

07/04/2011

Séminaire MaBioVis - LABRI

***Réseaux phylogénétiques
pour la visualisation de données
biologiques et textuelles***

Philippe Gambette



Plan

- Les réseaux phylogénétiques
- Méthodes de reconstruction
- Limites des méthodes combinatoires
- Illustration sur des données biologiques
- Utilisation sur des données textuelles
- Perspectives

Plan

- Les réseaux phylogénétiques
- Méthodes de reconstruction
- Limites des méthodes combinatoires
- Illustration sur des données biologiques
- Utilisation sur des données textuelles
- Perspectives

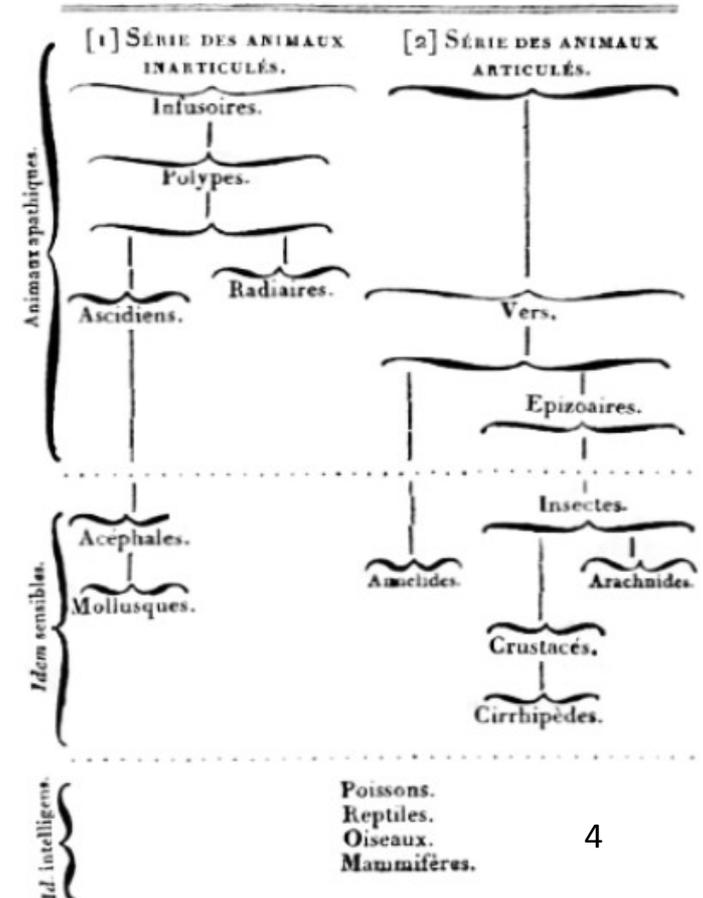
Les arbres phylogénétiques

Arbre phylogénétique d'un ensemble d'espèces :

- Les organiser en fonction de caractères communs
- Décrire leur évolution

classification

*ORDRE présumé de la formation des Animaux ,
offrant 2 séries séparées , subrameuses.*



*D'après Lamarck : Histoire naturelle des animaux
sans vertèbres (1815)*

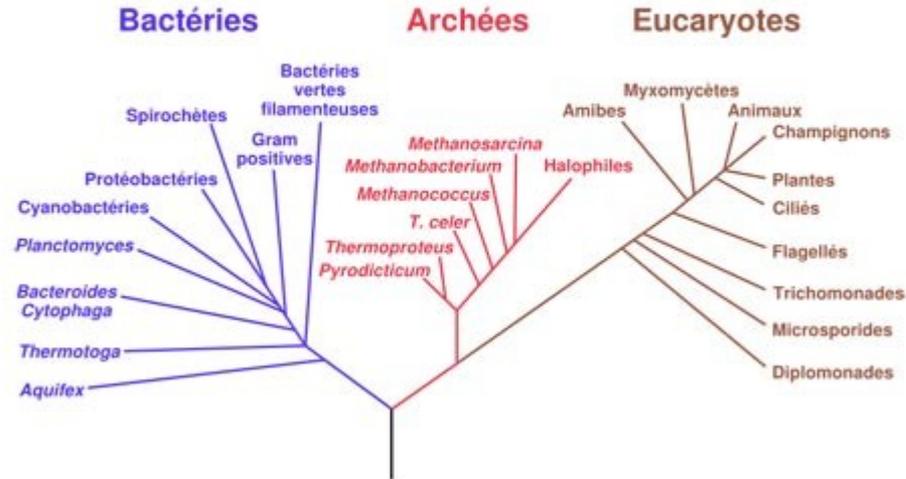
Les arbres phylogénétiques

Arbre phylogénétique d'un ensemble d'espèces :

- Les organiser en fonction de caractères communs
- Décrire leur **évolution**

modélisation

Arbre phylogénétique de la vie



D'après Woese, Kandler, Wheelis : Towards a natural system of organisms: proposal for the domains Archaea, Bacteria, and Eucarya, Proceedings of the National Academy of Sciences, 87(12), 4576–4579 (1990)

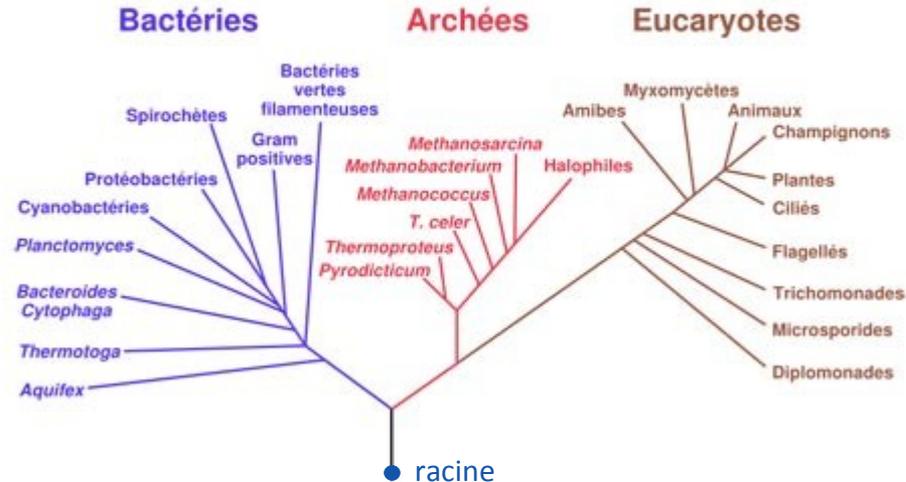
Les arbres phylogénétiques

Arbre phylogénétique d'un ensemble d'espèces :

- Les organiser en fonction de caractères communs
- Décrire leur **évolution**

modélisation

Arbre phylogénétique de la vie

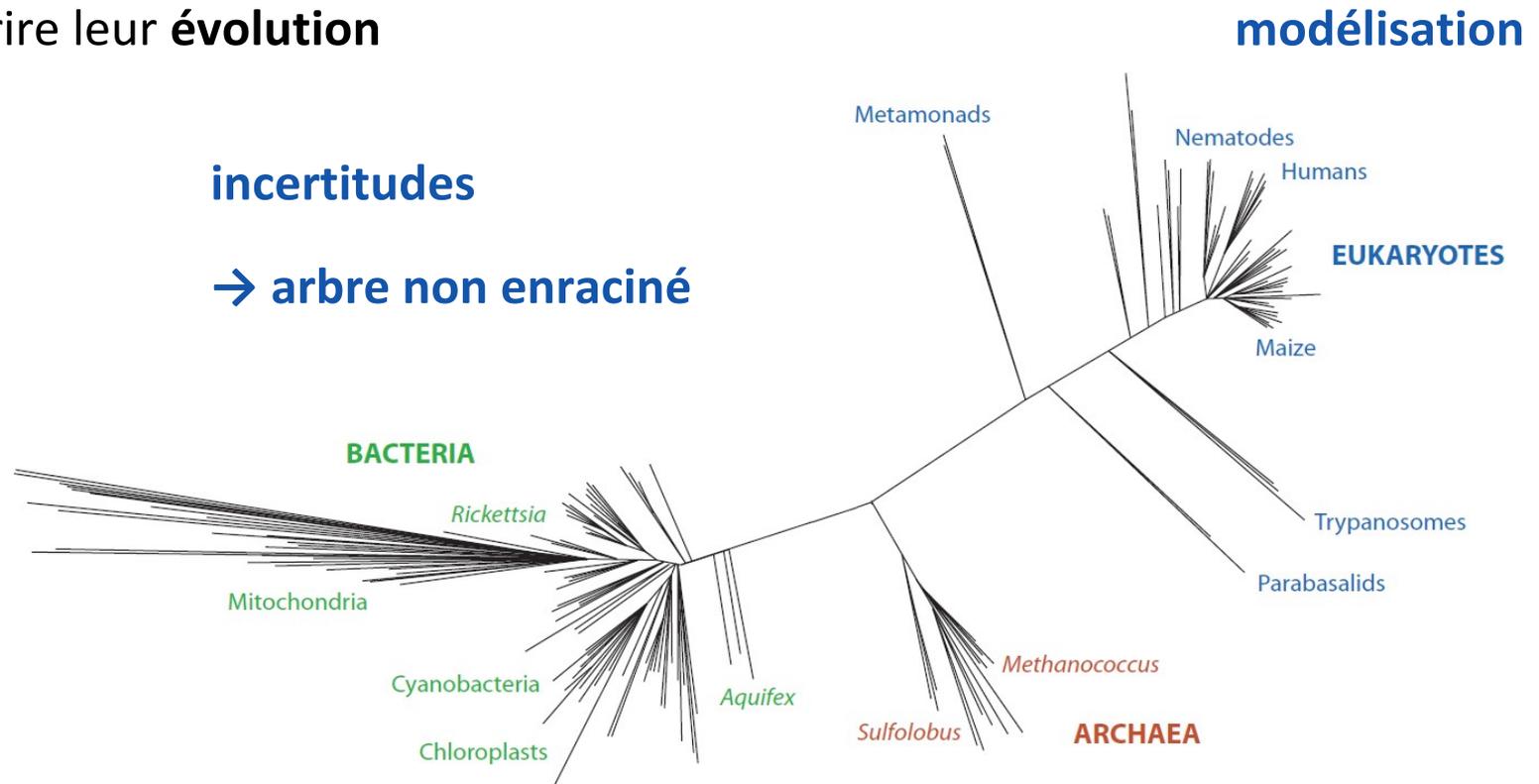


D'après Woese, Kandler, Wheelis : Towards a natural system of organisms: proposal for the domains Archaea, Bacteria, and Eucarya, Proceedings of the National Academy of Sciences, 87(12), 4576–4579 (1990)

Les arbres phylogénétiques

Arbre phylogénétique d'un ensemble d'espèces :

- Les organiser en fonction de caractères communs
- Décrire leur **évolution**

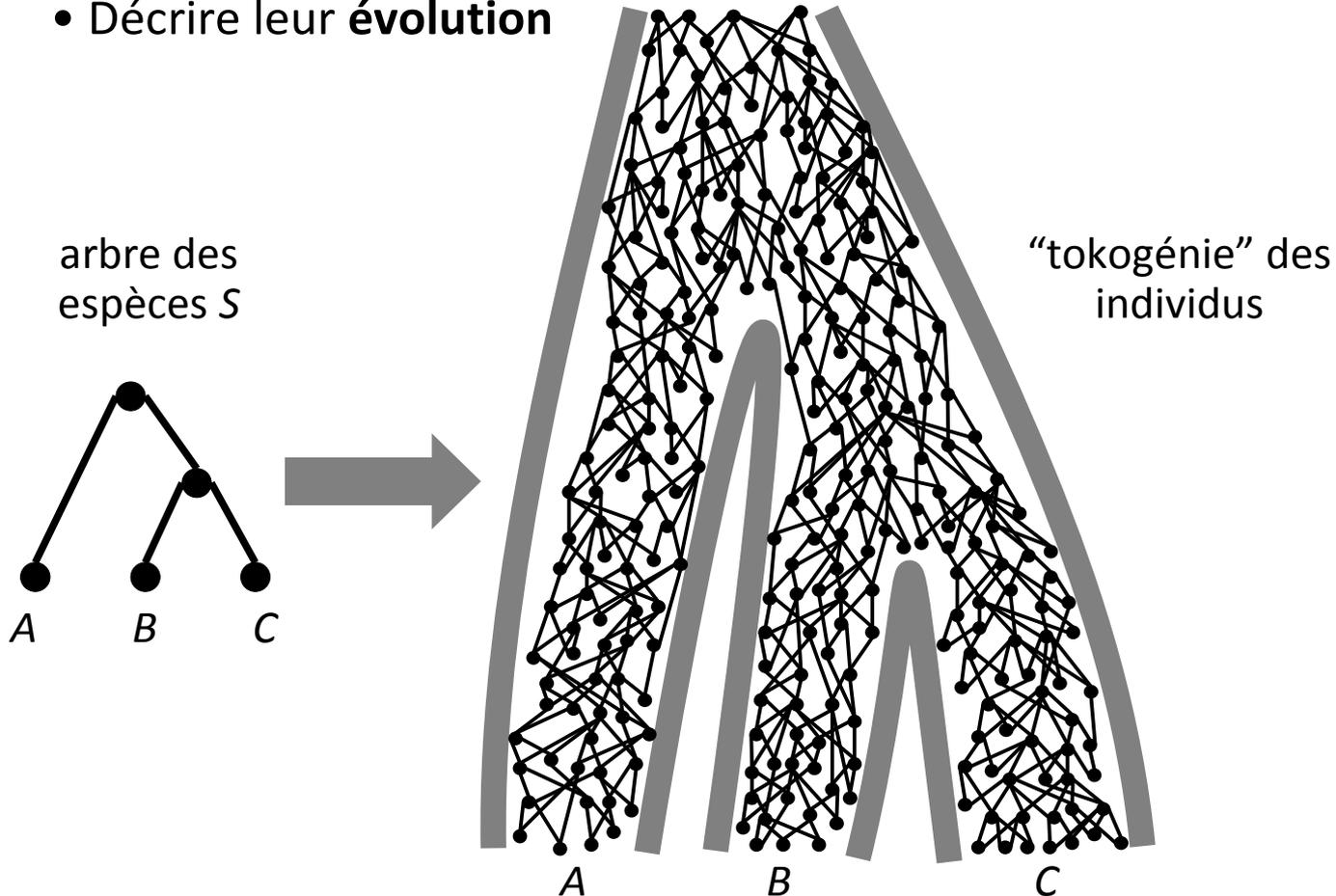


D'après Christophe Blumrich, David S. Spencer,
cité dans Doolittle : Uprooting the Tree of Life, Scientific American (Fév. 2000)

Les arbres phylogénétiques

Arbre phylogénétique d'un ensemble d'espèces :

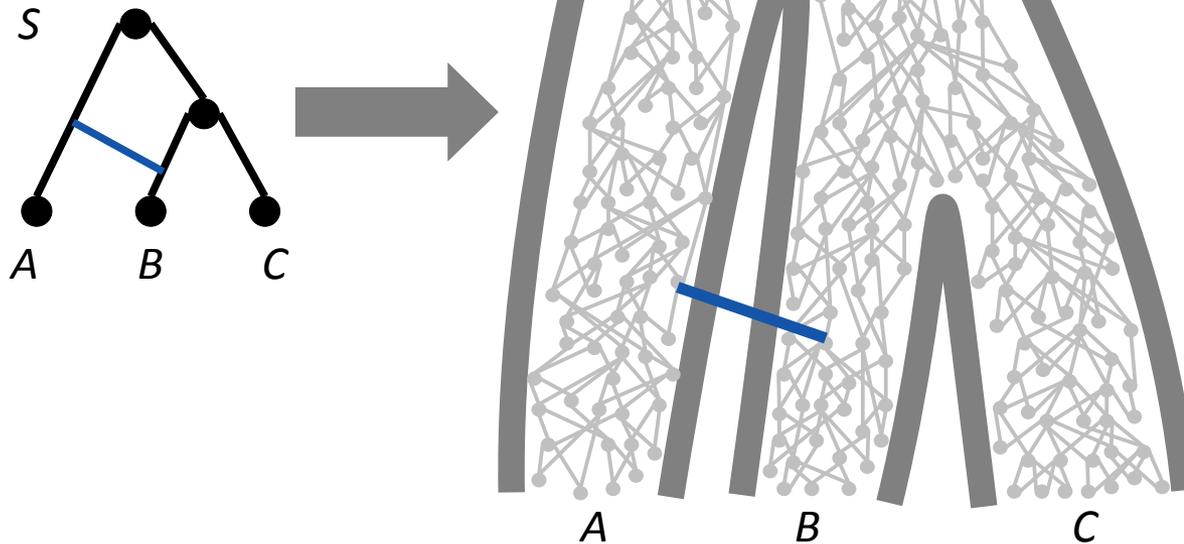
- Les organiser en fonction de caractères communs
- Décrire leur **évolution**



Transferts de matériel génétique

Transferts de matériel génétique entre espèces coexistantes :

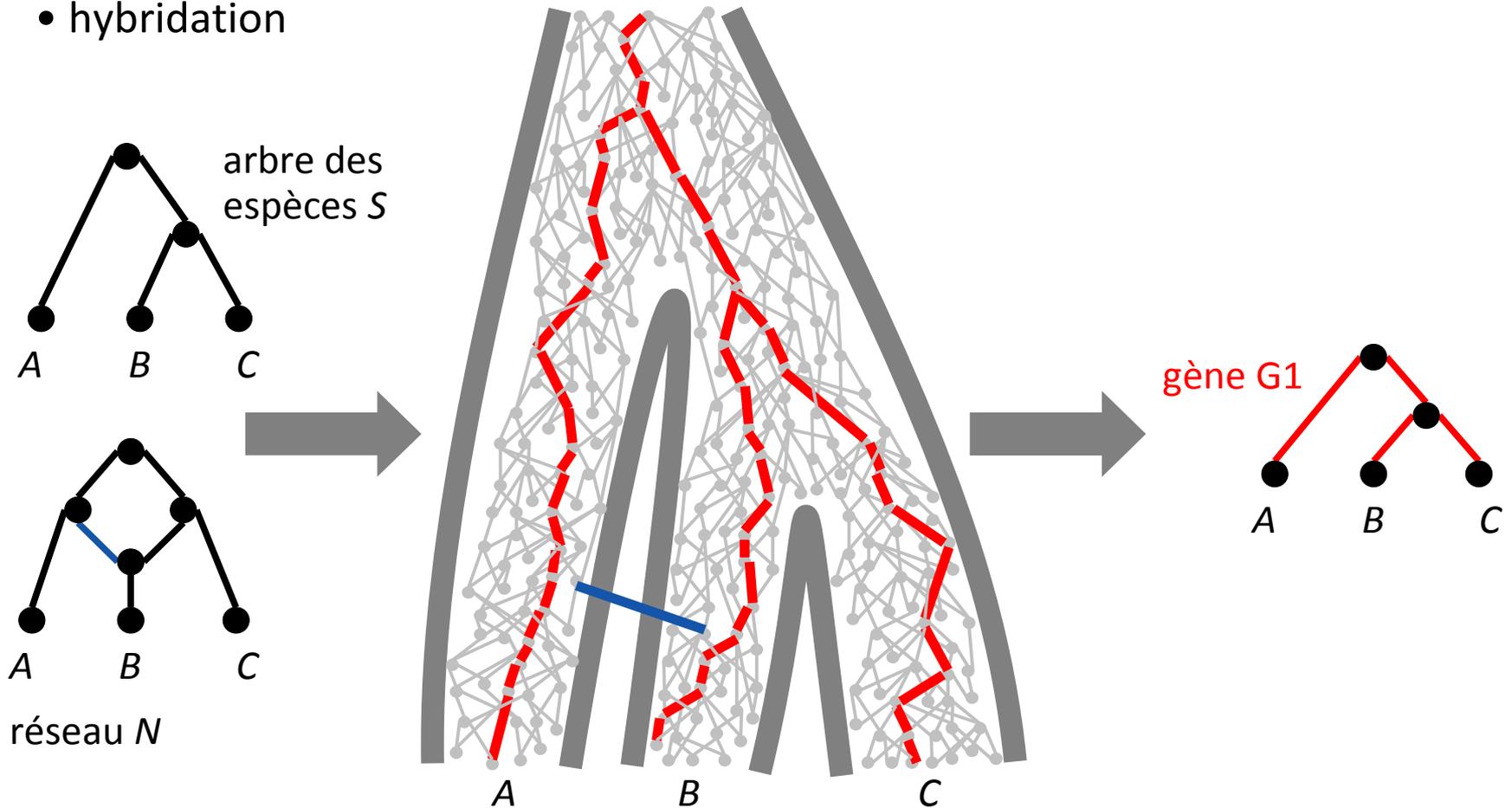
- transfert horizontal
- hybridation



Transferts de matériel génétique

Transferts de matériel génétique entre espèces coexistantes :

