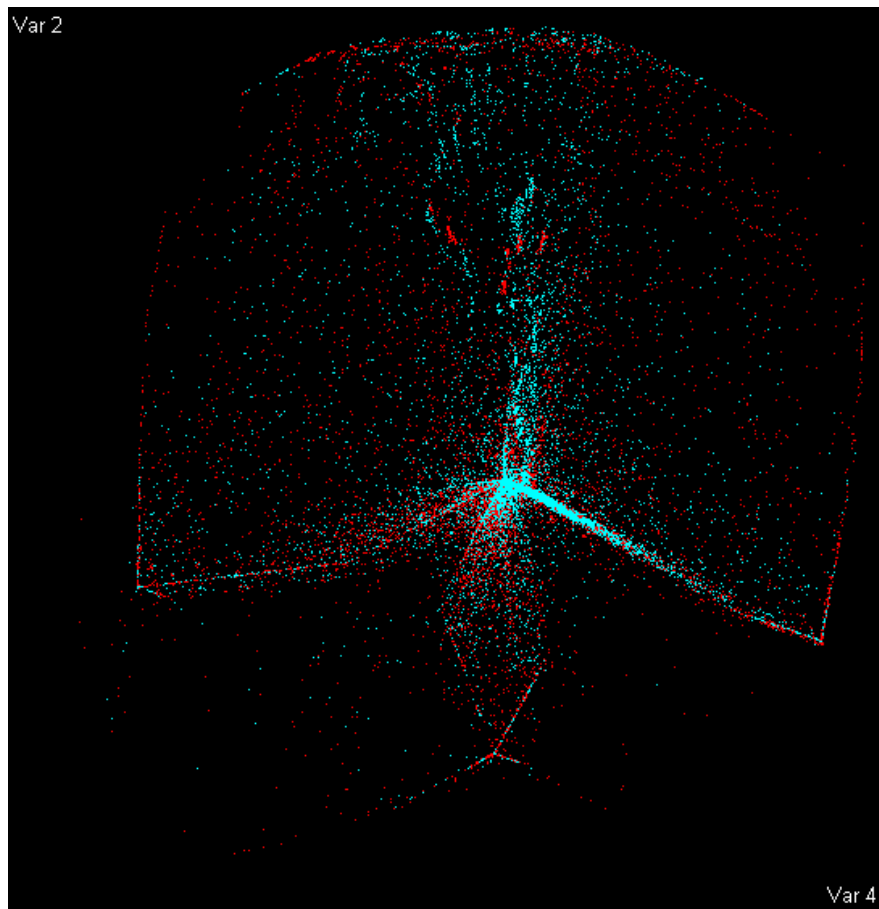


Figure IV-75 : Représentation graphique des variables 2 et 4 des états sensoriels observés au cours de l'expérience 37 (points rouges) et ceux observés lors de l'expérience 38 (points bleus).

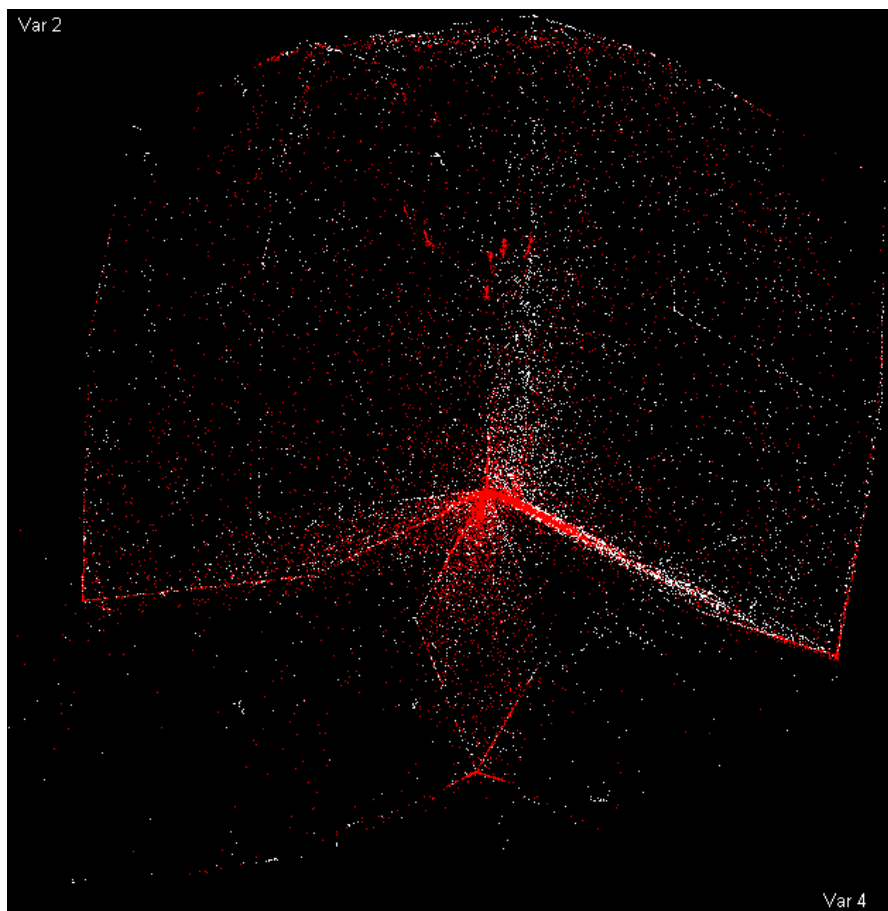


Figure IV-76 : Représentation graphique des variables 2 et 4 des états sensoriels observés au cours de l'expérience 37 (points rouges) et ceux observés lors de l'expérience 39 (points verts).

La Figure IV-76 montre que l'expérience 39 n'a pas formé des amas comme ceux observés dans les expériences 37 et 38. Une explication est que la présence d'obstacles aigus introduit des bifurcations dans la trajectoire du robot. Ainsi, le phénomène de comportement cyclique se trouve peu observé, cela se traduit par une répartition plus homogène des états sensoriels observés.

5.1.2. La conservation du comportement

L'étude comportementale du robot constitue la deuxième approche pour aborder les résultats de ces expériences. Contrairement aux conclusions motrices, les prémisses sensorielles évoluent, mais cette évolution conserve malgré tout le comportement d'évitement d'obstacle. Dans l'expérience 39, l'environnement se complexifie mais n'altère pas le comportement d'évitement du robot.

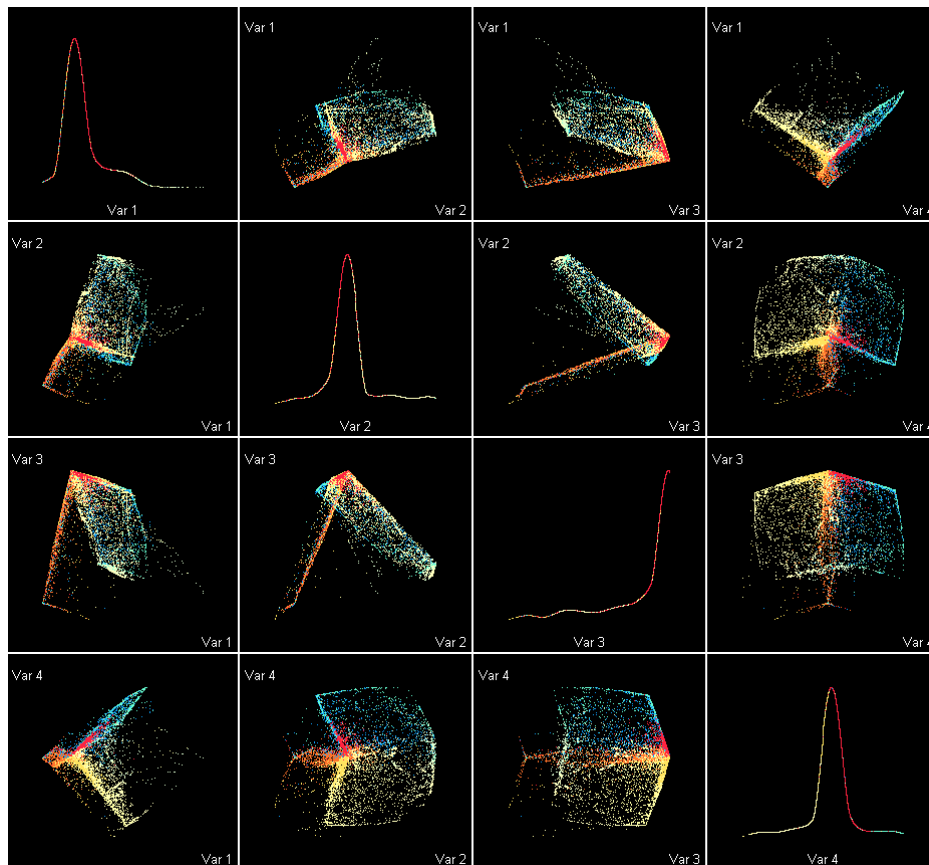


Figure IV-77 : Représentation graphique sensorielle observée au cours de l'expérience 37 : la couleur des points est fonction de la direction de l'effort moteur.

La Figure IV-77 affiche l'orientation de la force motrice suivant les commandes effectuées pour chaque état sensoriel rencontré lors de l'expérience 38. L'orientation de la force motrice est retrouvée en calculant l'arc tangent du rapport des commandes motrices. L'amplitude de la force n'est pas prise en compte dans ce type de représentation des consignes. La distinction de couleurs entre deux plans de l'espace sensoriel représente les actions pour lesquelles le robot tourne à gauche et celles où il s'oriente à droite. Le foyer de points de couleur rouge, correspondant à la non-détection d'obstacle, désigne l'action d'aller tout droit. La présence de ces aplats de couleurs montre qu'une règle peut glisser dans un sous-espace tout en restant cohérente avec le comportement d'évitement d'obstacle. Il en va de même pour la création de règles qui sera abordée ultérieurement.

Les prémisses sensorielles des règles initiales ne sont pas réparties de manière égale dans l'espace sensoriel. De part la relation qui les unie, ce déséquilibre se retrouve également dans les consignes motrices. Le robot tourne plus souvent d'un côté que de l'autre. La Figure IV-78 représente le moment de déclenchement de chaque règle, colorisé selon l'orientation de la motrice de commande induite par la règle. Au début de l'expérience, un grand nombre de règles se déclenchent et 80 % d'entre elles conduisent à une commande tournant dans une même direction. La fin de la première phase de la compétition à 580 s rééquilibre le nombre de règles s'orientant à droite et à gauche. Toutefois, le comportement global du robot semble privilégier une orientation, ce qui peut s'expliquer soit par la dissymétrie de la qualité du capteur, soit par le manque de règles à des

endroits critiques, soit par une courbe de réponse motrice trop inégale entre les deux moteurs.

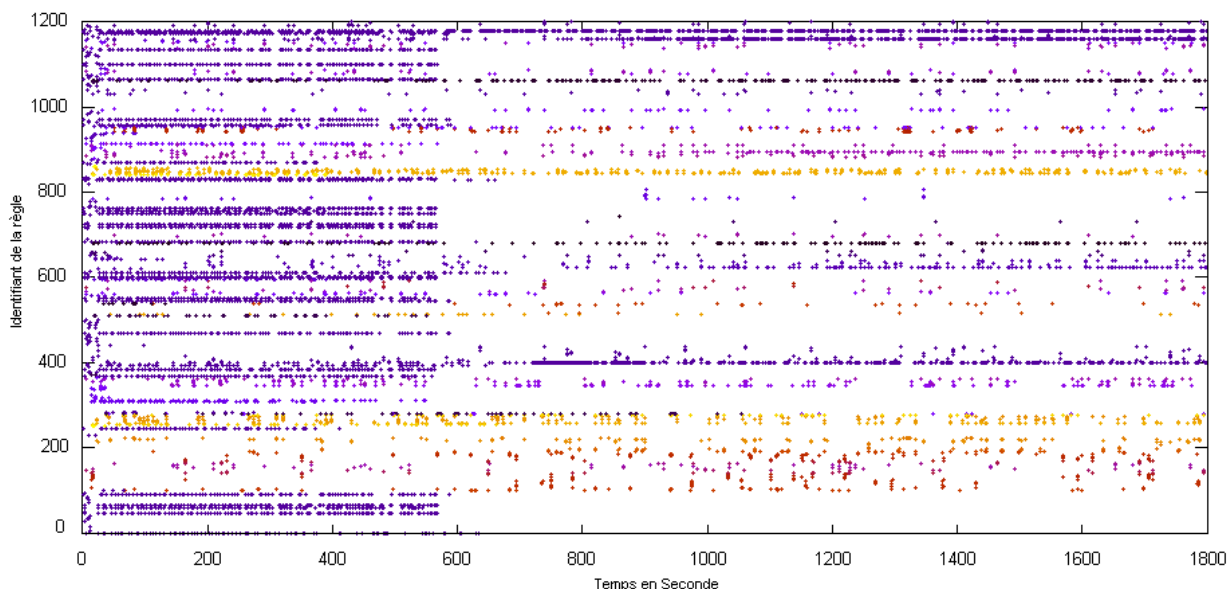


Figure IV-78 : Graphique indiquant la règle déclenchée au cours du temps de l'expérience 37 dont la couleur correspond à l'orientation de l'effort moteur.

L'expérience 38 montre la stabilité dans le déclenchement des règles et dans le maintien de l'équilibre entre les actions, Figure IV-79. Le graphique de l'expérience 39 possède une répartition temporelle des déclenchements identique à celle de l'expérience 38, Figure IV-80. L'environnement complexe semble toutefois participer à la rééquilibration des mouvements consistant à tourner soit à droite, soit à gauche.

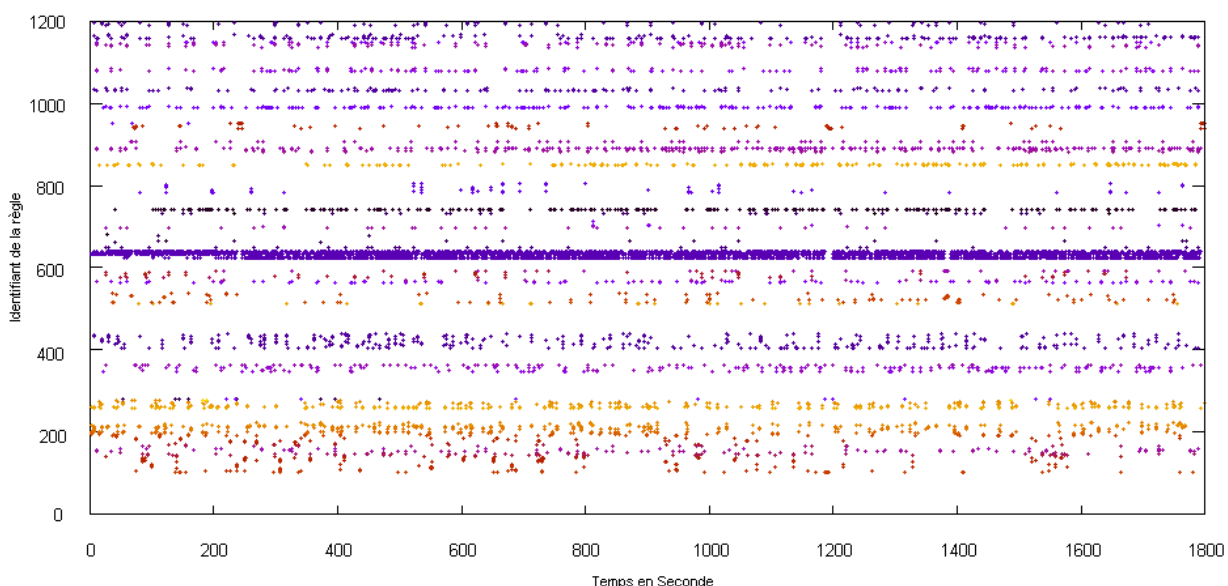


Figure IV-79 : Graphique indiquant la règle déclenchée au cours du temps de l'expérience 38 dont la couleur correspond à l'orientation de l'effort moteur.

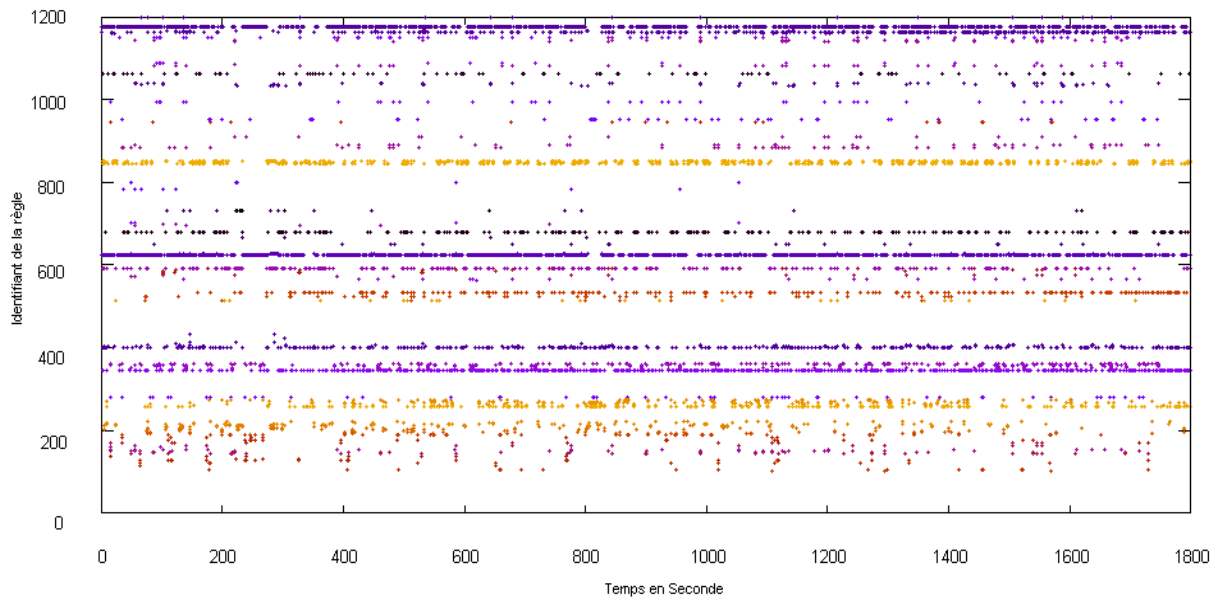


Figure IV-80 : Graphique indiquant la règle déclenchée au cours du temps de l'expérience 39 dont la couleur correspond à l'orientation de l'effort moteur.

La réduction des états moteurs possibles suite à l'élimination des règles touche principalement les états extrêmes (Figure IV-81). Ces états moteurs correspondent à des situations aberrantes pour un comportement d'évitement d'obstacle dans un environnement statique simple dont la position d'origine se trouve sans obstacle à proximité. Dans l'expérience 38, la répartition de l'espace des commandes motrices est conservée, bien que 30 états moteurs sur les 130 du départ soient éliminés (Figure IV-82).

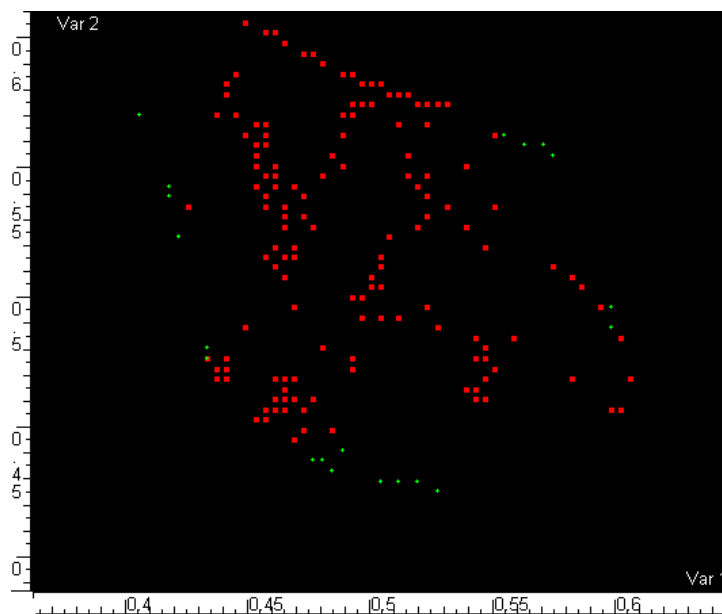


Figure IV-81 : Représentation des consignes motrices en 8 mm/s à la fin de l'expérience 37. Les croix vertes représentent les commandes éliminées et les carrés rouges sont les commandes conservées. Var2 : moteur gauche, Var1 : moteur droit.

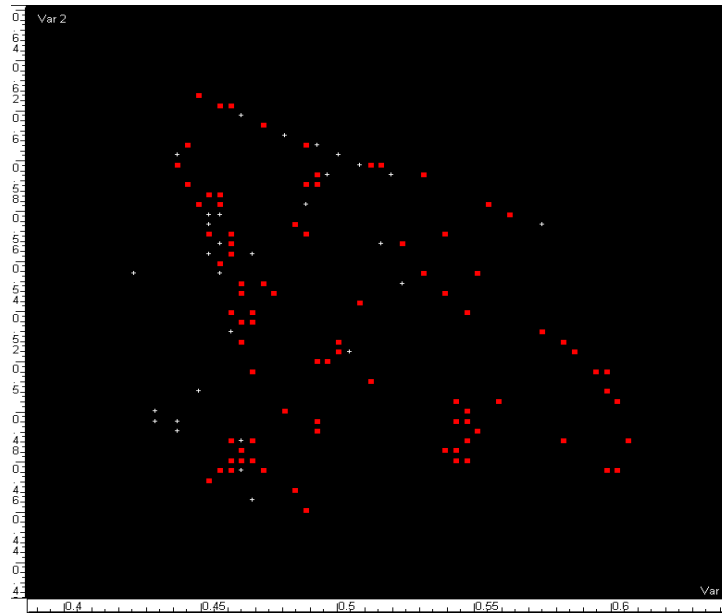


Figure IV-82 : Représentation des consignes motrices en 8 mm/s à la fin de l'expérience 38. Les croix vertes représentent les commandes et les carrés rouges sont les commandes conservées. Var2 : moteur gauche, Var1 : moteur droit.

5.1.3. L'évolution des populations de règles

Le troisième point de vue se penche sur l'évolution démographique de la base de règles et sur leur expressivité. Les cinq phases dégagées lors des expériences en environnement imposé se retrouvent dans la Figure IV-73. À la fin de l'expérience 37, il reste 230 règles dont 32 n'ayant jamais été déclenchées et 35 ayant été déclenchées entre 1 et 5 fois. Elles feront partie des 72 règles éliminées suite à la reconduite de l'expérience 37, soit l'expérience 38. Cependant, les règles qui se sont déclenchées jusqu'à 1000 fois dans l'expérience 37 ne se sont pas déclenchées au cours de l'expérience 38. Le détail de cette compétition sera davantage abordé lors de l'analyse de la dynamique des règles.

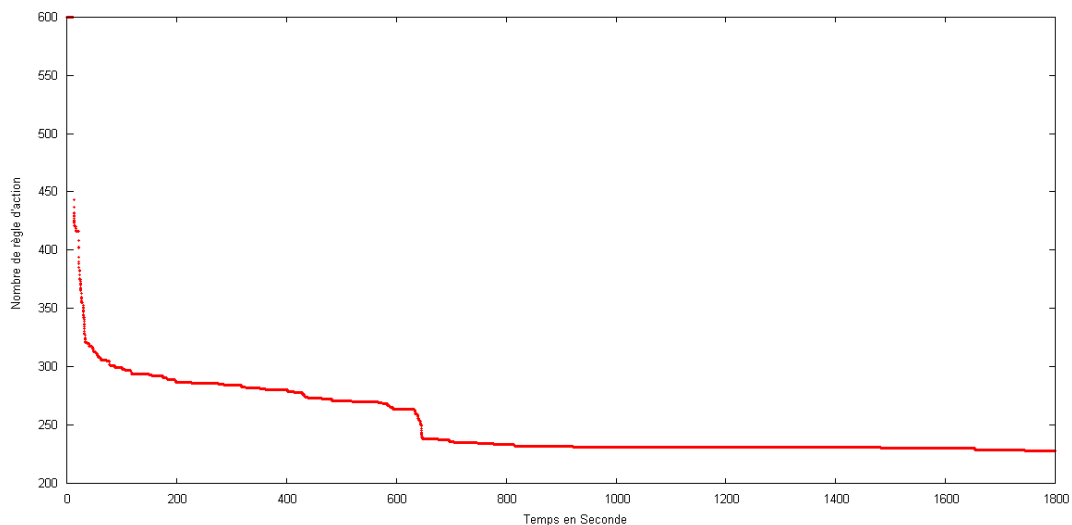


Figure IV-83 : Évolution du nombre de règles d'action au cours de l'expérience 37.

Au cours de l'expérience 39, 54 règles ont été éliminées et parmi les règles restantes à la fin, 10 ne se sont jamais déclenchées. La présence de règles qui ne se déclenchent plus dans les expériences 38 et 39 illustre la permanence de la compétition et la sensibilité du système dynamique sous-jacent.

Une manière d'appréhender l'expressivité des règles réside dans l'étude de leur répartition suivant le nombre de déclenchements. Pour l'expérience 37 (Figure IV-84), les deux règles ayant le plus grand nombre de déclenchements possèdent des prémisses sensorielles similaires, ce qui signifie qu'une seule prémisse correspond à plus d'un tiers des situations. Ainsi, les règles sensorimotrices adaptées peuvent être considérées comme un code optimal décrivant l'environnement en fonction des réactions suscitées. La notion d'entropie devient alors une mesure intéressante pour comparer les expériences.

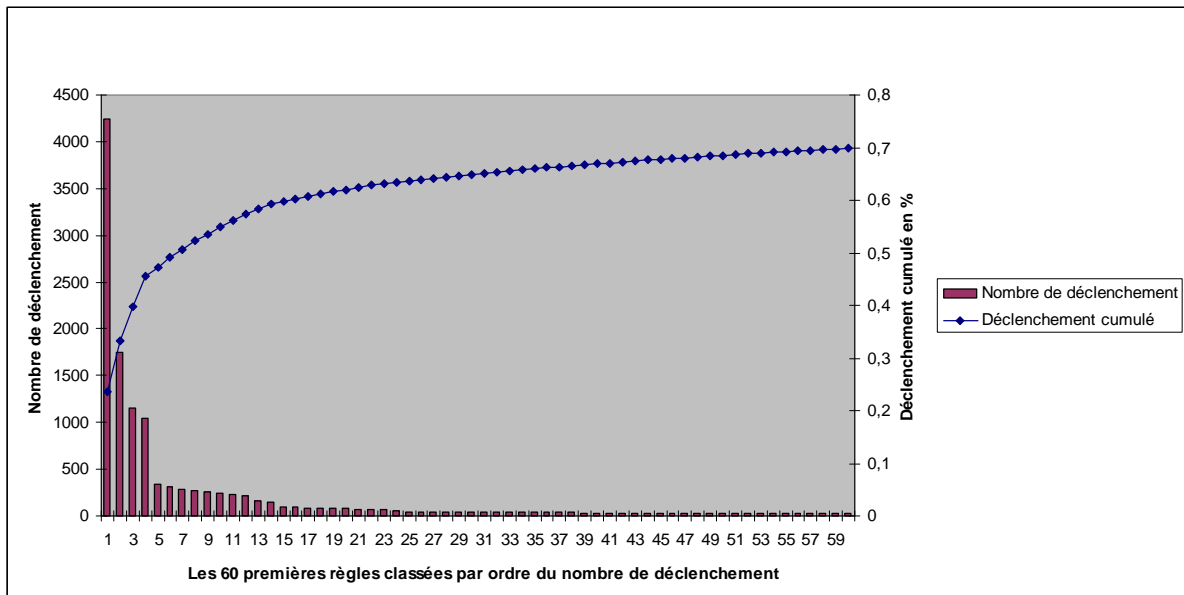


Figure IV-84 : Histogramme du nombre de déclenchements des 60 premières règles les plus actives et la courbe du pourcentage cumulé de leur participation au déclenchement global.

