

The Multiple Aspects of Tree Comparisons in Evolutionary Biology

Damien M. de Vienne
Laboratoire de Biométrie et Biologie Évolutive

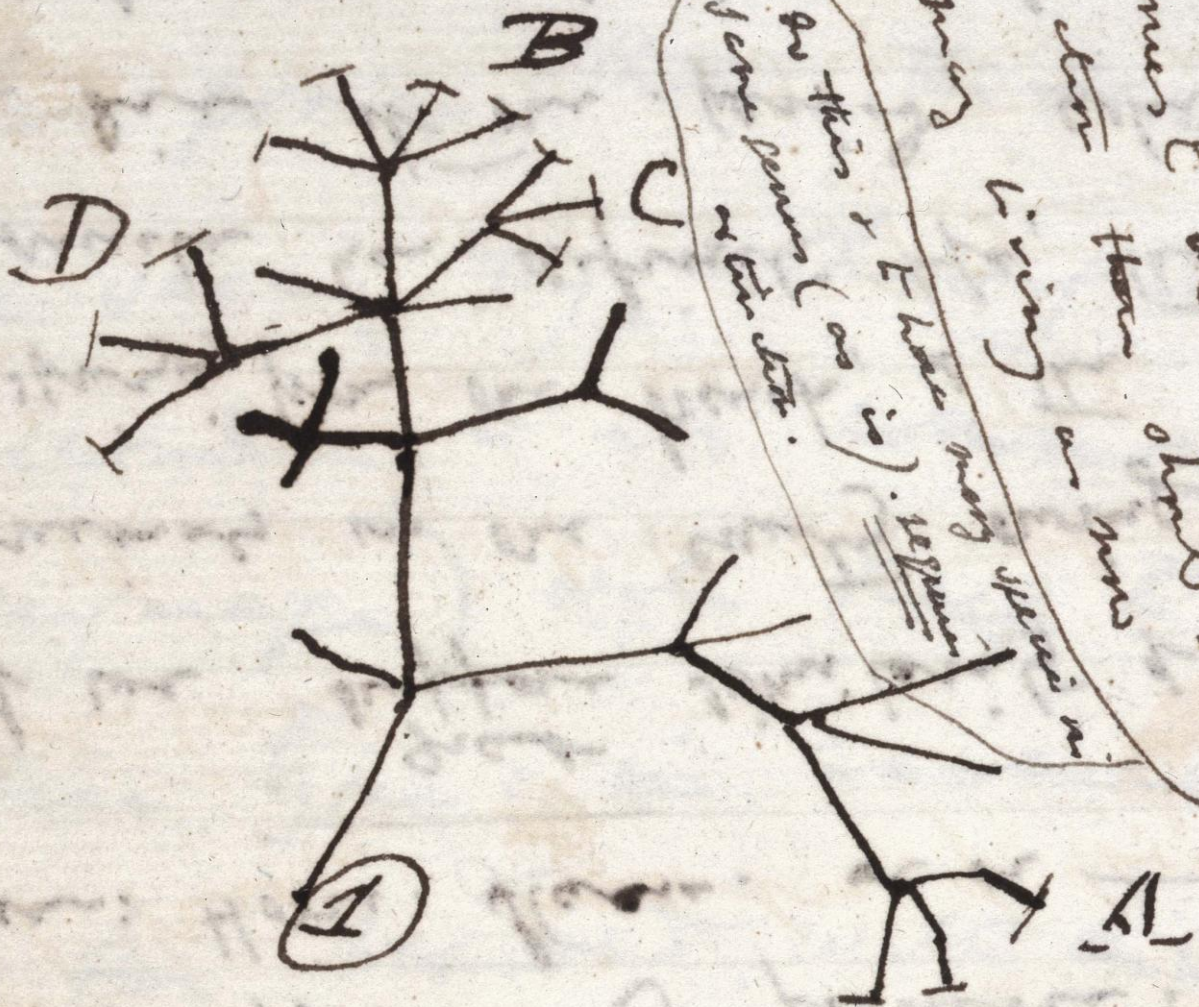


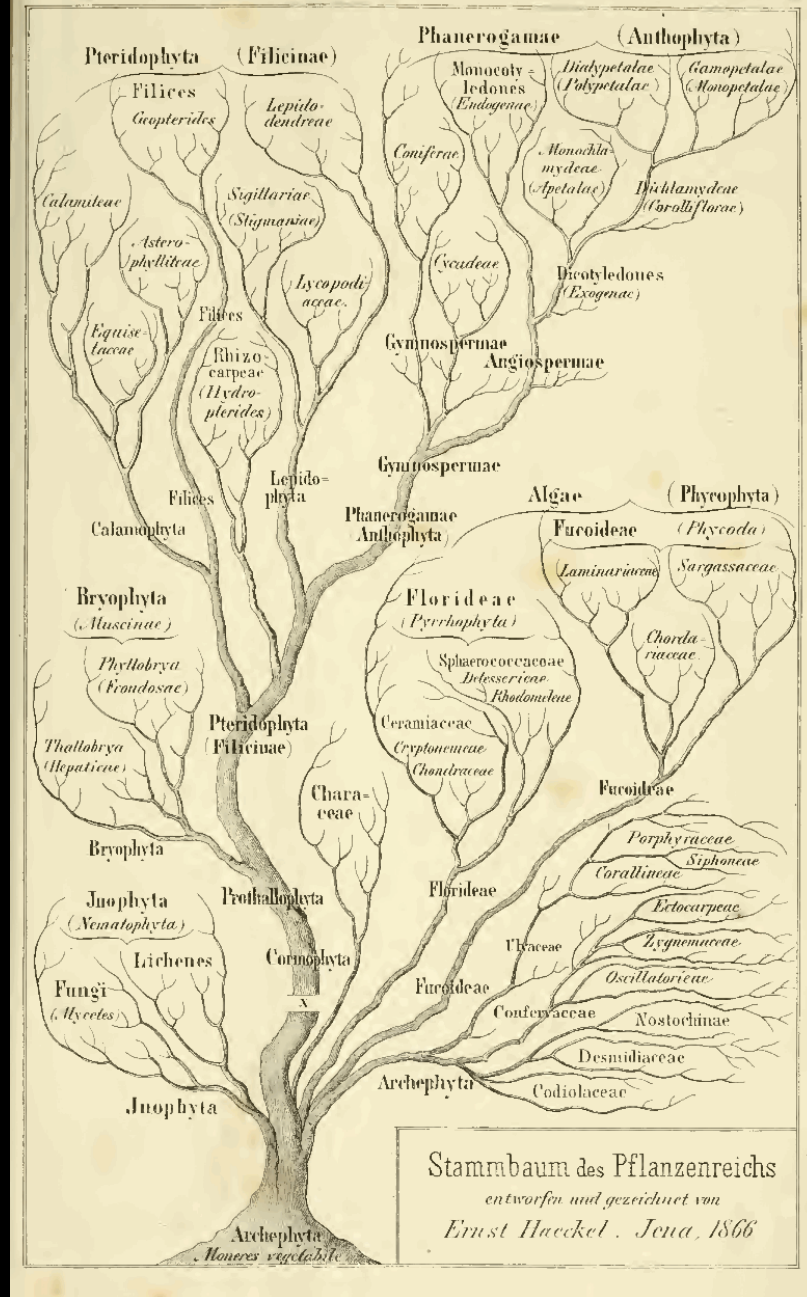
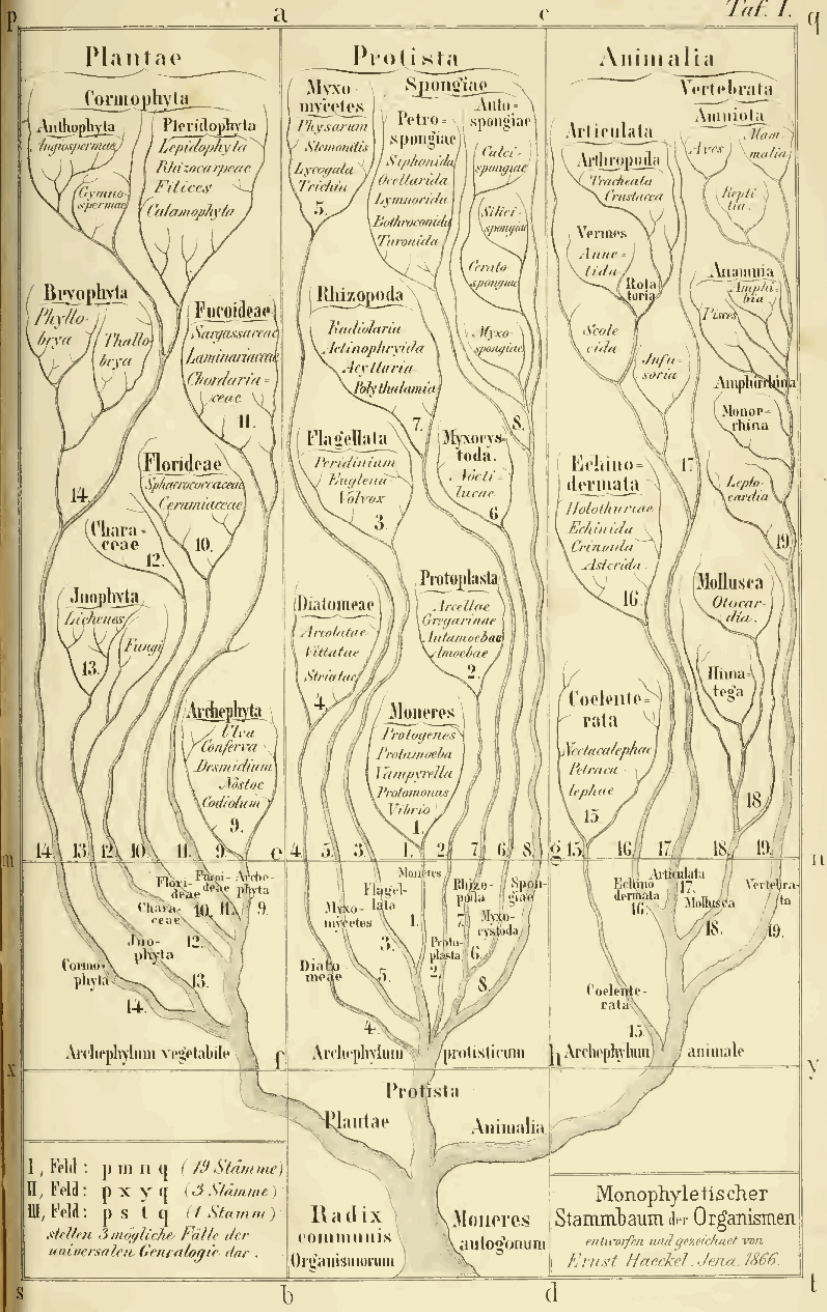
Aussois
7-10 avril 2015

phylogeny

- "genesis and evolution of a phylum," 1869, from German Phylogenie, coined 1866 by German biologist Ernst Heinrich Haeckel (1834-1919) from Greek phylon "race" (see phylo-) + -geneia "origin"

I think

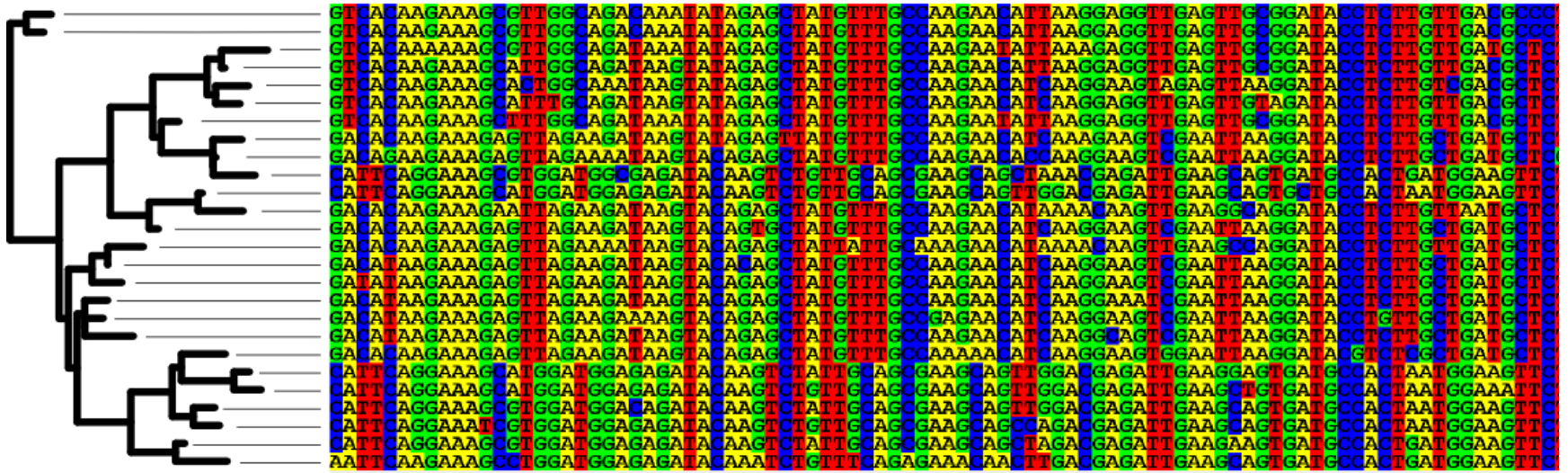




Molecular data

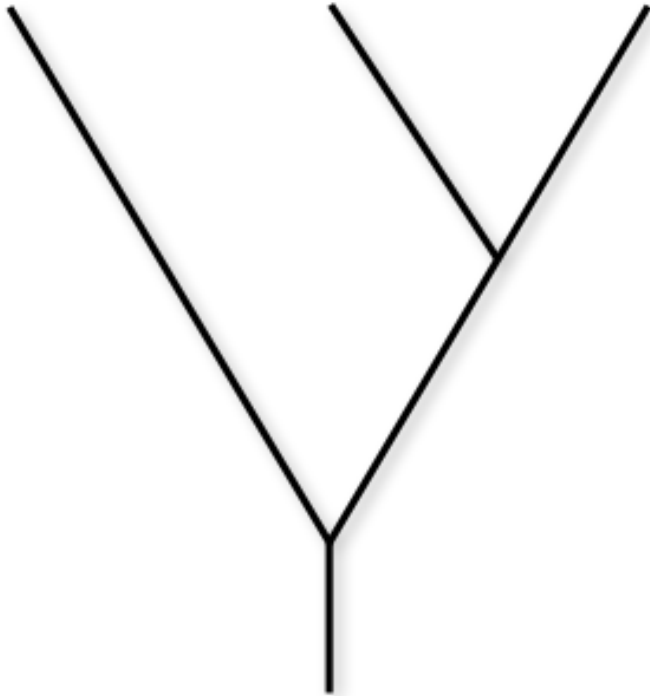
GT CACAAGAAAGCGTTGGCAGACAAATATAGAGCTATGTTTGCCA
GT CACAAGAAAGCTTTGGCAGATAAATATAGAGCTATGTTTGCCA
GT CACAAGAAAGAGTTAGAAGATAAGTATAGAGTTATGTTTGCCA
GT CAGAAGAAAGAGTTAGAAAATAAGTACAGAGCTATGTTTGCCA
CATT CAGGAAAGCGTGGATGGCAGATAACAAGTCTGTTGCAGCGA
CATT CAGGAAAGCATGGATGGAGAGATAACAAGTCTGTTGCAGCGA
GACACAAGAAAGAAATTAGAAGATAAGTACAGAGCTATGTTTGCCA
GACACAAGAAAGAGTTAGAAGATAAGTACAGTGCTATGTTTGCCA
GACACAAGAAAGAGTTAGAAAATAAGTACAGAGCTATTATTGCAA
GACATAAGAAAGAGTTAGAAGATAAGTACACAGCTATGTTTGCCA
GATATAAGAAAGAGTTAGAAGATAAGTACAGAGCTATGTTTGCCA
GACATAAGAAAGAGTTAGAAGATAAGTACAGAGCTATGTTTGCCA
GACATAAGAAAGAGTTAGAAGAAAAGTACAGAGCTATGTTTGCCG
GACATAAGAAAGAGTTAGAAGATAAGTACAGAGCTATGTTTGCCA
GACACAAGAAAGAGTTAGAAGATAAGTACAGAGCTATGTTTGCCA
CATT CAGGAAAGCATGGATGGAGAGATAACAAGTCTATTGCAGCGA
CATT CAGGAAAGCATGGATGGAGAGATAACAAGTCTGTTGCAGCGA
CATT CAGGAAAGCGTGGATGGACAGATAACAAGTCTATTGCAGCGA
CATT CAGGAAATCGTGGATGGAGAGATAACAAGTCTGTTGCAGCGA
CATT CAGGAAAGCGTGGATGGAGAGATAACAAGTCTATTGCAGCGA
AATT CAAGAAAGCGTGGATGGAGAGATACAAATCTGTTTCAGAGA

Molecular Phylogenies



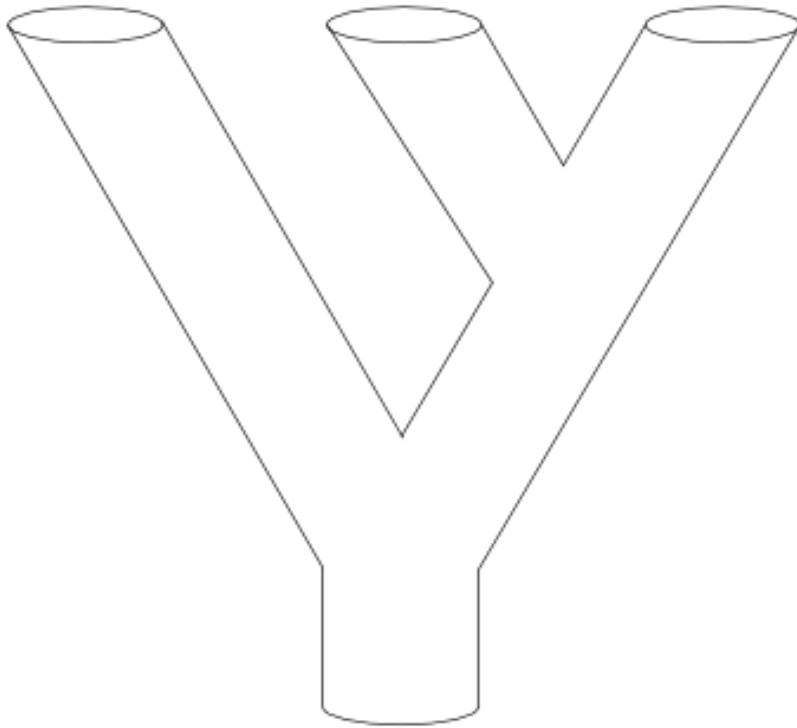
Species trees vs gene trees

Species A Species B Species C



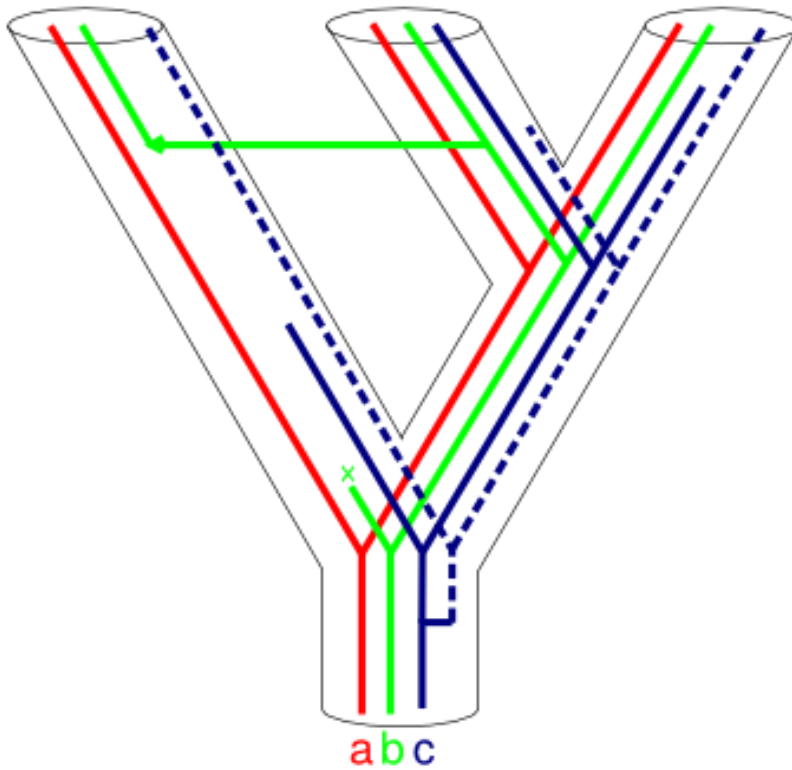
Species trees vs gene trees

Species A Species B Species C



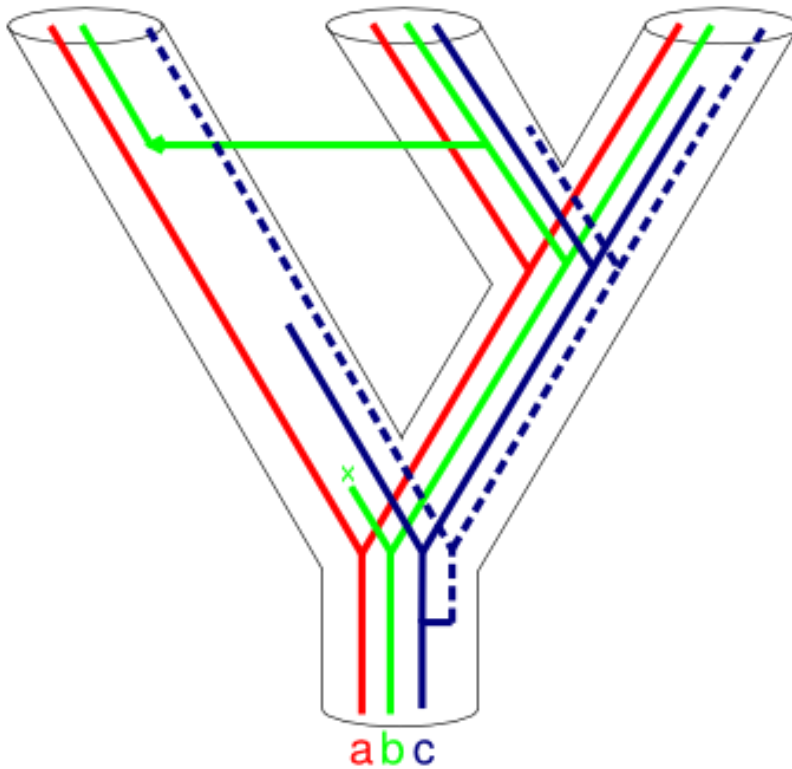
Species trees vs gene trees

Species A Species B Species C

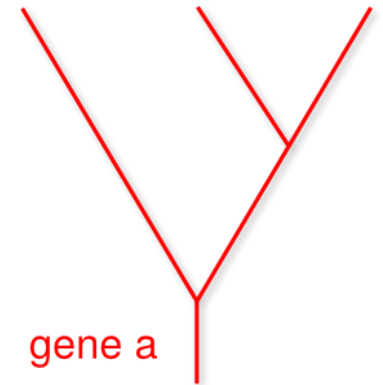


Species trees vs gene trees

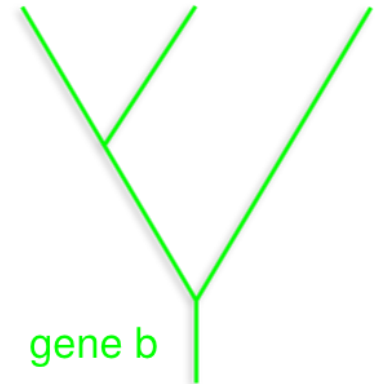
Species A Species B Species C



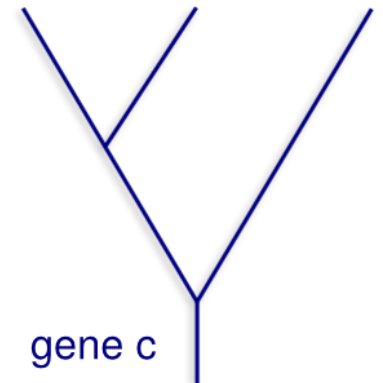
Species A Species B Species C

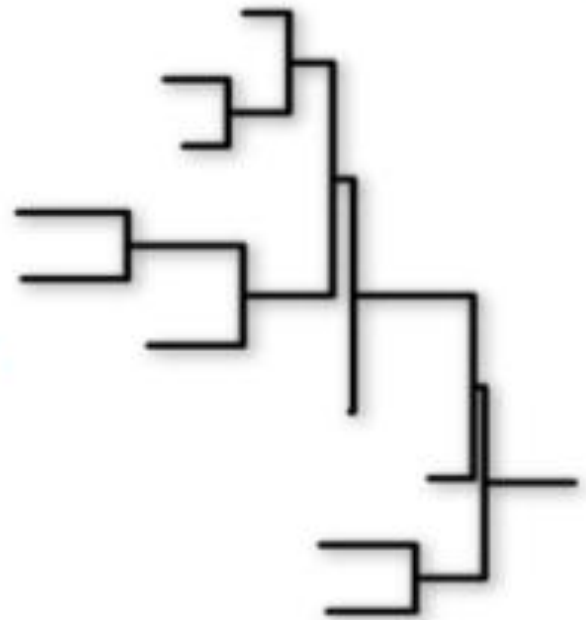
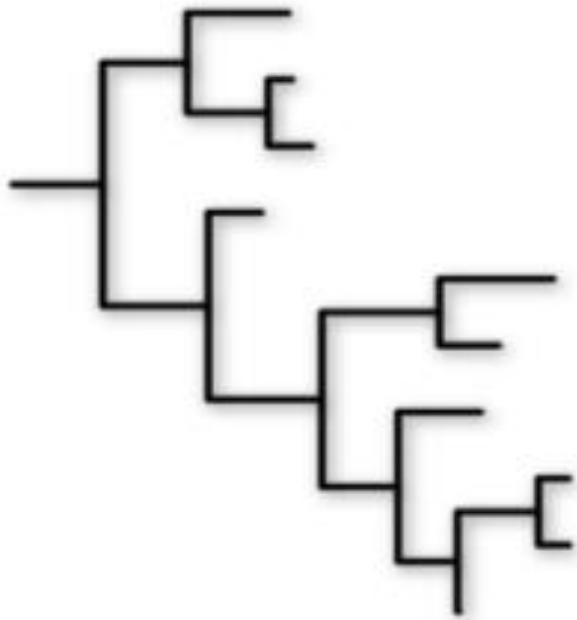


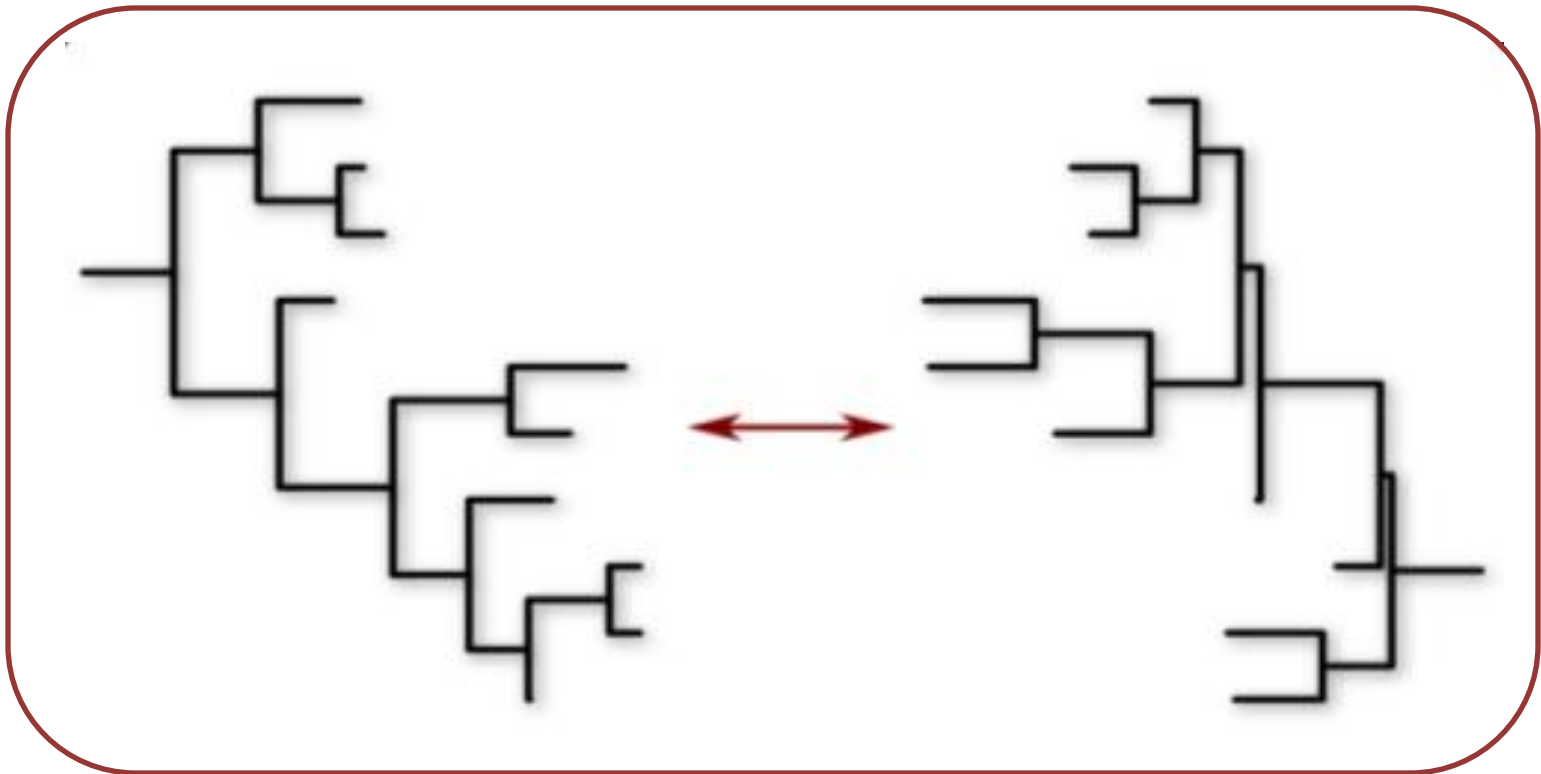
Species A Species B Species C



Species A Species C Species B







Species tree ↔ Species tree

Gene tree ↔ Gene tree

Species tree ↔ Gene tree

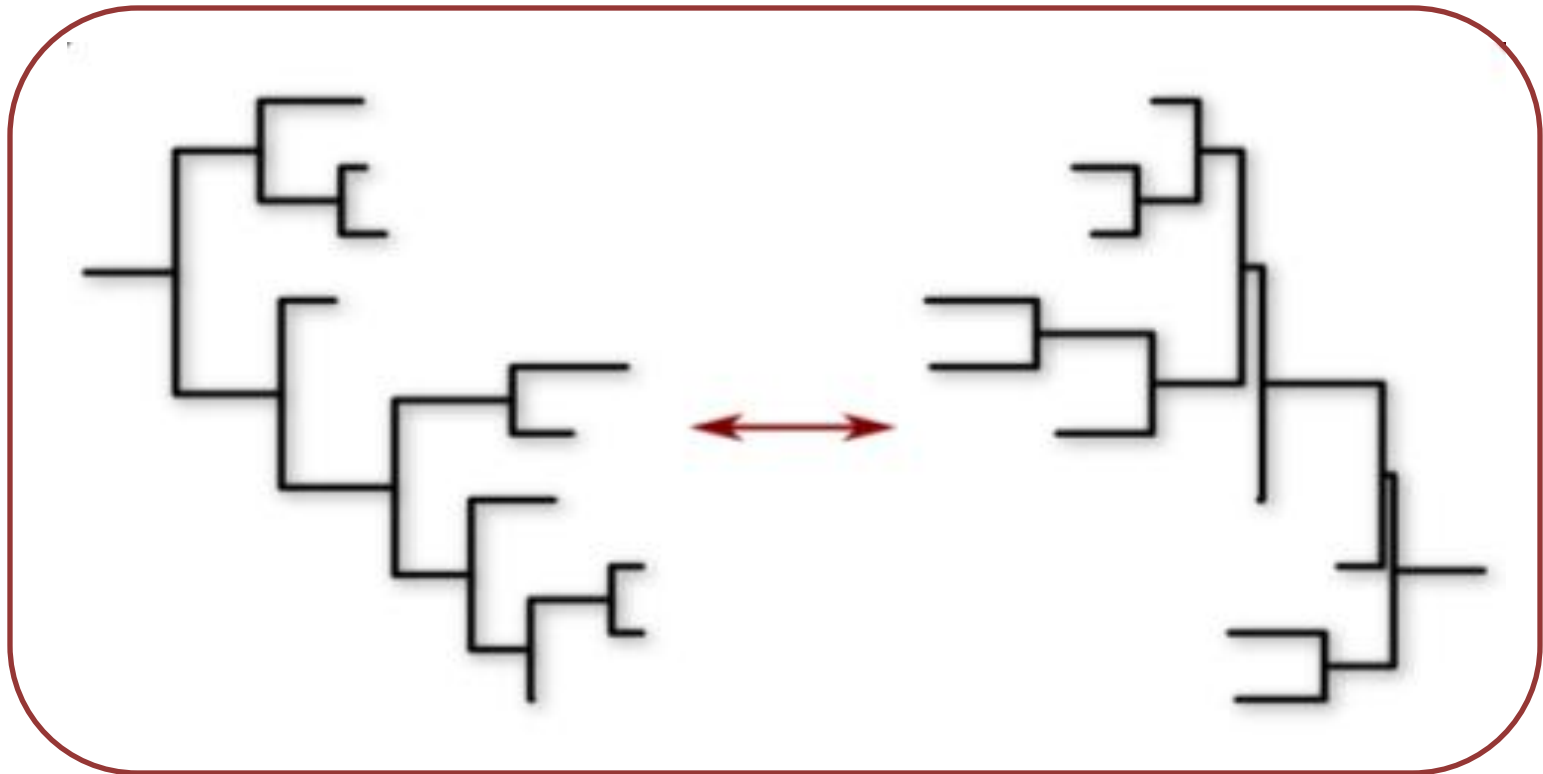
Multiple gene trees

Host parasite coevolutionary studies

Protein-Protein interaction detection

Reconciliation analyses

Phylogenomic studies



Species tree ↔ Species tree

Gene tree ↔ Gene tree

Species tree ↔ Gene tree

Multiple gene trees

Host parasite coevolutionary studies

Protein-Protein interaction detection

Reconciliation analyses

Phylogenomic studies



Fahrenholz, 1913



Fahrenholz's rule: "parasite phylogeny mirrors that of its host"