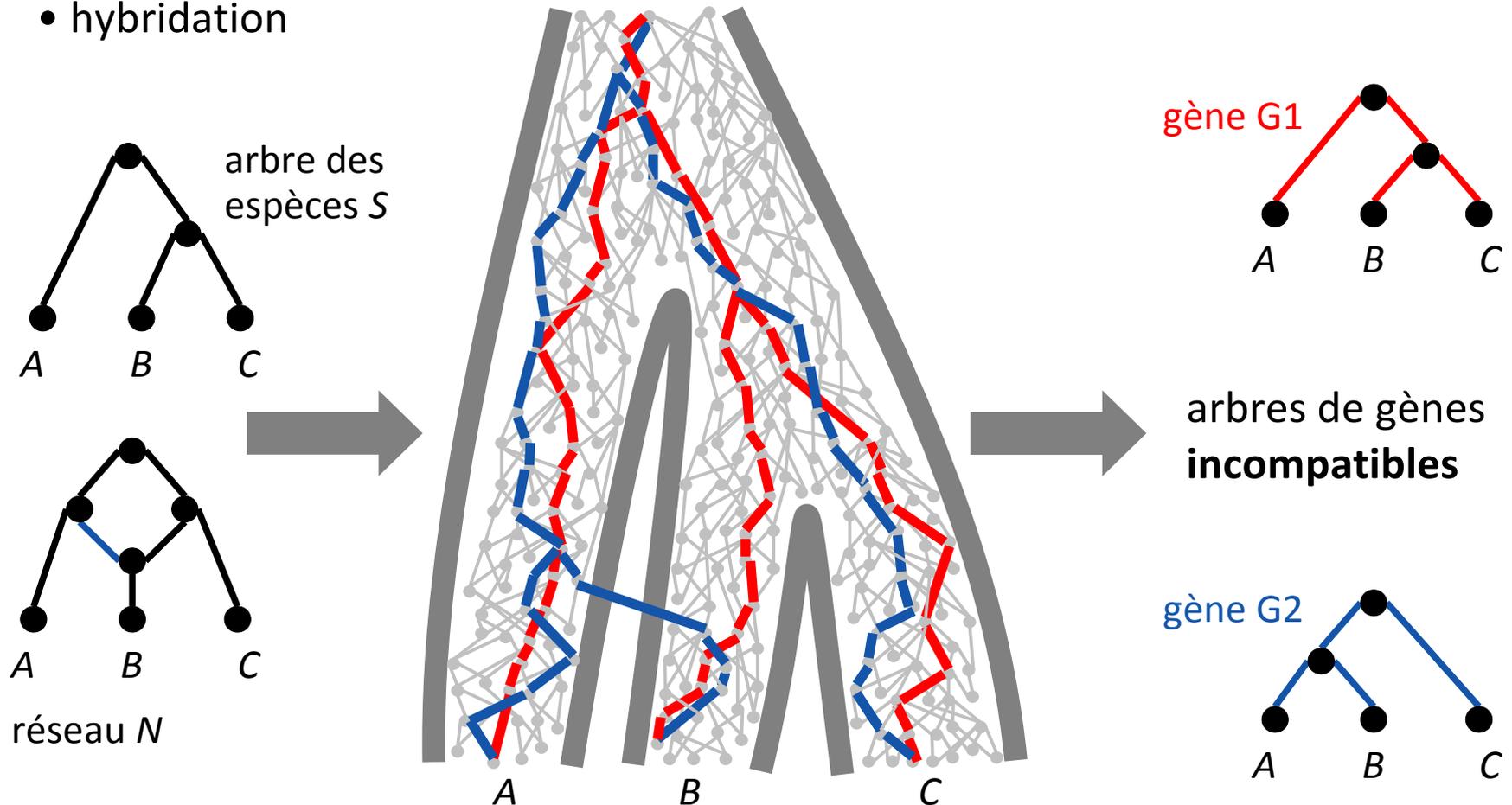
- transfert horizontal
- hybridation



Transferts de matériel génétique

Transferts de matériel génétique entre espèces coexistantes :

- transfert horizontal
- hybridation



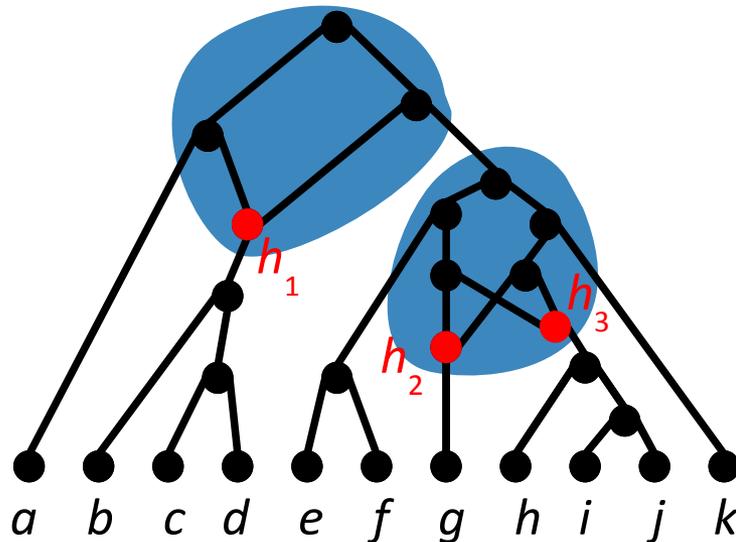
Réseaux phylogénétiques explicites

Réseau phylogénétique : réseau représentant des données d'évolution

- réseaux phylogénétiques **explicités**

modélisation de l'évolution

Réseau phylogénétique explicite enraciné :



Sommets à plus d'un parent :
réticulations

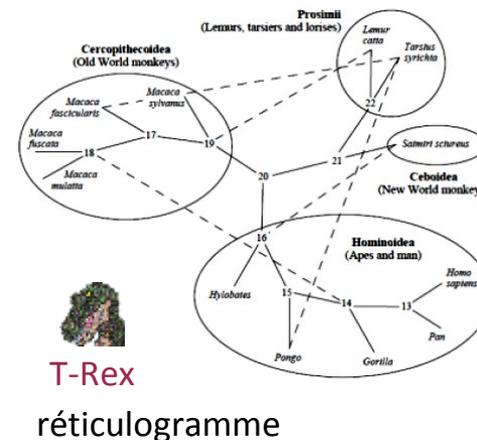
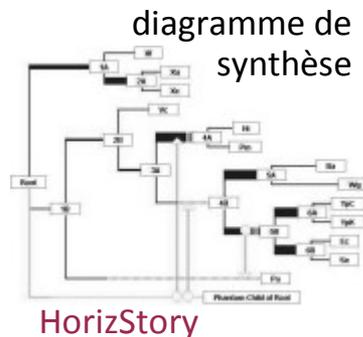
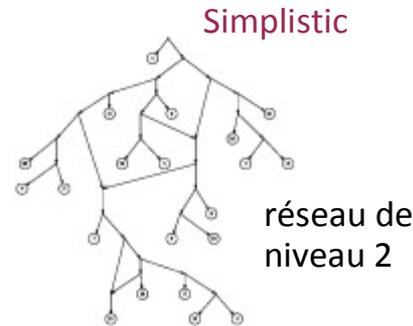
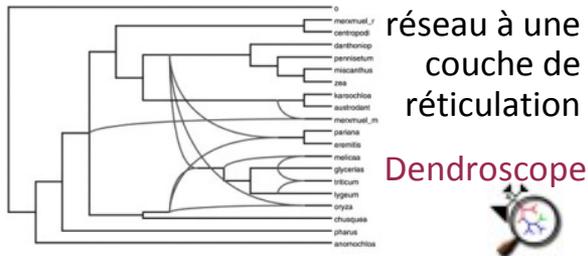
Partie non arborée : ***blob***.

Réseaux phylogénétiques explicites

Réseau phylogénétique : réseau représentant des données d'évolution

- réseaux phylogénétiques **explicites**

modélisation de l'évolution



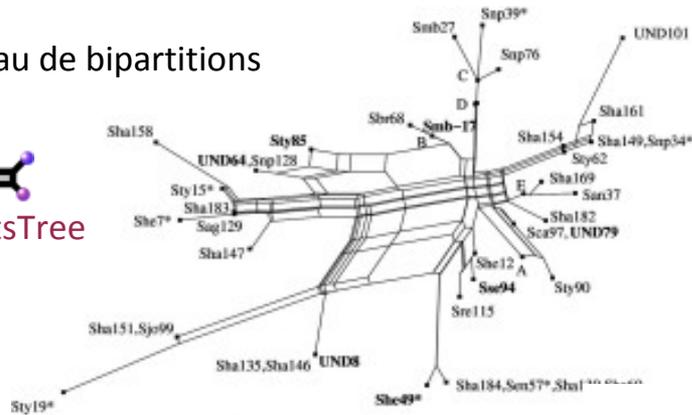
Réseaux phylogénétiques abstraits

Réseau phylogénétique : réseau représentant des données d'évolution

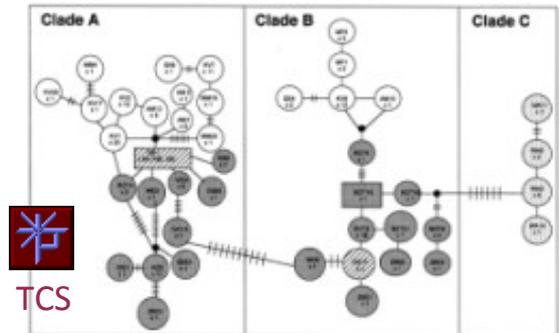
- réseaux phylogénétiques **abstraits**

classification, visualisation de données

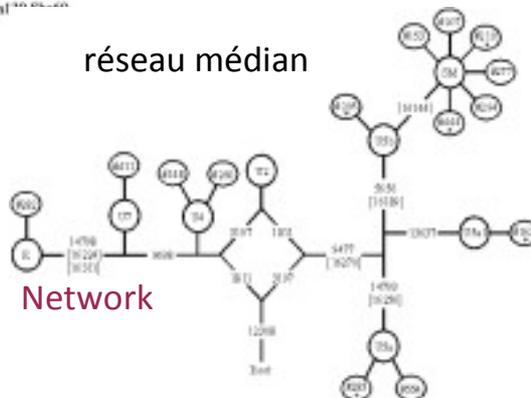
réseau de bipartitions



réseau couvrant minimum



réseau médian



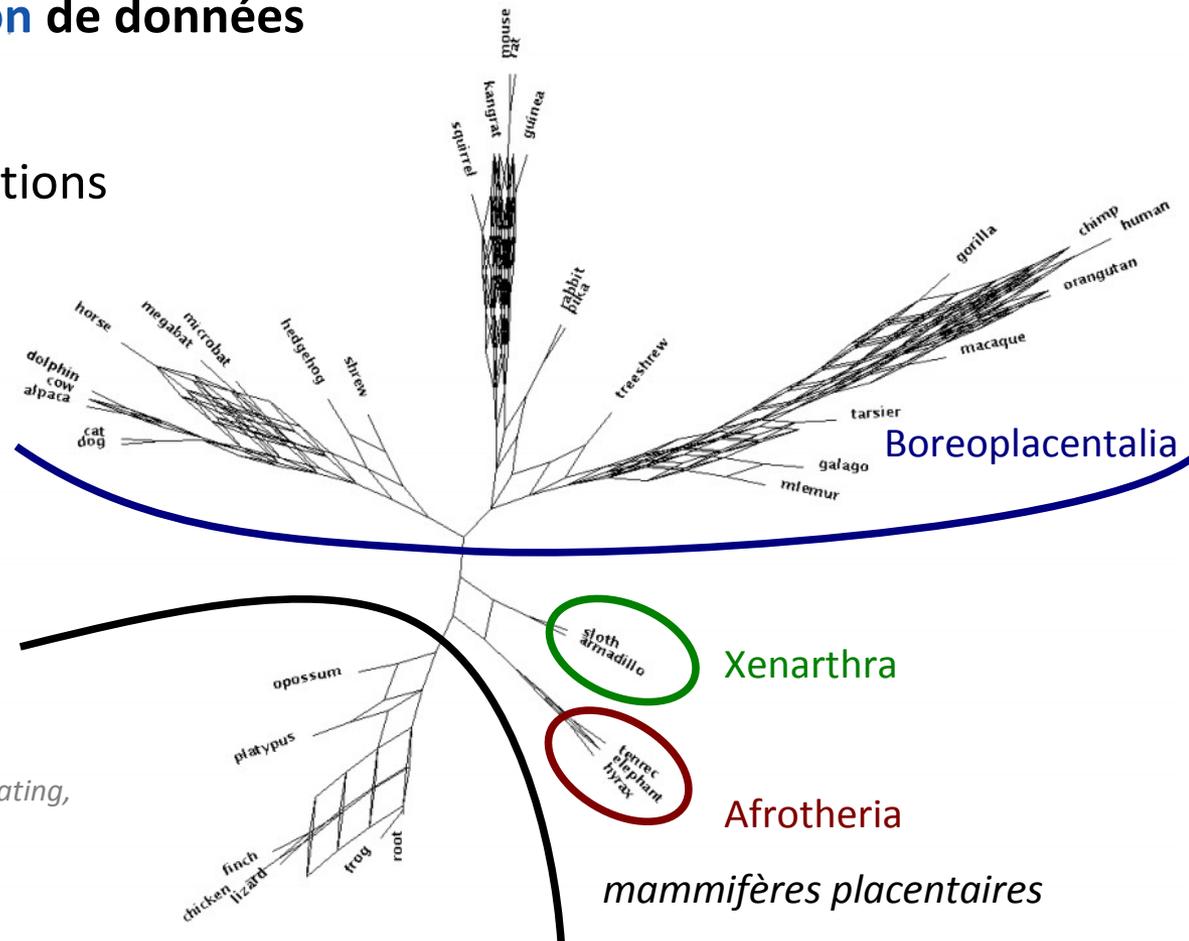
Réseaux phylogénétiques abstraits

Réseau phylogénétique : réseau représentant des données d'évolution

- réseaux phylogénétiques **abstraits**

classification, visualisation de données

Réseau abstrait de bipartitions
("consensus network") :



D'après Björn M. Hallström, Axel Janke -
Mammalian evolution may not be strictly bifurcating,
MBE, 2010

Réseaux phylogénétiques abstraits

Réseau phylogénétique : réseau représentant des données d'évolution

- réseaux phylogénétiques **abstraits**

classification, visualisation de données

Réseau abstrait de bipartitions
("consensus network") :



D'après Björn M. Hallström, Axel Janke -
Mammalian evolution may not be strictly bifurcating,
MBE, 2010

Réseaux phylogénétiques abstraits

Réseau phylogénétique : réseau représentant des données d'évolution

- réseaux phylogénétiques **abstraits**

classification, visualisation de données

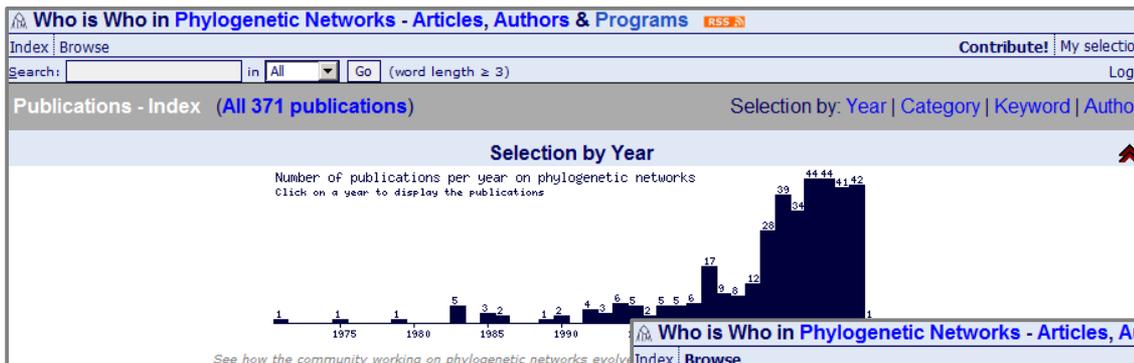
Réseau abstrait de bipartitions
("consensus network") :



D'après Björn M. Hallström, Axel Janke -
Mammalian evolution may not be strictly bifurcating,
MBE, 2010

Plate-forme bibliographique

Who is Who in Phylogenetic Networks, Articles, Authors & Programs



Selection by Year

Selection by Category

- Article (Journal) (216)
- Book (1)
- Misc (19)
- InProceedings (19)
- PhdThesis (19)
- Programs (52)

Selection by Key

abstract-network(46) approximation(8) APX-hard(2) ARG branch-and-bound(1) cactus-graph(1) characterization(7) circular-split-consistency(2) cophylogeny(1) distance-between-networks(21) divergence(1) explicit-network(93) exponential-algorithm(2) FPT(16) from-cluster(6) from-network(12) from-quartets(7) from-rooted-tree(26) from-splits(9) from-trees(6) from-triplets(17) from-ultra-metric(8) haplotype-network(2) haplotyping(1) heuristic(11) HMM programming(1) labeling(4) lateral-gene-transfer(35) level-k-sorting(5) MASN(4) median-network(15) MedianJoining(2) minisort(2) mu-distance(2) NeighborNet(11) nested-network(2) network-realization(2) parsimony(32) perfect(5) phylogenetic-network(46) polynomial(46) Program-Arlequin(5) Program-Beagle(3) Program-constNJ(1) Program-Dendroscope(7) Program-EEEEP(3) Program-GalledTree(4) Program-HybridInterleave(4) Program-HybridNET(1) Program-HybridNumber(3)

Publications related to 'Program Dendroscope': Dendroscope is an interactive viewer for large phylogenetic trees and networks. Available at www.dendroscope.org.

Associated keywords

abstract-network evaluation explicit-network FPT from-clusters from-rooted-trees galled-network level-k-phylogenetic-network NP-complete phylogenetic-network phylogeny polynomial Program-Bio-PhyloNetwork Program-Dendroscope Program-HybridInterleave Program-HybridNumber Program-NetGen Program-PhyloNet Program-SplitsTree Program-TCS reconstruction software split-network survey visualization

2010

1 Steven Kelk's k...
Leo van Iersel, Steven Kelk, Regula Rupp and Daniel H. Huson. Phylogenetic Networks Do not Need to Be Complex: Using Fewer Reticulations to Represent Conflicting Clusters. In *ISMB10*, Vol. 26(12):i124-i131 of *BIO*, 2010. [Comment]

Keywords: from clusters, level k phylogenetic network, Program Dendroscope, Program HybridInterleave, Program HybridNumber, reconstruction. Note: <http://dx.doi.org/10.1093/bioinformatics/btq202>.

2 Robert G. Beiko. Gene sharing and genome evolution: networks in trees and trees in networks. In *Biology and Philosophy*, 2010. [Comment]

Keywords: abstract network, explicit network, from rooted trees, galled network, phylogenetic network, phylogeny, Program Dendroscope, Program SplitsTree, reconstruction, split network, survey. Note: To appear, <http://dx.doi.org/10.1007/s10539-010-9217-3>.

Basé sur BibAdmin
par Sergiu Chelcea
+ nuages de mots, histogramme
des dates, liste des journaux,
graphes de co-auteurs,
définition des mots-clés.

Plan

- Les réseaux phylogénétiques
- **Méthodes de reconstruction**
- Limites des méthodes combinatoires
- Illustration sur des données biologiques
- Utilisation sur des données textuelles
- Perspectives

Reconstruction de réseaux phylogénétiques

espèce 1 : AATTGCAG TAGCCCAAAAT
espèce 2 : ACCTGCAG TAGACCAAT
espèce 3 : GCTTGCCG TAGACAAGAAT
espèce 4 : ATTTGCAG AAGACCAAAT
espèce 5 : TAGACAAGAAT
espèce 6 : ACTTGCAG TAGCACAAAAT
espèce 7 : ACCTGGTG TAAAAT

G1 G2

{séquences de gènes}

méthodes de distance

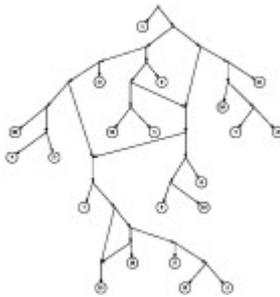
*Bandelt & Dress 1992 - Legendre & Makarenkov
2000 - Bryant & Moulton 2002*

méthodes de parcimonie

*Hein 1990 - Kececioglu & Gusfield 1994 - Jin,
Nakhleh, Snir, Tuller 2009*

méthodes de vraisemblance

*Snir & Tuller 2009 - Jin, Nakhleh, Snir, Tuller 2009 -
Velasco & Sober 2009*



réseau N

Reconstruction de réseaux phylogénétiques

espèce 1 : AATTGCAG TAGCCCAAAAT
espèce 2 : ACCTGCAG TAGACCAAT
espèce 3 : GCTTGCCG TAGACAAGAAT
espèce 4 : ATTTGCAG AAGACCAAAT
espèce 5 : TAGACAAGAAT
espèce 6 : ACTTGCAG TAGCACAAAAT
espèce 7 : ACCTGGTG TAAAAT

G1 G2

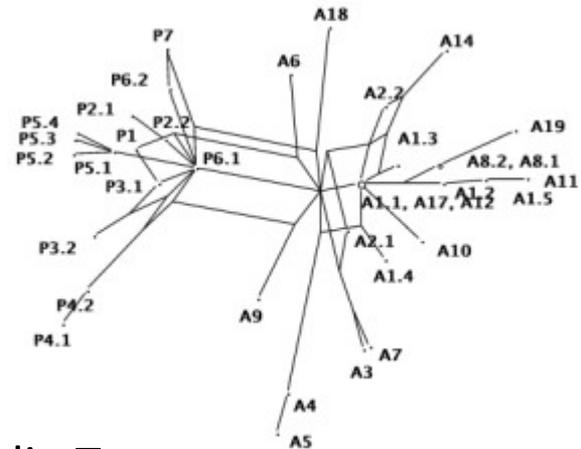
{séquences de gènes}

méthodes de distance

Bandelt & Dress 1992 - Legendre & Makarenkov
2000 - Bryant & Moulton 2002

"Split decomposition"

réseau N



Reconstruction de réseaux phylogénétiques

espèce 1 : AATTGCAG TAGCCCAAAAT
espèce 2 : ACCTGCAG TAGACCAAT
espèce 3 : GCTTGCCG TAGACAAGAAT
espèce 4 : ATTTGCAG AAGACCAAAT
espèce 5 : TAGACAAGAAT
espèce 6 : ACTTGCAG TAGCACAAAAT
espèce 7 : ACCTGGTG TAAAAT

G1 G2

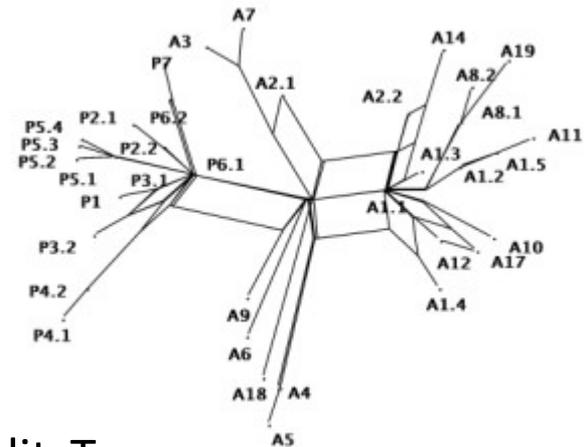
{séquences de gènes}

méthodes de distance

Bandelt & Dress 1992 - Legendre & Makarenkov
2000 - **Bryant & Moulton 2002**

“Neighbor Net”

réseau N



 SplitsTree

Reconstruction de réseaux phylogénétiques

espèce 1 : AATTGCAG TAGCCCAAAAT
espèce 2 : ACCTGCAG TAGACCAAT
espèce 3 : GCTTGCCG TAGACAAGAAT
espèce 4 : ATTTGCAG AAGACCAAAT
espèce 5 : TAGACAAGAAT
espèce 6 : ACTTGCAG TAGCACAAAAT
espèce 7 : ACCTGGTG TAAAAT

