

Enquête bioinformatique sur les origines de SARS-CoV-2

CM1 – Analyse des génomes de coronavirus

*Cours donné par **Jacques van Helden**, **Emese Meglécz** et **Gabriel Neve**
Sur base d'une enquête menée par Erwan Salard, José Haloy, Didier Casane,
Etienne Decroly et Jacques van Helden*

Nous mobiliserons une série d'outils bioinformatiques accessibles en ligne pour analyser les séquences de coronavirus et pour tenter de trouver des éléments informatifs concernant l'origine de SARS-CoV-2.

Supports de ce cours	Diapos, tuto, données	https://ivanheld.github.io/shnc-origines-sars-cov-2/
Uniprot	Base de donnée de séquences protéiques	https://www.uniprot.org/
NCBI Entrez	Bases de données biologiques	https://www.ncbi.nlm.nih.gov/
EMBOSS needle	Alignement de paires de séquences	https://www.ebi.ac.uk/Tools/psa/emboss_needle/
NCBI BLAST	Recherche de séquences par similarité	https://blast.ncbi.nlm.nih.gov/Blast.cgi
PIPprofiler	Profils de pourcentages de positions identiques	https://pipprofiler.france-bioinformatique.fr/
Clustal	Alignement de séquences multiples	https://www.ebi.ac.uk/Tools/msa/clustalo/
phylogeny.fr	Phylogénie moléculaire	https://www.phylogeny.fr/
AMU	page AMETICE de N&C3	https://ametice.univ-amu.fr/course/view.php?id=62928

Enquête bioinformatique sur les origines de SARS-CoV-2

CM2 – Inférer la phylogénie de SARS-CoV-2 à partir des séquences génomiques et protéiques

*Cours donné par **Jacques van Helden**, **Emese Megléc** et **Gabriel Neve**
Sur base d'une enquête menée par Erwan Salard, José Haloy, Didier Casane,
Etienne Decroly et Jacques van Helden*

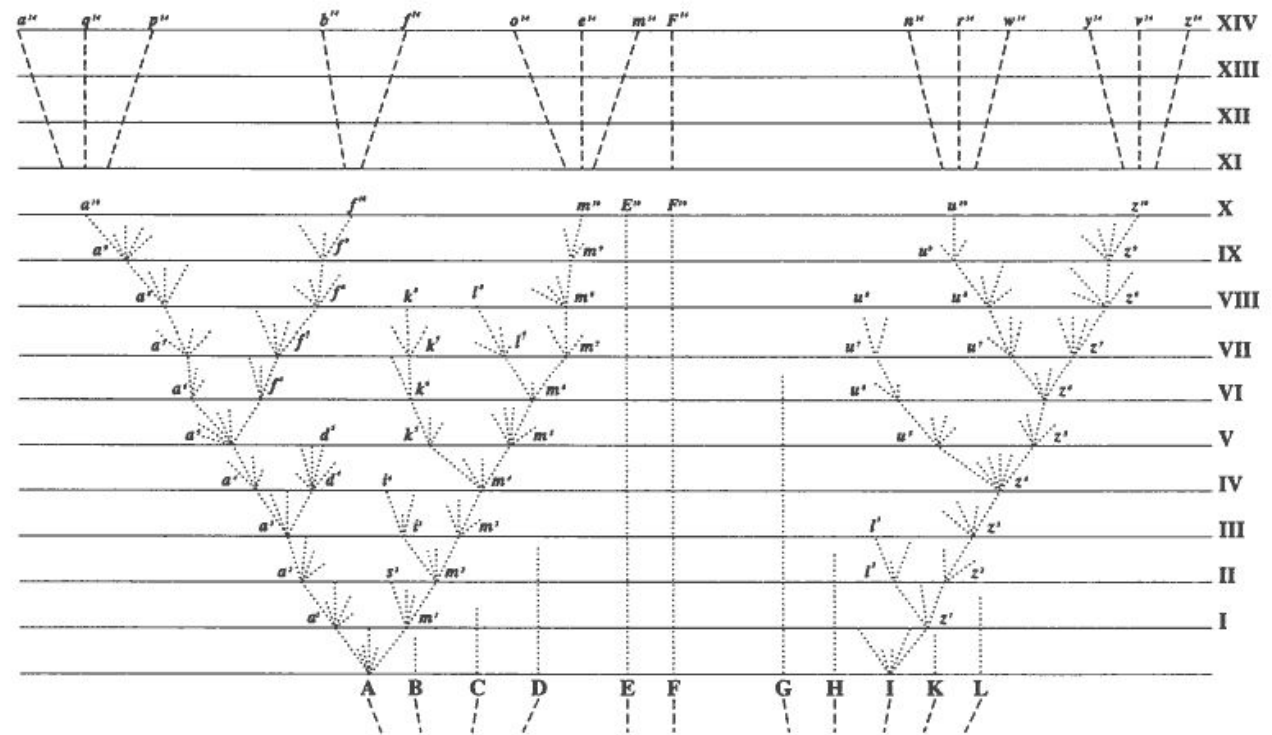
Inférence phylogénétique

*Cours donné par **Jacques van Helden**, **Emese Meglécz** et **Gabriel Neve**
Sur base d'une enquête menée par Erwan Salard, José Haloy, Didier Casane,
Etienne Decroly et Jacques van Helden*

Arbres de la vie

La divergence des caractères

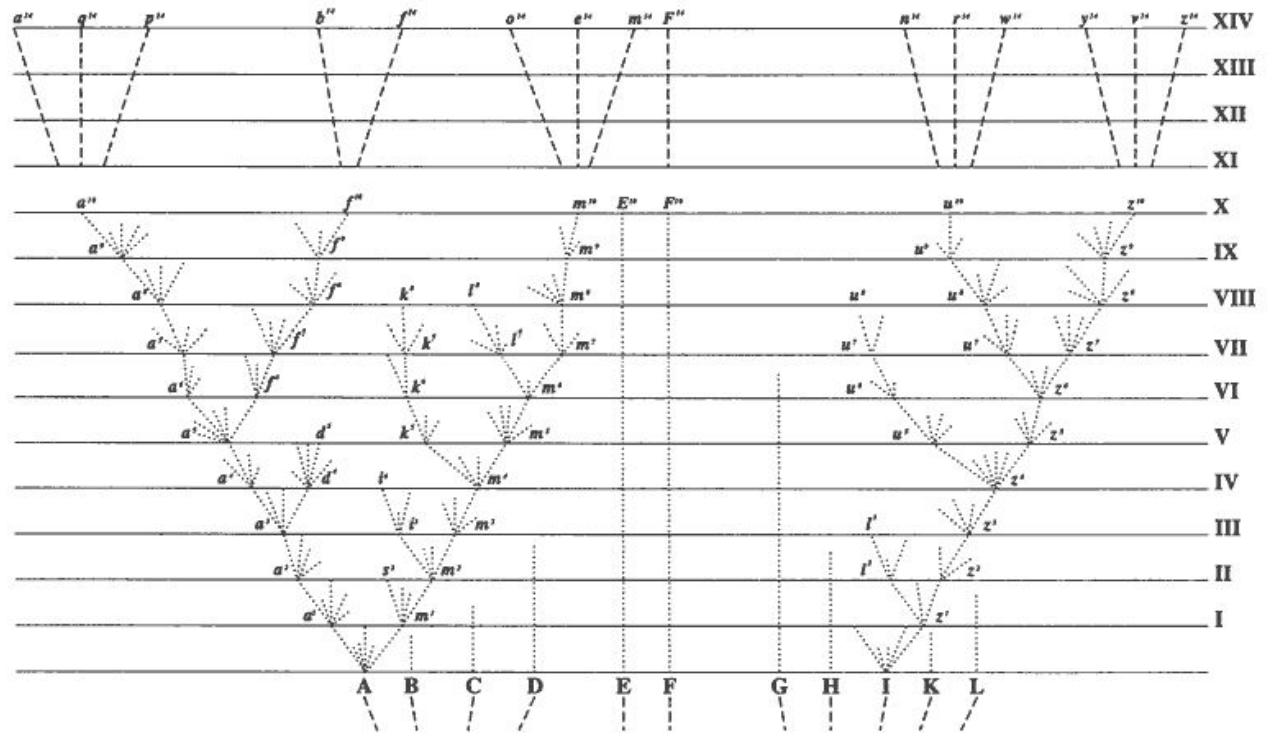
La seule figure de l'Origine des Espèces (C. Darwin, 1859) est une représentation conceptuelle de l'arbre de la vie.



La divergence des caractères

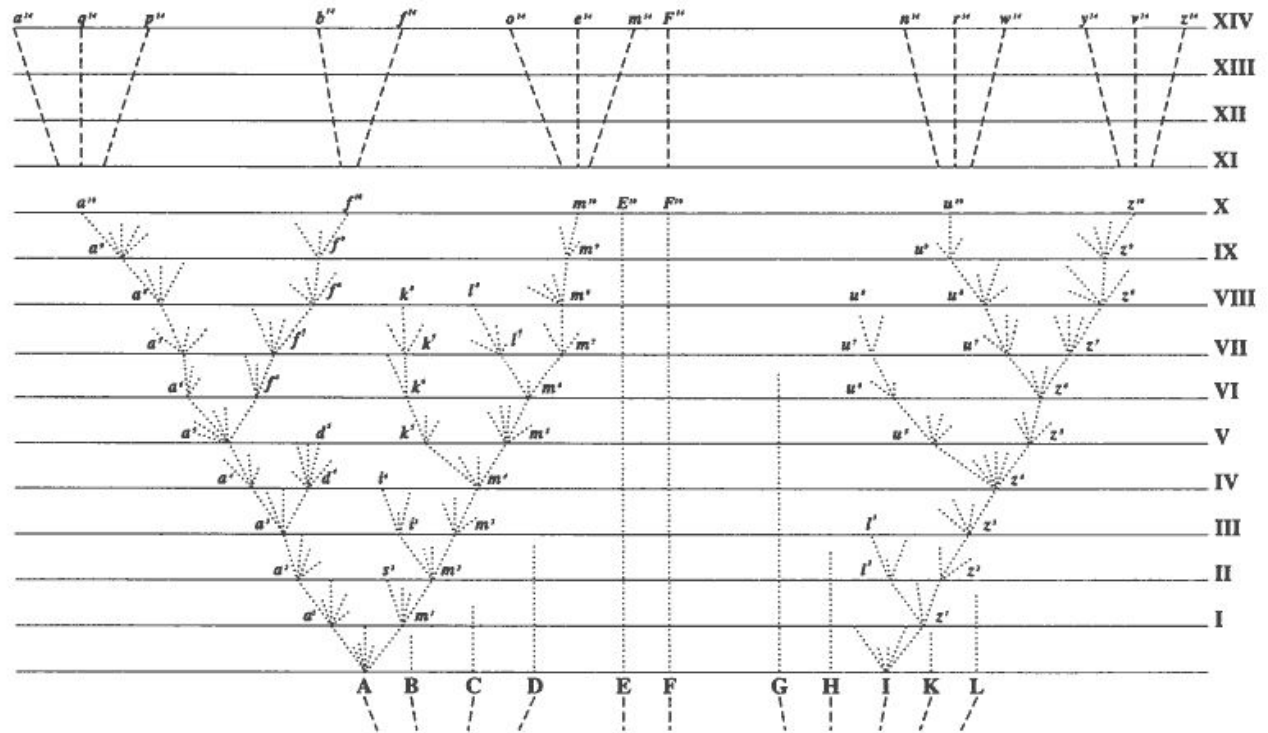
Il s'agit d'un arbre synchrone : chaque niveau horizontal représente un moment donné.

- La racine correspond aux époques les plus anciennes.
- Le niveau le plus élevé correspond au présent.
- A chaque époque on trouve des organismes de différents niveaux de complexité. La hauteur ne représente donc pas une complexité ou un "niveau d'évolution"



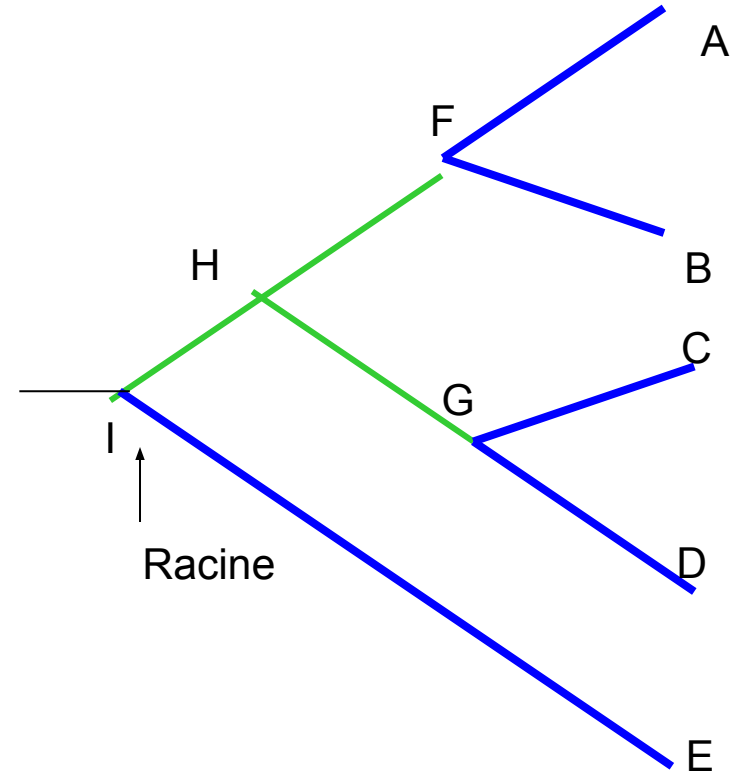
La divergence des caractères

- La plupart des branches sont abortives
- **Evolution graduelle** par accumulation de variations (mutations) le long des branches.
- Juste après un branchement, on a de très petites différences entre les variétés.
- Les observations dont on dispose sont généralement fragmentaires.
- Elles ne sont pas forcément placées sur une trajectoire linéaire depuis un ancêtre donné jusqu'aux espèces actuelles.



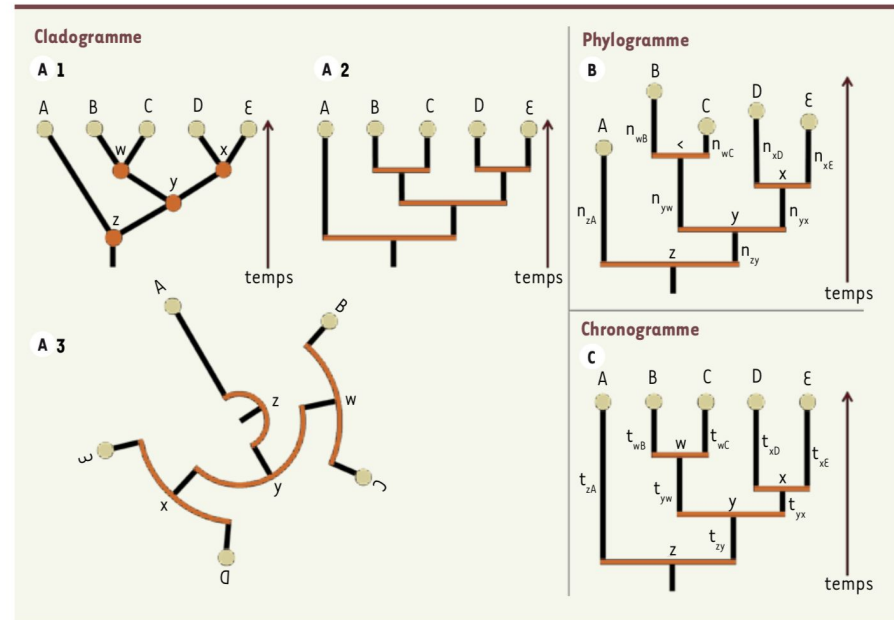
Unités taxonomiques opérationnelles (OTU) et hypothétiques (HTU)

- Les relations évolutives entre les objets étudiés (espèces, organes, séquences) sont représentées par des arbres phylogénétiques
- Les arbres sont des graphes composés de *noeuds* et de *branches*
 - Noeuds = unités taxonomiques
 - Feuilles ou **OTU = Unités Taxonomiques Opérationnelles** (A, B, C, D, E), espèces existantes.
 - **Noeuds internes ou HTU = Unités taxonomiques Hypothétiques** (F, G, H, I), correspondent aux espèces ancestrales.
 - Branches = relations de parenté(ancêtre/descendants) entre unités taxonomiques
 - Branches internes
 - Branches externes
- On appelle **topologie** l'ensemble des branchements de l'arbre.



Représentations arborescentes des histoires évolutives

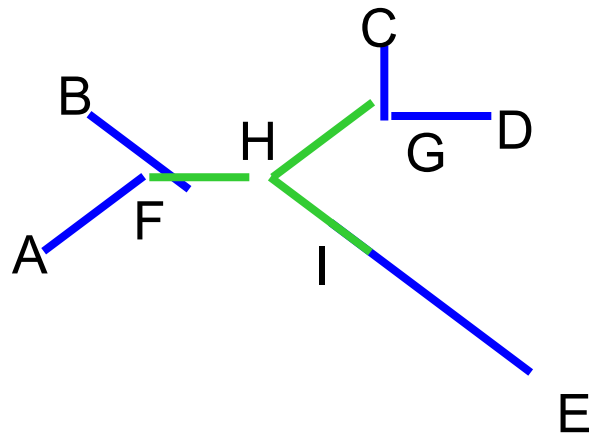
- On représente les histoires évolutives sous forme d'arbre
- Différents types de représentation peuvent être utilisés selon les cas.
 - Bifurcations triangulaires ou rectangulaires
 - Disposition radiale
- Selon les cas, les longueurs des branches représentent
 - Le nombre de divergences (cladogramme)
 - le nombre de différences entre deux espèces (phylogramme)
 - Le temps de divergence (chronogramme)



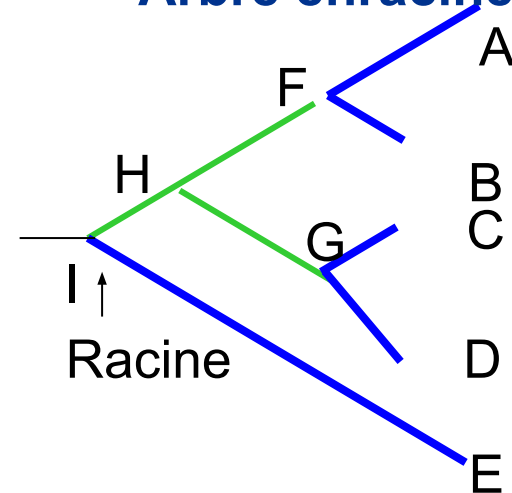
Arbres enracinés ou non enracinés

- Les arbres non-enracinés ne sont pas réellement des arbres phylogénétiques car ils n'ont pas de direction temporelle
-> indiquent les distances, mais pas les relations de parenté entre les noeuds.
- La **racine** définit une orientation de l'arbre, et donc un chemin évolutif unique vers chaque feuille.
- Elle symbolise le *dernier ancêtre commun* (i.e. le plus récent) de toutes les OTU.

