

École Doctorale INTERFACES
Approches interdisciplinaires: fondements, applications et
innovations

### Titre de la thèse

### **THÈSE**

présentée et soutenue publiquement le XXX

pour l'obtention du

#### Doctorat de l'Université Paris-Saclay

spécialité: Les sciences du vivant

par

M. David MIGNON

#### Composition du jury

Rapporteurs: Dr. Prénom1 Nom1 Rapporteur externe

Dr. Prénom2 Nom2 Rapporteur externe Pr. Prénom3 Nom3 Rapporteur interne

Examinateurs: Dr. Prénom4 Nom4 Examinateur

Dr. Prénom5 Nom5 Directeur de thèse Dr. Prénom6 Nom6 Directeur de thèse

## Remerciements

XXX

à XXX.

# Table des matières

| Liste des figures            | vii |
|------------------------------|-----|
| Liste des tables             | ix  |
| Abreviations                 | xi  |
| Introduction                 | 1   |
| 1 CPD                        | 5   |
| 2 proteus                    | 7   |
| 3 Comparaisons d'algorithmes | 9   |
| 4 PDZ                        | 11  |
| Conclusion                   | 13  |
| Ribliographie                | 91  |

# Liste des figures

## Liste des tables

| 1 | Les tests avec cinq positions actives   | 16 |
|---|---|----|
| 2 | Les tests avec dix positions actives    | 17 |
| 3 | Les tests avec vingt positions actives  | 18 |
| 4 | Les tests avec trente positions actives | 19 |

### Abreviations

 ${f H}$  algorithme heuristique

MC algorithme Monte-Carlo

**RE** algorithme "Replica Exchange";

 $\mathbf{GMEC}$ "'Global minimal energie cost"

Pfam "Protein family databank"

# Introduction

XXX

### Contexte

XXX

XXX

Citation entre crochets [??].

Citation dans le texte?.

# CPD

# proteus

### Comparaisons d'algorithmes

#### 3.1 Les méthodes pratiques

Nous cherchons maintenant à déterminer les performances et les qualités des différents algorithmes de proteus. Pour évaluer les différents algorithmes de proteus, comme pour leur établir un paramétrage, nous effectuons des séries de tests. Grâce à l'algorithme de type toulbar2 il est possible d'obtenir la séquence/conformation qui possède la plus haute énergie de dépliement. Cela constitue une information important qui va nous servir d'élément de comparaison. Le facteur temps est également un élément déterminant. Il est dans certain cas limitant, nous ne savons pas à l'avance quand toulbar2 termine. Et il apparaît d'emblée illusoire d'espérer voir ce programme converger dans toutes les situations intéressantes dans un temps raisonnable. D'autres métriques qui caractérisent les séquences d'acides aminés de meilleurs énergies obtenues seront également utilisées pour les évaluations et pour les paramétrages.

Dans la suite, on appelle «position active», une position pour laquelle, tous les types d'acides et tous les rotamères de chaque type d'acide aminé sont autorisés, au court de la recherche de proteus. On désigne «séquence/conformation» une séquence d'acides aminés munie à chaque position d'un rotamère (le backbone étant de toute façon fixé). Tandis ce que le terme simple «séquence» sans plus de précision désigne une séquence d'acides aminés.

### 3.1.1 les protéines

Les tests sont effectués sur neuf protéines choisies pour avoir des longueurs de backbone variées, plusieurs domaines représentés, mais aussi plusieurs structures pour chaque famille présente. Ainsi l'ensemble se décompose en deux protéines SH3 de 56 et 57 résidus, de trois protéines PDZ de longueur comprise entre 82 et 97 résidus et enfin de trois protéines

#### Chapitre 3. Comparaisons d'algorithmes

SH2 longues de 105 ou 109 résidus.L'ensemble a une moyenne, arrondie à l'unité inférieure, de quatre-vingt-neuf positions, voir les détails en table ??.

| Code PDB | résidus   | nombre de positions | domaine |
|----------|-----------|---------------------|---------|
| 1ABO     | 64-119    | 56                  | SH3     |
| 1CKA     | 134-190   | 57                  | SH3     |
| 1R6J     | 192 - 273 | 82                  | PDZ     |
| 1G9O     | 9-99      | 91                  | PDZ     |
| 2BYG     | 186 - 282 | 97                  | PDZ     |
| 1BM2     | 55 - 152  | 98                  | SH2     |
| 1O4C     | 1-105     | 105                 | SH2     |
| 1M61     | 4-112     | 109                 | SH2     |
| 1A81     | 9-117     | 109                 | SH2     |

Table 3.1 – Les protéines

| Protein | 1ABO        | 1CKA        |
|---------|-------------|-------------|
| 1ABO    | 100 (6e-42) | 26 (1e-07)  |
| 1CKA    | 26 (1e-07)  | 100 (2e-41) |

Table 3.2 – Pourcentage d'identité et e-value des alignements Blast native vs native pour nos protéine SH3.

| Protein      | 1R6J                     | 1G9O                     | 2BYG             |
|--------------|--------------------------|--------------------------|------------------|
| 1R6J<br>1G9O | 100(1e-59)<br>25 (3e-07) | 25 (3e-07)<br>100(2e-66) | no<br>35 (2e-11) |
| 2BYG         | no                       | 35 (2e-11)               | 100(7e-71)       |

Table 3.3 – Pourcentage d'identité et e-value des alignements Blast native vs native pour nos protéine PDZ (no= pas de touche avec une e-value inférieure à 10).

#### Alignements Blast croisés

| Protein | 1BM2        | 104C       | 1M61       | 1A81       |
|---------|-------------|------------|------------|------------|
| 1BM2    | 100 (7e-74) | 36 (2e-16) | 38(6e-10)  | 35 (1e-13) |
| 104C    | 36 (2e-16)  | 100(2e-79) | 27(3e-10)  | 33 (2e-12) |
| 1M61    | 38 (6e-10)  | 27 (3e-10) | 100(6e-81) | 57 (2e-47) |
| 1A81    | 35 (1e-13)  | 33 (2e-12) | 57(2e-47)  | 100(5e-83) |

Table 3.4 – Pour centage d'identité et e-value des alignements Blast native vs native pour nos protéine SH2.

