

# Formation Alignement/ Phylogénie/Motifs

Hélène Chiapello, Géraldine Pascal, Claire Hoede

INRA  
Unités GenPhyse et MIAT, Plateforme Genotoul-bioinfo

24 Février 2017

# Contexte de la formation

- Le point de départ : une demande d'une équipe de recherche du LISBP (E. Laville *et al.*)  
*“ Pouvoir caractériser fonctionnellement (fonction, substrat) des familles complexes d'enzymes de synthèse ou de dégradation de sucres” (CAZy)*
- C'est une question de recherche ;-)

# Tour de table

*Qui êtes vous ?*

*Quelle expérience avez-vous en bioinfo ?*

*Quelles attentes avez-vous pour cette formation ?*

# Objectifs de la formation

- Acquérir les **concepts de base en bioinformatique** pour
  - Réaliser des alignements de séquences multiples
  - Construire et manipuler des arbres phylogénétiques
  - Rechercher des motifs conservés dans les séquences
- Acquérir une **autonomie de pratique des outils bioinformatique** sans utiliser la ligne de commande
- Vous proposer une **démarche bioinformatique générique adaptée à votre question de recherche**

# Un point de départ



Applied and Environmental  
Microbiology



## Dividing the Large Glycoside Hydrolase Family 43 into Subfamilies: a Motivation for Detailed Enzyme Characterization

Keith Mewis,<sup>a</sup>  Nicolas Lenfant,<sup>b,c</sup> Vincent Lombard,<sup>b,c</sup> Bernard Henrissat<sup>b,c,d</sup>

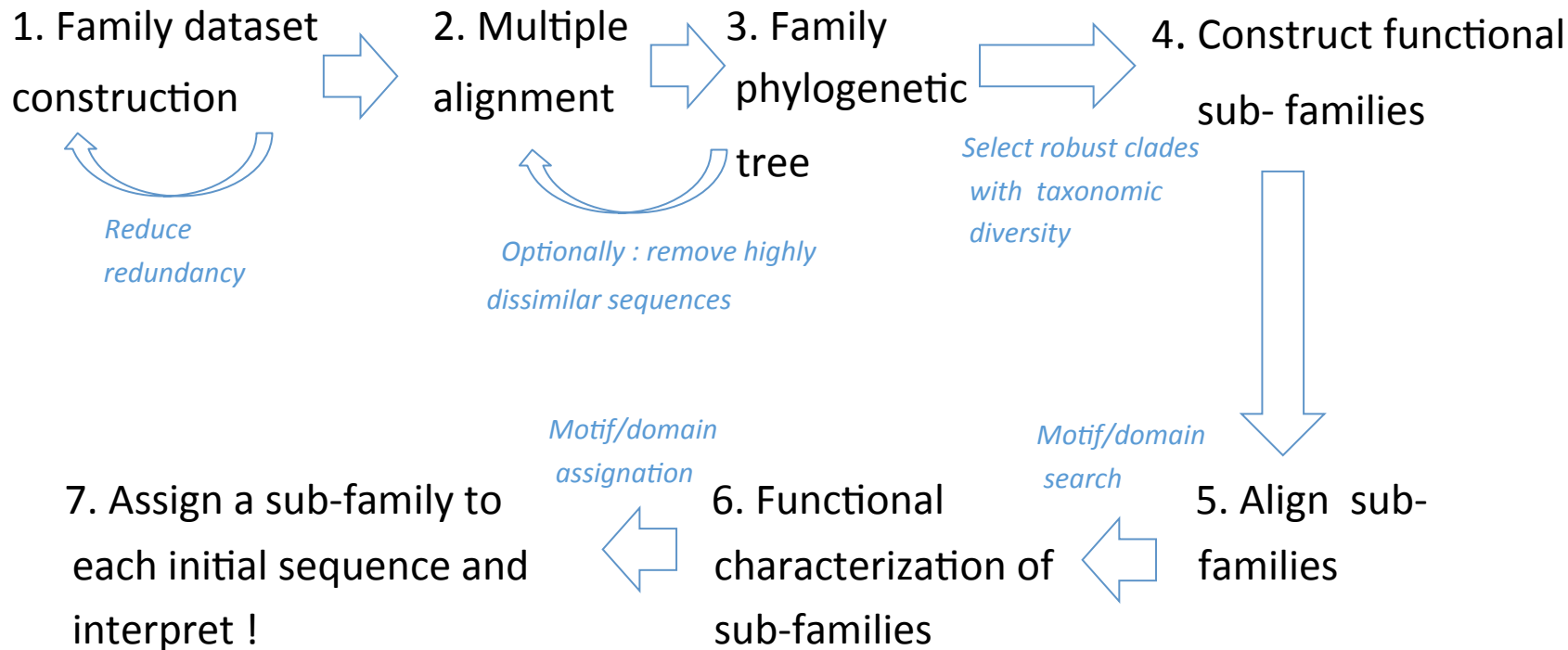
Genome Science and Technology Program, University of British Columbia, Vancouver, BC, Canada<sup>a</sup>; Architecture et Fonction des Macromolécules Biologiques, CNRS, Aix-Marseille Université, Marseille, France<sup>b</sup>; INRA, USC 1408 AFMB, Marseille, France<sup>c</sup>; Department of Biological Sciences, King Abdulaziz University, Jeddah, Saudi Arabia<sup>d</sup>

The rapid rise in DNA sequencing has led to an expansion in the number of glycoside hydrolase (GH) families. The GH43 family currently contains  $\alpha$ -L-arabinofuranosidase,  $\beta$ -D-xylosidase,  $\alpha$ -L-arabinanase, and  $\beta$ -D-galactosidase enzymes for the debranching and degradation of hemicellulose and pectin polymers. Many studies have revealed finer details about members of GH43 that necessitate the division of GH43 into subfamilies, as was done previously for the GH5 and GH13 families. The work presented here is a robust subfamily classification that assigns over 91% of all complete GH43 domains into 37 subfamilies that correlate with conserved sequence residues and results of biochemical assays and structural studies. Furthermore, cooccurrence analysis of these subfamilies and other functional modules revealed strong associations between some GH43 subfamilies and CBM6 and CBM13 domains. Cooccurrence analysis also revealed the presence of proteins containing up to three GH43 domains and belonging to different subfamilies, suggesting significant functional differences for each subfamily. Overall, the subfamily analysis suggests that the GH43 enzymes probably display a hitherto underestimated variety of subtle specificity features that are not apparent when the enzymes are assayed with simple synthetic substrates, such as pNP-glycosides.

GH43 family :

- Need to be divided into subfamilies
- A strategy proposed to identify functional differences between subfamilies

# Une démarche bioinfo générique



# Déroulé de la formation

3 parties :

- Partie 1 : Alignement de séquences - Géraldine Pascal
- Partie 2 : Construction d'arbres phylogénétiques - Hélène Chiapello
- Partie 3 : Recherche de motifs fonctionnels - Claire Hoede

Et de nombreux travaux pratiques sur un dataset de 62 séquences de CAZY GH130 (E. Laville)

# Partie 1 - la stratégie utilisée :

*Material and Methods - Mewis et al. App. And Env. Microb. 2016*

- *All complete GH43 domain sequences were taken from the CAZy database -> 62 sequences from Elisabeth*
- *To reduce redundancy and improve the processing time, sequences were clustered at 95% similarity by using CD-Hit*
- *In order to generate high-quality and relevant alignments, MAFFT was used to iteratively remove highly dissimilar sequences*
- *With these sequences, was used to generate a phylogenetic tree based on the midpoint root method -> NJ*



# Alignement de séquence, quelques points de théorie

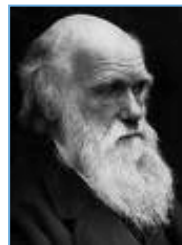
# La théorie de l'évolution

Les espèces se modifient au cours du temps et donnent naissance à de nouvelles espèces.

Selon Jean-Baptiste Lamarck (1744-1829), les **espèces évoluent** en adoptant des **caractères acquis** par les individus au **cours de leur vie**.



Charles Darwin (1809-1882) émet l'hypothèse de la **sélection du plus apte** (ou sélection naturelle) parmi des individus naturellement variant.



# La déduction par homologie, ou le «dogme central» de la bioinformatique

Si la bioinformatique «marche», c'est parce que l'évolution des gènes laisse une trace parfaitement visible lorsque l'on compare leur séquence

Évolution des gènes = mutations, insertions, délétion.

Les gènes des organismes modernes sont issus de remaniement de gènes ancestraux

On peut donc déduire la fonction de la plupart des gènes par comparaison avec les gènes «homologues» d'autres espèces.

(homologue = qui a un ancêtre commun)

## La déduction par homologie, ou le «dogme central» de la bioinformatique

Les **régions fonctionnelles** des gènes (sites catalytique, de fixation, etc.) **sont soumises à sélection**. Elles sont relativement **préservées par l'évolution** car des **mutations trop radicales sont désavantageuses**.

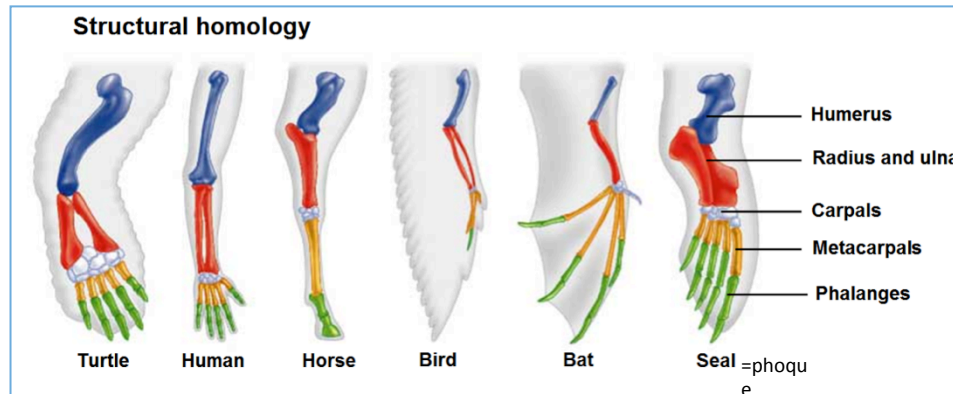
Les **régions non fonctionnelles** subissent peu de pressions de sélection et **divergent rapidement** au fur et à mesure que s'accumulent les mutations.

# L'homologie de séquence

En bioinformatique aussi **Homologie = parenté = ancêtre commun**

L'aile de l'aigle **est homologue** à l'aile du perroquet (oiseaux) et à la patte de la tortue (reptile)

L'aile de la chauve souris **est homologue** à la patte du cheval et au bras de l'homme (chiroptère, ongulé et primate - mammifères).



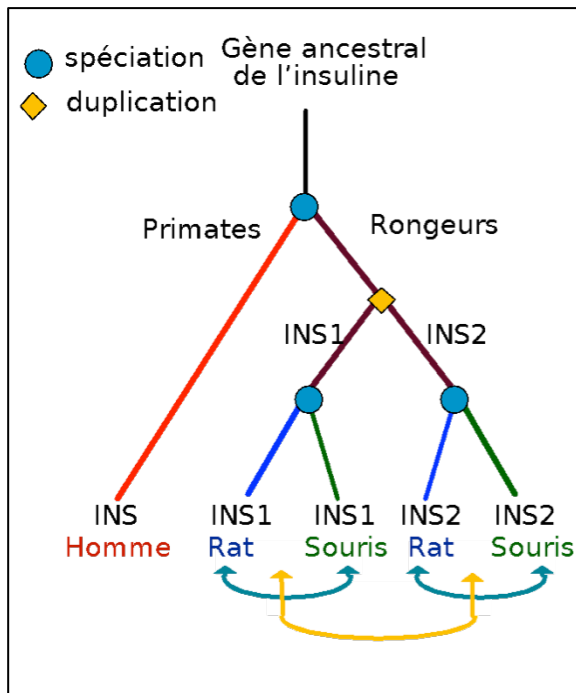
# L'homologie de séquence

On est homologue ou on ne l'est pas.

Donc **on ne dit pas**: "très homologue", "faible homologie", «22% d'homologie», etc.

Pour une notion **quantitative**, on parle de **similitude** ("très similaire", etc.) ou **d'identité** (28% d'identité)

# Orthologie et paralogie



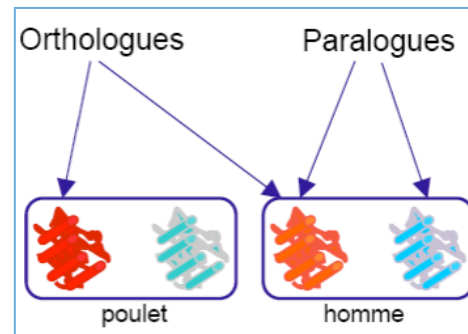
Homologie : 2 gènes sont homologues s'ils ont un ancêtre commun

Orthologie: 2 gènes sont orthologues s'ils ont divergé à la suite d'un événement de spéciation  
 ex: INS1 et INS2

Paralogie: 2 gènes sont paralogues s'ils ont divergé après un événement de duplication  
 ex: INS1 de rat et INS1 de souris

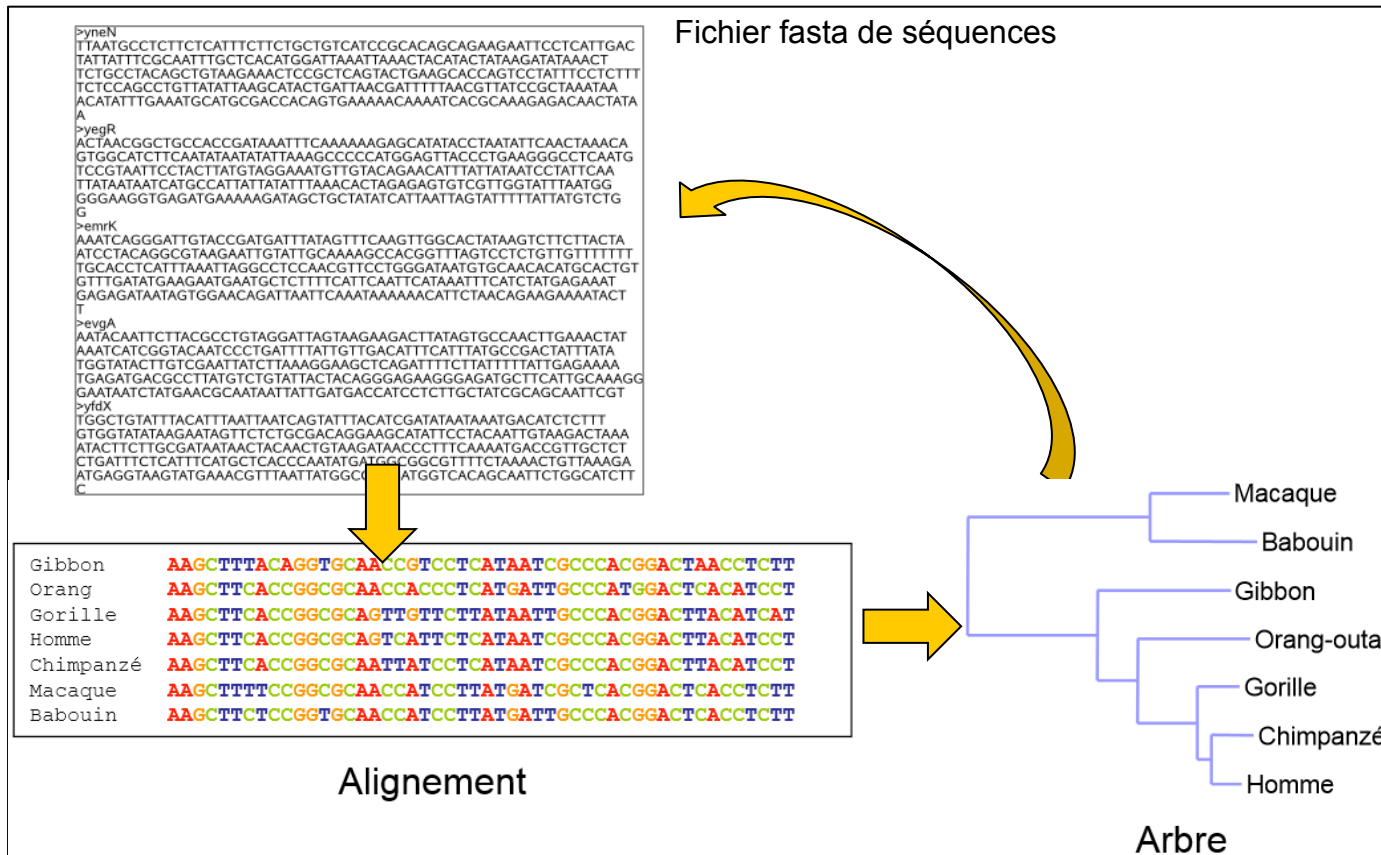
# Fonction et homologie

- **Homologie n'implique pas même fonction**: par exemple l'aile de l'oiseau et le bras humain n'ont pas la même fonction
- Des **orthologues rapprochés** (p. ex. homme/souris) ont le **plus souvent la même fonction** dans l'organisme.
- Des **orthologues distants** (p. ex. homme/mouche) ont **plus rarement le même rôle phénotypique**, mais peuvent exercer le même rôle dans une voie donnée.
- Les **paralogues acquièrent rapidement des fonctions différentes**





# Les différentes étapes de la reconstruction d'arbre phylogénétique



# Les différentes étapes de la reconstruction d'arbre phylogénétique

## Point de départ :

Un ensemble de séquences *homologues* alignées.

Avoir sa séquence et rechercher des séquences *similaires* dans les banques de données par blast par exemple

*OU*

Avoir *sa propre liste* de séquences déduite d'expériences biologiques

Séquences protéiques ou nucléiques

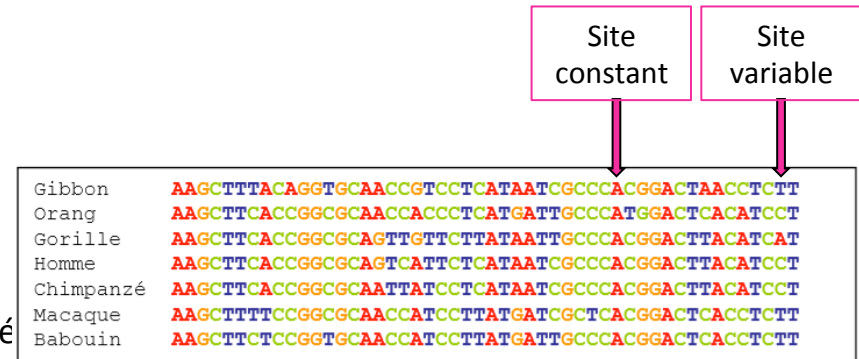
Format standard fasta ou phylip

## Alignement multiple

Chaque position dans l'alignement constitue un site.

## Résultat obtenu :

Un arbre décrivant les relations évolutives entre les séquences  
*i.e.* un arbre phylogénétique.

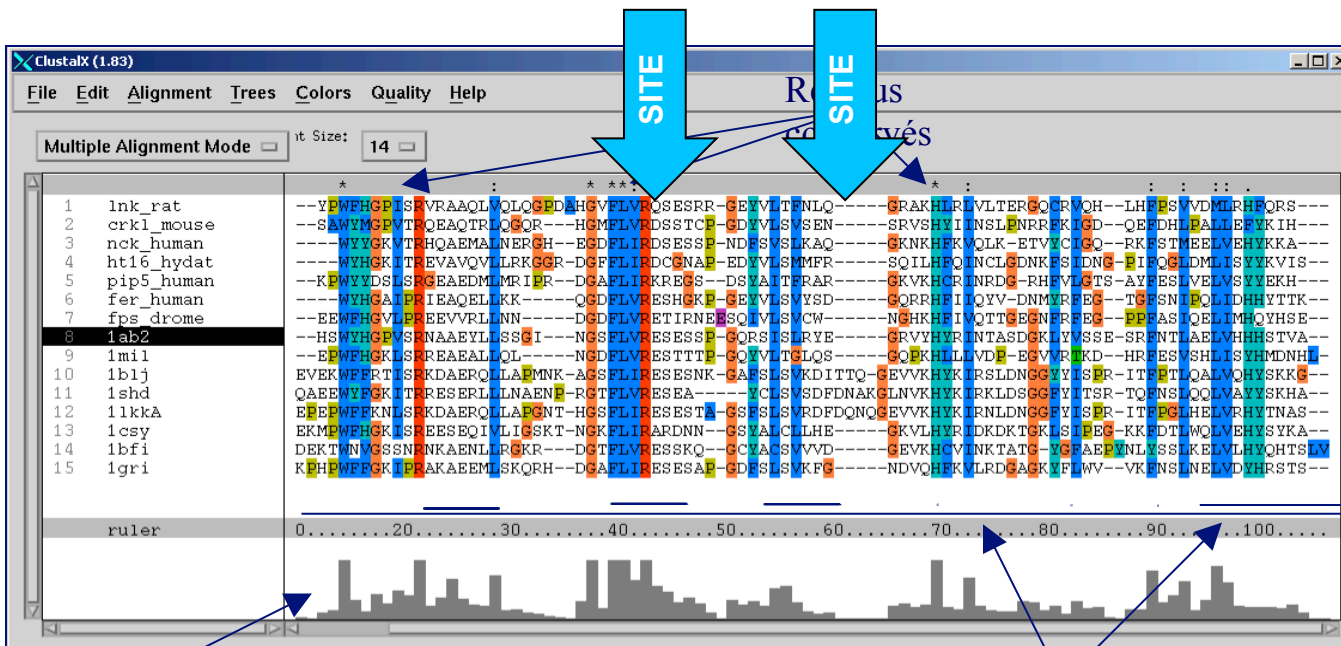


# L'alignement multiple



# Qu'est-ce qu'un alignement multiple ?

## Vue d'un alignement dans un éditeur d'alignement



Profil de conservation

Structure secondaire

# L'alignement multiple au format fasta

```

>danio
MKATSGPNGLRCLVLSCLFVLHIIFTRVAGNMA SP SHFGVAST----EVHRNRHPKRRNP
FRPPVESRTDDDGMRMLSLYRIAADADGRPKQHKIFGSNTIRLLQASTTEKHFPPTSSD
LQYTYTVKYELNDL-LLDKLVKASFMYLRSFMSRL-----PYICEASV
TSLQNFLEGDRI TMGFRSRWTE TDVTDHV-----SESKDGHVSFFARYWCTKPEHKR
SV----AHRKR-----PPQHHLRAPLLLLFLEENKHPVEWG-----
-KSFPPLSRPRTTR
>oryzias
MRATRGT RSLRLRTCALTCFLFLCSTCA-GAATRRNVASHFSTSKR RRS HQSAKHMGAH
HRPLTDEQKADQNLQFMLS YRSAAEPDGRPKQHRKFGSNTVRLRPTASSVHYRPTSGD
LCYTFVQYKLD TLP-SEQLVRASF IHLRTSSLTLSQT-----VEPPQCRAQI
ASF---GGKSLAVLEPHEQWTE TDITAHVSAHILQGRD INEEGHLTLTAQYWCTEPGKPD
-----VDERKRW-----NSEPHLEAPSLLLYLEEERENST--LELKDSFLDALN--
-SPTSS-----
>takifugu
-----MDDRKLSPSRRSSLR SRARLPDGGAH----
HRPLTDEQKADQNLRFMLNLYQSA AEPDGRPKQHRKFGSNTVRLKPSASSVRYLPAP--
-----TVQYHMDTLP-TEQLVRASFVHLRSSATSSSLNAT-----QGAKPPRC EARI
TSL---GQESLVTLEPHEQWMETDITAHVR---QDNQSPGKFLTLTGQYCA AEDALV
HSEEDVGLKWWTRLQERGRSGEHHLEVP SLLLYLEE QREEQR-----
-----RHR-
>tetraodon
-----
-RPLTDEQKADQNLQFMMSLYRSAAEPDGRPKQHRKFGSNTVRLKPSASSVSYLPASPD
HQYHFSVQYHLD TLP-SEQLVRASFVHLRSAP-APFNSSQG-----APPPRCR ARV
APL---GRESLVVLEPHQRWTE TDITAHVR---QRDQSPGGALTLTAQYRC TAPMAAQ
GGG---GLPRPW-----AQRGGQHLEVP SLLLYLEEERDGNVWMDLLG-----
-----PEQRRRRR
>gasterosteus
-----
-----QNLQFMLS YRSAAEPDGRPKQHRKFGSNTVRLRPSAASVRHLPASPD
HRYSFVQYNLD TVP-SEQLIRASF IHLRSAPSSSSARP-----LRPPRCRAQI
---PSLGKASLVTLEPHERWTE TDITAHVR---RGRSRGPGGHLTLTAQYWC TAWGGGF
-----PSTGGSP T SRRLFYTWRRS-----
-GRTPPQRTPRTRR

```

# Les outils les plus populaires

EMBL-EBI Services Research Training About us

## Multiple Sequence Alignment

Share Feedback

Tools > Multiple Sequence Alignment

**Multiple Sequence Alignment (MSA)** is generally the alignment of three or more biological sequences (protein or nucleic acid) of similar length. From the output, homology can be inferred and the evolutionary relationships between the sequences studied.

By contrast, [Pairwise Sequence Alignment](#) tools are used to identify regions of similarity that may indicate functional, structural and/or evolutionary relationships between two biological sequences.

**Clustal Omega** ?

New MSA tool that uses seeded guide trees and HMM profile-profile techniques to generate alignments. Suitable for medium-large alignments.

[Launch Clustal Omega](#)

**MUSCLE** ?

Accurate MSA tool, especially good with proteins. Suitable for medium alignments.

[Launch MUSCLE](#)

**ClustalW2** ?

Popular MSA tool that uses tree-based progressive alignments. Suitable for medium alignments.

[Launch ClustalW2](#)

**MView** ?

Transform a Sequence Similarity Search result into a Multiple Sequence Alignment or reformat a Multiple Sequence Alignment using the MView program.

[Launch MView](#)

**DbClustal** ?

Create a Multiple Sequence Alignment from a protein BLAST result using the DbClustal program.

[Launch DbClustal](#)

**T-Coffee** ?

Consistency-based MSA tool that attempts to mitigate the pitfalls of progressive alignment methods. Suitable for small alignments.

[Launch T-Coffee](#)

**Kalign** ?

Very fast MSA tool that concentrates on local regions. Suitable for large alignments.

[Launch Kalign](#)

**WebPRANK**

The EBI has a new phylogeny-aware multiple sequence alignment program which makes use of evolutionary information to help place insertions and deletions. Try it out at [WebPRANK](#).