G1 G2

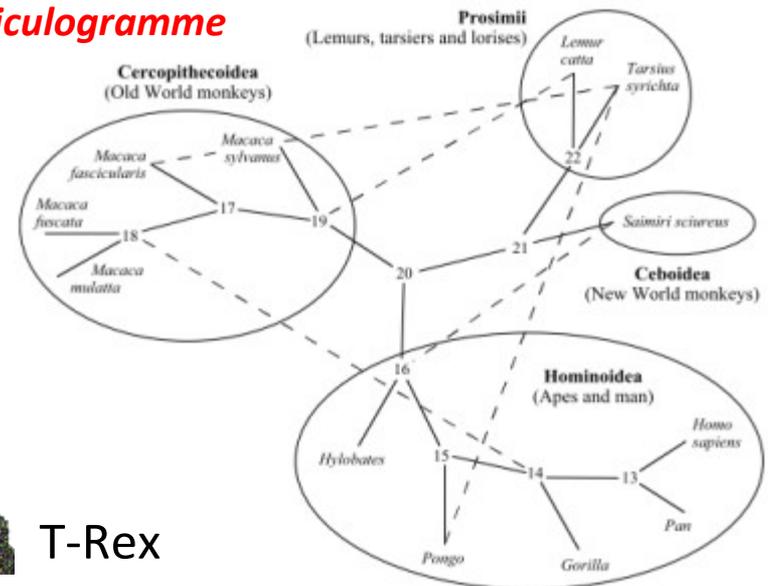
{séquences de gènes}

méthodes de distance

Bandelt & Dress 1992 - Legendre & Makarenkov
2000 - Bryant & Moulton 2002

réseau N

Réticulogramme



T-Rex

D'après Vladimir Makarenkov & Pierre Legendre,
Improving the additive tree representation of a
dissimilarity matrix using reticulations, DACRM, 2000

Reconstruction de réseaux phylogénétiques

**Problème : méthodes généralement lentes,
explosion du nombre de séquences.**

espèce 1 : AATTGCAG TAGCCCAAAAT
espèce 2 : ACCTGCAG TAGACCAAT
espèce 3 : GCTTGCCG TAGACAAGAAT
espèce 4 : ATTTGCAG AAGACCAAAT
espèce 5 : TAGACAAGAAT
espèce 6 : ACTTGCAG TAGCACAAAAT
espèce 7 : ACCTGGTG TAAAAT

G1 **G2**

{séquences de gènes}

méthodes de distance

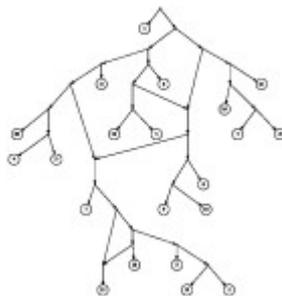
*Bandelt & Dress 1992 - Legendre & Makarenkov
2000 - Bryant & Moulton 2002*

méthodes de parcimonie

*Hein 1990 - Kececioglu & Gusfield 1994 - Jin,
Nakhleh, Snir, Tuller 2009*

méthodes de vraisemblance

*Snir & Tuller 2009 - Jin, Nakhleh, Snir, Tuller 2009 -
Velasco & Sober 2009*

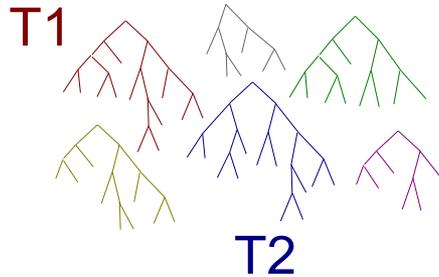


réseau *N*

Reconstruction de réseaux phylogénétiques

espèce 1 : AATTGCAG TAGCCCAAAAT
espèce 2 : ACCTGCAG TAGACCAAT
espèce 3 : GCTTGCCG TAGACAAGAAT
espèce 4 : ATTTGCAG AAGACCAAAT
espèce 5 : TAGACAAGAAT
espèce 6 : ACTTGCAG TAGCACAAAAT
espèce 7 : ACCTGGTG TAAAAT

G1 G2



{séquences de gènes}

Reconstruction d'un arbre pour chaque gène présent chez plusieurs espèces

Guindon & Gascuel, SB, 2003

{arbres}

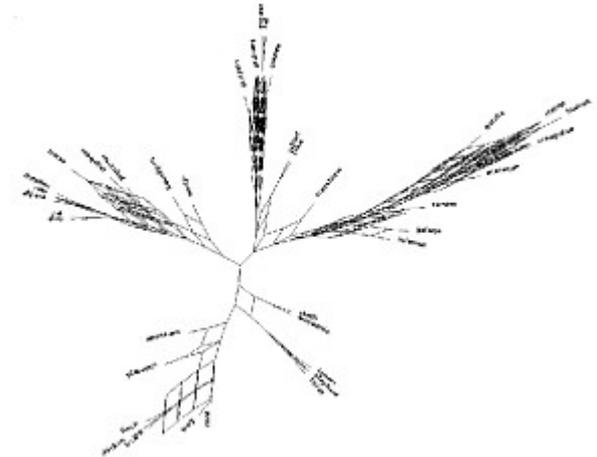
Base HOGENOM



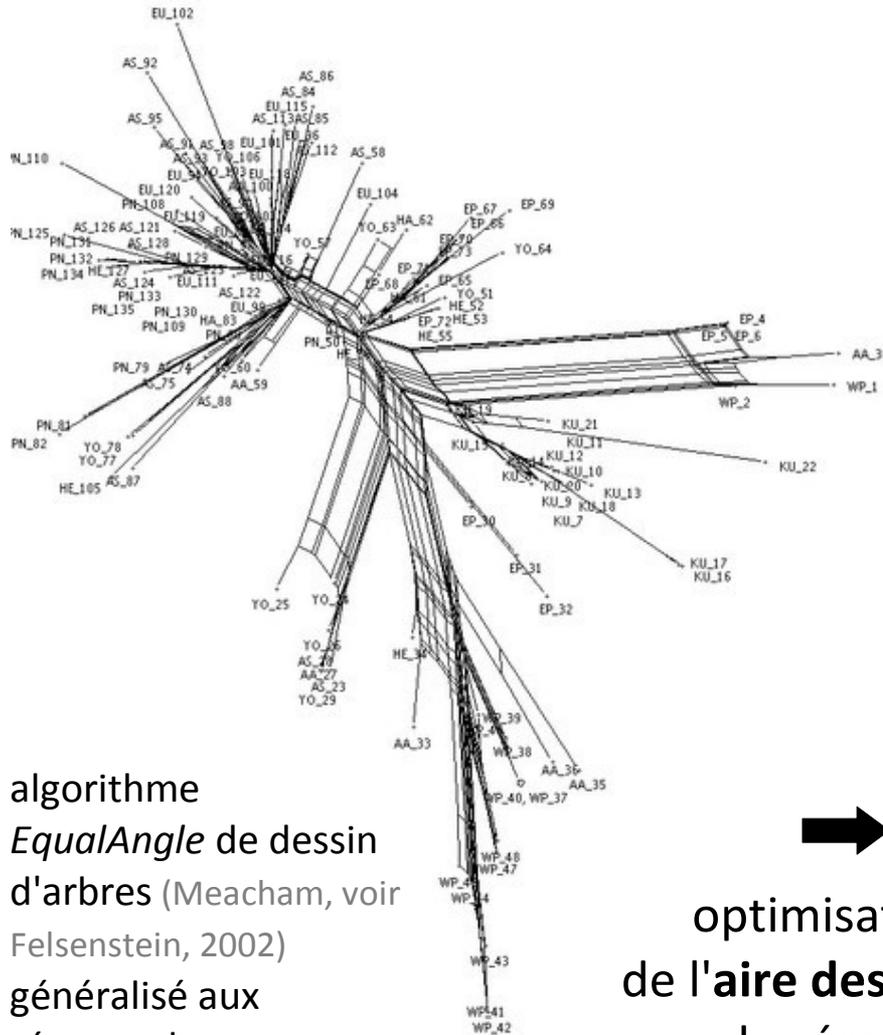
Dufayard, Duret, Penel, Gouy, Rechenmann & Perrière, BioInf, 2005

Réconciliation ou **consensus** d'arbres

super-réseau de consensus N



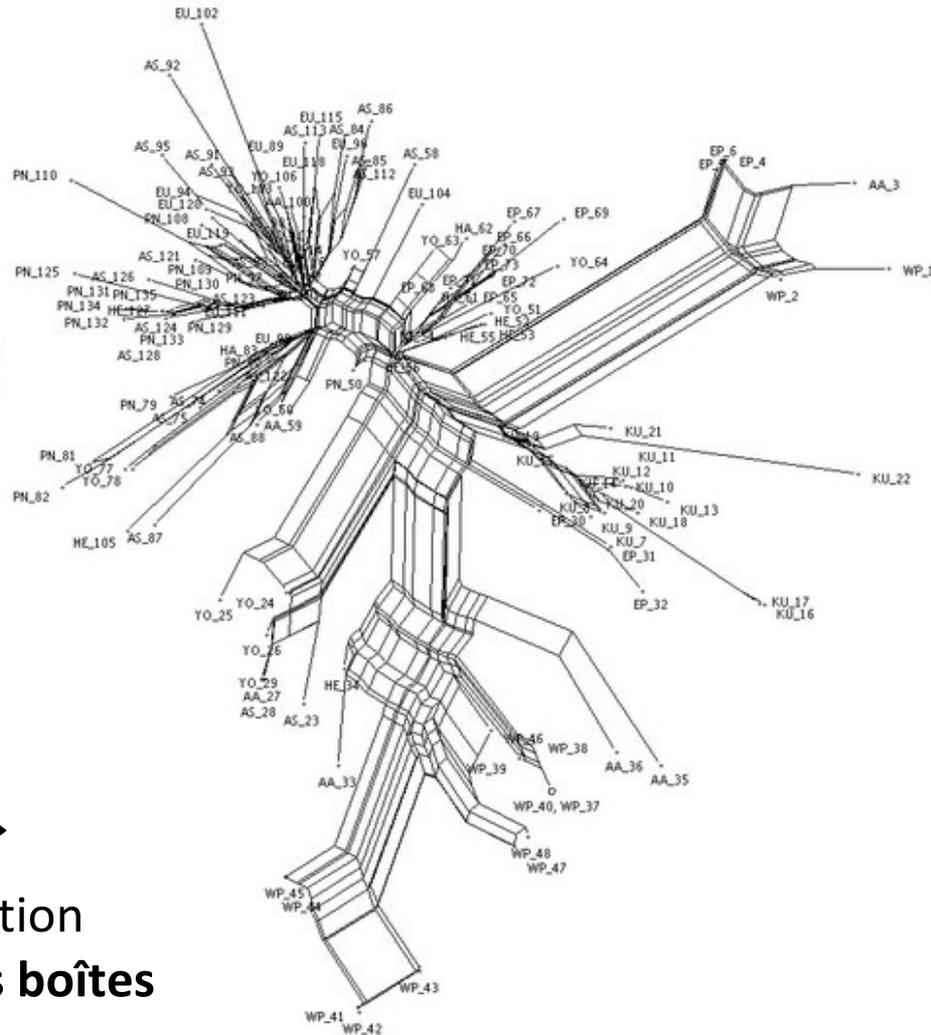
Visualisation des réseaux de bipartitions



algorithme
EqualAngle de dessin
d'arbres (Meacham, voir
Felsenstein, 2002)
généralisé aux
réseaux de
bipartitions.



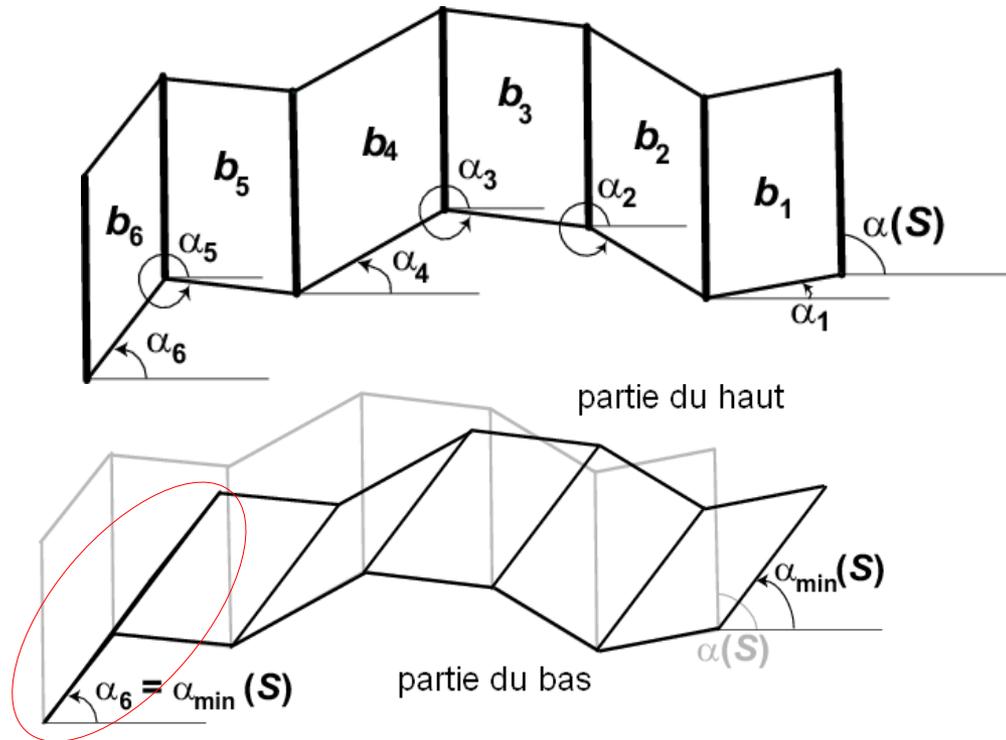
optimisation
de l'aire des boîtes
du réseau



Algorithme Box-opening

Collisions **locales** :

deux angles critiques $\alpha_{\min}(S)$ et $\alpha_{\max}(S)$ pour l'angle de la bipartition $\alpha(S)$

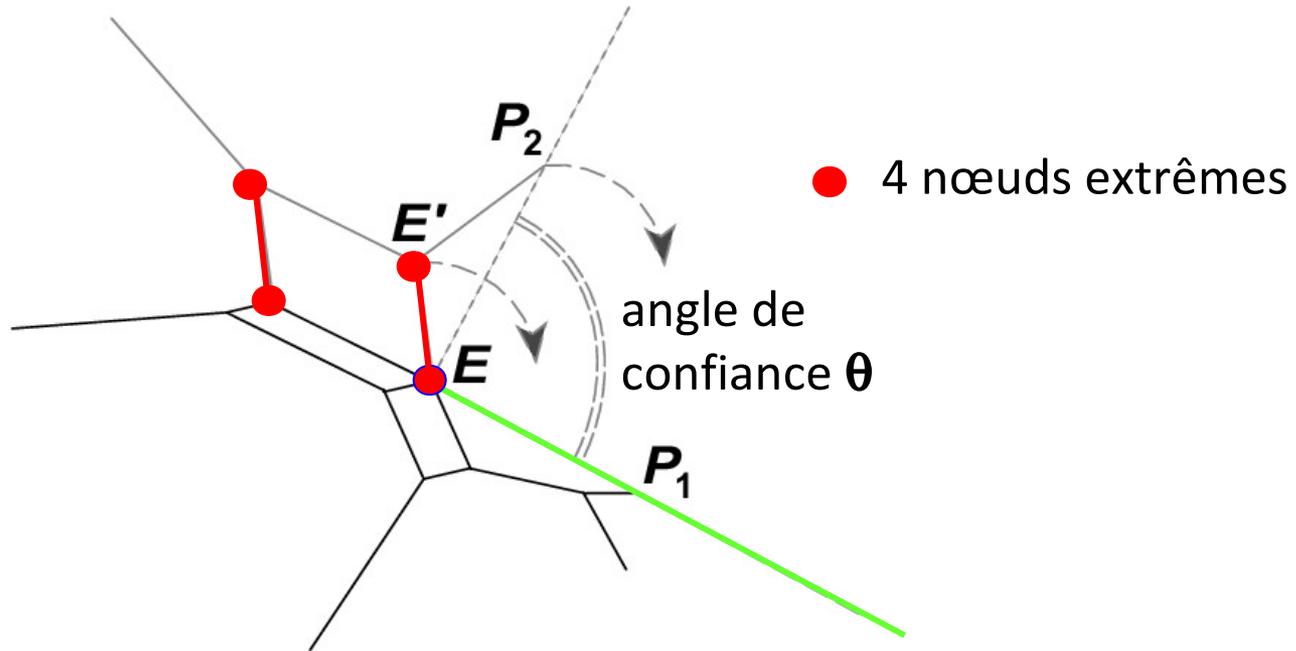


$$\alpha_{\max}(S) = \alpha(S) + \min_{\text{boîte } b_i} \{(\alpha_i - \alpha(S) - \pi) \bmod 2\pi\}$$

$$\alpha_{\min}(S) = \alpha(S) - \min_{\text{boîte } b_i} \{(\alpha(S) - \alpha_i) \bmod 2\pi\}$$

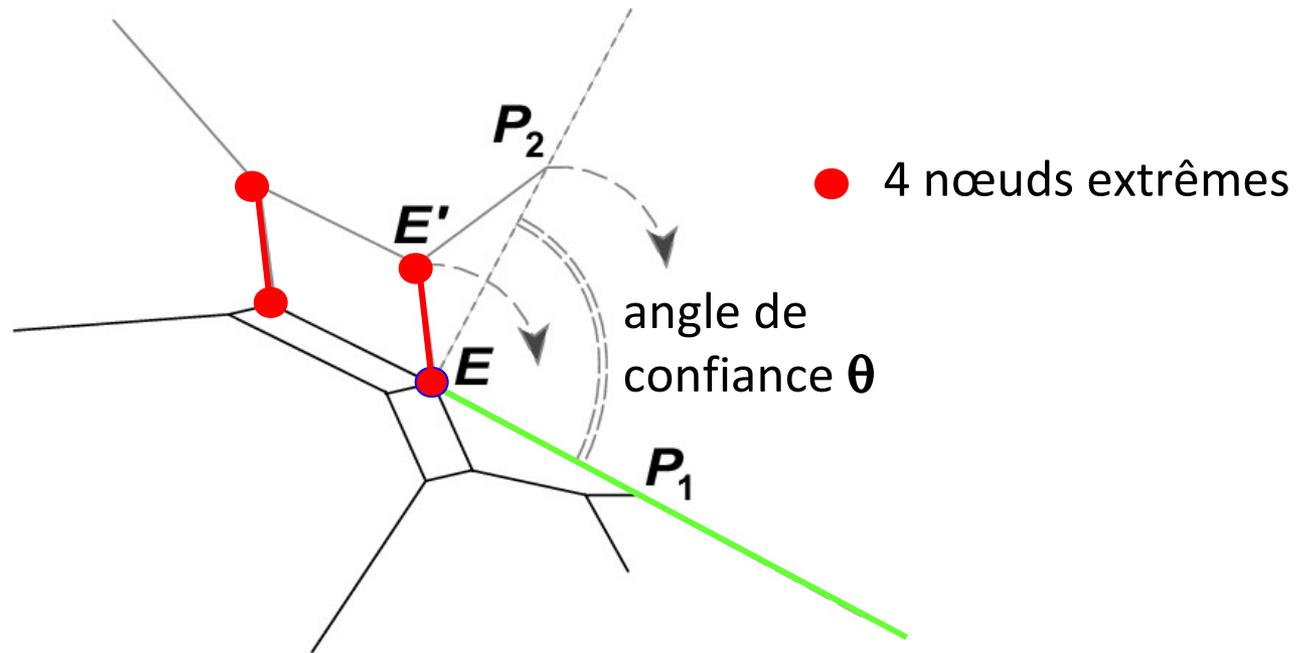
Algorithme Box-opening

Choix de l'angle $\alpha(S)$: collisions **globales**



Algorithme Box-opening

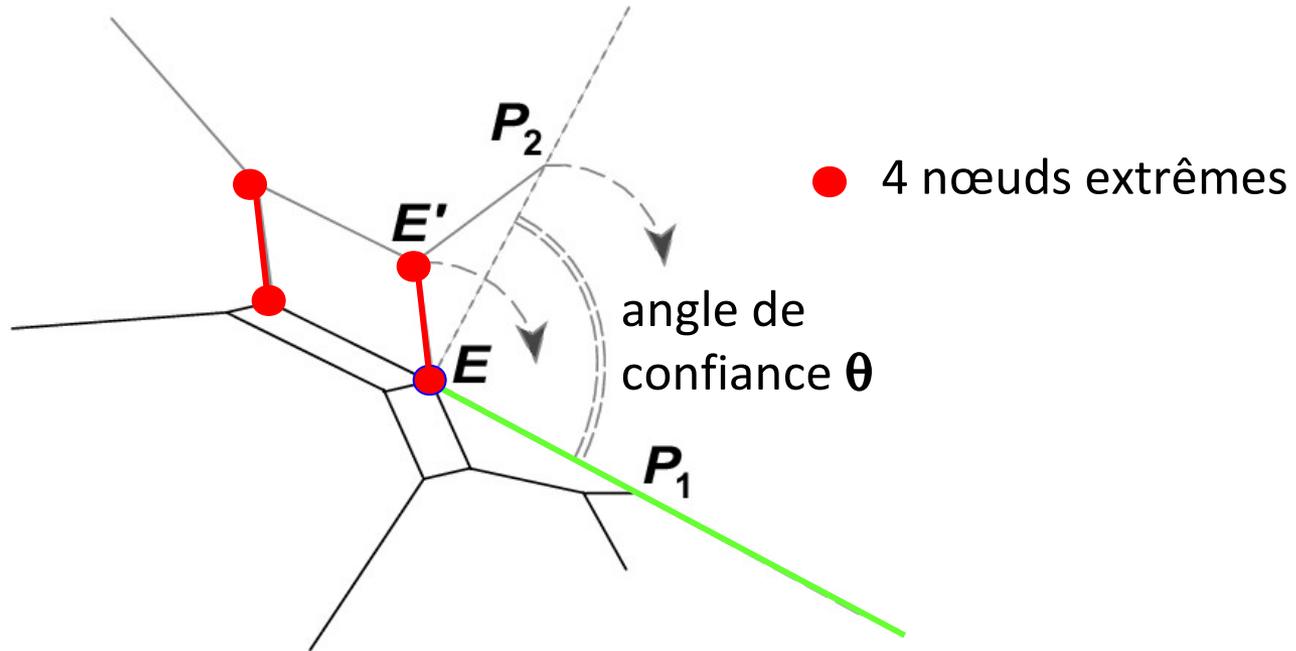
Choix de l'angle $\alpha(S)$: collisions **globales**



