

# IFT3295 -TP1

Louis-Andre Brassard  
Maggie Robert

October 2022

## Chevauchement de séquences

### 1)

La différence entre ce genre d'alignement et l'alignement global se divise en deux points

- Dans le cas de l'alignement local, nous initialisons la première rangée et la première colonne à 0 ce qui donne  $V(i,0) = 0$  et  $V(0,j) = 0$  contrairement à l'alignement globale où on initialise  $V(i,0) = i$  et  $V(0,j) = j$

- Pour l'alignement local, on trouve la valeur maximal dans tout le tableau et on part de cette valeur pour et on suit les pointeurs pour retourner au 0 le plus proche. Le trajet parcourus avec les pointeurs nous donne l'alignement local le donc le score est le plus élevé. Pour l'alignement globale, nous remplissons la table de programmation dynamique au complet et on part de la case  $V(i,j)$  pour retourner en suivant les pointeurs vers la case  $V(0,0)$ .

### 2)

Comme mentionné plus haut, les valeurs de  $V(i,0) \forall i$  et  $V(0,j) \forall j$  sont de 0. Ceci s'explique par le fait que l'alignement local optimal peut commencer et terminer n'importe où dans notre mot et pas seulement à l'origine  $V(0,0)$ .

### 3)

La fonction de récurrences pour ce genre d'alignement local est:

$$V(i, 0) = V(0, j) = 0 \quad \forall i, j$$
$$V(i, j) = \max \begin{cases} V(i-1, j-1) + \delta(v_i, w_j) \\ V(i-1, j) + \delta(v_i, -) \\ V(i, j-1) + \delta(-, w_j) \end{cases}$$

où

$$\delta = \begin{cases} 4 & \text{if } v_i = v_j, \\ -4 & \text{if } v_i \neq v_j, \\ -8 & \text{si } v_i \text{ ou } v_j \text{ est un indel} \end{cases}$$

4)

On commence par initialiser les pointeurs des conditions initiales. Pour chaque  $V(i, 0)$  les pointeurs point vers le haut. Pour chaque  $V(0, j)$  il n'y a pas de pointeurs. Ensuite on remplit la table de programmation dynamique en suivant les formules de récurrence montrées plus haut. Après avoir rempli la table au complet, on prend la case de la dernière colonne avec la plus grande valeur (excluant celle de la première ligne), qui se trouve à la ligne  $k$ .  $k$  représente le nombre de lettres du deuxième mot qui appartiennent à l'alignement. On ajoute à l'envers les  $len(mot2) - k$  lettres du deuxième mot et le même nombre d'espaces au premier mot. Ensuite on suit les pointeurs, en ajoutant à l'envers les lettres de l'alignement, jusqu'à ce que l'on croise une case sans un pointeur dans la première ligne qui se trouve à la colonne  $l$ . Finalement on ajoute à l'envers les  $l + 1$  premières lettres du premier mot et le même nombre d'espaces au deuxième mot. Ceci donne le chevauchement suffixe-préfixe optimal entre les deux séquences.

5)

