



ALGORITHMIQUE LU3IN003

RAPPORT

Alignement de séquences

Préparé par :
Rami BENELMIR
Erisa KOHANSAL

Année universitaire
2022/2023

Table des matières

1	Le problème d'alignement de séquences	3
1.1	Question 1	3
1.2	Question 2	3
2	Algorithmes pour l'alignement de séquences	4
2.1	Question 3	4
2.2	Question 4	4
2.3	Question 5	4
2.4	Question 6	5
2.5	Tâche A	5
3	Programmation dynamique	7
3.1	Question 7	7
3.2	Question 8	7
3.3	Question 9	7
3.4	Question 10	7
3.5	Question 11	7
3.6	Question 12	8
3.7	Question 13	8
3.8	Question 14	9
3.9	Question 15	9
3.10	Question 16	10
3.11	Question 17	10
3.12	Question 18	10
3.13	Tâche B	11
3.14	Question 19	11
3.15	Question 20	12
3.16	Tâche C	13
3.17	Question 21	14
3.18	Question 22	15
3.19	Question 23	15
3.20	Question 24	16
3.21	Question 25	17
3.22	Question 26	17
3.23	Question 27	18
3.24	Question 28	18
3.25	Tache D	18
3.26	Question 29	19

4	L'alignement local de séquences (BONUS)	19
4.1	Question 30	19
4.2	Question 31	20

1 Le problème d'alignement de séquences

1.1 Question 1

Soit (\bar{x}, \bar{y}) est un alignement de (x, y) et (\bar{u}, \bar{v}) est un alignement de (u, v) .

On veut démontrer que $(\bar{x}.\bar{u}, \bar{y}.\bar{v})$ est un alignement de $(x.u, y.v)$:

$$I : |\bar{x}.\bar{u}| = |\bar{x}| + |\bar{u}| = |\bar{y}| + |\bar{v}| = |\bar{y}.\bar{v}|$$

D'après les informations données, on sait que $\pi(\bar{x}) = \pi(x)$ et que $\pi(\bar{u}) = \pi(u)$. Or :

$$II : \pi(\bar{x}.\bar{u}) = \pi(\bar{x}).\pi(\bar{u}) = x.u$$

$$\text{De même : } \pi(\bar{y}.\bar{v}) = \pi(\bar{y}).\pi(\bar{v}) = y.v$$

$\forall i \in [1, \dots, |\bar{x}|], \bar{x}_i \neq _ , \bar{y}_i \neq _ \text{ et } \forall j \in [1, \dots, |\bar{u}|], \bar{u}_j \neq _ , \bar{v}_j \neq _ \text{ on a :}$

$$III : \forall i + j = k \in [1, \dots, |\bar{x}.\bar{u}|], (\bar{x}.\bar{u})_k \neq _ , (\bar{y}.\bar{v})_k \neq _$$

Conclusion : D'après *I, II* et *III*, cette propriété est vérifiée.

1.2 Question 2

La longueur maximale d'un alignement est alors majorée par $n + m$. Pour cela, on pourrait utiliser le caractère $_$ dans le but de créer un alignement de (x, y) .

C'est le cas de quatrième exemple dans la partie 2.2 :

$$\Sigma = A, T, G, C \quad x = ATTGTA \text{ et } y = ATCTTA.$$

L'alignement de longueur maximale de (x, y) correspond à :

$$\bar{x} = _ _ _ _ _ _ ATTGTA$$

$$\bar{y} = ATCTTA _ _ _ _ _ _$$

Donc on peut avoir jusqu'à $n + m$ de longueur maximum.

Conclusion : la longueur maximale d'un alignement de (x, y) est majorée par $n + m$

2 Algorithmes pour l'alignement de séquences

2.1 Question 3

On cherche le nombre de combinaisons possible de placer les gaps dans l'alignement de $(x, x + k)$. x est un mot de longueur n ($|x| = n$) or en ajoutant k gaps à x on obtient \bar{x} , c'est-à-dire : $|\bar{x}| = |x| + k = n + k$

Conclusion : $C_{\bar{x}+k}^{|x|} = C_{n+k}^n$ nous donne le nombre d'alignement de $(x, x + k)$ possible.

2.2 Question 4

$|x| = n$ et une fois qu'on ajoute k gaps à x , on aura $|\bar{x}| = |x| + k = n + k$

On peut obtenir le nombre de positionnement possible pour les k gaps à l'aide d'un coefficient binomial. Puisque qu'on n'a pas besoin de prendre en compte l'ordre des éléments, on utilise détermine le nombre de combinaison possibles : C_{n+k}^k

Une des conditions d'alignement est que les deux mots alignés n'aient pas de gap à la même position. D'après cette condition et d'après le fait que $n \geq m$ alors pour satisfaire ces deux conditions, k doit être plus petit que m ($k \leq m$).

Néanmoins, le nombre de gap dans y est égal à la différence de taille de x et y plus le k nombre de gaps ajoutés à x (c'est-à-dire $n - m + k$, qui est toujours positif) : C_n^{n-m+k}

En prenant en compte toutes ces contraintes, on aura : $C_{n+k}^k \cdot C_n^{n-m+k}$

Conclusion : $C_{n+k}^k \cdot C_n^{n-m+k}$ Sera le nombre total d'alignements possibles. Ainsi pour $|x| = 15$ et $|y| = 10$ on a :

$$\sum_{k=0}^{10} \frac{(k+15)!15!}{15!k!(10-k)!(5+k)!} = 298199265 \text{ alignements possibles.}$$

2.3 Question 5

La complexité temporelle, c'est le calcul de tous les alignements possibles, pour apres choisir celui qui à la distance la plus minimal.