**MAFFT** ?

MSA tool that uses Fast Fourier Transforms. Suitable for medium-large alignments.

[Launch MAFFT](#)

# Construire un alignement multiple



# Approches traditionnelles

- A. alignement multiple optimal
- B. alignement multiple progressif
- C. alignement multiple itératif

# Approches traditionnelles

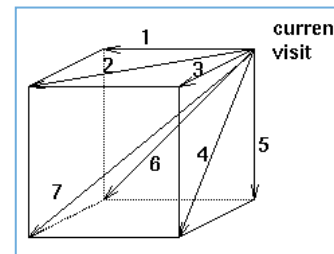
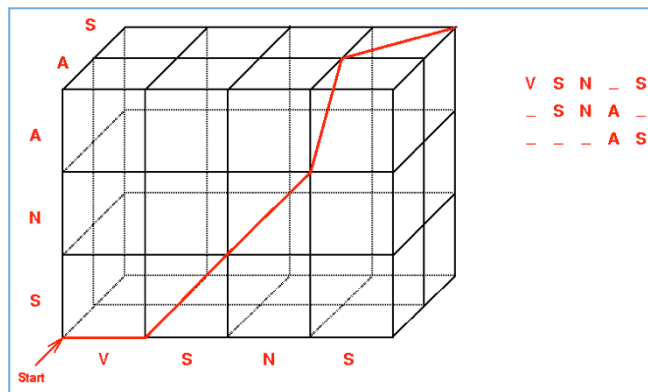
- A. alignement multiple optimal
- B. alignement multiple progressif
- C. alignement multiple itératif

# A. Alignement multiple optimal

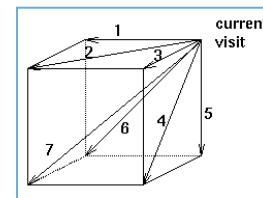
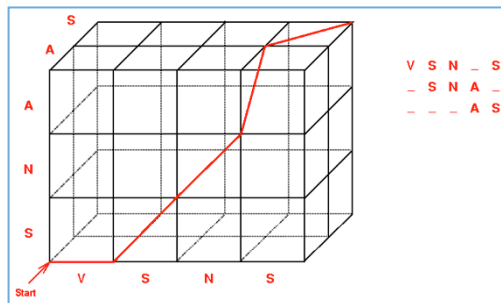
Extension directe des programmes dynamiques des alignements de séquences par paires à N dimensions (Sankoff, 1975).

Examine l'ensemble des alignements possibles afin de trouver l'alignement optimal

Exemple: alignement de 3 séquences



## A. Alignement multiple optimal



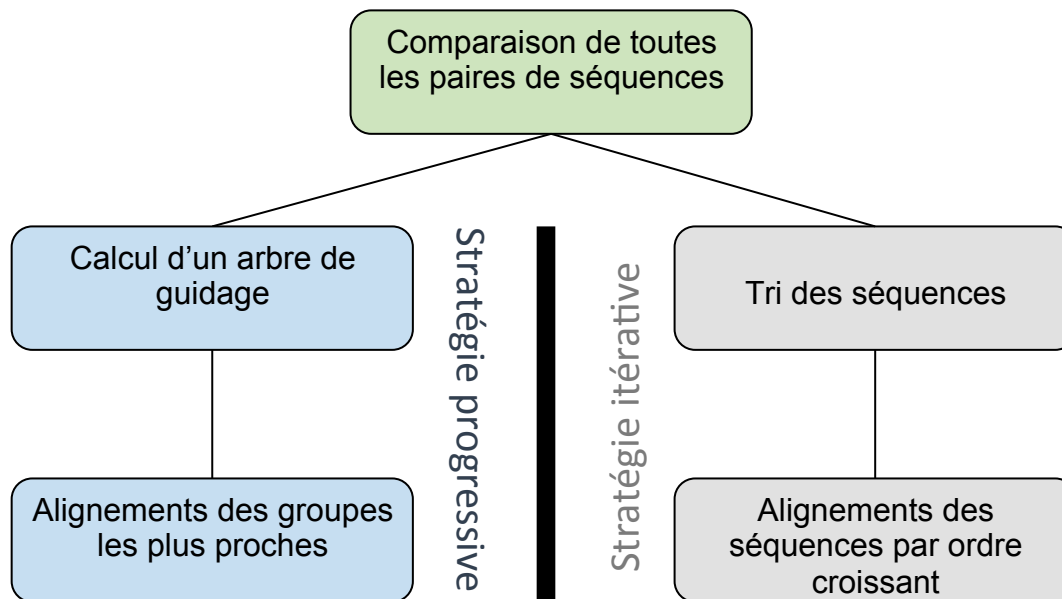
### Problème

L'alignement mathématique optimisé n'est pas nécessairement l'alignement biologique optimal.

Le temps CPU (temps de calcul sur ordinateur) et la mémoire requise sont prohibitifs pour un usage classique (temps requis est proportionnel à  $N^k$  avec  $k$  séquences de longueur  $N$ ).

En pratique, moins de 10 séquences peuvent être alignées.

# Stratégies d'alignements non-optimales



# Approches traditionnelles

- A. alignement multiple optimal
- B. alignement multiple progressif
- C. alignement multiple itératif

## B. Alignement multiple progressif

Évite le calcul de l'ensemble des alignements possibles

Non garantie d'obtenir l'alignement optimal

Principe :

Les séquences (ou groupe de séquences) sont alignées progressivement par paires

-> exemple Clustal, MUSCLE

## B. Alignement multiple progressif

Une approche alternative pour aligner des séquences multiples est de réaliser un **alignement progressif**.

L'algorithme procède en plusieurs étapes:

- Calculer une **matrice de distances**, qui indique la distance entre chaque paire de séquences.
- Construire un **arbre guide** qui regroupe en premier lieu les séquences les plus proches, et remonte en regroupant progressivement les séquences les plus éloignées.
- Utiliser cet arbre pour **aligner progressivement les séquences**.

Il s'agit d'une approche heuristique. Cette approche est possible pour un grand nombre de séquences, mais ne peut pas garantir de retourner l'alignement optimal.



## B. Alignement multiple progressif

1 ère étape: construction de la matrice de distance On effectue un alignement par paires entre chaque paire de protéines

A partir de chaque alignement par paire, calculer la distance entre les deux séquences.

## B. Alignement multiple progressif

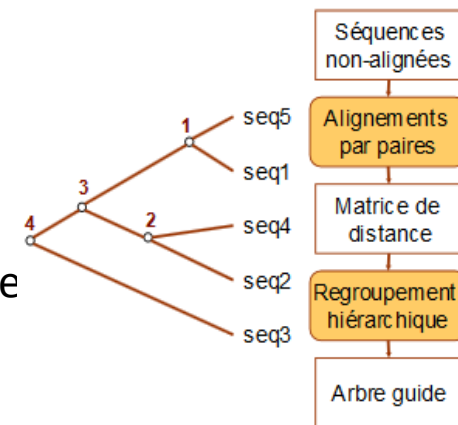
2ème étape : construire l'arbre guide

On peut calculer un arbre à partir d'une matrice de distance par regroupement hiérarchique.

On commence par regrouper les deux séquences les plus proches (groupe 1). On regroupe ensuite les groupes les plus proches. Les deux séquences les plus proches (groupe 2). Un groupe avec un groupe (groupe 3). Une séquence avec un groupe précédent (groupe 4).

Cet arbre sera ensuite utilisé comme guide pour déterminer l'ordre d'incorporation des séquences dans l'alignement multiple.

Attention ! Cet arbre ne doit pas être interprété comme un arbre phylogénétique. Il sert uniquement à identifier les similarités les plus fortes entre séquences pour construire l'alignement multiple.



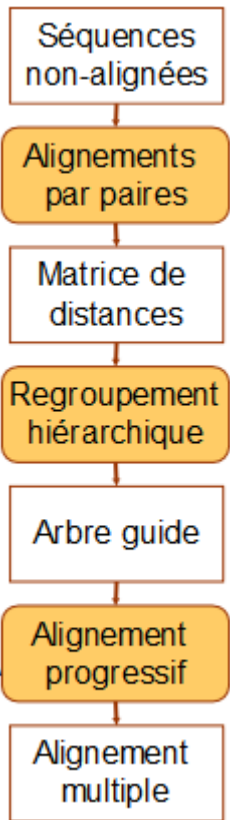
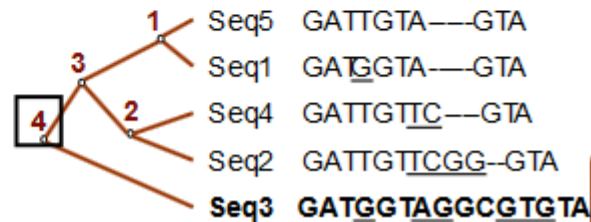
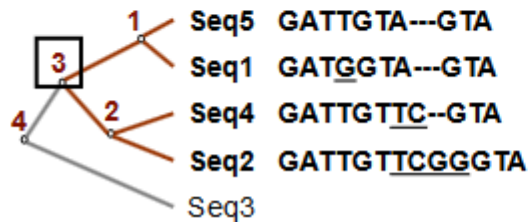
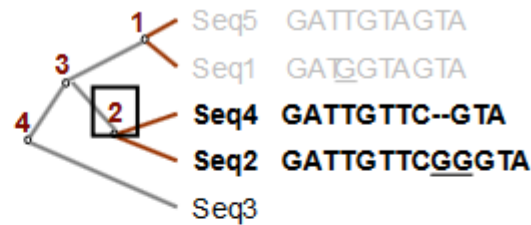
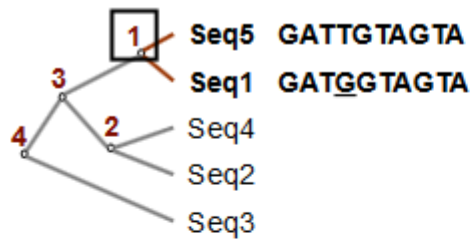
## B. Alignement multiple progressif

3 ème étape: alignement multiple

On construit un alignement multiple en incorporant progressivement les séquences selon leur ordre de branchement dans l'arbre guide, en remontant des plus proches aux plus éloignées.

## B. Aligement multiple progressif

3 ème étape: aligement multiple



# Approches traditionnelles

- A. alignement multiple optimal
- B. alignement multiple progressif
- C. alignement multiple itératif

## C. Alignement multiple itératif

Les méthodes itératives calculent une solution sous-optimale et continuent de la modifier intelligemment en utilisant la programmation dynamique ou d'autres méthodes jusqu'à ce que la solution converge. Contrairement aux méthodes progressives, les méthodes itératives peuvent corriger de façon dynamique les erreurs d'alignement.

## MAFFT: méthode itérative (Kato et al, 2002) d'alignement multiple

MAFFT = programmes de **nouvelle génération**

MAFFT a été écrit dans le but explicite d'**accélérer considérablement** le processus d'alignement multiple,

permettant ainsi d'aligner **un grand nombre de séquences** sans pour autant sacrifier à la qualité de l'alignement.

# MAFFT: méthode itérative (Kato et al, 2002) d'alignement multiple

## 3 grandes étapes :

### 1<sup>ère</sup> ETAPE

- Chaque **acide aminé** est décrit par sa **polarité** et son **volume**, les **séquences sont réécrites** dans ce système. Par exemple, les acides aminés hydrophobes I, L, M et V forment un seul groupe, de même que D, E, N et Q, etc.

suite de lettres (chaque séquence) → une suite de valeurs numériques.

- Les **nucléotides** sont recodées en utilisant les **fréquences locales des quatre bases**.
- Puis les **segments de similarité entre chaque paire** de séquences sont **repérés** au moyen d'un **algorithme de calcul** appelé transformée de Fourier rapide, ou **Fast Fourier Transform (FFT)** en anglais.
- Les **paires de séquences** sont ensuite **alignées** sur la base de ces segments de similarité. Sauf pour des séquences très divergentes, ce procédé permet d'**aligner toutes les paires de séquences environ 10 fois plus vite que ClustalW** (approche progressive).



# MAFFT: Alignement multiple

## 2<sup>ème</sup> ETAPE

Un **arbre de guidage** est ensuite calculé à partir des alignements précédents.

La **distance entre 2 séquences** est estimée à partir du nombre de **mots de 6 lettres** (*k-mers*) que ces **séquences partagent** dans ce nouvel alphabet

# MAFFT: Alignement multiple

## 3<sup>ème</sup> ETAPE

Les séquences sont ensuite alignées progressivement en suivant l'ordre indiqué par l'arbre de guidage.

## MAFFT: Alignement multiple

Plusieurs programmes sont proposés sur la page de garde du site Web de MAFFT.

Ce que nous venons de décrire correspond à l'option FFT-NS-1.

Ce nom un peu barbare signifie Fast Fourier Transform-New Scoring matrix-1 step.

# MAFFT: Alignement multiple

MAFFT peut occasionnellement procéder à un **deuxième passage** (approche itérative).

Dans ce dernier, l'alignement réalisé précédemment sert à re-calculer la distance entre chaque paire de séquences, un **nouvel arbre de guidage** et un nouveau alignement multiple.

Cette option s'appelle FFT-NS-2.

# MAFFT: Alignement multiple

MAFFT peut procéder à un **raffinement** de l'alignement.

Dans ce cas, **l'arbre de guidage est scindé en deux**, puis les deux moitiés sont réalignées.

On **recommence** ainsi tant que **le score d'alignement s'améliore**.

On procède alors à un nombre  $i$  d'itérations  $i$  étant inconnu *a priori*.

Cette option porte le nom de FFT-NS- $i$ .

On peut, sur la page de garde de MAFFT, limiter à deux le nombre d'itérations (« two cycles only »).

# MAFFT: Alignement multiple

Il faut par ailleurs noter que le site de MAFFT propose des programmes fondés non pas sur la transformées de Fourier rapide, mais sur l'algorithme de **programmation dynamique**.

Ainsi le programme nommé **G-INS-i** aligne les paires de séquences suivant l'**algorithme global de Needleman-Wunsch**, calcule un arbre de guidage, aligne toutes les séquences suivant cet arbre et procède enfin à un raffinement de l'alignement comme décrit ci-avant.

Les programmes **L-INS-i** et **E-INS-i** procèdent de la même façon, mais avec l'algorithme d'**alignement local de Smith-Waterman**.

Bien entendu, ces programmes, **nettement plus lents, ne conviennent pas pour un grand nombre de séquences.**

Enfin, le programme **Q-INS-i** est spécifiquement dédié à l'alignement de séquences d'**ARN**.

# MAFFT - bilan des paramètres

1. Mode basique, rapide — **juste progressif**

- a) FFTNS1 (fftns --retree 1)
- b) FFTNS2 (fftns) (same as mafft --retree 2)

OK jusqu'à 1 000 séquences facilement alignables

2. Mode intermédiaire — **progressif + itérations**

- a) FFTNSI (fftnsi) default two cycles, or e.g. fftnsi --maxiterate 1000
- b) NWNSI (nwnsi) same as FFTNSI, but no FFT, Needleman-Wunsch only.

OK entre 100 et 500 séquences

3. Mode avancé — **progressif + itérations + consistance (cf. T-Coffee)**

- a) EINSI (einsi) Smith-Waterman (plusieurs régions similaires même ordre)
- b) LINSI (linsi) Smith-Waterman stricte (1 région similaires)
- c) GINSI (linsi) global Needleman-Wunsch

# Publication MAFFT

- <https://academic.oup.com/mbe/article/30/4/772/1073398/MAFFT-Multiple-Sequence-Alignment-Software-Version>



# Visitions le site de MAFFT

- <http://mafft.cbrc.jp/alignment/server/>

[Download version](#)

[Mac OS X](#)

[Windows](#)

[Linux](#)

[Source](#)

**Online version**

**Alignment**

[mafft --add](#)

[Merge](#)

[Phylogeny](#)

[Rough tree](#)

[Merits / limitations](#)

[Algorithms](#)

[Tips](#)

[Benchmarks](#)

[Feedback](#)

[Follow](#)



For a large number of short sequences, try [an experimental service](#) (2016/Jul).

Multiple sequence alignment and NJ / UPGMA phylogeny

**Input:**

Paste protein or DNA sequences in fasta format. [Example](#)

A large empty text box for pasting sequences.

Au lieu d'utiliser les séquence brutes fournies par Elisabeth nous allons les clusteriser car vous aurez beaucoup de séquences à soumettre à MAFFT en réalité.

or upload a **plain text** file:  Aucun fichier choisi

Use structural alignment(s)

Allow unusual symbols (Selenocysteine "U", Inosine "I", non-alphabetical characters, etc.) [Help](#)

**UPPERCASE / lowercase:**

Same as input

Amino acid → UPPERCASE / Nucleotide → lowercase

**Direction of nucleotide sequences:** [Help](#)

Same as input

Adjust direction according to the first sequence (accurate enough for most cases)

Adjust direction according to the first sequence (only for highly divergent data; **extremely slow**)

**Output order:**

Same as input

Aligned

**Notify when finished** (optional; recommended when submitting large data):

Email address:

# Visitions le site de CD-HIT

- [http://weizhong-lab.ucsd.edu/cdhit\\_suite/cgi-bin/index.cgi?cmd=Server%20home](http://weizhong-lab.ucsd.edu/cdhit_suite/cgi-bin/index.cgi?cmd=Server%20home)

**CD-Hit** permet de regrouper les séquences très proche (ici 95%), limite la redondance et donc le nombre de séquences. Basiquement, CD-HIT est un algorithme progressif glouton qui débute avec la séquence la plus longue et sera donc la séquence représentative du premier cluster, puis traite les séquences restantes de la plus longue à la plus court pour classer chaque séquence comme une séquence redondante ou représentative basée sur ses similitudes avec la séquence existante.

# Visitions le site de CD-HIT

- [http://weizhong-lab.ucsd.edu/cdhit\\_suite/cgi-bin/index.cgi?cmd=Server%20home](http://weizhong-lab.ucsd.edu/cdhit_suite/cgi-bin/index.cgi?cmd=Server%20home)

CD-Hit -> Très utile dans le cadre de **très grand jeux de données**.

Input data :

<http://genoweb.toulouse.inra.fr/~formation/AlignClusteringLISBP/data>

# le site de CD-HIT

**CD-HIT Suite: Biological Sequence Clustering and Comparison**  
- this server at UCSD is not under regular maintenance, you may try [the server](#) if there is an issue

---

[Server home](#) | [cd-hit](#) | [cd-hit-est](#) | [h-cd-hit](#) | [h-cd-hit-est](#) | [cd-hit-2d](#) | [cd-hit-est-2d](#) | [result](#) | [calculated cluster](#)

**Sequence file and databases**

Load Query Fasta file from your computer:  all.fasta.txt

Incorporate annotation info at header line

**Sequence Identity Parameters**

Sequence identity cut-off

# Résultats de CD-Hit

You job 1487863706 is finished.

Program you ran: cd-hit

You input file is all.fasta.txt and we named it as [1487863706.fas.0](#)

Summary information for 1487863706.fas.0 included in [1487863706.fas.0.stat](#)

You required 1 runs for sequence clustering

1. Fasta file for representative sequences at 95% identity is [1487863706.fas.1](#)

Summary information for 1487863706.fas.1 included in [1487863706.fas.1.stat](#)

Corresponding cluster file is [1487863706.fas.1.clstr](#)

Sorted cluster file by size is [1487863706.fas.1.clstr.sorted](#)

Generated shell script is [run-1487863706.sh](#)

```
faa_stat.pl 1487863706.fas.0
```

```
cd-hit -i 1487863706.fas.0 -d 0 -o 1487863706.fas.1 -c 0.95 -n 5 -G 1 -g 1 -b 20 -s 0.0 -aL 0.0 -aS 0.0
```

```
faa_stat.pl 1487863706.fas.1
```

```
clstr_sort_by.pl no < 1487863706.fas.1.clstr > 1487863706.fas.1.clstr.sorted
```

```
clstr_list.pl 1487863706.fas.1.clstr 1487863706.clstr.dump
```

```
gnuplot1.pl < 1487863706.fas.1.clstr > 1487863706.fas.1.clstr.1; gnuplot2.pl 1487863706.fas.1.clstr.1 1487863706.fas.1.clstr.1.png
```

```
clstr_list_sort.pl 1487863706.clstr.dump 1487863706.clstr_no.dump
```

```
clstr_list_sort.pl 1487863706.clstr.dump 1487863706.clstr_len.dump len
```

```
clstr_list_sort.pl 1487863706.clstr.dump 1487863706.clstr_des.dump des
```

# 57 clusters de séquences

```

>BAC_FRA_CC_emb_CAH06518.1 cluster21 [Bacteroides fragilis NCTC 9343]
MSLFNDKVAKL LAGHEALLMRKNEPVEE GNGVITRYRYPVLTAAHTPVFWRDYLNEETNPF LMERIGMNA TLNAGAIKWDGKYLMLVRVEGAD
>BAC_OVA_CC_gb_EDO10988.1 cluster21 [Bacteroides ovatus ATCC 8483]
MNNMKSTFLFL LTTMTCTAYGQSSNHK ENKLPDWAFFGGFERPKVNPVISP IENTKFYCP LTKDSIAWESNDTFNPAATLYNGEIVVLYRA
>BAC_THE_CC_gb_AAO76140.1 cluster46 [Bacteroides thetaiotaomicron VPI-5482]
MNKIQIPWEERFV GCTDVMWRYSQNPVIGRYHIPSSNSIFNSAVVPFKDGFAGVFRCDNKAVQMNI FTGFSKDG IHWDISHEPIQFKAGNTEM
>BAC_THE_CC_gb_AAO78885.1 cluster26 [Bacteroides thetaiotaomicron VPI-5482]
MKSTFLFLVTT TMMTCTALGQPSNDKKNVLPDWAFFGGFERPQGANPVISP IENTKFYCPMTQDYVAVESNDTFNPAATLHDGKIVVLYRAEDK
>CAL_POL_CC_ref_WP_026485574.1 cluster1 [Caldanaerobius polysaccharolyticus]
MSKQLIVGEAL FNIPWQDRPAGCNDV VVRYQQNPVIPRDLIPSSNSIFNSAVVPNGE FAGVFRCDTKSRQMEIHSGRSKDGLNWEIDHERIK
>CEL_MIX_CC_gb_AAS19693.1 cluster16 [Cellvibrio mixtus]
MSSFREKAKALL QQHETLITKKNVAVKRDGNGIYDCYENPILTAEHAPVFWRYDLNEKTNP HLMERQGINAAFNFGAMYWNGKYILAVRVEGVD
>DTA_FER_CC_gb_ACT94389.1 cluster27 [Dyadobacter fermentans DSM 18053]
MKNFIAGFAICTS ILTGQAF AQTENKLPDWAFFGGFKRPAGVNPVSPDSTTRFFCFPMNKREVDWESNDTFNPAATMKNGKIVVLYRAEDKAGK
>LIS_INN_CC_emb_CAC96089.1 cluster31 [Listeria innocua Clip11262]
MNIYRYEENPLIT PLDVKPIHEGFVIGAFNGGVAEYNGEVLLLLVAEKPVSEDP EIVLAPVYNAKNKELEQSFRLDDENYDFEDPRMIRS
>RUM_ALB_CC_gb_ADU21379.1 cluster24 [Ruminococcus albus 7 = DSM 20455]
MIHEKYTEMRNEQE ALLSRKNTKTSFYNGIYDRYEHPVLTREH IPLHWRYDLNKTENPFFQERLGINAVFNAGAIKLNDRYCLVARVEGNDRK
>RUM_ALB_CC_gb_ADU20661.1 cluster38 [Ruminococcus albus 7 = DSM 20455]
MKTQIINGVSLP NIPWQDKPADCKDVIWRYDANPIIPRDQLPTSNSIFNSAVVPYSEKGRFAGVFRVDDKCRNMELHAGFSKDG IHWIDINPD
>THE_SP_CC_gb_ABY93074.1 cluster51 [Thermoanaerobacter sp. X514]
MFRRLRLSNKPIL SPIKEHEWEKEAVFNAAVIYEGNKFHLYFRASNNK FVLNTEKPEEKYKFKVSVSIGYAVSE DGINFERFDKPVLVGEIPQEA
>THE_SP_CC_gb_ABY93073.1 cluster54 [Thermoanaerobacter sp. X514]
MKLKRRLSDKPV LMPKAENEWERA AVFNTAAIYDNGLFH LIYRATDIGPHAKY GKYISRLGYAVSKDGINFMR LDKPVMNETEQELRGLDEP
>UNC_ORG_CC_gb_ADD61463.1 cluster3 [uncultured organism]
MSMSKVIIPWEERF AGCKDVLWRSVANPIIPRDL LPTSNSIFNSAVVPFGDGFAGVFRCDT SRRMRLHVGF SKDAINWNIKEEPLKFCDD
>BAC_XYL_emb_CBK68185.1 cluster28 [Bacteroides xylanisolvens XB1A]
MKRKLQNTAYLL MAAAFVASCSEKKQISEFPDWA WTDQRPPEGVNP IISPDTTTTFYCPMRQDSIAWEASD TFNPAATIYDGKVVVLYRAEDN
>CLO_SP_gb_AEY67872.1 cluster37 [Clostridium sp. BNL1100]
MSIIIGKTKLKN IFPWQDKPLGCS SVIWRHEGNPIIGWNPTPTRIARYNS SVVPWNSG FAGIFRADHKDGKARIHVGFSSDGVNWNVEDAPIVWHD
>PAE_POL_gb_ADO59098.1 cluster32 [Paenibacillus polymyxa SC2]
MSKLVNVDLVGN RIVGDSLSKMPWQDKPEGSEAPVWRHTENPVIGRN PVPGIARIFNSAVAPYEGRFVGI FRAETINGRPHLHLGWSDDGLA
>PRE_DNB_gb_AGB28392.1 cluster39 [Prevotella dentalis DSM 3688]
MNTLKIQGFALFG MPWEDRAENDKHMWRSHRNPIIPRDL LPTSNSIFNSAVVPFGDGYAGVFRCDT NRRMALHAGFSTNGIDWHINEAPLR
>LAC_PHY_gb_ABX42090.1 cluster47 [Lachnoclostridium phytofermentans ISDg]
MKNIPWEPRPVDC EDVVWRYSKNPIIHRNEIKRSNSIFNSAVVPFKDGYAGVFRCD DKRREMLLHAGFSVDGVKWNINPEPIEFQSEVEDSEP

```

Sauver le  
fichier fasta en  
57clusters.fasta

[Download version](#)

[Mac OS X](#)

[Windows](#)

[Linux](#)

[Source](#)

**Online version**

**Alignment**

[mafft --add](#)

[Merge](#)

[Phylogeny](#)

[Rough tree](#)

[Merits / limitations](#)

[Algorithms](#)

[Tips](#)

[Benchmarks](#)

[Feedback](#)

[Follow](#)



For a large number of short sequences, try [an experimental service](#) (2016/Jul).