Voir pièce jointe pour code

## Assemblage de fragments

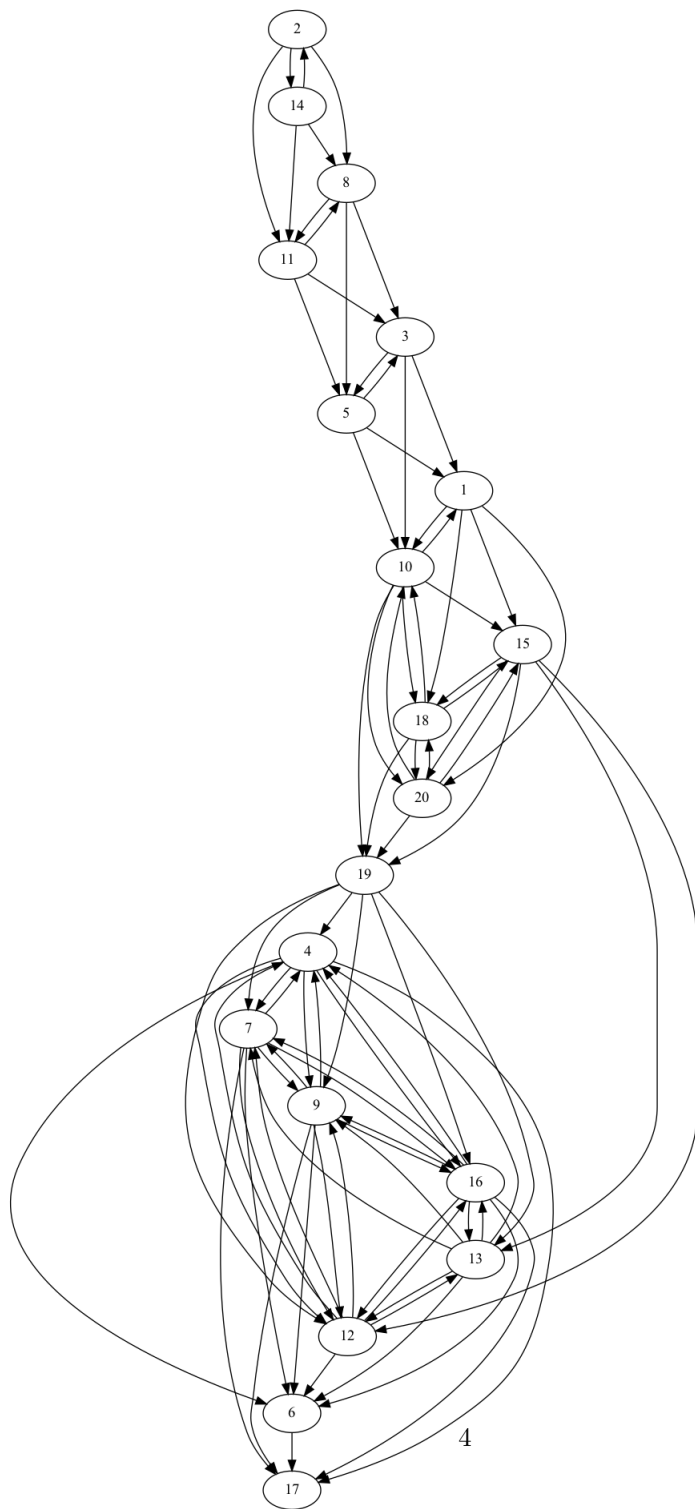
1)

La matrice des scores, avec les scores  $\geq 80$  en orange.

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
1	-	8	8	0	4	12	12	16	4	804	20	4	-4	4	396	4	20	652	12	608
2	4	-	8	12	44	8	4	416	8	12	424	16	0	1060	12	12	12	8	4	0
3	396	12	-	24	1016	16	48	0	36	128	4	20	40	8	20	20	16	8	8	8
4	28	8	28	-	12	640	1024	16	1068	12	12	556	16	24	12	848	152	16	8	16
5	500	4	856	28	-	24	52	8	52	208	16	20	36	4	8	24	36	64	4	12
6	52	4	4	16	0	-	4	4	12	28	4	16	4	0	28	8	672	16	12	20
7	32	12	32	832	16	720	-	4	908	12	12	300	24	24	24	560	232	12	4	0
8	8	12	644	4	540	16	8	-	8	4	1008	12	4	12	0	4	8	36	4	16
9	36	12	20	1012	16	668	1044	4	-	8	4	456	12	24	20	740	196	8	4	12
10	228	0	0	8	32	8	8	24	12	-	32	4	12	12	640	4	8	904	184	868
11	20	20	620	4	524	12	12	1072	4	12	-	24	28	24	8	8	16	20	12	12
12	28	28	48	900	16	460	820	0	856	4	16	-	452	24	16	996	20	16	8	4
13	68	32	32	744	24	272	656	0	692	16	8	860	-	36	16	848	32	24	4	28
14	16	908	12	12	52	20	8	452	12	16	492	16	8	-	4	12	20	16	8	16
15	20	8	28	12	16	12	16	12	8	4	12	80	260	28	-	28	20	168	632	452
16	20	16	24	1032	16	552	944	12	988	8	20	844	288	20	12	-	80	12	8	8
17	16	4	32	16	16	32	8	12	4	24	20	12	4	28	12	8	-	16	8	20
18	16	0	8	12	8	4	4	8	8	568	16	8	4	16	808	8	4	-	368	1004
19	48	20	16	388	8	20	316	8	368	16	20	520	732	16	24	492	20	20	-	4
20	8	4	0	4	0	8	4	0	0	356	4	0	64	8	908	4	16	820	428	-