➔ optimisations locales de l'aire

Algorithme Box-opening

Choix de l'angle $\alpha(S)$: collisions **globales**



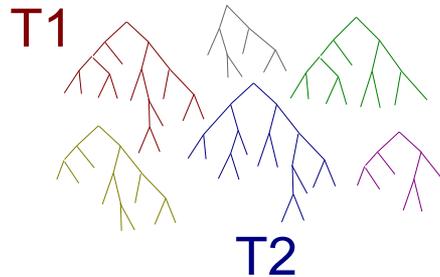
➡ **optimisations locales** de l'aire

➡ **optimisations globale** du réseau
métaheuristique (recuit simulé)

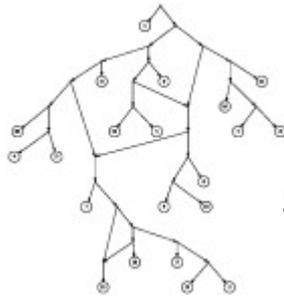
Reconstruction de réseaux phylogénétiques

espèce 1 : AATTGCAG TAGCCCAAAAT
espèce 2 : ACCTGCAG TAGACCAAT
espèce 3 : GCTTGCCG TAGACAAGAAT
espèce 4 : ATTTGCAG AAGACCAAAT
espèce 5 : TAGACAAGAAT
espèce 6 : ACTTGCAG TAGCACAAAAT
espèce 7 : ACCTGGTG TAAAAT

G1 G2



réseau
explicite



{séquences de gènes}

*Reconstruction d'un arbre pour chaque
gène présent chez plusieurs espèces*

Guindon & Gascuel, SB, 2003

{arbres}

Base HOGENOM



Dufayard, Duret, Penel, Gouy,
Rechenmann & Perrière, BioInf, 2005

Réconciliation ou consensus d'arbres

super-réseau optimal N

Reconstruction de réseaux phylogénétiques

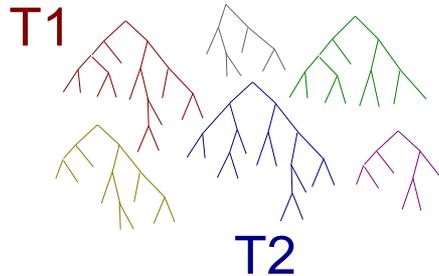
espèce 1 : AATTGCAG TAGCCCAAAAT
espèce 2 : ACCTGCAG TAGACCAAT
espèce 3 : GCTTGCCG TAGACAAGAAT
espèce 4 : ATTTGCAG AAGACCAAAT
espèce 5 : TAGACAAGAAT
espèce 6 : ACTTGCAG TAGCACAAAAT
espèce 7 : ACCTGGTG TAAAAT

G1 G2

{séquences de gènes}

Reconstruction d'un arbre pour chaque
gène présent chez plusieurs espèces

Guindon & Gascuel, SB, 2003



{arbres}

Base HOGENOM

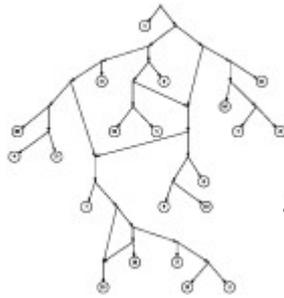


Dufayard, Duret, Penel, Gouy,
Rechenmann & Perrière, BioInf, 2005

> 500 espèces, >70 000 arbres

Réconciliation ou consensus d'arbres

réseau
explicite



super-réseau optimal N

Problème : la réconciliation d'arbres est un problème difficile

(NP-complet pour 2 arbres avec le minimum d'hybridations)

Bordewich & Semple, DAM, 2007

Triplets et quadruplets, clades et bipartitions

Problème :

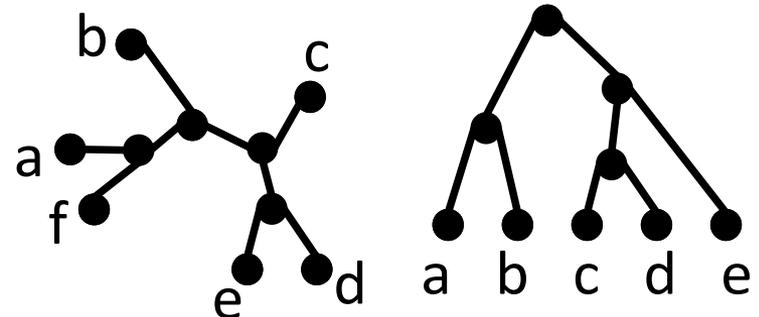
Reconstruire le **super-réseau** d'un ensemble d'arbres est
difficile.

Idée :

reconstituer un réseau contenant tous les :

triplets
quadruplets
clades
bipartitions

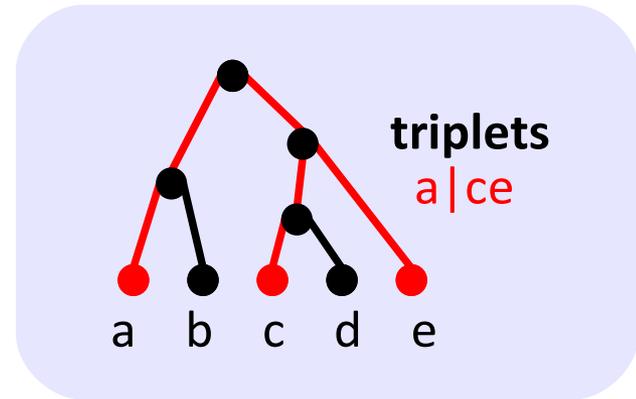
des arbres en entrée ?



Triplets et quadruplets, clades et bipartitions

Idée :

reconstituer un réseau contenant tous les :



des arbres en entrée ?

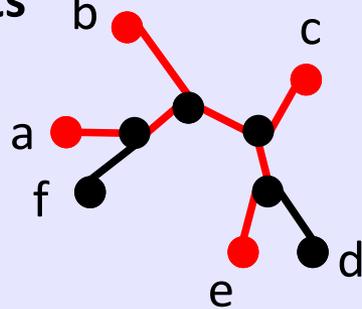
Triplets et quadruplets, clades et bipartitions

Idée :

reconstituer un réseau contenant tous les :

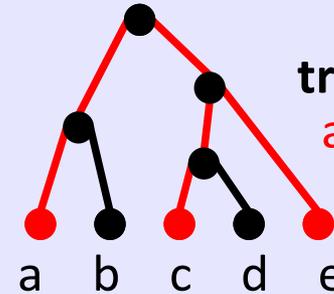
quadruplets

$ab|ce$



triplets

$a|ce$



des arbres en entrée ?

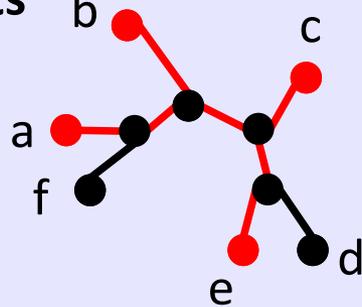
Triplets et quadruplets, clades et bipartitions

Idée :

reconstituer un réseau contenant tous les :

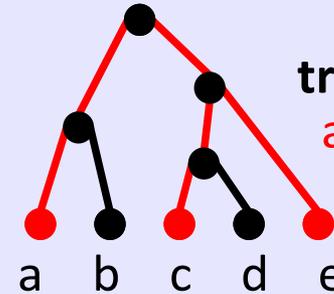
quadruplets

$ab|ce$



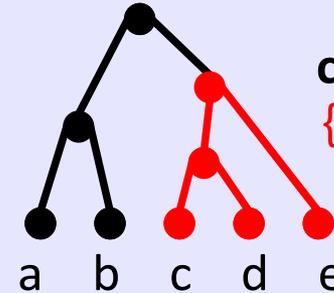
triplets

$a|ce$



clades

$\{c,d,e\}$



des arbres en entrée ?

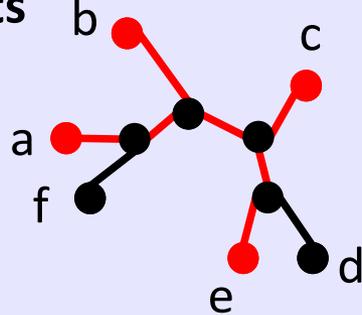
Triplets et quadruplets, clades et bipartitions

Idée :

reconstituer un réseau contenant tous les :

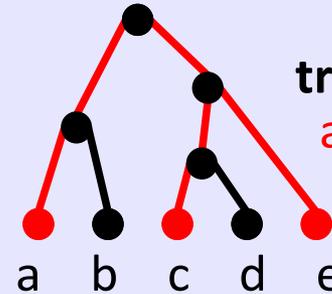
quadruplets

$ab|ce$



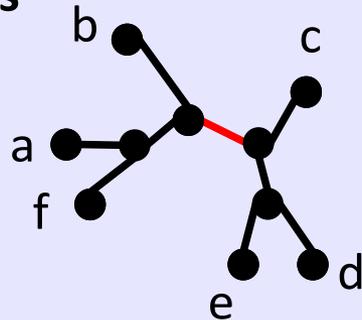
triplets

$a|ce$



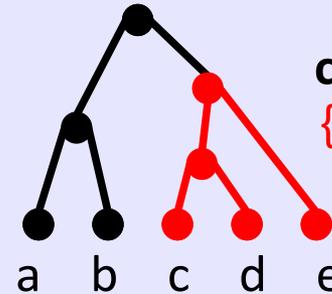
bipartitions

$\{a,b,f\}$
 $\{c,d,e\}$



clades

$\{c,d,e\}$



des arbres en entrée ?

Reconstruction depuis les triplets

{arbres}

Reconstruction d'un réseau de **niveau k** à partir d'un ensemble de **triplets**

Jansson, Nguyen & Sung, JOC, 2006 : NP-complet pour niveau 1,
Van Iersel, Kelk & Mních, JBCB, 2009 : NP-complet pour niveau k

{triplets}

niveau =
mesure de "complexité", d'éloignement par rapport à une
structure d'arbre.

Choy, Jansson, Sadakane & Sung, TCS, 2005



N'
réseau
de niveau k

Reconstruction depuis les triplets

{arbres}

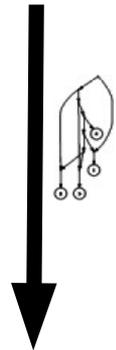
Reconstruction d'un réseau de **niveau k** à partir d'un ensemble de **triplets**

Jansson, Nguyen & Sung, JOC, 2006 : NP-complet pour niveau 1,
Van Iersel, Kelk & Mních, JBCB, 2009 : NP-complet pour niveau k

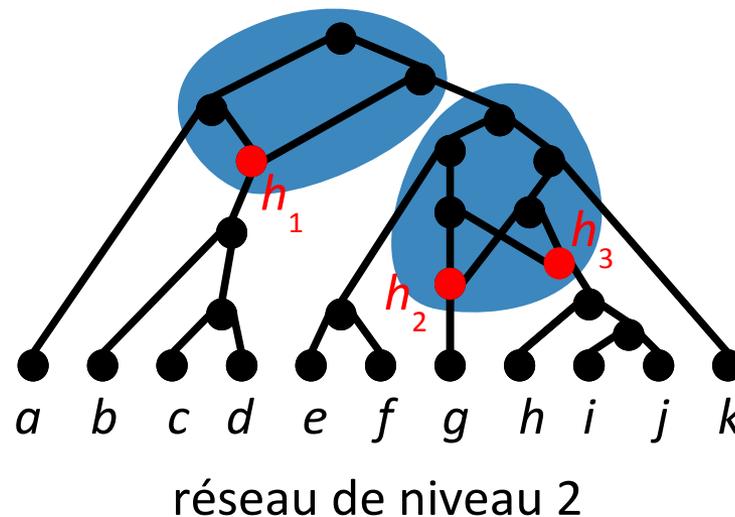
{ triplets}

niveau =
nombre maximum d'hybridations par partie non arborée (*blob*).

Choy, Jansson, Sadakane & Sung, TCS, 2005



N'
réseau
de niveau k



Reconstruction depuis les triplets

{arbres}

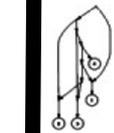
Reconstruction d'un réseau de **niveau k** à partir d'un ensemble de **triplets**

Jansson, Nguyen & Sung, JOC, 2006 : NP-complet pour niveau 1,
Van Iersel, Kelk & Mnich, JBCB, 2009 : NP-complet pour niveau k

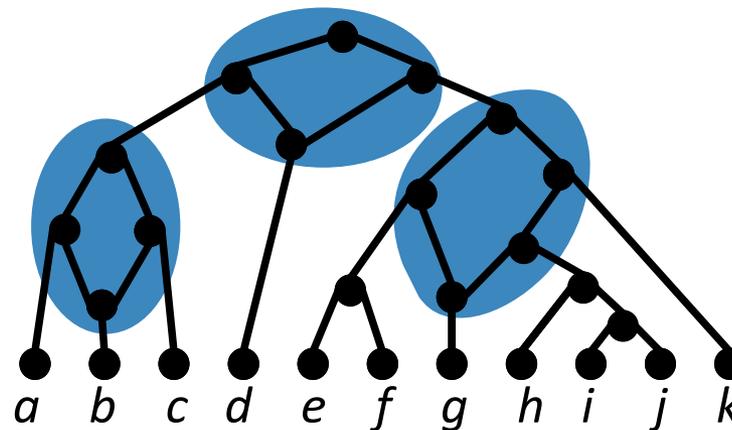
{triplets}

niveau =
nombre maximum d'hybridations par partie non arborée (*blob*).

Choy, Jansson, Sadakane & Sung, TCS, 2005



N'
réseau
de niveau k

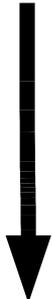


réseau de niveau 1 "galled tree"

Gusfield, Eddhu & Langley, IJOC, 2004

Reconstruction depuis les triplets

{arbres}



{triplets}



Méthodes exactes rapides pour reconstruire un **réseau de niveau 1 et 2** (s'il en existe un) à partir d'un ensemble **dense de triplets**

Jansson, Nguyen & Sung, SODA'05 : $O(n^3)$ pour niveau 1,

van Iersel, Kelk & al, RECOMB'08 : $O(n^8)$ pour niveau 2,

To & Habib, CPM'09 : $O(n^{5k+4})$ pour niveau k

Van Iersel & Kelk, J. Theor. Biol., 2011 : NP-complet de trouver le niveau minimal

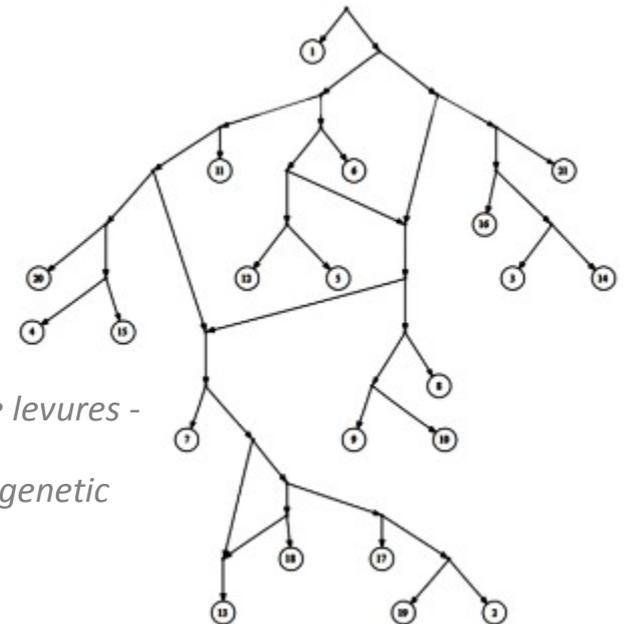
dense =

sur chaque ensemble de 3 feuilles, au moins 1 triplet existe dans T .

Programme Simplistic



N'
réseau
de niveau k



*Réseau phylogénétique de levures -
Van Iersel et al. :*

*Constructing level-2 phylogenetic
networks from triplets.*

RECOMB 2008

Reconstruction depuis les clades souples

{arbres}



{clades}



N'

réseau à 1
couche de
réticulation

Consensus de clades souples :

Dendroscope 

Huson et al., BMCB, 2007

Méthode exacte rapide de reconstruction de **réseaux à 1
couche de réticulation** à partir de **clades souples**

Huson, Rupp, Berry, Gambette & Paul, ISMB 2009

Méthode exacte de reconstruction de **réseaux de niveau k**
à partir de **clades souples**

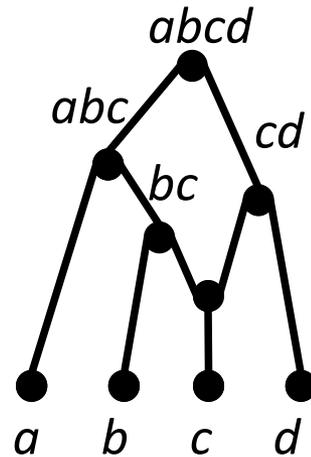
Iersel, Kelk, Rupp & Huson, ISMB 2010



meilleurs résultats mais plus lente pour niveau > 2 .
pour k fixé, certains ensembles de clades contenus
dans aucun réseau de niveau k .

Clades stricts et souples

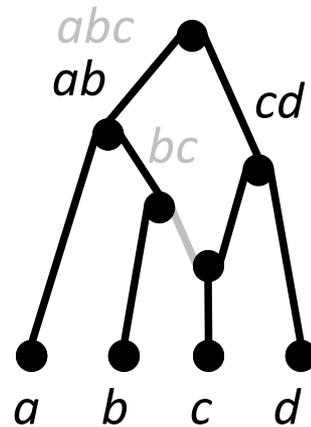
Clade "strict" : ensemble des feuilles sous un noeud du réseau



Clades stricts et souples

Clade “souple” : clade d'un arbre inclus dans le réseau

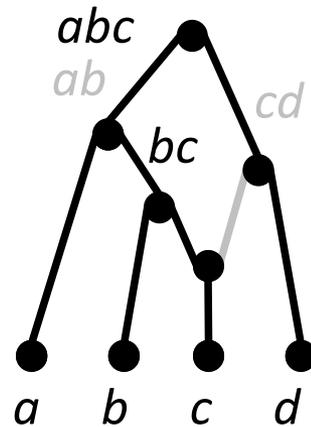
Modèle de **transmission arborée** des gènes
(gène transmis intégralement)



Clades stricts et souples

Clade “souple” : clade d'un arbre inclus dans le réseau

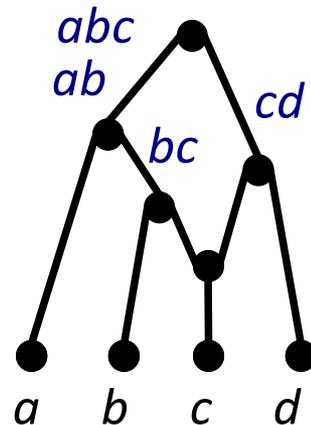
Modèle de **transmission arborée** des gènes
(gène transmis intégralement)



Clades stricts et souples

Modèle de **transmission arborée** des gènes
(gène transmis intégralement)

Clade “souple” : clade d'un arbre inclus dans le réseau



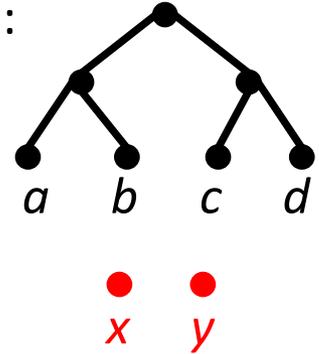
L'ensemble $S(N)$ de **tous les clades seulement compatibles** avec N peut être de taille **exponentielle**.

Tester si un **clade souple** appartient à un réseau : **NP-complet**.

Une approche de reconstruction en deux étapes

1- Trouver un **ensemble minimum de conflits** parmi les clades :

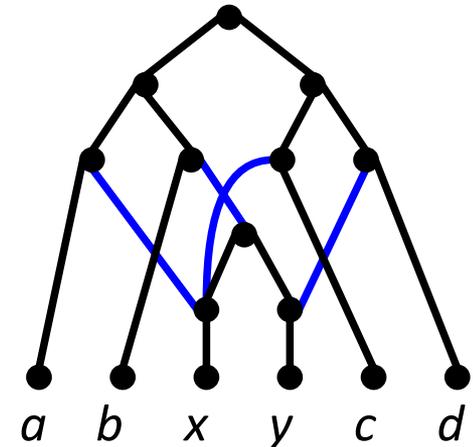
- partie sans conflits ➡ arbre,
- taxons impliqués dans des **conflits** ➡ sous les réticulations.