Multiple sequence alignment and NJ / UPGMA phylogeny

**Input:**

Paste protein or DNA sequences in fasta format. [Example](#)

input = données issus de CD-HIT  
57clusters.fasta

or upload a **plain text** file:  Aucun fichier choisi

Use structural alignment(s)

Allow unusual symbols (Selenocysteine "U", Inosine "I", non-alphabetical characters, etc.) [Help](#)

**UPPERCASE / lowercase:**

Same as input

Amino acid → UPPERCASE / Nucleotide → lowercase

**Direction of nucleotide sequences:** [Help](#)

Same as input

Adjust direction according to the first sequence (accurate enough for most cases)

Adjust direction according to the first sequence (only for highly divergent data; **extremely slow**)

**Output order:**

Same as input

Aligned

**Notify when finished** (optional; recommended when submitting large data):

Email address:



## Advanced settings

### Strategy

- Auto (FFT-NS-1, FFT-NS-2, FFT-NS-i or L-INS-i; depends on data size) **Updated**

### Progressive methods

- FFT-NS-1 (Very fast; recommended for >2,000 sequences; progressive method)
- FFT-NS-2 (Fast; progressive method)
- G-INS-1 (Slow; progressive method with an accurate guide tree)

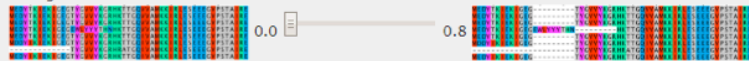
### Iterative refinement methods

- FFT-NS-i (Slow; iterative refinement method)
- E-INS-i (Very slow; recommended for <200 sequences with multiple conserved domains and long gaps) [Help](#) **Updated** (2015/Jun)
- L-INS-i (Very slow; recommended for <200 sequences with one conserved domain and long gaps) [Help](#)
- G-INS-i (Very slow; recommended for <200 sequences with global homology) [Help](#)
- Q-INS-i (Extremely slow; secondary structure of RNA is considered; recommended for a global alignment of highly divergent ncRNAs with <200 sequences  $\times$  <1,000 nucleotides; the number of iterative cycles is restricted to two, 2016/May) [Help](#)

### Align unrelated segments, too? *in Alpha Testing* (2014/Mar)

If the input data is expected to be globally conserved but locally contaminated by unrelated segments, try 'Unalignlevel>0' and possibly 'Leave gappy regions'.

#### Unalignlevel:

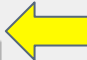


↑ Default

This feature is available only when G-INS-1 or G-INS-i is selected in the **Strategy** section above.

- Try to align gappy regions anyway
- Leave gappy regions (Not recommended for >~1,000 sequences)

**Parameters:**


Scoring matrix for amino acid sequences:  

Scoring matrix for nucleotide sequences:

↑ Switch it to '1PAM / κ=2' when aligning closely related DNA sequences.

Gap opening penalty:  (1.0 - 5.0)

Offset value:  (0.0 - 1.0)

**Score of N in nucleotide data:** [Example](#) 

↑ Long stretches of Ns tend to be gapped (excluded from the alignment).

- (nzero) N has no effect on the alignment score.
- (nwildcard) N is treated like a wildcard. *Experimental option* (2016/Apr/26)

↑ Try this if Ns should be aligned with usual letters.

 **Mafft-homologs** (Collects homologs from SwissProt by BLAST and performs profile-based alignments; Protein only): [Help](#) 

On

Show homologs (if any)

Number of homologs:  (5 - 200)

Threshold:  $E =$   ( $1e-5$  -  $1e-40$ )

**Plot LAST hits** (DNA only):

The top sequence vs the others  The longest sequence vs the others

Plot and alignment  Plot only  Alignment only

Threshold:



# Alignement -> score -> matrice de substitution

- **Distance d'édition**

- Selon ce concept, le bon alignement est celui qui minimise les opérations à réaliser pour passer d'une séquence à l'autre.
- Opérations: conservation, remplacement/mutation, délétion, insertion. Une pénalité peut être affectée à chaque opération, par exemple  $c=0$ ,  $m=1$ ,  $d=2$ ,  $i=2$ . La distance finale entre les deux séquences (distance d'édition) est la somme de ces pénalités.

	Seq 1	CAGTGGT-GC	
	Seq 2	CA-TCGTAGC	$c=0, m=1, d=2, i=2.$
Ou,	distance	ccicmccdcc = 0+0+2+0+1+0+0+2+0+0 = 5	
variante:	ressemblance	ccicmccdcc = 2+2-1+2-1+2+2-1+2+2 = 11	
			$c=2, m=-1, d=-1, i=-1.$

- **Une délétion à l'intérieur d'une séquence est considéré comme une insertion dans la séquence lui faisant face.**

# Matrices de Substitution

- Matrice 4X4 (nucléotides) ou 20x20 (acides aminés) décrivant la distance ou la similitude entre résidus.
  - Estiment le coût ou le taux de remplacement d'un résidu par un autre (distance).
  - Le choix d'une matrice affecte fortement le résultat de l'analyse. Chaque matrice de score représente implicitement une théorie évolutive donnée

## Matrices DNA

	<b>A</b>	<b>C</b>	<b>G</b>	<b>T</b>
<b>A</b>	2	-1	-1	-1
<b>C</b>	-1	2	-1	-1
<b>G</b>	-1	-1	2	-1
<b>T</b>	-1	-1	-1	2

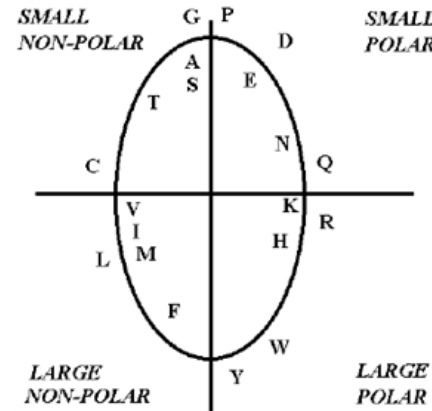
Matrice identité

	<b>A</b>	<b>C</b>	<b>G</b>	<b>T</b>
<b>A</b>	3	-1	1	-1
<b>C</b>	-1	3	-1	1
<b>G</b>	1	-1	3	-1
<b>T</b>	-1	1	-1	3

Matrice transition/transversion

# Matrices de Substitution

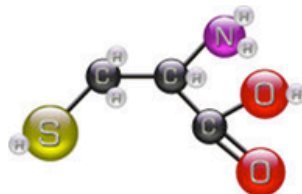
- Matrices fondées sur le code génétique
  - Les scores sont déterminés en fonction du nombre commun de nucléotides présents dans les codons des acides aminés, ce qui revient à considérer le minimum de changements nécessaires en bases pour convertir un acide aminé en un autre.
- Matrices fondées sur les propriétés physicochimiques
  - Les plus courantes sont celles basées sur le caractère hydrophile ou hydrophobe des protéines. Ces matrices sont peu utilisées.



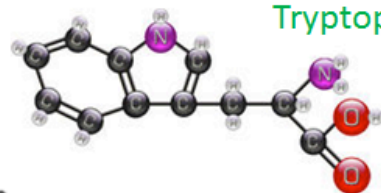
Une représentation bidimensionnelle des propriétés des aa calculée d'après la matrice de Dayhoff par G. Vriend, Centre for Molecular and Biomolecular Informatics, University of Nijmegen

# Matrice de Dayhoff (1979)

	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V	B	Z	X	
A	2																							
R	-2	6																						
N	0	0	2																					
D	0	-1	2	4																				
C	-2	-4	-4	-5	12																			
Q	0	1	1	2	-5	4																		
E	0	-1	1	3	-5	2	4																	
G	1	-3	0	1	-3	-1	0	5																
H	-1	2	2	1	-3	3	1	-2	6															
I	-1	-2	-2	-2	-2	-2	-3	-3	-2	5														
L	-2	-3	-3	4	-6	-2	-3	-4	-2	2	6													
K	-1	3	1	0	-5	1	0	-2	0	-2	-3	5												
M	-1	0	-2	-3	-5	-1	-2	-3	-2	2	4	0	6											
F	-4	-4	-4	-6	-4	-5	-5	-5	-2	1	2	-5	0	9										
P	1	0	-1	-1	-3	0	-1	-1	0	-2	-3	-1	-2	-5	6									
S	1	0	1	0	0	-1	0	1	-1	-1	-3	0	-2	-3	1	2								
T	1	-1	0	0	2	-1	0	0	-1	0	-2	0	-1	-3	0	1	3							
W	-6	2	-4	-7	-8	-5	-7	-7	-3	-5	-2	-3	-4	-0	-6	-2	-5	17						
Y	-3	-4	-2	4	0	-4	-4	-5	0	-1	-1	-4	-2	7	-5	-3	-3	0	10					
V	0	-2	-2	-2	-2	-2	-2	-1	-2	4	2	-2	2	-1	-1	-1	0	-6	-2	4				
B	0	-1	2	3	-4	1	2	0	1	-2	-3	1	-2	-5	-1	0	0	-5	-3	-2	2			
Z	0	0	1	3	-5	3	3	-1	2	-2	-3	0	-2	-5	0	0	-1	-6	-4	-2	2	3		
X	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	



Cystéine



Tryptophane

# Matrice de Blossum 62

	C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W	
C	9																				
S	-1	4																			
T	-1	1	5																		
P	-3	-1	-1	7																	
A	0	1	0	-1	4																
G	-3	0	-2	-2	0	6															
N	-3	1	0	-2	-2	0	6														
D	-3	0	-1	-1	-2	-1	1	6													
E	-4	0	-1	-1	-1	-2	0	2	5												
Q	-3	0	-1	-1	-1	-2	0	0	2	5											
H	-3	-1	-2	-2	-2	-2	1	-1	0	0	8										
R	-3	-1	-1	-2	-1	-2	0	-2	0	1	0	5									
K	-3	0	-1	-1	-1	-2	0	-1	1	1	-1	2	5								
M	-1	-1	-1	-2	-1	-3	-2	-3	-2	0	-2	-1	-1	5							
I	-1	-2	-1	-3	-1	-4	-3	-3	-3	-3	-3	-3	-3	1	4						
L	-1	-2	-1	-3	-1	-4	-3	-4	-3	-2	-3	-2	-2	2	2	4					
V	-1	-2	0	-2	0	-3	-3	-3	-2	-2	-3	-3	-2	1	3	1	4				
F	-2	-2	-2	-4	-2	-3	-3	-3	-3	-3	-1	-3	-3	0	0	0	-1	6			
Y	-2	-2	-2	-3	-2	-3	-2	-3	-2	-1	2	-2	-2	-1	-1	-1	-1	3	7		
W	-2	-3	-2	-4	-3	-2	-4	-4	-3	-2	-2	-3	-3	-1	-3	-2	-3	1	2	11	
C		S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W	

# Quelle matrice doit-on utiliser?

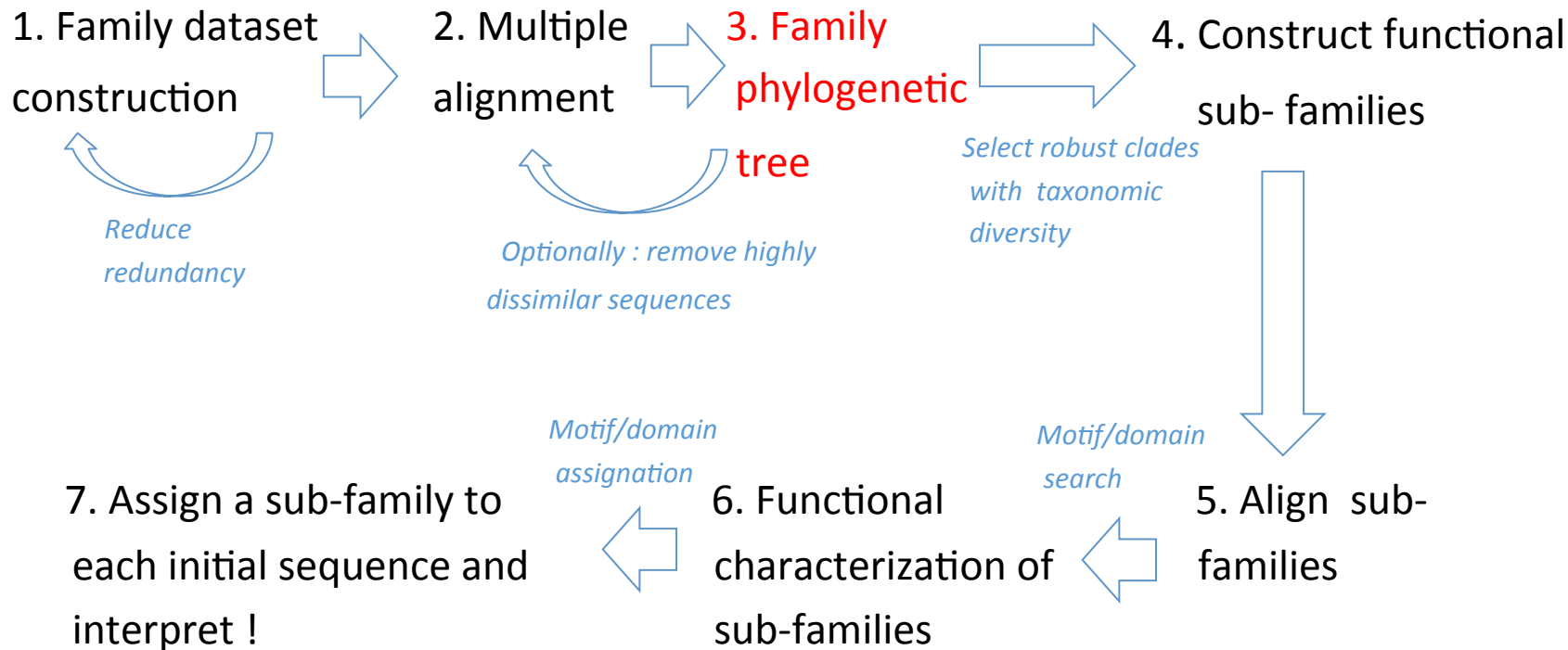
- Les matrices BLOSUM sont les plus souvent proposées comme matrices par défaut car les fréquences de substitution sont directement calculées à partir de l'alignement.
- La BLOSUM62 (ou PAM120) est utilisée comme matrice par défaut car elle offre un bon compromis quand les distances évolutives entre les séquences ne sont pas connues.
- La BLOSUM80 (ou PAM40) donnera de meilleurs résultats pour des séquences proches dans l'évolution. Elle tend à trouver des alignements courts fortement similaires.
- La BLOSUM30 (ou PAM350) donnera de meilleurs résultats pour des séquences éloignées dans l'évolution. Elle trouvera de plus longs alignements locaux de faible conservation.



# Partie 2 : Construction d'arbres phylogénétiques

Concepts de base

# Une démarche bioinfo générique



## Partie 2 - la stratégie utilisée:

*Material and Methods - Mewis et al. App. And Env. Microb. 2016*

*“Manual separation of subfamilies was decided based on phylogenetic distances in this reduced tree. Subfamilies were required to contain at least 5 sequences found in this reduced tree in order to generate a proper multiple-sequence alignment”*

Adaptation : define subfamilies using the refined **phylogenetic tree** and the three following criteria

- Sequences from robust and fully-resolved clades **MAFFT**
- Enough sequence number
- Large enough taxonomic diversity (at the class level) **NCBI**

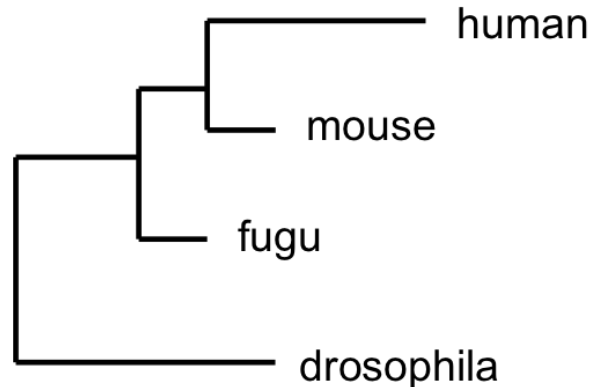
## La reconstruction d'arbre phylogénétique : plan

- Introduction et concepts de base (arbres, distances)
- Focus sur les méthodes de distance
  - UPGMA
  - NJ
- Interprétation
  - Test de la robustesse d'une topologie
  - Problèmes fréquents

# Qu'est-ce qu'un arbre phylogénétique ?

En **phylogénie moléculaire**, lorsqu'on analyse une famille de gènes

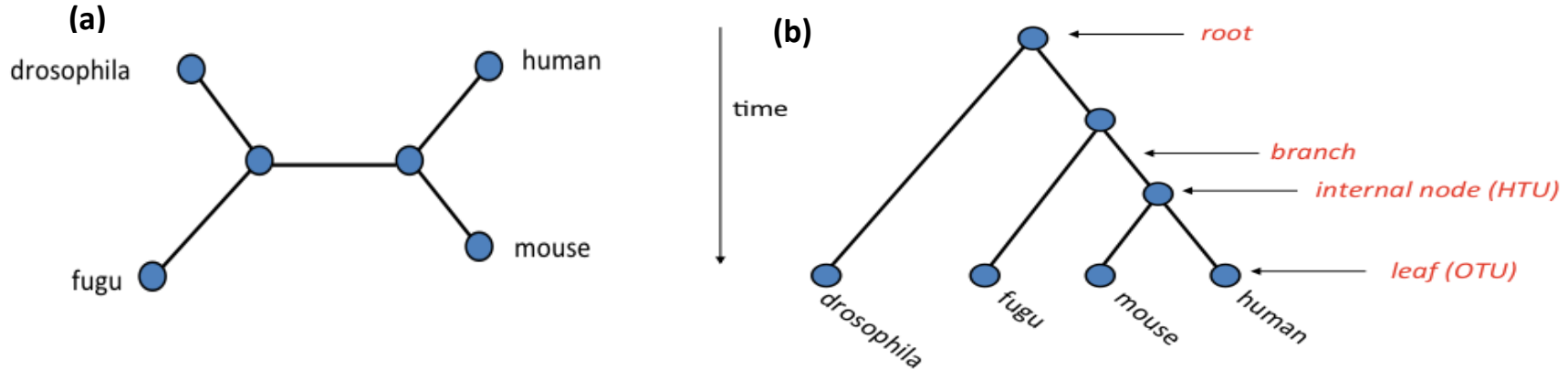
- un arbre phylogénétique est **une représentation de l'histoire évolutive de ces gènes** au cours du temps
- Cette représentation est **simplifiée**
- **L'histoire des gènes peut être différente de l'histoire évolutive des unités taxonomiques qui portent ces gènes**





# Les arbres phylogénétiques: notions de base

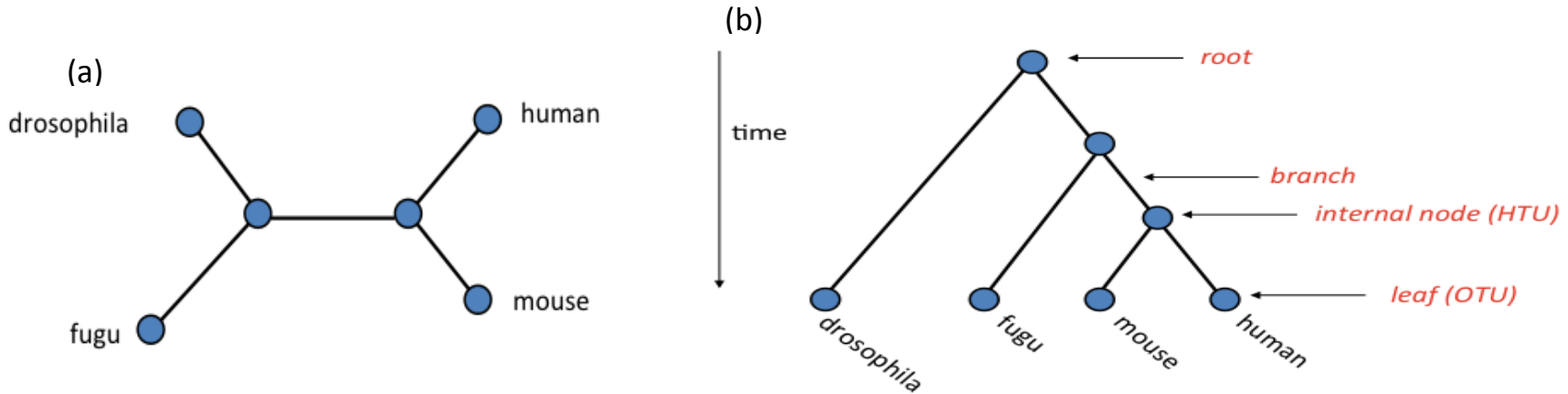
- Un arbre phylogénétique est défini par sa **topologie** et ses **longueurs de branche**
- Un arbre phylogénétique peut-être **non raciné (a)** ou **raciné (b)**



- Un arbre contient des **branches**, des **noeuds internes** et des **feuilles**
- Les termes usuels pour les taxa sont : **OTU**- Operational Taxonomic Units et **HTU**- Hypothetical Taxonomic Units

# Les arbres phylogénétiques: notions de base

- Un arbre phylogénétique peut-être *non raciné (a)* ou *raciné (b)*



- La très grande majorité des méthodes phylogénétiques s'intéressent aux arbres **bifurquants** : chaque nœud interne a un degré 3
- La plupart des méthodes phylogénétiques produisent des **arbres non-racinés**



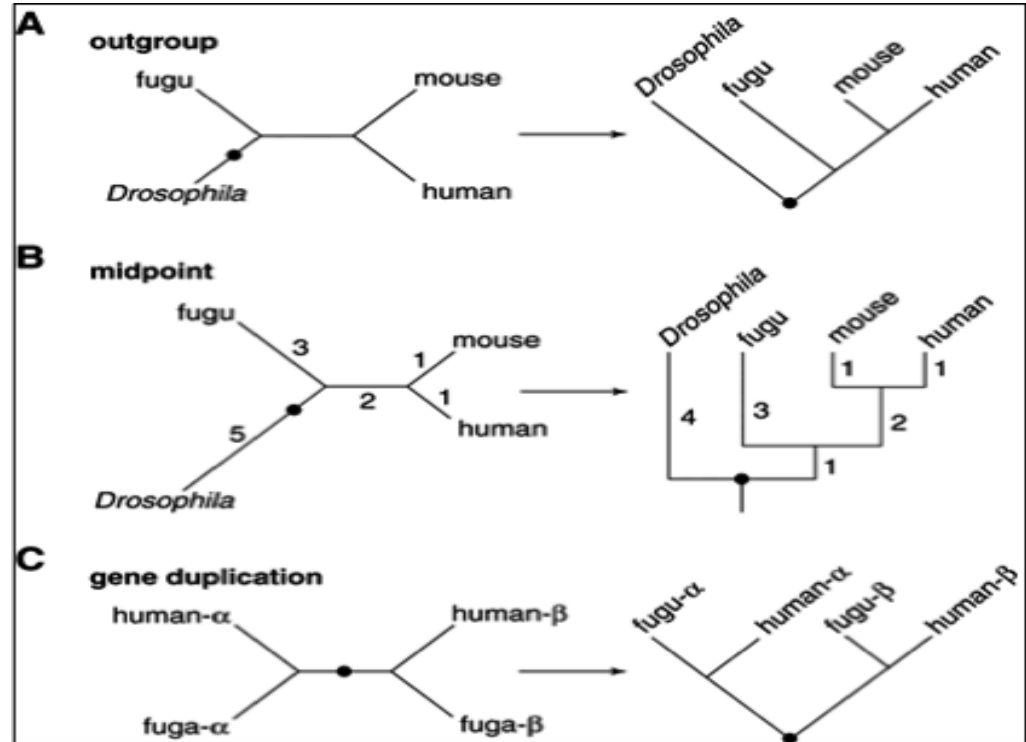
# Comment enraciner un arbre ?

