

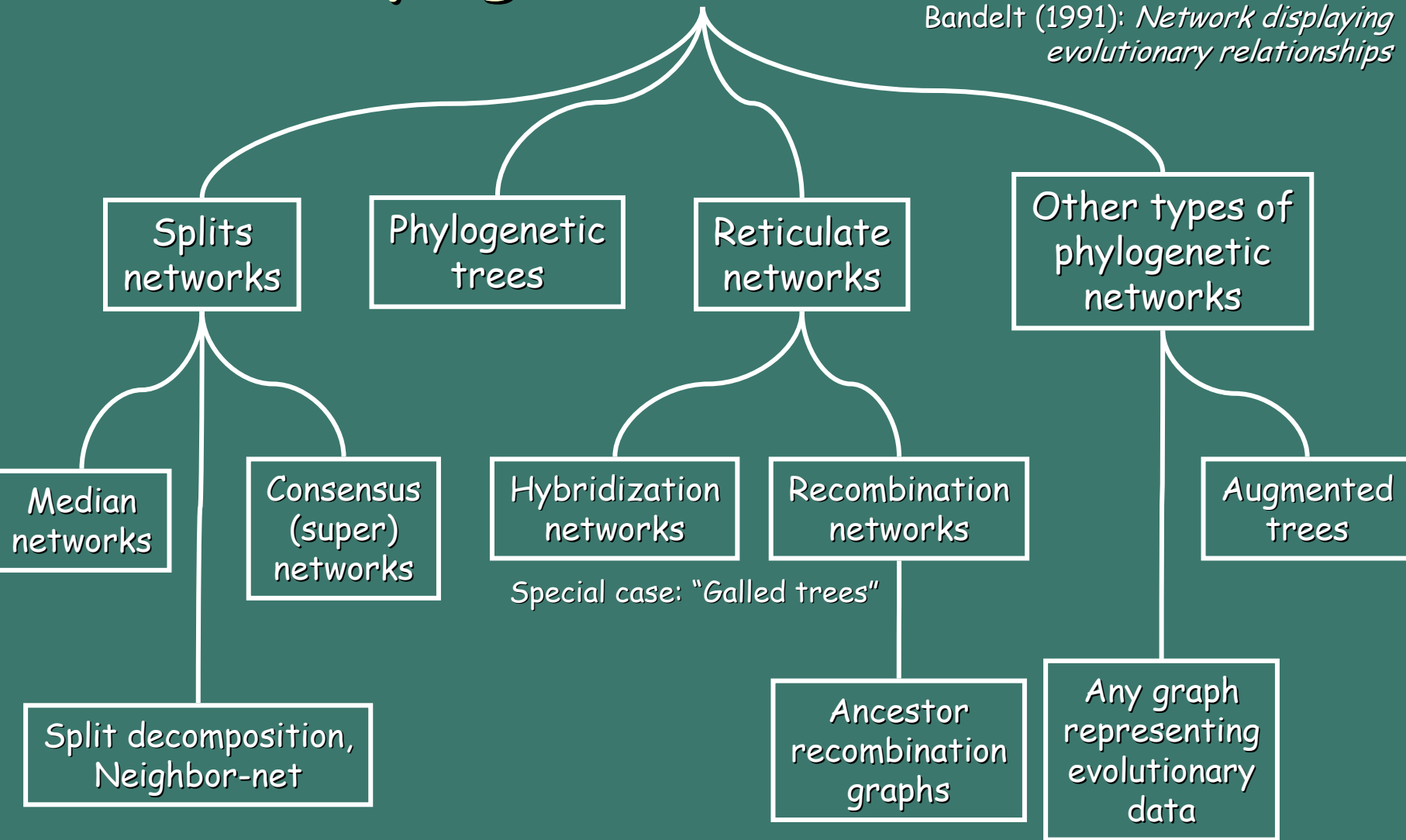
Contents

1. Phylogenetic trees
2. Splits networks
3. Consensus networks
4. Hybridization and reticulate networks
5. Recombination networks

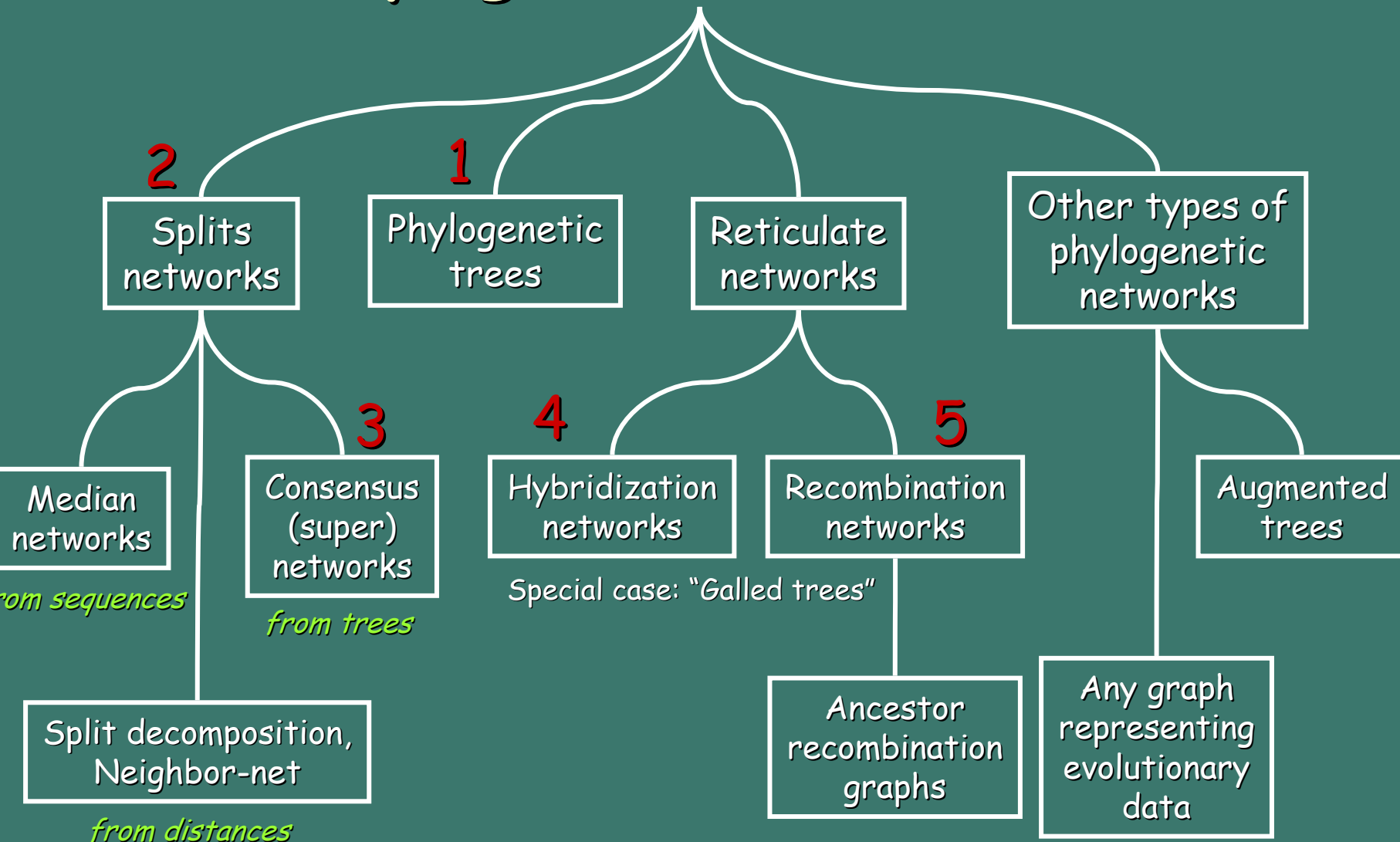


Phylogenetic Networks

Bandelt (1991): *Network displaying evolutionary relationships*

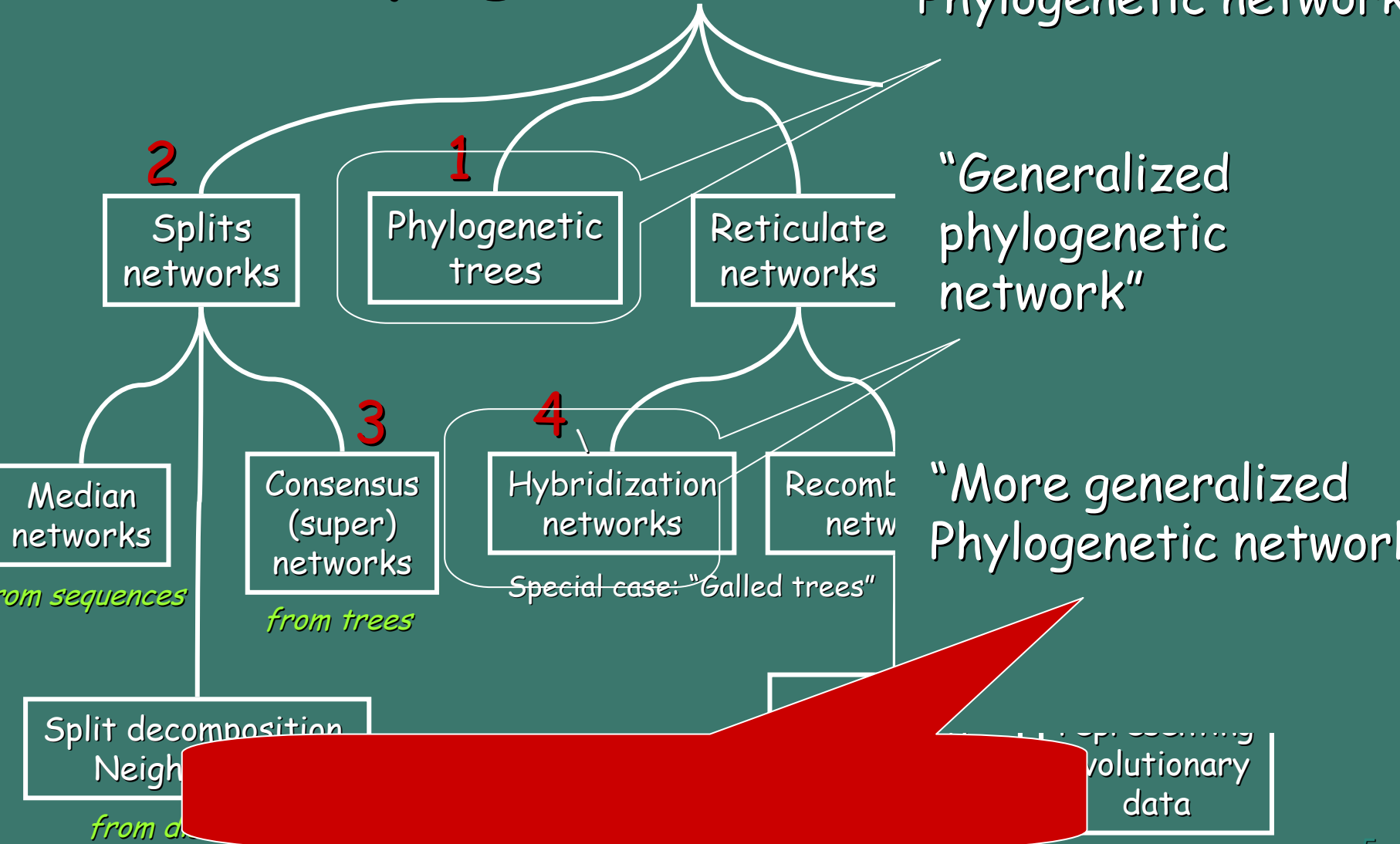


Phylogenetic Networks



Phylogenetic Net

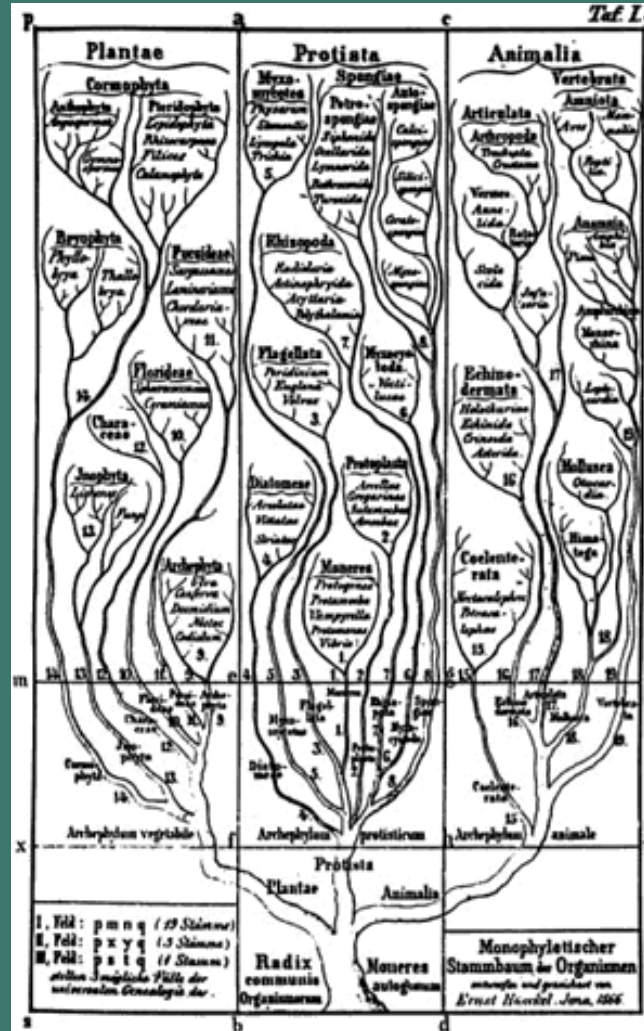
Dan Gusfield:
"Phylogenetic network"



Part I

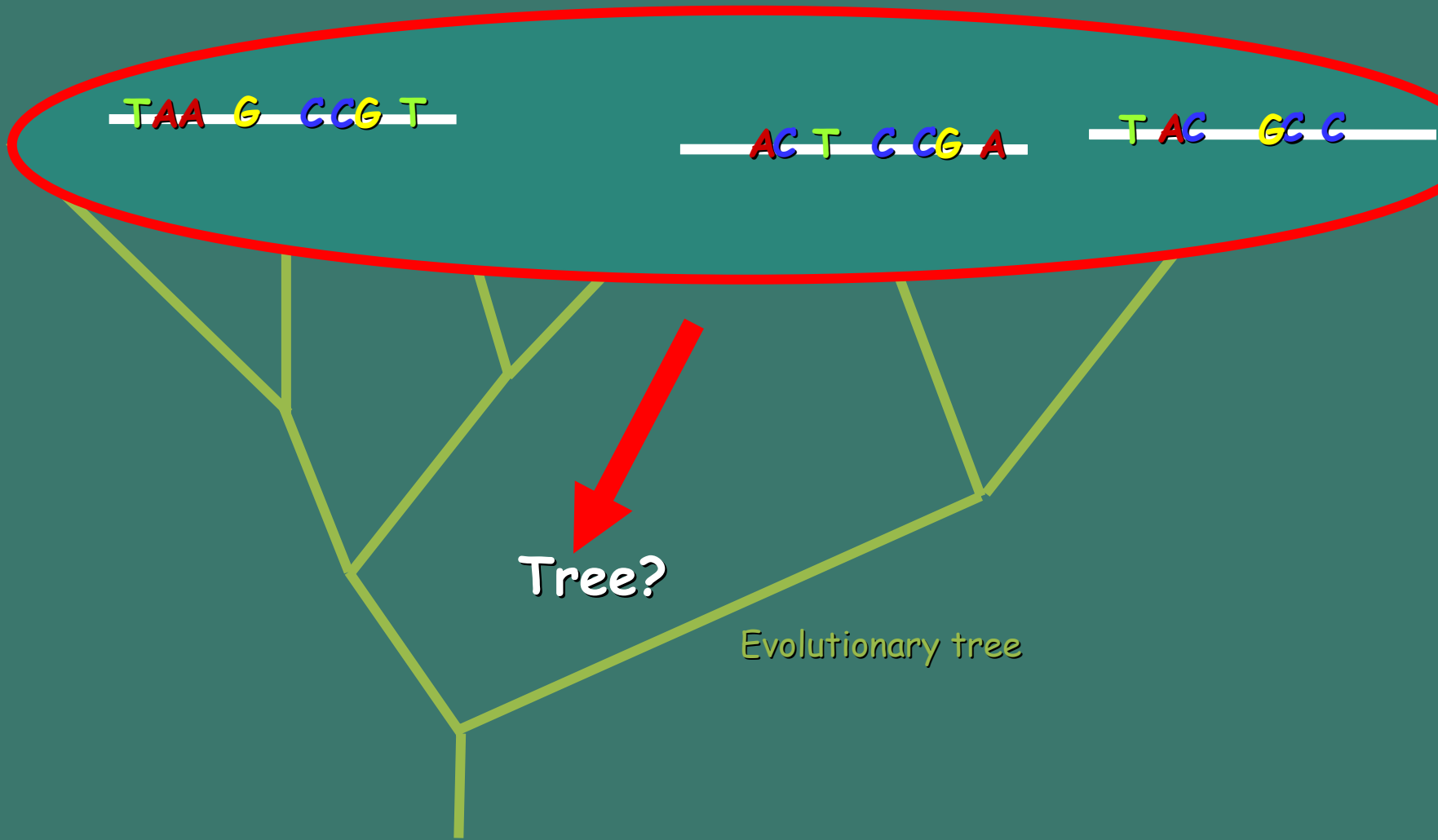
1. Phylogenetic trees
2. Splits networks
3. Consensus networks
4. Hybridization and reticulate networks
5. Recombination networks