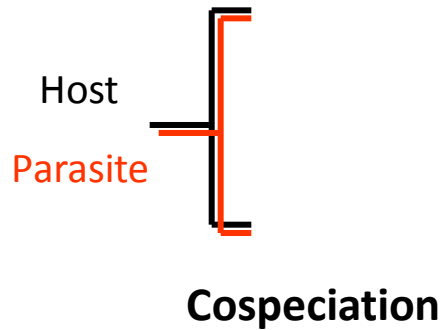


Fahrenholz, 1913

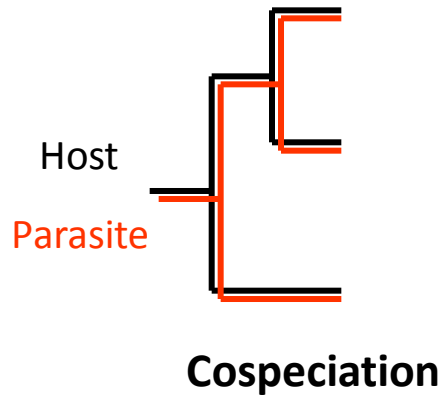


Fahrenholz's rule: "parasite phylogeny mirrors that of its host"

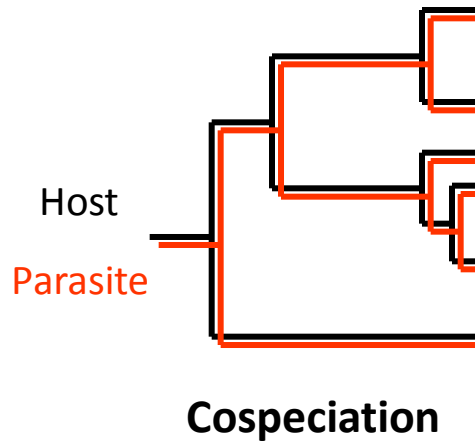
Cospeciation leads to congruent phylogenies



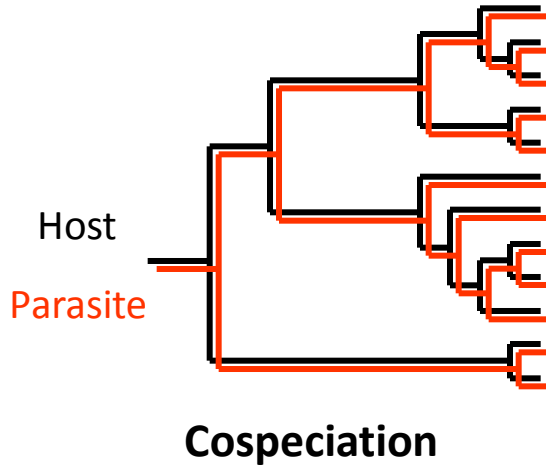
Cospeciation leads to congruent phylogenies



Cospeciation leads to congruent phylogenies

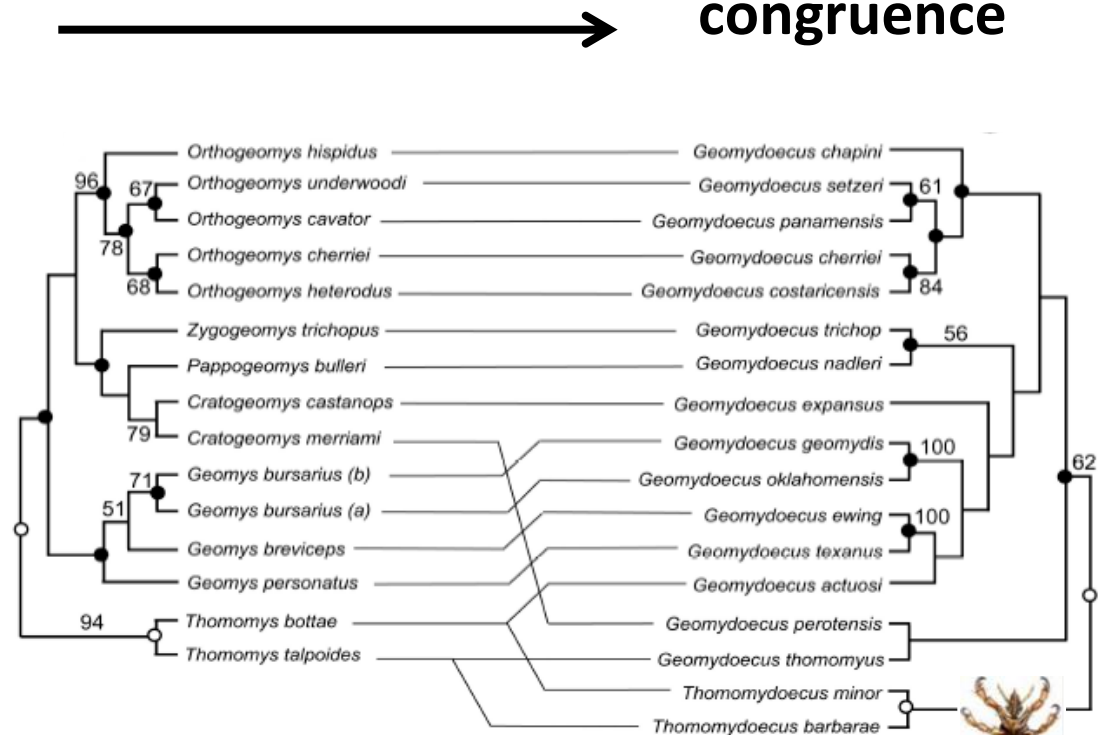
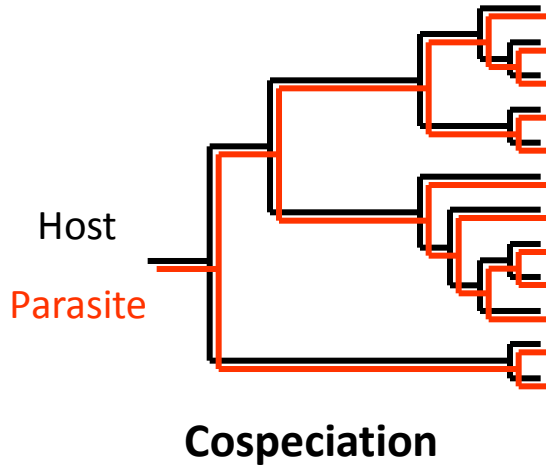


Cospeciation leads to congruent phylogenies

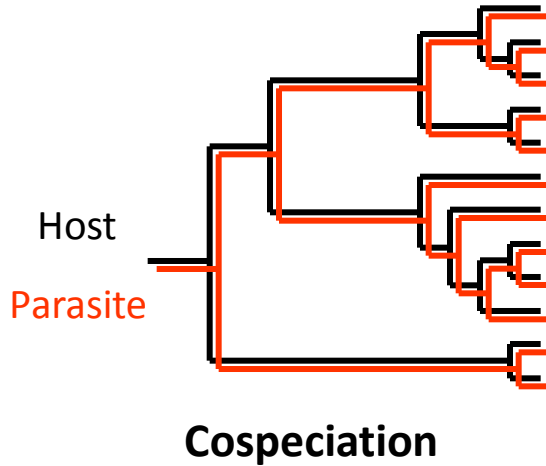


congruence

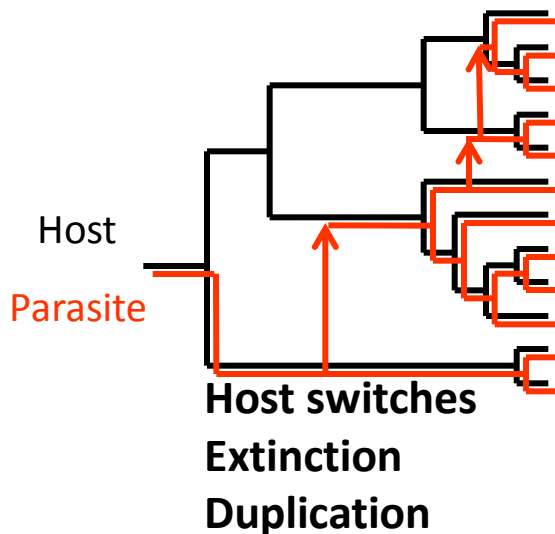
Cospeciation leads to congruent phylogenies



Other events lead to **incongruent** phylogenies

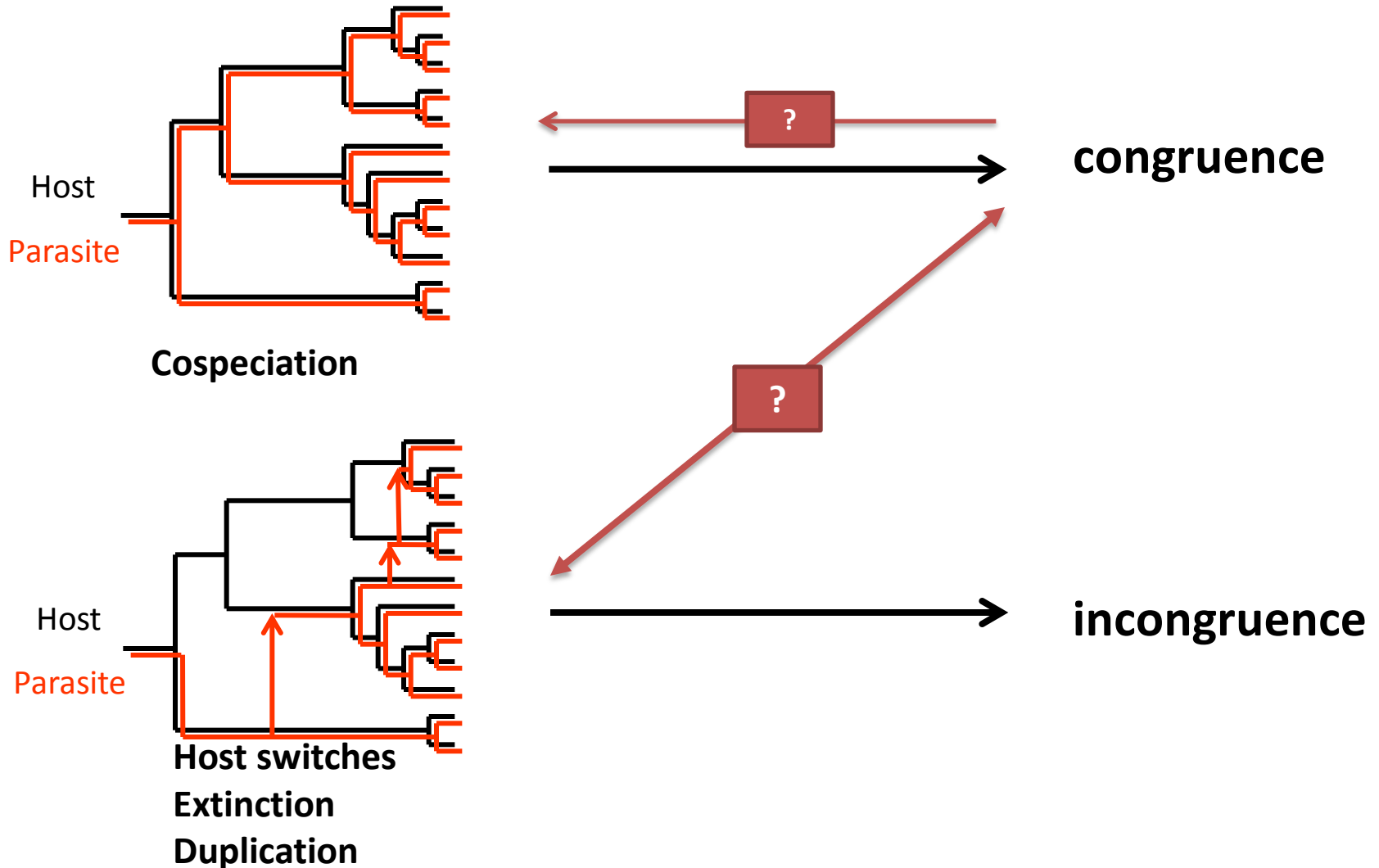


congruence



incongruence

Does congruence imply cospeciations?



Can host-switches produce congruent phylogenies?

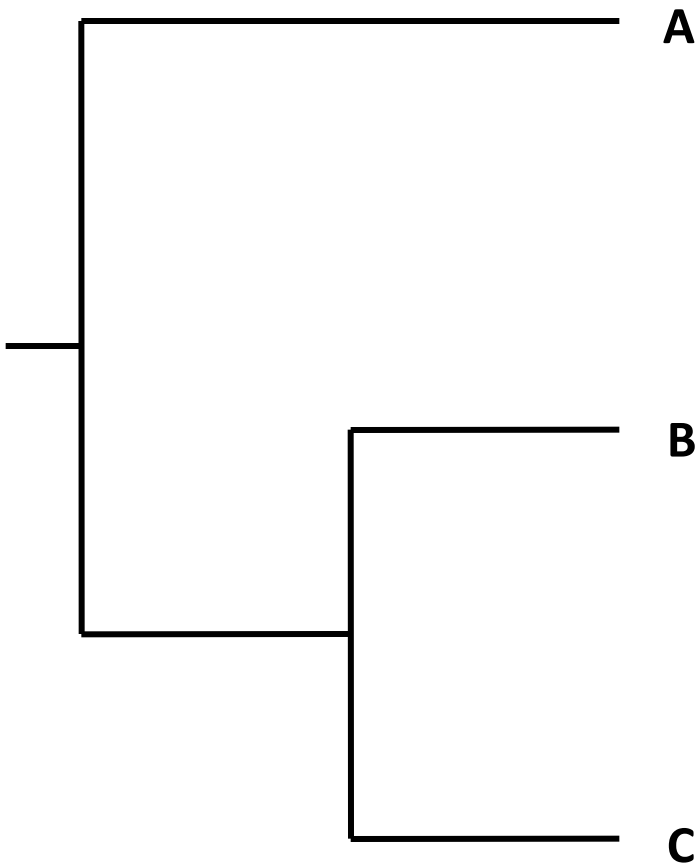
Can congruence be due to host-switches alone?

Adaptive radiation on a group of pre-existing host species

Under what conditions?

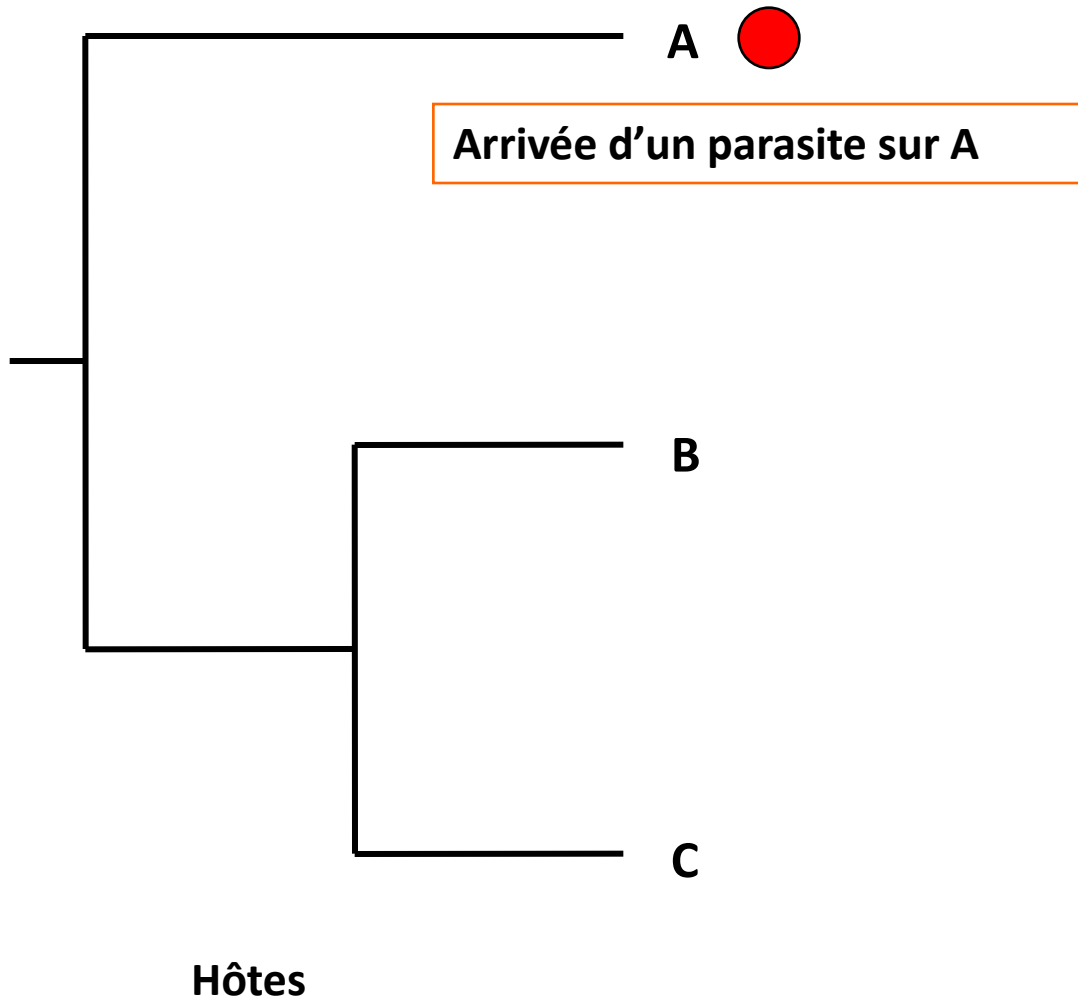
Host tree topology, first host parasitized, host-switch probability, time lag between switch and speciation as a function of the switch distance

Simulate adaptive radiation of a parasite on a group of pre-existing hosts

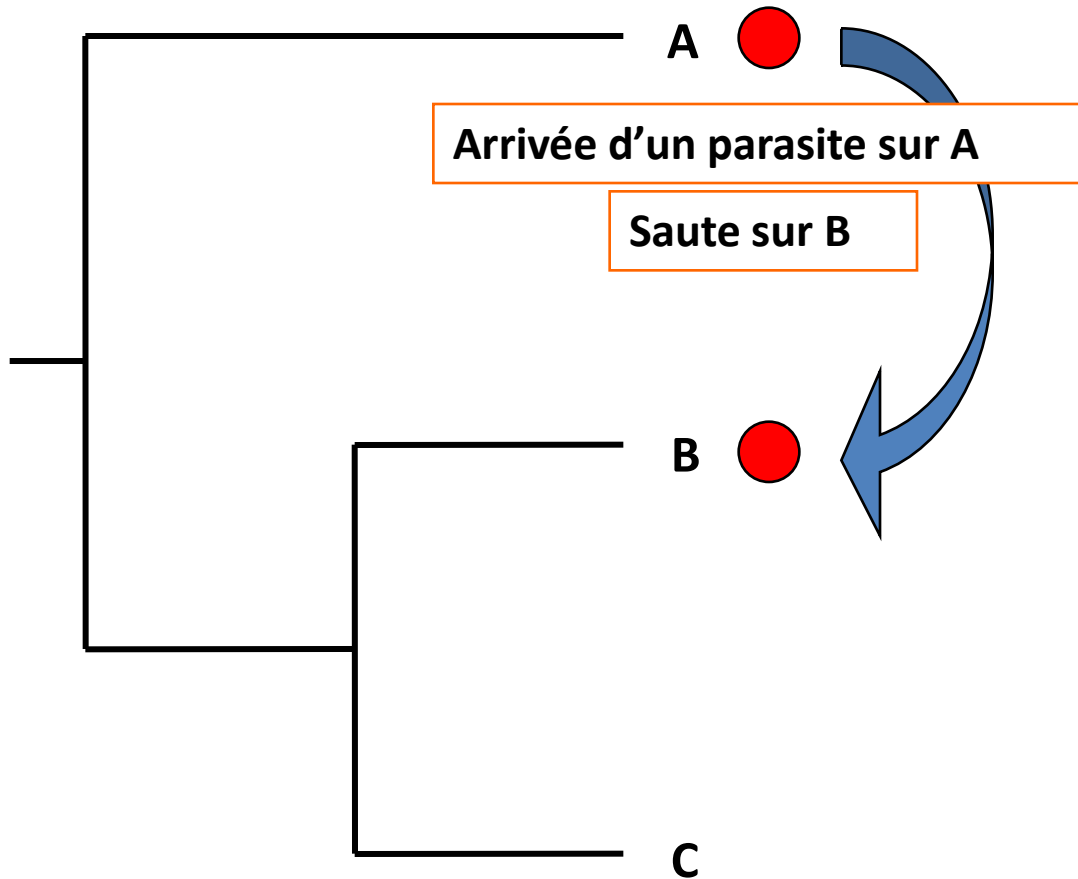


Hôtes

Simulate adaptive radiation of a parasite on a group of pre-existing hosts

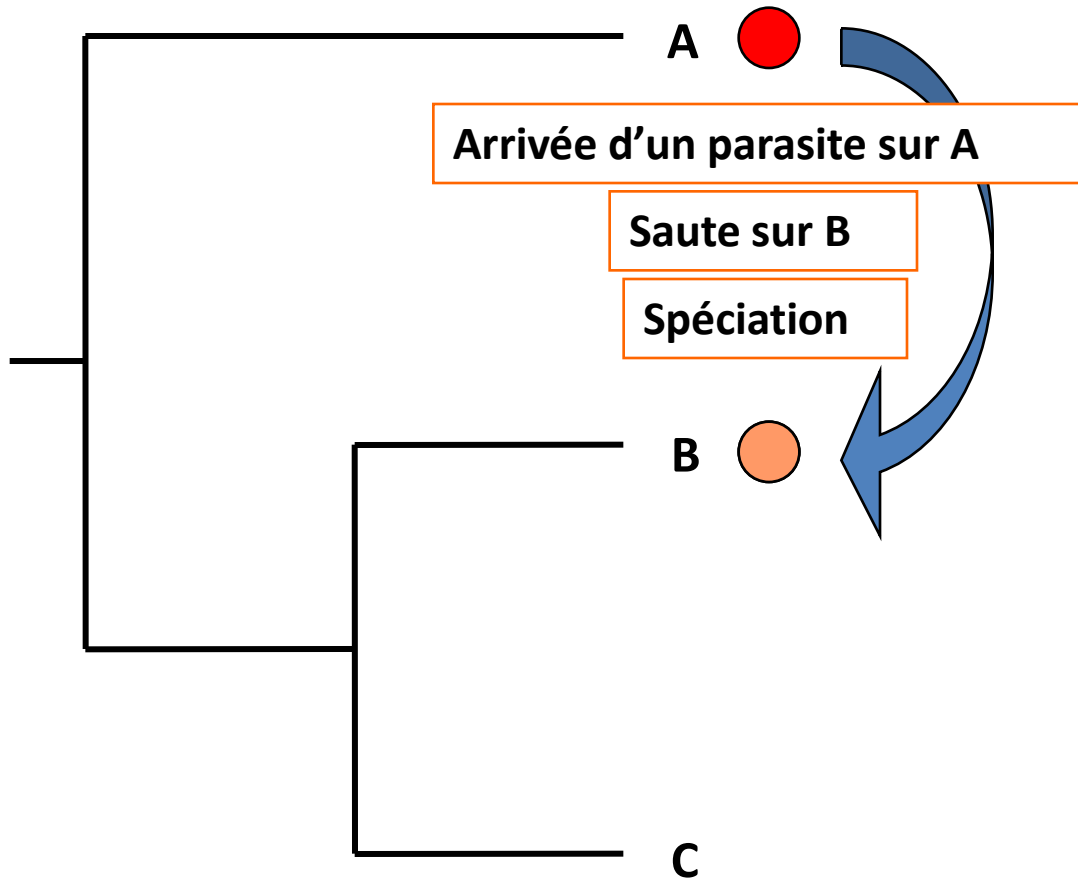


Simulate adaptive radiation of a parasite on a group of pre-existing hosts



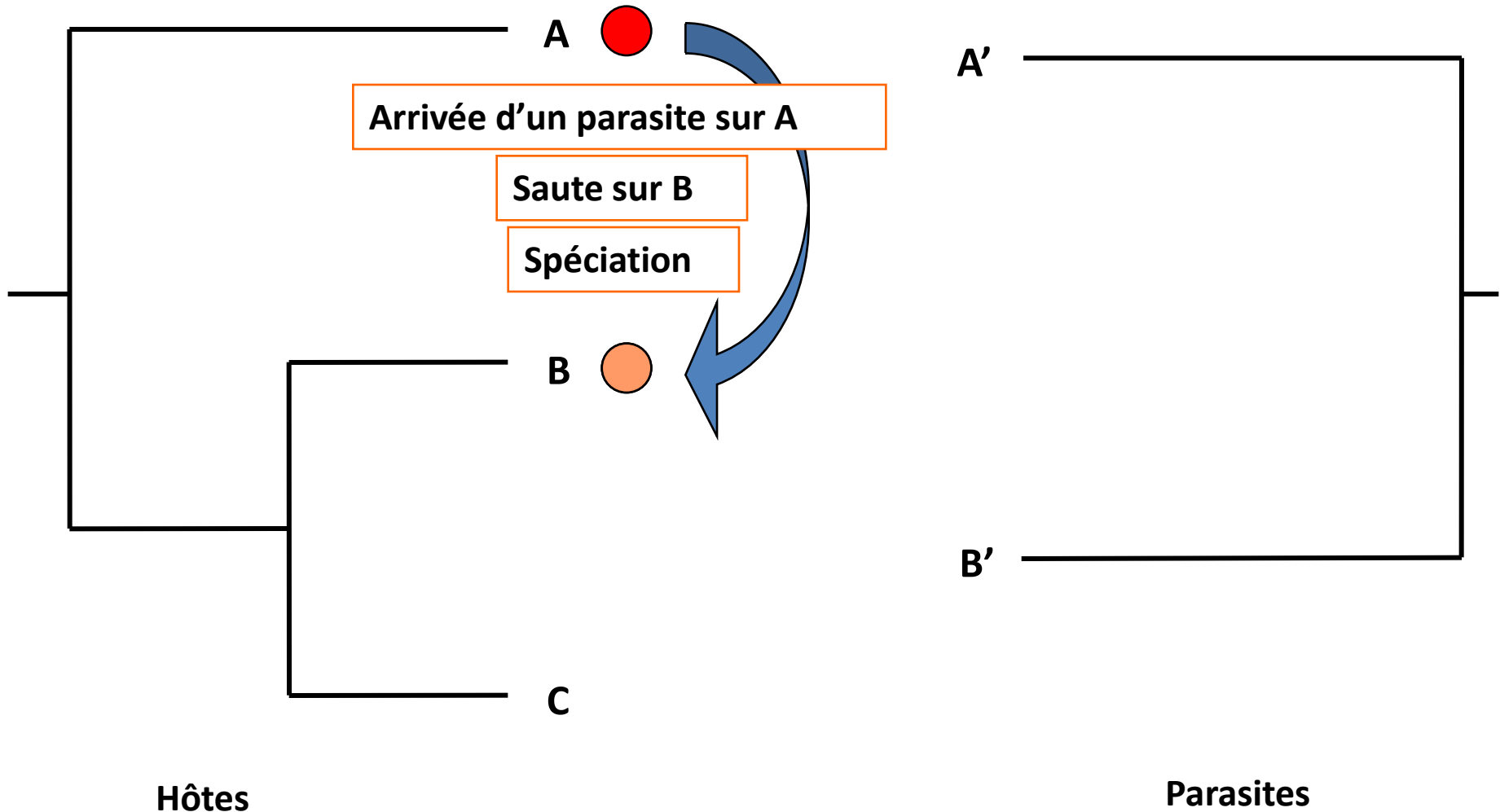
Hôtes

Simulate adaptive radiation of a parasite on a group of pre-existing hosts

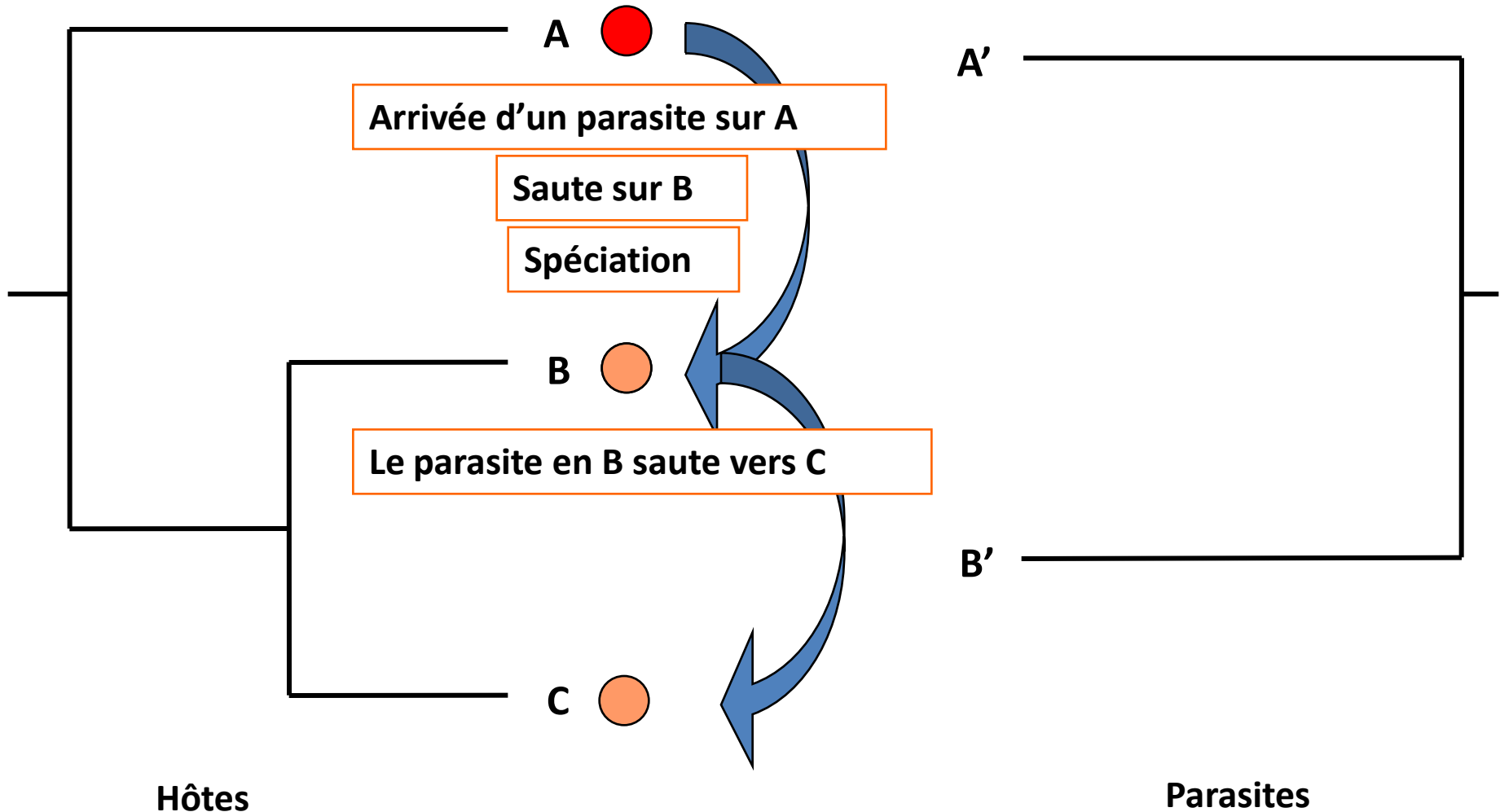


Hôtes

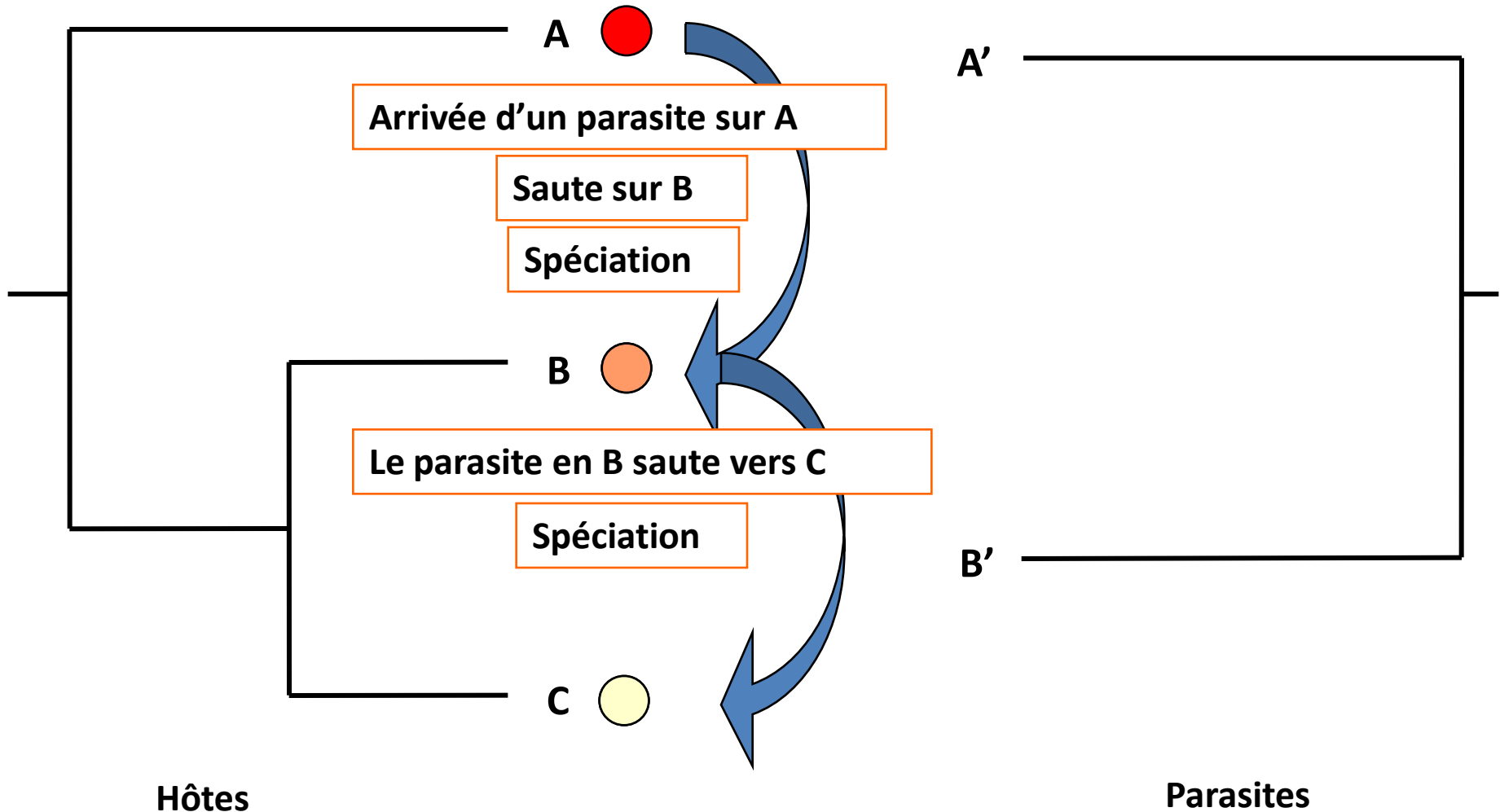
Simulate adaptive radiation of a parasite on a group of pre-existing hosts



