



République Algérienne Démocratique et Populaire



Ministère de l'Enseignement Supérieur et de la Recherche Scientifique

Université des Sciences et de la Technologie Houari Boumediene

Faculté d'Informatique
Département Informatique

Spécialité : Bio-Informatique

Rapport Module Architecture et Calcul Parallèle

Thème

Programmation parallèle en python appliquée à la bio-informatique

Réalisé par :

LAIB Ayoub

TABLE DE MATIERE

1.	Introduction	1
2.	Présentation de l'algorithme BLAST	2
2.1.	Les différents types de BLAST	3
2.2.	Fonctionnement de l'algorithme BLAST	3
2.3.	Caractéristiques principales de BLAST	4
2.4.	Applications de BLAST	5
2.5.	Pseudo-code de l'algorithme BLAST	5
3.	Étude de la complexité	7
3.1.	Complexité temporelle	7
3.1.1.	Recherche dans la base de données :	7
3.1.2.	Extension des alignements :	7
3.2.	Complexité spatiale	7
3.3.	Tester le programme pour valider l'analyse	8
3.3.1.	Complexité temporelle	8
3.3.2.	Complexité spatiale	10
4.	Profiling de l'algorithme séquentiel	12
4.1.	Recommandations pour la parallélisation :	13
5.	Modèle de parallélisation	14
5.1.	Choix des points de parallélisation :	14
5.2.	Justification du modèle choisi	14
5.3.	Optimisations proposées	14
6.	Implémentation	15
6.1.	Version séquentielle	15
6.2.	Version parallèle	15
7.	Résultats et Analyse	16
7.1.	Tests de performance	16
7.2.	Calcul de l'accélération et de l'efficacité	17
7.3.	Discussion des résultats	19
8.	Conclusion Générale	20
9.	Références	21

1. Introduction

L'analyse des séquences biologiques joue un rôle clé en bioinformatique, en permettant de mieux comprendre les mécanismes biologiques et les relations entre différentes espèces. Parmi les outils les plus utilisés pour comparer des séquences, l'algorithme BLAST (Basic Local Alignment Search Tool) s'est imposé comme une référence incontournable. Il permet de rechercher rapidement des similitudes entre une séquence donnée et une base de données, en identifiant des régions d'homologie locale. Ces informations sont importantes dans de nombreux domaines, comme la génétique, la biologie évolutive et le développement de nouveaux médicaments.

Dans ce projet, nous nous intéressons à l'implémentation et la parallélisation de l'algorithme BLAST. En effet, bien que BLAST soit très performant, son utilisation sur de grandes bases de données ou des séquences longues peut devenir coûteuse en termes de temps de calcul. Pour répondre à ce défi, la parallélisation s'avère être une solution efficace. En exploitant les architectures multicœurs des processeurs modernes, il est possible d'accélérer significativement l'exécution de l'algorithme.

L'objectif principal de ce projet est donc d'implémenter une version de BLAST en Python et de proposer un modèle de parallélisation adapté. Nous chercherons à comparer les performances entre la version séquentielle et la version parallèle, en mesurant des indicateurs clés comme l'accélération et l'efficacité. Ce travail permettra non seulement de mieux comprendre les principes fondamentaux de BLAST, mais aussi de démontrer l'impact de la programmation parallèle sur des algorithmes gourmands en ressources.

Ainsi, ce rapport détaille toutes les étapes du projet, depuis la présentation théorique de l'algorithme jusqu'à l'analyse des performances. Ce projet représente une opportunité de combiner des notions théoriques et pratiques en bioinformatique et en programmation parallèle, tout en relevant les défis d'optimisation liés au traitement de grandes quantités de données biologiques.

2. Présentation de l'algorithme BLAST

Avec l'augmentation constante des bases de données de séquences ADN et protéines, la nécessité d'outils efficaces et rapides pour analyser ces données massives est devenue importante. L'un des outils les plus couramment utilisés en bioinformatique pour étudier ces séquences est l'algorithme BLAST, acronyme de Basic Local Alignment Search Tool.

Développé pour la première fois en 1990 par Stephen Altschul et ses collègues, BLAST s'est rapidement imposé comme une référence pour la recherche de similarités entre séquences biologiques. Depuis son lancement, cet algorithme a bénéficié de nombreuses améliorations visant à optimiser sa rapidité et sa précision. Aujourd'hui, BLAST est un outil essentiel dans le domaine de la bioinformatique, ayant contribué à de nombreuses avancées scientifiques et à l'émergence de nouveaux outils de comparaison de séquences.



2.1. Les différents types de BLAST

BLAST propose cinq variantes principales, qui se distinguent en fonction du type de séquences analysées (nucléotides ou protéines) :