Phylogenetic Trees



Ernst Haeckel,
Tree of Life
1866

Tree Reconstruction Problem



Part II

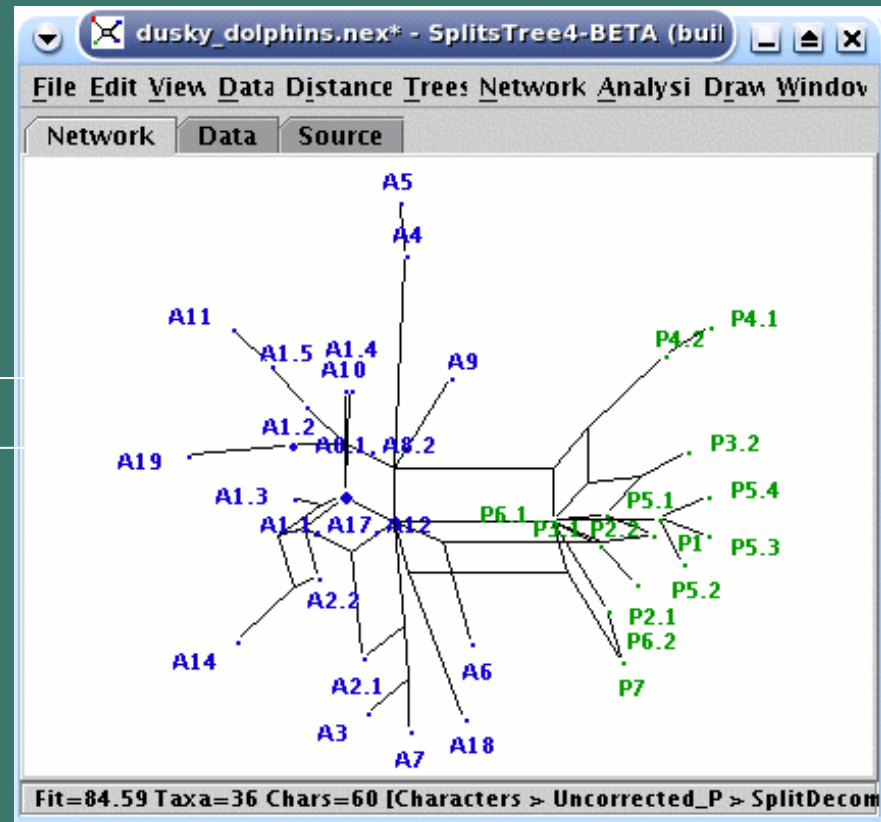
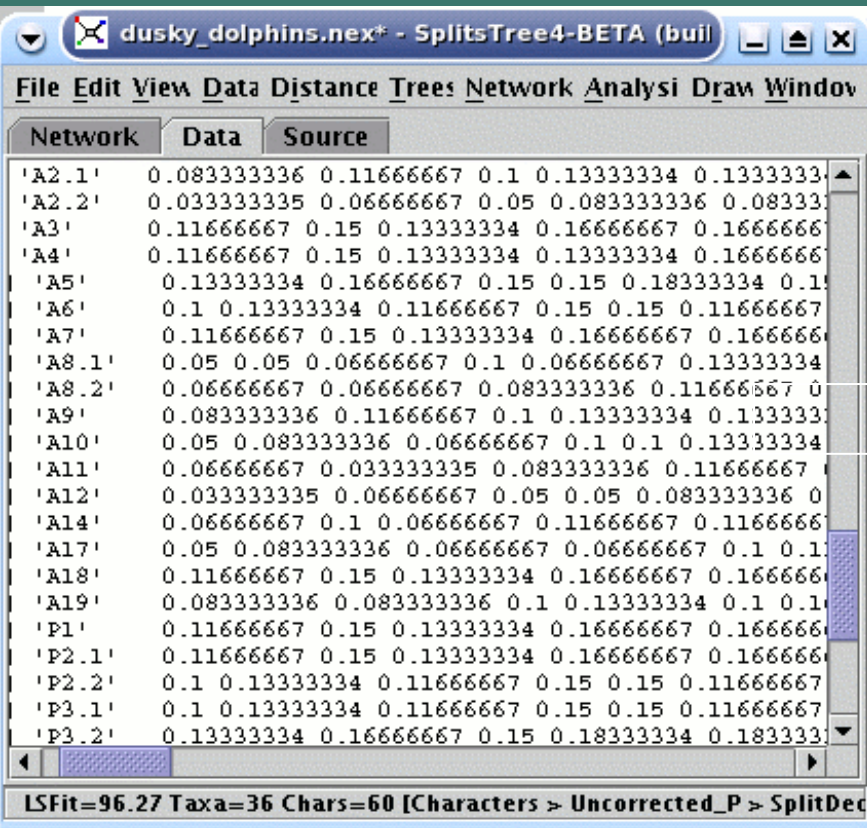
1. Phylogenetic trees
2. Splits networks
3. Consensus networks
4. Hybridization and reticulate networks
5. Recombination networks

Splits Networks

Represents incompatible signals in data, from:

- **Sequences, e.g.:**
 - Median network (Bandelt *et al*/1994)
 - Spectral analysis (Hendy and Penny 1993)
- **Distances, e.g.:**
 - E.g. Split decomposition (Bandelt and Dress 1992)
 - Neighbor-Net (Bryant and Moulton 2002)
- **Trees, e.g.:**
 - Consensus network (Holland and Moulton 2003)
 - Super network (H., Dezulian, Klopper and Steel 2004)
 - Bootstrap network (H., implemented in SplitsTree4)

Distances to Splits Network



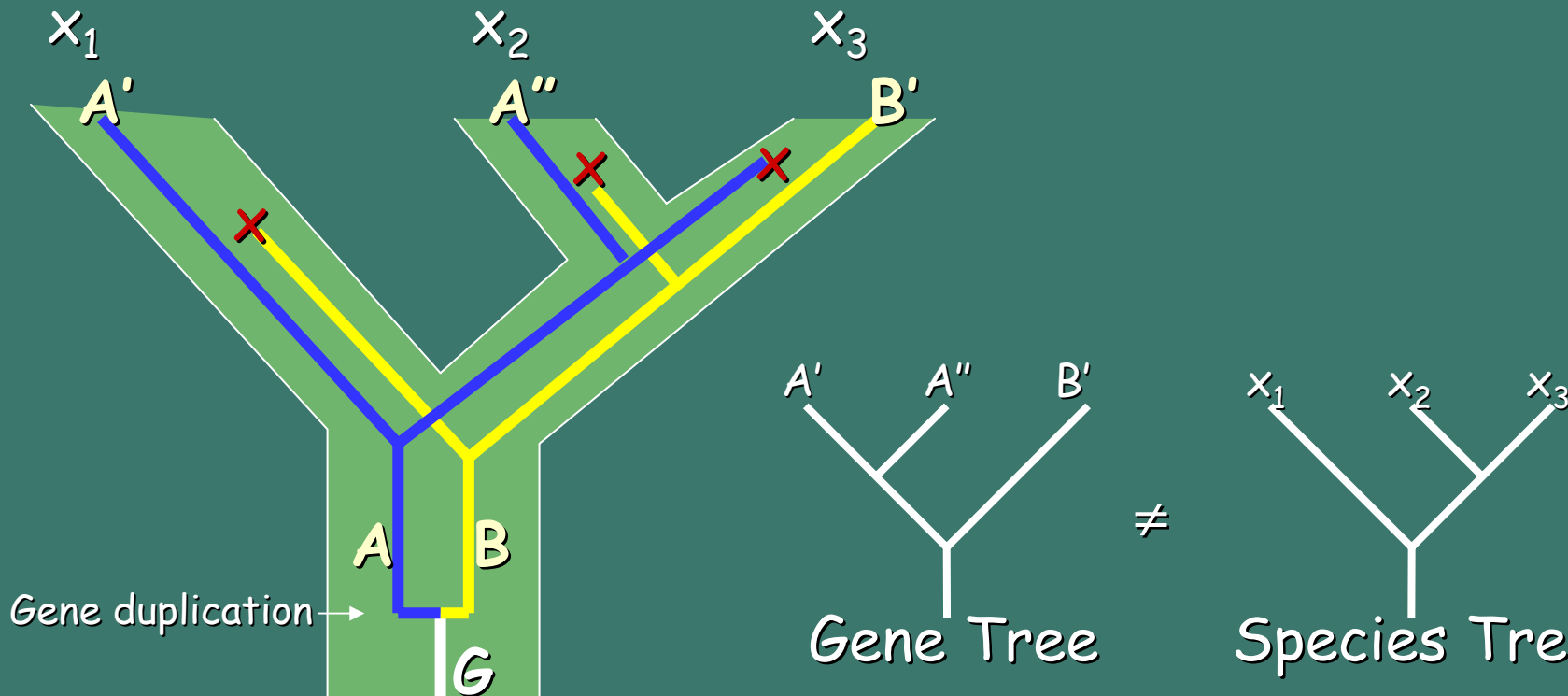
Split decomposition or Neighbor-Net produces network from distances

Part III

1. Phylogenetic trees
2. Splits networks
3. Consensus networks
4. Hybridization and reticulate networks
5. Recombination networks

Gene Trees Can Differ

Also allow gene duplication and loss:



Consensus of Different Gene Trees

- For a given set of species, different genes lead to different trees
- How to form a consensus of the trees?
 - Consensus trees
 - Consensus networks
 - Consensus super networks

The Split Encoding of a Tree

Tree T:



Split encoding $\Sigma(T)$:

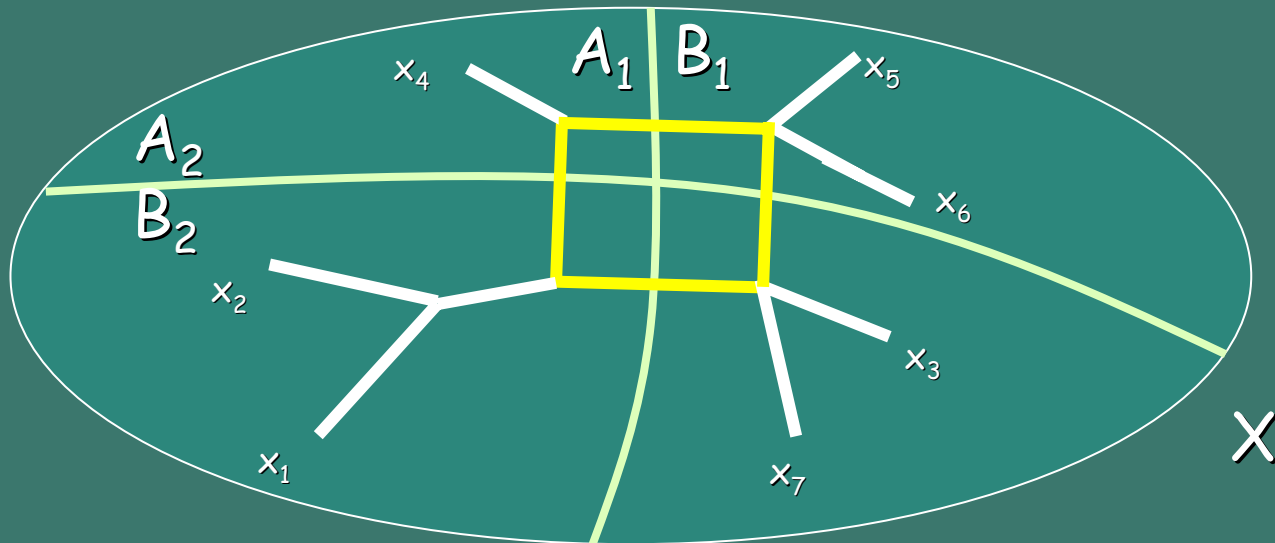
5 **trivial** splits: $\frac{\{a\}}{\{b, c, d, e\}}$, $\frac{\{b\}}{\{a, c, d, e\}}$, $\frac{\{c\}}{\{a, b, d, e\}}$, $\frac{\{d\}}{\{a, b, c, e\}}$ and $\frac{\{e\}}{\{a, b, c, d\}}$

2 **non-trivial** splits: $\frac{\{a, b, e\}}{\{c, d\}}$ and $\frac{\{a, b\}}{\{c, d, e\}}$.

Compatibility

- Two splits $A_1|B_1$ and $A_2|B_2$ of X are *compatible*, if
$$\emptyset \in \{A_1 \cap A_2, A_1 \cap B_2, B_1 \cap A_2, B_1 \cap B_2\}$$

Two **incompatible** splits:



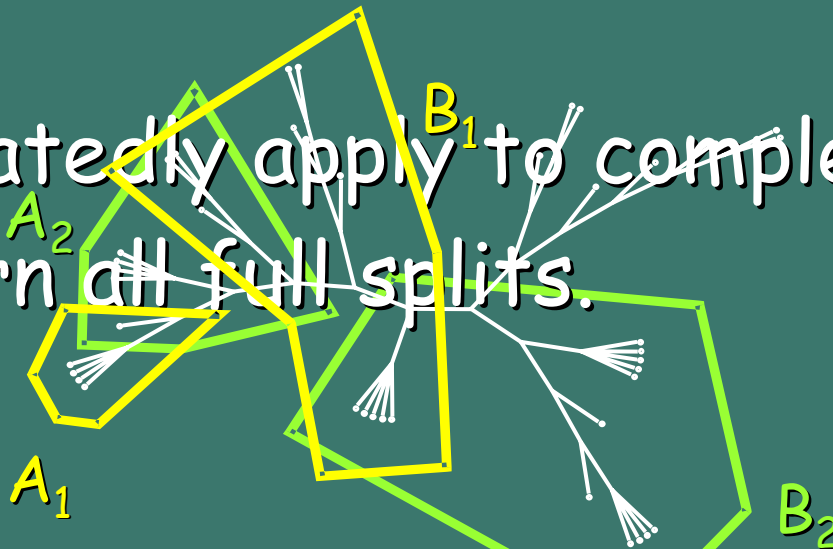
Z-Closure Method

- Idea: Extend partial splits.

- Z-rule:

$$\begin{array}{ccc}
 A_1 & \overset{n}{\text{---}} & A_2 \\
 \text{---} & \text{Z} & \text{---} \\
 B_1 & & B_2
 \end{array}
 \longrightarrow
 \begin{array}{cc}
 A_1 & A_1 \cup A_2 \\
 \text{---} & \text{---} \\
 B_1 \cup B_2 & B_2
 \end{array}$$

- Repeatedly apply to completion.
- Return all full splits.



Example

- Five fungal trees from [Pryor, 2000] and [Pryor, 2003]
- Trees:
 - ITS (two trees)
 - SSU (two trees)
 - Gpd (one tree)
- Numbers of taxa differ: "partial trees"

Example



SplitsTree 4
Beta

Written by:
Daniel H. Huson and David Bryant (2005)

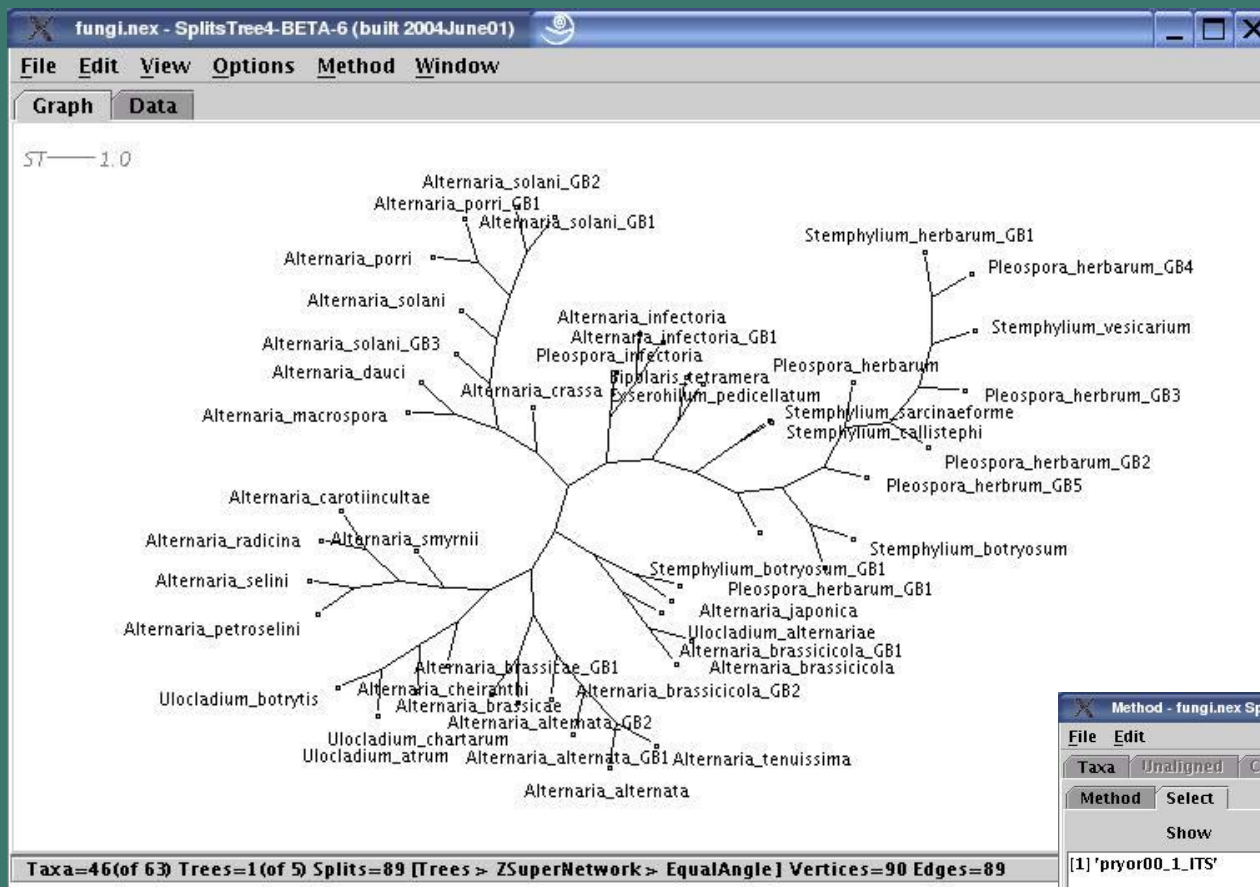
with contributions from:
Markus Franz, Migüel Jette', Tobias Klöpper
and Michael Schröder

www.splitstree.org

4beta25: User manual now available from www.splitstree.org!

Individual Gene Trees

ITS00



Method - fungi.nex SplitsTree4-BETA-6 (built 2004June01)

File Edit

Taxa Unaligned Characters Distances Quartets Trees Splits

Method Select

Show Hide

[1] 'pryor00_1_ITS'

< Show 'pryor03_1_ITS'

Hide > 'pryor03_2_Ssu'

'pryor03_3_gpd'

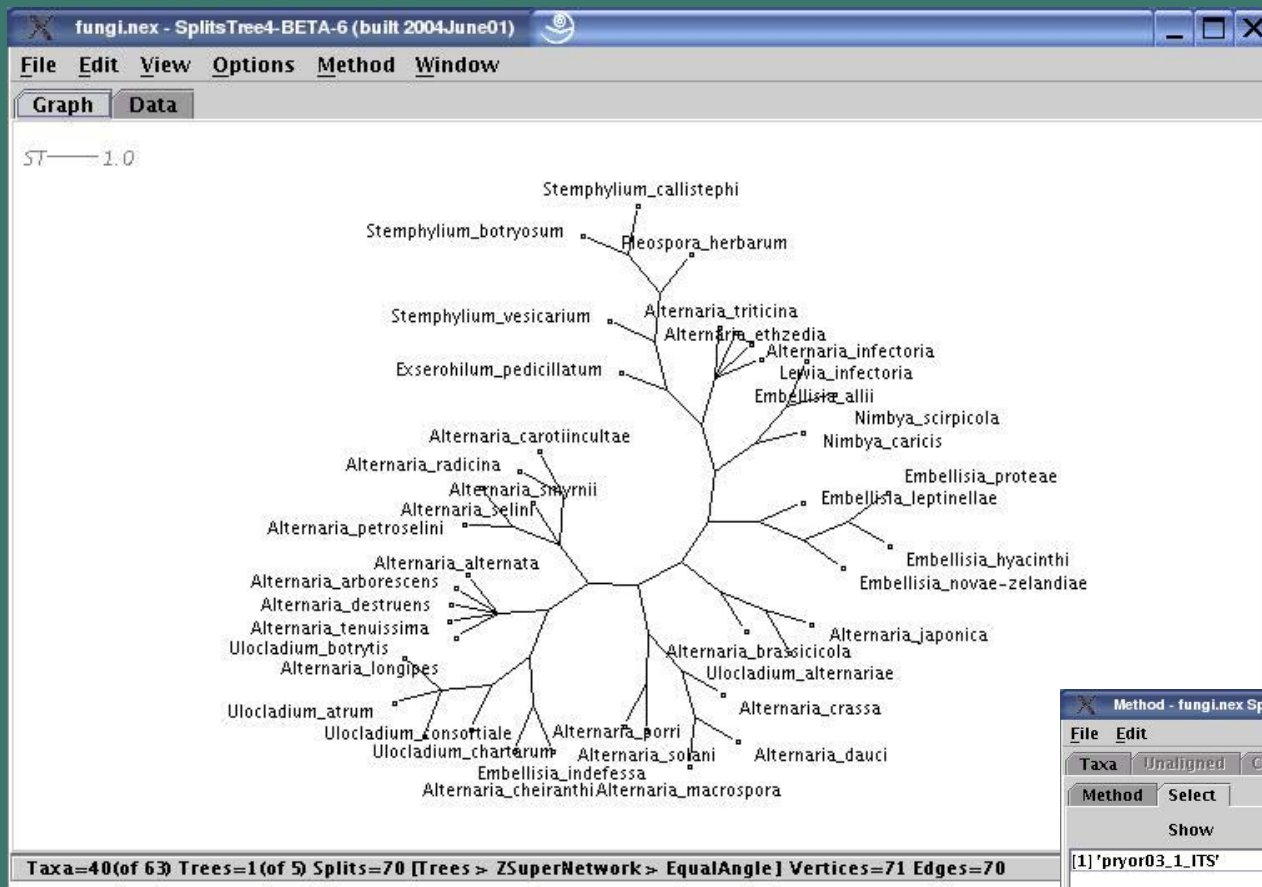
'pryor00_2_SsuDNA'

