



HAL
open science

La convergence des modularités structurelles et fonctionnelles des systèmes complexes

Nicolas Omont

► **To cite this version:**

Nicolas Omont. La convergence des modularités structurelles et fonctionnelles des systèmes complexes. Modélisation et simulation. Université d'Evry-Val d'Essonne, 2009. Français. NNT: . tel-00369892

HAL Id: tel-00369892

<https://theses.hal.science/tel-00369892>

Submitted on 22 Mar 2009

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.

Université d'Evry Val d'Essonne
Ecole doctorale des Génomes Aux Organismes
Programme Epigénomique

Thèse présentée pour le titre de docteur en informatique
par Nicolas OMONT

**La convergence des modularités structurelles et
fonctionnelles des systèmes complexes**

Soutenue le 12 janvier 2009 devant le jury suivant :

Jean-Loup Faulon,	président du jury
François Képès,	directeur de thèse
Jean-Pierre Nadal,	rapporteur
Olivier Teytaud,	directeur de thèse
Louis Wehenkel,	rapporteur
Jérôme Wojcik,	examineur

τί γὰρ ὠφελήσεται ἄνθρωπος ἐὰν τὸν κόσμον ὅλον κερδήσῃ
τὴν δὲ ψυχὴν αὐτοῦ ζημιωθῆ;

*Que sert à l'homme de gagner l'univers entier,
s'il vient à perdre son âme ?*

(Evangile selon Matthieu 16,26)

Remerciements

Si le doctorat est un diplôme attribué à une personne, il marque l'aboutissement d'un travail qui ne peut être accompli seul. Je tiens à remercier l'ensemble des personnes qui m'ont accompagné sur tout ou partie du chemin.

Passionné par le cours de bioinformatique (35) donné à l'Ecole Polytechnique par François Képès et Frédéric Dardel, grâce à la compréhension de Dominique Ventre, directeur des études initiales de Télécom Paris (Pardon, Telecom ParisTech), pour ce double cursus particulier, j'ai eu la possibilité de suivre le D.E.A. « Application des Mathématiques et de l'Informatique à la Biologie » de l'Université d'Evry. Tout indiquait à ce moment que la parenthèse de la bioinformatique se refermait avec la fin du D.E.A. lorsqu'Hiroaki Tanaka, directeur de la bioinformatique de Serono Genetic Institute, ex-Genset, contacte le laboratoire Statistique et Génome de Bernard Prum afin d'entamer une collaboration. Elle s'est développée par la mise en place de deux doctorats C.I.F.R.E. C'est de justesse que j'acceptais d'être l'un des deux doctorants. Sans toutes ces personnes, je n'aurais jamais commencé ce doctorat.

Rapidement après le début de la thèse, il devient évident que les méthodes classiques d'analyse des études d'association génétique ne suffiront pas pour les études à l'échelle du génome constituant le cœur des travaux de Serono Genetics Institute. Premier changement de programme : plutôt que d'étudier le croisement des résultats de l'analyse des études génétiques avec d'autres données, il est nécessaire de se concentrer sur l'obtention des résultats eux-mêmes. A peine la mise au point d'une nouvelle méthode achevée, la maison mère décide de fermer le centre de recherche. L'objet initial de la thèse nécessitant une forte collaboration avec les équipes de Serono Genetics Institute (gestionnaires de bases de données, curateurs, etc.), je décide d'arrêter le doctorat. Malgré tout, avec le soutien de François Képès et de Jérôme Wojcik, nous rédigeons un article sur la base de cette méthode.

L'histoire de ce doctorat aurait pu s'arrêter là. Cependant, Arnaud Renaud, Président-Directeur Général d'Artelys où je venais d'entrer, mais aussi professeur associé à la Sorbonne, décide d'explorer de nouvelles pistes pour que je finisse cette thèse. Grâce à cette nouvelle impulsion, avec le soutien ininterrompu de François Képès et le nouveau soutien d'Olivier Teytaud, chercheur à l'I.N.R.I.A. et collaborateur d'Artelys, nous montons un projet de poursuite de la thèse. Celui-ci prévoit l'audacieux

mélange de la bioinformatique et de l'optimisation des réseaux électriques ainsi qu'une partie minoritaire du travail effectué au sein d'un laboratoire académique. Je remercie Francis Quéfier, directeur de l'école doctorale « des Génomes aux Organismes » d'avoir fait confiance à ce projet sortant des sentiers battus en me permettant de m'inscrire en troisième année de thèse (au mois de juin !). Grâce à la passionnante collaboration d'Artelys avec Alain Hautot et Patrick Sandrin de R.T.E., j'étais introduit dans les arcanes de la tarification marginale des réseaux de transport électrique. La matière du dernier article était trouvée. Sans cette bienveillance de l'ensemble de ces personnes, la thèse n'aurait pas repris.

Commençait peut-être la phase la plus stimulante et la plus enrichissante de la thèse : celle de la synthèse de l'ensemble des travaux dans un cadre conceptuel commun. Je remercie ici particulièrement Olivier Teytaud pour avoir relu et enrichi le manuscrit dès les premiers jets les plus confus, et François Képès pour ses vues d'ensemble. Je remercie aussi tous les autres relecteurs pour leurs remarques et corrections de tous ordres : Emmanuel Avril, Laure Comby, Anne-Sophie de Courcy, Michel Funfschilling, Alain Hautot et Barbara Lucet. Je remercie aussi mes rapporteurs, Jean-Pierre Nadal et Louis Wehenkel pour leurs remarques judicieuses et leurs évaluations bienveillantes. Je remercie enfin Jean-Loup Faulon pour la curiosité dont il a fait preuve à la lecture de ce manuscrit et pour avoir accepté de présider le jury.

Si les encouragements d'Arnaud Renaud sont restés soutenus pendant cette phase, le volume de travail nécessaire à sa réalisation a largement débordé sur des heures normalement destinée à ma famille qui s'est bien agrandie depuis le début de la thèse. C'est pourquoi, la dernière personne que je souhaite remercier est Anne-Sophie, ma femme, qui m'a vu m'éclipser régulièrement le soir, les week-ends et pendant les vacances.

Résumé

L'objet de cette thèse est la *convergence structure-fonction* dans les *systèmes* complexes et ses applications aux systèmes vivants et aux systèmes technico-économiques.

Après avoir défini la *modularité* et identifié les difficultés associées à sa définition, cette thèse formalise le concept de convergence structure-fonction dans les *systèmes évolutifs et fonctionnels* et montre son intérêt pour l'évolutivité et la robustesse de ces systèmes. Ensuite, elle applique ce concept à des problématiques réelles de *systèmes évolutifs et fonctionnels* en biologie et en économie afin d'illustrer son utilité.

Ainsi, dans le cadre de la génomique, cette thèse montre que la longueur des *opérons* bactériens, qui sont à la fois des modules structurels et fonctionnels, est limitée du fait de contraintes dues à l'interaction des mécanismes de *transcription* et de *réplication*. Ensuite, elle fait l'hypothèse que la *modularité structurelle* des *points chauds de recombinaison* correspond au moins partiellement à la *modularité fonctionnelle* des *gènes*. Ceci permet de développer une nouvelle méthode d'analyse des études d'association génétique basée sur un découpage en régions géniques du génome dans le but de faciliter la compréhension du mécanisme fonctionnel de leur action sur le caractère étudié en analysant directement l'association de gènes ou de groupe de gènes avec ce caractère. Sur le plan structurel, les résultats sont d'une qualité comparable à ceux des méthodes classiques. En revanche, le découpage en régions devra encore être affiné afin d'obtenir une analyse fonctionnelle pleinement utile.

Enfin, dans le cadre de la libéralisation du marché européen de l'électricité, la correspondance effective entre structure et fonction de chaque acteur issu de la restructuration fait supposer que le principe de convergence structure-fonction y est bien appliqué. Cependant, des difficultés subsistent avant de parvenir à mettre en place des relations structurelles permettant d'atteindre l'optimum souhaité. Celui-ci inclut des échanges d'énergie à l'origine des contraintes couplantes entre les acteurs. A partir de la théorie de la décomposition par les prix, nous proposons un cadre permettant de définir des tarifs propres à les faciliter, en particulier celles liant producteurs et transporteurs.

En conclusion, cette thèse montre (a) la limite à la convergence structure-fonction que constitue la limite de la longueur des *opérons* bactériens, (b) la faisabilité de l'utilisation d'un découpage

basé sur les limites de gènes afin d'analyser des études d'association génétique à grande échelle et (c) l'importance d'améliorer « la grande boucle » des relations entre producteurs et transporteurs d'électricité afin d'assurer l'optimisation conjointe des investissements en capacité de production et de transport. Elle synthétise l'ensemble de ces résultats dans le cadre conceptuel de la convergence structure-fonction qui postule que la modularité structurelle des *systemes évolutifs et fonctionnels* tend à se superposer à leur modularité fonctionnelle afin de leur apporter robustesse et évolutivité.

Summary

The aim of this thesis is to investigate the structure-function convergence in complex systems by way of applications to living systems and technical-economical systems.

Once having both defined modularity and identified the difficulties coupled to its core definition, this thesis formalizes the structure-function convergence in evolutive and functional systems and illustrates the interest of this concept with regard to the evolvability and robustness of these systems. Furthermore, this concept is applied to open questions in real biological and economic systems.

In the field of genomics, this thesis establishes that the length of bacterial operons, which are structural and functional modules at the same time, is limited by the interactions of transcription and replication mechanisms. Then, this thesis makes the hypothesis that the modular structure defined by recombination hotspots at least partially corresponds to the functional modularity defined by genes. This enables to develop a new method to analyse genetic association studies. It is based on a partition of the genome into bins with boundaries based on gene boundaries. This method renders easier the understanding of the functional mechanism of their action on the studied character. Indeed, it analyses directly the association of individual or group of genes with this character. On the structural level, results are of a quality comparable to those obtained through standard methods. However, the gene based partition will need to be refined in order to obtain a fully useful functional analysis.

Finally, when considering the opening-up of the European electricity market, the correspondence between structure and function of actors issued from the reorganization suggests that the structure-function convergence principle is correctly applied. However, the present structural relationships between actors prevent the system from reaching the desired optimality. This optimality includes energy exchanges which impose coupling constraints on the system. Thanks to the price decomposition theory, we propose a framework to define tariffs useful to improve such relationships, particularly those linking production and transmission operators.

As a conclusion, this thesis shows (a) the limit of structure-function convergence implied by the length limit of bacterial operons, (b) the feasibility of a gene based bin analysis of genome-wide genetic association studies, (c) the importance of improving the relationships between production and transmission operators in order to assume a joint optimization of investments in production and

transmission capacities. This thesis sums up these results in the conceptual framework of structure-function convergence, which postulates that the modular structure of evolutive and functional systems tend to superimpose their functional modularity in order to give them robustness and evolvability.

Table des matières

I	Introduction générale	1
1	Définitions	3
1.1	Modularité structurelle	4
1.1.1	Interface	4
1.1.2	Modularité et complexité	5
1.2	Modularité structurelle des graphes	12
1.2.1	Modèles de complexité	13
1.2.2	Mesures de faible couplage	20
1.2.3	Algorithmes	22
1.2.4	Conclusion	26
1.3	Modularité des systèmes évolutifs et fonctionnels	27
1.3.1	Structure d'un système évolutif	28
1.3.2	Fonction d'un système évolutif fonctionnel	29
1.3.3	Convergence de la modularité structurelle et fonctionnelle	29
1.3.4	Fonction et modularité vis-à-vis de l'environnement	31
1.3.5	Conclusion	32
2	Applications	33
2.1	Exemples de convergence structure-fonction	33
2.1.1	Optimisation mathématique	34
2.1.2	Systèmes économiques	38
2.1.3	Evolution du vivant	40
2.1.4	Conclusion	43
2.2	Apport de la modularité des systèmes évolutifs et fonctionnels	43
2.2.1	Optimisation mathématique	44
2.2.2	Economie	46
2.2.3	Biologie	49

3	Objectifs de la thèse	55
3.1	Recombinaison de l'ADN	56
3.2	Distinction producteur-transporteur dans les systèmes électriques	57
 II Optimisation de la convergence structure-fonction des systèmes vivants par la recombinaison de l'ADN		59
4	Introduction	61
4.1	Le gène : une entité structurale et fonctionnelle	61
4.2	Mutation et gènes superposés	62
4.3	Recombinaison	63
4.3.1	Recombinaison homologue	64
4.3.2	Recombinaison hétérologue inter-espèce	66
4.3.3	Recombinaison hétérologue intra-espèce	67
4.4	Conclusion	67
5	Le biais d'orientation et de longueur des unités de transcription bactériennes	69
5.1	Introduction	69
5.2	Article	70
5.3	Conclusion	70
6	Les études d'association génétique	73
6.1	Introduction	73
6.2	Article	75
6.3	Conclusion	85
7	Conclusion	87
 III Convergence structure-fonction limitée de la distinction producteur-transporteur dans les systèmes électriques		89
8	Introduction	91
8.1	Généralités sur les systèmes électriques	91
8.1.1	Utilité des systèmes électriques	91
8.1.2	Caractéristiques techniques des systèmes électriques	92
8.1.3	Caractéristiques économiques des systèmes électriques	96
8.2	Système électrique européen	97

8.2.1	Caractéristiques techniques	97
8.2.2	Caractéristiques économiques	98
9	La tarification de long terme des réseaux électriques	101
9.1	Introduction	101
9.2	Article	102
9.3	Conclusion	110
9.3.1	Situations possibles et réalisées	110
9.3.2	Congestion et saturation	110
10	Conclusion	111
IV	Discussion et conclusion	113
11	Discussion	115
11.1	Pertinence du concept de modularité	115
11.2	Pertinence du concept de convergence structure-fonction	116
12	Conclusion	117
V	Glossaire et bibliographie	119
	Glossaire	121
	Bibliographie	129

Table des articles

- N. OMONT et F. KÉPÈS – « Developmental modularity and evolutionary canalization », 52–52
Revue de livre, *BioEssays* **27** (2005), p. 667–668.
- N. OMONT et F. KÉPÈS – « Transcription/replication collisions cause bacterial transcription units to be longer on the leading strand of replication », 70–70
Bioinformatics **20** (2004), p. 2729–2725.
- N. OMONT, K. FORNER, M. LAMARINE, G. MARTIN, F. KÉPÈS et J. WOJCIK – « Gene-based bin analysis of genome-wide association studies », 75–85
BMC Proceedings **2 Suppl 4** (2008), S6.
- N. OMONT, A. RENAUD et P. SANDRIN – « Long term nodal pricing and transmission costing », 102–110
Proceedings of the 16th Power Systems Computation Conference, Glasgow, 2008.

Contributions personnelles

Le cœur de cette thèse étant constitué de 4 articles, il est nécessaire de clarifier ma contribution à chacun d'entre eux. De plus, la partie introductive de cette thèse comporte des éléments originaux que je mentionne ici.

L'introduction est un état de l'art sur la relation entre structure et fonction sous la forme d'une synthèse des analyses issues de domaines différents (sous-additivité de l'entropie de Shannon, principe de subsidiarité en économie, modularité dans l'évolution naturelle, algorithmes de décomposition en optimisation). La synthèse en elle-même est fondée sur les travaux de Toussaint et de Schlosser. A ce titre, la seule originalité de cette introduction consiste en l'interprétation des mesures de modularités dans les graphes dans le cadre plus général des mesures de complexité sous-additives. Par une relecture dès les premiers brouillons, Olivier Teytaud a fortement influencé cette introduction. J'ai écrit la revue de livre qui y figure après avoir lu de manière approfondie l'ouvrage « Modularity in Development and Evolution ». François Képès en a amélioré le manuscrit.

Le premier article est le premier à montrer qu'il existe un biais de la longueur des opérons bactériens suivant leur orientation. J'ai eu l'idée d'étudier l'existence de ce biais. François Képès a collecté les données et m'a aidé lors de l'analyse puis lors de la rédaction de l'article.

Le second article ne présente pas d'originalité sur le plan de chaque élément de la méthode, ni sur les questions posés (Comment mesurer des associations entre ADN et maladie ? Comment traduire une association statistique de l'ADN avec une maladie en terme de mécanisme ?). Cependant, au début de la thèse en 2004, il n'existait pas de méthode rigoureuse d'analyse permettant d'obtenir une liste de gènes associés à une maladie à partir d'une étude d'association génétique. J'ai donc développé une méthode permettant d'aboutir à une telle liste en évaluant rigoureusement la qualité des résultats. Les éléments de méthodes existaient déjà (le découpage de l'ADN en régions appelées « bins » était classique chez Genset, le calcul du FDR – False Discovery Rate – était en cours d'introduction dans l'analyse « standard » univariée SNP par SNP). J'ai été l'initiateur de cet assemblage et j'ai réalisé l'essentiel des développements et des calculs. François Képès m'a aidé à identifier les bonnes questions. L'analyse des résultats a été faite par Jérôme Wojcik. J'ai rédigé l'article à l'exception de la discussion, écrite par Jérôme Wojcik. François Képès a amélioré le manuscrit obtenu.

Le troisième article s'inscrit dans une longue tradition des modèle d'optimisation des réseaux électriques. C'est le seul modèle que je connaisse dans lequel il y ait des pertes, pas de seconde loi, et des capacités de lignes à optimiser sur plusieurs scénarios simultanément. Si de tels modèles sont souvent utilisés pour obtenir des prix de court-terme différenciés spatialement et temporellement, c'est la seule fois où ils sont utilisés dans le but de construire des indicateurs tarifaires de long-terme. C'est l'article le moins personnel, dans le sens où les idées sont essentiellement celles d'Alain Hautot, avec des apports d'Arnaud Renaud et Patrick Sandrin. J'ai cependant une part prépondérante dans l'analyse en termes micro-économiques, dans la constitution et l'analyse de l'exemple ainsi que dans la rédaction de l'article.

Première partie

Introduction générale

Chapitre 1

Définitions

A l'origine du mot système, en grec ancien, σύστημα (*sústema*) signifie simplement « ensemble »¹, mais il tire son origine du verbe polysémique συνίστημι (*sunístèmi*) qui signifie au sens propre (I) « placer debout en même temps ». Lorsqu'il est transitif, le verbe peut aussi être synonyme de (II) « constituer, instituer », mais aussi (III) « rassembler, réunir pour une entreprise commune », (IV) « composer, réunir en un tout par l'assemblage des parties, d'où faire naître, créer, produire ». De même, lorsqu'il est introductif, il signifie (I) « se tenir ensemble », (II) « se rapprocher, se rencontrer, en venir aux mains », (III) « être engagé *ou* s'engager dans une liaison amicale, entrer en relation d'où être lié par mariage, être attaché comme ami *ou* disciple, *ou* partisan à », (IV) « se constituer ». A l'époque moderne, ainsi que le mentionne le dictionnaire de la langue française d'Émile Littré², un système est au sens propre un « composé de parties coordonnées entre elles ».

Ces définitions insistent beaucoup sur l'aspect structurel des systèmes : un système est un ensemble d'éléments reliés entre eux. Par exemple, le système solaire est composé d'un ensemble de planètes interagissant par l'intermédiaire des forces de gravitation. Cependant, une lecture attentive laisse poindre la deuxième caractéristique essentielle des systèmes : leur fonction. En effet, la coordination n'est pas un simple lien : c'est un lien pour un objectif. De même, « l'entreprise commune », même s'il s'agit en général « [d'un] combat, [d'une] guerre, [d'une] conjuration » met en valeur l'essence fonctionnelle de l'ensemble constitué.

En conclusion, l'histoire du mot nous révèle deux axes de la signification du mot « système » au développement inégal suivant les contextes :

- une structure : un système est constitué d'un ensemble d'éléments qui lui confèrent ses propriétés ;
- une fonction : un système peut posséder un objectif. On peut mesurer l'adéquation de sa configuration à celui-ci.

1. <http://home.scarlet.be/tabularium/bailly/844.htm>

2. <http://francois.gannaz.free.fr/Littré/xmlittré.php?requete=syst%E8me&submit=Rechercher>

Les sections suivantes explorent ces deux axes – structure dans les Sections (1.1) et (1.2), fonction dans la Section (1.3) – dans la perspective de la notion de modularité.

1.1 Modularité structurelle

Cette partie vise à préciser la notion de modularité structurelle et à développer des mesures permettant de la quantifier.

La structure d'un système est à l'origine de ses propriétés. Elle caractérise l'ensemble des états possibles du système, c'est-à-dire ses configurations. Dans la pratique, un système est un découpage de la réalité. Elle implique la définition :

- d'éléments constitutifs, aux propriétés connues. Ceux-ci pourront eux-mêmes être des systèmes à une échelle plus fine ;
- de relations entre éléments du système, associant leurs propriétés et donnant au système les siennes ;
- de relations avec l'environnement, associant les propriétés de certains éléments du système à celles de l'environnement. Ces relations particulières permettent de caractériser l'interface du système.

De manière générale, on dira qu'un système est modulaire s'il est partitionable en modules dont l'essentiel de la *complexité* n'est pas liée au reste du système. Ceci implique un faible couplage entre un module et le reste du système. Une notion renforcée de la modularité est popularisée dans la notion de « boîte noire » : un module qu'on peut ainsi qualifier possède une interface peu complexe comparée à sa complexité propre. Ceci implique nécessairement un faible couplage. De plus, une « boîte noire » peut être remplacée par une autre à partir du moment où son interface a les mêmes propriétés.

1.1.1 Interface

En absence de découpage, il n'existe qu'un seul système : l'univers entier. Le découpage est indispensable à l'utilité du concept de système. Ainsi que le mentionne Checkland (24), la limite d'un système est une distinction faite par un observateur qui marque la différence entre un système d'étude et son environnement. Ce découpage d'un système unique (l'univers) en deux sous-systèmes (le système d'étude et le reste de l'univers) est extrêmement asymétrique. Il fait nettement apparaître deux interfaces :

- Du point de vue du système d'étude, l'interface du reste de l'univers est ce qu'il est nécessaire d'en connaître pour l'étudier. Par exemple, dans le cas d'une transformation thermodynamique isotherme, il s'agit de la température de l'extérieur. Cette première interface est habituellement nommée sous le terme d'environnement.

- Symétriquement, il existe une seconde interface qui regroupe ce qu'il est nécessaire de connaître du système d'étude afin d'étudier son influence sur le reste du monde. Ce type d'interface se rencontre lors de changements d'échelle : par exemple, si la finesse du comportement d'un atome ne peut être comprise que par la mécanique quantique, les chimistes utilisent couramment des approximations décrivant les orbites possibles des électrons afin d'étudier la formation des molécules. De même, les biologistes se contentent en général de modéliser les *ARN* et les *protéines* comme de longues chaînes afin d'en étudier le repliement dans l'espace. Cette deuxième interface est de type « boîte noire » dans le sens où on ne s'intéresse pas au fonctionnement interne du système mais à ce qu'on en observe depuis l'extérieur et qui est utile pour comprendre son effet sur son environnement.

Poursuivant cette logique de découpage, un système d'étude peut être décomposé en sous-systèmes, aboutissant ainsi à une hiérarchie de systèmes. C'est l'approche réductionniste classique. Inversement, on peut aussi assembler des sous-systèmes dont l'interface est connue pour former un système de niveau intermédiaire (17; 185). C'est l'approche parfois retenue lorsque l'approche réductionniste a découpé trop finement un système en sous-systèmes (182; 2). En effet, dans ce cas, le système constitué des sous-systèmes est trop complexe pour être analysé. Il convient donc d'insérer un niveau de détail en regroupant les sous-systèmes à une échelle intermédiaire. Cette approche particulière est retenue en biologie des systèmes et en *science des systèmes* de manière générale.

1.1.2 Modularité et complexité

1.1.2.1 Modularité d'une partition

Un découpage d'un système par la définition d'interfaces est utile pour son appréhension par la pensée humaine s'il augmente sa simplicité apparente en l'organisant. Ceci constitue un *rasoir d'Occam* dans lequel est « nécessaire » la coupe qui réduit effectivement sa complexité apparente. Toutefois, il est extrêmement difficile d'évaluer la complexité d'un système (1) et encore plus sa complexité apparente car cela nécessite de mesurer notre degré d'ignorance d'un système.

Par exemple, on peut mesurer la complexité par la quantité d'information nécessaire pour décrire un état du système parmi l'ensemble des états auxquels il peut accéder, c'est-à-dire son *entropie* en théorie de l'information. Dans ce cas, la complexité apparente est plus élevée que la complexité réelle car elle intègre notre ignorance à propos du système. De fait, cette ignorance nous conduit à décrire avec plus d'incertitude qu'il n'en existe en réalité les états du système et donc à augmenter la quantité d'information nécessaire à leur description. Dans ce cadre, une coupe en deux parties indépendantes est utile car elle réduit effectivement l'écart entre complexité apparente et complexité réelle. En effet, il est d'autant plus aisé d'évaluer avec précision la complexité d'un système en l'étudiant de manière approfondie qu'il est petit. Il est alors possible d'obtenir une meilleure approximation de

la complexité d'un système pour lequel on a identifié deux parties indépendantes en sommant leurs complexités respectives.

Information mutuelle : Nous cherchons dans ce paragraphe à définir une mesure de la modularité de la coupe d'un système. Nous supposons pour cela que nous disposons d'une mesure de la complexité réelle d'un système S , par exemple à travers son *entropie* $H(S)$. Sans qu'il soit nécessaire que cette mesure de complexité soit une *entropie* au sens de la mécanique statistique ou de la théorie de l'information, il est suffisant pour la définition de modularité que nous exposons ici que la mesure de complexité satisfasse la propriété de sous-additivité pour toute partition \mathcal{P} de S en deux parties notées A et $S \setminus A$:

$$\max(H(A), H(S \setminus A)) \leq H(S) \leq H(A) + H(S \setminus A) \quad (1.1)$$

Cela signifie que, pour tout découpage en deux parties, la complexité de chacun des modules est inférieure à celle du système et que la somme des complexités des modules est supérieure à celle du système. La complexité de Kolmogorov-Chaitin, égale à la longueur du programme de longueur minimale décrivant un objet, vérifie aussi cette propriété (100; 23). Grâce à la sous-additivité, on définit l'*information mutuelle* $I(A, S \setminus A)$:

$$H(S) = H(A) + H(S \setminus A) - I(A, S \setminus A) \quad (1.2)$$

L'information mutuelle est positive et plus petite que la plus petite des complexités des deux parties :

$$0 \leq I(A, S \setminus A) \leq \min(H(A), H(S \setminus A)) \quad (1.3)$$

On note que le cardinal d'un ensemble (son nombre d'éléments dans le cas d'un ensemble fini) définit une structure analogue. En effet, la formule fondamentale pour le calcul des cardinaux est :

$$\text{Card}(C \cup D) = \text{Card}(C) + \text{Card}(D) - \text{Card}(C \cap D) \quad (1.4)$$

Ainsi, on peut représenter de telle mesures de complexité sous la forme de *diagrammes de Venn* où la complexité est analogue à la surface de l'aire correspondant à une courbe fermée. En effet, on peut établir la correspondance suivante :

$$\left\{ \begin{array}{l} H(S) \quad \leftrightarrow \quad C \cup D \\ H(A) \quad \leftrightarrow \quad C \\ H(S \setminus A) \quad \leftrightarrow \quad D \\ I(A, S \setminus A) \quad \leftrightarrow \quad C \cap D \end{array} \right. \quad (1.5)$$

On note que la correspondance ne se fait pas directement avec les ensembles A et $S \setminus A$. Par exemple, $\text{Card}(A \cap S \setminus A) = \text{Card}(\emptyset) = 0 \neq I(A, S \setminus A)$.

On peut aussi décomposer $H(S)$ sous la forme suivante :

$$H(S) = \underbrace{H(A|S \setminus A) + H(S \setminus A|A)}_{\text{Partie modulaire}} + \underbrace{I(A, S \setminus A)}_{\text{Partie anti-modulaire}} \quad (1.6)$$

Où $H(A|S \setminus A)$, dénommée complexité conditionnelle de A sachant $S \setminus A$, est définie par :

$$\underbrace{H(A)}_{\text{Entropie totale de } A} = \underbrace{H(A|S \setminus A)}_{\text{Entropie propre à } A} + \underbrace{I(A, S \setminus A)}_{\text{Entropie partagée avec } S \setminus A} \quad (1.7)$$

Effort d'étude : Voici une interprétation de l'Equation (1.2). Admettons que l'effort à appréhender la complexité réelle $H(S)$ d'un système soit une fonction *convexe*, croissante et positive de celle-ci notée $g(H(S))$. Ceci traduit la limitation de la complexité accessible à la pensée humaine. Cela traduit aussi le fait que la complexité apparente augmente plus vite que la complexité réelle. C'est une idée similaire qui conduit à la définition de la complexité de diamètre d par Chaitin (23) : elle est en effet une mesure de la complexité apparente d'un objet en ne s'autorisant à le décrire que comme une « somme » de parties de taille d au maximum. Cette complexité est plus grande que sa complexité réelle car on ne peut pas prendre en compte les corrélations entre ces parties de taille limitée.

On peut alors écrire :

$$g(H(S)) = g(H(A) + H(S \setminus A) - I(A, S \setminus A)) \quad (1.8)$$

$$\geq g(H(A)) + g(H(S \setminus A)) - g(I(A, S \setminus A)) \quad (1.9)$$

Après découpage, l'effort d'étude du système à travers les études indépendantes des deux parties A et $S \setminus A$ est $g(H(A)) + g(H(S \setminus A))$. On étudie la complexité associée à $I(A, S \setminus A)$ dans chaque partie, car $H(A) = H(A|S \setminus A) + I(A, S \setminus A)$ et $H(S \setminus A) = H(S \setminus A|A) + I(A, S \setminus A)$. Le surcroît d'effort correspondant est à faire deux fois. On ne peut en récupérer une partie à travers l'effort correspondant au terme $-g(I(A, S \setminus A))$. Ce n'est pas le cas si on étudie le système globalement. Le découpage oblige en quelque sorte à étudier deux fois la complexité non modulaire correspondant à l'information mutuelle. Ainsi, il n'est pas nécessairement utile. La coupe \mathcal{P} l'est si :

$$U_g(\mathcal{P}, S) = \frac{g(H(S))}{g(H(A)) + g(H(S \setminus A))} \geq 1 \quad (1.10)$$

Cette condition ne dépend pas seulement du système mais aussi de la forme exacte de g . Cependant, cette inégalité a d'autant plus de chance d'être vérifiée que l'information $I(A, S \setminus A)$ est faible. De plus,

à information mutuelle égale, si g est strictement *convexe*, le gain de simplicité $U_g(\mathcal{P}, S)$ sera d'autant plus important que $H(A)$ et $H(S \setminus A)$ sont égaux, c'est-à-dire que les deux parties sont de même taille.

Par exemple, supposons que l'*information mutuelle* soit nulle ($I(A, S \setminus A) = 0$) et que les deux modules soient de même complexité (ce qui implique $H(A) = H(S \setminus A) = H(S)/2$). Dans ce cas, on a :

$$g(H(S)) \geq g(H(A)) + g(H(S \setminus A)) = 2g\left(\frac{H(S)}{2}\right) \quad (1.11)$$

L'effort d'étude est donc effectivement réduit par un tel partitionnement. Le gain est :

$$U_g(\mathcal{P}) = \frac{g(H(S))}{2g\left(\frac{H(S)}{2}\right)} \geq 1 \quad (1.12)$$

A l'inverse, supposons que l'*information mutuelle* soit maximale :

$$I(A, S \setminus A) = \min(H(A), H(S \setminus A)) \quad (1.13)$$

Supposons aussi que les deux modules soient de même complexité (ce qui implique $H(A) = H(S \setminus A) = H(S)/2$). On a alors :

$$g(H(S)) \leq g(H(A)) + g(H(S \setminus A)) = 2g(H(S)) \quad (1.14)$$

Le découpage est inutile : l'étude de chaque module a la même complexité que l'étude du système. Le surcout dû au découpage est maximum et égal au coût de l'étude du système lui-même. Il n'y a aucun gain de simplicité :

$$U_g(\mathcal{P}) = \frac{1}{2} \leq 1 \quad (1.15)$$

Ces deux exemples extrêmes illustrent le fait que l'intérêt d'un partitionnement pour l'étude d'un système dépend du compromis entre réduction de l'*information mutuelle* et équilibre des partitions.

Modularité faible : Dans cette perspective de diminution de l'effort d'étude, nous nous intéressons d'abord à cet aspect de réduction de l'*information mutuelle*. On définit pour cela la modularité d'une partition d'un système par la fraction de complexité modulaire par rapport à la complexité totale du système. D'après l'Equation (1.6), pour deux parties, on a donc :

$$M^w(\mathcal{P}, S) = \frac{H(A|S \setminus A) + H(S \setminus A|A)}{H(S)} = 1 - \frac{I(A, S \setminus A)}{H(S)} \quad (1.16)$$

En généralisant à un découpage en I parties P_i , on peut définir la modularité par ce score³

$$M^w(\mathcal{P}, S) = \sum_{i \in I} M^w(P_i, S) \quad (1.17)$$

Avec :

$$M^w(P_i, S) = \frac{H(P_i|S \setminus P_i)}{H(S)} = \frac{H(P_i) - I(P_i, S \setminus P_i)}{H(S)} \quad (1.18)$$

Modularité forte : On constate que le numérateur de $M^w(P_i, S)$ dépend de S . On souhaite obtenir une mesure similaire dans laquelle le numérateur est une propriété intrinsèque de P_i . Dans ce but, de manière cohérente avec la Section (1.1.1), on définit l'interface de P_i notée $\Xi(P_i)$ comme la plus petite partie de P_i telle que $I(P_i, S \setminus P_i) = I(\Xi(P_i), S \setminus P_i)$. On peut alors écrire la décomposition suivante de la complexité :

$$H(S) = \underbrace{H(A \setminus \Xi(A) | \Xi(A)) + H((S \setminus A) \setminus \Xi(S \setminus A) | \Xi(S \setminus A))}_{\text{Partie modulaire}} + \underbrace{H(\Xi(A), \Xi(S \setminus A))}_{\text{Partie anti-modulaire}} \quad (1.19)$$

On peut voir l'étude des systèmes de complexité $H(A \setminus \Xi(A) | \Xi(A))$ pour chaque configuration de $\Xi(A)$ comme la caractérisation de la « relation entrée-sortie » du module A . L'étude des systèmes de complexité $H(A|S \setminus A)$ pour chaque configuration de $S \setminus A$ est aussi une caractérisation de cette relation, mais dans le contexte spécifique du système S .

A partir de cette décomposition, on peut aussi définir la modularité par ce score généralisé à un découpage en I parties :

$$M^s(\mathcal{P}, S) = \sum_{i \in I} M^s(P_i, S) \quad (1.20)$$

Avec :

$$M^s(P_i, S) = \frac{H(P_i \setminus \Xi(P_i) | \Xi(P_i))}{H(S)} = \frac{H(P_i) - H(\Xi(P_i))}{H(S)} \quad (1.21)$$

Comme $H(\Xi(P_i)) \geq I(P_i, S \setminus P_i)$, on a :

$$M^s(P_i, S) \leq M^w(P_i, S) \quad (1.22)$$

3. On note que cette généralisation n'est qu'une des généralisations possibles de l'*information mutuelle* à plus de deux ensembles. Il est connu que cette généralisation n'est pas triviale. Il en existe de nombreuses autres, telles la corrélation totale et l'information d'interaction.

Ainsi, on nomme $M^s(P_i, S)$ modularité forte et $M^w(P_i, S)$ modularité faible. Voici un exemple de système modulaire au sens faible, mais pas au sens fort : le système constitué de deux feuilles métalliques minces séparées par un isolant thermique moyen. Il y a un couplage faible entre la température de deux points se faisant face sur chaque feuille, cependant, il n'y pas à proprement parler d'interface entre les deux systèmes : l'étude de l'interaction des interfaces se confond avec celle du système. Le système n'est donc pas modulaire au sens fort.

1.1.2.2 Modularité d'un système

Système complexe – système modulaire : Simon (167) définit les systèmes complexes de la manière suivante :

« [A complex system is] one made up of a large number of parts that interact in a non-simple way. In such systems the whole is more than the sum of the parts, at least in the important pragmatic sense that, given the properties of the parts and the laws of their interaction, it is not a trivial matter to infer the properties of the whole. »

Cela correspond à l'absence d'un découpage tel que l'*information mutuelle* entre les deux systèmes soient la plus réduite possible (couplage faible). Ainsi, la complexité est liée à l'absence de faible couplage, ainsi que l'on nommera désormais la modularité au sens faible. En résumé :

$$\text{Modularité au sens fort} \implies \text{Faible Couplage} \iff \text{Non complexité} \quad (1.23)$$

Mesure perturbée de modularité : Il existe une différence entre la notion de complexité selon Simon et les définitions de modularité exposées jusqu'ici : pour Simon, la complexité est une propriété du système, alors que les modularités définies sont des propriétés du système S et d'une partition \mathcal{P} . On souhaiterait dire que la modularité d'un système est celle de sa coupe la plus fortement modulaire. Cependant, les mesures définies dans le paragraphe précédent présentent un inconvénient : le système « coupé » en un seul module le contenant intégralement (et éventuellement autant de modules vides que souhaités) est parfaitement modulaire d'après la mesure définie dans l'Equation (1.17). En effet, supposons que l'on souhaite une coupe en deux parties. La première est choisie vide : $A = \emptyset$. Sa complexité est $H(A) = 0$. La seconde est identique au système lui-même : $S \setminus A = S$. Ainsi, l'information mutuelle $I(A, S \setminus A)$ est nulle et la modularité est maximale et égale à 1.

C'est pourquoi, dans la pratique, on corrige les mesures de modularité afin d'obtenir des partitions de taille plus équilibrées, plus proche de la définition de Simon et cohérent avec la recherche de coupes du système permettant de réduire l'effort nécessaire pour son étude – cf. Equation (1.10). La taille d'un module A étant caractérisée par sa complexité $H(A)$, on perturbe les mesures de modularité à l'aide d'une fonction $f : [0, 1] \rightarrow \mathbb{R}$ pénalisant les modules représentant une grande proportion de

la complexité du système afin de privilégier les coupes où les modules sont de même taille :

$$M_f^w(\mathcal{P}, S) = \sum_i f\left(\frac{H(P_i)}{H(S)}\right) - \frac{I(P_i, S \setminus P_i)}{H(S)} \quad (1.24)$$

$$M_f^s(\mathcal{P}, S) = \sum_i f\left(\frac{H(P_i)}{H(S)}\right) - \frac{H(\Xi(P_i))}{H(S)}$$

On constate que $M_f^w(\mathcal{P}, S)$ possède certains points communs avec le score $U_g(\mathcal{P}, S)$ défini dans l'Equation (1.10). Pour une bipartition, on peut réécrire $M_f^w(\mathcal{P}, S)$:

$$M_f^w(\mathcal{P}, S) = f\left(\frac{H(A)}{H(S)}\right) + f\left(\frac{H(S \setminus A)}{H(S)}\right) - 2\frac{I(A, S \setminus A)}{H(S)} \quad (1.25)$$

Pour obtenir la meilleure coupe, il faut maximiser $M_f^w(\mathcal{P}, S)$ comme $U(\mathcal{P}, S)$. Entre deux coupes symétriques (dans lesquelles A et $S \setminus A$ sont de même complexité), les deux mesures privilégient celle dont l'information mutuelle est la plus faible. Entre deux coupes de même information mutuelle, les deux mesures privilégient la coupe la plus équilibrée. En revanche, en dehors de ces deux axes, le compromis entre ces deux objectifs est différent.

Voici une proposition de construire une fonction de perturbation f déjà utilisée dans la littérature dans le cadre de la modularité des graphes (129), mais sans être introduite de cette manière. Cette fonction f n'est pas définie de manière univoque, ainsi que le souligne Newman (128).

On souhaite trouver la fonction f qui pénalise les modules de grande taille tout en garantissant que la perturbation introduite soit négligeable dans le cas de modules très petits devant la taille du système. De manière plus formelle, on souhaite donc que f possède les propriétés suivantes :

- $f(0) = 0$ afin qu'un module « vide » ne pèse pas dans l'objectif ;
- $f(\cdot)$ est dérivable en 0 et $f'(0) = 1$ afin que, dans la limite d'un très grand nombre de petits modules, la fonction de régularisation ne perturbe pas l'objectif ;
- $\operatorname{argmax}(f(x)) \geq 1/2$ afin que la bipartition optimale d'un système composé de $2N$ sous-systèmes de même complexité sans aucune interaction soit bien symétrique (i.e. une partition en deux modules comprenant chacun N sous-systèmes) ;
- la fonction f a la *forte concavité* la plus élevée possible afin de favoriser l'obtention des modules les plus équilibrés possibles. En effet, cette *concavité* maximale permet de pénaliser au maximum la taille des modules.

L'ensemble de ces considérations implique :

- f est une parabole car toute autre courbe satisfaisant les conditions a une *forte concavité* moindre : $f(x) = ax^2 + bx + c$;
- $c = 0$ car $f(0) = 0$;
- $b = 1$ car $f'(0) = 1$;

– $a < -b$ car $f'(x) \geq 0$ pour $x \leq 1/2$.

Ainsi, la fonction f de forte concavité la plus élevée vérifiant les conditions est :

$$f(x) = x(1 - x) \quad (1.26)$$

Dans la suite on note :

$$M_r^w = M_{x(1-x)}^w \quad (1.27)$$

$$M_r^s = M_{x(1-x)}^s$$

Fortunato et Barthélemy (48) ont montré que, dans le cadre des graphes, ces mesures de modularité étaient limitées dans leur capacité à trouver les coupes les plus modulaires. Ceci ne semble pas étonnant dans la perspective de ce paragraphe puisque la coupe la plus modulaire de n'importe quel système est le système lui-même. Lors de l'utilisation de telles mesures, il est donc nécessaire d'inspecter chaque module trouvé et éventuellement en réitérant l'algorithme de recherche de partition sur chacun afin de rechercher de possibles sous-modules.

En conclusion, la perturbation des mesures de modularité pose des problèmes, particulièrement celui de définir un compromis entre taille et couplage. Toutefois, elle permet bien de définir utilement la modularité d'un système comme étant celle de sa coupe la plus modulaire.

1.2 Modularité structurelle des graphes

Les dépendances entre parties d'un système sont fréquemment modélisées par des graphes.

- En biologie, il s'agit des cartes d'interaction physique entre macromolécules, ou des réseaux de régulation génétique, ou bien encore des réseaux constitués de paires de gènes dont l'inactivation est létale (135).
- En science de l'ingénieur, les matrices de structure de design (Design Structure Matrix) sont utilisées de manière courante. Elles modélisent les contraintes entre les éléments d'un système (11).
- En optimisation non-linéaire, le graphe des « non zéros structurels » de la hessienne est aussi un objet d'étude.

Du fait de ce grand nombre d'applications pratiques lié au bon compromis entre simplicité de la modélisation sous forme de graphe et capacité à représenter les structures, nous illustrons les concepts de la section précédente sur des graphes non-orientés.

Soit un graphe $G(N, E)$ comptant N nœuds et E arêtes ayant éventuellement des arêtes partant et arrivant au même nœud et des arêtes multiples. On note par $E(X, Y)$ le nombre de liens de l'ensemble de nœuds X vers l'ensemble de nœuds Y .

On souhaite calculer la complexité et la modularité de ce graphe. Pour cela, sa topologie n'est pas suffisante. Nous proposons donc de faire des hypothèses supplémentaires afin de lui faire correspondre une distribution de probabilité. Cette distribution associée permet de calculer sa complexité puis la modularité de ses coupes. Nous passons en revue quatre hypothèses, respectivement nommées « sans cycles », « avec cycles », « basée sur le coefficient de clustering » et « basée sur une constante de diffusion ».

Dans un premier temps – cf. Section (1.2.1), on expose les mesures de complexité associés aux quatre hypothèses. Ensuite – cf. Section (1.2.2), on détaille les mesures de modularité, ou plus exactement de faible couplage associées. Enfin, on conclut – cf. Section (1.2.3) – par une analyse de certaines heuristiques classiques de partitionnement des réseaux, en cherchant à déterminer les définitions de complexité sur lesquelles elles se basent implicitement.

1.2.1 Modèles de complexité

1.2.1.1 Sans cycles :

Soit $d(n)$ le *degré* du nœud n . On attribue à chaque nœud n une variable binaire par arête incidente, notée $v_{e_i}(n) \in \{0, 1\}$ où e_i ($i \in [1, d(n)]$) est une arête dont une des extrémités est le nœud n . On note $V_e(n)$ la variable synthétique combinant l'ensemble des variables du nœud :

$$V_e(n) = \sum_{i \in [1, d(n)]} 2^{i-1} v_{e_i}(n) \quad (1.28)$$

On suppose que le système est tel que, pour tout arête e dont les extrémités sont n_1 et n_2 :

$$v_e(n_1) = v_e(n_2) \quad (1.29)$$

C'est-à-dire qu'il y a couplage complet entre $v_e(n_1)$ et $v_e(n_2)$. De plus, on suppose que toutes les variables binaires sont équiréparties (la probabilité de valoir 0 est égale à celle de valoir 1) :

$$p(v_e(n) = 0) = p(v_e(n) = 1) = \frac{1}{2} \quad (1.30)$$

Sous ces hypothèses, on constate que :

- l'*entropie* de la variable synthétique de chaque nœud est égale au *degré* du nœud ;
- l'*information mutuelle* entre deux nœuds est égale au nombre d'arêtes qui les relient ;
- l'*entropie* d'un groupe de nœuds est égale à la somme du nombre arêtes internes au groupe et du nombre d'arêtes vers le reste du graphe.

$$H(A) = E(A, A) + E(A, G \setminus A) \quad (1.31)$$

- l'*information mutuelle* entre deux groupes de nœuds A et $G \setminus A$ est égale au nombre d'arêtes entre les deux groupes.

$$I(A, G \setminus A) = E(A, G \setminus A) \quad (1.32)$$

Cette mesure de complexité qui est analogue à une *entropie* puisqu'elle est basée sur une distribution de probabilité, se calcule donc simplement en comptant le nombre d'arêtes présentes dans un graphe ou dans un sous-graphe.

1.2.1.2 Avec cycles :

Avec l'hypothèse « sans cycles », la complexité d'une *clique* de N nœuds est la même que la complexité de $N(N - 1)$ nœuds appariés deux à deux, chaque nœud étant relié uniquement au nœud auquel il est apparié. En pratique, la corrélation entre les variables des nœuds de la *clique* est beaucoup plus élevée, ce qui réduit d'autant la complexité. On corrige donc la complexité par le nombre de cycles dans le graphe G noté $C(G)$. On obtient une nouvelle mesure de complexité. On rappelle que, en notant $\kappa(G)$ le nombre de composantes connexes de G , on a :

$$C(G) = E(G, G) - N(G) + \kappa(G) \quad (1.33)$$

On obtient alors que :

- l'*information mutuelle* entre deux nœuds est égale au nombre d'arêtes qui les relie, comme précédemment ;
- la complexité d'un groupe de nœuds A est égale à la somme (a) du nombre de nœuds et (b) du nombre d'arêtes vers le reste du graphe (c) moins le nombre de composantes connexes du sous-graphe constitué des arêtes reliant les nœuds de A :

$$H(A) = N(A) + E(A, G \setminus A) - \kappa(A) \quad (1.34)$$

Où $E(A, G \setminus A)$ est le nombre d'arêtes entre A et le reste du graphe ;

- l'*information mutuelle* entre A et le reste du graphe est :

$$I(A, G \setminus A) = 2E(A, G \setminus A) + \kappa(G) - \kappa(A) - \kappa(G \setminus A) \quad (1.35)$$

Il s'agit bien d'une information mutuelle car on a :

$$I(A, G \setminus A) \geq E(A, G \setminus A) \geq 0 \quad (1.36)$$

En effet, chaque lien entre A et $G \setminus A$ peut au plus couper une seule composante connexe de G en deux lorsqu'on l'enlève. Ainsi, le nombre de composantes connexes du graphe privé des arêtes de A vers $G \setminus A$, c'est-à-dire $\kappa(A) + \kappa(G \setminus A)$, est inférieur à $E(A, G \setminus A) + \kappa(G)$. Ainsi, on a :

$$E(A, G \setminus A) + \kappa(G) \geq \kappa(A) + \kappa(G \setminus A) \quad (1.37)$$

Finalement, on note que l'Equation (1.36) signifie aussi que, pour la même coupe du même graphe, l'*information mutuelle* est plus grande que dans l'hypothèse « sans cycles », ce qui est bien l'objectif initial.

