

# bioinformatique



# Table des matières

|   |           |
|---|-----------|
| <b>I - chapitre2 Alignement pair</b>                        | <b>3</b>  |
| 1. objectifs spécifiques .....                              | 3         |
| 2. Introduction .....                                       | 3         |
| 2.1. Les similarités de séquences et score .....            | 3         |
| 2.2. La matrice d'identité .....                            | 5         |
| 2.3. Alignement multiple .....                              | 5         |
| 2.4. Exemple d'alignement de séquences par BLAST/NCBI ..... | 6         |
| <b>Glossaire</b>  | <b>8</b>  |
| <b>Abréviations</b>   | <b>9</b>  |
| <b>Références</b>   | <b>10</b> |
| <b>Bibliographie</b>  | <b>11</b> |

# I chapitre2 Alignement pair

## 1. objectifs spécifiques

1- Recherche d'alignement des séquences par Blast/NCBI

2- savoir les étapes des travaux par logiciel MEGA7

## 2. Introduction

Si une nouvelle séquence est obtenue à partir du séquençage génomique, la première étape est la recherche de similarités avec des séquences connues dans d'autres organismes.

Si la fonction/structure des séquences similaires/protéines est connue, très probablement

(highly likely) la nouvelle séquence correspond à une protéine avec la même fonction/structure. En effet, il a été trouvé que seulement à peu près 1% des gènes humains

n'ont pas de contrepartie dans le génome de souris et que la moyenne de similarité entre les gènes de la souris et de l'homme est de 85%.

Les similarités existent parce que toutes les cellules possèdent une cellule ancêtre commune (a mother cell). Donc, dans les différents organismes il pourrait avoir des mutations d'acides aminés dans certaines protéines parce que les acides aminés ne sont pas tous importants pour la fonction et peuvent être remplacés par des acides aminés qui ont des caractéristiques chimiques semblables sans changer la structure. Parfois les mutations sont tellement nombreuses qu'il est difficile de trouver des similarités.

La méthode du calcul des fonctions des gènes par similarités est appelée la génomique comparative ou la recherche d'homologie. Deux séquences sont homologues lorsqu'ils ont comme racine un ancêtre commun.(4)\*

### 2.1. Les similarités de séquences et score

Après le séquençage, les biologistes n'ont habituellement aucune idée de l'utilité des gènes trouvés. En espérant découvrir un indice sur leurs fonctions, ils tentent de trouver des similitudes entre des gènes nouvellement séquencés et d'autres déjà séquencés dont ils connaissent les fonctions.

Le jeu suivant, transformer un mot anglais en un autre mot en passant par une série de mots intermédiaires, dans laquelle chaque mot ne diffère du suivant que d'une seule lettre.

Pour transformer head en tail, on n'a besoin que de quatre intermédiaires :

head → heal → teal → tell → tall → tail.

Pour les séquences biologiques, il est connu comment une séquence peut mutée en une autre. Premièrement, il y'a les points de mutation ou un nucléotide ou acide aminé est changé en un autre. Deuxièmement, il y'a les suppressions ou un élément (nucléotide ou acide aminé) ou une subséquence entière d'un élément est supprimée de la séquence.

Troisièmement, il y'a les insertions ou un élément ou une subséquence est insérée dans la séquence. Un alignement peut s'interpréter comme le fruit d'un travail d'édition : trouver le nombre minimum d'opérations élémentaires d'édition qui permettent de transformer une séquence en une autre. On considère les trois opérations suivantes :

- (a) insertion : insertion d'une ou plusieurs lettres ;
- (b) délétion : suppression d'une ou plusieurs lettres ;
- (c) substitution : remplacement d'une lettre par une autre.

Dans une perspective évolutive ces trois opérations peuvent s'interpréter comme des mutations et le travail d'édition comme une tentative de reconstruction de l'histoire évolutive en considérant ces 3 mutations élémentaires. L'alignement suivant par exemple.

Le conte donne 12 lettres identiques sorties des 14 lettres de BIOINFORMATICS. Les mutations pourraient êtres :

- (1) suppression I BOINFORMATICS
- (2) insertion LI BOILINFORMATICS
- (3) insertion G BOILINGFORMATICS
- (4) changement de T en N BOILINGFORMANICS

Les deux textes semblent très similaires. Noter que l'insertion ou la suppression ne peuvent pas être distinguées si les deux séquences sont présentées (es que le I est supprimé de la première séquence ou inséré dans la seconde ?). Donc, les deux cas sont dénotées par “-”.