- Trois méthodes existent

## A. Outgroup rooting

## B. Midpoint rooting

## C. Usage of external knowledge (ex. ancestral **gene duplication**)



# Les formats d'arbres phylogénétique

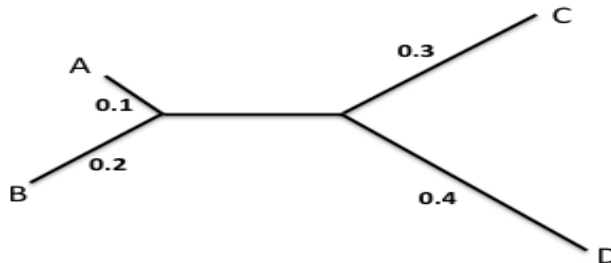
- Deux formats très utilisés: NEWICK and NEXUS



NEWICK: ((A,B), (C,D))

```
#NEXUS
BEGIN TAXA;
  TAXLABELS A B C D;
END;

BEGIN TREES;
  TREE tree1 = ((A,B),(C,D));
END;
```



NEWICK: ((A:0.1,B:0.2):0.2, (C:0.3,D:0.4))

```
#NEXUS:
Begin trees;
Translate
1 A,
2 B,
3 C,
4 D,
;
Tree tree2= [&U] ((1:0.1,2:0.2):0.2, (3:0.3,3:0.4));
End;
```

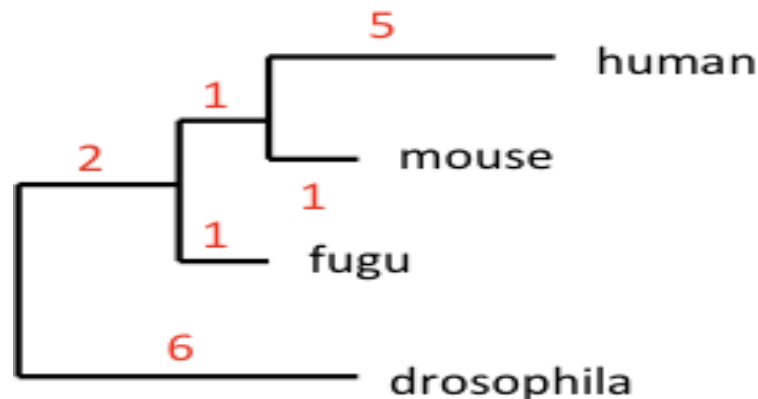
## Les distances génétiques (ou évolutives)

- Une distance génétique (évolutive) est **une mesure de la divergence entre deux séquences génétiques**
- Les **modalités de calcul de la distance génétique** entre deux séquences constituent une étape clé d'une analyse phylogénétique
  - C'est la première étape des méthodes basées sur des matrices de distance (UPGMA, NJ)
  - Dans les autres méthodes (maximum de vraisemblance, méthodes bayésiennes) les distances sont utilisées à travers les modèles de substitutions (nucléotides, AA) qui servent à calculer la vraisemblance d'un arbre

# Distances et arbres

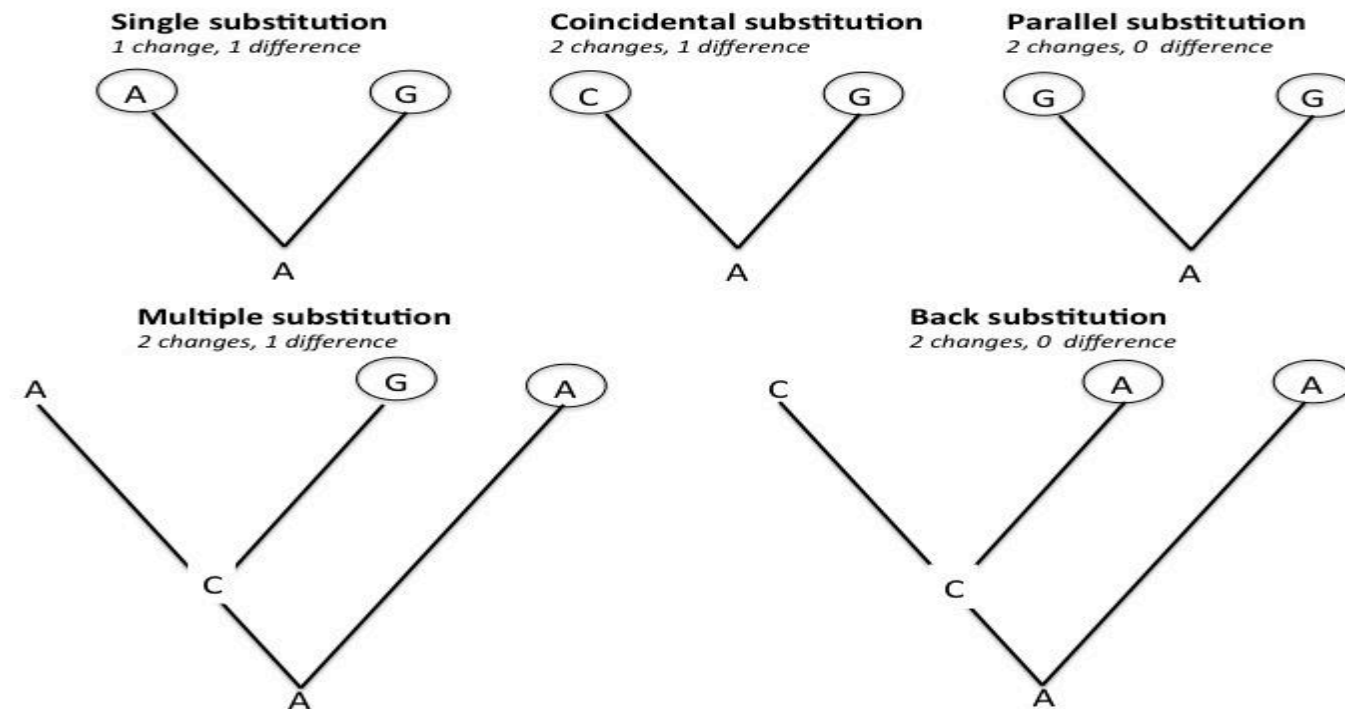
- Pour les séquences liées par un arbre, les **longueurs de branches** représentent la **distance entre les noeuds (séquences) dans l'arbre**
- Sous l'hypothèse d'horloge moléculaire, la **distance génétique est linéairement proportionnelle au temps écoulé**

human	x			
mouse	6	x		
fugu	7	3	x	
<i>drosophila</i>	14	10	9	x
	human	mouse	fugu	<i>drosophila</i>



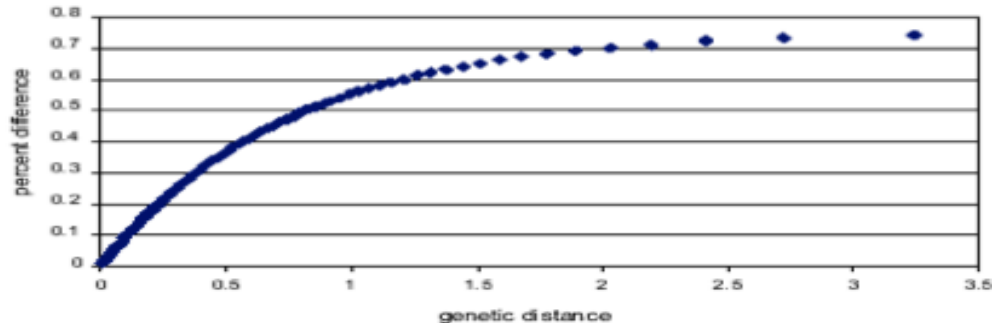
# Distances génétiques vs Distances observées

La distance basée sur le nombre de substitutions observées n'est pas toujours informative !



# Distances génétiques vs Distances observées

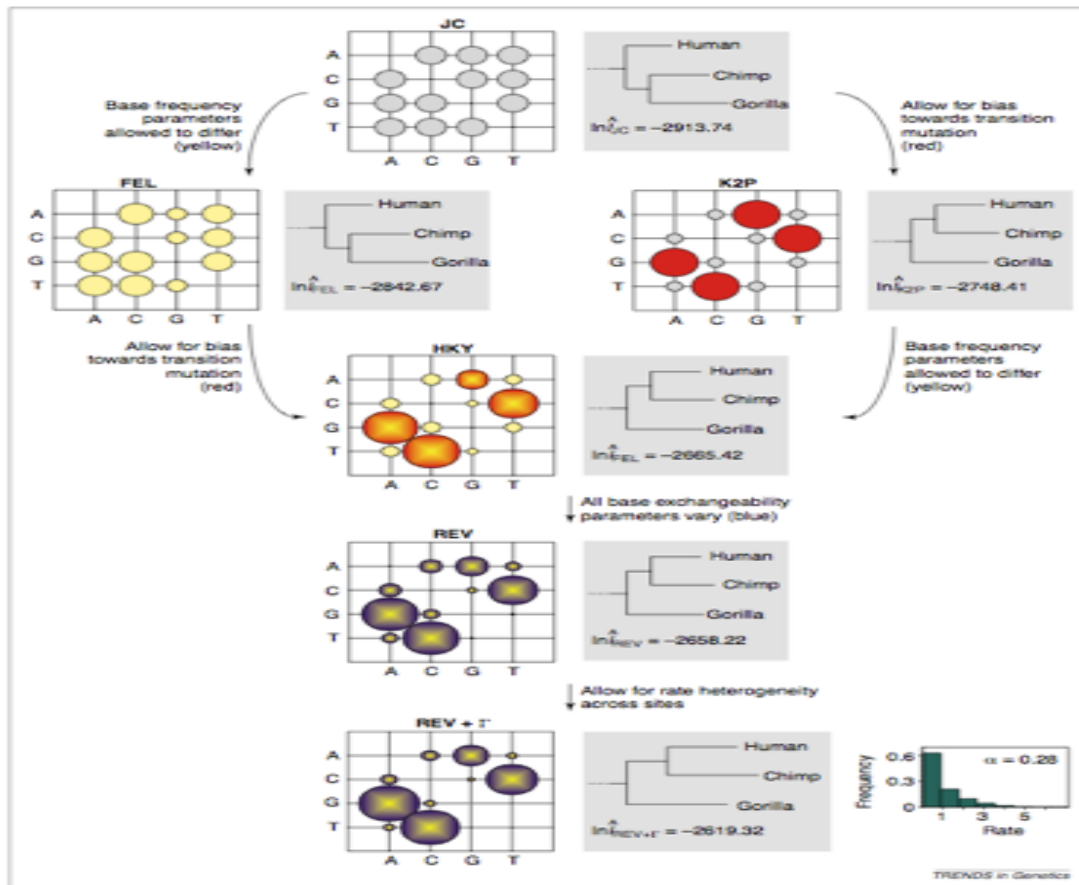
- La distance observée peut être calculée en comptant le nombre de sites où deux séquences diffèrent : elle est exprimée par le **nombre de différences nucléotidiques par site (p distance)**
- La distance observée est une sous-estimation de la distance génétique due aux substitutions multiples par site et à la saturation : des **modèles de substitution** sont utilisés.



## Les modèles de substitution nucléotidique

- Une substitution de nucléotide peut être modélisée comme un **événement aléatoire** ;
- Les **processus de Markov en temps continu** sont particulièrement adaptés pour décrire les taux de substitution à chaque site, avec les hypothèses suivantes :
  - Chaque site évolue indépendamment ;
  - La substitution d'un nucléotide  $i$  en un nucléotide  $j$  est indépendante du nucléotide présent avant  $i$  (propriété de Markov) ;
  - Les fréquences des 4 nucléotides sont à l'équilibre (stationnarité).
- Les substitutions à un site donné sont décrites par une chaîne de Markov dont les états sont les 4 nucléotides (A,T,C,G) et les changements de nucléotides sont décrits par la matrice des probabilités de transition  $P(t)$ .

# Les modèles de substitution nucléotidique : un aperçu





# Choix d'un modèle nucléotidique

- L'utilisation d'un modèle donnée détermine les **longueurs des branches** de l'arbre et peut parfois changer sa topologie
- **Conseils**
  - Aller du modèle le plus simple vers le plus complexe
  - le modèle **GTR (General Time Reversible)** est le plus complexe (9 paramètres)
  - On peut effectuer un test statistique pour choisir le modèle le plus pertinent (LRT, AIC, BIC)

Il est recommandé d'expérimenter !

# Les modèles de substitution protéiques

- **Concept similaire : les substitutions multiples d'acides aminés conduisent à une sous-estimation des distances d'évolution entre deux protéines homologues.**
- La fréquence de substitution des acides aminés dépend de l'AA: elle est plus élevée entre les acides aminés proches en terme de propriétés physiques (polarité, hydrophobicité, ...)
- Trop de paramètres (190) pour estimer les paramètres du modèle probabiliste => des **modèles empiriques** sont utilisés
- Le taux de transition entre les acides aminés sont estimés à partir de grands alignements de référence obtenus par concaténation de plusieurs protéines homologues

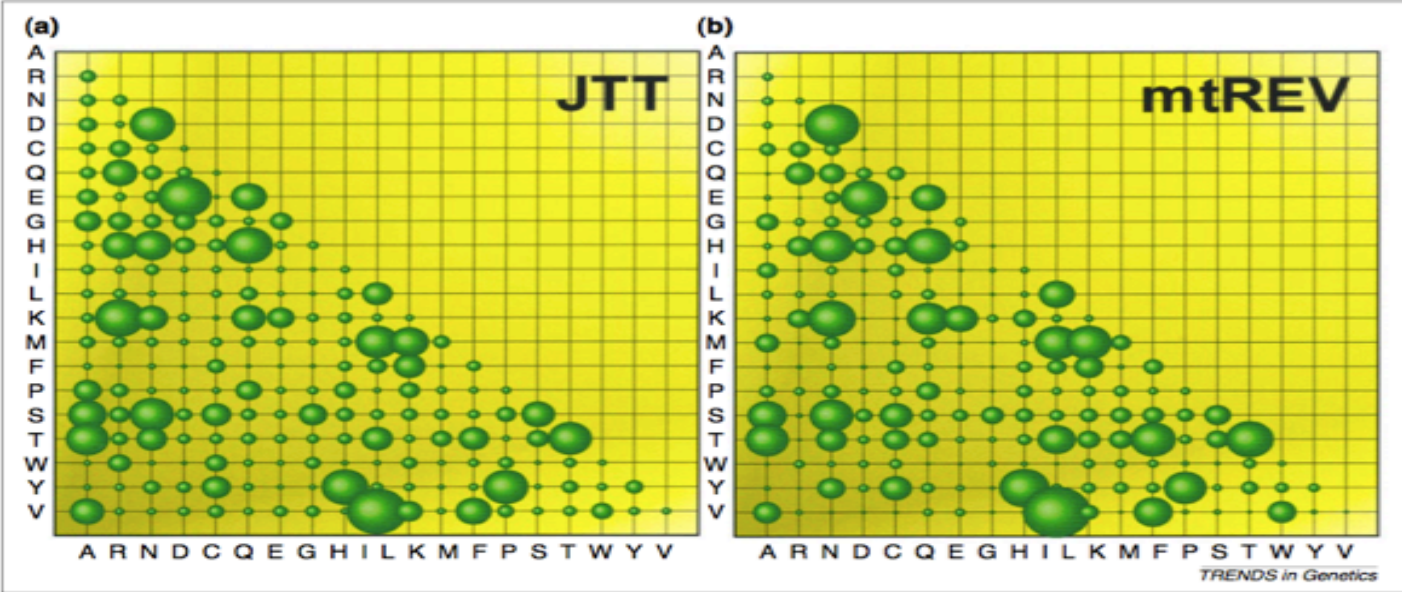
# Les principaux modèles d'évolution protéique

Model	Dataset	Ref
Poisson	Poisson process	Zuckerlandl, 1965
PAM	1300 protein sequences from 71 homolog families	Dayhoff 1978
Blosum	Extension of PAM dataset	Henikoff 1992
JTT	16 300 sequences	Jones 1992
mtREV	Mitochondrial DNA	Adachi 1996
<b>WAG &amp; LG</b>	<b>Likelihood methods</b>	<b>Whelan 2001</b>

- Les modèles WAG et LG sont les modèles les plus utilisés
- On peut effectuer un test statistique pour choisir le modèle le plus pertinent (LRT, AIC, BIC)

# Les modèles protéiques : exemple

JTT (1992, 16 300 sequences) vs mtREV (for mitochondrial proteins)

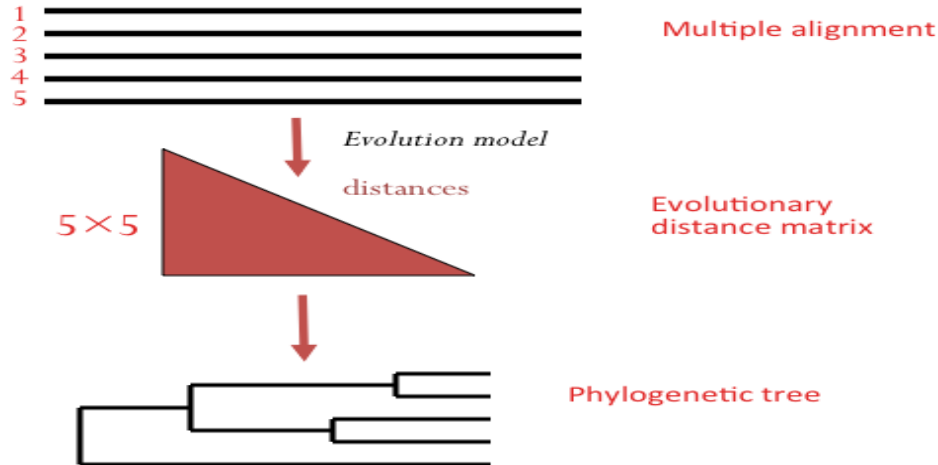


# Les méthodes d'inférence phylogénétique

Input data	Method	Principle of the algorithms
Distance matrix	Unweighted Pair Group Method (UPGMA)	clustering
	Neighbor-Joining (NJ)	clustering
Character state	Maximum Parsimony (MP)	Search for the tree(s) of minimum character changes
	Maximum Likelihood (ML)	Search for the tree(s) that maximizes the probability of observing the character states giving a tree topology and a model of evolution
	Bayesian Inference	Target a probability distribution of trees (set of possible trees for the data)

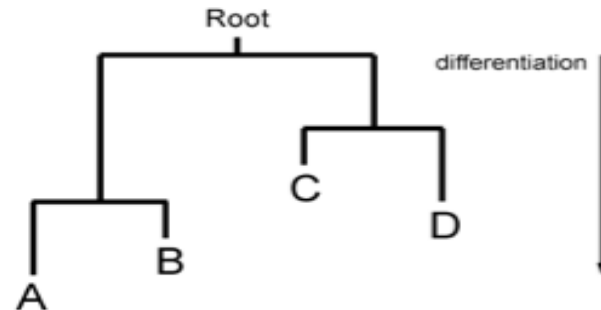
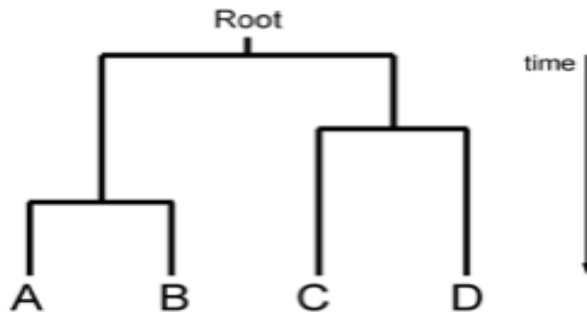
# Les méthodes de distances pour inférer un arbre phylogénétique

- Introduites en 1960, les plus anciennes !
- Objectif : représenter une matrice de distance génétique par un arbre
- Nécessitent un modèle d'évolution



# Les méthodes de distance

- Objectif des méthodes de distance : faire que les distances patristiques générées dans l'arbre représentent au mieux les distances de départ
- Deux méthodes principales
  - **UPGMA**: a clustering method that produced ultrametric trees
  - **Neighbor-Joining**: use a greedy algorithm to compute the Minimal Evolution tree *i.e.* the optimal topology is the one which minimizes the tree length



# La méthode UPGMA

## Unweighted Pair Group Method with Arithmetic Mean (Sneath & Sokal 1973)

Cette méthode est la plus simple, elle est aussi appelée **Classification ascendante hiérarchique**

Elle suppose que les taux de substitution entre séquences sont à peu près homogènes dans toutes les lignées (**hypothèse d'horloge moléculaire**)

Elle permet d'estimer des **arbres ultramétriques** à partir d'une matrice de distance

**L'algorithme est itératif** : il regroupe les séquences séparées par les distances les plus courtes, par ordre de similarité



# UPGMA : principe de l'algorithme

- **Données de départ** : matrice des distances  $D_{ij} = \{d_{ij}\}$  entre toutes les paires de séquences  $i$  et  $j$ .

	A	B	C	D	E
B	2				
C	4	4			
D	6	6	6		
E	6	6	6	4	
F	8	8	8	8	8

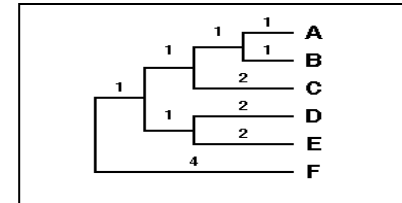
## Principe de l'algorithme

### Algorithme itératif:

1. Déterminer les deux clusters  $i$  et  $j$  pour lesquels  $d_{ij}$  est minimal.
2. Créer un nouveau cluster  $C_k = C_i \cup C_j$  et calculer  $d_{kl}$  pour tous les  $l$ .
3. Stocker  $C_i$  et  $C_j$  comme fils droit et gauche dans  $C_k$ . Stocker la hauteur de  $C_k$  :  $d_{ij} / 2$ .
4. Ajouter  $k$  à la liste des clusters à traiter et supprimer  $i$  et  $j$  de cette liste.

**Terminaison:** quand il ne reste plus d'un seul cluster

- **Résultat** : arbre ultramétrique

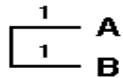


# UPGMA : exemple

- Si on a 6 UTOs séparées par une distance donnée (matrice) :

	A	B	C	D	E
B	2				
C	4	4			
D	6	6	6		
E	6	6	6	4	
F	8	8	8	8	8

- On commence par unir les 2 UTOs les plus similaires, A et B



$$\begin{aligned} \text{dist}(A,B),C &= (\text{dist}AC + \text{dist}BC) / 2 = 4 \\ \text{dist}(A,B),D &= (\text{dist}AD + \text{dist}BD) / 2 = 6 \\ \text{dist}(A,B),E &= (\text{dist}AE + \text{dist}BE) / 2 = 6 \\ \text{dist}(A,B),F &= (\text{dist}AF + \text{dist}BF) / 2 = 8 \end{aligned}$$

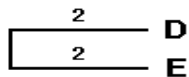
- Par la suite, A et B vont être considérées comme une unité (A,B)

# UPGMA : exemple

Les cycles suivants :

2ème cycle :

	A,B	C	D	E
C	4			
D	6	6		
E	6	6	4	
F	8	8	8	8

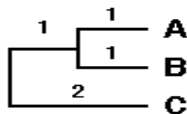


5ème cycle :

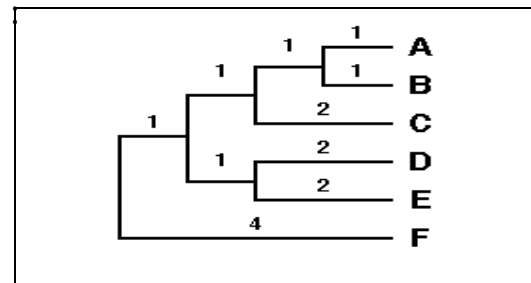
	ABC,DE
F	8

3ème cycle :

	A,B	C	D,E
C	4		
D,E	6	6	
F	8	8	8



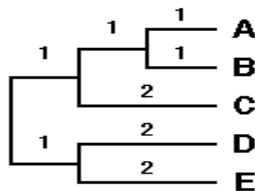
Fin



On obtient la topologie correcte

4ème cycle :

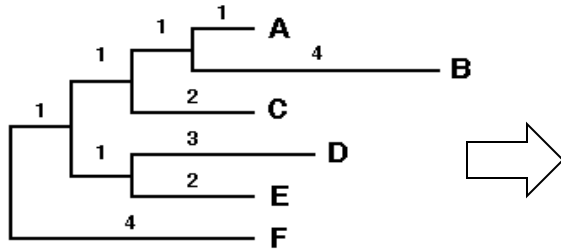
	AB,C	D,E
D,E	6	
F	8	8



# UPGMA : les limites

- La méthode UPGMA produit un arbre raciné
- Elle ne marche que pour les **distances ultramétriques** (principe d'horloge moléculaire stricte) : UPGMA échoue si le taux d'évolution varie entre les UTO

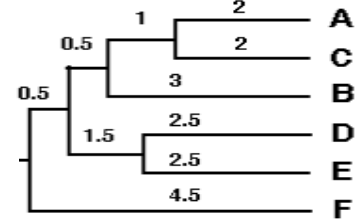
Arbre réel



Matrice de distances

	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

Arbre UPGMA



Topologie incorrecte

Conclusion : cette méthode est très rapide mais n'est quasiment plus employée en phylogénie moléculaire !

# L'algorithme du Neighbor-Joining

- La méthode a été proposée par **Saitou and Nei en 1987**
- Cette méthode utilise un **algorithme itératif**, elle est **très rapide**
- Cette approche produit des **arbres sans racine**
- La méthode est basée sur le **principe d'évolution minimale** :

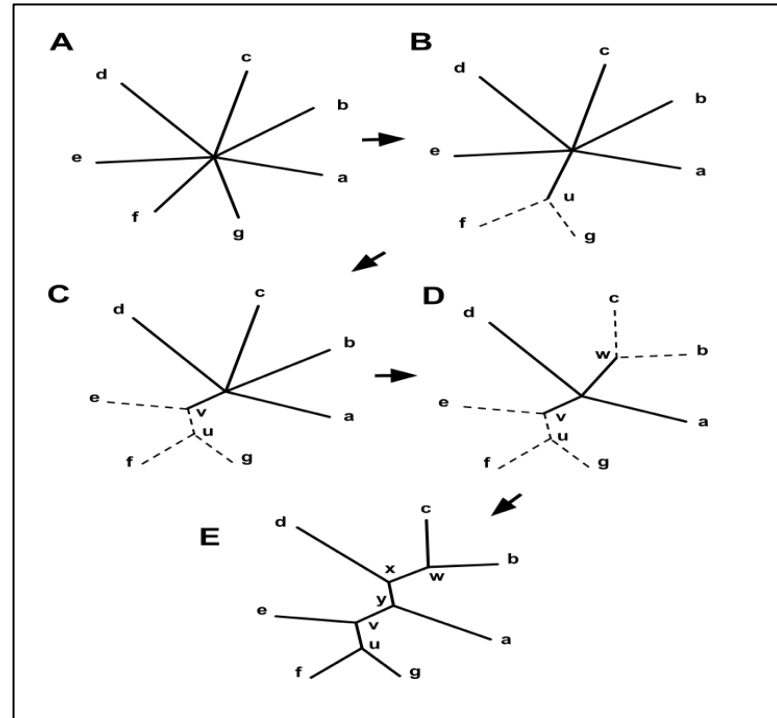
Si  $S$  = somme de toutes les branches de l'arbre avec  $T$  = nombre de branches et  $b_i$  = longueur de branche  $i$ , la meilleure topologie est celle qui correspond à  $S = \sum b_i$  minimum

- Contrairement à UPGMA, cette méthode **construit des arbres non ultramétriques**, traduisant une possible hétérogénéité dans le nombre de mutations accumulées

# Le Neighbor-Joining (NJ)

## Principe de l'algorithme :

- Start with a **star tree** (A)
- Compute the **matrix  $Q_{ij}^*$**  and find the pair of taxa with lowest value (here f and g)
- Join f and g and **create a new internal node, u**, as shown in (B)
- Compute the distances from node u to the nodes a-e
- **Repeat the process** : u and e are joined to the newly created v, as shown in (C).
- Two more iterations lead first to (D), and then to (E).



