
Evolution moléculaire et phylogénie

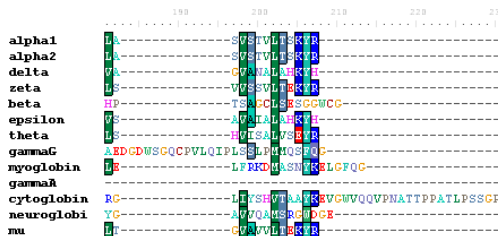
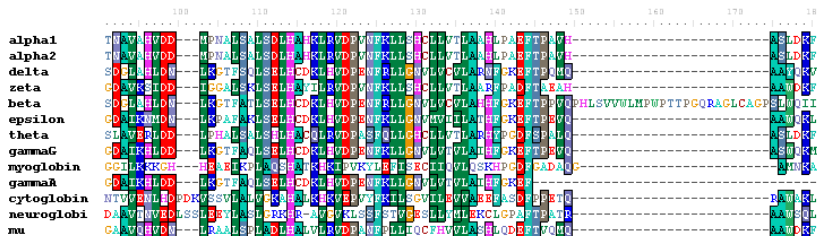
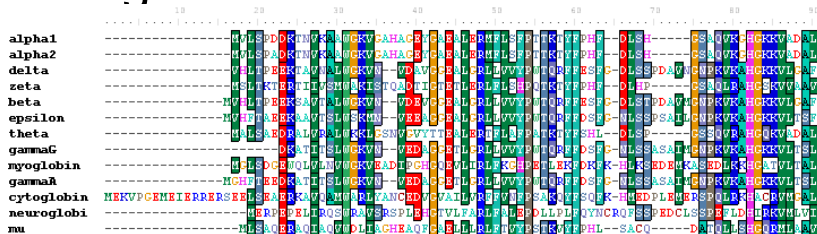
M2 – EcoSciences / BIM INSA

Méthodes de reconstruction phylogénétique

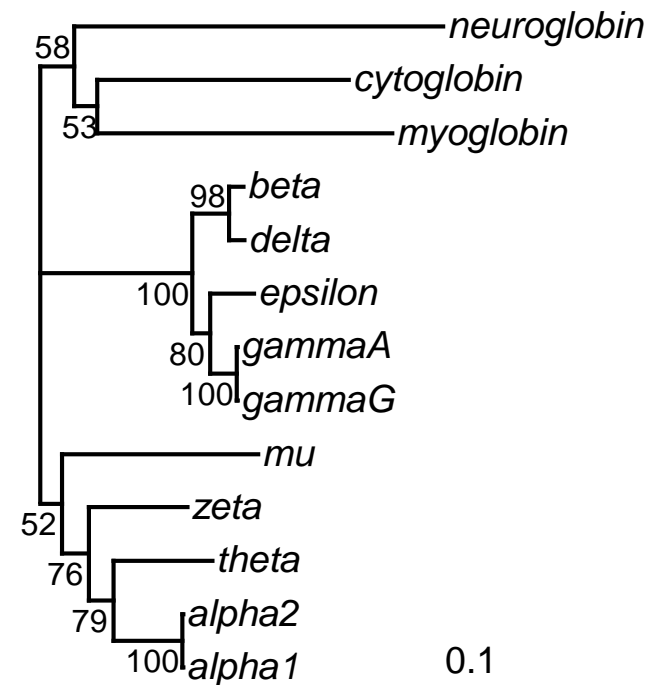
- Quatre grandes familles de méthodes
 - Parcimonie
 - Méthodes de distance
 - Maximum de vraisemblance
 - Méthodes bayésiennes
-

Données utilisées en phylogénie moléculaire

- Point de départ = alignement de séquences homologues
- Arrivée = arbre décrivant les liens évolutifs entre les séquences de l'alignement



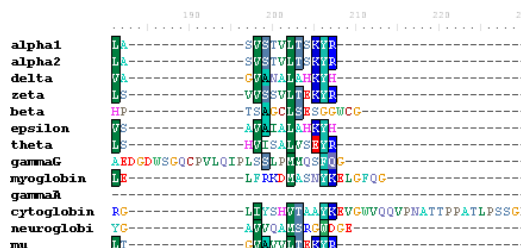
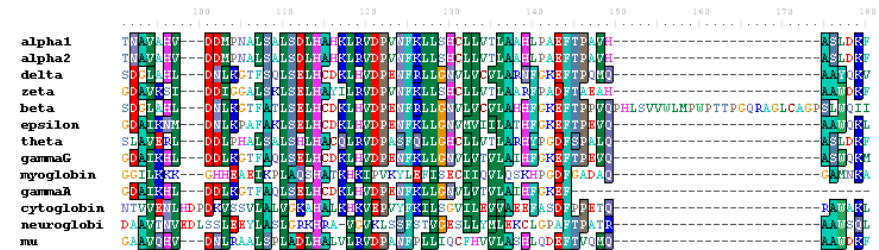
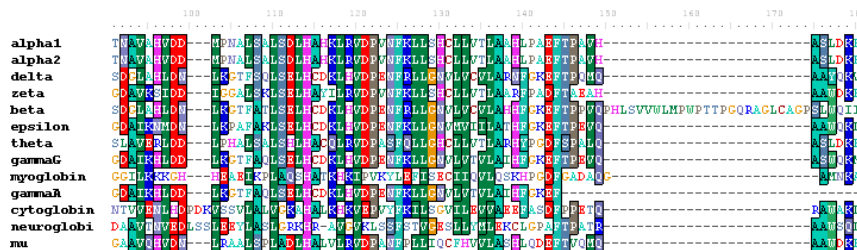
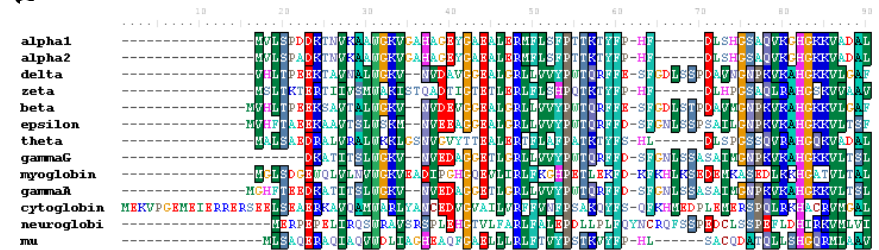
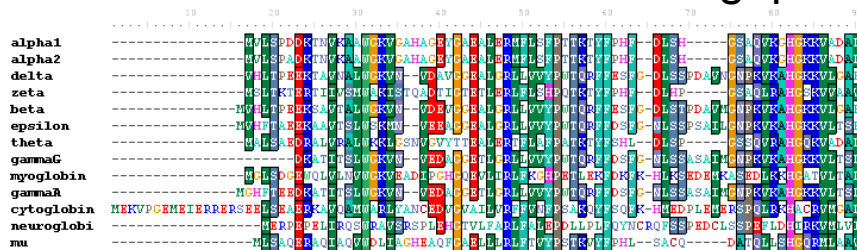
109 / 230 positions
conservées pour l'analyse



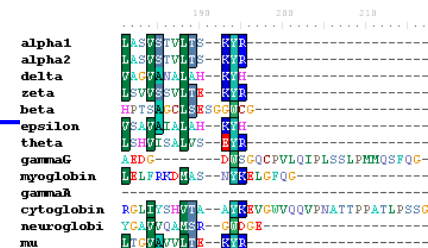
(Alignement des 13 globines humaines réalisé avec clustalW (http://www.frangun.org/HSglobin_A.fasta), arbre construit avec Seaview (BioNJ, 100 réplicats de bootstrap))

