

Correction des TP1 et TP2 :

Jeux de données

Analyser le tableau 1

Nous pouvons observer que le nombre de systèmes homologues à ComED est très variable d'un génome à l'autre. Ces systèmes sont absents de *S. agalactiae* et *S. suis*, sont présents à un nombre variable d'exemplaires même dans le même groupe taxonomique. Cette versatilité suggère des événements de gains/perdes de gènes récents au cours de l'évolution. Si nous utilisons la proximité chromosomique pour reconstruire les systèmes, nous observons que le nombre de partenaires histidine kinase (HK) peut varier de un à trois (*S. uberis*), un et deux partenaire HK étant ce qu'il est trouvé de plus fréquent. Nous pouvons observer que l'annotation fonctionnelle des séquences donne peu d'information sur leur fonction biologique, de même les noms de gènes/protéines utilisés sont très peu fiables.

Alignement multiples des séquences homologues a ComE de *S. pneumoniae*

L'alignement obtenu est de très bonne qualité avec très peu d'insertions/délétions (indels) en dehors des régions Nter et Cter. Une variabilité au niveau de la partie Nter des protéines est souvent observée en raison de la difficulté de prédire correctement les débuts des gènes. Ces régions peuvent être éditées pour supprimer les indels mais cela aura peu de répercussions sur les reconstructions d'arbres. En effet, par défaut les méthodes basées sur une distance éliminent les colonnes comportant au moins une délétion et les méthodes basées sur le maximum de vraisemblance prennent en compte efficacement les délétions.

Construction des arbres en utilisant la méthode de distance BioNJ

Remarques : seule la longueur des branches horizontales est significative et chaque valeur de bootstrap est associée à une bi-partition de l'arbre.

La rotation des branches autour des nœuds (swap) ne change pas les bi-partitions et les longueurs de branches, l'arbre conserve sa topologie. Les arbres obtenus ne sont pas enraciné. Par défaut, le logiciel utilise la méthode du point médian. Comme nous disposons d'un groupe externe, nous allons l'utiliser pour enraciner tous nos arbres. Le nœud ancêtre sera sur la branche reliant ce groupe externe aux autres séquences. Cet enracinement permet d'orienter l'arbre (distinguer les nœuds pères des nœuds fils) et donc les différents événements qui se sont produits au cours de l'évolution.

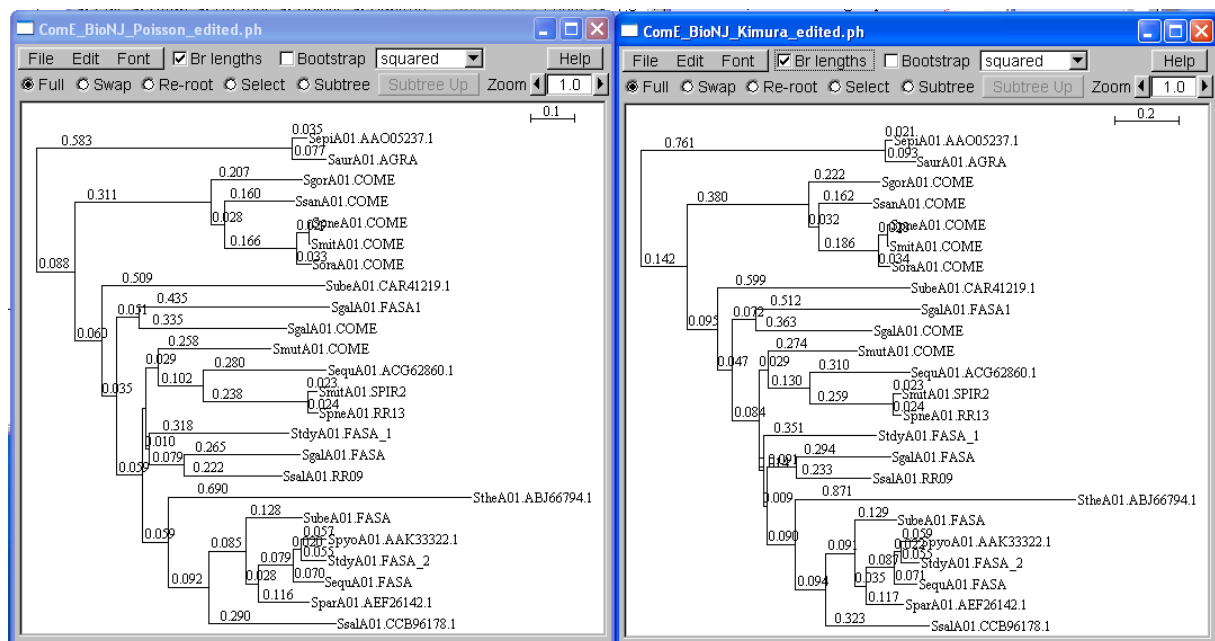
Comparaison topologie obtenue avec distance de Poisson et Kimura (approximation de la distance PAM)

Nous pouvons remarquer que les branches menant aux feuilles (branches externes) ont des longueurs similaires alors qu'elles ont tendance à être plus longues avec la méthode de distance Kimura pour les branches les plus profondes (noter la différence d'échelle entre les deux topologies).

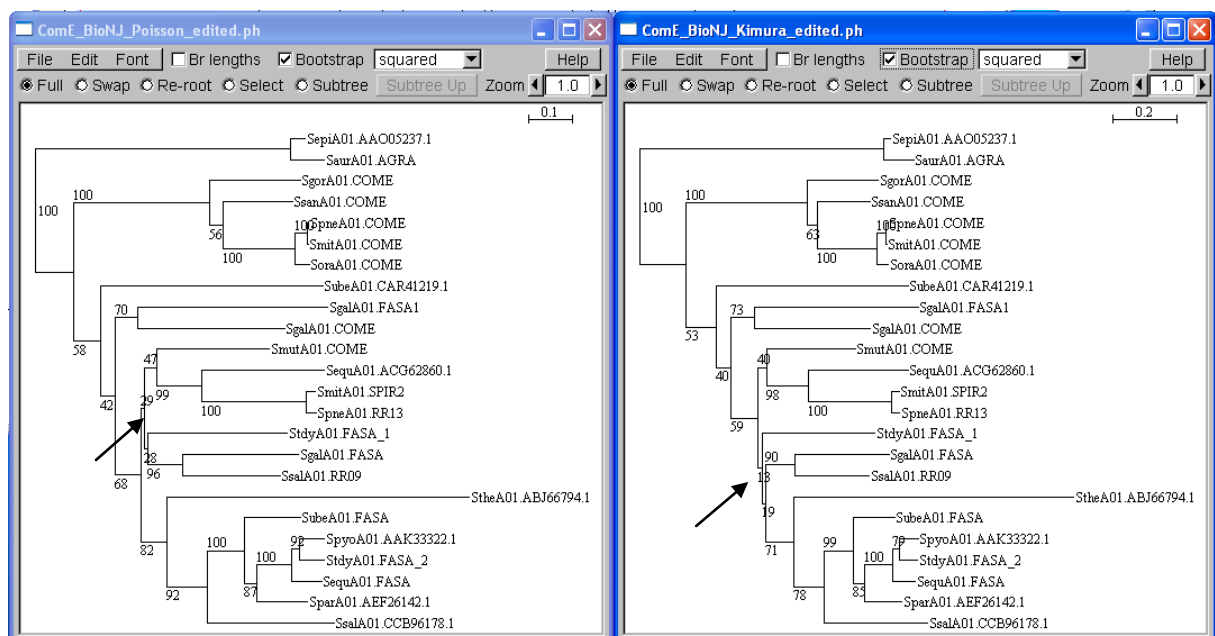
Les valeurs de bootstrap sont un peu meilleures avec Kimura.

Nous pouvons remarquer que les séquences semblent évoluer à peu près à la même vitesse (elles sont alignées verticalement) sauf la séquence StheA01.ABJ66794.1 qui montre clairement une accélération.

