

# Construction d'arbres phylogénétiques de virus

Quentin Moreau

2021-2022

La pandémie actuelle, d'après IPBES [1], sera probablement suivie d'autres épidémies. Mieux connaître les virus devient alors un impératif scientifique. Les arbres phylogénétiques permettent justement aux biologistes de mieux connaître les espèces.

**Problème.** Comment construire un arbre phylogénétique optimisé à partir d'un ensemble de séquences ADN ?

## 1 Modélisation et notations

**Définition 1.0.1.** L'ADN est une macromolécule formée elle-même d'une chaîne de molécules appelées nucléotides qui sont l'adénine (A), la cytosine (C), la guanine (G) et la thymine (T). Chaque être vivant est caractérisé par son ADN.

### Notations :

- $\Sigma$  désigne l'alphabet des nucléotides  $\{A, C, G, T\}$  ;
- Le langage de l'ADN, noté  $\mathcal{L}$  est  $\Sigma^*$  ;
- Pour tout mot  $w$  de  $\mathcal{L}$ ,  $\mathcal{P}(w)$ ,  $\mathcal{F}(w)$  et  $\mathcal{S}(w)$  désignent respectivement les ensembles des préfixes, des facteurs et des suffixes de  $w$  ;
- Pour tout mot  $w$  de  $\mathcal{L}$ ,  $|w|$  désigne la longueur de  $w$  ;
- Pour tout mot  $w$  de  $\mathcal{L}$ , pour tout entier  $i$  de  $\llbracket 0, |w| - 1 \rrbracket$ ,  $w[i]$  désigne la  $(i + 1)$ ème lettre de  $w$ .
- Nous identifions les lettres  $A$ ,  $C$ ,  $G$  et  $T$  aux nombres 1, 2, 3 et 4. Nous écrivons par exemple, pour toute matrice  $M$  de  $\mathcal{M}_4(\mathbf{R})$ ,  $M[A, G]$  le coefficient de la première ligne et de la troisième colonne de  $M$ .

## 2 Recherche des régions codantes

### 2.1 Les régions codantes

**Définition 2.1.1.** Certaines régions de l'ADN évoluent plus lentement que d'autres (subissent moins de modifications génétiques d'une génération à l'autre). Elles sont appelées régions codantes, tandis que celles qui évoluent plus rapidement sont appelées régions non codantes par opposition.

*Remarque.* Les régions non codantes évoluent particulièrement rapidement chez les virus et peuvent être considérées comme aléatoires. Il est donc important de ne travailler que sur les régions codantes, plus représentatives des ressemblances entre espèces.

**Problème.** Comment déterminer les régions codantes ?

Soit  $w$ , une séquence ADN de  $\mathcal{L}$ . Les régions codantes de  $w$  sont des facteurs  $f$  de  $\mathcal{F}(w)$  de la forme suivante :

- (1)  $ATG \in \mathcal{P}(f)$  ;
- (2)  $f$  est suivi d'un élément de  $\{TAA, TAG, TGA\}$  ;
- (3)  $\forall c \in \mathcal{F}(f), c \notin \{TAA, TAG, TGA\}$  ;
- (4)  $|f| \in 3\mathbb{Z}$ .

Cependant, les conditions (1), (2), (3) et (4) ne sont pas suffisantes mais seulement nécessaires. Nous ajoutons donc, en accord avec [2] la condition :

- (5)  $|f| \geq 300$ .

*Remarque.* La condition heuristique (5) n'est pas nécessaire. Nous sommes donc susceptibles de trouver des faux positifs et de faux négatifs dans nos tests sur les régions codantes.

**Définition 2.1.2.** Le mot  $ATG$  est nommé codon start et les éléments de  $\{TAA, TAG, TGA\}$  sont nommés codons stop. L'ensemble des codons start et stop est noté  $\mathcal{C}$ .

## 2.2 Algorithmes de recherche de motif

Soit  $w$ , une séquence ADN de  $\mathcal{L}$ , et  $c$ , un codon de  $\mathcal{C}$ .

**Définition 2.2.1.** Une occurrence de  $c$  dans  $w$  est un indice  $i$  de  $\llbracket 0, |w| - |c| \rrbracket$  tel que :  $\forall j \in \llbracket 0, |c| - 1 \rrbracket, w[i + j] = c[j]$ .

**Problème.** Comment trouver toutes les occurrences de  $c$  dans  $w$  le plus efficacement possible ?

**Définition 2.2.2.** Posons  $\mathcal{A} = (\mathcal{P}(c), \{\varepsilon\}, \{c\}, \delta)$  où  $\delta$  est définie comme l'application de  $\mathcal{P}(c) \times \Sigma$  dans  $\mathcal{P}(c)$  qui à tout couple  $(q, a)$  de  $\mathcal{P}(c) \times \Sigma$  associe le plus long mot de  $\mathcal{P}(c) \cap \mathcal{S}(q.a)$ . L'automate  $\mathcal{A}$  est appelé automate des suffixes de  $c$ .

**Proposition 2.2.1.** L'automate des suffixes de  $c$  reconnaît le langage  $\mathcal{L}.c$ .

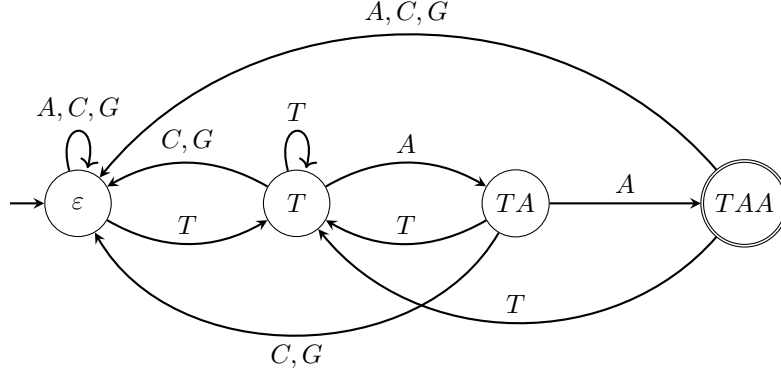


FIGURE 1 – Exemple : automate des suffixes de  $TAA$

---

**Algorithme 1** : Construction de l'automate des suffixes

---

**Entrée** :  $c \in \mathcal{C}$   
 $A.\delta \leftarrow$  table de longueur  $(|c| + 1) \times |\Sigma|$ ;  
*// les préfixes de  $c$  sont identifiés à leurs longueurs, donc les états de  $A$  sont des entiers*  
**Pour**  $i \leftarrow 0$  à  $i \leq |c|$  **faire**  
    **Pour**  $j \leftarrow 0$  à  $j < |\Sigma|$  **faire**  
         $A.\delta[i][j] \leftarrow 0$ ;  
    **Fin pour**  
**Fin pour**  
**Pour**  $i \leftarrow 1$  à  $i \leq |c|$  **faire**  
    ancienne transition  $\leftarrow A.\delta[i - 1][c[i - 1]]$ ;  
     $A.\delta[i - 1][c[i - 1]] \leftarrow i$ ;  
    **Pour**  $j \leftarrow 0$  à  $j < |\Sigma|$  **faire**  
         $A.\delta[i][j] \leftarrow A.\delta[\text{ancienne transition}][j]$ ;  
    **Fin pour**  
**Fin pour**  
**Renvoyer**  $A$

---

**Proposition 2.2.2.** *L'algorithme 1 admet une complexité temporelle en  $O(|c| \times |\Sigma|)$ .*

*Démonstration.* L'algorithme présente, par deux fois, deux boucles pour imbriquées qui réalisent  $|c| + 1$  tours et  $|\Sigma|$  tours.  $\square$

**Proposition 2.2.3.** *L'algorithme 2 admet une complexité temporelle en  $O(|w| + |c| \times |\Sigma|)$  et une complexité spatiale en  $O(|c| \times |\Sigma|)$ .*

