

Bio-informatique : comparaison de séquences et calcul d'arbres phylogénétiques

Jean-Stéphane Varré

Université Lille 1 - CRIStAL - Inria Lille-Nord Europe

JEIA - 24 février 2016

Qu'est-ce que la bio-informatique ?

L'ordinateur est une extension de la paillasse pour mener des expériences en biologie moléculaire.

*C'est de la biologie **in silico**.*

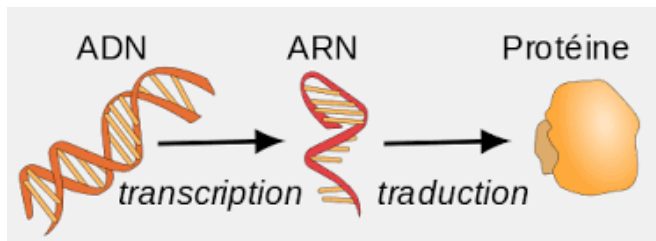
- ▶ acquisition et organisation des données biologiques
- ▶ conception de logiciels pour l'analyse, la comparaison et la modélisation des données
- ▶ analyse des résultats produits par les logiciels

Qu'est-ce que la bio-informatique ?

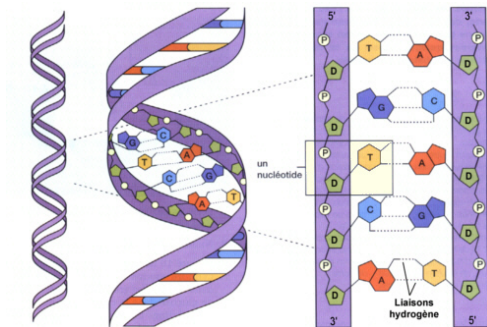
en anglais, distinction entre :

- ▶ « **Bioinformatics** » : applique des algorithmes, modèles statistiques dans l'objectif d'interpréter, classer et comprendre des données biologiques,
- ▶ « **Computational Biology** » : développer des modèles mathématiques et outils associés pour résoudre des problèmes biologiques.

Dogme central



L'ADN

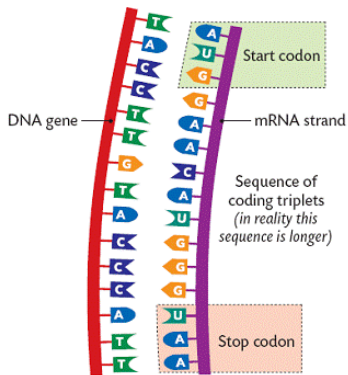


- ▶ c'est l'enchaînement de nucléotides le long d'une macromolécule d'ADN
- ▶ une séquence d'ADN est représentée par un texte écrit sur les 4 lettres A, T, C et G
- ▶ le brin complémentaire est la correspondance exacte de l'autre brin :
A ↔ T et C ↔ G
- ▶ ayant un sens de lecture (5' vers 3')

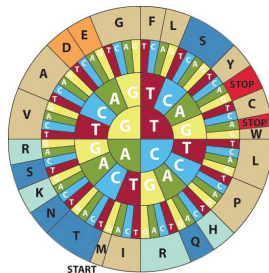
La séquence du gène ABO

```
>gi|58331215|ref|NM_020469.2| Homo sapiens ABO blood group (transferase A, alpha)
GGAGGCCGAGACCAGACGCGGAGCCATGGCCGAGGTGTTGCGGACGCTGGCCGAAAACAAAATGCCAC
GCACCTTCGACCTATGATCCTTTTCCTAATAATGCTTGTCTTGGTCTTGTGGTTACGGGGTCCCTAAGCC
CCAGAAGTCTAATGCCAGGAAGCCTGGAACGGGGGTTCTGCATGGCTGTTAGGGAACCTGACCATCTGCA
GCGCGTCTCGTTGCCAAGGATGGTCTACCCCCAGCCAAAGGTGCTGACACCGTGTAGGAAGGATGTCCTC
GTGGTGACCCCTTGGCTGGCTCCCATTGTCTGGGAGGGCACATTCAACATCGACATCCTCAACGAGCAGT
TCAGGCTCCAGAACACCACCATTGGGTTAACTGTGTTTGCCATCAAGAAATACGTGGCTTTCCTGAAGCT
GTTCTCGAGACGCGGAGAAGCACTTCATGGTGGGCCACCGTGTCCACTACTATGTCTTCACCGACCAG
CCGGCCGCGGTGCCCCGCGTGACGCTGGGGACCGGTCGGCAGCTGTCAGTGCTGGAGGTGCGCGCCTACA
AGCGCTGGCAGGACGTGTCCATGCGCCGCATGGAGATGATCAGTGACTTCTGCGAGCGGCGCTTCCTCAG
CGAGGTGGATTACCTGGTGTGCGTGGACGTGGACATGGAGTTCGCGACCACGTGGGCGTGGAGATCCTG
ACTCCGCTGTTCCGGCACCCCTGCACCCCGGCTTCTACGGAAGCAGCCGGGAGGCCTTACCTACGAGCGCC
GGCCCCAGTCCCAGGCCTACATCCCCAAGGACGAGGGCGATTTCTACTACCTGGGGGGGTTCTTCGGGGG
GTCGGTGCAAGAGGTGCAGCGGCTCACAGGGCCTGCCACCAGGCCATGATGGTGCACCAGGCCAACGGC
ATCGAGGCCGTGTGGCACGACGAGAGCCACCTGAACAAGTACCTGCTGCGCCACAAACCCACCAAGGTGC
TCTCCCCCGAGTACTTGTGGGACCAGCAGCTGCTGGGCTGGCCCGCCGTCCTGAGGAAGCTGAGGTTAC
TGCGGTGCCCAAGAACCACCAGGCGGTCCGGAACCCGTGAGCGGCTGCCAGGGGCTCTGGGAGGGCTGCC
GGCAGCCCCGTCCCCCTCCCGCCCTTGGTTTTAGCAGAACGGGTAACCTCTGTTTTCTTTGTCCGTCTG
TTGTGAGTAACTGAAGCCTAGGCCCCGTCCCCACCTCAAATCACACACACCCCTCCCCACCACAGAGAC
ACCATTACATACACAGACACACACAGAAAGACACACACAGACACAAAATCACACACACACCCCTCCCCGCC
ACAGAGACACCATTACATACACAGACACACACAGAAAGACACAGACACAAAATCACACACACACCCCTCCC
CGCCACAGAGACACACCATTACATACACAGACACGCAATCGCAGATACGCCCTTCCGGCCACAGAAACAC
ACCATTACACACACATACACAGAAAGACACACACAGACACACAATCACACGCAGCCCCTCCCCGCCACAG
AGACACACCATTACATACACAGACACACACAGAAAGACAC
```

Du gène à la protéine



- ▶ c'est la suite de nucléotides au sein du gène qui détermine la suite d'**acides aminés** qui compose une protéine
- ▶ décryptage grâce au code génétique



Exemple

133,275,214

GGAGGCCGAGACAGACGCGGAGCCATGGCCGAGGTGTTGCCGACGCTGGCCG
gtgagtgcaggcctcgccccgggt.....tgatttttctactcctgttttcag
GAAAACCAAAATGCCACGCACTTCGACCTATGATCCTTTTCTAATAATGCTTGTCTGG
TCTTGTITGG
gtaagacacatttgaccatcgaggc.....gtcttggcacacttcctttctgcag
TTACGGGGTCTTAAGCCCCAGAAGTCTAATGCCAGGAAGCTGGAACGGGGGTTCTG
gtgagtgcagggaagagcaggtgga.....catctcctgtgtttctattctgcag
CATGGCTGTTAGGAAACCTGACCATCTGCAGCGCGTCTCGTTGCCAAG
gtataatgtcagtgcctcccttcag.....acgtggcggcgctttgctgcttcgag
GATGGTCTACCCCGAGCCAAAGGTGCTGACACCGTG
gtgagtaagtactgcactgaaa.....accgcacgcctctctccatgtgcag
TAGGAAGGATGCTCTGTG

gt
ACCCCTTGGCTGGCTCCCAATTGTCTGGGAGGGCACATTCAACATCGACATCCTCAACGAG
CAGTTCCAGGCTCCAGAAACACCAACATTGGGTTAACTGTGTTGCCATCAAGAA
gtaagtcaagtgcaggtggccgagggt.....cagccccgtccgctgcttcgag
ATACGTGGCTTTCCTGAAGCTGTTCTCGGAGACGCGGGAGAAGCACTTCATGGTGGGCCA
CGGTGTCCACTACTATGTCTTACCAGACAGCCGGCGGGTCCCGCGTGACCGTGGG
GACCGGTCCGACGCTGCACTGCTGGAGGTGCCGCGCTACAAGCGCTGGCAGGACGTGTC
CATGGCCGCGATGGAGATGATCAGTGACTTCTCGGAGCGCGCTTCTCAGCGAGGTGGA
TTACCTGGTGTGGCTGGACGTGGACATGGAGTTCGCGGACACGCTGGGCGTGGAGATCCT
GACTCCGCTGTTCCGGCACCTGCACCCCGGCTTCTACGGAAGCAGCGGGAGGCCCTTAC
CTACAGCGCCCGCCCACTCCAGGCTACATCCCAAGGACGAGGGCGATTCTACTA
CTGGGGGGGTTCTTCCGGGGTCCGTGCAAGAGGTGCAGCGCTCACCAGGGCCTGCCA
CCAGCCATGATGGTGCAGCCAGGCCAACGGCATCGAGGCGGTGGCCAGCAGAGAGCCA
CCTGAACAAGTACTCTGCTGCCACAAACCCACCAAGGTGCTCTCCCGAGTACTTGTG
GGACAGCAGCTGCTGGCTGGCCCGCGCTCCTGAGGAAGCTGAGGTTCACTGCGGTGCC
CAAGAACCCAGCGCGGTCCGGAACCCGTGAGCGGCTGCCAGGGCTCTGGAGGGCTGC
CGGCAGCCCGTCCCGCTCCCGCTTGGTTTTAGCAGAACGGTAAACTCTGTTTCCTT
TGTCGCTCCTGTTGTGAGTAACTGAAGCCTAGGCCCGTCCCACTCAAATCACACACA
CCCCCTCCCAACACAGAGACACATTACATACAGACACACAGAAAGACACACACA
GACACAAATCACACACACAACCTCCCGCCACAGACACCAATTACATACAGACACA
CACAGAAAGACACAGACACAATACACACACACCCCTCCCGCCACAGAGACACACCAT
TACATACAGACAGCAATCGCAGATAGCCCTTCCGGCCACAGAAACACACACATTAGA
CACACATACAGAAAGACACACAGACACACAATACACAGCAGCCCTCCCGCCACA
GAGACACACATTACATACAGACACACAGAAAGACAC