Il y a peu d'incongruences entre ces deux arbres, elles sont généralement associées à des branches courtes supportées par de faibles valeurs de bootstrap (indiquées par une flèche sur les arbres ci-dessous)



Distance de Poisson et Kimura, longueurs de branches

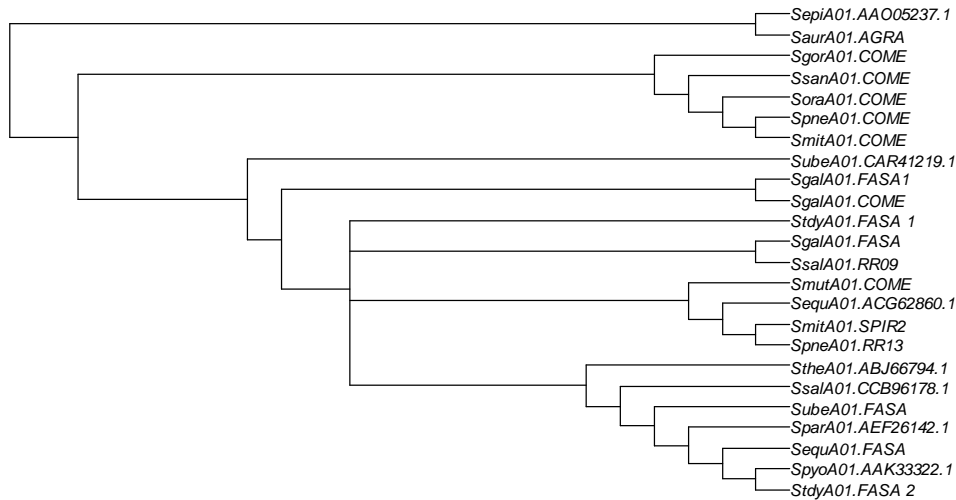


Distance de Poisson et Kimura, valeurs de Bootstraps

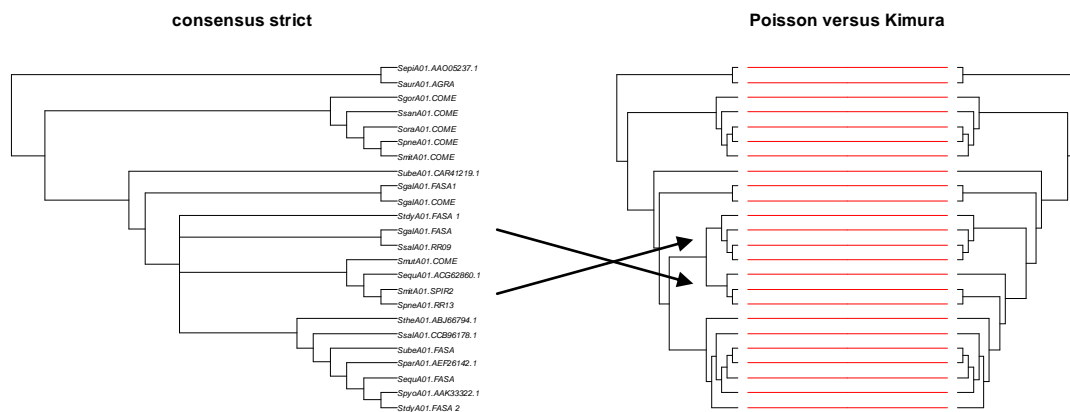
Comparaison des bipartitions des arbres obtenues avec les distances de Poisson et Kimura. Construction de l'arbre consensus.

Attention, les longueurs de branches associées à l'arbre consensus n'ont pas de signification phylogénétique. Nous observons une bonne résolution de l'arbre consensus ce qui traduit une très grande majorité de bipartitions communes entre les deux arbres. Il y a une seule région où les bipartitions ont été fusionnées (multifurcation ou polytomie, nœud dans un arbre qui connecte plus de trois branches).

Consensus strict between Poisson and Kimura



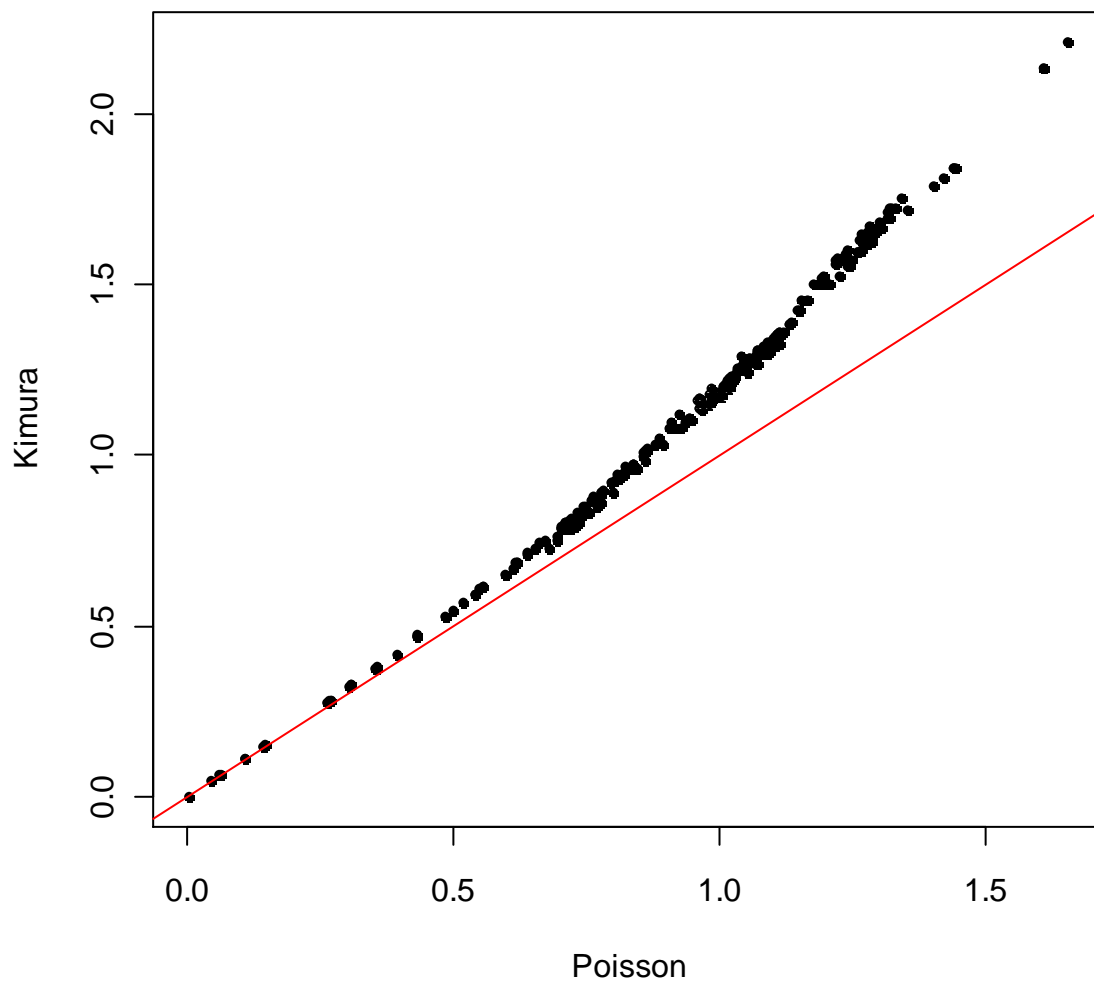
Congruence des arbres obtenues avec les distances de Poisson et Kimura



Attention, les deux arbres ne sont pas exactement alignés ! Comme précédemment, les longueurs de branches n'ont pas de signification phylogénétique.

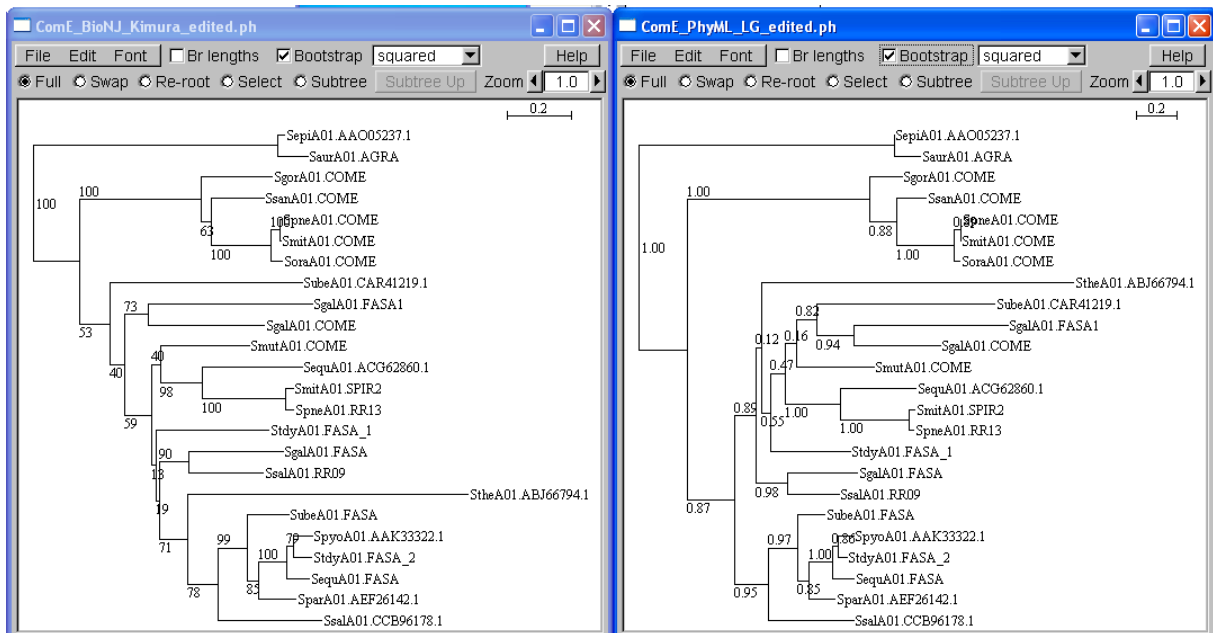
Relation entre les distances d'arbres Poisson/ Kimura

On observe une bonne corrélation entre les distances obtenues avec les deux modèles pour les petites distances. Par contre, il y a un décrochage très net pour les distances > 0.5 . Cela montre que la distance de Poisson sous estime les distances 'réelles' par rapport à la distance de Kimura quand la divergence augmente entre les paires de séquences.