Arbre non-enraciné



Arbre enraciné



Combien d'arbres ?

- Le nombre d'arbres possibles augmente de façon vertigineuse en fonction du nombre d'éléments terminaux (qu'ils représentent des molécules ou des espèces).
- Un seul de ces arbres correspond à l'histoire évolutive réelle.
- Puisqu'on ne dispose pas a priori de cet arbre, on doit l'*inférer* à partir des éléments actuels (les unités taxonomiques opérationnelles, UTO).

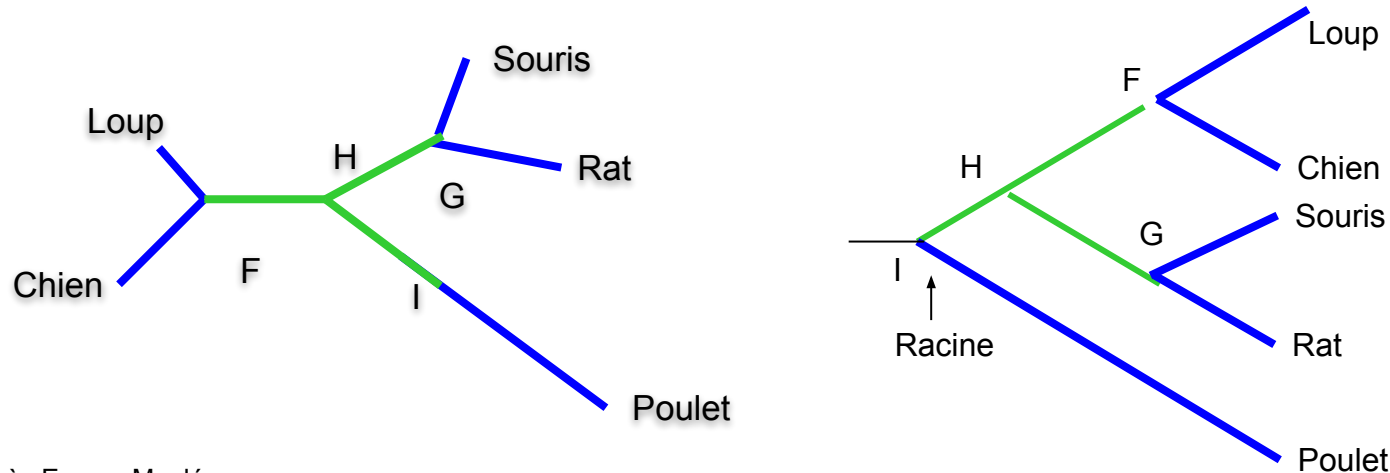
n	Nb arbres enracinés	Nb arbres non-enracinés
2	1	1
3	3	1
4	15	3
5	105	15
6	945	105
7	10,395	945
8	135,135	10,395
9	2,027,025	135,135
10	3.45E+07	2,027,025
11	6.55E+08	3.45E+07
12	1.37E+10	6.55E+08
13	3.16E+11	1.37E+10
14	7.91E+12	3.16E+11
15	2.13E+14	7.91E+12
16	6.19E+15	2.13E+14
17	1.92E+17	6.19E+15
18	6.33E+18	1.92E+17
19	2.22E+20	6.33E+18
20	8.20E+21	2.22E+20

$$N_R = \frac{(2n! 3)!}{2^{n!2} (n! 2)!}$$

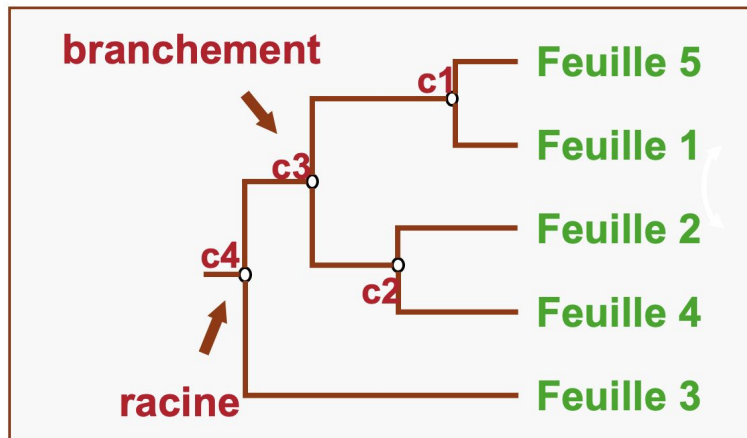
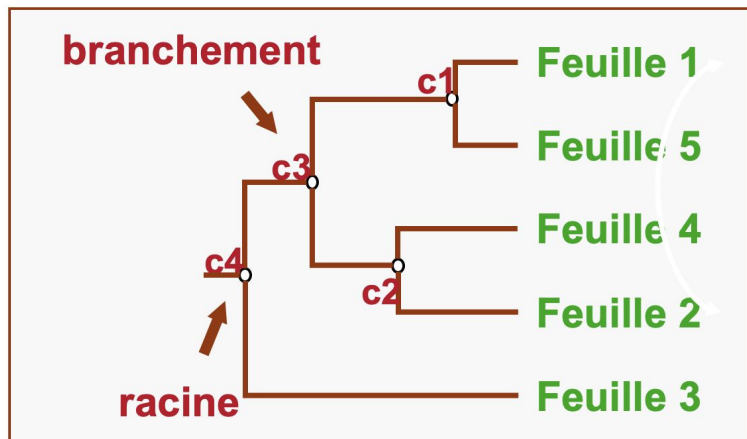
$$N_U = \frac{(2n! 5)!}{2^{n!3} (n! 3)!}$$

Comment enraciner un arbre phylogénétique ?

- Connaissance *a priori* de la feuille la plus externe parmi les OTU étudiées (« *outgroup* »)
 - Exemple: chien, loup, souris, rat et poulet
 - Sur base des connaissances biologiques, on décide que le **Groupe extérieur** est le poulet
- Sans connaissance *a priori* du OTU les plus externes parmi les OTU étudiées
 - Enracinement au poids moyen: on enracine l'arbre sur la branche qui minimise la moyenne des distances aux feuilles.
 - Ceci implique une hypothèse d'**horloge moléculaire**: on considère que le taux de mutation est constant au cours de l'évolution, et égal entre les branches.
 - Cette hypothèse n'est généralement pas très réaliste, il s'agit d'une approximation.



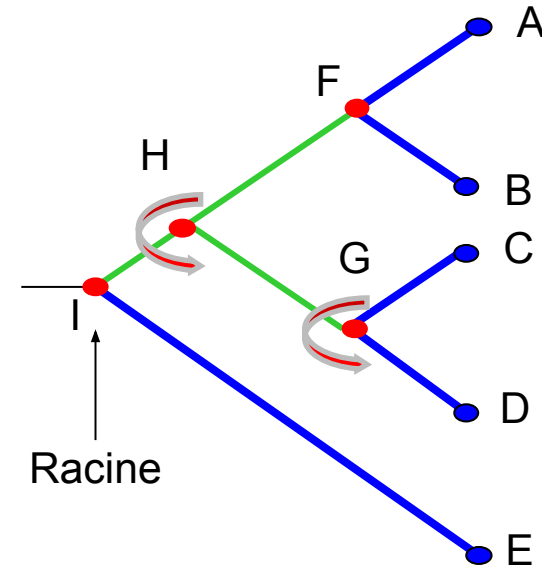
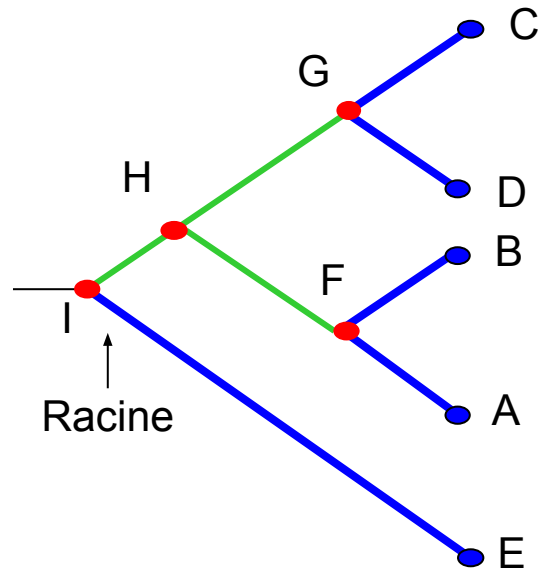
Isomorphisme sur un arbre



- Dans un arbre, les deux enfants de chaque branche peuvent être interchangeés.
- Le résultat est un arbre **isomorphe**, considéré équivalent à l'arbre initial.
- Les deux arbres de gauche sont équivalents.
- Cependant
 - Arbre du dessus: les feuilles 1 et 2 sont très éloignées.
 - Arbre du dessous: les feuilles 1 et 2 sont voisines.
- Les distances verticales entre deux nœuds ne reflètent pas leur distance réelle !
- La distance entre deux nœuds est la somme des longueurs des branches qui les séparent.

L'isomorphisme des arbres phylogénétiques

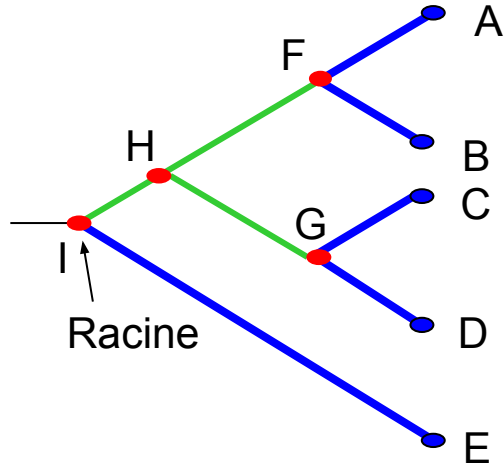
- Il faut éviter le piège d'évaluer les distance entre feuilles sur base de leur proximité verticale.
 - Les structures ci-dessous sont absolument identiques.
 - Pourtant les feuilles B et D semblent voisines sur le graphe de gauche, et éloignées sur celui de droite.
- Pour évaluer la distance entre deux nœuds d'un arbre , il faut prendre en compte la longueur totale du chemin le plus court pour les rejoindre (somme des longueurs de branches).



Que représente l'échelle d'un arbre phylogénétique ?

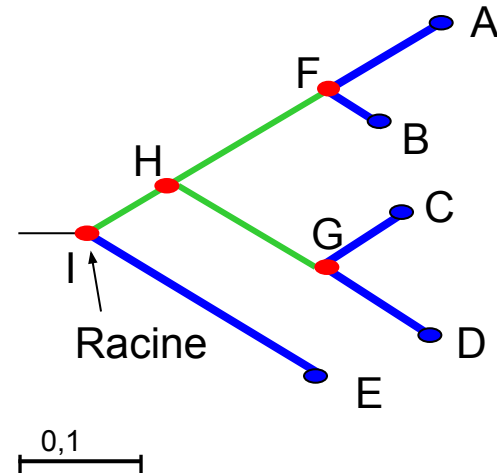
■ Cladogramme

- Représentation sans échelle
- L'arbre indique uniquement l'ordre des branchements.
- Les longueurs de branches ne sont pas proportionnelles au nombre de changements évolutifs.



■ Phylogramme

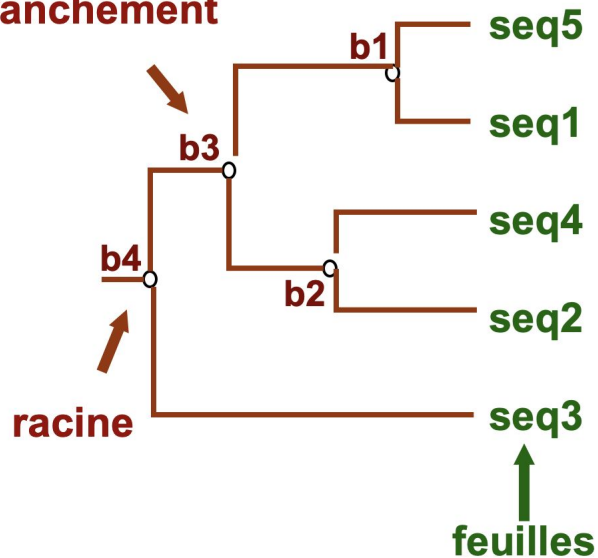
- Représentation avec échelle
- L'arbre indique les distances évolutives entre nœuds.
- Les longueurs de branches sont proportionnelles au nombre d'événements évolutifs (substitutions ou substitution/sites).



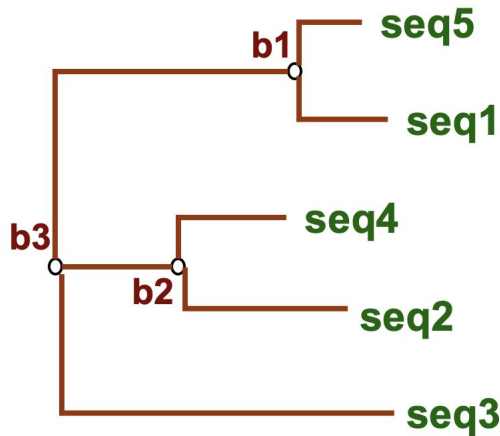
Calcul de la distance sur un arbre

Arbre enraciné

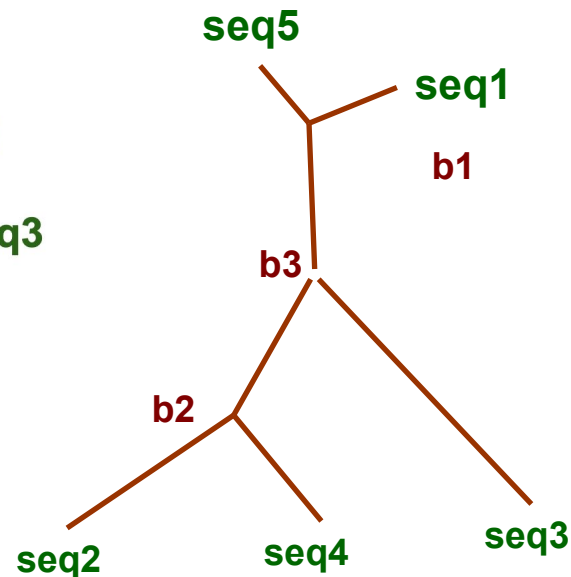
branchement



Arbre non-enraciné



Arbre non-enraciné



- La distance entre deux nœuds est la somme des longueurs des branches qui les séparent.

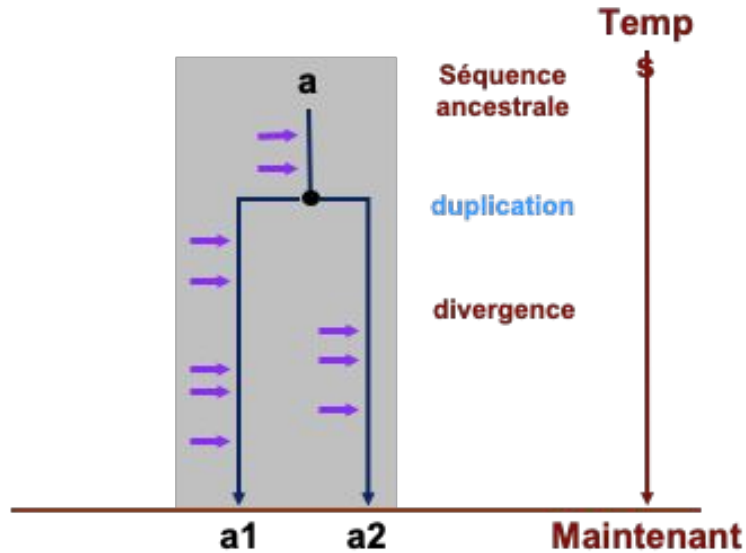
Concepts: homologie, paralogie, orthologie

- Pour l'analyse de la phylogénie moléculaire, nous porterons un intérêt tout particulier à deux événements évolutifs particuliers: duplication et spéciation.
- **Duplication**
 - Une duplication est une mutation qui génère un dédoublement d'une partie de l'ADN génomique. La duplication peut recouvrir l'ensemble du génome (formation de polyploïdes), un chromosome entier, ou un fragment de chromosome de taille plus ou moins grande.
 - Les duplications peuvent éventuellement entraîner l'apparition de copies multiples d'un ou plusieurs gènes, provoquant ainsi une certaine redondance de l'information génétique.
 - Dans certains cas, l'une des copies dupliquées du gène acquiert, par accumulation de mutations, de nouvelles caractéristiques qui lui permettent d'assumer une nouvelle fonction. Ce mécanisme, appelé duplication divergence, est à l'origine de la diversification des fonctions biologiques.
- **Spéciation**
 - Processus évolutif qui résulte en la formation d'espèces distinctes à partir d'une espèce unique.
- Les événements de duplication et spéciation suscitent l'apparition de copies multiples à partir d'une seule séquence, soit au sein d'une même espèce (duplication), soit au sein des espèces distinctes dérivées de la spéciation. Ces séquences, dont la similarité résulte d'une séquence ancestrale commune, sont dites **homologues**

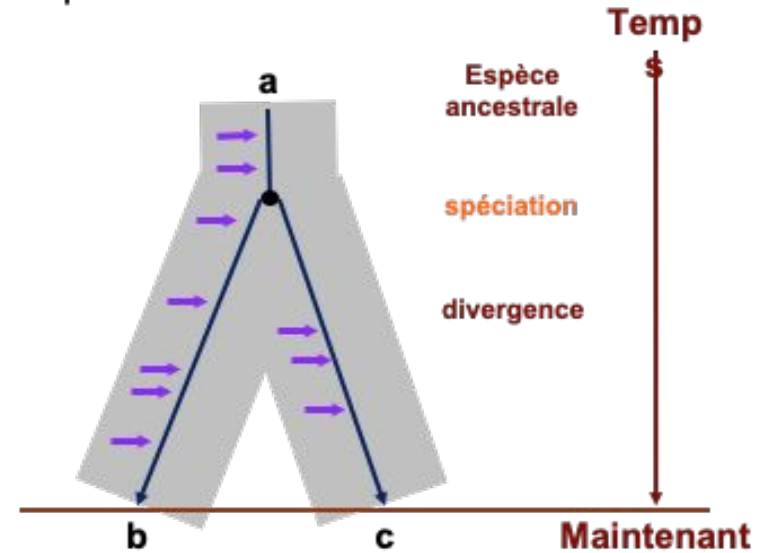
Scénarios évolutifs

- Nous disposons de deux séquences, et nous supposons qu'elles divergent d'un ancêtre commun.
- La divergence peut résulter
 - d'une **duplication** (création de deux copies du gène dans le même génome)
 - ou d'une **spéciation** (formation d'espèces séparées à partir d'une espèce unique).
- Les **flèches violettes** indiquent les mutations (substitutions, délétions, insertions) qui s'accumulent au sein d'une séquence particulière au cours de son histoire évolutive. Ces mutations sont à l'origine de la diversification des séquences, des structures et des fonctions.