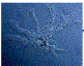
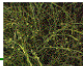







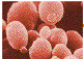

133,255,176

$$53 + 70 + 57 + 48 \\ + 36 + 19 + 113 + 1181 = 1580$$

354 acides aminés

maevlrtrlag kpkchalrpm ilflimlvlv
lfgygvlspr slmpgslerg fcmavrepdh
lqrvsilprmv yppqkvltpc rkdvlvtpw
lapivwegtf nidilneqfr lqnttigitv
faikkyvaf1 klfletaekh fmvghrvhy
vftdqaavp rvtlgtgrql svlevraykr
wqdvsmrrme misdfcerrf lsevdy1vcv
dvdmeifrdhv gveiltplfg tlhpgfygss
reaftyerpp qsqa1ypkde gdfy1ylgff
ggsvqevqrl trachqammv dqang1eavw
hdeshlnkyl lrhkptkvl1s peylwdq1ll
gwpav1rklr ftavpkn1ha vrn1p

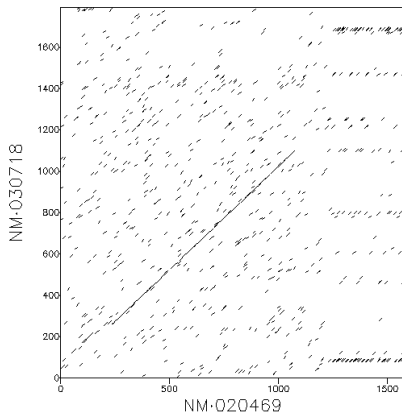
Taille des objets manipulés

		Taille du génome (nucléotides)	Nbre de gènes (protein-coding)	
	<i>Amoeba dubia</i>	~ 670 000 000 000	?	
	<i>Psilotum nudum</i>	~ 250 000 000 000	?	
	<i>Fritillaria assyriaca</i>	~ 100 000 000 000	?	
	<i>Necturus lewisi</i>	~100 000 000 000	?	
	<i>Homo sapiens</i>	2 900 000 000	23 000	
	<i>Vitis vinifera</i>	487 000 000	30 400	
	<i>Drosophila melanogaster</i>	160 000 000	14 000	
	<i>Arabidopsis thaliana</i>	115 000 000	28 000	
	<i>Caenorhabditis elegans</i>	98 000 000	19 400	
	<i>Saccharomyces cerevisiae</i>	12 500 000	5 800	
	<i>Escherichia coli</i>	4 600 000	4 300	

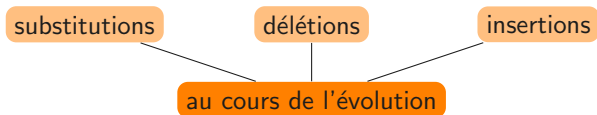
Activités

- A Etant donné une séquence, calculer son taux de GC.
- B Etant donné une séquence, calculer son complémentaire inversé
- C Etant donné une séquence que l'on sait codant pour un gène, calculer la séquence protéique associée
- D Etant donné une séquence, extraire toutes les ORF (Open Reading Frame) : toutes les sous-séquences débutant par un codon START, se terminant par un codon STOP, de taille supérieure à 3 codons, sur les deux brins (au premier STOP rencontré, la traduction s'arrête)

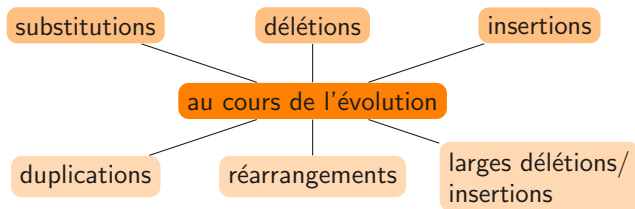
Comparaison de séquences



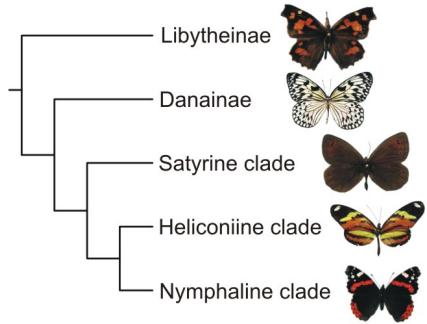
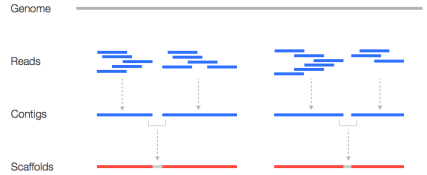
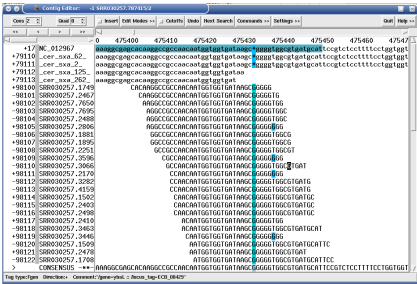
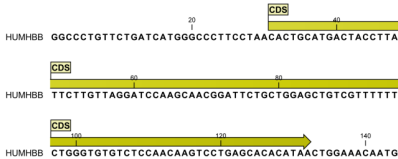
Pourquoi existe-t-il des séquences similaires ?



Pourquoi existe-t-il des séquences similaires ?



De l'utilité de savoir comparer



Examples

```
G A C T C C G
|         | | |
G C T A C C G
```

3 substitutions

Examples

```
G A C T C C G
|           | | |   3 substitutions
G C T A C C G
```

```
G A C T - C C G
|  | |   | | |   2 indels
G - C T A C C G
```


Exemples

G A C T C C G

| | | |

G C T A C C G

3 substitutions

G A C T - C C G

| | | | | |

G - C T A C C G

2 indels

- ▶ calculer la similarité = aligner
- ▶ mesurer la similarité = obtenir l'alignement optimal
- ▶ nécessite de calculer un score

Exemples

```
G A C T C C G
|         | | |
G C T A C C G
```

3 substitutions

```
G A C T - C C G
| | | | |
G - C T A C C G
```

2 indels

- ▶ calculer la similarité = aligner
- ▶ mesurer la similarité = obtenir l'alignement optimal
- ▶ nécessite de calculer un score

A rapprocher de ...

distance de Hamming entre deux textes : nombre minimum de substitutions

distance de Levenshtein entre deux textes : nombre minimum de substitutions et d'indels

Comparaison du gène ABO homme/souris

NM_020469	0	-----	0
NM_030718	1	gtgttcagagctgtgtattatctcccctgggggtgagttctcctgtgtgc	50
NM_020469	1	-----ggaggccgagaccagacgaggagcca	26
		
NM_030718	51	ctgagacctggcctgtgcctaacagctgtgtg-ca-cagacactgaacca	98
NM_020469	27	tggcc-gaggtgt-tg-c-ggacgct-ggcccgaaaaccaaattgccac-	70
		
NM_030718	99	t--cctg-ggt-tctgacatgaatctcag-aggaaagaccgaaatgcaact	143
NM_020469	71	gcacttcgacctatg--atccttttcc-taataatgcttgtcttggctctt	117
		
NM_030718	144	tc-cttc-acct-tggaatcc-ttccttcgagtgttgtcttagtctt	189
NM_020469	118	gtttggttacggggtcctaagcccagaagtctaatagccaggaagcctgg	167
		
NM_030718	190	ctttggctacctgttcctaagc-----t-t-----ca-gaagcc-ag	223
NM_020469	168	aacgggggttctgcatggctgttagggaacctgaccatctgcagcgcgtc	217
NM_030718	224	aac-----t--tg--gg-t-----c---acc-----cag-gag-c	243
NM_020469	218	tcgttgccaagg-atggtctacccccagcaaaagtgctgacaccgtgta	266
NM_030718	244	t-g-tgactaggaat-gcctatctgcagccaagggtgctaaaacccta	290
NM_020469	267	ggaagatgtctcctcgtggtagcccttggtgctcccattgtctgggag	316
		
NM_030718	291	ggaaagatgttcttgtcttgactccttggtggcgcacctcatctgggag	340
NM_020469	317	ggcaccattcaacatcgacatcctcaacgagcagttcaggctccagaacac	366
		
NM_030718	341	gggaccttcaacatcgacatattgaatgagcagttcaggattcggaaatac	390
NM_020469	367	caccattgggttaactgtgtttgccatcaagaaa-tacgtggct-ttcct	414
		
NM_030718	391	tacgattggaactgactgtatttgctatcaa-aaagtatgtgg-tgttcct	438

Alignement

- ▶ données :
 - ▶ une paire de séquences (ADN / protéine)
 - ▶ une schéma de score : comment compter ce qui se ressemble ?

