

A short (?) introduction to phylogenetic networks

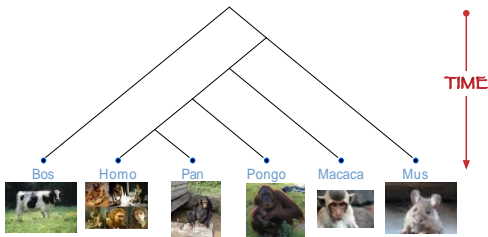
CÉLINE SCORNAVACCA

ISE-M, Equipe Phylogénie & Evolution Moléculaires
Montpellier, France

Rooted species trees ...

... are oriented connected and acyclic graphs, where terminal nodes are associated to a set of species:

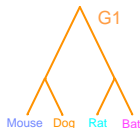
- the leaves or **taxa** represent extant organisms
- internal nodes represent hypothetical ancestors
- each **internal node** represents the lowest common ancestor of all taxa below it (**clade**)
- the only node without ancestor is called **root**



Gene trees

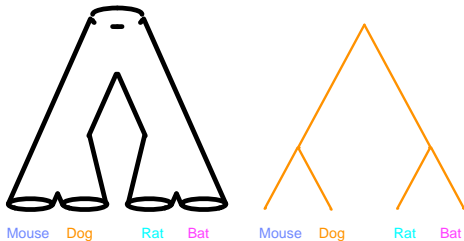
- Gene trees are built by analyzing a **gene family**, i.e., **homologous** molecular sequences appearing in the genome of different organisms.

Mouse	G	G	A	G	C	T	T	G	A	G	C	C	G	G	A	A	T	A	G	T	A	G	G	A	A	C	A	T	C	T	T	A	A	G	A	A	T	T	T	A	A	T	T	C	G	A	G	C	
Dog	G	G	A	A	T	C	T	G	A	A	C	A	G	G	C	T	T	A	G	T	A	G	C	C	A	C	T	A	G	A	A	T	A	A	G	A	C	T	T	T	A	A	T	T	C	G	A	G	C
Bat	G	G	A	A	T	T	T	G	A	A	C	A	G	G	T	T	T	A	G	T	A	G	C	C	A	C	T	A	G	A	A	T	A	A	G	A	C	T	C	T	A	A	T	T	C	G	A	G	C
Rat	G	G	A	A	T	T	T	G	A	A	C	C	G	G	C	T	C	G	T	A	G	C	A	A	C	A	A	G	A	A	T	A	A	G	C	T	T	A	A	T	T	C	G	T	G	C			



Gene trees

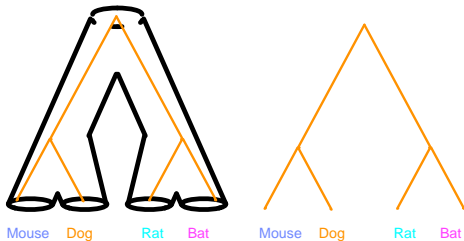
- Gene trees are built by analyzing a gene family.



- Used, among other things, to estimate species trees.

Gene trees

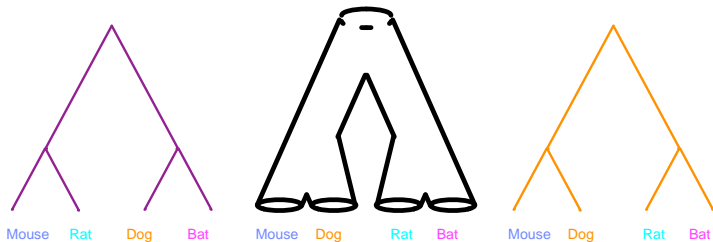
- Gene trees are built by analyzing a gene family.



- Used, among other things, to estimate species trees.

Gene trees

- Gene trees are built by analyzing a gene family.



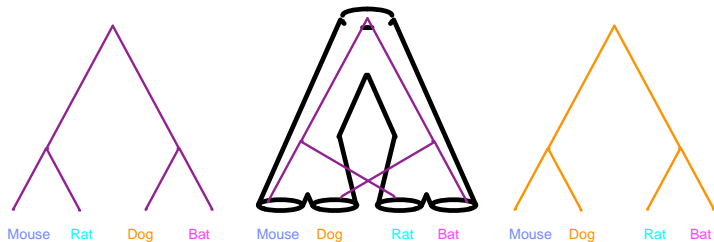
- Used, among other things, to estimate species trees.

Gene trees can significantly differ from the species tree for:

- methodological reasons
- biological reasons

Gene trees

- Gene trees are built by analyzing a gene family.



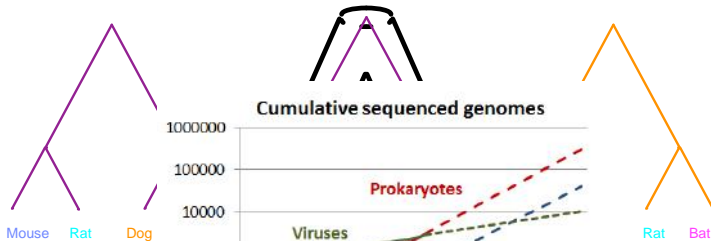
- Used, among other things, to estimate species trees.

Gene trees can significantly differ from the species tree for:

- methodological reasons
- biological reasons
- We usually use several gene families...

Gene trees

- Gene trees are built by analyzing a gene family.



- Used, among o

Gene trees can sig

- methodological
- biological reasons

ee for:

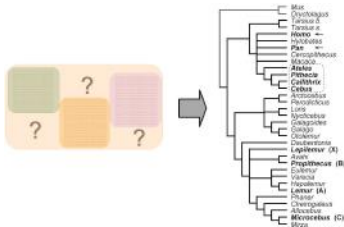
- We usually use several gene families...

<http://sulab.org/2013/06/sequenced-genomes-per-year/>

Reconstruction of phylogenies for multiple datasets

The two main *classic* approaches:

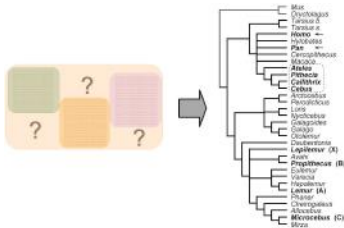
- Supermatrix approach: assembling primary data



Reconstruction of phylogenies for multiple datasets

The two main *classic* approaches:

- Supermatrix approach: **assembling primary data**

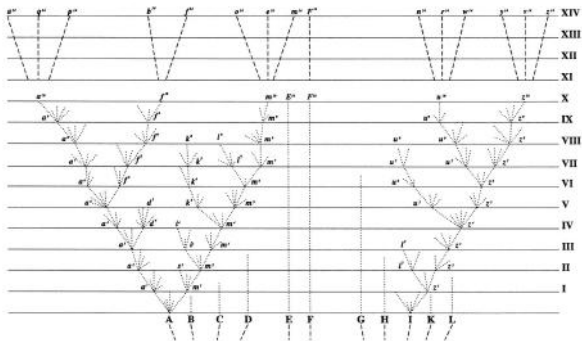


- Supertree approach: **assembling trees**



An implicit assumption

The implicit assumption of using trees is that, at a macroevolutionary scale, each (current or extinct) species or gene only descends from one ancestor. Darwin described evolution as "descent with modification", a phrase that does not necessarily imply a tree representation...

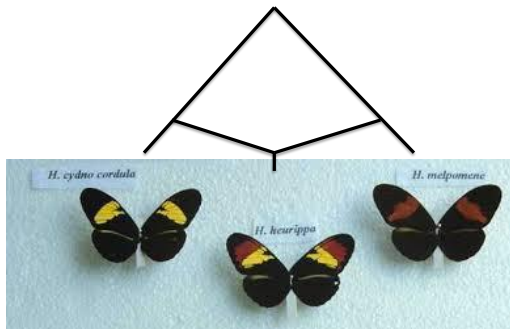


A *new* approach: building phylogenic networks

Why do we need them? Due to reticulate evolutionary phenomena (hybridization, recombination, horizontal gene transfer) the evolution of a set of species sometimes cannot be described using phylogenetic trees.

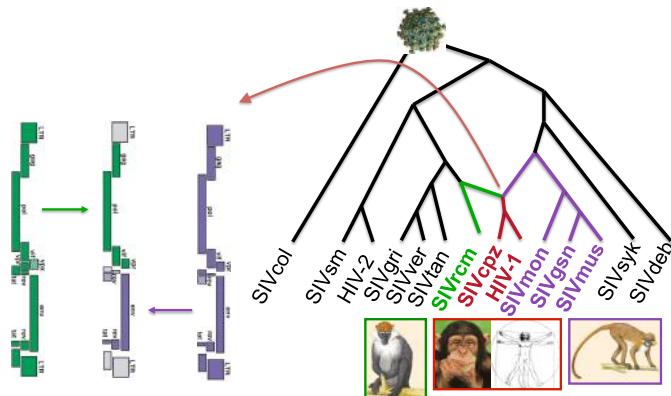
A new approach: building phylogenetic networks

Why do we need them? Due to reticulate evolutionary phenomena (hybridization, recombination, horizontal gene transfer) the evolution of a set of species sometimes cannot be described using phylogenetic trees.



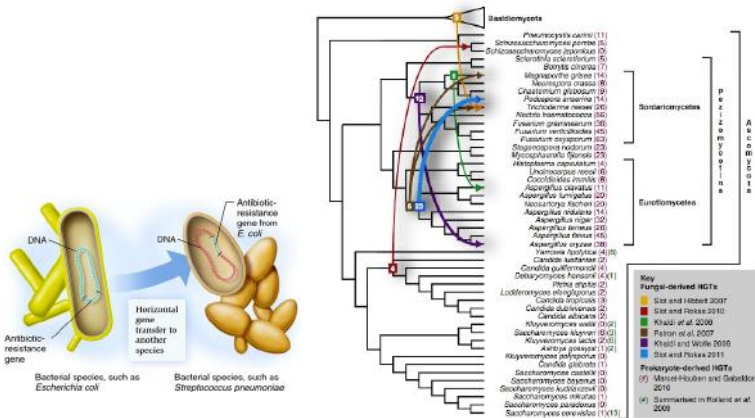
A new approach: building phylogenetic networks

Why do we need them? Due to reticulate evolutionary phenomena (hybridization, recombination, horizontal gene transfer) the evolution of a set of species sometimes cannot be described using phylogenetic trees.



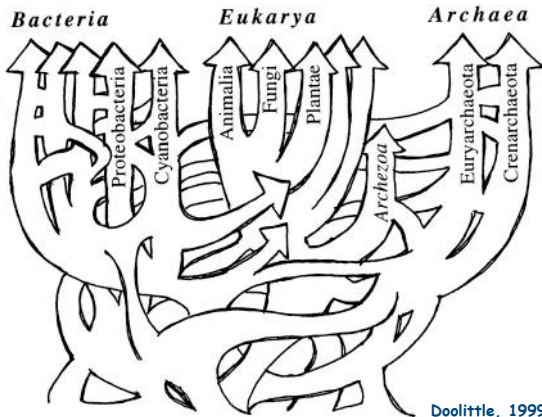
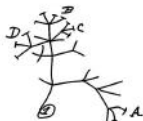
A new approach: building phylogenetic networks

Why do we need them? Due to reticulate evolutionary phenomena (hybridization, recombination, horizontal gene transfer) the evolution of a set of species sometimes cannot be described using phylogenetic trees.