---

**Algorithme 2** : Recherche avec l'automate des suffixes

---

**Entrée** :  $w \in \mathcal{L}$ ,  $c \in \mathcal{C}$   
file  $\leftarrow$  file vide ;  
A  $\leftarrow$  automate des suffixes de  $c$  ;  
etat  $\leftarrow$  0;  
**Pour**  $i \leftarrow 0$  **à**  $i < |w|$  **faire**  
    etat  $\leftarrow$  A. $\delta$ [etat][w[i]];  
    **Si** etat =  $|c|$  **alors**  
        file.enfiler (i);  
    **Fin si**  
**Fin pour**  
**Renvoyer** file

---

*Démonstration.* Le calcul de l'automate des suffixes de  $c$  nécessite un  $O(|c| \times |\Sigma|)$  d'opérations élémentaires d'après la proposition 2.2.2. De plus, l'algorithme 2 parcourt ensuite la chaîne  $w$  à travers une boucle pour.  $\square$

Durées d'executions des algorithmes sur 200000 recherches d'un motif de 3 lettres dans une séquence de 10000 lettres

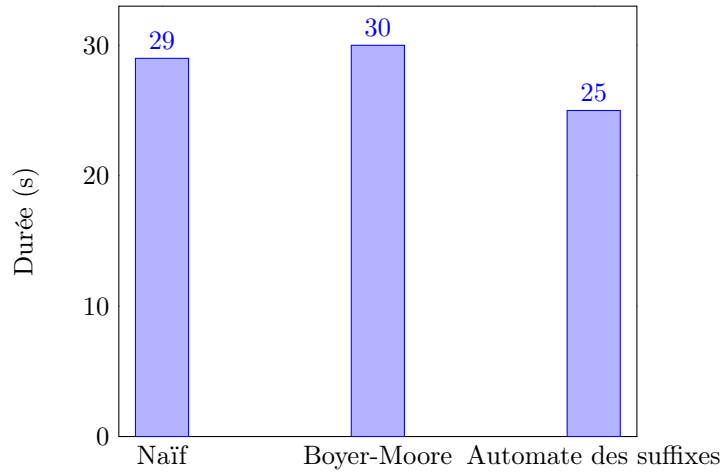


FIGURE 2 – Comparaison des différents algorithmes de recherche

Dans la suite, nous travaillerons exclusivement sur les régions codantes. Ainsi, lorsque l'on parlera de la séquence ADN  $w$ , on entendra la concaténée de ses régions codantes.

### 3 Comparaison des espèces

#### 3.1 Distance de Levenshtein

**Problème.** Comment quantifier la distance entre les espèces ?

**Définition 3.1.1.** La distance de Levenshtein pondérée  $d$  est définie de  $\mathcal{L}^2$  dans  $\mathbf{R}$  récursivement par :

$$\begin{aligned} & \forall (S_1, S_2) \in \mathcal{L}^2, \forall (X, Y) \in \Sigma^2, \\ d(S_1.X, S_2.Y) = & \min(d(S_1, S_2) + M[X, Y], d(S_1.X, S_2) + \alpha, d(S_1, S_2.Y) + \alpha), \\ & \text{et } d(\varepsilon, S_1) = d(S_1, \varepsilon) = \alpha|S_1| \end{aligned}$$

$$\text{où } M = \begin{matrix} & \begin{matrix} A & C & G & T \end{matrix} \\ \begin{matrix} A \\ C \\ G \\ T \end{matrix} & \begin{pmatrix} 0 & \gamma & \beta & \gamma \\ \gamma & 0 & \gamma & \beta \\ \beta & \gamma & 0 & \gamma \\ \gamma & \beta & \gamma & 0 \end{pmatrix} \end{matrix} \text{ avec } \alpha, \beta \text{ et } \gamma \text{ des réels.}$$

*Remarque.* Dans cette modélisation,  $\alpha$  représente le coût d'insertion ou de déletion, et  $\beta$  et  $\gamma$  représentent des coûts de deux types de substitutions appelés respectivement transition et transvection.

#### 3.2 Algorithme de Needleman-Wunsch

Pour calculer la distance de Levenshtein, on utilise un algorithme de programmation dynamique dû à Needleman et Wunsch [3] :

---

**Algorithme 3** : Algorithme de Needleman-Wunsch

---

**Entrée** :  $(w_1, w_2) \in \mathcal{L}^2$ ,  $\alpha \in \mathbf{R}$ ,  $M \in \mathcal{M}_{|\Sigma|}(\mathbf{R})$   
distances  $\leftarrow$  table de longueur  $(|w_1| + 1) \times (|w_2| + 1)$  ;  
**Pour**  $j \leftarrow 0$  à  $j \leq |w_2|$  **faire**  
    **Pour**  $i \leftarrow 0$  à  $i \leq |w_1|$  **faire**  
        **Si**  $i = 0$  **alors**  
            distances[ $i$ ][ $j$ ]  $\leftarrow j \times \alpha$ ;  
        **Fin si**  
        **Sinon si**  $j = 0$  **alors**  
            distances[ $i$ ][ $j$ ]  $\leftarrow i \times \alpha$ ;  
        **Fin si**  
        **Sinon**  
            distances[ $i$ ][ $j$ ]  $\leftarrow \min(\text{distances}[i-1][j-1] + M[w_1[i-1]][w_2[j-1]], \text{distances}[i][j-1] + \alpha, \text{distances}[i-1][j] + \alpha)$ ;  
        **Fin**  
    **Fin pour**  
**Fin pour**  
**Renvoyer** distances[ $w_1$ ][ $w_2$ ]

---

**Proposition 3.2.1.** *L'algorithme 3 admet une complexité spatiale et temporelle en  $O(|w_1| \times |w_2|)$ .*

*Démonstration.* L'algorithme 3 introduit un tableau bidimensionnel de taille  $|w_1| \times |w_2|$  et admet deux boucles pour imbriquées réalisant  $|w_2|$  et  $|w_1|$  tours.  $\square$

*Remarque.* L'algorithme a été linéarisé dans le programme réel en n'enregistrant que deux colonnes.

## 4 Création de l'arbre phylogénétique

L'algorithme UPGMA permet enfin de construire l'arbre phylogénétique à proprement parler :

---

**Algorithme 4 : UPGMA**

---

```
Entrée : distances, arbres phylogénétiques, // arbres phylogénétiques est un
tableau d'arbres phylogénétiques (chacun est représenté par un couple
d'une chaîne de caractère et d'un pointeur) déjà créés initialisé par
les noms des différentes espèces
n // n désigne le nombre d'espèces à traiter

Si  $n \leq 1$  alors
| Renvoyer arbres phylogénétiques[0]
Fin si

// recherche du minimum des coefficients non diagonaux
minimum  $\leftarrow$  distances[1][0];
 $i_1 \leftarrow 1$ ;
 $i_2 \leftarrow 0$ ;
Pour  $i \leftarrow 1$  à  $i < n$  faire
| Pour  $j \leftarrow 0$  à  $j < i$  faire
| | Si distances[i][j] < minimum alors
| | |  $i_1 \leftarrow i$ ;
| | |  $i_2 \leftarrow j$ ;
| | | minimum  $\leftarrow$  distances[i][j];
| | Fin si
| Fin pour
Fin pour

// calcul des nouvelles distances
Pour  $i \leftarrow 0$  à  $i < n$  faire
| distances[i][ $i_1$ ]  $\leftarrow$  (distances[i][ $i_2$ ] + distances[i][ $i_1$ ])/2;
| distances[ $i_1$ ][i]  $\leftarrow$  distances[i][ $i_1$ ];
| distances[i][ $i_2$ ]  $\leftarrow$  distances[i][ $i_1$ ];
| distances[ $i_2$ ][i]  $\leftarrow$  distances[i][ $i_1$ ];
Fin pour

// échange des lignes  $i_1$  et  $n-1$ 
Pour  $i \leftarrow 0$  à  $i < n$  faire
| distances[ $i_1$ ][i]  $\leftarrow$  distances[n - 1][i];
| distances[i][ $i_1$ ]  $\leftarrow$  distances[i][n - 1];
Fin pour

nouvel arbre phylogénétique  $\leftarrow$  ("ancêtre
commun", [arbres phylogénétiques[ $i_1$ ], arbres phylogénétiques[ $i_2$ ]]);
arbres phylogénétiques[ $i_2$ ]  $\leftarrow$  nouvel arbre phylogénétique;
arbres phylogénétiques[ $i_1$ ]  $\leftarrow$  arbres phylogénétiques[n - 1];

Renvoyer UPGMA(distances, arbres phylogénétiques, n - 1)
```