L'entropie $E(X)$ de Shannon correspond à une mesure de la dispersion de l'information, avec p_i la probabilité d'un événement désigné par l'indice i :

$$E(X) = -\sum_{i=1}^n p_i * \frac{\ln(p_i)}{\ln(2)}$$

L'entropie basée sur la probabilité de déclenchement d'une règle par rapport à l'ensemble des règles initiales est de 5,65 bits. Toutefois, ce calcul sur l'ensemble des règles ne correspond pas à la probabilité effective du déclenchement des règles. En effet, une fois éliminée, une règle a une probabilité nulle de se déclencher. En calculant l'entropie uniquement sur la population finale, elle descend à 4,53 bits. Le désordre maximal d'un code correspond à l'apparition équiprobable des mots qui le composent, l'entropie vaut alors :

$$E(X) = -\frac{\ln(n)}{\ln(2)}$$

La Figure IV-85 indique l'entropie calculée à partir de l'ensemble des règles initiales et de l'ensemble des règles finales, ainsi que leur entropie maximale respective. Concernant l'expérience 37, la baisse de l'entropie doit provenir davantage de la diminution du nombre de règles plus que d'une réelle structuration de l'information. Il sera intéressant de comparer la mesure de l'entropie et son évolution avec les expériences commençant directement par l'emploi d'environnements complexes.

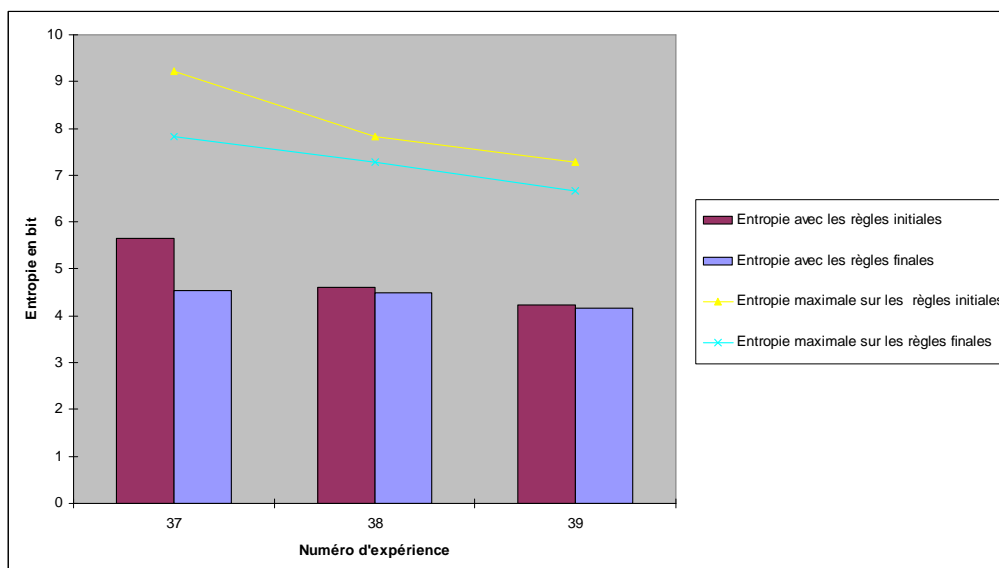


Figure IV-85 : Évolution de l'entropie de l'ensemble des règles suivant les expériences 37, 38 et 39.

5.1.4. Les dynamiques

Le quatrième point de vue porte sur la dynamique du système au travers des valeurs des règles déclenchées. La dynamique de l'alpha montre la spécialisation de certaines règles à partir de 650 s. Suite à la seconde transition, l'augmentation de la force des règles les plus souvent utilisées monte à 1. Cette observation conforte l'idée que le système n'est pas voué à l'élimination totale des règles. Le score augmente globalement, illustrant ainsi l'ajustement des prémisses comme pour les expériences avec l'environnement simulé. La transition à 650 s correspond à l'effondrement démographique qui doit par ailleurs faciliter la spécialisation des règles restantes.

La dynamique de l'expérience 38, illustrée par la Figure IV-87, se stabilise et aucune transition brusque n'est observée contrairement à la Figure IV-86. De même, la convergence de l'évaluation et du score à partir de 1200 s montre une stabilisation de l'ajustement des prémisses. La force moyenne, après la décroissance au cours de l'expérience 37, stagne autour de 0,18 puis vers 1200 s elle croît de nouveau avec approximativement un coefficient directeur de $6,6 \cdot 10^{-4}$. Cela confirme, en plus du fait que la force de certaines règles atteint une valeur de 1, que le système ne va pas vers l'élimination totale ou vers un état stable fragile dû à une force moyenne faible.

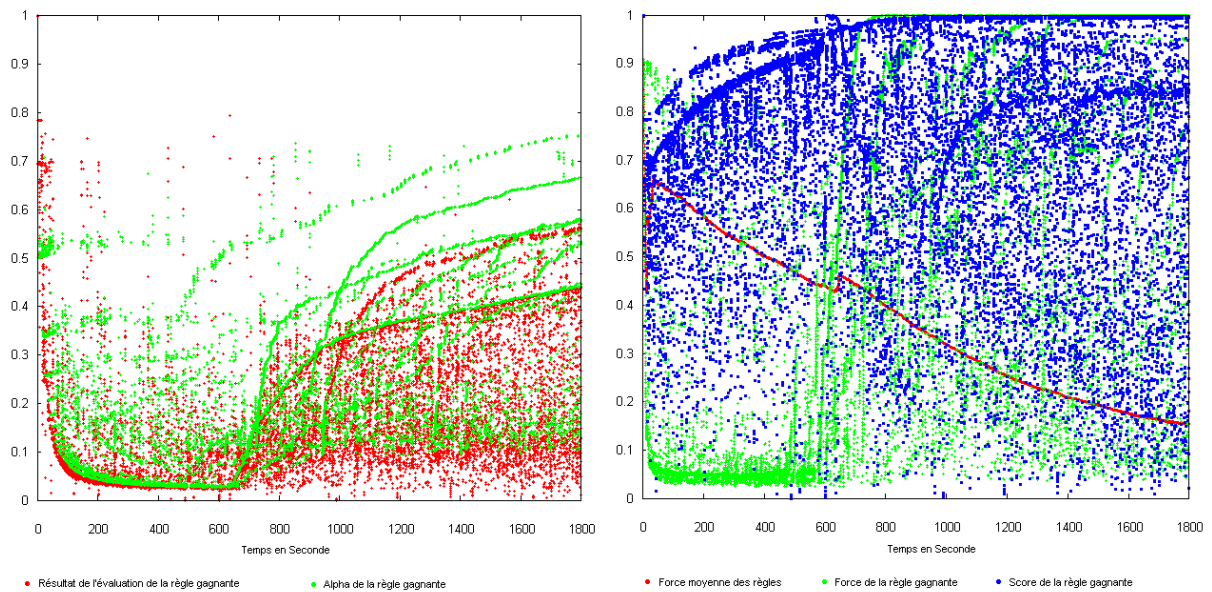


Figure IV-86 : Représentation graphique à droite de l'alpha et du résultat de l'évaluation de la règle déclenchée, à gauche représentation de la force et du score de la règle déclenché ainsi que la force moyenne des règles de l'expérience 37.

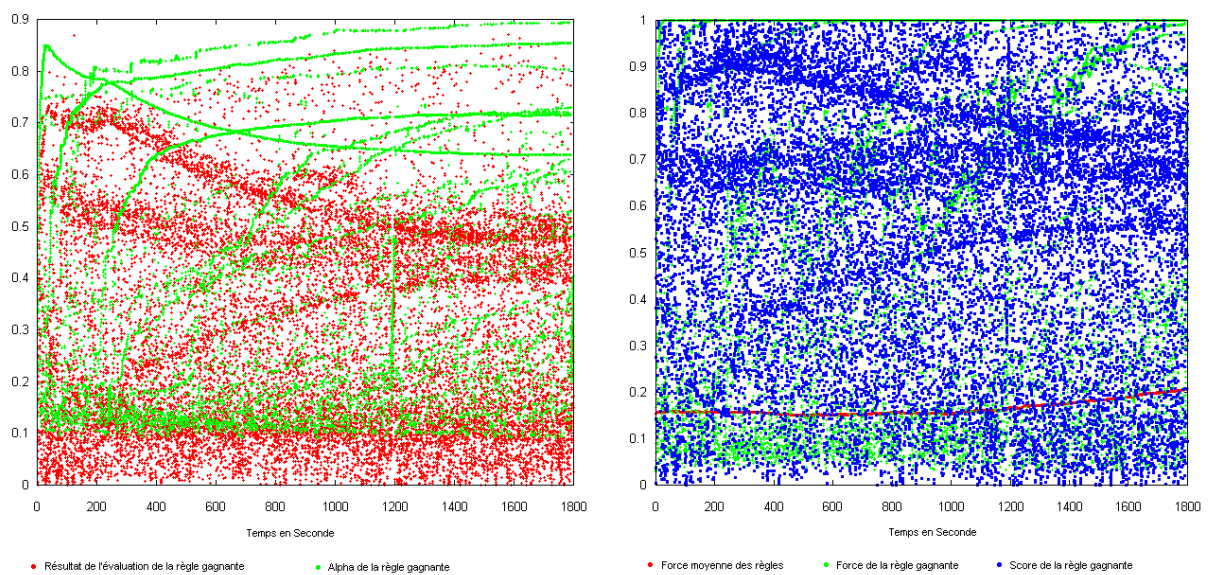


Figure IV-87 : Représentation graphique à droite de l'alpha et du résultat de l'évaluation de la règle déclenchée, à gauche représentation de la force et du score de la règle déclenché ainsi que la force moyenne des règles de l'expérience 38.

Concernant l'expérience 39, la complexification de l'environnement a légèrement perturbé la dynamique, dans la mesure où une règle arrive à une position très dominante, contrairement aux deux autres expériences. L'augmentation continue et rapide des caractéristiques de cette règle transparait dans la Figure IV-88. Après vérification, cette règle a été déclenchée 266 fois au cours de l'expérience 37 et 4 fois au cours de l'expérience 38. Sur le même espace sensoriel, une dizaine de règles dans l'expérience 37 subsistent. La fréquence d'apparition des états sensoriels correspondant permet le maintien de toutes ces règles, cependant, trois d'entre elles se déclenchent 100 à 500 fois plus que les autres. Dans

l'expérience 38, ces trois règles n'ont quasiment pas été déclenchées et trois autres ont pris leur place en terme de nombre de déclenchements. Au cours de la troisième expérience, seules deux règles se déclenchent pour un même type d'état sensoriel : l'une est issue de l'expérience 38 et l'autre, issue de l'expérience 37, redevient compétitive dans l'expérience 39. La complexité de la compétition illustrée par cet exemple n'empêche pas la convergence du système au-delà des 90 min, soit 54000 sélections. Par ailleurs, cette convergence se trouve fortement suggérée par l'augmentation continue de la force moyenne et les scores élevés observés (Figure IV-88).

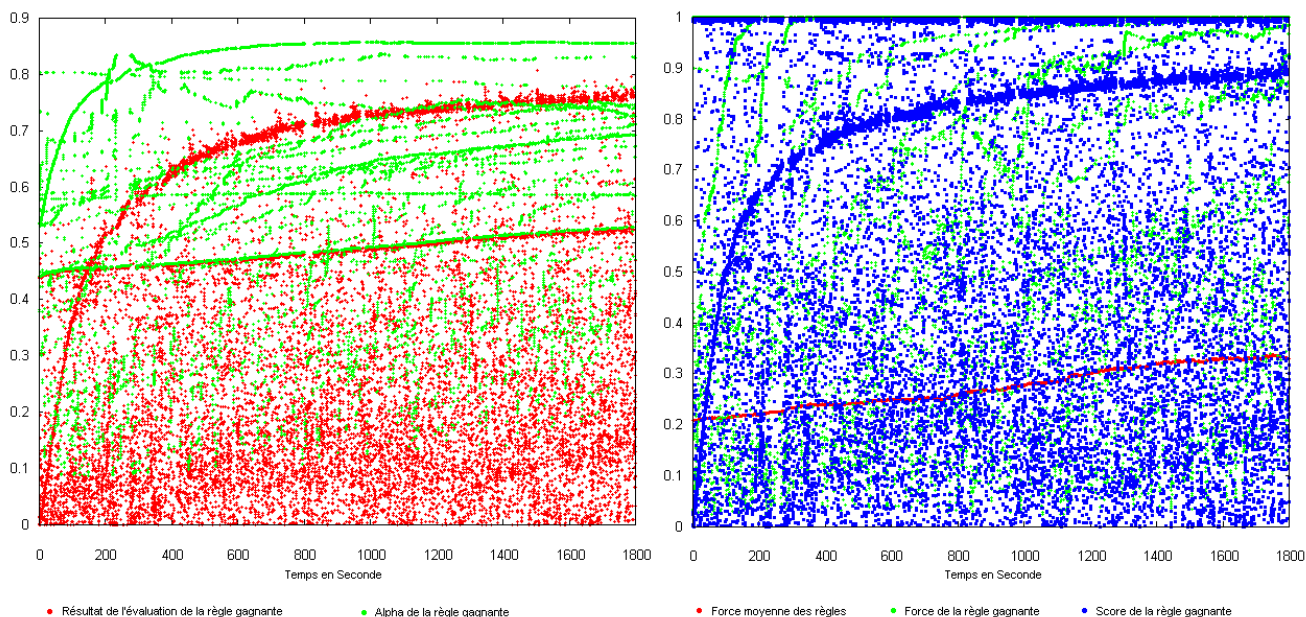


Figure IV-88 : Représentation graphique à droite de l'alpha et du résultat de l'évaluation de la règle déclenchée, à gauche représentation de la force et du score de la règle déclenchée ainsi que la force moyenne des règles de l'expérience 39.

5.1.5. Le positionnement des règles finales

La répartition des prémisses des règles à la fin de l'expérience 37 couvre l'espace sensoriel observé et les règles initiales situées en dehors de cet espace ont été éliminées (Figure IV-89).

La Figure IV-90 illustre la migration des prémisses de certaines règles vers des états sensoriels plus fréquents, notamment en bas à gauche. Elle traduit également la forte compétition dans la zone à densité maximale. En effet 274 prémisses se trouvent dans un rayon de 0,015 du foyer, sachant que dans l'espace sensoriel les valeurs pour l'axe de la première composante (Var1) vont de 0 à 0,75 et pour la seconde (Var2) de 0,15 à 1. Cette petite zone contient environ 46,5% des règles initiales. À la fin de l'expérience, le nombre de prémisses se trouvant dans ce même périmètre s'élève à 4, soit moins de 0,02% des règles finales. La compétition permet bien d'éliminer les règles redondantes.

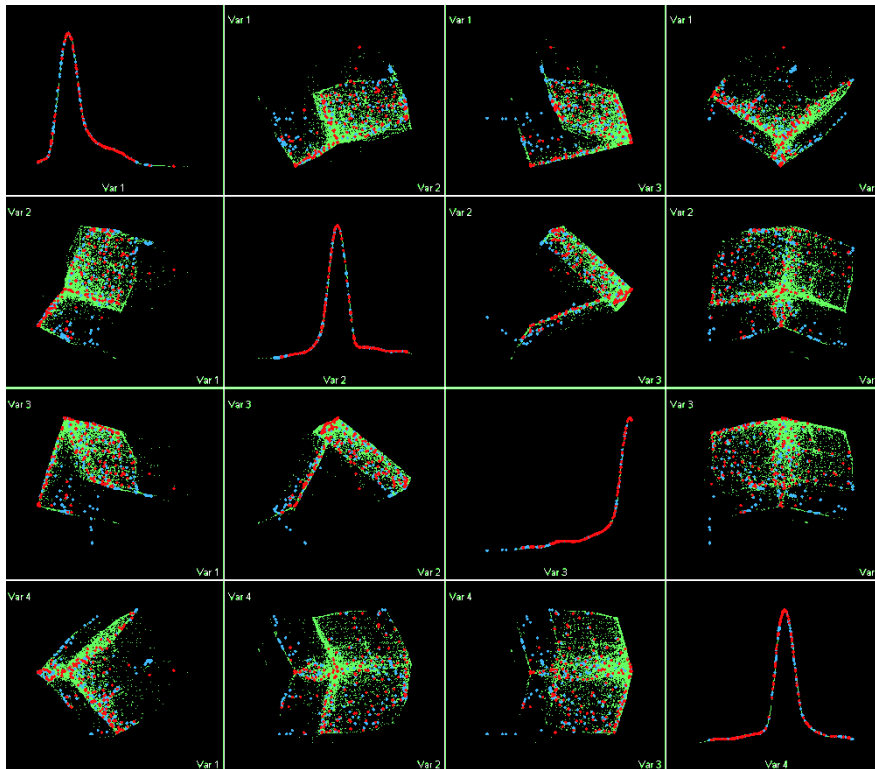


Figure IV-89 : Représentation de l'environnement (croix vertes), des prémisses appartenant aux règles initiales (croix bleus) et des prémisses appartenant aux règles finales (ronds rouges) de l'expérience 37.

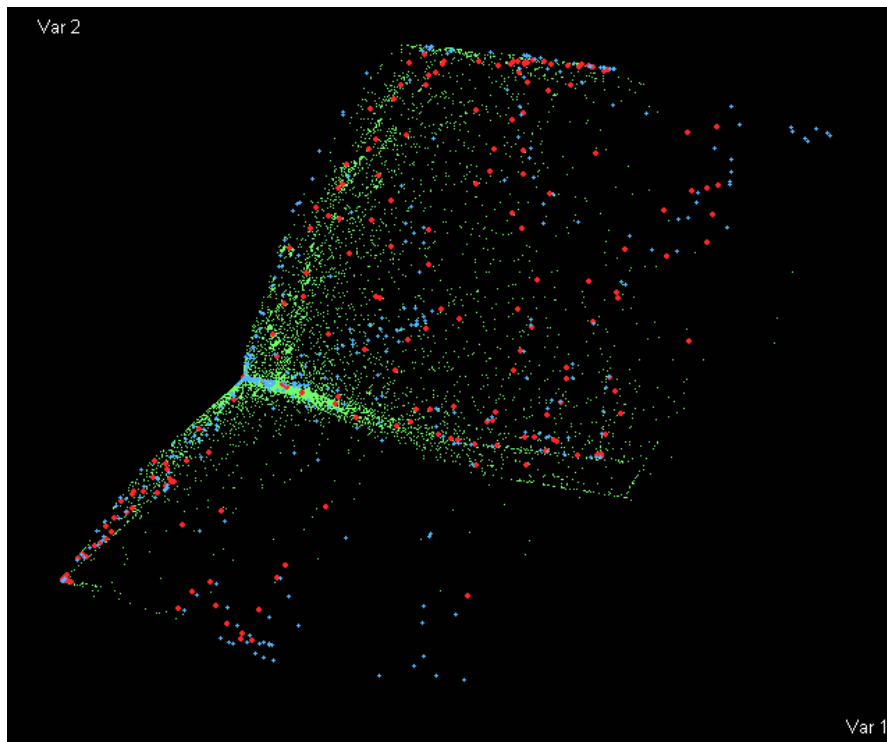


Figure IV-90 : Représentation de l'environnement (croix vertes), des prémisses appartenant aux règles initiales (croix bleus) et des prémisses appartenant aux règles finales (ronds rouges) de l'expérience 37.

Des états sensoriels forment des amas et des prémisses de règles se situent à proximité ou sur ces amas. La Figure IV-91 colorie tous les états sensoriels ayant déclenché une même règle au cours de l'expérience 37. Le champ des prémisses est varié et l'incertitude sur les valeurs des prémisses peut fluctuer d'un facteur 100. Une même prémisses peut couvrir plusieurs amas. Une capture plus fine des prémisses doit être possible grâce à un maillage plus précis des règles initiales et un paramètre d'enchère plus faible. L'architecture proposée permet bien l'adaptation des prémisses à un type de stimulation, même en environnement réel, avec de fortes variabilité et incertitude.

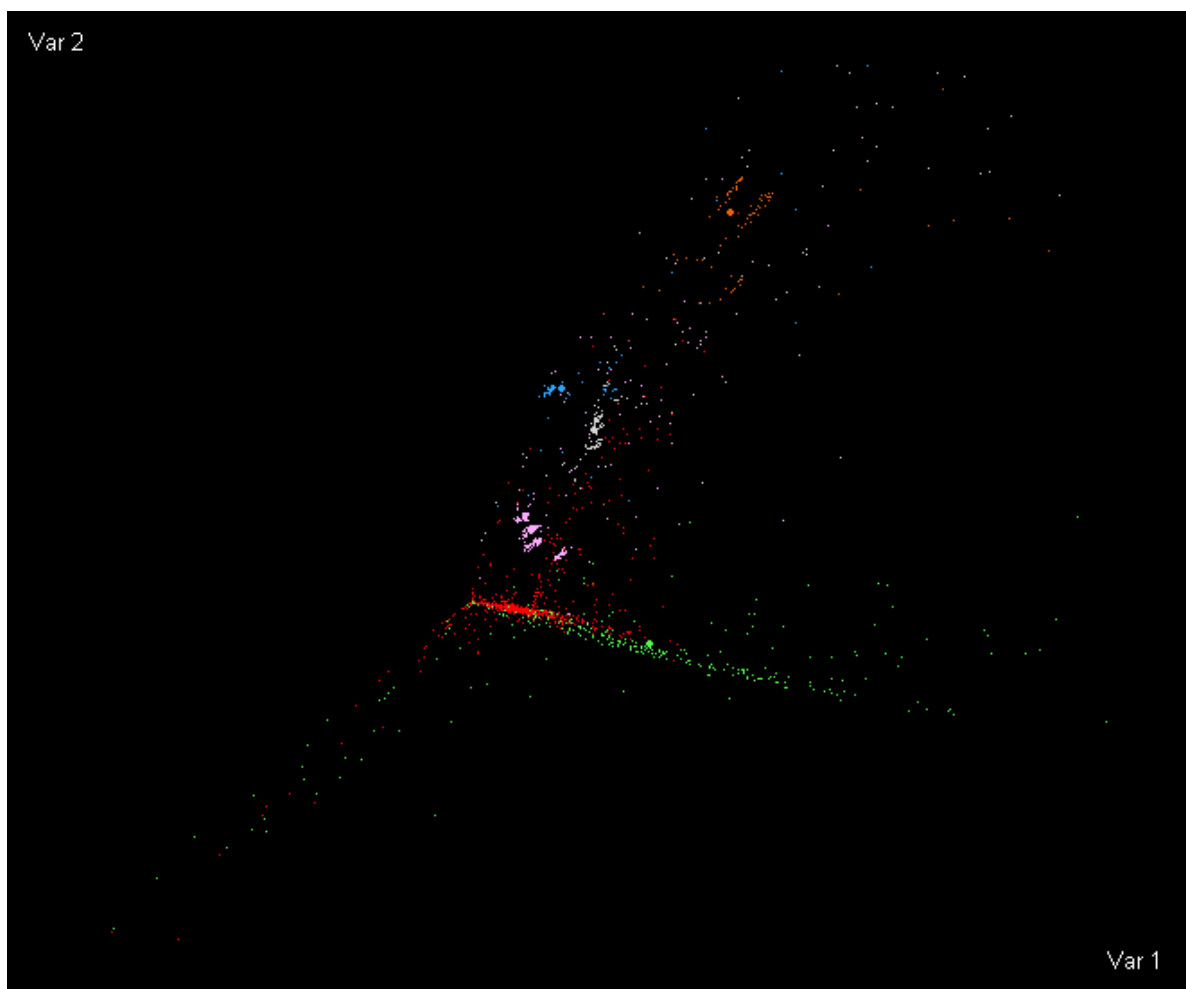


Figure IV-91 : Représentation sur les deux premières composantes avec la couleur associée des états sensoriels (points) déclenchant une des 6 règles sélectionnées (ronds) parmi les règles finales de l'expérience 37.

La comparaison des règles finales des expériences 38 et 39 montre une diminution des prémisses ainsi qu'une répartition homogène de l'espace sensoriel observé. Cette diminution est particulièrement visible sur les projections comprenant la quatrième composante sur la Figure IV-92. L'environnement complexe a pu jouer un rôle dans l'élimination d'une certaine redondance, maintenue par la dissymétrie de l'espace sensoriel stimulé puisque le nombre de situations variées était plus important.

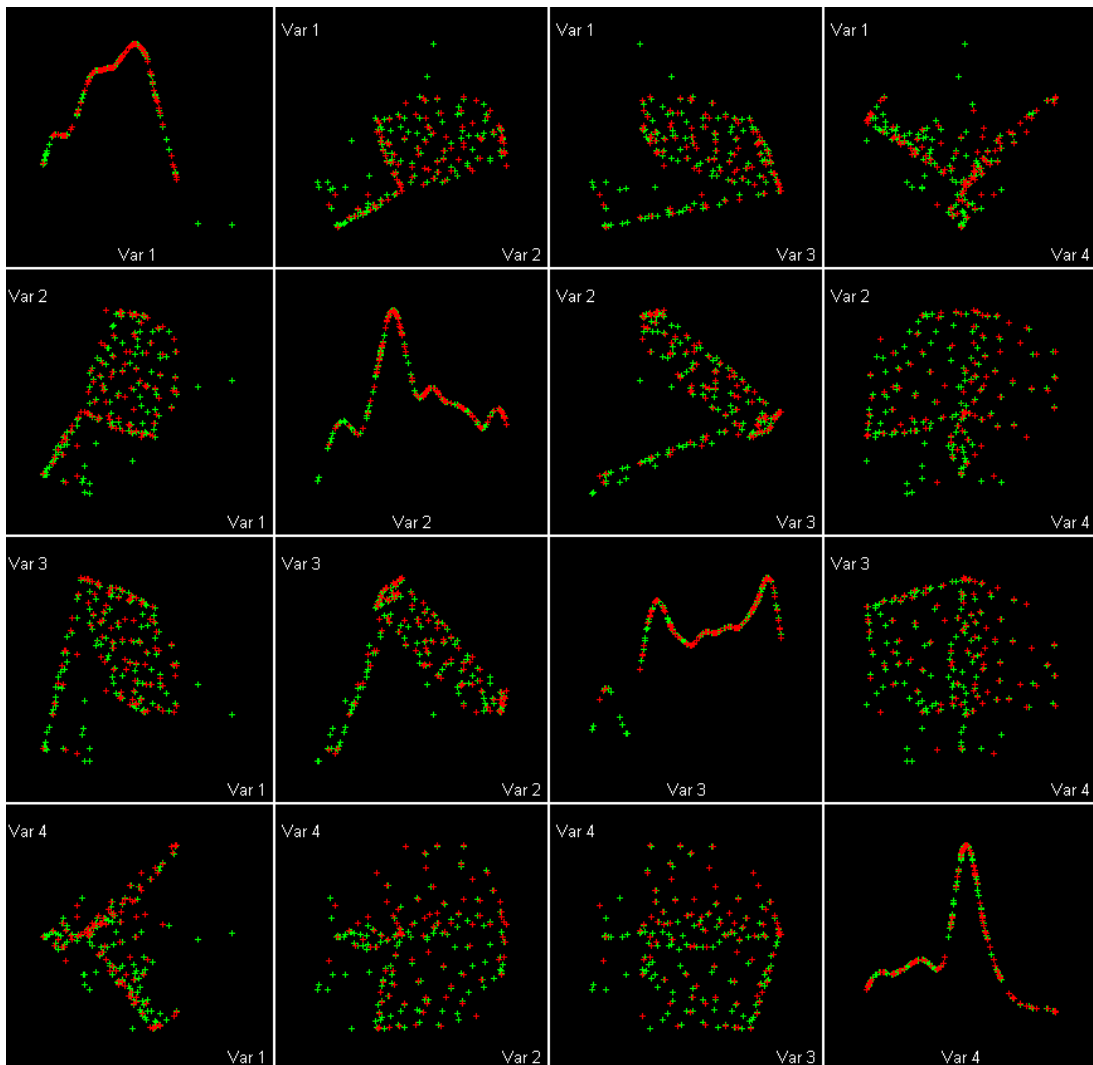


Figure IV-92 : Représentation des prémisses appartenant aux règles finales de l'expérience 38 (croix vertes) et celles de l'expérience 39 (croix rouges).

Les observations des états sensoriels de l'expérience 39 ne présentent pas d'amas, comme précédemment. Les prémisses sensorielles s'étendent davantage dans l'espace sensoriel (Figure IV-93 et Figure IV-94) et la répartition des zones d'influence se révèle moins nette, tout en restant visible.