Show all

Hide all

46 taxa

Individual Gene Trees



ITS03

Method - fungi.nex SplitsTree4-BETA-6 (built 2004.June01)

File Edit

Taxa Unaligned Characters Distances Quartets Trees Splits

Method Select

Show Hide

[1] 'pryor03_1_ITS'

< Show 'pryor03_2_SSU'
'pryor03_3_gpd'
'pryor00_1_ITS'
'pryor00_2_SSurDNA'

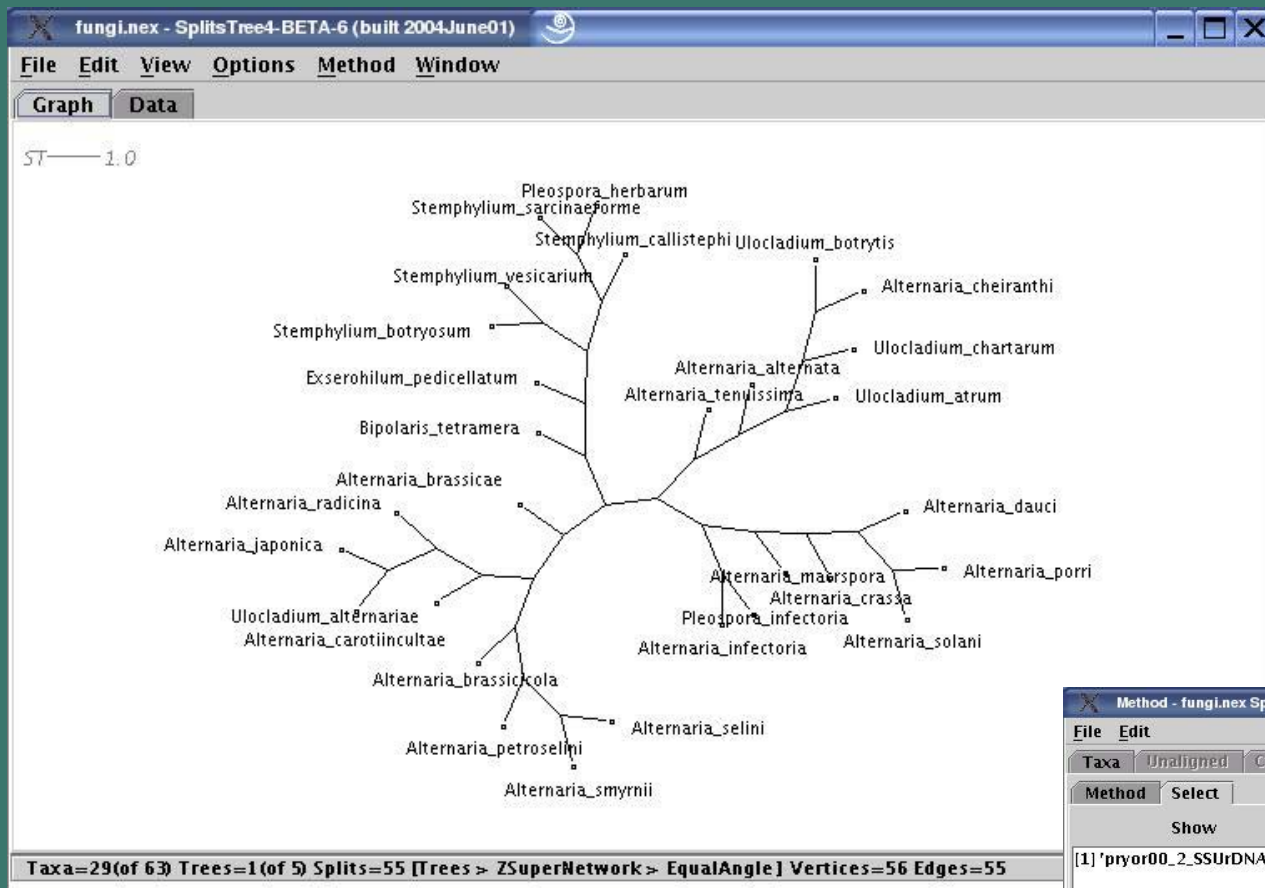
Hide >

Show all

Hide all

40 taxa

Individual Gene Trees



SSU00

Method - fungi.nex SplitsTree4-BETA-6 (built 2004June01)

File Edit

Taxa Unaligned Characters Distances Quartets Trees Splits

Method Select

Show Hide

[1] 'pryor00_2_SsUrDNA'

< Show 'pryor03_1_ITS'

Hide > 'pryor03_2_SsU'

'pryor03_3_gpd'

'pryor00_1_ITS'

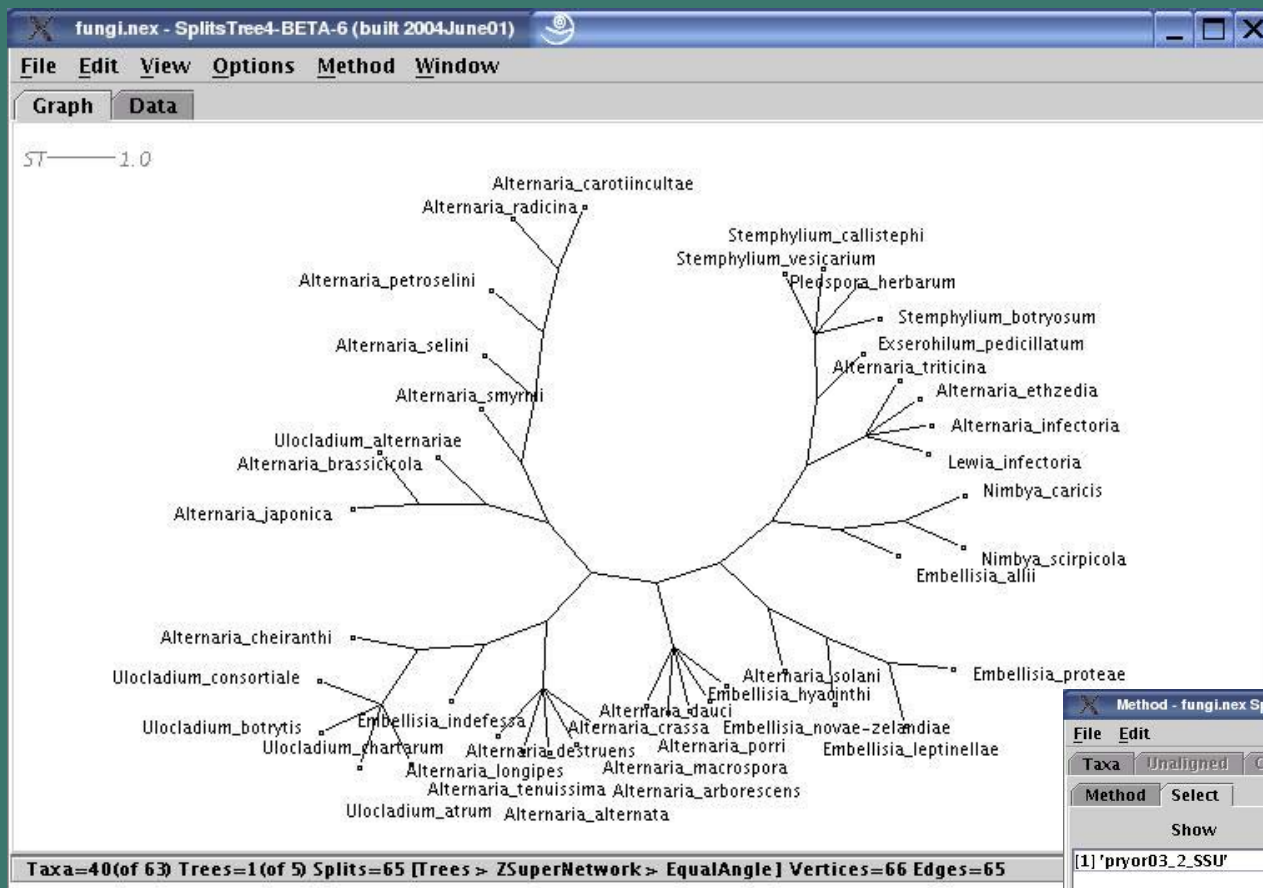
Show all

Hide all

29 taxa

Individual Gene Trees

SSU03



Method - fungi.nex SplitsTree4-BETA-6 (built 2004June01)

File Edit

Taxa Unaligned Characters Distances Quartets Trees Splits

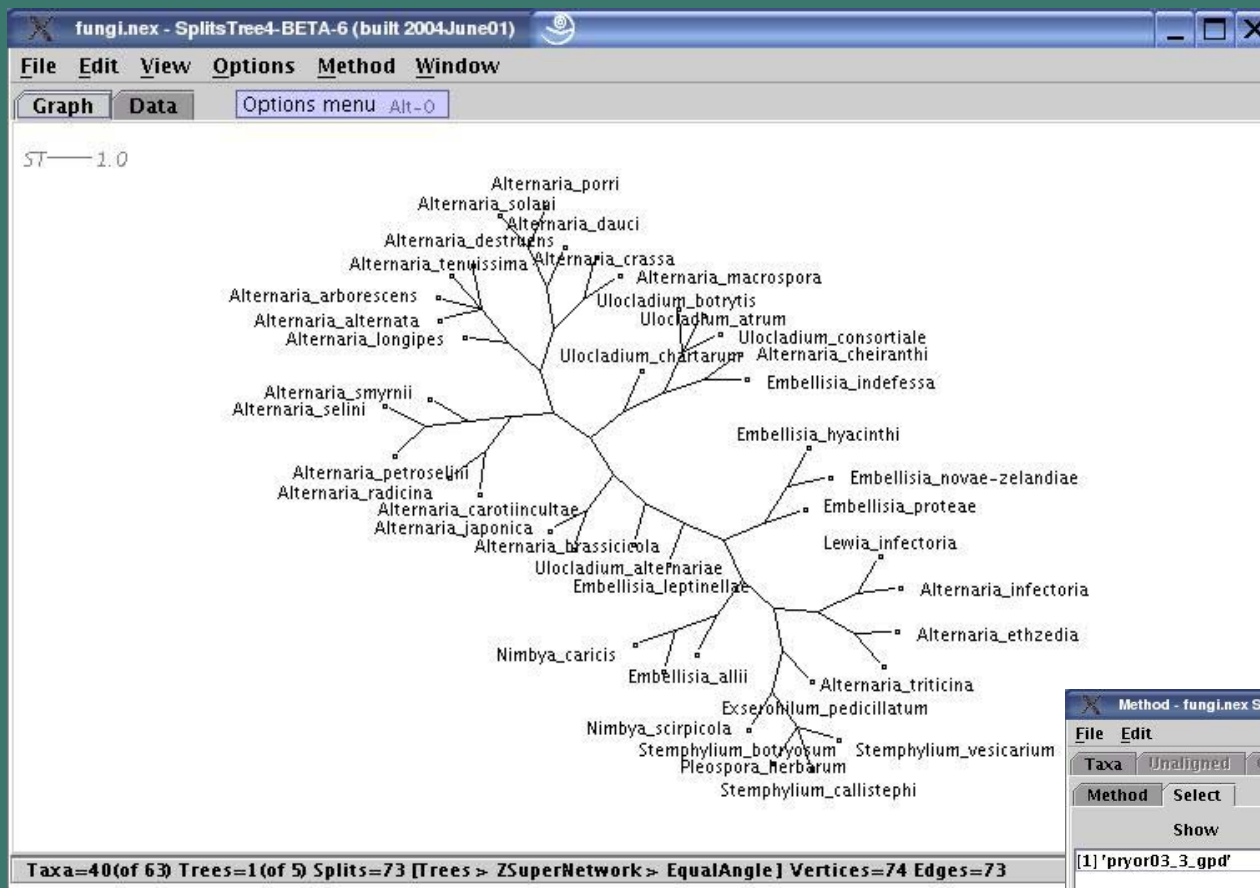
Method Select

Show Hide

[1] 'pryor03_2_Ssu'	< Show	'pryor03_1_ITS'
	Hide >	'pryor03_3_gpd'
	Show all	'pryor00_1_ITS'
	Hide all	'pryor00_2_SsuDNA'

40 taxa

Individual Gene Trees



Gpd03

40 taxa

Method - fungi.nex SplitsTree4-BETA-6 (built 2004June01)

File Edit

Taxa Unaligned Characters Distances Quartets Trees Splits

Method Select

Show	Hide	App
[1] 'pryor03_3_gpd'	'pryor03_1_ITS'	
	'pryor03_2_SSU'	
	'pryor00_1_ITS'	
	'pryor00_2_SSurDNA'	

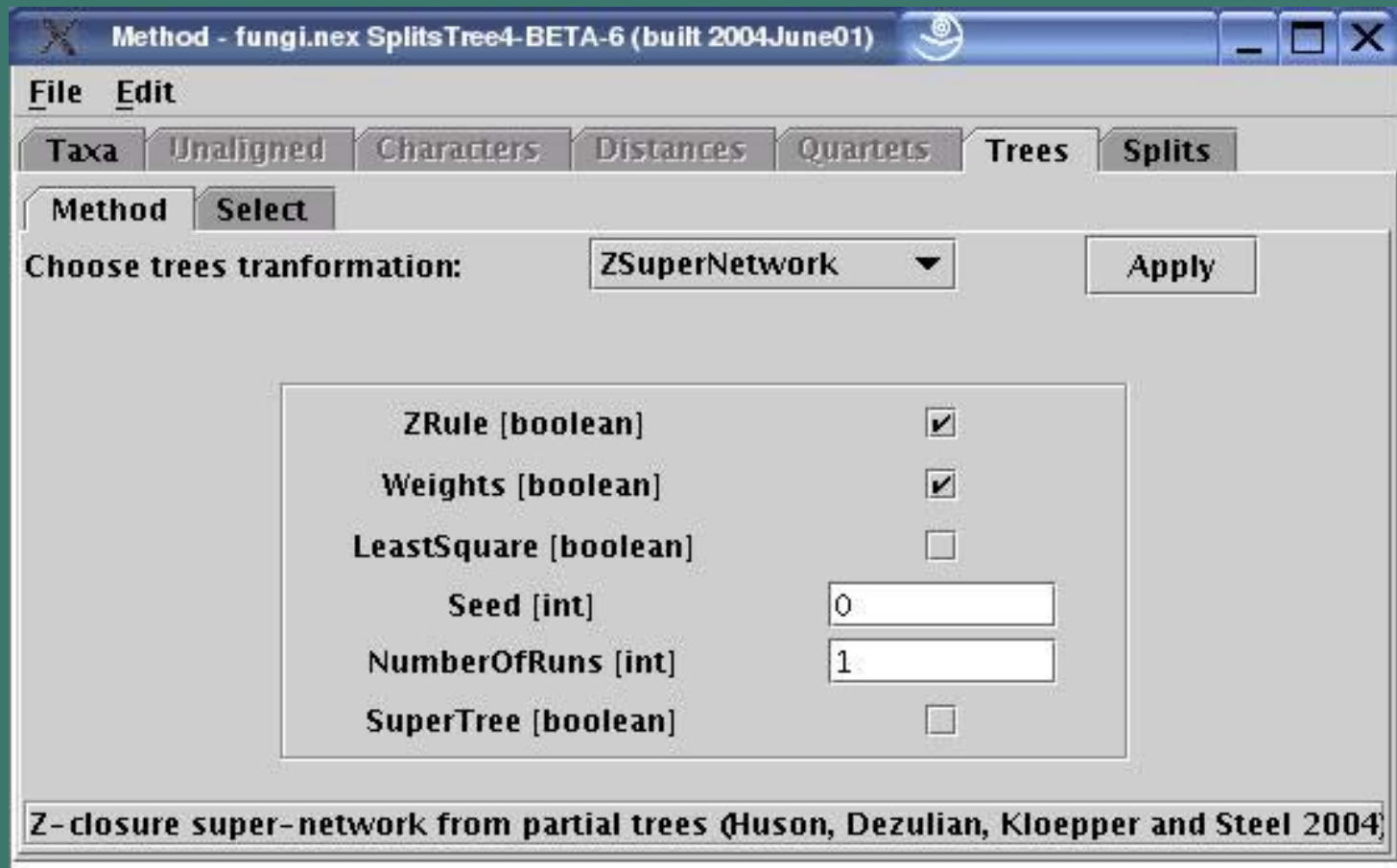
< Show

Hide >

Show all

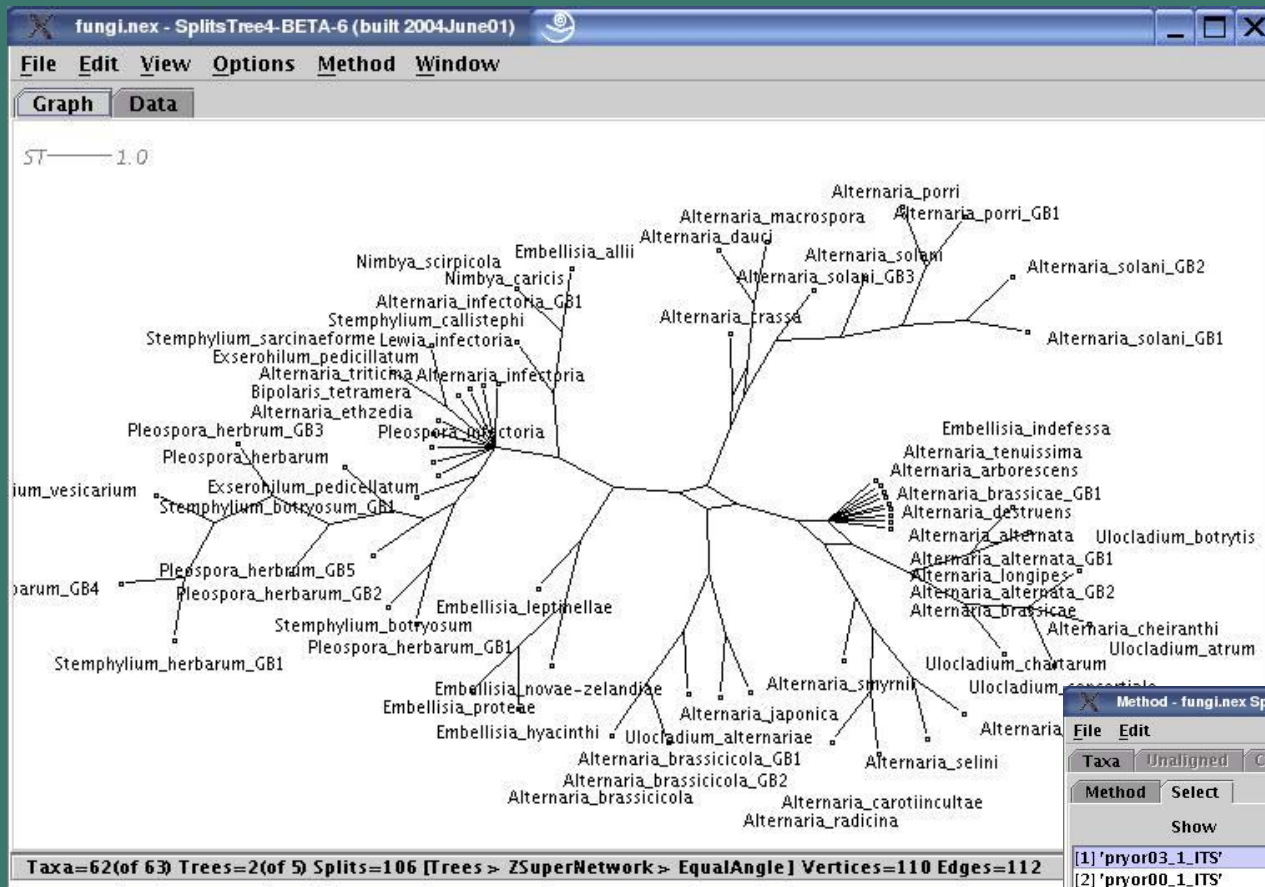
Hide all

Gene Trees as Super Network



Z-closure: a fast super-network method

Gene Trees as Super Network



ITS00+
ITS03

Method - fungi.nex SplitsTree4-BETA-6 (built 2004June01)

