



HAL
open science

Estimateurs de biodiversité robustes pour le métabarcoding

Sylvain Moinard

► **To cite this version:**

Sylvain Moinard. Estimateurs de biodiversité robustes pour le métabarcoding. Biodiversité. Université Grenoble Alpes [2020-..], 2023. Français. NNT : 2023GRALV107 . tel-04624750

HAL Id: tel-04624750

<https://theses.hal.science/tel-04624750>

Submitted on 25 Jun 2024

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.

THÈSE

Pour obtenir le grade de

DOCTEUR DE L'UNIVERSITÉ GRENOBLE ALPES

École doctorale : CSV- Chimie et Sciences du Vivant

Spécialité : Biodiversité-Ecologie-Environnement

Unité de recherche : Laboratoire d'ECologie Alpine

Estimateurs de biodiversité robustes pour le métabarcoding

Robust biodiversity estimators for metabarcoding

Présentée par :

Sylvain MOINARD

Direction de thèse :

Eric COISSAC
Université Grenoble Alpes

Directeur de thèse

Christelle GONINDARD
Université Grenoble Alpes

Co-directrice de thèse

Rapporteurs :

Nigel Gilles YOCCOZ
FULL PROFESSOR, Universitetet i Tromsø

Guillaume ACHAZ
PROFESSEUR DES UNIVERSITES, Université Paris Cité

Thèse soutenue publiquement le **13 décembre 2023**, devant le jury composé de :

Eric COISSAC
MAITRE DE CONFERENCES HDR, Université Grenoble Alpes

Directeur de thèse

Christelle GONINDARD
MAITRESSE DE CONFERENCES HDR, Université Grenoble Alpes

Co-directrice de thèse

Nigel Gilles YOCCOZ
FULL PROFESSOR, Universitetet i Tromsø

Rapporteur

Guillaume ACHAZ
PROFESSEUR DES UNIVERSITES, Université Paris Cité

Rapporteur

François POMPANON
PROFESSEUR DES UNIVERSITES, Université Grenoble Alpes

Président

Lucie ZINGER
MAITRESSE DE CONFERENCES, ENS Paris

Examinatrice

Purificación LOPEZ-GARCIA
DIRECTRICE DE RECHERCHE, CNRS délégation Ile-de-France Sud

Examinatrice

Invités :

Didier Piau
PROFESSEUR, Université Grenoble Alpes



Avant-propos

J'ai réalisé ma thèse d'octobre 2020 à décembre 2023 au Laboratoire d'Écologie Alpine (LECA) à l'Université Grenoble-Alpes, au sein de l'École Doctorale de Chimie et des Sciences du Vivant. Elle a été financée par une bourse AMX de l'École polytechnique obtenue en 2020.

Ma thèse s'inscrit à l'interface entre l'écologie et les mathématiques appliquées. J'ai suivi le Cycle ingénieur de l'École polytechnique puis ai intégré le Master Mathématiques pour les Sciences du Vivant de l'Université Paris-Saclay. J'ai acquis dans ces deux cursus mes premières connaissances en écologie. J'ai effectué mon stage de Master 2 au LECA dans mon équipe de thèse.

Durant ma thèse, j'ai dispensé 135 heures de Travaux Dirigés en Licence Biotechnologie pour la santé (UFR de Pharmacie) en biostatistiques, informatique et mathématiques. J'ai encadré deux stages de Licence 3 et un projet de recherche bibliographique de Master 1. J'ai également été représentant élu des doctorants et post-doctorants du LECA de 2022 à 2023.

Résumé

Du fait de l'activité humaine, la biodiversité connaît un bouleversement rapide à l'échelle mondiale : extinction d'espèces, diminution des populations, introduction d'espèces invasives... Pour mettre en place des mesures de conservation, il faut être capable d'évaluer l'état de la biodiversité sur un site donné. De nombreux indicateurs existent mais l'observation directe des espèces pour les calculer est laborieuse. L'étude de l'ADN environnemental par le métabarcoding permet de dépasser cette limite et fournit une source d'information prometteuse pour améliorer la gestion des écosystèmes. Cependant, l'estimation des abondances relatives des espèces est encore mal établie pour les données de métabarcoding à cause de biais introduits au cours de l'expérience et de l'incertitude de l'identification des espèces présentes.

Pendant ma thèse, j'ai cherché à améliorer les aspects quantitatifs du métabarcoding sur le plan expérimental tout en développant de nouveaux outils théoriques d'analyse. Ma thèse comporte trois chapitres. Le premier est consacré à un nouvel algorithme d'inférence de paramètres pour modèles aléatoires appelé *Fixed Landscape Inference MethOd* (*flimo*). Il est applicable à une large gamme de modèles utilisés en écologie, pour le métabarcoding mais aussi en génétique ou en dynamique des populations. Il fonctionne dans le même cadre que les algorithmes d'*Approximate Bayesian Computation* (ABC) en procédant par simulations du modèle sans considérer sa vraisemblance. Sur les exemples étudiés, les résultats de *flimo* sont obtenus beaucoup plus vite que pour les algorithmes utilisés en comparaison, avec une précision similaire. Le deuxième chapitre présente mes travaux sur le métabarcoding quantitatif. Deux des biais du processus ont été mesurés expérimentalement et je propose un protocole simple pour les corriger. Le premier biais est lié à la variabilité interspécifique du nombre de copies des portions d'ADN ciblées par le métabarcoding et a été mesuré par digital droplet PCR. Le second biais est dû à l'amplification plus ou moins efficace des différents marqueurs au cours de la PCR. Ce biais a été mesuré à la fois par qPCR Taqman et retrouvé numériquement à partir d'un mélange d'ADN de composition connue. Ces travaux sont appuyés par le développement d'un modèle mathématique de PCR qui décrit plus précisément le phénomène que les modèles couramment utilisés en métabarcoding. Le troisième chapitre aborde le sujet de l'attribution des séquences observées dans les données de métabarcoding à des espèces réelles. Je montre certaines limites des algorithmes communs de traitement des séquences et propose des pistes d'amélioration pour un algorithme existant, *obiclean*, qui trie les séquences sur la base d'un graphe de mutation. Les séquences sont attribuées aux espèces de manière probabiliste sur la base de la vraisemblance d'un modèle d'erreur et en observant les cooccurrences des différents variants au sein d'un jeu de données avec de multiples échantillons indépendants. Les mesures de biodiversité sont ensuite adaptées à ce contexte aléatoire.

Mes travaux ont apporté plusieurs contributions à la recherche en écologie. J'ai émis plusieurs idées nouvelles pour interpréter les données de métabarcoding, en affinant la compréhension des processus qui affectent leur qualité et en proposant des solutions réalistes pour les analyser. Mes travaux méthodologiques de modélisation et de développement algorithmique me semblent aussi importants pour améliorer les outils statistiques utilisés en écologie.

Abstract

Due to human activity, biodiversity is undergoing rapid change on a global scale : species extinction, population decline, introduction of invasive species... To implement conservation measures, we need to be able to assess the state of biodiversity on a given site. Numerous indicators exist, but direct observation of species to compute them is laborious. The analysis of environmental DNA using metabarcoding overcomes this limitation and provides a promising source of information to improve ecosystem monitoring. However, the estimation of relative abundances of species is still poorly established for metabarcoding data, due to biases introduced during the experiment and to the uncertainty of identifying the species present.

During my PhD, I aimed to improve the quantitative aspects of metabarcoding both experimentally and by developing new theoretical tools for analysis. My PhD thesis is divided into three chapters. The first chapter is devoted to a new parameter inference algorithm for random models, called *Fixed Landscape Inference MethOd* (*flimo*). It is applicable to a wide range of models used in ecology, for metabarcoding but also in population genetics or population dynamics. It operates within the same framework as Approximate Bayesian Computation (ABC) algorithms, by simulating the model without considering its likelihood. On the examples studied, *flimo* results are obtained much faster than for the algorithms used in comparison, with similar accuracy. The second chapter presents my work on quantitative metabarcoding. Two biases in the process have been measured experimentally, and I propose a simple protocol to correct them. The first bias is linked to interspecific variability in the copy number of DNA portions targeted by metabarcoding, and was measured by digital droplet PCR. The second bias is due to the more or less efficient amplification of the species barcodes during PCR. This bias was measured both by Taqman qPCR and recovered numerically from a DNA mixture of known composition. This work is supported by the development of a mathematical model of PCR that describes more precisely the process than the models commonly used in metabarcoding. The third chapter deals with the attribution of sequences observed in metabarcoding data to real species. I show some limitations of sequence processing algorithms and propose improvements to an existing algorithm, *obiclean*, which classifies sequences on the basis of a mutation graph. Sequences are assigned to species probabilistically on the basis of the likelihood of an error model and by observing cooccurrences of different variants within a dataset with multiple independent samples. Biodiversity measurements are then adapted to this stochastic context.

My work has brought several contributions to ecological research. I have put forward new ideas for interpreting metabarcoding data, refining our understanding of the processes affecting their quality and proposing realistic solutions for analyzing them. My methodological work in modeling and algorithmic development is also important for improving the statistical tools used in ecology.

Liste des abréviations

Entre parenthèses figure l’acronyme en anglais lorsque celui-ci est utilisé dans les manuscrits (Chapitres 1 et 2).

- ABC : *Approximate Bayesian Computation*
- ADN (DNA) : Acide DésoxyriboNucléique
- ADNe (eDNA) : ADN environnemental
- ARN : Acide RiboNucléique
- ASV : *Amplicon Sequence Variant*
- bp : *base pair*
- (B)SL : *(Bayesian) Synthetic Likelihood*
- ddPCR : digital droplet PCR
- dNTP : désoxyriboNucléotide TriPhosphate
- EM : Espérance-Maximisation
- flimo : *Fixed Landscape Inference MethOd*
- HMM : *Hidden Markov Model*
- iid : indépendant et identiquement distribué
- LAMP : *Loop-mediated isothermal AMPlification*
- MCMC : *Markov Chain Monte Carlo*
- (M)OTU : *(Molecular) Operational Taxonomic Unit*
- PCR : *Polymerase Chain Reaction*
- qPCR : PCR quantitative
- RRA : *Raw Reads Abundances*
- RFU : *Relative Fluorescence Unit*
- RMSE : *Root-Mean-Square Error*
- SAEM : *Stochastic Approximation Expectation Maximisation*
- SMC : *Sequential Monte Carlo*

Liste des notations mathématiques

Les notations utilisées de manière récurrente dans le manuscrit sont consignées ici :

- a, b : paramètres de la loi Beta dans le Chapitre 3
- C_t : cycle de franchissement du seuil (*threshold*) en qPCR
- qD : nombre de Hill d'ordre q
- F_n : fluorescence au cycle n en qPCR
- \mathcal{F}_n : filtration d'une chaîne de Markov (information du système à l'étape n)
- qH : entropie HCDT d'ordre q
- J : fonction objectif à minimiser
- K : capacité de charge, quantité maximale de molécules synthétisables dans un milieu réactionnel
- l, l_θ : log-vraisemblance d'un modèle
- M_n^s : nombre de molécules de l'espèce s au cycle n au cours de la PCR
- \mathcal{M} : communauté artificielle (*mock community*) d'ADN environnemental
- N_v : nombre de variants observés dans des données de métabarcoding
- p, p_θ : vraisemblance d'un modèle
- q : paramètre d'importance des espèces rares pour l'entropie HCDT et les nombres de Hill
- R_s : nombre de lectures (*reads*) de l'espèce s après le séquençage
- S : nombre d'espèces d'une communauté
- α, β, γ : échelles de la biodiversité
- θ, Θ : paramètres et espace des paramètres des modèles étudiés
- λ_n^s, Λ_s : efficacité d'amplification PCR de l'espèce s apparente au cycle n et intrinsèque
- μ : taux de mutation
- σ : écart-type

Remerciements

Avant tout, merci à mon jury de thèse : Guillaume Achaz, Purificación López-García, François Pompanon, Gilles Yoccoz et Lucie Zinger d'avoir accepté d'évaluer mes travaux.

Bien sûr, un grand merci à mes encadrants de m'avoir accompagné depuis plus de trois ans dans cette grande aventure. Je vois derrière moi de beaux projets que nous avons menés à bien malgré les embûches.

Merci à Éric pour votre confiance et la richesse de vos idées qui commencent par "je dis peut-être une bêtise mais..." et sont souvent brillantes.

Merci à Christelle pour ta lucidité, ton optimisme et aussi le plaisir de pouvoir râler sur nos petits Biotech. Tu as toujours veillé à ce que ma thèse garde le bon cap.

Merci aussi à Didier Piau pour l'attention précise et exigeante que vous avez accordée à mon travail et pour votre soutien durant les longs premiers mois de ma thèse.

Merci à Frédéric Boyer pour ta sympathie, ton suivi au début de ma thèse et pour avoir nourri mes réflexions scientifiques de tes nombreux questionnements.

Merci ensuite aux membres de mon Comité de Suivi Individuel : Adeline Leclercq-Samson, Eric Marcon, François Munoz et Pierre Pudlo. Nos échanges ont été fructueux et ont marqué des étapes importantes pour que je reste motivé et aille de l'avant dans ma thèse.

Merci à Frédéric Laporte pour ta bonne humeur et ton indéfectible soutien pendant nos nombreuses heures de PCR. L'année prochaine, vous gagnez le tournoi interlabo !

Merci à celles et ceux qui font vivre le LECA et dont le soutien m'a été précieux, en particulier Delphine Rioux pour les manips, Christian Miquel et Florence Sagnimorte sans qui la logistique et l'administratif n'auraient pas été aussi simples.

Merci à Édouard Oudet, Pierre Taberlet et Stefaniya Kamenova, entre autres, pour nos interactions scientifiques stimulantes qui ont marqué mon entrée dans le monde de la recherche.

J'ai aussi une pensée chaleureuse pour des professeurs marquants, en particulier Alex Petitcolin et Denis Choimet. Vous m'avez apporté, au-delà du goût des mathématiques, l'envie et l'énergie de découvrir ce qui se cache plus loin.

Je me dois bien sûr de remercier les ministres de Darwinie, des Nodes ou d'ailleurs.

Merci à Rémy, camarade depuis la boue de la Courtine, berger philosophe et esprit frappeur des alpages et des couloirs du LECA (sûrement une âme volée à un clown).

Merci à Tiphaine, complice de la première heure, de m'avoir supporté jusque-là malgré un catapultage à Sassenage dont le récit légendaire reste encore sujet à caution.

Merci à Marie, infatigable Pharmadoc², grâce à qui j'ai été promu Éminent Professeur Junior... Les feuilletés du pot te sont dédiés! (too soon?)

Merci à Laura, montagnarde survoltée et esprit libre, la preuve en personne qu'on n'est pas vieux à 32 ans (c'est mon manuscrit, j'ai le droit d'être taquin).

Merci à Alexis, mécano-nageur et imitateur hors pair, fin punchlineur malgré le sur-apprentissage sur une filmographie douteuse.

Merci à Stephen, Dr Aculeca selon des sources occultes, pour ton humour corrosif et ta vigilance à ce que les petits doctorants survivent aux affres de la thèse.

Merci à Ségolène, nos croisements à vélo sont une horloge infaillible pour mesurer mon retard le matin. Merci à Gaspard le musico-tronicien au succès annoncé à la prochaine Eurovision! Et merci à Marion la botaniste aux cheveux roses (un jour peut-être?). En outre, je crois m'y être engagé :

Louis Majesté des moustiques
Gratitude d'un voisin
Artifice post-prandial

Merci aussi à Camille à qui je lègue le meilleur bureau du labo, Matthias et Marianne les roublards des buffets, à Matthieu, Manu, Marie, Naiara, Louise, Caroline, Titouan, Clara, Matthew, Lara...

Je réserve le dernier mot à mes proches, bien présents à mes côtés. Merci à ma famille, à Lucie, à mes amis de toujours : Eloi, Baptiste, Guillaume, Nicolas, Valentin, encore Valentin, Héloïse, Jeanne... Vous savez que je vous dois beaucoup. Ces trois années, celles d'avant et toutes celles à venir n'auraient pas la même saveur sans vous!

Table des matières

Introduction générale	1
0.1 Pourquoi mesurer la biodiversité?	1
0.1.1 Biodiversité à l'ère de l'Anthropocène	1
0.1.2 Définitions de la biodiversité	2
0.1.3 Le problème de l'espèce	3
0.2 Indices de biodiversité	3
0.2.1 Composantes de la biodiversité	3
0.2.2 Niveaux d'étude	4
0.2.3 Notations	5
0.2.4 Mesures de la diversité neutre α et γ	5
0.2.5 Entropie et nombres de Hill	6
0.2.6 Spectre de diversité	7
0.2.7 Diversité β	7
0.2.8 Décomposition de la diversité neutre	8
0.2.9 Mesures de biodiversité en pratique	8
0.3 Métabarcoding et ADN environnemental	10
0.3.1 ADN environnemental	10
0.3.2 Émergence du métabarcoding	11
0.3.3 Applications du métabarcoding et de l'ADNe	11
0.3.4 Protocole du métabarcoding	12
0.3.5 Interprétation des données	16
0.4 Réaction en Chaîne par Polymérase (PCR)	16
0.4.1 Procédé	17
0.4.2 Mécanismes de saturation au cours de la PCR	19
0.4.3 Applications de la PCR	19
0.4.4 Techniques de PCR quantitatives	19
0.4.5 Inhibition de la PCR	23
0.5 Modélisation de la PCR	24
0.5.1 Notations	24
0.5.2 Intérêt et limites du modèle exponentiel	25
0.5.3 Modèles analytiques	25
0.5.4 Modèles mécanistiques	26
0.5.5 Modèles aléatoires	29
0.5.6 Limites de la modélisation	30
0.5.7 Représentation de quelques modèles	30

0.6	Métabarcoding quantitatif	32
0.6.1	Enjeux	32
0.6.2	Biais successifs	32
0.6.3	Biais biologiques	33
0.6.4	Biais techniques	34
0.6.5	Méthodes de correction	36
0.7	Limites et perspectives du métabarcoding	39
0.7.1	Détection d'espèces	39
0.7.2	Détermination des communautés	39
0.7.3	Alternatives au métabarcoding	40
0.8	Outils mathématiques en écologie et biologie	42
0.8.1	Motivations	42
0.8.2	Inférence de paramètres pour modèles stochastiques	44
0.9	Plan de la thèse	52
1	<i>Fixed Landscape Inference MethOd (flimo)</i>	53
1.1	Introduction	53
1.1.1	Contexte de développement de <i>flimo</i>	53
1.1.2	Simulation de variables aléatoires	54
1.1.3	Couplage de variables aléatoires	58
1.1.4	Algorithmes d'optimisation déterministe	60
1.2	Résumé en langue française	68
1.3	Manuscrit	68
1.4	Résultats complémentaires	90
1.4.1	Couplage entre <i>flimo</i> et un algorithme EM ?	90
1.5	Conclusion	90
2	Métabarcoding quantitatif	91
2.1	Introduction	91
2.2	Résumé en langue française	92
2.3	Manuscrit	92
2.4	Résultats complémentaires	116
2.4.1	Modélisation mécanistique de PCR	116
2.4.2	Modélisation de l'amplification par PCR avec mismatch	120
2.4.3	Mesure de la compétition par un contrôle interne	126
2.4.4	Mesure de la concentration d'ADN cible selon le tissu chez plusieurs espèces de plantes arctiques	130
2.5	Conclusion	130
3	Assigination probabiliste des séquences	131
3.1	Introduction	131
3.1.1	Variabilité des séquences	132
3.1.2	Modèles de mutation	134
3.1.3	Construction d'unités taxonomiques	136
3.1.4	Algorithmes de regroupement des séquences	137
3.1.5	Détection des chimères	140

3.1.6	Algorithmes de traitement post-clustering	140
3.1.7	Plan du chapitre	140
3.2	Un premier modèle de mutation pendant la PCR	141
3.2.1	Développement du modèle	141
3.2.2	Résultats	141
3.3	Assignation taxonomique : Matériel et Méthodes	142
3.3.1	Données étudiées	142
3.3.2	Étude des mutations ponctuelles	143
3.3.3	Cooccurrences des variants	146
3.3.4	Matrice de probabilité d'assignation	148
3.3.5	Adaptation des indices de biodiversité	150
3.4	Assignation taxonomique : Résultats	152
3.4.1	Mutations ponctuelles	152
3.4.2	Ratio de mutants	152
3.4.3	Cooccurrences des variants	155
3.4.4	Similarité génétique des séquences cooccurrentes	159
3.4.5	Matrice de probabilité	159
3.4.6	Estimation de la diversité α	160
3.5	Collaboration : Géotypage de loups	162
3.6	Conclusion	162
4	Discussion générale et perspectives	165
4.1	<i>Flimo</i> : <i>Fixed Landscape Inference MethOd</i> (Chapitre 1)	166
4.1.1	Bilan de l'étude	166
4.1.2	Perspectives	167
4.2	Métabarcoding quantitatif (Chapitre 2)	168
4.2.1	Bilan de l'étude	168
4.2.2	Perspectives	170
4.3	Assignation probabiliste des séquences (Chapitre 3)	173
4.3.1	Bilan de l'étude	173
4.3.2	Perspectives	174
4.4	Indices de biodiversité	176
4.5	Conclusion	177
	Bibliographie	178
A	Manuscrit <i>The Fixed Landscape Inference MethOd</i>	203
B	Manuscrit <i>Quantitative metabarcoding</i>	211
C	Modèle de PCR avec mismatch (Chapitre 2)	215
D	Calculs du modèle de mutation (Chapitre 3)	217
E	Marqueurs pour métabarcoding	219

Introduction générale

La conservation de la biodiversité est un enjeu fondamental du XXI^e siècle. Les changements rapides que connaît la Terre depuis 150 ans nécessitent une action urgente et globale pour préserver le patrimoine naturel de la planète. Pour protéger la biodiversité, il est crucial de bien la connaître et la comprendre. Mes travaux s’inscrivent dans ce vaste projet. Au cours de ma thèse, j’ai cherché à améliorer les outils de mesure de la biodiversité à partir de l’ADN environnemental analysé par métabarcoding. Mes travaux ont aussi donné lieu à des développements méthodologiques en statistique. Avant de présenter mes différents projets, je dresse dans cette introduction un état de l’art des sujets qui m’ont intéressé pendant les trois années de ma thèse. J’aborde deux thèmes principaux : d’une part, les mesures de biodiversité en lien avec l’ADN environnemental étudié par le métabarcoding ; d’autre part, les outils mathématiques utilisés en écologie à un sens plus large.

0.1 Pourquoi mesurer la biodiversité ?

0.1.1 Biodiversité à l’ère de l’Anthropocène

Depuis la révolution industrielle du XIX^e siècle, la Terre connaît des changements radicaux dont le moteur principal est l’action humaine. Ces bouleversements marquent une ère nouvelle caractérisée entre autres par un changement climatique global et rapide et une modification profonde des écosystèmes entraînant une extinction de masse d’espèces. Bien que la notion d’Anthropocène en tant qu’ère géologique fasse débat¹, elle illustre un tournant dans la place qu’occupe l’humanité dans son environnement.

L’état de la biodiversité est documenté dans le rapport de 2019 de l’IPBES (*Intergovernmental Science-Policy Platform on Biodiversity and Ecosystem Services*) (IPBES, 2019)². L’IPBES un groupe d’experts de la biodiversité à l’interface entre sciences et politique dont le but est de renforcer les actions de conservation dans le monde entier. Une de ses observations est que “dans la plupart des régions du monde, la nature a aujourd’hui été altérée de manière significative par de multiples facteurs humains, et la grande majorité des indicateurs relatifs aux écosystèmes et à la biodiversité montrent

1. <http://quaternary.stratigraphy.org/working-groups/anthropocene/>
2. https://www.ipbes.net/sites/default/files/2020-02/ipbes_global_assessment_report_summary_for_policymakers_fr.pdf

un déclin rapide.” De nombreuses mesures de la biodiversité viennent appuyer cette conclusion et montrent une accélération du déclin depuis 1970. La Figure 0.1.1 illustre l'érosion de la biodiversité à l'échelle mondiale.

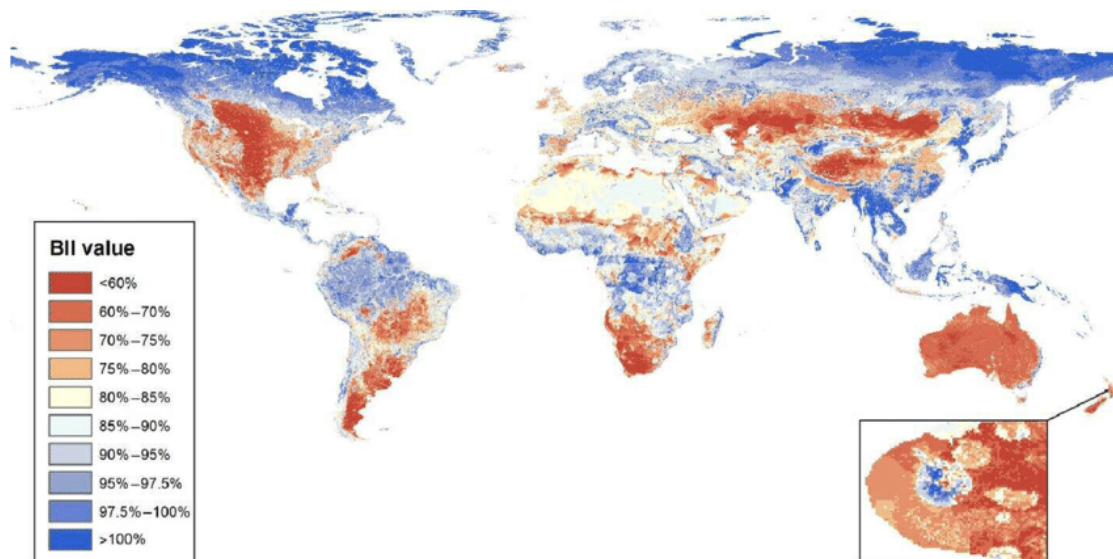


FIGURE 0.1.1 – Indice d'intégrité de la biodiversité (*Biodiversity Intactness Index*, BII) estimé par Newbold et al. (2016). Cet indice désigne la part restante aujourd'hui de la biodiversité naturelle d'origine du site. La zone en bas à droite est une réserve naturelle de Nouvelle-Zélande où les mesures de conservation ont des conséquences visibles. Source de la figure : Purvis et al. (2018).

Une des conséquences de cette perte de biodiversité est la dégradation des contributions de la nature aux populations. Celles-ci regroupent l'ensemble des apports matériels et immatériels nécessaires à une bonne qualité de vie comme les ressources, la régulation du climat ou encore le bien-être psychologique. Là encore, ces contributions sont quantifiées dans le rapport de l'IPBES par un grand nombre d'indicateurs.

0.1.2 Définitions de la biodiversité

Ces enjeux nécessitent de définir la biodiversité pour préciser les objectifs de conservation. Le terme "biodiversité", contraction de "diversité biologique", a été formulé en 1986 par Walter Rosen (Wilson, 1999). La biodiversité désigne donc la variabilité du vivant à différentes échelles dont les plus communes sont les écosystèmes, les espèces et les gènes. De très nombreuses définitions existent, la plus restrictive étant le nombre d'espèces différentes vivant dans un écosystème donné. Marc Loreau propose une définition (Barbault et al., 2005) que je trouve pertinente et complète :

“La Terre abrite une extraordinaire diversité biologique, qui inclut non seulement les espèces qui habitent notre planète, mais aussi la diversité de leurs gènes, la multitude des interactions écologiques entre elles et avec leur environnement physique, et la variété des écosystèmes complexes qu'elles constituent. Cette biodiversité, qui est le produit de plus de 3 milliards d'années

d'évolution, constitue un patrimoine naturel et une ressource vitale dont l'humanité dépend de multiples façons."

0.1.3 Le problème de l'espèce

Une difficulté émerge avant de mesurer la biodiversité. L'espèce est l'unité de base pour caractériser le vivant, mais ses contours sont flous (Hey, 2001). Ce caractère élémentaire n'a d'ailleurs rien d'évident chez certains groupes comme les micro-organismes. Mayden (1997) recense vingt-deux définitions qui mènent à des décomptes d'espèces variables. Une définition biologique générale est donnée par Mayr (1942) :

"Une espèce est un groupe de populations naturelles au sein duquel les individus peuvent, réellement ou potentiellement, échanger du matériel génétique ; toute espèce est séparée des autres par des mécanismes d'isolement reproductif."

Ce sont les notions d'interfécondabilité et d'isolement reproductif qui caractérisent ici l'espèce. Mais d'autres approches sont possibles, par exemple une définition phylogénétique de Cracraft (1983) :

"le plus petit ensemble d'organismes individuels contenant leurs ancêtres et leurs descendants et ayant en commun une diagnose qui les différencie des autres ensembles."

Cette définition est reformulée par Wheeler (1999) :

"le plus petit agrégat de populations (sexué) ou de lignées (asexué) diagnostiqué par une combinaison unique d'états de caractères chez des individus comparables."

Ce sont alors les caractères communs qui permettent de caractériser l'espèce.

J'introduirai plus loin la notion d'Unité Taxonomique Opérationnelle (OTU), qui est une approximation pragmatique de l'espèce selon différents critères de similarité. Les OTU permettent de regrouper les observations afin de calculer des indices de biodiversité, notamment pour l'étude de l'ADN environnemental.

0.2 Indices de biodiversité

Il est bien sûr nécessaire de pouvoir mesurer la biodiversité pour mettre en place des mesures de conservation. Les outils de mesure visent à fournir un résumé simple de la composition d'un écosystème. Mes travaux sur les mesures de biodiversité ont été guidés par l'ouvrage très complet de Marcon (2015).

0.2.1 Composantes de la biodiversité

Trois composantes permettent de décrire la biodiversité de manière complète (Stirling, 2007). La première est la richesse, c'est-à-dire simplement le nombre d'espèces

différentes dans l'écosystème considéré.

Vient ensuite l'équitabilité qui prend en compte l'abondance relative de chaque espèce. La motivation est qu'une espèce très présente ne joue pas le même rôle dans la biodiversité qu'une espèce représentée par un seul individu. La diversité est considérée comme maximale quand les espèces sont toutes présentes dans les mêmes proportions. C'est une mesure "neutre" : cela ne signifie pas que l'écosystème est en meilleure santé dans ces conditions d'équilibre. Les mesures de biodiversité se traduisent souvent par un nombre "équivalent" (dans un sens à définir) d'espèces dans un écosystème où l'observation de chaque espèce est équiprobable.

La troisième composante est la diversité des espèces. La distance entre les espèces présentes est mesurée soit par leur proximité évolutive (diversité phylogénétique) (Webb et al., 2006), soit par leurs traits fonctionnels ou niches écologiques (diversité fonctionnelle) (Tilman, 1997). La diversité fonctionnelle est utile pour la conservation si l'on décide de se concentrer sur les groupes fonctionnels plutôt que sur les espèces individuellement (Walker, 1992). La diversité phylogénétique permet, elle, de lever le problème de l'espèce : regrouper ou non deux espèces proches ne change qu'à la marge la mesure de la biodiversité.

Au cours de ma thèse, je me suis concentré sur les deux premières composantes qui définissent la diversité neutre : chaque espèce est considérée comme intrinsèquement unique et équidistante de chacune autre. La Figure 0.2.2 illustre le rôle de la richesse et de l'équitabilité dans la perception de la biodiversité.

0.2.2 Niveaux d'étude

L'échelle spatiale varie en fonction du contexte d'étude. Trois niveaux de mesure se démarquent (Whittaker, 1960) :

- la diversité α mesure la diversité locale, à l'échelle d'une communauté et dans un habitat homogène ;
- la diversité β mesure la similarité entre des communautés locales ;
- la diversité γ mesure la diversité à l'échelle de la métacommunauté, c'est-à-dire d'une communauté de communautés.

À une échelle très globale, il peut être pertinent d'utiliser des *indicateurs* de biodiversité dans le sens donné par Balmford et al. (2003) : on ne cherche alors pas à décrire exhaustivement la biodiversité mais à se concentrer sur certains groupes typiques ou cruciaux des écosystèmes.

Mes travaux d'étude de l'ADN environnemental ont pour application principale la diversité α .

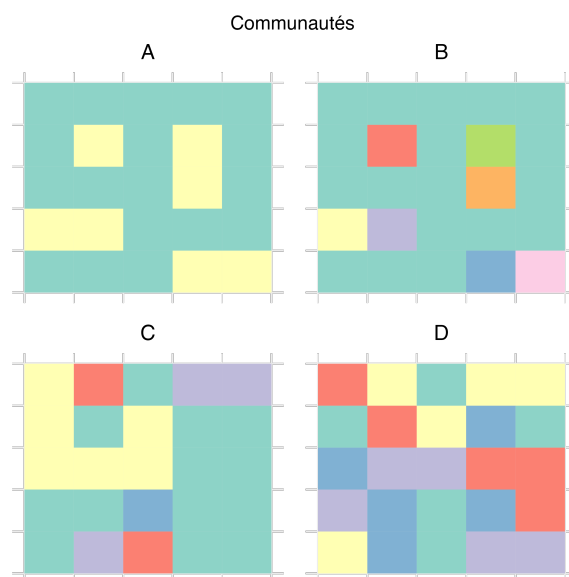


FIGURE 0.2.2 – Rôle de la richesse (en haut) et de l'équitabilité (en bas) dans la perception de la diversité pour quatre populations de 25 individus dont la couleur définit l'espèce. Toute chose étant égale par ailleurs, cette diversité semble plus riche à droite qu'à gauche : soit parce qu'il y a plus d'espèces différentes, soit parce que leurs abondances sont mieux réparties. En revanche, la communauté B est plus ou moins diverse que C et D selon le critère retenu. Figure inspirée de Marcon (2015).

0.2.3 Notations

On considère un écosystème (local ou global) comportant S espèces. Chaque espèce s représente une fraction p_s de l'ensemble des individus. Ces abondances relatives peuvent aussi être définies par unité de biomasse.

0.2.4 Mesures de la diversité neutre α et γ

Les diversités α et γ sont intrinsèques à un milieu donné, elles peuvent être définies de manière absolue. On peut les mesurer avec les mêmes indices, choisis pour être invariants par effet d'échelle. Voici trois exemples classiques (Broms et al., 2015) qui accordent un rôle différent aux abondances relatives.

- La richesse spécifique, tout simplement S . Cet indice, très simple, fait abstraction des abondances.
- L'indice de Shannon correspond à une mesure d'entropie : $H = -\sum_{s=1}^S p_s \ln p_s$. Cet indice donne le même poids à chaque observation (Marcon, 2015).
- L'indice de Gini-Simpson, $E = 1 - \sum_{s=1}^S p_s^2$, s'interprète (en population infinie) comme la probabilité que deux individus tirés au hasard ne soient pas de la même espèce.

0.2.5 Entropie et nombres de Hill

Ces trois indices peuvent être réunis dans un cadre commun par une mesure d'entropie empruntée à la théorie de l'information. L'entropie quantifie la biodiversité par la surprise générée lors d'une nouvelle observation. Elle s'exprime sous la forme générale suivante :

$$\sum_{s=1}^S p_s I(p_s) \quad (1)$$

où $p_s \mapsto I(p_s)$ est une fonction d'information définie sur $[0, 1]$ qui vérifie :

- I est décroissante : la surprise augmente lorsqu'on observe une espèce plus rare ;
- $\lim_{q \rightarrow 0} I(q) = +\infty$: observer une espèce rare apporte une grande information ;
- $I(1) = 0$: une observation certaine n'apporte aucune information.

Dans ce formalisme, on retrouve la richesse spécifique (moins un) et les indices de Shannon et Gini-Simpson pour $I(p_s) = \frac{1-p_s}{p_s}$; $I(p_s) = -\ln(p_s)$ et $I(p_s) = 1 - p_s$ respectivement.

Une forme commune est l'entropie HCDDT, du nom de ses (ré)inventeurs successifs (Havrda et al., 1967; Daróczy, 1970; Tsallis, 1988), définie par :

$${}^q H = \frac{1}{1-q} \left(1 - \sum_{s=1}^S p_s^q \right) \quad (2)$$

pour un paramètre $q \geq 0$ qui caractérise l'importance des abondances relatives : plus q est grand, moins les espèces rares sont prises en compte. On retrouve aisément les indices classiques : $S - 1 = {}^0 H$, $H = {}^1 H$ (par continuité en $q = 1$), $E = {}^2 H$.

On aimerait donner une interprétation facile des valeurs de ces indices de biodiversité. Pour cela, un choix consensuel (Ellison, 2010) consiste à transformer l'entropie HCDDT en nombres de Hill (Hill, 1973; Chao et al., 2014), définis par :

$${}^q D = \left(\sum_{s=1}^S p_s^q \right)^{\frac{1}{1-q}} \quad (3)$$

Ceux-ci sont interprétables comme "le nombre d'espèces équiprobables dont l'entropie est la même que celle de la communauté étudiée" (Marcon, 2015). Une fois de plus, cette interprétation dépend de l'importance donnée aux espèces rares décrite par le paramètre q . Les indices classiques sont liés aux nombres de Hill :

$${}^0 D = S, \quad {}^1 D = e^H, \quad {}^2 D = \frac{1}{1-E}$$

Pour établir l'équivalence de l'entropie HCDT et des nombres de Hill, on définit un logarithme et une exponentielle déformés, réciproques l'un de l'autre (Tsallis, 1994) :

$$\ln_q(x) = \frac{x^{1-q} - 1}{1 - q} \quad (4)$$

$$e_q^x = (1 + (1 - q)x)^{\frac{1}{1-q}} \quad (5)$$

On a alors

$${}^qH = \ln_q({}^qD) \quad (6)$$

$${}^qD = e_q^{{}^qH} \quad (7)$$

0.2.6 Spectre de diversité

Afin de comparer des communautés, il est recommandé (Kindt et al., 2006; Leinster and Cobbold, 2012) d'observer le profil de diversité ou spectre de diversité α , c'est-à-dire la courbe ${}^qD = f(q)$. Celle-ci décroît de la richesse spécifique $S = {}^0D$ à ${}^{+\infty}D = \frac{1}{p_{\max}}$ où p_{\max} est la proportion de l'espèce la plus abondante³. Les spectres de diversité ne fournissent pas une relation d'ordre totale des communautés selon leur diversité. Des exemples de spectres de diversité de Hill sont donnés sur les Figures 3.4.12 et 3.4.13.

0.2.7 Diversité β

La diversité β mesure la dissimilarité entre des communautés locales. Il existe de nombreuses manières de définir cette dissimilarité qui ont mené à de multiples définitions de la diversité β (Anderson et al., 2011). Ce sujet n'est pas abordé dans ma thèse où je me suis concentré sur la diversité α . À titre d'exemple, je cite la dissimilarité de Bray-Curtis (Bray and Curtis, 1957) qui utilise les abondances observées $n_{s,k}$ (nombre d'individus de l'espèce s au site k) :

$$D = \frac{\sum_s |n_{s,1} - n_{s,2}|}{\sum_s (n_{s,1} + n_{s,2})} \quad (8)$$

Une généralisation est proposée par Anderson et al. (2006). Cette dissimilarité dite de Gower modifiée accorde la même importance à un changement d'un ordre de grandeur qu'au passage de zéro à un individu. Elle incorpore aussi une pondération w_s pour chaque espèce.

$$D_{MG} = \frac{\sum_s w_s |x_{s,1} - x_{s,2}|}{\sum_s w_s} \quad (9)$$

avec $x_{s,k} = \log_{10}(n_{s,k}) + 1$ si $n_{s,k} > 0$ et 0 sinon. Cette approche logarithmique est pertinente quand les disparités d'abondance sont grandes comme c'est le cas pour les données d'ADN environnemental.

3. Dans le cas où l'espèce s vérifiant $p_s = p_{\max}$ est unique.

0.2.8 Décomposition de la diversité neutre

Une fois les trois niveaux de diversité définis, il est possible de décomposer la diversité γ en une composante α (la diversité moyenne des sites locaux) et une composante β (la différence entre les sites), idéalement indépendantes. C'est un vaste sujet qui sort du cadre de ma thèse, mais je donne une décomposition dont je trouve l'interprétation élégante. Jost (2007) et Chao et al. (2012) montrent qu'une décomposition adaptée pour les nombres de Hill est multiplicative :

$${}^qD_\gamma = {}^qD_\beta {}^qD_\alpha \quad (10)$$

où ${}^qD_\alpha$ et ${}^qD_\gamma$ sont les nombres équivalents d'espèces présentés plus haut, et ${}^qD_\beta$ est un nombre équivalent de communautés totalement distinctes. Des décompositions additives existent aussi (Marcon et al., 2014).

0.2.9 Mesures de biodiversité en pratique

Mes travaux ne portent pas sur les méthodes traditionnelles de mesure de biodiversité. Je m'intéresse plutôt à l'ADN environnemental, présenté plus loin (Chapitres 2 et 3) et pour lequel l'analyse est conceptuellement différente.

0.2.9.1 Acquisition des données

Dans les études dites traditionnelles, il faut procéder à un échantillonnage des zones étudiées pour estimer quelles espèces sont présentes et en quelles proportions. Chez les animaux, l'échantillonnage consiste souvent en la capture des individus (pêche, piégeage...) qui sont ensuite identifiés et éventuellement marqués en vue d'une recapture ultérieure. Il est aussi possible d'observer les individus à distance par piège vidéo ou par des techniques acoustiques. Le choix des sites d'échantillonnage (aléatoire, stratifié, systématique...) dépend des objectifs de l'étude.

0.2.9.2 Effort d'échantillonnage

Naturellement, le nombre d'espèces observées et la précision d'estimation de leurs abondances dépendent de l'effort d'échantillonnage. Celui-ci peut être représenté par la courbe d'accumulation qui donne le nombre d'espèces observées selon le nombre d'individus observés ou la superficie de la zone étudiée, comme sur la Figure 0.2.3.

L'effort peut ensuite être quantifié par différents indices dont je donne deux exemples. Un premier est le taux de couverture (Good, 1953), défini comme la proportion des espèces découvertes par :

$$\sum_{s=1}^S \mathbb{1}_{n_s > 0} p_s \quad (11)$$

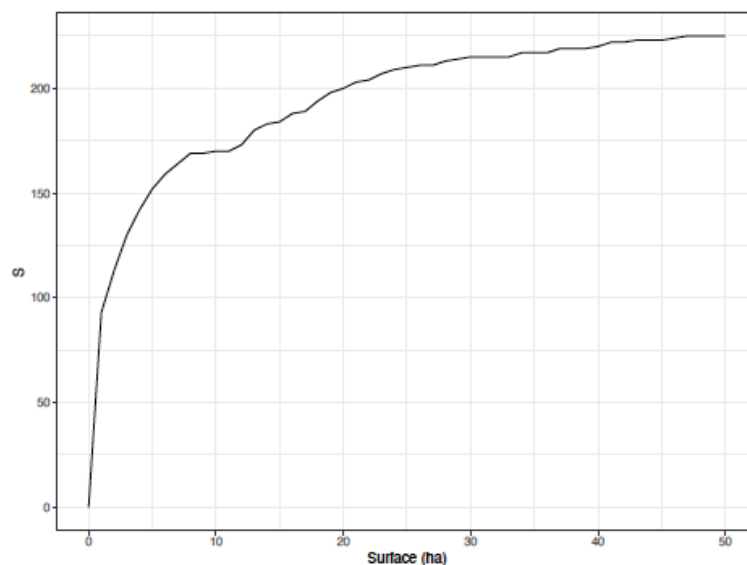


FIGURE 0.2.3 – Courbe d’accumulation des espèces d’arbres du dispositif de l’île de Barro Colorado (Panama), une parcelle forestière tropicale de 50 hectares (Condit et al., 2012). Source de la figure : Marcon (2015).

où n_s est le nombre d’observations de l’espèce s .

Un autre indice, le taux de complétude, ne prend pas en compte les abondances mais simplement le nombre d’espèces découvertes :

$$\sum_{s=1}^S \mathbb{1}_{n_s > 0} \quad (12)$$

La comparaison de jeux de données entre eux est un vrai enjeu car les mesures dépendent de l’effort d’échantillonnage. Il est commun de traiter les données par raréfaction (réduction aléatoire du jeu de données) et/ou extrapolation (prédiction d’observation avec un effort plus grand) pour prendre en compte le biais induit par des tailles de jeux de données différentes (Broms et al., 2015; Hsieh et al., 2016).

0.2.9.3 Correction des indices classiques

Le biais principal des mesures de biodiversité traditionnelles provient du risque de ne pas observer certaines espèces. Un grand nombre d’estimateurs empiriques ont été construits pour limiter ce biais (Marcon, 2015), notamment en utilisant le nombre d’espèces dont on a observé un ou deux individus pour estimer le nombre d’espèces non observées (Chao and Jost (2015) par exemple).

Par ailleurs, des méthodes de correction sont proposées lorsque l’attribution taxonomique est incertaine, par exemple dans les forêts tropicales où les noms vernaculaires ne

désignent pas toujours une espèce d'arbre unique (Guitet et al., 2014). Pour contourner cette difficulté, il est aussi possible de changer d'échelle taxonomique en établissant une richesse générique à la place d'une richesse spécifique (Williams and Gaston, 1994).

0.2.9.4 Limites classiques

Malgré ces corrections, les méthodes de suivi traditionnelles sont limitées. Elles sont chères (en temps de travail) et des différences méthodologiques affectent la comparaison de groupes d'espèces distincts (Pawlowski et al., 2020). Des erreurs d'identification sont aussi inévitables malgré l'expertise des naturalistes (Deiner et al., 2017). À titre d'exemple, Lawton et al. (1998) dresse un inventaire de huit groupes d'animaux d'une forêt tropicale au Cameroun. Les auteurs estiment que 10 000 heures de travail scientifique ont été nécessaires à l'inventaire de 2000 espèces et que celui-ci est très incomplet, avec entre 10 et 100 fois plus d'espèces présentes estimées.

0.3 Métabarcoding et ADN environnemental

De nouvelles techniques de biologie moléculaire ont émergé il y a une vingtaine d'années. Leur but est de remplacer l'observation directe des espèces par une analyse de l'information génétique qu'elles laissent dans leur environnement, et ainsi de contourner les limites des méthodes classiques. Cette information génétique est appelée l'ADN environnemental (ADNe). Alberdi and Gilbert (2019) atteste de la pertinence des données d'ADNe pour établir des mesures de biodiversité, en particulier pour calculer des nombres de Hill.

Le métabarcoding, une de ces techniques, est au cœur de ma thèse et fait partie des grandes thématiques étudiées au LECA. Je présente dans le détail son usage, ses applications et ses limites actuelles. Je fais ensuite un détour par une étape du protocole, l'amplification par PCR, sur laquelle j'ai concentré une partie de mes analyses. Enfin, je me penche sur le sujet principal qui a motivé ma thèse : les outils disponibles pour obtenir des données quantitatives à partir d'ADN environnemental.

0.3.1 ADN environnemental

L'ADN environnemental est défini par Taberlet et al. (2018) comme "un mélange complexe d'ADN génomique, provenant de nombreux organismes différents, trouvé dans un échantillon environnemental". Les auteurs distinguent deux types d'échantillons. Les échantillons environnementaux à proprement parler sont constitués de sol, de sédiment, d'eau, de fèces ou encore d'air Ruppert et al. (2019). Le second type désigne les échantillons "en vrac" (*bulk samples*) contenant des organismes entiers, par exemple des insectes entiers capturés par un piège Malaise (Elbrecht et al., 2021). L'ADN est généralement de bonne qualité dans les échantillons "en vrac" mais pas forcément dans les échantillons environnementaux. Ceux-ci contiennent de l'ADN intracellulaire provenant des cellules vivantes collectées dans l'échantillon ou de l'ADN extracellulaire de cellules mortes et soumises à un niveau de dégradation

variable. L'ADNe désigne des échantillons récents mais l'étude de l'ADN ancien repose sur les mêmes technologies. Les échantillons ADN ancien proviennent de substrats particuliers comme les sédiments de lac (Giguet-Covex et al., 2014) ou le permafrost gelé (Willerslev et al., 2014).

Les études de l'ADNe poursuivent deux buts : la détection d'une espèce ciblée et l'identification de l'ensemble des espèces à l'échelle d'une communauté (Seymour, 2019). Ces deux approches reposent sur l'utilisation de codes-barres moléculaires spécifiques à chaque espèce. Il s'agit d'une portion du génome avec peu de variabilité intraspécifique, de la variabilité interspécifique (même pour deux groupes taxonomiques proches). Elle doit être flanquée de deux portions d'ADN conservées autant que possible entre les différents groupes, de sorte à pouvoir amplifier cette partie précise du génome.

0.3.2 Émergence du métabarcoding

La première mention de l'ADNe est faite par Ogram et al. (1987). Elle est suivie par une première étude de biodiversité à partir d'ARN présent dans de l'eau de mer (Giovannoni et al., 1990). Le barcoding, qui tire son nom des codes-barres moléculaires, a d'abord été développé pour détecter une espèce ciblée (Hebert et al., 2003). Cette technique nécessite de l'ADN isolé et de bonne qualité, ce qui n'est pas toujours réaliste pour des échantillons d'ADNe. De plus, le barcoding ne permet pas de répondre à certaines questions écologiques, par exemple sur la composition de communautés. C'est pourquoi le métabarcoding a été développé (Taberlet et al., 2012), ou plus précisément l'*eDNA metabarcoding* (Pawlowski et al., 2020). Son essor a été rendu possible grâce au développement du séquençage haut-débit dit de nouvelle génération (Deiner et al., 2017). Grâce au métabarcoding, il n'est plus nécessaire d'observer directement ou d'isoler les espèces de la communauté à étudier.

La métagénomique permet aussi l'étude de l'ADNe en séquençant sans amplification (méthode dite *shotgun*) l'ensemble de l'ADNe. Je présente cette technique en 0.7.3. La Figure 0.3.4 illustre les trois techniques de biologie moléculaire de mesure de la biodiversité : le barcoding pour la détection d'espèces ciblées, le métabarcoding et la métagénomique pour analyser une communauté. J'en profite pour évoquer la métatranscriptomique qui étudie l'ARN trouvé dans un échantillon environnemental, à la manière de l'ADNe pour la métagénomique.

0.3.3 Applications du métabarcoding et de l'ADNe

Actuellement, les études reposant sur l'ADNe (barcoding, métabarcoding, métagénomique) connaissent un vrai succès dans une vaste gamme de domaines, comme le montre la Figure 0.3.5 : le suivi d'espèces invasives (Schneider et al., 2016), la détermination de régimes alimentaires (Pompanon et al., 2012), la gestion des écosystèmes (Beng and Corlett, 2020), la détection de parasites (Mulerio et al., 2021), la conservation d'espèces (Nordstrom et al., 2022), l'observation de réseaux d'interactions plantes-pollinisateurs (Pornon et al., 2016), la supervision de systèmes agricoles (Kestel et al.,

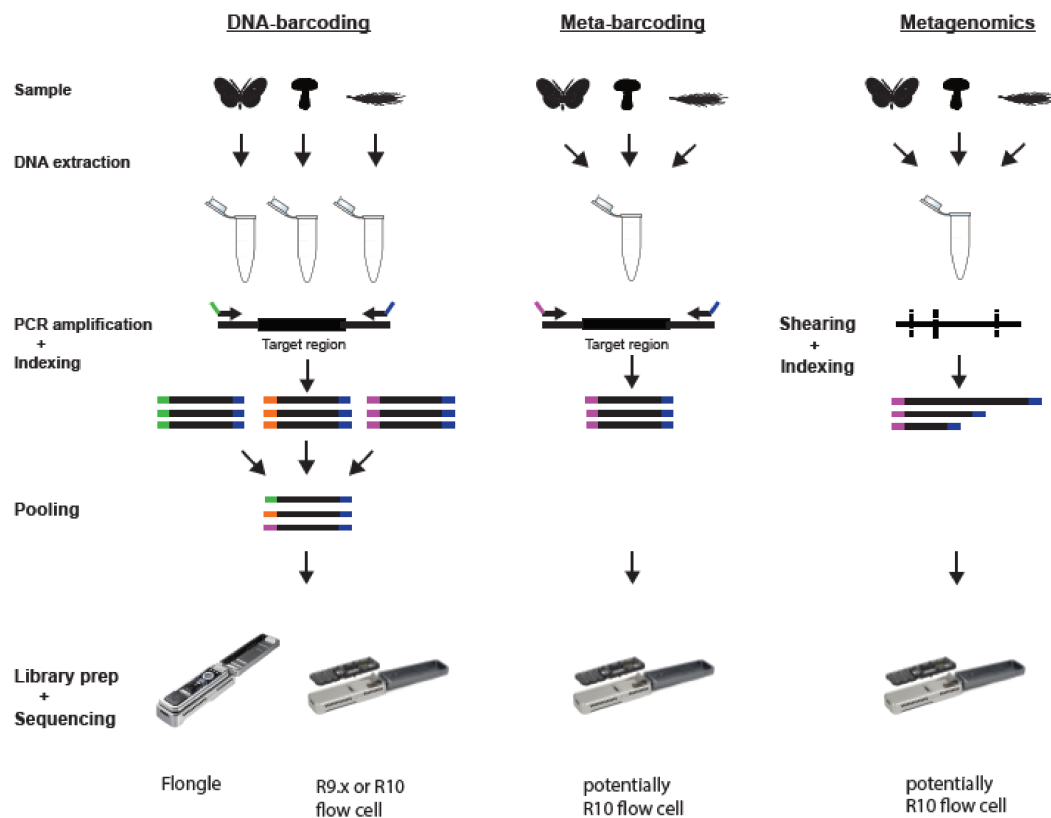


FIGURE 0.3.4 – Trois techniques de biologie moléculaire pour mesurer la biodiversité, ici avec un séquençage Nanopore. Source de la figure : Krehenwinkel et al. (2019).

2022), ou encore la génétique des populations (Adams et al., 2019)...

De son côté, l'ADN ancien est utilisé entre autres pour reconstituer la biodiversité d'écosystèmes disparus, comme celle du Groenland il y a deux millions d'années (Kjær et al., 2022), ou encore pour attester de la présence d'espèces préhistoriques quand les fossiles font défaut (Haile et al., 2009).

0.3.4 Protocole du métabarcoding

Le protocole expérimental du métabarcoding peut se résumer en six étapes (Zinger et al., 2019) :

1. design expérimental, choix du marqueur ;
2. échantillonnage et préservation du matériel collecté ;
3. extraction de l'ADN ;
4. amplification par Réaction de Polymérisation en Chaîne (*Polymerase Chain Reaction*, PCR) des codes-barres moléculaires grâce à un thermocycleur ;
5. séquençage haut-débit des amplicons (molécules créées pendant la PCR) pour former une librairie ;
6. analyse bioinformatique des séquences obtenues appelées lectures.

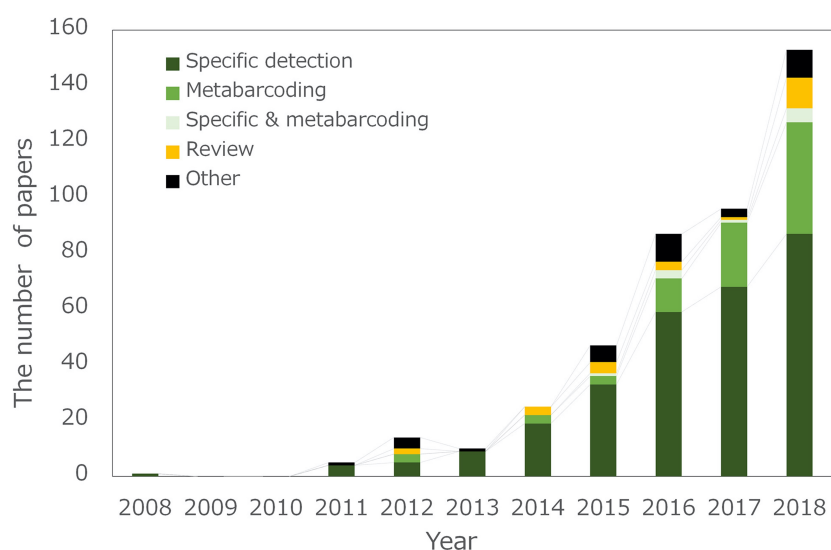


FIGURE 0.3.5 – Nombres d’articles portant sur le barcoding (*specific*) et le métabarcoding publiés entre 2008 et 2018. Source de la figure : Tsuji et al. (2019).

0.3.4.1 Design expérimental pour le métabarcoding

Choix du marqueur Le choix du marqueur est une étape cruciale pour répondre aux questions écologiques de manière pertinente et limiter les biais expérimentaux (Taberlet et al., 2012; Fahner et al., 2016; Krehenwinkel et al., 2017). Le marqueur désigne la région du génome où se trouvent les codes-barres moléculaires. Le guide Taberlet et al. (2018) en propose plusieurs dizaines. Ces codes-barres sont généralement courts pour être observables dans de l’ADN dégradé. Des amorces dites universelles ciblent des groupes taxonomiques larges, par exemple le gène mitochondrial cytochrome c oxidase I (COI) (Hebert et al., 2003) pour les animaux et les gènes chloroplastiques *rbcl* (Hollingsworth, 2011) ou encore ITS2 (Chen et al., 2010) pour les plantes. La spécificité et la sensibilité du marqueur sont des critères importants : on cherche à amplifier et à identifier le maximum d’espèces du groupe ciblé mais à ne pas amplifier les autres groupes. Le consortium *International Barcode Of Life* (IBOL) propose des standards pour générer des données de métabarcoding dans les meilleures conditions.

Un exemple d’amorces est donné sur la Figure 0.3.6. Deux autres exemples sont donnés en Annexe E : *Sper01* pour les Spermatophytes (aussi appelé *g-h*, ciblant la boucle P6 de l’intron *trnL* (UAA) du chloroplaste (Taberlet et al., 2007)), et *Euka03* pour les eukaryotes (ciblant le gène 18S).

Les amorces peuvent être désignées par le logiciel *ecoPrimers* (Riaz et al., 2011). Celui-ci cherche dans une base de génomes un potentiel marqueur qui maximise le nombre d’espèces qui peuvent être amplifiées pour un groupe donné (sensibilité et spécificité) et parmi celles-ci le nombre qui peuvent être différenciées les unes des autres (résolution). Les performances des amorces sont testées grâce à des logiciels de PCR *in silico* comme *ecoPCR* (Ficetola et al., 2010), *Primer-Blast*⁴ ou *FastPCR* (Kalendar,

4. <https://www.ncbi.nlm.nih.gov/tools/primer-blast/index.cgi>

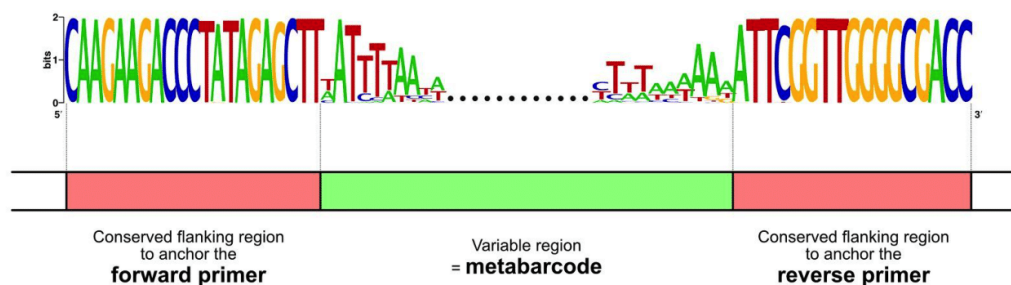


FIGURE 0.3.6 – Exemple d’amorce utilisé en métabarcoding : *Lumb01* ciblant les vers de terre (Bienert et al., 2012). Les régions initiales et finales (*forward* et *reverse primers*) sont conservées, la région intermédiaire est variable. La hauteur pour chaque base correspond au taux de conservation. Source de la figure : Taberlet et al. (2018).

2022). Ceux-ci évaluent le nombre de mismatches possibles pour l’ensemble des taxons considérés. Les mismatches sont des erreurs de complémentarité entre les amorces et le code-barres.

Bases de référence Les bases de référence des génomes sont obtenues par *genome skimming* (Coissac et al., 2016; Alsos et al., 2020). On peut citer par exemple Phyloalps⁵ qui contient la référence de plusieurs milliers de plantes alpines.

Nombre de réplicats et contrôles En général, trois réplicats techniques (amplifications PCR indépendantes) sont produits à partir de chaque échantillon afin de pouvoir identifier un réplicat défectueux tout en minimisant le coût de l’expérience. Ces réplicats techniques doivent être accompagnés de réplicats biologiques (échantillons indépendants) pour fournir des résultats complets (Ruppert et al., 2019).

Il est important de standardiser les protocoles utilisés et d’incorporer des contrôles pour rendre l’étude robuste. Ces contrôles peuvent être positifs (échantillon de composition connue) et négatifs (réplicats sans incorporation d’ADN ajoutés avant l’extraction et avant l’amplification) pour estimer les biais quantitatifs (présentés en 0.6) et pour détecter des contaminations (Zinger et al., 2019; Calderón-Sanou et al., 2020).

0.3.4.2 Extraction de l’ADN

Il est nécessaire d’extraire l’ADN avant l’analyse en détruisant les cellules et les protéines présentes. L’ADN est extrait à partir de kits commerciaux ou selon le protocole CTAB. Ali et al. (2017) présente une revue des techniques communes. À cette étape, l’ADN cible (les codes-barres moléculaires) est mélangé à une grande quantité d’ADN non ciblé. Il faut donc amplifier l’ADN cible pour améliorer la détection.

5. <https://data.phyloalps.org/browse/>

0.3.4.3 Amplification par PCR

L'amplification par PCR est au cœur de mes travaux de thèse. Je détaille ce processus en 0.4 et dans le Chapitre 2. À l'issue de la PCR, les molécules d'ADN cible ont été multipliées de plusieurs ordres de grandeur : elles fournissent un signal assez clair pour être étudié après le séquençage.

0.3.4.4 Séquençage

Le séquençage permet enfin de lire une fraction des molécules amplifiées qui peut être analysée. Pervez et al. (2022) présente une revue des techniques de séquençage existantes. Le séquençage utilisé en métabarcoding est de deuxième génération, en particulier Illumina qui séquence des fragments courts (jusqu'à 300 paires de bases). Les méthodes de seconde génération ont remplacé le séquençage par la méthode historique de Sanger et ont permis le développement massif du métabarcoding. Actuellement, de nouvelles méthodes (Nanopore) se mettent en place (section 0.7.3). La Figure 0.3.7 donne un aperçu de l'évolution des capacités de séquençage depuis son invention.

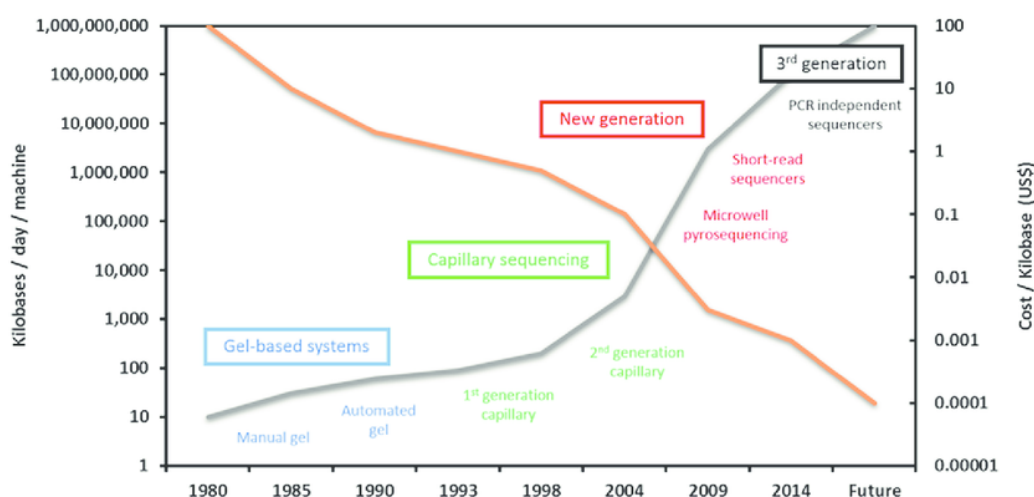


FIGURE 0.3.7 – Évolution des capacités de séquençage en termes de coût et de débit. Les séquenceurs de *Next generation* sont aussi appelés de seconde génération. Source de la figure : Goss-Souza et al. (2016).

0.3.4.5 Analyse bioinformatique

Les données de séquençage brutes sont ensuite analysées. Les étapes préliminaires usuelles sont :

- Aligner les lectures des données de séquençage *paired end* au cours duquel les brins complémentaires ont été séquencés ;
- démultiplexer les données dans le cas où plusieurs échantillons sont réunis dans la même librairie ;

- dérépliquer les lectures en regroupant les variants, c'est-à-dire les séquences identiques ;
- filtrer les séquences selon des critères simples comme la longueur des fragments ou en retirant les singletons (variants observés une seule fois) ;
- regrouper les variants proches en unités taxonomiques, proxys des espèces réelles ;
- Attribuer les unités taxonomiques à des espèces référencées dans une base de données.

Des logiciels bioinformatiques comme les OBITools (Boyer et al., 2016)⁶) permettent de réaliser ces différentes étapes. Le sujet du regroupement des variants est traité en détail dans le Chapitre 3.

0.3.5 Interprétation des données

Le métabarcoding n'est pas exempt de limites pour produire des données exploitables du point de vue écologique. Deux difficultés majeures existent.

La première est l'attribution des séquences observées à des espèces connues. Même avec des bases de données complètes, cette tâche n'est pas évidente car le nombre de variants observés dépasse de plusieurs ordres de grandeur le nombre estimé d'espèces. Ces nombreux variants ont pour origine des erreurs de PCR et de séquençage et de la variabilité biologique. Il est nécessaire de traiter cette variabilité car elle affecte les mesures de biodiversité. Une première étape consiste donc à extraire des unités taxonomiques de ces données par clustering pour limiter le nombre de fausses détections (Chapitre 3).

La seconde difficulté a trait aux abondances des espèces présentes. La pratique consistant à comparer la fréquence de lectures de chaque espèce à leur abondance relative dans l'écosystème (en biomasse ou en nombre d'individus) est contestable du fait des nombreux biais affectant ces valeurs. Je reviens sur ce point en 0.6 et dans le Chapitre 2.

0.4 Réaction en Chaîne par Polymérase (PCR)

Après avoir présenté globalement le métabarcoding et ces enjeux, je m'arrête ici sur l'amplification par PCR sur laquelle se concentre le Chapitre 2. La PCR est identifiée comme une limite majeure au développement du métabarcoding quantitatif? Un objectif de ma thèse est de mesurer et de corriger les biais dus à la PCR. J'en présente d'abord le fonctionnement et certaines des applications.

La PCR est une technique massivement employée en biologie moléculaire qui multiplie le nombre de molécules d'ADN présentes dans un échantillon par réplifications successives à partir d'une séquence d'ADN matrice (*template*). Cela permet de détecter plus facilement des séquences spécifiques dans un mélange d'ADN. La PCR a été inventée aux environs de 1983 par l'Américain Kary Mullis (publié dans Mullis and Faloona (1987)) qui reçoit le Prix Nobel de Chimie pour cette invention en 1993.

6. <https://pythonhosted.org/OBITools/welcome.html>

0.4.1 Procédé

La PCR est basée sur une succession de 30 à 45 cycles de température. À chaque cycle, les molécules d'ADN matrice présentes peuvent être répliquées et les nouvelles molécules participent à l'amplification aux cycles suivants. Les principaux réactifs de la PCR sont :

- des séquences d'ADN matrice ;
- une enzyme, la Taq polymérase, qui permet la réplication de l'ADN. Celle-ci est thermorésistante pour supporter les hautes températures du cycle de PCR. Certaines polymérases ont une activité exonucléase dans le sens 3' → 5' dite de *proofreading* pour corriger les erreurs de réplication (Khare and Eckert, 2002; Gohl et al., 2021).
- Deux types d'amorces (*primers*) complémentaires des brins 3' "sens" et "anti-sens". Ces amorces permettent d'amplifier exclusivement l'ADN matrice. L'amplification peut fonctionner dans le cas où la complémentarité entre l'amorce et le brin 3' n'est pas parfaite, mais ces mismatches réduisent l'efficacité de la PCR.
- Des nucléotides ou dNTP (désoxyribonucléotide triphosphate) qui servent de briques élémentaires pour la réplication de l'ADN ;
- une solution tampon pour maintenir le pH à un niveau optimal pour l'action de la polymérase ;
- des cofacteurs facilitant l'activité de la polymérase comme des ions Mg^{2+} .

La Figure 0.4.8 illustre le déroulement de la PCR. Au départ, l'ADN présent est double-brin. Une étape de dénaturation initiale (10 minutes à 95°C⁷) sépare ces doubles-brins.

Au cours de chaque cycle, trois étapes se répètent. Chaque cycle commence par une étape de dénaturation (30 secondes à 95°C⁷) qui permet de dés-hybrider les brins d'ADN et d'homogénéiser le milieu réactionnel.

Ensuite, une phase d'hybridation (30 secondes à 50-60°C⁷) permet aux amorces de se fixer sur les brins d'ADN matrice. Cette température de fusion, notée T_m , dépend des amorces (*melting temperature*). Il s'agit de la température à laquelle 50 % de l'ADN double-brin est dissocié.

Enfin, une phase d'élongation ou extension (1 minute à 72°C⁷) a lieu pendant laquelle la polymérase synthétise un nouveau brin à partir de l'ADN matrice. Les molécules d'ADN créées sont appelées amplicons.

En théorie, le nombre de molécules d'ADN matrice devrait doubler à chaque cycle. En pratique, l'efficacité de PCR est un peu inférieure à 1 car certaines molécules ne sont pas répliquées au cours de chaque cycle. Cette efficacité varie d'une PCR à l'autre : elle dépend principalement de la séquence matrice et de l'échantillon biologique (comportant des inhibiteurs) mais a priori pas de la quantité initiale d'ADN (Karlen et al., 2007). L'amplification est exponentielle à un taux inférieur à 2 pendant plusieurs cycles puis ralentit jusqu'à atteindre un plateau. Une cinétique typique de PCR est représentée sur la Figure 0.4.9.

7. Les durées et températures sont indicatives.

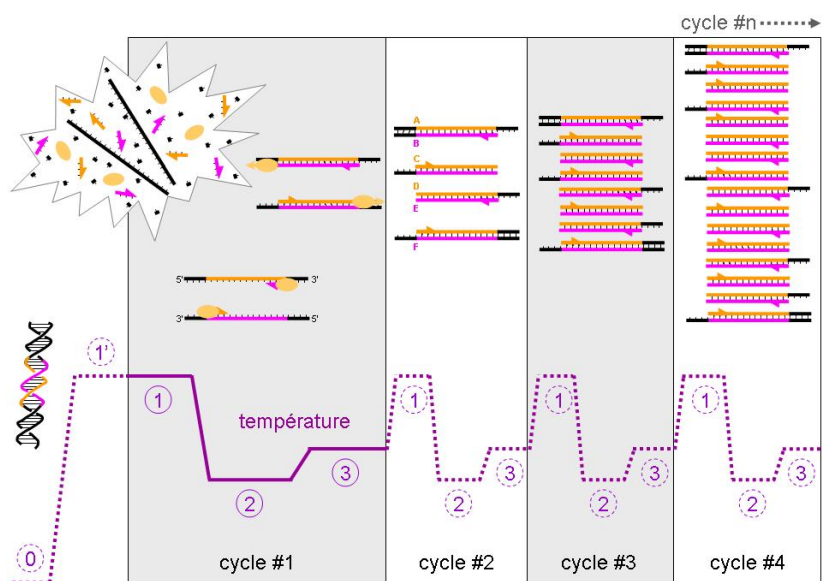


FIGURE 0.4.8 – Déroulement des premiers cycles d’une amplification par PCR. L’état initial (0) est suivi de la dénaturation initiale (1’) puis d’une répétition de cycles : dénaturation (1), hybridation (2), élongation (3). Source de la figure : https://fr.wikipedia.org/wiki/Réaction_en_chaîne_par_polymérase.

Cinétique de PCR

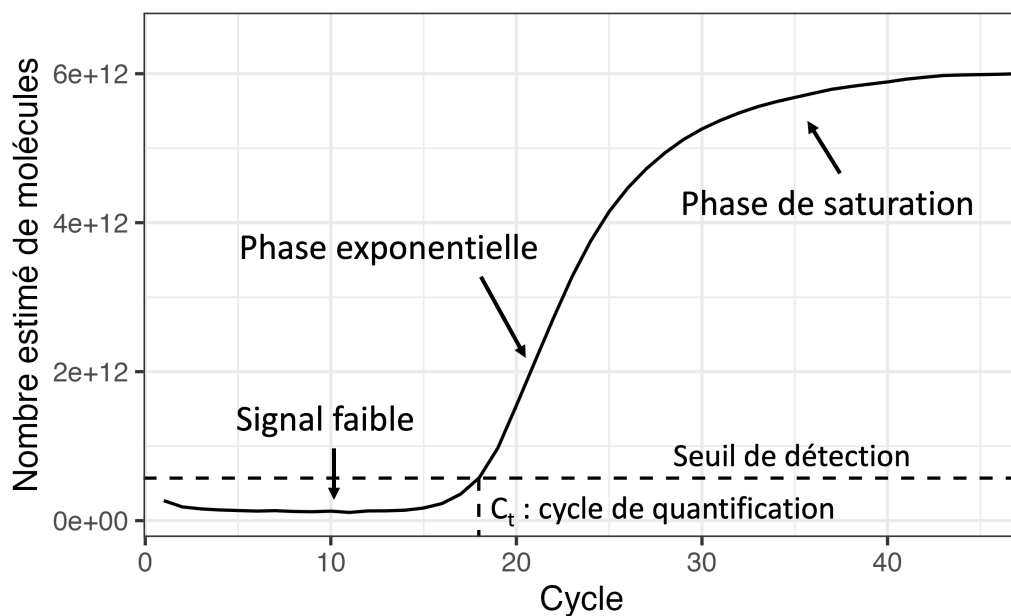


FIGURE 0.4.9 – Cinétique typique de PCR, avec une phase exponentielle puis une phase de plateau. Cette cinétique est observée par qPCR à partir de mesures de fluorescence (section 0.4.4).

0.4.2 Mécanismes de saturation au cours de la PCR

Différentes hypothèses ont été émises pour expliquer la saturation responsable du plateau observé (Van Der Graaf and Schoemaker, 1999; Gottschalk and Dunn, 2005). La première hypothèse est l'épuisement des réactifs (dNTP, amorces) qui deviennent limitants (Hayward, 1998; Carr and Moore, 2012). La seconde hypothèse, qui n'exclut pas la première, est que la saturation est due à l'auto-hybridation : les brins simples d'ADN se lient entre eux et ne participent donc pas à l'amplification (Mathieu-Daude, 1996; Gevertz et al., 2005; Mehra and Hu, 2005). Ce phénomène est supposé d'autant plus fort que les brins d'ADN sont nombreux, ce qui induit un biais contre les espèces les plus abondantes dans l'amplification (Mathieu-Daude, 1996). Ces différentes hypothèses sont étudiées dans les modèles de PCR présentés en 0.5.

0.4.3 Applications de la PCR

Le protocole de PCR connaît de nombreuses variantes. Une application majeure de la PCR est le diagnostic médical pour détecter la présence d'un virus ou d'une bactérie chez un patient. Le cas le plus connu ces dernières années est bien sûr le virus SARS-CoV-2 responsable du Covid-19. Sa présence est détectée chez le patient par Reverse-Transcriptase PCR (RT-PCR) où l'ARN est transcrit en ADN (transcription inverse) avant l'amplification par PCR (Tahamtan and Ardebili, 2020).

Des techniques de PCR quantitatives sont également très utilisées pour estimer une charge virale, un niveau d'expression de gène ou encore la quantité d'ADN dans un échantillon environnemental. L'ADN matrice présent dans un échantillon est quantifié soit de manière absolue (en copies de matrice dans l'échantillon) soit de manière relative à une séquence de référence. Ces techniques, utiles pour l'étude de l'ADNe, sont détaillées dans la section suivante.

0.4.4 Techniques de PCR quantitatives

La principale technique est la PCR quantitative (qPCR) à proprement parler ou PCR en temps réel. Elle est apparue au début des années 1990 (Higuchi et al., 1993). Elle consiste en une PCR classique où l'on fait intervenir un agent fluorescent dont l'activité dépend de la quantité d'ADN dans le milieu réactionnel. Son intensité est mesurée à chaque cycle en Unité Relative de Fluorescence (*Relative Fluorescence Unit*, RFU). Deux classes principales de qPCR existent.

0.4.4.1 Détection non spécifique

La première catégorie de qPCR permet de détecter l'ensemble des molécules d'ADN matrice du milieu réactionnel. L'agent fluorescent (*dye*), en général le SybrGreen, n'est activé que s'il est lié à une molécule d'ADN double-brin. Ce principe est illustré sur la Figure 0.4.10. Pour contrôler la spécificité de l'amplification, on observe la courbe de fusion qui montre la fluorescence (RFU) en fonction de la température (T) lors de la phase de dénaturation. L'analyse de la variation de RFU en fonction de T (plus

précisément des pics de $-\frac{dRFU}{dT}$) permet de détecter la présence de différentes séquences d'ADN car la température de fusion T_m varie d'une séquence à une autre.

0.4.4.2 Détection d'une séquence spécifique

La seconde catégorie de qPCR repose sur des sondes qui ciblent une séquence spécifique, par exemple une seule espèce dans un échantillon environnemental. La technique la plus commune est la qPCR Taqman, qui incorpore une sonde Taqman complémentaire de la séquence d'intérêt. Cette sonde est composée d'un fluorophore relié à un inhibiteur (*quencher*), fonctionnant sur le principe du transfert d'énergie entre molécules fluorescentes (*Fluorescence Resonance Energy Transfer*, FRET). Tant que ces deux composants sont en contact, aucune fluorescence n'est émise. La sonde se fixe sur l'ADN simple-brin pendant l'hybridation, puis le fluorophore est séparé de l'inhibiteur durant la phase d'élongation, ce qui déclenche la fluorescence. Ce principe est aussi illustré sur la Figure 0.4.10. En incorporant plusieurs sondes avec des spectres d'émission différents, il est possible de suivre l'évolution de plusieurs espèces simultanément au cours d'une même PCR. Des alternatives à la qPCR Taqman existent, par exemple en utilisant une hybridation de deux sondes, des balises moléculaires (*Molecular Beacon*), ou encore des amorces scorpions. Wang and Yang (2013) en présente une revue.

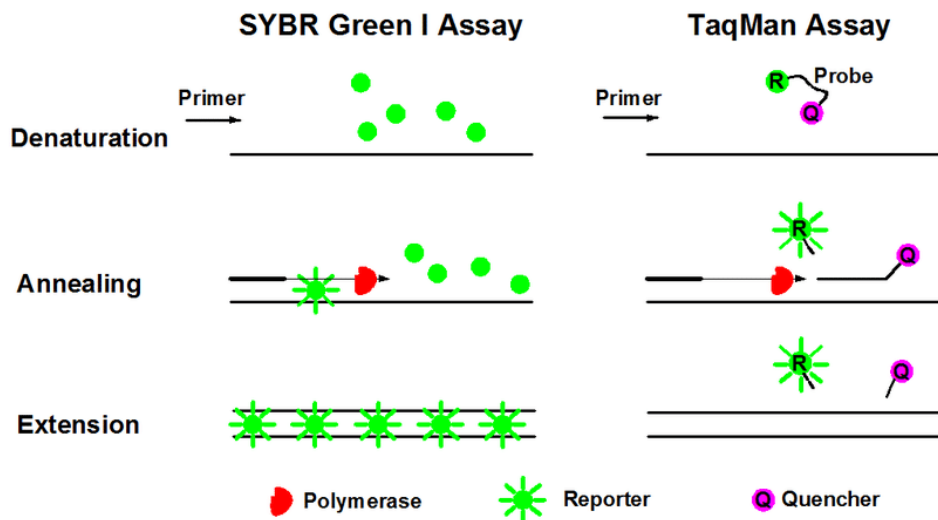


FIGURE 0.4.10 – Principe de la qPCR SybrGreen (à gauche) et Taqman (à droite). Dans le premier cas, le fluorophore est activé quand il se lie à une molécule d'ADN double-brin. Dans le second, la synthèse d'un nouveau brin d'ADN sépare le fluorophore de son inhibiteur. Source de la figure : Cao et al. (2020).

0.4.4.3 Principe de quantification par qPCR

Pour les deux approches de qPCR, la quantification repose sur la détection du cycle C_t (*threshold cycle*) où un niveau seuil de fluorescence est atteint (Figure 0.4.9). Plus ce cycle est tardif, moins l'échantillon comportait initialement d'ADN. La fluorescence est considérée comme proportionnelle au nombre de molécules (Gill et al., 2022). Ce cycle est atteint durant la phase exponentielle de la PCR, ce qui rend la mesure plus précise que les études en "point final" (Gál et al., 2006). Les conditions expérimentales peuvent être ajustées pour optimiser l'efficacité de la PCR (Zhang et al., 2019). La quantification est relative à une espèce de référence ou absolue en se basant sur des contrôles positifs (Peirson, 2003).