Pour le calcul de de tout les alignements possibles ce fait en : $\sum_0^m C_{n+k}^k \cdot C_n^{n-m+k}$.

Et pour enumerer tous les alignements pour trouver l'alignements du cout minimal le cout sera négligé, parce que c'est trop petit « par rapport au cout ddu calcul des alignement.

Conclusion : On est sur une complexité temporelle de $O(\sum_0^m C_{n+k}^k \cdot C_n^{n-m+k})$, et c'est très couteux.

2.4 Question 6

Pour trouver l'alignement minimal on doit avoir une variable qui stock la distance, et un tableau de taille max d'alignement qui est $n+m$.

Conclusion : On est sur une complexité spatiale de $O(n+m)$

2.5 Tâche A

On a eu les bons résultats lors d'implémentation de "DIST_NAIF" sur les trois instances "Inst_0000010_44.adn", "Inst_0000010_7.adn" et "Inst_0000010_8.adn".

```

3
4 print("Instances_genome\Inst_0000010_44.adn")
5 _, _, x, y = readFile("Instances_genome\Inst_0000010_44.adn")
6 print("distance d'édition : ", DIST_NAIF(x, y))
7
8 print("Instances_genome\Inst_0000010_7.adn")
9 _, _, x, y = readFile("Instances_genome\Inst_0000010_7.adn")
10 print("distance d'édition : ", DIST_NAIF(x, y))
11
12 print("Instances_genome\Inst_0000010_8.adn")
13 _, _, x, y = readFile("Instances_genome\Inst_0000010_8.adn")
14 print("distance d'édition : ", DIST_NAIF(x, y))
15
PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL JUPYTER
Instances_genome\Inst_0000010_44.adn
distance d'édition : 10
Instances_genome\Inst_0000010_7.adn
distance d'édition : 8
Instances_genome\Inst_0000010_8.adn
distance d'édition : 2

```

FIGURE 1 – Les résultats rendus par "DIST_NAIF" qui est exécutée sur les instances "Inst_0000010_44.adn", "Inst_0000010_7.adn" et "Inst_0000010_8.adn"

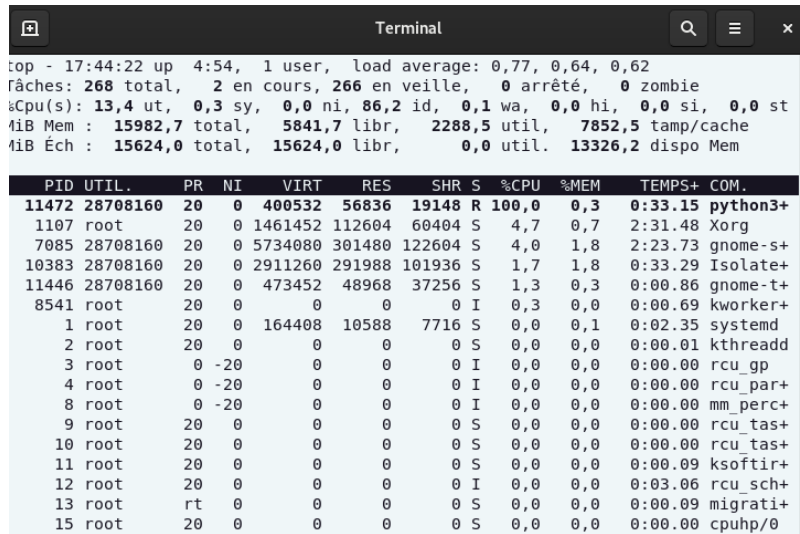
Pour déterminer l'instance pour laquelle le temps d'exécution de la fonction "DIST_NAIF" va dépasser une minute, on lance la fonction "trace_courbe_tacheA" situé dans le fichier "traceFunc.py". Cette fonction se termine lorsque le temps d'exécution d'une instance dépasse une minute, qui est l'instance "Inst_0000012_56.adn".

```

25 tailleA, tempsA = trace_courbe_tacheA("Instances_genome\\")
26
PROBLEMS 1 OUTPUT DEBUG CONSOLE TERMINAL JUPYTER
l'exécution de Instances_genome\Inst_0000012_56.adn a pris 83.41903805732727 seconds qui est plus d'une minute

```

FIGURE 2 – L'exécution de la fonction "DIST_NAIF" prend plus d'une minute pour l'instance "Inst_0000012_56.adn"



Terminal

top - 17:44:22 up 4:54, 1 user, load average: 0,77, 0,64, 0,62
  ches: 268 total, 2 en cours, 266 en veille, 0 arr t , 0 zombie
 %Cpu(s): 13,4 ut, 0,3 sy, 0,0 ni, 86,2 id, 0,1 wa, 0,0 hi, 0,0 si, 0,0 st
 Mem: 15982,7 total, 5841,7 libr, 2288,5 util, 7852,5 tamp/cache
 Mem: 15624,0 total, 15624,0 libr, 0,0 util, 13326,2 dispo Mem

