

Gene tree-species tree methods in RevBayes

Bastien Boussau

LBBE, CNRS, Université de Lyon



What I spend my time doing

- Developing methods for:
 - molecular evolution
 - phylogenetic reconstruction
 - genomics, genome scale phylogenetic reconstruction
- Using these methods to investigate:
 - Genome evolution (gene family evolution, drift)
 - Evolution of life on Earth (species trees, ancestral trait reconstruction)

Why the tautological title?

“Phylogenomic” has been used to describe attempts at reconstructing species trees based on 10-100 genes

We are interested in the species tree, but also in genome evolution

A genome can contain >20,000 genes

Genome-scale: thousands of genes, dozens of species

Genomes, Genes, gene families, gene trees and species trees

Genome



Genomes, Genes, gene families, gene trees and species trees

Genome
Gene



Homo sapiens; GeneA: ACTGGTGATGACATGAC...

Genomes, Genes, gene families, gene trees and species trees

Genome
Gene



Gene family
alignment

Homo sapiens; GeneA: ACTGGTGATGACATGAC...

Homo sapiens; GeneA: ACTGGTGATGACATAAC...

Homo sapiens; GeneB: ACTGTTGATGACATGAC...

Mus musculus; GeneC: ACTGATGATGACAAGAC...

Mus musculus; GeneD: ACTGGTGA--CCATGAC...

Bison bison; GeneE: ACTGGTGATGACACGAC...

Canis lupus; GeneF: ACT--TCATGAAACGAC...

Genomes, Genes, gene families, gene trees and species trees

Genome
Gene



Gene family
alignment

Homo sapiens; GeneA: ACTGGTGATGACATGAC...

Homo sapiens; GeneA: ACTGGTGATGACATAAC...

Homo sapiens; GeneB: ACTGTTGATGACATGAC...

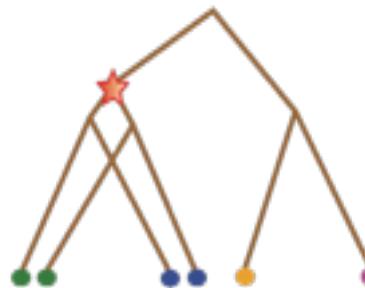
Mus musculus; GeneC: ACTGATGATGACAAGAC...

Mus musculus; GeneD: ACTGGTGA--CCATGAC...

Bison bison; GeneE: ACTGGTGATGACACCGAC...

Canis lupus; GeneF: ACT--TCATGAAACGAC...

Gene tree



Genomes, Genes, gene families, gene trees and species trees

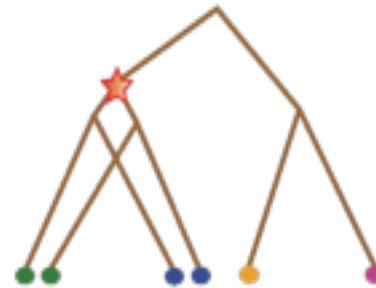
Genome
Gene



Gene family
alignment

Homo sapiens; GeneA: ACTGGTGATGACATGAC...
Homo sapiens; GeneB: ACTGTTGATGACATGAC...
Mus musculus; GeneC: ACTGATGATGACAAGAC...
Mus musculus; GeneD: ACTGGTGA--CCATGAC...
Bison bison; GeneE: ACTGGTGATGACACCGAC...
Canis lupus; GeneF: ACT--TCATGAAACGAC...

Gene tree



Species trees

A graphical model for phylogenomics

Species tree

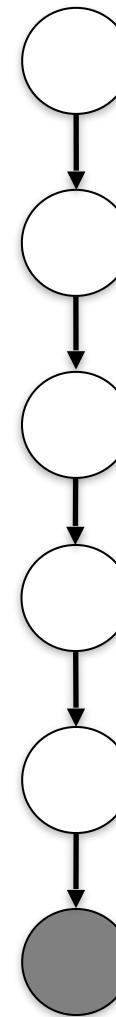
Gene tree

Gene family alignment

Gene

Genome sequences

Raw sequences



A graphical model for phylogenomics

Species tree

Species tree construction

Gene tree

Gene tree construction

Gene family alignment

*Homology prediction
Alignment*

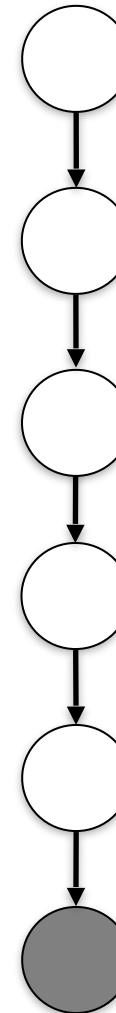
Gene

Gene prediction

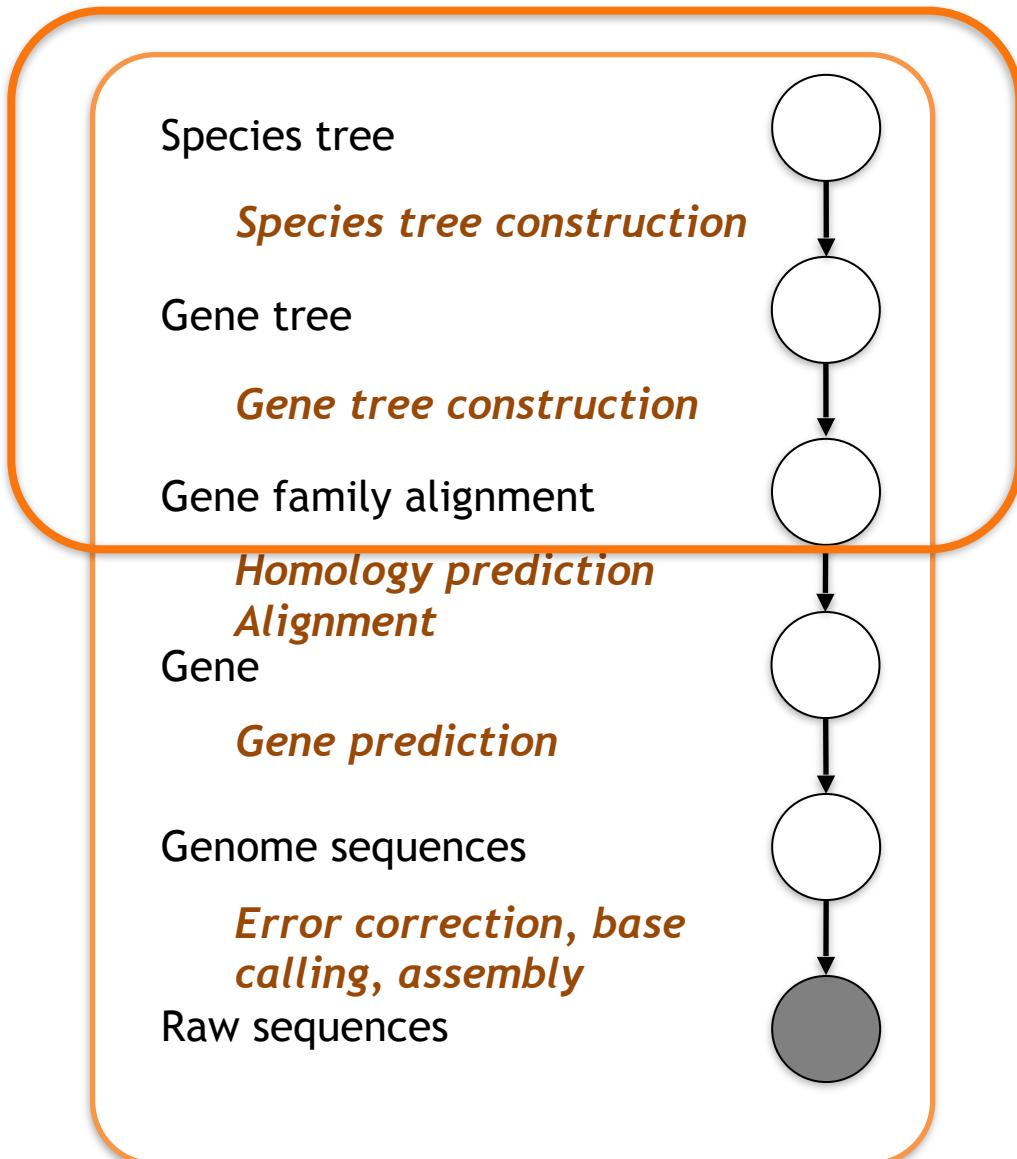
Genome sequences

*Error correction, base
calling, assembly*

Raw sequences



A graphical model for phylogenomics



The usual approach to reconstructing phylogenetic trees

Homo sapiens;GeneA: ACTGGTGATGACATAAC...

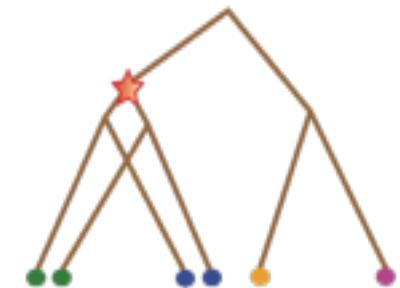
Homo sapiens;GeneB: ACTGTTGATGACATGAC...

Mus musculus;GeneC: ACTGATGATGACAAGAC...

Mus musculus;GeneD: ACTGGTGA--CCATGAC...

Bison bison; GeneE: ACTGGTGATGACACCGAC...

Canis lupus; GeneF: ACT--TCATGAAACGAC...



The usual approach to reconstructing phylogenetic trees

Homo sapiens;GeneA: ACTGGTGATGACATAAC...

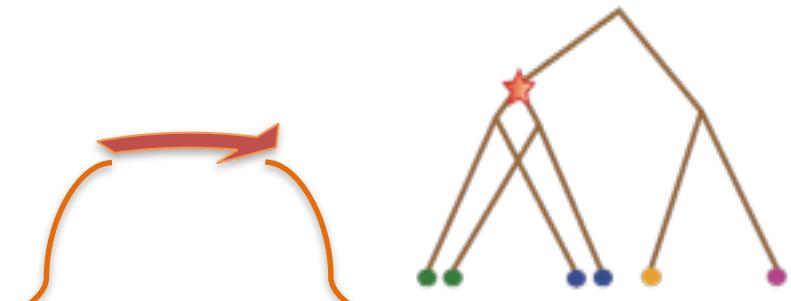
Homo sapiens;GeneB: ACTGTTGATGACATGAC...

Mus musculus;GeneC: ACTGATGATGACAAGAC...

Mus musculus;GeneD: ACTGGTGA--CCATGAC...

Bison bison; GeneE: ACTGGTGATGACACCGAC...

Canis lupus; GeneF: ACT--TCATGAAACGAC...



- Parsimony
- Model-based approaches: e.g. Felsenstein pruning algorithm (1981) to compute $P(\text{alignment} | \text{Gene tree})$

The usual approach to reconstructing phylogenetic trees

Homo sapiens;GeneA: ACTGGTGATGACATAAC...

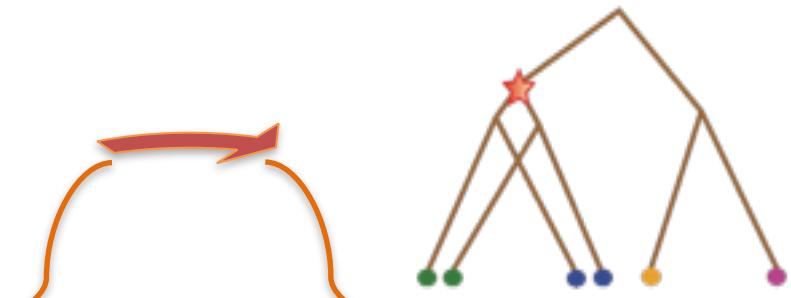
Homo sapiens;GeneB: ACTGTTGATGACATGAC...

Mus musculus;GeneC: ACTGATGATGACAAGAC...

Mus musculus;GeneD: ACTGGTGA--CCATGAC...

Bison bison; GeneE: ACTGGTGATGACACCGAC...

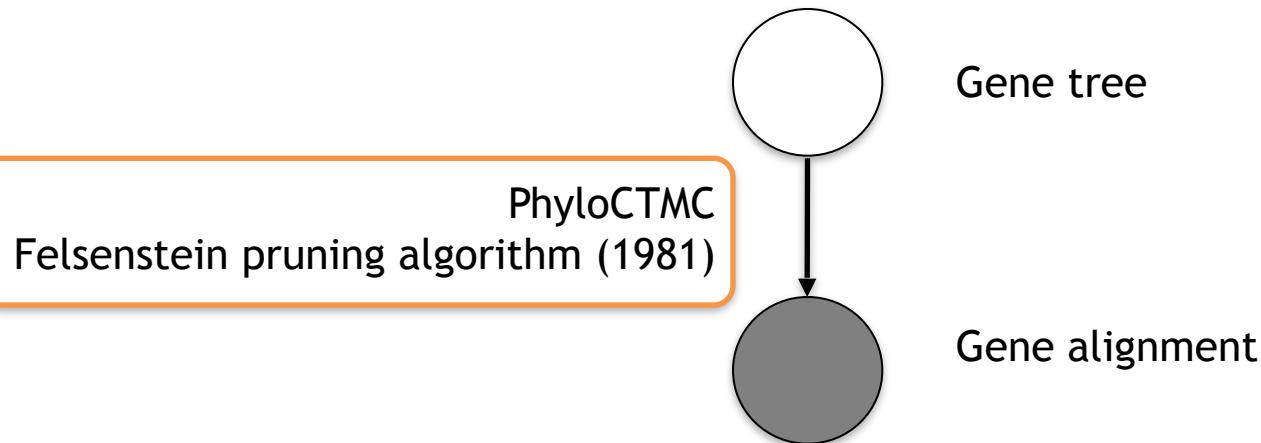
Canis lupus; GeneF: ACT--TCATGAAACGAC...



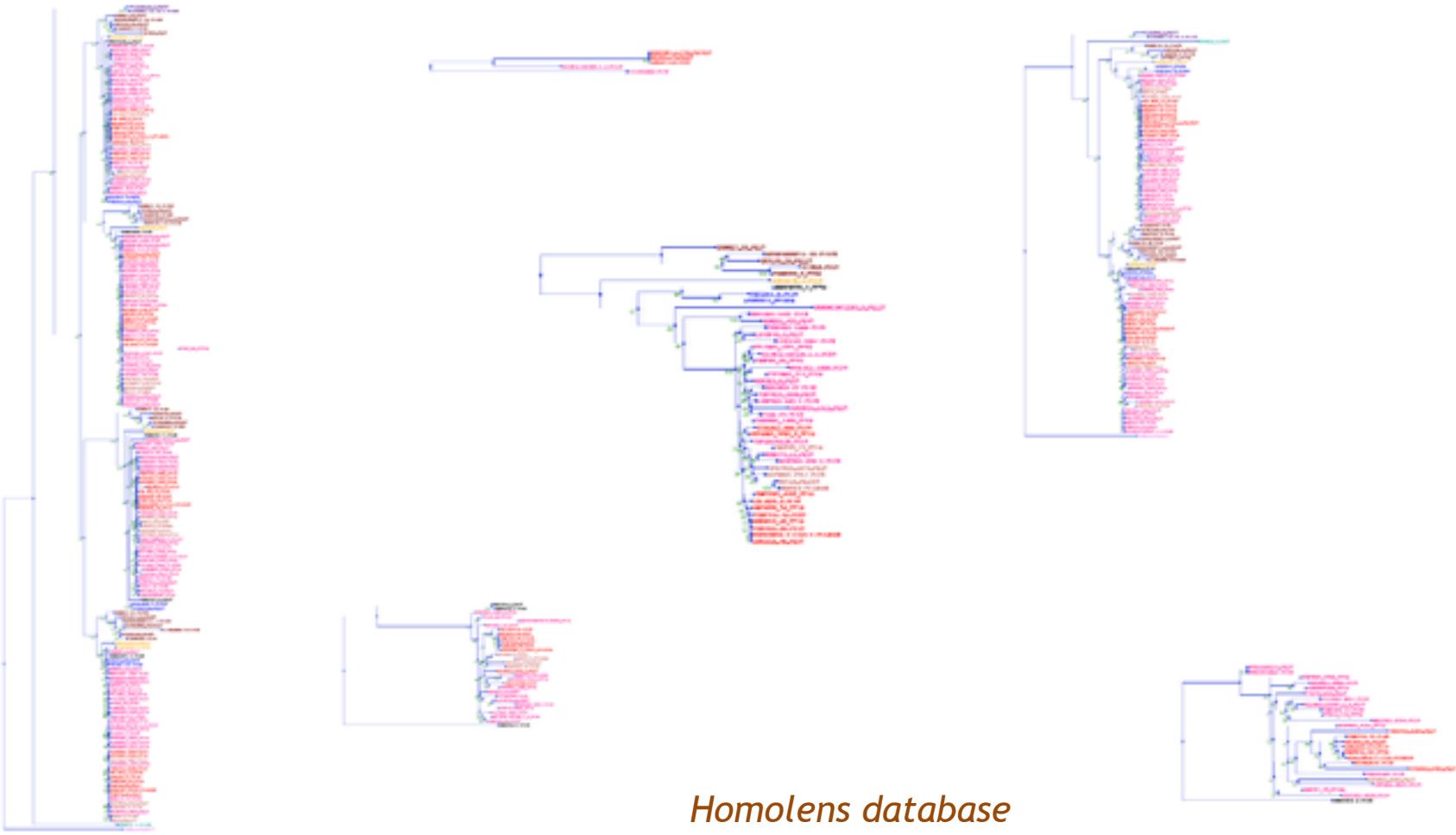
- Parsimony
- Model-based approaches: e.g. Felsenstein pruning algorithm (1981) to compute $P(\text{alignment} | \text{Gene tree})$

No species tree object in the usual approach

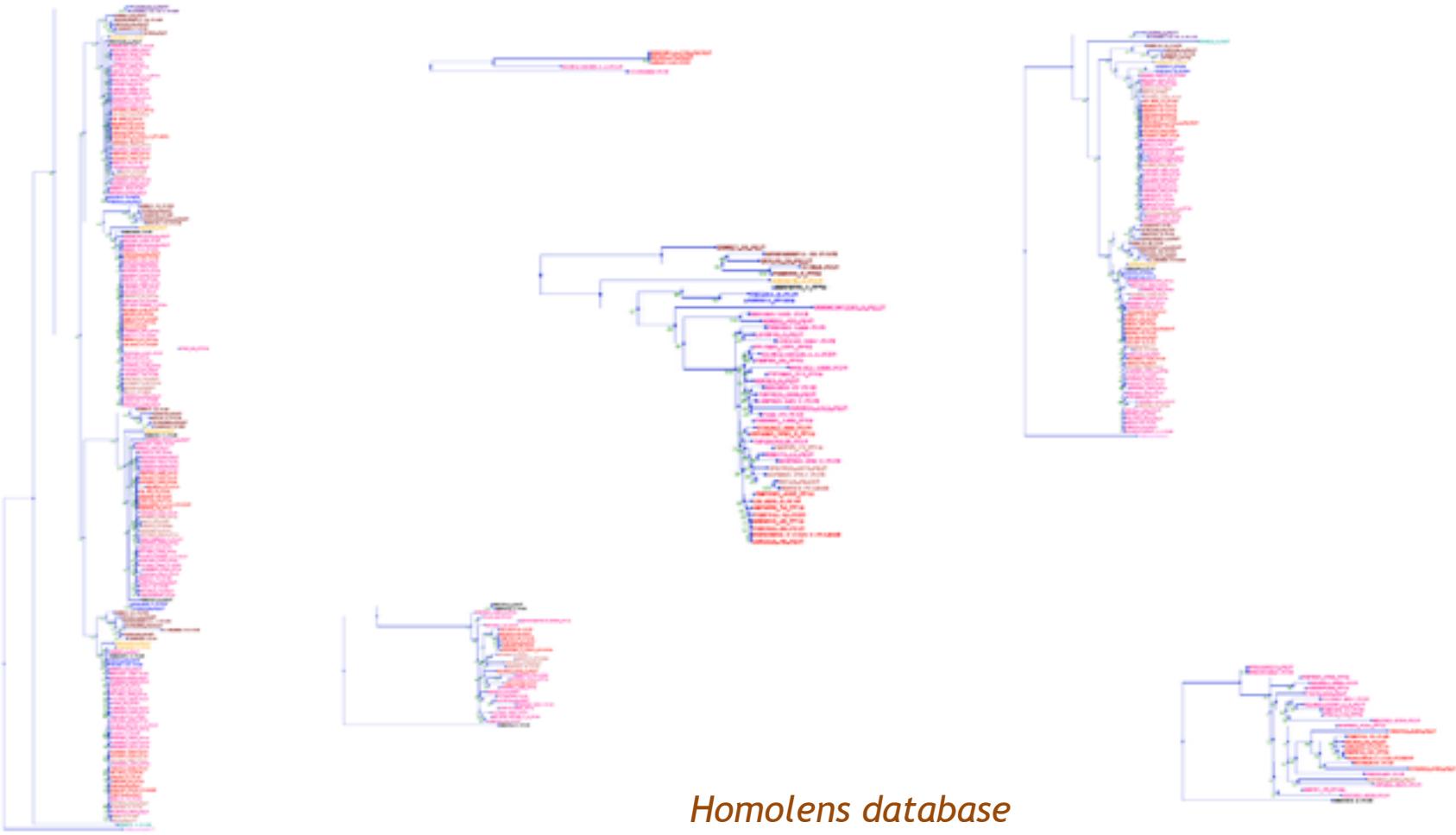
The gene tree graphical model



Gene trees inferred from sequences alone



Gene trees inferred from sequences alone



Complex, messy, not to be trusted

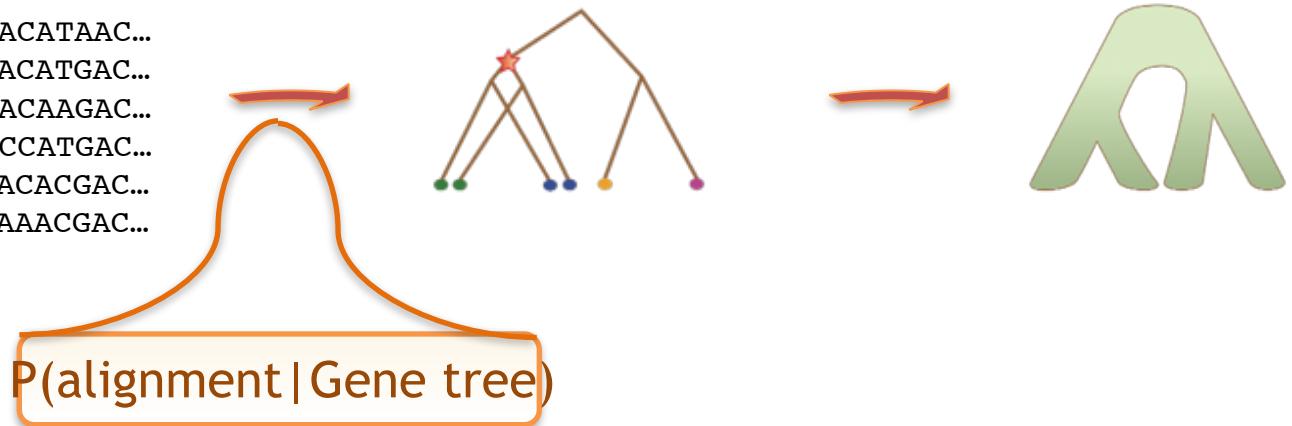
The species tree-gene tree approach to reconstructing phylogenetic trees

Homo sapiens GeneA: ACTGGTGATGACATAAC...
Homo sapiens GeneB: ACTGTTGATGACATGAC...
Mus musculus GeneC: ACTGATGATGACAAGAC...
Mus musculus GeneD: ACTGGTGA--CCATGAC...
Bison bison; GeneE: ACTGGTGATGACACGAC...
Canis lupus; GeneF: ACT--TCATGAAACGAC...



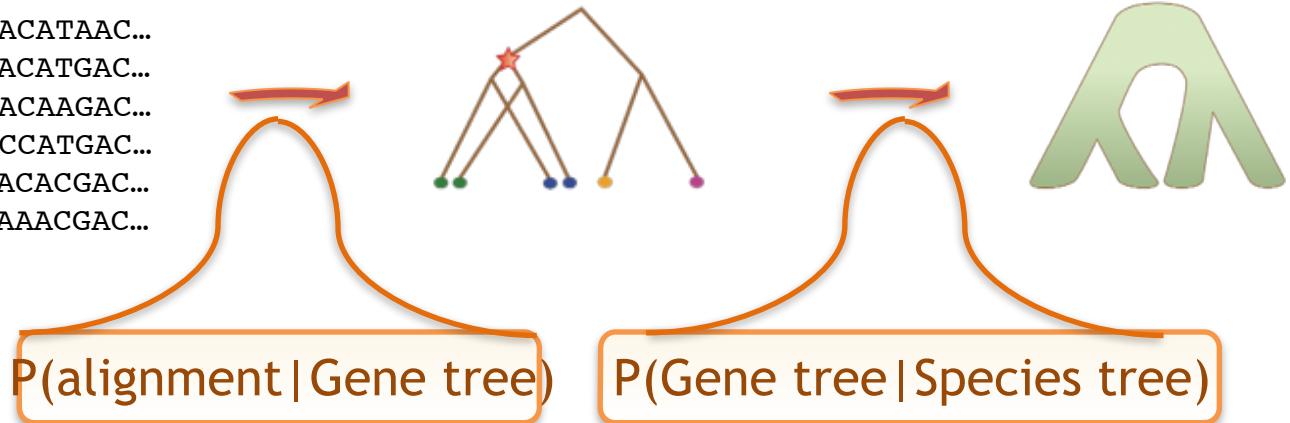
The species tree-gene tree approach to reconstructing phylogenetic trees

Homo sapiens GeneA: ACTGGTGATGACATAAC...
Homo sapiens GeneB: ACTGTTGATGACATGAC...
Mus musculus GeneC: ACTGATGATGACAAGAC...
Mus musculus GeneD: ACTGGTGA--CCATGAC...
Bison bison; GeneE: ACTGGTGATGACACGAC...
Canis lupus; GeneF: ACT--TCATGAAACGAC...



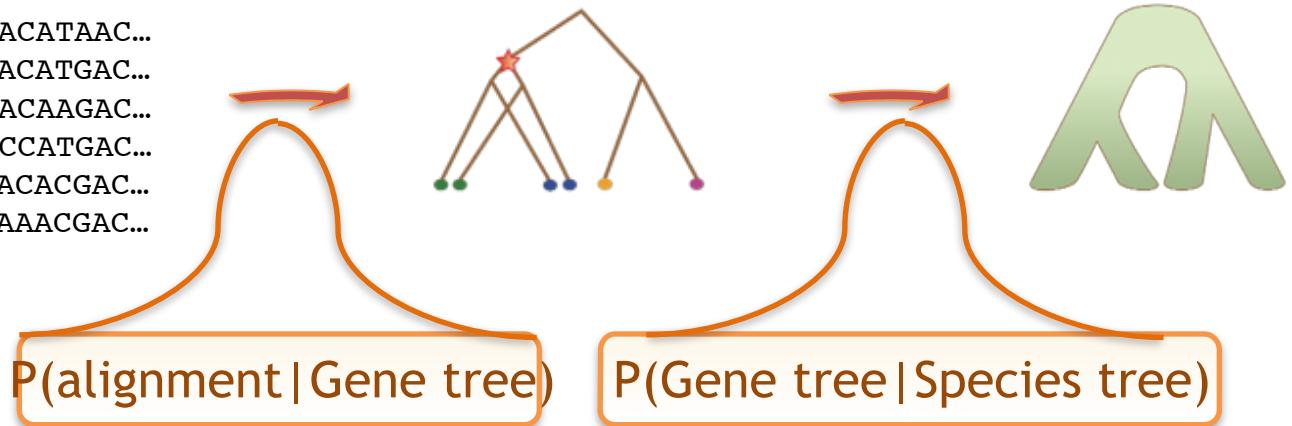
The species tree-gene tree approach to reconstructing phylogenetic trees

Homo sapiens GeneA: ACTGGTGATGACATAAC...
Homo sapiens GeneB: ACTGTTGATGACATGAC...
Mus musculus GeneC: ACTGATGATGACAAGAC...
Mus musculus GeneD: ACTGGTGA--CCATGAC...
Bison bison; GeneE: ACTGGTGATGACACGAC...
Canis lupus; GeneF: ACT--TCATGAAACGAC...



The species tree-gene tree approach to reconstructing phylogenetic trees

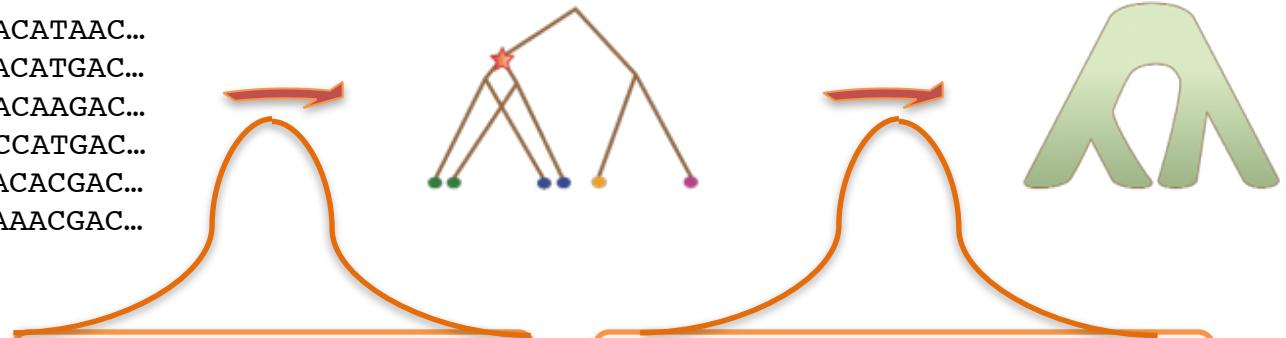
Homo sapiens GeneA: ACTGGTGATGACATAAC...
Homo sapiens GeneB: ACTGTTGATGACATGAC...
Mus musculus GeneC: ACTGATGATGACAAGAC...
Mus musculus GeneD: ACTGGTGA--CCATGAC...
Bison bison; GeneE: ACTGGTGATGACACGAC...
Canis lupus; GeneF: ACT--TCATGAAACGAC...



$$P(\text{alignment, Gene tree} | \text{Species tree}) = P(\text{alignment} | \text{Gene tree}) \times P(\text{Gene tree} | \text{Species tree})$$

The species tree-gene tree approach to reconstructing phylogenetic trees

Homo sapiens GeneA: ACTGGTGATGACATAAC...
Homo sapiens GeneB: ACTGTTGATGACATGAC...
Mus musculus GeneC: ACTGATGATGACAAGAC...
Mus musculus GeneD: ACTGGTGA--CCATGAC...
Bison bison; GeneE: ACTGGTGATGACACGAC...
Canis lupus; GeneF: ACT--TCATGAAACGAC...