Duplication



Spéciation



- La similarité entre deux traits (organes, séquences) peut s'interpréter par deux hypothèses alternatives: homologie et analogie.
- **Homologie**
 - La similarité s'explique par le fait que les deux séquences divergent d'un ancêtre commun.
 - Les différences entre les deux caractères homologues résultent de l'accumulation de mutations à partir de l'ancêtre commun. Il s'agit donc d'une évolution par **divergence évolutive**.
- **Analogie**
 - Ressemblance entre deux traits (organes, séquence) qui ne résulte pas d'une origine ancestrale commune (par opposition à l'homologie).
 - Les traits similaires sont apparus de façon **indépendante**. Leur ressemblance peut éventuellement manifester l'effet d'une pression évolutive qui a sélectionné les mêmes propriétés.
 - Dans ce cas, on parle de **convergence évolutive**.

- Inférence
 - Avant d'affirmer que deux séquences sont homologues, nous devrions pouvoir retracer leur histoire jusqu'à leur ancêtre commun.
 - Nous ne pouvons malheureusement pas disposer des séquences de toutes les espèces disparues. Il est donc impossible de démontrer formellement l'homologie.
 - Cependant, nous pouvons appuyer l'hypothèse d'homologie sur une analyse de la vraisemblance d'un scénario évolutif (taux de mutations, niveaux de similarités).
 - L'inférence d'homologie est toujours attachée à un certain **risque de faux positifs**. Les modèles évolutifs nous permettent d'estimer ce risque.

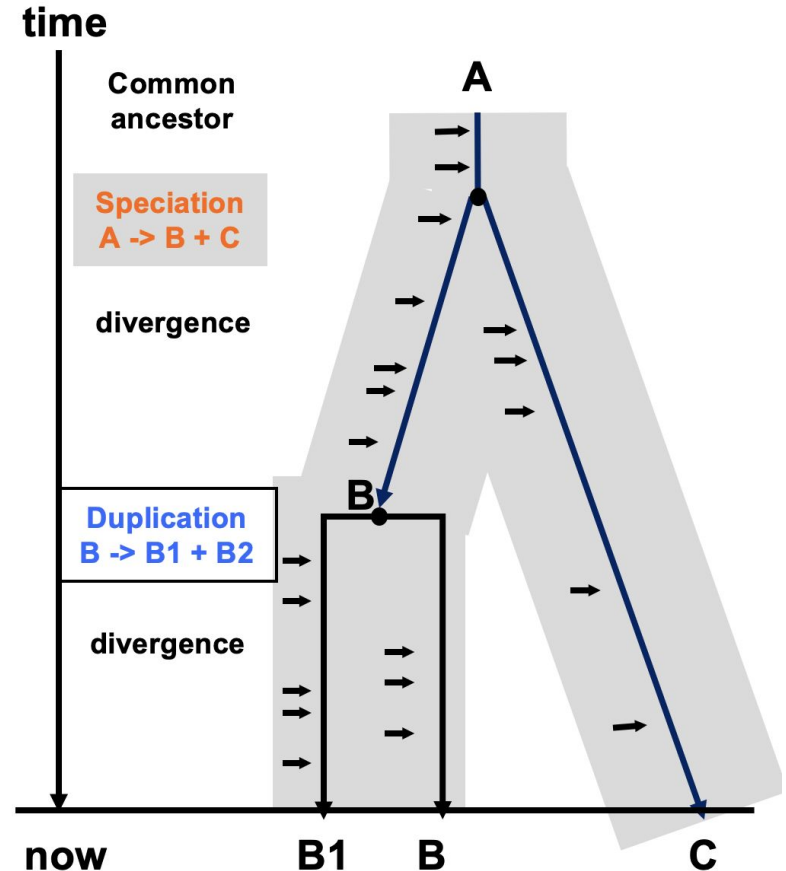


L'homologie est une relation logique (soit vraie, soit fausse).

- Deux séquences sont homologues (possèdent des caractères communs parce qu'elles dérivent d'un ancêtre commun) ou elles ne le sont pas.
- Il est donc complètement inapproprié de parler de « niveau d'homologie » ou « pourcentage d'homologie ».
- La formulation correcte
 - On observe un certain niveau de similarité entre deux séquences (pourcentages de résidus identiques, pourcentages de résidus « similaires »).
 - Sur cette base, on évalue deux scénarios évolutifs: cette similarité peut provenir d'une évolution convergente (analogie) ou divergente à partir d'un ancêtre commun (homologie).
 - Si la deuxième hypothèse est la plus vraisemblable, on *infère* que les séquences sont homologues.

Orthologie versus paralogie

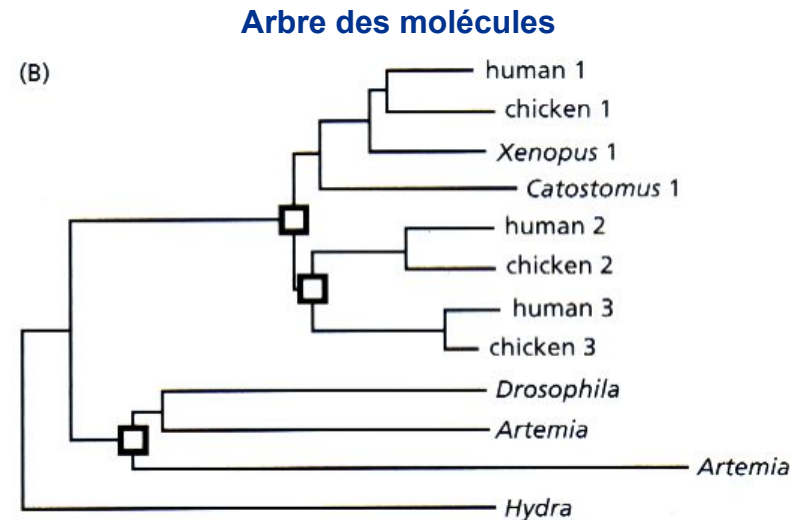
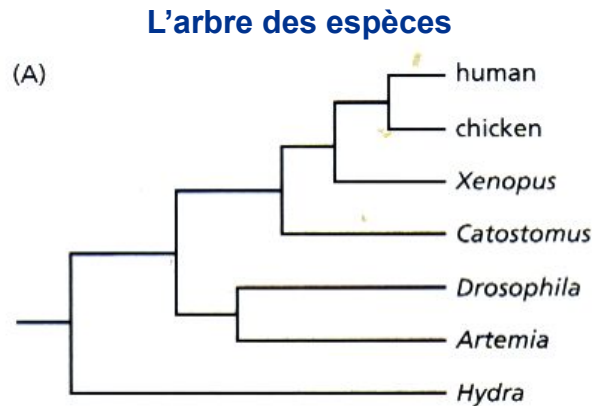
- Zvelebil & Baum (2000) fournissent une définition claire et opérationnelle des concepts d'orthologie et paralogie.
 - **Orthologues**: séquences dont le dernier ancêtre commun précède immédiatement un événement de spéciation.
 - **Paralogues** séquences dont le dernier ancêtre commun précède immédiatement un événement de duplication
- Exemples:
 - B et C sont **orthologues**, car leur dernier ancêtre commun (A) précède un événement de **spéciation** ($A \rightarrow B + C$).
 - B1 et B2 sont **paralogues** car le premier événement évolutif qui succède à leur dernier ancêtre commun (B) est une **duplication** ($B \rightarrow B1 + B2$).



Approches d'inférence phylogénétique

Inférence phylogénique à partir de séquences moléculaires

- En partant d'une famille de séquences macromoléculaires (ADN, ARN, protéines), on peut construire des arbres phylogéniques.
- En comparant l'arbre des molécules et l'arbre des espèces, on peut inférer l'histoire évolutive de cette famille de séquences.
- Nous reviendrons plus tard sur cet exemple, en expliquant les méthodes bioinformatiques permettant d'inférer des arbres moléculaires à partir de séquences, et les façons d'interpréter ces arbres en tenant compte de la filiation des espèces.



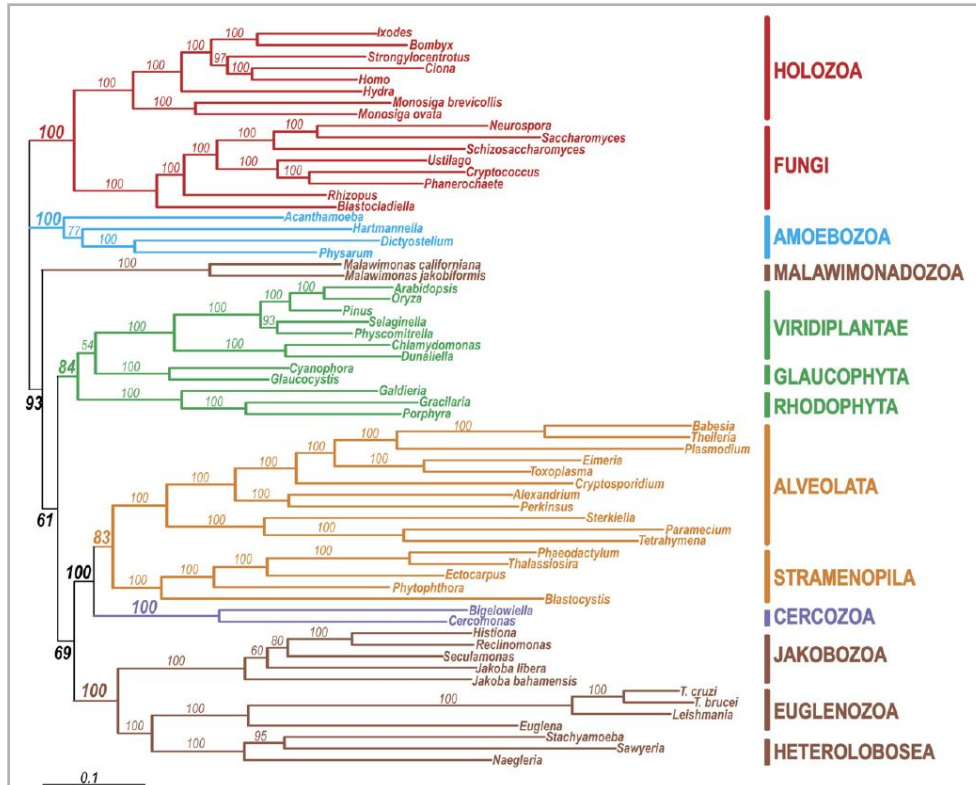


Figure 1. Maximum-Likelihood Tree of Eukaryotes

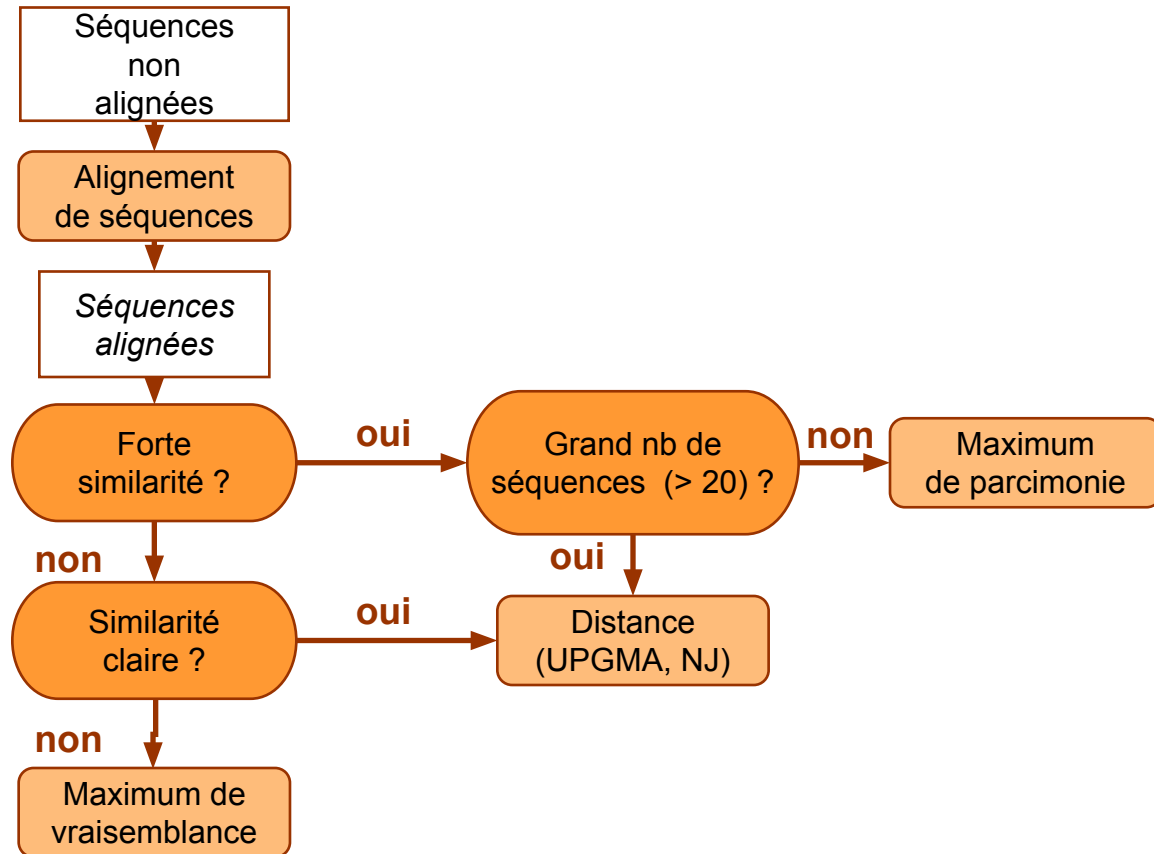
The tree includes 64 species and is based on 143 concatenated nucleus-encoded proteins (31,604 amino acid positions). Numbers indicate support values of RaxML analysis (100 replicates) with the WAG + F + Γ model. Posterior probabilities obtained in the Bayesian Inference with MrBayes are 1.0 for all branches. The scale bar denotes the estimated number of amino acid substitutions per site. The tree was rooted according to a gene fusion [13, 16].

- En phylogénie moléculaire, une approche classique consiste à se concentrer sur un gène considéré comme représentatif, et à construire un arbre sur base de la divergence de séquence de ce gène.
- Ces approches peuvent maintenant être généralisées en comparant les séquences de plusieurs centaines de gènes (ci-contre, arbre basé sur 143 familles de protéines).
- Elles permettent d'inférer des phylogénies entre organismes très éloignés (règnes différents), et d'établir ainsi des scénarios concernant les premières étapes de la diversification des êtres vivants.

- Source: Rodríguez-Ezpeleta et al. Curr Biol (2007) vol. 17 (16) pp. 1420-5
- Toward resolving the eukaryotic tree: the phylogenetic positions of jakobids and cercozoans

Inférence phylogénétique par comparaison de séquences

- Il existe plusieurs méthodes pour inférer un arbre évolutif à partir de séquences.
 - Maximum de parcimonie
 - Distance
 - Maximum de vraisemblance
- On part toujours d'un jeu de séquences alignées (alignement multiple).
- Le choix de la méthode dépend du nombre de séquences, et de leur degré de similarité.