---

## 5 Arbres créés

### 5.1 Arbre de référence

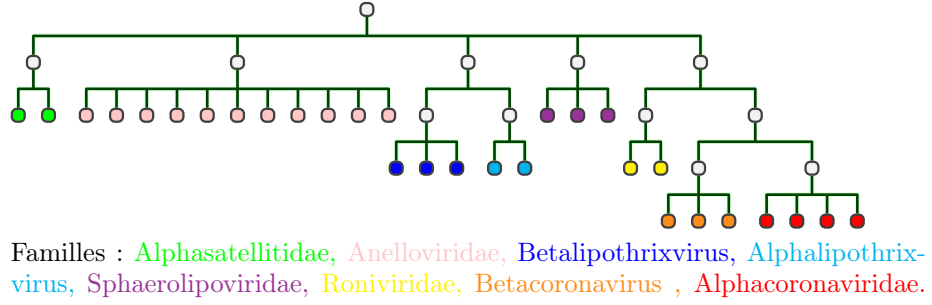


FIGURE 3 – Arbre phylogénétique du NCBI [4]

### 5.2 Arbres obtenus à partir des algorithmes

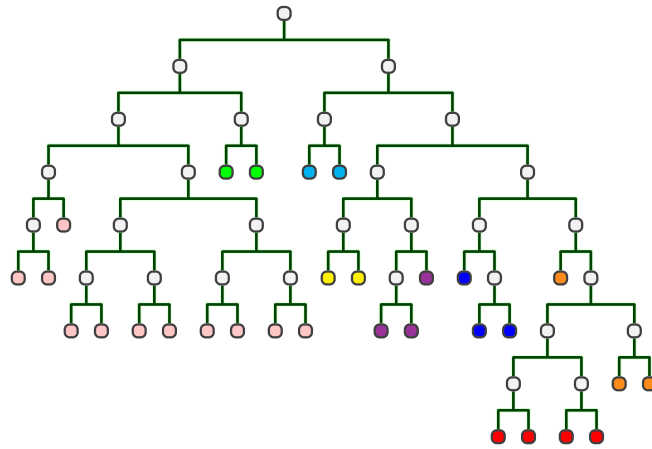


FIGURE 4 – Arbre pour  $\alpha = \beta = \gamma = 1$



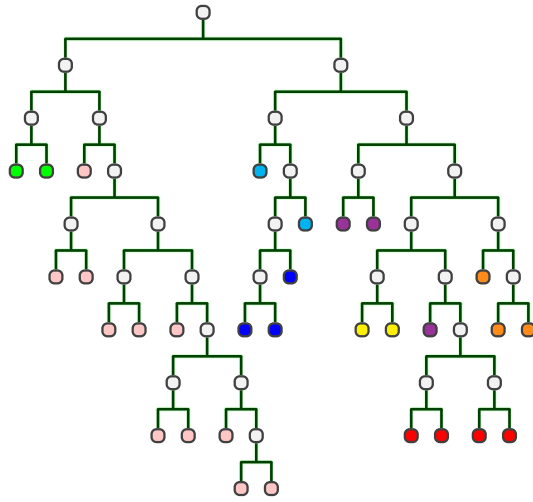


FIGURE 5 – Arbre pour  $\alpha \in \{10, 11, 100\}$ ,  $\beta = 1$  et  $\gamma = 3$

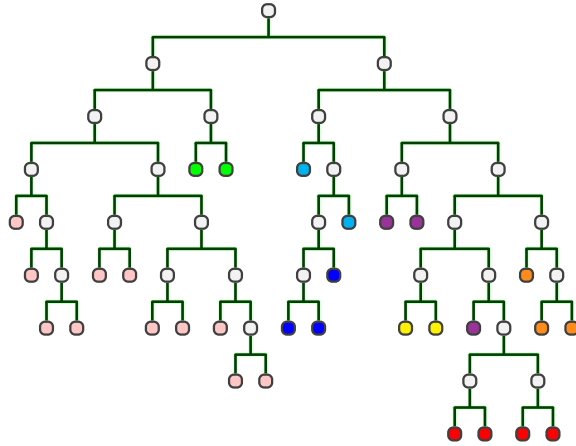


FIGURE 6 – Arbre pour  $\alpha = 10$ ,  $\beta = 1$  et  $\gamma = 2$

## 6 Annexe

```

1 /* Quentin MOREAU
2 * TIPE Bioinformatique
3 * Reconstruction d'arbres phylogenetiques
4 *
5 * fichier d'entete main
6 *
7 * 2021–2022
8 */
9

```

```

10
11
12 #ifndef MAIN_H_INCLUDED
13 #define MAIN_H_INCLUDED
14
15 #include "construction_arbre_phylogenetiques.h"
16
17
18
19 typedef enum Nucleotide Nucleotide;
20 enum Nucleotide
21 {
22     A, C, G, T
23 };
24
25
26
27 typedef struct Element Element;
28 struct Element
29 {
30     int element;
31     Element* precedent;
32     Element* suivant;
33 };
34
35
36 typedef struct File File;
37 struct File
38 {
39     Element* premierElement;
40     Element* dernierElement;
41 };
42
43
44
45 void enfiler(File* file, int element);
46 int defiler(File* file);
47 void supprimer(File* file);
48 Arbre* arbrePhylogenetique(char* nomsFichiers[], char* noms[], int
    nombreSequence);
49
50
51
52 #endif // MAIN_H_INCLUDED

```

main.h

```

1 /* Quentin MOREAU
2  * TIPE Bioinformatique
3  * Reconstruction d'arbres phylogenetiques
4  *
5  * fichier source main
6  *
7  * 2021–2022
8  */
9
10

```

```

11
12 #include <stdio.h>
13 #include <stdlib.h>
14 #include <time.h>
15
16 #include "main.h"
17 #include "fichiers.h"
18 #include "recherche_motif.h"
19 #include "prediction_regions_codantes.h"
20 #include "alignement.h"
21 #include "construction_arbre_phylogenetique.h"
22
23
24
25
26 int main()
27 {
28
29
30     char* virusFichiers[30] = { "sequences/Coronaviridae/Middle
    East respiratory syndrome-related coronavirus.fasta",
31     "sequences/Coronaviridae/Human coronavirus 229E.fasta",
32     "sequences/Coronaviridae/Bat coronavirus HKU10.fasta",
33     "sequences/Coronaviridae/Bat coronavirus CDPHE15.fasta",
34     "sequences/Coronaviridae/Human coronavirus NL63.fasta",
35     "sequences/Coronaviridae/Human coronavirus HKU1.fasta",
36     "sequences/Coronaviridae/Severe acute respiratory syndrome-
    related coronavirus.fasta",
37     "sequences/Roniviridae/Gill-associated virus.fasta",
38     "sequences/Roniviridae/Yellow head virus.fasta",
39     "sequences/Sphaerolipoviridae/Haloarcula hispanica
    icosahedral virus 2.fasta",
40     "sequences/Sphaerolipoviridae/Haloarcula hispanica virus
    PH1.fasta",
41     "sequences/Sphaerolipoviridae/Haloarcula virus HCIV1.fasta"
    ,
42     "sequences/Adnaviria/Acidianus filamentous virus 3.fasta",
43     "sequences/Adnaviria/Acidianus filamentous virus 6.fasta",
44     "sequences/Adnaviria/Acidianus filamentous virus 8.fasta",
45     "sequences/Adnaviria/Sulfolobales Beppu filamentous virus
    2.fasta",
46     "sequences/Adnaviria/Sulfolobus filamentous virus 1.fasta",
47     "sequences/Alphasatellitidae/Ageratum yellow vein Singapore
    alphasatellite.fasta",
48     "sequences/Alphasatellitidae/Coconut foliar decay
    alphasatellite.fasta",
49     "sequences/Anelloviridae/Simian torque teno virus 30.fasta"
    ,
50     "sequences/Anelloviridae/Simian torque teno virus 31.fasta"
    ,
51     "sequences/Anelloviridae/Simian torque teno virus 32.fasta"
    ,
52     "sequences/Anelloviridae/Simian torque teno virus 34.fasta"
    ,
53     "sequences/Anelloviridae/Torque teno virus 10.fasta",
54     "sequences/Anelloviridae/Torque teno virus 11.fasta",
55     "sequences/Anelloviridae/Torque teno virus 12.fasta",

