



HAL
open science

Classification et caractérisation de familles enzymatiques à l'aide de méthodes formelles

Gaëlle Garet

► **To cite this version:**

Gaëlle Garet. Classification et caractérisation de familles enzymatiques à l'aide de méthodes formelles. Informatique [cs]. Université de Rennes 1, 2014. Français. NNT: . tel-01096916v1

HAL Id: tel-01096916

<https://inria.hal.science/tel-01096916v1>

Submitted on 18 Dec 2014 (v1), last revised 2 Feb 2015 (v2)

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 / UNIVERSITÉ DE RENNES 1
sous le sceau de l'Université Européenne de Bretagne

pour le grade de
DOCTEUR DE L'UNIVERSITÉ DE RENNES 1

Mention : Informatique
École doctorale Matisse

présentée par
Gaëlle GARET

préparée à l'unité de recherche Inria/Irisa – UMR6074
Institut de Recherche en Informatique et Système Aléatoires
Composante universitaire : ISTIC

**Classification et
caractérisation
de familles enzy-
matiques à l'aide
de méthodes for-
melles**

**Thèse à soutenir à Rennes
le 16 décembre 2014**

devant le jury composé de :

Jean-Christophe JANODET

Professeur à l'Université d'Evry-Val-d'Essonne / *Rapporteur*

Amedeo NAPOLI

Directeur de recherche au Loria, Nancy / *Rapporteur*

Colin DE LA HIGUERA

Professeur à l'Université de Nantes / *Examineur*

Olivier RIDOUX

Professeur à l'Université de Rennes 1 / *Examineur*

Mirjam CZJEK

Directrice de recherche CNRS, Roscoff / *Examinatrice*

Jacques NICOLAS

Directeur de recherche à Inria, Rennes / *Directeur de thèse*

François COSTE

Chargé de recherche à Inria, Rennes / *Co-directeur de thèse*

*Ainsi en était-il depuis toujours. Plus les hommes accumulaient des connaissances,
plus ils prenaient la mesure de leur ignorance.*
Dan Brown, Le Symbole perdu

Tu me dis, j'oublie. Tu m'enseignes, je me souviens. Tu m'impliques, j'apprends.
Benjamin Franklin

Remerciements

Je remercie tout d'abord la région Bretagne et Inria qui ont permis de financer ce projet de thèse.

Merci à Jean-Christophe Janodet et Amedeo Napoli qui ont accepté de rapporter cette thèse et à Olivier Ridoux, Colin De La Higuera et Mirjam Czjzek pour leur participation au jury.

J'aimerais aussi dire un grand merci à mes deux directeurs de thèse : Jacques Nicolas et François Coste, qui m'ont toujours apporté leur soutien tant dans le domaine scientifique que personnel. Jacques, merci pour ta patience, ta gentillesse et ton soutien à toute épreuve. Un grand merci surtout pour tout ce que tu m'a apporté au niveau scientifique, entre autre pour formaliser les problèmes rencontrés, et pour ta pile de bibliographie amassée pendant des années :) François, merci de m'avoir permis de faire mon stage de master qui m'a conduit ici, merci pour les discussions, constructives ou non, qui m'ont aidé à mieux comprendre les grammaires et les concepts qu'elles impliquent. Merci aussi de m'avoir donné l'occasion de découvrir le monde durant ces trois ans aux travers des voyages, je me souviendrais longtemps des voyages, Washington à l'occasion d'ICGI, Cordoba, les visites, le parc des condors ;)

Merci à tous ceux qui ont collaborés de près ou de loin à mes travaux, je pense notamment aux biologistes de la station de Roscoff qui m'ont toujours accueillis avec beaucoup de patience (et il en faut pour expliquer les concepts biologiques à une informaticienne !). En vrac, Myrjam, Gurvan, Thierry, Agnès, Catherine et tout ceux que j'oublie sur le moment.

Merci aussi à ceux qui m'ont accueilli dans leur labo en Argentine ou en Allemagne lors de mes différents séjours.