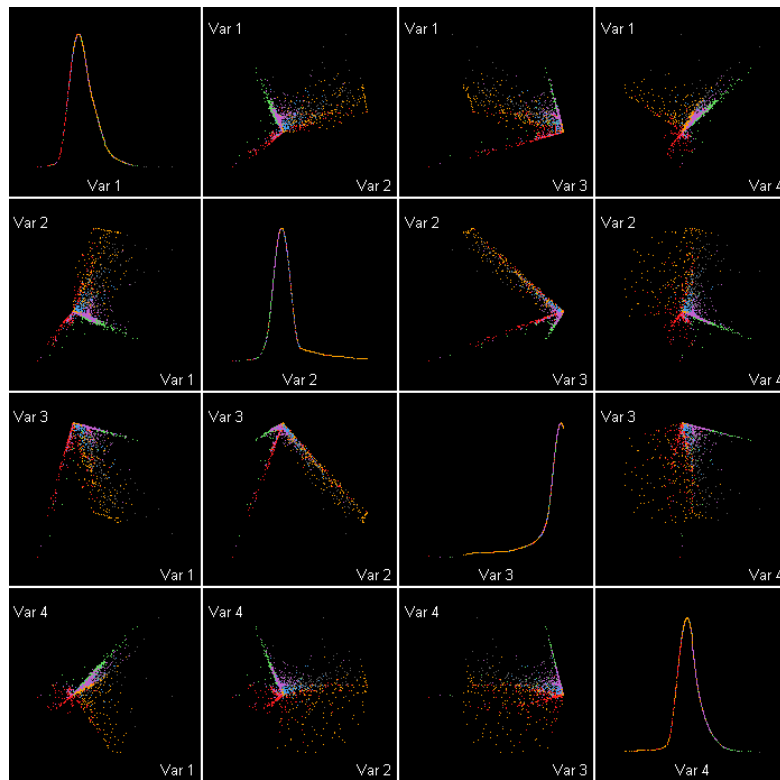


Figure IV-93 : Représentation des états sensoriels observés déclenchant l'une des 5 règles sélectionnées.

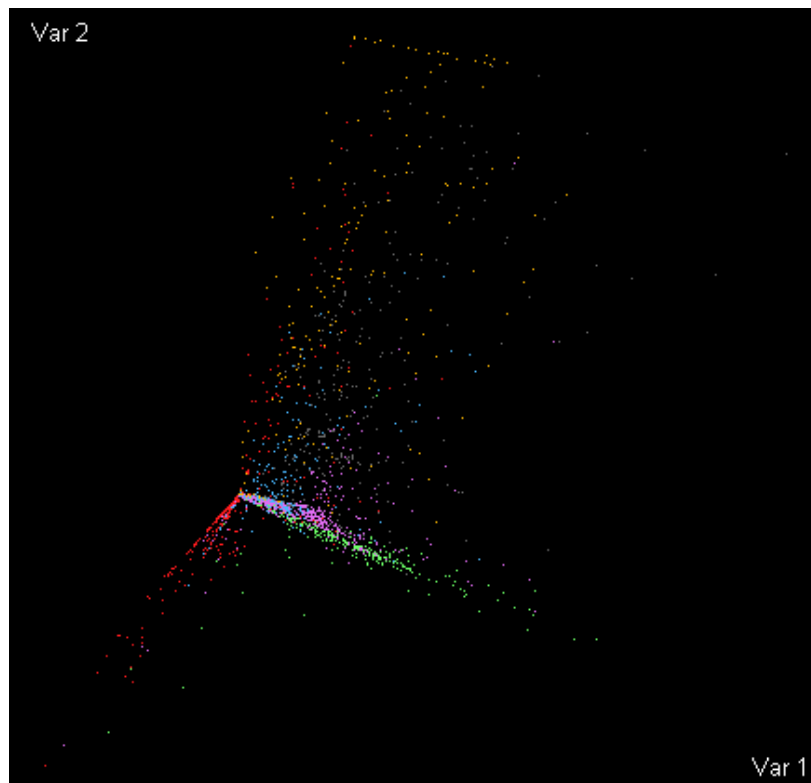


Figure IV-94 : Représentation des états sensoriels observés déclenchant l'une des 5 règles sélectionnées selon les deux premières composantes.

5.1.6. Les structures temporelles

La dernière analyse des résultats de cette première série d'expériences porte sur l'enchaînement des règles. Les amas observés donnent des indices sur la périodicité du comportement mais sans précisément déceler de véritables structures temporelles, comme par exemple une séquence d'évitement. La représentation des déclenchements des règles, comme celle de la Figure IV-80, ne facilite pas non plus l'observation de séquences d'actions. L'élaboration de bandes chronologiques des règles déclenchées devient alors nécessaire. Ces bandes se colorisent en fonction des états sensoriels, des commandes effectuées ou de l'identifiant de la règle déclenchée.

La compétition se traduit par l'évolution du nombre de franges de couleurs différentes au début de l'expérience vers un nombre plus réduit et surtout correspondant à la variation de la stimulation sensorielle (Figure IV-95). Le grand aplatissement de couleur vert clair illustre la suprématie d'une règle qui s'explique par la fréquence d'apparition de l'état sensoriel correspondant. Les quatre premiers graphiques illustrent l'adéquation entre les états sensoriels et les prémisses des règles déclenchées.

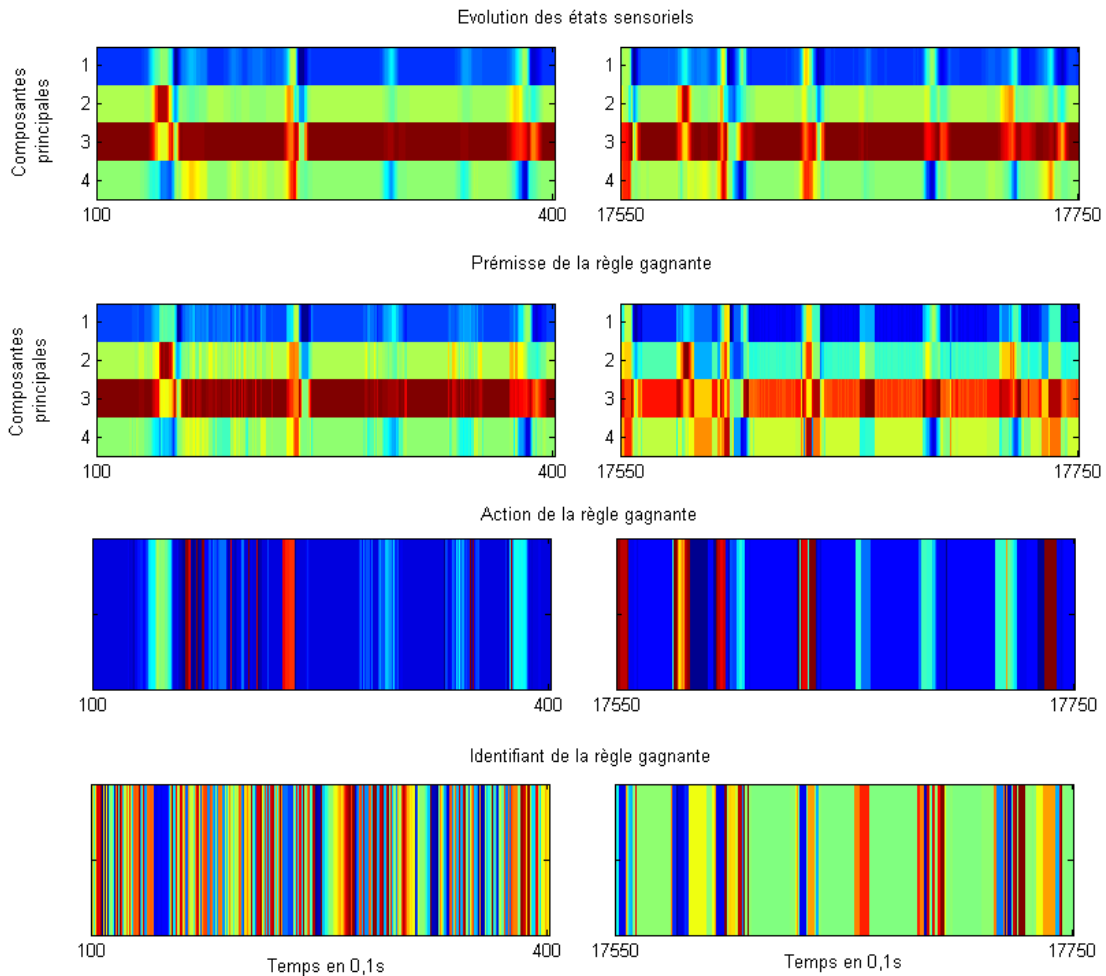


Figure IV-95 : Bandes chronologiques de l'expérience 37 coloriées en fonction des états sensoriels observés, des prémisses déclenchées, des commandes effectuées et de l'identifiant de la règle déclenchée.

En dilatant encore l'échelle du temps (Figure IV-96), une transition correspond à un enchaînement d'états sensoriels distincts avec une valeur motrice associée particulière. Cet enchaînement conforte les choix effectués concernant l'amplitude des consignes de vitesse, la fréquence d'acquisition sensorielle et la fréquence de sélection. La transition visualisée dure 1,4 s, soit 14 sélections qui se répartissent sur une dizaine de règles différentes. L'évaluation de la durée de ces transitions déterminera les paramètres fixés pour la création de règles.

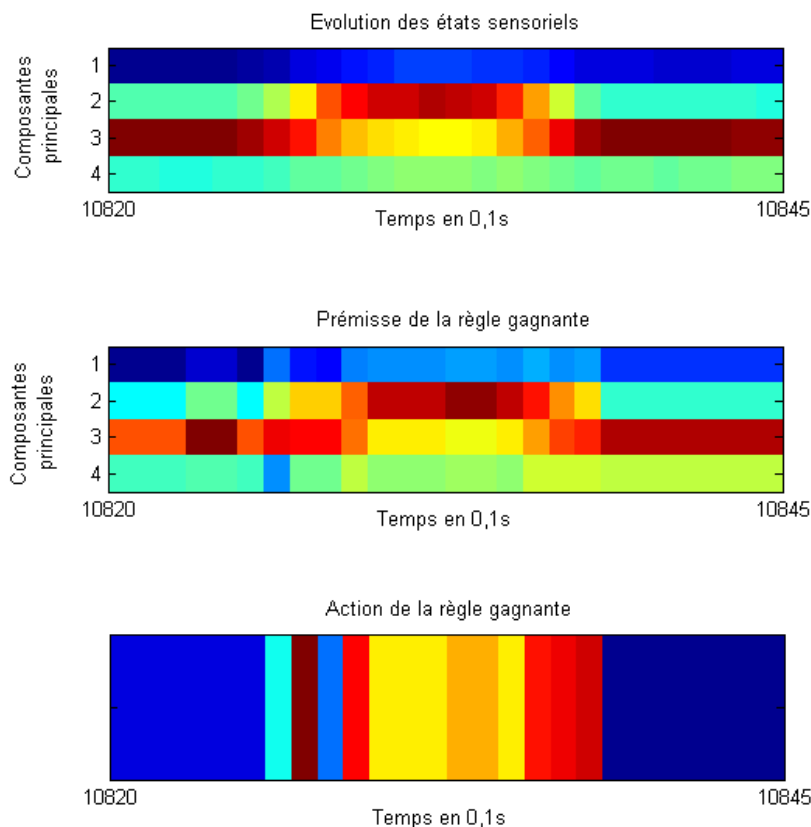


Figure IV-96 : Bandes chronologiques de l'expérience 37 coloriées en fonction des états sensoriels observés, des prémisses déclenchées et des commandes effectuées.

5.2. Influence de la complexité de l'environnement

5.2.1. Les états sensoriels observés

Précédemment, il a été évoqué que l'environnement complexe pouvait être à l'origine de la disparition des amas d'états sensoriels observés dans l'expérience 37 et 38. Mais cette hypothèse est contredite par la Figure IV-97 qui superpose les états sensoriels observés au cours des expériences 37 et 40. Les observations de l'expérience 40 indiquent la présence d'amas deux fois plus nombreux et de taille similaire. Bien que le nombre et la position des amas observés soient différents, ils restent dans les mêmes cônes de l'espace sensoriel. De la Figure IV-97 à la Figure IV-99, les figures montrent la projection des états sensoriels sur les composantes 2 et 4, qui représentent au mieux le phénomène d'amas dans l'espace sensoriel.

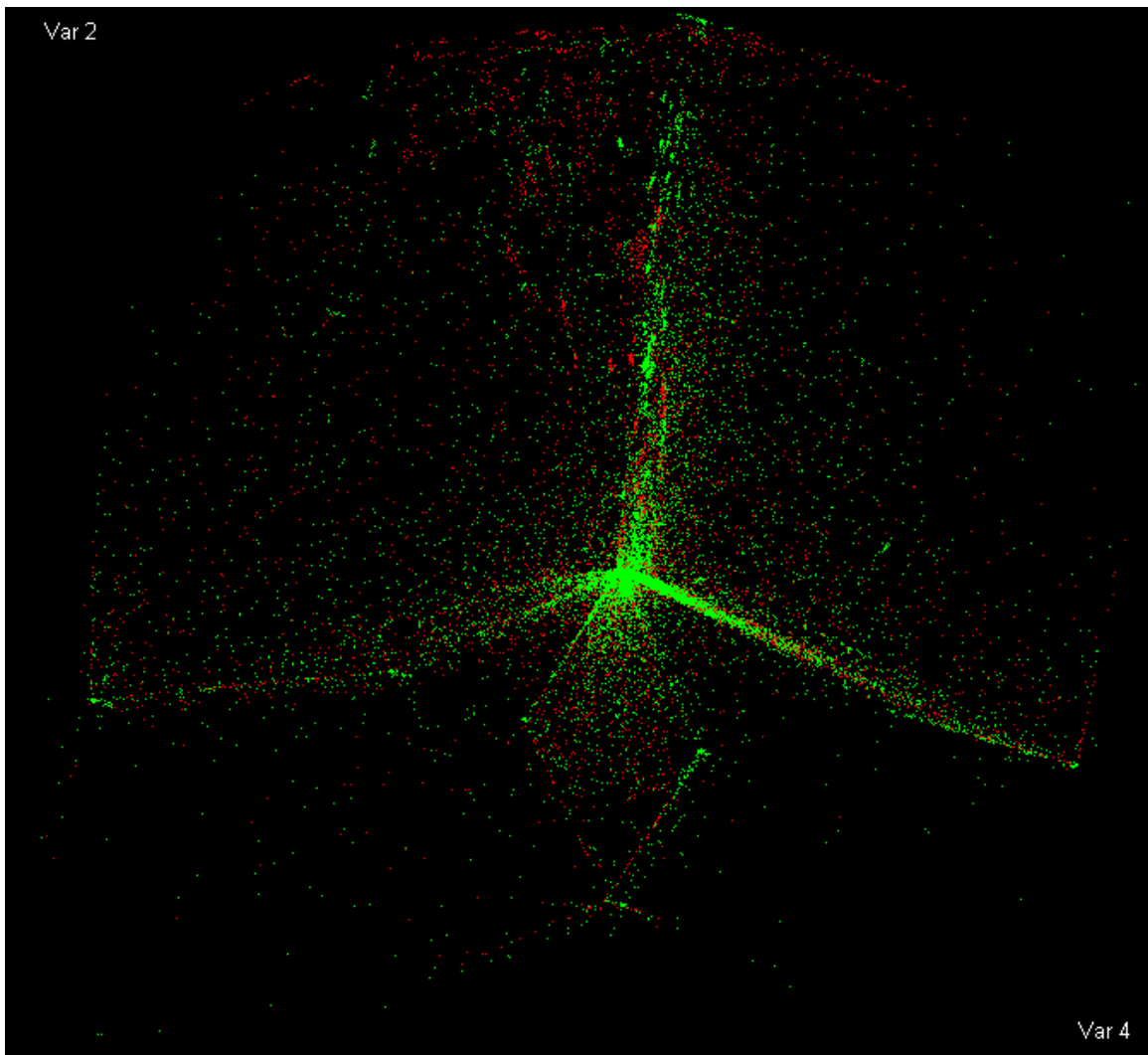


Figure IV-97 : Sur les composantes 2 et 4, représentation des états sensoriels observés au cours de l'expérience 37 (points verts) et de l'expérience 40 (points rouges).

La comparaison entre les expériences 40 et 42 de la distribution des états sensoriels observés (Figure IV-98) indique la disparition des amas au cours de cette dernière qui intervient ainsi dans les deux séries d'expériences à la troisième évolution (les expériences 42 et 39). Les états sensoriels observés au cours de l'expérience 42 dessinent davantage les arêtes de l'hypercube qui confinent tous les états sensoriels. Le dessin de ces arêtes signifie que les capteurs prennent des valeurs élevées voire maximales et par conséquent que le robot s'approche davantage des obstacles.

Dans l'expérience 43, qui débute avec les règles initiales originelles et un bruit infrarouge, ne contient pas d'amas de points sensoriels aussi nettement que dans l'expérience 40 ou 37. Bien que le bruit rajouté ne nuise pas au comportement d'évitement d'obstacle, celui-ci ne favorise pas l'établissement d'un comportement cyclique ou le déclenchement de règles générant un déplacement lent. En revanche, le regroupement de points au centre de l'axe vertical dessiné par les autres points semble être un amas dont l'amplitude correspond à celle du bruit émis (Figure IV-99). Par ailleurs, le bruit amplifiant les signaux, les arêtes de l'hypercube apparaissent davantage que lors de l'expérience 40.

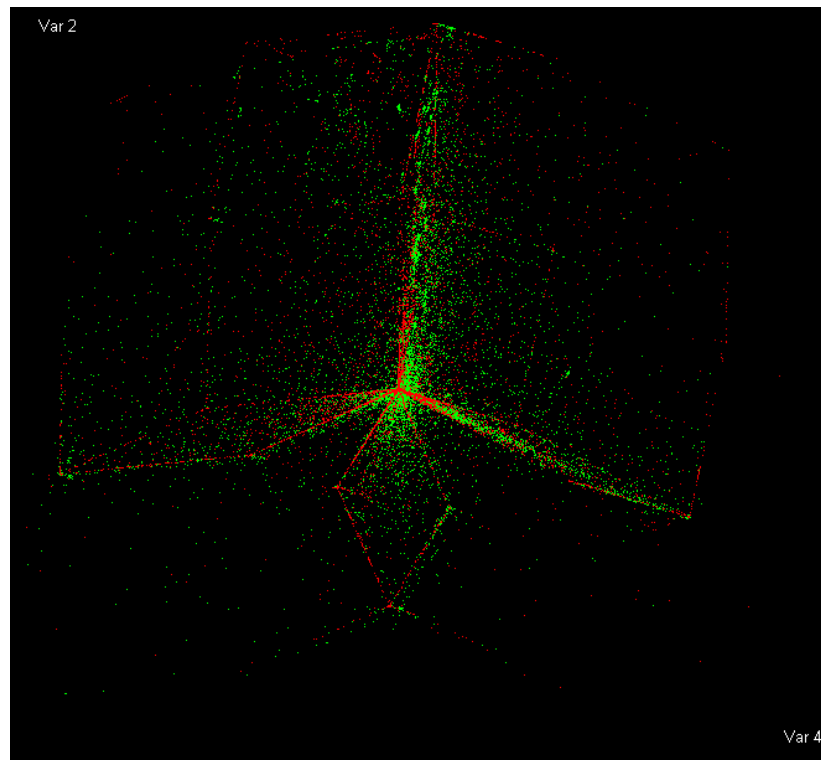


Figure IV-98 : Sur les composantes 2 et 4, représentation des états sensoriels observés au cours de l'expérience 40 (points verts) et de l'expérience 42 (points rouges).

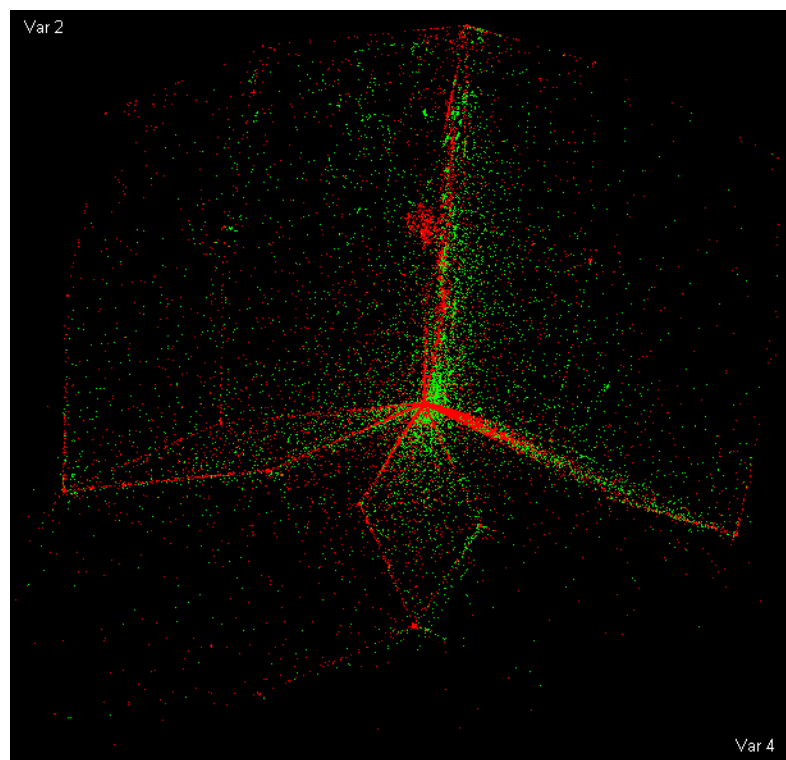


Figure IV-99 : Sur les composantes 2 et 4, représentation des états sensoriels observés au cours de l'expérience 40 (points verts) et de l'expérience 43 (points rouges).

5.2.2. La conservation du comportement

Le comportement d'évitement d'obstacle est conservé dans les quatre expériences : 40, 41, 42 et 43. Dans les expériences 41 et 42, le robot s'approche davantage des obstacles et ses actions deviennent plus sèches. La répartition des consignes dans l'espace des commandes (Figure IV-100) ressemble à celle observée lors de la première série d'expériences : 37, 38 et 39.

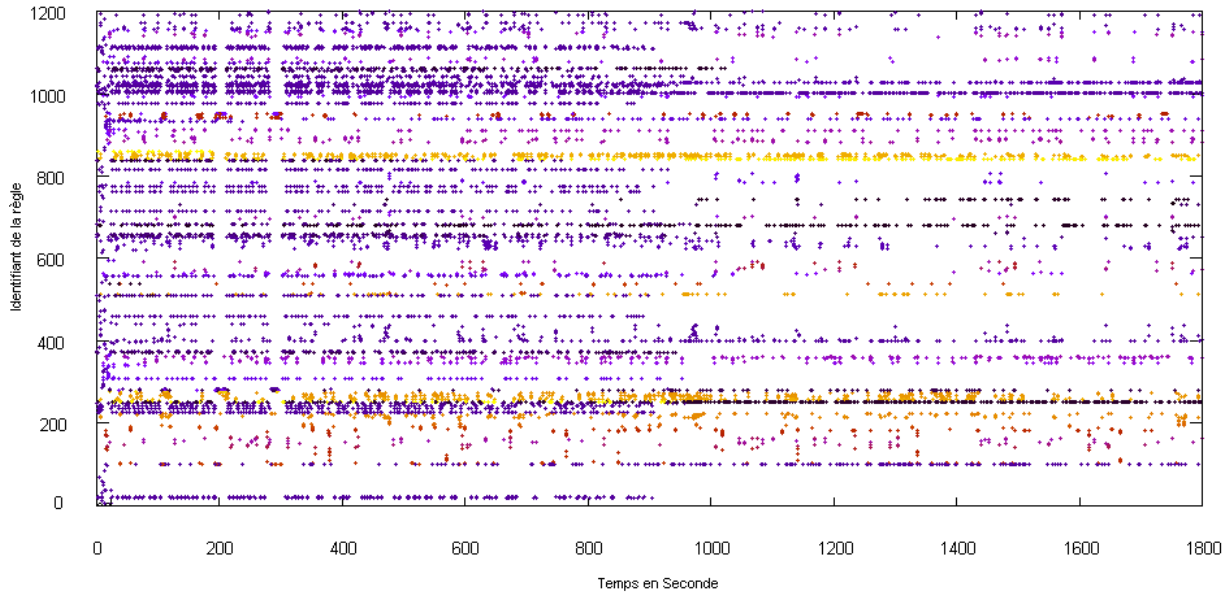


Figure IV-100 : Graphique indiquant le numéro de la règle gagnante à chaque élection de l'expérience 40 : la couleur des points est fonction de la direction de l'effort moteur.

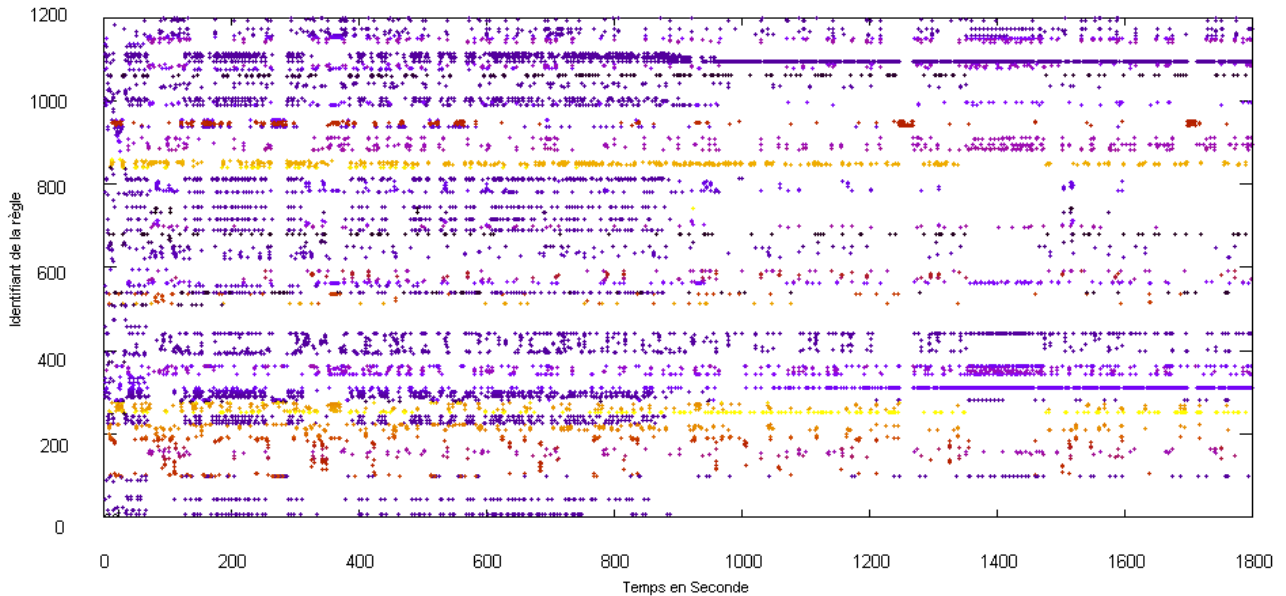


Figure IV-101 : Graphique indiquant le numéro de la règle gagnante à chaque élection de l'expérience 40 : la couleur des points est en fonction de la direction de l'effort moteur.

La Figure IV-100 qui affiche l'indice de la règle déclenchée et l'orientation de l'action décrit une transition moins brusque que celle observée au cours de l'expérience 37 (Figure IV-78). Le début de cette transition se trouve décalé dans le temps de 250 s environ, et ce décalage se retrouve également dans l'expérience 43. Par ailleurs, sur la Figure IV-101, une douzaine de règles se sont déclenchées successivement pendant 160 s. Le mouvement rotatif induit durant cette période explique ainsi l'amas central remarqué dans la Figure IV-99.

Par la suite, dans les expériences 41 et 42, le déclenchement des règles présente une certaine régularité, comme pour les expériences 38 et 39.