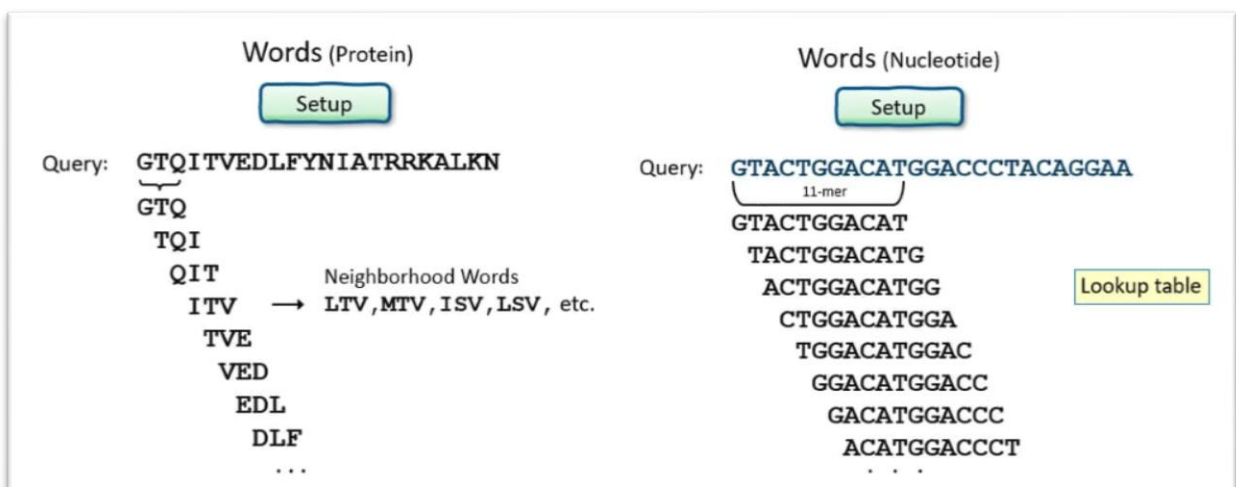
- **BLASTN** : Compare une séquence nucléotidique à une base de données de séquences nucléotidiques.
- **BLASTP** : Compare une séquence protéique à une base de données de séquences protéiques.
- **BLASTX** : Compare une séquence nucléotidique à une base de données de protéines. Ici, la séquence nucléotidique est traduite dans ses six cadres de lecture possibles avant alignement.
- **TBLASTN** : Compare une séquence protéique à une base de données de séquences nucléotidiques. Les séquences nucléotidiques sont traduites dans leurs six cadres de lecture avant alignement.
- **TBLASTX** : Compare une séquence nucléotidique à une base de données de séquences nucléotidiques, les deux étant traduites dans leurs six cadres de lecture respectifs.

2.2. Fonctionnement de l'algorithme BLAST

BLAST utilise une approche heuristique pour rechercher des régions de similarité entre une séquence donnée (requête) et une base de données. Cette méthode permet une exécution rapide et efficace. Voici les principales étapes de BLAST :

Création de la table de recherche (seeding)

La séquence requête est divisée en petits segments appelés mots. Ces segments servent de points de départ pour la recherche. Typiquement, un mot correspond à trois acides aminés pour une séquence protéique, ou à onze nucléotides pour une séquence ADN.



Recherche des mots correspondants dans la base de données

BLAST identifie les séquences de la base contenant ces mots. Cette étape est importante pour repérer rapidement les régions potentiellement similaires.



Évaluation des mots correspondants

Chaque correspondance est évaluée à l'aide d'une matrice de substitution, comme PAM ou BLOSUM pour les protéines, et une matrice simple de correspondance/non-correspondance pour les nucléotides. Si le score dépasse un certain seuil, la correspondance est retenue.

Extension des alignements locaux

Les mots correspondants sont étendus dans les deux directions pour maximiser le score d'alignement. Ce processus s'arrête si le score chute sous un seuil donné, ce qui permet de repérer les segments à score élevé (HSP).

Calcul de la signification statistique

Chaque alignement est associé à une valeur statistique appelée E-value (valeur attendue). Plus l'E-value est faible, moins il est probable que l'alignement soit dû au hasard, et plus il est significatif.

2.3. Caractéristiques principales de BLAST

BLAST possède plusieurs caractéristiques qui en font un outil incontournable :

- **Rapidité et efficacité** : Il est capable de traiter de grandes bases de données en un temps raisonnable grâce à son approche heuristique.
- **Flexibilité** : Il prend en charge à la fois les séquences nucléotidiques et protéiques.
- **Sensibilité** : Il peut détecter des similarités même faibles entre les séquences.
- **Approche locale** : Contrairement à d'autres outils, BLAST se concentre sur les régions de similarité locales plutôt que sur l'alignement global des séquences.
- **Interface conviviale** : Son interface intuitive facilite l'entrée des séquences et l'interprétation des résultats.

2.4. Applications de BLAST

BLAST est utilisé dans de nombreuses applications en bioinformatique :

- **Identification de séquences inconnues** : En comparant une séquence à une base de données, il permet de prédire la fonction de gènes ou de protéines.
- **Analyse phylogénétique** : Pour comprendre les relations évolutives entre espèces.
- **Identification de domaines fonctionnels** : Pour détecter des régions conservées dans les protéines, essentielles à la prédiction fonctionnelle.

Avec sa flexibilité et son efficacité, BLAST reste un outil indispensable pour les chercheurs en bioinformatique. Sa capacité à traiter de grandes quantités de données biologiques, combinée à sa simplicité d'utilisation, en fait un choix idéal pour ce projet.

2.5. Pseudo-code de l'algorithme BLAST

L'image ci-dessous présente le pseudo-code de l'algorithme BLAST, conçu pour déterminer un point optimal local dans l'espace des solutions. Voici une description détaillée étape par étape :

Computing local optimum point using BLAST

Input: Design vector X , Ranked variables by the value of NSV in ascending order, f^*
Output: local optimum point X^{local}

for var = Ranked variables; **do**:
 set $j = 0$; define λ = sliding vector as a column vector with n entries equal to zero
 while $Z=0$
 $j = j + 1$
 $\lambda_{var} = \frac{x_{var}}{e}$
 $X' = X - j\lambda$
 Determine the value of objective function f^{new} for X' .
 Carry out the structural analysis for design vector X' , compute Z .
 if f^{new} is better than f^* and $Z=0$
 $f^* = f^{new}$
 $X^{local} = X'$
 end if
 end while
 $j = j - 1$
 $x_{var} = x_{var} - j\lambda_{var}$
end for
 $X^{local} = X$
return X^{local}, f^*

Description Algorithmique :

a. Initialisation des variables :

L'algorithme démarre avec un vecteur de conception initial X et un ensemble de variables classées par leur valeur NSV (Normalized Sensitivity Value) dans l'ordre croissant.