Un grand merci au groupe Symbiose (Dyliss, Genscale et Genouest), surtout ne changez rien !! Malgré les multiples départs et les arrivées, l'ambiance est toujours restée la même : excellente, ces 3 ans n'auraient pas été les mêmes sans chacun d'entre vous. Je me rappellerai des séminaires au vert, des midi activités, notamment les tournois de ping pong, l'escalade, le badminton (si si des fois j'y allais), les jeudi soirs amarillys, les afters de folie ou les samedi matins au marché, le voyage en Roumanie et j'en passe ... Tellement de bons souvenirs !

Merci à mathilde de m'avoir supportée dans le bureau pendant ces 3 ans (enfin plutôt 2 ;)), les tours de smarties, les aprem entre filles, ...

Merci à ceux qui m'ont permis de squatter leur bureau sur la fin (Julie, Charles et Ayme-

rick) et à Claudia qui m'a supportée pendant la fin de la rédaction. Merci à Catherine, Coline ou encore Anaïs pour les pauses potins quand la rédaction devenait difficile. Un merci particulier à ceux qui ont toujours été là à n'importe quelle heure du jour ou de la nuit, mathilde, vincent, sylvain : le stage de M2, les discussions mc do, les conférences, les séjours à Roscoff, ...

Merci à tout ceux qui m'ont fait découvrir la vulgarisation scientifique et qui nous ont aidé dans l'organisation de Sciences en Cour[t]s, une aventure à la fois scientifique et humaine qui a permis de penser à autre chose de temps en temps. Notamment les orgas, parmi eux Sylvain, Coraline, Charles, Marie, Kévin, Nico, ... mais aussi tous les réals ! Et je souhaite bon courage aux nouveaux (entre autre Clovis, Fanny et Paulin) qui ont accepté de reprendre le flambeau pour continuer l'aventure !

Merci à mes amis de Quimper qui m'ont écouté pendant un certain nombre d'heure parler de cette thèse sans rien comprendre 'Je pense notamment à Camille, Elise, Céline, JR, Maureen, Kane, ... Merci à vous d'être là depuis nos années déjantées du lycée.

Et merci à toute ma famille qui m'a soutenue et encouragée dans cette aventure qu'est la thèse (même quand je n'étais pas toujours facile à vivre !), maman, Nico, Pierre, Audrey, Caro,

Enfin, merci à tout ceux que je n'ai pas cité mais qui se reconnaîtront.

Remerciements		5
Table des matières		8
Introduction		13
1 Modélisation de familles d'enzymes		17
1.1 Problème biologique : reconnaissance de classes d'enzymes		17
1.1.1 Un peu d'histoire		17
1.1.2 Qu'est-ce qu'une enzyme ?		18
1.1.3 Des familles de séquences		20
1.1.4 Le problème de l'annotation d'enzymes à partir de séquences		22
1.2 Modélisation d'une famille à partir d'un ensemble de séquences		24
1.2.1 Alignement de séquences		24
1.2.2 Découverte de modèles réguliers à partir d'un alignement multiple		27
2 Inférence grammaticale sur des séquences biologiques		31
2.1 Apprendre un langage		31
2.1.1 Langages et grammaires, quelques définitions		32
2.1.2 Cadre théorique de l'apprentissage d'un langage		34
2.1.3 La généralisation comme recherche dans un espace d'hypothèses		34
2.1.4 Apprenabilité : le cadre de l'identification à la limite		36
2.1.5 Validation du langage appris		37
2.2 Inférence de grammaires régulières		38
2.2.1 État de l'art		38
2.2.2 Application aux séquences de protéines		39
2.2.3 Protomata-Learner		40
2.2.4 Discussion		42