PID	UTIL.	PR	NI	VIRT	RES	SHR	S	%CPU	%MEM	TEMPS+	COM.
11472	28708160	20	0	400532	56836	19148	R	100,0	0,3	0:33.15	python3+
1107	root	20	0	1461452	112604	60404	S	4,7	0,7	2:31.48	Xorg
7085	28708160	20	0	5734080	301480	122604	S	4,0	1,8	2:23.73	gnome-s+
10383	28708160	20	0	2911260	291988	101936	S	1,7	1,8	0:33.29	Isolate+
11446	28708160	20	0	473452	48968	37256	S	1,3	0,3	0:00.86	gnome-t+
8541	root	20	0	0	0	0	I	0,3	0,0	0:00.69	kworker+
1	root	20	0	164408	10588	7716	S	0,0	0,1	0:02.35	systemd
2	root	20	0	0	0	0	S	0,0	0,0	0:00.01	kthreadd
3	root	0	-20	0	0	0	I	0,0	0,0	0:00.00	rcu_gp
4	root	0	-20	0	0	0	I	0,0	0,0	0:00.00	rcu_par+
8	root	0	-20	0	0	0	I	0,0	0,0	0:00.00	mm_perc+
9	root	20	0	0	0	0	S	0,0	0,0	0:00.00	rcu_tas+
10	root	20	0	0	0	0	S	0,0	0,0	0:00.00	rcu_tas+
11	root	20	0	0	0	0	S	0,0	0,0	0:00.09	ksoftir+
12	root	20	0	0	0	0	I	0,0	0,0	0:03.06	rcu_sch+
13	root	rt	0	0	0	0	S	0,0	0,0	0:00.09	migrati+
15	root	20	0	0	0	0	S	0,0	0,0	0:00.00	cpuhp/0

FIGURE 3 – Consommation de m moire par la fonction "DIST_NAIF"

Sur la figure 2, la premi re ligne correspond   l'ex cution de "DIST_NAIF" sur l'instance "Inst_0000012_56.adn". Cette fonction consomme 0.3% de la m moire. Vu que la taille de cette instance n'est en effet pas si grande, cette consommation de m moire est  lev e.

3 Programmation dynamique

3.1 Question 7

Si $\bar{u}_l = _$, d'après les conditions d'alignement on ne peut pas avoir un gap sur la même position dans les deux mots (\bar{u}_l étant le l-ième élément de \bar{u}). Or cela implique que $\bar{v}_l = y_j$ (y_j est le j-ième élément de y). De même pour le cas où on a $\bar{v}_l = _$, alors $\bar{u}_l = x_i$ (x_i est le i-ième élément de x). Maintenant si $\bar{u}_l \neq _$ et $\bar{v}_l \neq _$, cela veut dire que $\bar{u}_l = x_i$ et $\bar{v}_l = y_j$.

3.2 Question 8

$$C(\bar{u}, \bar{v}) = C(\bar{u}_{[1, \dots, l-1]}, \bar{v}_{[1, \dots, l-1]}) + \begin{cases} C_{ins}, & \text{si } \bar{u}_l = 0 \\ C_{del}, & \text{si } \bar{v}_l = 0 \end{cases}$$

3.3 Question 9

On sait que la distance est calculée récursivement, on ajoute à chaque fois le coût minimal des opérations. D'après les deux questions précédentes, on peut conclure :

$$D(i, j) = \min \begin{cases} D(i-1, j) + C_{ins} \\ D(i, j-1) + C_{del} \\ D(i-1, j-1) + C_{sub}(x_i, y_j) \end{cases}$$

3.4 Question 10

$D(0, 0) = d(\emptyset, \emptyset) = 0$ c'est la distance entre deux mots vides.

3.5 Question 11

— $D(0, j)$ consiste en l'ajout de gap dans x . C'est-à-dire qu'on va répéter l'opération d'insertion j fois (un alignement de longueur j) :

$$\begin{aligned} D(0, j) &= D(0, j-1) + C_{ins} \\ &= [D(0, j-2) + C_{ins}] + C_{ins} = D(0, j-2) + 2 \times C_{ins} \\ &= [D(0, j-3) + C_{ins}] + 2 \times C_{ins} = D(0, j-3) + 3 \times C_{ins} \\ &= \dots \\ &= D(0, j-j) + j \times C_{ins} = D(0, 0) + j \times C_{ins} = j \times C_{ins} \end{aligned}$$

D'après la question précédente $D(0, 0) = 0$.

Alors on a : $D(0, j) = j \times C_{ins}$

- $D(i, 0)$ consiste en l'ajout de gap dans y . C'est-à-dire qu'on va répéter l'opération de suppression i fois de suite (un alignement de longueur i) :

$$\begin{aligned}
 D(i, 0) &= D(i - 1, 0) + C_{del} \\
 &= [D(i - 2, 0) + C_{del}] + C_{del} = D(i - 2, 0) + 2 \times C_{del} \\
 &= [D(i - 3, 0) + C_{del}] + 2 \times C_{del} = D(i - 3, 0) + 3 \times C_{del} \\
 &= \dots \\
 &= D(i - i, 0) + i \times C_{del} = D(0, 0) + i \times C_{del} = i \times C_{del}
 \end{aligned}$$

D'après la question précédente $D(0, 0) = 0$.