5.2.3. L'évolution des populations de règles

La complexification de l'environnement n'a pas altéré le comportement du robot qui évite toujours les obstacles. En revanche, cette complexification, soit par le bruit soit par l'ajout d'éléments, a eu pour conséquence de modifier la dynamique du système. La Figure IV-102 montre que la deuxième transition souffre d'un décalage de 350 s entre l'expérience 37 et l'expérience 40 ou l'expérience 43. Ce décalage s'explique par un grand nombre de configurations sensorielles différentes ou du moins par une meilleure uniformisation de leurs fréquences d'apparition qui aurait pour effet de diminuer la fréquence de l'enchère. Une différence concernant la pente de la seconde transition apparaît entre les expériences 40 et 43. Le bruit provoque le déclenchement d'un ensemble de règles plus grand et de manière très variable. Cette disparité doit alors introduire des perturbations dans la dynamique de la compétition et ainsi désynchroniser l'élimination des règles. À la fin de ces trois expériences, le nombre de règles restantes ne diffère que de deux ou trois règles.

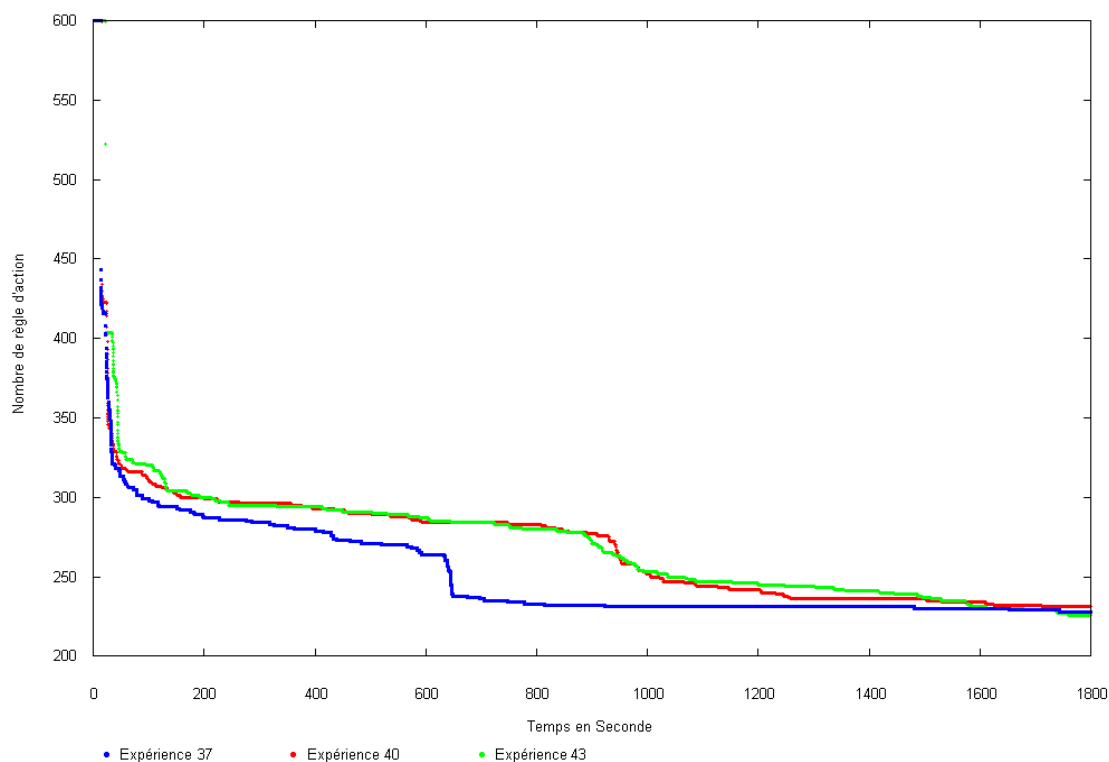


Figure IV-102 : Évolution du nombre de règles d'action au cours de l'expérience 37, 40 et 43.

En dehors de ce décalage dans l'expérience 40, l'évolution de la démographie des expériences suivantes (41 et 42) reste similaire à celle des expériences 38 et 39 (Figure IV-103).

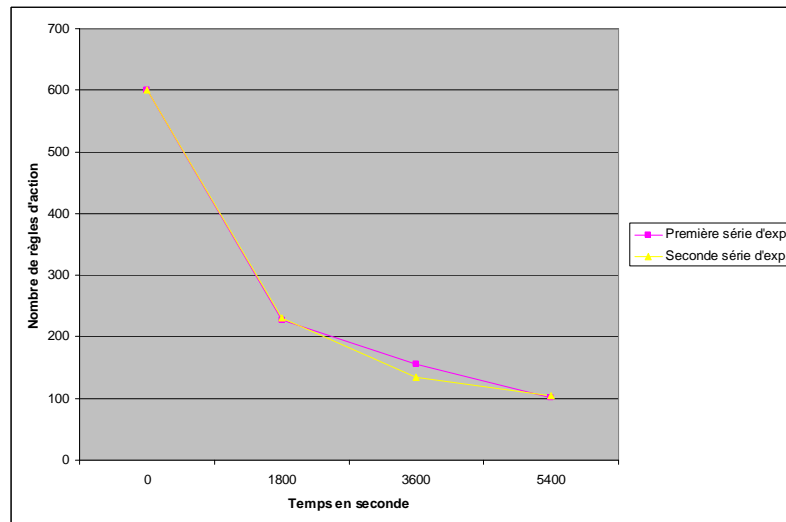


Figure IV-103 : Évolution du nombre de règles d'action finales au cours de la première série d'expériences (37, 38, 39) et de la seconde série d'expériences (40, 41 et 42).

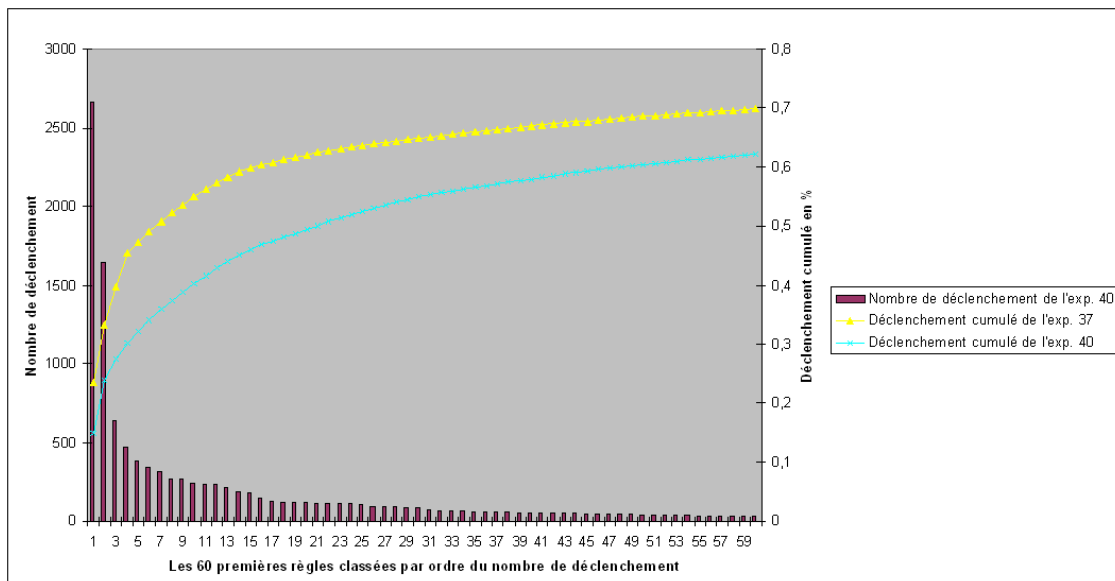


Figure IV-104 : Histogramme du nombre de déclenchements des 60 règles les plus activées et la courbe du pourcentage cumulé de leur participation au déclenchement global.

L'historgramme du nombre de déclenchements des 60 règles les plus activées (Figure IV-104) suggère toutefois une différence qualitative plutôt que quantitative. Dans le cas de l'environnement simple, les 60 premières règles couvrent 70% des déclenchements contre 60% dans le cas de l'environnement complexe. En d'autres termes, l'homogénéisation de la distribution des états sensoriels provoquée par la complexification de l'environnement se répercute dans la répartition des règles déclenchées.

La différence du nombre de déclenchements moyens se retrouve également entre les expériences 39 et 42. En enlevant les deux règles à fort taux de déclenchements, la moyenne de déclenchement d'une règle est de 34 avec un écart type de 113 à la fin de l'expérience 39 et une moyenne de 58 avec un écart type de 86 à la fin de l'expérience 42.

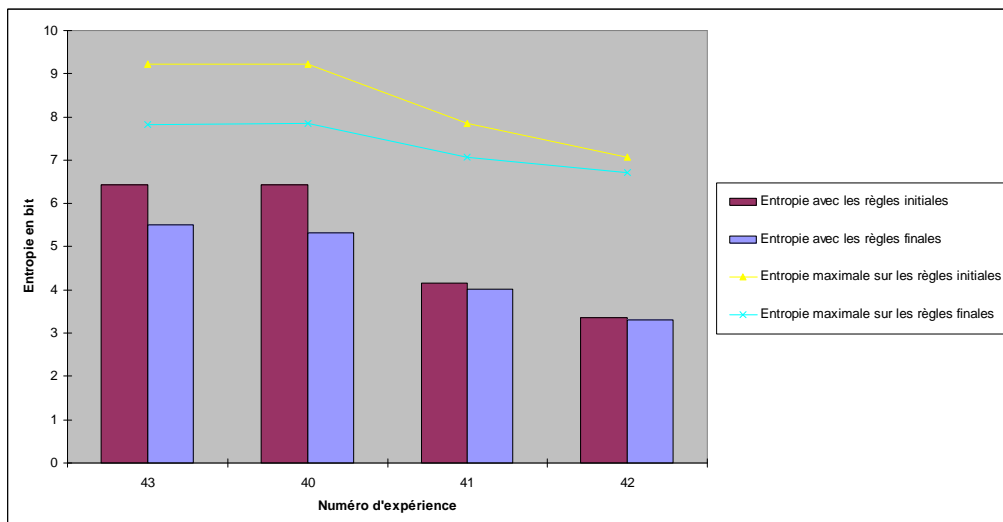


Figure IV-105 : Évolution de l'entropie de l'ensemble des règles suivant les expériences 43, 40, 41 et 42.

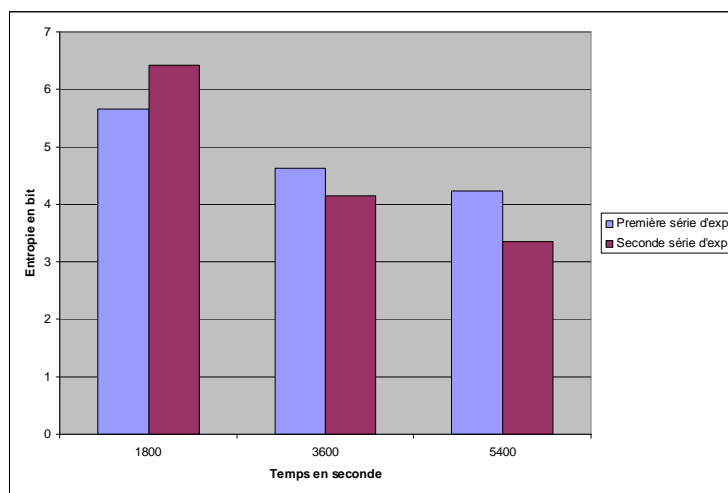


Figure IV-106 : Évolution de l'entropie au cours de la première série d'expériences (37, 38, 39) et de la seconde série d'expériences (40, 41 et 42).

La présence du bruit dans l'expérience 43 n'augmente pas l'entropie des règles finales comparée à celle de l'expérience 40. La prise en compte des quatre premières composantes principales uniquement explique certainement cette indifférence puisque le bruit se situe dans les dernières composantes calculées par l'ACP. La diminution de l'entropie, visible sur la Figure IV-105 et la Figure IV-106, présente une pente de -1,5 alors qu'elle est de -0,71 pour la première série d'expériences. Une interprétation serait de considérer qu'une stimulation sensorielle variée structure plus rapidement l'espace sensoriel décrit par les règles. Par exemple, lors de la première série d'expériences, la compétition entre 6 règles pour la représentativité de l'espace sensoriel le plus dense a duré jusqu'à la fin de la troisième

expérience. En revanche, pour la seconde série, la compétition pour cette même zone n'a duré que le temps de la première expérience. La règle dominante au terme de l'expérience 41 s'est déclenchée plus de 8000 fois, soit environ la somme des déclenchements des règles en compétition pour ce sous-espace dans l'expérience 40.

A la fin de l'expérience 40, seules 15 règles ne se sont pas déclenchées et 50 se sont déclenchées moins de 5 fois au lieu de respectivement 32 et 35 fois pour l'expérience 37. Par la suite, dans les expériences 41 et 42, toutes les règles se sont déclenchées plus d'une fois, contrairement aux expériences 38 et 39. Autrement dit, l'environnement simple a permis une augmentation de la force de certaines règles suffisante pour subir la taxe pendant 30 min sans passer en dessous du seuil, et leur non déclenchement s'explique par la redondance de ces règles avec d'autres qui, elles, s'activent.

5.2.4. Les dynamiques

Les évolutions des caractéristiques de la règle déclenchée lors de l'expérience 40 ressemblent aux évolutions constatées dans l'expérience 39, avec un décalage temporel d'environ 350 s (Figure IV-107). Ce décalage correspond à celui observé sur la Figure IV-102 concernant l'évolution du nombre de règles.

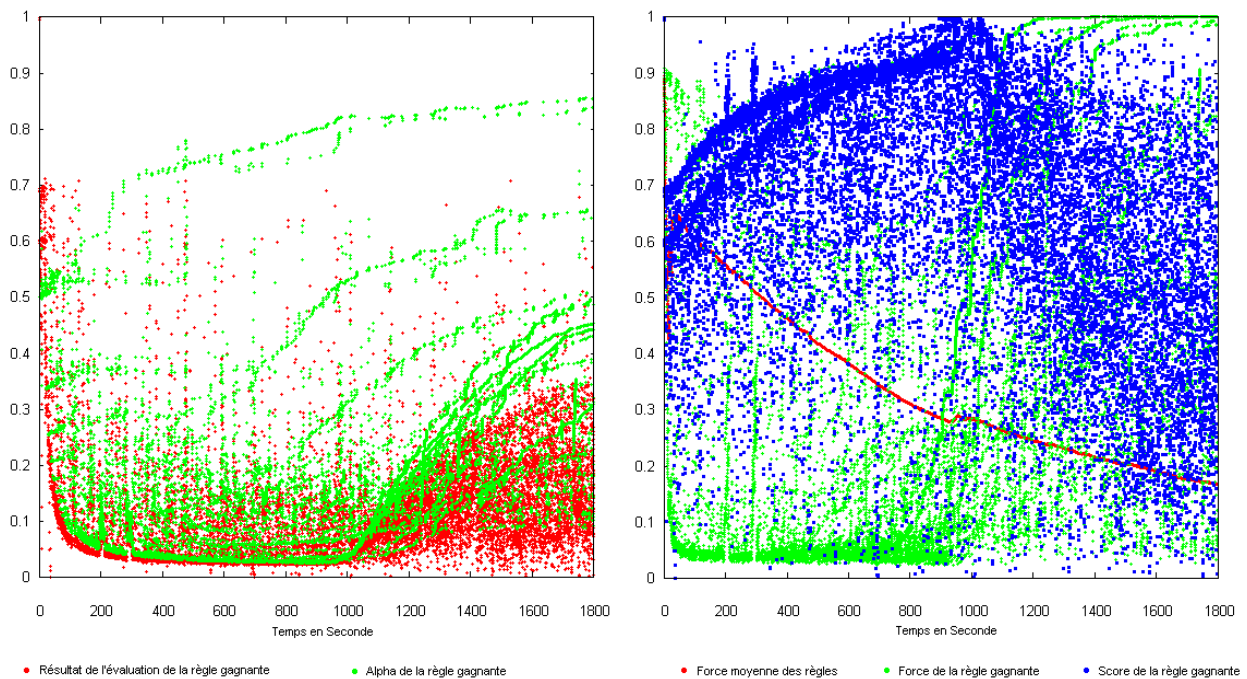


Figure IV-107 : Représentation graphique à droite de l'alpha et du résultat de l'évaluation de la règle déclenchée, à gauche représentation de la force et du score de la règle déclenchée ainsi que la force moyenne des règles de l'expérience 40.

Les expériences 41 et 42 présentent une évolution linéaire (Figure IV-108). À la fin de l'expérience 42, la force moyenne atteint, 0,3, soit 0,1 de plus qu'à la fin de l'expérience 39. De nombreuses règles possèdent à la fin de l'expérience une force élevée ou maximale comme dans la première série d'expériences. Toutefois, la richesse de l'environnement généralise cette augmentation de la force à davantage de règles.

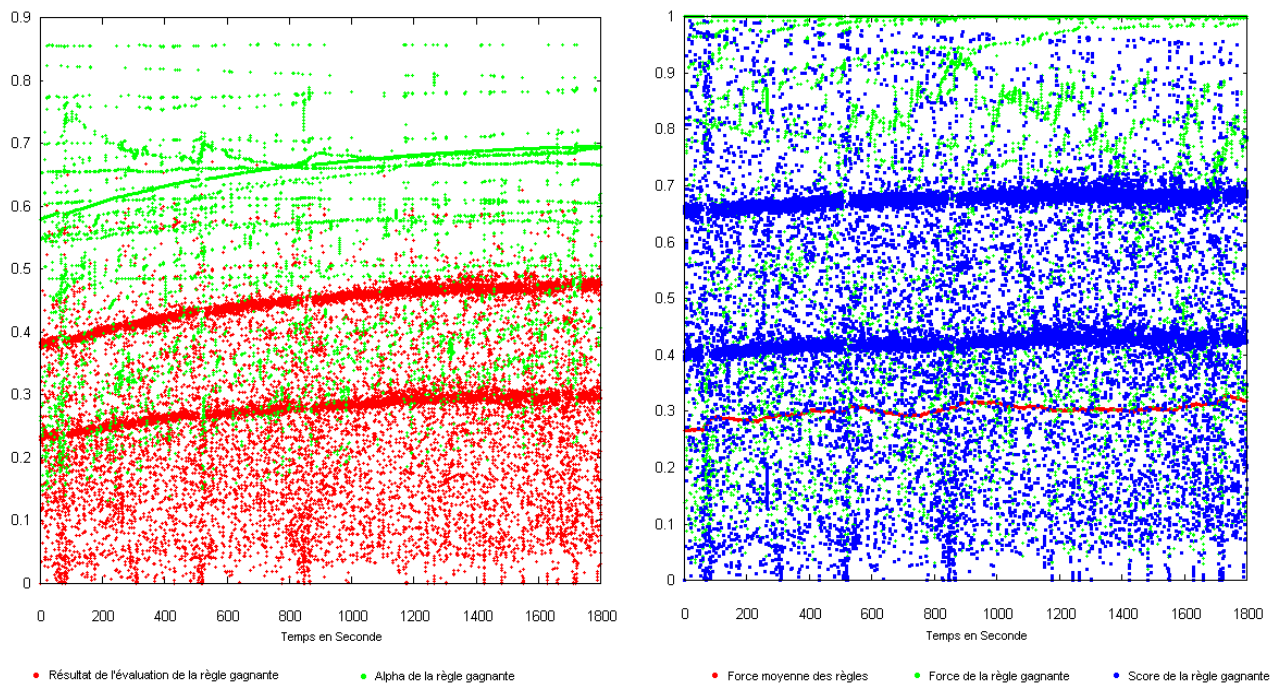


Figure IV-108 : Représentation graphique à droite de l'alpha et du résultat de l'évaluation de la règle déclenchée, à gauche représentation de la force et du score de la règle déclenchée ainsi que la force moyenne des règles de l'expérience 42.

Le décalage temporel observé dans l'expérience 40 se retrouve dans l'expérience 43. Le bruit n'affecte pas significativement l'allure des différentes variables (Figure IV-109).

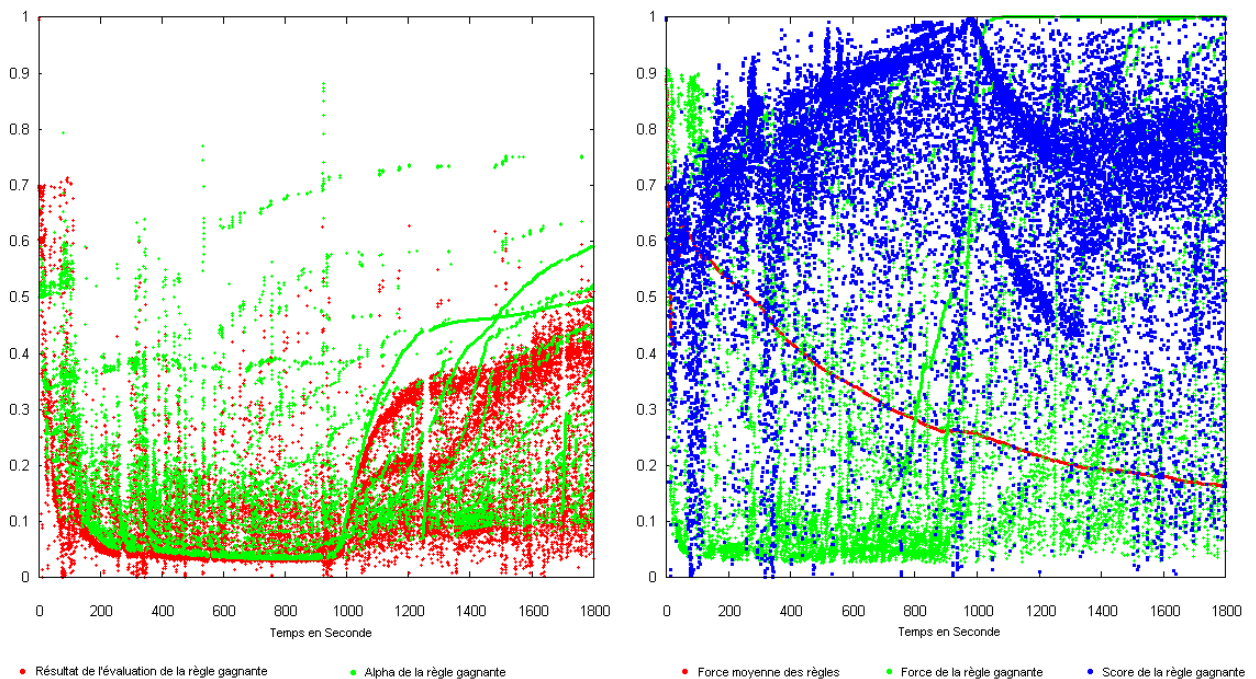


Figure IV-109 : Représentation graphique à droite de l'alpha et du résultat de l'évaluation de la règle déclenchée, à gauche représentation de la force et du score de la règle déclenchée ainsi que la force moyenne des règles de l'expérience 43.

5.2.5. Le positionnement des règles finales

A l'issue de l'expérience 41, la répartition des prémisses des règles ressemble à celle observée à la fin de l'expérience 38, soit au terme de 36000 déclenchements. Toutefois, la répartition des prémisses finales de l'expérience 41 semble plus symétrique et homogène que celle de l'expérience 38 (Figure IV-110). En revanche, la comparaison entre les expériences 42 et 39 ne montre pas les mêmes signes de distinction (Figure IV-111).

La désertion de certaines zones de l'espace sensoriel de règles finales entre les expériences 38 et 39 se retrouve également entre les expériences 40 et 42. Par ailleurs, la distinction entre les règles finales des premières expériences et celles des troisièmes expériences pour les deux séries d'expériences demeure perceptible par l'élimination de certaines règles. Cette ressemblance (Figure IV-110 et Figure IV-111) se trouve mise en évidence par la visualisation des couples de variantes suivantes : var4 et var1 ; var4 et var2 ; var4 et var3 ; var1 et var2 ; var1 et var3.

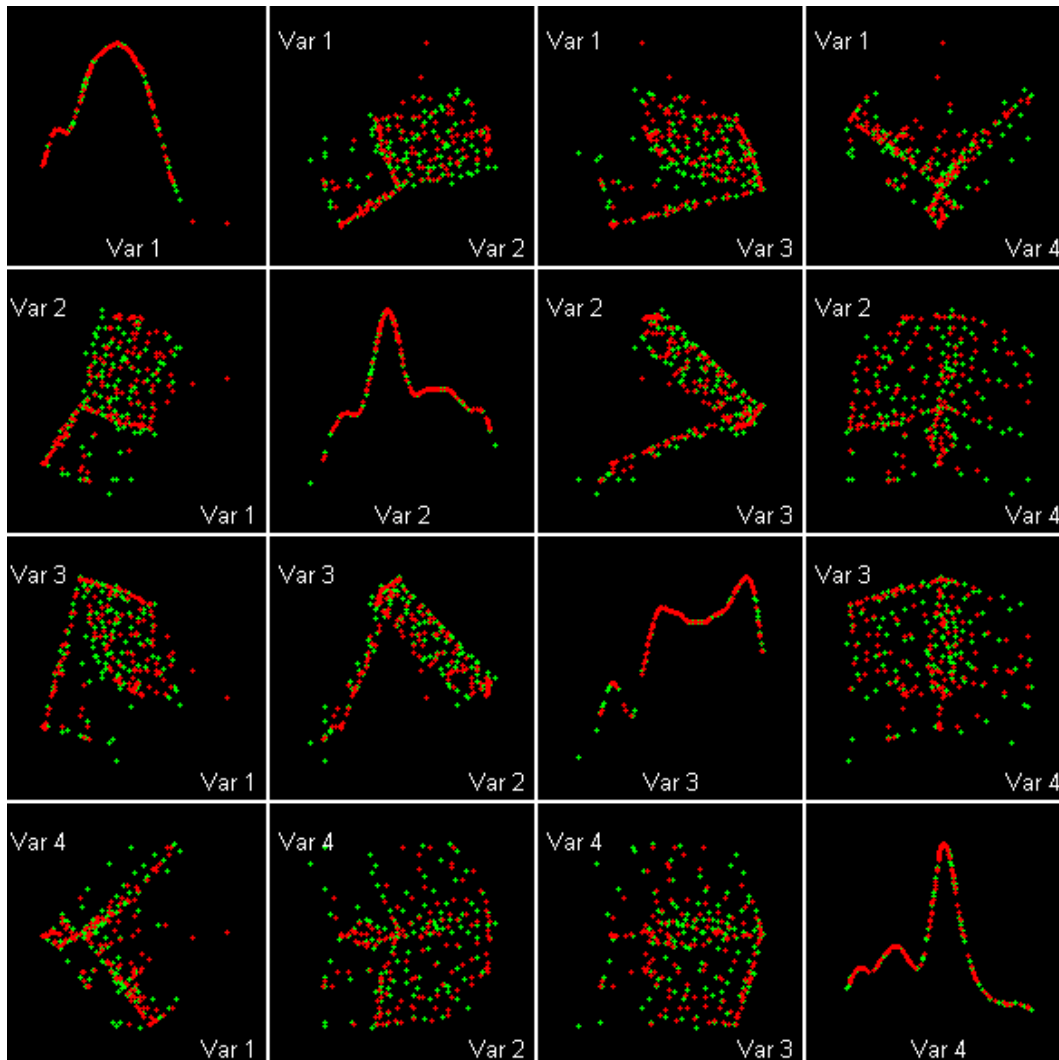


Figure IV-110 : Représentation des prémisses appartenant aux règles finales de l'expérience 41 (croix vertes) et celles de l'expérience 38 (croix rouges).