Exemple : la famille des opsines

- Pour inférer un arbre phylogénétique à partir d'une famille de séquences, on part toujours d'un alignement multiple.
- La figure ci-dessous montre la première partie d'un alignement multiple entre 50 opsines de mammifère.
- A l'œil nu, on distingue déjà 2 groupes évidents.
 - Dessus: opsines sensibles aux ondes moyennes (vert) ou longues (rouge)

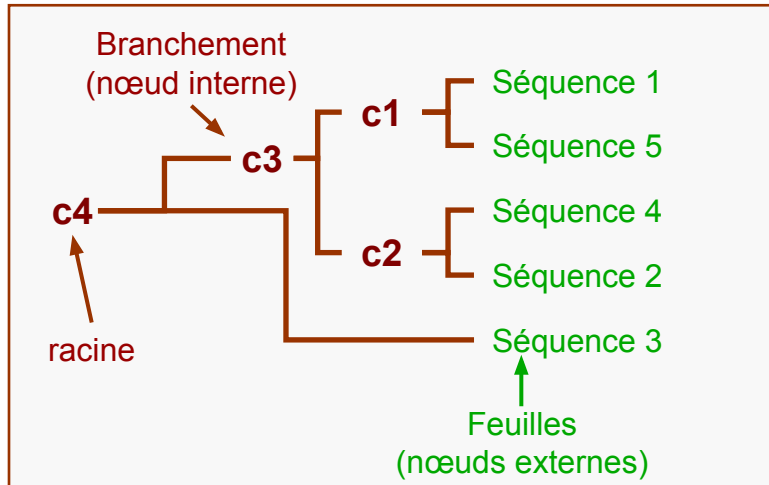


Principe de la construction de l'arbre

Matrice de distance

	séquence 1	séquence 2	séquence 3	séquence 4	séquence 5
séquence 1	0.00	4.00	6.00	3.50	1.00
séquence 2	4.00	0.00	6.00	2.00	4.50
séquence 3	6.00	6.00	0.00	5.50	6.50
séquence 4	3.50	2.00	5.50	0.00	4.00
séquence 5	1.00	4.50	6.50	4.00	0.00

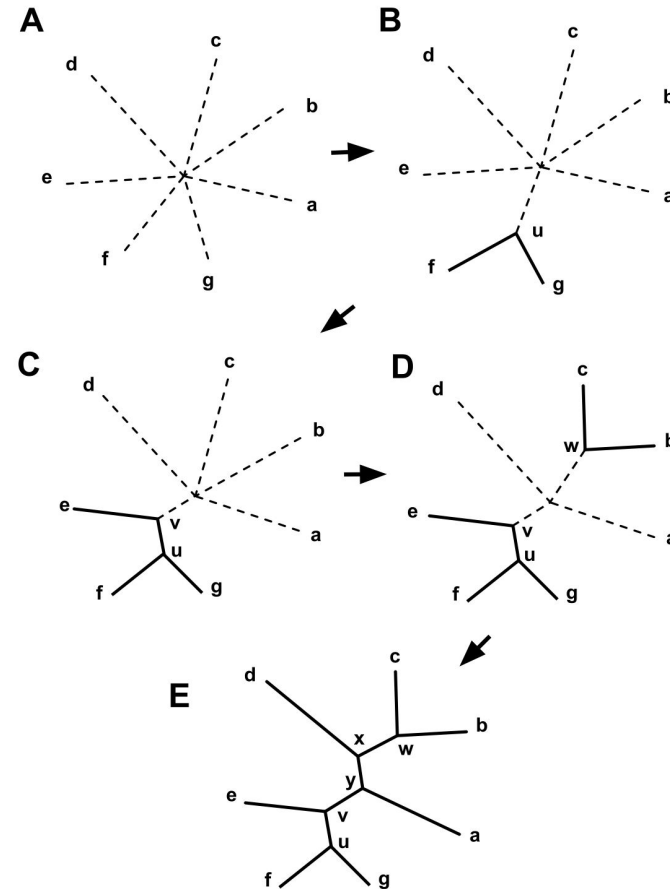
Arbre



- Le clustering hiérarchique est une méthode de clustering agrégative.
 - Prend une matrice de distance en entrée
 - Regroupe progressivement les objets en allant des plus proches aux plus distants.
- Il existe plusieurs possibilités pour établir une règle d'agglomération, qui définit la distance entre deux groupes.
 - Liaison simple (**single linkage**): distance entre groupes A et B est la distance entre les plus proches de leurs éléments respectifs.
 - Liaison moyenne (**average linkage**): distance moyenne entre tous les objets des deux groupes (=UPGMA).
 - Liaison complète (**complete linkage**): distance entre les éléments les plus éloignés des groupes A et B.
- Algorithmes
 - 1. Assigner chaque objet à un cluster séparé.
 - 2. Identifier la paire de clusters les plus proches, et les regrouper en un seul.
 - 3. Répéter la seconde étape jusqu'à ce qu'il ne

Neighbour joining (NJ) - Méthode

- Développé par Saitou et Nei (1987) est une approximation de l'algorithme pour trouver l'arbre le plus court (minimum évolution)
- Principe:
 - A chaque étape, rechercher le couple d'UTO qui minimise la longueur totale de l'arbre

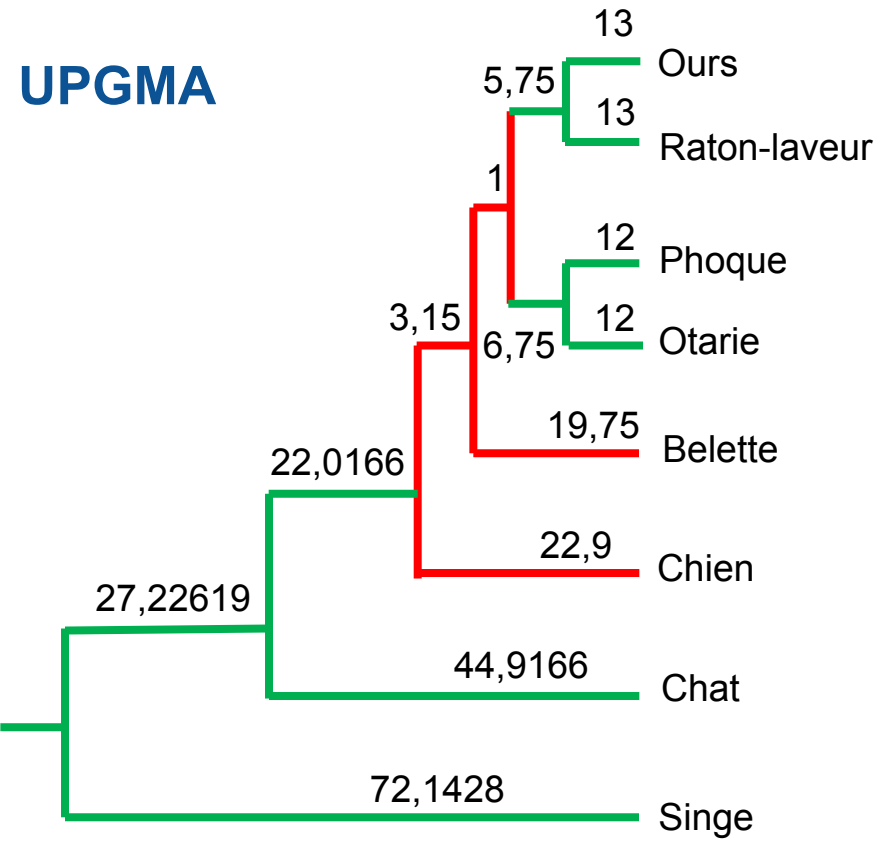


Propriétés de la méthode NJ

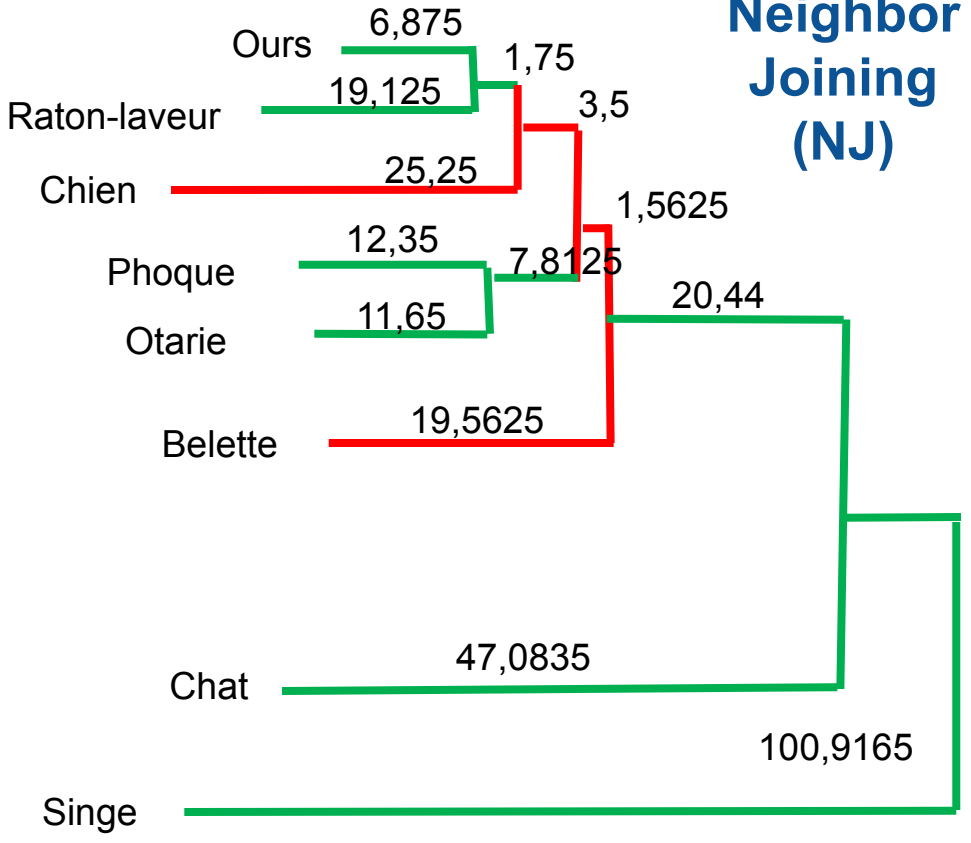
- Méthode rapide et simple qui permet de travailler avec un très grand nombre de taxons.
- Les arbres ne sont pas enracinés.
- Les longueurs des branches sont informatives (phylogramme).
- Bonne approximation de la méthode du minimum d'évolution (l'arbre le plus court).
- Retrouve l'arbre vrai si la matrice de distances est un reflet exact des distances évolutives (malheureusement ce n'est pas souvent le cas).
- Ne dépend pas d'hypothèse de l'horloge moléculaire, donc la méthode est applicable dans les cas où le taux d'évolution varie entre les lignées.

Comparison UPGMA - Neighbour Joining

UPGMA



Neighbor Joining (NJ)

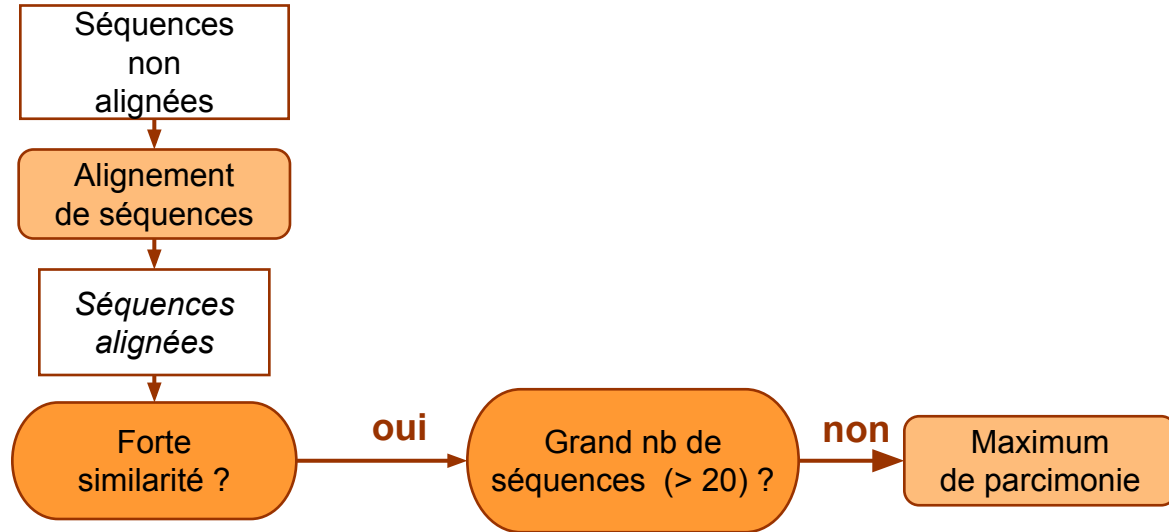


*Inférence phylogénétique par la méthode
du maximum de parcimonie*

Inférence phylogénétique par comparaison de séquences

■ Approches alternatives

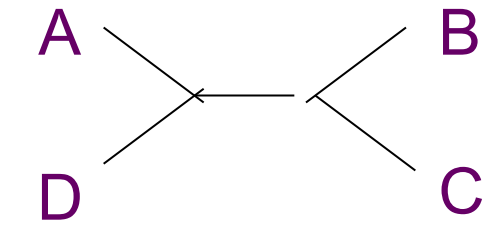
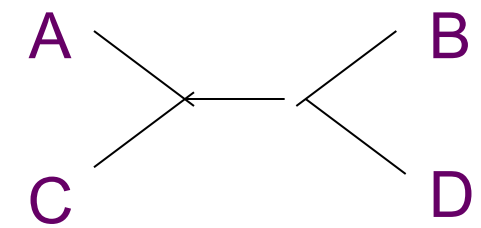
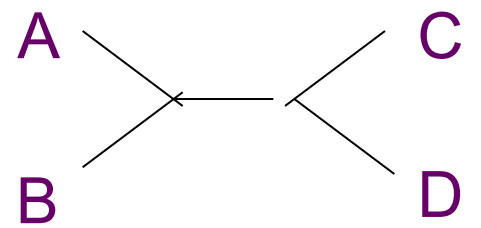
- Maximum de parcimonie
- Distance
- Maximum de vraisemblance



- Principe:
 - ❑ Identifier la topologie T qui implique le plus petit nombre de changements évolutifs suffisant à rendre compte des différences observées entre les OTU étudiées.
 - ❑ Utilise des états de caractères discrets => L'arbre le plus parcimonieux => plus court chemin conduisant aux états de caractères observés
- Algorithme
 - ❑ Construction de tous les arbre possibles.
 - ❑ Pour chaque site (position de l'alignement), on compte le nombre de substitutions nécessaires pour expliquer chaque arbre.
 - ❑ On retient l'arbre qui nécessite le plus petit nombre de substitutions au total (en tenant compte de tous les sites).
- Caractéristique des arbres obtenus
 - ❑ Solutions multiples => plusieurs arbres avec le même nombre minimum de changements peuvent être obtenus.
 - ❑ La longueur des branches ne reflète par la distance évolutive (arbre sans échelle).
 - ❑ Arbres non enracinés.

Déterminer toutes les topologies possibles
4 UTO => 3 arbres non racinés

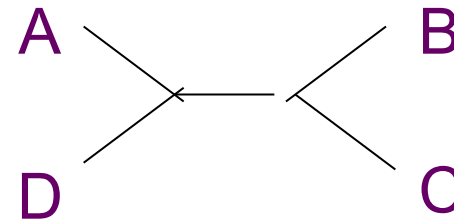
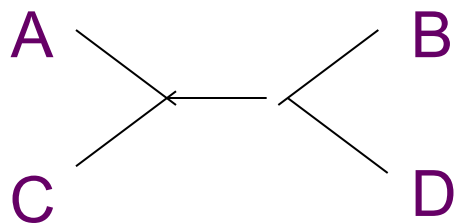
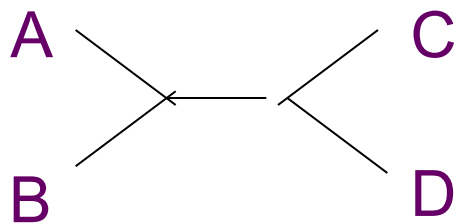
Espèces	Séquences								
	1	2	3	4	5	6	7	8	9
A	A	A	G	A	G	T	T	C	A
B	A	G	C	C	G	T	T	C	T
C	A	G	A	T	A	T	C	C	A
D	A	G	A	G	A	T	C	C	T



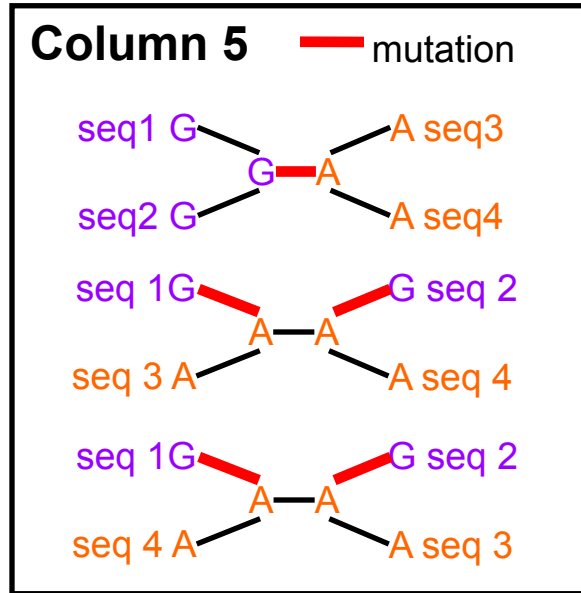
Maximum de parcimonie – classification des sites

- Caractère **invariant**: toutes les OTU possèdent le même état de caractères pour un site donné.
- Caractère **variable**
 - **Non informatif** si les états de caractères à ce site ne favorisent aucune topologie parmi l'ensemble des topologies possibles
 - **Informatif** si les états de caractères à ce site favorise une (ou plusieurs) topologie(s) parmi l'ensemble des topologies possibles

Espèces	Séquences								
	1	2	3	4	5	6	7	8	9
A	A	A	G	A	G	T	T	C	A
B	A	G	C	C	G	T	T	C	T
C	A	G	A	T	A	T	C	C	A
D	A	G	A	G	A	T	C	C	T



position	1	2	3	4	5	6	7	8	9
seq1	A	A	G	A	G	T	G	C	A
seq2	A	G	C	C	G	T	G	C	G
seq3	A	G	A	T	A	T	C	C	A
seq4	A	G	A	G	A	T	C	C	G

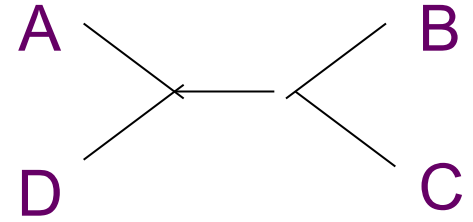
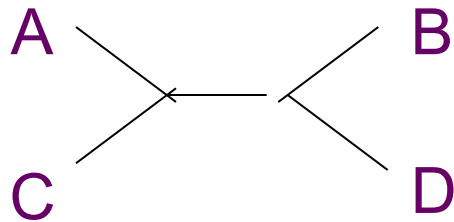
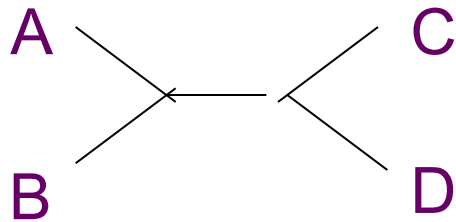


Adapted from Mount (2000)

- Pour chacune des colonnes de l'alignement, tous les arbres possibles sont évalués.
- Pour chaque colonne informative, l'arbre qui présente le plus petit nombre de mutations est retenu.
- On retient ensuite l'arbre qui correspond au nombre le plus élevé de colonnes (consensus)
- Note: cette approche peut éventuellement retourner plusieurs arbres *ex aequos*.