1.2.1.3 Basée sur le coefficient de clustering :

Nous définissons dans ce paragraphe une complexité basée sur le coefficient de clustering, très utilisé dans la littérature (200). Pour calculer ce coefficient, on compte pour chaque nœud k le nombre de paires de voisins qu'il possède. Ainsi, un nœud possédant $d(k)$ voisins compte $d(k)(d(k) - 1)$ paires de voisins. Parmi ces paires de voisins, il en existe un certain nombre $c(k)$ qui sont reliés par une arête. Le coefficient de clustering du nœud est alors :

$$C(k) = \frac{c(k)}{d(k)(d(k) - 1)} \quad (1.38)$$

Un triangle étant un ensemble de 3 nœuds reliés deux à deux (une *clique* de 3 nœuds), il s'agit somme toute du ratio du nombre de triangles ayant leur pointe en k par rapport au nombre de triangles pourraient exister vu les paires existantes de voisins de k . Pour un graphe G , on appelle coefficient de clustering la moyenne pondérée de ces coefficients :

$$C(G) = \frac{\sum_k c(k)}{\sum_k d(k)(d(k) - 1)} \quad (1.39)$$

Par extension, pour un sous-graphe A , le coefficient de clustering est défini de la manière suivante : en incluant uniquement les nœuds du sous-graphe et en conservant $c(k)$ et $d(k)$ mesurés sur le graphe complet.

A partir de là, on définit la complexité suivante :

$$H(A) = \frac{1}{3} \sum_k d(k)(d(k) - 1) - c(k) \quad (1.40)$$

On normalise par $1/3$ afin que $\sum_k c(k)$ soit égal au nombre de triangles dans le graphe. L'*information mutuelle* entre deux sous-graphes est alors égale au nombre de triangles « cassés » par la séparation, car ces triangles sont comptés dans la complexité des deux sous-graphes. Il s'agit d'une complexité

proche de la complexité « avec cycles », car l'élimination d'un triangle élimine aussi un cycle. En revanche, l'élimination d'un cycle n'élimine pas nécessairement un triangle. On note que cette complexité est sensible à l'insertion de nœuds fictifs au milieu d'un lien (en cassant des triangles), alors que ce n'est pas le cas pour la complexité « avec cycles ».

1.2.1.4 Basée sur une constante de diffusion :

D'un tout autre point de vue, imaginons une marche aléatoire sur un graphe : à chaque instant, un marcheur saute d'un nœud à un de ses voisins directs de manière équiprobable. Autrement dit, chaque arête représente une transition possible et toutes les transitions ont la même probabilité.

Après un petit nombre de pas de temps, la probabilité de présence en un nœud du graphe dépend fortement du nœud de départ. Après un très grand nombre, elle n'en dépend plus du tout. Dans ce cas, la probabilité de présence en un nœud x est l'inverse de son *degré* soit $1/d(x)$. Entre les deux, la structure modulaire d'un graphe peut se révéler : en effet, le marcheur aura tendance à rester dans une sous-partie du graphe fortement interconnectée et faiblement connectée avec l'extérieur. On peut chercher à caractériser les ordres de grandeur de ces phénomènes de diffusion ou leurs constantes de temps, qui seront des indicateurs de complexité utiles pour caractériser la modularité de partitionnement.

Valeurs propres d'un graphe : Pour un graphe G connexe, la constante de diffusion est traditionnellement présentée comme étant la première valeur propre non nulle du *laplacien* du graphe⁴. De manière équivalente, on peut la définir de la manière suivante :

$$\lambda_G = \inf_f \frac{\sum_{u \text{ voisin de } v} (f(u) - f(v))^2}{\sum_v f(v)^2} \quad (1.41)$$

$$\sum_v f(v)d(v) = 0$$

On peut généraliser la définition de cette valeur propre à un sous-graphe A . On note \hat{A} l'ensemble des nœuds de A et les nœuds directement voisins de A . On définit la valeur propre de Neumann λ_A du sous-graphe A :

$$\lambda_A = \inf_f \frac{\sum_{u \in \hat{A} \text{ voisin de } v \in A} (f(u) - f(v))^2}{\sum_{v \in A} f(v)^2} \quad (1.42)$$

$$\sum_{v \in A} f(v)d(v) = 0$$

4. Pour une vision plus complète des concepts mathématiques utilisés dans cette section, le lecteur est invité à se référer aux articles de Chung, en particulier la référence (26).

Cette formulation sous la forme d'une minimisation permet de comprendre ce que représente λ_G . L'objectif minimisé pour obtenir λ_G mesure les différences de la valeur de la fonction f entre nœuds voisins. Toutefois, deux contraintes la régissent :

- d'une part, sa norme ne peut être nulle sous peine d'annuler le dénominateur ;
- d'autre part, sa moyenne doit être nulle, ce qui, combiné avec la première contrainte, interdit la solution constante sur l'ensemble des nœuds.

Ainsi, la fonction f de l'optimum sera telle qu'il existera au moins deux parties dans le graphe : une première partie dans laquelle les valeurs de f seront négatives et une deuxième dans laquelle les valeurs de f seront positives. Du fait de l'objectif défini dans l'Equation (1.41), les deux parties sont en général faiblement connectées. On peut même se servir directement de ces vecteurs propres pour rechercher des modules (20).

De plus, on montre que cette valeur propre λ_G permet de borner le temps de diffusion sur le graphe (27). Ainsi, si P est la matrice d'adjacence d'un graphe connexe (telle que $p(i, j) = 1/d(i)$ si i voisin de j et 0 sinon), on sait que la vitesse de convergence vers le vecteur de probabilité stationnaire π du vecteur de probabilité de présence μP^s d'un marcheur aléatoire en chaque nœud du graphe après s pas connaissant son vecteur de probabilité initial μ est bornée : quel que soit ϵ strictement positif, si le nombre de pas s est supérieur à $(2/\lambda_G) \ln(1/\epsilon)$, alors $\|\mu P^s - \pi\| \leq \epsilon$. L'inverse de λ_G est donc bien un temps de diffusion caractéristique du graphe.

Enfin, on montre qu'on peut borner inférieurement cette constante (26) (et supérieurement son inverse) :

$$\frac{1}{\lambda_G} \leq 8d(G)D(G)^2 \quad (1.43)$$

où $d(G)$ désigne le *degré* maximum des nœuds de G et $D(G)$ est le *diamètre* de G , c'est-à-dire la plus grande distance entre deux nœuds de G . On peut généraliser à certains sous-graphes dits « convexes » une borne similaire.

Constante approchée : Les inverses des valeurs propres de Neumann pourraient être de bons candidats pour définir la complexité d'un graphe, toutefois la sous-additivité de ces valeurs propres n'est pas garantie dans le cas général alors qu'elle est nécessaire pour définir une complexité selon l'Equation (1.1). Nous proposons donc de définir une mesure de complexité basée sur la borne supérieure de l'inverse de λ_G donnée dans l'Equation (1.43).

A cette fin, on généralise le *diamètre* aux sous-graphes en définissant $D(A)$ comme la distance maximale entre toutes les paires de nœuds de \hat{A} (nœuds de A et nœuds directement voisins), sachant que la distance entre deux nœuds, mesurée par le chemin le plus court, peut sortir du sous-graphe :

On remarque que le *diamètre* et le *degré* vérifient tous deux la propriété de sous-additivité permettant de les utiliser comme complexité :

- $D(A) \leq D(G)$ car le maximum sur un sous-ensemble est inférieur au maximum sur un ensemble.
- $D(G) \leq D(A) + D(G \setminus A)$ car, soit le chemin le plus long appartient à un des deux sous-graphes (A par exemple) et $D(G) = D(A)$, soit il n'y appartient pas. Dans ce cas, le chemin le plus long passe au moins par une arête frontière. Chaque demi-chemin est lui-même le chemin le plus court entre nœud extrême et nœud frontière dans son sous-graphe. Le *diamètre* des sous-graphes étant supérieur ou égal aux longueurs de ces chemins, la somme des *diamètres* est strictement supérieure au *diamètre* du graphe.
- De même $d(A) \leq d(G)$ et $d(G) = \max(d(A), d(G \setminus A)) \leq d(A) + d(G \setminus A)$.

On peut donc considérer la complexité suivante, qui donne une borne supérieure au temps de diffusion sur A :

$$H(A) = \ln(d(A)D(A)^2) \quad (1.44)$$

On remarque que la complexité d'un graphe non connexe est infinie. Ceci ne constitue pas une limitation importante : la modularité d'un graphe non connexe est maximale, quelle que soit la mesure non-perturbée retenue. Dans ce cas, une étude d'intérêt porterait sur chacun des sous-graphes connexes et non sur le graphe entier.

1.2.1.5 Conclusions

Ces mesures définissent de manière très différente la complexité des graphes. A titre d'exemple, si on considère un graphe de complexité fixée selon l'une des définitions, on peut rechercher les graphes de complexité minimale et maximale selon les autres définitions de complexité. On constate que l'écart entre les minimums et les maximums est souvent très important. Les paragraphes suivants passent en revue la comparaison de ces extremums pour les complexités sans cycle, avec cycle et basée sur la diffusion.

Considérons d'abord un graphe connexe de complexité avec cycle fixée (nombre de nœuds N donné). On peut rechercher les graphes de complexité minimale et maximale selon trois des définitions de complexité données précédemment :

- Sans cycle : le graphe de complexité minimale selon cette définition est un fil de $N - 1$ arêtes (complexité $N - 1$) et le graphe de complexité maximale est une *clique* de $N(N - 1)/2$ arêtes (complexité $N(N - 1)/2$).

- Basée sur la diffusion : supposons N divisible par 3. Un graphe de complexité très élevée⁵, probablement maximale, est une « étoile filante » dont le nœud central est de *degré* $N/3$ et dont une seule des branches est de longueur $2N/3$ quand les autres sont de longueur 1. Sa complexité de diffusion est :

$$\log \left(\frac{N}{3} \left(\frac{2N}{3} + 1 \right)^2 \right) \quad (1.45)$$

Le minimum est une *clique* (complexité de diffusion $\log(N)$).

Supposons maintenant un graphe connexe de E arêtes. De même que précédemment, on peut calculer les graphes de complexité minimale et maximale pour chacune des définitions :

- Avec cycle : le graphe de complexité avec cycle minimale est une *clique* qui compte approximativement \sqrt{E} nœuds (complexité \sqrt{E}). Le graphe de complexité avec cycle maximale est un fil qui compte donc $E + 1$ nœuds (complexité $E + 1$).
- Basée sur la diffusion : le graphe de complexité minimale est inchangé par rapport au cas sans cycle (*clique*). Sa complexité de diffusion est approximativement $\log(\sqrt{E}) = \log(E)/2$. Un graphe de complexité très élevée⁶, probablement maximale, est aussi une « étoile filante ». Supposons E divisible par 3. La queue de l'étoile compte $2E/3$ arêtes et le *degré* du nœud central est $E/3 + 1$. La complexité de diffusion de l'étoile est donc :

$$\log \left(\left(\frac{E}{3} + 1 \right) \left(\frac{2E}{3} + 1 \right)^2 \right) \quad (1.46)$$

Supposons enfin un graphe de produit $\delta = d(G)D(G)^2$ fixé. De même que précédemment, on a :

- Sans cycle : un graphe de complexité faible⁷ est une « étoile filante » comptant un nœud central avec $(\delta/4)^{\frac{1}{3}}$ voisins et une queue de longueur $(2\delta)^{\frac{1}{3}} - 1$. Sa complexité sans cycle est donc approximativement :

$$\left((2\delta)^{\frac{1}{3}} - 1 \right) + \left(\left(\frac{\delta}{4} \right)^{\frac{1}{3}} - 1 \right) = \delta^{\frac{1}{3}} \left(2^{\frac{1}{3}} + 2^{-\frac{2}{3}} \right) - 2 \quad (1.47)$$

Le graphe de complexité sans cycle maximale est une *clique* comptant $\delta/2$ arêtes.

- Avec cycle : un graphe de complexité très faible est une « étoile filante » comme précédemment. En effet, par rapport au cas de la complexité sans cycle, on peut associer un nœud à chaque

5. Supposons un graphe du type « étoile filante ». On maximise le produit $d(G)D(G)^2$ sous la contrainte que la somme du nombre de nœuds dans la queue et dans le cœur de l'étoile soit constante. Le résultat correspond à une approximation dans la mesure où les variables de nombre de nœuds sont considérées comme continues.

6. Raisonnement similaire à celui de la Note (5).

7. Raisonnement similaire à celui de la Note (5), mais en minimisant la somme des arêtes des queues et de l'étoile à δ fixé.

arête, sauf pour le nœud central. Sa complexité avec cycle (en nombre de nœuds) est donc :

$$\left((2\delta)^{\frac{1}{3}} - 1 \right) + \left(\left(\frac{\delta}{4} \right)^{\frac{1}{3}} - 1 \right) + 1 = \delta^{\frac{1}{3}} \left(2^{\frac{1}{3}} + 2^{-\frac{2}{3}} \right) - 1 \quad (1.48)$$

Enfin, un graphe de grande complexité avec cycle est un graphe acyclique où chaque nœud excepté les feuilles a k voisins. On le constitue en partant d'un nœud central, sur lequel on greffe k nœuds, sur lesquels on greffe $k - 1$ nœuds, puis on répète la procédure P fois, obtenant ainsi un *diamètre* $2(P + 1)$. On cherche à maximiser le nombre de nœuds, égal à $k(k - 1)^P$. Comme $k(2(P + 1))^2$ est égal à δ , on substitue k par $\delta / (2(P + 1))^2$ pour obtenir le nombre de nœuds en fonction du nombre de pas effectués P :

$$\frac{\delta}{(2(P + 1))^2} \left(\frac{\delta}{(2(P + 1))^2} - 1 \right)^P \quad (1.49)$$

Cette fonction possède un minimum. En effet, le nombre de pas P est borné car le nombre de nœuds est un entier strictement positif.

Malgré l'existence d'une tendance générale, qui s'exprime par le fait que le minimum et le maximum de la complexité selon une hypothèse tendent vers l'infini lorsque la complexité selon une autre hypothèse tend vers l'infini, on constate la grande hétérogénéité des mesures. Ainsi, on comprend que l'utilisation d'une ou l'autre définition de complexité influe fortement sur les définitions de modularité qui en découlent.

Enfin, il existe certainement d'autres complexités possibles sur les graphes, ce qui illustre parfaitement les difficultés liées à la définition de la complexité en général. Par exemple, on note aussi les quatre modèles exposés constituent des extrêmes. Entre ces extrêmes, on trouve un continuum de modèles de complexité : en effet, toute combinaison linéaire à facteurs positifs des complexités définies est aussi une mesure de complexité.

1.2.2 Mesures de faible couplage

A partir des complexités définies dans la section précédente, on peut construire les mesures de modularité suivantes sur les graphes :

1.2.2.1 Sans cycles :

La définition du faible couplage donnée dans l'Equation (1.24) devient :

$$M_r^{w,n}(P, G) = \sum_i \frac{\sum_j E(P_i, P_j)}{E} \left(1 - \frac{\sum_j E(P_i, P_j)}{E} \right) - \sum_{j \neq i} \frac{E(P_i, P_j)}{E} \quad (1.50)$$

Ce qu'on peut réécrire :

$$M_r^{w,n}(P, G) = \sum_i \frac{E(P_i, P_i)}{E} - \left(\frac{\sum_j E(P_i, P_j)}{E} \right)^2 \quad (1.51)$$

Cette mesure de modularité faible correspond exactement à celle présentée par Girvan et Newman (129). Le cadre des graphes permet une interprétation plus intuitive de cette mesure. En effet, chaque terme est la différence entre la fraction des arêtes internes d'un module et cette fraction après reconnexion aléatoire des arêtes. Cette reconnexion aléatoire conserve constant le nombre d'arêtes du module (plus exactement la somme des degré des nœuds du module) et autorise les arêtes multiples et les arêtes partant et arrivant sur le même nœud. $M_r^{w,n}(P, G)$ varie de $-2/|P|$ pour un graphe comprenant $|P|$ composantes connexes comptant chacune 2 nœuds et 1 lien que l'on découpe en autant de modules que de nœuds (et non de 0 comme on le lit dans Ziv et al. (210)) à $1 - 1/|Z|$ pour un réseau comptant $|Z|$ composantes connexes de même taille. $M_r^w(P, G) = 0$ correspond à une coupe dont le niveau de modularité ne varie pas après une reconnexion aléatoire de ce type. L'interprétation de ce niveau 0 de la modularité perturbée n'est pas aussi évidente dans le cas général d'un système quelconque. Il caractérise le compromis entre équilibre et modularité du partitionnement que réalise cette mesure.

Celle-ci a été ensuite étendue par Radicchi et al. (147) en contraignant tous les nœuds d'un module à avoir plus de liens vers l'intérieur du module que vers l'extérieur.

1.2.2.2 Avec cycles :

La définition du faible couplage donnée dans l'Equation (1.24) devient :

$$M_r^{w,c}(P, G) = \sum_i \frac{N(P_i) - \kappa(P_i)}{N(G) - \kappa(G)} - \left(\frac{N(P_i) - \kappa(P_i) + \sum_{j \neq i} E(P_i, P_j)}{N(G) - \kappa(G)} \right)^2 \quad (1.52)$$

En général, on s'intéresse uniquement à la décomposition de graphes connexes en modules connexes :

$$M_r^{w,c}(P, G) = 1 - \frac{|P| - 1}{N(G) - 1} - \sum_i \left(\frac{N(P_i) - 1 + \sum_{j \neq i} E(P_i, P_j)}{N(G) - 1} \right)^2 \quad (1.53)$$

Cette mesure de modularité favorise la recherche de modules plus « denses » en arêtes que la mesure de modularité n'utilisant pas les cycles.

1.2.2.3 Basée sur le coefficient de clustering :

De manière analogue à la précédente, cette mesure de modularité favorise la recherche de modules denses en triangles et séparés par peu de triangles. L'utilité de cette mesure est liée à la réalité de

l'existence des triangles dans un graphe. Par exemple, il n'y a presque pas de triangles dans un réseau ferroviaire où chaque nœud est une gare. En effet, il y a généralement plusieurs gares entre chaque point de rencontres de lignes. Au contraire, dans un réseau social, les *cliques* ont un sens fort car « les amis de mes amis sont mes amis ». L'étude des triangles peut donc y être pertinente.

1.2.2.4 Basée sur la diffusion :

L'*information mutuelle* entre deux sous-graphes d'une partition est d'autant plus grande que les *diamètres* des sous-graphes sont petits devant le *diamètre* du graphe G . On définit donc pour tout graphe connexe :

$$M_r^{w,d}(P, G) = \sum_i \frac{\ln(d(P_i)D(P_i)^2)}{\ln(d(G)D(G)^2)} \left(1 - \frac{\ln(d(P_i)D(P_i)^2)}{\ln(d(G)D(G)^2)} \right) - \frac{I(P_i, G \setminus P_i)}{\ln(d(G)D(G)^2)} \quad (1.54)$$

Soit :

$$M_r^{w,d}(P, G) = |P| + \sum_i - \left(\frac{\ln(d(P_i)D(P_i)^2)}{\ln(d(G)D(G)^2)} \right)^2 - \frac{\ln(d(G \setminus P_i)D(G \setminus P_i)^2)}{\ln(d(G)D(G)^2)} \quad (1.55)$$

Cette mesure de modularité favorise la recherche de modules de même constante de diffusion (partie quadratique), tout en cherchant à ce que les constantes de diffusion des modules soient les plus élevées possibles (partie linéaire).

1.2.3 Algorithmes

Les hétérogénéités entre mesures de modularité illustrent la difficulté à définir de manière théorique une décomposition modulaire utilisable, c'est-à-dire dans laquelle les modules ont approximativement la même complexité, même sur le modèle simple des graphes. Cependant, cette difficulté est masquée par celle de résolution du problème informatique associé à l'objectif. En effet, les problèmes de coupe dans des graphes, même sous leur forme la plus simple (bisection en deux parties de même taille) sont NP-complets (55).

Ainsi, de nombreuses heuristiques ont été développées depuis les années 1970 afin de trouver des décompositions modulaires de graphes. Dans le vaste domaine des heuristiques de recherche de modules, nous donnons les trois approches suivantes à titre d'exemple pour illustrer la difficulté pratique du problème et l'écart qui existe entre les heuristiques et les objectifs définis dans ce paragraphe. Pour une revue plus complète, on peut se référer utilement à la revue de da F. Costa et al. (44).

1.2.3.1 Centralité des liens

Introduit par Girvan et Newman (129), mais déjà mentionné par Freeman (51), cette classe d'algorithmes basée sur la centralité des liens est divisive : partant du graphe entier, ils enlèvent une à une les arêtes les plus centrales jusqu'à dissocier le graphe en plusieurs composantes connexes. A chaque itération, l'arête la plus centrale est enlevée. Le calcul de la centralité est réévalué après chaque itération. Dans la version initiale, cette centralité se mesure par le nombre de chemins les plus courts entre paires de nœuds qui passent par une arête donnée. De manière alternative, dans un réseau doté d'une loi de répartition des flux comme un réseau électrique, la proportion du flux entre deux points passant par chaque lien peut être utilisée à la place des plus courts chemins. On montre que cette loi spécifique aux réseaux électrique est équivalente à utiliser la probabilité de passage par un lien d'un marcheur aléatoire allant d'un nœud à un autre. Le résultat de l'algorithme est un *dendrogramme* indiquant les coupes successives du graphe, jusqu'à ce que chaque nœud soit isolé. Il faut évaluer par des critères non fournis par l'algorithme (variations de la modularité, etc.) la pertinence des coupes afin de ne retenir que la partie pertinente du *dendrogramme*.

L'algorithme réalise bien un compromis entre équilibre de la bipartition (le nombre maximum de chemins qui peut passer par l'interface est égal au produit du nombre de nœuds dans chaque partition, qui est maximum quand les deux partitions sont de même taille) et taille de l'interface. En effet, la centralité moyenne des liens d'une interface de I liens entre deux partitions de taille N_1 et N_2 est :

$$\bar{c} = \frac{N_1 N_2 / 2}{E} = \frac{N(N-1)\tilde{I}}{E I} \quad (1.56)$$

avec \tilde{I} défini comme le nombre de liens entre les deux modules après reconnexion entièrement aléatoire des arêtes (sans conservation du nombre d'arêtes dans chaque module ni des *degrés* des nœuds). L'algorithme recherche donc des partitions selon des lignes de faiblesses par rapport à ce que l'on attendrait dans un *graphe aléatoire d'Erdős-Rényi* (pour ces graphes, si les arêtes d'un nœud vers lui-même sont autorisées, la probabilité que deux nœuds soient connectés est $N(N-1)/2E$).

Ceci ne correspond pas à la définition de la modularité faible perturbée basée sur un modèle sans cycle $M_r^{s,n}(P, G)$ de la Section (1.2.2.1), car la reconnexion entièrement aléatoire ne tient pas compte du nombre de liens présents initialement dans le module. Ainsi, un module comportant peu de liens internes et une interface vers l'extérieur réduite à un seul lien sera aussi bien détecté que le même module, comportant plus de liens internes, alors que sa modularité est plus faible. Malgré cet écart, Newman propose d'utiliser cette mesure de modularité conjointement à son algorithme. Il nous semblerait plus adapté d'utiliser la modularité faible perturbée basée sur un modèle avec cycle $M_r^{s,c}(P, G)$ de la Section (1.2.2.2), qui cherche aussi à équilibrer les tailles de modules évaluées en nombre de nœuds.

Enfin, on note que l'algorithme utilise la propriété de « modularité » des mesures de modularité : une fois qu'un système est coupé en plusieurs parties, la partition ultérieure d'une partie n'influe pas sur la contribution à la modularité du système des autres parties. En particulier, on peut descendre dans chaque branche de l'arbre et s'arrêter lorsque la modularité de la branche atteint son maximum. On est alors certain d'avoir conservé la partie de l'arbre donnant la modularité la plus grande.

1.2.3.2 Network Information Bottleneck

Introduit par Ziv et al. (210), l'algorithme NIB (« Network Information Bottleneck ») cherche à trouver les modules qui permettent de connaître au mieux la probabilité de présence d'un marcheur aléatoire au temps T sans connaître son point de départ, mais seulement son module de départ. Cet algorithme cherche à regrouper en un nombre de modules fixé a priori les nœuds qui donnent un profil de probabilité de présence semblable au temps T . Il est donc sensible au temps de diffusion sur le graphe et privilégie les structures à faible temps de diffusion (*cliques*). On constate que le choix du temps T est important. En effet, il faut que le temps T soit suffisamment élevé afin que le marcheur aléatoire ait parcouru tout le module plusieurs fois afin d'avoir « oublié » son nœud de départ, mais suffisamment faible pour qu'il n'ait pas parcouru plusieurs fois tout le graphe en ayant « oublié » sa position de départ. Le temps T caractérise donc a priori la taille des modules potentiels. Il doit donc être cohérent avec le nombre de modules recherchés et la taille du graphe. Ziv et al. proposent d'utiliser l'inverse de la première valeur propre non nulle du *laplacien* du graphe – cf. Section (1.2.1.4).

En variant le nombre de modules et à l'aide du calcul de la modularité de la partition obtenue, selon la mesure de modularité faible perturbée sur des graphes sans cycles, on peut trouver un nombre optimal de modules. La mesure utilisée par Ziv (210) est basée sur le modèle « sans-cycle », alors qu'on peut penser qu'une mesure basée sur les temps de diffusion est plus adaptée, même si celle présentée dans la Section (1.2.2.4) peut manquer de finesse pour un tel usage.

Toujours en utilisant la théorie de l'information, l'algorithme de Rosvall et al (155) cherche à maximiser l'*information mutuelle* entre un graphe et sa représentation sous forme modulaire. Dans leur modélisation, ceci revient à minimiser l'*entropie* de la représentation modulaire du graphe. Celle-ci est définie comme étant l'*entropie* d'une variable aléatoire définie sur l'espace des graphes « compatibles » avec la description modulaire. Moins il existe de graphes dans cet espace, plus l'*entropie* de la description modulaire est faible et mieux elle représente le graphe initial. Cependant, afin de trouver des modules, il lui est nécessaire de pénaliser les partitions qui ne sont pas modulaires. En effet, l'algorithme a tendance à regrouper des nœuds qui ont les mêmes voisins plutôt que des nœuds qui appartiennent au même module. Cette pénalisation montre que l'objectif utilisé de « compresser » au mieux la représentation du graphe n'est pas adapté. L'information utilisée à compresser est plutôt celle de l'algorithme NIB : la distribution de probabilité d'un marcheur aléatoire après un certain

nombre de pas de temps. Récemment, Rosvall et al ont d'ailleurs présenté (156) une nouvelle version de leur algorithme cherchant effectivement à compresser la description d'une marche aléatoire sur le graphe d'origine à l'aide d'un découpage en modules.

1.2.3.3 Attaques ciblées

A côté de ces algorithmes visant à identifier des modules, de nombreuses études ont été réalisées afin d'identifier la vitesse à laquelle un graphe connexe se décompose en graphes non connexes lorsqu'on enlève progressivement des nœuds et/ou des liens. Ces études visent aussi à étudier une forme de modularité des graphes.

Le constat initial est le suivant : par rapport à un *graphe aléatoire d'Erdős-Rényi* dans lequel la probabilité qu'une paire de nœuds soit connectée par un lien est uniforme, les réseaux observés dans de très nombreux domaines (12) présentent les caractéristiques suivantes :

- *diamètre* comparable, à peine plus grand que celui d'un graphe aléatoire ayant le même nombre de nœuds et de liens, autrement dit l'effet « small world » ;
- fort coefficient de clustering, bien supérieur à celui d'un graphe aléatoire, indice d'une modularité importante.

Cette concomitance est inattendue. Ainsi, les réseaux réguliers peuvent avoir un coefficient de clustering important, mais ils ont alors un *diamètre* élevé. Par exemple, le graphe correspondant à un pavage du plan avec des triangles équilatéraux possède un coefficient de clustering tendant vers 1, mais un *diamètre* élevé. Dans la plupart des réseaux réels, ces deux caractéristiques se retrouvent ensemble car il suffit qu'il existe quelques nœuds de *degrés* très élevés (« hubs ») reliant des paquets de nœuds très connectés entre eux (« clusters ») (149; 150). Toutefois, il existe d'autres réseaux qui ne possèdent pas cette structure. Par exemple, les réseaux de transport électrique sont beaucoup plus réguliers (ils ont une distribution du *degré* des nœuds plus concentrée), car ils se développent dans un plan (comme il existe peu de croisements de lignes sans nœuds, il s'agit de graphes essentiellement planaires) (154). De même, les réseaux de régulation génétique, qui sont orientés, ont une distribution de *degré* entrant concentrée (61). Ceci signifie que chaque gène est régulé par peu de *protéines*.

Cette structure « hub-cluster » est à l'origine d'un grand nombre d'études visant à caractériser la vitesse à laquelle le graphe se décompose en sous-graphes non-connexes lorsqu'on enlève les « hubs », en la comparant avec cette même vitesse lorsqu'on enlève des nœuds au hasard. Dans un *graphe aléatoire d'Erdős-Rényi*, les deux vitesses sont identiques. Lorsque cela n'est pas le cas, il existe potentiellement une modularité. Un algorithme d'attaque ciblée permet de la mettre en évidence⁸.

8. Il existe cependant un problème technique à résoudre : lorsqu'on enlève un nœud, si le graphe se décompose en plusieurs parties non connexes, on ne sait pas à laquelle des parties rattacher le nœud enlevé. Une solution consiste à remplacer ce nœud par autant de nœuds que de parties qu'il déconnecte et d'ajouter des arêtes de manière à former une *clique* avec l'ensemble de ces nœuds. L'algorithme d'attaque ciblée donne des partitions du graphe modifié de la sorte.

Il a ainsi été montré une grande robustesse des réseaux réels à des attaques aléatoires, mais une robustesse faible aux attaques ciblés. Ce résultat est aussi vrai pour les réseaux de transport électriques, même s'ils ne possèdent pas de structure « hub-cluster » à proprement parler (3; 154). Ce résultat est par exemple très important pour l'étude de la propagation des épidémies, car il stipule qu'il est essentiel de déconnecter un nombre limité de « hubs » de contacts humains pour bloquer la propagation des épidémies.

1.2.4 Conclusion

Malgré la formalisation proposée de la définition de la modularité sur la base de l'*information mutuelle*, le problème théorique est repoussé dans la difficulté à définir puis à mesurer la complexité (1). En effet, comme nous l'avons vu à propos des graphes, il existe de nombreuses manières de la définir pour un même système. Cela conduit à des usages extrêmement variés de la modularité et de la complexité dont il n'est pas évident de comprendre l'unité, même s'il existe une direction générale exprimant le fait que la complexité est liée à la taille du graphe en nombre de nœuds et/ou de liens. Ainsi, il semble que ces deux problèmes de définition, centraux pour la *science des systèmes* soient intrinsèquement liés.

De ce fait, ainsi que Mitchell le souligne à la vue de la diversité des usages de la notion de modularité, il est peu probable de formuler des lois générales à son propos (122). Schlosser (162) fait la même remarque en notant :

« “[M]odularity” is in danger of becoming a buzzword applicable to everything but lacking conceptual force. Modularity needs to be more precisely defined and operational criteria for its application need to be given if the concept is to play an important theoretical role in a new evolutionary-developmental synthesis. »

Face à ce bilan mitigé, nous relevons malgré tout deux points positifs :

- Le constat de fragilité du fonctionnement des réseaux réels face à des attaques ciblés. Par son caractère fonctionnel (la fonction des réseaux est d'assurer la connexité de leurs nœuds), ce résultat anticipe ceux de la section suivante où nous renforcerons la notion de modularité dans les systèmes en l'étudiant dans le cadre des *systèmes évolutifs et fonctionnels*.
- Les nombreux résultats « pratiques » de découpage modulaires de graphes. Le travail de Newman (128) est à cet égard exemplaire. Dans sa recherche de « communautés » dans des graphes (généralement des graphes dits « sociaux » dans lesquels les nœuds sont des personnes et les arêtes des liens entre ces personnes), il se base en ultime ressort sur son intuition pour qualifier ce qu'est une coupe modulaire, jugeant « à l'œil » de sa modularité (selon ses propres mots). Cette approche, éminemment pratique, est efficace pour s'adapter à la connaissance intuitive

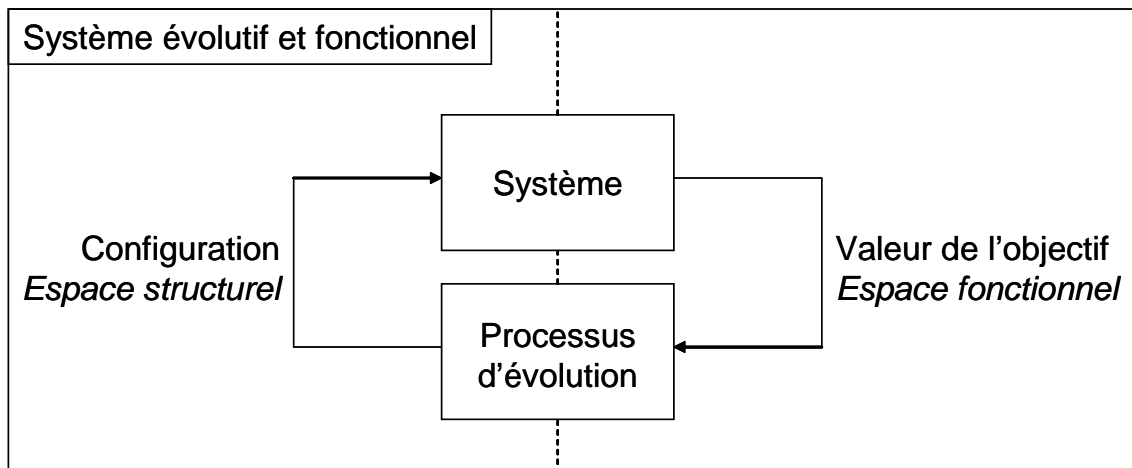


FIGURE 1.1 – Diagramme d'un système évolutif et fonctionnel

des « experts » du système étudié (sociologues, biologistes, etc.), même si elle peine à définir les points communs et les différences entre les modularités recherchées dans chaque système.

Ces approches pratiques sont enrichies par un cadre conceptuel commun tel que nous l'avons ébauché. En effet, il permet la comparaison des différentes notions de modularité et de complexité utilisées, l'évaluation des différents algorithmes face aux déclinaisons de ces notions et une meilleure compréhension de ce qui fait fondamentalement l'unité du concept de modularité.

1.3 Modularité des systèmes évolutifs et fonctionnels

Face aux limites de la définition de la modularité dans le cadre général et surtout face à la difficulté de trouver des propriétés communes à tous les systèmes modulaires, nous proposons d'étudier la modularité d'une classe de systèmes particuliers que nous appelons *systèmes évolutifs et fonctionnels*.

Voici une définition succincte de cette classe de systèmes : supposons qu'un système soit muni d'une fonction à optimiser et d'un processus susceptible de changer sa « configuration » afin d'atteindre ce but – Figure (1.1). Ce nouveau système englobant le système d'origine, le processus d'évolution ainsi que leurs liens par le biais des configurations et de l'objectif est alors à la fois évolutif et fonctionnel. Le principal critère pour décider si un système est bien évolutif et non simplement adaptatif réside certainement dans le fait qu'un système adaptatif possède un domaine de fonctionnement bien défini et étudié à l'avance, tandis qu'un système évolutif est susceptible de « découvrir » des configurations nouvelles. S'inspirant de de Jong (84), cette partie cherche à définir formellement les modularités structurelles et fonctionnelles de tels *systèmes évolutifs et fonctionnels*.

Dans la Section (1.3.1), nous expliciterons une manière de structurer les différentes configurations d'un système de manière à pouvoir y discerner des caractères stables dans le temps et pouvant évoluer de manière indépendante, c'est-à-dire des modules structurels. Dans la Section (1.3.2), nous donnerons la définition de la fonction d'un système et de la modularité fonctionnelle. Ensuite, nous exposerons dans la Section (1.3.3) le point central de cette introduction, à savoir la définition de la convergence structure-fonction et la raison fondamentale de son omniprésence dans les *systèmes évolutifs et fonctionnels*. Nous en tirerons immédiatement les conséquences en termes de limite du concept de fonction dans la Section (1.3.4) avant de conclure.

1.3.1 Structure d'un système évolutif

Un système évolutif est un système qui peut varier, tant du point de vue de la relation entre ses éléments que du point de vue de l'existence même des éléments dans l'ensemble des configurations. En effet, certains éléments peuvent n'exister que dans un sous-ensemble des configurations. Face à de telles possibilités de variation, il convient de formaliser la définition d'un tel système afin de comprendre où sont les continuités dans son évolution.

Comme le font Stadler et al. (172), supposons que le système soit caractérisé par des configurations. Sans structure, un observateur peut seulement dire si deux observations du système correspondent au même état ou non. On suppose donc l'existence d'une pré-topologie sur cet espace des configurations possibles, c'est-à-dire qu'on est capable de définir quelles configurations sont proches d'une configuration donnée (la relation n'est pas nécessairement symétrique). Cette relation permet de rechercher des factorisations locales (194) valables pour un sous-ensemble de l'ensemble des configurations. En substance, une factorisation locale est l'identification de variables ou de caractères pouvant varier indépendamment (ce qui ne veut pas dire qu'ils ne sont pas corrélés) permettant de décrire le système. Les caractères qui ne peuvent être décomposés en caractères de rang inférieur sont dits primitifs. Ce sont ses éléments. Dans un espace vectoriel, la notion de factorisation est équivalente à celle de choix d'une base. Dans un espace de taille finie, la proximité se représente par un graphe dans lequel il existe une arête orientée entre une configuration et les configurations proches. Il existe alors une unique décomposition en facteurs primaires du graphe, c'est-à-dire que le graphe est égal au *produit cartésien des graphes* de ses facteurs primaires.

Ceci définit donc simultanément ce que peuvent être les éléments d'un système évolutif (facteurs primaires) et ses sous-systèmes (regroupements de facteurs).

Enfin, supposons :

- que l'évolution du système soit à temps discret ;
- que la relation de proximité entre configurations soit définie à partir de probabilité de transition du système d'une configuration à l'autre ;
- qu'il existe une probabilité de présence en chaque configuration ;

- qu’il existe une factorisation commune à l’ensemble des configurations.

Dans ce cas, la modularité structurelle d’une partition du système évolutif en sous-systèmes peut donc se mesurer grâce à l’*entropie* de la distribution de probabilité de présence en chaque configuration et aux mesures développées dans la Section (1.1).

1.3.2 Fonction d’un système évolutif fonctionnel

Supposons qu’on définisse une fonction objectif f de l’espace des configurations p du système dans \mathbb{R} . Supposons aussi, comme Toussaint (178), qu’on puisse lui associer sa distribution de probabilité de Boltzmann :

$$F(p) = \frac{\exp f(p)}{Z} \quad (1.57)$$

Pour cela, il faut que la constante de normalisation Z soit finie ($Z = \sum_p \exp f(p) < +\infty$). On constate que non seulement l’ordre mais aussi les valeurs de la fonction objectif sont importantes pour définir cette distribution.

Supposons que l’on parvienne à exprimer la fonction objectif à partir d’une partie des variables issues d’une factorisation du système évolutif. On peut alors former des partitions à partir de ces variables et mesurer la modularité de ces partitions vis-à-vis de la distribution de Boltzmann de la fonction objectif. On parle alors de modularité fonctionnelle.

1.3.3 Convergence de la modularité structurelle et fonctionnelle

Toussaint et von Seelen (179) donnent un cadre conceptuel pour comprendre pourquoi il est important que la modularité structurelle et la modularité fonctionnelle se recourent dans les *systèmes évolutifs et fonctionnels*.

Soit g une configuration du système. Soit p la configuration « fonctionnelle » correspondante c’est-à-dire que p se déduit de g en ne retenant que des variables qui sont nécessaires à l’évaluation de la fonction objectif. Par conséquent, il existe plusieurs configurations g pour une seule configuration p . Autrement dit, il existe des évolutions « neutres » du point de vue de la fonction objectif, mais non du point de vue de l’évolutivité du système, comme nous allons le voir. Ce type d’évolution est fréquent dans les systèmes reproductifs vivants (159). La première *théorie de l’évolution neutre* a été écrite en 1968 par Kimura (95). Certains suggèrent qu’elle s’applique aussi pour les systèmes technologiques (114).

La valeur de la fonction objectif en g est $f(p)$. On peut définir d’autres fonctions objectifs associées à g . En particulier, étant donnée la distribution de probabilité $\sigma_g(p)$ de transition depuis une configuration structurelle quelconque g vers une configuration fonctionnelle quelconque p (variabilité dite fonctionnelle), l’espérance de la valeur de la fonction objectif après un pas d’évolution (*fonction*

de fitness effective), notée $\langle f \rangle_{\sigma_g}(g)$, qui caractérise l'évolutivité de g s'écrit :

$$\langle f \rangle_{\sigma_g}(g) = \sum_p \sigma_g(p) f(p) \quad (1.58)$$

On constate que cette fonction objectif dépend maintenant des probabilité de transition en g . Ainsi, toutes les configurations du système qui avaient pour même valeur d'objectif $f(p)$ n'ont plus nécessairement la même valeur d'objectif effectif.

En utilisant la distribution de Boltzmann de la fonction f , on peut réécrire cette équation :

$$\begin{aligned} \langle f \rangle_{\sigma_g}(g) &= \sum_p \sigma_g(p) \ln(F(p)) + \ln(Z) \\ &= -D(\sigma_g \| F) - H(\sigma_g) + \text{const} \end{aligned} \quad (1.59)$$

Dans cette équation, $D(\sigma_g \| F)$ représente la divergence de Kullback-Leibler entre les distributions σ_g et F et $H(\sigma_g)$ est l'entropie de la distribution σ_g ⁹.

De manière triviale, on peut vérifier que cette fonction de fitness effective est maximale si on est certain d'atteindre la valeur maximum de f en quittant g . En effet, dans ce cas, $H(\sigma_g) = 0$ et $-D(\sigma_g \| F) = \max_{g_1} (f(g_1))$. La convergence structure-fonction n'a alors aucune importance. Cela peut arriver lorsque le processus d'évolution est parfaitement adapté à la fonction à optimiser (régulation, etc.). En effet, a contrario, un processus évolutionnaire ne peut garantir d'aboutir à l'optimum en une seule étape que s'il est conçu pour résoudre un seul problème d'optimisation. Même s'il est déterministe ($H(\sigma_g) = 0$), la transition vers la configuration suivante est limitée (par des contraintes technologiques par exemple) et/ou choisie de manière à obtenir la plus grande fitness effective en moyenne sur l'ensemble des problèmes qu'il est susceptible de résoudre et qu'il est incapable de différencier avec l'information dont il dispose. La transition ne consiste donc pas à atteindre immédiatement et à coup sur l'optimum. Pour cette raison, dans le second cas, la connaissance des modularités de la fonction objectif peut déjà l'aider à réduire au minimum l'ensemble des problèmes qu'il est encore incapable de différencier à un instant donné. Par exemple, en optimisation non-linéaire, connaître la hessienne (et sa structure modulaire) en plus du gradient permet de calculer un pas qui sera en moyenne plus efficace qu'en connaissant le gradient seul – cf. Section (2.1.1.2).

Cependant, c'est lorsque la variabilité fonctionnelle σ_g des systèmes évolutifs et fonctionnels n'est pas nulle que la convergence structure-fonction joue le plus grand rôle. Ainsi que le mentionne Tous-saint, plusieurs raisons non exclusives l'une de l'autre peuvent expliquer que la variabilité fonctionnelle d'un système soit positive :

9. $D(p \| q) + H(p)$ est aussi appelée *entropie croisée* de p et q . Elle est aussi définie dans le cas continu, ce qui assure la généralisation de cette équation au cas où il existe une continuité de variations possibles à partir de g .

- Variabilité structurelle irréductible : le processus d'évolution n'est pas parfait. Par exemple, des erreurs s'introduisent dans la copie du matériel génétique. Le système global doit être robuste, c'est-à-dire conserver une fonction de fitness effective élevée malgré cette variabilité structurelle subie. Cette variabilité induit une variabilité fonctionnelle σ_g qu'il faut structurer au mieux.
- Variabilité mesurée de la fonction objectif : le système n'est pas déterministe. A configuration égale, la valeur de l'objectif est variable, mais selon des probabilités connues. Le système global doit être robuste à cette incertitude intrinsèque. La variabilité fonctionnelle est une stratégie de diversification pour y parvenir.
- Variabilité complète de la fonction objectif : le système doit rester flexible. La fonction objectif est elle aussi variable. Le système global doit être robuste à ces variations de fonction objectif, par exemple dues à des changements environnementaux. Le système doit pouvoir évoluer rapidement vers une configuration de *fitness* élevée pour la nouvelle fonction objectif. La variabilité fonctionnelle est un moyen privilégié d'y parvenir en explorant de nouvelles solutions. Ce moyen est d'autant plus efficace que la variation de l'objectif est modulaire. C'est par exemple le cas lorsque seule une partie de la fonction objectif évolue.

Admettons que cette variabilité fonctionnelle soit d'entropie fixée $H(\sigma_g)$. Dans ce cas, la « bonne » variabilité est celle qui possède la moindre divergence avec la distribution de probabilité associée à l'objectif, c'est-à-dire $D(\sigma_g||F) = 0$. Le processus d'évolution sera d'autant plus efficace que la divergence est faible, ce qui implique en particulier de calquer la modularité structurelle sur la modularité fonctionnelle, quel que soit la manière dont celle-ci est définie. En particulier on note qu'une séparabilité de la fonction objectif en plusieurs ensembles de variables (fonction objectif décomposable en sommes telle que chaque variable ne soit présente que dans une somme) correspond à un produit de distribution de Boltzmann. C'est pourquoi les probabilités de variation structurelle de ces mêmes ensembles de variables sont idéalement indépendantes les unes des autres. C'est ce que nous appellerons la convergence structure-fonction.

Celle-ci permet de garantir un ensemble de robustesses du système, en particulier une certaine invariance fonctionnelle face à des variations structurelles ou une certaine invariance structurelle face à une variation fonctionnelle. De fait, comme nous le développerons sur des exemples choisis dans la Section (2.2), elle joue un rôle essentiel en permettant des évolutions indépendantes par modules qui évitent les effets d'un module sur un autre et permettent d'introduire des nouvelles solutions sans déstabiliser l'existant puisque la fonction objectif possède une modularité identique.

1.3.4 Fonction et modularité vis-à-vis de l'environnement

On peut déjà souligner un point d'importance dans le débat sur la possibilité d'attribuer une fonction à un système : sans cette *convergence*, il n'est pas possible d'attribuer une fonction à une

structure. Ainsi, l'existence même d'un système d'étude muni d'une fonction objectif suppose que ce système soit modulaire structurellement et fonctionnellement vis-à-vis de son environnement.