3	Apprentissage de grammaires par substituabilité	45
3.1	Apprentissage de grammaires algébriques	45
3.1.1	Les difficultés de l'apprentissage de grammaires algébriques	46
3.1.2	Apprentissage de la structure d'un langage	48
3.1.3	Concepts formels et formalisation du principe de substituabilité	51
3.1.4	Discussion	53
3.2	Apprendre des langages substituables	54
3.2.1	Langages formels et substituabilité locale	54
3.2.2	Comparaison des classes de langage substituables	56
3.2.3	Propriétés de clôture	58
3.2.4	La substituabilité comme principe de généralisation	60
3.2.5	Un premier algorithme générique d'apprentissage pour les langages substituables	62
3.3	Apprentissage de grammaires réduites	64
3.3.1	Grammaires réduites	64
3.3.2	Apprentissage d'une grammaire réduite : l'algorithme ReGLiS	65
3.3.3	Complexité de l'algorithme ReGLiS	68
3.3.4	Exemple de réduction d'une grammaire pour le langage naturel	70
3.4	Expérimentations	72
3.4.1	Comparaison des temps d'exécution sur des données simulées	72
3.4.2	Processus d'apprentissage sur les séquences de protéines	73
3.4.3	Résultats d'apprentissage sur des familles de protéines	76
4	Classification de séquences par analyse de concepts formels	79
4.1	Analyse de concepts formels à partir d'un PLMA	80
4.1.1	Codage des séquences d'enzymes	80
4.1.2	Observation d'un lien séquence/structure sur une superfamille multifonctionnelle	80
4.2	Annotation via l'analyse de concepts formels	83
4.2.1	Formalisation du problème de classification	83
4.2.2	Classification supervisée	88
4.2.3	Classification non supervisée	90
4.3	Expérimentation sur des superfamille d'enzymes	92
4.3.1	Expérimentation sur des jeux de données connus	92
4.3.2	Expérimentations sur des données réelles	97
4.3.3	Apprentissage de grammaire sur une superfamille de séquences	99
4.4	Conclusion de cette approche	100
5	Conclusions et perspectives	103
5.1	Les contributions apportées	103
5.2	Perspectives	105
	Bibliographie	107

Annexes	119
A. Alignement partiel local partiel de la superfamille des GH16	119
B. Résultats de classement obtenus sur les séquences HAD d' <i>Ectocarpus</i> . .	121
C. Grammaire obtenue la superfamille des GH16	129
Table des figures	147

Introduction

Prédire l'activité d'une enzyme à partir de sa séquence en acides aminés est une tâche d'une importance majeure pour la compréhension des réactions biochimiques avec de nombreuses retombées dans le domaine des biotechnologies. Cette thèse considère le problème de l'apprentissage de signatures caractéristiques de familles d'enzymes. Les enzymes sont des protéines particulières et on dispose déjà d'un certain nombre d'outils pour caractériser des familles de protéines et découvrir de nouveaux membres à ces familles. Le problème général est celui de l'inférence de modèles à partir d'un ensemble de séquences partageant une fonction commune. Parmi les modèles les plus utilisés, on trouve les modèles de Markov cachés ou encore les langages réguliers. Dans cette thèse nous nous intéressons à des modèles algébriques plus expressifs qui sont capables de capturer des interactions entre éléments éloignés sur la séquence de la protéine.

Dans ce but, nous nous sommes intéressés dans un premier temps à l'inférence de grammaires hors-contexte. Il s'agit d'un problème pour lequel il existe un certain nombre d'algorithmes heuristiques pour le traitement automatique des langues. Tous sont basés sur le principe de substituabilité de Harris pour lequel le travail de A. Clark [CE07] offre un cadre théorique et des résultats d'apprenabilité en introduisant la classe des langages hors-contexte substituables. Nous avons étendu ce principe en introduisant de nouvelles classes de langages et de nouveaux critères de généralisation basés sur des classes de substituabilité locales et/ou contextuelles. Ces nouveaux critères permettent ainsi de traiter des ensembles de séquences de protéines pour lesquels les exemples sont moins nombreux et les séquences beaucoup plus longues que pour les langues naturelles.

Afin de pouvoir utiliser ce modèle en pratique sur des ensembles de données réels, nous proposons un algorithme d'apprentissage efficace permettant de réduire la grammaire tout au long de l'apprentissage jusqu'à l'obtention d'une grammaire canonique non redondante du langage substituable cible. Nous avons ainsi obtenu de bons résultats sur des données réelles avec une haute spécificité et une bonne sensibilité grâce à une généralisation basée sur la substituabilité locale.

En pratique, la disponibilité des familles d'enzymes n'est pas immédiate. Il existe des groupes de séquences phylogénétiquement liées, appelés superfamilles, qui possèdent des motifs communs pour les repérer. Au sein de ces superfamilles, il existe un certain nombre de familles avec des activités diverses déterminées expérimentalement. Le coût et la difficulté des expérimentations ne permettent pas actuellement d'avoir des séquences d'apprentissage pour chaque famille existante.