Alors on a : $D(i, 0) = i \times C_{del}$

Conclusion : $D(0, j) = j \times C_{ins}$ et $D(i, 0) = i \times C_{del}$

3.6 Question 12

Algorithm 1 DIST_1

Require: Deux mots x, y

Ensure: Un tableau T

```

1:  $n = \text{longueur}(x)$ 
2:  $m = \text{longueur}(y)$ 
3:  $T[n][m] = 0$ 
4: for  $j = [1 \dots n]$  do
5:    $T[j][0] \leftarrow j * C_{del}$ 
6: end for
7: for  $j = [1 \dots m]$  do
8:    $T[0][j] \leftarrow j * C_{ins}$ 
9: end for
10: for  $i = [1 \dots n]$  do
11:   for  $j = [1 \dots m]$  do
12:      $T[i][j] \leftarrow \min(T[i - 1][j - 1] + C_{sub}(x[i - 1], y[j - 1])); T[i][j - 1] + C_{ins}; T[i - 1][j] + C_{del})$ 
13:   end for
14: end for
15: return  $T$ 

```

3.7 Question 13

La complexité spatiale de DIST_1 est en $\theta(n \times m)$ car on stocke le coût dans chaque case de notre matrice T .

3.8 Question 14

La complexité temporelle de `DIST_1` est en $\theta(n \times m)$ car on a deux boucles imbriquées allant de 0 à $n-1$ et de 0 à $m-1$. Toutes les autres opérations sont des opérations élémentaires.

3.9 Question 15

On cherche à démontrer que :

Si $i > 0$ et $D(i, j) = D(i - 1, j) + C_{del}$, alors $\forall (\bar{s}, \bar{t}) \in Al^*(i - 1, j)$,
 $(\bar{s}.x_i, \bar{t}._) \in Al^*(i, j)$

Commençons par

$$\begin{aligned} D(i, j) &= D(i - 1, j) + C_{del} \\ &= d(x_{[1, \dots, i]}, y_{[1, \dots, j]}) + C_{del} \\ &= C(\bar{s}.x_i, \bar{t}._) \end{aligned}$$

Or $C(\bar{s}.x_i, \bar{t}._) \in Al^*(i, j)$

3.10 Question 16

Algorithm 2 SOL_1

Require: Deux mots x , y et un tableau T

Ensure: Un tableau T

```

1:  $i \leftarrow \text{longueur}(x)$ 
2:  $j \leftarrow \text{longueur}(y)$ 
3:  $\bar{x} = [ ]$ 
4:  $\bar{y} = [ ]$ 
5: while  $i > 0$  and  $j > 0$  do
6:   if  $i > 0$  and  $j > 0$  and  $T(i,j) = T(i-1,j-1) + \text{Csub}(x[i-1], y[j-1])$  then
7:      $\bar{x} = x_i.\bar{x}$ 
8:      $\bar{y} = y_j.\bar{y}$ 
9:      $i \leftarrow i - 1$ 
10:     $j \leftarrow j - 1$ 
11:   else
12:     if  $i > 0$  and  $T(i,j) = T(i,j-1) + \text{Cdel}$  then
13:        $\bar{x} = x_i.\bar{x}$ 
14:        $\bar{y} = -.\bar{y}$ 
15:        $i \leftarrow i - 1$ 
16:     else
17:       if  $j > 0$  and  $T(i,j) = T(i,j-1) + \text{Cins}$  then
18:          $\bar{x} = -.\bar{x}$ 
19:          $\bar{y} = y_j.\bar{y}$ 
20:          $j \leftarrow j - 1$ 
21:       end if
22:     end if
23:   end if
24: end while
25: return  $(\bar{x}, \bar{y})$ 

```

3.11 Question 17

Notre problème ALI se résolve avec deux fonctions :

— DIST_1 se calcule en $O(nm)$.

— SOL_1 se calcule en $O(n+m)$.

Conclusion : On a une complexité temporelle de $O(nm)$.

3.12 Question 18

Notre problème ALI se résolve avec deux fonctions :

- DIST_1 qui alloue un tableau de taille $n \times m$, ce qui a une complexité spatiale de $O(nm)$.
- SOL_1 qui alloue deux tableaux de taille n et m , ce qui donne une complexité spatiale de $O(n+m)$.

Conclusion : On a une complexité spatiale de $O(nm)$.

3.13 Tâche B

On utilise la fonction "trace_courbe_tacheB" qu'on a implémentée dans le fichier "traceFunc.py" afin d'obtenir la courbe de consommation de temps CPU en fonction de la taille $|x|$: La première ligne de la figure 3 correspond à la