\* $Q_{ij}$  : matrice des distances corrigées

# Le Neighbor-Joining (NJ) : bilan

## Avantages :

- L'algorithme est constructif : **si les distances de départ sont patristiques, l'arbre obtenu est correct**
- Méthode très rapide et performante
- Mesure globale – l'arbre de longueur totale minimale

## Inconvénients :

- Donne des mauvais résultats en cas de fort taux de substitution ou de variations de taux entre les branches
- Pas recommandé pour des séquences très divergentes
- Problème d'attraction des branches longues

# Le Neighbor-Joining en pratique

## Conseils :

- **Tester plusieurs modèles d'évolution** pour construire la matrice de distances (en particulier pour les séquences assez divergentes)
- Lorsque les taux de substitutions sont élevés ou variables entre les UTO, utiliser **l'algorithme BioNJ\***, une variante de l'algorithme du NJ qui donne une meilleure précision de topologie

## Logiciels :

- **MAFFT** (UPGMA, NJ) : <http://mafft.cbrc.jp/alignment/server/>
- **NJ and BioNJ** <http://www.atgc-montpellier.fr/fastme/> or [http://www.phylogeny.fr/one\\_task.cgi?task\\_type=bionj](http://www.phylogeny.fr/one_task.cgi?task_type=bionj)
- **Seaview** (NJ and BioNJ, standalone program) <http://doua.prabi.fr/software/seaview>



# Comment évaluer un arbre phylogénétique ?

## Problème de confiance

- Quelle confiance avoir en l'arbre inféré?
- Quelles parties de l'arbre sont fiables / non fiables?
- Comment pouvons-nous valider l'arbre ?

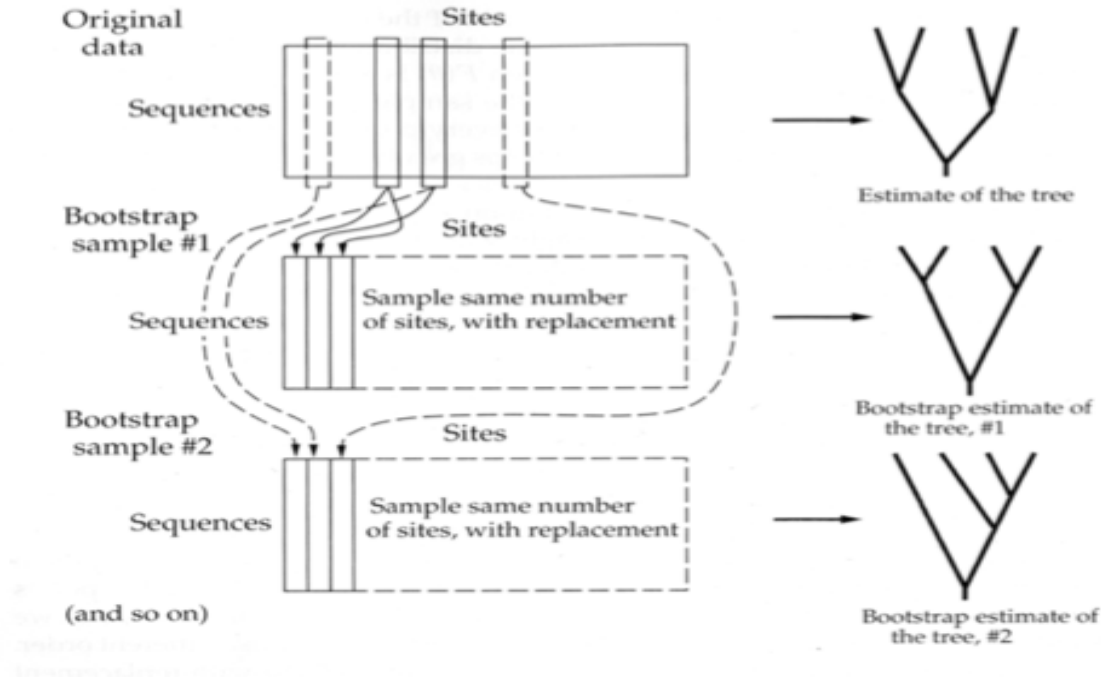
**Problème: le vrai arbre est inconnu !**

## Solution :

- Utilisez le bootstrap pour évaluer la fiabilité de l'arbre inféré et des clades
- Combiner les sous-échantillons et les arbres de consensus pour obtenir des valeurs de soutien sur les branches

# Évaluer la topologie d'un arbre

**Principe du Bootstrap:** on ré-échantillonne les positions d'un alignement



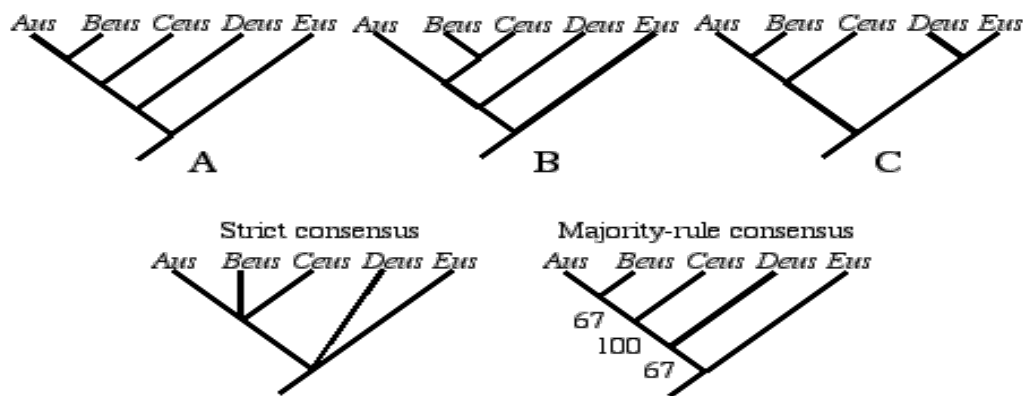
# Principe du bootstrap et de l'arbre consensus

- Inférer plusieurs ( $\sim 100$ ) arbres en utilisant des techniques de **rééchantillonnage**;
- Identifier et conserver uniquement les clades présents dans de nombreux arbres ;
- Combinez les différents arbres pour produire un **arbre de consensus** qui est compatible avec la plupart les arbres .
- En général, l'arbre de consensus n'a pas de longueur de branche et une résolution inférieure à celle de l'arbre d'origine.
- Superposer les valeurs de bootstrap sur l'arbre d'origine

# Comment construire un arbre consensus ?

## Règles de consensus :

- **Strict Consensus:** clades présents dans tous les arbres;
- **Règle de majorité:** les clades sont présents dans au moins la moitié des arbres;
- **Règle de majorité étendue:** clades sont présents dans au moins la moitié des arbres et un peu plus jusqu'à ce que l'arbre soit résolu.



# Le bootstrap en pratique

## Attention à l'interprétation des valeurs de bootstrap :

- Les valeurs de bootstrap n'ont pas d'interprétation statistique claire;
- Une valeur bootstrap de 95% ne signifie pas que le clade correspondant a 95% de chance d'être "vrai";
- Les valeurs de bootstrap sont difficiles à interpréter quantitativement.

Cependant, les valeurs de Bootstrap sont (assez) faciles à interpréter **qualitativement**:

- **Plus la valeur bootstrap est élevée, plus vous pouvez être confiant dans votre clade;**
- 95%, 90% et 66% constituent le seuil traditionnel pour être confiants dans un clade.

## Quelques derniers conseils

- **Attraction des longues branches:** Les longues branches ont tendance à se regrouper dans l'arbre
- Solution : "décomposer" de longues branches en ajoutant quelques UTO à l'analyse;
- **Saturation:** Les caractères ont évolué depuis si longtemps qu'ils sont presque aléatoires
- Solution : Supprimer les sites saturés et / ou taxons; Lorsqu'elles sont disponibles, utiliser des séquences protéiques au lieu des séquences nucléiques;
- **Filtrage des alignements** : pour les phylogénies basées sur un seul gène / protéine, le filtrage des régions avec "gaps" n'améliore pas la reconstruction des arbres (ref)
- **Solution : Ne pas filtrer les alignements single gene/protein !**

# Travaux pratiques

<http://mafft.cbrc.jp/alignment/server/>

• **Exercice 1 : Construire un arbre UPGMA et un arbre de Neighbor-Joining à partir de l'alignement des 57 séquence de la famille GH131**

*Commenter les différences*

Multiple sequence alignment by MAFFT ver.7

Clustal format | FastA format | MAFFT result | View | Tree | Refine dataset | Return to home

NJ or UPGMA tree (β)

57 sequences, 1839 total sites, 216 gap-free sites, 55 conserved sites

Get Reset

Settings

**Method:**

- NJ -- Conserved sites (55 AAs)
- NJ -- All gap-free sites (216 AAs)
- Average linkage (UPGMA) -- alignment scores (for up to 50,000 sequences)
- Minimum linkage -- alignment scores (for up to 50,000 sequences)
- Memory-saving tree -- alignment scores (for larger data)

**Substitution model** (valid when NJ is selected):

- JTT
- WAG **Alpha**
- Poisson

**Heterogeneity among sites** (valid when JTT or WAG is selected):

- Ignore (α = -)
- Estimate
- Specify; α = 1.00 (0.10 - 5.0)

**Bootstrap** (valid for NJ):

On

Number of resampling: 100 (5 - 1000)  
(The number of sequences must be <1000 for Poisson model, or <100 for other models.)

Get Reset

Alignment id = 1702012351s24814826CQ3URJ1cpv8W8zMG7AF  
Posted at Wed Feb 1 23:51:28 JST 2017  
[Feedback form](#)

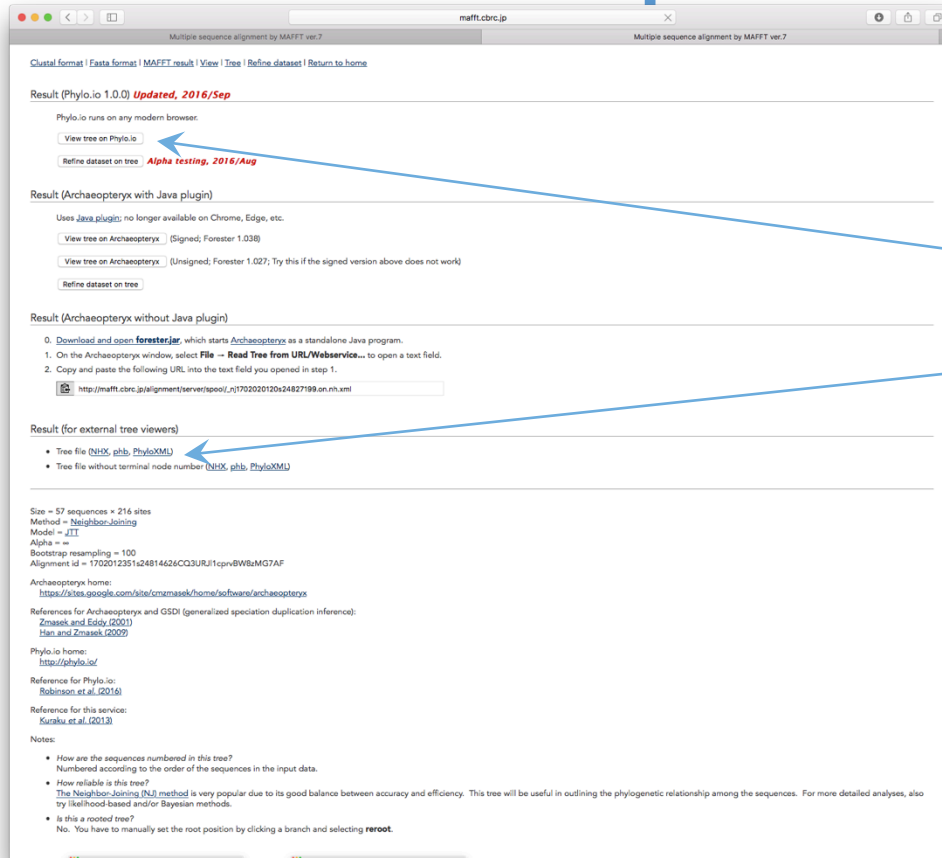
Pour le NJ utiliser les paramètres suivant :

- All gap-free sites
- **Substitution model : WAG**
- Heterogeneity among sites : Ignore (default)
- Bootstrap : yes

Visualiser l'arbre et interpréter



# Arbre des 57 séquences



Multiple sequence alignment by MAFFT ver.7

Clustal format | Fasta format | MAFFT result | View | Tree | Refine dataset | Return to home

Result (Phylo.io 1.0.0) **Updated, 2016/Sep**

Phylo.io runs on any modern browser:

[View tree on Phylo.io](#)

[Refine dataset on tree](#) **Alpha testing, 2016/Aug**

Result (Archaeopteryx with Java plugin)

Uses [Java plugin](#): no longer available on Chrome, Edge, etc.

[View tree on Archaeopteryx](#) (Signed; Forester 1.036)

[View tree on Archaeopteryx](#) (Unsigned; Forester 1.027; Try this if the signed version above does not work)

[Refine dataset on tree](#)

Result (Archaeopteryx without Java plugin)

0. Download and open [forester.jar](#) which starts [Archaeopteryx](#) as a standalone Java program.

1. On the Archaeopteryx window, select **File** → **Read Tree from URL/Webservice...** to open a text field.

2. Copy and paste the following URL into the text field you opened in step 1.

Result (for external tree viewers)

- Tree file ([NEXUS](#), [PHIL](#), [PhyloXML](#))
- Tree file without terminal node number ([NEXUS](#), [PHIL](#), [PhyloXML](#))

Size = 57 sequences \* 216 sites  
 Method = Neighbor-joining  
 Model = JTT  
 Alpha = -  
 Bootstrap resampling = 100  
 Alignment id = 1702012351a24814a26c303URJ1cprvBWBtMG7AF

Archaeopteryx home:  
<https://sites.google.com/site/crmamasek/home/software/archaeopteryx>

References for Archaeopteryx and GSDI (Generalized speciation duplication inference):  
[Zmasek and Eddy \(2001\)](#)  
[Han and Zmasek \(2009\)](#)

Phylo.io home:  
<http://phylo.io/>

Reference for Phylo.io:  
[Robinson et al. \(2016\)](#)

Reference for this service:  
[Kuraku et al. \(2013\)](#)

Notes:

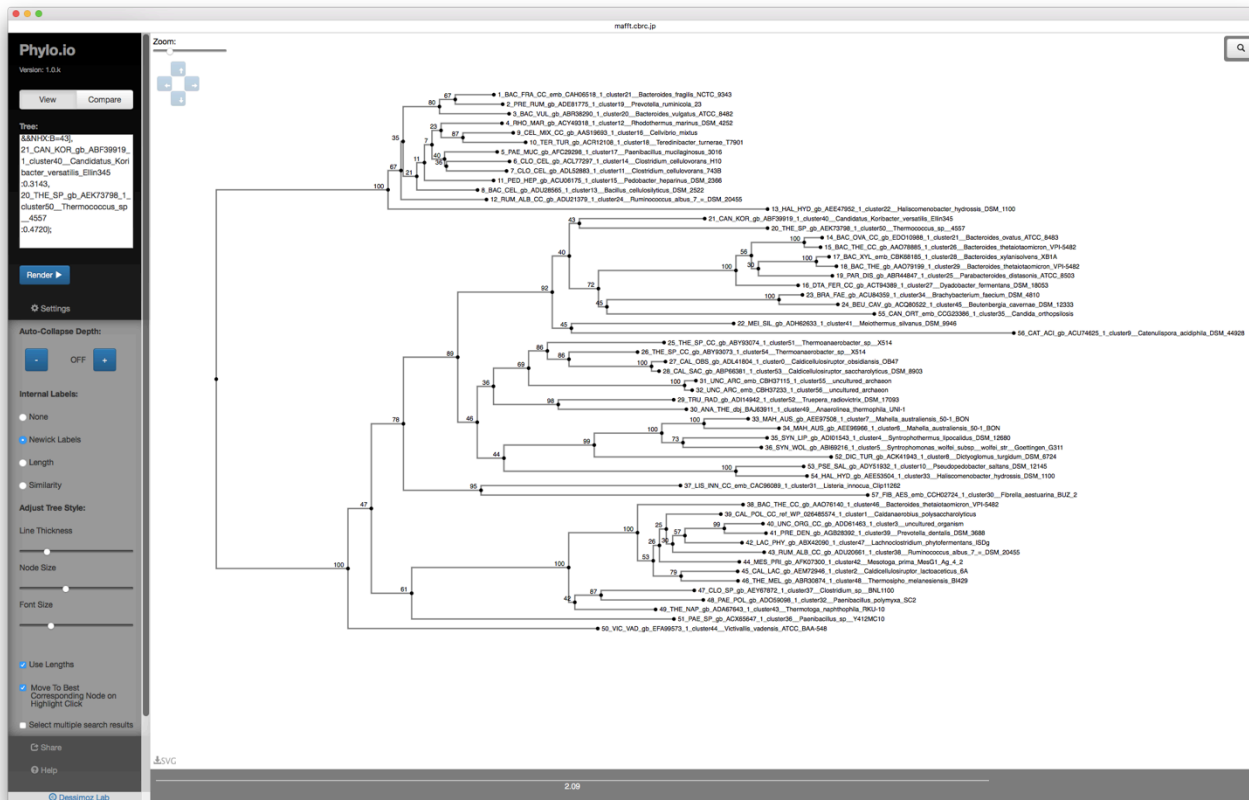
- How are the sequences numbered in this tree?  
 Numbered according to the order of the sequences in the input data.
- How reliable is this tree?  
 The Neighbor-joining (NJ) method is very popular due to its good balance between accuracy and efficiency. This tree will be useful in outlining the phylogenetic relationship among the sequences. For more detailed analyses, also try likelihood-based and/or Bayesian methods.
- Is this a rooted tree?  
 No. You have to manually set the root position by clicking a branch and selecting **reroot**.

Explore the tree using Phylo.io

Select “view tree on Phylo.io”

Export/save tree in a standard format

# Arbre des 57 séquences

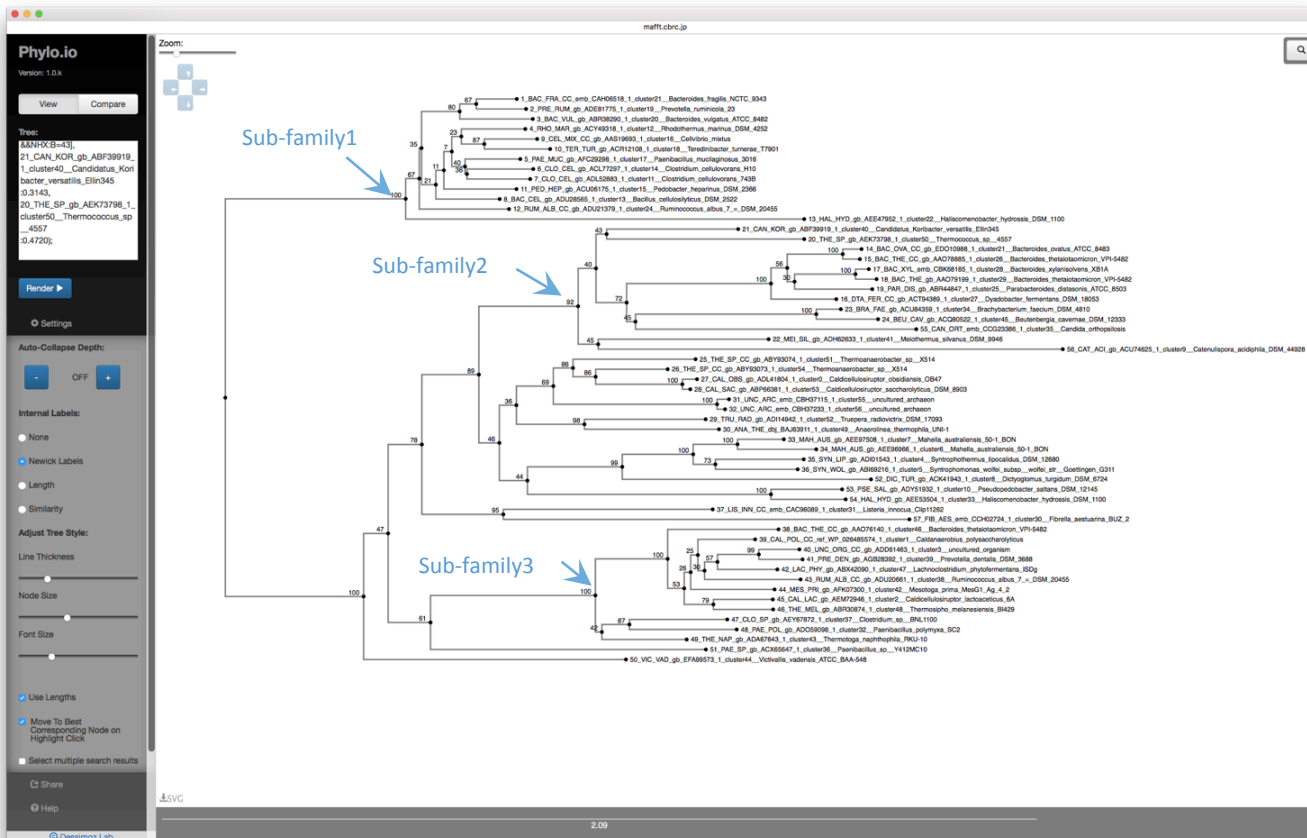


## Exercise 2: select 3 sub-families meeting the 3 following criteria:

1. Bootstrap clade values  $\geq 90$
2. At least 12 sequences
3. Large enough taxonomic diversity (at the class level)

DOMAIN	Mnemonic (memory aid)
Kingdom	D
Phylum	K
Class	P
Order	C
Family	O
Genus	F
species	G
subspecies	S

# Arbre des 57 séquences



## Subfamilies definition

- Sub\_family 1 : sequence 1 to 13
- Sub\_family 2 : sequence 14 to 24 + 55 + 56
- Sub\_family 3 : sequence 38 to 49

For taxonomic diversity, use <https://www.ncbi.nlm.nih.gov/Taxonomy/Browser/wwwtax.cgi> with species name

# Définition des sous-familles

## Sub-Family 1 : 13 seqs

>BAC\_FRA\_CC\_emb\_CAH06518.1 cluster21 [Bacteroides fragilis NCTC 9343]

[Bacteroidia](#);

>PRE\_RUM\_gb\_ADE81775.1 cluster19 [Prevotella rumincola 23]

Class : [Bacteroidia](#);

>BAC\_VUL\_gb\_ABR38290.1 cluster20 [Bacteroides vulgatus ATCC 8482]

[Bacteroidia](#);

>RHO\_MAR\_gb\_ACY49318.1 cluster12 [Rhodothermus marinus DSM 4252]

Class : ?

>PAE\_MUC\_gb\_AFC29298.1 cluster17 [Paenibacillus mucilaginosus 3016]

>CLO\_CEL\_gb\_ACL77297.1 cluster14 [Clostridium cellulovorans H10]

>CLO\_CEL\_gb\_ADL52883.1 cluster11 [Clostridium cellulovorans 743B]

>BAC\_CEL\_gb\_ADU28565.1 cluster13 [Bacillus cellulosilyticus DSM 2522]

>CEL\_MIX\_CC\_gb\_AAS19693.1 cluster16 [Cellvibrio mixtus]

Class : [Gammaproteobacteria](#);

>TER\_TUR\_gb\_ACR12108.1 cluster18 [Teredinibacter turnerae T7901]

[Gammaproteobacteria](#);

>PED\_HEP\_gb\_ACU06175.1 cluster15 [Pedobacter heparinus DSM 2366]

Phylum : [Bacteroidetes](#) Class :

Phylum : [Bacteroidetes](#)

Phylum : [Bacteroidetes](#) Class :

Phylum : [Bacteroidetes](#)

Phylum : [Bacilli](#); Class : [Bacillales](#)

Phylum : [Firmicutes](#); Class : [Clostridia](#)

Phylum : [Firmicutes](#); Class : [Bacilli](#)

Phylum : [Firmicutes](#); Class : [Bacilli](#)

Phylum : [Proteobacteria](#);

Phylum : [Proteobacteria](#); Class :

108  
Phylum : [Bacteroidetes](#); Class

# Définition des sous-familles

## Sub-Family 2 : 13 seqs

>BAC\_OVA\_CC\_gb\_EDO10988.1 cluster21 [Bacteroides ovatus ATCC 8483]  
>BAC\_THE\_CC\_gb\_AAO78885.1 cluster26 [Bacteroides thetaiotaomicron VPI-5482]  
>DTA\_FER\_CC\_gb\_ACT94389.1 cluster27 [Dyadobacter fermentans DSM 18053]  
>BAC\_XYL\_emb\_CBK68185.1 cluster28 [Bacteroides xylanisolvens XB1A]  
>BAC\_THE\_gb\_AAO79199.1 cluster29 [Bacteroides thetaiotaomicron VPI-5482]  
>PAR\_DIS\_gb\_ABR44847.1 cluster25 [Parabacteroides distasonis ATCC 8503]  
>THE\_SP\_gb\_AEK73798.1 cluster50 [Thermococcus sp. 4557]  
>CAN\_KOR\_gb\_ABF39919.1 cluster40 [Candidatus Koribacter versatilis Ellin345]  
>MEI\_SIL\_gb\_ADH62633.1 cluster41 [Meiothermus silvanus DSM 9946]  
>BRA\_FAE\_gb\_ACU84359.1 cluster34 [Brachybacterium faecium DSM 4810]  
>BEU\_CAV\_gb\_ACQ80522.1 cluster45 [Beutenbergia cavernae DSM 12333]  
>CAN\_ORT\_emb\_CCG23386.1 cluster35 [Candida orthopsilosis]  
>CAT\_ACI\_gb\_ACU74625.1 cluster9 [Catenulispora acidiphila DSM 44928]

Combien de phyla ? Combien de classes ?