$P(\text{alignment} | \text{Gene tree})$

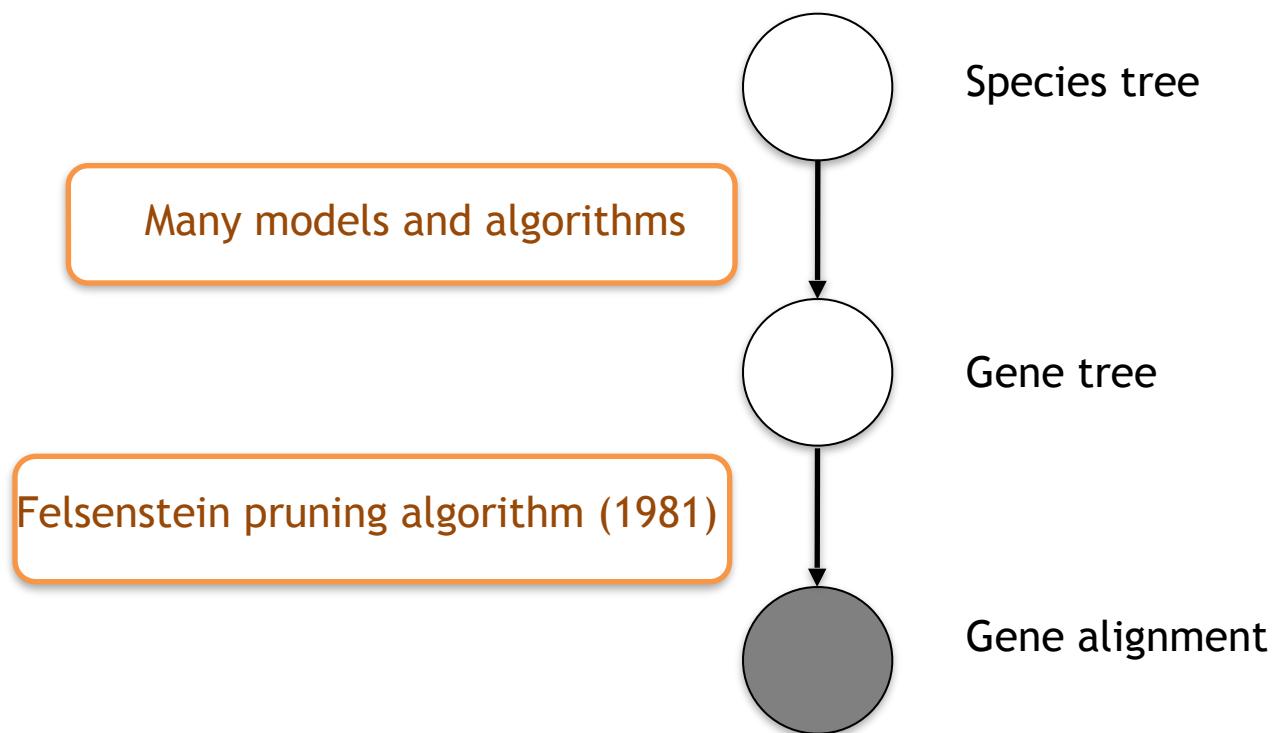
$P(\text{Gene tree} | \text{Species tree})$

Felsenstein pruning algorithm (1981)

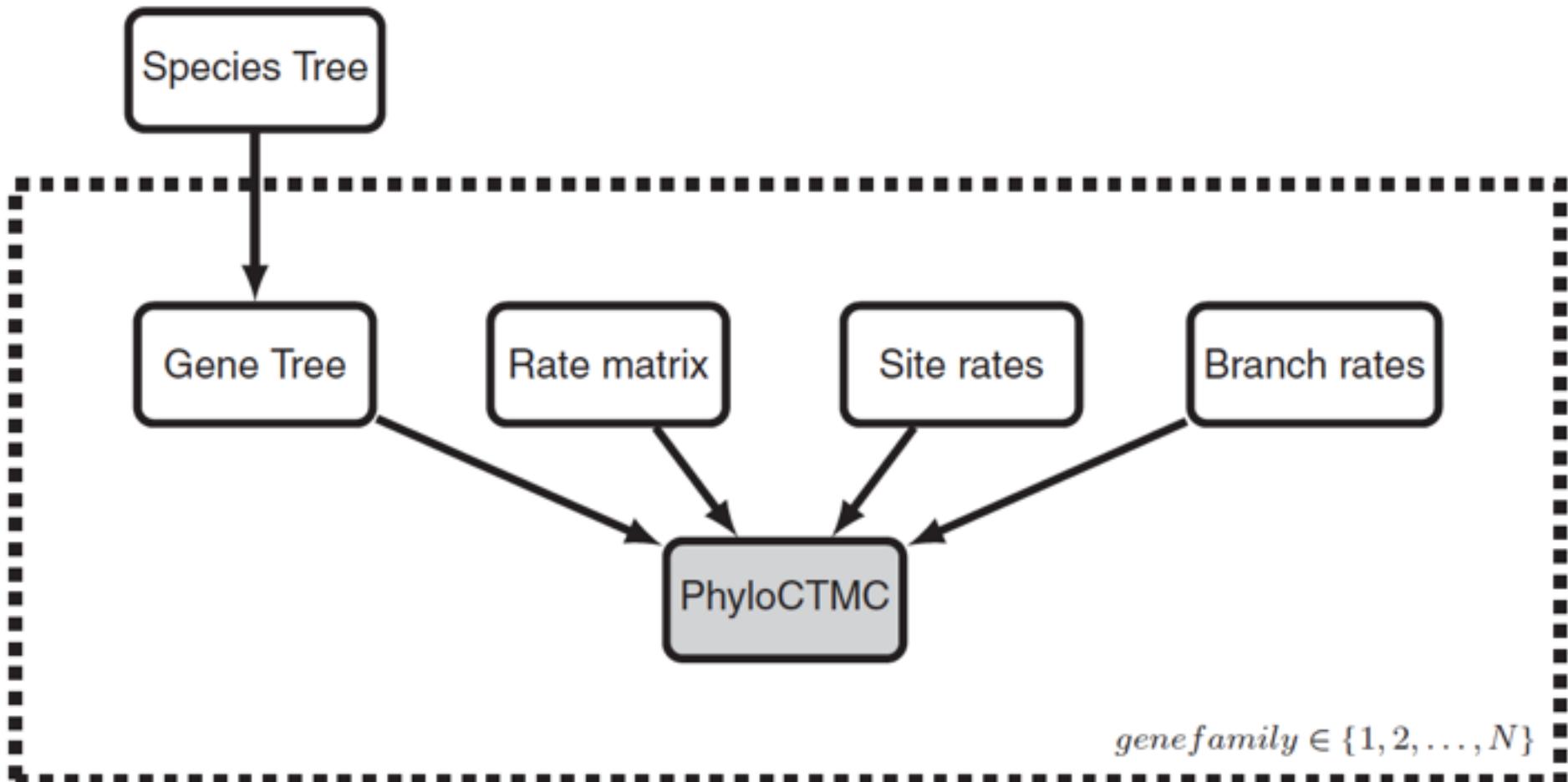
Many models and algorithms

$$P(\text{alignment}, \text{Gene tree} | \text{Species tree}) = P(\text{alignment} | \text{Gene tree}) \times P(\text{Gene tree} | \text{Species tree})$$

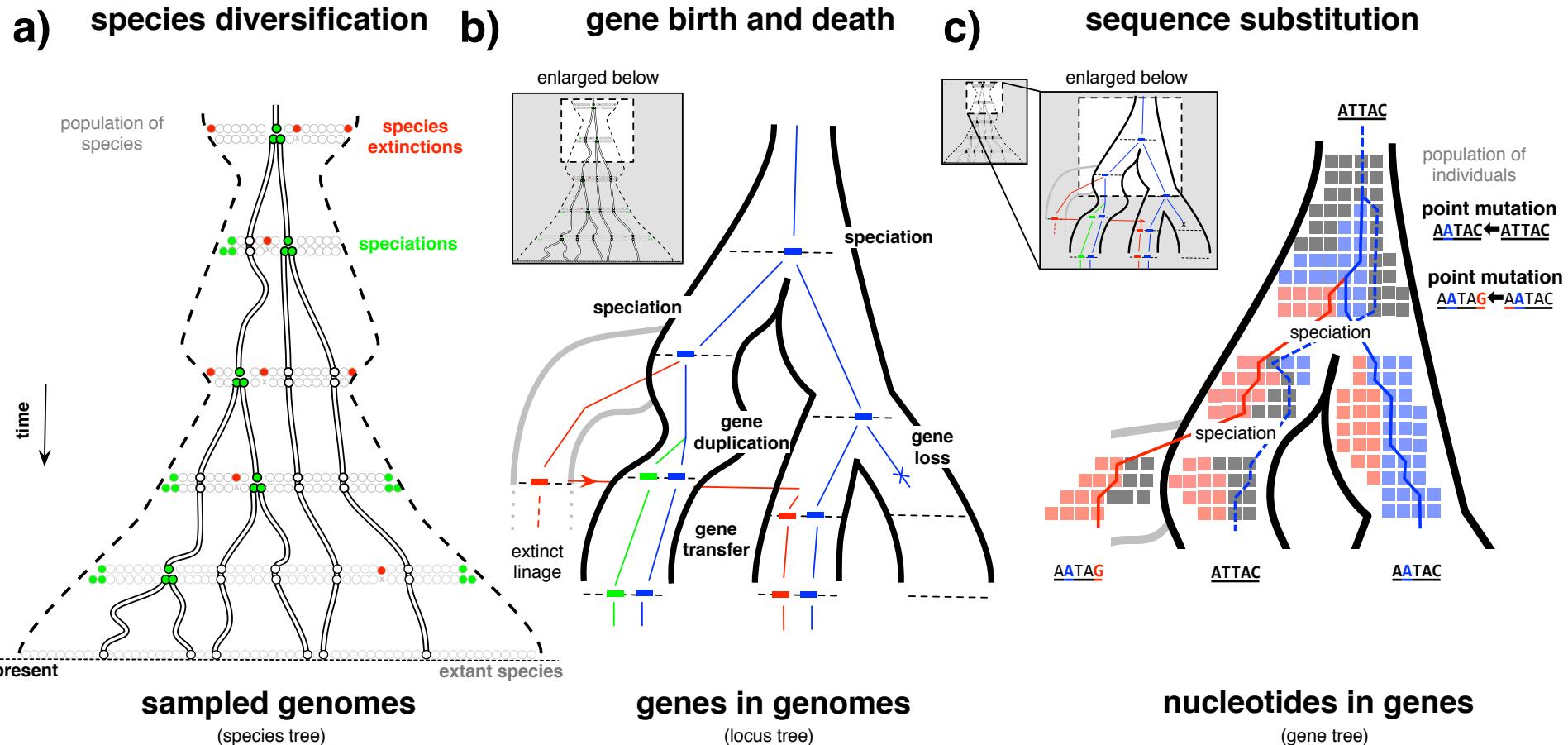
The species tree-gene tree graphical model



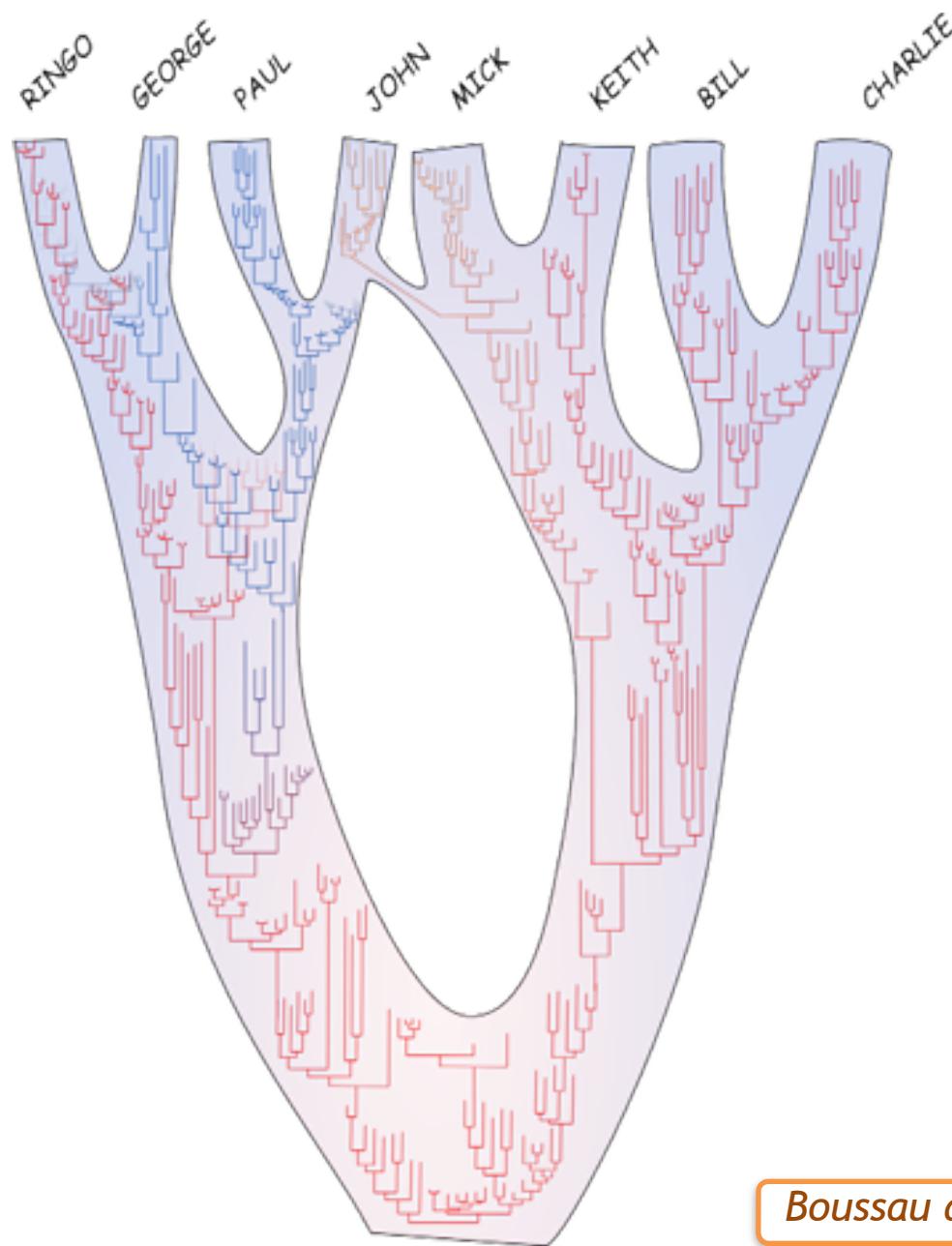
The species tree-gene tree graphical model



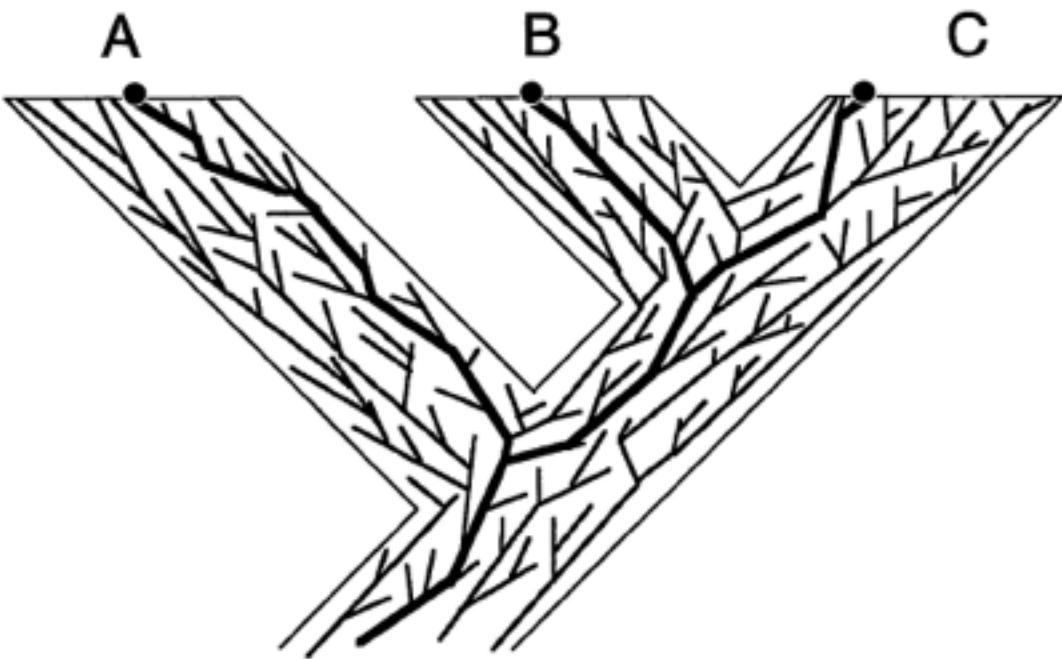
A hierarchy of processes



Gene tree and species tree are linked

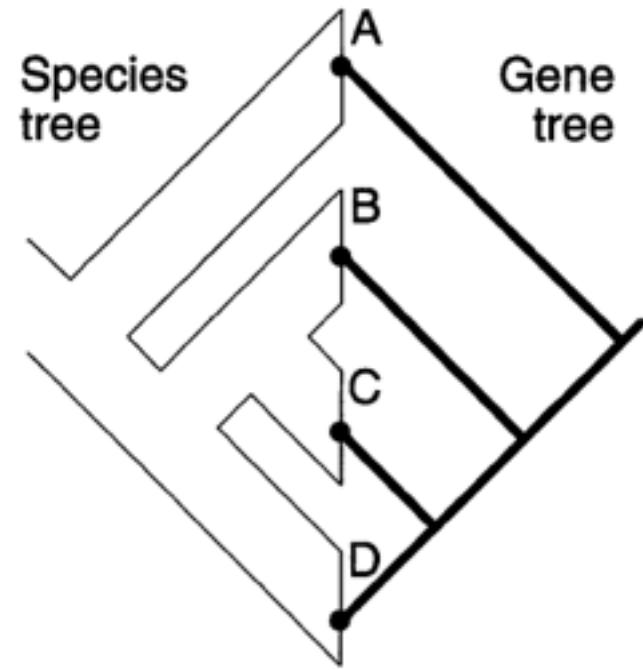
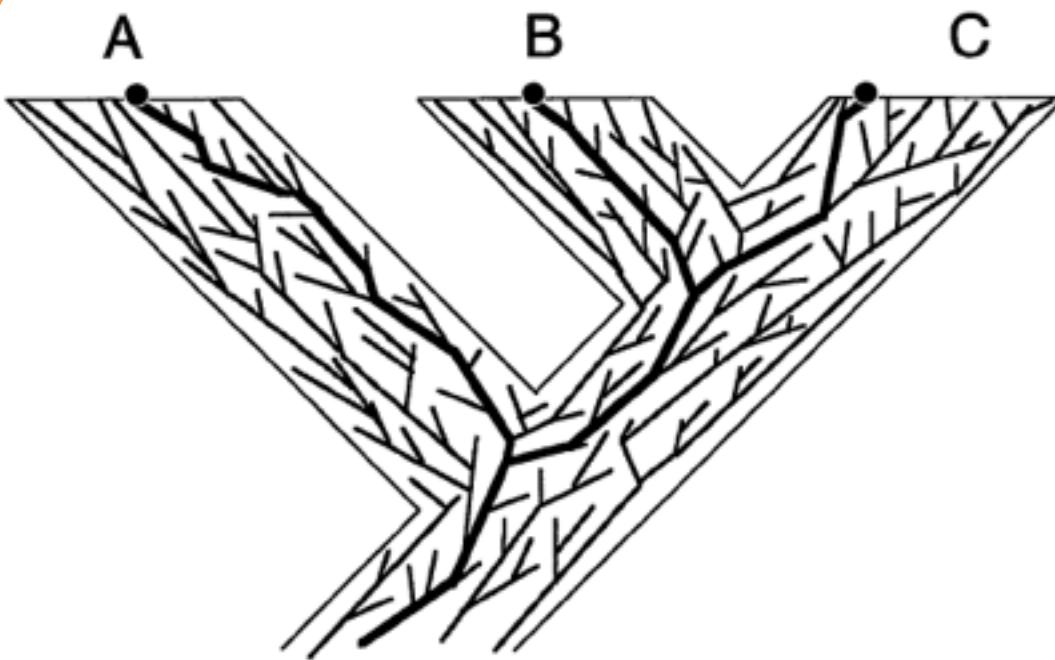


Gene trees in species trees



Maddison, *Syst. Biol.* 1997

Gene trees in species trees



Maddison, *Syst. Biol.* 1997

Challenges

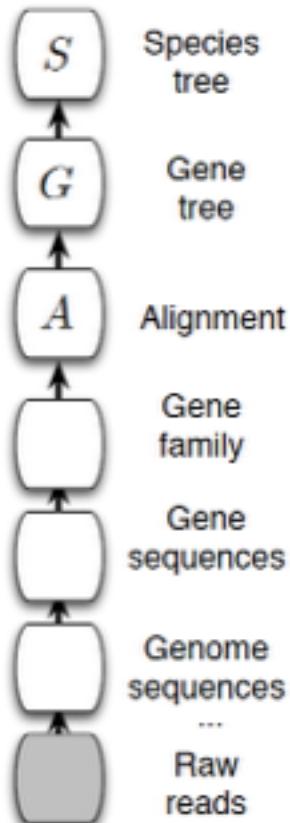
Given gene trees: reconstruct species tree

Given species tree and alignment: reconstruct gene tree

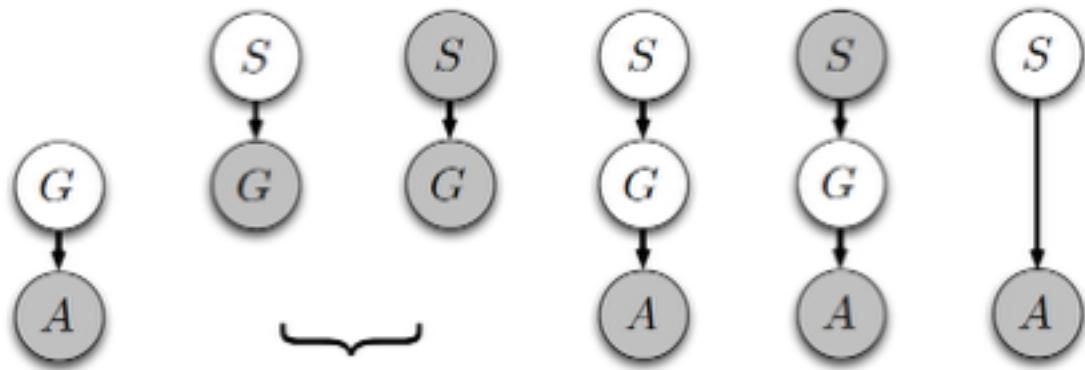
Given alignments: reconstruct genes and species trees

The landscape of gene tree-species tree methods

Phylogenomics Inference pipeline



Gene tree-species tree models published in the literature



$P(A|G)$

PhyML [Guindon & Gascuel, 2003];
RAxML [Stamatakis et al., 2005];
MrBayes [Ronquist et al., 2012];
BEAST [Drummond et al., 2012].

$P(G|S)$

BEST [Edwards et al., 2007]; [Rasmussen & Kellis, 2012];
MP-EST [Liu et al., 2010];
ODT [Szöllősi et al., 2012];

$P(A|G, S)$

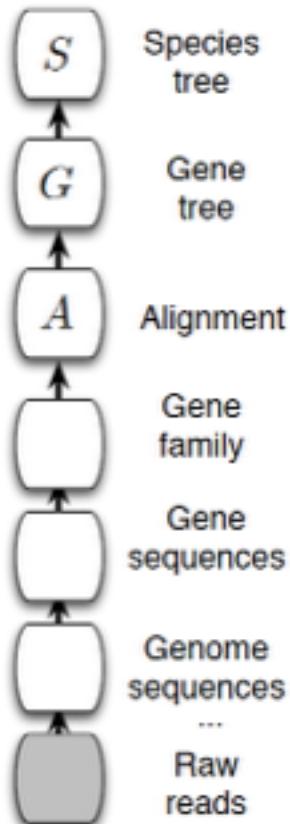
*BEAST [Heled & Drummond, 2008];
PHYLDOG [Roussau et al., 2013].

$P(A|S)$

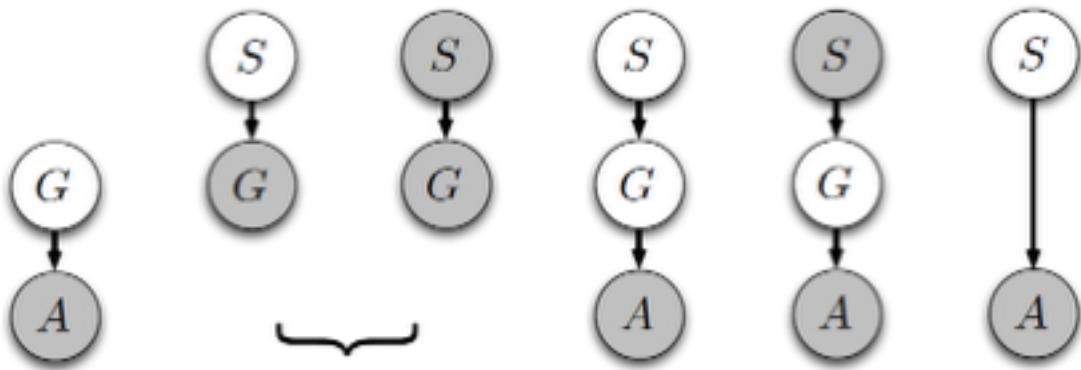
Prime-GSR [Akerborg et al., 2009]; [Bryant et al., 2012];
IMa2 [Hey, 2011];
Prime-DLRS [Sjöstrand et al., 2012];
exODT [Szöllősi et al., 2013a];
POMO (De Maio et al., 2013)

The landscape of gene tree-species tree methods

Phylogenomics Inference pipeline



Gene tree-species tree models published in the literature



$P(A|G)$

PhyML
[Guindon & Gascuel, 2003];
RAxML
[Stamatakis et al., 2005];
MrBayes
[Ronquist et al., 2012];
BEAST
[Drummond et al., 2012].

$P(G|S)$

BEST
[Edwards et al., 2007]; [Rasmussen & Kellis, 2012];
MP-EST
[Liu et al., 2010];
ODT
[Szöllősi et al., 2012];

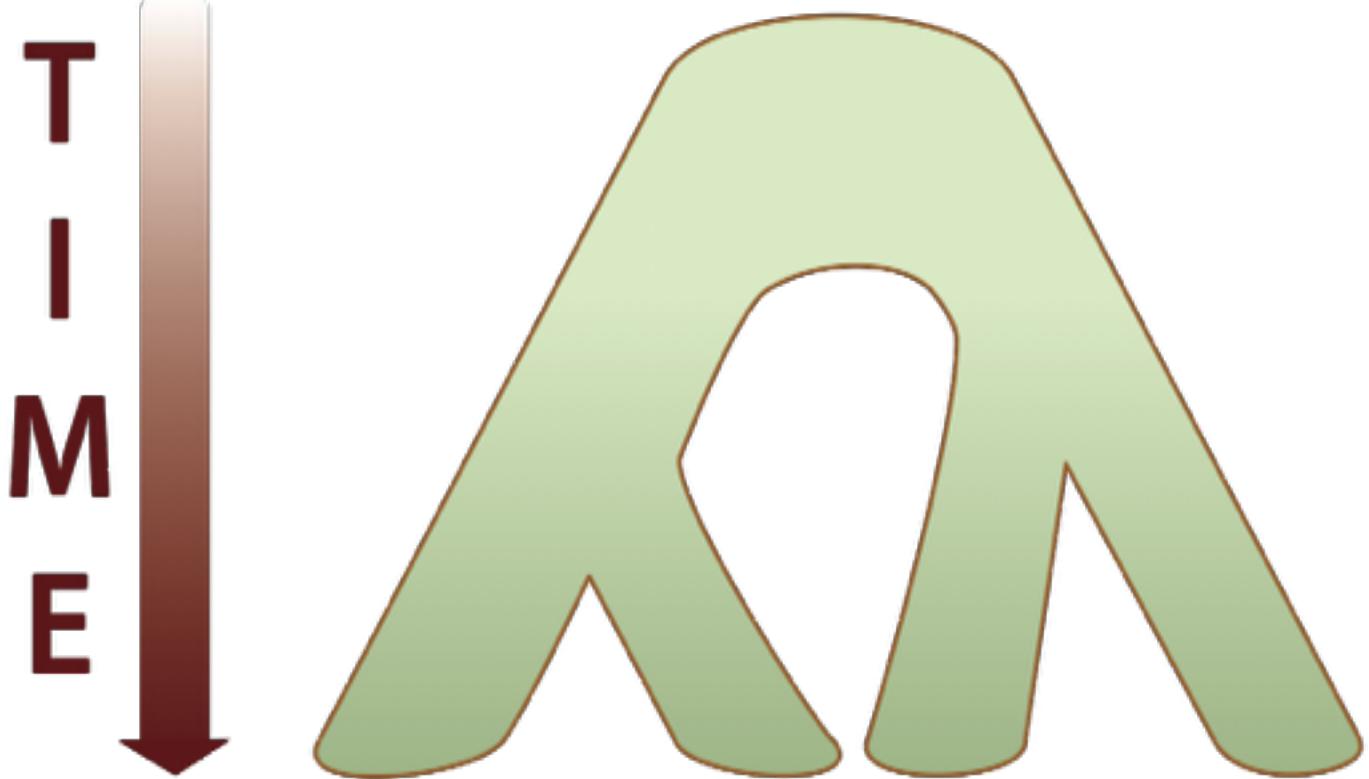
$P(A|G, S)$

*BEAST
[Heled & Drummond, 2008];
PHYLDOG
[Roussau et al., 2013].