Construction des arbres en utilisant une méthode du maximum de vraisemblance

Nous observons des changements topologiques importants entre ces deux arbres. Le plus important concerne la position de la séquence StheA01.ABJ66794.1 qui est radicalement différente. Nous pouvons également observer une perturbation générale qui conduit à un décalage des groupes de séquences par rapport à la verticale, ce qui traduit des vitesses relatives d'évolution différentes pour ces groupes. Les valeurs de bootstrap sont un peu meilleures pour la méthode PhyML LG mais restent faibles pour les régions incompatibles entre les deux arbres.



BioNJ distance de Kimura

PhyML : matrice LG

Informations données lors du déroulement du programme PhyML :

269 patterns found (out of a total of 272 sites) : 269

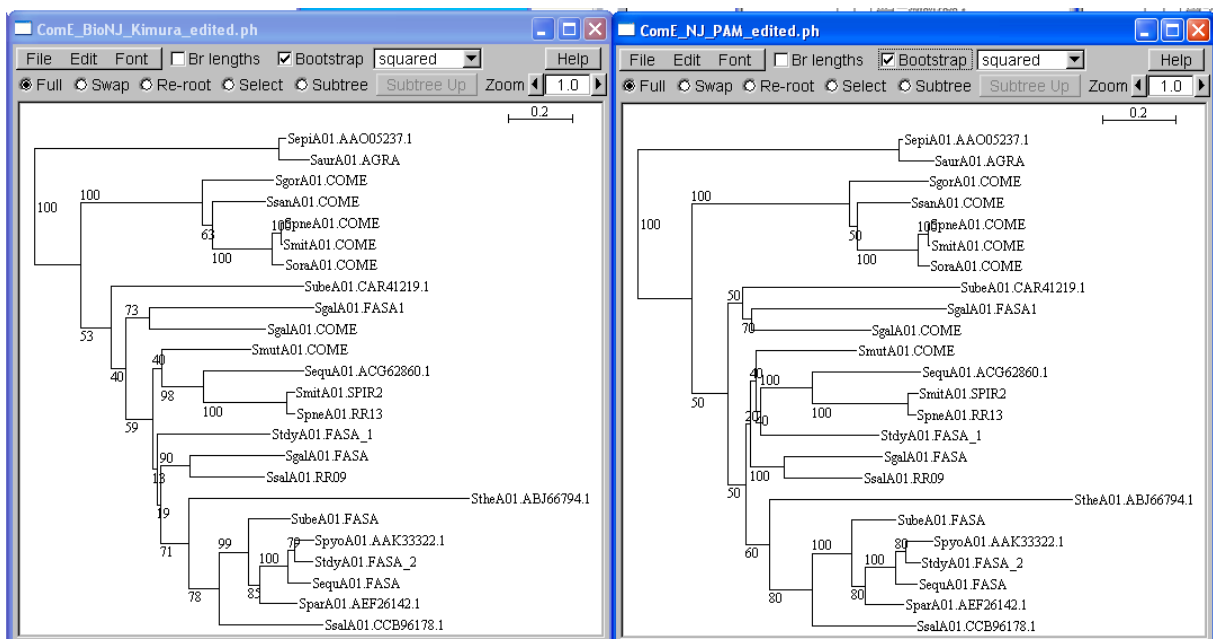
sites ont été utilisés pour le calcul de l'arbre

30 sites without polymorphism (11.03%) : 30

positions invariantes dans l'alignement

Log likelihood of PhyML LG 4 tree: -7848.270726

Construction des arbres en utilisant la méthode de distance NJ et une distance PAM

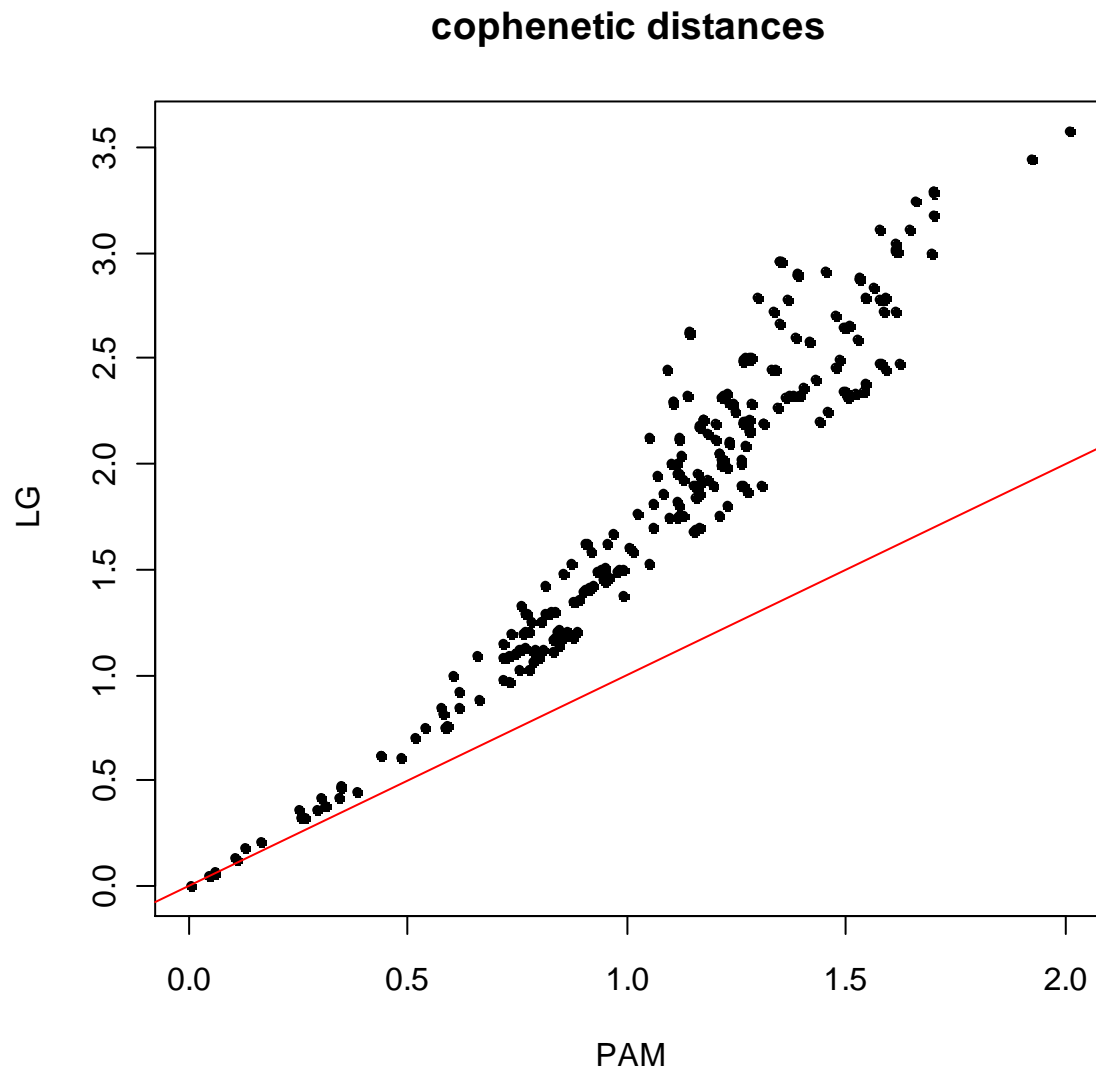


BioNJ distance Kimura

NJ distance PAM

L'arbre obtenu avec la distance PAM est très proche de celui obtenu avec les distances de Poisson et Kimura.