File Edit

Taxa Unaligned Characters Distances Quartets Trees Splits

Method Select

Show	Hide	App
[1] 'pryor03_1_ITS'	'pryor03_2_Ssu'	
[2] 'pryor00_1_ITS'	'pryor03_3_gpd'	
	'pryor00_2_SsuDNA'	

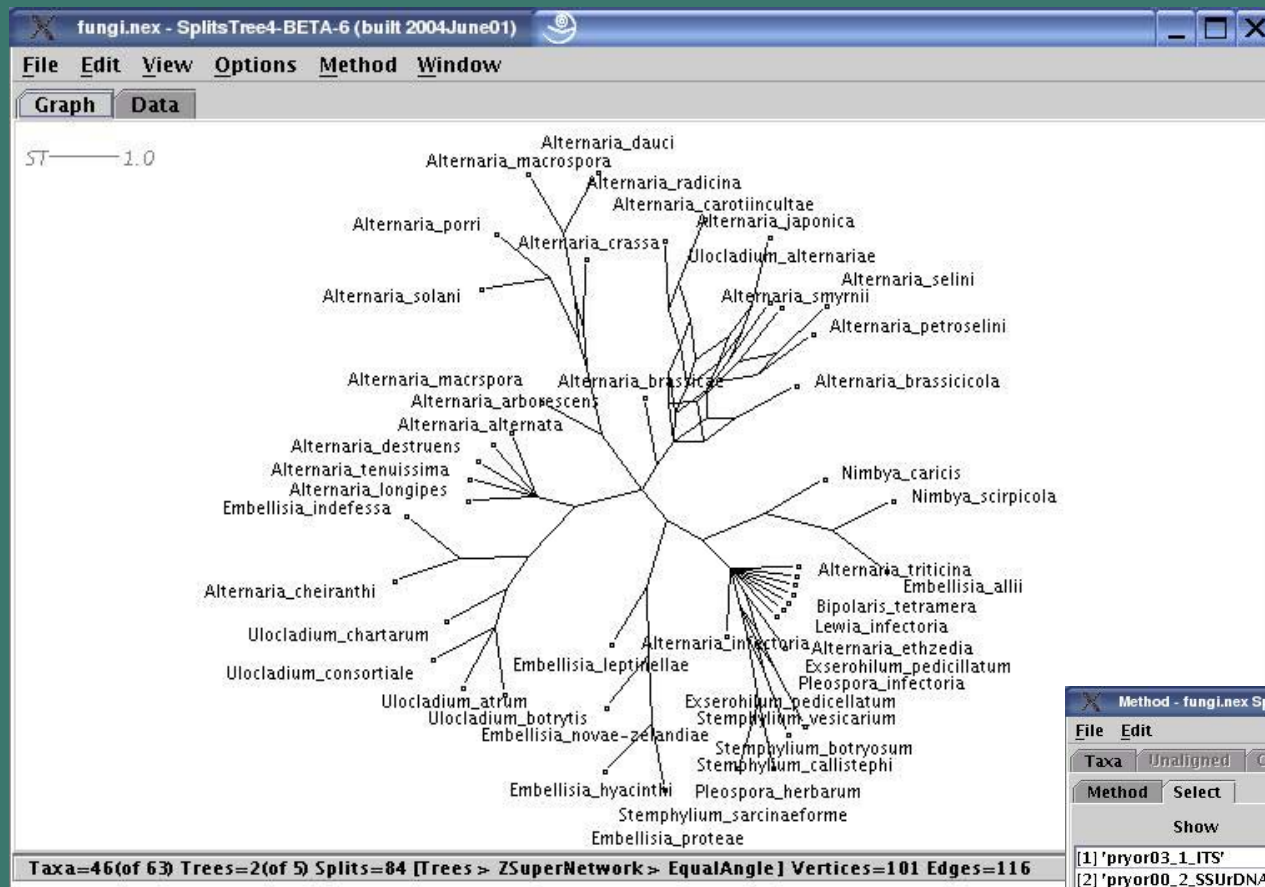
< Show

Hide >

Show all

Hide all

Gene Trees as Super Network



ITS03+
SSU00

Method - fungi.nex SplitsTree4-BETA-6 (built 2004June01)

File Edit

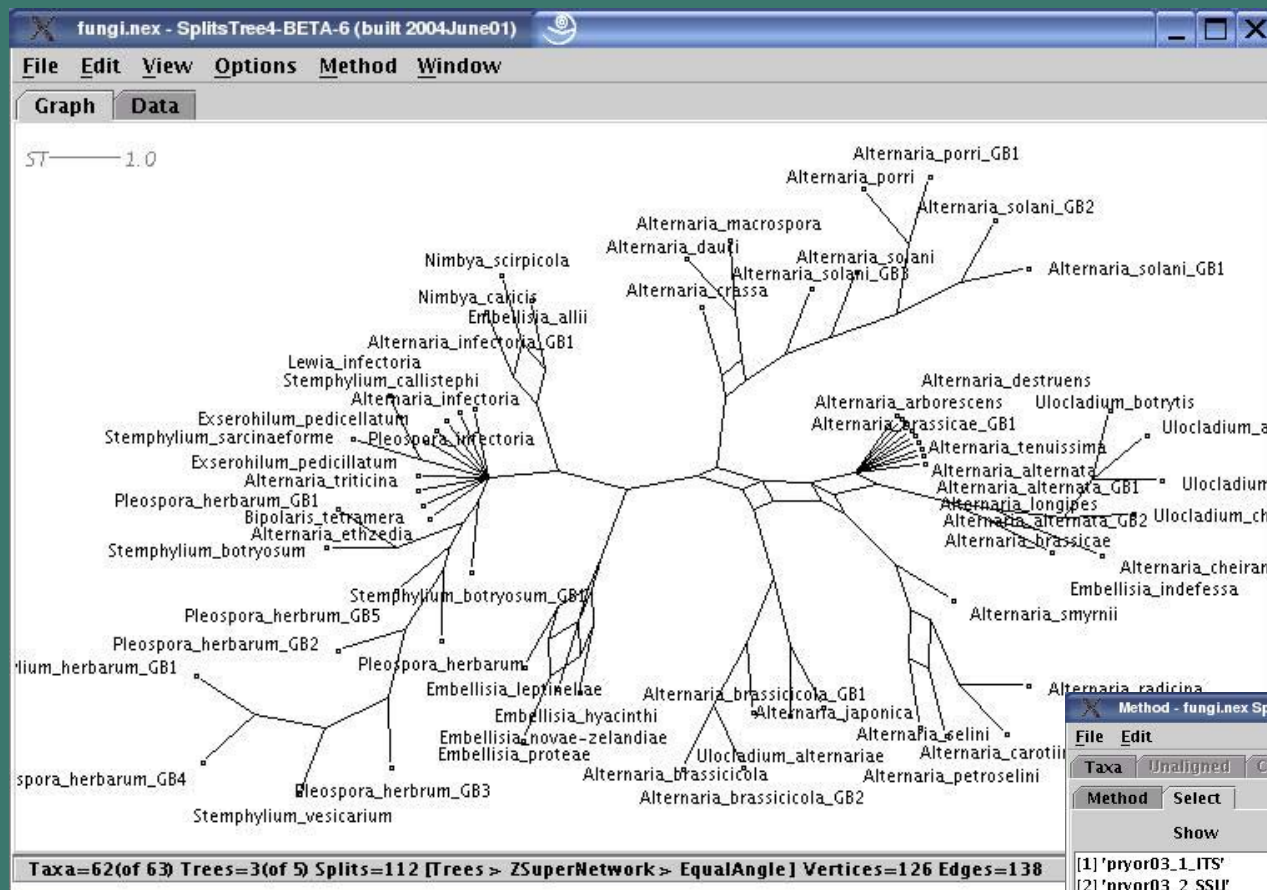
Taxa Unaligned Characters Distances Quartets Trees Splits

Method Select

Show Hide

[1] 'pryor03_1_ITS'	< Show	'pryor03_2_SSU'
[2] 'pryor00_2_SSuDNA'	Hide >	'pryor03_3_gpd'
	Show all	'pryor00_1_ITS'
	Hide all	

Gene Trees as Super Network



ITS00+
ITS00+
SSU03

Method - fungi.nex SplitsTree4-BETA-6 (built 2004June01)

File Edit

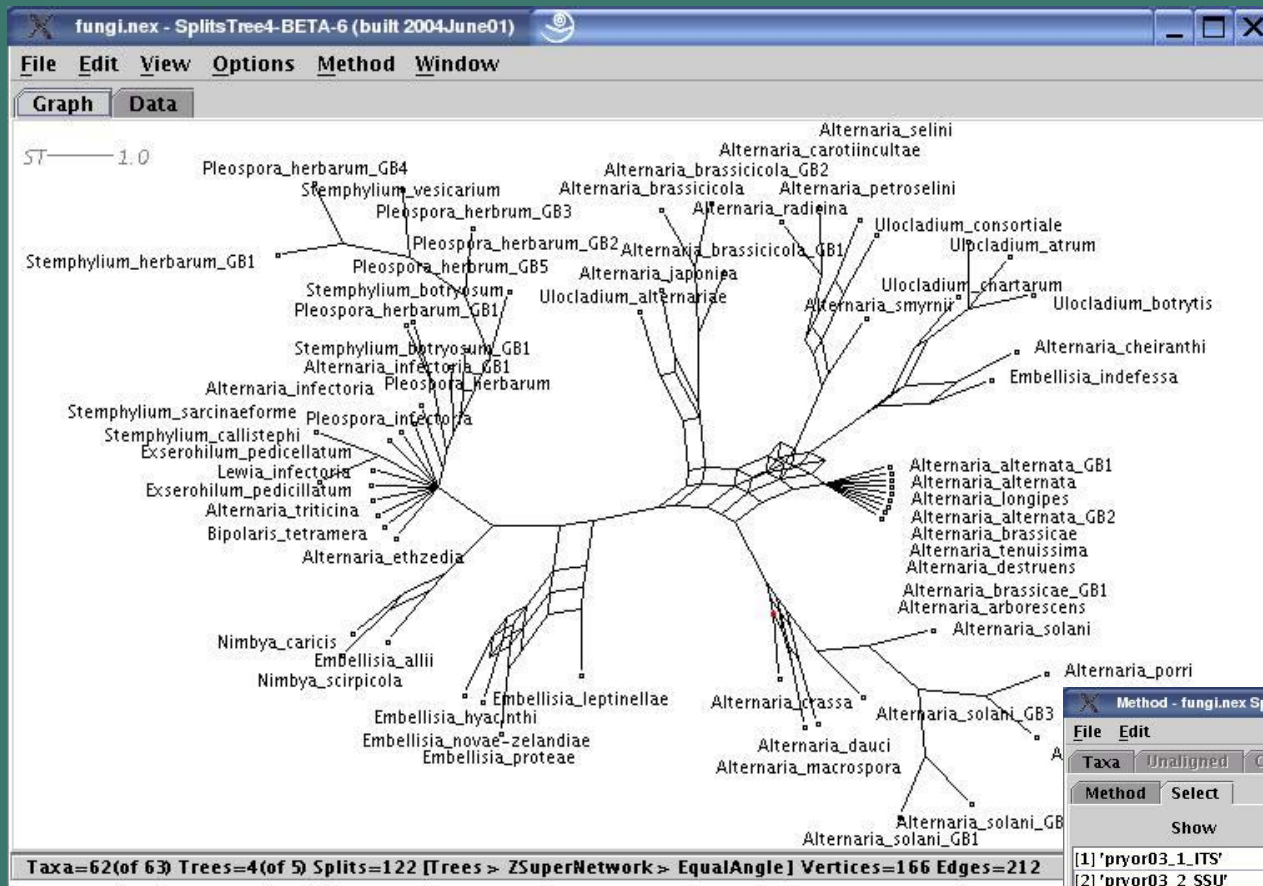
Taxa Unaligned Characters Distances Quartets Trees Splits

Method Select

Show	Hide
[1] 'pryor03_1_ITS'	'pryor03_3_gpd'
[2] 'pryor03_2_SSU'	'pryor00_2_SSuDNA'
[3] 'pryor00_1_ITS'	

< Show Hide > Show all Hide all

Gene Trees as Super Network



ITS00+
ITS03+
SSU03+
Gpd03

Method - fungi.nex SplitsTree4-BETA-6 (built 2004June01)

File Edit

Taxa Unaligned Characters Distances Quartets Trees Splits

Method Select

Show	Hide	App
[1] 'pryor03_1_ITS'	'pryor00_2_SSuDNA'	
[2] 'pryor03_2_SSU'		
[3] 'pryor03_3_gpd'		
[4] 'pryor00_1_ITS'		

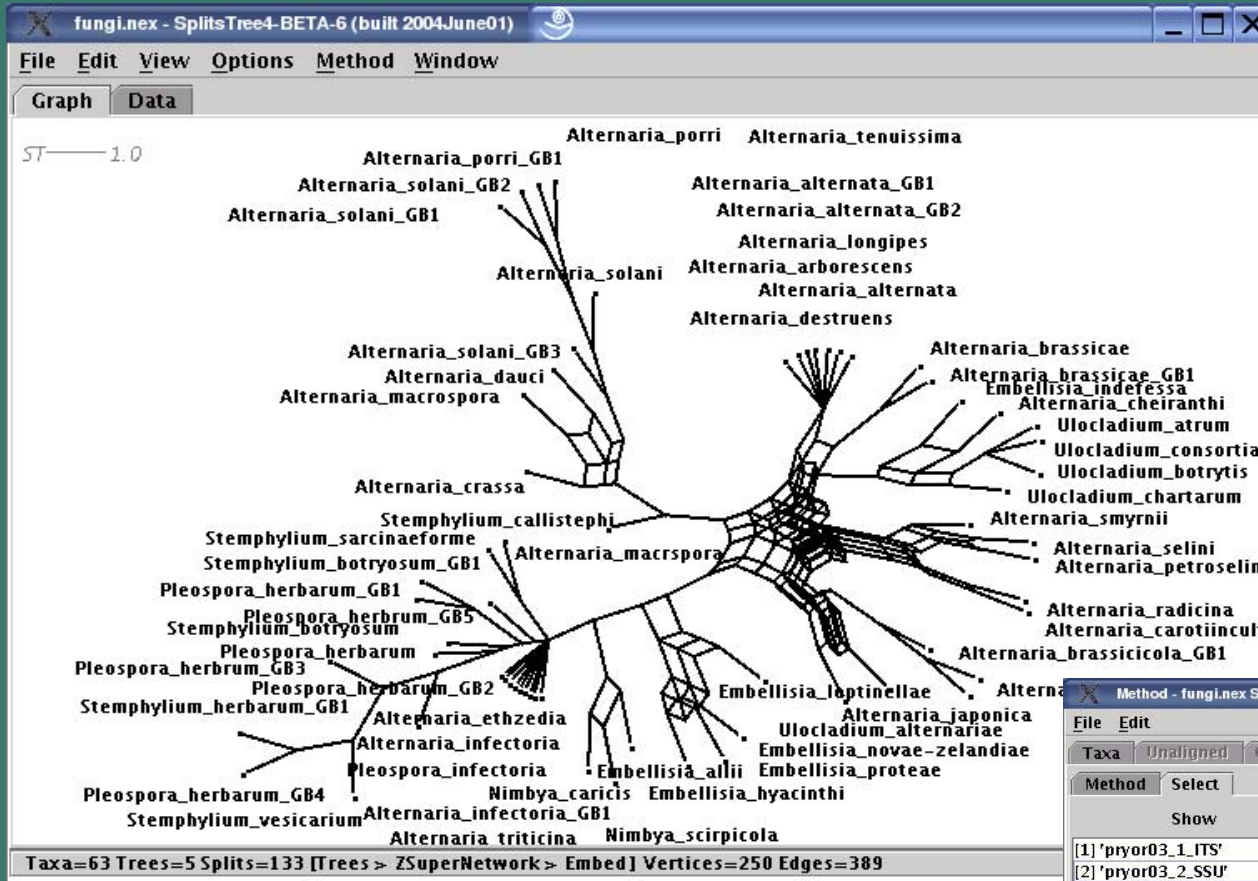
< Show

Hide >

Show all

Hide all

Gene Trees as Super Network



ITS00+
ITS03+
SSU00+
SSU03+
Gpd03



Part IV

1. Phylogenetic trees
2. Splits networks
3. Consensus networks
4. Hybridization and reticulate networks
5. Recombination networks

Hybridization

- Occurs when two organisms from different species interbreed and combine their genomes



Copyright © 2003 University of Illinois

Water hemp



Copyright © 2003 University of Illinois

Hybrid

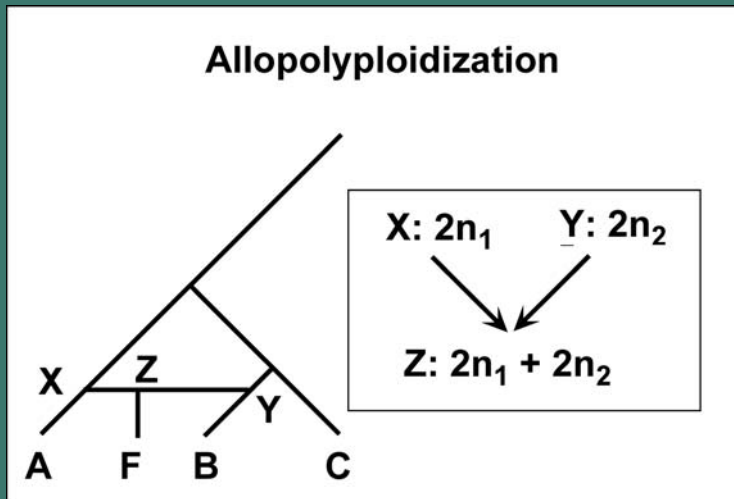


Copyright © 2003 University of Illinois

Pigs weed

Speciation by Hybridization 1

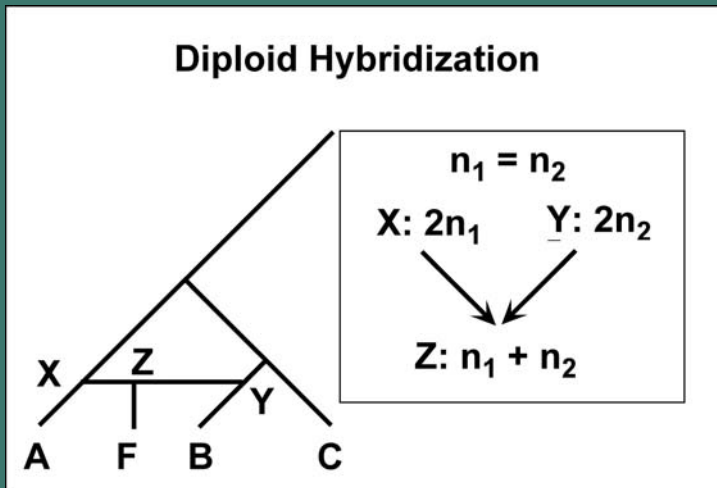
- In allopolyploidization, two different lineages produce a new species that has the complete nuclear genomes of both parental species:



Linder *et al.* 2004

Speciation by Hybridization 2

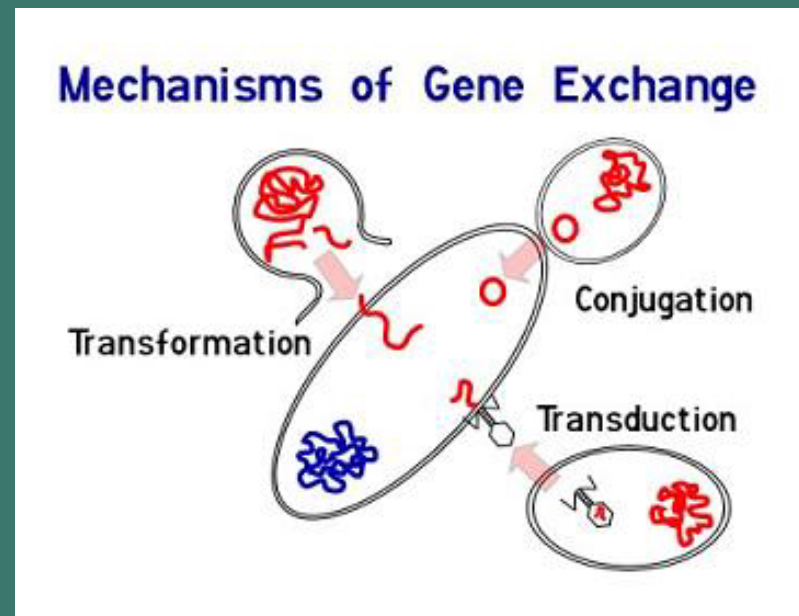
- In diploid (or homoploid) hybrid speciation, each of the parents produces normal gametes (haploid) to produce a normal diploid hybrid:



Linder *et al.*, 2004

Horizontal Gene Transfer

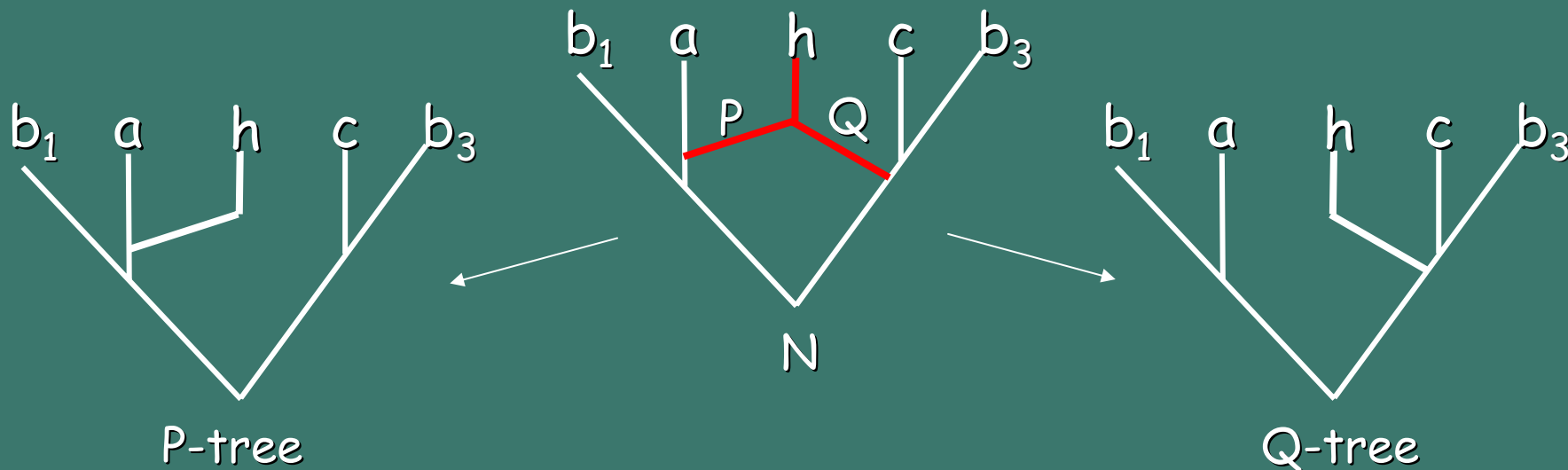
- There are a number of known mechanisms by which bacteria can exchange genes
 - Transformation
 - Conjugation
 - transduction



<http://www.pitt.edu/~heh1/research.html>

Reticulate Networks and Trees

- The evolutionary history associated with any given gene is a **tree**
- A network N with k reticulations gives rise to 2^k different gene trees



Rooted Reticulate Network

Definition Let X be a set of taxa. A rooted *reticulate network* N on X is a connected, directed acyclic graph with:

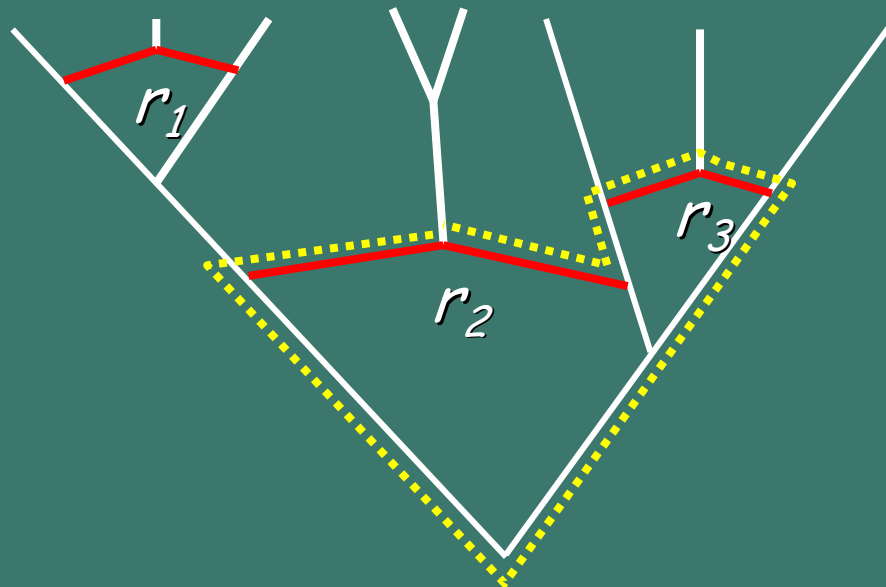
- precisely one node of indegree 0, the *root*,
- all other nodes are *tree nodes* of indegree 1, or *reticulation nodes* of indegree 2,
- every edge is a *tree edge* joining two tree nodes, or a *reticulation edge* from a tree node to a reticulation node, and
- the set of leaves consists of tree nodes and is labeled by X .

Most Parsimonious Network Problem:

- Given a set of trees \mathcal{T} , determine a reticulate network N such that
$$\mathcal{T} \subseteq \mathcal{T}(N)$$
and N contains a minimum number of reticulation nodes.
- In fully generality, this is known to be a computationally hard problem
[Wang *et al* 2001, Bordewich and Semple 2004].

Independent Reticulations

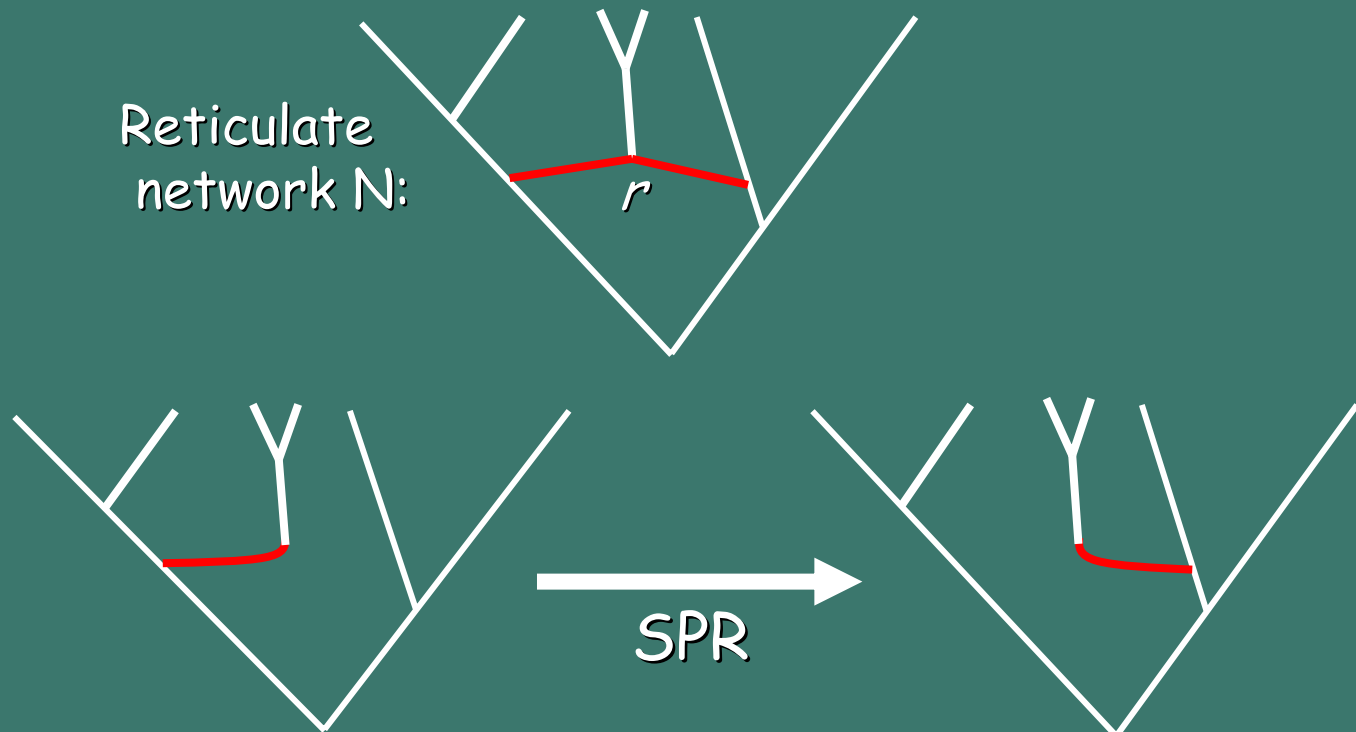
- Reticulation nodes $r_i, r_j \in N$ are *independent*, if they are not contained in a common cycle:



- Independent reticulations also called *galls* and a network only containing galls is also called a *galled tree* [Gusfield et al. 2003]

SPR's and Reticulations

Observation [Maddison 1997]: If N contains only one reticulation r , then it corresponds to a "sub-tree prune and regraft" operation:

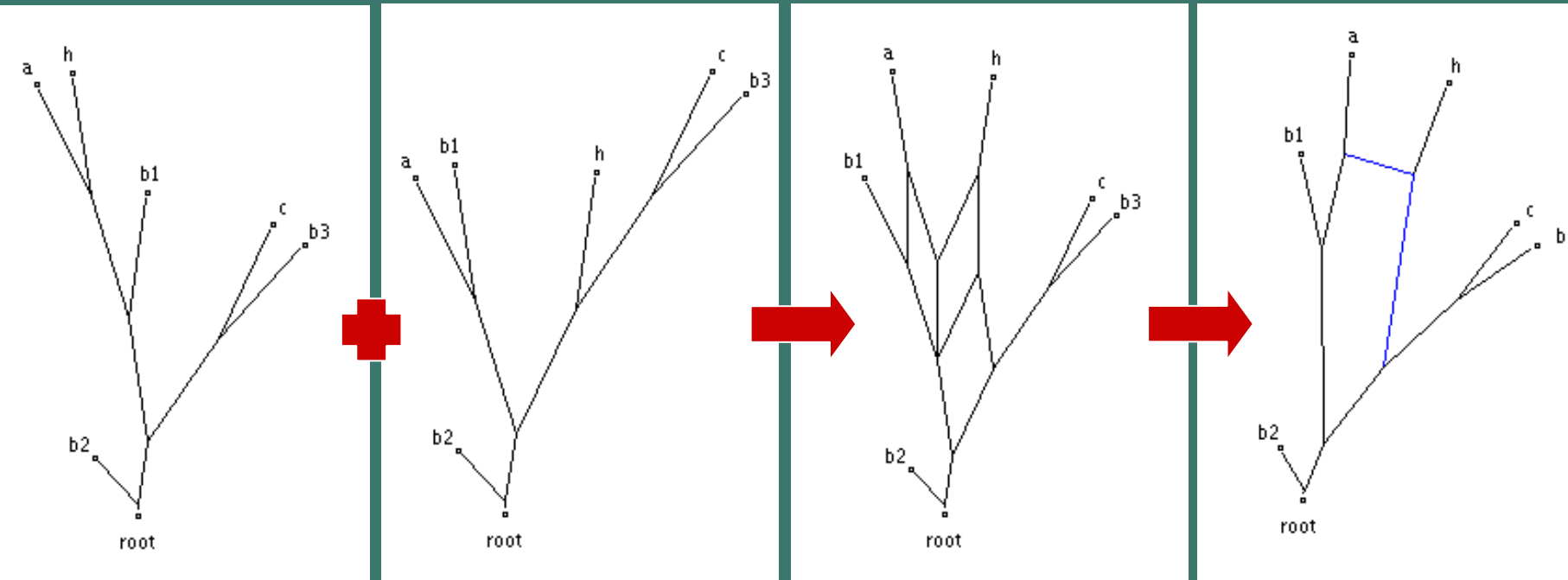


SPR-Based Algorithm

- Given two bifurcating trees, compute their SPR distance:
 - If $f = 0$, return the tree
 - If $f = 1$, return the reticulate network
 - Else, return fail
- Generalized to networks with multiple independent reticulations [Nakhleh *et al*/2004]
- Maximum agreement forest approach (Semple *et al*/2005)

Splits-Based Approach

A new splits-based approach [Huson, Kloepper, Lockhart and Steel 2005]



gene tree1

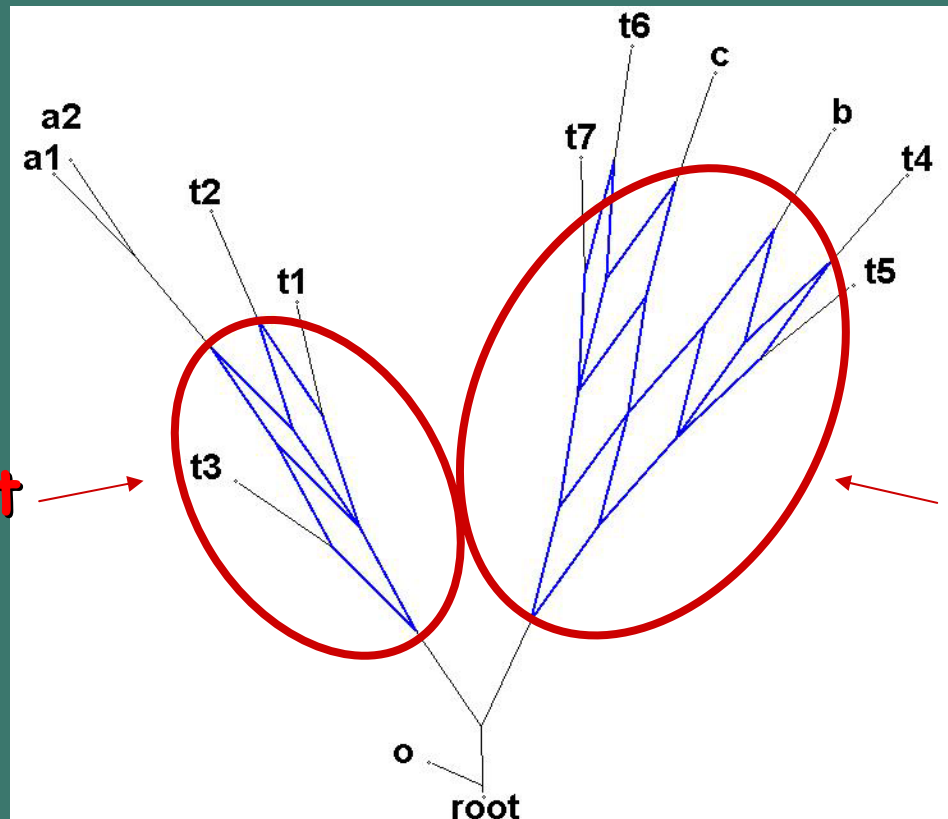
gene tree2

splits network
of all splits

reticulate
network

Decomposition Theorem

- Each incompatibility component can be considered independently:

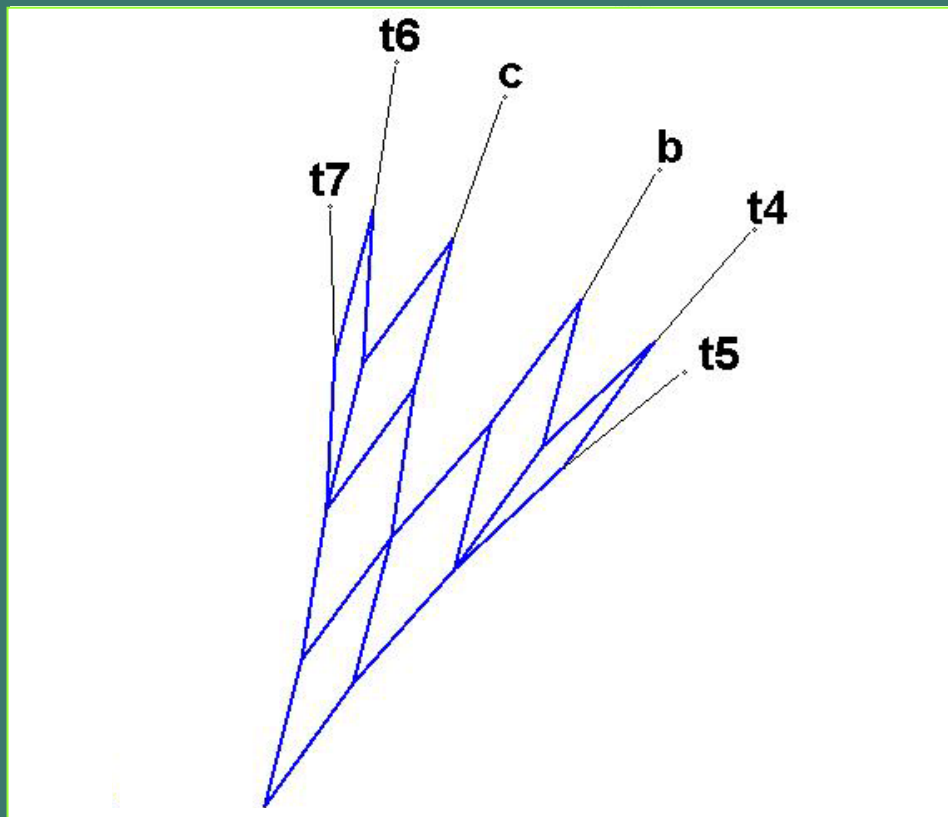


1. component

2. component

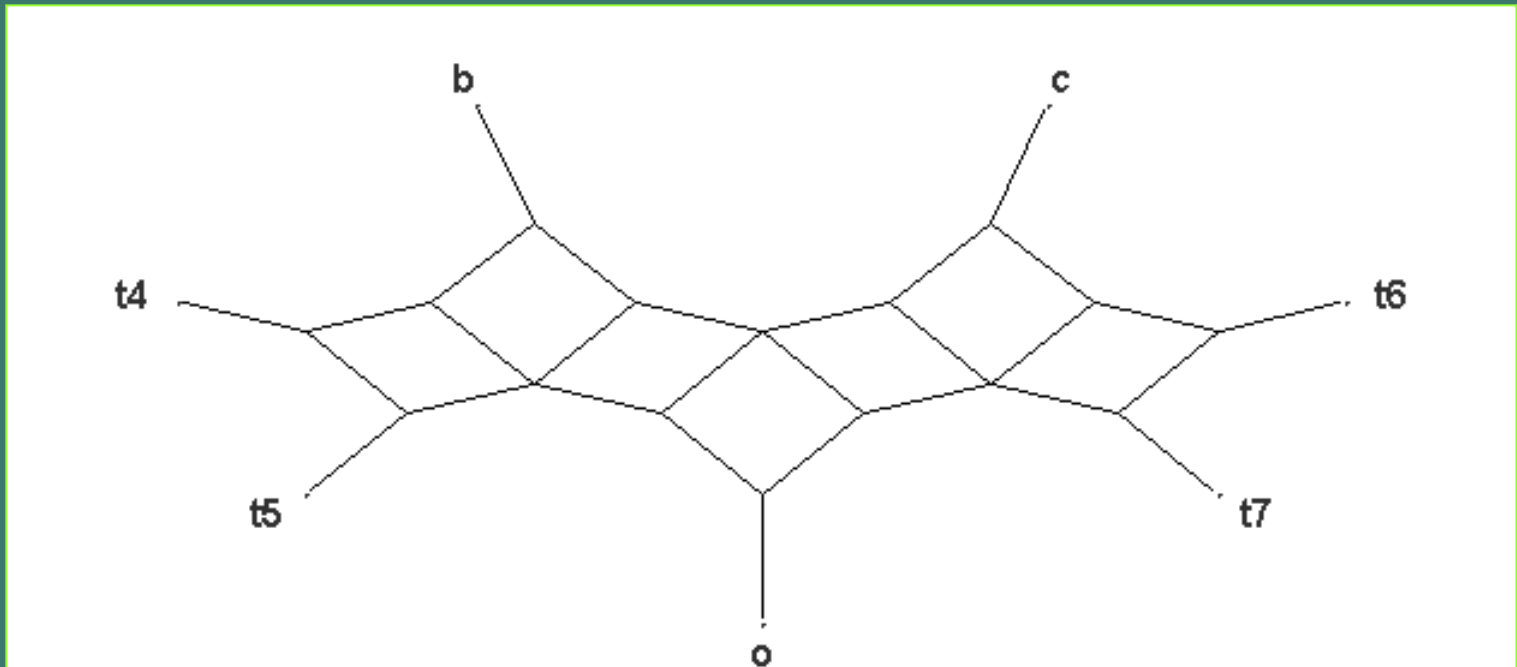
Decomposition Theorem

- Consider a component:



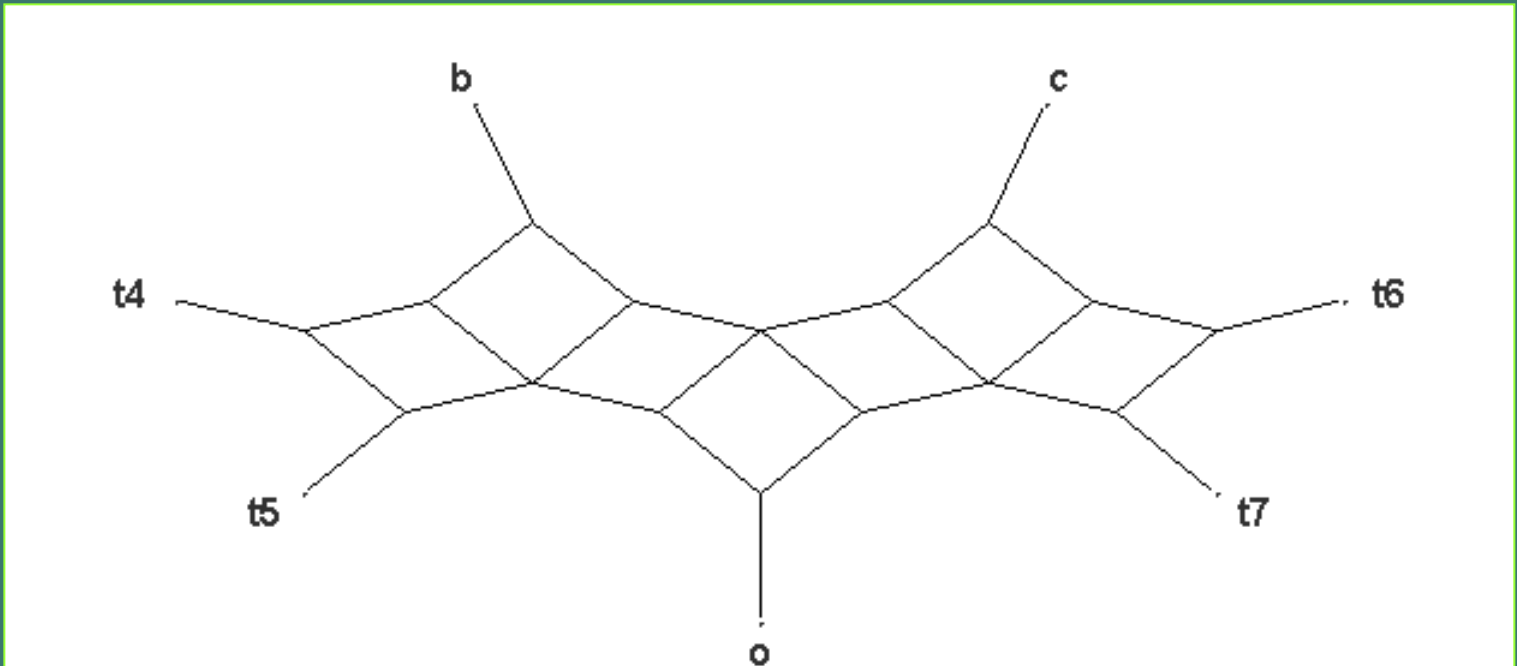
Algorithm

- Find decomposition $R \cup B$ as a set of reticulate taxa and backbone taxa



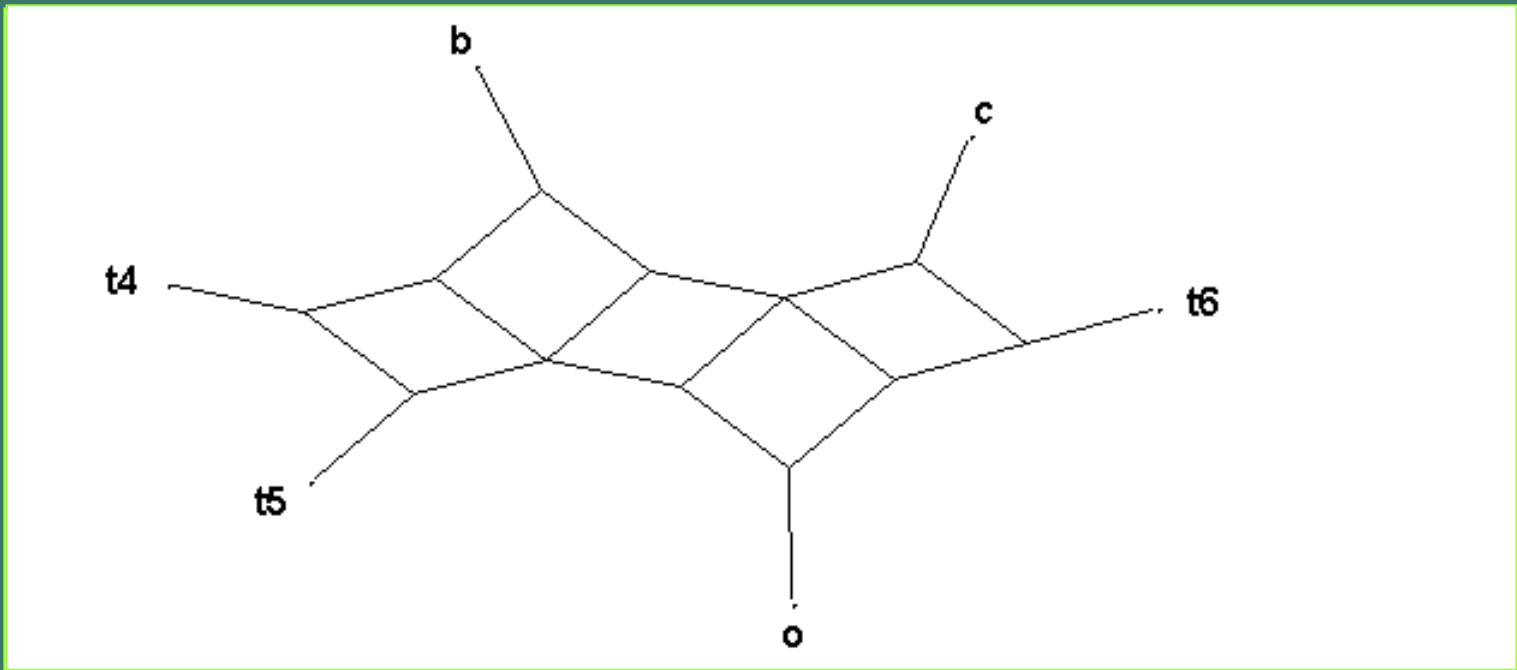
Algorithm

- Necessary condition: splits restricted to B must correspond to a tree



Algorithm

- Consider all choices for R of size 1
[Gusfield *et al.*, 2003, 2004]:

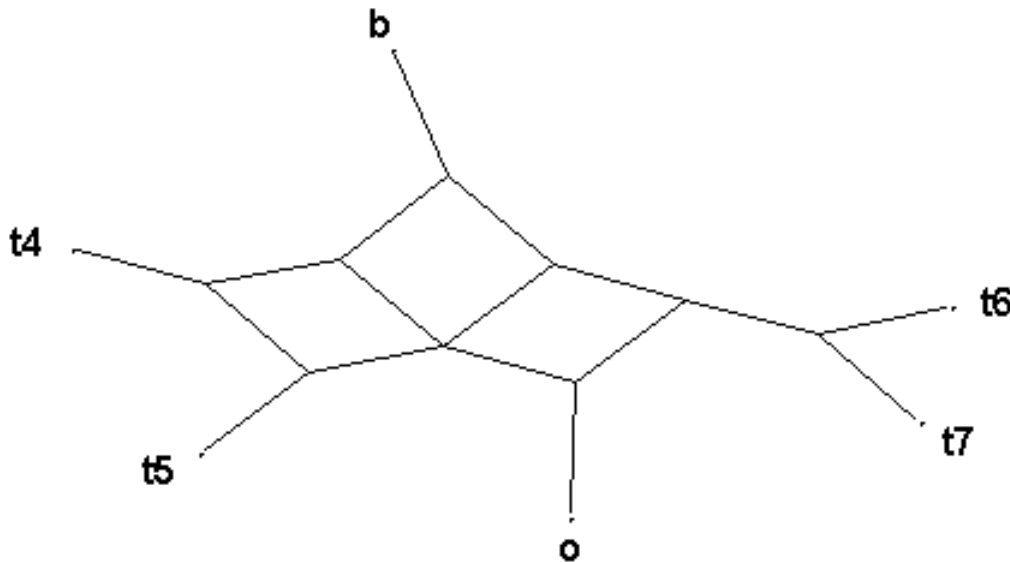


$R = \{t_7\}$

☹️ ...not a tree, R not good

Algorithm

- Consider all choices for R of size 1
[Gusfield *et al.*, 2003, 2004]:

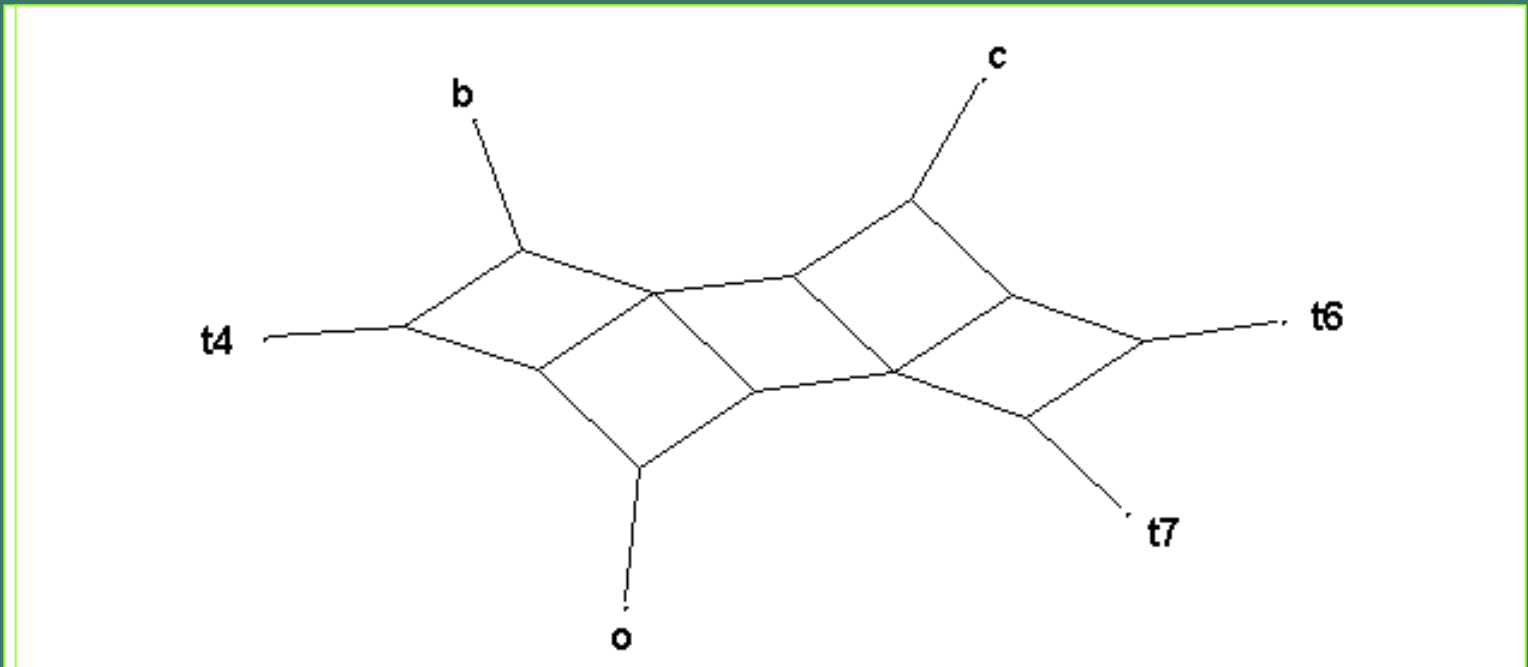


$R=\{c\}$

☹️ ...not a tree, R not good

Algorithm

- Consider all choices for R of size 1
[Gusfield *et al.*, 2003, 2004]:

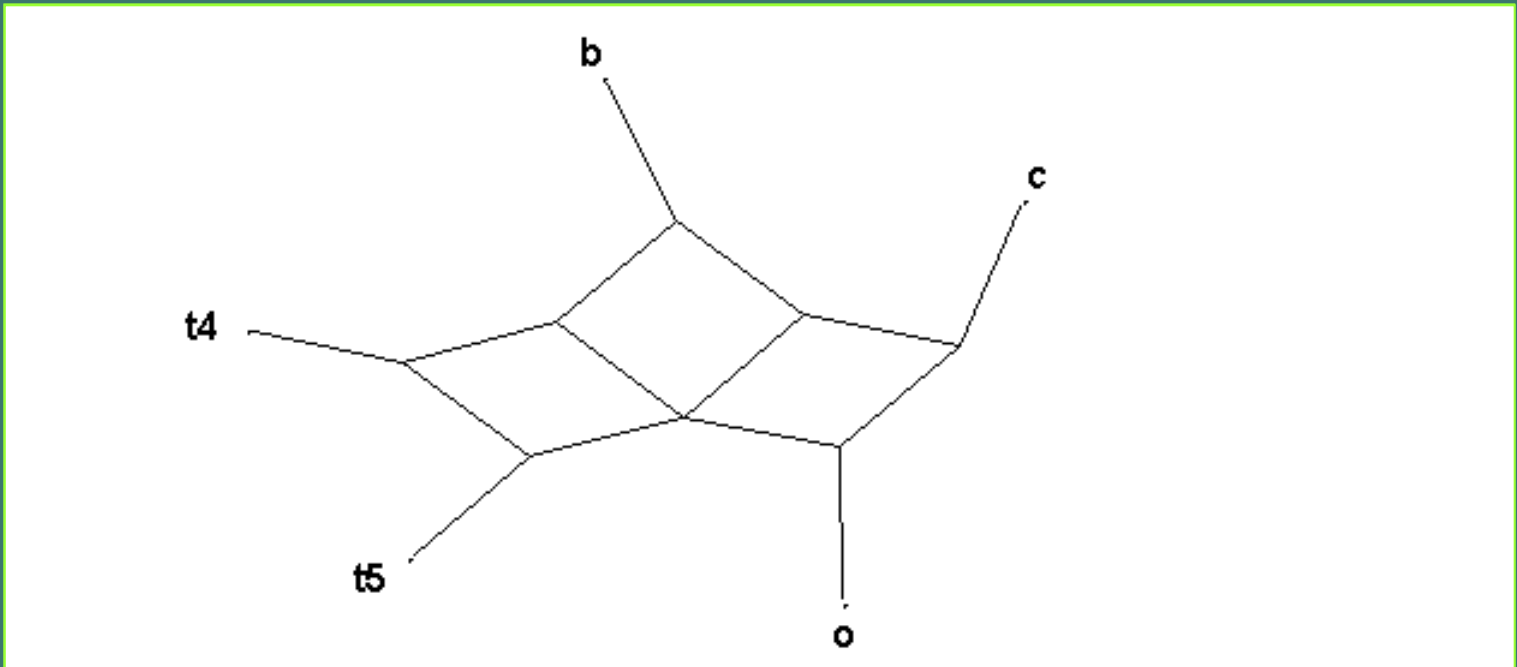


$R = \{t_5\}$

☹ ...not a tree, R not good

Algorithm

- Consider all choices for R of size 2:
[H., Klopper, Lockhart and Steel, 2005]

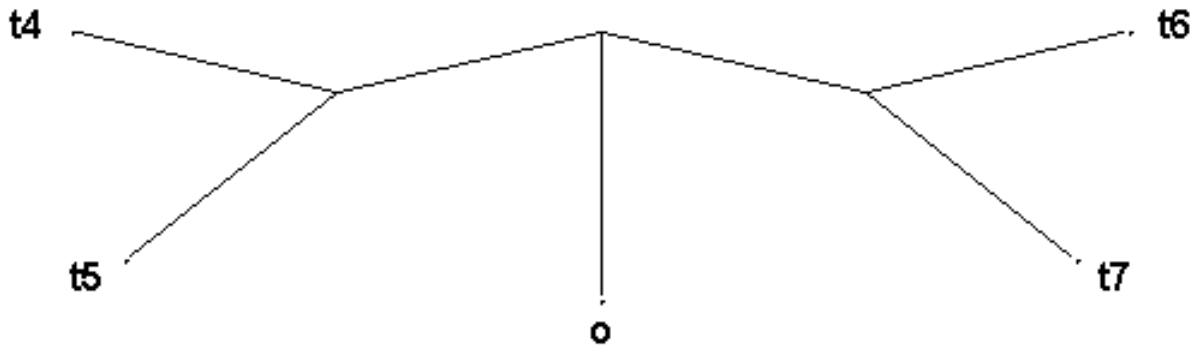


$R = \{t_6, t_7\}$

☹️ ...not a tree, R not good

Algorithm

- Consider all choices for R of size 2:
[H., Klopper, Lockhart and Steel, 2005]