Ce problème avait déjà été identifié, par exemple par Langlois (106) ou par Watson (199) : il y a des limites à la définition de la fonction d'un système, car celle-ci suppose l'indépendance du système par rapport à son environnement, ce qui est rarement parfaitement le cas. De manière imagée, on peut voir l'ours blanc comme un organisme adapté à l'Arctique tandis que l'Arctique est un environnement constant. Il en est souvent ainsi, c'est-à-dire qu'il est généralement possible de considérer qu'un système possède une fonction objectif lié à un environnement constant car le système a peu d'influence sur son environnement. Cependant, il arrive que « l'ours blanc modifie l'Arctique ». C'est-à-dire que le système influe sur son environnement de manière sensible. Cela signifie qu'il est impossible de définir une fonction objectif du système indépendamment des variations de l'environnement : il est nécessaire d'inclure une représentation des interactions de l'environnement et du système dans le système étudié. Il est donc nécessaire d'étudier le système non modulaire constitué du système initial et de son environnement. Ceci illustre parfaitement la réflexion théorique de Toussaint (179) sur la nécessité d'une modularité à la fois structurelle et fonctionnelle d'un système vis-à-vis de son environnement : il est impossible d'optimiser un système sur la base de sa seule fonction objectif si elle ne dépend pas seulement de sa configuration. De fait, la fonction objectif partielle définie par rapport à un environnement donné n'a pas de sens pour l'évolution du système en cas d'influence mutuelle.

A l'inverse, un système isolé ne peut pas avoir de fonction : La fonction n'est pas intrinsèque au système. Elle dépend de son intégration dans un système fonctionnel plus vaste. Par exemple, un livre épais constituera un bon moyen de caler un meuble bancal. On ne peut cependant pas dire que le livre « contienne » sa fonction : elle lui est donnée par la personne qui l'a utilisé comme une simple cale. Ainsi, tant du point de vue structurel que fonctionnel, la recherche de modules dans un *système évolutif et fonctionnel* est nécessairement un équilibre entre intégration fonctionnelle et indépendance structurelle (64).

1.3.5 Conclusion

Cette partie a montré l'existence d'un lien profond entre variabilité de la structure et séparabilité de la fonction et modularité dans les *systèmes évolutifs et fonctionnels*. Elle a aussi montré les limites à la définition de la fonction dans ces systèmes. Dans les parties suivantes, on déclinera le principe général de convergence structure-fonction pour différents *systèmes évolutifs et fonctionnels*. Après avoir décrit ces systèmes et mis en lumière cette *convergence*, on s'intéressera particulièrement à ses conséquences en termes de robustesse et de flexibilité.

Chapitre 2

Applications

Dans le Chapitre (1), nous avons explicité de manière formelle le lien entre modularité structurelle et la modularité fonctionnelle. Ce chapitre vise à l'illustrer à travers des exemples de *systèmes évolutifs et fonctionnels*. Nous nous intéressons à trois classes de tels systèmes :

- les méthodes mathématiques d'optimisation, qui peuvent parfois s'inspirer des systèmes suivants, ou bien être utilisées pour les modéliser ;
- les systèmes techniques et économiques ;
- l'évolution du vivant (et non le développement d'un organisme vivant).

Dans un premier temps, ce chapitre explicite dans quelle mesure il est possible d'identifier la fonction et les mécanismes d'évolution des systèmes des deux premiers types, c'est-à-dire dans quelle mesure ce sont des *systèmes évolutifs et fonctionnels*. On montre aussi dans quelle mesure ces mécanismes assurent la convergence structure-fonction. Dans un deuxième temps, ce chapitre montre comment cette *convergence*, généralement alliée à une variabilité fonctionnelle non nulle, apporte robustesse et flexibilité à ces systèmes.

2.1 Exemples de convergence structure-fonction

Les *systèmes évolutifs et fonctionnels* ont tendance à posséder des modularités à la fois structurelles et fonctionnelles. On peut donner par des exemples un aperçu des développements des sections suivantes :

- Certains groupement de gènes apportent une contribution partiellement séparable à la fonction de fitness de l'organisme qui les porte. Par exemple, les gènes à l'origine de l'établissement des segments de polarité chez la drosophile forment un module évoluant de manière relativement indépendante (115; 36).

- Dans les sociétés humaines, la monnaie est l'instrument par excellence de modularisation de l'économie. En effet, elle permet de découpler la gestion des limitations de ressources entre les différents acteurs d'une société. Son objectif est de décomposer le problème d'optimisation de l'utilisation des ressources de façon à ce que la résolution, par chacun, de son sous-problème consistant à obtenir « le plus » de l'argent dont il dispose, conduise à l'optimum global.
- Les mathématiciens et les informaticiens cherchent toujours à « exploiter la structure » du problème d'optimisation qu'ils essaient de résoudre. Ils parlent en fait d'exploiter structurellement la modularité fonctionnelle du problème. En général, cela commence par choisir des variables les plus indépendantes possibles les unes des autres. Ceci revient souvent à décomposer le problème en problèmes qu'il est possible de résoudre de manière relativement indépendante, par exemple liés par un petit nombre de contraintes couplantes. Ainsi, les méthodes de décomposition par les prix utilisent la modularité en s'inspirant du rôle de la monnaie dans l'économie : elles cherchent à trouver le « prix » de chaque contrainte couplante afin que la résolution de chaque sous-problème intégrant les « prix » de ces contraintes mène à l'optimum global. De même les algorithmes dit « génétiques » cherchent à découvrir la modularité de la fonction à optimiser (65; 198; 75) en utilisant des méthodes s'inspirant notamment de l'évolution naturelle (68; 183). En particulier, les opérateurs de *recombinaison* s'inspirent de la *recombinaison* génétique, qui est un des processus évolutifs de gestion de la modularité. En effet, allié à la sélection, ils permettent de regrouper sur un même chromosome puis de transmettre ensemble les informations couplées du système tout en séparant les informations indépendantes du point de vue de l'objectif par des *recombinaisons*. Ainsi, l'évolution est biaisée vers l'exploration des solutions accessibles par des *recombinaisons* de segments d'information (modules structurels) indépendants fonctionnellement.

Cependant, même si cela relève plus de l'exception que de la règle, tous les *systèmes évolutifs et fonctionnels* ne sont pas modulaires. Ainsi dans la plupart des réseaux de neurones artificiels, il est impossible de donner un sens aux calculs réalisés par les couches intermédiaires, qui dépendent du bruit stochastique de l'apprentissage (165; 7). De même en neurologie, il existe des cas où il est impossible de lier une structure et une fonction. Par exemple, on montre que la modularité fonctionnelle n'est pas nécessairement exploitée sous forme de modularité structurelle (140).

2.1.1 Optimisation mathématique

En optimisation mathématique, on s'intéresse essentiellement à des algorithmes qui font évoluer une ou plusieurs solutions à chaque pas de temps dans la mémoire d'un ordinateur. Nous nous concentrons sur le cas où la fonction objectif est explicite.

Dans ce cadre, l'utilisation de la « structure » du problème (qui est en fait sa modularité fonctionnelle) est souvent la clé des méthodes de résolutions efficaces. Il s'agit en général de trouver une

représentation (structurelle) du problème ajustée à sa fonction objectif. Même pour les méthodes dans lesquelles il est difficile d'introduire des modularités structurelles, tels les réseaux de neurones, des moyens de le faire ont été développés, en particulier dans le cas où la fonction à apprendre possède des sorties intermédiaires (14) ou lorsque le réseau est dynamique et que ses entrées au pas de temps $t + 1$ dépendent de ses sorties au pas de temps t (132; 66).

Ainsi que le mentionne de Jong et al. (85), la séparabilité définit la modularité fonctionnelle du problème d'optimisation, tandis que la représentation des solutions définit la modularité structurelle de sa résolution. En effet, le choix des variables du problème et éventuellement de sa décomposition en sous-problèmes revient à définir ce qui peut varier indépendamment. Le problème de *convergence* est donc intrinsèquement lié à celui de représentation : certaines représentations auront la modularité structurelle du système calquée sur la modularité fonctionnelle de l'objectif qui lui est assigné, d'autres non.

Par exemple, en optimisation combinatoire, l'existence de modules optimisables indépendamment de leur contexte permet de réduire fortement la complexité du problème. De même, en optimisation continue, que ce soit par décomposition en sous-problèmes ou bien dans la prise en compte des zéros de la hessienne, on cherche à représenter dans la structure du système à optimiser les séparabilités de la fonction objectif, c'est-à-dire à obtenir une hessienne qui compte le moins possible d'éléments hors diagonale et que ceux qui le sont permettent de former une structure diagonale par bloc. Nous détaillons ces deux cas dans les sections suivantes.

2.1.1.1 Optimisation combinatoire

Concepts : A cheval entre l'optimisation informatique et la théorie évolutionnaire, de Jong et al. (85) définissent le module fonctionnel par rapport à sa contribution à la fonction du système. Un module est caractérisé par le fait qu'on peut l'optimiser de manière relativement indépendante de son environnement. En effet, dans le cas contraire, sa configuration optimale serait différente pour chaque configuration de l'environnement. Dans un système séparable, chaque module peut être optimisé séparément : l'optimum global est trouvé à partir des optimums de chaque module. Cette condition est trop forte et de Jong et al. l'assouplissent en définissant la décomposabilité : l'optimisation de chaque module peut en partie dépendre de la configuration des autres modules du système. Ceci laisse la place pour la cooptimisation ou la coévolution des modules, afin de rechercher un optimum couplant plusieurs modules.

Cette définition ne fait pas d'hypothèses sur les mécanismes d'évolution. Toutefois, pour vérifier que cette modularité correspond bien à une convergence structure-fonction au sens de Toussaint –Section (1.3), il faut expliciter un processus d'évolution typique des algorithmes combinatoires. Par exemple, supposons que l'algorithme soit conçu de telle façon qu'il soit possible de remplacer entièrement la configuration d'un module par une autre en une seule itération. Supposons de plus que

l'algorithme soit à même de sélectionner, parmi les configurations d'un module, seulement celles qui sont optimales dans au moins un contexte (une configuration des autres modules). S'il existe effectivement des configurations qui ne sont optimales dans aucun contexte, l'algorithme augmente sa probabilité d'atteindre l'optimum en considérant uniquement des évolutions consistant à changer la configuration de modules. Les algorithmes génétiques avec opérateurs de recombinaison sont conçus dans cet objectif. En effet, supposons que, grâce à la sélection, les configurations optimales d'un module soient surreprésentées dans la population. Dans ce cas, l'opérateur de recombinaison peut assurer le remplacement d'un module par un autre. Il agira de manière imparfaite, n'ayant pas la connaissance explicite des modules, mais cela peut cependant être suffisant pour orienter la sélection.

Mesure : supposons le processus évolutif suivant : Soit un système à N_S configurations et O_S optimums. Supposons une évolution entièrement aléatoire. Dans ce cas, la probabilité d'aboutir à une solution optimale depuis l'état courant est uniforme : O_S/N_S . Avec un objectif égal à 1 si la configuration est optimale et 0 sinon, la fonction de fitness effective est $O_S/N_S \leq 1$.

Supposons maintenant que l'évolution n'explore aléatoirement que les configurations où l'ensemble des modules sont dans une configuration optimale dans au moins un contexte. Dans un cas à deux modules A et $S \setminus A$ comptant respectivement O_A et $O_{S \setminus A}$ optimums, on obtient la fonction de fitness effective suivante : $O_S / (O_A O_{S \setminus A}) \leq 1$.

De plus, on a : $N_S = N_A N_{S \setminus A}$. Notons $H(S) = -\log(O_S/N_S)$ le logarithme de la fonction de fitness effective. Comme $N_A \geq O_A$ et $N_{S \setminus A} \geq O_{S \setminus A}$, on a :

$$H(S) \leq H(A) + H(S \setminus A) \quad (2.1)$$

De plus, $H(S) \geq \min(H(A), H(S \setminus A))$. La fonction H définit donc une mesure de complexité au sens de l'Equation (1.1).

De cette définition, on en déduit aussitôt les mesures de modularité développée dans la Section (1.1). Ainsi, un sous-système est d'autant plus complexe qu'il a peu de solutions optimales (du type « aiguille dans une botte de foin ») et beaucoup d'interactions avec le reste du système. Enfin, supposons que le problème d'optimisation S ne soit pas dégénéré ($\log(O_S) = 0$). Dans ce cas, $I(A, S \setminus A) = \log(O_A) + \log(O_{S \setminus A})$. Le système est donc d'autant plus modulaire que le nombre d'optimums partiels de A et de $S \setminus A$ est petit.

2.1.1.2 Optimisation non-linéaire

Concepts : en optimisation non-linéaire, on recherche un minimum local d'une fonction objectif en partant d'un point initial. Les algorithmes cherchent en général un chemin strictement décroissant

vers l'objectif¹. Ils font en général face à une difficulté majeure : choisir le pas à faire (norme et direction) de manière à maximiser l'espérance de diminution de l'objectif.

Lorsque la fonction est différentiable, la direction de plus grande pente est en général une direction de descente adéquate. En effet, l'existence d'un gradient non nul implique l'existence d'un voisinage où la fonction objectif a une valeur inférieure. Le problème de non-modularité survient au niveau du choix de la norme du vecteur de descente. En effet, afin de minimiser le nombre de pas pour atteindre l'objectif, il faut faire les plus grands pas possibles. D'un autre côté, l'information contenue dans le gradient n'est valable que dans un voisinage du point actuel. La hessienne de l'objectif permet de déterminer cette norme. En particulier, l'existence d'une hessienne bloc-diagonale indique l'existence de modules fonctionnels correctement représentés structurellement (au moins localement). L'algorithme en tire parti car il peut calculer indépendamment les normes optimales du pas dans les sous-espaces de chaque module, s'éloignant ainsi de la stricte direction de plus grande pente (algorithmes de gradient conjugué par exemple (158)). Au contraire, à l'intérieur d'un module, il lui faut tenir compte de la force des interactions de chaque variable pour déterminer le voisinage de validité de son gradient. Enfin, si les zéros sont globaux, le calcul de la hessienne s'en trouve accéléré.

L'existence de zéros locaux ou globaux de la hessienne est donc l'expression de la convergence structure-fonction pour un problème non-linéaire. Elle permet, comme en optimisation combinatoire d'avancer indépendamment la résolution des différentes parties du problème.

Mesure : à moins de calculer la distribution de probabilité de succès d'un pas partant d'un gradient et d'une hessienne donnée pour l'ensemble des fonctions à laquelle peut s'appliquer un algorithme, il semble difficile d'utiliser la mesure de la fonction de fitness effective pour calculer la modularité d'un problème.

Frimannslund (52) définit le graphe de covariation qui permet de généraliser la notion de hessienne creuse à des problèmes non différentiables. Le graphe de covariation étant un graphe non orienté, on peut appliquer les mesures de modularité définies dans la section (1.2). La question se pose de savoir quelle est la mesure la plus adaptée. Le gain en temps de calcul d'une hessienne est d'autant plus grand que le nombre de cases non-vides hors diagonale est grand devant le nombre total de cases. Ainsi, la modularité basée sur le nombre de liens (« sans cycles ») semble être une mesure adaptée si cet effet est effectivement la source principale de gain.

1. Ceci est faux pour la recherche d'optimums globaux. Par exemples, les algorithmes évolutionnaires possédant des stratégies d'évolutions telles qu'on autorise une croissance de l'objectif sont réputées plus robustes que ceux utilisant les stratégies ordinaires (83).

2.1.2 Systèmes économiques

En économie, on étudie les systèmes de productions de biens et de services. Ceux-ci évoluent à chaque instant suite aux décisions des acteurs et aux événements exogènes. Ils voient leur fonction évoluer au cours du temps suite à des ruptures technologiques ou à des modifications de préférences des acteurs.

Dans un premier temps, nous montrerons en quoi la différenciation entre structure et fonction demande une attention particulière et quelles en sont les limites. Ensuite, nous montrerons comment les *externalités* correspondent à des absences de convergence structure-fonction. Nous terminerons avec l'impact de changement de fonction objectif sur la structure du système à travers des exemples de ruptures technologiques.

2.1.2.1 Structure et fonction

Pour la plupart des systèmes conçus par les hommes, la fonction est évidente car leur but est de répondre à un objectif préexistant. L'*analyse fonctionnelle* a même été théorisée par les ingénieurs et est maintenant enseignée dès le lycée². Cependant, définir une finalité est controversé dans les systèmes économiques. En effet, si chaque acteur peut éventuellement ordonner ses préférences sur la configuration de la partie qui le concerne, l'agrégation de l'ensemble de ces préférences constitue un sujet en soi car il s'agit d'arbitrer entre préférences contradictoires voire incompatibles des acteurs. Par exemple, la science ne peut donner qu'un éclairage sur l'équilibre entre richesse moyenne et écarts de richesse.

En tout état de cause, la structure d'un système économique est l'ensemble de ce qui régit son évolution (acteurs et règles attachées) tandis que sa fonction évalue la production réalisée compte tenu des ressources disponibles. Ces deux éléments ne sont pas entièrement stables dans le temps. En effet, les contraintes technologiques évoluent, augmentant (ou diminuant) le domaine du possible (changeant ainsi la fonction du système) et les capacités de gestion de l'information (donc les formes possibles d'organisation et ainsi finalement les structures possibles). De plus, les organisations participent aussi à cette modification de leur environnement, ce qui limite la possibilité de définir leur fonction, ainsi que nous le mentionnons de manière générale dans la Section (1.3.4). Par exemple, les organisations ne s'adaptent pas toujours à la technologie. Ainsi, Ford, en démodularisant l'industrie automobile a permis les sauts technologiques aboutissants à la Ford T, sans avoir à se concerter avec des partenaires. Au contraire, du point de vue de la modularisation mais toujours à titre d'exemple d'une action d'un système sur son environnement, en concevant un PC modulaire, IBM a permis malgré lui la naissance de l'industrie informatique modulaire d'aujourd'hui. Dans les deux cas, ils ont fortement modifié leur environnement technologique par leurs décisions. Ils ont ensuite dû changer

2. <http://www.education.gouv.fr/bo/2000/hs8/ing.htm>

de stratégie car, si leur objectif était fixe (gagner de l'argent), leur « configuration » optimale pour y parvenir avait changé. Ainsi, après Ford, l'industrie automobile s'est à nouveau démodularisée avec l'apparition de nombreux sous-traitants, et IBM a perdu son rôle de leader de l'industrie des PC (11).

Pour les besoins de cette thèse, on considérera que la structure regroupe les organisations et les règles qui permettent de prendre les décisions sur l'évolution du système. Au contraire, la fonction relève d'une évaluation de l'adéquation de la production du système aux besoins compte tenu des contraintes. Par exemple, la fonction peut être un surplus social (somme des bénéfices des acteurs), tandis que la structure est l'organisation en acteurs du système ainsi que leurs règles d'interactions.

2.1.2.2 Externalités et non-modularité fonctionnelle

Autant il est fréquent en optimisation mathématique qu'un seul acteur possède l'ensemble de l'information sur le système et l'ensemble du pouvoir sur l'évolution du système, autant ceci n'a même pas de sens pour l'organisation des sociétés humaines.

En effet, un grand nombre d'acteurs possède chacun un pouvoir de décision partiel sur le système productif et une information partielle sur son état. L'objectif de l'organisation d'une société (de sa structuration) est de parvenir à ce que les décisions de chaque acteur permettent à l'ensemble du système d'atteindre son optimum fonctionnel, quel qu'il soit.

En économie, les *externalités* sont des exemples de *non-convergence entre structure et fonction*. En effet, de la manière la plus générale, une *externalité* apparaît lorsque l'activité de production ou de consommation d'un agent engendre des coûts ou des bénéfices pour d'autres agents n'ayant aucune influence directe sur l'activité. Un moyen d'éliminer une *externalité* consiste à établir une compensation financière inverse afin que le bilan de l'activité soit neutre pour les agents impliqués malgré eux. Dans le cadre d'un marché où ce type de compensation est fixée par les prix, l'absence d'*externalité* signifie que les acteurs prennent des décisions en étant « informés » par le biais des prix de l'ensemble des conséquences de leurs décisions sur l'économie. L'information dont il s'agit ici est donc bien celle qui pourrait leur faire changer de comportement à travers du prix à payer pour une action donnée (ou à recevoir pour sa réalisation). Par exemple, la tendance actuelle est d'essayer d'intégrer l'objectif global de limiter le réchauffement climatique dans les objectifs de chaque acteur de l'économie. Il y a non-modularité structurelle (chaque voiture et chaque cheminée participe au réchauffement climatique), qui ne possède pas de flux d'information correspondant (prix de l'émission de dioxyde de carbone).

2.1.2.3 Contraintes technologiques et modularité structurelle

A l'inverse des *externalités*, il existe des cas où la progression des technologies a fait apparaître des modularités fonctionnelles. C'est-à-dire que des liens sont devenus suffisamment faibles pour pouvoir organiser de manière modulaire l'appareil de production.

Ainsi, la dissociation technologique en différents composants des fonctions d'un ordinateur qui est à la base de l'invention du PC par IBM a permis la modularisation fonctionnelle de l'organisation de la conception et la production des PC. Ceci s'est manifesté par l'apparition d'un grand nombre d'entreprises intervenant dans leur conception et la fabrication. Les interactions entre ces entreprises ont été rendues possibles grâce à l'existence de standards. Ces standards sont des interfaces technologiques parfaitement spécifiées précisant le comportement attendu de chaque module vis-à-vis du reste du système.

Il existe d'autres cas où les contraintes technologiques sont plus prégnantes. Par exemple, si une automobile est constituée de sous-systèmes modulaires, produits dans des usines différentes, il n'en reste pas moins que le constructeur automobile garde la maîtrise du design de l'ensemble des pièces. Au contraire, la conception même des PC est modulaire : des interfaces ont été spécifiées sous la forme de standards. Dans le premier cas, il s'agit de modularité de la production et dans le second de la modularité du design. En effet, chacun peut faire évoluer une partie du PC (processeur, écran, disque dur, système d'exploitation) indépendamment des autres (11).

Ainsi, il existe en économie un lien complexe entre structure et fonction car une structure adéquate peut créer les conditions de sauts technologiques et les sauts technologiques peuvent en retour fortement modifier la structure d'un secteur.

2.1.3 Evolution du vivant

Les biologistes évolutionnaires s'intéressent à l'évolution des systèmes reproductifs vivants. Bien qu'il existe une unité de la vie sur Terre, on étudie généralement l'évolution d'une population homogène d'organismes vivants, par exemple d'une *espèce*. Cette évolution des systèmes vivants est marquée par un processus d'exploration-sélection dans laquelle l'essentiel de l'information sur la configuration du système est transmise de génération en génération par le biais de l'information génétique généralement portée par de l'*ADN*. A chaque génération, la duplication d'*ADN* n'est pas complètement fidèle, ce qui permet l'exploration. La sélection des individus les plus aptes à se reproduire intervient ensuite. La biologie est certainement le domaine où la question de la différence entre structure et fonction et son corollaire sur la convergence structure-fonction se pose de la manière la plus aiguë aujourd'hui. Cette section est donc dédiée à donner un aperçu des réponses actuelles.

2.1.3.1 Différence entre structure et fonction

La biologie est toute orientée vers l'évolution³. Si Lander (105) partage ce point de vue, il rappelle que la biologie moléculaire échappe à cet aspect téléologique, tandis que les généticiens se limitent volontairement à caractériser la fonction sans référence à un objectif, mais en mesurant des perturbations (sans nécessairement de notion de fonction de fitness de l'organisme) suite à une altération génétique (une variation structurelle). Hormis ces exceptions, la plupart des biologistes se réfère implicitement à la fonction de fitness d'un organisme qu'un processus d'évolution-sélection semble optimiser. Celle-ci comporte deux axes indissociables : la capacité de *reproduction* de l'organisme et son interaction avec l'environnement (78). Ainsi que nous le mentionnions déjà pour le cas général dans la Section (1.3.4), la fonction ne peut pas être définie lorsque l'environnement est lié à l'organisme. Par exemple, les organismes vivants ont produit le dioxygène de l'atmosphère, fortement nocif pour un organisme qui ne sait le séquestrer. Ils ont ainsi dû évoluer face à un environnement qu'ils avaient eux même transformés, ce qui rend impossible la définition d'une seule fonction sur l'ensemble de la période.

En tout état de cause, l'omniprésence du concept d'évolution, la *convergence* souvent forte des modules structurels et fonctionnels et les limites aux expériences praticables créent une difficulté théorique et pratique à bien séparer les deux concepts. Par exemple, l'évolution trouve souvent des solutions analogues à des problèmes analogues : l'œil, par exemple, a été « inventé » de nombreuses fois de manière différente (78). Dans ce cas, il y avait un avantage sélectif à disposer d'un organe (un module qu'on peut enlever) pour capter la lumière. Il est en effet dans la pratique extrêmement délicat de séparer la modularité fonctionnelle, liée au fait que l'impact du fait de voir sur la fonction de fitness est indépendant de celui des autres caractéristiques de l'organisme, de la modularité structurelle liée aux possibilités d'évolution de l'information héritable associée à l'œil indépendamment de l'information associée aux autres caractéristiques. Dans les deux cas, cela peut conduire à la coévolution de caractères, soit qu'ils soient liés fonctionnellement soit qu'ils le soient structurellement pour des raisons historiques (205; 25), mais observer la différence est difficile car les biologistes évolutionnaires sont essentiellement restreints à observer un seul résultat d'un processus évolutionnaire.

Ainsi, Qi et al. (145) propose d'appeler un motif un module structurel, et de réserver le nom de module aux modules fonctionnels. Pour lui, un motif est caractérisé par une structure et une dynamique, mais n'a pas « d'autonomie fonctionnelle ». Cette limite semble peu marquée. Wolf (204) voit plutôt un module comme un groupement de motifs, mais souligne aussitôt la frontière floue entre les deux concepts, qui sont parfois confondus sous le même nom (67). Nous retiendrons qu'un motif est

3. « Une théorie, un système, une hypothèse, l'Evolution ?... Non point : mais, bien plus que cela, une condition générale à laquelle doivent se plier et satisfaire désormais, pour être pensables et vrais, toutes les théories, toutes les hypothèses, tous les systèmes. Une lumière éclairant tous les faits, une courbure que doivent épouser tous les traits : voilà ce qu'est l'Evolution », *Œuvres de Pierre Teilhard de Chardin, tome I : Le phénomène humain* (1956), cité par Dobzhansky (40).

une structure relativement petite à laquelle on peut associer une dynamique déterminée (boucle de rétroaction ou autre), sans préjuger de sa fonction (120; 166; 121; 131; 136).

Schlosser (162) possède lui une vision beaucoup plus claire de cette distinction. Il propose une définition de la modularité basée sur des processus :

« A subprocess of any process qualifies as a module when its components are likely to interact as an integrated and context-insensitive unit to produce a particular IOR [Input-Output Relationship] in the face of various perturbations. »

Il introduit ensuite la notion d'unité d'évolution :

« [I]t will be a unit of those constituents, which tend to coevolve by recurrently constraining each other evolutionary modification. ».

A première vue, il ne semble pas évident de séparer ce qui relève de la structure et de la fonction dans les deux types de modules. En effet, la première définition semble orientée « traitement du signal », sans notion de fonction, tandis que la seconde ne semble pas distinguer les contraintes issues de limitations de l'exploration de configurations de celles provenant de la sélection basée sur la fonction de fitness. Cependant, précisant ce qu'il entend par une relation entrée-sortie particulière indépendante, il donne aussi la définition suivante pour les modules de développement :

« Modules of *development* are units of interacting elements that make a relatively invariant contribution to the development of an organism. »

Enfin, il associe les unités d'évolutions aux « patrons de couplages des constituants d'un système reproductif » et les modules aux « contributions interdépendantes à la fonction de fitness » (161). Les processus modulaires de Schlosser sont donc des modules fonctionnels, du fait même qu'il est extrêmement peu probable d'observer un processus dans un organisme vivant qui n'influence pas sa fonction de fitness. Par opposition, les unités d'évolutions sont des modules structurels, car elles définissent les possibilités d'évolutions (hors sélection) d'un système reproductif.

Schlosser s'intéresse alors aux conditions dans lesquels cette identification entre processus de développement et unité d'évolution peut être faite, autrement dit aux conditions de convergence structure-fonction dans le cadre particulier des systèmes reproductifs. Nous mentionnons ces conditions dans le paragraphe suivant.

2.1.3.2 Conditions de convergence

Schlosser (161) formalise l'existence d'un lien entre module (à laquelle on peut attribuer une fonction) et unité d'évolution (structure) dans les organismes vivants par la satisfaction partielle de ces trois conditions :

- Le système reproducteur peut se passer de l'unité d'évolution pour se reproduire (existence de perturbations non létales de l'unité d'évolution). En effet, dans le cas contraire, il ne sera

pas possible d'observer l'existence de l'unité d'évolution, puisque toutes ses variations seront éliminées par la sélection.

- Il y a convergence entre les effets d'une perturbation d'un module et d'une perturbation de l'unité d'évolution correspondante. Ceci tend à limiter la pléiotropie, c'est-à-dire la réutilisation des mêmes supports d'information héritable dans plusieurs processus. En lien avec la première propriété, ceci signifie aussi que les perturbations du module sont non-létales.
- Les deux propriétés précédentes doivent persister pendant plusieurs itérations de reproduction.

Ces conditions sont une interprétation spécialisée de la *convergence* entre les distributions de probabilité d'évolution et de la fonction de fitness, cherchant à caractériser des ensembles de variables structurelles ayant des impacts indépendants sur la fonction de fitness.

2.1.4 Conclusion

Nous avons montré que la convergence structure-fonction est un phénomène important pour trois domaines d'étude de *systèmes évolutifs et fonctionnels* :

- En optimisation mathématique, la *convergence* est réalisée lorsque la « structure » d'un problème (c'est-à-dire la structure de son objectif, i.e. sa modularité fonctionnelle) est reflétée dans sa représentation (c'est-à-dire ce qui est susceptible de varier indépendamment, i.e. sa modularité structurelle).
- En économie, la *convergence* est réalisée lorsqu'il n'existe pas d'*externalités*. Autrement dit, les acteurs du système, agissant sur sa structure, peuvent alors prendre des décisions indépendamment les uns des autres sur la base d'une information partielle de manière à ce que la somme de leurs décisions mène à l'optimum global. Ceci n'est possible qu'en calquant leurs domaines d'action sur les modularités fonctionnelles (notamment liées aux limitations technologiques) du système de production.
- En biologie, la *convergence* est réalisée lorsque l'information codant pour un processus apporte une contribution indépendante à la fonction de fitness est elle-même indépendante des informations codant pour d'autres processus.

Il reste à étudier dans la dernière section de ce chapitre les explications plausibles pour cette *convergence* en détaillant notamment les gains d'efficacité qui en découlent.

2.2 Apport de la modularité des systèmes évolutifs et fonctionnels

Si la convergence structure-fonction est étudiée dans de nombreux domaines pour des *systèmes évolutifs et fonctionnels* variés, ses apports sont multiples. Pour l'essentiel, elle apporte robustesse et évolutivité, généralement en structurant une variabilité fonctionnelle nécessaire – cf. Section (1.3.3).

- En optimisation mathématique, certains algorithmes sont stochastiques. Ils possèdent une variabilité fonctionnelle qu'il convient de structurer pour un maximum d'efficacité, notamment face à des changements de fonction objectif au cours de l'optimisation. D'autres algorithmes sont déterministes et ne possèdent pas de variabilité fonctionnelle à proprement parler. Néanmoins, la prise en compte de la modularité fonctionnelle de l'objectif leur permet d'être plus efficaces en résolvant de manière relativement découplée des sous-problèmes plus simples.
- En économie, le contexte est changeant et jamais entièrement connu, soit du fait d'actions d'autres acteurs, soit du fait d'évolution des contraintes. Le système doit être suffisamment flexible pour y faire face. Pour cela, il ne peut se permettre d'explorer une seule solution. Une variabilité fonctionnelle est nécessaire afin d'essayer plusieurs solutions.
- En biologie, l'évolution est un processus non borné où de meilleures solutions (adaptations à des niches existantes) sont toujours envisageables et où la découverte ou l'apparition de nouvelles niches sont fréquentes. La partie exploration du processus d'exploration-sélection que constitue l'évolution naturelle est donc essentielle.

2.2.1 Optimisation mathématique

Pour illustrer l'apport de la convergence structure-fonction en optimisation, nous commençons par exposer l'exemple de la décomposition par les prix d'un problème d'optimisation convexe puis nous mentionnons des résultats d'expériences basées sur des algorithmes génétiques.

2.2.1.1 Optimisation convexe et décomposition par les prix

Supposons qu'un problème d'optimisation convexe soit tel qu'il existe une partition des variables du système telle que la fonction objectif s'écrive comme une somme de fonction dépendants chacune uniquement des variables d'un seul module de la partition.

La décomposition de Wolfe-Dantzig (34) vise à utiliser la structure modulaire de ce problème afin de faciliter sa résolution : il s'agit de résoudre les sous-problèmes de chaque module à contexte global fixé (par les variables duales des contraintes couplantes) puis de calculer le nouveau contexte à partir des nouveaux optimums. On peut montrer que, dans le cas convexe, sous certaines hypothèses, notamment de *forte convexité*, cet algorithme converge.

Comme dans tous les algorithmes d'optimisation, il existe une incertitude sur la qualité des itérations possibles à partir d'un état lié au fait que l'algorithme est conçu pour résoudre une classe de problème et non une instance spécifique. Ainsi, il doit à chaque étape explorer une nouvelle configuration qui ne sera pas optimale pour le problème en cours de résolution (atteindre directement l'optimum serait le plus efficace), mais optimale vis-à-vis de la distribution de gain associé à l'en-

semble des problèmes de la classe pour une exploration donnée (espérance de gain maximale, gain positif même dans des cas extrêmes, etc.).

L'utilisation de la modularité fonctionnelle par la décomposition de Wolfe-Dantzig permet cependant de limiter cette incertitude en tenant compte des indépendances entre parties du problème de manière à accroître à la fois la vitesse de convergence pour une large classe de problèmes (robustesse) et l'évolutivité (changement d'objectifs). Nous étudions en détail ces deux aspects dans les paragraphes suivants.

Robustesse - vitesse de convergence : Un découpage modulaire limite le nombre de contraintes couplantes. Dans une certaine mesure, un découpage modulaire apporte une robustesse à la décomposition car elle permet de prendre de bonnes directions à chaque itération, surtout si le problème est régularisé (152). On peut en effet supposer que la convergence est d'autant plus difficile que le nombre de contraintes couplantes est élevé.

De plus, on peut utiliser des algorithmes plus performants adaptés au type de chaque sous-problème, réduisant ainsi la durée d'une itération.

Évolutivité - changement de l'objectif Un avantage d'une telle décomposition réside dans le fait que le détail de la résolution de chaque sous-problème est indifférent à l'algorithme coordinateur qui donne les prix correspondant à un état du système. Ceci introduit une flexibilité plus grande par rapport à une évolution de la fonction objectif : il est possible de faire évoluer la fonction objectif d'un sous-problème sans remettre en cause la fonction objectif des autres.

En revanche, la décomposition ne permet pas de faire face à l'apparition de nouvelles contraintes couplantes ou de couplage des objectifs.

2.2.1.2 Optimisation combinatoire

Robustesse Pour les algorithmes génétiques, une bonne représentation du problème à résoudre permet souvent une plus grande robustesse notamment face à des mutations délétères.

Par exemple, considérons les problèmes de design topologique optimal dans lesquels on souhaite trouver la forme optimale d'une structure (répartition de la matière dans le plan ou dans l'espace) répondant à certaines contraintes. Dans ces problèmes, les lois de la mécanique sont primordiales. A titre d'illustration, si le but est de former un pont entre deux rives, la continuité de la forme est essentielle pour assurer sa fonction de jonction des rives comme pour assurer sa résistance à son propre poids. Dans ce cadre, la représentation la plus évidente, à savoir une matrice basée sur un maillage de l'espace et stockant l'information sur la présence ou à l'absence de matière dans chaque cellule de base du maillage, n'est pas la plus efficace. En effet, cette représentation ne tient pas compte de la continuité de l'espace et de la continuité des solutions à obtenir. Ainsi, une mutation ponctuelle

enlevant la matière permettant de faire « tenir » un barreau pourra créer un trou et sera suffisante pour rendre l'objet entier inutilisable. Au contraire, des représentations plus compactes, basées par exemple sur un *diagramme de Voronoï* (63), permettent d'assurer des transformations liées pour des points voisins dans l'espace.

Ainsi, la connaissance empirique des modularités fonctionnelles (la présence de matière en un endroit donnée n'est efficace que si de la matière est aussi présente dans le voisinage) est exploitée structurellement en assurant que les variations (mutations et recombinaisons) conservent cette continuité du problème initial.

Flexibilité Un certain nombre d'expériences ont été réalisées afin de montrer comment la convergence structure-fonction induit une flexibilité accrue face à des changements d'objectif. En effet, les résultats de Toussaint et al. exposés dans la Section (1.3.3) montrent l'importance de la fonction de fitness effective d'un état. Or, dans le cas où le système peut être soumis arbitrairement à plusieurs fonctions objectif, il se peut qu'il n'existe aucun état capable d'avoir une *fitness* importante avec les deux fonctions, cependant il se peut qu'il existe des états tels que leur *fitness* effective soit grande dans les deux cas (par exemple, s'il existe deux états voisins, l'un étant optimal avec une fonction et l'autre étant optimal avec l'autre).

Avant Toussaint et al., Lipson et al. avaient développé cette idée de promouvoir la modularité en changeant l'environnement ou les contraintes au cours de l'évolution du système (112; 113), même si certaines critiques peuvent être faites sur leurs travaux (54). Kashtan et al. illustrent cet effet de promotion en faisant évoluer des circuits logiques (86). Dans le cas qu'ils présentent, il se trouve que, pour parvenir à ce que les deux états optimaux soient proches, du fait d'une modularité dans le changement de fonction objectif, il faut que l'organisation du système reflète cette modularité pour y parvenir. Lorsqu'elle la reflète, il est possible de passer de l'état optimal pour un objectif à l'état optimal pour le second en seulement deux itérations. Il s'agit d'un exemple de convergence structure-fonction qui favorise la flexibilité du système pour lui permettre de s'adapter rapidement à un objectif changeant.

Ainsi, on peut penser que, s'il existe une modularité commune aux différentes fonctions objectifs rencontrées par un problème d'optimisation lors de sa résolution, les algorithmes globalement les plus efficaces seront ceux qui auront intégré cette modularité dans leur représentation du problème.

2.2.2 Economie

En économie, la variabilité fonctionnelle est extrêmement utile. En effet, elle permet de s'adapter aux changements, notamment technologiques et de limiter les risques d'orienter le système dans une mauvaise direction par l'exploration de plusieurs pistes simultanément. Dans un premier temps, nous verrons comment le *principe de subsidiarité* permet de structurer cette modularité en assurant la

convergence structure-fonction. Ensuite, nous verrons comment la théorie libérale classique pousse ce principe à l'extrême, dans une sorte de transposition à l'économie de la décomposition par les prix utilisée en optimisation mathématique – cf. Section (2.2.1.1).

2.2.2.1 Modularité et principe de subsidiarité

On pourrait aisément éliminer toute *externalité* en considérant un seul acteur qui déciderait de l'évolution de l'ensemble du système. Cependant, pour différentes raisons (information imparfaite sur l'état actuel et futur du système, absence de consensus sur l'objectif, etc.), les décisions de cet acteur ne mènent pas nécessairement à l'optimum. A l'opposé, l'existence d'un grand nombre d'acteurs rend difficile la mise à plat de l'ensemble des *externalités*. Ainsi, il est nécessaire que le système conserve une certaine variabilité fonctionnelle afin d'explorer simultanément un ensemble de solutions, puisque la meilleure, pour autant qu'elle existe, est inatteignable.

C'est ce compromis entre difficulté à éliminer les *externalités* lorsque les acteurs sont nombreux et difficulté à obtenir l'information croissante avec la taille du système qui fonde le *principe de subsidiarité*. Celui-ci postule que la responsabilité des décisions doit être allouée à la plus petite entité capable de résoudre le problème d'elle-même.

Intérêts : Détaillant l'intérêt de cette *subsidiarité* (sans qu'il emploie le nom), Baldwin (11) donne trois raisons à la modularité dans l'industrie (mais on peut généraliser à un système économique quelconque) :

- Rendre la complexité gérable : aucun acteur ne doit connaître l'ensemble du système : chacun peut se contenter du module sur lequel il travaille et de ses interfaces. Ayant une information plus complète sur la configuration de la partie du système qui le concerne, chaque acteur prend de meilleures décisions.
- Permettre le travail en parallèle : chaque module peut être développé séparément.
- S'accommoder de l'incertitude du futur : on peut expérimenter le développement d'un nouveau module (d'une nouvelle implémentation de ces interfaces) sans impact sur le reste du système. Il s'agit d'améliorer l'évolutivité du système. De plus, un module au comportement inadéquat peut être remplacé par un autre. Ceci apporte au système une robustesse à l'encontre des mauvaises décisions qui aurait pu mener à ce module. On trouve un tel exemple de remplacement en informatique : La transition technologique entre les tubes à vide et les transistors n'a pas remis en cause les principes de l'informatique. On a changé le module correspondant au matériel, mais la logique des circuits et de leur utilisation n'a pas varié (67).

Ainsi, la modularité induit une tolérance à l'incertitude et est propices aux expérimentations. Baldwin et al (11) montrent que la modularisation couplée à l'expérimentation réduit les risques liés au développement du système et augmente les chances de profit.

Limites : À l'inverse, la décomposition modulaire a ses limites. Baldwin et al (11) indiquent qu'elle a des coûts, liés à la définition des interfaces et aux tests d'intégration des modules. En effet, plus il y a de couplage entre systèmes, plus il faut définir des interfaces puis tester l'assemblage des modules.

De plus, Marshall et al. (117) et Langlois (106) précisent qu'un système modulaire peut être sous-optimal (par exemple, le just-in-time dans l'automobile augmente l'interdépendance en supprimant les stocks intermédiaires). Le just-in-time est ainsi très délicat à mettre en place car il suppose l'échange d'informations rapidement et de manière exacte. Dans le cas contraire, les états atteints sont très sous-optimaux. C'est pourquoi, on peut préférer un système avec des stocks afin d'apporter de la robustesse.

La dernière limite de la modularité réside dans la fixation des interfaces : si les objectifs évoluent à l'intérieur de chaque module, le système s'adapte très bien. En revanche, si une partie couplante de l'objectif évolue, son changement sera très difficile (160). Il existe un cas connu de tous : la disposition des touches sur les claviers occidentaux diminue la vitesse de frappe. En effet, à l'origine, les marguerites des machines à écrire étaient telles que la frappe simultanée ou immédiatement successive de deux touches voisines les bloquaient. Ainsi, les claviers furent conçus dans chaque langue afin de minimiser le « risque » de frappe successive de touches voisines, ce qui a pour effet de ralentir la cadence de frappe. Cette contrainte a disparu et il existe des dispositions de clavier, par exemple les claviers Dvorak (42) qui permettent de taper la majorité des mots sans quitter la rangée centrale. Cependant, leur adoption est difficile car le coût d'apprentissage d'une nouvelle configuration de clavier est élevé. Du fait de leur rôle d'interface entre l'homme et la machine, leur évolution n'est possible qu'en faisant évoluer simultanément les utilisateurs et la technologie, ce qui s'est révélé impossible.

Langlois (106) précise que le design des interfaces est très difficile car il faut bien connaître ou anticiper le fonctionnement interne des modules et même anticiper la manière dont ils pourront évoluer. Il s'agit bien là d'un problème de *fitness* effective réduite : le système est certes optimal, mais il ne possède pas l'évolutivité nécessaire pour s'adapter à un nouvel objectif.

2.2.2.2 Modularité et théorie libérale classique

Poussant à un extrême l'application du *principe de subsidiarité*, la théorie libérale souhaite donner le maximum de pouvoir aux acteurs individuels d'une société. Ainsi Langlois (106) écrit :

« We can think of Adam Smith's "obvious and simple system of natural liberty" as among the earliest proposals for how a complex modern society might be made more productive through a modular design of social and economic institutions. In separating mine from thine, rights of private property modularize social interaction, which is then mediated through the interface of voluntary exchange, all under the governance of the systems architecture of common law. »

L'instance centrale régulatrice implicite dans cette citation ne collecte pas d'informations : elle se contente de définir le cadre et les règles des échanges d'informations. Ce faisant, elle tente de contraindre les échanges afin que la somme des intérêts particuliers soit équivalente à son idée de l'intérêt général. Ainsi, il est de la responsabilité de chaque acteur de prendre les décisions qui le concerne, lui et sa propriété. De la sorte, la théorie libérale espère promouvoir la robustesse (si les prix sont représentatifs et l'information correcte, les décisions prises par chacun seront optimales) et l'évolutivité (chacun est libre d'essayer une nouvelle solution pour s'adapter à un nouveau contexte).

Il s'agit d'une sorte de décomposition de Wolfe-Dantzig – cf. Section (2.2.1.1) – du système afin d'organiser l'économie d'une société. La théorie libérale se restreint donc à atteindre des optimums définissables comme étant la somme des intérêts particuliers sous certaines règles. Les conditions d'existence de ces optimums sont proches de ceux d'une décomposition par les prix. C'est pourquoi la loi des rendements décroissants est très souvent utilisée dans la théorie libérale et fait partie des conditions de *concurrence* pure et parfaite, car elle assure que le problème d'optimisation correspondant soit convexe. Au contraire, s'il existe plusieurs optimums locaux, la « main invisible » ne peut assurer la convergence vers l'optimum global. Finalement, on note qu'une telle décomposition réduit effectivement la complexité du problème au sens de Simon (1967) pour qui, dans un système complexe, « the whole is more than the sum of the parts ». En effet, par une telle décomposition, l'optimisation de l'ensemble est équivalente à l'optimisation conjointe de chacune des parties.

2.2.2.3 Conclusion

Le *principe de subsidiarité* se trouve bien dans la perspective d'optimiser la *fitness* effective du système économique en ajustant sa structure à sa fonction. En effet, il assure d'une part des prises de décisions efficaces à contexte inchangé et d'autre part des initiatives pertinentes pour s'adapter à un nouveau contexte par un compromis adapté entre difficultés de maîtrise des *externalités* et difficultés de centralisation.

En contrepartie, son application est difficile et généralement subjective. La théorie libérale classique, par la position extrême retenue vis-à-vis de ce principe, a le mérite de la simplicité. Cependant, elle aussi est délicate à appliquer car il est souvent difficile de créer des conditions suffisamment proches des conditions de *concurrence* pure et parfaite nécessaires à la validité de ses conclusions en termes d'optimalité.

2.2.3 Biologie

En biologie, il a souvent semblé contradictoire de sélectionner la robustesse, c'est-à-dire de conserver le niveau d'optimisation en minimisant l'impact de mutations, et l'évolutivité, c'est-à-dire la capacité à explorer de nouvelles formes. En effet, la robustesse semble nécessiter un faible taux de

mutation et l'évolutivité un taux élevé (103; 202). Après avoir exposé ce paradoxe, nous verrons comment la modularité permet de structurer les mutations pour le résoudre. Enfin, comme la question de l'existence d'augmentations de modularité résultant de la sélection soit controversée, nous donnerons quelques exemples pour l'illustrer.

2.2.3.1 Le paradoxe de la robustesse et de l'évolutivité

D'une part, pour les systèmes évolutionnaires vivants, les biologistes se posent la question de la robustesse aux mutations génétiques dans le sens du maintien des fonctionnalités acquises. Il s'agit en quelque sorte d'une vision restreinte de la notion de vitesse de convergence : supposons qu'une espèce ait atteint un optimum par rapport à sa niche, dans quelle mesure est-elle capable d'y rester ? En effet, contrairement à un algorithme qui possède des conditions d'arrêt, l'évolution ne s'arrête pas.