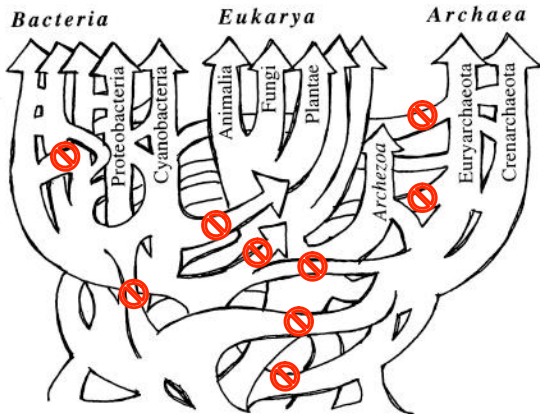
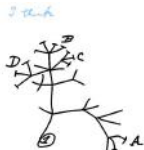
The network of life

3-10-11



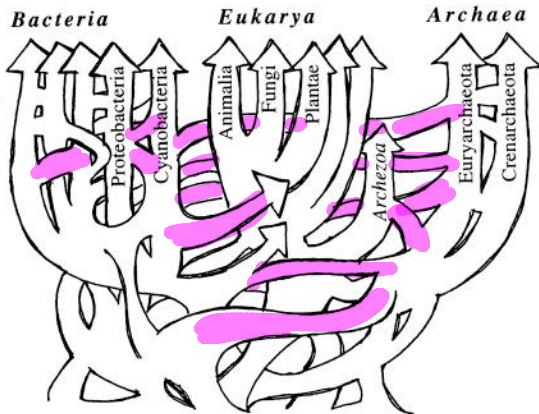
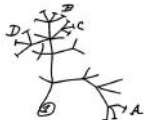
Three different paradigms

We (want to) see only the tree



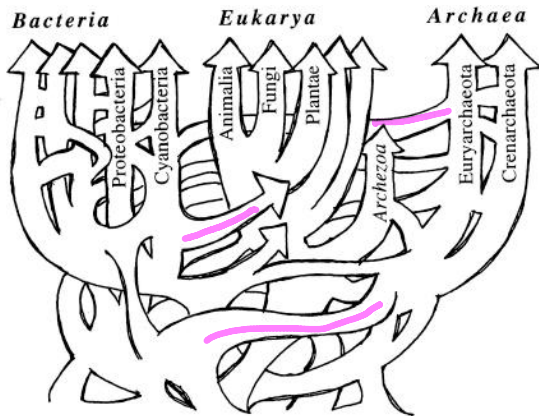
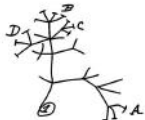
Three different paradigms

It is a big mess, no chance to retrieve the past

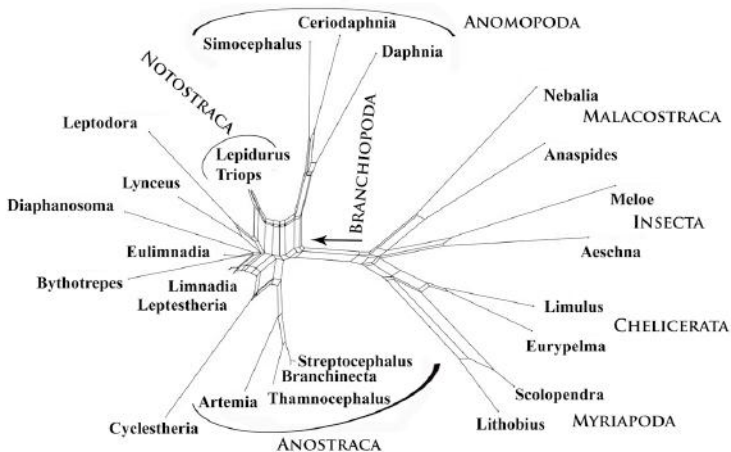


Three different paradigms

There is an underlying tree structure, with some reticulate events

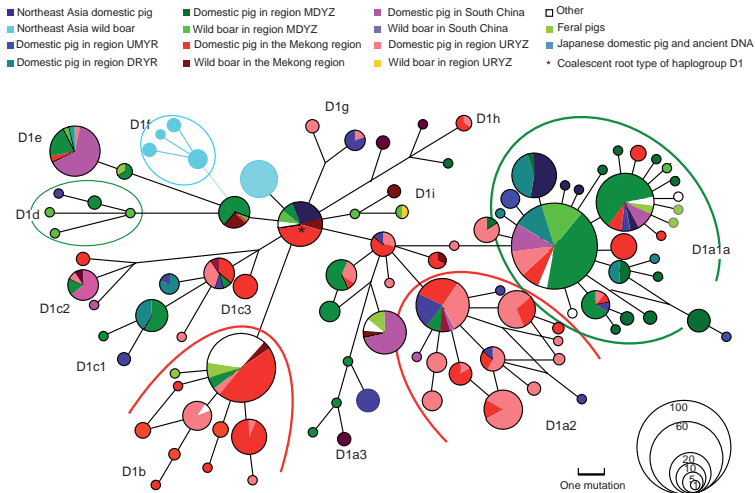


An example - a split network



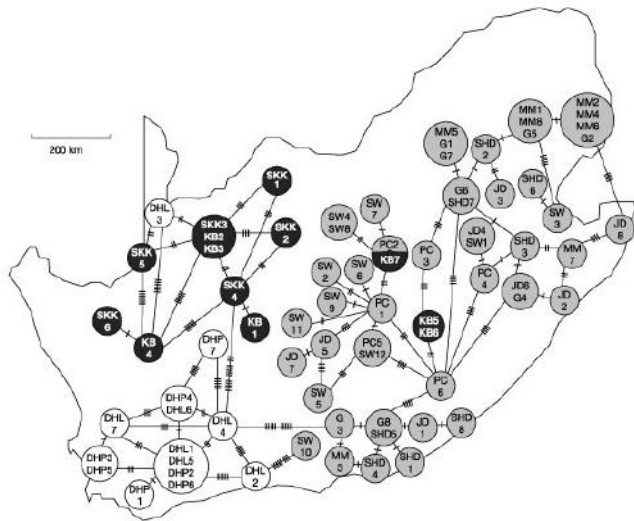
J. Wagele and C. Mayer. Visualizing differences in phylogenetic information content of alignments and distinction of three classes of long-branch effects. BMC Evolutionary Biology, 7(1):147, 2007

An example - a reduced median network



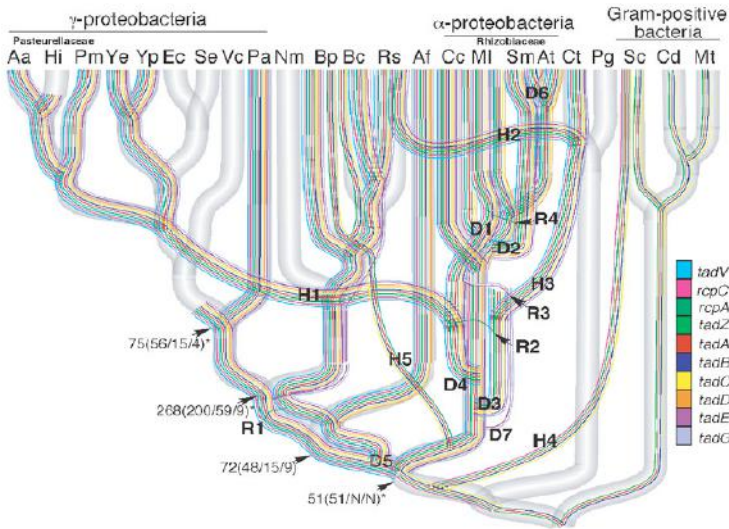
G.-S. Wu, Y.-G. Yao, K.-X. Qu, Z.-L. Ding, H. Li, M. Palanichamy, Z.-Y. Duan, N. Li, Y.-S. Chen, and Y.-P. Zhang. Population phylogenomic analysis of mitochondrial DNA in wild boars and domestic pigs revealed multiple domestication events in East Asia. *Genome Biology*, 8(11):R245, 2007

An example - a minimum spanning network



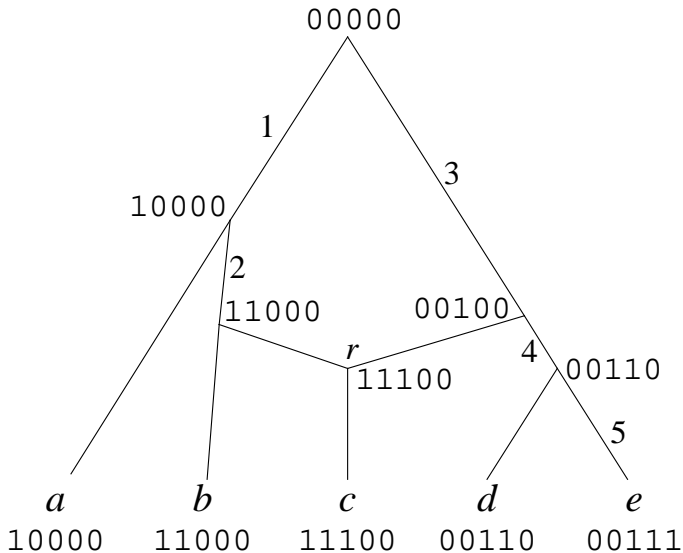
C. M. Miller-Butterworth, D. S. Jacobs, and E. H. Harley. Strong population sub- structure is correlated with morphology and ecology in a migratory bat. *Nature*, 424(6945):187-191, 2003

An example - a DTLR network

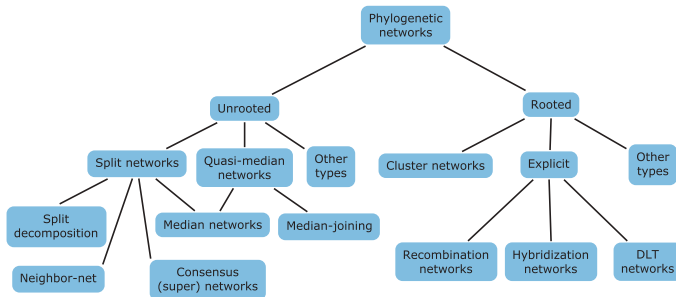
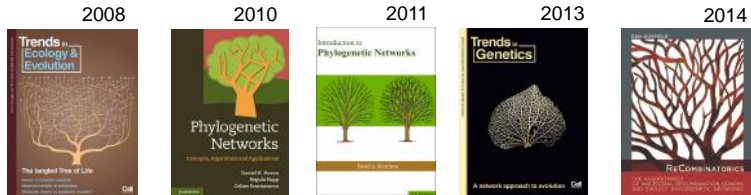


P.J. Planet, S.C. Kachlany, D.H. Fine, R. DeSalle, and D.H. Figurski. The wide spread colonization island of *actinobacillus actinomycetemcomitans*. *Nature Genetics*, 34:193–198, 2003.