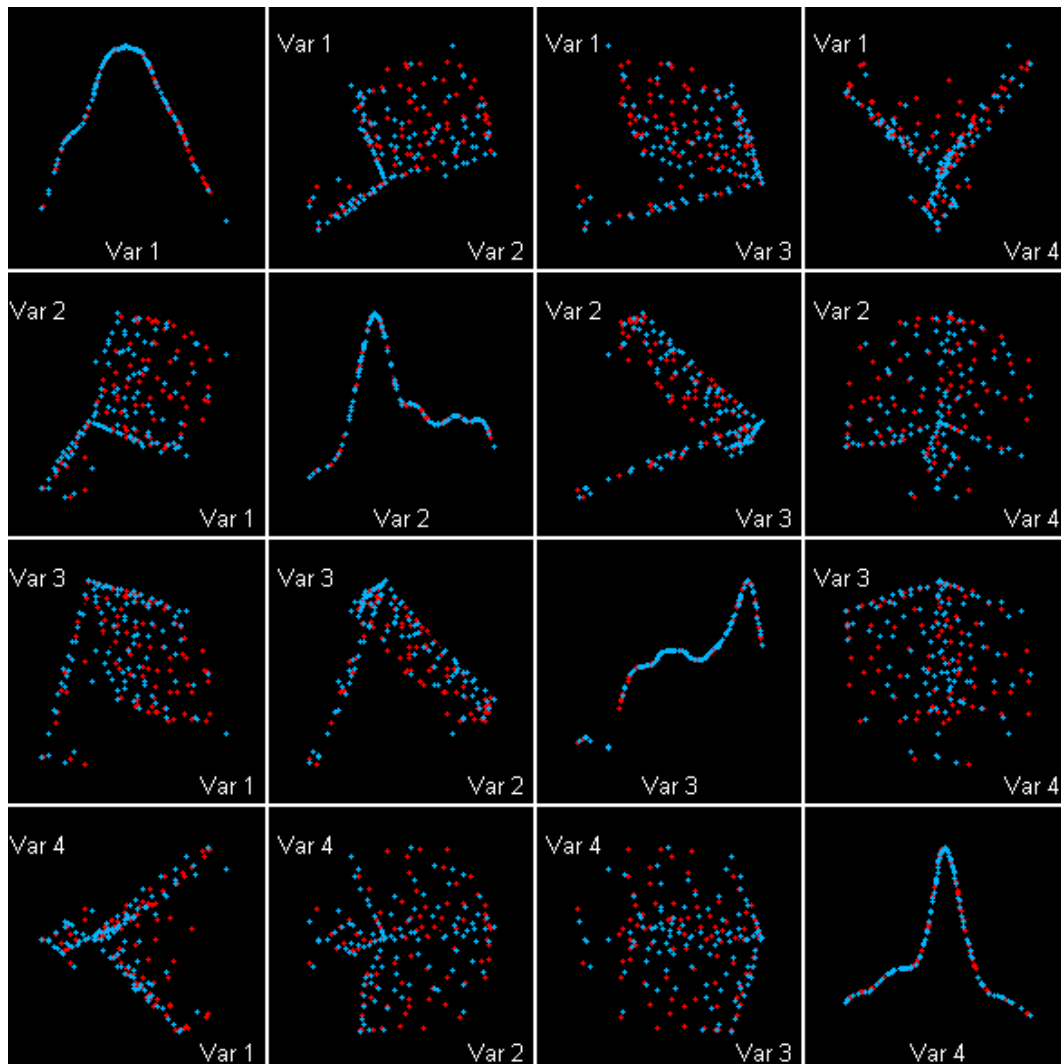


Figure IV-111 : Représentation des prémisses appartenant aux règles finales de l'expérience 42 (croix vertes) et celles de l'expérience 39 (croix rouges).

L'élimination de ces règles provient du manque d'états sensoriels dans leur champ de réception comme l'illustre la Figure IV-112. Cette réduction de l'espace sensoriel exploré montre que l'auto-organisation des règles résulte d'un couplage sensorimoteur et non uniquement sensoriel.

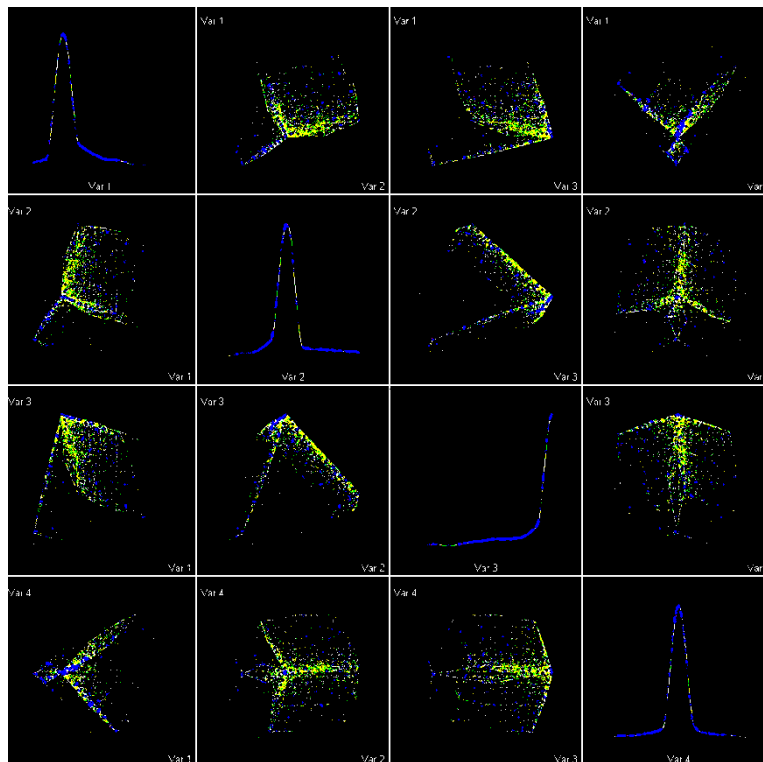


Figure IV-112 : Représentation des prémisses appartenant aux règles finales (ronds bleus) et les états sensoriels observés (points verts) de l'expérience 42.

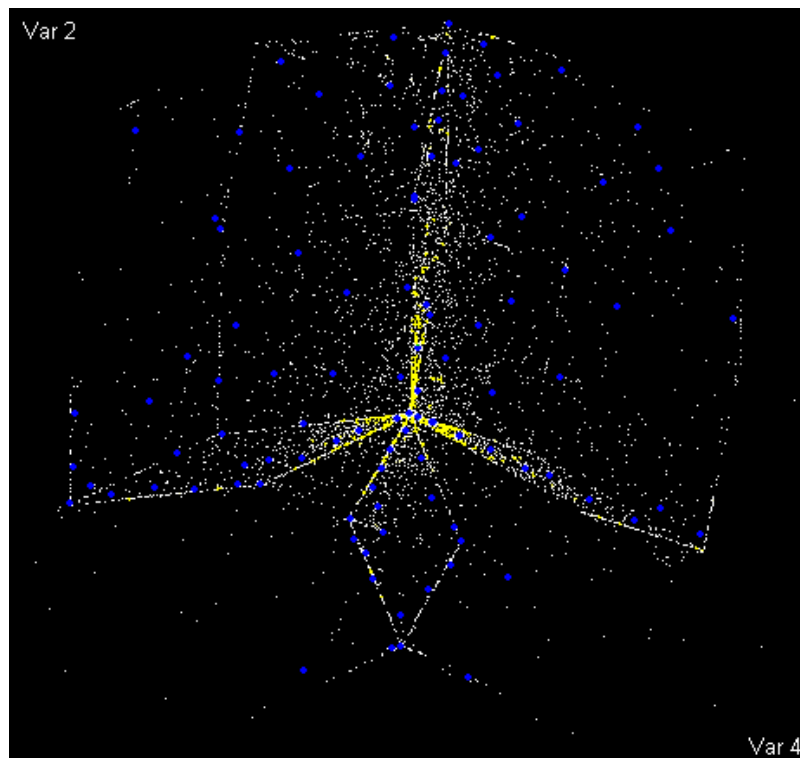


Figure IV-113 : Représentation des prémisses appartenant aux règles finales (ronds bleus) et les états sensoriels observés (points verts) de l'expérience 42 selon les composantes 2 (Var2) et 4 (Var4).

Les prémisses se situent préférentiellement sur les arêtes et les sommets de l'hypercube de l'espace sensoriel dessiné par les états observés (Figure IV-113).

Afin de mettre en relation le déclenchement des règles et la situation du robot, une expérience supplémentaire a été filmée pendant 30 s. Ce délai relativement court provient des difficultés à maintenir une juste synchronisation entre le déroulement du programme et l'acquisition. Toutefois, la Figure IV-114 représente la position du robot à chaque déclenchement par une flèche dont la taille correspond au diamètre de celui-ci et dont la pointe symbolise l'avant du robot. Les flèches rouges indiquent le déclenchement de la règle 399. La localisation dans l'enceinte du robot lors du déclenchement de cette règle montre que le système a pu identifier une configuration sensorielle spécifique à son environnement. Ainsi, il est possible d'attribuer une sémantique à la troisième personne à la règle 399 au-delà de sa simple position dans l'espace sensoriel, comme l'indique la Figure IV-115 par une croix rouge.

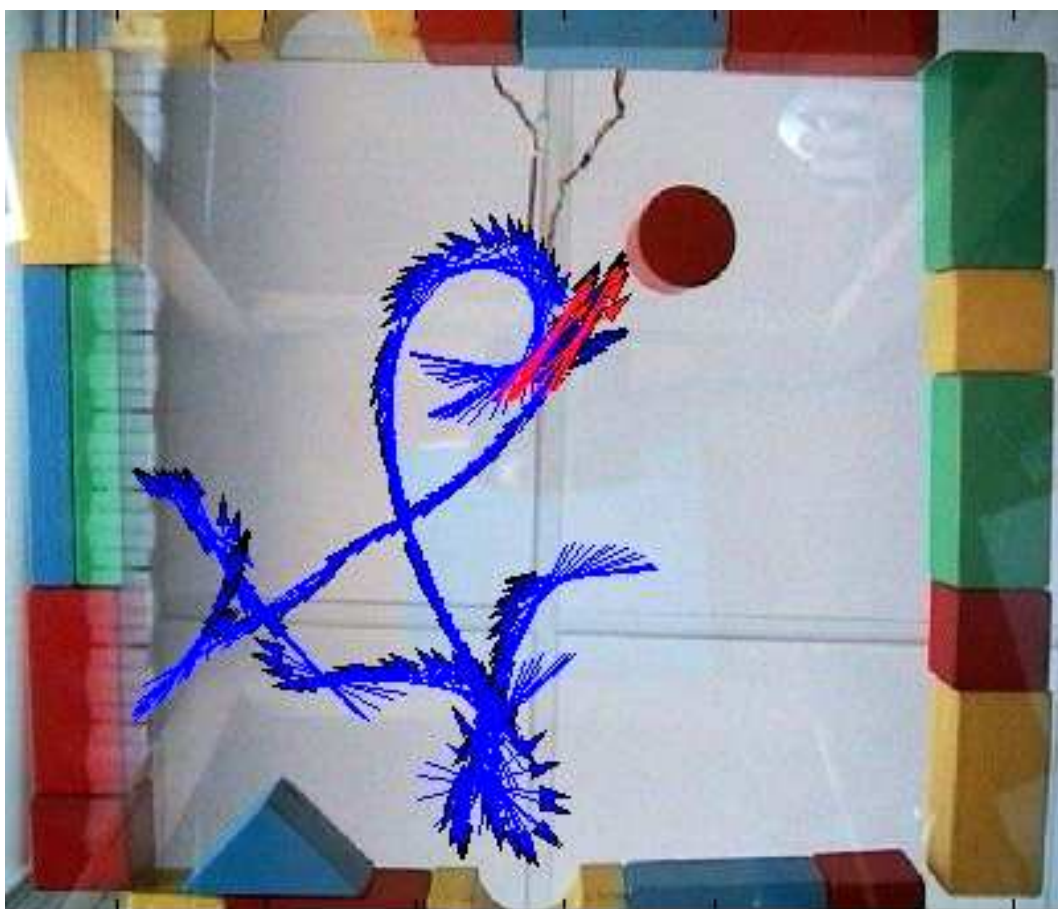


Figure IV-114 : Visualisation à l'aide de flèches des positions du robot lors d'une séquence de 30s, les flèches rouges indiquent le déclenchement de la règle 399.

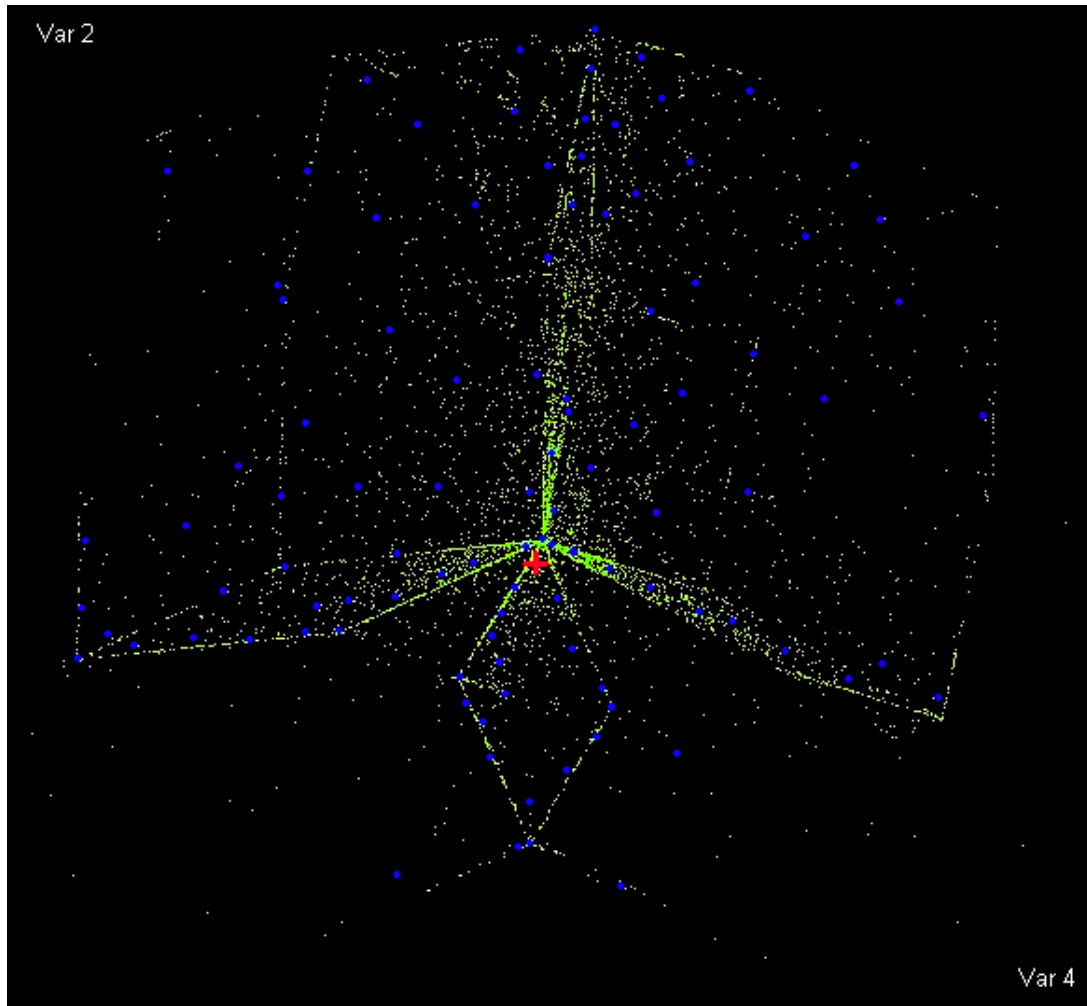


Figure IV-115 : Représentation des prémisses appartenant aux règles finales (ronds bleus) ainsi que les états sensoriels observés (points verts) de l'expérience 42 et la croix rouge correspond à la prémisse de la règle 399 selon les composantes 2 (Var2) et 4 (Var4).

5.3. Extension de la base de règles par création spontanée

Les règles initiales utilisées lors de la troisième série d'expériences sont les règles finales de l'expérience 38. Dès le début des expériences, des règles sont créées à une fréquence de 0,02 Hz, soit toutes les 50 s. Les nouvelles règles sensorimotrices possèdent un message sensoriel en prémisse et un message moteur en conclusion. Les règles de gestion associées ont leur message de conclusion avec une étiquette temporelle à 1, comme les règles de gestion des règles initiales. La troisième série d'expériences se compose des expériences 44 et 45. La première expérience se déroule dans l'environnement complexe et la seconde dans l'environnement simple mais bruité par la présence d'une lampe émettant des infrarouges. Les règles initiales des expériences 44 et 45 correspondent aux règles finales de l'expérience 38, soit 156 règles.

Pour ces deux expériences, les dynamiques observées à travers les caractéristiques de la règle déclenchée ressemblent à celles observées dans l'expérience 39. La création régulière de règles surcompense le nombre de règles éliminées ; ainsi le nombre de règles finales atteint 236 pour l'expérience 44 et 219 pour l'expérience 45 (Figure IV-116). À titre

de comparaison, l'expérience 39, qui est identique à l'expérience 44 mais sans création de règles, possède 102 règles au terme de celle-ci.

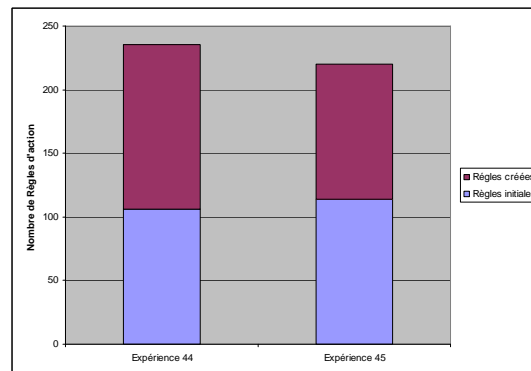


Figure IV-116 : Nombre de règles d'action restantes à la fin des expériences 44 et 45.

Dans l'expérience 39, 77% des règles finales se sont déclenchées plus de 20 fois. Dans les expériences 44 et 45, les règles créées sont initialisées avec une force à $5 \cdot 10^{-1}$, un alpha à 1, une incertitude $1 \cdot 10^{-2}$ pour les états sensoriels et $1 \cdot 10^{-4}$ pour l'étiquette temporelle. Comme pour l'environnement imposé, ces trois derniers paramètres favorisent le déclenchement des règles nouvellement créées. La Figure IV-117 et la Figure IV-118 offrent un moyen de comparaison entre les expériences 44 et 45, selon un critère de pertinence basé sur le nombre de déclenchements d'une règle.

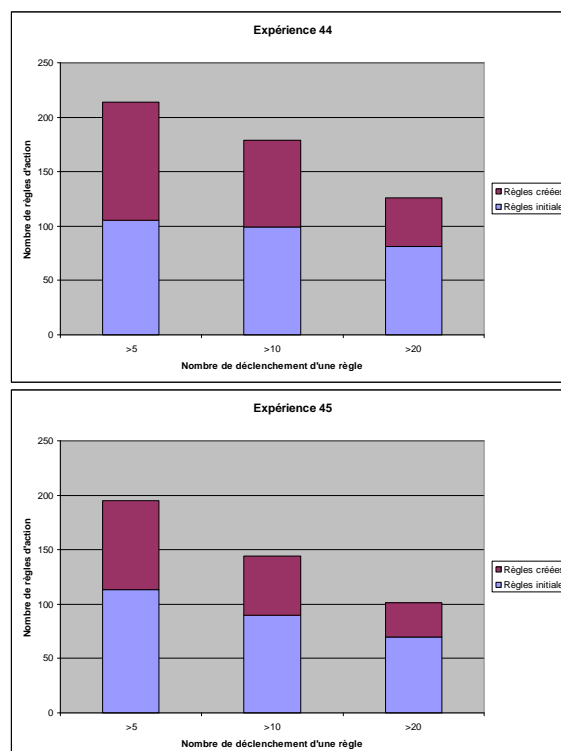


Figure IV-117 : Comparaison entre les expériences 44 et 45 du nombre de règles finales suivant le nombre minimal de déclenchements d'une règle pour la prise en compte de celle-ci.

La Figure IV-117 exprime une équivalence quantitative alors que la Figure IV-118 révèle une différence d'ordre qualitative. En effet, les règles qui se sont déclenchées plus de 20 fois participent à 82,5% au nombre total de déclenchements dans l'expérience 44 contre 45% dans l'expérience 45. Cette différence s'explique par la présence du bruit qui favorise le déclenchement d'un plus grand nombre de règles pour une seule configuration topologique du robot. La diminution du taux de participation n'est pas équivalente entre celle des règles initiales (un facteur de 1,7) et celle des règles créées (un facteur de 2,7). L'instabilité nuit davantage aux nouvelles règles qu'aux anciennes en terme de pertinence liée au nombre de déclenchements.

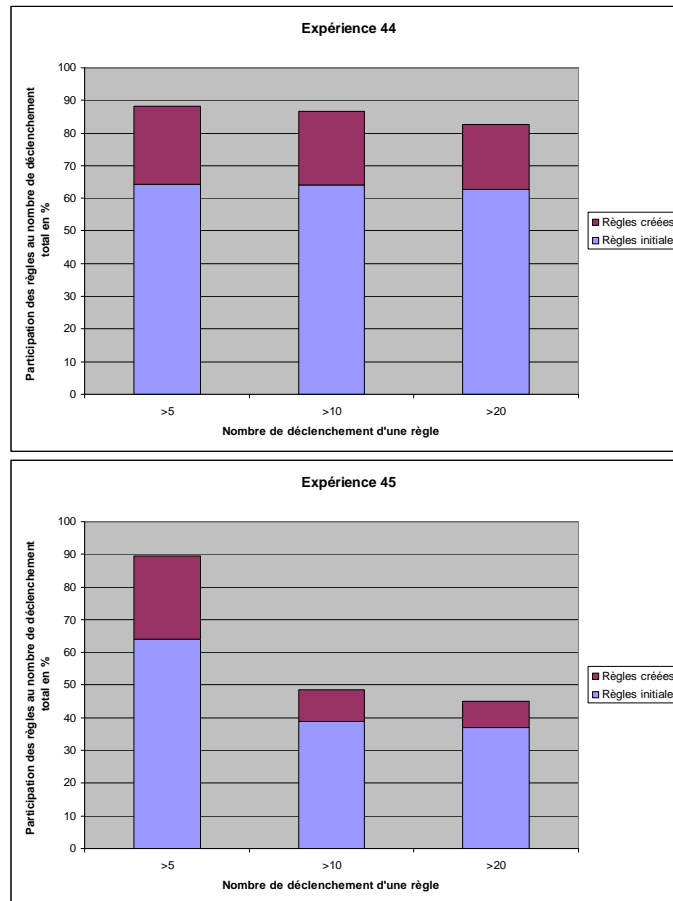


Figure IV-118 : Comparaison entre les expériences 44 et 45 du pourcentage de déclenchements des règles finales par rapport au nombre de déclenchements total, suivant le nombre de déclenchement minimal d'une règle pour la prise en compte de celle-ci.

La Figure IV-119 représente l'âge de la règle déclenchée au cours des expériences 44 et 45 ainsi que l'âge moyen des règles. Cette visualisation confirme la propension des nouvelles règles à se déclencher.

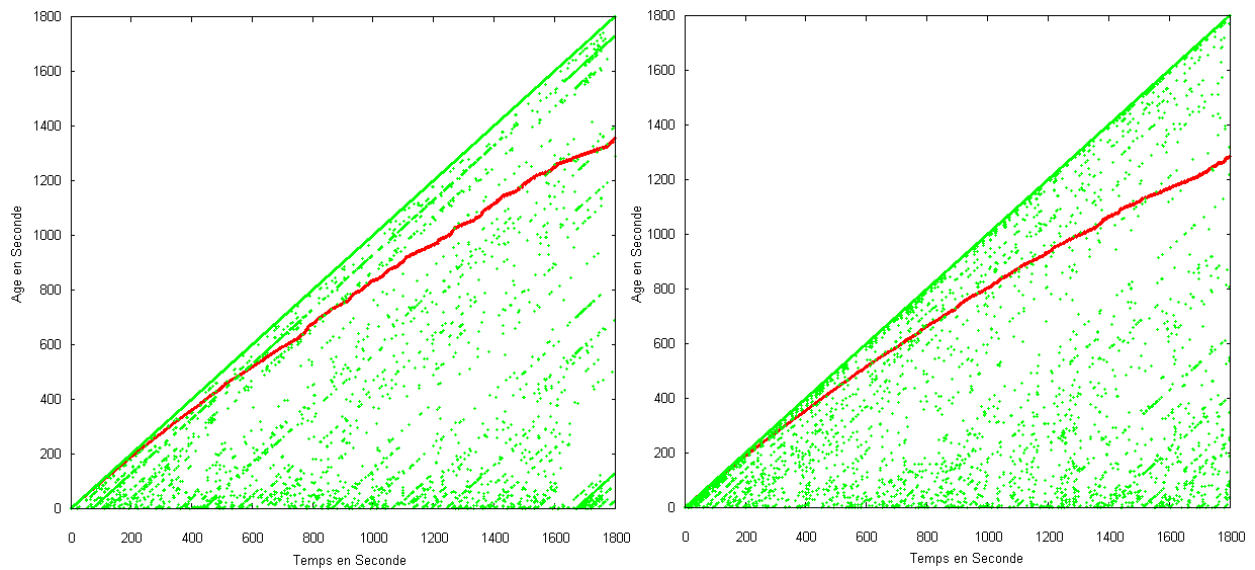


Figure IV-119 : Représentation de l'âge de la règle déclenchée (croix verte) et la moyenne d'âge des règles ; à droite pour l'expérience 44 et à gauche pour l'expérience 45.

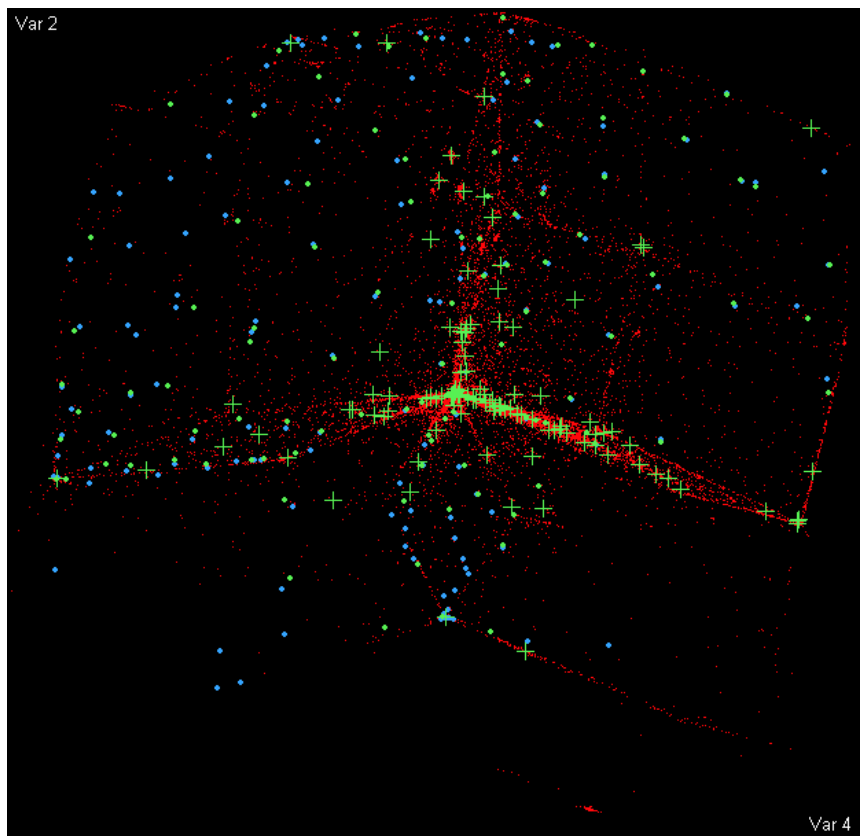


Figure IV-120 : Représentation des prémisses appartenant aux règles finales issues des règles initiales (ronds verts), les règles initiales (ronds bleus) et des règles créées déclenchées plus de 5 fois (croix vertes) et les états sensoriels observés (points rouges) de l'expérience 44 selon les composantes 2 (Var2) et 4 (Var4).

La Figure IV-120 et la Figure IV-121 donnent un aperçu de la répartition des règles nouvelles symbolisées par des croix vertes. La première figure représente les règles dont le

nombre de déclenchements est supérieur à 5, la seconde représente celles dont le nombre de déclenchements est supérieur à 20. Dans les deux cas, la majorité des règles créées et persistantes se trouve dans le sous-espace sensoriel à forte densité d'états sensoriels observés. Malgré tout, les règles initiales restent présentes au sein de ce sous-espace sensoriel, ce qui signifie que la compétition demeure perpétuelle à cause d'une stimulation sensorielle trop intense.