Déterminer toutes les topologies possibles
4 UTO => 3 arbres non racinés

Espèces	Séquences								
	1	2	3	4	5	6	7	8	9
A	A	A	G	A	G	T	T	C	A
B	A	G	C	C	G	T	T	C	T
C	A	G	A	T	A	T	C	C	A
D	A	G	A	G	A	T	C	C	T



Maximum de parcimonie - Méthode

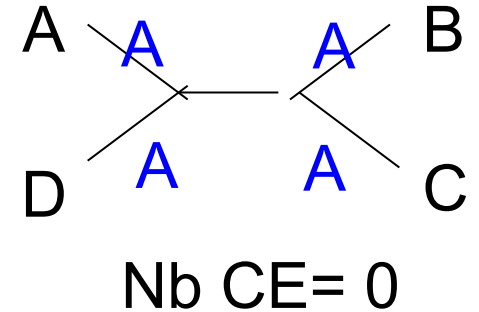
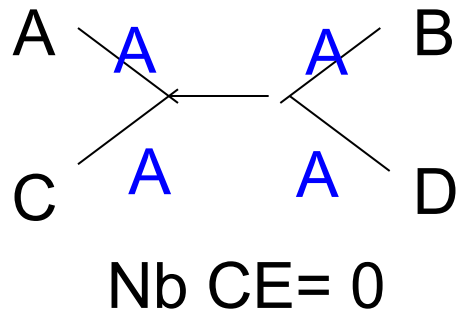
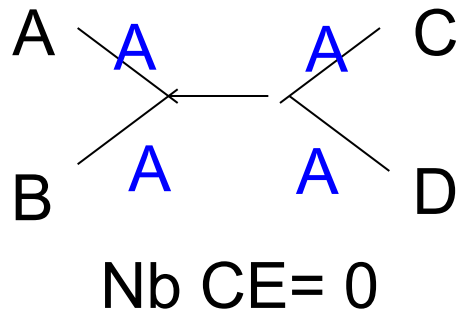
Pour un caractère donné, on compte le nombre de changements évolutifs (CE) pour chaque topologie possible.

Étude du caractère n°1

	1	2	3	4	5	6	7	8	9
A	A	A	G	A	G	T	T	C	A
B	A	G	C	C	G	T	T	C	T
C	A	G	A	T	A	T	C	C	A
D	A	G	A	G	A	T	C	C	T

Caractère constant (même état de caractère à tous les sites).

Caractère **non informatif** : ne favorise aucune topologie par rapport à une autre.

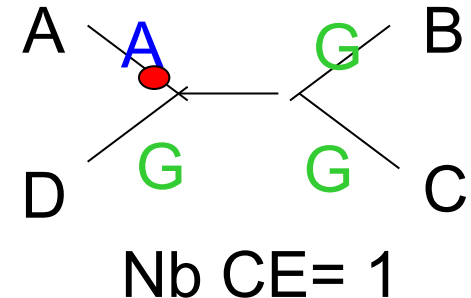
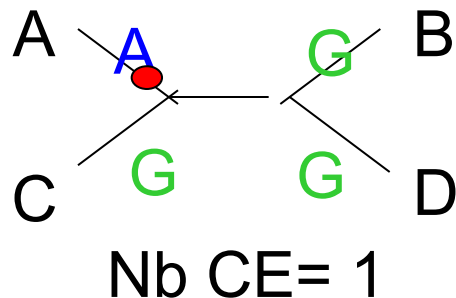
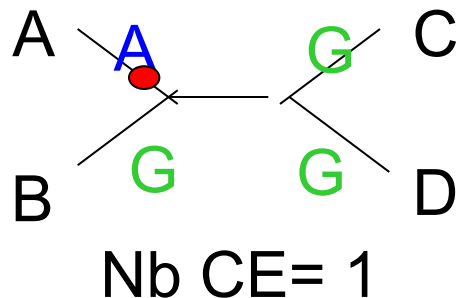


Étude du caractère n°2

	1	2	3	4	5	6	7	8	9
A	A	A	G	A	G	T	T	C	A
B	A	G	C	C	G	T	T	C	T
C	A	G	A	T	A	T	C	C	A
D	A	G	A	G	A	T	C	C	T

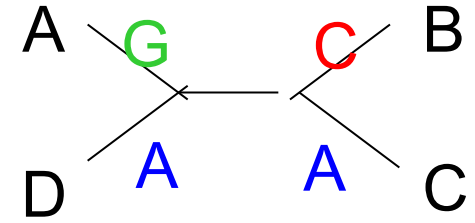
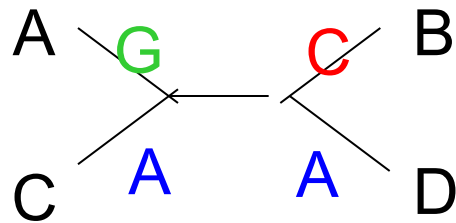
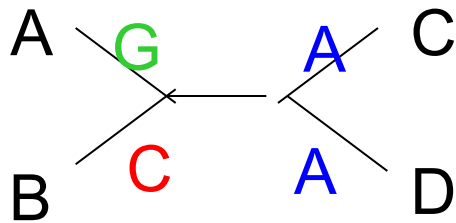
Caractère **variable** mais **non informatif**.

Caractère ne favorisant aucune topologie par rapport à une autre.

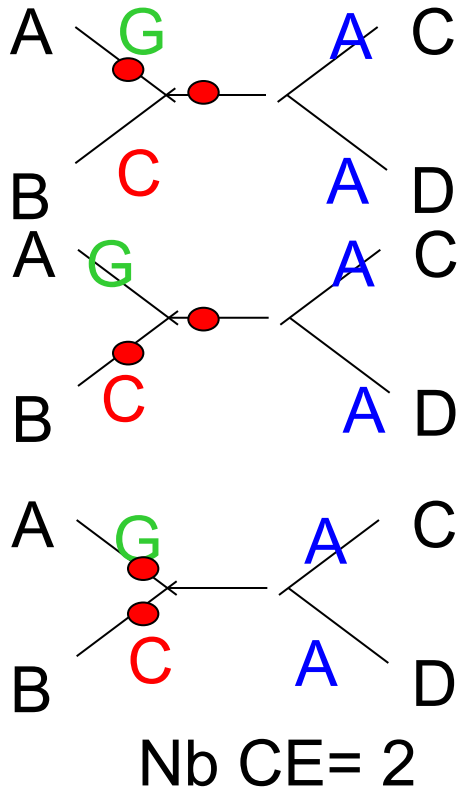


Étude du caractère n°3

	1	2	3	4	5	6	7	8	9
A	A	A	G	A	G	T	T	C	A
B	A	G	C	C	G	T	T	C	T
C	A	G	A	T	A	T	C	C	A
D	A	G	A	G	A	T	C	C	T

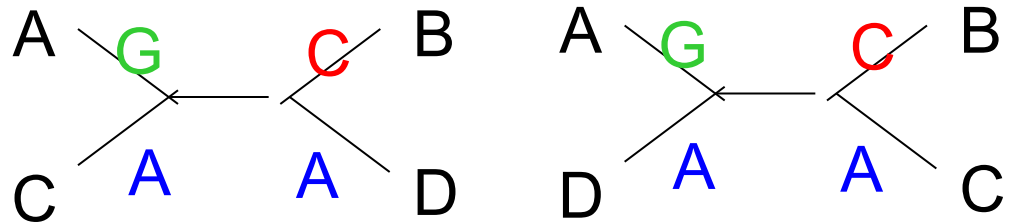


Étude du caractère n°3



Arbre 1

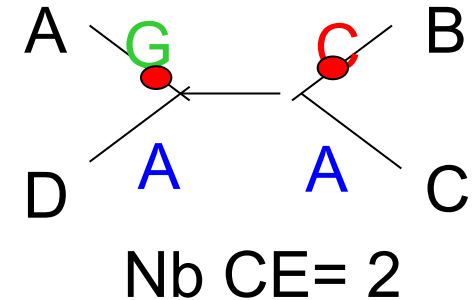
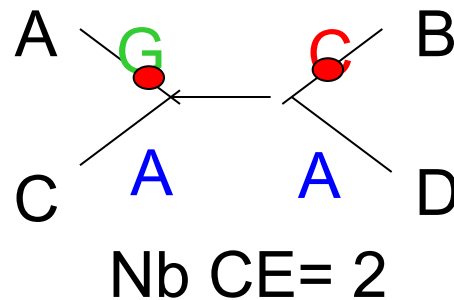
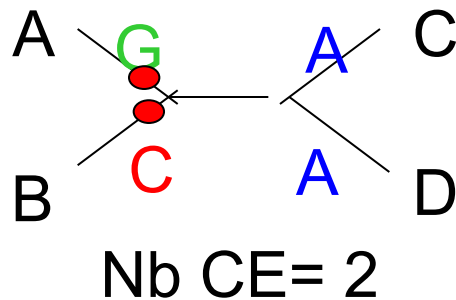
	1	2	3	4	5	6	7	8	9
A	A	A	G	A	G	T	T	C	A
B	A	G	C	C	G	T	T	C	T
C	A	G	A	T	A	T	C	C	A
D	A	G	A	G	A	T	C	C	T



Étude du caractère n°3

	1	2	3	4	5	6	7	8	9
A	A	A	G	A	G	T	T	C	A
B	A	G	C	C	G	T	T	C	T
C	A	G	A	T	A	T	C	C	A
D	A	G	A	G	A	T	C	C	T

Caractère **variable** mais **non informatif**: tous les scénarios « coûtent » 2 CE.
 Caractère ne favorisant aucune topologie par rapport à une autre.

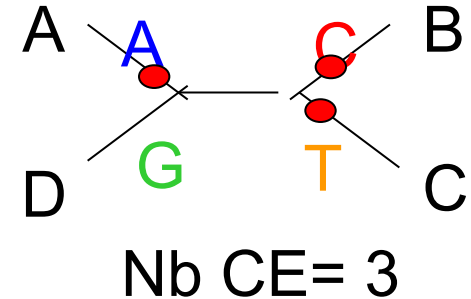
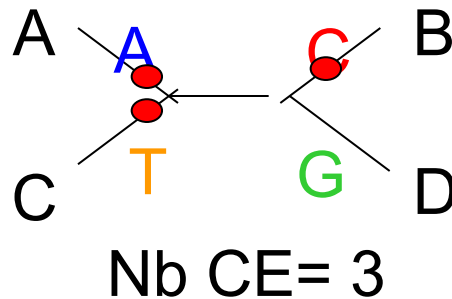
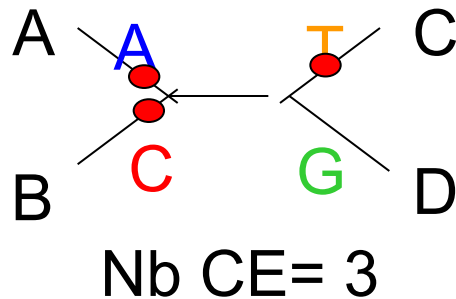


Étude du caractère n°4

	1	2	3	4	5	6	7	8	9
A	A	A	G	A	G	T	T	C	A
B	A	G	C	C	G	T	T	C	T
C	A	G	A	T	A	T	C	C	A
D	A	G	A	G	A	T	C	C	T

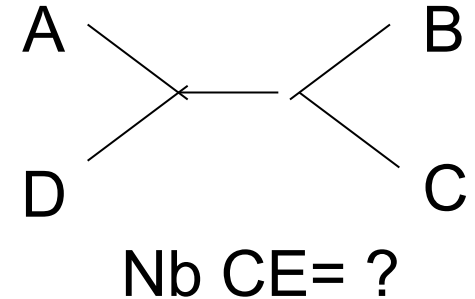
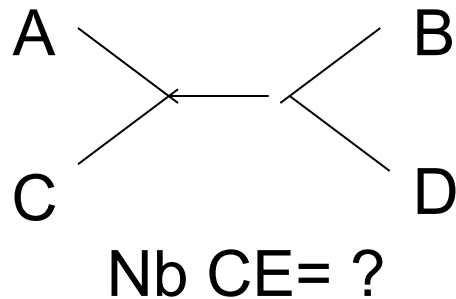
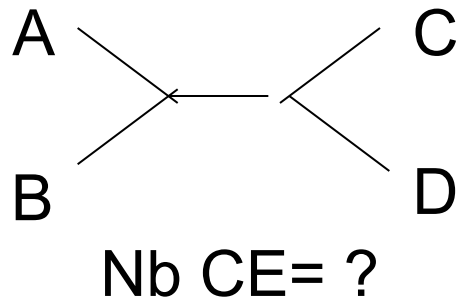
Caractère **variable** mais **non informatif**.

Caractère ne favorisant aucune topologie par rapport à une autre.



Étude du caractère n°5

	1	2	3	4	5	6	7	8	9
A	A	A	G	A	G	T	T	C	A
B	A	G	C	C	G	T	T	C	T
C	A	G	A	T	A	T	C	C	A
D	A	G	A	G	A	T	C	C	T

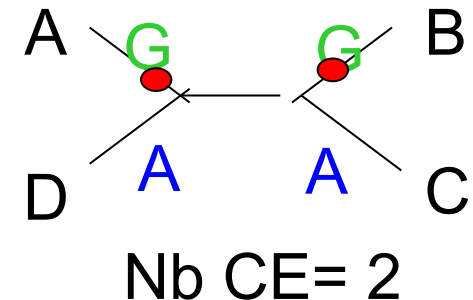
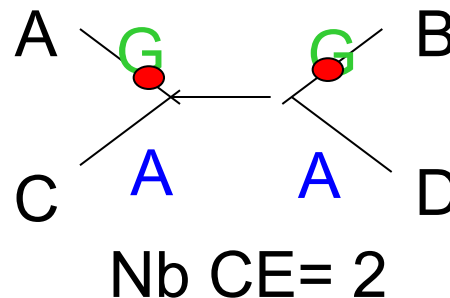
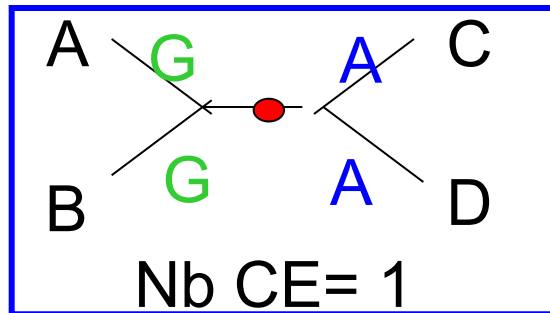


Étude du caractère n°5

	1	2	3	4	5	6	7	8	9
A	A	A	G	A	G	T	T	C	A
B	A	G	C	C	G	T	T	C	T
C	A	G	A	T	A	T	C	C	A
D	A	G	A	G	A	T	C	C	T

Caractère variable et **informatif** (au moins 2 états de caractère sont partagés par au moins 2 OTU).

Caractère favorisant la première topologie par rapport aux deux autres.



- Caractère ***invariant***: toutes les OTU possèdent le même état de caractères pour un site donné
- Caractère ***variable***
 - ▣ ***Non informatif*** si les états de caractères à ce site ne favorisent aucune topologie parmi l'ensemble des topologies possibles
 - ▣ ***Informatif*** si les états de caractères à ce site favorise une (ou plusieurs) topologie(s) parmi l'ensemble des topologies possibles

Maximum de parcimonie - Méthode

On compte ensuite le nombre total de mutations nécessaires pour chaque topologie.