Relation entre les distances d'arbres PAM / LG

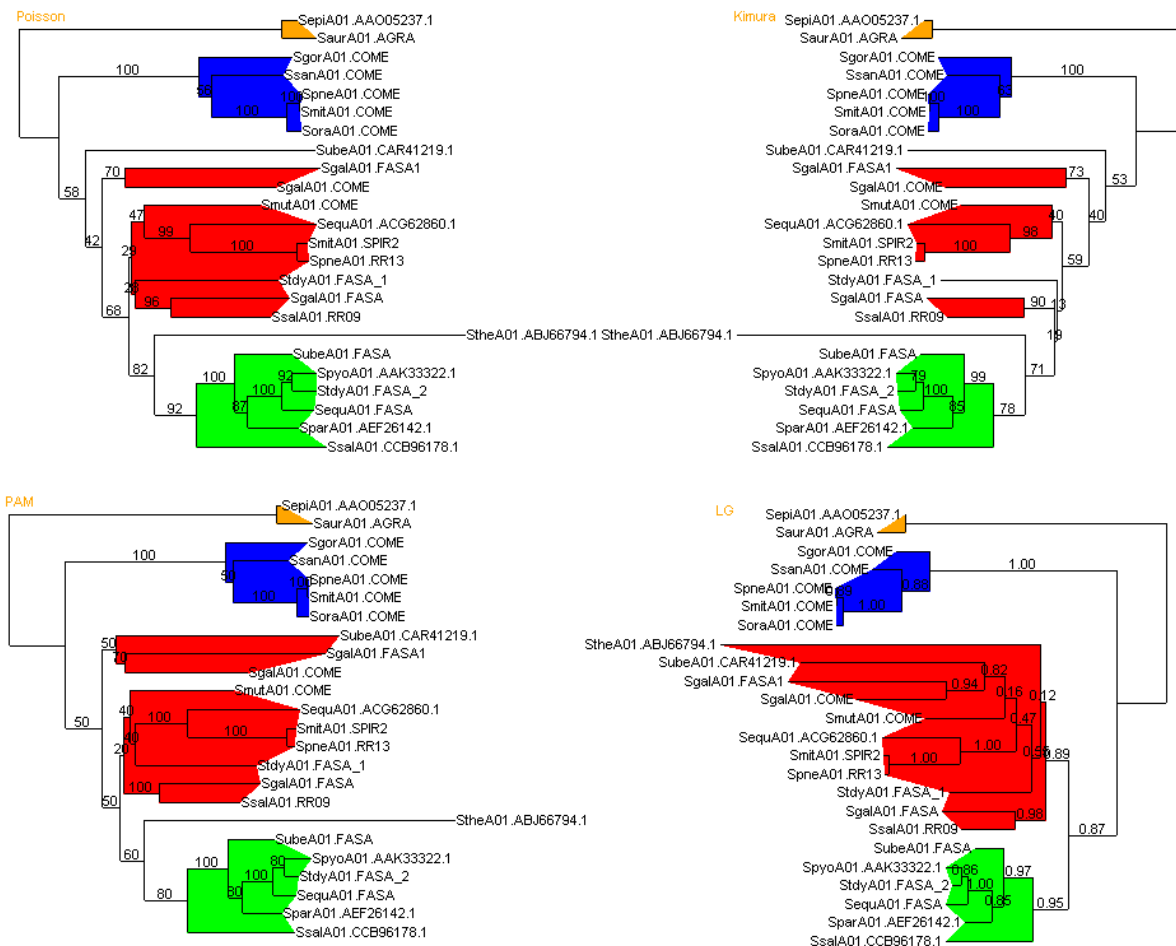


La corrélation est beaucoup moins bonne que celle observée précédemment. Nous observons une sous estimation des distances par la méthode NJ PAM en regard de la méthode PhyML LG. De plus, il y a une dispersion importante des points suggérant un traitement différent des substitutions observées entre les paires de séquences.

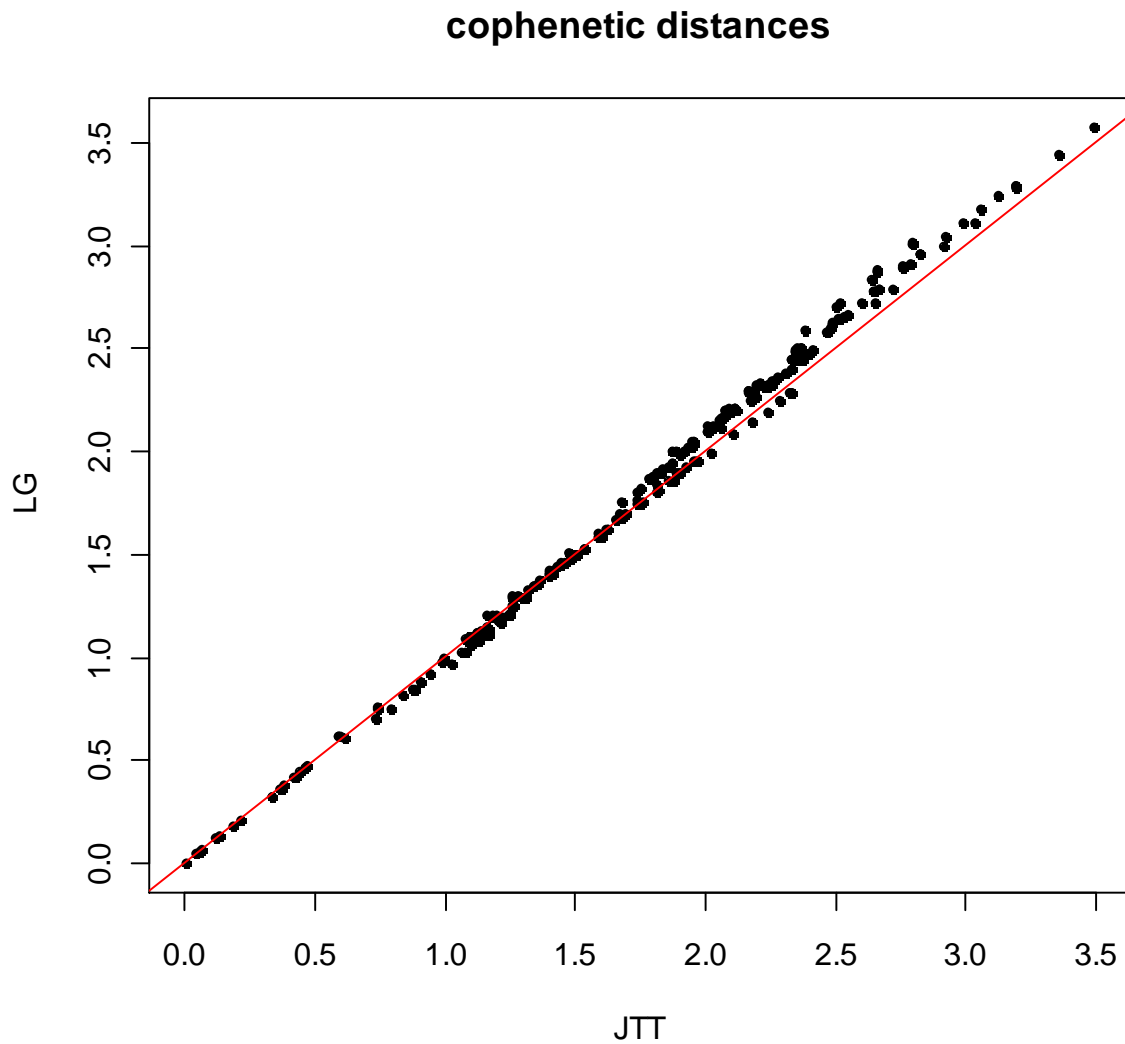
Edition et annotation des arbres

Nous pouvons faire les mêmes remarques que précédemment. Les quatre groupes de séquences apparaissent très clairement. Nous avons confirmation que les incongruences entre les arbres sont imputables aux feuilles du groupe rouge. Il est à noter que ce groupe est monophylétique avec la méthode PhyML matrice LG (support de aLTR de 0.89).

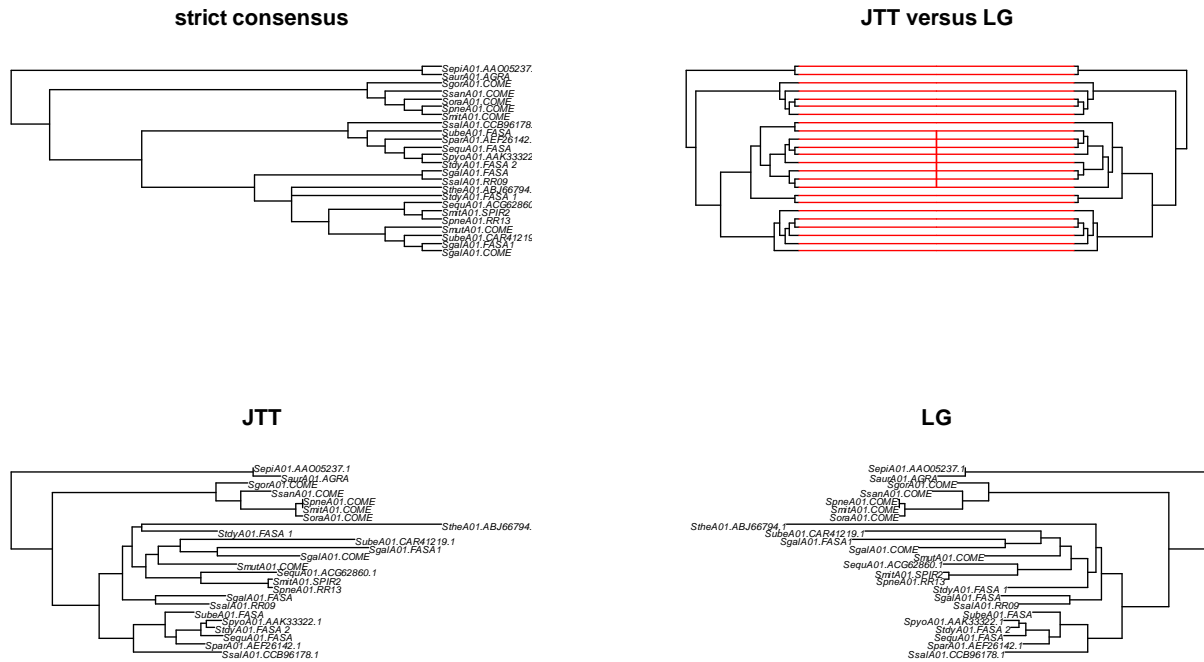
En rapport avec la question biologique que nous sommes posée, à savoir si la différence entre les temps de latence observée chez *S. pneumoniae* et *S. mutans* entre le moment où le CSP a été ajouté et celui où la transcription des gènes précoces est observée peut s'expliquer en analysant les protéines impliquées dans la régulation du processus, nous pouvons observer que la séquence de ComE de *S. mutans* appartient au groupe rouge, comme les séquences BlpR de *S. pneumoniae* (SpneA01.RR13) et non au groupe bleu renfermant la séquence de ComE de *S. pneumoniae*. Les gènes *comE* de *S. mutans* et *S. pneumoniae* sont donc paralogues et non pas orthologues, ce qui suggère des différences fonctionnelles.