La méthode du $\Delta\Delta C_t$ (Livak and Schmittgen, 2001) est une approche classique pour convertir la valeur de C_t en une différence d'abondance (*fold-change*) entre deux séquences. Au sein d'un échantillon, on calcule la différence de C_t entre la séquence d'intérêt et une séquence de référence (équation 13). Cette différence est ensuite comparée à la différence observée dans un contrôle dans lequel les deux séquences sont présentes en quantités connues (équation 14). La différence d'expression est alors estimée par $2^{-\Delta\Delta C_t}$. Cela suppose que les deux séquences ont une efficacité de PCR toutes deux environ égales à 1 (rendement parfait), ce qui est contestable. Dans ce cas, un décalage de C_t d'un cycle entre deux échantillons signifie que l'un contient deux fois plus d'ADN que l'autre. Cette différence de $1C_t$ (et donc une différence d'abondance de 2) est considéré comme la précision limite pour différencier deux échantillons (d'après la documentation technique de Thermo-Fisher⁸).

$$\Delta C_t(\text{échantillon}) = C_t(\text{espèce à doser}) - C_t(\text{référence}) \quad (13)$$

$$\Delta\Delta C_t = \Delta C_t(\text{échantillon}) - \Delta C_t(\text{contrôle positif}) \quad (14)$$

Une autre méthode de quantification repose sur une gamme de dilution sérielle d'un échantillon. On observe une relation linéaire entre C_t et le logarithme de la concentration en ADN de l'échantillon (Ramakers et al., 2003). Un exemple est montré sur la Figure 0.4.11. La pente de la régression linéaire

$$C_t = a \log(\text{Concentration initiale}) + b \quad (15)$$

permet d'estimer l'efficacité de PCR Λ par la relation (Svec et al., 2015) :

$$\Lambda = 10^{-1/a} - 1 \quad (16)$$

Cette équation découle directement du modèle de PCR exponentiel présenté en 0.5. La mesure d'efficacité est affectée par l'inhibition au cours de la PCR (voir plus bas). Je discute de la précision de cette méthode dans le Chapitre 2.

8. <https://www.thermofisher.com/fr/fr/home/life-science/pcr/real-time-pcr/real-time-pcr-learning-center/gene-expression-analysis-real-time-pcr-information/precision-qpcr.html>

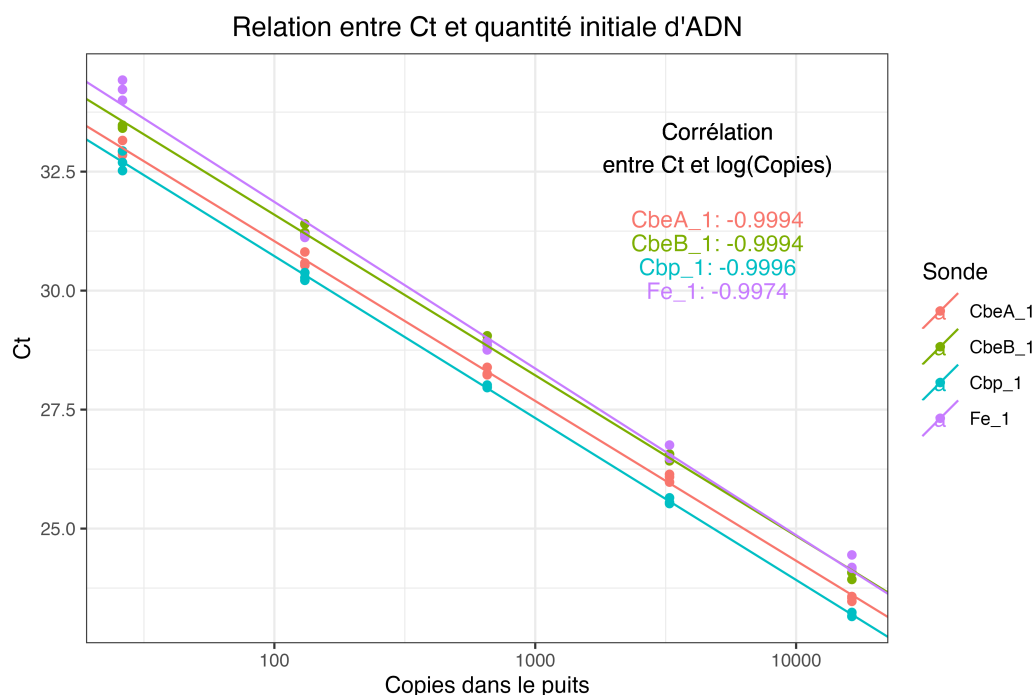


FIGURE 0.4.11 – Cycles de quantification C_t en fonction du nombre de copies initialement présentes dans l'échantillon constitué de trois espèces, pour quatre sondes Taqman différentes. Chaque concentration fait l'objet de trois réplicats. L'axe des abscisses est logarithmique.

La mesure d'efficacité est généralement utilisée pour comparer l'espèce à doser dans l'échantillon à une espèce de référence en prenant en compte leurs différences d'amplification (Pfaffl, 2001). Le ratio entre les deux est donné par :

$$\text{Ratio} \frac{\text{espèce à doser}}{\text{référence}} = \frac{\Lambda_{\text{espèce à doser}}^{\Delta C_t(\text{échantillon} - \text{contrôle})}}{\Lambda_{\text{référence}}^{\Delta C_t(\text{échantillon} - \text{contrôle})}} \quad (17)$$

Cette mesure est plus précise que la méthode $\Delta\Delta C_t$ (Pfaffl, 2001) et les deux méthodes peuvent mener à des conclusions biologiques différentes (Skern et al., 2005).

0.4.4.4 Droplet digital PCR (ddPCR)

La droplet digital PCR (ddPCR) est une autre technique de PCR quantitative (parfois appelée simplement digital PCR (dPCR)). C'est une technologie plus récente : la première machine commercialisée est le *QX100 ddPCR System* en 2011 par Bio-Rad (Morley, 2014). Elle ne repose pas sur une détection en temps réel de la quantité d'ADN mais en fournit une quantification absolue. Le principe, illustré sur la Figure 0.4.12, consiste à former environ 20 000 gouttelettes contenant au plus quelques molécules d'ADN à partir de l'échantillon. Une PCR est ensuite réalisée pour chaque gouttelette. Si la gouttelette contenait de l'ADN, un signal fluorescent est détecté grâce à

l'agent EvaGreen. Les gouttelettes positives et négatives sont comptées et un modèle de Poisson permet de déterminer le nombre probable de molécules dans l'échantillon. Cette technique est plus précise que la qPCR (Doi et al., 2015; Uchii et al., 2019; Brys et al., 2021). Plusieurs raisons sont avancées : le signal est binaire (présence ou absence) et non continu (comparaison à un seuil), et la ddPCR effectue un grand nombre de répliquats techniques (chaque gouttelette). La différence d'expression détectable est d'environ 1.2 *fold-change* (d'après la documentation technique de Bio-Rad⁹).

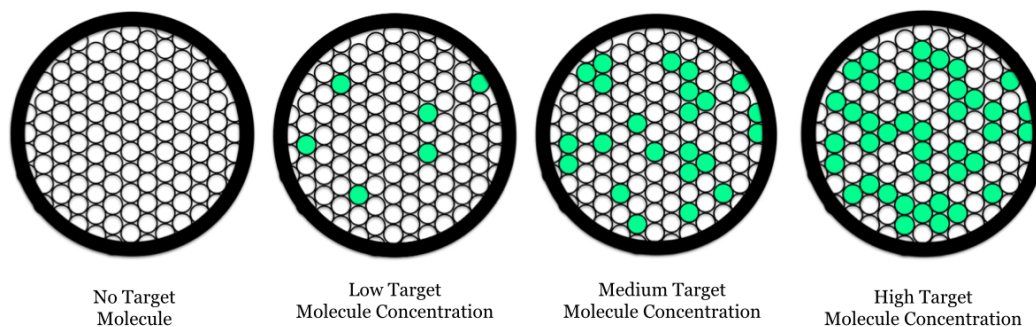


FIGURE 0.4.12 – Principe de la ddPCR. Les puits fluorescents sont comptés pour quantifier l'ADN présent. Source de la figure : https://en.wikipedia.org/wiki/Digital_polymerase_chain_reaction.

0.4.5 Inhibition de la PCR

L'inhibition au cours de la PCR est due à des substances qui diminuent l'amplification en affectant l'activité de la polymérase. Il est nécessaire de la prendre en compte car elle affecte les résultats de quantifications.

Les inhibiteurs peuvent se lier à la polymérase ou aux brins d'ADN, interférer avec le complexe polymérase-ADN pendant l'élongation ou encore interférer avec les sondes fluorescentes en qPCR (Bar et al., 2012; Sidstedt et al., 2020). Il existe trois types principaux d'inhibition : l'inhibition compétitive (l'inhibiteur se fixe sur l'enzyme et exclut le substrat, c'est-à-dire la molécule interagissant normalement avec l'enzyme), incompétitive (l'inhibiteur se fixe sur le complexe enzyme-substrat) et non compétitive (les sites de fixation du substrat et de l'inhibiteur sont distincts).

Les inhibiteurs peuvent provenir de l'échantillonnage ou de l'extraction de l'ADN (Schrader et al., 2012). L'inhibition fausse l'estimation du rendement par gamme de dilution en qPCR car les répliquats les plus concentrés sont soumis à des inhibiteurs qui sont eux aussi plus concentrés. La dilution a un double effet sur la concentration en ADN et en inhibiteur, ce qui "aplatit" la pente de régression de $C_t = a \log \text{Concentration} + b$ et augmente artificiellement l'efficacité mesurée ? Celle-ci

9. <https://www.bio-rad.com/fr-fr/life-science/digital-pcr/digital-pcr-assays/ddpcr-copy-number-determination-assays>

dépasse parfois le maximum théorique de 1 (Svec et al., 2015). C'est une des raisons pour lesquelles la ddPCR est plus précise que la qPCR.

L'inhibition est une limite connue à la quantification des échantillons environnementaux (Uchii et al., 2019; Takasaki et al., 2021). En particulier, des faux négatifs sont reportés : des espèces présentes ne sont pas détectées car la PCR est totalement inhibée (Jane et al., 2015; Fujii et al., 2019). Les solutions proposées consistent généralement à diluer plus fortement les échantillons et à utiliser des kits de purification de l'ADN. En qPCR, des outils statistiques de *Kinetic Outlier Detection* visent à détecter les cinétiques de PCR présentant des signes d'inhibition (Bar et al., 2012).

0.5 Modélisation de la PCR

L'émergence de la qPCR a conduit à de nombreux travaux de modélisation depuis la fin des années 1990. Le but est de décrire les cinétiques observées afin d'améliorer la quantification de l'ADN présent dans un échantillon avant l'amplification. Il peut s'agir d'améliorer la précision de la mesure ou d'obtenir une quantification absolue tout en réduisant la quantité de réplicats nécessaires. L'idée historique, non réalisée à ce jour, était de quantifier automatiquement un échantillon à partir d'une seule mesure et sans biais lié à la nature de l'échantillon ni aux paramètres techniques du protocole utilisé.

Voici une revue des modèles existants. La Figure 0.5.13 montre les cinétiques établies par les plus réalistes. Le package R `qpcR` (Ritz and Spiess, 2008) permet d'en manipuler plusieurs¹⁰. Il existe deux grandes catégories de modèles. Les premiers sont analytiques : ils ont pour but de représenter fidèlement les cinétiques observées sans justification biologique mais en utilisant une forme mathématique appropriée pour être facilement ajustables aux données. Les seconds sont mécanistiques et visent à mettre en équation les mécanismes biochimiques de la PCR. Par ailleurs, certains modèles sont déterministes tandis que d'autres sont aléatoires pour étudier l'incertitude observée dans les données. Ce sujet sera abordé dans le Chapitre 2.

0.5.1 Notations

Dans ce qui suit, M_n désigne le nombre d'amplicons et F_n la fluorescence (en RFU) au cycle n . L'efficacité de PCR est notée Λ . Lorsque celle-ci est variable, elle est notée λ_n pendant le cycle n . Afin de simplifier l'écriture, les modèles sont décrits avec un seul type d'ADN matrice présent, mais ceux-ci sont généralisables à un nombre arbitraire d'espèces sans difficulté.

Pour une quantification absolue, le lien entre le nombre de molécules et la fluorescence observée est établi de manière proportionnelle grâce à une courbe de référence.

10. voir page 73 de la documentation pour une liste des modèles implémentés, <https://cran.r-project.org/web/packages/qpcR/qpcR.pdf>

Ce lien n'a pas besoin d'être explicité pour une quantification relative.

0.5.2 Intérêt et limites du modèle exponentiel

Le modèle élémentaire est le modèle exponentiel (Suzuki and Giovannoni, 1996). Il a l'avantage d'être très simple à manipuler et pertinent en début d'amplification. En effet, l'efficacité de PCR est à peu près constante pendant plusieurs cycles car la réaction n'est pas encore saturée. À chaque cycle, ce modèle postule qu'une molécule présente une probabilité $\Lambda \in [0, 1]$ constante d'être répliquée. Si $\Lambda = 0$, l'amplification n'a pas lieu ; si $\Lambda = 1$, l'amplification est parfaite. En moyenne, le processus est donc décrit par :

$$\text{Pour tout cycle } n, \quad M_n = M_0 \cdot (1 + \Lambda)^n \quad (18)$$

Ce modèle exponentiel est compatible avec la mesure du cycle de quantification C_t . En revanche, il diverge rapidement quand le nombre de cycles augmente et accroît artificiellement les écarts d'abondance dus à des efficacités de PCR différentes. Par ailleurs, les mesures basées sur les valeurs de C_t n'exploitent qu'un seul point par courbe d'amplification. Un modèle plus précis permettrait théoriquement d'exploiter l'ensemble de la dynamique de PCR pour extraire de l'information sur l'échantillon.

0.5.3 Modèles analytiques

Je présente d'abord quelques modèles analytiques parmi les plus répandus. Ceux-ci décrivent explicitement la fluorescence en fonction du cycle par une sigmoïde avec un nombre limité de paramètres. Le premier est un modèle logistique à quatre paramètres (Liu and Saint, 2002) appelé b_4 dans le package qpcR (équation 19). Les paramètres sont liés au niveau de fluorescence de base (F_{base}), au niveau de fluorescence final sans le signal de base (F_{max}), à l'efficacité de PCR (k) et au cycle de demi-amplification ($n_{1/2}$). La quantité initiale d'ADN est donnée par l'équation 20 (en fluorescence, à convertir en nombre de molécules).

$$F_n = F_{\text{base}} + \frac{F_{\text{max}}}{1 + \exp((n - n_{1/2})/k)} \quad (19)$$

$$\text{donc } F_0 = F_{\text{base}} + \frac{F_{\text{max}}}{1 + \exp(n_{1/2}/k)} \quad (20)$$

Sa forme récurrente est définie par l'équation 21 ce qui permet d'expliciter le rendement cycle par cycle :

$$F_n - F_{\text{base}} = (F_{n-1} - F_{\text{base}}) \times \frac{1 + \exp(-(n-1 - n_{1/2})/k)}{1 + \exp(-(n - n_{1/2})/k)} \quad (21)$$

Ce modèle montre qu'une fluorescence finale plus élevée ne traduit pas une quantité initiale plus grande. Il a été étudié à des fins de quantification par un certain nombre

d'auteurs (Tichopad, 2003; Rutledge, 2004; Zhao and Fernald, 2005; Goll et al., 2006; Chervoneva et al., 2007; Rutledge and Stewart, 2008). Ce modèle a fait l'objet de plusieurs variantes par la suite. Un modèle log-logistique a été établi (l_4 dans qpcR, équation 22) avec une meilleure qualité d'ajustement (Zhao and Fernald, 2005).

$$F_n = F_{\text{base}} + \frac{F_{\text{max}}}{1 + \exp((\log(n) - \log(n_{1/2}))/k)} \quad (22)$$

Les modèles b_4 et l_4 ont été améliorés avec l'ajout d'un cinquième paramètre (Spiess et al., 2008) afin de prendre en compte l'asymétrie de l'amplification (Van Der Graaf and Schoemaker, 1999; Gottschalk and Dunn, 2005) : l'accélération de l'amplification est plus rapide que sa décélération, alors que les modèles précédents sont symétriques par rapport au point d'inflexion, c'est-à-dire le point de demi-amplification. Les modèles b_5 et l_5 figurent en équations 23 et 24. L'asymétrie est caractérisée par le paramètre f . Le modèle l_5 parvient au meilleur ajustement selon les tests statistiques réalisés par Spiess et al. (2008).

$$F_n = F_{\text{base}} + \frac{F_{\text{max}}}{(1 + \exp((n - n_{1/2})/k))^f} \quad (23)$$

$$F_n = F_{\text{base}} + \frac{F_{\text{max}}}{(1 + \exp((\log(n) - \log(n_{1/2}))/k))^f} \quad (24)$$

Je cite deux autres modèles analytiques : le modèle de Gompertz (Van Der Graaf and Schoemaker, 1999) et le modèle de Chapman (Zhao and Fernald, 2005).

Des essais ont été menés avec ces modèles pour obtenir automatiquement une quantification, sans résultat probant à ma connaissance. L'intérêt est de ne pas utiliser de seuil arbitraire de fluorescence comme c'est le cas pour déterminer C_t et de s'affranchir des différences de fluorescence des sondes ou des fluorophores. Une possibilité est de considérer le cycle de plus grande accélération caractérisé par le maximum de dérivée seconde de la cinétique (*Second Derivative Maximum*) (Zhao and Fernald, 2005).

0.5.4 Modèles mécanistiques

Les modèles mécanistiques mettent en équation les mécanismes biochimiques de la PCR de manière plus ou moins simplifiée. Ils sont définis par récurrence : l'amplification au cycle $n + 1$ dépend de l'état du système au cycle n . La forme commune est donnée par l'équation 25 (éventuellement en décrivant la fluorescence F_n au lieu du nombre de molécules M_n) : le modèle décrit une probabilité λ_n de chaque molécule d'être répliquée et aucune molécule ne disparaît. Dans certains cas, les auteurs choisissent les réactions chimiques qu'ils jugent pertinentes pour décrire chaque cycle de PCR et calculent explicitement les cinétiques chimiques associées pour établir leur modèle.

$$M_n = M_{n-1}(1 + \lambda_n) \quad (25)$$

$$\text{soit } \lambda_n = \frac{M_n - M_{n-1}}{M_{n-1}}$$

Les premiers modèles décrivaient surtout la phase d'élongation en modélisant la probabilité d'ajout de chaque nucléotide (Stolovitzky and Cecchi, 1996; Velikanov and Kapral, 1999). Cette représentation assez lourde est évitée dans la plupart des travaux ultérieurs. D'autres modèles prennent en compte l'ensemble du cycle de manière plus ou moins exhaustive, comme Gevertz et al. (2005) ou encore Mehra and Hu (2005) qui incorpore une vingtaine de réactions chimiques différentes. L'inconvénient est que ces modèles se prêtent mal à l'ajustement aux données observées.

Un des premiers modèles mécanistiques simplifiés a été établi par Schnell and Mendoza (1997) en décrivant la réaction enzymatique de la PCR à partir de l'équation de Michaelis-Menten (équation 26). Les entités en présence sont les molécules d'ADN simple brin (M), l'ADN double-brin (D) et la polymérase (P). Les amorces et dNTP sont considérés comme en excès.



Le modèle de PCR s'écrit alors selon l'équation 27. Ce modèle est valable en début de PCR. Il prend en compte la saturation comme un accroissement linéaire (et non un plateau) remplaçant la phase exponentielle. Ce modèle est adapté par Lalam (2006) pour faciliter l'estimation du paramètre K_M .

$$M_n = M_{n-1} \left(1 + \frac{K_M}{K_M + M_{n-1}} \right) \quad (27)$$

avec $K_M = \frac{k_2 + k_{-1}}{k_1}$ la constante de Michaelis-Menten

De manière un peu différente, le modèle MAK (*Mass Action Kinetics*) (Boggy and Woolf, 2010) ne représente que la phase exponentielle de la cinétique de PCR. Son objectif est là encore de quantifier l'ADN à partir d'un seul réplicat. Les réactifs sont tous considérés en excès et toutes les réactions sont décrites comme totales. Deux réactions sont alors en compétition au cours de chaque cycle : la synthèse d'une nouvelle molécule double-brin à partir d'un brin simple (équation 28) et la réhybridation (équation 29) qui consiste pour deux brins à se lier, à la manière des dimères d'amorces.



Je profite de ce modèle simple pour montrer le genre de système qu'il faut résoudre pour établir la formule de récurrence. La cinétique chimique correspondant aux réactions 28 et 29 au cycle n s'écrit selon le système 30 qu'il faut résoudre en temps. Le résultat quand le temps de réaction tend vers $+\infty$ donne le modèle de l'équation 31, en comptant les molécules d'ADN simple-brin M_n :

$$\begin{cases} \frac{dM}{dt} = -k_a M - k_b M^2 & , \quad M(0) = M_{n-1} \\ \frac{dD}{dt} = k_a M + \frac{1}{2} k_b M^2 & , \quad D(0) = 0 \end{cases} \quad (30)$$

$$M_n = M_{n-1} + k \log \left(1 + \frac{M_{n-1}}{k} \right) \quad \text{avec } k = \frac{k_a}{2k_b} \quad (31)$$

Les modèles précédents ne permettent pas de reproduire la saturation sous forme d'un plateau. Hayward (1998) y parvient en considérant la saturation sous forme de l'épuisement des réactifs. Une quantité maximale de molécules K peut être créée dans le milieu d'amplification. Le rendement décroît à mesure que le nombre d'amplicons s'approche de ce seuil, selon un paramètre c (équation 32). Ce modèle correspond au modèle logistique du Chapitre 2 avec $c = \Lambda$.

$$M_n = M_{n-1} \left(1 + \Lambda - c \frac{M_{n-1}}{K} \right) \quad (32)$$

Enfin, le modèle de Carr and Moore (2012) (*cm3* dans qpcR) prend en compte les deux phénomènes de saturation pour expliquer l'asymétrie de la cinétique (équation 33). D'une part, les réactifs sont limitants comme dans Hayward (1998). D'autre part, les amplicons eux-mêmes génèrent de l'inhibition (Boggy and Woolf, 2010). Dans le cas d'une réaction totale, le milieu réactionnel permet de créer M_{max} molécules mais les auteurs estiment qu'une petite fraction des amorces et des dNTP sont consommées au moment de la saturation. Celle-ci est induite par un mécanisme d'inhibition reposant sur la loi d'action de masse : l'activité enzymatique décroît lorsque la quantité de produit augmente. Dans le modèle, les amplicons agissent comme des inhibiteurs d'affinité avec l'enzyme caractérisée par le paramètre K_d .

$$M_n = M_{n-1} \left(2 - \frac{M_{n-1}}{M_{max}} - \frac{M_{n-1}}{K_d - M_{n-1}} \right) \quad (33)$$

Pour inclure la variabilité d'efficacité, il faudrait remplacer ce modèle par :

$$M_n = M_{n-1} \left(1 + \Lambda \left(1 - \frac{M_{n-1}}{M_{max}} - \frac{M_{n-1}}{K_d - M_{n-1}} \right) \right) \quad (34)$$

0.5.5 Modèles aléatoires

La littérature est assez vaste sur les aspects stochastiques de la PCR. La représentation naturelle de l'amplification par PCR est un processus de branchement de type Galton-Watson : chaque molécule suit une évolution indépendante avec une probabilité d'être répliquée à chaque cycle dépendant de la taille de population (Stolovitzky and Cecchi, 1996; Velikanov and Kapral, 1999; Jagers and Klebaner, 2003; Lalam, 2006). L'adaptation commune de l'équation 25 est simplement donnée par la loi conditionnelle d'évolution donnée par l'équation 35, dont l'espérance correspond bien à l'équation 25 :

$$M_n | \mathcal{F}_{n-1} \sim M_{n-1} + \text{Binomial}(M_{n-1}, \lambda_n) \quad (35)$$

\mathcal{F}_{n-1} est une filtration adaptée à M_{n-1} , c'est-à-dire l'information observable au cycle $n - 1$ ¹¹. Le conditionnement ne dépend pas seulement de M_{n-1} dans le cas où plusieurs espèces sont présentes, ce qui n'est pas explicite ici mais l'est dans le Chapitre 2.

Les travaux sur les modèles aléatoires de PCR visent principalement à décrire la variabilité de la mesure, notamment lorsque le nombre de molécules au départ est faible (Piau, 2005). En général, les auteurs constatent que l'incertitude dépend de la taille de population initiale, par exemple sous la forme (Peccoud and Jacob, 1996) :

$$\frac{\text{écart-type}}{\text{espérance}}(M_n) \sim \frac{1}{\sqrt{M_0}} \quad (36)$$

mais qu'une approximation déterministe (en considérant une population de taille infinie) est pertinente du fait de la loi des grands nombres (Peccoud and Jacob, 1996; Jagers and Klebaner, 2003; Piau, 2004).

Dans mes travaux (Chapitre 2), j'ai utilisé des modèles aléatoires à des fins de simulation mais je n'ai pas cherché à estimer par le calcul les incertitudes générées.

Ces modèles aléatoires présentent l'intérêt d'incorporer naturellement des modèles de mutation (Sun, 1995; Wang et al., 2000; Pritchard et al., 2005; Piau, 2005) : chaque molécule nouvellement créée a une probabilité d'être différente de la séquence d'origine. Ces modèles de mutation sont abordés dans le Chapitre 3. Les mutations ont donné lieu à une approche par coalescent (Weiss, 1997) pour décrire le cas (avéré en pratique) où seulement un sous-échantillon de molécules est séquencé après la PCR. Les mutations sont négligées pour définir la dynamique de réaction car celles-ci sont rares (Boggy and Woolf, 2010) et il est raisonnable de penser que deux amplicons très similaires (une séquence mutée et une séquence souche) ont une efficacité d'amplification similaire.

11. Cela correspond à un formalisme mathématique précis mais que je ne détaille pas. La filtration naturelle $\mathcal{F}_{n-1} = \sigma(M_k^s, k \leq n-1, 1 \leq s \leq S)$ tenant compte de chacune des S espèces présentes est le choix logique.

0.5.6 Limites de la modélisation

Malgré l'intérêt théorique de ces nombreux modèles, les méthodes basées sur l'ajustement des données n'ont pas remplacé les méthodes de calcul d'efficacité par dilution sérielle car elles sont généralement moins précises (Karlen et al., 2007; Boggy and Woolf, 2010). Voici quelques éléments d'explication. Tout d'abord, la phase de plateau, connue pour être plus fluctuante que la phase exponentielle, peut induire des ajustements de modèle erronés (Rutledge and Stewart, 2008). Les tentatives de correction en tronquant la fin d'amplification ou en pondérant les données selon la phase de l'amplification n'ont pas été convaincants (Goll et al., 2006). Du fait du caractère exponentiel de la PCR, les modèles induisent une erreur faible sur l'estimation des efficacités Λ mais ces erreurs deviennent importantes lorsque Λ est convertie en quantité initiale M_0 (Tellinghuisen and Spiess, 2015). De plus, les efficacités estimées par les différents modèles sont parfois inconsistantes (Tellinghuisen and Spiess, 2014). Certains paramètres peuvent être non identifiables, différentes combinaisons de paramètres produisant la même cinétique. Cela rend plus complexe la quantification à partir des paramètres estimés.

0.5.7 Représentation de quelques modèles

La Figure 0.5.13 représente des données réelles de qPCR ainsi que sept modèles implémentés dans le package `qpcR` ou par mes soins : b_4 , l_4 , b_5 , l_5 pour les modèles analytiques ; $cm3$, le modèle logistique (Hayward, 1998) et le modèle mécanistique que nous avons développé ("nôtre" sur la courbe, Chapitre 2). Comme on le constate, tous les modèles s'ajustent bien à la courbe, à l'exception du modèle logistique en début de saturation. En revanche, il est difficile de donner un sens biologique cohérent aux différents paramètres. Cette figure est à mettre en relation avec la Figure 2.4.1 du Chapitre 2 où je représente les trois modèles comparés dans mes travaux sur le métabarcoding.

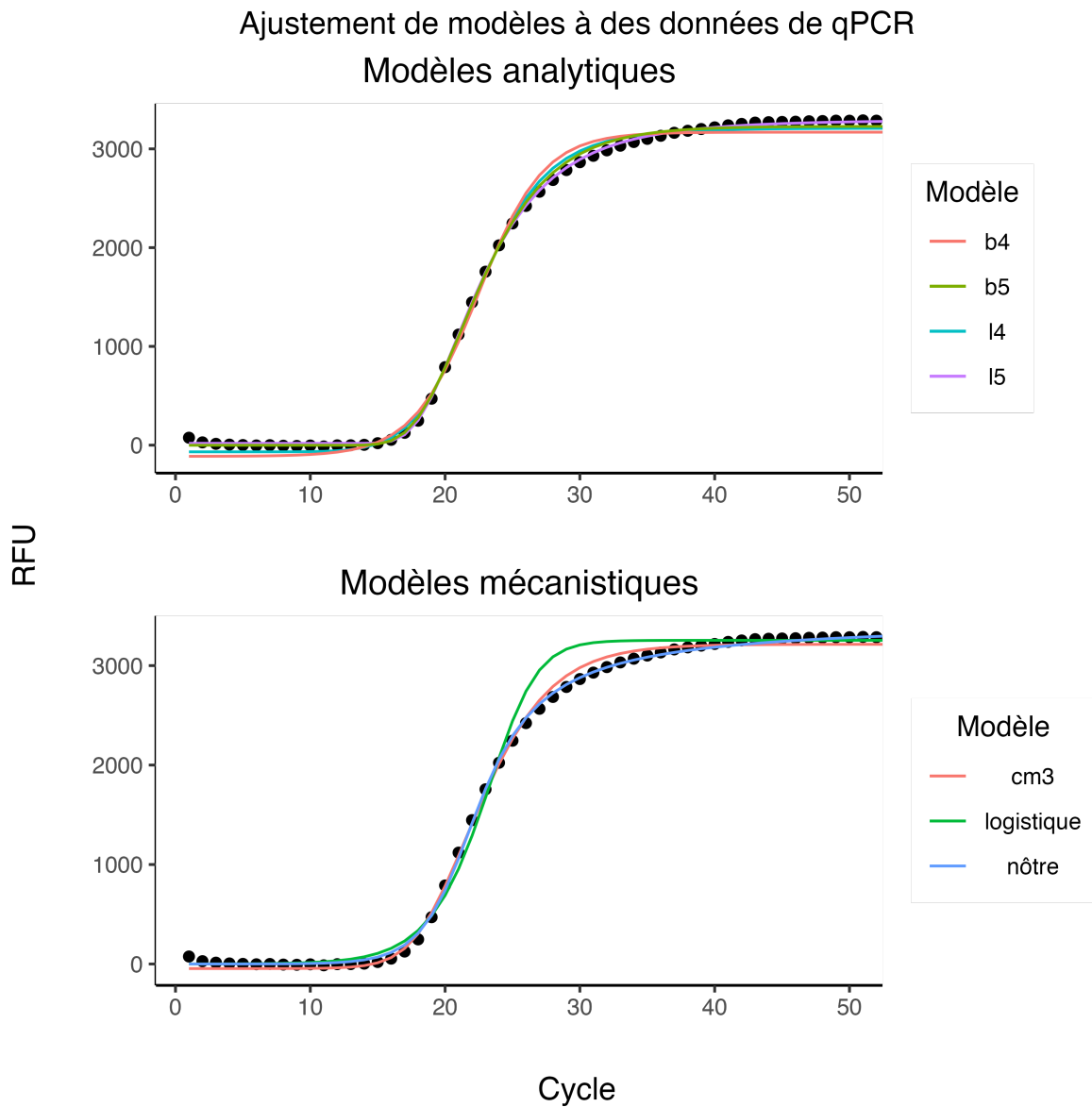


FIGURE 0.5.13 – Ajustement numérique de modèles analytiques (en haut) et mécanistiques (en bas) à des données de qPCR. L'ajustement a été optimisé pour accorder plus d'importance au début de l'amplification (vers le cycle C_t). L'asymétrie de la qPCR est apparente en comparaison du modèle logistique (courbe verte, panneau bas) qui est symétrique par rapport au point de demi amplification.

0.6 Métabarcoding quantitatif

Les indices de biodiversité reposent sur la richesse de l'écosystème et sur les abondances relatives de chaque espèce. Le métabarcoding a rapidement eu pour ambition de fournir une estimation de ces abondances relatives. Cette question a été la première motivation de ma thèse. Cette section établit un état de l'art des méthodes de quantification utilisées. Le Chapitre 2 présente le protocole expérimental et d'analyse que j'ai développé pour parvenir au métabarcoding quantitatif.

0.6.1 Enjeux

L'analyse naturelle consistant à assimiler les proportions des lectures des différents variants aux proportions réelles des espèces dans l'environnement s'avère mitigée. La corrélation entre lectures séquencées et biomasse ou nombre d'individus ne fait pas consensus (van der Loos and Nijland, 2021) et est souvent variable d'une étude à une autre (Lamb et al., 2019). Certaines études rapportent des corrélations positives (Rourke et al., 2022), d'autres des corrélations faibles ou l'absence de corrélation entre des marqueurs différents (Bell et al., 2017; Ershova et al., 2021).

Par ailleurs, de la variabilité est observée entre les réplicats biologiques, même si les échantillons sont prélevés à proximité (Levi et al., 2019; Steinke et al., 2021). De même, la comparaison inter-études mène parfois à des interprétations écologiques incohérentes entre différents types d'échantillon (piège Malaise et sol par exemple (Marquina et al., 2019)).

Ces différences sont dues à de nombreux biais qui interviennent au cours de l'expérience dont un des plus importants est le biais d'amplification par PCR. Au cours de ma thèse, j'ai cherché à mesurer et à corriger ce biais (entre autres) pour améliorer les résultats de métabarcoding quantitatif. Ces biais font l'objet d'un grand nombre d'études, dont Fonseca (2018), Lamb et al. (2019) ou Piñol et al. (2019) présentent des revues.

0.6.2 Biais successifs

Les biais sont soit d'origine biologique et donc inhérents à l'étude de l'ADNe, soit d'origine technique et sont créés au cours de l'expérience en laboratoire. Les quantités concernées successives sont :

- La biomasse,
- l'ADN déposé dans l'environnement,
- l'ADN environnemental échantillonné,
- l'ADN extrait,
- l'ADN cible amplifié par PCR,
- les lectures obtenues par séquençage.

Les erreurs aléatoires ou systématiques se propagent au cours de l'expérience (Kelly et al., 2019). Ces dernières sont induites par des méthodes dites de *high precision, low accuracy*, comme la PCR. La Figure 0.6.14 reprend le déroulement d'une expérience de métabarcoding et indique les biais à chaque phase.

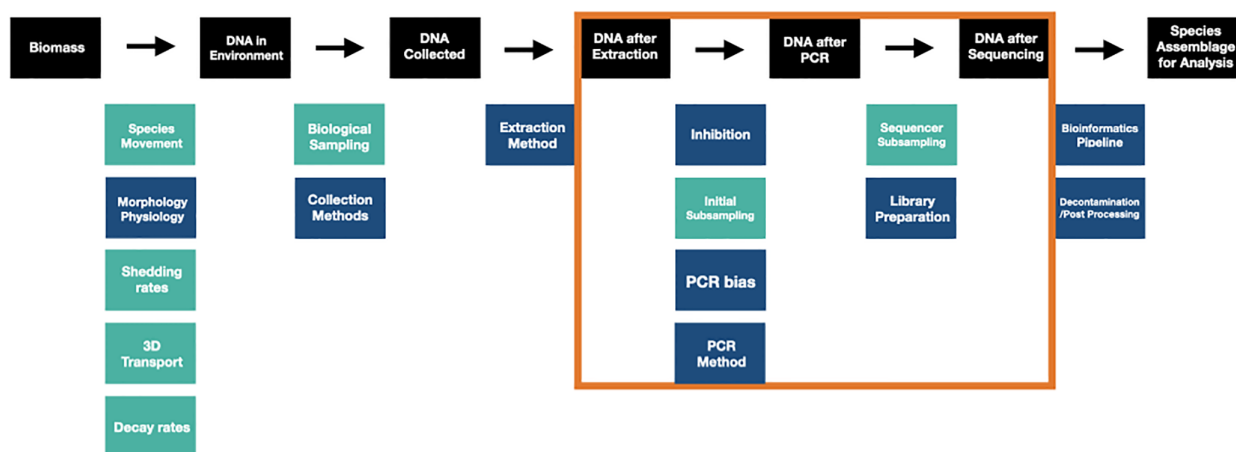


FIGURE 0.6.14 – Différentes quantités successivement traitées au cours du métabarcoding (en noir) et biais affectant la quantification des espèces. En bleu sont représentés les biais considérés comme systématiques par les auteurs et en vert les biais aléatoires. Gold et al. (2023b) et le Chapitre 2 de cette thèse proposent des modèles (indépendants) de la partie encadrée du processus. Source de la figure : Gold et al. (2023b).

0.6.3 Biais biologiques

Les biais biologiques ont déjà affecté les abondances relatives au moment où l'ADNe est échantillonné et doivent donc être traités en complément du métabarcoding.

0.6.3.1 Taux de dépôt et transport de l'ADN

La quantité d'ADN déposée dans l'environnement dépend sans surprise de la biomasse des individus (Yates et al., 2021) ou de la taille de l'organisme (Takahara et al., 2012).

Mais le taux de dépôt d'ADN dans l'environnement (*shedding rate*) dépend de nombreux autres facteurs dont Andruszkiewicz Allan et al. (2021); Wilder et al. (2023) établissent un inventaire détaillé. Par exemple, il est supposé que ce taux est plus faible pour les reptiles que pour les amphibiens du fait de leurs types de peau différents, ce qui complexifie l'analyse conjointe des deux groupes taxonomiques (Nordstrom et al., 2022). Chez les organismes aquatiques, ce taux dépend du stade de vie (Goldberg et al., 2011), du stress (Pilliod et al., 2014) et de la température de l'eau (Lacoursière-Roussel et al., 2016). Les études sur les insectes capturés par piège Malaise montrent aussi que les espèces ne sont pas toutes capturées avec la même probabilité (deWaard et al., 2019).

Il faut également prendre en compte que l'ADNe n'a pas forcément été déposé sur le site d'échantillonnage. Il est parfois transporté, par exemple dans les cours d'eau où l'on retrouve de l'ADNe déposé en amont (Jane et al., 2015; Deiner et al., 2017).

0.6.3.2 Taux de dégradation de l'ADN

Une fois déposé dans l'environnement, le taux de dégradation de l'ADN (*decay rate*) dépend de facteurs détaillés dans Andruszkiewicz Allan et al. (2021); Wilder et al. (2023). Il dépend du type d'ADN (extra ou intracellulaire) (Turner et al., 2014), de nombreux facteurs abiotiques (exposition au soleil, température, pH, salinité...) (Eichmiller et al., 2016; Barnes et al., 2014; Strickler et al., 2015; Seymour et al., 2018) et de facteurs biotiques comme la présence d'enzymes ou de micro-organismes (Barnes et al., 2014). Après la collecte, les conditions de conservation des échantillons influent aussi le taux de dégradation.

Le biais de dégradation est moins important que le biais d'amplification par PCR (Krehenwinkel et al., 2018).

0.6.3.3 Digestion

Les études de régimes alimentaires à partir de fèces ajoutent un biais : l'effet de la digestion sur l'ADN en fonction du tissu ingéré et de l'espèce (Thomas et al., 2014). Par ailleurs, les échantillons contiennent de l'ADN du prédateur en grande quantité, ce qui peut limiter la détection de certaines espèces (Cuff et al., 2023). L'utilisation d'un oligonucléotide bloquant l'amplification de la séquence du prédateur peut affecter d'autres codes-barres et induire un nouveau biais (Piñol et al., 2015).

0.6.3.4 Capacité de détection

Ces facteurs induisent un biais difficile à caractériser, notamment pour comparer plusieurs études. Ils influent la capacité de détection qui varie dans le temps (saisonnalité) et spatialement (Laramie et al., 2015; Li et al., 2019). Sassoubre et al. (2016) recense des études (principalement en eau douce) montrant que la quantité d'ADN passe sous le seuil de détection en 4 à 54 jours. Cet intervalle est à considérer avec précaution, comme l'attestent les études sur l'ADN ancien. Ces taux de dépôt, de dégradation et de détection d'ADN ont surtout été mesurés chez les macro-organismes aquatiques (Sassoubre et al., 2016; Klymus et al., 2015; Pilliod et al., 2014; Maruyama et al., 2014).

0.6.4 Biais techniques

Les biais techniques sont induits par le métabarcoding lui-même et peuvent être atténués par le design d'un protocole expérimental adéquat et lors du traitement des données.

0.6.4.1 Extraction

Des travaux font état de différences de qualité d'extraction entre taxons et entre différentes méthodes d'extraction (Schiebelhut et al., 2017; Dopheide et al., 2019; Martoni et al., 2022; Iwaszkiewicz-Eggebrecht et al., 2023). Pour les échantillons "en vrac", un facteur identifié est le ratio surface/volume des insectes (Marquina et al., 2019).

0.6.4.2 Concentration en ADN cible

L'ADN amplifié est une petite région du génome souvent issue d'organites : les chloroplastes chez les plantes (marqueurs rBCL, *Sper01*) et les mitochondries chez les animaux (marqueur COI). Le nombre de copies par cellule est variable d'une espèce à l'autre (Zoschke et al., 2007; Krehenwinkel et al., 2017; Garrido-Sanz et al., 2022) et d'un type de tissu à l'autre chez un même individu (Wiesner et al., 1992). Il peut aussi varier au cours du développement de la plante (Golczyk et al., 2014; Sakamoto and Takami, 2018). Le même résultat est observé chez les bactéries pour le gène 16S (Kembel et al., 2012).

Ainsi, la masse d'ADN dosé sur l'ensemble de la cellule est un indicateur biaisé de la quantité d'ADN effectivement amplifiée. Ce biais est étudié dans le Chapitre 2.

0.6.4.3 PCR

Le biais de PCR est identifié comme une limite majeure du métabarcoding (Taberlet et al., 2012). Il fait l'objet d'une étude approfondie au Chapitre 2. Les codes-barres moléculaires ne sont pas tous amplifiés avec la même efficacité pour un marqueur donné.

D'abord, des amorces présentent des mismatches avec certains codes-barres qui limitent l'hybridation et donc l'efficacité d'amplification par PCR (Krehenwinkel et al., 2017). C'est le cas notamment pour le marqueur COI (Elbrecht and Leese, 2015; Piñol et al., 2015). Ce marqueur est dégénéré, c'est-à-dire que certaines de ses bases sont variables, ce qui induit des mismatches récurrents. Les mismatches induisent un biais phylogénétique car certains groupes sont systématiques affectés (Liu et al., 2023).

La phase d'élongation de la PCR engendre aussi des biais. La longueur de la séquence, la structure secondaire de l'ADN (Pawluczyk et al., 2015), la répétition d'homopolymères et le taux de GC (Nichols et al., 2018) sont des facteurs de biais identifiés. Ceux-ci peuvent varier selon la polymérase utilisée.

Le biais d'hybridation intervient dans les premiers cycles où les molécules d'ADN matrice ont du mal à se lier aux amorces, alors que les nouveaux amplicons sont complémentaires de celles-ci. Le biais d'élongation, lui, affecte l'ensemble de la PCR et est fortement marqué pendant la phase exponentielle de l'amplification. L'effet du biais de PCR est important : le facteur entre les abondances avant et après PCR peut valoir au moins 4 sans mismatch (Silverman et al., 2021) et jusqu'à 10 en présence de mismatches (Parada et al., 2016).

0.6.4.4 Séquençage

À ma connaissance, le séquençage n'induit pas un biais contre certaines espèces. Porazinska et al. (2010) montre qu'il existe de la variabilité entre réplicats de séquençage mais que ceux-ci mènent à des résultats comparables. Le séquençage n'étant effectué que sur une petite fraction des amplicons, cela peut limiter la détection d'espèces rares.

0.6.5 Méthodes de correction

De nombreux travaux visent à améliorer l'estimation des abondances à partir d'ADNe. L'analyse peut porter sur l'abondance interspécifique ("l'espèce A est plus abondante que l'espèce B") ou intraspécifique ("l'espèce A est plus abondante dans l'échantillon 1 que dans le 2") (Luo et al., 2022). Certains protocoles adaptent le protocole expérimental tandis que d'autres se concentrent sur l'analyse bioinformatique. Mes travaux du Chapitre 2 portent sur la variabilité interspécifique.

0.6.5.1 Utilisation de PCR quantitative

Les techniques de PCR quantitative sont des outils communs pour quantifier la variabilité intraspécifique. Les études comparatives de la qPCR, de la ddPCR et du métabarcoding ont des conclusions variables. Un certain nombre observe que la qPCR (Harper et al., 2018; Chandelier et al., 2021) ou la ddPCR (Wood et al., 2019; Espinosa-Prieto et al., 2023) permettent une meilleure détection des espèces rares que le métabarcoding, mais d'autres études concluent l'inverse (McCarthy et al., 2022). Une corrélation entre les abondances établies par qPCR, ddPCR et métabarcoding est généralement observée (Rourke et al., 2022).

Plusieurs auteurs prônent donc des approches jointes utilisant la qPCR (Català et al., 2017) ou la ddPCR (Picard et al., 2022; Pont et al., 2023) pour améliorer les conclusions écologiques établies par métabarcoding.

Dans mes travaux, nous avons utilisé la ddPCR pour améliorer la quantification interspécifique en mesurant le nombre de copies de cible de différentes espèces à concentration d'ADN total ou à biomasse égale (Chapitre 2).

0.6.5.2 Modélisation

De nombreuses modélisations du métabarcoding ou des données d'ADNe ont été développées de sorte à décrire les différentes étapes et leurs biais et à faciliter l'analyse des données. Kelly et al. (2019) modélise le dépôt d'ADN dans l'environnement (biais systématique selon l'espèce), l'amplification par PCR (idem) et le séquençage (effet de sous-échantillonnage aléatoire). Le but de cette étude est d'estimer par des simulations l'importance relative des différents paramètres et d'observer dans quelle mesure les lectures sont représentatives de la biomasse. Gold et al. (2023b) ajoute une étape aléatoire caractérisant le nombre de molécules d'ADN extraites et applique ces travaux

pour évaluer le risque de non-détection d'espèces dans certains réplicats.

En général, le modèle utilisé pour décrire la PCR est le modèle exponentiel présenté en 0.5, équation 18 (Kelly et al. (2019); Gold et al. (2023b); Silverman et al. (2021); Shelton et al. (2022)...) Ce modèle s'interprète grâce à une régression linéaire, ce qui rend son utilisation simple (équation 15). Cet aspect est discuté en détail dans le Chapitre 2.

Pour décrire la quantité d'ADN déposé dans l'environnement en fonction de la biomasse, Yates et al. (2021) propose une relation allométrique supposant que le taux de dépôt par unité de masse diminue quand la masse augmente, selon la relation :

$$\text{ADN} = C \cdot \text{Masse}^b \quad (37)$$

où C est une constante et le paramètre b est compris entre 0 et 1. La valeur de b est de l'ordre de 0.75 mais doit être inférée et varie dans un intervalle assez grand.

En écologie des communautés, de nombreux modèles d'occupation estiment la probabilité de présence réelle des espèces à partir des données d'ADNe (Griffin et al., 2020; Jurburg et al., 2021; Burian et al., 2021). J'évoque ce sujet en 4.2.2.6.

0.6.5.3 Communautés de composition connue

Une idée récurrente est l'utilisation de communautés artificielles (*mock communities*) où les espèces d'intérêt sont mélangées en proportions connues pour reproduire un échantillon environnemental qui est ensuite traité par métabarcoding. Les modèles évoqués précédemment sont utilisés pour établir des facteurs correctifs.

Une première expérience est réalisée par Thomas et al. (2014) à partir de trois espèces de poisson présentes dans le régime alimentaire du phoque commun.

Thomas et al. (2016) complète ces travaux en construisant des communautés artificielles avec une espèce à tester et une espèce de référence en proportions de biomasse fraîche égales. L'étude établit des facteurs correctifs (ratio d'abondance entre l'espèce d'intérêt et la référence) applicables à des expériences ultérieures.

Des protocoles équivalents sont proposés en utilisant plusieurs communautés artificielles d'insectes de composition variable (Krehenwinkel et al., 2017) ou avec des communautés de deux espèces en proportions variables de biomasse de racine (20%-80%, 50%-50%, 80%-20%) pour analyser des communautés de plantes (Matesanz et al., 2019). McLaren et al. (2019) reprend une méthode similaire en décomposant la correction en une succession de facteurs correctifs (extraction, PCR, séquençage) pour des communautés bactériennes.

Shelton et al. (2022) va plus loin en utilisant des communautés artificielles (à partir de dosage d'ADN total par qPCR) avec différentes espèces aquatiques pour estimer les efficacités de PCR et leurs incertitudes à partir du modèle exponentiel. Cette approche

corrective a plusieurs similitudes avec mon projet du Chapitre 2, même si la mise en œuvre est différente.

Enfin, Gold et al. (2023a) reprend le protocole d'inférence de paramètres de Shelton et al. (2022) en couplant les données de métabarcoding à analyser avec un autre jeu de données de décomptes morphologiques établis sur les mêmes sites et les mêmes espèces.

Les résultats montrent une amélioration de la quantification grâce à une correction simple. En revanche, les communautés artificielles sont lourdes à réaliser et portent sur un nombre limité de taxons. Dans certains cas, la constitution de communautés artificielles est impossible, par exemple pour certaines bactéries (Browne et al., 2016). De plus, la correction est sensible aux erreurs expérimentales commises lors de l'élaboration de ces communautés.

0.6.5.4 Modification du protocole expérimental

Pour éviter ces inconvénients, Silverman et al. (2021) propose une approche sans communauté artificielle, dans une étude du microbiote intestinal avec un marqueur sans mismatch. Un même échantillon est amplifié dans plusieurs réplicats pendant 10 à 35 cycles. Les abondances observées pour chaque nombre de cycles sont ensuite utilisées pour estimer les efficacités d'amplification dans le modèle exponentiel (équation 18). Shelton et al. (2022) ne parvient pas à reproduire de résultats par cette méthode.

0.6.5.5 Contrôle interne (*spike-in*)

D'autres protocoles consistent à incorporer dans les échantillons des contrôles internes sous la forme de molécules d'ADN en quantité connue appelées *spike-in*. Ces molécules sont amplifiées par le même marqueur que le groupe taxonomique d'intérêt (Harrison et al., 2021), mais ces contrôles existent aussi en métagénomique (où il n'y a pas d'amplification) (Ji et al., 2020). Il peut s'agir de séquences synthétiques ou d'échantillons d'espèces réelles (*spike-in* biologiques). Le *spike-in* est généralement introduit avant l'extraction de l'ADN pour être traité de la même façon que l'échantillon.

En l'absence de tout biais, on s'attend à obtenir le même nombre de lectures de *spike-in* pour tous les réplicats. Pour corriger les biais expérimentaux, le nombre de lectures des espèces est divisé par le nombre de lectures du *spike-in* (Ji et al., 2020). Cela permet de quantifier la variabilité intraspécifique mais n'apporte pas d'information sur les biais propres à chaque espèce. Au cours de ma thèse, j'ai cherché à caractériser la variabilité interspécifique grâce un *spike-in* synthétique, mais ces travaux n'ont pas abouti à un résultat probant (section 2.4.3).

Harrison et al. (2021) propose une quantification absolue des abondances en multipliant les abondances initiales relatives estimées par le nombre de molécules de *spike-in* qui est mesuré au préalable (par cytométrie en flux (*flow cytometry*) Props et al. (2016)). Cet article fournit également une revue d'études ayant utilisé de tels contrôles internes.

Des *spikes-in* biologiques peuvent être utilisés dans les communautés artificielles pour estimer des facteurs correctifs d'abondance : c'est le cas des espèces de référence dans Thomas et al. (2016) et Iwaszkiewicz-Eggebrecht et al. (2023).

Une approche différente consiste à ajouter aux molécules d'ADN des tags aléatoires (de 7 à 12 bases quelconques) avant la PCR (Hoshino and Inagaki, 2017). À la fin de l'amplification, le nombre de tags différents par espèce est corrélé à l'abondance initiale et permet une correction similaire à un dosage par ddPCR (Hoshino et al., 2021).

0.7 Limites et perspectives du métabarcoding

Je dresse un bilan des performances actuelles du métabarcoding pour détecter et quantifier les espèces dans un écosystème et les alternatives proposées pour dépasser les limites identifiées.

0.7.1 Détection d'espèces

Les études basées sur l'ADNe dépassent souvent les méthodes de suivi traditionnelles en termes de détection d'espèces. Ces meilleures performances sont attestées pour la détection d'espèces (Valentini et al., 2016; Schneider et al., 2016) et pour la résolution (i.e., la différenciation entre plusieurs espèces) (Kraaijeveld et al., 2015). En revanche, le métabarcoding peut ne pas être suffisant pour répondre aux questions écologiques et est alors complémentaire de méthodes traditionnelles (Nordstrom et al., 2022), par exemple pour prendre en compte le transport de l'ADN dans l'environnement.

Par ailleurs, les méthodes d'analyse de l'ADNe sont réputées moins chères que les méthodes traditionnelles car les inventaires morphologiques sont longs à réaliser.

0.7.2 Détermination des communautés

Comme je l'ai évoqué plus haut, l'acquisition de données quantitatives fiables par métabarcoding ne fait pas encore l'objet d'un protocole unifié.

Une des ambitions actuelles du métabarcoding est la détermination exhaustive de la composition des communautés. Pour cela, toutes les espèces doivent être détectées, idéalement avec des abondances représentatives de la réalité. Dans ce but, Ficetola and Taberlet (2023) établit une liste de méthodes disponibles :

- la combinaison de plusieurs marqueurs au sein d'une même étude pour couvrir différents groupes taxonomiques ;
- l'utilisation de marqueurs universels permettant d'étudier toutes les espèces simultanément (au risque que certains groupes soient moins bien amplifiés) ;
- l'utilisation conjointe de marqueurs spécifiques et universels pour avoir un aperçu global et une étude précise des groupes clés ;

- l'utilisation simultanée de plusieurs amorces ;

ou encore la suppression de l'amplification par PCR, pointée du doigt depuis les débuts du métabarcoding (Taberlet et al., 2012) :

- l'usage de la métagénomique (séquençage shotgun) ;
- l'enrichissement par capture précédant le séquençage shotgun.

Ainsi, les solutions proposées ici aux limites du métabarcoding sont expérimentales : le protocole est modifié ou optimisé pour améliorer la qualité des données. La section 0.6 présentant les approches quantitatives du métabarcoding, aborde aussi les améliorations proposées pour l'analyse bioinformatique des données.

0.7.3 Alternatives au métabarcoding

Plusieurs techniques alternatives ont été développées, entre autres, pour éviter l'étape d'amplification par PCR.

0.7.3.1 Amplification isotherme de l'ADN

Les techniques d'amplification isotherme de l'ADN sont des alternatives à la PCR qui ne recourent pas à des cycles de températures, et ne nécessitent donc pas de thermocycleur. Cela les rend plus simples et plus abordables (Bartholomew et al., 2015). Une des méthodes les plus connues est l'Amplification Isotherme Médinée par les Boucles (*Loop-mediated isothermal amplification*, LAMP) (Notomi, 2000). Elle peut être quantitative (Hardinge and Murray, 2020) à l'instar des méthodes de PCR quantitative présentées en 0.4, sans être exempte de biais (Liu et al., 2017). D'autres méthodes existent, dont Gill and Ghaemi (2008) présente une revue : l'Amplification par Réplication Circulaire (*Rolling Circle Amplification*, RCA) (Ali et al., 2014), l'Amplification par Polymérase Recombinase (*Recombinase Polymerase Amplification*, RPA) (Piepenburg et al., 2006), l'Amplification de Séquence d'ADN (*Nucleic acid sequence-based amplification*, NASBA) (Hønsvall and Robertson, 2017)... L'amplification isotherme est compatible avec les nouvelles méthodes de séquençage qui remplacent Illumina.

0.7.3.2 Nouvelles méthodes de séquençage

Les techniques de séquençage évoluent encore, ce qui permet d'envisager de nouveaux protocoles pour le métabarcoding. En particulier, le séquenceur de poche Nanopore MinION d'Oxford Technology, commercialisé en 2014, permet d'effectuer un séquençage rapide sur le terrain pour des expériences de dimensionnement variable (Krehenwinkel et al., 2019). Le fait que le MinION soit portatif permet de mener des expériences de métabarcoding dans des régions reculées et d'améliorer l'accès à ces techniques dans des pays où peu de moyens de séquençage sont disponibles. Trois jours sont nécessaires pour établir une librairie, ce qui est bien moins que les délais usuels pour le séquençage Illumina pour lequel les échantillons sont envoyés à des prestataires tels que l'entreprise Fasteris. Le débit du MinION est moindre par rapport à Illumina mais il peut séquencer de longues molécules (plusieurs milliers de paires de bases). Son

inconvenient principal est son fort taux d'erreur (quelques pour cent). Son utilisation est souvent couplée à des méthodes d'amplification isotherme (LAMP notamment) afin d'améliorer la qualité des séquences consensuelles établies. Des expériences ont été menées sur des fragments d'ADN très courts (moins de 100 bases), à l'image des codes-barres moléculaires (Wilson et al., 2019). Des travaux récents confirment que cette technologie est adaptée au métabarcoding (Baloglu et al., 2021).

0.7.3.3 Métagénomique

La métagénomique s'affranchit de l'amplification par PCR en séquençant l'ensemble de l'ADN collecté, selon la méthode dite shotgun. Le séquençage de l'ensemble de l'ADN permet d'obtenir une information plus importante que les seuls codes-barres moléculaires sur lesquels se concentre le métabarcoding. Taberlet et al. (2012) estime la part informative de l'ADN global à environ 1 à 10 %. Le séquençage shotgun est particulièrement intéressant pour l'ADN ancien (Murchie et al., 2021) car il permet d'authentifier l'âge de l'ADN (Eisenhofer and Weyrich, 2018). En revanche, l'analyse des données est complexe et dépend grandement des bases de données de génomes, qui se concentrent aujourd'hui sur quelques espèces. La métagénomique est également très coûteuse en comparaison du métabarcoding (Ficetola and Taberlet, 2023) mais des résultats montrent une amélioration de la quantification par rapport au métabarcoding (Lang et al., 2019; Pierella Karlusich et al., 2022).

Un inconvenient du séquençage shotgun brut est son absence de spécificité, puisque l'ensemble de l'ADN collecté est séquençé. Des méthodes ont été développées pour enrichir l'échantillon en ADN d'intérêt et faciliter l'analyse en réduisant le coût d'acquisition de l'information génétique pertinente.

0.7.3.4 Enrichissement par capture

L'enrichissement par capture précède le séquençage shotgun et permet d'augmenter la spécificité de l'étude. La méthode classique est celle de la capture hybride. Elle consiste à fragmenter l'ADN puis à ajouter dans l'échantillon des sondes (une par séquence cible) qui se lient aux fragments ciblés. Le prélèvement est ensuite nettoyé pour ne garder que l'ADN lié à une sonde qui est alors séquençé (Singh, 2022). En revanche, cette technique nécessite des bases de données de qualité et le design de sondes pour chaque espèce d'intérêt.

0.8 Outils mathématiques en écologie et biologie

Même si les problématiques liées au métabarcoding quantitatif sont la motivation première de ma thèse, j'ai été amené à réfléchir à des sujets plus généraux relevant des statistiques appliquées à l'écologie au sens large.

Du fait de mon cursus, mon premier réflexe pour traiter un problème en écologie est de me tourner vers les mathématiques. C'est un positionnement commun et l'interface entre ces disciplines est aujourd'hui très développée. Il me semble important d'aborder les interactions entre mathématiques et écologie dans cette introduction. Le Chapitre 1 de ma thèse aborde le développement d'un nouvel outil statistique applicable à un large ensemble de problèmes d'écologie. De même, mes travaux relatifs au métabarcoding (Chapitres 2 et 3) reposent en bonne partie sur des modèles et des algorithmes que j'ai étudiés ou développés. Je ne présente pas une revue des outils mathématiques utilisés en écologie — c'est un sujet bien trop vaste — mais me concentre sur les algorithmes d'inférence pour modèles stochastiques pour donner le contexte d'étude du Chapitre 1. Je détaillerai dans l'introduction de ce chapitre plusieurs aspects techniques de mathématiques appliquées utiles à la compréhension de ces travaux.

0.8.1 Motivations

Avant de présenter les méthodes d'inférence, je donne quelques motivations et repères historiques sur l'usage des mathématiques en écologie ou biologie, en particulier au sujet de la modélisation qui est un thème récurrent de ma thèse. Ces réflexions sont inspirées entre autres de Legay (1997) et de l'introduction du cours *Modèles aléatoires et Évolution* de Sylvie Méléard (École polytechnique).

Une définition possible d'un modèle est un substitut à un objet réel qu'il n'est pas possible d'étudier directement. Les raisons peuvent en être multiples, surtout pour l'étude du vivant : le système étudié peut être trop complexe, trop lent à évoluer ou encore soumis à des interdits éthiques ou légaux (en anatomie par exemple). Le chercheur ou la chercheuse utilise donc une imitation du réel qui permet de contourner ces difficultés. Tous les modèles ne sont d'ailleurs pas mathématiques : la drosophile est un organisme modèle en génétique.

Une spécificité de la biologie et de l'écologie est la nature essentiellement aléatoire des systèmes qu'elles décrivent. Ceux-ci résultent de comportements globaux d'une population, d'un ensemble de cellules... dont chaque unité élémentaire évolue de manière aléatoire, au sens que son évolution est trop complexe pour être prise en compte précisément à l'échelle macroscopique. Cette variabilité individuelle est un enjeu majeur et le travail de modélisation consiste en un compromis entre une représentation fidèle de l'objet d'étude et la simplicité d'analyse.

Les modèles poursuivent des buts variés : ils servent à mieux comprendre, expliquer ou généraliser un phénomène, à tester des hypothèses, à prédire des comportements...

Les modèles biologiques ont émergé plus tard que les modèles physiques du fait de la grande diversité et complexité des systèmes étudiés, ce qui rendait difficile l'établissement de théories générales. Dans le *Discours de la méthode* (1637), Descartes rejette l'analyse des systèmes complexes par manque de moyens pour les étudier. Il y suggère de décomposer ceux-ci autant que nécessaire en problèmes simples pour que chacun réponde à un critère d'évidence, c'est-à-dire de vérité incontestable. On cherche alors à décrire les problèmes sous la forme : "une cause \Rightarrow un effet". C'est loin d'être évident en biologie car les systèmes sont rarement décomposables. Et de fait, les premiers modèles en écologie sont très simples, comme celui de Malthus dans *An Essay on the Principle of Population* (1798) pour décrire la dynamique de population selon la production de ressources. Plus tard, le développement des statistiques a apporté un nouveau point de vue en acceptant de ne considérer que quelques facteurs contrôlés à côté de facteurs non contrôlés qui peuvent être innombrables et complexes. Le schéma "une cause \Rightarrow un effet" est ainsi abandonné et le critère d'évidence de Descartes est remplacé par un réseau de cohérences : la conviction se base désormais sur un faisceau d'indices et non plus sur une vérité allant de soi. Un exemple typique est le développement de l'analyse de variance (ANOVA) par Fisher en 1925 (Fisher, R. A., 1935) pour évaluer la productivité de pommiers. Dans cette étude, la variété des arbres est un facteur contrôlé sur des parcelles randomisées pour distinguer son effet du type de sol, de l'exposition au soleil, de l'exploitation agricole locale, etc. Plus récemment, l'informatique a accéléré le développement des modèles parallèlement à l'acquisition de données massives, par exemple en génomique. L'usage du *machine learning* est une révolution dans la mesure où un modèle peut être développé automatiquement sans qu'un expérimentateur n'ait formulé une problématique à laquelle les données doivent apporter des éléments de réponse.

Une fois qu'un modèle est développé, il est courant de l'utiliser pour simuler le phénomène imité. Deux axes d'étude se présentent alors : la justification du modèle par des procédés statistiques (tests d'adéquation...) et l'action sur le modèle pour le rendre le plus fidèle possible à la réalité, c'est-à-dire à des données observées. C'est surtout ce second objectif d'optimisation des modèles qui m'a intéressé durant ma thèse, c'est pourquoi je présente un tour d'horizon des algorithmes qui répondent à cette problématique.

0.8.2 Inférence de paramètres pour modèles stochastiques

0.8.2.1 Cadre général

En statistique, l'inférence désigne les méthodes qui permettent d'acquérir de l'information sur un phénomène à partir d'observations partielles, typiquement pour estimer des caractéristiques d'une population à partir d'un sous-échantillon de quelques individus.

Dans ma thèse, je m'intéresse au problème de l'inférence de paramètres pour des modèles stochastiques dont je présente d'abord le formalisme mathématique. En guise d'illustration, considérons l'exemple d'un modèle aléatoire de PCR présenté en 0.5 : une réalisation du modèle, sachant les paramètres θ de la PCR, est le nombre de molécules de deux espèces présentes.

On considère un ensemble mesurable (Ω, \mathcal{A}) , où Ω est un univers et \mathcal{A} la tribu des événements. Les éléments de Ω sont les réalisations possibles du modèle, par exemple : "l'espèce 1 compte 10 molécules et l'espèce 2, 5 molécules". Dans ce texte, Ω sera toujours un sous-ensemble de \mathbb{R}^k , $k \geq 1$. Dans mon exemple, $k = 2$. Les éléments de \mathcal{A} , les événements possibles, sont des sous-ensembles de Ω , par exemple : "Toutes les issues possibles du modèle où l'espèce 1 est plus abondante que l'espèce 2".

Le modèle est défini par une famille de distributions de probabilité \mathbb{P}_θ . Les paramètres θ sont choisis dans un ensemble $\Theta \subset \mathbb{R}^m$, $m \geq 1$.

Une réalisation du modèle sachant les paramètres est simplement une variable aléatoire X_θ de loi \mathbb{P}_θ . Ici, les distributions considérées sont discrètes ou à densité. La fonction de masse ou de densité de ces lois, respectivement, est notée p_θ ou p selon le contexte. Je me contenterai de parler de densité par la suite mais ces mentions incluent la fonction de masse des lois discrètes.

Dans l'exemple, on peut choisir $m = 4$: les quatre paramètres seraient les quantités initiales et les efficacités de PCR des deux espèces. Le nombre de paramètres m est quelconque mais plus ceux-ci sont nombreux, plus le modèle est difficile à analyser. C'est le fléau de la dimension (ou *curse of dimensionality*) qui fait références aux difficultés qui émergent lorsque la dimension d'un problème augmente.

Le problème d'inférence est alors : ayant observé des variables $X_\theta^{(1)}, \dots, X_\theta^{(k)}$ de paramètres θ inconnu, quelle est la valeur de θ ?

La difficulté du contexte aléatoire est qu'un même jeu de paramètres peut générer toute une gamme de réalisations. On cherche donc à construire un estimateur $\hat{\theta}$ de θ le plus "juste" possible. L'estimateur idéal n'est pas biaisé, a une faible variabilité autour de la vraie valeur, converge vers celle-ci lorsque le nombre d'observations augmente et est robuste face aux données extrêmes.

Un choix apprécié est $\hat{\theta}$ maximisant la vraisemblance du modèle car cet estimateur

vérifie un certain nombre de garanties théoriques. La fonction de vraisemblance du modèle sachant une donnée observée $X = x$ est une fonction des paramètres θ , définie par la densité de la loi et notée $p_\theta(x)$ ou $p(x; \theta)$. Elle donne la plausibilité que les paramètres du modèle soient θ sachant que la donnée $X = x$ a été observée. Il est souvent utile de considérer la log-vraisemblance $l = \log p$.

Lorsque cette vraisemblance n'est pas accessible (incalculable, inexploitable mathématiquement ou numériquement du fait de sa complexité), d'autres stratégies doivent être mises en place. Une autre quantité peut être maximisée (on parle de M-estimateur), ou on procède par la méthode des moments. Celle-ci consiste à choisir les paramètres qui permettent d'égaliser les moments théoriques et empiriques de la distribution étudiée (espérance, variance...).

De nombreux algorithmes permettent de traiter cette classe de problèmes et le Chapitre 1 présente le développement d'une nouvelle méthode. Je dresse un état de l'art des principales méthodes d'inférence utilisées en statistiques. Cette revue n'est pas exhaustive et mon objectif est de donner les intuitions de ces méthodes, sans rentrer dans un haut niveau de détails théoriques et d'implémentation.

0.8.2.2 Catégories d'algorithmes

Dans ce qui suit, les données observées sont des réalisations indépendantes d'un modèle, notées $\mathbf{x} = (x_1, \dots, x_n)$, $n \geq 1$.

Fréquentiste ou bayésien ? Les méthodes d'inférence sont divisibles en deux grandes catégories : les statistiques fréquentistes et bayésiennes (Shoemaker et al., 1999). Les premières considèrent une probabilité comme la fréquence de l'événement étudié lorsque l'expérience est répétée un grand nombre de fois. Les secondes traitent une probabilité comme une mesure de l'incertitude.

Les méthodes fréquentistes visent à construire un estimateur ponctuel des paramètres (maximum de vraisemblance par exemple). Les méthodes bayésiennes, elles, supposent que le paramètre est une variable aléatoire dont on estime une densité de probabilité a posteriori en incluant une connaissance a priori de l'expérimentateur (le prior). L'inférence bayésienne repose sur la formule de Bayes pour déterminer cette densité postérieure :

$$p(\theta|\mathbf{x}) = \frac{p(\mathbf{x}|\theta) p(\theta)}{p(\mathbf{x})} \quad (38)$$

où $p(\mathbf{x}|\theta)$ est la vraisemblance, $p(\theta)$ le prior et $p(\mathbf{x})$ l'évidence (une constante de normalisation en pratique).

Un premier embranchement dans le choix de la méthode d'inférence est la question de la vraisemblance : celle-ci est-elle accessible, peut-on facilement la calculer et la

manipuler pour extraire de l'information sur le modèle ? Si c'est le cas, les algorithmes construisant un estimateur du maximum de vraisemblance sont la meilleure solution. Sinon, des alternatives doivent être mises en place. Une option commune est d'étudier le modèle sur la base de simulations, dans l'idée de construire une distribution empirique ou d'extraire des caractéristiques du modèle pour des paramètres donnés. Hartig et al. (2011) expose les différences conceptuelles des deux approches.

0.8.2.3 Enjeux en pratique

L'implémentation de ces algorithmes doit certes fournir des estimateurs robustes, mais le temps de calcul est aussi un critère crucial, surtout pour les problèmes complexes ou de grande dimension. Pour cela, les choix d'implémentation (distributions approchées, hyperparamètres, statistiques résumées...) font partie intégrante de la tâche d'inférence.

Il se peut aussi que les modèles ne soient pas identifiables, c'est-à-dire que des jeux de paramètres différents donnent des résultats identiques. Cette problématique fait l'objet d'études spécifiques (Gustafson, 2014).

0.8.2.4 Méthodes avec vraisemblance

Dans la plupart des applications, la vraisemblance n'est connue sous une forme analytique explicite. Un cas fréquent est sa connaissance à une constante de normalisation près, typiquement l'évidence du modèle $p(\mathbf{x})$ dans l'équation 38. Je décris d'abord deux méthodes de Monte Carlo qui simulent des données à partir de la connaissance partielle de la vraisemblance. Elles construisent ainsi des estimateurs du maximum de vraisemblance (approche fréquentiste) ou du maximum de la distribution postérieure (approche bayésienne). Luengo et al. (2020) établit une revue très complète de ces méthodes de Monte Carlo.

Méthodes de Monte-Carlo par chaînes de Markov (MCMC) Les méthodes de Monte-Carlo par chaînes de Markov (*Markov chain Monte Carlo*, MCMC) sont un outil majeur en statistique bayésienne. Elles permettent d'échantillonner dans une distribution lorsqu'une approche directe n'est pas possible en simulant une chaîne de Markov dont la distribution stationnaire est la distribution à simuler.

L'algorithme de Metropolis-Hastings (Metropolis et al., 1953; Hastings, 1970) est à l'origine des MCMC. On suppose que l'on connaît la distribution visée π à une constante multiplicative près. On choisit un noyau de transition g de la chaîne de Markov tel que $g(x'|x)$ est la probabilité de passer de x à x' .

À l'itération k , la chaîne de Markov est dans un état θ_k . On tire un paramètre θ' selon $g(\cdot|\theta_k)$. θ' est accepté comme valeur de θ_{k+1} avec une probabilité $\frac{\pi(\theta')g(\theta_k|\theta')}{\pi(\theta_k)g(\theta'|\theta_k)}$. Dans ce cas, θ_{k+1} prend la valeur θ' , sinon $\theta_{k+1} = \theta_k$. On obtient ainsi un ensemble de paramètres acceptés dont on peut extraire une valeur la plus vraisemblable.

L'autre grande méthode historique de MCMC est l'échantillonnage de Gibbs (Geman and Geman, 1987) qui est un cas particulier de l'algorithme de Metropolis-Hastings adapté en grande dimension. Dans ce cas, les coordonnées de la chaîne de Markov sont mises à jour une à une, successivement.

Ces algorithmes sont très généraux et connaissent de nombreuses variations passées en revue dans Luengo et al. (2020). La dépendance à la vraisemblance des méthodes MCMC est contournée dans un certain nombre d'algorithmes, comme les méthodes MCMC-ABC présentées plus loin.

Échantillonnage préférentiel L'échantillonnage préférentiel (*importance sampling*) permet d'estimer les moments d'une distribution. Le principe est de simuler des données selon une loi simple f (dont le choix a une grande importance) et de donner un poids à chaque simulation θ égal à $\frac{f(\theta)}{\pi(\theta)}$, où π est encore la distribution cible. Cette méthode permet de réduire la variance des estimateurs. Là encore, Luengo et al. (2020) propose un état de l'art détaillé des méthodes actuelles.

J'en profite pour introduire les méthodes de Monte Carlo Séquentielles (*Sequential Monte Carlo*, SMC), qui consistent en un échantillonnage préférentiel itératif avec ré-échantillonnage successif de valeurs selon leurs poids à l'itération donnée (Del Moral et al., 2006). Cette approche est utilisée par les méthodes sans vraisemblance (voir les SMC-ABC plus loin).

0.8.2.5 Cas des modèles de Markov cachés

Une classe de modèles est particulièrement fréquente : les modèles de Markov cachés, où une partie de l'information n'est pas observée. Il s'agit de chaînes de Markov dont la vraisemblance est difficilement accessible à moins d'avoir accès à des variables cachées (ou latentes) $\mathbf{z} = (z_1, \dots, z_n)$, $n \geq 1$. Le calcul de la log-vraisemblance jointe $l(\mathbf{x}, \mathbf{z}; \theta)$ est alors plus facile. Ces variables cachées décrivent généralement des états du système, discrets ou continus. La Figure 0.8.15 illustre ces modèles. Un cas typique est une population dans laquelle plusieurs groupes sont mélangés. Étudier les caractéristiques d'un individu sans connaître son groupe d'origine est difficile mais est trivial si on le connaît. Je présente deux méthodes pour traiter ces modèles.

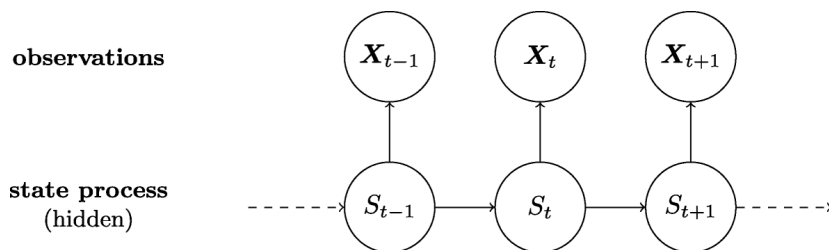


FIGURE 0.8.15 – Représentation générique d'un modèle de Markov caché. Des observations X_t sont émises à partir des états cachés S_t (notés z_t dans mon texte). Les flèches indiquent une dépendance. Source de la figure : Popov et al. (2019).

Algorithme d'Espérance-Maximisation (EM) L'algorithme Espérance-Maximisation (EM) est un algorithme itératif de recherche du maximum de log-vraisemblance d'un modèle de Markov caché (Dempster et al., 1977) :

$$\begin{aligned}\hat{\theta} &= \operatorname{argmax}_{\theta \in \Theta} l(\mathbf{x}; \theta) \\ &= \operatorname{argmax}_{\theta \in \Theta} \sum_{i=1}^n l(x_i; \theta)\end{aligned}\quad (39)$$

L'algorithme consiste en une répétition de deux étapes. À l'itération k :

- **Étape E** : À partir de l'estimation courante θ_k , calculer l'expression de :

$$Q(\theta; \theta_k) = \mathbb{E}_{\theta_k}[l(\mathbf{x}, \mathbf{z}; \theta)] = \int_{\mathcal{Z}} \log p(\mathbf{x}, \mathbf{z}; \theta) p(\mathbf{z}|\mathbf{x}; \theta_k) d\mathbf{z} \quad (40)$$

où \mathcal{Z} est l'ensemble de définition des variables latentes \mathbf{z} et où p est utilisée indifféremment pour désigner les deux densités.

- **Étape M** : Mettre à jour θ_k en maximisant

$$\theta_{k+1} = \operatorname{argmax}_{\theta \in \Theta} Q(\theta; \theta_k) \quad (41)$$

L'algorithme EM converge vers un maximum (local) de la vraisemblance (Wu, 1983). Plusieurs variantes existent, notamment lorsque l'étape de maximisation est compliquée ou que $Q(\theta; \theta_k)$ est difficile à calculer. J'introduis l'algorithme *Stochastic Approximation Expectation-Maximisation* (SAEM) (Kuhn and Lavielle, 2004, 2005) qui permet d'approcher $Q(\theta; \theta_k)$. L'étape M est inchangée et l'étape E est remplacée par :

- Simuler des variables latentes $\mathbf{z}^{(k+1)}$ selon $p(\cdot|\mathbf{x}; \theta_k)$;
- Mettre à jour l'estimation de Q :

$$\widehat{Q}(\theta, \theta_{k+1}) = \widehat{Q}(\theta, \theta_k)(1 - \gamma_k) + \gamma_k \log p(\mathbf{z}^{(k+1)}, \mathbf{x}; \theta) \quad (42)$$

où (γ_k) est une suite décroissante telle que $\sum \gamma_k = +\infty$ et $\sum \gamma_k^2 < +\infty$.

L'algorithme SAEM est abordé dans le Chapitre 1 du fait de son utilisation de simulations.

Algorithme forward Un autre algorithme, l'algorithme *forward* est utilisé pour estimer la probabilité d'un état (la variable \mathbf{z}) à un temps donné. Plus précisément, il permet de calculer la densité jointe des données jusqu'à l'étape k et de l'état caché du système à k :

$$\alpha_k(z_k) = p_\theta(x_1, \dots, x_k, z_k) \quad (43)$$

Les valeurs de α_k sont établies par récurrence :

$$\alpha_0(z_0) = p(z_0) \text{ à déterminer selon le modèle} \quad (44)$$

$$\alpha_{k+1}(z_{k+1}) = p(x_{k+1}|z_{k+1}) \int_{\mathcal{Z}} p(z_{k+1}|z) \alpha_k(z) dz \quad (45)$$

où l'intégrale (ou somme pour des états discrets) est calculée sur tous les états possibles z_k . Cette formule permet de substituer une intégrale à m dimensions, le nombre d'étapes de la chaîne de Markov, par m intégrales à une dimension, ce qui est nettement moins coûteux en termes de calcul.

Dans le cas de l'inférence de paramètres, le calcul de la vraisemblance des observations est a priori suffisant. Mais il est possible d'estimer la probabilité de l'ensemble des états cachés du modèle à partir de l'algorithme *forward-backward* qui inclut une autre récurrence en partant de l'état final du système.

0.8.2.6 Méthodes sans vraisemblance basées sur des simulations

Lorsque la vraisemblance n'est pas disponible, plusieurs méthodes alternatives, présentées dans cette section, la remplace par des simulations aléatoires. Cranmer et al. (2020) expose ce concept et liste les méthodes actuelles. La méthode présentée dans le Chapitre 1 appartient à ce cadre.

L'information des simulations peut être utilisée de deux manières (Drovandi and Frazier, 2022) : soit l'ensemble des données produites est utilisé, soit celles-ci sont agrégées en statistiques résumées (*summary statistics*), comme les moments empiriques, pour simplifier l'analyse. On parle alors de réduction de l'information. La première approche peut être coûteuse en temps de calcul et est plus sensible au fléau de la dimension mais n'induit pas de perte d'information. La seconde évite ces limites mais le choix des statistiques résumées est une étape cruciale et complexe.