Bilan:

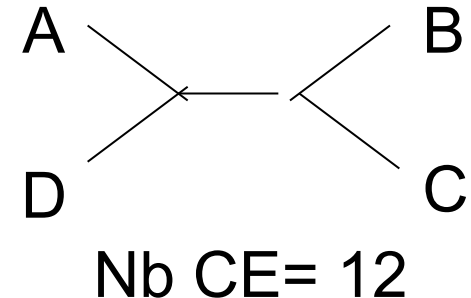
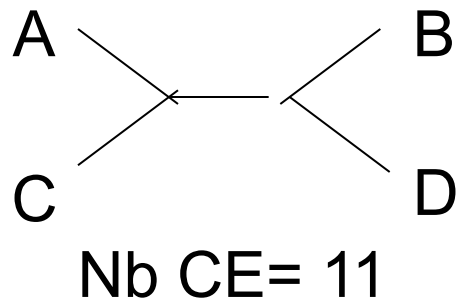
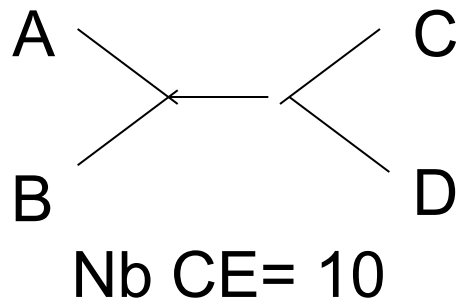
$$T1 = 0+1+2+3+1+0+1+0+2=10$$

$$T2 = 0+1+2+3+2+0+2+0+1=11$$

$$T3 = 0+1+2+3+2+0+2+0+2=12$$

	1	2	3	4	5	6	7	8	9
A	A	A	G	A	G	T	T	C	A
B	A	G	C	C	G	T	T	C	T
C	A	G	A	T	A	T	C	C	A
D	A	G	A	G	A	T	C	C	T

L'arbre le plus parcimonieux = arbre 1



Maximum de parcimonie - désavantages

- Le nombre d'arbres possibles augmente rapidement avec le nombre d'UTOs (séquences).
 - Dans les exemples qui précèdent nous avons analysé 4 séquences.
 - Pour analyser ne fût-ce que 20 séquences, on se trouve confronté à un nombre astronomique de possibilités.
- La parcimonie repose intrinsèquement sur une hypothèse de l'horloge moléculaire => suppose que toutes les branches ont évolué avec la même vitesse.
- Cette méthode fonctionne seulement avec les séquences très conservées.

n	Nb arbres enracinés	Nb arbres non-enracinés
2	1	1
3	3	1
4	15	3
5	105	15
6	945	105
7	10,395	945
8	135,135	10,395
9	2,027,025	135,135
10	3.45E+07	2,027,025
11	6.55E+08	3.45E+07
12	1.37E+10	6.55E+08
13	3.16E+11	1.37E+10
14	7.91E+12	3.16E+11
15	2.13E+14	7.91E+12
16	6.19E+15	2.13E+14
17	1.92E+17	6.19E+15
18	6.33E+18	1.92E+17
19	2.22E+20	6.33E+18
20	8.20E+21	2.22E+20

$$N_R = \frac{(2n! 3)!}{2^{n!2} (n! 2)!}$$

$$N_U = \frac{(2n! 5)!}{2^{n!3} (n! 3)!}$$

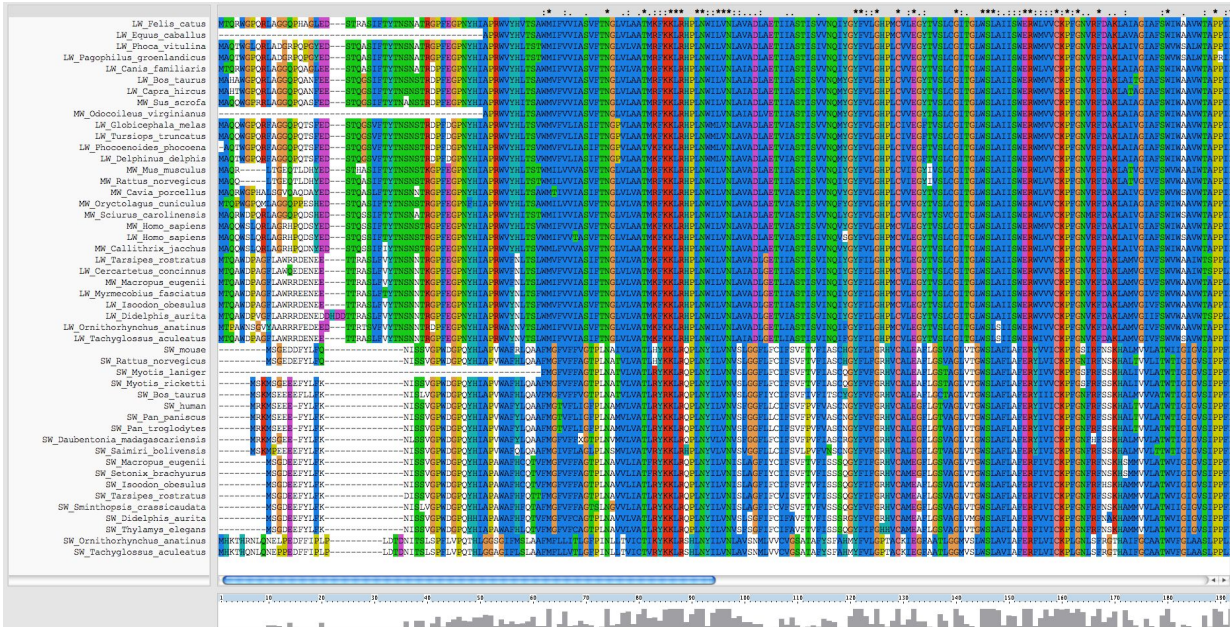
- Résumé des méthodes d'inférence d'arbre implémentées dans PHYLIP.
- Note: le temps de calcul augmente drastiquement quand on passe de méthodes de voisinage (NJ, UPGMA: temps quadratique) aux méthodes de kitch ou fitch (puissance 4 de la longueur des séquences).

Phylip program	method	rooted tree	time	accuracy	remarks
fitch	Fitch-Margoliah	no	$O(n^4)$	higher	loss of accuracy when the tree contains long branches
kitsch	Fitch-Margoliah	yes	$O(n^4)$	higher	
neighbor	neighbour-joining	no	$O(n^2)$	lower	suitable when rate of evolution varies among branches
neighbor	UPGMA	yes	$O(n^2)$	lower	assumes constant rate of evolution along the branches

Evaluation de la robustesse de l'inférence: le "bootstrap"

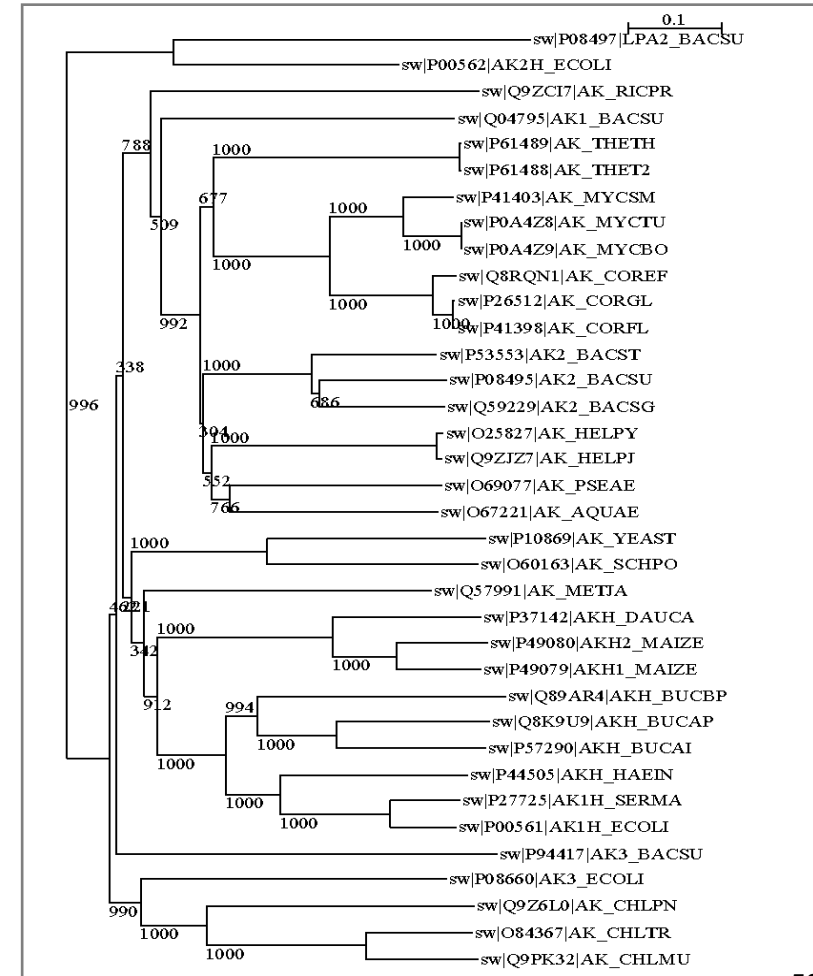
Quelle est la fiabilité d'un arbre inféré ?

- On se base sur les colonnes d'un alignement multiple pour inférer un arbre phylogénétique cohérent avec les différences entre groupes de séquences.
- Cependant, selon les colonnes choisies on peut observer des variations de séquence qui touchent des sous-groupes différents.
- Comment évaluer la robustesse de l'inférence par rapport aux particularité des échantillons disponibles ?



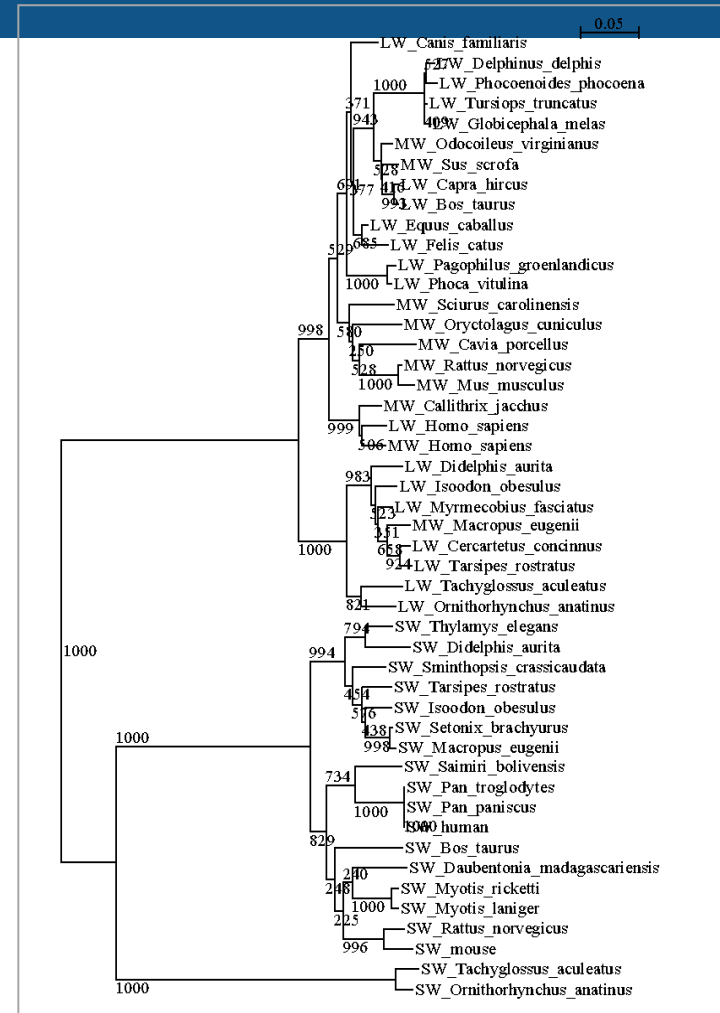
Bootstrapping

- Dans certains cas, les données ne permettent pas d'inférer la phylogénie de façon fiable.
- Pour évaluer la fiabilité de l'inférence, on peut appliquer la méthode du **bootstrapping**.
 - Etant donné un alignement de N séquences et M colonnes, on effectue une sélection aléatoire de M colonnes **avec remise**. Certaines colonnes sont donc tirées plusieurs fois, et d'autres aucune fois.
 - On calcule un arbre avec les colonnes échantillonnées.
 - On répète l'opération un bon nombre de fois (1000), et on compte le nombre de fois où chaque branchement de l'arbre original se reproduit.

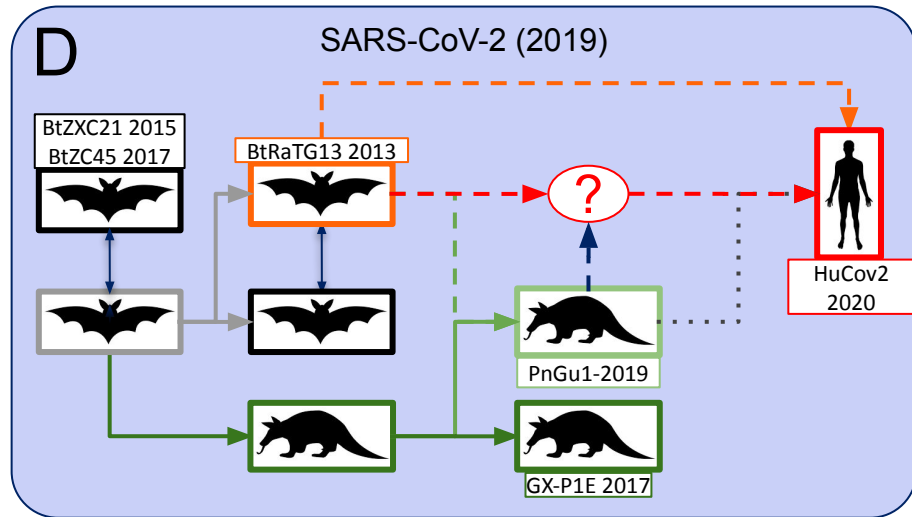
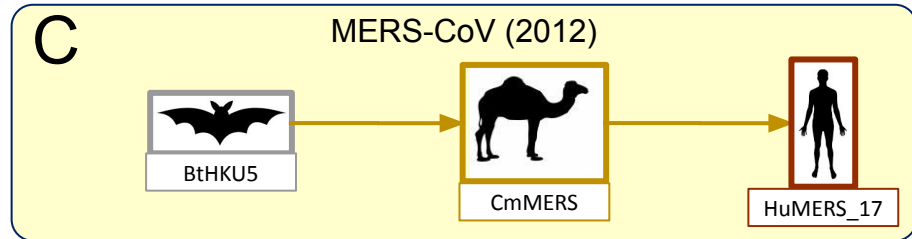
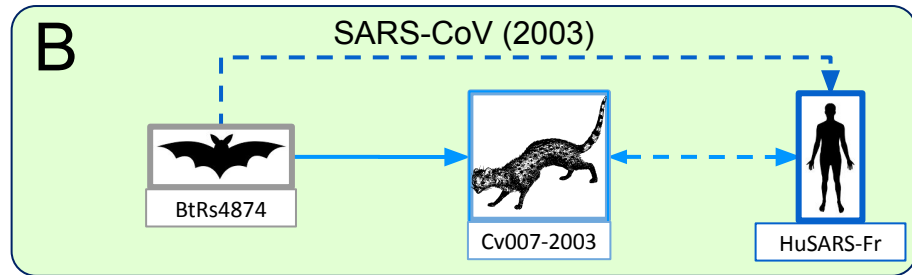
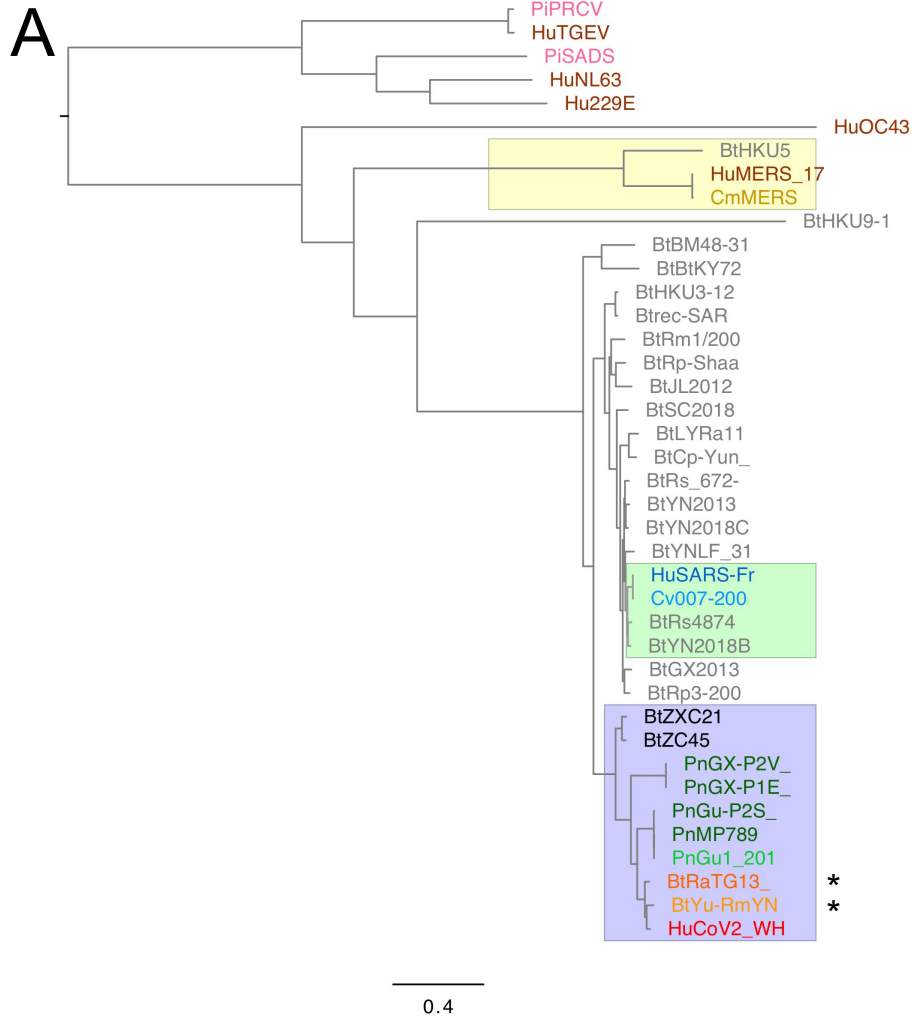


Bootstrapping

- Le phylogramme permet d'identifier les relations entre longueurs des branches et valeurs de bootstrap.
- Les valeurs de bootstrap sont cependant moins faciles à lire que sur un cladogramme (où toutes les branches ont la même longueur).