Des règles finales se positionnent sur des amas éloignés de l'origine des règles initiales, notamment les trois croix situées en haut au centre de la Figure IV-120 et les deux croix à droite. Hormis les règles créées se trouvant dans la partie à forte concentration d'états sensoriels, il n'apparaît pas de règles créées supplantant des règles initiales. La création de règles a permis ici de coloniser à la fois des zones précises entourées de règles initiales moins spécialisées ainsi que des zones plus larges non couvertes par celles-ci.

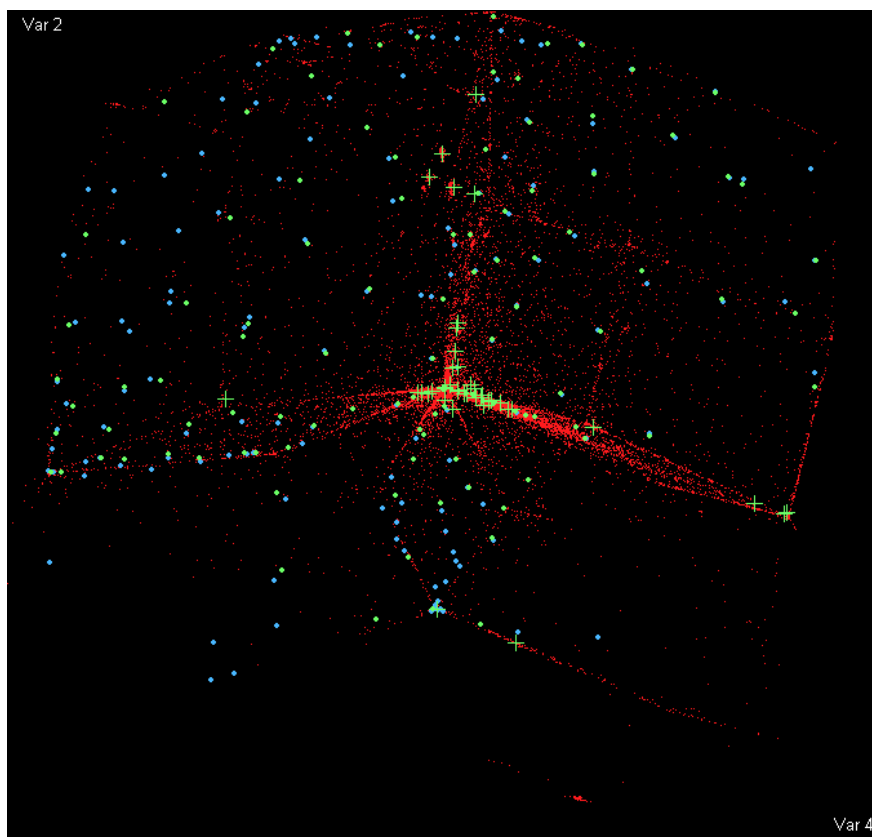


Figure IV-121 : Représentation des prémisses appartenant aux règles finales issues des règles initiales (ronds verts), les règles initiales (ronds bleus) et des règles créées déclenchées plus de 20 fois (croix vertes) et les états sensoriels observés (points rouges) de l'expérience 44 selon les composantes 2 (Var2) et 4 (Var4).

La répartition des règles créées dans l'expérience 45 reste similaire à celle de l'expérience 44. En majorité, les règles créées et persistantes se trouvent dans le sous-espace contenant le plus d'états sensoriels. Les autres règles créées se maintiennent dans les régions peu couvertes par les règles initiales (Figure IV-122 et Figure IV-123). Comme pour l'expérience 43, les états sensoriels observés lors de l'expérience 45 ne constituent pas d'amas. Cette absence d'amas s'explique par la variabilité qu'induit le bruit infrarouge bien qu'il soit en parti filtré en ne prenant que les quatre premières composantes.

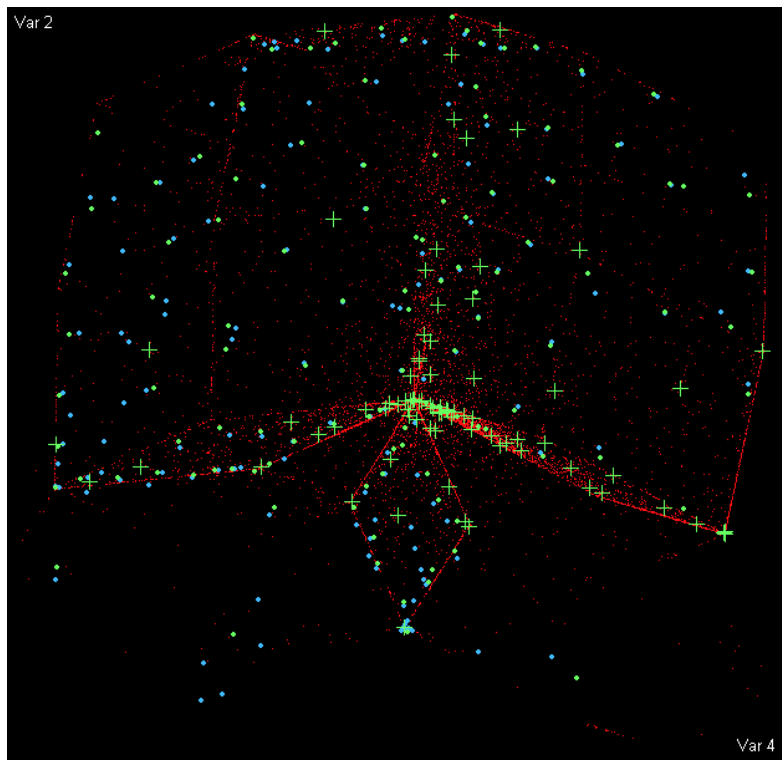


Figure IV-122 : Représentation des prémisses appartenant aux règles finales (ronds verts) issues des règles initiales (ronds bleus) et des règles créées ayant été déclenchées plus de 5 fois (croix verte) et les états sensoriels observés (points rouges) de l'expérience 45 selon les composantes 2 (Var2) et 4 (Var4).

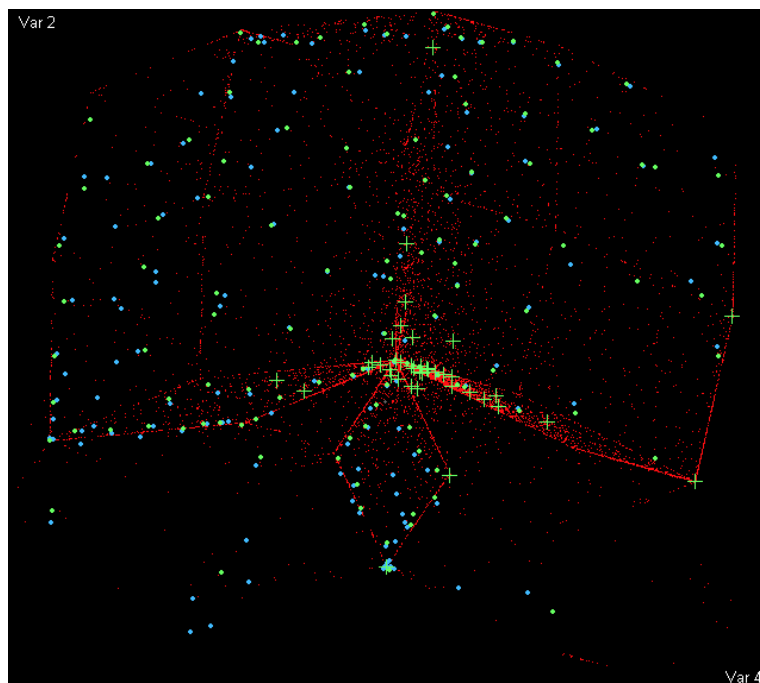


Figure IV-123 : Représentation des prémisses appartenant aux règles finales (ronds verts) issues des règles initiales (ronds bleus) et des règles créées ayant été déclenchées plus de 20 fois et les états sensoriels observés (points rouges) de l'expérience 45 selon les composantes 2 (Var2) et 4 (Var4).

En conclusion, la création de règles permet de découvrir de nouvelles règles pérennes c'est-à-dire de s'intégrer dans le couplage opérationnel. En effet, aucune règle finale de l'expérience 39 ne se situe sur le sommet de l'hypercube en bas à droite de la Figure IV-124 dessiné par les états sensoriels observés. Les règles finales issues des règles initiales n'atteignent pas non plus ce sommet, seules les règles créées y accèdent (Figure IV-121 et Figure IV-123). Par ailleurs, le fait que des règles créées se fixent au centre des amas aperçus dans l'expérience 45 montre que les nouvelles règles peuvent se spécialiser même si des règles initiales couvrent en partie cette zone.

Toutefois, dans le cas d'une forte densité d'un sous-ensemble d'états sensoriels, la création de règles génère une compétition perpétuelle qui ne favorise pas la stabilisation des règles dans ce sous-espace. Le mécanisme d'élimination des règles redondantes est trop long par rapport à la fréquence de création de règles similaires. Une solution serait alors de détecter la nouveauté pour en déduire l'opportunité de créer de nouvelles règles. La prédiction peut être l'un de ces dispositifs.

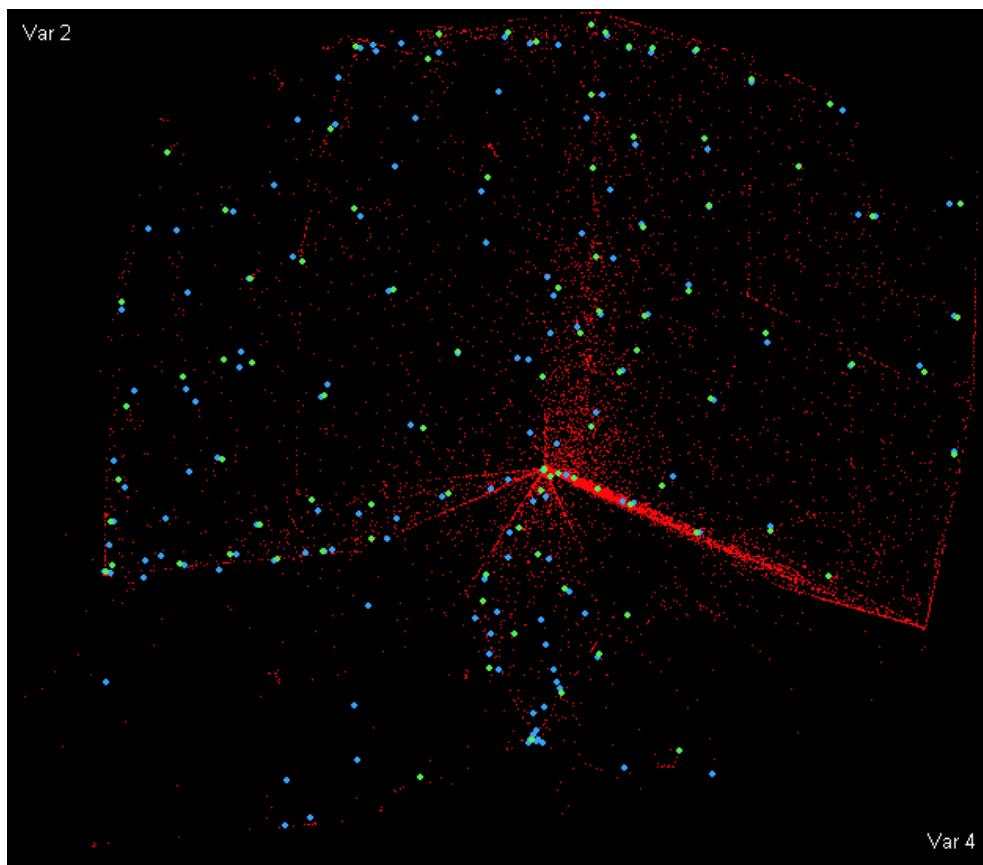


Figure IV-124 : Représentation des prémisses appartenant aux règles finales (ronds verts), issues des règles initiales (ronds bleus) et les états sensoriels observés (points rouges) de l'expérience 39 selon les composantes 2 (Var2) et 4 (Var4).

5.4. Étude d'un mécanisme prédictif

Le processus de création de règles de la troisième série conservait un gabarit de règles sensorimotrices identiques à celles déjà existantes, c'est-à-dire un seul message sensoriel en guise de prémisses et un seul message de commande en guise de conclusion. En revanche, dans la quatrième série d'expériences, afin d'évaluer la capacité de la boucle prédictive à

dégager des structures spatio-temporelles, les règles sensorimotrices générées comportent 4 messages sensoriels pour la prémisse et 4 messages moteurs pour la conclusion. Ainsi, la prédiction porte sur un temps de 500 ms environ. L'ordonnement des règles de ce gabarit constituant la boucle prédictive a également été vérifié dans un environnement imposé.

L'expérience 46 a duré 1800 s avec une fréquence de création de règles à 0,05 Hz soit 90 règles sensorimotrices créées au total. Les étiquettes temporelles des messages sensoriels appartenant aux prémisses permettent de construire des trajectoires dans l'espace sensoriel. Ces trajectoires peuvent parcourir la moitié de l'espace sensoriel, suivant la vitesse du robot et son environnement. La Figure IV-125 illustre quelques trajectoires observées. Le nombre moyen de déclenchements d'une règle prédictive vaut 19 avec un écart type de 14, le nombre maximal de déclenchements atteint 84. Ce faible nombre de déclenchements, comparé aux 18000 déclenchements pour les règles sensorimotrices, s'explique par la taille des prémisses qui raréfie l'apparition d'une configuration adaptée dans la mémoire événementielle. Au total, il y a eu 1676 déclenchements de règles prédictives générées, soit 9,31% des déclenchements. Ce chiffre confirme les résultats des expériences dans un environnement imposé concernant la viabilité des règles possédant une prémisse ayant plusieurs messages sensoriels.

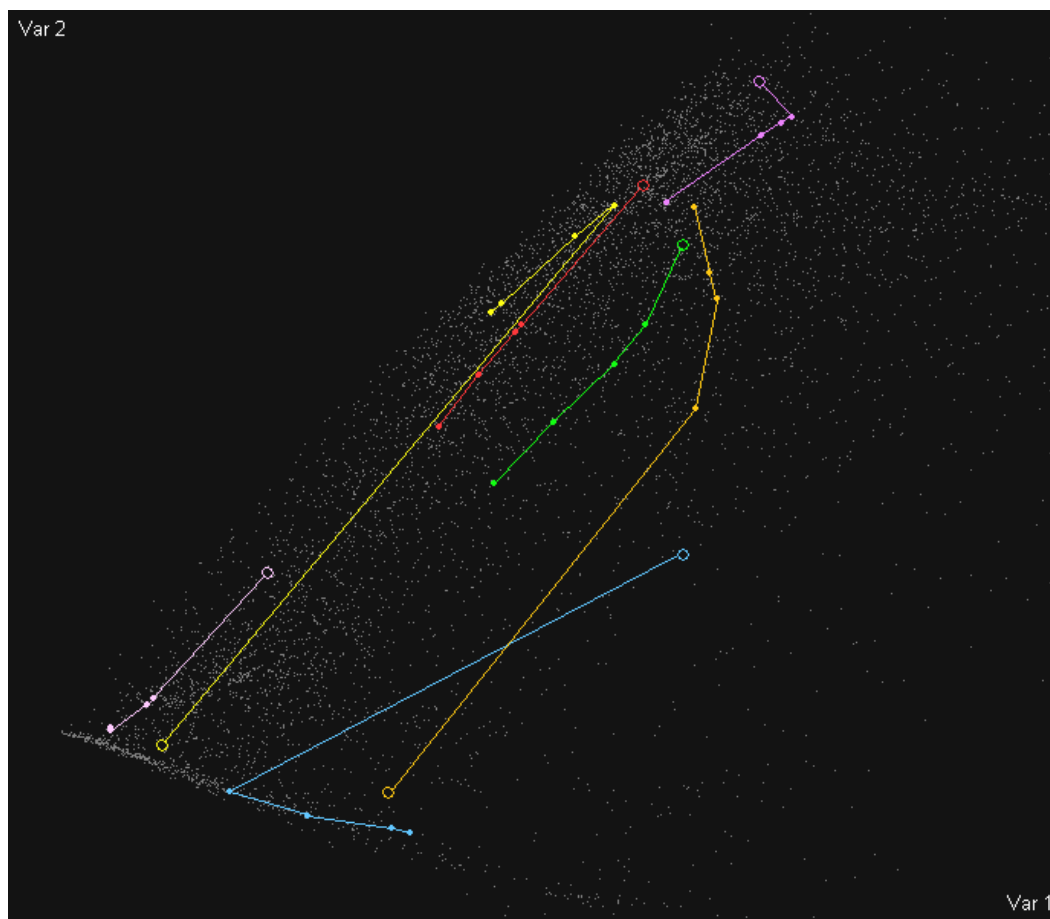


Figure IV-125 : Représentation de 7 prémisses appartenant aux règles créées de l'expérience 46 selon les composantes 1 et 2. Les messages sensoriels des prémisses sont reliés par ordre chronologique et le cercle au bout du segment représente l'état sensoriel prédit.

L'étude des structures temporelles dans les parties précédentes a révélé qu'entre les transitions s'intercale un état sensoriel stationnaire durant trois à cinq fois le temps moyen d'une transition. Ainsi, la probabilité de générer une règle prédictive portant sur un état stationnaire se révèle beaucoup plus grande que la probabilité de saisir une transition. De même, la probabilité de déclenchement pour les règles prédictives portant sur un état stationnaire se trouve supérieure à celle des secondes. Par conséquent, la majorité des règles prédictives déclenchées traduit l'attente de la stabilité, soit une confirmation d'une absence de changement. La Figure IV-126 représente la dérivée en chaque point des états sensoriels observés. La moyenne générale des différences entre deux points successifs est de 0,016 avec un écart type de 0,022.

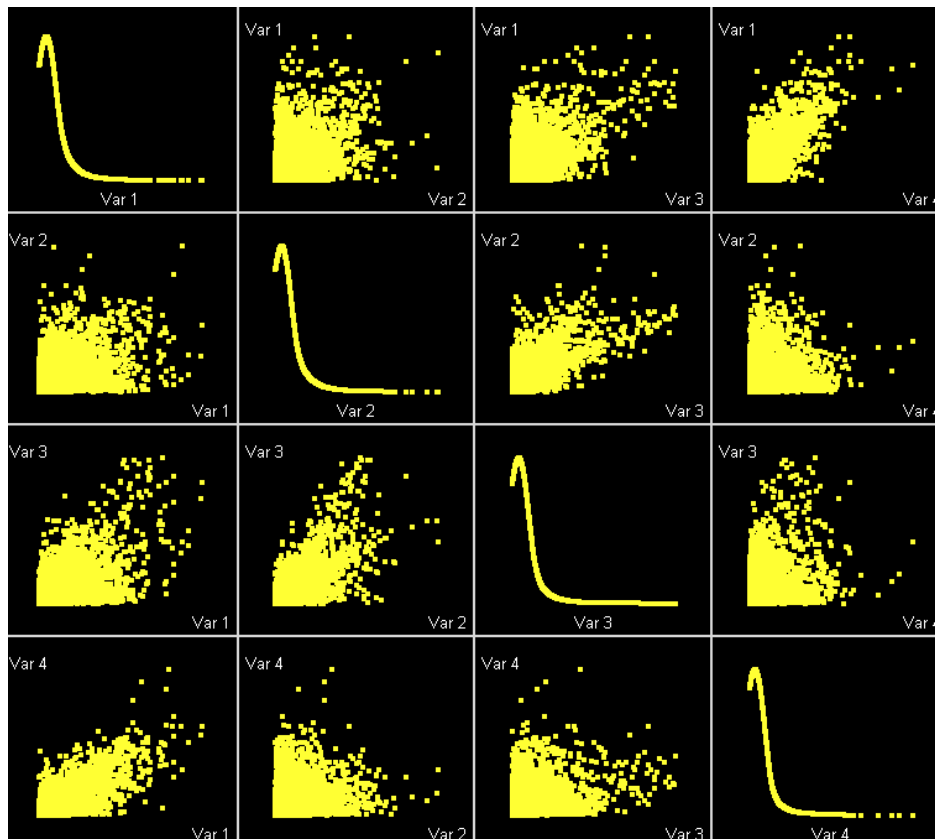


Figure IV-126 : Représentation de la dérivée en chaque point des états sensoriels observés lors de l'expérience 46.

A chaque anticipation se trouve associé un taux de rétribution proportionnel au succès de la prédiction. La valeur de la rétribution calculée s'ajoute ensuite aux forces de toute règle déclenchée entre le début de la prédiction et sa vérification, ainsi qu'aux règles déclenchées juste avant la prédiction. La Figure IV-127 représente la distribution des valeurs du taux de rétribution. Cet histogramme possède une bosse à chacune des bornes de l'intervalle 0 et 1. La première bosse contient les valeurs inférieures à 0,18 soit 18% des récompenses, la seconde, plus large, recouvre les valeurs au-delà de 0,6 soit 57% des taux de rétribution.

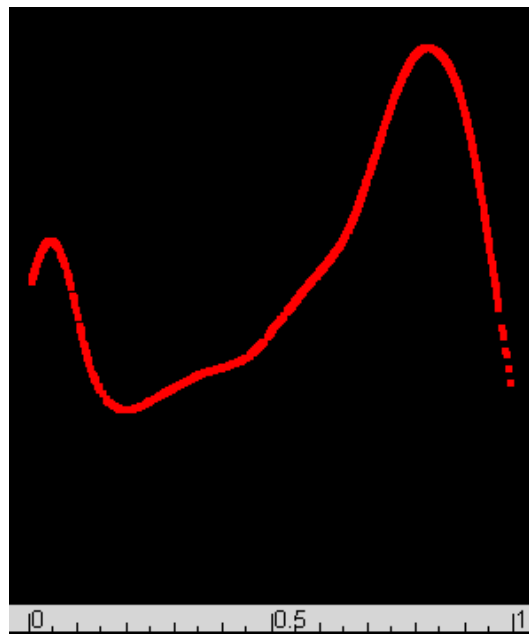


Figure IV-127 : Histogramme des récompenses réalisées avec la méthode de « Averaged Shifted Histograms » à partir de 20 histogrammes décalés de l'expérience 46.

Globalement, cette répartition montre qu'une boucle prédictive d'une épaisseur temporelle de 500ms environ est suffisante pour identifier des états sensoriels imprévus. Cette distinction pourra être alors le point de départ d'une réflexion sur l'élaboration de modalités liées à des situations environnementales ou sur l'opportunité de la création de nouvelles règles.

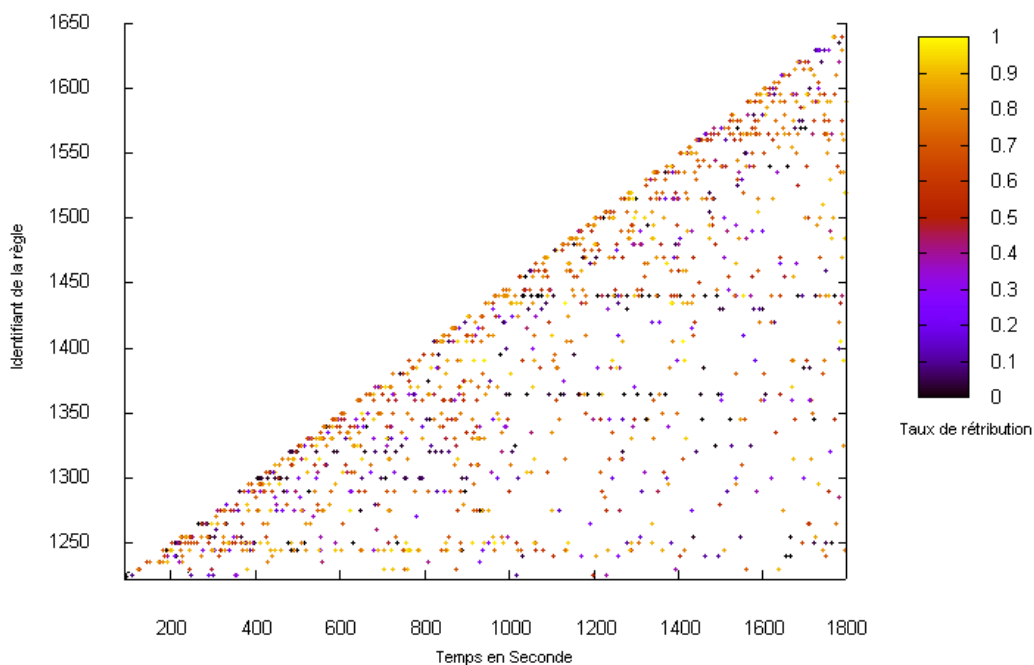


Figure IV-128 : Graphique indiquant le numéro de la règle créée gagnante au cours de l'expérience 46 : la couleur des points est en fonction du taux de rétribution.

Pour finir, comme pour les expériences 44 et 45, le déclenchement des règles créées ce concentre lors des 50 premières secondes de leur apparition, puis seules quelques-unes conservent un déclenchement régulier (Figure IV-128). Les commandes des règles créées gardent une diversité analogue à celle des règles initiales (Figure IV-129).

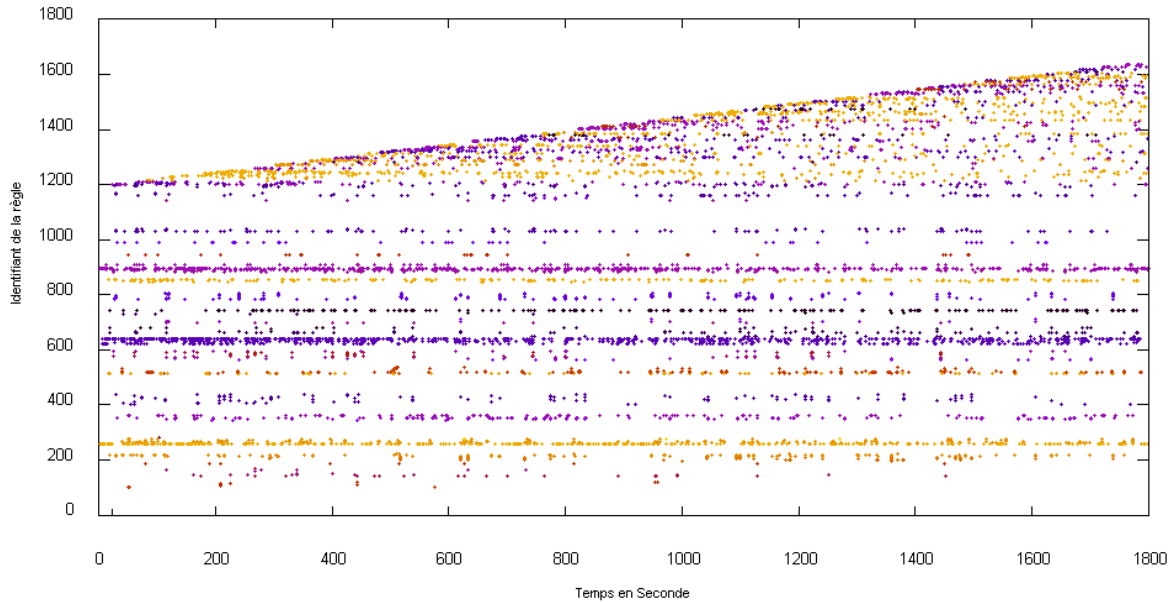


Figure IV-129 : Graphique indiquant le numéro de la règle gagnante à chaque élection de l'expérience 46 : la couleur des points est en fonction de la direction de l'effort moteur.

6. Conclusion

La caractérisation de l'environnement et des primitives associées aux différentes manières de visualiser l'état du système cognitif (comportement du robot, évolution des caractéristiques de la population de règles et les messages sensoriels des prémisses) a permis de dégager les propriétés de ce système rétroactif complexe en (A) boucle ouverte et en (B) boucle fermée.

A - Études en boucle ouverte

La première étude reposant sur 36 expériences montre que le mode propre du système cognitif possède principalement quatre propriétés préalables à une autopoïèse sémiotique. Les deux premières propriétés se rattachent au premier stade de la cognition alors que les deux dernières se rapportent au deuxième stade.