MAXIMUM COMPATIBLE SUBSET

2- Attacher à l'arbre les taxons impliqués dans des conflits avec un **nombre minimal d'arcs** :

MINIMUM ATTACHMENT PROBLEM



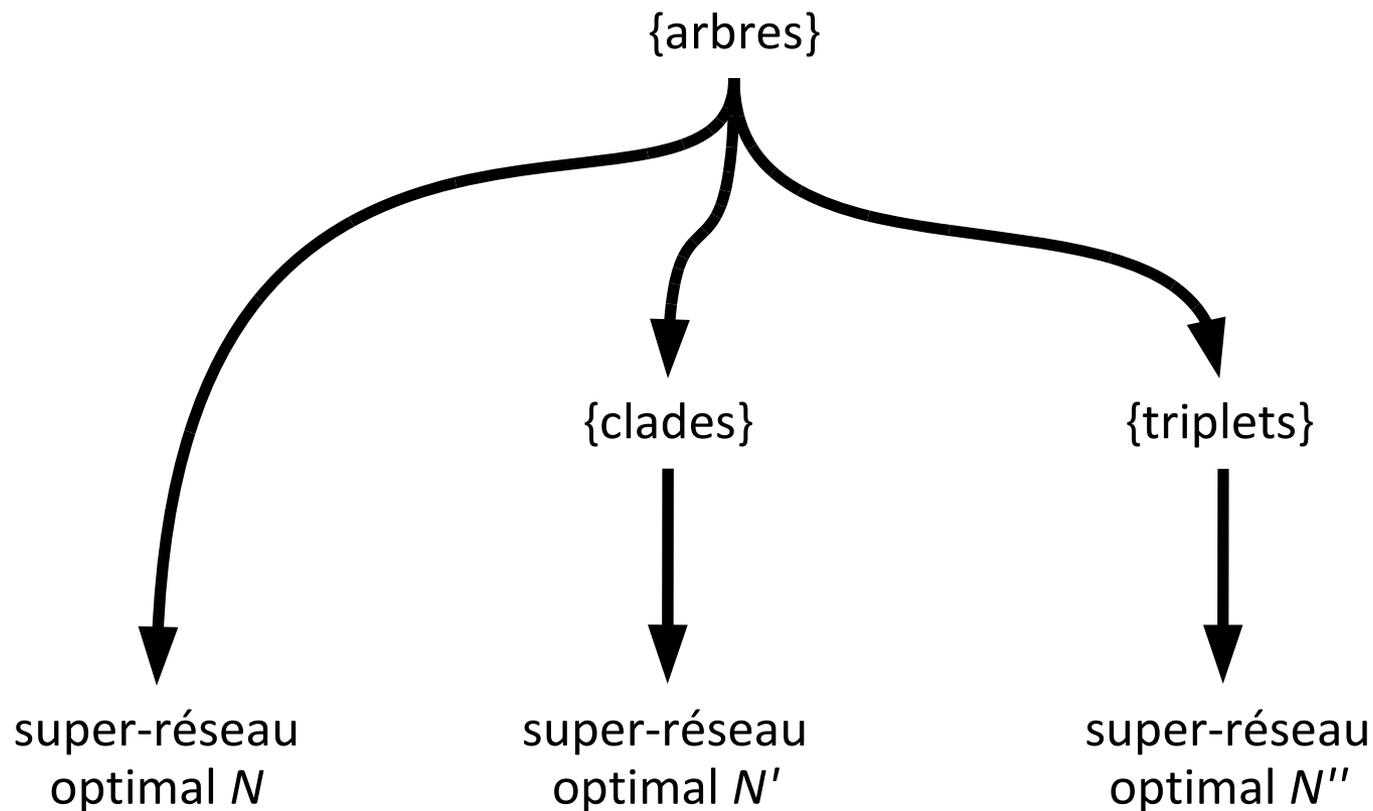
Plan

- Les réseaux phylogénétiques
- Méthodes de reconstruction
- **Limites des méthodes combinatoires**
- Illustration sur des données biologiques
- Utilisation sur des données textuelles
- Perspectives

Reconstruction combinatoire de réseaux phylogénétiques

Idée :

modifier le type de données à traiter

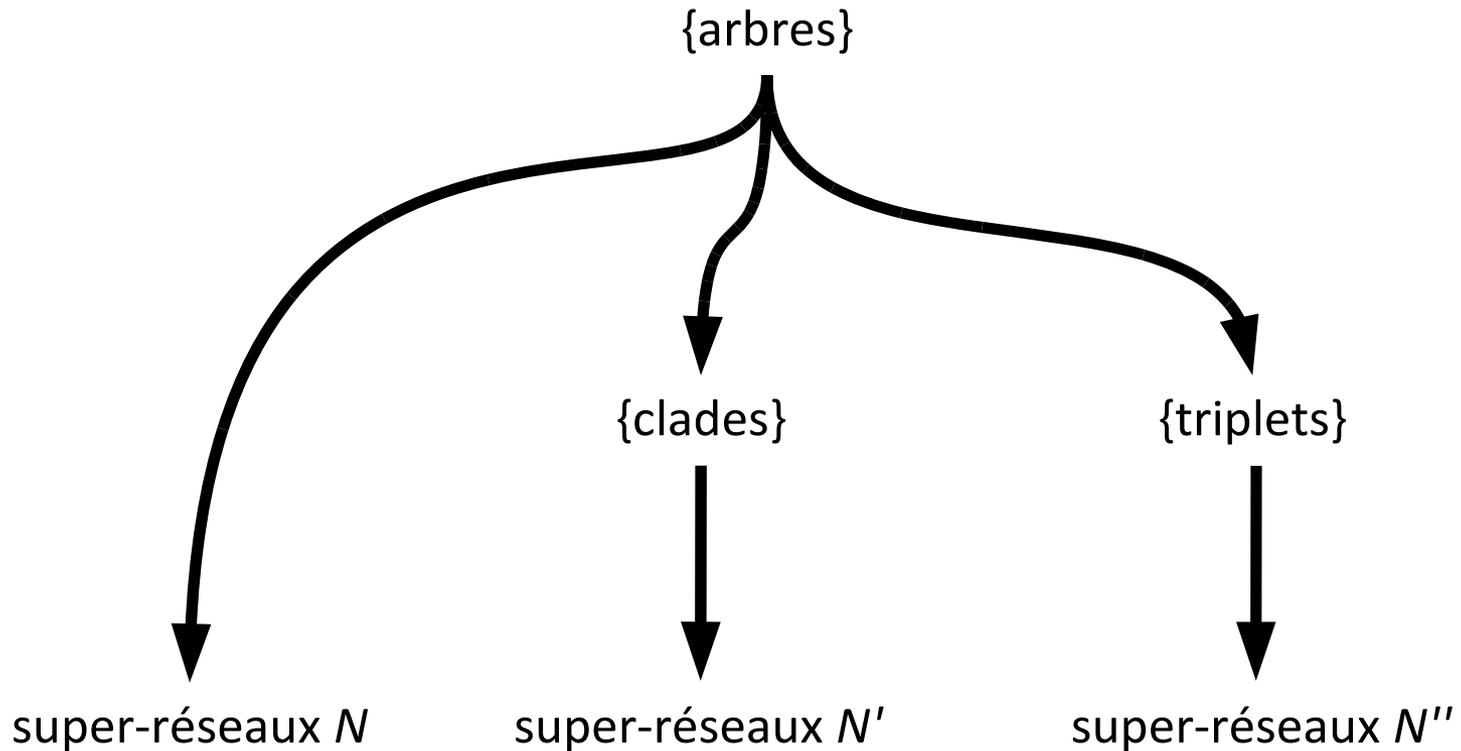


$$N = N' = N'' ?$$

Reconstruction combinatoire de réseaux phylogénétiques

Idée :

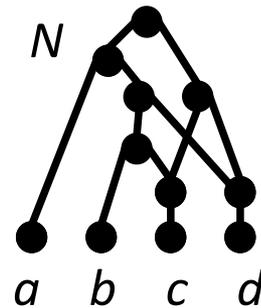
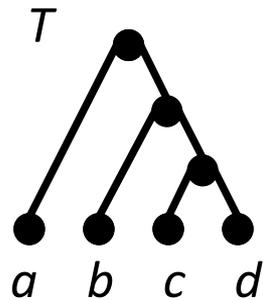
modifier le type de données à traiter



$$\{N\} \subseteq \{N'\} \subseteq \{N''\}$$

Reconstruction combinatoire de réseaux phylogénétiques

Un réseau qui contient l'ensemble de **tous les triplets ou clades d'un arbre T** ne contient **pas forcément T** .

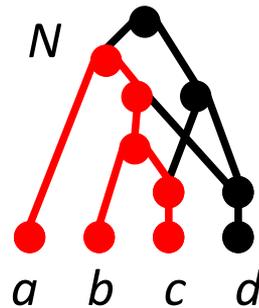
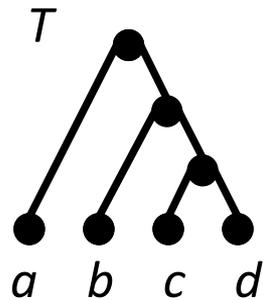


contient $\{a|bc, a|bd, a|cd, b|cd\}$
mais pas T

contient $\{abcd, bcd, cd, a, b, c, d\}$
mais pas T

Reconstruction combinatoire de réseaux phylogénétiques

Un réseau qui contient l'ensemble de **tous les triplets ou clades d'un arbre T** ne contient **pas forcément T** .

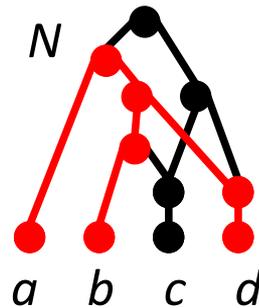
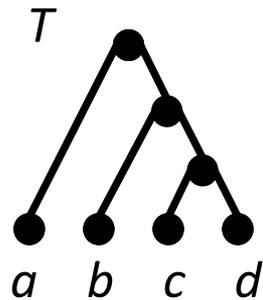


contient $\{a|bc, a|bd, a|cd, b|cd\}$
mais pas T

contient $\{abcd, bcd, cd, a, b, c, d\}$
mais pas T

Reconstruction combinatoire de réseaux phylogénétiques

Un réseau qui contient l'ensemble de **tous les triplets ou clades d'un arbre T** ne contient **pas forcément T** .

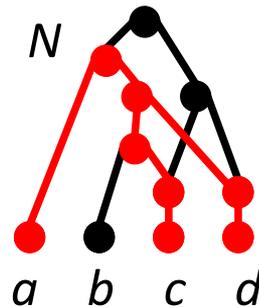
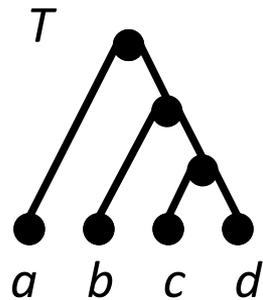


contient $\{a|bc, a|bd, a|cd, b|cd\}$
mais pas T

contient $\{abcd, bcd, cd, a, b, c, d\}$
mais pas T

Reconstruction combinatoire de réseaux phylogénétiques

Un réseau qui contient l'ensemble de **tous les triplets ou clades d'un arbre T** ne contient **pas forcément T** .

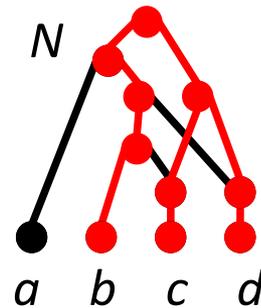
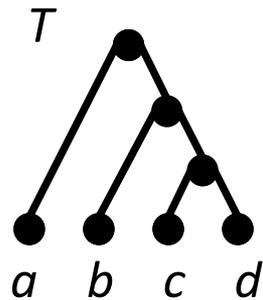


contient $\{a|bc, a|bd, a|cd, b|cd\}$
mais pas T

contient $\{abcd, bcd, cd, a, b, c, d\}$
mais pas T

Reconstruction combinatoire de réseaux phylogénétiques

Un réseau qui contient l'ensemble de **tous les triplets ou clades d'un arbre T** ne contient **pas forcément T** .

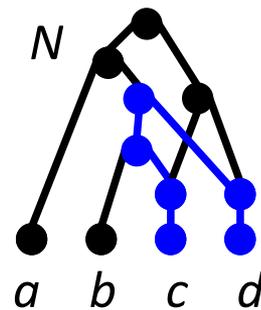
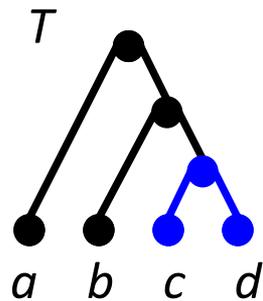


contient $\{a|bc, a|bd, a|cd, b|cd\}$
mais pas T

contient $\{abcd, bcd, cd, a, b, c, d\}$
mais pas T

Reconstruction combinatoire de réseaux phylogénétiques

Un réseau qui contient l'ensemble de **tous les triplets ou clades d'un arbre T** ne contient **pas forcément T** .

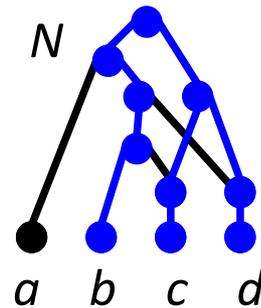
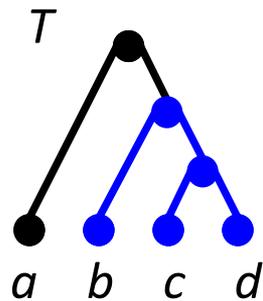


contient $\{a|bc, a|bd, a|cd, b|cd\}$
mais pas T

contient $\{abcd, bcd, cd, a, b, c, d\}$
mais pas T

Reconstruction combinatoire de réseaux phylogénétiques

Un réseau qui contient l'ensemble de **tous les triplets ou clades d'un arbre T** ne contient **pas forcément T** .

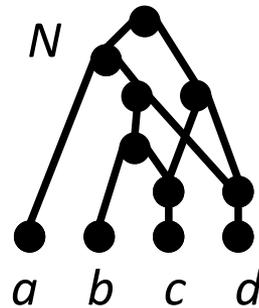
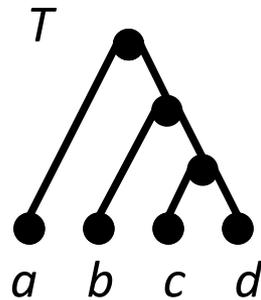


contient $\{a|bc, a|bd, a|cd, b|cd\}$
mais pas T

contient $\{abcd, bcd, cd, a, b, c, d\}$
mais pas T

Reconstruction combinatoire de réseaux phylogénétiques

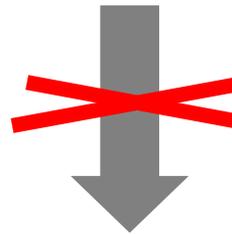
Un réseau qui contient l'ensemble de **tous les triplets ou clades d'un arbre T** ne contient **pas forcément T** .



contient $\{a|bc, a|bd, a|cd, b|cd\}$
mais pas T

contient $\{abcd, bcd, cd, a, b, c, d\}$
mais pas T

contient les clades / triplets d'un arbre T

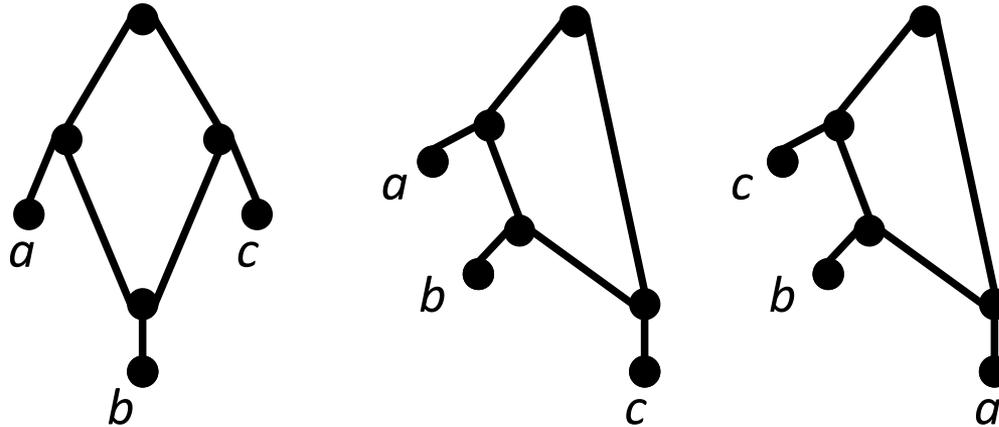


contient T .

Ambiguïté des solutions

Ambiguïté de la reconstruction, même à partir de données complètes et correctes

Plusieurs réseaux minimaux **distincts** ont exactement le **même ensemble** d'arbres, de triplets, de clades.

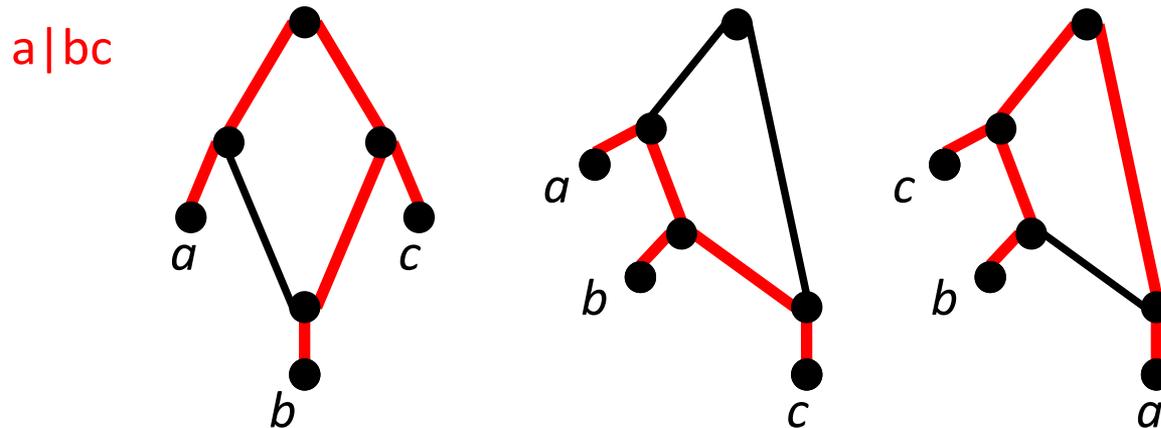


Caractérisation pour les réseaux de niveau 1 :
les seuls cas ambigus sont les blobs ci-dessus (< 5 sommets)

Ambiguïté des solutions

Ambiguïté de la reconstruction, même à partir de données complètes et correctes

Plusieurs réseaux minimaux **distincts** ont exactement le **même ensemble** d'arbres, de triplets, de clades.

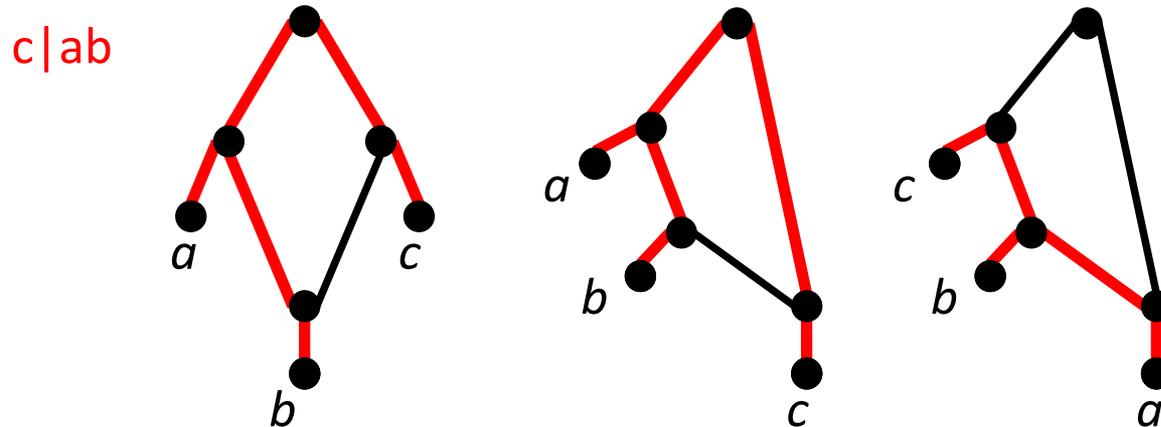


Caractérisation pour les réseaux de niveau 1 :
les seuls cas ambigus sont les blobs ci-dessus (< 5 sommets)

Ambiguïté des solutions

Ambiguïté de la reconstruction, même à partir de données complètes et correctes

Plusieurs réseaux minimaux **distincts** ont exactement le **même ensemble** d'arbres, de triplets, de clades.

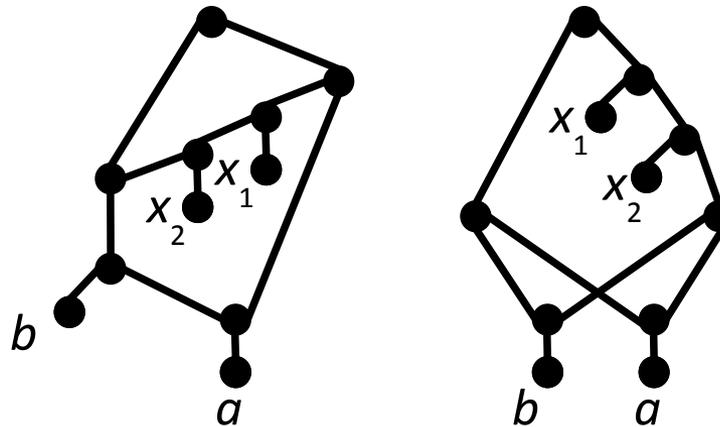


Caractérisation pour les réseaux de niveau 1 :
les seuls cas ambigus sont les blobs ci-dessus (< 5 sommets)

Ambiguïté des solutions

Ambiguïté de la reconstruction, même à partir de données complètes et correctes

Plusieurs réseaux minimaux **distincts** ont exactement le **même ensemble** d'arbres, de triplets, de clades.