# Définition des sous-familles

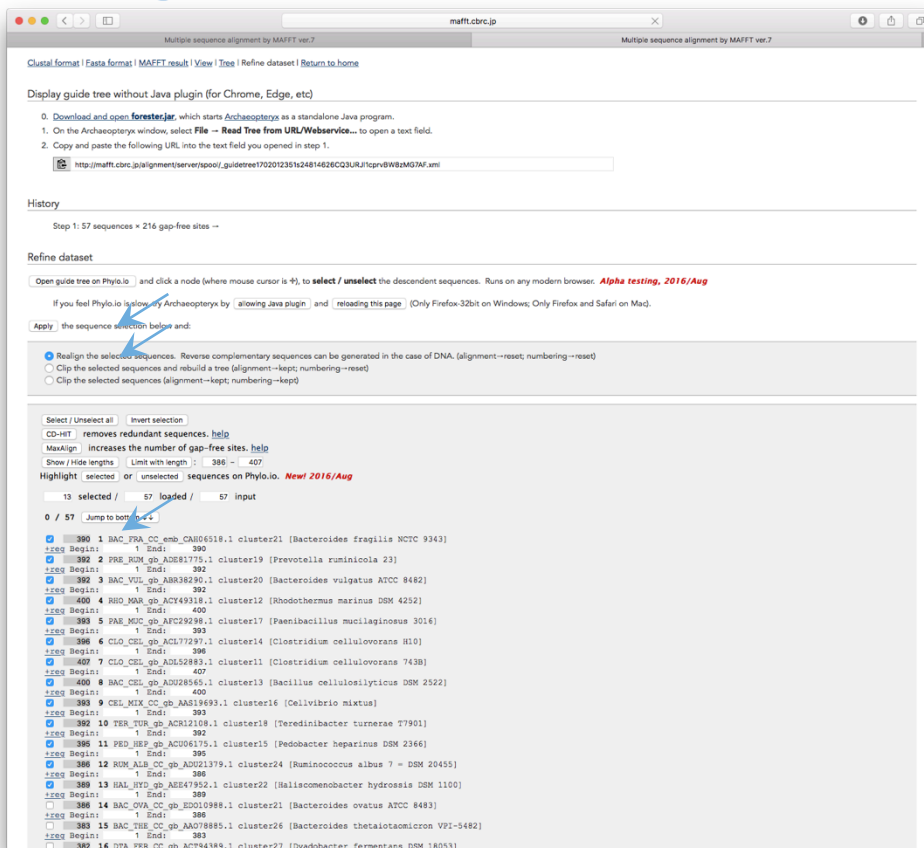
## Sub-Family 3 : 12 seqs

>BAC\_THE\_CC\_gb\_AAO76140.1 cluster46 [Bacteroides thetaiotaomicron VPI-5482]  
>CAL\_POL\_CC\_ref\_WP\_026485574.1 cluster1 [Caldanaerobius polysaccharolyticus]  
>RUM\_ALB\_CC\_gb\_ADU20661.1 cluster38 [Ruminococcus albus 7 = DSM 20455]  
>UNC\_ORG\_CC\_gb\_ADD61463.1 cluster3 [uncultured organism]  
>PRE\_DEN\_gb\_AGB28392.1 cluster39 [Prevotella dentalis DSM 3688]  
>LAC\_PHY\_gb\_ABX42090.1 cluster47 [Lachnospirillum phytofermentans ISDg]  
>MES\_PRI\_gb\_AFK07300.1 cluster42 [Mesotoga prima MesG1.Ag.4.2]  
>CAL\_LAC\_gb\_AEM72946.1 cluster2 [Caldicellulosiruptor lactoaceticus 6A]  
>THE\_MEL\_gb\_ABR30874.1 cluster48 [Thermosiphon melanesiensis BI429]  
>CLO\_SP\_gb\_AEY67872.1 cluster37 [Clostridium sp. BNL1100]  
>PAE\_POL\_gb\_ADO59098.1 cluster32 [Paenibacillus polymyxa SC2]  
>THE\_NAP\_gb\_ADA67643.1 cluster43 [Thermotoga naphthophila RKU-10]

Combien de phyla ? Combien de classes ?

# Alignement des sous-familles

- **Exercise 3:** Refine the dataset to build three sub-alignments corresponding to the three subfamilies



Multiple sequence alignment by MAFFT ver.7

Clustal format | FastA format | MAFFT result | View | Tree | Refine dataset | Return to home

Display guide tree without Java plugin (for Chrome, Edge, etc)

- Download and open **forester.jar**, which starts **Archoopteryx** as a standalone Java program.
- On the Archoopteryx window, select **File** → **Read Tree from URL/Webservice...** to open a text field.
- Copy and paste the following URL into the text field you opened in step 1.

History

Step 1: 57 sequences × 216 gap-free sites →

Refine dataset

Open guide tree on Phylo.io and click a node (where mouse cursor is ☞), to **select / unselect** the descendant sequences. Runs on any modern browser. **Alpha testing, 2016/Aug**

If you feel Phylo.io is slow, try Archoopteryx by allowing Java plugin and reloading this page (Only Firefox-32bit on Windows; Only Firefox and Safari on Mac).

Apply the sequence selection below and:

- Realign the selected sequences. Reverse complementary sequences can be generated in the case of DNA. (alignment--reset; numbering--reset)
- Clip the selected sequences and rebuild a tree (alignment--keep; numbering--reset)
- Clip the selected sequences (alignment--keep; numbering--keep)

Select / Unselect all | Invert selection

CD-HT removes redundant sequences. help

MaxAlign increases the number of gap-free sites. help

Show / Hide lengths | Limit with length : 386 - 407

Highlight selected or unselected sequences on Phylo.io. **New! 2016/Aug**

13 selected / 57 loaded / 57 input

0 / 57 | Jump to bottom

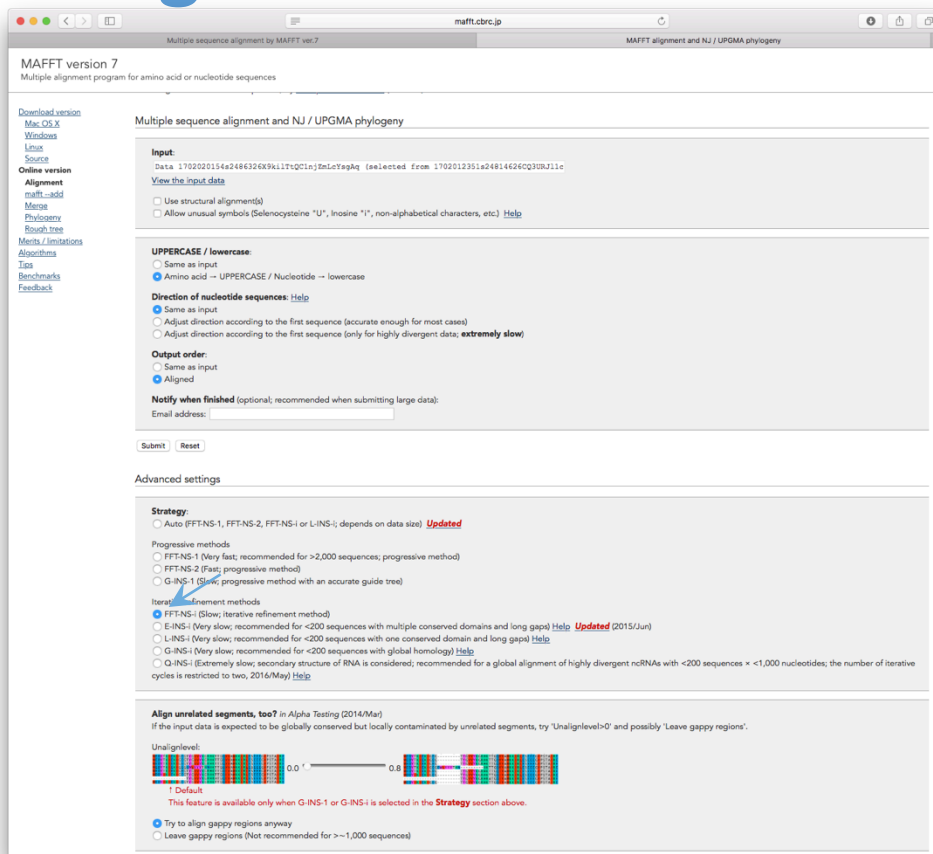
- 390 1 BAC\_FTA\_CC\_emb\_CNH06518.1 cluster21 [Bacteroides fragilis NCTC 9343]
- 392 2 FRB\_RDM\_gb\_ADE81715.1 cluster19 [Prevotella ruminicola 23]
- 392 3 BAC\_VIG\_gb\_ABR38290.1 cluster20 [Bacteroides vulgatus ATCC 8482]
- 400 4 RH0\_P042\_gb\_ACY4919.1 cluster12 [Rhodothermus marinus DSM 4232]
- 393 5 PAR\_MUC\_gb\_AFC29036.1 cluster17 [Paenibacillus muclilaginosus 3016]
- 396 6 CLO\_CES\_gb\_ACL17297.1 cluster14 [Clostridium cellulovorans H10]
- 407 7 CLO\_CES\_gb\_ACL52883.1 cluster11 [Clostridium cellulovorans 7438]
- 400 8 BAC\_CES\_gb\_AJ028565.1 cluster13 [Bacillus cellululositicus DSM 2522]
- 393 9 CEE\_MIX\_CC\_gb\_AAS19693.1 cluster16 [Cellivibrio mixtus]
- 392 10 TER\_TUR\_gb\_ACR12108.1 cluster18 [Teredinibacter turnerae T7901]
- 396 11 PED\_HEP\_gb\_AJ008175.1 cluster15 [Pedobacter heparinus DSM 2366]
- 386 12 RUM\_ALB\_CC\_gb\_AD021379.1 cluster24 [Ruminococcus albus 7 = DSM 20455]
- 389 13 HAL\_RYD\_gb\_AEE47952.1 cluster22 [Haloscomenobacter hydroxilis DSM 1100]
- 386 14 BAC\_OVA\_CC\_gb\_EC010988.1 cluster21 [Bacteroides ovatus ATCC 8483]
- 383 15 BAC\_TES\_CC\_gb\_AJ078885.1 cluster26 [Bacteroides thetaiotaomicron VPI-5482]
- 382 16 DTA\_FER\_CC\_gb\_ACT94389.1 cluster27 [Dyadobacter fermentans DSM 18053]

## Subfamilies definition

- Example : Sub\_family 1

- Select sequences
- Select 'Realign the selected sequences'

# Alignement des sous-familles



MAFFT version 7  
Multiple alignment program for amino acid or nucleotide sequences

Multiple sequence alignment and NJ / UPGMA phylogeny

**Input:**  
Data: 1702020134e2486326X9k11TtQC1eJEnLeYagk (selected from 1702012351a248144264029URJ11c)  
[View the input data](#)

Use structural alignment(s)  
 Allow unusual symbols (Selenocysteine "U", Inosine "I", non-alphabetical characters, etc.) [Help](#)

**UPPERCASE / lowercase**  
 Same as input  
 Amino acid → UPPERCASE / Nucleotide → lowercase

**Direction of nucleotide sequences:** [Help](#)  
 Same as input  
 Adjust direction according to the first sequence (accurate enough for most cases)  
 Adjust direction according to the first sequence (only for highly divergent data; **extremely slow**)

**Output order:**  
 Same as input  
 Aligned

**Notify when finished** (optional; recommended when submitting large data):  
Email address:

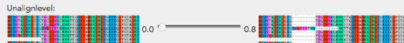
**Advanced settings**

**Strategy:**  
 Auto (FFT-NS-1, FFT-NS-2, FFT-NS-I or L-INS-i; depends on data size) **Updated**

**Progressive methods**  
 FFT-NS-1 (Very fast; recommended for >2,000 sequences; progressive method)  
 FFT-NS-2 (Fast; progressive method)  
 G-INS-1 (Slow; progressive method with an accurate guide tree)

**Iterative refinement methods**  
 FFT-NS-I (Slow; iterative refinement method)  
 E-INS-i (Very slow; recommended for <200 sequences with multiple conserved domains and long gaps) [Help](#) **Updated (2015/Jun)**  
 L-INS-i (Very slow; recommended for <200 sequences with one conserved domain and long gaps) [Help](#)  
 G-INS-i (Very slow; recommended for <200 sequences with global homology) [Help](#)  
 Q-INS-i (Extremely slow; secondary structure of RNA is considered; recommended for a global alignment of highly divergent ncRNAs with <200 sequences × <1,000 nucleotides; the number of iterative cycles is restricted to two, 2016/May) [Help](#)

**Align unrelated segments, too?** in Alpha Testing (2014/Mar)  
If the input data is expected to be globally conserved but locally contaminated by unrelated segments, try 'Unalignlevel>0' and possibly 'Leave gappy regions'.

Unalignlevel:  
  
1 Default  
This feature is available only when G-INS-1 or G-INS-i is selected in the Strategy section above.

Try to align gappy regions anyway  
 Leave gappy regions (Not recommended for >~1,000 sequence)

Subfamily 1 alignment

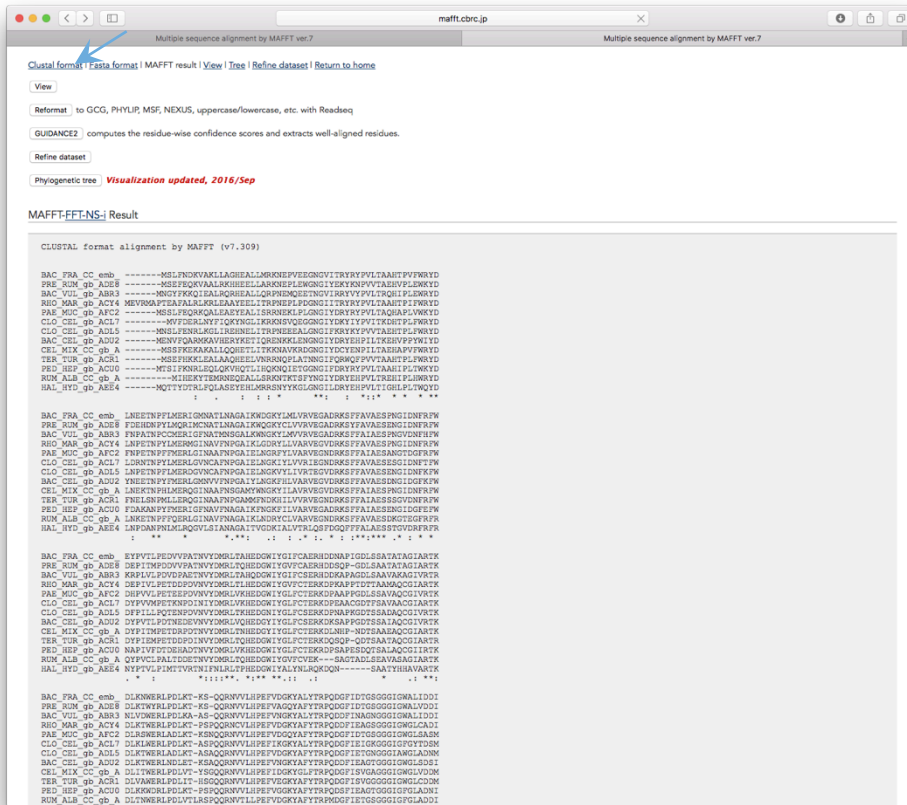
- Example : Sub\_family 1

Strategy :

Select "FFT-NS-I iterative refinement method "



# Alignement des sous-familles



Multiple sequence alignment by MAFFT ver.7

Clustal format | Fasta format | MAFFT result | View | Tree | Refine dataset | Return to home

View

Reformat to GCG, PHYLIP, MSF, NEXUS, uppercase/lowercase, etc. with Readseq

GUIDANCE2 computes the residue-wise confidence scores and extracts well-aligned residues.

Refine dataset

Phylogenetic tree Visualization updated, 2016/Sep

MAFFT-FFT-NS-i Result

```

CLUSTAL format alignment by MAFFT (v7.309)

BAC_FRA_cc emb -----MSLFDKRVAKLAGEHALMKRKEPVEEGNVIYTRYVYVPLTAAITPFWRYD
PRE_RUM_gb_ADE8 -----MSEFDKRVAKLAGEHALMKRKEPVEEGNVIYTRYVYVPLTAAITPFWRYD
BAC_VUL_gb_ABR3 -----MNDYFKQIEALQGRHALLQRFPMQDETNQVIRRYVYVPLTQRIPLERWYD
RHO_MAR_gb_ACT4 MEVWNAFZAFALALQKLEALAEELIETRESEIPLFQSHIITRYVYVPLTAAITPFWRYD
PAE_MOC_gb_AFC2 -----MSSLFDQRQALEAEYTEALISSRNEKPLGNGIYDRYVYVPLTAQAAPLWRYD
CLO_CEL_gb_ACI1 -----MVTDEKLNLYIQKYNGLKRNMSVDEGNGIYDRYIYVPIYKCHTVPFWRYD
CLO_CEL_gb_ADU2 -----MENVFAQRKKAVERFETETIQRENKLENGNGIYDRYRHPILTKRHPVFPYIYD
CEL_MIX_cc_gb_A -----MSEFDKRVAKLAGEHALMKRKEPVEEGNVIYTRYVYVPLTAAITPFWRYD
TER_TUR_gb_ACR1 -----MSEFHKKLEALAAQHEILVNSRNLATNNGIIFQRWQFPVYVYVPLTAAITPFWRYD
PED_REF_gb_ACH9 -----MSTFRNKLELQKQVGLIHLQKQIETGONGIYDRYVYVPLTAAITPFWRYD
RUM_ALB_cc_gb_A -----MHEETFRHGGKQALSRQNTSTFNGYVYVPLTAAITPFWRYD
HAL_HTD_gb_AE4 -----MOTTYDPLQLASYTEHLMKRSNTYKGLGNGILDRYRHPILTKRHPILTWQYD
: * * * * * :
: * * * * * :

BAC_FRA_cc emb IAEENFFMLERIDGNATINAGALFEGQKTIAMLVVVEADKRSFFVAASEGNGINFRFW
PRE_RUM_gb_ADE8 FDEHDFYIMQRIMCNATINAGALFKWQKICLVVVEADKRSFFVAASEGNGINFRFW
BAC_VUL_gb_ABR3 FWFATNPKDCEIETQATNSGALKRQKIMVYVVEADKRSFFVAASEGNGINFRFW
RHO_MAR_gb_ACT4 LRFENFPLMERIDGNATINAGALFKGQKTIAMLVVVEADKRSFFVAASEGNGINFRFW
PAE_MOC_gb_AFC2 FWFETNFFMERLGINAAMPFGALINRQFVLAVVEGNDKRSFFVAASEGNGITGRFR
CLO_CEL_gb_ACI1 LGQNFNPLMERIDGNATINAGALFKQKTIAMLVVVEADKRSFFVAASEGNGINFRFW
CLO_CEL_gb_ADU2 IWFENFFMLERIDGNATINAGALINRQFVLAVVEGNDKRSFFVAASEGNGINFRFW
CEL_MIX_cc_gb_A LRFENFPLMERIDGNATINAGALFKQKTIAMLVVVEADKRSFFVAASEGNGITGRFR
TER_TUR_gb_ACR1 FRELINPMLERQGINAAMPFGALINRQFVLAVVEGNDKRSFFVAASEGNGITGRFR
PED_REF_gb_ACH9 FQKAMPYFMRIGTNAVYVVEADKRSFFVAASEGNGITGRFR
RUM_ALB_cc_gb_A LKKEINFFQERLGINAVYVVEADKRSFFVAASEGNGITGRFR
HAL_HTD_gb_AE4 IWFQANPMLRQGVISINAGALFEGQKTIAMLVVVEADKRSFFVAASEGNGITGRFR
: * * * * * :
: * * * * * :

BAC_FRA_cc emb EYFVTLFEDVFAVAVYVQMLTAEHDGQYVGFCAERDDNAPIGDLSAATAGIARTK
PRE_RUM_gb_ADE8 DEPIITPDQVFAVAVYVQMLTAEHDGQYVGFCAERDDSDP-GDLSAATAGIARTK
BAC_VUL_gb_ABR3 KRPLVLPDQVFAVAVYVQMLTAEHDGQYVGFCAERDDKAPAGDLSAATAGIARTK
RHO_MAR_gb_ACT4 DEPIVLPDQVFAVAVYVQMLTAEHDGQYVGFCAERDDKAPAGDLSAATAGIARTK
PAE_MOC_gb_AFC2 DHPVLPDQVFAVAVYVQMLTAEHDGQYVGFCAERDDKAPAGDLSAATAGIARTK
CLO_CEL_gb_ACI1 DYPVQVYVQVQVYVQMLTAEHDGQYVGFCAERDDKAPAGDLSAATAGIARTK
CLO_CEL_gb_ADU2 DFYFLLPQVFAVAVYVQMLTAEHDGQYVGFCAERDDKAPAGDLSAATAGIARTK
CEL_MIX_cc_gb_A DYPVQVYVQVQVYVQMLTAEHDGQYVGFCAERDDKAPAGDLSAATAGIARTK
TER_TUR_gb_ACR1 DYPVQVYVQVQVYVQMLTAEHDGQYVGFCAERDDKAPAGDLSAATAGIARTK
PED_REF_gb_ACH9 NAFVYVQVYVQVQVYVQMLTAEHDGQYVGFCAERDDKAPAGDLSAATAGIARTK
RUM_ALB_cc_gb_A QYVFLPALDEYVQVYVQMLTAEHDGQYVGFCAERDDKAPAGDLSAATAGIARTK
HAL_HTD_gb_AE4 NTFVFLPMTVQVYVQMLTAEHDGQYVGFCAERDDKAPAGDLSAATAGIARTK
: * * * * * :
: * * * * * :

BAC_FRA_cc emb DLKNWRLPDILK-KS-QQRNVVLPFVQGYAYVTRPDQGITDGGGGIGMALDID
PRE_RUM_gb_ADE8 DLKTVWRLPDILK-KS-QQRNVVLPFVQGYAYVTRPDQGITDGGGGIGMALDID
BAC_VUL_gb_ABR3 NLVWRLPDILK-KS-QQRNVVLPFVQGYAYVTRPDQGITDGGGGIGMALDID
RHO_MAR_gb_ACT4 DLKTVWRLPDILK-PSPQQRNVVLPFVQGYAYVTRPDQGITDGGGGIGMALDID
PAE_MOC_gb_AFC2 DLKTVWRLPDILK-KSQQRNVVLPFVQGYAYVTRPDQGITDGGGGIGMALDID
CLO_CEL_gb_ACI1 DLKTVWRLPDILK-ASQQRNVVLPFVQGYAYVTRPDQGITDGGGGIGMALDID
CLO_CEL_gb_ADU2 DLKTVWRLPDILK-ASQQRNVVLPFVQGYAYVTRPDQGITDGGGGIGMALDID
CEL_MIX_cc_gb_A DLTVWRLPDILV-YSQQRNVVLPFVQGYAYVTRPDQGITDGGGGIGMALDID
TER_TUR_gb_ACR1 DLKTVWRLPDILK-DSQQRNVVLPFVQGYAYVTRPDQGITDGGGGIGMALDID
PED_REF_gb_ACH9 DLKTVWRLPDILK-PSPQQRNVVLPFVQGYAYVTRPDQGITDGGGGIGMALDID
RUM_ALB_cc_gb_A DLTVWRLPDILVLSQQRNVVLPFVQGYAYVTRPDQGITDGGGGIGMALDID

```

Subfamily 1 alignment

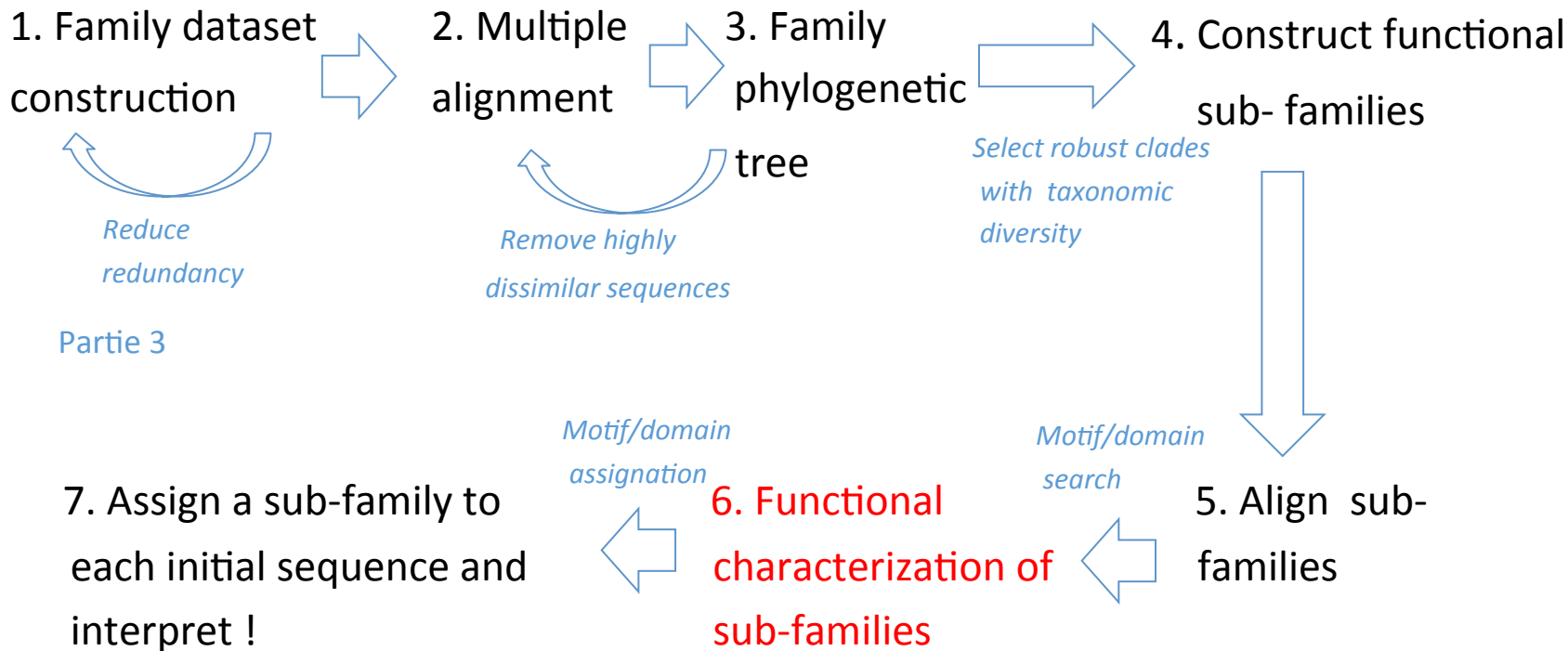
- Example : Sub\_family 1 : sequence 1 to 13

Export alignment of subfamily 1 in clustal format

Same for subfamilies 2 and 3!

# Motifs dans les séquences

# Et une démarche bioinfo générique



## Partie 3 - la stratégie utilisée:

*Material and Methods - Mewis et al. App. And Env. Microb. 2016*

- Hidden Markov models (HMMs) were created for each subfamily, as well as for the complete GH43 subfamily, by using HMMer3 **Hmmer2 sans ligne de commande**
- All GH43 sequences were compared to these HMMs by use of HMMer3 to assign a subfamily to each sequence. Each sequence was compared to all FIG 1 Phylogenetic tree of the GH43 family, showing the 37 subfamilies as colored branches. **Hmmer2 sans ligne de commande**

Ajout de contenu :

- détecter *de novo* des motifs et les rechercher dans d'autres séquences **MEME, hmmer ...**
- Annoter une séquence avec des bases de données de motifs et de domaines connus **Interpro**

# Recherche de profils de novo : MEME

## Travaux pratiques

<http://meme-suite.org/tools/meme>

<http://alternate.meme-suite.org/>

• **Exercice 1 : Rechercher des motifs conservés dans le jeu de données** <http://genoweb.toulouse.inra.fr/~formation/AlignClusteringLISBP/data/GH130InputRenamed.fasta>.