$P(A|S)$

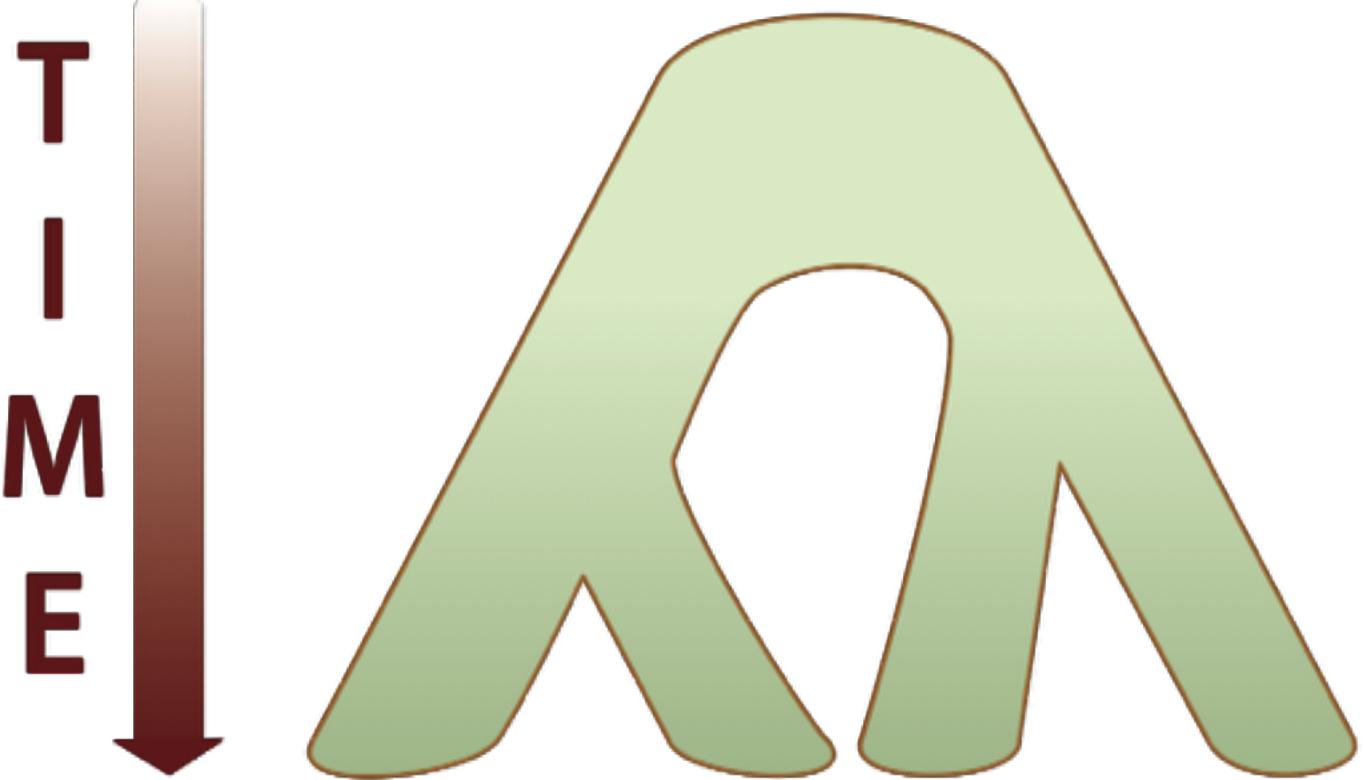
Prime-GSR
[Akerborg et al., 2009]; [Bryant et al., 2012];
IMa2
[Hey, 2011];
Prime-DLRS
[Sjöstrand et al., 2012];
exODT
[Szöllősi et al., 2013a];
POMO
(De Maio et al., 2013)
ALE
[Szöllősi et al., 2013b];

RevBayes



Discrete character: a
Continuous character: 0.1
Species: A

a 0.2 b 0.2 a 0.4
Species: B C D

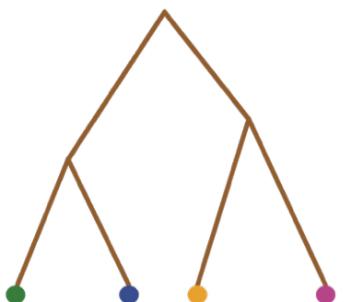
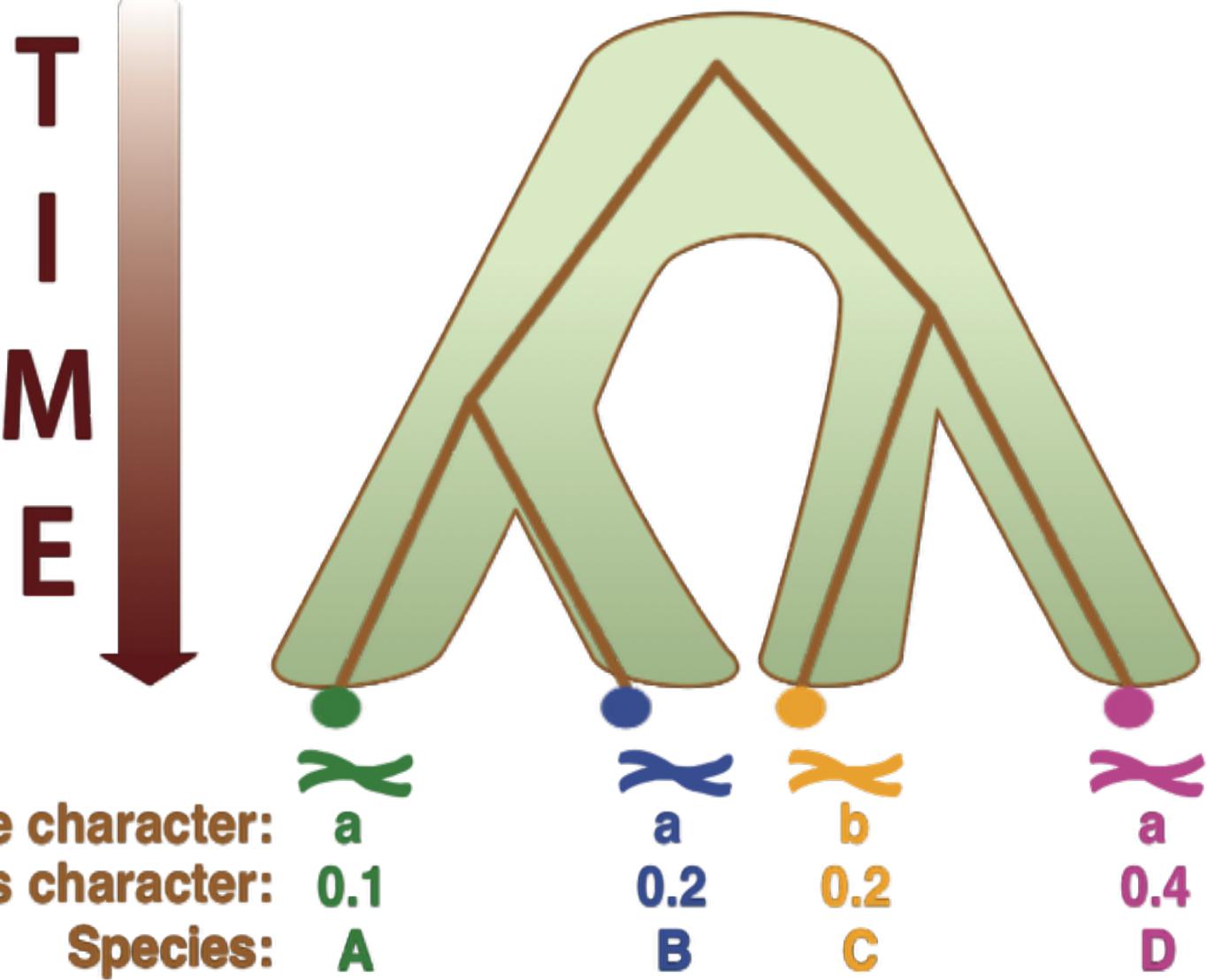


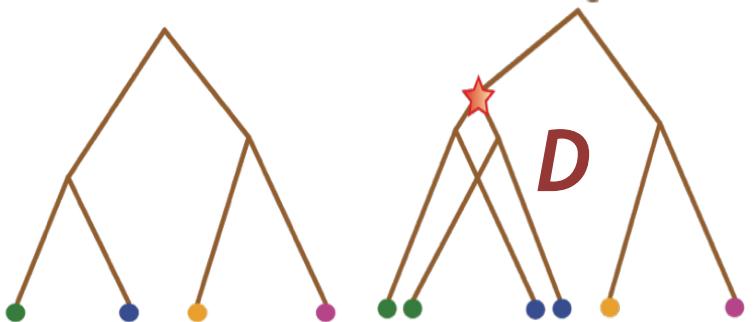
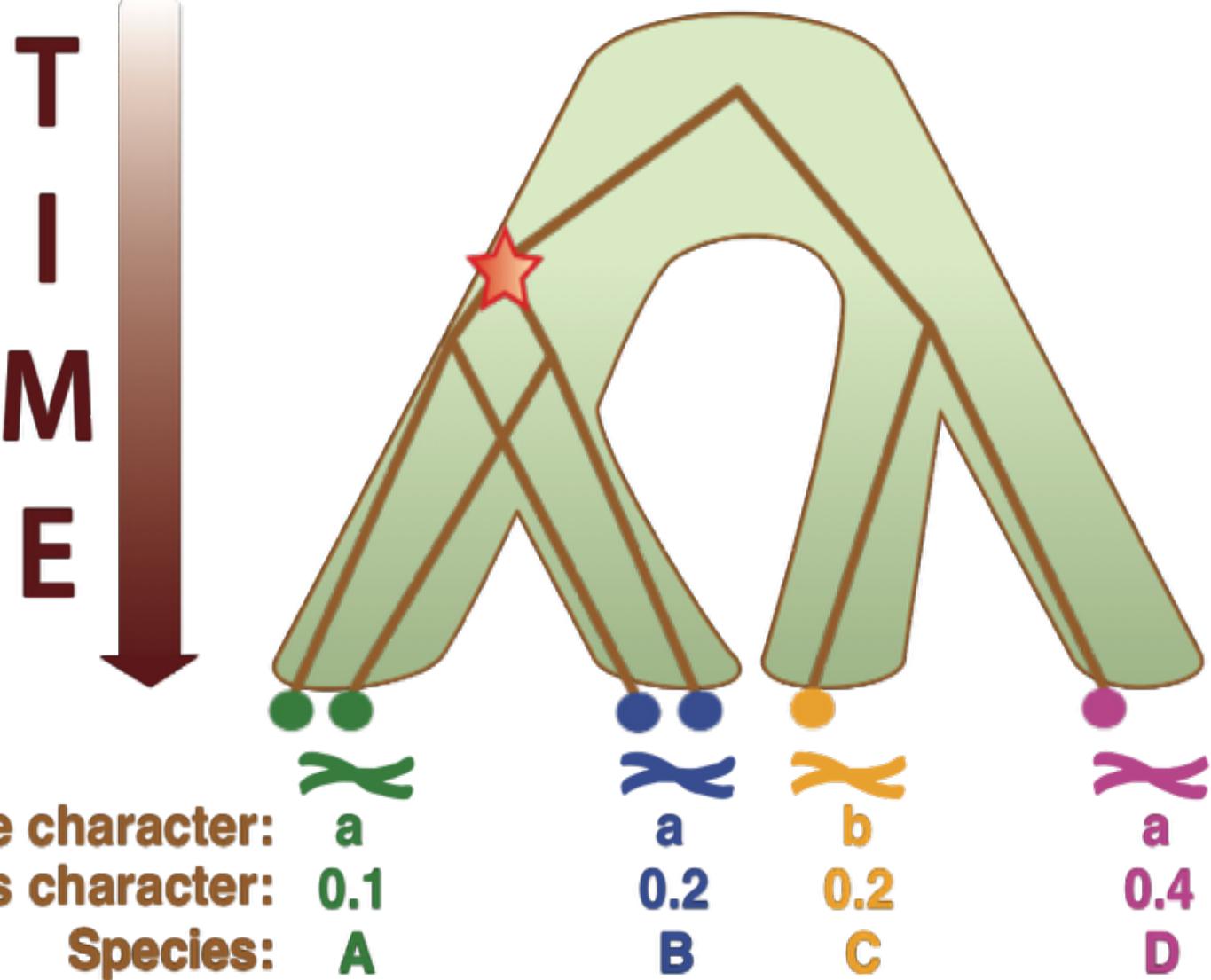
Discrete character: a
Continuous character: 0.1
Species: A

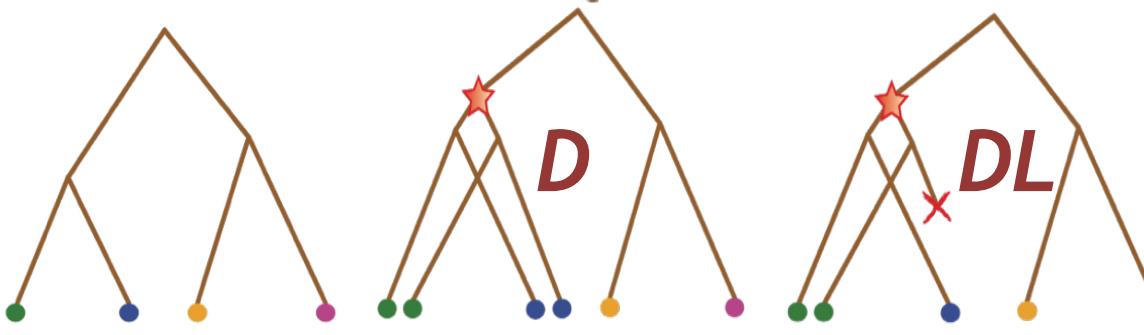
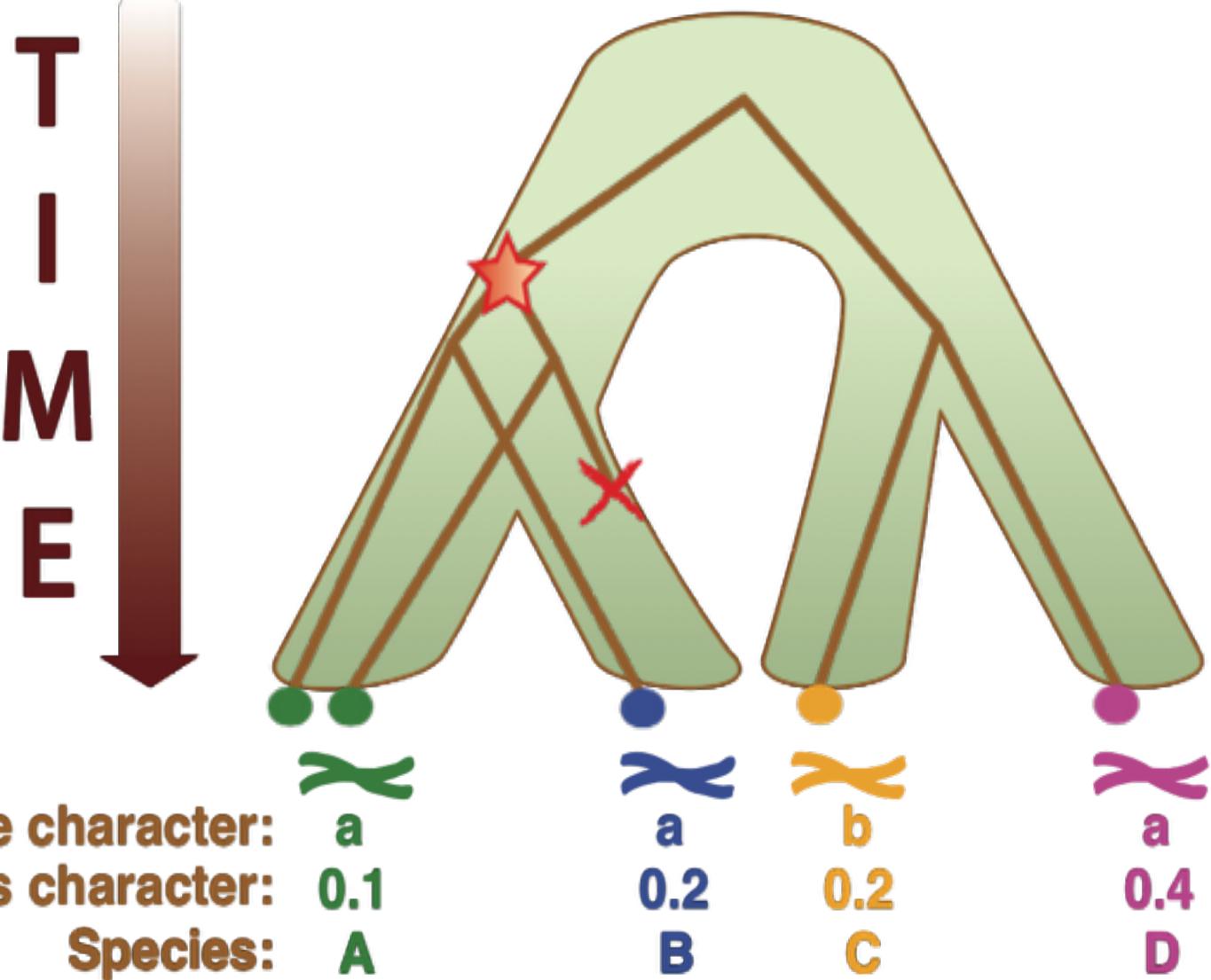
a
0.2
B

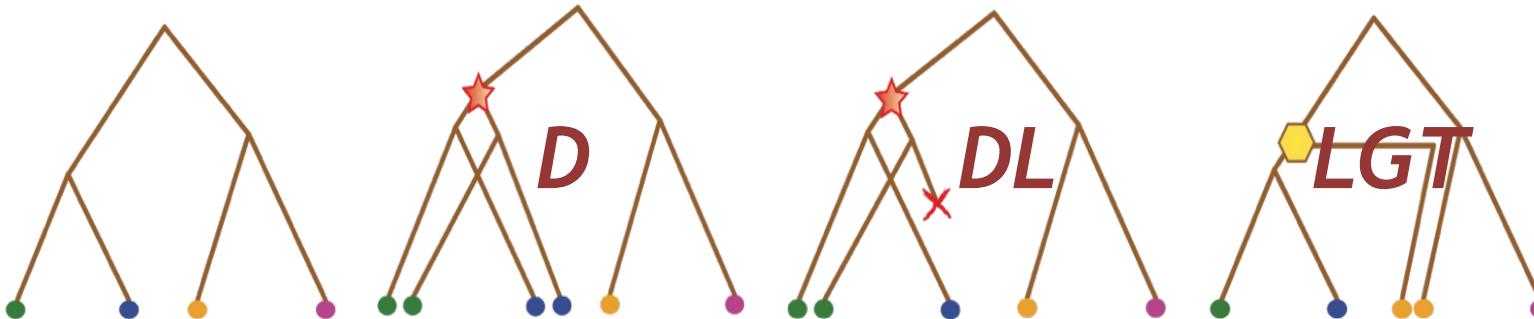
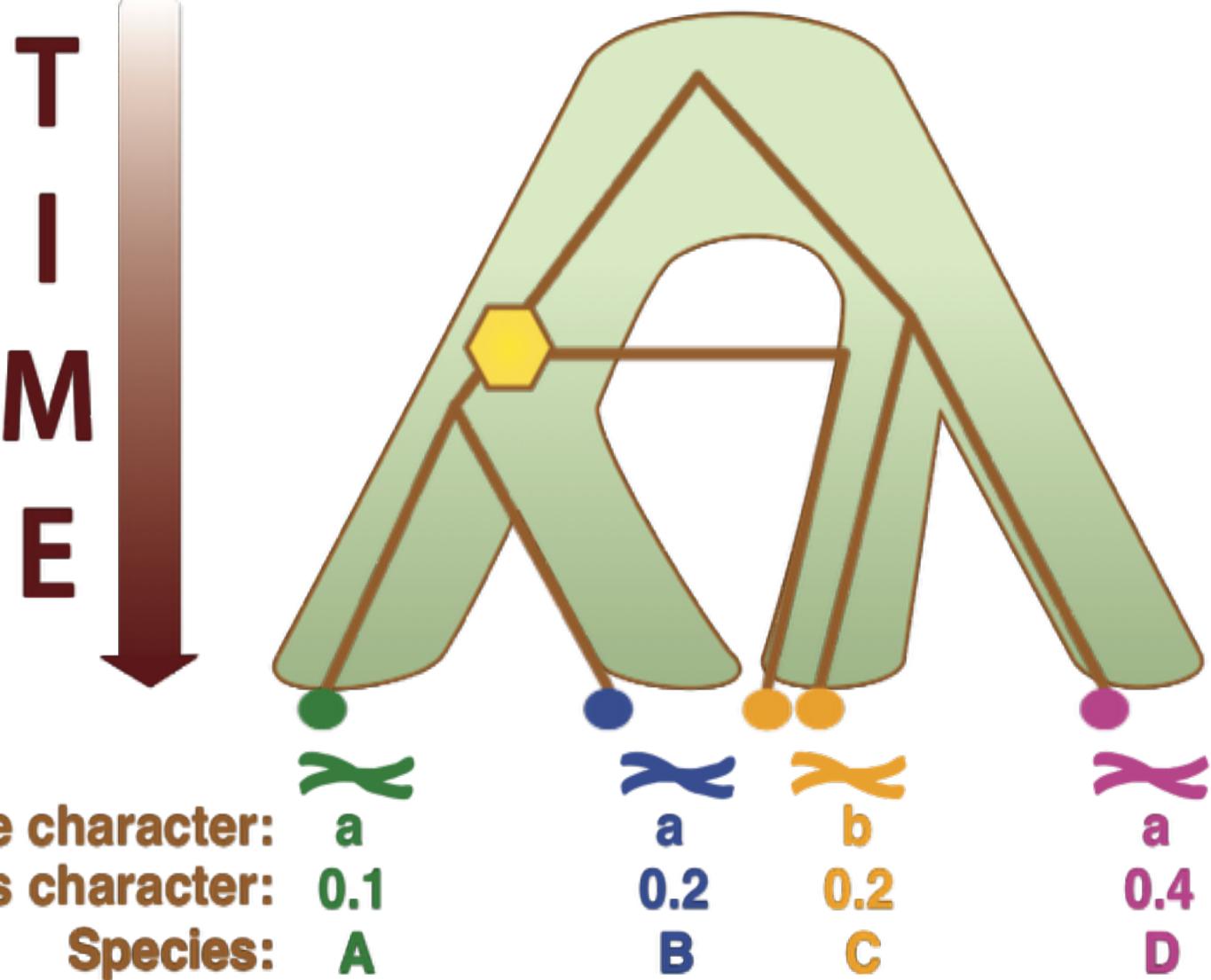
b
0.2
C

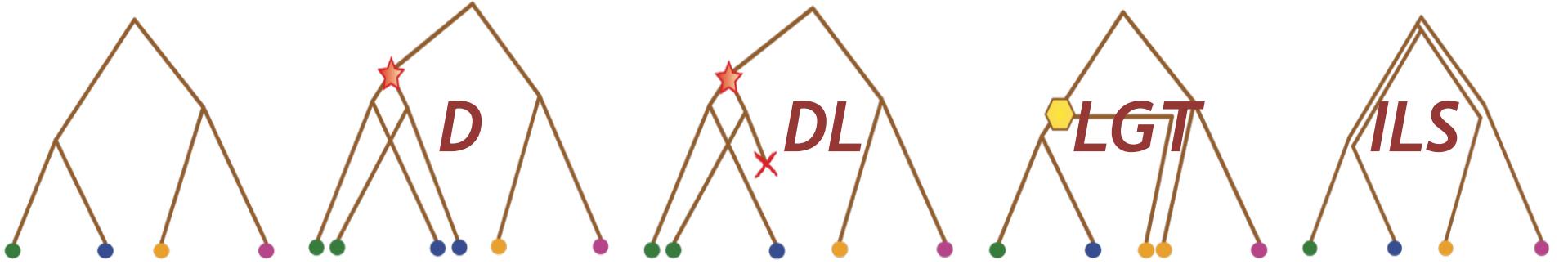
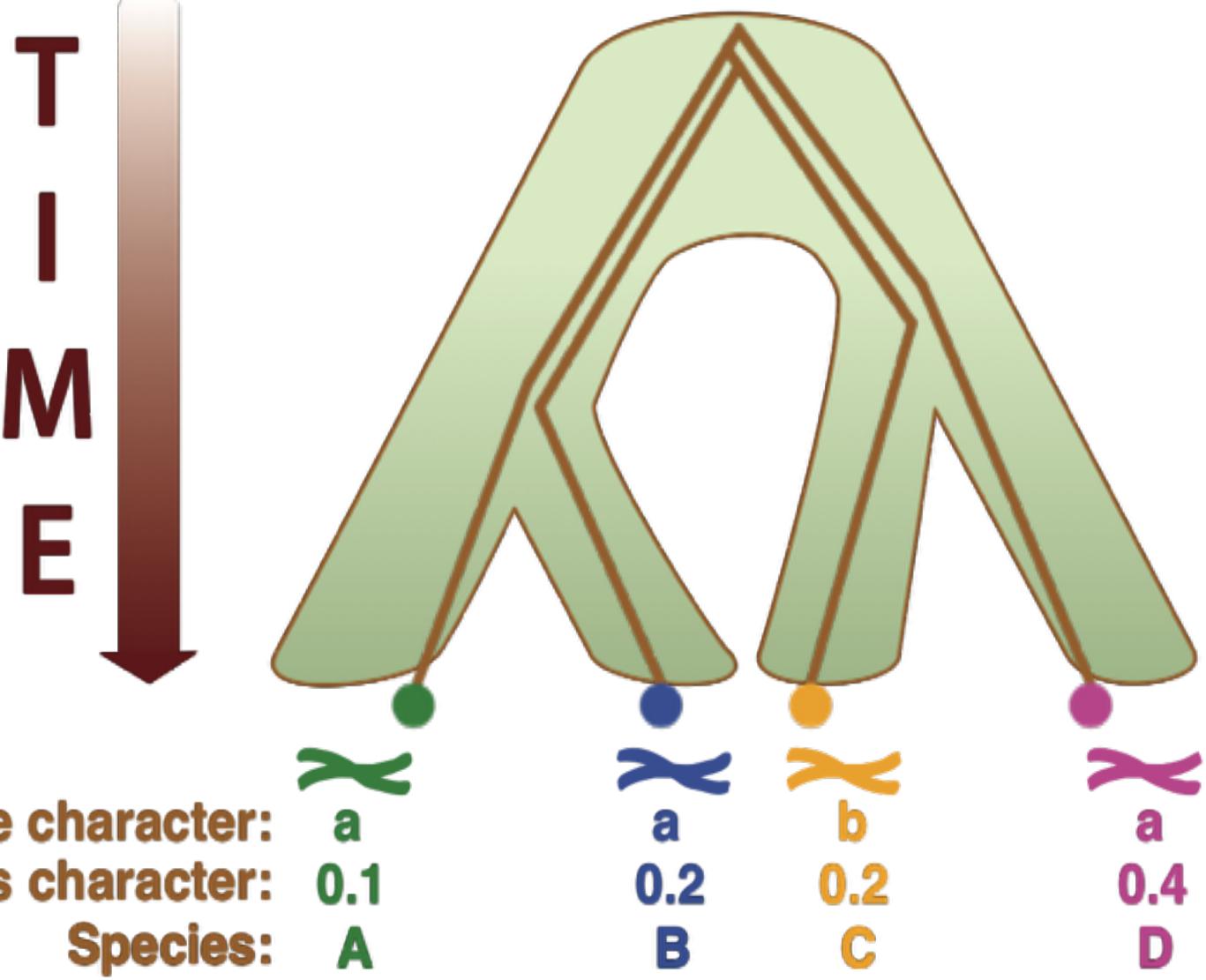
a
0.4
D

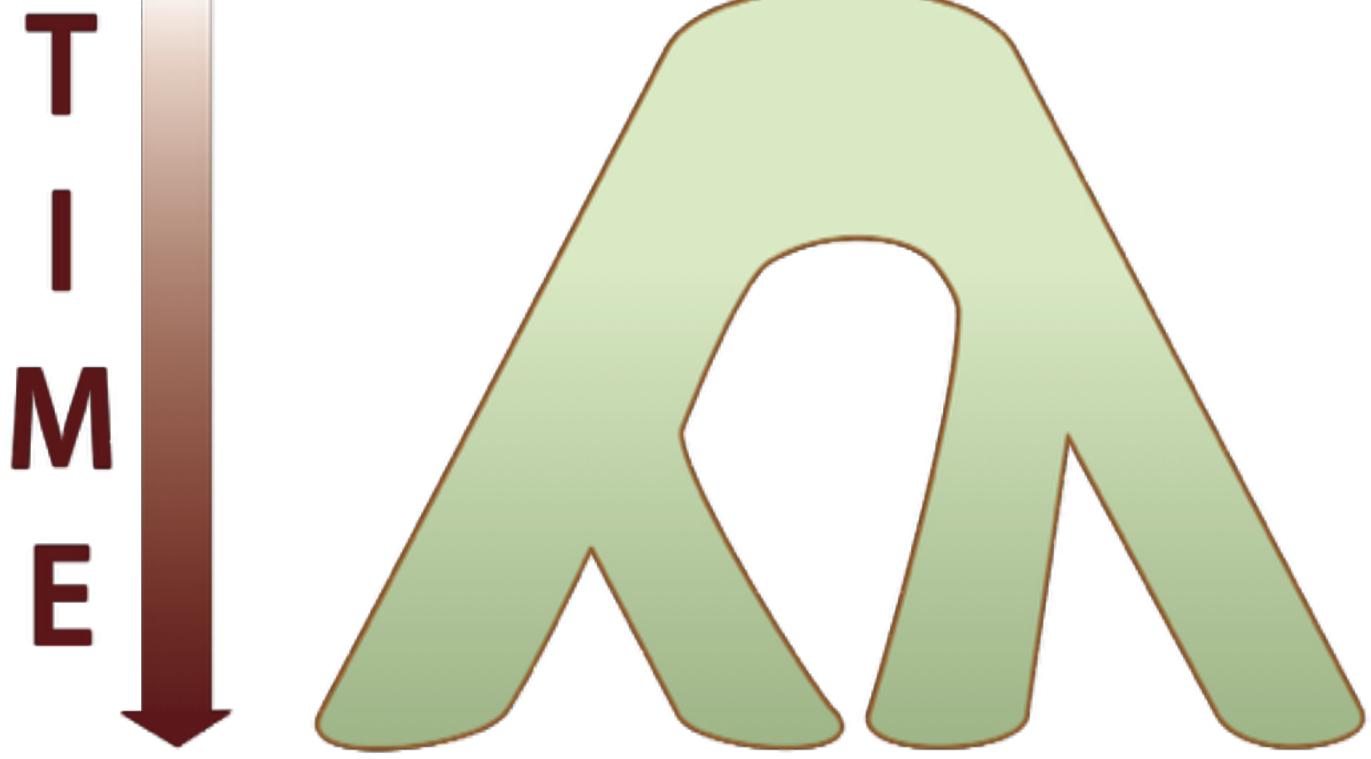






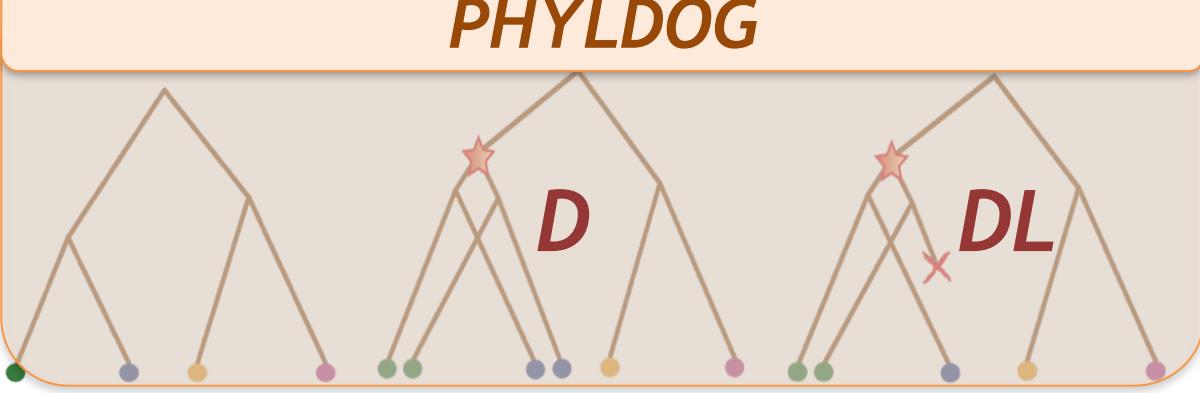






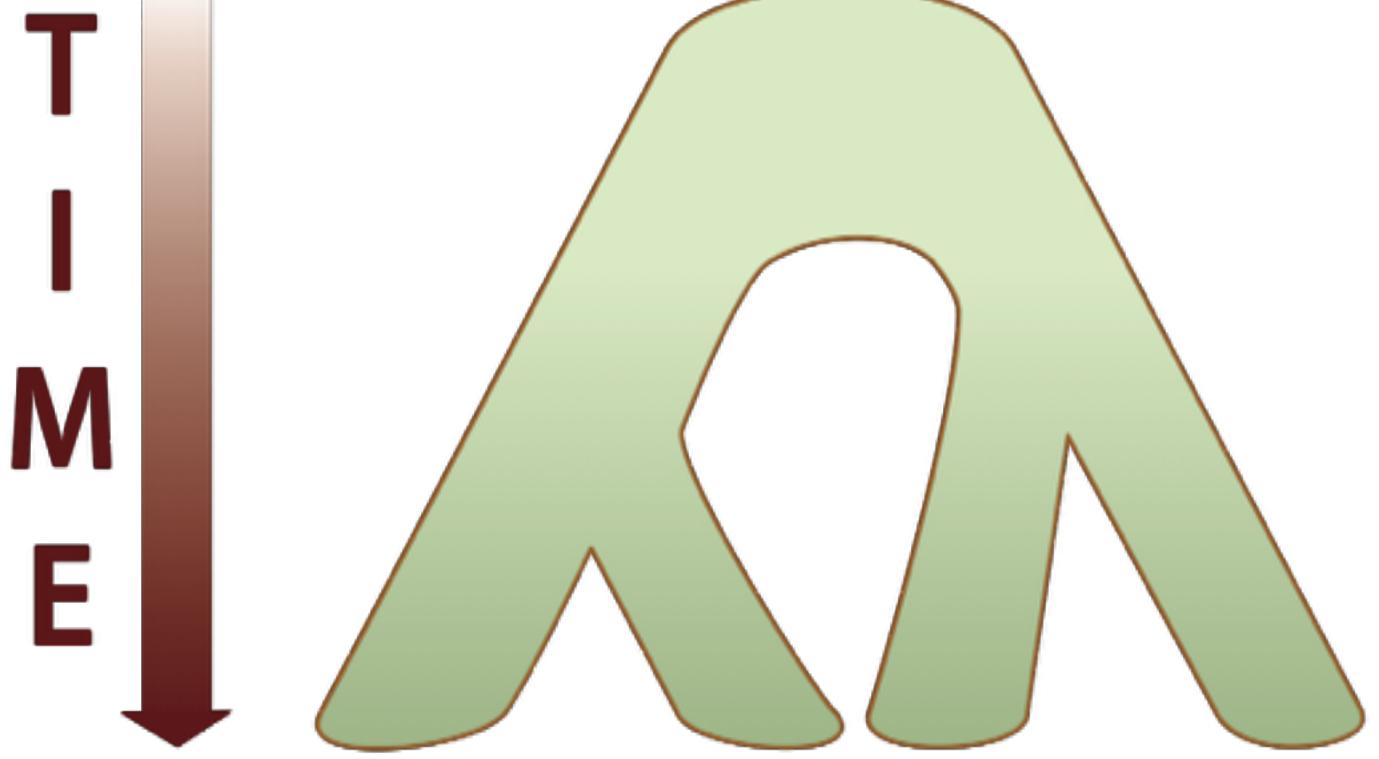
Discrete character: a
 Continuous character: 0.1