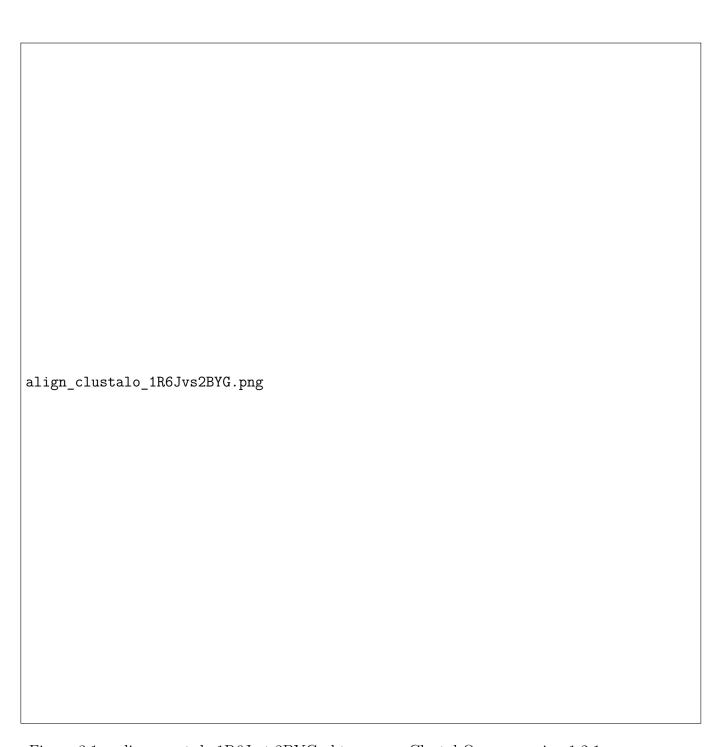


Figure 3.1 – alignement de 1R6J et 2BYG obtenu avec Clustal Omega version  $1.2.1\,$ 

#### 3.1.2 Description des tests

Les tests sont répartis en deux ensembles :

- 1. un ensemble de tests où toutes les positions de la séquence sont actives (cela correspond aux situations de design complet de protéines)
- 2. un ensemble de tests où le nombre de positions actives est gardé sous contrôle de façon à maîtriser la taille de l'espace d'exploration

Ensemble «Tout actif» Pour le premier ensemble de tests, la totalité de la matrice d'énergie est exploitée et pour chaque position l'espace d'exploration correspond à l'espace d'état déclaré dans le fichier ".bb". C'est-à-dire que tous les types de résidu et tous les rotamères sont possibles à chaque position. Comme l'espace des séquences/conformations à explorer est gigantesque, nous ne faisons pas de tentatives de recherche du GMEC par méthode exacte.

Nous effectuons des recherches avec les algorithmes suivants :

- heuristique, noté H par la suite;
- Monte-Carlo, noté MC;
- «Replica Exchange», noté RE);

L'ensemble «nombre d'actifs limité» L'ensemble «Nombre d'actifs limité» est composé de six groupes de tests avec un nombre de positions actives fixe définit de la façon suivante :

- 1. aucune position active
- 2. une position active
- 3. cinq positions
- 4. dix positions
- 5. vingt positions
- 6. trente positions

Lorsqu'une position n'est pas active, l'acide aminé de la position est fixé en utilisant l'acide aminé de la séquence native. La chaîne latérale est, elle, laissée libre. Il n'y a donc jamais dans nos tests de position où l'état est complètement fixé.

Le groupe «aucune position active» n'est constitué que d'un test par algorithme pour chaque protéine. Il y a donc neuf tests par algorithme. Ce sont les tests pendant lesquels la séquence d'acides aminés est fixe et correspond à la séquence native de la protéine.

Pour les tests avec une seule position active, comme des temps de calcul le permettent, nous décidons d'être exhaustifs : Toutes les positions sont testées, il y a alors huit cent quatre tests par algorithme. Pour tous les autres groupes de tests (cinq,dix,vingt et trente positions actives), cinq tests sont effectués par protéine, c'est-à-dire quarante-cinq tests par algorithme.

le choix des positions actives Pour définir complètement les tests, il reste maintenant à décrire le choix des positions actives pour les groupes de numéro trois jusqu'à six. Il y a peu d'intérêt à tester des situations avec des positions actives sans interaction entre-elles. En effet, s'il existe une position active P dont chaque résidu est sans interaction avec tous les résidus possibles des autres positions actives, déterminer le meilleur état pour P est proche du test du groupe 2 avec P comme position active. Toutefois, cela n'est pas exactement la même question, parce que les positions actives différentes de P peuvent influencer la position de la chaîne latérale de positions inactives qui à leur tour peuvent influencer l'état de P. Ainsi, le choix des positions actives ne se fait non pas par tirage aléatoire, car le risque d'obtenir des positions avec peu d'interactions est trop grand. Il se fait sous contrainte d'interaction.

**positions en interactions** Pour cela, nous utilisons la notion de voisinage de proteus. Elle se définit de la façon suivante : Deux positions P et Q sont en interactions s'il existe un rotamère  $r_P$  de P et un rotamère  $r_Q$  de Q tels que :

$$|E(r_P,r_O)| > S_{Vois}$$

avec  $S_{Vois}$  un seuil donné par l'utilisateur à la configuration de proteus (voir chap.?? pour les détails).

Alors on appelle «n-uplet en interaction» la donnée de n positions avec  $n \in \{5,10,20,30\}$  et d'un seuil  $S_{Vois}$  tels que pour toute paire de positions (P,Q) du n-uplet, P et Q sont en interactions.

choix des positions actives Pour définir les positions actives, nous exécutons proteus en mode verbeux, sans effectuer d'optimisation. Pour cela, il existe plusieurs façons de procéder, ici nous utilisons le mode Monte-Carlo avec une trajectoire de zéro pas. Ces exécutions produisent en sortie standard la liste des voisins pour chaque position au seuil donnée en paramètre. Pour chacune des neuf protéines, nous exécutons proteus avec  $S_{Vois}$  égal à dix, cinq et un à tour de rôle; trois listes de voisins sont obtenues. Ensuite, un script dédié recherche dans ces listes, les n-uplets en interaction, en partant de la liste de voisins

au sens le plus fort, c'est-à-dire dix, vers celle au sens le plus faible (0.1).La recherche s'arrête lorsque cinq n-uplets au moins sont trouvés.

Nous obtenons quarante-cinq n-uplets pour le groupe à cinq (respectivement dix, vingt et trente) positions actives pour un seuil  $S_{Vois}$  égal à dix (respectivement dix, un et un). Les positions actives de tous les tests sont en annexe ??). Pour chaque n-uplet, un fichier de configuration de proteus est créé dans lequel la balise <Space\_Constraints> fixe les positions inactives en utilisant le type d'acide aminé présent dans la séquence native.

#### 3.1.3 Définition de protocole comparable

Nous voulons comparer les algorithmes très différents. Un algorithme peut garantir l'obtention du minimum global en énergie (GMEC) si l'exécution se termine, mais ne garantit pas qu'elle se termine. Un autre permet un contrôle très fin du temps d'exécution sans garantie du GMEC, et d'autres enfin ont des objectifs plus large que la seule obtention du GMEC. Mais le GMEC reste le meilleur point de commun. Nous allons donc y concentrer une part importante des comparaisons.