Alignement

- ▶ données :
 - ▶ une paire de séquences (ADN / protéine)
 - ▶ une schéma de score : comment compter ce qui se ressemble ?
- ▶ but :
 - ▶ déterminer le degré de similarité (meilleur score)
 - ▶ montrer la similarité (meilleur alignement)

Alignement

- ▶ données :
 - ▶ une paire de séquences (ADN / protéine)
 - ▶ un schéma de score : comment compter ce qui se ressemble ?
- ▶ but :
 - ▶ déterminer le degré de similarité (meilleur score)
 - ▶ montrer la similarité (meilleur alignement)
- ▶ décrit la ressemblance grâce à 3 opérations (**mutations ponctuelles**)
 - ▶ **insertion**
 - ▶ **délétion**
 - ▶ **identité/substitution**

Alignement

- ▶ données :
 - ▶ une paire de séquences (ADN / protéine)
 - ▶ une schéma de score : comment compter ce qui se ressemble ?
- ▶ but :
 - ▶ déterminer le degré de similarité (meilleur score)
 - ▶ montrer la similarité (meilleur alignement)
- ▶ décrit la ressemblance grâce à 3 opérations (**mutations ponctuelles**)
 - ▶ **insertion**
 - ▶ **délétion**
 - ▶ **identité/substitution**
- ▶ mesure la ressemblance en donnant un poids à chaque opération
 - ▶ poids positif (“récompense”) aux *bonnes parties* de l’alignement
e.g appariement de deux lettres identiques ou proches
 - ▶ poids négatif ou nul (“pénalité”) aux *mauvaises parties* de l’alignement
e.g appariement de deux lettres non relatés, non-appariement

Composantes d'un schéma de scores

- ▶ score (ou poids) pour une identité/substitution : matrice s de similarité
 - ▶ exemple : $s(a, b)$ = score d'alignement des nucléotides a et b

$$\begin{pmatrix} +1 & -1 & -1 & -1 \\ -1 & +1 & -1 & -1 \\ -1 & -1 & +1 & -1 \\ -1 & -1 & -1 & +1 \end{pmatrix}$$

Composantes d'un schéma de scores

- ▶ score (ou poids) pour une identité/substitution : matrice s de similarité

- ▶ exemple : $s(a, b)$ = score d'alignement des nucléotides a et b

$$\begin{pmatrix} +1 & -1 & -1 & -1 \\ -1 & +1 & -1 & -1 \\ -1 & -1 & +1 & -1 \\ -1 & -1 & -1 & +1 \end{pmatrix}$$

- ▶ score (ou poids) d'un *indel* (insertion/délétion)

- ▶ exemple : score unitaire = -2 par *indel*

Composantes d'un schéma de scores

- ▶ score (ou poids) pour une identité/substitution : matrice s de similarité

- ▶ exemple : $s(a, b)$ = score d'alignement des nucléotides a et b

$$\begin{pmatrix} +1 & -1 & -1 & -1 \\ -1 & +1 & -1 & -1 \\ -1 & -1 & +1 & -1 \\ -1 & -1 & -1 & +1 \end{pmatrix}$$

- ▶ score (ou poids) d'un *indel* (insertion/délétion)

- ▶ exemple : score unitaire = -2 par *indel*

- ▶ **score de l'alignement** = somme des scores des événements élémentaires

- ▶ exemple :

A	A	C	G	T	A	C	G	A	T	A
A	A	C	G	T	A	-	A	A	G	A

1	1	1	1	1	1	-2	-1	1	-1	1 = 4

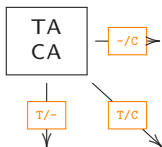
Comment calculer le meilleur alignement ?

- ▶ prenons 2 séquences de longueur n toutes les deux : alignement de longueur maximale $2n$
- ▶ exemple avec les séquences TA et CA

TA
CA

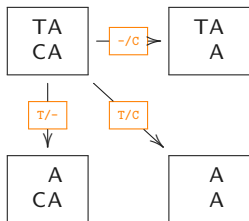
Comment calculer le meilleur alignement ?

- ▶ prenons 2 séquences de longueur n toutes les deux : alignement de longueur maximale $2n$
- ▶ exemple avec les séquences TA et CA



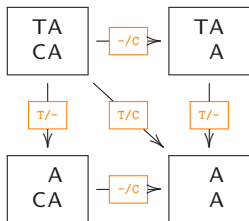
Comment calculer le meilleur alignement ?

- ▶ prenons 2 séquences de longueur n toutes les deux : alignement de longueur maximale $2n$
- ▶ exemple avec les séquences TA et CA



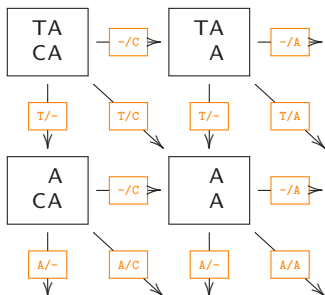
Comment calculer le meilleur alignement ?

- ▶ prenons 2 séquences de longueur n toutes les deux : alignement de longueur maximale $2n$
- ▶ exemple avec les séquences TA et CA



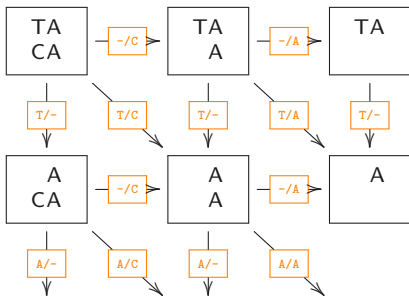
Comment calculer le meilleur alignement ?

- ▶ prenons 2 séquences de longueur n toutes les deux : alignement de longueur maximale $2n$
- ▶ exemple avec les séquences TA et CA



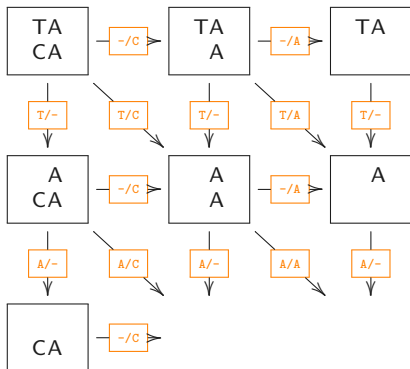
Comment calculer le meilleur alignement ?

- ▶ prenons 2 séquences de longueur n toutes les deux : alignement de longueur maximale $2n$
- ▶ exemple avec les séquences TA et CA



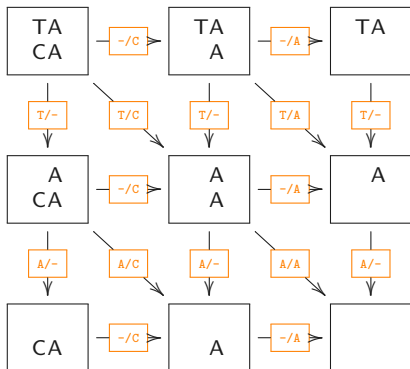
Comment calculer le meilleur alignement ?

- ▶ prenons 2 séquences de longueur n toutes les deux : alignement de longueur maximale $2n$
- ▶ exemple avec les séquences TA et CA



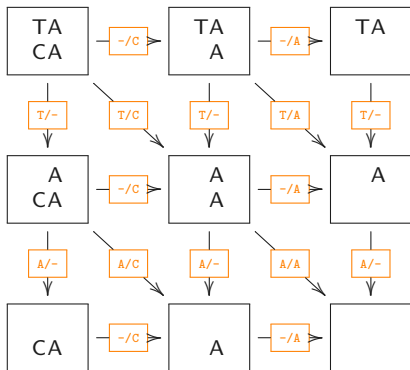
Comment calculer le meilleur alignement ?

- ▶ prenons 2 séquences de longueur n toutes les deux : alignement de longueur maximale $2n$
- ▶ exemple avec les séquences TA et CA



Comment calculer le meilleur alignement ?

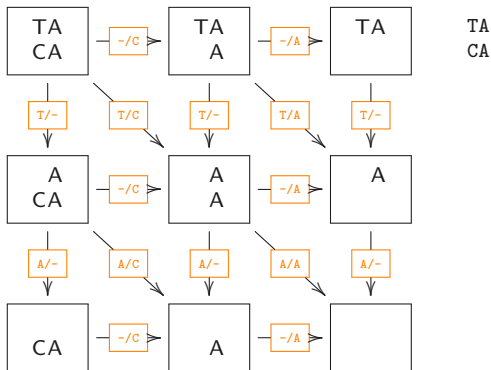
- ▶ prenons 2 séquences de longueur n toutes les deux : alignement de longueur maximale $2n$
- ▶ exemple avec les séquences TA et CA



- ▶ une suite de couples de lettres en suivant les flèches donne un alignement

Comment calculer le meilleur alignement ?

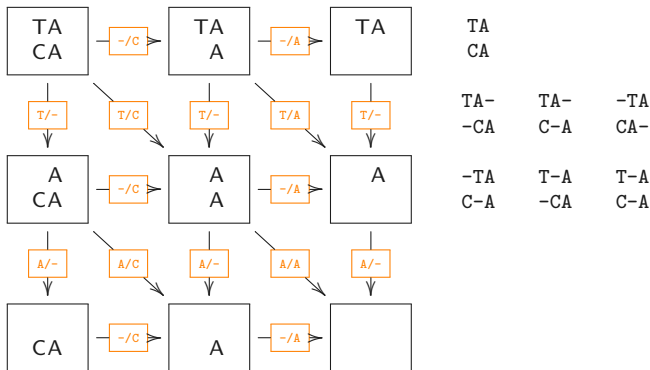
- ▶ prenons 2 séquences de longueur n toutes les deux : alignement de longueur maximale $2n$
- ▶ exemple avec les séquences TA et CA



- ▶ une suite de couples de lettres en suivant les flèches donne un alignement

Comment calculer le meilleur alignement ?