PHYLDODG



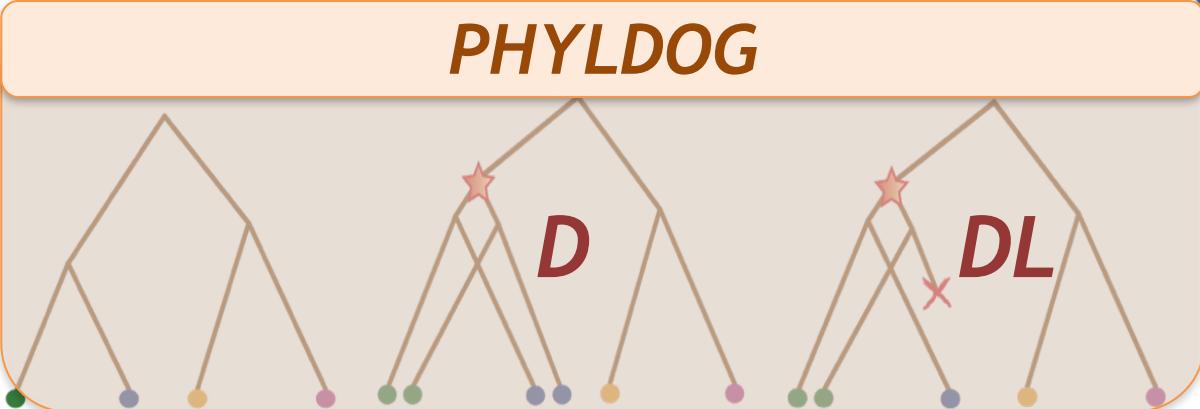
a 0.2 B
 b 0.2 C
 a 0.4 D





Discrete character: a
 Continuous character: 0.1

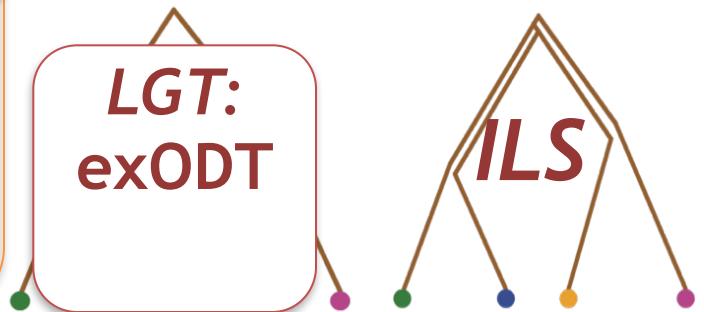
PHYLDOD

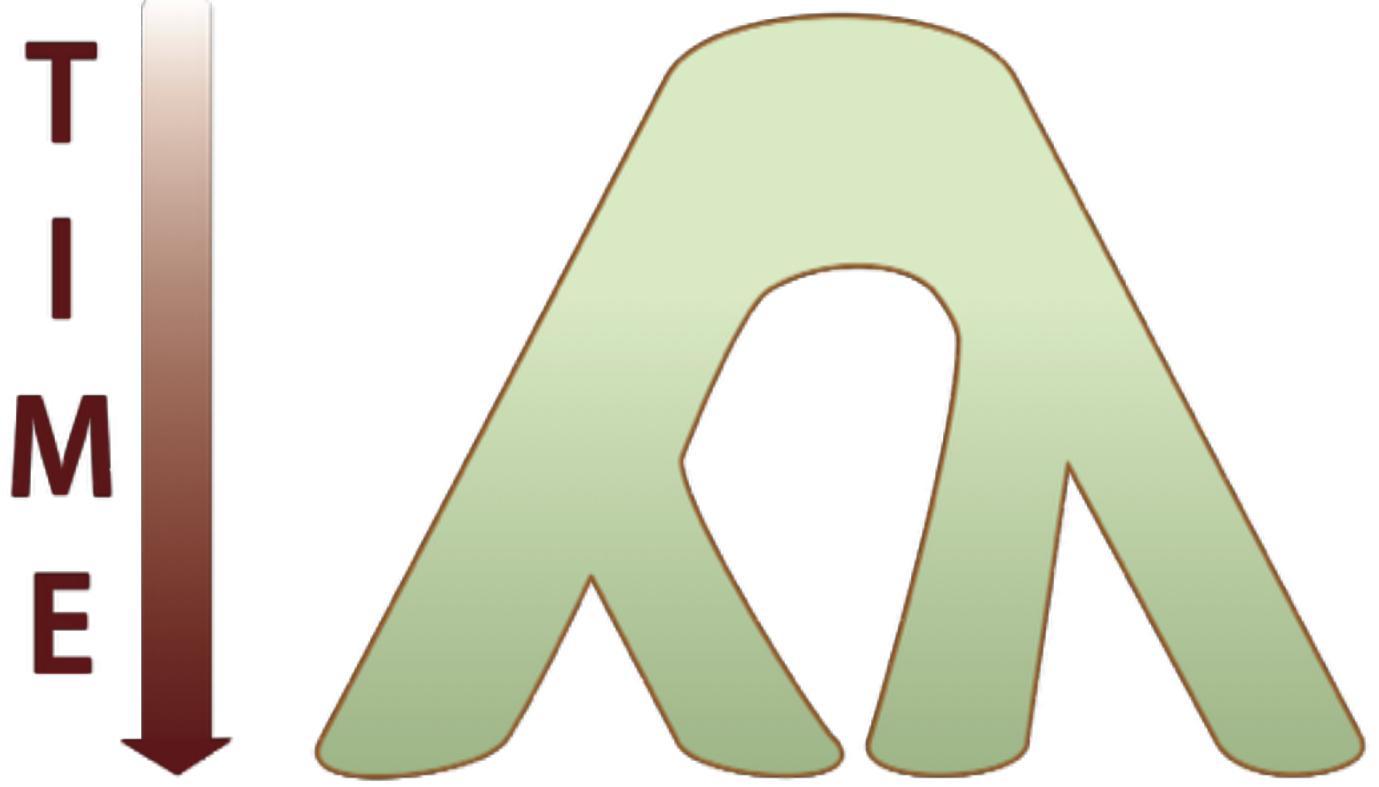


a
 b
 0.2
 B C

a
 0.4
 D

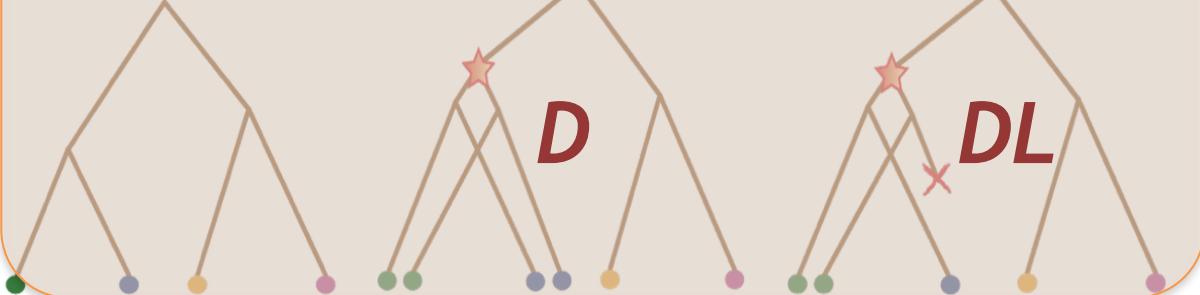
LGT:
exODT





Discrete character: a
 Continuous character: 0.1

PHYLDOD



a
 b
 0.2
 B C

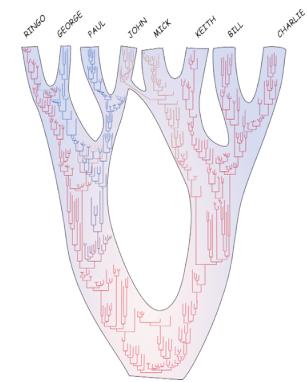
a
 0.4
 D

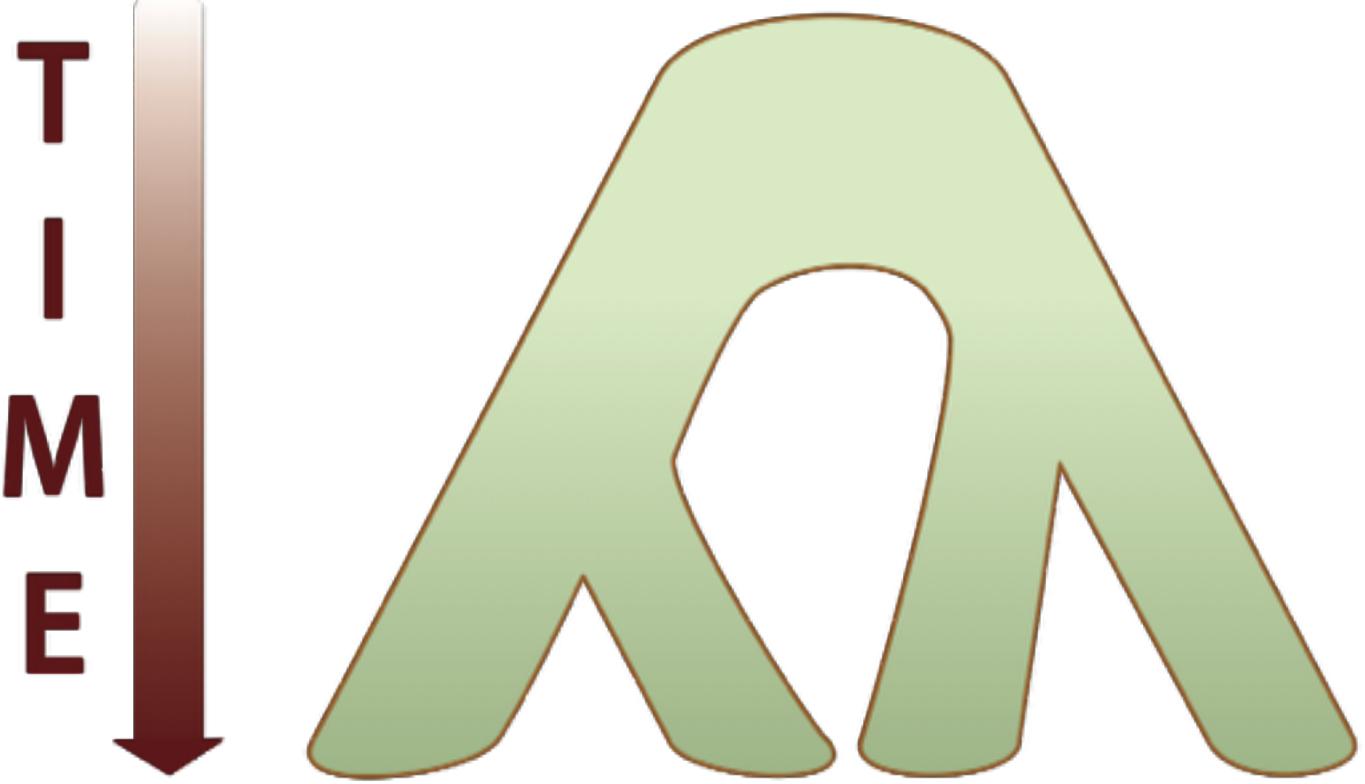
LGT:
exODT

ILS:
MSC

Plan

1. Modeling the relationship between species tree and gene tree
 - coalescent models
 - models of gene duplication and loss
 - models of gene transfer
 - models that combine the above





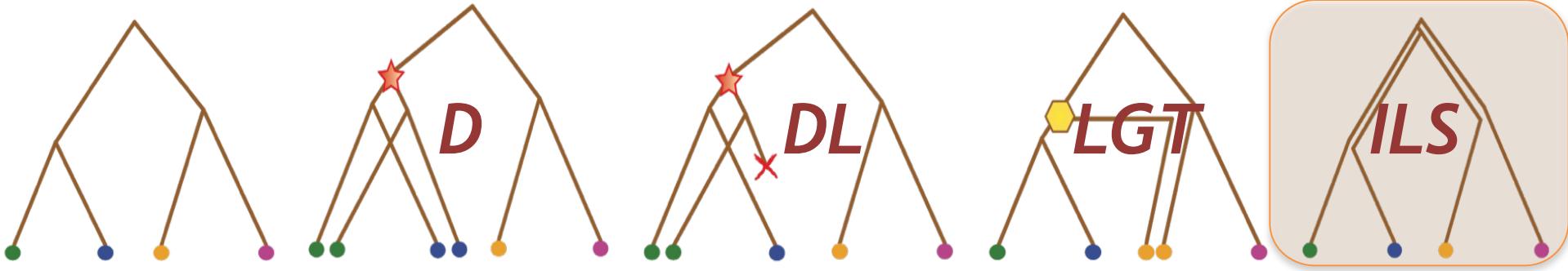
Discrete character:
Continuous character:
Species:

a
0.1
A

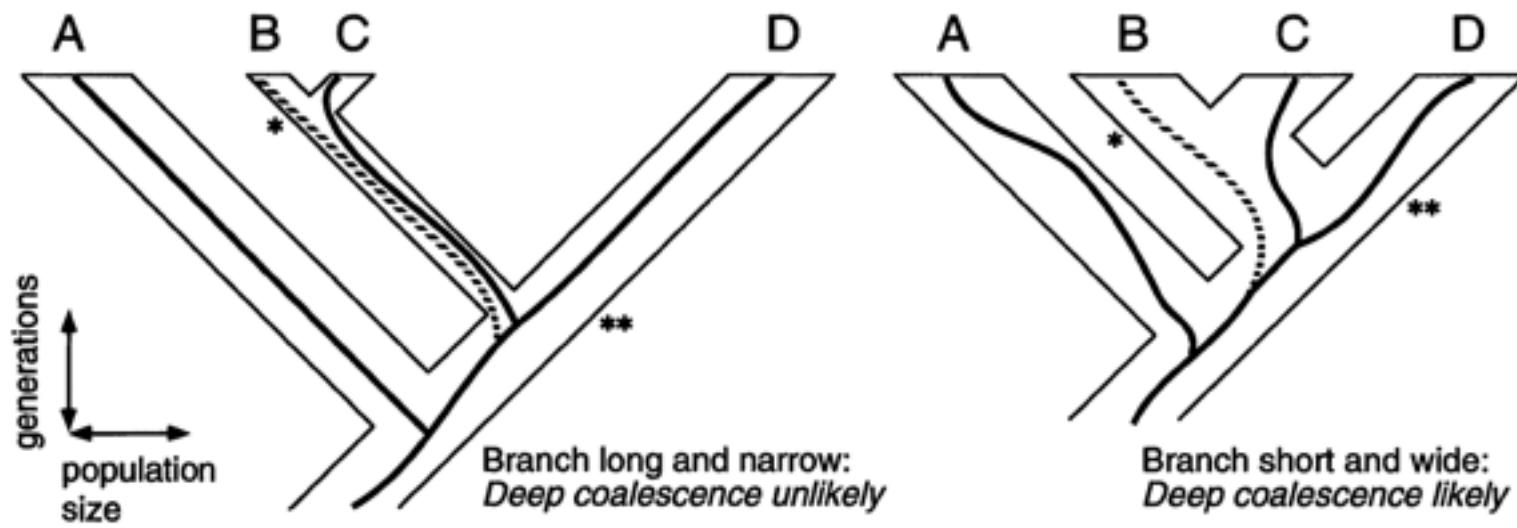
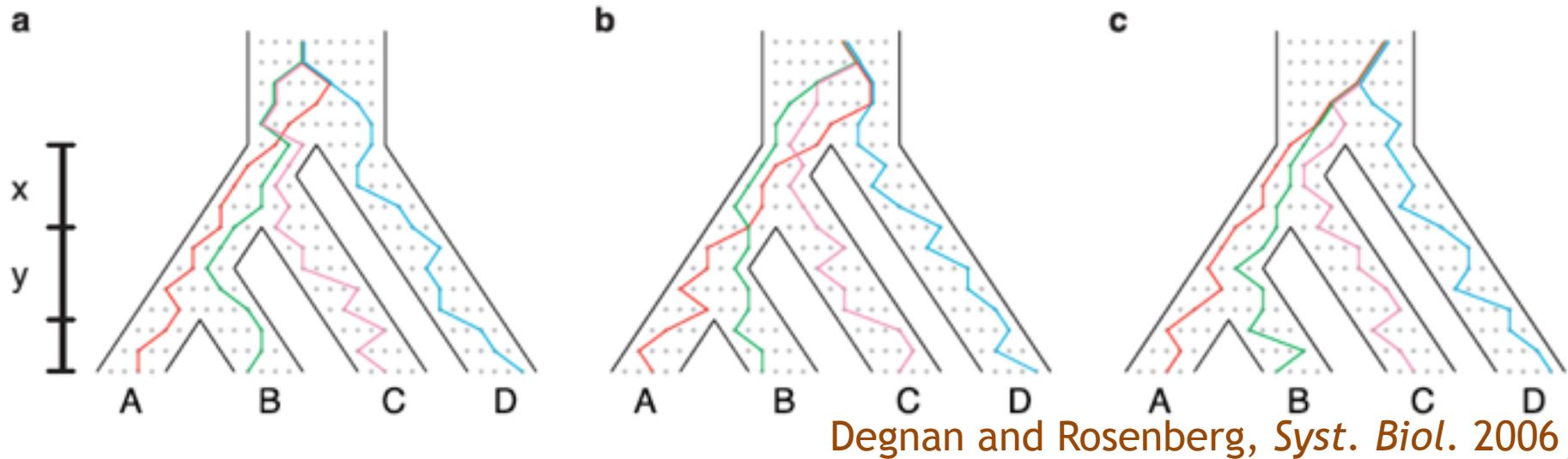
a
0.2
B

b
0.2
C

a
0.4
D



Incomplete lineage sorting

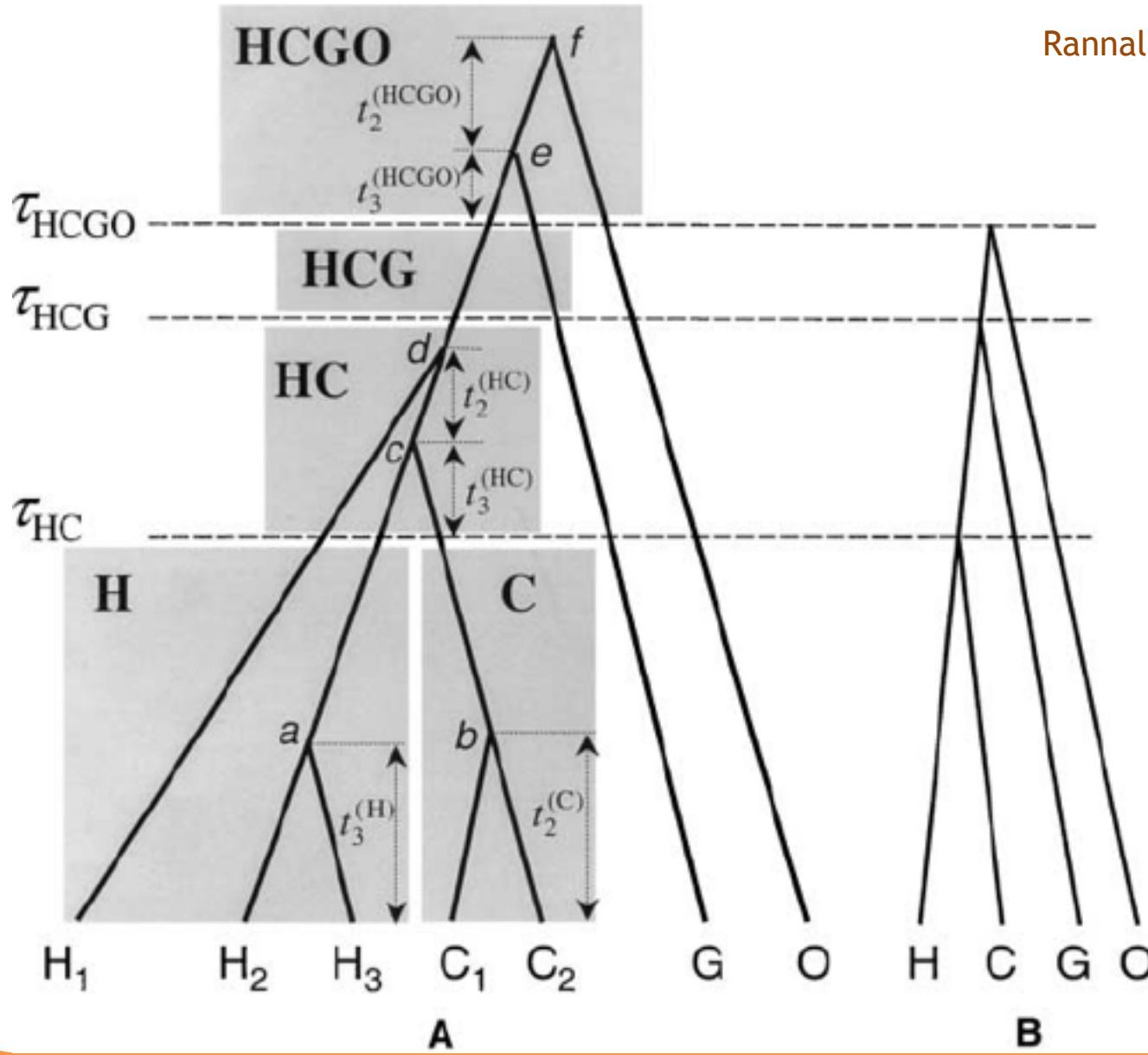


Maddison, *Syst. Biol.* 1997

Reconstructing species trees given ILS

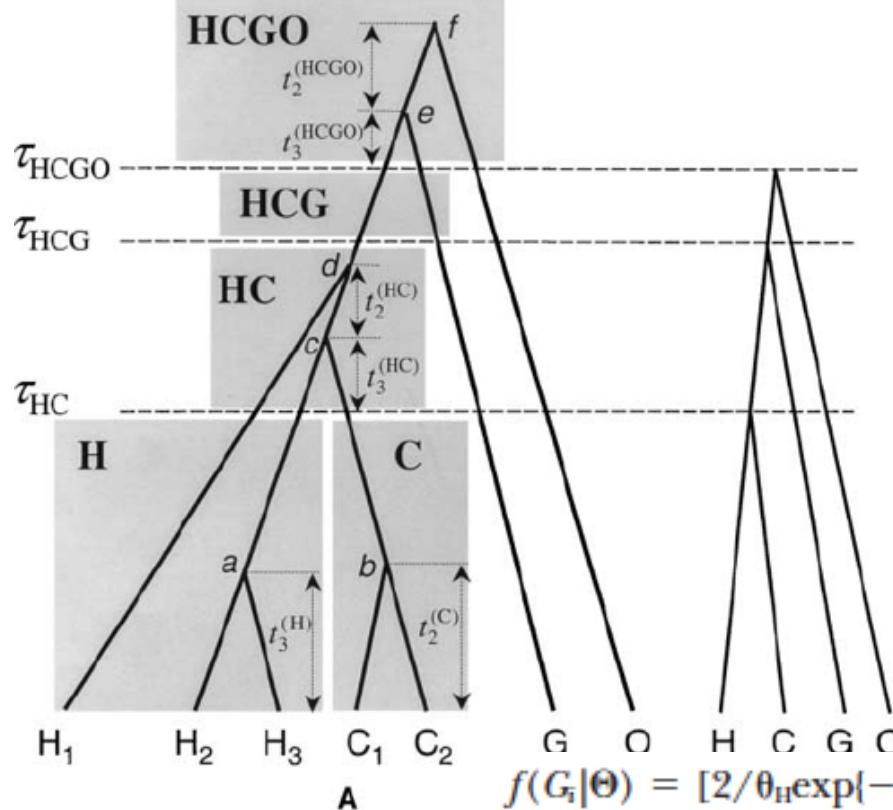
- Parsimony method: minimizing deep coalescences (Maddison 1997)
- Distance method: NJst (Liu and Yu 2011)
- Summary statistic methods: STAR (Liu et al, 2009), GLASS (Mossel and Roch 2010), iGLASS (Jewett and Rosenberg 2012), ASTRAL (Mirarab et al., 2014)
- Maximum Likelihood methods: STEM (Kubatko et al. 2009), MP-EST (Liu et al., 2010)
- Bayesian methods: BeST (2007), *BEAST (Heled and Drummond, 2009)
- Bypassing the gene trees with math tricks: SNAPP (Bryant et al., 2012), POMO (DeMaio et al., 2013), SVDQuartets (Chifman and Kubatko, 2014)

Likelihood of a gene tree given a species tree under the multi-species coalescent



Rannala and Yang, *Genetics* 2003

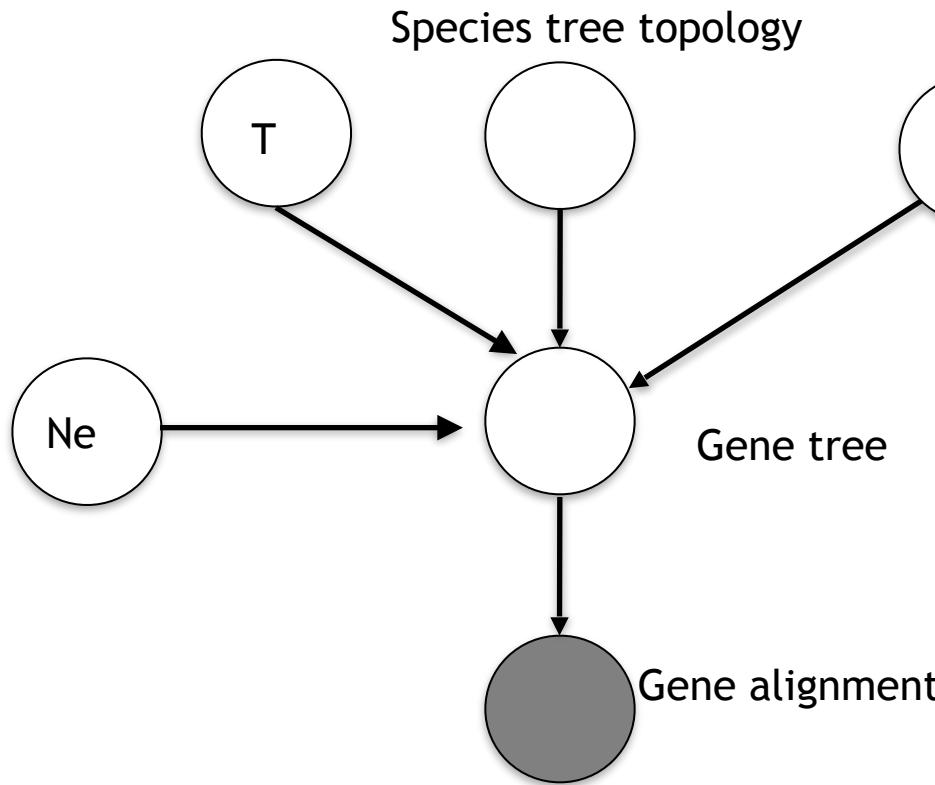
Likelihood of a gene tree given a species tree under the multi-species coalescent



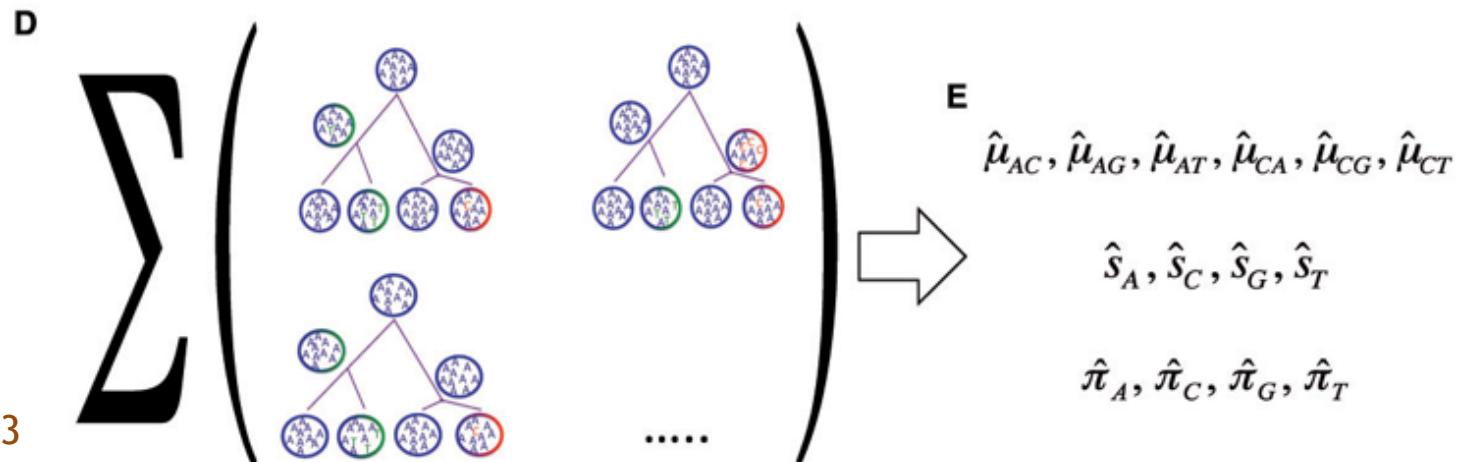
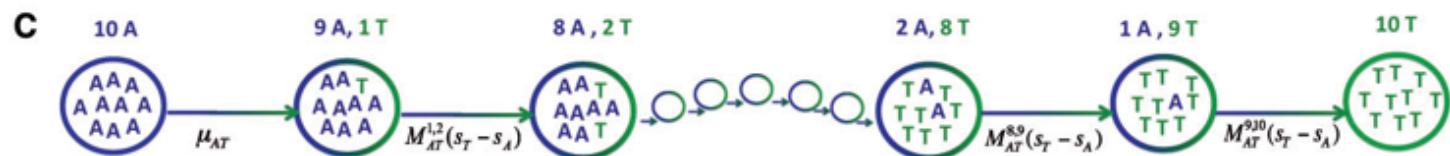
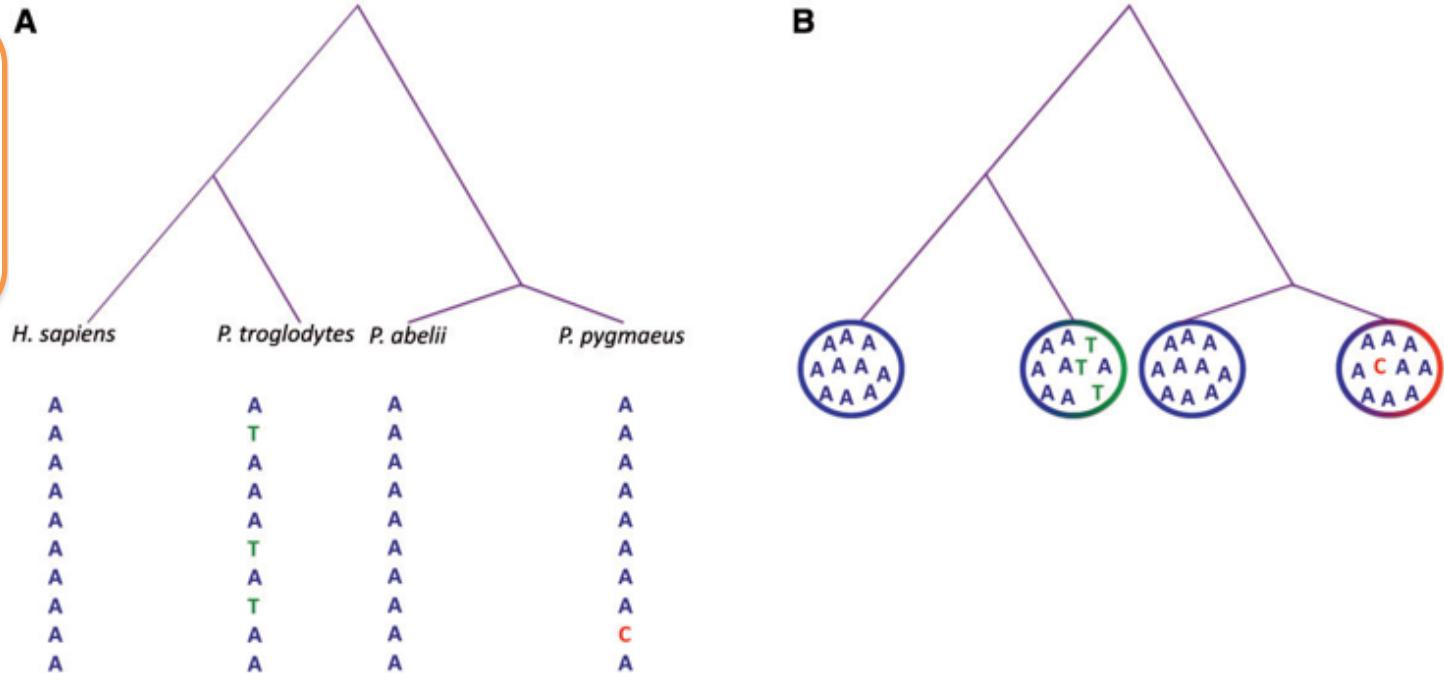
Rannala and Yang, *Genetics* 2003

$$\begin{aligned}
 f(G_i | \Theta) = & [2/\theta_H \exp\{-6t_3^{(H)}/\theta_H\} \exp\{-2(\tau_{HC} - t_3^{(H)})/\theta_H\}] \\
 & \times [2/\theta_C \exp\{-2t_2^{(C)}/\theta_C\}] \\
 & \times [2/\theta_{HC} \exp\{-6t_3^{(HC)}/\theta_{HC}\} \times 2/\theta_{HC} \exp\{-2t_2^{(HC)}/\theta_{HC}\}] \\
 & \times [\exp\{-2(\tau_{HCG} - \tau_{HC} - (t_3^{(HC)} + t_2^{(HC)}))/\theta_{HCG}\}] \\
 & \times [2/\theta_{HCGO} \exp\{-6t_3^{(HCGO)}/\theta_{HCGO}\} \times 2/\theta_{HCGO} \exp\{-2t_2^{(HCGO)}/\theta_{HCGO}\}].
 \end{aligned}$$