Nous devons noter également que l'obtention du GMEC est théorique, en pratique nous n'avons pas de preuve que le code de l'algorithme exact que nous utilisons n'a pas de bogue. Cependant, nous mettons de côté cette éventualité et dans toute la suite GMEC désigne aussi bien le minimum global en énergie que le résultat de toulbar2 lorsqu'il se termine. Le Monte-Carlo et le «Replica exchange» possèdent de nombreux paramètres de configuration, ce qui rend l'ensemble des protocoles possibles très grand. Se pose alors la question de l'optimisation du protocole. L'objectif fixé ici, n'est pas la recherche d'un protocole optimal pour chacun des tests, mais d'évaluer, avec les tests, un protocole optimisé par algorithme. Nous allons alors dans un premier temps, recherche les meilleurs paramétrages pour le Monte-Carlo et le «Replica Exchange» sur l'ensemble de tests «tout actif». Puis, sur la base des résultats obtenus, les protocoles seront fixés pour effectuer les comparaisons sur l'ensemble «tout actif» et celui à «nombre d'actifs limité». Le programme toulbar2 possède aussi de nombreuses options. Deux paramétrages différents seront utilisés.

Pour rendre les protocoles comparables, le temps d'exécution maximum est fixé à vingtquatre heures pour tous les exécutions. Toulbar2 donne sa meilleure séquence/conformation en dernier, il n'y a donc pas post-traitement nécessaire. C'est également le cas pour le Monte-Carlo à condition de configurer l'impression de la trajectoire avec la balise  $Print\_Threshold = 0$ . dans le fichier de configuration. Pour le "Replica Exchange" et l'heuristique, un tri des séquences selon l'énergie est nécessaire. Mais il n'y a pas beaucoup de séquences :

- 1. L' Heuristique fournit une séquence/conformation à chaque cycle.
- 2. Le "Replica Exchange avec  $Print\_Threshold = 0$  produit autant de fichiers de séquences/rotamères que de marcheurs. Chacun ne contenant pas plus de quelques dizaines de séquences/rotamères.

Nous pouvons donc négliger la durée du tri dans le temps total d'exécution.

Protocole heuristique Pour l'algorithme heuristique, il n'y a dans notre situation qu'un seul paramètre à renseigner : le nombre de cycles à effectuer. Quelques essais préliminaires sur la plus grosse protéine (Table ??) avec toute les positions actives, montre que la version utilisée de proteus peut effectuer jusqu'à environ 110000 cycles sur nos machines de calculs en l'espace de vingt-quatre heures. Ainsi, le protocole H est défini comme le protocole qui utilise le mode heuristique de proteus et qui effectue cent dix mille cycles. Sont également définis les variantes H-, H+ et H++ comme des protocoles plus courts ou plus longs à facteur entier près (Table ??). Par ailleurs, certaines comparaisons de l'heuristique avec le Monte-Carlo ont été faites avec une version précédente du programme proteus. Ce protocole sera noté h. Il diffère aussi de H par le fait que l'option d'optimisation du compilateur Intel utilisé est -O2 contre -O3 pour H.

| Nom | nombre de cycles |
|-----|------------------|
| Н   | 110000           |
| Н-  | 1100             |
| H+  | 330000           |
| H++ | 990000           |
| h   | 100000           |

Table 3.5 – Les protocoles heuristiques

**Protocoles Monte-Carlo** On distingue deux ensembles de protocoles Monte-Carlo. Dans le premier, les noms sont de la forme "mc\*". Il rassemble les protocoles utilisés pour le paramétrage du Monte-Carlo. Le second est constitué des protocoles utilisés lors des comparaisons.

Les éléments à paramétrer pour l'algorithme Monte-Carlo sont les suivants :

- 1. la température
- 2. le nombre de pas (avec le nombre de trajectoires et la longueur de trajectoire )
- 3. Le seuil de voisinage
- 4. Les probabilités de changements de la séquence/conformation

Ce qui représente un ensemble de protocoles trop grand pour une approche exhaustive. Pour l'essentiel, nous allons faire varier les paramètres un par un, en prenant comme point de départ un protocole qui rend le comportement de marcheur Monte-Carlo «proche» de l'heuristique.

La température est le paramètre principal du Monte-Carlo, c'est elle qui contrôle le taux d'acceptation du critère de Metropolis. Alors, la première étape de cette optimisation va consister à faire varier la température, entre 0.001 et 0.5, en conservant les autres paramètres fixés (protocoles de mc0 à mc5). Le nombre de pas total effectué est le produit de deux paramètres, le nombre de trajectoires et la longueur de trajectoire. Les protocoles mc1b et mc2b testent l'effet d'une augmentation du nombre de pas. Tandis que mc2c et mc2d testent l'effet de la variation du nombre de trajectoires par rapport à la longueur. Le protocole mc2e s'intéresse aux probabilités de changement de la trajectoire. Il y a cinq balises dans proteus qui contrôle ces changements:

- <Prot> donne la probabilité de modifications de rotamère à une position.
- <Prot\_Prot> donne la probabilité de modifications de rotamère à deux positions.
- <Mut> donne la probabilité de modifications de type de résidu à une position.
- <Mut\_Prot> donne la probabilité de modifications de rotamères à deux positions.
- <Mut\_Mut> donne la probabilité de modifications de type de résidu à deux positions.

La table ?? donne les probabilités utilisées par ces cinq paramètres dans l'ordre de la liste précédente.

Enfin, mc4b se distingue des autres par un seuil de voisinage plus grand ((Table ??)).

Seconde version de proteus Pour la partie comparaison avec les autres algorithmes, quatre protocoles sont utilisés. Les protocoles MCa et MCb s'inspirent fortement de mc2d et mc2e, en étant adapté à la contrainte du temps de calcul de la comparaison et en utilisant la nouvelle version de proteus (les lettres capitales dans le nom des protocoles signifient l'utilisation de la dernière version de proteus). MCa- est une variante de MCa avec une trajectoire six fois plus courte. Enfin, MC0 s'inspire de mc0 dans le sens où la température est suffisamment froide pour que nous puissions considérer qu'il n'y a pas de baisse de l'énergie au cours d'une trajectoire.