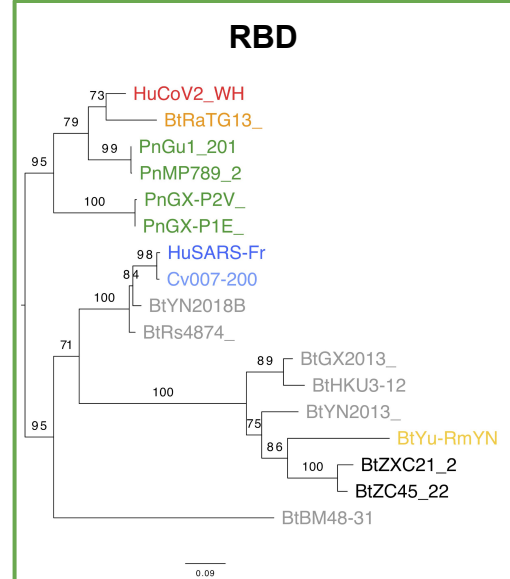
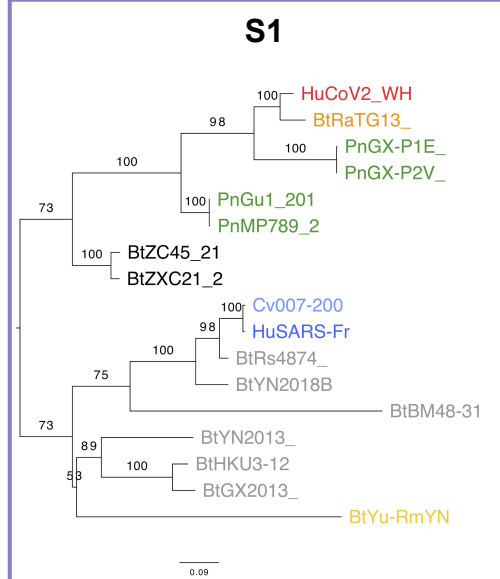
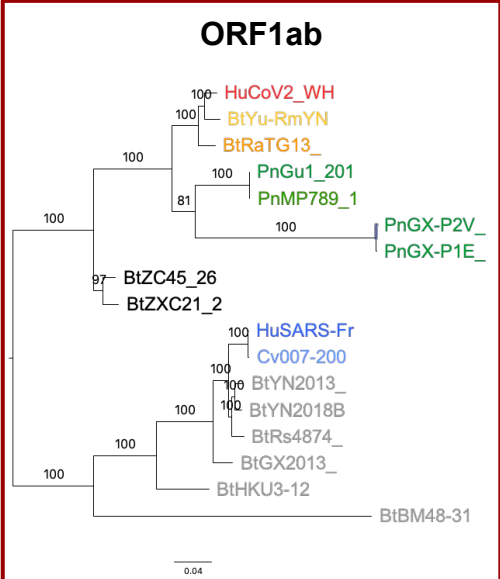
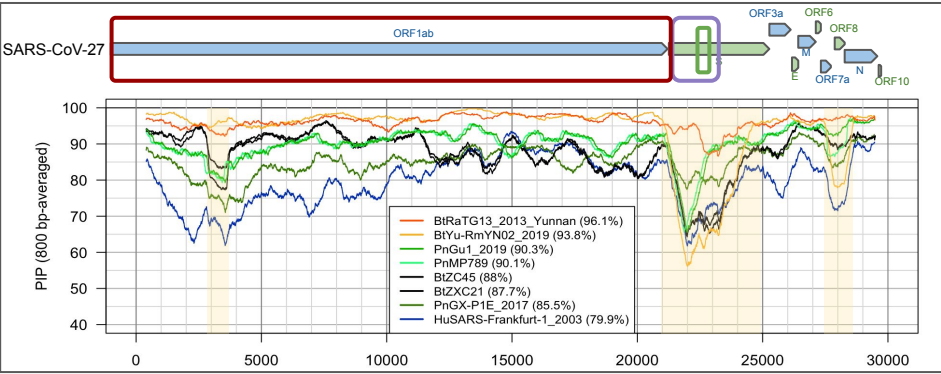


Phylogénie des coronavirus

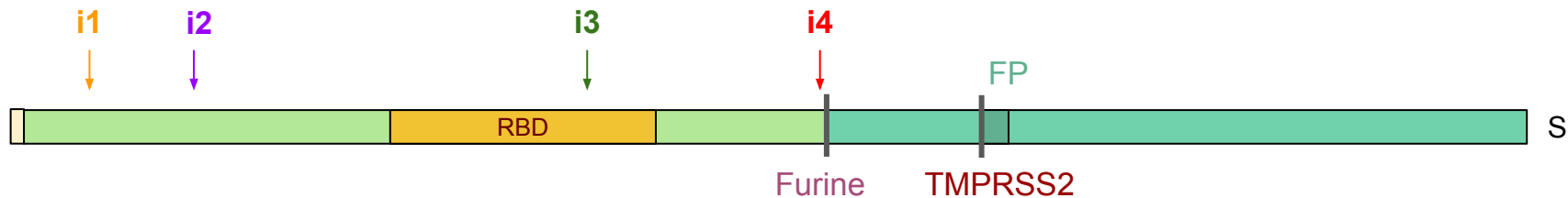


Comparaison entre coronavirus - gène S

Profils de pourcentages de positions identiques (PPI) entre différents génomes de coronavirus (d'humain, de chauve-souris et de pangolin) et SARS-CoV-2 (la référence qui correspond à 100% sur toute la largeur).

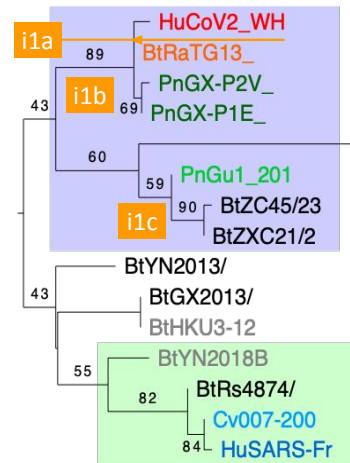


Insertion 1



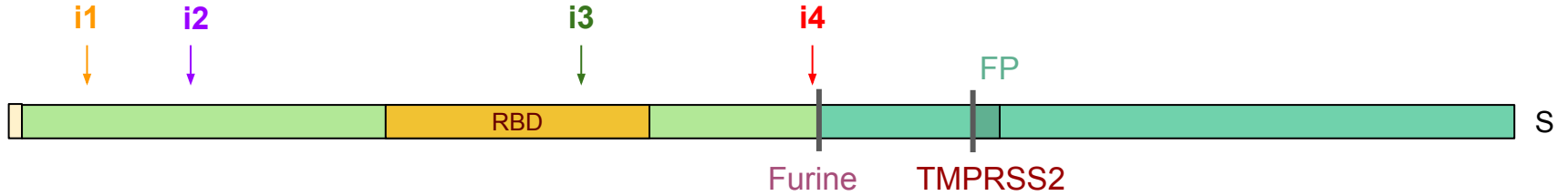
Strain	Sequence
HuCoV2_WH01_2019	WFHAIHVSGTNGTKRFDNP
BtRaTG13_2013_Yunnan	WFHAIHVSGTNGIKRFDNP
PnGX-P1E_2017	WFNTI--NYQGGFKKFDNP
PnGX-P2V_2018	WFNTIHLNYQGGFKKFDNP
PnGu1_2019	WYYAL-TKTNSAEKRVDNP
BtZC45	WYYSL-TTNNAATKRFDNP
BtZXC21	WYYSL-TTNNAATKRFDNP
BtYu-RmYN02_2019	WYNFW-----NQAYTSR
BtYN2013	QFFTQ-----GTNIDNP
BtHKU3-12	QYFSL-NVSDRYTYFDNP
BtGX2013	QYFSL-NVSDRYTYFDNP
BtYN2018B	RFITF-----GLNFDNP
BtRs4874	GFHTI-----NHRFDNP
Cv007-2004	GFHTI-----NHTFDNP
HuSARS-Frankfurt-1_2003	GFHTI-----NHTFGNP

.....80.....90

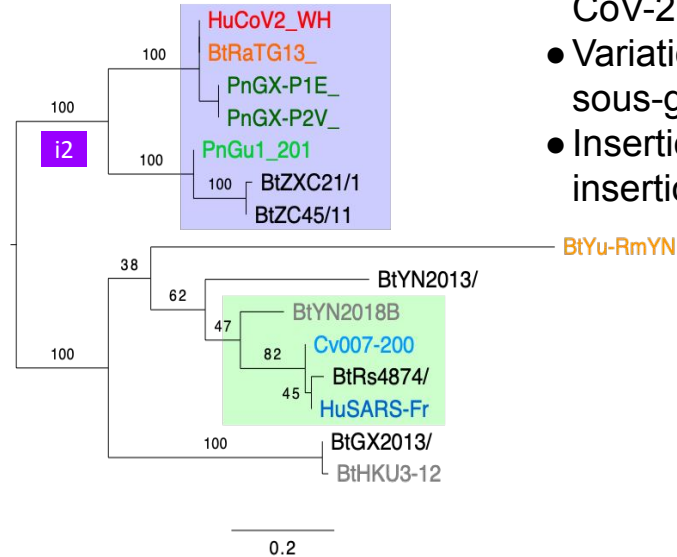


- Insertion retrouvée dans tous les génomes de la branche CoV-2.
- Variations selon les sous-groupes.
- Insertions indépendantes ou insertion suivie de mutations ?

Insertion 2

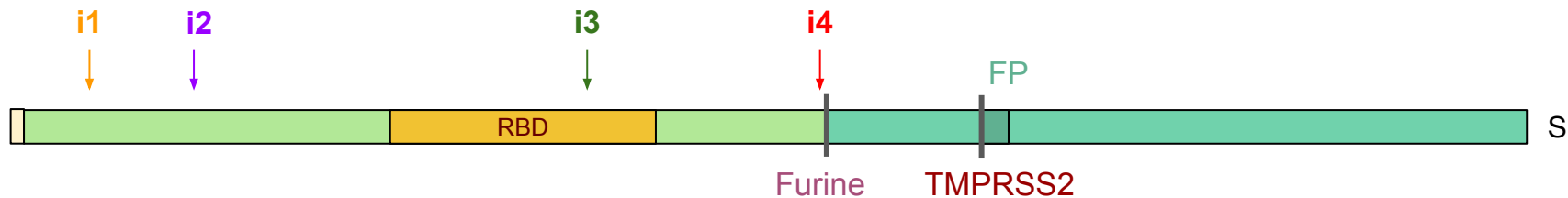


HuCoV2_WH01_2019	YYHKNKNSWMESEFRVYSS	
BtRaTG13_2013_Yunnan	YYHKNKNSWMESEFRVYSS	
PnGX-P1E_2017	YYHNNKNTWVENEFRVYSS	
PnGX-P2V_2018	YYHNNKNTWVENEFRVYSS	
PnGu1_2019	YYH-NNKTSSTREFAVYSS	i2
BtZC45	YYH-NNKTSIREFAVYSS	
BtZXC21	YYH-NNKTSIREFAVYSF	
BtYu-RmYN02_2019	AGGQOTSAA-----VYIS	
BtYN2013	FKSNNSQLSH-----LFS	
BtHKU3-12	SRGTQONAW-----VYQS	
BtGX2013	SRGTQONS-----VYQS	
BtYN2018B	LRSNNTQIPSY----IFNN	
BtRs4874	SKPTGTQTHM----IFDN	
Cv007-2004	SKPMGTQTHM----IFDN	
HuSARS-Frankfurt-1_2003	SKPMGTQTHM----IFDN	
160.....170	

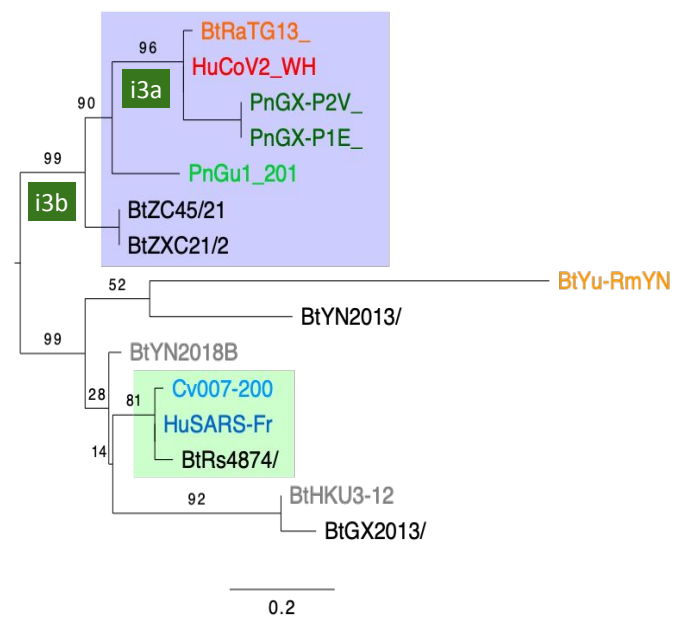


- Insertion retrouvée dans tous les génomes de la branche CoV-2.
- Variations selon les sous-groupes.
- Insertions indépendantes ou insertion suivie de mutations ?

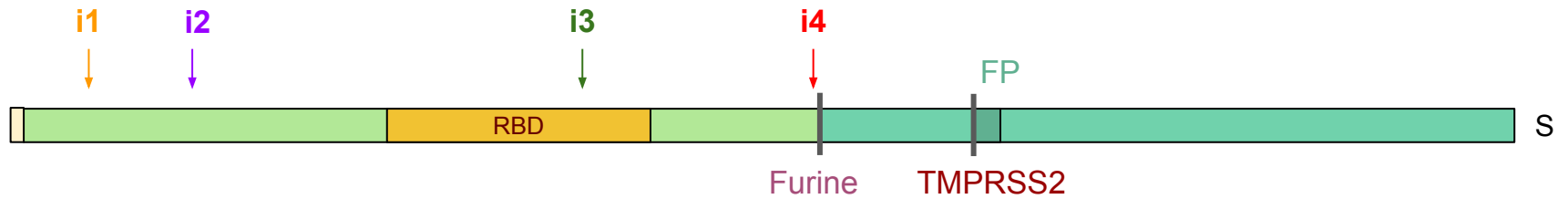
Insertion 3



HuCoV2_WH01_2019	LLALHRSYLTPGDS	SSGWTA
BtRaTG13_2013_Yunnan	LLALHRSYLTPGDS	SSGWTA
PnGX-P1E_2017	LLALHRSYLTPGKLES	GWTT
PnGX-P2V_2018	LLALHRSYLTPGKLES	GWTT
PnGu1_2019	LLTIHRGDPMP---	NNGWTV
BtZC45	LLTIHRGDPMP---	NNGWTA
BtZXC21	LLTIHRGDPMS---	NNGWTA
BtYu-RmYN02_2019	VLTF-----	RSNSQP
BtYN2013	FLAVYRVA-----	AGSISV
BtHKU3-12	VMAMFSQT-----	TSNFLP
BtGX2013	VMAMFSQS-----	TSNFLP
BtYN2018B	LLTAFPPN-----	PGYWGT
BtRs4874	ILTAFSPA-----	QDTWGT
Cv007-2004	ILTAFSPA-----	QGTWGT
HuSARS-Frankfurt-1_2003	ILTAFSPA-----	QDIWGT
260.....	270..

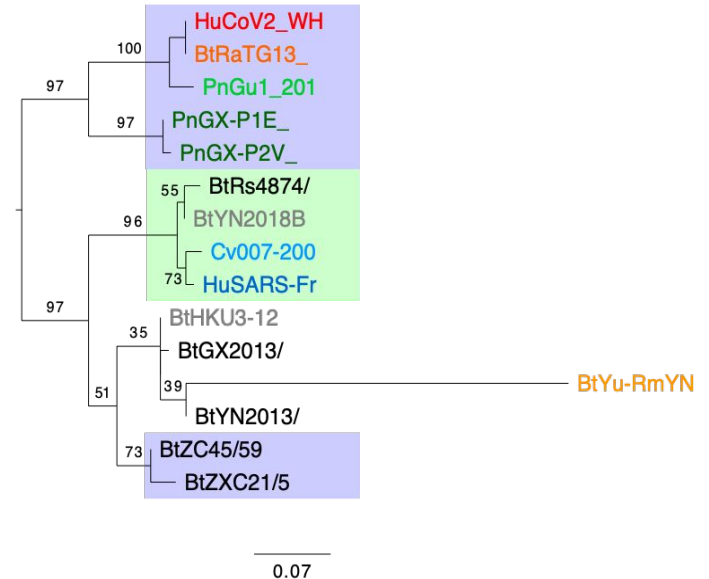


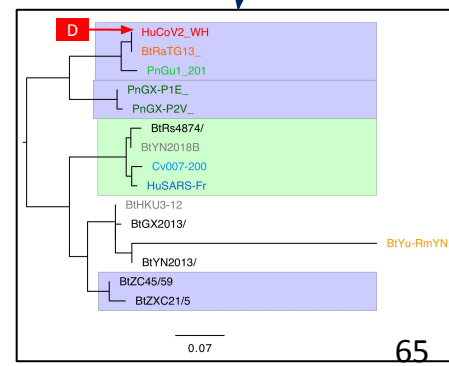
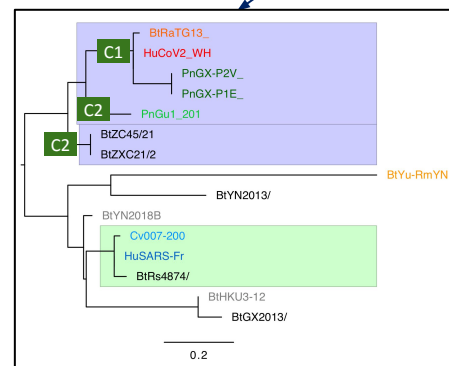
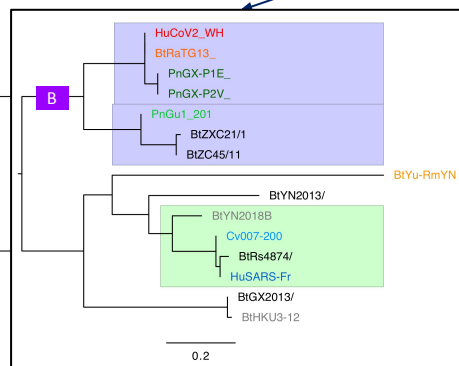
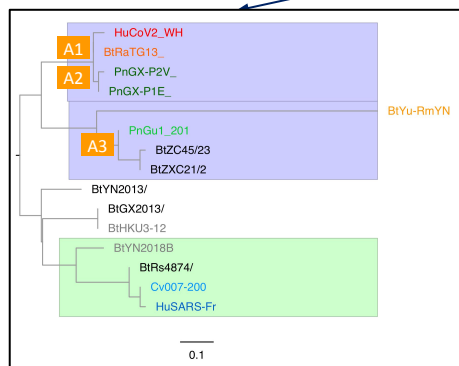
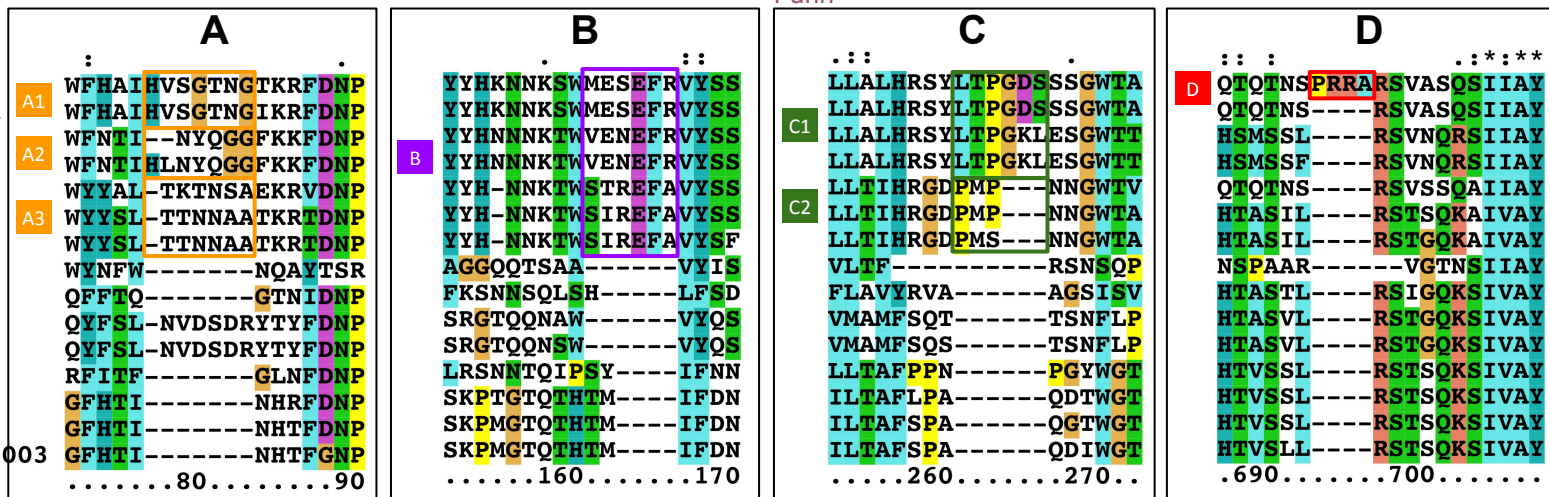
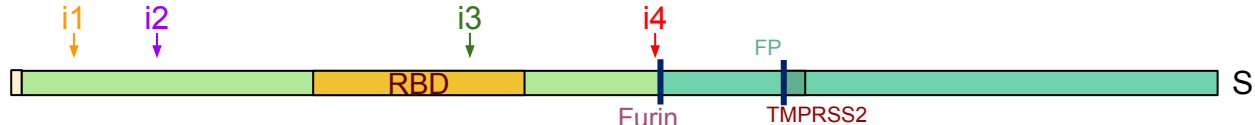
Insertion 4