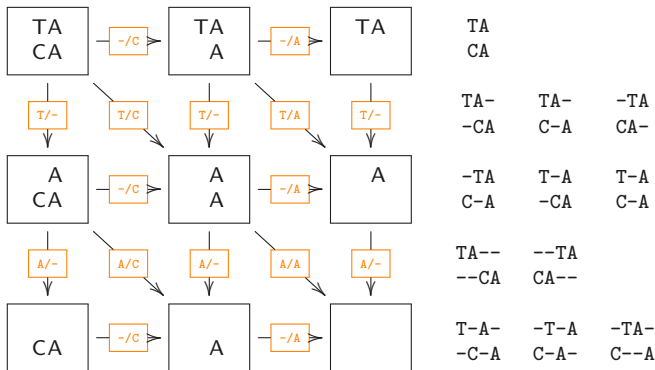
- ▶ prenons 2 séquences de longueur n toutes les deux : alignement de longueur maximale $2n$
- ▶ exemple avec les séquences TA et CA



- ▶ une suite de couples de lettres en suivant les flèches donne un alignement

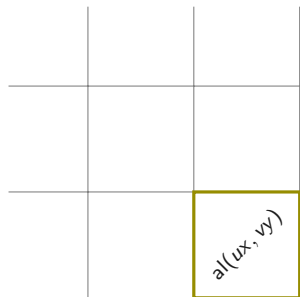
Comment calculer le meilleur alignement ?

- ▶ prenons 2 séquences de longueur n toutes les deux : alignement de longueur maximale $2n$
- ▶ exemple avec les séquences TA et CA



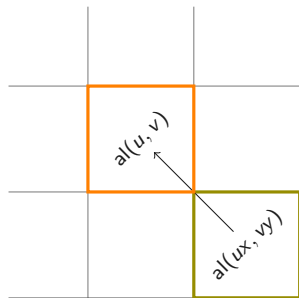
- ▶ une suite de couples de lettres en suivant les flèches donne un alignement

Une expression récursive du meilleur score



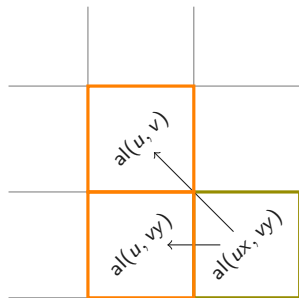
$$a(u_x, v_y) = \max \left\{ \right.$$

Une expression récursive du meilleur score



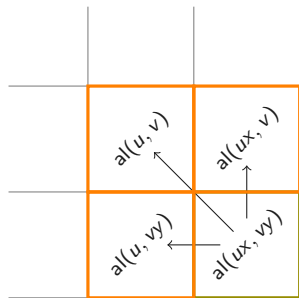
$$al(ux, vy) = \max \left\{ \begin{array}{l} al(u, v) + \text{score}(x, y) \end{array} \right.$$

Une expression réursive du meilleur score



$$al(ux, vy) = \max \begin{cases} al(u, v) + \text{score}(x, y) \\ al(u, vy) + \text{score}(x, -) \end{cases}$$

Une expression récursive du meilleur score



$$al(ux, vy) = \max \begin{cases} al(u, v) + \text{score}(x, y) \\ al(u, vy) + \text{score}(x, -) \\ al(ux, v) + \text{score}(-, y) \end{cases}$$

Programmation dynamique

Needleman-Wunsch

match = 1 mismatch = -1 gap = -1

		G	C	A	T	G	C	U	
		0	-1	-2	-3	-4	-5	-6	-7
G		-1	1	0	-1	-2	-3	-4	-5
A		-2	0	0	1	0	-1	-2	-3
T		-3	-1	-1	0	2	1	0	-1
T		-4	-2	-2	-1	1	1	0	-1
A		-5	-3	-3	-1	0	0	0	-1
C		-6	-4	-2	-2	-1	-1	1	0
A		-7	-5	-3	-1	-2	-2	0	0

GCATG-CU

| | | |
G-ATTACA

GCA-TGCU

| | | |
G-ATTACA

GCA-TGCU

| | | |
G-ATTACA

source : Wikipedia

Difficulté du problème

- ▶ nombre max d'alignements (séqu. de lg n)

$$\binom{2n}{n} = \frac{(2n)!}{(n!)^2} \sim \frac{2^{2n}}{\sqrt{2\pi n}}$$

pour deux séquences de longueur 100 : $2 \cdot 10^{57}$ alignements

- ▶ grâce à la représentation en tableau : complexité en temps et en espace $\mathcal{O}(n^2)$
(proportionnel au produit de la longueur des séquences)

pour deux séquences de longueur 100 : 10000 opérations

Un modèle plus fin pour les indels

- ▶ ne plus compter une succession d'indels comme plusieurs évènements mais comme un seul : un **gap**
- ▶ changement de la fonction de score :
 - ▶ coût associé à l'ouverture de gap
 - ▶ coût associé à l'extention de gap
- ▶ exemple (ouv = -2, ext = -1) :

A	T	C	G	G	C	A	T	G	C	C	G
A			G	G	A	A	T	G	C	-	G
1	-2	-1	1	1	-1	1	1	1	1	-2	1 = 2

Conséquence sur l'algorithme

La récurrence s'écrit alors :

$$al(n, m) = \max \left\{ \right.$$

Conséquence sur l'algorithme

La récurrence s'écrit alors :

$$al(n, m) = \max \left\{ \begin{array}{l} al(n-1, m-1) + \text{score}(u[n], v[m]) \\ \end{array} \right.$$

Les conditions initiales (première ligne, première colonne) sont :

$$al(0, i) = al(i, 0) = \text{gap}(i)$$

Conséquence sur l'algorithme

La récurrence s'écrit alors :

$$\text{al}(n, m) = \max \left\{ \begin{array}{l} \text{al}(n-1, m-1) + \text{score}(u[n], v[m]) \\ \max_i(\text{al}(n-i, m) + \text{gap}(i)) \end{array} \right.$$

Les conditions initiales (première ligne, première colonne) sont :

$$\text{al}(0, i) = \text{al}(i, 0) = \text{gap}(i)$$

Inconvénient : la complexité de l'algorithme augmente (devient en $O(n^3)$)
Mais on peut faire mieux ...

Conséquence sur l'algorithme

La récurrence s'écrit alors :

$$\text{al}(n, m) = \max \begin{cases} \text{al}(n-1, m-1) + \text{score}(u[n], v[m]) \\ \max_i(\text{al}(n-i, m) + \text{gap}(i)) \\ \max_i(\text{al}(n, m-i) + \text{gap}(i)) \end{cases}$$

Les conditions initiales (première ligne, première colonne) sont :

$$\text{al}(0, i) = \text{al}(i, 0) = \text{gap}(i)$$

Inconvénient : la complexité de l'algorithme augmente (devient en $O(n^3)$)
Mais on peut faire mieux ...

Activités

- A Calculer le score d'un alignement donné
- B Modifier un alignement à la main pour obtenir un meilleur score
- C Produire deux jeux de score « pertinents » tels que ce n'est pas le même meilleur alignement produit
- D ★ Programmer le calcul du meilleur alignement
- E ★★ Produire un meilleur alignement
- F ★★★ Produire les meilleurs alignements

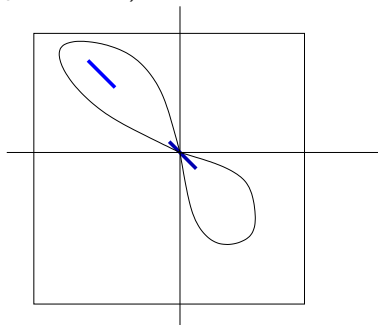
Histoire des algorithmes d'alignement

- ▶ Needleman & Wunsch, 1970 pour identifier une ressemblance globale de deux séquences
- ▶ Smith & Waterman, 1981 pour identifier une ressemblance locale de deux séquences (il suffit d'ajouter max avec 0 à l'éq. de rec.)
- ▶ variation pour ne pas compter les gaps en début/fin
- ▶ algorithme du K-band pour limiter la recherche aux alignements avec un petit nombre d'erreurs
- ▶ recherche des alignements co-optimaux
- ▶ recherche des n meilleurs alignements locaux

Et pour aller plus loin ?

- ▶ algorithmes inadaptés pour traiter de grands volumes de données
- ▶ mise en place d'heuristiques (telle que BLAST)

idée : indexer les séquences pour connaître toutes les positions de mots courts (de longueur 11 par exemple) qui formeront des candidats potentiels pour un alignement



Alignement multiple

- ▶ extension de l'alignement 2 à 2

entrée k séquences

```
C A T G C G A G T A G T A G
C A T G G T A G T A G
C C T G G A G T A C G T A G
C A T G A G C G T A G
```

sortie un tableau contenant les k séquences, avec des indels

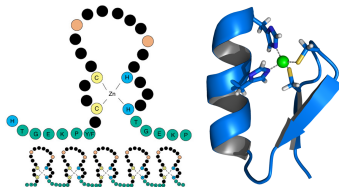
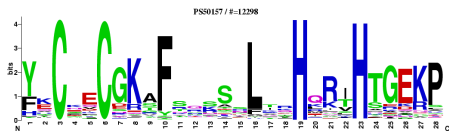
```
C A T G C G A G T A - G T A G
C A T G - - - G T A - G T A G
C C T G - G A G T A C G T A G
C A T G - - A G - - C G T A G
```

- ▶ algorithmiquement, on peut appliquer la même méthode
- ▶ biologiquement, à quelle conservation **syntaxique** correspond la famille de séquences portant la même fonction ?