**Protocoles "Replica Exchange"** L'algorithme «Replica Exchange» (RE) est une extension du Monte-Carlo. Les paramètres d'un protocole RE sont ceux d'un protocole Monte-Carlo plus trois autres :

| Nom  | Temp  | Long. de trajectoire(mega) | Nb de trajectoires | Voisin | Proba           |
|------|-------|----------------------------|--------------------|--------|-----------------|
| mc0  | 0.001 | 3                          | 1000               | 10     | 0; 1; 0.1; 0;0  |
| mc1  | 0.1   | 3                          | 1000               | 10     | 0; 1; 0.1; 0; 0 |
| mc2  | 0.2   | 3                          | 1000               | 10     | 0; 1; 0.1; 0; 0 |
| mc3  | 0.3   | 3                          | 1000               | 10     | 0; 1; 0.1; 0; 0 |
| mc4  | 0.5   | 3                          | 1000               | 10     | 0; 1; 0.1; 0; 0 |
| mc5  | 0.7   | 3                          | 1000               | 10     | 0; 1; 0.1; 0; 0 |
| mc1b | 0.1   | 6                          | 1000               | 10     | 1; 1; 1; 1; 0   |
| mc2b | 0.2   | 6                          | 1000               | 10     | 0; 1; 0.1; 0;0  |
| mc2c | 0.2   | 3                          | 10000              | 10     | 0; 1; 0.1; 0;0  |
| mc2d | 0.2   | 3000                       | 1                  | 10     | 0; 1; 0.1; 0;0  |
| mc2e | 0.2   | 3                          | 1000               | 10     | 1;0;0.1;0;0     |
| mc4b | 0.5   | 10                         | 100                | 10     | 0;1;0;1;0       |
| MC0  | 0.01  | 1000                       | 1                  | 10     | 1;0;0.1;0;0     |
| MCa  | 0.2   | 6000                       | 1                  | 10     | 1;0;0.1;0;0     |
| MCa- | 0.2   | 1000                       | 1                  | 10     | 1;0;0.1;0;0     |
| MCb  | 0.2   | 6000                       | 1                  | 10     | 0; 1; 0.1; 0; 0 |

Table 3.6 – Les protocoles Monte-Carlo

- le nombre de marcheurs
- la température pour chaque marcheur
- la période de «swap», c'est-à-dire la période (en nombre de pas) à laquelle le test de Hasting sur l'échange de température est effectué.

Pour avoir des exécutions en parallèle avec au plus un marcheur par coeur du processeur, nous limiter nos tests à quatre ou huit marcheurs. La distribution des températures est un élément déterminant dans le comportement des marcheurs, car c'est elle qui pilote en grande partie le taux d'acceptation des échanges de températures. Nous suivons l'idée proposée par Kofke de lui faire suivre une progression géométrique ( $\frac{T_i}{T_{i+1}} = C$ , avec C une constante) [???]. Ceci garantie alors que le taux d'acceptation d'échange entre  $T_ietT_{i+1}$  soit égale pour tout nos i.De plus, nous souhaitons centrer approximativement, nos distributions sur la température ambiante (environ 0.6 kcal/mol). Dans toute la suite, les températures et les énergies sont exprimées en kcal/mol.

Voici les températures pour le RE quatre marcheurs :

- 10, 1, 0.1 et 0.01
- 2, 1, 0.5 et 0.25
- -1, 0.5, 0.25 et 0.125

et celles pour le RE huit marcheurs :

-3, 2, 1.333, 0.888, 0.592, 0.395, 0.263 et 0.175

— 10, 3.16, 1, 0.316, 0.1, 0.0316, 0.01 et 0.00316

Ici les protocoles ne se font qu'avec une seule trajectoire par marcheur. Et la contrainte du temps de calcul se comprend comme vingt-quatre heures de calculs cumulées sur tous les marcheurs. Ainsi les longueurs de trajectoire sont définit pour le RE à quatre marcheurs comme le quart d'une trajectoire MC, pour le RE à huit marcheurs comme le huitième.

La table ?? donne les probabilités utilisées par les cinq balises qui contrôlent les modifications de la séquence/conformation à chaque pas, dans l'ordre de la liste de la section ??.

| Nom   | marcheurs | Temp         | Traj (mega) | seuil voisin | Proba          | swap period (me |
|-------|-----------|--------------|-------------|--------------|----------------|-----------------|
| RE4a  | 4         | 10<->0.01    | 1500        | 10           | 1;0;0.1;0;0    | 7.5             |
| RE4b  | 4         | 1 < -> 0.125 | 1500        | 10           | 1;0;0.1;0;0    | 7.5             |
| RE4c  | 4         | 2 < -> 0.25  | 1500        | 10           | 1;0;0.1;0;0    | 7.5             |
| RE8a1 | 8         | 10<->0.00316 | 750         | 0            | 1;0;0.1;0;0    | 2.5             |
| RE8a2 | 8         | 10<->0.00316 | 750         | 10           | 1;0;0.1;0;0    | 2.5             |
| RE8b1 | 8         | 3 < -> 0.175 | 750         | 10           | 0; 1; 0.1; 0;0 | 7.5             |
| RE8b2 | 8         | 3 < -> 0.175 | 750         | 10           | 1;0;0.1;0;0    | 7.5             |
| RE8b3 | 8         | 3 < -> 0.175 | 750         | 10           | 1;0;0.1;0;0    | 1               |
|       |           |              |             |              |                |                 |

Table 3.7 – Les protocoles «Replica Exchange»

Protocoles Toulbar2 Après avoir converti nos matrices au format «wcsp» grâce à un script dédié,nous pouvons utiliser toulbar2. Le protocole de recherche du GMEC est le suivant : L'exécutable toulbar2 de version 0.9.7.0 est lancé avec les options « -l=3 -m -d : -s», ce qui correspond au paramétrage conseillé dans la documentation CDP [??]. Si l'exécution se termine en moins de vingt-quatre heures, le protocole est achevé. Sinon le programme est arrêté et une seconde version (la 0.9.6.0) est lancée avec les options «-l=1 -dee=1 -m -d : -s». Au bout de vingt-quatre heures si le programme n'est pas terminé, il est arrêté. La dernière séquence/conformation imprimée en sortie est collectée. Le choix de la seconde version et du paramétrage fait suite à une discussion avec monsieur Seydou Traoré.

Toulbar2 offre également la possibilité de fournir la liste des séquences/conformations dont l'énergie est comprise entre celle qui correspond au GMEC,  $E_{GMEC}$  et une autre  $E_{upper\_bound}$  donnée en paramètre. Pour utiliser cette fonctionnalité nous utilisons le paramétrage : «-d : -a -s -ub= $E_{upper\_bound}$  ».Cependant, il s'avère que cette utilisation peut utiliser une quantité de mémoire vive importante. Alors, pour eviter tout plantage de nos machines, la mémoire que toulbar2 peut allouer est limité à 30 Go.

#### 3.1.4 Outils d'analyse des données

Superfamily/SCOP Superfamily [?] est un ensemble composé :

- D'une base de données de modèles de Markov cachés, où chaque modèle représente une structure 3D d'un domaine de la classification SCOP.
- D'une série de scripts qui annotent à partir des informations de la base,les séquences données en entrée. Ici, nous utilisons uniquement l'association au modèle 3D le plus vraisemblable.

Nous travaillons avec la base de données à la version 1.75, et en conjonction, nous utilisons SAM (version 3.5) [?] et HMMER (version 3.0) [?] recommandés par l'équipe de Superfamily. Le paramétrage utilisé est celui par défaut.

Taux d'identité de séquences Soient S et N deux séquences d'acides aminés de même longueur l.