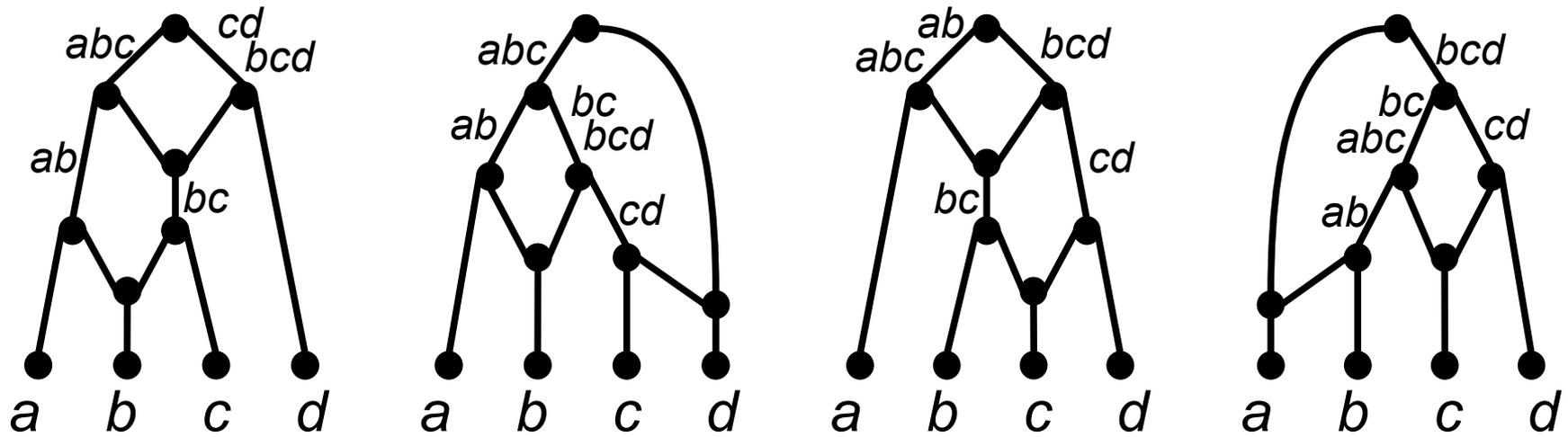


2 réseaux de niveau 2 avec le même ensemble de triplets

Ambiguïté des solutions

Ambiguïté de la reconstruction, même à partir de données complètes et correctes

Plusieurs réseaux minimaux **distincts** ont exactement le **même ensemble** d'arbres, de triplets, de clades.



4 réseaux de niveau 2 avec le même ensemble de triplets et de clades

Utilisation pratique

méthodes existantes
à faire ou améliorer

Conditions d'utilisation	Données disponibles	Traitements possibles
arbres enracinés	arbres non enracinés	<i>Enracinement à partir d'un arbre d'espèces de référence ou des contraintes topologiques</i>
arbres de gènes simple copie	“MUL-trees” (avec gènes dupliqués)	Traitement des MUL-trees Scornavacca, Berry & Ranwez, 2009
Clades et triplets corrects	données bruitées	Nettoyage d'arbre PhySIC_IST, 2008 Filtrage des données (clades avec support élevé, présent dans plus de x% des arbres) <i>Edition de données : solution contenant le maximum des données d'entrée</i>
Données complètes (ensembles denses de triplets, clades complets)	données partielles, gènes supprimés	<i>Sélection d'un grand nombre d'arbres avec un grand nombre d'espèces en commun</i> <i>Sélection du nombre maximum de taxons avec densité des triplets</i> problèmes NP-complets

Plan

- Les réseaux phylogénétiques
- Méthodes de reconstruction
- Limites des méthodes combinatoires
- **Illustration sur des données biologiques**
- Utilisation sur des données textuelles
- Perspectives

Illustrations

16 arbres sur 47 taxons de la base HOGENOM

(proteobacteria)

24 Enterobacteriales

2 Pasteurellales

1 Aeromonadales

9 Alteromonadales

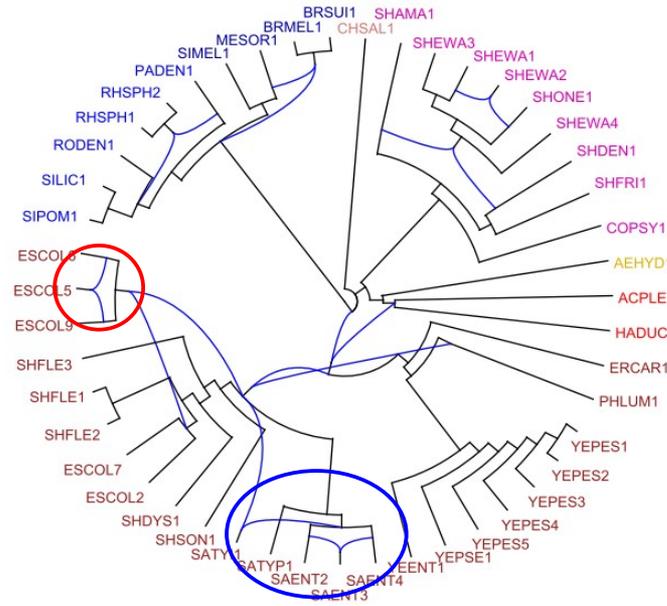
1 Oceanospirillales

6 Rhodobacterales

4 Rhizobiales

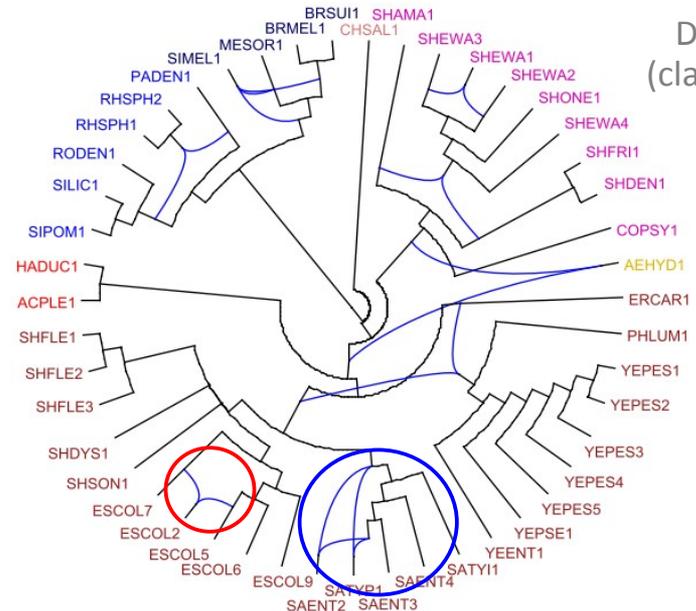
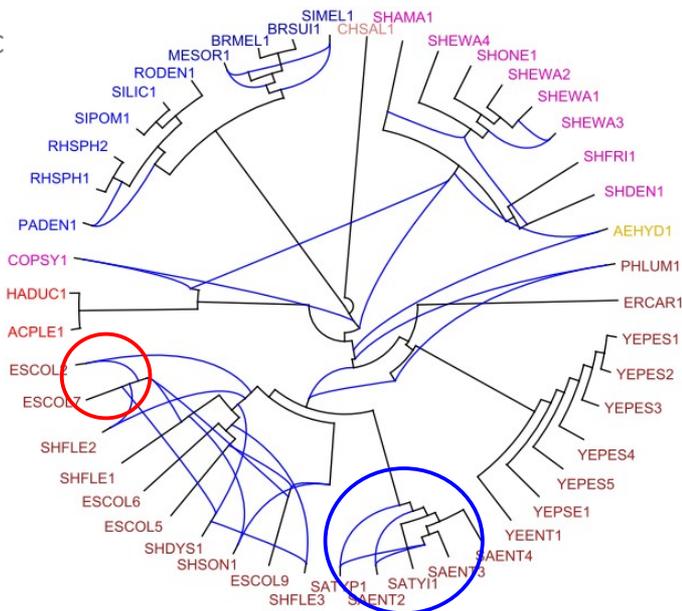


Réseaux contenant les triplets, clades souples, présents dans au moins 20% des arbres



Lev1athan
(heuristique triplets, niveau 1)
24 sec.

Simplistic
(triplets, niveau 7)
63 sec.



Dendroscope
(clades, "galled network")
<1 sec.

Illustrations

16 arbres sur 47 taxons de la base HOGENOM

(proteobacteria)

24 Enterobacteriales

2 Pasteurellales

1 Aeromonadales

9 Alteromonadales

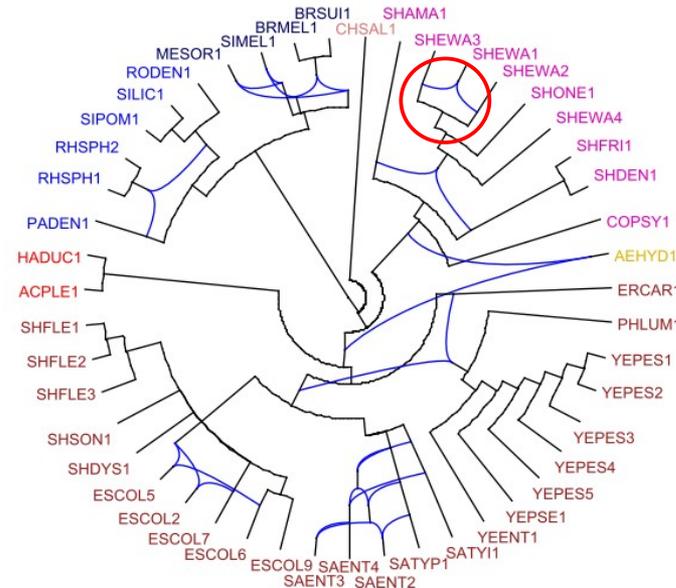
1 Oceanospirillales

6 Rhodobacterales

4 Rhizobiales

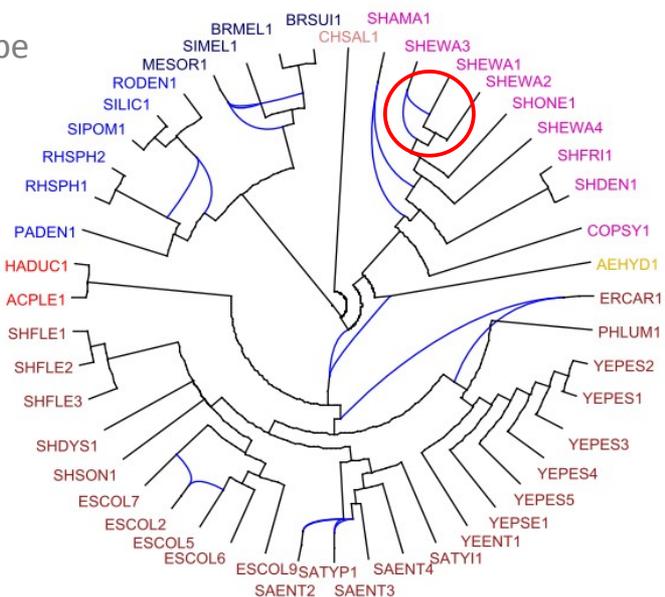


Réseaux contenant les triplets, clades souples, présents dans au moins 20% des arbres

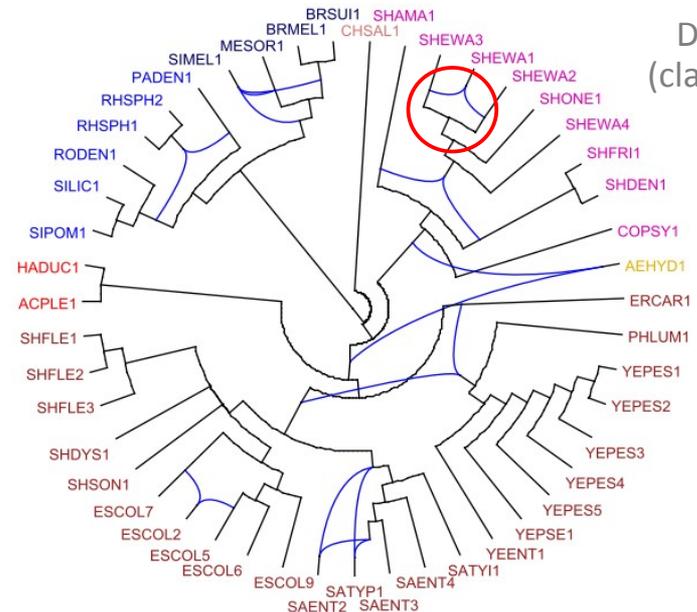


Dendroscope
(clades
stricts,
"cluster
network",
level 1)
<1 sec.

Dendroscope
(clades,
niveau 2)
2 sec.



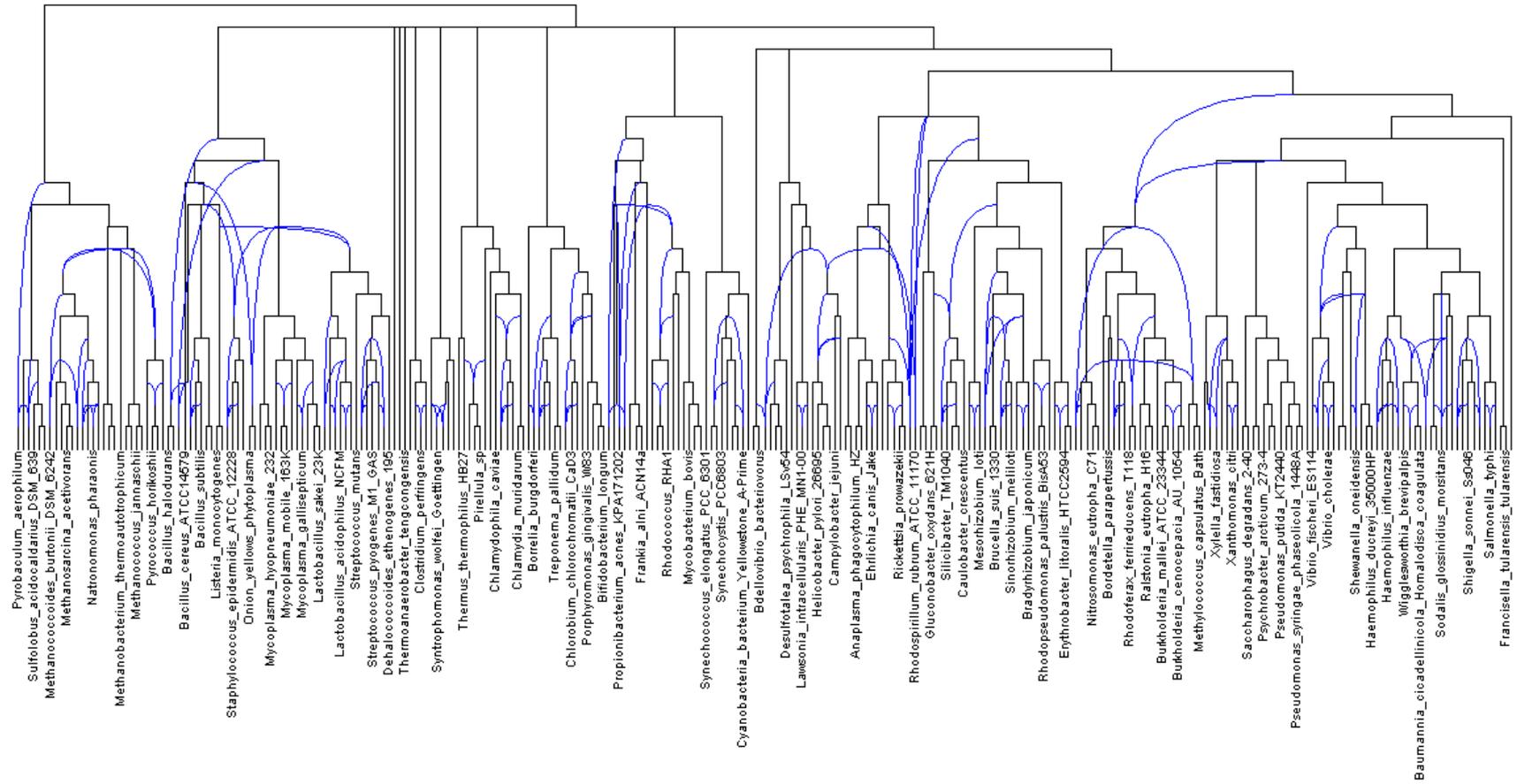
Dendroscope
(clades, "galled
network")
<1 sec.



Illustrations

9 arbres sur 279 espèces de procaryotes Clades dans au moins 2 arbres

Auch, Steigle, Huson & Henz, 2009

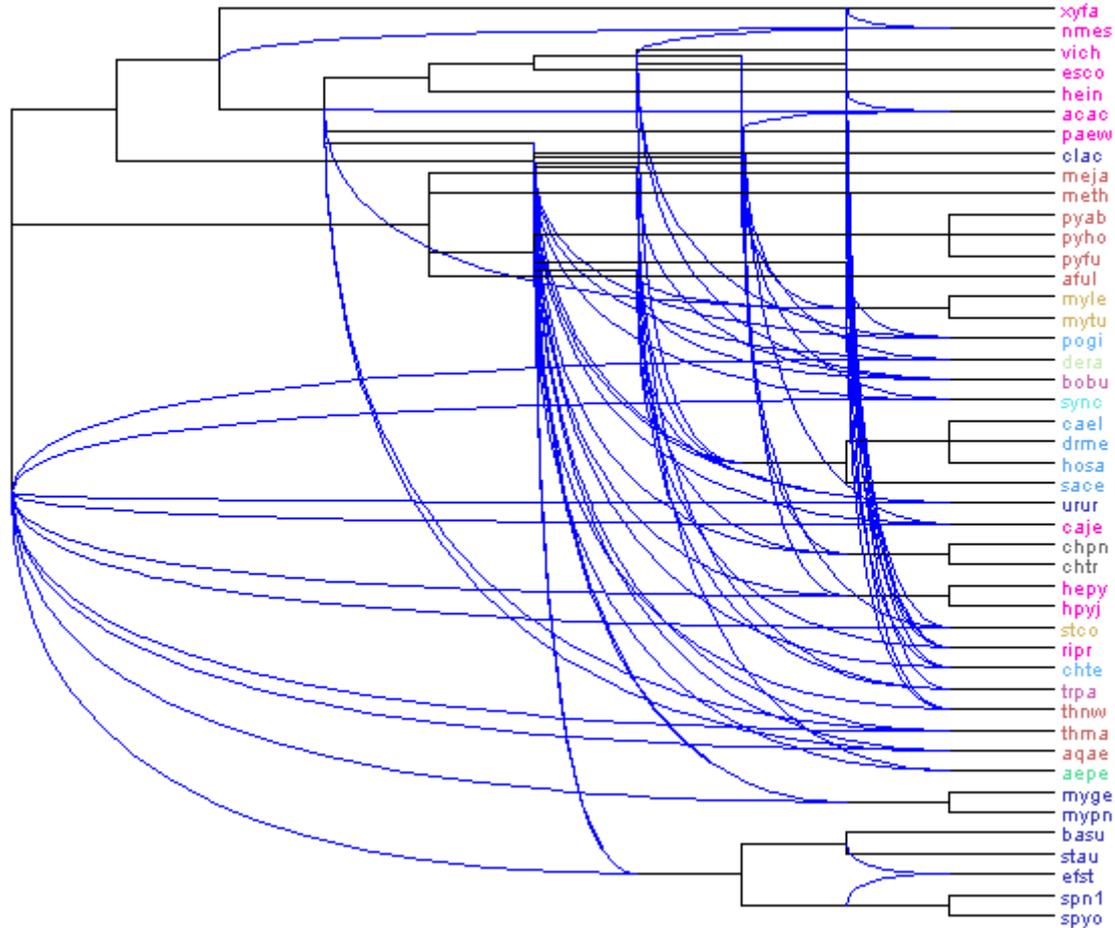


Dendroscope
(clusters, galled network)
2 sec.

Illustrations

23 arbres, 45 espèces des 3 domaines du vivant
clades avec 99% de confiance (bootstrap)

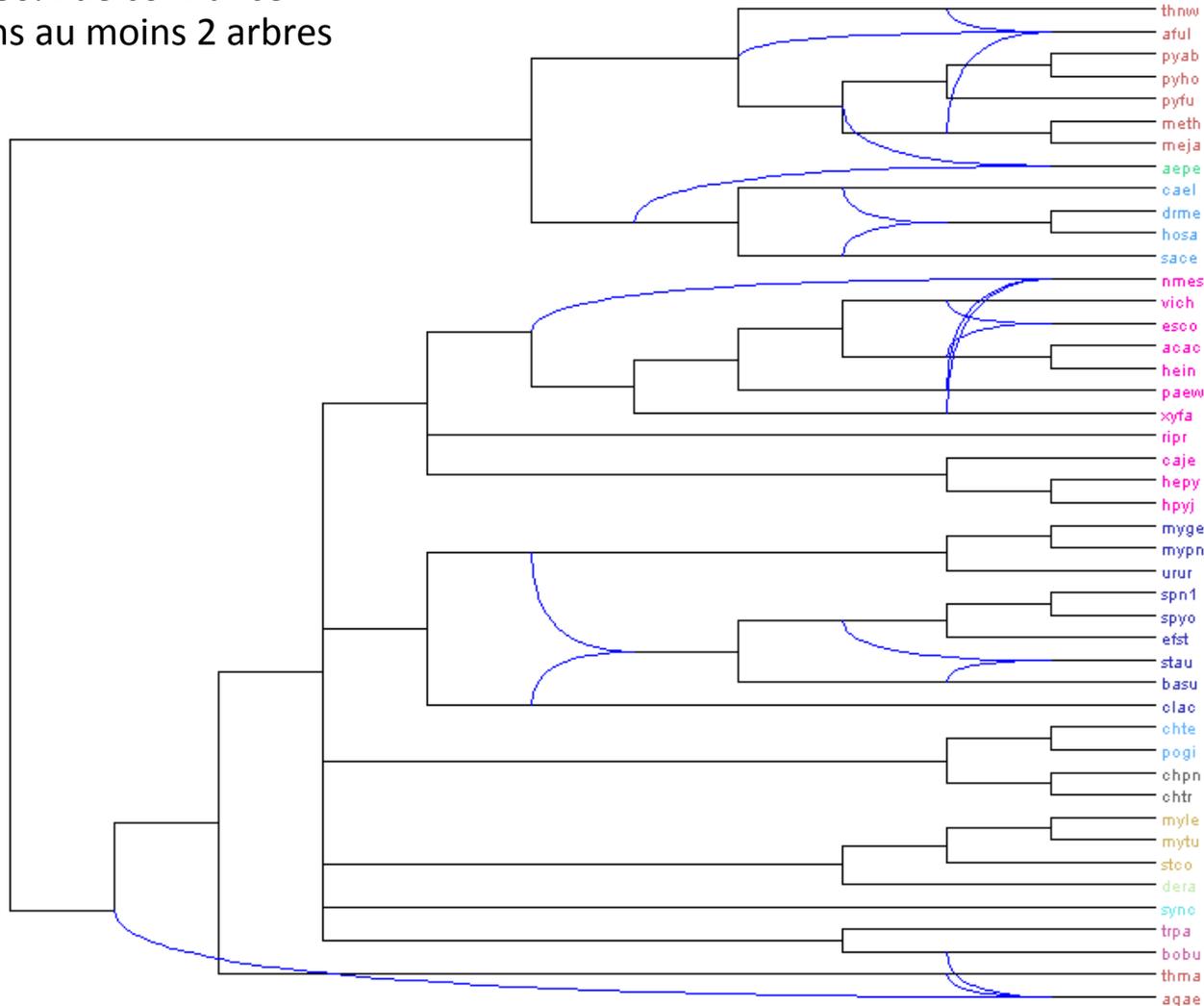
Dendroscope
("galled
network")
4 sec.