Motif doigt de zinc (C2H2-type)

TYY1_HUMAN YVCPFDG**C**NKKFAQSTNLKSHILT--H
 YKQ8_CAEEL YK**C**T--V**C**RKDISSESLE**R**THMFKQH
 BASO_HUMAN FQ**C**D--I**C**KKTFKNACSVKI**H**HKN-MH
 ZG2-9_XENL FV**C**T--V**C**GKTYKYKHGLN**T**HLHS--H
 P43_XENBO LK**C**SVPG**C**KRSFRKKRALRI**H**VSE--H
 IKAR_MOUSE FEC**N**--M**C**GYSQDRYEFSS**H**ITRGEH
 TRA1_CAEEL YK**C**E**F**AD**C**E**K**A**F**S**N**AS**D**R**A**K**H**Q**N**R-T**H**
 ZN10_HUMAN YK**C****N**--Q**C**G**I**IFSQNSPFIV**H**Q**I**A--H
 XFIN_XENLA F**R****C****S**--E**C**SR**S**F**T**H**N**SD**L**T**A**H**M**R**K**--H
 TF3A_BUFAM **C**K**C**E**T**E**N****C**N**L**A**F**T**T**A**S**N**M**R**L****H**F**K**R-A**H**
 ZG58_XENLA FV**C**T--E**C**N**L**S**F**A**G**L**A**N**L**R**S**H**Q**H**L**--H
 P43_XENBO Y**R****C**S**Y**E**D****C**Q**T**V**S**P**T**W**T**A**L**Q**T**H**L**K**K**--H
 TSH_DROME F**R****C**V--W**C**K**Q**S**F**P**T**L**E**A**L**T**H**M**K**D**S**K**H**
 ZN76_HUMAN F**R****C**G**Y**K**G****C**G**R**L**Y**T**T**A**H**H**L**K**V**H**E**R**A**--H
 TF3A_BUFAM Y**R****C**P**R**E**N****C**D**R**T**Y**T**T**K**F**N**L**K**S**H**I**L**T**-**F**H
 SUHW_DROAN Y**A****C****K**--I**C**G**K**D**F**T**R**S**Y**H**L**K**R**H**Q**K**Y**S**S****C**
 ZN76_HUMAN Y**T****C**P**E**P**H****C**G**R**G**F**T**S**A**T**N**Y**K**N**H**V**R**I**--H
 SRYC_DROME F**K****C****N**--Y**C**P**R**D**F**T**N**F**N**W**L**K**H**T**R**R-R**H**
 EVI1_HUMAN Y**R****C****K**--Y**C**D**R**S**F**S**I**S**S**N**L**Q**R**H**V**R**N**-I**H**

modélisation : motif **Prosite**



C-x(2,4)-**C**-x(3)-[LIVMFYWC]-x(8)-**H**-x(3,5)-**H**

Calculer un alignement multiple

avec un arbre guide

① TACCATGA ② TACCATA ③ GACGACCA ④ GACCATCTCA

Calculer un alignement multiple

avec un arbre guide

① TACCATGA ② TACCATA ③ GACGACCA ④ GACCATCTCA

- calcul des meilleurs alignements 2 à 2 :
scores (Match = 1, Mismatch = -1, Indel = -1)

	①	②	③	④
①	.	6	0	2
②	.	.	0	3
③	.	.	.	4
④

Calculer un alignement multiple

avec un arbre guide

① TACCATGA ② TACCATA ③ GACGACCA ④ GACCATCTCA

- calcul des meilleurs alignements 2 à 2 :
scores (Match = 1, Mismatch = -1, Indel = -1)
- construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④
①	.	6	0	2
②	.	.	0	3
③	.	.	.	4
④

Calculer un alignement multiple

avec un arbre guide

① TACCATGA ② TACCATA ③ GACGACCA ④ GACCATCTCA

- calcul des meilleurs alignements 2 à 2 :
scores (Match = 1, Mismatch = -1, Indel = -1)
- construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④
①	.	6	0	2
②	.	.	0	3
③	.	.	.	4
④

Calculer un alignement multiple

avec un arbre guide

① TACCATGA ② TACCATA ③ GACGACCA ④ GACCATCTCA

- calcul des meilleurs alignements 2 à 2 :
scores (Match = 1, Mismatch = -1, Indel = -1)
- construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④
①	.	6	0	2
②	.	.	0	3
③	.	.	.	4
④



Calculer un alignement multiple

avec un arbre guide

① TACCATGA ② TACCATA ③ GACGACCA ④ GACCATCTCA

- calcul des meilleurs alignements 2 à 2 :
scores (Match = 1, Mismatch = -1, Indel = -1)
- construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④
①	.	6	0	2
②	.	.	0	3
③	.	.	.	4
④



① TACCATGA
② TACCAT-A

Calculer un alignement multiple

avec un arbre guide

① TACCATGA ② TACCATA ③ GACGACCA ④ GACCATCTCA

1. calcul des meilleurs alignements 2 à 2 :

scores (Match = 1, Mismatch = -1, Indel = -1)

2. construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④		①②	③	④
①	.	6	0	2	①②	.	0	2.5
②	.	.	0	3	③	.	.	4
③	.	.	.	4	④	.	.	.
④				



① TACCATGA

② TACCAT-A

Calculer un alignement multiple

avec un arbre guide

① TACCATGA ② TACCATA ③ GACGACCA ④ GACCATCTCA

- calcul des meilleurs alignements 2 à 2 :
scores (Match = 1, Mismatch = -1, Indel = -1)
- construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④		①②	③	④
①	.	6	0	2	①②	.	0	2.5
②	.	.	0	3	③	.	.	4
③	.	.	.	4	④	.	.	.
④				



① TACCATGA
② TACCAT-A

Calculer un alignement multiple

avec un arbre guide

① TACCATGA ② TACCATA ③ GACGACCA ④ GACCATCTCA

- calcul des meilleurs alignements 2 à 2 :
scores (Match = 1, Mismatch = -1, Indel = -1)
- construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④		①②	③	④
①	.	6	0	2	①②	.	0	2.5
②	.	.	0	3	③	.	.	4
③	.	.	.	4	④	.	.	.
④				



① TACCATGA
② TACCAT-A

Calculer un alignement multiple

avec un arbre guide

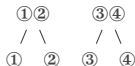
① TACCATGA ② TACCATA ③ GACGACCA ④ GACCATCTCA

1. calcul des meilleurs alignements 2 à 2 :

scores (Match = 1, Mismatch = -1, Indel = -1)

2. construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④		①②	③	④
①	.	6	0	2	①②	.	0	2.5
②	.	.	0	3	③	.	.	4
③	.	.	.	4	④	.	.	.
④				



① TACCATGA

② TACCAT-A

Calculer un alignement multiple

avec un arbre guide

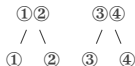
① TACCATGA ② TACCATA ③ GACGACCA ④ GACCATCTCA

1. calcul des meilleurs alignements 2 à 2 :

scores (Match = 1, Mismatch = -1, Indel = -1)

2. construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④		①②	③	④
①	.	6	0	2	①②	.	0	2.5
②	.	.	0	3	③	.	.	4
③	.	.	.	4	④	.	.	.
④				



① TACCATGA
② TACCAT-A

③ GACGA-C-CA
④ GACCATCTCA

Calculer un alignement multiple

avec un arbre guide

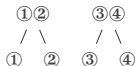
① TACCATGA ② TACCATA ③ GACGACCA ④ GACCATCTCA

1. calcul des meilleurs alignements 2 à 2 :

scores (Match = 1, Mismatch = -1, Indel = -1)

2. construction d'un arbre de clustering (et de l'alignement) :

①	②	③	④	①②	③	④	①②	③④
①	.	6	0	2	①②	.	0	2.5
②	.	.	0	3	③	.	.	4
③	.	.	.	4	④	.	.	.
④				



① TACCATGA
② TACCAT-A

③ GACGA-C-CA
④ GACCATCTCA

Calculer un alignement multiple

avec un arbre guide

① TACCATGA ② TACCATA ③ GACGACCA ④ GACCATCTCA

1. calcul des meilleurs alignements 2 à 2 :

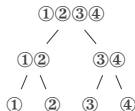
scores (Match = 1, Mismatch = -1, Indel = -1)

2. construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④
①	.	6	0	2
②	.	.	0	3
③	.	.	.	4
④

	①②	③	④
①②	.	0	2.5
③	.	.	4
④	.	.	.

	① ②	③④
①②	.	1.25
③④	.	.



① TACCATGA
② TACCAT-A

③ GACGA-C-CA
④ GACCATCTCA

Calculer un alignement multiple

avec un arbre guide

① TACCATGA ② TACCATA ③ GACGACCA ④ GACCATCTCA

1. calcul des meilleurs alignements 2 à 2 :

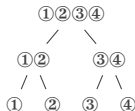
scores (Match = 1, Mismatch = -1, Indel = -1)

2. construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④
①	.	6	0	2
②	.	.	0	3
③	.	.	.	4
④

	①②	③	④
①②	.	0	2.5
③	.	.	4
④	.	.	.

	① ②	③④
①②	.	1.25
③④	.	.



① TACCATGA
② TACCAT-A

③ GACGA-C-CA
④ GACCATCTCA

① TACCAT--GA
② TACCAT--A
③ GACGA-C-CA
④ GACCATCTCA

Calculer un alignement multiple

avec un arbre guide

① TACCATGA ② TACCATA ③ GACGACCA ④ GACCATCTCA

1. calcul des meilleurs alignements 2 à 2 :

scores (Match = 1, Mismatch = -1, Indel = -1)

2. construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④
①	.	6	0	2
②	.	.	0	3
③	.	.	.	4
④

	①②	③	④
①②	.	0	2.5
③	.	.	4
④	.	.	.

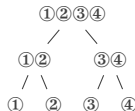
	① ②	③④
①②	.	1.25
③④	.	.