D'autre part, les biologistes s'interrogent sur la capacité des organismes complexes à évoluer, c'est-à-dire à s'adapter à une fonction de fitness changeante. Kirschner et Gerhart (96) définissent l'évolutivité comme la capacité à générer des variations fonctionnelles héréditaires. En effet, ne connaissant pas à l'avance les changements de fonction de fitness, les organismes se doivent d'être capables de générer des variations fonctionnelles. Bonner – cité en exergue par Wagner et Altenberg (195) – se demande même comment les organismes complexes parviennent à évoluer tout court, vu la faible probabilité des mutations bénéfiques.

En première approche, il existe une opposition entre robustesse et évolutivité, car la robustesse diminue avec la variabilité et l'évolutivité augmente (110). D'ailleurs, on sait que les bactéries contrôlent le taux de mutation afin qu'il soit plus élevé pendant les périodes de stress où l'exploration de nouvelles configurations est essentielle (146; 16; 15). Cependant, ce n'est pas un paradoxe de réunir à la fois robustesse et évolutivité car la robustesse est évaluée en moyenne sur l'ensemble des événements tandis que l'évolutivité l'est uniquement sur des événements peu probables. Ainsi, Wagner (193) mesure la robustesse par la probabilité de conserver la même valeur de fonction de fitness tandis qu'il définit l'évolutivité par le nombre de valeurs différentes de la fonction de fitness accessibles à partir de l'état actuel. D'une certaine manière, il suppose deux scénarios différents :

- Supposons que la fonction de fitness soit inchangée depuis suffisamment longtemps. Dans ce cas, l'organisme est largement optimisé pour sa niche. Les mutations qui se produisent seront délétères voire létales avec une très grande probabilité. En première approximation, elles mènent donc à des configurations de *fitness* égale à $-\infty$. L'optimum serait donc qu'elles aient une probabilité nulle de se produire. Pour ne pas diminuer la *fitness*, il faut donc limiter au maximum la probabilité des mutations et la variabilité du système. La robustesse semble donc associée à une faible variabilité.
- Supposons maintenant que la fonction de fitness change et que la valeur de la fonction de fitness associé à l'état courant devienne progressivement égale à $-\infty$ (survie impossible), tandis que

cette valeur passe de $-\infty$ à une valeur réelle pour un des états accessibles à partir de l'état courant. Par exemple, ceci se produit si la niche d'un système évolue progressivement. Dans ce cas, la survie de l'*espèce* passe par la « découverte » de ce nouvel état. Il est alors nécessaire que le système possède une variabilité fonctionnelle élevée afin que la probabilité de découvrir la nouvelle configuration où la survie soit possible. Même si cette variabilité est associée à l'existence de nombreuses mutations délétères, la *fitness* effective sera élevée car le gain associé à la découverte de cette nouvelle configuration est très important.

- A mi-chemin entre robustesse et évolutivité, Carlson et Doyle développent le concept de tolérance hautement optimisée (21; 22; 33; 99; 97; 98). De portée plus large, elle s'applique en biologie pour décrire l'optimisation d'un organisme vis-à-vis d'événements arrivant rarement à l'échelle d'une génération. Il ne s'agit ni d'être dans un état robuste à une seule fonction de fitness ni évolutif face à un vaste ensemble de fonctions de fitness jamais rencontrées. Il s'agit que l'état soit robuste face à un ensemble de fonctions de fitness que le système évolutif rencontre régulièrement avec des probabilités variables. Le poids de chaque fonction de fitness dans la fonction de fitness synthétique résultante est proportionnel à ces probabilités. Ainsi, plus un événement est fréquent, plus on acceptera une dégradation de la fonction de fitness globale pour y faire face.

Ainsi, si l'évolutivité augmente avec la variabilité tandis que la robustesse diminue, il ne s'agit que d'un facteur les opposants. En effet, à variabilité fixée, l'évolutivité n'est pas nécessairement contradictoire avec la robustesse car l'évolutivité repose sur un grand nombre d'états accessibles avec une probabilité qui peut être très faible, tandis que la robustesse est liée à la forte probabilité de stabilité de l'état courant à fonction de fitness inchangée. Une structuration adéquate de cette variabilité, notamment en utilisant les évolutions « neutres » – cf. Section (1.3.3), permet d'assurer les deux simultanément.

2.2.3.2 Apport de la modularité

Dans ce cadre, la modularité joue un rôle important. En effet, par l'optimisation séparée de différents aspects, elle permet :

- d'améliorer l'évolutivité en permettant l'évolution d'un caractère – introduction, modification, voire suppression (73) – sans impact sur les autres. Au contraire, s'il n'y a pas de nécessité de conserver une variabilité fonctionnelle, par exemple en présence de sélection stabilisante (i.e. dans un environnement non changeant), Wagner (195) note que la modularité n'apparaît pas. De même, il ne suffit pas que l'environnement change afin d'introduire la modularité : il faut que les changements soient modulaires sur le plan fonctionnel. Ainsi, Hintze (71) n'observe pas de modularité dans les réseaux métaboliques évolués en l'absence de certains précurseurs à tour de rôle. En effet, dans ce cas, la meilleure stratégie consiste à créer des liens entre

les chaînes de réactions utilisant chaque précurseur afin de pouvoir remplacer un précurseur absent. Ainsi, le système le plus efficace n'est pas nécessairement le plus modulaire, même dans un environnement changeant ;

- d'améliorer la robustesse en limitant l'impact des mutations délétères à un seul module (176; 204; 5; 173; 53; 144). Ainsi, l'expérience a été réalisée de supprimer un par un les gènes de la levure (82). Seule une faible proportion des gènes sont indispensables à la survie dans un environnement de laboratoire. Par ailleurs, on peut tracer un graphe où deux gènes sont reliés par un lien si les *protéines* qu'ils produisent interagissent. On remarque alors que les gènes les plus connectés sont beaucoup plus souvent indispensables. Or, conformément à la théorie, le réseau est plus rapidement déconnecté par des attaques sur ces nœuds que sur les autres. On peut donc considérer que la conservation de la connexité du réseau est une approximation utile bien que grossière de la fonction de fitness d'un organisme vivant (6) qu'il est possible d'étudier en utilisant un algorithme d'attaque ciblée – cf. Section (1.2.3.3). Il ne faut pas pour autant considérer que la centralité cause l'essentialité. En effet, l'existence de variables cachées causant à la fois la centralité et l'essentialité est probable (211). Toutefois, on peut voir cette fragilité de l'attaque des gènes les plus connectés comme le pendant de l'impact relativement faible de l'attaque des autres nœuds. En effet, il est nécessaire d'assurer la cohérence de l'ensemble des modules. Minimiser le nombre de nœuds cruciaux par les liens qu'ils assurent permet de réduire le risque qu'une mutation délétère les touche.

En conclusion, Lipson (111) pense même que la capacité à représenter les configurations de manière modulaire, régulière et hiérarchique est cruciale pour les systèmes évolutionnaires non bornés. Sans cette capacité, ils ne pourraient pas « passer à l'échelle », c'est-à-dire rester évolutif et robuste malgré leur taille.

Enfin, nous reproduisons dans les pages suivantes une revue du livre sur la modularité dans le développement et l'évolution édité par Schlosser et Wagner (163). En effet, on retrouve dans l'ensemble de l'ouvrage cette idée que la modularité façonne l'évolutivité des organismes.

<http://dx.doi.org/10.1002/bies.20243>

2.2.3.3 Exemples

Parmi les sciences des systèmes évolutifs, la biologie évolutionnaire présente une difficulté supplémentaire. En effet, il est impossible d'observer l'évolution naturelle, mais seulement son résultat (sauf cas particulier d'organismes se multipliant rapidement). Par delà les simples constats généralisés de coévolution – cf. Section (2.1.3.1), voici donc des preuves parmi les plus saillantes de l'impact positif d'une modularité structurelle convergente avec la modularité fonctionnelle sur l'évolutivité.

Blocage et modularité insuffisante Il existe un exemple de blocage de l'évolution pour cause de non convergence structure-fonction : celui de la transition à la multicellularité des *volvocacées* (algues vertes) mentionné par Nedelcu et al. (126).

Un organisme multicellulaire doit contrôler la prolifération des cellules qui le composent. En particulier, il existe certaines cellules qui participent à la *reproduction*, dites cellules germinales. Ces cellules doivent pouvoir se diviser pour permettre la *reproduction* de l'organisme. Les autres cellules sont dites somatiques. Initialement, elles doivent se diviser pour former l'organisme. Ultérieurement, celui-ci peut subsister sans qu'elles ne se divisent, mais ce n'est pas optimal car l'organisme peut subir des atteintes telles que certaines cellules meurent. Une division des cellules somatiques est alors nécessaire pour les remplacer. Au contraire, une prolifération anarchique s'apparente à un cancer. C'est là un risque crucial pour les organismes multicellulaires qui doivent donc absolument contrôler la division des cellules. Pour bloquer tout risque de prolifération anarchique, les algues vertes ont évolué vers une solution où la croissance en volume même des cellules somatiques est stoppée, et, par conséquent, leur division car celle-ci est conditionnée par une augmentation initiale de leur volume (comme chez la grande majorité des organismes). Du fait de cette maîtrise de la prolifération, il existe des algues vertes multicellulaires. La multicellularité leur procure des avantages sélectifs (mobilité en particulier), mais la solution trouvée pour maîtriser la prolifération a limité le potentiel d'évolution. En effet, cela serait un avantage que la croissance en volume et la division ne se limitent pas à la phase de développement et qu'il existe des cellules souches somatiques capables de reconstituer ou renouveler un tissu. Entre autre, la durée de vie des algues vertes multicellulaires est limitée à 5 jours du fait de la mort programmée de ses cellules somatiques. Les avantages liés à ce découplage sont inatteignables par les algues vertes. Au contraire, chez les organismes multicellulaires évolués (champignons, plantes et *animaux*), il existe des cellules somatiques capables de se diviser, ce qui est probablement une des raisons expliquant l'explosion de complexité dans ces groupes.

La fonction de fitness associée à la maîtrise jointe de la multicellularité et à la capacité de croissance des cellules somatiques est séparable, puisqu'il existe des organismes multicellulaires possédant des cellules souches somatiques. Ces organismes possèdent deux unités d'évolution afin que la sélection puisse opérer relativement indépendamment sur les deux aspects de la fonction de fitness que sont la formation d'un tissu multicellulaire et son renouvellement. Cependant, dans le contexte évolutif accessible aux algues vertes, il n'existe pas de telle solution séparée. En effet, elles ne possèdent qu'une seule unité d'évolution portant sur la multicellularité mais interdisant la croissance en volume des cellules somatiques. Autrement dit, elles possèdent un seul module structurel (gène *regA* responsable du blocage de la croissance des cellules) pour deux modules fonctionnels (assurer une croissance multicellulaire non-anarchique d'une part et renouveler les tissus d'autre part). Il faudrait d'abord qu'elles perdent la multicellularité pour l'acquérir à nouveau sous une forme moins contraignante pour l'évolution future.

Taux de diversification et modularité En taxinomie, les *espèces* sont regroupées en *taxons* d'après leurs caractéristiques. Ces *taxons* sont dénommés suivant leur rang. Les *insectes* forment un *taxon* du rang de classe. Il existe un *taxon* inclus dans les *insectes* dont une caractéristique commune est le fait que la larve est radicalement différente de l'adulte. Les individus des *espèces* de ce *taxon* du rang de sous-classe connaissent une métamorphose complète. C'est pourquoi on dit qu'ils sont *holométaboliques*. Par opposition, cette métamorphose est dite incomplète chez les autres *insectes*. La plupart d'entre eux sont dit *hémimétaboliques* car le mode de vie de la larve est comparable à celui de l'adulte et le passage de l'un à l'autre est progressif. On pense que les *insectes holométaboliques* sont un *taxon* monophylétique ou *clade*, c'est-à-dire qu'il regroupe l'ensemble des *espèces* descendant d'un même ancêtre commun. C'est aussi le cas des *insectes*, mais pas du groupe constitué des *insectes* privé de ceux qui sont *holométaboliques*.

L'étude des fossiles permet de tracer le nombre d'*espèces* apparues dans un *taxon*. Un grand nombre d'entre elles a disparu par la suite. En considérant que les *espèces* sont une « solution » de qualité suffisante au « problème » d'optimisation de la fonction de fitness, on peut supposer que le *taxon* le plus évolutif est celui qui a vu apparaître le plus grand nombre d'*espèces* par unité de temps. Ceci définit le taux d'apparition de nouvelles *espèces*. Afin de le pondérer par la taille du *taxon*, on définit le taux de diversification d'un *taxon* comme le taux d'apparition de nouvelles *espèces* rapporté à sa taille.

Yang (207) montre que les *insectes holométaboliques* ont un taux de diversification supérieur aux autres *insectes*. De manière plus précise, en moyenne sur une longue période de plus de 250 millions d'années, proportionnellement à leur taille respective, plus de nouveaux *taxons* du rang de familles apparaissent chez les *holométaboles* que chez les *hémimétaboles*. On peut supposer que cela est lié au fait que la larve et l'adulte peuvent s'adapter chacun indépendamment à des niches différentes, comme dans le cas de l'oursin (191). Ainsi, il existe une modularité fonctionnelle (l'environnement de la larve peut être différent de celui de l'adulte et chacun peut changer indépendamment) exploitée structurellement par les *holométaboles*, d'où leur plus grande évolutivité. Au contraire, les autres *insectes* sont incapables d'adapter les formes de la larve et l'adulte à deux environnements radicalement différents, par exemple eau et air, ce qui limite le nombre de niches environnementales potentielles dans lesquelles une nouvelle *espèce* pourrait apparaître.

Chapitre 3

Objectifs de la thèse

Un des moyens classiques d'aborder l'unité des systèmes a été de les représenter sous forme de graphe et de montrer les points communs dans leurs caractéristiques topologiques (134). Du point de vue fonctionnel, on a ainsi montré que la plupart des réseaux réels sont peu sensibles à des attaques sur des nœuds quelconques mais que les attaques sur des nœuds de *degré* élevé les déconnectaient rapidement. Ceci suggère l'importance de la modularité dans les réseaux réels car ils comptent généralement peu de nœuds assurant la cohérence globale de l'ensemble, appelés « hubs », et beaucoup de nœuds formant des modules plus ou moins nets appelés « clusters » et faiblement connectés entre eux par les « hubs » (200). De plus, des mécanismes d'évolution basés sur l'attachement préférentiel permettent de rendre compte de cette structure (3; 127; 39). Dans des domaines variés, il est donc raisonnable de penser que la structure « hubs-clusters » résultant d'une évolution corresponde en partie à la fonction du système. En effet, celle-ci repose généralement sur la conservation de la connexité du réseau (153). Au simple niveau topologique, il y a donc convergence structure-fonction.

Pour généraliser cette approche, constatant ses limites (92) notamment pour les réseaux électriques (154) qui sont mal décrits par une structure « hub-cluster », nous avons étudié et mis en lumière le rôle central de la convergence structure-fonction pour la robustesse et l'évolutivité des *systèmes évolutifs et fonctionnels* – cf. Section (2.2). Afin d'utiliser ce concept dans les deux grands domaines d'existence de *systèmes évolutifs et fonctionnels*, à savoir les systèmes vivants d'une part et les systèmes créés par les hommes de l'autre, nous avons choisi d'apporter des éléments de réponses à deux problématiques ouvertes :

- Le rôle de la *recombinaison* de l'*ADN* dans la convergence structure-fonction des systèmes vivants reproductifs pour lesquels l'évolution est essentiellement combinatoire. La Section (3.1) expose cette problématique à laquelle est consacrée la Partie (II) de cette thèse.
- La nécessaire coordination des producteurs et des transporteurs d'électricité face à la non-modularité fonctionnelle des systèmes électriques, systèmes technico-économiques dont la mo-

délisation relève essentiellement de l'optimisation non-linéaire. La Section (3.2) expose cette problématique à laquelle est consacrée la Partie (III) de cette thèse.

Ceci permet de vérifier d'un part que les concepts développés dans l'introduction ne sont ni trop abstraits pour être utiles ni trop spécifiques pour apporter une vision d'ensemble, d'autre part qu'ils forment un cadre conceptuel intéressant pour analyser la complexité des *systèmes évolutifs et fonctionnels* et qu'ils sont complémentaires des approches existantes basées sur la topologie de réseaux.

3.1 Recombinaison de l'ADN

Les mutations des génomes peuvent être classées dans différentes catégories, qui se complètent les unes les autres (125). Les plus simples sont les mutations ponctuelles qui ne modifient qu'une seule position dans l'ADN. D'autres insèrent ou éliminent une seule position. Si de telles mutations sont suffisantes pour découvrir de nouvelles fonctions (164), les mutations de plus grande ampleur raccourcissent considérablement le chemin d'une configuration à une autre (46).

Par exemple, certaines insèrent ou éliminent une portion complète d'ADN. De telles insertions peuvent conduire à la duplication de gènes, souvent suivie par une différenciation de leurs fonctions (102). On peut penser que de tels phénomènes de duplication jouent un rôle essentiel dans la formation de structure « hub-cluster » (13). En effet, on sait que de telles structures peuvent émerger si, lors de l'évolution du système, de nouveaux liens sont créés préférentiellement sur les nœuds qui sont déjà les plus connectés. Même si la probabilité de duplication est la même pour l'ensemble des nœuds est identique, les nœuds qui sont les plus connectés ont plus de chance d'avoir un voisin qui se duplique, d'où le fait que les duplications donnent l'apparence d'une croissance avec attachement préférentiel. C'est ainsi que, même en tenant compte des pertes de liens redondants successifs aux duplications (74; 47; 189), on montre que la duplication des gènes joue un rôle important dans la modularité structurelle du graphe des interactions entre gènes, sans que l'existence d'une modularité fonctionnelle correspondante ne soit évidente (169; 171; 197).

De même, les *recombinaisons* génétiques que nous qualifierons d'homologues (aboutissant au remplacement d'une portion de l'ADN par une autre comparable, généralement par échange entre deux ADN semblables, par exemple par *enjambement* lors de la *méiose*) constituent un opérateur de mutation qui façonne fortement la modularité structurelle de l'information génétique. En effet, du fait des *recombinaisons*, plus deux positions de l'ADN sont proches, plus elles ont de chance d'être transmises ensemble. Nous montrerons que cette propriété est effectivement exploitée pour structurer le génome (par exemple en éloignant les gènes les uns des autres pour favoriser les *recombinaisons* entre gènes). La question se pose alors de savoir si cette modularité structurelle est bien calquée sur une modularité fonctionnelle, comme le laisse entendre l'introduction de cette thèse.

Après avoir répondu à cette question générale, nous nous poserons deux questions particulières :

- Chez les bactéries, existe-t-il des contraintes structurelles, notamment de taille des unités d'information génétique qui interdisent la formation de modules de trop grande taille ? Autrement dit, y-a-t-il une limite à la non-modularité de l'information génétique ? Ceci impliquerait un fort biais de l'évolution vers les structures modulaires, indépendamment de la faiblesse de la modularité fonctionnelle correspondante.
- Chez l'homme, est-il possible d'utiliser une éventuelle convergence structure-fonction pour analyser de manière fonctionnelle des études d'association génétique, alors que les méthodes classiques donnent des résultats structurels ? Ceci permettrait de faciliter leur analyse et leur croisement avec des données fonctionnelles que les derniers développements technologiques ont permis d'accumuler en quantité et ainsi d'atteindre l'objectif principal de telles études : la meilleure compréhension des mécanismes à l'origine de maladies.

3.2 Distinction producteur-transporteur dans les systèmes électriques

Si la *recombinaison* permet d'illustrer la convergence structure-fonction dans les systèmes vivants, l'étude de la restructuration du secteur de l'énergie électrique en Europe montre la pertinence de ce concept dans une autre grande classe de *systèmes évolutifs et fonctionnels*, à savoir ceux créés par l'homme.

Cette restructuration consiste en une redéfinition des acteurs et de leurs relations. Cette introduction a illustré le fait que cette modification de la structure n'est bénéfique que si elle ajuste encore mieux la structure à la fonction – cf. Section (2.2.2). Nous nous interrogerons sur l'existence de complexité et de non-modularités fonctionnelles dans les systèmes électriques afin de mettre en évidence les difficultés liées à cette restructuration. Au passage, nous verrons en quoi cette complexité échappe en grande partie à l'analyse topologique des graphes correspondant aux réseaux de transport électrique.

A partir de là, nous nous poserons la question suivante : comment ajuster la relation entre les gestionnaires de réseaux de transport et les autres acteurs du réseau, en particulier les producteurs, afin que la dissociation de la production et du transport ne mène pas à une sous-optimisation des investissements ?

Deuxième partie

Optimisation de la convergence structure-fonction des systèmes vivants par la recombinaison de l'ADN

Chapitre 4

Introduction

4.1 Le gène : une entité structurale et fonctionnelle

Pour les organismes vivants, la principale formalisation de la *convergence* entre structure et fonction de fitness réside dans le concept de gène. En effet, un gène est à la fois :

- une séquence d'*ADN* ;
- l'ensemble des macromolécules (*ARN* non messagers et *protéines*) qu'il permet de produire ;
- la caractérisation de sa régulation au sens large. Par exemple, pour un gène *codant* une *protéine*, il peut s'agir des conditions de *transcription* en *ARN* messenger, d'*épissage* ou de dégradation de l'*ARN* messenger.

Ainsi, le gène est à la fois structurel, car la donnée de sa séquence d'*ADN* suffit à le caractériser, et fonctionnel, car il est aussi caractérisé par le rôle de l'information qu'il contient dans la cellule, c'est-à-dire, plus ou moins indirectement, son impact sur la fonction de fitness de l'organisme qui le porte. Comme le montre ce double sens structurel et fonctionnel, la *convergence* est admise comme une évidence. Cependant l'étude de son origine n'est pas sans intérêt. Elle réside fondamentalement dans la robustesse et évolutivité permise par le fait de calquer la modularité structurelle sur la modularité fonctionnelle, ainsi que l'explique la Section (2.2.3.2).

Afin d'illustrer de manière détaillée ce propos, et dans le but de poser les concepts utiles à la compréhension des articles de cette partie, nous explorons maintenant les implications de cette *convergence* par rapport aux deux grands modes de variations de l'information héritable contenue dans l'*ADN*, à savoir les mutations ponctuelles et les *recombinaisons*.

4.2 Mutation et gènes superposés

On tient en général pour acquis le fait que l'information génétique des gènes ne se superpose pas. En effet, dans la perspective de la convergence structure-fonction, il est évident que l'évolutivité est d'autant plus grande que les variations structurelles ponctuelles n'ont d'impact que sur un module fonctionnel, c'est-à-dire un gène (au sens des *protéines* qu'il permet de produire).

Toutefois, il existe des exceptions, particulièrement mais pas exclusivement, dans les organismes où il existe une pression de sélection pour diminuer la taille du génome :

- Les *procaryotes*. En effet, contrairement aux *eucaryotes*, ils ne possèdent qu'une origine de *réplication*. Même si une *réplication* peut être initiée avant que la précédente ne soit achevée, la vitesse de division de la cellule diminue avec la taille du génome (138). En effet, lors de la *réplication*, la machinerie cellulaire lit l'*ADN* à une vitesse dépendante de l'organisme. Afin d'accélérer la vitesse de division, l'organisme doit donc accélérer la vitesse de *réplication* (mesurée en bases par seconde) et/ou raccourcir la longueur du génome (mesurée en bases). Ceci est particulièrement sensible dans un environnement riche en nutriments et en espace où un écart de 1% dans la vitesse de division de deux organismes suffit à marginaliser le plus lent en moins de 1000 générations (soit 10 jours avec une division tous les quarts d'heures). Il y a alors moins d'un individu du type le moins rapide pour 1000 de l'autre.
- Les virus. En effet, leur matériel génétique doit loger dans une enveloppe de taille limitée.
- Les mitochondries. En effet, ces *organites* responsables de la production d'énergie dans les cellules *eucaryotes* ont une origine *procaryote*.
- Les *eucaryotes*, même si cela reste un phénomène rare dont la fonction est encore mal connue (116) car la pression de sélection de la taille des génomes *eucaryotes* est faible.

La superposition crée des contraintes qui se traduisent par des possibilités réduites de variations de la portion concernée (123). En effet, toute mutation dans la région de superposition peut avoir des conséquences sur la fonction de chacun des gènes superposés. Même si elle est favorable pour un gène, elle pourra être fortement défavorable pour l'autre et ne sera donc pas sélectionnée. Les mutations ne peuvent donc être sélectionnées que si leur effet global est positif. Toutefois, Krakauer (102) montre aussi que la superposition peut être stable et même bénéfique car la superposition crée un lien entre les conditions de *transcription* des deux gènes, à la manière des gènes inclus dans le même *opéron* dont nous discutons dans la Section (4.3). La transcription simultanée et dans les mêmes proportions peut être effectivement être un avantage évolutif dans différentes situations dès lors qu'il est nécessaire d'assurer que les protéines codées par les deux gènes soient produites dans une proportion bien définie (gènes codant des enzymes pour des réactions qui s'enchaînent, des protéines qui s'assemblent deux par deux pour former un complexe, etc.).

4.3 Recombinaison

Parallèlement aux mutations ponctuelles, les *recombinaisons* sont un processus majeur de l'évolution de l'*ADN* des systèmes reproductifs. Celles-ci ont lieu à l'occasion de la coupure des deux brins d'une séquence d'*ADN*. Cette coupure n'a généralement pas lieu au même endroit sur chaque brin. Il apparaît alors deux extrémités assez courtes (une dizaine de bases) simple brin. Il existe alors des *enzymes* spécialisées dans la recherche d'extrémités d'*ADN* simple brin dont l'une est le complémentaire de l'autre susceptibles d'être connectées. Il s'agit d'une heuristique qui permet en général de reconnecter les extrémités qui viennent d'être séparées. Cependant, il arrive que la connexion ne répare pas la cassure en trouvant une extrémité simple brin de séquence complémentaire qui n'est pas l'extrémité initiale. Il peut s'agir :

- D'une autre portion de l'*ADN* de l'organisme. S'il y a eu deux cassures, cela peut aboutir à la suppression pure et simple de la séquence d'*ADN* située entre les deux cassures. Chez un organisme *diploïde*, comme l'homme, l'*enjambement* des chromosomes lors de la *méiose* entre dans cette catégorie. Comme l'enjambement ne fait qu'échanger des séquences similaires entre chromosomes, ses conséquences sont généralement moins dramatiques que les *translocations* de séquences d'*ADN* entre chromosomes non homologues ou les suppressions.
- D'une séquence d'*ADN* provenant d'une autre *espèce*. Ce cas a surtout été étudié chez *procaryotes*. Il s'agit du phénomène de *transfert latéral ou horizontal* de gènes, par opposition au transfert vertical du matériel génétique à la descendance. La probabilité de réussite d'une telle intégration est rare à l'échelle d'un seul individu, mais importante au regard de l'évolution, comme en témoignent les nombreuses traces de tels transferts.
- D'une séquence d'*ADN* provenant d'un autre organisme au génome similaire. Ce cas peut se produire chez les *procaryotes*. L'effet est alors proche de celui de la *reproduction sexuée* : le génome de l'organisme est alors une mosaïque entre le génome initial et celui de la bactérie dont provient la séquence intégrée.

Les trois sections suivantes explorent successivement ces cas. A proprement parler, toutes ces *recombinaisons* peuvent être qualifiées d'homologues, car le mécanisme reconnecte uniquement les extrémités qui sont complémentaires l'une de l'autre. Par abus de langage, dans cette thèse, on appellera *recombinaison* homologue uniquement les *recombinaisons* aboutissant au remplacement d'une séquence par une autre comparable. Les autres *recombinaisons* seront qualifiées d'hétérologues, dans le sens où l'*ADN* formé à la suite de la reconnexion n'est pas homologue à celui qui existait avant les coupures. Ainsi, l'*enjambement* des chromosomes humains sera qualifié de *recombinaison* homologue. En revanche, le déplacement, l'insertion ou la suppression d'information génétique ne le seront pas, même à l'intérieur du même organisme.

4.3.1 Recombinaison homologue

Du point de vue de l'évolutivité, la *recombinaison* homologue apporte un avantage certain : pourvu qu'elle respecte la division en gènes de la séquence d'*ADN*, elle permet d'échanger la configuration d'un module par une autre, quelle qu'elle soit, sans avoir à passer par le long chemin des mutations point à point qui mène de l'un à l'autre (46). Ainsi, si on suppose qu'il existe dans une population une diversité de configurations pour plusieurs modules, chacune optimale dans des conditions particulières, l'évolution explorera plus rapidement la combinaison de ces configurations localement optimales, réduisant ainsi la complexité de l'optimisation, telle que définie pour l'optimisation combinatoire dans la Section (2.1.1.1).

On s'attend donc à ce que la sélection pour l'évolutivité ait conduit à moduler la probabilité de *recombinaison* homologue afin de superposer la modularité structurelle que les *recombinaisons* créent à la modularité fonctionnelle associée, c'est-à-dire en premier lieu celle des gènes. En particulier, on peut penser que les *recombinaisons* ont lieu préférentiellement en dehors des gènes afin de préserver leur cohérence (non-modularité) interne. Chez l'homme, il existe des observations de deux types :

- La première consiste à comparer l'*ADN* d'individus puis à évaluer ensuite le taux de *recombinaison* tout au long de l'*ADN* à partir de l'identification de séquences entières identiques chez un grand nombre d'individus (donc jamais séparées par une *recombinaison* dans l'histoire des populations humaines). C'est une étude de *déséquilibre de liaison* ;
- La seconde consiste à observer directement les *recombinaisons* homologues lors de la *méiose*. Dans le premier cas, on observe un résultat après sélection. On ne peut donc pas déterminer si un taux de *recombinaison* localement faible correspond :
 - au fait qu'il y ait effectivement moins de *recombinaisons* dans la région concernée ;
 - au fait qu'il y ait autant de *recombinaisons* mais que celles-ci conduisent à une diminution de la fonction de fitness de l'organisme possédant le matériel recombiné telle que les traces de ces *recombinaisons* ne sont plus visibles dans les populations contemporaines.

Seule l'observation directe d'une modulation du taux de *recombinaison* au long de l'*ADN* permet de déduire qu'il y a eu sélection sur la fonction de fitness effective afin de calquer la modularité structurelle (indépendance relative des variations de séquences d'*ADN* dues à la probabilité élevée d'une *recombinaison* entre les deux séquences) sur la modularité fonctionnelle (indépendance fonctionnelle relative des gènes). Par exemple, on peut penser que les *recombinaisons* situées à l'intérieur des gènes sont plus souvent délétères (impact négatif sur la fonction de fitness) car un *ADN* possède une cohérence interne élevée. Ainsi, il se peut qu'un gène ne soit fonctionnel que si deux bases relativement éloignées l'une de l'autre dans sa *séquence codante* sont identiques. Toute *recombinaison* entre deux variantes du gène toutes deux fonctionnelles car portant des bases identiques aux deux positions mais différentes d'une variante à l'autre produirait des variantes non fonctionnelles. Au contraire, on peut

penser que ce cas est plus rares entre des bases de deux gènes différents, d'où la possibilité de les séparer par des *recombinaisons*.

Les études directes montrent que, chez de nombreuses *espèces*, les *recombinaisons* homologues ne se produisent que dans des régions très étroites de l'*ADN*, appelées *points chauds de recombinaison*, séparées par de larges régions sans *recombinaison* (56; 90). Il est possible que cette distribution presque discrète soit liée à l'augmentation de la fonction de fitness effective qui peut apparaître lorsque les séquences entre deux *points chauds* ne sont jamais recombinaisonnées. Qu'un tel intervalle contienne un gène ou non, la question de l'existence d'une modularité fonctionnelle associée se pose donc. Toutefois, la réponse n'est pas évidente car il semble que, du singe ou de la souris à l'homme, la localisation des *points chauds* ne soit pas conservée, ce qui semble indiquer que l'effet de la localisation des *points chauds* sur la fonction de fitness effective n'est pas essentiel et évolue rapidement (31).

D'un autre côté, les études génétiques montrent aussi une grande variation du *déséquilibre de liaison* (124) et un *déséquilibre de liaison* fort est en général le signe d'un taux de *recombinaison* faible. Les observations directes et indirectes concluent toutes deux à l'existence de ces *points chauds* (80; 81). De plus on constate une augmentation du *déséquilibre de liaison* (ainsi qu'une diminution du polymorphisme) dans les *régions codantes (exons)* de l'*ADN* (180; 70; 88). A une échelle moyenne (de l'ordre de la longueur d'un gène), on observe au contraire que les gènes se trouvent dans des régions de faible *déséquilibre de liaison*, ce qui est cohérent avec le fait que le *déséquilibre de liaison* doit couvrir l'information correspondant à un seul module fonctionnel, soit en général un seul gène. Ce phénomène est anticipé depuis longtemps puisque Hill et Robertson ont montré comment la sélection effective était réduite pour des variations structurelles fortement liées (trop proches pour être séparées par une *recombinaison*) (69). Certains (187; 138) avancent même l'hypothèse que l'existence d'une quantité importante d'*ADN non-codant* dans les gènes, les *introns*, permet de créer des biais encore plus importants dans les probabilités de *recombinaison* afin de calquer au mieux la modularité structurelle due à la *recombinaison* sur la modularité fonctionnelle. Dans le champ de la programmation génétique (141) (l'art de construire des programmes en faisant évoluer une population de programmes par mutation et recombinaison), une observation similaire a déjà été faite : les codes tendent à devenir très longs (phénomène dit de « bloat », littéralement de bouffissure) et une grande partie n'est jamais exécutée (équivalent des séquences d'*ADN* qui ne sont pas codantes et ne sont pas utiles à l'initiation de la transcription). On n'a cependant jamais prouvé que cela fût lié à l'amélioration de l'efficacité des *recombinaisons*.

A une plus large échelle, sans que cela soit contradictoire avec l'observation précédente, on constate même qu'il existe plus de gènes dans les régions à haut et bas *déséquilibre de liaison* que dans les régions de *déséquilibre* moyen (168). Ceci indique que, pour certaines fonctions, un haut niveau de brassage des gènes associés est bénéfique (système immunitaire, etc.), tandis que c'est le contraire pour d'autres (mécanismes de gestion de l'*ADN* et de l'*ARN*, etc.). Ce résultat est confirmé

par l'analyse des taux de *recombinaison* de Kato et al. (87). Ceux-ci avancent l'hypothèse que la sélection naturelle réduit la diversité génétique apparente en sélectionnant fortement les *recombinaisons* dans certaines régions.

Enfin, les études des séquences d'*ADN* montrent l'existence de séquences spécifiques associées aux *points chauds de recombinaison* (168).

En guise de conclusion, nous faisons donc les deux remarques suivantes :

- La structure génique et celle des *points chauds* se superposent partiellement, indiquant clairement que les séquences des gènes sont des modules structurels correspondant à une modularité fonctionnelle.
- L'information concernant la localisation des *points chauds* est partiellement *héritable* car en partie déterminée par les caractéristiques de l'*ADN*. Il est donc possible que l'évolutivité soit ainsi soumise à sélection.

4.3.2 Recombinaison hétérologue inter-espèce

De même, la *recombinaison* hétérologue inter-espèce impliquant l'intégration de matériel étranger à l'*espèce* permet elle aussi de calquer la modularité structurelle sur la modularité fonctionnelle. En effet, un organisme peut acquérir en une seule étape une nouvelle fonctionnalité sans avoir à la faire évoluer *ab initio* (148). Ce phénomène appelé *transfert horizontal* a surtout été étudié chez les *procaryotes* (18), mais il commence à être mieux connu chez les *eucaryotes* (91). Woese pense même qu'il a eu un rôle plus important que la transmission de l'information par héritage (transfert dit « vertical ») aux origines de la vie (203).

Même si un gène possède une fonction, elle n'est en général pas suffisamment indépendante de celle des autres gènes pour apporter à elle toute seule une contribution indépendante à la fonction de fitness de l'organisme qui le porte. En général, plusieurs *protéines* et donc le plus souvent plusieurs gènes sont nécessaires pour apporter une telle contribution indépendante. En accord avec la *convergence* entre structure et fonction de fitness, le transfert d'une fonctionnalité entre deux *espèces* d'organismes a d'autant plus de chances d'être fructueux que l'ensemble des gènes nécessaires à l'accomplissement de la fonction est transmis simultanément, c'est-à-dire dans la même séquence d'*ADN*. Ce mécanisme, dit « des *opérons* égoïstes » (108), est supposé être à l'origine de l'existence d'*opérons* chez les bactéries, c'est-à-dire d'enchaînements des séquences des gènes assurant ensemble une fonction sur l'*ADN* d'une bactérie. Au contraire, il n'existe pas de tels *opérons* chez les *eucaryotes*, ce qui ne facilite pas ces transferts. Aujourd'hui, il existe un grand nombre d'*opérons* bactériens dont on sait qu'ils ont été intégrés de la sorte (18). Ainsi, la vitesse de l'apparition des résistances aux antibiotiques est liée à l'existence même de ces *transferts horizontaux de gènes* (133).

4.3.3 Recombinaison hétérologue intra-espèce

Enfin, Wagner (194) fait l'hypothèse que ces *recombinaisons* hétérologues intra-espèce participent d'une forme de gestion du risque par les populations bactériennes. En effet, une fonctionnalité peut être nécessaire uniquement pour survivre à un événement rare. Dans ce cas, porter l'information génétique est un handicap la plupart du temps. Grâce au *transfert latéral* de gènes, il est envisageable que seule une partie de la population bactérienne conserve l'information génétique d'une situation rare. En quelque sorte, il y a mutualisation des coûts de stockage et de reproduction de l'information génétique tout en assurant qu'une partie de la population est capable de survivre à un événement rare.

Il existe une forme analogue de gestion collective du risque chez l'homme : certaines maladies possédant des facteurs génétiques existeraient précisément du fait que ces mêmes facteurs de prédisposition procurent d'autres avantages à l'individu ou à la population dans son ensemble. Par exemple, les schizophrènes semblent avoir une résistance accrue aux blessures et infections, ce qui a pu avoir une grande importance dans le passé (38). De même, il existe une prédisposition au diabète qui se développe dans les conditions d'une alimentation riche en énergie et d'un style de vie sédentaire. Cette prédisposition est héritée du passé où l'alimentation était irrégulière et les infections nombreuses (99) et serait probablement un avantage dans de telles conditions.

Dans les deux cas, pour reprendre les termes de Doyle (33), il s'agit de tolérance hautement optimisée : de manière optimale, plus l'événement est rare, plus la fraction de la population susceptible d'y résister est faible.

4.4 Conclusion

A côté des duplications de gènes, nous avons montré l'importance des différentes formes de *recombinaisons* dans la convergence structure-fonction.

Certains cherchent à augmenter artificiellement les *recombinaisons* pour optimiser des bactéries (32; 209). De manière plus théorique, des modèles montrent même l'impossibilité de l'évolution en l'absence d'un taux de *recombinaison* suffisant. Xia et Levitt (206) développent par exemple un modèle d'évolution artificielle de *protéines* qui échoue à évoluer vers la forme la plus robuste en absence de *recombinaison*, du fait de l'énorme espace des possibles. Ce genre de « catastrophe de complexité », ainsi que le formule Kauffman (89), dans laquelle la sélection n'est plus efficace, ne peut être évitée que par une sélection permanente de l'évolutivité, c'est-à-dire par une amélioration continue de l'adéquation des variations structurelles aux variations fonctionnelles. De manière imagée, on peut dire que le « bricolage » évolutionnaire si cher à François Jacob (78) et présent à tous les niveaux d'observation du vivant (49) ne serait probablement pas efficace sans un « bricolage du bricolage » favorisant l'apparition d'assemblages bénéfiques.

Insistant sur l'importance de cette organisation du génome liée aux *recombinaisons*, les deux chapitres suivants présentent deux études en montrant une nouvelle facette pour la première et utilisant des résultats connus à son propos pour la seconde.

La première étude montre que les *opérons* sont plus longs lorsqu'ils sont *transcrits* dans la même direction que la *réplication*. Ceci suggère que la longueur des *opérons* est soumise à une pression de sélection. Ainsi, il pourrait exister une pression de sélection à la modularité, du simple fait que les modules de grande complexité sont défavorables pour l'organisme qui les porte.

La seconde étude se base sur l'existence d'une convergence structure-fonction entre les points de *recombinaison* qui délimitent les modules structurels de l'*ADN* et les gènes qui marquent une des plus fortes modularités fonctionnelles des organismes vivants, les premiers se situant majoritairement entre les seconds. Sur la base de cette hypothèse, nous construisons une nouvelle méthode d'analyse des études d'association génétique se focalisant sur l'obtention de résultats fonctionnels (association potentielle de gènes avec un caractère) plus que structurels (association potentielle de mutations localisées de l'*ADN* avec un caractère).

Chapitre 5

Le biais d'orientation et de longueur des unités de transcription bactériennes

5.1 Introduction

Si on peut lister la liste des bénéfices potentiels des *opérons* pour les organismes susceptibles de les intégrer et de les exploiter, il est nécessaire de constater certaines limites à leur existence.

- Tout d'abord, il faut que l'organisme hôte puisse survivre, au moins dans certaines conditions sans posséder l'*opéron*. Ainsi, les gènes correspondant aux *enzymes* du métabolisme central (ensemble des réactions chimiques liées à la production d'énergie par les cellules) ne sont pas de bons candidats pour former un *opéron*, car aucun organisme ne peut s'en dispenser.
- Ensuite, il ne suffit pas d'intégrer un *opéron*, encore faut-il qu'il soit *transcrit* de manière adéquate (au bon moment, au bon endroit). D'un côté, le fait de recevoir toute l'information génétique en une seule fois et d'avoir la possibilité de réaliser un seul *transcrit* (ARN messager) pour l'ensemble des gènes réduit le nombre de régions régulatrices fonctionnelles nécessaires : il suffit d'obtenir par mutation une seule région régulatrice en amont de l'*opéron* (car il est peu probable qu'une région régulatrice adaptée à un organisme fonctionne dans un organisme un tant soit peu éloigné).
- De l'autre côté, la force de l'*opéron* (la modularité et la mutualisation de la région régulatrice) est aussi une faiblesse. En effet, il se trouve – mais est-ce une coïncidence ? – que, en très grande majorité, les organismes acceptant les *opérons* sont aussi ceux qui n'ont qu'une origine de *réplication*. Pour ces organismes, il est fréquent que la vitesse de *réplication* soit un facteur limitant de la division. Or, la plupart des bactéries résolvent mal les collisions entre les *transcriptases* et la ou les *réplicases*. En règle générale, vu la pression afin d'accélérer la *réplication*, la *transcription* est purement et simplement arrêtée par la *réplicase*. Si le transcrit

(ARN messager) n'est pas complet, dans le meilleur des cas il se dégradera rapidement, dans le pire des cas, il produira une *protéine* incomplète qui pourra avoir des effets délétères. Il existe donc une pression de sélection visant :

- à raccourcir la longueur des gènes et des *opérons* ;
- à les orienter dans le sens de la *réplication*. Ainsi, ils sont lus dans le même sens à la transcription et à la réplication. La réplacase entre donc en collision moins souvent avec la transcriptase du fait d'une vitesse relative moins élevée.

Les *opérons* ont donc une taille limite, non pas tant par la capacité d'intégration d'ADN nouveau par la bactérie, mais simplement par son incapacité à transcrire efficacement les derniers gènes d'un *opéron* trop long lors des périodes de *réplication* rapide. De plus, on a montré que cette pression est d'autant plus grande que le gène est important pour la fonction de fitness de l'organisme (143) et que le gène est fortement exprimé (186).

Ainsi, à la manière de la fable de Tempus et Hora de Herbert Simon ¹, il existe une pression pour obtenir une forte modularité structurelle de l'organisation du génome : il est nécessaire que des *transcrits* relativement courts puissent être utilisés par l'organisme, autrement dit qu'ils aient une fonction.

C'est l'objet de l'article suivant que de montrer cette limitation intrinsèque de la longueur des gènes par un modèle liant le biais d'orientation, qui était déjà connu et dont on a approfondi depuis l'analyse (8), au biais en longueur, que nous avons mis en lumière. Ce biais en longueur suggère l'importance de la contrainte de longueur pour les unités de *transcription* bactériennes. Si cet article montrait qu'il n'était pas nécessaire que la *réplication* soit davantage ralentie par les collisions frontales que par les collisions colinéaires, des études postérieures à l'article et y faisant référence suggèrent que la pression de sélection est majoritairement sur la vitesse de *réplication*, fortement ralentie par les collisions frontales, plutôt que sur la diminution des *transcrits* incomplets (143).

5.2 Article

<http://dx.doi.org/10.1093/bioinformatics/bth317>

5.3 Conclusion

Il semble donc que la longueur des unités de *transcription* soit limitée chez les organismes à centre unique de *réplication*. Au contraire, chez les *eucaryotes*, de multiples centres de *réplication* existent. Ceci a relâché la contrainte sur la taille du génome, mais aussi sur la longueur des gènes, ce qu'on observe chez les *eucaryotes* et encore plus les *eucaryotes* supérieurs qui comptent des gènes très longs du fait de l'existence de *séquences non codantes* incluses dans les gènes, les *introns* (138).

1. <http://polaris.gseis.ucla.edu/pagre/simon.html>

Toutefois, les *eucaryotes* ont eux une capacité beaucoup plus limitée d'intégrer des *opérons*, probablement du fait que leur évolutivité ne réside plus dans leur capacité à intégrer des fonctions essentiellement métaboliques (diversification des capacités d'alimentation, diversification des milieux de vie, etc.) dans leur génome. Par exemple, les *animaux* s'appuient plutôt sur leur capacité à changer de morphologie (142) grâce, en particulier, à de réseaux de régulation beaucoup plus développés (93; 28; 79).

Ainsi, on peut se poser la question de l'évolutivité des organismes possédant à la fois la capacité d'intégrer des *opérons* et la capacité de gérer plusieurs points de *réplication*, à la manière de ce qui a été fait pour montrer que les *insectes à métamorphose complète* (plus modulaire fonctionnellement que ceux à *métamorphose incomplète* car capable de s'adapter indépendamment à leur milieu larvaire et à leur milieu adulte) ont connu un taux de diversification plus rapide (207) – cf. Section (2.2.3.3). Les *nématodes* (vers ronds) possèdent ces deux capacités, sans que l'on ait montré pour le moment une évolutivité supérieure de cet embranchement.

Enfin, si les biais de longueur et de position ont maintenant été bien étudiés, l'étude des contraintes structurelles portant sur le génome reste ouverte, notamment sur l'organisation en 3 dimensions de l'*ADN*. Par exemple, on a découvert une périodicité des localisations des gènes exprimés simultanément (101) sur l'*ADN*, dont on ne comprend pas encore avec précision la fonction, probablement liée à l'amélioration de leur régulation commune.

Chapitre 6

Les études d'association génétique

6.1 Introduction

Le but des études génétiques est d'identifier les variations de l'information génétique à l'origine de variations de caractères des individus qui les portent. Chez l'homme, ces caractères sont généralement des maladies ou des résistances à des médicaments. Ces études nécessitent de lire l'information génétique de plusieurs individus et d'analyser leurs différences. Les études les plus simples du point de vue conceptuel cherchent à corréliser les différences d'information génétique interindividuelles au caractère d'intérêt. Ces études sont dites « études d'association ». Des études plus compliquées nécessitant de modéliser la transmission de l'information cherchent à corréliser le caractère d'intérêt à l'information génétique d'un individu et de ses parents simultanément. Ces études sont dites de liaison. Dans tous les cas, la connaissance de la structure de l'information génétique est essentielle, comme nous allons le montrer dans ce paragraphe.