Effet des paramètres de PhyML sur la reconstruction des arbres : matrice JTT versus matrice LG

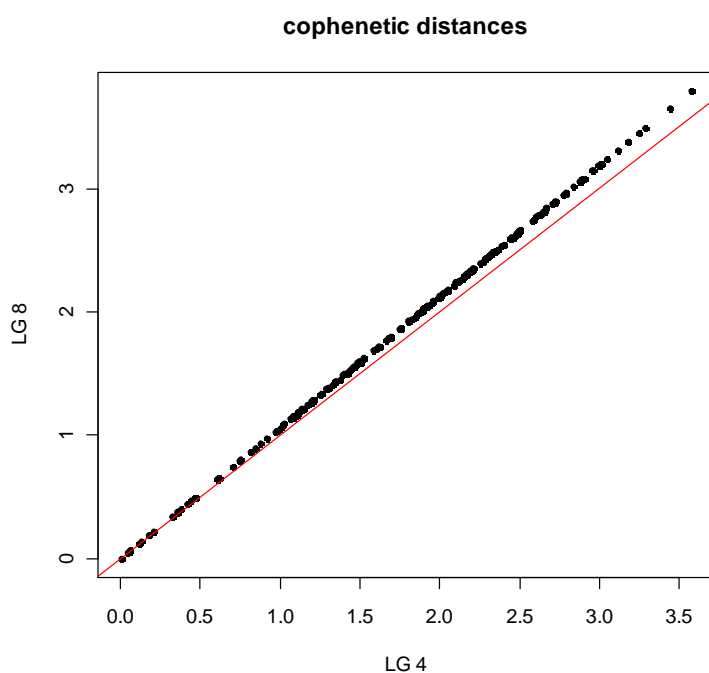


La corrélation est très bonne entre les résultats des deux modèles évolutifs. Les différences n'apparaissent que pour les distances les plus grandes.

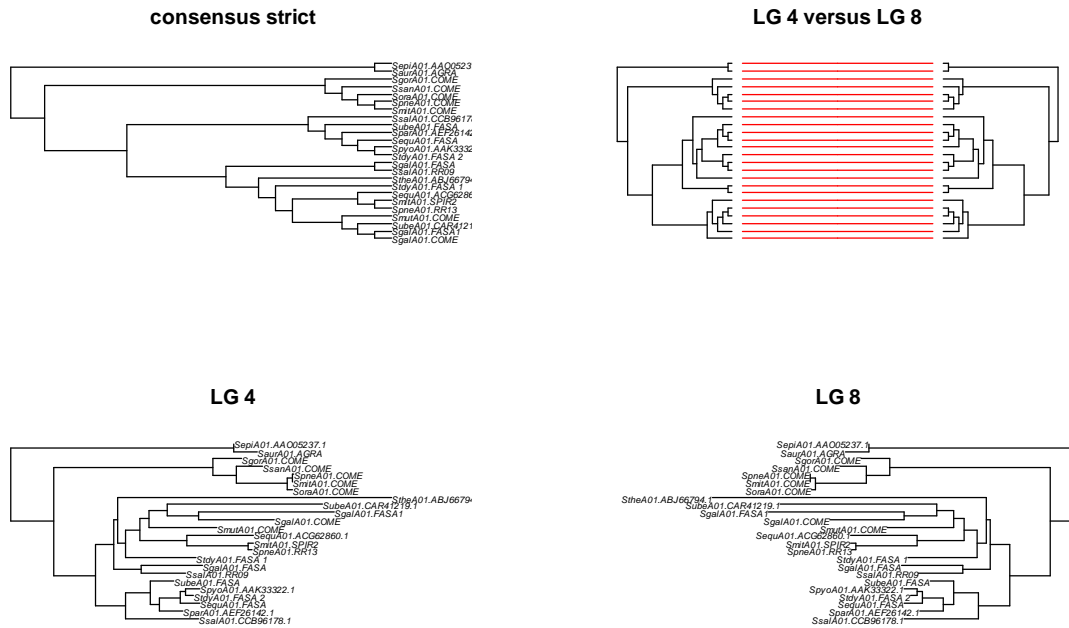


La seule différence topologique entre ces deux arbres est la position de la feuille StdyA01.FASA_1.
 Log likelihood of PhyML with LG 4 tree: -7848.270726
 Log likelihood of PhyML with JTT 4 tree: -7945.734999.
 L'arbre obtenu avec le model LG est le plus vraisemblable.

Vitesse différentes d'évolution des sites : nombre de classes de sites (paramètre de la gamma distribution (4 classes versus 8 classes))



La corrélation est quasiment parfaite suggérant que les distances ne diffèrent que par un facteur multiplicatif très faible.



Il n'est donc pas étonnant que les deux arbres présentent les mêmes topologies avec des valeurs de bootstrap très proches.

Log likelihood of PhyML LG 4 tree: -7848.270726

Log likelihood of PhyML LG 8 tree: -7845.940913

Le modèle avec 8 catégories est légèrement plus vraisemblable que le modèle à 4 catégories. Comme nous avons un petit jeu de données (272 sites), il n'est pas conseillé de choisir les modèles les plus compliqués car il y a plus de paramètres à estimer (ici 8 taux à la place de 4). Dans ce cas, ces paramètres pourront être mal estimés car les données ne sont pas assez nombreuses et donc conduire à des arbres moins bons qu'avec des modèles plus simples.

Recherche du modèle évolutif le plus adapté à l'alignement ComE

```
*****
MODEL OPTIMIZATION
*****

ProtTest options
-----
Alignment file..... :
C:\Users\fichant\Documents\Phylogenie\Phylogenomic\cours-
TD\TD_phylo\2012\TP1\SpneA01.COME_CleanUp_muscle.phy
Tree..... : BioNJ
StrategyMode..... : Fixed BIONJ JTT
Candidate models..... :
Matrices..... : JTT LG WAG Blosum62
Distributions..... : +G Uniform
Observed frequencies... : true
*****
Model..... : LG
Number of parameters..... : 47 (0 + 47 branch length estimates)
-lnL..... = 8568.97
(seconds)
```

```

Model..... : LG+F
  Number of parameters..... : 66 (19 + 47 branch length
estimates)
    aminoacid frequencies..... = observed (see above)
  -lnL..... = 8543.93
    (seconds))

Model..... : LG+G
  Number of parameters..... : 48 (1 + 47 branch length estimates)
    gamma shape (4 rate categories).. = 1.0
  -lnL..... = 8272.77
    (seconds))

Model..... : LG+G+F
  Number of parameters..... : 67 (20 + 47 branch length
estimates)
    gamma shape (4 rate categories).. = 1.0
    aminoacid frequencies..... = observed (see above)
  -lnL..... = 8235.34

```

```

*****
                        AKAIKE INFORMATION CRITERION
*****

```

```

*****
Best model according to AIC: LG+G+F
Sample Size:      272.0
Confidence Interval: 100.0
*****

```

Model	deltaAIC	AIC	AICw	-lnL
LG+G+F	0.00	16604.68	1.00	-8235.34
LG+G	36.86	16641.54	0.00	-8272.77
WAG+G+F	84.02	16688.70	0.00	-8277.35
JTT+G+F	109.99	16714.68	0.00	-8290.34
WAG+G	189.53	16794.21	0.00	-8349.11
Blosum62+G+F	214.22	16818.90	0.00	-8342.45
JTT+G	248.35	16853.03	0.00	-8378.52
Blosum62+G	281.14	16885.83	0.00	-8394.91
WAG+F	574.31	17179.00	0.00	-8523.50
LG+F	615.17	17219.86	0.00	-8543.93
WAG	626.67	17231.35	0.00	-8568.68
LG	627.26	17231.94	0.00	-8568.97
Blosum62+F	662.16	17266.85	0.00	-8567.42
JTT+F	698.27	17302.96	0.00	-8585.48
Blosum62	698.71	17303.39	0.00	-8604.70
JTT	806.17	17410.85	0.00	-8658.43

Les modèles +G et +G+F avec la matrice LG sont les plus vraisemblables. Nous pouvons remarquer que pour chaque matrice les modèles +G+F devancent les modèles plus simples.