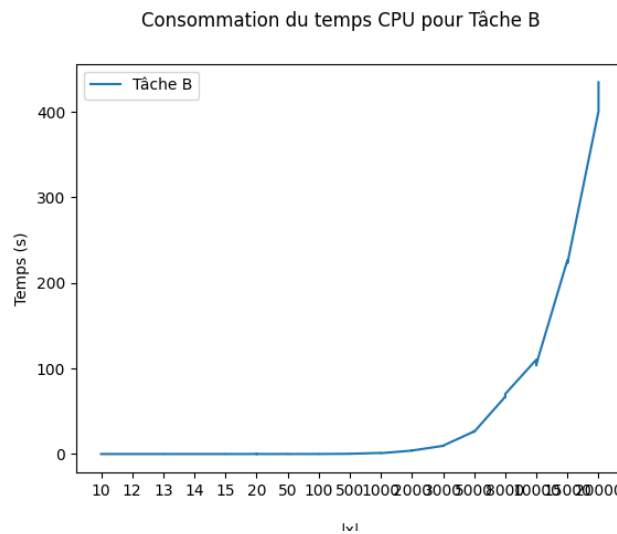
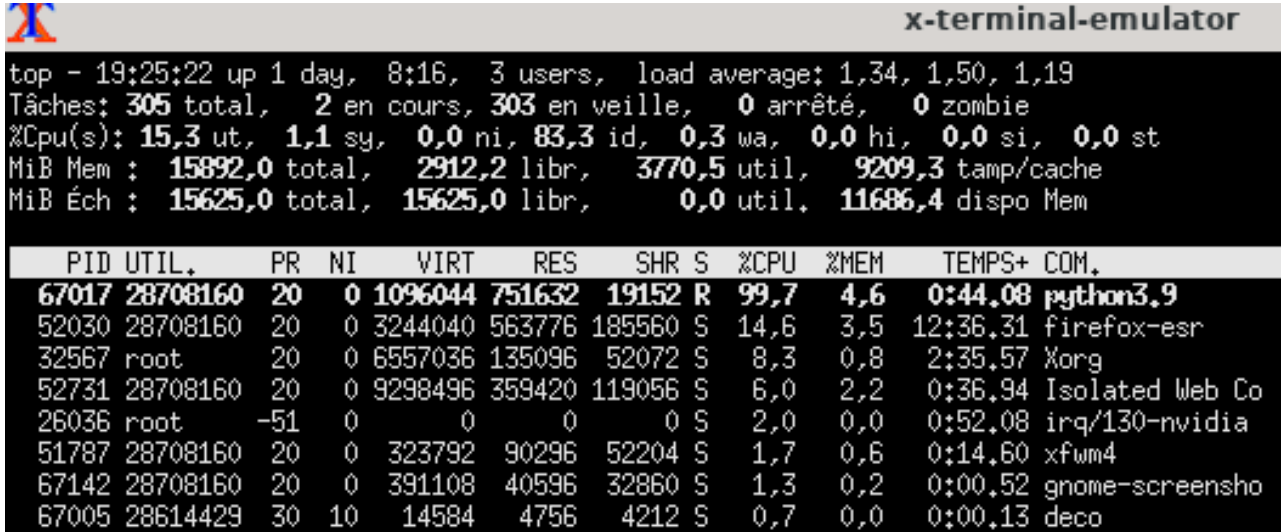


FIGURE 4

consommation mémoire de "PROG_DYN" exécutée sur une instance de grande taille "Inst_0010000_7.adn". Cette consommation est très importante, vu que ça consomme 4.6% de la mémoire.

3.14 Question 19

Car à chaque tour de boucle la valeur minimum de $D(i, j)$ est calculée à partir de $D(i-1, j-1)$, $D(i-1, j)$, $D(i, j-1)$ or on a besoin que de ces deux lignes i et $i-1$. Alors, on pourra utiliser une matrice à deux lignes afin d'améliorer la complexité spatiale.



PID	UTIL.	PR	NI	VIRT	RES	SHR	S	%CPU	%MEM	TEMPS+	COM.
67017	28708160	20	0	1096044	751632	19152	R	99,7	4,6	0:44.08	python3.9
52030	28708160	20	0	3244040	563776	185560	S	14,6	3,5	12:36.31	firefox-esr
32567	root	20	0	6557036	135096	52072	S	8,3	0,8	2:35.57	Xorg
52731	28708160	20	0	9298496	359420	119056	S	6,0	2,2	0:36.94	Isolated Web Co
26036	root	-51	0	0	0	0	S	2,0	0,0	0:52.08	irq/130-nvidia
51787	28708160	20	0	323792	90296	52204	S	1,7	0,6	0:14.60	xfwm4
67142	28708160	20	0	391108	40596	32860	S	1,3	0,2	0:00.52	gnome-screensho
67005	28614429	30	10	14584	4756	4212	S	0,7	0,0	0:00.13	deco

FIGURE 5 – Consommation mémoire de la fonction "PROG_DYN" exécutée sur l'instance "Inst_0010000_7.adn"

3.15 Question 20

Algorithm 3 DIST_2

Require: Deux mots x , y

Ensure: Distance entre x et y

```

1:  $n = \text{longueur}(x)$ 
2:  $m = \text{longueur}(y)$ 
3:  $T[2][m] = 0$ 
4: for  $j = [1 \dots m]$  do
5:    $T[0][j] \leftarrow i * C_{ins}$ 
6: end for
7: for  $i = [0 \dots n]$  do
8:   for  $j = [0 \dots m]$  do
9:      $T[1][j] \leftarrow \min(T[0][j - 1] + C_{sub}(x[i - 1], y[j - 1])); T[1][j - 1] +$ 
        $C_{ins}; T[0][j] + C_{del})$ 
10:   end for
11:   for  $j = [0 \dots m]$  do
12:      $T[0][j] \leftarrow T[1][j]$ 
13:   end for
14: end for
15: return  $T$ 

```

3.16 Tâche C

Pour le traçage de la courbe de consommation de temps CPU en fonction de la taille $|x|$, on a implémenté la fonction "trace_courbe_DIST_2" dans le fichier "traceFunc.py".