Nous avons précédemment montré l'importance de la *recombinaison* homologue dans l'origine de la structure actuelle du génome humain. En effet, que ce soit à partir d'observations directes ou d'inférence à partir de données génétiques, on observe des *points chauds de recombinaison* séparés par de longs espaces sans *recombinaison*. Nous avons vu l'intérêt pour l'évolutivité du contrôle précis des emplacements des *points chauds* afin que la modularité structurelle qu'ils créent se superpose à la modularité fonctionnelle associée à l'information qui les sépare (essentiellement de gènes ou de portions de gènes).

Une conséquence pratique de ce fort *déséquilibre de liaison* réside dans le fait qu'il n'est pas nécessaire de lire entièrement l'information génétique d'un individu (l'*ADN* contenu dans ses 23 paires de chromosomes, soit environ 2×3 milliards de bases, contenant l'équivalent de 12 gigaoctets d'information) pour pouvoir la caractériser. En effet, si on découpe le génome en chaque point où une ou plusieurs *recombinaisons* ont eu lieu dans l'histoire des populations humaines, on peut lister

les séquences possibles entre chacun de ces points. A la fois le nombre de points où il y a eu des *recombinaisons* et le nombre de séquences entre chacun de ces points est probablement relativement restreint. Ainsi, on peut espérer qu'un échantillonnage suffisamment dense (typiquement, lire une base sur 1000) permette d'associer de manière univoque les valeurs échantillonnées entre deux points de *recombinaison* à la seule séquence complète existant dans les populations humaines compatible avec ces valeurs. Ainsi, on peut espérer caractériser les différences d'information génétique entre individus en lisant uniquement 2×3 millions de positions dans chacun des *ADN*.

Cette idée est à la base des études d'association génétique à l'échelle du génome (57; 9; 196). Leur but est d'identifier des variants de l'*ADN* associés à la présence d'un caractère spécifique dans un groupe d'individus. Par exemple, la base T à la place de la base A en position 5 204 808 du chromosome 11 est connue pour être associée à certaines formes de drépanocytoses¹. La mise en œuvre d'études d'association demande beaucoup de précautions pour tenir compte des spécificités de l'information génétique (72; 60). On peut les résumer en deux points, que l'expérience accumulée par les premières études a permis de préciser (118) :

- un bon choix des individus. Il faut que l'échantillon soit d'une taille suffisante afin d'avoir la puissance statistique nécessaire pour détecter les biais escomptés. Il faut aussi une homogénéité de l'arrière plan génétique afin d'éliminer les risques de variables cachées, par exemple si une maladie est surreprésentée dans un sous-groupe de la population et l'homogénéité des caractères (tous les individus atteints doivent l'être d'une maladie ayant une cause unique) ;
- un bon échantillonnage de l'information génétique. En effet, en échantillonnant, on ne pense pas lire directement les bases associées au caractère, comme dans le cas de la drépanocytose : on espère lire suffisamment près pour qu'il n'y ait jamais eu de *recombinaison* entre la position lue et la séquence à l'origine du changement de caractère. En particulier, l'échantillonnage doit être réparti tout au long de l'*ADN* en tenant compte des points de *recombinaison*.

Supposons que la maladie étudiée soit effectivement lié à l'information génétique (on peut pour cela étudier le caractère chez des jumeaux n'ayant pas été élevés ensemble par exemple. S'ils développent toujours le même caractère, on peut supposer qu'il y a une origine génétique. Dans le cas contraire, les conditions de vie jouent certainement un rôle essentiel). Supposons que l'étude soit bien réalisée. A quelles conditions l'analyse peut-elle détecter le ou les variants génétiques associés à la maladie ? Une situation favorable est celle des caractères mendéliens. En effet, dans ce cas, la convergence structure-fonction est complète : dans le cas extrême de la drépanocytose, une simple mutation d'une base présente sur un chromosome suffit à provoquer l'apparition de la maladie sous une forme atténuée et sur les deux chromosomes sous une forme sévère. L'étude de ces cas est maintenant ancienne et a connu de grands succès. On pourrait imaginer de les étudier avec les moyens d'aujourd'hui en réalisant une étude d'association génétique. Il est probable que l'on parviendrait à

1. http://www.ncbi.nlm.nih.gov/entrez/dispmim.cgi?id=141900&a=141900_AllelicVariant0243

identifier des effets aussi forts. En revanche, il reste aujourd'hui à élucider des associations génétiques avérées (d'après les études de jumeaux séparés), mais qui ne sont pas mendéliennes. Il n'y a pas dans ce cas de convergence structure-fonction forte pour faciliter le travail : de multiples variants dispersés dans l'ensemble du génome interagissent très probablement entre eux et avec l'environnement pour déclencher une maladie (177). Toutefois, les techniques d'analyses multivariées adaptées à ces situations sont encore à développer (72). En effet, elles font face à une difficulté théorique liée à la faiblesse des effets à détecter par rapport à la taille de l'espace des génomes possibles.

Nous nous concentrons donc sur le développement d'une analyse univariée de l'association de variants génétiques et du caractère étudié. Nous souhaitons orienter l'analyse vers la mesure directe de l'association des gènes avec le caractère étudié. En effet, ceci a un grand intérêt pour faciliter l'analyse fonctionnelle qui découle d'une étude génétique car celle-ci vise en effet à transformer une information portant sur une corrélation en mécanisme causal permettant de reconstituer l'impact d'une mutation sur une maladie (50). Elle consiste souvent à caractériser le rôle d'un gène dans une maladie. Nous connaissons la tendance des points de *recombinaison* à se situer entre les gènes, qui est caractéristique d'une convergence structure-fonction. Nous savons que cette convergence structure-fonction n'est pas parfaite et qu'il n'existe pas encore de description détaillée et valable pour toutes les populations humaines des *points chauds de recombinaison*. Nous avons donc développé une méthodologie d'analyse des associations des gènes à des caractères discrets, en particulier de type cas/contrôle (malade/sain) en considérant que les *recombinaisons* se situent entre les gènes. Nous en présentons les principes et les résultats dans l'article suivant.

6.2 Article

Gene-based bin analysis of genome-wide association studies

Nicolas Omont^{1,2}, Karl Forner¹, Marc Lamarine¹, Gwendal Martin¹, François Képès² and Jérôme Wojcik*¹

Address: ¹Merck Serono International S.A., 9 chemin des Mines, 1202 Geneva, Switzerland and ²Epigenomics Project, Genopole®, 523 Terrasses de l'Agora, 91034 Évry cedex, France

Email: Jérôme Wojcik* - jerome.wojcik@merckserono.net

* Corresponding author

from Machine Learning in Systems Biology: MLSB 2007
Evry, France. 24–25 September 2007

Published: 17 December 2008

BMC Proceedings 2008, 2(Suppl 4):S6

This article is available from: <http://www.biomedcentral.com/1753-6561/2/S4/S6>

© 2008 Omont et al; licensee BioMed Central Ltd.

This is an open access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/2.0>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Abstract

Background: With the improvement of genotyping technologies and the exponentially growing number of available markers, case-control genome-wide association studies promise to be a key tool for investigation of complex diseases. However new analytical methods have to be developed to face the problems induced by this data scale-up, such as statistical multiple testing, data quality control and computational tractability.

Results: We present a novel method to analyze genome-wide association studies results. The algorithm is based on a Bayesian model that integrates genotyping errors and genomic structure dependencies. p -values are assigned to genomic regions termed bins, which are defined from a gene-biased partitioning of the genome, and the false-discovery rate is estimated. We have applied this algorithm to data coming from three genome-wide association studies of Multiple Sclerosis.

Conclusion: The method practically overcomes the scale-up problems and permits to identify new putative regions statistically associated with the disease.

Background

The last years have shown a tremendous increase in the number of markers available for association studies. Previous studies were dealing either with the whole genome at a very low resolution (for instance 5 264 microsatellites in [1]) or with a carefully chosen region of few millions of base pairs [2,3]. Recent technologies allow the genome-wide genotyping of hundred of thousands SNPs [4]. This has arisen the need of new methodological developments to overcome different issues, such as the multiple-testing

problem, gene biases, data quality analysis and the computational tractability.

Firstly, the multiple testing problem seems to cause association studies ability to detect associations to decrease as the number of markers increases. The classical analysis strategy, based on an association test for each marker [5], encounters increasing difficulties as more than one million of markers are available: Increasing the number of markers prevents from the detection of the mild genetic effects expected in complex diseases, as only strong effects

emerges from the huge noise generated by the increased quantity of data.

Methods like False Discovery Rate (FDR) [6] computation allow to control the error rigorously, but do not increase the statistical power. Better strategies based on haplotype blocks are being developed, the first step being gathering such block data (see the HapMap project, [7]). The gain of such strategies is two-folded: (i) the number of tests is independent of the number of markers (ii) the statistical power may be increased if markers of the same haplotype block are not fully correlated.

Secondly, a genetic association of a given SNP is a statistical feature and does not explain by itself a phenotype. To biologically interpret an associated marker, its haplotype block should first be delimited. Then, the association can be refined by fine-scale genotyping technologies or ideally by full resequencing. This eventually allows to identify functional mutations. Most of the time, these mutations impact relatively close genes. This is a first argument to bias association analysis towards genes. Moreover, even if haplotype blocks are unreachable, DNA might be cut into distinct regions (called *bins*) on another basis, so as to limit the multiple-testing problem and make it independent of the number of markers. Combining these two arguments leads to choose one bin for each gene, and to create "desert" bins in large unannotated regions. It allows to associate a list of genes with a test, which simplifies the analysis of results. The drawbacks are (i) that it makes more difficult the study of these "deserts", however the goal is here to maximize, not the chance of finding an association, but the chance of elucidating a mechanism of a complex disease given the current knowledge (ii) that a bin might contain several haplotype blocks, resulting in a dilution of the association signal if only one block is associated. Reciprocally, neighbor bins are not independent because they may share a haplotype block. However, with the classical strategy, correlated neighbor SNPs would also be tested separately.

Thirdly, genome-wide genotyping data are obtained by high-throughput experiments which encompass limitations requiring careful statistical methodology. Especially, with *Affy. technology*, the trade-off between the call rate (i.e. errors detected by the genotyping process and resulting in missing genotypes in the data set) and the error rate (i.e. errors left in the data) is difficult to adjust. Obtaining unbiased statistical results is then conditioned to good pre-processing filters. Indeed spurious markers must be eliminated and missing data correctly managed.

In addition, for most of SNPs used in this study, some genotypes are held by less than few percents of patients, which, given the usual collection size of a few hundreds,

(i) is not enough for good asymptotic approximations and (ii) should be considered with care given possible high error rate.

Finally, whatever algorithmic solution is developed, because the number of markers available will probably quickly reach a few millions, creating a scalability problem, it has to be linear in the number of markers.

In this paper we present a novel Bayesian algorithm developed to easily analyze genome-wide association studies. This algorithm is based on a gene-based partitioning of DNA into regions, called bins. A p -value of association is computed for each bin. The model takes into account genotyping errors and missing data and tries to detect simple differences in the haplotype block structure between cases and controls. The study of different collections is allowed. The multiple testing problem is addressed by estimation of FDR. The method has been applied to analyze the results of three genome-wide case-control association studies of the complex disease Multiple Sclerosis (MS). It identifies putatively associated bins, containing genes previously described to be linked to MS (see [8] for review) as well as new candidate genes.

Materials

Three association studies dealing with Multiple Sclerosis (MS) in three independent collections have been realized. Around 600 patients have been recruited for each study, half of them as cases affected by the disease, half of them as controls (Table 1). Genotypes of the 116 204 SNPs have been determined for each patient using Affymetrix GeneChip® human mapping 100 K technology (*Affy. technology*).

Methods

Notations

Stochastic variables are noted with a round letter (\mathcal{V}), a realization is noted in lower case (v). Indices are noted in lower case (k), ranging from 1 to the corresponding upper case letter (K). Unless needed, this range of indices ($k \in [1, K]$) is omitted. The number of different values is noted $\#(\mathcal{V})$. The n -dimensional table of the number of individuals having the same combination of values for given var-

Table 1: Genome-wide association multiple sclerosis collections.

Coll.	Origin	#Cases	#Controls	%Females
A	French	314	352	69
B	Swedish	279	301	71
C	American	289	289	85

variables $\mathcal{V}^k, k \in [1, K]$ (the contingency table) is noted $n(\mathcal{V}^1, \dots, \mathcal{V}^K)$. The marginalization of such a contingency table over one variable, for example \mathcal{V}^1 , is noted $n(\oplus, \mathcal{V}^2, \dots, \mathcal{V}^K) = \sum_{v \in \#(\mathcal{V}^1)} n(v, \mathcal{V}^2, \dots, \mathcal{V}^K)$. Estimation of a probability distribution $P(\mathcal{V})$ is noted with hatted letter, $\hat{P}(\mathcal{V})$. Each bin $b \in [1, B]$ contains J_b genetic markers \mathcal{G}_b^j with $j \in [1, J_b]$. Each patient $i \in [1, I]$ has a phenotype value $s(i)$ (in case-control studies, $\#(S) = 2$), discrete co-variable values $v_m(i), m \in [1, M]$ (gender: $m = 1$, or collection of origin: $m = 2$), and a genotype value for each marker $g_b^j(i)$ (with SNPs, $\#(\mathcal{G}_b^j) = 3$). A patient i is represented by this vector:

$$i = [s(i), v_m(i), g_b^j(i)] \text{ with } : m \in [1, M], b \in [1, B], j \in [1, J_b] \quad (1)$$

The data set is noted $D = \{i\}_{i \in [1, I]}$. A first level of the method aggregates predictors at the bin level. The "restriction" of a patient to a bin is noted i_b , the corresponding data set being $D_b = \{i_b\}_{i \in [1, I]}$:

$$i_b = [s(i), v_m(i), g_b^j(i)] \text{ with } : m \in [1, M], j \in [1, J_b] \quad (2)$$

Data preprocessing

Due to *Affy. technology* (the D.M. calling algorithm), errors on heterozygotic genotypes are more frequent. It can be detected through the deviation of a SNP from the Hardy-Weinberg equilibrium, which basically states that, noting $P(a) = P(aa) + P(Aa)/2$ and $P(A) = P(AA) + P(Aa)/2$:

$$\begin{cases} P(aa) = P(a)^2 \\ P(Aa) = 2P(a)P(A) \\ P(AA) = P(A)^2 \end{cases} \quad (3)$$

Therefore, the following pre-processing filters are applied: SNPs are discarded (i) if the number of missing genotypes is higher than 5% because the genotyping process quality was low for this SNP, (ii) if the minimum allele frequency in controls $MAF = \min(P(a), P(A))$ is lower than 1%, because the SNP holds no information, or (iii) if the probability that the SNP follows the Hardy-Weinberg equilibrium in controls is lower than 0.02.

Bin definition

Bins are defined on DNA from protein genes as defined in the version 35.35 of Ensembl [9] of the human DNA sequence. The basic region of a gene lie from the begin-

ning of its first exon to the end of its last exon. Overlapping genes are clustered in the same bin. If two consecutive genes or clusters of overlapping genes are separated by less than 200 kbp, the bin limit is fixed in the middle of the interval. Otherwise, the limit of the upstream bin is set 50 kbp downstream its last exon, the limit of the downstream bin is set 50 kbp upstream its first exon, and a special bin corresponding to a *desert* is created in between the two bins. With these rules, desert bins have a minimum length of 100 kbp (Figure 1).

Assessing bin association

General model, hypotheses and statistics

We assume that each bin constitutes an independent data set. The following ideal probability distribution is defined:

$$\forall b \in [1, B], P(I_b) = P(S, \mathcal{V}_m, \mathcal{G}_b^j) \quad (4)$$

As experimenters choose cases and controls (phenotypes) each individual subset of the study is a realization of the conditional distributions $P(\mathcal{G}_b^j | S, \mathcal{V}_m)$. Estimations of probability distribution are possible from contingency tables:

$$\hat{P}(\mathcal{G}_b^j | S, \mathcal{V}_m) = \frac{n(S, \mathcal{V}_m, \mathcal{G}_b^j)}{n(S, \mathcal{V}_m, \oplus)} \quad (5)$$

On the contrary, due to the experimental design, estimations of $P(S, \mathcal{V}_m)$ are impossible.

A general way to assess the association of a bin b is to estimate whether $(\mathcal{G}_b^j)_{j \in [1, J_b]}$ is independent from the pheno-

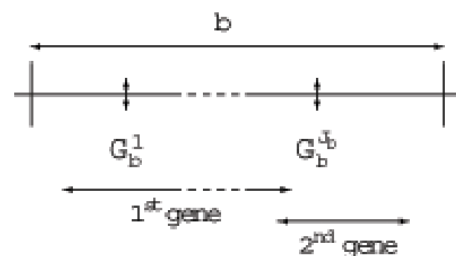


Figure 1
Representation of a bin containing two genes and J_b markers.

type S , i.e., whether $P(\mathcal{G}_b^j|\mathcal{V}_m)$ is "far" from $P(\mathcal{G}_b^j|S, \mathcal{V}_m)$.

$$H_0^b : P(\mathcal{G}_b^j|S, \mathcal{V}_m) = P(\mathcal{G}_b^j|\mathcal{V}_m) \quad (6)$$

However, as only $P(\mathcal{G}_b^j|\mathcal{V}_m, S)$ is estimable, estimation of $P(\mathcal{G}_b^j|\mathcal{V}_m)$ is not possible. Therefore, one estimates $\widehat{P}_{H_0^b}(\mathcal{G}_b^j|\mathcal{V}_m)$ assuming H_0^b , as indicated by the subscript:

$$\widehat{P}_{H_0^b}(\mathcal{G}_b^j|\mathcal{V}_m) = \frac{n(\oplus, \mathcal{V}_m, \mathcal{G}_b^j)}{n(\oplus, \mathcal{V}_m, \oplus)} \quad (7)$$

We have chosen likelihood ratio LR as a statistic to estimate the "distance" between estimations of $P(\mathcal{G}_b^j|S, \mathcal{V}_m)$ and $P(\mathcal{G}_b^j|\mathcal{V}_m)$. For each patient, the LR is:

$$LR(i_b) = \frac{p(g_b^j(i)|s(i), v_m(i))}{\widehat{p}_{H_0^b}(g_b^j(i)|v_m(i))} \quad (8)$$

As all patients are considered to be independently chosen, the LR of the set of patients available is:

$$LR(D_b) = \prod_{i \in [1, I]} LR(i_b) \quad (9)$$

p-value estimation and FDR

To assess estimation errors due to randomness and sample size, the probability that H_0^b is true given the observation, i.e. the p -value π_b needs to be computed. This is theoretically achieved by enumerating all possible outcomes $D_b(\sigma)$ of the experiment that lead to the observed data $D_b(\sigma_0)$ (σ is a enumeration parameter to be defined. The following notation simplification is done: $D_b(\sigma_0) = D_b$). Then the probability $p(D_b(\sigma))$ of each outcome assuming that H_0^b is true is computed as well as its LR. Finally, the p -value is:

$$\begin{aligned} \pi_b &= p(LR(D_b(\sigma)) \geq LR(D_b)) \\ &= \sum_{\sigma | LR(D_b(\sigma)) \geq LR(D_b)} p(D_b(\sigma)) \end{aligned} \quad (10)$$

In this article, estimation of p -values is based on permutations: possible outcomes are obtained through patient phenotype permutations σ and σ_0 is the identity permutation. The probability of each permutation is uniform. The

denominator of equation (8) is constant with respect to such permutations, therefore it is omitted. Sampling this space is possible: random permutations of the phenotypes are drawn and used to compute a LR. This is a Monte-Carlo procedure, for which we propose an optimized implementation that guarantees the precision required for FDR estimation:

For each bin b , compute LR for new permutations of phenotypes until the number of permutations realized N_b satisfies the following equation, noting $\hat{\pi}_b$ the estimation of the bin p -value:

$$N_b \geq \left(B\theta \frac{\gamma}{\delta} \right)^2 \min \left(\frac{1-\theta}{\theta}, \frac{1-\pi_b}{\pi_b} \right) \quad (11)$$

θ and γ/δ control the quality of the method: θ is an upper bound of the threshold that is expected to be used to select bins. γ/δ controls the error due to the randomness of the process: Assuming that two consecutive p -values $\pi_{b1} < \pi_{b2} \approx \theta$ are sufficiently spaced (probability $p_s = e^{-\delta}$), $\hat{\pi}_{b1} < \hat{\pi}_{b2}$ with a confidence $c = \text{cdf}(\mathcal{N}(0, 1), \gamma)$ (standard normal cumulative distribution function). In this article, $B = 11264$, $\theta = 0.001$, $\delta = 1$ and $\gamma = \sqrt{2}$ thus $N_b = 507003$, $p_s = 0.37$ and $c = 0.92$.

To address multiple testing, the method uses an FDR estimation defined as in [10]:

$$FDR(\theta) = \frac{\widehat{\Pi}_0 \theta B}{\#\{b | \pi_b < \theta\}} \quad (12)$$

The numerator is an estimation of the expectation of the number of false-positive with $\pi_b \leq \theta$. $\widehat{\Pi}_0$ is an estimation of the proportion of bins under the null hypothesis. Given that it is expected to be very high in current study, it is (conservatively) fixed at its upper bound: $\widehat{\Pi}_0 = 1$. The denominator is the number of tests with p -values below. The ratio is therefore an estimation of the proportion of false negatives in the set of bins with a p -value below θ . Because we want to analyze thoroughly the FDR for around the 10 bins with the lowest p -values, the FDR is not controlled at a specified threshold as in [6] but only estimated.

This estimation relies on two main hypothesis: (i) tests are independent or positively correlated [11], (ii) p -values are continuously and uniformly distributed in $[0, 1]$. Assuming that sharing of haplotype block by neighbor bins is the

only source of correlation between tests, the positive correlation seems reasonable. Indeed, if the p -value of a not associated bin decreases, the p -values of bins sharing the same haplotype block are more than likely to decrease too. The uniform distribution is less obvious, because the number of possible contingency tables is finite so that even the null distribution is not uniform. However, the sample size is one to two order of magnitude higher than in other applications of FDR to discrete data in which the problem is acute [12].

Model of linkage disequilibrium and error

Correlation between markers induced by LD is modelled with an inhomogeneous hidden Markov chain of order 1. Indeed, as a rough approximation, for each marker, most information is found on its first neighbor on each direction of DNA. In a directed graphical model, independence assumptions consist in:

$$P(\mathcal{G}_b^j | \mathcal{G}_b^l)_{l \neq j} = \begin{cases} P(\mathcal{G}_b^j | \mathcal{G}_b^{j-1}) & \text{if } j \neq 1 \\ P(\mathcal{G}_b^j) & \text{if } j = 1 \end{cases} \quad (13)$$

Finally, this assumptions also allow to obtain correct estimations because corresponding contingency tables are sufficiently filled. They implies that contingency tables are computed for 2 SNPs ($\#(\mathcal{G}_b^j) = 3$), the phenotype ($\#(\mathcal{S}) = 2$) and the co-variables together. The gender co-variable is not be used. It requires the hypothesis that the SNP distribution is independent from it. The only co-variable is the study patients belong to (Table 1, $\#(\mathcal{V}_2) = 3$). As collection sizes for a given study are around 600, the average number of patients in each cell of contingency tables is then $\bar{n} = 33$.

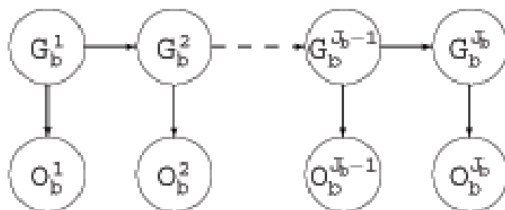


Figure 2
Error and LD model of bin b .

An error model (Figure 2) is introduced linking observed genotypes O_b^j with real ones ($O_b^j \in \{aa, Aa, AA, \emptyset\}$, where \emptyset means that the observed genotype is missing):

$$P(O_b^j | (\mathcal{G}_b^l)_{l \in [1, j_b]}) = P(O_b^j | \mathcal{G}_b^j) \quad (14)$$

Since \mathcal{G}_b^j are hidden variables, estimation of a priori probabilities of $P(\mathcal{G}_b^j | \mathcal{G}_b^{j-1})$ and $P(O_b^j | \mathcal{G}_b^j)$ is not straightforward. Usual strategy is to use an Expectation-Maximization (E.-M.) algorithm to infer the state of hidden variables. However, it is not required in order to assess bin associations. Therefore, an alternative strategy is developed. $P(\mathcal{G}_b^j | \mathcal{G}_b^{j-1})$ and $P(\mathcal{G}_b^1)$ are estimated through the removal of patients with missing genotypes:

$$\hat{P}(\mathcal{G}_b^j | \mathcal{G}_b^{j-1}) = \frac{n(O_b^j, O_b^{j-1}) + C}{n(\oplus, O_b^{j-1}) - n_{\emptyset} + mC} \quad (15)$$

Where n_{\emptyset} is the number of patients with either O_b^j or O_b^{j-1} missing and m is the number of cells. To obtain more regular estimates, a constant is added to all cell counts. It is a Dirichlet prior on parameters. This constant is chosen to be $C = \alpha_0 \bar{n}$, where α_0 is the chosen error rate and \bar{n} is the mean number of individuals per cell. This constant means that uncertainty on low cell counts is high, not only because of randomness, but also because of genotyping errors.

On the other hand, given the previously developed structure of errors, the following model of $P(O_b^j | \mathcal{G}_b^j)$ is chosen:

$$P(O_b^j | \mathcal{G}_b^j) = \begin{pmatrix} O_b^j \backslash \mathcal{G}_b^j & aa & Aa & AA \\ aa & 1-\beta & 1-2\beta & 0 \\ & 1-\alpha & \alpha & \\ Aa & 1-\beta & 1-2\beta & 1-\beta \\ & \alpha & 1-2\alpha & \alpha \\ AA & 0 & 1-2\beta & 1-\beta \\ & & \alpha & 1-\alpha \\ \emptyset & \beta & 2\beta & \beta \end{pmatrix} \quad (16)$$

The missing rate β is estimated for each marker through the resolution of the non-linear system drawn from the preceding model. The maximum error rate α_0 is estimated

during external comparison of *Affy. technology* and other technologies. In this study, the error rate is chosen to be $\alpha_0 = 0.05$. The error rate is $\alpha = \min(\alpha_0, P(O_b^j = Aa)/(1 - P(O_b^j = \emptyset)))$ in order that the system always have a solution for β .

Likelihood computation

With the current model, the likelihood of a patient is the sum of the likelihoods over all possible combinations of real genotypes:

$$L_1(i_b) = p(o_b^j(i)) = \sum_{g_b^j \in \{1, \dots, \mathcal{G}_b^j\}} \left(\prod_{j>1} p(o_b^j(i) | g_b^j) p(g_b^j | g_b^{j-1}(i)) p(o_b^j(i) | g_b^j) p(g_b^j) \right) \quad (17)$$

This is a computation in $O(\prod \#(\mathcal{G}_b^j)) \sim O(3^{J_b})$. Some approximations in the model are required to obtain computations linear with the number of markers. The following one is based on two-marker sliding windows and corresponds to the model of Figure 3:

$$L_2(i_b) = \prod_{j \geq 2} \sum_{g_b^{j-1}, g_b^j} \left(p(o_b^j(i) | g_b^j) p(g_b^j, g_b^{j-1}) p(o_b^{j-1}(i) | g_b^{j-1}) \right) \quad (18)$$

This equation considers information coming from two neighbor markers together. Compared to the full model, information flow is limited to pair of markers. The likelihood could be falsely increased in this extreme situation: suppose that a missing genotype is inferred *aa* from its left neighbor and *AA* from its right neighbor, the merging of this two inferences would results in a contradiction and thus a low resulting likelihood. On the contrary, the approximated likelihood does not detect this contradiction and is falsely increased. This likelihood is named thereafter "two-marker" likelihood.

Simplifying further leads to consider markers one by one. There is no model of linkage disequilibrium anymore, but

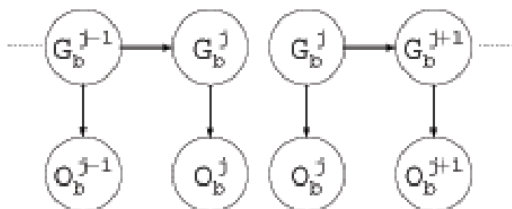


Figure 3
Simplified model of two-marker likelihood computation.

noise is reduced as cells are better filled. This likelihood is named thereafter "naive likelihood" because it corresponds to a naive Bayesian model:

$$L_3(i_b) = \prod_j \sum_{g_b^j} p(o_b^j(i) | g_b^j) p(g_b^j) \quad (19)$$

Results

The method has been applied to each of the three collections A, B, C (Table 1) as well as to the three collections at once (ABC), considering the collection of origins as a co-variable. The overall computation time is about 10 days on a single processor.

The pre-processing filters discard around 20% of SNP: for collection A (resp. B and C), out of 112 463 SNP, 84 430 (resp. 93 548 and 86 652) SNP remains. If all SNP satisfied the Hardy-Weinberg equilibrium, 2 249 SNP are

Table 2: Associated bins at FDR 5% threshold (top), at FDR 50% threshold before (middle) and after exclusion of MHC region bins (bottom). A, B, C, ABC: collection designs, L₂: two-marker likelihood, L₃: naive likelihood.

FDR 5% with MHC	L ₃	L ₂
A	3	2
B	3	6
C	2	2
ABC	4	6
<hr/>		
FDR 50% with MHC	L ₃	L ₂
A	6	6
B	14	7
C	6	28
ABC	20	33
<hr/>		
FDR 50% w/o MHC	L ₃	L ₂
A	2	0
B	1	1
C	0	0
ABC	8	10

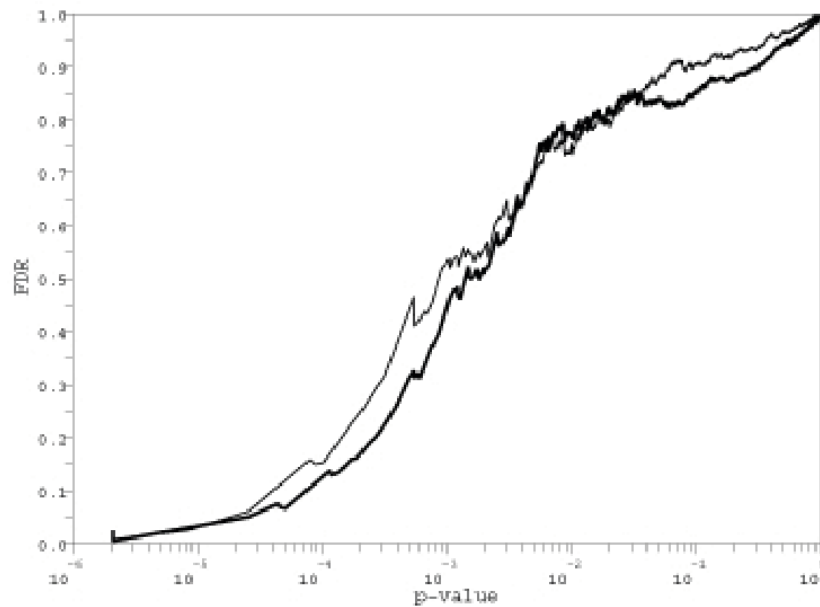


Figure 4
FDR versus p -values of bins sorted in increasing order for the three collections design (ABC). Thick line: two-marker likelihood L_2 , thin line: naive L_3 .

expected to be discarded. 9 422 were for collection A. It can be explained (i) by artifacts of DM calling algorithm which has a higher error rate on heterozygotic genotypes (ii) by deviations from the assumptions underlying this theoretical equilibrium. The bin partitioning algorithm divides the genome into 19 556 gene bins and 1 993 desert bins. Out of these 21 549 bins, only 11 264 (52%) contain one SNP or more after pre-processing in at least one collection and are considered for further analysis. Before pre-processing, out of 12 512 SNP with one bin or more, 2 781 have only one SNP, and 2 188 bins 10 SNP or more. The maximum is 210.

Figure 4 shows the FDR plotted against p -values computed using the two-marker L_2 or the naive L_3 likelihood for the three collection design. Two-marker FDR remains below naive FDR until a p -value level of 0.01 and both increase slowly towards 1. FDR against the number of selected SNP plots are detailed by collection in Figure 5. As observed in other studies [13], the FDR is not monotonous with the p -value. The oscillations are less important for the three collection design, maybe because of the three time increase of sample size. With a FDR threshold of 5%, only between 2 and 6 bins are selected depending on the collections and likelihood considered (Table 2, top). Most of them are located, in the Major Histocompatibility Complex (MHC) region, mainly in the class III subregion. The class II subregion is known to be associated with MS [14]. The three collection design selects more associated bins than one

collection designs, independently on the likelihood. Results with a less stringent FDR threshold of 50% (Table 2, middle) shows a greater power of L_2 over L_3 for the three collection design. However, FDR is misleading in this study because the MHC region is known to be associated with MS. It leads to an overestimation of the FDR at which bins outside of this region are selected. It contains 12 of the 33 bins selected by L_2 on the three collection design. As a result, only 10 and not 21 bins are selected (Table 2, bottom).

Discussion

We have developed a new method to practically analyze genome-wide association studies data. Our algorithm is based on a bin partitioning of the genome, takes advantage of studying several collections simultaneously, takes into account genotyping errors and local genomic structure (LD), and handles the multiple testing problem through FDR estimation while staying computationally tractable. The method has been applied to analyze three association studies in Multiple Sclerosis.

The FDR threshold is chosen according to the desired application. To conduct expensive further experiments with putatively associated genes, a very low rate of false-positives is required. A FDR threshold of 5% seems reasonable. On the contrary, if one wants to minimize the false-negative rate, a FDR of 50% is acceptable.

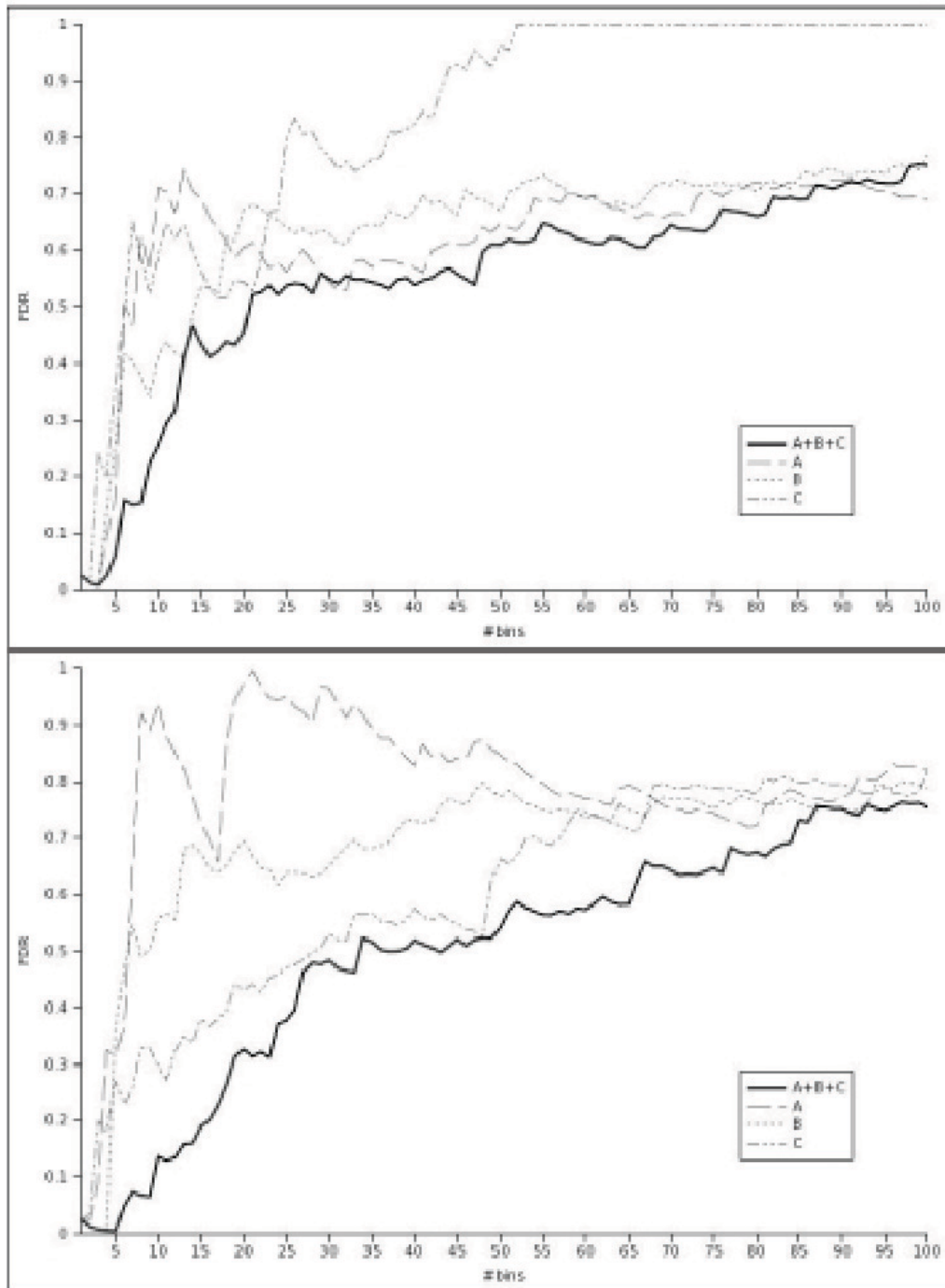


Figure 5
FDR versus number of bins selected using L_3 naive likelihood (top) and L_2 two-marker likelihood (bottom). A: solid, B: dash, C: dash dot, ABC: thick.

Applying the method to experimental genome-wide association data on three collections permits (i) to assess the algorithm and evaluate the different parameters and design and (ii) to identify genes potentially associated to Multiple Sclerosis. We have evidenced that the three collection design outperforms the one-study design in terms of expected number of true-positives, despite differences between the studied collections, especially on the severity of the disease. Furthermore, with this three collection design, the two-marker likelihood L_2 seems to be more efficient thanks to the additional information used. With this configuration, a FDR threshold of 5% gives 6 associated bins. Four of them are located in the MHC region, known to be linked to Multiple Sclerosis [14]. It is a validation of the method. The two others are bins containing olfactory receptor genes *OR2T2* and *OR4A47*. The biological meaning of such association is unclear but the extended MHC regions contain many other olfactory genes [14] and olfactory dysfunction has already been reported in Multiple Sclerosis [15]. At FDR threshold of 50% and after exclusion of bins from MHC, the method selects ten bins. They open the perspective of insights to explain Multiple Sclerosis.

Competing interests

The authors declare that they have no competing interests.

Authors' contributions

NO, FK and JW conceived and designed the model. NO, KF, ML and GM wrote the analysis tool. NO, FK and JW wrote the manuscript.

Acknowledgements

We are grateful to the Serono Genetics Institute banking, genotyping and genetic analysis team for producing high-quality data. The work has been significantly made easier by the Serono Genetics Institute Research Knowledge Management team and we acknowledge them particularly. We also thank Pierre-Yves Bourguignon for the idea of the hidden Markov chain and Jean Duchon for the distribution of the distance between two p -values. This article has also been greatly improved thanks to the comments of many reviewers.

This article has been published as part of *BMC Proceedings* Volume 2 Supplement 4, 2008: Selected Proceedings of Machine Learning in Systems Biology: MLSB 2007. The full contents of the supplement are available online at <http://www.biomedcentral.com/1753-6561/2?issue=S4>.

References

- Dib C, Fauré S, Fizames C, Samson D, Drouot N, Vignal A, Millasseau P, Marc S, Hazan J, Seboun E, Lathrop M, Gyapay G, Morissette J, Weissenbach J: **A comprehensive genetic map of the human genome based on 5,264 microsatellites.** *Nature* 1996, **380(6570)**:152-4.
- Cardon LR, Bell JL: **Association study designs for complex diseases.** *Nature reviews Genetics* 2001, **2(2)**:91-99.
- Lewis CM: **Genetic association studies: design, analysis and interpretation.** *Briefings in bioinformatics* 2002, **3(2)**:146-153.
- Kennedy GC, Matsuzaki H, Dong S, min Liu W, Huang J, Liu G, Su X, Cao M, Chen W, Zhang J, Liu W, Yang G, Di X, Ryder T, He Z, Surti U, Phillips MS, Boyce-Jacino MT, Fodor SPA, Jones KW: **Large-scale**

genotyping of complex DNA. *Nature biotechnology* 2003, **21(10)**:1233-7.

- Klein RJ, Zeiss C, Chew EY, Tsai JY, Sackler RS, Haynes C, Henning AK, SanGiovanni JP, Mane SM, Mayne ST, Bracken MB, Ferris FL, Ott J, Barnstable C, Hoh J: **Complement factor H polymorphism in age-related macular degeneration.** *Science* 2005, **308(5720)**:385-9.
- Benjamini Y, Hochberg Y: **Controlling the false discovery rate: a practical and powerful approach to multiple testing.** *Journal of the Royal Statistical Society. Series B* 1995, **57(1)**:289-300.
- International HapMap Consortium T: **A haplotype map of the human genome.** *Nature* 2005, **437(7063)**:1299-320.
- Dyment DA, Ebers GC, Sadovnick AD: **Genetics of multiple sclerosis.** *Lancet neurology* 2004, **3(2)**:104-110.
- Birney E, Andrews D, Caccamo M, Chen Y, Clarke L, Coates G, Cox T, Cunningham F, Curwen V, Cutts T, Down T, Durbin R, Fernandez-Suarez XM, Flicek P, Gräf S, Hammond M, Herrero J, Howe K, Iyer V, Jekosch K, Kähäri A, Kasprzyk A, Keefe D, Kokocinski F, Kulesha E, London D, Longden I, Melsopp C, Meidl P, Overduin B, Parker A, Proctor G, Prlic A, Rae M, Rios D, Redmond S, Schuster M, Sealy I, Searle S, Severin J, Slater G, Smedley D, Smith J, Stabenau A, Stalker J, Trevanion S, Ureta-Vidal A, Vogel J, White S, Woodwork C, Hubbard TJP: **Ensembl 2006.** *Nucleic Acids Research* 2006:D556-D561.
- Storey JD, Tibshirani R: **Statistical significance for genomewide studies.** *Proceedings of the National Academy of Sciences of the United States of America* 2003, **100(16)**:9440-5.
- Benjamini Y, Yekutieli D: **The control of the false discovery rate in multiple testing under dependency.** *Annals of Statistics* 2001, **29(4)**:1165-1188.
- Pounds S, Cheng C: **Improving false discovery rate estimation.** *Bioinformatics* 2004, **20(11)**:1737-45.
- Pounds SB: **Estimation and control of multiple testing error rates for microarray studies.** *Briefings in bioinformatics* 2006, **7**:25-36.
- Horton R, Wilming L, Rand V, Lovering RC, Bruford EA, Khodiyar VK, Lush MJ, Povey S, Talbot CC, Wright MW, Wain HM, Trowsdale J, Ziegler A, Beck S: **Gene map of the extended human MHC.** *Nature reviews Genetics* 2004, **5(12)**:889-899.
- Zivadinov R, Zorzon M, Bragadin LM, Pagliaro G, Cazzato G: **Olfactory loss in multiple sclerosis.** *Journal of the neurological sciences* 1999, **168(2)**:127-130.

Publish with **BioMed Central** and every scientist can read your work free of charge

"BioMed Central will be the most significant development for disseminating the results of biomedical research in our lifetime."

Sir Paul Nurse, Cancer Research UK

Your research papers will be:

- available free of charge to the entire biomedical community
- peer reviewed and published immediately upon acceptance
- cited in PubMed and archived on PubMed Central
- yours — you keep the copyright

Submit your manuscript here:
http://www.biomedcentral.com/info/publishing_adv.asp



Erratum

A la place de l'Equation (13), il faut lire :

$$P(\mathcal{G}_b^j | \mathcal{G}_b^l)_{l < j} = \begin{cases} P(\mathcal{G}_b^j | \mathcal{G}_b^{j-1}) & \text{if } j > 1 \\ P(\mathcal{G}_b^j) & \text{if } j = 1 \end{cases}$$

$$P(\mathcal{G}_b^j | \mathcal{G}_b^l)_{l > j} = \begin{cases} P(\mathcal{G}_b^j | \mathcal{G}_b^{j+1}) & \text{if } j < J_b \\ P(\mathcal{G}_b^j) & \text{if } j = J_b \end{cases}$$

6.3 Conclusion

En conclusion, si l'application de la méthode d'étude d'association des gènes aux études génétiques de Serono Genetics Institute sur la *sclérose en plaques* donne des résultats satisfaisants comparés à la méthode standard – elle identifie des régions situés dans le *Complexe Majeur d'Histocompatibilité (CMH)*, comme cela était déjà connu et a été confirmé depuis (174) –, il faut en constater les limites, pour l'essentiel liées au concept même de l'étude d'association génétique chez l'homme, mais aggravées par le biais fonctionnel introduit par l'étude directe de l'association des gènes.

A titre d'illustration, nous trouvons que la séquence génétique du gène *NOTCH4* est associée à la *sclérose en plaques*. Toutefois, le détail de la zone, tel que présenté dans la figure (6.1), montre un fort *déséquilibre de liaison* avec le début de la sous-région 2 du *CMH* (76). Ainsi, on ne peut en fait être aussi précis sur l'association de *NOTCH4*. La mutation ayant réellement un effet fonctionnel peut se trouver dans l'ensemble de la région en *déséquilibre de liaison*. Il s'agit là d'une myopie inhérente aux études d'associations génétiques : sans *recombinaison*, il est impossible de préciser la localisation de l'association. Cette myopie est aggravée par le découpage fonctionnel : l'association structurelle découverte ne peut simplement être attribuée fonctionnellement au gène *NOTCH4*. D'autres études (non génétiques) seraient nécessaires pour comprendre les mécanismes causant cette association de certains variants du *CMH* à la *sclérose en plaques*.

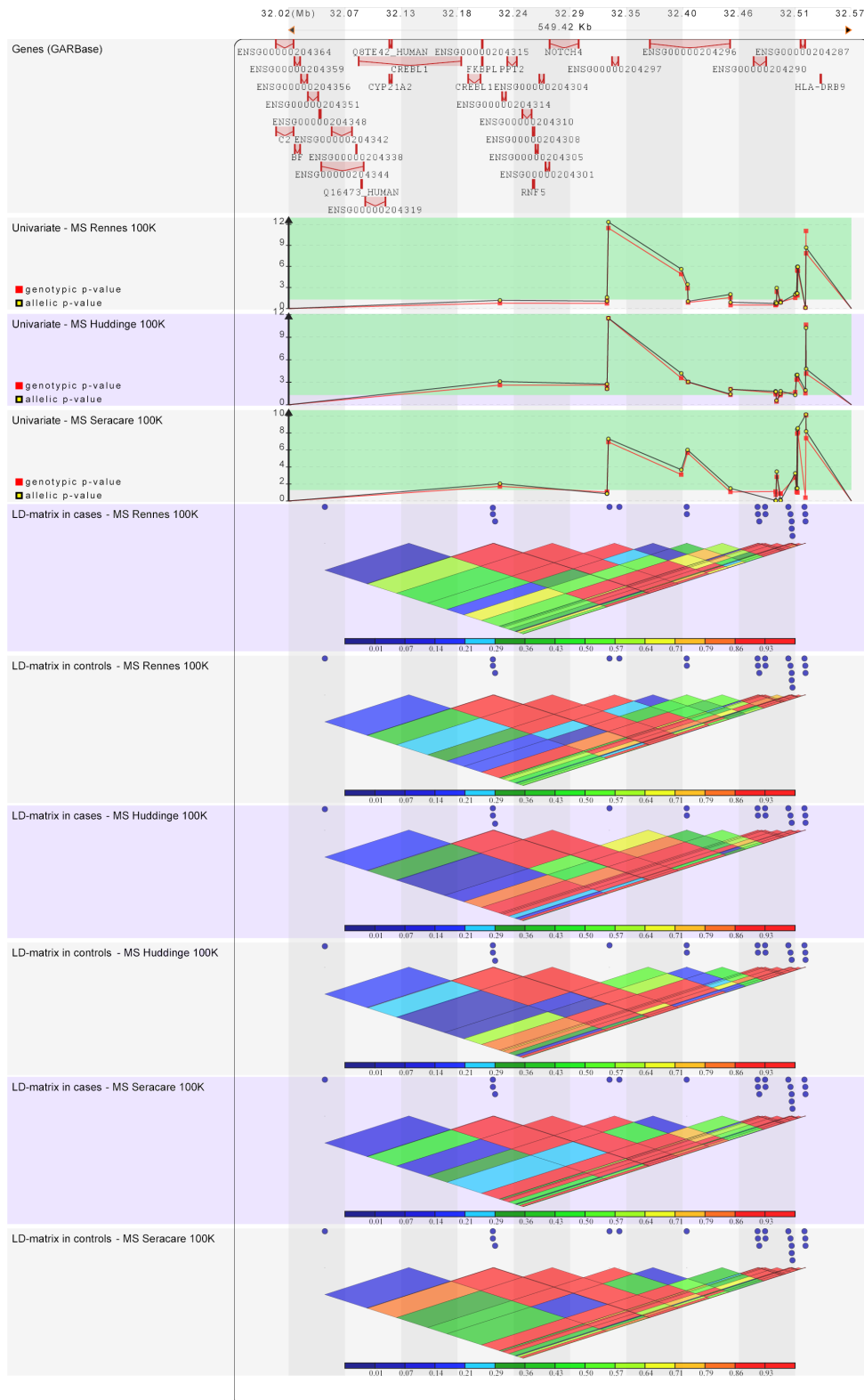


FIGURE 6.1 – Représentation graphique de l'association des marqueurs génétiques autour de *NOTCH4* et des déséquilibres de liaison dans les études génétiques sur la sclérose en plaques de Serono Genetics Institute.