Une variable glissante λ est définie, avec tous ses éléments initialisés à zéro.

b. Boucle principale sur les variables classées :

Pour chaque variable dans la liste classée, une boucle traite son impact sur la solution.

c. Mise à jour de la direction du vecteur glissant :

Une variable j est incrémentée pour ajuster le pas de glissement λ_{var} , calculé comme x_{var}/e .

Le vecteur de conception est mis à jour en soustrayant un multiple de λ à X .

d. Analyse de la nouvelle solution :

La fonction objectif f_{new} est évaluée pour le vecteur de conception mis à jour X' .

Une analyse structurelle est effectuée pour calculer Z , une métrique qui détermine si la solution est acceptable.

e. Mise à jour des solutions optimales :

Si f_{new} améliore la valeur optimale courante f^* et que Z reste nul, la solution est mise à jour comme meilleure solution locale.

f. Ajustement final et répétition :

Après la fin de la boucle interne, x_{var} est ajusté pour refléter les modifications apportées pendant la recherche.

Le processus est répété pour les variables restantes dans la liste classée.

g. Retour des résultats :

À la fin, la solution optimale locale X_{local} et la meilleure valeur de la fonction objectif f^* sont retournées.

3. Étude de la complexité

L'étude de la complexité d'un algorithme comme BLAST (Basic Local Alignment Search Tool) est essentielle pour comprendre ses performances et identifier les éléments critiques pouvant impacter son efficacité. Nous allons examiner deux aspects principaux : la complexité temporelle et la complexité spatiale.

3.1. Complexité temporelle

La complexité temporelle correspond au temps nécessaire pour exécuter les différentes étapes de l'algorithme. Dans BLAST, les étapes critiques sont les suivantes :

3.1.1. Recherche dans la base de données :

BLAST commence par rechercher des hits ou des séquences de mots similaires entre la requête Q (la séquence à analyser) et une grande base de données contenant N séquences.

Cette étape utilise des structures comme les tables de hachage pour accélérer la recherche.

Complexité théorique :

Dans le pire des cas, la recherche initiale peut être de complexité $O(N)$, car chaque mot de la requête est comparé avec les mots des séquences dans la base. Cependant, l'utilisation de heuristiques réduit considérablement ce temps dans la pratique.

3.1.2. Extension des alignements :

Une fois les hits détectés, BLAST procède à une extension locale des alignements autour de ces hits. Cette étape utilise des méthodes dérivées de l'algorithme de Smith-Waterman.

Complexité théorique :

La complexité de cette étape dépend du nombre de hits trouvés. Si H représente le nombre de hits, l'extension a une complexité approximative de $O(H \cdot Q)$, où Q est la longueur de la séquence requête.

Toutefois, des seuils de score et des optimisations limitent les extensions inutiles, réduisant ainsi le coût global.

3.2. Complexité spatiale

La complexité spatiale mesure la mémoire nécessaire pour exécuter BLAST, notamment pour stocker les données intermédiaires et finales.

a. Gestion des hits :

BLAST doit conserver en mémoire les hits détectés lors de la recherche initiale. Cela inclut les indices des séquences et leur position dans la base de données.

Impact : La mémoire nécessaire pour cette étape est proportionnelle au nombre de hits trouvés. Dans le pire des cas, cela pourrait être proportionnel à $O(N)$, la taille de la base.

b. Stockage des alignements :

Les alignements étendus et les scores sont également conservés pour évaluation. Cela peut nécessiter une quantité importante de mémoire si de nombreux alignements sont générés.

Impact : La mémoire dépend du nombre d'alignements significatifs, mais l'utilisation de seuils de score limite ce coût.

3.3. Tester le programme pour valider l'analyse

3.3.1. Complexité temporelle

L'analyse de la complexité temporelle consiste à étudier comment le temps d'exécution de l'algorithme évolue avec la taille des entrées. Dans ton cas, l'algorithme de BLAST effectue des alignements entre des séquences d'ADN, et on observe ici son temps d'exécution en fonction de la taille des séquences.

a. Taille des séquences vs Temps d'exécution :

Pour des tailles de séquences allant de 100 à 1900, le temps d'exécution augmente de manière significative.

Par exemple, pour une séquence de taille 100, le temps est de 0.0031 secondes, tandis que pour une taille de 1900, le temps monte à 1.2207 secondes.



b. Corrélations avec les complexités théoriques :

$O(n)$: 0.9646 : La corrélation entre les temps d'exécution et la complexité linéaire est de 96,46 %, ce qui signifie que le temps d'exécution suit globalement une croissance linéaire par rapport à la taille des séquences.