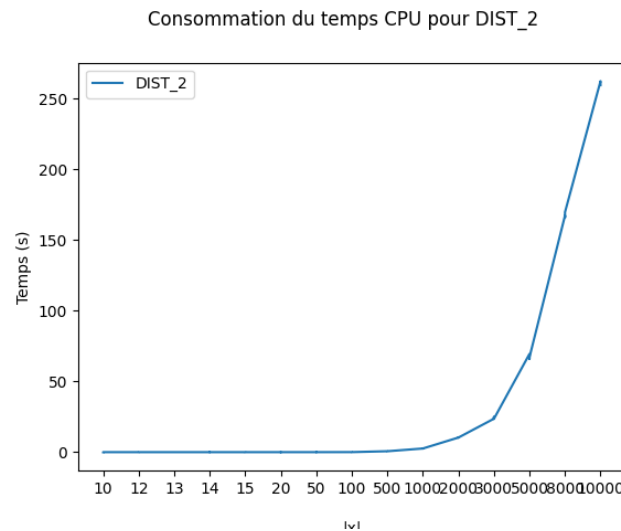


FIGURE 6

On peut alors comparer les deux fonctions "DIST_1" et "DIST_2" :

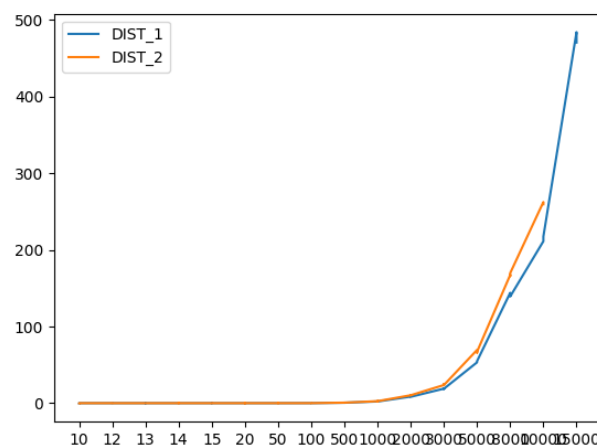
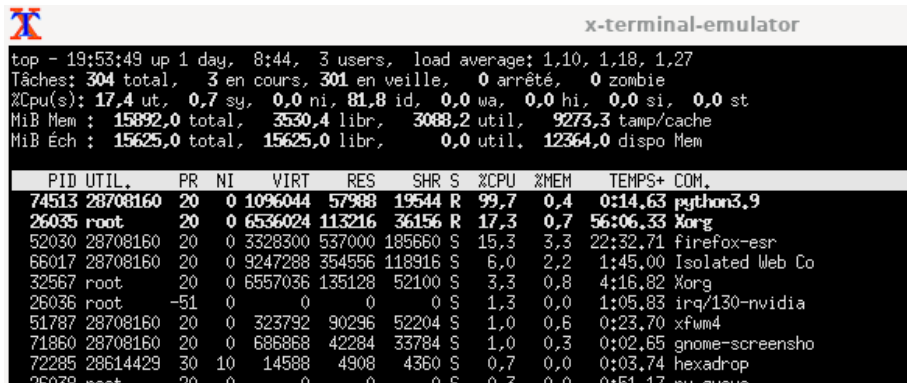


FIGURE 7

D'après la figure 7, on peut voir que la courbe de "DIST_2" monte moins vite par rapport à celle de "DIST_1". Cela veut dire qu'on a réussi à améliorer la complexité



```

top - 19:53:49 up 1 day, 8:44, 3 users, load average: 1.10, 1.18, 1.27
Tâches: 304 total, 3 en cours, 301 en veille, 0 arrêté, 0 zombie
%Cpu(s): 17.4 ut, 0.7 sy, 0.0 ni, 81.8 id, 0.0 wa, 0.0 hi, 0.0 si, 0.0 st
MiB Mem : 15892.0 total, 3530.4 libr, 3088.2 util, 9273.3 tam/cache
MiB Éch : 15625.0 total, 15625.0 libr, 0.0 util, 12364.0 dispo Mem

  PID  UTIL.  PR  NI  VIRT  RES  SHR  S  %CPU  %MEM  TEMPS+  COM.
 74513 28708160 20 0 1096044 57988 19544 R 99.7 0.4 0:14.63 python3.9
 26035 root 20 0 6536024 113216 36156 R 17.3 0.7 56:06.33 Xorg
 52030 28708160 20 0 3328300 537000 185660 S 15.3 3.3 22:32.71 firefox-esr
 66017 28708160 20 0 9247288 354556 118916 S 6.0 2.2 1:45.00 Isolated Web Co
 32567 root 20 0 6557036 135128 52100 S 3.3 0.8 4:16.82 Xorg
 26036 root -51 0 0 0 0 0 S 1.3 0.0 1:05.83 irq/130-nvidia
 51787 28708160 20 0 323792 90296 52204 S 1.0 0.6 0:23.70 xfwm4
 71860 28708160 20 0 686868 42284 33784 S 1.0 0.3 0:02.65 gnome-screensho
 72285 28614429 30 10 14588 4908 4360 S 0.7 0.0 0:03.74 hexadrop
 95028 root 20 0 0 0 0 0 S 0.7 0.0 0:51.17 su - su

```

FIGURE 8 – Consommation de la mémoire par la fonction "DIST_2" pour une instance d'une grande taille, ici on a utilisé "Inst_0010000_7.adn"

temporelle de "DIST_2" par rapport à celle de "DIST_1".

La première ligne sur la figure 8 correspond aux détails de l'exécution de la fonction "DIST_2" sur l'instance "Inst_0010000_7.adn". Cette exécution occupe 0.4% de la mémoire. En comparant ce résultat avec celui de "DIST_1" dans les questions précédentes, on conclut que la consommation de mémoire a diminué de $4.6\% - 0.4\% = 4.4\%$.

Conclusion : En conclusion, "DIST_2" a une complexité spatiale (resp. une complexité temporelle) qui est plus petite que celle de "DIST_1".