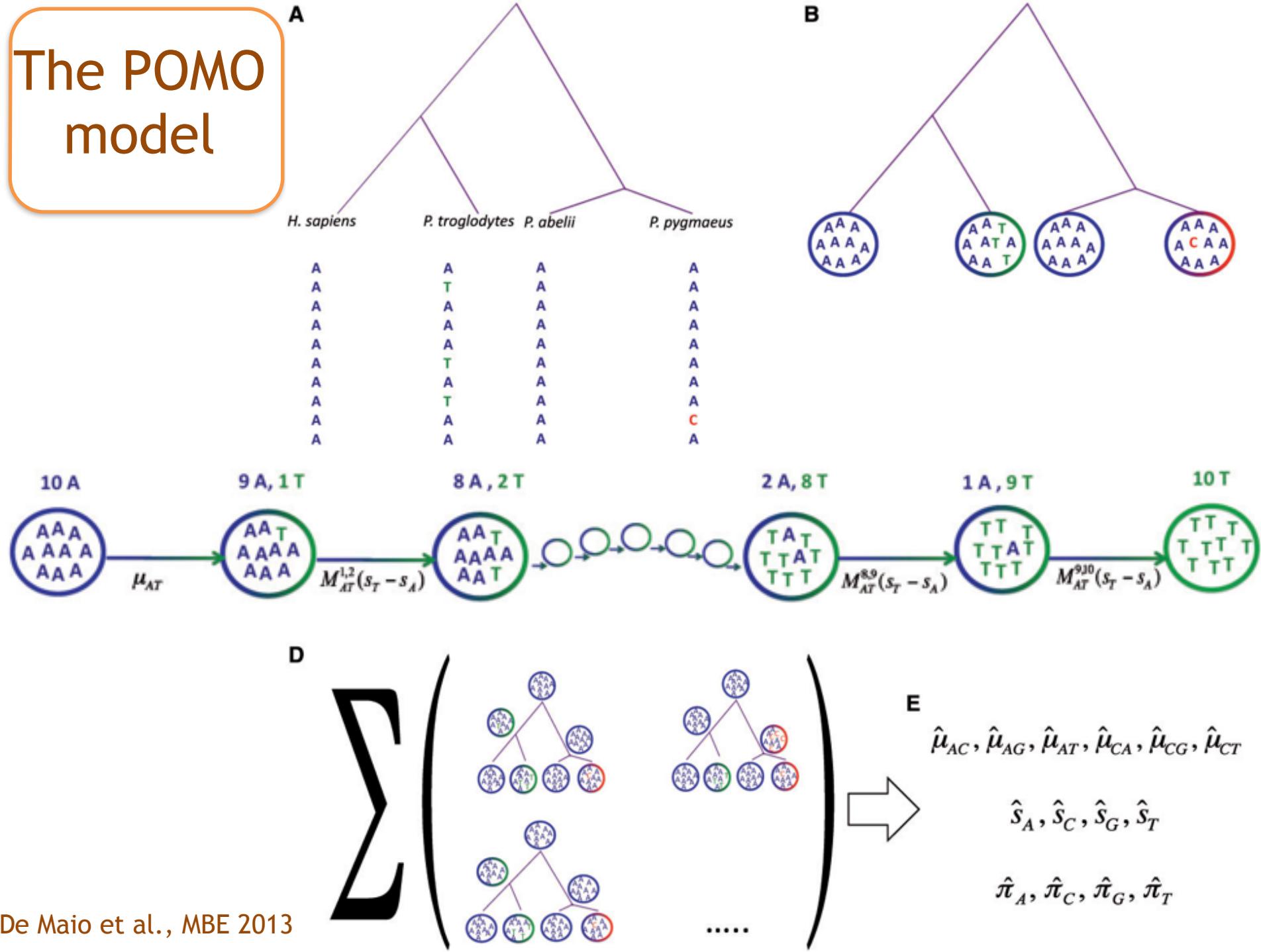
The species tree-gene tree graphical model with ILS



The POMO model



The POMO model



The POMO model

$$\mathbf{A} \left(\begin{array}{c} \binom{1}{N-1} \\ \binom{2}{N-2} \\ \binom{3}{N-3} \\ \binom{4}{N-4} \\ \vdots \\ \binom{N-4}{1} \\ \binom{N-3}{2} \\ \binom{N-2}{3} \\ \binom{N-1}{2} \\ \binom{1}{1} \end{array} \right) \left(\begin{array}{ccccccccc} \binom{1}{N-1} & \binom{2}{N-2} & \binom{3}{N-3} & \binom{4}{N-4} & \cdots & \binom{N-4}{4} & \binom{N-3}{3} & \binom{N-2}{2} & \binom{N-1}{1} \\ * & N \cdot M_{IJ}^{1,2} & * & N \cdot M_{IJ}^{2,3} & & & & & \\ N \cdot M_{IJ}^{2,1} & * & N \cdot M_{IJ}^{3,2} & * & N \cdot M_{IJ}^{3,4} & & & & 0 \\ & N \cdot M_{IJ}^{3,1} & N \cdot M_{IJ}^{4,2} & * & * & & & & \\ & & N \cdot M_{IJ}^{4,3} & * & * & & & & \\ & & & * & * & & & & \\ & & & & * & N \cdot M_{IJ}^{N-3,N-2} & N \cdot M_{IJ}^{N-2,N-3} & N \cdot M_{IJ}^{N-2,N-1} & \\ & & & & N \cdot M_{IJ}^{N-1,N-2} & * & N \cdot M_{IJ}^{N-1,N-1} & * & N \cdot M_{IJ}^{N-1,N} \\ & & & & & N \cdot M_{IJ}^{N,N-1} & & & * \end{array} \right)$$

$$\mathbf{B}_A \left(\begin{array}{c} A \\ C \\ G \\ T \\ \vdots \\ \binom{N-1}{1} \\ \binom{N-2}{1} \\ \binom{N-3}{1} \\ \binom{N-4}{1} \\ \vdots \\ \binom{N-1}{T} \\ \binom{N-2}{T} \\ \binom{N-3}{T} \\ \binom{N-4}{T} \\ \vdots \\ \binom{N-1}{1} \end{array} \right) \left(\begin{array}{cccc|cccc|cccc|cccc|} A & C & G & T & \binom{1}{N-1} & \cdots & \binom{N-1}{N-1} & A & C & G & T & \binom{1}{N-1} & \cdots & \binom{N-1}{N-1} & A & C & G & T \\ * & 0 & 0 & 0 & 0 & \cdots & N^2 \cdot \mu_{AC} & 0 & 0 & 0 & 0 & N^2 \cdot \mu_{AG} & 0 & 0 & 0 & 0 & N^2 \cdot \mu_{AT} & 0 & \cdots & \cdots \\ 0 & * & 0 & 0 & 0 & \cdots & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \cdots \\ 0 & 0 & * & 0 & 0 & \cdots & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \cdots \\ 0 & 0 & 0 & * & 0 & \cdots & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \cdots \\ N \cdot M_{AC}^{1,0} & 0 & 0 & 0 & \hline & & & & B_{AC}^N & & & 0 & & & & 0 & & & 0 & & \cdots \\ 0 & N \cdot M_{AC}^{N,N+1} & 0 & 0 & \hline & & & & 0 & & & B_{AG}^N & & & 0 & & & 0 & & \cdots \\ N \cdot M_{AG}^{1,0} & 0 & 0 & 0 & \hline & & & & 0 & & & 0 & & & B_{AT}^N & & & 0 & & & \cdots \\ 0 & 0 & N \cdot M_{AG}^{N,N+1} & 0 & \hline & & & & 0 & & & 0 & & & 0 & & & B_{AT}^N & & \cdots \\ N \cdot M_{AT}^{1,0} & 0 & 0 & 0 & \hline & & & & 0 & & & 0 & & & \vdots & & & \vdots & & & \vdots & & \ddots \end{array} \right)$$

Pros and cons of the Pomo model

1. Pros:

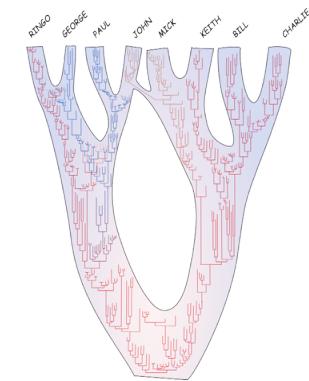
1. bypasses gene tree estimation (easy to combine with DL or DTL models...)
2. does not assume linkage between sites

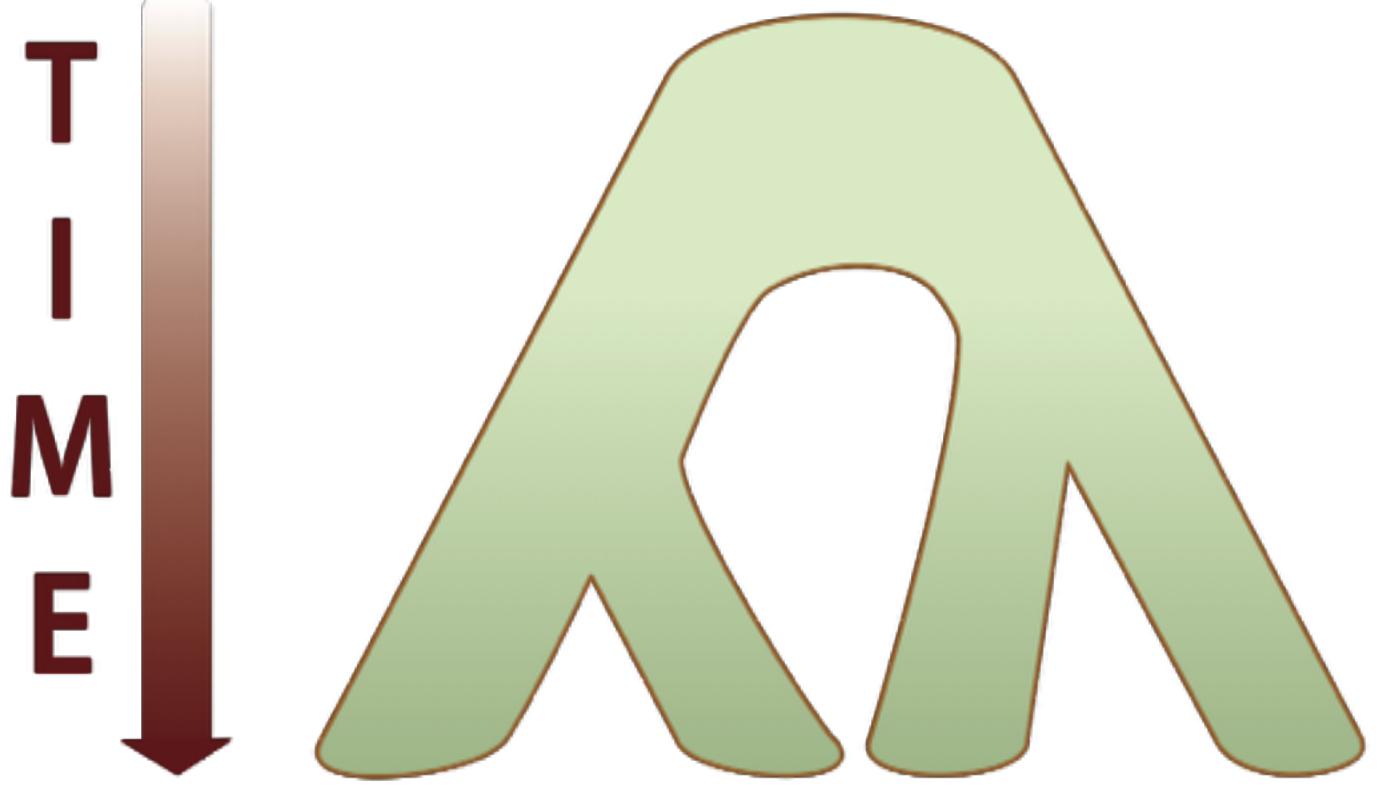
2. Cons:

1. You don't get gene trees
2. It's not sure what branch lengths mean anymore
3. matrix can become big
4. Only bimorphic sites are considered: 3- or 4-morphic sites need to be handled differently
5. Not sure how it behaves in practice

Plan

1. Modeling the relationship between species tree and gene tree
 - coalescent models
 - models of gene duplication and loss
 - models of gene transfer
 - models that combine the above





Discrete character:
Continuous character:
Species:



a

0.1

A



a

0.2

B



b

0.2

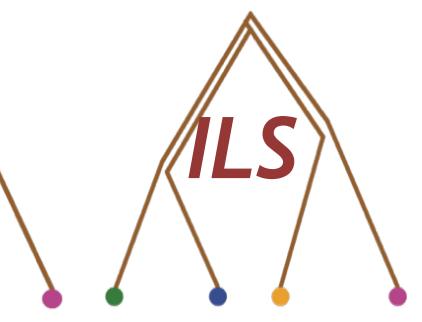
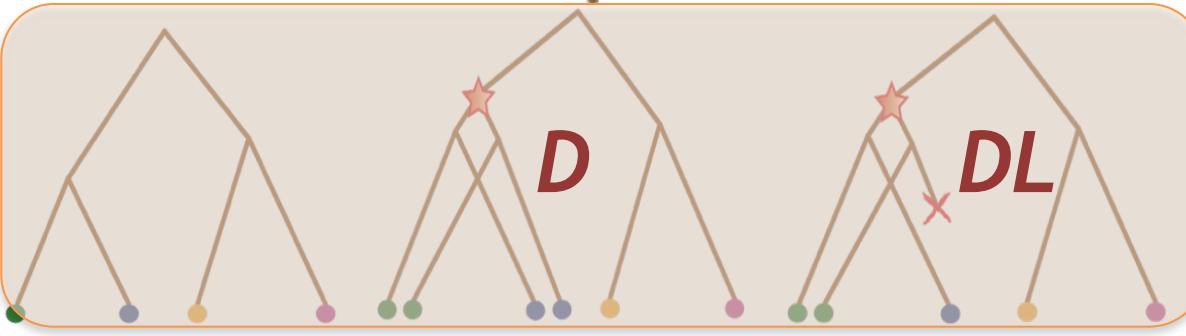
C



a

0.4

D

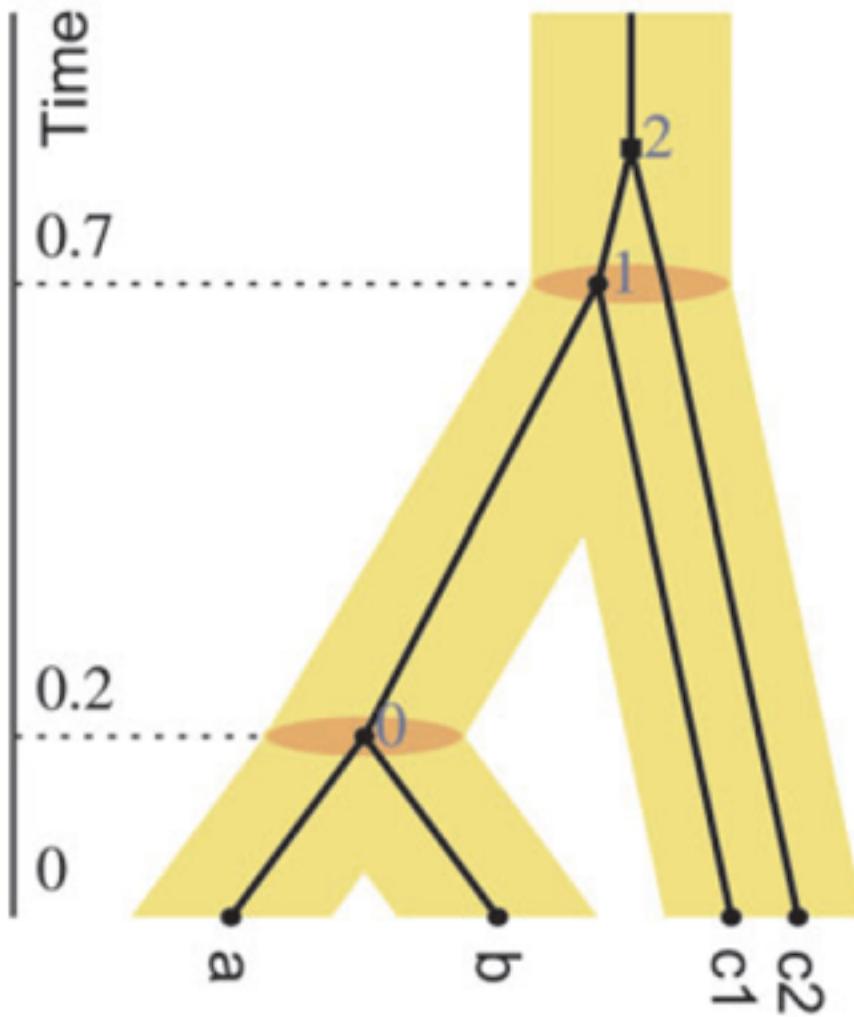


DL models

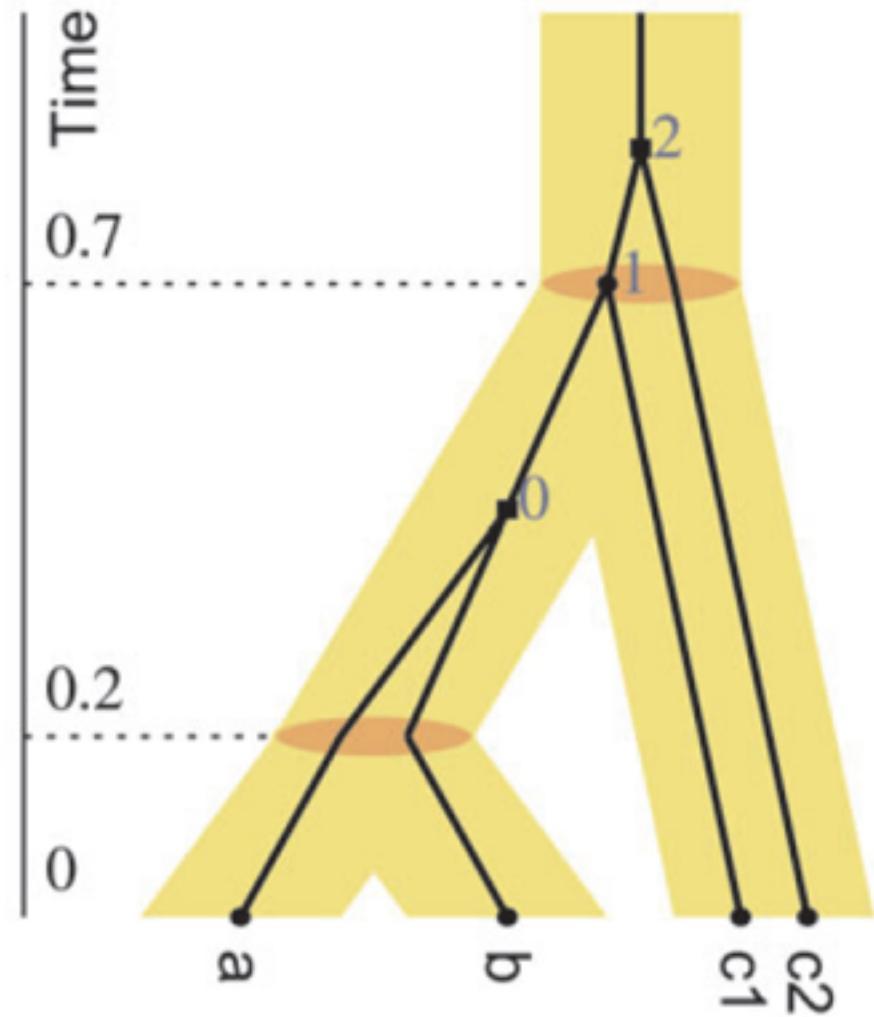
- Reconciliation between a gene tree and a species tree:
 - Parsimony method: minimizing numbers of duplications (Goodman et al. 1979; Page 1994; Zmasek and Eddy 2001)
 - Model-based method: Arvestad et al. 2003, 2004, 2009; Görecki et al. 2011
- Reconstruction of a gene tree given a species tree and alignment:
 - Parsimony method: Chen et al. 2000; Durand et al. 2006
 - Model-based method: Rasmussen and Kellis 2007, 2011; Akerborg et al. 2009; Sjöstrand et al. 2012
- Reconstruction of a species tree given gene trees:
 - Parsimony methods: Page and Charleston 1997, Bansal et al. 2007; Wehe et al. 2008; Bansal et al. 2010; Chang et al. 2011
- Joint reconstruction of gene and species trees given alignments:
 - Model-based method: Boussau et al., 2013

Many possible reconciliations

a)

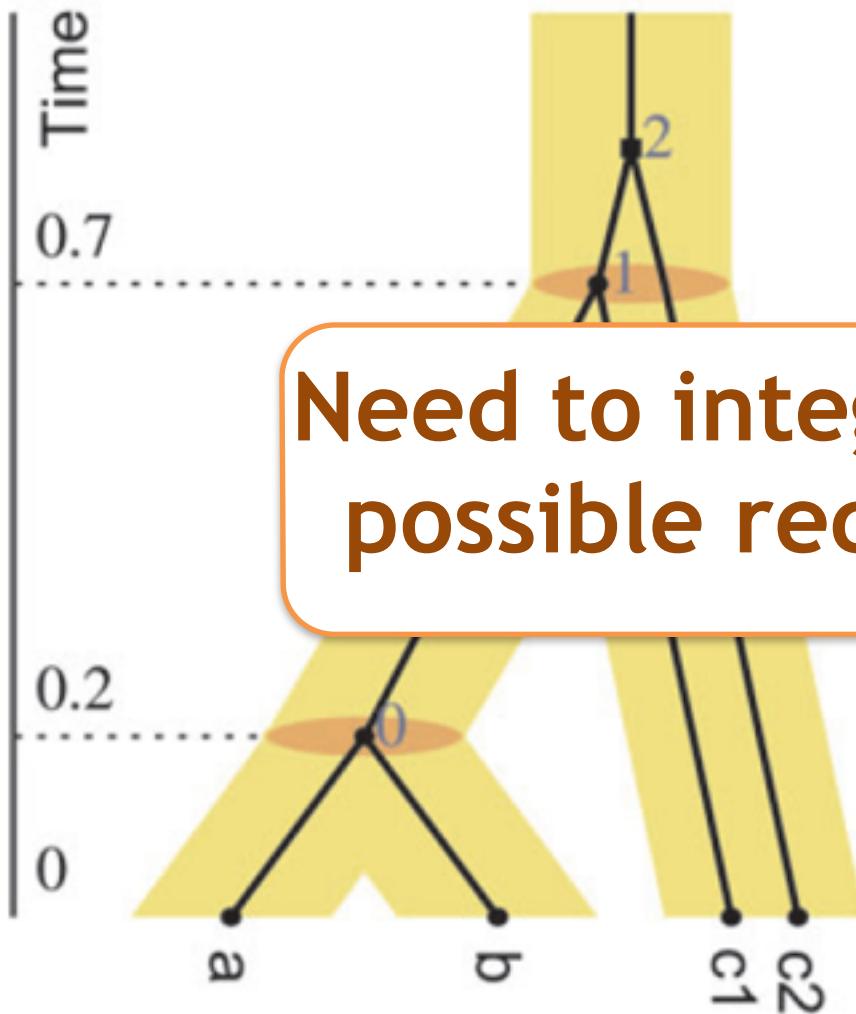


b)

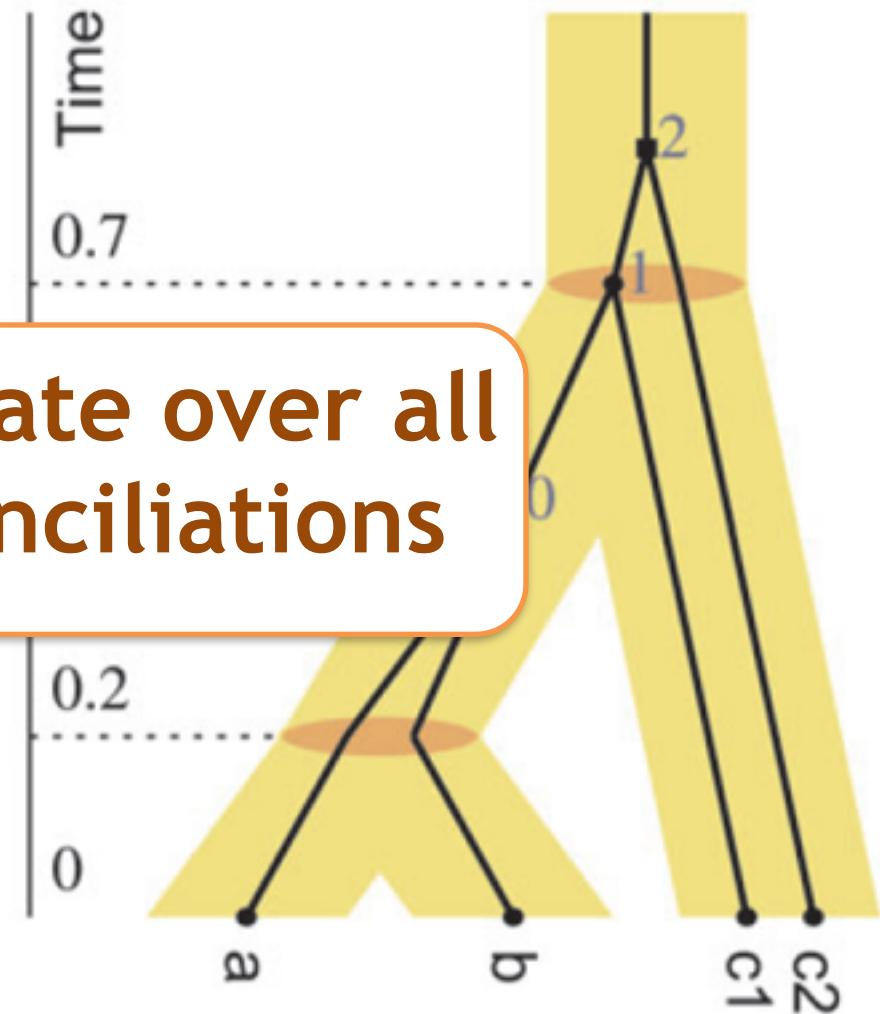


Many possible reconciliations

a)

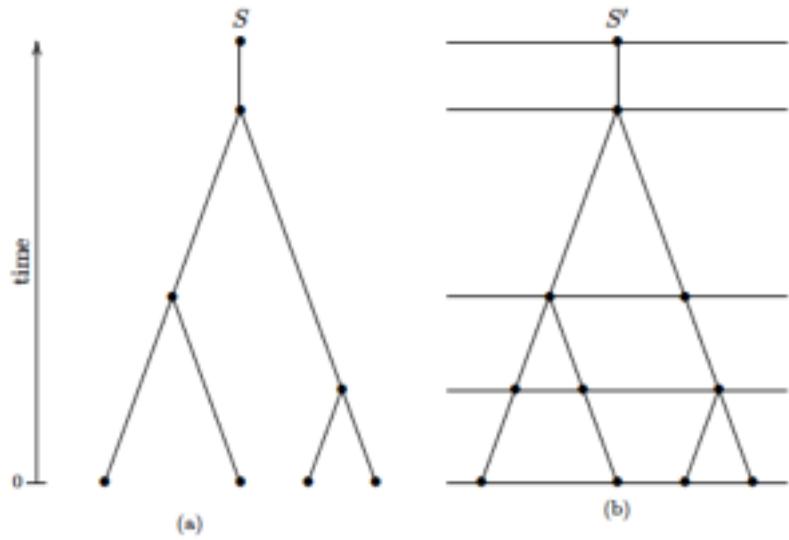


b)

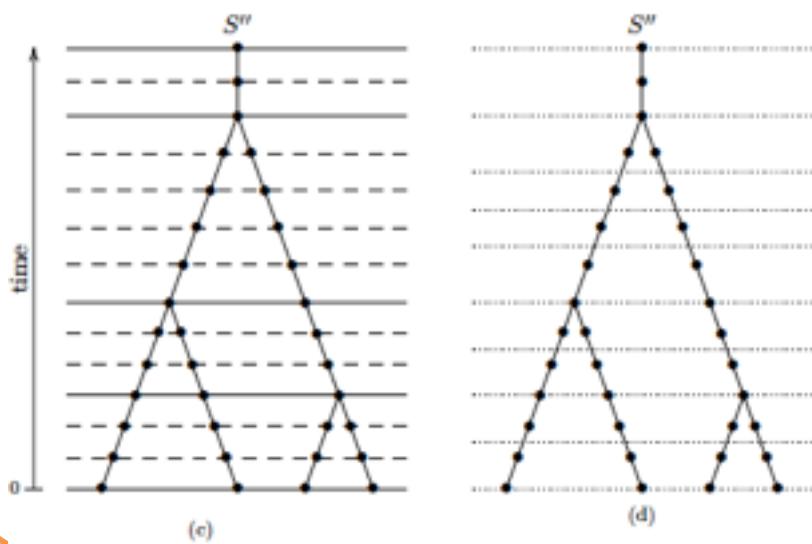


Need to integrate over all possible reconciliations

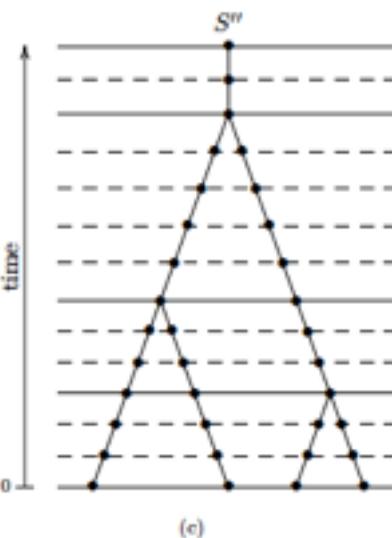
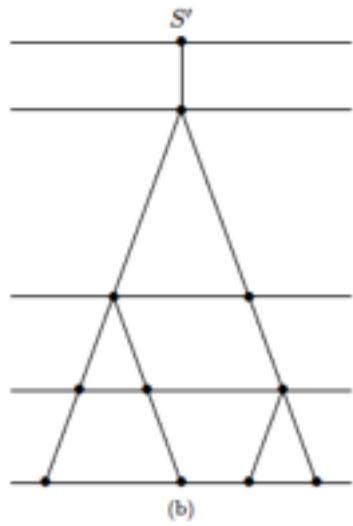
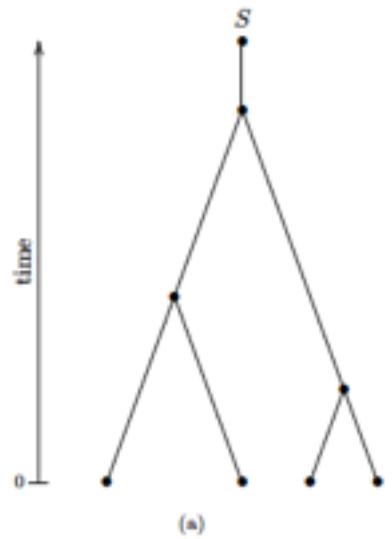
Considering all mappings: The species tree discretization approximation