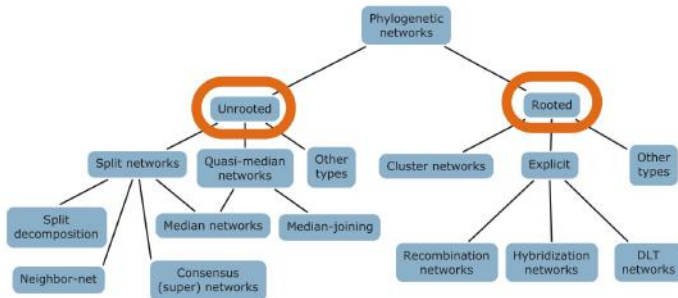
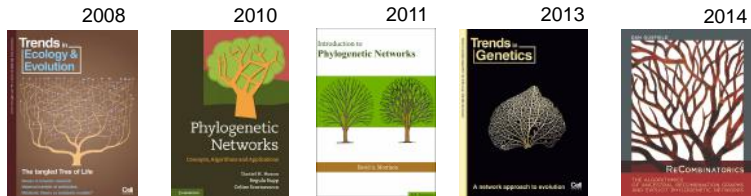
An example - a recombination network



Phylogenetic networks

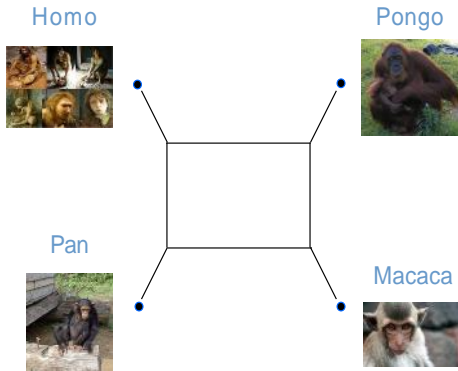


Phylogenetic networks



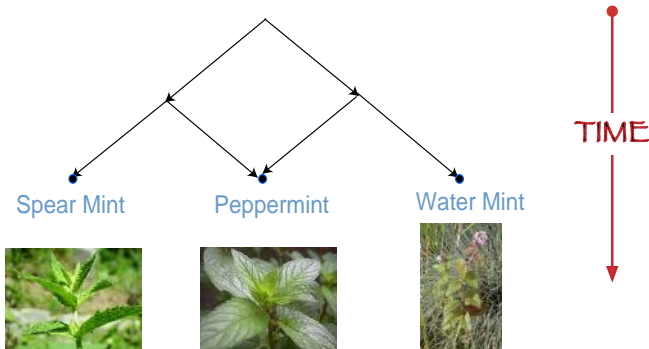
A phylogenetic network ...

... is any connected graph, where terminal nodes are associated to a set of species.

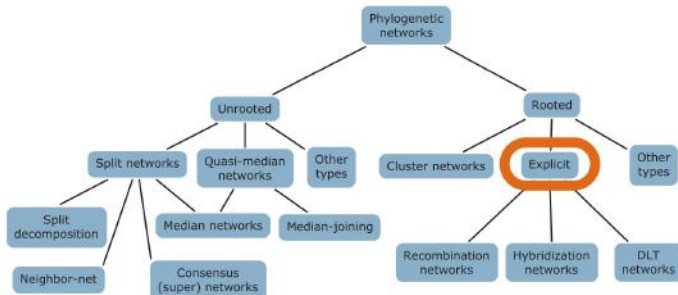
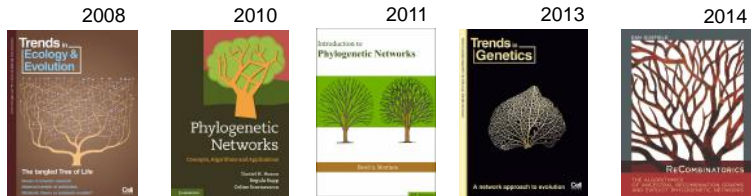


A rooted phylogenetic network ...

... is any single-rooted directed acyclic graph, where terminal nodes are associated to a set of species.

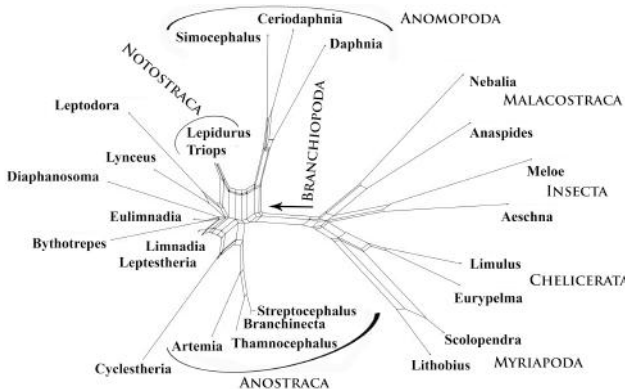


Phylogenetic networks



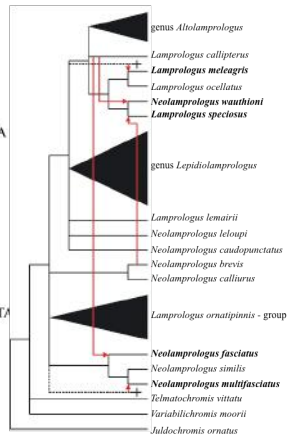
Abstract VS explicit phylogenetic networks

Split network:



Shows conflicting placement of taxa

Hybridization network:



Shows putative hybridization history

The plan of the survey

- ① combinatorial and distance methods not accounting for ILS
 - unrooted networks
 - rooted networks (explicit or not)
- ② methods accounting for ILS (always explicit)

Unrooted phylogenetic networks



SplitsTree4

by Daniel Huson and David Bryant

with contributions from Markus Franz, Migüel Jette',
Tobias Kloepper and Michael Schröder

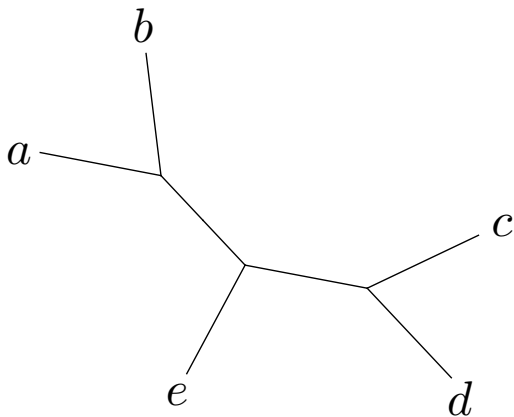
www.splitstree.org

Reconstruction of unrooted phylogenetic networks

- from splits
- from distances (via splits or not)
- from trees (via splits)
- from sequences (via splits or not)

Splits

A *split* $A \mid B$ on \mathcal{X} is a partition of a taxon set \mathcal{X} into two non-empty sets.



Compatible splits

Definition (Compatible splits)

Two splits $S_1 = A_1|B_1$ and $S_2 = A_2|B_2$ on \mathcal{X} are called **compatible**, if one of the following four possible intersections of their split parts is empty: $A_1 \cap A_2$, $A_1 \cap B_2$, $B_1 \cap A_2$ or $B_1 \cap B_2$. Otherwise, the two splits are called incompatible. A set of splits \mathcal{S} is called compatible if all pairs of splits in \mathcal{S} are compatible.

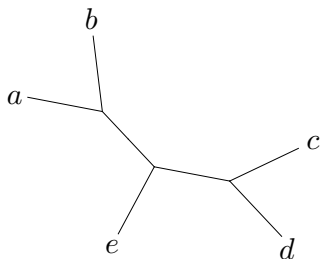
Example

$$\begin{array}{l} S_1 = \begin{array}{l} \{a\}|\{b, c, d, e\} \\ \{b\}|\{a, c, d, e\} \\ \{c\}|\{a, b, d, e\} \\ \{d\}|\{a, b, c, e\} \\ \{e\}|\{a, b, c, d\} \\ \{a, b\}|\{c, d, e\} \\ \{a, b, e\}|\{c, d\} \end{array} \end{array} \quad \begin{array}{l} S_2 = \begin{array}{l} \{a, b, d, e, h\} | \{c, f, g\} \\ \{a, c, d, e, g, h\} | \{b, f\} \\ \{a, c, e, g\} | \{b, d, f, h\} \\ \{a, c, g\} | \{b, d, e, f, h\} \\ \{a, c, e, f, g\} | \{b, d, h\} \\ \{a, e, h\} | \{b, c, d, f, g\} \end{array} \end{array}$$

Compatible splits

Theorem (Compatible splits)

Let S be a set of splits on \mathcal{X} and assume that S contains all trivial splits on \mathcal{X} . There exists a unique unrooted phylogenetic tree T that realizes S , that is, with $\mathcal{S}(T) = S$, if and only if S is compatible.



(a) Unrooted tree T

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

(b) Split encoding of T

Circular splits

Definition (Circular splits)

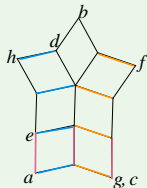
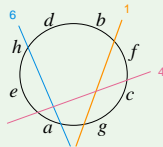
A set of splits \mathcal{S} on \mathcal{X} is called *circular*, if there exists a linear ordering $\pi = (x_1, \dots, x_n)$ of the elements of \mathcal{X} for \mathcal{S} such that each split $S \in \mathcal{S}$ is *interval-realizable*, that is, has the form

$$S = \{x_p, x_{p+1}, \dots, x_q\} \mid (\mathcal{X} \setminus \{x_p, x_{p+1}, \dots, x_q\}),$$

for appropriately chosen $1 < p \leq q \leq n$.

Example

$\{a, b, d, e, h\} \mid \{c, f, g\}$
 $\{a, c, d, e, g, h\} \mid \{b, f\}$
 $\{a, c, e, g\} \mid \{b, d, f, h\}$
 $\{a, c, g\} \mid \{b, d, e, f, h\}$
 $\{a, c, e, f, g\} \mid \{b, d, h\}$
 $\{a, e, h\} \mid \{b, c, d, f, g\}$



Circular splits

Problem (Consecutive Ones problem)

Let M be a binary matrix. Does there exist a permutation of the columns of the matrix M such that in each row, all ones in the row occur in a single consecutive run?

Example

$\{a, b, d, e, h\} \mid \{c, f, g\}$
 $\{a, c, d, e, g, h\} \mid \{b, f\}$
 $\{a, c, e, g\} \mid \{b, d, f, h\}$
 $\{a, c, g\} \mid \{b, d, e, f, h\}$
 $\{a, c, e, f, g\} \mid \{b, d, h\}$
 $\{a, e, h\} \mid \{b, c, d, f, g\}$

a	b	c	d	e	f	g	h
0	0	1	0	0	1	1	0
0	1	0	0	0	1	0	0
0	1	0	1	0	1	0	1
0	1	0	1	1	1	0	1
0	1	0	1	0	0	0	1
0	1	1	1	0	1	1	0

(a) Input matrix

a	e	h	d	b	f	c	g
0	0	0	0	0	1	1	1
0	0	0	0	1	1	0	0
0	0	1	1	1	1	0	0
0	1	1	1	1	1	0	0
0	0	1	1	1	0	0	0
0	0	0	1	1	1	1	1

(b) Permuted matrix

Circular splits

Problem (Consecutive Ones problem)

Let M be a binary matrix. Does there exist a permutation of the columns of the matrix M such that in each row, all ones in the row occur in a single consecutive run?

Example

$\{a, b, d, e, h\} \mid \{c, f, g\}$
 $\{a, c, d, e, g, h\} \mid \{b, f\}$
 $\{a, c, e, g\} \mid \{b, d, f, h\}$
 $\{a, c, g\} \mid \{b, d, e, f, h\}$
 $\{a, c, e, f, g\} \mid \{b, d, h\}$
 $\{a, e, h\} \mid \{b, c, d, f, g\}$

a	b	c	d	e	f	g	h
0	0	1	0	0	1	1	0
0	1	0	0	0	1	0	0
0	1	0	1	0	1	0	1
0	1	0	1	1	1	0	1
0	1	0	1	0	0	0	1
0	1	1	1	0	1	1	0

(a) Input matrix

a	e	h	d	b	f	c	g
0	0	0	0	0	1	1	1
0	0	0	0	1	1	0	0
0	0	1	1	1	1	0	0
0	1	1	1	1	1	0	0
0	0	1	1	1	0	0	0
0	0	0	1	1	1	1	1

(b) Permuted matrix

Note

decision problem is polynomial solvable
finding an ordering of X that minimizes the number of runs of ones in the matrix M (Optimal Consecutive Ones problem) is NP-hard.

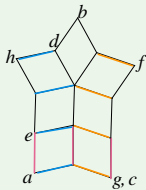
Circular splits

Theorem (Circular implies outer-labeled planar)

A set of splits \mathcal{S} on $\mathcal{X} = \{x_1, \dots, x_n\}$ is circular if and only if it can be represented by a split network N that is outer-labeled planar.