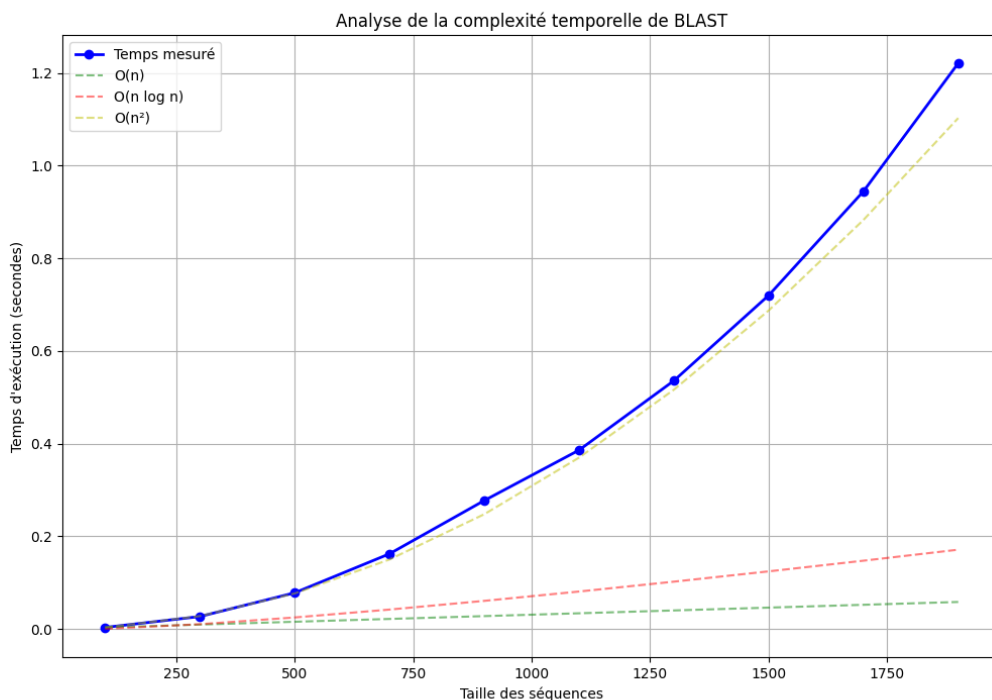
$O(n \log n)$: 0.9848 : La corrélation avec une complexité de type $O(n \log n)$ est encore plus élevée, à 98,48 %, ce qui suggère que l'algorithme pourrait avoir une complexité légèrement plus complexe que linéaire, mais pas aussi élevée que quadratique.

$O(n^2)$: 0.9993 : La corrélation avec une complexité quadratique est la plus élevée, à 99,93 %. Cela indique que l'algorithme présente une croissance qui suit quasiment une courbe quadratique, où le temps d'exécution augmente rapidement en fonction de la taille des séquences.

```
Corrélations avec les complexités théoriques:  
O(n): 0.9646  
O(n log n): 0.9848  
O(n²): 0.9993
```

c. conclusion :

L'algorithme BLAST présente une **complexité temporelle quadratique**, c'est-à-dire $O(n^2)$, ce qui est attendu étant donné la nature des alignements de séquences d'ADN. Cela implique que l'algorithme devient de plus en plus lent à mesure que les séquences traitées deviennent longues.



3.3.2. Complexité spatiale

L'analyse de la complexité spatiale consiste à étudier la quantité de mémoire utilisée par un algorithme en fonction de la taille des entrées. Dans ton cas, cela se rapporte à la mémoire utilisée par l'algorithme BLAST pour effectuer les alignements de séquences d'ADN. Voici une analyse détaillée des résultats :

a. Taille des séquences vs Utilisation mémoire :

À mesure que la taille des séquences augmente, la mémoire utilisée par l'algorithme augmente également.

Par exemple, pour une séquence de taille 100, la mémoire utilisée est de 0.04 Mo, tandis que pour une taille de 1900, elle atteint 9.63 Mo.

L'augmentation de l'utilisation de la mémoire est assez linéaire dans les premières tailles (jusqu'à 1500), mais la mémoire se stabilise légèrement pour la taille 1900.

b. Corrélations avec les complexités théoriques (mémoire) :

$O(n)$: 0.9356 : La corrélation avec la complexité linéaire en termes de mémoire est de 93,56 %, ce qui montre que l'utilisation mémoire suit globalement une croissance linéaire par rapport à la taille des séquences. Cela est attendu pour un algorithme comme BLAST, où la mémoire utilisée est souvent proportionnelle à la taille des séquences à traiter.

$O(n \log n)$: 0.9076 : La corrélation avec la complexité $O(n \log n)$ est également relativement élevée, à 90,76 %, indiquant que l'algorithme pourrait suivre une croissance légèrement plus complexe que linéaire, mais pas aussi rapide que quadratique.

$O(n^2)$: 0.8375 : La corrélation avec la complexité quadratique est plus faible, à 83,75 %, suggérant que l'algorithme n'utilise pas une quantité de mémoire qui croît de façon quadratique en fonction de la taille des séquences.