La première propriété provient de l'étude sur le rôle des trois principaux paramètres régissant la dynamique de l'évaluation globale : la taxe qui entraîne l'élimination des règles passives, l'enchère qui élimine les règles redondantes et le remboursement qui compense la taxe et les enchères pour les règles fréquemment actives. La conjonction de ces trois notions induit une dynamique qui converge vers un état stable. Un état stable représente l'équilibre de l'auto-organisation des règles, soit de la clôture opérationnelle du système. L'existence de ces équilibres dynamiques est primordiale puisque l'autopoïèse représente un équilibre dynamique qui résiste aux perturbations.

La deuxième propriété correspond à l'auto-organisation des règles et l'ajustement des prémisses qui conduisent à une classification de l'espace sensoriel. L'ajustement des messages sensoriels des prémisses permet aux règles de se spécialiser, de s'élargir ou de migrer vers des stimulations sensorielles plus fréquentes et moins couvertes. L'évolution du nombre de classes et leur contenu dépendent à la fois de la valeur du taux employé et de l'environnement. Cette double dépendance montre l'existence d'un couplage sensoriel entre les perturbations environnementales imposées et la dynamique interne du système. L'ensemble des règles à chaque instant représente le couplage ponctuel du robot avec son environnement. Cependant, le terme couplage employé pour décrire l'auto-organisation des règles lors des essais avec un environnement imposé se comprend au sens faible puisque les réactions du système restent internes et sans influence sur l'observation de l'environnement. Le couplage entre le robot et l'environnement devient une réelle interaction lorsque les actions du robot conditionnent l'observation de l'environnement et que ces observations conditionnent la forme des règles aboutissant à l'action.

La troisième propriété de cette auto-organisation consiste à pouvoir intégrer de nouvelles règles générées à partir des signes contenus dans la mémoire événementielle. Cette méthode permet d'accéder plus aisément à ces points extrêmes sans mettre en péril les règles anciennes et adaptées. En effet, la compétition élimine les nouvelles règles redondantes. La stabilité de l'auto-organisation et sa propension à couvrir au mieux l'espace des stimulations offre une robustesse indispensable pour la mise en œuvre des schèmes cognitifs générant de nouvelles règles au deuxième stade de la cognition.

La quatrième propriété importante repose sur l'utilisation de l'opérateur de reconnaissance attendue associé à une rétribution. L'utilisation de ce schème cognitif permet de maintenir des règles que la dynamique initiale n'aurait pas maintenues. Les mécanismes subcognitifs sous-tendent l'auto-organisation des règles vers un équilibre « par défaut » mais ces règles peuvent elles-mêmes traduire des mécanismes cognitifs afin d'orienter cet équilibre. La faculté à s'auto-rétribuer se trouve au cœur du deuxième stade *puisque'elle représente la phase de dressage.

B - Études en boucle fermée

L'étude du système cognitif en boucle fermée, c'est-à-dire dans un environnement réel, repose sur une dizaine d'expériences. Ces dernières ont permis de retrouver les propriétés dégagées lors de la première étude. Plus particulièrement, trois points indiquent que le système cognitif développé représente une autopoïèse sémiotique minimale.

Le premier point concerne le comportement du système. Les règles s'ajustent et s'éliminent de telle sorte que l'entropie de l'ensemble de règles diminue et tend vers un équilibre. Sans consigne, l'évolution de règles conduit à un couplage entre l'ensemble des règles et l'environnement sans dénaturer son comportement global « d'évitement d'obstacle ». Contrairement aux autres approches classiques élaborant un couplage sensorimoteur, celui-ci ne repose pas sur un couplage par entrée/sortie à l'aide d'un critère extérieur mais sur un couplage par clôture. Le robot n'apprend pas à éviter les obstacles, il n'optimise pas non plus l'évitement. L'évolution des règles correspond à l'auto-organisation de l'espace sensorimoteur de second ordre en dehors de toute finalité explicite.

Le deuxième point montre que la dynamique des règles permet de rendre compte de l'alternance entre l'assimilation et l'accommodation (Piaget, 1964). Le premier processus d'adaptation correspond à l'équilibre de la dynamique entre les règles (la sémiotique), le second, l'accommodation, correspond à l'intégration d'une perturbation entraînant un

déséquilibre dans le système. L'origine de cette perturbation peut être endogène ou exogène. Dans le cadre du système cognitif, la création régulière de règles représente une perturbation endogène, toutefois elle ne génère pas un déséquilibre dans un environnement déjà connu. En revanche, lorsque de nouvelles situations apparaissent et perturbent le système cognitif, les deux perturbations se rencontrent pour former de nouvelles règles pertinentes et ainsi accommoder l'ensemble des règles pérennes. L'emploi d'une nouvelle règle amène alors son assimilation. Ainsi, le système cognitif construit ses règles en fonction de l'équilibre du couplage par clôture et du déséquilibre induit par des perturbations et non en fonction d'un objectif explicite, de sorte que le système représente un système constructiviste.

Le dernier point souligne la faisabilité d'exprimer en terme de règles les conditions de leur maintien ainsi que la capacité à dégager des structures temporelles. La conservation de nouvelles règles de prédiction qui auraient dû disparaître si elles n'avaient pas été rétribuées montre que la sémiotique possède les moyens d'orienter sa propre dynamique. Par ailleurs, le déclenchement irrégulièrement alterné entre les règles sensorimotrices associées à une prédiction et celles qui ne le sont pas illustre l'alternance entre les phases d'instinct et de dressage.

En somme, le système cognitif possède toutes les propriétés identifiées au premier et au deuxième stade de la généalogie cognitive. Par conséquent, le système cognitif développé à partir de l'architecture cognitive proposée représente une autopoïèse sémiotique minimale.

CONCLUSIONS ET PERSPECTIVES

« Et en poussant jusqu'à son terme cette logique absurde, je dois reconnaître que cette lutte suppose l'absence totale d'espoir (qui n'a rien à voir avec le désespoir), le refus continu (qu'on ne doit pas confondre avec le renoncement) et l'insatisfaction consciente (qu'on ne saurait assimiler à l'inquiétude juvénile). »

Le mythe de Sisyphe p. 52 d'Albert Camus (1942).

« Un beau jour le monde se refroidira et tout mourra ; mais c'est loin, et sa valeur présente à dépréciation constante est presque nulle. Pas plus que le présent n'a moins de valeur parce que le futur sera vide. L'humanité, qui remplit l'avant de la scène de mon tableau, je la trouve intéressante, et dans l'ensemble admirable. »

Logique philosophique et probabilité p. 337 de Frank Ramsey (1925).

Dans l'objectif de mieux dégager les principales conclusions, (A) la problématique sera rappelée et une synthèse de chaque chapitre sera présentée. (B) Les perspectives de ce travail seront ensuite évoquées, concernant à la fois l'étude de l'architecture cognitive et le développement des schèmes cognitifs.

A - Synthèse générale

Les technologies actuelles permettent de concevoir des robots performants pour des tâches complexes lorsque celles-ci sont suffisamment bien définies pour spécifier les étapes de leur réalisation et les heuristiques associées. Les robots ainsi conçus ne possèdent pas d'autonomie dans la mesure où ils ne spécifient pas eux-mêmes leurs objectifs ainsi que les moyens d'y parvenir. L'autonomie suppose l'aptitude à modifier ses objectifs ou à en ajouter de nouveaux en fonction des contingences environnementales. Ainsi, cette capacité devient cruciale lorsque l'environnement se révèle imprévisible à cause de la méconnaissance du milieu ou de la présence d'agents autonomes comme les humains par exemple. En définitive, l'autonomie se trouve sous-tendue par les capacités cognitives de l'individu.

Les progrès technologiques survenus au cours des cinquante dernières années apparaissent ainsi davantage dans le cadre de la réalisation mécanique du robot ou des traitements d'informations spécifiques que dans celui du développement de capacités cognitives. Les nombreuses approches existantes en robotique cognitive montrent que ce déficit ne provient pas d'un désintéressement de la communauté scientifique mais d'une impasse générale.

i - Synthèse du premier chapitre

Afin de déterminer la nature cette impasse, une grille d'analyse philosophique et épistémologique a été établie pour mener une étude transversale des différentes approches en sciences cognitives. Les résultats de cette étude présentés dans le premier chapitre montrent que toutes les approches possèdent une limite intrinsèque pour le développement de systèmes cognitifs et souligne leur diversité tant sur les hypothèses ontologiques ou cognitives que sur les méthodes scientifiques employées. Néanmoins, toutes les approches se réclament d'au moins une des deux hypothèses ontologiques. L'interactionnisme

n'échappe pas à cette critique dans la mesure où il introduit indirectement une de ces hypothèses lors de la réalisation pratique. L'acceptation d'une hypothèse ontologique conduit à définir a priori la notion d'objet (même au plus bas niveau et quelle que soit sa nature, physique ou idéale) et à réduire la cognition à l'extraction et à la manipulation de ces objets. En d'autres termes, la définition a priori de la notion d'objet délimite, dans le cadre de la cognition artificielle, ce qui est connaissable et les moyens d'y accéder, de sorte que la définition de la notion d'objet implique une notion de vérité.

ii - Synthèse du deuxième chapitre

La notion de vérité représente alors la clé pour comprendre les raisons conduisant à cette impasse générale. Le caractère philosophique de l'impasse explique par ailleurs sa difficulté à transparaître dans la mesure où les positions métaphysiques (conscientes ou inconscientes) des concepteurs apparaissent rarement dans les approches en robotique cognitive. L'analyse de la notion de vérité révèle l'impossibilité à concevoir un concept de vérité sur des hypothèses ontologiques sans introduire des tensions intenable. Dans ce cas, seul le pragmatisme offre une philosophie cognitive alternative, sans hypothèses ontologiques, qui permette de s'affranchir d'une définition a priori de la notion d'objet.

Dans le cadre du pragmatisme, « la vérité est une chose qui se fait... Le vrai consiste simplement dans ce qui est avantageux pour notre pensée » (James, 1907). L'application d'une croyance se trouve évaluée par les actions et les croyances auxquelles elle conduit. Contrairement à la notion d'utilité de Mill (1861), la détection de l'avantageux n'est pas une donnée objective. L'idée subjective de ce qui est avantageux se construit et évolue dynamiquement au cours de l'usage. Autrement dit, le pragmatisme repose sur l'existence d'une rétribution interne dont les conditions subjectives évolueraient dynamiquement.

Dans un premier temps, le pragmatisme définit la sémiologie (Peirce, 1934) comme la production de relations et de signes afin de maintenir une satisfaction interne de l'individu toujours en interaction avec l'environnement. Percevoir l'activité cognitive comme une résistance à la perturbation que produisent de nouvelles expériences a incité à effectuer un parallèle entre la notion d'autopoïèse de Maturana et Varela (1989) et celle de sémiologie. Les principaux points de convergence révélés par ce parallèle ont conduit alors à définir la cognition comme une autopoïèse sémiotique.

Dans un second temps, l'interprétation des remarques éthologiques de Lorenz (1937) à travers la définition pragmatiste de la cognition proposée a permis d'établir une généalogie de la cognition en quatre stades au cours de l'évolution des espèces. Ces derniers représentent autant d'états dans le développement de la robotique cognitive, mais la mise en perspective évolutionniste permet surtout de distinguer d'une part la phylogénèse et l'ontogénèse physique, et d'autre part la « culturogénèse » et l'ontogénèse sémiotique.

iii - Synthèse du troisième chapitre

La distinction entre l'autopoïèse physique et l'autopoïèse sémiotique a montré que la concrétude revendiquée dans la définition originale de l'autopoïèse (Varela, 1989) renvoie en définitive à la mise à l'épreuve de la structure du système. Pour l'autopoïèse sémiotique, la mise à l'épreuve se traduit par l'existence d'un mécanisme d'oubli que la dynamique des croyances tente de compenser et par le fait que l'application d'une croyance entraîne un effet conditionnant l'application d'autres croyances.

Cette précision sur la définition de l'autopoïèse sémiotique apporte deux conclusions. La première souligne que le concept d'autopoïèse sémiotique se révèle indépendant du

concept d'autopoièse physique. Par conséquent, la réalisation d'un système cognitif à l'aide d'un système robotique traditionnel est envisageable. Toutefois, bien qu'il existe une indépendance conceptuelle entre ontogénèse physique et ontogénèse sémiotique, ces deux phénomènes se trouvent intimement liés dans le cadre d'un individu cognitif naturel.

La seconde conclusion insiste sur le fait qu'un système logique décrit des relations logiques alors qu'une autopoièse sémiotique réalise des relations effectives. La description logique de l'activité cognitive du système se révèle alors intrinsèquement limitée, néanmoins la description logique de l'architecture cognitive demeure réalisable et correspond à la spécification. L'étude sur des systèmes logiques a permis par ailleurs d'identifier les propriétés minimales à incorporer dans l'architecture cognitive afin de pouvoir retrouver les principes logiques qui sous-tendent les descriptions logiques.

En s'appuyant à la fois sur les remarques phénoménales élémentaires concernant la perception et la conception, et sur la réinterprétation de la triade sémiotique peircienne, l'autopoièse sémiotique a été formalisée en termes d'architecture cognitive représentant l'ensemble des mécanismes subcognitifs nécessaires au déroulement de la sémiose et de schèmes cognitifs représentant la structure des croyances engendrant la sémiose. Cette formalisation a permis d'insister sur la nécessité d'une part de définir les croyances de manière holistique et d'autre part de relier la sémiose à des primitives sensorimotrices, c'est-à-dire d'incarner le système cognitif. Les primitives sensorielles et motrices représentent des processus à part entière dont le principe reste fixe. Ces primitives peuvent ainsi contenir des boucles sensorimotrices primaires sur lesquelles se fondent les boucles sensorimotrices secondaires définies par des schèmes cognitifs. Par ailleurs, l'expression des conditions de rétribution interne en termes de croyances permet l'auto-orientation de la dynamique sémiotique qui, associée à la génération de croyances, représente la notion d'autonomie et le caractère pragmatiste du système.

Ainsi formalisée, l'architecture cognitive ne repose sur aucune hypothèse ontologique tout en offrant la possibilité d'en concevoir et d'en réaliser. En s'inspirant des systèmes de classeurs considérés comme des structures dissipatives subsymboliques, une architecture cognitive a été spécifiée et implémentée afin de produire une autopoièse sémiotique minimale.

iv - Synthèse du quatrième chapitre

Considérer la cognition comme une autopoièse sémiotique implique que la sémiose débute à partir d'un ensemble de croyances précédant toute expérience cognitive : les *proto-croyances*. Par ailleurs, la nécessité des proto-croyances exprime l'impossibilité d'éliminer le normatif (Putnam, 1983). L'évolution des capacités cognitives d'un individu dépend, en plus de l'environnement, de la complexité de ces schèmes cognitifs initiaux. Pour l'implémentation robotique, les schèmes cognitifs initiaux ont représenté une boucle sensorimotrice secondaire simple et un détecteur de régularité par prédiction. L'utilisation de ces deux schèmes a montré, d'une part que l'architecture cognitive proposée possède des auto-organisations permettant de classifier son environnement sans a priori, et d'autre part que cette auto-organisation native peut être orientée par les croyances elles-mêmes. Autrement dit, l'architecture cognitive proposée répond à l'exigence définie par les deux premiers stades de la généalogie de la cognition.

Sur le plan technologique, en jouant le rôle de coordonnateur entre les primitives, l'architecture cognitive peut être perçue comme un intergiciel et les schèmes cognitifs comme la programmation cognitive de la coordination et de son évolution par une

autopoïèse sémiotique. Par rapport aux méthodes de classification traditionnelle, l'architecture cognitive offre la possibilité de fusionner de nombreux critères de classification sans remettre en cause l'auto-organisation native.

*

L'adoption d'une hypothèse ontologique étant la source de l'impasse générale en robotique cognitive, il a été proposé d'intégrer la notion d'autopoïèse (Varela, 1989) au pragmatisme (James, 1907 ; Peirce, 1878), toujours en étroite relation avec les observations éthologiques (Lorenz, 1939) et phénoménales sur la perception et la conception. Dans ce cadre, la cognition s'identifie à l'autopoïèse sémiotique interagissant avec les primitives sensorielles et motrices. En approximant, les approches traditionnelles peuvent être perçues comme des approches qui conçoivent des primitives complexes afin de pallier les limites de leur architecture cognitive qui se réduit à un multiplexeur plus ou moins adaptatif. Cette conception interdit alors toute modification des schèmes comportementaux. En revanche, la démarche proposée souhaite diminuer la complexité des primitives afin de la transférer en termes de schèmes cognitifs dans l'architecture cognitive, c'est-à-dire de concevoir une programmation cognitive. La traduction des tâches complexes en termes de schèmes cognitifs permet de les coordonner et de les optimiser en fonction des autres tâches, mais surtout de modifier potentiellement leur objectif, en somme de rendre l'individu autonome. La majorité des difficultés du développement de la robotique cognitive pragmatiste réside dans la définition des primitives et des proto-croyances, perspectives ouvrant de nombreuses pistes de recherche.

B - Perspectives de recherches

L'approche d'une robotique cognitive pragmatiste a été définie et a permis d'élaborer une formalisation de l'architecture cognitive résultante ainsi qu'une spécification. Toutefois, les études proposées sur l'architecture cognitive en elle-même et sur l'expressivité des schèmes se veulent des études préliminaires dont l'objectif est de vérifier la cohérence globale de la démarche. Cette cohérence ayant été montré, deux types d'investigations doivent être menés afin de poursuivre le projet de l'artificialisation de la cognition, (i) l'un concernant les aspects théoriques de l'architecture cognitive, (ii) l'autre concernant la modélisation de capacités cognitives à l'aide de l'architecture cognitive.

i - Recherches sur l'architecture cognitive

L'étude à mener sur les propriétés de l'architecture cognitive prend principalement trois formes différentes. La première représente les études théoriques sur la correspondance entre les systèmes logiques d'ordre supérieur et l'architecture cognitive afin de déterminer si la notion de recopie est suffisante pour accéder aux logiques d'ordre supérieur. Par ailleurs, la modélisation des logiques traditionnelles ou modales en termes de schèmes cognitifs devrait permettre d'entrevoir les schèmes cognitifs fondamentaux et les expériences indispensables qui se trouve à l'origine de ces schèmes logiques.

La deuxième forme d'étude doit se concentrer sur la signification des différents paramètres initiaux par rapport aux algorithmes de classification traditionnels. Dans ce cadre, le système est étudié en boucle ouverte, ce qui facilitera une interprétation de la dynamique initiale de l'architecture cognitive à travers la théorie de l'information.

Le troisième type d'étude consiste à modéliser les divers types de classeurs et à comprendre la traduction des différents paramètres. En effet, les mécanismes subcognitifs

des classeurs traditionnels deviennent des schèmes cognitifs à part entière. Plus particulièrement, la modélisation des CSA qui intègre le *Q-Learning* permettra d'explorer la mise en œuvre d'un système de rétribution interne complexe.

L'ensemble de ces études conditionnera l'interprétation des résultats issus de la construction de schèmes cognitifs inspirés par les sciences cognitives.

ii - Recherches avec l'architecture cognitive

Les études précédentes sur l'architecture cognitive étant essentiellement comparatives, seuls quelques axes de visualisation sont privilégiés, mais dans le cadre de l'étude de l'évolution des schèmes inspirés des modèles cognitifs, la détection des phénomènes liés à l'autopoïèse sémiotique passe par la visualisation et le recoupement de l'ensemble des données enregistrées. Dans ce cas, avant d'entreprendre ces études, il devient indispensable de regrouper les modes de visualisation accessibles à partir d'une base de données globale pour chaque expérience (vidéo, primitives sensorielles et primitives motrices ainsi que les enregistrements concernant les règles). Quatre axes de recherche se dégagent ensuite prioritairement et chacun d'eux doit être étudié en fonction des différents stades de la cognition considérés.

Le premier axe porte sur la constitution des boucles sensorimotrices primaires et secondaires : Comment distinguer leurs caractéristiques ? Comment intégrer la notion d'effort et quel est son rôle ? Comment les boucles sensorimotrices se traduisent en incorporant la notion d'intensité dans les primitives sensorimotrices via la légitimité des signes ? Quelles sortes d'espaces sensorimoteurs peuvent être élaborés et comment les utiliser pour l'orientation vers un objectif ? Comment construire des boucles sensorimotrices dans un environnement dynamique ?

Le second axe de recherche consiste à concevoir des schèmes cognitifs traduisant des capacités attentionnelles : comment définir les primitives sensorielles pour les capteurs de type matriciel pour la construction de tels schèmes ? Comment moduler la fenêtre attentionnelle et la diriger ? Comment établir un contexte et segmenter le monde en fonction de celui-ci, autrement dit dégager une forme du fond ? Comment évaluer et modifier cette forme ?

Le troisième axe explore les différentes manières de construire des schèmes cognitifs reproduisant les différents conditionnements tels que le conditionnement classique ou le conditionnement opérant. Comment comparer les comportements des robots avec ceux des animaux ? Comment s'auto-organise un ensemble de règles traduisant plusieurs types de conditionnement ?

Le quatrième axe de recherche s'intéresse au développement des conventions langagières lors de l'interaction entre des agents autonomes. Comment leurs échanges peuvent contribuer à l'évolution de leur sémiotique ? Quels schèmes cognitifs initiaux doivent être imaginés pour que la rétribution interne de chacun conduise à la convergence des objectifs, soit à la coopération ? Comment amener un agent à modéliser un autre agent pour anticiper cette convergence ?

L'avantage d'identifier la cognition à une autopoïèse sémiotique est d'offrir un cadre explicatif pour des domaines divers et vastes traditionnellement isolés dans leur tentative d'artificialisation. Cette transversalité constitue la clé pour résoudre les difficultés à concevoir un système cognitif artificiel qui intègre nécessairement les propriétés recherchées par les quatre axes de recherche.

*

L'interactionnisme a relevé toutes les difficultés de la dualité entre une description logique et une description physique du monde mais sans fournir une alternative suffisamment précise pour unifier les problèmes liés à la cognition. La formalisation du pragmatisme, grâce à la notion d'autopoïèse, a permis d'élaborer en définitive une théorie de la connaissance fondée sur la cognition et non sur le monde, qui dépasse cette dualité au point de proposer des pistes pour concevoir une entité cognitive. La compréhension de la portée de cette théorie cognitive passe alors par un examen philosophique et épistémologique de ses implications dans tous les domaines se réclamant des sciences cognitives.

BIBLIOGRAPHIE

- "World Robotics 2005"(2005), in: International Federation of Robotics.
- Afonso, A. (2006), *Propriétés analogiques des représentations mentales de l'espace : étude comparative auprès des personnes voyantes et non-voyantes*, Orsay: Université Paris-Sud.
- Anderson, J. R. (1974), "Retrieval of prepositional information from long-term memory", *Cognitive Psychology* 6:451-474.
- Asch, S. E. (1951), "Effects of group pressure upon the modification and distortion of judgments", in H. Guetzkow (ed.), *Groups, Leadership and Men: Research in Human Relations* Pittsburgh, 177-190.
- Atick, J. J., Z. Li, and A. N. Redlich (1992), "Understanding Retinal Color Coding from First Principles", *Neural Computation* 4 (4):559-572.
- Auvray, M., S. Hanne-ton, J. K. O'Regan, and C. Lenay (2004), "Distal attribution in sensory substitution", Paper read at ASSC8, at Anvers.
- Auvray, M., and J. K. O'Regan (2001), "Influence of semantic factors on blindness to progressive changes in visual scenes", *Perception* 30.
- Bachelard, G. (2002), *La philosophie du non*. Edited by quadrige. puf ed. Original edition, 1940.
- Bach-y-Rita, P., C. Collins, F. Saunders, B. White, and L. Scadden (1969), "Vision Substitution by Tactile Image Projection", *Nature* 221:963-964.
- Bacon, F. (1986), *Novum Organum*. Edited by Épipiméthée. PUF ed. Original edition, 1620.
- Ballard, D. H., and R. P. N. Rao (1995), "Deictic codes for the embodiment of cognition", in: National Resource Laboratory for the Study of Brain and Behavior.
- Barlow, H. B. (1961), "Possible principles underlying the transformation of sensory messages", in W. Rosenblith (ed.), *Sensory Communication*, Cambridge, MA: The MIT Press, 217-234.
- Barlow, H. B., and W. R. Levick (1965), "The mechanism of directionally selective units in rabbit's retina", *Journal of Physiology* 178 (3):477-504.
- Barto, A. G. (1995), "Adaptive critic and the basal ganglia", in J. L. Davis J. C. Houk, & D. G. Beiser (ed.), *Models of information processing in the basal ganglia* Cambridge: MIT Press, 215-232.
- Bayes, T. (1763), "An Essay towards solving a Problem in the Doctrine of Chances", *Philosophical Transactions of the Royal Society of London* 53:370-418.
- Bell, A. J., and T. J. Sejnowski (1997), "The "independent components" of natural scenes are edge filters", *Vision Research* 37 (23):3327-3338.
- Berthoz, A. (1997), *Le sens du mouvement*. Edited by Sciences: Odile Jacob. Original edition, 1997.
- Bethe, Beer, and Uexküll (1899), "Vorschläge zu einer objectivirenden Nomenklatur in der Physiologie des Nervensystems", *Biologisches Centralblatt* 5 (19):517.

- Blodgett, H. C. (1929), "The effect of the introduction of reward upon the maze performance of rats", *University of California Publications in Psychology* 4:114–133.
- Bonelli, P., and A. Parodi (1991), "An Efficient Classifier System and its Experimental Comparison with two Representative learning methods on three medical domains", Paper read at Proceedings of the 4th International Conference on Genetic Algorithms at San Mateo.
- Bourgine, P., and J. Stewart (2004), "Autopoiesis and Cognition", *Artificial Life* 10 (3):327-345.
- Braitenberg, V. (1984), *Vehicles: Essays in synthetic psychology*. Cambridge, MA: The MIT Press.
- Brewka, G., J. Dix, and K. Konolige (1997), "Nonmonotonic Reasoning: An Overview", *Center for the Study of Language and Information*.
- Broca, P. (1861), "Remarques sur le siège de la faculté du langage articulé, suivies d'une observation d'aphémie (perte de la parole)", *Bulletin de la Société Anatomique* 6:330-357.
- Brooks, R. (2001), "The relationship between matter and life", *Nature* 409:409 à 411.
- Brooks, R. A. (1986), "A Robust Layered Control System for a Mobile Robot", *IEEE Journal of Robotics and Automation* RA-2 (1):14-23.
- Brooks, R. A., and P. A. Viola (1990), "Network Based Autonomous Robot Motor Control: From Hormones to Learning", in R. Eckmiller (ed.), *Advanced Neural Computers*, Amsterdam: North Holland, 341-348.
- Brouwer, L. E. J. (1927), "Über Definitionsbereiche von Funktionen (Intuitionistic reflections on formalism)", *Mathematische Annalen* XVII.
- Bruner, J. (2000), *Culture et modes de pensée. L'esprit humain dans ses œuvres*. Paris: Retz.
- Buche, C., R. Querrec, P. De Loor, and P. Chevaillier (2004), "MASCARET : A Pedagogical Multi-Agent System for Virtual Environment for Training", *International Journal of Distance Education Technologies* 2 (4):41-61.
- Buche, C., C. Septseault, and P. De Loor (2006), "Les systèmes de classeurs. une présentation générale", *Revue des Sciences et Technologies de l'Information (RSTI-TSI)* 20:963-990.
- Cajal, R. y. (1893), "La rétines des vertébrés", *La cellule* 9:119-225.
- Camus, A. (1982), *Discours de Suède*. Edited by nrf. Gallimard ed. Original edition, 1958.
- Camus, A. (1998), *L'homme révolté*. Edited by folio. Gallimard ed, *essais*. Original edition, 1951.
- Carew, T. J., R. D. Hawkins, T. W. Abrams, and E. R. Kandel (1984), "A test of Hebb's postulate at identified synapses which mediate classical conditioning in *Aplysia*", *Journal of Neuroscience* 4:1217-1224.
- Carnap, R. (2002), *La construction logique du monde*. Translated by T. Rivain. Edited by Mathésis. Vrin ed. Original edition, 1928.
- Chaitin, G. J. (2002), *The Limits of Mathematics: A Course on Information Theory and the Limits of Formal Reasoning* par Springer.
- Changeux, J.-P. (1983), *L'homme neuronal*. Paris: Fayard.