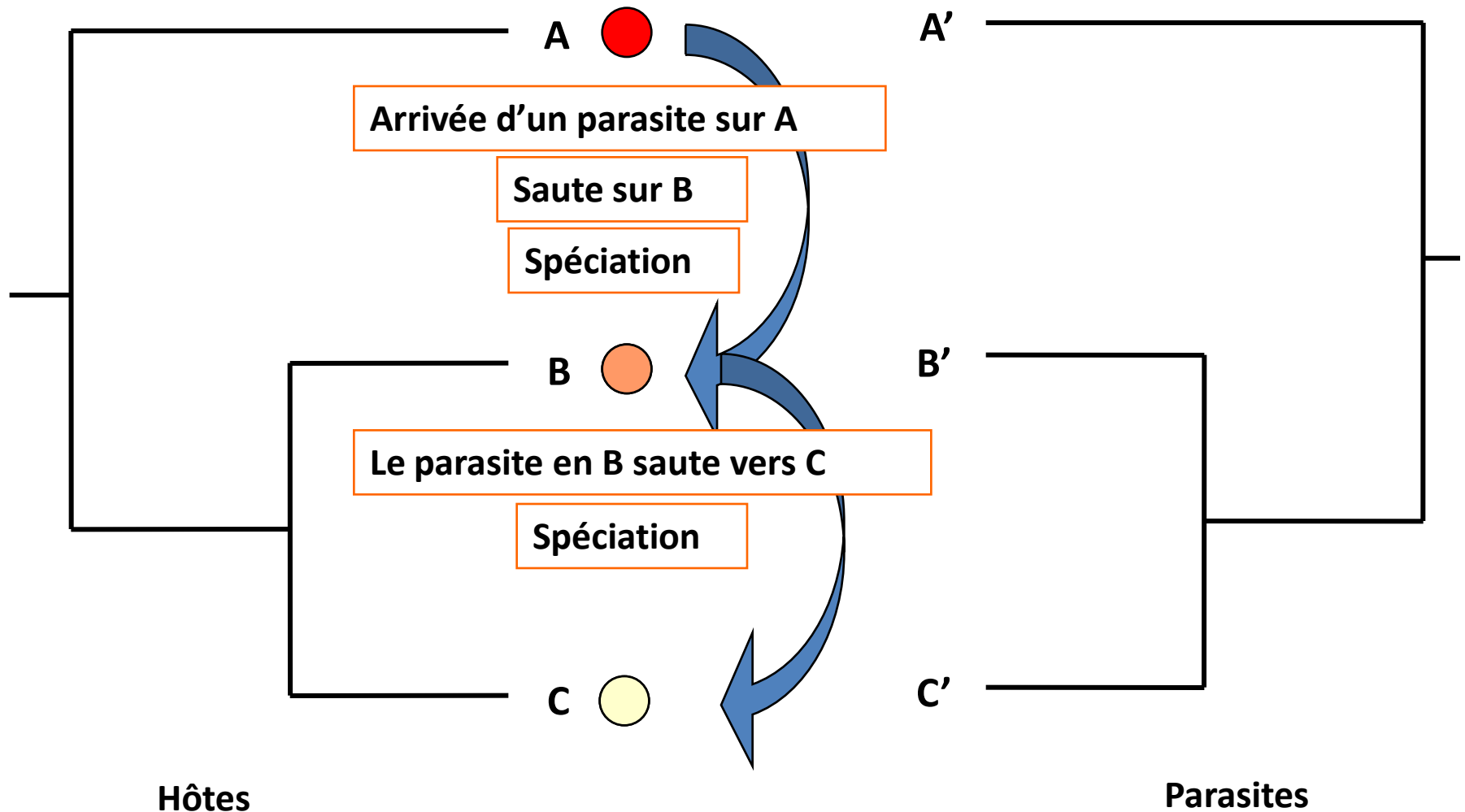
Simulate adaptive radiation of a parasite on a group of pre-existing hosts



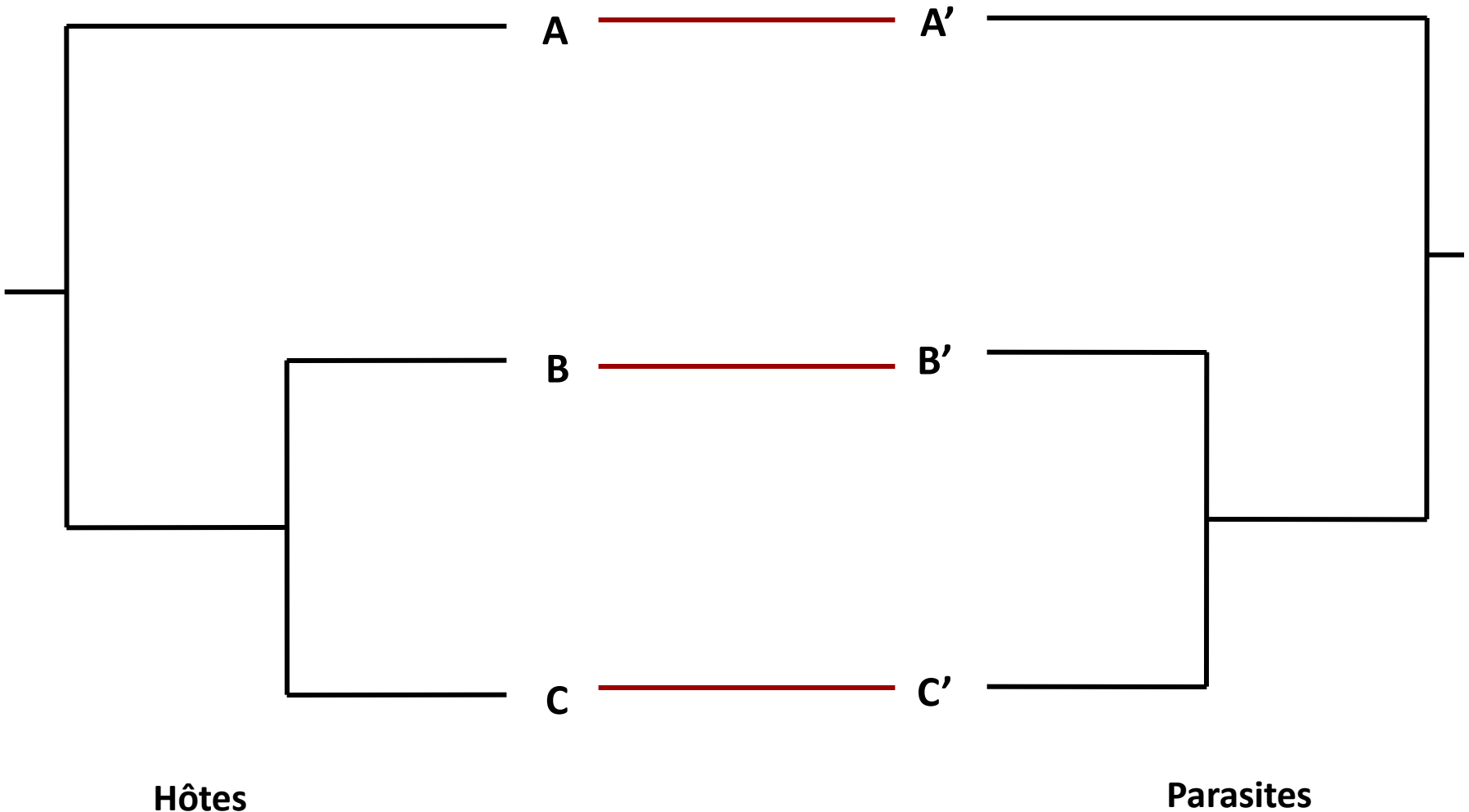
Simulate adaptive radiation of a parasite on a group of pre-existing hosts



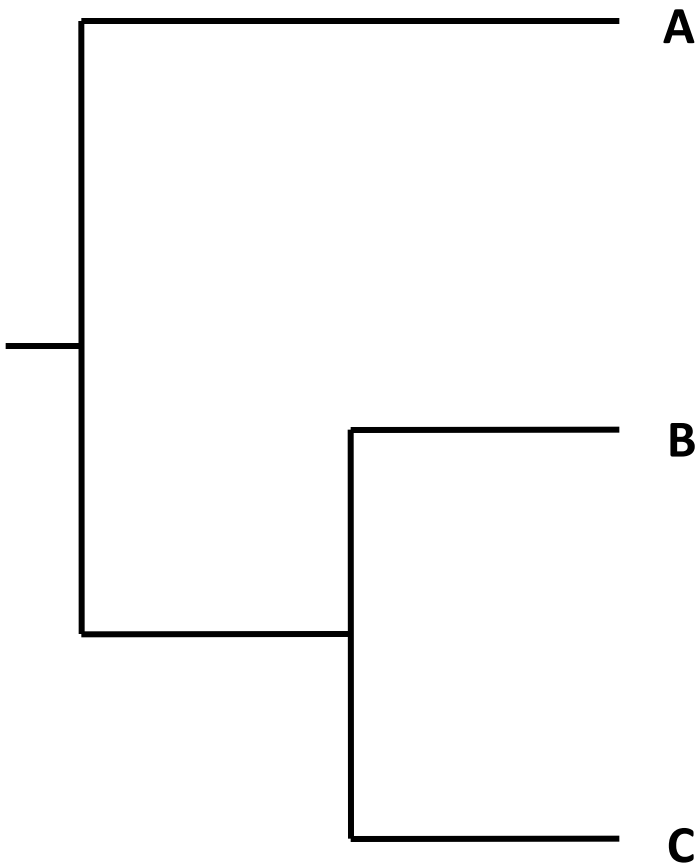
Simulate adaptive radiation of a parasite on a group of pre-existing hosts



Simulate adaptive radiation of a parasite on a group of pre-existing hosts

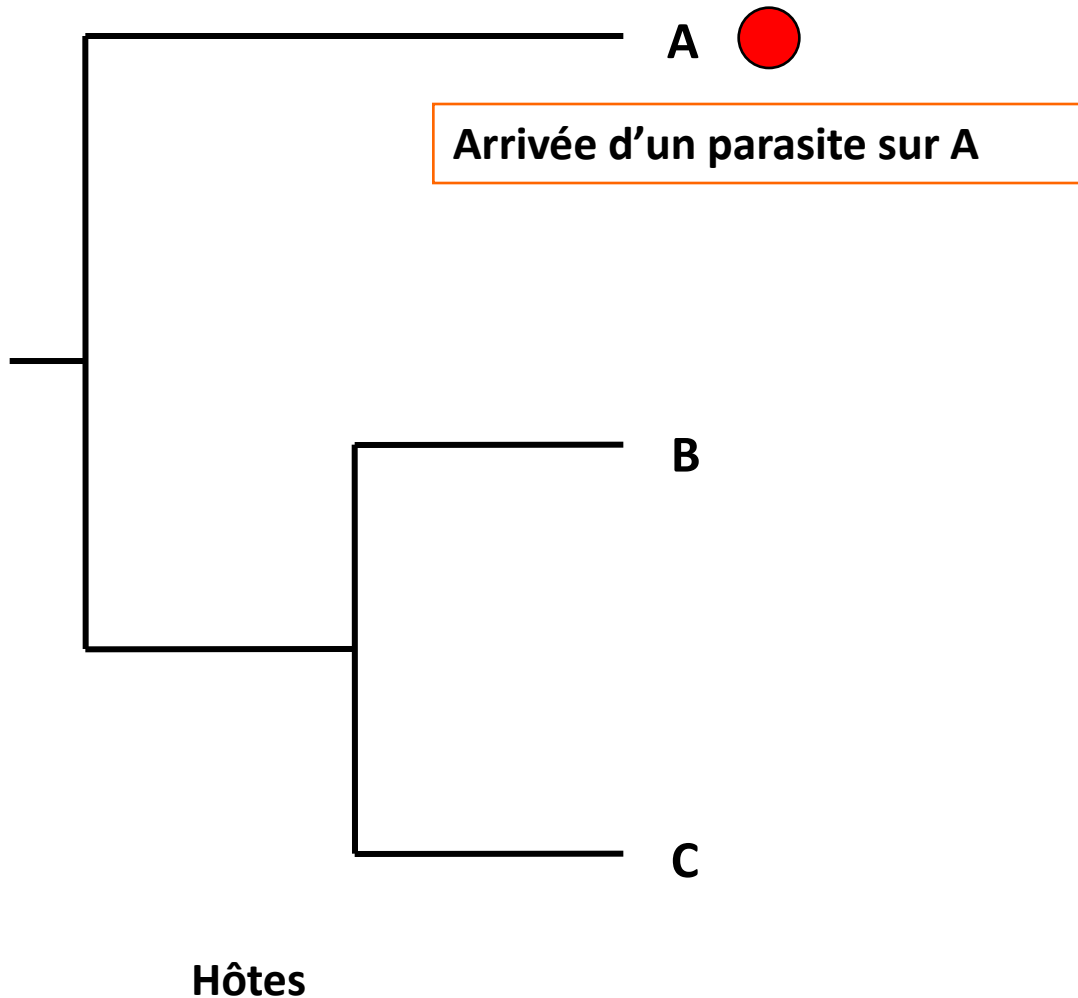


Simulate adaptive radiation of a parasite on a group of pre-existing hosts

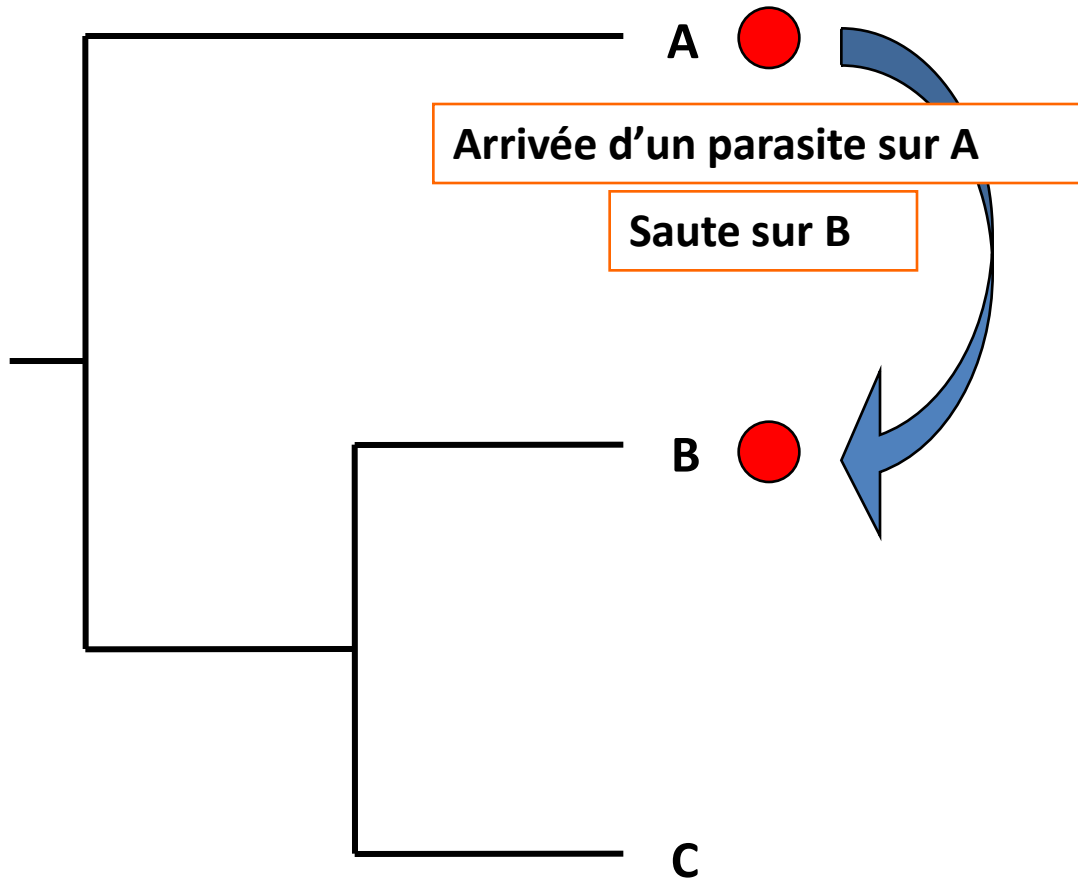


Hôtes

Simulate adaptive radiation of a parasite on a group of pre-existing hosts

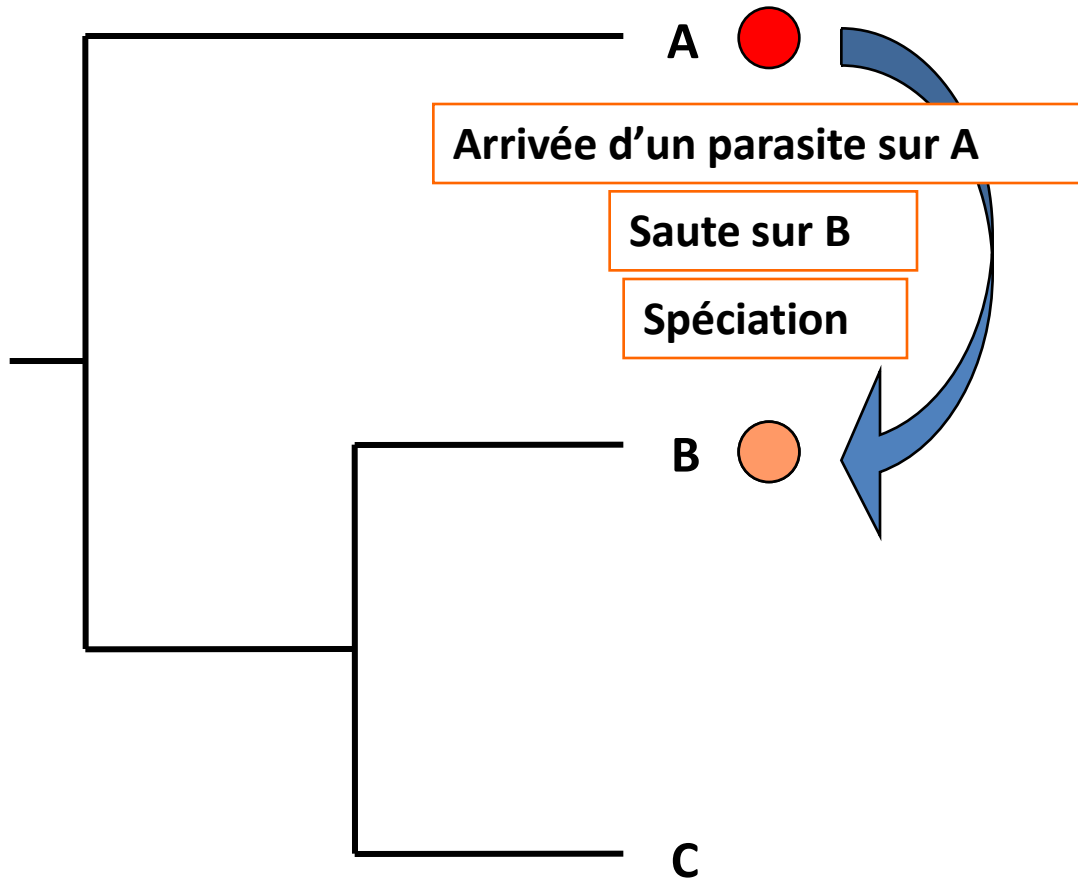


Simulate adaptive radiation of a parasite on a group of pre-existing hosts



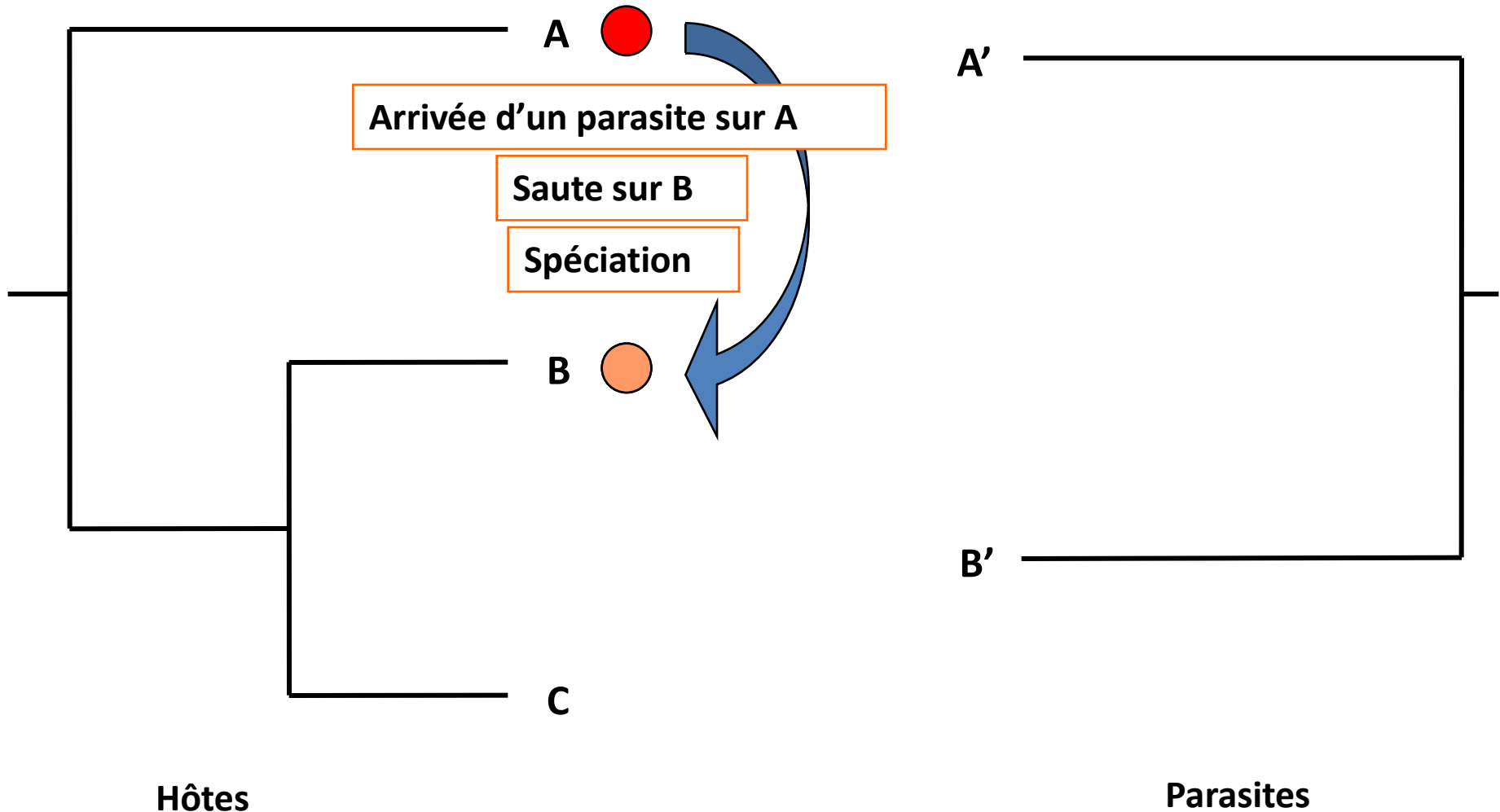
Hôtes

Simulate adaptive radiation of a parasite on a group of pre-existing hosts

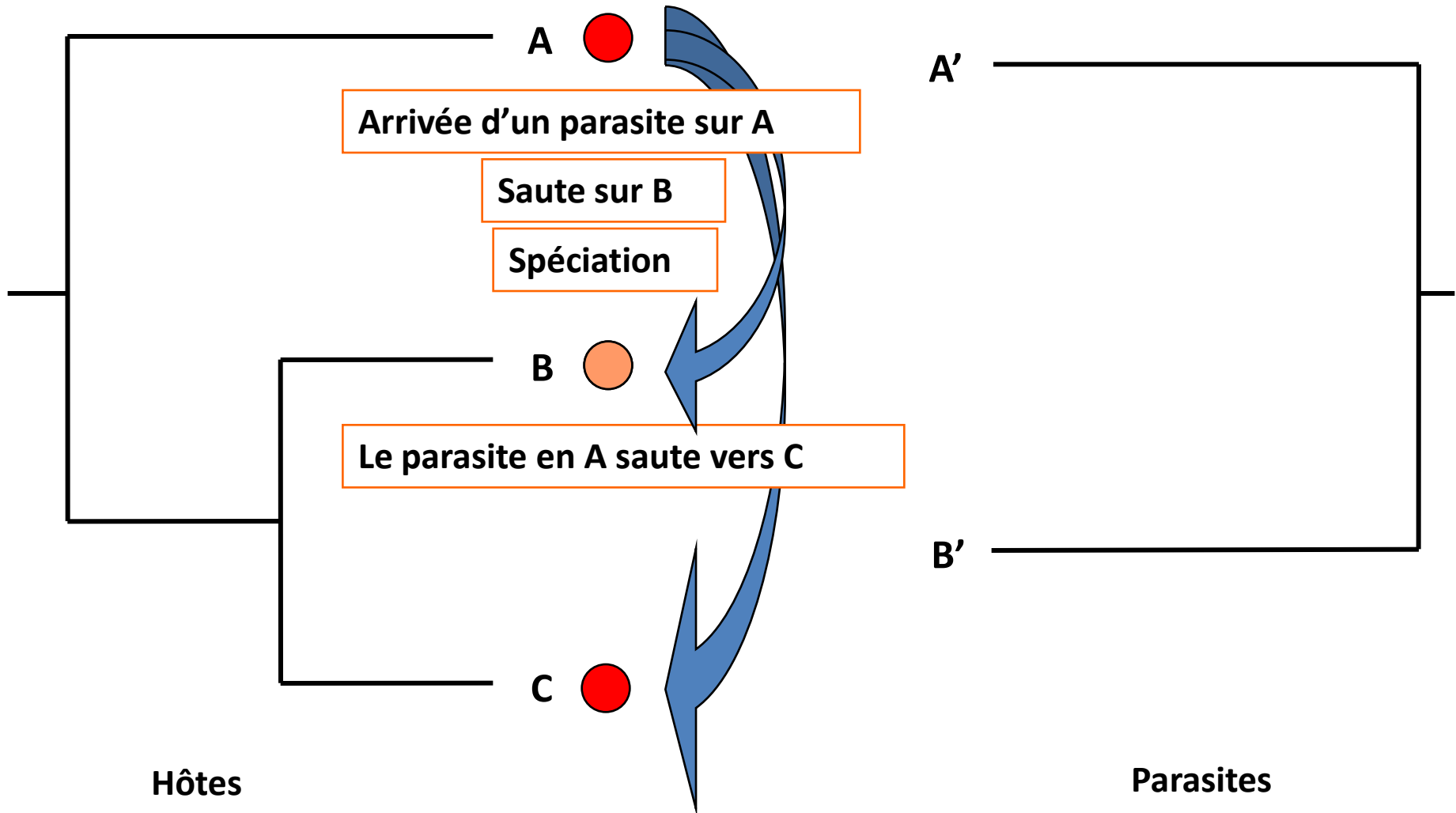


Hôtes

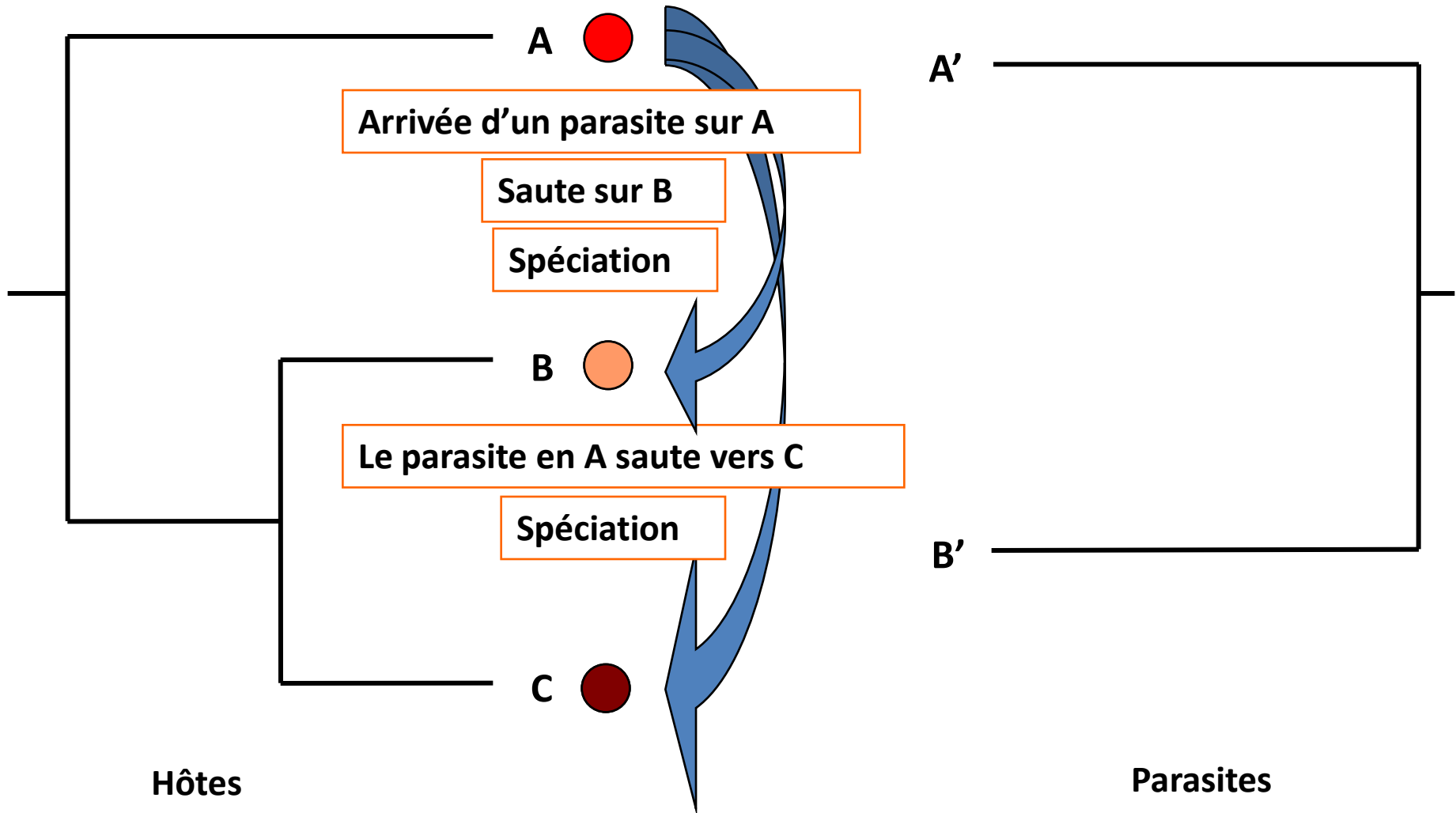
Simulate adaptive radiation of a parasite on a group of pre-existing hosts