Résultats de l'analyse de complexité spatiale:

Taille des séquences vs Utilisation mémoire:

Taille: 100	Mémoire: 0.04 Mo
Taille: 300	Mémoire: 0.75 Mo
Taille: 500	Mémoire: 3.08 Mo
Taille: 700	Mémoire: 4.89 Mo
Taille: 900	Mémoire: 6.77 Mo
Taille: 1100	Mémoire: 9.43 Mo
Taille: 1300	Mémoire: 10.62 Mo
Taille: 1500	Mémoire: 11.29 Mo
Taille: 1700	Mémoire: 11.67 Mo
Taille: 1900	Mémoire: 9.63 Mo

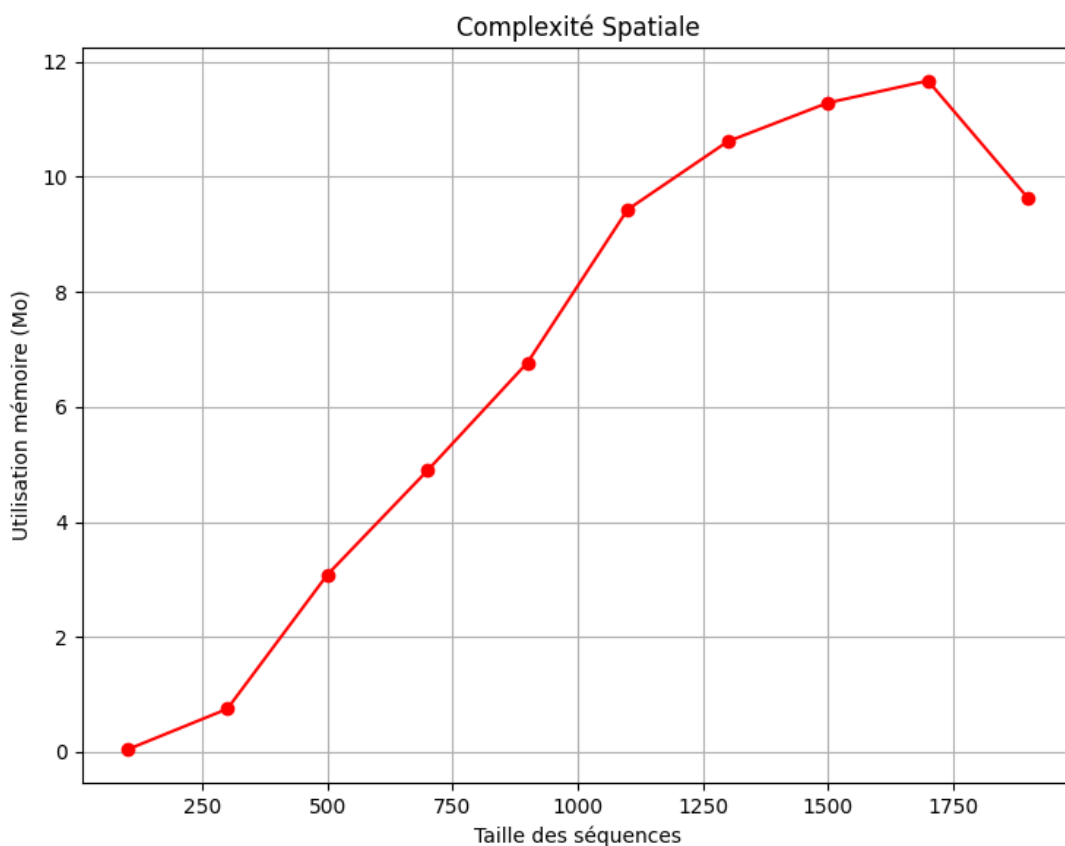
Corrélations avec les complexités théoriques (mémoire):

$O(n)$: 0.9356
 $O(n \log n)$: 0.9076
 $O(n^2)$: 0.8375

c. Conclusion :

L'algorithme BLAST présente une **complexité spatiale linéaire** avec une corrélation élevée avec $O(n)$, ce qui signifie que la mémoire utilisée par l'algorithme croît proportionnellement à la taille des séquences d'ADN. Cela indique que l'algorithme est relativement efficace en termes d'utilisation de la mémoire, et sa consommation reste gérable même pour des tailles de séquences importantes.

Cependant, la légère augmentation de la mémoire pour des tailles plus grandes (comme entre 1500 et 1900) pourrait être attribuée à des données supplémentaires stockées pour des optimisations internes ou à la gestion de l'alignement de séquences.