- Chomsky, N. (1992), *Aspects of the theory of syntax*. 7th printing ed. Cambridge, MA: The MIT Press.
- Churchland, P. (1981), "Eliminative Materialism and the Propositional Attitudes", *Journal of Philosophy* 78 (2):67-90.
- Churchland, P. S. (1986), *Neurophilosophy. Toward a Unified Science of the Mind/Brain*. Cambridge, MA: The MIT Press.
- Chwistek, L. (1921), "Antynomje logiki formalnej", *Przegląd Filozoficzny* 24:164-171.
- Conway, J. (1970), "Mathematical games", *Scientific American* October:120-127.
- Cornoldi, C., R. De Beni, and A. Pra Baldi (1989), "Generation and retrieval of general, specific and autobiographic images representing concrete nouns", *Acta Psychologica* 72:25-39.
- Cornuéjol, A., and L. Miclet (2002), *Apprentissage artificiel, Concepts et algorithmes* Dunod.
- Damasio, A. R. (1995), *L'erreur de Descartes*. Translated by Marcel Blanc: Odile Jacob. Original edition, 1994.
- Darwin, C. (1999), *On the Origin of Species*. Edited by Flammarion. Original edition (1859).
- Dawkins, R. (1976), *The Selfish Gene*. Oxford: Oxford University Press.
- De Finetti, B. (1937), "La prévision : ses lois logiques, ses sources subjectives", *Annales de l'institut Henri Poincaré* 7 (1):1-68.
- Dempster, A. P., N. M. Laird, and D. B. Rubin (1977), "Maximum likelihood from incomplete data via EM algorithm", *Journal of the Royal Statistical Society* 39:1-38.
- Denis, M., and M. Cocude (1997), "On the metric properties of visual images generated from verbal descriptions: Evidence for the robustness of the mental scanning effect", *European Journal of Cognitive Psychology*, 9:353-379.
- Dennett, D. (1993), *La conscience expliquée*. Paris: Odile Jacob.
- Denquive, N., and P. Tarroux (2002), "Multi-resolution codes for scene categorization", Paper read at European Symposium on Artificial Neural Networks ESANN 2002, April 23-26, at Bruges, Be.
- Descartes (1995), *Discours de la méthode: Le livre de poche*. Original edition, 1637.
- Dewey, J. (1896), "The reflex arc concept in psychology", *Psychological Review* 3:357:370.
- Dokic, J. (1999), "L'action située et le principe de Ramsey ", in M. de Fornel and L. Quéré (dir.) (eds.), *La logique des situations. Nouveaux regards sur l'écologie des activités sociales, Raisons Pratiques*, Paris: Editions de l'Ecole des Hautes Etudes en Sciences Sociales, 31-155.
- Dokic, J. (2004), *Qu'est ce que la perception, Chemins Philosophiques*: Vrin. Original edition, 2004.
- Dokic, J., and P. Engel (2001), *Ramsey, Vérité et succès, philosophie*. puf. Original edition, 2001.
- Dorigo, M. (1995), "ALECSYS and the AutoNOMouse: Learning to Control a Real Robot by Distributed Classifier Systems", *Machine learning journal* 19 (3):209-240.
- Dorigo, M., and H. Bersini (1994), "A Comparison of Q-Learning and Classifier Systems", Paper read at From Animals to Animats 3 : Third International Conference on Simulation of Adaptive Behavior - SAB.

- Drogoul, A. (1993), *De la Simulation Multi-Agents à la Résolution Collective de Problèmes*, Paris: Université Pierre et Marie Curie.
- Eco, U. (1987), *Le signe*. Edited by Biblioessais: Le Livre de Poche.
- Edelman, G. M. (2000), *Biologie de la conscience*. Translated by Ana Gerschenfeld. Edited by poches: Odile Jacob. Original edition, 1992.
- Engel, P. (1994), *Introduction à la philosophie de l'esprit, sciences cognitives: La découverte*. Original edition, 1994.
- Engel, P. (1998), *La vérité*. Edited by Optique philosophie: Hatier. Original edition, 1998.
- Enroth-Cugell, C., and J. Robson (1966), "The contrast sensitivity of retinal ganglion cells of the cat", *Journal of Physiology* 187:517–552.
- Faucheux, C., and S. Moscovici (1971), *Psychologie sociale théorique et expérimentale*. Mouton ed. Paris.
- Ferber, J. (2005), *Multiagent Systems, an Introduction to Distributed Artificial Intelligence*. Pearson Education ed.
- Feyerabend, P. (1988), *Contre la méthode*. Translated by Baudoin Jurdant. Edited by Points and Agnès Schlumberger, *sciences*: Seuil. Original edition, 1975.
- Field, D. J. (1987), "Relations between the statistics of natural images and the response properties of cortical cells", *Journal of the Optical Society of America A* 4:2379-2394.
- Finke, R. A., and S. Pinker (1982), "Spontaneous mental image scanning in mental extrapolation", *Journal of Experimental Psychology: Learning, Memory, and Cognition* 8:142-147.
- Floreano, D., and S. Nolfi (1998), "Co-evolving predator and prey robots: Do 'arm races' arise in artificial evolution?" *Artificial Life* 4 (4):311-335.
- Fox, M., M. Ghallab, G. Infantes, and D. Long (2005), "Robot introspection through learned hidden Markov models", *Artificial Intelligence* 169.
- Franceschini, N., and S. Viollet (2002), "Visual servo system based on a biologically-inspired scanning sensor", in MIT Press (ed.), *Neurotechnology for Biomimetic Robots*, Cambridge, 57-71.
- Frege, G. (1994), *Ecrits logiques et philosophiques*. Translated by Claude Imbert. Edited by Points, *Essais*: Seuil. Original edition, 1879-1925.
- Gardenfors, P. (1988), *Knowledge in Flux*. MIT Press ed. Cambridge.
- Gaussier, P., A. Revel, J. Banquet, and V. Babeau (2002), "From view cells and place cells to cognitive map learning : processing stages of the hippocampal system", *Biological Cybernetics* 86:15-28.
- Gaussier, P., and S. Zrehen (1994), "A Topological Neural Map for On-Line Learning. Emergence of Obstacle Avoidance in a Mobile Robot", Paper read at Third International Conference on Simulation of Adaptive Behavior, 8-12 August, 1994, at Cambridge, MA.
- Gell-Mann, M. (1995), *Le quark et le jaguar*. Translated by Gilles Minot: Albin Michel. Original edition, 1994.
- Gérard, P. (2002), *Systèmes de classeurs : étude de l'apprentissage latent* Université Paris VI.

- Gerstner, W., and W. Kistler (2002), "Mathematical formulations of Hebbian Learning." *Biological Cybernetics* 87:404-415.
- Gibson, E., and N. Rader (1979), "The perceiver as performer", in G.A. Hale and M. Lewis (eds.), *Attention and cognitive development*, New York, NY: Plenum.
- Grassé, P.-P. (1959), "La reconstruction du nid et les coordinations inter-individuelles chez *bellicositermes natalensis* et *cubitermes* sp. la théorie de la stigmergie : essai d'interprétation du comportement des termites constructeurs", *Insectes Sociaux* 6:41–80.
- Grice, H. P. (1962), "Some Remarks about the Senses", *Analytical Philosophy*, 133-153.
- Ha-Duong, M. (2005), *Modèles de précaution en économie: introduction aux probabilités imprécises*. Université de Paris I.
- Hebb, D. O. (1949), *The organization of behavior*. New York: Wiley.
- Hegel, G. W. F. (2003), *Principes de la philosophie du droit*. Edited by Quadrigue. PUF ed. Original edition, 1821.
- Heidegger, M. (1995), *Etre et Temps*. Translated by François Vezin. Edited by nrf: Gallimard. Original edition, 1927.
- Hempel, C. G. (1960), *Philosophy of natural science*. Edited by Prentice Hall. Englewood Cliffs.
- Heudin, J.-C. (1998), *L'évolution au bord du chaos*. Hermès ed. Paris.
- Hjelmslev, L. (1957), "Dans quelle mesure les significations des mots peuvent-elles être considérées comme formant une structure", Paper read at Reports of the Eighth International Congress of Linguists, at Oslo.
- Hoffmann, J. (1993), *Vorhersage und Erkenntnis*. Göttingen: Hogrefe.
- Hofstadter, D. (1985), *Godel, Escher, Bach. Les Brins d'une Guirlande Eternelle*. Francaise ed: InterEditions.
- Holland, J. (1976), "Adaptation", in R. Rosen and F. Snell (eds.), *Progress in theoretical biology*, New York: Plenum.
- Holland, J., and J. Reitman (1978), "Cognitive systems based on adaptive algorithms", in D. A. Waterman and F. Hayes-Roth (eds.), *Pattern-directed inference systems*, New York: Academic Press.
- Holland, J. H. (1975), *Adaptation in natural and artificial systems*. The University of Michigan Press.
- Holland, J. H. (1980), "Adaptive algorithms for discovering and using general patterns in growing knowledge bases", *International Journal of Policy Analysis and Information Systems* 4:245-268.
- Hopfield, J. J. (1982), "Neural networks and physical systems with emergent collective computational abilities", *Proceedings of the National Academy of Sciences of the United States* 79:2554-2558.
- Hotelling, H. (1933), "Analysis of a complex of statistical variables into principal components", *Journal of educational psychology* 24:417-441.
- Hume (1983), *Enquête sur l'entendement humain*. Edited by Gf: Flammarion. Original edition, 1748.

- Husserl, E. (1985), *Idées directrices pour une phénoménologie*. Edited by tel. Gallimard ed. Original edition, 1913.
- Husserl, E. (1996), *Logique formelle et Logique transcendantale*. Edited by Épiméthée. PUF ed. Original edition, 1929.
- Ilya Prigogine, I. S. (1986), *La nouvelle alliance*. Edited by Folio. Gallimard ed.
- Ising, E. (1924), *Beitrag zur Theorie des Ferro- und Paramagnetismus* Thèse de doctorat de l'Université d'Hambourg.
- Itti, L., and C. Koch (2001), "Computational Modelling of visual attention", *Nature reviews neuroscience* 2:194-203.
- James, W. (1968), *Le Pragmatisme*. Translated by E. Le Brun: Flammarion. Original edition, 1907.
- Johnson-Laird, P. (1983), *Mental models : Toward a cognitive science of language, inference, and consciousness*. Cambridge: Harvard University Press.
- Kandel, E. R. (1988), *Cellular Biology of Neurons*. Amer Physiological Society ed.
- Kant, E. (1987), *Critique de la raison pure*. Translated by Jules Barni. Edited by GF: Flammarion. Original edition, 1781.
- Kaplan, F. (1999), "Dynamiques de l'auto-organisation lexicale: simulations multi-agents et Têtes parlantes", *In Cognito : Revue internationale francophone en Sciences Cognitives* 15:3-23.
- Khun, T. S. (2001), *La structure des révolutions scientifiques*. Translated by Laure Meyer. Edited by Champs: Flammarion. Original edition, 1970.
- Koffka, K. (1935), *Principles of Gestalt Psychology*. Lund Humphries ed. London.
- Kohonen, T. (1982), "Self-organized formation of topographically correct feature maps", *Biological Cybernetics* 43:59-69.
- Korzybski, A. (1933), *Science and Sanity, an Introduction to Non-Aristotelian Systems and General Semantics*
- Kosslyn, S. (1980), *Image and mind*. Mass.: Harvard University Press.
- Kripke, S. (1959), "A Completeness Theorem in Modal Logic", *Journal of Symbolic Logic* 24:1-14.
- Krivine, J.-L. (2004), "Wigner, Curry et Howard - La déraisonnable efficacité des mathématiques", Paper read at ARCo'04, at Université de Compiègne.
- Laplace, P.-S. D. (1814), *Essai philosophique sur les probabilités*. Paris.
- Latour, B. (1997), *Nous n'avons jamais été modernes*. Edited by Poche: La découverte. Original edition, 1991.
- Latour, B. (2001), *Le métier de chercheur regard d'un anthropologue*. Edited by Sciences en questions: INRA. Original edition, 1994.
- Lehar, S. (2003), *The World in Your Head. A Gestalt View of the Mechanism of Conscious Experience*. Erlbaum ed.
- Lewin, K. (1959), *Psychologie dynamique*. Edited by Les relations humaines. PUF ed. Paris.

- Linsker, R. (1988), "Self-organization in a perceptual network", *Computer Magazine* 21:105-117.
- Locke, J. (2000), *Essai sur l'entendement humain*. Translated by M. Coste. Vrin ed. Original edition, 1689.
- Mach, E. (1996), *L'analyse des sensations. Le rapport du physique au psychique*. Translated by F. Eggers et J.-M. Monnoyer. Jacqueline Chambon ed. Original edition, 1900.
- Machrouh, Y., J. S. Lienard, and P. Tarroux (2001), "Multiscale feature extraction from visual environment in an active vision system", Paper read at International Workshop on Visual Form 4, at Capri, It.
- Mandelbrot, B. (1995), *Les Objets fractals, survol du langage fractal*. Edited by Champs. Flammarion ed.
- Mariotte, E. (1668), *Nouvelle découverte touchant la vue*. Paris.
- Marr, D., and E. Hildreth (1980), "Theory of edge detection", *Proceedings of the Royal Society of London [Biol]* 207:187-217.
- Mataric, M. (1992), "Integration of representation into goal-driven behavior-based robots", *IEEE Journal of Robotics and Automation*, 8 (3):304-312.
- Maturana, H. (1978), "Biology of language: The epistemology of reality", in Miller George and Elizabeth Lenneberg (eds.), *Psychology and Biology of Language and Thought: Essays in Honor of Eric Lenneberg*. Academic Press, 27-63.
- McCarthy, J. (1986), "Circumscription - a form of non-monotonic reasoning", *Artificial Intelligence* 13:27-39.
- McCulloch, W., and W. Pitts (1943), "A logical calculus of the ideas immanent in nervous activity", *Bulletin of Mathematical Biophysics* 5:115-133.
- Merleau-ponty, M. (1942), *La Structure du comportement*. Paris.
- Merleau-Ponty, M. (1948), *Sens et Non-Sens*.
- Merleau-Ponty, M. (1964), *L'oeil et l'esprit*. Edited by nrf: Gallimard. Original edition, 1964.
- Métivier, M. (2004), *Méthodes évolutionnaires et apprentissage : Apprentissage par imitation dans le cadre des systèmes de classeurs* Université Paris V.
- Mettrie, J. O. d. L. (2004), *L'homme plus que machine*. Edited by Rivages poche: Payot & Rivages. Original edition, 1764.
- Mill, J. S. (1988), *L'utilitarisme*. Translated by Georges Tanesse. Edited by Champs: Flammarion. Original edition, 1861.
- Mill, J. S. (2003), *La Nature*. Translated by Estiva Reus. Edited by Poche: La Découverte. Original edition, 1874.
- Nadal, J. P., and N. Parga (1997), "Redundancy reduction and independent component analysis: Conditions on cumulants and adaptive approaches", *Neural Comput* 9 (7):1421-1456.
- Nadal, J.-P., and N. Parga (1994), "Duality Between Learning Machines: A Bridge Between Supervised and Unsupervised Learning", *Neural Computation* 6:491-508.

- Nagel, E., J. R. Newman, K. Gödel, and J.-Y. Girard (2000), *Le théorème de Gödel*. Translated by Jean-Baptiste Scherrer. Edited by Points, sciences: Seuil. Original edition, 1989.
- Nagel, T. (1974), "What is it like to be a bat?" *Philosophical Review* 83 (4):435-451.
- Nagel, T. (1986), *The View from Nowhere*. Oxford: Oxford University Press.
- Nash, J. (1951), "Equilibrium points in n-person games", *Proceedings of the National Academy of Sciences* 36:48-49.
- Neumann, J. V. (1996), *L'ordinateur et le cerveau*. Edited by Champs. Flammarion ed. Original edition, 1957.
- Newell, A. (1980), "Physical symbol system", *Cognitive Science* 4:135-183.
- Newell, A., J. Shaw, and H. Simon (1960), "A variety of intelligent learning in a General Problem Solver", in, *Self-Organizing Systems*: Pergammon Press, 153-189.
- O'Regan, J. K. (1992), "Solving the "Real" Mysteries of Visual Perception: The World as an Outside Memory", *Canadian Journal of Psychology* 46 (3):461-488.
- Pearson, K. (1901), "On lines and planes of closest fit to systems of points in space", *Phil. Mag* 2:559-572.
- Peirce, C. S. (1960), *Collected Papers - Principles of philosophy, Elements of logic*. Cambridge Harvard University Press.
- Penfield, W., and E. Boldrey (1937), "Somatic motor and sensory representation in the cerebral cortex of man as studied by electrical stimulation", *Brain* 60:389-443.
- Penrose, R. (1994), *Shadows of the mind* Oxford Press ed. Oxford.
- Philipona, D., J. K. O Regan, and J.-P. Nadal (2003), "Is there something out there? Inferring space from sensorimotor dependencies", *Neural Computation* 15 (9):2029-2049.
- Philipona, D., and J. K. O'Regan (2005), "Perception sensorimotrice de l'espace", *arobase* 1:71-74.
- Piaget, J. (1992), *Sagesse et illusions de la philosophie*. Edited by quadrige: puf. Original edition, 1965.
- Piaget, J. (2001), *Six études de psychologie*. Edited by Folio essais: Denoël. Original edition, 1964.
- Platon (1964), *Le banquet*. Translated by Emile Chambry. Edited by GF: Flammarion. Original edition, 380 av. J.-C.
- Poggio, G. F. (1995), "Stereoscopic processing in monkey visual cortex: A review", in T.V. Pappathomas, C. Chubb, A. Gorea and E. Kowler (eds.), *Early vision and beyond*, Cambridge, MA: The MIT Press, 43-54.
- Poincaré, H. (1999), *Science et méthode*. Paris: Kimé. Original edition, 1908.
- Poincaré, H. (2000), *Science et Méthode*. Edited by Philosophia scientiae. Kimé ed. Original edition, 1908.
- Popper, K. (2003), *Logique de la découverte scientifique*. Payot ed. Original edition, 1934.
- Putnam, H. (1975), "The Meaning of Meaning", *Philosophical Papers* 2:228.

- Putnam, H. (1992), *Pourquoi ne peut-on pas "naturaliser" la raison*. Translated by Christian Bouchindhome. Edited by tiré à part, *définitions: l'éclat*. Original edition, 1983.
- Putnam, H. (1996), *Philosophie de la logique*. Translated by Patrick Peccatte. Edited by tiré à part: l'éclat. Original edition, 1971.
- Pylyshyn, Z. (1984), *Cognition and computation*. Cambridge: MIT Press.
- Quine, W. v. O. (1977), *Le Mot et la Chose*. Translated by P. Gochet J. Dopp. Edited by Champs. Flammarion ed.
- Ramsey, F. P. (1930), "On a problem of formal logic", *Proceedings of the London Mathematical Society* 48:264-286.
- Ramsey, F. P. (1978), *Foundations. Essays in Philosophy, Logic, Mathematics and Economics*. D. H. Mellor ed. Londres.
- Reiter, R. (1978), "On closed world data bases", in H. Gallaire and J. Minker (eds.), *Logic and Data Bases*, New York: Plenum Press, 55-76.
- Reiter, R. (1980), "A logic for default reasoning", *Artificial Intelligence* 13:81-132.
- Rensik, J. A., J. K. O'Regan, and J. J. Clark (1997), "To see or not to see: The need for attention to perceive changes in scenes", *Psychological Science* 8 (5):368-373.
- Robert, G., and A. Guillot (2006), "MHICS, une architecture de sélection de l'action adaptative pour joueurs artificiels", in T. Cazenave (ed.), *Intelligence Artificielle et Jeu*, Paris: Hermès, 47-80.
- Robert, G., P. Portier, and A. Guillot (2002), "Classifier systems as 'Animat' architectures for action selection in MMORPG", Paper read at 3rd International Conference on Intelligent Games and Simulation, at London.
- Roll, J. P. (2003), "Physiologie de la kinesthèse. La proprioception musculaire : sixième sens ou sens premier ?" *Intellectica* (36-37):49-66.
- Rosenblatt, F. (1958), "The perceptron: a probabilistic model for information storage and organization in the brain", *Psychological Review* 65:386-408.
- Ruderman, D. L. (1994), "The statistics of natural images", *Network : Computation in Neural Systems* 5:517-548.
- Russell, B. (1989), *Écrits de logique philosophique*. Translated by Jean-Michel Roy. PUF ed. Original edition, 1918.
- Rutherford, E. (1911), "The Scattering of a and b Particles by Matter and the Structure of the Atom", *Philosophical Magazine* 21:669-688.
- Sakarovitch, J. (2003), *Éléments de théorie des automates*. Vuibert ed.
- Sanchez, S. (2004), *Mécanismes évolutionnistes pour la simulation comportementale d'acteurs virtuels*. Université des Sciences Sociales Toulouse I.
- Sanza, C. (2001), *Évolution d'Entités Virtuelles Coopératives par Système de Classifieurs*, Toulouse: Université Paul Sabatier.
- Sartre, J.-P. (1970), *L'existentialisme est un humanisme*. Edited by Collection Pensée: Nagel. Original edition, 1946.
- Sartre, J.-P. (1996), *L'être et le néant*. Edited by Tel: Gallimard. Original edition, 1943.

- Savage, L. (1954), *The foundations of statistics*. New-York: John Wiley.
- Schilpp, P. A. (1963), *The philosophy of Rudolf Carnap*: La Salle, Open Court.
- Schöner, G., and M. Dose (1992), "A dynamical systems approach to task-level system integration used to plan and control autonomous vehicle motion", *Robotics and Autonomous Systems* 10:253-267.
- Schulenburg, S., and P. Ross (2001), "Strength and Money: An LCS Approach to Increasing Returns", in Pier Lanzi, Wolfgang Stolzmann and Stewart Wilson (eds.), *Advances in Learning Classifier Systems*, Berlin: Springer-Verlag, 114-137.
- Schultz, W., and R. Romo (1990), "Dopamine neurons of the monkey midbrain, contingencies of responses to stimuli eliciting immediate behavioral reactions", *Journal of Neuroscience* 63:607-624.
- Scott, D. W. (1992), *Multivariate Density Estimation: Theory, Practice, and Visualization*. New York: John Wiley.
- Searle, J. (1984), "Intentionality in Its Place in Nature", *Synthese* 61:3-16.
- Searle, J. R. (1980), "Minds, brains and programs", *Behavioral and Brain Sciences* 3:417-424.
- Shepard, R. (1975), "Form, formation and transformation of internal representations", in R. Solso (ed.), *Information processing and cognition : the Loyola Symposium*, New York: Lawrence Erlbaum Associates, 87-112.
- Sherif, M. (1935), "A study of some social factors in perception", *Archives of Psychology* 187:1-60.
- Sigaud, O. (2001), *Automatisme et subjectivité : l'anticipation au coeur de l'expérience* Université Paris I.
- Simon, H. A. (2004), *Les sciences de l'artificiel*. Translated by Jean-Louis Le Moigne. Edited by folio, *essais*: Gallimard. Original edition, 1996.
- Simoncelli, E. P., and B. A. Olshausen (2001), "Natural Image Statistics and Neural Representation", *Annual Review of Neuroscience* 24:1193-1216.
- Simons, D. J. (1996), "In sight, out of mind : When object representation fail", *Psychological Science* 7 (5):301-305.
- Sims, K. (1994), "Evolving Virtual Creatures. Computer Graphics", Paper read at SIGGRAPH '94.
- Smith, S. F. (1980), *Learning System Based on Genetic Adaptive Algorithms*: University of Pittsburgh.
- Smithson, M. (1989), *Ignorance and Uncertainty - Emerging Paradigms*. Springer-Verlag ed. New York.
- Soler, L. (2000), *Introduction à l'épistémologie*. Edited by philo: ellipses. Original edition, 2000.
- Steels, L. (1999), "The Spontaneous Self-organization of an Adaptive Language", in S. Muggleton (ed.), *Machine Intelligence*, Oxford: Oxford University Press, 205-224.
- Stolzmann, W. (1998), "Anticipatory classifier systems", Paper read at Third Annual Genetic Programming Conference.

- Suri, R., and W. Schultz (1998), "Learning of sequential movements by neural network model with dopamine-like reinforcement signal", *Experimental Brain Research* 121 (3):350–354.
- Thorpe, S. J., K. R. Gegenfurtner, M. Fabre-Thorpe, and H. H. Bülthoff (2001), "Detection of Animals in Natural Images Using Far Peripheral Vision", *European Journal of Neuroscience* 14 (5):869-876.
- Tolman, E. C. (1948), "Cognitive maps in rat and men", *The Psychological Review* 55 (4):189-208.
- Troadeç, B. (1998), *Psychologie du développement cognitif*. Edited by Synthèse, *Psychologie*: Armand Colin. Original edition, 1998.
- Troadeç, B., and C. Martinot (2003), *Le développement cognitif. Théories actuelles de la pensée en contextes*. Paris Belin.
- Turing, A., and J.-Y. Girard (2000), *La machine de Turing*. Translated by Julien Basch. Edited by Points and Patrice Blanchard, *Sciences*: Seuil. Original edition, 1995-1936.
- Valenzuela-Rendon, M. (1991), "The Fuzzy Classifier System : a Classifier System for Continuously Varying Variables", Paper read at Proceedings of the Fourth International Conference on Genetic Algorithms.
- Varela, F. (1989), *Autonomie et connaissance*. Translated by Paul Bourguine and Paul Dumouchel. Edited by La couleur des idées: Seuil. Original edition, 1982.
- Varela, F., E. Thompson, and E. Rosch (1991), *The Embodied Mind : Cognitive Science and Human Experience*: MIT Press.
- Varela, F. J. (1996), *Invitation aux sciences cognitives*. Translated by Pierre Lavoie. Edited by Points, *Sciences*: Seuil. Original edition, 1988.
- Varela, F. J. (2004), *Quel savoir pour l'éthique ?* Translated by Franz Regnot. Edited by poche: La Découverte. Original edition, 1992.
- Viterbi, A. (1967), "Error bounds for convolutional codes and an asymptotically optimum decoding algorithm", *IEEE Transactions on Information Theory* 13 (2):260-269.
- Volken (2003), "« Je le vois, mais je ne le crois pas... » Preuves et vérités dans les sciences formelles", Paper read at Raisons et rationalités, at Arolla
- Waddington, C. (1956), *Principles of Embryology*. London: Allen and Unwin.
- Walter, G. (1950), "An Electromechanical Animal", *dialectica* 4 (3):206-213.
- Watkins, C. J. C. (1989), *Learning from delayed rewards*. Psychology Department, Cambridge, England.: Cambridge University.
- Watson, L. (1979), *Lifetide: The biology of the unconscious*. New York: Simon Schuster
- Weismann, A. (1883), *Ueber die Vererbung*.
- Weizenbaum, J. (1966), "ELIZA - A Computer Program for the Study of Natural Language Communication between Man and Machine", *Communications of the Association for Computing Machinery* 9:36-45.
- Wiener, N. (1948), *Cybernetics or Control and Communication in the Animal and the Machine*. MIT Press ed. Cambridge.

- Wilson, W. (1990), "The Animat Path to AI", Paper read at First International Conference on Simulation of Adaptive Behavior, at Cambridge, MA.
- Wilson, W. (1994), "ZCS : A zeroth level classifier system", *Evolutionary Computation* 2 (1):1-18.
- Wilson, W. (1995), "Classifier Fitness Based on Accuracy", *Evolutionary Computation* 3 (2):149-175.
- Winograd, T. (1972), *Understanding natural language*. Edinburgh Academic Press ed.
- Wittgenstein, L. (2001), *Tractatus logico-philosophicus*. Translated by Gilles-Gaston Granger. Gallimard ed. Original edition, 1921.
- Wolfram, S., ed. (1986), *Theory and Applications of Cellular Automata*. Singapore: World Scientific.
- Yarbus, A. F. (1967), *Eye Movements and Vision*. New York: Plenum Press.
- Zadeh, L. (1978), "Fuzzy sets", *Information and Control* 1:3-28.
- Ziemke, T. (2001), "The construction of 'reality' in the robot: constructivist perspectives on situated artificial intelligence and adaptive robotics", in A. Riegler (ed.), *Foundations of Science*.
- Zuse, K. (1967), "Rechnender raum", *Elektronische Datenverarbeitung* 8:336-344.
- Zykov, V., E. Mytilinaios, B. Adams, and H. Lipson (2005), "Self-reproducing machines", *Nature* 435 (7038):163-164.