Approximate Bayesian Computation (ABC) Les méthodes d'*Approximate Bayesian Computation* (ABC) ont pour but d'estimer la vraisemblance postérieure en simulant des réalisations du modèle (Sisson et al., 2018). L'idée est d'établir une distribution postérieure empirique à partir de simulations.

La première implémentation d'ABC est un algorithme de rejet (Tavaré et al., 1997; Pritchard et al., 1999). Pour chaque simulation, un paramètre θ est tiré aléatoirement selon le prior. Pour ce paramètre, on simule une réalisation du modèle X_θ . La simulation est acceptée selon sa similarité aux données mesurée par ρ , i.e. si $\rho(X_\theta, \mathbf{x}) \leq \epsilon$ pour une tolérance $\epsilon \geq 0$ donnée. L'ensemble des valeurs θ acceptées

forme la distribution postérieure.

Lorsque le modèle n'est pas basique, il est fastidieux de produire des simulations proches des données. On utilise donc des statistiques résumées s et le critère d'acceptation devient :

$$\rho(s(X_\theta), s(\mathbf{x})) \leq \epsilon \quad (46)$$

Ces statistiques peuvent être des moments, des statistiques d'ordre... choisies de sorte à être informatives et faciles à calculer. Leur nombre doit être limité sous peine d'augmenter le temps de calcul nécessaire à produire des simulations proches des données, en conséquence du fléau de la dimension. Leur choix est un facteur important pour les performances de l'algorithme (Fearnhead and Prangle, 2012). Des travaux plus récents reconsidèrent la possibilité de traiter des données complètes au lieu de statistiques résumées (Drovandi and Frazier, 2022).

Beaumont et al. (2002) améliore l'algorithme en remplaçant le rejet des simulations par une pondération selon $\rho(s(X_\theta), s(\mathbf{x}))$ et en ajustant les paramètres acceptés par régression. Cette seconde idée est typique des améliorations proposées dont Marin et al. (2012) fournit une revue. À partir de cette première version, deux grandes classes d'ABC ont émergé (Blum and François, 2010) : les *Monte Carlo Markov Chain ABC* (MCMC-ABC) et les *Sequential Monte Carlo ABC* (SMC-ABC). Les méthodes ABC les plus efficaces actuellement sont des SMC-ABC.

Les méthodes MCMC-ABC ont été développées pour contourner une limite des ABC : lorsque le prior est peu informatif, une part importante des paramètres évalués se situe dans une région de faible vraisemblance de l'ensemble Θ . Les MCMC-ABC adaptent donc les MCMC sans utiliser la vraisemblance. La première version est due à Marjoram et al. (2003) puis a été améliorée par Wegmann et al. (2009). L'algorithme utilise le même cadre que l'algorithme de Metropolis-Hastings présenté plus haut. À partir d'un paramètre θ_k , un paramètre θ' est tiré selon $g(\cdot|\theta_k)$ et utilisé pour simuler une réalisation du modèle x' . Celle-ci est comparée aux données selon $\rho(s(x'), s(\mathbf{x})) \leq \epsilon$, comme précédemment.. Si cette condition est vérifiée, θ' est accepté avec une probabilité $\frac{p(\theta')g(\theta_k|\theta')}{p(\theta_k)g(\theta'|\theta_k)}$, où p est un prior.

D'autre part, les méthodes de *Sequential Monte Carlo ABC* (SMC-ABC) consistent en un raffinement successif de la densité postérieure. À chaque étape, des paramètres θ sont évalués selon la méthode ABC pour une tolérance ϵ donnée. Ensuite, cette tolérance est diminuée et une nouvelle population de paramètres est sélectionnée par échantillonnage préférentiel dans la distribution précédente (Sisson et al., 2007; Del Moral et al., 2012).

Vraisemblance synthétique Les méthodes de vraisemblance synthétique, utilisées dans un contexte bayésien ou non (*Bayesian Synthetic Likelihood*, (B)SL) sont une autre approche lorsque la vraisemblance n'est pas accessible. Elles remplacent le calcul

de la vraisemblance exacte par une vraisemblance des statistiques résumées approchée par une loi normale selon la procédure suivante (Wood, 2010; Fasiolo et al., 2016; Price et al., 2018; An et al., 2020) :

- Des simulations \mathbf{x}_θ sont réalisées selon le paramètre θ évalué.
- Leurs statistiques résumées $s(\mathbf{x}_\theta)$ sont calculées.
- On suppose que les statistiques résumées suivent une loi normale multivariée : $s(\mathbf{x}_\theta) \sim \mathcal{N}(\mu_\theta, \Sigma_\theta)$. Les paramètres μ_θ et Σ_θ sont estimés à partir des simulations.
- La vraisemblance approchée des statistiques résumées des données $p(s(\mathbf{x}); \theta)$ est calculée à partir de la loi normale multivariée.

Cette vraisemblance peut ensuite être maximisée pour obtenir un estimateur de θ .

0.8.2.7 Autres approches

D'autres approches existent notamment dans le cas de modèles de Markov cachés. Fasiolo et al. (2016) en présente plusieurs reposant sur les filtres de particules (Gordon et al., 1993; Doucet and Johansen, 2008). Ceux-ci consistent à estimer la vraisemblance en simulant des trajectoires du modèle et en ré-échantillonnant parmi elles selon une pondération dépendant de la vraisemblance des observations sachant les états cachés. On trouve les méthodes MCMC particulières (PMCMC) (Andrieu et al., 2010), les filtres itérés (Ionides et al., 2011), l'estimation des paramètres en cascade (*Parameter Cascading*) (Ramsay et al., 2007)...

Il existe aussi des méthodes qui ne reposent pas sur des simulations mais effectuent un calcul approché de la vraisemblance (sous forme d'intégrale), par exemple grâce à des méthodes asymptotiques ou des quadratures multiples (Evans and Swartz, 1995).

Pour conclure, j'évoque les flux normalisants (*normalizing flows*). Ce sont des procédés de transformation de lois de probabilité simples, typiquement des lois normales, de sorte à reproduire des distributions de plus grande complexité (Papamakarios et al., 2019; Kobyzev et al., 2021). Cela permet d'approcher des distributions complexes avec une procédure simple à la fois pour effectuer des simulations et estimer la densité. Les flux normalisants sont utilisés en inférence variationnelle qui s'intéresse à l'estimation de variables cachées d'un modèle (Rezende and Mohamed, 2015) plutôt qu'aux paramètres.

0.9 Plan de la thèse

J'ai présenté dans cette introduction un état de l'art des thématiques principales traitées dans ma thèse. Les limites et perspectives soulevées ont motivé mes travaux de recherche des trois années écoulées. Mon manuscrit de thèse est composé de trois chapitres.

Le premier chapitre, relativement indépendant des deux suivants, présente un algorithme d'inférence pour modèles stochastiques que j'ai développé, appelé *Fixed Landscape Inference MethOd (flimo)*. Celui-ci a vu le jour face à des difficultés techniques pour analyser des données de métabarcoding. Mais son potentiel m'a conduit à lui consacrer une partie de mes travaux pour l'étudier plus en profondeur et l'appliquer à des problématiques écologiques variées. Son objectif premier est de répondre à la question : Comment inférer efficacement les paramètres d'un modèle aléatoire dont la vraisemblance n'est pas accessible mais dont il est facile de produire des simulations? L'idée sous-jacente était d'améliorer les méthodes disponibles, notamment l'approche des ABC. La solution technique à ce problème repose sur une gestion spécifique de l'aléa des simulations qui permet d'accélérer le processus d'inférence. Sur les applications étudiées, les résultats sont aussi précis que ceux obtenus par d'autres méthodes mais sont établis plus rapidement.

Le deuxième chapitre de ma thèse traite du métabarcoding quantitatif en répondant à plusieurs questions : 1) Comment mesurer le biais d'élongation de l'amplification PCR affectant les données de métabarcoding? 2) Comment mesurer le biais de concentration d'ADN cible? 3) Quel protocole expérimental et informatique permet de corriger efficacement ces biais? J'ai utilisé plusieurs techniques de PCR quantitative pour quantifier les biais observés. Ceux-ci ont été corrigés avec succès pour des communautés artificielles de plantes alpines. La correction repose sur un protocole informatique qui tire profit d'un modèle de PCR plus élaboré que le modèle exponentiel utilisé habituellement en métabarcoding.

Le troisième chapitre aborde le lien entre métabarcoding et mesures de biodiversité : 1) Comment attribuer les nombreuses séquences observées aux différentes unités taxonomiques? 2) Comment adapter les indices de biodiversité à l'incertitude de cette attribution? À partir d'un graphe de mutation établi par un algorithme existant, *obiclean*, j'ai évalué la vraisemblance des mutations, ponctuelles ou non, pour attribuer les séquences à différentes unités taxonomiques avec une certaine probabilité. Ces probabilités ont ensuite été utilisées pour établir des indices de biodiversité. Une étude des cooccurrences entre les variants souches et leurs mutants apparaît comme une piste très prometteuse pour caractériser les mutations de manière non paramétrique.

Enfin, la discussion générale dresse un bilan de mes contributions à la recherche et ouvre des perspectives sur de nouveaux projets ou des questions qui ont émergé de mes travaux.

Chapitre 1

Fixed Landscape Inference MethOd (flimo)

1.1 Introduction

1.1.1 Contexte de développement de *flimo*

Ce chapitre présente un nouvel algorithme d'inférence de paramètres pour modèles aléatoires appelé *Fixed Landscape Inference MethOd (flimo)*, ou "Méthode d'inférence par fixation du paysage [aléatoire]". Ces travaux relèvent des statistiques plus que de l'écologie. Mais la question de l'inférence de paramètres "décrivant au mieux des données observées" revient fréquemment en écologie ou biologie, où les phénomènes étudiés sont souvent décrits par des modèles probabilistes complexes.

Le développement de *flimo* n'était pas prévu dans mon projet de thèse. J'ai d'abord développé cet algorithme pour faciliter l'étude d'un modèle de PCR que nous voulions "inverser" pour estimer les quantités initiales des espèces présentes à partir de leur abondance finale (Chapitre 2). La vraisemblance de ce modèle n'était pas calculable. J'ai donc développé une méthode reposant sur des simulations en recourant à des statistiques résumées.

De manière générale, la difficulté de l'inférence provient du caractère aléatoire du processus étudié qui a deux sources : d'une part, les données sont aléatoirement bruitées ; d'autre part, chaque ensemble de simulations pour évaluer un paramètre θ est aléatoire. Cette seconde source d'aléa peut être réduite en agrégeant un grand nombre de simulations, mais l'évaluation de deux paramètres proches est toujours discontinuée du fait que des simulations distinctes sont utilisées pour l'un et l'autre.

Avec mes encadrants et Edouard Oudet (Laboratoire Jean Kuntzmann, Grenoble), nous avons réfléchi à un moyen de supprimer l'aléa des simulations à la racine, en modifiant la manière dont les tirages aléatoires sont réalisés numériquement. C'est l'idée au cœur de la méthode *flimo* : en rendant les simulations déterministes, il devient possible d'évaluer les paramètres de manière plus régulière et donc d'accélérer

l'exploration de l'espace des paramètres, que l'on désigne comme un paysage aléatoire "figé" par *flimo*. Cet algorithme, conçu pour étudier le modèle de PCR au départ, s'est montré particulièrement performant et nous avons décidé de l'étendre à un cadre plus large. La partie principale de ce chapitre est le manuscrit que nous avons écrit à son sujet.

J'introduis d'abord des concepts mathématiques utiles à la compréhension de *flimo*. En premier lieu, j'explique comment les simulations aléatoires sont réalisées en informatique (section 1.1.2). Ensuite, j'évoque les couplages de variables aléatoires (section 1.1.3) dont les simulations produites par *flimo* sont un cas particulier. Enfin, je dresse un rapide état de l'art des méthodes d'optimisation déterministes (section 1.1.4) qui sont utilisées par *flimo* pour inférer les paramètres du modèle étudié.

J'inclus ensuite le manuscrit de l'article présentant la méthode *flimo*.

1.1.2 Simulation de variables aléatoires

En introduction générale, j'ai présenté plusieurs méthodes d'inférence reposant sur des simulations aléatoires. Une question importante a été tue : comment simuler numériquement une variable aléatoire ? Dans la vie quotidienne, le hasard est facile à simuler, en lançant un dé par exemple. De nombreux phénomènes physiques sont également imprévisibles par essence. Mais en informatique, il n'est a priori pas possible de générer un nombre réellement aléatoire : les instructions transmises à la machine suivent une procédure déterministe. À la place, on crée des suites de nombres "pseudo-aléatoires". Leur tirage est déterministe mais elles vérifient des propriétés attendues de suites de variables aléatoires (indépendance, distribution...). Cette section présente les méthodes les plus communes, en particulier la méthode de la transformée inverse utilisée par *flimo*.

1.1.2.1 Génération de nombres pseudo-aléatoires selon une loi uniforme

Une idée récurrente est de simuler des variables aléatoires élémentaires selon la loi uniforme $\mathcal{U}([0, 1])$ puis de les transformer pour représenter n'importe quelle distribution. La méthode historique de simulation, appelée générateur congruentiel linéaire, est illustrée sur la Figure 1.1.1. Une suite de nombres (x_n) est simulée puis convertie en une suite (u_n) qui correspond aux tirages uniformes (équation 1.1).

$$\begin{aligned} x_0 &\in \mathbb{N} \\ x_{n+1} &\equiv ax_n + b \text{ modulo } T \\ \text{et ainsi } u_n &= \frac{x_n}{T} \overset{\text{approx.}}{\sim} \mathcal{U}([0, 1]) \end{aligned} \tag{1.1}$$

où a est un "grand" entier (historiquement $7^5 = 16807$), b est un entier quelconque et T est la période du générateur (historiquement $2^{31} - 1$), c'est-à-dire le nombre d'itérations avant de régénérer la même suite de nombres.

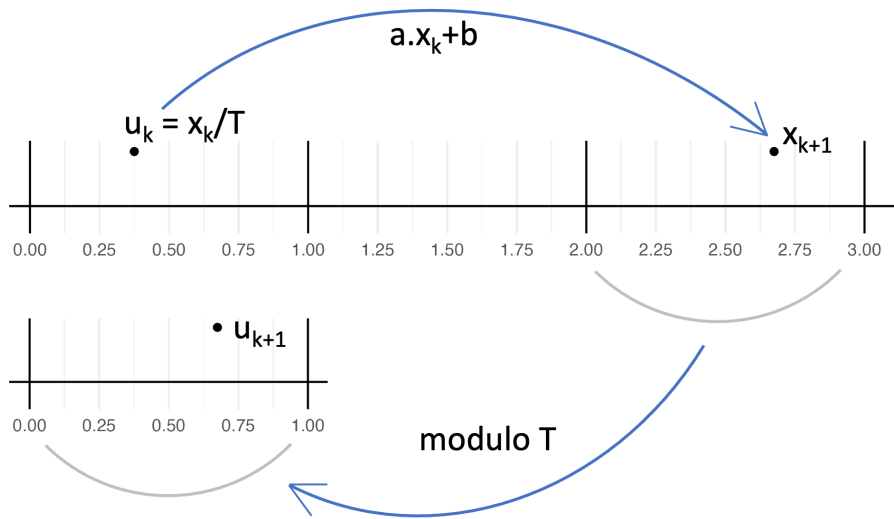


FIGURE 1.1.1 – Illustration du fonctionnement des générateurs congruents linéaires.

La méthode utilisée aujourd’hui est celle de Mersenne Twister (Matsumoto and Nishimura, 1998) pour générer des nombres flottants sur 32 ou 64 bits. Sa période est de $2^{19937} - 1 \simeq 4.3 \times 10^{6001}$.

La graine aléatoire x_0 (*seed*) est déterminée par l’horloge de l’ordinateur pour changer à chaque utilisation. Il est souvent utile de fixer la graine à une valeur arbitraire pendant la phase de développement d’un algorithme.

Une fois que l’on dispose d’un générateur pseudo-aléatoire, on peut générer des variables aléatoires quelconques par transformation.

1.1.2.2 Méthode de la transformée inverse

La méthode la transformée inverse est fondamentale dans le fonctionnement de *flimo*. Elle permet de simuler une variable aléatoire X en inversant sa fonction de répartition $F : x \mapsto \mathbb{P}(X \leq x)$. Pour cela, on utilise la fonction quantile :

$$F^{-1} : u \in]0, 1[\mapsto \inf \{x \mid F(x) \geq u\} \quad (1.2)$$

qui est la fonction inverse généralisée à gauche de F . Elle agit comme la fonction réciproque de F car $Q(F(X)) = X$ presque sûrement. On montre facilement que

$$\text{pour } U \sim \mathcal{U}([0, 1]), F^{-1}(U) \stackrel{\text{loi}}{=} X \quad (1.3)$$

c’est-à-dire que la variable aléatoire $F^{-1}(U)$ suit la loi que l’on désire simuler. La Figure 1.1.2 illustre cette méthode pour une loi normale centrée réduite, de fonction quantile $Q(r) = \sqrt{2} \operatorname{erf}^{-1}(2r - 1)$, $r \in]0, 1[$ avec erf la fonction d’erreur usuelle. Cette méthode fonctionne pour tout type de variable à partir du moment où l’on peut inverser F analytiquement ou numériquement.

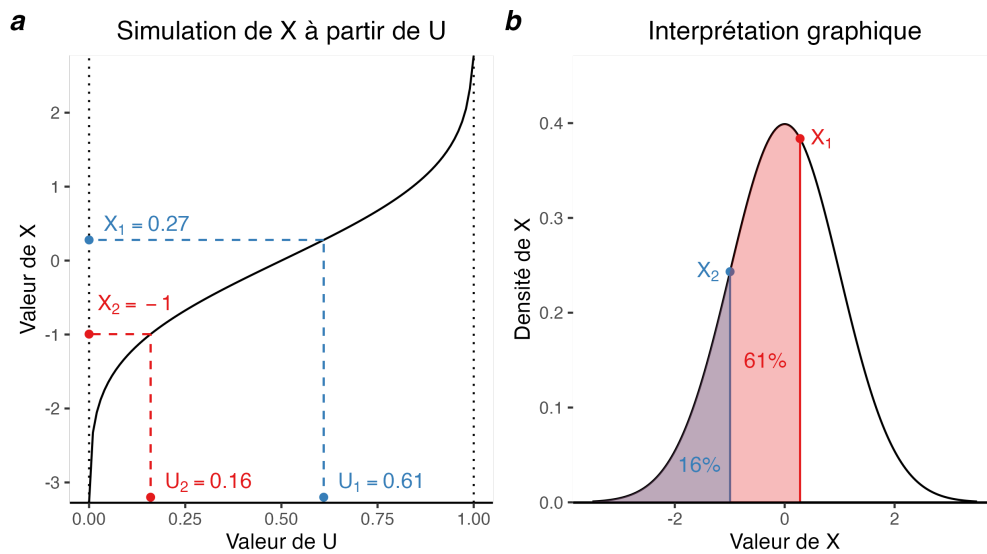


FIGURE 1.1.2 – Méthode de la transformée inverse pour simuler deux variables selon une loi normale centrée réduite $\mathcal{N}(0, 1)$. Deux variables uniformes sont simulées : on obtient $U_1 = 0.61$ et $U_2 = 0.16$. Le panneau **a** montre la détermination de X_1 et X_2 à partir de la fonction quantile de la loi $\mathcal{N}(0, 1)$. Le panneau **b** montre la densité de cette loi. Le lien entre U_1 et X_1 (resp. U_2 et X_2) s'interprète graphiquement : X_1 est la valeur de la loi telle que $U_1 = 61\%$ de la masse de la densité est située avant X_1 (resp. $U_2 = 16\%$ de la masse avant X_2).

1.1.2.3 Simulation de variables gaussiennes : transformation de Box-Muller

Dans le cas de la loi normale, la méthode de la transformée inverse implique une approximation de la fonction F^{-1} . La transformation de Box-Muller est une méthode connue pour simuler deux variables iid dans $\mathcal{N}(0, 1)$ sans approximation (Box and Muller, 1958) en utilisant une transformation polaire de la densité de la loi normale.

Soit $R \geq 0$ tel que $R^2 \sim \text{Exp}(\frac{1}{2})$ et $\theta \sim \mathcal{U}([0, 2\pi])$. Ces variables sont obtenues par transformée inverse :

$$\begin{aligned} U, V &\sim \mathcal{U}([0, 1]) \\ R^2 &= -2 \log U \text{ et } \theta = 2\pi V \end{aligned} \quad (1.4)$$

Les variables

$$X = R \cos(\theta) \text{ et } Y = R \sin(\theta) \quad (1.5)$$

suivent alors la loi $\mathcal{N}(0, 1)$ et sont indépendantes. Cette propriété est démontrée en utilisant un théorème d'identification de la densité¹.

1. Présenté ici : https://perso.univ-rennes2.fr/system/files/users/fromont_m/ProbasE nsai_Magalie.pdf, Théorème 5.

1.1.2.4 Méthode de rejet

La méthode de rejet est une alternative à la transformée inverse. On suppose que l'on connaît la densité f de la loi désirée mais qu'il n'est pas évident de simuler directement une variable aléatoire X suivant cette loi. Pour simuler X , on simule une variable aléatoire Y de densité g que l'on accepte comme une réalisation de X avec une probabilité proportionnelle à $\frac{f(Y)}{g(Y)}$ (Forsythe, 1972).

Soit $U \sim \mathcal{U}([0, 1])$ et $c \geq 1$ tel que $c g(x) \geq f(x)$ presque partout. Alors

$$Y|(U c g(Y) \leq f(Y)) \stackrel{\text{loi}}{=} X \quad (1.6)$$

On simule donc un couple de variables indépendantes (Y, U) jusqu'à ce que la condition $U c g(Y) \leq f(Y)$ soit vérifiée. La valeur de Y est alors une réalisation de même loi que X . Cette méthode est illustrée sur la Figure 1.1.3.

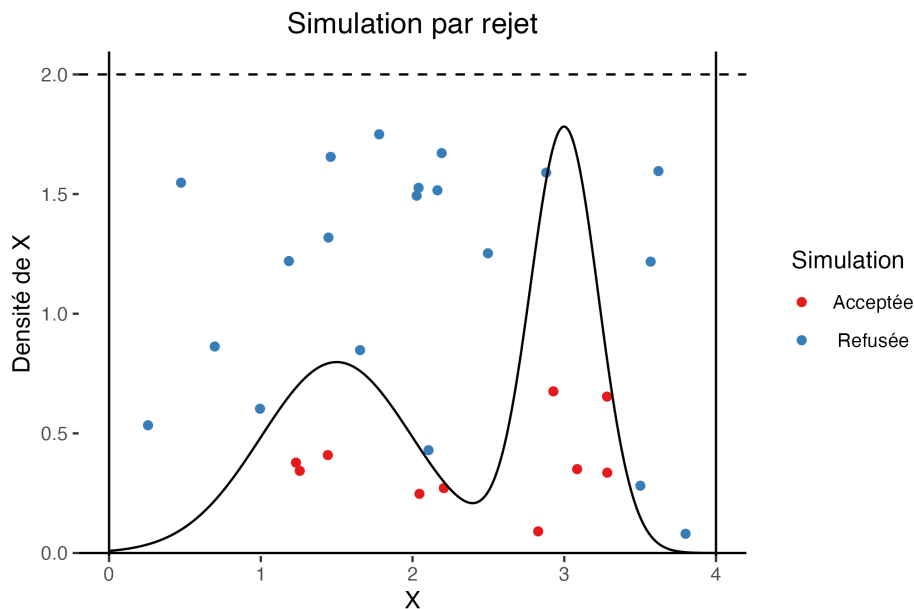


FIGURE 1.1.3 – Méthode du rejet pour une densité quelconque f (courbe continue). La densité g choisie est celle d'une loi uniforme entre 0 et 4 (ligne pointillée) de densité $g(y) = \frac{1}{4}$ sur $[0, 4]$. Ici, $c = 8$ de sorte que $c g(y) > \max_{0 \leq y \leq 4} f(y)$. Dans cet exemple, il a fallu 36 simulations pour obtenir 10 valeurs de la loi de X .

1.1.2.5 Méthodes approchées

Les algorithmes précédents permettent de simuler des nombres qui suivent théoriquement la loi exacte désirée. D'autres méthodes permettent de simuler des distributions approchées dans le cas où la distribution exacte n'est pas accessible pour diverses raisons. C'est ce que font les méthodes d'échantillonnage préférentiel et MCMC présentées en 0.8.2. Les flux normalisants peuvent aussi servir à des fins de simulation (Papamakarios et al., 2019).

1.1.3 Couplage de variables aléatoires

En théorie des probabilités, le couplage est une technique de comparaison de deux distributions en considérant une distribution jointe dont chaque marginale est une des distributions en question. J'en fais une brève présentation car *flimo* repose sur un couplage dit monotone établi entre les simulations numériques.

Soit μ et ν deux distributions définies sur un ensemble mesurable (Ω, \mathcal{A}) . Par définition, un couplage de μ et ν est une distribution γ sur $(\Omega \times \Omega, \mathcal{A} \times \mathcal{A})$ telle que les distributions marginales de γ sont μ et ν .

Le couplage monotone est un cas particulier lié à la notion de dominance stochastique. Pour deux mesures de probabilité μ et ν sur \mathbb{R} , on dit que μ domine stochastiquement ν si pour tout $x \in \mathbb{R}$, $\mu([x, +\infty[) \geq \nu([x, +\infty[)$. En notant F_μ et F_ν les fonctions de répartition de μ et ν , cela est équivalent à $F_\mu(x) \geq F_\nu(x)$ pour tout $x \in \mathbb{R}$.

Pour deux variables aléatoires réelles X et Y , un théorème assure que : X domine stochastiquement Y si et seulement s'il existe un couplage (\hat{X}, \hat{Y}) de X et Y tel que $\mathbb{P}(\hat{X} \geq \hat{Y}) = 1$. (\hat{X}, \hat{Y}) est alors appelé couplage monotone de X et Y .

De plus, pour toute fonction f croissante, $(f(\hat{X}), f(\hat{Y}))$ est un couplage monotone de $f(X)$ et $f(Y)$ et $\mathbb{E}[f(X)] \geq \mathbb{E}[f(Y)]$, si les espérances existent. La preuve de ces assertions ainsi qu'une étude bien plus vaste des couplages sont disponibles dans *Modern Discrete Probability : An Essential Toolkit*, Chapitre 4, de Sébastien Roch (en préparation de publication par *Cambridge University Press*²).

La manière classique de construire un couplage monotone est d'utiliser la méthode de la transformée inverse (section 1.1.2) en utilisant les fonctions quantiles des lois de X et Y de fonctions de répartition F et G et une variable aléatoire uniforme U :

$$U \sim \mathcal{U}([0, 1]) \text{ et } (\hat{X}, \hat{Y}) = (F^{-1}(U), G^{-1}(U)) \quad (1.7)$$

L'algorithme *flimo* compare des jeux de paramètres $\theta_1, \dots, \theta_n$ en produisant des simulations couplées $(X_{\theta_i})_{1 \leq i \leq n}$. Dans le cas où θ est un paramètre de position, par exemple l'espérance de la loi, le couplage est monotone :

$$\begin{aligned} &\text{si } \theta_1 \leq \dots \leq \theta_n, \\ &\text{alors } X_{\theta_1} \leq \dots \leq X_{\theta_n} \end{aligned} \quad (1.8)$$

Pour d'autres paramètres, comme l'écart-type d'une loi normale, il n'y a pas de résultat de monotonie.

2. <https://people.math.wisc.edu/~roch/mdp/>

La Figure 1.1.4 illustre le processus de simulation par couplage pour une loi de Poisson.

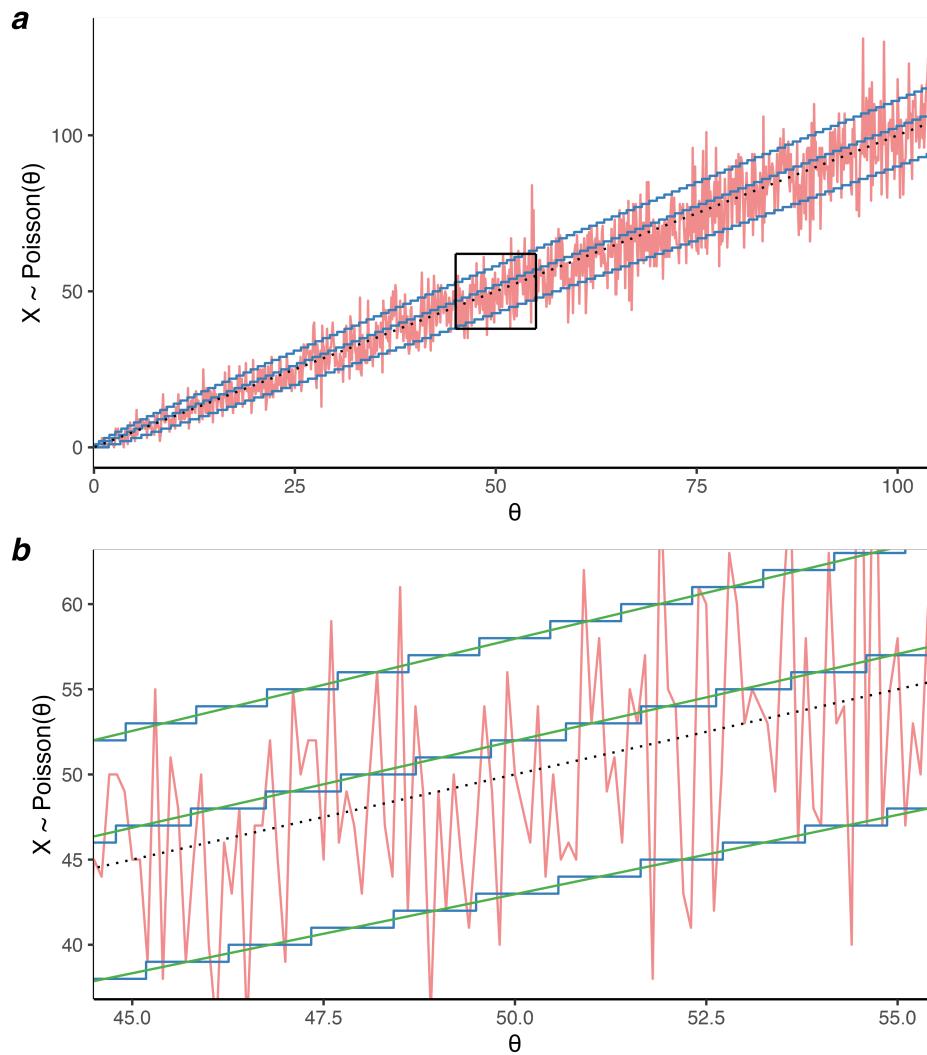


FIGURE 1.1.4 – Simulations selon une loi de Poisson. Le panneau **a** montre une simulation indépendante pour chaque valeur de θ (en rouge) ou des simulations couplées (en bleu) selon l'équation 1.7 pour $U \in \{0.16, 0.61, 0.87\}$, de bas en haut. La ligne pointillée est l'espérance du tirage. Le panneau **b** est un grossissement du panneau **a** (selon le rectangle noir). Ce panneau illustre le fait que la loi de Poisson est discrète : sa fonction quantile est constante par morceaux. Les courbes vertes sont des simulations couplées de la loi normale approchant la loi de Poisson, $\mathcal{N}(\theta, \theta)$. Celle-ci présente l'avantage d'être continue.

1.1.4 Algorithmes d'optimisation déterministe

L'optimisation est une problématique récurrente qui touche tous les domaines liés de près ou de loin à la modélisation. Une fois qu'un processus a été observé et modélisé, il est naturel de chercher à l'optimiser pour l'ajuster au mieux au phénomène observé.

Les méthodes d'inférence présentées en 0.8.2 ont aussi pour but d'optimiser des modèles aléatoires, mais j'aborde ici le cas où la fonction à optimiser est déterministe. Les outils utilisés sont assez différents malgré des problématiques similaires. Ce sujet est important pour *flimo*. En effet, même si elle traite des modèles stochastiques, la gestion spécifique des simulations conduit à employer un algorithme d'optimisation déterministe.

Cette section présente les processus d'optimisation communs, leur implémentation et le calcul des différentielles. Je ne développe ni un cadre théorique complet ni une revue exhaustive et technique de ces algorithmes mais je m'efforce de mettre en lumière les enjeux importants pour construire un problème d'optimisation. Ma référence principale est le cours *Approximation numérique et optimisation* de Grégoire Allaire (École polytechnique, édition 2017). Nocedal and Wright (2006) est un ouvrage de référence complet.

Je ne parle ici que d'algorithmes déterministes d'optimisation locale. Toute une gamme d'algorithmes a pour objectif une optimisation globale, souvent sous forme stochastique : recuit simulé, algorithmes génétiques, optimisation bayésienne (Shahriari et al., 2016)... Ces algorithmes sont moins efficaces mais nécessitent moins d'hypothèses sur la fonction à optimiser. J'évoquerai leurs avantages en Discussion générale.

1.1.4.1 Cadre mathématique

Plusieurs éléments sont nécessaires pour définir un problème d'optimisation. Dans ce texte, l'espace de définition des paramètres du modèle est $\Theta \subset \mathbb{R}^n$, $n \geq 1$ quelconque. Θ est appelé ensemble admissible. Cet espace est de dimension finie mais le nombre de paramètres n'est pas limité.

On considère une fonction objectif $J : \mathbb{R}^n \rightarrow \mathbb{R}$ continue dont on cherche le minimum (ou le maximum) sur Θ . Le problème s'écrit :

$$\inf_{\theta \in \Theta} f(\theta) \tag{1.9}$$

On écrit inf et non min car il n'y a pas de garantie que ce minimum soit atteint (i.e., qu'il existe $\theta^* \in \Theta$ tel que $f(\theta^*) = \min_{\theta \in \Theta} f(\theta)$).

1.1.4.2 Enjeux

Dans le cas général, il faut d'abord montrer l'existence de minima locaux (des résultats généraux sont établis en dimension finie), puis en donner des conditions nécessaires

(mais pas suffisantes) d'optimalité. La convexité est une propriété cruciale pour les problèmes d'optimisation car elle fournit une caractérisation des minima nettement plus simple. Par ailleurs, on cherche généralement à optimiser des fonctions différentiables pour profiter des outils efficaces dont on dispose dans ce cadre.

1.1.4.3 Cas des fonctions différentiables

Lorsque les fonctions sont régulières, on peut extraire des informations sur leurs optima à partir de leurs différentielles (ou dérivées) d'ordre 1 et 2. Dans la suite, on considère que la fonction J est aussi régulière que nécessaire pour l'ensemble des méthodes présentées³.

Je traite plus bas du cas des fonctions non différentiables. Pour celles-ci, les performances algorithmiques sont nettement moins bonnes et il n'est pas facile de donner des conditions d'optimalité.

Gradient La différentielle d'ordre 1 est caractérisée en chaque point $\theta \in \Theta$ par un élément $\nabla J(\theta) \in \mathbb{R}^n$ appelé gradient qui vérifie, dans un voisinage de θ ,

$$J(\theta + h) = J(\theta) + \langle \nabla J(\theta), h \rangle + o(h) \quad \text{avec} \quad \lim_{h \rightarrow 0} \frac{o(h)}{\|h\|} = 0 \quad (1.10)$$

Son existence et son unicité sont garanties par le théorème de représentation de Riesz.

Hessienne La différentielle d'ordre 2 est caractérisée par une matrice $H_J(\theta) \in \mathbb{R}^{n \times n}$, qui vérifie :

$$J(\theta + h) = J(\theta) + \langle \nabla J(\theta), h \rangle + \langle H_J(\theta)h, h \rangle + o(\|h\|^2) \quad \text{avec} \quad \lim_{h \rightarrow 0} \frac{o(\|h\|^2)}{\|h\|^2} = 0 \quad (1.11)$$

A certaines conditions exposées plus loin, la hessienne permet de distinguer les minima locaux des maxima locaux.

1.1.4.4 Analyse convexe

La convexité est une propriété importante pour caractériser les minima d'une fonction. Les problèmes convexes, impliquant des fonctions convexes sur des ensembles convexes, sont donc une catégorie particulièrement favorable à étudier. Voici quelques définitions.

3. La plus exigeante est la méthode de Newton qui impose que J soit de classe C^3 sur \mathbb{R}^n , c'est-à-dire différentiable (au sens de Fréchet) continûment trois fois.

Ensemble convexe Un ensemble K est dit convexe si, pour tous $x, y \in K$, le segment $[x, y]$ est entièrement contenu dans K , c'est-à-dire :

$$\forall x, y \in K, \forall t \in [0, 1], \quad tx + (1 - t)y \in K \quad (1.12)$$

La Figure 1.1.5 illustre cette définition.

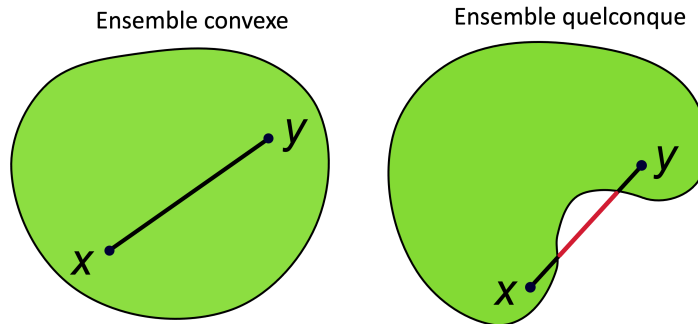


FIGURE 1.1.5 – Ensemble convexe (à gauche) et non convexe (à droite) dans \mathbb{R}^2 . Source de la figure : https://fr.wikipedia.org/wiki/Ensemble_convexe.

Fonction convexe Une fonction J est dite convexe si son graphe est en dessous de toutes ses cordes (Figure 1.1.6), c'est-à-dire :

$$\forall x, y \in I, \forall t \in [0, 1], \quad f(tx + (1 - t)y) \leq tf(x) + (1 - t)f(y) \quad (1.13)$$

Pour une fonction convexe, chaque minimum local est un minimum global, et si l'inégalité est stricte dans l'équation 1.13 (pour $x \neq y$, $0 < t < 1$), alors il existe au plus un minimum (global).

La convexité est aussi caractérisée par la dérivée seconde : J est convexe si et seulement si $\langle H_J(\theta)h, h \rangle \geq 0$, $\forall \theta, h \in \mathbb{R}^n$, ce qui traduit le fait que la pente augmente.

1.1.4.5 Conditions d'optimalité

Dans ce texte, je ne présente que le cas où l'ensemble admissible Θ est convexe et même plus précisément le cas où chaque paramètre est compris sur un intervalle (potentiellement infini) $[a, b] \subset \mathbb{R}$. Je n'aborde pas ici les cas où les paramètres optimaux sont des bornes de ces intervalles.

Pour d'autres formes de contraintes sur les paramètres impliquant des égalités et des inégalités ($F(\theta) \leq 0$ pour une fonction F quelconque), la résolution passe par des Lagrangiens. C'est un sujet très intéressant mais qui n'a pas sa place ici.

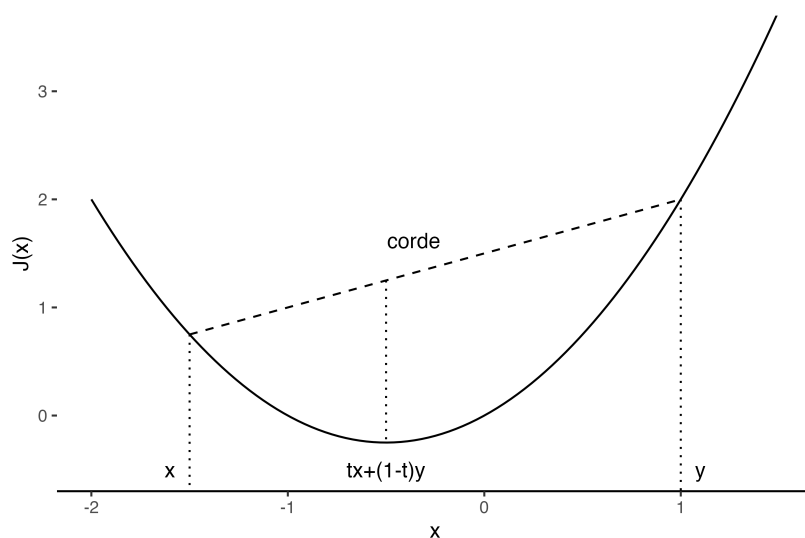


FIGURE 1.1.6 – Graphe d’une fonction convexe : la corde est au-dessus du graphe de la fonction. L’épigraphe, c’est-à-dire la partie au-dessus du graphe, est un ensemble convexe, d’où la dénomination commune.

Condition d’optimalité du premier ordre Supposons donc que Θ est un ensemble convexe. L’inéquation d’Euler (1.14) donne une condition nécessaire d’optimalité. Si θ^* intérieur à Θ est un minimum local de J ,

$$\langle \nabla J(\theta^*), \theta - \theta^* \rangle \geq 0 \quad , \quad \forall \theta \in \Theta \quad (1.14)$$

Réciproquement, si J est convexe et si θ^* vérifie 1.14, alors u est un minimum global de J sur Θ . Pour les points intérieurs de Θ , on retrouve la condition nécessaire : $\nabla J(\theta^*) = 0$. Elle n’est pas suffisante, comme l’illustre la Figure 1.1.7.

Condition d’optimalité du second ordre Les différentielles d’ordre 2 fournissent une autre condition d’optimalité. Ici, on suppose $\Theta = \mathbb{R}^n$. Si θ^* est un minimum local de J , alors

$$\nabla J(\theta^*) = 0 \text{ et } \langle H_J(\theta^*)\theta, \theta \rangle \geq 0 \quad , \quad \forall \theta \in \mathbb{R}^n \quad (1.15)$$

Réciproquement, si pour tout θ dans un voisinage de θ^* ,

$$\nabla J(\theta^*) = 0 \text{ et } \langle H_J(\theta)x, x \rangle \geq 0 \quad , \quad \forall x \in \mathbb{R}^n \quad (1.16)$$

alors θ^* est un minimum local de J . La fonction $J : x \mapsto x^3$ en $x = 0$ illustre que l’équation 1.15 est une condition nécessaire mais pas suffisante : J a un point d’inflexion en 0 qui n’est ni un minimum ni un maximum sur \mathbb{R} .

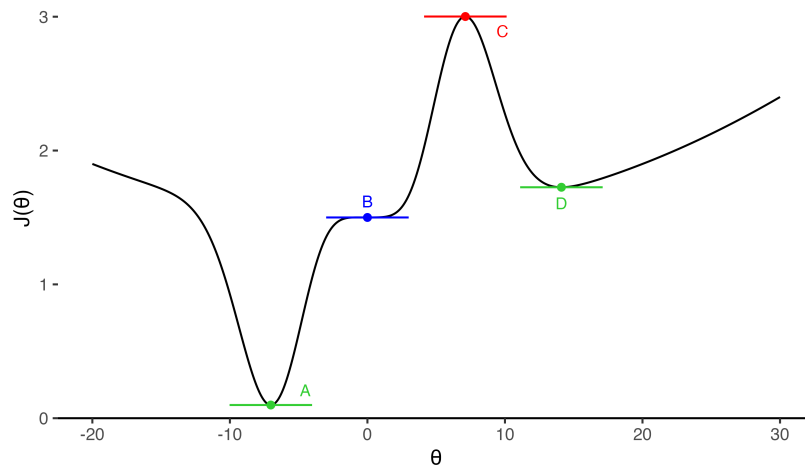


FIGURE 1.1.7 – Graphe d’une fonction avec quatre points de dérivée nulle (tangente à la courbe plate). A est le minimum global, B un point d’inflexion, C un maximum local, D un minimum local. Selon l’algorithme retenu et la condition initiale, l’optimisation pourrait converger vers l’un de ces quatre points.

C’est pour vérifier l’une ou l’autre des conditions 1.14 et 1.15 que l’on cherche à trouver les solutions du problème $\nabla J(\theta) = 0$ (souvent noté $J'(\theta) = 0$ en une dimension) sur l’ensemble admissible pour trouver des solutions candidates.

1.1.4.6 Quelques algorithmes d’optimisation

Je présente maintenant différents algorithmes itératifs qui utilisent ces propriétés pour approcher numériquement le minimum d’une fonction. Le fonctionnement général est le suivant : on fixe une condition initiale $\theta_0 \in \Theta$ puis on construit une suite $(\theta_n)_n$ de Θ convergeant vers le minimum de la fonction J dans Θ :

$$\theta_{n+1} = \theta_n - \mu_n w_n \quad (1.17)$$

où $w_n \in \Theta$ est un vecteur unitaire donnant la direction d’évolution et $\mu_n \in \mathbb{R}$ est le pas effectué dans cette direction. Si le pas est trop grand, l’algorithme risque de diverger, s’il est trop petit, l’algorithme est lent.

Je ne donne pas de détail sur les preuves de convergence qui impliquent une propriété de convexité plus forte. Si celle-ci n’est pas vérifiée, il existe un risque que la suite oscille entre des minima locaux, par exemple dans une "vallée" optimale. De même, il est difficile de donner une idée des vitesses de convergence réelles car celles-ci sont très dépendantes du problème à traiter. Je me contenterai d’une description qualitative.

Le critère de convergence dépend généralement de la différence entre θ_{n+1} et θ_n ou entre $J(\theta_{n+1})$ et $J(\theta_n)$, ou encore de la valeur de $\|\nabla J(\theta_n)\|$.

Un mot sur l'implémentation Les packages pour la méthode *flimo* ont été développés dans les langages R (largement diffusé chez les écologues) et Julia (pour ses performances numériques). J'utilise des algorithmes d'optimisation déjà implémentés. En R, j'utilise des méthodes disponibles dans la fonction `base::optim`. En Julia, j'utilise le package `Optim.jl`. Je citerai au fil du texte les méthodes retenues. Ce ne sont pas les mêmes dans les deux langages, du fait des implémentations disponibles.

1.1.4.7 Méthodes d'optimisation différentielles

Descente de gradient La première catégorie utilise le gradient pour donner la direction d'évolution : $w_n = \frac{\nabla J(\theta_n)}{\|\nabla J(\theta_n)\|}$: la suite des θ_n évolue selon la plus forte pente. Il existe ensuite différentes possibilités pour choisir le pas μ_n , généralement en résolvant le problème d'optimisation en une dimension : $\inf_{\mu \in \mathbb{R}} J(\theta_n - \mu w_n)$ ("recherche linéaire"). Lorsque l'ensemble admissible Θ est formé d'intervalles, la contrainte $\theta_{n+1} \in \Theta$ est facile à satisfaire par projection sur Θ .

Ces méthodes ont une vitesse de convergence linéaire. La méthode du gradient conjugué permet d'améliorer la vitesse de convergence. L'algorithme du gradient stochastique est une alternative dans le cas où le calcul de ∇J est compliqué.

Méthode de (quasi-)Newton La deuxième catégorie utilise l'information fournie par la différentielle de second ordre. Elle repose sur la méthode de Newton. Pour une fonction $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$, celle-ci permet de trouver un point θ vérifiant $F(\theta) = 0$ (et de différentielle d'ordre 1 $F'(\theta)$ inversible). La relation de récurrence est donnée par :

$$\theta_{n+1} = \theta_n - (F'(\theta_n))^{-1} F(\theta_n) \quad (1.18)$$

Cette méthode a une vitesse de convergence quadratique (mais chaque itération est plus longue à calculer que pour la descente de gradient). Pour résoudre un problème de minimisation, on remplace F par ∇J , ce qui donne :

$$\theta_{n+1} = \theta_n - (H_J(\theta_n))^{-1} \nabla J(\theta_n) \quad (1.19)$$

On ne calcule pas l'inverse de la hessienne car cela est coûteux en temps, mais on résout à la place le système linéaire $H_J(\theta_n)X = \nabla J(\theta_n)$.

L'inconvénient de cette méthode est qu'elle nécessite un calcul de hessienne à chaque itération. Les méthodes de quasi-Newton contournent cette limite en approchant $H_J(\theta_n)$ par une matrice S_n symétrique définie positive. Un choix typique est fait par l'algorithme BFGS à partir de la méthode de la sécante dont l'idée est de trouver une matrice vérifiant :

$$S_n (\theta_n - \theta_{n-1}) = (\nabla J(\theta_n) - \nabla J(\theta_{n-1})) \quad (1.20)$$

Pour *flimo* implémenté en R, j'utilise une adaptation appelée L-BFGS-B (Byrd et al., 1995) qui requière moins de mémoire et qui permet d'imposer des contraintes sur l'ensemble admissible sous formes d'intervalles.

Méthode de point intérieur Pour *flimo* en Julia, j'utilise la méthode IPNewton (*Interior-Point Newton*) (Wächter and Biegler, 2006; Nocedal and Wright, 2006) qui est une méthode de quasi-Newton dite primale-duale de point intérieur, utilisant des pénalisations successives pour imposer les contraintes sur Θ . Cela consiste à modifier le problème d'optimisation en remplaçant l'objectif $J(\theta_n)$ par $J(\theta_n) + B(\theta_n)$ où B est une fonction barrière. La fonction barrière est diminuée au cours de l'optimisation qui se rapproche ainsi du problème exact.

1.1.4.8 Méthodes d'optimisation non différentielles

Les méthodes d'optimisation non différentielles sont connues pour converger moins vite et pour fournir des solutions variant (parfois fortement) selon les conditions initiales. En revanche, ils s'appliquent à des problèmes où la fonction objectif est complexe ou irrégulière. Ils fonctionnent par itération mais en conservant en mémoire plusieurs points qui définissent la région de Θ où chercher un minimum.

Voici deux méthodes que j'utilise en R et en Julia selon la dimension du problème.

Algorithme de Brent L'algorithme de Brent permet de minimiser une fonction à une dimension $J : \mathbb{R} \rightarrow \mathbb{R}$ sur un segment $[a, b]$. Elle consiste à tester plusieurs méthodes (dichotomie, méthode de la sécante, interpolation quadratique inverse) et à effectuer une itération avec celle qui donne le meilleur résultat pour actualiser l'intervalle de recherche (Brent, 2002).

Algorithme de Nelder-Mead L'algorithme de Nelder-Mead permet de minimiser une fonction non différentiable à plusieurs dimensions. Elle consiste à actualiser un simplexe ($n + 1$ points dans un espace à n dimensions, soit un segment en 1D ou un triangle en 2D) qui définit l'espace de recherche (Nelder and Mead, 1965; Gao and Han, 2012). À chaque itération, le point le moins bon est remplacé par un meilleur.

1.1.4.9 Différentiation automatique

Enfin, j'aborde un dernier point : le calcul des différentielles. Par défaut, les gradients (1.21) et les hessiennes (1.22) sont calculés de manière approchée par différences finies. En une dimension, la méthode la plus simple est, avec h "petit" :

$$f'(x) \simeq \frac{f(x+h) - f(x)}{h} \tag{1.21}$$

$$f''(x) \simeq \frac{f(x+h) - 2f(x) + f(x-h)}{h^2} \tag{1.22}$$

Ces calculs peuvent être longs, notamment quand la fonction J est complexe et que chaque appel est coûteux en temps. De plus, les approximations numériques atténuent la précision du calcul.

La différentiation automatique permet de réduire ces limites. Elle part du constat que les fonctions, même complexes, sont généralement constituées d'une suite d'opérations élémentaires (additions, multiplications...) et d'applications de fonctions simples (exponentielle, cosinus...) que l'on sait dériver. La différentiation automatique utilise cette décomposition pour calculer les différentielles de fonctions en appliquant la règle de la chaîne (*chain rule*) pour les fonctions composées $x \mapsto g(f(x))$:

$$\frac{d(g \circ f)}{dx} = \frac{dg}{df} \frac{df}{dx} \quad (1.23)$$

La différentiation automatique est utilisée pour *flimo* en Julia pour améliorer la vitesse d'inférence, grâce au package ForwardDiff.jl (Revels et al., 2016). Les fonctions quantiles de quasiment toutes les lois continues considérées ont pu être différenciées automatiquement, alors même que leur expression est non triviale.

Je reprends l'exemple proposé sur Wikipédia⁴ qui est une bonne illustration. On considère une fonction à deux variables f que l'on décompose successivement.

$$f(x, y) = xy + \sin(x) \quad (1.24)$$

$$= w_1 w_2 + \sin(w_1) \quad (1.25)$$

$$= w_3 + w_4 \quad (1.26)$$

$$= w_5 \quad (1.27)$$

Pour calculer la dérivée de f par rapport à x , on écrit

$$\frac{\partial w_1}{\partial x} = \frac{\partial x}{\partial x} = 1 \quad (1.28)$$

$$\frac{\partial w_2}{\partial x} = \frac{\partial y}{\partial x} = 0 \quad (1.29)$$

$$\frac{\partial w_3}{\partial x} = \frac{\partial(w_1 w_2)}{\partial x} = \frac{\partial w_1}{\partial x} w_2 + w_1 \frac{\partial w_2}{\partial x} \quad (1.30)$$

$$\frac{\partial w_4}{\partial x} = \frac{\partial \sin(w_1)}{\partial x} = \cos(w_1) \frac{\partial w_1}{\partial x} \quad (1.31)$$

$$\frac{\partial w_5}{\partial x} = \frac{\partial w_3 + w_4}{\partial x} = \frac{\partial w_3}{\partial x} + \frac{\partial w_4}{\partial x} \quad (1.32)$$

et donc en recombinaison :

4. https://fr.wikipedia.org/wiki/Diff%C3%A9rentiation_automatique

$$\frac{\partial f}{\partial x}(x, y) = \frac{\partial w_5}{\partial x} = \dots = y + \cos(x) \quad (1.33)$$

et de la même façon,
$$\frac{\partial f}{\partial y}(x, y) = x \quad (1.34)$$

Selon que l'on commence à calculer les premiers termes (w_1 puis $w_2\dots$) ou les derniers termes (w_5 puis $w_4\dots$), on parle d'accumulation directe ou inverse dont le choix dépend du contexte. Ce ne sont d'ailleurs pas les seules méthodes disponibles pour la différentiation automatique.

Après avoir traité ces différents points techniques, je présente le manuscrit décrivant *flimo* dans la version soumise auprès de PeerJ en septembre 2023. La section *Supplementary Information* mentionnée dans le texte est placée en Annexe A.

1.2 Résumé en langue française

La modélisation en biologie doit être adaptée à des données de plus en plus complexes et massives. L'efficacité des algorithmes d'inférence utilisés pour estimer les paramètres des modèles est donc un enjeu de développement majeur. Beaucoup d'entre eux sont basés sur des processus d'optimisation stochastiques qui nécessitent un temps de calcul important. Nous présentons ici l'algorithme *Fixed Landscape Inference Method (flimo)*, une nouvelle méthode d'inférence sans vraisemblance pour modèles aléatoires continus. Cette méthode utilise un algorithme d'optimisation déterministe de type quasi-Newton pour obtenir une estimation ponctuelle des paramètres du modèle. Pour cela, la différence de statistiques résumées est minimisée entre les données observées et des simulations réalisées avec une gestion particulière des tirages aléatoires. Cet usage de statistiques résumées rapproche *flimo* de l'algorithme *Approximate Bayesian Computation (ABC)*. Comme les ABC, *flimo* peut aussi estimer une distribution empirique des paramètres. Trois applications sont présentées ici : un exemple théorique habituel, l'inférence de paramètres d'une distribution g-et-k ; un problème de génétique des populations, l'inférence d'une valeur sélective à partir de séries temporelles dans une population suivant le modèle de Wright-Fisher ; et l'étude d'un modèle de Ricker, simulant des dynamiques de population chaotiques. Pour les deux premières applications, les résultats montrent une réduction drastique du temps d'inférence en comparaison d'autres méthodes usuelles, pour une précision similaire. Même quand la vraisemblance du modèle étudié est disponible, la simplicité et l'efficacité de *flimo* en font une alternative pertinente. Une implémentation en Julia et en R est disponible sur <https://git.metabarcoding.org/lecasofts/flimo>. Pour utiliser *flimo*, la seule condition est de pouvoir simuler le processus étudié.

1.3 Manuscrit

The Fixed Landscape Inference MethOd (*flimo*):
a versatile alternative to Approximate Bayesian
Computation, faster by several orders of
magnitude

Sylvain Moinard^{1*}, Edouard Oudet², Didier Piau³,
Eric Coissac^{1†}, Christelle Gonindard-Melodelima^{1†}

^{1*}Univ. Grenoble Alpes, Univ. Savoie Mont Blanc, CNRS, LECA,
FR-38000, Grenoble, France.

²Univ. Grenoble-Alpes, CNRS, Laboratoire Jean Kuntzmann (LJK),
FR-38000, Grenoble, France.

³Univ. Grenoble-Alpes, CNRS, Institut Fourier, FR-38000, Grenoble,
France.

*Corresponding author(s). E-mail(s):

sylvain.moinard@univ-grenoble-alpes.fr;

Contributing authors: edouard.oudet@univ-grenoble-alpes.fr;
didier.piau@univ-grenoble-alpes.fr; eric.coissac@metabarcoding.org;
christelle.gonindard@univ-grenoble-alpes.fr;

†These authors contributed equally to this work.

Abstract

Modelling in biology must adapt to increasingly complex and massive data. The efficiency of the inference algorithms used to estimate model parameters is therefore questioned. Many of these are based on stochastic optimization processes that require significant computing time. We introduce the Fixed Landscape Inference MethOd (*flimo*), a new likelihood-free inference method for continuous state-space stochastic models. It applies deterministic gradient-based optimization algorithms to obtain a point estimate of the parameters, minimizing the difference between the data and some simulations according to some prescribed summary statistics. In this sense, it is analogous to Approximate Bayesian Computation (ABC). Like ABC, it can also provide an approximation of the distribution of the parameters. Three applications are proposed: a usual theoretical example,

namely the inference of the parameters of g-and-k distributions; a population genetics problem, not so simple as it seems, namely the inference of a selective value from time series in a Wright-Fisher model; and simulations from a Ricker model, representing chaotic population dynamics. In the two first applications, the results show a drastic reduction of the computational time needed for the inference phase compared to the other methods, despite an equivalent accuracy. Even when likelihood-based methods are applicable, the simplicity and efficiency of *flimo* make it a compelling alternative. Implementations in Julia and in R are available on <https://metabarcoding.org/flimo>. To run *flimo*, the user must simply be able to simulate data according to the chosen model.

Keywords: Approximate Bayesian Computation, Likelihood-free inference, Model optimization, Simulation-based inference, Stochastic models

MSC Classification: 6204 , 62F10 , 62M09 , 9208

1 Introduction

Modelling in biology and ecology presents some important conceptual challenges, due to the increasing complexity and size of the available data. Even for the simplest models, the likelihood is often intractable, especially in population genetics problems (Stephens, 2004). Bayesian sampling methods (Shoemaker et al., 1999) are often used to study such models. In addition to the optimal solution, these methods provide the posterior distribution of the parameters. Among these, Markov chain Monte Carlo (MCMC) methods offer the advantage of yielding convergent estimators of the maximum likelihood of the data (Luengo et al., 2020) but at the cost of some large computational times, and, sometimes, of a complicated preliminary analysis to determine the likelihood function of the model. However, in cases where it is possible to simulate the process under study for a given set of parameters, other methods simply compare data to simulations (Cranmer et al., 2020; Hartig et al., 2011), either through summary statistics (Nunes and Balding, 2010) or by using the full information (Drovandi and Frazier, 2022). Approximate Bayesian Computation (ABC) methods (Sisson et al., 2018) belong to this class of algorithms. The original rejection-sampling method (Beaumont et al., 2002; Pritchard et al., 1999; Tavaré et al., 1997) has been replaced by new approaches combining ABC with iterative Monte Carlo methods as in Sequential Monte Carlo ABC (Dean et al., 2014; Del Moral et al., 2012), where the distributions are refined step by step or coupling ABC and MCMC (Marjoram et al., 2003; Wegmann et al., 2009). Bayesian Synthetic Likelihood methods are another important class of simulation-based algorithms, which estimate the likelihood of summary statistics using a multivariate normal distribution (Price et al., 2018; Wood, 2010). Other Bayesian methods exist, such as particle filtering (Fasiolo et al., 2016).

Some non-Bayesian methods are also used to estimate the maximum likelihood of a data set with respect to a given model. In this category, one can mention among others Expectation-Maximization algorithms (Dempster et al., 1977) or forward algorithms (Bollback et al., 2008) for Hidden Markov Models (HMM) . While being efficient

when the model is adapted, these methods also require in general some significant computation times or a substantial preliminary theoretical analysis.

Like ABC methods, but out of the Bayesian framework, the Fixed Landscape Inference MethOd (*flimo*), which we propose in the present paper, adjusts some summary statistics of the simulations to those of the data. But the idea behind *flimo* is to replace the time-consuming rejection-sampling approach by an efficient gradient descent phase. To increase its efficiency, *flimo* relies on algorithms that are usually only applicable in a deterministic framework. Many efficient local optimization algorithms for deterministic functions (Nocedal and Wright, 2006) exist, such as quasi-Newton algorithms that require Hessian computation. However, these methods need a smooth solution landscape with a limited number of local optima since, when this condition is not met, the convergence cannot be guaranteed. This explains why these methods are not suitable for the optimization of stochastic functions. By definition, for a given set of parameters, a stochastic function may return different values. The non-constancy of the value returned during different optimization cycles induces an instability of the landscape and, thus, many spurious local optima, preventing the correct estimation of the gradient. To overcome this limitation, in *flimo*, the solution landscape is stabilized by fixing the randomness of the simulation beforehand, drawing all needed random values from a unique random seed uniformly distributed on $[0, 1]$. Later on, these values will be reused for each optimization cycle, by transforming them into the appropriate distribution using quantile functions, a common approach to generate random values with a prescribed distribution. Thus, the simulations become deterministic, and the objective function to be minimized becomes stable. To our knowledge, these works are the first to use this deterministic quantile approach for random model inference problems.

To apply *flimo*, all one needs to do is to simulate the process and to choose some appropriate summary statistics to compare the data and the simulations. The adaptation of existing simulators to *flimo* is thus straightforward. Although *flimo* was developed to provide some point estimates of the unknown parameters, it can also be used to approximate their full distributions.

To illustrate the workings of *flimo*, we present three applications: a theoretical benchmark, namely the inference of the parameters of a g-and-k distribution; a population genetics problem, namely the estimation of the allelic selection parameter from some time series in a Wright-Fisher model; and the estimation of parameters of a Ricker model, that gives an example of a possible application in population dynamics. In the two first cases, we compare *flimo* to other existing methods, to highlight the advantages of each approach in terms of bias and precision of estimates as well as computation time.

2 Material and Methods

2.1 Description of the algorithm

The *flimo* algorithm relies on the construction of a regular deterministic objective function which is built from random simulations of the process to be modelled and which is then efficiently minimized. Contrarily to classical stochastic methods, *flimo*

needs a single random draw to evaluate all the candidate parameters θ and to select the best one according to the chosen summary statistics. This selection is performed using a deterministic algorithm. The steps of the *flimo* algorithm are described below and summarized in Algorithm 1.

In the studied probability space, denoted $(\Omega, \mathcal{A}, \mathbb{P})$, let X_θ be the random variable defining a drawing in the considered model, with parameters θ . X_θ is defined as an application $X_\theta : \omega \in \Omega \mapsto X(\omega; \theta)$. From a mathematical point of view, *flimo* works by fixing the outcome $\omega_0 \in \Omega$ and by considering a deterministic function called *simQ* : $\theta \mapsto X(\omega_0; \theta)$ which gives a drawing in the model for any θ parameters.

Algorithm 1 Fixed Landscape Inference MethOd

Input:

#Data
 y^{obs}
 #Defined by model
 n_{draw} #random draws for one simulation
 $simQ$ #adequate quantile simulator
 #Chosen by user
 n_{sim} #simulations to perform
 s #summary statistics
 d #distance between summary statistics

Output: $\hat{\theta} = \operatorname{argmin} J$

#Define quantiles matrix
 1: $\mathbf{R} \leftarrow (r_{i,j})$ with $r_{i,j} \sim \mathcal{U}([0, 1])$ i.i.d.
 #Define objective function
 2: $J : \theta \mapsto d(s(y^{obs}), s(simQ(\theta, \mathbf{R})))$
 #Run optimization
 3: $\hat{\theta} \leftarrow \operatorname{argmin} J$

2.1.1 Preliminary drawing of the randomness

A random simulation of a model is based on a certain number of draws of some simple random variables. For example, in the Wright-Fisher model (Ewens, 2004), a binomial distribution is drawn at each generation and another one at each sampling. To run *flimo*, the user needs to determine an upper bound of the number of these draws for a simulation. This value is noted n_{draw} . Then the user decides the number of simulations n_{sim} to perform to estimate the typical summary statistics of parameters θ . These simulations can be averaged for example. This choice is based on a trade-off between computation time and dispersion of the estimators. We suggest to test several values of n_{sim} , typically between 10 and 1000. Once n_{draw} and n_{sim} are fixed, the randomness is drawn. A matrix R of dimension $n_{sim} \times n_{draw}$ is set such that each entry $r_{i,j} \sim \mathcal{U}([0, 1])$ and the entries are independent. These values will be considered as the quantiles of each random draw involved in the n_{sim} simulations and converted

on purpose to a realization of the desired random distribution thanks to its quantile function parameterized using θ .

2.1.2 Special framework to obtain an empirical distribution

To obtain a convenient empirical distribution, the user must set n_{sim} to 1 and run several independent inferences. To improve performance, it is useful to set the initial condition of inference $n + 1$ to the inferred value of inference n .

2.1.3 Use of the randomness to carry out simulations

In the chosen model with parameter θ , the k^{th} draw of the process can be written $Z_k \sim \mathcal{L}_k(\theta)$ where \mathcal{L}_k is a probability distribution parameterized by θ . This draw may depend on the state of the system at step $k - 1$.

The cumulative distribution function (CDF) of \mathcal{L}_k is denoted by F_k^θ . Recall that the quantile function Q_k^θ is defined by $Q_k^\theta(q) = \inf\{x \mid F_k^\theta(x) \geq q\}$ for every $q \in [0, 1]$. Thus, Q_k^θ acts as the inverse function of F_k^θ since $Q_k^\theta(q) \leq x$ if and only if $q \leq F_k^\theta(x)$. In the most common cases, $F_k^\theta(Q_k^\theta(q)) = q$ for every $q \in]0, 1[$. The models considered in biology or ecology can be sophisticated but they are often based on compositions of usual distributions (normal distribution, Poisson distribution...). It is the quantile functions of these elementary blocks that must be known, which is generally verified in practice. Under *flimo*, each step $Z_k \sim \mathcal{L}_k(\theta)$ of simulation i is replaced by equation 1.

$$Z_k^i = Q_k^\theta(r_{i,k}) \quad (1)$$

By construction, $Z_k^i \sim \mathcal{L}_k(\theta)$. Indeed, if $U \sim \mathcal{U}([0, 1])$, then $Q_k^\theta(U)$ is a random variable with cumulative distribution function F_k^θ . This method of generation of pseudo-random numbers is called inverse transform sampling. Moreover, once the matrix R is fixed, each run, and thus the whole set of simulations, becomes deterministic for any value of θ . This also yields a global monotone coupling since larger position parameters of the distribution yield larger values of the drawn random variables.

Let $(\theta, \mathbf{R}) \mapsto simQ(\theta, \mathbf{R})$ denote the simulator using quantiles instead of random calls. It is crucial to underline that once the matrix R is fixed, $simQ(\theta, \mathbf{R})$ produces exactly n_{sim} independent simulations of the model with respect to θ , as a classical simulator would do. We note $y^\theta = simQ(\theta, \mathbf{R})$ the simulations performed for a set of parameters θ .

2.1.4 Building the objective function

The user must then choose some summary statistics s , appropriate for the model studied. Many studies have been carried out on this subject (Nunes and Balding, 2010). As for the ABC methods, this choice plays a major role in the quality of inference. One must also choose a distance d to compare the summary statistics of the data, noted y^{obs} , and the simulations y^θ . The Euclidean norm and the Mean Absolute Deviation are two reasonable options. Once these components are chosen, the objective function J is simply defined by equation 2.

$$J(\theta) = d(s(y^{obs}), s(y^\theta)) \quad (2)$$

The function J is deterministic with the same smoothness as the implied quantile functions (again for some fixed matrix R). The usual continuous probability distributions have a smooth quantile function; on the other hand for discrete distributions the quantile function is piecewise constant.

2.1.5 Deterministic optimization algorithm and automatic differentiation

A deterministic local optimization algorithm is then used to estimate $\underset{\theta}{\operatorname{argmin}} J(\theta)$. If the stochastic process involves only draws from continuous distributions (e.g. normal distributions), it is possible to use a gradient-based second-order, e.g. quasi-Newton-type algorithm. In the case of discrete distributions, two routes are available. If θ is low dimensional, a gradient-free method may be suitable. Otherwise, each discrete distribution can be replaced by an adequate continuous distribution, such as a normal distribution with same mean and variance.

When the (transformed) probability distributions are continuous, the J function is differentiable almost everywhere. It is then possible to accelerate the inference by using an Automatic Differentiation module (Bartholomew-Biggs et al., 2000; Revels et al., 2016). Thus, the gradient and the Hessian of J are computed automatically by *chain rule* and not by finite difference which is the standard method of estimating differentials. Automatic Differentiation reduces both the risks of numerical errors and the computational times.

2.2 Implementation

2.2.1 Packages overview

Implementations of *flimo* are freely available in the Julia package `Jflimo.jl` and in the R package `f1imo`, both in <https://metabarcoding.org/flimo>. The Julia implementation takes advantage of the good numerical performances of this language (Bezanson et al., 2017). This version is used for the three applications presented in the present paper. The R implementation is deposited on the CRAN. While the language R (R Core Team, 2021) is probably the data analysis language that biologists use the most, it suffers from some performance limitations when compared to Julia.

In Julia, the optimization functions come from the package `Optim.jl` (K Mogensen and N Riseth, 2018). The adequate application framework (when the objective function J is differentiable) allows us to use the *IPNewton* method (Wächter and Biegler, 2006), an interior-point Newton algorithm solving constrained optimization problems. One may use this method in combination with the Automatic Differentiation module `ForwardDiff.jl` (Revels et al., 2016). In the case of a non-differentiable problem, the Brent method (Brent, 2002) is adequate in the one-dimensional case and an implementation with the Nelder-Mead optimization method (Gao and Han, 2012) is provided for multidimensional problems.

Our R package proposes two modes: `flimoR` and `flimoRJ`. In `flimoR` mode, *flimo* is implemented with the `optim` function of the R `stats` package with the L-BFGS-B method (Byrd et al., 1995), a modification of the BFGS quasi-Newton method in the differentiable case. The `flimoRJ` mode uses the Julia functions implemented in `Jflimo.jl` thanks to the R package `JuliaConnector` (Lenz et al., 2022).

2.2.2 Building adequate simulator

Any classical simulator of the studied process can be used, simply adapting it according to the following basic procedure. Each time random draws are performed with a random function, this must be replaced by the associated quantile function. To ensure the independence of the draws, each quantile must be used only once for each simulation.

In the R case, one replaces the random functions (`rpois`, `rnorm`...) by their quantile version (`qpois`, `qnorm`...). In the Julia case, the procedure is the same with the packages `Random.jl` and `Distributions.jl`: the `rand` calls are replaced by some `quantile` calls. In both cases, the number of random drawings has to be replaced by the adequate submatrix of R . A tutorial is provided on the git page of the project.

2.3 Comparison to other inference algorithms

We present in detail the two applications used to study *flimo*. One can reproduce our results, using the scripts available on the git page of the project (<https://metabarcoding.org/flimo/flimo>). The main work was performed on a laptop MacBook Air (2017, 2.2 GHz Intel Core i7 Dual Core Processor). The results presented in supplementary material were parallelized on the GRICAD infrastructure servers.

2.3.1 Estimate of the parameters of g-and-k distributions

Definition

The family of g-and-k distributions can be viewed as an asymmetric generalization of the normal distributions $\mathcal{N}(\mu, \sigma^2)$, using two additional shape parameters g and k . Estimating the parameters of a g-and-k distribution from a random sample is a classical example used to evaluate ABC methods (Sisson et al., 2018). The density of a general g-and-k distribution is not explicit, but its quantile function is: equation 3 holds for every $q \in [0, 1]$. Here, $B > 0$, $k > -1/2$, and $z(\cdot)$ is the quantile function of the standard normal distribution. Normal distributions are the $g = k = 0$ case of this family.

$$Q(q | A, B, g, k) = A + B \frac{\left(1 + 0.8 \frac{1 - \exp(-gz(q))}{1 + \exp(-gz(q))}\right)}{(1 + z(q)^2)^k z(q)} \quad (3)$$

Compared inference methods

As is usually done, we choosed $\theta_{\text{true}} = (A, B, g, k)_{\text{true}} = (3, 1, 2, 0.5)$ (Drovandi and Pettitt, 2011). We generated 100 independent data sets y^{obs} of 1000 draws of the

distribution with parameters θ_{true} . One typical distribution is shown in Supplementary Figure S1. Other parameter sets have been studied under the same conditions.

The relative efficiencies of the MCMC, ABC and *flimo* methods used with different objective functions and computational efforts were measured, our goal being to disentangle the effects of these aspects of the optimization process on the quality of the results (Table 1). Parameter inferences are performed by each method for each of the 100 simulated data sets.

The MCMC method (available online: <https://github.com/pierrejacob/winference> (Bernton et al., 2019)) is considered the gold standard. The Sequential Monte Carlo ABC algorithm implemented in the same package `winference` and called here *wABC*, uses as summary statistics the Wasserstein distance of order 1 (Drovandi and Frazier, 2022), thus comparing the complete collection of empirical quantiles of the two distributions y^{obs} and y^θ . The objective function J is then defined as an average absolute deviation (equation 4).

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n |Y_{(i)}^{\text{obs}} - Y_{(i)}^\theta| \quad (4)$$

Here, n denotes the common observation size and the $Y_{(i)}$ are the order statistics $Y_{(1)} \leq \dots \leq Y_{(n)}$ of the sample $Y = (Y_1, \dots, Y_n)$ under consideration.

This statistics has been implemented to build the so-called *wflimo* method. We then explored two strategies to reduce the computation of the ABC and *flimo* methods. The first method aims to reduce the computational effort of the optimization algorithm, allowing it less time to converge (*wabc-short* method). The second aims to reduce the computational complexity of the objective function by replacing the Wasserstein distance with four summary statistics based on the empirical octiles (Drovandi and Pettitt, 2011), which significantly reduces the amount of information used (*oABC*, *oflimo* and *oflimo-short* methods). These summary statistics, called Moment Estimates and denoted by $s(y) = (S_A(y), S_B(y), S_g(y), S_k(y))$ for a set of realizations y , characterize the parameters of the distribution. With E_i for $1 \leq i \leq 8$ denoting the empirical octiles of the samples, the summary statistics are defined in equations 5.

$$\begin{aligned} S_A &= E_4 \\ S_B &= E_6 - E_2 \\ S_g &= \frac{E_6 + E_2 - 2E_4}{E_6 - E_2} \\ S_k &= \frac{E_7 - E_5 + E_3 - E_1}{E_6 - E_2} \end{aligned} \quad (5)$$

Following the recommendations of (Fearnhead and Prangle, 2012), the simulation procedure is accelerated. Indeed, it is possible, using once again the inverse transform simulation method, to simulate uniform distribution order statistics $(U_{(i)})_{1 \leq i \leq 7}$ with

the exponential spacings method (Ripley, 1987) (equation 6) for a data set of size 1000 and to convert them to realizations of a g-and-k distribution with respect to θ .

$$V_i \sim \Gamma(1000/8) \text{ independent,} \quad U_{(i)} = \frac{\sum_{j=1}^i V_j}{\sum_{k=1}^7 V_k} \quad (6)$$

Then, the objective function J chosen to evaluate the parameter sets θ is defined by equation 7.

$$J(\theta) = \left(\frac{S_A(y^{obs}) - S_A(y^\theta)}{S_A(y^{obs})} \right)^2 + \left(\frac{S_B(y^{obs}) - S_B(y^\theta)}{S_B(y^{obs})} \right)^2 + \left(\frac{S_g(y^{obs}) - S_g(y^\theta)}{S_g(y^{obs})} \right)^2 + \left(\frac{S_k(y^{obs}) - S_k(y^\theta)}{S_k(y^{obs})} \right)^2 \quad (7)$$

Method	Package: Implementation	Summary statistics	Bounds or Prior	Computation time control
<i>MCMC</i>	winference : <i>metropolishastings</i>	-	Prior : $\mathcal{U}([0, 10]^4)$ IC : inferred by <i>wABC</i>	8000 iterations (burn-in : 2000)
<i>oABC</i>	gk : 100 best simulations	Moment Estimates		$n_{sim} = 5 \times 10^6$
<i>wABC</i> <i>wABC-short</i>	winference : <i>wsmc</i>	1-Wasserstein distance	Prior : $\mathcal{U}([0, 10]^4)$	$max_{time} = 180s$ 1024 particles $max_{time} = 18s$ 100 particles
<i>oflimo</i> <i>oflimo-short</i>	Jflimo : IPNewton with AD	Moment Estimates	Bounds : $[0, 10]^4$	$n_{sim} = 1000$ $n_{sim} = 10$
<i>wflimo</i>	IPNewton without AD	1-Wasserstein distance	IC : $\mathcal{U}([0, 10]^4)$	$n_{sim} = 1$ best of 20 inferences

Table 1 Application framework for 100 inferences of the g-and-k distribution. IC stands for Initial Condition. For each non-*flimo* method, default setup are used (available online for **gk** (Prangle, 2017) and **winference** https://github.com/pierrejacob/winference/blob/master/inst/tutorials/tutorial_gandk.pdf). Initial Condition for *MCMC* is the mean value of each parameter inferred by *wABC*. For *wABC* inferences, the running time exceeds the time limit set in parameter which is treated as a stop condition if it has been exceeded at the previous iteration.

Estimation of parameters distribution

The *flimo* algorithm can also be used to approximate the distribution of the model parameters, as Bayesian methods provide posterior distributions. The *MCMC* method started with the true parameters is again used as a reference and the obtained posterior

marginal distributions are compared to those obtained by *wABC* and *wflimo*. To compare the posterior densities, a Kolmogorov-Smirnov test is used for each of the four marginal distributions of the parameters. This test is overpowered for these sample sizes, so we focused only on the test statistics D and used it as a distance between distributions. Recall that D is a distance of type L^∞ between the CDFs, thus, for CDFs F_1 and F_2 , D is defined by equation 8.

$$D = \sup_{x \in \mathbb{R}} |F_1(x) - F_2(x)| \quad (8)$$

2.3.2 Estimate of the selection value in a Wright-Fisher model

Model definition

Flimo has then been applied to infer the strength of selection in a Wright-Fisher model (Ewens, 2004) from some time series of allele frequencies (Paris et al., 2019). This is a typical study of population genetics problems. The distribution of two alleles A_0 and A_1 of a locus in a population of size N_e is simulated over several generations, with a selective value s applied on the allele A_1 . At regular intervals, the allele frequencies are estimated by sampling a part of the population. Noting $X(t)$ the actual proportion of the alternative allele A_1 at generation t , the model is written with a binomial distribution (equations 9 and 10). The function f is defined on $[0, 1]$, where $s \geq -1$ is the selective value of A_1 and $h \in [0, 1]$ is the dominance parameter.

$$X(t+1) | X(t) \sim \frac{1}{N_e} \mathcal{B}(N_e, f(X(t))) \quad (9)$$

$$f(x) = \frac{x(1 + sh + s(1 - h)x)}{1 + 2shx + s(1 - 2h)x^2} \quad (10)$$

The available data is sampled at times $t_1 = 0 < t_2 < \dots < t_n = T$. Let $X_k = X(t_k)$. At each observation time, n_k alleles ($n_k/2$ individuals) are sampled. The number of A_1 alleles sampled is denoted by Y_k , hence, conditionally on X_k , the distribution of Y_k is binomial (equation 11).

$$Y_k | X_k \sim \mathcal{B}(n_k, X_k) \quad (11)$$

Approach

We compared the *flimo* method with the *compareHMM* method (Paris et al., 2019). Thanks to the different approximations implemented in its model, this method is one of the fastest available today, while having a higher accuracy than e.g. the *WFABC* method (Foll et al., 2015) used in the same context. *CompareHMM* is non-Bayesian and relies on Maximum Likelihood Estimation. It follows a previous approach (Bollback et al., 2008) which uses a forward algorithm on the Hidden

Markov Model to compute the likelihood of the model with either the exact binomial model (*compareHMM-Bin*) or an approximate model using a Beta with spikes (*compareHMM-Bws*) distribution (Tataru et al., 2015) where transitions from one generation to the next are represented by a mixture model with a probability that the allele frequency is fixed at 0 or 1, and a Beta distribution conditional on non-fixation otherwise. *CompareHMM-Bin* has an algorithmic complexity of $\mathcal{O}(N_e^3)$ so its computation time makes it prohibitive for large populations. It was used as a reference for $N_e = 10^2$ and $N_e = 10^3$. The results of (Paris et al., 2019) have been reproduced from the python code available online (<https://github.com/CyrielParis/compareHMM/>).