2)

a)



b)

En choisissant les arêtes avec les scores les plus élevés, on obtient les chevauchements suivants:



c)

La séquence finale est ci-dessous. Sa longueur est 1581 caractères.

```
CTAGGGACTTGGAGAACACAAGTATTATGAAAAGTACTGATGAAA
GTTATTAACAGGTTTCGAAAAATAACTTTACTATCTGACGTGTTGC
TTCTGCCGAGGATGACCGTTATTCTGGTTTTGCATTTATATTTCA
CGTATGGTTAAATGTGCCAGCGTTGTGGTTTAAAACTAATAGTAAT
AATATGCTTCTTTGTTTCAGTTGGCTAGAGATTTACTACATCCGTCC
TTGGAAGAGGAAAAGAAAAAACATAACAAGAAACGCCTAGTACAA
AGTCCAAATTCTTACTTTTATGGATGTAAAATGTCCCGGTCAAATTT
GAAATCTTAATTCTTTTACTAAAGAAAATTTCTGAAGGGATTGCTA
GTGTGGTGTGTATAGTTAAGATACATTTCGAATCCTCTGTTGAGTAG
AAGTGGGATTACAGAATTGGAAATGTGAGGGACATTTTCATAATAA
ATGTACTGAAACCATTGTAGAAACATTTTCGGGGTAGTGAAAACAGG
CCTTCAAAGGAGATATGCCATTTGACTGGACTTGTGGATAATCAAA
ATGTGGATACACAAAATAGAGATGTTCTTTAATTGTAAAATACTCA
CTGGATTTTGTACGATGTAGCACAGAAAAAAAATACATTGATTACAC
GTTTTTAAAAATTTTTGGGTGCTGCTAGAAAAGTTTATGTTACAC
ATGGGGCTTGCTGTTTTCATAGCACTGAAGTTAATGATTTTTTTTACA
TATTACCTGAAATCTCGAACAGGTCCTGTTTTCTCTGCTTTTCAT
TTTTAACATTGCCTTTTTTTTTTTTTTAGGTTGCTACAAGATAACCA
CGGTTTTTCAGCCATGCTCAGACAGTGGTTCTTTGTGTAGGTTGTTT
AACAGTGTTGTGCCAGCCTACAGGAGGAAAGGCCAGACTCACAGAA
GGTATACATTTGTCATTCTCCAACCCAGTGATGAGATTGATGATTA
TAAATGTCTCTATCTTCACTGAAAAGTTTACAGAAATCTTAATGAT
TCCCAAAATAACTTATCTCACACTGGAAGAGTTCAAGTGGATTGGC
AGCAAATCTGAGATCTATTTGGTGTGACCTGGTGAGATCTAAATAT
GGAGTCAGCACATGATTTTTTAAAGAGTAATATTGCTAAGTAATATT
GCTAAGTATAGTCTGAAAATACCTCTAATCAAAATTTTTTACTTGA
GAAAAGTATTCAGTATAGTTCCTAAAAATTAAGAGTATATTTCTGG
TATAAAAGGATAAATATTCTGTATATGAGTATTAATCCAATATGCT
TAAAACTTCAGTATTTTACTTAAAAGTACTGTTTGTGCTATAAAATT
ATACCAAAGTAGAATGCACTTGTTTAATATACTCTCATGATTCTTT
TGCAGGGTGTTTCAATTTAGAAGAAAGCAACACTAATGATTCAAACAG
CTTCCTGAATTTTAATTTTGTGTTGTCTCACAGAAAGCCTTATCAT
AAATTCCATAATTCTAATTAATTTACCAAGATAATGTCATTACATT
TGGTTATGTAAGTTATACAGCAGTAATCTCCTATTTTGGTGTGCTAGT
TTTTCACTAAAGTTTTTAT
```