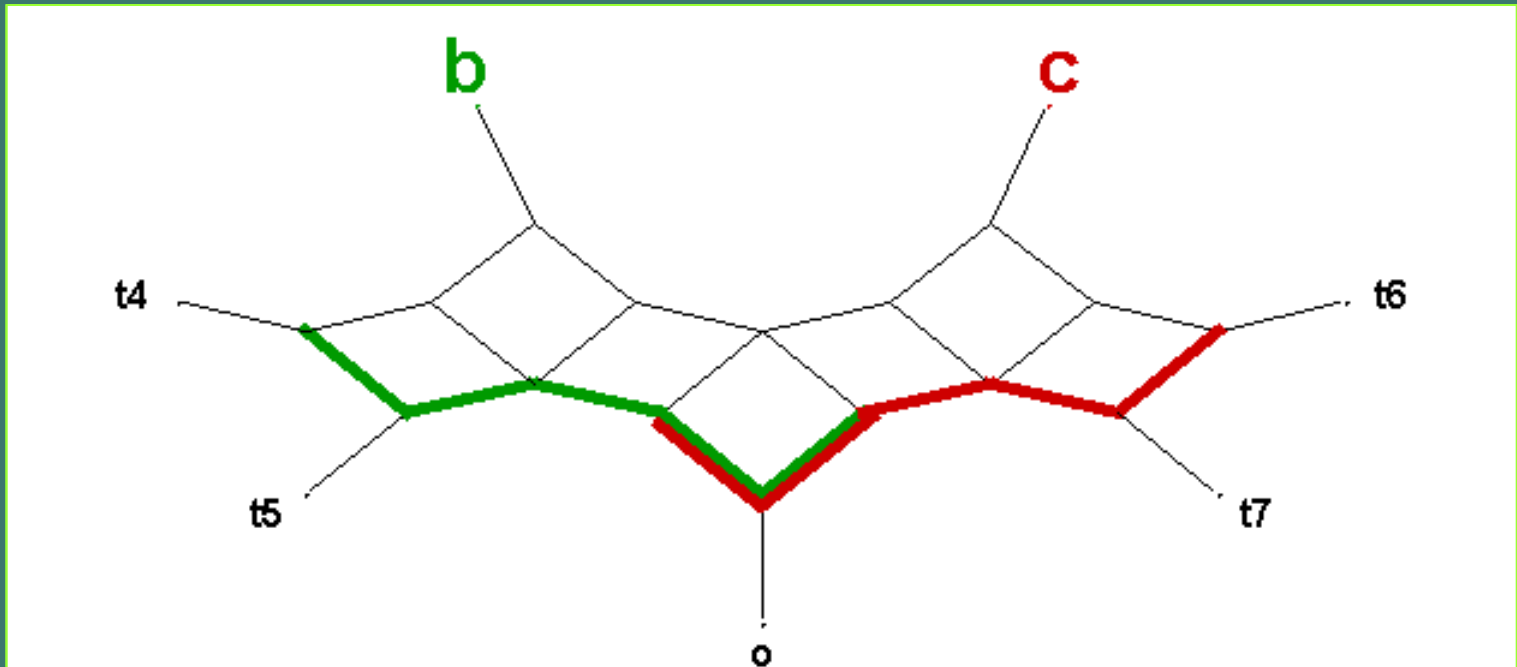


$R = \{b, c\}$

😊 ...is a **tree**, R is a candidate

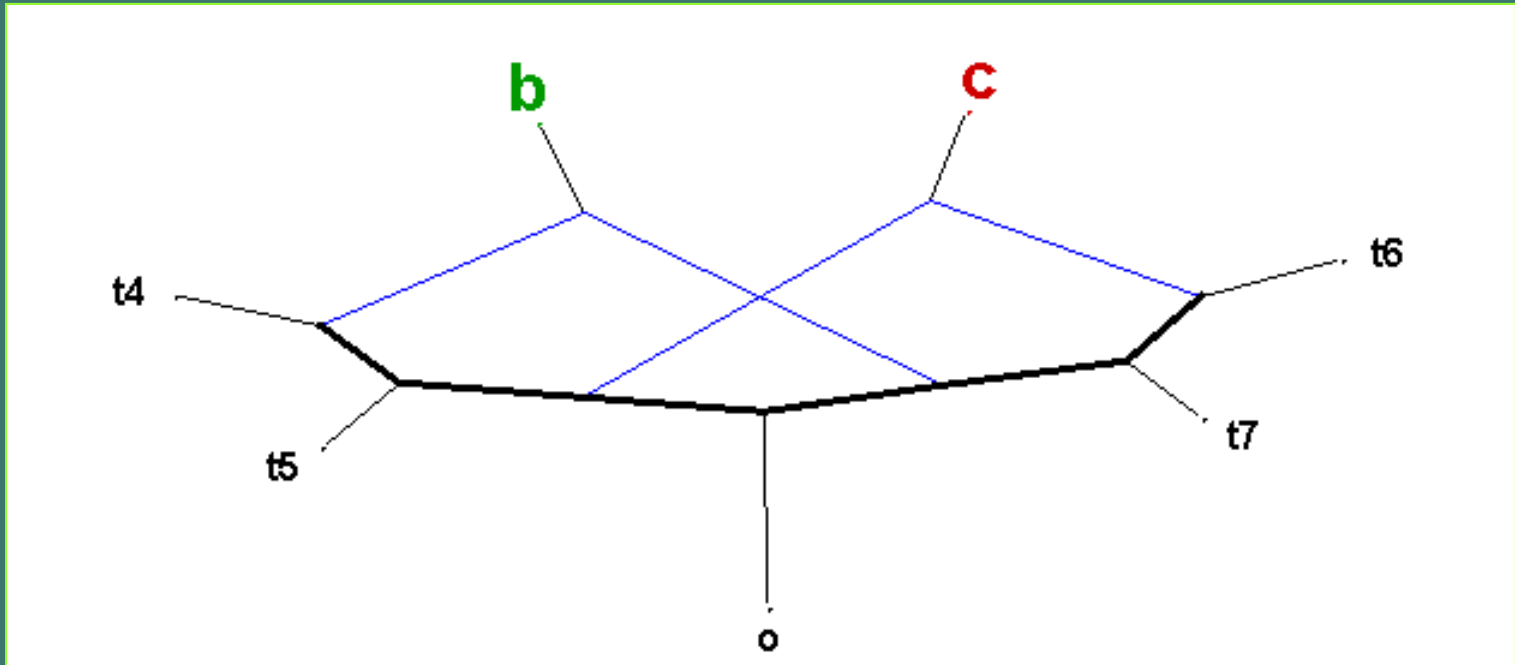
Check Candidate

- For $R=\{b,c\}$, check that reticulation cycles overlap correctly along a path:



Network Construction

- Modify splits network to represent reticulations:

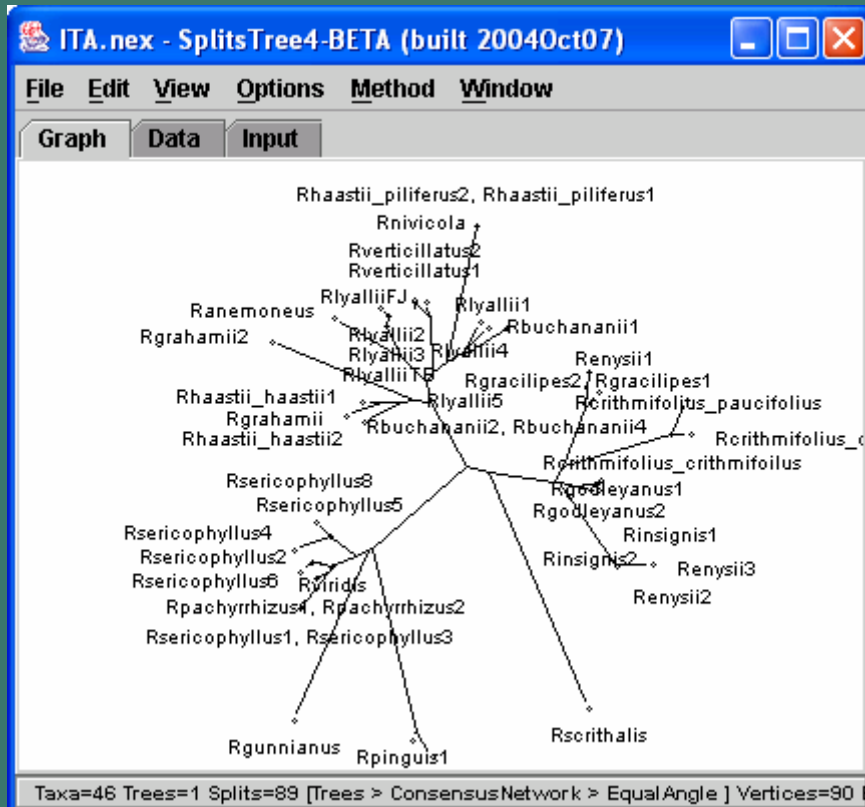


Splits-Based Algorithm

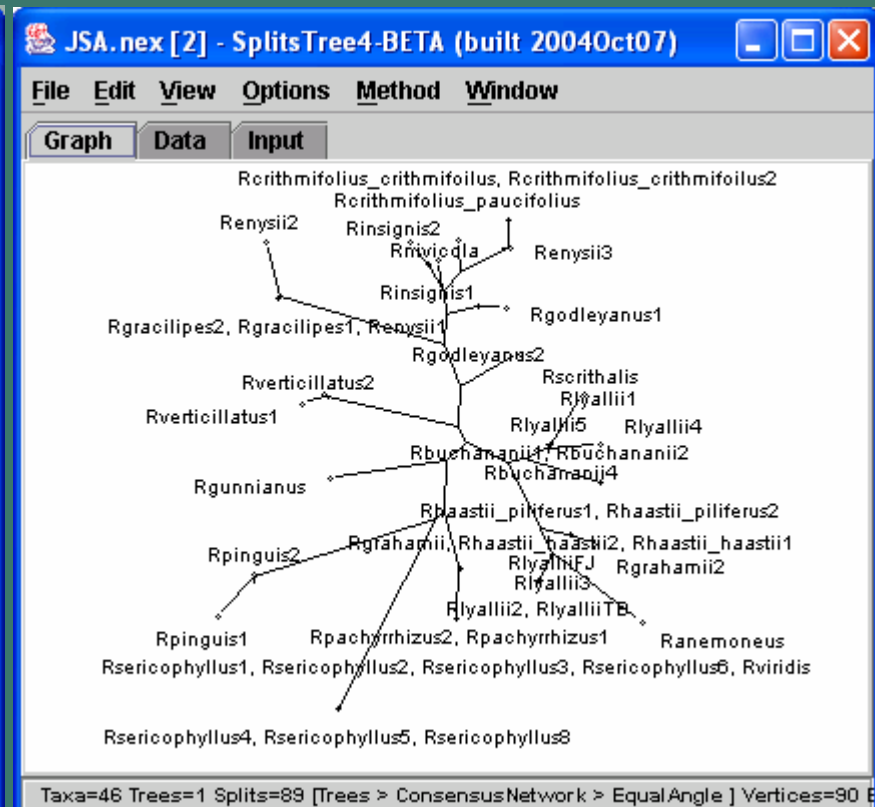
- Input:
 - Set of trees \mathcal{T} , not necessarily bifurcating, can be partial trees
 - Parameter k
- Output:
 - All reticulate networks N for which every incompatibility component can be explained by at most k overlapping reticulations
- Complexity: polynomial for fixed k

Application to Real Data

New Zealand *Ranunculus* (buttercup) species



Nuclear ITS region



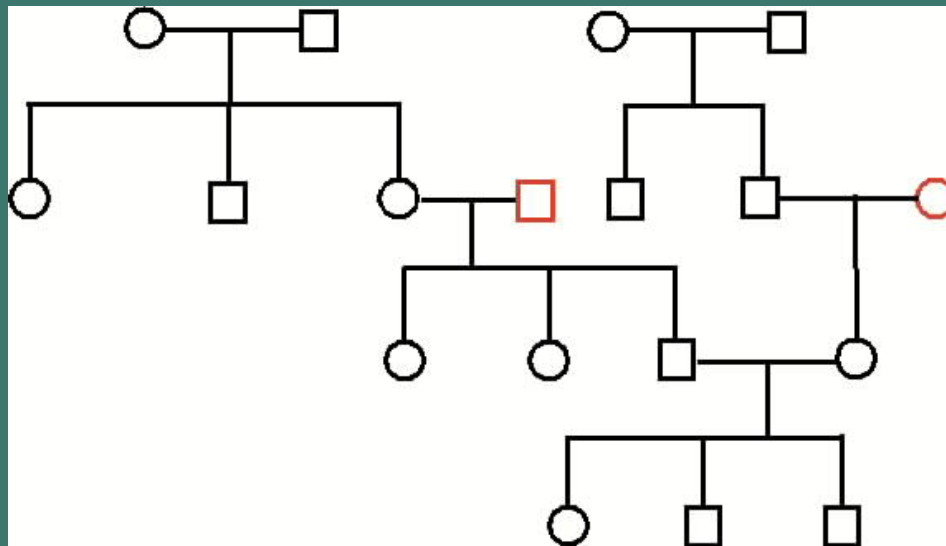
Chloroplast J_{SA} region

Part V

1. Phylogenetic trees
2. Splits networks
3. Consensus networks
4. Hybridization and reticulate networks
5. Recombination networks

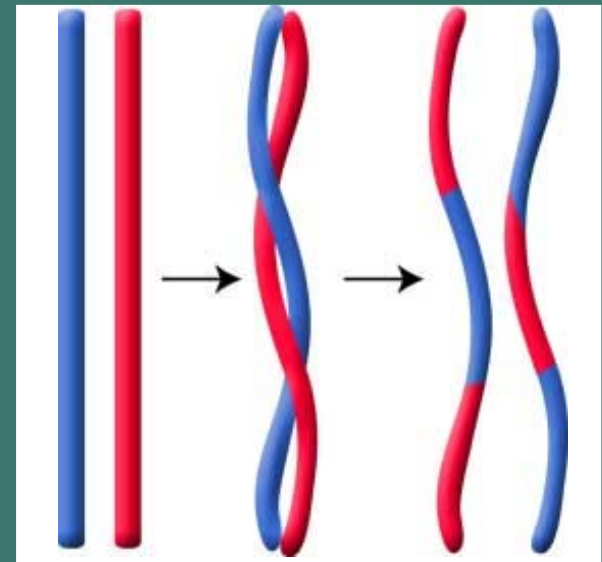
Recombination

- Recombination is studied in population genetics [24, 20, 16, 46, 47, 48] and there *ancestor recombination graphs* (ARGs) are used for statistical purposes.



Chromosomal Recombination

- We will study the combinatorial aspects of chromosomal (meiotic) recombination and will consider *recombination networks* rather than ARGs.
- Simplifying assumptions:
 - all sequences have a common ancestor, and
 - any position can mutate at most once.



Recombination Network

Tree-based approach [Gusfield *et al.* 2003]
for computing galled trees:

- For each component:
 - Determine whether removing one taxon produces a perfect phylogeny
 - If so, arrange taxa in gall
 - Return description of network

Recombination Network

Splits-based approach [Huson & Klopper 2005]
for computing overlapping networks:

- Determine a reticulate network as described above.
- Compute the labeling of nodes and edges.

Recombination Network

- [Lungso, Sun and Hein, to appear in WABI 2005]:
Branch and bound approach to
computing unrestricted recombination
network

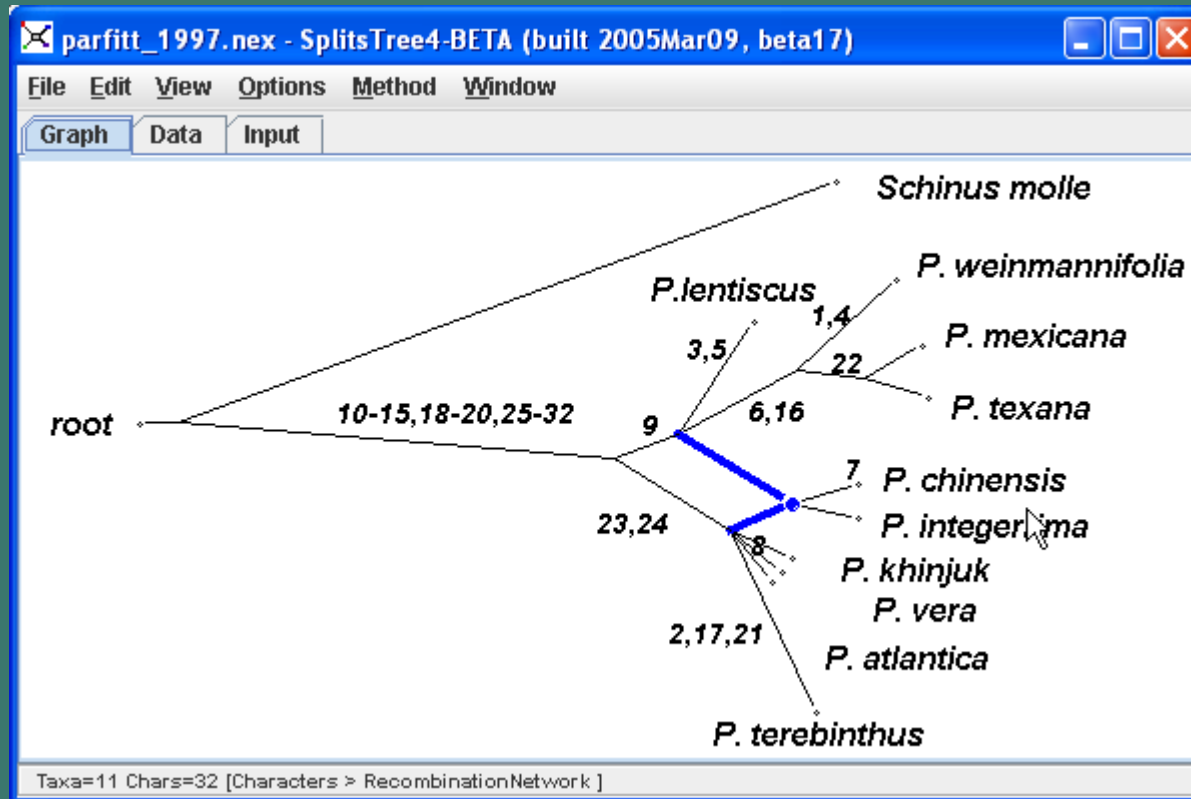
Example 1, Data

- Input: Presence (0) or absence (1) of a given restriction site in a 3.2kb region of variable chloroplast DNA in *Pistacia* [Parfitt & Badenes 1997]:

■ <i>P.lentiscus</i>	0111001010000000111000000010000
■ <i>P.weinmannifolia</i>	11001110100000010111000000010000
■ <i>P.chinensis</i>	01011000100000000111001100010000
■ <i>P.integerrima</i>	01011010100000000111001100010000
■ <i>P.terebinthus</i>	00011010000000001111101100010000
■ <i>P.atlantica</i>	01011011000000000111001100010000
■ <i>P.mexicana</i>	01011110100000010111010000010000
■ <i>P.texana</i>	01011110100000010111010000010000
■ <i>P.khinjuk</i>	01011010000000000111001100010000
■ <i>P.vera</i>	01011010000000000111001100010000
■ <i>Schinus molle</i>	01011010011111100000000011101111

Example 1, Recombination Network

- Load this data in to SplitsTree4 and select RecombinationNetwork to obtain:





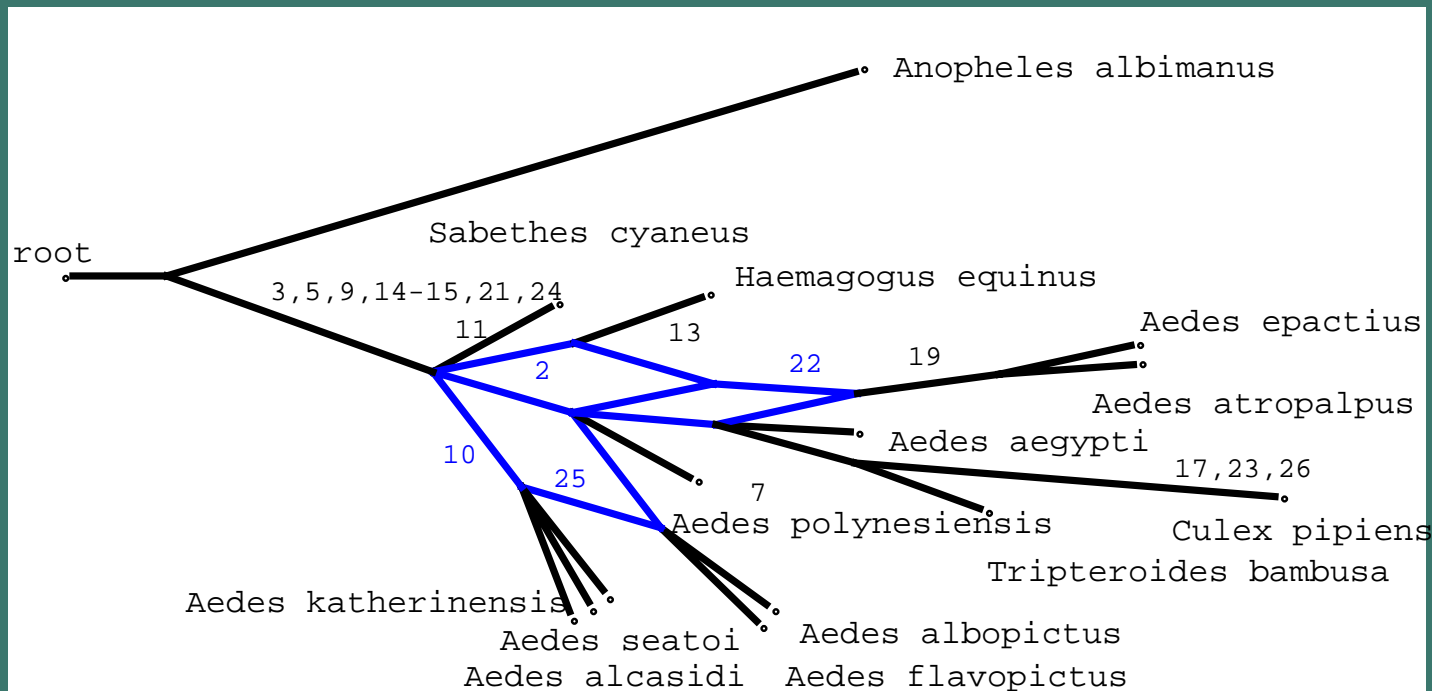
Example 2, Data

- Input: Restriction maps of the rDNA cistron (length \approx 10kb) of twelve species of mosquitoes using eight 6bp recognition restriction enzymes [Kumar *et al*, 1998]:

<i>Aedes albopictus</i>	11110101010100010101010010
<i>Aedes aegypti</i>	11110101000100010101000010
<i>Aedes seatoi</i>	11110101010100010101010000
<i>Aedes avopictus</i>	11110101010100010101010010
<i>Aedes alcalasi</i>	11110101010100010101010000
<i>Aedes katherinensis</i>	11110101010100010101010000
<i>Aedes polynesiensis</i>	11110101000100010101010010
<i>Aedes triseriatus</i>	10110101000110010101000000
<i>Aedes atropalpus</i>	10110101000100010111000010
<i>Aedes epactius</i>	10110101000100010111000010
<i>Haemagogus equinus</i>	10110101000110010101010000
<i>Armigeres subalbatus</i>	10110101000100010101000000
<i>Culex pipiens</i>	11110111000100011101001011
<i>Tripteroides bambusa</i>	11110111000100010101000010
<i>Sabethes cyaneus</i>	11110101001100010101010000
<i>Anopheles albimanus</i>	11011101100101110101110100

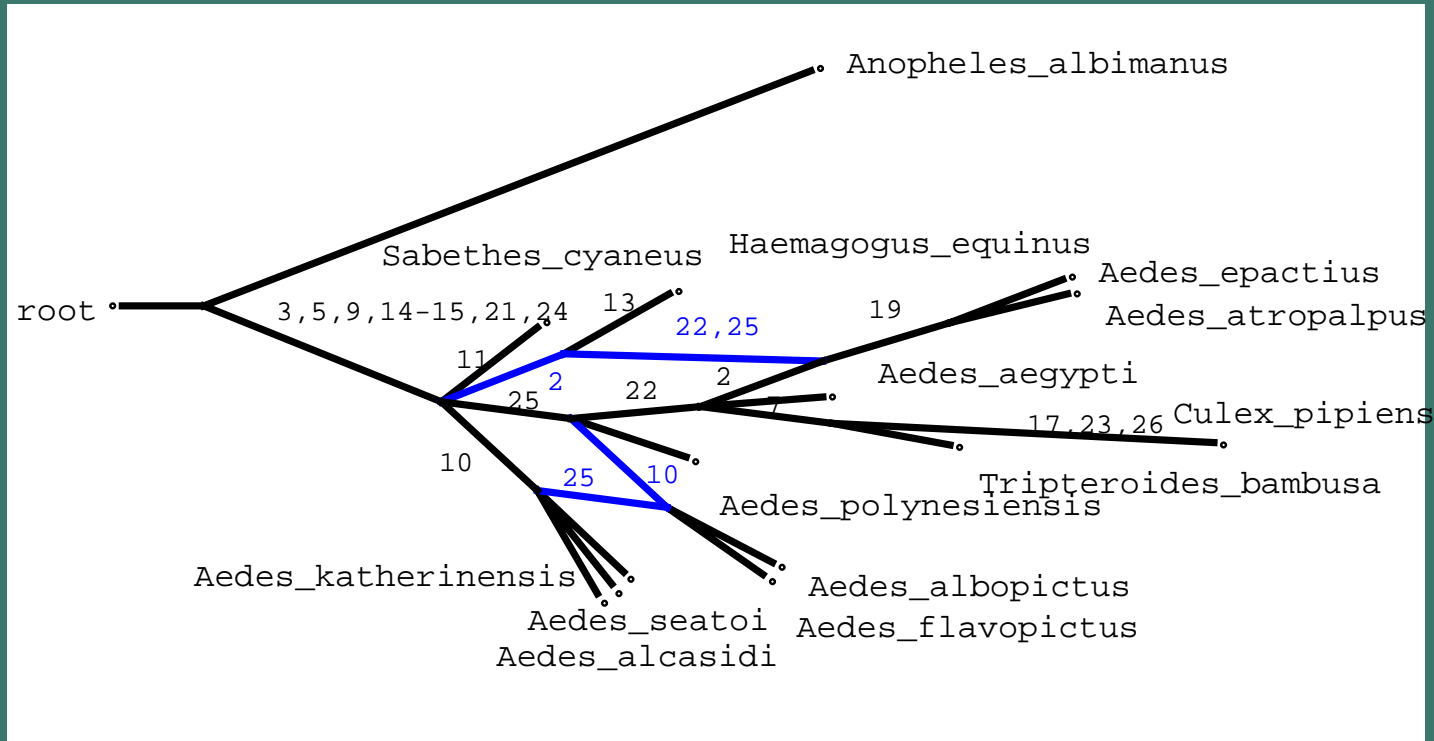
Example 2, Subset

- Recombination scenarios based on the complete data set look unconvincing. However, trial-and-error removal of two taxa *Aedes triseriatus* and *Armigeres subalbatus* gives rise to a simpler splits network:



Example 2, Recombination Network

- A possible recombination scenario is given by:



- Here, *Haemagogus equinus* appears to arise by a single-crossover recombination, and a second such recombination leads to *A. albopictus* and *A. avopictus*.



Example 3, Data

- 19 restriction endonucleases were used to analyze patterns of cleavage site variation in the mtDNA of *Zonotrichia*.
- 7 taxa, 122 characters
- [Zink *et al*, 1991]

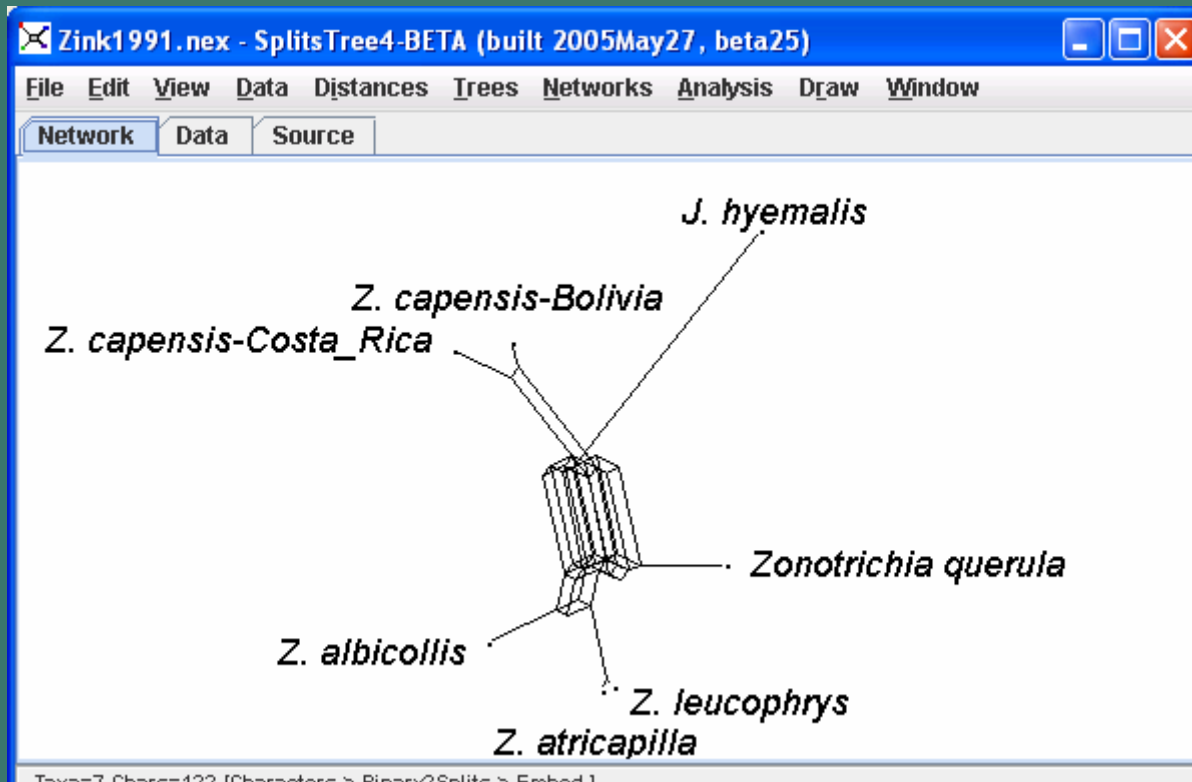
```

notrichia_querula'      1110001111110001111001111111100001111000111010101100011111001111111000100111110011111101111001101111
_atricapilla'          1110001111110000110000111111110000111110011100010111001111100111111100010011111101111101111101101101
_leucophrys'           1110001111110000110000111111110000111110011110010111001111100111111100010011111101111101111101101101
_albicollis'           1110001111110000110000101111110100001110101110001011110111111011111110010011111001111100111100111001101101
_capensis--Bolivia'    0111001110010000100000111111100011011000001110101011110111110011011110010110111110001111100111011011101
_capensis--Costa_Rica' 111001111001011010000011111110001101100000111000101111011111001101111001011111110001101100111011011101
_hyemalis'             1110101110010001100011001111100000111100011110111111001110010101010111001100111100000110010111001101101
    
```

- However: recombination of mtDNA unlikely

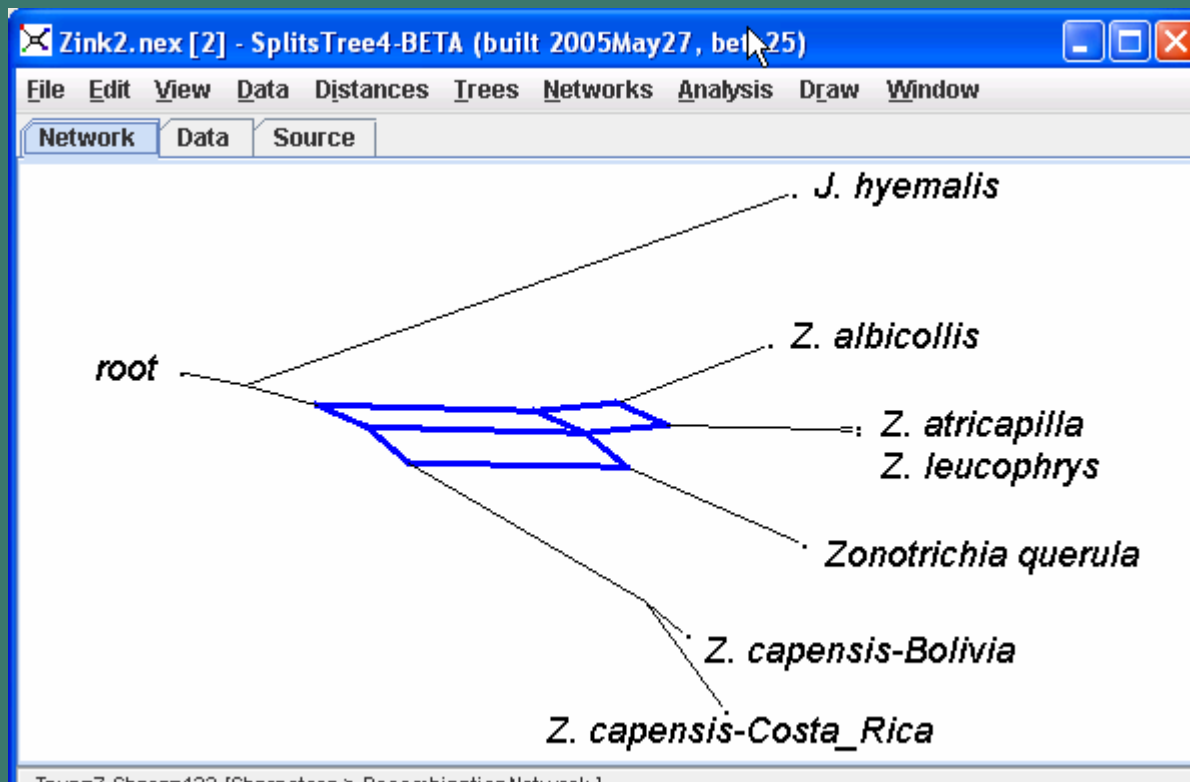
Example 3, Median Network

- The unrooted splits network for a dataset of restriction sites in the mtDNA of *Zonotrichia*:



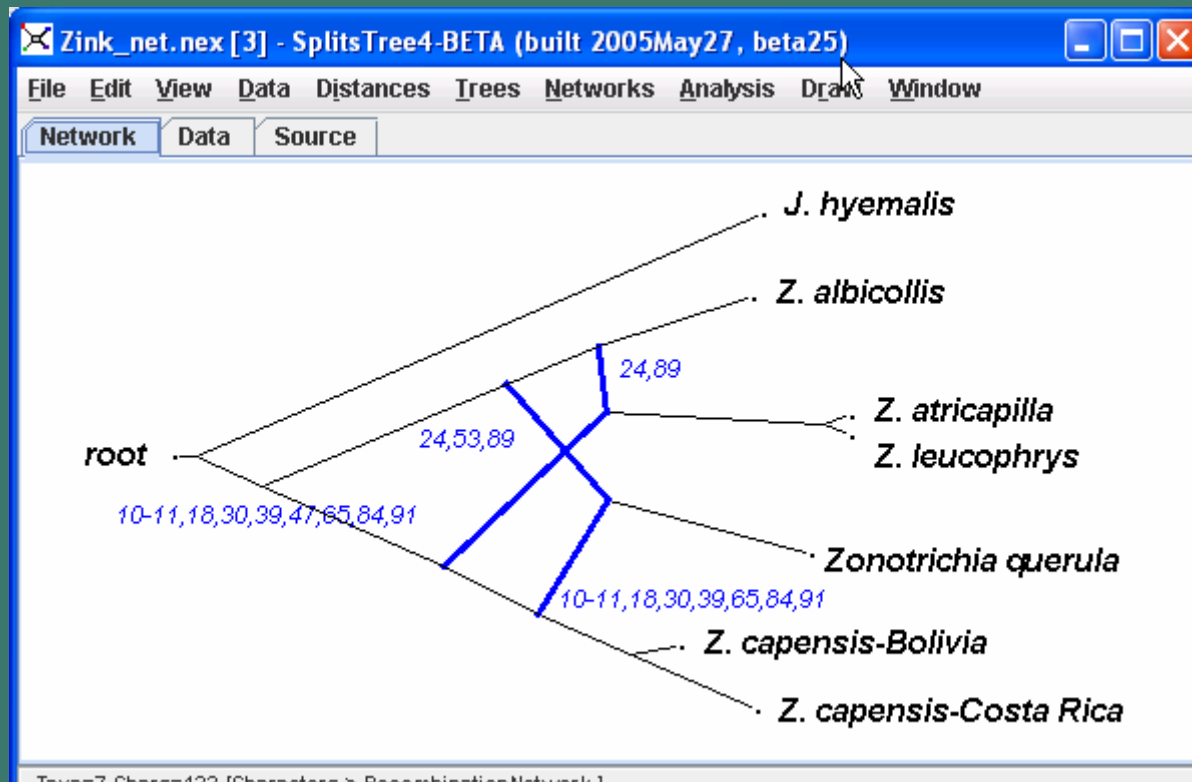
Example 3, Significant Differences

- The rooted splits network for the same data set, but suppressing all splits that are only supported by one site in the data:



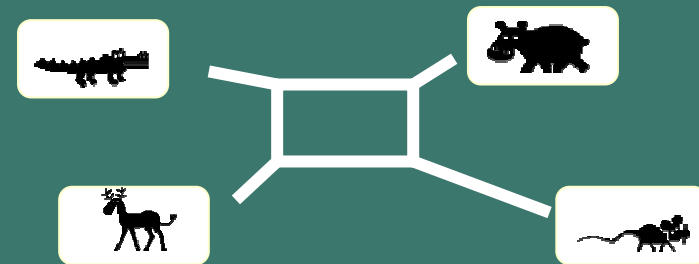
Recombination Network

- Possible recombination scenario involving two non-independent reticulations:



Summary

- Incompatible signals in gene trees can be usefully displayed using splits networks
- A reticulate network may be extracted by combinatorial analysis of individual components
- Implementations of many tree and network methods are available in **SplitsTree4**



3rd RECOMB Comparative Genomics Satellite Workshop

Trinity College Dublin

18th-20th September, 2005



Bringing together researchers from molecular biology, computer science and mathematics to discuss the latest developments in comparative genomics. This is the third workshop in this highly successful series.

Scientific organising committee

Aoife McLysaght, Trinity College Dublin

Daniel Huson, Tübingen University, Germany

Jens Lagergren, Stockholm Bioinformatics Centre and KTH

David Sankoff, University of Ottawa

http://www.gen.tcd.ie/recomb_cg