For *flimo*, the unknown initial value of $X(0)$ is estimated by $\widehat{X}(0) = \frac{Y_0}{n_0}$. The unknown parameter is then just the A_1 selective value $\theta = s$ while the other quantities are assumed to be known. The objective to minimize is the mean absolute deviation around the median, as defined in equation 12.

$$J(\theta) = Mean(|Y_k^{\text{obs}} - Median(Y_k(\theta))|)_{1 \leq k \leq n} + 10^{-2} |\theta| \quad (12)$$

Here, $Median(Y_k(\theta))$ is the median of the n_{sim} simulations done with parameter θ . The correction term is present to avoid wrong convergence results because J is almost constant over a wide interval close to the lower bound. We used three different adaptations of the Wright-Fisher model. The first model is the Beta with spikes (Bws) approximated model as in *compareHMM-Bws* used with a gradient-based optimization. This model does not work well with *flimo* as the quantiles of the Beta distribution have no closed analytical form. One needs to inverse numerically the CDF $F(x; \alpha, \beta)$, a step which takes a substantial time. The samplings are simulated with approximated normal distributions instead of binomial distributions. The second model is the Nicholson Gaussian (NG) approximated model (Nicholson et al., 2002) with the same optimization process. Classically, each binomial distribution is replaced by the normal distribution with same mean and variance, with absorbing states at allele frequencies 0 and 1. The third model is the original binomial model with a gradient-free optimization. The number of simulations n_{sim} (10 and 200) was chosen for a compromise between efficiency and robustness. The exact model is more difficult to optimize due to its piecewise constant objective function. This is why it leads to worse performances and the approximation of the model is relevant.

Simulated data

We considered populations with three effective sizes $N_e = 10^2$, 10^3 , and 10^4 and the initial proportion $X(0) = 0.2$. The value $N_e = 10^2$ allows to test the robustness of the inferences despite the strong stochastic variations linked to the small population size, while $N_e = 10^4$ provides a very robust data set. In accordance with a range of tests performed in (Paris et al., 2019), we set the dominance parameter to $h = 0.5$ and we simulated over $T = 45$ generations with a sampling every $\Delta t = 5$ generations of $n_k = 0.3N_e$ alleles. In the two first scenarios, 100 data sets were simulated to compare the performance of the *flimo* method with Paris et al. (2019). For the main scenario studied, declined in two cases, $N_e = 10^2$ and $N_e = 10^3$, the selective value is set to

$s = 0.1$ (slight selective advantage for A_1). Data is shown in Supplementary Figure S3. Then we used a single data set with $N_e = 10^4$ and $s = 0.1$ to study the influence of the number of simulations in the dispersion of the estimate, with n_{sim} chosen in 10, 20, 50, 100, 200, 500 and 1000. Two extreme scenarios were also analyzed with $N_e = 10^3$, and with $s = 0.01$ and $s = 1$.

2.3.3 Chaotic population dynamics under Ricker model

Flimo is also applicable to more chaotic models. We present here an application based on the Ricker model which can describe erratic population dynamics. This model is known for the high non-linearity of its likelihood (Wood, 2010). This model is written according to equation 13. The observable data are $(Y_t)_{0 \leq t \leq T}$ for a given time limit T . It has three parameters: $\theta = (r, \sigma, \Phi) \in \mathbb{R}_+^3$. A typical realisation of this model is shown in Figure 1.

$$\begin{aligned} Y_t &\sim \text{Poisson}(\Phi N_t) \\ N_{t+1} &= r N_t e^{-N_t + e_t} \\ \text{with } e_t &\sim \mathcal{N}(0, \sigma^2) \text{ iid and } N_0 = 1 \end{aligned} \tag{13}$$

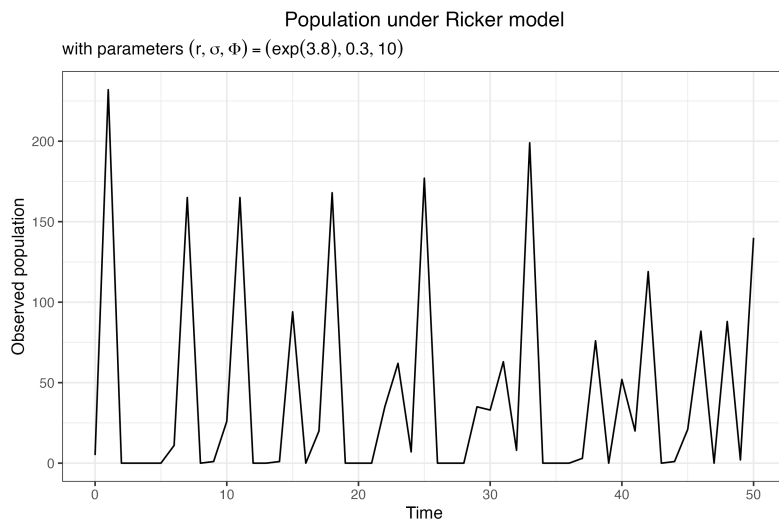


Fig. 1 Typical population simulated from Ricker model with parameters $(r, \sigma, \Phi) = (\exp(3.8), 0.3, 10)$ and $T = 50$.

The inference methods developed for this model fall into two main categories: an information-reduction approach (ABC or synthetic likelihood), or state space methods (particle MCMC, iterated filtering, parameter cascading) (Fasiolo et al., 2016). We used Wood (2010)'s work to develop an inference protocol adapted to *flimo*. The

summary statistics he used were studied and six of those that vary most regularly with the parameters, in the context of random simulations carried out by fixing the randomness beforehand were kept (s_1, \dots, s_6). Those are the 4 parameters of the cubic regression of the ordered $Y_{t+1} - Y_t$ values, the average population and the number of observed zeros. The objective function to be minimized is given by equation 14.

$$J(\theta) = \sum_{k=1}^6 \left(\frac{s_k(\text{sim}Q(\theta)) - s_k(\text{data})}{s_k(\text{data})} \right)^2 \quad (14)$$

Flimo was applied to 100 simulated data sets for seven different parameter sets, around the values $(r, \sigma, \Phi) = (\exp(3.8), 0.3, 10)$ studied by Wood (2010). The optimization method chosen is the Nelder-Mead gradient-free method, with either 100 or 500 (as in Wood (2010)) simulations used by *flimo*.

3 Results

3.1 Parameter inference for g-and-k distributions

3.1.1 Point estimate of the parameters

The *wABC* and *wflimo* methods provide estimates of all four parameters of g-and-k distributions, which are consistent with those obtained by *MCMC*. The averages of the estimates over the 100 simulated data sets yield p-values greater than 0.09 for all the different optimization methods hence, according to Wilcoxon tests with the usual significance level 0.05, they are not statistically distinguishable (Figure 2a-d).

The accuracy of the three methods is also comparable (Figure 2a-d). For *A* and *B*, the variances of the estimates are not significantly different. When comparing *wABC* to *MCMC*, the variance is multiplied by a factor of 1.06 for *A* and by a factor of 0.99 for *B* (the p-values of Fisher tests being 0.756 and 0.984, respectively). For the comparison between *wflimo* and *MCMC*, the increase factors are 1.14 and 1.03, respectively (with p-values 0.522 and 0.898). The shape parameters are estimated with less precision. To wit, the estimates of *g* are significantly more dispersed with the *wABC* and *wflimo* methods than with *MCMC*, with a variance ratio to *MCMC* of 2.38 and 1.97, respectively (p-values < 0.002 for both methods). The dispersion of the estimates of *k* is not significantly different between *wABC* or *wflimo*, and *MCMC*, the respective variance ratio 1.10 and 1.50 has a p-value close to the 5% threshold for the *wflimo* versus *MCMC* comparison (p-value 0.060). Finally, while no significant difference can be shown between the three methods in terms of bias and accuracy, *wflimo* runs 26.4 times faster than *MCMC* and 23.6 times faster than *wABC* (Figure 2e).

When the *wABC* parameters are adjusted to achieve a computation time comparable to *wflimo* (*wABC-short*), the method ceases to be reliable and strong biases on the estimates of the parameters *A* and *g* appear (Figure 2a and c), as well as an increase in the variance for each parameter compared to *wABC* (with variance ratios 2.58 for *A*, 1.38 for *B*, 83.7 for *g*, and 2.39 for *k*). If the second strategy, based on a

less complex objective function, is applied to the ABC method (*oABC*), it allows a correct estimation of the parameters A and g but the estimates of B and k become strongly biased. It is also less efficient in terms of computation time (Figure 2e). As regards the effects of these methods of reduction of the computational times, *flimo* is more robust. The use of the simplified objective function (*oflimo*) reduces the computational time compared to *wflimo* by a factor of 5.3 without introducing any bias on the estimates (the p-values comparing the estimates to the simulated parameters are 0.631, 0.139, 0.666, and 0.143). Only the variance of the estimates increases relative to *wflimo*, by ratios of 1.14, 1.67, 1.84, 4.72 respectively (with p-values 0.5, 0.01, 0.003, < 0.001 respectively). When the computational effort reduction is applied in conjunction with *oflimo* (*oflimo-short*), no significant estimation bias occurs (p-values: 0.803, 0.963, 0.758, 0.497) and no significant increase in the variances of the estimates compared to *oflimo* is observable (variance ratios: 1.10, 1.10, 1.05, 1.18; p-values: 0.632, 0.622, 0.880, 0.448). This last method, combining the two optimization procedures, runs 5157 times faster than *MCMC*, without introducing any significant bias on the parameters estimation and only increasing the variance by a factor of 1.45, and 1.91 for A and B , respectively. In contrast, for the two shape parameters, the increase in variance is much larger: 3.79 and 8.36 for g and k , respectively.

The inference results for other parameter sets, with comparable conclusions, are shown in Supplementary Figure S2. The time have not been included, due to different parallelisation processes between the different methods.

The different versions of *flimo* were compared for the *oflimo-short* method. The `flimoR` version is about 30 times slower than `Jflimo.jl` (median time of 0.50s instead of 0.016s for one inference). The `flimoRJ` mode has a similar computation time to `Jflimo.jl` with an additional fixed cost of about 2s corresponding to the switch from R to Julia, regardless of the number of inferences.

3.1.2 Estimation of the parameters distributions

The quantile function of the density obtained by *wflimo* is comparable to those obtained by the *MCMC* and *wABC* methods. (Figure 3a-d). For the parameters A and B , the three densities are closely related: D (equation 8) equals 0.076 versus 0.11 for A , 0.082 versus 0.039 for B , for the comparisons *wflimo* versus *MCMC* and *wABC* versus *MCMC*. For g , the distribution obtained by *wflimo* is closer to that obtained by *MCMC* than by *wABC* (0.25 versus 0.45). For k , the densities obtained by *wflimo* and *wABC* are equivalent (0.28 versus 0.21). However, the current implementation of *flimo* is not optimized to work in such a mode. Indeed, *wflimo* takes 252s to perform 1024 inferences, while *wABC* takes only 192s.

3.2 Estimation of a Wright-Fisher selection value

For $s = 0.1$ and for $N_e = 10^2$ and $N_e = 10^3$, both methods exhibit highly correlated s values over the 100 simulated data sets (namely, $R \in [0.88, 94]$ for $N_e = 10^2$ and $R \in [0.82, 87]$ for $N_e = 10^3$, Figure 4a-b and Table 2). Compared to *compareHMM-Bin*, for $N_e = 10^2$, *flimo* systematically overestimated s by about

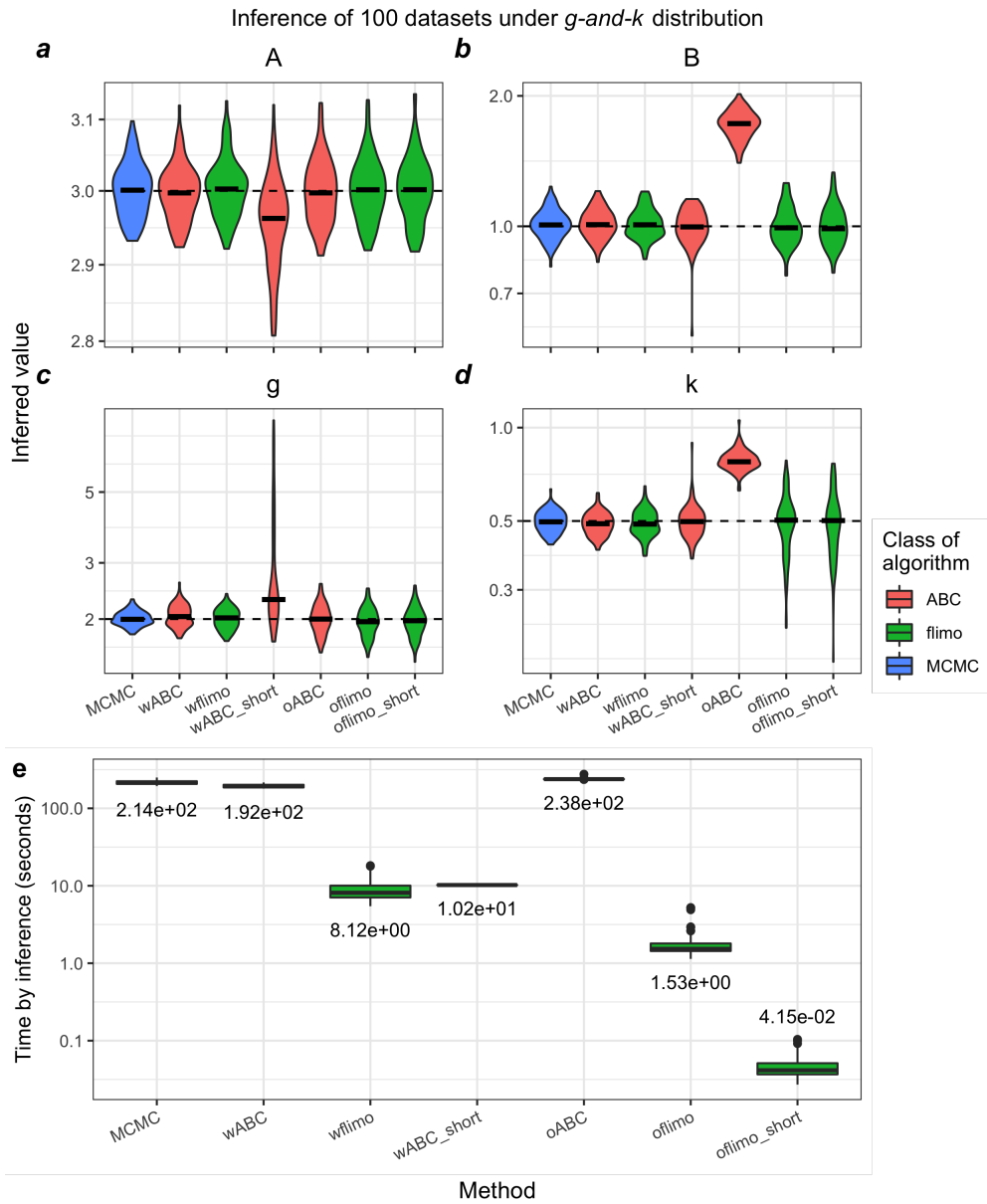


Fig. 2 Inference of the parameters of g -and- k distributions for 100 independent data sets of 1000 observations and different methods (see Table 1). Each y-scale is logarithmic. **Panels a-d**: Inference results for each parameter. Medians are plotted as bold black lines. Horizontal dashed lines correspond to the simulated values of the parameters. **Panel e**: Box plot of the 100 running times in seconds.

10 – 15%. This systematic over-estimation almost disappears for $N_e = 10^3$, with a difference of estimation between the methods ranging from 0.03% to 3% (Table 2). For

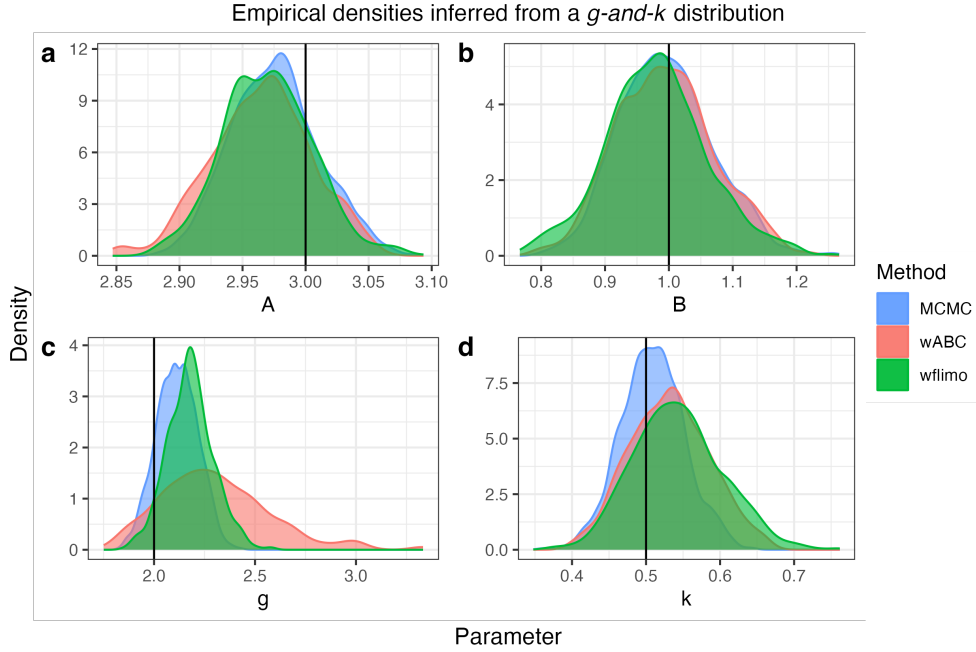


Fig. 3 The empirical distributions of the four parameters of g -and- k distributions as estimated by different methods. The estimation of the distributions by the method *MCMC*, used as reference (in grey), is compared to two approximations, the first one provided by the *wABC* method (in blue), the second one by the *wflimo* method (in red).

$N_e = 10^2$, the inferred mean value of s is 0.883 for *compareHMM-Bin* and 0.896 for *compareHMM-Bws* versus $\hat{s} \in [0.104, 0.114]$ for the different implementations of *flimo*. For $N_e = 10^3$, we have $\hat{s}_{compareHMM-Bin} = 0.0991$, $\hat{s}_{compareHMM-Bws} = 0.992$ and $\hat{s}_{flimo} \in [0.1001, 0.1013]$. This shows that the values inferred by *flimo* are not further from the expected value $s = 0.1$ than the values inferred by *compareHMM*. On average, an inference by *compareHMM-Bin* lasted 1.2s for $N_e = 10^2$ and 1.0 10^3 s for $N_e = 10^3$.

For $N_e = 10^4$, a single population was simulated. As expected, the differences between the s values estimated using *flimo* and *compareHMM-Bws* are very small (less than 1%) even with $n_{sim} = 10$, a very small number of simulations. The only noticeable effect of increasing the number of simulations used by the *flimo* algorithm is to reduce the standard deviation of the estimates by a factor close to $10^{-2}/\sqrt{n_{sim}}$ (Figure 4c). For the two extreme scenarios $s = 0.01$ and $s = 1$, *flimo* behaves similarly to what was observed for $s = 0.1$ (Supplementary Table S1).

3.3 Chaotic population dynamics under Ricker model

The inference results are shown in Table 3. The computation time for one inference is of the order of 15-20 seconds for $n_{sim} = 100$ and of 80-90 seconds for $n_{sim} = 500$ (the number used in Wood (2010)). The estimators are globally unbiased, with a rather large dispersion. In the case where the parameter r is maximal (89), the inference fails: it might be necessary to include other summary statistics to deal with this extreme

Population size	Criterion	<i>compareHMM</i>		<i>flimo</i>				
		<i>Bws</i>	<i>Binomial</i>		<i>Bws</i>		<i>NG</i>	
			10	200	10	200	10	200
$N_e = 10^2$	Correlation	0.9994	0.89	0.91	0.93	0.93	0.88	0.94
	Median difference $\hat{s}_{flimo} - \hat{s}_{compareHMM-Bws}$		0.011	0.011	0.0096	0.014	0.016	0.016
	Seconds by inference	1.4	0.016	0.29	0.17	3.3	0.014	0.11
$N_e = 10^3$	Correlation	0.9997	0.85	0.87	0.82	0.88	0.84	0.87
	Median difference $\hat{s}_{flimo} - \hat{s}_{compareHMM-Bws}$		$2.5 \cdot 10^{-3}$	$4.5 \cdot 10^{-4}$	$2.2 \cdot 10^{-4}$	$1.5 \cdot 10^{-3}$	$8.8 \cdot 10^{-3}$	$1.3 \cdot 10^{-3}$
	Seconds by inference	1.4	0.037	0.69	0.34	7.1	0.013	0.10

Table 2 Inference results, based on 100 simulated data sets, using *compareHMM-Bin*, *compareHMM-Bws* or *flimo* for two different numbers of simulations and three versions of the Wright-Fisher model. Three quantities are presented: the Pearson correlation coefficient, the median of the difference between the values inferred by *compareHMM-Bin* and the other methods, and the median of the computation times.

case. It is often necessary to run several inferences with different starting conditions, because of bad convergences that are detected by their high error score. Some remain in the results after 20 tries, especially in scenarios where the parameters r or ϕ are lower. These have a large impact on the standard deviation of the estimators. The passage from $n_{sim} = 100$ to $n_{sim} = 500$ does not seem to improve significantly the accuracy of the inference: increasing the number of possible tries rather than the number of simulations might be relevant.

The interest of *flimo* in this framework is that it is not necessary to consider the likelihood of the model, nor to evaluate the likelihood of the summary statistics as the synthetic likelihood approach do. Neither is it necessary to approximate the model to achieve a reasonable point inference of the parameters. The choice of suitable summary statistics is again crucial here: [Wood \(2010\)](#)'s study has been very useful for our work.

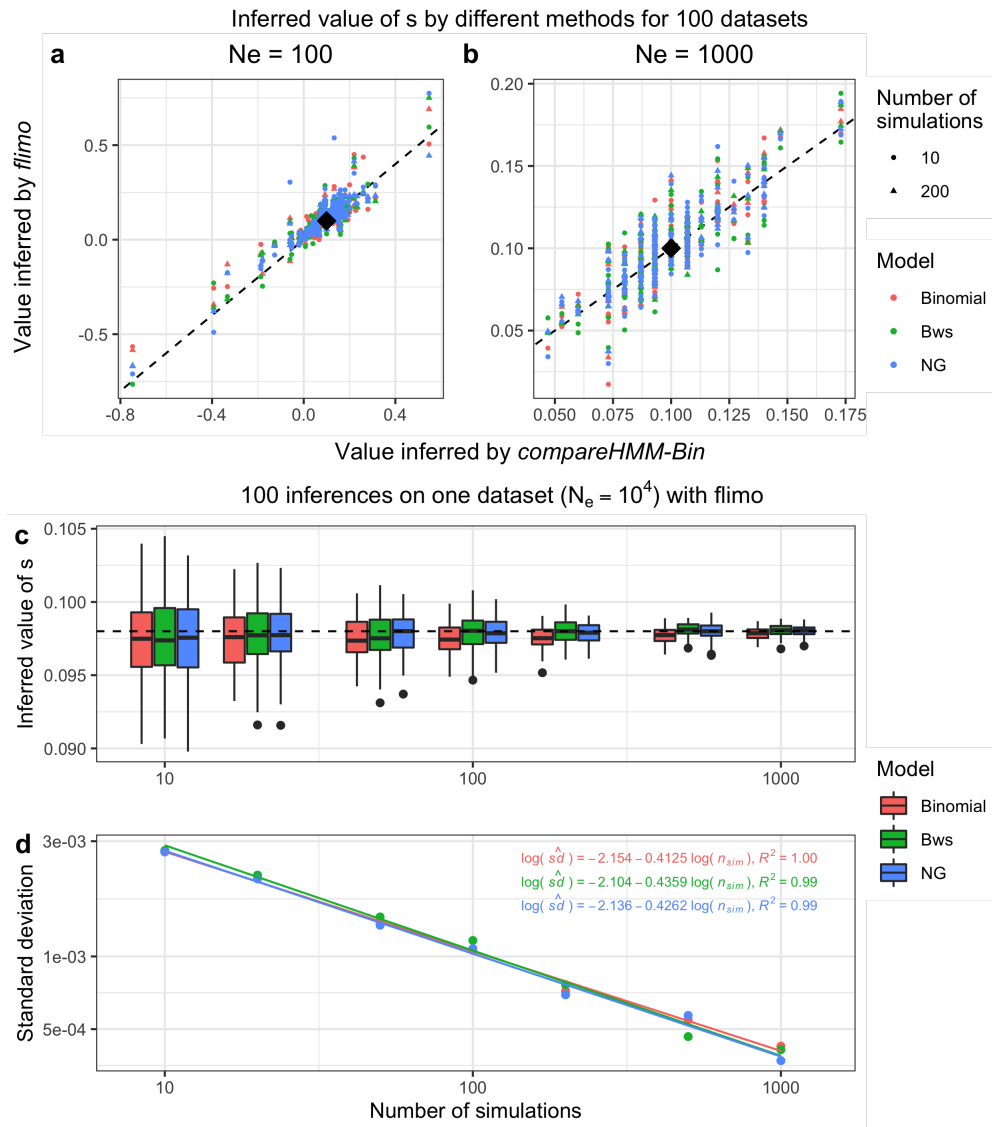


Fig. 4 Comparison of inferred selection values s by the *compareHMM* method and different implementations of *flimo*. **Panels a and b**: compared inferred values for one hundred data sets with effective population size $N_e = 10^2$ (**a**) and $N_e = 10^3$ (**b**). Dashed line corresponds to identity. Black diamond is the simulated value $s = 0.1$. **Panels c and d**: Influence of the number of simulations on *flimo* inferences for a population of size $N_e = 10^4$. One hundred inferences were run for every number of simulations. **Panel c**: inferred values of s . The dashed line represents the value inferred by *compareHMM-Bws*, $\hat{s} = 0.098$. **Panel d**: standard deviation of the inferred values of s with linear regression curve.

r^*	σ^*	Φ^*	n_{sim}	Median inferred r (1 <i>s.d.</i>)	Median inferred σ (1 <i>s.d.</i>)	Median inferred Φ (1 <i>s.d.</i>)	Detected Outliers
$\simeq 45$	0.3	10	100	44 (8.5)	0.30 (0.13)	10 (0.56)	0
			500	45 (8.8)	0.27 (0.14)	10 (0.58)	0
$\simeq 22$	0.3	10	100	21 (7.4)	0.31 (0.14)	10 (23)	12
			500	21 (7.9)	0.30 (0.15)	10 (25)	15
$\simeq 89$	0.3	10	100	82 (17)	0.41 (0.20)	10 (0.87)	0
			500	85 (15)	0.35 (0.18)	10 (0.73)	0
$\simeq 45$	0.15	10	100	42 (7.6)	0.15 (0.13)	10 (9.0)	1
			500	42 (7.7)	0.15 (0.13)	10 (9.0)	2
$\simeq 45$	0.6	10	100	44 (15)	0.54 (0.21)	9.9 (0.91)	0
			500	44 (15)	0.56 (0.21)	9.9 (0.89)	0
$\simeq 45$	0.3	5	100	42 (12)	0.30 (0.19)	5.0 (21)	5
			500	42 (13)	0.31 (0.20)	5.0 (23)	6
$\simeq 45$	0.3	20	100	44 (10)	0.26 (0.14)	20 (1.4)	0
			500	45 (10)	0.29 (0.15)	20 (1.4)	0

Table 3 Inference results by *flimo* on several populations simulated under the Ricker model, for different parameter values and different computational effort. The exact values of r^* are $exp(3.8)$, $exp(3.8)/2$ and $2exp(3.8)$. This value of $exp(3.8)$ is used in [Wood \(2010\)](#).

4 Discussion

Concerning the point inference of parameters, *flimo* has the advantage of being considerably faster than the other methods with comparable accuracy: up to 5,000 times faster for the g-and-k example, and up to 100 times faster on the Wright-Fisher model. Obviously, the efficiency of an algorithm strongly depends on its implementation, especially on the programming languages used. The algorithm *flimo* is implemented in Julia, a language known for its good numerical performances. The other algorithms used here are implemented either in C++ for *MCMC* and *wABC*, or in Python with *Numpy* for *compareHMM*, which are comparably efficient programming languages (Aruoba and Fernández-Villaverde, 2015). This ensures that its computation speed is an intrinsic property of the *flimo* algorithm. Like in ABC methods, the choice of the summary statistics is important. However, *flimo* seems to be less sensitive than the ABC implementations tested here. This lower sensitivity allows to select some summary statistics that can be calculated quickly, thus decreasing the optimization time of *flimo*. This is illustrated by the use of the Wasserstein distance which is necessary to have good estimates of the parameters of g-and-k distributions using the ABC method, while it only slows down *flimo* with low benefit on the accuracy of the results. The behaviour of *flimo* in high dimension (in terms of number of parameters) is the purpose of a further investigation (Supplementary File 1), based on a toy example from Li et al. (2017). It seems realistic to infer a few dozens of parameters.

Flimo also allows to obtain suitable empirical estimates of the distribution of the inferred parameters, but at the cost of its computation efficiency. Obtaining a distribution of parameters rather than a point value is relevant, in particular for estimating the uncertainties of these estimated parameters. An interesting strategy could be to implement some hybrid methods, using *flimo* to circumscribe a region of interest in the search space, which could then be used as a fine prior for a Bayesian method.

As with ABC, the model-based approach of *flimo* has its limitations. In particular, the models selected by the user bias the analysis, since this prior choice influences the conclusions that will be drawn from the algorithm's results (Csilléry et al., 2010). However, *flimo* allows for the comparison of several different scenarios, provided that the same summary statistics are used for the different models. The way in which *flimo* works does, however, imply some constraints that other simulation-based methods do not have. First of all, the models studied must verify some properties. On the one hand, it is necessary to be able to calculate efficiently the quantiles of the concerned distributions. This is not always the case, as for Beta distributions. In these cases, it is advantageous to approximate the distribution, for example by a normal distribution. Moreover, *flimo* needs that the states of the model do not have large discontinuities, otherwise the optimization algorithms tend to fail. For example, *flimo* is not suitable for building a most likely phylogenetic tree. In addition, the usual constraints of optimization problems apply here: there must not be many local minima for the considered objective function, that should ideally be convex. Unlike Bayesian methods which can explore the parameters space widely, *flimo* tends to revolve around a single area of the search space, which depends on the initial conditions of the optimization. It is therefore necessary in practice to perform several inferences to ensure good convergence. Resampling methods, like jackknife, can be used to make the inference more robust.

In general, it does not seem possible to prove the convergence of the point estimators established by *flimo* to the theoretical parameters of the model. This is mainly due to the fact that summary statistics do not have a priori any particular properties, and therefore matching summary statistics from data and simulations does not prove that the parameter values are similar. A simple example of a calculation where convergence can be shown is given in Supplementary File 2. This illustrates the fact that the *flimo* approach (both in the simulation and minimisation process) does not in itself induce an asymptotic bias.

This study shows that *flimo* is a solid alternative to other inference methods, especially ABC, with a simple implementation which, thanks to the Julia and R packages that we developed. It can easily be adapted to a wide range of models from different fields, for example in population dynamics or population genetics. By rethinking the role of randomness in stochastic model inference problems, *flimo* makes it possible to use efficient deterministic gradient-based optimization algorithms to infer the parameters of probabilistic models whose likelihood or moments are intractable. All these qualities make *flimo* a particularly simple and effective method of inference.

5 Statement and Declarations

5.1 Competing interests

This work was supported by the Alpalga project (ANR-20-CE02-0020). The authors have no relevant financial or non-financial interests to disclose.

5.2 Data and Code availability

All the data and scripts are available on <https://metabarcoding.org/flimo>.

6 Author contributions statement

All authors conceived the algorithm. S.M., C.G. and E.C. wrote the manuscript. E.O. and D.P. contributed to writing the manuscript. S.M. developed the packages and performed computational experiments with assistance and guidance from C.G., E.C. and E.O. C.G. and E.C. supervised the project.

7 Acknowledgments

We thank Pierre Pudlo and Adeline Leclercq-Samson for their help on some theoretical aspects of the paper, and Simon Boitard for sharing with us his work on the Wright-Fisher model.

Some of the computations presented in this paper were performed using the GRICAD infrastructure (<https://gricad.univ-grenoble-alpes.fr>), which is supported by Grenoble research communities.

1.4 Résultats complémentaires

1.4.1 Couplage entre *flimo* et un algorithme EM ?

Une limite de *flimo* évoquée par Adeline Leclercq-Samson (Laboratoire Jean Kuntzmann, Grenoble) est que l'algorithme repose sur une heuristique, sans garantie théorique de convergence. Nous avons donc émis l'idée d'adapter l'approche de *flimo* à un algorithme d'Espérance-Maximisation (EM), et plus particulièrement à la version SAEM qui repose sur des simulations aléatoires. Comme expliqué en section 0.8.2 (dont je reprends les notations), il est parfois compliqué d'estimer l'espérance de la log-vraisemblance jointe du modèle à chaque itération k : $Q(\theta, \theta_k)$. L'algorithme SAEM contourne cette difficulté en estimant Q par des simulations. Puisque l'étape M consiste à maximiser Q , *flimo* semble adéquat dans ce cadre. Malgré cela, je ne suis pas parvenu à établir de résultat dans cette direction. Le calcul suivant montre la limite à laquelle j'ai été confronté. L'itération k de l'inférence par SAEM est :

- **Simulation** des variables latentes \mathbf{x}_k selon $p(\cdot | \mathbf{y}; \theta_{k-1})$;
- **Mise à jour de Q** : $Q(\theta, \theta_k) = Q(\theta, \theta_{k-1}) + \gamma_k (\log p(\mathbf{x}_k, \mathbf{y}; \theta) - Q(\theta, \theta_{k-1}))$ où (γ_k) est une suite décroissante telle que $\sum \gamma_k = +\infty$ et $\sum \gamma_k^2 < +\infty$;
- **Maximisation** : $\theta_k = \operatorname{argmax}_{\theta \in \Theta} Q(\theta; \theta_{k-1})$

On peut expliciter le calcul de la mise à jour de Q (équation 1.35). Cela explicite le fait que les simulations $(\mathbf{x}_i)_i$ dépendent des paramètres successifs θ_{i-1} . Or l'intérêt de *flimo* est de maximiser une fonction reposant sur des simulations selon le paramètre à optimiser, en l'occurrence θ . Mais ici, au moment de maximiser $Q(\theta, \theta_k)$, les simulations sont déjà fixées, et ne dépendent dans tous les cas pas de θ .

$$\begin{aligned}
 Q(\theta, \theta_k) &= Q(\theta, \theta_{k-1})(1 - \gamma_k) + \gamma_k \log p(\mathbf{x}_k, \mathbf{y}; \theta) \\
 &= (Q(\theta, \theta_{k-2})(1 - \gamma_{k-1}) + \gamma_{k-1} \log p(\mathbf{x}_{k-1}, \mathbf{y}; \theta))(1 - \gamma_k) + \gamma_k \log p(\mathbf{x}_k, \mathbf{y}; \theta) \\
 &= \dots \\
 &= \sum_{i=1}^k \gamma_i \prod_{j=i+1}^k (1 - \gamma_j) \log p(\mathbf{y}, \mathbf{x}_i(\theta_{i-1}); \theta)
 \end{aligned} \tag{1.35}$$

Pour cette raison, cette piste de développement n'a pas été poursuivie.

1.5 Conclusion

Ce chapitre a présenté la méthode *Fixed Landscape Inference MethOd* (*flimo*) pour inférer des paramètres dans des modèles aléatoires. Cet algorithme est particulièrement efficace pour établir des estimateurs ponctuels des paramètres. Les domaines d'application sont vastes et *flimo* pourra devenir à l'avenir une alternative pertinente aux autres méthodes d'inférence sans vraisemblance comme les ABC.

Un bilan plus global et des perspectives sont exposés dans la Discussion générale du manuscrit.

Chapitre 2

Métabarcoding quantitatif par mesure des biais

2.1 Introduction

Ce chapitre s'intéresse à un grand enjeu du métabarcoding : l'acquisition de données quantitatives fiables pour calculer des indices de biodiversité. C'est un champ de recherche très actif mais les biais identifiés sont encore mal pris en compte. J'ai présenté en 0.6 des méthodes utilisées actuellement pour les corriger. Le manuscrit de ce chapitre présente nos travaux pour mesurer et corriger deux biais du métabarcoding : le biais d'amplification PCR et le biais de concentration de l'ADN cible par rapport à la concentration d'ADN total. Il a été soumis en octobre 2023 auprès du journal *Molecular Ecology Resources*. La section *Supplementary Information* mentionnée dans le manuscrit est placée en Annexe B.

Je présente ensuite plusieurs résultats études complémentaires. Le premier concerne un nouveau modèle mécanistique de PCR que nous avons développé. Je montre après des résultats de modélisation de PCR en présence de mismatches d'amorces. Ensuite, je présente un autre protocole expérimental testé pendant ma thèse pour mesurer la compétition entre espèces durant la PCR grâce à un contrôle interne. La dernière partie est consacrée à une étude de dosage de l'ADN cible par ddPCR en collaboration avec Stefaniya Kamenova (University of Oslo, Norvège).

2.2 Résumé en langue française

Les analyses par métabarcoding connaissent actuellement un grand succès du fait de leurs performances pour le suivi de la biodiversité. Cependant, il est encore difficile de tirer des conclusions quantitatives précises sur les écosystèmes étudiés, principalement à cause de biais inhérents à l'ADN environnemental ou introduit durant le protocole expérimental. Ces biais altèrent le lien entre la quantité d'ADN observée et la biomasse ou le nombre d'individus des espèces détectées. Deux des biais inhérents au métabarcoding ont été mesurés : le ratio entre les concentrations d'ADN total et d'ADN cible, et le biais d'amplification par PCR. Une méthode de correction est proposée. Tous les tests expérimentaux ont été réalisés sur une communauté artificielle de plantes alpines avec le marqueur *Sper01*, qui est supposé induire un faible biais d'amplification du fait de la bonne conservation de ses sites d'amorçage. Notre approche combine des techniques de PCR quantitatives standards (qPCR et digital droplet PCR) et un modèle aléatoire de PCR réaliste prenant en compte la saturation. Le modèle a été utilisé pour estimer les efficacités de PCR de chaque espèce et pour inférer leurs vraies proportions dans la communauté artificielle, à partir de la fréquence de leur lectures. Les corrections sont faciles à mettre en place et peuvent être appliquées à des données de métabarcoding produites auparavant. Ces travaux montrent l'importance relative des deux biais considérés et sont une porte ouverte au métabarcoding quantitatif, même si de nombreux autres biais doivent encore être considérés.

2.3 Manuscrit

Towards quantitative DNA Metabarcoding: A method to overcome PCR amplification bias

Sylvain Moinard^{1*}, Didier Piau², Frédéric Laporte¹, Delphine Rioux¹, Pierre Taberlet¹, Christelle Gonindard-Melodelima^{1*†} and Eric Coissac^{1*†}

¹Univ. Grenoble Alpes, Univ. Savoie Mont Blanc, CNRS, LECA, FR-38000, Grenoble, France.

²Univ. Grenoble-Alpes, CNRS, Institut Fourier, FR-38000, Grenoble, France.

*Corresponding authors. E-mails: sylvain.moinard@univ-grenoble-alpes.fr; christelle.gonindard@univ-grenoble-alpes.fr; eric.coissac@metabarcoding.org;

Contributing authors: didier.piau@univ-grenoble-alpes.fr; frederic.laporte@univ-grenoble-alpes.fr; delphine.rioux@univ-grenoble-alpes.fr; pierre.taberlet@univ-grenoble-alpes.fr;

†These authors contributed equally to this work.

Abstract

Metabarcoding analyses have recently undergone significant development due to the power of this technique in biodiversity monitoring. However, it is still difficult to draw accurate quantitative conclusions about the ecosystems studied, mainly because of biases inherent in the environmental DNA or introduced during the experimental process. These biases alter the relationship between the amount of DNA observed and the biomass or number of individuals of the species detected. Two of the biases inherent in metabarcoding have been measured: the ratio between total DNA and target DNA concentrations, and the PCR amplification bias. A method for their correction is proposed. All experimental tests were performed on mock alpine plant communities using the marker *Sper01*, which is expected to have low amplification bias due to its highly conserved priming sites. Our approach combines standard quantitative PCR techniques (qPCR and digital droplet PCR) with

2 *Correcting PCR bias in metabarcoding data*

a realistic stochastic model of PCR dynamics that accounts for PCR saturation. The model was used to estimate PCR efficiencies for each species and to infer the true species proportions of the mock communities from the read relative frequencies. The corrections are easy to implement and can be applied to previously generated DNA metabarcoding data. This work demonstrates the relative importance of the two biases considered and is an open door to quantitative metabarcoding data, although many other biases remain to be considered.

Keywords: Amplification bias, droplet digital PCR, PCR model, Quantitative metabarcoding, Taqman qPCR

Introduction

In the context of mass species extinction (Barnosky et al., 2011), biodiversity assessment is currently a major challenge. Classically, biodiversity inventories consist not only of a list of species occurring at a site, but also of quantitative data assessing the abundance of each species. Traditional approaches based on direct observation by taxonomists may be unrealistic in terms of available skills and costs, given the enormous effort required to conduct such a survey on a global scale and across the tree of life. Therefore, high-throughput methods, including DNA metabarcoding (Taberlet, Coissac, Pompanon, Brochmann, & Willerslev, 2012), are the only chance to achieve such a goal. DNA metabarcoding has been used for more than a decade in many areas of ecology, such as biodiversity monitoring (e.g. Bohmann et al., 2014), detection of invasive species (e.g. Klymus, Marshall, & Stepien, 2017), or tracking animal diets (e.g. Pompanon et al., 2012). It is now part of the basic toolbox of ecologists, if we consider more than a thousand articles published annually based on this technique. While metabarcoding provides a satisfactory species inventory (Beng & Corlett, 2020; Ficetola & Taberlet, 2023; Taberlet et al., 2012) with some insight into their relative abundance (Pornon et al., 2016), the quality of quantitative data produced is questionable (Krehenwinkel et al., 2017; Yang et al., 2021).

The relationship between the abundance of a species in the field and the number of sequence reads measured in a DNA metabarcoding experiment is far from straightforward. Many reasons can lead to biased abundance estimates. Biases arise from both natural properties and technical issues (Luo, Ji, Warton, & Yu, 2022; van der Loos & Nijland, 2021). At least three natural biases can be considered. First, if the amount of DNA shed into the environment depends on the biomass of individuals (Elbrecht & Leese, 2015; Elbrecht, Peinert, & Leese, 2017; Lamb et al., 2019), it is also a function of shedding rates specific to each DNA source (Wilder, Farrell, & Green, 2023). Second, the relationship between the eDNA sampled, and the DNA actually shed depends on its decay rate, which in turn depends on the ecosystem studied (Andruszkiewicz Allan,

Zhang, Lavery, & Govindarajan, 2021; Krehenwinkel et al., 2018). Third, the number of copies of the DNA marker targeted by metabarcoding per unit of biomass or per individual varies from species to species (Garrido-Sanz, Senar, & Piñol, 2022; Krehenwinkel et al., 2017; Zoschke, Liere, & Börner, 2007), and may also vary among tissues, during development or according to phenology. Two main sources can be considered for technical biases. First, the DNA extraction method, whose efficiency depends on the extracted substrate and varies between taxonomic groups (Dopheide, Xie, Buckley, Drummond, & Newcomb, 2019). Second, the PCR amplification can incur species-specific amplification biases (Pawluczyk et al., 2015) related to the annealing step (Piñol, Mir, Gomez-Polo, & Agustí, 2015) or to the PCR extension step, which may depend, among other things, on the GC content of the metabarcodes (Nichols et al., 2018). Thus, the sum of all these biases obscures the relationship between the abundance of the sequenced reads and the relative abundance of the species in terms of biomass or number of individuals.

Metabarcoding thus requires an appropriate statistical approach to robustly estimate species abundances (Alberdi & Gilbert, 2019; Mächler, Walser, & Altermatt, 2021). For a long time, that quantification problem has been considered. Authors have proposed improvements by optimizing the choice of primers (Krehenwinkel et al., 2017), by varying the number of PCR cycles for different replicates (Silverman et al., 2021) or by creating mock communities to infer correction factors with one species of interest and one control species (Thomas, Deagle, Eveson, Harsch, & Trites, 2016), with two species of interest in different quantities (Matesanz et al., 2019) or by comparing several mock communities of more complex composition (Krehenwinkel et al., 2017); or to infer PCR efficiencies (Shelton et al., 2022). Internal controls can be used, but these do not allow measuring amplification bias (Smets et al., 2016; Ushio et al., 2018).

The present paper examines the biases introduced by the most commonly criticized step of DNA metabarcoding, the PCR amplification. The strength of the amplification bias and its impact on the estimated abundances of metabarcoding are assessed. This study is based on a new mathematical model of PCR amplification that is applicable to the simulation of DNA metabarcoding experiments. Several models exist to describe PCR dynamics (*e.g.* Carr & Moore, 2012; Hayward, 1998; Mehra & Hu, 2005) but have not been linked to metabarcoding. The model developed from existing models considers the amplification bias between species in conjunction with the saturation phase of PCR amplification, with a minimum number of parameters. A usual model in quantitative metabarcoding is the exponential model, also called log-ratio linear model (*e.g.* Gold et al., 2023; Kelly, Shelton, & Gallego, 2019; Shelton et al., 2022), where the abundance of each species increases geometrically during the PCR. The non-treatment of saturation is not a problem in quantitative real-time PCR (qPCR) because the amplification starts with an exponential phase, but is incompatible with metabarcoding PCR, which relies on the final state of the system.

4 *Correcting PCR bias in metabarcoding data*

The impact of low priming site conservation on species detection and quantification of COI markers has been widely discussed. These biases are related to the annealing phase of PCR cycles due to primer mismatches (Clarke, Soubrier, Weyrich, & Cooper, 2014; Piñol et al., 2015; Pompanon et al., 2012). To specifically target the biases induced by the extension step of PCR, we assessed them on three mock alpine plant communities using the *Sper01* marker (Taberlet et al., 2007). This marker is widely used in many ecological studies: soil biodiversity (Yoccoz et al., 2012), paleoecology based on ancient eDNA (Willerslev et al., 2014) or diet (Valentini et al., 2009). Although there is very little variation at the *Sper01* priming sites, no strong annealing bias can be assumed for this marker. However, the length of the metabarcodes and the complexity of its sequence (length and frequency of homopolymers) varies from species to species, making it an appropriate candidate to study extension bias. PCR efficiencies for three species were accurately estimated using Taqman qPCR to calibrate our model and then to infer the pre-PCR eDNA proportions of each species. Combined with precise estimates of target DNA concentrations in each species by droplet digital PCR (ddPCR), the results of this experiment demonstrate the benefit of handling PCR extension bias and the variation of target DNA concentration among taxa to correctly estimate taxa abundance from DNA metabarcoding results. Although only a single marker was studied here on a limited number of species, the presented protocol is easily generalizable and opens perspectives for quantitative DNA metabarcoding (qMetabarcoding).

Material and Methods

Metabarcoding experiment

Quantification biases were investigated using three mock communities composed of thirteen alpine plants belonging to the *Spermatophyta* clade (Supplementary Table 1), using the *Sper01* primer (Taberlet, Bonin, Zinger, & Coissac, 2018; Taberlet et al., 2007) targeting the P6 loop of the *trnL* of the chloroplast genome. Plant species were selected for having no mismatches at their priming sites with the *Sper01* primers.

Plant sampling

Plants leaves were collected in Chartreuse and Belledonne massif in the French Alps during Spring 2021 (Supplementary Table 1). Freshly collected material was stored in silica gel before DNA extraction.

DNA Extraction

Plant DNA was extracted using the CTAB protocol (Doyle, 1990), except for *Carpinus betulus*, for which a *DNeasy Plant Mini Kit* (Qiagen) was used after unsuccessful CTAB extractions.

Quantification of target DNA

The total DNA concentration for each plant sample was determined using Qubit (ThermoFisher). The amount of DNA targeted by the *Sper01* primer is not proportional to the total DNA concentration, as the number of chloroplasts per cell is expected to vary between different species and tissues and during plant development (Golczyk et al., 2014; Sakamoto & Takami, 2018; Zoschke et al., 2007). ddPCR was used to provide absolute quantification of the *Sper01* target DNA. ddPCR was preferred over qPCR because it is much less affected by inhibition than qPCR, which varies from sample to sample. (Sjostedt, Rådström, & Hedman, 2020). This quantification was performed using serial dilutions of total DNA concentrations ranging from $6.25 \times 10^{-2} \text{ ng}/\mu\text{l}$ to $6.25 \times 10^{-5} \text{ ng}/\mu\text{l}$ with one or two replicates for each condition. The reaction mixtures had a total volume of $20 \mu\text{l}$ ($5 \mu\text{l}$ of DNA solution, $10 \mu\text{l}$ of Master Mix EvaGreen, $0.6 \mu\text{l}$ of primers (forward and reverse) at $10 \mu\text{M}$, $4.4 \mu\text{l}$ of milliQ water). The *QX200 Droplet Digital System* (Bio-Rad) was used to generate droplets (*QX200 Droplet Generator*) and to analyze them after PCR amplification (*QX200 Droplet Reader* with the *QuantaSoft Software*). Thermocycler conditions with optimized annealing temperature for the *Sper01* primer (52°C) were set (30 seconds at 95°C , 30 seconds at 52°C , one minute at 72°C). Replicates identified as incorrect by the reader and the most diluted replicate in cases where this concentration was outside the expected detection range were removed.

In this study, the concentration index chosen to compare the samples is the expected number of target copies per *ng* of total DNA. It is calculated from each assay as in the equation 1. The number of copies per μl (in target DNA) is the value measured by ddPCR. $C(\text{Total DNA})_{\text{replicate}}$ is the total DNA concentration of the sample in the reaction mix. The average concentration for each species is used for the rest of the protocol.

$$\text{Concentration}(\text{Copies}/\text{ng DNA}) = \frac{(\text{Copies}/\mu\text{l})_{\text{ddPCR}}}{C(\text{Total DNA})_{\text{replicate}}} \quad (1)$$

This choice of index was made because mock communities are composed of purified plant DNA extracts. This correction includes two bias factors: the number of target copies per genome and the genome size. In the case of an eDNA sample, the amount of target DNA detected does not depend on the genome size. The aim is to quantify the number of sampled cells, so the only factor to correct is the difference in copy numbers per genome. Two choices are then possible. Either, the concentration is established in copies per gram of tissue (equation 2), and then $(\text{Copies}/\mu\text{l})_{\text{ddPCR}}$ is calculated for a known mass of tissue $m(\text{Tissue})_{\text{replicate}}$. Or the concentration is established by taking into account the genome size using the C-value (equation 3, Supplementary Table 1) with the same ddPCR assays as in equation 1. The Kew C-value database

6 Correcting PCR bias in metabarcoding data

(<https://cvalues.science.kew.org/>) can be used to implement this correction, possibly using a parent species as an approximation.

$$\text{Concentration(Copies/g tissue)} = \frac{(\text{Copies}/\mu\text{l})_{\text{ddPCR}}}{m(\text{Tissue})_{\text{replicate}}} \quad (2)$$

$$\text{Concentration(Copies/genome)} = \frac{(\text{Copies}/\mu\text{l})_{\text{ddPCR}}}{C(\text{Total DNA})_{\text{replicate}}} \times C_{\text{value}} \quad (3)$$

Mock communities

Three mock communities were constructed after the ddPCR assays: (i) a uniform community (\mathcal{M}_U) where each plant has the same concentration of target DNA, (ii) a community where each plant has the same concentration of total DNA (\mathcal{M}_T), and (iii) a community where the concentrations of target DNA are distributed according to a geometric sequence of common ratio 1/2 (concentrations of 1, 1/2, 1/4...) (\mathcal{M}_G). The species used are described in Table 1. The metabarcode sequences are given in the Supplementary Table 1 and the exact composition of each community is given in the Supplementary Table 2. The comparison between \mathcal{M}_U and \mathcal{M}_T communities allows determining the bias introduced by variation in the number of chloroplast genomes per unit of total DNA. The \mathcal{M}_U and \mathcal{M}_G comparison allows the estimation of relative PCR extension step efficiencies.

Species	Short form	Length	GC content (%)	Total DNA concentration (ng/ μ l)	Rank (\mathcal{M}_G)
<i>Briza media</i>	Bme	53	39.6	183	1
<i>Rosa canina</i>	Rca	51	31.4	50.8	2
<i>Lotus corniculatus</i>	Lco	55	38.2	65.2	3
<i>Populus tremula</i>	Ptr	68	25.0	31.4	4
<i>Salvia pratensis</i>	Spr	46	26.1	24.4	5
<i>Lonicera xylosteum</i>	Lxy	46	32.6	45.8	6
<i>Fraxinus excelsior</i>	Fex	39	33.3	22.4	7
<i>Acer campestre</i>	Aca	56	39.3	12.2	8
<i>Capsella bursa-pastoris</i>	Cbp	48	45.8	38.8	9
<i>Geranium robertianum</i>	Gro	53	34.0	15.0	10
<i>Carpinus betulus</i>	Cbe	61	27.9	9.14	11
<i>Abies alba</i>	Aal	47	44.7	3.58	12
<i>Rhododendron ferrugineum</i>	Rfe	46	30.4	3.90	13

Table 1: Plants used for the three mock communities and their characteristics for the *Sper01* marker. Total DNA concentrations are assayed in the samples after extraction by Qubit. Rank stands for decreasing abundance in the \mathcal{M}_G community.

DNA metabarcoding PCR amplification 212

For each community, 20 replicates ($2\mu\text{l}$ of DNA) and one PCR negative control ($2\mu\text{l}$ of milliQ water) are made. Three wells are left blank (sequencing controls). Each well was individually tagged. 40 PCR cycles were run with an optimized annealing temperature for *Sper01* (30 seconds at 95°C , 30 seconds at 52°C , one minute at 72°C).

Metabarcoding DNA Sequencing 218

High-throughput sequencing was performed on NextSeq (Illumina) by Fasteris (Plan-les-Ouates, Switzerland; <https://www.fasteris.com/>). One library was constructed per community following the Metafast protocol (as proposed by Fasteris).

Bioinformatic pipeline 223

All the bioinformatic work was performed on a laptop MacBook Air (2017, 2.2 GHz Intel Core i7 Dual Core Processor). The data and analysis scripts are available on the project's git page, <https://github.com/LECA-MALBIO/metabar-bias>. Raw data was processed with OBITools (version 4 aka OBITools4; Boyer et al., 2016, <https://metabarcoding.org/obitools4>). Unless otherwise stated, the further analyses were carried out using R.

A DNA metabarcoding experiment model 230

The goal of this part is to estimate the initial relative abundances p_s of each species s , from the number of reads R_s among the S different species in the considered environmental sample. To achieve this, a simulation model was used to generate virtual metabarcoding data that are compared with observed data.

The model integrates the three steps involved in the production of a DNA metabarcoding result from a DNA extract, as in Gold et al. (2023): i) the sampling of a portion of the DNA extract, ii) the PCR amplification, iii) the sampling of a portion of the PCR reaction for sequencing.

Sampling of a portion of the DNA extract 239

The initial number of molecules in a replicate r , $M_0^s(r)$, is modeled by a Poisson distribution with expectation m_0^s .

On the basis of simulations, it was observed that the final observed proportions had a standard deviation about 25 times higher than the proportions in equivalent simulated data. Such a standard deviation can be obtained in simulations by replacing the Poisson distribution with a negative binomial distribution with parameters $r = \frac{m_0^s}{\delta^2}$ and $p = \frac{1}{\delta^2}$, so that its expectation is m_0^s and its variance is $m_0^s \delta^2$, with δ the overdispersion of the standard deviation obtained with an initial Poisson distribution (with here $\delta \simeq 25$). Our choice for the Poisson distribution simplifies the model and the mean value remains unchanged.

$$M_0^s \sim \text{Poisson}(m_0^s) \tag{4}$$

so that $\mathbb{E}[M_0^s] = m_0^s$ and $\text{Var}(M_0^s) = m_0^s$

PCR amplification

The used PCR model, here called logistic model, accounts for the different amplification efficiencies and the saturation phase. The model is an adaptation of the deterministic model in [Hayward \(1998\)](#), which also features a logistic form with an additional parameter. It is also a simplified form of the model in [Carr and Moore \(2012\)](#). Compared with a conventional exponential model, the logistic model accounts for saturation phase at the end of the PCR. Both are parametric stochastic models. Figure 1 illustrates kinetics described by these two parametric stochastic models. The models are fitted to qPCR data generated from a sample of *Capsella bursa-pastoris*. The fit is a non-linear regression with more weight given to the acceleration cycles (15th-25th cycles).

The models considered describe the evolution of the number of DNA molecules of each species cycle by cycle, denoted M_k^s for each species s at PCR cycle k . Each molecule already present is maintained and has a probability λ_k^s of being replicated again, modeled by a binomial distribution (equation 5) depending on the state of the system after cycle $k - 1$. More precisely, amplification of species s depends on M_{k-1}^s and on $M_k = \sum_t M_{k-1}^t - M_0^t$ because the saturation depends on all the molecules that were created before cycle k (equation 6).

$$M_k^s | M_{k-1}^s, M_k \sim M_{k-1}^s + \text{Bin}(M_{k-1}^s, \lambda_k^s) \tag{5}$$

Let K be a charge capacity K , ie the total number of DNA molecules that can be created during the amplification. This quantity is not known a priori. Due to saturation, the effective PCR efficiency of each species, λ_k^s , decreases during the PCR. The logistic saturation has been chosen for its simple shape (equation 6).

$$\lambda_k^s = \begin{cases} \Lambda_s \left(1 - \frac{M_k}{K}\right) & \text{if } X_k \leq 1 \\ 0 & \text{otherwise} \end{cases} \tag{6}$$

The purely exponential model is a special case with no saturation where $\lambda_k^s = \Lambda_s$ at each cycle k . Under the exponential model, the expected number of molecules at cycle k is given by equation 7.

$$M_k^s = M_0^s (1 + \Lambda_s)^k \tag{7}$$

Sampling of a portion of the PCR reaction for sequencing

All the molecules created by the PCR are not sequenced: only a fraction constitutes the observed data, denoted R_s for each species s . At the end of a PCR amplification with n cycles, the sequencing step is described as a multinomial sub-sampling of R_{total} molecules (equation 8).

$$R_1, \dots, R_S | M_n^1, \dots, M_n^S \sim \text{Multinom} \left(R_{\text{total}}, \left(\frac{M_n^1}{\sum_t M_n^t}, \dots, \frac{M_n^S}{\sum_t M_n^t} \right) \right) \quad (8)$$

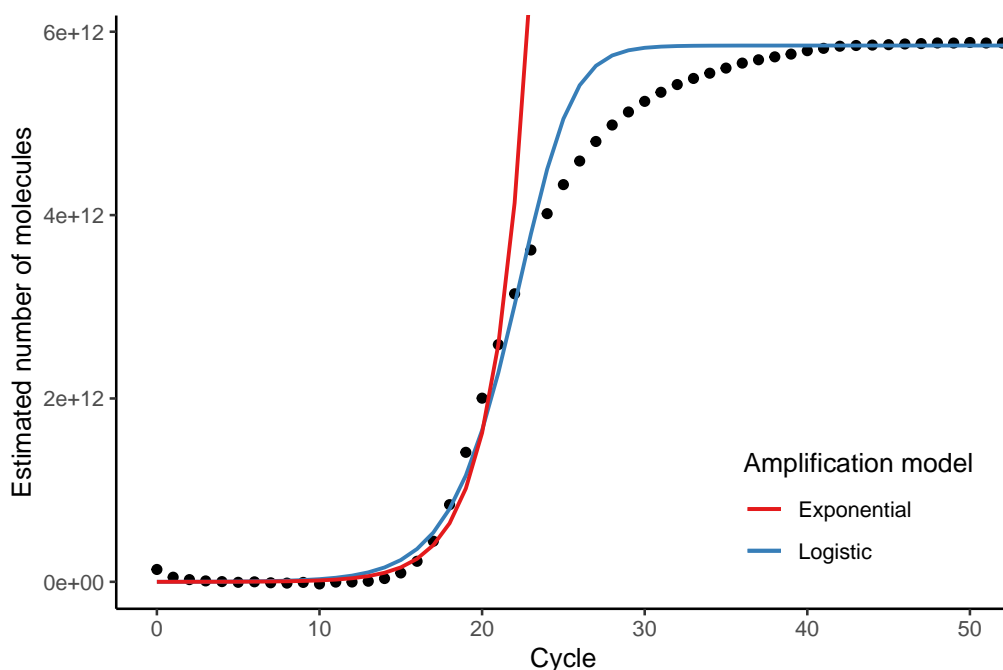


Fig. 1: Observed qPCR kinetics for a sample of *Capsella bursa-pastoris* (black dots) compared to two PCR models fitted to the data. Blue curve: logistic model; red curve: exponential model. An asymmetry of amplification is observed around the inflection point, which creates here a gap between the 25th cycle and the 40th cycle for the logistic model. This asymmetry is taken into account in more sophisticated PCR models (Gottschalk & Dunn, 2005). It is known that the first qPCR cycles correspond to a RFU background noise (Rao et al., 2013).

Measure of the amplification efficiencies

283

Using Taqman qPCR assay

284

PCR amplification efficiencies Λ_s were measured by qPCR for three of the plant species present in our mock communities: *Carpinus betulus*, *Capsella bursa-pastoris* and *Fraxinus excelsior*. These three species were chosen because their metabarcodes differ widely in sequence length and GC content. This makes it possible to expect different amplification efficiencies and to design specific Taqman internal probes that allow individual PCR efficiency measurements within a mixture of the three plant DNAs. Two different probes were designed for *Carpinus betulus* to evaluate the influence of the probe itself on the measurement. The four probes used are described in the Supplementary Table 3. The assay was performed using Taqman qPCR on a uniform community composed of these three species. A 5-fold serial dilution from 654 to 1 copies/ μl in the reaction mix (25 μl with 5 μl of DNA) was performed for each probe, with three replicates per concentration. Taqman qPCR was chosen to measure PCR efficiency because it allows measurement from a mixture of the three plant DNAs. This ensures the same inhibitory effect for each species. Since each individual DNA extract has its own pool of inhibitors that interfere with qPCR assays, independent measurement on pure extract would not be realistic (Svec, Tichopad, Novosadova, Pfaffl, & Kubista, 2015). This approach uses the quantification cycles C_t measured by qPCR.

285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303

The exponential model (equation 7), which is valid before the PCR saturation phase, can be used to estimate apparent PCR efficiencies. Estimated efficiencies are referred to as apparent efficiencies because inhibition is always present. For this study, however, only the relative values of the efficiencies are important. A commonly used formula (equation 9, Gill, Bleka, & Fonnelløp, 2022) can be derived from the exponential model to estimate amplification efficiencies from a series of qPCRs performed on successive dilutions. However, a major limitation of this formula that has been identified here is that the estimation of the slope is very sensitive to small variations in C_t , resulting in a large variance of the estimator of the efficiency Λ .

304
305
306
307
308
309
310
311
312
313

$$\begin{aligned} \text{Linear regression: } C_t &= -\frac{\log_{10}(m_0)}{\log_{10}(1 + \Lambda)} + \frac{\log_{10}(M_{C_t})}{\log_{10}(1 + \Lambda)} \\ &= a \log_{10}(m_0) + b + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma^2) \text{ iid} \\ \Lambda &= 10^{-1/a} - 1 \\ \text{and } M_{C_t} &= 10^{-b/a} \end{aligned} \tag{9}$$

To estimate efficiencies more precisely, this approach is modified. This linear regression approach was used only to estimate M_{C_t} , the number of molecules present at C_t . This quantity was assumed to be shared by all species. The value retained was the average of the estimated value for each species:

314
315
316
317

$M_{C_t} = 1.5 : 10^{11}$ (1.8×10^{11} , 2.3×10^{11} , 1.1×10^{11} and 1.3×10^{11} for probes CbeA, CbeB, Cbp and Fex, resp.) (equation 9).

From this value of M_{C_t} , other quantities can be estimated: K , the total number of molecules that can be created during amplification, and Λ_s , the PCR efficiencies.

The estimate of K is given by equation 10. This equation is established from observed relative fluorescence unit (RFU) values, assuming within-replicate proportionality between RFU and DNA copy number (Gill et al., 2022), although RFU values are not standardized, and fluorescence saturates at the end of amplification and depends on many experimental factors (Svec et al., 2015). This approach only provides an estimate: $K = 7.6 \times 10^{12}$ here (see Discussion). For the inference of initial abundances, the value of K used was estimated numerically (see below).

$$K \simeq M_{C_t} \times \sum_{s=1}^3 \frac{RFU_{End}(s)}{RFU_{C_t}} \quad (10)$$

Then, the efficiencies Λ_s were estimated for each replicate, for which the initial quantities M_0^s are known (equation 11). For subsequent analyses, the average Λ_s over all replicates is used.

$$M_{C_t} = M_0^s (1 + \Lambda_s)^{C_t(s)}$$

$$\text{so } \Lambda_s = \left(\frac{M_{C_t}}{M_0^s} \right)^{1/C_t(s)} - 1 \quad (11)$$

If the value of M_{C_t} chosen is not the average but the minimum value 1.1×10^{11} (resp. maximum value 2.3×10^{11}), the inferred values of Λ_s are multiplied by a factor of 0.980 (resp. 1.03).

Using the \mathcal{M}_U community

PCR efficiencies Λ_s were also inferred by fitting the logistic PCR model presented above to the experimental data. K is also inferred at the same time in this approach. For the \mathcal{M}_U community, the known quantities are the read numbers R_s and the initial quantities m_0^s of each species. The Fixed Landscape Inference MethOd (*flimo*, Moinard, Oudet, Piau, Coissac, & Gonindard-Melodelima, 2022) implemented in Julia was used for this purpose. This algorithm is based on model simulations compared with observed data using summary statistics, in the same way as ABC methods. It has the advantage of being faster. The *flimo* method minimizes an objective function in the form of a χ^2 statistic (equation 12).

$$\begin{aligned} & \operatorname{argmin}_{0 \leq \Lambda_1, \dots, \Lambda_S \leq 1, K > 0} J_{m_0}((\Lambda_s)_s, K) \\ \text{with } & J_{m_0}(\Lambda_1, \dots, \Lambda_S, K) = \sum_{s=1}^S \frac{(\overline{p}_s(\text{data}) - \widehat{p}_s)^2}{\overline{p}_s(\text{data})} \end{aligned} \quad (12)$$

where $\widehat{p}_s = \frac{R_s}{R_{\text{total}}}$ is the average proportion of species s in a replicate, estimated over $n_{\text{sim}} = 190$ simulations for given $(m_0^s)_s$, $(\Lambda_s)_s$ and K , and $\overline{p}_s(\text{data})$ is the average proportion of species s in the data.

The inferred efficiencies are relative, as the model can produce similar results for different ranges of Λ_s , especially when the value of K changes. The maximum efficiency value has been set at 1. 101 parameter inferences were performed. For the inference of initial proportions in \mathcal{M}_T and \mathcal{M}_G , the values retained for Λ_s and K are those associated with the median value of K .

Correction of relative abundances of a MOTU

Figure 2 summarizes the additional pipeline recommended for correcting amplification bias in a metabarcoding experiment. The PCR amplification efficiency of each species is estimated from samples of species characteristic of the ecosystem studied that are assayed by ddPCR. There are two ways of doing this: Taqman qPCR or a mock community study. These efficiencies are then used to infer the initial proportions of each species.

Using the Ratio method

Previous works (*e.g.* Shelton et al., 2022; Silverman et al., 2021) showed that a reference mock community can be used to correct abundances in another community composed of the same species. Although this was not the main objective of our work, this result was verified using the three communities studied. The \mathcal{M}_U community was used as a reference to correct abundances in the \mathcal{M}_T and \mathcal{M}_G communities. In the \mathcal{M}_U community, each species had a starting relative frequency of $1/13 \simeq 7.7\%$, which should have been observed in the final read proportions in the absence of amplification bias. The correction factor for each species c_s is therefore simply the median ratio between the expected and the observed reads frequencies over all replicates in the \mathcal{M}_U community (equation 13).

$$c_s = \operatorname{Median} \left(\frac{\text{Observed reads frequency}}{\text{Expected reads frequency}}(s) \right) \quad (13)$$

For the \mathcal{M}_T and \mathcal{M}_G communities, this correction factor is applied to estimate the initial proportions \widehat{p}_s for each species s with R_s reads (equation 14).

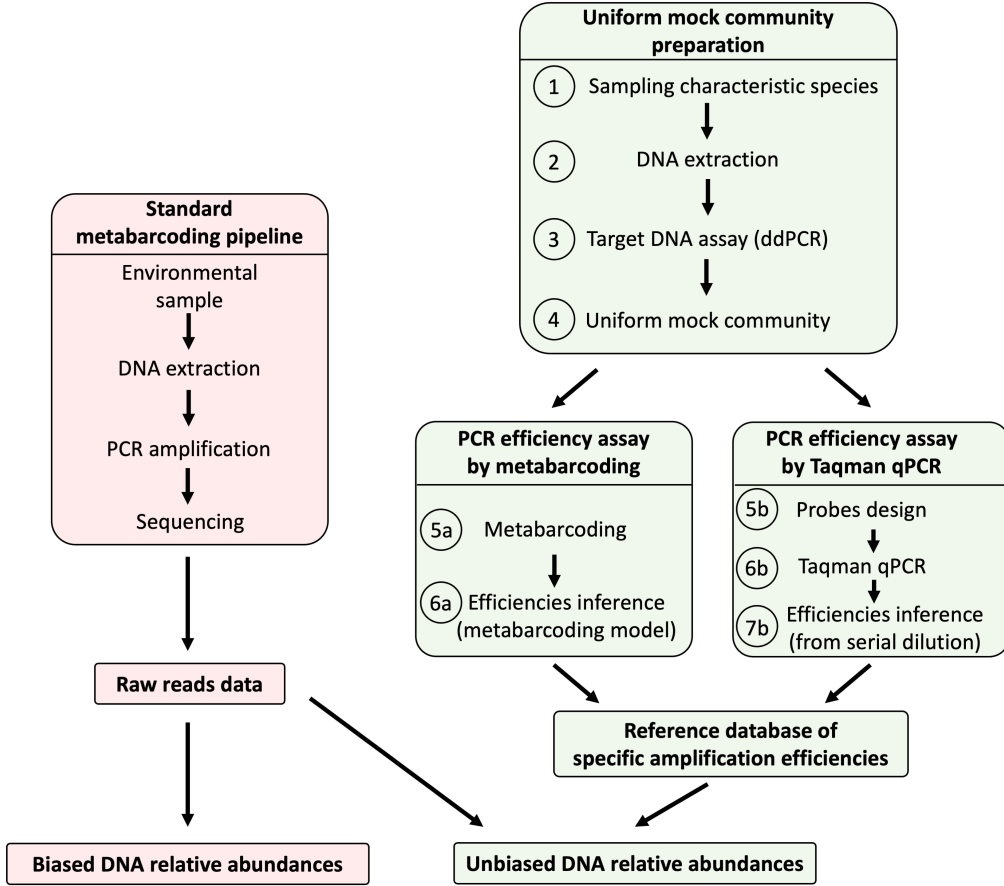


Fig. 2: Additional pipeline recommended for correcting amplification bias in a metabarcoding experiment as presented in this study.

$$\begin{aligned}
 R'_s &= \frac{R_s}{c_s} \\
 \hat{p}_s &= \frac{R'_s}{\sum_t R'_t}
 \end{aligned} \tag{14}$$

Using the estimated amplification efficiencies

378

The inference of the actual proportions of eDNA from the relative read abundances (RRA) measured after DNA metabarcoding sequencing is achieved by the same algorithmic method presented above, but this time the Λ_s efficiencies are assumed to be known, either measured by Taqman qPCR or inferred from the model fit for the \mathcal{M}_U community. The parameters to be inferred are the initial quantities $(m_0^s)_s$ for the \mathcal{M}_T and \mathcal{M}_G communities (equation 15).

379

380

381

382

383

384

$$\operatorname{argmin}_{m_0^1, \dots, m_0^S > 0} J_\Lambda((m_0^s)_s)$$

$$\text{with } J_{\Lambda}(m_0^1, \dots, m_0^S) = \sum_{s=1}^S \frac{(\overline{p_s}(\text{data}) - \hat{p}_s)^2}{\overline{p_s}(\text{data})} \quad (15)$$

An estimate of these proportions can be obtained using the exponential model, but this requires knowledge of the PCR effective number of “exponential cycles” at saturation (equation 17) which cannot be inferred simultaneously to Λ_s and K with the exponential model. In average, the final proportion p_s of each species s is:

$$p_s = \frac{R_s}{R_{\text{total}}} = \frac{m_0^s(1 + \Lambda_s)^{n_{\text{eff}}}}{K} \quad (16)$$

Taking the logarithm and summing all the species :

$$n_{\text{eff}} = \frac{\log\langle p_s \rangle + \log K - \log\langle m_0^s \rangle}{\log(1 + \Lambda_s)} \quad (17)$$

where $\langle p_s \rangle$, $\langle m_0^s \rangle$ and $\langle 1 + \Lambda_s \rangle$ are the geometric mean of the p_s , m_0^s and $1 + \Lambda_s$ ($\langle p_s \rangle = \left(\prod_{s=1}^S p_s\right)^{1/S}$, etc.).

Criteria for measuring quantification errors

The distance between the observed or corrected proportions $(\hat{p}_s)_s$, median over all the replicates) and the initial theoretical proportions (p_s^{th}) is measured by two RMSE (*Root-Mean-Square Error*) criteria. The error measured is either absolute (equation 18) or relative (normalized by the theoretical proportions, equation 19).

$$\text{Absolute Error: AbsErr}((\hat{p}_s)_s) = \sqrt{\frac{1}{S} \sum_{s=1}^S (\hat{p}_s - p_s^{\text{th}})^2} \quad (18)$$

and

$$\text{Relative Error: RelErr}((\hat{p}_s)_s) = \sqrt{\frac{1}{S} \sum_{s=1}^S \left(\frac{\hat{p}_s - p_s^{\text{th}}}{p_s^{\text{th}}}\right)^2} \quad (19)$$

Ecological conclusions: biodiversity indices

To compare theoretical, observed and inferred compositions, biodiversity indices were computed for \mathcal{M}_T and \mathcal{M}_G . Hill numbers (Hill, 1973) (equation 20), interpretable as an effective number of species in the community, were chosen with $q = 1$ (linked to Shannon entropy) and $q = 2$ (linked to Gini-Simpson index).

$${}^q D = \left(\sum_{s=1}^S p_s^q \right)^{\frac{1}{1-q}} \quad (20)$$

Results

405

ddPCR assay

406

The concentrations of each plant sample measured by ddPCR are shown in Figure 3. For the same total DNA concentration, there was a wide variability in average target concentration, ranging from 3.7×10^4 copies per *ng* for *Rhododendron ferrugineum* to 2.5×10^5 copies per *ng* for *Populus tremula* with an average of 1.1×10^5 copies per *ng* among the thirteen species. The factor between the extremes is thus 6.6.

412

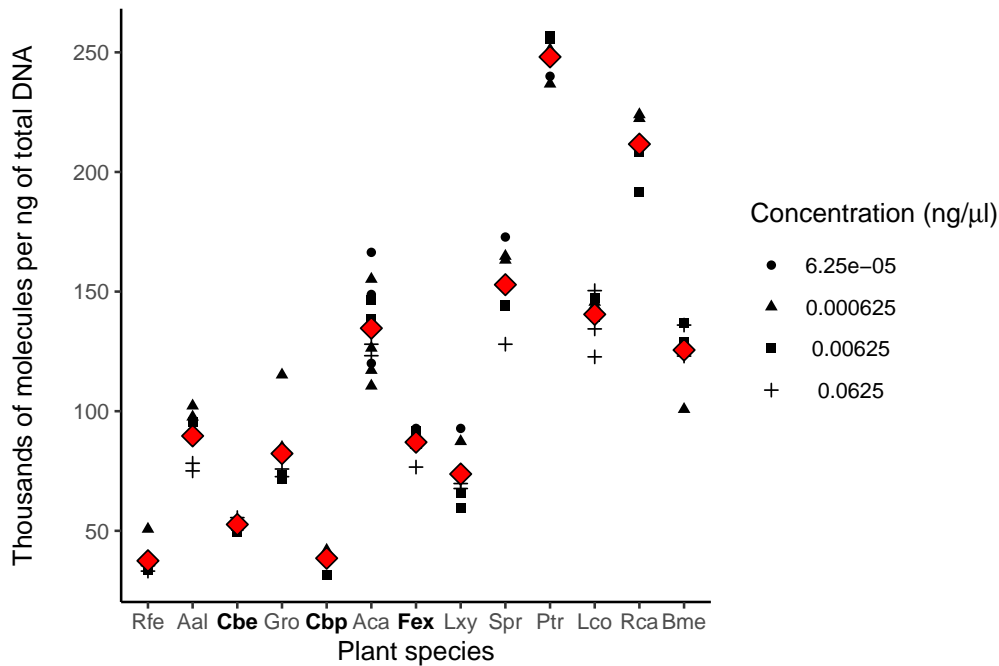


Fig. 3: Number of target DNA molecules (thousands) per *ng* of total DNA for thirteen alpine plants, computed with the index used in equation 1. Each black dot is a replicate, for different total DNA concentrations. The red diamonds correspond to the mean for each species.

Metabarcoding experiment

413

Raw sequencing data

414

After processing with the OBITools, an average of 37,000 reads per non-negative replicate was obtained with a standard deviation of 27,000 reads (first

416

and third quartiles : 14,000 and 56,000 reads). Negative controls showed negligible contamination. For each community out of the 20 PCR replicates, one replicate with fewer than 5,000 reads was discarded from further analysis.

Reads proportions

The comparison of observed and expected read proportions is shown in Figure 4. Significant differences can be observed: at most, between the observed and expected proportions, there is a factor of 3.0 for *Geranium robertianum* in the \mathcal{M}_U community, 4.2 for *Abies alba* in \mathcal{M}_T and 9.0 for *Abies alba* in \mathcal{M}_G .

Comparing the observed proportions with the expected proportions allows to visualize the two biases under study. For example, *Rosa canina* species has both good efficiency and a high target concentration: the two biases add up. Conversely, *Geranium robertianum* is penalized by both biases. *Salvia pratensis* has a higher-than-average concentration, but poor efficiency. *Capsella bursa-pastoris* is well amplified, but its target concentration is low.

The joint effect of the double bias is visible for \mathcal{M}_T , with median proportions comprised between 1.5% and 26%, and between 2.6% and 17% for \mathcal{M}_U .

Inter-replicate variability is significant in some species, such as *Populus tremula* (in \mathcal{M}_U : mean proportion : 8.6%, varying from 3.3% to 14%, standard deviation of 2.7%).

Inferring PCR efficiencies and abundances

The apparent PCR efficiencies for the three species tested (**Fex**, **Cbe**, **Cbp**) measured using the Taqman qPCR method for the four probes have a relative difference of the order of 5%. That can be considered low, but due to the exponential nature of PCR, it has a real impact on the final proportions in the community due to the exponential nature of PCR amplification.

Table 2 shows the abundances in the reference mock community \mathcal{M}_U and the efficiencies inferred from the Taqman qPCR assay and from the model fit to the \mathcal{M}_U community, with the *flimo* method. The standard deviation of the PCR efficiencies established for each probe by Taqman qPCR varies between 0.0049 (Cbp) and 0.011 (Fex). When inferred from \mathcal{M}_U , the standard deviation of Λ_s is between 0.0027 (Cbp) and 0.010 (Gro). The lowest efficiency is around 15% lower than the maximum. The absolute values determined by Taqman qPCR are overestimated in relation to these values, but once normalized, they are broadly similar, even though more values would be required for a rigorous comparison. Because of this similarity and the fact that the assay involves only three species, the results are based on efficiencies measured in \mathcal{M}_U .

The inferred median K value is 2.36×10^{12} (s.d. = 2.61×10^{12}). The choice of the value of K has a negligible influence on the inferred values of Λ_s .

Table 3 shows the proportions in the \mathcal{M}_T and \mathcal{M}_G communities, as well as the errors compared to the theoretical proportions and the biodiversity indices.