*Paramètres : laisser les paramètres par défauts sauf monter à **13** le nombre de motifs recherché.*

***Cela va être un peu long...***

# A quoi ça sert ?

- **rechercher des éléments fonctionnels**
- rechercher des sites de restriction
- **déterminer la signature d'une famille ou d'une sous-famille protéique pour caractériser ses membres**

# A quoi ça sert ?

- Exemples : la TATA box dans les séquences promotrices des gènes.
- Chez les protéines les motifs sont constitués de résidus conservés pas nécessairement consécutifs. Ex dans prosite : FN2\_2, [PS51092](#); Fibronectin type-II collagen-binding domain profile

```

BSPH1_HUMAN/40-84
BSPH1_HUMAN/85-132
BSPH1_MOUSE/40-84
BSPH1_MOUSE/85-133
BSPH2_MOUSE/35-79
BSPH2_MOUSE/80-128
ESPB1_CANLF/46-90
ESPB1_CANLF/91-139
ESPB1_CANLF/146-192
ESPB1_CANLF/199-245
ESPB1_HUMAN/24-68
ESPB1_HUMAN/69-117
ESPB1_HUMAN/124-170
ESPB1_HUMAN/177-223
ESPB1_PIG/24-68
ESPB1_PIG/69-117
ESPB1_PIG/124-170
ESPB1_PIG/177-223
FA12_BOVIN/42-90
FA12_CAVPO/41-89
FA12_HUMAN/42-90
FA12_MOUSE/42-90

```

```

VTDGEVFPFH KNGTYDGIKSKA--RHKNGLNKTVEG--YKKGCSA
EDFANGVFPFW RRLIYWECTDD EAFGKKNCSL KFNKDRINKYCE-
TEDGAVFPFL RSEIFYDCVNFNL--KHKNGLNKTVEG--YKKGCSA
SDYAPCAFPFW RHMIYWDCTED EVFGKKNCSL PNNKQVQKYGIE
ISTDSCVFPFV ADGFHYSCISLHS--DYDNGSLDFQVQ--RRIYGT
QDPPKCIFFFO KQKLIKKTKE YILNRSNCSL TENYNDGKIKQCS
DQKDSVFPFV KGSYFSCIKTNS--FSPNCA TRAVYNG--QKKGCSA
DDYPRCIFPFI RKGSHNSCITE SFLRRLNCSV SSF DENQKYGCT
SFSKPCIFPSI RNSTIFECMEDE--NKLNCPT ENMDEGKISLQAD
VPGFPCHPFS KNKNYNGIGK TKENLWGCATSYN DQDHTVYV--
GMHEECVFPFT KGSVYFTGTHIHS--LSPNCA TRAVYNG--QKKGCS
EDYPRCIFPFI RKGAYNSCISQ SFLGSLNCSV SVF DEKQKYGCT
SLRKPCHPFI RNNVSDCMEDES--NKLNCPT ENMDKGGKISLQAD
VPGFPCHPFN KNKNYFNCTNE SKENLWGCATSYN DQDHTVYV--
DTKDSVFPFH KGFTYFSCIRTNS--LSPNCA TRAVYNG--QKKGCS
EDYPRCIFPFI RGRSHRNCIVE SFFGKLNCSV SSF DEKQKYGCT
SFSKPCIFPSI RNHIISECLEDES--NKLNCPT ENMDMGGKISLQAD
VPGFPCHPFN KNKNYFNCTNE SKENLWGCATSYN DQDHTVYV--
VTGEPCHPFO HRQLHHKGIHR RPPGPRPGAT PNF EKDQRYAYGLE
VTGEPCHPFO NRQLYHHGIHK RPPGPRPGAT PNF DQDQRYAYGLE
VTGEPCHPFO HRQLYHKGIHK RPPGPRPGAT PNF DQDQRYAYGLE
VDGRLCHPFO HRQLHHKGIHKRRPGRPGAT PNF DEDQRYAYGLE

```

# Comment les décrit-on ?



Exemple d'un motif dans des séquences de maltose binding proteins :

Yvfk_Bs	PT <b>P</b> NIPEMNEI <b>W</b>
Yvfk_Bs	PT <b>P</b> NIPEMAEV <b>W</b>
MalX_Sp	PL <b>P</b> NISQMSAV <b>W</b>
MalE_Sc	PR <b>P</b> ALPEYSSL <b>W</b>
MalE_Tm	PM <b>P</b> NVPEMAPV <b>W</b>
CymE_Ko	AM <b>P</b> SIPEMGYL <b>W</b>
MalE_Ea	IM <b>P</b> NIIPQMSAF <b>W</b>
MalE_Sy	IM <b>P</b> NIIPQMSAF <b>W</b>
MalE_Ec	IM <b>P</b> NIIPQMSAF <b>W</b>
consens :	PM <b>P</b> NIPEMSAX <b>W</b>

Expression régulière ou patterns (format Prosite)

Syntaxe:

- : séparation des éléments

x : n'importe quel acide aminé

(i,j) : nombre d'occurrences entre i et j avec  $i < j$

[NHG] : alternative entre N H et G pour un même site

PROSITE: [PAI]-[TLRM]-P-[NAS]-[ILV]-[PS]-[EQ]-[MY]-[NASG]...



# Comment les décrit-on ?



Exemple d'un motif dans des séquences de maltose binding proteins :

Yvfk_Bs	PT <b>P</b> NIPEMNEI <b>W</b>
Yvfk_Bs	PT <b>P</b> NIPEMAEV <b>W</b>
MalX_Sp	PL <b>P</b> NISQMSAV <b>W</b>
MalE_Sc	PR <b>P</b> ALPEYSSL <b>W</b>
MalE_Tm	PM <b>P</b> NVPEMAPV <b>W</b>
CymE_Ko	AM <b>P</b> SIPEMGYL <b>W</b>
MalE_Ea	IM <b>P</b> NIPOMSAF <b>W</b>
MalE_Sy	IM <b>P</b> NIPOMSAF <b>W</b>
MalE_Ec	IM <b>P</b> NIPOMSAF <b>W</b>
consensus :	PM <b>P</b> NIPEMSAX <b>W</b>

**Convient bien pour des motifs très conservés**

Pour tenir compte de la répartition des différents acides aminés à chaque position : profils ou matrices de fréquences

PROSITE: [PAI] - [TLRM] - P - [NAS] - [ILV] - [PS] - [EQ] - [MY] - [NASG] . . .

# Comment les décrit-on ?

Les matrices de fréquences ou de probabilité ou profils :

AATAGTCGC

GGTAGTCTA

ATTAGTCGA

GCTAGTCGG

position frequency matrix (PFM)  
position probability matrix (PPM)

le profil correspondant :

	1	2	3	4	5	6	7	8	9
A:	0.5	0.25	0	1	0	0	0	0	0.5
T:	0	0.25	1	0	0	1	0	0.25	0
G:	0.5	0.25	0	0	1	0	0	0.75	0.25
C:	0	0.25	0	0	0	0	1	0	0.25

# Comment les décrit-on ?

**Les matrices de fréquences ou de probabilités ou profils :**

	1	2	3	4	5	6	7	8	9
A:	0.5	0.25	0	1	0	0	0	0	0.5
T:	0	0.25	1	0	0	1	0	0.25	0
G:	0.5	0.25	0	0	1	0	0	0.75	0.25
C:	0	0.25	0	0	0	0	1	0	0.25

Ces matrices de fréquences peuvent être modifiées pour tenir compte de la fréquence de chaque résidus : elles deviennent ainsi des matrices de poids.

De plus une transformation  $(\log_2((f(x)+\epsilon)/f(b)))$  est très souvent utilisée pour enlever les 0 dans la matrice et donc autoriser de nouvelles séquences qui peuvent correspondre au profil.

# Comment les décrit-on ?

On obtient alors des pseudo-comptages.

Ces matrices sont appelées PWM : position weight matrix, PSWM : position-specific weight matrix ou PSSM : position-specific scoring matrix.

Le score d'un mot dans le modèle est la somme des poids puisqu'on est en log.

# Comment les décrit-on ?

## Les matrices de poids :

le profil correspondant :

	1	2	3	4	5	6	7	8	9
A:	0.5	0.25	0	1	0	0	0	0	0.5
T:	0	0.25	1	0	0	1	0	0.25	0
G:	0.5	0.25	0	0	1	0	0	0.75	0.25
C:	0	0.25	0	0	0	0	1	0	0.25

PWM : position weight matrix

PSWM : position-specific weight matrix

PSSM : position-specific scoring matrix

$\log_2((f(x) + 0.05) / 0.25)$

	1	2	3	4	5	6	7	8	9
A:	1.14	0.26	-2.32	2.07	-2.32	-2.32	-2.32	-2.32	1.14
T:	-2.32	0.26	2.07	-2.32	-2.32	2.07	-2.32	0.26	-2.32
G:	1.14	0.26	-2.32	-2.32	2.07	-2.32	-2.32	1.68	0.26
C:	-2.32	0.26	-2.32	-2.32	-2.32	-2.32	2.07	-2.32	0.26

# Score d'un mot dans ma séquence ?

Calcul d'un score d'une séquence :

	1	2	3	4	5	6	7	8	9
A:	1.14	0.26	-2.32	2.07	-2.32	-2.32	-2.32	-2.32	1.14
T:	-2.32	0.26	2.07	-2.32	-2.32	2.07	-2.32	0.26	-2.32
G:	1.14	0.26	-2.32	-2.32	2.07	-2.32	-2.32	1.68	0.26
C:	-2.32	0.26	-2.32	-2.32	-2.32	-2.32	2.07	-2.32	0.26

Score : CAGTGACCC ????

# Score d'un mot dans ma séquence ?

Calcul d'un score d'une séquence :

	1	2	3	4	5	6	7	8	9
A:	1.14	0.26	-2.32	2.07	-2.32	-2.32	-2.32	-2.32	1.14
T:	-2.32	0.26	2.07	-2.32	-2.32	2.07	-2.32	0.26	-2.32
G:	1.14	0.26	-2.32	-2.32	2.07	-2.32	-2.32	1.68	0.26
C:	-2.32	0.26	-2.32	-2.32	-2.32	-2.32	2.07	-2.32	0.26

Score : CAGTGACCC ????

$$-2.32 + 0.26 - 2.32 - 2.32 + 2.07 - 2.32 + 2.07 + 0.26 = -4.62$$

score > 0 : la séquence a plus de chance de correspondre au motif qu'une séquence aléatoire,

score < 0 : séquence plus probablement aléatoire

## Ce score est-il significatif ?

Il est nécessaire ensuite de comparer le score calculé aux scores max et min possibles avec le profil et de choisir le seuil de significativité.

⇒ pas de p-value.



## Autres raffinements possibles du profil

Il est possible de tenir compte des taux de substitution / conservation observées en générale dans les séquences protéiques (matrice PAM / BLOSUM) et de rajouter des indels (en rajoutant une colonne pour les gaps).

**Ce type de formalisme convient donc pour des motifs moins bien conservés et permet d'utiliser des connaissances biologiques pour améliorer le profil.**

# Recherche de profils de novo : MEME

Le programme MEME de la suite du même nom découvre des motifs non gappés dans un ensemble de séquence.

- Représente les motifs par une PPM (Position Probability Matrix) mais rajoute aussi l'information de la fréquence des nucléotides observées.
- Randomise les données en entrée (ou prend un set de données contrôle) pour trouver les motifs significativement enrichis dans les séquences.
- Plutôt pour de l'ADN (ex : promoteurs) mais certains logiciels fonctionnent aussi sur des séquences protéiques

# Recherche de profils de novo : MEME

Utilise :

- des méthodes statistiques (algorithme EM) pour choisir automatiquement la longueur, le nombre d'occurrences et la description de chaque motif.
- Un modèle de background (fréquence des résidus, des mots de taille 2, 3 selon le paramétrage...) pour calculer la vraisemblance du motif, sa significativité (e-value) et la PSPM.

# Recherche de profils de novo : MEME

## Les scores :

- Le **log likelihood ratio** (LLR) du motif =  $\log (\text{Pr}(\text{occurrence} \mid \text{motif}) / \text{Pr}(\text{occurrence} \mid \text{back}))$ . C'est une mesure de la différence du locus avec le modèle de background.
- **E-value** du motif : estimation du nombre de motif de même taille et ayant le même nombre d'occurrences qui aurait obtenu un LLR égal ou supérieur si les séquences avaient été obtenues aléatoirement selon le background modèle d'ordre 0.

# Recherche d'un ensemble de motifs : MAST

MAST recherche un ensemble de motifs non gappés dans des séquences non alignées.

- calcul de la p-value pour chaque occurrence de motif
- calcul de la p-value globale (produit des p-values des occurrences de chaque motifs dans la séquence)
- Les résultats sont triés sur la e-value globale (p-value corrigée pour les tests multiples).
- Les p-values sont calculées en utilisant un modèle de background obtenu à partir d'une vieille version de la base de données nr du NCBI.

[http://meme-suite.org/doc/mast.html?man\\_type=web](http://meme-suite.org/doc/mast.html?man_type=web)

# Travaux pratiques

## En attendant ....

<http://emboss-genotoul.toulouse.inra.fr/>

[http://meme-suite.org/doc/overview.html?man\\_type=web](http://meme-suite.org/doc/overview.html?man_type=web)

- **jetons un coup d'oeil aux différents outils de la suite MEME et au manuel**
- **Jetons aussi un coup d'oeil à la suite emboss.**

*En particulier les sections EDIT et PROTEINS PROFILE*

# Travaux pratiques

<http://meme-suite.org/tools/meme>

## •Exercice 1 : Rechercher des motifs conservés dans le jeu de données GH130InputRenamed.fasta

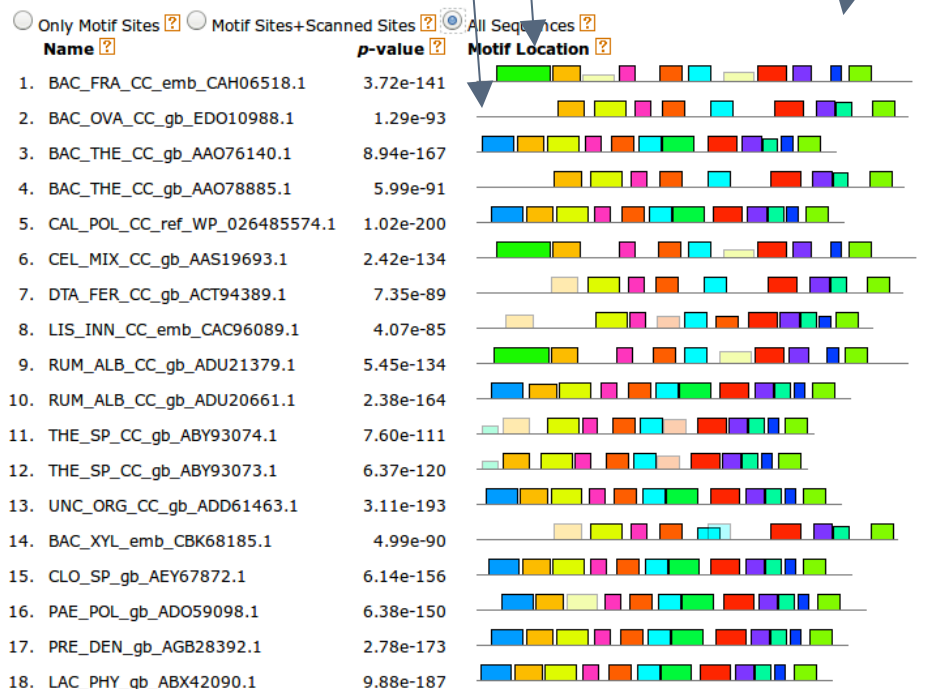
Utilisez les pages de résultats HTML de MEME et de MAST

Comment ces motifs sont ils répartis ?

- dans l'arbre ?
- selon les fonctions connues des protéines connues ?

Regardons les motifs 7, 11 et 13 de plus près

### MOTIF LOCATIONS



# Travaux pratiques

<http://meme-suite.org/tools/meme>; <http://emboss-genotoul.toulouse.inra.fr/>

## •Exercice 2 : télécharger et explorer ces motifs

**Faire 3 groupes : un sur le motif 7, un autre sur le 11 et le dernier sur le 13**

- 1) Télécharger le logo et le fasta des occurrences du motif
- 2) Passer dans emboss : Utiliser Prophecy pour générer une matrice de fréquence à partir de ce fichier. La sauvegarder en utilisant un nom explicite. Réutiliser Prophecy pour générer une matrice de poids de Grisbkov. La sauvegarder sous un autre nom explicite.
- 3) Utiliser Profit pour rechercher le motif sous la forme d'une matrice de fréquence simple dans les séquences *GH130InputRenamed.fasta* et Prophet pour chercher le motif sous forme de matrice de poids, augmenter la pénalité d'ouverture de gap à 3.
- 4) Comparer. Que remarquez-vous ? Choisir un seuil de significativité pour Prophet.



# Travaux pratiques

**En option : si vous avez le temps**

[http://www.cazy.org/GH130\\_characterized.html](http://www.cazy.org/GH130_characterized.html)

## •Exercice 2 : télécharger et explorer ces motifs

***Diviser la salle en 3 groupes : un sur le motif 7, un autre sur le 11 et le dernier sur le 13***

*5) Une des séquences présentant ces motifs a-t-elle une structure pdb ? Si oui, ouvrir une des fiches pdb correspondante et regarder où se situe le motif dans la structure secondaire et dans la structure tertiaire.*

# Comment les décrit-on ?

Il y a trois grands types de formalismes pour représenter un motifs :



# Comment les décrit-on ?

- Il y a trois grands types de formalismes pour représenter un motifs :
- les patterns (expression régulières ou format Prosite)
  - les profils ou matrices de fréquence ou de poids
  - les HMM

# Comment les décrit-on ? Les HMM

Les HMM : Chaînes de markov cachées

Basiquement, il s'agit d'un modèle statistique représentant le motif en utilisant la probabilité d'avoir un résidu après un autre.

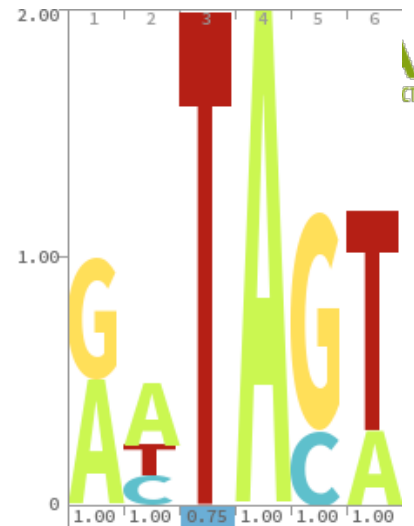
Ex :

AATACT

GA-AGT

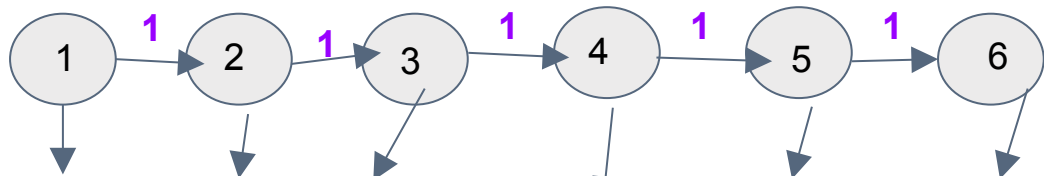
ATTAGA

GCTAGT



**Si pas de gap**

Probabilités de transition sans gap = 1



Probabilités d'émission

A = 0.5  
G = 0.5

A = 0.5  
T = 0.25  
C = 0.25

T = 0.75  
- = 0.25

A = 1

G = 0.75  
C = 0.25

T = 0.75  
A = 0.25

# Comment les décrit-on ? Les HMM

Les HMM : Chaînes de markov cachées

Basiquement, il s'agit d'un modèle statistique représentant le motif en utilisant la probabilité d'avoir un résidu après un autre.

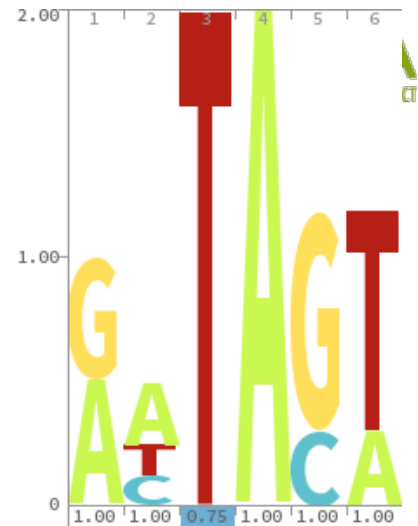
Ex :

AATACT

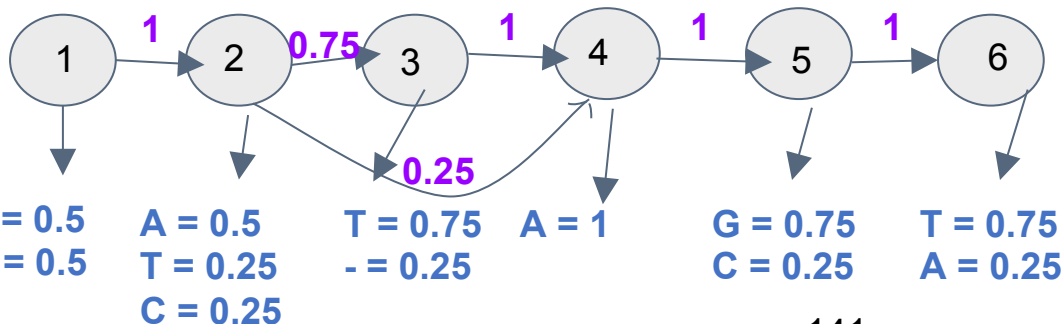
GA-AGT

ATTAGA

GCTAGT



Probabilités de transition sans gap = 1

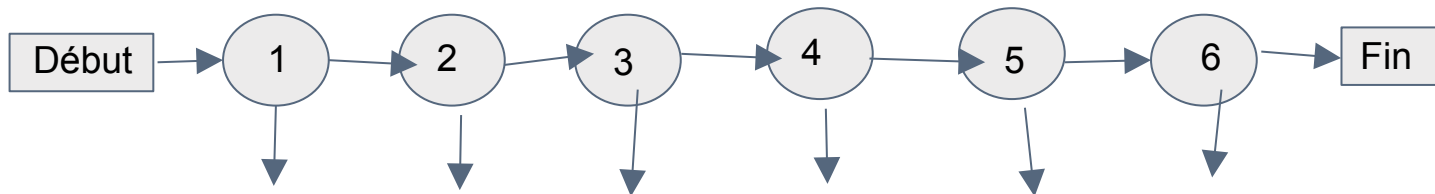


# Comment les décrit-on ? Les HMM

Les HMM : Chaînes de markov cachées  
 Si on généralise l'exemple précédent :

Types de probabilités de transitions :  
 match -> match  
 begin -> match  
 match -> end

Match / mismatch



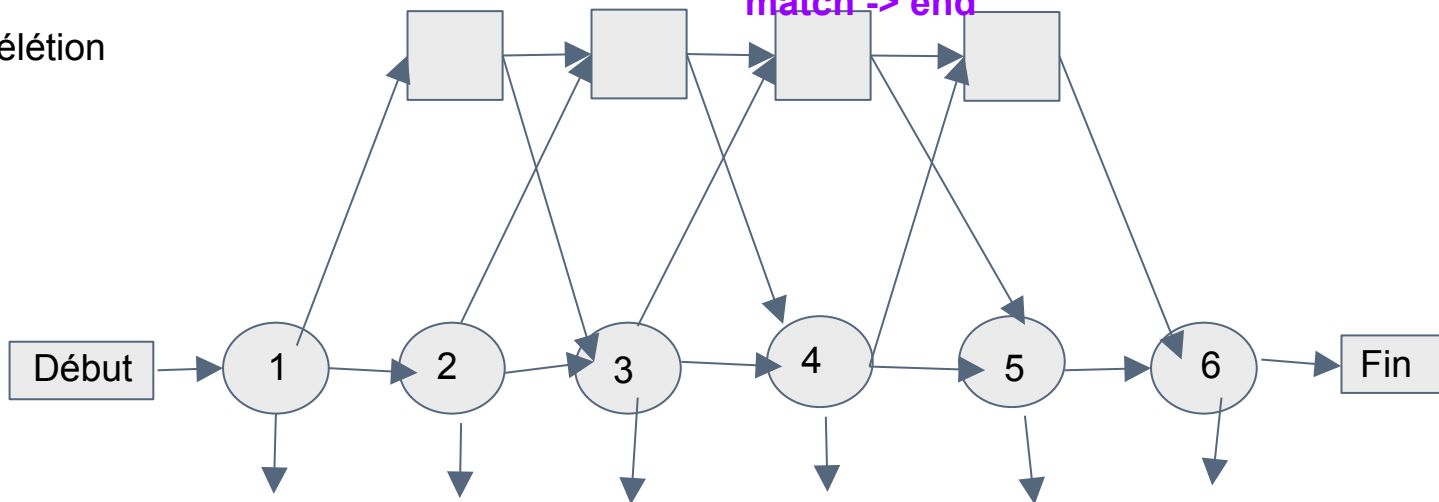
Probabilités d'émission

# Comment les décrit-on ? Les HMM

Les HMM : Chaînes de markov cachées  
 Si on généralise l'exemple précédent :

Types de probabilités de transitions :  
 match -> match  
 delete -  
 > delete  
 match -> delete  
 match  
 match -> end  
 begin ->

Événements de délétion



Probabilités d'émission

Match / mismatch

# Comment les décrit-on ? Les HMM

Les HMM : Chaînes de markov cachées  
 Si on généralise l'exemple précédent :

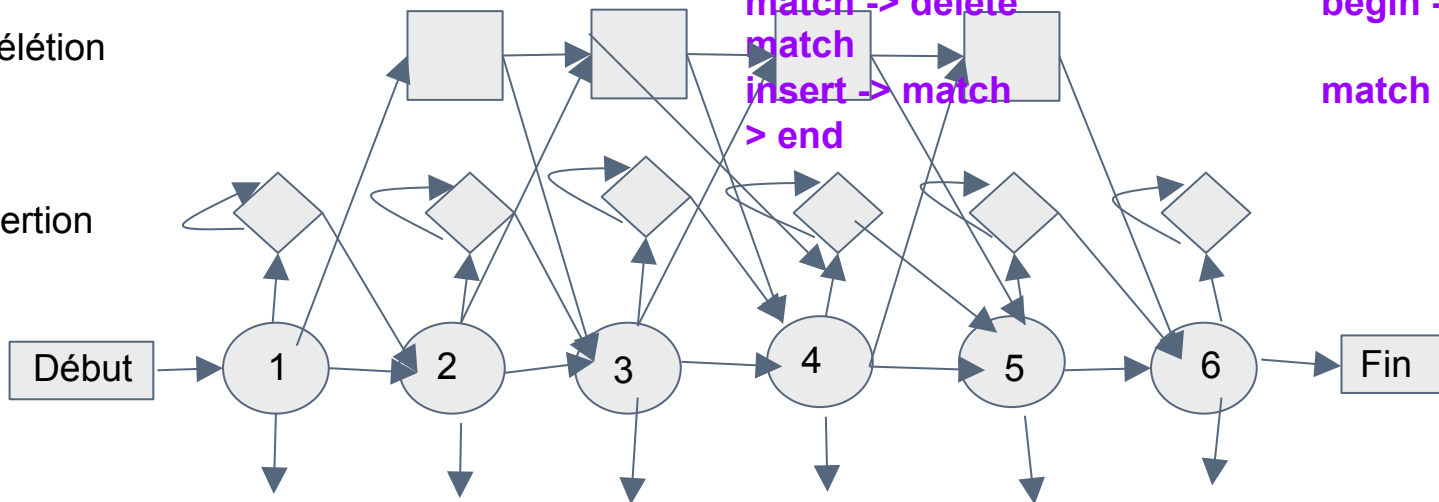
Types de probabilités de transitions :

- match -> match
- insert ->
- insert
- match -> insert
- delete -
- > delete
- match -> delete
- begin ->
- match
- insert -> match
- match -
- > end

Événements de délétion

Evénements d'insertion

Match / mismatch



Probabilités d'émission



# Comment les décrit-on ? Les HMM

Les HMM sont donc utiles pour :

- représenter des motifs complexes.
- permettre l'utilisation d'algorithmes utilisés dans des problèmes similaires.
- comme pour les profils, il est possible de rajouter de la “connaissance biologique” dans le motif en modifiant les probabilités par exemple en étant plus tolérant aux substitutions entre acides aminés fréquemment rencontrés dans les protéines en générale.
- retrouver des occurrences de motifs dans des séquences très divergentes.