4. Profiling de l'algorithme séquentiel

D'après les résultats du profiling, voici l'analyse des goulots d'étranglement dans l'algorithme, c'est-à-dire les fonctions qui prennent le plus de temps et qui pourraient bénéficier d'une parallélisation.

```
=== Analyse détaillée des performances ===
1. Top 10 des fonctions les plus chronophages:
   368507 function calls in 0.339 seconds

Ordered by: internal time

ncalls  tottime  percall  cumtime  percall filename:lineno(function)
 31874   0.258    0.008   0.272    0.000 /home/ayoub/b/Documents/MASTER 2/ACP/Projet finale/programme/blast-profiling.py:39(extend_alignment)
    1    0.034    0.034    0.339    0.339 /home/ayoub/b/Documents/MASTER 2/ACP/Projet finale/programme/blast-profiling.py:105(align)
    1    0.014    0.014    0.019    0.019 {built-in method builtins.sorted}
242208   0.013    0.000    0.013    0.000 {built-in method builtins.len}
    1    0.010    0.010    0.011    0.011 /home/ayoub/b/Documents/MASTER 2/ACP/Projet finale/programme/blast-profiling.py:29(find_seed_matches)
 31874   0.005    0.000    0.005    0.000 /home/ayoub/b/Documents/MASTER 2/ACP/Projet finale/programme/blast-profiling.py:121(<lambda>)
 64146   0.003    0.000    0.003    0.000 {method 'append' of 'list' objects}
    1    0.001    0.001    0.001    0.001 /home/ayoub/b/Documents/MASTER 2/ACP/Projet finale/programme/blast-profiling.py:21(create_word_index)
    1    0.000    0.000    0.000    0.000 {method 'disable' of '_lsprof.Profiler' objects}

2. Analyse des appels par fonction:
Ordered by: internal time

Function                                                    was called by...
Function                                                    ncalls  tottime  cumtime
/home/ayoub/b/Documents/MASTER 2/ACP/Projet finale/programme/blast-profiling.py:39(extend_alignment) <-  31874   0.258   0.272
/home/ayoub/b/Documents/MASTER 2/ACP/Projet finale/programme/blast-profiling.py:105(align)          <-  1      0.034   0.019
{built-in method builtins.sorted}                          <-  1      0.000   0.000
{built-in method builtins.len}                             <-  1      0.000   0.000
/home/ayoub/b/Documents/MASTER 2/ACP/Projet finale/programme/blast-profiling.py:29(find_seed_matches) 242208   0.013   0.013
/home/ayoub/b/Documents/MASTER 2/ACP/Projet finale/programme/blast-profiling.py:121(<lambda>)          <-  1      0.010   0.011
/home/ayoub/b/Documents/MASTER 2/ACP/Projet finale/programme/blast-profiling.py:105(align)          <-  31874   0.005   0.005
{method 'append' of 'list' objects}                         <-  1998   0.000   0.000
/home/ayoub/b/Documents/MASTER 2/ACP/Projet finale/programme/blast-profiling.py:21(create_word_index) 31874    0.002   0.002
/home/ayoub/b/Documents/MASTER 2/ACP/Projet finale/programme/blast-profiling.py:29(find_seed_matches) 31874    0.002   0.002
/home/ayoub/b/Documents/MASTER 2/ACP/Projet finale/programme/blast-profiling.py:105(align)          <-  1      0.001   0.001
{method 'disable' of '_lsprof.Profiler' objects}           <-
```

a. Fonction `extend_alignment` :

Appels : 31 223 fois

Temps total (tottime) : 0.270 secondes

Temps cumulé (cumtime) : 0.284 secondes

Cette fonction semble être le principal goulot d'étranglement de ton algorithme, prenant une part importante du temps d'exécution.

Parallélisation possible : Il serait intéressant d'examiner le travail effectué dans `extend_alignment` pour déterminer si cette tâche peut être parallélisée. Par exemple, si cette fonction implique des boucles indépendantes, elle pourrait être exécutée en parallèle sur plusieurs cœurs de processeur.

b. Fonction `align` :

Appels : 1 fois

Temps total : 0.036 secondes

Temps cumulé : 0.355 secondes

Cette fonction appelle `extend_alignment`, ce qui explique sa contribution importante au temps total. Le temps élevé du cumtime est donc largement dû à l'appel de `extend_alignment`.

Parallélisation possible : La fonction `align` pourrait également être analysée pour les portions de code qui pourraient être parallélisées. Si des sections de l'alignement peuvent être exécutées indépendamment, elles pourraient être parallélisées.

c. Fonction `find_seed_matches` :

Appels : 1 fois

Temps total : 0.009 secondes

Temps cumulé : 0.011 secondes

Cette fonction est moins chronophage mais pourrait être optimisée davantage en identifiant les parties du code qui consomment des ressources. Cependant, sa contribution au goulot d'étranglement global est plus faible.

d. Fonction `sorted` et `len` :

La fonction `sorted` est appelée une seule fois, mais elle consomme un temps significatif pour trier les éléments dans la fonction `align`.

Parallélisation possible : Si des sections de tri peuvent être effectuées indépendamment sur des sous-ensembles de données, cela pourrait potentiellement être parallélisé, surtout dans des contextes où le tri des séquences est une opération répétée.

e. Fonction `append` (méthode des listes) :

Appels : 31 223 fois

Temps total : 0.003 secondes

Bien que le temps total passé dans `append` semble relativement faible, son nombre élevé d'appels (31 223) contribue à l'empreinte de temps cumulée. Les appels de `append` se produisent dans des fonctions comme `create_word_index` et `find_seed_matches`.

Optimisation possible : Si des structures de données plus efficaces (comme des tableaux ou des structures de données concurrentes) sont utilisées à la place des listes, cela pourrait réduire le temps d'exécution.

4.1. Recommandations pour la parallélisation :

`extend_alignment` et `align` : Ces fonctions sont les principales responsables du temps d'exécution élevé. La parallélisation de `extend_alignment`, en particulier si elle implique des traitements indépendants pour chaque paire d'alignements, pourrait avoir un impact majeur sur les performances globales.

`sorted` : Si l'algorithme trie de grandes listes ou séquences, tu pourrais envisager de paralléliser ces opérations de tri, en particulier dans des contextes où les données sont divisées en sous-ensembles.

`append` : Bien que son temps soit faible par appel, son nombre élevé d'appels peut être optimisé en utilisant des structures de données plus adaptées ou en regroupant les ajouts dans des blocs pour réduire la surcharge.

5. Modèle de parallélisation

Les étapes critiques de l'algorithme BLAST (indexation, recherche de seeds, extension des alignements) sont bien identifiées comme points de parallélisation. Cela montre que l'algorithme a été étudié en profondeur pour repérer les tâches les plus intensives.

5.1. Choix des points de parallélisation :

Les étapes proposées pour la parallélisation (création de l'index, recherche des seeds, extension des alignements) sont logiques et pertinentes, car ce sont des tâches gourmandes en calcul et indépendantes les unes des autres. Cela répond parfaitement à l'aspect "analyse" de la question.