Tofiq, *PhD thesis*



Considering all mappings: The species tree discretization approximation



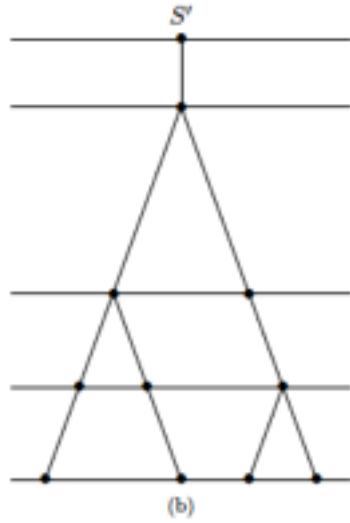
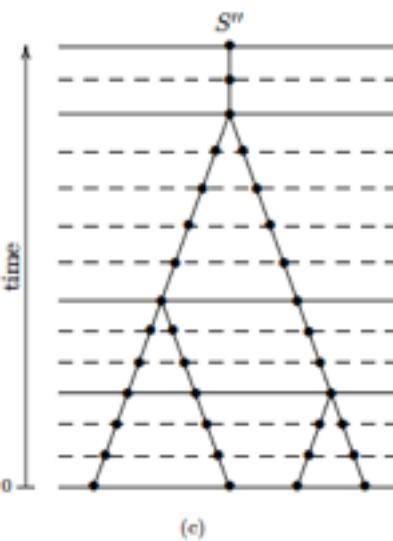
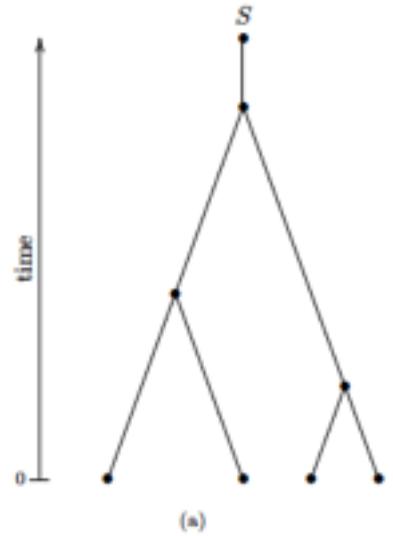
Tofigh, *PhD thesis*

Mapping between a gene tree
and species tree



Mapping between nodes of
the gene tree and nodes of
the augmented species tree

Considering all mappings: The species tree discretization approximation



Tofigh, PhD thesis

Mapping between a gene tree
and species tree



Mapping between nodes of
the gene tree and nodes of
the augmented species tree

Integrating over all mappings:
done through dynamic
programming from the leaves
of the gene tree to its root

Birth-death models for DL

Model-based

approaches:

often use a **birth-death**

process (Arvestad et al.,

2003, 2004, 2009; Akerborg et

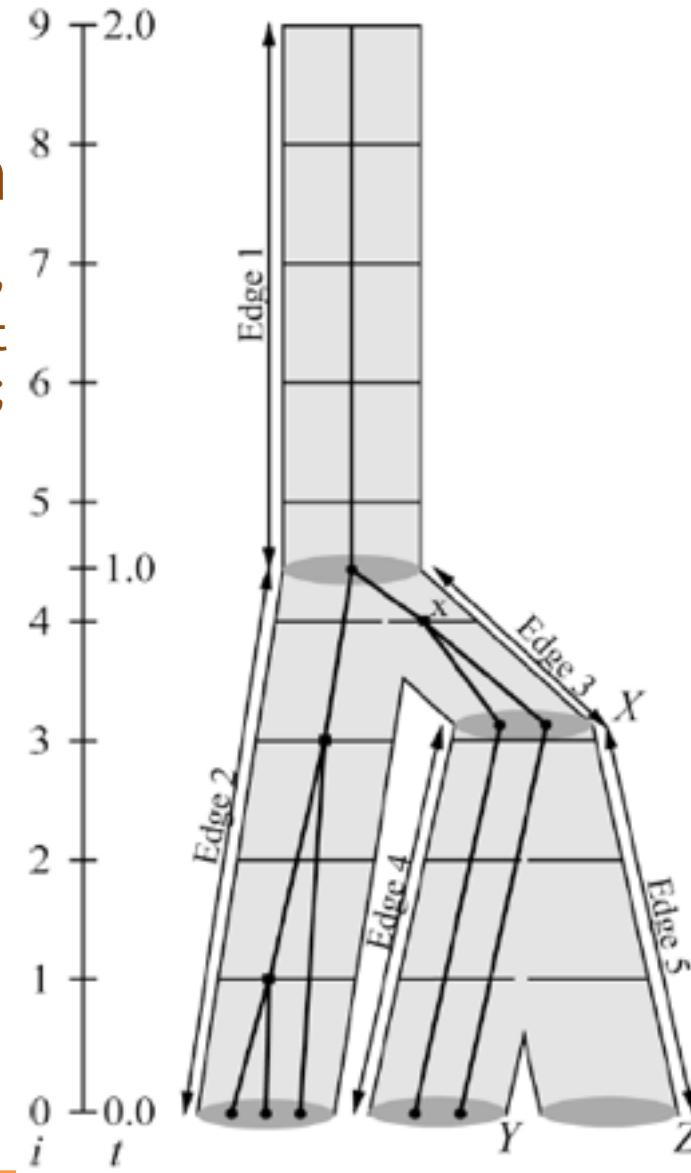
al., 2009; Sjöstrand et al., 2012;

Boussau et al., 2013)

Akerborg et al., PNAS 2009

Birth=Duplication: λ

Death=Loss: μ



The Akerborg et al. 2009 model

- Input: Species tree with time, rooted gene tree, λ and μ assumed to be known
- Output: reconciliation between species tree and gene tree
- Relaxed-clock model
- Gene trees generated according to a birth-death model

$$\Pr [D, G|S] = \int \Pr [D|G, l = rt] p[r|G] p[G, t|S] dt dr$$

The Akerborg et al. 2009 model

- Input: Species tree with time, rooted gene tree, λ and μ assumed to be known
- Output: reconciliation between species tree and gene tree
- Relaxed-clock model
- Gene trees generated according to a birth-death model

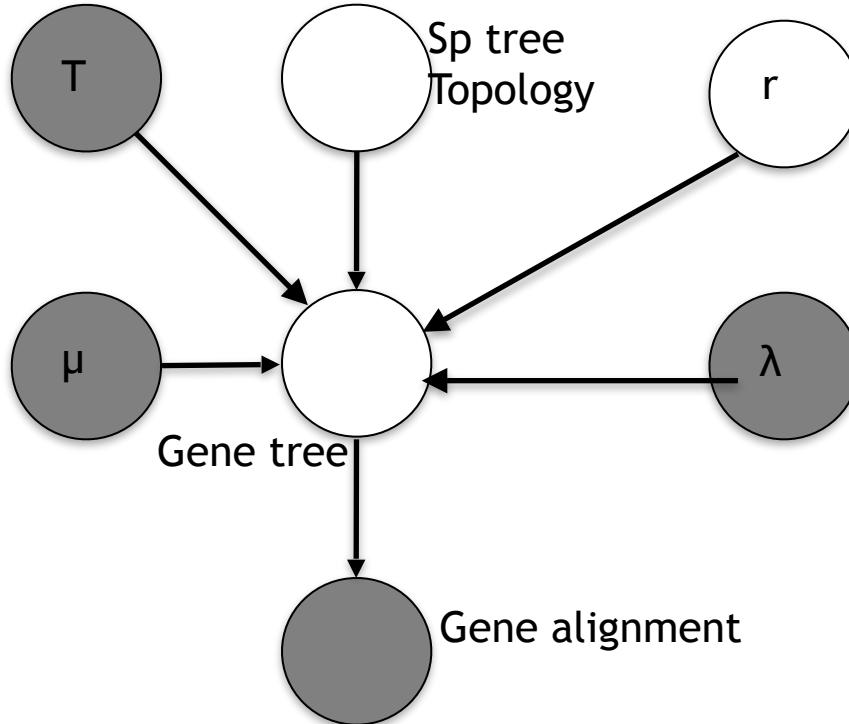
$$\Pr [D, G|S] = \int \Pr [D|G, l = rt] p[r|G] p[G, t|S] dt dr$$

Felsenstein 1981

Gamma prior

Dynamic
programming
to integrate
over all
mappings with
BD model

The species tree-gene tree graphical model with DL (Akerborg 2009)



The Akerborg et al. 2009 model vs PHYLDOG

- Akerborg et al. 2009:
 - Generative model
 - Integrates over a large space of mappings
 - Integrates over all possible scenarios of duplications and losses

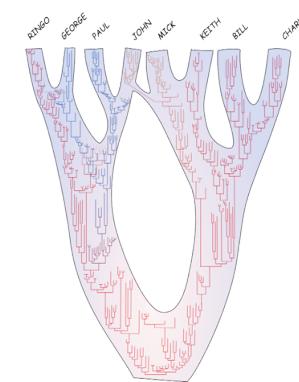
PHYLDOG

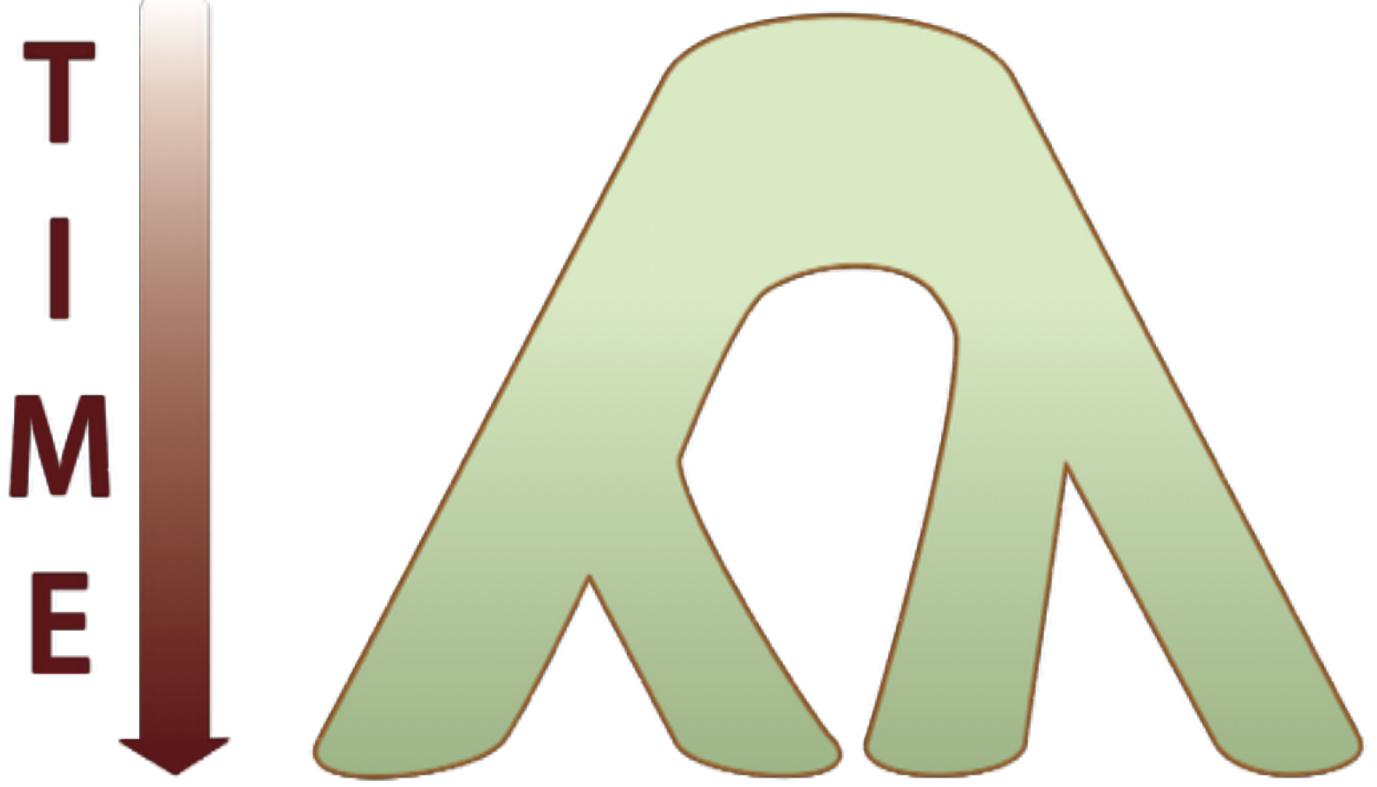
- use unrooted, non-clock gene trees
- Use the most parsimonious mapping
- Use a double-recursive tree traversal
- Integrates over a large subset of possible scenarios of duplications and losses (Boussau et al., Genome Research 2013)

Plan

1. Modeling the relationship between species tree and gene tree

- coalescent models
- models of gene duplication and loss
- models of gene transfer
- models that combine the above





Discrete character:
Continuous character:
Species:



a

0.1

A



a

0.2

B



b

0.2

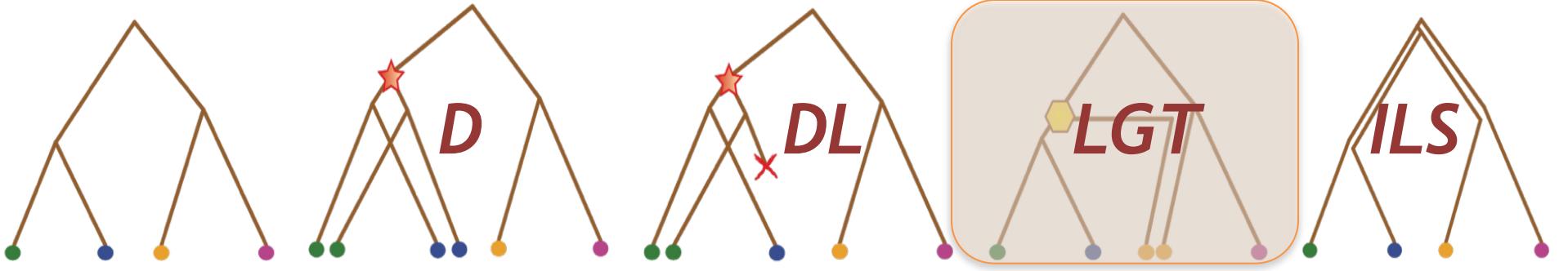
C



a

0.4

D



Pure transfer models

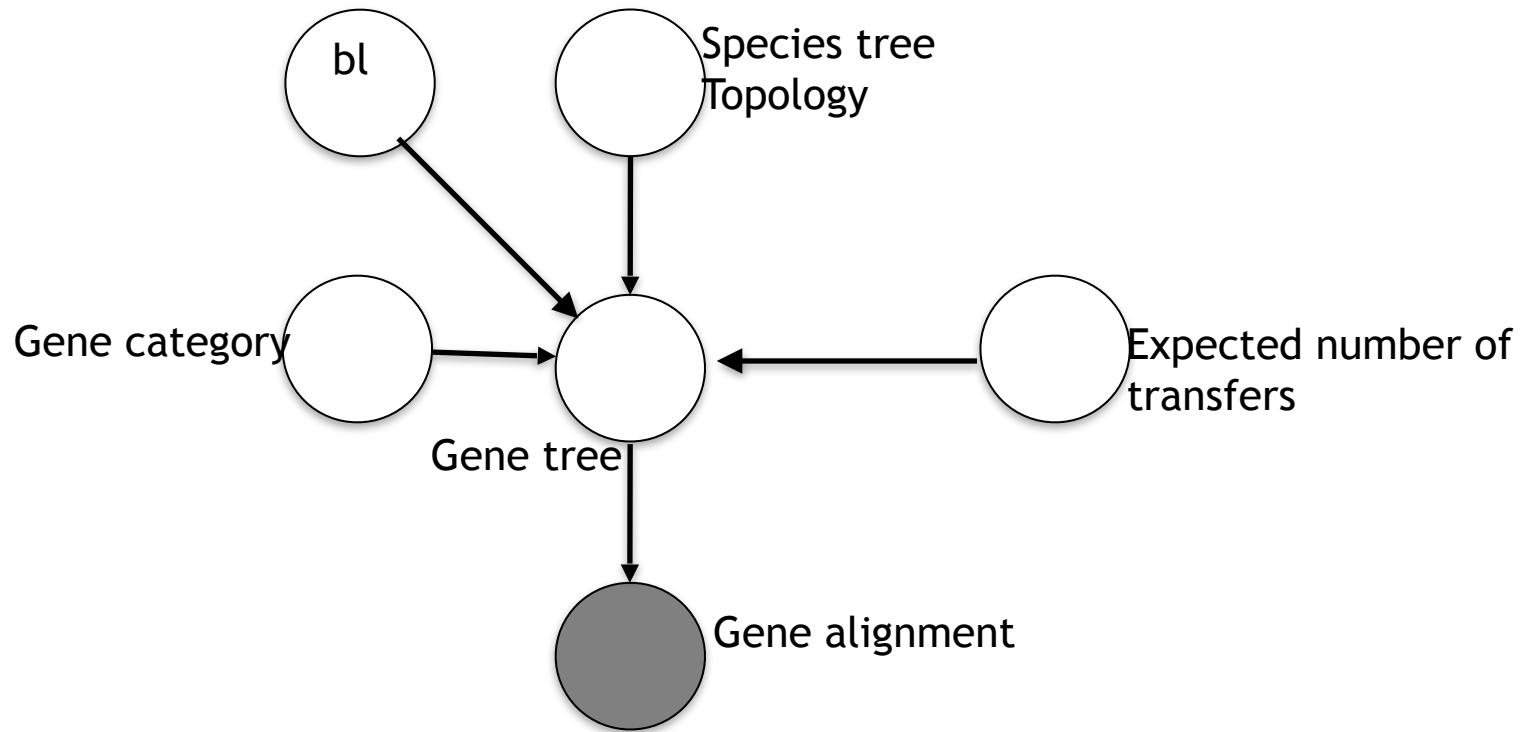
- Identification of transfers (not probabilistic):
 - Nakleh et al., 2005; Beiko and Hamilton 2006; Abby et al., 2010; Hill et al. 2010
- Joint reconstruction of gene and species trees given alignments:
 - Model-based method: Suchard, *Genetics* 2005

Suchard “random walk through tree space” model

- Input: gene alignments
- Output: species tree, gene trees
- Each gene tree is a Poisson number of transfers away from the species tree
- Random walk through gene trees
- C classes of genes, with different propensities to be transferred
- Has been used with at most 6 species, 144 orthologous gene alignments, but could be used with 8 species or even more

Suchard, *Genetics* 2005

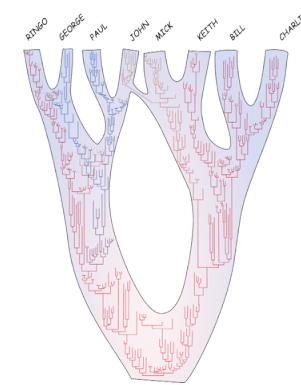
The species tree-gene tree graphical model with transfers (Suchard 2005)



Plan

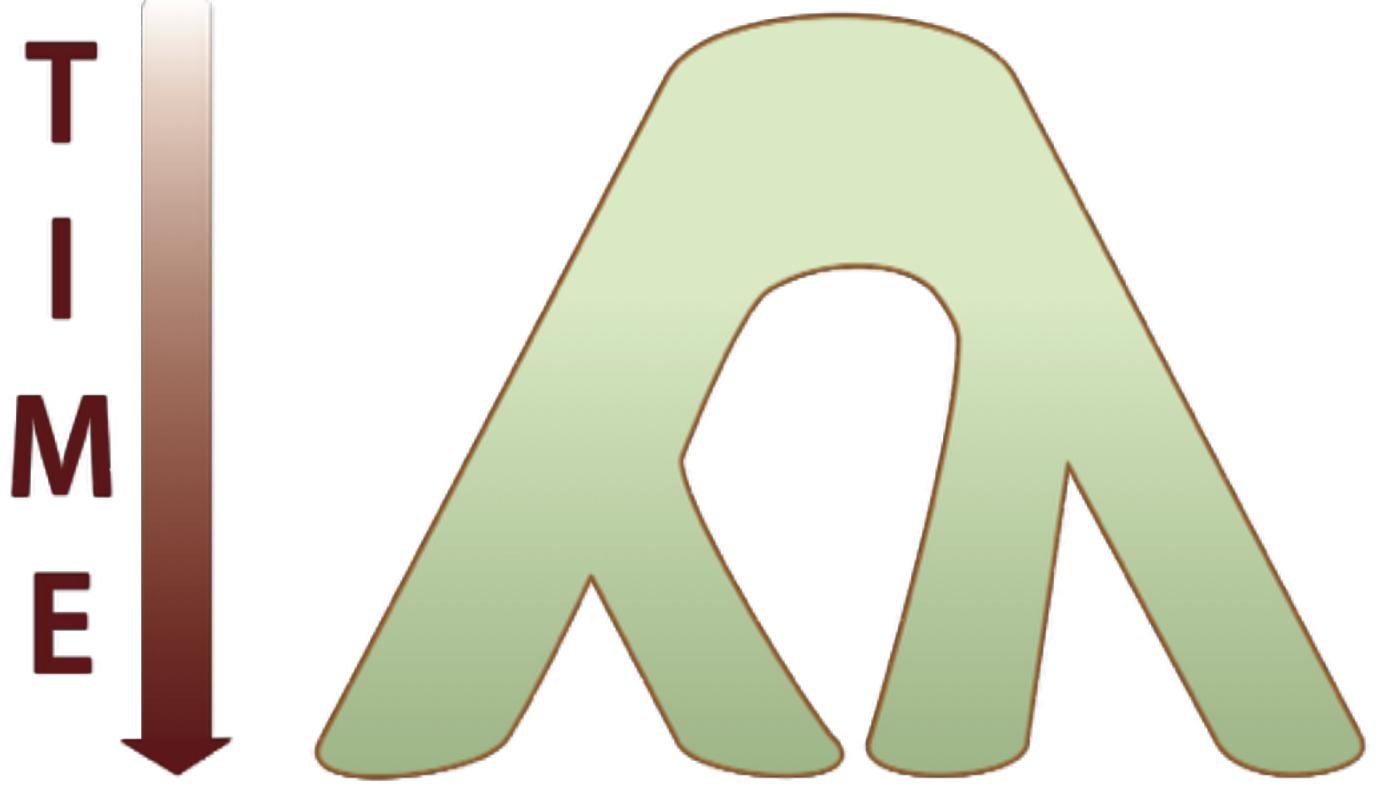
1. Modeling the relationship between species tree and gene tree

- coalescent models
- models of gene duplication and loss
- models of gene transfer
- models that combine the above



Model that combines DL and ILS

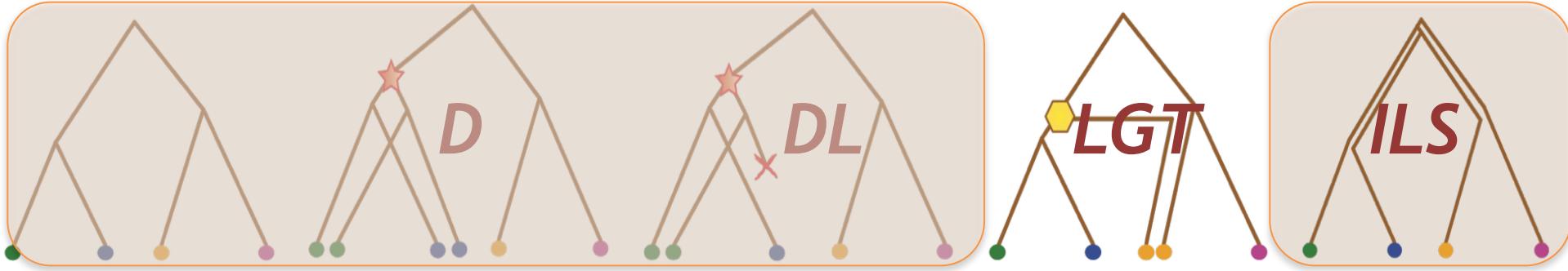
- Reconciliation between gene tree and species tree:
 - Model-based method: Rasmussen and Kellis, *Genome Research* 2012
- Input: dated and rooted species tree, DL rates, Ne
- Output: Locus tree and DL+ILS scenario



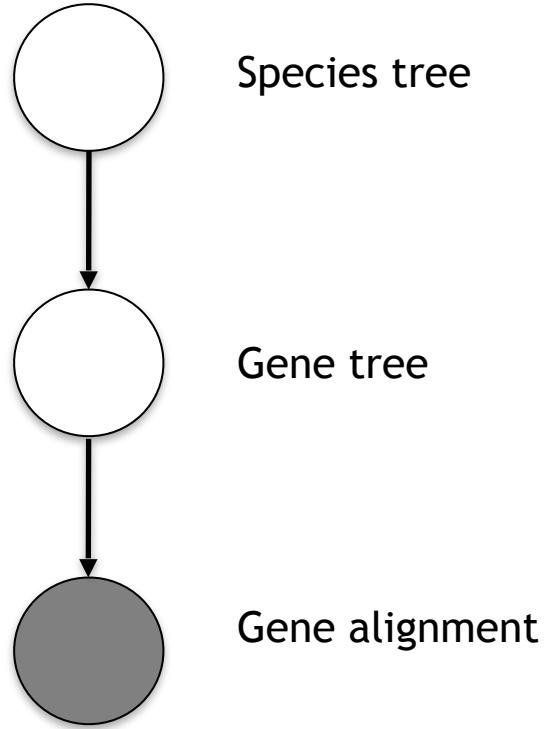
Discrete character: a
 Continuous character: 0.1
 Species: A

a 0.2
 b 0.2
 Species: B C

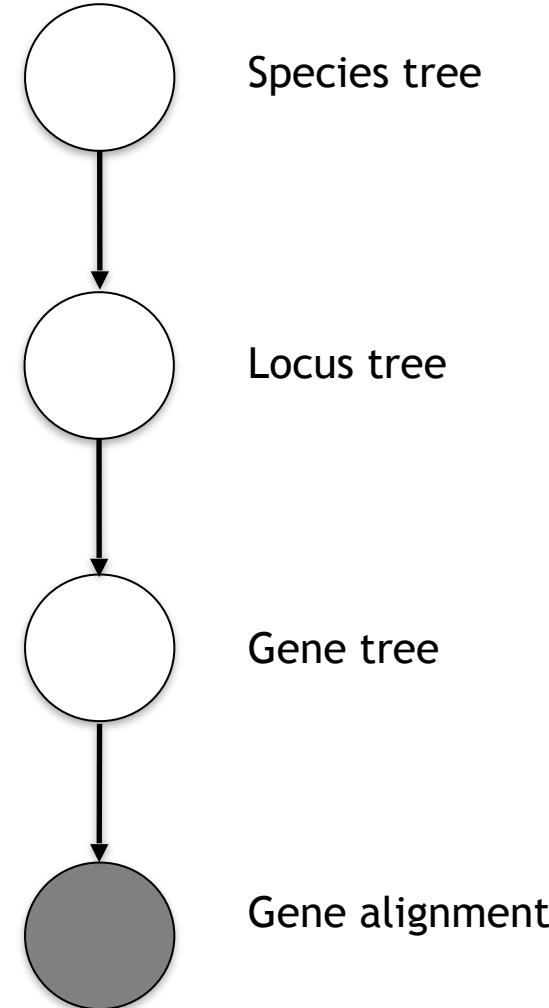
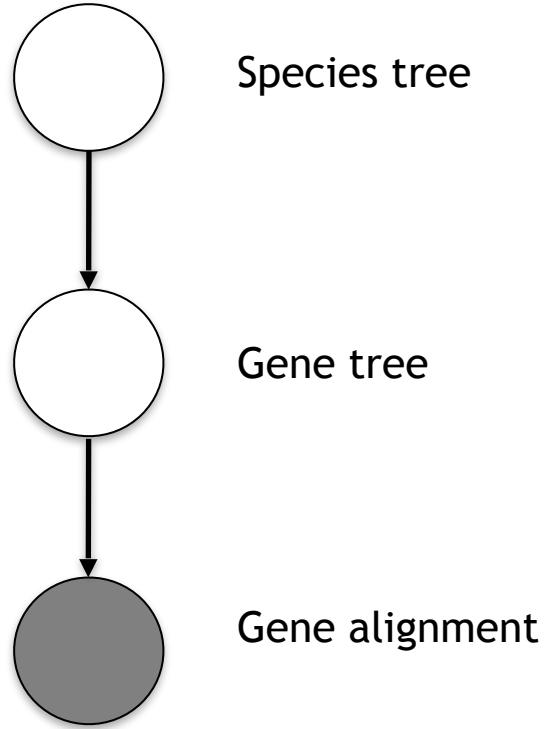
a 0.4
 Species: D



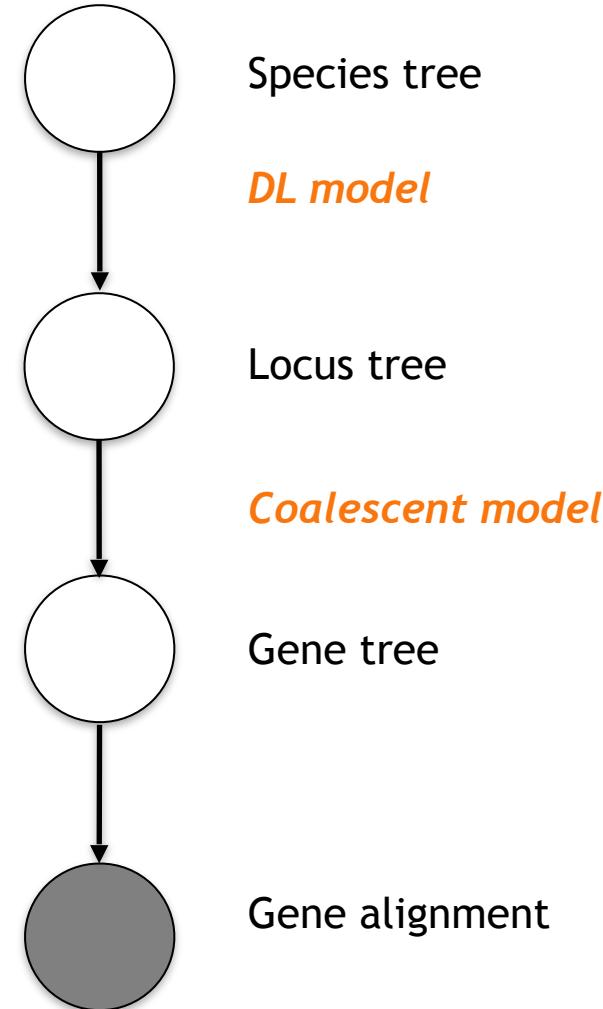
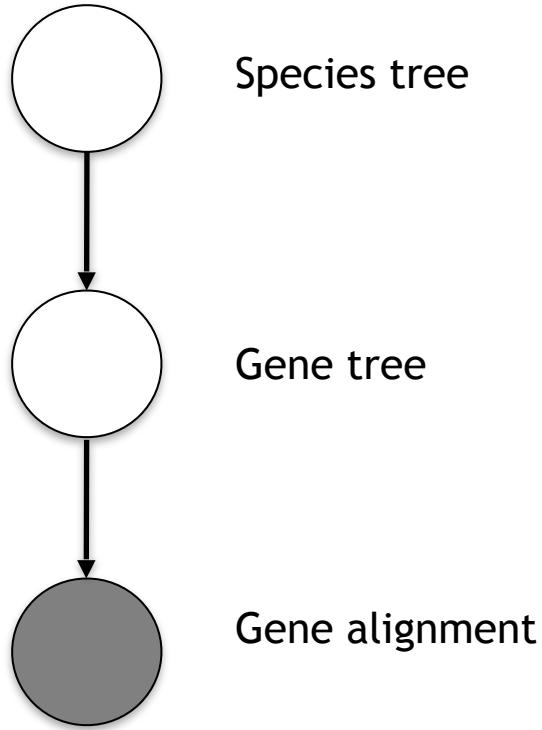
The Rasmussen and Kellis model (2012)



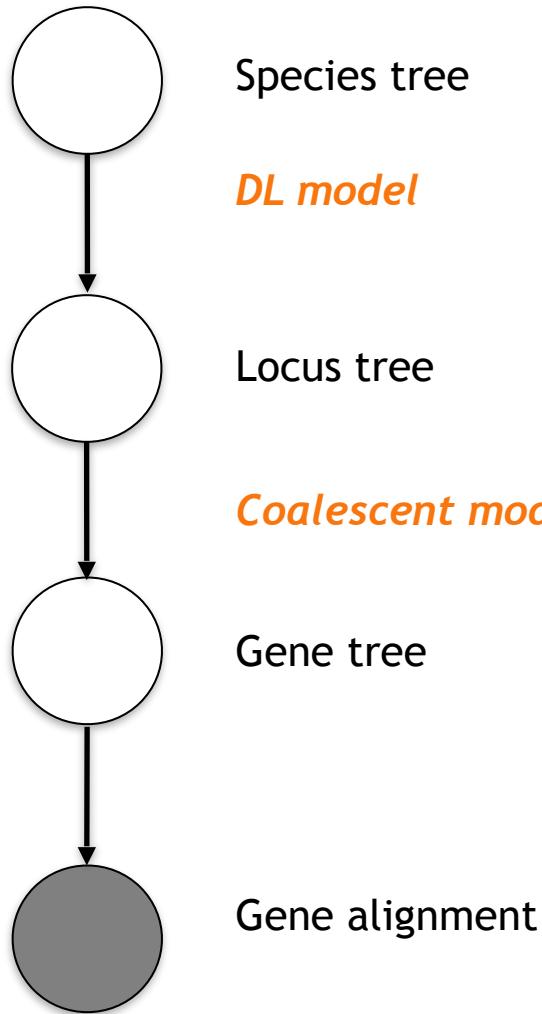
The Rasmussen and Kellis model (2012)