Sequence	Conservation	Insertion Site
HuCoV2_WH01_2019	QTQTNSPRRRARSVASQSIIAY	i4
BtRaTG13_2013_Yunnan	QTQTNS----RSVASQSIIAY	
PnGX-P1E_2017	HSMSSL----RSVNORSIIAY	
PnGX-P2V_2018	HSMSSF----RSVNORSIIAY	
PnGu1_2019	QTQTNS----RSVSSQAI IAY	
BtZC45	HTASIL----RSTSOKAIVAY	
BtZXC21	HTASIL----RSTGOKAIVAY	
BtYu-RmYN02_2019	NSPAAR-----VGTNSIIAY	
BtYN2013	HTASTL----RSIGOKSIVAY	
BtHKU3-12	HTASVL----RSTGOKSIVAY	
BtGX2013	HTASVL----RSTGOKSIVAY	
BtYN2018B	HTVSLL----RSTSOKSIVAY	
BtRs4874	HTVSLL----RSTSOKSIVAY	
Cv007-2004	HTVSLL----RSTSOKSIVAY	
HuSARS-Frankfurt-1_2003	HTVSLL----RSTSOKSIVAY	

.690.....700.....



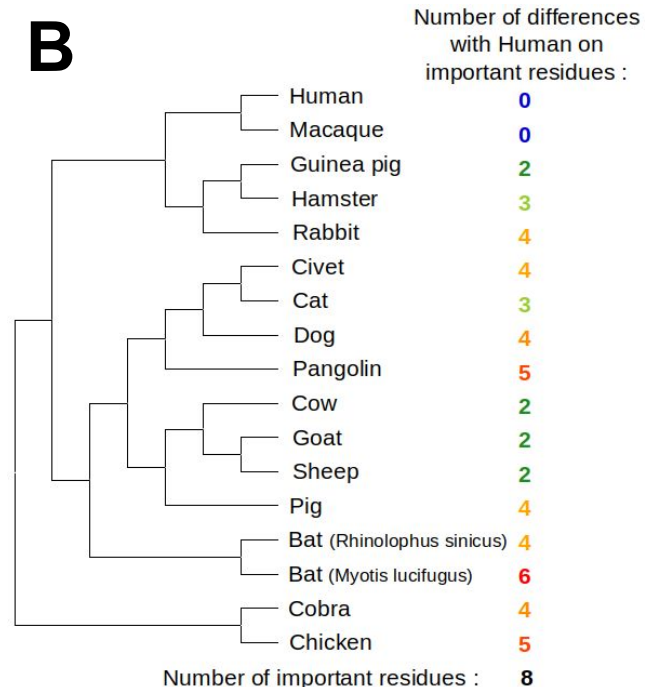


A

ACE2 receptor



S protein

B

Que peut-on conclure à ce stade ? Discussion en classe (virtuelle)

Que peut-on conclure à ce stade ?

- Les éléments disponibles ne permettent pas de démontrer si le génome est d'origine naturelle ou artificielle
- Comment progresser sur cette question, selon les différentes hypothèses ?
 - Zoonose récente d'origine naturelle : intensifier la collecte d'échantillons animaux (sites naturels, élevages, marchés, laboratoires).
 - Circulation à bas bruit dans les populations humaines: caractériser les virus dans les échantillons prélevés depuis 2013 chez des patients atteints de pneumonies atypiques (en particulier les mineurs de 2013).
 - Virus développé en laboratoire : contrôle du laboratoire par une commission indépendante.

Retracer par analyse des génomes la propagation du virus dans les populations humaines

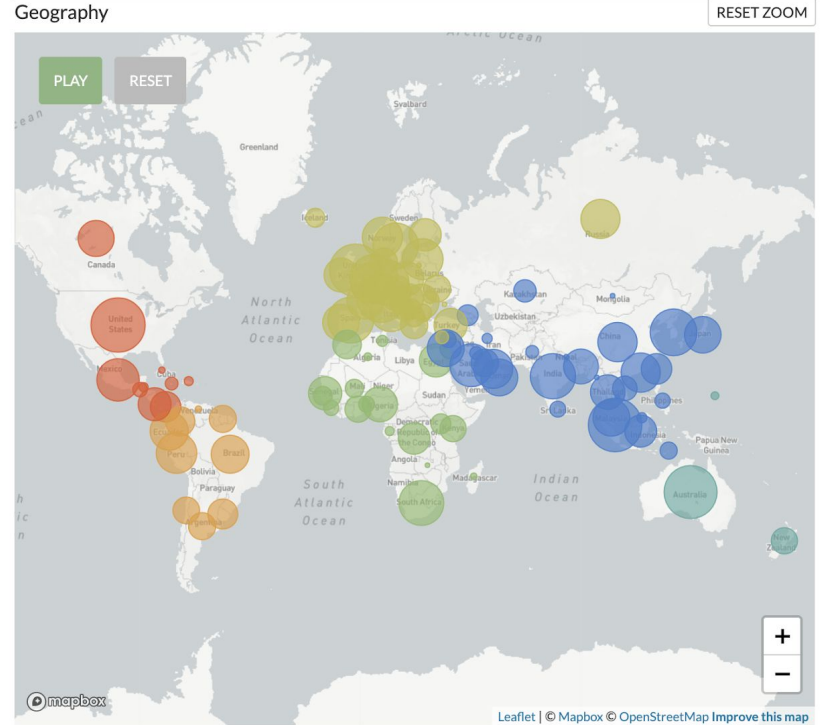
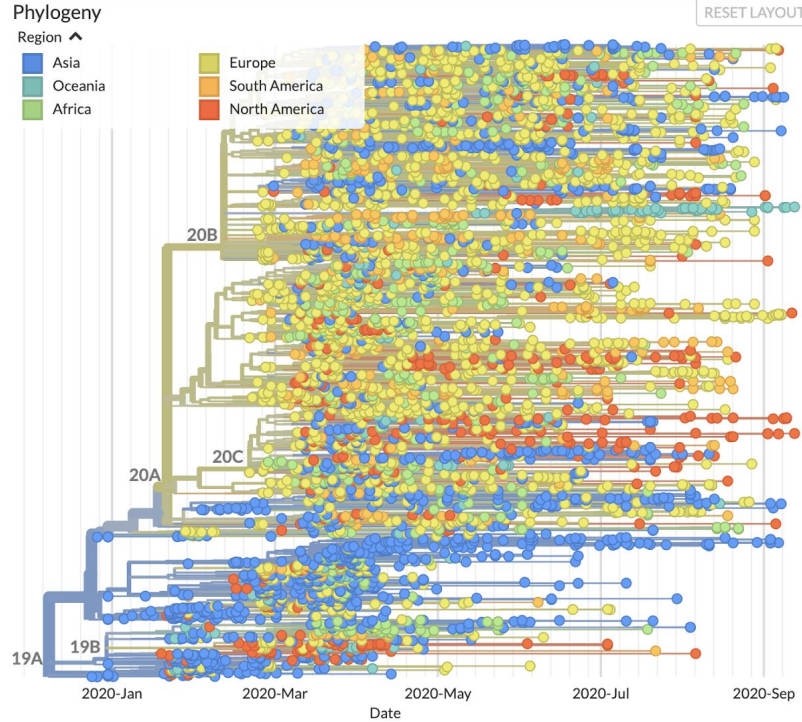
Jacques van Helden

Un outil informatique pour retracer la progression de l'épidémie

Genomic epidemiology of novel coronavirus - Global subsampling

Maintained by the Nextstrain team. Enabled by data from **GISAID**

Showing 4678 of 4678 genomes sampled between Dec 2019 and Sep 2020.



Evolution du virus pendant sa propagation

Review

SARS-CoV-2 and COVID-19: A genetic, epidemiological, and evolutionary perspective



Manuela Sironi^a, Seyed E. Hasnain^b, Benjamin Rosenthal^j, Tung Phan^c, Fabio Luciani^d, Marie-Anne Shaw^e, M. Anice Sallum^f, Marzieh Ezzaty Mirhashemi^g, Serge Morand^h, Fernando González-Candelas^{i,k}, on behalf of the Editors of *Infection, Genetics and Evolution*

M. Sironi, et al.

Infection, Genetics and Evolution 84 (2020) 104384

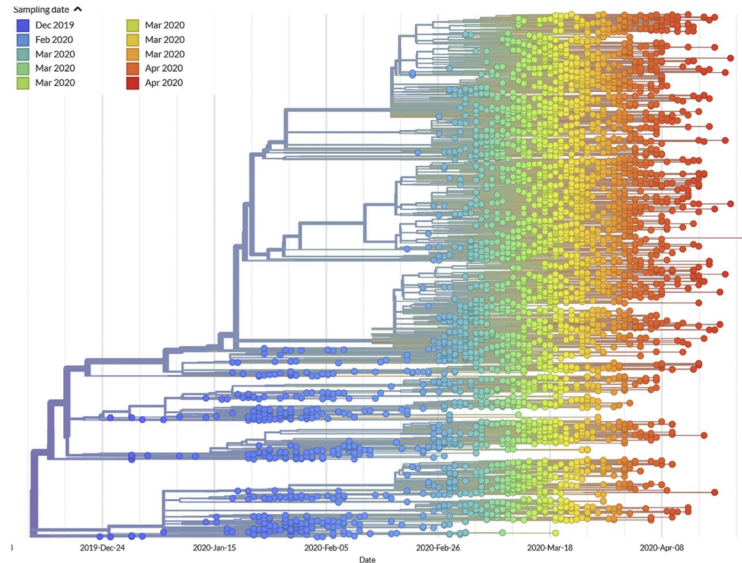


Fig. 3. Time-stamped maximum likelihood phylogenetic reconstruction of SARS-CoV-2 isolates deposited to GISAID.org and rendered by <https://nextstrain.org/ncov>. Isolates are represented by colored circles with the color code corresponding to time of sampling as detailed in the legend.

M. Sironi, et al.

Infection, Genetics and Evolution 84 (2020) 104384

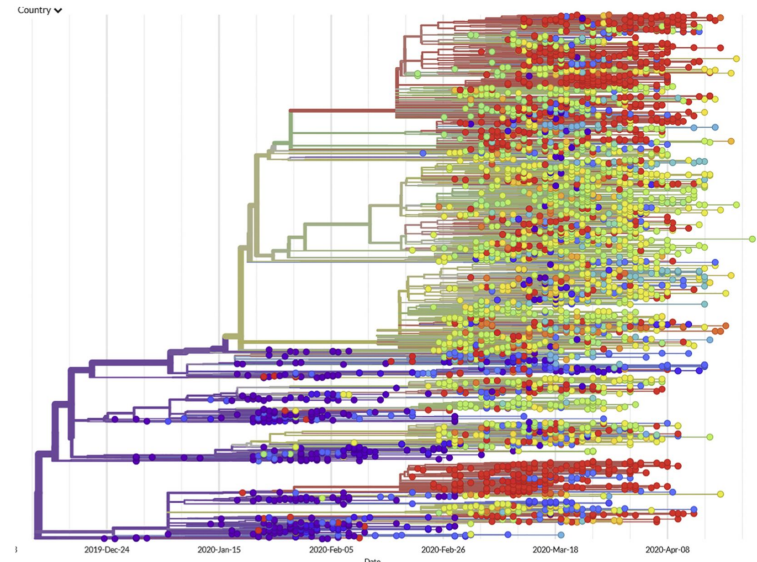


Fig. 4. The same reconstruction of SARS-CoV-2 phylogeny, now denoted by geography. Isolates originated and initially diversified in China (purple), followed by multiple and independent introductions to Oceania (blue), Europe (green and yellow), and North America (red). Less information is known about Africa, India, South America, and other populations of major concern in the Global South. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

Special Section: SARS-CoV-2

On the origin and continuing evolution of SARS-CoV-2

Xiaolu Tang^{1,†}, Changcheng Wu^{1,†}, Xiang Li^{2,3,4,†}, Yuhe Song^{2,5,†}, Xinmin Yao¹,
Xinkai Wu¹, Yuange Duan¹, Hong Zhang¹, Yirong Wang¹, Zhaohui Qian⁶,
Jie Cui^{ID 2,3,*} and Jian Lu^{ID 1,*}

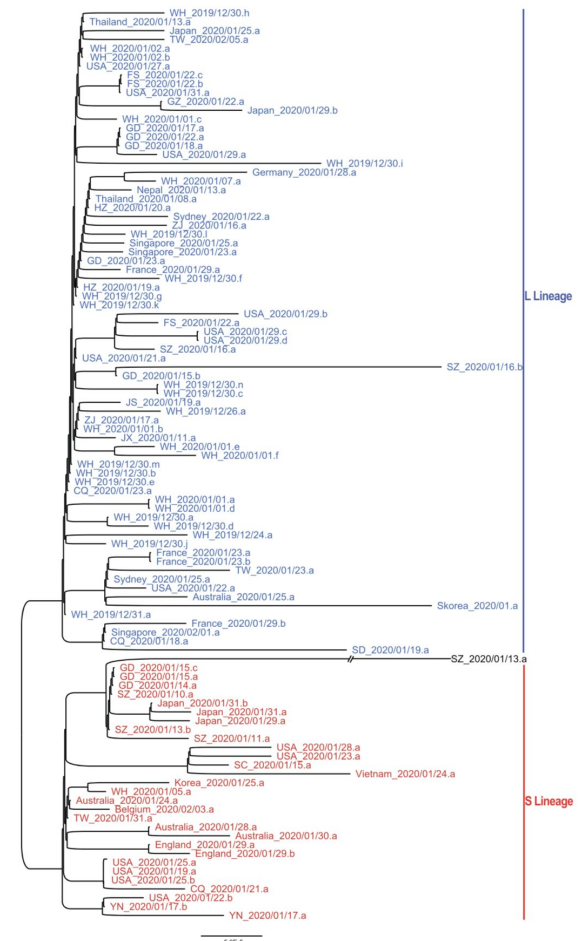
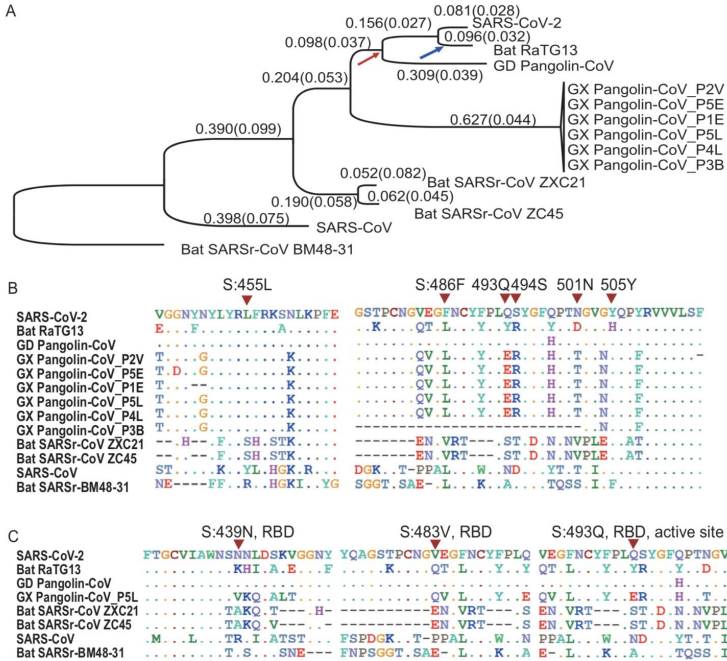


Figure 5. The unrooted phylogenetic tree of the 103 SARS-CoV-2 genomes. The ID of each sample is the same as in Fig. 4A. Note WH/2019/12/31.a represents the reference genome (NC_045512). Note SZ/2020/01/13.a had C at both positions 8,782 and 28,144 in the genome, belonging to neither L nor S lineage.