Chapitre 7

Conclusion

Si les gènes et les *opérons* sont le signe d'une forte convergence structure-fonction nécessaire à l'évolutivité des systèmes reproductifs vivants, on constate qu'il existe des limites de différentes sortes à cette *convergence*.

- D'une part les *opérons* ont une taille limitée du fait de l'incapacité des bactéries à résoudre les collisions *réplicases-transcriptases* sans interrompre la *transcription* et/ou ralentir la *réplication*. Ainsi, supposons qu'il soit nécessaire d'avoir un *opéron* très grand afin d'assurer une fonction non modularisable (très complexe), les bactéries n'évolueraient certainement pas vers une telle solution. Elles exploreraient plutôt les solutions sous-optimales comportant plusieurs modules structurels (gènes et *opérons*) pour un seul module fonctionnel. Par exemple, il existe de nombreuses protéines qui ne sont actives que lorsqu'elles s'assemblent pour former un complexe. Même si ce mécanisme d'assemblage fournit d'autres avantages (notamment de pouvoir réguler l'activité du complexe en régulant son assemblage), on peut faire l'hypothèse que le fait de conserver des *opérons* courts a aussi un rôle dans l'explication de l'omniprésence de tels complexes.
- D'autre part, la disposition des *points chauds de recombinaison* montre bien l'existence d'un biais vers les *recombinaisons* conservant la structure génique, mais ce biais est relativement faible chez l'homme : il existe des points de *recombinaison* dans les gènes, c'est-à-dire plusieurs modules structurels pour un seul fonctionnel, et des gènes non séparés par des points de *recombinaison*, c'est-à-dire plusieurs modules fonctionnels pour un seul structurel. Ainsi, les régions à fort *déséquilibre de liaison* constituant des modules structurels regroupant de nombreux gènes qui n'ont aucun lien fonctionnel entre eux sont très nombreuses. A ce titre, on peut noter le caractère renforcé de cette observation dans les populations européennes car elles descendent d'un groupe restreint (57; 151; 208) qui a vécu il y a relativement peu de temps (entre 15 000 et 30 000 ans auparavant). Par conséquent, seules des recombinaisons qui ont eu

lieu depuis sont visibles dans les populations contemporaines. Ainsi, le fait que deux gènes non apparentés fonctionnellement mais voisins sur l'*ADN* n'aient pas été séparés par des recombinaisons est encore plus fréquent dans les populations européennes que dans les populations africaines.

Ainsi, si la première étude a été concluante pour montrer l'existence d'un moyen détourné d'une sélection amplifiée de la modularité par limitation de la taille des *opérons* (même en cas de non-modularité, des solutions modulaires avec des *opérons* courts sont privilégiées), la seconde a montré qu'il était nécessaire d'améliorer les découpages du génome humain permettant d'inférer une association fonctionnelle d'un gène à partir de l'association structurelle d'une séquence. En effet, il est pour cela indispensable de lier l'association d'une séquence à celle d'un gène (et donc à sa fonction), ce qui n'est pas possible lorsque plusieurs gènes n'ont jamais été séparés par des *recombinaisons*. Une piste d'amélioration serait de découper le génome en utilisant à la fois les limites de gènes et les *points chauds de recombinaison*, en regroupant les gènes qui ne sont pas séparés par de tels *points chauds*. Ceci permettrait d'obtenir des régions structurellement indépendantes car séparées par des points chauds, tout en regroupant ensemble les gènes pour lequel l'analyse sera dans tous les cas « myope » puisque ces gènes ont toujours été liés dans l'histoire de la population étudiée.

Troisième partie

Convergence structure-fonction limitée de la distinction producteur-transporteur dans les systèmes électriques

Chapitre 8

Introduction

8.1 Généralités sur les systèmes électriques

8.1.1 Utilité des systèmes électriques

L'activité des sociétés modernes est conditionnée à leur consommation continue d'énergie. Du fait de sa facilité d'utilisation (lumière, conversion en énergie mécanique ou thermique, utilisation pour l'alimentation des systèmes d'information), l'énergie électrique a joué très tôt un rôle central dans leur fonctionnement.

En absence de contraintes, l'énergie serait toujours produite là où elle est consommée, car son transport a un coût. Dans la pratique, cette situation est rare. C'est pourquoi le transport d'énergie est devenu une activité importante. A cette fin, de multiples réseaux ont été mis en place.

Dans le cas de l'électricité, le transport de l'énergie répond à plusieurs nécessités :

- La centralisation de la production : la production a tendance à être centralisée lorsque la productivité d'une grosse centrale est meilleure que celle d'un grand nombre de petits groupes. En contrepartie, la distance moyenne entre le producteur et les consommateurs augmente.
- Les contraintes géographiques : si la centralisation de la production était le seul facteur conduisant au transport de l'électricité, les producteurs d'électricité seraient proches des lieux de consommation. C'est d'ailleurs le cas d'un certain nombre de centrales. Cependant, d'autres facteurs peuvent conduire à éloigner producteurs et consommateurs :
 - La présence d'une source d'énergie (énergie potentielle de l'eau, soleil, vent), ou d'un pré-requis technique (eau de refroidissement des centrales nucléaires) ;
 - L'arbitrage entre coût de transport de l'énergie électrique par rapport à une autre (centrales installées dans les ports).
- La mutualisation des moyens de production : la capacité de production est dimensionnée pour faire face à la consommation la plus élevée. En agrandissant la taille d'un système électrique,

on diminue la variabilité de la consommation et de la production. Par exemple, l'existence d'un grand nombre de consommateurs « lisse » et rend prévisible la courbe de consommation, alors même que la consommation de chacun ne l'est pas. De même, pour continuer à desservir la demande en cas d'arrêt de plusieurs groupes de production, ce qui est un événement courant, il faut disposer de groupes de secours en quantité suffisante. En l'absence de réseau, ces groupes devraient être présents sur tous les lieux de consommation ; cela pourrait amener à en construire un très grand nombre. En construisant des lignes, on peut mutualiser ces groupes et en construire ainsi beaucoup moins, les mêmes groupes pouvant être utilisés pour « assurer » différents cas d'indisponibilité de groupes. Ceci est d'autant plus rentable que la distance de transport de l'électricité est courte, ce qui est le cas en Europe. En effet, le coût d'une ligne de transport est alors plus faible que celui d'un groupe de production de même puissance.

8.1.2 Caractéristiques techniques des systèmes électriques

8.1.2.1 Un stockage extrêmement limité

Comme pour toutes les énergies, le premier principe de la thermodynamique s'applique à l'électricité : il y a conservation de l'énergie électrique. Toute l'énergie produite est perdue, consommée ou stockée pour utilisation ultérieure.

Au contraire de la plupart des énergies, l'électricité est techniquement très difficile à stocker. A proprement parler, seuls les condensateurs stockent l'énergie électrique. D'autres dispositifs de stockage existent, mais ils sont en fait capables de convertir l'électricité en une énergie stockable (chimique pour les piles, potentielle pour les systèmes couplant des réservoirs d'eau à des altitudes différentes) puis de reconstituer l'électricité à partir du stock. Dans les systèmes actuels, ces dispositifs sont peu répandus.

Ainsi, dans un système électrique, à chaque seconde, le flux d'énergie électrique produite (la puissance produite) est égal au flux d'énergie consommée (la puissance consommée, pertes incluses). Ce n'est pas le cas des autres systèmes énergétiques (ceux des combustibles fossiles en particulier), dans lequel de nombreux stockages sont présents (des réservoirs souterrains de gaz naturels aux réservoirs des automobiles).

Un système électrique est donc essentiellement constitué de producteurs, de consommateurs et d'un réseau de transport. Il comporte très peu de dispositifs de stockage. Dans cette thèse, les quelques dispositifs de stockage seront considérés comme des consommateurs ou des producteurs suivant qu'ils stockent ou rendent de l'énergie au système.

La principale conséquence est que la valeur financière d'une même quantité d'énergie est très variable au cours du temps, d'une heure à l'autre par exemple. Ce n'est pas le cas pour d'autres sources d'énergie, comme les combustibles fossiles, car leur stockage est aisé.

8.1.2.2 Un contrôle restreint

Transport Pour l'essentiel, la répartition des flux dans les lignes d'un réseau n'est pas explicitement contrôlée : elle suit les lois de la physique. D'autres réseaux se comportent de la sorte (gaz, eau), mais, dans le cas de l'électricité, les moyens de contrôle des flux sont restreints et/ou indirects.

En particulier, il n'existe pas de moyen direct de limiter le flux dans une ligne à la capacité de la ligne, c'est-à-dire de « robinet » en quelque sorte. A défaut, si le flux dépasse la capacité, la ligne est déconnectée automatiquement du réseau afin de la préserver. Dans ce cas, plus aucun flux ne la traverse, alors que l'optimum serait que le flux soit bloqué à la valeur maximum admissible. Ce genre de situation est inconnu de la plupart des réseaux. En télécommunications, par exemple, la saturation se traduit soit par l'impossibilité d'une nouvelle connexion (réseau téléphonique commuté) soit par des données perdues (Internet). Pour un réseau de fluide, on peut limiter le débit dans l'ouvrage saturé à l'aide de vannes.

Dans un premier temps, une telle déconnexion de ligne n'a que peu de conséquences sur la production et la consommation sur le réseau (seules les pertes peuvent être légèrement modifiées). L'électricité est donc acheminée par les lignes voisines selon les lois de la physique. Si celles-ci se trouvent à leur tour en surcharge, elles se déconnectent les unes après les autres, jusqu'à entraîner la séparation du réseau en plusieurs parties non connexes. Apparaît alors un déséquilibre entre production et consommation dans chaque partie qui est résolu comme l'indique le paragraphe suivant.

Production Il existe des dispositifs de contrôle de l'équilibre production-consommation dans tous les systèmes électriques. Certains permettent aussi un contrôle indirect des flux. En voici une description rapide.

On appelle zone synchrone une partie d'un réseau électrique dans lequel tous les alternateurs de production tournent à la même fréquence et en même temps. Il est impossible de connecter deux réseaux – et donc d'échanger de l'énergie – avec des lignes ordinaires s'ils ne forment pas une zone synchrone. De plus, comme nous allons le voir, tous les producteurs d'une zone synchrone sont solidaires pour assurer l'équilibre production-consommation.

Connaissant les lois de la physique et grâce à des prévisions de consommation, il est possible de construire des plans de production et de prévoir les flux dans une zone synchrone à partir d'une situation de production et de consommation équilibrée donnée. En particulier, en faisant varier la production (déplacement de puissance d'une région à une autre) ou en délestant une consommation, on contrôle indirectement les flux dans les lignes.

Ce système comprenant entre autre le réglage dit secondaire est hautement centralisé. En effet, il est nécessaire de coordonner des actions à l'échelle d'une zone synchrone qui peut couvrir tout un continent afin d'assurer à la fois l'équilibre production-consommation et l'admissibilité des flux dans les ouvrages. Le réglage secondaire n'est pas modulaire, car il nécessite qu'un organisme centralisé

puisse donner une consigne à chaque producteur du réseau. Toutefois, face à la difficulté d'une telle centralisation, une organisation basée sur plusieurs zones de réglages à l'intérieur d'une même zone synchrone a été retenue. Malgré tout, une coordination conséquente des actions entre zones de réglage reste nécessaire.

Cependant, le délai associé à ce contrôle dit secondaire est de l'ordre de la minute. Ceci n'est pas suffisamment rapide en regard de la faible inertie du système. Il existe donc un système entièrement automatique réagissant en quelques secondes à un déséquilibre entre production et consommation. Ce système dit de réglage primaire est décentralisé : si la production n'est plus suffisante, la fréquence du système, ordinairement stabilisée à 50 Hz en Europe, baisse globalement. En effet, les alternateurs ralentissent car ils doivent faire face à une charge plus élevée. Ce signal est détecté par certains producteurs qui augmentent instantanément leur production. Si ce système est modulaire au sens où les actions sont décentralisées, cette modularité ne se superpose pas à celle des zones de réglages, car les actions sont toujours réparties dans toute la zone synchrone. Il se crée donc une solidarité de fait entre les producteurs de toute la zone.

Même si ces deux contrôles permettent de garantir en permanence l'équilibre production-consommation, l'un ne tient pas compte des contraintes du réseau et l'autre possède un délai d'action important. Il demeure donc des situations à risque dont voici un exemple. Si le déséquilibre est trop important (suite à la séparation du réseau en parties non-connexes par exemple), les producteurs ne parviennent pas à ajuster la production à la consommation grâce aux réglages primaires et/ou secondaires. Dans ce cas, afin de préserver leurs groupes de production, ils les déconnectent. En les déconnectant, ils aggravent le déséquilibre et entraînent d'autres déconnexions. En quelques secondes, le système s'arrête : il ne subsiste aucun flux d'énergie.

Ces effets dominos, évoquant la fable de Tempus et Hora de Herbert Simon¹, sont d'une ampleur impressionnante. Ainsi, la coupure volontaire mais malencontreuse d'une ligne dans le Nord de l'Allemagne a entraîné des coupures de courant jusqu'à Tunis dans l'heure qui a suivie, la phase finale de propagation du réseau régional allemand à l'ensemble de la zone s'effectuant elle-même en moins d'une minute (77).

Enfin, dans la Section (1.2.3.3), nous avons exposé un algorithme général de décomposition de graphes par retranchement successif des nœuds les plus connectés. Il semble en effet que la fonction de la plupart des graphes réels comporte entre autre celle de garantir la connexité de ces éléments. Par ce paragraphe, nous constatons que cela n'est qu'une vision très approximative de la réalité. Un réseau qui semble connexe à la suite du retranchement d'un nœud ne le serait pas nécessairement en réalité du fait des pannes en cascades. De plus, dans chaque zone connexe, le service de l'ensemble de la demande n'est pas assuré non plus. Les systèmes électriques sont donc moins robustes aux attaques ciblées sur des nœuds que l'étude de la connexité ne le montre. Malgré ces limites, une étude met en

1. <http://polaris.gseis.ucla.edu/pagre/simon.html>

évidence que les réseaux qui ont le moins d'indisponibilité sont ceux qui ont le plus de boucles, c'est-à-dire un rapport du nombre de liens sur le nombre de nœuds élevé (170). Ainsi, comme en biologie – cf. Section (2.2.3.2), la simple topologie du réseau apporte des informations utiles sur sa fonction.

Finalement, si le contrôle de l'équilibre production-consommation est complexe et si ses défaillances peuvent entraîner séparation du réseau en parties non-connexes, la section suivante expose les moyens mis en œuvre pour assurer son bon fonctionnement.

8.1.2.3 Stabilisation des systèmes électriques

Ainsi il existe des avantages à transporter l'électricité, mais ils sont contrebalancés par la difficulté à exploiter des grands systèmes, difficulté qui croît avec leur taille. Afin de la surmonter, deux grandes classes de moyens ont été mis en œuvre, le premier agissant sur la structure et le second sur la fonction telles que nous les avons définies dans la section (2.1.2.1) :

- à technologie égale, une meilleure prévision des risques de pannes et leurs conséquences pour éviter de placer le système dans un état à risque. Cela nécessite une vision globale du système, qui se traduit en général par une centralisation des prises de décision ;
- l'introduction de nouvelles technologies telles les lignes à courant continu afin d'augmenter la modularité fonctionnelle du système.

Anticipation des situations à risque Malgré des leviers de contrôles restreints et indirects, les écroulements de réseaux électriques sont rares. Il existe donc des moyens permettant de les éviter. Le principal consiste à anticiper à court terme comme à long terme les états possibles du réseau et à vérifier par la simulation numérique qu'ils sont stables.

Cette anticipation passe en premier lieu par la prévision de la consommation et des flux à l'avance, notamment du jour pour le lendemain. Fort heureusement, celle-ci est généralement prévisible avec une grande précision (157).

Ensuite, parmi les règles les plus connues, celle dite du « N-1 » consiste à vérifier par la simulation que l'état du réseau sera compatible avec les flux prévus dans les éléments disponibles du système (producteurs, consommateurs, réseau de transport), mais aussi avec les flux qui se réaliseraient si n'importe lequel d'entre eux venait de surcroît à tomber en panne.

Comme une panne locale (de production ou de transport) peut avoir un effet global, cette vérification de la compatibilité du réseau avec les flux prévus nécessite d'avoir une vision globale de son état. Ceci a naturellement poussé à centraliser l'exploitation des systèmes électriques.

Autrement dit, la fonction même des systèmes électriques est complexe car leur stabilisation n'est pas modulaire : elle est extrêmement difficile sans une vision et des moyens d'actions globaux. De ce fait, afin d'assurer la convergence structure-fonction, on lui a fait correspondre une structure de décision complexe, c'est-à-dire centralisée et intégrée (non-modulaire).

Lignes asynchrones à courant continu Dans un système électrique, les déconnexions en cascade amplifient l'effet d'une panne initiale. L'absence de contrôle du flux circulant dans une ligne du réseau est une cause souvent déterminante dans un tel effet.

Il existe cependant des lignes électriques ne présentant pas ce défaut : ces lignes sont à courant continu. Il est possible de régler le flux d'énergie sur une ligne de ce type de manière directe. Ainsi, le risque de déconnexion de la ligne suite à une surcharge disparaît. De ce fait, elles ne participent pas aux cascades de déconnexion de lignes. En revanche, elles sont beaucoup plus onéreuses que des lignes ordinaires. Enfin, elles permettent les échanges d'énergie entre zones synchrones.

Ainsi deux zones électriques constituées de lignes ordinaires (à courant alternatif) reliées par des lignes à courant continu peuvent être beaucoup plus aisément exploitées de manière indépendante. En effet, elles ont des réglages primaires indépendants et s'isolent mutuellement de la propagation en cascade des pannes de lignes. Le changement de technologie permet de modulariser la fonction des systèmes électriques : la fonction de stabilisation du système (réglage primaire) se décompose en une fonction de stabilisation par zone, mais le couplage fonctionnel souhaité de transfert d'énergie reste possible.

8.1.3 Caractéristiques économiques des systèmes électriques

La section précédente montre que les difficultés propres aux systèmes électriques poussent à la centralisation de leur exploitation, du moins tant que les lignes à courant continu ne sont pas utilisées pour connecter différentes parties du réseau. Dans la pratique, la centralisation implique généralement qu'une seule entité ait le contrôle, au moins à court terme et au moins en dernier ressort, sur la production (quelle quantité est produite, à quel endroit, etc.), la consommation (quel consommateur est alimenté) et la topologie du réseau (quelle ligne est sous tension, etc.).

On peut ajouter que, comme de nombreux réseaux (réseaux de transports routiers et ferroviaires par exemple) notamment, les réseaux de transports électriques constituent des monopoles naturels, dans lesquels le coût marginal est inférieur au coût moyen. En effet, les coûts marginaux de développement de capacités de transport sont décroissants : il existe des économies d'échelle, notamment liées à l'existence de coûts « fixes » (non liés à sa capacité) en particulier lors de la construction d'une ligne (139). On observe ce résultat dans le fait que, plus la capacité d'une ligne est petite, plus son taux d'occupation est faible (94). De plus, il est aujourd'hui très difficile d'obtenir l'autorisation de construire de nouvelles lignes en Europe. Cependant, ce n'est pas cette situation de monopole naturel qui explique la centralisation, c'est bien la difficulté d'exploitation. Ainsi, avant la création de la SNCF, il existait de nombreuses compagnies ferroviaires en France, possédant chacune un réseau. L'exploitation indépendante était possible car l'interconnexion de deux réseaux ferroviaires est relativement simple. Par exemple, il était possible d'organiser des correspondances piétonnes entre deux gares proches appartenant à deux compagnies différentes.

Ainsi, la conjonction des difficultés techniques d'exploitation (ensemble du système) et du caractère de monopole naturel (transport uniquement) ont conduit naturellement à l'existence de monopoles de production et de transport d'électricité en Europe. La section suivante détaille ce phénomène.

8.2 Système électrique européen

8.2.1 Caractéristiques techniques

Historiquement, du fait des contraintes techniques d'exploitation et de la nature de monopole naturel – cf. Section (8.1.3), des monopoles de production et de transport de l'électricité se sont constitués, dans chaque pays d'Europe.

Cependant, du fait du caractère stratégique de l'électricité, la centralisation s'arrêta au niveau national en Europe occidentale. Ainsi, on obtint une juxtaposition de systèmes complexes (zones synchrones) gérés de manière centralisée. L'intérêt du transport électrique à l'échelle internationale n'avait pas pour autant diminué. C'est ainsi que se fit dans l'après deuxième guerre mondiale la première connexion Suisse-France afin d'exploiter pleinement le potentiel hydroélectrique suisse (175). Depuis cette époque, il est nécessaire de trouver les moyens de gérer un système complexe et fonctionnellement non modulaire par une structure modulaire, essentiellement calquée sur les pays. En effet, les moyens de contrôle du système ne sont pas modulaires – cf. Section (8.1.2.3). A cette fin, les deux axes ont été exploités :

- la coordination de l'exploitation des réseaux, par la mise en place de processus standardisés d'exploitation et par l'échange d'information sur l'état des réseaux. Cette coordination a lieu à l'intérieur de chaque zone synchrone ;
- la connexion de zones synchrones par des liens asynchrones (lignes à courant continu rattachées aux deux zones synchrones par des stations de conversion alternatif/continu).

C'est ainsi que le système a évolué vers les zones synchrones de la figure (8.1). On constate l'existence de vastes zones synchrones non connectées (UCTE et IPS/UPS par exemple) ou connectées par des lignes à courant continu (NORDEL, Grande Bretagne et Irlande). Afin d'illustrer les conséquences de ce découpage, on peut relever que, sur la carte, l'Ukraine est coupée en deux. Cela signifie qu'il ne peut pas y avoir de transfert d'énergie entre l'ouest de l'Ukraine, connecté à l'UCTE et l'est, connecté à l'IPS/UPS.

De plus, à l'intérieur des zones synchrones, on trouve une ou plusieurs organisations chargées du service d'équilibrage entre production et consommation – il existe en général une zone de réglage par pays (responsable entre autre du réglage secondaire de sa zone) –, alors même que les déséquilibres éventuels ont un impact sur l'ensemble des acteurs de la zone (producteurs, transporteurs, consommateurs). Cette *non-convergence* des modularités structurelles et fonctionnelles est compensée par une coordination importante, notamment par la mise en place de standards (181).

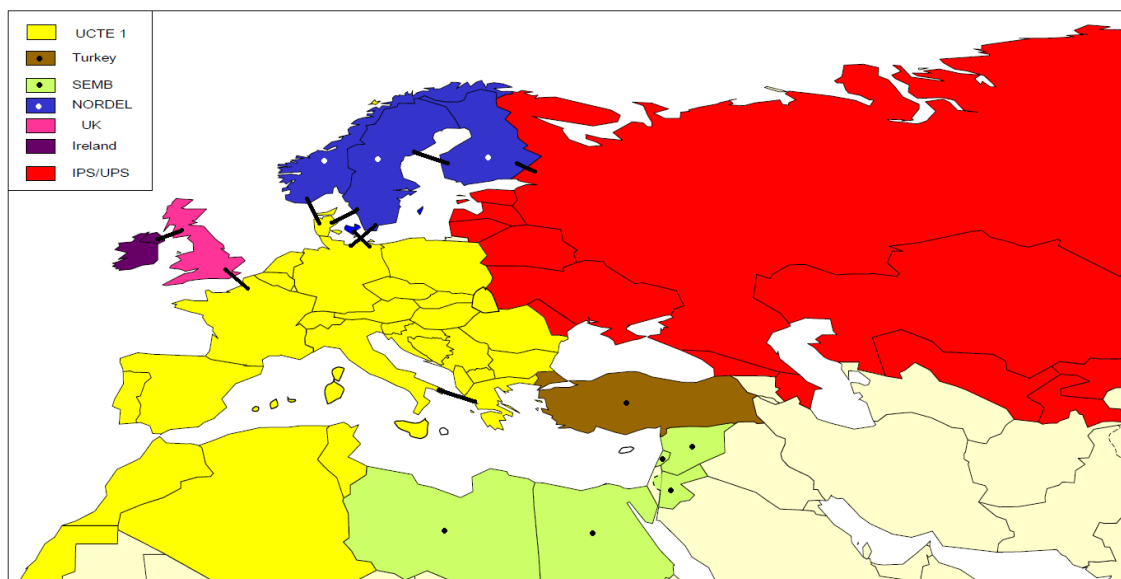


FIGURE 8.1 – Zones synchrones européennes en 2006 (extrait de « European, cis and mediterranean interconnection : State of play 2006 » (175)). Les lignes noires ajoutées représentent les principales connexions à courant continu existant en 2006.

8.2.2 Caractéristiques économiques

Dans le but d'accroître la productivité de l'économie européenne (41), l'Union Européenne a décidé de créer un marché européen de l'électricité où, idéalement, chaque consommateur serait libre de choisir son fournisseur d'électricité dans toute l'Europe (29). Dans la perspective de cette thèse, on peut dire que cette restructuration a pour but de changer la forme de la variabilité structurelle du système électrique européen. En effet, comme illustré dans la Section (2.2.2), le changement de structure du système correspond à modifier ses évolutions possibles afin de lui permettre d'évoluer vers plus d'efficacité (en particulier dans sa prise en compte des besoins des consommateurs) et plus réactif aux changements fonctionnels (technologies, coûts des énergies).

Dans la pratique, les contraintes techniques placent d'emblée des limites à une possible libre concurrence car le service rendu aux consommateurs n'est pas la simple fourniture d'énergie, mais un service complexe alliant :

- la fourniture d'énergie ;
- en un lieu déterminé ;

mais aussi :

- la possibilité de consommer plus ou moins dans certaines limites en temps réel ;
- un niveau de certitude sur la continuité de la fourniture d'énergie.

Ceci rend extrêmement difficile la mise en place de marchés où les prix représentent bien l'ensemble des coûts correspondants aux services rendus. C'est-à-dire qu'une des difficultés de la réforme consiste à ne pas créer d'*externalités* là où il n'en existait pas, ce qui aurait pour résultat de « désoptimiser » un système que la centralisation pouvait permettre d'optimiser.

Afin de simplifier le service soumis à *concurrence*, l'Union Européenne décida en 1996 d'orienter la libéralisation vers une libéralisation de la production, chaque consommateur pouvant choisir son producteur et tous les producteurs ayant les mêmes conditions d'accès au « réseau » (137). En revanche, dans cette vision, l'ensemble des services associés au réseau (notamment transport et continuité de la fourniture d'énergie) restent des monopoles centralisés.

C'est ainsi que, pour le transport et les services auxiliaires (gestion de l'équilibre, etc.), du fait des contraintes techniques, l'Union Européenne a mis en place des gestionnaires de réseaux, indépendants des producteurs, à la fois responsables de l'exploitation et des investissements dans le réseau de transport. S'il est possible de séparer l'exploitation et les investissements dans le réseau de transport comme le fait l'Ecosse par exemple (119), la solution privilégiée par l'Union Européenne est bien celle de gestionnaires de réseaux contrôlant l'exploitation et les investissements. En effet, il apparaît en pratique extrêmement délicat de mener les investissements adéquats pour une entreprise qui n'exploite pas le réseau (30). Ainsi, au moins pour l'exploitation, ces gestionnaires sont des monopoles sur un territoire donné. Ils sont régulés par des commissions publiques indépendantes dans chaque pays. Elles sont chargées de contrôler qu'ils n'abusent pas de leur situation de monopole.

Les relations entre processus de transport et production, auparavant intégrés au sein des mêmes entreprises sont donc à réinventer. L'objet du seul et unique chapitre de cette partie est de modéliser le système de transport européen dans le cadre de la théorie libérale classique – cf. Section (2.2.2.2) – afin que les gestionnaires de réseaux puissent donner des signaux tarifaires de long terme aux producteurs et aux consommateurs. Ceci permettrait de réduire au maximum les *externalités* que pourrait créer la dissociation de la production et du transport d'électricité malgré l'impossibilité de créer un marché du transport d'électricité.

Chapitre 9

La tarification de long terme des réseaux électriques

9.1 Introduction

Dans le contexte européen actuel, une réflexion fondamentale permet d'orienter et d'éclairer les discussions sur le problème de création d'un marché efficace de l'électricité. En se basant sur la théorie classique, en accord à la fois avec l'objectif de l'Union Européenne d'améliorer la productivité par la création de marché et avec la longue tradition d'optimisation des coûts par les monopoles intégrés, nous considérerons que la fonction d'un système électrique est de maximiser le surplus social de l'ensemble des acteurs du système (59), à savoir les producteurs, le ou les transporteurs et les consommateurs, tout en assurant l'équilibre des comptes du ou des transporteurs (201). Ce surplus social correspond au « bénéfice » que la société dans son ensemble retire du système. Il est défini comme étant la différence entre le prix que les consommateurs auraient été prêts à payer l'énergie électrique consommée d'un côté et les coûts de production et de transport de l'autre.

L'article ci-dessous présente un tel modèle mathématique des systèmes électriques qui promet d'être utile aux gestionnaires de réseau pour donner des indications tarifaires de long terme aux producteurs et aux consommateurs. D'une part, celles-ci peuvent orienter leur décisions de manière optimale. D'autre part, utilisées pour construire des tarifs, elles permettent aux gestionnaires de réseau de disposer des moyens financiers pour réaliser les investissements nécessaires. En effet, dans le cadre d'un marché européen de la production électrique, les investissements dans de nouvelles capacités de production se feront lorsqu'on anticipera des prix de marché suffisamment élevés. Sous condition d'un marché bien organisé, ceci devrait conduire à la maximisation du surplus social annoncée. Au contraire, pour le transport, comme les monopoles régulés échappent au marché, il convient d'élaborer des stratégies tarifaires permettant d'investir au mieux au regard du surplus social et de donner

les signaux tarifaires adéquats aux acteurs du système afin qu'ils puissent en tenir compte dans leurs décisions.

Très concrètement et très simplement, vu que les investissements, aussi bien en transport qu'en production ont des durées de vie de plusieurs dizaines d'années et relèvent de processus de décisions longs et complexes, il revient aux acteurs d'échanger toute l'information possible le plus en amont possible, notamment à travers des signaux tarifaires de long terme afin d'éviter des situations contre-productives très simples telles que :

- la construction de lignes vers un lieu où aucune capacité de production ne s'installe ;
- réciproquement, la construction de capacité de production dans un lieu où les capacités de transport ne sont pas disponibles. Ceci est d'autant plus vraisemblable qu'il faut aujourd'hui plus de temps pour construire une ligne à partir du moment où la décision est prise que pour construire certains types de génération (éoliennes, comme cela se produit actuellement aux Etats-Unis ¹, ou centrales à gaz).

Ainsi, la coordination des investissements de production et de transport, si elle n'était pas nécessairement aisée dans un monopole intégré (107; 184), doit trouver de nouveaux moyens d'être optimale dans le contexte européen actuel (45; 19). Cet article donne les bases pour y répondre en ébauchant ce que pourrait être un signal tarifaire de long terme.

9.2 Article

1. <http://www.iht.com/articles/2008/08/26/business/grid.php>

Long term nodal pricing and transmission costing

Nicolas Omont, Arnaud Renaud

Artelys
Paris, France
<http://www.artelys.com>

Patrick Sandrin

Réseau de Transport d'Electricité
Paris, France
<http://www.rte-france.com>

Abstract - This article questions costing principles through their links with long term nodal pricing and discusses remaining steps to the implementation of the associated model. The long term nodal pricing model differs from the standard nodal pricing model in that it considers transmission capacities as decision variables.

We show that a tariff based on long term nodal prices on an optimal network is equivalent to a tariff based on the compensation of each network element development costs proportionally to their usage by network transactions (balanced production–demand sets), i.e. based on the proportional costing principle.

In a multi–situation framework, we show that development costs of each line are mostly allocated to situations in which the line is saturated. However, we show that cost allocation between saturating situations is complex as it entangles spatial and temporal effects. Intrinsically linked with the fact that the network is dimensioned to face contingencies, the analysis of multiple saturating situation cost allocation is the key remaining step to the model implementation.

Keywords - Long Term Nodal Pricing, Marginal Pricing, Proportional Transmission Costing, Multi-period Transmission Cost Allocation.

1 Introduction

The revenues of Transmission System Operators (TSOs) are usually based on two paradigms. On the one hand, market mechanisms lead to fix prices for the usage of network elements. On the other hand, tariffs are built from cost evaluations. They are adjusted so as to ensure economic sustainability. Each TSO uses each paradigm to a different extent, based on their respective advantages and drawbacks. Few links have been established between them.

Many market structures rely on the theory of nodal pricing [1]. This theory asserts that the maximization of the social welfare resulting from the network usage is reached by choosing a market structure that leads to fix energy prices at each network node (nodal prices) at the value of specific dual variables of an optimal power flow problem ([2], See [3] for an introduction). This advocates the fact that such nodal prices would provide the right incentives to each market actor. However, in its usual formulation in which transmission capacities are parameters, it computes short term prices which do not allow TSOs to recover their expenses [4]. Therefore, it has to be completed by other mechanisms like access or “club membership”[5].

Conversely, costing principles used to establish tariffs

are numerous (See [6] for review) partly because arguments in favor of one or another often lack of theoretical grounds. In particular, two questions are regularly debated: (a) Transaction based power flow methods face the problem of counter flow pricing. Should flows whose direction is opposite to the line net flow get credit for lightening the flow, pay according to their absolute value, or pay 0? (b) How should periodic (annual) costs be allocated to each time slot? Should they be allocated to the peak one, to each of them (if yes, with which weights?), or based on each line peak state?

In this article, we introduce a new framework in order to compute long term nodal prices. Its properties allow to show links with specific costing principles, theoretically supporting them while pointing out the non–optimality of others. As developed below, the model differs from the short–term nodal price framework on two points: (a) Network expansion; (b) Simultaneous optimization over several situations. It induces a peak load pricing [7] in a model with fixed production and variable transmission capacities.

Firstly, considering network expansion is intrinsically linked with long term nodal prices and full cost recovery. Indeed, the model minimizes development costs, which includes both investment and operation costs. These development costs would be the prices paid by a TSO if it rented the assets to third party owners. On the one hand, these costs would reflect the expected asset usage during its full economical life, i.e. they would be long term costs. On the other hand, at the beginning of an optimization period, capacity renting would effectively be a decision variable in hands of the TSO. We show that a nodal pricing based on such an optimization scheme leads to the full recovery of development costs. This reconciliation property induced by long term transmission expansion planning has already been noted by [8]. However, given the choice of a continuous model to the expense of realism, interpretation of dual variables as nodal prices is possible. Thanks to this property, we show that this nodal pricing is equivalent to compensation of each network element development cost proportionally to its usage by transactions occurring on the network. This founds costing principles based on: (a) Transactions; (b) Cost allocation proportional to network usage; (c) Credit to counter flows. To the best of our knowledge, this result is original.

Secondly, several situations need to be considered because networks are dimensioned to face operating conditions resulting from: (a) Normal variations of demand and production, both in quantity and location; (b) Contingencies [9]. In the standard framework, optimization over sev-

eral normal situations is equivalent to optimization over each situation separately because all variables are specific to each situation. However, in our model, as line capacities are chosen once for all situations, separate optimization on each of them is impossible. Besides, multi-situation long term nodal pricing leads to allocate development costs to each situation. As a primary result, if only normal situations are considered, we show the usual result that most of each element development costs are allocated to the situations which saturate it. It backs cost allocation to the peak situation of each line. However, as suggested in [4], taking into account contingencies will probably result in less concentrated allocations. Indeed, the cost allocated to contingency situations derived from a normal situation should be reported on it. Consequently, even if a line is not saturated in a normal situation, it might be saturated in some associated contingency situations, leading to a significant cost allocation to the normal situation.

The paper is organized as follows. After the model presentation in Section 2, we will analyze optimality conditions, focusing on consequences on conditions of transmission capacity development in Section 3. In Section 4, the equivalence between nodal pricing and proportional costing is developed. In Section 5, the first steps towards establishing a sound situation cost allocation are presented through a three node three situation example.

2 Model

2.1 System Representation

The graph underlying the power system network is modeled as a set of nodes $n \in \mathcal{N}$ and a set of edges $e \in \mathcal{E}$.

Edges represent either physical lines or abstraction of them. They are arbitrarily oriented. An edge belongs to the downstream edge set \mathcal{E}_n^- of its upstream node n_e^+ and to the upstream edge set \mathcal{E}_n^+ of its downstream node n_e^- :

$$\begin{cases} e \in \mathcal{E}_n^- \Leftrightarrow n_e^+ = n \\ e \in \mathcal{E}_n^+ \Leftrightarrow n_e^- = n \end{cases} \quad (1)$$

For each edge e , the cost of developing one or several lines of total capacity X_e is $K_e(X_e) \stackrel{\text{def}}{=} k_e X_e$, i.e. $\partial K_e(X_e)/\partial X_e = k_e$. To simplify equations, the load level t_e of a line of capacity X_e transited by a flow Z_e is defined such that:

$$t_e X_e \stackrel{\text{def}}{=} Z_e \quad (2)$$

The availability level $f_e \in [0, 1]$ of the line represents the upper limit of the load level ($|t_e| < f_e$). It can be strictly below 1 due to technical constraints, or even 0 if the line is out of order.

The loss coefficient α_e is defined as the proportion of power lost if the line is saturated ($|t_e| = 1$). The existence of this loss coefficient independent of the capacity requires that the resistance is inversely proportional to it. Assuming a fixed phase shift, Pouillet's law asserts that the resistance is inversely proportional to cable section, while the capacity is proportional to it. These two relationships allow to define this independent loss coefficient.

A set of producers $g \in \mathcal{G}_n$ and a set of consumers $d \in \mathcal{D}_n$ are attached to each node n . For each producer, the cost $c_g^g(G_g)$ of producing a given quantity of energy G_g is defined as a convex function without fixed costs ($c_g^g(0) = 0$). The production is limited by Y_g . Equivalently, the demand D_d and the cost $c_d^u(U_d)$ of unserved demand U_d are given for each consumer. The set of producers and consumers are respectively noted $\mathcal{G} \stackrel{\text{def}}{=} \bigcup_n \mathcal{G}_n$ and $\mathcal{D} \stackrel{\text{def}}{=} \bigcup_n \mathcal{D}_n$. Finally, the net injection at each node is defined as:

$$I_n \stackrel{\text{def}}{=} \sum_{g \in \mathcal{G}_n} G_g - \sum_{d \in \mathcal{D}_n} (D_d - U_d) \quad (3)$$

2.2 Objective function

Given demands D_d , the objective function is the sum of production and transmission costs and of the cost of unserved energy:

$$\sum_{e \in \mathcal{E}} k_e X_e + \sum_{d \in \mathcal{D}} c_d^u(U_d) + \sum_{g \in \mathcal{G}} c_g^g(G_g) \quad (4)$$

The minimization of the criterium with respect to X_e , t_e , U_d and G_g is submitted to the following constraints, defined along with their corresponding dual variables (Greek letters are used to denote them).

Energy conservation: Given the assumption on losses, they are equal to $\alpha_e t_e^2 X_e$. Assuming losses are split between each end node, the energy balance that should be satisfied at each node n is:

$$- \sum_{e \in \mathcal{E}_n^-} X_e t_e \left(1 + \frac{\alpha_e}{2} t_e\right) - \sum_{e \in \mathcal{E}_n^+} X_e t_e \left(1 - \frac{\alpha_e}{2} t_e\right) - I_n = 0 \leftarrow \lambda_n \quad (5)$$

The voltage law is not taken into account due to the non-convexity implied by its coupling with capacity development [10]. The model is therefore a flow model as in [11], except that it includes losses.

Load level: The available capacity of each line should not be exceeded.

$$t_e^2 \leq f_e^2 \quad \leftarrow \theta_e \quad (6)$$

Capacity: The capacity of each line is positive.

$$X_e \geq 0 \quad \leftarrow \xi_e \quad (7)$$

Unserved demand: The unserved demand of each consumer is positive and limited by its demand.

$$0 \leq U_d \leq D_d \quad \leftarrow \underline{\delta}_d, \bar{\delta}_d \quad (8)$$

Generation: The generation of each producer is positive and limited by its capacity.

$$0 \leq G_g \leq Y_g \quad \leftarrow \underline{\gamma}_g, \bar{\gamma}_g \quad (9)$$

Overall, the Lagrangian of the system is:

$$\begin{aligned} \mathcal{L} \left(X_e, t_e, U_d, G_g, \lambda_n, \theta_e, \xi_e, \underline{\gamma}_g, \bar{\gamma}_g, \bar{\delta}_d, \underline{\delta}_d \right) = & \quad (10) \\ & \sum_{e \in \mathcal{E}} k_e X_e + \sum_{d \in \mathcal{D}} c_d^u(U_d) + \sum_{g \in \mathcal{G}} c_g^g(G_g) \\ & + \sum_{n \in \mathcal{N}} \lambda_n \left(- \sum_{e \in \mathcal{E}_n^-} X_e t_e \left(1 + \frac{\alpha_e}{2} t_e \right) \right. \\ & \quad \left. - \sum_{e \in \mathcal{E}_n^+} X_e t_e \left(1 - \frac{\alpha_e}{2} t_e \right) - I_n \right) \\ & + \sum_{e \in \mathcal{E}} \theta_e (t_e^2 - f_e^2) - \sum_{e \in \mathcal{E}} \xi_e X_e \\ & + \sum_{d \in \mathcal{D}} \bar{\delta}_d (U_d - D_d) - \sum_{d \in \mathcal{D}} \underline{\delta}_d U_d \\ & + \sum_{g \in \mathcal{G}} \bar{\gamma}_g (G_g - Y_g) - \sum_{g \in \mathcal{G}} \underline{\gamma}_g G_g \end{aligned}$$

2.3 Multi-situation extension

Networks are dimensioned to face different operating conditions, resulting from hourly/seasonal demand variations and contingencies. The following extension of the previous model allows to take into account both types of situations: Demands D_d^s , unserved demand costs $c_d^{u,s}(U_d^s)$, production limits Y_g^s , production costs $c_g^{g,s}(G_g^s)$, and availability levels f_e^s are defined for each situation $s \in \mathcal{S}$. Consequently, load levels t_e^s , productions G_g^s and unserved demand U_d^s are also indexed by the situations, as are associated constraints and dual variables. On the contrary, development costs k_e , developed capacities X_e , and positive capacity constraint dual variables ξ_e are the same for all situations. Each situation is assigned a probability p^s representing its expected duration during the optimization period. The modified objective has the following shape :

$$\sum_{e \in \mathcal{E}} k_e X_e + \sum_{s \in \mathcal{S}} p^s \left(\sum_{d \in \mathcal{D}} c_d^{u,s}(U_d^s) + \sum_{g \in \mathcal{G}} c_g^{g,s}(G_g^s) \right) \quad (11)$$

3 Optimality condition analysis

3.1 Mono-situation optimality analysis

Let us introduce the following notations:

$$\begin{cases} \Delta \lambda_e &= \lambda_{n_e^-} - \lambda_{n_e^+} \\ \lambda_e &= (\lambda_{n_e^+} + \lambda_{n_e^-})/2 \end{cases} \quad (12)$$

Given Equations (1) and (10), the associated first-order K.K.T. (Karush–Kuhn–Tucker) necessary optimality conditions are:

$$\frac{\partial \mathcal{L}}{\partial X_e} = 0 = -\Delta \lambda_e t_e + \lambda_e \alpha_e t_e^2 + k_e - \xi_e \quad (13)$$

$$\frac{\partial \mathcal{L}}{\partial t_e} = 0 = X_e (-\Delta \lambda_e + 2\lambda_e \alpha_e t_e) + 2\theta_e t_e \quad (14)$$

$$\frac{\partial \mathcal{L}}{\partial U_d} = 0 = \frac{\partial c_d^u}{\partial U_d}(U_d) + \bar{\delta}_d - \underline{\delta}_d - \lambda_n \quad (15)$$

$$\frac{\partial \mathcal{L}}{\partial G_g} = 0 = \frac{\partial c_g^g}{\partial G_g}(G_g) - \underline{\gamma}_g + \bar{\gamma}_g - \lambda_n \quad (16)$$

with the complementary slackness conditions:

$$\begin{cases} \theta_e (t_e^2 - f_e^2) &= 0 \\ \xi_e X_e &= 0 \\ \bar{\delta}_d (U_d - D_d) &= \underline{\delta}_d U_d = 0 \\ \bar{\gamma}_g (G_g - Y_g) &= \underline{\gamma}_g G_g = 0 \end{cases} \quad (17)$$

and the positivity of $\theta_e, \xi_e, \bar{\delta}_d, \underline{\delta}_d, \bar{\gamma}_g$ and $\underline{\gamma}_g$.

3.1.1 Node analysis

From Equation (16), for an active unsaturated producer ($G_g \in]0, Y_g[$):

$$\frac{\partial c_g^g}{\partial G_g}(G_g) = \lambda_n \quad (18)$$

Therefore λ_n is equal to the marginal cost of producing energy at node n . If the λ_n is below the initial marginal production cost ($\lambda_n < \frac{\partial c_g^g}{\partial G_g}(0)$), then the producer is inactive ($G_g = 0$). If the production reached its maximum ($G_g = Y_g$), λ_n is above this marginal cost ($\lambda_n \geq \frac{\partial c_g^g}{\partial G_g}(Y_g)$). For a more comprehensive discussion of optimality for producers and consumers, see [3]. It hints at interpreting λ_n as the energy price at node n .

3.1.2 Edge analysis

Equations (13) and (14) allow to discuss whether a capacity X_e is developed on edge e , how it is loaded (t_e) and how its development cost k_e is related to energy price difference $\Delta \lambda_e$ and mean energy price λ_e .