## Recherche d'introns et Blast

a)

En utilisant le code standard et en décortiquant le brin codant selon les trois cadre nous obtenons les suite d'acides aminé que voici

cadre 1 :

```
RTKFALPYVS*SCESPFLGTWRTQVL*KVLMKVINRFRKITLLSDVLLLP  
RMTVIPVFACIFHVWLNVPALWFKTNSNNMLLCSVG*RFTTSVLGRGK  
EKT*KETPSTKSKFLLYGCKMSR*NLKS*FLY*RKFL*GLLVWCV*LRYI  
RTLCL*VEVGLQNWKCQGHFHNKCTETIVETFRGK*KQAFKGDMPFDC  
TCG*SKCGYSK*RCSLIVKYSLDFDDVAQKKNTLITLFLKNFVSLLEKFM  
LHMWLAVS*H*SY*FFYILPEMSNRCCFLLFIFNIAFFFF*FATRSRPSA  
MLRQWFFV*VVQQCCASLQEERPDSQKDYHLAFSNPVMRLMIINVSIFT  
EKFKEILMITKITYLSLEEFKWIGSKSEIYLV*PGEI*IWSQHMIFLRVILLS  
NIAKYSKLIPLIKIHYLRKVFRIVPKN*EYISGIKG*IICI*VLIQYS*NFSILLK  
STFCH*NYSKGRMHLFNILS*FFCRLFI*KKATLMIQTAS*ILILCCLTESLII  
NSIILINLPR*CNYIWFCKVYSSNLLFWCQFFNKVLIMGK
```

Plus longue séquence entre un codon start et un codon stop pour le cadre 1 :

```
MKVINRFRKITLLSDVLLLP
```

cadre 2 :

```
GQSLPYHMFDPDRAKALF*GLGEHKYYEKY**KLLTGFEK*LYYLTCFC  
RG*PLFLFLHVFYFTYG*MCQRCGLKLIVICFFVQLARDLLHPSLEEEKK  
KHKKKRLVQSPNSYFMDVKCPGKI*NLNSFTKENFCRDC*CGVYS*DTL  
EPSVE*KWDYRIGNVRDIFIINVLPK*KHFGVSENRRHSKEICHLTALVD  
NQNVDTQNRDVL*L*NTHWILTM*HRKKIH*LHCF*KILCRC*KSLCYTC  
GLLFHSTEVTDFFTTYQKCRGTGAVFLCFSFLTLPFFFFSLLQDHHGFQP  
CSDSGSLCRLFNSVVPAYRRKGQTHRRIIIWHSTQ**D**L*MSLSSLKSL  
KKS**LPK*LISHWKSSSGLAANLRSIWCDLVRSKYGVST*FF*E*YC*VIL  
LSIV*KYL*SKLFT*EKYSE*FLKIKSIFLV*KDK*SVYEY*SNILKTSVFYL  
KVLVFIKIIAKVECTCLYSHDSFADCSFRRKQH**FKQLPEF*FCVVSQK  
ALS*IP*F*LIYQDNVITFGFVRYTAVISYFGVSFSIKF*LWA
```

Plus longue séquence entre un codon start et un codon stop pour le cadre 2 :

```
MCQRCGLKLIVICFFVQLARDLLHPSLEEEKKHKKKRLVQSPNSYFMDVKCPGKI
```

cadre 3 :

```
DKVCLTICFLIVRKPFSDLENTSIMKSTDESY*QVSKNNFTI*RVASAED  
DRYSCFCMYISRMVKCASVVV*N****YASLFSWLEIYYIRPWKRKRKNI  
KRNA*YKVQILTLWM*NVQVKFEILIPLLKKISVGIASVVCIVKIH*NPLLS  
RSGITELEMSGTFS**MY*NHCRNISC*VKTGIQRRYAI*LHLWIIKMWIL  
KIEMFFNCKILTGF*RCSTEKKYIDYTVFKKFCVAARKVYVTHVACCFI  
ALKLLIFLHITRNVEQVLFSSAFHF*HCLFFFLVCYKITTVFSHAQTVVLC
```