Le Taux d'identité Id(S,N) de S par rapport N est égal au pourcentage de position où l'acide aminé est identique dans S et N. C'est-à-dire

$$Id(S,N) = \frac{\sum_{1 \le i \le l} \mathbb{1}(\hat{s_i}, n_i)}{l} \times 100$$

avec  $s_i$  et  $n_i$  l'acide animé en i de S et de N respectivement, et  $\mathbb{1}(x,y)$  la fonction qui vaut 1 lorsque x=y et 0 sinon.

Taux d'identité par position Le taux d'identité d'un alignement  $A_S$  à la position i par rapport à une séquence N de même longueur se définit comme :

$$Id(A_S,i) = \frac{\sum_{1 < j < m} \mathbb{1}(s_i^j, n_i)}{m} \times 100$$
, avec m le nombre de séquences de  $A_S$ .

Alignements Pfam Ce taux d'identité donne une mesure de la ressemblance entre un alignement et une séquence. Cela nous permet de comparer nos séquences calculées à la séquence native. Mais cela n'est pas notre seule objectif. Et nous voulons les évaluer par rapport à l'ensemble des séquences du domaine protéique de la native. La base de données Pfam (Protein families database) [?] regroupe les domaines protéiques connus en famille. Chaque famille étant représentée par des alignements multiples de séquences et des profiles de modèles de Markov cachés [?]. Dans la suite, nous n'utiliserons l'alignement dit « seed» qui se base sur un petit ensemble de membres représentatifs de la famille et l'alignement « full» , plus large, qui est généré par modèle de Markov caché à partir de l'alignement « seed». Les alignements correspondent pour nous aux familles PF00017 (domaine SH2), PF00018 (domaine SH3) et PF00595 (domaine PDZ).

Score BLOSUM Pour tenir compte des ressemblances et des différences entre les acides aminés lors d'une substitution, nous avons besoin d'une matrice de coût. Nous utilisons la matrice BLOSUM62 (BLOcks SUbstitution Matrix) [?] qui est construite à partir de blocs d'alignement très conservés (ici plus de 62% d'identités). Les fréquences des mutations y sont calculées. Le score BLOSUM d'une substitution est alors le logarithme de la fréquence de la mutation correspondante. À cela est ajouté un score de pénalités pour l'insertion d'un gap (c'est-à-dire un saut dans l'alignement).

On définit alors simplement un score de similarité de deux séquences de même longueur comme la somme des scores BLOSUM62 sur toutes les positions. De même le score de similarité d'un alignement par rapport à une séquence sera défini comme la moyenne des scores de similarité sur ensemble des séquences de l'alignement. Et enfin un score de similarité de deux ensembles de séquences alignés entre eux comme la moyenne des scores de similarité du premier ensemble par rapport aux séquences du second.

similarité d'un ensemble à une famille Pfam Afin de calculer un score de similarité d'un ensemble de nos séquences par rapport à une famille Pfam, il faut commencer par aligner nos séquences avec l'alignement de la famille. Pour cela nous utilisons le programme d'alignement BLAST [?]. Il implémente une heuristique qui recherche puis étend les meilleurs alignements locaux. Nous procédons comme suit :

- 1. La commande blastpgp est utilisée avec comme database (paramètre -d ) l'alignement Pfam et comme séquence en entrée ( paramètre -i ) la séquence native.
- 2. Dans la sortie blast, la séquence qui produit l'alignement le plus significatif avec la native est collectée, notons-la  $S_0$ .
- 3. L'alignement blast est alors utilisé pour positionner la native par rapport à  $S_0$  et les gaps nécessaires pour aligner la native à  $S_0$  sont ajoutés.
- 4. Le positionnement et les gaps sont alors appliqués tels quels à la liste de nos séquences.

Répartition de l'énergie selon les centiles Pour étudier différentes distributions d'ensemble de séquences/conformations selon l'énergie, nous déterminons les centiles de la façon suivante :

- 1. L'ensemble de séquences/conformations est trié selon l'énergie.
- 2. L'intervalle entre la meilleure énergie et la moins bonne est divisé en cent intervalles consécutifs contenant le même nombre de séquences/conformations (un centième du cardinal de l'ensemble).

#### Chapitre 3. Comparaisons d'algorithmes

3. Les quatre-vingt-dix-neuf valeurs d'énergie obtenues par ce découpage sont les centiles.

## PDZ

## Conclusion

XXX

Annexe 1 :Liste des positions actives pour chaque test

| Nom       | $S_{Vois}$ | positions actives      |
|-----------|------------|------------------------|
| 1A81 1    | 10         | 10 13 16 84 86         |
| 1A81 2    | 10         | 20 21 24 27 116        |
| 1A81 3    | 10         | $35\ 38\ 56\ 105\ 107$ |
| 1A81 4    | 10         | $44\ 47\ 52\ 65\ 67$   |
| 1A81 5    | 10         | 82 84 86 87 90         |
| 1ABO 1    | 10         | 64 66 90 93 100        |
| 1ABO 2    | 10         | 72 74 80 104 111       |
| 1ABO 3    | 10         | 79 82 102 111 115      |
| 1ABO 4    | 10         | 83 86 104 105 106      |
| 1ABO 5    | 10         | 93 100 102 113 116     |
| $1BM2\ 1$ | 10         | 101 106 140 141 146    |
| 1BM22     | 10         | 120 128 131 132 135    |
| 1BM23     | 10         | 58 61 127 128 129      |
| 1BM24     | 10         | 74 75 98 100 105       |
| 1BM25     | 10         | 85 87 95 110 128       |
| 1CKA 1    | 10         | 136 138 158 175 190    |
| 1CKA 2    | 10         | 149 166 169 171 181    |
| 1CKA 3    | 10         | 151 153 157 159 172    |
| 1CKA 4    | 10         | 164 170 172 184 187    |
| 1CKA 5    | 10         | 172 174 182 186 187    |
| 1G9O 1    | 10         | 10 13 54 57 92         |
| 1G9O 2    | 10         | 15 39 42 54 57         |
| 1G9O 3    | 10         | 24 26 28 39 42         |
| 1G9O 4    | 10         | 48 53 57 59 88         |
| 1G9O5     | 10         | 75 78 79 86 88         |
| $1M61\ 1$ | 10         | 12 20 23 24 27         |
| $1M61\ 2$ | 10         | 17 20 24 37 49         |
| $1M61\ 3$ | 10         | 27 33 51 100 102       |
| $1M61\ 4$ | 10         | 5 8 10 11 36           |
| $1M61\ 5$ | 10         | 59 71 84 87 94         |
| 104C 1    | 10         | 20 21 32 34 46         |
| 1O4C 2    | 10         | 2 71 79 81 82          |
| 104C 3    | 10         | 33 45 63 71 73         |
| 104C 4    | 10         | 43 45 63 71 85         |
| 104C 5    | 10         | 8 33 82 83 86          |
| 1R6J 1    | 10         | 194 237 239 270 272    |
| 1R6J 2    | 10         | 199 201 211 218 232    |
| 1R6J 3    | 10         | 213 218 227 232 238    |
| 1R6J 4    | 10         | 221 227 232 267 269    |
| 1R6J 5    | 10         | 241 254 258 267 269    |
| 2BYG 1    | 10         | 189 191 221 244 246    |
| 2BYG 2    | 10         | 205 224 239 245 248    |
| 2BYG 3    | 10         | 232 233 265 272 274    |
| 2BYG 4    | 10         | 238 240 243 276 278    |
| 2BYG 5    | 10         | 253 261 264 265 274    |