La tache des algorithmes bioinformatiques est de trouver à partir de deux séries (la partie à gauche dans l'exemple au-dessus) l'alignement optimal (la partie à droite dans l'exemple au-dessus). L'alignement optimal est l'arrangement des deux séries d'une manière ou le nombre de mutations est minimal.

L'alignement peut être global (sur toute la longueur de la séquence) ou local (sur les parties les mieux conservées), selon la relation présumée entre les séquences. On définit un score d'alignement qui permet de définir le meilleur alignement de deux séquences et de quantifier leur ressemblance.

|                    |   |                   |
|--------------------|---|-------------------|
| BIOINFORMATICS     | → | BIOI-N-FORMATICS  |
| BOILING FOR MANICS |   | B-OILINGFORMANICS |

## 2.2. La matrice d'identité

La matrice d'identité ou matrice de dot (Dot Matrix) est un outil de représentation des alignements, où une séquence est écrite horizontalement en haut et l'autre verticalement à gauche. Ce qui donne une matrice où chaque lettre de la première séquence est couplée avec chaque lettre de la deuxième séquence. Pour chaque correspondance de lettres un point (dot) est inscrit dans la position concordante dans la matrice. Quelles paires apparaissent dans l'alignement optimal ? On va voir ci-après que chaque chemin à travers la matrice correspond à un alignement(6)\*

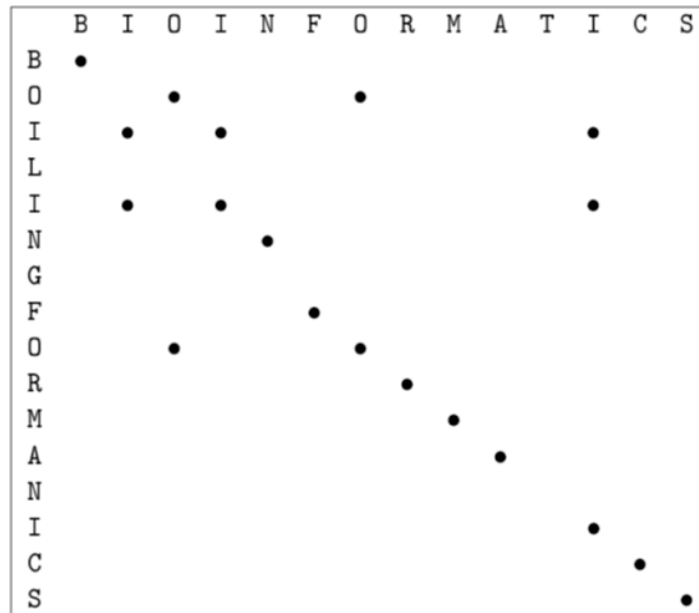


Figure 2a. Principe opérationnel de la matrice d'identité

## 2.3. Alignement multiple

Le but de la comparaison des séquences protéiques est de découvrir des similitudes « biologiques » (i.e. structurales ou fonctionnelles) parmi les protéines. Des protéines biologiquement similaires peuvent ne pas exhiber une forte similitude de séquences et l'on aimerait reconnaître la ressemblance structurale / fonctionnelle, même lorsque les séquences sont très différentes.

La comparaison simultanée de nombreuses séquences permet souvent de trouver des similitudes invisibles dans la comparaison de séquences par paires « l'alignement par paires chuchote... l'alignement multiple crie ».

L'alignement multiple est la base de l'étude de familles de protéines et de domaines fonctionnels. Son but est de révéler des similarités de séquence ou de structure dans une famille de séquences voisines dans l'évolution ou par la fonction.

Il convient de bien analyser le résultat de l'alignement multiple avant de passer à la construction de l'arbre phylogénétique et de bien régler les paramètres du logiciel. Nous allons procéder à l'alignement multiple du jeu de séquences en utilisant l'outil ClustalW.

Ces séquences appartiennent à la famille des facteurs de transcription du type "Basic Leucine Zipper". Ce sont des gènes qui codent pour des protéines qui régulent la transcription des ARNm.