VGCSTVLCQPTGGKARLTEGLSFGILQPSDEIDDYKCLYLH\*KV\*RNLN  
 DYQNNLSLTGRVQVDWQQI\*DLFGVTW\*DLNMESAHDFKSNIAK\*YC\*  
 V\*SENTSNQNYLLEKSIQNSS\*KLRVYFWYKRINNLYMSINPIFLKLQYFT  
 \*KYFLSLKL\*QR\*NALV\*YTMILLQIVHLEESNTNDSNSFLNFNFLVLSHR  
 KPYHKFHNSN\*FTKIM\*LHLVL\*GIQQ\*SPILVSVFQ\*SFDYGQ

Plus longue séquence entre un codon start et un codon stop pour le cadre 3 :  
 MILLQIVHLEESNTNDSNSFLNFNFLVLSHRKPYHKFHNSN

À l'aide d'un algorithme qui scan les cadres pour trouver codons de départs "M", on cherche le "\*" suivant le codon de départ qui sont séparés par un nombre suffisant d'acide aminé soit 30, ce qui correspond à 90 nucléotides. Le résultat de cet algorithme nous donne une suite d'acides aminés dans chaque cadre qui nous indique la protéine que l'on cherche.

Lorsqu'on isole les séquences commençant par un "M" et finissant par un "\*" on peut voir que le début de notre gèneX se trouve dans le cadre 2

**b)**

Pour résoudre ce genre de problème, nous croyons que le meilleur algorithme de programmation dynamique disponible est celui de pondération affine.

Cet algorithme a les conditions initiales et relations de récurrences suivantes :

Conditions initiales

$$\begin{aligned} V(i, 0) &= F(i, 0) = -\rho - i \cdot \sigma \\ V(0, j) &= E(0, j) = -\rho - j \cdot \sigma \end{aligned}$$

Récurrences

$$V(i, j) = \max \begin{cases} V(i-1, j-1) + \delta(v_i, w_j) \\ E(i, j) \\ F(i, j) \end{cases}$$

$$E(i, j) = \max \begin{cases} E(i, j-1) - \sigma \\ V(i, j-1) - \rho - \sigma \end{cases}$$

$$F(i, j) = \max \begin{cases} F(i-1, j) - \sigma \\ V(i-1, j) - \rho - \sigma \end{cases}$$



Où  $\rho$  est le coût d'ouverture d'un gap et  $\sigma$  le coût de continuation d'un gap

Pour initier le tableau de  $(n+1 \times m+1)$  cases, on doit garder en tête 3 valeurs par case soit  $V(i,j)$ ,  $E(i,j)$  et  $F(i,j)$ . Cela nous servira à calculer les cases suivantes en sachant si nous avons un match, mismatch, que nous commençons un gap ou que nous le poursuivons. Les valeurs de  $V$ ,  $E$  et  $F$  sont initialisées selon les règles de conditions initiales. Si elles ne figurent pas dans les conditions initiales elles sont initialisées à  $-\infty$ . Une fois que le tableau est bien initialisé, on descend dans le tableau en suivant les règles de récurrences montrées plus haut en prenant bien soin de mémoriser les pointeurs qui ont permis de calculer la valeur observée. Une fois que le tableau est rempli, on part de la case  $(n+1, m+1)$  et on regarde la valeur la plus élevée parmi  $V, E, F$  et on suit son pointeur jusqu'à ce qu'on arrive à  $V(0,0)$ ,  $E(0,0)$  ou  $F(0,0)$ .

Nous avons choisi cette façon de faire puisque nous voulons minimiser les introns. Les introns peuvent être représentés par des gaps dans la séquence et la pondération affine vise à réduire le nombre de gaps dans notre alignement en lui donnant un poids élevé pour l'ouverture d'un gap, mais en ne pénalisant pas autant sévèrement la continuation d'un gap. Donc nous avons pensé que cet algorithme serait parfait pour cette tâche.

c)

La protéine X du Dr Osbourne porte le nom de ribosomal protein S27 like [Homo sapiens]. Selon le National Library of Science américain, cette protéine sert à "selectively regulates the expression and alternative splicing of inflammatory and immune response genes in thyroid cancer cells"