Example

$\{a, b, d, e, h\} \mid \{c, f, g\}$
 $\{a, c, d, e, g, h\} \mid \{b, f\}$
 $\{a, c, e, g\} \mid \{b, d, f, h\}$
 $\{a, c, g\} \mid \{b, d, e, f, h\}$
 $\{a, c, e, f, g\} \mid \{b, d, h\}$
 $\{a, e, h\} \mid \{b, c, d, f, g\}$



Weakly compatible splits

Definition (Weakly compatible splits)

A set of splits \mathcal{S} on \mathcal{X} is called **weakly compatible**, if any *three* distinct splits in \mathcal{S} are weakly compatible. Three such splits $S_1 = \frac{A_1}{B_1}$, $S_2 = \frac{A_2}{B_2}$, and $S_3 = \frac{A_3}{B_3}$ are called *weakly compatible*, if

- ① at least one of the following four intersections is empty:
 $A_1 \cap A_2 \cap A_3$, $A_1 \cap B_2 \cap B_3$, $B_1 \cap A_2 \cap B_3$ and $B_1 \cap B_2 \cap A_3$,
- ② at least one of the following four intersections is empty:
 $B_1 \cap B_2 \cap B_3$, $B_1 \cap A_2 \cap A_3$, $A_1 \cap B_2 \cap A_3$ and $A_1 \cap A_2 \cap B_3$.

Example

$$S_1 = \begin{array}{l} \{a, b, d, e, h\} \mid \{c, f, g\} \\ \{a, c, d, e, g, h\} \mid \{b, f\} \\ \{a, c, e, g\} \mid \{b, d, f, h\} \\ \{a, c, g\} \mid \{b, d, e, f, h\} \\ \{a, c, e, f, g\} \mid \{b, d, h\} \\ \{a, e, h\} \mid \{b, c, d, f, g\} \end{array}$$

$$S_2 = \begin{array}{l} \{a, b, d, e, h\} \mid \{c, f, g\} \\ \{a, c, d, e, g, h\} \mid \{b, f\} \\ \{a, c, e, g\} \mid \{b, d, f, h\} \\ \{a, c, d, e\} \mid \{b, f, g\} \\ \{a, b\} \mid \{c, d, e, f, g\} \\ \{a, e, f\} \mid \{b, c, d, g\} \end{array}$$

Weakly compatible splits

Definition (Weakly compatible splits)

A set of splits \mathcal{S} on \mathcal{X} is called **weakly compatible**, if any *three* distinct splits in \mathcal{S} are weakly compatible. Three such splits $S_1 = \frac{A_1}{B_1}$, $S_2 = \frac{A_2}{B_2}$, and $S_3 = \frac{A_3}{B_3}$ are called *weakly compatible*, if

- 1 at least one of the following four intersections is empty:

Why all this?!?

- 2 Phylogenetic networks reconstructed from weakly compatible are easier than the ones

Ex reconstructed from generic splits

$$S_1 = \begin{array}{l} \{a, b, d, e, h\} \mid \{c, f, g\} \\ \{a, c, d, e, g, h\} \mid \{b, f\} \\ \{a, c, e, g\} \mid \{b, d, f, h\} \\ \{a, c, g\} \mid \{b, d, e, f, h\} \\ \{a, c, e, f, g\} \mid \{b, d, h\} \\ \{a, e, h\} \mid \{b, c, d, f, g\} \end{array}$$

$$S_2 = \begin{array}{l} \{a, b, d, e, h\} \mid \{c, f, g\} \\ \{a, c, d, e, g, h\} \mid \{b, f\} \\ \{a, c, e, g\} \mid \{b, d, f, h\} \\ \{a, c, d, e\} \mid \{b, f, g\} \\ \{a, b\} \mid \{c, d, e, f, g\} \\ \{a, e, f\} \mid \{b, c, d, g\} \end{array}$$

UPN from distances

or “how to get the splits from distances”

From weighted splits to distances

Any split $S \in \mathcal{S}$ can be used to define a distance matrix D_S on \mathcal{X} , by setting:

$$d_S(x, y) = \begin{cases} 1 & \text{if } S \text{ separates } x \text{ and } y, \\ 0 & \text{else.} \end{cases}$$

for all taxa x and y in \mathcal{X} .

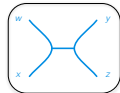
Assume we are given a weighting of a set of splits, $\omega : \mathcal{S} \rightarrow \mathbb{R}^{\geq 0}$, then we define the weighted split metric $D_{(\mathcal{S}, \omega)}$ as:

$$d_{(\mathcal{S}, \omega)}(x, y) = \sum_{S \in \mathcal{S}} \omega(S) \times d_S(x, y)$$

for all taxa x and y in \mathcal{X} .

PN from distances: the split decomposition

The **split decomposition** algorithm [Bandelt and Dress, 1992]: Given a distance matrix D on $\mathcal{X} = \{x_1, \dots, x_n\}$ we start by computing the **isolation index** for quartets and splits :



- for any four taxa w, x, y and z with $\{w, x\} \cap \{y, z\} = \emptyset$, but not necessarily $w \neq x$ or $y \neq z$:

$$\hat{\alpha}_D\left(\frac{\{w, x\}}{\{y, z\}}\right) = \frac{1}{2}(\max\{d(w, x) + d(y, z), d(w, y) + d(x, z), d(w, z) + d(x, y)\} - d(w, x) - d(y, z)).$$

- for any (partial) split S : $\alpha_D(S) = \min\{\hat{\alpha}_D\left(\frac{\{w, x\}}{\{y, z\}}\right) \mid w, x \in A, y, z \in B\} \geq 0$.



Then, we set $\mathcal{X}_0 = \emptyset$ and $\mathcal{S}_0 = \emptyset$. Now, assume that the set of splits \mathcal{S}_i on the first i taxa $\mathcal{X}_i = \{x_1, \dots, x_i\}$. To obtain \mathcal{S}_{i+1} on $\mathcal{X}_{i+1} = \{x_1, \dots, x_{i+1}\}$, for each split $\frac{A}{B} \in \mathcal{S}_i$ do:

- 1 Consider $S = \frac{A \cup \{x_{i+1}\}}{B}$. If $\alpha_D(S) > 0$, set $\omega(S) = \alpha_D(S)$ and add S to \mathcal{S}_{i+1} .
- 2 Consider $S = \frac{A}{B \cup \{x_{i+1}\}}$. If $\alpha_D(S) > 0$, set $\omega(S) = \alpha_D(S)$ and add S to \mathcal{S}_{i+1} .
- 3 Consider $S = \frac{\mathcal{X}_i}{\{x_{i+1}\}}$. If $\alpha_D(S) > 0$, set $\omega(S) = \alpha_D(S)$ and add S to \mathcal{S}_{i+1} .

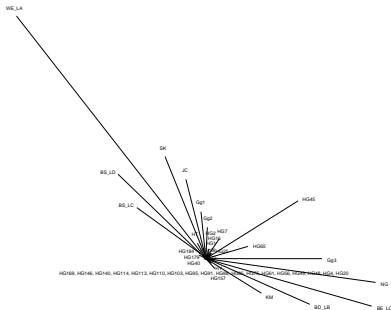
The result is given by \mathcal{S}_n .

PN from distances: the split decomposition

- A split S whose isolation index $\alpha_D(S)$ is greater than 0 is called a D -split. D -splits are always weakly compatible.
- It follows from this that the split decomposition always computes a set of weakly compatible splits

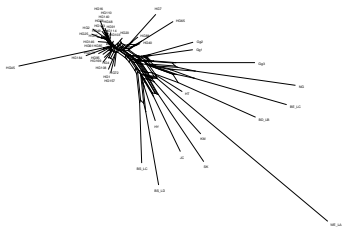
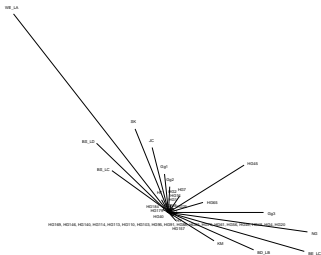
PN from distances: the split decomposition

- A split S whose isolation index $\alpha_D(S)$ is greater than 0 is called a D -split. D -splits are always weakly compatible.
- It follows from this that the split decomposition always computes a set of weakly compatible splits
- The SD is a conservative method
- It can be used for small number of taxa or low divergence



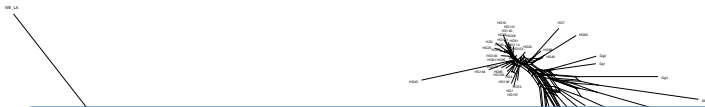
PN from distances: Neighbor-Net

- Given a distance matrix D on \mathcal{X} , the Neighbor-Net algorithm [Bryant and Moulton, 2004] computes a circular ordering π of \mathcal{X} from D and then a set of weighted splits \mathcal{S} that are interval-realizable with respect to π :
 - produces circular splits
 - uses together with circular network algorithm to get planar networks
 - can be used for large number of taxa and high divergence



PN from distances: Neighbor-Net

- Given a distance matrix D on \mathcal{X} , the Neighbor-Net algorithm [Bryant and Moulton, 2004] computes a circular ordering π of \mathcal{X} from D and then a set of weighted splits \mathcal{S} that are interval-realizable with respect to π :
 - produces circular splits
 - uses together with circular network algorithm to get planar networks
 - can be used for large number of taxa and high divergence



Theorem (Neighbor-Net consistency)

Given a circular distance matrix D on \mathcal{X} , the Neighbor-Net algorithm produces a circular ordering π and a set of weighted splits \mathcal{S} that are interval-realizable with respect to π such that $D = D(\mathcal{S})$.

PN from distances

Other algorithms from distances:

- Minimum spanning network
- T-Rex
- ...


A great source of information:

<http://phylnet.univ-mlv.fr/>

Who is Who in Phylogenetic Networks

🏠 Authors Community Keywords Publications **Software** Browse Basket Account Contribute! About Help 🔍

Programs and their Input Data

How do I interact with the graph 

Below, you can find all programs present at least **1 time(s)** in [Who is who in phylogenetic networks](#), as well as the links with the data they use as input.

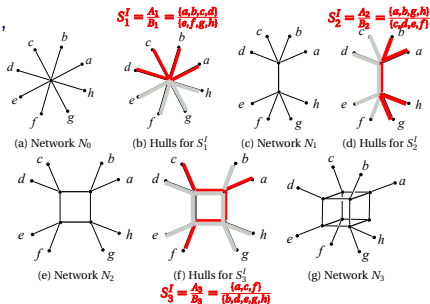
UPN from splits

or “what to do with the splits?”

PN from splits: the Convex hull algorithm

Let \mathcal{S} be a set of splits on \mathcal{X} comprising all trivial ones. Assume that we have already computed a split network N for \mathcal{S} . To obtain a split network for $\mathcal{S} \cup S$, where $S = \frac{A}{B}$ is a new non-trivial split, modify N as follows:

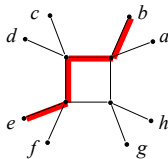
- 1 Compute the two convex hulls $H(A)$ and $H(B)$ in N and let M be the split graph induced in N by the set of nodes $H(A) \cap H(B) \neq \emptyset$.
- 2 Create a copy M' of M and for each node v and edge e in M let v' and e' denote their copies in M' , respectively.
- 3 For every edge f that leads from some node u in $H(B) \setminus H(A) \neq \emptyset$ to a node v in M , redirect the edge f so that it leads from u to v' .
- 4 Connect each pair of nodes v in M and v' in M' by a new edge $f = (v, v')$ and set $\sigma(f) = S$



PN from splits: the circular network algorithm

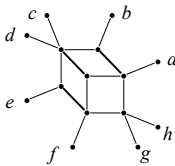
Let \mathcal{S} be a set of splits on \mathcal{X} comprising all trivial ones. Assume that we have already computed a split network N for \mathcal{S} . To obtain a split network for $\mathcal{S} \cup S$, where $S = \frac{\{x_p, \dots, x_q\}}{\mathcal{X} \setminus \{x_p, \dots, x_q\}}$ is a new non-trivial split, modify N as follows: (splits have to be considered in a certain order)

- 1 Determine the path $M(x_p, x_q)$ and let \dot{M} denote the path obtained by removing the first and last (leaf) edges from $M(x_p, x_q)$.
- 2 Create a copy \dot{M}' of \dot{M} and for each node v and edge e in \dot{M} let v' and e' denote their copies, respectively.
- 3 Redirect every edge f that leads from some node $u = \lambda(x_i)$ to some node v in \dot{M} so that it leads from u to v' , for all $i = p, \dots, q$.
- 4 Connect each pair of nodes v in \dot{M} and v' in \dot{M}' by a new edge $f = (v, v')$ and set $\sigma(f) = S_{t+1}$.