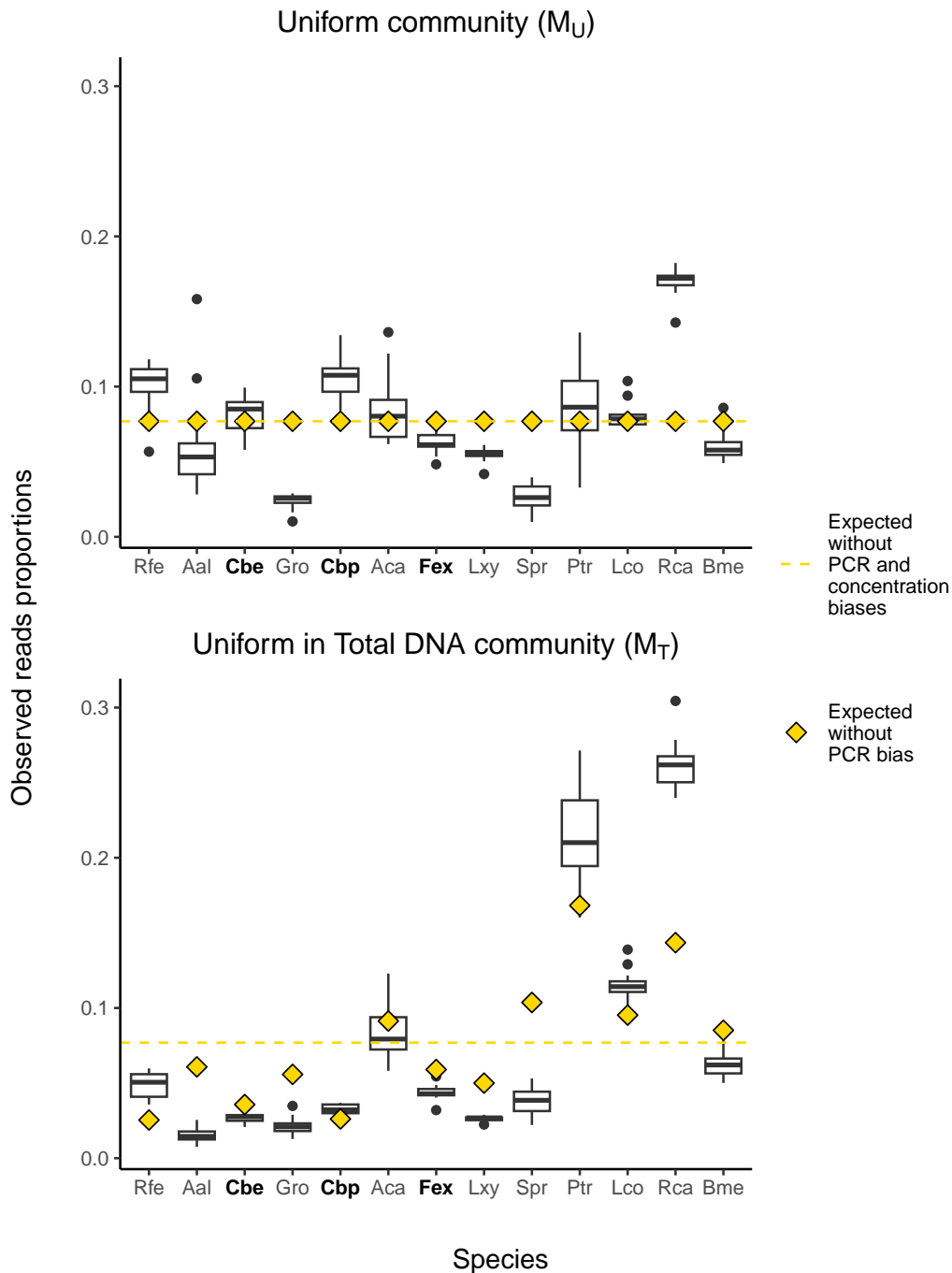


Fig. 4: Observed relative proportions of reads of thirteen plant species for the mock communities M_U and M_T . Gold lines indicate proportions expected in the absence of target concentration and amplification bias. Gold diamonds are the proportions expected in the absence of amplification bias. For the M_U community, the deviation of the boxplots from the diamonds shows the amplification bias alone. For the M_T community, both biases are present. Concentration bias is visible as the difference between the diamond and the line.

The standard deviation of inferred proportions varies from 0.28% (Lco) to 2.0% (Gro) of the mean inferred proportion for M_T (resp. from 0.11% (Lco) 459 460

Species	Average proportion in \mathcal{M}_U (%)		PCR Efficiency inferred from	
	Theoretical	Observed	Taqman	\mathcal{M}_U
Bme	7.7	6.1		0.918
Rca	7.7	17		1.00
Lco	7.7	8.0		0.939
Ptr	7.7	8.6		0.945
Spr	7.7	2.7		0.855
Lxy	7.7	5.5		0.910
Fex	7.7	6.3	0.920	0.920
Aca	7.7	8.3		0.942
Cbp	7.7	11	0.968	0.961
Gro	7.7	2.4		0.847
Cbe	7.7	8.1	0.951 (CbeA) 0.927 (CbeB)	0.940
Aal	7.7	5.8		0.914
Rfe	7.7	10		0.958

Table 2: Proportions in \mathcal{M}_U and relative PCR amplification efficiencies measured for the four Taqman qPCR probes and inferred from the \mathcal{M}_U community. The maximum efficiency was set at 1 for *Rosa canina*. Efficiencies inferred were normalized so that Fex has the same efficiencies with both methods.

to 7.0% (Aal) of the mean inferred proportion for \mathcal{M}_G). The results of the two corrections are comparable and both improve the RMSE criteria, as expected. The corrected biodiversity indices also seem to better approximate the real biodiversity than the observed values.

461
462
463
464

Species	Average proportion in \mathcal{M}_T (%)				Average proportion in \mathcal{M}_G (%)			
	Theoretical	Observed	Inferred with \mathcal{M}_U	Inferred with Λ_s	Theoretical	Observed	Inferred with \mathcal{M}_U	Inferred with Λ_s
Bme	8.5	6.2	8.4	8.3	50	36	54	53
Rca	14	26	12	12	25	40	20	20
Lco	9.5	11	11	12	13	15	16	16
Ptr	17	21	20	20	6.3	5.2	5.5	5.6
Spr	10	3.9	12	11	3.1	0.63	2.0	2.1
Lxy	5.0	2.7	3.8	3.9	1.6	0.94	1.5	1.5
Fex	5.9	4.3	5.7	5.7	0.78	0.68	0.97	0.95
Aca	9.1	7.9	7.8	8.2	0.39	0.16	0.17	0.17
Cbp	2.6	3.2	2.4	2.5	0.20	0.19	0.15	0.16
Gro	5.6	2.1	6.5	7.1	0.098	0.019	0.064	0.085
Cbe	3.6	2.8	2.6	2.7	0.049	0.030	0.031	0.034
Aal	6.1	1.5	2.3	2.1	0.024	0.0045	0.0072	0.0049
Rfe	2.5	5.1	3.7	3.9	0.012	0.014	0.012	0.014
AbsErr		0.045	0.017	0.018		0.057	0.020	0.019
RelErr		0.53	0.26	0.28		0.50	0.34	0.33
1D	11	8.9	11	11	4.0	3.7	3.7	3.8
2D	10	6.7	9.1	9.1	3.0	3.1	2.8	2.8

Table 3: Proportions of species in \mathcal{M}_T and \mathcal{M}_G . Inferred with \mathcal{M}_U means corrected by the ratios. Proportions inferred with Λ_s are obtained by fitting the PCR model using the efficiencies inferred previously.

PCR bias importance: comparison of model simulations and observed data

To illustrate the effect of small differences in efficiency, PCR kinetics were simulated for two species with equal initial quantities. Figure 5 shows the final proportions of the two species according to the difference in PCR efficiency. These simulations are compared with the proportions observed in the \mathcal{M}_U community when comparing *Rosa canina* (the most efficiently amplified species) and the other species individually. These two proportion series are very close to each other.

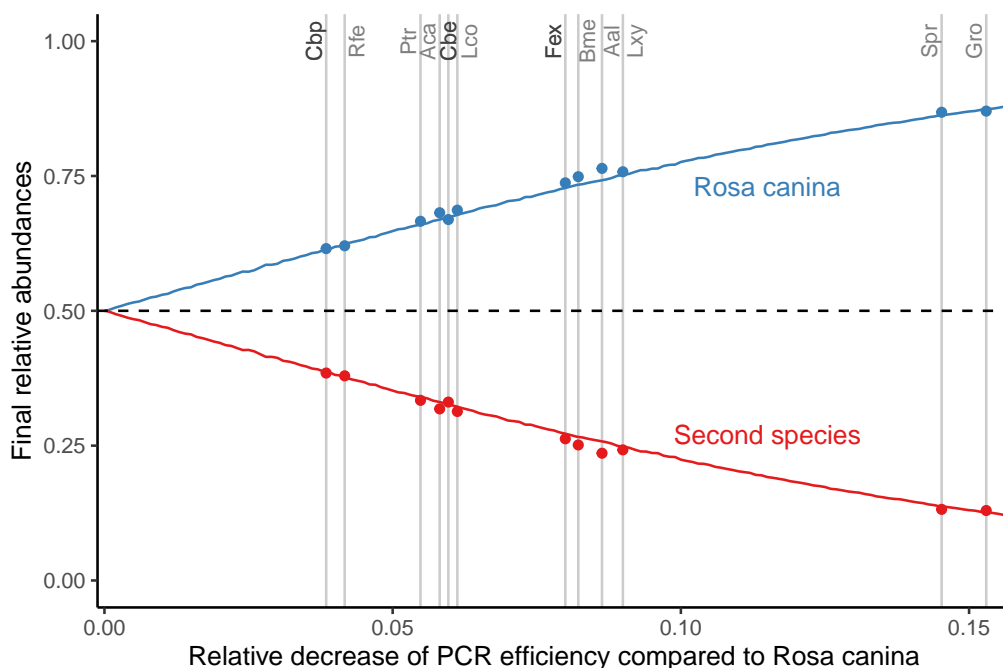


Fig. 5: Relative abundances in a mock community of two initially evenly distributed species simulated with the logistic model (lines) and observed in the \mathcal{M}_U community (dots) considering only *Rosa canina* and the other species individually. The first species has an efficiency of $\Lambda_1 = 1$. The second has a variable efficiency, of value $\Lambda_2 = \Lambda_1(1 - x)$ along the x -axis ($\Lambda_2 \in [0.85, 1.0]$).

Discussion

The quantitative aspect of DNA metabarcoding is regularly questioned by ecologists. Here, two potential biases were considered and their relative effects quantified.

The first is well known. It has long been discussed by microbial ecologists (Kembel, Wu, Eisen, & Green, 2012; Milivojević et al., 2021) and has been identified for macro-organisms (Garrido-Sanz et al., 2022; Krehenwinkel et al., 2017). It can be summarized by a simple question: how many copies of the target gene marker are present per genome in each species under consideration?

In macro-organisms such as plants and animals, most of the targeted markers are carried by the chloroplast or mitochondrial genome, but the same question remains: how many copies of the organelle genome are there per cell? This amount can be estimated by ddPCR. In this study, the communities studied were composed of purified DNA, assayed by ddPCR in marker copies per ng of total DNA. In the case of an environmental DNA sample, this measurement is easily adapted by taking into account the C-value (equation 3) to assay the target in marker copies per genome. Another option is to measure the target in copies per unit of biomass, as in [Thomas et al. \(2016\)](#) or ([Matesanz et al., 2019](#)) (equation 2) Such a reference database makes it possible to convert molecule proportions into species biomass proportions. Here is an example of the bias in marker copies per genome. Among the 13 plants tested, the one more concentrated in chloroplast DNA, *Populus tremula* (Ptr), has 6.6 times more copies per unit of nuclear DNA than the one less concentrated, *Rhododendron ferrugineum* (Rfe). According to the Kew C-value database, the 1C value of Ptr is 0.45 pg ([Siljak-Yakovlev et al., 2010](#)) and that of *Rhododendron ponticum*, the only *Rhododendron* measured, is 0.74 pg ([Bou Dagher-Kharrat et al., 2013](#)) (See Supplementary Table 1). Both together allow to estimate that the bias in chloroplast abundance (in copies per genome) can lead to a 4-fold overestimation of Ptr abundances relative to Rfe (equation 2).

The second type of bias is amplification bias, which our study quantifies precisely. The amplification efficiency of a marker for the species s (Λ_s) is an intrinsic property of the sequence. It does not depend on co-amplified sequences. In this study, we propose two methods to measure it. Both provide similar values, and the choice between them depends on practical convenience. The values obtained can be used to correct the composition of any community, as long as differences in amplifiability between the species present do not cause one or more to disappear. The proposed correction method combines the generation of a reference base for the amplifiability and a mathematical model of the PCR. It does not require any modification of the metabarcoding protocol. Therefore, it can be applied to already generated results and is easy to implement. In particular, experimental parameters do not need to be adjusted to keep amplification within the exponential regime. The logistic model is valid even if saturation is not reached, *i.e.* if $\sum_s M_n^s < K$. However, it is necessary to achieve the saturation plateau while amplifying the reference community (\mathcal{M}_U in this case) in order to infer the efficiencies of Λ_s and K . Without obtaining the saturation plateau, one cannot determine the stage at which the PCR was terminated. Estimating K from Taqman qPCR data is imprecise, and numerical inference along with Λ_s seems more appropriate. The values of Λ_s are relative and vary with K linearly: this is why a precise value of K is not necessary to establish a correction.. This second method suggests that a precise K value is not necessary to determine the Λ_s . The order of magnitude is the same for both approaches and is comparable to the estimate of [Newton and Graham \(2000\)](#) (1.81×10^{12} molecules at the end of the PCR).

The amplification bias is accumulated over each PCR cycle. Thus, the final bias on the observed read relative frequencies is a function of the amplifiability per cycle and the number of amplification cycles. In PCR, the actual number of amplification cycles is not necessarily the number of cycles programmed into the PCR instrument. This number may be lower because the total amount of DNA that can be synthesized is limited by the nucleotide concentration. It is therefore possible that the plateau will be reached before the programmed number of cycles has been reached, with the last cycles not corresponding to any amplification (Figure 1). Correcting for bias using the ratio method (*e.g.* Shelton et al., 2022; Silverman et al., 2021) requires that each sample, including the reference mock community used to estimate it, be amplified with the same effective number of PCR cycles n_{eff} . For the exponential model, which is an assumption of the ratio method (Luo et al., 2022), the change in n_{eff} can lead to significant variations in the inferred proportions. This means that each sample must contain the same total number of target DNA molecules at the start of the PCR. In our study, each mock community was prepared with close total amounts of target DNA, thus respecting the ideal condition for using the ratio method. Therefore, as shown in Table 3, the corrections made by the ratio method and our PCR model-based approach are strictly equivalent. When samples contain different amounts of target DNA, the efficiency of the ratio method decreases because the number of effective PCR cycles varies from sample to sample. To illustrate, using our logistic PCR models, DNA metabarcoding experiments were simulated for two uniform communities containing 10 times less or 10 times more initial DNA molecules than \mathcal{M}_U . The maximum final ratio observed from simulation (Rca to Gro) is 9.5 in the former case (abundances of 19% versus 2.0%) versus 5.4 in the latter (abundances of 16% versus 2.8%), while it was 6.7 in the original \mathcal{M}_U (abundances of 17% versus 2.6%). Fortunately, oppositely to ratio method, our PCR model-based correction allows us to estimate the effective number of PCR cycles for each sample, thereby accounting for sample heterogeneity, making our correction method more robust.

When the two species with the most different amplifiability, *Rosa canina* ($\Lambda_{Rca} = 1.000$) and *Geranium robertianum* ($\Lambda_{Gro} = 0.847$) are co-amplified, with equal amounts of initial target DNA in the extract, the ratio between the RRA observed after sequencing can be up to 6.7 (Figure 5), leading to a strong overestimation of Rca abundance relative to Gro. This initial assessment shows that due to the exponential nature of PCR, even a small difference in amplifiability, as little as 15% between Rca and Gro, the two extreme species tested, can have as strong an effect on the observed RRA as the bias observed due to chloroplast richness. Sometimes the two biases studied push in the same direction, as in *Populus tremula* (Ptr), which has a high chloroplast concentration and a high amplifiability, or *Capsella bursa-pastoris* (Cbe), which combines both a low chloroplast concentration and a low amplifiability (Fig. 4). Sometimes, by chance, both biases partially compensate, as in *Salvia pratensis* (Spr).

Even if the abundances observed by traditional surveys and those of metabarcoding reads are correlated (Yoccoz et al., 2012), it is necessary to be cautious when analyzing DNA metabarcoding data in terms of quantitative information. If we consider the estimation of biodiversity indices, the worst situation is the estimation of α -diversity. Because of all the biases acting simultaneously on DNA metabarcoding measures, but their good reproducibility, the information they provide is inherently relative. Relative in terms of abundances, DNA metabarcoding can at best provide relative abundances, but also relative because the values provided are biased. Therefore, only changes between measures are truly meaningful. Although it has been shown that α -diversity of plant communities can be correctly estimated from DNA metabarcoding data (Calderón-Sanou, Münkemüller, Boyer, Zinger, & Thuiller, 2020), the limited condition under which this is true, Hill numbers computed for $q = 1$, indicates that this is because at this level of weighting of rare versus abundant species by chance most of the biases are compensated. This phenomenon can also be observed in our results (Table 3), where 1D and 2D values estimated from raw RRA and corrected abundances do not strongly differ, while the error between RRA and theoretical composition decreases by a factor of two when using corrected abundances. This discrepancy between the decrease in error due to the correction and the not so good increase in the quality of the α -diversity estimates can be at least partially explained in \mathcal{M}_G by the abundances of the two most abundant species, *Briza media* (Bme) and *Rosa canina* (Rca), which have inverted abundances when estimated from RRA. For any study analyzing changes in diversity across time or ecological gradients, because metabarcoding measures are biased but accurate, the true β -diversity patterns can be easily detected using metabarcoding. In fact, because the biases are repeatable between measures, they often amplify the pattern because the errors correlate with the ecological signal. The problem of all these biases only arises when trying to disentangle the observed pattern from changes in specific species. Therefore, we can strongly encourage people to be very cautious when interpreting the observed pattern, and to be careful not to over-interpret changes in the abundance of a few species in the community as an ecological cause.

Conclusion

We investigated two of the biases that prevent proper quantification of relative eDNA abundances in metabarcoding data. Despite their importance, these biases are far from being corrected or even considered in most current studies. In this study, we measure the two studied biases and propose a simple method to correct the amplification biases in the limit of extreme cases where some species are so strongly disadvantaged that they disappear from the raw results. The advantage of our method compared to the previous ones is that it is more robust to sample variability, while compared to the spiking-based method it does not require any change in metabarcoding protocols. This also allows the reanalysis of previously obtained results, providing the opportunity

for a better ecological interpretation of them. By combining relative abundance
correction and ddPCR to estimate the amount of target DNA in each sample,
we can even consider the possibility of having access to an absolute quantifi-
cation of DNA in the analyzed DNA extracts for each species instead of only
relative abundances. This opens the possibility to increase the robustness of
the quantitative interpretation of DNA metabarcoding results, although other
biases still need to be assessed and modeled in a similar way to fully achieve
the goal of truly quantitative metabarcoding.

Acknowledgments

The authors thank Christian Miquel for the logistic support, Frédéric Boyer
and Clément Lionnet for their help with the bioinformatics pipeline. This work
was supported by the Alpalga project (ANR-20-CE02-0020).

Data Accessibility and Benefit Sharing statement

Data Accessibility statement

The data and analysis scripts are available on the project's git page, <https://github.com/LECA-MALBIO/metabar-bias>.

Competing interests

The authors declare no competing interests.

Authors' contributions

SM, EC, CG and DP studied the PCR models. SM, EC, CG and PT
designed the associated experimental protocol. SM, EC, CG and PT wrote
the manuscript. DP contributed to the writing of the manuscript. EC and PT
sampled the plants. DR and SM performed the extractions and metabarcod-
ing PCRs. FL and SM performed the qPCR and ddPCR assays. SM wrote the
analysis script. EC and CG supervised the project.

2.4 Résultats complémentaires

2.4.1 Modélisation mécanistique de PCR

Nous avons développé un nouveau modèle mécanistique de PCR en même temps que nous établissions le protocole expérimental présenté dans le manuscrit pour pallier les approximations du modèle logistique. L'auteur principal de ce modèle est Didier Piau mais sa conception a émergé de discussions d'équipe. J'ai ensuite manipulé numériquement le modèle.

Dans ce modèle, la saturation est causée par la diminution de la quantité de nucléotides et d'amorces au cours de l'amplification. L'hybridation entre brins d'ADN n'est pas prise en compte. L'intérêt de ce modèle est de reproduire fidèlement la cinétique de PCR avec un nombre limité de paramètres. Il est à comparer aux modèles présentés en section 0.5.

2.4.1.1 Description du modèle

Au cours de chaque cycle n , on note M_n le nombre de molécules d'ADN simple-brin, D_n les molécules d'ADN double-brin, $P_n = P_0$ la polymérase considérée en excès et constante au cours de l'amplification, N_n les dNTP et A_n les amorces. On ne tient pas compte du fait que les molécules d'ADN et les amorces ont une version en sens direct et une version complémentaire. Les étapes modélisées sont les suivantes :

Hybridation L'hybridation forme des complexes AM et AMP selon les réactions 2.1 et 2.2. Pour simplifier, on calcule l'équilibre pour la formation des complexes AM puis des complexes AMP .



Complexes AM Le nombre de complexes AM formés à cette étape est noté C_n . À l'équilibre, les quantités de molécules sont données par :

$$\begin{aligned} k_1 AM &= k_{-1} [AM] \\ \text{soit } k_1 (A_n - C_n)(M_n - [AM]_n) &= k_{-1} C_n \end{aligned} \quad (2.3)$$

que l'on résout de sorte à avoir $[AM]_n < \min(A_n, M_n)$:

$$\begin{aligned}
C_n &= \frac{A_n + M_n + \frac{k_{-1}}{k_1} - \sqrt{(A_n + M_n + \frac{k_{-1}}{k_1})^2 - 4A_nM_n}}{2} \\
&= \frac{2A_nM_n}{A_n + M_n + \frac{k_{-1}}{k_1} + \sqrt{(A_n + M_n + \frac{k_{-1}}{k_1})^2 - 4A_nM_n}}
\end{aligned} \tag{2.4}$$

Complexes AMP À l'équilibre, on a (la quantité de polymérase étant constante) :

$$\begin{aligned}
k_2[AM]P &= k_{-2}[AMP] \\
\text{soit } k_2(C_n - [AMP]_n)P_0 &= k_{-2}[AMP]_n
\end{aligned} \tag{2.5}$$

et donc

$$\begin{aligned}
[AMP]_n &= \frac{k_2P_0}{k_2P_0 + k_{-2}}C_n \\
\text{et } [AM]_n &= C_n - [AMP]_n
\end{aligned} \tag{2.6}$$

Élongation Ensuite, l'élongation a lieu pendant un temps T selon la réaction 2.7. l est le nombre de bases d'une molécule d'ADN.



qu'on modélise de manière approchée par

$$\begin{aligned}
D_n &= (1 - e^{-k_3T/l})[AMP]_n \frac{N_n}{N_0} \\
&= (1 - e^{-k_3T/l}) \frac{k_2P_0}{k_2P_0 + k_{-2}} C_n \frac{N_0 - l(M_n - M_0)}{N_0} \\
\text{avec } N_n &= N_0 - l(M_n - M_0)
\end{aligned} \tag{2.8}$$

Dénaturation Enfin, la dénaturation est irréversible et complète (réaction 2.9).



À la fin de la dénaturation, on a :

$$\begin{aligned}
M_{n+1} &= M_n + D_n \\
(\text{et } A_{n+1} &= A_n - D_n)
\end{aligned} \tag{2.10}$$

Cycle complet En prenant en compte les différents résultats ainsi que $A_n + M_n = A_0 + M_0 =: K$, on obtient :

$$\begin{aligned}
M_{n+1} &= M_n + (1 - e^{-k_3/T}) \frac{k_2 P_0}{k_2 P_0 + k_{-2}} C_n \frac{N_0 - l(M_n - M_0)}{N_0} \\
&= M_n \left(1 + 2r \frac{A_0 + M_0 - M_n}{A_0 + M_0 + \frac{k_{-1}}{k_1} + \sqrt{(A_0 + M_0 + \frac{k_{-1}}{k_1})^2 - 4(A_0 + M_0 - M_n)M_n}} \right. \\
&\quad \left. \left(\frac{N_0 + lM_0}{lN_0} - \frac{M_n}{K} \right) \right) \\
&= M_n \left(1 + 2rc \frac{K - M_n}{1 + \sqrt{1 - 4c^2(K - M_n)M_n}} \left(e - \frac{M_n}{K} \right) \right) \tag{2.11}
\end{aligned}$$

$$\text{avec } r = (1 - e^{-k_3 T/L}) \frac{P_0}{P_0 + k_{-2}/k_2} \frac{lK}{N_0} \tag{2.12}$$

$$c = \frac{K}{K + k_{-1}/k_1} \tag{2.13}$$

$$e = \frac{N_0 + lM_0}{lK} \tag{2.14}$$

et en introduisant $\gamma_c(x) = \frac{2}{1 + \sqrt{1 - 4c^2(1-x)x}}$ pour $x \in]0, 1[$, il vient :

$$M_{n+1} = M_n \left(1 + rc\gamma_c \left(\frac{M_n}{K} \right) \left(1 - \frac{M_n}{K} \right) \left(e - \frac{M_n}{K} \right) \right) \tag{2.15}$$

En pratique, $K = A_0 + M_0 \simeq A_0$ car la quantité de molécules d'ADN en fin de PCR est très supérieure à la quantité initiale. De même, $e \simeq \frac{N_0}{lK}$. La dépendance en la quantité d'ADN initiale est donc négligeable. Au début de PCR, le rendement caractéristique en phase exponentielle est $\Lambda = r c e$. Le facteur γ_c est inclus entre 1 et 2 : c'est une légère accélération de la cinétique par rapport au modèle logistique.

Modèle aléatoire Nous avons adapté ce modèle en un modèle aléatoire de la même manière que le modèle logistique du manuscrit. La variabilité inter-réplicats et le séquençage sont aussi pris en compte, comme expliqué dans le manuscrit.

– **Variabilité inter-réplicats.** On note m_0 la quantité moyenne de molécules parmi l'ensemble des réplicats. On a deux choix possibles, la variabilité initiale peut être représentée par une loi de Poisson (plus simple) ou par une loi binomiale négative de plus grande variance. Celle-ci est plus réaliste mais nécessite d'estimer un paramètre de dispersion supplémentaire.

$$\begin{aligned}
M_0 &\sim \text{Poisson}(m_0) \\
\text{ou } M_0 &\sim \text{NegBin}(r = \frac{m_0}{\delta - 1}, p = 1 - \frac{1}{\delta}) \\
\text{et ainsi } \mathbb{E}[M_0] &= m_0 \text{ et } \text{Var}(M_0) = m_0 \text{ ou } \delta m_0, \delta > 1
\end{aligned} \tag{2.16}$$

– **Amplification cycle par cycle.** \mathcal{F}_n est une filtration décrivant l'information du système à la fin du cycle n . On a :

$$M_{n+1} | \mathcal{F}_n = M_n + \text{Binomial} \left(M_n, r c \gamma_c \left(\frac{M_n}{K} \right) \left(1 - \frac{M_n}{K} \right) \left(e - \frac{M_n}{K} \right) \right) \tag{2.17}$$

Cela impose que $r c e \leq 1$.

– **Séquençage** Le nombre de lectures R est un sous-échantillon de facteur d (supposé connu) des amplicons au cycle final n :

$$R | M_n \sim \text{Binomial} \left(K.d, \frac{M_n}{K} \right) \tag{2.18}$$

2.4.1.2 Résultats

La dynamique d'amplification par PCR est observée par qPCR. Des données expérimentales sont comparées aux trois modèles étudiés sur la Figure 2.4.1, similaire à la Figure 0.5.13 en introduction. L'intérêt du modèle mécanistique par rapport au modèle logistique est qu'il prend en compte l'asymétrie par rapport au point d'inflexion de la courbe. En revanche, il faut inférer plus de paramètres : M_0, r, c, e, K . La cinétique résultante est convaincante, notamment pour prendre en compte l'asymétrie de l'amplification. Notre modèle mécanistique n'a pas été inclus dans le manuscrit car les résultats d'inférence avec ce modèle et avec le modèle logistique étaient identiques.

L'hypothèse que les dNTP et les amorces sont tous deux limitants est contestable : dans les protocoles habituels, il semblerait que les amorces soient en excès. La ré-hybridation des brins d'ADN pourrait être incorporée dans le modèle à la place ou en supplément.

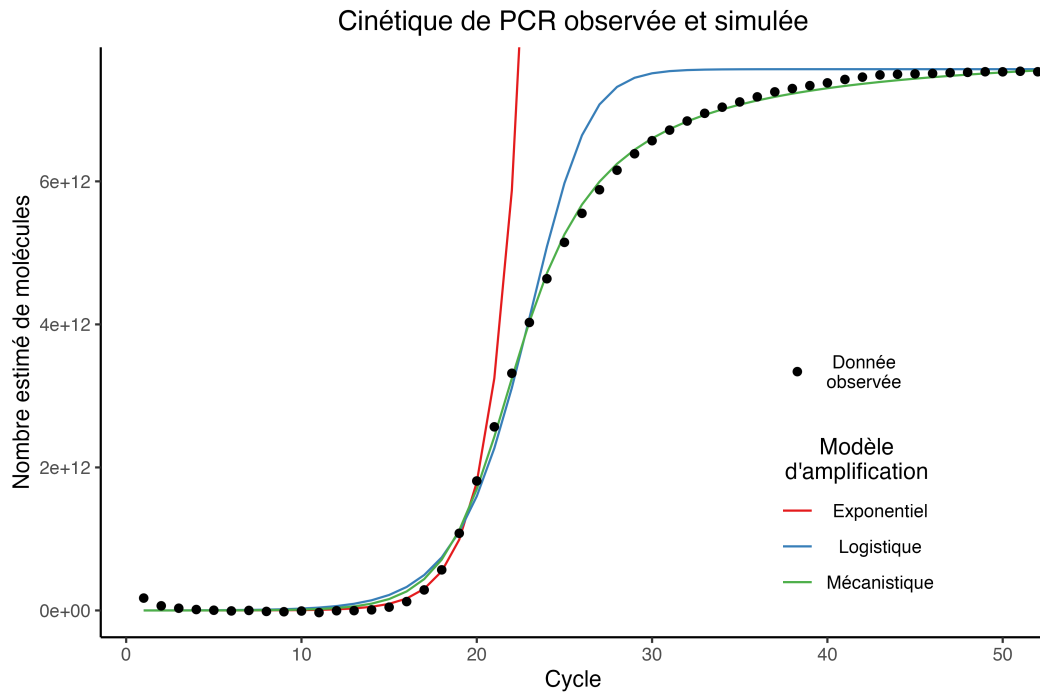


FIGURE 2.4.1 – Cinétique de PCR observée par qPCR pour un échantillon de *Capsella bursa-pastoris* (points) comparée à trois modèles de PCR ajustés. En rouge : modèle exponentiel ; en bleu : modèle logistique ; en vert : modèle mécanistique.

2.4.2 Modélisation de l'amplification par PCR avec mismatch

Une extension possible des modèles de PCR concerne les mismatches d'amorces. Ceux-ci sont connus pour avoir un effet important sur l'amplification (Piñol et al., 2015). En présence de mismatch, deux choix expérimentaux sont possibles : soit l'amplification est faite avec des amorces dégénérées (différentes versions sont incorporées dans la réaction), soit une seule amorce est retenue, compatible avec certaines espèces mais pas parfaitement avec d'autres.

Je propose ici une adaptation des modèles de PCR pour étudier ces deux cas. Le modèle mécanistique est modifié dans sa forme mathématique finale, pas sur le plan biochimique. Pour le modèle exponentiel, il est possible d'explicitier les résultats du nouveau modèle. Pour les autres modèles, je procède par simulation. Dans tous les cas, des résultats sont établis à partir des modèles déterministes associés qui donnent l'espérance des nombres de molécules.

On étudie d'abord une PCR avec une seule espèce de séquence s , comptant M_n molécules au cycle n . Deux amorces existent : p , complémentaire de s , et p' , qui a un mismatch avec s . Les molécules créées à partir d'une amorce p (resp. p') sont de séquence s (resp. s'), quelle que soit la séquence à l'origine de la réplication. On note M'_n le nombre de molécules de séquence s' ($M'_0 = 0$).

2.4.2.1 Cas 1 : Amorces non dégénérées

Dans le premier cas, la seule amorce utilisée est p' . Son mismatch avec la séquence s induit une réduction de son amplification d'un facteur $\rho \in [0, 1]$. Le modèle d'amplification général devient :

$$\begin{aligned} M_{n+1} | \mathcal{F}_n &= M_n = M_0 \\ M'_{n+1} | \mathcal{F}_n &= M'_n + \text{Binomial}(M'_n, \lambda_n) + \text{Binomial}(M_n, \rho \lambda_n) \end{aligned} \quad (2.19)$$

Pour le modèle exponentiel déterministe, l'évolution de M'_n est :

$$M'_{n+1} = M'_n + \Lambda M'_n + \rho \Lambda M_0 \quad (2.20)$$

C'est une suite arithmético-géométrique dont on peut expliciter la forme :

$$M'_n = (1 + \Lambda)^n \rho M_0 + \Lambda \rho M_0 \quad (2.21)$$

Le nombre total de molécules de l'espèce étudiée est donc :

$$M_n + M'_n = ((1 + \Lambda)^n \rho + 1 + \Lambda \rho) M_0 \simeq (1 + \Lambda)^n \rho M_0 \quad (2.22)$$

L'espèce est pénalisée sur sa quantité initiale apparente, pas sur son efficacité d'amplification.

2.4.2.2 Cas 2 : Amorces dégénérées

On considère maintenant un système où les deux amorces p et p' sont présentes, compatibles avec les séquences s et s' , respectivement, et avec un mismatch avec l'autre séquence. On étudie deux modes de saturation : soit les deux amorces s'épuisent indépendamment, soit la saturation est commune, par exemple car les nucléotides ont un caractère plus limitant que les amorces. Dans le premier cas, le milieu contient assez d'amorces pour créer K molécules de s (resp. K' molécules de s'). On suppose $K = K'$. Dans le second cas, il y a une seule capacité de charge K .

On émet l'hypothèse que pendant la phase d'hybridation, une molécule s se lie à une amorce p avec une probabilité π_n^{ps} et une amorce p' avec une probabilité $\pi_n^{p's}$, et de manière analogue pour s' . En quelque sorte, les amorces avec mismatch parasitent les paires complémentaires ps et $p's'$. Dans le cas de deux saturations indépendantes, ces probabilités dépendent de l'affinité de s pour l'amorce p (et de s' pour p'), $\pi > 0.5$, et des quantités d'amorces déjà consommées, $1 - \frac{M_n - M_0}{K}$ et $1 - \frac{M'_n - M'_0}{K'}$:

$$\pi_n^{ps} = \frac{\pi(1 - \frac{M_n - M_0}{K})}{\pi(1 - \frac{M_n - M_0}{K}) + (1 - \pi)(1 - \frac{M'_n - M'_0}{K'})} \quad (2.23)$$

$$\pi_n^{p's} = 1 - \pi_n^{ps} \quad (2.24)$$

Dans le second cas, les pondérations sont simplement $\pi_n^{ps} = \pi_n^{p's'} = \pi$ et $\pi_n^{p's} = \pi_n^{ps'} = 1 - \pi$.

Au cycle n , les nombres de paires amorces-séquences M_n^{ps} , $M_n^{p's}$, $M_n^{ps'}$ et $M_n^{p's'}$ sont tirés selon une loi multinomiale :

$$\begin{aligned} (M_n^{ps}, M_n^{p's}) &\sim \text{Multinomial}(M_n, (\pi_n^{ps}, 1 - \pi_n^{ps})) \\ (M_n^{ps'}, M_n^{p's'}) &\sim \text{Multinomial}(M_n, (1 - \pi_n^{p's'}, \pi_n^{p's'})) \end{aligned} \quad (2.25)$$

Les succès de l'élongation sont ensuite λ_n (paire complémentaire) et $\rho\lambda_n$ (paire avec mismatch). Ici, la saturation est reportée sur la phase d'élongation. La loi d'évolution devient :

$$\begin{aligned} M_{n+1}|\mathcal{F}_n &= M_n + \text{Binomial}(M_n^{ps}, \lambda_n) + \text{Binomial}(M_n^{p's'}, \rho\lambda_n) \\ M'_{n+1}|\mathcal{F}_n &= M'_n + \text{Binomial}(M_n^{p's}, \rho\lambda_n) + \text{Binomial}(M_n^{p's'}, \lambda_n) \end{aligned} \quad (2.26)$$

Si l'on considère le modèle exponentiel déterministe associé (avec $K = K' = +\infty$), on peut écrire :

$$\begin{pmatrix} M_n \\ M'_n \end{pmatrix} = \begin{pmatrix} 1 + \pi\Lambda & (1 - \pi)\rho\Lambda \\ (1 - \pi)\rho\Lambda & 1 + \pi\Lambda \end{pmatrix}^n \begin{pmatrix} M_0 \\ M'_0 \end{pmatrix} \quad (2.27)$$

La matrice de transition est diagonalisable. Les calculs sont présentés en Annexe C. On obtient, avec $M'_0 = 0$:

$$\begin{pmatrix} M_n \\ M'_n \end{pmatrix} = \frac{M_0}{2} \begin{pmatrix} (1 + (\pi + \rho - \pi\rho)\Lambda)^n + (1 + (\pi - \rho + \pi\rho)\Lambda)^n \\ (1 + (\pi + \rho - \pi\rho)\Lambda)^n - (1 + (\pi - \rho + \pi\rho)\Lambda)^n \end{pmatrix} \quad (2.28)$$

La quantité d'intérêt est l'ensemble des molécules de l'espèce en question :

$$M_n + M'_n = (1 + (\pi + \rho - \pi\rho)\Lambda)^n M_0 \quad (2.29)$$

Dans ce cas, l'efficacité d'amplification est affectée par les amorces dégénérées.

2.4.2.3 Résultats de simulation

On considère désormais un mélange de deux espèces et on souhaite quantifier l'effet du choix d'amorce(s) sur les abondances finales, pour vérifier les comportements établis par le calcul.

Les modèles logistique et mécanistique sont étudiés par simulation à partir des équations 2.19 pour le cas 1 et 2.26 pour le cas 2.

Cas 1 : Amorces non dégénérées Plusieurs scénarios sont testés avec différentes quantités initiales et efficacités de PCR des deux espèces en faisant varier le paramètre ρ entre 0 et 1. On observe parfaitement la relation établie par l'équation 2.22 : l'espèce avec un mismatch et une quantité initiale M_0 a une proportion finale exactement équivalente à celle d'une espèce sans mismatch dont la quantité initiale aurait été ρM_0 , avec la même efficacité de PCR.

Cas 2 : Amorces dégénérées Deux amorces sont présentes dans le milieu réactionnel. Chaque espèce est compatible avec l'une des deux et a un mismatch avec l'autre. Des scénarios sont étudiés comme pour le premier cas.

En fonction des paramètres et de la saturation (indépendante ou commune), on observe différents régimes. J'illustre ces résultats pour un scénario où l'espèce 1 est mieux amplifiée que l'espèce 2 : $\Lambda_1 = 0.92 > \Lambda_2 = 0.9$. L'espèce 1 (resp. 2), compatible avec l'amorce p (resp. p'), a pour quantité initiale M_0^1 variable (resp. $M_0^{2'} = 1000$ fixée). La quantité d'intérêt est le ratio d'accroissement défini par l'équation 2.30 qui vaut 1 en l'absence de biais.

$$\text{Ratio d'accroissement} = \frac{\text{Quantité finale espèce 1}}{\text{Quantité initiale espèce 1}} / \frac{\text{Quantité finale espèce 2}}{\text{Quantité initiale espèce 2}} \quad (2.30)$$

La Figure 2.4.2 montre ce ratio dans différents cas. Si la saturation est commune, on retrouve le comportement prédit par le modèle exponentiel (équation 2.29). Si les saturations sont indépendantes pour les deux amorces, ce constat n'est plus vrai, en particulier lorsque π est grand et ρ est petit. Par exemple, pour $\pi = 0.9$ et $M_0^1 = 100$, l'espèce 1 est environ trois fois "trop" abondante pour les petites valeurs de ρ . On peut interpréter ce résultat. Lorsque l'amplification avec mismatch est fortement pénalisée, l'espèce minoritaire (moins abondante au départ ou en cas de quantités initiales égales, avec une moindre efficacité) est sur-représentée à la fin car elle est amplifiée plus longtemps avec un bon rendement, comme le montre la Figure 2.4.3. L'espèce dominante est moins à même de puiser dans les deux stocks de ressources.

2.4.2.4 Discussion des modèles avec mismatch

Cette étude théorique montre l'importance du choix des amorces lors de l'amplification de plusieurs espèces : les comportements prédits sont assez différents en présence d'une ou deux amorces.

D'un point de vue purement mathématique, le premier cas, avec une seule amorce, est plus facile à analyser. En pratique, selon la valeur du paramètre ρ , i.e. selon la pénalisation du mismatch, il se peut que l'abondance finale d'une espèce chute drastiquement au point de ne pas être détectée. Le risque de non-détection n'augmente visiblement pas en présence d'amorces dégénérées, en revanche le biais est difficile à quantifier.

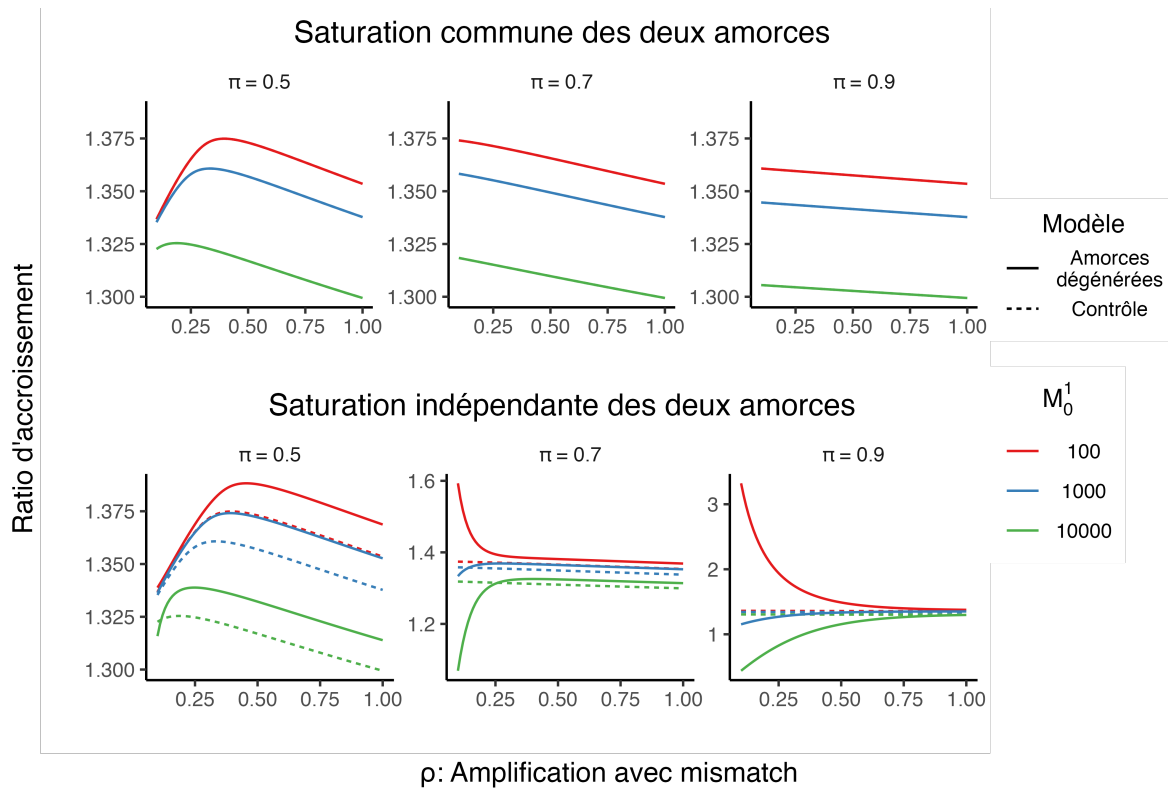


FIGURE 2.4.2 – Ratio d'accroissement pour plusieurs scénarios. En haut, la saturation est commune, en bas, deux saturations ont lieu. Le modèle "Contrôle" correspond au cas de deux espèces sans mismatch avec un rendement de $(\pi + \rho - \pi\rho)\Lambda_1$ et $(\pi + \rho - \pi\rho)\Lambda_2$, selon les calculs établis pour le modèle exponentiel.

La valeur réelle des paramètres liés au mismatch (pénalité ρ et affinité π) n'est pas connue dans cette étude. Des travaux expérimentaux sont donc nécessaires pour déterminer dans quel régime d'amplification est vérifié en pratique et pour valider les résultats des modèles, notamment pour déterminer les mécanismes principaux de saturation en présence d'amorces dégénérées.

Il serait intéressant d'adapter les équations de réaction à l'origine du modèle mécanistique pour vérifier si les approximations de calcul sont correctes. Ensuite, il faudrait répartir la saturation entre la phase d'hybridation (raréfaction des amorces) et d'élongation (raréfaction des dNTP), par exemple en modifiant l'équation 2.26. L'influence des paramètres c , e et K n'a pas été étudiée spécifiquement ici. Enfin, le biais d'amplification en présence de mismatches augmente avec la température de fusion T_m (Sipos et al., 2007) : ce facteur pourrait être pris en compte d'une manière ou d'une autre.

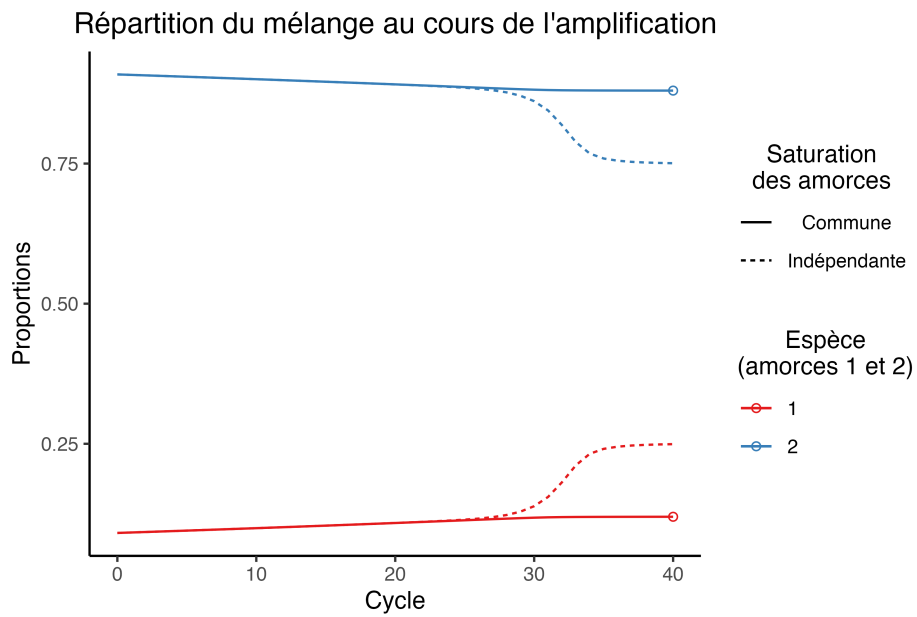


FIGURE 2.4.3 – Proportions des deux espèces au cours de la PCR pour $M_0^1 = 100$, $M_0^2 = 1000$, $\rho = 0.1$, $\pi = 0.9$. Lorsque les saturations sont indépendantes, on observe un changement de régime. Autour du cycle 30, l'espèce 2 atteint sa saturation pour l'amorce compatible p' . Cette saturation n'intervient qu'autour du cycle 35 pour l'espèce 1. Les cercles désignent les proportions attendues avec les efficacités $(\pi + \rho - \pi\rho)\Lambda_s$ et les mêmes quantités initiales dans un cas sans aucun mismatch.

2.4.3 Mesure de la compétition par un contrôle interne

Mes premiers travaux sur le métabarcoding quantitatif reposaient sur un autre protocole expérimental que celui présenté dans le manuscrit. L'objectif était de mesurer le biais d'amplification directement dans le protocole de métabarcoding grâce à une molécule synthétique appelée *Spike*, introduite en différentes quantités dans l'échantillon à analyser. Son rôle était conçu différemment des contrôles internes présentés en 0.6.5 car l'objectif n'est pas d'établir une référence pour une quantification absolue de l'ADN ou la mesure de la variabilité intraspécifique.

2.4.3.1 Hypothèse initiale

Notre hypothèse était que l'ajout d'une espèce compétitrice devait faire varier la compétition interspécifique pendant la PCR. Plus l'échantillon contient de molécules, plus la saturation sera vite atteinte pour l'ensemble des espèces. Cette saturation fige les proportions, qui varient donc avec la quantité de *Spike*. De manière simpliste, on peut considérer que la PCR suit une cinétique exponentielle coupée net quand un certain seuil de molécules est atteint, toutes espèces confondues. Plus il y a de molécules de départ, et moins le nombre de cycles avant cet arrêt est important, et donc moins les écarts dus aux différences d'efficacité n'auront eu le temps de se creuser.

L'objectif de ce projet est de mesurer cet effet avec le marqueur *Sper01* pour les trois communautés artificielles présentées dans l'article. Ensuite, le but est d'inférer les proportions initiales et les efficacités d'amplification des espèces présentes à partir d'un modèle de PCR grâce à l'algorithme *fimo* (Chapitre 1).

2.4.3.2 Design de la molécule *Spike*

Les molécules de *Spike* sont amplifiées par le même marqueur que les plantes étudiées. Sa séquence doit être clairement différente des codes-barres des espèces ciblées. Ce *Spike* (50 paires de bases, 56.0% de GC) a été synthétisé par Sigma-Aldrich. Ces informations sont à comparer à la *Supplementary Table 1* du manuscrit (Annexe B). Sa séquence est :

$$\text{accctcagcctcgcccaaggttgaattatgaaacctgtgacggtcgggtc} \quad (2.31)$$

2.4.3.3 Protocole expérimental

À partir d'un échantillon d'ADNe, quatre compositions sont faites. Chacune contient la même quantité d'échantillon. La première composition ne contient pas de *Spike* et les trois suivantes en contiennent une quantité croissante (Table 2.1). Chaque composition a fait l'objet de vingt réplicats ($2 \mu\text{l}$ d'ADN), un contrôle négatif ($2 \mu\text{l}$ d'eau pure) et trois puits vides. Nous voulions déterminer a posteriori un nombre raisonnable de réplicats pour appliquer ce protocole.

Communauté	Composition	Copies par puits	Part de <i>Spike</i>	Influence du <i>Spike</i>
\mathcal{M}_U	Plantes	2.5×10^5		
	<i>Spike-0</i>	0	0% (0 esp.)	
	<i>Spike-I</i>	1.9×10^4	7.1 % (1 esp.)	-7.1%
	<i>Spike-II</i>	9.5×10^4	28 % (5 esp.)	-28%
	<i>Spike-III</i>	1.9×10^5	43 % (10 esp.)	-43%
\mathcal{M}_T	Plantes	1.1×10^5		
	<i>Spike-0</i>	0	0% (0 esp.)	
	<i>Spike-I</i>	6.3×10^3	5.3 % (0.33 esp.)	-5.3%
	<i>Spike-II</i>	3.2×10^4	22% (1.7 esp.)	-22%
	<i>Spike-III</i>	6.3×10^4	36 % (3.3 esp.)	-36%
\mathcal{M}_G	Plantes	2.5×10^5		
	<i>Spike-0</i>	0	0%, 0 esp.)	
	<i>Spike-I</i>	1.1×10^4	4.2 % (0.088 esp.)	-4.2%
	<i>Spike-II</i>	5.4×10^4	18 % (0.44 esp.)	-18%
	<i>Spike-III</i>	1.1×10^4	30 % (0.88 esp.)	-30%

TABLE 2.1 – Nombre de molécules de *Spike* dans les trois communautés artificielles et les quatre compositions. "Plantes" désigne l'ensemble des copies des espèces de plantes de la communauté artificielle. Les puits ont un volume de $2 \mu l$. "esp." signifie "équivalent de l'espèce la plus abondante". L'influence du *Spike* correspond à la diminution de la proportion relative de chaque espèce par rapport à la composition *Spike-0*.

Comme les autres espèces, le *Spike* a été dosé par ddPCR pour *Sper01* pour des concentrations d'ADN allant de $7.0 \times 10^{-9} ng/\mu l$ et $7.0 \times 10^{-5} ng/\mu l$. Sa concentration mesurée est de 3.5×10^8 copies/ng d'ADN. Tout l'ADN de *Spike* est ciblé par le marqueur, d'où la différence d'ordre de grandeur par rapport aux espèces de plantes (Figure 3 du manuscrit). L'efficacité de PCR du *Spike*, notée Λ_{Spike} , est inconnue.

Pour un échantillon contenant S espèces, le modèle comporte $2S + 1$ paramètres inconnus : efficacité d'amplification et quantité initiale des S espèces présentes, et efficacité d'amplification du *Spike*. Il faut au moins trois compositions car chacune possède S degrés de liberté. En effet, la quantité finale de *Spike* est fixée par la saturation : $R_{\text{Spike}} = \overline{R_{\text{total}}} - \sum_{s=1}^S R_s$, où R_s est le nombre de lectures de l'espèce s .

2.4.3.4 Implémentation

Le protocole est similaire à celui du manuscrit. Nous utilisons ici aussi la méthode *flimo*. Afin d'accélérer les calculs, l'inférence est réalisée en Julia. Les lois du modèle (logistique dans la version finale) sont approchées par des lois normales de même moyenne et variance de manière à bénéficier d'un module de Différentiation Automatique. Cette implémentation fonctionne efficacement, alors même que la fonction à différentier est complexe.

La fonction objectif à minimiser est définie par l'équation 10 du manuscrit, mais

cette fois tous les paramètres, m_0^s et Λ_s , sont inconnus. Il faut résoudre :

$$\operatorname{argmin}_{m_0^1, \dots, m_0^S > 0; 0 \leq \Lambda_1, \dots, \Lambda_S, \Lambda_{\text{Spike}} \leq 1} J((m_0^s)_s, (\Lambda_s)_s, \Lambda_{\text{Spike}}) \quad (2.32)$$

$$\text{avec } J((m_s^{(0)})_s, (\Lambda_s)_s, \Lambda_{\text{Spike}}) = \sum_{\text{composition}} \sum_{c=1}^{S+1} \frac{\left(\overline{p_s(c)}_{\text{données}} - \widehat{p}_s(c) \right)^2}{\overline{p_s(c)}_{\text{données}}} \quad (2.33)$$

Le protocole a été testé avec succès sur des données de PCR simulées pour des échantillons de 2 à 13 espèces avec des rendements et des quantités initiales variables mais avec une variabilité intraspécifique visiblement trop faible.

2.4.3.5 Résultats expérimentaux

Les proportions relatives de chaque espèce sont visibles sur la Figure 2.4.4. Les résultats d'inférence n'ont pas été inclus ici car l'algorithme ne converge pas. On constate deux points importants : la variabilité inter-réplicats est très grande pour certaines espèces et on n'observe pas de différence claire liée à la quantité de *Spike*. Les données ne se prêtent donc pas au protocole envisagé qui tolère mal cette incertitude et la faiblesse du signal que l'on cherche à observer. L'inférence fonctionne dans des cas plus simples, comme dans le manuscrit où les efficacités d'amplification sont connues.

2.4.3.6 Discussion

La faiblesse du signal est due au fait que le *Spike* n'occupe qu'une place relative dans les communautés (au plus 43% attendu dans la communauté \mathcal{M}_U *Spike-III*, Table 2.1). Mais nous voulions éviter d'amplifier et de séquencer majoritairement des molécules de *Spike* car cela n'apporte pas d'information sur le contenu de l'échantillon. De plus, les efficacités d'amplification varient de quelques pour cent seulement. Il est difficile de mesurer cette variation alors qu'elle a un vrai effet sur les abondances finales.

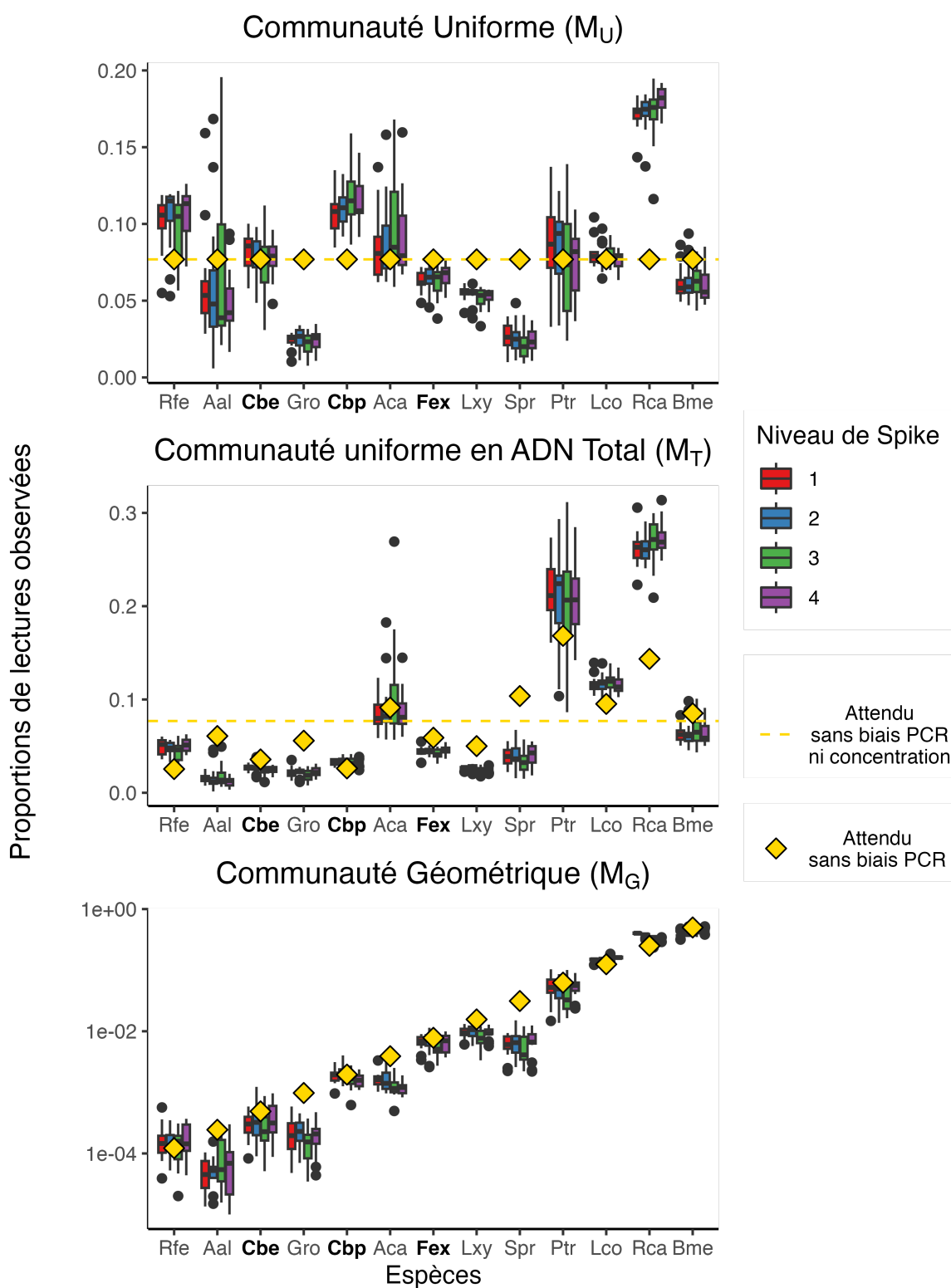


FIGURE 2.4.4 – Proportions relatives de lectures des treize espèces étudiées dans trois communautés artificielles différentes pour le marqueur *Sper01* en excluant le *Spike*. Les losanges jaunes correspondent aux proportions attendues en l'absence de biais d'amplification. Les lignes pointillées vertes correspondent aux proportions attendues en l'absence de biais de concentration et d'amplification.

2.4.4 Mesure de la concentration d'ADN cible selon le tissu chez plusieurs espèces de plantes arctiques

En complément de mes travaux sur le métabarcoding quantitatif, j'ai collaboré avec Stefaniya Kamenova (Université d'Oslo, Norvège) sur un projet visant à mesurer la concentration en ADN cible pour différentes espèces de plantes et différents tissus : racines, feuilles, fleurs et graines. Ces travaux sont complémentaires de la section *Quantification of target DNA* de l'article et confirment le biais de concentration étudié.

2.4.4.1 Ma contribution au projet

Dans ce projet, ma contribution principale a été l'analyse des résultats de ddPCR. J'ai calculé les rapports entre le nombre de copies mesuré par ddPCR et la biomasse sèche de chaque espèce et ai codé les scripts d'analyse. J'ai également participé au design expérimental pour les mesures par ddPCR des différents échantillons (gamme de dilution, nombre de réplicats...) à partir de mes propres travaux avec la ddPCR.

2.5 Conclusion

Dans ce chapitre, j'ai présenté plusieurs travaux visant à améliorer l'interprétation quantitative des données de métabarcoding. Le projet principal a montré comment mesurer deux biais grâce à des techniques de PCR quantitatives et à les corriger en utilisant un modèle de PCR réaliste. Les résultats complémentaires ont porté sur des extensions des modèles de PCR et d'autres approches expérimentales. Ces travaux permettront de tirer des conclusions écologiques plus robustes à partir des données d'ADN environnemental.

Là encore, un bilan global et des perspectives sont exposés dans la Discussion générale du manuscrit.

Chapitre 3

Assignment probabiliste des variants de métabarcoding

3.1 Introduction

Le métabarcoding a pour but de fournir des informations sur la biodiversité à partir de l'ADN environnemental. Mais pour cela, il faut retrouver les espèces présentes à partir d'un grand nombre de séquences observées, les variants. Ce nombre est colossal en comparaison du nombre d'espèces attendues, comme l'illustre la Figure 3.1.1. Dans un réplicat d'un contrôle positif composé de 13 espèces de plantes alpines, 17 043 variants sont observés dont 13 108 singletons, c'est-à-dire avec une seule lecture.

Sans surprise, la variabilité des séquences a une influence sur les mesures de biodiversité. Elle peut être d'origine biologique mais aussi technique ou méthodologique. Avant même d'établir des abondances comme au chapitre précédent, il faut identifier la liste des espèces présentes de manière la plus fiable possible en limitant le nombre de faux positifs et de faux négatifs (Ficetola et al., 2015). Les faux positifs sont limités en regroupant les variants en unités taxonomiques selon des critères de similarité et d'abondance. Les faux négatifs peuvent être des espèces rares difficilement détectables ou des espèces dont les codes-barres sont proches et qui sont agrégées à tort dans une même unité taxonomique.

Dans ce chapitre, mon objectif est d'adapter les mesures de biodiversité à cette incertitude sur l'origine des séquences observées. Pour cela, je m'intéresse aux algorithmes de traitement de ces erreurs et en particulier à l'algorithme *obiclean* des OBITools (Boyer et al., 2016), développés au LECA, pour rendre son approche plus robuste. Avant de présenter mes recherches, je reviens sur les différents types d'erreur qui affectent les données puis sur les méthodes existantes pour constituer des unités taxonomiques.

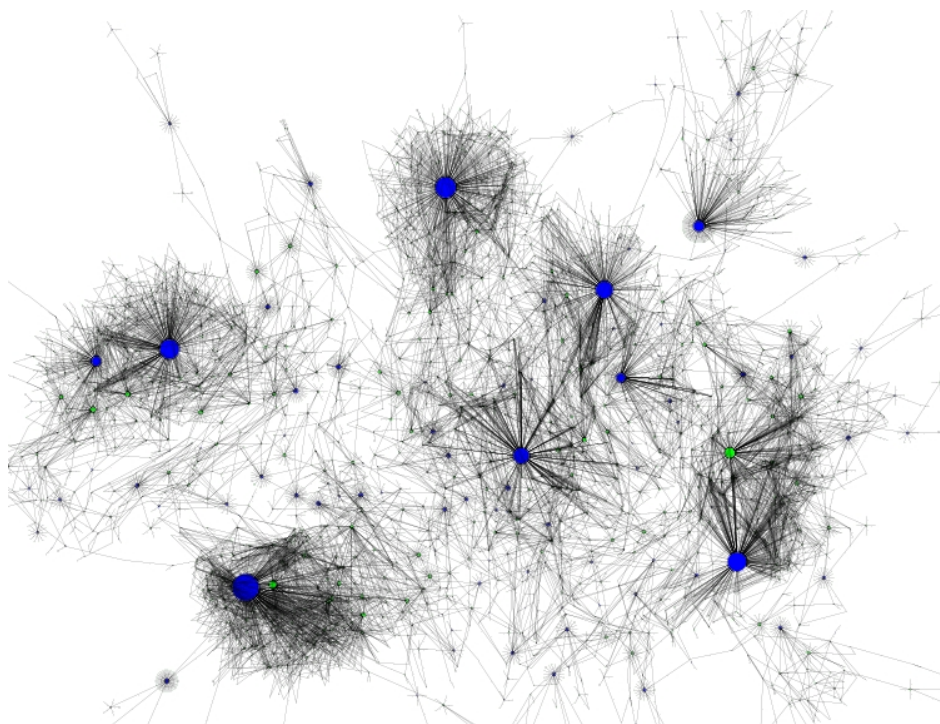


FIGURE 3.1.1 – Graphe des variants d’un réplicat PCR pour le marqueur *Sper03*. Chaque nœud est un variant. Une arête est présente entre deux nœuds à condition qu’une mutation ponctuelle permette de passer de l’un à l’autre. Parmi les variants présents, 336 sont des souches potentielles, 12 374 sont reliés à un variant plus abondant et 4 333 ne sont reliés à aucun autre variant. Graphe établi par *obiclean* et visualisé avec *yEd* (<https://www.yworks.com/products/yed>).

3.1.1 Variabilité des séquences

3.1.1.1 Causes

Il est important de comprendre l’origine de la variabilité pour la prendre en compte de manière pertinente dans les mesures de biodiversité.

La première cause de variabilité est biologique, soit intraspécifique (Estensmo et al., 2021) soit intragénomique (Pereira et al., 2020) : au sein d’une même espèce ou chez un même individu, respectivement, différentes versions du marqueur utilisé cohabitent. Pour la plupart des marqueurs communs en métabarcoding, la similarité des codes-barres au sein d’une espèce est de 96%-99%, parfois moins (Bonin et al., 2023).

La contamination extérieure est une autre forme de variabilité biologique. Elle désigne la détection d’espèces absentes de l’échantillon environnemental mais qui y ont été incluses par erreur lors de la collecte ou en laboratoire. La contamination est détectée à partir de contrôles négatifs (Calderón-Sanou et al., 2020).

Le protocole technique du métabarcoding est la cause de la variabilité dite technique ou méthodologique. Klepke et al. (2022) montre que la richesse avant filtration des séquences augmente avec le nombre de cycles de PCR tandis que la richesse après filtration diminue. La technologie de séquençage utilisée et le protocole de construction des bibliothèques conduit également à des taux d'erreurs différents (Schirmer et al., 2015). La richesse observée augmente avec la profondeur de séquençage, c'est-à-dire le nombre d'amplicons séquencés par réplicats PCR (Shirazi et al., 2021; Klepke et al., 2022). Le traitement des différents réplicats PCR et les choix d'analyse informatique ont aussi des conséquences sur la variabilité observée (Alberdi et al., 2018).

3.1.1.2 Types d'erreur

La variabilité méthodologique est responsable un grand nombre d'artefacts moléculaires, i.e., de variants qui ne correspondent pas à une séquence ayant un sens écologique. Ces artefacts sont de différents types (Zinger et al., 2019; Calderón-Sanou et al., 2020) :

- mutations ponctuelles durant la PCR : insertions, délétions, substitutions ;
- modification de la longueur des homopolymères (bases répétées) (Huntley and Golding, 2006) ou des microsatellites (répétition d'un motif dans la séquence d'ADN) (Leclercq et al., 2010) appelées *slippages*. Ces mutations apparaissent lorsque la polymérase arrête provisoirement de répliquer une région répétée puis reprend alors que le nouveau brin s'est refixé à un motif identique mais décalé ;
- formation de chimères (séquences obtenues par recombinaison de deux séquences existantes) (Schnell et al., 2015; Bjørnsgaard Aas et al., 2017) ;
- présence de séquences d'ADN sans lien avec des séquences connues ;
- présence de contaminants extérieurs ;
- *tag jumps*, c'est-à-dire lectures erronées issues de chimères pendant le séquençage. Les *tags* identifiant les réplicats forment alors une mauvaise combinaison.

3.1.1.3 Abondance des erreurs

L'abondance des erreurs est caractérisée par le nombre de variants différents observés et par l'abondance (relative) des lectures de mutants par rapport aux lectures des séquences connues. Du point de vue de la modélisation, il est commun de les décrire par leur taux d'apparition. Ce taux peut être donné en nombre de mutations par base et par réplication (ou par cycle), ou bien en nombre de mutations par séquence et par réplication pour une représentation plus globale.

Ces abondances varient beaucoup d'une expérience à l'autre et il est difficile d'en donner une estimation précise. En guise d'exemple, j'ai observé dans mes travaux environ 10%-15% de lectures non attribuées à des espèces connues, tandis que les variants inconnus représentent de l'ordre de 99.9% des variants observés. Parmi ces variants, un grand nombre correspond à des singletons, c'est-à-dire des séquences observées une seule fois dans toute l'expérience. Bjørnsgaard Aas et al. (2017) et

Potapov and Ong (2017) estiment qu'entre 30 et 60 % des variants sont des chimères, généralement rares (quelques lectures).

Les taux de mutation varient en fonction des paramètres techniques, des séquences... Il est donc pertinent d'inférer ces taux à partir des données sans supposer de valeur a priori lorsque l'on analyse les données. Je cite tout de même quelques ordres de grandeur, en erreurs par base et par réplication. Potapov and Ong (2017) indique un taux de substitution de 1.5×10^{-4} , un taux de délétion de 1 à 5×10^{-6} et un taux d'insertion de $10^{-7} - 10^{-6}$. Les *slippages* ont plutôt un taux de 10^{-3} (Leclercq et al., 2010) et les recombinaisons créant des chimères de 10^{-5} à 2×10^{-4} (Potapov and Ong, 2017).

Je glisse une remarque sur les corrections classiques des indices de biodiversité. Celles-ci utilisent les espèces rares (observées une ou deux fois) pour estimer le taux de faux négatifs, comme exposé en 0.2.9.3. Dans le cas des données de métabarcoding, cette correction mène évidemment à un résultat aberrant. À l'inverse, le nombre de singletons est extrêmement élevé. Chiu and Chao (2016) utilise un raisonnement analogue aux corrections classiques pour établir des nombres de Hill, avec ou sans rarefaction/extrapolation, et comparer la diversité α de communautés microbiennes. Ils procèdent en estimant le nombre réel d'espèces "avec une seule lecture" à partir des variants observés 2, 3 et 4 fois.

3.1.2 Modèles de mutation

Des modèles de mutation sont utilisés pour décrire ces erreurs dans les algorithmes de traitement des données.

3.1.2.1 Modèles de substitution

Les modèles les plus répandus traitent des seules substitutions car les insertions et délétions sont plus difficiles à représenter dans un seul modèle. Les modèles de substitutions ont d'abord été développés pour construire les arbres phylogénétiques et pour fournir une estimation de la durée d'évolution ou une datation des événements de spéciation.

Ces modèles, markoviens, décrivent la probabilité de substitution des bases d'ADN au cours du temps (continu) sous la forme d'un processus de Poisson. Les hypothèses principales de ces modèles sont la réversibilité en temps, c'est-à-dire que le modèle n'admet pas une séquence ancestrale en particulier dont découleraient les autres séquences ; l'indépendance des bases ; la neutralité des mutations et le nombre fini de sites.

Les modèles sont paramétrés par un vecteur de fréquence des bases à l'équilibre, un taux global μ qui définit le rythme temporel des mutations et 6 paramètres μ_{AC} , μ_{AG} , μ_{AT} , μ_{CG} , μ_{CT} et μ_{TG} . Ceux-ci décrivent le taux instantané de chaque mutation de la première base vers la seconde. La réversibilité en temps induit la symétrie $\mu_{AC} = \mu_{CA}$,

etc. Je cite trois exemples représentatifs des modèles de substitution. Celui de Jukes and Cantor (1969) postule que tous les μ_{ij} sont constants; celui de Kimura (1980) distingue les transitions ($\mu_{AG} = \mu_{CT}$) des transversions ($\mu_{AC} = \mu_{AT} = \mu_{CG} = \mu_{TG}$) et Tavaré (1986) laisse l'ensemble des paramètres libres. La Figure 3.1.2 rappelle les deux catégories de substitution. Je ne présente pas le formalisme précis de ces modèles car le traitement des données de métabarcoding diverge à partir de ce cadre.

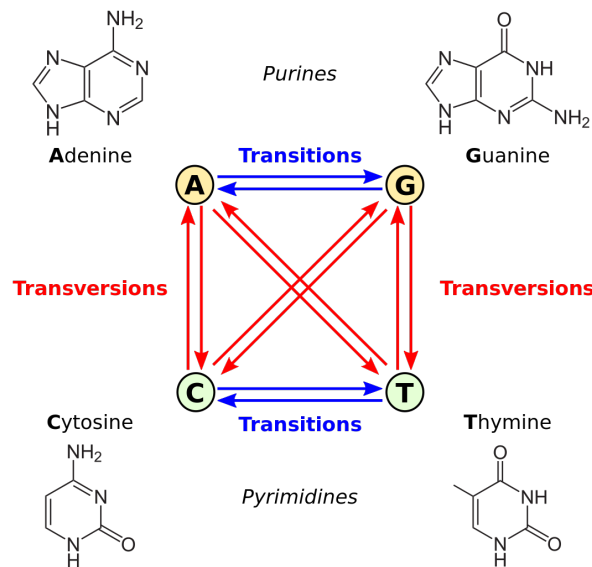


FIGURE 3.1.2 – Deux types de substitution. Les transitions sont plus probables que les transversions. Source de la figure : <https://fr.wikipedia.org/wiki/Transversion>.

Ces modèles de mutation sont faciles à coupler avec les modèles de PCR (Sun, 1995; Wang et al., 2000; Pritchard et al., 2005) où chaque molécule nouvellement créée a une certaine probabilité d'être un mutant qui dépend du type de mutation. Ces modèles de PCR permettent d'ailleurs d'inclure indifféremment tout type de mutation non ponctuelle, car la séquence est considérée dans sa globalité et non base par base.

Pour le métabarcoding, je conserve la notation des taux de mutation μ_{ij} même si le sens mathématique est un peu différent. Les modèles de mutation au cours de la PCR ne vérifient pas la réversibilité en temps car une séquence est facilement identifiable comme la séquence souche par observation des abondances. Cela induit que $\mu_{ij} \neq \mu_{ji}$ ($i \neq j$) a priori. De plus, le temps est discret et seulement 30 à 45 générations permettent les événements de mutation, selon le nombre de cycles de PCR.

Dans les données de métabarcoding, on observe une symétrie des ratios de mutants pour les substitutions (Potapov and Ong, 2017). Il ne s'agit pas des mutations réciproques mais pour les transitions : $\mu_{AG} = \mu_{TC}$ et $\mu_{CT} = \mu_{GA}$; et pour les transversions : $\mu_{AC} = \mu_{TG}$, $\mu_{AT} = \mu_{TA}$, $\mu_{CA} = \mu_{GT}$, $\mu_{CG} = \mu_{GC}$. Le phénomène est facile à expliquer. Je note \bar{i} la base complémentaire de i . Une mutation $i \rightarrow j$ peut avoir deux origines. Soit cette mutation s'est effectivement produite sur le brin direct, soit la mutation $\bar{i} \rightarrow \bar{j}$ a eu lieu sur le brin complémentaire puis a été recopiée lors de la réplication de ce brin complémentaire. On obtient alors un brin direct avec la

mutation $i \rightarrow j$. Cela ne veut pas dire que les deux mutations ont la même probabilité d'apparaître mais que les deux phénomènes sont indiscernables pour ce type de données.

Les insertions et délétions peuvent être traitées en considérant une base "-". $i \rightarrow -$ est donc une délétion d'une base i et $- \rightarrow i$ une insertion d'un i . Ces mutations ont des comportements moins prédictibles mais cela permet d'unifier le traitement des mutations ponctuelles. Considérer tous les types de délétion et insertion ajoute 8 paramètres aux 12 précédents (6 par symétrie).

3.1.2.2 Modèles de mutation non ponctuelle

Des modèles prennent en compte à la fois les mutations ponctuelles et les *slippages* (Leclercq et al., 2010) paramétrés par un taux d'extension (ajout d'un motif) ou de contraction (retrait d'un motif). Les *slippages* ne sont pas considérés dans le projet principal de ce chapitre ; je les évoquerai en 3.5.

3.1.3 Construction d'unités taxonomiques

La détection des erreurs permet de corriger les données pour améliorer les conclusions écologiques. Pour cela, une partie des variants est simplement retirée du jeu de données, comme les singletons, les séquences trop courtes, les séquences éloignées de toute référence connue... Les autres variants sont ensuite regroupés en unités taxonomiques qui rassemblent les séquences issues de la variabilité biologique et méthodologique. On dresse ainsi une liste de séquences biologiquement significatives ou correctes.

Ces unités taxonomiques sont construites à partir des (dis)similarités de séquences mesurées de différentes manières. Le choix de la métrique a une influence sur la composition des unités construites. Parmi les distances d'édition habituelles, on trouve la distance de Levenshtein qui compte le nombre minimal de mutations ponctuelles (substitutions, insertions, délétions) pour passer d'une séquence à une autre et la distance de Hamming prenant en compte seulement les substitutions. Cette dernière est plus rapide à calculer mais ne fonctionne qu'avec des séquences de même longueur. Un autre exemple est la distance de Kimura à deux paramètres (K2P), construite à partir du modèle de mutation du même auteur (Kimura, 1980), qui tient compte de la plus grande fréquence des transitions par rapport aux transversions.

Le premier type de regroupement de séquences est l'Unité Taxonomique Opérationnelle Moléculaire (MOTU), souvent appelée OTU. Une OTU vise à agréger les séquences selon les deux sources de variabilité (biologique et technique). Le critère le plus commun est de regrouper les séquences ayant au moins 97% de similarité. Ce seuil cherche à exploiter le *barcoding gap* représenté sur la Figure 3.1.3, c'est-à-dire le creux entre les distributions de similarité intraspécifique et interspécifique. La pertinence de ce *gap* est toutefois contestée (Wiemers and Fiedler, 2007). Les OTU sont construites par clustering à partir de bases de référence ou de manière non supervisée (*de novo*).

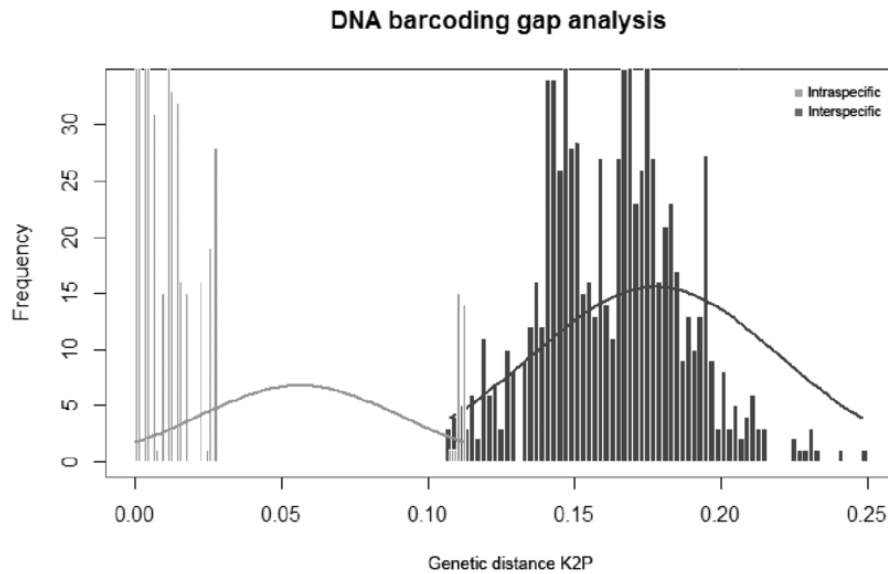


FIGURE 3.1.3 – Distribution de distance entre les codes-barres pour le marqueur COI chez les tipules (insectes). Le *barcoding gap* désigne la zone qui sépare les deux distributions de distance entre les codes-barres selon que la variabilité est intraspécifique ou interspécifique. Cette zone est bien marquée mais on constate une ambiguïté pour une partie de la variabilité intraspécifique. Source de la figure : Rodrigues et al. (2019).