Table 1 – Les tests avec cinq positions actives  $\,$ 

| Nom              | $S_{Vois}$ | positions actives   |
|------------------|------------|---|
| 1A81 1           | 10         | 13 15 39 41 53 86 89 90 93 103                                  |
| 1A81 2           | 10         | 39 41 53 55 64 66 76 89 92 103                                  |
| 1A81 3           | 10         | 51 53 64 66 68 74 76 82 88 92                                   |
| 1A81 4           | 10         | 76 82 87 88 90 91 92 95 97 99                                   |
| 1A81 5           | 10         | 9 10 11 16 41 51 53 66 88 89                                    |
| 1ABO 1           | 10         | 64 72 74 79 89 91 101 103 108 111                               |
| $1ABO\ 2$        | 10         | 66 68 80 82 88 90 100 102 104 111                               |
| 1ABO 3           | 10         | 69 70 72 74 80 81 106 113 114 115                               |
| 1ABO 4           | 10         | 71 78 83 84 94 99 101 104 105 106                               |
| 1ABO5            | 10         | 72 79 82 94 99 102 104 106 111 115                              |
| $1BM2\ 1$        | 10         | 119 120 121 122 123 125 131 134 135 140                         |
| 1BM2 2           | 10         | 125 126 127 129 130 133 134 136 137 147                         |
| 1BM23            | 10         | 83 99 101 106 108 135 140 141 146 148                           |
| 1BM24            | 10         | 85 95 97 110 118 120 125 128 131 132                            |
| 1BM25            | 10         | 99 101 106 139 140 141 142 143 144 146                          |
| 1CKA 1           | 10         | 134 135 160 161 162 173 174 175 176 179                         |
| 1CKA 2           | 10         | 137 139 143 151 153 157 159 172 182 186                         |
| 1CKA 3           | 10         | 138 140 147 149 150 155 166 169 181 188                         |
| 1CKA 4           | 10         | 140 141 153 154 155 157 174 175 184 186                         |
| 1CKA 5           | 10         | 151 153 157 166 168 173 174 176 178 179                         |
| 1G9O 1           | 10         | 10 11 13 14 15 16 53 54 57 92                                   |
| 1G9O 2           | 10         | 15 17 24 26 39 42 48 51 53 88                                   |
| 1G9O 3           | 10         | 26 28 39 42 48 53 57 59 88 90                                   |
| 1G9O 4           | 10         | 34 35 58 60 68 70 74 75 89 91                                   |
| 1G9O 5           | 10         | 71 73 74 77 80 81 82 83 84 85                                   |
| 1M61 1           | 10         | 10 12 20 23 24 27 35 49 102 104                                 |
| 1M61 2           | 10         | 17 20 21 24 37 39 40 47 49 58                                   |
| 1M61 3           | 10         | 34 36 46 48 59 61 71 83 84 87                                   |
| 1M61 4<br>1M61 5 | 10         | 5 6 11 36 46 48 61 69 83 84                                     |
| 104C 1           | 10         | 59 61 70 71 75 77 83 86 87 92<br>31 33 45 47 61 63 73 86 89 100 |
| 104C 1<br>104C 2 | 10<br>10   | 50 51 52 53 63 72 73 77 85 89                                   |
| 104C 2<br>104C 3 | 10         | 61 62 63 71 72 73 79 85 88 89                                   |
| 104C 3<br>104C 4 | 10         | 73 74 75 76 77 89 92 94 96 101                                  |
| 104C 4<br>104C 5 | 10         | 90 91 93 96 98 99 101 102 103 104                               |
| 1R6J 1           | 10         | 193 194 195 197 199 218 232 236 267 269                         |
| 1R6J 2           | 10         | 199 209 211 213 218 227 232 238 265 267                         |
| 1R6J 3           | 10         | 201 204 205 209 211 218 241 258 265 267                         |
| 1R6J 4           | 10         | 209 211 213 218 227 238 241 258 265 267                         |
| 1R6J 5           | 10         | 238 240 241 242 246 257 258 261 265 267                         |
| 2BYG 1           | 10         | 194 196 203 205 224 233 239 245 274 276                         |
| 2BYG 2           | 10         | 203 205 207 224 227 233 239 243 245 276                         |
| 2BYG 3           | 10         | 206 207 222 245 248 253 256 261 264 265                         |
| 2BYG 4           | 10         | 221 222 245 248 251 253 256 261 264 265                         |
| 2BYG 5           | 10         | 247 248 249 250 251 252 259 262 263 275                         |