(a) Network N_2

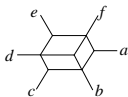
$$\frac{A}{B} = \frac{\{a, f, g, h\}}{\{b, c, d, e\}}$$



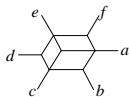
(b) Network N_3

PN from splits: **attention!!!**

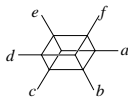
All four different split networks shown below represent the same set of splits. The split networks shown in (a) and (b) were computed using the circular network algorithm processing the splits and taxa in two different orders. The one shown in (c) was constructed using the convex hull algorithm. The split network shown in (d) can be obtained by deleting superfluous edges in any of the first three.



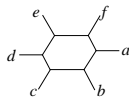
(a)



(b)



(c)



(d)

UPN from trees

or “how to get splits from a bunch of trees”

PN from trees: Consensus split networks

Consensus splits [Holland et al, 2004]

- Input: Trees on identical taxon sets
- Determine splits in more than X% of trees
- For >50%, the result is compatible

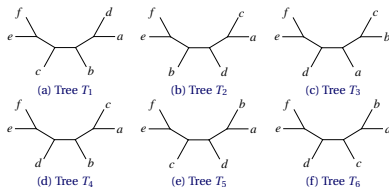


Figure 11.1 (a)– (f) Six different phylogenetic trees T_1, \dots, T_6 on $\mathcal{X} = \{a, b, c, d, e, f\}$. (g) Their majority consensus tree and (h) their consensus split network for $d = 2$, representing all splits that are present in more than $\frac{1}{3}$ of the trees. Note that in this case the network is still a tree, but more resolved than the majority consensus tree. (i) The consensus split network for $d = 5$ and (j) the split network representing all splits present in the six trees.

PN from trees: Consensus split networks

Consensus splits [Holland et al, 2004]

- Input: Trees on identical taxon sets
- Determine splits in more than X% of trees
- For >50%, the result is compatible

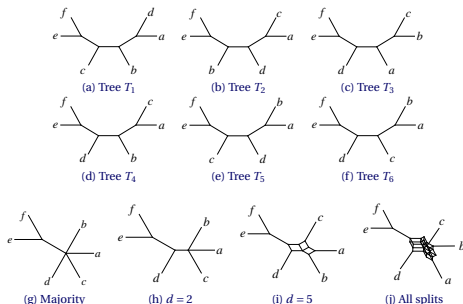


Figure 11.1 (a)–(f) Six different phylogenetic trees T_1, \dots, T_6 on $\mathcal{X} = \{a, b, c, d, e, f\}$. (g) Their majority consensus tree and (h) their consensus split network for $d = 2$, representing all splits that are present in more than $\frac{1}{3}$ of the trees. Note that in this case the network is still a tree, but more resolved than the majority consensus tree. (i) The consensus split network for $d = 5$ and (j) the split network representing all splits present in the six trees.

PN from trees: Consensus super splits networks

Consensus super splits [Huson et al, 2004, Whitfield et al 2008].

Input: Trees on overlapping taxon sets

- Use the Z-closure to complete partial splits
- Use the “distortion” values to filter splits

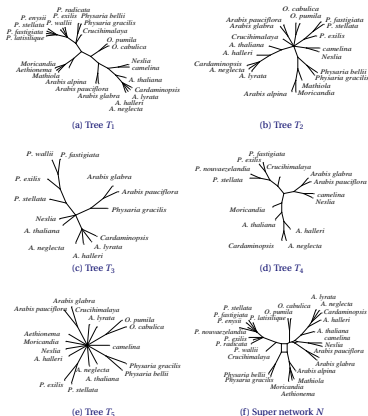


Figure 11.3 (a)–(e) Five partial gene trees T_1, \dots, T_5 on 13–25 plant taxa. (f) The corresponding super split network N on all 26 taxa, computed using the Z-closure method. The edges in N are scaled to represent the number of input trees that contain the edge. The network N shows that the placement of the pair of taxa *Physaria bellii* and *Physaria gracilis* differs in the five trees.

The Z-closure

The *Z-closure* method takes a set of **partial splits** on \mathcal{X} as input and infers a set of **complete splits** on \mathcal{X} as output.

- Two partial splits $S_1 = \frac{A_1}{B_1} \in \mathcal{S}$ and $S_2 = \frac{A_2}{B_2} \in \mathcal{S}$ are said to be in *Z-relation* to each other, if exactly one of the four intersections $A_1 \cap A_2$, $A_1 \cap B_2$, $B_1 \cap A_2$ or $B_1 \cap B_2$ is empty. In this case, if the empty intersection is $A_1 \cap B_2$, say, then we write $\frac{A_1}{B_1} Z \frac{A_2}{B_2}$.
- The *Z-operation* applied to S_1 and S_2 is defined as the creation of two new splits

$$S'_1 = \frac{A_1}{B_1 \cup B_2} \text{ and } S'_2 = \frac{A_1 \cup A_2}{B_2}.$$

- If at least one of the two new splits contains more taxa than its predecessor, the pair of splits is called *productive*.

We obtain a set of complete splits from a set partial splits $\mathcal{S} = \{S_1, \dots, S_m\}$ on \mathcal{X} as follows: While \mathcal{S} contains a productive pair of splits $\{S_i, S_j\}$, apply the Z-operation to obtain two new splits $\{S'_i, S'_j\}$ and then replace the former pair by the latter pair in \mathcal{S} . Finally, add all trivial splits on \mathcal{X} .

The distortion values

- Let $\mathcal{T} = (T_1, \dots, T_k)$ be a set of partial trees on \mathcal{X} . For any complete split S on \mathcal{X} we define the distortion of S (with respect to \mathcal{T}) as $\sigma(S) = \sum_{i=1}^k \sigma(T_i, S)$
- $\sigma(T_i, S)$ denotes the minimal homoplasy score for S on the input tree T_i , i.e. the smallest number of (SPR or TBR) branch-swapping operations required to transform some refinement of T_i into a tree that contains the split S
- The distortion of a split can be efficiently computed using dynamic programming

UPN from sequences

Median networks

For a condensed¹ multiple alignment M of binary sequences on \mathcal{X} , its **median network** is a phylogenetic network $N = (V, E, \sigma, \lambda)$ whose node set is given by the **median closure** $V = \bar{M}$ and in which any two nodes a and b are connected by an edge e of color $\sigma(e) = i \in E$, if and only if they differ in exactly in their i -th position (as haplotypes). An associated taxon labeling $\lambda : X \rightarrow V$ maps each taxon x onto the node $\lambda(x)$ that represents the corresponding sequence.

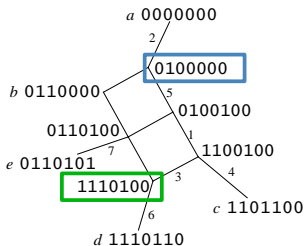
¹ each set of identical sequences is pooled into a single haplotype, then all constant columns are removed and finally, every set of columns that have the same pattern is replaced by a single column i that is assigned a weight $\omega(i)$ that equals the number of columns in the represented set.

Median networks

For a condensed¹ multiple alignment M of binary sequences on \mathcal{X} , its **median network** is a phylogenetic network $N = (V, E, \sigma, \lambda)$ whose node set is given by the **median closure** $V = \bar{M}$ and in which any two nodes a and b are connected by an edge e of color $\sigma(e) = i \in E$, if and only if they differ in exactly in their i -th position (as haplotypes). An associated taxon labeling $\lambda : X \rightarrow V$ maps each taxon x onto the node $\lambda(x)$ that represents the corresponding sequence.

a	0000000
b	0110000
c	1101100
d	1110110
e	0110101

(a) Alignment M



(b) Median network N

¹ each set of identical sequences is pooled into a single haplotype, then all constant columns are removed and finally, every set of columns that have the same pattern is replaced by a single column i that is assigned a weight $\omega(i)$ that equals the number of columns in the represented set.

Quasi median networks

a A A A A A
 b B B A A A
 c A B A B B
 d A A B B C
 e A A C B C

(a) Input M

a 0 0 0 0 0 0 0 0
 b 1 1 0 0 0 0 0 0
 c 0 1 0 0 0 1 1 0
 d 0 0 1 1 0 1 1 0
 e 0 0 1 0 1 1 1 0

(b) Binary expansion M_1

a 0 0 0 0 0 0 0
 b 1 1 0 0 0 0 0
 c 0 1 0 0 0 1 1
 d 0 0 1 1 0 1 0
 e 0 0 1 0 1 1 0

(c) Condensed M_1

0 0 0 0 0 0 0
 1 1 0 0 0 0 0
 0 1 0 0 0 1 1
 0 0 1 1 0 1 0
 0 0 1 0 1 1 0
 0 0 1 0 0 1 0
 0 0 0 0 0 1 0
 0 1 0 0 0 0 0
 0 1 0 0 0 1 0

(d) Median closure M_2

0 0 0 0 0 0 0 0
 1 1 0 0 0 0 0 0
 0 1 0 0 0 1 1 0
 0 0 1 1 0 1 1 0
 0 0 1 0 1 1 1 0
 0 0 1 0 0 1 1 0
 0 0 0 0 0 1 1 0
 0 1 0 0 0 0 0 0
 0 1 0 0 0 1 1 0

(e) Expanded M_2

A A A A A
 B B A A A
 A B A B B
 A A C B C
 A A B B C
 A A * B C
 A A A B *
 A B A A A
 A B A B *

(f) Multi-states M_3

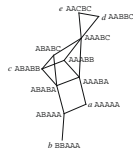
A 0 0 0
 B 1 1 0
 C 1 0 1

$AA * BC = \begin{cases} A A A B C \\ A A B B C \\ A A C B C \end{cases}$
 $AAAB * = \begin{cases} A A A B A \\ A A A B B \\ A A A B C \end{cases}$
 $ABAB * = \begin{cases} A B A B A \\ A B A B B \\ A B A B C \end{cases}$

(g) Expansion of virtual medians

A A A A A
 B B A A A
 A B A B B
 A A C B C
 A A B B C
 A A A B C
 A A A B A
 A A A B B
 A B A A A
 A B A B A
 A B A B C

(h) Final matrix M_4



(i) Quasi-median network N

How to keep the complexity of the network down...

The number of nodes of the quasi-median network can be very large, even for a small number of short sequences. Thus, the quasi-median network is rarely useful in practice. There exist two alternative methods:

- median-joining algorithm, which aims at computing an UPN that is as informative as a quasi-median network, but usually much smaller. The algorithm has a parameter Δ that is used to control how complex the resulting phylogenetic network will be.
- geodesically-pruned quasi-median networks: a method that aims at computing a pruned version of the full quasi-median network by considering only those sequences that lie on a geodesic between two of the original input sequences.

How to keep the complexity of the network down...