Alignement multiples des séquences homologues a ComD de *S. pneumoniae*

Recherche du modèle évolutif le plus adapté à l'alignement ComD

```
*****
MODEL OPTIMIZATION
*****

ProtTest options
-----
Alignment file..... :
D:\Enseignement\TD3\2012\Results\ComD\SpneA01.COMD_CleanUp_muscle.fst
Tree..... : BioNJ
StrategyMode..... : Fixed BIONJ JTT
Candidate models..... :
  Matrices..... : JTT LG WAG Blosum62
  Distributions..... : +G Uniform
  Observed frequencies... : true
*****
AKAIKE INFORMATION CRITERION
*****

*****
Best model according to AIC: LG+G+F
Sample Size: 475.0
Confidence Interval: 100.0
*****
Model deltaAIC AIC AICw -lnL
-----
LG+G+F 0.00 49641.79 1.00 -24741.89
JTT+G+F 148.20 49789.98 0.00 -24815.99
WAG+G+F 213.60 49855.39 0.00 -24848.70
LG+G 412.43 50054.22 0.00 -24967.11
Blosum62+G+F 466.35 50108.14 0.00 -24975.07
Blosum62+G 864.39 50506.17 0.00 -25193.09
WAG+G 901.37 50543.15 0.00 -25211.58
JTT+G 928.68 50570.47 0.00 -25225.24
WAG+F 1302.05 50943.83 0.00 -25393.92
Blosum62+F 1357.13 50998.92 0.00 -25421.46
JTT+F 1365.50 51007.29 0.00 -25425.64
LG+F 1392.76 51034.55 0.00 -25439.27
Blosum62 1627.57 51269.35 0.00 -25575.68
LG 1663.27 51305.06 0.00 -25593.53
WAG 1880.42 51522.20 0.00 -25702.10
JTT 2080.67 51722.46 0.00 -25802.23
-----
```

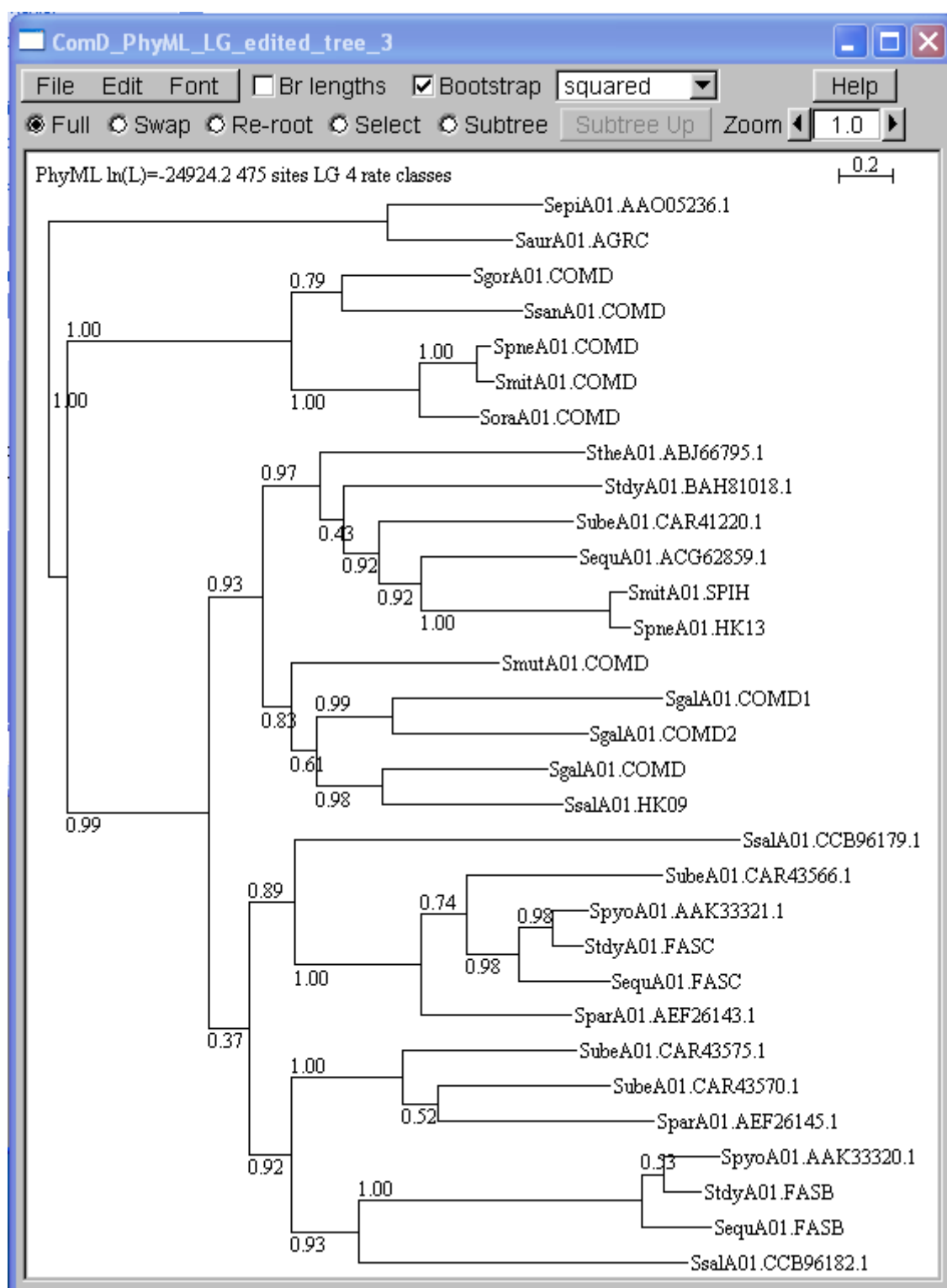
Le modèle LG+G+F est le plus vraisemblable comme dans le cas de ComE. Par contre, les modèles suivants ne changent que la matrice, ainsi pour chaque matrice les modèles +G+F devancent les modèles plus simples.

Construction de l'arbre des protéines ComD en utilisant le modèle le plus adapté et une méthode du maximum de vraisemblance

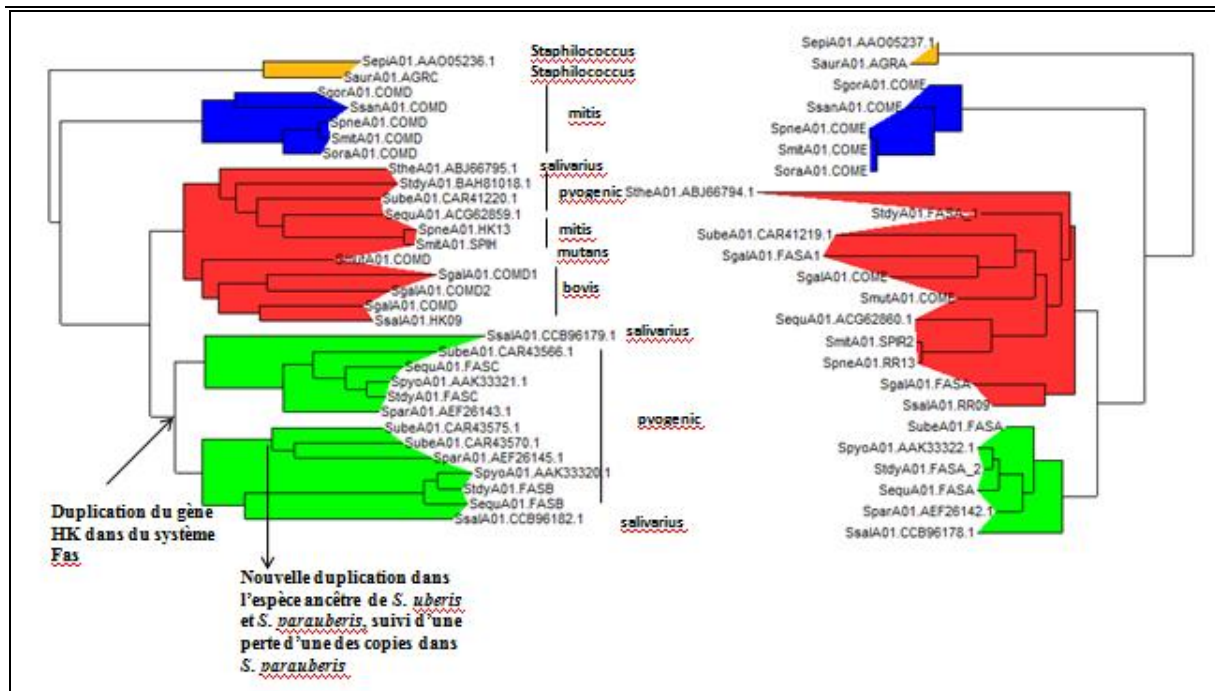
72 patterns found. (out of a total of 475 sites)

14 sites without polymorphism (2.95%).

Log likelihood of the current tree: -24924.205478.



Comparaison des arbres obtenus sur les séquences homologues à ComD et à ComE: interprétation biologique



Les séquences de *Staphylococcus* (Saur et Sepi) servent de groupe externe et permettent donc de connaître le nœud correspondant au nœud ancêtre hypothétique de l'ensemble de nos séquences de streptocoques. A partir de ce nœud, nous remarquons qu'une branche conduit à un nœud interne regroupant un certain nombre de séquences (coloriées en bleu sur chacun des deux arbres) parmi lesquelles ComD et ComE de *S. pneumoniae*. Parmi ces séquences il n'y a pas de paralogie (pas deux séquences appartenant à la même espèce) donc nous avons un groupe de séquences orthologues. Tous les génomes appartiennent au group mitis. Le système ComDE de *S. mutans* n'appartient pas à ce sous-arbre, il ne forme donc pas un système orthologue au système ComDE de *S. pneumoniae* et n'a donc probablement pas la même fonction. Ceci pourrait expliquer les différences de temps de latence observées lors de l'ajout du CSP avant le déclenchement de l'état de compétence chez *S. mutans*. Son système ComDE ne doit pas intervenir de la même manière dans la régulation de la compétence que le système ComDE de *S. pneumoniae*. Il est également possible qu'il ne soit pas impliqué dans cette régulation. C'est ce que nous savons aujourd'hui. La compétence étant régulée chez *S. mutans* par un autre système ComRS.