3.17 Question 21

Algorithm 4 mot_gaps

Require: Une entier k

Ensure: Un mot de k gaps

```

1: mot = [ ]
2:  $mot \leftarrow -' -' * k$ 
3: return (mot)

```

3.18 Question 22

Algorithm 5 align_lettre_mot

Require: Deux mots x, y **Ensure:** Alignement de x et y

```

1: for  $i = [0 \dots \text{longueur}(y) - 1]$  do
2:   if  $x[i] == y[i]$  then
3:     Sortir de la boucle
4:   end if
5: end for
6:  $x \leftarrow \text{mot\_gaps}(i) + x + \text{mot\_gaps}(\text{longueur}(y) - i - 1)$ 
7: return  $(x, y)$ 

```

3.19 Question 23

On découpe x et y au milieu : $x_1 : \text{BAL}$, $x_2 : \text{LON}$, $y_1 : \text{RO}$, $y_2 : \text{ND}$.

Alignement de $(x_1, y_1) : (\text{B}, \text{A}, \text{L})$, $(\text{R}, \text{O}, -)$ avec une distance de 13.

Alignement de $(x_2, y_2) : (\text{L}, \text{O}, \text{N}, -)$, $(-, -, \text{N}, \text{D})$ avec une distance de 9.

Alignement de $(x_1.x_2, y_1.y_2) : (\text{B}, \text{A}, \text{L}, \text{L}, \text{O}, \text{N}, -)$, $(\text{R}, \text{O}, -, -, -, \text{N}, \text{D})$ avec une distance de 22.

Alors l'alignement optimal est : $(\text{B}, \text{A}, \text{L}, \text{L}, \text{O}, \text{N}, -)$, $(\text{R}, -, -, -, \text{O}, \text{N}, \text{D})$ qui a une distance de 17.

Conclusion : par absurde $(\text{B}, \text{A}, \text{L}, \text{L}, \text{O}, \text{N}, -)$, $(\text{R}, \text{O}, -, -, -, \text{N}, \text{D})$ n'est pas un alignement optimal.

3.20 Question 24

Algorithm 6 SOL_2_REC

Require: Deux mots x, y
Ensure: Distance entre x et y

```

1:  $n = \text{longueur}(x)$ 
2:  $m = \text{longueur}(y)$ 
3: if  $n == 0$  then
4:   return  $(\text{mot\_gaps}(m), y)$ 
5: end if
6: if  $m == 0$  then
7:   return  $(x, \text{mot\_gaps}(n))$ 
8: end if
9: if  $n = 1$  then
10:  return  $\text{align\_lettre\_mot}(x, y)$ 
11: else
12:   $\text{index\_i} = n \text{ div } 2$ 
13:   $\text{index\_j} = \text{coupure}(x, y)$ 
14:   $(x1, y1) \leftarrow \text{SOL\_REC\_2}(x_{[0..\text{index\_i}]}, y_{[0..\text{index\_j}]})$ 
15:   $(x2, y2) \leftarrow \text{SOL\_REC\_2}(x_{[\text{index\_i}+1..n]}, y_{[\text{index\_j}+1..m]})$ 
16: end if
17: return  $(x1 + x2, y1 + y2)$ 

```

3.21 Question 25

Algorithm 7 Coupure

Require: Deux mot x et y

Ensure: L'indice j^*

```

1: T : Une liste de deux case pour les distances
2:  $c_{li}$  : Une liste de deux cases pour stocker la colonne sur la ligne i
3: for  $j = [0 \dots longueur(y)]$  do
4:   T[0].append( j * Cins )
5:   T[1].append( -1 )
6:    $c_{li}[0].append( j )$ 
7:    $c_{li}[0].append( -1 )$ 
8: end for
9: for  $i = [1 \dots longueur(n) \text{ div } 2]$  do
10:  T[1][0]  $\leftarrow i * Cdel$ 
11:  for  $j = [1 \dots longueur(y)]$  do
12:    T[1][j]  $\leftarrow \min (T[0][j-1] + Csub(x[i-1], y[j-1]) ; T[1][j-1] + Cins ; T[0][j] + Cdel )$ 
13:  end for
14: end for
15: for  $i = [longueur(n) \text{ div } 2 + 1 \dots longueur(n)]$  do
16:  T[1][0]  $\leftarrow i * Cdel$ 
17:  for  $j = [1 \dots longueur(y)]$  do
18:    T[1][j]  $\leftarrow \min (T[0][j-1] + Csub(x[i-1], y[j-1]) ; T[1][j-1] + Cins ; T[0][j] + Cdel )$ 
19:    If  $T[1][j] = T[1][j-1] + Cins$  then
20:       $c_{li}[1][j] \leftarrow c_{li}[1][j-1]$ 
21:    If  $T[1][j] = T[0][j] + Cdel$  then
22:       $c_{li}[1][j] \leftarrow c_{li}[0][j]$ 
23:    If  $T[1][j] = T[0][j-1] + Csub(x[i-1], y[j-1])$  then
24:       $c_{li}[1][j] \leftarrow c_{li}[0][j-1]$ 
25:    end for
26:     $c_{li}[0] \leftarrow c_{li}[1]$ 
27:  end for
28: return  $c_{li}[0][m]$ 

```

3.22 Question 26

La complexité spatiale de notre fonction coupure est de $O(m)$, étant donné qu'on deux tableau (liste pour notre cas) de deux cases de taille m.

3.23 Question 27

On sait que la complexité de coupure est en $O(m)$, on a : $n > m$ donc : $O(n)$. Et à chaque appel récursif on décompose les deux mots x et y en deux et on lance deux appels. En utilisant le théorème maître on trouve : $T(n) = 2 T(\frac{n}{2}) + O(n)$. SOL_2_REC est de complexité : $O(n \log(n))$.