① TACCATGA
② TACCAT-A

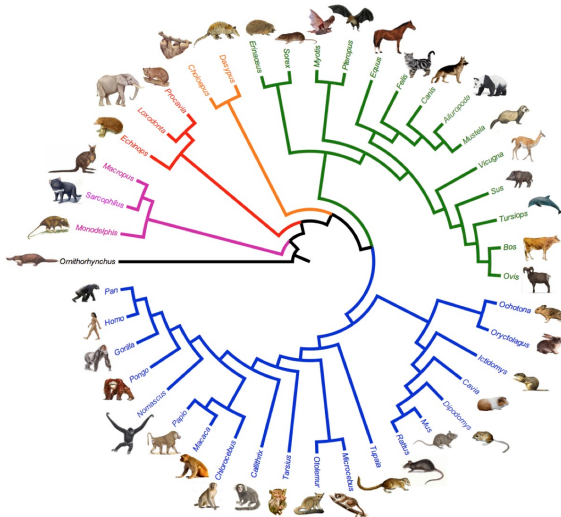


③ GACGA-C-CA
④ GACCATCTCA



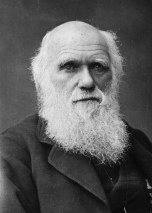
① TACCAT--GA
② TACCAT---A
③ GACGA-C-CA
④ GACCATCTCA

Reconstruction phylogénétique

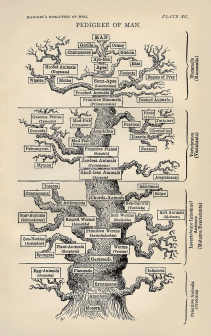


Reconstruction phylogénétique

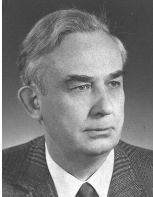
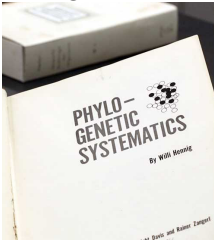
Darwin, 1859



Haeckel, 1879

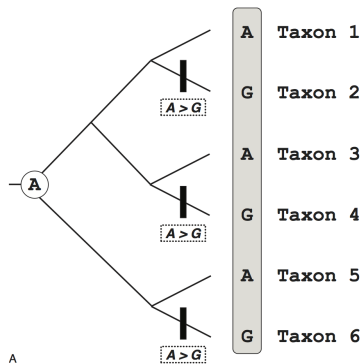


Hennig, 1950



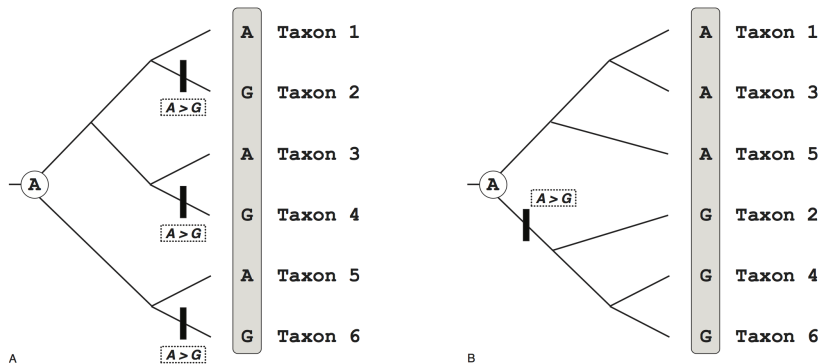
Sous l'éclairage de la phylogénie

« Rien en biologie n'a de sens, si ce n'est à la lumière de l'évolution. »
(Theodosius Dobzhansky)

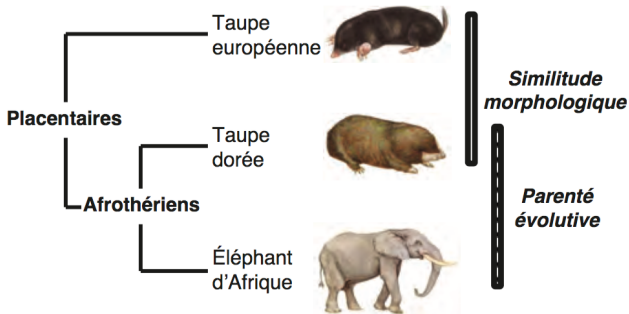


Sous l'éclairage de la phylogénie

« Rien en biologie n'a de sens, si ce n'est à la lumière de l'évolution. »
(Theodosius Dobzhansky)

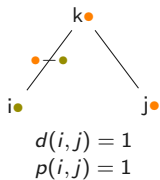


De la ressemblance à l'homologie

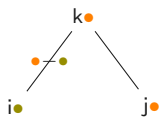


homologie = un caractère partagé par deux espèces ayant un ancêtre commun

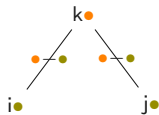
L'évolution nous joue des tours



L'évolution nous joue des tours



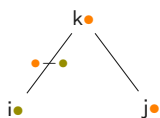
$$d(i, j) = 1$$
$$\rho(i, j) = 1$$



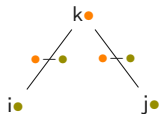
$$d(i, j) = 0$$
$$\rho(i, j) = 2$$

convergence

L'évolution nous joue des tours

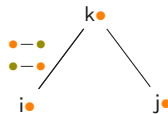


$$d(i, j) = 1$$
$$p(i, j) = 1$$



$$d(i, j) = 0$$
$$p(i, j) = 2$$

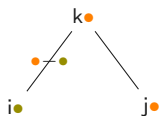
convergence



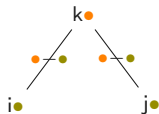
$$d(i, j) = 0$$
$$p(i, j) = 2$$

réversion

L'évolution nous joue des tours

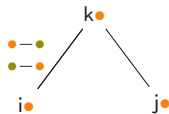


$$d(i, j) = 1$$
$$p(i, j) = 1$$



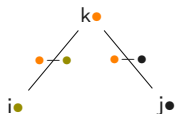
$$d(i, j) = 0$$
$$p(i, j) = 2$$

convergence



$$d(i, j) = 0$$
$$p(i, j) = 2$$

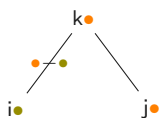
réversion



$$d(i, j) = 1$$
$$p(i, j) = 2$$

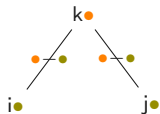
mutations multiples

L'évolution nous joue des tours



$$d(i, j) = 1$$

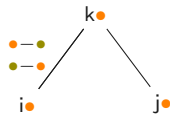
$$p(i, j) = 1$$



$$d(i, j) = 0$$

$$p(i, j) = 2$$

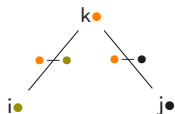
convergence



$$d(i, j) = 0$$

$$p(i, j) = 2$$

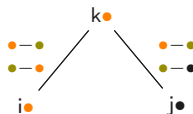
réversion



$$d(i, j) = 1$$

$$p(i, j) = 2$$

mutations multiples



$$d(i, j) = 1$$

$$p(i, j) = 4$$

convergence + réversion + mutations multiples

Méthode de distances

- ▶ l'alignement des séquences fournit un score
- ▶ sert d'estimateur de la distance évolutive entre deux séquences
- ▶ pour corriger les problèmes distance observée/distance réelle : on pose un modèle d'évolution permettant de corriger
- ▶ on procède à une classification hiérarchique

*On construit l'arbre à partir des feuilles en regroupant progressivement les noeuds 2 à 2 pour former des **clusters**.*

Activité

- A.1 Poser un modèle simple : chaque position évolue indépendamment, des substitutions surviennent aléatoirement selon un processus de Poisson
modèle de Jukes et Cantor, probabilité de substitution s et de conservation $1 - 3s$
- A.2 Simuler l'évolution de séquences, conserver le vrai nombre de substitutions (i.e. la distance observée)
- A.3 Calculer la distance observée puis la distance corrigée
$$d = -\frac{3}{4} \ln \left(1 - \frac{4}{3}p \right)$$
- A.4 En faire une représentation graphique

Un exemple idéal

matrice de distances :

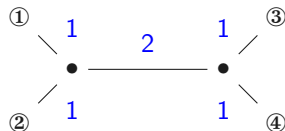
	①	②	③	④
①	0	...		
②	2	0	...	
③	4	4	0	...
④	4	4	2	0

Un exemple idéal

matrice de distances :

	①	②	③	④
①	0	...		
②	2	0	...	
③	4	4	0	...
④	4	4	2	0

arbre calculé :

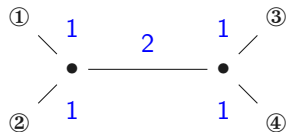


Un exemple idéal

matrice de distances :

	①	②	③	④
①	0	...		
②	2	0	...	
③	4	4	0	...
④	4	4	2	0

arbre calculé :

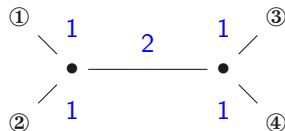


Un exemple idéal

matrice de distances :

	①	②	③	④
①	0	...		
②	2	0	...	
③	4	4	0	...
④	4	4	2	0

arbre calculé :



arbre vrai :



Un exemple plus réaliste

matrice de distances :

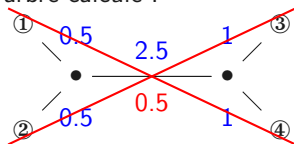
	①	②	③	④
①	0	...		
②	1	0	...	
③	4	3	0	...
④	2	3	2	0

Un exemple plus réaliste

matrice de distances :

	①	②	③	④
①	0	...		
②	1	0	...	
③	4	3	0	...
④	2	3	2	0

arbre calculé :

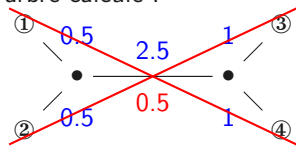


Un exemple plus réaliste

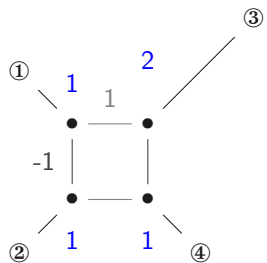
matrice de distances :

	①	②	③	④
①	0	...		
②	1	0	...	
③	4	3	0	...
④	2	3	2	0

arbre calculé :



arbre représentant la matrice :

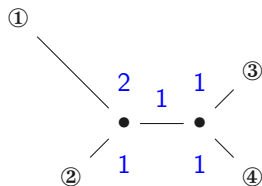


Un exemple plus réaliste

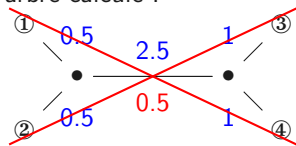
matrice de distances :