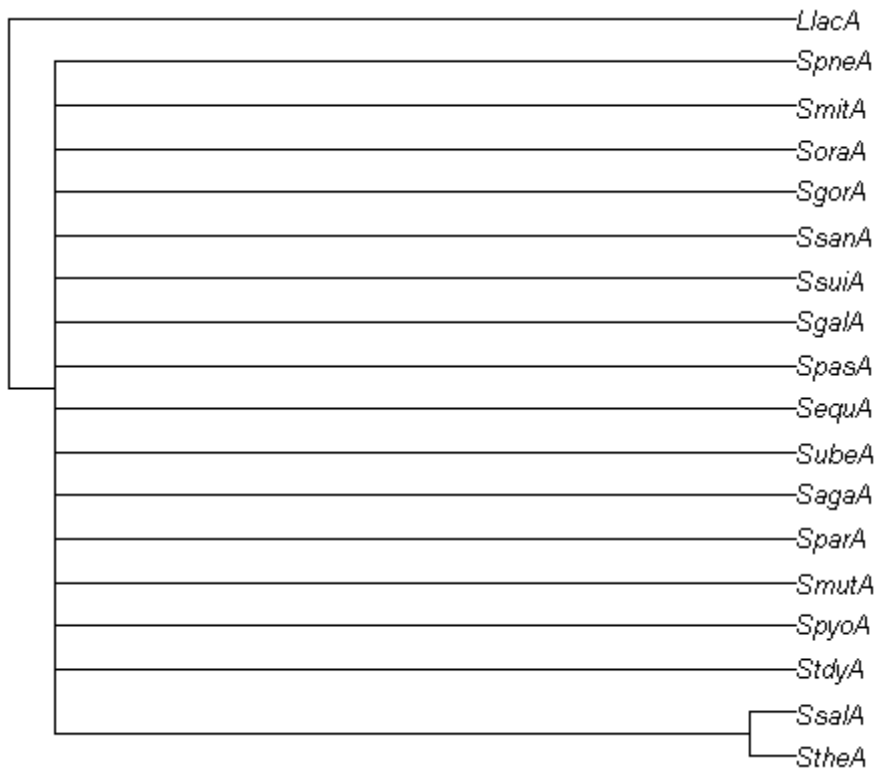
Pour les deux sous-arbres bleus, les topologies sont identiques. On dit que les deux sous-arbres sont **congruents**. En termes évolutifs, cela indique que les deux partenaires du système, ComD et ComE, ont **coévolué**.

La deuxième partie des arbres est plus complexe à analyser. Première remarque : des séquences homologues à BlpR et BlpH de *S. pneumoniae* (groupe rouge) présentent la plus grande distribution taxonomique avec 9 espèces représentées appartenant à 5 groupes taxonomiques (salivarius, mitis, pyogenic, mutans, bovis). Ceci suggère que les gènes codant pour ce système étaient présents dans l'ancêtre commun aux streptocoques et certaines espèces les auraient perdus. Le génome de *S. gallolyticus* se distingue par l'occurrence de trois copies du système (paralogues). On remarquera aussi que les séquences de *S. thermophilus* (StheA01) et *S. salivarius* (SsalaA01) (groupe salivarius) ne sont pas regroupées ce qui suggère des transferts horizontaux de gènes.

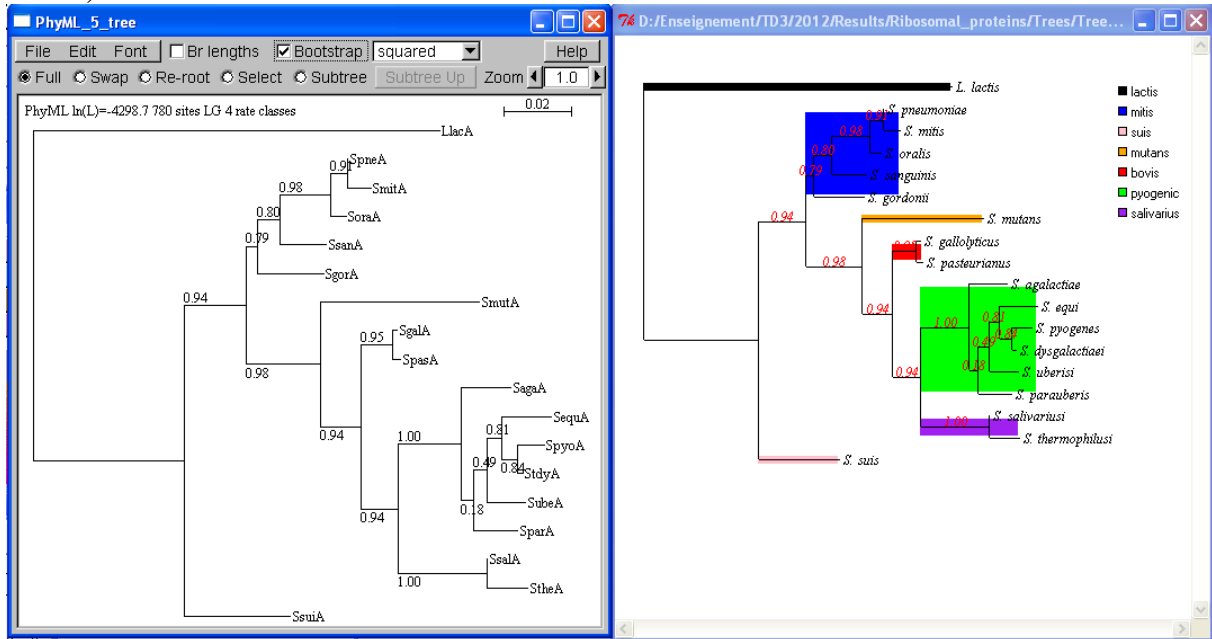
Le groupe vert est décomposé en deux sous arbres sur les ComD et un seul sous-arbre sur les ComE. Ce groupe renferme les séquences du système Fas qui posséderait deux HK par RR. Des signaux différents pourraient être "sentis" par chacun des senseurs et activer le même régulateur et donc activer les mêmes gènes. Ils sont trouvés majoritairement dans le groupe pyogenic.

Arbre consensus

Un des cinq arbres ne possède que 17 feuilles (tips) alors que les autres en possèdent 18. On ne peut calculer un arbre consensus que si les différents arbres ont exactement les mêmes feuilles, d'où le refus de la méthode quand nous demandons le consensus avec les 5 arbres. En supprimant l'arbre incriminé, nous obtenons l'arbre suivant qui confirme une totale incongruence entre les différents arbres (que des multifurcations) car aucune bifurcation commune, excepté pour *SsalA* et *StheA*.

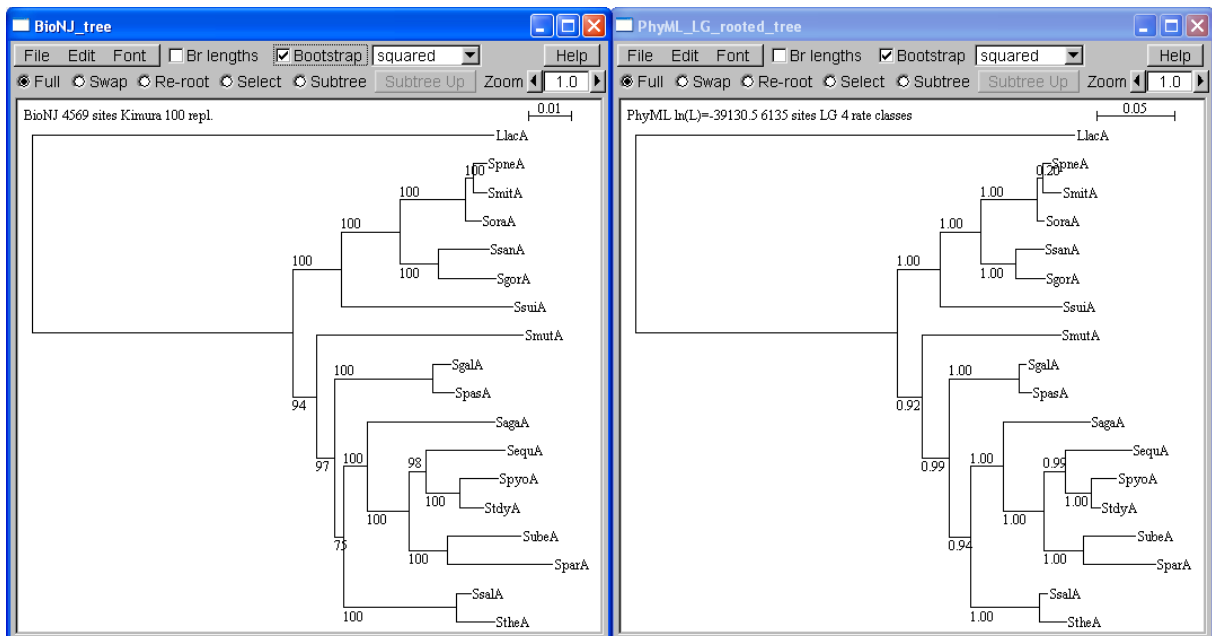


Arbre obtenu en concaténant les 5 alignement (méthode PhyML, matrice LG paramètres par défaut).