Table 2 – Les tests avec dix positions actives  $\,$ 

| Nom       | $S_{Vois}$ | positions actives   |
|-----------|------------|---|
| 1A81 1    | 1          | 9 11 12 13 15 16 17 19 20 25 28 39 40 41 42 43 44 45 114 117                    |
| 1A81 2    | 1          | 9 11 12 13 15 16 17 19 20 25 28 39 40 41 42 43 44 45 47 117                     |
| 1A81 3    | 1          | 9 11 12 13 15 16 17 19 19 41 43 48 51 68 74 84 86 109 114 117                   |
| 1A81 4    | 1          | 12 13 15 16 17 19 20 25 28 39 40 41 42 43 44 45 47 86 114 117                   |
| $1A81\ 5$ | 1          | 13 15 16 19 41 43 48 51 60 64 68 70 71 74 84 86 87 88 109 114 117               |
| 1ABO 1    | 1          | 64 66 67 68 82 86 87 88 89 90 91 101 102 102 103 103 108 111 113 116            |
| 1ABO 2    | 1          | 64 65 65 66 67 84 87 88 89 90 91 93 100 101 102 103 108 111 113 116             |
| 1ABO 3    | 1          | 65 66 67 87 88 89 90 91 93 94 95 100 101 102 103 106 108 111 113 116            |
| 1ABO 4    | 1          | 64 65 66 67 69 87 88 89 90 91 93 100 101 102 103 106 108 111 113 116            |
| 1ABO 5    | 1          | 66 67 68 82 86 87 88 89 90 91 93 100 101 102 103 106 108 111 113 116            |
| 1BM2 1    | 1          | 55 56 60 61 62 83 84 85 86 87 95 97 99 110 118 125 127 133 150 152              |
| 1BM22     | 1          | 55 56 60 61 62 83 84 85 86 87 95 97 99 110 118 125 127 128 129 152              |
| 1BM23     | 1          | 55 56 58 60 61 62 64 67 69 73 83 84 85 86 87 129 132 133 150 152                |
| 1BM24     | 1          | 55 56 60 61 62 69 83 84 85 86 87 95 97 99 110 129 132 133 150 152               |
| 1BM2 5    | 1          | 58 60 60 61 61 62 64 67 69 73 75 83 84 85 86 129 132 133 150 152                |
| 1CKA 1    | 1          | 134 135 136 137 138 139 150 151 160 161 162 163 164 170 171 172 173 179 189 190 |
| 1CKA 2    | 1          | 134 135 136 137 139 150 151 153 160 161 162 163 164 170 171 172 173 179 189 190 |
| 1CKA 3    | 1          | 134 136 137 139 150 151 157 158 160 161 162 163 164 170 171 172 173 179 189 190 |
| 1CKA 4    | 1          | 136 137 139 150 151 153 158 159 160 161 162 163 164 170 171 172 173 179 189 190 |
| 1CKA 5    | 1          | 137 139 150 151 153 158 160 161 162 163 164 170 171 172 173 174 175 179 189 190 |
| 1G9O 1    | 1          | 9 10 11 13 14 15 31 34 38 54 57 58 60 68 90 91 92 94 95 96                      |
| 1G9O 2    | 1          | 9 11 13 14 15 16 31 34 38 54 57 58 60 68 90 91 92 94 95 96                      |
| 1G9O 3    | 1          | 9 11 13 14 15 31 34 38 54 55 57 58 60 68 90 91 92 94 95 96                      |
| 1G9O 4    | 1          | 9 11 13 15 16 17 54 57 58 59 60 61 68 89 90 91 92 94 95 96                      |
| 1G9O 5    | 1          | 10 11 13 15 16 17 54 57 58 60 61 68 89 90 90 91 92 94 95 96                     |
| 1M61 1    | 1          | 34 35 36 37 38 39 42 46 47 48 49 50 61 63 69 71 77 78 81 82                     |
| 1M61 2    | 1          | 35 36 37 38 39 42 46 47 48 49 50 61 63 69 71 77 78 81 82 83                     |
| 1M61 3    | 1          | 38 39 42 46 47 48 49 50 61 63 69 71 77 78 81 82 83 84 85 87                     |
| 1M61 4    | 1          | 42 46 47 48 49 50 61 63 69 71 77 78 81 82 83 84 85 87 88 98                     |
| 1M61 5    | 1          | 5 7 50 61 63 69 71 77 78 81 82 83 84 85 87 88 98 103 104 109                    |
| 104C 1    | 1          | 32 33 34 35 43 45 63 65 71 73 79 80 81 82 83 84 85 86 87 89                     |
| 104C 2    | 1          | 3 4 5 7 8 9 11 17 31 32 33 34 35 43 45 63 65 71 73 79                           |
| 104C 3    | 1          | 1 3 4 5 6 7 8 9 11 17 31 32 33 35 43 45 65 81 82 83                             |
| 104C 4    | 1          | 1 2 3 4 5 6 7 8 9 11 12 13 14 17 19 35 65 81 82 83                              |
| 104C 5    | 1          | 1 3 4 5 6 7 8 9 11 12 17 31 32 33 34 35 65 81 82 83                             |
| 1R6J 1    | 1          | 193 194 195 197 214 215 217 218 233 235 236 237 239 240 241 242 247 269 270 273 |
| 1R6J 2    | 1          | 193 194 197 198 199 217 233 235 236 237 238 239 240 241 242 247 268 270 272 273 |
| 1R6J 3    | 1          | 193 195 197 217 233 235 236 239 240 241 242 244 245 247 268 269 270 270 272 273 |
| 1R6J 4    | 1          | 193 195 197 217 233 235 236 237 239 241 242 244 245 247 268 269 270 272 273 273 |
| 1R6J 5    | 1          | 193 194 197 198 199 233 236 237 239 239 240 241 247 268 268 269 270 270 272 273 |
| 2BYG 1    | 1          | 186 187 188 189 190 191 192 215 216 219 244 246 270 271 273 274 278 280 281 282 |
| 2BYG 2    | 1          | 186 187 188 189 190 215 216 219 221 223 240 243 270 271 273 274 278 280 281 282 |
| 2BYG 3    | 1          | 186 187 188 189 190 215 216 219 221 223 240 243 244 270 271 273 278 280 281 282 |
| 2BYG 4    | 1          | 186 187 188 189 190 215 216 219 221 223 240 243 244 270 274 276 278 280 281 282 |
| 2BYG 5    | 1          | 187 189 190 191 192 215 216 219 221 223 240 243 244 270 274 276 278 280 281 282 |