```

```

56     "sequences/Anelloviridae/Torque teno virus 13.fasta",
57     "sequences/Anelloviridae/Torque teno virus 14.fasta",
58     "sequences/Anelloviridae/Torque teno virus 15.fasta",
59     "sequences/Anelloviridae/Torque teno virus 16.fasta"
60 };
61
62 char* virusNoms[30] = {"Middle East respiratory syndrome-
related coronavirus",
63     "Human coronavirus 229E",
64     "Bat coronavirus HKU10",
65     "Bat coronavirus CDPHE15",
66     "Human coronavirus NL63",
67     "Human coronavirus HKU1",
68     "Severe acute respiratory syndrome-related coronavirus",
69     "Gill-associated virus",
70     "Yellow head virus",
71     "Haloarcula hispanica icosahedral virus 2",
72     "Haloarcula hispanica virus PH1",
73     "Haloarcula virus HCIV1",
74     "Acidianus filamentous virus 3",
75     "Acidianus filamentous virus 6",
76     "Acidianus filamentous virus 8",
77     "Sulfolobales Beppu filamentous virus 2",
78     "Sulfolobus filamentous virus 1",
79     "Ageratum yellow vein Singapore alphasatellite",
80     "Coconut foliar decay alphasatellite",
81     "Simian torque teno virus 30",
82     "Simian torque teno virus 31",
83     "Simian torque teno virus 32",
84     "Simian torque teno virus 34",
85     "Torque teno virus 10",
86     "Torque teno virus 11",
87     "Torque teno virus 12",
88     "Torque teno virus 13",
89     "Torque teno virus 14",
90     "Torque teno virus 15",
91     "Torque teno virus 16"
92 };
93
94 arbrePhylogenetique(virusFichiers, virusNoms, 30);
95
96
97
98
99
100
101
102     return 0;
103 }
104
105
106
107
108
109
110
111

```

```

112 Arbre* arbrePhylogenetique(char* nomsFichiers[], char* noms[], int
    nombreSequence)
113 {
114     Nucleotide** sequencesLues = malloc(sizeof(Nucleotide*) *
    nombreSequence);
115     Nucleotide** sequences = malloc(sizeof(Nucleotide*) *
    nombreSequence);
116     int* tailles = malloc(sizeof(int) * nombreSequence);
117
118     if (sequencesLues == NULL || sequences == NULL || tailles ==
    NULL)
119     {
120         printf("echec (arbrePhylogenetique)");
121         exit(EXIT_FAILURE);
122     }
123
124     Nucleotide codonStop0[] = { T, A, A };
125     Nucleotide codonStop1[] = { T, A, G };
126     Nucleotide codonStop2[] = { T, G, A };
127     Nucleotide codonStart[] = { A, T, G };
128
129
130     AutomateSuffixes* automates[4] = { constructionAutomateSuffixes
    (codonStop0, 3),
131                                     constructionAutomateSuffixes(
    codonStop1, 3),
132                                     constructionAutomateSuffixes(
    codonStop2, 3),
133                                     constructionAutomateSuffixes(
    codonStart, 3)
134     };
135
136
137     for (int i = 0; i < nombreSequence; i++)
138     {
139         sequencesLues[i] = lireSequence(nomsFichiers[i]);
140
141         Nucleotide* sequenceCodanteLue = extraireSequencesCodantes(
    sequencesLues[i] + 1, sequencesLues[i][0], automates);
142         sequences[i] = sequenceCodanteLue + 1;
143         tailles[i] = sequenceCodanteLue[0];
144
145         if (tailles[i] == 0)
146         {
147             printf("pas de region codante trouvee");
148             exit(EXIT_FAILURE);
149         }
150     }
151
152     int** distances = distancesLevenshtein(sequences, tailles,
    nombreSequence);
153
154     Arbre* arbre = constructionArbreUPGMA(distances, noms,
    nombreSequence);
155
156     afficherArbre(arbre);
157

```

```

158     return NULL;
159 }
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176  ////////////////////////////////////// FILES
177  //////////////////////////////////////
178
179
180
181
182 void enfiler(File* file , int element)
183 {
184
185     Element* dernierElement = file->dernierElement;
186     Element* nouvelElement = malloc(sizeof(Element));
187
188     if (nouvelElement == NULL)
189     {
190         printf("echec (enfiler)");
191         exit(EXIT_FAILURE);
192     }
193
194     nouvelElement->element = element;
195     nouvelElement->precedent = NULL;
196     nouvelElement->suivant = NULL;
197
198     if (dernierElement != NULL)
199     {
200         dernierElement->suivant = nouvelElement;
201         nouvelElement->precedent = dernierElement;
202     }
203     else
204     {
205         file->premierElement = nouvelElement;
206     }
207
208     file->dernierElement = nouvelElement;
209 }
210
211
212 int defiler(File* file)
213 {

```

```

214     Element* premierElement = file->premierElement;
215     if (premierElement == NULL)
216     {
217         return -1;
218     }
219     int element = premierElement->element;
220     file->premierElement = premierElement->suivant;
221     if (premierElement->suivant != NULL)
222     {
223         premierElement->suivant->precedent = NULL;
224     }
225     free(premierElement);
226     return element;
227 }
228
229
230
231 void supprimer(File* file)
232 {
233     while (file->premierElement != NULL)
234     {
235         defiler(file);
236     }
237 }

```

main.c

```

1  /* Quentin MOREAU
2  * TIPE Bioinformatique
3  * Reconstruction d'arbres phylogenetiques
4  *
5  * fichier d'entete fichiers
6  *
7  * 2021-2022
8  */
9
10
11
12 #ifndef FICHIERS_H_INCLUDED
13 #define FICHIERS_H_INCLUDED
14
15
16
17 Nucleotide* lireSequence(char* nom);
18 void afficherSequence(Nucleotide* sequence, int taille);
19
20
21
22
23
24
25 #endif // FICHIERS_H_INCLUDED

```

fichiers.h

```

1  /* Quentin MOREAU
2  * TIPE Bioinformatique

```

```

3  * Reconstruction d'arbres phylogenetiques
4  *
5  * fichier source fichier
6  *
7  * 2021-2022
8  */
9
10
11 #include <stdlib.h>
12 #include <stdio.h>
13 #include <string.h>
14
15 #include "main.h"
16 #include "fichiers.h"
17
18
19
20 Nucleotide* lireSequence(char* nom)
21 {
22     FILE* fichier = NULL;
23     fopen_s(&fichier, nom, "r");
24
25     if (fichier == NULL)
26     {
27         printf("impossible d'ouvrir %s", nom);
28         exit(EXIT_FAILURE);
29     }
30
31     fseek(fichier, 0, SEEK_END);
32     int tailleFichier = ftell(fichier);
33     rewind(fichier);
34
35     char test = 'a';
36     int curseur = 0;
37     while (test != '\n')
38     {
39         curseur++;
40         test = fgetc(fichier);
41     }
42
43     int lettre = 0;
44
45     int taille = (tailleFichier - 1 - curseur) - (tailleFichier - 1 -
46     curseur) / 61; // 60 caract par ligne + 1 \n
47     Nucleotide* sequence = malloc(sizeof(Nucleotide) * (taille + 1)
48     ); // on enregistre aussi la taille
49
50     if (sequence == NULL)
51     {
52         printf("echec (lireSequence 1)");
53         exit(EXIT_FAILURE);
54     }
55
56     sequence[0] = taille;
57     int i = 1;