L'arbre obtenu en concaténant 5 fichiers est remarquablement cohérent avec la classification en groupe des Streptocoques. Nous pouvons également remarquer de bonne valeurs de aLRT.

Arbre obtenu sur l'alignement concaténé des 43 familles de protéines



Arbre obtenu avec BioNJ et la matrice Kimura à gauche et avec PhyML LG 4 classes à droite. Dans l'arbre de gauche, les sites comportant des indels ont été éliminés de l'analyse.

Résultat de PhyML :

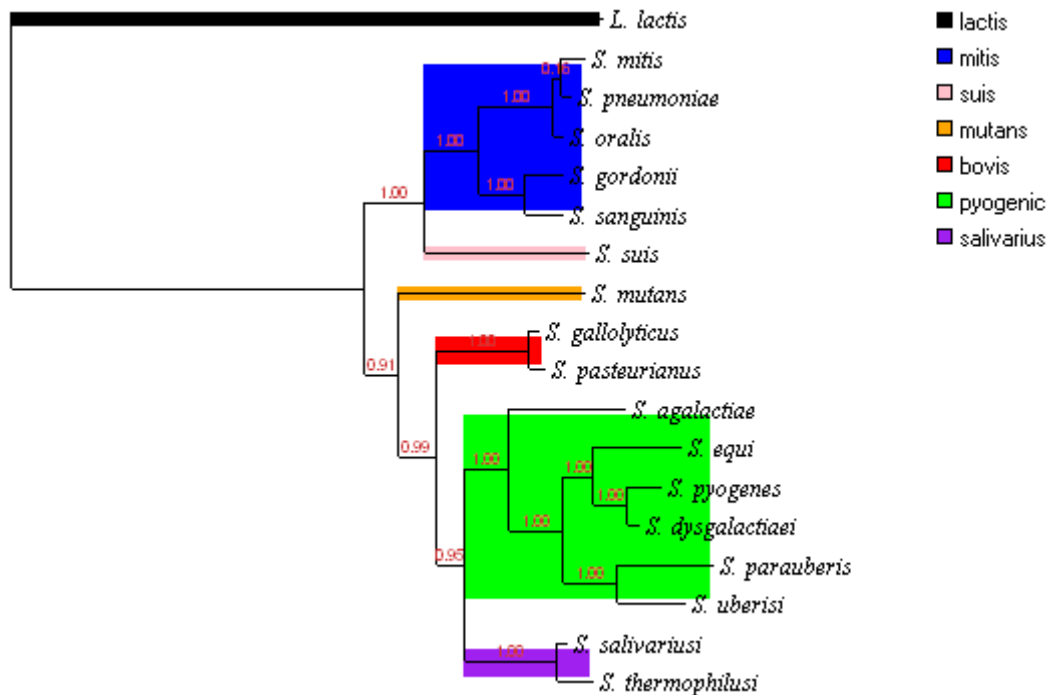
1529 patterns found. (out of a total of 6135 sites)

4304 sites without polymorphism (70.15%).

LG 4 classes : Log likelihood of the current tree: -38869.518584.

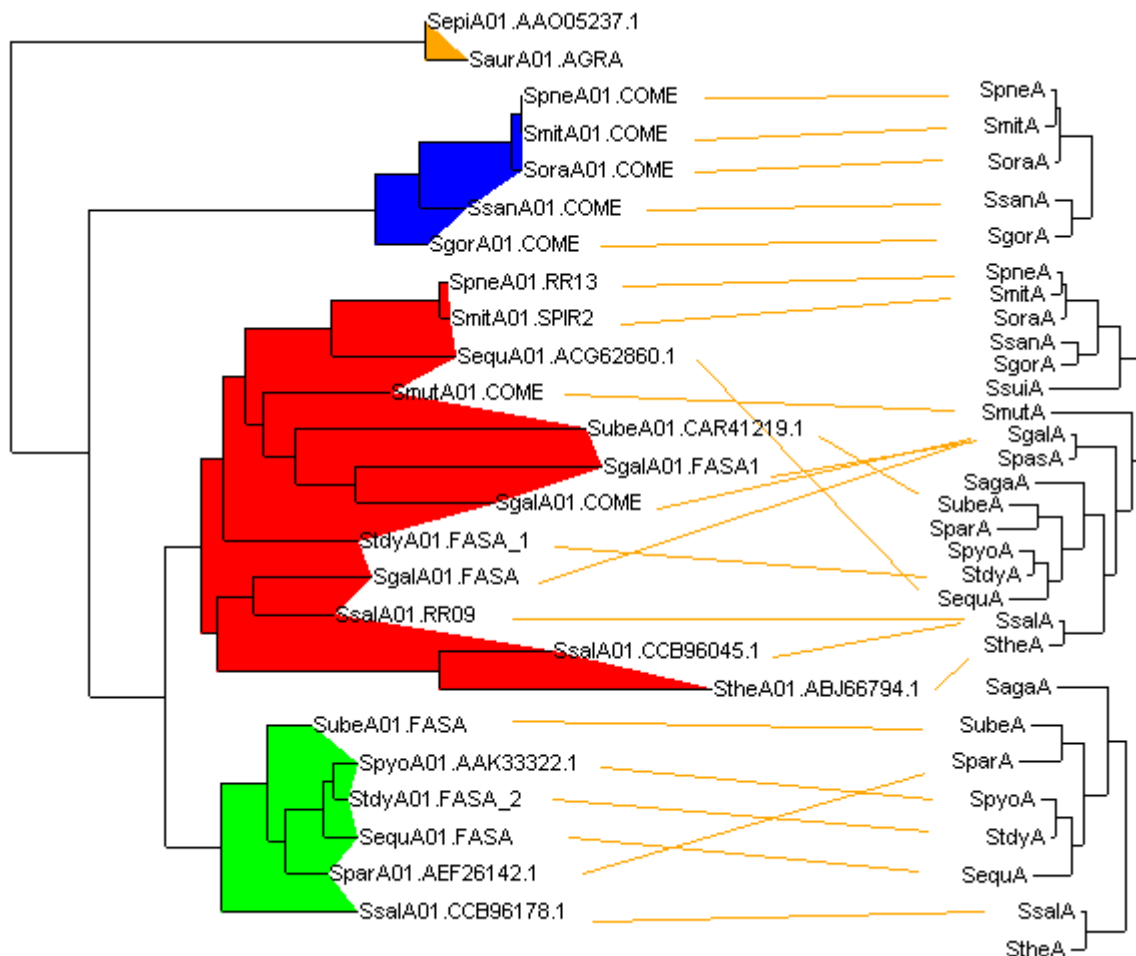
LG 8 classes : Log likelihood of the current tree: -38776.454864.

Quelle que soit la méthode utilisée, les arbres ont exactement la même topologie avec de très bonnes valeurs de bootstrap, ce qui suggère que l'alignement des 43 protéines ribosomiques contient suffisamment d'information phylogénétique pour résoudre les relations phylogénétiques entre ces espèces. Nous pouvons néanmoins observer que la branche menant à *S. pneumoniae*, et *S. mitis* à une très faible valeur de aLRT avec PhyML. Ce résultat pourrait être dû à la faible divergence des protéines ribosomiques qui n'apporteraient pas suffisamment d'information pour les espèces qui auraient divergées récemment.



Comparaison de la topologie de l'arbre obtenu sur ComE avec l'arbre des espèces

Quand nécessaire le sous-arbre des espèces a été extrait pour le mettre en face de la topologie des sous-groupes de ComE (bleu et vert). Il y a une très bonne congruence entre les deux pour le sous-arbre bleu des séquences ComE (et donc ComD). Nous voyons cependant pour les autres sous-groupes, quelques différences avec l'arbre des espèces. Pour le sous-arbre rouge, *S. equi* est à une place non attendue par rapport à la phylogénie des espèces. Sgal possède 3 paralogues dont deux sont correctement placés avec Sube et pourraient correspondre à une duplication dans Sgal. Par contre, le sous-groupe formé par la troisième copie de Sgal et SsalA1.RR9 pose un problème. Dans le sous-arbre vert, seule la séquence de Spar pose un petit problème car elle ne possède pas un ancêtre commun avec Sube mais cependant branche juste après Sube en groupe externe de Spyo, Stdy et Sequi.



Comme nous l'avons déjà remarqué, les séquences orthologues à ComE et ComD de *S. pneumoniae* ne sont présentes que dans le groupe mitis (groupe bleu). L'arbre obtenu avec les 43 protéines ribosomiques supposé représenter la phylogénie des espèces de Streptocoques montre que les différentes espèces du groupe mitis descendent bien d'une espèce ancêtre commune. Nous pouvons donc émettre l'hypothèse que le système ComDE a été acquis par cette espèce ancêtre et hérité ensuite par spéciation par les espèces actuelles. Cependant, que ce soit dans les arbres obtenus avec une méthode de distance (BioNJ ou NJ) ou avec une méthode de maximum de vraisemblance (PhyML), les sous-arbres correspondant aux systèmes ComDE branchent à l'extérieur du sous-arbre comportant les autres séquences. Ceci indique que ce système n'a probablement pas été acquis par duplication mais par transfert horizontal par l'ancêtre commun. En effet une duplication aurait du se traduire par la topologie d'arbre suivante (idem pour ComE) :