Table 3 – Les tests avec vingt positions actives

| Nom       | $S_{Vois}$ | positions actives   |
|-----------|------------|---|
| 1A81 1    | 1          | 9 11 12 13 15 16 17 19 20 25 26 27 28 29 36 38 39 40 41 42 43 48 51 68 74 84 86 109 114 117                             |
| $1A81\ 2$ | 1          | 9 10 11 12 13 15 16 17 19 20 25 28 39 41 43 48 51 68 74 83 84 86 87 88 90 91 93 109 114 117                             |
| $1A81\ 3$ | 1          | 9 11 12 13 15 16 17 19 20 25 27 28 36 38 39 40 41 41 42 43 43 48 51 68 74 84 86 109 114 117                             |
| $1A81\ 4$ | 1          | 9 10 11 12 13 15 16 17 19 20 25 28 36 39 40 41 42 43 43 44 45 48 51 68 74 84 86 109 114 117                             |
| $1A81\ 5$ | 1          | 9 10 11 12 13 15 16 17 19 20 25 28 39 40 41 42 43 44 45 47 48 51 52 68 74 84 86 109 114 117                             |
| 1ABO 1    | 1          | 64 65 66 67 68 70 71 72 75 78 79 80 81 82 83 86 87 88 89 90 91 93 100 101 102 103 108 111 113 116                       |
| $1ABO\ 2$ | 1          | 64 65 66 67 68 72 75 78 80 81 82 83 84 86 87 88 89 90 91 93 94 100 101 102 103 104 108 111 113 116                      |
| 1ABO 3    | 1          | 64 66 67 68 70 71 72 78 82 86 87 88 89 90 91 93 94 95 96 99 100 101 102 103 104 105 108 111 113 116                     |
| 1ABO 4    |            | 64 65 66 67 70 71 72 68 82 86 87 88 89 90 91 93 94 95 96 99 100 101 102 103 104 105 108 111 113 116                     |
| 1ABO 5    | 1          | 65 66 67 70 71 72 75 78 80 82 86 87 88 89 90 91 93 94 95 96 99 100 101 102 103 108 111 113 116                          |
| $1BM2\ 1$ | 1          | 55 56 58 60 61 62 83 84 85 86 87 95 97 99 110 118 119 120 121 122 125 127 128 129 130 131 132 133 150 152               |
| 1BM2 2    | 1          | 56 58 60 61 62 83 84 85 86 87 95 97 99 110 118 119 120 121 122 123 125 127 128 129 130 131 132 133 150 152              |
| 1BM23     | 1          | 58 60 61 62 83 84 85 86 87 95 97 99 108 109 110 118 120 121 122 123 125 127 128 129 132 133 134 135 150 152             |
| 1BM24     | 1          | 55 56 58 60 61 62 83 84 85 86 87 95 97 99 108 109 110 118 120 121 125 127 128 129 130 131 132 133 150 152               |
| 1BM25     | 1          | 56 58 60 61 62 67 83 84 85 86 87 95 97 99 110 111 112 113 115 118 125 127 128 129 130 131 132 133 150 152               |
| 1CKA 1    | 1          | 134 135 136 137 139 140 141 142 143 144 146 147 148 149 150 151 158 159 160 161 162 163 164 170 171 172 173 179 189 190 |
| 1CKA 2    |            | 134 135 136 137 139 143 144 146 147 148 149 150 151 158 159 160 161 162 163 164 170 171 172 173 179 186 187 188 189 190 |
| 1CKA 3    | 1          | 135 136 137 139 144 146 147 148 149 150 151 157 158 159 160 161 162 163 163 164 170 171 172 173 179 186 187 188 189 190 |
| 1CKA 4    | 1          | 136 137 139 140 141 142 143 144 146 147 148 151 158 159 160 161 162 163 164 170 171 172 173 179 184 186 187 188 189 190 |
| 1CKA 5    | 1          | 134 136 137 139 140 141 142 143 144 146 147 148 151 158 159 160 161 162 163 164 170 171 172 173 179 182 187 188 189 190 |
| 1G9O 1    | 1          | 9 10 11 13 15 24 31 34 38 40 41 42 43 46 48 49 50 51 54 57 58 60 68 89 90 91 92 94 95 96                                |
| 1G9O2     | 1          | 9 11 13 15 31 32 34 38 40 41 42 43 46 48 49 50 51 54 57 58 60 68 89 90 90 91 92 94 95 96                                |
| 1G9O 3    |            | 9 10 11 13 15 31 32 34 38 40 41 42 43 46 48 49 50 51 54 57 58 60 68 89 90 91 92 94 95 96                                |
| 1G9O4     | 1          | 10 11 13 14 15 31 32 34 38 40 41 42 43 46 48 49 50 51 54 57 58 60 61 68 89 90 91 92 94 95 96                            |
| 1G9O5     | 1          | 10 11 13 14 15 31 32 40 41 42 43 46 48 49 50 51 54 57 58 60 61 62 68 87 89 90 91 92 94 95 96                            |
| $1M61\ 1$ | 1          | 12 14 15 20 34 35 36 37 38 39 42 46 47 48 49 50 61 63 69 71 77 78 81 82 83 84 85 87 88 98                               |
|           | 1          | 6 7 8 10 11 12 14 15 20 21 34 35 36 37 38 39 42 46 47 48 49 50 61 63 69 71 77 78 81 82                                  |
| $1M61\ 3$ | 1          | 5 7 35 36 37 38 39 42 46 47 48 49 50 61 63 69 71 77 78 81 82 83 84 85 87 88 98 103 104 109                              |
|           | 1          | 7 8 10 11 12 14 15 20 34 35 36 37 38 39 42 46 47 48 49 50 61 63 69 71 77 78 81 82 83 84                                 |
| 1M61 5    | 1          | 8 10 11 12 14 15 20 34 35 36 37 38 39 42 46 47 48 49 50 61 63 69 71 77 78 81 82 83 84 85                                |
|           | 1          | 1 2 3 4 5 6 7 8 9 11 17 31 32 33 34 35 43 45 63 65 71 81 82 83 90 91 92 93 94 96  |
|           | 1          | 1 3 4 5 7 8 9 11 17 31 32 33 34 35 43 45 63 65 71 73 79 80 81 82 83 90 91 92 93 96                                      |
|           | 1          | 1 3 4 5 6 7 8 9 11 17 31 32 33 34 35 43 45 63 65 71 73 79 80 81 82 83 90 91 92 93 96                                    |
|           | 1          | 1 3 4 5 7 8 9 11 17 31 32 33 34 35 43 45 63 65 71 73 79 80 81 82 83 84 91 92 93 96                                      |
| 104C 5    | 1          | 1 3 4 5 6 7 8 9 11 17 31 32 33 34 35 43 45 63 65 71 73 79 80 81 82 83 84 92 93 96                                       |
| 1R6J 1    | 1          | 193 194 195 197 198 199 217 218 219 220 221 222 223 224 225 226 227 228 229 233 235 236 237 239 247 268 269 270 272 273 |
| 1R6J 2    | 1          | 193 194 195 197 198 199 217 220 221 222 223 224 225 226 227 228 229 233 235 235 236 237 239 240 247 268 269 270 272 273 |
| 1R6J 3    | 1          | 193 194 195 197 198 199 208 217 220 221 222 223 224 225 226 227 228 229 230 233 235 236 237 239 247 268 269 270 272 273 |
| 1R6J 4    | 1          | 193 194 195 197 198 199 217 220 221 222 223 224 225 226 227 228 229 233 235 236 237 239 240 247 249 268 269 270 272 273 |
| 1R6J 5    | 1          | 194 195 197 198 199 217 220 221 222 223 224 225 226 227 228 229 233 235 236 237 239 240 247 249 250 268 269 270 272 273 |
| 2BYG 1    |            | 186 187 188 189 190 191 192 215 216 219 221 223 240 243 244 246 250 251 252 253 254 255 256 257 259 260 278 280 281 282 |
| 2BYG 2    |            | 186 187 188 189 190 191 192 198 215 216 219 221 223 240 243 244 251 252 253 254 255 256 257 259 260 278 278 280 281 282 |
| 2BYG 3    |            | 186 187 188 189 190 190 191 192 215 216 219 221 223 240 243 244 251 252 253 254 255 256 257 259 260 278 246 280 281 282 |
| 2BYG 4    |            | 186 187 188 189 190 191 192 215 216 219 221 223 240 243 244 245 246 251 252 253 254 255 256 257 259 260 278 278 281 282 |
| 2BYG 5    | 1          | 186 187 188 189 190 191 192 215 216 219 221 223 240 243 244 251 252 253 254 255 256 257 259 260 278 246 278 280 281 282 |

Table 4 – Les tests avec trente positions actives

### Résum

### Titre de la thèse

XXX

Mots-clés : motclé1, motclé2, motclé3

### Abstract

Thesis title

XXX

 $\textbf{Keywords:} \ \text{keyword1}, \ \text{keyword2}, \ \text{keyword3}$