UPN from ...

quartets ... QNet

<http://www2.cmp.uea.ac.uk/~vlm/qnet/>

<http://phylnet.univ-mlv.fr/>

Rooted phylogenetic networks



Dendroscope 3

by Daniel H. Huson

with contributions from Benjamin Albrecht,
Philippe Gambette, Leo van Iersel,
Celine Scornavacca and others.

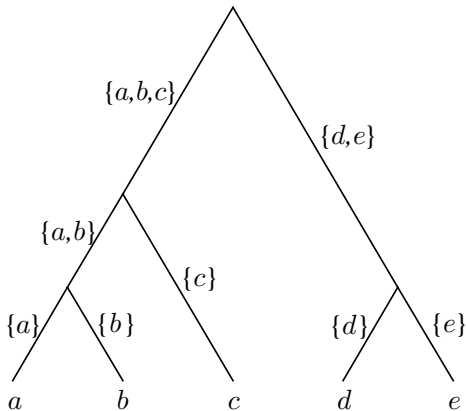
www-ab.informatik.uni-tuebingen.de/software/dendroscope

Reconstruction of rooted phylogenetic networks

- from clusters
- from trees (via clusters or not)
- from sequences
- from distances

Clusters

A *cluster* C on \mathcal{X} is a subset of the taxon set \mathcal{X} .



Compatible clusters

Definition (Compatible splits)

Two clusters C_1 and C_2 on X are called **compatible**, if $C_2 \subseteq C_1$ or $C_1 \subseteq C_2$ or $C_1 \cap C_2 = \emptyset$. Otherwise, the two clusters are called **incompatible**. A set of clusters \mathcal{C} is called **compatible** if all pairs of clusters in \mathcal{C} are compatible.

Example

$$\begin{array}{cc} \begin{array}{c} \{a\} \\ \{b\} \\ \{c\} \\ \{d\} \\ \{e\} \\ \{a, b\} \\ \{d, e\} \\ \{a, b, c\} \\ \{a, b, c, d, e\} \end{array} & \begin{array}{c} \{a\} \\ \{b\} \\ \{c\} \\ \{d\} \\ \{e\} \\ \{a, b\} \\ \{d, e\} \\ \{a, b, e\} \\ \{a, b, c, d, e\} \end{array} \\ C_1 = & C_2 = \end{array}$$

Compatible clusters

Definition (Compatible splits)

Two clusters C_1 and C_2 on X are called **compatible**, if $C_2 \subseteq C_1$ or $C_1 \subseteq C_2$ or $C_1 \cap C_2 = \emptyset$. Otherwise, the two clusters are called **incompatible**. A set of clusters \mathcal{C} is called **compatible** if all pairs of clusters in \mathcal{C} are compatible.

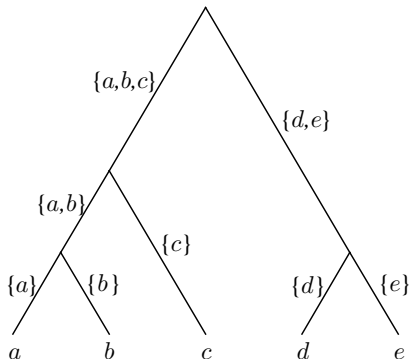
Example

$$\begin{array}{cc} \mathcal{C}_1 = & \mathcal{C}_2 = \\ \begin{array}{c} \{a\} \\ \{b\} \\ \{c\} \\ \{d\} \\ \{e\} \\ \{a, b\} \\ \{d, e\} \\ \{a, b, c\} \\ \{a, b, c, d, e\} \end{array} & \begin{array}{c} \{a\} \\ \{b\} \\ \{c\} \\ \{d\} \\ \{e\} \\ \{a, b\} \\ \{d, e\} \\ \{a, b, e\} \\ \{a, b, c, d, e\} \end{array} \end{array}$$

Compatible clusters

Theorem (Compatible clusters)

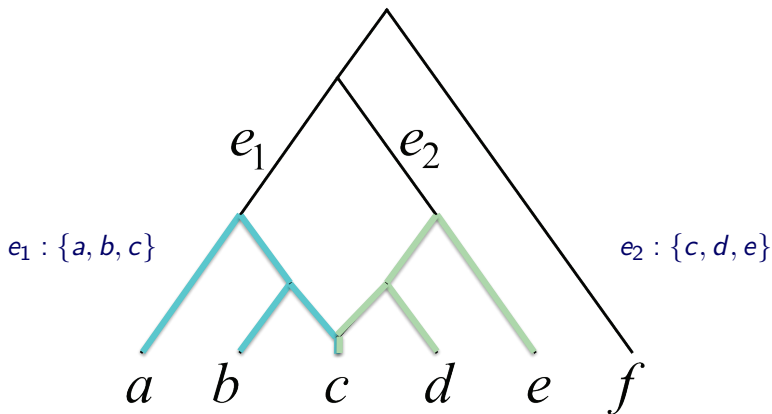
Let \mathcal{C} be a set of clusters on \mathcal{X} and assume that \mathcal{C} contains all trivial splits on \mathcal{X} . There exists a unique rooted phylogenetic tree T that realizes \mathcal{C} , that is, with $\mathcal{C}(T) = \mathcal{C}$, if and only if \mathcal{C} is compatible.



RPN from clusters

When a phyl. network N represents a cluster C ?

HARDWIRED SENSE : if there exists a tree edge of N such that the set of all taxa below the edge equals C



RPN from clusters - Cluster networks

The [Cluster-popping algorithm](#) [Huson & Rupp, 2008]:

for each cluster $C \in \mathcal{C}$ create a node u and define $\nu(C) = u$ and $\nu^{-1}(u) = C$.

Create a root node ρ and set $\nu^{-1}(\rho) = \mathcal{X}$

for each $C \in \mathcal{C}$ in order of non-increasing cardinality **do**

Unmark all nodes

Push ρ onto a stack S and mark ρ

while S is not empty **do**

Pop v off S

Set *isBelow* = *false*

for each child w of v **do**

if $C \subset \nu^{-1}(w)$ **then**

Set *isBelow* = *true*

if w is unmarked **then** Mark w and push w onto S

if *isBelow* = *false* **then** Create a new edge $(v, \nu(C))$

for each node v with indegree ≥ 2 and outdegree $\neq 1$ **do**

Create a new node v'

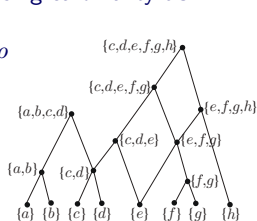
Redirect all in-edges of v to v'

Create a new edge (v', v)

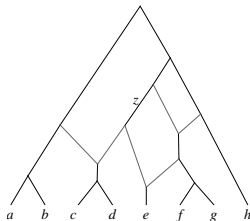
for each taxon $a \in \mathcal{X}$ **do**

Set $\lambda(a) = \nu^{-1}(\{a\})$

end



(a) Hasse diagram



(b) Cluster network

RPN from clusters - Hardwired sense

The [Cluster-popping algorithm](#) [Huson & Rupp, 2008]:

for each cluster $C \in \mathcal{C}$ create a node u and define $\nu(C) = u$ and $\nu^{-1}(u) = C$.

Create a root node ρ and set $\nu^{-1}(\rho) = \mathcal{X}$

for each $C \in \mathcal{C}$ in order of non-increasing cardinality **do**

Unmark all nodes

Push ρ onto a stack S and mark ρ

while S is not empty **do**

Pop v off S

Set *isBelow* = *false*

for each child w of v **do**

if $C \subset \nu^{-1}(w)$ **then**

Set *isBelow* = *true*

if w is unmarked **then** Mark w and push w onto S

if *isBelow* = *false* **then** Create a new edge $(v, \nu(C))$

for each node v with indegree ≥ 2 and outdegree $\neq 1$ **do**

Create a new node v'

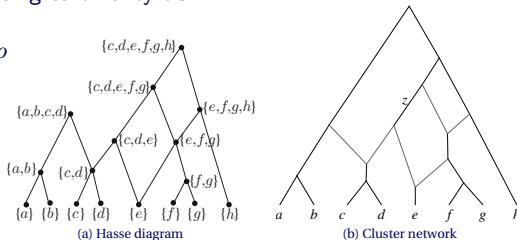
Redirect all in-edges of v to v'

Create a new edge (v', v)

for each taxon $a \in \mathcal{X}$ **do**

Set $\lambda(a) = \nu^{-1}(\{a\})$

end

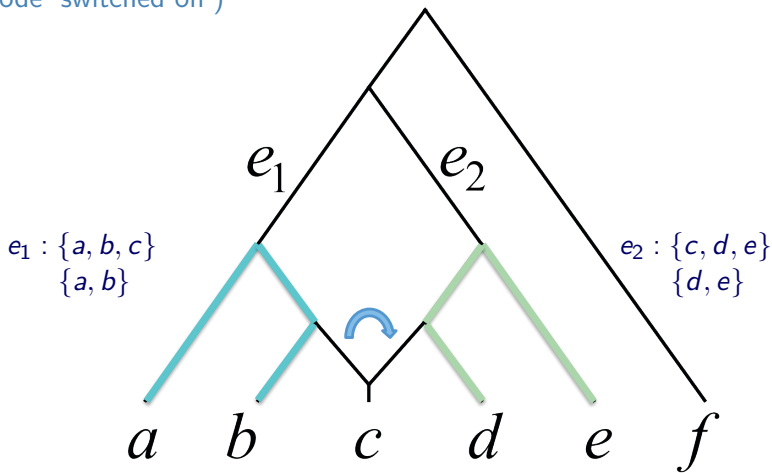


Problem

The networks obtained in this way are too complex

When a phyl. network N represents a cluster C ?

SOFTWARED SENSE : if there exists a tree edge of N such that the set of all taxa below the edge equals C (with one edge per reticulation node "switched on")

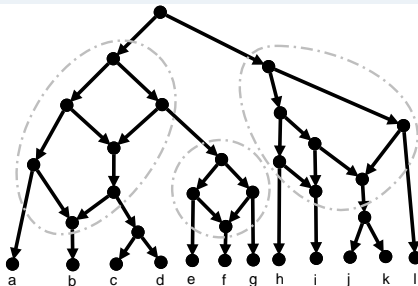


Constructing minimal softwired networks

- cluster containment: NP-hard
- minimization NP-hard, APX-hard

A possible solution ... topological constraints:

- galled trees (if every non-trivial biconnected component of N properly contains exactly one reticulation) (it does not always exist)
- galled networks (if every reticulation in N has a tree cycle) (still NP-hard)
- level- k networks (maximum reticulation number among biconnected components of N is k) (still NP-hard)



RPN from clusters - Softwired sense

CASS algorithm [van Iersel et al., 2010], guaranteed to construct minimal networks if level is 1 or 2:

Algorithm 8.5.1 (Level- k networks for tangled clusters) *Let $k \geq 0$ be a fixed number. Input is a set of clusters \mathcal{C} on \mathcal{X} , which is assumed to be tangled in the initial call to the algorithm. Recursively construct a set of level- k networks \mathcal{N} for \mathcal{C} in the following steps (if any such network exists):*

Set $\mathcal{N} = \emptyset$
if $k = 0$ **then**
 if \mathcal{C} **is compatible then**
 Compute the rooted phylogenetic tree T for \mathcal{C} (and all trivial clusters on \mathcal{X})
 Insert an auxiliary edge that separates the root from the rest of the tree T
 Set $\mathcal{N} = \{T\}$
 else (comment: $k > 0$)
 Create a new auxiliary taxon z
 for each taxon $x \in \mathcal{X} \cup \{z\}$ **do**
 Set $\mathcal{C}' = \mathcal{C}|_{\mathcal{X} - \{x\}}$, remove all trivial clusters and then collapse to obtain \mathcal{C}'' on \mathcal{X}''
 Recursively compute the set \mathcal{N}'' of level- $(k-1)$ networks for \mathcal{C}'' on \mathcal{X}''
 for each network N'' in \mathcal{N}'' **do**
 Replace each leaf of N'' , which is labeled by a composite taxon $\tilde{A} \in \mathcal{X}''$, by the rooted phylogenetic tree $T(\mathcal{C}'|_{\tilde{A}})$ that represents $\mathcal{C}'|_{\tilde{A}}$
 Let N' be the resulting network
 for each pair of (not necessarily distinct) edges e_1 and e_2 in a new copy of N' **do**
 Create two nodes r and u , add an edge from r to u and set $\lambda(x) = u$
 Insert a new node v_1 into e_1 and connect v_1 to r
 Insert a new node v_2 into e_2 and connect v_2 to r
 if the resulting network N represents \mathcal{C} (disregarding all auxiliary taxa) **then**
 add N to \mathcal{N}
return \mathcal{N}

If the returned set \mathcal{N} is empty, then the algorithm reports fail.