```



```

58
59 while (lettre != EOF)
60 {
61     lettre = getc(fichier);
62     if (lettre != '\n' && lettre != EOF)
63     {
64         switch (lettre)
65         {
66             case 'A':
67                 sequence[i] = A;
68                 break;
69             case 'C':
70                 sequence[i] = C;
71                 break;
72             case 'G':
73                 sequence[i] = G;
74                 break;
75             case 'T':
76                 sequence[i] = T;
77                 break;
78             default:
79                 printf("%c", lettre);
80                 printf("\n\nERROR\n\n");
81             }
82             i += 1;
83         }
84     }
85
86     fclose(fichier);
87
88
89
90     return sequence;
91 }
92
93
94
95 void afficherSequence(Nucleotide* sequence, int taille)
96 {
97     for (int i = 0; i < taille; i++)
98     {
99         switch (sequence[i])
100         {
101             case A:
102                 printf("A");
103                 break;
104             case C:
105                 printf("C");
106                 break;
107             case G:
108                 printf("G");
109                 break;
110             case T:
111                 printf("T");
112                 break;
113             default:
114                 printf(" erreur ");

```

```

115         break;
116     }
117 }
118 }

```

#### fichiers.c

```

1  /* Quentin MOREAU
2  * TIPE Bioinformatique
3  * Reconstruction d'arbres phylogenetiques
4  *
5  * fichier entete recherche_motif
6  *
7  * 2021–2022
8  */
9
10
11
12 #ifndef RECHERCHE_MOTIF_H_INCLUDED
13 #define RECHERCHE_MOTIF_H_INCLUDED
14
15 #include "main.h"
16
17
18
19
20
21
22 typedef struct AutomateSuffixes AutomateSuffixes;
23 struct AutomateSuffixes
24 {
25     int etatFinal;
26     int** transitions;
27 };
28
29
30
31
32
33 File* rechercheNaive(Nucleotide sequence[], Nucleotide motif[], int
    tailleSequence, int tailleMotif);
34 File* boyerMoore(Nucleotide sequence[], Nucleotide motif[], int
    tailleSequence, int tailleMotif);
35 AutomateSuffixes* constructionAutomateSuffixes(Nucleotide motif[],
    int taille);
36 File* rechercheAutomateSuffixes(AutomateSuffixes* automate,
    Nucleotide sequence[], int tailleSequence);
37
38
39
40 #endif // RECHERCHE_MOTIF_H_INCLUDED

```

#### recherche\_motif.h

```

1  /* Quentin MOREAU
2  * TIPE Bioinformatique
3  * Reconstruction d'arbres phylogenetiques

```

```

4  *
5  * fichier source recherche_motif
6  *
7  * 2021-2022
8  */
9
10
11
12 #include <stdlib.h>
13 #include <stdio.h>
14
15 #include "recherche_motif.h"
16
17
18
19
20
21
22
23
24
25 //////////////////////////////////////// RECHERCHE NAIVE
26 ////////////////////////////////////////
27
28
29
30
31 File* rechercheNaive(Nucleotide sequence[], Nucleotide motif[], int
    tailleSequence, int tailleMotif)
32 {
33     File* indices = malloc(sizeof(File));
34
35     if (indices == NULL)
36     {
37         exit(EXIT_FAILURE);
38     }
39
40     indices->dernierElement = NULL;
41     indices->premierElement = NULL;
42
43     for (int i = 0; i < tailleSequence - tailleMotif; i++)
44     {
45         int indiceMotif = 0;
46         int test = 1;
47
48         while (test)
49         {
50             if (indiceMotif < tailleMotif)
51             {
52                 test = sequence[i + indiceMotif] == motif[
indiceMotif];
53                 indiceMotif++;
54             }
55             else
56             {
57                 test = 0;

```

```

58         }
59     }
60
61     if (indiceMotif >= tailleMotif)
62     {
63         enfiler(indices, i);
64     }
65 }
66
67
68 return indices;
69 }
70
71
72
73
74
75
76
77
78
79 //////////////////////////////////////////////////// BOYER MOORE
80 ////////////////////////////////////////////////////
81
82
83
84
85
86 File* boyerMoore(Nucleotide sequence[], Nucleotide motif[], int
87     tailleSequence, int tailleMotif)
88 {
89     File* indices = malloc(sizeof(File));
90
91     if (indices == NULL)
92     {
93         exit(EXIT_FAILURE);
94     }
95
96     indices->dernierElement = NULL;
97     indices->premierElement = NULL;
98
99     int table[4] = { -1, -1, -1, -1 };
100
101     for (int i = 0; i < tailleMotif; i++)
102     {
103         table[motif[i]] = i;
104     }
105
106     for (int i = 0; i < tailleSequence + 1 - tailleMotif; i++)
107     {
108         int indiceMotif = tailleMotif - 1;
109         int test = 1;
110
111         while (test)
112         {

```

```

113         if (indiceMotif >= 0)
114         {
115             test = sequence[i + indiceMotif] == motif[
116                 indiceMotif];
117             indiceMotif--;
118         }
119         else
120         {
121             test = 0;
122         }
123     }
124     if (indiceMotif < 0) // le motif est apparu
125     {
126         enfiler(indices, i);
127         if (i + tailleMotif < tailleSequence)
128         {
129             i += tailleMotif - table[sequence[i + tailleMotif]]
130             - 1;
131         }
132     }
133     else // le motif n'est pas apparu
134     {
135         i += max(0, indiceMotif - table[sequence[i + indiceMotif
136             ]] - 1);
137     }
138 }
139
140
141
142 return indices;
143 }
144
145
146
147
148
149
150
151
152
153
154
155
156 //////////////////////////////////////// AUTOMATE DES SUFFIXES
157 ////////////////////////////////////////
158
159
160
161
162
163
164
165 AutomateSuffixes* constructionAutomateSuffixes(Nucleotide motif[],

```

```

166     int taille)
167 {
168     AutomateSuffixes* automate = malloc(sizeof(AutomateSuffixes));
169     int** transitions = malloc(sizeof(int*) * (taille+1));
170
171     if (automate == NULL || transitions == NULL)
172     {
173         exit(EXIT_FAILURE);
174     }
175
176     automate->etatFinal = taille;
177     automate->transitions = transitions;
178
179     for (int i = 0; i <= taille; i++)
180     {
181         transitions[i] = malloc(sizeof(int) * 4);
182         if (transitions[i] == NULL)
183         {
184             exit(EXIT_FAILURE);
185         }
186     }
187
188
189     for (int i = 0; i <= taille; i++)
190     {
191         for (int j = 0; j < 4; j++)
192         {
193             transitions[i][j] = 0;
194         }
195     }
196
197
198
199     for (int i = 1; i <= taille; i++)
200     {
201
202         int ancienneTransition = transitions[i - 1][motif[i - 1]];
203         // enregistrement de l'ancienne transition
204
205         transitions[i - 1][motif[i - 1]] = i;
206
207         for (int j = 0; j < 4; j++)
208         {
209             transitions[i][j] = transitions[ancienneTransition][j];
210         }
211     }
212
213
214
215
216     return automate;
217 }
218
219
220

```

```

221
222
223
224
225
226
227 File* rechercheAutomateSuffixes(AutomateSuffixes* automate,
    Nucleotide sequence[], int tailleSequence)
228 {
229
230     File* indices = malloc(sizeof(File));
231
232     if (indices == NULL)
233     {
234         exit(EXIT_FAILURE);
235     }
236
237     indices->dernierElement = NULL;
238     indices->premierElement = NULL;
239
240
241
242
243
244     int etat = 0;
245
246     for (int i = 0; i < tailleSequence; i++)
247     {
248         etat = automate->transitions[etat][sequence[i]];
249
250         if (etat == automate->etatFinal)
251         {
252             enfiler(indices, i + 1 - automate->etatFinal); //
automate->etatFinal = taille du motif
253         }
254     }
255
256     return indices;
257 }

```

recherche\_motif.c

```

1  /* Quentin MOREAU
2  * TIPE Bioinformatique
3  * Reconstruction d'arbres phylogenetiques
4  *
5  * fichier d'entete prediction_regions_codantes
6  *
7  * 2021-2022
8  */
9
10
11
12 #ifndef PREDICTION_REGIONS_CODANTES_H_INCLUDED
13 #define PREDICTION_REGIONS_CODANTES_H_INCLUDED
14
15 #define TAILLE 300