Le seuil arbitraire de similarité est supprimé par la seconde catégorie d'unité appelée *Amplicon Sequence Variant* (ASV) ou encore *Exact Sequence Variant* (ESV), *zero-radius OTU* (zOTU), *sub-OTU* (sOTU)... Les ASV ne cherchent à regrouper que les erreurs de PCR/séquençage autour de leur séquence d'origine mais laissent apparente la variabilité biologique pour conserver l'information la plus complète possible. Ils sont construits à partir de la similarité des séquences et des nombres de lectures grâce à une analyse du graphe de similarité des séquences et/ou de modèles d'erreurs. Les algorithmes de construction des ASV sont dits de *denoising*. Les ASV sont considérés comme des OTU avec un seuil de similarité de 100 %.

3.1.4 Algorithmes de regroupement des séquences

De nombreux algorithmes permettent de construire les unités taxonomiques, avec chacun de multiples combinaisons de paramètres possibles. Plusieurs comparaisons de leurs performances ont été réalisées, par exemple Pauvert et al. (2019); Prodan et al. (2020); Porter and Hajibabaei (2020); Mathon et al. (2021). Ces travaux concluent à une influence des choix d'analyse sur l'interprétation biologique, à la fois en termes d'espèces retrouvées (et de fausses espèces détectées à tort) mais aussi de mesures de biodiversité. Ils attestent qu'il n'y a pas actuellement de consensus sur le meilleur traitement des données. Globalement, les ASV ont une meilleure sensibilité et spécificité que les OTU. Mais dans certains contextes, la variabilité biologique n'est pas informative sur la biodiversité et il est alors pertinent de construire des OTU, parfois à partir des ASV. La communauté débat de l'avantage d'une technique par rapport à une

autre : ASV devant remplacer les OTU (Callahan et al., 2017) ou au contraire les deux approches étant complémentaires (Glassman and Martiny, 2018; Estensmo et al., 2021).

Je présente ensuite les principaux algorithmes utilisés. La répartition entre les différentes catégories est poreuse, par exemple pour l’algorithme dbOTU3 que je classe comme algorithme de *denoising* du fait de sa construction, malgré son nom. Je fais cette remarque pour expliciter le continuum entre tous ces algorithmes et rejeter une opposition conceptuelle franche entre les ASV et les OTU.

3.1.4.1 Algorithmes de clustering

Kopylova et al. (2016) fait une revue des algorithmes de clustering et distingue le clustering *closed-reference*, le clustering *de novo* et le clustering *open-reference*.

Les algorithmes de clustering *closed-reference* établissent les OTU par comparaison à une base de référence, comme SortMeRNA (Kopylova et al., 2012). Porter and Hajibabaei (2020) présente une liste de bases de référence.

À l’inverse, les algorithmes de clustering *de novo* utilisent la similarité des séquences du jeu de données sans incorporer d’information extérieure. Le clustering *de novo* est souvent considéré comme plus performant (Westcott and Schloss, 2015). Parmi eux, on trouve entre autres UCLUST/USEARCH (Edgar, 2010), UPARSE (Edgar, 2013), VSEARCH (Rognes et al., 2016) ou OptiClust (Westcott and Schloss, 2017).

Le clustering *open-reference* est simplement une combinaison des deux approches précédentes (Rideout et al., 2014).

Des logiciels incluent plusieurs algorithmes de clustering pour une utilisation flexible selon le contexte d’étude. C’est le cas de QIIME2 (Bolyen et al., 2019) et mothur (Schloss et al., 2009) (clustering hiérarchique essentiellement).

3.1.4.2 Algorithmes de *denoising*

Les algorithmes de *denoising* sont en général plus sophistiqués que les précédents. Peng and Dorman (2021) dresse une liste des méthodes les plus utilisées.

Modèle probabiliste d’erreurs Une première classe d’algorithmes définit le nombre d’erreurs attendues par séquence à partir d’un modèle prenant en compte, entre autres, l’abondance et la similarité des séquences.

Un des plus connus est DADA2 (Callahan et al., 2016). Il construit un modèle d’erreurs pour chaque valeur de qualité de séquençage possible en utilisant les abondances et la similarité mesurée par l’algorithme de Needleman-Wunsch. AmpliCI (Peng and Dorman, 2021) est similaire mais étudie la qualité base par base et non à l’échelle de la séquence.

Deblur (Amir et al., 2017) et UNOISE3 (Edgar, 2016b) reposent sur des modèles d'erreurs dont les paramètres ne sont pas inférés sur les données mais choisis a priori. Deblur retire à chaque variant le nombre de lectures erronées attendues de ses voisins (déterminés par la distance de Hamming) plus abondants. UNOISE3 construit des clusters pour lesquels chaque variant est accepté ou non selon sa similarité (distance de Levensthein) et son ratio d'abondance par rapport au centroïde du cluster.

DADA2, Deblur et UNOISE3 sont les trois algorithmes les plus utilisés aujourd'hui mais il en existe de nombreux autres. J'en cite deux à titre d'exemples. L'algorithme *Minimum Entropy Decomposition* (MED) (Eren et al., 2015) construit itérativement des groupes de séquences. Si l'entropie de Shannon d'un cluster est trop grande, celui-ci est partitionné. dbOTU3 (Olesen et al., 2017) identifie les séquences similaires et émet l'hypothèse que la somme des séquences d'une OTU et une séquence individuelle qui en fait partie suivent des lois de Poisson de paramètres proportionnels. La classification est établie selon un test du rapport de vraisemblance.

Partition de graphe Une autre classe d'algorithme repose sur la partition d'un graphe de mutations. L'algorithme *swarm2* (Mahé et al., 2015) partitionne le graphe au niveau des séquences les moins abondantes de part et d'autre de pics d'abondances qui sont considérés comme des vraies séquences. L'algorithme *obiclean* des OBITools (Boyer et al., 2016) fait partie de cette catégorie, j'en détaille le fonctionnement dans la section suivante.

3.1.4.3 Algorithme *obiclean*

L'algorithme *obiclean* établit un graphe des mutations ponctuelles et identifie les séquences comme mutantes en fonction du ratio d'abondance entre elles et leurs souches potentielles. Ce graphe de mutations est l'élément de départ des travaux de ce chapitre, j'en donne donc les détails.

Une arête est présente dans le graphe si une seule mutation ponctuelle permet de passer d'une séquence à l'autre. Chaque arête est orientée de la plus abondante à la moins abondante, i.e., de la souche présumée vers le mutant. Ensuite, le ratio de lectures mutant/souche est étudié pour chaque mutant. Si ce ratio est inférieur à un certain seuil (1 par défaut), le nombre de lectures du mutant est ajouté à celui de la souche pour définir une nouvelle variable appelée poids (*obiclean-weight*). Il est alors possible de ne conserver que les racines du graphe, appelées *heads*. Les nœuds internes du graphe sont dit *internal* tandis que les variants sans aucune connexion sont appelés *singletons* (à ne pas confondre avec les variants observés une seule fois).

La quatrième version des OBITools¹ développée par Eric Coissac est implémentée dans le langage Go. J'ai principalement codé en R pour faciliter l'exploration des données. La motivation de mon projet est de modifier l'algorithme d'*obiclean* de sorte

1. <https://git.metabarcoding.org/obitools/obitools4>

à remplacer le seuil fixe appliqué au ratio par un critère donnant la probabilité pour chaque variant d'être un mutant issu d'un autre variant.

3.1.5 Détection des chimères

Certains algorithmes sont conçus pour détecter les chimères, comme UCHIME2 (Edgar, 2016a). Je n'en fais pas l'inventaire ici.

3.1.6 Algorithmes de traitement post-clustering

Comme évoqué plus haut, il peut être intéressant de former des OTU à partir des ASV pour améliorer l'interprétation écologique (Estensmo et al., 2021).

Des algorithmes interviennent une fois les ASV/OTU construits pour détecter de potentiels groupes vraisemblablement erronés, comme LULU (Frøslev et al., 2017). LULU considère les cooccurrences des ASV/OTU avec une similarité suffisante. Si un cluster est systématiquement moins abondant qu'un cluster voisin, il est traité comme une erreur du cluster abondant. Dans ce chapitre, j'étudie également les cooccurrences des séquences pour détecter des erreurs.

3.1.7 Plan du chapitre

J'ai illustré en introduction de ce chapitre les enjeux autour des erreurs de séquence dans les données de métabarcoding. Ce chapitre traite principalement deux questions.

La première concerne l'origine des erreurs observées. Des algorithmes de *denoising* construisent leurs modèles de correction principalement autour de l'étape de séquençage, comme DADA2. Peut-on négliger l'influence de la PCR dans les modèles de correction ? Pour traiter ce point, j'ai effectué un travail de modélisation mis en relation avec des résultats observés dans la littérature.

Ensuite, le cœur de l'étude consiste à définir une matrice de probabilité décrivant les origines possibles de chaque variant sans attribuer une lecture à une classe ou une autre avec un seuil arbitraire. Pour cela, je suis parti du graphe de mutations fourni par *obiclean*. J'ai construit un modèle de mutation ponctuelle et étudié les cooccurrences des variants pour distinguer les mutants vraisemblables. Mon but est d'établir le modèle le plus simple possible qui puisse représenter la diversité des erreurs observées. Cette étude m'a ensuite permis de redéfinir les indices de biodiversité dans le cadre de cette attribution probabiliste des variants aux différentes unités taxonomiques.

Pour conclure, je décrirai une collaboration réalisée avec Frédéric Boyer (LECA) sur le génotypage à partir de données de métabarcoding.

3.2 Un premier modèle de mutation pendant la PCR

Je présente d'abord un modèle qui illustre le rôle de la PCR dans la création de mutants. Ce modèle déterministe est inspiré de celui de Sun (1995) et Wang et al. (2000) qui est un processus de branchement de type Galton-Watson. L'objectif est de montrer que le taux de mutation dépend du nombre de cycles d'amplification. Conceptuellement, ce résultat est contradictoire avec l'approche de DADA2 qui construit une matrice de paramètres pour chacun des 40 scores de qualité de séquençage sans prendre en compte les paramètres de la PCR.

3.2.1 Développement du modèle

Le modèle décrit les molécules d'ADN simple brin selon deux catégories : souches et mutantes. Leurs nombres moyens à la fin du cycle k sont notés M_k^s et M_k^m , respectivement. Les hypothèses du modèle sont les suivantes. Initialement, $M_0^s > 0$ molécules de souche sont présentes et $M_0^m = 0$. Puis à chaque cycle :

- chaque molécule est répliquée avec une probabilité Λ (il s'agit du modèle d'amplification exponentielle, équation 18) ;
- chaque molécule nouvellement créée à partir de la souche a une probabilité μ d'être un mutant.

μ est un taux de mutation par séquence et par cycle, supposé faible.

Le nombre total de molécules au cycle k est $M_k = M_k^s + M_k^m = (1 + \Lambda)^k M_0^s$, d'après le modèle exponentiel. Les formules de récurrence pour M_k^m et M_k^s sont :

$$\begin{cases} M_{k+1}^s = (1 + \Lambda(1 - \mu))M_k^s \\ M_{k+1}^m = \Lambda\mu M_k^s + (1 + \Lambda)M_k^m \end{cases} \quad (3.1)$$

3.2.2 Résultats

Le calcul d'une expression simplifiée est donné en Annexe D. On obtient :

$$\begin{pmatrix} M_n^s \\ M_n^m \end{pmatrix} = M_0^s \begin{pmatrix} (1 + \Lambda(1 - \mu))^n & \\ (1 + \Lambda)^n - (1 + \Lambda(1 - \mu))^n & (1 + \Lambda)^n \end{pmatrix} \quad (3.2)$$

La part de séquences mutées après n cycles est :

$$\frac{M_n^m}{M_n} = \frac{(1 + \Lambda)^n - (1 + \Lambda(1 - \mu))^n}{(1 + \Lambda)^n} \quad (3.3)$$

On effectue un développement limité avec $\mu \ll 1$. Il vient

$$\frac{M_n^m}{M_n} = \frac{n\Lambda}{1 + \Lambda}\mu + o(\mu) \simeq \frac{n\Lambda}{1 + \Lambda}\mu \quad (3.4)$$

On peut interpréter ce résultat. Au cycle n , les mutants proviennent de n événements de mutation différents (avec chacun la même abondance $(1 + \Lambda)^{n-1}\Lambda\mu$) : amplification de la séquence souche pendant $k - 1$ cycles, mutation au cycle k puis amplification du mutant pendant $n - k$ cycles.

Ce résultat est cohérent avec Klepke et al. (2022) qui observe bien une dépendance du nombre d'ASV au nombre de cycles de PCR alors que le nombre d'espèces inféré ne dépend pas de ce facteur.

Chaque amplification réelle a un "nombre équivalent de cycles d'amplification exponentielle" n_{exp} qui correspond au nombre de cycles du modèle exponentiel aboutissant au même nombre de molécules finales K (la capacité de charge du milieu d'amplification). On peut réécrire simplement ce modèle :

$$M_0^s(1 + \Lambda)^{n_{exp}} = K \quad (3.5)$$

$$\text{soit } n_{exp} = \frac{\log K - \log M_0^s}{\log(1 + \Lambda)} \quad (3.6)$$

Ainsi, le ratio de mutants dépend de la quantité initiale d'ADN :

$$\frac{M_n^m}{M_n} \simeq \frac{n_{exp}\Lambda}{1 + \Lambda} \mu = \frac{(\log K - \log M_0^s)\Lambda}{(1 + \Lambda) \log(1 + \Lambda)} \mu \quad (3.7)$$

Ce ratio est d'autant plus grand que le nombre de molécules initial est faible. Cette hypothèse est à vérifier expérimentalement. Expliciter la valeur de Λ selon M_0^s , C_t et M_{C_t} par l'équation 9 du manuscrit du Chapitre 2 n'apporte pas d'information supplémentaire car il subsiste toujours un degré de liberté.

3.3 Assignment taxonomique : Matériel et Méthodes

3.3.1 Données étudiées

Dans ce chapitre, deux jeux de données ont été considérés. La Table 3.1 en fournit quelques informations. Le premier est celui produit pour mes travaux sur le métabarcoding quantitatif avec la molécule *Spike* (Chapitre 2, section 2.4.3) pour le marqueur *Sper01*. Dans ces données, 252 répliquats (dont 12 contrôles négatifs) sont répartis en trois communautés artificielles avec chacune quatre abondances de *Spike* différentes.

Le second jeu de données est issu d'un projet du consortium ORCHAMP sur la distribution des algues dans les Alpes françaises (Stewart et al., 2021). Il est constitué de 195 échantillons de sol prélevés sur cinq sites puis analysés avec le marqueur *Euka03*. Les données comportent 706 répliquats (dont 17 contrôles négatifs et 20 contrôles positifs). Ici, le graphe de mutations est construit pour les mutations dont les souches ont au moins 100 lectures attribuées.

Les descriptifs des marqueurs dans Taberlet et al. (2018) sont inclus en Annexe E.

Marqueur	Variants	Part de singletons	Lectures	Part de singletons
<i>Euka03</i>	1 909 504	79%	12 887 189	12%
<i>Sper01</i>	98 670	71 %	1 390 244	0.68%

TABLE 3.1 – Données utilisées pour l'étude d'*obiclean*.

3.3.2 Étude des mutations ponctuelles

Dans cette étude, j'ai d'abord traité les mutations ponctuelles (substitutions, insertions, délétions) puis toutes les autres erreurs en bloc. Les mutations ponctuelles sont plus faciles à détecter grâce à des méthodes d'alignement. Toutes les mutations d'une base i vers une base j (différente) sont regroupées en une seule catégorie de mutation. Les bases possibles sont A, T, G, C et "-", même si les insertions et délétions montrent un comportement moins prévisible que les substitutions. La symétrie évoquée en 3.1.2.1 permet de ne considérer que 6 taux de substitutions et 8 taux d'insertions et délétions. Ces taux de mutation sont caractérisés par le ratio de lectures mutant/souche.

3.3.2.1 Modélisation du ratio de mutants

Chaque séquence est traitée individuellement. Le ratio des lectures du mutant par rapport aux lectures de la souche est modélisé par une combinaison de deux lois que l'on appelle beta-Poisson. D'abord, le ratio théorique d'un mutant m (une séquence en particulier par rapport à une séquence souche donnée) est tiré selon son type de mutation ($i \rightarrow j$) :

$$\mu_m \sim \text{Beta}(a_{ij}, b_{ij}) \quad (3.8)$$

La loi Beta a deux paramètres à déterminer qui définissent l'espérance et la variance de la loi, données par :

$$\mu_{ij} := \mathbb{E}[\mu_m] = \frac{a_{ij}}{a_{ij} + b_{ij}} \quad , \quad \text{Var}(\mu_m) = \frac{a_{ij}b_{ij}}{(a_{ij} + b_{ij})^2(a_{ij} + b_{ij} + 1)} \quad (3.9)$$

On peut ainsi définir le ratio moyen de mutants pour un type donné, μ_{ij} . Une fois que ce ratio théorique de mutants est déterminé, le nombre effectif de lectures du mutant m dépend du nombre de lectures R_s de la souche selon :

$$R_m | \mu_m \sim \text{Poisson}(\mu_m R_s) \quad (3.10)$$

Au total, il faut donc estimer 12 paramètres pour les substitutions et 16 paramètres pour les insertions/délétions donnés par la Table 3.3.2.1. Ce modèle permet de prendre en compte la variance observée du ratio de mutants pour un type de mutation donné. Celle-ci est en général bien plus grande que la variance de la simple loi de Poisson.

Base d'origine	Base de destination				
	A	C	G	T	-
A		$(a, b)_{AC}$	$(a, b)_{AG}$	$(a, b)_{AT}$	$(a, b)_{A-}$
C	$(a, b)_{CA}$		$(a, b)_{CG}$	$(a, b)_{CT}$	$(a, b)_{C-}$
G	$(a, b)_{GT}$	$(a, b)_{CG}$		$(a, b)_{CA}$	$(a, b)_{G-}$
T	$(a, b)_{AT}$	$(a, b)_{AG}$	$(a, b)_{AC}$		$(a, b)_{T-}$
-	$(a, b)_{-A}$	$(a, b)_{-C}$	$(a, b)_{-G}$	$(a, b)_{-T}$	

TABLE 3.2 – Paramètres à inférer pour le modèle de ratio de mutants.

3.3.2.2 Inférence des paramètres a et b

Les paramètres sont inférés pour chaque échantillon (en regroupant les réplicats techniques) à partir des mutations potentielles dont la souche est une racine (*head*) du graphe. Ce choix est fait pour estimer les paramètres sur des données robustes. L'inférence est effectuée pour les réplicats avec moins de 1000 lectures pour *Euka03* et 5000 pour *Sper01*. On considère ainsi 2274 combinaisons type \times échantillon pour *Euka03* et 168 pour *Sper01*.

La procédure d'inférence de a et b est itérative. Pour chaque type et échantillon, on accepte initialement toutes les mutations puis on répète les étapes suivantes :

- i calcul de la moyenne et de la variance du ratio de mutants en prenant en compte les mutants non observés (ratios égaux à 0) ;
- ii calcul de a et b par la méthode des moments ;
- iii pour chaque mutant, calcul de la valeur de μ_m maximisant sa vraisemblance ;
- iv pour chaque mutant, calcul de la log-vraisemblance sachant μ_m ;
- v mise à jour des mutants acceptés selon le test : vraisemblance $\geq 10^{-3}$. Dans l'algorithme d'attribution probabiliste, ce critère n'est pas utilisé pour classifier les variants à assigner.

3.3.2.3 Précisions sur l'implémentation

i) Calcul des moments Pour calculer l'espérance et la variance du ratio de mutants, on prend en compte les ratios valant 0. Pour cela, on comptabilise les mutants observés au moins une fois dans le jeu de données mais pas retrouvés dans l'échantillon considéré.

Une possibilité envisagée est de compter toutes les mutations possibles pour une séquence donnée (par exemple, pour $A \rightarrow T$, compter les bases A). Cela retire la dépendance au jeu de données global mais accentue l'excès de 0 (l'effet *zero-inflated*). Je justifie mon choix de manière pragmatique : l'absence systématique d'un mutant sur une grande base de données est peut-être due à une raison autre que le hasard.

ii) Inférence de a et b La méthode des moments est utilisée pour estimer a et b . L'espérance du ratio du mutant m par rapport à sa souche s est donnée par :

$$\begin{aligned}\mathbb{E}\left[\frac{R_m}{R_s}\right] &= \frac{1}{R_s}\mathbb{E}[\mathbb{E}[R_m|\mu]] \\ &= \frac{1}{R_s}\mathbb{E}[R_s\mu] = \mathbb{E}[\mu] \\ &= \frac{a}{a+b}\end{aligned}\tag{3.11}$$

Elle est donc indépendante de R_s .

La variance du ratio est donnée par :

$$\begin{aligned}Var\left(\frac{R_m}{R_s}\right) &= \mathbb{E}\left[Var\left(\frac{R_m}{R_s}\right)|\mu\right] + Var\left(\mathbb{E}\left[\frac{R_m}{R_s}\right]|\mu\right) \\ &= \frac{\mathbb{E}[\mu]}{R_s} + Var(\mu) \\ &= \frac{1}{R_s}\frac{a}{a+b} + \frac{ab}{(a+b)^2(a+b+1)}\end{aligned}\tag{3.12}$$

La variance dépend donc de R_s . Pour s'affranchir de ce facteur, on calcule le terme correspondant à la seule variance de la loi Beta :

$$V = Var\left(\frac{R_m}{R_s}\right) - \frac{1}{R_s}\mathbb{E}\left[\frac{R_m}{R_s}\right] = \frac{ab}{(a+b)^2(a+b+1)}\tag{3.13}$$

À partir de ces valeurs, on calcule, avec $\mu = \mathbb{E}\left[\frac{R_m}{R_s}\right]$:

$$a = \mu\left(\frac{\mu(1-\mu)}{V} - 1\right)\tag{3.14}$$

$$b = (1-\mu)\left(\frac{\mu(1-\mu)}{V} - 1\right)\tag{3.15}$$

iii) μ maximisant la vraisemblance On note $\theta = (a, b, R_s)$. La vraisemblance totale du modèle pour un nombre de lectures de mutant R_m est :

$$\begin{aligned}p_\theta(R_m; \mu) &= p_\theta^{Poisson}(R_m|\mu) \cdot p_\theta^{Beta}(\mu) \\ &= \frac{(\mu R_s)^{R_m}}{R_m!} e^{-\mu R_s} \cdot \frac{\mu^{a-1}(1-\mu)^{b-1}}{B(a, b)}\end{aligned}\tag{3.16}$$

où B est la fonction Beta.

On cherche μ maximisant la log-vraisemblance $l_\theta = \log p_\theta$.

$$\begin{aligned} l_\theta(\mu|R_m) &= l_\theta(R_m; \mu) - l_\theta(R_m) \\ &= (R_m + a - 1) \log \mu + (b - 1) \log(1 - \mu) - \mu R_s + C \end{aligned} \quad (3.17)$$

où C est un terme qui ne dépend pas de μ . On annule la dérivée :

$$\frac{\partial l_\theta(\mu|R_m)}{\partial \mu} = \frac{R_m + a - 1}{\mu} - \frac{b - 1}{1 - \mu} - R_s = 0 \quad (3.18)$$

Les solutions sont données par les racines d'un polynôme de degré 2 en μ :

$$\mu_{\pm} = \frac{(R_s + R_m + a + b - 2) \pm \sqrt{(R_s + R_m + a + b - 2)^2 - 4R_s(R_m + a - 1)}}{2R_s} \quad (3.19)$$

La seule solution sur $[0, 1]$ est μ_- (observé numériquement).

3.3.3 Cooccurrences des variants

Nous avons adopté une autre approche pour traiter les erreurs qui ne correspondent pas à des mutations ponctuelles. On attend une "corrélation" entre une séquence souche et ses erreurs. Les erreurs ne doivent pas être plus abondantes que la souche et doivent être absentes quand la souche est absente. L'analyse des corrélations nécessite des échantillons indépendants car deux espèces ont a priori un ratio constant dans des réplicats techniques d'un même échantillon. Cela est pertinent pour les données de *Euka03*. Pour *Sper01*, seulement trois échantillons (\mathcal{M}_U , \mathcal{M}_T et \mathcal{M}_G) présentent des compositions différentes (modulo le Spike) : les résultats doivent être interprétés avec plus de précautions.

Pour comparer deux variants, nous avons étudié plusieurs mesures reposant sur l'ensemble des échantillons pour lesquels au moins un des deux variants est présent. En premier lieu, les corrélations du nombre de lectures attribuées, du logarithme du nombre de lectures ($\log(1+\text{lectures})$) et du rang du nombre de lectures (Spearman) sont calculées.

Nous observons aussi l'écart-type et la déviation absolue à la moyenne (MAD) du ratio de lectures entre le mutant potentiel m et sa souche potentielle s parmi les échantillons concernés :

$$\sigma_m = \sigma \left(\frac{R_m}{R_s + \epsilon} \right) \quad (3.20)$$

$$MAD_m = \mathbb{E} \left[\left| \frac{R_m}{R_s + \epsilon} - \mathbb{E} \left[\frac{R_m}{R_s + \epsilon} \right] \right| \right] \quad (3.21)$$

avec $\epsilon = 10^{-3}$

Si m est une erreur de s , σ_m et MAD_m sont censés être faibles. En revanche, pour deux séquences indépendantes, ces valeurs sont supposées grandes, par exemple car en fonction des échantillons, le variant le plus abondant de la paire n'est pas toujours le même.

La mesure de cooccurrence retenue dans la suite de l'étude est l'écart-type du ratio. Ce choix sera justifié dans la section Résultats.

On sélectionne 45 366 (resp. 1 205) variants pour *Euka03* (resp. *Sper01*) considérés comme des vrais mutants car leur vraisemblance est toujours supérieure à 10^{-3} selon le modèle de mutation ponctuelle. Les paires de variants de cette base de référence permettent de construire une distribution empirique de l'écart-type du ratio à laquelle on compare les paires de variants à évaluer. La fonction de répartition empirique de cette distribution est :

$$x \mapsto \mathbb{P}(\sigma_m \leq x) = \frac{\text{Nombre de paires de référence où } \sigma_m \leq x}{\text{Nombre de paires de référence}} \quad (3.22)$$

3.3.3.1 Similarité génétique des variants

On vérifie ensuite la similarité génétique des paires de séquences détectées avec une forte cooccurrence. Cette étape est nécessaire pour s'assurer que deux espèces ne soient pas traitées comme une paire souche-mutant alors qu'elles sont conjointement présentes pour des raisons écologiques.

Pour cela, la distance de Levensthein entre les séquences est calculée. Certaines paires de séquences sont étudiées individuellement pour chercher de potentielles chimères ; cette approche doit être généralisée au terme du projet.

3.3.3.2 Stratégie de comparaison

Il n'est pas raisonnable d'étudier les cooccurrences des $N_v(N_v - 1)/2$ paires de séquences pour un jeu de données à N_v variants. Pour *Euka03*, cela ferait environ 10^{11} paires sans considérer les variants singletons. Il faut donc décider des paires de séquences à considérer. Une option est d'observer la similarité des séquences mais ce critère n'est pas adéquat pour détecter les chimères.

Pour chaque réplicat, on considère le graphe de mutations et les cooccurrences de ses racines (*heads*). Pour celles-ci, les souches potentielles considérées sont les variants avec au moins 100 lectures dans le réplicat, et les mutants présumés ceux avec au moins 10 lectures. On fait cela pour limiter le nombre de combinaisons à tester environ 10^4 pour *Euka03* (quelques centaines pour *Sper01*). Ce choix m'a paru pragmatique pour que le traitement des séquences améliore la qualité des mesures de biodiversité en corrigeant la classification des séquences les plus abondantes.

3.3.4 Matrice de probabilité d'assignation

Les analyses du modèle de substitution et des cooccurrences sont réunies pour établir une matrice de probabilité P . Celle-ci a pour coefficients la probabilité qu'un variant i soit un mutant d'un variant j : $P_{ij} = \mathbb{P}(i \text{ est issu de } j)$.

Cette matrice est de dimension N_v^2 (avec N_v le nombre de variants) et la somme de chaque ligne i , $\sum_j P_{ij}$, vaut 1. Cette matrice est gigantesque mais quasiment tous les coefficients sont nuls car seulement quelques parents potentiels sont considérés.

On utilise une classification bayésienne. La probabilité que le variant i soit de catégorie C_i (détaillée après) sachant l'information contenue dans les données \mathcal{D} (abondance et écriture de chaque séquence pour chaque réplicat) est, d'après la formule de Bayes :

$$p(C_i|\mathcal{D}) = \frac{p(C_i) p(\mathcal{D}|C_i)}{p(\mathcal{D})} \text{ proportionnel à } p(C_i) p(R_i|C_i) \quad (3.23)$$

$p(\mathcal{D})$ est une constante (\mathcal{D} étant donnée). $p(C_i)$ correspond au prior et $p(\mathcal{D}|C_i)$ à la vraisemblance. La proportionnalité indiquée est due au fait que $\mathcal{D} \setminus \{R_i\}$ ne dépend pas de C_i : en clair, la modélisation est faite de sorte que le choix de la catégorie du variant i (dans le réplicat étudié) n'affecte pas la vraisemblance du reste des données.

Les catégories possibles sont :

- $C_i = \text{vrai}$: le variant i correspond à une vraie espèce ;
- $C_i = \text{ponct}_{ij}$: le variant i est issu d'une mutation ponctuelle de j ;
- $C_i = \text{coocc}_{ij}$: la variant j est issu de j par un autre type de mutation détectée par cooccurrence.

Le prior choisi est $\mathbb{P}(C_i = \text{vrai}) = 10^{-3}$, en utilisant l'ordre de grandeur observé pour des données de contrôle positif. Ensuite, les autres catégories se voient attribuer le prior $\mathbb{P}(C_i) = \frac{1 - \mathbb{P}(C_i = \text{vrai})}{\text{Nombre de parents}}$.

On choisit la vraisemblance d'une vraie espèce selon la loi suivante :

$$\frac{1}{R_i} | C_i = \text{vrai} \sim \text{Exp}(10) \quad (3.24)$$

$$\text{soit } p(R_i | C_i = \text{vrai}) = \frac{1}{10} \exp\left(-\frac{1}{10 R_i}\right) \quad (3.25)$$

Ce choix est retenu car il est simple et prend en compte l'abondance des séquences. En revanche, le choix du paramètre 10 est arbitraire. Il permet de pénaliser les séquences rares tandis que les autres ont une densité à peu près constante. Ce choix est discuté dans la Discussion générale du manuscrit.

Pour le modèle de substitution, on utilise la vraisemblance calculée plus haut :

$$p(R_i|C_i = \text{punct}_{ij}) = p_{\theta_i}(R_i; \mu_-) \quad (3.26)$$

Pour les mutants détectés par cooccurrence, on utilise une densité empirique $\hat{p}_{\sigma(\text{ratio})}$ estimée par noyau (*kernel density*) à partir des mutants de référence :

$$p(R_i|C_i = \text{coocc}_{ij}) = \hat{p}_{\sigma(\text{ratio})} \left(\frac{R_i}{\epsilon + R_j} \right) \quad (3.27)$$

À partir de ces vraisemblances, chaque probabilité est calculée par l'équation 3.23. La constante de normalisation est obtenue en sommant toutes les catégories c possibles pour i (vraie espèce et parentèle estimée) :

$$p(C_i|\mathcal{D}) = \frac{p(C_i) p(R_i|C_i)}{\sum_c p(c) p(R_i|c)} \quad (3.28)$$

L'attribution des probabilités se fait par ordre croissant du nombre de lectures. Lorsqu'un variant est traité, ses enfants sont réattribués aux parents du parent.

L'abondance exprimée en nombre de lectures attribuées à chacun des variants est :

$$R' = P \cdot R \quad (3.29)$$

P est la matrice de probabilité de dimension $N_v \times N_v$ et R est le vecteur des nombres de lectures (*reads*) de dimension N_v . L'abondance attribuée à chaque classe j est donc :

$$R'_j = \sum_{i=1}^{N_v} \mathbb{P}(i \text{ est un mutant de } j) R_i = \sum_{i=1}^{N_v} P_{ij} R_i \quad (3.30)$$

3.3.5 Adaptation des indices de biodiversité

Les indices de biodiversité reposent sur la connaissance de l'abondance relative de chaque espèce s , p_s . Dans Marcon (2015), p_s est définie comme "la probabilité qu'un individu tiré au hasard appartienne à l'espèce s et dont l'estimateur, \hat{p}_s est la fréquence observée". Cette définition s'adapte bien à une attribution probabiliste des variants aux différentes classes.

3.3.5.1 Estimation de p_s

Les données sont composées de R_{total} lectures. Dans le cas déterministe, chaque observation i correspond à une espèce \mathcal{S}_i . L'abondance relative de l'espèce s est :

$$\hat{p}_s = \frac{1}{R_{total}} \sum_{i=1}^{R_{total}} \mathbb{1}_{\mathcal{S}_i=s} \quad (3.31)$$

Lorsque l'espèce d'une observation est incertaine, on décide de la partager entre les différentes classes possibles. On définit donc :

$$\hat{p}_s = \frac{1}{R_{total}} \sum_{i=1}^{R_{total}} \mathbb{P}(\mathcal{S}_i = s) \quad (3.32)$$

$$= \frac{1}{R_{total}} \sum_{i=1}^{N_v} \mathbb{P}(\mathcal{S}_i = s) R_i \quad (3.33)$$

On ne traite évidemment pas les lectures une par une mais variant par variant. Entre les deux lignes, la sommation change : on passe de R_{total} observations (lectures) à N_v variants regroupant les lectures identiques.

3.3.5.2 Diversité α

À partir de cette définition de p_s , on peut utiliser directement les nombres de Hill pour calculer la diversité α :

$${}^q D = \left(\sum_{s=1}^{N_v} \hat{p}_s^q \right)^{1/1-q} \quad (3.34)$$

La diversité α d'un échantillon est calculée en moyennant les vecteurs d'abondances attribuées établis pour les différents réplicats. La diversité γ s'obtient de la même manière. Je ne l'étudie pas dans ce texte.

3.3.5.3 Estimation de la richesse à partir de $q=0$

Cette définition n'est pas satisfaisante quand $q = 0$. En effet, cela revient à comptabiliser tous les variants observés comme des espèces, alors qu'on en attend beaucoup moins. On ne souhaite pas comptabiliser les mutants probables. Une redéfinition possible de la richesse est :

$$S^* = \sum_{i=1}^{N_v} \mathbb{P}(\text{Au moins une lecture est attribuée au variant } i) \quad (3.35)$$

Il faut donner du sens à cette probabilité. En utilisant la matrice d'abondance R' définie en 3.29, on choisit :

$$S^* = \sum_{i=1}^{N_v} \mathbb{P}(R'_i \geq 1) := \sum_{i=1}^{N_v} \min(R'_i, 1) \quad (3.36)$$

Les variants avec au moins une lecture attribuée sont comptabilisés comme une espèce, les autres le sont avec une certaine pondération.

La première limite de cette formule est que la richesse dépend de l'abondance absolue : doubler le nombre de lectures en remplaçant R par $2R$ conduit à une estimation de S^* plus grande car les abondances R'_s sont remplacées par $2R'_s$.

Ensuite, cette formule induit une discontinuité du spectre de Hill. On a :

$$p_s = \frac{1}{R_{total}} \sum_{i=1}^{N_v} \mathbb{P}(\mathcal{S}_i = s) \cdot R_i = \frac{R'_s}{R_{total}} \quad (3.37)$$

On retrouve donc :

$$\lim_{q \rightarrow 0} {}^q D = N_v > S^* \quad (3.38)$$

La différence étant qu'un variant est comptabilisé pleinement dès que $R'_i > 0$ pour ${}^0 D$, et non de manière progressive jusqu'à $R'_i = 1$ pour S^* définie par 3.36.

3.4 Assignation taxonomique : Résultats

3.4.1 Mutations ponctuelles

La Table 3.3 donne la répartition des séquences selon leur statut dans le graphe de mutations d'*obiclean*. La Table 3.4 donne la répartition des mutations potentielles par type de mutation. On constate que les deux jeux de données ne semblent pas avoir la même structure de mutations : d'une part, *obiclean* ne classe pas les mutants de la même manière, d'autre part, les types de mutation n'ont pas des abondances comparables.

Marqueur	Variants			Lectures		
	<i>head</i>	<i>internal</i>	<i>singleton</i>	<i>head</i>	<i>internal</i>	<i>singleton</i>
<i>Euka03</i>	10%	65%	25%	59%	33%	8%
<i>Sper01</i>	2%	92%	6%	87%	13%	0.3%

TABLE 3.3 – Répartition des variants et des lectures selon le statut donné par *obiclean* dont le fonctionnement est présenté en 3.1.4.3. Je rappelle que les *singletons* définis par *obiclean* sont les variants sans connexion dans le graphe.

Marqueur	Transitions	Transversions	Délétions	Insertions
<i>Euka03</i>	60%	30%	9%	0.2%
<i>Sper01</i>	48%	16%	24%	12%

TABLE 3.4 – Part de chaque type de mutation ponctuelle détectée par *obiclean*. Le fort taux d'insertion/délétion pour *Sper01* est en fait dû à des homopolymères de certains codes-barres : il s'agit plutôt d'erreurs de type *slippage*.

3.4.2 Ratio de mutants

Même si le type de mutation est un facteur important pour le ratio de mutants, la variabilité est grande au sein d'une même classe. En effet, la mutation dépend d'autres facteurs, par exemple la position dans la séquence ou les bases voisines, dont on ne souhaite pas décrire la dépendance. La Figure 3.4.4 montre la distribution de ces ratios pour le marqueur *Euka03*. Le pic autour de 10^{-2} est un artefact dû au critère des 100 lectures de la séquence souche dans la construction du graphe d'*obiclean*. La Figure 3.4.5 montre les résultats en séparant les types de substitution symétriques pour illustrer leur similarité. Cette similarité provient de la PCR, pas du séquençage.

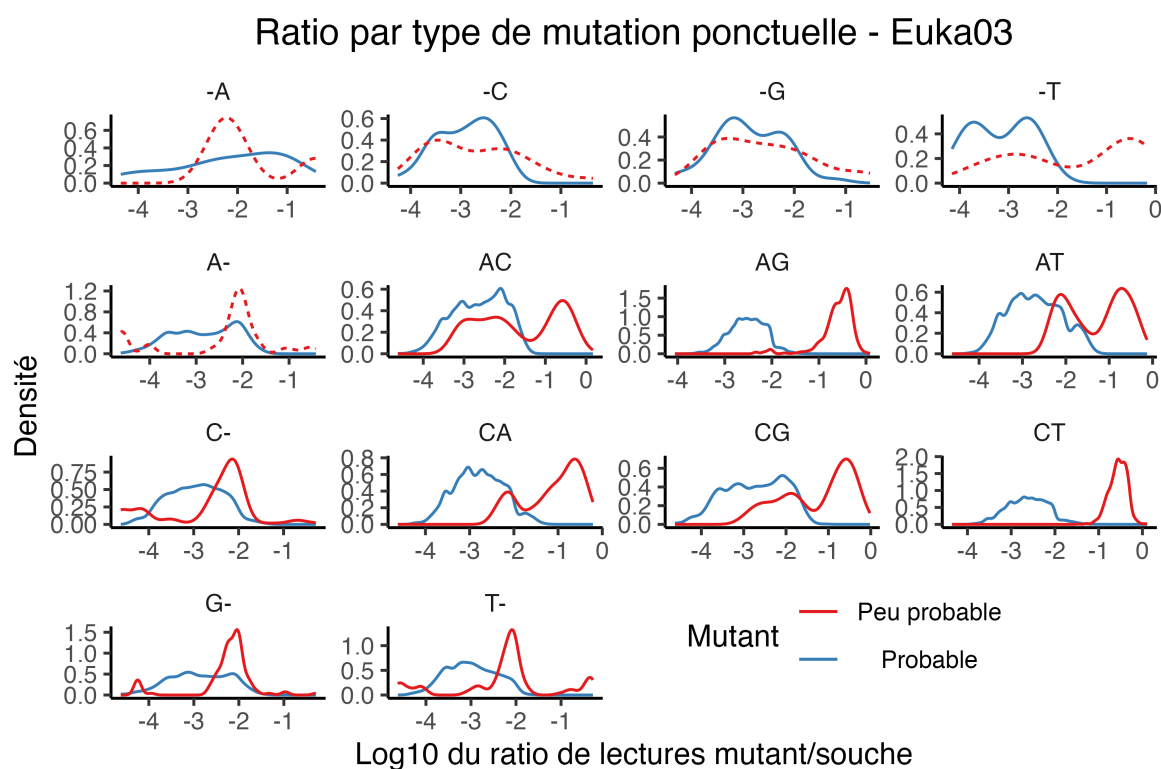


FIGURE 3.4.4 – Densité du logarithme du ratio de mutants pour les différents types de mutations ponctuelles dans les données *Euka03* et selon que la vraisemblance du mutant est supérieure à 10^{-3} (765 121 occurrences) ou non (6 622 occurrences). Les lignes pointillées correspondent aux cas où la densité est calculée à partir de moins de 100 valeurs.

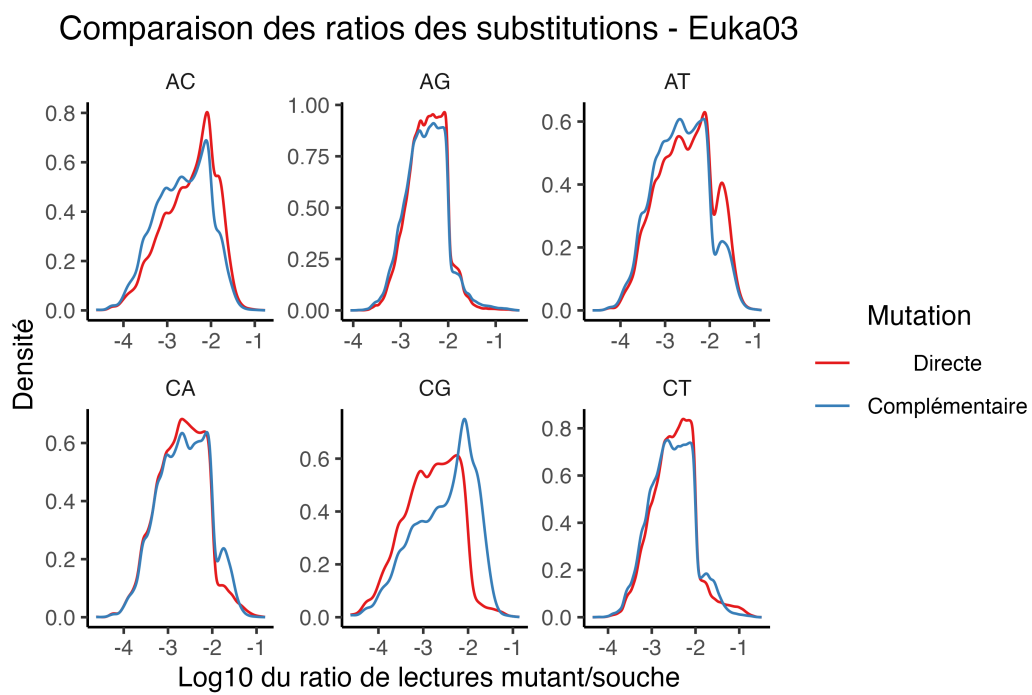


FIGURE 3.4.5 – Densité du logarithme du ratio de mutants probables. Les types complémentaires (TG pour AC, etc., voir Table 3.3.2.1) sont affichés séparément pour faire apparaître la similarité observée.

3.4.3 Cooccurrences des variants

Les Figures 3.4.6 et 3.4.7, montrent les critères considérés pour caractériser les paires mutant-souche pour *Euka03*. Sur chacune, trois jeux de données sont représentés. Le premier est constitué des 45 366 mutants considérés comme vraisemblables par le modèle de mutation. Le deuxième regroupe les 924 mutants rejetés dans au moins un réplicat. Le troisième comprend 10 000 paires de séquences quelconques où la "souche" est un variant racine du graphe avec au moins 10 occurrences et le "mutant" est une séquence choisie au hasard, avec une plus grande probabilité donnée aux séquences cooccurrent avec la "souche". Pour les trois corrélations, les lignes verticales représentent les quantiles 10^{-2} , 10^{-3} et 10^{-4} de la distribution des vrais mutants. Pour l'écart-type et la déviation absolue à la moyenne du ratio, les lignes verticales correspondent aux quantiles 0.99, 0.999 et 0.9999 de la distribution des vrais mutants.

Pour les trois corrélations, les distributions se superposent et les quantiles de la distribution de référence sont assez espacés. C'est nettement moins le cas pour l'écart-type du ratio qui ne génère d'ailleurs aucun "faux négatif" parmi les mutants avérés. Ce critère nous paraît donc préférable. La distribution de la déviation à la moyenne du ratio semble équivalente à celle de l'écart-type. Elle n'a pas été considérée dans la suite. La densité de l'écart-type du ratio pour les deux jeux de données est montrée sur la Figure 3.4.8. L'écart-type médian pour *Sper01* est de 3.8×10^{-4} et de 1.9×10^{-3} pour *Euka03*.

Pour certains variants, la distribution des lectures est compatible avec le modèle de mutation ponctuelle mais les corrélations sont mauvaises (inférieures à 0), comme illustré sur la Figure 3.4.9. Les (nombreuses) absences du mutant doivent être prises en compte car elles font partie de l'information, mais elles ont une influence importante sur le calcul des corrélations. Cela n'exclut pas que le variant mutant soit aussi issu d'une autre souche présente dans seulement certains échantillons.

Pour conclure ces résultats de cooccurrence, les écarts-type du ratio sont comparés à la moyenne des ratios sur la Figure 3.4.10 pour les mutants avérés. La dépendance est claire.

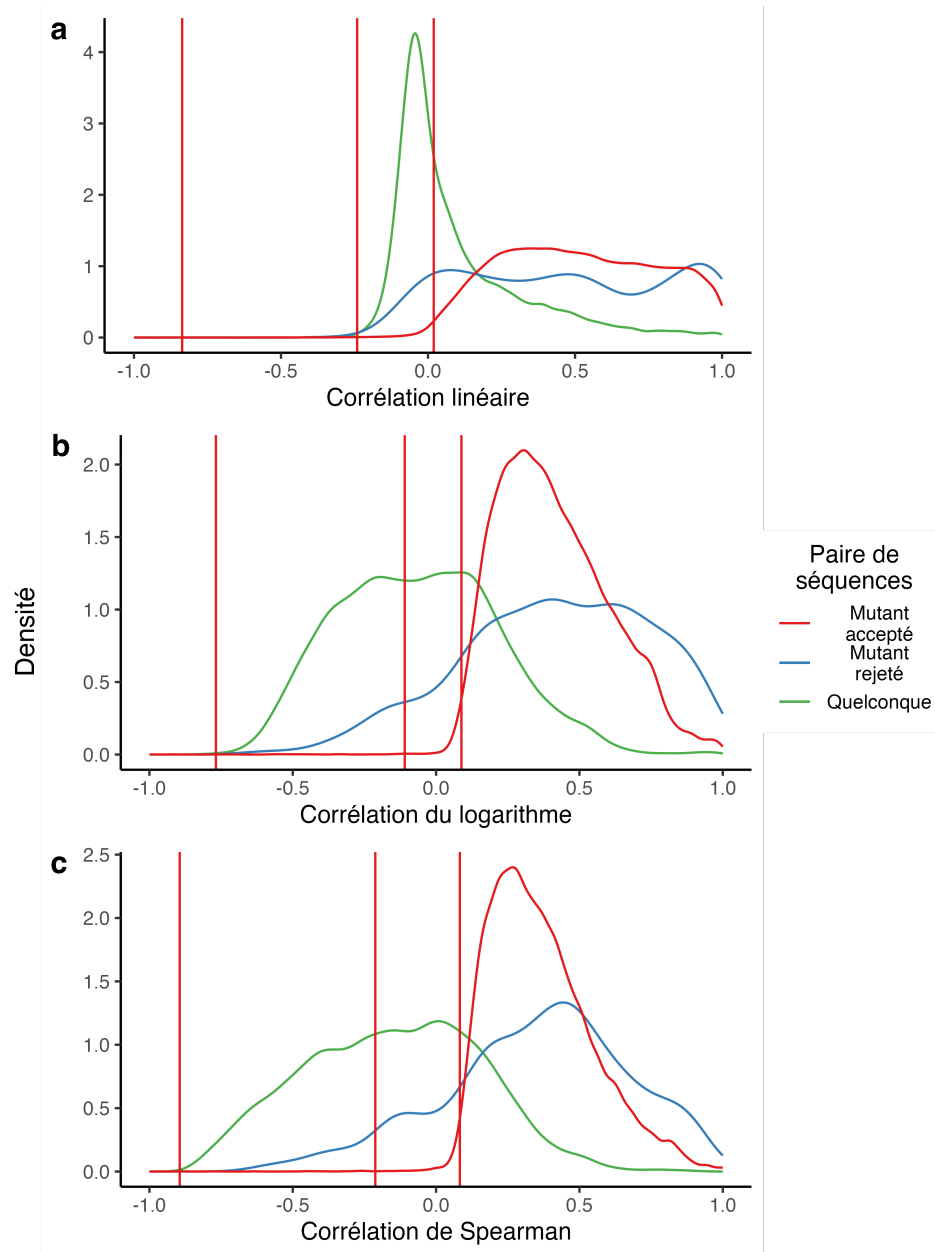


FIGURE 3.4.6 – Corrélations des lectures de la souche et du mutant potentiel. Panneau **a** : corrélation linéaire. Panneau **b** : corrélation linéaire du logarithme des lectures. Panneau **c** : corrélation de Spearman.

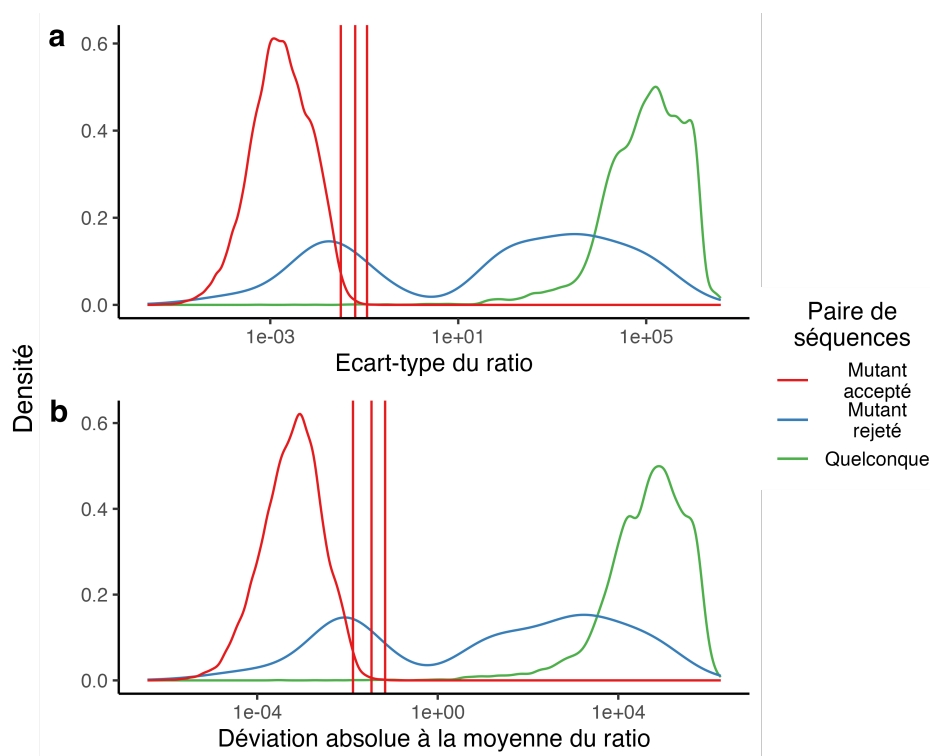


FIGURE 3.4.7 – Densité du logarithme de l'écart-type (panneau **a**) et de la déviation absolue à la moyenne (panneau **b**) du ratio des lectures des mutants par rapport à leurs souches. Selon ces critères, une partie des mutants rejetés par le modèle de mutation ponctuelle semble acceptable (premier mode de la distribution bleue).

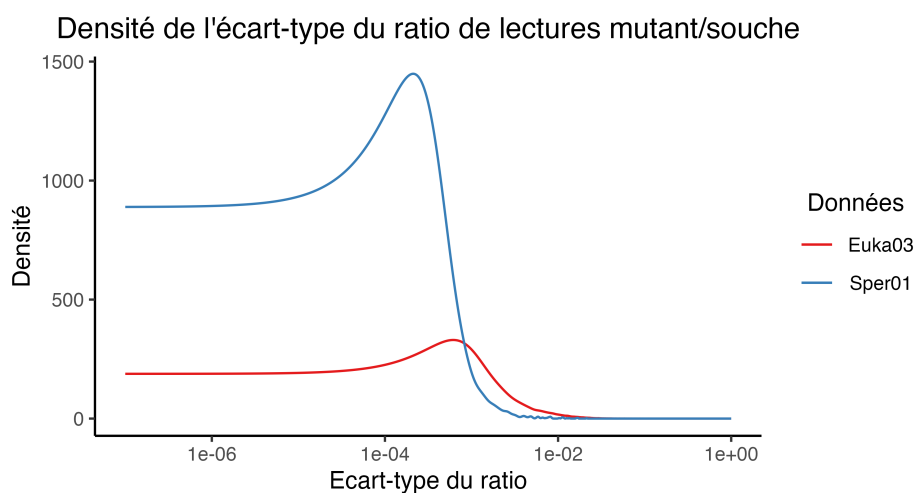


FIGURE 3.4.8 – Densité de l'écart-type du ratio de lectures mutant/souche observée chez des mutants avérés pour les deux jeux de données. La différence d'amplitude des courbes est un artefact de l'échelle logarithmique en abscisse : la distribution pour *Sper01* est piquée autour d'une valeur plus faible (2.1×10^{-4}) qu'*Euka03* (6.2×10^{-4}).

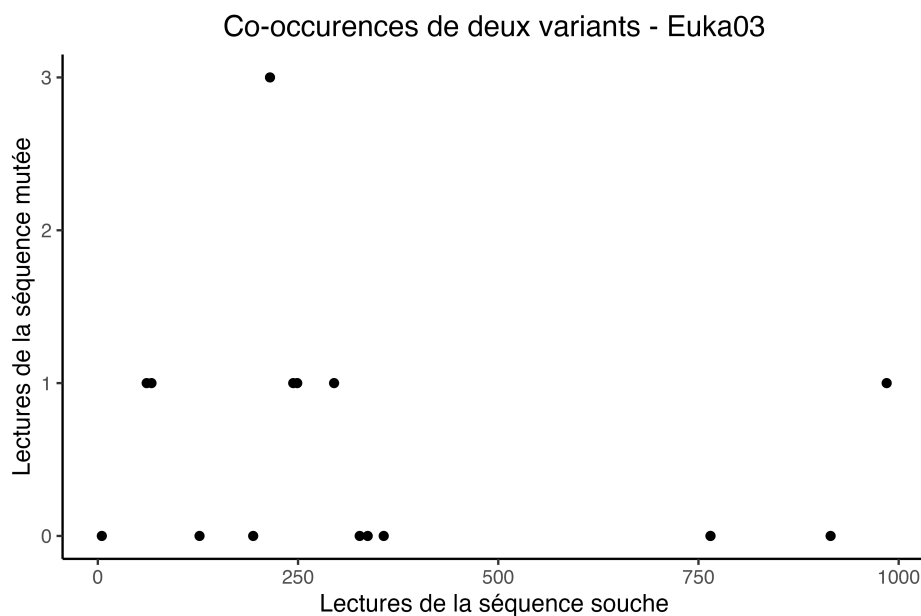


FIGURE 3.4.9 – Lectures attribuées à deux séquences souche et mutant (substitution $A \rightarrow G$) dans seize réplicats. Ce mutant n'a pas d'autres parents identifiés. Il est présent sept fois et absent neuf fois. La corrélation linéaire est de -0.20, la corrélation des logarithmes de -0.059 et la corrélation des rangs de -0.24. En revanche, l'écart-type du ratio est de 5.9×10^{-3} , ce qui le classe comme mutant probable d'après la base de référence.

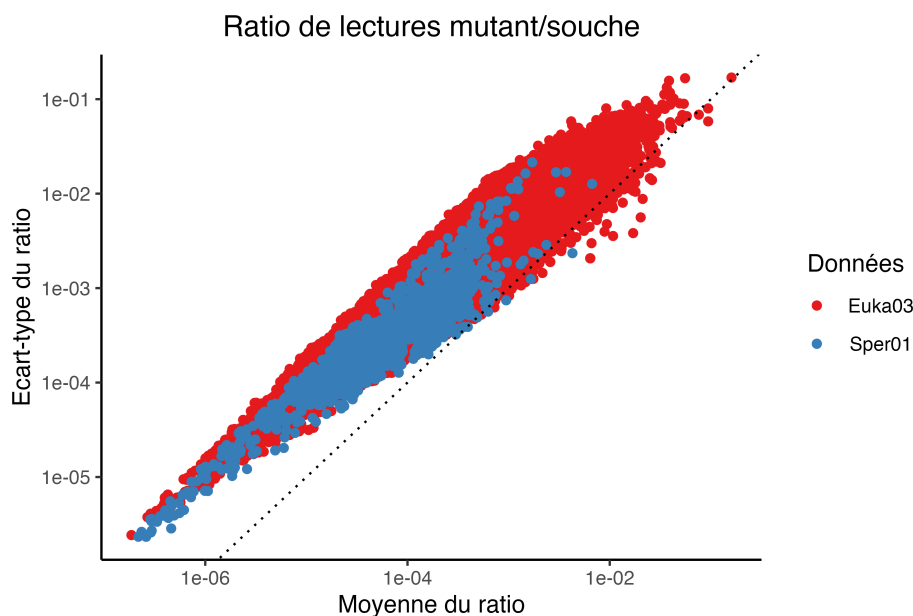


FIGURE 3.4.10 – Lien entre l'écart-type et la moyenne du ratio de lectures pour les deux jeux de données. L'échelle log-log écrase un peu la dispersion : la corrélation des deux grandeurs est de 0.70 pour *Sper01* et de 0.77 pour *Euka03*.

3.4.4 Similarité génétique des séquences cooccurrentes

Certains variants cooccurrent de manière significative avec un variant plus abondant en l'absence de mutation ponctuelle. Pour ces paires de séquences, la distance de Levensthein ne révèle pas de proximité particulière. L'observation des séquences concernées montre parfois des chimères mais aussi des paires de séquences sans lien apparent. Cela montre que le critère de l'écart-type du ratio n'est pas très spécifique même si sa sensibilité est bonne.

Un prochain objectif du projet est donc de mettre au point une approche qui résout ce problème, notamment lors du choix des paires de séquences testées.

3.4.5 Matrice de probabilité

Voici quelques éléments d'analyse des matrices de probabilité d'assignation taxonomique. La Figure 3.4.11 montre le nombre de variants attribués à chacun des variants (sans prendre en compte les abondances), selon :

$$\text{Nombre de variants attribués à } j = \sum_{i=1}^{N_v} P_{ij} \quad (3.39)$$

Un variant avec une valeur de 1 est a priori un variant isolé. Si la valeur est inférieure, ce variant est partiellement attribué à un autre variant dont il est une erreur possible. À l'inverse, les variants avec plus d'une attribution ont absorbé des erreurs.

Pour *Euka03*, 60 à 70 % des variants ne se trouvent liés à aucun autre. Cela est dû à une raison technique. Seulement une partie du graphe de mutations est utilisée ici puisqu'on ne considère que les souches avec au moins cent lectures attribuées.

De 3% à 5% des variants ont des parents "par cooccurrence". Ce résultat est assez faible et indique qu'il faudra améliorer la stratégie de comparaison des paires de séquence, notamment pour identifier des chimères.

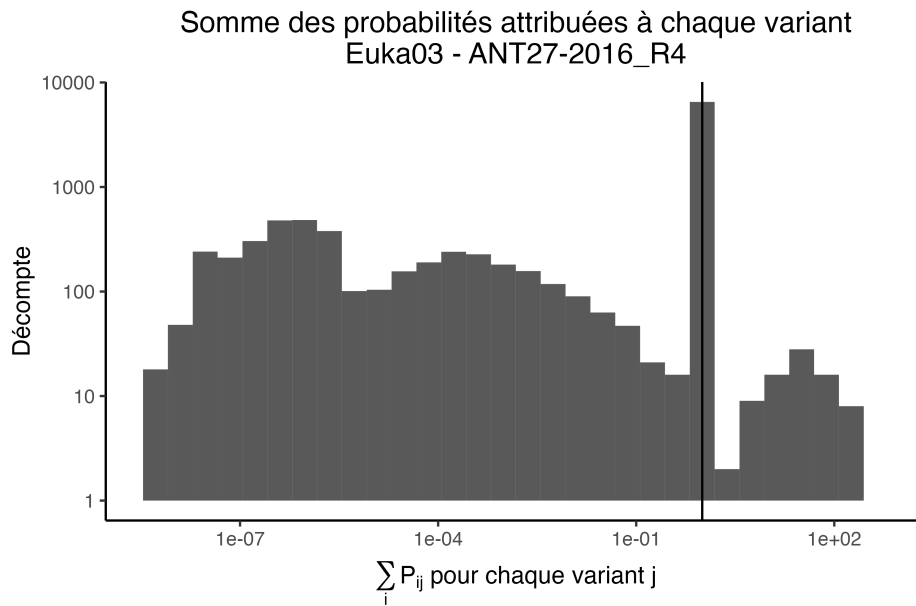


FIGURE 3.4.11 – Histogramme du nombre de variants attribué à chacun des variants pour un réplicat de *Euka03*. Ici, 85 variants se sont vu attribuer plus d’un variant.

3.4.6 Estimation de la diversité α

La Figure 3.4.12 montre le spectre de Hill pour les quatre réplicats de l’échantillon ANT27-2016 des données *Euka03*. Mon estimation de la biodiversité se situe systématiquement entre celle des lectures brutes et celle d’*obiclean*. C’est logique : chaque variant garde une fraction de son abondance du fait de l’incertitude d’attribution. Je reviens sur ce point dans en Discussion Générale du manuscrit.

La Figure 3.4.13 montre un résultat semblable pour *Sper01*. La Figure agrège les estimations de biodiversité de la communauté artificielle \mathcal{M}_U pour deux concentrations de *Spike* (Table 2.1). Puisqu’il s’agit d’un contrôle positif, il est possible de comparer les spectres obtenus au spectre réel. Pour les petites valeurs de q , le tri plus drastique d’*obiclean* est plus performant. On observe que pour les grandes valeurs de q , toutes les procédures sous-estiment la biodiversité.

Dans les deux cas, la plus-value de notre approche est difficile à établir : pour les petites valeurs de q , elle souffre des mêmes limites de surabondance des variants que les données brutes. Ensuite, la correction est similaire à celle effectuée par *obiclean*. Ce résultat est concordant avec celui de Calderón-Sanou et al. (2020) qui propose d’établir les diversités α , β et γ à partir des nombres de Hill pour un paramètre q compris entre 1 et 2. Le paramètre q est alors utilisé comme un filtre, ce qui le détourne néanmoins de son rôle écologique.

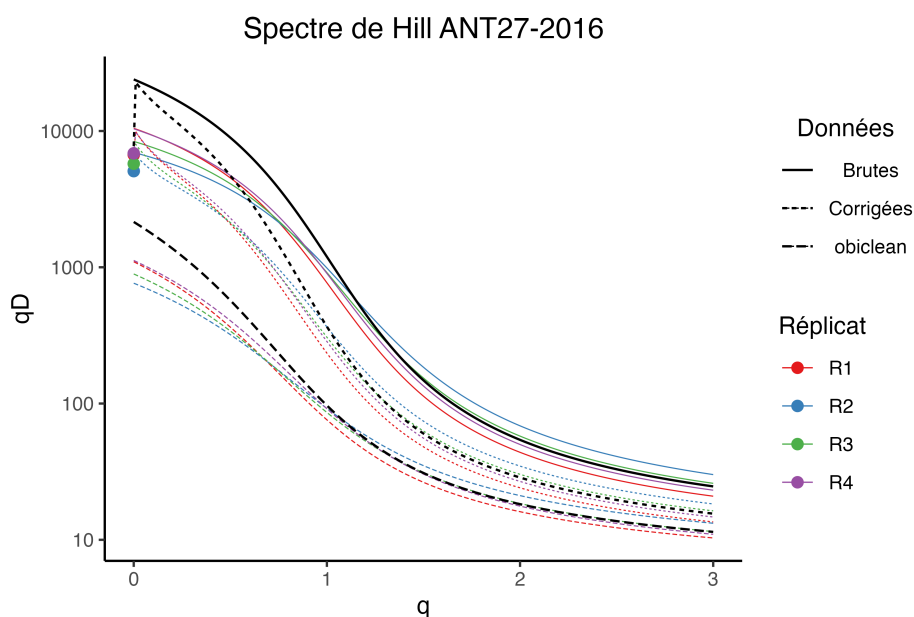


FIGURE 3.4.12 – Spectre de Hill pour un échantillon (quatre réplicats) pour trois procédures de (non-)tri des variants. "Corrigées" signifie : traitées par la procédure de ce texte. Les lignes noires représentent la diversité α estimée en moyennant les abondances des réplicats. Les points représentent l'estimation du nombre d'espèces S^* définie en 3.3.5.3.

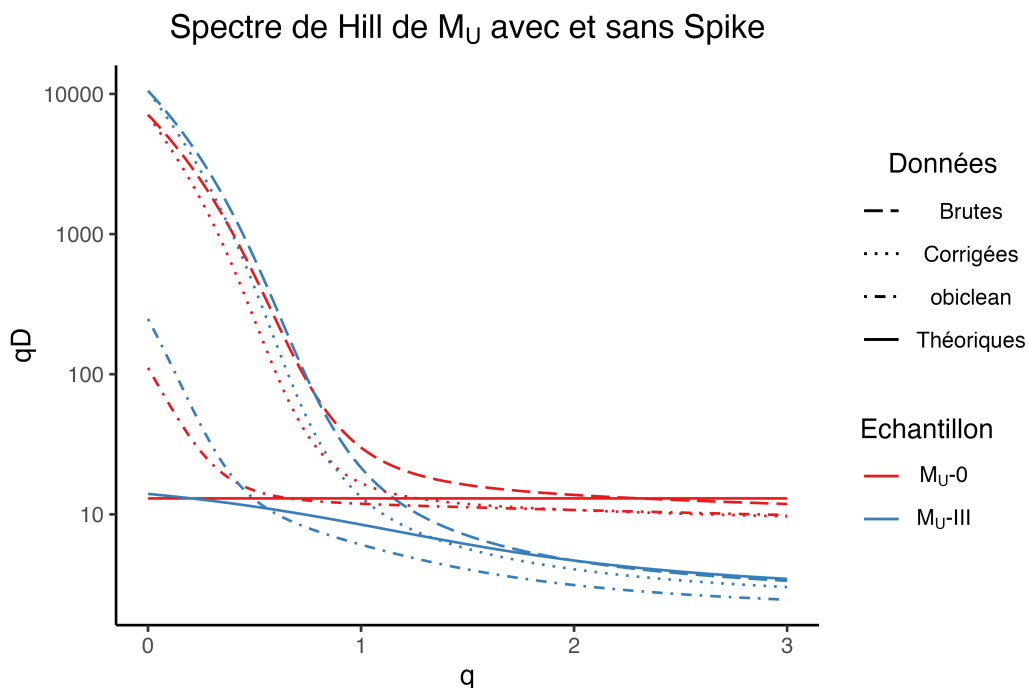


FIGURE 3.4.13 – Spectre de Hill de la communauté \mathcal{M}_U de *Sper01* pour trois procédures de (non-)tri des variants et le spectre théorique. "Corrigées" signifie : traitées par la procédure de ce texte.

3.5 Collaboration : Génotypage de loups

Les travaux présentés dans ce chapitre sont en lien avec un projet de l'équipe de Tomaž Skrbinšek (Université de Ljubljana, Slovénie) auquel j'ai contribué sur invitation de Frédéric Boyer (LECA). L'étude vise à définir le profil génétique de loups à partir d'une série d'allèles. L'ADN est prélevé à partir de traces de morsures. L'ADN cible est composé de microsatellites (répétition d'une séquence d'ADN).

Ma contribution a consisté à développer avec Frédéric Boyer un algorithme de détermination du génotype le plus probable. Pour un échantillon de huit réplicats, les données se présentent sous la forme de variants formant un graphe dont la Figure 3.5.14 donne un exemple. L'objectif est de déterminer quelles séquences sont les vrais allèles, et en particulier si l'individu est homozygote ou hétérozygote. L'ADN étant très peu concentré, l'un ou l'autre des allèles peut ne pas être détecté dans certains réplicats (*dropouts*), comme l'illustre la Figure 3.5.15.

Nous avons évalué les génotypes en prenant en compte la disparition possible d'un allèle et en établissant un modèle d'erreur traitant les mutations ponctuelles et les *slippages*. Une vraisemblance est ainsi attribuée aux différents génotypes considérés ("allèles = séquence i et séquence j "). Cette modélisation, dont je ne donne pas le détail, est proche du modèle de mutation que je présente dans ce chapitre. Un manuscrit est en préparation.

3.6 Conclusion

Dans ce chapitre, j'ai étudié les erreurs observées dans les données de métabarcoding pour comprendre leurs origines et leur structure ainsi que les approches existantes pour les corriger. J'ai ensuite étudié plusieurs nouvelles pistes pour les traiter et extraire un signal écologique plus clair à partir de ces données. Je discute en détail des contributions et des perspectives de ce projet dans la Discussion générale.

WCH07719 – Lup21 – dropout/ratio

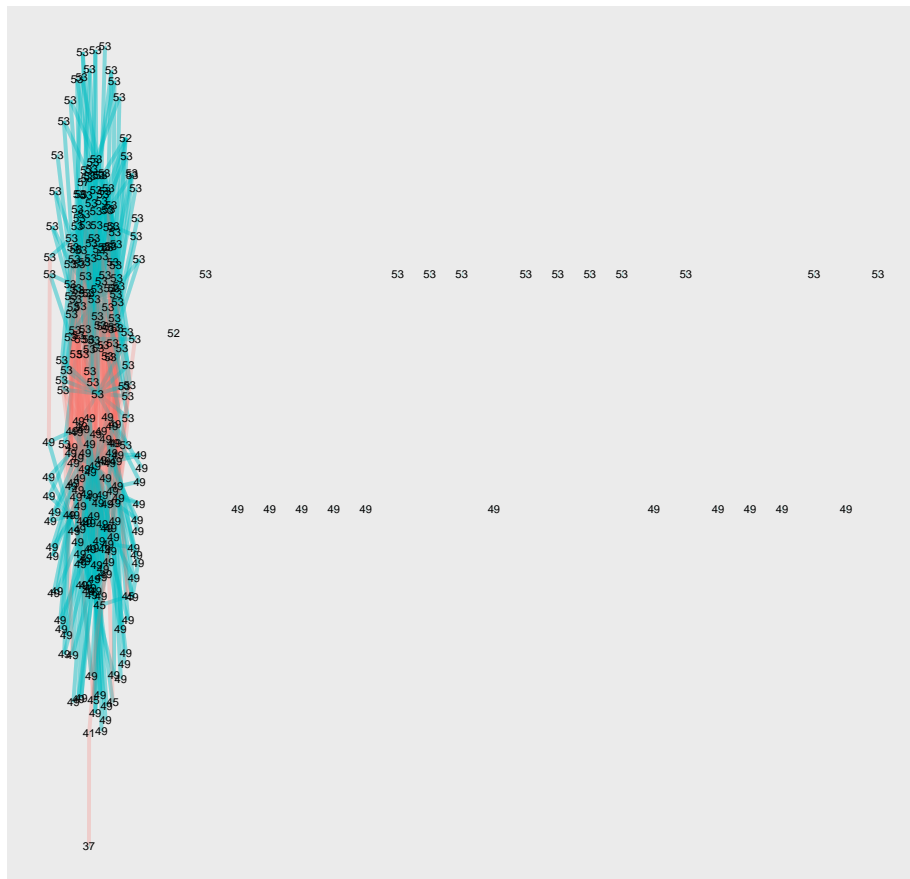


FIGURE 3.5.14 – Graphe des variants observés pour un échantillon analysé dans huit réplicats. Ici, l’individu est hétérozygote. Les nombres correspondent au nombre de répétition du motif des microsatellites. Les nœuds du graphe sont des variants. Les arêtes indiquent qu’il est possible de passer de l’un à l’autre avec une mutation ponctuelle ou un *slippage*. Figure : Frédéric Boyer.

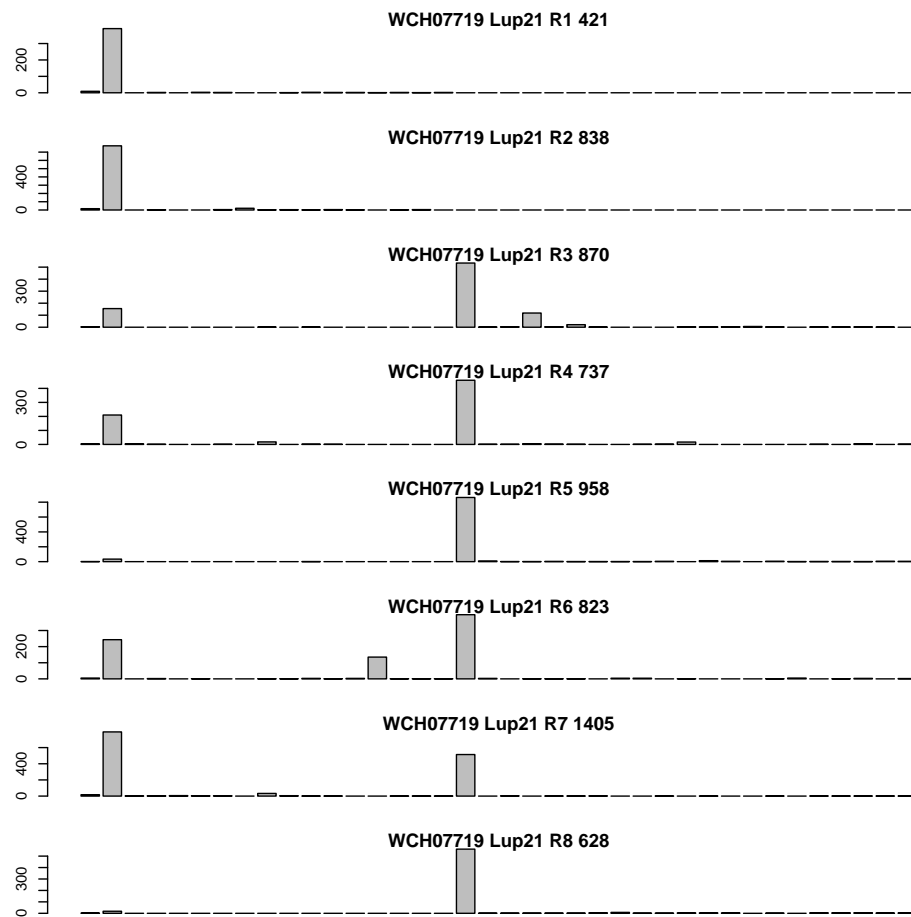


FIGURE 3.5.15 – Nombre de lectures selon la longueur du variant pour chacun des huit répliqués d'un échantillon. On observe des *dropouts* dans les répliqués 1, 2, 5 et 8 : l'un ou l'autre des allèles manque. Figure : Frédéric Boyer.

Chapitre 4

Discussion générale et perspectives

Dans cette dernière partie, je résume les principales contributions de ma thèse et les perspectives qui s'ouvrent à la suite de mes projets.

La motivation de ma thèse était la conception d'indicateurs de biodiversité robustes pour le métabarcoding. Pour y parvenir, j'ai consacré deux de mes projets à l'étude des données de métabarcoding. Le Chapitre 2 a abordé une première limite de ces données : le signal écologique est biaisé par le protocole du métabarcoding car les abondances des différentes espèces sont modifiées au cours de l'expérience. Le Chapitre 3 a traité le fait que le signal écologique est bruité : avant de tirer des conclusions sur la biodiversité représentée dans l'échantillon d'ADNe, il faut agréger les séquences d'ADN obtenues.

Ma thèse a aussi été l'occasion de développer de nouveaux outils d'analyse utilisés en écologie. Le développement méthodologique est un champ de recherche transversal passionnant : les outils conçus contribuent à répondre à des problématiques variées et à remettre en question les pratiques de recherche actuelles. C'est l'objet du Chapitre 1 décrivant l'algorithme d'inférence *flimo* que j'ai développé et appliqué à différents problèmes en lien avec l'écologie. Les deux chapitres suivants partagent cet objectif méthodologique en proposant des protocoles expérimentaux et informatiques applicables à l'ensemble des expériences de métabarcoding.

Ces travaux m'ont amené à jongler entre plusieurs disciplines : j'ai passé du temps à développer et à implémenter des modèles mathématiques mais j'ai aussi mené par moi-même (avec un grand soutien des équipes techniques du LECA) des expérimentations en biologie moléculaire, de la conception à la paillasse. Ma formation en biologie avant ma thèse était minime mais ce défi a été enrichissant à relever.

J'établis maintenant un bilan de chacun de mes chapitres et les recherches que je suggère de mener pour poursuivre mes travaux.

4.1 Chapitre 1 : une nouvelle méthode d'inférence pour modèles aléatoires

4.1.1 Bilan de l'étude

Dans le premier chapitre, je me suis intéressé à la problématique commune en écologie de l'inférence de paramètres pour des modèles aléatoires. La méthode *Fixed Landscape Inference MethOd* (*flimo*) occupe une place à part dans mes travaux car elle n'a pas de lien direct avec le métabarcoding et les thématiques traitées au LECA. Son développement a été motivé par une difficulté technique dans mon projet sur le métabarcoding quantitatif. La puissance de cette méthode m'a ensuite conduit à en faire un projet indépendant qui a pris de l'ampleur au cours de ma thèse.

Comme je l'ai exposé en introduction, la littérature est riche au sujet des méthodes d'inférence pour les modèles aléatoires. Il s'agit d'un enjeu technique important, en particulier en écologie et en biologie où les phénomènes étudiés sont représentés par des modèles complexes. Les méthodes reposant sur des simulations sont intéressantes car elles évitent des développements analytiques lourds. En revanche, les approches bayésiennes nécessitent un temps de calcul important pour construire une distribution postérieure des paramètres du modèle. La méthode *flimo*, elle, est conçue pour établir un estimateur ponctuel des paramètres. Cette approche limite l'information à fournir et permet d'accélérer le processus d'inférence. Pour cela, l'aléa des simulations est supprimé par un couplage de sorte qu'elles conservent les propriétés des variables aléatoires représentant le modèle étudié. Les paramètres sont ensuite estimés grâce à des algorithmes d'optimisation déterministes. Les applications présentées ont montré la grande efficacité de *flimo* sans concession sur la précision des résultats. Son cadre d'application est vaste et ne se limite a priori pas à l'écologie.

Bien sûr, la méthode *flimo* n'est pas exempte de limites. Certaines limites techniques sont communes à d'autres algorithmes comme les ABC. Pour chaque nouveau problème, il faut choisir des statistiques résumées les plus informatives possibles tout en étant faciles à manipuler : c'est la contrepartie des méthodes basées sur des simulations. J'ai émis des hypothèses sur les limites d'application de *flimo* en lien avec la régularité des modèles étudiés mais celles-ci sont relativement spéculatives et doivent être étudiées plus en détail. A priori, *flimo* ne devrait pas être approprié pour un modèle avec des discontinuités marquées (par exemple le choix de positionnement de branches dans un arbre) mais il est possible que d'autres classes de modèles soient inaccessibles. Cela pourrait être le cas des modèles qui ne sont pas constitués de combinaisons de variables aléatoires simples. Un certain nombre de perspectives concrètes ont été identifiées et sont évoquées plus loin.

L'usage des statistiques résumées restreint aussi les garanties théoriques sur les résultats obtenus. Rien ne prouve que les valeurs trouvées par *flimo* sont justes : le fonctionnement est heuristique. Je suis convaincu que la procédure de *flimo* n'induit pas en soi de biais dans les résultats d'inférence mais cela n'a pas été démontré

formellement. Ces aspects théoriques sont importants pour la poursuite du projet et nécessiteraient une collaboration avec des spécialistes du domaine, pour prendre du recul par rapport au développement actuel de *flimo*.

Un autre enjeu est la diffusion de *flimo*. Dans cette optique, la publication du manuscrit est bien sûr l'échéance principale. De nouvelles applications permettront aussi d'identifier plus finement le besoin auquel répond *flimo* ("résoudre rapidement un problème d'inférence complexe") dans des contextes où les algorithmes actuels ne sont pas satisfaisants.