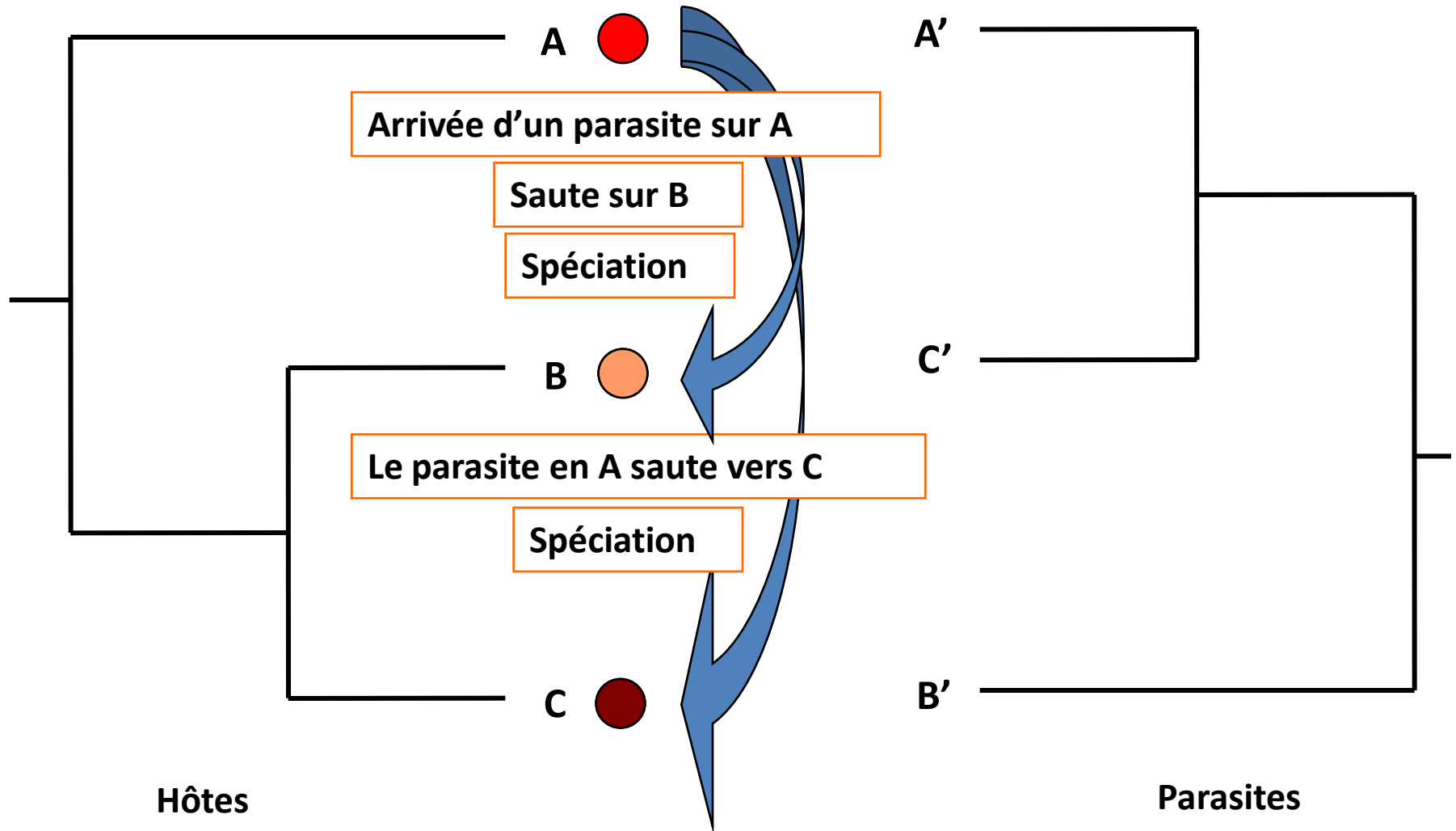
Simulate adaptive radiation of a parasite on a group of pre-existing hosts



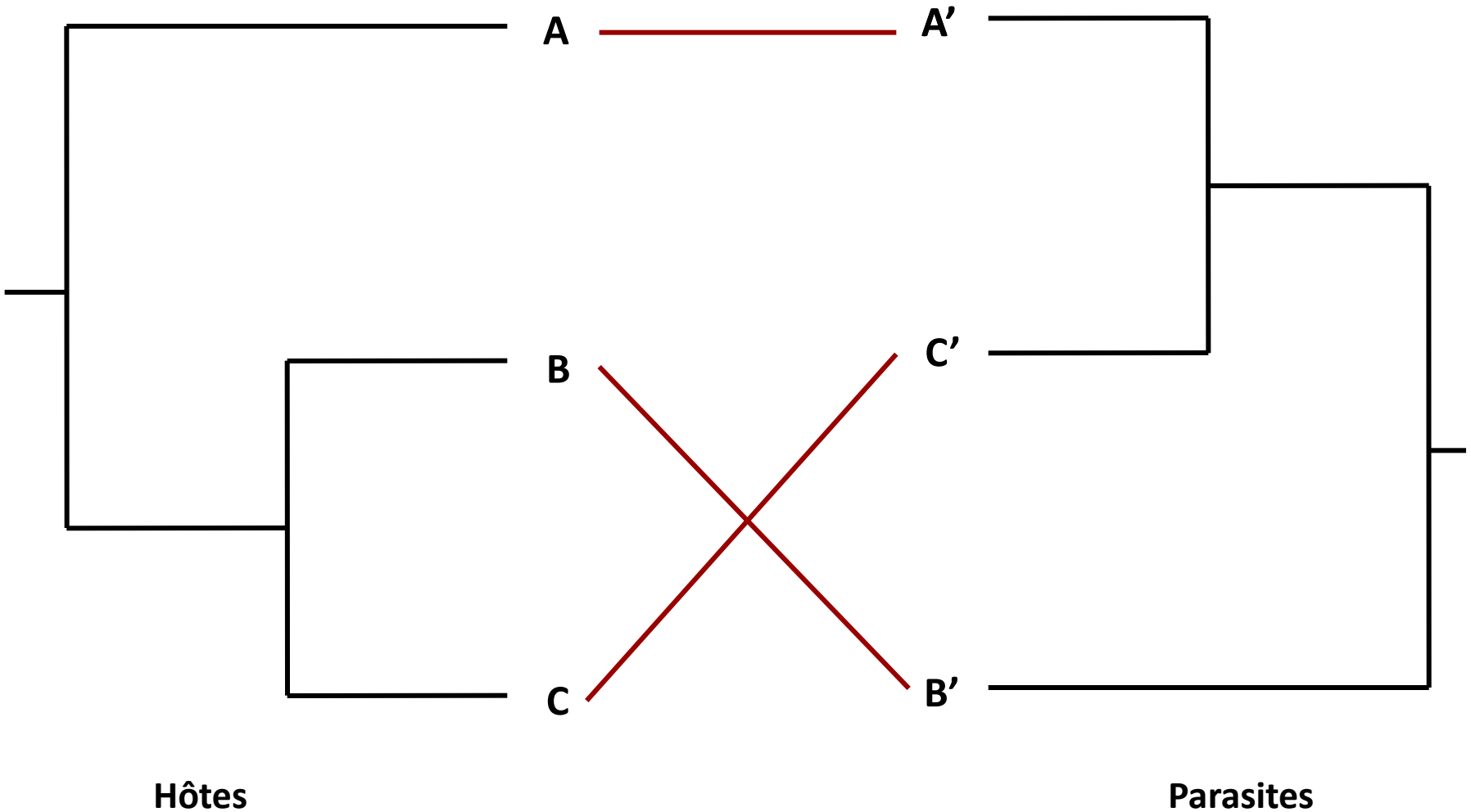
Simulate adaptive radiation of a parasite on a group of pre-existing hosts



Simulate adaptive radiation of a parasite on a group of pre-existing hosts



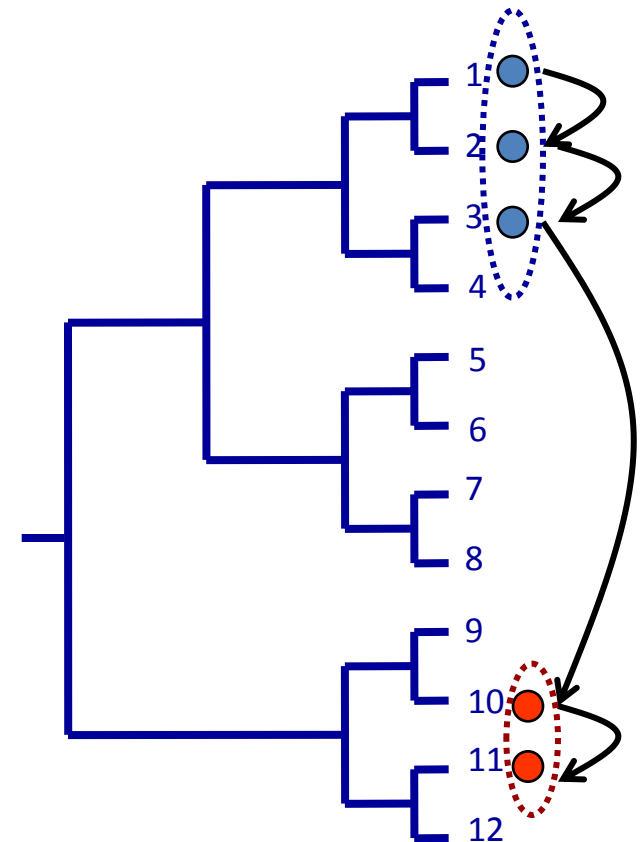
Simulate adaptive radiation of a parasite on a group of pre-existing hosts



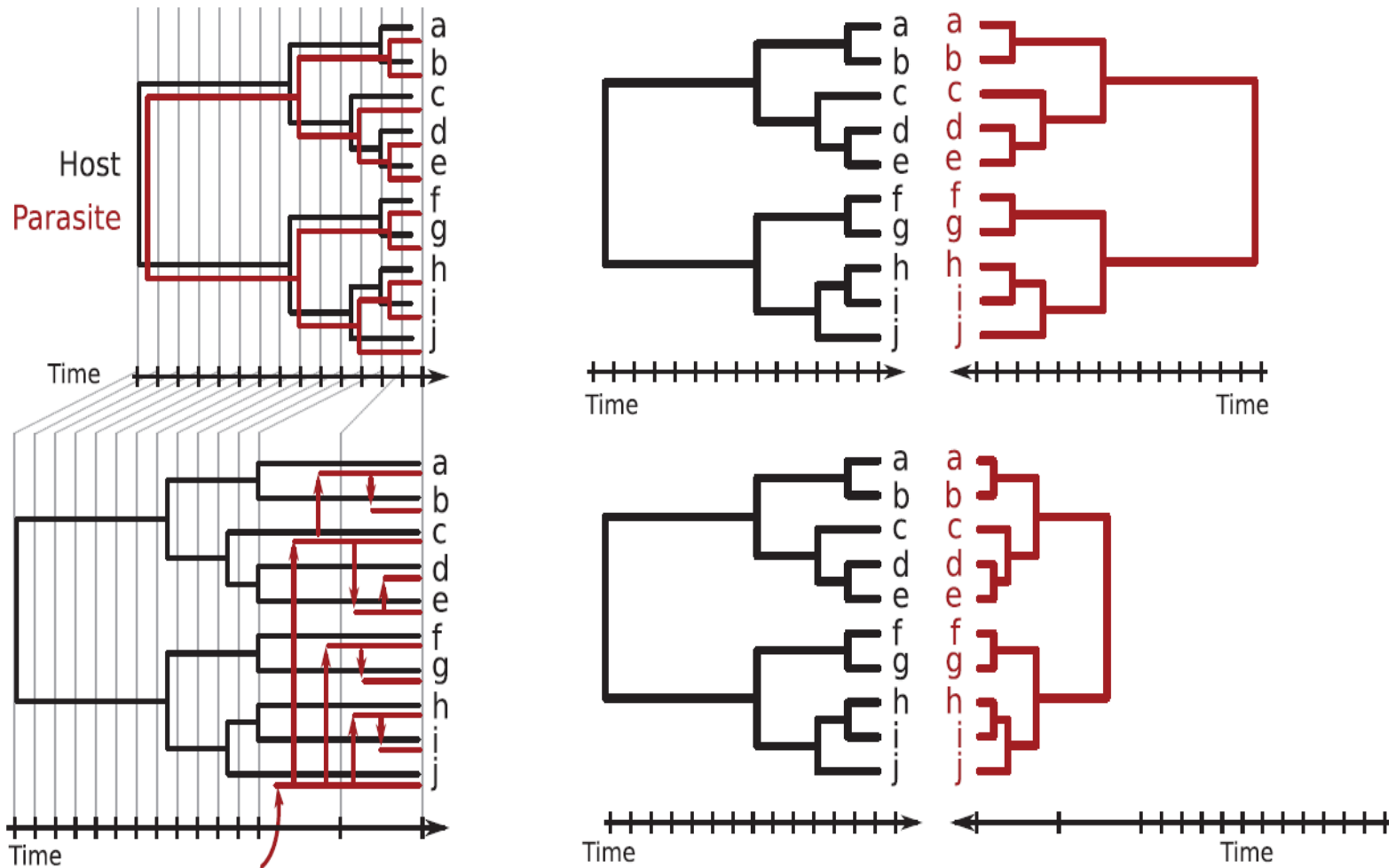
High congruence between host and parasite phylogenies can be obtained without cospeciation under plausible conditions

Highest probability for switches to closely related hosts

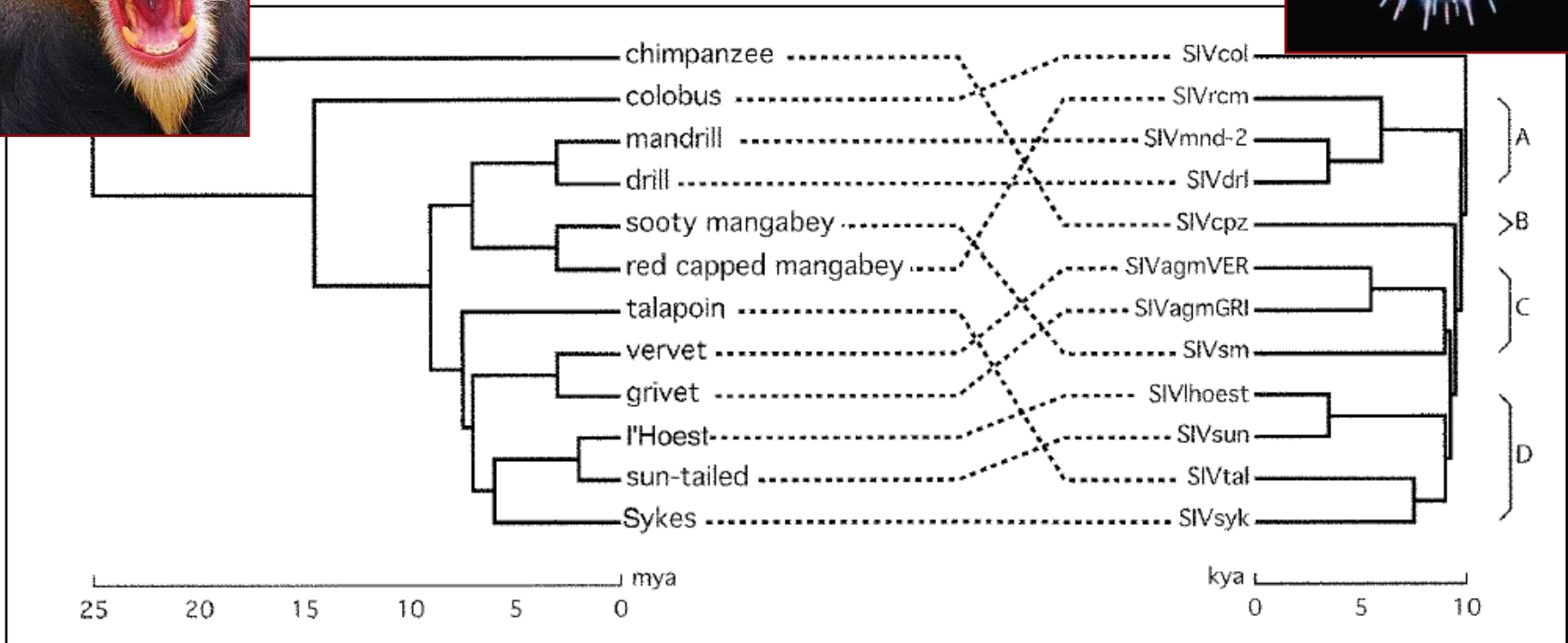
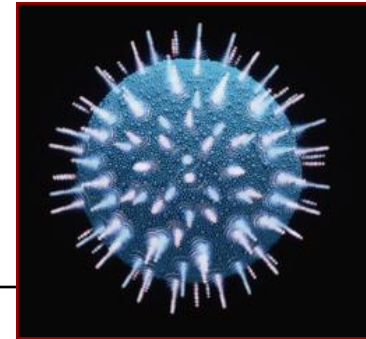
Faster speciation after distant switches than after close ones



The main difference between cospeciation-based and host-shift based congruence is the age of nodes in the trees



High congruence between host and parasite phylogenies can be obtained without cospeciation under plausible conditions



Co-phylogenetic methods

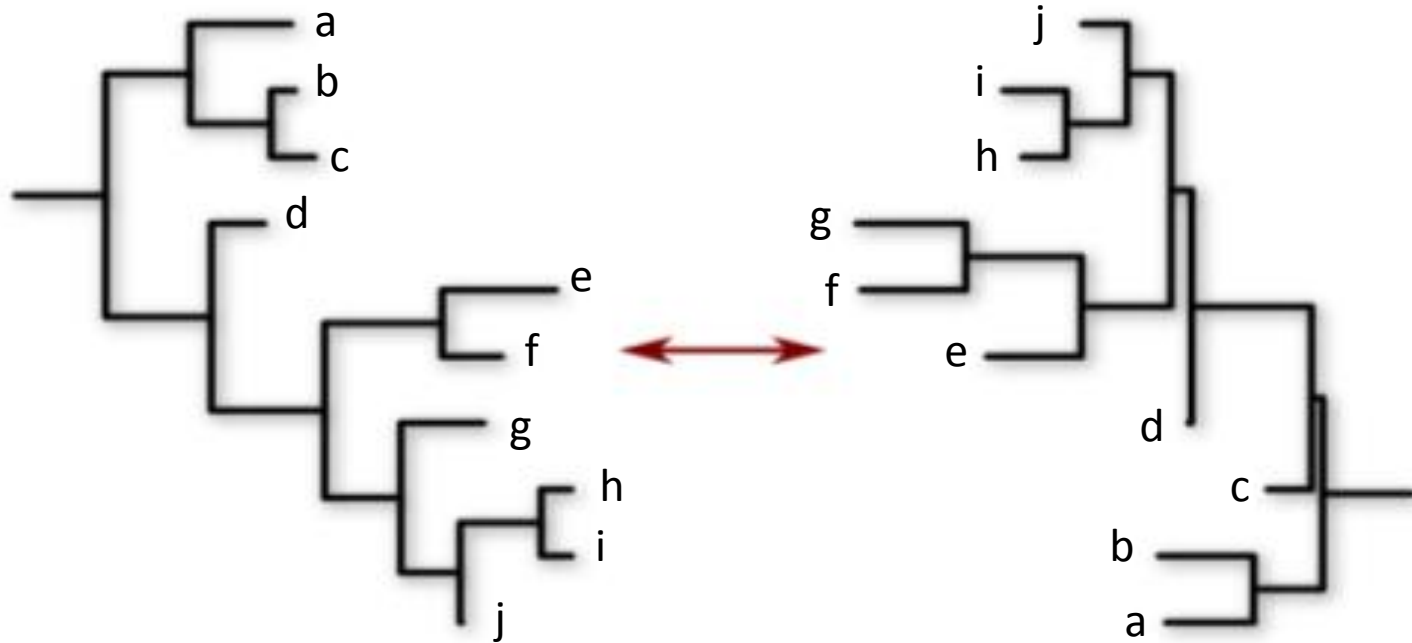
- Event- or cost-based
 - Estimate a scenario
 - Costs are associated to each event (cospeciation/duplication/host-shift/extinction)
 - Try to minimize cospeciation
 - Cospeciation is always less costly
- Topology- and distance-based methods
 - Based on comparison of a score with its distribution after permutations
 - No *a priori* on the reasons for the overall congruence (no events)
 - Often, *a posteriori* interpretation that congruence = cospeciation

Co-phylogenetic methods

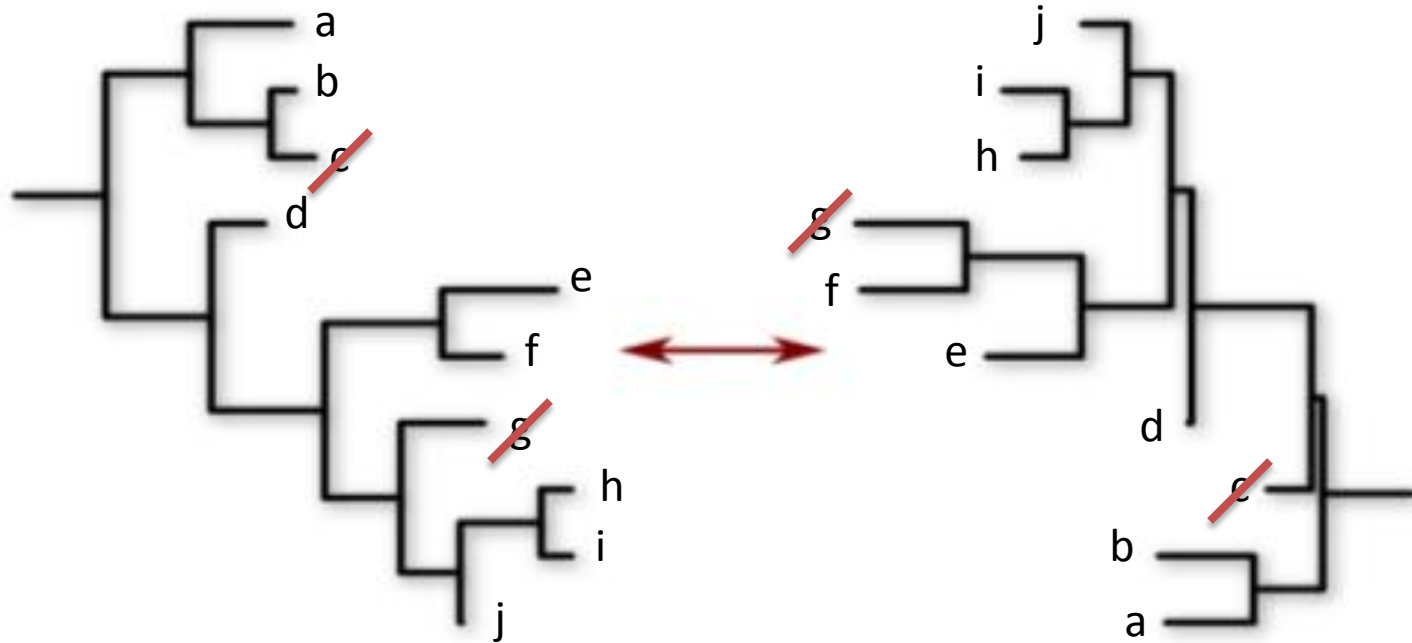
- Event- or cost-based
 - Estimate a scenario
 - Costs are associated to each event (cospeciation/duplication/host-shift/extinction)
 - Try to minimize cospeciation
 - Cospeciation is always less costly
- Topology- and distance-based methods
 - Based on comparison of a score with its distribution after permutations
 - No *a priori* on the reasons for the overall congruence (no events)
 - Often, *a posteriori* interpretation that congruence = cospeciation

I_{cong} : New topological distance based on expected size of the MAST (no permutations)

Maximum Agreement Subtree (MAST)



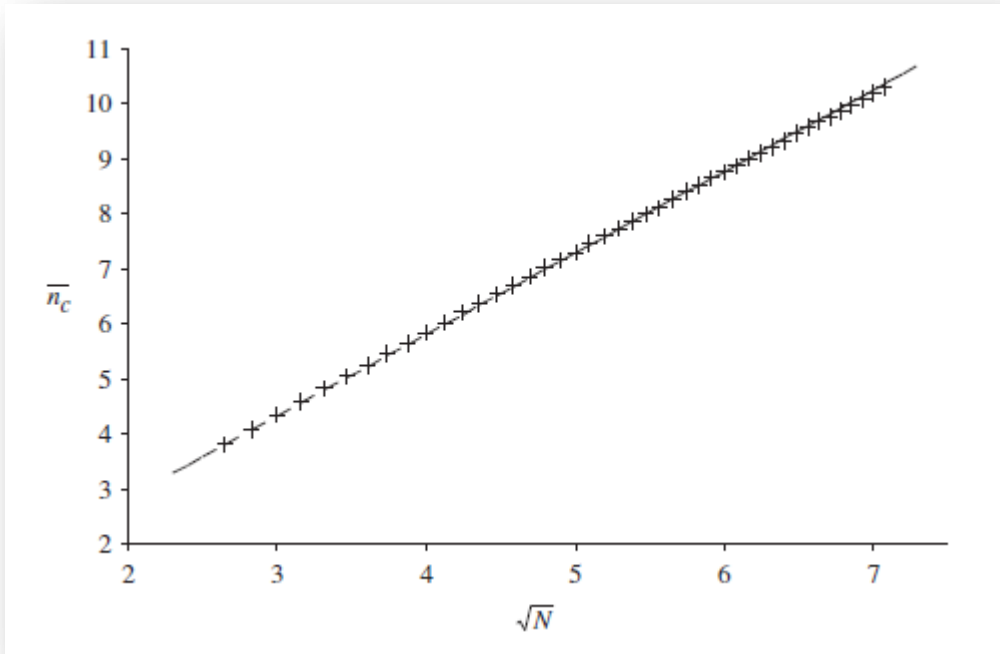
Maximum Agreement Subtree (MAST)



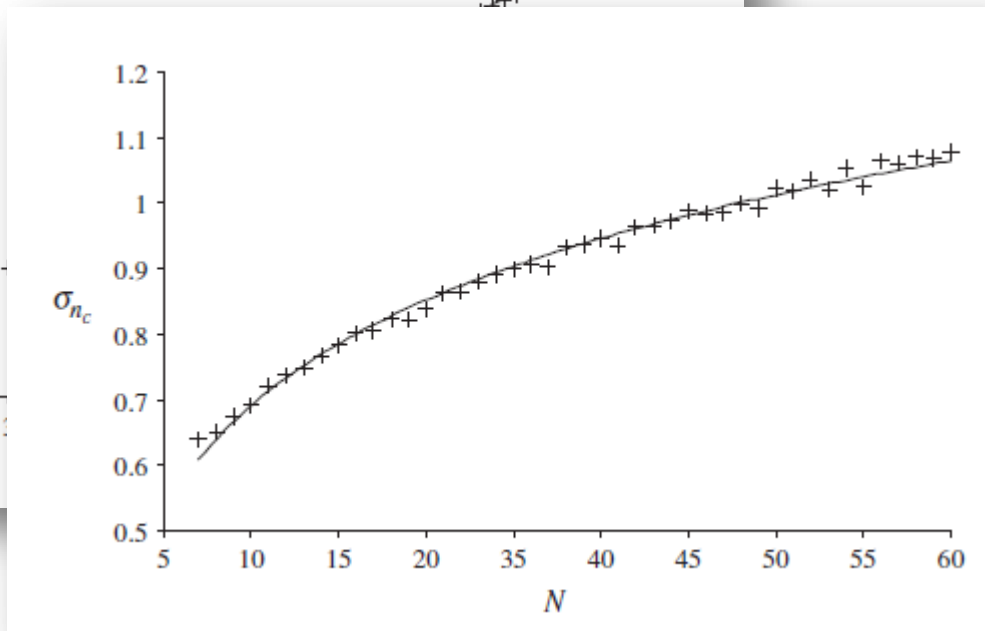
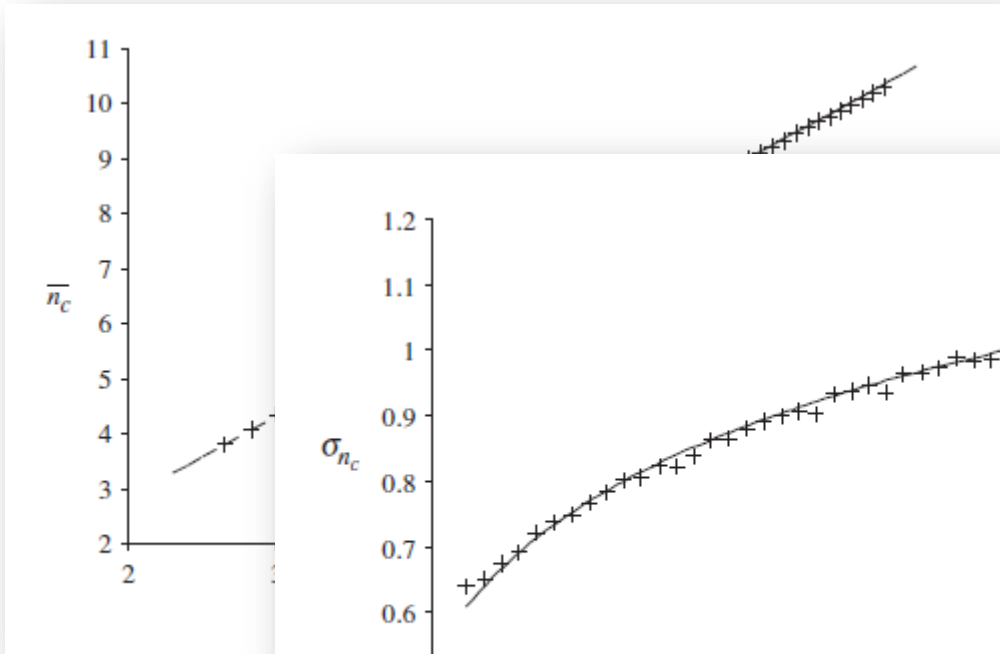
$N = 10$