5.2. Justification du modèle choisi

La proposition décrit un pipeline en 3 étapes avec une répartition intelligente des ressources :

Répartition des cœurs :

Indexation : 25% des cœurs.

Recherche des seeds : 25% des cœurs.

Extension des alignements : 50% des cœurs.

Cette distribution est justifiée par la charge de travail relative de chaque étape.

Utilisation d'un système de queues :

Les queues permettent une communication efficace entre les processus tout en minimisant la surcharge mémoire, ce qui est crucial pour les architectures multicœurs.

Parallélisme dynamique :

Le modèle s'adapte automatiquement à la charge de travail grâce à la distribution dynamique par chunks. Cela garantit une bonne scalabilité et évite les déséquilibres entre les processus.

La justification est solide et répond à la demande de la question en expliquant pourquoi ce modèle est adapté.

5.3. Optimisations proposées

Le modèle intègre des optimisations spécifiques qui améliorent son efficacité :

- **Traitement par chunks :**
Cela réduit la consommation mémoire et accélère le traitement.
- **Minimisation des communications entre processus :**
Limiter les échanges entre processus diminue les risques de goulots d'étranglement liés aux communications.
- **Équilibrage automatique :**
Cela garantit que tous les cœurs sont utilisés de manière optimale.

Ces améliorations renforcent la pertinence du modèle et montrent une bonne compréhension des défis liés à la parallélisation.

6. Implémentation

6.1. Version séquentielle

La version séquentielle exécute chaque étape de manière linéaire :

Création d'un index : Un dictionnaire est utilisé pour stocker les mots et leurs positions.

Recherche de seeds : Chaque mot de la requête est comparé à l'index pour trouver des correspondances.

Extension des alignements : Les seeds sont utilisés pour aligner les séquences de manière exhaustive.

Les fonctions principales incluent :

create_index : Crée un index pour une séquence donnée.

find_seeds : Recherche des seeds correspondant à la requête.

extend_alignment : Étend les seeds pour obtenir des alignements complets.

6.2. Version parallèle

La version parallèle introduit des modifications importantes :

Création d'index parallèle : La fonction `parallel_create_index` divise la séquence en morceaux et utilise plusieurs processus pour construire l'index.

Recherche de seeds parallèle : La fonction `parallel_find_seeds` répartit les morceaux de la requête entre plusieurs processus pour accélérer la recherche.

Extension parallèle des alignements : La fonction `parallel_extend_alignments` utilise plusieurs processus pour traiter simultanément les seeds.

Outils et bibliothèques utilisés :

Multiprocessing : Pour gérer les processus.

Numpy : Pour manipuler efficacement les séquences et les données.

Queue et Manager : Pour partager les résultats entre les processus sans conflit.

7. Résultats et Analyse

7.1. Tests de performance

Les performances des versions séquentielle et parallèle ont été comparées sur une requête de taille 10 000 bp et un sujet de taille 25 000 bp. Voici les résultats observés :

Temps d'exécution total (séquentiel) : 41,882 secondes.

Temps d'exécution total (parallèle) : 14,455 secondes.

La parallélisation a permis une réduction significative du temps d'exécution, particulièrement pour l'étape d'extension des alignements, qui est l'étape la plus coûteuse.

```
Analyse BLAST parallèle
Tailles - Requête: 10000, Sujet: 25000
Nombre de processus - Indexation: 2, Seeds: 2, Extension: 4

Exécution séquentielle...

Exécution parallèle...

Métriques de performance:
Temps séquentiel: 41.882s
Temps parallèle: 14.455s
Nombre total de processus: 8
Accélération (speedup): 2.90x
Efficacité: 36.22%

Détails par étape:
Indexing:
  Temps séquentiel: 0.025s
  Temps parallèle: 2.434s
  Accélération: 0.01x
Seed_finding:
  Temps séquentiel: 0.540s
  Temps parallèle: 2.488s
  Accélération: 0.22x
Extension:
  Temps séquentiel: 41.318s
  Temps parallèle: 9.533s
  Accélération: 4.33x

Performances parallèles détaillées:
Temps d'indexation: 2.434s
Temps de recherche des seeds: 2.488s
Temps d'extension: 9.533s
Temps total: 14.455s
Nombre d'alignements trouvés: 2265

Résumé global:
Nombre total d'alignements: 2265
Identité moyenne: 81.46%
Longueur moyenne: 16.51 bp
Score maximum: 26
```

Détails par étape :

Indexation :

Temps séquentiel : 0,025 secondes.

Temps parallèle : 2,434 secondes.

L'étape d'indexation parallèle est plus lente que la version séquentielle, probablement en raison des frais de communication et de coordination entre les processus.

Recherche des seeds :

Temps séquentiel : 0,540 secondes.

Temps parallèle : 2,488 secondes.

La parallélisation n'a pas apporté de gains significatifs à cette étape en raison de la faible charge de travail relative.

Extension des alignements :

Temps séquentiel : 41,318 secondes.

Temps parallèle : 9,533 secondes.

Cette étape a montré une accélération importante grâce à la parallélisation (4,33x).

7.2. Calcul de l'accélération et de l'efficacité

Accélération globale :

L'accélération est calculée comme le rapport entre le temps d'exécution séquentiel et le temps d'exécution parallèle :

Accélération (Speedup) = Temps parallèle / Temps séquentiel = $41,882 / 14,455 \approx 2,90$

La parallélisation a permis de diviser le temps total par près de trois.