4.1.2 Perspectives

J'ai identifié plusieurs sujets concrets qui justifieraient des travaux complémentaires.

4.1.2.1 Nouvelles applications

Un premier objectif est la généralisation de la méthode à travers de nouvelles applications. Deux ont déjà été envisagées. D'une part, François Munoz (LiPhy, Grenoble) m'a proposé d'adapter *flimo* pour le package *ecolottery* (Munoz et al., 2018)¹. Ce package modélise l'assemblage de communautés soumises à des filtres environnementaux. D'autre part, j'ai été sollicité par Ramon Grima (Université d'Edinburgh, Royaume-Uni) pour tester *flimo* sur un modèle de réseaux de réactions biochimiques (Ocal et al., 2020). Ces deux études reposent actuellement sur des approches bayésiennes et présentent a priori un cadre favorable pour *flimo*.

4.1.2.2 Optimisation globale

L'implémentation de *flimo* peut d'ores et déjà être améliorée. Edouard Oudet a suggéré d'utiliser des algorithmes d'optimisation globale pour l'inférence. Actuellement, la fonction objectif est minimisée localement avec des algorithmes de type quasi-Newton. L'optimisation locale est efficace à certaines conditions difficiles à garantir en pratique.

L'optimisation globale n'assume pas de propriété de la fonction objectif. Dans ce cadre, diverses stratégies existent pour explorer l'espace des paramètres. Des méthodes aléatoires heuristiques sont utilisées, comme l'algorithme du recuit simulé. Lors du développement de *flimo*, j'ai rapidement exclu ces algorithmes stochastiques car nous cherchions justement à faire disparaître l'aléa des problèmes étudiés, mais je pense que cette décision devrait être reconsidérée.

Le rôle de *flimo* est de lisser et de figer le paysage des paramètres, qui est aléatoire dans le problème originel. Mais la manière d'explorer cet espace figé n'est pas imposée par *flimo*. La fonction objectif ne vérifie pas forcément les hypothèses de prédilection des

1. <https://github.com/frmunoz/ecolottery>

algorithmes d'optimisation déterministes, en particulier la convexité. L'optimisation globale pourrait donc être une alternative pertinente, notamment pour les problèmes multidimensionnels non différentiables traités par la méthode de Nelder-Mead.

4.1.2.3 Accessibilité de l'implémentation

Enfin, l'interface utilisateur de *fimo* dans les packages R et Julia pourrait être reprise en main. C'est un point clé pour promouvoir notre méthode auprès d'utilisateurs qui hésiteraient entre plusieurs approches. Actuellement, je n'ai pas eu de retour sur les fonctionnalités, la documentation, les exemples... La mise en place de *fimo* exige de modifier le code du simulateur pour remplacer les tirages aléatoires par des appels aux fonctions quantiles. Ça n'est pas le cas des implémentations d'ABC par exemple. Cette conversion peut être effectuée de manière mécanique (en R, remplacer les *rnorm* par des *qnorm*, etc.). Mais la gestion de la matrice des quantiles aléatoires (quels quantiles utiliser à quelle étape du modèle) empêche pour le moment une automatisation du processus.

4.2 Chapitre 2 : Analyse quantitative des données de métabarcoding

4.2.1 Bilan de l'étude

Au cours de ma thèse, de nombreux travaux ont été publiés sur le métabarcoding quantitatif. Cette problématique est en pleine effervescence et la communauté scientifique est loin d'avoir répondu à toutes les questions soulevées. Cela conforte la pertinence du projet présenté dans le Chapitre 2. Le métabarcoding quantitatif vise à calculer des indices de biodiversité plus représentatifs que ceux calculés à partir des données brutes. Au début de ce projet, le biais d'amplification par PCR était déjà connu et étudié mais sa mesure précise n'avait pas été réalisée. Des facteurs à l'origine des biais avaient déjà été identifiés comme les mismatches et le taux de GC des codes-barres. Par ailleurs, j'ai l'impression que le biais de concentration de l'ADN cible était souvent négligé dans les études sur les macro-organismes. Même si le métabarcoding est une technique très utilisée, la littérature montre que le traitement des biais est souvent superficiel. Les abondances des lectures sont souvent assimilées aux abondances des espèces. C'est pour cela que les comparaisons avec les données d'abondance établies par des méthodes traditionnelles sont parfois contradictoires, que les conclusions varient en fonction du marqueur utilisé, etc.

Mes contributions sont à considérer selon deux points de vue. D'une part, j'ai apporté de nouvelles connaissances sur la manière dont sont produites les données de métabarcoding et en particulier la nature et l'amplitude des biais qui les affectent. D'autre part, j'ai proposé un protocole réaliste pour mesurer et corriger lesdits biais. Ce protocole peut être utilisé dès à présent pour les expériences de métabarcoding réalisées pour suivre la biodiversité.

Mon premier résultat est la mesure des concentrations d'ADN ciblé par le marqueur choisi pour étudier un groupe d'espèces. Ces concentrations varient de manière non négligeable, d'un facteur maximal de plus de 6 dans mes travaux. Ce biais est un problème pour la composition de communautés artificielles, régulièrement utilisées à des fins de quantification. Ces communautés ont un sens écologique si elles sont constituées à masse de tissu égale mais dans certains cas, elles sont composées à quantité égale d'ADN total dosé par Qubit (Clarke et al. (2014) ou Piñol et al. (2015) par exemple). Cela introduit un biais dans l'étude. La ddPCR est la technique appropriée pour mesurer ces concentrations. Nous avons d'ailleurs effectué plusieurs dosages par qPCR avant de réaliser que l'inhibition variable des différents échantillons empêchait de mesurer précisément les concentrations d'ADN cible.

Ensuite, mes résultats portent sur le biais de PCR. Cette étude est divisée en plusieurs composantes. La première est théorique et concerne les modèles de PCR. La modélisation a plusieurs intérêts pratiques. Elle permet de mesurer les biais en estimant les efficacités d'amplification des espèces et en évaluant les conséquences de ces différences de rendement sur les abondances finales. Mais les modèles permettent aussi d'estimer les proportions initiales dans la communauté en corrigeant ce biais. Comme l'illustrent mes développements sur les mismatches d'amorces, les modèles permettent aussi de simuler des réactions PCR pour étudier des hypothèses avant de mettre en place une validation expérimentale. Les modèles que j'ai développés ne sont pas beaucoup plus compliqués que le modèle exponentiel usuel et représentent plus fidèlement la dynamique d'amplification. Cela rend les paramètres plus facilement interprétables et limite les choix arbitraires comme le nombre de cycles effectifs du modèle exponentiel. Les deux modèles étudiés ont des avantages différents : le modèle mécanistique décrit les mécanismes biochimiques de la PCR, le modèle logistique est plus facile à manipuler.

Sur le plan expérimental, la qPCR Taqman a permis de mesurer les efficacités d'amplification de trois espèces. Cette approche est prometteuse bien qu'une généralisation à plus grande échelle soit nécessaire pour la valider. Un grand enjeu sera de déterminer s'il est possible d'établir des bases de référence d'efficacité de PCR, au moins pour des conditions expérimentales données (disons au sein d'un même laboratoire).

Le traitement de ces deux biais a abouti à un protocole complet complémentaire du metabarcoding qui peut dès à présent être utilisé en conditions expérimentales réelles. Toutefois, j'ai le sentiment que certaines limites à la quantification ne seront pas dépassées avec les outils disponibles. Je pense notamment aux facteurs qui altèrent l'ADN environnemental avant qu'il ne soit collecté : taux de dépôt, de dégradation... Il est admis que le biais entre deux espèces est constant (par exemple dans Luo et al. (2022), "*it is reasonable to assume that the ratio of the biases of every pair of species is fixed*"), mais ce postulat me paraît contestable. Notre hypothèse est justement que l'estimation des rendements d'amplification puis des abondances initiales est moins dépendante des conditions expérimentales que les facteurs correctifs établis

directement à partir d'une communauté de référence.

Pour traiter de manière cohérente l'ensemble des biais, chaque échantillon devrait faire l'objet de nombreuses analyses complémentaires, au moins pour estimer l'incertitude de la mesure. Sinon, quel est le sens de corriger finement les biais de PCR si d'autres erreurs relativement importantes sont négligées? Cette question me paraît primordiale pour établir des mesures de biodiversité pertinentes à partir du métabarcoding. Cela conforte l'idée que cette technique doit être conçue comme complémentaire des suivis de biodiversité traditionnels.

J'aborde enfin les techniques d'étude de l'ADN environnemental sans PCR. Celles-ci, comme le séquençage shotgun (avec capture ou non), sont décrites comme des solutions d'avenir pour s'affranchir des biais du métabarcoding. Mes travaux sur la PCR sont-ils déjà obsolètes? Je pense que non, pour plusieurs raisons. D'une part, les techniques alternatives évoquées n'ont pas remplacé le métabarcoding avec PCR car elles imposent d'autres contraintes, par exemple le design de sondes de capture, et car le métabarcoding est simple à mettre en place et peu coûteux. D'autre part, la PCR est une technologie utilisée bien au-delà du métabarcoding. Mes résultats contribuent donc plus largement à la compréhension de ce procédé majeur en biologie moléculaire.

4.2.2 Perspectives

Mes travaux ont apporté des résultats mais aussi beaucoup de nouvelles questions. Voici des pistes importantes pour prolonger mon projet. J'aborde des sujets expérimentaux puis des problématiques de modélisation.

4.2.2.1 Utilisation de la droplet digital PCR

Dans mes travaux, la ddPCR a joué un rôle important mais la finalité du projet était plutôt la mesure des biais de PCR. Cette technique me semble appropriée pour mesurer les autres biais qui affectent les données de métabarcoding, du fait de sa précision et de son protocole standardisé. Par exemple, le taux de dépôt d'ADN dans l'environnement pourrait être analysé selon les conditions biotiques et abiotiques. De même, il est possible d'étudier la dégradation de l'ADN ou la concentration selon la phénologie et le type de tissu. L'étude menée par Stefaniya Kamenova va de ce sens. D'autres travaux ont été réalisés (par exemple Andruszkiewicz Allan et al. (2021); Wilder et al. (2023)) mais essentiellement sur des milieux aquatiques et par qPCR, qui est moins précise que la ddPCR. Il faudrait les généraliser pour d'autres groupes taxonomiques et les intégrer aux protocoles de correction pour le métabarcoding. Idéalement, il faudrait établir des bases de référence pour corriger ces biais de manière systématique et à moindres frais. Cette possibilité dépendra de la variabilité intraspécifique de ces mesures pour une condition ou un type de tissu donné.

4.2.2.2 Utilisation de la qPCR Taqman

J'ai proposé un protocole pour mesurer les efficacités de PCR par la qPCR Taqman pour s'affranchir de l'effet variable des inhibiteurs. Ces travaux devraient être complétés par une série de tests. Il faudrait réaliser d'autres communautés artificielles de composition variable avec des sondes pour chaque espèce pour s'assurer de la reproductibilité des mesures. Il faudrait aussi comparer les mesures Taqman sur des séquences naturelles et sur les mêmes séquences synthétiques pour évaluer l'influence de l'inhibition dans les échantillons. Cette étude a été envisagée avec l'entreprise Argaly² en Savoie.

4.2.2.3 Compréhension des mécanismes de la PCR

J'ai l'impression que les mécanismes de saturation de la PCR ne sont pas parfaitement compris. Est-on capable de quantifier le rôle de l'auto-hybridation et des réactifs limitants (dNTP, amorces)? Cette question est importante pour comprendre les biais de PCR car la compétition interspécifique pendant l'amplification n'est pas la même selon l'une ou l'autre des hypothèses. Mon étude des modèles de PCR avec mismatch illustre ce point : les comportements sont différents si les espèces ont une saturation indépendante (épuisement d'une ressource propre ou auto-hybridation) ou une saturation commune.

Cela conduit à une autre question : dans quelle mesure les corrections établies à partir des ratios (avec une communauté de référence) ou des efficacités estimées sont-elles valides lorsque la communauté et les conditions expérimentales changent? En particulier, notre hypothèse que les efficacités d'amplification sont une propriété intrinsèque de la séquence est-elle vérifiée? En d'autres termes, le paramètre Λ_s a-t-il un sens? Il paraît raisonnable de supposer que l'écriture de la séquence détermine son efficacité d'amplification mais celle-ci peut-elle être influencée par des interactions avec le code-barres d'une autre espèce? De manière similaire, on peut se demander si deux espèces avec une efficacité identique (si cela est possible) sont toujours affectées de la même façon si les conditions d'amplification changent.

De manière subsidiaire, l'expérience menée avec la molécule *Spike* pour générer de la compétition n'a pas donné de résultat exploitable. Est-ce uniquement dû à la faiblesse du signal ou bien l'ajout de *Spike* n'a-t-il pas modifié la compétition au sein des échantillons? Il m'est difficile d'émettre des hypothèses pour répondre à ces questions, il me faudrait des connaissances en biologie moléculaire plus complètes.

Enfin, si cela est possible, il serait intéressant de prédire l'efficacité de PCR à partir de l'écriture de la séquence amplifiée. Une option serait de construire une base de données avec une multitude de séquences (synthétiques?) et leur efficacité d'amplification puis d'étudier cette base de données avec des méthodes statistiques adéquates. Nichols et al. (2018) amorce ces travaux en comparant l'amplification selon le taux de GC pour douze séquences synthétiques.

2. <https://www.argaly.com>

4.2.2.4 Variabilité intraspécifique entre les réplicats

La question de la variabilité intraspécifique entre les réplicats d'un même échantillon reste aussi en suspens à la fin de ma thèse. J'ai observé une forte variabilité d'abondance pour certaines espèces et ce constat a été observé par ailleurs (Iwaszkiewicz-Eggebrecht et al., 2023). Je n'ai pas identifié de raison pour laquelle cette dispersion est très différente d'une espèce à l'autre, comme on l'observe sur la Figure 2.4.4 (résultats de métabarcoding avec le *Spike*). Est-ce dû à la qualité de l'échantillon ? À une propriété particulière du marqueur ? À un simple effet de sous-échantillonnage de l'ADN extrait puis plus tard séquencé ? Des travaux suggèrent que les espèces rares ont une plus grande variabilité relative mais je n'ai pas observé ce phénomène dans mes données.

Cette variabilité a une influence sur les conclusions écologiques à partir du métabarcoding. Quelle précision sur les abondances corrigées peut-on espérer sachant la variabilité des données brutes ? Il est possible d'augmenter le nombre de réplicats pour estimer les abondances de manière plus robuste mais il serait intéressant d'imaginer d'autres solutions pour contrôler cette variabilité.

4.2.2.5 Modélisation de la PCR

J'ai consacré une partie importante de mes travaux à la modélisation de la PCR mais le sujet n'est pas clos. Nous avons envisagé de publier le modèle mécanistique assorti d'une revue des modèles existants et de leurs spécificités ; l'étude des mismatches pourrait aussi en faire partie.

J'ai identifié plusieurs limites dans mes travaux actuels. La première est la variabilité intraspécifique que je viens d'évoquer. J'ai développé des modèles aléatoires et non déterministes pour prendre en compte cette variabilité. Elle y intervient essentiellement à l'étape initiale de constitution des réplicats, où les nombres de molécules sont tirées dans une loi de Poisson ou binomiale négative. On constate ensuite que l'amplification et le séquençage génèrent très peu d'aléa (selon le ratio espérance/variance). J'avais aussi travaillé sur un modèle où l'amplification cycle par cycle était simulée par une loi Beta-binomiale pour augmenter sa variance mais cela impliquait une amplification moins régulière qui ne correspond pas aux données observées.

Une autre limite est l'estimation concrète des paramètres du modèle, notamment la capacité de charge du milieu réactionnel ou les paramètres biochimiques du modèle mécanistique. Les valeurs utilisées sont assez spéculatives car elles n'ont pas été déterminées sur des critères biologiques mais par ajustement numérique.

Par ailleurs, il est intéressant de disposer d'un modèle dont on peut calculer analytiquement des quantités comme les moments. Pour les modèles logistique et mécanistique, cela n'est pas possible. Si l'on cherche à calculer par récurrence l'espérance du nombre de molécules M_k , on observe que l'espérance $\mathbb{E}[M_{k+1}]$ dépend de $\mathbb{E}[M_k]$ et de $\mathbb{E}[M_k^2]$. Mais le calcul de $\mathbb{E}[M_{k+1}^2]$ dépend à son tour de $\mathbb{E}[M_k^3]$ et $\mathbb{E}[M_k^4]$, etc. Cette piste de

calcul n'est donc pas réaliste. Une adaptation du modèle avec des quantités d'intérêt explicites faciliterait la tâche d'inférence des paramètres.

4.2.2.6 Lien avec les modèles d'occupation

Les modèles d'occupation sont un outil commun en écologie des communautés pour décrire la répartition des espèces dans les écosystèmes. Ils reposent sur le fait que la détection des espèces est imparfaite et visent à donner des probabilités de présence selon l'ensemble des observations. Mes travaux sur le métabarcoding quantitatif pourraient compléter ces modèles en prenant en compte les différences d'amplification des espèces de la métacommunauté.

Si une espèce est rare dans un écosystème, le risque qu'elle ne soit pas détectée est d'autant plus grand que son efficacité d'amplification est faible pour le marqueur utilisé. Comment déterminer la probabilité de non-détection sachant les espèces détectées, les efficacités d'amplification des espèces de la métacommunauté et l'abondance (réelle mais inconnue) de l'espèce manquée? Cette approche pourrait fournir un nouveau champ d'application à mes travaux.

4.3 Chapitre 3 : Attribution probabiliste des séquences pour le métabarcoding

4.3.1 Bilan de l'étude

Le troisième chapitre portait sur l'analyse des nombreux variants observés dans les données de métabarcoding et les procédés utilisés pour les répartir en unités taxonomiques ayant un sens écologique. De nombreuses idées ont été explorées pour étudier les mécanismes d'apparition des séquences erronées et leur détection.

Tous les objectifs de ce projet n'ont pas été atteints. En particulier, l'idée d'une attribution probabiliste des variants était élégante mais paraît erronée, en tout cas dans son application actuelle. La motivation était de rendre plus robuste le traitement des erreurs en rejetant les paramètres arbitraires des approches existantes. L'attribution probabiliste supprime ces seuils et crée un continuum entre les unités taxonomiques. Mais l'hypothèse sous-jacente est que chaque variant observé peut être issu d'une espèce réellement présente dans l'échantillon. Or, au vu des ordres de grandeur connus (de l'ordre de 1000 fois plus de variants que d'espèces), accorder une probabilité, même infime, à chaque variant revient en fait à tolérer une part importante du bruit des données. C'est ce qu'on observe sur les spectres de Hill, par exemple sur la Figure 3.4.12. Dans ce sens, mes travaux ne permettent pas encore de remplacer l'algorithme de *obiclean*.

Je pense que l'intérêt du projet ne réside là. Les différentes composantes développées fournissent en soi des outils pertinents qui méritent d'être valorisés. Le projet d'attribution taxonomique est divisé en trois parties : l'étude des mutations

ponctuelles, l'étude des cooccurrences et l'assignation probabiliste des variants.

Le modèle de mutation ponctuelle répond à l'objectif de prendre en compte la spécificité des données (différents taux selon la mutation, variabilité entre les mutants et les échantillons) avec un nombre de paramètres raisonnables. Il me semble assez facile de construire un algorithme de *denoising* "classique" à partir de ce modèle, à la manière de DADA2 ou de UNOISE3.

Les cooccurrences entre séquences souches et séquences mutées sont une contribution importante pour la détection des erreurs. Elles ont déjà été étudiées (Frøslev et al., 2017; Olesen et al., 2017) mais pas de manière aussi précise que dans mon projet. L'écart-type du ratio des lectures est un indicateur avec une bonne sensibilité mais sa spécificité doit être améliorée et la stratégie de test doit être affinée. Cette approche est très prometteuse pour détecter n'importe quelle erreur car elle n'assume aucun modèle particulier : on peut envisager une nouvelle procédure algorithmique avec seulement ce critère de cooccurrence pour déterminer la vraisemblance de chaque variant d'être erroné ou issu d'une vraie espèce. Dans ce cadre, le modèle de mutation ponctuelle est utile a posteriori pour vérifier la qualité de la détection par cooccurrence.

Enfin, la construction de la matrice d'assignation puis le calcul des indices de biodiversité permettent de revenir à la problématique initiale de ma thèse. Comme je l'ai écrit plus haut, cette approche n'a pas abouti à un résultat concret satisfaisant. Mais la redéfinition des abondances relatives pour les indices de biodiversité présente un intérêt conceptuel : il n'y a pas d'obstacle théorique à leur calcul dans un cas où les observations sont incertaines. Cette idée pourrait être reprise dans le cadre de suivis traditionnels de la biodiversité pour lesquels les données ont une tout autre forme que celles de métabarcoding.

Pour finir, je reviens sur la distinction entre algorithmes de clustering et de *denoising*, ou entre OTU et ASV. Dans ce chapitre, j'ai illustré le fait que les erreurs ne provenaient pas seulement du séquençage mais aussi de l'amplification par PCR. Ce résultat n'est pas nouveau mais permet de relativiser les conclusions des algorithmes de *denoising* qui assument parfois l'inverse. Une partie de la communauté estime que le *denoising* devrait remplacer les approches par clustering. Ce débat ne me semble pas spécialement pertinent et je rejoins Antich et al. (2021) : le choix de la méthode dépend des questions écologiques posées. Les OTU sont intéressants pour agréger l'information à l'échelle de l'espèce tandis que les ASV fournissent une plus grande précision qui n'est pas toujours utile. Les performances en termes de détection et de spécificité du *denoising* sont meilleures mais les ASV produits restent des estimations des espèces présentes, ils ne sont pas "plus vrais" en soi que les OTU.

4.3.2 Perspectives

Ce projet est encore en évolution et de nombreuses pistes peuvent encore être étudiées.

4.3.2.1 Modèle de mutation ponctuelle

Le modèle de mutation ponctuelle a été conçu pour prendre en compte les mutants en considérant simplement la substitution observée (ou l'insertion/délétion). La variabilité des ratios observés est due au contexte de ces mutations, par exemple les bases voisines ou la position dans la séquence. Certains mutants sont systématiquement surabondants par rapport à l'attendu du modèle. Leur étude serait intéressante pour mieux comprendre les mécanismes de mutation.

4.3.2.2 Choix des cooccurrences testées

Dans la version actuelle de mon algorithme, très peu de paires de séquences testées montrent un lien significatif de cooccurrence. C'est prévisible dans la mesure où le choix de ces paires est relativement arbitraire ici. On pourrait se baser sur la similarité des séquences mais cela conduirait sûrement à omettre un grand nombre de chimères.

Deux options sont envisageables : soit trouver un critère pertinent pour déterminer quelles paires de séquences sont testées, soit établir un critère de cooccurrence assez rapide à calculer pour pouvoir tester des paires à plus grande échelle.

Par ailleurs, l'écart-type du ratio semble pertinent mais peu de critères de cooccurrence ont été étudiés à l'heure actuelle. Une exploration plus complète pourrait être menée, en s'inspirant par exemple des algorithmes de regroupement de séquences (*binning*) en métagénomique qui étudient les co-occurrences des séquences, comme Arisdakessian et al. (2021).

4.3.2.3 Graphe des variants à partir des cooccurrences

L'idée a été émise de remplacer les modèles de mutation par l'étude des cooccurrences. Dans ce cadre, un graphe des variants pourrait être construit à partir d'une métrique basée sur ces cooccurrences. Ce graphe serait ensuite comparé au graphe de mutations d'*obiclean*. Plusieurs questions se posent : les graphes sont-ils similaires ? Dans l'un et l'autre, quelles sont les propriétés des séquences souches, des séquences mutées ? J'évoque deux premières idées concrètes.

D'abord, l'arité (le nombre d'arêtes d'un nœud) pourrait être étudiée. Elle est a priori révélatrice de la nature du variant : on s'attend à ce qu'une vraie séquence soit reliée à un grand nombre de variants (ses erreurs) dépendant de son abondance.

Ensuite, on peut se demander quelle serait la structure des chimères dans ce graphe. On s'attend à une dépendance à deux séquences souches : est-ce facile à caractériser ? L'écart-type du ratio est peut-être insuffisant puisque la chimère peut être présente à condition que ses deux parents soient présents. Le même problème se pose pour les mutants (de substitution par exemple) qui peuvent provenir de différentes séquences souches indépendantes : la cooccurrence du mutant avec une première souche est faussée par l'absence de prise en compte des autres souches.

4.3.2.4 Attribution probabiliste

Cette étape de la procédure de classification ne donne pas un résultat satisfaisant pour le moment. Toutefois, l'approche actuelle présente des limites potentiellement faciles à ajuster pour obtenir de meilleures performances. Je pense notamment aux distributions utilisées pour construire la matrice d'assignation. La vraisemblance choisie pour les lectures des "vrais" variants ($p(R_i|C_i = \text{vrai})$, équation 3.25) a été choisie pour sa simplicité mais elle pourrait être modifiée pour attribuer des probabilités égales à zéro en dessous d'une certaine abondance. Mais cela revient à réintroduire un seuil arbitraire d'abondance : un problème global subsiste.

4.3.2.5 Comparaison à des algorithmes de *denoising*

Dans ce projet, je me suis concentré sur l'algorithme *obiclean*. Ses performances sont reconnues mais de nombreux autres algorithmes répondent à la même problématique. À terme, il serait intéressant de comparer ma procédure de classification aux algorithmes les plus répandus, comme DADA2 ou UNOISE3. Certes, j'ai montré que les modèles reposant sur le séquençage seulement sont théoriquement incomplet mais plusieurs études montrent les bonnes performances pratiques de DADA2 notamment. Peut-être que ce défaut de modélisation induit un biais à identifier. La vitesse de calcul serait alors un autre critère de performance à considérer.

4.4 Indices de biodiversité

Pour conclure cette discussion, je fais un bref retour à la problématique de ma thèse : Comment améliorer les indices de biodiversité établis à partir des données de métabarcoding ? Comme évoqué en préambule de cette discussion, le Chapitre 2 améliore le traitement des biais tandis que le Chapitre 3 vise à réduire le bruit des données. J'ai essentiellement travaillé sur la diversité α qui correspond à l'échelle locale décrite par un échantillon d'ADN environnemental. L'analyse quantitative du Chapitre 2 a traité la variabilité interspécifique au sein d'un tel échantillon.

Mais les indices de biodiversité ne se résument pas à cette seule échelle α . L'adaptation de mes travaux à la diversité γ , à l'échelle de la métacommunauté, ne devrait pas poser de problème : les deux mesurent intrinsèquement la biodiversité dans un espace donné.

En revanche, la diversité β , importante en écologie des communautés, est mesurée de manière différente. Elle n'est pas intrinsèque mais relative à un ensemble de sites considérés. Dans ce cadre, l'étude de la variabilité intraspécifique dans les données de métabarcoding devient une information importante, alors que celle-ci était plutôt perçue comme une source de bruit dans mes données (section "Variabilité intraspécifique entre les réplicats" ci-dessus). Adapter mes résultats pour la diversité β serait une perspective intéressante.

Enfin, je rappelle que mes travaux portaient sur la diversité neutre qui postule que chaque espèce est intrinsèquement unique et équidistante de toutes les autres. En pratique, la diversité des espèces doit être prise en compte pour comprendre le fonctionnement des écosystèmes. La considération de la diversité fonctionnelle (taux de dépôt d'ADN selon les traits ?) ou phylogénétique (amplification préférentielle pour certains groupes phylogénétiques, par exemple à cause de mismatches systématiques (Liu et al., 2023)) offrirait de nouvelles perspectives d'étude pour le métabarcoding.

4.5 Conclusion

Ce manuscrit a présenté mes travaux de recherche sur les indicateurs de biodiversité établis par métabarcoding et sur les outils statistiques utilisés en écologie. De nombreux autres sujets auraient pu être abordés, bien sûr, mais je pense que ma thèse présente aujourd'hui un ensemble cohérent. Je conclus donc ici les projets menés ces trois dernières années. J'espère avoir suscité chez le lecteur ou la lectrice de la curiosité et de nouveaux questionnements !

Bibliographie

- Adams, C. I., Knapp, M., Gemmell, N. J., Jeunen, G.-J., Bunce, M., Lamare, M. D., and Taylor, H. R. (2019). Beyond Biodiversity : Can Environmental DNA (eDNA) Cut It as a Population Genetics Tool? *Genes*, 10(3) :192.
- Alberdi, A., Aizpurua, O., Gilbert, M. T. P., and Bohmann, K. (2018). Scrutinizing key steps for reliable metabarcoding of environmental samples. *Methods in Ecology and Evolution*, 9(1) :134–147.
- Alberdi, A. and Gilbert, M. T. P. (2019). A guide to the application of Hill numbers to DNA-based diversity analyses. *Molecular Ecology Resources*, 19(4) :804–817.
- Ali, M. M., Li, F., Zhang, Z., Zhang, K., Kang, D.-K., Ankrum, J. A., Le, X. C., and Zhao, W. (2014). Rolling circle amplification : a versatile tool for chemical biology, materials science and medicine. *Chemical Society Reviews*, 43(10) :3324.
- Ali, N., Rampazzo, R. D. C. P., Costa, A. D. T., and Krieger, M. A. (2017). Current Nucleic Acid Extraction Methods and Their Implications to Point-of-Care Diagnostics. *BioMed Research International*, 2017 :1–13.
- Alsos, I. G., Lavergne, S., Merkel, M. K. F., Boleda, M., Lammers, Y., Alberti, A., Pouchon, C., Denoeud, F., Pitelkova, I., Puşcaş, M., Roquet, C., Hurdu, B.-I., Thuiller, W., Zimmermann, N. E., Hollingsworth, P. M., and Coissac, E. (2020). The Treasure Vault Can be Opened : Large-Scale Genome Skimming Works Well Using Herbarium and Silica Gel Dried Material. *Plants*, 9(4) :432.
- Amir, A., McDonald, D., Navas-Molina, J. A., Kopylova, E., Morton, J. T., Zech Xu, Z., Kightley, E. P., Thompson, L. R., Hyde, E. R., Gonzalez, A., and Knight, R. (2017). Deblur Rapidly Resolves Single-Nucleotide Community Sequence Patterns. *mSystems*, 2(2) :e00191–16.
- An, Z., Nott, D. J., and Drovandi, C. (2020). Robust Bayesian synthetic likelihood via a semi-parametric approach. *Statistics and Computing*, 30(3) :543–557.
- Anderson, M. J., Crist, T. O., Chase, J. M., Vellend, M., Inouye, B. D., Freestone, A. L., Sanders, N. J., Cornell, H. V., Comita, L. S., Davies, K. F., Harrison, S. P., Kraft, N. J. B., Stegen, J. C., and Swenson, N. G. (2011). Navigating the multiple meanings of beta diversity : a roadmap for the practicing ecologist : Roadmap for beta diversity. *Ecology Letters*, 14(1) :19–28.
- Anderson, M. J., Ellingsen, K. E., and McArdle, B. H. (2006). Multivariate dispersion as a measure of beta diversity. *Ecology Letters*, 9(6) :683–693.
- Andrieu, C., Doucet, A., and Holenstein, R. (2010). Particle Markov Chain Monte Carlo Methods. *Journal of the Royal Statistical Society Series B : Statistical Methodology*, 72(3) :269–342.
- Andruszkiewicz Allan, E., Zhang, W. G., Lavery, A., and Govindarajan, A. (2021). Envi-

- ronmental DNA shedding and decay rates from diverse animal forms and thermal regimes. *Environmental DNA*, 3(2) :492–514.
- Antich, A., Palacin, C., Wangenstein, O. S., and Turon, X. (2021). To denoise or to cluster, that is not the question : optimizing pipelines for COI metabarcoding and metaphylogeography. *BMC Bioinformatics*, 22(1) :177.
- Arisdakessian, C. G., Nigro, O. D., Steward, G. F., Poisson, G., and Belcaid, M. (2021). CoCoNet : an efficient deep learning tool for viral metagenome binning. *Bioinformatics*, 37(18) :2803–2810.
- Aruoba, S. B. and Fernández-Villaverde, J. (2015). A comparison of programming languages in macroeconomics. *Journal of Economic Dynamics and Control*, 58 :265–273.
- Balmford, A., Green, R. E., and Jenkins, M. (2003). Measuring the changing state of nature. *Trends in Ecology & Evolution*, 18(7) :326–330.
- Baloğlu, B., Chen, Z., Elbrecht, V., Braukmann, T., MacDonald, S., and Steinke, D. (2021). A workflow for accurate metabarcoding using nanopore MinION sequencing. *Methods in Ecology and Evolution*, 12(5) :794–804.
- Bar, T., Kubista, M., and Tichopad, A. (2012). Validation of kinetics similarity in qPCR. *Nucleic Acids Research*, 40(4) :1395–1406.
- Barbault, R., Le Duc, J.-P., and Barbeau, P. (2005). *Biodiversité, science et gouvernance : actes de la conférence internationale, Paris, 24-28 janvier 2005*. Muséum national d’histoire naturelle, Paris.
- Barnes, M. A., Turner, C. R., Jerde, C. L., Renshaw, M. A., Chadderton, W. L., and Lodge, D. M. (2014). Environmental Conditions Influence eDNA Persistence in Aquatic Systems. *Environmental Science & Technology*, 48(3) :1819–1827.
- Barnosky, A. D., Matzke, N., Tomiya, S., Wogan, G. O. U., Swartz, B., Quental, T. B., Marshall, C., McGuire, J. L., Lindsey, E. L., Maguire, K. C., Mersey, B., and Ferrer, E. A. (2011). Has the Earth’s sixth mass extinction already arrived? *Nature*, 471(7336) :51–57.
- Bartholomew, R. A., Hutchison, J. R., Straub, T. M., and Call, D. R. (2015). PCR, Real-Time PCR, Digital PCR, and Isothermal Amplification. In Yates, M. V., Nakatsu, C. H., Miller, R. V., and Pillai, S. D., editors, *Manual of Environmental Microbiology*, pages 2.3.2–1–2.3.2–13. ASM Press, Washington, DC, USA.
- Bartholomew-Biggs, M., Brown, S., Christianson, B., and Dixon, L. (2000). Automatic differentiation of algorithms. *Journal of Computational and Applied Mathematics*, 124(1-2) :171–190.
- Beaumont, M. A., Zhang, W., and Balding, D. J. (2002). Approximate Bayesian Computation in Population Genetics. *Genetics*, 162(4) :2025–2035.
- Bell, K. L., Fowler, J., Burgess, K. S., Dobbs, E. K., Gruenewald, D., Lawley, B., Morozumi, C., and Brosi, B. J. (2017). Applying Pollen DNA Metabarcoding to the Study of Plant–Pollinator Interactions. *Applications in Plant Sciences*, 5(6) :1600124.
- Beng, K. C. and Corlett, R. T. (2020). Applications of environmental DNA (eDNA) in ecology and conservation : opportunities, challenges and prospects. *Biodiversity and Conservation*, 29(7) :2089–2121.
- Bernton, E., Jacob, P. E., Gerber, M., and Robert, C. P. (2019). Approximate Bayesian computation with the Wasserstein distance. *Journal of the Royal Statistical Society : Series B (Statistical Methodology)*, 81(2) :235–269.
- Bezanson, J., Edelman, A., Karpinski, S., and Shah, V. B. (2017). Julia : A Fresh Approach

- to Numerical Computing. *SIAM Review*, 59(1) :65–98.
- Bienert, F., De Danieli, S., Miquel, C., Coissac, E., Poillot, C., Brun, J.-J., and Taberlet, P. (2012). Tracking earthworm communities from soil dna. *Molecular Ecology*, 21(8) :2017–2030.
- Bjørnsgaard Aas, A., Davey, M. L., and Kauserud, H. (2017). ITS all right mama : investigating the formation of chimeric sequences in the ITS2 region by DNA metabarcoding analyses of fungal mock communities of different complexities. *Molecular Ecology Resources*, 17(4) :730–741.
- Blum, M. G. B. and François, O. (2010). Non-linear regression models for Approximate Bayesian Computation. *Statistics and Computing*, 20(1) :63–73.
- Boggy, G. J. and Woolf, P. J. (2010). A Mechanistic Model of PCR for Accurate Quantification of Quantitative PCR Data. *PLoS ONE*, 5(8) :e12355.
- Bohmann, K., Evans, A., Gilbert, M. T. P., Carvalho, G. R., Creer, S., Knapp, M., Yu, D. W., and de Bruyn, M. (2014). Environmental DNA for wildlife biology and biodiversity monitoring. *Trends in Ecology & Evolution*, 29(6) :358–367.
- Bollback, J. P., York, T. L., and Nielsen, R. (2008). Estimation of 2Ne s From Temporal Allele Frequency Data. *Genetics*, 179(1) :497–502.
- Bolyen, E., Rideout, J. R., Dillon, M. R., Bokulich, N. A., Abnet, C. C., Al-Ghalith, G. A., Alexander, H., Alm, E. J., Arumugam, M., Asnicar, F., Bai, Y., Bisanz, J. E., Bittinger, K., Brejnrod, A., Brislawn, C. J., Brown, C. T., Callahan, B. J., Caraballo-Rodríguez, A. M., Chase, J., Cope, E. K., Da Silva, R., Diener, C., Dorrestein, P. C., Douglas, G. M., Durall, D. M., Duvallet, C., Edwardson, C. F., Ernst, M., Estaki, M., Fouquier, J., Gauglitz, J. M., Gibbons, S. M., Gibson, D. L., Gonzalez, A., Gorlick, K., Guo, J., Hillmann, B., Holmes, S., Holste, H., Huttenhower, C., Huttley, G. A., Janssen, S., Jarmusch, A. K., Jiang, L., Kaehler, B. D., Kang, K. B., Keefe, C. R., Keim, P., Kelley, S. T., Knights, D., Koester, I., Kosciulek, T., Kreps, J., Langille, M. G. I., Lee, J., Ley, R., Liu, Y.-X., Loftfield, E., Lozupone, C., Maher, M., Marotz, C., Martin, B. D., McDonald, D., McIver, L. J., Melnik, A. V., Metcalf, J. L., Morgan, S. C., Morton, J. T., Naimey, A. T., Navas-Molina, J. A., Nothias, L. F., Orchanian, S. B., Pearson, T., Peoples, S. L., Petras, D., Preuss, M. L., Priesse, E., Rasmussen, L. B., Rivers, A., Robeson, M. S., Rosenthal, P., Segata, N., Shaffer, M., Shiffer, A., Sinha, R., Song, S. J., Spear, J. R., Swafford, A. D., Thompson, L. R., Torres, P. J., Trinh, P., Tripathi, A., Turnbaugh, P. J., Ul-Hasan, S., Van Der Hooft, J. J. J., Vargas, F., Vázquez-Baeza, Y., Vogtmann, E., Von Hippel, M., Walters, W., Wan, Y., Wang, M., Warren, J., Weber, K. C., Williamson, C. H. D., Willis, A. D., Xu, Z. Z., Zaneveld, J. R., Zhang, Y., Zhu, Q., Knight, R., and Caporaso, J. G. (2019). Reproducible, interactive, scalable and extensible microbiome data science using QIIME 2. *Nature Biotechnology*, 37(8) :852–857.
- Bonin, A., Guerrieri, A., and Ficetola, G. F. (2023). Optimal sequence similarity thresholds for clustering of molecular operational taxonomic units in DNA metabarcoding studies. *Molecular Ecology Resources*, 23(2) :368–381.
- Bou Dagher-Kharrat, M., Abdel-Samad, N., Douaihy, B., Bourge, M., Fridlender, A., Siljak-Yakovlev, S., and Brown, S. C. (2013). Nuclear DNA C-values for biodiversity screening : Case of the Lebanese flora. *Plant Biosystems - An International Journal Dealing with all Aspects of Plant Biology*, 147(4) :1228–1237.
- Box, G. E. P. and Muller, M. E. (1958). A Note on the Generation of Random Normal Deviates. *The Annals of Mathematical Statistics*, 29(2) :610–611.

- Boyer, F., Mercier, C., Bonin, A., Le Bras, Y., Taberlet, P., and Coissac, E. (2016). obitools : a unix-inspired software package for DNA metabarcoding. *Molecular Ecology Resources*, 16(1) :176–182.
- Bray, J. R. and Curtis, J. T. (1957). An Ordination of the Upland Forest Communities of Southern Wisconsin. *Ecological Monographs*, 27(4) :325–349.
- Brent, R. P. (2002). *Algorithms for minimization without derivatives*. Dover books on mathematics. Dover Publications, Mineola, NY, unabridged republication of the work publ. by prentice-hall ... 1973 edition.
- Broms, K. M., Hooten, M. B., and Fitzpatrick, R. M. (2015). Accounting for imperfect detection in Hill numbers for biodiversity studies. *Methods in Ecology and Evolution*, 6(1) :99–108.
- Browne, H. P., Forster, S. C., Anonye, B. O., Kumar, N., Neville, B. A., Stares, M. D., Goulding, D., and Lawley, T. D. (2016). Culturing of ‘unculturable’ human microbiota reveals novel taxa and extensive sporulation. *Nature*, 533(7604) :543–546.
- Brys, R., Halfmaerten, D., Neyrinck, S., Mauvisseau, Q., Auwerx, J., Sweet, M., and Mergeay, J. (2021). Reliable eDNA detection and quantification of the European weather loach (*Misgurnus fossilis*). *Journal of Fish Biology*, 98(2) :399–414.
- Burian, A., Mauvisseau, Q., Bulling, M., Domisch, S., Qian, S., and Sweet, M. (2021). Improving the reliability of eDNA data interpretation. *Molecular Ecology Resources*, 21(5) :1422–1433.
- Byrd, R. H., Lu, P., Nocedal, J., and Zhu, C. (1995). A Limited Memory Algorithm for Bound Constrained Optimization. *SIAM Journal on Scientific Computing*, 16(5) :1190–1208.
- Calderón-Sanou, I., Münkemüller, T., Boyer, F., Zinger, L., and Thuiller, W. (2020). From environmental DNA sequences to ecological conclusions : How strong is the influence of methodological choices? *Journal of Biogeography*, 47(1) :193–206.
- Callahan, B. J., McMurdie, P. J., and Holmes, S. P. (2017). Exact sequence variants should replace operational taxonomic units in marker-gene data analysis. *The ISME Journal*, 11(12) :2639–2643.
- Callahan, B. J., McMurdie, P. J., Rosen, M. J., Han, A. W., Johnson, A. J. A., and Holmes, S. P. (2016). DADA2 : High-resolution sample inference from Illumina amplicon data. *Nature Methods*, 13(7) :581–583.
- Cao, Y., Yu, M., Dong, G., Chen, B., and Zhang, B. (2020). Digital PCR as an Emerging Tool for Monitoring of Microbial Biodegradation. *Molecules*, 25(3) :706.
- Carr, A. C. and Moore, S. D. (2012). Robust Quantification of Polymerase Chain Reactions Using Global Fitting. *PLoS ONE*, 7(5) :e37640.
- Català, S., Berbegal, M., Pérez-Sierra, A., and Abad-Campos, P. (2017). Metabarcoding and development of new real-time specific assays reveal *Phytophthora* species diversity in holm oak forests in eastern Spain. *Plant Pathology*, 66(1) :115–123.
- Chandelier, A., Hulin, J., San Martin, G., Debode, F., and Massart, S. (2021). Comparison of qPCR and Metabarcoding Methods as Tools for the Detection of Airborne Inoculum of Forest Fungal Pathogens. *Phytopathology*®, 111(3) :570–581.
- Chao, A., Chiu, C.-H., and Hsieh, T. C. (2012). Proposing a resolution to debates on diversity partitioning. *Ecology*, 93(9) :2037–2051.
- Chao, A., Chiu, C.-H., and Jost, L. (2014). Unifying Species Diversity, Phylogenetic Diversity, Functional Diversity, and Related Similarity and Differentiation Measures Through Hill

- Numbers. *Annual Review of Ecology, Evolution, and Systematics*, 45(1) :297–324.
- Chao, A. and Jost, L. (2015). Estimating diversity and entropy profiles via discovery rates of new species. *Methods in Ecology and Evolution*, 6(8) :873–882.
- Chen, S., Yao, H., Han, J., Liu, C., Song, J., Shi, L., Zhu, Y., Ma, X., Gao, T., Pang, X., Luo, K., Li, Y., Li, X., Jia, X., Lin, Y., and Leon, C. (2010). Validation of the ITS2 Region as a Novel DNA Barcode for Identifying Medicinal Plant Species. *PLoS ONE*, 5(1) :e8613.
- Chervoneva, I., Li, Y., Iglewicz, B., Waldman, S., and Hyslop, T. (2007). Relative quantification based on logistic models for individual polymerase chain reactions. *Statistics in Medicine*, 26(30) :5596–5611.
- Chiu, C.-H. and Chao, A. (2016). Estimating and comparing microbial diversity in the presence of sequencing errors. *PeerJ*, 4 :e1634.
- Clarke, L. J., Soubrier, J., Weyrich, L. S., and Cooper, A. (2014). Environmental metabarcodes for insects : *in silico* PCR reveals potential for taxonomic bias. *Molecular Ecology Resources*, 14(6) :1160–1170.
- Coissac, E., Hollingsworth, P. M., Lavergne, S., and Taberlet, P. (2016). From barcodes to genomes : extending the concept of DNA barcoding. *Molecular Ecology*, 25(7) :1423–1428.
- Condit, R., Chisholm, R. A., and Hubbell, S. P. (2012). Thirty Years of Forest Census at Barro Colorado and the Importance of Immigration in Maintaining Diversity. *PLoS ONE*, 7(11) :e49826.
- Cracraft, J. (1983). Species concepts and speciation analysis. In Johnston, R. F., editor, *Current Ornithology*, pages 159–87. Plenum Press.
- Cranmer, K., Brehmer, J., and Louppe, G. (2020). The frontier of simulation-based inference. *Proceedings of the National Academy of Sciences*, 117(48) :30055–30062.
- Csilléry, K., Blum, M. G., Gaggiotti, O. E., and François, O. (2010). Approximate Bayesian Computation (ABC) in practice. *Trends in Ecology & Evolution*, 25(7) :410–418.
- Cuff, J. P., Kitson, J. J. N., Hemprich-Bennett, D., Tercel, M. P. T. G., Browett, S. S., and Evans, D. M. (2023). The predator problem and PCR primers in molecular dietary analysis : Swamped or silenced ; depth or breadth ? *Molecular Ecology Resources*, 23(1) :41–51.
- Daróczy, Z. (1970). Generalized information functions. *Information and Control*, 16(1) :36–51.
- Dean, T. A., Singh, S. S., Jasra, A., and Peters, G. W. (2014). Parameter Estimation for Hidden Markov Models with Intractable Likelihoods : Estimating intractable HMMs. *Scandinavian Journal of Statistics*, 41(4) :970–987.
- Deiner, K., Bik, H. M., Mächler, E., Seymour, M., Lacoursière-Roussel, A., Altermatt, F., Creer, S., Bista, I., Lodge, D. M., Vere, N., Pfrender, M. E., and Bernatchez, L. (2017). Environmental DNA metabarcoding : Transforming how we survey animal and plant communities. *Molecular Ecology*, 26(21) :5872–5895.
- Del Moral, P., Doucet, A., and Jasra, A. (2006). Sequential Monte Carlo Samplers. *Journal of the Royal Statistical Society Series B : Statistical Methodology*, 68(3) :411–436.
- Del Moral, P., Doucet, A., and Jasra, A. (2012). An adaptive sequential Monte Carlo method for approximate Bayesian computation. *Statistics and Computing*, 22(5) :1009–1020.
- Dempster, A. P., Laird, N. M., and Rubin, D. B. (1977). Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society : Series B*, 39 :1–38.
- deWaard, J. R., Levesque-Beaudin, V., deWaard, S. L., Ivanova, N. V., McKeown, J. T., Miskie, R., Naik, S., Perez, K. H., Ratnasingham, S., Sobel, C. N., Sones, J. E., Steinke,

- C., Telfer, A. C., Young, A. D., Young, M. R., Zakharov, E. V., and Hebert, P. D. (2019). Expedited assessment of terrestrial arthropod diversity by coupling Malaise traps with DNA barcoding. *Genome*, 62(3) :85–95.
- Doi, H., Uchii, K., Takahara, T., Matsushashi, S., Yamanaka, H., and Minamoto, T. (2015). Use of Droplet Digital PCR for Estimation of Fish Abundance and Biomass in Environmental DNA Surveys. *PLOS ONE*, 10(3) :e0122763.
- Dopheide, A., Xie, D., Buckley, T. R., Drummond, A. J., and Newcomb, R. D. (2019). Impacts of DNA extraction and PCR on DNA metabarcoding estimates of soil biodiversity. *Methods in Ecology and Evolution*, 10(1) :120–133.
- Doucet, A. and Johansen, A. M. (2008). A tutorial on particle filtering and smoothing : Fifteen years later.
- Doyle, J. J. (1990). Isolation of plant dna from fresh tissue.
- Drovandi, C. and Frazier, D. T. (2022). A comparison of likelihood-free methods with and without summary statistics. *Statistics and Computing*, 32(3) :42.
- Drovandi, C. C. and Pettitt, A. N. (2011). Likelihood-free Bayesian estimation of multivariate quantile distributions. *Computational Statistics & Data Analysis*, 55(9) :2541–2556.
- Edgar, R. C. (2010). Search and clustering orders of magnitude faster than BLAST. *Bioinformatics*, 26(19) :2460–2461.
- Edgar, R. C. (2013). UPARSE : highly accurate OTU sequences from microbial amplicon reads. *Nature Methods*, 10(10) :996–998.
- Edgar, R. C. (2016a). UCHIME2 : improved chimera prediction for amplicon sequencing. preprint, Bioinformatics.
- Edgar, R. C. (2016b). UNOISE2 : improved error-correction for Illumina 16S and ITS amplicon sequencing. preprint, Bioinformatics.
- Eichmiller, J. J., Best, S. E., and Sorensen, P. W. (2016). Effects of Temperature and Trophic State on Degradation of Environmental DNA in Lake Water. *Environmental Science & Technology*, 50(4) :1859–1867.
- Eisenhofer, R. and Weyrich, L. (2018). Proper Authentication of Ancient DNA Is Still Essential. *Genes*, 9(3) :122.
- Elbrecht, V., Bourlat, S. J., Hörren, T., Lindner, A., Mordente, A., Noll, N. W., Schäffler, L., Sorg, M., and Zizka, V. M. (2021). Pooling size sorted Malaise trap fractions to maximize taxon recovery with metabarcoding. *PeerJ*, 9 :e12177.
- Elbrecht, V. and Leese, F. (2015). Can DNA-Based Ecosystem Assessments Quantify Species Abundance ? Testing Primer Bias and Biomass—Sequence Relationships with an Innovative Metabarcoding Protocol. *PLOS ONE*, 10(7) :e0130324.
- Elbrecht, V., Peinert, B., and Leese, F. (2017). Sorting things out : Assessing effects of unequal specimen biomass on DNA metabarcoding. *Ecology and Evolution*, 7(17) :6918–6926.
- Ellison, A. M. (2010). Partitioning diversity ¹. *Ecology*, 91(7) :1962–1963.
- Eren, A. M., Morrison, H. G., Lescault, P. J., Reveillaud, J., Vineis, J. H., and Sogin, M. L. (2015). Minimum entropy decomposition : Unsupervised oligotyping for sensitive partitioning of high-throughput marker gene sequences. *The ISME Journal*, 9(4) :968–979.
- Ershova, E. A., Wangensteen, O. S., Descoteaux, R., Barth-Jensen, C., and Præbel, K. (2021). Metabarcoding as a quantitative tool for estimating biodiversity and relative biomass of marine zooplankton. *ICES Journal of Marine Science*, 78(9) :3342–3355.

- Espinosa-Prieto, A., Beisel, J., Verschuren, P., and Hardion, L. (2023). Toward freshwater plant diversity surveys with eDNA barcoding and metabarcoding. *Environmental DNA*, page edn3.407.
- Estensmo, E. L. F., Maurice, S., Morgado, L., Martin-Sanchez, P. M., Skrede, I., and Kause-rud, H. (2021). The influence of intraspecific sequence variation during DNA metabar-coding : A case study of eleven fungal species. *Molecular Ecology Resources*, 21(4) :1141–1148.
- Evans, M. and Swartz, T. (1995). Methods for Approximating Integrals in Statistics with Special Emphasis on Bayesian Integration Problems. *Statistical Science*, 10(3).
- Ewens, W. J. (2004). *Mathematical Population Genetics*, volume 27 of *Interdisciplinary Ap-plied Mathematics*. Springer New York, New York, NY.
- Fahner, N. A., Shokralla, S., Baird, D. J., and Hajibabaei, M. (2016). Large-Scale Monitoring of Plants through Environmental DNA Metabarcoding of Soil : Recovery, Resolution, and Annotation of Four DNA Markers. *PLOS ONE*, 11(6) :e0157505.
- Fasiolo, M., Pya, N., and Wood, S. N. (2016). A Comparison of Inferential Methods for Highly Nonlinear State Space Models in Ecology and Epidemiology. *Statistical Science*, 31(1).
- Fearnhead, P. and Prangle, D. (2012). Constructing summary statistics for approximate Baye-sian computation : semi-automatic approximate Bayesian computation : Semi-automatic Approximate Bayesian Computation. *Journal of the Royal Statistical Society : Series B (Statistical Methodology)*, 74(3) :419–474.
- Ficetola, G., Coissac, E., Zundel, S., Riaz, T., Shehzad, W., Bessièrè, J., Taberlet, P., and Pompanon, F. (2010). An In silico approach for the evaluation of DNA barcodes. *BMC Genomics*, 11(1) :434.
- Ficetola, G. F., Pansu, J., Bonin, A., Coissac, E., Giguët-Covex, C., De Barba, M., Gielly, L., Lopes, C. M., Boyer, F., Pompanon, F., Rayé, G., and Taberlet, P. (2015). Replication levels, false presences and the estimation of the presence/absence from eDNA metabarcoding data. *Molecular Ecology Resources*, 15(3) :543–556.
- Ficetola, G. F. and Taberlet, P. (2023). Towards exhaustive community ecology via DNA metabarcoding. *Molecular Ecology*, page mec.16881.
- Fisher, R. A. (1935). *Statistical Methods for Research Workers* : By R. A. Fisher. Edinburgh : Oliver ... Boyd. Ed. 5. XIII+319 pages and supplementary tables. Illus. 1934. 15 / net. *Agronomy Journal*, 27(1) :76–76.
- Foll, M., Shim, H., and Jensen, J. D. (2015). WFABC : a Wright-Fisher ABC-based approach for inferring effective population sizes and selection coefficients from time-sampled data. *Molecular Ecology Resources*, 15(1) :87–98.
- Fonseca, V. G. (2018). “Pitfalls in relative abundance estimation using eDNA metabarcoding”. *Molecular Ecology Resources*, 18(5) :923–926.
- Forsythe, G. E. (1972). Von Neumann’s Comparison Method for Random Sampling from the Normal and Other Distributions. *Mathematics of Computation*, 26(120) :817.
- Frøslev, T. G., Kjølner, R., Bruun, H. H., Ejrnæs, R., Brunbjerg, A. K., Pietroni, C., and Hansen, A. J. (2017). Algorithm for post-clustering curation of DNA amplicon data yields reliable biodiversity estimates. *Nature Communications*, 8(1) :1188.
- Fujii, K., Doi, H., Matsuoka, S., Nagano, M., Sato, H., and Yamanaka, H. (2019). Environ-mental DNA metabarcoding for fish community analysis in backwater lakes : A comparison of capture methods. *PLOS ONE*, 14(1) :e0210357.
- Gao, F. and Han, L. (2012). Implementing the Nelder-Mead simplex algorithm with adaptive

- parameters. *Computational Optimization and Applications*, 51(1) :259–277.
- Garrido-Sanz, L., Senar, M. A., and Piñol, J. (2022). Relative species abundance estimation in artificial mixtures of insects using mito-metagenomics and a correction factor for the mitochondrial DNA copy number. *Molecular Ecology Resources*, 22(1) :153–167.
- Geman, S. and Geman, D. (1987). Stochastic Relaxation, Gibbs Distributions, and the Bayesian Restoration of Images. In *Readings in Computer Vision*, pages 564–584. Elsevier.
- Gevertz, J. L., Dunn, S. M., and Roth, C. M. (2005). Mathematical model of real-time PCR kinetics. *Biotechnology and Bioengineering*, 92(3) :346–355.
- Giguet-Covex, C., Pansu, J., Arnaud, F., Rey, P.-J., Griggo, C., Gielly, L., Domaizon, I., Coissac, E., David, F., Choler, P., Poulénard, J., and Taberlet, P. (2014). Long livestock farming history and human landscape shaping revealed by lake sediment DNA. *Nature Communications*, 5(1) :3211.
- Gill, P., Bleka, O., and Fonnelløp, A. E. (2022). Limitations of qPCR to estimate DNA quantity : An RFU method to facilitate inter-laboratory comparisons for activity level, and general applicability. *Forensic Science International : Genetics*, 61 :102777.
- Gill, P. and Ghaemi, A. (2008). Nucleic Acid Isothermal Amplification Technologies—A Review. *Nucleosides, Nucleotides and Nucleic Acids*, 27(3) :224–243.
- Giovannoni, S. J., Britschgi, T. B., Moyer, C. L., and Field, K. G. (1990). Genetic diversity in Sargasso Sea bacterioplankton. *Nature*, 345(6270) :60–63.
- Glassman, S. I. and Martiny, J. B. H. (2018). Broadscale Ecological Patterns Are Robust to Use of Exact Sequence Variants versus Operational Taxonomic Units. *mSphere*, 3(4) :e00148–18.
- Gohl, D. M., Auch, B., Certano, A., LeFrançois, B., Bouevitch, A., Doukhanine, E., Fragel, C., Macklaim, J., Hollister, E., Garbe, J., and Beckman, K. B. (2021). Dissecting and tuning primer editing by proofreading polymerases. *Nucleic Acids Research*, 49(15) :e87–e87.
- Golczyk, H., Greiner, S., Wanner, G., Weihe, A., Bock, R., Börner, T., and Herrmann, R. G. (2014). Chloroplast DNA in Mature and Senescing Leaves : A Reappraisal. *The Plant Cell*, 26(3) :847–854.
- Gold, Z., Kelly, R. P., Shelton, A. O., Thompson, A. R., Goodwin, K. D., Gallego, R., Parsons, K. M., Thompson, L. R., Kacev, D., and Barber, P. H. (2023a). Archived dna reveals marine heatwave-associated shifts in fish assemblages. *Environmental DNA*, page edn3.400.
- Gold, Z., Shelton, A. O., Casendino, H. R., Duprey, J., Gallego, R., Van Cise, A., Fisher, M., Jensen, A. J., D’Agnese, E., Andruszkiewicz Allan, E., Ramón-Laca, A., Garber-Yonts, M., Labare, M., Parsons, K. M., and Kelly, R. P. (2023b). Signal and noise in metabarcoding data. *PLOS ONE*, 18(5) :e0285674.
- Goldberg, C. S., Pilliod, D. S., Arkle, R. S., and Waits, L. P. (2011). Molecular Detection of Vertebrates in Stream Water : A Demonstration Using Rocky Mountain Tailed Frogs and Idaho Giant Salamanders. *PLoS ONE*, 6(7) :e22746.
- Goll, R., Olsen, T., Cui, G., and Florholmen, J. (2006). Evaluation of absolute quantitation by nonlinear regression in probe-based real-time PCR. *BMC Bioinformatics*, 7(1) :107.
- Good, I. J. (1953). The population frequencies of species and the estimation of population parameters. *Biometrika*, 40(3-4) :237–264.
- Gordon, N., Salmond, D., and Smith, A. (1993). Novel approach to nonlinear/non-Gaussian Bayesian state estimation. *IEE Proceedings F Radar and Signal Processing*, 140(2) :107.

- Goss-Souza, D., Borges, C. D., Mendes, L. W., Aparecido Navarrete, A., Baretta, D., and Siu Mui, T. (2016). 5. Exploring Diversity of Soil Microorganisms : A Multidimensional Approach. In Gheler-Costa, C., Lyra-Jorge, M. C., and Martins Verdade, L., editors, *Biodiversity in Agricultural Landscapes of Southeastern Brazil*, pages 66–86. De Gruyter Open.
- Gottschalk, P. G. and Dunn, J. R. (2005). The five-parameter logistic : A characterization and comparison with the four-parameter logistic. *Analytical Biochemistry*, 343(1) :54–65.
- Griffin, J. E., Matechou, E., Buxton, A. S., Bormpoudakis, D., and Griffiths, R. A. (2020). Modelling Environmental DNA Data ; Bayesian Variable Selection Accounting for False Positive and False Negative Errors. *Journal of the Royal Statistical Society Series C : Applied Statistics*, 69(2) :377–392.
- Guitet, S., Sabatier, D., Brunaux, O., Hérault, B., Aubry-Kientz, M., Molino, J.-F., and Baraloto, C. (2014). Estimating tropical tree diversity indices from forestry surveys : A method to integrate taxonomic uncertainty. *Forest Ecology and Management*, 328 :270–281.
- Gustafson, P. (2014). Bayesian inference in partially identified models : Is the shape of the posterior distribution useful? *Electronic Journal of Statistics*, 8(1).
- Gál, A. B., Carnwath, J. W., Dinnyes, A., Herrmann, D., Niemann, H., and Wrenzycki, C. (2006). Comparison of real-time polymerase chain reaction and end-point polymerase chain reaction for the analysis of gene expression in preimplantation embryos. *Reproduction, Fertility and Development*, 18(3) :365.
- Haile, J., Froese, D. G., MacPhee, R. D. E., Roberts, R. G., Arnold, L. J., Reyes, A. V., Rasmussen, M., Nielsen, R., Brook, B. W., Robinson, S., Demuro, M., Gilbert, M. T. P., Munch, K., Austin, J. J., Cooper, A., Barnes, I., Möller, P., and Willerslev, E. (2009). Ancient DNA reveals late survival of mammoth and horse in interior Alaska. *Proceedings of the National Academy of Sciences*, 106(52) :22352–22357.
- Hardinge, P. and Murray, J. A. H. (2020). Full Dynamic Range Quantification using Loop-mediated Amplification (LAMP) by Combining Analysis of Amplification Timing and Variance between Replicates at Low Copy Number. *Scientific Reports*, 10(1) :916.
- Harper, L. R., Lawson Handley, L., Hahn, C., Boonham, N., Rees, H. C., Gough, K. C., Lewis, E., Adams, I. P., Brotherton, P., Phillips, S., and Hänfling, B. (2018). Needle in a haystack? A comparison of eDNA metabarcoding and targeted qPCR for detection of the great crested newt (*Triturus cristatus*). *Ecology and Evolution*, 8(12) :6330–6341.
- Harrison, J. G., John Calder, W., Shuman, B., and Alex Buerkle, C. (2021). The quest for absolute abundance : The use of internal standards for DNA-based community ecology. *Molecular Ecology Resources*, 21(1) :30–43.
- Hartig, F., Calabrese, J. M., Reineking, B., Wiegand, T., and Huth, A. (2011). Statistical inference for stochastic simulation models - theory and application : Inference for stochastic simulation models. *Ecology Letters*, 14(8) :816–827.
- Hastings, W. K. (1970). Monte Carlo sampling methods using Markov chains and their applications. *Biometrika*, 57(1) :97–109.
- Havrda, J., Charvát, F., Havrda, J., and Charvát, F. (1967). Quantification method of classification processes : Concept of structural a-entropy. *Kybernetika*.
- Hayward, A. (1998). Modeling and analysis of competitive RT-PCR. *Nucleic Acids Research*, 26(11) :2511–2518.
- Hebert, P. D. N., Cywinska, A., Ball, S. L., and deWaard, J. R. (2003). Biological identi-

- fications through DNA barcodes. *Proceedings of the Royal Society of London. Series B : Biological Sciences*, 270(1512) :313–321.
- Hey, J. (2001). The mind of the species problem. *Trends in Ecology & Evolution*, 16(7) :326–329.
- Higuchi, R., Fockler, C., Dollinger, G., and Watson, R. (1993). Kinetic PCR Analysis : Real-time Monitoring of DNA Amplification Reactions. *Nature Biotechnology*, 11(9) :1026–1030.
- Hill, M. O. (1973). Diversity and Evenness : A Unifying Notation and Its Consequences. *Ecology*, 54(2) :427–432.
- Hollingsworth, P. M. (2011). Refining the DNA barcode for land plants. *Proceedings of the National Academy of Sciences*, 108(49) :19451–19452.
- Hoshino, T. and Inagaki, F. (2017). Application of Stochastic Labeling with Random-Sequence Barcodes for Simultaneous Quantification and Sequencing of Environmental 16S rRNA Genes. *PLOS ONE*, 12(1) :e0169431.
- Hoshino, T., Nakao, R., Doi, H., and Minamoto, T. (2021). Simultaneous absolute quantification and sequencing of fish environmental DNA in a mesocosm by quantitative sequencing technique. *Scientific Reports*, 11(1) :4372.
- Hsieh, T. C., Ma, K. H., and Chao, A. (2016). iNEXT : an R package for rarefaction and extrapolation of species diversity (H numbers). *Methods in Ecology and Evolution*, 7(12) :1451–1456.
- Huntley, M. A. and Golding, G. B. (2006). Selection and Slippage Creating Serine Homopolymers. *Molecular Biology and Evolution*, 23(11) :2017–2025.
- Hønsvall, B. K. and Robertson, L. J. (2017). From research lab to standard environmental analysis tool : Will NASBA make the leap? *Water Research*, 109 :389–397.
- Ionides, E. L., Bhadra, A., Atchadé, Y., and King, A. (2011). Iterated filtering. *The Annals of Statistics*, 39(3).
- IPBES (2019). Summary for policymakers of the global assessment report on biodiversity and ecosystem services. Technical report, Zenodo. Version Number : summary for policy makers.
- Iwazkiewicz-Eggebrecht, E., Granqvist, E., Buczek, M., Prus, M., Kudlicka, J., Roslin, T., Tack, A. J. M., Andersson, A. F., Miraldo, A., Ronquist, F., and Łukasik, P. (2023). Optimizing insect metabarcoding using replicated mock communities. *Methods in Ecology and Evolution*, 14(4) :1130–1146.
- Jagers, P. and Klebaner, F. (2003). Random variation and concentration effects in PCR. *Journal of Theoretical Biology*, 224(3) :299–304.
- Jane, S. F., Wilcox, T. M., McKelvey, K. S., Young, M. K., Schwartz, M. K., Lowe, W. H., Letcher, B. H., and Whiteley, A. R. (2015). Distance, flow and PCR inhibition : eDNA dynamics in two headwater streams. *Molecular Ecology Resources*, 15(1) :216–227.
- Ji, Y., Huotari, T., Roslin, T., Schmidt, N. M., Wang, J., Yu, D. W., and Ovaskainen, O. (2020). SPIKEPIPE : A metagenomic pipeline for the accurate quantification of eukaryotic species occurrences and intraspecific abundance change using DNA barcodes or mitogenomes. *Molecular Ecology Resources*, 20(1) :256–267.
- Jost, L. (2007). Partitioning diversity into independent alpha and beta components. *Ecology*, 88(10) :2427–2439.
- Jukes, T. H. and Cantor, C. R. (1969). Evolution of Protein Molecules. In *Mammalian Protein*

- Metabolism*, pages 21–132. Elsevier.
- Jurburg, S. D., Keil, P., Singh, B. K., and Chase, J. M. (2021). All together now : Limitations and recommendations for the simultaneous analysis of all eukaryotic soil sequences. *Molecular Ecology Resources*, 21(6) :1759–1771.
- K Mogensen, P. and N Riseth, A. (2018). Optim : A mathematical optimization package for Julia. *Journal of Open Source Software*, 3(24) :615.
- Kalendar, R. (2022). A Guide to Using FASTPCR Software for PCR, In Silico PCR, and Oligonucleotide Analysis. In Basu, C., editor, *PCR Primer Design*, volume 2392, pages 223–243. Springer US, New York, NY. Series Title : Methods in Molecular Biology.
- Karlen, Y., McNair, A., Perseguers, S., Mazza, C., and Mermod, N. (2007). Statistical significance of quantitative PCR. *BMC Bioinformatics*, 8(1) :131.
- Kelly, R. P., Shelton, A. O., and Gallego, R. (2019). Understanding PCR Processes to Draw Meaningful Conclusions from Environmental DNA Studies. *Scientific Reports*, 9(1) :12133.
- Kembel, S. W., Wu, M., Eisen, J. A., and Green, J. L. (2012). Incorporating 16S Gene Copy Number Information Improves Estimates of Microbial Diversity and Abundance. *PLoS Computational Biology*, 8(10) :e1002743.
- Kestel, J. H., Field, D. L., Bateman, P. W., White, N. E., Allentoft, M. E., Hopkins, A. J., Gibberd, M., and Nevill, P. (2022). Applications of environmental DNA (eDNA) in agricultural systems : Current uses, limitations and future prospects. *Science of The Total Environment*, 847 :157556.
- Khare, V. and Eckert, K. A. (2002). The proofreading 3' -> 5' exonuclease activity of DNA polymerases : a kinetic barrier to translesion DNA synthesis. *Mutation Research/Fundamental and Molecular Mechanisms of Mutagenesis*, 510(1-2) :45–54.
- Kimura, M. (1980). A simple method for estimating evolutionary rates of base substitutions through comparative studies of nucleotide sequences. *Journal of Molecular Evolution*, 16(2) :111–120.
- Kindt, R., Van Damme, P., and Simons, A. J. (2006). Tree Diversity in Western Kenya : Using Profiles to Characterise Richness and Evenness. *Biodiversity and Conservation*, 15(4) :1253–1270.
- Kjær, K. H., Winther Pedersen, M., De Sanctis, B., De Cahsan, B., Korneliussen, T. S., Michelsen, C. S., Sand, K. K., Jelavić, S., Ruter, A. H., Schmidt, A. M. A., Kjeldsen, K. K., Tesakov, A. S., Snowball, I., Gosse, J. C., Alsos, I. G., Wang, Y., Dockter, C., Rasmussen, M., Jørgensen, M. E., Skadhauge, B., Prohaska, A., Kristensen, J. A., Bjerager, M., Allentoft, M. E., Coissac, E., PhyloNorway Consortium, Alsos, I. G., Coissac, E., Rouillard, A., Simakova, A., Fernandez-Guerra, A., Bowler, C., Macias-Fauria, M., Vinner, L., Welch, J. J., Hidy, A. J., Sikora, M., Collins, M. J., Durbin, R., Larsen, N. K., and Willerslev, E. (2022). A 2-million-year-old ecosystem in Greenland uncovered by environmental DNA. *Nature*, 612(7939) :283–291.
- Klepke, M. J., Sigsgaard, E. E., Jensen, M. R., Olsen, K., and Thomsen, P. F. (2022). Accumulation and diversity of airborne, eukaryotic environmental DNA. *Environmental DNA*, page edn3.340.
- Klymus, K. E., Marshall, N. T., and Stepien, C. A. (2017). Environmental DNA (eDNA) metabarcoding assays to detect invasive invertebrate species in the Great Lakes. *PLOS ONE*, 12(5) :e0177643.
- Klymus, K. E., Richter, C. A., Chapman, D. C., and Paukert, C. (2015). Quantification of

- eDNA shedding rates from invasive bighead carp *Hypophthalmichthys nobilis* and silver carp *Hypophthalmichthys molitrix*. *Biological Conservation*, 183 :77–84.
- Kobyzev, I., Prince, S. J., and Brubaker, M. A. (2021). Normalizing Flows : An Introduction and Review of Current Methods. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(11) :3964–3979.
- Kopylova, E., Navas-Molina, J. A., Mercier, C., Xu, Z. Z., Mahé, F., He, Y., Zhou, H.-W., Rognes, T., Caporaso, J. G., and Knight, R. (2016). Open-Source Sequence Clustering Methods Improve the State Of the Art. *mSystems*, 1(1) :e00003–15.
- Kopylova, E., Noé, L., and Touzet, H. (2012). SortMeRNA : fast and accurate filtering of ribosomal RNAs in metatranscriptomic data. *Bioinformatics*, 28(24) :3211–3217.
- Kraaijeveld, K., De Weger, L. A., Ventayol García, M., Buermans, H., Frank, J., Hiemstra, P. S., and Den Dunnen, J. T. (2015). Efficient and sensitive identification and quantification of airborne pollen using next-generation DNA sequencing. *Molecular Ecology Resources*, 15(1) :8–16.
- Krehenwinkel, Pomerantz, and Prost (2019). Genetic Biomonitoring and Biodiversity Assessment Using Portable Sequencing Technologies : Current Uses and Future Directions. *Genes*, 10(11) :858.
- Krehenwinkel, H., Fong, M., Kennedy, S., Huang, E. G., Noriyuki, S., Cayetano, L., and Gillespie, R. (2018). The effect of DNA degradation bias in passive sampling devices on metabarcoding studies of arthropod communities and their associated microbiota. *PLOS ONE*, 13(1) :e0189188.
- Krehenwinkel, H., Wolf, M., Lim, J. Y., Rominger, A. J., Simison, W. B., and Gillespie, R. G. (2017). Estimating and mitigating amplification bias in qualitative and quantitative arthropod metabarcoding. *Scientific Reports*, 7(1) :17668.
- Kuhn, E. and Lavielle, M. (2004). Coupling a stochastic approximation version of EM with an MCMC procedure. *ESAIM : Probability and Statistics*, 8 :115–131.
- Kuhn, E. and Lavielle, M. (2005). Maximum likelihood estimation in nonlinear mixed effects models. *Computational Statistics & Data Analysis*, 49(4) :1020–1038.
- Lacoursière-Roussel, A., Rosabal, M., and Bernatchez, L. (2016). Estimating fish abundance and biomass from eDNA concentrations : variability among capture methods and environmental conditions. *Molecular Ecology Resources*, 16(6) :1401–1414.
- Lalam, N. (2006). Estimation of the reaction efficiency in polymerase chain reaction. *Journal of Theoretical Biology*, 242(4) :947–953.
- Lamb, P. D., Hunter, E., Pinnegar, J. K., Creer, S., Davies, R. G., and Taylor, M. I. (2019). How quantitative is metabarcoding : A meta-analytical approach. *Molecular Ecology*, 28(2) :420–430.
- Lang, D., Tang, M., Hu, J., and Zhou, X. (2019). Genome-skimming provides accurate quantification for pollen mixtures. *Molecular Ecology Resources*, 19(6) :1433–1446.
- Laramie, M. B., Pilliod, D. S., and Goldberg, C. S. (2015). Characterizing the distribution of an endangered salmonid using environmental DNA analysis. *Biological Conservation*, 183 :29–37.
- Lawton, J. H., Bignell, D. E., Bolton, B., Bloemers, G. F., Eggleton, P., Hammond, P. M., Hodda, M., Holt, R. D., Larsen, T. B., Mawdsley, N. A., Stork, N. E., Srivastava, D. S., and Watt, A. D. (1998). Biodiversity inventories, indicator taxa and effects of habitat modification in tropical forest. *Nature*, 391(6662) :72–76.

- Leclercq, S., Rivals, E., and Jarne, P. (2010). DNA Slippage Occurs at Microsatellite Loci without Minimal Threshold Length in Humans : A Comparative Genomic Approach. *Genome Biology and Evolution*, 2(0) :325–335.
- Legay, J.-M. (1997). *L'expérience et le modèle*. Editions Quæ.
- Leinster, T. and Cobbold, C. A. (2012). Measuring diversity : the importance of species similarity. *Ecology*, 93(3) :477–489.
- Lenz, S., Hackenberg, M., and Binder, H. (2022). The JuliaConnectoR : A Functionally-Oriented Interface for Integrating *Julia* in *R*. *Journal of Statistical Software*, 101(6).
- Levi, T., Allen, J. M., Bell, D., Joyce, J., Russell, J. R., Tallmon, D. A., Vulstek, S. C., Yang, C., and Yu, D. W. (2019). Environmental DNA for the enumeration and management of Pacific salmon. *Molecular Ecology Resources*, 19(3) :597–608.
- Li, J., Lawson Handley, L. J., Harper, L. R., Brys, R., Watson, H. V., Di Muri, C., Zhang, X., and Hänfling, B. (2019). Limited dispersion and quick degradation of environmental DNA in fish ponds inferred by metabarcoding. *Environmental DNA*, 1(3) :238–250.
- Li, J., Nott, D., Fan, Y., and Sisson, S. (2017). Extending approximate Bayesian computation methods to high dimensions via a Gaussian copula model. *Computational Statistics & Data Analysis*, 106 :77–89.
- Liu, M., Burridge, C. P., Clarke, L. J., Baker, S. C., and Jordan, G. J. (2023). Does phylogeny explain bias in quantitative DNA metabarcoding? *Metabarcoding and Metagenomics*, 7 :e101266.
- Liu, N., Zou, D., Dong, D., Yang, Z., Ao, D., Liu, W., and Huang, L. (2017). Development of a multiplex loop-mediated isothermal amplification method for the simultaneous detection of *Salmonella* spp. and *Vibrio parahaemolyticus*. *Scientific Reports*, 7(1) :45601.
- Liu, W. and Saint, D. A. (2002). Validation of a quantitative method for real time PCR kinetics. *Biochemical and Biophysical Research Communications*, 294(2) :347–353.
- Livak, K. J. and Schmittgen, T. D. (2001). Analysis of relative gene expression data using real-time quantitative pcr and the 2-delta delta ct method. *Methods*, 25(4) :402–408.
- Luengo, D., Martino, L., Bugallo, M., Elvira, V., and Särkkä, S. (2020). A survey of Monte Carlo methods for parameter estimation. *EURASIP Journal on Advances in Signal Processing*, 2020(1) :25.
- Luo, M., Ji, Y., Warton, D., and Yu, D. W. (2022). Extracting abundance information from DNA-based data. *Molecular Ecology Resources*, pages 1755–0998.13703.
- Mahé, F., Rognes, T., Quince, C., De Vargas, C., and Dunthorn, M. (2015). Swarm v2 : highly-scalable and high-resolution amplicon clustering. *PeerJ*, 3 :e1420.
- Marcon, E. (2015). Mesures de la Biodiversité. Master. Kourou, France. cel-01205813.
- Marcon, E., Scotti, I., Hérault, B., Rossi, V., and Lang, G. (2014). Generalization of the Partitioning of Shannon Diversity. *PLoS ONE*, 9(3) :e90289.
- Marin, J.-M., Pudlo, P., Robert, C. P., and Ryder, R. J. (2012). Approximate Bayesian computational methods. *Statistics and Computing*, 22(6) :1167–1180.
- Marjoram, P., Molitor, J., Plagnol, V., and Tavaré, S. (2003). Markov chain Monte Carlo without likelihoods. *Proceedings of the National Academy of Sciences*, 100(26) :15324–15328.
- Marquina, D., Esparza-Salas, R., Roslin, T., and Ronquist, F. (2019). Establishing arthropod community composition using metabarcoding : Surprising inconsistencies between soil

- samples and preservative ethanol and homogenate from Malaise trap catches. *Molecular Ecology Resources*, 19(6) :1516–1530.
- Martoni, F., Piper, A. M., Rodoni, B. C., and Blacket, M. J. (2022). Disentangling bias for non-destructive insect metabarcoding. *PeerJ*, 10 :e12981.
- Maruyama, A., Nakamura, K., Yamanaka, H., Kondoh, M., and Minamoto, T. (2014). The Release Rate of Environmental DNA from Juvenile and Adult Fish. *PLoS ONE*, 9(12) :e114639.
- Matesanz, S., Pescador, D. S., Pías, B., Sánchez, A. M., Chacón-Labela, J., Illuminati, A., Cruz, M., López-Angulo, J., Marí-Mena, N., Vizcaíno, A., and Escudero, A. (2019). Estimating belowground plant abundance with DNA metabarcoding. *Molecular Ecology Resources*, 19(5) :1265–1277.
- Mathieu-Daude, F. (1996). DNA rehybridization during PCR : the 'Cot effect' and its consequences. *Nucleic Acids Research*, 24(11) :2080–2086.
- Mathon, L., Valentini, A., Guérin, P., Normandeau, E., Noel, C., Lionnet, C., Boulanger, E., Thuiller, W., Bernatchez, L., Mouillot, D., Dejean, T., and Manel, S. (2021). Benchmarking bioinformatic tools for fast and accurate eDNA metabarcoding species identification. *Molecular Ecology Resources*, 21(7) :2565–2579.
- Matsumoto, M. and Nishimura, T. (1998). Mersenne twister : a 623-dimensionally equidistributed uniform pseudo-random number generator. *ACM Transactions on Modeling and Computer Simulation*, 8(1) :3–30.
- Mayden, R. L. (1997). A hierarchy of species concepts : The denouement in the saga of the species problem. In Claridge, M. F., Dawah, H. A., and Wilson, M. R., editors, *Species : The units of diversity*, pages 381–423. Chapman & Hall.
- Mayr, E. (1942). Systematics and the origin of species from the viewpoint of a zoologist. *New York : Columbia University Press*.
- McCarthy, A., Rajabi, H., McClenaghan, B., Fahner, N. A., Porter, E., Singer, G. A. C., and Hajibabaei, M. (2022). Comparative analysis of fish environmental DNA reveals higher sensitivity achieved through targeted sequence-based metabarcoding. *Molecular Ecology Resources*, pages 1755–0998.13732.
- McLaren, M. R., Willis, A. D., and Callahan, B. J. (2019). Consistent and correctable bias in metagenomic sequencing experiments. *eLife*, 8 :e46923.
- Mehra, S. and Hu, W.-S. (2005). A kinetic model of quantitative real-time polymerase chain reaction. *Biotechnology and Bioengineering*, 91(7) :848–860.
- Metropolis, N., Rosenbluth, A. W., Rosenbluth, M. N., Teller, A. H., and Teller, E. (1953). Equation of State Calculations by Fast Computing Machines. *The Journal of Chemical Physics*, 21(6) :1087–1092.
- Milivojević, T., Rahman, S. N., Raposo, D., Siccha, M., Kucera, M., and Morard, R. (2021). High variability in SSU rDNA gene copy number among planktonic foraminifera revealed by single-cell qPCR. *ISME Communications*, 1(1) :63.
- Moinard, S., Oudet, E., Piau, D., Coissac, E., and Gonindard-Melodelima, C. (2022). The Fixed Landscape Inference MethOd (flimo) : an alternative to Approximate Bayesian Computation, faster by several orders of magnitude. *arXiv*. Publisher : arXiv Version Number : 1.
- Morley, A. A. (2014). Digital PCR : A brief history. *Biomolecular Detection and Quantification*, 1(1) :1–2.