# Comment les décrit-on ? Les HMM

- ATTENTION :

Un profil HMM peut être aligné sur une séquence globalement ou localement (seulement une partie du profil) mais cela est choisi lors de la construction du profil et non lors de son utilisation.

# Comment utilise-t-on les motifs HMM ?

## Ma séquence contient-elle ce motif ?

Pour chercher si une séquence contient un profil HMM connu, le programme va chercher le chemin le plus probable de la séquence à travers le modèle en utilisant les probabilités de transition et d'émission pour calculer le score ( $S = P_{\text{Obs}} / P_{\text{attendu}}$ ).

# Comment utilise-t-on les motifs HMM ?

A partir du score  $S$ , et grâce au modèle probabiliste, le logiciel calcule la p-value.

**p-value** : probabilité d'obtenir par chance un score au moins égal au score  $S$ .

# Comment utilise-t-on les motifs HMM ?

**Quelles sont les séquences contenant mon profil HMM dans cet ensemble de séquences (database) ?**

- la p-value doit être corrigée pour les tests multiples  $\Rightarrow$  e-value
- la e-value est le nombre de hits qu'on aurait obtenus avec un aussi bon score contre une database de même taille contenant des séquences aléatoires.
- le score est un score log transformé appelé aussi log odds score ou bit score. Il ne dépend pas de la taille de la database mais seulement du modèle et de la séquence cible.

# Pourquoi utilise-t-on les motifs HMM ?

- **Créer un profil HMM représentatif** d'une famille de protéines d'intérêt ou d'un domaine protéique d'intérêt. Puis utiliser ce profil pour rechercher d'autres séquences appartenant à cette famille ou contenant ce domaine. ⇒ hmmerBuild, HmmerSearch (Hmmer2), hmmbuild, hmmsearch (hmmer3)
- **Annoter une séquence** : chercher des motifs connus ou des domaines fonctionnels qu'on ne trouverait pas avec un blast standard car trop peu similaires ⇒ Interpro, hmmerPfam (hmmer2), hmmScan (hmmer3)
- **Construire un alignement multiple** très rapidement en utilisant le profil HMM comme graine. ⇒ hmmerAlign (hmmer2), hmmlalign (hmmer3)

# Travaux pratiques

[http://sequenceconversion.bugaco.com/convert/biology/sequences/fasta\\_to\\_stockholm.php](http://sequenceconversion.bugaco.com/convert/biology/sequences/fasta_to_stockholm.php)

## •Exercice 3 : caractérisons nos sous-familles par un profil HMM

*Diviser la salle en 3 : un sur chaque sous-famille*

1) *Convertir l'alignement multiple clustal en alignement au format stockholm*

*Si besoin : <http://genoweb.toulouse.inra.fr/~formation/AlignClusteringLISBP/data/SubFamilies>*

# Travaux pratiques

<http://emboss-genotoul.toulouse.inra.fr/>

## •Exercice 3 : caractérisons nos sous-familles par un profil HMM

*Diviser la salle en 3 : un sur chaque sous-famille*

*2) Construire un profil HMM global représentant la sous-famille*

*Utiliser en entrée l'alignement multiple en format stockholm*

*Attention : mettre obligatoirement un nom pour le HMM*

*Laisser le reste par défaut*

*Sauvegarder le fichier hmmfile avec un nom explicite et jeter un coup d'oeil au fichier.*

*Que sont les lignes et les colonnes de la matrice ?*



# Travaux pratiques

<http://skylign.org>

• **Exercice 3 : caractérisons nos sous-familles par un profil HMM**

*Diviser la salle en 3 : un sur chaque sous-famille*

3) *Générer un logo pour ce profil HMM, le sauvegarder avec un nom explicite.*

# Travaux pratiques

<http://emboss-genotoul.toulouse.inra.fr/>

## •Exercice 3 : caractérisons nos sous-familles par un profil HMM

***Diviser la salle en 3 : un sur chaque sous-famille***

4) Calibrer le profil HMM avec ehmmcalibrate.

*Ceci permet d'augmenter la sensibilité de la recherche de ce motif contre un ensemble de séquence. Pour cela il génère un grand nombre de séquences aléatoires et compare le profil HMM original contre elles. Il calcule ainsi des paramètres de valeurs extrêmes qu'il écrit dans le profil HMM modifié qu'il génère.*

*Sauvegarder le fichier outhmmfile avec un nom explicite*

# Travaux pratiques

<http://emboss-genotoul.toulouse.inra.fr/>

• **Exercice 4** : recherchons les membres de la sous-famille dans le dataset d'origine (GH130InputRenamed.fasta) en utilisant le motif HMM.

*Diviser la salle en 3 : un sur chaque sous-famille*

*5) utiliser le bon programme pour rechercher le motif dans l'ensemble des séquences, en utilisant les paramètres par défaut.*

*Pourquoi y a-t-il deux sections de résultats ?*

*Qu'en concluez-vous ?*

# Travaux pratiques

**En option : si vous avez le temps**

<https://www.ebi.ac.uk/Tools/hmmer/search/hmmsearch>

• **Exercice 5 : utilisons hmmer3 pour rechercher des protéines similaires dans des banques**

*Diviser la salle en 3 : un sur chaque sous-famille*

*1) utiliser l'alignement de la sous-famille en entrée  
faire la recherche contre "swissprot",*

*dans les paramètres avancés, customised results cliquer sur tout.*

*Explorer les résultats. Comparer les scores avec ceux que nous avons obtenus.*

*N'oubliez pas la possibilité de cliquer sur ">".*

# Travaux pratiques

**En option : si vous avez le temps**

<https://www.ebi.ac.uk/Tools/hmmer/search/hmmsearch>

**•Exercice 5 : utilisons hmmer3 pour rechercher des protéines similaires dans des banques**

*Diviser la salle en 3 : un sur chaque sous-famille*

*2) utiliser l'alignement de la sous-famille en entrée*

*faire la même chose que précédemment contre la banque "ensembl genome bacteria"*

*Que remarquez-vous ?*

# Petit bilan intermédiaire

OK

- **Créer un profil HMM représentatif** d'une famille de protéines d'intérêt ou d'un domaine protéique d'intérêt. Puis utiliser ce profil pour rechercher d'autres séquences appartenant à cette famille ou contenant ce domaine. ⇒ hmmerBuild, HmmerSearch (Hmmer2), hmmbuild, hmmsearch (hmmer3)

Conseil : Penser à enrichir le modèle avec les nouvelles protéines trouvées (processus itératif)

- **Annoter une séquence** : chercher des motifs connus ou des domaines fonctionnels qu'on ne trouverait pas avec un blast standard car trop peu similaires ⇒ Interpro, hmmerPfam (hmmer2), hmmScan (hmmer3)

# Analyse fonctionnelle et structurale de séquences protéiques

Interpro : <http://www.ebi.ac.uk/interpro/>

- classifie en famille (si elles sont connues)
- prédit la présence de domaine (au sens de interpro : unité structurale et / ou fonctionnelle)
- de sites annotés comme importants (site actif, site de liaison...)
- et de domaines structuraux (structure super secondaire comme l'hélice - boucle - hélice des facteurs de transcription)
- pour cela il intègre de nombreux outils et base de données.

# Travaux pratiques

<http://www.ebi.ac.uk/interpro/>

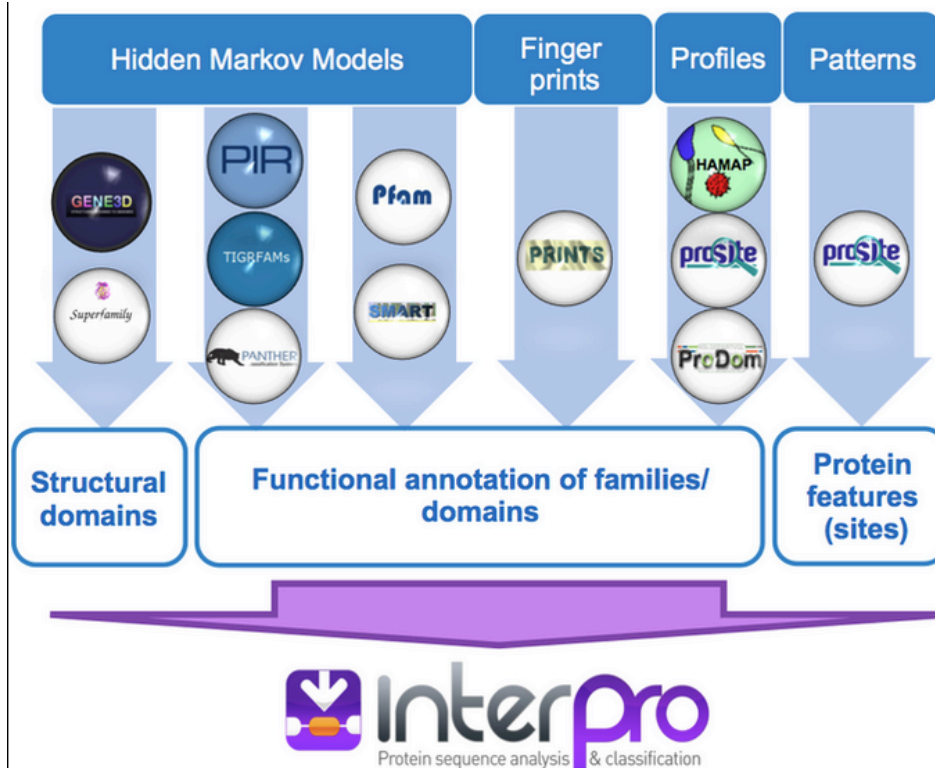
- **Exercice 7 : utiliser Interpro**

*1) Choisir la séquence de votre choix dans la sous-famille étudiée pour interroger interpro*

***Cela va être un peu long...***

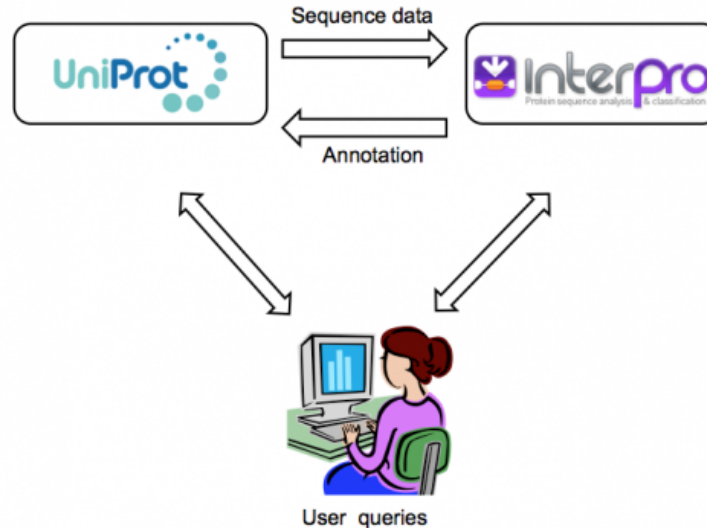


# Analyse fonctionnelle et structurale de séquences protéiques



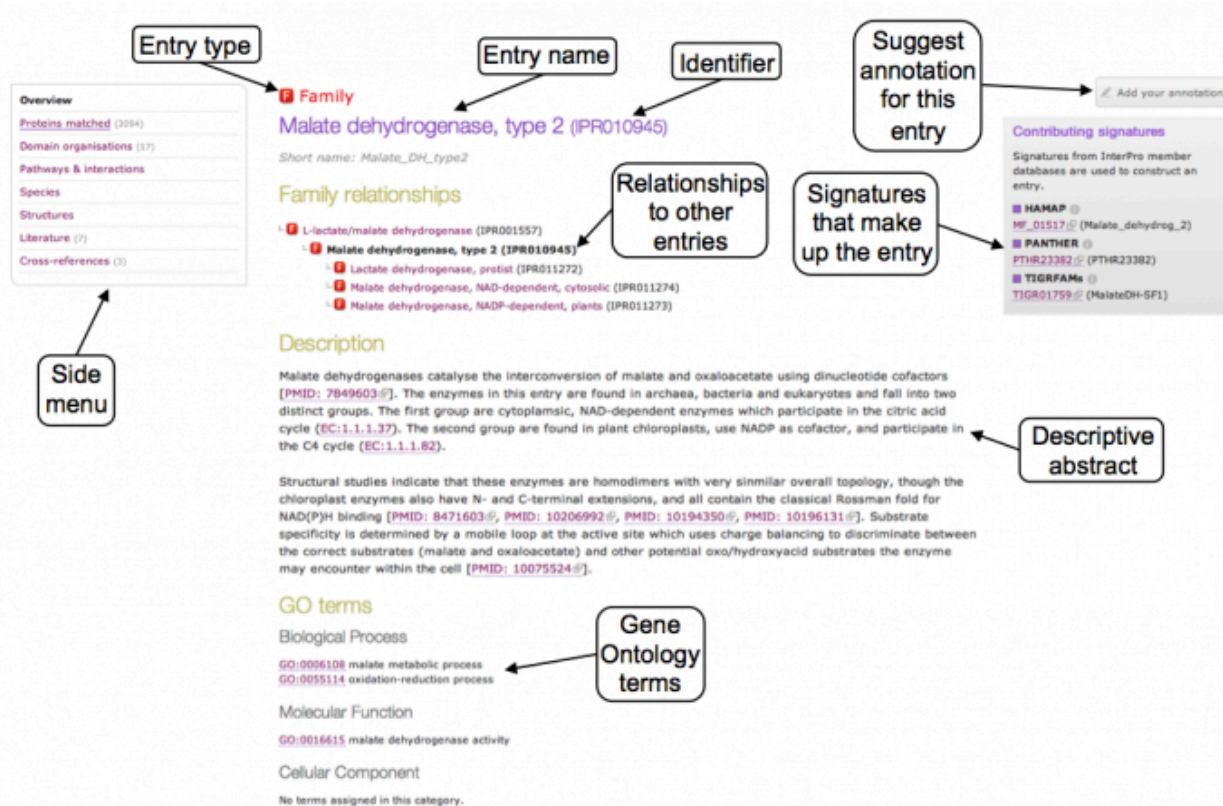
# Analyse fonctionnelle et structurale de séquences protéiques

La base de donnée UniproKB est régulièrement analysée avec Interpro pour en améliorer l'annotation.



# Analyse fonctionnelle et structurale de séquences protéiques

Intégration de différentes informations caractéristiques de la même famille protéique, du même domaine ou de la même caractéristique de séquence.  
 Ex de la fiche Interpro de la famille des malate déshydrogénase de type 2



The screenshot shows the InterPro entry for Malate dehydrogenase, type 2 (IPR010945). Annotations highlight the following elements:

- Entry type:** Family
- Entry name:** Malate dehydrogenase, type 2 (IPR010945)
- Identifier:** IPR010945
- Side menu:** Overview, Proteins matched (2094), Domain organisations (17), Pathways & Interactions, Species, Structures, Literature (7), Cross-references (3)
- Relationships to other entries:** L-lactate/malate dehydrogenase (IPR01557), Malate dehydrogenase, type 2 (IPR010945), Lactate dehydrogenase, protist (IPR01272), Malate dehydrogenase, NAD-dependent, cytosolic (IPR01274), Malate dehydrogenase, NADP-dependent, plants (IPR01273)
- Suggest annotation for this entry:** Add your annotation
- Signatures that make up the entry:** HAMAP (MF\_01517 (Malate\_dehydrog\_2)), PANTHER (PTHR23382 (PTHR23382)), TIGRFAMs (TIGR01759 (MalateDH-SF1))
- Descriptive abstract:** Malate dehydrogenases catalyse the interconversion of malate and oxaloacetate using dinucleotide cofactors [PMID: 7849603]. The enzymes in this entry are found in archaea, bacteria and eukaryotes and fall into two distinct groups. The first group are cytoplasmic, NAD-dependent enzymes which participate in the citric acid cycle (EC:1.1.1.37). The second group are found in plant chloroplasts, use NADP as cofactor, and participate in the C4 cycle (EC:1.1.1.82).
- Gene Ontology terms:** Biological Process: GO:0006108 malate metabolic process, GO:0055114 oxidation-reduction process; Molecular Function: GO:0016615 malate dehydrogenase activity; Cellular Component: No terms assigned in this category.

# Analyse fonctionnelle et structurale de séquences protéiques

Les 4 catégories figurant dans Interpro



# Analyse fonctionnelle et structurale de séquences protéiques

Les hiérarchies entre les entrées interpro des familles et des domaines.  
Pas d'overlap entre les deux.

## **F** Family

Malate dehydrogenase, type 2 (IPR010945)

Short name: Malate\_DH\_type2

## Family relationships

- ↳ **F** L-lactate/malate dehydrogenase (IPR001557)
  - ↳ **F** Malate dehydrogenase, type 2 (IPR010945)
    - ↳ **F** Lactate dehydrogenase, protist (IPR011272)
    - ↳ **F** Malate dehydrogenase, NAD-dependent, cytosolic (IPR011274)
    - ↳ **F** Malate dehydrogenase, NADP-dependent, plants (IPR011273)

## **D** Domain

Death-like domain (IPR011029)

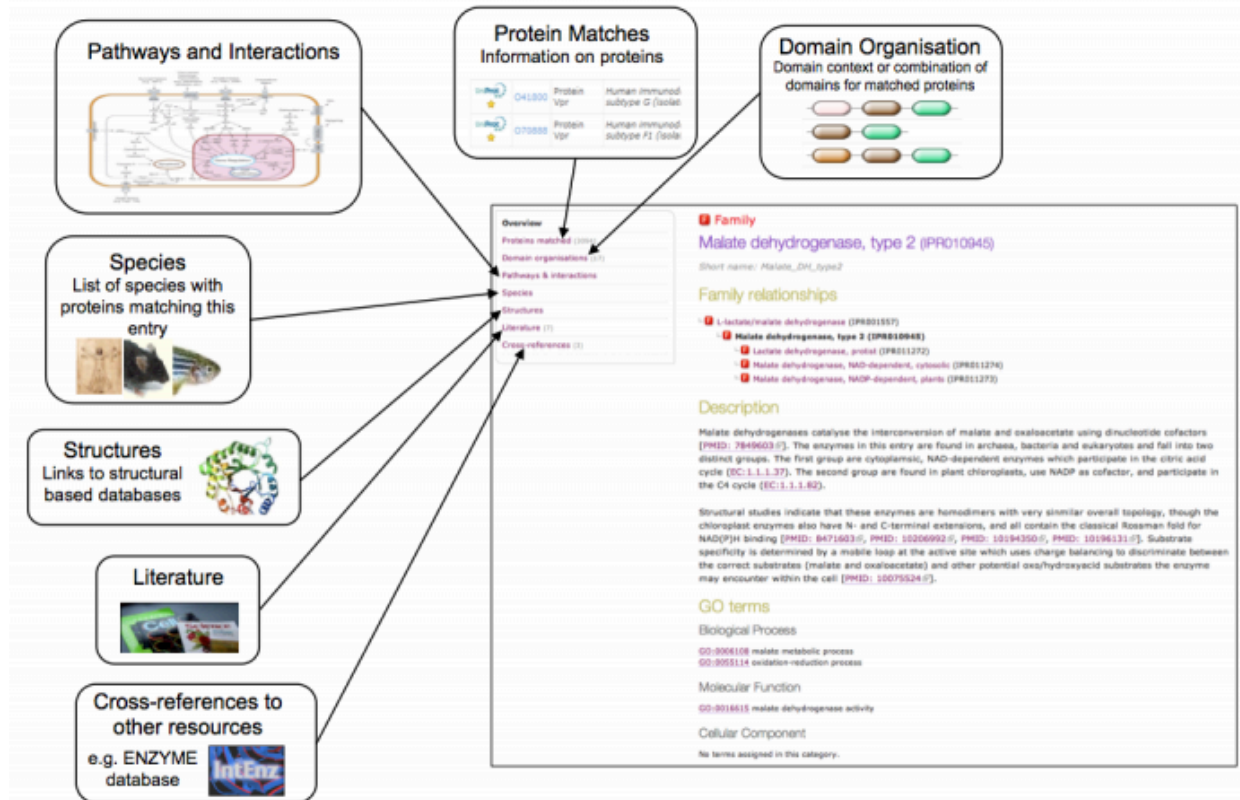
Short name: DEATH-like\_dom

## Domain relationships

- ↳ **D** Death-like domain (IPR011029)
  - ↳ **D** Caspase Recruitment (IPR001315)
  - ↳ **D** DAPIN domain (IPR004020)
  - ↳ **D** Death domain (IPR000488)
  - ↳ **D** Death effector domain (IPR001875)

# Analyse fonctionnelle et structurale de séquences protéiques

Menu en haut à gauche d'une entrée Interpro



# Analyse fonctionnelle et structurale de séquences protéiques

Protein

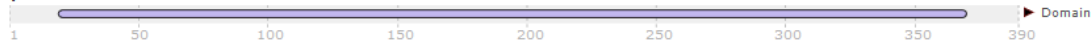
CCGI|60491760|EMBL|CAH06518.1|

Length 390 amino acids

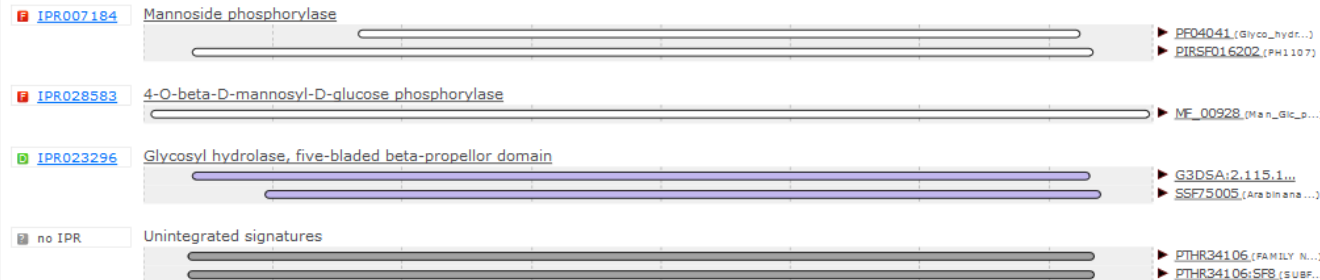
## Protein family membership

- Mannoside phosphorylase (IPR007184)
  - 4-O-beta-D-mannosyl-D-glucose phosphorylase (IPR028583)

## Domains and repeats



## Detailed signature matches



## GO Term prediction

Biological process  
None predicted.

Molecular function  
None predicted.

Cellular component  
None predicted.

Famille prédite par Interpro

Domaines et répétitions prédits comme contenus dans la séquence

Information détaillée concernant les motifs trouvés sur la séquence

Les GO (Gene Ontology) termes prédits à partir des informations précédentes

# Travaux pratiques

<http://www.ebi.ac.uk/interpro/>

- **Exercice 6 : utiliser Interpro**

2) *Explorer les résultats. Passer la souris sur les différents motifs trouvés. Ouvrir les fiches Interpro et des autres bases de données liées.*



# Conclusion

Les motifs peuvent :

- représenter des zones conservées dans des familles de protéines et servir ainsi de “portrait robot” de la famille.
- décrire des domaines ou sites fonctionnels ou structuraux.

Nous avons manipulé plusieurs outils permettant de détecter *de novo* des motifs ou de rechercher des motifs connus dans des séquences.

Identifier des domaines dans une séquence est une étape essentielle de l’annotation fonctionnelle surtout si notre séquence est loin des séquences présentes dans la banque de données curées manuellement.

Pour cela la meilleure stratégie est d’interroger et d’intégrer plusieurs bases de données complémentaires.

# Quelques liens utiles....

<http://skylign.org/>

<http://eddylab.org/software/hmmer3/3.1b2/Userguide.pdf>

# Conclusion Générale

Aujourd'hui nous avons appris à faire des alignements, des arbres et à rechercher des motifs conservés sans ligne de commande.

Mais vous aurez sans doute besoin de l'aide d'un bioinfo pour modifier les header en input du fichier fasta.

Il faudra enlever la redondance et peut-être sous-échantillonner pour certaines interfaces.

Tous les outils présentés existent en ligne de commande et permettent ainsi une utilisation à grande échelle.

# MERCI !

Feedback :

<https://enquetes.inra.fr/index.php?sid=84236&newtest=Y>