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

(a) Set of tangled input clusters \mathcal{C}

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

(b) Set $\mathcal{C}' = \mathcal{C}|_{\mathcal{X} - \{c\}}$

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

(c) Set \mathcal{C}'' obtained by collapsing \mathcal{C}'

$\{a, g\}, \{d, e\}$

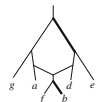
(d) Result of removing $\{b, f\}$ from \mathcal{C}''



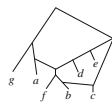
(e) Rooted tree T for set of clusters in (d)



(f) Level-1 network for set of clusters \mathcal{C}'' in (c)



(g) Level-1 network for set of clusters \mathcal{C}' in (b)

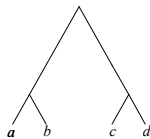


(h) Final level-2 network for set \mathcal{C} in (a)

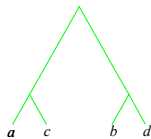
RPN from trees (via clusters)

RPN from trees - option 1

- decompose the trees in clusters
- apply the cluster methods to (a part of) the clusters



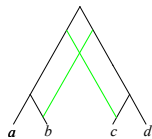
(a) Tree T_1



(b) Tree T_2

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

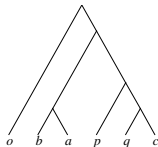
(c) The clusters



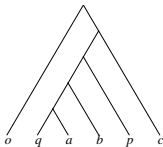
(d) Network N

RPN from trees - option 1

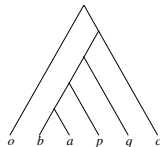
- decompose the trees in clusters
- apply the cluster methods to (a part of) the clusters



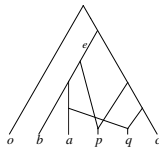
(a) Tree T_1



(b) Tree T_2



(c) Tree T_3

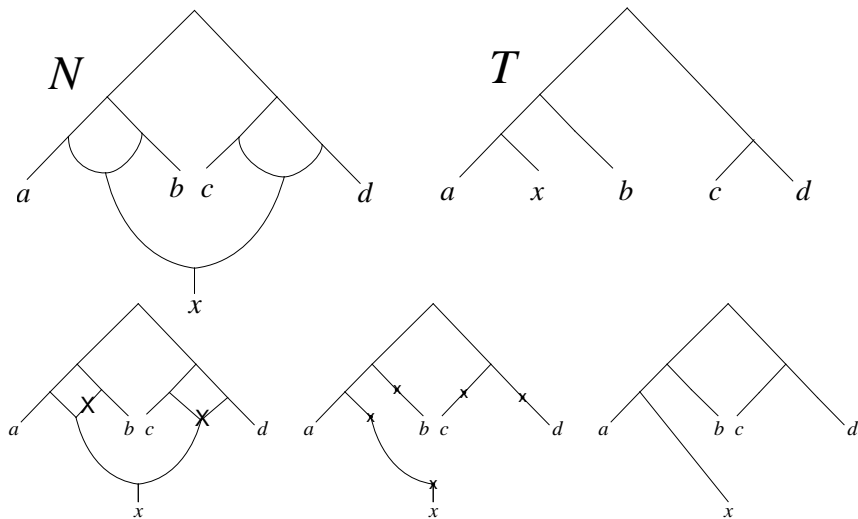


(d) Network N

RPN from trees (NOT via clusters)

When a phyl. network N embeds a tree T ?

if T can be obtained from N by performing a series of node deletions, edge deletions and node suppressions



Problem

Given:

A set of (binary) trees (with same taxa set) and **different** topology.

Question:

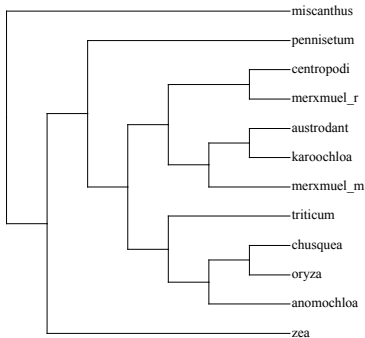
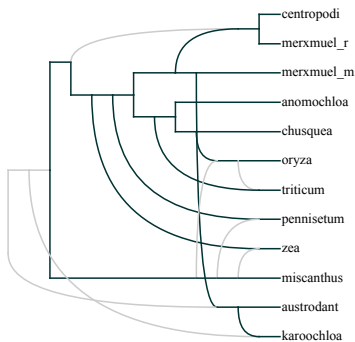
What is the **most probable** evolutionary history?

Assumptions: Difference is caused by hybridizations, parsimony

Answer:

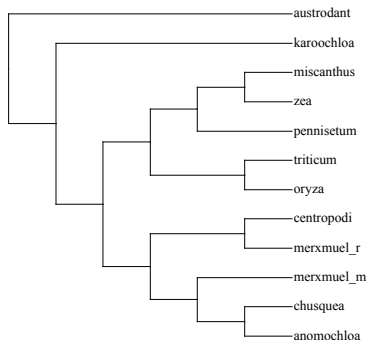
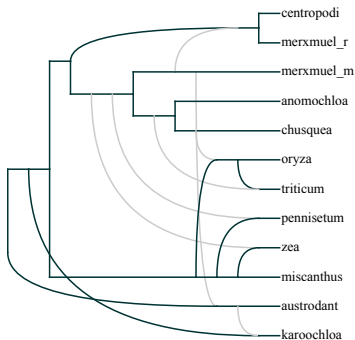
Network embedding the trees with a **minimal** number of hybridization nodes.

Tree embedding



Hybridization network (left) highlighting the embedding of the phylogenetic tree based on gene *rbcL* (right).

Tree embedding



Hybridization network (left) highlighting the embedding of the phylogenetic tree based on gene *waxy* (right).

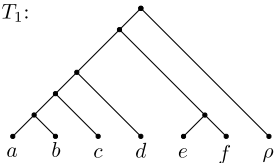
Agreement forests

An **agreement forest** for two rooted bifurcating phylogenetic trees T_1 and T_2 on $\mathcal{X} \cup \rho$ is a set of components $\mathcal{F} = \{F_\rho, F_1, \dots, F_n\}$ on $\mathcal{X} \cup \rho$ such that...

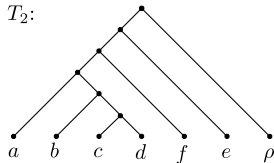
- 1 ... the taxon ρ is contained in F_ρ .
- 2 ... each component F_i is a **restricted subtree** of T_1 and T_2 .
- 3 ... the trees in $\{T_1(\mathcal{X}_i | i = \rho, 1, \dots, n)\}$ and $\{T_2(\mathcal{X}_i | i = \rho, 1, \dots, n)\}$ are **node disjoint subtrees** of T_1 and T_2 , respectively.

Acyclic agreement forests

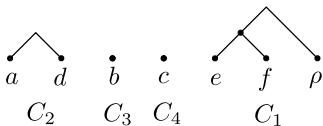
T_1 :



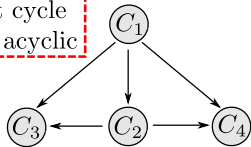
T_2 :



$\mathcal{F}(T_1, T_2)$:

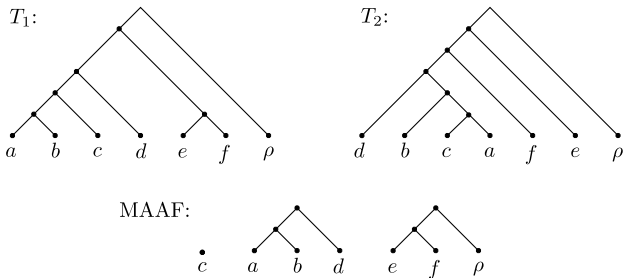


No direct cycle
 \Rightarrow AF is acyclic

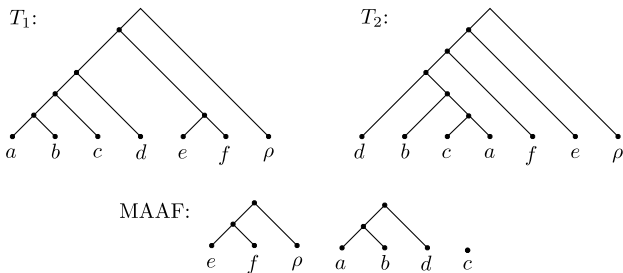


A **maximal acyclic** agreement forest, denoted by *MAAF*, is any acyclic agreement forest of **minimal** size.

Using MAAF to construct hybridization networks

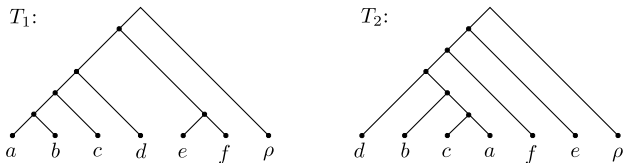


Using MAAF's to construct hybridization networks

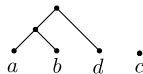
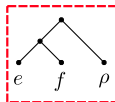


Compute acyclic ordering.