```

```

16
17 #define STOP0 0
18 #define STOP1 1
19 #define STOP2 2
20 #define START 3
21
22
23 File* rechercheCodonsStart(Nucleotide sequence[], int taille,
    AutomateSuffixes** automates);
24 File* rechercheCodonsStop(Nucleotide sequence[], int taille,
    AutomateSuffixes** automates);
25 File* predictionORF(Nucleotide sequence[], int taille,
    AutomateSuffixes** automates);
26 Nucleotide* extraireSequencesCodantes(Nucleotide sequence[], int
    taille, AutomateSuffixes** automates);
27
28
29 #endif // PREDICTION_REGIONS_CODANTES_H_INCLUDED

```

---

prediction\_regions\_codantes.h

---

```

1  /* Quentin MOREAU
2  * TIPE Bioinformatique
3  * Reconstruction d'arbres phylogenetiques
4  *
5  * fichier source prediction_regions_codantes
6  *
7  * 2021–2022
8  */
9
10
11
12 #include <stdio.h>
13 #include <stdlib.h>
14
15 #include "main.h"
16 #include "recherche_motif.h"
17 #include "prediction_regions_codantes.h"
18
19
20
21
22 File* rechercheCodonsStart(Nucleotide sequence[], int taille,
    AutomateSuffixes** automates)
23 {
24
25     if (automates != NULL)
26     {
27         return rechercheAutomateSuffixes(automates[3], sequence,
            taille);
28     }
29     else
30     {
31         Nucleotide codonStart[] = { A,T,G };
32         return boyerMoore(sequence, codonStart, taille, 3);
33     }
34 }

```



```

35
36
37
38 File* rechercheCodonsStop(Nucleotide sequence[], int taille ,
AutomateSuffixes** automates)
39 {
40
41     File* codonsStop0 = NULL;
42     File* codonsStop1 = NULL;
43     File* codonsStop2 = NULL;
44
45     if (automates != NULL)
46     {
47         codonsStop0 = rechercheAutomateSuffixes(automates[STOP0],
sequence, taille);
48         codonsStop1 = rechercheAutomateSuffixes(automates[STOP1],
sequence, taille);
49         codonsStop2 = rechercheAutomateSuffixes(automates[STOP2],
sequence, taille);
50     }
51     else
52     {
53         Nucleotide codonStop0[] = { T, A, A };
54         Nucleotide codonStop1[] = { T, A, G };
55         Nucleotide codonStop2[] = { T, G, A };
56         codonsStop0 = boyerMoore(sequence, codonStop0, taille, 3);
57         codonsStop1 = boyerMoore(sequence, codonStop1, taille, 3);
58         codonsStop2 = boyerMoore(sequence, codonStop2, taille, 3);
59     }
60
61
62     File* codonsStop = malloc(sizeof(File));
63
64     if (codonsStop == NULL)
65     {
66         exit(EXIT_FAILURE);
67     }
68
69     codonsStop->dernierElement = NULL;
70     codonsStop->premierElement = NULL;
71
72
73     while (codonsStop0->premierElement != NULL)
74     {
75         if (codonsStop1->premierElement != NULL)
76         {
77             if (codonsStop2->premierElement != NULL)
78             {
79                 if (codonsStop0->premierElement->element <
codonsStop1->premierElement->element)
80                 {
81                     if (codonsStop0->premierElement->element <
codonsStop2->premierElement->element)
82                     {
83                         enfiler(codonsStop, defiler(codonsStop0));
84                     }
85                     else

```

```

86         {
87             enfiler(codonsStop, defiler(codonsStop2));
88         }
89     }
90     else
91     {
92         if (codonsStop1->premierElement->element <
codonsStop2->premierElement->element)
93         {
94             enfiler(codonsStop, defiler(codonsStop1));
95         }
96         else
97         {
98             enfiler(codonsStop, defiler(codonsStop2));
99         }
100     }
101     } // codonStopActuel2 = NULL
102
103     else
104     {
105         if (codonsStop0->premierElement->element <
codonsStop1->premierElement->element)
106         {
107             enfiler(codonsStop, defiler(codonsStop0));
108         }
109         else
110         {
111             enfiler(codonsStop, defiler(codonsStop1));
112         }
113     }
114     } // codonStopActuel1 = NULL
115
116     else
117     {
118         if (codonsStop2->premierElement != NULL)
119         {
120             if (codonsStop0->premierElement->element <
codonsStop2->premierElement->element)
121             {
122                 enfiler(codonsStop, defiler(codonsStop0));
123             }
124             else
125             {
126                 enfiler(codonsStop, defiler(codonsStop2));
127             }
128         } // codonStopActuel2 = NULL
129         else
130         {
131             enfiler(codonsStop, defiler(codonsStop0));
132         }
133     }
134     } // codonStopActuel0 = NULL
135
136     while (codonsStop1->premierElement != NULL)
137     {
138         if (codonsStop2->premierElement != NULL)
139         {

```

```

140         if (codonsStop1->premierElement->element < codonsStop2
->premierElement->element)
141         {
142             enfiler(codonsStop, defiler(codonsStop1));
143         }
144         else
145         {
146             enfiler(codonsStop, defiler(codonsStop2));
147         }
148     }
149     else
150     {
151         enfiler(codonsStop, defiler(codonsStop1));
152     }
153 }
154
155 while (codonsStop2->premierElement != NULL)
156 {
157     enfiler(codonsStop, defiler(codonsStop2));
158 }
159
160 return codonsStop;
161 }
162
163
164
165
166
167
168
169 File* predictionORF(Nucleotide sequence[], int taille,
AutomateSuffixes** automates)
170 {
171     File* codonsStart = rechercheCodonsStart(sequence, taille,
automates);
172     File* codonsStop = rechercheCodonsStop(sequence, taille,
automates);
173
174
175     File* orfs = malloc(sizeof(File));
176     orfs->dernierElement = NULL;
177     orfs->premierElement = NULL;
178
179
180     int tailleSequenceCodante = 0;
181
182
183
184     while (codonsStart->premierElement != NULL && codonsStop->
premierElement != NULL)
185     {
186         // recherche de deux codons stops en phase successifs
distants d'au moins TAILLE
187         int indiceCodonStop0 = defiler(codonsStop); // premier
codon stop
188         int indiceCodonStop1 = defiler(codonsStop); // deuxieme
codon stop

```

```

189     int invalide = (indiceCodonStop1 - indiceCodonStop0) % 3 !=
190     0 || (indiceCodonStop1 - indiceCodonStop0) <= TAILLE; // = 1 !
191     int fileVide = 0;
192
193     while (invalide)
194     {
195         if (codonsStop->premierElement == NULL)
196         {
197             invalide = 0;
198             fileVide = 1;
199         }
200         else if ((indiceCodonStop1 - indiceCodonStop0) % 3 == 0
201         && (indiceCodonStop1 - indiceCodonStop0) < TAILLE)
202         {
203             int indiceCodonStopSuivant = defiler(codonsStop);
204             indiceCodonStop0 = indiceCodonStop1;
205             indiceCodonStop1 = indiceCodonStopSuivant; // il ne
206             doit pas y avoir de codon stop entre les deux !
207         }
208         else if ((indiceCodonStop1 - indiceCodonStop0) % 3 !=
209         0)
210         {
211             indiceCodonStop1 = defiler(codonsStop);
212         }
213         else
214         {
215             invalide = 0;
216         }
217     }
218
219     if (!fileVide) // on a trouve deux codons stops ideaux
220     {
221         if (codonsStart->premierElement != NULL)
222         {
223             int indiceCodonStart = defiler(codonsStart);
224
225             invalide = indiceCodonStart > indiceCodonStop0 &&
226             indiceCodonStart < indiceCodonStop1 && (indiceCodonStart -
227             indiceCodonStop0) % 3 == 0 && indiceCodonStop1 -
228             indiceCodonStart >= TAILLE;
229             fileVide = 0;
230
231             while (invalide)
232             {
233                 if (codonsStart->premierElement == NULL)
234                 {
235                     invalide = 0;
236                     fileVide = 1;
237                 }
238                 else if (indiceCodonStart > indiceCodonStop1)
239                 {
240                     invalide = 0; // le codon start doit etre
241                     entre les deux codons stop donc on sort
242                 }
243             }
244         }
245     }