Alignements et gaps

- Chaque colonne de l'alignement représente une position (ou site) composée de résidus homologues, cad dérivant d'un même site ancêtre
- La qualité des alignements est essentielle
 - ⇒ Les régions où l'alignement est ambigu doivent être retirées (automatiquement ou manuellement) avant l'analyse phylogénique
- La plupart des méthodes de reconstruction ne prend en compte que les substitutions et non les événements d'insertions/délétions
 - ⇒ Les sites contenant des gaps sont ignorés



230
(ClustalW)



218
(Muscle)

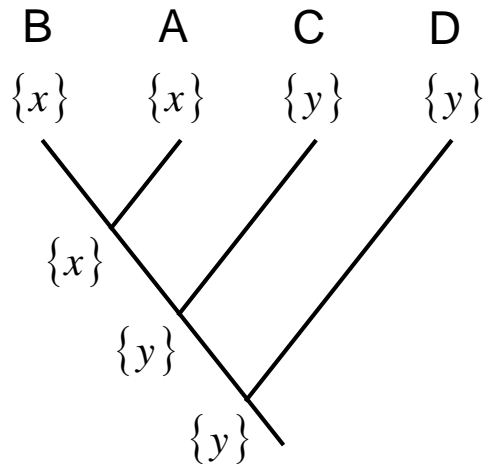
La parcimonie

- Fondement: rasoir d'Occam
 - « Les multiples ne doivent pas être utilisés sans nécessité. »
(pluralitas non est ponenda sine necessitate) ou sous une forme plus moderne « les hypothèses les plus simples sont les plus vraisemblables »

Le critère de parcimonie

- Soit un caractère relevé dans 4 espèces {A,B,C,D} (dont on connaît la phylogénie) et présentant les états de caractères { x, x, y, y }

⇒ Quelle histoire a pu conduire à cet état final?



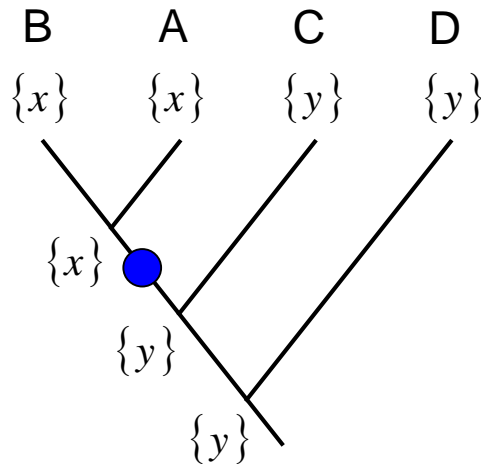
Le critère de parcimonie

- Soit un caractère relevé dans 4 espèces {A,B,C,D} (dont on connaît la phylogénie) et présentant les états de caractères {x, x, y, y}

⇒ Quelle histoire a pu conduire à cet état final?

● Substitution $y \Rightarrow x$

● Substitution $x \Rightarrow y$



$$N_C = 1$$

Similarité par
ascendance commune

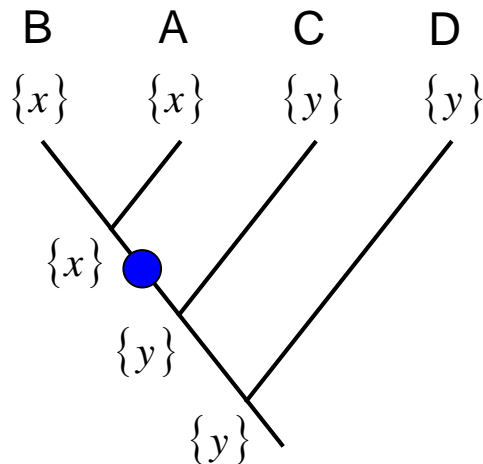
Le critère de parcimonie

- Soit un caractère relevé dans 4 espèces {A,B,C,D} (dont on connaît la phylogénie) et présentant les états de caractères {x, x, y, y}

⇒ Quelle histoire a pu conduire à cet état final?

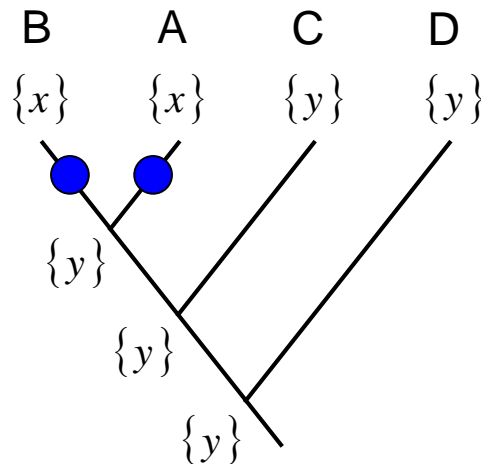
● Substitution $y \Rightarrow x$

● Substitution $x \Rightarrow y$



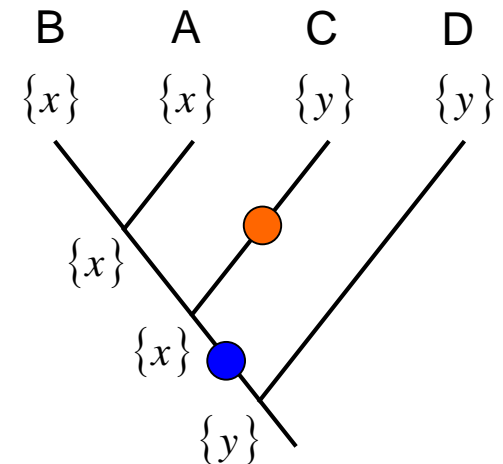
$N_c = 1$

Similarité par
ascendance commune



$N_c = 2$

Similarité par
convergence



$N_c = 2$

Similarité par
réversion

Les scénarios homoplasiques demandent plus de changements évolutifs. L'emploi du critère de parcimonie en phylogénie moléculaire n'est justifié que si les convergences et les réversions sont rares.

Le maximum de parcimonie

- Principe: rechercher parmi l'espace des arbres définissant les liens entre n séquences la topologie qui minimise le nombre de changements évolutifs
 - ⇒ Quelle est la topologie qui implique le moins de changements d'état de caractères pour rendre compte des différences observées entre les UTO étudiées
 - Procédure:
 - 1) pour une topologie T fixée et pour un site donné de l'alignement, calculer (N_C) le nombre de changements évolutifs nécessaires pour expliquer les états de caractères observés
 - 2) calculer (N_C) pour chaque site de l'alignement $\Rightarrow L$, la longueur de l'arbre
 - 3) calculer L pour toutes les topologies T possibles \Rightarrow retenir l'arbre le plus parcimonieux (cad l'arbre le plus court)
-

Parcimonie: Etape 1

- Pour une topologie T fixée et pour un site donné de l'alignement, calculer (N_C) le nombre de changements évolutifs nécessaires pour expliquer les états de caractères observés
-

Algorithme de Fitch: calcul du nombre minimal de changements évolutifs

- Soit une topologie T fixée et racinée de manière arbitraire, soit V l'ensemble de ses nœuds
 - Pour tout $p \in V$ on définit:
 - C_p , le nombre minimal de changements dans le sous-arbre dont p est la racine
 - S_p , l'état de p , cad l'ensemble des résidus en p compatibles avec C_p changements évolutifs dans le sous-arbre raciné par p .
- Soit q et r les deux nœuds fils de p

Algorithme 1 Nombre de substitutions avec Fitch

si $S_q \cap S_r \neq \emptyset$ **alors**

$S_p \leftarrow S_q \cap S_r$

$c_p \leftarrow c_q + c_r$

sinon

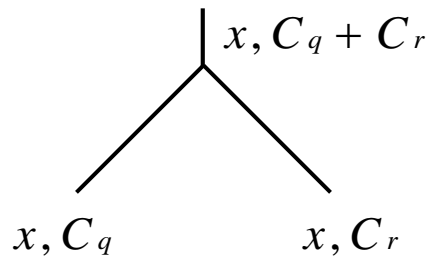
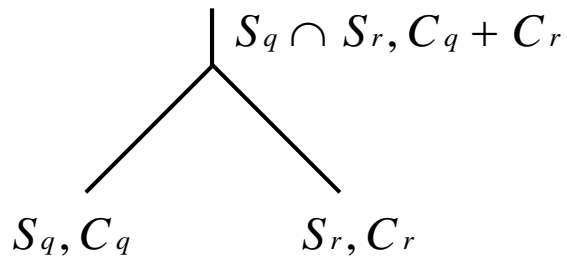
$S_p \leftarrow S_q \cup S_r$