Size of the MAST (n_c) = 8

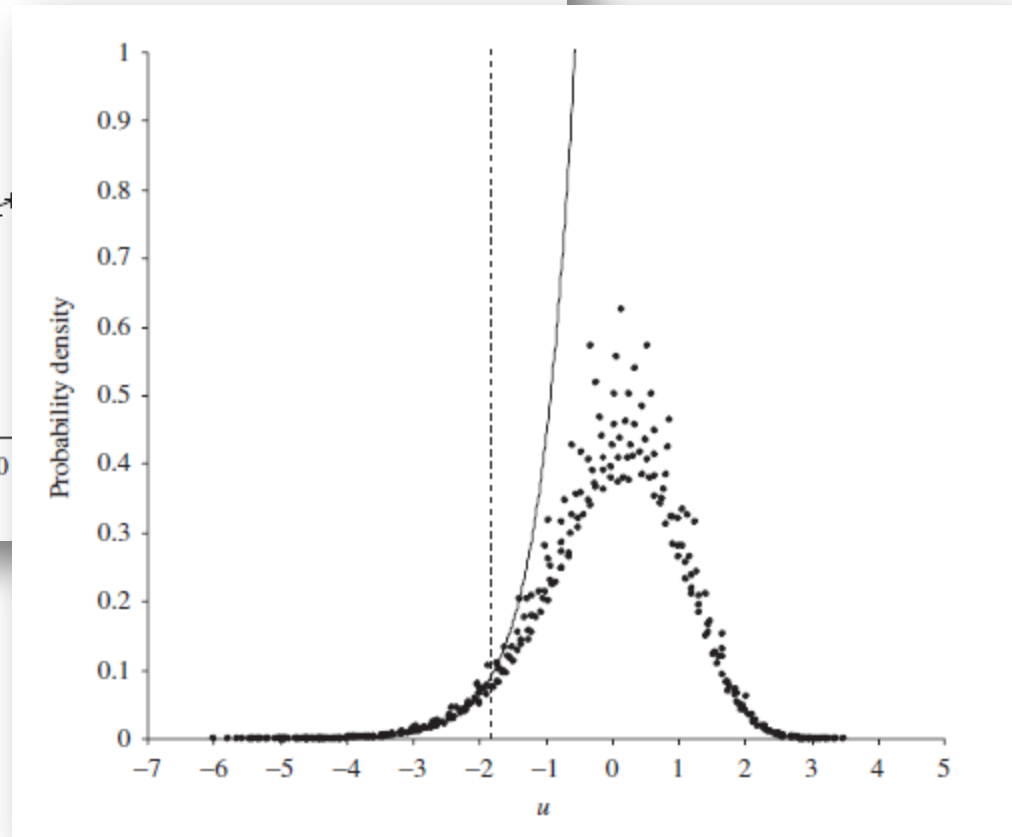
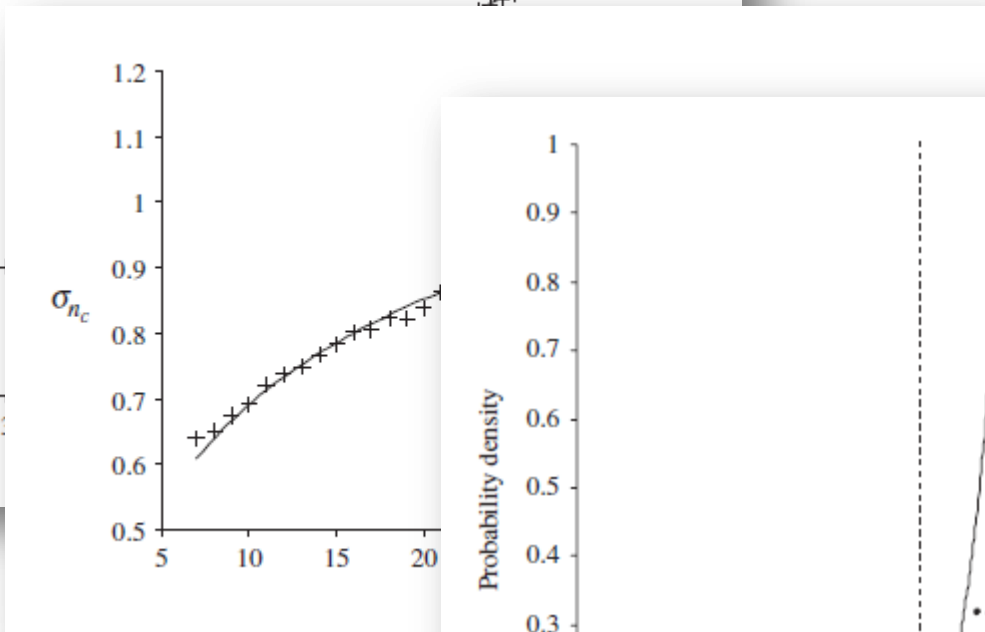
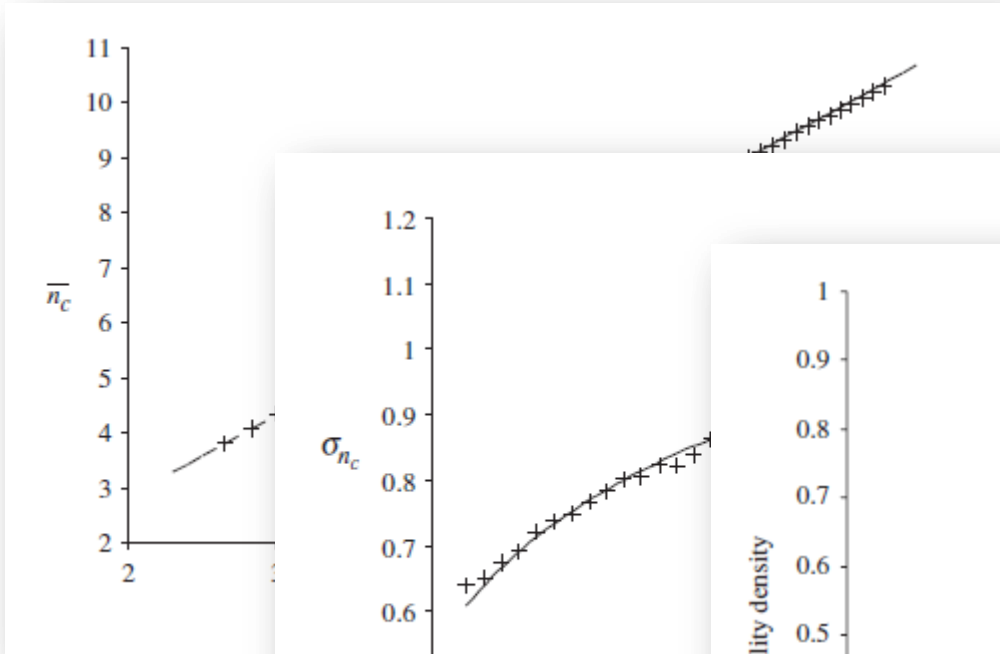
I_{cong}



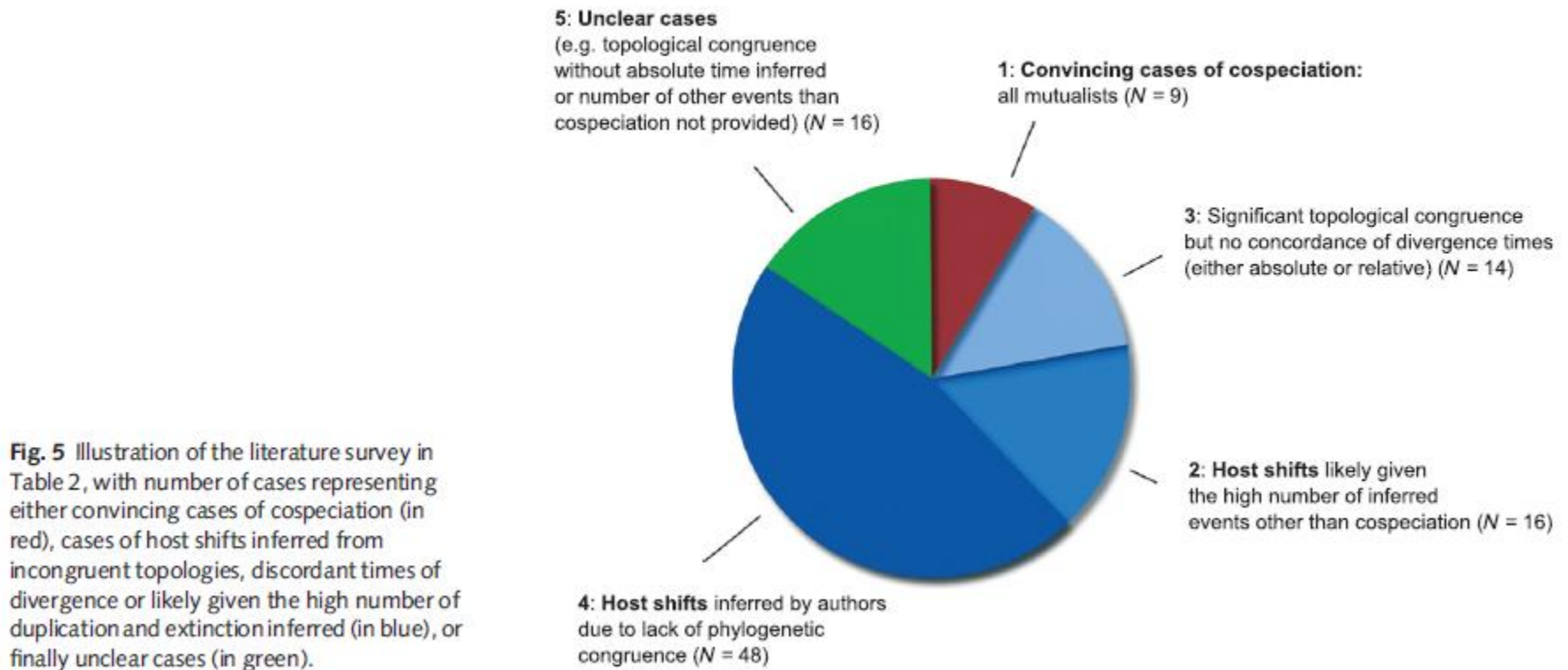
I_{cong}



I_{cong}



Literature survey reveals prevalence of host-shift speciation over cospeciation



Link between [coevolution, specialization, speciation] and [cospeciation or host-shift speciation]

- Coevolution -> specialization -> speciation

(reviewed in Summers et al. 2003)

Speed depends on parasite and host generation time, dispersal rates, effective pop size, etc... (Huyse et al. 2005)

→ Specialization of two parasite lineages on sister host species may result in a cospeciation event.

BUT is it what prevails in the long term?

MAYBE NOT / APPARENTLY NOT

Link between [coevolution, specialization, speciation] and [cospeciation or host-shift speciation]

When new species are formed:

loss of associated parasites (Enemy release, Kean & Crowley 2002, Genton et al. 2005)

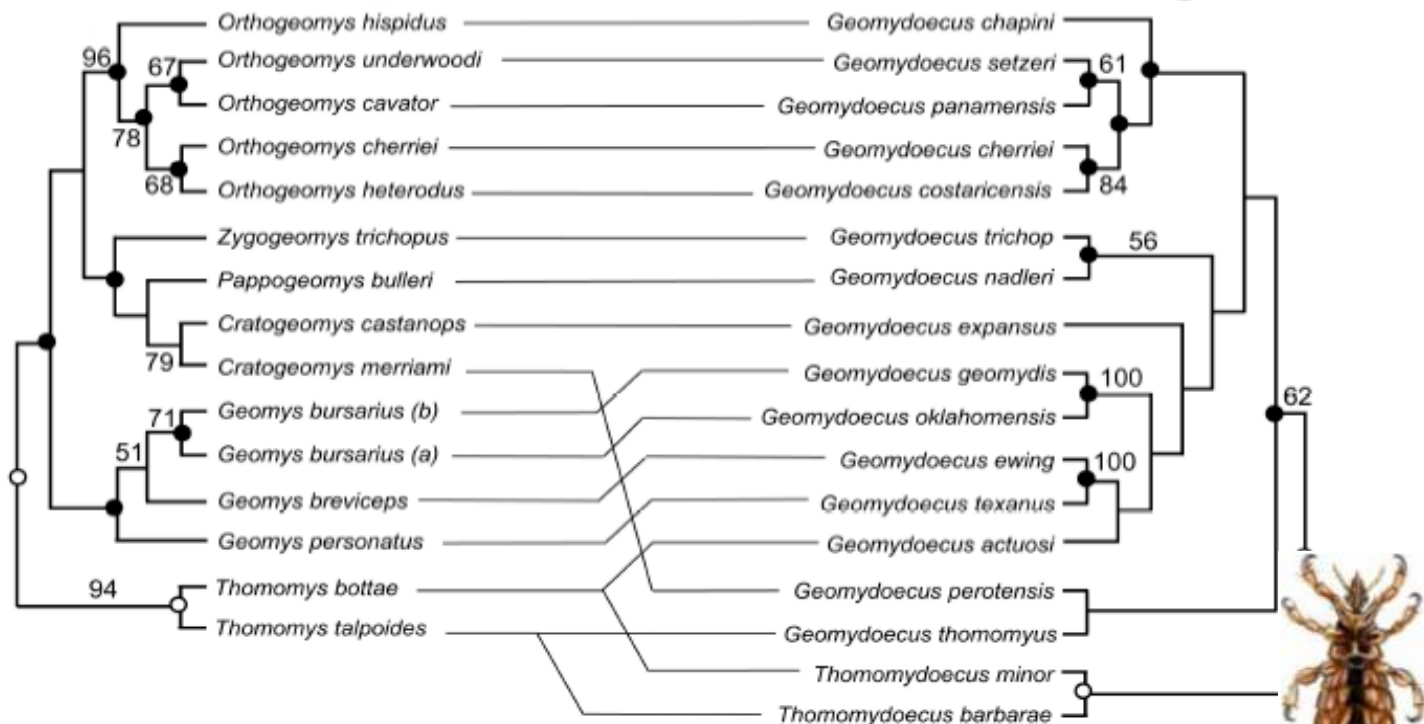
smaller population -> not compatible with specialist parasites (de Castro & Bolker, 2005)

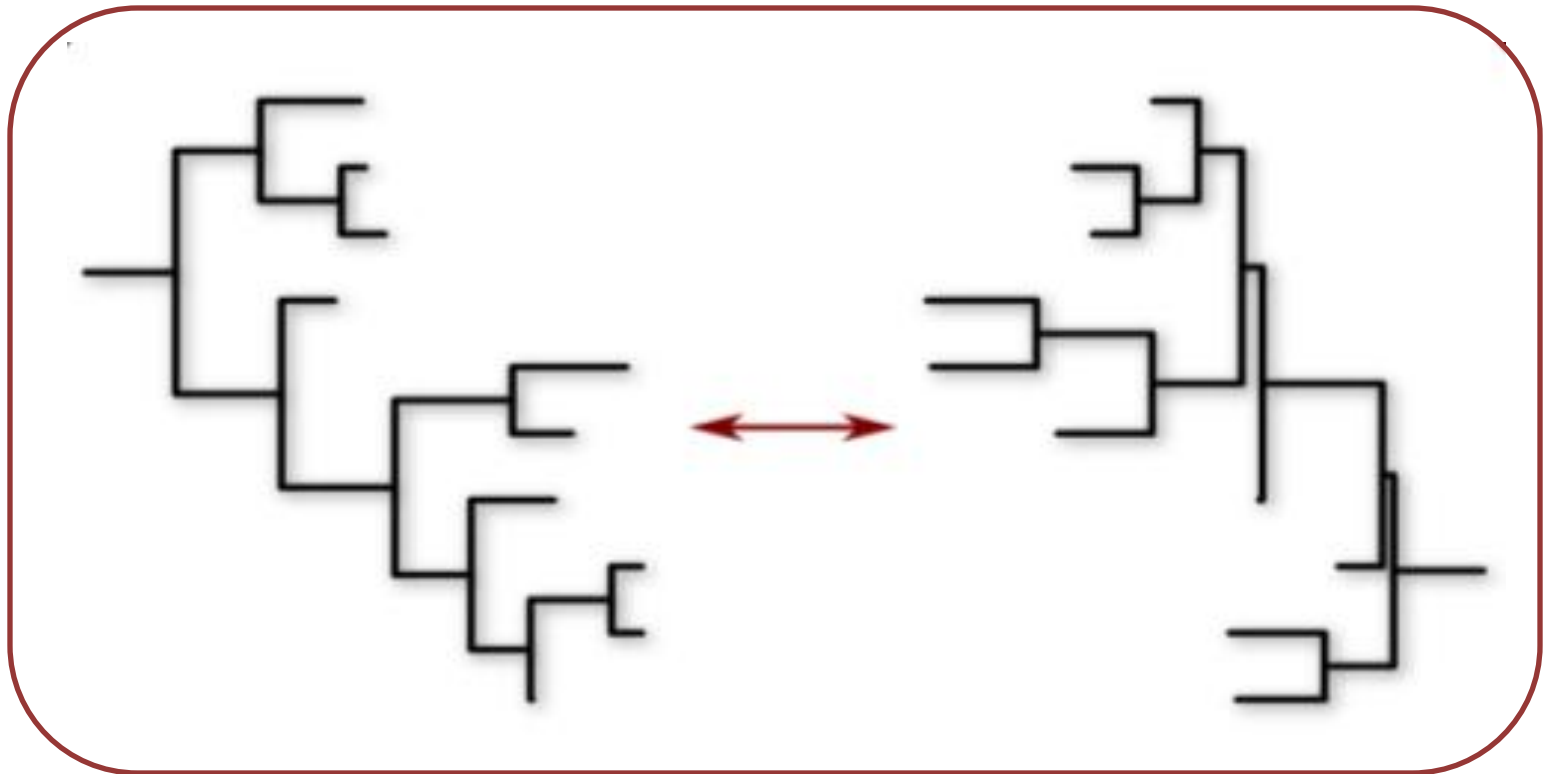
but coevolution hinders the persistence of generalists...

Against the idea coevolution leading to cospeciation.

MODEL? Under what conditions does host-parasite coevolution leads to cospeciation or to speciation by host-shifts?

Link between [coevolution, specialization, speciation] and [cospeciation or host-shift speciation]





Species tree \leftrightarrow Species tree

Gene tree \leftrightarrow Gene tree

Species tree \leftrightarrow Gene tree

Multiple gene trees

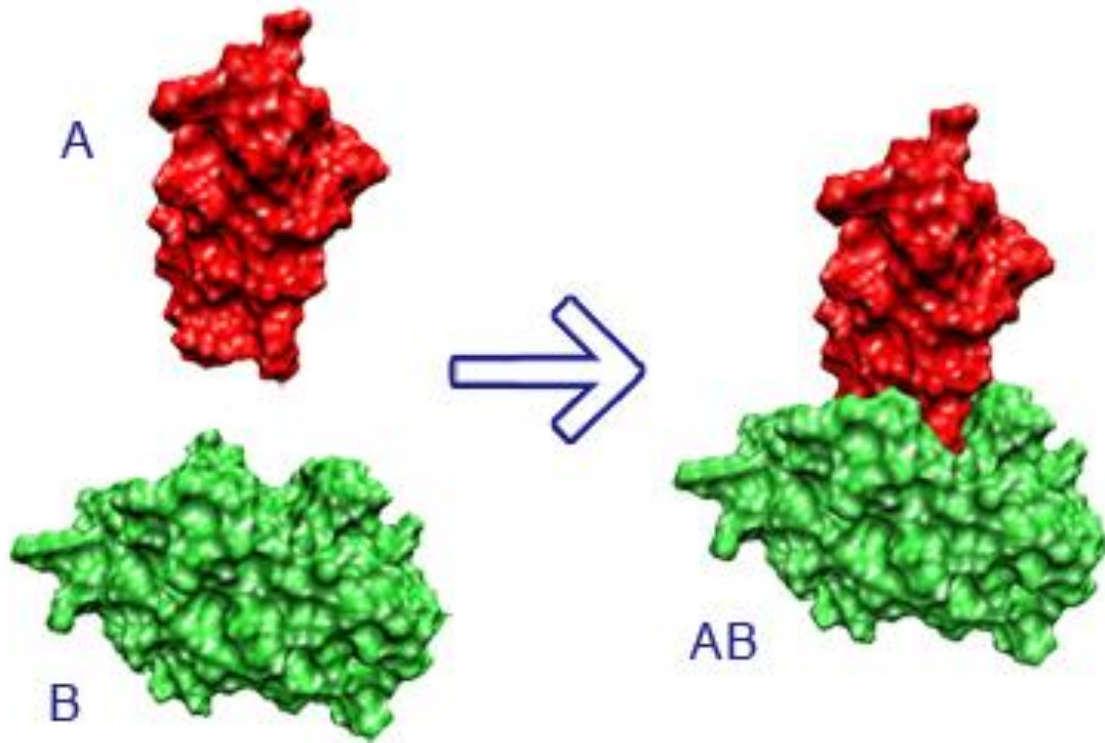
Host parasite coevolutionary studies

Protein-Protein interaction detection

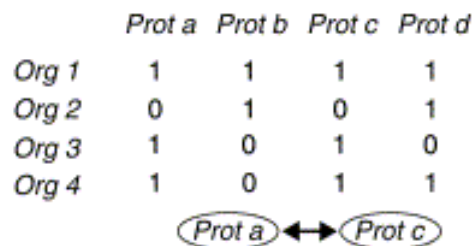
Reconciliation analyses

Phylogenomic studies

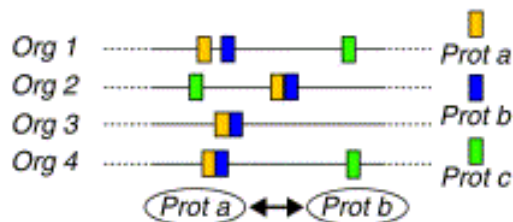
Two proteins involved in the same complex/pathway are expected to be coevolving



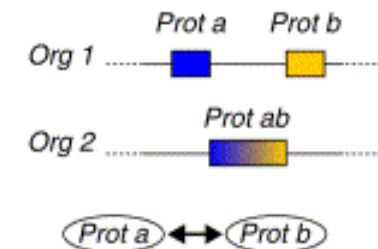
(a) Phylogenetic profiles



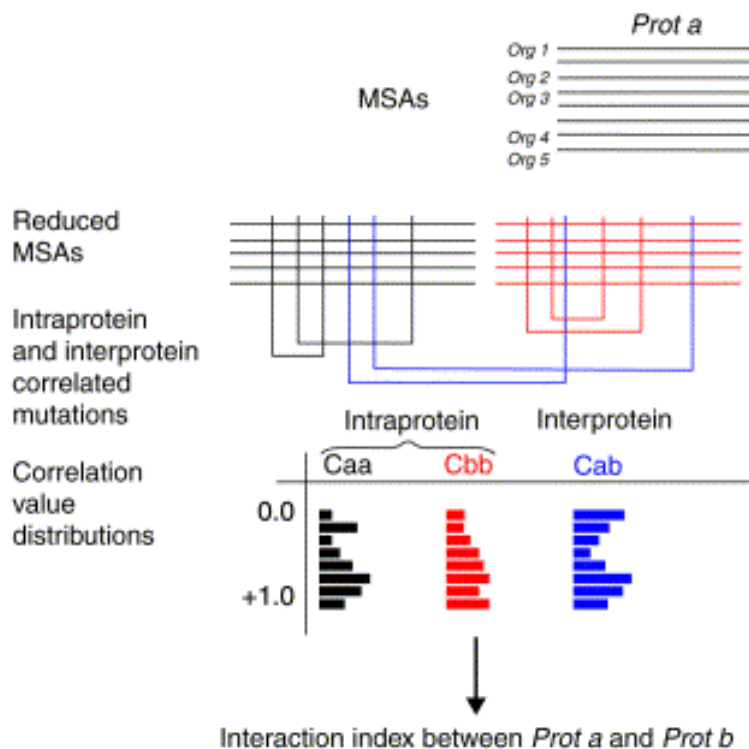
(b) Conservation of gene neighborhood



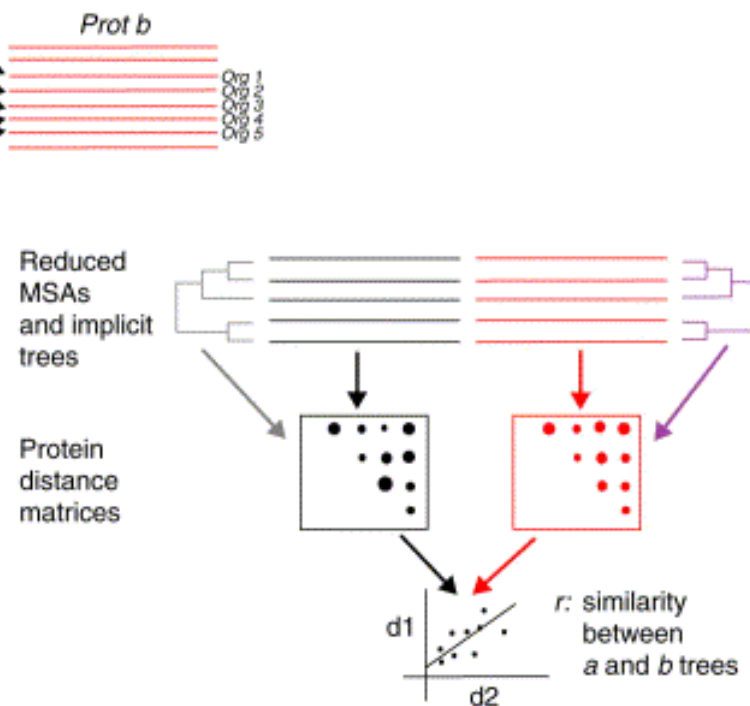
(c) Gene fusion



(e) Correlated mutations



(d) Similarity of phylogenetic trees

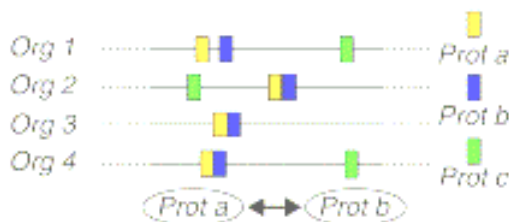


(a) Phylogenetic profiles

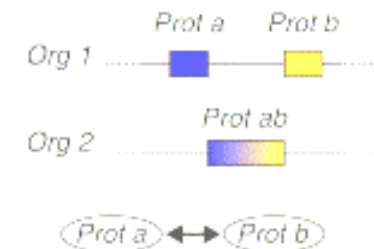
	Prot a	Prot b	Prot c	Prot d
Org 1	1	1	1	1
Org 2	0	1	0	1
Org 3	1	0	1	0
Org 4	1	0	1	1

(Prot a) ↔ (Prot c)

(b) Conservation of gene neighborhood



(c) Gene fusion



Still many false positive and false negative when predicting PPI on the interactome of *E. coli*.

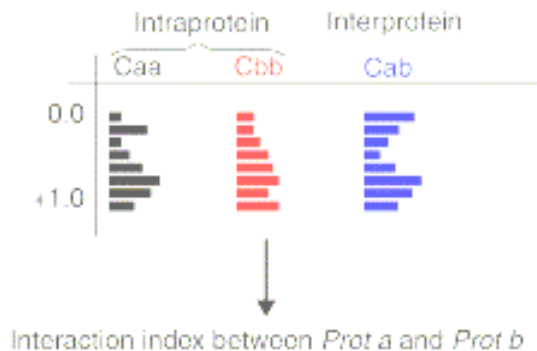
DATA: 2177 *E. coli* proteins and their orthologs in 115 other prokaryotic genomes

(e) Correlate

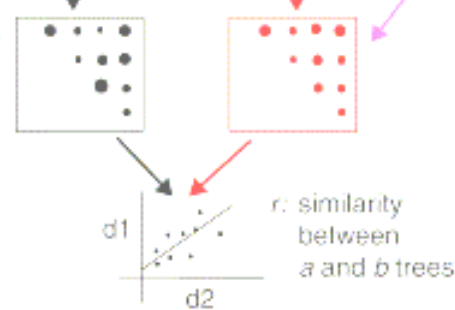
Reduced MSAs

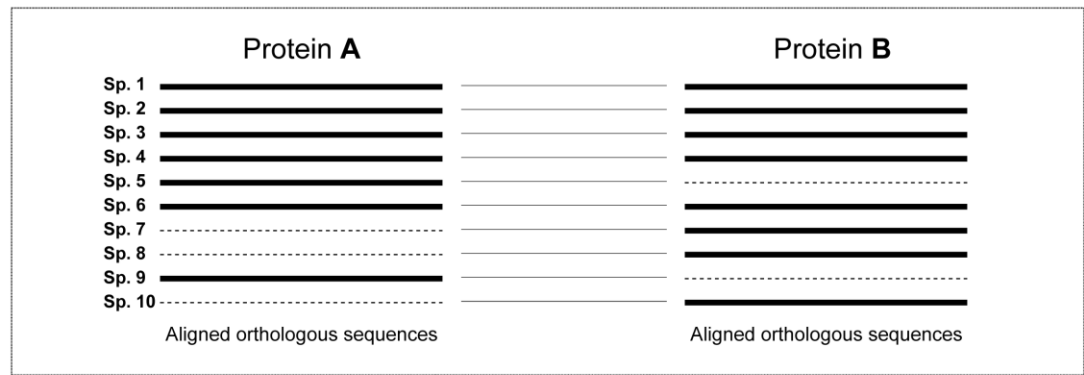
Intraprotein and interprotein correlated mutations

Correlation value distributions



Protein distance matrices

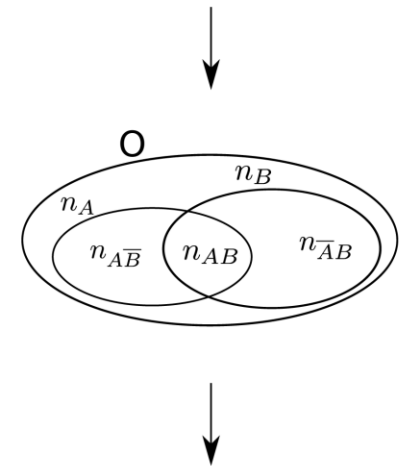
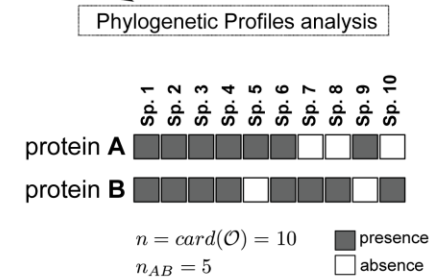
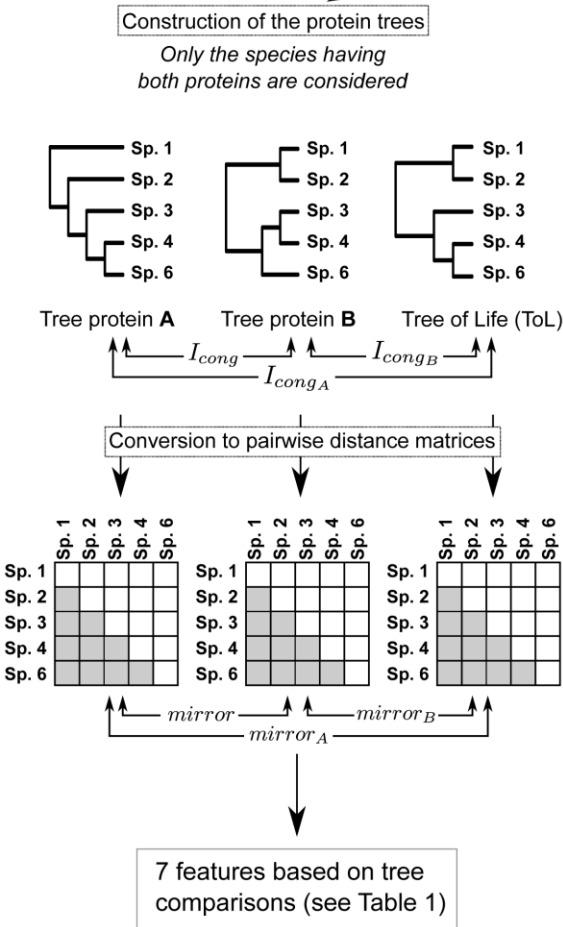