Nous avons donc abordé le problème de la prédiction de l'appartenance d'une séquence à une famille. Nous présentons un classifieur basé sur l'identification de blocs de sous-séquences communes entre des séquences de familles connues et inconnues appartenant à la même superfamille. Nous nous appuyons pour cela sur la recherche de concepts formels construits sur le produit des blocs et des séquences. Contrairement à la plupart des classifieurs dans ce domaine, nous avons introduit les classes (familles) comme des objets à part entière. Nous traitons le problème non supervisé de la détection de nouvelles familles dans les séquences non étiquetées comme un problème d'optimisation en minimisant le

nombre de nouvelles familles tout en maximisant le support d'une nouvelle famille en terme de blocs caractéristiques. Nous avons utilisé ce classifieur sur des données provenant du génome d'une algue brune récemment séquencée, *Ectocarpus siliculosus*, relativement éloignée des espèces habituellement étudiées.

La dernière étape de l'étude est la création de grammaires hors-contexte pour chaque famille repérée au sein d'une superfamille. Ces caractérisations expressives permettent potentiellement de repérer des interactions intéressantes au sein des protéines.

Dans un premier temps, le chapitre 1 est consacré à l'introduction des concepts biologiques nécessaires à la compréhension de cette thèse. Il présentera les notions d'enzymes et de familles de séquences, ainsi que certaines méthodes utilisées pour caractériser un ensemble de séquences protéiques. Nous verrons ainsi les limites des approches utilisées actuellement.

puis, le chapitre 2 introduit les concepts mathématiques et informatiques utilisés pour modéliser les familles d'enzymes. Nous présentons dans ce but la théorie des langages et l'inférence grammaticale, et insistons sur l'état de l'art utilisant des langages réguliers.

Dans le chapitre 3, nous cherchons à obtenir des grammaires de la classe des grammaires hors-contexte, plus expressives que les grammaires régulières classiquement utilisées. Dans ce but, nous introduisons les concepts de substituabilité locale et contextuelle comme critères sur lesquels s'appuiera la généralisation des séquences d'apprentissage.

Le chapitre 4 est consacré à l'apprentissage supervisé et non supervisé des propriétés des différentes sous-familles présentes dans l'échantillon d'apprentissage grâce à l'analyse de concepts formels. Il sera ainsi possible de discriminer des ensembles de séquences présentant des activités inconnues. Cela permettra d'introduire une phase de sélection dans la méthode présentée au chapitre 2.

Enfin, le chapitre 5 présentera les perspectives et conclusions de ce travail.

Résumé

Cette thèse propose une nouvelle approche de découverte de signatures de familles (et superfamilles) d'enzymes. Dans un premier temps, étant donné un échantillon aligné de séquences appartenant à une même famille, cette approche infère des grammaires algébriques caractérisant cette famille. Pour ce faire, de nouveaux principes de généralisation et de nouvelles classes de langages ont été introduites sur la base de la substituabilité locale. Un algorithme a également été développé à cet effet qui produit une grammaire réduite, conservant la structuration des exemples, d'un langage substituable. Dans un second temps, ce manuscrit présente une méthode de classification des séquences d'une superfamille en familles à l'aide d'une analyse de concepts formels basée sur l'alignement des séquences qui permet la détection de nouvelles familles et la découverte des motifs fonctionnels pour améliorer les signatures précédentes.

Mots clés

bioinformatique, enzyme, famille, inférence grammaticale, grammaire algébrique, substituabilité, analyse de concepts formels

Abstract

This thesis proposes a new approach to discover signatures of families (and superfamilies) enzymes. At first, given a sample of aligned sequences belonging to the same family, this approach infers context-free grammars characteristic of this family. To do this, new principles of generalization and new classes have been introduced based on substitutability. An algorithm has also been developed for this purpose, which produces a reduced grammar able to retain the structure of examples. In a second step, this manuscript presents a method for classification of a superfamily sequences into families with a formal concept analysis based on alignment sequences allowing detection of new families and the discovery of patterns to improve functional previous signatures.

Keywords

bioinformatics, enzyme, family, grammatical inference, context-free grammar, substitutability, formal concept analysis