If a capacity is installed ($X_e > 0$), the expression of the relationship between k_e and X_e depends on whether the development cost k_e is higher than the cost of losses at saturation $\lambda_e \alpha_e f_e^2$ or not. If it is, the line is saturated, i.e. its load level is maximum ($|t_e| = f_e$). It may be pointed out that saturations resulting from this modeling are different from short-term congestions since the line capacities are always optimal. The relationship between $\Delta \lambda_e$ and k_e is linear:

$$\begin{cases} k_e = |\Delta \lambda_e| f_e - \lambda_e \alpha_e f_e^2 \\ t_e = \text{sgn}(\Delta \lambda_e) f_e \end{cases} : k_e \geq \lambda_e \alpha_e f_e^2 \quad (19)$$

If the development cost is strictly lower than the cost of losses, the line is unsaturated ($|t_e| < f_e$). The relationship between $\Delta \lambda_e$ and k_e is quadratic:

$$\begin{cases} k_e = \Delta \lambda_e^2 / 4\lambda_e \alpha_e \\ t_e = \Delta \lambda_e / 2\lambda_e \alpha_e \end{cases} : k_e < \lambda_e \alpha_e f_e^2 \quad (20)$$

In both cases, the load level t_e and the energy price difference $\Delta \lambda_e$ are oriented in the same direction, which means that the energy is transmitted from low price nodes to high price nodes:

$$t_e \Delta \lambda_e \geq 0 \quad (21)$$

If $X_e = 0$, the first order conditions are degenerated because Equation (14) is easily satisfied. Summarizing

Equations (19) and (20), let us define the development cost at which a capacity can be developed:

$$K_e(\Delta\lambda_e, \lambda_e) = \quad (22)$$

$$\begin{cases} |\Delta\lambda_e|f_e - \lambda_e\alpha_e f_e^2 & : |\Delta\lambda_e| \geq 2\lambda_e\alpha_e f_e \\ \Delta\lambda_e^2/4\lambda_e\alpha_e & : |\Delta\lambda_e| < 2\lambda_e\alpha_e f_e \end{cases} \quad (23)$$

The study of the stability of solution shows that the development cost is higher than this cost :

$$X_e = 0 \implies k_e \geq K_e(\Delta\lambda_e, \lambda_e) \quad (24)$$

Finally, at optimum, the development cost k_e cannot be lower than $K_e(\Delta\lambda_e, \lambda_e)$. The graph of Figure (1) sums up this discussion by representing the relationship in the $(\Delta\lambda_e, k_e)$ plane.

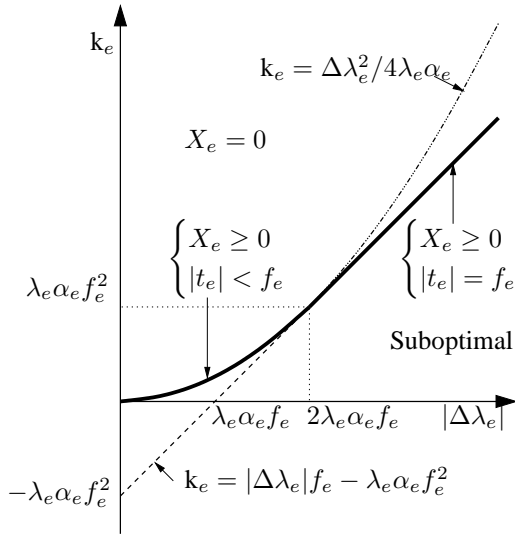


Figure 1: Existence of capacity at optimality given the nodal price difference $|\Delta\lambda_e|$ and the unit capacity development cost k_e . The black thick curve represents the relationship $K(\Delta\lambda_e, \lambda_e)$ existing when capacities are developed ($X_e > 0$) between the development cost k_e , the energy price difference between line extremities $\Delta\lambda_e$, and the mean energy price λ_e . If the development cost k_e is lower than the cost of losses at saturation $\lambda_e\alpha_e f_e^2$, the developed capacities are unsaturated (lower part of the graph, parabolic part of the thick curve), else they are saturated (upper part of the graph, linear part of the curve). Above the thick curve, no capacities are developed because they are too expensive. Below this curve, the situation is suboptimal.

A direct consequence of the previous discussion, particularly of the existence of a development cost $K_e(\Delta\lambda_e, \lambda_e)$ at which a capacity can be developed as a function of dual variables at edge extremities, is the following theorem:

Theorem 1. *Excluding degeneration in parameters, the optimal network is a tree. If parameters are degenerated, there is always a tree in the set of solutions.*

Without losses, the theorem is a classical result of linear programming [12]. Introduction of losses does not change the result because the model defines the resistance as inversely proportional to the capacity.

3.1.3 Price and cost analysis

Classically, the dual variable λ_n is interpreted as the price of energy at node n obtained as a market equilibrium in which transmission of energy is charged at marginal price [1]. It casts a new light on the previous discussion on edge existence: Let us assume that the capacity owner of edge e is able to buy or sell energy at the price $\lambda_{n_e^+}$ and $\lambda_{n_e^-}$ respectively to each end node. It can transmit energy provided it purchases half of its losses at each end node. This paragraph discusses the cost and the revenue associated to a small variation around the optimum of the transmitted energy. It shows how it implies that the total revenue of transmission compensates not only losses but also development costs, even if the line is unsaturated.

First of all, due to Equation (21) and to the arbitrary orientation of edges, we assume without loss of generality that $\Delta\lambda_e \geq 0$ and $t_e \geq 0$. At optimum, the opportunity to transmit a small additional quantity of energy δZ_e is null. Let us assume that the edge is unsaturated. This energy could have been transmitted either by increasing the edge capacity X_e or the load level t_e .

In the first case ($\delta Z_e = t_e \delta X_e$), the opportunity is null because the transmission revenue $\Delta\lambda_e \delta Z_e$ is exactly compensated by two costs: (a) The development cost increase $k_e \delta X_e$; (b) The cost of losses due to the additional flow $\lambda_e \alpha_e t_e \delta Z_e$. As the load level is constant, the mean cost of losses remains $\lambda_e \alpha_e t_e$.

On the contrary, in the second case ($\delta Z_e = X_e \delta t_e$), the opportunity is null because transmission revenue $\Delta\lambda_e \delta Z_e$ is exactly compensated by two costs: (a) The increase of the cost of losses for the original flow $\lambda_e \alpha_e Z_e \delta t_e$, due to the increase of the load level; (b) The cost of losses due to the additional flow $\lambda_e \alpha_e t_e \delta Z_e$. As already noted in [4], due to the quadratic nature of losses, both costs are coincidentally equal: $\lambda_e \alpha_e Z_e \delta t_e = \lambda_e \alpha_e t_e \delta Z_e$. This discussion is summed up in the two following equations. The first one, derived from Equation (13), is related to an increase of the capacity. The second one, derived from Equation (14), is related to an increase of load level:

$$\Delta\lambda_e \delta Z_e = \Delta\lambda_e t_e \delta X_e = k_e \delta X_e + \lambda_e \alpha_e t_e \delta Z_e \quad (25)$$

$$\Delta\lambda_e \delta Z_e = \Delta\lambda_e X_e \delta t_e = \lambda_e \alpha_e Z_e \delta t_e + \lambda_e \alpha_e t_e \delta Z_e \quad (26)$$

As development costs are linear ($K_e(X_e) \stackrel{\text{def}}{=} k_e X_e$), it is possible to integrate previous equations in order to obtain the financial balance associated to edge e :

$$\begin{cases} \Delta\lambda_e Z_e = k_e X_e + \alpha_e t_e \lambda_e Z_e \\ \Delta\lambda_e Z_e = 2\alpha_e t_e \lambda_e Z_e \end{cases} \quad (27)$$

Consistently with the opportunity analysis, the revenue $\Delta\lambda_e Z_e$ is equally split into development costs $k_e X_e$ and loss purchase $\alpha_e t_e \lambda_e Z_e$. Indeed, assuming that the capacity is given, the energy is transmitted at marginal cost $2\alpha_e t_e \lambda_e$, whereas the first units transmitted generate less losses. The difference is allocated to development costs.

Finally, let us consider a saturated edge ($|t_e| = f_e$). A similar financial analysis leads to:

$$\begin{cases} \Delta\lambda_e Z_e = k_e X_e + \alpha_e f_e \lambda_e Z_e \\ \Delta\lambda_e Z_e = 2\alpha_e f_e \lambda_e Z_e + 2\theta_e \end{cases} \quad (28)$$

In the saturated as in the unsaturated case, the energy transmission does not generate any additional economic rent because the development costs are linear and developable capacities unlimited. Overall, this paragraph sums up in the following theorem:

Theorem 2. *In the long term nodal pricing model, the revenue of transmission recovers exactly the development costs and the cost of losses. The development cost share is at least equal to cost of losses one. It is equal if the capacity is unsaturated and higher if it is.*

3.2 Multi-situation optimality analysis

With the multi-situation extension, the optimality conditions are similar to the previous ones, except the derivative with respect to the edge capacity:

$$\frac{\partial \mathcal{L}}{\partial X_e} = 0 = k_e - \xi_e + \sum_{s \in \mathcal{S}} \left(-\Delta \lambda_e^s t_e^s + \lambda_e^s \alpha_e t_e^{s2} \right) \quad (29)$$

Using Equation (22) which defines the nodal price-development cost relationship in the mono-situation case, let us define k_e^s and rewrite Equation (29) when $X_e > 0$:

$$\begin{aligned} k_e^s &\stackrel{\text{def}}{=} K_e(\Delta \lambda_e^s, \lambda_e^s) \\ k_e &= \sum_{s \in \mathcal{S}} k_e^s \end{aligned} \quad (30)$$

This cost decomposition allows to assess the value of each edge developed capacities in each situation. As in the mono-situation case, if the edge is unsaturated, the development costs allocated to the situation are equal to the cost of losses. It means that every used edge ($Z_e^s > 0$) has a value, even if it is unsaturated. It is false in the short term nodal price framework in which nodal price difference covers only losses on unsaturated edges. The fraction of development costs allocated to a situation cannot exceed the cost of losses unless the edge is saturated. In usual networks, the cost of losses represent only a fraction of the development costs (around one fifth). As a result, most of development costs are allocated to the saturating situation or to saturating situations.

The topology of the network is not a tree anymore, even if availability constraints are not used ($f_e^s = 1$). This fact is illustrated by the example developed in Section 5.

4 Nodal pricing and proportional costing equivalence

Despite many contributions, the question of counter flow pricing in transaction based power flow costing methods remains open. In this section, we show the original result that the long term nodal pricing is equivalent to a transaction based costing method in which counter flows get credit for lightening the flow. This method is based on the proportional costing principle, i.e. transactions compensate the development costs and the losses of each edge proportionally to their usage of it (Theorem 3). Moreover development costs are exactly recovered on the optimal network.

Let us assume that, for a given network, a node potential λ_n representing a price of energy at node n is given. Let us assume that a loss function $L_e(Z_e)$ satisfying $L_e(0) = 0$ is given for each edge. In the previously exposed model, $L_e(Z_e) = \alpha_e Z_e^2 / X_e$. Let us also define the mean loss factor as :

$$\ell_{Z_e} = \begin{cases} L_e(Z_e)/Z_e & : Z_e \neq 0 \\ 0 & : Z_e = 0 \end{cases} \quad (31)$$

Let us define a network usage u through sets of productions G_g^u , demands D_d^u , lost loads U_d^u and the corresponding set of net injections I_n^u , partially loading lines of the network at level Z_e^u such that:

$$\begin{cases} I_n^u \stackrel{\text{def}}{=} - \sum_{e \in \mathcal{E}_n^-} Z_e^u + \ell_{Z_e} Z_e^u / 2 \\ X_e = 0 \implies Z_e^u = 0 \end{cases} \quad (32)$$

If the network is a tree, there is only one network usage set Z_e^u per injection set I_n^u . On the contrary, if the network has cycles, there can be several usage sets for one injection set, especially if losses are not taken into account. Injection sets correspond to usual transactions, including losses at each node. Due to this definition, the revenue from a nodal tariff of network usage u can be rewritten in the form of edge compensations:

$$\begin{aligned} \sum_{g \in \mathcal{G}} -\lambda_n G_g^u + \sum_{d \in \mathcal{D}} \lambda_n (D_d^u - U_d^u) \\ = \sum_{e \in \mathcal{E}} Z_e^u (\Delta \lambda_e - \ell_{Z_e} \lambda_e) \end{aligned} \quad (33)$$

This equivalence is valid for every node potential, every network topology, and every loss function satisfying Equation (31).

Additionally, in the long term nodal pricing model, due to Equation (13), this relationship between development cost and energy prices is valid for edges with $X_e > 0$:

$$k_e X_e = Z_e (\Delta \lambda_e - \ell_{Z_e} \lambda_e) \quad (34)$$

Combining Equations (33) and (34) leads to the following equivalence theorem:

Theorem 3. *For any given network usage u , the nodal tariff based on the optimal dual variables λ_n of the long term nodal pricing model is equivalent to a tariff based on the proportional compensation of development costs on the optimal network:*

$$\sum_{g \in \mathcal{G}} -\lambda_n G_g^u + \sum_{d \in \mathcal{D}} \lambda_n (D_d^u - U_d^u) = \sum_{e \in \mathcal{E} | t_e > 0} Z_e^u \frac{k_e}{t_e} \quad (35)$$

In particular, given that $t_e = 0$ and $X_e > 0$ implies $k_e = 0$, the overall financial balance is:

$$\sum_{g \in \mathcal{G}} -\lambda_n G_g + \sum_{d \in \mathcal{D}} \lambda_n (D_d - U_d) = \sum_{e \in \mathcal{E}} k_e X_e \quad (36)$$

(A)

(B)

Edge	X_e	k_e	α_e
e_1	0.51	1	2%
e_2	1.52	1	2%
e_3	1.01	1	2%
Tot.	3.03	-	-

(C)

Node	$p^1 \frac{\partial c_g^{g,1}}{\partial G_n}(G_g^1)$	$p^2 \frac{\partial c_g^{g,2}}{\partial G_n}(G_g^2)$	$p^3 \frac{\partial c_g^{g,3}}{\partial G_n}(G_g^3)$
n_A	10	-	10
n_B	10	10	-
n_C	-	10	-

(D)

Sit.	1				2				3				Tot.
N.	G_g^1	D_d^1	λ_n^1	pay.	G_g^2	D_d^2	λ_n^2	pay.	G_g^3	D_d^3	λ_n^3	pay.	pay.
n_A	1.53	0	10	15.30	0	2.5	10.63	-26.57	1.53	0	10	15.34	4.07
n_B	0.51	0	10	5.10	1.02	0	10	10.19	0	1.5	10.79	-16.18	-0.89
n_C	0	2	10.76	-21.53	1.53	0	10	15.31	0	0	10.14	0	-6.21
Tot.	2.04	2	-	-1.13	2.55	2.5	-	-1.06	1.53	1.5	-	-0.84	-3.03

Sit.	1				2				3				Tot.
E.	Z_e^1	loss.	$\Delta\lambda_e^1$	comp.	Z_e^2	loss.	$\Delta\lambda_e^2$	comp.	Z_e^3	loss.	$\Delta\lambda_e^3$	comp.	comp.
e_1	-0.51	0.01	-0.76	0.28	0	0	0	0	0.51	0.01	0.65	0.22	0.51
e_2	1.52	0.03	0.76	0.84	-1.52	0.03	-0.63	0.64	0.52	0.01	0.14	0.03	1.52
e_3	0	0	0	0	1.01	0.02	0.63	0.42	-1.01	0.02	-0.79	0.59	1.01
Tot.	-	0.04	-	1.13	-	0.05	-	1.06	-	0.03	-	0.84	3.03

Table 1: Three node three situation example. Data is bold characters and results are in standard characters. Figure (A) sketches the network with three nodes n_A , n_B , and n_C and three potential lines e_1 , e_2 , and e_3 . Arrows on edges indicate the arbitrary direction of positive flow. Table (B) presents capacity development results X_e along with development costs k_e and loss coefficient α_e . Table (C) presents active producers and their production costs in each situation. The production cost (10) is chosen so that the cost of losses is below the development costs so that each line is saturated at least once. Table (D) presents detailed situation results associated to nodes in the upper part and to edges in the lower part. The production G_g^s , demand D_d^s , nodal price λ_n^s , and associated TSO payment (pay. = $G_g^s \lambda_n^s$) are given for each node. The flow Z_e^s , losses (loss. = $\alpha_e t_e Z_e^s$), nodal price difference $\Delta\lambda_e$, and associated TSO compensation after loss purchase (pay. = $\Delta\lambda_e Z_e^s - \lambda_e^s \cdot \text{loss}$) for each edge. N. = Node; E. = Edge; Sit. = Situation; Tot. = Total.

In other words, the revenues from a nodal tariff (or from the equivalent proportional costing tariff) are exactly equal to development costs. Indeed, as already mentioned in Theorem 2, there is no additional economic rent for capacity owner. This theorem can be easily extended to the multi-situation extension.

Finally, the reverse of Theorem (3) is true when the network is a tree:

Theorem 4. *The proportional compensation of development costs is equivalent to a nodal tariff if the network is a tree and provided that all edges with capacities and no flow have null development costs (if $X_e > 0$ and $Z_e = 0$, i.e. $t_e = 0$, then $k_e = 0$).*

5 Multi-situation cost allocation

The multi-situation model entangles spatial and temporal aspects so that theoretical analysis have not yet been thoroughly explored. In particular, the underlying principle of edge development cost allocation to each situation is not straightforward. In order to illustrate it, let us con-

sider the three node three situation example described in Table (1-A). Each node n_A , n_B , n_C holds one producer and one consumer. In each situation numbered from 1 to 3, only one consumer is active: The positive demands are D_C^1 , D_A^2 and D_B^3 {Table (1-D, Col. D_d^s)}. So as to forbid local production, producers are disabled ($Y_g^s = 0$) when the corresponding consumer is active {Table (1-D) Col G_g^s }. Furthermore, in situation 3, producer in n_C is also deactivated so that the only active producer is located in n_A . Capacities may be installed between each node on three edges noted e_1 , e_2 , and e_3 . Capacity development costs and losses are uniform on the three edges {Table (1-B) Col. k_e and α_e }. Active producer marginal energy cost is chosen so that each edge with developed capacity will have at least one saturating situation ($p^s \frac{\partial c_g^{g,s}}{\partial G_n}(G_g^s) = 10$) {Table (1-C)}, i.e. the development costs are sufficiently higher than the costs of losses.

The optimization results in developing network capacity on each edge {Table (1-B, Col. X_e)}. The network is not a tree, even if all edges are always fully available.

The sub graph used in each situation is not necessarily a tree either: In the last situation, all edges are used to send energy from n_A to n_B {Table (1-D, Col. Z_e^3)}. In the two other situations, the edge linking producing nodes is not used ($Z_3^1 = Z_1^2 = 0$). Besides, the analysis of situation development cost allocation confirms the theoretical analysis and shows that usual cost allocation techniques are not optimal with respect to the current model. First of all, as expected, development costs are exactly compensated by the nodal price tariff {Tables (1-D, Col. *Tot. comp.*) and (1-B Line *Tot.*)}. All three edges are saturated in two of the three situations. As expected, for each edge, most development costs are allocated to these saturated situations and no costs are allocated to unused edges {Table (1-D, Col. *comp.*)}.

However, the allocation of costs between saturated situations is not related to usual allocation principles. For example, let us consider three allocation principles: Uniform (U); Proportional to the total demand in MW (D); Proportional to the network load in MW.km, assuming that all edges have the same length (L); According to the current model (N). Table (2) presents the resulting cost allocation results and confirms that all schemes are different.

Sit. \ Meth.	U	D	L	N
1	1.01	1.01	0.93	1.13
2	1.01	1.26	1.16	1.06
3	1.01	0.76	0.94	0.84
Tot.	3.03	3.03	3.03	3.03

Table 2: Development cost allocation to situation 1,2 and 3 using different methods. U = Uniform; D = Proportional to the total demand in MW; L = Proportional to the network load in MW.km assuming that all edges have the same length; N = Long term nodal pricing.

Surprisingly, the current model does not even allocate costs monotonously with the total demand in MW or with the network load in MW.km: costs allocated to situation 2 are lower than costs allocated to situation 1. This is linked with the fact that node n_C cannot produce in situation 3. Indeed, if this constraint is relaxed, the cost allocation resulting from the model (N) is equivalent to the cost allocation proportional to the total demand (D). Applying the constraint rises the compensation of edge e_3 in situation 3 because it directly links the producer node n_A to the consumer node n_B . As a counter effect, it relieves the compensation of the same edge in situation 2, thus lowering the global situation cost allocation.

6 Conclusion

In this article, we show that, provided a long term multi-situation framework is used, nodal pricing theory gives a useful theoretical perspective on the variety of existing costing principles. Inside this perspective, we show that a tariff based on long term nodal prices is equivalent to a tariff based on the proportional compensation of development costs on the optimal network. We also show that the development costs are mostly allocated to the saturating situations. However, among saturating situations, the allocation of costs does not seem to follow simple rules.

The remaining steps before a full implementation of the model are related to security modeling and multiple saturating situations. Indeed, the network is dimensioned to face contingency situations, derived from normal operating situation. A reasonable scheme consists in transferring to the normal situation the costs allocated to derived unobserved contingency situations. In normal situations, all demand is served, so that it is unlikely that an edge is saturated more than once. However, in contingency situations, due to their low probability, there may be unserved demand. As a result, multiple saturating situations might frequently occur. The analysis of cost allocation in such cases is therefore the main milestone towards implementation.

References

- [1] J. Boucher and Y. Smeers. Alternative models of restructured electricity systems, part 1: No market power. *Operations Research*, 49:821–838, 2001.
- [2] F.C. Schweppe, M.C. Caramanis, R.D. Tabors, and R.E. Bohn. *Spot Pricing of Electricity*. Kluwer Academic Publishers, London, 1988.
- [3] M. Hsu. An introduction to the pricing of electric power transmission. *Utilities Policy*, 6:257–270, 1997.
- [4] I.J. Perez-Arriaga, F.J. Rubio, J.F. Puerta, J. Arceluz, and J. Marin. Marginal pricing of transmission services: an analysis of cost recovery. *IEEE Transactions on Power Systems*, 10:546–553, 1995.
- [5] W.W. Hogan. Contract networks for electric power transmission. *Journal of Regulatory Economics*, 4:211–242, 1992.
- [6] Z. Jing, X. Duan, F. Wen, Y. Ni, and F.F. Wu. Review of transmission fixed costs allocation methods. *IEEE Power Engineering Society General Meeting*, 4:2585–2592, 2003.
- [7] M. Boiteux and P. Stasi. Sur la détermination des prix de revient de développement dans un système interconnecté de production–distribution. *Document VI.10 du Congrès de Rome de l’U.N.I.P.E.D.E.*, reprinted in G. Morlat et F. Bessière (eds.) *Vingt-cinq ans d’économie électrique*, Dunod, Paris, 361–400, 1952.
- [8] A.S.D. Braga and J.T. Saraiva. Long term marginal prices – solving the revenue reconciliation problem of transmission providers. In *Proceedings of the 15th Power Systems Computation Conference*, 2005.
- [9] F. Bouffard, F.D. Galiana, and J.M. Arroyo. Umbrella contingencies in security-constrained optimal power flow. In *Proceedings of the 15th Power System Computation Conference*, 2005.
- [10] L. Bahiense, G.C. Oliveira, M. Pereira, and S. Granville. A mixed integer disjunctive model for transmission network expansion. *IEEE Transactions on Power Systems*, 16:560–565, 2001.
- [11] R. Villasana, L. L. Garver, and S.J. Salon. Transmission network planning using linear programming. *IEEE Transactions on Power Apparatus and Systems*, 104:349–356, 1985.
- [12] R.K. Ahuja, T.L. Magnanti, and J.B. Orlin. *Network Flows. Theory, Algorithms, and Applications*. Prentice Hall, 1993.

9.3 Conclusion

Cet article, bien que laissant un grand nombre de questions ouvertes, met en lumière plusieurs points :

9.3.1 Situations possibles et réalisées

L'établissement d'une tarification des réseaux électriques fait face à la difficulté suivante : à chaque instant, le système est configuré pour faire face à un grand nombre d'aléas, et pourtant un seul se réalise. Le fait de pouvoir faire face à cet ensemble d'aléas a un coût (A court terme, c'est la non utilisation de la capacité maximum des ouvrages pour faire face à des situations « N-1 ». A moyen terme, c'est l'anticipation d'hiver rigoureux. A long terme, c'est l'incertitude sur le niveau de consommation et les coûts de l'énergie). De nouveau apparaît la complexité du service rendu par la fourniture d'électricité : à côté du coût propre de l'énergie (ou de son prix, plus exactement dans un système de marché), se juxtapose le coût de la sécurité de l'approvisionnement. Par exemple, le système est conçu pour faire face à des aléas extrêmes en termes de charge du réseau, néanmoins le système peut apparaître comme surdimensionné si ceux-ci ne se réalisent pas.

Le modèle de long terme que nous avons exposé inclut ces coûts d'aléas non réalisés, cependant il ne résout pas le problème de leur intégration dans une grille tarifaire pratique (par exemple sous la forme classique d'une prime fixe de puissance maximale et d'une prime proportionnelle d'énergie consommée).

9.3.2 Congestion et saturation

Le modèle permet de distinguer la congestion d'une ligne de sa simple saturation. Dans les deux cas, la ligne est utilisée à sa capacité maximale. Dans le cas d'une congestion, le surplus social aurait été plus grand si la ligne avait été de capacité plus importante. Dans le cas d'une saturation, la capacité de la ligne n'est pas supérieure car le surplus social se trouverait diminué par les coûts de développement de capacité supplémentaire. Seul le second cas se rencontre dans le modèle. Ceci permet de préciser l'objectif d'un marché européen de l'électricité : il est vraisemblable que, dans un système européen maximisant le surplus social, il demeure ponctuellement des saturations entraînant des différences de prix d'électricité entre régions européennes. Des capacités de transport plus grandes coûteraient plus qu'elles ne seraient utiles.

Ainsi, il n'est pas nécessairement efficaces de développer le système jusqu'à ce que les capacités de transport jusqu'à ce que n'importe quel consommateur européen puisse choisir n'importe quel producteur européen. En effet, il se peut que les capacités de transport optimales ne permettent pas cette configuration.

Chapitre 10

Conclusion

En observant le modèle théorique exposé, on constate qu'il s'agit d'un modèle sur lequel une décomposition par les prix peut être pratiquée – cf. Section (2.2.1.1), même s'il n'est pas nécessairement convexe à cause des pertes (elles sont proportionnelles au produit de deux variables, ce qui n'est pas le cas dans un modèle de court-terme où les capacités des lignes sont fixées). On peut ainsi décomposer le problème en un sous problème d'optimisation par producteur, un par ligne et un par consommateur. Il ne reflète donc pas directement la complexité des systèmes électriques annoncée dans l'introduction.

Toutefois, une analyse précise montre que l'existence même d'un grand nombre de situations dimensionnantes pour le système, dont seul un petit nombre se réalise, demande à l'ensemble des acteurs de se mettre d'accord sur les situations auxquelles le système doit répondre, sachant qu'il ne s'agit pas uniquement des situations réalisées. A court terme, cela signifie une coordination accrue des acteurs afin d'être certain que le système puisse répondre à un ensemble de pannes données, ce qui nécessite pratiquement l'exposition du détail des actifs de chaque acteur, probabilité de défaillance incluse. A long terme, la construction d'un marché sur la base de cette décomposition du problème de maximisation du surplus social entre producteurs, transporteurs et consommateurs semble compromise par la longueur et la complexité des processus d'investissement, même si la publication d'indicateurs tarifaires de long terme dont nous ébauchons la construction pourrait faciliter un tel processus.

En revanche, autant la modularité apparaît faible pour le transport, ce que l'Union Européenne a reconnu en créant des monopoles régulés de transport, autant la relative modularité des producteurs et surtout des consommateurs semble réelle. Traditionnellement et du fait de limitations technologiques, on a attribué à ces derniers une courbe de demande particulièrement simple : ils seraient prêts à payer un prix donné très élevé jusqu'à un certain niveau de consommation après lequel l'utilité marginale deviendrait immédiatement nulle. La technologie des « compteurs intelligents » s'alliant à la libérali-

sation, on verra peut-être apparaître des solutions pratiques permettant au consommateur de montrer la réalité de son élasticité au prix en choisissant des fournisseurs aux offres innovantes. Il s'agit bien là de la force attendue des systèmes modulaires évolutifs : leur permettre de s'adapter aux changements (nouvelles technologies, raréfaction de l'énergie fossile, etc.) malgré une connaissance très partielle de la fonction à optimiser, en l'occurrence le surplus social.

Quatrième partie

Discussion et conclusion

Chapitre 11

Discussion

De manière additionnelle aux discussions spécifiques à chaque étude développée dans cette thèse, ce chapitre vise à discuter de la pertinence des concepts de modularité et de convergence structure-fonction d'après leur utilité dans chacune de ces études. Ceci permet tout à la fois de mettre en lumière leurs apports dans chaque étude et les limites à leurs applications pratiques et, à partir de ces constats, de dégager des perspectives d'approfondissement théoriques et pratiques de ces concepts.

11.1 Pertinence du concept de modularité

Dès l'introduction, nous avons constaté que la définition même de la modularité était extrêmement délicate.

D'abord, la modularité n'a pas pu être définie comme la propriété d'un système, mais comme la propriété d'un partitionnement d'un système. La notion de modularité d'une distribution de probabilité est donc mal définie. Malgré cela, les résultats sur la convergence structure-fonction – cf. Section (1.3.3) – portent sur de telles distributions et nous avons alors extrapolé en parlant de superposition des modularités structurelles et fonctionnelles, quelle que soit la manière exacte de les définir. En effet, puisqu'il y a superposition des distributions de probabilité, il y a aussi correspondance des partitionnements modulaires, quels qu'ils soient. Toutefois, il reste un travail théorique important pour préciser le concept de modularité, probablement en lien avec celui de complexité sur lequel il repose. Ceci renforcera l'utilité de l'étude de la convergence structure-fonction.

Ensuite, la modularité telle qu'elle est pratiquée fait intervenir un compromis entre taille des modules et faiblesse des connexions. En effet, en l'absence de compromis, le meilleur partitionnement est toujours celui qui laisse entier le système – cf. Section (1.1.2.2). Ce compromis reste à étudier de manière théorique, même si les difficultés algorithmiques liées à la recherche de modules sont déjà nombreuses et masquent souvent cette question. Le travail en ce domaine est à la fois théorique (pour

étudier l'impact des compromis) et pratique (pour étudier leur réalisme dans différents domaines). Par exemple, plutôt que la pénalisation quadratique de la complexité utilisée, on pourrait imaginer l'impact de placer une limite inférieure stricte en termes de complexité à la taille des modules, ce qui correspondrait par exemple à des contraintes de longueur minimale d'un gène ou bien de taille minimale d'une entreprise.

Au final, la modularité n'est pas le concept creux annoncé par certains, mais son utilisation pratique pour la recherche de modules fait face à de nombreuses difficultés théoriques et pratiques. Fort heureusement, il existe un certain nombre de cas où la modularité tant structurelle que fonctionnelle des systèmes est suffisamment claire (gènes, *opérons* et acteurs économiques), ce qui permet de dépasser ces difficultés et de poursuivre l'analyse de ces systèmes.

11.2 Pertinence du concept de convergence structure-fonction

A travers une analyse en optimisation mathématique, en économie et en biologie, l'introduction de cette thèse a montré la pertinence du concept de convergence structure-fonction. Ceci implique en particulier que les modularités structurelles et fonctionnelles se superposent dans les *systèmes évolutifs et fonctionnels*. Les résultats de cette thèse en biologie et en économie ont permis de vérifier l'utilité de ce concept, tout en percevant des limites à ses applications pratiques.

Ainsi, si les *opérons* sont bien des unités structurelles calquées sur une modularité fonctionnelle, nos résultats suggèrent que leur longueur est soumise à sélection du fait de limitations des mécanismes de *réplication* et de *transcription*. En quelque sorte, ces limitations forcent le trait de la modularité structurelle, posant ainsi une limite à la convergence structure-fonction.

D'autre part, nos travaux montrent que l'utilisation de la seule modularité fonctionnelle des gènes pour analyser fonctionnellement des études d'associations génétiques, bien que donnant des résultats encourageants, devrait être complétée d'informations structurelles sur les *points chauds de recombinaison* pour être réellement efficace. En effet, dans le cas des populations humaines, la force de la *convergence* n'est pas suffisante pour assurer la superposition des *points chauds* et des limites de gènes.

Enfin, le modèle économique que nous avons utilisé pour représenter les systèmes électriques permet de constater que la décomposition structurelle en acteurs souhaitée par la commission européenne (producteurs, transporteurs et consommateurs) correspond bien à une modularité fonctionnelle de l'objectif retenu dans la modélisation. Toutefois, l'analyse précise des relations induites par le découpage met en lumière certaines difficultés de coordination entre les acteurs. Ceci illustre que le principe de *convergence* ne se réduit pas à la superposition des modularités : le comportement des acteurs et leurs relations sont essentiels afin de garantir la meilleure superposition des probabilités de transition et de la fonction objectif.

Chapitre 12

Conclusion

On pense souvent que, plus un concept est général, moins il est applicable. Il est vrai qu'il s'abstrait alors d'un grand nombre de caractéristiques propres à l'ensemble des problèmes qu'il tente de généraliser au point de rendre triviale son application à un système pratique. C'est certainement un des reproches que l'on peut faire à la modularité, un concept si général que l'on peine à le définir, à le mesurer et encore plus à trouver des propriétés générales des systèmes modulaires même si la découverte de la structure « hub-cluster » de la plupart des systèmes réels est certainement une étape importante dans la compréhension de l'unité des mécanismes de gestion de la complexité .

Toutefois, la généralité du phénomène convergence structure-fonction dans les *systèmes évolutifs et fonctionnels* laisse à penser que son étude approfondie est nécessaire à la compréhension de ce qui pourrait faire leur unité. Elle possède une caractérisation mathématique précise par le biais de la fonction de fitness effective, même si difficilement applicable à des systèmes réels du fait de la taille incommensurable de l'espace de leurs évolutions possibles. Grâce à cette caractérisation mathématique, il est possible d'utiliser cette *convergence* comme cadre d'interprétation de réalités très variées afin de constater à quelle point celle-ci est réalisée ou non dans les systèmes d'études.

Dans cette thèse, nous avons privilégié l'étude de trois systèmes en utilisant ce cadre conceptuel pour obtenir des résultats pratiques :

- Nous avons montré qu'il existait une pression de sélection sur la longueur des *opérons* bactériens afin de diminuer l'impact des collisions des *réplicases* et des *transcriptases*. Cette pression pousse à ce qu'il existe plusieurs modules structurels (*opérons*) même lorsque la fonction qu'ils assurent est unique (par exemple formation d'un complexe).
- Nous avons aussi illustré les difficultés liées à l'existence d'une modularité structurelle due aux *points chauds de recombinaison* qui ne se superpose pas parfaitement à la modularité fonctionnelle associée aux gènes pour la réalisation d'études d'associations génétiques chez l'homme. En effet, cette *non-convergence* diminue les possibilités de formuler des hypothèses fonction-

nelles sur le mécanisme de la maladie étudiée à partir du constat d'associations structurelles de séquences d'*ADN* avec la maladie.

- Enfin, nous avons proposé une méthodologie pour calculer un tarif de long terme pour le transport de l'électricité. En se rapprochant du résultat d'une décomposition par les prix du problème d'optimisation associé au système, ce tarif permettrait de limiter les *externalités*. En effet, l'optimalité serait atteinte malgré la *non-convergence* entre le découpage structurel du système entre producteurs et transporteurs et sa réalité fonctionnelle qui n'est pas modulaire selon ce même découpage. Par exemple, il est important de coordonner les investissements dans des capacités de production et dans des capacités de transport afin d'éviter de construire des lignes inutiles ou des centrales qui ne peuvent fonctionner par manque de lignes.

Ces études pratiques laissent penser que des études théoriques plus nombreuses des espaces d'évolution des systèmes fonctionnels (190), s'inspirant par exemple de celle réalisée sur le réseau de mutations de l'*ARN* (130; 192) permettront de mieux se représenter les formes de ces espaces, leurs caractéristiques communes et leurs différences. Grâce à cette compréhension approfondie, une nouvelle étape pourra être franchie en *science des systèmes*, dont on peut penser qu'elle irriguera en retour un grand nombre d'études pratiques.

Cinquième partie

Glossaire et bibliographie

Glossaire

- ADN** Acide désoxyribonucléique. Support de l'information héréditaire dans les organismes vivants. L'ADN stable est double-brin. C'est-à-dire qu'il est formé de deux molécules linéaires appelées brins portant une information complémentaire et appariées à la manière d'une fermeture éclair. A partir d'un seul brin, on peut reconstituer le second. L'ADN est une chaîne dont un élément est appelé base. Voir aussi la référence (4). 40, 55, 56, 61–66, 68, 70, 71, 73, 74, 88, 118
- analyse fonctionnelle** Analyse dédiée à la détermination des fonctions d'un produit en cours de conception. Voir aussi la référence (10). 38
- ARN** Acide ribonucléique. Macromolécule de structure semblable à l'ADN mais ne possédant qu'un brin. Elle est utilisée soit directement dans les mécanismes cellulaires (souvent des mécanismes apparus avant ceux assurés par des protéines) soit comme intermédiaire dans la production des protéines à partir de l'ADN (ARN messager). Les ARN messagers sont dégradés en quelques minutes. Voir aussi la référence (4). 5, 61, 65, 69, 70, 118
- clade** En taxinomie, entité conceptuelle regroupant tous les organismes vivants descendant du même ancêtre commun. Voir aussi la référence (109). 54
- clique** En théorie des graphes, une clique est un ensemble de nœuds deux à deux adjacents. Autrement dit, toutes les arêtes possibles entre nœuds de la clique existent. Voir aussi la référence (58). 14, 15, 18, 19, 22, 24, 25
- CMH** En immunologie, le Complexe Majeur d'Histocompatibilité désigne à la fois un système de reconnaissance du soi présent chez la plupart des vertébrés et certaines protéines très variables d'un individu à l'autre qui sont impliquées dans cette fonction. 85
- complexité** Mesure conjointe de la taille d'un système et de la force du couplage de ses différentes parties caractérisant la difficulté d'analyse d'un système par l'homme. La complexité d'un système ayant un nombre fini d'états possible et décrit par une distribution de probabilité peut être mesurée par son entropie de Shannon. Les complexités de Kolmogorov et de Chaitin sont ba-

sées sur la longueur du programme de longueur minimale permettant de décrire un objet. 4–11, 13–20, 22, 26, 27, 35, 36, 47, 49, 53, 56, 57, 64, 67, 68, 110, 111, 115–117

concavité Une fonction f est concave si son opposé $-f$ est convexe. 11

concurrence Situation dans laquelle les acteurs d'un système économique ont le choix entre plusieurs alternatives pour satisfaire un même besoin. Voir aussi la référence (37). 49, 98, 99

convergence structure-fonction Dans les systèmes évolutifs et fonctionnels, tendance de la distribution de probabilité des variations structurelles à être calquée sur la variabilité de la fonction objectif. En particulier, la modularité structurelle de ces systèmes tend à être superposable à leur modularité fonctionnelle. vii, viii, 28, 30–33, 35, 37–44, 46, 47, 53, 55, 57, 61, 62, 66–68, 74, 75, 87, 95, 97, 115–118

convexité Une fonction f de \mathbb{X} dans \mathbb{R} est convexe si, quel que soit $\alpha \in [0, 1]$, quel que soit x et y de \mathbb{X} : $f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y)$. La fonction est dite strictement convexe si cette inégalité est strictement satisfaite. 7, 8

degré En théorie des graphes, le degré d'un nœud est le nombre d'arêtes connectées à un nœud. S'il n'existe pas d'arêtes multiples et d'arêtes partant et arrivant au même nœud, c'est aussi le nombre de nœuds voisins du nœud. Dans un graphe orienté, le degré sortant est le nombre d'arêtes partant d'un nœud. Le degré entrant est le nombre d'arêtes arrivant en un nœud. Voir aussi la référence (58). 13, 16–19, 23, 25, 55

dendrogramme Représentation graphique qui comporte des ramifications semblables à celles d'un arbre. 23

diagramme de Venn Un diagramme de Venn est une représentation schématiques d'ensembles. Voir aussi http://fr.wikipedia.org/wiki/Diagramme_de_Venn. 6

diagramme de Voronoï Décomposition particulière d'un espace métrique déterminée par les distances à un ensemble discret d'objets de l'espace, en général un ensemble discret de points. Voir aussi http://fr.wikipedia.org/wiki/Diagramme_de_Voronoï. 46

diamètre En théorie des graphes, le diamètre est la distance maximale séparant deux nœuds d'un graphe non orienté. La distance entre deux nœuds est égale à la longueur en nombre d'arêtes d'un des plus court chemins les reliant. 17, 18, 20, 22, 25

diploïdie Caractéristique des espèces, telle l'espèce humaine, qui possèdent des paires de chromosomes homologues dont un seul est transmis à chaque descendant lors d'une reproduction sexuée. L'existence de paires homologues est à l'origine du phénomène d'enjambement lors de la méiose. 63

déséquilibre de liaison Corrélation de l'information génétique trouvée en deux positions distinctes de l'ADN. Les origines de cette corrélation sont multiples, mais la proximité des positions en

est généralement responsable, du fait de la faible probabilité d'une recombinaison séparant les deux positions. Voir aussi l'introduction de (60). 64, 65, 73, 85, 87

enjambement Recombinaison homologue propre aux organismes diploïdes. Elle se déroule lors de la méiose et consiste en l'échange de séquences d'ADN entre chromosomes homologues. Voir aussi la référence (4). 56, 63

entropie En théorie de l'information, l'entropie de Shannon mesure de la quantité d'information contenue dans une source d'information. En physique statistique, l'entropie mesure du degré de désordre d'un système. Dans les deux cas, pour une variable aléatoire X discrète, elle peut s'écrire $H(X) = -\sum_x p(x) \ln(p(x))$. L'entropie est sous-additive ($\max(H(X), H(Y)) \leq H(\{X, Y\}) \leq H(X) + H(Y)$). 5, 6, 13, 14, 24, 29–31

entropie croisée Terme de théorie de l'information. Soit deux distributions discrètes p et q . Leur entropie croisée est définie par $\sum_x p(x) \ln(q(x))$. 30

enzyme ensemble formé de protéines et/ou d'ARN capable de catalyser (faciliter) une réaction chimique. 63, 69

espèce En taxinomie, l'espèce est l'unité, ou taxon de base. Le concept d'espèce pose de nombreux problèmes dans sa définition et dans son application. Voir aussi la référence (62). 40, 50, 51, 54, 63, 65, 66

eucaryote Taxon regroupant l'ensemble des organismes vivants dont la cellule possède un noyau dans lequel est confiné l'ADN. Voir aussi la référence (109) et <http://tolweb.org/Eukaryotes/3>. 62, 66, 70, 71

externalité Une externalité désigne une situation dans laquelle une activité économique cause des coûts ou des bénéfices à des tiers qui n'ont aucune influence sur son déroulement. Voir aussi Voir aussi la référence (37) . 38–40, 43, 47, 49, 99, 118

fonction de fitness En biologie, modélisation du mécanisme de l'évolution naturelle dans laquelle on considère que le « problème » d'un organisme est d'optimiser une fonction appelée fonction de fitness. Cette fonction est caractérisées par deux axes indissociables : (1) assurer sa reproduction (et donc sa survie jusqu'à la reproduction) (2) en étant adapté à son environnement (sa « niche »). Par analogie, en informatique, on qualifie ainsi la fonction objectif d'un algorithme évolutionnaire. 30, 31, 33, 36, 37, 41–43, 46, 48–54, 61, 64–66, 70, 117

forte concavité Une fonction f est fortement concave si son opposé $-f$ est fortement convexe. 11, 12

forte convexité Une fonction f de \mathbb{X} dans \mathbb{R} est fortement convexe avec la constante $b > 0$ si la fonction $x \mapsto f(x) - b\|x\|^2/2$ est convexe. L'introduction de (152) approfondit cette notion issue de l'optimisation convexe. 44

- graphe aléatoire d'Erdős-Rényi** Modèle de graphe aléatoire dans lequel la probabilité qu'une arête existe entre deux nœuds est uniforme. Voir aussi l'introduction de (150). 23, 25
- gène** Séquence d'ADN qui spécifie la synthèse d'une macromolécule dite « fonctionnelle », c'est-à-dire une protéine ou un ARN non messager. Les gènes sont généralement séparés sur l'ADN. On connaît avec précision les limites de leurs séquences codantes spécifiant la synthèse à proprement parler. En revanche, on ne connaît pas précisément l'étendue des séquences servant à initier leur transcription. Voir aussi la référence (4). vii, viii, 12, 25, 33, 52, 53, 56, 61–71, 73, 75, 85, 87, 88, 116, 117
- holométabole** Les holométaboles sont des insectes dont la larve diffère radicalement de l'adulte (chenille/papillon, asticot/mouche, etc.). Voir aussi <http://www.tolweb.org/Endopterygota/8243>. 54, 71
- hémimétabole** Les hémimétaboles sont des insectes caractérisés par un développement progressif, sans stade immobile (chrysalide) entre la larve et l'adulte. Au contraire des holométaboles, ils ne constituent pas une clade. 54, 71
- information mutuelle** En théorie de l'information, l'information mutuelle de deux variables mesure la dépendance mutuelle de deux variables. Elle est définie à partir de l'entropie de Shannon. Dans le cadre de cette thèse, par abus de langage, on parle d'information mutuelle à propos de la grandeur définie de manière similaire à partir de toute mesure de complexité vérifiant la propriété de sous-additivité. 6, 8–10, 13–15, 22, 24, 26
- insecte** Les insectes sont des animaux formant la classe des hexapodes. Voir aussi la référence (109) et <http://www.tolweb.org/Insecta/8205>. 54
- laplacien** En théorie des graphes, les éléments $l(i, j)$ de la matrice laplacienne L d'un graphe de n nœuds sont définis de la manière suivante : sur la diagonale, le degré du nœud i ($l(i, i) = d(i)$) ; hors diagonale, 0 partout sauf si il existe un lien entre i et j (Dans ce dernier cas, $l(i, j) = -1$). Voir aussi la référence (58). 16, 24
- modularité** La modularité de la coupe d'un système en modules est la fraction de complexité qui demeure dans les modules et non entre eux. Plus cette fraction proche de 1, plus la modularité est élevée. Par abus de langage, on parle de modularité d'un système lorsqu'il en existe des coupes modulaires. vii, 4, 6, 8–13, 16, 18, 20–27, 29–37, 40, 42, 46–48, 50–52, 55, 57, 64, 68, 69, 88, 94, 111, 115–117
- modularité fonctionnelle** La modularité d'une fonction objectif se mesure par la possibilité de la séparer en plusieurs fonctions objectifs optimisables indépendamment. vii, 28, 29, 31, 33–35, 40, 41, 43–46, 52, 54–57, 61, 64–66, 68, 73, 95, 116, 117