	①	②	③	④
①	0	...		
②	3	0	...	
③	4	3	0	...
④	4	3	2	0

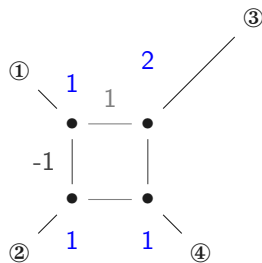
arbre vrai :



arbre calculé :



arbre représentant la matrice :



Neighbor-Joining

Satou et Nei, 1987

- ▶ autorise un taux de mutation différent sur les différentes branches,
- ▶ à partir des données initiales, calcule une matrice qui donne un arbre en étoile basé sur la divergence des taxons,
- ▶ la topologie de l'arbre est obtenu à partir de cette nouvelle matrice de distances,
- ▶ les longueurs des branches sont corrigées avec la divergence

Etapas du Neighbor-Joining

1. calcul des divergences

$$r_i = \sum_k d_{ik}$$

2. calcul de la nouvelle matrice (représentation en étoile)

$$M_{ij} = d_{ij} - \frac{r_i + r_j}{n-2}$$

3. choix des voisins les plus proches a et b , fusion en U ,
calcul des branches de l'arbre

$$l_{aU} = \frac{d_{ab}}{2} + \frac{r_a - r_b}{2n-4}, \quad l_{bU} = d_{ab} - l_{aU}$$

4. calcul des nouvelles distances pour les autres noeuds k

$$d_{kU} = \frac{d_{ak} + d_{bk} - d_{ab}}{2}$$

5. recommencer à partir de cette nouvelle matrice

Neighbor-Joining

Example

d_{ij}	A	B	C	D	E
B	5				
C	4	7			
D	7	10	7		
E	6	9	6	5	
F	8	11	8	9	8

Neighbor-Joining

Example

d_{ij}	A	B	C	D	E
B	5				
C	4	7			
D	7	10	7		
E	6	9	6	5	
F	8	11	8	9	8

1.

$$r_i = \sum_k d_{ik}$$

i	A	B	C	D	E	F
r_i	30	42	32	38	34	44

Neighbor-Joining

Example

d_{ij}	A	B	C	D	E
B	5				
C	4	7			
D	7	10	7		
E	6	9	6	5	
F	8	11	8	9	8

1.

$$r_i = \sum_k d_{ik}$$

r_i	A	B	C	D	E	F
	30	42	32	38	34	44

2.

$$M_{ij} = d_{ij} - \frac{r_i + r_j}{n - 2}$$

M_{ij}	A	B	C	D	E	F
B	-13					
C	-11.5	-11.5				
D	-10	-10	-10.5			
E	-10	-10	-10.5	-13		
F	-10.5	-10	-11	-11.5	-11.5	

Neighbor-Joining

Example

d_{ij}	A	B	C	D	E
B	5				
C	4	7			
D	7	10	7		
E	6	9	6	5	
F	8	11	8	9	8

1.

$$r_i = \sum_k d_{ik}$$

r_i	A	B	C	D	E	F
	30	42	32	38	34	44

2.

$$M_{ij} = d_{ij} - \frac{r_i + r_j}{n - 2}$$

M_{ij}	A	B	C	D	E	F
B	-13					
C	-11.5	-11.5				
D	-10	-10	-10.5			
E	-10	-10	-10.5	-13		
F	-10.5	-10	-11	-11.5	-11.5	

Neighbor-Joining

Example

d_{ij}	A	B	C	D	E
B	5				
C	4	7			
D	7	10	7		
E	6	9	6	5	
F	8	11	8	9	8

1.

$$r_i = \sum_k d_{ik}$$

i	A	B	C	D	E	F
r_i	30	42	32	38	34	44

2.

$$M_{ij} = d_{ij} - \frac{r_i + r_j}{n - 2}$$

M_{ij}	A	B	C	D	E	F
B	-13					
C	-11.5	-11.5				
D	-10	-10	-10.5			
E	-10	-10	-10.5	-13		
F	-10.5	-10	-11	-11.5	-11.5	

3.

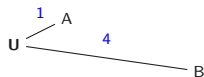
$$I_{AU} = \frac{d_{AB}}{2} + \frac{r_A - r_B}{2(n - 2)} = \frac{5}{2} + \frac{30 - 42}{2(6 - 2)} = 1$$

$$I_{BU} = d_{AB} - I_{AU} = 5 - 1 = 4$$

Neighbor-Joining

Example

d_{ij}	A	B	C	D	E
B	5				
C	4	7			
D	7	10	7		
E	6	9	6	5	
F	8	11	8	9	8



1.

$$r_i = \sum_k d_{ik}$$

i	A	B	C	D	E	F
r_i	30	42	32	38	34	44

2.

$$M_{ij} = d_{ij} - \frac{r_i + r_j}{n - 2}$$

M_{ij}	A	B	C	D	E	F
B	-13					
C	-11.5	-11.5				
D	-10	-10	-10.5			
E	-10	-10	-10.5	-13		
F	-10.5	-10	-11	-11.5	-11.5	

3.

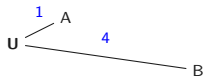
$$I_{AU} = \frac{d_{AB}}{2} + \frac{r_A - r_B}{2(n - 2)} = \frac{5}{2} + \frac{30 - 42}{2(6 - 2)} = 1$$

$$I_{BU} = d_{AB} - I_{AU} = 5 - 1 = 4$$

Neighbor-Joining

Example

d_{ij}	A	B	C	D	E
B	5				
C	4	7			
D	7	10	7		
E	6	9	6	5	
F	8	11	8	9	8



4.

$$d'_{CU} = \frac{d_{AC} + d_{BC} - d_{AB}}{2} = \frac{4 + 7 - 5}{2}$$

d'_{ij}	U	C	D	E	F
C	3				
D					
E					
F					

Neighbor-Joining

Example

d_{ij}	A	B	C	D	E
B	5				
C	4	7			
D	7	10	7		
E	6	9	6	5	
F	8	11	8	9	8



4.

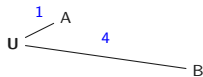
$$\forall k \quad d'_{kU} = \frac{d_{Ak} + d_{Bk} - d_{AB}}{2}$$

d'_{ij}	U	C	D	E	F
C	3				
D	6	7			
E	5	6	5		
F	7	8	9	8	

Neighbor-Joining

Example

d_{ij}	U	C	D	E	F
C	3				
D	6	7			
E	5	6	5		
F	7	8	9	8	



Neighbor-Joining

Example

d_{ij}	U	C	D	E	F
C	3				
D	6	7			
E	5	6	5		
F	7	8	9	8	



1.

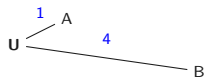
$$r_i = \sum_k d_{ik}$$

i	U	C	D	E	F
r_i	21	24	27	24	32

Neighbor-Joining

Example

d_{ij}	U	C	D	E	F
C	3				
D	6	7			
E	5	6	5		
F	7	8	9	8	



1.

$$r_i = \sum_k d_{ik}$$

i	U	C	D	E	F
r_i	21	24	27	24	32

2.

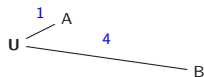
$$M_{ij} = d_{ij} - \frac{r_i + r_j}{n - 2}$$

M_{ij}	U	C	D	E	F
C	-12				
D	-10	-10			
E	-10	-10	-12		
F	-10.6	-10.6	-10.6	-10.6	

Neighbor-Joining

Example

d_{ij}	U	C	D	E	F
C	3				
D	6	7			
E	5	6	5		
F	7	8	9	8	



1.

$$r_i = \sum_k d_{ik}$$

i	U	C	D	E	F
r_i	21	24	27	24	32

2.

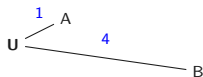
$$M_{ij} = d_{ij} - \frac{r_i + r_j}{n - 2}$$

M_{ij}	U	C	D	E	F
C	-12				
D	-10	-10			
E	-10	-10	-12		
F	-10.6	-10.6	-10.6	-10.6	

Neighbor-Joining

Example

d_{ij}	U	C	D	E	F
C	3				
D	6	7			
E	5	6	5		
F	7	8	9	8	



1.

$$r_i = \sum_k d_{ik}$$

i	U	C	D	E	F
r_i	21	24	27	24	32

2.

$$M_{ij} = d_{ij} - \frac{r_i + r_j}{n - 2}$$

M_{ij}	U	C	D	E	F
C	-12				
D	-10	-10			
E	-10	-10	-12		
F	-10.6	-10.6	-10.6	-10.6	

3.