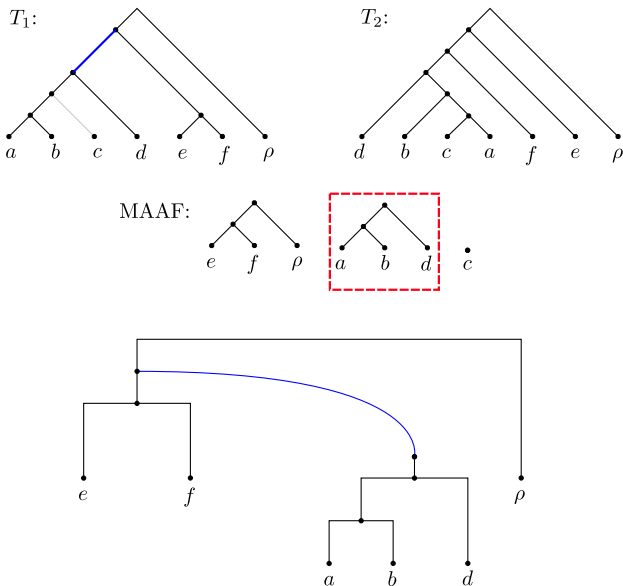
Using MAAF to construct hybridization networks



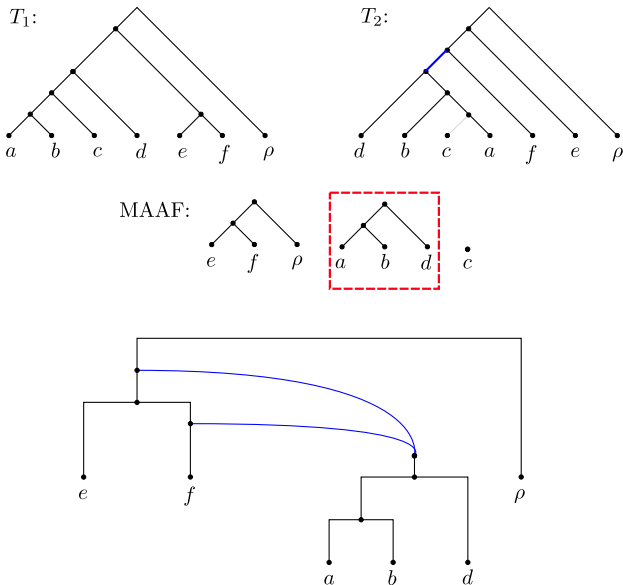
MAAF:



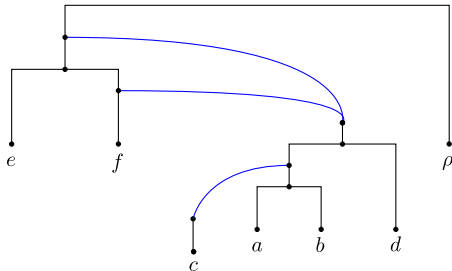
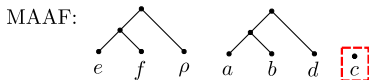
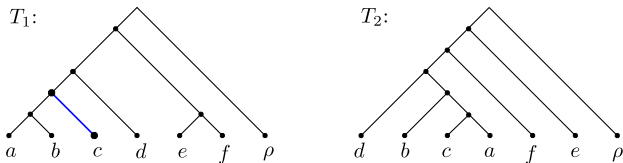
Using MAAF to construct hybridization networks



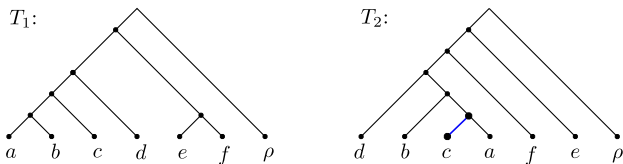
Using MAAF to construct hybridization networks



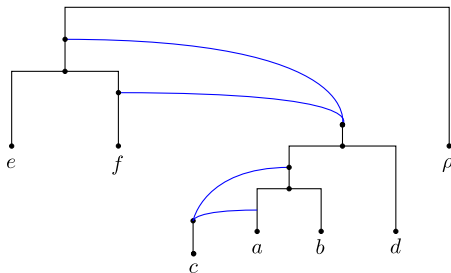
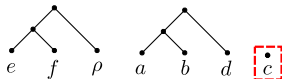
Using MAAF to construct hybridization networks



Using MAAF to construct hybridization networks



MAAF:

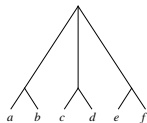


- ① HybridNet : <http://www.cs.cityu.edu.hk/~lwang/software/Hn/treeComp.html>
- ② UltraNet: <http://rnc.r.dendai.ac.jp/ultraNet.html>
- ③ HybridInterleave: <http://www.math.canterbury.ac.nz/~c.semple/software.shtml>
- ④ NonbinaryCycleKiller: the two trees are not necessarily binary
<http://homepages.cwi.nl/~iersel/cyclekiller/>
- ⑤ Dendroscope: the two trees are not necessarily binary and not necessarily on the same taxon set.
- ⑥ Hybroscale: multiply trees, not necessarily binary and not necessarily on the same taxon set
www.bio.ifi.lmu.de/software/services/hybroscale
- ⑦ ...

RPN from triplets

Triplets

A rooted triple is given by a bifurcating, rooted phylogenetic tree on a set of three taxa x, y, z .

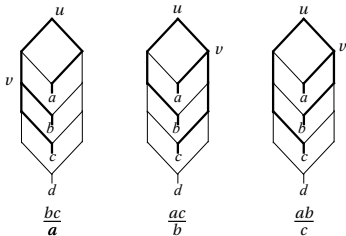


$\frac{ab}{c}$	$\frac{ab}{d}$	$\frac{ab}{e}$	$\frac{ab}{f}$
$\frac{cd}{a}$	$\frac{cd}{b}$	$\frac{cd}{e}$	$\frac{cd}{f}$
$\frac{ef}{a}$	$\frac{ef}{b}$	$\frac{ef}{c}$	$\frac{ef}{d}$

Triplets

A rooted triple $yz|x$ is said to be contained in a rooted phylogenetic network N , if there exist two nodes u and v in N such that:

- 1 There exists a directed path from u to the node labeled x .
- 2 There exists a directed path from u to v .
- 3 There exists a directed path from v to the node labeled y .
- 4 There exists a directed path from v to the node labeled z .
- 5 All four paths are node-disjoint (except at their end-points).



Software for networks from triplets

<https://leovaniersel.wordpress.com/software/>

LEV1ATHAN: A Practical Algorithm for Reconstructing Level-1 Phylogenetic Networks. Combines any set of phylogenetic trees into a level-1 phylogenetic network (a galled tree) that is consistent with a large number of the triplet topologies of the input trees. [Paper](#). [Download](#).

SIMPLISTIC: Constructs level- k phylogenetic networks from triplets. This program always returns a phylogenetic network consistent with all input triplets. Partly based on the SL- k and MINPITS algorithms in [this paper](#). [Download](#).

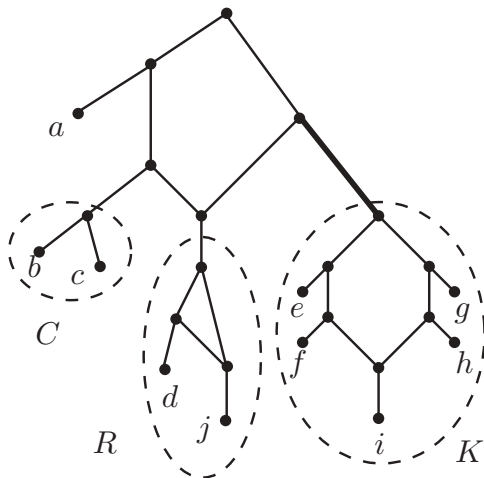
MARLON: Constructs a level-1 phylogenetic network with a minimum number of reticulations consistent with a dense set of triplets, if such a network exists. [Paper](#). [Download](#).

LEVEL2: Constructs a level-2 phylogenetic network consistent with a dense set of triplets, if such a network exists. [Paper](#). [Download](#).

Software for networks from binets and trinets

TriLoNet

<https://www.uea.ac.uk/computing/TriLoNet>



RPN from sequences

– ARGs –

Recombination networks

Definition (Recombination network)

Let M be a multiple alignment of binary sequences of length L , on \mathcal{X} . A *recombination network* N representing M is given by a bicomposing rooted phylogenetic network on \mathcal{X} , together with two additional labellings:

- 1 Each node v of N is labeled by a binary sequence $\sigma(v)$ of length L .
- 2 Each tree edge e is labeled by a set of positions $\delta(e) \subseteq \{1, \dots, L\}$.

These two labellings must fulfill the following compatibility conditions:

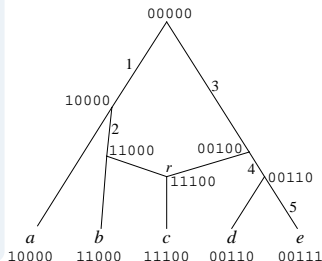
- A The sequence $\sigma(v)$ assigned to any leaf v must equal the sequence in M that is given for the taxon associated with v .
- B If r is a reticulate node (often called a recombination node in this context) with parents v and w , then the sequence $\sigma(r)$ must be obtainable from the two sequences $\sigma(v)$ and $\sigma(w)$ by a crossover.
- C If $e = (v, w)$ is a tree edge, then the set of positions at which the two sequences $\sigma(v)$ and $\sigma(w)$ differ must equal $\delta(e)$.

For computational reasons, the following condition is usually also required

- D Any given position may mutate at most once in the network. In other words, for any given position i there exists at most one edge e with $i \in \delta(e)$.

This condition is usually referred to as the infinite sites assumption because for sequences of infinite length it holds that the probability of the same site being hit by a mutation more than once is zero, under a uniform distribution.

a 10000
 b 11000
 c 11100
 d 00110
 e 00111



Recombination networks

Definition (Recombination network)

Let M be a multiple alignment of binary sequences of length L , on \mathcal{X} . A *recombination network* N representing M is given by a bicomposing rooted phylogenetic network on \mathcal{X} , together with two additional labellings:

- 1 Each node v of N is labeled by a binary sequence $\sigma(v)$ of length L .
- 2 Each tree edge e is labeled by a set of positions $\delta(e) \subseteq \{1, \dots, L\}$.

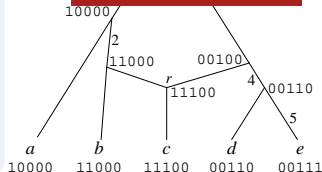
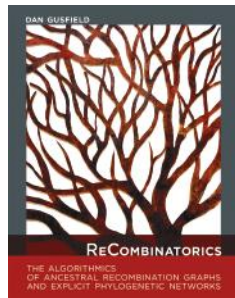
These two labellings must fulfill the following compatibility conditions:

- A The sequence $\sigma(v)$ assigned to any leaf v must equal the sequence in M that is given for the taxon associated with v .
- B If r is a reticulate node (often called a recombination node in this context) with parents v and w , then the sequence $\sigma(r)$ must be obtainable from the two sequences $\sigma(v)$ and $\sigma(w)$ by a crossover.
- C If $e = (v, w)$ is a tree edge, then the set of positions at which the two sequences $\sigma(v)$ and $\sigma(w)$ differ must equal $\delta(e)$.

For computational reasons, the following condition is usually also required

- D Any given position may mutate at most once in the network. In other words, for any given position i there exists at most one edge e with $i \in \delta(e)$.

This condition is usually referred to as the infinite sites assumption because for sequences of infinite length it holds that the probability of the same site being hit by a mutation more than once is zero, under a uniform distribution.



Recombination networks

Definition (Recombination network)

Let M be a multiple alignment of binary sequences of length L , on \mathcal{X} . A *recombination network* N representing M is given by a bicomining rooted phylogenetic network on \mathcal{X} , together with two additional labellings:

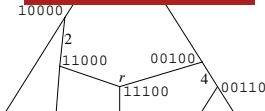
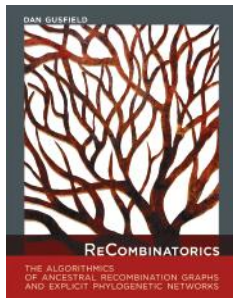
- 1 Each node v of N is labeled by a binary sequence $\sigma(v)$ of length L .
- 2 Each tree edge e is labeled by a set of positions $\delta(e) \subseteq \{1, \dots, L\}$.

These two labellings must fulfill the following compatibility conditions:

- A The sequence $\sigma(v)$ assigned to any leaf v must equal the sequence in M that is given for the taxon associated with v .
- B If r is a reticulate node (often called a recombination node in this context) with parents v and w , then the sequence $\sigma(r)$ must be obtainable from the two sequences $\sigma(v)$ and $\sigma(w)$ by a crossover.
- C If $e = (v, w)$ is a tree edge, then the set of positions at which the two sequences $\sigma(v)$ and $\sigma(w)$ differ must equal $\delta(e)$.

For computational reasons, the following condition is usually also required

- D Any given position may mutate at most once in the network. In other words, for any given position i there exists at most one edge e with $i \in \delta(e)$.



the approach has to solve two NP-hard problems

Bacter: <http://tgvaughan.github.io/bacter/>, a package of BEAST 2

ARGweaver: <http://mdrasmus.github.io/argweaver/doc/>

RPN from sequences

- Other approaches –