- Development of new descriptors of coevolution

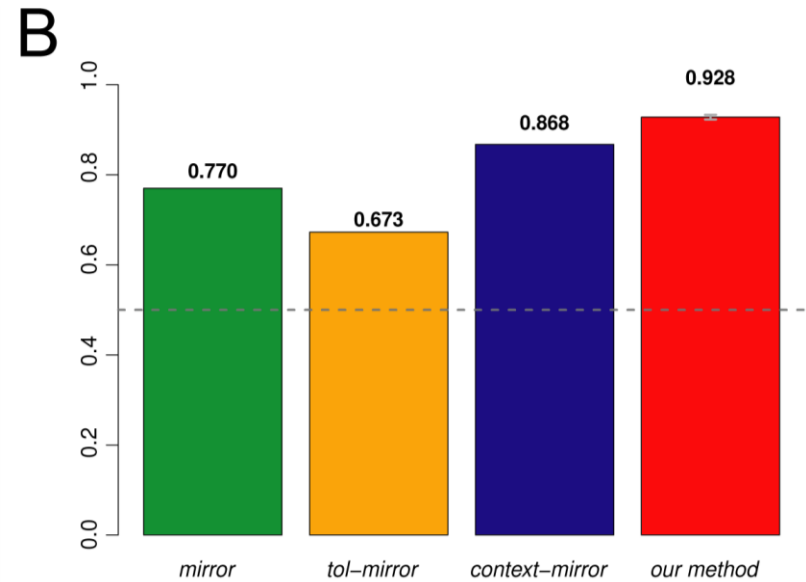
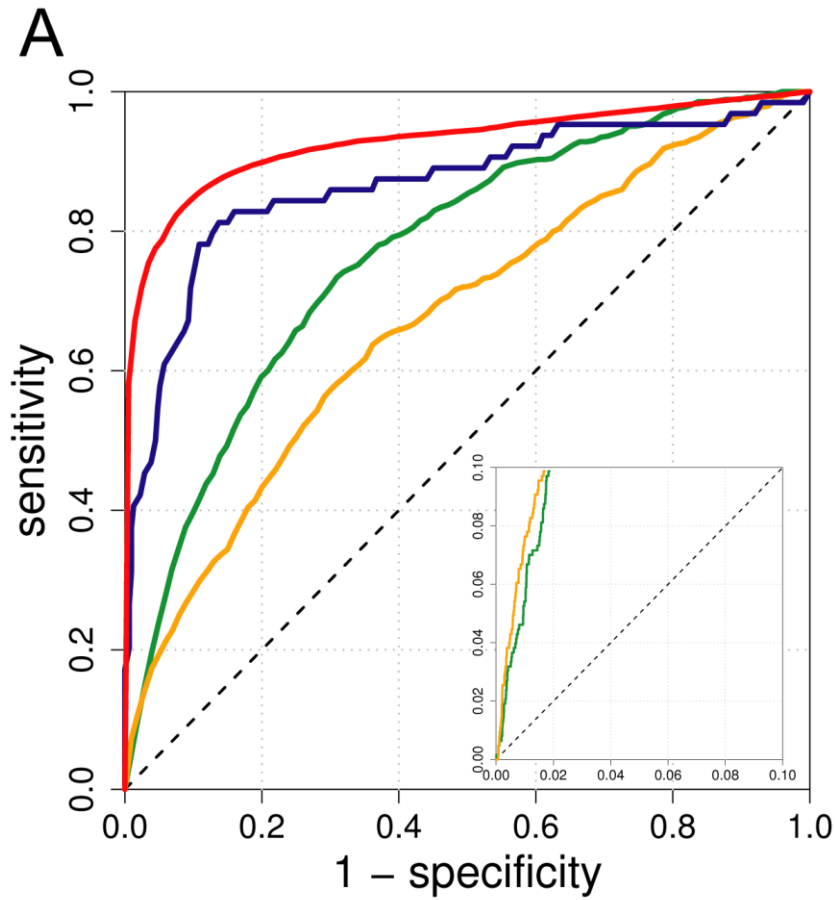
-Combination of all these features in a Machine Learning framework

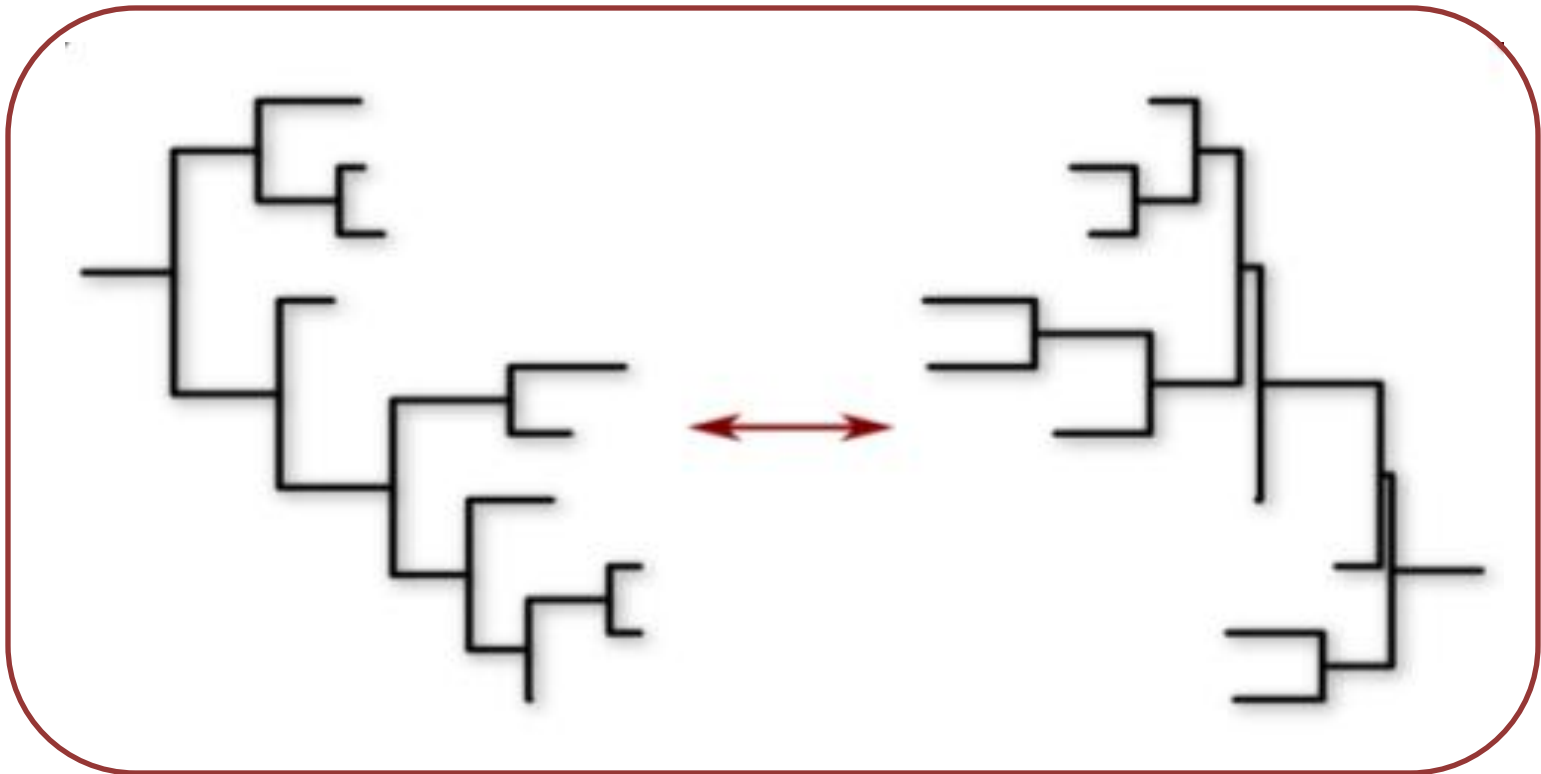
-Look at the capacity to correctly sort positive and negative pairs



25 features based on quality measures from Phylogenetic Profiles (see Table 2)

Combining multiple classifiers allows improving PPI prediction





Species tree ↔ Species tree

Gene tree ↔ Gene tree

Species tree ↔ Gene tree

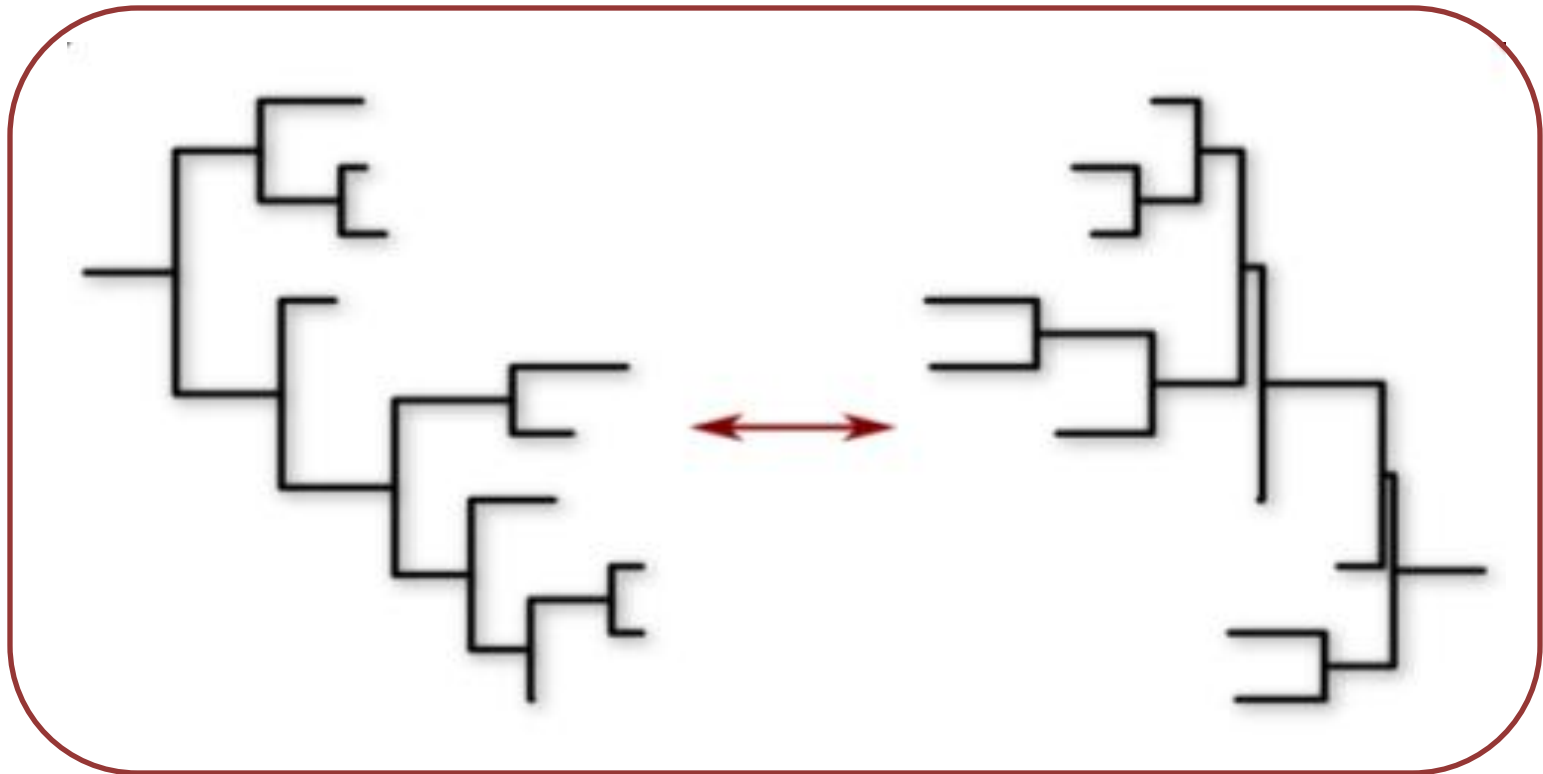
Multiple gene trees

Host parasite coevolutionary studies

Protein-Protein interaction detection

Reconciliation analyses

Phylogenomic studies



Species tree \leftrightarrow Species tree

Gene tree \leftrightarrow Gene tree

Species tree \leftrightarrow Gene tree

Multiple gene trees

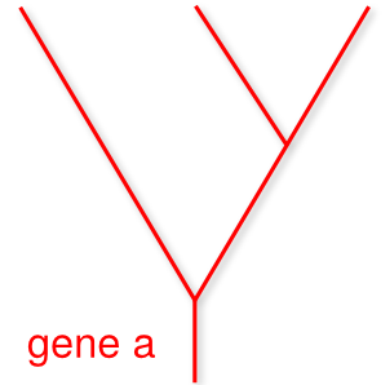
Host parasite coevolutionary studies

Protein-Protein interaction detection

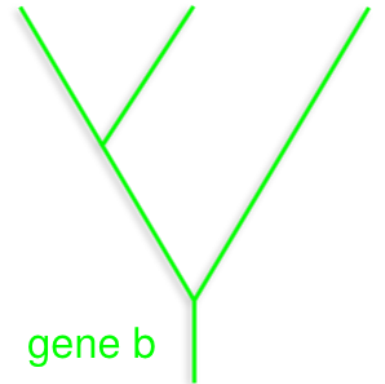
Reconciliation analyses

Phylogenomic studies

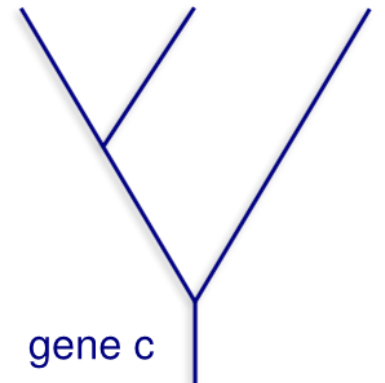
Species A Species B Species C



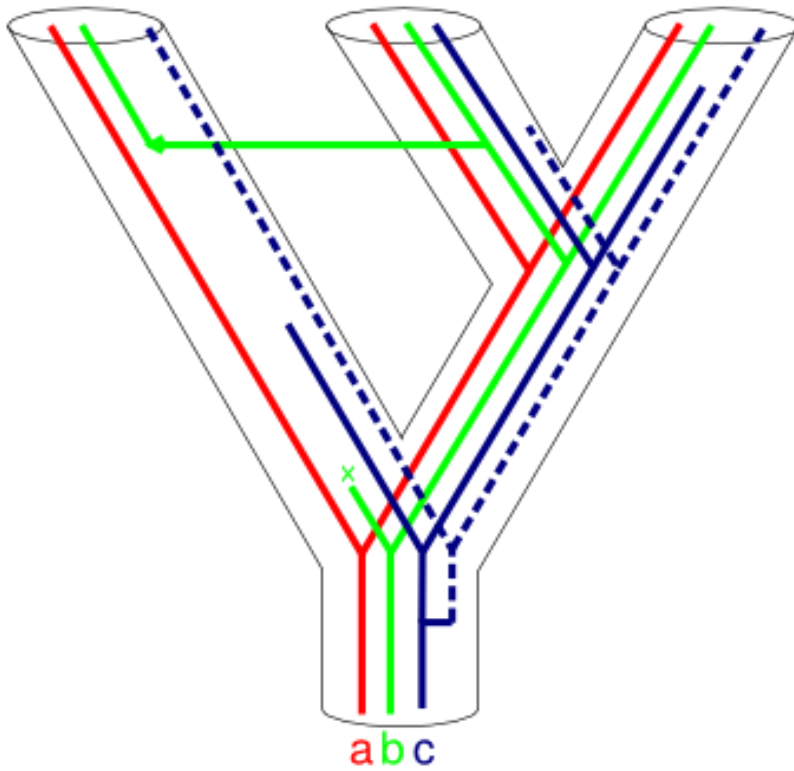
Species A Species B Species C



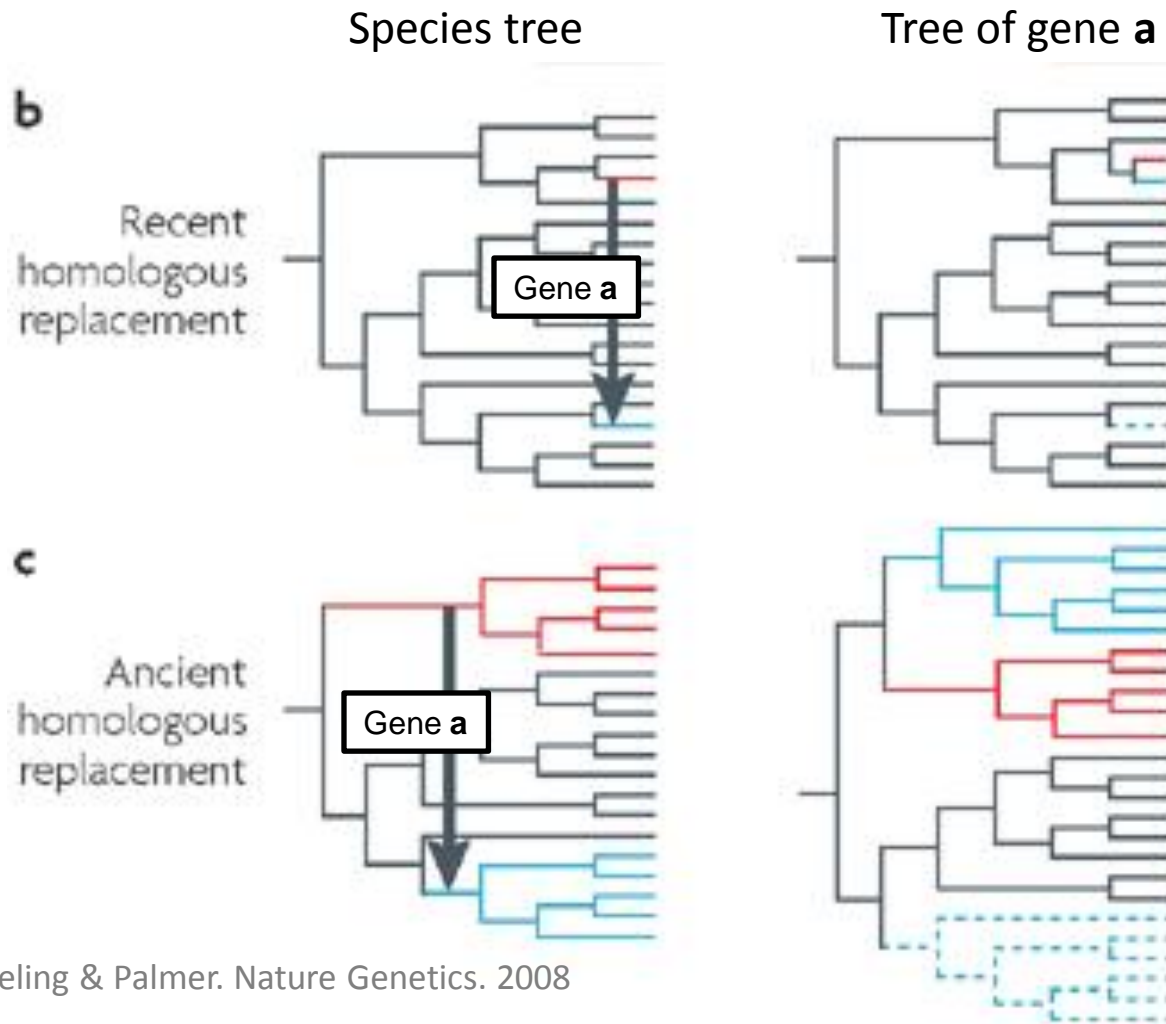
Species A Species C Species B



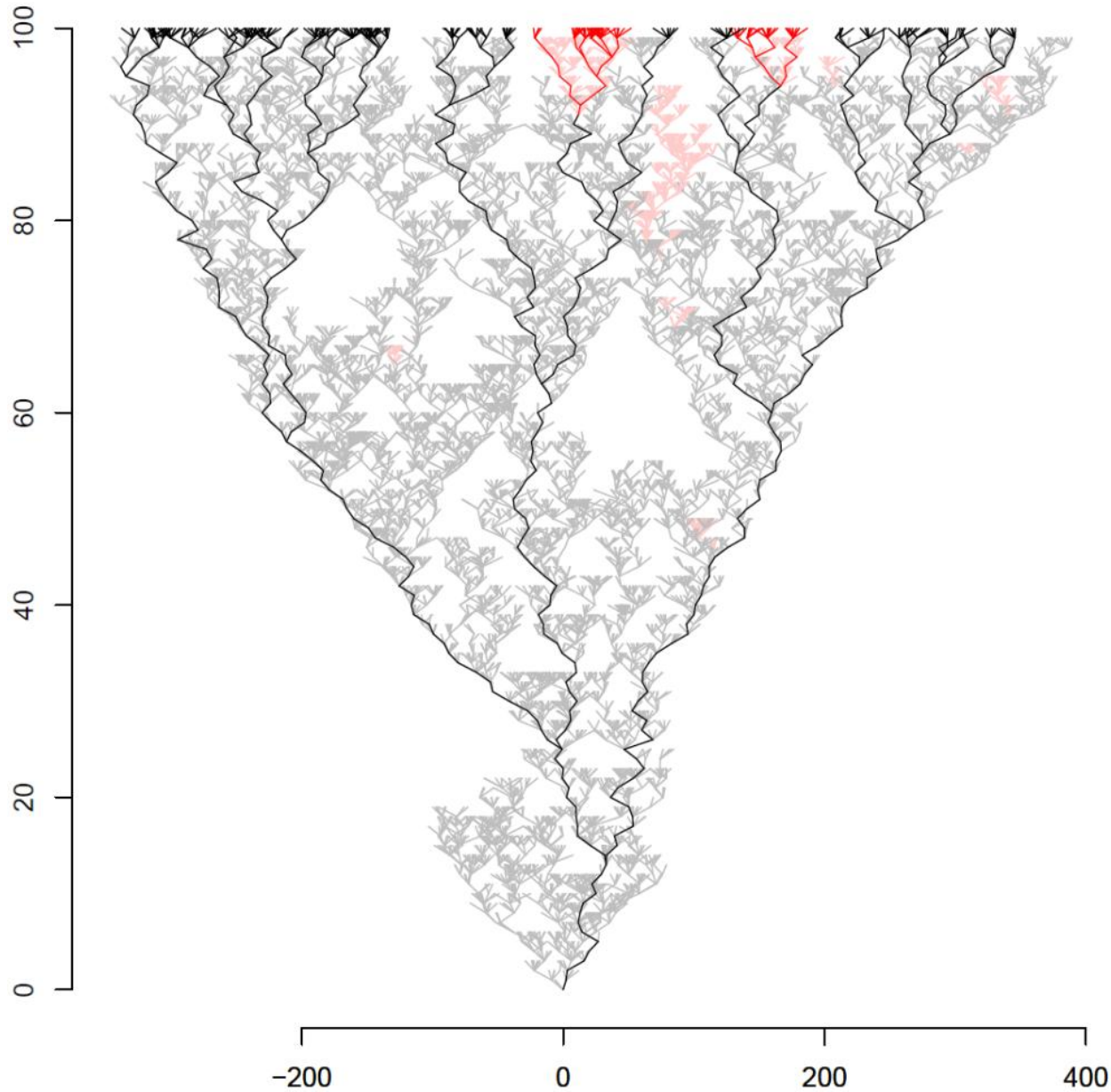
Species A Species B Species C



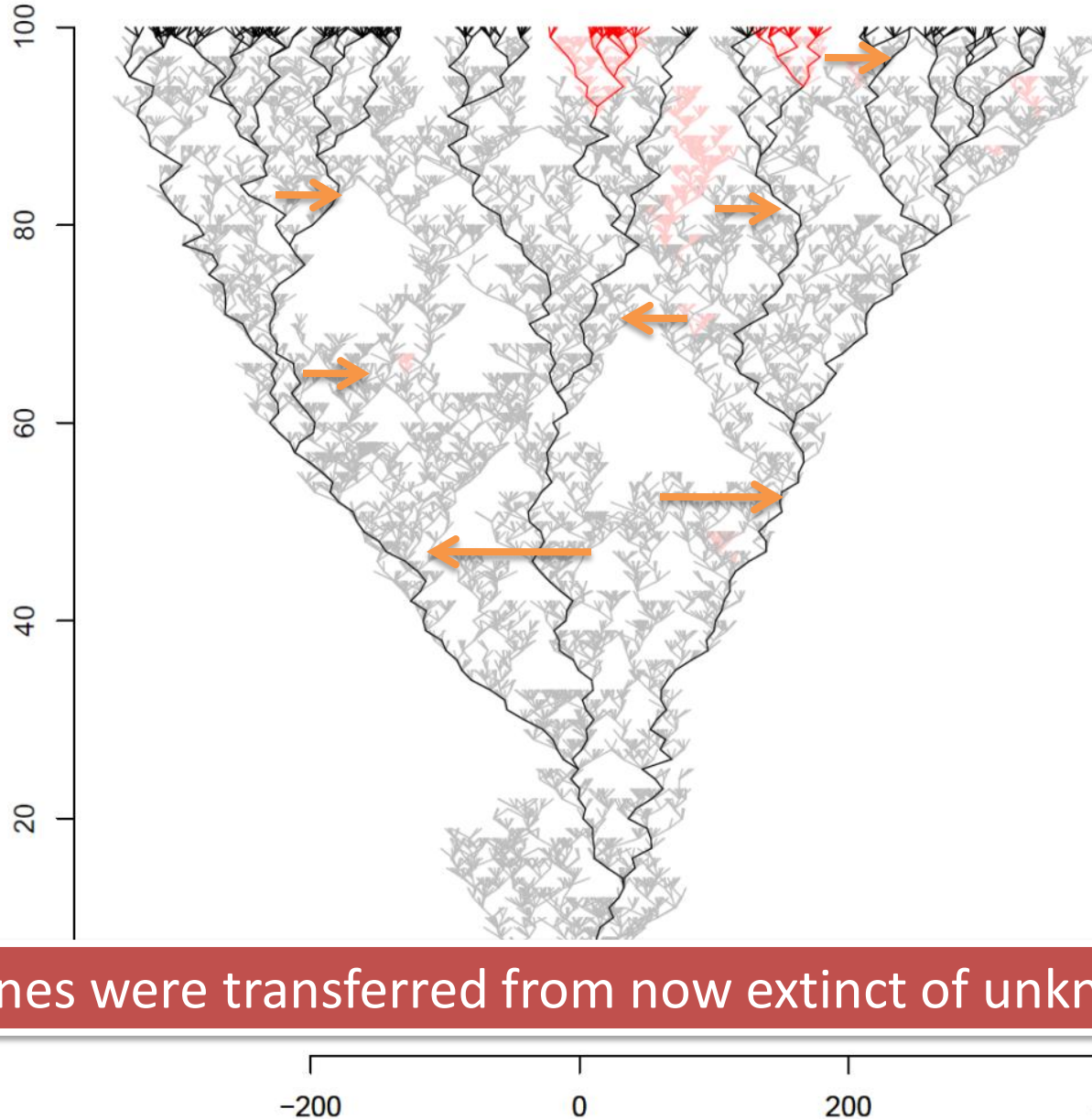
Consequence of horizontal gene transfer on gene tree topology



“To a first approximation, all species are extinct”

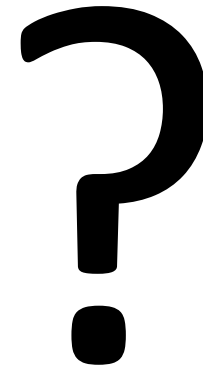
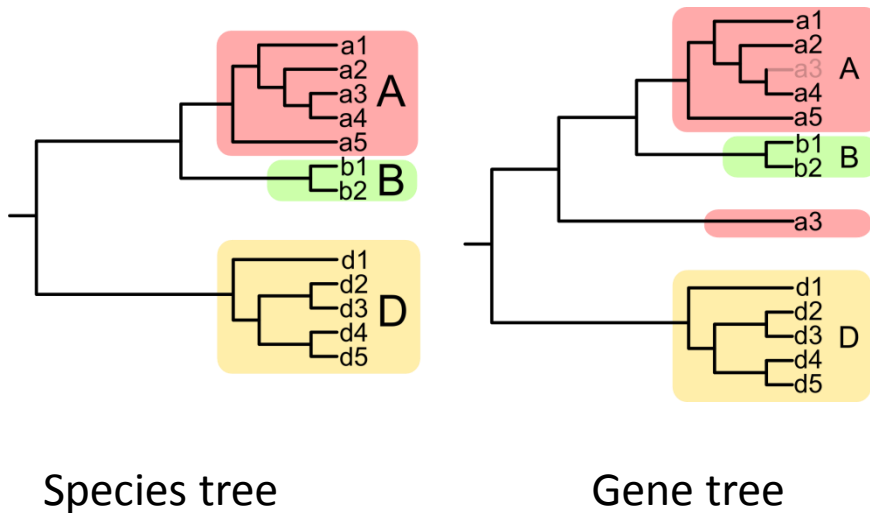


“To a first approximation, all species are extinct”

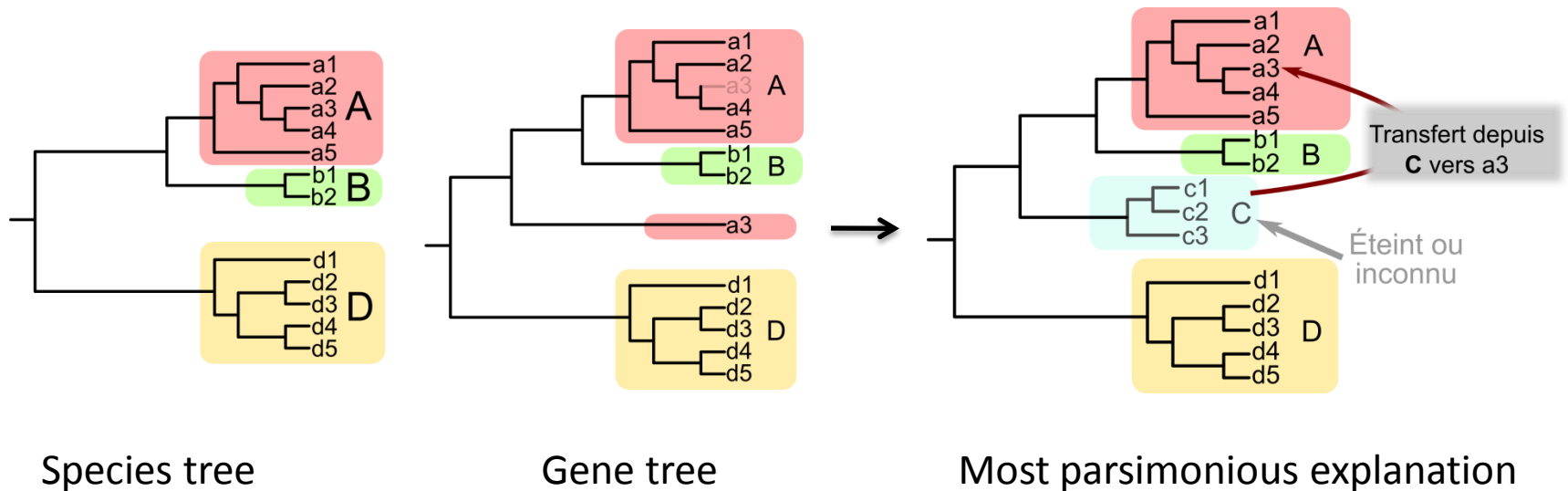


MANY genes were transferred from now extinct of unknown species

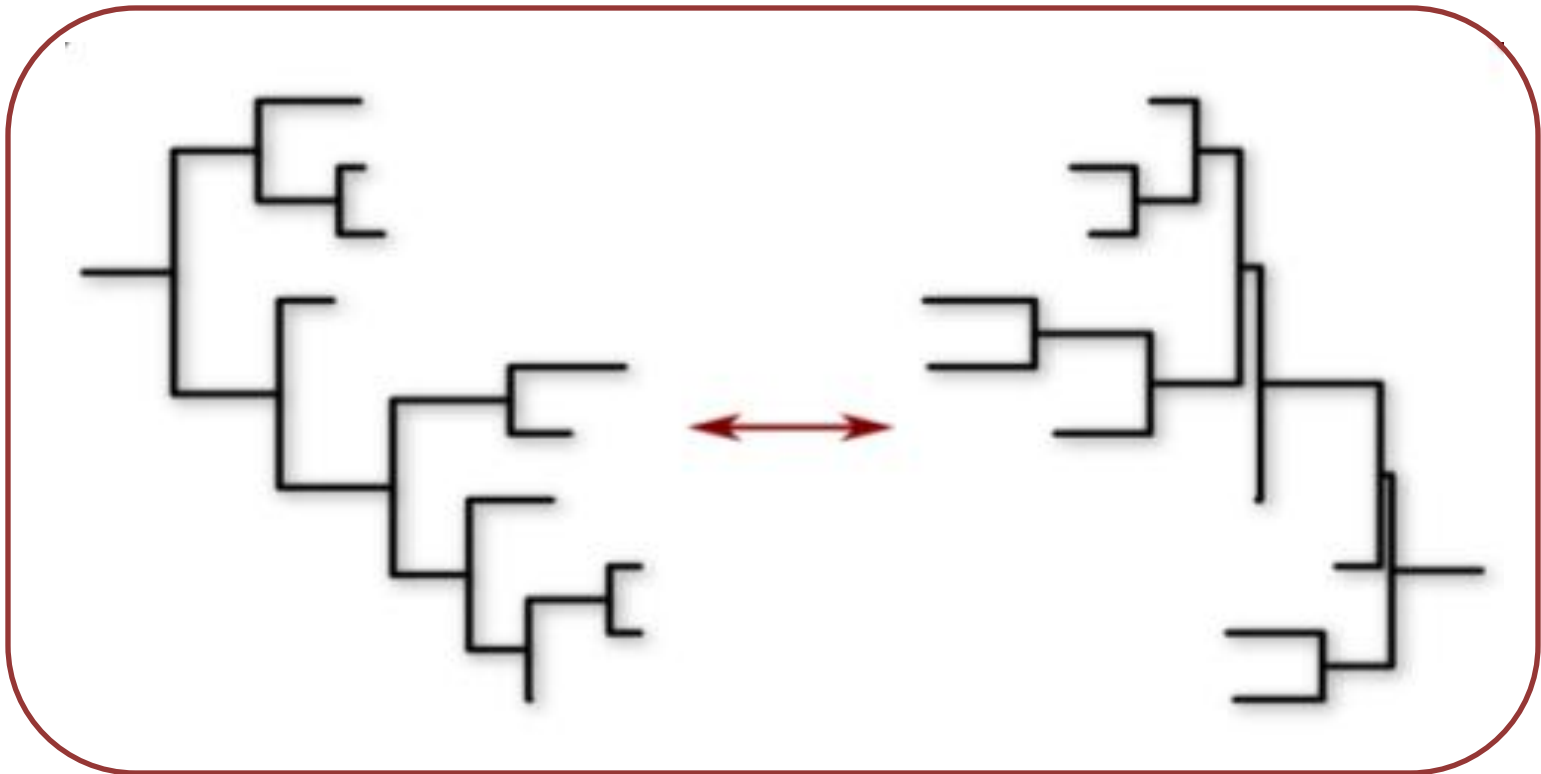
CAN we use HGT detection to explore extinct or unknown diversity?