3.24 Question 28

Etant donné qu'on doit faire $n*m$ itérations pour trouver la bonne coupure, la complexité temporelle de coupure est de $O(nm)$.

3.25 Tache D

La fonction "trace_courbe_tacheD" dans le fichier "traceFunc.py" nous permet de tracer la courbe de consommation de temps CPU de la tâche D en fonction de $|x|$. La première ligne de la figure ci-dessus contient les détails de l'exécution de la

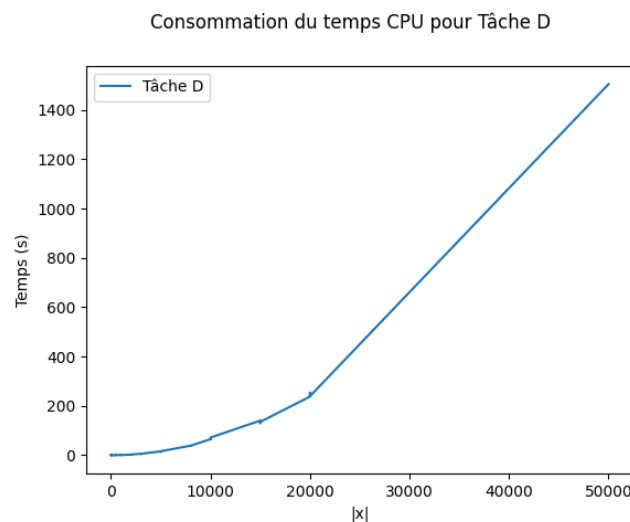


FIGURE 9

fonction "coupure" pour une grande instance, "Inst_0010000_7.adn". On a obtenu cette figure grâce à la commande "top" de Linux. La consommation mémoire étant égale à 0.4%, cette valeur n'est pas une grande consommation par rapport à la taille de notre instance.

```

top - 20:06:37 up 1 day, 8:57, 4 users, load average: 2.13, 1.76, 1.53
Tâches: 316 total, 4 en cours, 312 en veille, 0 arrêté, 0 zombie
%Cpu(s): 24,3 ut, 2,3 sy, 1,2 ni, 71,8 id, 0,4 wa, 0,0 hi, 0,0 si, 0,0 st
MiB Mem : 15892,0 total, 2786,4 libr, 3776,1 util, 9329,5 temp/cache
MiB Éch : 15625,0 total, 15625,0 libr, 0,0 util, 11614,8 dispo Mem

```

PID	UTIL.	PR	NI	VIRT	RES	SHR	S	%CPU	%MEM	TEMPS+	COM.
78203	28708160	20	0	1097004	58712	19016	R	100,3	0,4	0:10,34	python3.9
26035	root	20	0	6536024	113216	36156	R	69,7	0,7	57:44,20	Xorg
52030	28708160	20	0	3657868	748828	211532	S	16,0	4,6	27:58,19	firefox-esr
77637	28614429	30	10	24796	15108	4496	S	15,0	0,1	0:18,79	metaballs
32567	root	20	0	6574604	152320	67564	S	9,7	0,9	5:53,60	Xorg
75250	28708160	20	0	9221952	360080	119240	S	6,3	2,2	0:20,85	Isolated Web Co
26036	root	-51	0	0	0	0	S	2,0	0,0	1:14,83	irq/130-nvidia

FIGURE 10 – La consommation mémoire de la fonction "coupure"

3.26 Question 29

Oui, en voulant améliorer la complexité spatiale on a perdu en complexité temporelle. On remarque bien que la complexité temporelle de SOL_2_REC est bien supérieur de celle de SOL_1.

4 L'alignement local de séquences (BONUS)

4.1 Question 30

En s'inspirant des instances fournies, on a créé trois nouvelles instances "Instance_0010000_5.adn", "Instance_0015000_9.adn" et "Instance_0020000_12.adn" situés dans le dossier "Question30". Dans le fichier "main.py", on a exécuté la fonction "DIST_2" sur ces trois instances et on a obtenu les résultats ci-dessous :

```

25  _, _ , x, y = readFile("Question30\\Instance_0010000_5.adn")
26  print(DIST_2(x,y))
27  _, _ , x, y = readFile("Question30\\Instance_0015000_9.adn")
28  print(DIST_2(x,y))
29  _, _ , x, y = readFile("Question30\\Instance_0020000_12.adn")
30  print(DIST_2(x,y))

```

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL JUPYTER

```

19990
29982
39976

```

FIGURE 11

Ces résultats sont cohérents avec notre hypothèse :

$$|x| = 10000, |y| = 5, (|x| - |y|) \times C_{del} = 19990$$

$$|x| = 15000, |y| = 9, (|x| - |y|) \times C_{del} = 29982$$

$$|x| = 20000, |y| = 12, (|x| - |y|) \times C_{del} = 39976$$

Conclusion : On conclut que le coût d'un alignement global (x, y) quand $|x| \gg |y|$ relativement à $|\sum|$ vaut $|x| - |y| \times C_{del}$

4.2 Question 31

Cela semble être une bonne idée pour les alignements petits, mais dans les pire des cas, c'est-à-dire quand la taille d'alignement est $n + m$ ($|x| = n, |y| = m$), dans ce cas là on doit d'abord parcourir tout l'alignement afin d'enlever les gaps, donc une boucle en $n + m$. Puis, une deuxième boucle $n + m$.

On aura une complexité de $O(n + m)^2$. ce qu'est coûteux.