Le résultat d'une partie l'alignement multiple de cette série de séquences est le suivant(5)\*

```

Solanum.tuberosum1466pb      -GGCTGCAC----ACCAAT-CAGCT-----CAGGGTC-----TCC 1172
Triticum.monococcum1062pb    TGACCACAG----GC-AGT-CTGCC-----CGTGCAC-----TTC 931
Rattus.norvegicus1785pb     GGGCAGCCC----ACCAG--CAGCTG-----CAGGAAGCTGATATCC 1427
Zea.mays1236pb              TGGTAGGG----TC--AT-CAGCCC-----CGAGCGCACGGGTGTAC 1047
Oryza.satival272pb         TGGTAG-AA---GCTAG--AGCTT-----AGCTAGC-----1099
Xenopus.laevis1188pb       CGACAGCAACGACTGCTAA--AGTTGC-----CGAAAGC-----1049
Arabidopsis.thaliana1489pb  TAACCAGAA---AAA-GATTCAT-----TGGTTTT-----1281
Triticum.aestivum1585pb     TTGTAGAAGAAGGATCCATCTCTCCCTTCTCTCAGACATAGTCATGCA 1324
                               *
Solanum.tuberosum1466pb     TT-----GOCITAGG-----AGAGT----ACTTTAAACGTC- 1199
Triticum.monococcum1062pb   TT-----GTGATAAG-----TGATT---ACTCATCCCGGC- 958
Rattus.norvegicus1785pb    TTAAACTGAGTCAGGCATCAAGA---CTAAGC---ACTCAGCAAGTG- 1468
Zea.mays1236pb            ATA-----GCTTTCAG-----TAGATCG--AATCCAGGCATG- 1078
Oryza.satival272pb        -----TAGCGAG-----AGAGT--AGCTCAGCTAAGC- 1125
Xenopus.laevis1188pb      -----GCAGCAGA-----GATCCCTAATACTATAAAAAG- 1077
Arabidopsis.thaliana1489pb -----GTGATT---TTGATTG---AGGTAACATATG- 1306
Triticum.aestivum1585pb    TCATGCT-----CCTCGAGAGTCTCTGATATGAGCACATGATCCATGG 1366
                               *
Solanum.tuberosum1466pb     TTGG----TGCTCTTA----GCTCACTTTGGGC-----TGGTCGT 1231
Triticum.monococcum1062pb   TTGG----TGCCCTAA----GTTCTCTTTGG-C-----T--TTGC 987
Rattus.norvegicus1785pb    CTGGA---CTGGTTTGACTCTCGATTGCCCCAAGCCAGCAGAAGTGGTAGT 1515
Zea.mays1236pb            TCCA-----TCAACAAGCAGTTTCTTC-----TCGTGAT 1107
Oryza.satival272pb        TTAATTAGCTGGCTTGAT---TGCTTGGCTTTG-----TGGCTGG 1161
Xenopus.laevis1188pb      TAGG-----GAT---GTCTTTTGATA-----CGTCAC 1102
Arabidopsis.thaliana1489pb TCTG----TATTTTTAT-----TTACTGTATGACTCAGCGCAGGTTAAA 1345
Triticum.aestivum1585pb    TTAATTAACAGGATCTAC----ATCCTCCTG-----TGCTCAT 1400
                               *

```

Figure 3 Cet alignement présente beaucoup de gap qui faussent l'interprétation. Ceci est dû au fait que nos séquences appartiennent à des individus dont la taxonomie est totalement différente. Nous avons aligné des séquences de grenouille, de blé, etc.

## 2.4. Exemple d'alignement de séquences par BLAST/NCBI

La Figure 4 représente le résultat d'un alignement de la séquence partiel du gène ARNr16S

d'Aeromonas veronii obtenue sur GenBank, via le programme BlastN.