$c_p \leftarrow c_q + c_r + 1$

fin si

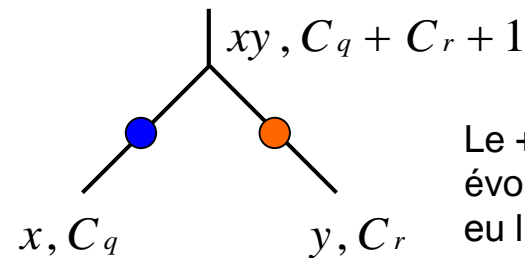
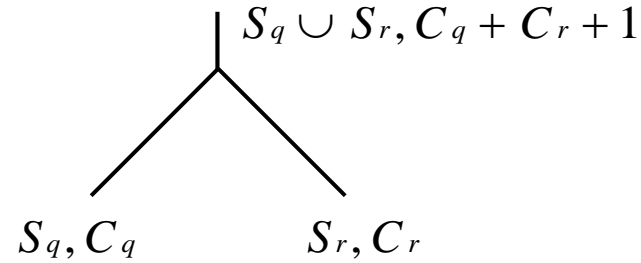
Algorithme de Fitch: Illustration

Cas n°1: $S_q \cap S_r \neq \emptyset$



$\forall x \in S_q \cap S_r$

Cas n°2: $S_q \cap S_r = \emptyset$



$\forall x \in S_q$

$\forall y \in S_r$

Le +1 reflète le changement évolutif qui a nécessairement eu lieu dans l'une des deux branches filles

soit $x \Rightarrow y$ ●

Soit $y \Rightarrow x$ ●

Algorithme de Fitch: Application

Algorithme 1 Nombre de substitutions avec Fitch

si $S_q \cap S_r \neq \emptyset$ alors

$$S_p \leftarrow S_q \cap S_r$$

$$c_p \leftarrow c_q + c_r$$

sinon

$$S_p \leftarrow S_q \cup S_r$$

$$c_p \leftarrow c_q + c_r + 1$$

fin si

La racine est placée de manière arbitraire et n'a aucune influence sur le nombre de changements évolutifs inférés

Les états de caractères inférés aux nœuds ne représentent pas des caractères ancestraux, ni tous les états de caractères possibles !

Initialisation du calcul récursif aux feuilles de l'arbre

-P = {x} = résidu présent à cette feuille
-C_p = 0

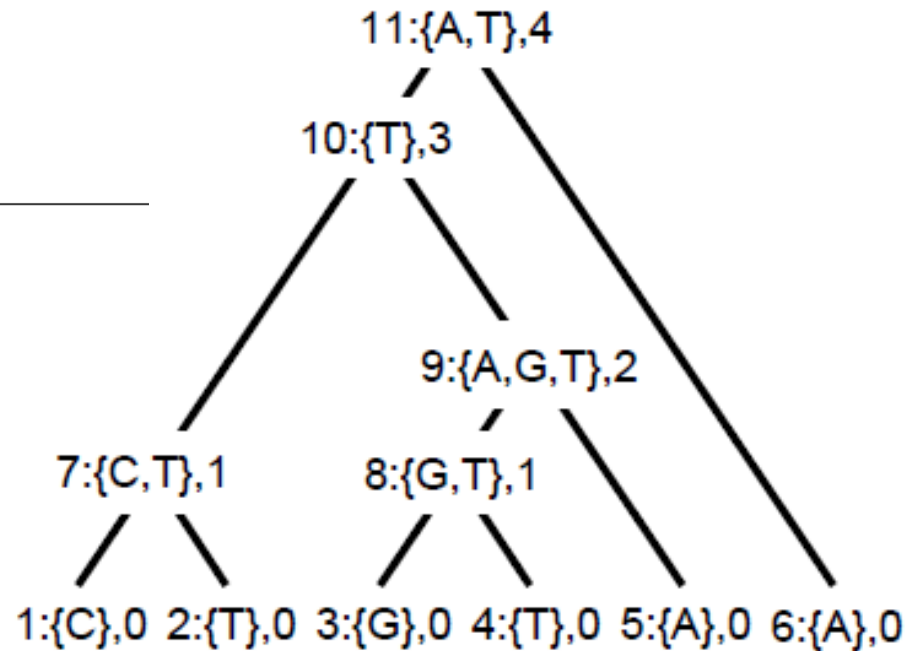
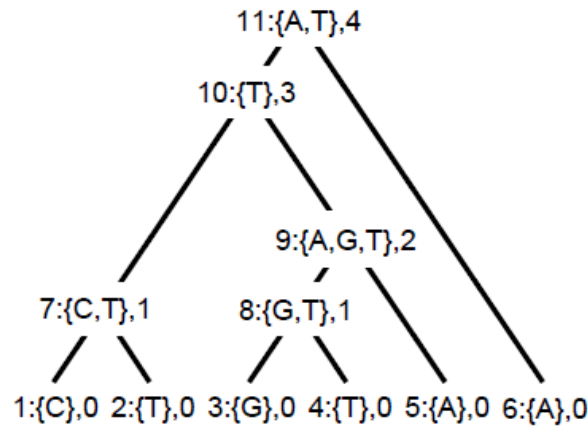


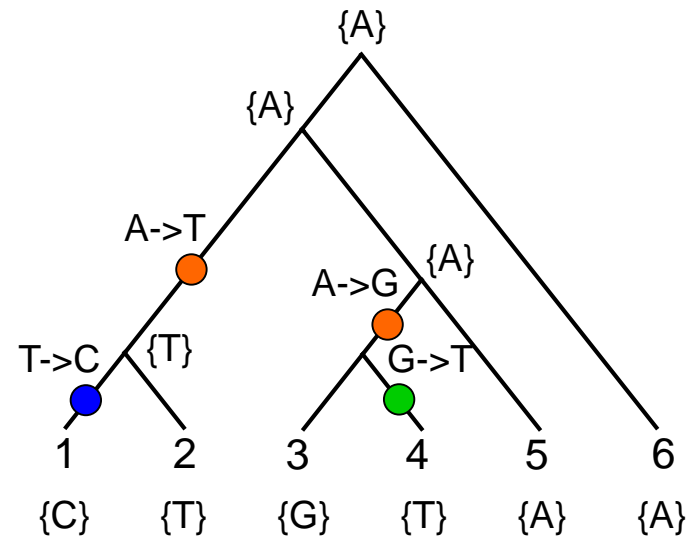
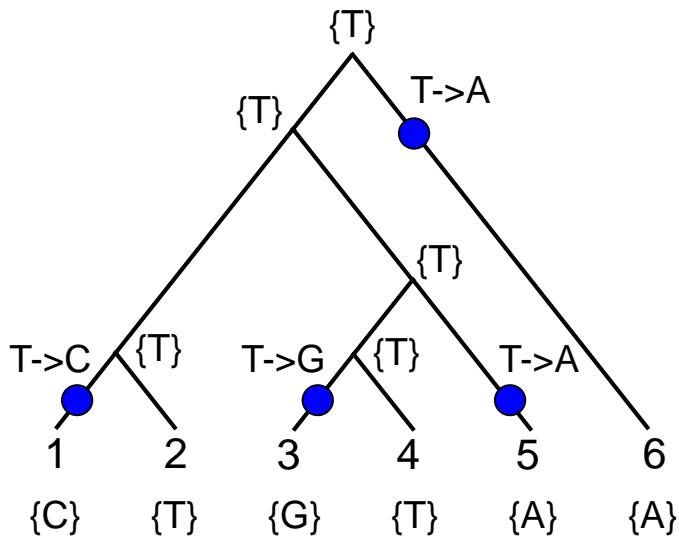
Figure 3.1 – Calcul du nombre de changements pour un site avec l'algorithme de Fitch. Les nœuds sont numérotés de 1 à 11 et les ensembles générés ainsi que le nombre de substitutions inférées sont indiqués. Dans cet exemple, le nombre minimum de changements est égal à 4.