- modularité structurelle** Un module structurel est une partie de la configuration d'un système susceptible de varier indépendamment du reste de la configuration. vii, 27, 29, 31, 33–35, 39, 41, 43, 52, 56, 61, 64–66, 70, 73, 97, 115–117
- méiose** Division cellulaire particulière propre aux espèces possédant une reproduction sexuée dans lequel une cellule portant des paires de chromosomes homologues (diploïde) se divise en deux cellules appelées gamètes portant un chromosome de chaque paire (haploïde). La fusion de deux gamètes provenant d'individus généralement différents donne naissance à une cellule à nouveau diploïde pouvant se développer en un nouvel individu. 56, 63, 64
- métazoaire** Nom moderne du taxon constitué par les animaux. Voir aussi la référence (109) et <http://tolweb.org/Animals/2374>. 53, 71
- nématode** Vers ronds. Voir aussi la référence (109) et <http://www.tolweb.org/nematoda>. 71
- opéron** Groupe de gènes s'enchaînant en une séquence d'ADN transcrits simultanément en un seul ARN messenger. Les opérons existent essentiellement chez les procaryotes. Voir aussi la référence (4). vii, 62, 66, 68–71, 87, 88, 116, 117
- organite** structure délimitée par une membrane incluse dans une cellule (elle-même délimitée de l'extérieur par une membrane). 62
- point chaud de recombinaison** Région très peu étendue du génome d'une espèce à reproduction sexuée dans laquelle la probabilité d'enjambement est très élevée comparée à celle de son voisinage. vii, 65, 66, 73, 75, 87, 88, 116, 117
- principe de subsidiarité** Principe philosophique qui postule que la responsabilité des décisions doit être allouée à la plus petite entité capable de résoudre le problème d'elle-même. Voir aussi le traité de Rome (43). 46–49
- procaryote** Taxon regroupant l'ensemble des organismes vivants dont la cellule ne possède pas de noyau. L'ADN est donc accessible à l'intérieur de celles-ci et non confiné dans le noyau. Les bactéries sont les représentants les plus connus de ce taxon.. 62, 63, 66
- produit cartésien de graphes** Généralisation aux graphes du produit cartésien. Voir aussi http://en.wikipedia.org/wiki/Cartesian_product_of_graphs. 28
- protéine** Macromolécule constituée d'acides aminés produite par traduction de l'information génétique contenue dans l'ADN, par l'intermédiaire des ARN messagers. Les protéines assurent la plus grande partie des mécanismes cellulaires. Voir aussi la référence (4). 5, 25, 52, 61, 62, 66, 67, 70
- rasoir d'Occam** Principe philosophique selon lequel « Les multiples ne doivent pas être utilisés sans nécessité ». 5

- recombinaison** Mécanisme de reconnexion d'extrémités libres et similaires (homologues) d'ADN. Ce mécanisme permet de « réparer » des coupures de la molécule d'ADN. Dans cette thèse, par abus de langage, on parle de recombinaison homologue uniquement quand la réparation aboutit à la formation d'un ADN similaire à celui d'origine, c'est-à-dire que la reconnexion a permis de restaurer l'ADN initial ou a remplacé une séquence par une autre homologue, par exemple issue d'un chromosome homologue. Dans le cas contraire, l'ADN obtenu n'est pas homologue à celui de départ et nous parlerons de recombinaison hétérologue. Voir aussi la référence (4). 34, 55–57, 61, 63–68, 73–75, 85, 87, 88
- reproduction** ensemble des mécanismes par lesquels une espèce se perpétue. La reproduction implique la transmission de l'information génétique, support principal de l'hérédité. La reproduction sexuée implique la fusion du matériel génétique issu de deux individus pour en former un nouveau. Ce mécanisme implique l'existence de paires de chromosomes dits homologues porteurs d'une information génétique semblable. Chaque parent transmet un des chromosomes de chaque paire. L'existence de paires homologues est à l'origine du phénomène d'enjambement lors de la méiose. La reproduction asexuée est basée sur la division cellulaire. Elle n'implique pas l'existence de chromosomes homologues. 41, 53, 63, 66
- réplication** Processus de duplication de l'ADN dans une cellule en vue de sa division. Ce processus se déroule de manière linéaire en partant d'un point donné de l'ADN. Chez les procaryotes, ce processus est initié en un seul point de l'ADN, appelé origine de réplication. Chez les eucaryotes, il existe de nombreux points d'initiation. Les enzymes acteurs de ce mécanisme, en particulier celles qui lisent la séquence d'ADN, sont appelées de manière générique répliques. Voir aussi la référence (4). vii, 62, 68–71, 87, 116, 117
- science des systèmes** La science des systèmes ou systémique est un cadre conceptuel permettant de décrire ou d'analyser l'organisation de tout système. 5, 26, 118
- sclérose en plaques** Maladie auto-immune (liée à l'activité anormale de certains anticorps dirigés contre la gaine de myéline des fibres nerveuses) neurologique chronique souvent invalidante. 85
- système** Un système est un objet d'étude au sens général. Il possède des éléments reliés entre eux. Ces éléments caractérisent sa configuration et lui donnent ses propriétés. Il possède aussi une limite et des relations avec le reste de l'univers, ou environnement. vii, 3–13, 21, 24, 26–36, 38–40, 42–52, 55–57, 63, 65, 87, 91–99, 101, 102, 110–112, 115–118
- système évolutif et fonctionnel** Un système évolutif est un système qui possède une configuration susceptible d'évoluer. Un système fonctionnel est un système dont on peut mesurer l'adéquation de ses propriétés à un but. Dans cette thèse, un système évolutif et fonctionnel possède les deux propriétés et est en plus muni d'un processus d'évolution capable de faire évoluer sa

configuration afin d'améliorer son adéquation à son but. vii, viii, 26–30, 32–34, 43, 55–57, 116, 117

séquence codante Séquence d'ADN porteuse d'une information codant pour une ou plusieurs protéines, après traduction de cette information en utilisant le code génétique. Par opposition, les séquences non-codantes n'ont pas pour fonction d'être traduites en protéines. Voir aussi la référence (4). 61, 64, 65, 70

taxon En taxinomie, entité conceptuelle regroupant tous les organismes vivants possédant en commun certains caractères bien définis. Voir aussi la référence (109). 54

théorie de l'évolution neutre Théorie de l'évolution selon laquelle la plupart des mutations n'ont aucun impact direct sur la valeur de la fonction de fitness. Voir aussi la référence (95). 29

transcription Mécanisme de formation d'ARN messager (transcrit) contenant la même information qu'une séquence d'ADN. La transcription d'un gène est une condition sine qua non de son activation. Les enzymes acteurs de ce mécanisme, en particulier celles qui lisent la séquence d'ADN, sont appelées de manière générique transcriptases. L'initiation de la transcription d'un gène est régulée par des protéines régulatrices qui se fixent sur des séquences spécifiques dites elles aussi régulatrices. Voir aussi la référence (4). vii, 61, 62, 68–70, 87, 116, 117

transfert horizontal Le transfert horizontal (ou latéral) désigne le transfert d'information génétique entre deux organismes vivants en dehors des mécanismes de reproduction, dans lesquels l'information génétique est transmise « verticalement » aux descendants. 63, 66, 67

translocation Echange de matériel génétique entre deux chromosomes non homologues. Certaines translocations sont responsables de cancers (104). 63

volvocacée Les volvocacées sont un ordre d'algues vertes unicellulaires ou coloniales dans la division des chlorophytes. Les individus sont mobiles avec deux, quatre ou rarement huit flagelles en forme de fouet. Voir aussi http://www.tolweb.org/Green_plants. 53

épissage Chez les eucaryotes, les gènes sont souvent constitués d'une succession de séquences codantes appelées exons et de séquences non codantes appelées introns. L'ensemble du gène est transcrit en ARN, introns inclus. L'excision des introns et le recollement des exons voisins, appelée épissage, a lieu avant la traduction. Voir aussi la référence (4). 61, 65, 70

Bibliographie

- [1] C. ADAMI – « What is complexity ? », *Bioessays* **24** (2002), no. 12, p. 1085–1094.
- [2] A. C. AHN, M. TEWARI, C.-S. POON et R. S. PHILLIPS – « The limits of reductionism in medicine : could systems biology offer an alternative ? », *PLoS Med* **3** (2006), no. 6, p. e208.
- [3] R. ALBERT et A.-L. BARABASI – « Statistical mechanics of complex networks », *Reviews of Modern Physics* **74** (2002), p. 47–97.
- [4] B. ALBERTS, A. JOHNSON, J. LEWIS, M. RAFF, K. ROBERTS et P. WALTER – *Molecular biology of the cell*, Garland, 2007.
- [5] M. ALDANA et P. CLUZEL – « A natural class of robust networks. », *Proc Natl Acad Sci U S A* **100** (2003), no. 15, p. 8710–8714.
- [6] E. ALMAAS – « Biological impacts and context of network theory. », *J Exp Biol* **210** (2007), no. Pt 9, p. 1548–1558.
- [7] U. ALON – « Biological networks : the tinkerer as an engineer. », *Science* **301** (2003), no. 5641, p. 1866–1867.
- [8] K. ARAKAWA et M. TOMITA – « Selection effects on the positioning of genes and gene structures from the interplay of replication and transcription in bacterial genomes », *Evolutionary Bioinformatics* **3** (2007), p. 279–286.
- [9] K. G. ARDLIE, L. KRUGLYAK et M. SEIELSTAD – « Patterns of linkage disequilibrium in the human genome. », *Nat Rev Genet* **3** (2002), no. 4, p. 299–309.
- [10] R. BACHELET – *Cours d'analyse fonctionnelle en conception*, Ecole Centrale de Lille, 2007.
- [11] C. Y. BALDWIN et K. B. CLARK – « Complex engineered systems. », ch. Modularity in the design of complex engineering systems, p. 175–205, Springer, Berlin / Heidelberg, 2006.

- [12] A.-L. BARABÁSI, Z. DEZSŐ, E. RAVASZ, S.-H. YOOK, et Z. N. OLTVAI – « Scale-free and hierarchical structures in complex networks », *AIP Conf. Proc.* **661** (2003), p. 1–16.
- [13] A.-L. BARABÁSI et Z. N. OLTVAI – « Network biology : understanding the cell's functional organization. », *Nat Rev Genet* **5** (2004), no. 2, p. 101–113.
- [14] Y. BENGIO – « Using a financial training criterion rather than a prediction criterion », CIRANO Working Papers 98s-21, CIRANO, 1998.
- [15] A. BERGMAN et M. L. SIEGAL – « Evolutionary capacitance as a general feature of complex gene networks. », *Nature* **424** (2003), no. 6948, p. 549–552.
- [16] I. BJEDOV, O. TENAILLON, B. GÉRARD, V. SOUZA, E. DENAMUR, M. RADMAN, F. TADDEI et I. MATIC – « Stress-induced mutagenesis in bacteria. », *Science* **300** (2003), no. 5624, p. 1404–1409.
- [17] S. BORNHOLDT – « Systems biology. Less is more in modeling large genetic networks. », *Science* **310** (2005), no. 5747, p. 449–451.
- [18] Y. BOUCHER, C. J. DOUADY, R. T. PAPKE, D. A. WALSH, M. E. R. BOUDREAU, C. L. NESBØ, R. J. CASE et W. F. DOOLITTLE – « Lateral gene transfer and the origins of prokaryotic groups. », *Annu Rev Genet* **37** (2003), p. 283–328.
- [19] M. BUYGI, G. BALZER, H. SHANECHI et M. SHAHIDEHPOUR – « Market-based transmission expansion planning », *IEEE Transactions on Power Systems* **19** (2004), p. 2060–2067.
- [20] A. CAPOCCI, V. D. SERVEDIO, G. CALDARELLI et F. COLAIORI – « Communities detection in large networks. », *Algorithms and Models for the Web-Graph* (S. Leonardi, éd.), Lecture Notes in Computer Science, vol. 3243/2004, Springer Berlin / Heidelberg, 2004, p. 181–187.
- [21] J. M. CARLSON et J. DOYLE – « Highly optimized tolerance : A mechanism for power laws in designed systems », *Physical Review E* **60** (1999), no. 2, p. 1412+.
- [22] J. M. CARLSON et J. DOYLE – « Complexity and robustness. », *Proc Natl Acad Sci U S A* **99 Suppl 1** (2002), p. 2538–2545.
- [23] G. J. CHAITIN – « The maximum entropy formalism », ch. Toward a mathematical definition of "life", p. 477–498, MIT Press, 1979.
- [24] P. CHECKLAND – *Systems thinking, systems practice.*, Wiley, Chichester, U.K., 1981.
- [25] Y. CHEN et N. V. DOKHOLYAN – « The coordinated evolution of yeast proteins is constrained by functional modularity. », *Trends Genet* **22** (2006), no. 8, p. 416–419.

- [26] F. R. K. CHUNG – « Eigenvalues of graphs. », *Proceedings of the International Congress of Mathematicians, Zürich*, Birkhäuser Verlag, Berlin, 1994, p. 1333–1342.
- [27] F. R. K. CHUNG et S.-T. YAU – « Eigenvalue inequalities for graphs and convex subgraphs », *Communications on Analysis and Geometry* **5** (1997), p. 575–623.
- [28] M. J. COHN – « A review of “from DNA to diversity : molecular genetics and the evolution of animal design” by Sean B. Carroll, Jennifer K. Grenier, and Scott D. Weatherbee », *Evolution & Development* **3** (2001), p. 364–365.
- [29] COMMISSION DES COMMUNAUTÉS EUROPÉENNES – *Exposé des motifs du 3ème paquet législatif relatif au marché européen du gaz et de l’électricité.*, 2007.
- [30] — , *Impact assessment accompanying the 3rd legislative package on EU electricity and gas markets.*, 2007.
- [31] G. COOP – « Can a genome change its (hot)spots? », *Trends Ecol Evol* **20** (2005), no. 12, p. 643–645.
- [32] A. CRAMERI, G. DAWES, E. RODRIGUEZ, S. SILVER et W. P. STEMMER – « Molecular evolution of an arsenate detoxification pathway by dna shuffling. », *Nat Biotechnol* **15** (1997), no. 5, p. 436–438.
- [33] M. E. CSETE et J. C. DOYLE – « Reverse engineering of biological complexity. », *Science* **295** (2002), no. 5560, p. 1664–1669.
- [34] G. B. DANTZIG et P. WOLFE – « Decomposition principle for linear programs. », *Operations Research* **8** (1960), no. 1, p. 101–111.
- [35] F. DARDEL et F. KÉPÈS – *Bioinformatique. génomique et post-génomique*, Les Editions de l’Ecole Polytechnique, 2002.
- [36] G. VON DASSOW, E. MEIR, E. M. MUNRO et G. M. ODELL – « The segment polarity network is a robust developmental module. », *Nature* **406** (2000), no. 6792, p. 188–192.
- [37] H. DEFALVARD – *Fondements de la microéconomie*, vol. 2, L’équilibre des marchés, De Boeck Université, 2003.
- [38] A. DEMARET – *Ethologie et psychiatrie*, Mardaga, 1995.
- [39] Z. DEZSŐ – « The topology and dynamics of complex networks », Thèse, University of Notre Dame, 2005.

- [40] T. DOBZHANSKY – « Nothing in biology makes sense except in the light of evolution. », *The American Biology Teacher* **35** (1973), p. 125–129.
- [41] J. M. DURÃO BARROSO – *Travaillons ensemble pour la croissance et l'emploi. un nouvel élan pour la stratégie de lisbonne.*, 2005, Communication du président de la commission des communautés européennes.
- [42] A. DVORAK et AL. – *Typewriter keyboard*, 1936, U.S. Patent #2, 040, 248.
- [43] Etats fondateurs – Rome, *Traité instituant la communauté européenne*, 1957, Version consolidée dans le journal officiel n°C 325 du 24 décembre 2002.
- [44] L. DA F. COSTA, F. A. RODRIGUES, G. TRAVIESO et P. R. V. BOAS – « Characterization of complex networks : A survey of measurements », *Advances In Physics* **56** (2007), p. 167.
- [45] R. FANG et D. HILL – « A new strategy for transmission expansion in competitive electricity markets », *IEEE Transactions on Power Systems* **18** (2003), p. 374–380.
- [46] C. FLAMM, I. HOFACKER, B. STADLER et P. STADLER – « Saddles and barrier in landscapes of generalized search operators », *Foundations of Genetic Algorithms*, Lecture Notes in Computer Science, vol. 4436, 2007, p. 194–212.
- [47] A. FORCE, W. A. CRESKO, F. B. PICKETT, S. R. PROULX, C. AMEMIYA et M. LYNCH – « The origin of subfunctions and modular gene regulation. », *Genetics* **170** (2005), no. 1, p. 433–446.
- [48] S. FORTUNATO et M. BARTHÉLEMY – « Resolution limit in community detection. », *Proc Natl Acad Sci U S A* **104** (2007), no. 1, p. 36–41.
- [49] H. B. FRASER – « Modularity and evolutionary constraint on proteins. », *Nat Genet* **37** (2005), no. 4, p. 351–352.
- [50] — , « Coevolution, modularity and human disease. », *Curr Opin Genet Dev* **16** (2006), no. 6, p. 637–644.
- [51] L. C. FREEMAN – « A set of measures of centrality based on betweenness. », *Sociometry* **40** (1977), no. 1, p. 35–41.
- [52] L. FRIMANNSLUND – « On curvature and separability in unconstrained optimisation. », Thèse, University of Bergen, Norway, 2006.
- [53] A. GARDNER et A. T. KALINKA – « Recombination and the evolution of mutational robustness. », *J Theor Biol* **241** (2006), no. 4, p. 707–715.

- [54] A. GARDNER et W. ZUIDEMA – « Is evolvability involved in the origin of modular variation ? », *Evolution* **57** (2003), no. 6, p. 1448–1450.
- [55] M. R. GAREY et D. S. JOHNSON – *Computers and intractability : A guide to the theory of NP-completeness*, W. H. Freeman, 1979.
- [56] J. L. GERTON, J. DERISI, R. SHROFF, M. LICHTEN, P. O. BROWN et T. D. PETES – « Inaugural article : global mapping of meiotic recombination hotspots and coldspots in the yeast *saccharomyces cerevisiae*. », *Proc Natl Acad Sci U S A* **97** (2000), no. 21, p. 11383–11390.
- [57] D. B. GOLDSTEIN et M. E. WEALE – « Population genomics : linkage disequilibrium holds the key. », *Curr Biol* **11** (2001), no. 14, p. R576–R579.
- [58] M. GONDRAN et M. MINOUX – *Graphes et algorithmes*, Eyrolles, 1979.
- [59] R. J. GREEN – *Electricity transmission pricing : How much does it cost to get it wrong ?*, Cambridge Working Paper in Economics, 2004.
- [60] M. GUEDJ – « Méthodes statistiques pour l’analyse des données génétiques d’association à grande Échelle. », Thèse, Université d’Evry Val d’Essonne, 2007.
- [61] N. GUELZIM, S. BOTTANI, P. BOURGINE et F. KÉPÈS – « Topological and causal structure of the yeast transcriptional regulatory network. », *Nat Genet* **31** (2002), no. 1, p. 60–63.
- [62] H. L. GUYADER – « Doit-on abandonner le concept d’espèce ? », *Le Courrier de l’environnement de l’INRA* **46** (2002), p. 51–64.
- [63] H. HAMDA, F. JOUVE, E. LUTTON, M. SCHOENAUER et M. SEBAG – « Compact unstructured representations for evolutionary design », *Applied Intelligence* **16** (2002), p. 139–155.
- [64] T. F. HANSEN – « Is modularity necessary for evolvability ? remarks on the relationship between pleiotropy and evolvability. », *Biosystems* **69** (2003), no. 2-3, p. 83–94.
- [65] G. HARIK – *Linkage learning via probabilistic modeling in the ECGA*, Illinois Genetic Algorithms Laboratory Technical Report, 1999.
- [66] C. HARTLAND – *Mécanisme d’anticipation pour la robustesse de contrôleurs insitu*, Mémoire, Université Paris XI, 2005.
- [67] L. H. HARTWELL, J. J. HOPFIELD, S. LEIBLER et A. W. MURRAY – « From molecular to modular cell biology. », *Nature* **402** (1999), no. 6761 Suppl, p. C47–C52.

- [68] M. HUISKEN, C. IGEL et M. TOUSSAINT – « Task-dependent evolution of modularity in neural networks », *Connection Science* **14** (2002), p. 2002.
- [69] W. G. HILL et A. ROBERTSON – « The effect of linkage on limits to artificial selection. », *Genet Res* **8** (1966), no. 3, p. 269–294.
- [70] D. A. HINDS, L. L. STUVE, G. B. NILSEN, E. HALPERIN, E. ESKIN, D. G. BALLINGER, K. A. FRAZER et D. R. COX – « Whole-genome patterns of common dna variation in three human populations. », *Science* **307** (2005), no. 5712, p. 1072–1079.
- [71] A. HINTZE et C. ADAMI – « Evolution of complex modular biological networks. », *PLoS Comput Biol* **4** (2008), no. 2, p. e23.
- [72] J. N. HIRSCHHORN et M. J. DALY – « Genome-wide association studies for common diseases and complex traits. », *Nat Rev Genet* **6** (2005), no. 2, p. 95–108.
- [73] C. T. HITTINGER, A. ROKAS et S. B. CARROLL – « Parallel inactivation of multiple gal pathway genes and ecological diversification in yeasts. », *Proc Natl Acad Sci U S A* **101** (2004), no. 39, p. 14144–14149.
- [74] S. D. HOOPER et O. G. BERG – « On the nature of gene innovation : duplication patterns in microbial genomes. », *Mol Biol Evol* **20** (2003), no. 6, p. 945–954.
- [75] G. HORNBY – « Measuring, enabling and comparing modularity, regularity and hierarchy in evolutionary design », *Proceedings of the Genetic and Evolutionary Computation Conference (GECCO-2005)*, 2005, p. 1729–1736.
- [76] R. HORTON, L. WILMING, V. RAND, R. C. LOVERING, E. A. BRUFORD, V. K. KHODIYAR, M. J. LUSH, S. POVEY, C. C. TALBOT, M. W. WRIGHT, H. M. WAIN, J. TROWSDALE, A. ZIEGLER et S. BECK – « Gene map of the extended human MHC. », *Nat Rev Genet* **5** (2004), no. 12, p. 889–899.
- [77] INVESTIGATION COMMITTEE – *System disturbance on 4 november 2006.*, UCTE, 2007.
- [78] F. JACOB – « Evolution and tinkering. », *Science* **196** (1977), no. 4295, p. 1161–1166.
- [79] — , « Complexity and tinkering. », *Ann N Y Acad Sci* **929** (2001), p. 71–73.
- [80] A. J. JEFFREYS, L. KAUPPI et R. NEUMANN – « Intensely punctate meiotic recombination in the class II region of the major histocompatibility complex. », *Nat Genet* **29** (2001), no. 2, p. 217–222.
- [81] A. J. JEFFREYS, R. NEUMANN, M. PANAYI, S. MYERS et P. DONNELLY – « Human recombination hot spots hidden in regions of strong marker association. », *Nat Genet* **37** (2005), no. 6, p. 601–606.

- [82] H. JEONG, S. P. MASON, A. L. BARABÁSI et Z. N. OLTVAI – « Lethality and centrality in protein networks. », *Nature* **411** (2001), no. 6833, p. 41–42.
- [83] J. JÄGERSKÜPPER et T. STORCH – « How comma selection helps with the escape from local optima », *Parallel Problem Solving from Nature - PPSN IX*, Lecture Notes in Computer Science, vol. 4193, Springer Berlin / Heidelberg, 2006, p. 52–61.
- [84] E. D. DE JONG – « Representation development from pareto-coevolution », *Proceedings of the Genetic and Evolutionary Computation Conference (GECCO-2003)*, vol. 2723/2003, 2003.
- [85] E. D. DE JONG, D. THIERENS et R. A. WATSON – « Defining modularity, hierarchy, and repetition. », *Proceedings of the Genetic and Evolutionary Computation Conference (GECCO-2004)*, 2004.
- [86] N. KASHTAN et U. ALON – « Spontaneous evolution of modularity and network motifs. », *Proc Natl Acad Sci U S A* **102** (2005), no. 39, p. 13773–13778.
- [87] M. KATO, F. MIYA, Y. KANEMURA, T. TANAKA, Y. NAKAMURA et T. TSUNODA – « Recombination rates of genes expressed in human tissues. », *Hum Mol Genet* **17** (2008), no. 4, p. 577–586.
- [88] M. KATO, A. SEKINE, Y. OHNISHI, T. A. JOHNSON, T. TANAKA, Y. NAKAMURA et T. TSUNODA – « Linkage disequilibrium of evolutionarily conserved regions in the human genome. », *BMC Genomics* **7** (2006), p. 326.
- [89] S. A. KAUFFMAN – *The origins of order : Self-organization and selection in evolution.*, Oxford University Press, 1993.
- [90] L. KAUPPI, A. J. JEFFREYS et S. KEENEY – « Where the crossovers are : recombination distributions in mammals. », *Nat Rev Genet* **5** (2004), no. 6, p. 413–424.
- [91] P. J. KEELING et J. D. PALMER – « Horizontal gene transfer in eukaryotic evolution. », *Nat Rev Genet* **9** (2008), no. 8, p. 605–618.
- [92] E. F. KELLER – « Revisiting "scale-free" networks. », *Bioessays* **27** (2005), no. 10, p. 1060–1068.
- [93] D. N. KEYS, D. L. LEWIS, J. E. SELEGUE, B. J. PEARSON, L. V. GOODRICH, R. L. JOHNSON, J. GATES, M. P. SCOTT et S. B. CARROLL – « Recruitment of a hedgehog regulatory circuit in butterfly eyespot evolution. », *Science* **283** (1999), no. 5401, p. 532–534.
- [94] D.-H. KIM et A. E. MOTTER – « Fluctuation-driven capacity distribution in complex networks », *New Journal of Physics* **10** (2008), no. 5, p. 053022 (19pp).
- [95] M. KIMURA – « Evolutionary rate at the molecular level. », *Nature* **217** (1968), no. 5129, p. 624–626.

- [96] M. KIRSCHNER et J. GERHART – « Evolvability. », *Proc Natl Acad Sci U S A* **95** (1998), no. 15, p. 8420–8427.
- [97] H. KITANO – « Biological robustness. », *Nat Rev Genet* **5** (2004), no. 11, p. 826–837.
- [98] — , « Towards a theory of biological robustness. », *Mol Syst Biol* **3** (2007), p. 137.
- [99] H. KITANO, K. ODA, T. KIMURA, Y. MATSUOKA, M. CSETE, J. DOYLE et M. MURAMATSU – « Metabolic syndrome and robustness tradeoffs. », *Diabetes* **53 Suppl 3** (2004), p. S6–S15.
- [100] A. N. KOLMOGOROV – « Three approaches for defining the concept of information quantity », *Information Transmission* **1** (1965), p. 3–11.
- [101] F. KÉPÈS – « Periodic transcriptional organization of the *E.coli* genome. », *J Mol Biol* **340** (2004), no. 5, p. 957–964.
- [102] D. C. KRAKAUER – « Stability and evolution of overlapping genes. », *Evolution* **54** (2000), no. 3, p. 731–739.
- [103] D. C. KRAKAUER – *Robustness in biological systems – a provisional taxonomy*, Santa Fe Institute Working Paper, 2003.
- [104] R. KURZROCK, H. M. KANTARJIAN, B. J. DRUKER et M. TALPAZ – « Philadelphia chromosome-positive leukemias : from basic mechanisms to molecular therapeutics. », *Ann Intern Med* **138** (2003), no. 10, p. 819–830.
- [105] A. D. LANDER – « A calculus of purpose. », *PLoS Biol* **2** (2004), no. 6, p. e164.
- [106] R. N. LANGLOIS – « Modularity in technology and organization. », *Journal of Economic Behavior & Organization* **49** (2002), no. 1, p. 19–37.
- [107] G. LATORRE, R. CRUZ, J. AREIZA et A. VILLEGAS – « Classification of publications and models on transmission expansion planning », *IEEE Transactions on Power Systems* **18** (2003), p. 938–946.
- [108] J. G. LAWRENCE et J. R. ROTH – « Selfish operons : horizontal transfer may drive the evolution of gene clusters. », *Genetics* **143** (1996), no. 4, p. 1843–1860.
- [109] G. LECOINTRE et H. L. GUYADER – *Classification phylogénétique du vivant*, Editions Belin, 2006.
- [110] R. E. LENSKI, J. E. BARRICK et C. OFRIA – « Balancing robustness and evolvability. », *PLoS Biol* **4** (2006), no. 12, p. e428.

- [111] H. LIPSON – « Principles of modularity, regularity, and hierarchy for scalable systems. », *Journal of Biological Physics and Chemistry* **7** (2007), p. 125–128.
- [112] H. LIPSON, J. B. POLLACK et N. P. SUH – « Promoting modularity in evolutionary design. », *Proceedings of the 2001 ASME Design Engineering Technical Conferences*, 2001.
- [113] — , « On the origin of modular variation. », *Evolution* **56** (2002), p. 1549–1556.
- [114] J. LOBO, J. H. MILLER et W. FONTANA – *Neutrality in technological landscapes*, Sante Fe Institute Working Paper, 2004.
- [115] W. MA, L. LAI, Q. OUYANG et C. TANG – « Robustness and modular design of the drosophila segment polarity network. », *Mol Syst Biol* **2** (2006), p. 70.
- [116] I. MAKALOWSKA, C.-F. LIN et W. MAKALOWSKI – « Overlapping genes in vertebrate genomes. », *Comput Biol Chem* **29** (2005), no. 1, p. 1–12.
- [117] R. MARSHALL et P. LEANEY – « A systems engineering approach to product modularity. », *Proceedings of the Institution of Mechanical Engineers Part B.*, 1999.
- [118] M. I. MCCARTHY, G. R. ABECASIS, L. R. CARDON, D. B. GOLDSTEIN, J. LITTLE, J. P. A. IOANNIDIS et J. N. HIRSCHHORN – « Genome-wide association studies for complex traits : consensus, uncertainty and challenges. », *Nat Rev Genet* **9** (2008), no. 5, p. 356–369.
- [119] P. MICHAEL – « The arguments for and against ownership unbundling of energy transmission networks », *Energy policy* **36** (2008), p. 704–713.
- [120] R. MILO, S. SHEN-ORR, S. ITZKOVITZ, N. KASHTAN, D. CHKLOVSKII et U. ALON – « Network motifs : simple building blocks of complex networks. », *Science* **298** (2002), no. 5594, p. 824–827.
- [121] R. MILO, S. ITZKOVITZ, N. KASHTAN, R. LEVITT, S. SHEN-ORR, I. AYZENSHTAT, M. SHEFFER et U. ALON – « Superfamilies of evolved and designed networks. », *Science* **303** (2004), no. 5663, p. 1538–1542.
- [122] S. D. MITCHELL – « Modularity – more than a buzzword ? », *Biological Theory* **1** (2006), no. 1, p. 98–101.
- [123] M. MIZOKAMI, E. ORITO, K. OHBA, K. IKEO, J. Y. LAU et T. GOJOBORI – « Constrained evolution with respect to gene overlap of hepatitis b virus. », *J Mol Evol* **44 Suppl 1** (1997), p. S83–S90.
- [124] S. MYERS, L. BOTTOLO, C. FREEMAN, G. McVEAN et P. DONNELLY – « A fine-scale map of recombination rates and hotspots across the human genome. », *Science* **310** (2005), no. 5746, p. 321–324.

- [125] M. W. NACHMAN, V. L. BAUER, S. L. CROWELL et C. F. AQUADRO – « DNA variability and recombination rates at X-linked loci in humans. », *Genetics* **150** (1998), no. 3, p. 1133–1141.
- [126] A. M. NEDELCOU et R. E. MICHOD – « Modularity in development and evolution. », ch. Evolvability, Modularity, and Individuality during the Transition to Multicellularity in Volvocalean Green Algae, p. 466–489, The University of Chicago Press, 2004.
- [127] M. E. J. NEWMAN – « The structure and function of complex networks », *SIAM Review* **45** (2003), p. 167–256.
- [128] M. E. J. NEWMAN – « Finding community structure in networks using the eigenvectors of matrices », *Physical Review E* **74** (2006), p. 036104.
- [129] M. E. J. NEWMAN et M. GIRVAN – « Finding and evaluating community structure in networks. », *Phys Rev E Stat Nonlin Soft Matter Phys* **69** (2004), no. 2 Pt 2, p. 026113.
- [130] E. VAN NIMWEGEN, J. P. CRUTCHFIELD et M. HUYNEN – « Neutral evolution of mutational robustness. », *Proc Natl Acad Sci U S A* **96** (1999), no. 17, p. 9716–9720.
- [131] Y. D. NOCHOMOVITZ et H. LI – « Highly designable phenotypes and mutational buffers emerge from a systematic mapping between network topology and dynamic output. », *Proc Natl Acad Sci U S A* **103** (2006), no. 11, p. 4180–4185.
- [132] S. NOLFI, J. L. ELMAN et D. PARISI – « Learning and evolution in neural networks », *Adaptive Behavior* **3** (1994), p. 5–28.
- [133] H. OCHMAN, J. G. LAWRENCE et E. A. GROISMAN – « Lateral gene transfer and the nature of bacterial innovation. », *Nature* **405** (2000), no. 6784, p. 299–304.
- [134] Z. N. OLTVAI et A.-L. BARABÁSI – « Systems biology. life's complexity pyramid. », *Science* **298** (2002), no. 5594, p. 763–764.
- [135] S. L. OOI, X. PAN, B. D. PEYSER, P. YE, P. B. MELUH, D. S. YUAN, R. A. IRIZARRY, J. S. BANDER, F. A. SPENCER et J. D. BOEKE – « Global synthetic-lethality analysis and yeast functional profiling. », *Trends Genet* **22** (2006), no. 1, p. 56–63.
- [136] S. R. PALADUGU, V. CHICKARMANE, A. DECKARD, J. P. FRUMKIN, M. McCORMACK et H. M. SAURO – « In silico evolution of functional modules in biochemical networks. », *Syst Biol (Stevenage)* **153** (2006), no. 4, p. 223–235.
- [137] PARLEMENT ET CONSEIL EUROPÉEN – *Directive 96/92/CE concernant des règles communes pour le marché intérieur de l'électricité.*, 1996.

- [138] L. PATHY – « Modular assembly of genes and the evolution of new functions. », *Genetica* **118** (2003), no. 2-3, p. 217–231.
- [139] I. PEREZ-ARRIAGA, F. RUBIO, J. PUERTA, J. ARCELUZ et J. MARIN – « Marginal pricing of transmission services : an analysis of costrecovery », *IEEE Transactions on Power Systems* **10** (1995), p. 546–553.
- [140] D. C. PLAUT – « Double dissociation without modularity : evidence from connectionist neuropsychology. », *J Clin Exp Neuropsychol* **17** (1995), no. 2, p. 291–321.
- [141] R. POLI, W. B. LANGDON et N. F. MCPHEE – *A field guide to genetic programming*, Lulu Enterprises, 2008.
- [142] A. M. POOLE, M. J. PHILLIPS et D. PENNY – « Prokaryote and eukaryote evolvability. », *Biosystems* **69** (2003), no. 2-3, p. 163–185.
- [143] M. N. PRICE, E. J. ALM et A. P. ARKIN – « Interruptions in gene expression drive highly expressed operons to the leading strand of DNA replication. », *Nucleic Acids Res* **33** (2005), no. 10, p. 3224–3234.
- [144] S. R. PROULX, S. NUZHIDIN et D. E. L. PROMISLOW – « Direct selection on genetic robustness revealed in the yeast transcriptome. », *PLoS ONE* **2** (2007), no. 9, p. e911.
- [145] Y. QI et H. GE – « Modularity and dynamics of cellular networks. », *PLoS Comput Biol* **2** (2006), no. 12, p. e174.
- [146] C. QUEITSCH, T. A. SANGSTER et S. LINDQUIST – « *Hsp90* as a capacitor of phenotypic variation. », *Nature* **417** (2002), no. 6889, p. 618–624.
- [147] F. RADICCHI, C. CASTELLANO, F. CECCONI, V. LORETO et D. PARISI – « Defining and identifying communities in networks. », *Proc Natl Acad Sci U S A* **101** (2004), no. 9, p. 2658–2663.
- [148] P. B. RAINEY et T. F. COOPER – « Evolution of bacterial diversity and the origins of modularity. », *Res Microbiol* **155** (2004), no. 5, p. 370–375.
- [149] E. RAVASZ, A. L. SOMERA, D. A. MONGRU, Z. N. OLTVAI et A. L. BARABÁSI – « Hierarchical organization of modularity in metabolic networks. », *Science* **297** (2002), no. 5586, p. 1551–1555.
- [150] E. RAVASZ – « Evolution, hierarchy and modular organization in complex networks », Thèse, University of Notre Dame, 2004.

- [151] D. E. REICH, M. CARGILL, S. BOLK, J. IRELAND, P. C. SABETI, D. J. RICHTER, T. LAVERY, R. KOUYOUMJIAN, S. F. FARHADIAN, R. WARD et E. S. LANDER – « Linkage disequilibrium in the human genome. », *Nature* **411** (2001), no. 6834, p. 199–204.
- [152] A. RENAUD – « Algorithmes de régularisation et décomposition pour les problèmes variationnels monotones. », Thèse, Ecole Nationale Supérieure des Mines de Paris, 1993.
- [153] C. RODRIGUEZ-CASO, M. A. MEDINA et R. V. SOLÉ – « Topology, tinkering and evolution of the human transcription factor network. », *FEBS J* **272** (2005), no. 24, p. 6423–6434.
- [154] M. ROSAS-CASALS, S. VALVERDE et R. V. SOLE – « Topological vulnerability of the european power grid under errors and attacks », *Int J Bifurcat Chaos Appl Sci Eng* **17** (2007), p. 2465–2476.
- [155] M. ROSVALL et C. T. BERGSTROM – « An information-theoretic framework for resolving community structure in complex networks. », *Proc Natl Acad Sci U S A* **104** (2007), no. 18, p. 7327–7331.
- [156] — , « Maps of random walks on complex networks reveal community structure. », *Proc Natl Acad Sci U S A* **105** (2008), no. 4, p. 1118–1123.
- [157] RTE FORECAST TEAM – *Electricity consumption in France : Characteristics and forecast method*, http://www.rte-france.com/htm/an/vie/telecharge/prevconsoelec_an.pdf.
- [158] A. RUSZCZYŃSKI – *Nonlinear optimization*, Princeton University Press, 2006.
- [159] S. L. RUTHERFORD – « From genotype to phenotype : buffering mechanisms and the storage of genetic information. », *Bioessays* **22** (2000), no. 12, p. 1095–1105.
- [160] M. SAKO – « The business of systems integration », ch. Modularity and outsourcing : the nature of co-evolution of product and organisation architecture in the global automotive industry, p. 229–53, Oxford University Press, 2003.
- [161] G. SCHLOSSER – « Modularity and the units of evolution. », *Theory in Biosciences* **121** (2002), no. 1, p. 1–80.
- [162] — , « Modularity in development and evolution. », ch. The Role of Modules in Development and Evolution, p. 519–582, in Schlosser et Wagner (163), 2004.
- [163] G. SCHLOSSER et G. P. WAGNER (éds.) – *Modularity in development and evolution*, The University of Chicago Press, 2004.

- [164] E. A. SCHULTES et D. P. BARTEL – « One sequence, two ribozymes : implications for the emergence of new ribozyme folds. », *Science* **289** (2000), no. 5478, p. 448–452.
- [165] T. J. SEJNOWSKI et C. R. ROSENBERG – « Parallel networks that learn to pronounce english text. », *Complex Systems* **1** (1987), p. 145–168.
- [166] S. S. SHEN-ORR, R. MILO, S. MANGAN et U. ALON – « Network motifs in the transcriptional regulation network of escherichia coli. », *Nat Genet* **31** (2002), no. 1, p. 64–68.
- [167] H. A. SIMON – *The sciences of the artificial.*, 1996 éd., MIT Press, Cambridge, MA, 1969.
- [168] A. V. SMITH, D. J. THOMAS, H. M. MUNRO et G. R. ABECASIS – « Sequence features in regions of weak and strong linkage disequilibrium. », *Genome Res.* **15** (2005), p. 1519–1534.
- [169] R. V. SOLÉ et P. FERNANDEZ – *Modularity “for free” in genome architecture ?*, Santa Fe Institute Working Paper, 2003.
- [170] R. V. SOLÉ, M. ROSAS-CASALS, B. COROMINAS-MURTRA et S. VALVERDE – « Robustness of the european power grids under intentional attack. », *Phys Rev E Stat Nonlin Soft Matter Phys* **77** (2008), no. 2 Pt 2, p. 026102.
- [171] R. V. SOLÉ et S. VALVERDE – « Spontaneous emergence of modularity in cellular networks. », *J R Soc Interface* **5** (2008), no. 18, p. 129–133.
- [172] B. M. STADLER, P. F. STADLER, G. P. WAGNER et W. FONTANA – « The topology of the possible : formal spaces underlying patterns of evolutionary change. », *J Theor Biol* **213** (2001), no. 2, p. 241–274.
- [173] J. STELLING – « Mathematical models in microbial systems biology. », *Curr Opin Microbiol* **7** (2004), no. 5, p. 513–518.
- [174] A. SVEJGAARD – « The immunogenetics of multiple sclerosis. », *Immunogenetics* **60** (2008), no. 6, p. 275–286.
- [175] SYSTINT WORKGROUP – *European, cis and mediterranean interconnection : State of play 2006.*, UCTE–EURELECTRIC, 2007.
- [176] D. THIEFFRY et D. ROMERO – « The modularity of biological regulatory networks. », *Biosystems* **50** (1999), no. 1, p. 49–59.
- [177] T. A. THORNTON-WELLS, J. H. MOORE et J. L. HAINES – « Genetics, statistics and human disease : analytical retooling for complexity. », *Trends Genet* **20** (2004), no. 12, p. 640–647.

- [178] M. TOUSSAINT – « On the evolution of phenotypic exploration distributions. », *Foundations of Genetic Algorithms*, Lecture Notes in Computer Science, vol. 7, Springer Berlin / Heidelberg, 2003, p. 169–182.
- [179] M. TOUSSAINT et W. VON SEELEN – « Complex adaptation and system structure. », *BioSystems* **90** (2007), p. 769–782.
- [180] T. TSUNODA, G. M. LATHROP, A. SEKINE, R. YAMADA, A. TAKAHASHI, Y. OHNISHI, T. TANAKA et Y. NAKAMURA – « Variation of gene-based SNPs and linkage disequilibrium patterns in the human genome. », *Hum Mol Genet* **13** (2004), no. 15, p. 1623–1632.
- [181] UCTE – *Operation handbook.*, 2004.
- [182] M. H. V. VAN REGENMORTEL – « Reductionism and complexity in molecular biology. Scientists now have the tools to unravel biological and overcome the limitations of reductionism. », *EMBO Rep* **5** (2004), no. 11, p. 1016–1020.
- [183] E. A. VARIANO, J. H. MCCOY et H. LIPSON – « Networks, dynamics, and modularity. », *Phys Rev Lett* **92** (2004), no. 18, p. 188701.
- [184] S. VASSENA, P. MACK, P. ROUSSEAU, C. DRUET et L. WEHENKEL – « A probabilistic approach to power system network planning under uncertainties », *IEEE Bologna Power Tech Conference Proceedings*, 2003.
- [185] J. M. G. VILAR – « Modularizing gene regulation. », *Mol Syst Biol* **2** (2006), p. 2006.0016.
- [186] J. VIÑUELAS, F. CALEVRO, D. REMOND, J. BERNILLON, Y. RAHBÉ, G. FEBVAY, J.-M. FAYARD et H. CHARLES – « Conservation of the links between gene transcription and chromosomal organization in the highly reduced genome of *Buchnera aphidicola*. », *BMC Genomics* **8** (2007), p. 143.
- [187] C. A. VOIGT, C. MARTINEZ, Z.-G. WANG, S. L. MAYO et F. H. ARNOLD – « Protein building blocks preserved by recombination. », *Nat Struct Biol* **9** (2002), no. 7, p. 553–558.
- [188] WAGNER – *Robustness and evolvability in living systems.*, Princeton University Press, 2005.
- [189] A. WAGNER – « How the global structure of protein interaction networks evolves. », *Proc Biol Sci* **270** (2003), no. 1514, p. 457–466.
- [190] — , « Circuit topology and the evolution of robustness in two-gene circadian oscillators. », *Proc Natl Acad Sci U S A* **102** (2005), no. 33, p. 11775–11780.

- [191] — , « Robustness and evolvability in living systems. », ch. The many ways of building the same body, p. 173–191, in (188), 2005.
- [192] — , « Robustness and evolvability in living systems. », ch. RNA structure, p. 39–61, in (188), 2005.
- [193] — , « Robustness and evolvability : a paradox resolved. », *Proc Biol Sci* **275** (2008), no. 1630, p. 91–100.
- [194] G. P. WAGNER et P. F. STADLER – « Quasi-independence, homology and the unity of type : a topological theory of characters. », *J Theor Biol* **220** (2003), no. 4, p. 505–527.
- [195] G. P. WAGNER et L. ALTENBERG – « Complex adaptations and the evolution of evolvability. », *Evolution* **50** (1996), no. 3, p. 967–976.
- [196] W. Y. S. WANG, B. J. BARRATT, D. G. CLAYTON et J. A. TODD – « Genome-wide association studies : theoretical and practical concerns. », *Nat Rev Genet* **6** (2005), no. 2, p. 109–118.
- [197] Z. WANG et J. ZHANG – « In search of the biological significance of modular structures in protein networks. », *PLoS Comput Biol* **3** (2007), no. 6, p. e107.
- [198] R. A. WATSON – « Analysis of recombinative algorithms on a non-separable building-block problem », *Foundations of Genetic Algorithms* (W. Martin et W. Spears, éd.), Lecture Notes in Computer Sciences, vol. 6, Morgan Kaufmann, 2001, p. 69–89.
- [199] — , « Compositional evolution : Interdisciplinary investigations in evolvability, modularity, and symbiosis. », Thèse, Brandeis University, 2002.
- [200] D. J. WATTS et S. H. STROGATZ – « Collective dynamics of “small-world” networks. », *Nature* **393** (1998), no. 6684, p. 440–442.
- [201] R. WIJAYATUNGA, B. CORY et M. SHORL – « Security and revenue reconciliation in optimal transmission pricing », *IEE proceedings in Generation, Transmission and Distribution* **146** (1999), p. 355–359.
- [202] C. O. WILKE et C. ADAMI – « Evolution of mutational robustness. », *Mutat Res* **522** (2003), no. 1-2, p. 3–11.
- [203] C. R. WOESE – « On the evolution of cells. », *Proc Natl Acad Sci U S A* **99** (2002), no. 13, p. 8742–8747.
- [204] D. M. WOLF et A. P. ARKIN – « Motifs, modules and games in bacteria. », *Curr Opin Microbiol* **6** (2003), no. 2, p. 125–134.

- [205] S. WUCHTY, Z. N. OLTVAI et A.-L. BARABÁSI – « Evolutionary conservation of motif constituents in the yeast protein interaction network. », *Nat Genet* **35** (2003), no. 2, p. 176–179.
- [206] Y. XIA et M. LEVITT – « Roles of mutation and recombination in the evolution of protein thermodynamics. », *Proc Natl Acad Sci U S A* **99** (2002), no. 16, p. 10382–10387.
- [207] A. S. YANG – « Modularity, evolvability, and adaptative radiations : a comparison of hemi- and holometabolous insects. », *Evolution and Development* **3** (2001), p. 59–72.
- [208] W. ZHANG, A. COLLINS, J. GIBSON, W. J. TAPPER, S. HUNT, P. DELOUKAS, D. R. BENTLEY et N. E. MORTON – « Impact of population structure, effective bottleneck time, and allele frequency on linkage disequilibrium maps. », *Proc Natl Acad Sci U S A* **101** (2004), no. 52, p. 18075–18080.
- [209] Y.-X. ZHANG, K. PERRY, V. A. VINCI, K. POWELL, W. P. C. STEMMER et S. B. DEL CARDAYRÉ – « Genome shuffling leads to rapid phenotypic improvement in bacteria. », *Nature* **415** (2002), no. 6872, p. 644–646.
- [210] E. ZIV, M. MIDDENDORF et C. H. WIGGINS – « Information-theoretic approach to network modularity. », *Phys Rev E Stat Nonlin Soft Matter Phys* **71** (2005), no. 4 Pt 2, p. 046117.
- [211] E. ZOTENKO, J. MESTRE, D. P. O’LEARY et T. M. PRZYTYCKA – « Why do hubs in the yeast protein interaction network tend to be essential : reexamining the connection between the network topology and essentiality. », *PLoS Comput Biol* **4** (2008), no. 8, p. e1000140.