```

```

238         else if ((indiceCodonStart - indiceCodonStop0)
% 3 != 0 || (indiceCodonStop1 - indiceCodonStart) < TAILLE ||
indiceCodonStart < indiceCodonStop0)
239         {
240             indiceCodonStart = defiler(codonsStart);
241         }
242         else
243         {
244             invalide = 0;
245         }
246     }
247
248
249     // ajout
250     if (!fileVide)
251     {
252
253         enfiler(orfs , indiceCodonStart);
254         enfiler(orfs , indiceCodonStop1);
255
256         tailleSequenceCodante += indiceCodonStop1 - 3 -
indiceCodonStart;
257
258
259
260         // decalage de codonStartActuel pour que l'
indice du codon stop soit superieur a celui du codon start
261
262         int codonStartInvalide = codonsStart->
premierElement != NULL;
263
264         while (codonStartInvalide)
265         {
266             if (codonsStart->premierElement->element >
indiceCodonStop1) // le codon start suivant fonctionne
267             {
268                 codonStartInvalide = 0;
269             }
270             else if (codonsStart->premierElement->
suivant != NULL)
271             {
272                 if (codonsStart->premierElement->
suivant->element > indiceCodonStop1)
273                 {
274                     defiler(codonsStart);
275                     codonStartInvalide = 0;
276                 }
277                 else
278                 {
279                     defiler(codonsStart);
280                 }
281             }
282             else
283             {
284                 defiler(codonsStart);
285                 codonStartInvalide = 0;
286             }

```

```

287     }
288     }
289 }
290 }
291
292 }
293
294 // ajout de la taille en tete de file
295
296 Element* premierElement = orfs->premierElement;
297 Element* nouvelElement = malloc(sizeof(Element));
298 nouvelElement->element = tailleSequenceCodante;
299 nouvelElement->precedent = NULL;
300 nouvelElement->suivant = NULL;
301
302 if (premierElement != NULL)
303 {
304     premierElement->precedent = nouvelElement;
305     nouvelElement->suivant = premierElement;
306 }
307 else
308 {
309     orfs->dernierElement = nouvelElement;
310 }
311
312 orfs->premierElement = nouvelElement;
313
314 return orfs;
315 }
316
317
318
319 Nucleotide* extraireSequencesCodantes(Nucleotide sequence[], int
    taille, AutomateSuffixes** automates)
320 {
321     File* orfsFile = predictionORF(sequence, taille, automates);
322
323     int tailleSequenceCodante = defiler(orfsFile);
324
325     Nucleotide* orfs = malloc(sizeof(int) * (tailleSequenceCodante
    + 1));
326     int indice = 1;
327
328     if (orfs == NULL)
329     {
330         exit(EXIT_FAILURE);
331     }
332
333     orfs[0] = tailleSequenceCodante;
334
335
336
337     while (orfsFile->premierElement != NULL)
338     {
339         int indiceStartCodon = defiler(orfsFile);
340         int indiceStopCodon = defiler(orfsFile);
341

```

```

342     for (int i = indiceStartCodon; i < indiceStopCodon - 3; i
      ++){
343     {
344         orfs[indice] = sequence[i];
345         indice++;
346     }
347 }
348
349
350
351
352     return orfs;
353 }

```

#### prediction\_regions\_codantes.c

```

1  /* Quentin MOREAU
2  * TIPE Bioinformatique
3  * Reconstruction d'arbres phylogenetiques
4  *
5  * fichier d'entete alignement
6  *
7  * 2021-2022
8  */
9
10
11 #ifndef ALIGNEMENT_H_INCLUDED
12 #define ALIGNEMENT_H_INCLUDED
13
14 #define TRANSITION 1
15 #define TRANSVERSION 3
16 #define COUT_INSERTION_DELETION 10
17
18
19 int alignementRecuratif(Nucleotide sequence1[], Nucleotide sequence2
    [], int tailleSequence1, int tailleSequence2);
20 int coutRecuratif(Nucleotide sequence1[], Nucleotide sequence2[],
    int indice1, int indice2);
21 int needlemanWunsch(Nucleotide sequence1[], Nucleotide sequence2[],
    int tailleSequence1, int tailleSequence2);
22 int** distancesLevenshtein(Nucleotide** sequences, int* tailles,
    int nombreSequence);
23
24
25 #endif // ALIGNEMENT_H_INCLUDED

```

#### alignement.h

```

1  /* Quentin MOREAU
2  * TIPE Bioinformatique
3  * Reconstruction d'arbres phylogenetiques
4  *
5  * fichier source alignement
6  *
7  * 2021-2022
8  */
9

```

```

10
11 #include <stdlib.h>
12 #include <stdio.h>
13
14 #include "main.h"
15 #include "alignement.h"
16
17
18
19 int alignementRekursif(Nucleotide sequence1[], Nucleotide sequence2
20 [], int tailleSequence1, int tailleSequence2)
21 {
22     return coutRekursif(sequence1, sequence2, tailleSequence1 - 1,
23     tailleSequence2 - 1);
24 }
25
26 int coutRekursif(Nucleotide sequence1[], Nucleotide sequence2[],
27 int indice1, int indice2)
28 {
29     int const COUTS_SUBSTITUTIONS[4][4] = { {0, TRANSVERSION,
30     TRANSITION, TRANSVERSION},
31     {TRANSVERSION, 0,
32     TRANSVERSION, TRANSITION},
33     {TRANSITION,
34     TRANSVERSION, 0, TRANSVERSION},
35     {TRANSVERSION,
36     TRANSITION, TRANSVERSION, 0} };
37
38     if (indice1 == 0 && indice2 == 0) return 0;
39     if (indice1 == 0) return indice2 * COUT_INSERTION_DELETION;
40     if (indice2 == 0) return indice1 * COUT_INSERTION_DELETION;
41
42     return min(coutRekursif(sequence1, sequence2, indice1 - 1,
43     indice2 - 1) + COUTS_SUBSTITUTIONS[sequence1[indice1]][
44     sequence2[indice2]],
45     min(coutRekursif(sequence1, sequence2, indice1, indice2 -
46     1) + COUT_INSERTION_DELETION,
47     coutRekursif(sequence1, sequence2, indice1 - 1, indice2
48     ) + COUT_INSERTION_DELETION));
49 }
50
51
52
53
54 int** distancesLevenshtein(Nucleotide** sequences, int* tailles,
55 int nombreSequence)
56 {
57     int** distances = malloc(sizeof(int*) * nombreSequence);
58
59     if (distances == NULL)
60     {
61         exit(EXIT_FAILURE);
62     }
63 }

```



```

55
56     for (int i = 0; i < nombreSequence; i++)
57     {
58         distances[i] = malloc(sizeof(int) * nombreSequence);
59         if (distances[i] == NULL)
60         {
61             exit(EXIT_FAILURE);
62         }
63
64         for (int j = 0; j < nombreSequence; j++)
65         {
66             distances[i][j] = 0;
67         }
68     }
69
70     for (int i = 1; i < nombreSequence; i++)
71     {
72         for (int j = 0; j < i; j++)
73         {
74             //distances[i][j] = alignementRecurcif(sequences[i],
75             sequences[j], tailles[i], tailles[j]);
76             distances[i][j] = needlemanWunsch(sequences[i],
77             sequences[j], tailles[i], tailles[j]);
78             distances[j][i] = distances[i][j];
79         }
80     }
81     return distances;;
82 }
83
84
85
86
87
88
89
90
91 int needlemanWunsch(Nucleotide sequence1[], Nucleotide sequence2[],
92     int tailleSequence1, int tailleSequence2)
93 {
94     if (tailleSequence2 < tailleSequence1)
95     {
96         return needlemanWunsch(sequence2, sequence1,
97         tailleSequence2, tailleSequence1);
98     }
99
100     int* colonnes[2] = { NULL, NULL };
101
102     colonnes[0] = malloc(sizeof(int) * (tailleSequence1 + 1));
103     colonnes[1] = malloc(sizeof(int) * (tailleSequence1 + 1));
104
105     int const COUTS_SUBSTITUTIONS[4][4] = { {0, TRANSVERSION,
106     TRANSITION, TRANSVERSION},