Données de Brown JR, Douady CJ, Italia MJ, Marshall WE, Stanhope MJ:
Universal trees based on large combined protein sequence data sets. Nat Genet 2001, 28:281--285

Illustrations

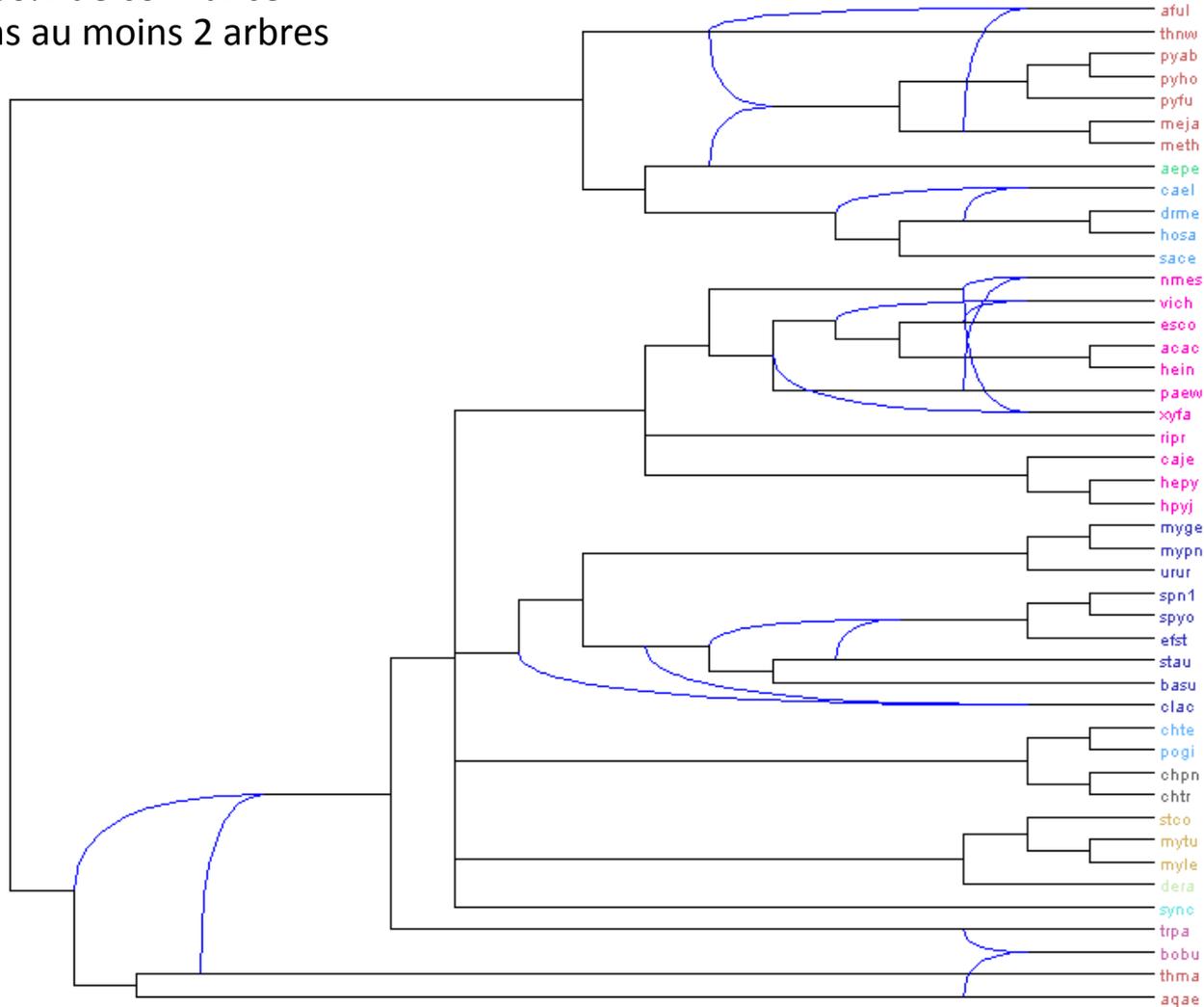
23 arbres, 45 espèces des 3 domaines du vivant
clades avec 80% de confiance
présents dans au moins 2 arbres



Dendroscope
("galled
network")
<1 sec.

Illustrations

23 arbres, 45 espèces des 3 domaines du vivant
clades avec 80% de confiance
présents dans au moins 2 arbres

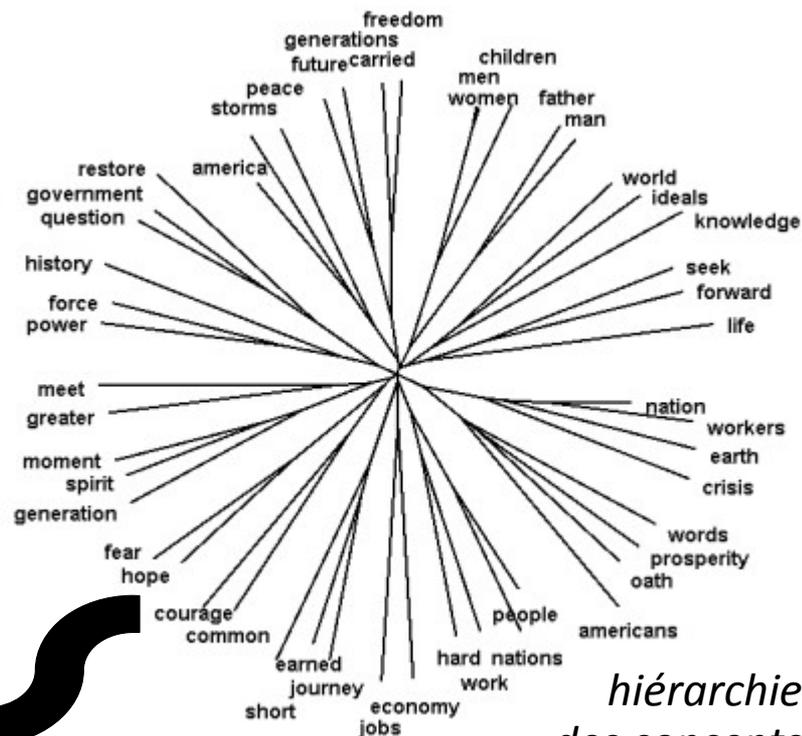
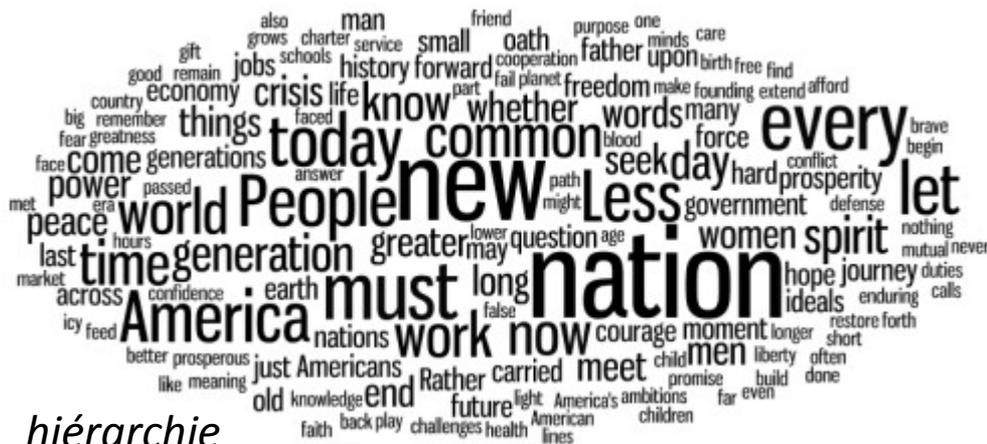


Dendroscope
(réseau de
niveau 3)
<1 sec.

Plan

- Les réseaux phylogénétiques
- Méthodes de reconstruction
- Limites des méthodes combinatoires
- Illustration sur des données biologiques
- **Utilisation sur des données textuelles**
- Perspectives

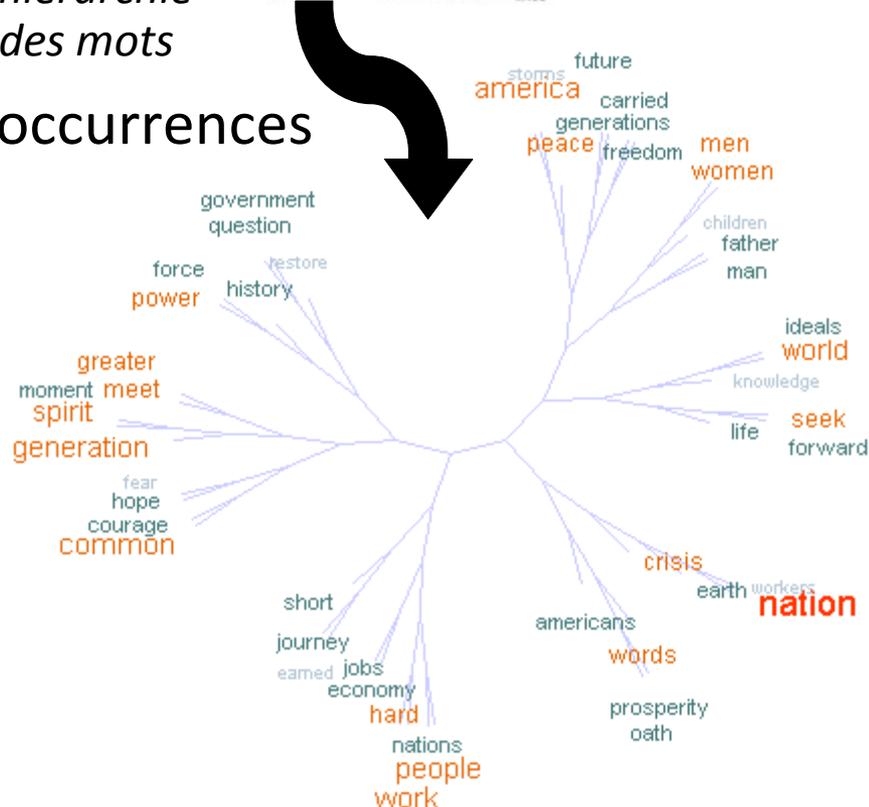
Nuage arboré, une information double



hiérarchie
des mots
occurrences

cooccurrences

hiérarchie
des concepts



Nuages de tags

- Construits depuis un ensemble de tags
- Taille de police liée à la fréquence



Ce qui est habituellement cité comme le premier nuage de tags, dans *Microserfs* de D. Coupland, HarperCollins, Toronto, 1995

Extraire l'information sémantique d'un texte

- analyse arborée Brunet (Hyperbase)
- graphe de cooccurrence Brunet (Hyperbase)
- graphe sémantique Grimmer (Wordmapper)
- lexicogramme récursif Martinez (Coocs)
- désambiguïsation lexicale Véronis (Hyperlex)
- réseau Phrasenet Viegas et al. (IBM Many Eyes)
- projection géodésique Viprey (Astartex)

Extraire l'information sémantique d'un texte

- analyse arborée
- graphe de cooccurrence
- graphe sémantique
- lexicogramme récursif
- désambiguïstation lexicale
- réseau Phrasenet
- projection géodésique

Brunet (Hyperbase)

Brunet (Hyperbase)

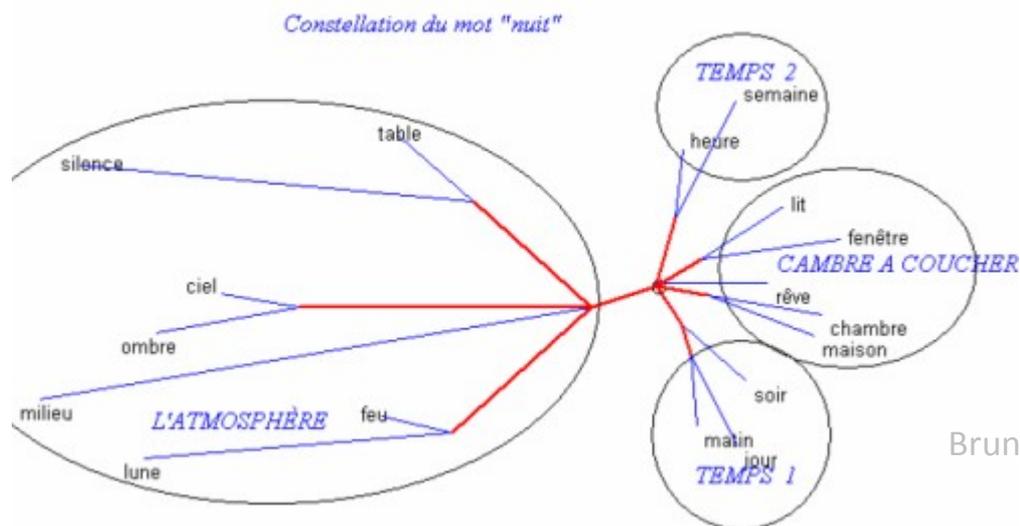
Grimmer (Wordmapper)

Martinez (Coocs)

Véronis (Hyperlex)

Viegas et al. (IBM Many Eyes)

Viprey (Astartex)



Brunet, *Les séquences (suite)*,
JADT'08

Extraire l'information sémantique d'un texte

- analyse arborée
- graphe de cooccurrence
- graphe sémantique
- lexicogramme récursif
- désambiguïsation lexicale
- réseau Phrasenet
- projection géodésique

Brunet (Hyperbase)

Brunet (Hyperbase)

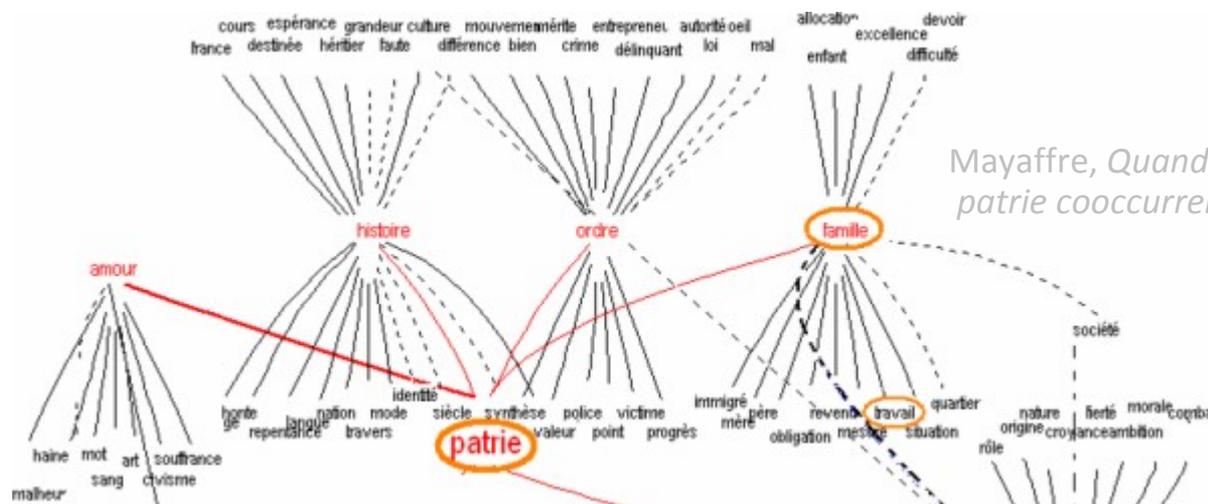
Grimmer (Wordmapper)

Martinez (Coocs)

Véronis (Hyperlex)

Viegas et al. (IBM Many Eyes)

Viprey (Astartex)



Mayaffre, Quand travail, famille, et patrie cooccurrent dans le discours de Nicolas Sarkozy, JADT'08

Extraire l'information sémantique d'un texte

- analyse arborée
- graphe de cooccurrence
- graphe sémantique
- lexicogramme récursif
- désambiguïstation lexicale
- réseau Phrasenet
- projection géodésique

Brunet (Hyperbase)

Brunet (Hyperbase)

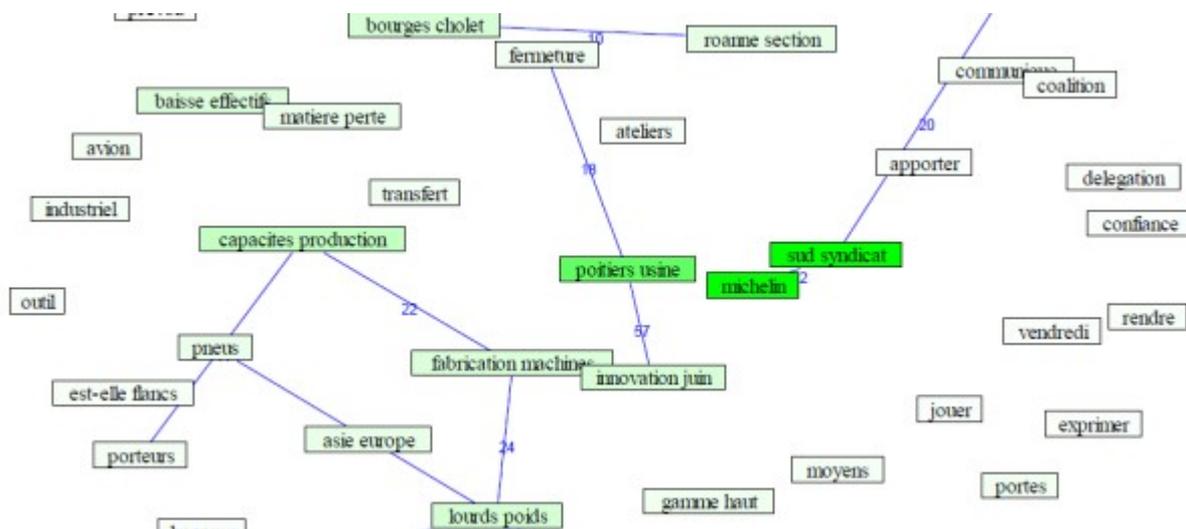
Grimmer (Wordmapper)

Martinez (Coocs)

Véronis (Hyperlex)

Viegas et al. (IBM Many Eyes)

Viprey (Astartex)



Peyrat-Guillard,
*Analyse du discours
syndical sur l'entreprise, JADT'08*

Extraire l'information sémantique d'un texte

- analyse arborée
- graphe de cooccurrence
- graphe sémantique
- lexicogramme récursif
- désambiguïstation lexicale
- réseau Phrasenet
- projection géodésique

Brunet (Hyperbase)

Brunet (Hyperbase)

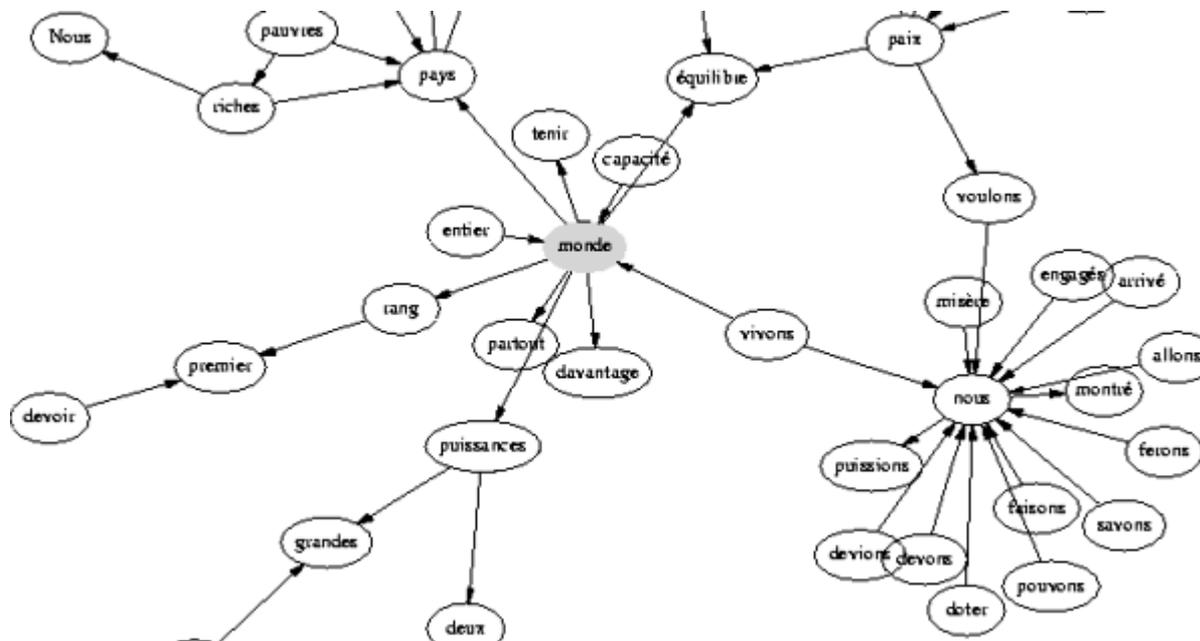
Grimmer (Wordmapper)

Martinez (Coocs)

Véronis (Hyperlex)

Viegas et al. (IBM Many Eyes)

Viprey (Astartex)



Leblanc, Martinez
*L'analyse contrastive des
réseaux de cooccurrence*
JADT 2006.