CAN we use HGT detection to explore extinct or unknown diversity?



- Simulate species evolution with speciation, extinction, duplication
- Estimate sets of parameters allowing detection of extinct/unknown diversity
- Test on real bacterial dataset



Species tree ↔ Species tree

Gene tree ↔ Gene tree

Species tree ↔ Gene tree

Multiple gene trees

Host parasite coevolutionary studies

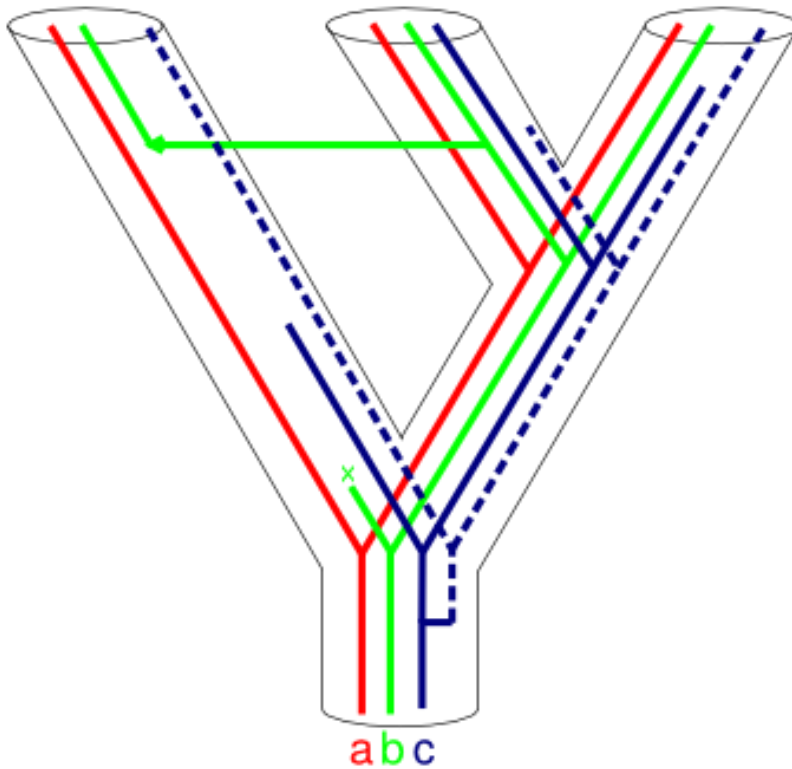
Protein-Protein interaction detection

Reconciliation analyses

Phylogenomic studies

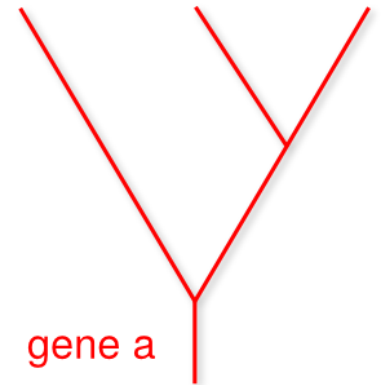
Incongruences between gene trees are common in phylogenomics

Species A Species B Species C

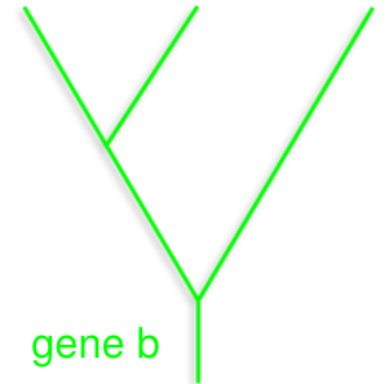


Incongruences

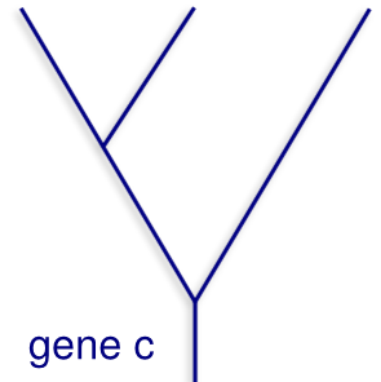
Species A Species B Species C



Species A Species B Species C



Species A Species C Species B



Existing methods for dealing with this variability

- Supermatrices
- Supertree
- Bayesian sampling
- Agreement subtrees
- Networks

Existing methods for dealing with this variability

- Supermatrices
- Supertree
- Bayesian sampling
- Agreement subtrees
- Networks

However, by concatenating the multilocus data, or by summarizing or obtaining a consensus of their individual gene trees, the latter methods lose a wealth of potentially interesting information, especially by removing outlier data or trees.

Phylo-MCOA

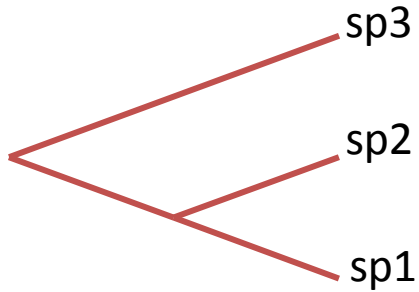


The goal

- Quickly compare a large number of gene trees
- Visualize the overall evolutionary history of the group analyzed.
- Find gene trees that tell the same story
- Find gene trees that tell different stories (produce discordant topologies)
 - Identify the species responsible for the discordance
 - Identify candidates for interesting biological processes

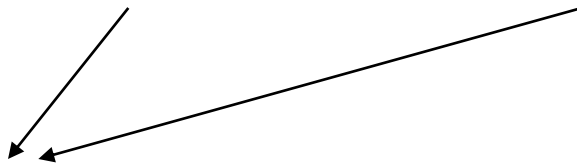
Principle

Gene tree G1



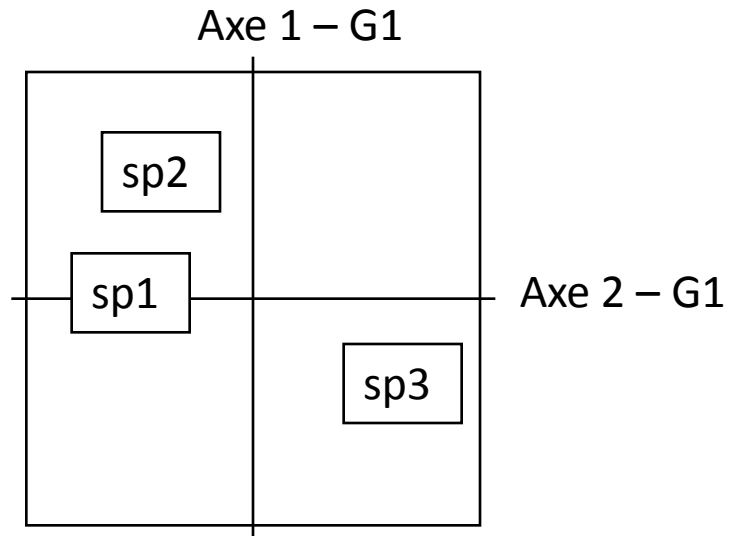
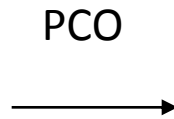
Alignment for gene G1

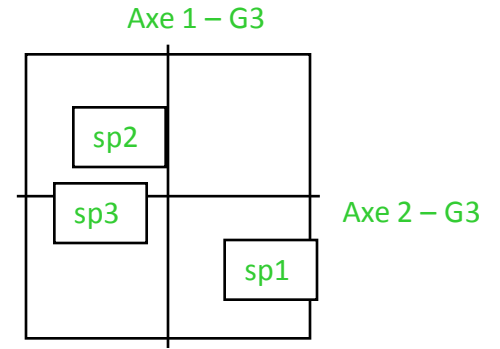
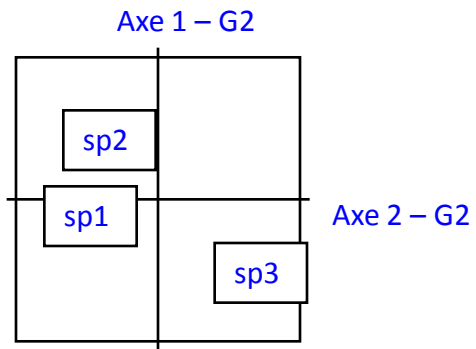
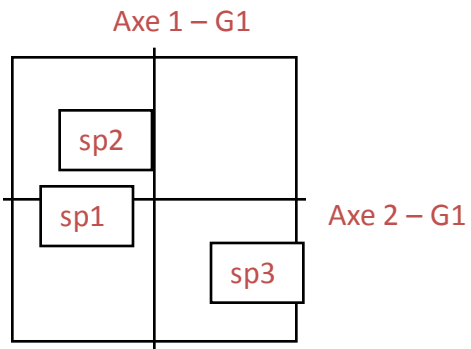
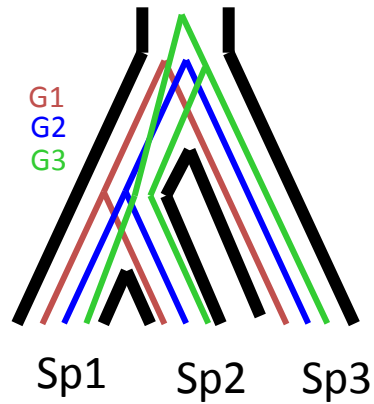
Sp1	AATGTGCATAC
Sp2	AATCTGCATAC
Sp3	AATGTCATTAC



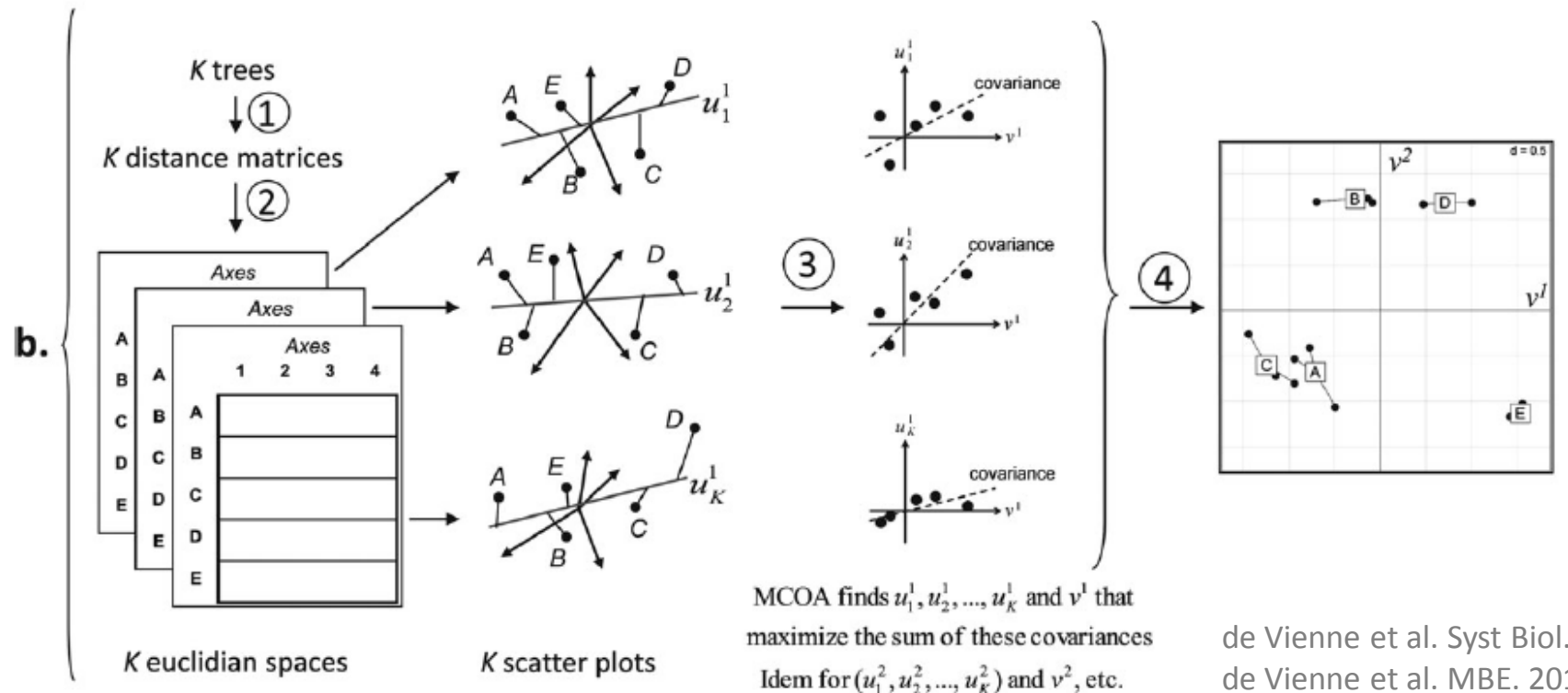
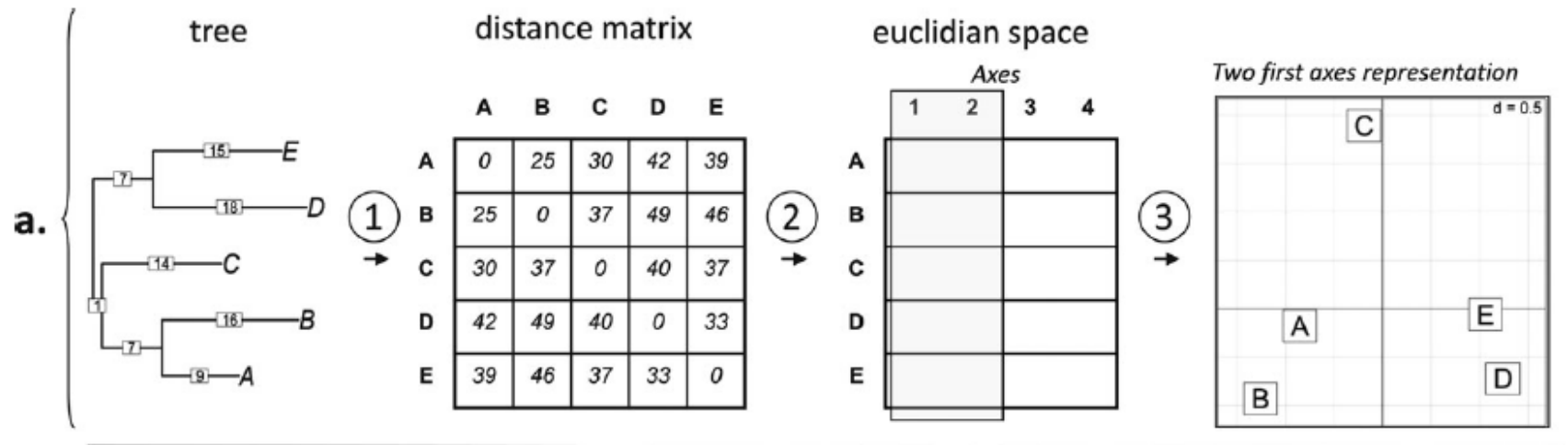
	Sp1	Sp2	Sp3
Sp1	0	x	y
Sp2	x	0	z
Sp3	y	z	0

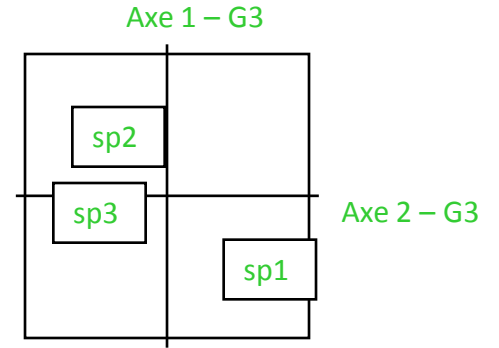
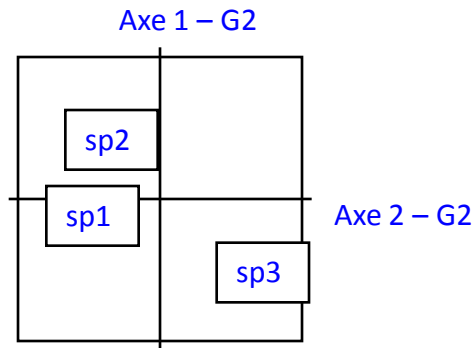
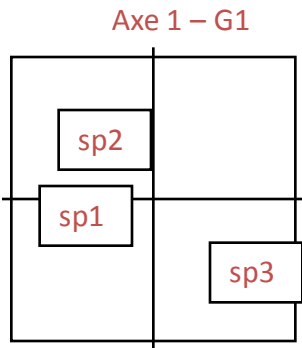
Pairwise distance matrix



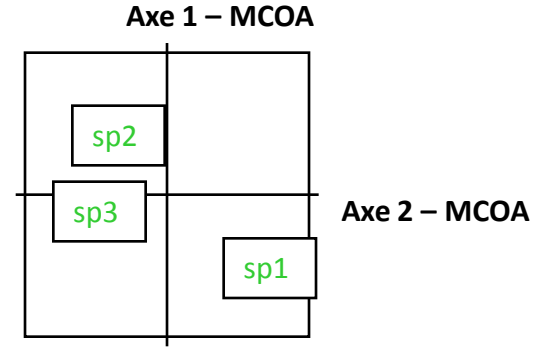
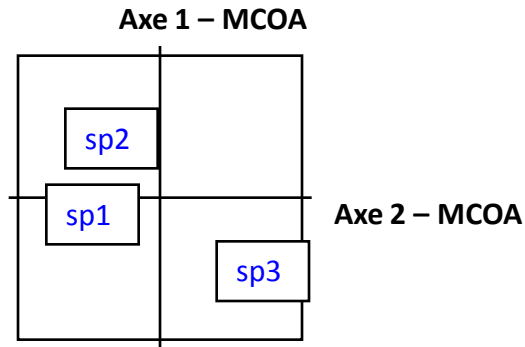
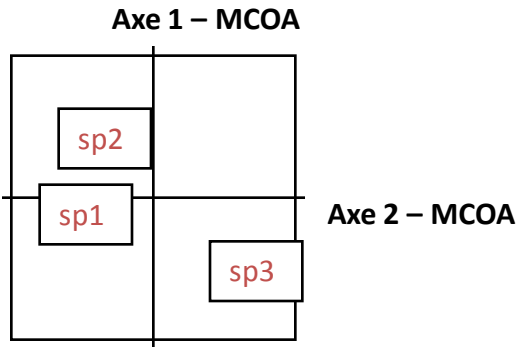


We want to visualize if the three genes tell the same story.
 So we want to compare the individual PCOs.
 So we need common axes. That's what MCOA does.

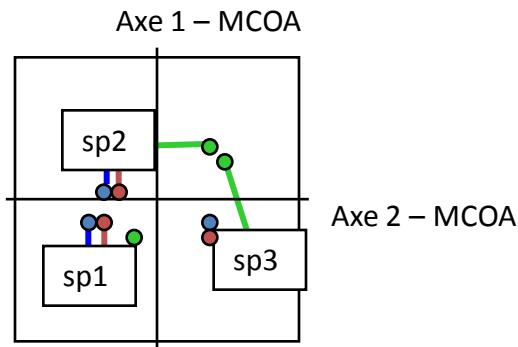




MCOA coordinates the individual axes

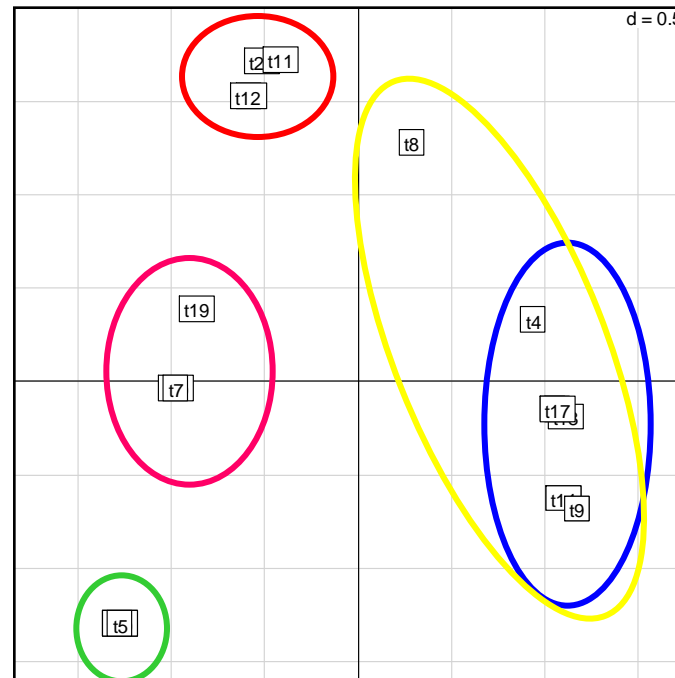
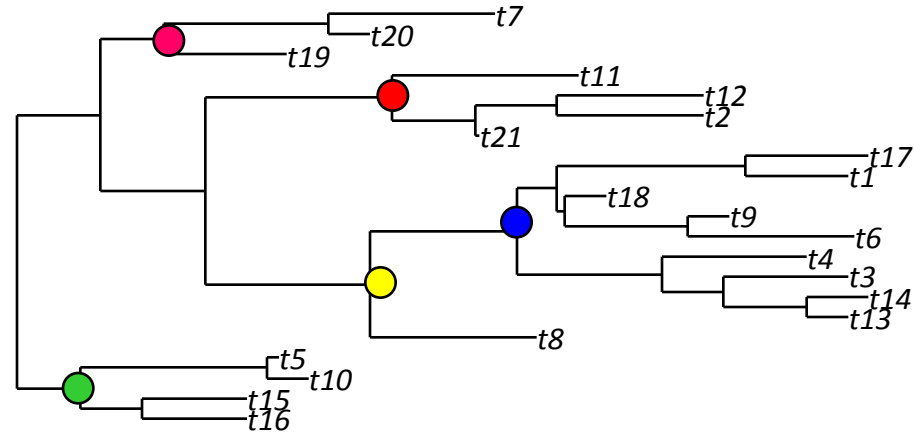


We can produce a single plot to visualize the concordance

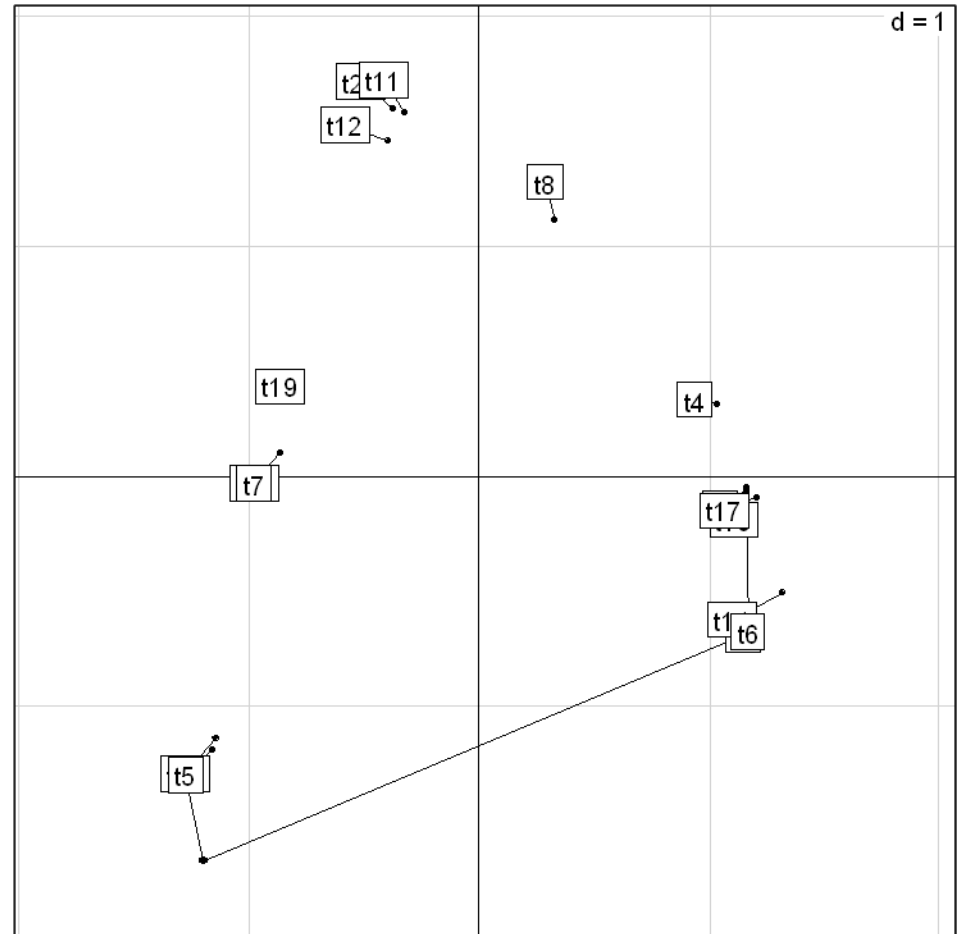
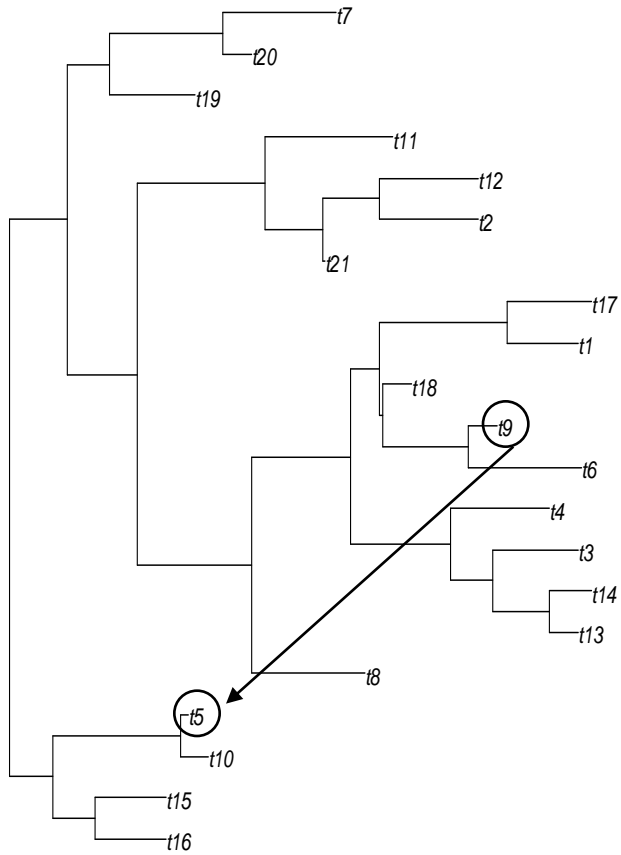


Cohesion plot

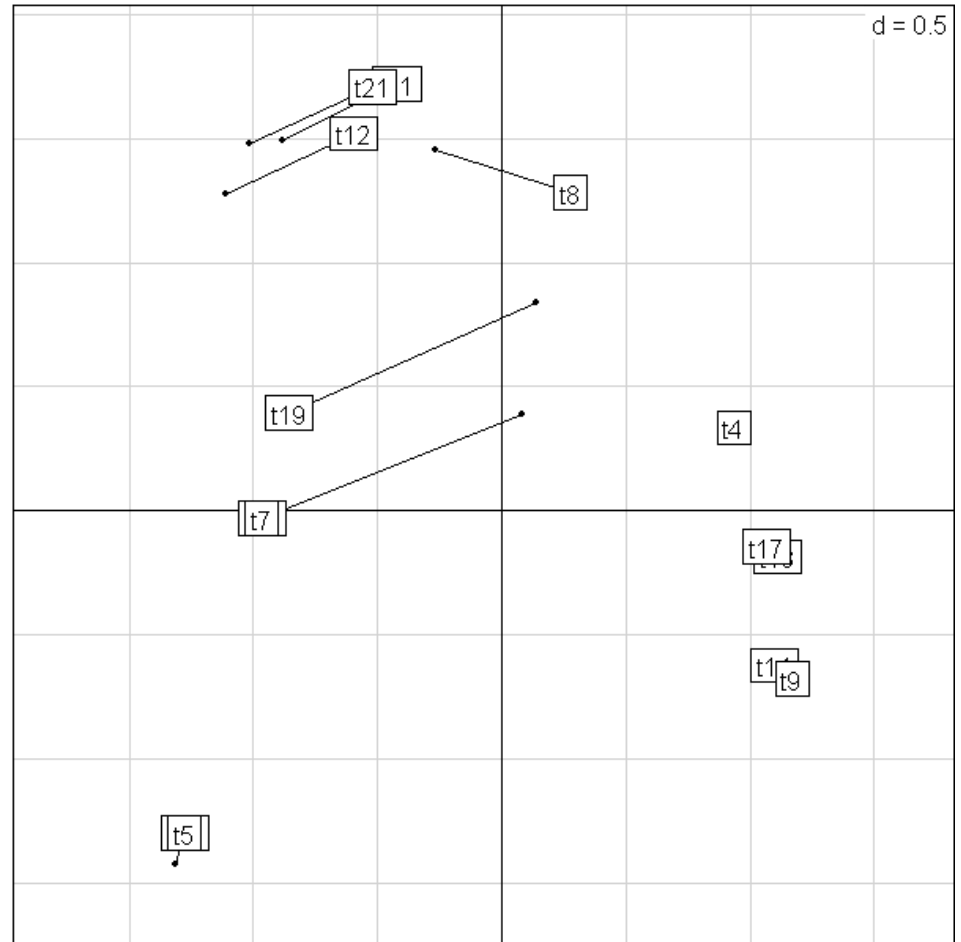
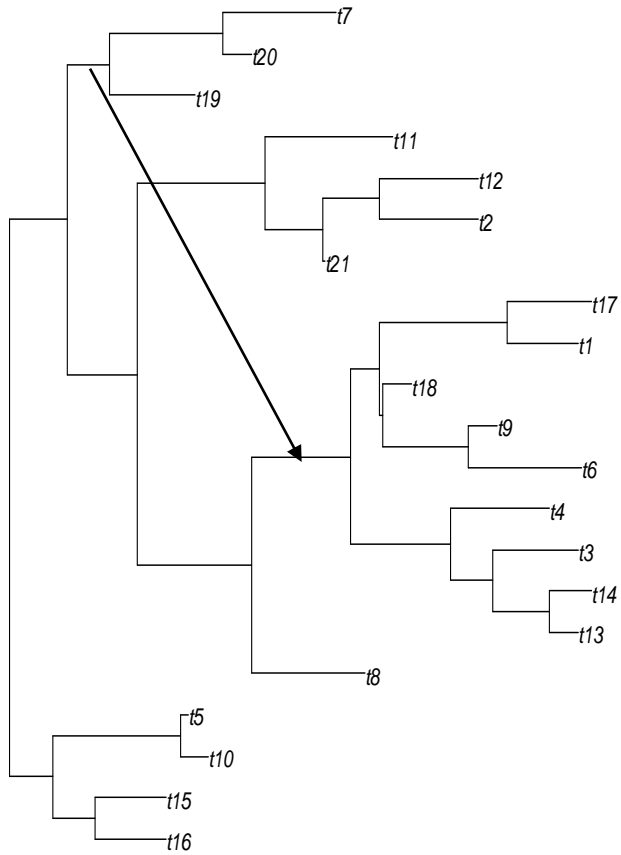
All the genes tell the same story