$$N_C = 4$$

Des scénarios multiples



Il existe plusieurs scénarios impliquant $N_C = 4$ changements évolutifs



Parcimonie: Etapes 2 et 3

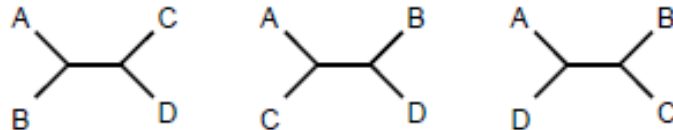
- Etape 2:
 - Calculer N_C pour chaque site de l'alignement
 - Sommer tous les valeurs de N_C pour l'ensemble des sites
 - Calculer L , la longueur totale de l'arbre
 - Etape 3:
 - Répéter l'étape 2 pour chaque topologie T composant l'espace des arbres possibles à n feuilles
 - Retenir l'arbre de longueur L minimale \Leftrightarrow arbre le plus parcimonieux
-

Tous les sites ne sont pas équivalents

- Tous les sites ne contiennent pas une information permettant de discriminer les topologies
 - Les sites constants (1 seul état de caractère)
 - Ne sont pas informatifs
 - Sites variables (au moins 2 états de caractères)
 - Informatifs: présentent au moins deux états de caractères chacun partagés par au moins deux séquences
 - Non informatifs: tous les autres
-

Tous les sites ne sont pas équivalents

- Soit A, B, C et D quatre séquences d'ADN homologues alignées
 - Il existe 3 topologies non racinées possibles
 - Il existe 4 états de caractères {A,T,C,G}
 - Il existe $4^4 = 256$ motifs différents observables à une position
 - ⇒ Seuls 36 sont informatifs, et sont tous du type $\{x,x,y,y\}$, $\{x,y,x,y\}$ ou $\{x,y,y,x\}$ (avec $x \neq y$ et $x,y \in \{A,T,C,G\}$)



ABCD	Topology 1 (A-C, B-D)	Topology 2 (A-B, C-D)	Topology 3 (A-D, B-C)
AAAA	0	0	0
AAAC	1	1	1
AAAG	1	1	1
AAAT	1	1	1
AACA	1	1	1
AACC	1	2	2
AACG	2	2	2
AACT	2	2	2
AAGA	1	1	1
AAGC	2	2	2
AAGG	1	2	2
AAGT	2	2	2
AATA	1	1	1
AATC	2	2	2
AATG	2	2	2
AATT	1	2	2
ACAA	1	1	1
...
TTTT	0	0	0

Figure 3.2 – Sous-ensemble des 256 motifs possibles pour un site dans le cas d'un arbre à quatre UTO (A, B, C et D) et mesure du nombre de substitutions inférées pour chacune des trois topologies non racinées possibles. En gras figurent les cas où ce nombre est différent entre plusieurs topologies.

Algorithme de Fitch: Reconstruction de séquences ancestrales

- La reconstruction de séquences ancestrales par l'algorithme de Fitch part de la racine et descend vers les feuilles
- Les notations $p, q, r, S_p, S_q, S_r, C_p, C_q,$ et C_r sont identiques à celles de l'algorithme de Fitch permettant de calculer le nombre de changements évolutifs
- Soit F_a , l'ensemble des états de caractères possibles du nœud a qui est l'ancêtre immédiat du nœud p , alors il est possible de déterminer F_p , l'ensemble des états de caractères ancestraux au nœud p .

Algorithme 2 Séquences ancestrales avec Fitch

```

 $F_p \leftarrow S_p \cap F_a$ 
si  $F_p \neq F_a$  alors
  si  $S_q \cap S_r \neq \emptyset$  alors
     $F_p \leftarrow ((S_q \cup S_r) \cap F_a) \cup S_p$ 
  sinon
     $F_p \leftarrow S_p \cup F_a$ 
  fin si
fin si

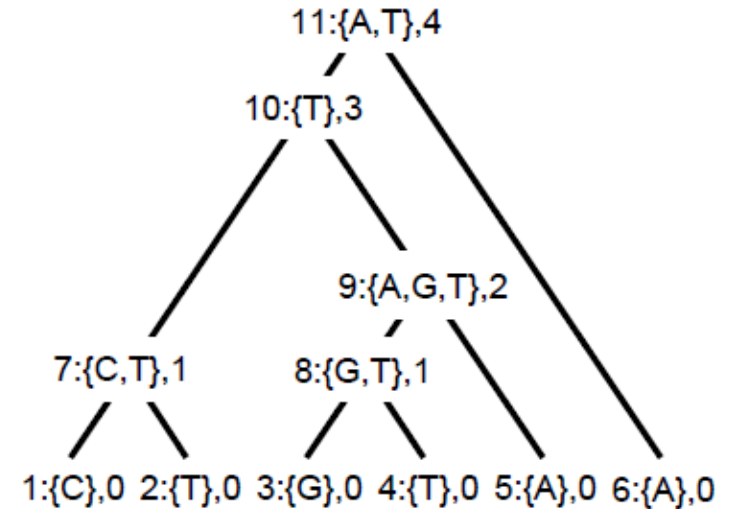
```

Si $F_p = F_a$, alors aucun changement évolutif ne s'est produit le long de la branche reliant a et p
 Si F_a présente des états de caractères qui ne sont pas présents en F_p , alors il peut être nécessaire d'ajouter des états possibles à F_p

Algorithme de Fitch: Reconstruction de séquences ancestrales

Algorithme 2 Séquences ancestrales avec Fitch

$F_p \leftarrow S_p \cap F_a$
si $F_p \neq F_a$ **alors**
 si $S_q \cap S_r \neq \emptyset$ **alors**
 $F_p \leftarrow ((S_q \cup S_r) \cap F_a) \cup S_p$
 sinon
 $F_p \leftarrow S_p \cup F_a$
 fin si
fin si



Exemple 1:

-nœud 11, $F_{11} = \{A, T\}$
 -nœud 10 (fils du nœud 11), $S_{10} = \{T\}$
 $\Rightarrow F_p \leftarrow S_p \cap F_a \Leftrightarrow F_{10} = \{T\}$
 S_{10} ne contient pas toutes les valeurs de F_{11} , donc $F_{11} \neq F_{10}$
 Or $S_7 = \{C, T\}$ et $S_9 = \{A, G, T\}$
 $\Rightarrow S_7 \cap S_9 \neq \emptyset$
 $\Rightarrow F_{10} \leftarrow (((S_7 \cup S_9) \cap F_{11}) \cup S_{10})$
 $\Rightarrow F_{10} \leftarrow (((\{C, T\} \cup \{A, G, T\}) \cap \{A, T\}) \cup \{T\})$
 $\Rightarrow F_{10} = \{A, T\}$

Exemple 2:

-nœud 10, $F_{10} = \{A, T\}$
 -nœud 9 (fils du nœud 10), $S_9 = \{A, G, T\}$
 $\Rightarrow F_p \leftarrow S_p \cap F_a \Leftrightarrow F_9 = \{A, T\}$
 S_9 contient toutes les valeurs de F_{10} ,
 donc $F_{10} = F_9$