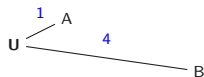
$$l_{DV} = \frac{d_{DE}}{2} + \frac{r(D) - r(E)}{2(n-2)} = \frac{5}{2} + \frac{27-24}{2(5-2)} = 3$$

$$l_{EV} = d_{DE} - l_{DV} = 5 - 3 = 2$$

Neighbor-Joining

Example

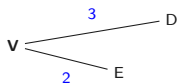
d_{ij}	U	C	D	E	F
C	3				
D	6	7			
E	5	6	5		
F	7	8	9	8	



1.

$$r_i = \sum_k d_{ik}$$

i	U	C	D	E	F
r_i	21	24	27	24	32



2.

$$M_{ij} = d_{ij} - \frac{r_i + r_j}{n - 2}$$

M_{ij}	U	C	D	E	F
C	-12				
D	-10	-10			
E	-10	-10	-12		
F	-10.6	-10.6	-10.6	-10.6	

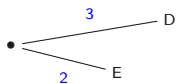
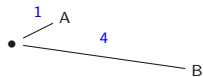
3.

$$l_{DV} = \frac{d_{DE}}{2} + \frac{r(D) - r(E)}{2(n-2)} = \frac{5}{2} + \frac{27 - 24}{2(5-2)} = 3$$

$$l_{EV} = d_{DE} - l_{DV} = 5 - 3 = 2$$

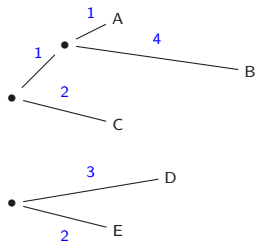
Neighbor-Joining

Example



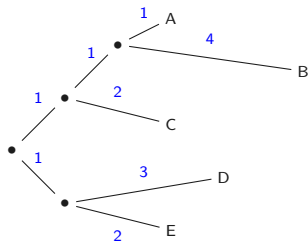
Neighbor-Joining

Example



Neighbor-Joining

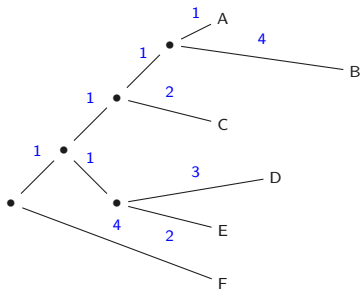
Example



Neighbor-Joining

Example

d_{ij}	A	B	C	D	E
B	5				
C	4	7			
D	7	10	7		
E	6	9	6	5	
F	8	11	8	9	8



Il existe de nombreuses autres méthodes de RP

bien plus performantes

- ▶ les méthodes probabilistes
 - ▶ maximum de vraisemblance : quel est l'arbre le plus vraisemblable dans le modèle d'évolution

calcul de $p(\mathcal{D}|\mathcal{T})$, la probabilité des données \mathcal{D} (i.e. de l'alignement multiple) sachant l'arbre \mathcal{T}

$$p(\mathcal{D}|\mathcal{T}) = \prod_{i=1}^N p(\mathcal{D}_i|\mathcal{T})$$

- ▶ inférence bayésienne :

calcul de $p(\mathcal{T}|\mathcal{D})$, la probabilité que l'arbre \mathcal{T} soit le vrai sachant les données \mathcal{D} (i.e. de l'alignement multiple)

$$p(\mathcal{T}|\mathcal{D}) = \frac{p(\mathcal{D}|\mathcal{T})p(\mathcal{T})}{\sum_{\mathcal{T}} p(\mathcal{D}|\mathcal{T})p(\mathcal{T})}$$

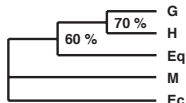
Puis-je croire en mon arbre ?

La méthode de bootstrap

Matrice originale

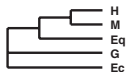
	0	1	4	6	7	8	9
Glaucomyx	C	G	T	A	T	A	A
Mus	C	A	C	T	G	C	G
Homo	C	A	T	A	A	C	A
Equus	T	A	C	A	A	C	A
Echinops	T	G	C	A	G	T	G

Consensus de bootstrap



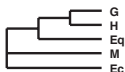
Réplication N°1

000114669
799015556
Glaucomyx CGGTTATTA
Mus CAACGCGG
Homo CAATAACCA
Equus TAACAACCA
Echinops TGGCAGTTG



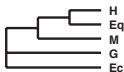
N°2

001144477
790055555
Glaucomyx CGTTAAAAA
Mus CACCGGGG
Homo CATTAAAAA
Equus TACCAAAAA
Echinops TGCCGGGGG



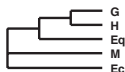
N°3

000118999
999116666
Glaucomyx GGGTTAAAA
Mus AAATTAGGG
Homo AAAAAACAA
Equus AAAAGAAA
Echinops GGAAGGGG



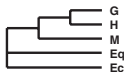
N°4

111114448
000115556
Glaucomyx TTTTAAAAA
Mus CCCTTGGGA
Homo TTTAAAAAC
Equus CCCAAAAAG
Echinops CCCAAGGGG



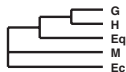
N°5

001488889
790566666
Glaucomyx CGTAAAAAA
Mus CACGAAAAG
Homo CATACCCCA
Equus TACAGGGGA
Echinops TGCGGGGGG



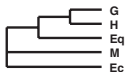
Réplication N°6

014448999
705556666
Glaucomyx CTAAAAAAA
Mus CCGGAGGG
Homo CTAACAATA
Equus TCAAAGAAA
Echinops TCGGGGGG



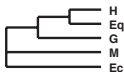
N°7

001114699
790015566
Glaucomyx CGTTTAAAA
Mus CACCTGGGG
Homo CATTAAAAA
Equus TACCAAAAA
Echinops TGCCAGGGG



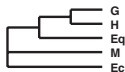
N°8

014489999
715566666
Glaucomyx CTAAAAAAA
Mus CTGGAGGG
Homo CAAACAAAA
Equus TAAAGAAAA
Echinops TAGGGGGG



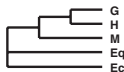
N°9

011167899
900155666
Glaucomyx GTTTTAAAA
Mus ACCTCGAGG
Homo ATTACACAA
Equus ACCACAGAA
Echinops GCCATGGGG



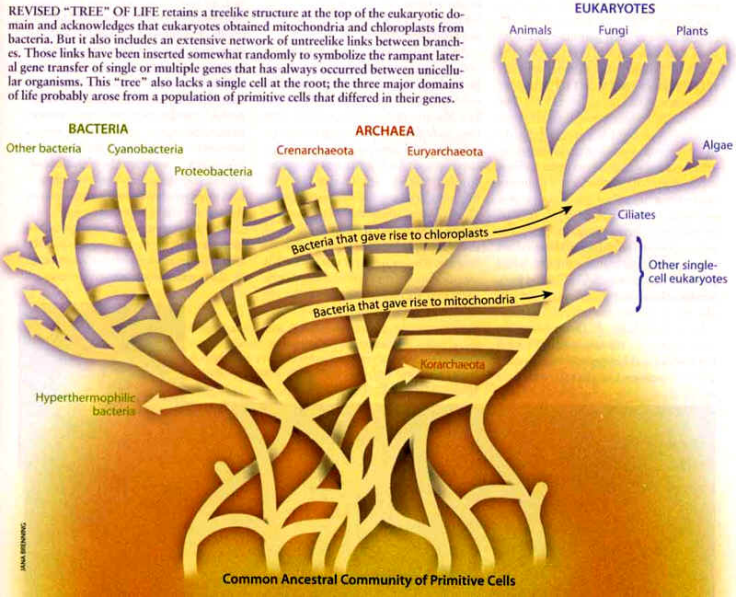
N°10

001111689
770001566
Glaucomyx CCTTTTTAA
Mus CCCCTCAG
Homo CCTTTACCA
Equus TTCCACGA
Echinops TTCCCATGG

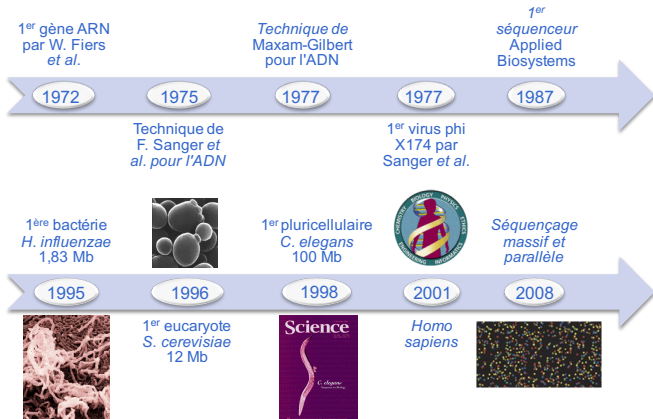


Une histoire bien plus complexe ...

REVISED "TREE" OF LIFE retains a treelike structure at the top of the eukaryotic domain and acknowledges that eukaryotes obtained mitochondria and chloroplasts from bacteria. But it also includes an extensive network of untreelike links between branches. Those links have been inserted somewhat randomly to symbolize the rampant lateral gene transfer of single or multiple genes that has always occurred between unicellular organisms. This "tree" also lacks a single cell at the root; the three major domains of life probably arose from a population of primitive cells that differed in their genes.



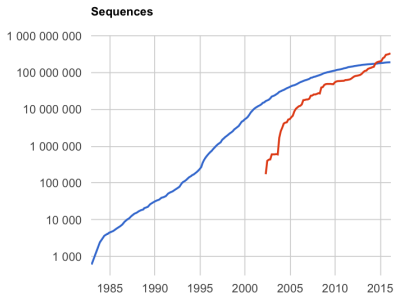
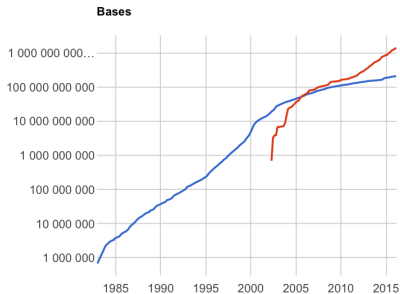
Petite histoire du séquençage



Genbank (release 212, février 2016)

GenBank		WGS	
Bases	Sequences	Bases	Sequences
207,018,196,067	190,250,235	1,399,865,495,608	333,012,760

Croissance de la taille des banques de données



WGS / Sequences

Banques de données

- ▶ EMBL (European Bioinformatics Institute, depuis 1992)
- ▶ GenBank (National Center for Biotechnology Information, depuis 1988)
- ▶ Ensembl (génomes eukaryotes)
- ▶ PubMed (banques de publications scientifiques)



EMBL



e!Ensembl

PubMed

Remerciements

Le matériel est issu principalement des cours de bio-informatique donnés à l'Université Lille 1 (auteurs : Sylvain Legrand, Laurent Noé, Hélène Touzet, Maude Pupin, Jean-Stéphane Varré).

Remerciements à Samuel Blanquart pour la nouvelle version du chapitre de livre « Phylogénie Moléculaire » du livre « Biologie Evolutive ».