```

```

{TRANSVERSION, 0,

```

```

TRANSVERSION,TRANSITION},
107                                     {TRANSITION,
TRANSVERSION,0,TRANSVERSION},
108                                     {TRANSVERSION,
TRANSITION,TRANSVERSION,0}  };
109
110
111     if (colonnes[0] == NULL || colonnes[1] == NULL)
112     {
113         exit(EXIT_FAILURE);
114     }
115
116     for (int i = 0; i <= tailleSequence1; i++)
117     {
118         colonnes[0][i] = i * COUT_INSERTION_DELETION;
119     }
120     for (int i = 0; i <= tailleSequence1; i++)
121     {
122         colonnes[1][i] = i * COUT_INSERTION_DELETION;
123     }
124
125
126
127     int k = 0;
128
129     while (k < tailleSequence2)
130     {
131
132         colonnes[k % 2][0] = (k+1) * COUT_INSERTION_DELETION;
133
134         for (int i = 0; i < tailleSequence1; i++)
135         {
136             colonnes[k % 2][i + 1] = min(colonnes[(k + 1) % 2][i] +
137             COUTS_SUBSTITUTIONS[sequence1[i]][sequence2[k]], /* \ */
138             min(colonnes[k % 2][i] + COUT_INSERTION_DELETION,
139             /* | */
140             colonnes[(k + 1) % 2][i + 1] +
141             COUT_INSERTION_DELETION)); /* - */
142         }
143
144         k++;
145     }
146
147     int cout = colonnes[(k + 1) % 2][tailleSequence1];
148
149
150
151     free(colonnes[0]);
152     free(colonnes[1]);
153
154
155     return cout;
156
157 }

```

---

alignement.c

---

```
1  /* Quentin MOREAU
2  * TIPE Bioinformatique
3  * Reconstruction d'arbres phylogenetiques
4  *
5  * fichier d'entete construction_arbre_phylogenetique
6  *
7  * 2021–2022
8  */
9
10
11
12 #ifndef CONSTRUCTION_ARBRE_PHYLOGENETIQUE_H_INCLUDED
13 #define CONSTRUCTION_ARBRE_PHYLOGENETIQUE_H_INCLUDED
14
15
16 #include "main.h"
17
18
19 typedef struct Arbre Arbre;
20 struct Arbre
21 {
22     char* nom;
23     Arbre** arbres;
24 };
25
26
27 void afficherArbre(Arbre* arbre);
28 void afficherArbreRecuratif(Arbre* arbre, int espace);
29 Arbre* constructionArbreUPGMA(int** distances, char** noms, int
    taille);
30 Arbre* constructionArbreUPGMARecursif(int** distances, Arbre**
    arbres, int taille);
31
32
33
34 #endif // CONSTRUCTION_ARBRE_PHYLOGENETIQUE_H_INCLUDED
```

---

construction\_arbre\_phylogenetique.h

---

```
1  /* Quentin MOREAU
2  * TIPE Bioinformatique
3  * Reconstruction d'arbres phylogenetiques
4  *
5  * fichier source construction_arbre_phylogenetique
6  *
7  * 2021–2022
8  */
9
10
11 #include <stdio.h>
12 #include <stdlib.h>
13 #include <string.h>
14
```

```

15 #include "main.h"
16 #include "construction_arbre_phylogenetique.h"
17
18
19
20
21
22
23 //////////////////////////////////////////////////// AFFICHAGE ARBRE
24 ////////////////////////////////////////////////////
25
26 void afficherArbre(Arbre* arbre)
27 {
28     afficherArbreRecuratif(arbre, 0);
29 }
30
31 void afficherArbreRecuratif(Arbre* arbre, int espace)
32 {
33     if (arbre != NULL)
34     {
35         for (int i = 0; i < espace; i++) printf(" ");
36         printf(" -> ");
37         printf(arbre->nom);
38
39         if (arbre->arbres != NULL)
40         {
41             printf("\n");
42             afficherArbreRecuratif(arbre->arbres[0], espace + 4 +
43             strlen(arbre->nom));
44             printf("\n");
45             afficherArbreRecuratif(arbre->arbres[1], espace + 4 +
46             strlen(arbre->nom));
47             printf("\n");
48         }
49     }
50 }
51
52
53
54
55
56
57
58
59
60 //////////////////////////////////////////////////// UPGMA
61 ////////////////////////////////////////////////////
62
63
64
65 Arbre* constructionArbreUPGMA(int** distances, char** noms, int
66     taille)
67 {

```

```

67     Arbre** arbres = malloc(sizeof(Arbre*) * taille);
68
69     if (arbres == NULL)
70     {
71         exit(EXIT_FAILURE);
72     }
73
74     for (int i = 0; i < taille; i++)
75     {
76         arbres[i] = malloc(sizeof(Arbre));
77         arbres[i]->arbres = NULL;
78         arbres[i]->nom = noms[i];
79     }
80
81     return constructionArbreUPGMARecursif(distances, arbres, taille);
82 }
83
84
85 Arbre* constructionArbreUPGMARecursif(int** distances, Arbre**
86     arbres, int taille)
87 {
88     if (taille <= 1)
89     {
90         return arbres[0];
91     }
92
93     int indice1 = 1, indice2 = 0;
94     int minimum = distances[1][0];
95
96     // recherche de minimum
97     for (int i = 1; i < taille; i++)
98     {
99         for (int j = 0; j < i; j++)
100         {
101             if (distances[i][j] < minimum) // on a indice2<indice1
102             {
103                 indice1 = i;
104                 indice2 = j;
105                 minimum = distances[i][j];
106             }
107         }
108
109         // distances recalculées
110
111         for (int i = 0; i < taille; i++)
112         {
113             distances[i][indice1] = (distances[i][indice2] + distances[
114 i][indice1]) / 2;
115             distances[indice1][i] = distances[i][indice1];
116             distances[i][indice2] = distances[i][indice1];
117             distances[indice2][i] = distances[i][indice1];
118         }
119
120         // insertion de la ligne taille-1 dans la ligne indice1
121         for (int i = 0; i < taille; i++)

```

```

121     {
122         distances[indice1][i] = distances[taille - 1][i];
123         distances[i][indice1] = distances[indice1][i];
124     }
125
126
127
128     if (arbres[indice1] == NULL || arbres[indice2] == NULL)
129     {
130         exit(EXIT_FAILURE);
131     }
132
133     Arbre** sousArbres = malloc(2 * sizeof(Arbre*));
134     Arbre* nouvelArbre = malloc(sizeof(Arbre));
135
136     if (sousArbres == NULL || nouvelArbre == NULL)
137     {
138         exit(EXIT_FAILURE);
139     }
140
141     sousArbres[0] = arbres[indice1];
142     sousArbres[1] = arbres[indice2];
143
144     nouvelArbre->nom = "ancetre commun";
145     nouvelArbre->arbres = sousArbres;
146
147     arbres[indice2] = nouvelArbre;
148     arbres[indice1] = arbres[taille - 1];
149
150
151     return constructionArbreUPGMARecursif(distances, arbres, taille
152         - 1);
153
154
155 }

```

construction\_arbre\_phylogenetique.c

## Références

- [1] IPBES. *Échapper à l'«ère des pandémies» : Les experts mettent en garde contre de pires crises à venir ; Options proposées pour réduire les risques.* URL : [https://ipbes.net/sites/default/files/2020-11/20201029%20Media%20Release%20IPBES%20Pandemics%20Workshop%20Report%20FR\\_Final\\_0.pdf](https://ipbes.net/sites/default/files/2020-11/20201029%20Media%20Release%20IPBES%20Pandemics%20Workshop%20Report%20FR_Final_0.pdf).
- [2] François RECHENMANN. *À la recherche des régions codantes.* URL : <https://interstices.info/a-la-recherche-de-regions-codantes/>.
- [3] Christian D. Wunsch SAUL B. NEEDLEMAN. "A general method applicable to the search for similarities in the amino acid sequence of two proteins". In : *Journal of Molecular Biology* 48 (1970), p. 443-453. ISSN : 0022-2836.

- [4] NCBI. *Taxonomy browser (Viruses)*. URL : <https://www.ncbi.nlm.nih.gov/Taxonomy/Browser/wwwtax.cgi?mode=Undef&name=Viruses&lvl=3&srchmode=1&keep=1&unlock>.