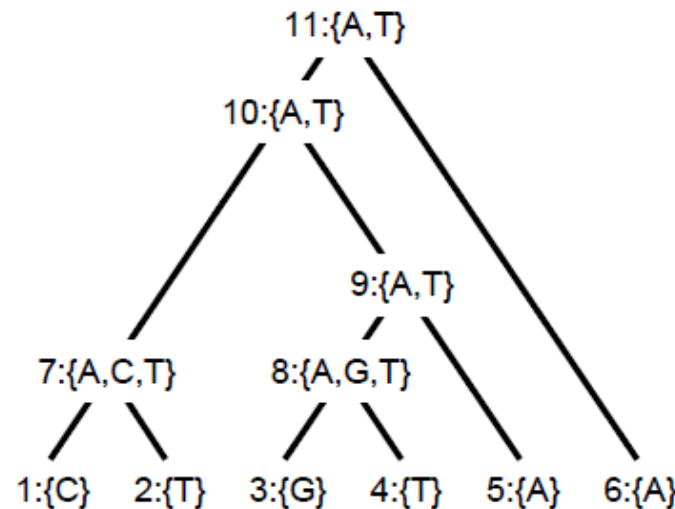
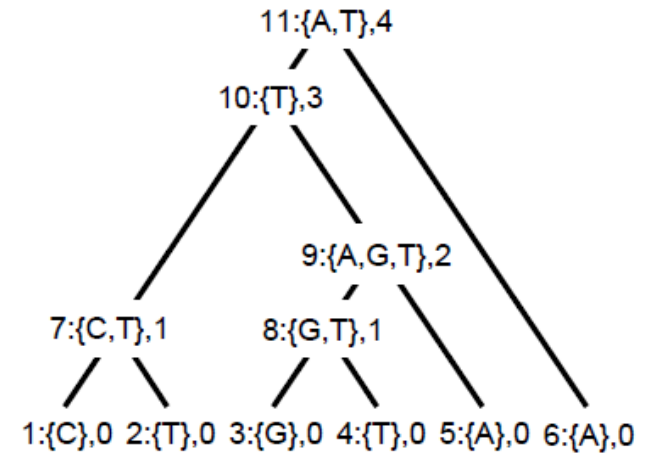
Algorithme de Fitch: Reconstruction de séquences ancestrales

Algorithme 2 Séquences ancestrales avec Fitch

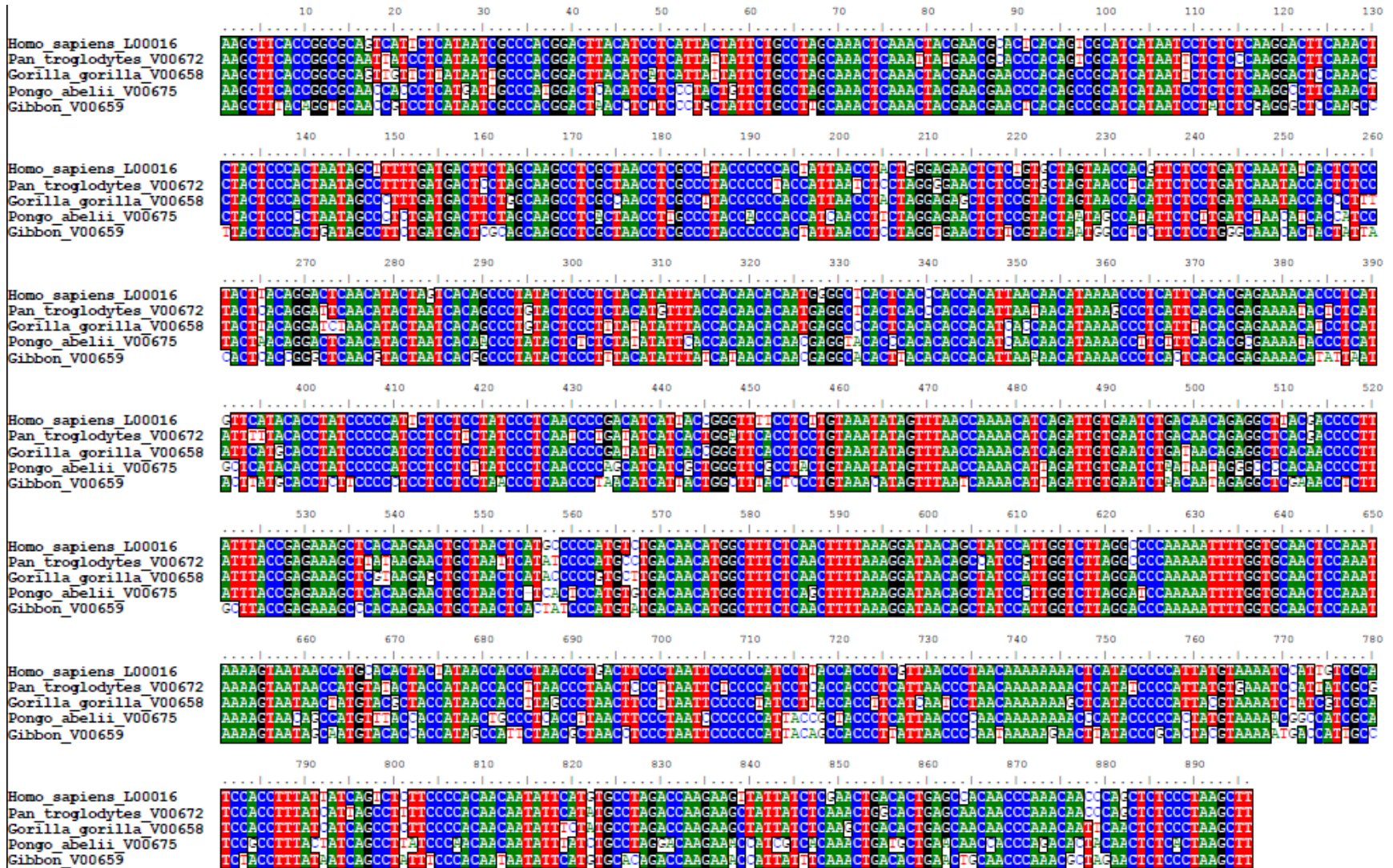
```

 $F_p \leftarrow S_p \cap F_a$ 
si  $F_p \neq F_a$  alors
  si  $S_q \cap S_r \neq \emptyset$  alors
     $F_p \leftarrow ((S_q \cup S_r) \cap F_a) \cup S_p$ 
  sinon
     $F_p \leftarrow S_p \cup F_a$ 
  fin si
fin si

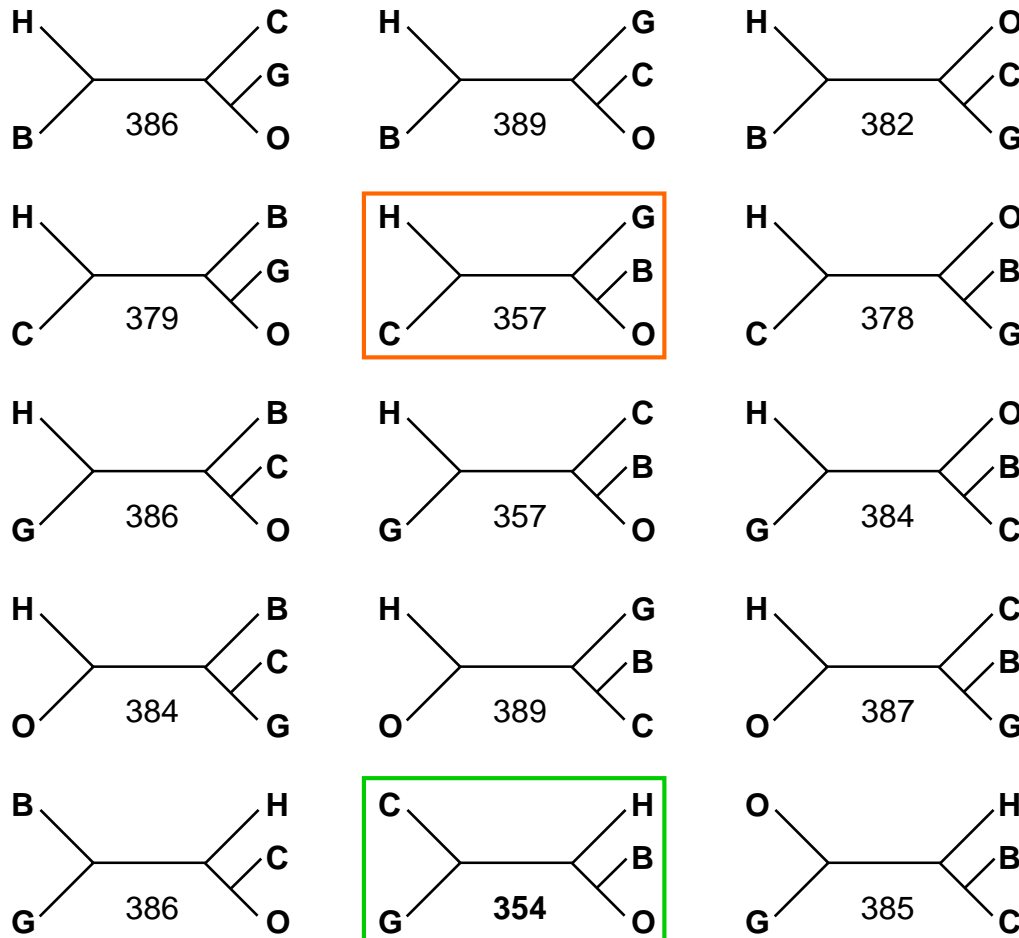
```