Extraire l'information sémantique d'un texte

- analyse arborée
- graphe de cooccurrence
- graphe sémantique
- lexicogramme récursif
- désambiguïstation lexicale
- réseau Phrasenet
- projection géodésique

Brunet (Hyperbase)

Brunet (Hyperbase)

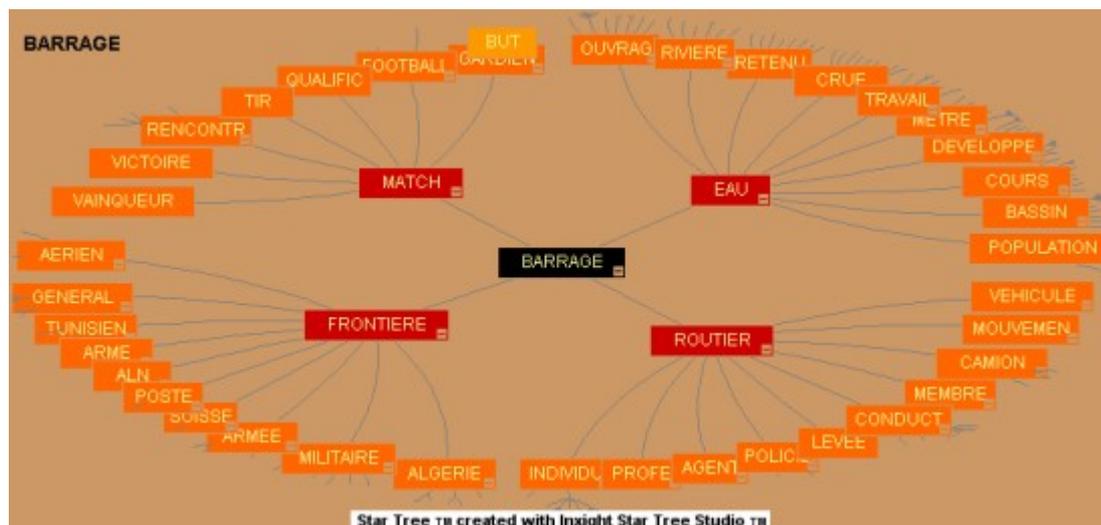
Grimmer (Wordmapper)

Martinez (Coocs)

Véronis (Hyperlex)

Viegas et al. (IBM Many Eyes)

Viprey (Astartex)



Désambiguïstation du mot
"barrage".

Véronis, *HyperLex: Lexical Cartography for Information Retrieval*, 2004

Extraire l'information sémantique d'un texte

- analyse arborée
- graphe de cooccurrence
- graphe sémantique
- lexicogramme récursif
- désambiguïstation lexicale
- réseau Phrasenet
- projection géodésique

Brunet (Hyperbase)

Brunet (Hyperbase)

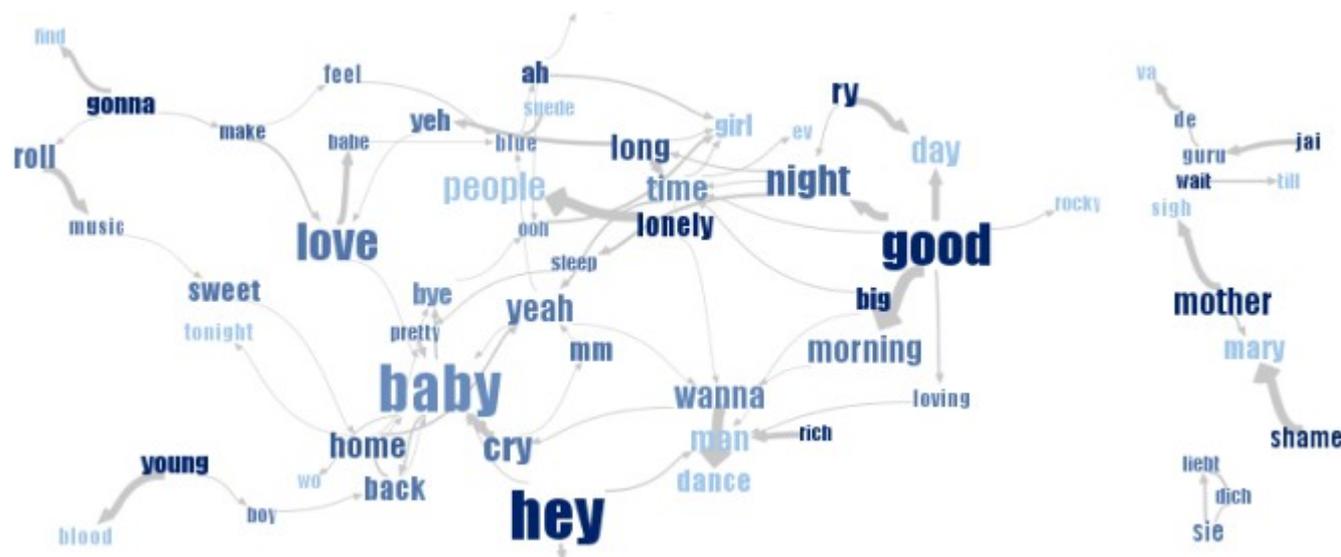
Grimmer (Wordmapper)

Martinez (Coocs)

Véronis (Hyperlex)

Viegas et al. (IBM Many Eyes)

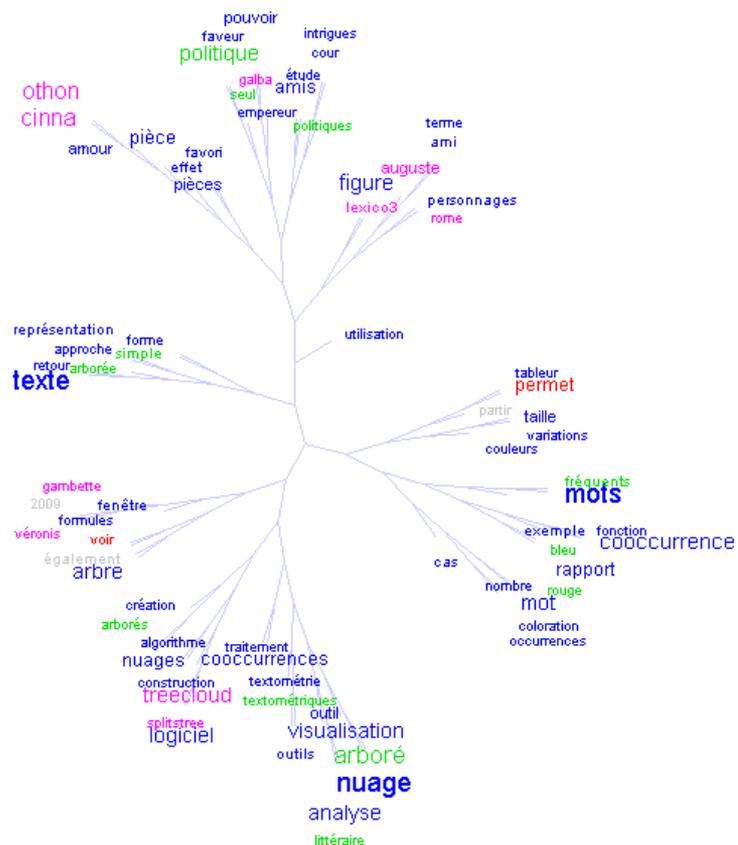
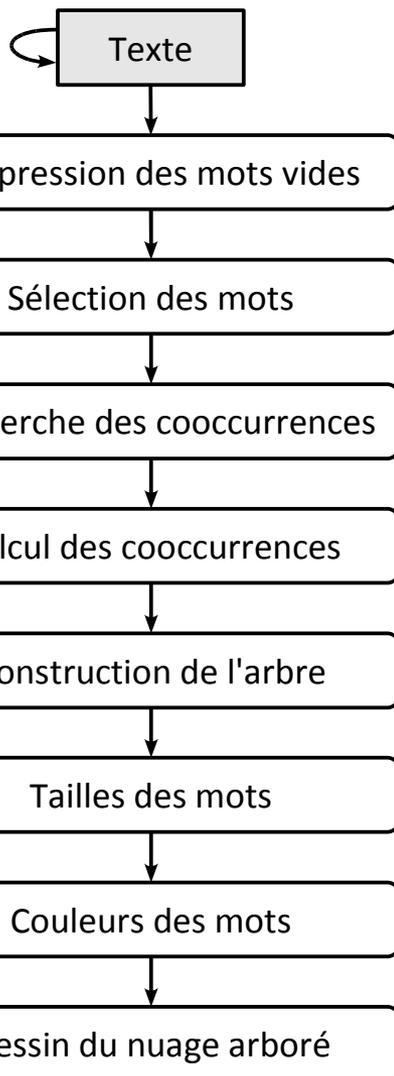
Viprey (Astartex)



Visualisation PhraseNet
de paroles des Beatles
créé avec Many Eyes (IBM)
<http://many-eyes.com>

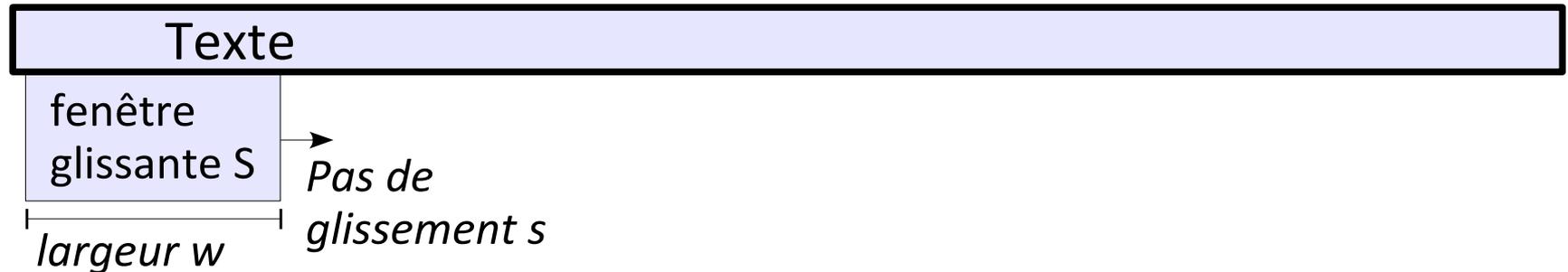
Construction des nuages arborés

Concordance d'un mot, lemmatisation
ou remplacements divers...



Construction de la matrice de cooccurrences

Formules de distance sémantique à base de cooccurrence



matrices de cooccurrence

$O_{11}, O_{12}, O_{21}, O_{22}$

	$v \in S$	$v \notin S$
$u \in S$	O_{11}	O_{12}
$u \notin S$	O_{21}	O_{22}



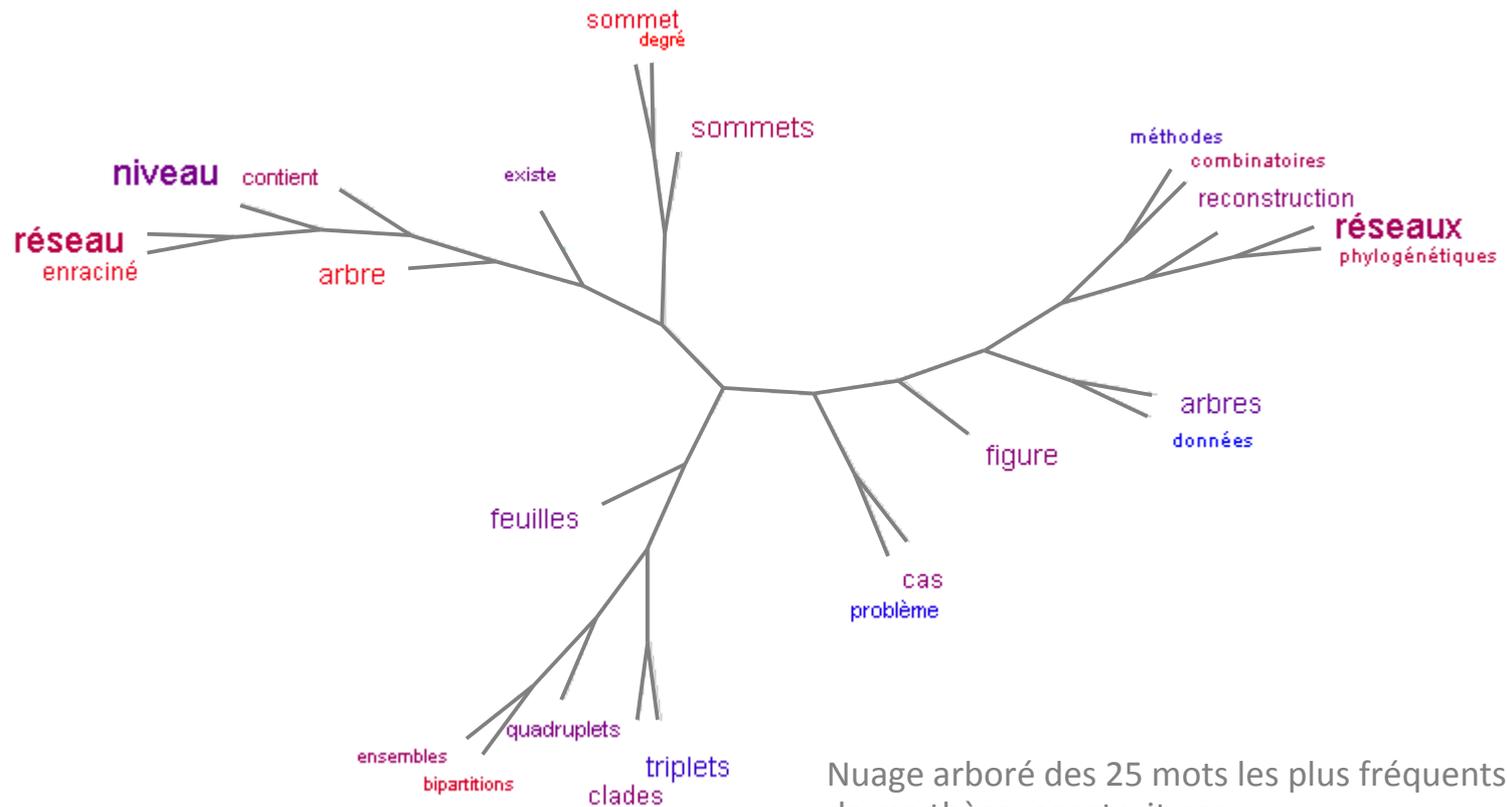
matrice de dissimilarité sémantique

chi squared, mutual information, liddel, dice, jaccard, gmean, hyperlex, minimum sensitivity, odds ratio, zscore, log likelihood, poisson-stirling...

Evert, *Statistics of words cooccurrences*, thèse, 2005
Gambette, *User manual for TreeCloud*, 2009

Construction du réseau phylogénétique

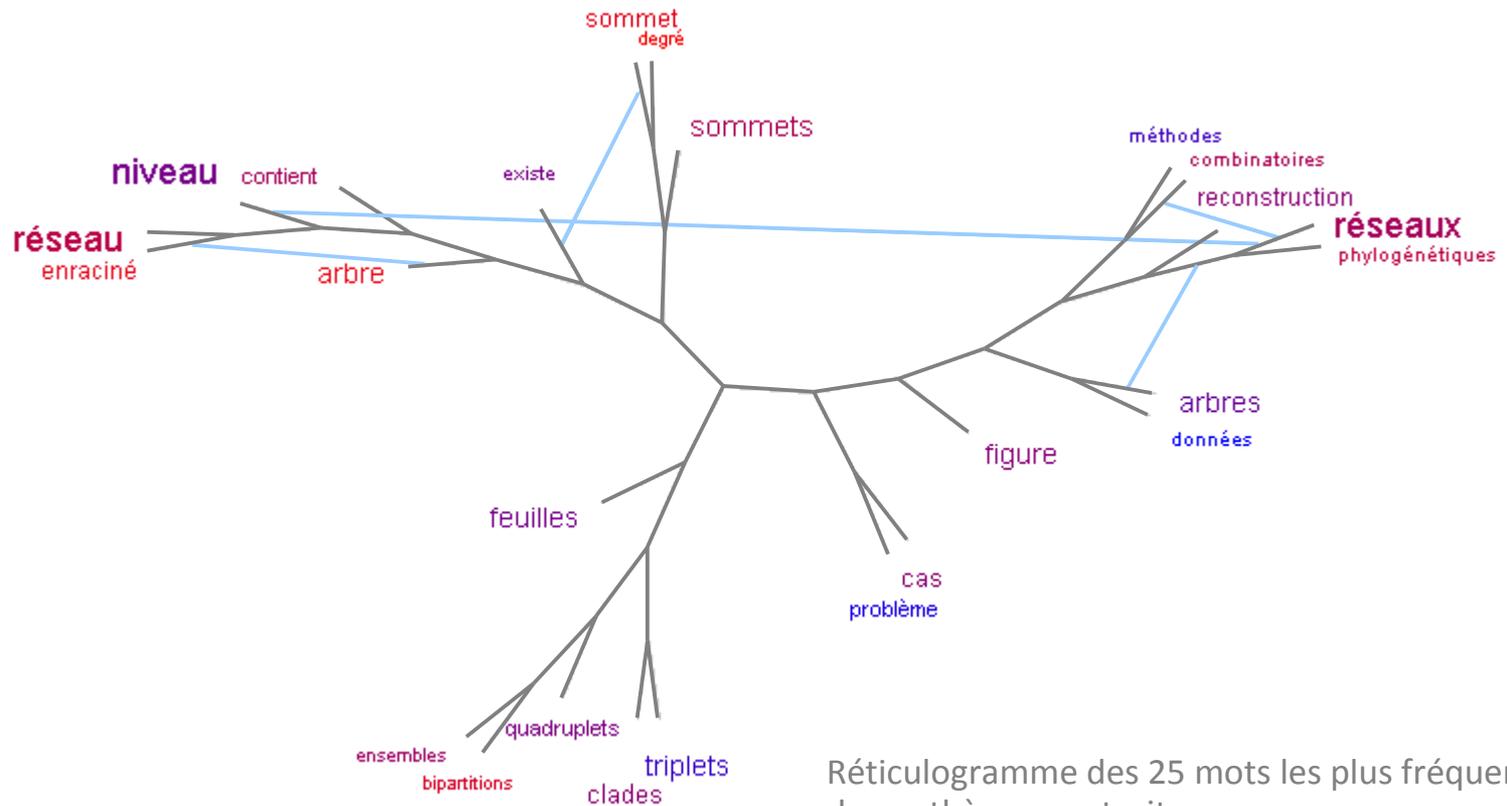
Méthode T-Rex bien adaptée



Nuage arboré des 25 mots les plus fréquents de ma thèse, construit par  TreeCloud,  SplitsTree et  T-Rex
Coloration : rouge au début, bleu à la fin

Construction du réseau phylogénétique

Méthode T-Rex bien adaptée



Réticulogramme des 25 mots les plus fréquents de ma thèse, construit par  TreeCloud,  SplitsTree et  T-Rex
Coloration : rouge au début, bleu à la fin

Plan

- Les réseaux phylogénétiques
- Méthodes de reconstruction
- Limites des méthodes combinatoires
- Illustration sur des données biologiques
- Utilisation sur des données textuelles
- Perspectives

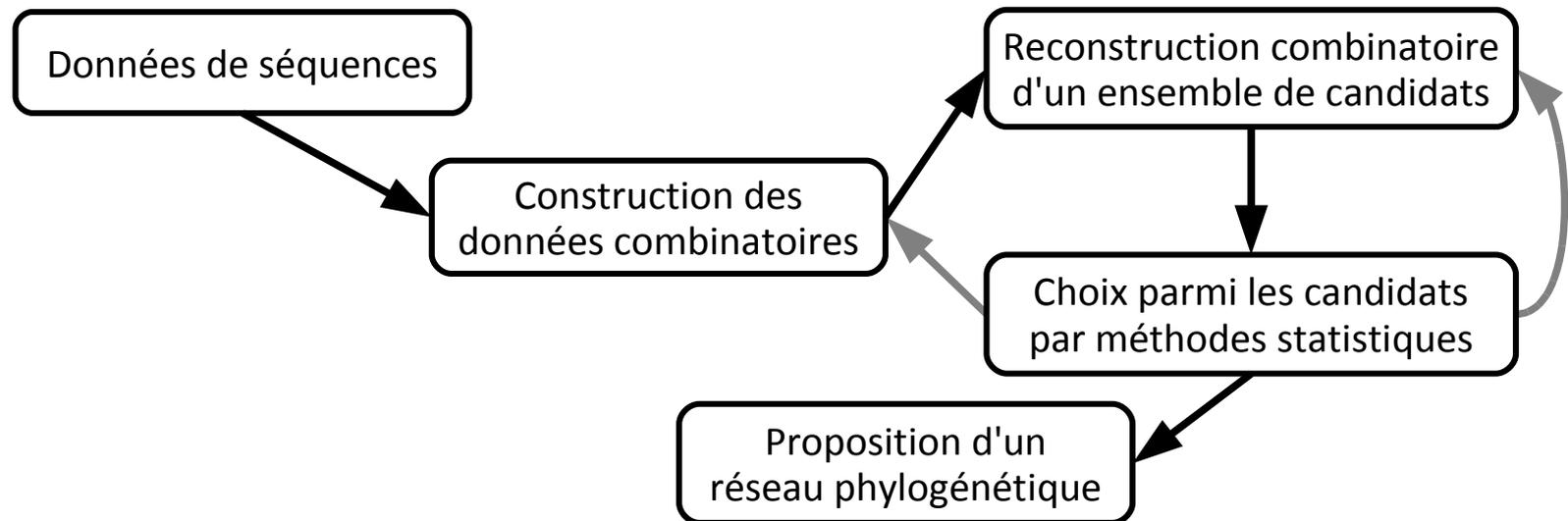
Perspectives de recherche

Combinatoire :

- Meilleure connaissance des réseaux de faible niveau, enracinés ou non : dénombrement, caractérisations...
- Mise à jour ou modification d'un réseau face à de nouvelles données

Bioinformatique :

- Fonction des gènes transférés (“autoroutes de transfert”)
- Intégration des méthodes combinatoires dans une approche statistique



Autres applications des réseaux phylogénétiques :

- visualiser la polysémie dans les nuages arborés

Merci !

Who is who in Phylogenetic Networks



<http://atgc.lirmm.fr/phylnet>

Dendroscope



<http://www.dendroscope.org>

TreeCloud



<http://www.treecloud.org>

SplitsTree



<http://www.splitstree.org>