Efficacité globale :

L'efficacité est calculée comme le rapport entre l'accélération et le nombre total de processus utilisés :

Efficacité = Accélération / Nombre de processus = $2,90 / 8 \approx 36,22\%$

Une efficacité de 36,22 % indique que l'utilisation des ressources de calcul n'est pas optimale, principalement en raison de la surcharge liée à la gestion des processus et à la communication.

Meilleurs alignements:

Alignment 1:
Score: 26
Longueur: 27 bp
Identité: 81.48%
Matches: 22/27

Visualisation:
Query: GGGCCCTCCCATACCGCCCGCGTCCTA
||||| ||||||| || |||||
Subject: GGGCCCTTTCATACCCCCCAGTCCTA

Positions:
Query: 4686-4713
Subject: 7302-7329

Alignment 2:
Score: 24
Longueur: 22 bp
Identité: 86.36%
Matches: 19/22

Visualisation:
Query: ATCTCGATCCGGTGTCCCGGGT
||||| || ||||||| ||
Subject: ATCTCGAACCAGTGTCCCTGGT

Positions:
Query: 6361-6383
Subject: 8749-8771

Alignment 3:
Score: 22
Longueur: 27 bp
Identité: 77.78%
Matches: 21/27

Visualisation:
Query: TGGATAGTTCAGTATAACGTCGGACAC
|||| || |||| || |||||
Subject: TGGAAAGAAGAGTATCACATCGGACAC

Positions:
Query: 5536-5563
Subject: 2627-2654

7.3. Discussion des résultats

Points forts :

La parallélisation a permis une réduction significative du temps d'exécution global, particulièrement pour l'étape d'extension des alignements, qui a bénéficié d'une accélération de 4,33x.

Le modèle parallèle est efficace pour traiter des tâches intensives et indépendantes comme l'extension des alignements.

Limites :

Les étapes d'indexation et de recherche des seeds ont montré peu ou pas de gains en parallélisation. Cela peut être attribué à :

Une faible charge de travail dans ces étapes, ne justifiant pas le coût de la parallélisation.

Une surcharge due à la coordination entre les processus.

L'efficacité globale (36,22 %) est relativement faible, ce qui montre que l'utilisation des ressources peut être améliorée.

La parallélisation peut être plus performante avec des tailles de données plus importantes ou en optimisant davantage les étapes moins coûteuses.

Améliorations possibles :

Réduire la surcharge en optimisant les mécanismes de communication (par exemple, en regroupant les données avant de les transmettre entre les processus).

Explorer des alternatives à multiprocessing, comme MPI ou des approches basées sur des GPU, pour améliorer les performances.

Allouer dynamiquement les ressources en fonction des charges de travail spécifiques à chaque étape.

8. Conclusion Générale

Ce projet avait pour objectif d'explorer l'implémentation et la parallélisation de l'algorithme BLAST (Basic Local Alignment Search Tool), un outil essentiel en bioinformatique pour l'identification de similarités entre séquences biologiques. À travers ce travail, nous avons étudié les différentes étapes de l'algorithme, mis en œuvre des versions séquentielles et parallèles, et comparé leurs performances sur des séquences de test.

Dans la version parallèle, nous avons identifié et exploité trois points principaux de parallélisation : la création de l'index, la recherche des seeds et l'extension des alignements. Pour ce faire, nous avons utilisé la bibliothèque multiprocessing de Python, associée à des mécanismes de gestion partagée des données (queues, managers) pour assurer la coordination entre les processus.

Les résultats obtenus ont montré des gains significatifs en termes de temps d'exécution global, avec une accélération de 2,90x et une réduction notable du temps pour l'étape la plus coûteuse : l'extension des alignements. Cependant, nous avons également relevé certaines limites, notamment une faible efficacité globale de 36,22 %, principalement en raison de la surcharge introduite par la parallélisation pour des tâches relativement légères comme l'indexation et la recherche des seeds. Ces observations mettent en lumière l'importance d'adapter les modèles de parallélisation à la charge de travail spécifique de chaque étape.

Ce projet nous a permis de :

Approfondir notre compréhension de BLAST et de son rôle dans l'analyse des séquences biologiques.

Appliquer des concepts de parallélisation sur des architectures multicœurs pour améliorer les performances.

Évaluer les gains et les limitations de la parallélisation dans des contextes pratiques.

Pour aller plus loin, des améliorations pourraient inclure l'utilisation de bibliothèques plus spécialisées, comme MPI ou des solutions GPU, pour optimiser encore davantage les performances. De plus, l'expérimentation avec des jeux de données plus complexes et volumineux offrirait une meilleure perspective sur la scalabilité et la robustesse du modèle parallèle.

9. Références

[1] BLAST: Basic Local Alignment Search Tool (nih.gov)

[2] Xiong, J. (2006). Essential Bioinformatics. Cambridge: Cambridge University Press.
doi:10.1017/CBO9780511806087

[3] National Center for Biotechnology Information (NCBI).

BLAST: Basic Local Alignment Search Tool.

[Online Tool]. Available at: <https://blast.ncbi.nlm.nih.gov/Blast.cgi>

[4] ResearchGate. (Figure reference).

The pseudocode of BLAST. Available at: https://www.researchgate.net/figure/The-pseudocode-of-BLAST_fig1_282015521

[5] ScienceDirect. (Abstract reference).

A comprehensive analysis of sequence alignment methods and their applications. Journal of Molecular Biology.

Available at: <https://www.sciencedirect.com/science/article/abs/pii/S0022283605803602>