Application à la phylogénie des Hominoïdes



Application à la phylogénie des Hominoïdes

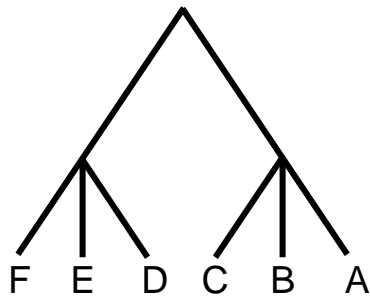
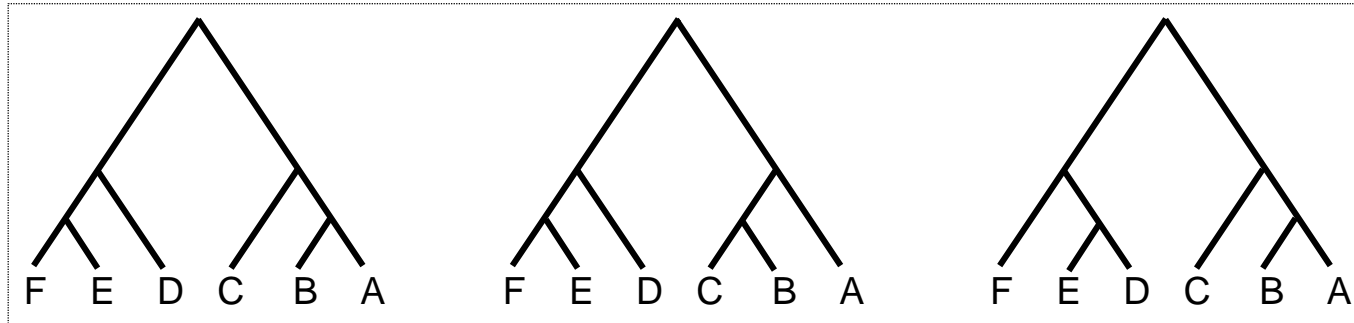


B= Gibbon, H = Homme, C = Chimpanzé, G = Gorille, O = Orang-outang

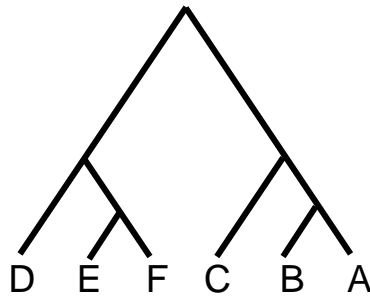
Parcimonie: Récapitulatif & propriétés

- Produit des arbres non racinés
- Le positionnement des changements dans un arbre n'est pas unique
 - ne permet pas d'inférer des longueurs de branches de manière unique
- Plusieurs arbres équiparcimonieux peuvent être trouvés
 - Inférence de consensus
- Le nombre d'arbre croissant de manière rapide avec le nombre de séquences, seule un sous-ensemble des topologies est testé pour identifier l'arbre le plus parcimonieux
 - Utilisation d'heuristiques pour explorer l'espace des arbres de manière rationnelle
 - Aucune certitude d'identifier l'arbre le plus parcimonieux à la fin de l'analyse
- Absence de critères pour discriminer le(les) arbre(s) le(s) plus parcimonieux des arbres légèrement moins parcimonieux
 - ex. est-ce qu'un arbre comptant 2504 pas est significativement meilleur que les 20 arbres comptant 2506 pas ?
- La parcimonie classique (algorithme de Fitch) considère toutes les substitutions comme équivalentes
 - Parcimonie pondérée (algorithme de Sankoff) permet de pondérer les types de changements

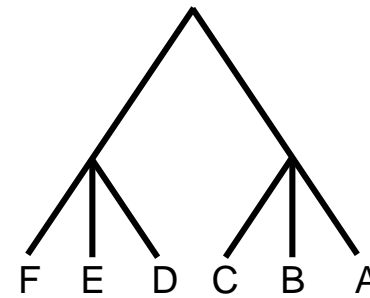
Consensus d'arbres



Strict



Maj. 50%



Maj. 80%

Explorer l'espace des topologies

- $n < 12$: Exploration exhaustive
 - $n < 20$: branch-and-bound
 - $n > 20$: heuristiques

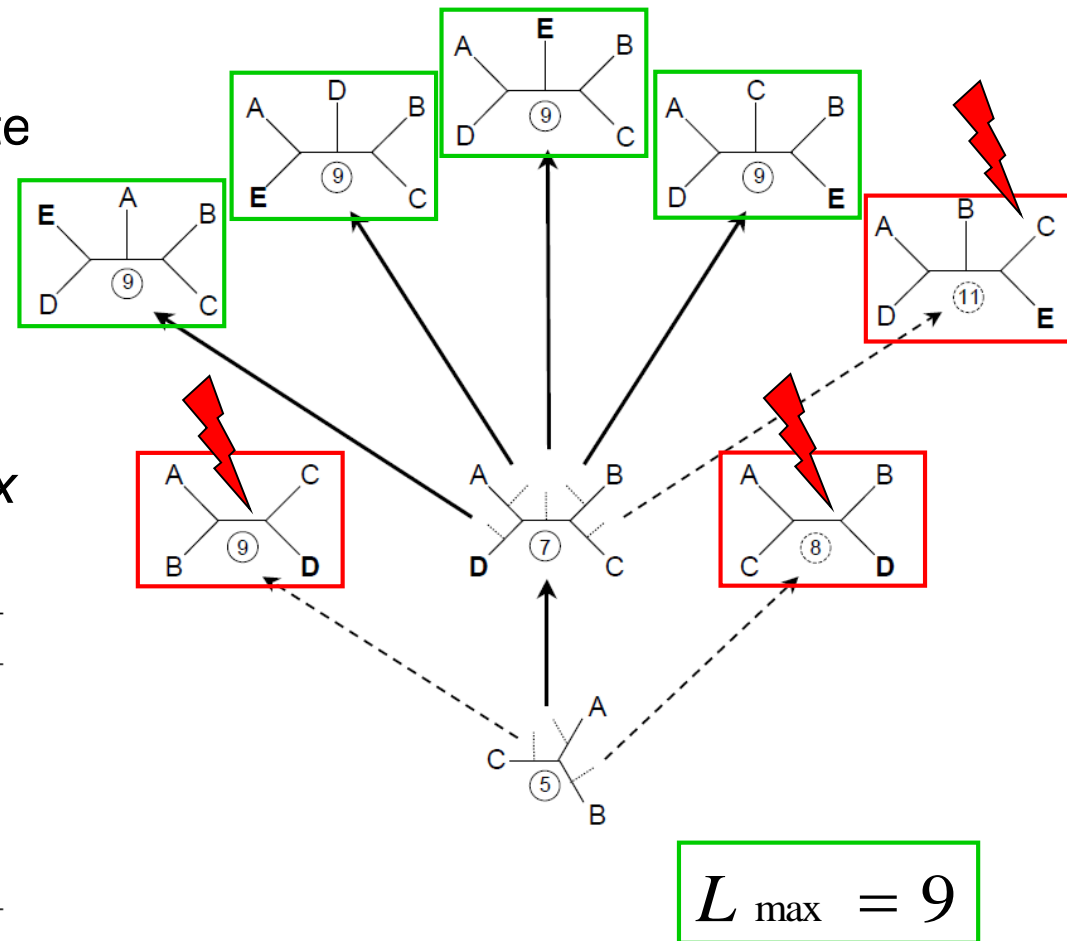
 - Utilisé pour la parcimonie, mais aussi les moindres carrés, le maximum de vraisemblance, etc.