- Mulero, S., Toulza, E., Loisier, A., Zimmerman, M., Allienne, J.-F., Foata, J., Quilichini, Y., Pointier, J.-P., Rey, O., and Boissier, J. (2021). Malacological survey in a bottle of water : A comparative study between manual sampling and environmental DNA metabarcoding approaches. *Global Ecology and Conservation*, 25 :e01428.
- Mullis, K. B. and Faloona, F. A. (1987). Specific synthesis of DNA in vitro via a polymerase-catalyzed chain reaction. In *Methods in Enzymology*, volume 155, pages 335–350. Elsevier.
- Munoz, F., Grenié, M., Denelle, P., Taudière, A., Laroche, F., Tucker, C., and Violle, C. (2018). *ecolottery* : Simulating and assessing community assembly with environmental filtering and neutral dynamics in R. *Methods in Ecology and Evolution*, 9(3) :693–703.
- Murchie, T. J., Kuch, M., Duggan, A. T., Ledger, M. L., Roche, K., Klunk, J., Karpinski, E., Hackenberger, D., Sadoway, T., MacPhee, R., Froese, D., and Poinar, H. (2021). Optimizing extraction and targeted capture of ancient environmental DNA for reconstructing past environments using the PalaeoChip Arctic-1.0 bait-set. *Quaternary Research*, 99 :305–328.
- Mächler, E., Walser, J., and Altermatt, F. (2021). Decision-making and best practices for taxonomy-free environmental DNA metabarcoding in biomonitoring using Hill numbers. *Molecular Ecology*, 30(13) :3326–3339.
- Nelder, J. A. and Mead, R. (1965). A Simplex Method for Function Minimization. *The Computer Journal*, 7(4) :308–313.
- Newbold, T., Hudson, L. N., Arnell, A. P., Contu, S., De Palma, A., Ferrier, S., Hill, S. L. L., Hoskins, A. J., Lysenko, I., Phillips, H. R. P., Burton, V. J., Chng, C. W. T., Emerson, S., Gao, D., Pask-Hale, G., Hutton, J., Jung, M., Sanchez-Ortiz, K., Simmons, B. I., Whitmee, S., Zhang, H., Scharlemann, J. P. W., and Purvis, A. (2016). Has land use pushed terrestrial biodiversity beyond the planetary boundary? A global assessment. *Science*, 353(6296) :288–291.
- Newton, C. R. and Graham, A. (2000). *PCR. The Introduction to biotechniques series*. BIOS Scientific Publ, Oxford, 2. ed., repr edition.
- Nichols, R. V., Vollmers, C., Newsom, L. A., Wang, Y., Heintzman, P. D., Leighton, M., Green, R. E., and Shapiro, B. (2018). Minimizing polymerase biases in metabarcoding. *Molecular Ecology Resources*, 18(5) :927–939.
- Nicholson, G., Smith, A. V., Jonsson, F., Gustafsson, O., Stefansson, K., and Donnelly, P. (2002). Assessing population differentiation and isolation from single-nucleotide polymorphism data. *Journal of the Royal Statistical Society : Series B (Statistical Methodology)*, 64(4) :695–715.
- Nocedal, J. and Wright, S. J. (2006). *Numerical Optimization*. Springer Series in Operations Research and Financial Engineering. Springer New York.
- Nordstrom, B., Mitchell, N., Byrne, M., and Jarman, S. (2022). A review of applications of environmental DNA for reptile conservation and management. *Ecology and Evolution*, 12(6).
- Notomi, T. (2000). Loop-mediated isothermal amplification of DNA. *Nucleic Acids Research*, 28(12) :63e–63.
- Nunes, M. A. and Balding, D. J. (2010). On Optimal Selection of Summary Statistics for Approximate Bayesian Computation. *Statistical Applications in Genetics and Molecular Biology*, 9(1).
- Ocal, K., Grima, R., and Sanguinetti, G. (2020). Parameter estimation for biochemical reaction networks using Wasserstein distances. *Journal of Physics A : Mathematical and Theo-*

- retical*, 53(3) :034002.
- Ogram, A., Sayler, G. S., and Barkay, T. (1987). The extraction and purification of microbial DNA from sediments. *Journal of Microbiological Methods*, 7(2-3) :57–66.
- Olesen, S. W., Duvallat, C., and Alm, E. J. (2017). dbOTU3 : A new implementation of distribution-based OTU calling. *PLOS ONE*, 12(5) :e0176335.
- Papamakarios, G., Nalisnick, E., Rezende, D. J., Mohamed, S., and Lakshminarayanan, B. (2019). Normalizing Flows for Probabilistic Modeling and Inference. *arXiv*. Publisher : arXiv Version Number : 2.
- Parada, A. E., Needham, D. M., and Fuhrman, J. A. (2016). Every base matters : assessing small subunit rRNA primers for marine microbiomes with mock communities, time series and global field samples : Primers for marine microbiome studies. *Environmental Microbiology*, 18(5) :1403–1414.
- Paris, C., Servin, B., and Boitard, S. (2019). Inference of Selection from Genetic Time Series Using Various Parametric Approximations to the Wright-Fisher Model. *G3 Genes/Genomes/Genetics*, 9(12) :4073–4086.
- Pauvert, C., Buée, M., Laval, V., Edel-Hermann, V., Fauchery, L., Gautier, A., Lesur, I., Vallance, J., and Vacher, C. (2019). Bioinformatics matters : The accuracy of plant and soil fungal community data is highly dependent on the metabarcoding pipeline. *Fungal Ecology*, 41 :23–33.
- Pawlowski, J., Apothéoz-Perret-Gentil, L., and Altermatt, F. (2020). Environmental DNA : What’s behind the term ? Clarifying the terminology and recommendations for its future use in biomonitoring. *Molecular Ecology*, 29(22) :4258–4264.
- Pawluczyk, M., Weiss, J., Links, M. G., Egaña Aranguren, M., Wilkinson, M. D., and Egea-Cortines, M. (2015). Quantitative evaluation of bias in PCR amplification and next-generation sequencing derived from metabarcoding samples. *Analytical and Bioanalytical Chemistry*, 407(7) :1841–1848.
- Peccoud, J. and Jacob, C. (1996). Theoretical uncertainty of measurements using quantitative polymerase chain reaction. *Biophysical Journal*, 71(1) :101–108.
- Peirson, S. N. (2003). Experimental validation of novel and conventional approaches to quantitative real-time PCR data analysis. *Nucleic Acids Research*, 31(14) :73e–73.
- Peng, X. and Dorman, K. S. (2021). AmpliCI : a high-resolution model-based approach for denoising Illumina amplicon data. *Bioinformatics*, 36(21) :5151–5158.
- Pereira, T. J., De Santiago, A., Schuelke, T., Hardy, S. M., and Bik, H. M. (2020). The impact of intragenomic rRNA variation on metabarcoding-derived diversity estimates : A case study from marine nematodes. *Environmental DNA*, 2(4) :519–534.
- Pervez, M. T., Hasnain, M. J. U., Abbas, S. H., Moustafa, M. F., Aslam, N., and Shah, S. S. M. (2022). A Comprehensive Review of Performance of Next-Generation Sequencing Platforms. *BioMed Research International*, 2022 :1–12.
- Pfaffl, M. W. (2001). A new mathematical model for relative quantification in real-time RT-PCR. *Nucleic Acids Research*, 29(9) :45e–45.
- Piau, D. (2004). Immortal branching Markov processes : Averaging properties and PCR applications. *The Annals of Probability*, 32(1A) :337–364.
- Piau, D. (2005). Confidence intervals for nonhomogeneous branching processes and polymerase chain reactions. *The Annals of Probability*, 33(2) :674–702.

- Picard, M., Pochon, X., Atalah, J., Pearman, J. K., Rees, A., Howarth, J. D., Moy, C. M., Vandergoes, M. J., Hawes, I., Khan, S., and Wood, S. A. (2022). Using metabarcoding and droplet digital PCR to investigate drivers of historical shifts in cyanobacteria from six contrasting lakes. *Scientific Reports*, 12(1) :12810.
- Piepenburg, O., Williams, C. H., Stemple, D. L., and Armes, N. A. (2006). DNA Detection Using Recombination Proteins. *PLoS Biology*, 4(7) :e204.
- Pierella Karlusich, J. J., Pelletier, E., Zinger, L., Lombard, F., Zingone, A., Colin, S., Gasol, J. M., Dorrell, R. G., Henry, N., Scalco, E., Acinas, S. G., Wincker, P., Vargas, C., and Bowler, C. (2022). A robust approach to estimate relative phytoplankton cell abundances from metagenomes. *Molecular Ecology Resources*, pages 1755–0998.13592.
- Pilliod, D. S., Goldberg, C. S., Arkle, R. S., and Waits, L. P. (2014). Factors influencing detection of eDNA from a stream-dwelling amphibian. *Molecular Ecology Resources*, 14(1) :109–116.
- Piñol, J., Mir, G., Gomez-Polo, P., and Agustí, N. (2015). Universal and blocking primer mismatches limit the use of high-throughput DNA sequencing for the quantitative metabarcoding of arthropods. *Molecular Ecology Resources*, 15(4) :819–830.
- Piñol, J., Senar, M. A., and Symondson, W. O. C. (2019). The choice of universal primers and the characteristics of the species mixture determine when DNA metabarcoding can be quantitative. *Molecular Ecology*, 28(2) :407–419.
- Pompanon, F., Deagle, B. E., Symondson, W. O. C., Brown, D. S., Jarman, S. N., and Taberlet, P. (2012). Who is eating what : diet assessment using next generation sequencing. *Molecular Ecology*, 21(8) :1931–1950.
- Pont, D., Meulenbroek, P., Bammer, V., Dejean, T., Erős, T., Jean, P., Lenhardt, M., Nagel, C., Pekarik, L., Schabuss, M., Stoeckle, B. C., Stoica, E., Zornig, H., Weigand, A., and Valentini, A. (2023). Quantitative monitoring of diverse fish communities on a large scale combining eDNA metabarcoding and qPCR. *Molecular Ecology Resources*, 23(2) :396–409.
- Popov, V., Ellis-Robinson, A., and Humphris, G. (2019). Modelling reassurances of clinicians with hidden Markov models. *BMC Medical Research Methodology*, 19(1) :11.
- Porazinska, D. L., Sung, W., Giblin-Davis, R. M., and Thomas, W. K. (2010). Reproducibility of read numbers in high-throughput sequencing analysis of nematode community composition and structure. *Molecular Ecology Resources*, 10(4) :666–676.
- Pornon, A., Escaravage, N., Burrus, M., Holota, H., Khimoun, A., Mariette, J., Pellizzari, C., Iribar, A., Etienne, R., Taberlet, P., Vidal, M., Winterton, P., Zinger, L., and Andalo, C. (2016). Using metabarcoding to reveal and quantify plant-pollinator interactions. *Scientific Reports*, 6(1) :27282.
- Porter, T. M. and Hajibabaei, M. (2020). Putting COI Metabarcoding in Context : The Utility of Exact Sequence Variants (ESVs) in Biodiversity Analysis. *Frontiers in Ecology and Evolution*, 8 :248.
- Potapov, V. and Ong, J. L. (2017). Examining Sources of Error in PCR by Single-Molecule Sequencing. *PLOS ONE*, 12(1) :e0169774.
- Prangle, D. (2017). `gk` : An R Package for the g-and-k and generalised g-and-h Distributions. *arXiv :1706.06889 [stat]*. arXiv : 1706.06889.
- Price, L. F., Drovandi, C. C., Lee, A., and Nott, D. J. (2018). Bayesian Synthetic Likelihood. *Journal of Computational and Graphical Statistics*, 27(1) :1–11.
- Pritchard, J. K., Seielstad, M. T., Perez-Lezaun, A., and Feldman, M. W. (1999). Population

- growth of human Y chromosomes : a study of Y chromosome microsatellites. *Molecular Biology and Evolution*, 16(12) :1791–1798.
- Pritchard, L., Corne, D., Kell, D., Rowland, J., and Winson, M. (2005). A general model of error-prone PCR. *Journal of Theoretical Biology*, 234(4) :497–509.
- Prodan, A., Tremaroli, V., Brolin, H., Zwinderman, A. H., Nieuwdorp, M., and Levin, E. (2020). Comparing bioinformatic pipelines for microbial 16S rRNA amplicon sequencing. *PLOS ONE*, 15(1) :e0227434.
- Props, R., Monsieurs, P., Mysara, M., Clement, L., and Boon, N. (2016). Measuring the biodiversity of microbial communities by flow cytometry. *Methods in Ecology and Evolution*, 7(11) :1376–1385.
- Purvis, A., Newbold, T., De Palma, A., Contu, S., Hill, S. L., Sanchez-Ortiz, K., Phillips, H. R., Hudson, L. N., Lysenko, I., Börger, L., and Scharlemann, J. P. (2018). Modelling and Projecting the Response of Local Terrestrial Biodiversity Worldwide to Land Use and Related Pressures : The PREDICTS Project. In *Advances in Ecological Research*, volume 58, pages 201–241. Elsevier.
- R Core Team (2021). R : A language and environment for statistical computing. Vienna, Austria.
- Ramakers, C., Ruijter, J. M., Deprez, R. H., and Moorman, A. F. (2003). Assumption-free analysis of quantitative real-time polymerase chain reaction (PCR) data. *Neuroscience Letters*, 339(1) :62–66.
- Ramsay, J. O., Hooker, G., Campbell, D., and Cao, J. (2007). Parameter Estimation for Differential Equations : a Generalized Smoothing Approach. *Journal of the Royal Statistical Society Series B : Statistical Methodology*, 69(5) :741–796.
- Rao, X., Lai, D., and Huang, X. (2013). A New Method for Quantitative Real-Time Polymerase Chain Reaction Data Analysis. *Journal of Computational Biology*, 20(9) :703–711.
- Revels, J., Lubin, M., and Papamarkou, T. (2016). Forward-Mode Automatic Differentiation in Julia. *arXiv :1607.07892 [cs]*. arXiv : 1607.07892.
- Rezende, D. J. and Mohamed, S. (2015). Variational Inference with Normalizing Flows. *arXiv*. Publisher : arXiv Version Number : 6.
- Riaz, T., Shehzad, W., Viari, A., Pompanon, F., Taberlet, P., and Coissac, E. (2011). eco-Primers : inference of new DNA barcode markers from whole genome sequence analysis. *Nucleic Acids Research*, 39(21) :e145–e145.
- Rideout, J. R., He, Y., Navas-Molina, J. A., Walters, W. A., Ursell, L. K., Gibbons, S. M., Chase, J., McDonald, D., Gonzalez, A., Robbins-Pianka, A., Clemente, J. C., Gilbert, J. A., Huse, S. M., Zhou, H.-W., Knight, R., and Caporaso, J. G. (2014). Subsampled open-reference clustering creates consistent, comprehensive OTU definitions and scales to billions of sequences. *PeerJ*, 2 :e545.
- Ripley, B. D. (1987). Stochastic Models. In *Wiley Series in Probability and Statistics*, pages 96–117. John Wiley & Sons, Inc., Hoboken, NJ, USA.
- Ritz, C. and Spiess, A.-N. (2008). *qpcR* : an R package for sigmoidal model selection in quantitative real-time polymerase chain reaction analysis. *Bioinformatics*, 24(13) :1549–1551.
- Rodrigues, L., Ortega, I., Vieira, R., Carrasco, D., and Proietti, M. (2019). Crane flies (Diptera, Tipuloidea) from southern Neotropical salt marshes : survey with DNA barcoding. *Iheringia. Série Zoologia*, 109 :e2019013.

- Rognes, T., Flouri, T., Nichols, B., Quince, C., and Mahé, F. (2016). VSEARCH : a versatile open source tool for metagenomics. *PeerJ*, 4 :e2584.
- Rourke, M. L., Fowler, A. M., Hughes, J. M., Broadhurst, M. K., DiBattista, J. D., Fielder, S., Wilkes Walburn, J., and Furlan, E. M. (2022). Environmental DNA (eDNA) as a tool for assessing fish biomass : A review of approaches and future considerations for resource surveys. *Environmental DNA*, 4(1) :9–33.
- Ruppert, K. M., Kline, R. J., and Rahman, M. S. (2019). Past, present, and future perspectives of environmental DNA (eDNA) metabarcoding : A systematic review in methods, monitoring, and applications of global eDNA. *Global Ecology and Conservation*, 17 :e00547.
- Rutledge, R. G. (2004). Sigmoidal curve-fitting redefines quantitative real-time PCR with the prospective of developing automated high-throughput applications. *Nucleic Acids Research*, 32(22) :e178–e178.
- Rutledge, R. G. and Stewart, D. (2008). A kinetic-based sigmoidal model for the polymerase chain reaction and its application to high-capacity absolute quantitative real-time PCR. *BMC Biotechnology*, 8(1) :47.
- Sakamoto, W. and Takami, T. (2018). Chloroplast DNA Dynamics : Copy Number, Quality Control and Degradation. *Plant and Cell Physiology*, 59(6) :1120–1127.
- Sassoubre, L. M., Yamahara, K. M., Gardner, L. D., Block, B. A., and Boehm, A. B. (2016). Quantification of Environmental DNA (eDNA) Shedding and Decay Rates for Three Marine Fish. *Environmental Science & Technology*, 50(19) :10456–10464.
- Schiebelhut, L. M., Abboud, S. S., Gómez Daglio, L. E., Swift, H. F., and Dawson, M. N. (2017). A comparison of DNA extraction methods for high-throughput DNA analyses. *Molecular Ecology Resources*, 17(4) :721–729.
- Schirmer, M., Ijaz, U. Z., D’Amore, R., Hall, N., Sloan, W. T., and Quince, C. (2015). Insight into biases and sequencing errors for amplicon sequencing with the Illumina MiSeq platform. *Nucleic Acids Research*, 43(6) :e37–e37.
- Schloss, P. D., Westcott, S. L., Ryabin, T., Hall, J. R., Hartmann, M., Hollister, E. B., Lesniewski, R. A., Oakley, B. B., Parks, D. H., Robinson, C. J., Sahl, J. W., Stres, B., Thallinger, G. G., Van Horn, D. J., and Weber, C. F. (2009). Introducing mothur : Open-Source, Platform-Independent, Community-Supported Software for Describing and Comparing Microbial Communities. *Applied and Environmental Microbiology*, 75(23) :7537–7541.
- Schneider, J., Valentini, A., Dejean, T., Montarsi, F., Taberlet, P., Glaizot, O., and Fumagalli, L. (2016). Detection of Invasive Mosquito Vectors Using Environmental DNA (eDNA) from Water Samples. *PLOS ONE*, 11(9) :e0162493.
- Schnell, I. B., Bohmann, K., and Gilbert, M. T. P. (2015). Tag jumps illuminated - reducing sequence-to-sample misidentifications in metabarcoding studies. *Molecular Ecology Resources*, 15(6) :1289–1303.
- Schnell, S. and Mendoza, C. (1997). Enzymological Considerations for the Theoretical Description of the Quantitative Competitive Polymerase Chain Reaction (QC-PCR). *Journal of Theoretical Biology*, 184(4) :433–440.
- Schrader, C., Schielke, A., Ellerbroek, L., and Johne, R. (2012). PCR inhibitors - occurrence, properties and removal. *Journal of Applied Microbiology*, 113(5) :1014–1026.
- Seymour, M. (2019). Rapid progression and future of environmental DNA research. *Communications Biology*, 2(1) :80.
- Seymour, M., Durance, I., Cosby, B. J., Ransom-Jones, E., Deiner, K., Ormerod, S. J., Col-

- bourne, J. K., Wilgar, G., Carvalho, G. R., De Bruyn, M., Edwards, F., Emmett, B. A., Bik, H. M., and Creer, S. (2018). Acidity promotes degradation of multi-species environmental DNA in lotic mesocosms. *Communications Biology*, 1(1) :4.
- Shahriari, B., Swersky, K., Wang, Z., Adams, R. P., and De Freitas, N. (2016). Taking the Human Out of the Loop : A Review of Bayesian Optimization. *Proceedings of the IEEE*, 104(1) :148–175.
- Shelton, A. O., Gold, Z. J., Jensen, A. J., D’Agnese, E., Andruszkiewicz Allan, E., Van Cise, A., Gallego, R., Ramón-Laca, A., Garber-Yonts, M., Parsons, K., and Kelly, R. P. (2022). Toward quantitative metabarcoding. *Ecology*.
- Shirazi, S., Meyer, R. S., and Shapiro, B. (2021). Revisiting the effect of pcr replication and sequencing depth on biodiversity metrics in environmental dna metabarcoding. *Ecology and Evolution*, 11(22) :15766–15779.
- Shoemaker, J. S., Painter, I. S., and Weir, B. S. (1999). Bayesian statistics in genetics : a guide for the uninitiated. *Trends in Genetics*, 15(9) :354–358.
- Sidstedt, M., Rådström, P., and Hedman, J. (2020). PCR inhibition in qPCR, dPCR and MPS—mechanisms and solutions. *Analytical and Bioanalytical Chemistry*, 412(9) :2009–2023.
- Siljak-Yakovlev, S., Pustahija, F., Oli, E. M., Boguni, F., Muratovi, E., Bai, N., Catrice, O., and Brown, S. C. (2010). Towards a Genome Size and Chromosome Number Database of Balkan Flora : C-Values in 343 Taxa with Novel Values for 242. *Advanced science letters*, 3(2) :190–213.
- Silverman, J. D., Bloom, R. J., Jiang, S., Durand, H. K., Dallow, E., Mukherjee, S., and David, L. A. (2021). Measuring and mitigating PCR bias in microbiota datasets. *PLOS Computational Biology*, 17(7) :e1009113.
- Singh, R. R. (2022). Target Enrichment Approaches for Next-Generation Sequencing Applications in Oncology. *Diagnostics*, 12(7) :1539.
- Sipos, R., Székely, A. J., Palatinszky, M., Révész, S., Márialigeti, K., and Nikolausz, M. (2007). Effect of primer mismatch, annealing temperature and PCR cycle number on 16S rRNA gene-targetting bacterial community analysis : PCR parameters influencing quantitative bias. *FEMS Microbiology Ecology*, 60(2) :341–350.
- Sisson, S. A., Fan, Y., and Beaumont, M. A., editors (2018). *Handbook of Approximate Bayesian Computation*. Chapman and Hall/CRC, Boca Raton, Florida : CRC Press, [2019], 1 edition.
- Sisson, S. A., Fan, Y., and Tanaka, M. M. (2007). Sequential Monte Carlo without likelihoods. *Proceedings of the National Academy of Sciences*, 104(6) :1760–1765.
- Skern, R., Frost, P., and Nilsen, F. (2005). Relative transcript quantification by quantitative pcr : Roughly right or precisely wrong? *BMC Molecular Biology*, 6(1) :10.
- Smets, W., Leff, J. W., Bradford, M. A., McCulley, R. L., Lebeer, S., and Fierer, N. (2016). A method for simultaneous measurement of soil bacterial abundances and community composition via 16S rRNA gene sequencing. *Soil Biology and Biochemistry*, 96 :145–151.
- Spieß, A.-N., Feig, C., and Ritz, C. (2008). Highly accurate sigmoidal fitting of real-time PCR data by introducing a parameter for asymmetry. *BMC Bioinformatics*, 9(1) :221.
- Steinke, D., Braukmann, T. W., Manerus, L., Woodhouse, A., and Elbrecht, V. (2021). Effects of Malaise trap spacing on species richness and composition of terrestrial arthropod bulk samples. *Metabarcoding and Metagenomics*, 5 :e59201.

- Stephens, M. (2004). Inference Under the Coalescent. In Balding, D. J., Bishop, M., and Cannings, C., editors, *Handbook of Statistical Genetics*, page bbc22. John Wiley & Sons, Ltd, Chichester.
- Stewart, A., Rioux, D., Boyer, F., Gielly, L., Pompanon, F., Saillard, A., Thuiller, W., Valay, J.-G., Maréchal, E., and Coissac, E. (2021). Altitudinal Zonation of Green Algae Biodiversity in the French Alps. *Frontiers in Plant Science*, 12 :679428.
- Stirling, A. (2007). A general framework for analysing diversity in science, technology and society. *Journal of The Royal Society Interface*, 4(15) :707–719.
- Stolovitzky, G. and Cecchi, G. (1996). Efficiency of DNA replication in the polymerase chain reaction. *Proceedings of the National Academy of Sciences*, 93(23) :12947–12952.
- Strickler, K. M., Fremier, A. K., and Goldberg, C. S. (2015). Quantifying effects of UV-B, temperature, and pH on eDNA degradation in aquatic microcosms. *Biological Conservation*, 183 :85–92.
- Sun, F. (1995). The Polymerase Chain Reaction and Branching Processes. *Journal of Computational Biology*, 2(1) :63–86.
- Suzuki, M. T. and Giovannoni, S. J. (1996). Bias caused by template annealing in the amplification of mixtures of 16S rRNA genes by PCR. *Applied and Environmental Microbiology*, 62(2) :625–630.
- Svec, D., Tichopad, A., Novosadova, V., Pfaffl, M. W., and Kubista, M. (2015). How good is a PCR efficiency estimate : Recommendations for precise and robust qPCR efficiency assessments. *Biomolecular Detection and Quantification*, 3 :9–16.
- Taberlet, P., Bonin, A., Zinger, L., and Coissac, E. (2018). *Environmental DNA*, volume 1. Oxford University Press.
- Taberlet, P., Coissac, E., Pompanon, F., Brochmann, C., and Willerslev, E. (2012). Towards next-generation biodiversity assessment using DNA metabarcoding. *Molecular Ecology*, 21(8) :2045–2050.
- Taberlet, P., Coissac, E., Pompanon, F., Gielly, L., Miquel, C., Valentini, A., Vermet, T., Cortier, G., Brochmann, C., and Willerslev, E. (2007). Power and limitations of the chloroplast trnL (UAA) intron for plant DNA barcoding. *Nucleic Acids Research*, 35(3) :e14–e14.
- Tahamtan, A. and Ardebili, A. (2020). Real-time RT-PCR in COVID-19 detection : issues affecting the results. *Expert Review of Molecular Diagnostics*, 20(5) :453–454.
- Takahara, T., Minamoto, T., Yamanaka, H., Doi, H., and Kawabata, Z. (2012). Estimation of Fish Biomass Using Environmental DNA. *PLoS ONE*, 7(4) :e35868.
- Takasaki, K., Aihara, H., Imanaka, T., Matsudaira, T., Tsukahara, K., Usui, A., Osaki, S., and Doi, H. (2021). Water pre-filtration methods to improve environmental DNA detection by real-time PCR and metabarcoding. *PLOS ONE*, 16(5) :e0250162.
- Tataru, P., Bataillon, T., and Hobolth, A. (2015). Inference Under a Wright-Fisher Model Using an Accurate Beta Approximation. *Genetics*, 201(3) :1133–1141.
- Tavaré, S. (1986). Some probabilistic and statistical problems in the analysis of dna sequences. In *Lectures on Mathematics in the Life Sciences*.
- Tavaré, S., Balding, D. J., Griffiths, R. C., and Donnelly, P. (1997). Inferring Coalescence Times From DNA Sequence Data. *Genetics*, 145(2) :505–518.
- Tellinghuisen, J. and Spiess, A.-N. (2014). Statistical uncertainty and its propagation in the analysis of quantitative polymerase chain reaction data : Comparison of methods. *Analytical*

- Biochemistry*, 464 :94–102.
- Tellinghuisen, J. and Spiess, A.-N. (2015). Bias and Imprecision in Analysis of Real-Time Quantitative Polymerase Chain Reaction Data. *Analytical Chemistry*, 87(17) :8925–8931.
- Thomas, A. C., Deagle, B. E., Eveson, J. P., Harsch, C. H., and Trites, A. W. (2016). Quantitative DNA metabarcoding : improved estimates of species proportional biomass using correction factors derived from control material. *Molecular Ecology Resources*, 16(3) :714–726.
- Thomas, A. C., Jarman, S. N., Haman, K. H., Trites, A. W., and Deagle, B. E. (2014). Improving accuracy of DNA diet estimates using food tissue control materials and an evaluation of proxies for digestion bias. *Molecular Ecology*, 23(15) :3706–3718.
- Tichopad, A. (2003). Standardized determination of real-time PCR efficiency from a single reaction set-up. *Nucleic Acids Research*, 31(20) :122e–122.
- Tilman, D. (1997). The Influence of Functional Diversity and Composition on Ecosystem Processes. *Science*, 277(5330) :1300–1302.
- Tsallis, C. (1988). Possible generalization of Boltzmann-Gibbs statistics. *Journal of Statistical Physics*, 52(1-2) :479–487.
- Tsallis, C. (1994). What are the numbers that experiments provide. *Química Nova*, 17 :468–471.
- Tsuji, S., Takahara, T., Doi, H., Shibata, N., and Yamanaka, H. (2019). The detection of aquatic macroorganisms using environmental DNA analysis—A review of methods for collection, extraction, and detection. *Environmental DNA*, 1(2) :99–108.
- Turner, C. R., Barnes, M. A., Xu, C. C. Y., Jones, S. E., Jerde, C. L., and Lodge, D. M. (2014). Particle size distribution and optimal capture of aqueous microbial eDNA. *Methods in Ecology and Evolution*, 5(7) :676–684.
- Uchii, K., Doi, H., Okahashi, T., Katano, I., Yamanaka, H., Sakata, M. K., and Minamoto, T. (2019). Comparison of inhibition resistance among PCR reagents for detection and quantification of environmental DNA. *Environmental DNA*, 1(4) :359–367.
- Ushio, M., Murakami, H., Masuda, R., Sado, T., Miya, M., Sakurai, S., Yamanaka, H., Minamoto, T., and Kondoh, M. (2018). Quantitative monitoring of multispecies fish environmental dna using high-throughput sequencing. *Metabarcoding and Metagenomics*, 2 :e23297.
- Valentini, A., Miquel, C., Nawaz, M. A., Bellemain, E., Coissac, E., Pompanon, F., Gielly, L., Cruaud, C., Nascetti, G., Wincker, P., Swenson, J. E., and Taberlet, P. (2009). New perspectives in diet analysis based on DNA barcoding and parallel pyrosequencing : the *trnL* approach. *Molecular Ecology Resources*, 9 :51–60.
- Valentini, A., Taberlet, P., Miaud, C., Civade, R., Herder, J., Thomsen, P. F., Bellemain, E., Besnard, A., Coissac, E., Boyer, F., Gaboriaud, C., Jean, P., Poulet, N., Roset, N., Copp, G. H., Geniez, P., Pont, D., Argillier, C., Baudoin, J.-M., Peroux, T., Crivelli, A. J., Olivier, A., Acqueberge, M., Le Brun, M., Møller, P. R., Willerslev, E., and Dejean, T. (2016). Next-generation monitoring of aquatic biodiversity using environmental DNA metabarcoding. *Molecular Ecology*, 25(4) :929–942.
- Van Der Graaf, P. and Schoemaker, R. (1999). Analysis of asymmetry of agonist concentration–effect curves. *Journal of Pharmacological and Toxicological Methods*, 41(2-3) :107–115.
- van der Loos, L. M. and Nijland, R. (2021). Biases in bulk : DNA metabarcoding of marine communities and the methodology involved. *Molecular Ecology*, 30(13) :3270–3288.
- Velikanov, M. V. and Kapral, R. (1999). Polymerase Chain Reaction : A Markov Process

- Approach. *Journal of Theoretical Biology*, 201(4) :239–249.
- Walker, B. H. (1992). Biodiversity and Ecological Redundancy. *Conservation Biology*, 6(1) :18–23.
- Wang, C. and Yang, C. J. (2013). Application of Molecular Beacons in Real-Time PCR. In Yang, C. J. and Tan, W., editors, *Molecular Beacons*, pages 45–59. Springer Berlin Heidelberg, Berlin, Heidelberg.
- Wang, D., Zhao, C., Cheng, R., and Sun, F. (2000). Estimation of the Mutation Rate During Error-prone Polymerase Chain Reaction. *Journal of Computational Biology*, 7(1-2) :143–158.
- Webb, C. O., Losos, J. B., and Agrawal, A. A. (2006). Integrating phylogenies into community ecology. *Ecology*, 87(sp7) :S1–S2.
- Wegmann, D., Leuenberger, C., and Excoffier, L. (2009). Efficient Approximate Bayesian Computation Coupled With Markov Chain Monte Carlo Without Likelihood. *Genetics*, 182(4) :1207–1218.
- Weiss, G. (1997). A coalescent approach to the polymerase chain reaction. *Nucleic Acids Research*, 25(15) :3082–3087.
- Westcott, S. L. and Schloss, P. D. (2015). De novo clustering methods outperform reference-based methods for assigning 16S rRNA gene sequences to operational taxonomic units. *PeerJ*, 3 :e1487.
- Westcott, S. L. and Schloss, P. D. (2017). OptiClust, an Improved Method for Assigning Amplicon-Based Sequence Data to Operational Taxonomic Units. *mSphere*, 2(2) :e00073–17.
- Wheeler, Q. D. (1999). Why the phylogenetic species concept?—Elementary. *Journal of Nematology*, 31(2) :134–141.
- Whittaker, R. H. (1960). Vegetation of the Siskiyou Mountains, Oregon and California. *Ecological Monographs*, 30(3) :279–338.
- Wiemers, M. and Fiedler, K. (2007). Does the DNA barcoding gap exist? – a case study in blue butterflies (Lepidoptera : Lycaenidae). *Frontiers in Zoology*, 4(1) :8.
- Wiesner, R. J., Rüegg, J., and Morano, I. (1992). Counting target molecules by exponential polymerase chain reaction : Copy number of mitochondrial DNA in rat tissues. *Biochemical and Biophysical Research Communications*, 183(2) :553–559.
- Wilder, M. L., Farrell, J. M., and Green, H. C. (2023). Estimating edna shedding and decay rates for muskellunge in early stages of development. *Environmental DNA*, 5(2) :251–263.
- Willerslev, E., Davison, J., Moora, M., Zobel, M., Coissac, E., Edwards, M. E., Lorenzen, E. D., Vestergård, M., Gussarova, G., Haile, J., Craine, J., Gielly, L., Boessenkool, S., Epp, L. S., Pearman, P. B., Cheddadi, R., Murray, D., Bräthen, K. A., Yoccoz, N., Binney, H., Cruaud, C., Wincker, P., Goslar, T., Alsos, I. G., Bellemain, E., Brysting, A. K., Elven, R., Sønstebo, J. H., Murton, J., Sher, A., Rasmussen, M., Rønn, R., Mourier, T., Cooper, A., Austin, J., Möller, P., Froese, D., Zazula, G., Pompanon, F., Rioux, D., Niderkorn, V., Tikhonov, A., Savvinov, G., Roberts, R. G., MacPhee, R. D. E., Gilbert, M. T. P., Kjær, K. H., Orlando, L., Brochmann, C., and Taberlet, P. (2014). Fifty thousand years of Arctic vegetation and megafaunal diet. *Nature*, 506(7486) :47–51.
- Williams, P. H. and Gaston, K. J. (1994). Measuring more of biodiversity : Can higher-taxon richness predict wholesale species richness? *Biological Conservation*, 67(3) :211–217.
- Wilson, B. D., Eisenstein, M., and Soh, H. T. (2019). High-Fidelity Nanopore Sequencing of

- Ultra-Short DNA Targets. *Analytical Chemistry*, 91(10) :6783–6789.
- Wilson, E. O., editor (1999). *Biodiversity : papers from the National Forum on Biodiversity held September 21 - 25, 1986, in Washington, D.C.* Number 1 in Biodiversity / E. O. Wilson, ed. National Academy Press, Washington, D.C, 14. print edition.
- Wood, S. A., Pochon, X., Laroche, O., Ammon, U., Adamson, J., and Zaiko, A. (2019). A comparison of droplet digital polymerase chain reaction (PCR), quantitative PCR and metabarcoding for species-specific detection in environmental DNA. *Molecular Ecology Resources*, 19(6) :1407–1419.
- Wood, S. N. (2010). Statistical inference for noisy nonlinear ecological dynamic systems. *Nature*, 466(7310) :1102–1104.
- Wu, C. F. J. (1983). On the Convergence Properties of the EM Algorithm. *The Annals of Statistics*, 11(1).
- Wächter, A. and Biegler, L. T. (2006). On the implementation of an interior-point filter line-search algorithm for large-scale nonlinear programming. *Mathematical Programming*, 106(1) :25–57.
- Yang, C., Bohmann, K., Wang, X., Cai, W., Wales, N., Ding, Z., Gopalakrishnan, S., and Yu, D. W. (2021). Biodiversity Soup II : A bulk-sample metabarcoding pipeline emphasizing error reduction. *Methods in Ecology and Evolution*, 12(7) :1252–1264.
- Yates, M. C., Glaser, D. M., Post, J. R., Cristescu, M. E., Fraser, D. J., and Derry, A. M. (2021). The relationship between eDNA particle concentration and organism abundance in nature is strengthened by allometric scaling. *Molecular Ecology*, 30(13) :3068–3082.
- Zhang, H., Li, H., Zhu, H., Pekárek, J., Podešva, P., Chang, H., and Neužil, P. (2019). Revealing the secrets of PCR. *Sensors and Actuators B : Chemical*, 298 :126924.
- Zhao, S. and Fernald, R. D. (2005). Comprehensive Algorithm for Quantitative Real-Time Polymerase Chain Reaction. *Journal of Computational Biology*, 12(8) :1047–1064.
- Zinger, L., Bonin, A., Alsos, I. G., Bálint, M., Bik, H., Boyer, F., Chariton, A. A., Creer, S., Coissac, E., Deagle, B. E., De Barba, M., Dickie, I. A., Dumbrell, A. J., Ficitola, G. F., Fierer, N., Fumagalli, L., Gilbert, M. T. P., Jarman, S., Jumpponen, A., Kausrud, H., Orlando, L., Pansu, J., Pawlowski, J., Tedersoo, L., Thomsen, P. F., Willerslev, E., and Taberlet, P. (2019). DNA metabarcoding—Need for robust experimental designs to draw sound ecological conclusions. *Molecular Ecology*, 28(8) :1857–1862.
- Zoschke, R., Liere, K., and Börner, T. (2007). From seedling to mature plant : Arabidopsis plastidial genome copy number, RNA accumulation and transcription are differentially regulated during leaf development : Plastome copy number in Arabidopsis leaf development. *The Plant Journal*, 50(4) :710–722.

Annexe A

Annexes au manuscrit *The Fixed Landscape Inference MethOd (flimo)*

Je reproduis ici les fichiers de la section *Supplementary Information* tels qu'ils sont inclus dans le manuscrit *The Fixed Landscape Inference MethOd (flimo)* du Chapitre 1.

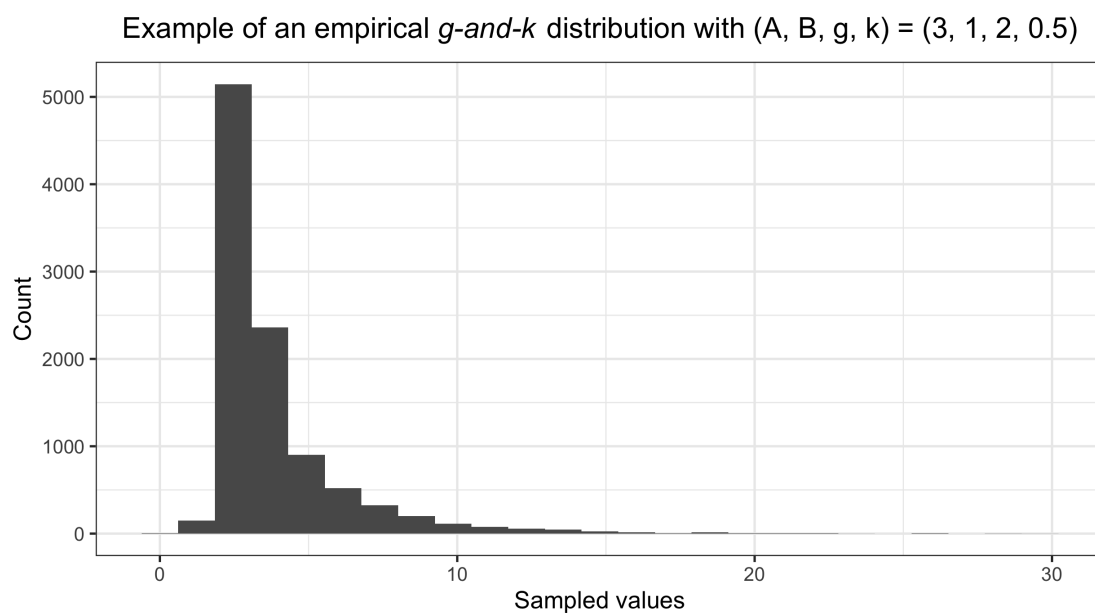


FIGURE S1 – Empirical histogram of a *g-and-k* distribution for parameters $(A, B, g, k) = (3, 2, 1, 0.5)$.

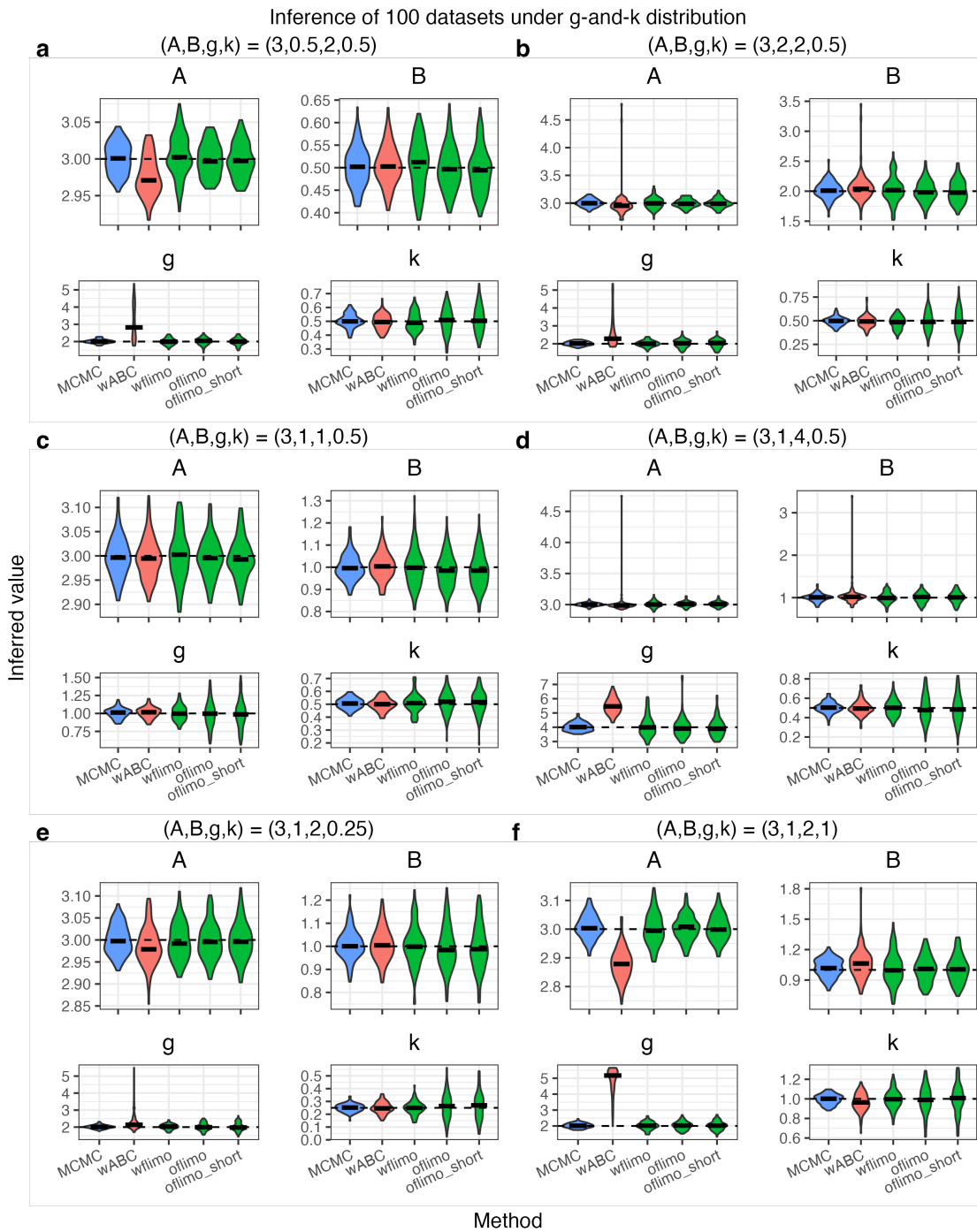


FIGURE S2 – Estimate of the parameters of g-and-k distributions for different parameter sets. The *oABC* and *wABC-short* methods have not been included due to their poor performance on the main example. In the case where $g = 4$ (panel d), i.e. when the peak of the g-and-k distribution is sharper, some *flimo* inferences tend to overestimate g . An attempt to improve this inference by jackknifing was made. For the cases studied, the median of one hundred inferred values with jackknifing is closer to the theoretical value. The average of the inferred parameters is a weaker estimator, due to a certain number of outliers. The inference results of g by *wABC* are biased on panels d and f. This is not the case for *flimo*. We did not try to correct the outliers isolated for A and B for the *wABC* method (panels b and d in particular). The times were not reported because the different inference methods were not used under the same conditions.

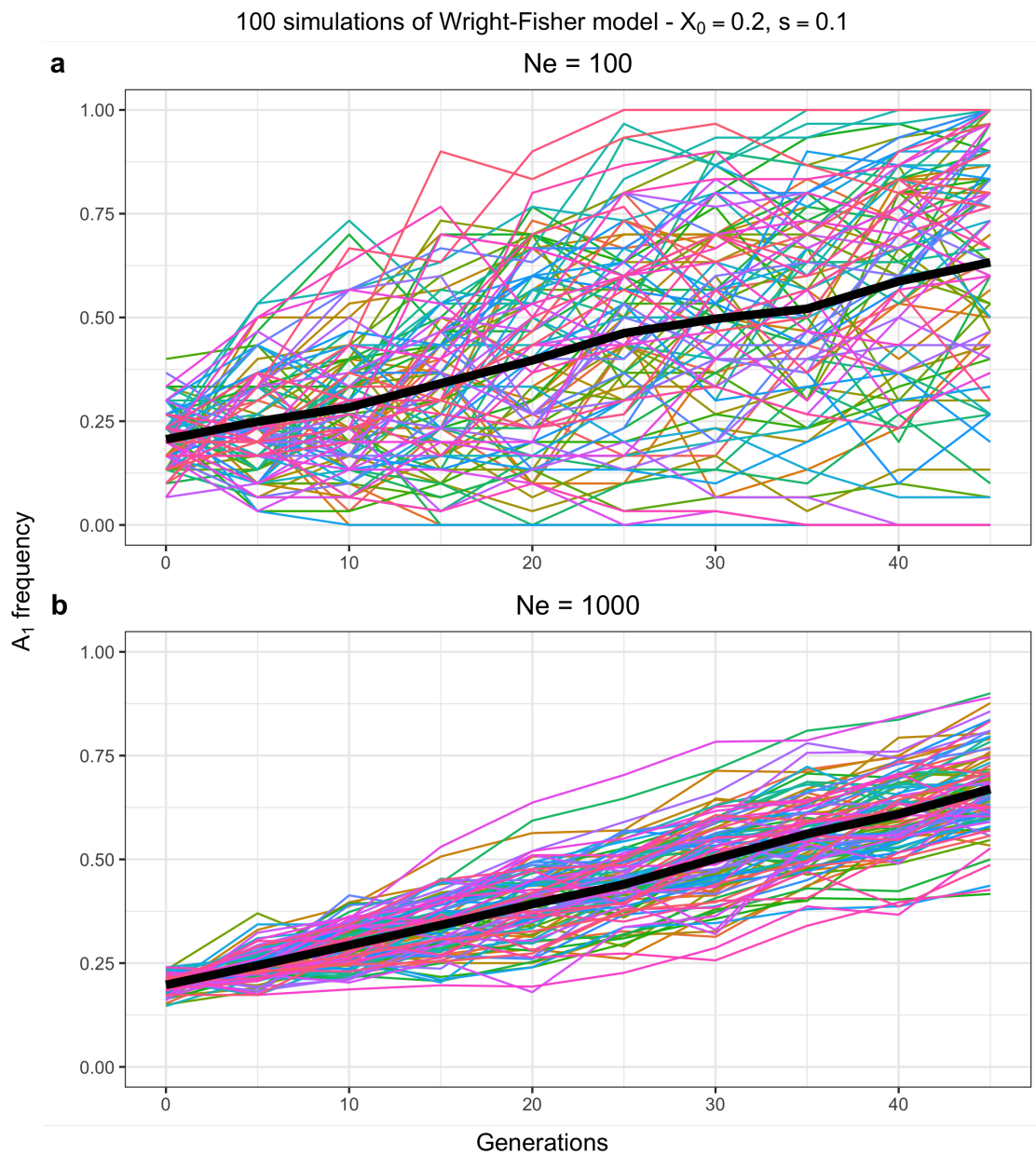


FIGURE S3 – One hundred simulated data sets of the Wright-Fisher model for $s = 0.1$ and two population sizes (**Panel a** : $N_e = 10^2$; **Panel b** : $N_e = 10^3$). Thick black line is the average A_1 proportion at each sampled time. **Panel a** : the dispersion from one simulation to another is important with a non-negligible probability of allele fixation.

Selective value	Criterion	<i>compareHMM</i>		<i>fimo</i>				
		<i>Bws</i>	<i>Binomial</i>		<i>Bws</i>		<i>NG</i>	
			10	200	10	200	10	200
$s = 0.01$	Correlation	0.998	0.86	0.89	0.85	0.89	0.88	0.90
	Median difference $\hat{s}_{fimo} - \hat{s}_{compareHMM-Bin}$	0.0 (0 %)	3.2×10^{-3} (32%)	2.5×10^{-3} (25%)	-5.0×10^{-4} (-5.0%)	2.8×10^{-3} (28%)	4.3×10^{-4} (4.3%)	9.5×10^{-4} (9.5%)
	Seconds by inference	1.4	0.038	0.76	0.34	6.4	0.013	0.10
	Correlation	0.98	0.75	0.72	0.65	0.69	0.64	0.69
$s = 1$	Median difference $\hat{s}_{fimo} - \hat{s}_{compareHMM-Bin}$	-6.7×10^{-3} (-0.67%)	-2.1×10^{-2} (-2.1%)	-1.5×10^{-2} (-1.5%)	-3.0×10^{-2} (-3.0%)	-1.8×10^{-2} (-1.8%)	-2.5×10^{-2} (-2.5%)	-2.2×10^{-2} (-2.2%)
	Seconds by inference	3.0	0.015	0.29	0.22	3.8	0.015	0.12
	Correlation	0.98	0.75	0.72	0.65	0.69	0.64	0.69

TABLE S1 – Estimation of a Wright-Fisher selection value - alternative scenarios. Inference results, based on 100 simulated data sets with $N_e = 10^3$ and $s \in \{0.01, 1\}$, using *compareHMM* or different implementations of the *fimo* method. Three quantities are presented : the Pearson correlation coefficient, the median of the difference between the values inferred by the *fimo* or *compareHMM-Bws*, and *compareHMM-Bin* methods, and the median of the computation times. For several inferences with $s = 0.01$, there is however a systematic overestimation of *fimo* compared to *compareHMM-Bin*. Note that in the latter case, the mean of the estimated s by *fimo* is closer to the simulated value $s = 0.01$ than with *compareHMM-Bin*. The large deviation from the theoretical value is due to important random fluctuations under these simulation conditions. For $s = 1$, the mean values inferred by *compareHMM* are more distant from the true value than those inferred by *fimo*.

Supplementary File 1 : Performance of *flimo* for a high dimension problem

We adapt here the toy example from Li et al. (2017). This basic example allows to evaluate the limits of the inference algorithms. In Li et al. (2017), the aim is to illustrate the performance in the construction of the posterior distribution for different ABC implementations. In our case, it is more a question of testing the limits in numbers of parameters reasonably inferable by *flimo*. The model considered is given by the equation A.1.

$$\begin{aligned} y &\sim \mathcal{N}_p(\theta, \Sigma) \quad p \geq 2 & (A.1) \\ \text{with } y &= (y_1, \dots, y_p)^T \\ \theta &= (\theta_1, \dots, \theta_p)^T \\ \Sigma &= \text{diag}(\sigma_0, \dots, \sigma_0) \end{aligned}$$

We use the prior chosen by the authors to determine the initial condition of the inference (equation A.2). As in their study, we fix $\sigma_0 = 1$ and $b = 0.1$.

$$\begin{aligned} \theta &\sim \mathcal{N}(0, A) \text{ with } A = \text{diag}(100, 1, \dots, 1) & (A.2) \\ \text{and then transforming } \theta_2 &\leftarrow \theta_2 + b\theta_1^2 - 100b \end{aligned}$$

The data consist of a single observation $y_{obs} = (10, 0, \dots, 0)^T$ and the summary statistics are the data vector itself : $s(y) = y$. The distance between the summary statistics is the Euclidean distance. The objective function for *flimo* is therefore given by the equation A.3. To evaluate *flimo*, we choose $n_{sim} = 10$.

$$J(\theta) = \sqrt{\sum_{i=1}^p (\overline{s(y_\theta)_i} - s(y_{obs})_i)^2} \quad (A.3)$$

where $\overline{s(y_\theta)}$ is the mean of the n_{sim} simulations produced by *flimo*.

In our case, *flimo* does not construct a posterior distribution. Therefore the evaluation criterion chosen is not a Kullback-Leibler divergence for the distribution of (θ_1, θ_2) but the Euclidean distance between the data and the inferred values for the first two components (equation A.4).

$$\text{Error}(\hat{\theta}) = (y_{obs,1} - \hat{\theta}_1)^2 + (y_{obs,2} - \hat{\theta}_2)^2 \quad (A.4)$$

The values of p tested range from 2 to 50. Table S2 shows the results. The median inference time seems to be proportional to p^a with $2 \leq a \leq 3$ and longer extreme cases are common. The error criterion is stable as p increases.

p	Mean Error	Maximum Error	Median Time	Maximum Time
2	0.44	0.98	0.0049	0.13
5	0.42	0.84	0.019	0.29
10	0.40	1.0	0.086	4.1
20	0.40	1.4	0.40	16
30	0.40	1.3	1.2	47
40	0.34	0.97	2.4	85
50	0.38	0.91	9.8	162

TABLE S2 – Inference accuracy and computation time (in seconds) as a function of the number of parameters p of the multivariate normal model with $n_{sim} = 10$, for 100 replicates.

Supplementary File 2 : Assessment of convergence in a toy example

In some very simple cases, it is possible to show analytically that the parameters inferred by *flimo* are convergent estimators of the model parameters. In this example, the observed data are drawn in a normal distribution (equation A.5). The summary statistics are the empirical mean and variance (equation A.6).

The goal is to infer $\theta^* = (\mu^*, \sigma^*)$. In practice, we want to show that the estimators built by *flimo* converge to the empirical moments (\bar{y}, σ_y^2) . In this case, it is of course matching with the moment estimators.

$$Y^{obs} = (y_1, \dots, y_n) \sim \mathcal{N}(\mu^*, \sigma^*) \text{ iid}, \quad n \geq 2 \quad (\text{A.5})$$

$$\bar{y} = \text{Mean}(Y^{obs}) \quad ; \quad \sigma_y^2 = \text{Var}(Y^{obs}) \quad (\text{A.6})$$

In the following, $\theta = (\mu, \sigma)$, $\mu \in \mathbb{R}$, $\sigma > 0$. Let $m \geq 2$ be the number of simulations to generate. There is no need for $n = m$. Let $z_\theta(q) = \mu + \sigma\sqrt{2}\text{erf}^{-1}(2q - 1)$, $q \in]0, 1[$ be the quantile function of the distribution $\mathcal{N}(\mu, \sigma)$. The quantiles used for the simulations are fixed : let $Q = (q_1, \dots, q_m) \sim \mathcal{U}([0, 1]) \text{ iid}$ and $\text{sim}Q : \theta \mapsto (z_\theta(q_1), \dots, z_\theta(q_m))$ be the simulation function. $\text{sim}Q$ entries are independent draws in $\mathcal{N}(\mu, \sigma)$ by inverse transform sampling. The chosen objective function to minimize to find optimal parameters is given by equation A.7.

$$J_Q(\theta) = \left(\bar{y} - \overline{\text{sim}Q(\theta)} \right)^2 + \left(\sigma_y^2 - \text{Var}(\text{sim}Q(\theta)) \right)^2 \quad (\text{A.7})$$

$$\text{where } \overline{\text{sim}Q(\theta)} = \frac{1}{m} \sum_{i=1}^m z_\theta(q_i) = \frac{1}{m} \sum_{i=1}^m \mu + \sigma\sqrt{2}\text{erf}^{-1}(2q_i - 1)$$

$$= \mu + \frac{\sigma\sqrt{2}}{m} \sum_{i=1}^m \text{erf}^{-1}(2q_i - 1) = \mu + \sigma \overline{e(Q)}$$

$$\text{with } \overline{e(Q)} = \frac{\sqrt{2}}{m} \sum_{i=1}^m \text{erf}^{-1}(2q_i - 1)$$

$$\text{and } \text{Var}(\text{sim}Q(\theta)) = \frac{1}{m} \sum_{i=1}^m \left(z_\theta(q_i) - \overline{\text{sim}Q(\theta)} \right)^2 = \frac{1}{m} \sum_{i=1}^m z_\theta(q_i)^2 - (\mu + \sigma \overline{e(Q)})^2$$

The computation with the unbiased estimator of the variance (normalized by $\frac{1}{m-1}$) does not lead to a close form of the solution. The formula is developed in equation A.8.

$$\begin{aligned}
\frac{1}{m} \sum_{i=1}^m z_{\theta}(q_i)^2 &= \frac{1}{m} \sum_{i=1}^m (\mu + \sigma\sqrt{2}\text{erf}^{-1}(2q_i - 1))^2 = \mu^2 + \sigma^2 \overline{e^2(Q)} + 2\mu\sigma \overline{e(Q)} \\
\text{with } \overline{e^2(Q)} &= \frac{2}{m} \sum_{i=1}^m \text{erf}^{-1}(2q_i - 1)^2 \\
\text{so } \text{Var}(\text{sim}Q(\theta)) &= \mu^2 + \sigma^2 \overline{e^2(Q)} + 2\mu\sigma \overline{e(Q)} - (\mu + \sigma \overline{e(Q)})^2 \\
&= \sigma^2 \left(\overline{e^2(Q)} - \overline{e(Q)}^2 \right) \\
J(\theta) &= \bar{y}^2 + \mu^2 + \sigma^2 \overline{e^2(Q)} + 2\mu\sigma \overline{e(Q)} - 2\bar{y}(\mu + \sigma \overline{e(Q)}) + \\
&\quad \sigma_y^4 + \sigma^4 \left(\overline{e^2(Q)} - \overline{e(Q)}^2 \right)^2 - 2\sigma_y^2 \sigma^2 \left(\overline{e^2(Q)} - \overline{e(Q)}^2 \right) \tag{A.8}
\end{aligned}$$

Let $\hat{\theta} = \text{argmin}J(\theta)$. Its explicit value can be obtain by equation A.9.

$$\begin{aligned}
\frac{\partial J}{\partial \mu}(\theta) &= 2\mu - 2\bar{y} + 2\sigma \overline{e(Q)} \\
\frac{\partial J}{\partial \sigma}(\theta) &= 2\sigma \overline{e^2(Q)} + 2\mu \overline{e(Q)} - 2\bar{y} \overline{e(Q)} + 4\sigma^3 \left(\overline{e^2(Q)} - \overline{e(Q)}^2 \right)^2 \\
&\quad - 4\sigma_y^2 \sigma \left(\overline{e^2(Q)} - \overline{e(Q)}^2 \right)
\end{aligned}$$

Solving $\frac{\partial J}{\partial \mu}(\hat{\theta}) = 0$ and $\frac{\partial J}{\partial \sigma}(\hat{\theta}) = 0$ we get $\hat{\mu} = \bar{y} - \hat{\sigma} \overline{e(Q)}$ and then

$$\hat{\sigma} = \frac{\sigma_y}{\sqrt{\left(\overline{e^2(Q)} - \overline{e(Q)}^2 \right)}} \quad \text{and} \quad \hat{\mu} = \bar{y} - \sigma_y \frac{\overline{e(Q)}}{\sqrt{\left(\overline{e^2(Q)} - \overline{e(Q)}^2 \right)}} \tag{A.9}$$

We can compute the limit when $m \rightarrow +\infty$. Recall that

$$\overline{e(Q)} = \frac{\sqrt{2}}{m} \sum_{i=1}^m \text{erf}^{-1}(2q_i - 1) \quad \text{and} \quad \overline{e^2(Q)} = \frac{2}{m} \sum_{i=1}^m \text{erf}^{-1}(2q_i - 1)^2$$

with $(\sqrt{2}\text{erf}^{-1}(2q_i - 1))_i$ being iid draws in a distribution $\mathcal{N}(0, 1)$. So $\overline{e(Q)}$ is a convergent estimator of the mean of $\mathcal{N}(0, 1)$: $\lim_{m \rightarrow \infty} \overline{e(Q)} = 0$. In the same way, $\overline{e^2(Q)} - \overline{e(Q)}^2$ is a convergent estimator of the variance of $\mathcal{N}(0, 1)$, so $\lim_{m \rightarrow \infty} \overline{e^2(Q)} - \overline{e(Q)}^2 = 1$. To conclude,

$$\lim_{m \rightarrow \infty} \hat{\mu} = \bar{y} \quad , \quad \lim_{m \rightarrow \infty} \hat{\sigma} = \sigma_y$$

Annexe B

Annexes au manuscrit *Towards quantitative DNA Metabarcoding*

Je reproduis ici les fichiers de la section *Supplementary Information* tels qu'ils sont inclus dans le manuscrit *Towards quantitative DNA Metabarcoding* du Chapitre 2.

Species	Barcode	1C value (pg)
<i>Briza media</i> (Bme)	atccgtgtttgagaaaacaaggggttctcgaa ctagaatacaaaggaaaag	6.35 ¹
<i>Rosa canina</i> (Rca)	atcccgttttatgaaaacaacaaggtttcagaa agcgagaataaataaag	1.40
<i>Lotus corniculatus</i> (Lco)	atcctgctttacgaaaacaagggaagttcagtt aagaaagcgacgagaaaaatg	0.87 ¹
<i>Populus tremula</i> (Ptr)	atcctatTTTTcgaaaacaacaaaaaacaac aaaggttcataaagacagaataagaatacaaaag	0.45
<i>Salvia pratensis</i> (Spr)	atcctgttttctcaaaaacaaggttcaaaaaacg aaaaaaaaaag	0.46
<i>Lonicera xylosteum</i> (Lxy)	atccagtttccgaaaacaaggtttagaaagca aaaatcaaaaag	0.70
<i>Fraxinus excelsior</i> (Fex)	atcctgttttcccaaaaacaaggttcagaaagaa aaaag	0.84
<i>Acer campestre</i> (Aca)	atcctgttttacgagaataaaaacaagcaaaaca gggttcagaaagcgagaaaggg	1.02 ¹
<i>Capsella bursa-pastoris</i> (Cbp)	atcctggtttacggaacacaccggagtttaca agcgagaaaaaagg	0.40
<i>Geranium robertianum</i> (Gro)	atcctttttacgaaaataaagaggggctcaca agcgagaatagaaaaaag	1.76 ²
<i>Carpinus betulus</i> (Cbe)	atcctgttttcccaaaaacaataaaaacaattta aggggttcataaagcgagaataaaaaag	1.03
<i>Abies alba</i> (Aal)	atccggttcataagaaaagggtttctctccttc tcctaaggaaagg	17.29
<i>Rhododendron ferrugineum</i> (Rfe)	atcctttttcgcaaaaacaagaattccgaaa gctaaaaaaaag	0.74 ²

SUPPLEMENTARY TABLE S1 – Metabarcodes and genome sizes of the plants used for the three mock communities for the *Sper01* marker. In our data, the sequence of *Salvia pratensis* has an insertion (a, position 35) compared to the reference sequence. 1C value characterises the genome size and can be found in the Kew database (<https://cvalues.science.kew.org/>). ¹ : average value of different assays. ² : average C-value of species in same genus when species C-value is missing.

Species	Concentration (ng/ μ l)			Copies per well (2 μ l of DNA)			Rank
	\mathcal{M}_U	\mathcal{M}_T	\mathcal{M}_G	\mathcal{M}_U	\mathcal{M}_T	\mathcal{M}_G	
<i>Briza media</i>	0.076	0.038	4.9×10^{-1}	1.9×10^4	9.7×10^3	1.2×10^5	1
<i>Rosa canina</i>	0.045	0.038	1.5×10^{-1}	1.9×10^4	1.6×10^4	6.2×10^4	2
<i>Lotus corniculatus</i>	0.068	0.038	1.1×10^{-1}	1.9×10^4	1.1×10^4	3.1×10^4	3
<i>Populus tremula</i>	0.038	0.038	3.1×10^{-2}	1.9×10^4	1.9×10^4	1.6×10^4	4
<i>Salvia pratensis</i>	0.062	0.038	2.5×10^{-2}	1.9×10^4	1.2×10^4	7.8×10^3	5
<i>Lonicera xylosteum</i>	0.13	0.038	2.6×10^{-2}	1.9×10^4	5.7×10^3	3.9×10^3	6
<i>Fraxinus excelsior</i>	0.11	0.038	1.1×10^{-2}	1.9×10^4	6.7×10^3	1.9×10^3	7
<i>Acer campestre</i>	0.071	0.038	3.6×10^{-3}	1.9×10^4	1.0×10^4	9.7×10^2	8
<i>Capsella bursa-pastoris</i>	0.25	0.038	6.3×10^{-3}	1.9×10^4	3.0×10^3	4.8×10^2	9
<i>Geranium robertianum</i>	0.12	0.038	1.2×10^{-3}	1.9×10^4	6.3×10^3	2.4×10^2	10
<i>Carpinus betulus</i>	0.18	0.038	1.2×10^{-3}	1.9×10^4	4.1×10^3	1.2×10^2	11
<i>Abies alba</i>	0.11	0.038	3.4×10^{-4}	1.9×10^4	6.9×10^3	6.1×10^1	12
<i>Rhododendron ferrugineum</i>	0.25	0.038	4.0×10^{-4}	1.9×10^4	2.9×10^3	3.0×10^1	13

SUPPLEMENTARY TABLE S2 – Total DNA concentrations (in the samples) and number of molecules per well of plants used for the three mock communities. The rank for \mathcal{M}_G stands for the decreasing abundance in terms of target DNA.

Species	Probe name	Probe sequence	Positions	Tm (salt-adjusted)
<i>Fraxinus excelsior</i>	Fex	ttttccaaaacaaggttcagaaagaaaa	7 - 36	63.9°C
<i>Capsella bursa-pastoris</i>	Cbp	aacacaccggagtttacaagcgag	16 - 40	65.8°C
<i>Carpinus betulus</i>	CbeA	tcctgtttccaaaacaataaaacaaat	1 - 30	62.5°C
	CbeB	ttaagggttcataaagcgagaataaaaaag	32 - 61	63.9°C

SUPPLEMENTARY TABLE S3 – Taqman probes designed for three different species for the *Sper01* marker. For *Carpinus betulus*, two probes were designed. The positions designate the bases of the original metabarcodes.

Annexe C

Modèle de PCR avec mismatch (Chapitre 2)

On cherche à simplifier l'écriture du nombre de molécules s et s' donnée par le système en 2.27 :

$$\begin{pmatrix} M_n \\ M'_n \end{pmatrix} = \begin{pmatrix} 1 + \pi\Lambda & (1 - \pi)\rho\Lambda \\ (1 - \pi)\rho\Lambda & 1 + \pi\Lambda \end{pmatrix}^n \begin{pmatrix} M_0 \\ M'_0 \end{pmatrix} \quad (\text{C.1})$$

Pour simplifier l'écriture, notons provisoirement

$$A = \begin{pmatrix} 1 + \pi\Lambda & (1 - \pi)\rho\Lambda \\ (1 - \pi)\rho\Lambda & 1 + \pi\Lambda \end{pmatrix} = \begin{pmatrix} a & b \\ b & a \end{pmatrix} \quad (\text{C.2})$$

La matrice A étant symétrique réelle, elle est diagonalisable. Pour $v \in \mathbb{R}$, on calcule :

$$\det(A - vI) = \begin{vmatrix} a - v & b \\ b & a - v \end{vmatrix} = (a - v)^2 - b^2 \quad (\text{C.3})$$

L'équation $\det(A - vI) = 0$ admet donc deux solutions distinctes : $v_- = a - b$ et $v_+ = a + b$. On trouve aisément des vecteurs propres de A :

$$A \begin{pmatrix} 1 \\ 1 \end{pmatrix} = (a + b) \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad \text{et} \quad A \begin{pmatrix} 1 \\ -1 \end{pmatrix} = (a - b) \begin{pmatrix} 1 \\ -1 \end{pmatrix} \quad (\text{C.4})$$

On peut donc écrire

$$A = P \begin{pmatrix} a + b & 0 \\ 0 & a - b \end{pmatrix} P^{-1} \quad (\text{C.5})$$

$$\text{avec } P = \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix} \quad \text{et son inverse } P^{-1} = \frac{1}{2} \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix} \quad (\text{C.6})$$

$$\begin{aligned}
\text{Donc } A^n &= PA^n P^{-1} \\
&= \frac{1}{2} \begin{pmatrix} (a+b)^n + (a-b)^n & (a+b)^n - (a-b)^n \\ (a+b)^n - (a-b)^n & (a+b)^n + (a-b)^n \end{pmatrix} \quad (\text{C.7})
\end{aligned}$$

avec :

$$\begin{cases} (a+b)^n + (a-b)^n = (1 + \pi\Lambda + (1 - \pi)\rho\Lambda)^n + (1 + \pi\Lambda - (1 - \pi)\rho\Lambda)^n \\ (a+b)^n - (a-b)^n = (1 + \pi\Lambda + (1 - \pi)\rho\Lambda)^n - (1 + \pi\Lambda - (1 - \pi)\rho\Lambda)^n \end{cases} \quad (\text{C.8})$$

On peut ainsi expliciter le nombre de molécules de s et s' avec $M'_0 = 0$:

$$\begin{pmatrix} M_n \\ M'_n \end{pmatrix} = \frac{M_0}{2} \begin{pmatrix} (1 + (\pi + \rho - \pi\rho)\Lambda)^n + (1 + (\pi - \rho + \pi\rho)\Lambda)^n \\ (1 + (\pi + \rho - \pi\rho)\Lambda)^n - (1 + (\pi - \rho + \pi\rho)\Lambda)^n \end{pmatrix} \quad (\text{C.9})$$

Annexe D

Calculs du modèle de mutation (Chapitre 3)

On cherche à simplifier l'écriture du nombre de molécules souches et mutantes la matrice donnée par le système en 3.1 :

$$\begin{pmatrix} M_{k+1}^s \\ M_{k+1}^m \end{pmatrix} = \begin{pmatrix} 1 + \Lambda(1 - \mu) & 0 \\ \Lambda\mu & 1 + \Lambda \end{pmatrix} \begin{pmatrix} M_k^s \\ M_k^m \end{pmatrix} \quad (\text{D.1})$$

soit par récurrence, avec $M_0^m = 0$:

$$\begin{pmatrix} M_n^s \\ M_n^m \end{pmatrix} = \begin{pmatrix} 1 + \Lambda(1 - \mu) & 0 \\ \Lambda\mu & 1 + \Lambda \end{pmatrix}^n \begin{pmatrix} M_0^s \\ 0 \end{pmatrix} \quad (\text{D.2})$$

Pour simplifier l'écriture, notons provisoirement

$$A = \begin{pmatrix} 1 + \Lambda(1 - \mu) & \Lambda\mu \\ \Lambda\mu & 1 + \Lambda(1 - \mu) \end{pmatrix} = \begin{pmatrix} a & 0 \\ b - a & b \end{pmatrix} \quad (\text{D.3})$$

On diagonalise la matrice A . Pour $v \in \mathbb{R}$, on calcule :

$$\det(A - vI) = \begin{vmatrix} a - v & 0 \\ b - a & b - v \end{vmatrix} = (a - v)(b - v) \quad (\text{D.4})$$

L'équation $\det(A - vI) = 0$ admet donc deux solutions distinctes strictement positives $v_- = a$ et $v_+ = b$. On trouve aisément des vecteurs propres de A :

$$A \begin{pmatrix} 1 \\ -1 \end{pmatrix} = a \begin{pmatrix} 1 \\ -1 \end{pmatrix} \quad \text{et} \quad \begin{pmatrix} 0 \\ 1 \end{pmatrix} = b \begin{pmatrix} 0 \\ 1 \end{pmatrix} \quad (\text{D.5})$$

On peut donc écrire

$$A = P \begin{pmatrix} a & 0 \\ 0 & b \end{pmatrix} P^{-1} \quad (\text{D.6})$$

$$\text{avec } P = \begin{pmatrix} 1 & 0 \\ -1 & 1 \end{pmatrix} \text{ et son inverse } P^{-1} = \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix} \quad (\text{D.7})$$

$$\begin{aligned} \text{Donc } A^n &= P A^n P^{-1} \\ &= \begin{pmatrix} a^n & 0 \\ b^n - a^n & b^n \end{pmatrix} \end{aligned} \quad (\text{D.8})$$

donc en reprenant les notations initiales, on obtient :

$$\begin{aligned} \begin{pmatrix} M_n^s \\ M_n^m \end{pmatrix} &= \frac{1}{2} \begin{pmatrix} (1 + \Lambda(1 - \mu))^n & 0 \\ (1 + \Lambda)^n - (1 + \Lambda(1 - \mu))^n & (1 + \Lambda)^n \end{pmatrix} \begin{pmatrix} M_0^s \\ 0 \end{pmatrix} \\ &= M_0^s \begin{pmatrix} (1 + \Lambda(1 - \mu))^n \\ (1 + \Lambda)^n - (1 + \Lambda(1 - \mu))^n \end{pmatrix} \end{aligned} \quad (\text{D.9})$$

Annexe E

Marqueurs pour métabarcoding

Marqueurs *Sper01* et *Euka03* décrits par Taberlet et al. (2018).

Sper01

Target taxonomic group: Spermatophyta (seed plants)

NCBI taxid: 58024

Forward primer: GGGCAATCCTGAGCCAA

Reference: Taberlet *et al.* (2007)

Reverse primer: CCATTGAGTCTCTGCACCTATC

Reference: Taberlet *et al.*

(2007)

Recommended annealing temperature: 52°C

Target gene: P6 loop of the *trnL* intron, chloroplast DNA

Coverage for the target group: 98.8% (739 species amplified in silico out of 748)

Min. length: 10 bp

Mean length: 48 bp

Max. length: 220 bp

Taxonomic resolution in the target group:

Species

Genus

Family

Order

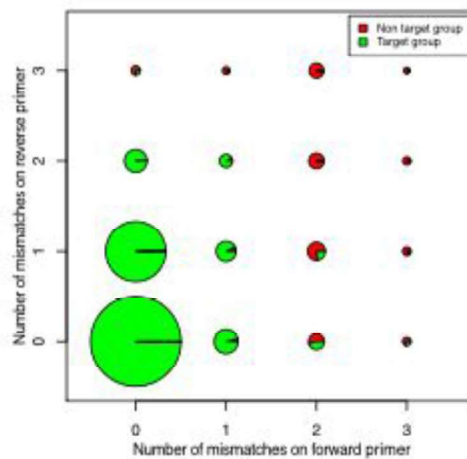
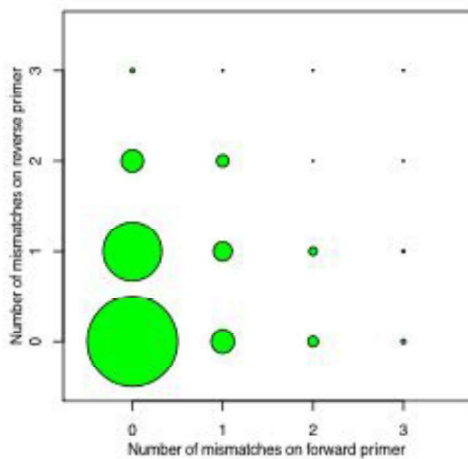
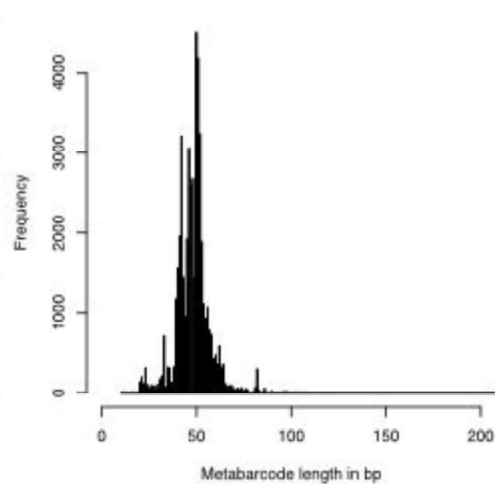
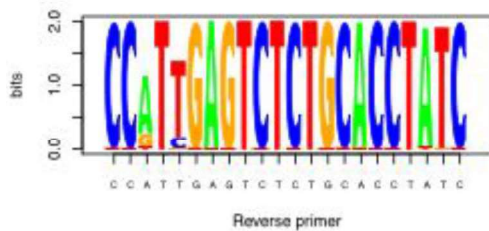
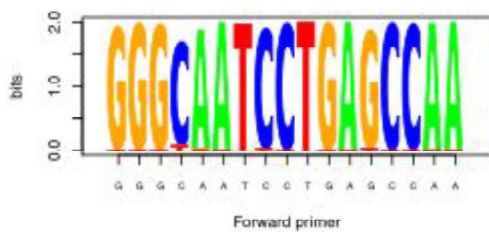
21.5% (48494)

36.9% (8417)

77.4% (371)

89.6% (67)

Comments: Widely used for analyzing degraded template. Correspond to the g/h primers of Taberlet *et al.* (2007).



Euka03

Target taxonomic group: Eukaryota

Forward primer: CCCTTTGTACACACCGCC

Reverse primer: CTTCYGCAGGTTACCTAC

Recommended annealing temperature: 55°C

Target gene: 18S nuclear ribosomal DNA (V9)

Coverage for the target group: NA

Min. length: 49 bp

Mean length: 133 bp

NCBI taxid: 2759

Reference: this book

Reference: this book

Max. length: 264 bp

Taxonomic resolution in the target group:

Species	Genus	Family	Order
58.1% (11815)	74.3% (5160)	84.0% (2343)	88.6% (673)

Comments: Highly specific of eukaryotes. Do not amplify properly amphipods/isopods (see Peca02). Amplify the same region than the primers Euka04 described in Amaral-Zettler *et al.* (2009), but Euka03 primers are optimized to amplify a wider range of eukaryotes. The main advantage of this Euka03 when compared to Euka01 or Euka02 is the more homogenized size of the metabarcode, and a higher taxonomic resolution, but the V9 (Euka03) region has been less sequenced than the V7 region (Euka01 and Euka02), leading to a less comprehensive reference database.