>Aeromonas veronii

GenBank, via le programme BlastN réalise un alignement en utilisant ses propres séquences et propose celle qui présente la meilleure identité avec la nôtre en calculant un score qui correspond au nombre de nucléotides identiques chez les deux séquences. Ce score peut être traduit sous forme de pourcentage d'identité (%id). La valeur calculée de E-value indique la probabilité que le résultat de cet alignement a eu lieu par hasard. Donc plus cette valeur est proche du zéro et mieux c'est. Or tous les alignements ont abouti à des valeurs nulles de la E-value ; ce qui exprime que les identités retrouvées entre nos séquences et celles proposées par GenBank ne sont pas dues au hasard.

```

Query 1 TACTTTTCCGGCGAGCGCGGACGGGTGAGTAATGCCTGGGGATCTGCCAGTCGAGGG 60
      |||
Sbjct 61 TACTTTTCCGGCGAGCGCGGACGGGTGAGTAATGCCTGGGGATCTGCCAGTCGAGGG 120

Query 61 GGATAACTACTGGAAACGGTAGCTAATACCGCATACGCCCTACGGGGAAAGCAGGGGAC 120
      |||
Sbjct 121 GGATAACTACTGGAAACGGTAGCTAATACCGCATACGCCCTACGGGGAAAGCAGGGGAC 180

Query 121 CTTGGGGCCTTGGCGATTGGATGAACCCAGGTGGGATTARCTAGTTGGTGAGGTAATGG 180
      |||
Sbjct 181 CTTGGGGCCTTGGCGATTGGATGAACCCAGGTGGGATTAGCTAGTTGGTGAGGTAATGG 240

Query 181 CTCACCAAGGCGACGATCCCTARCTGGTCTGAGAGGATGATCAGCCACTGGAAGTGG 240
      |||
Sbjct 241 CTCACCAAGGCGACGATCCCTARCTGGTCTGAGAGGATGATCAGCCACTGGAAGTGG 300

Query 241 ACACGGTCCAGACTCCTACGGGAGGCGAGCTGGGGAAATATTGCACAATGGGGAAACCC 300
      |||
Sbjct 301 ACACGGTCCAGACTCCTACGGGAGGCGAGCTGGGGAAATATTGCACAATGGGGAAACCC 360

Query 301 TGATGCMCCATGCCGGGTGTGTGAAGAAGCCCTTCGGGTTGTAAAGCACTTCAGCGAG 360
      |||
Sbjct 361 TGATGCMCCATGCCGGGTGTGTGAAGAAGCCCTTCGGGTTGTAAAGCACTTCAGCGAG 420

Query 361 GAGGAAAGGTTGGTAGCTAATAACTGCCAGCTGTGACGTTACTCGCAGAAGCACCAG 420
      |||
Sbjct 421 GAGGAAAGGTTGGTAGCTAATAACTGCCAGCTGTGACGTTACTCGCAGAAGCACCAG 480

Query 421 CTAACCTCCGTGCCAGCAGCCGGTAATACGGAGGTTGCAAGCGTTAATCGGAATTACTG 480
      |||
Sbjct 481 CTAACCTCCGTGCCAGCAGCCGGTAATACGGAGGTTGCAAGCGTTAATCGGAATTACTG 540

Query 481 GCGGTAAAGCGCACGACGGCGGTTGGATAAGTTAGATGTGAAAGCCCCGGGCTCAACTG 540
      |||
Sbjct 541 GCGGTAAAGCGCACGACGGCGGTTGGATAAGTTAGATGTGAAAGCCCCGGGCTCAACTG 600

Query 541 GGAATTGCATTAAAACTGTCCAGCTAGAGTCTTGTAGAGGGGGTAGAATCCAGGTGT 600
      |||
Sbjct 601 GGAATTGCATTAAAACTGTCCAGCTAGAGTCTTGTAGAGGGGGTAGAATCCAGGTGT 660

Query 601 AGCGGTGAAATGCGTAGAGATCTGGAGGAATACCGGTGGCGAAGGCGGCCCC 653
      |||
Sbjct 661 AGCGGTGAAATGCGTAGAGATCTGGAGGAATACCGGTGGCGAAGGCGGCCCC 713

```

Figure 4. Analyse bioinformatique des séquences d'ADNr16s sur GenBank, via le programme BlastN.

# Glossaire

## **bioinformatique**

est la science de l'utilisation de l'ordinateur dans l'acquisition, le traitement et l'analyse de l'information biologique



# Abréviations

**NCBI** : National Center for Biotechnology Information

**PDB** : Protein Data Bank

# Références

- référence* Imbs D., Hassan M.S. (2000) Bioinformatique Travail d'étude. Université de Nice Sophia Antipolis. 23p.
- référence* Tamura K., Stecher G., Peterson D., Filipski A., Kumar S. (2013) MEGA6: Molecular Evolutionary Genetics Analysis version 6.0. Molecular Biology and Evolution: 30 2725-2729.
- référence* - Saitou N., Nei M. (1987). The Neighbor-joining Method: a new method for reconstructing phylogenetic trees. Mol. Biol. Evol. 4 (4): 406-425.
- référence* Darlu P., Tassy P., (2004). La reconstruction phylogénétique. Concepts et méthodes. Masson. Paris. ISBN : 2-225-84229-9. 241 p.
- référence* Corpet F., Chevalet C. (2000) Génétique moléculaire : principes et application aux populations animales. INRA Prod. Anim. Numéro hors série : 191-195.
- référence* Hochreiter S. (2013) Bioinformatics I, Sequence Analysis and Phylogenetics. Institute of Bioinformatics Johannes Kepler University Linz, Austria. 166p.

# Bibliographie

Boubendir A. (2019) Cours de Bioinformatique Centre Universitaire de Mila, Algérie.