 - Topologie de départ?
 - Topologie aléatoire
 - Meilleure topologie issue d'une recherche séquentielle
-

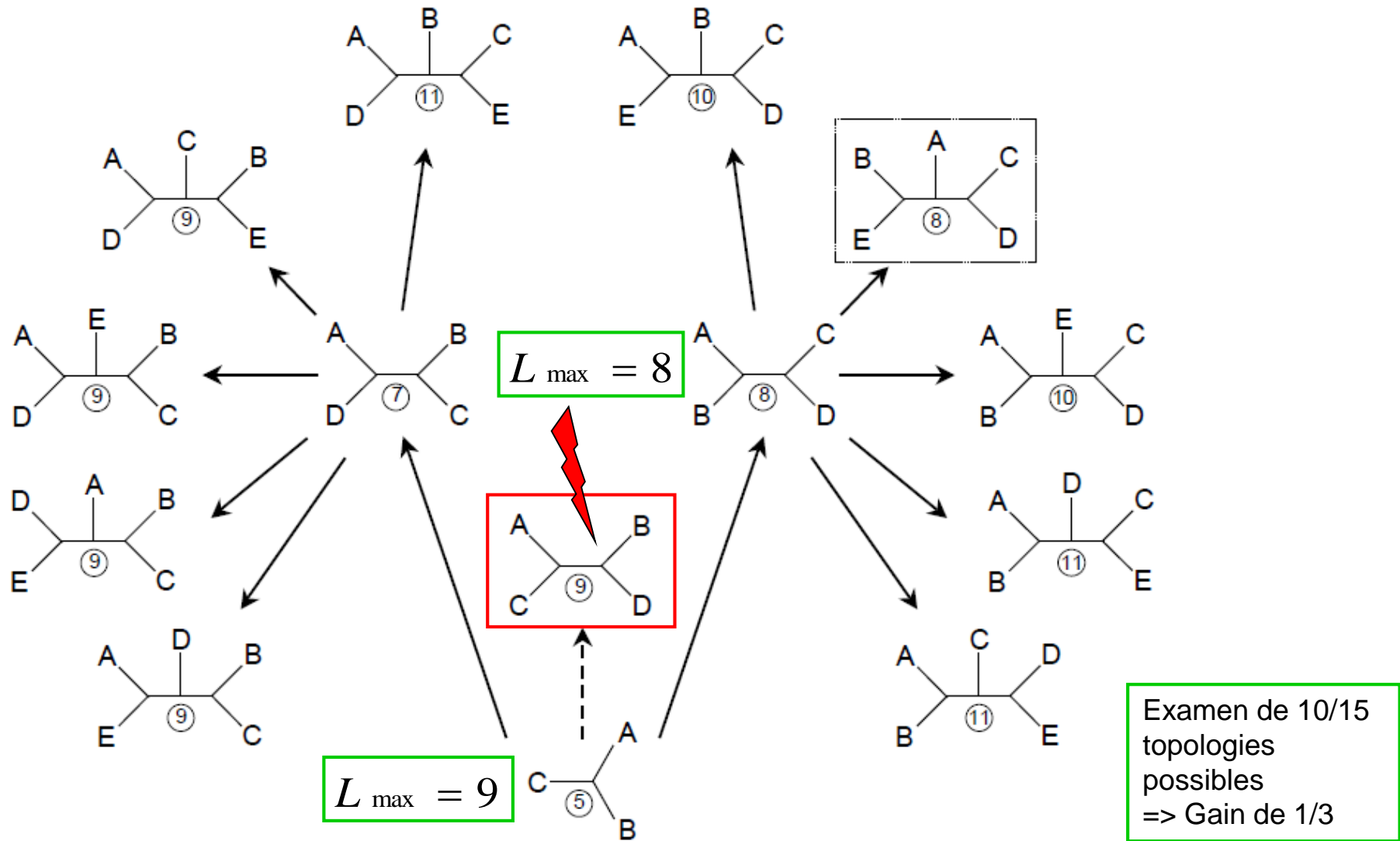
Recherche séquentielle

- Arbre à 3 feuilles
- Choix du 4^{ième} taxon à ajouter
 - ordre des taxa dans l'alignement
 - aléatoirement
 - maximum du minimum (taxon qui induit un L_{max} minimal)