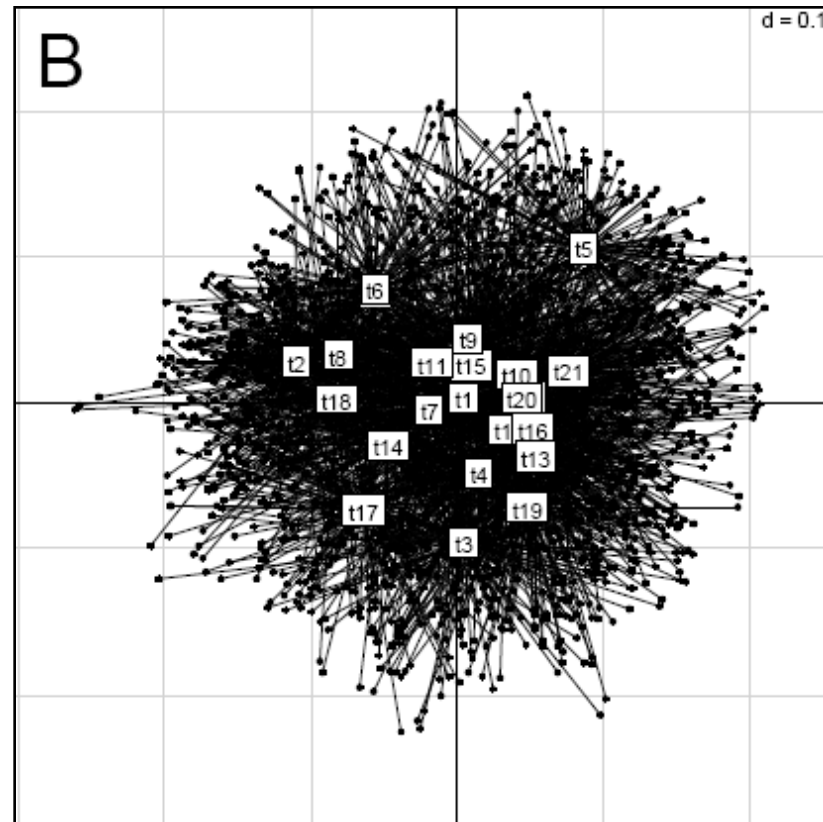
All the genes tell the same story + 1 recent HGT



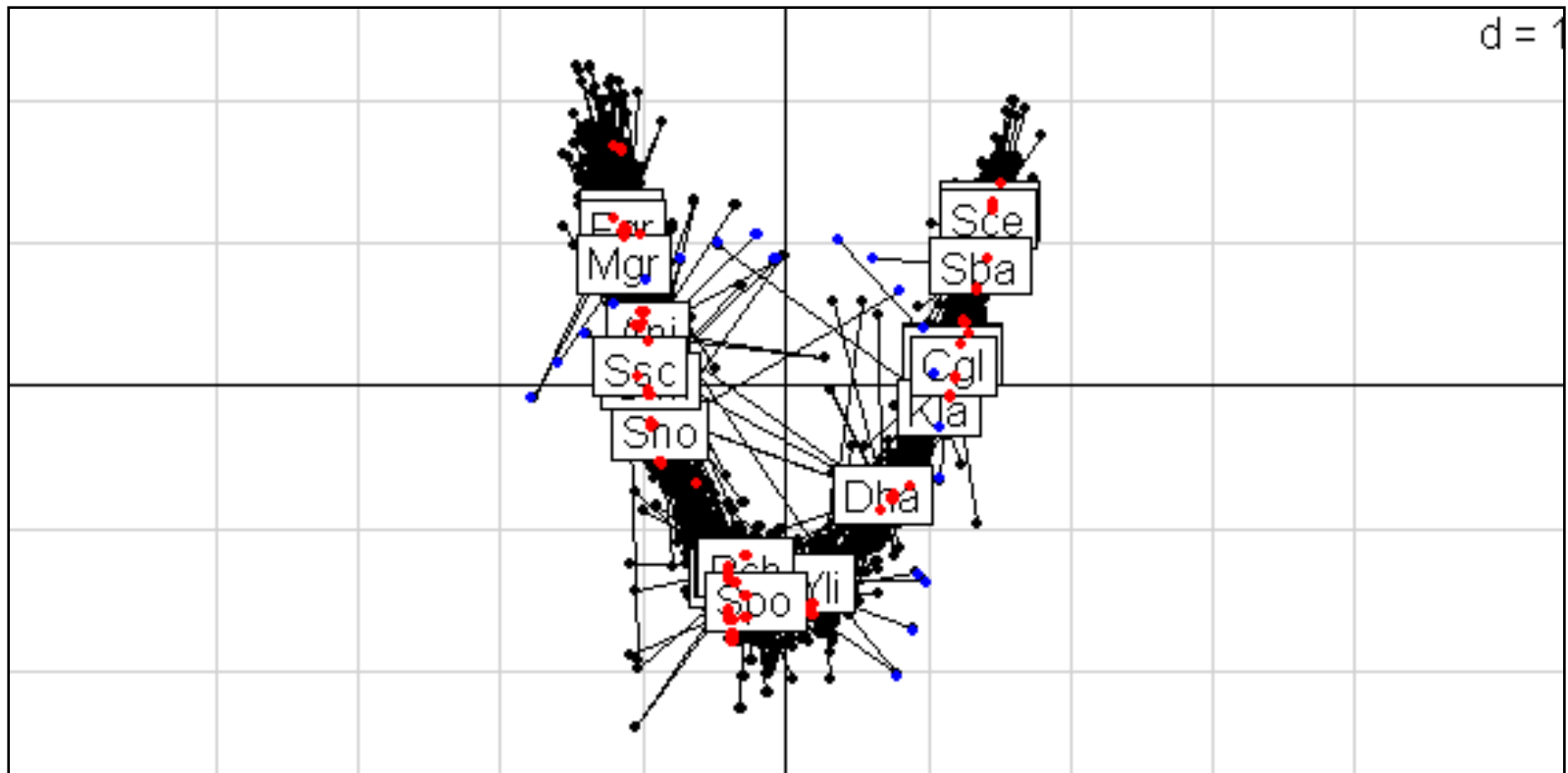
All the genes tell the same story + 1 ancient HGT



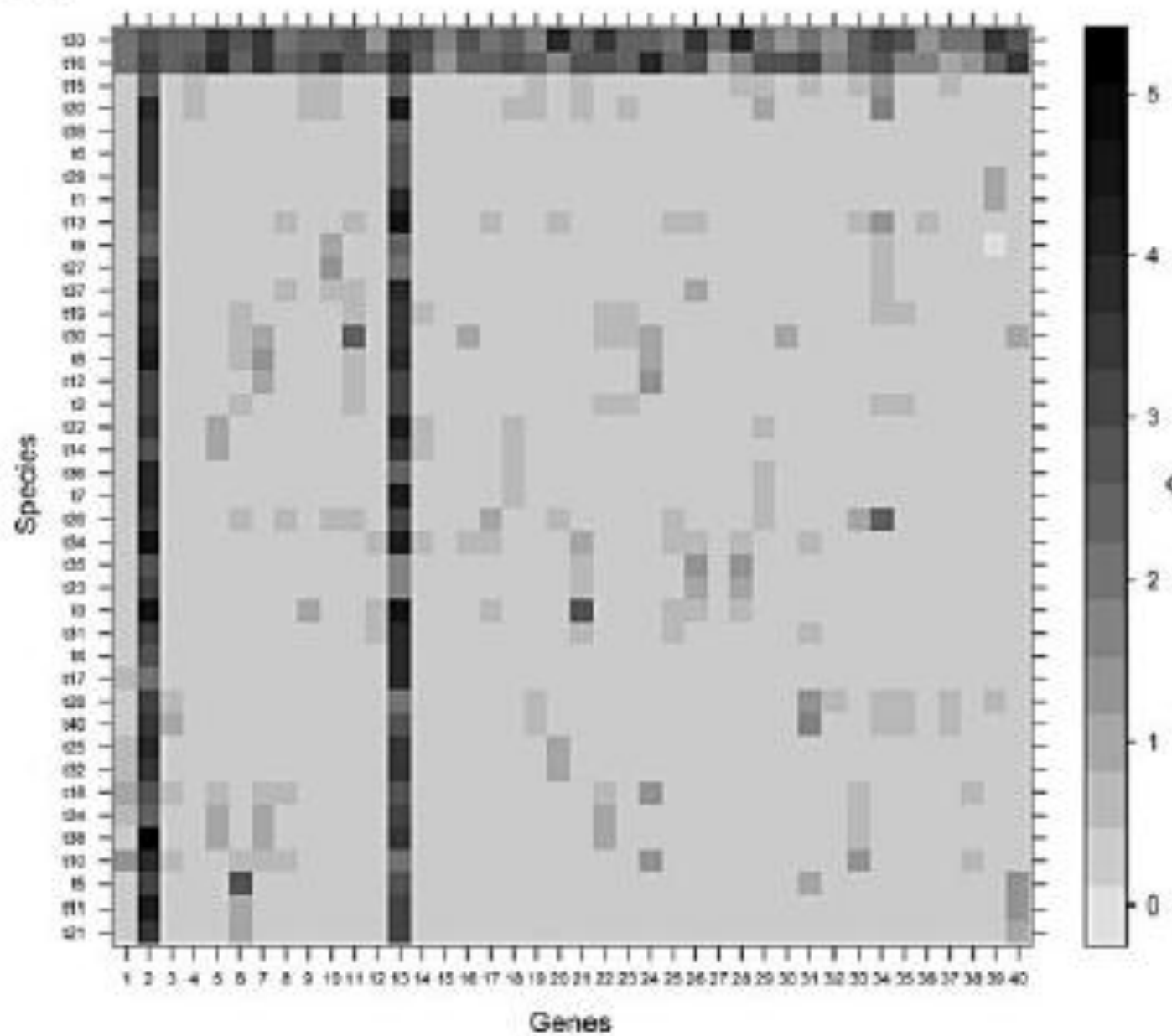
All the genes tell different stories



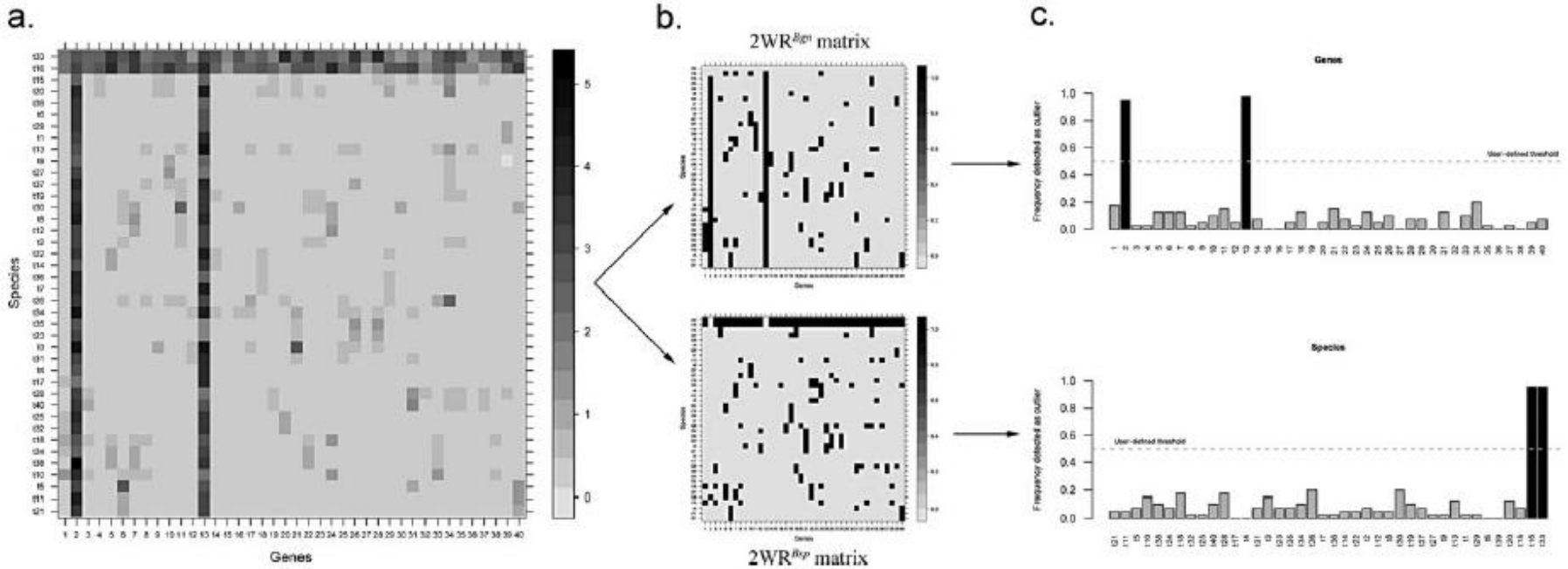
Real example (21 fungal species)



a.

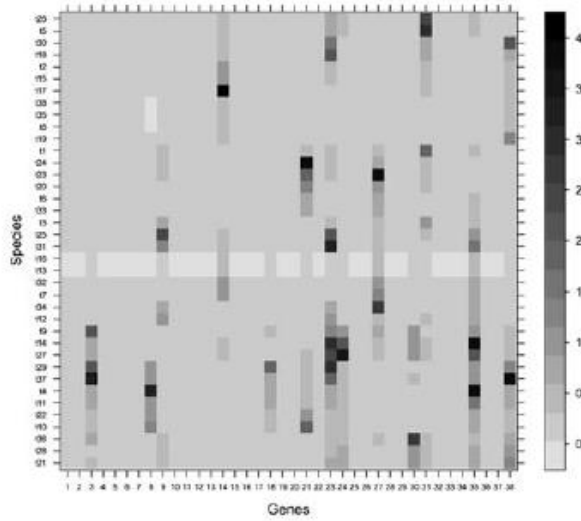


From cohesion plot to 2-way reference matrix



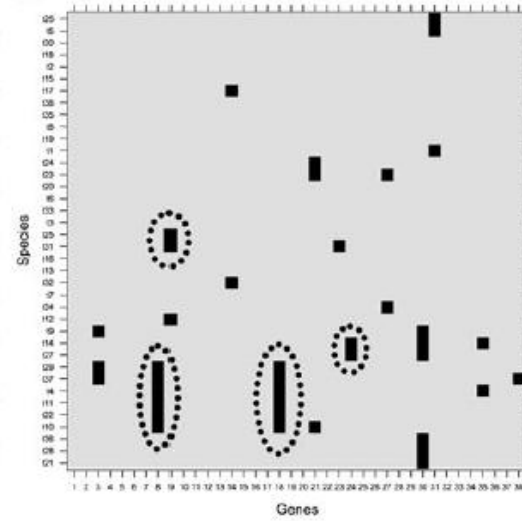
→ Identify outliers

a.



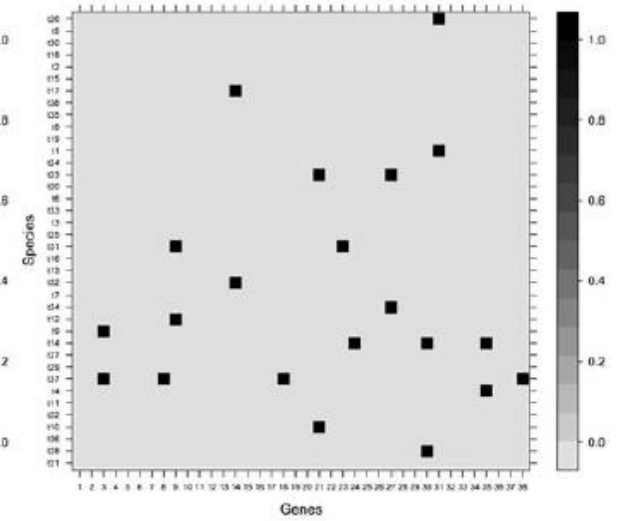
2WR matrix

b.



2WR^{Bgnsp} matrix
with "islands"

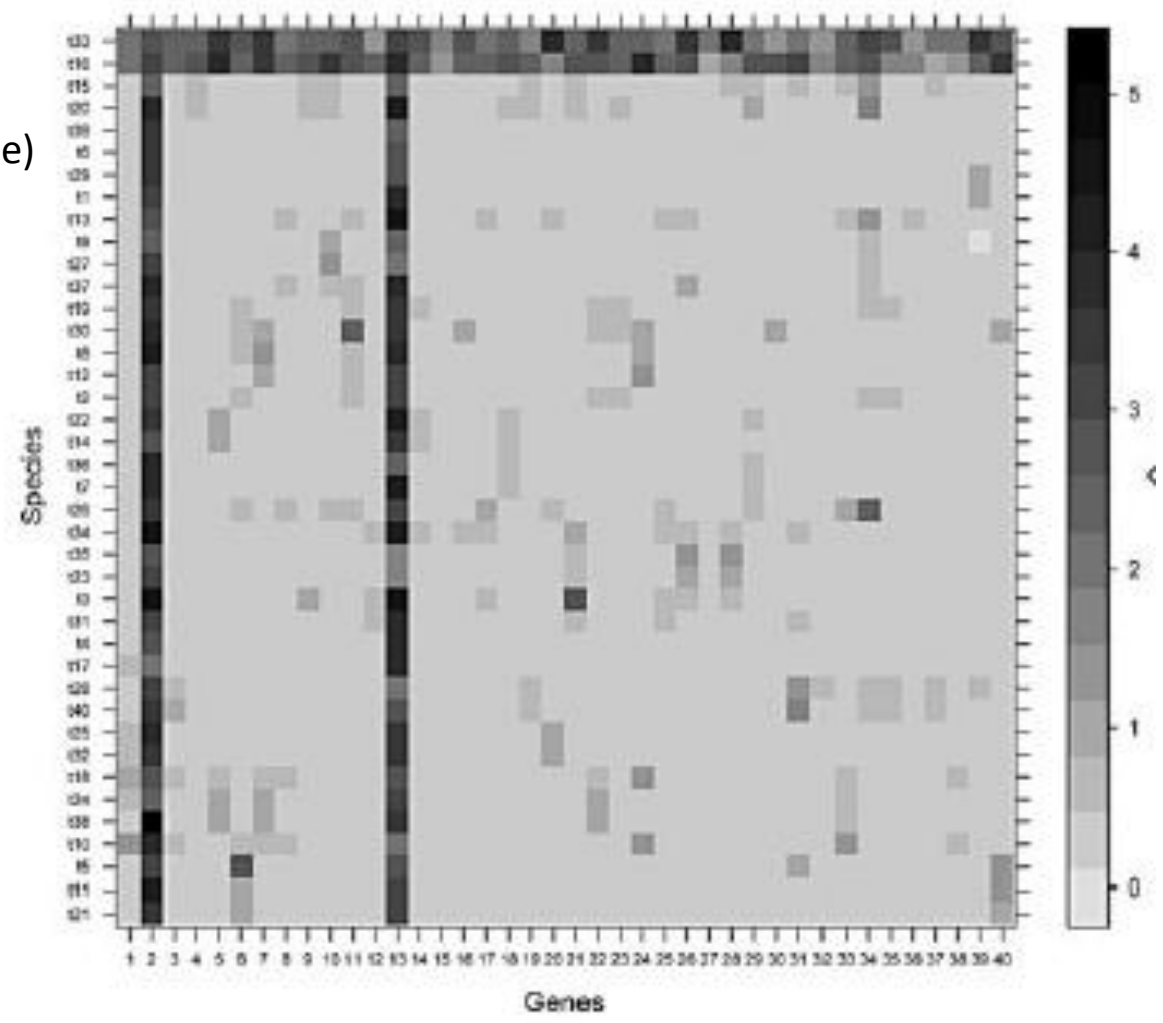
c.



2WR^{Bgnsp} matrix
without "islands"

→ Identify outliers

- Any problem (species identification...)
- “Volatile” species
(acceptors of genes by HGT, for example)



What it means

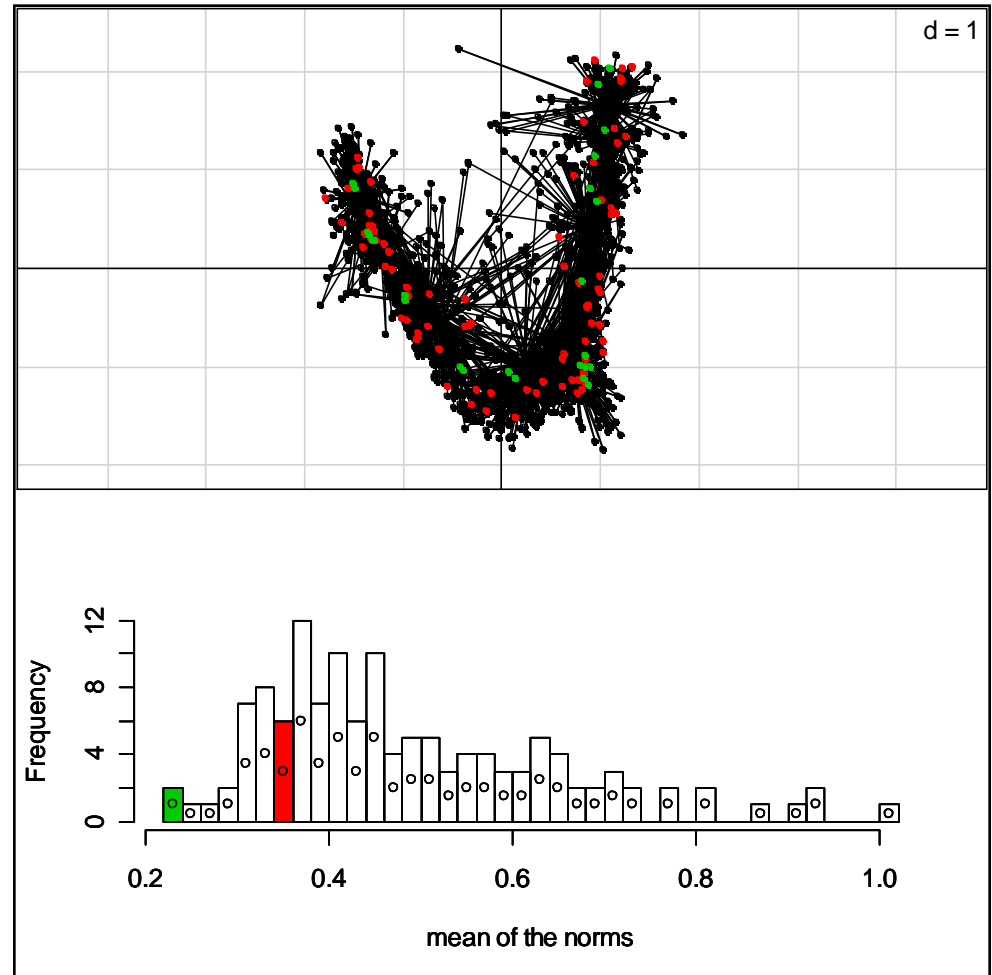
- Any problem (sequencing, paralogy, tree reconstruction...)
- “highly transferable genes /shared genes

Thanks

- HP interaction
 - Tatiana Giraud
 - Guislaine Refrégier
 - Michael Hood
 - ...
- Icong
 - Olivier Martin
- PPI
 - Jérôme Azé
- Phylo-MCOA
 - Gabriela Aguilera
 - Sebastien Ollier

Our method, Phylo-MCOA

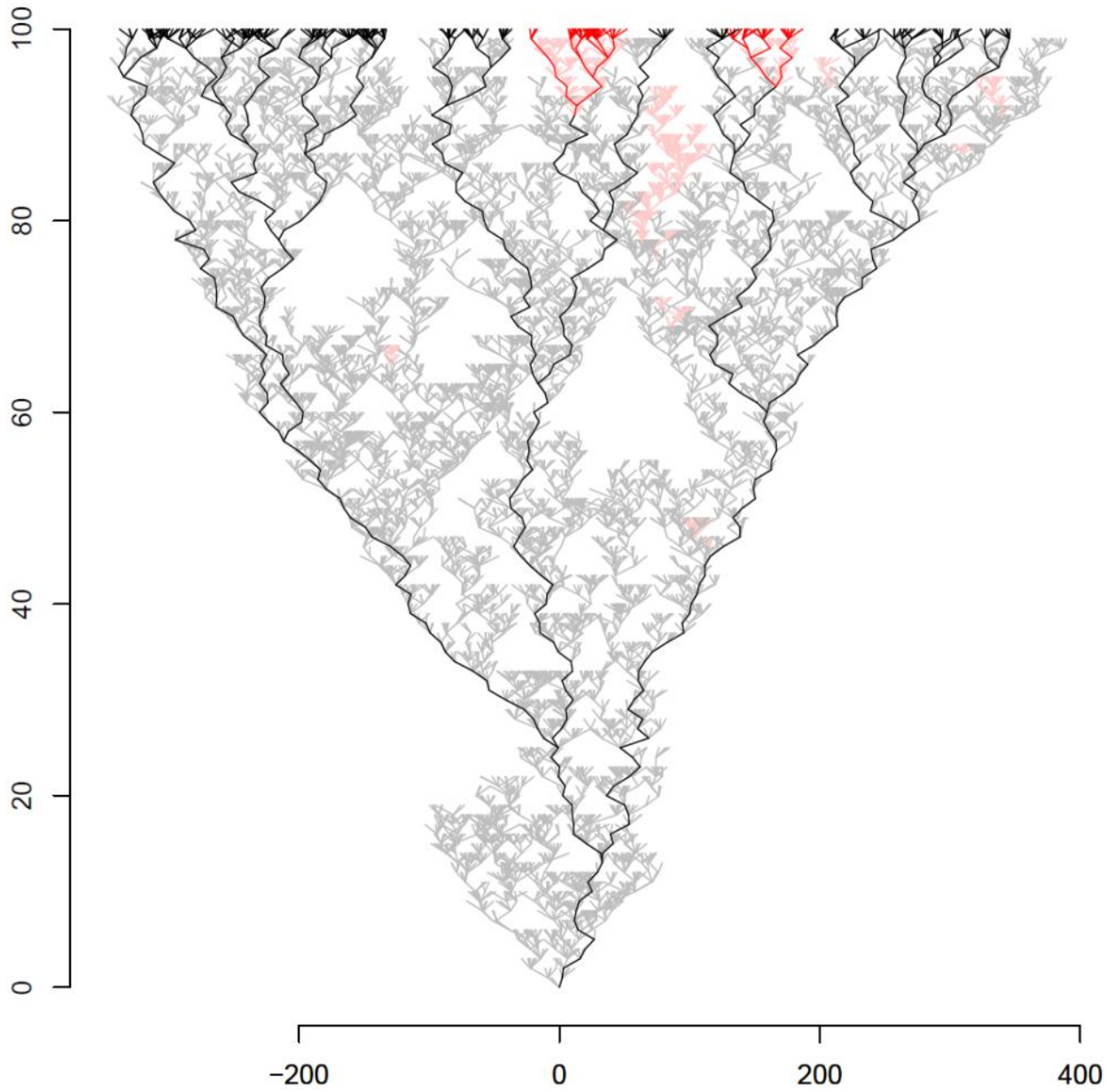
- Is fast
- Is user friendly
- Is concordant with other methods
- Is colorful
- Is free

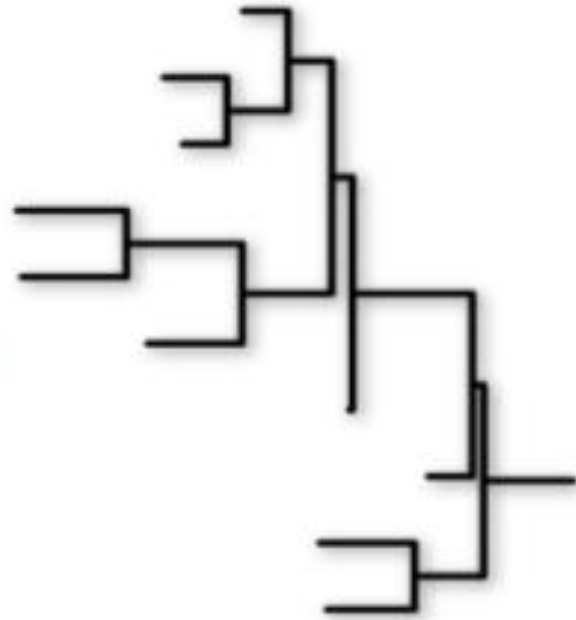
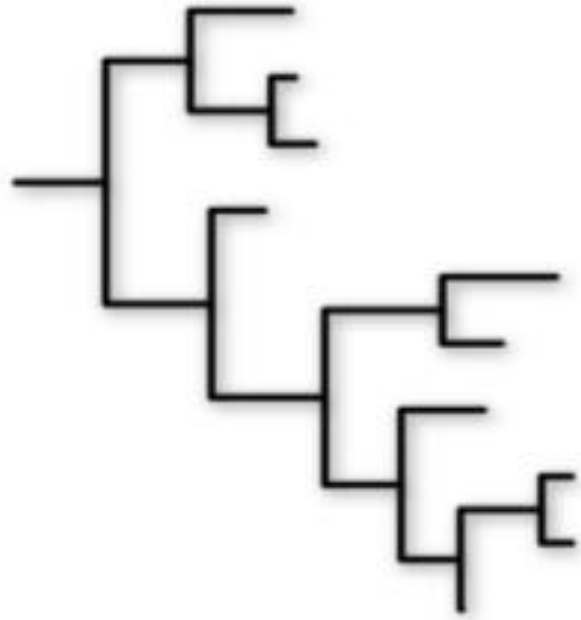


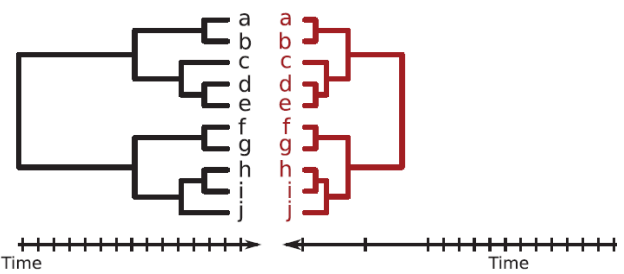
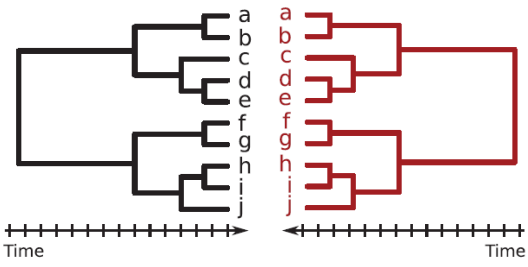
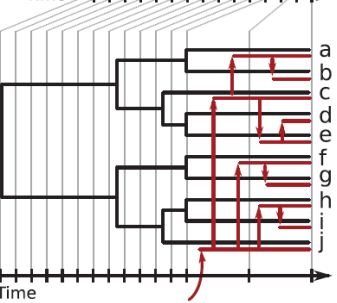
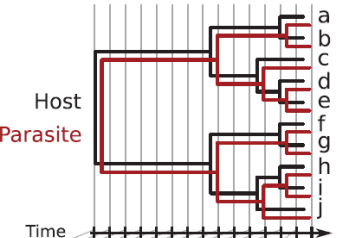
Examples with simulated data

- All genes tell different evolutionary histories
- All genes tell the same evolutionary history
- Recent HGT
- Ancient HGT

- Tree comparisons and coevolution
 - Between host and parasites (JEB)
 - Between genes/proteins (PloS ONE)
- Tree comparisons and species tree/gene tree “reconciliation”
 - Multiple gene trees in comparative genomics







I_{cong}

ON-LINE CALCULATION OF THE CONGRUENCY INDEX I_{cong}
TO ASSESS THE EXISTENCE OF TOPOLOGICAL CONGRUENCE
BETWEEN TWO TREES

NEW version, goes faster !!



If you need help for the formatting of your trees or if you want to see an example, [click here](#).

If you use this index for a publication, please cite
de Vienne D.M., Giraud, T. and Marlin, O.C. 2007. A Congruence Index for Testing
Topological Similarity between Trees. *Bioinformatics* **23** (23): 3119-3124.

Type or paste the first tree here:

Type or paste the second tree here:

Calculate I_{cong} and its P-value | Reset