UTO	1	2	3	4	5	6
A	A	T	T	A	A	T
B	T	T	A	T	T	T
C	A	A	T	T	T	T
D	A	A	T	A	A	A
E	T	T	A	A	A	T

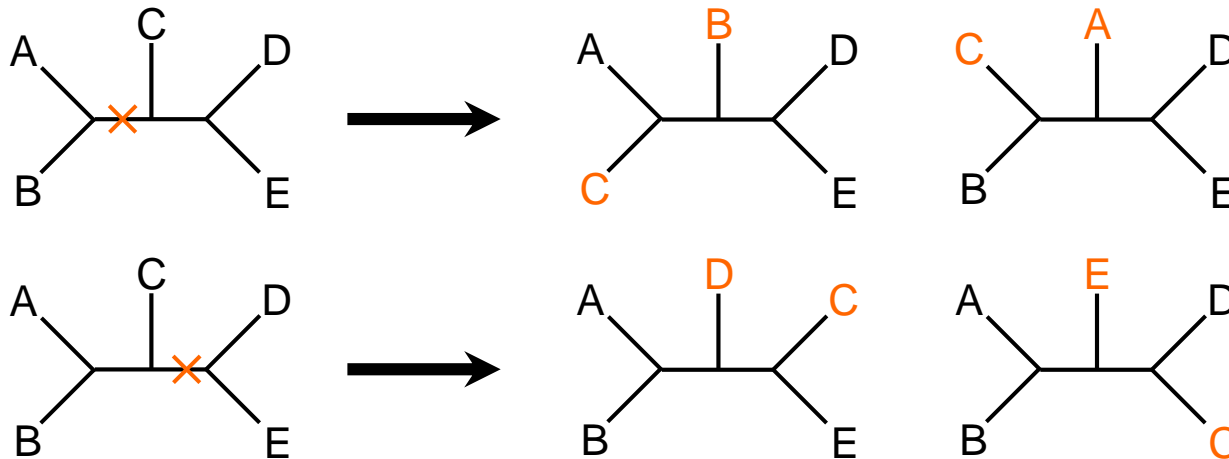


Branch-and-bound

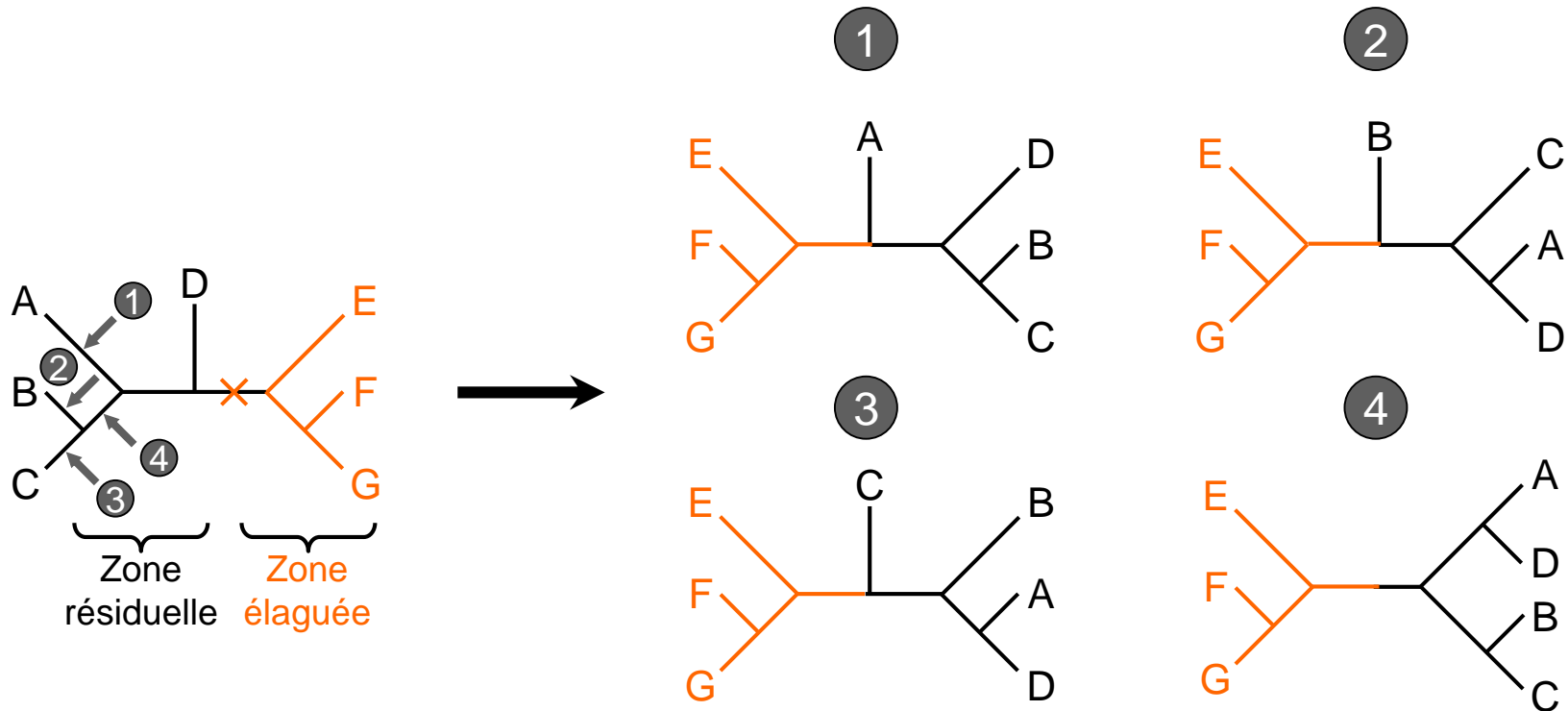


Nearest Neighbor Interchange (NNI)

- Examen des topologies se situant à une distance topologique $d_T = 2$ de l'arbre de départ
- $2(n - 3)$ arbres situés à une distance topologie $d_T = 2$

Complexité en $O(n)$ 

Subtree pruning and regrafting (SPR)



Si coupure au niveau d'une branche interne: $(2n - 8)$ arbres voisins

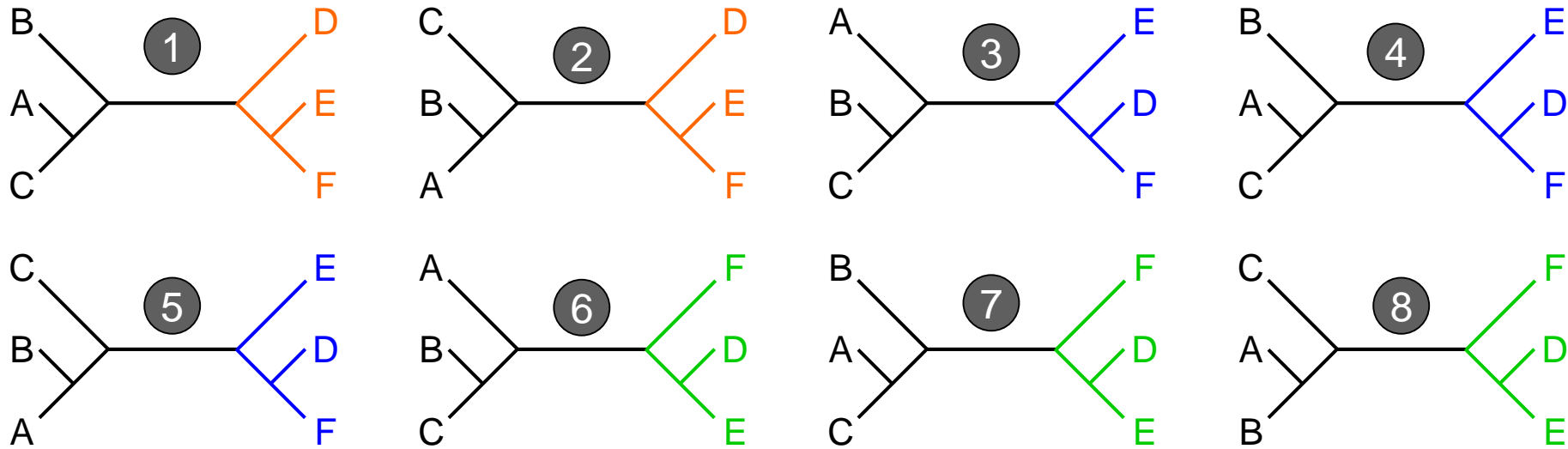
Si coupure au niveau d'une branche externe: $(2n - 6)$ arbres voisins

Un arbre non raciné compte: $(n - 3)$ branches internes et n branches externes

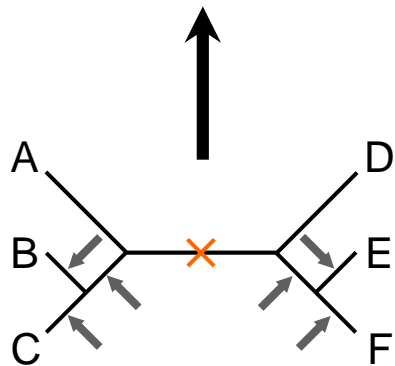
$$\Rightarrow \text{Nombre de voisins explorables: } \begin{aligned} & nx(2n - 6) + (n - 3)(2n - 8) \\ & = 4(n - 3)(n - 2) \end{aligned}$$

Complexité en $O(n^2)$

Tree Bisection and Reconnection (TBR)



$(2n - 3)(n - 3)^2$ Réarrangements maximum possibles



Complexité en $O(n^3)$

Lectures conseillées

- « Biologie évolutive » Thomas, Lefèvre, Raymond (2010) de boeck
 - « Evolution » Barton, Briggs, Eisen, Goldstein, Patel (2007) Cold Spring Harbor Laboratory Press
 - « Concepts et méthodes en phylogénie moléculaire » Perrière & Brochier-Armanet (2010) Springer collection IRIS
-