The Rasmussen and Kellis model (2012)

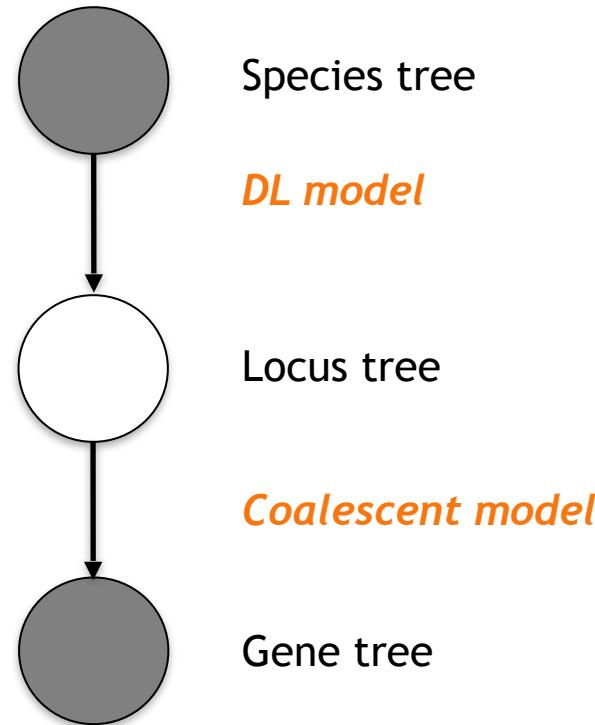
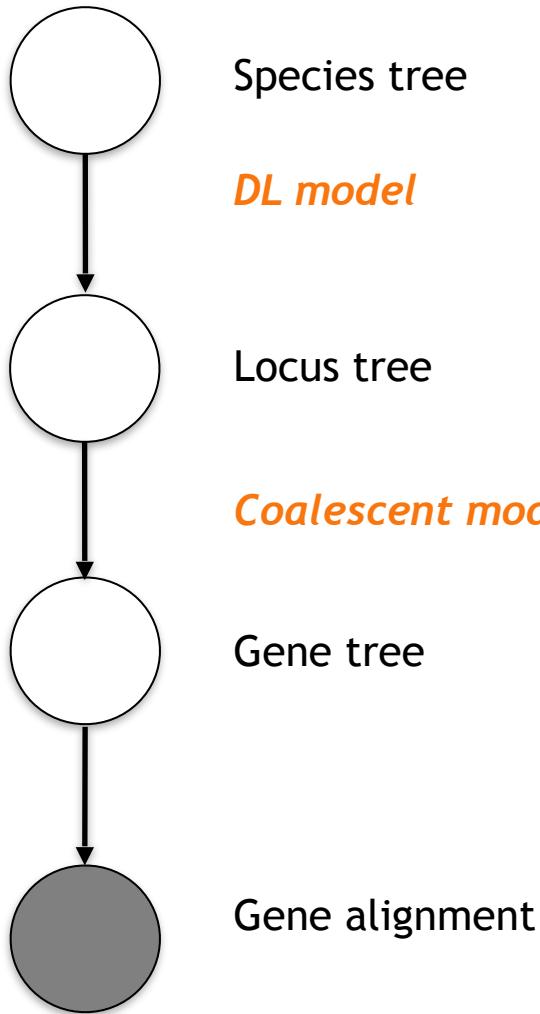


The Rasmussen and Kellis model (2012)



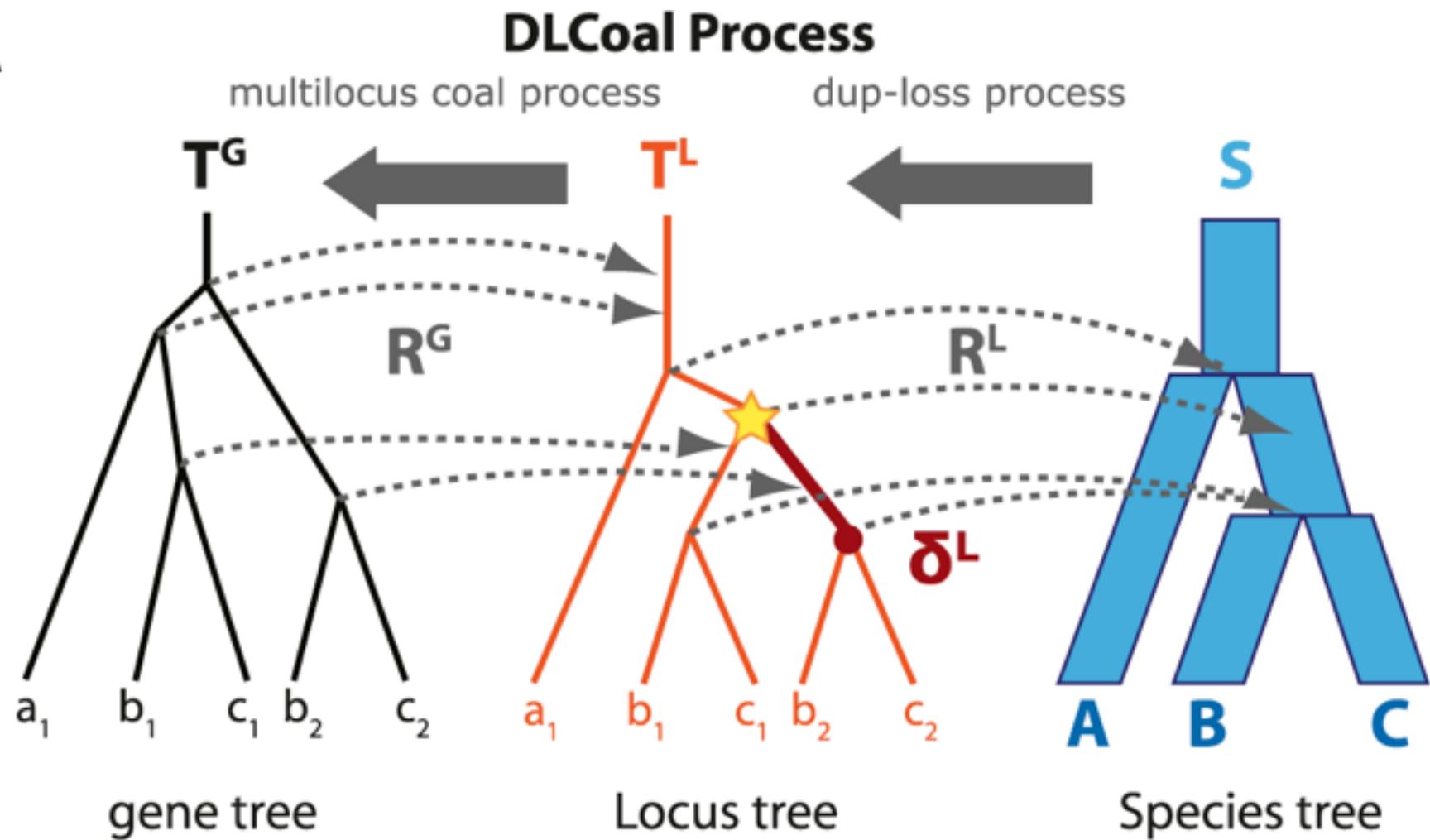
The Rasmussen and Kellis model (2012)

In practice:

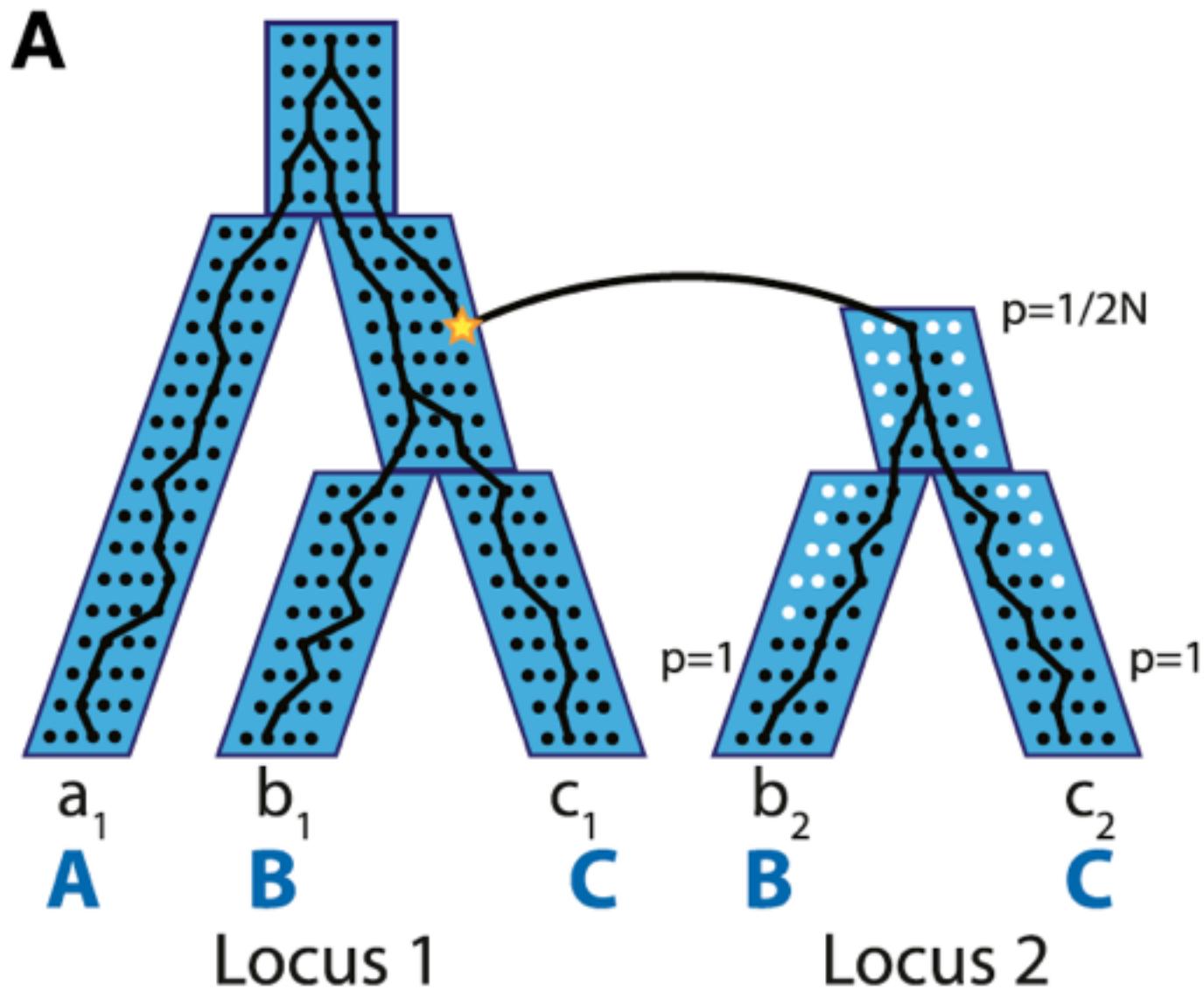


The Rasmussen and Kellis model (2012)

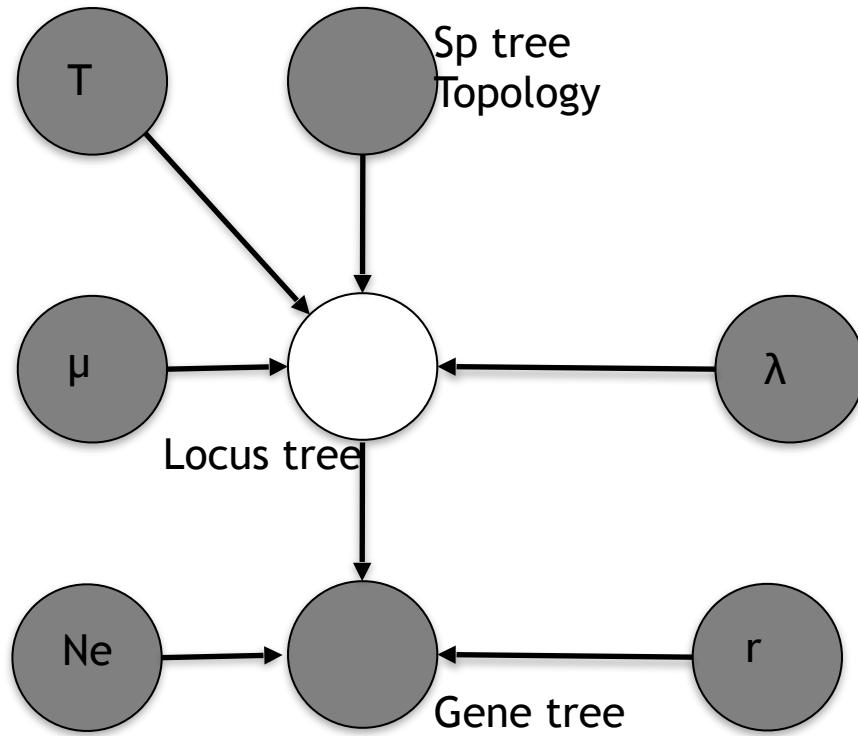
A

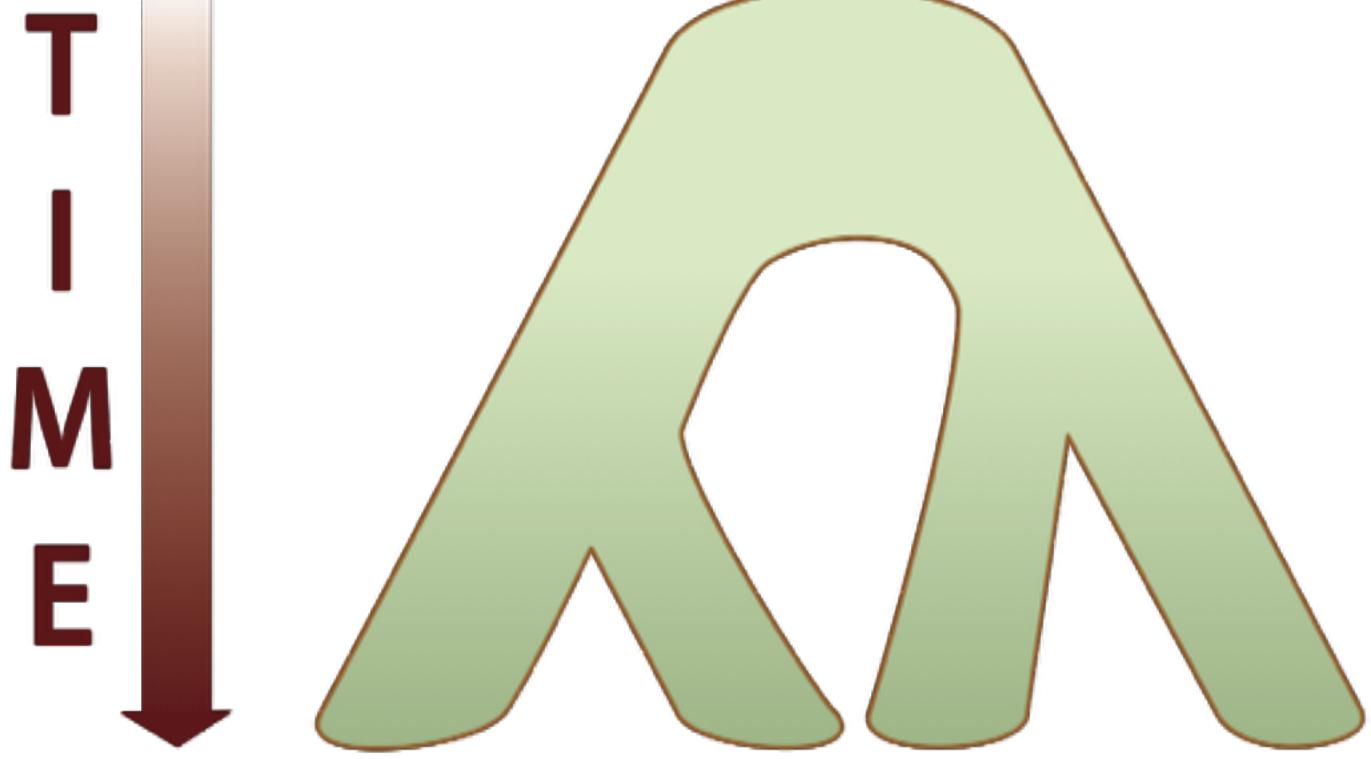


The Rasmussen and Kellis model (2012)



The species tree-gene tree graphical model with DL + ILS (Rasmussen and Kellis 2012)





Discrete character:
Continuous character:
Species:



a

0.1

A



a

0.2

B



b

0.2

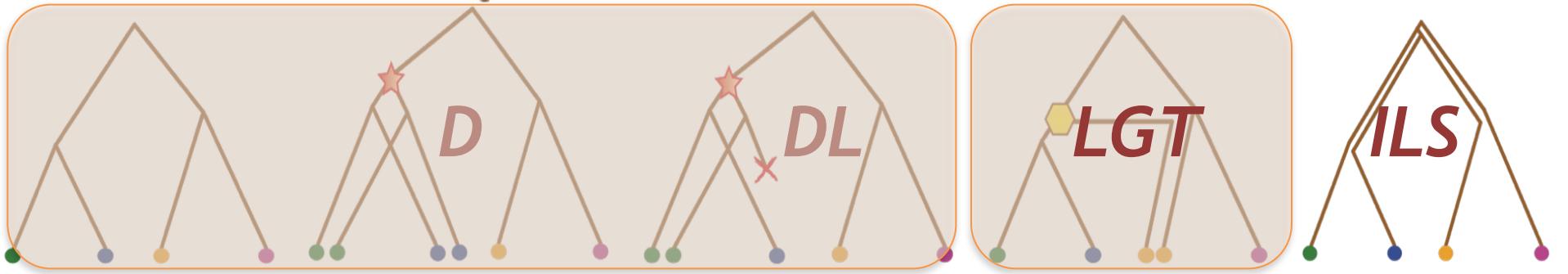
C



a

0.4

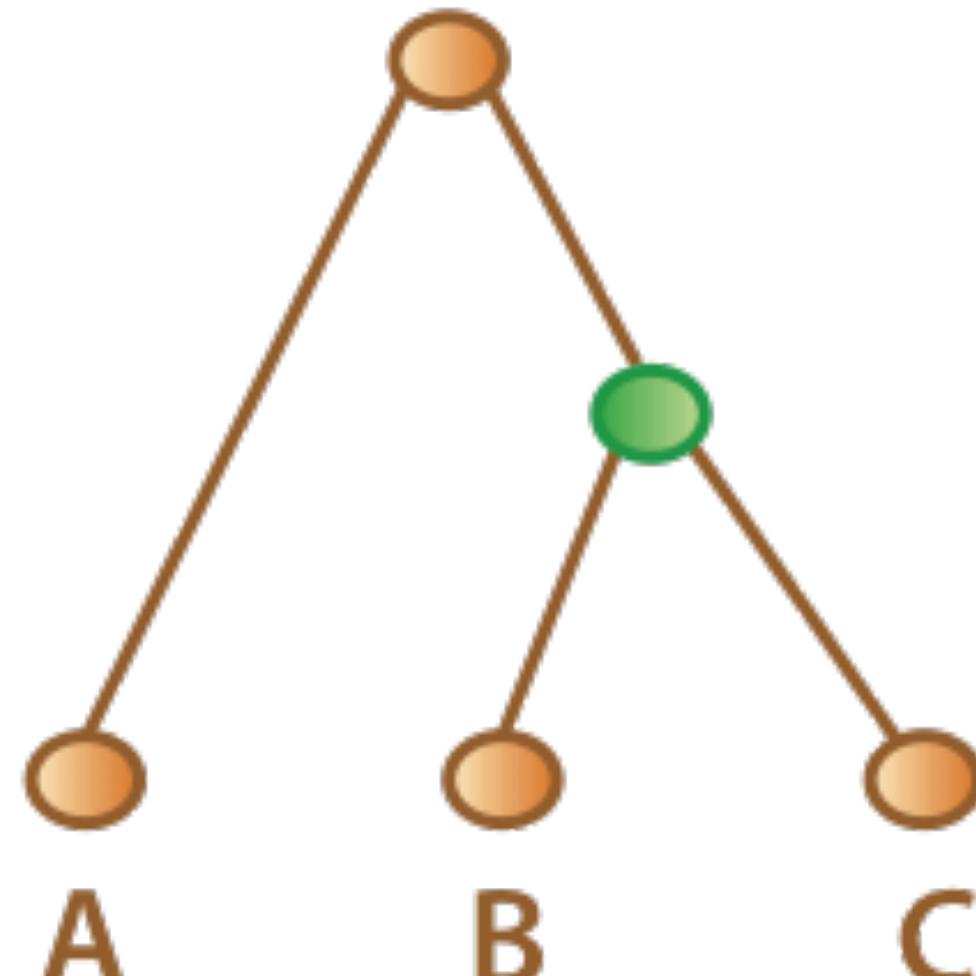
D



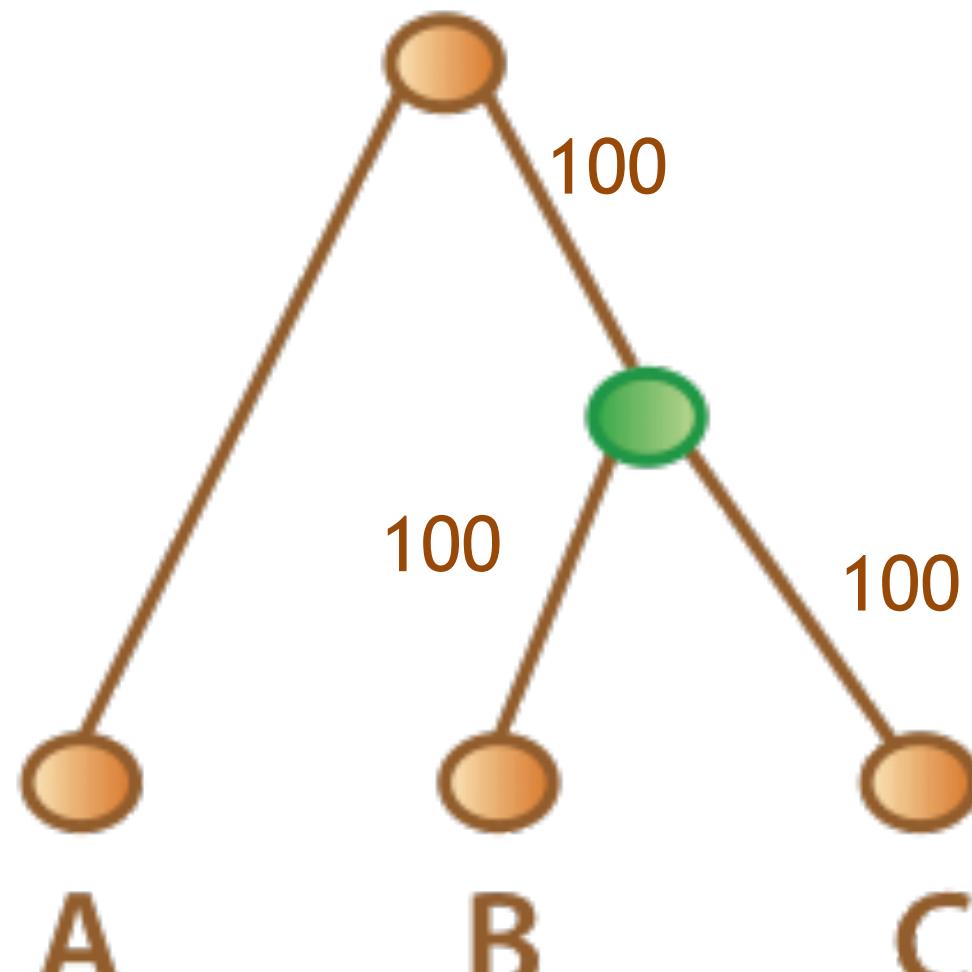
Model that combines DL and T

- Reconciliation between a species tree and a gene tree:
Parsimony setting: Gorecki, RECOMB 2004; Gorbunov and Lyubetski Mol. Biol. (Mosk), 2009; Libeskind-Hadas and Charleston, JCB 2009; Tofigh et al., IEEE ACM 2011; Doyon et al., Comparative Genomics 2011
- Reconstruction of a species tree given gene trees:
 - Model-based method: Szöllősi et al., PNAS 2012; Szöllősi et al., Systematic Biology 2013a, 2013b; Sjöstrand et al., Syst Biol 2014

Dating the tree of life is difficult

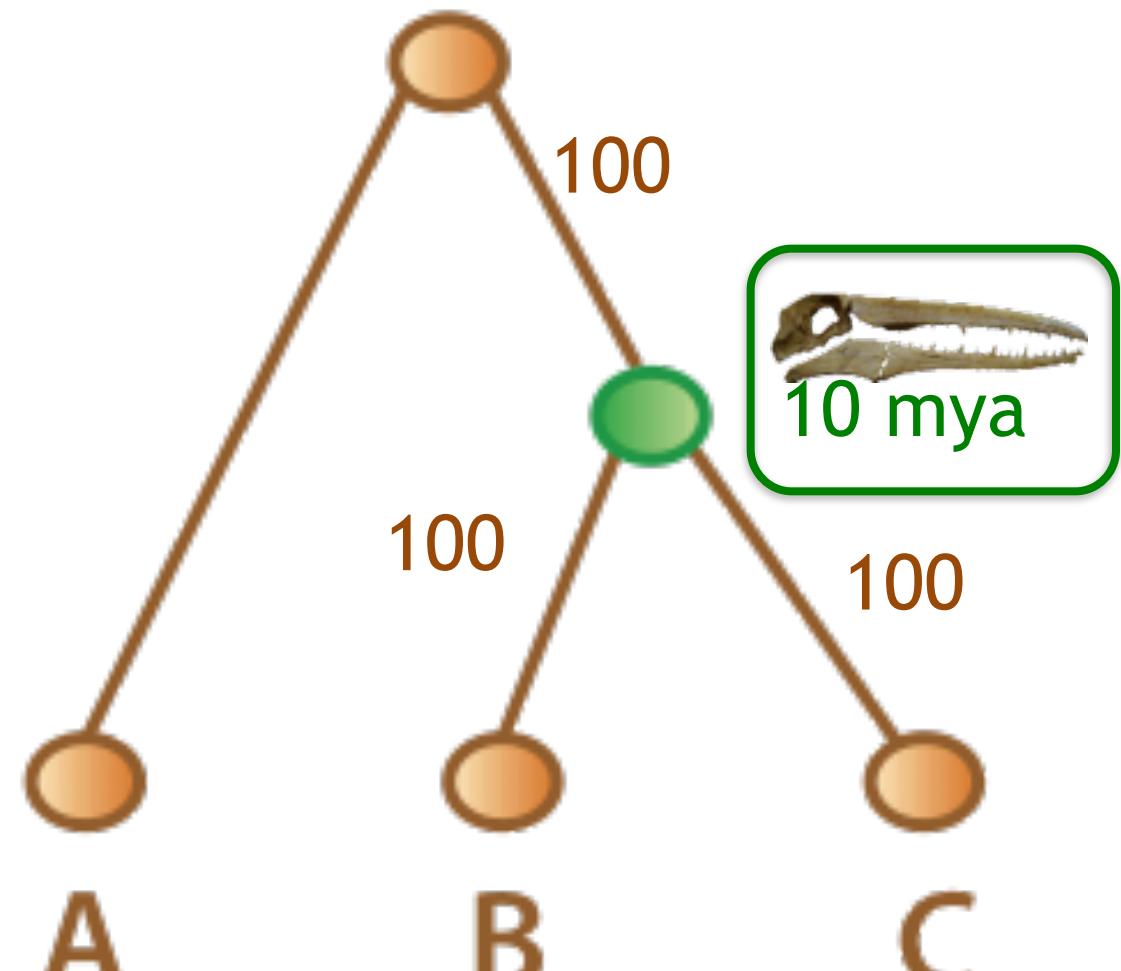


Dating the tree of life is difficult

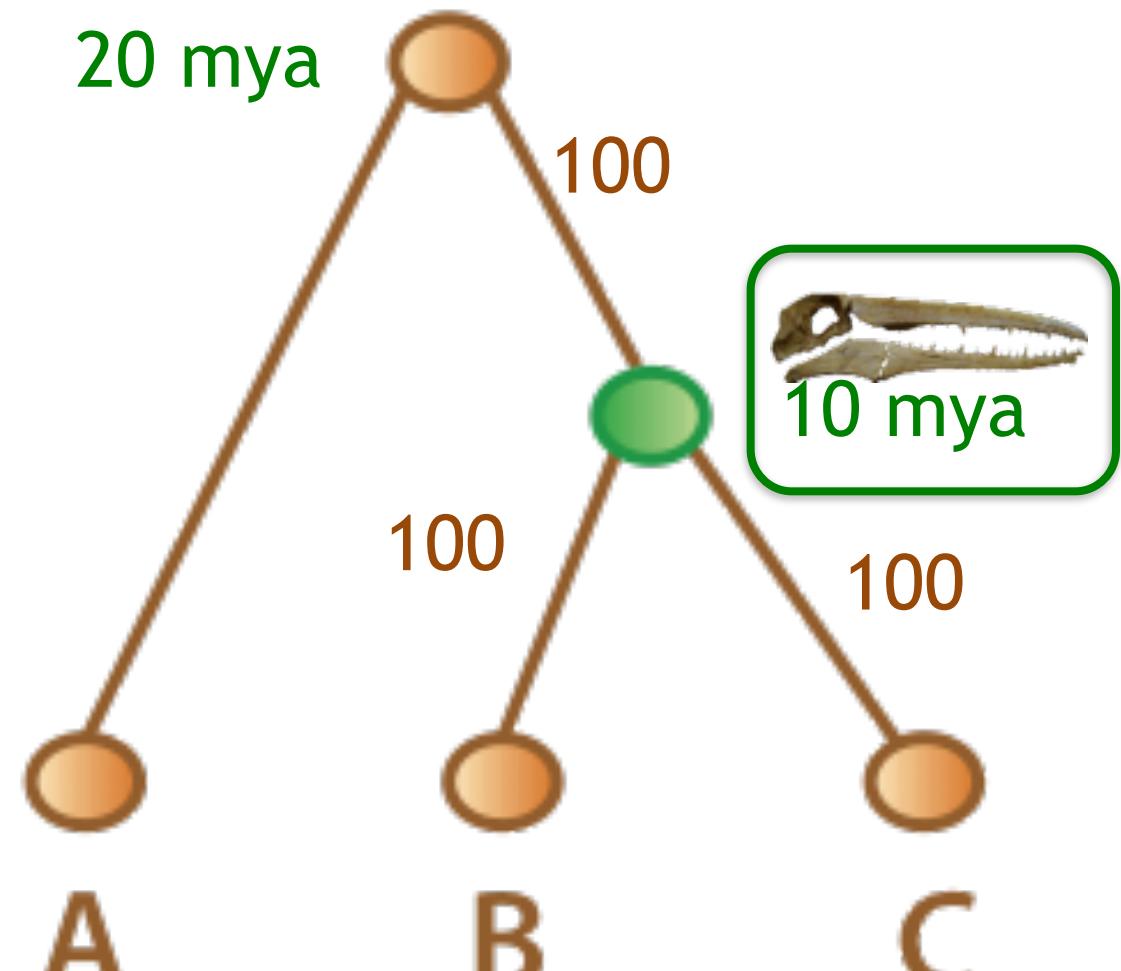


Species:

Dating the tree of life is difficult

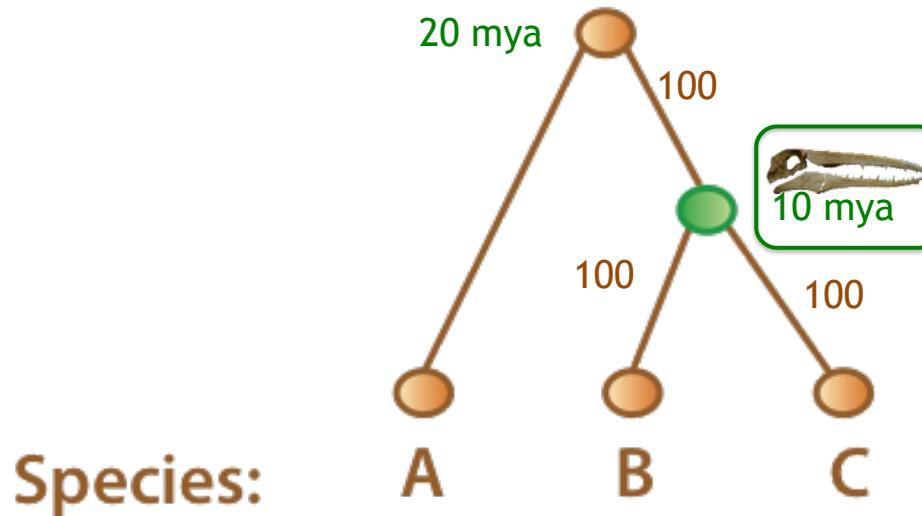


Dating the tree of life is difficult



Species:

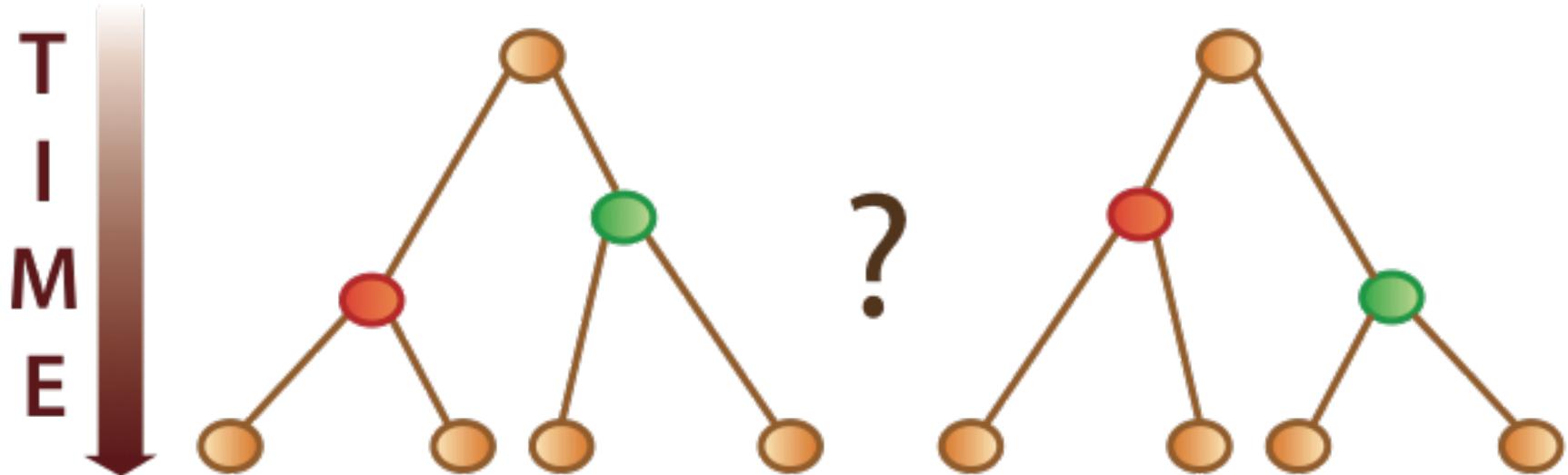
Dating the tree of life is difficult



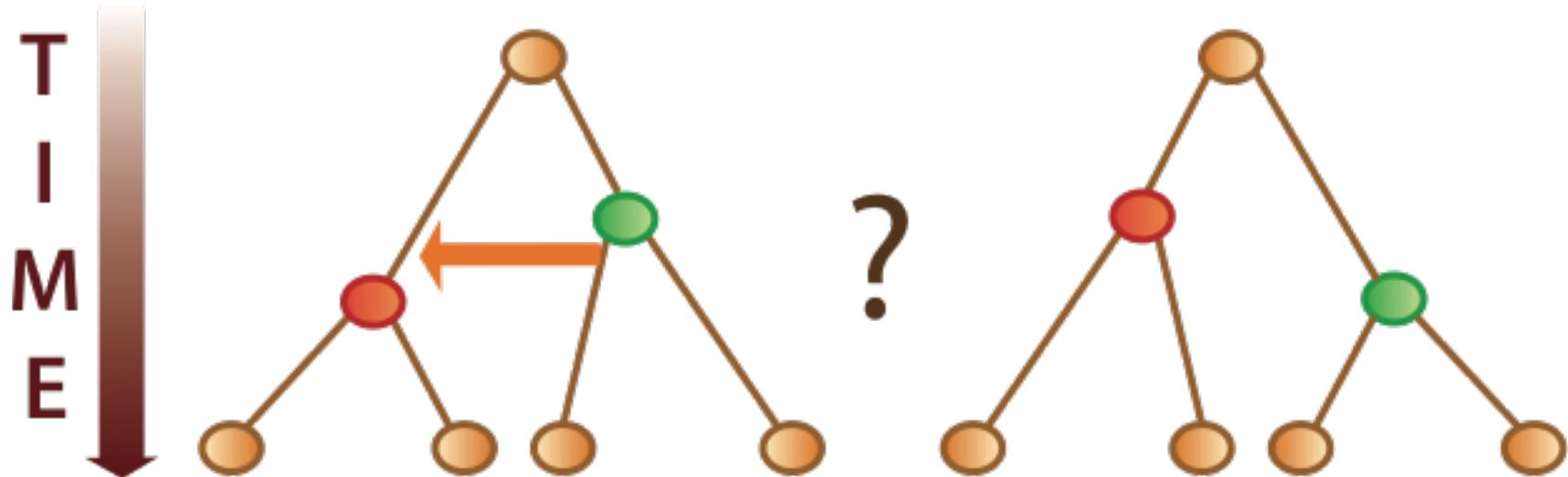
Very few fossils for the first 3 billion years of evolution on Earth

Other constraints are needed to date the tree of life

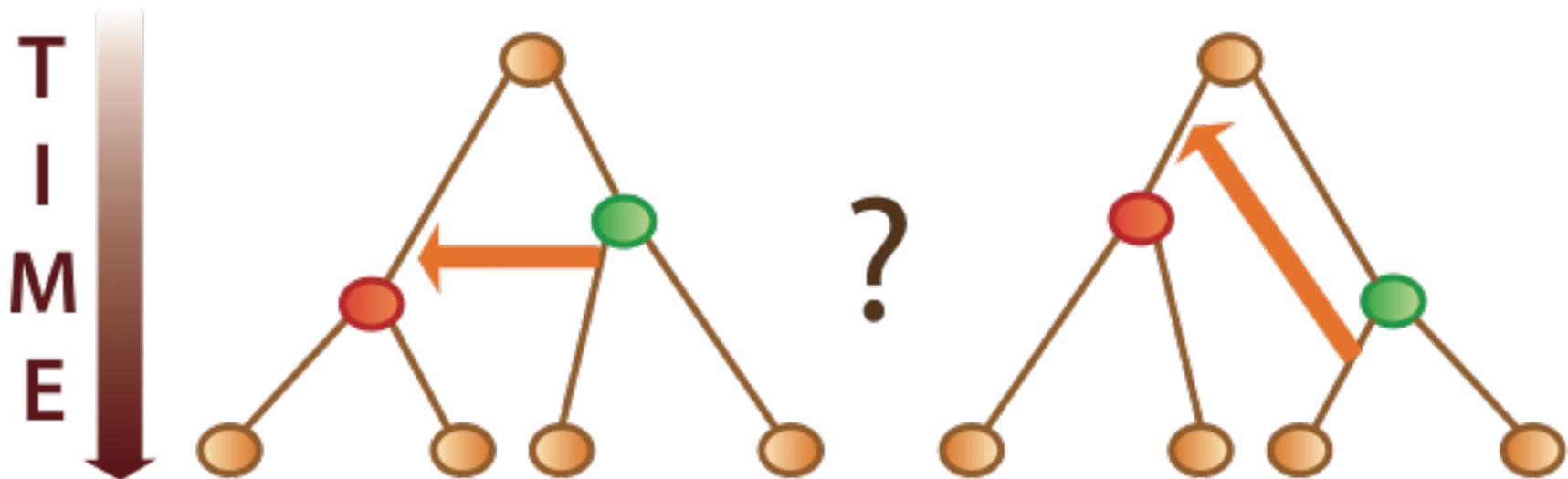
Using transfers to date clades



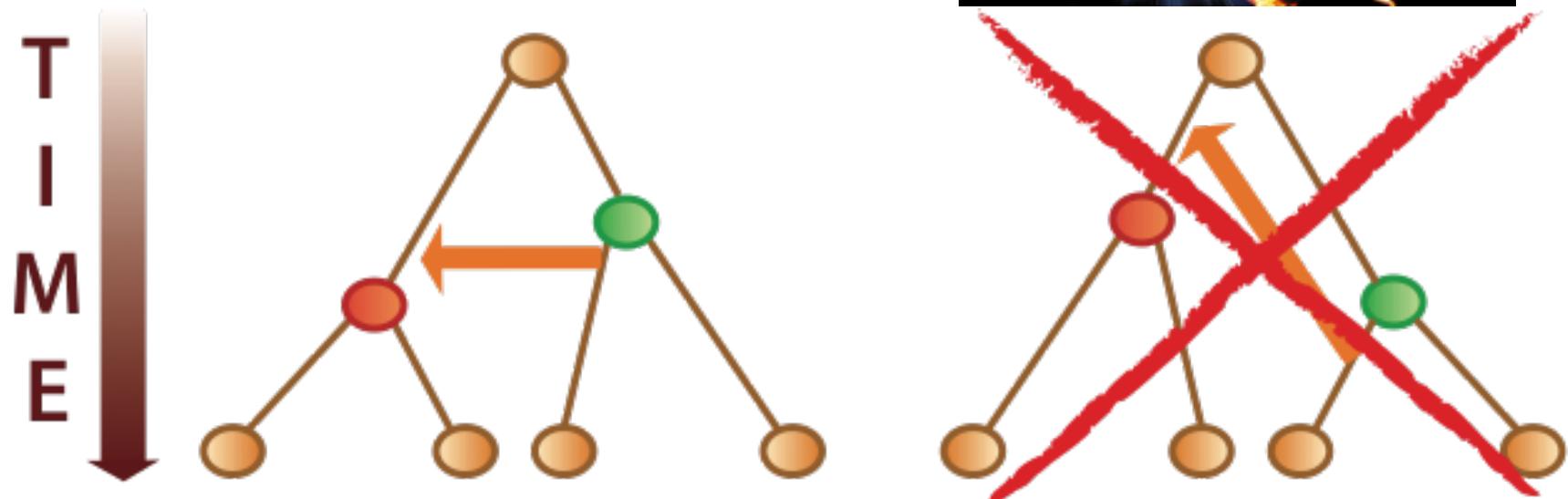
Using transfers to date clades



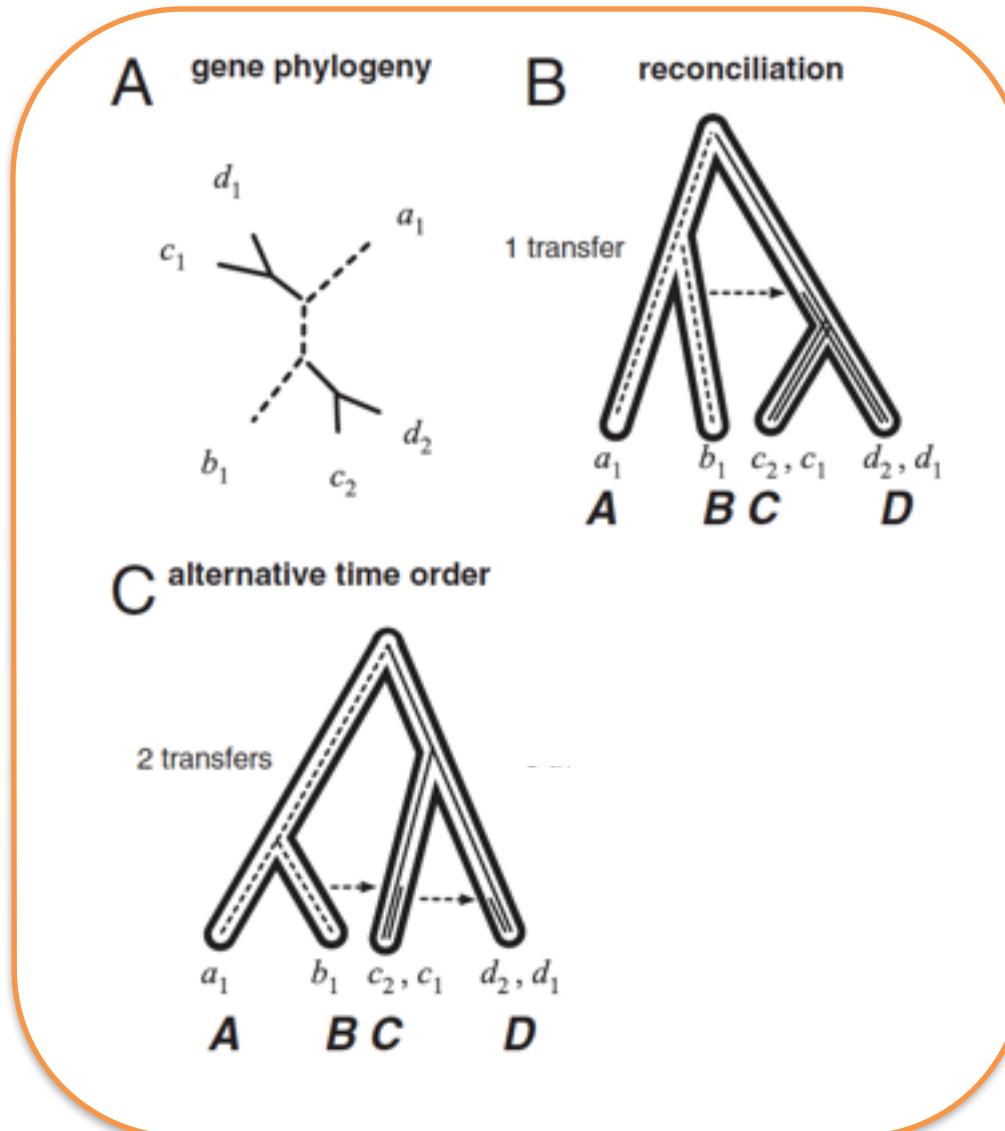
Using transfers to date clades



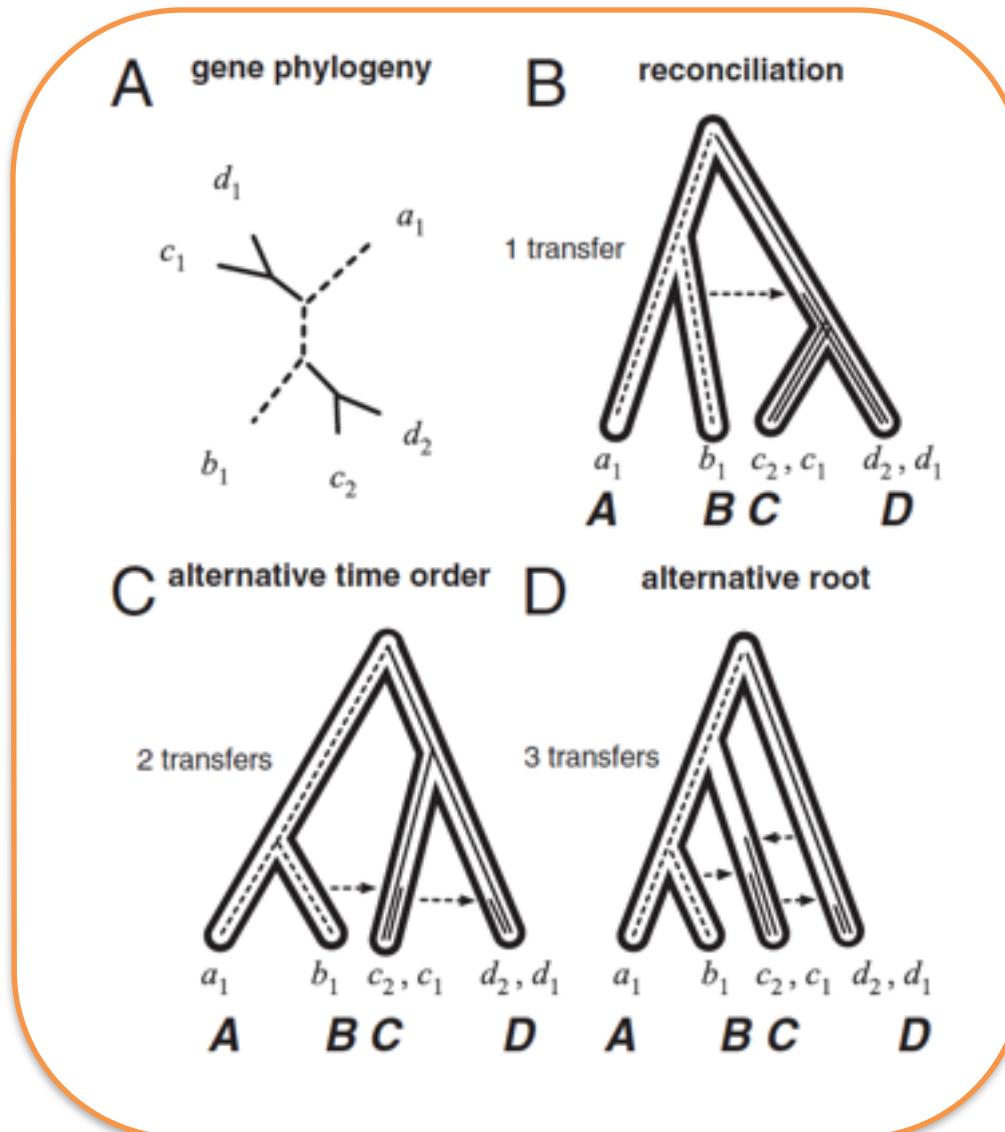
Using transfers to date clades



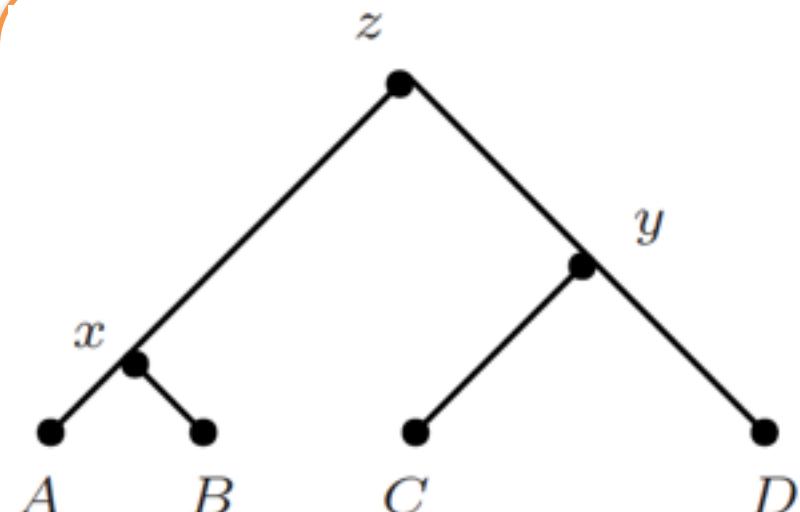
Gene transfers provide means to root and date species trees



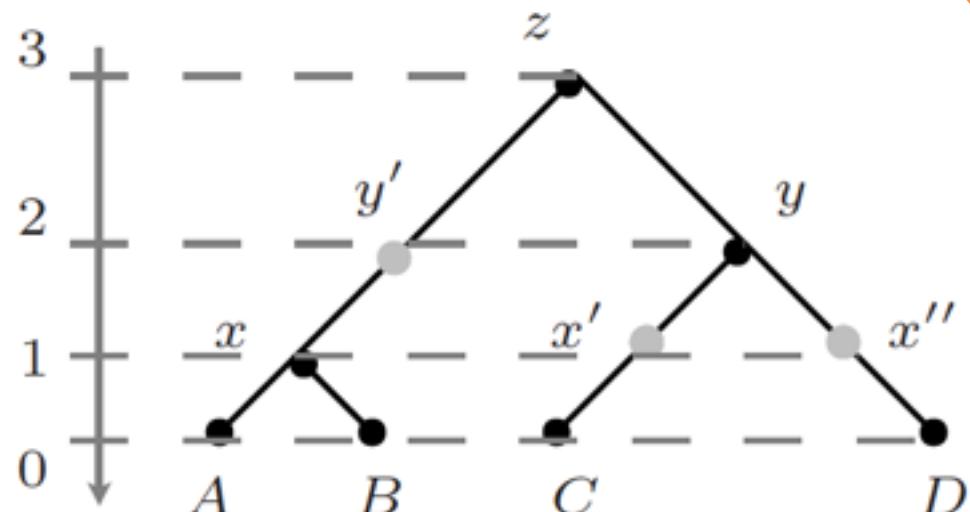
Gene transfers provide means to root and date species trees



Avoiding back to the future scenarios

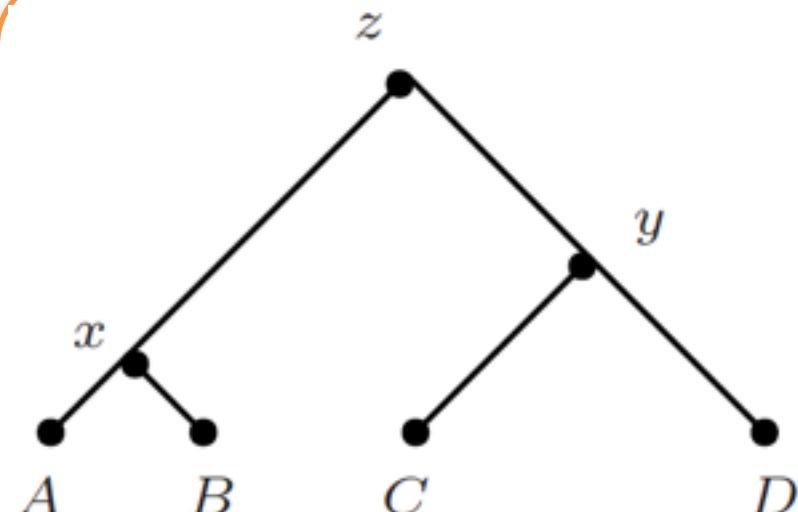


(a) A species tree S

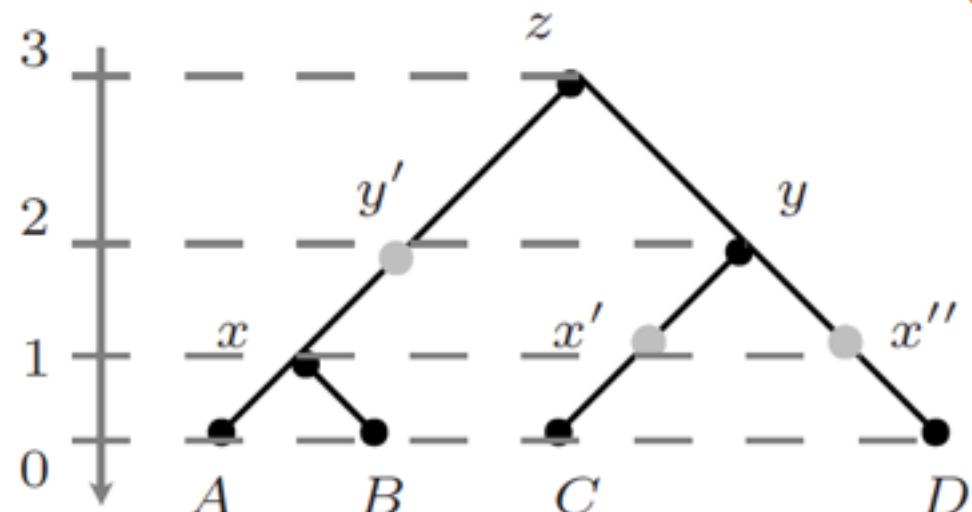


(b) The subdivision S' of S

Avoiding back to the future scenarios



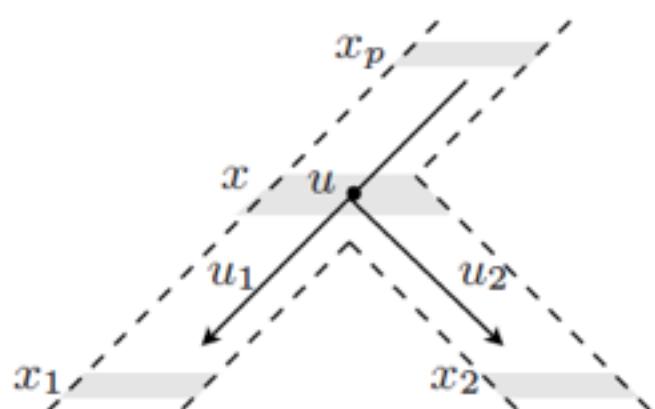
(a) A species tree S



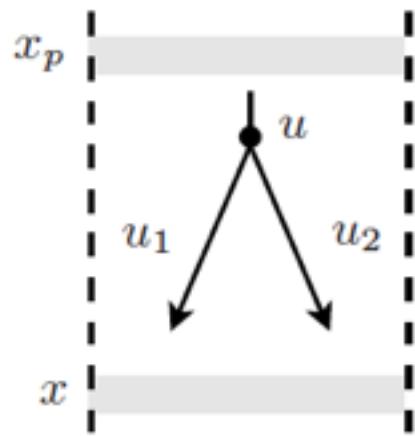
(b) The subdivision S' of S

→ Algorithm in $O(\#NodesInS' \cdot \#NodesInG)$

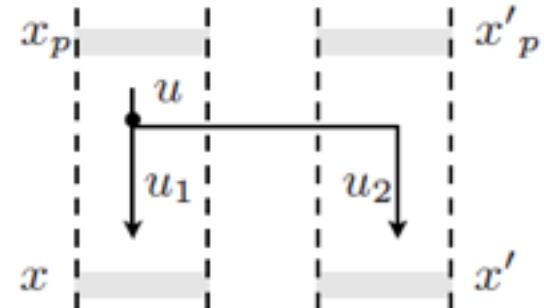
6 possible events are enough



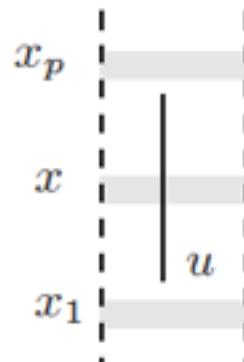
(a) Speciation (S)



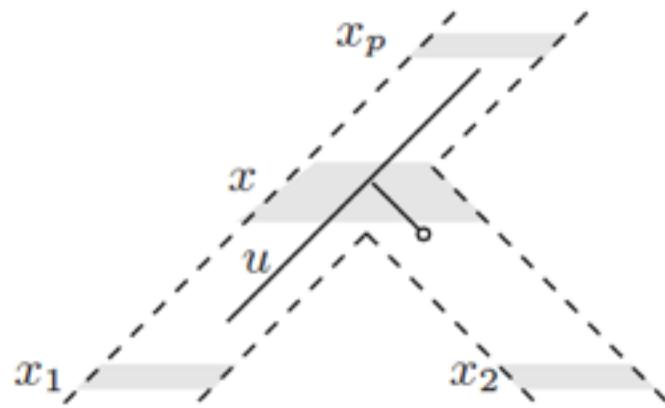
(b) Duplication (D)



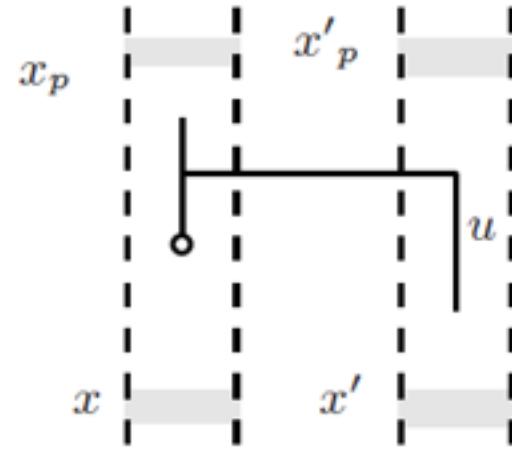
(c) Transfer (T)



(d) Ø event



(e) Speciation + Loss (SL)



(f) Transfer+Loss (TL)

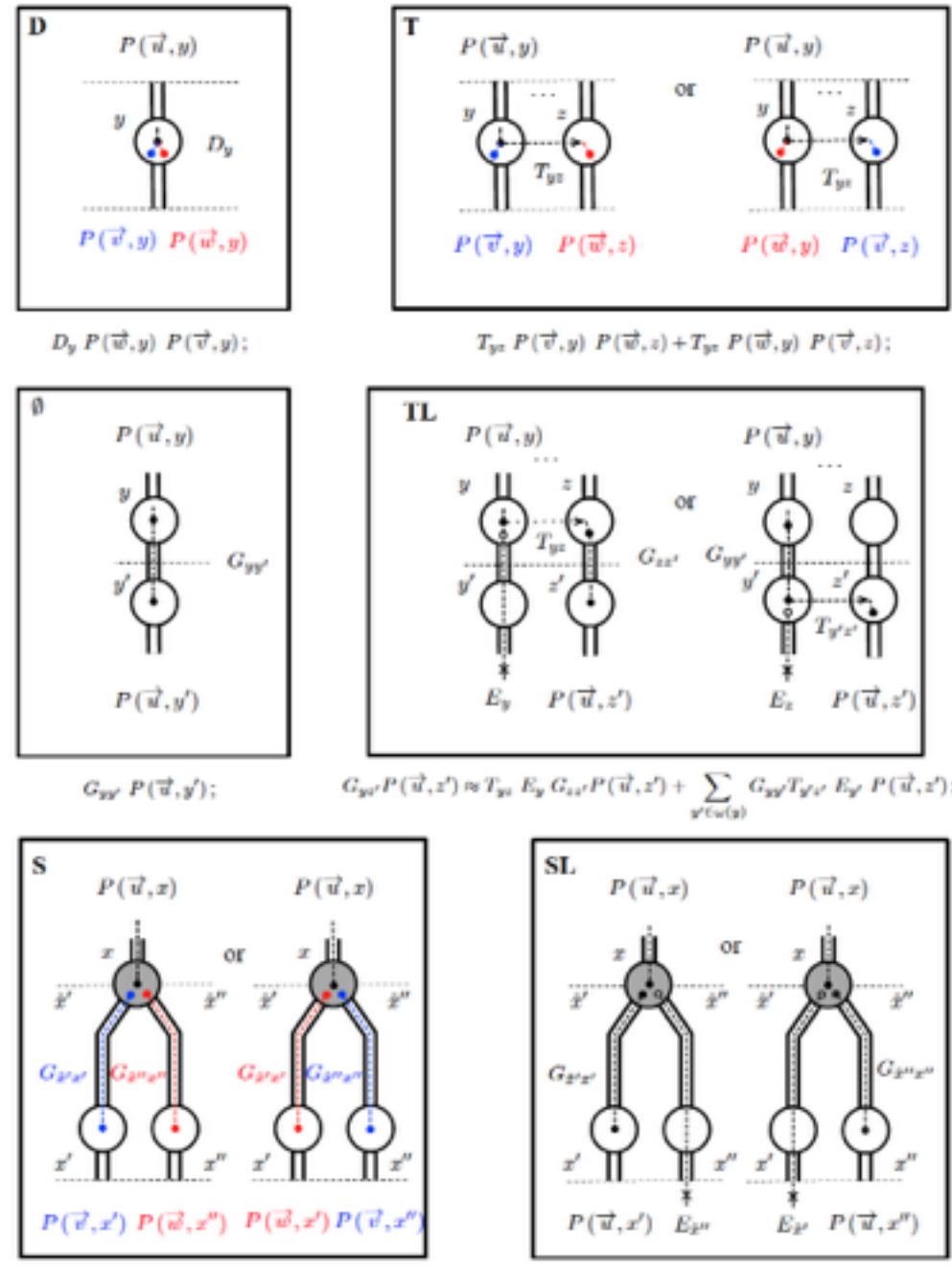
Dynamic programming in the Doyon et al. 2011 approach

- Input: rooted augmented species tree S' , rooted gene tree G
- Output: reconciliation scenario between S' and G
- Finding scenarios can be done in polynomial time using DP:
 - Fill the cost matrix c ($\# \text{NodesIn}G$, $\# \text{NodesIn}S'$)
 - Post-order traversal of the gene tree
 - For each node i of G :
 - Map it to each node j of S'
 - Consider all of 6 events, and set $c(i, j) = \min(\text{cost(possible events)})$
 - At the root, return the minimum cost found

From parsimony to model-based

- Input: gene trees
- Output: reconciled rooted gene trees, rooted species tree with ordered nodes
- Birth-death process to model gene evolution:
 - Birth parameters:
 - D
 - T: rate of receiving a gene through transfer
 - Death parameter:
 - L
- One or more sets of D, T, L parameters
- Double recursive DP for computing $P(\text{gene tree} \mid \text{species tree})$
- Integrating over all possible scenarios
- Sliced species tree

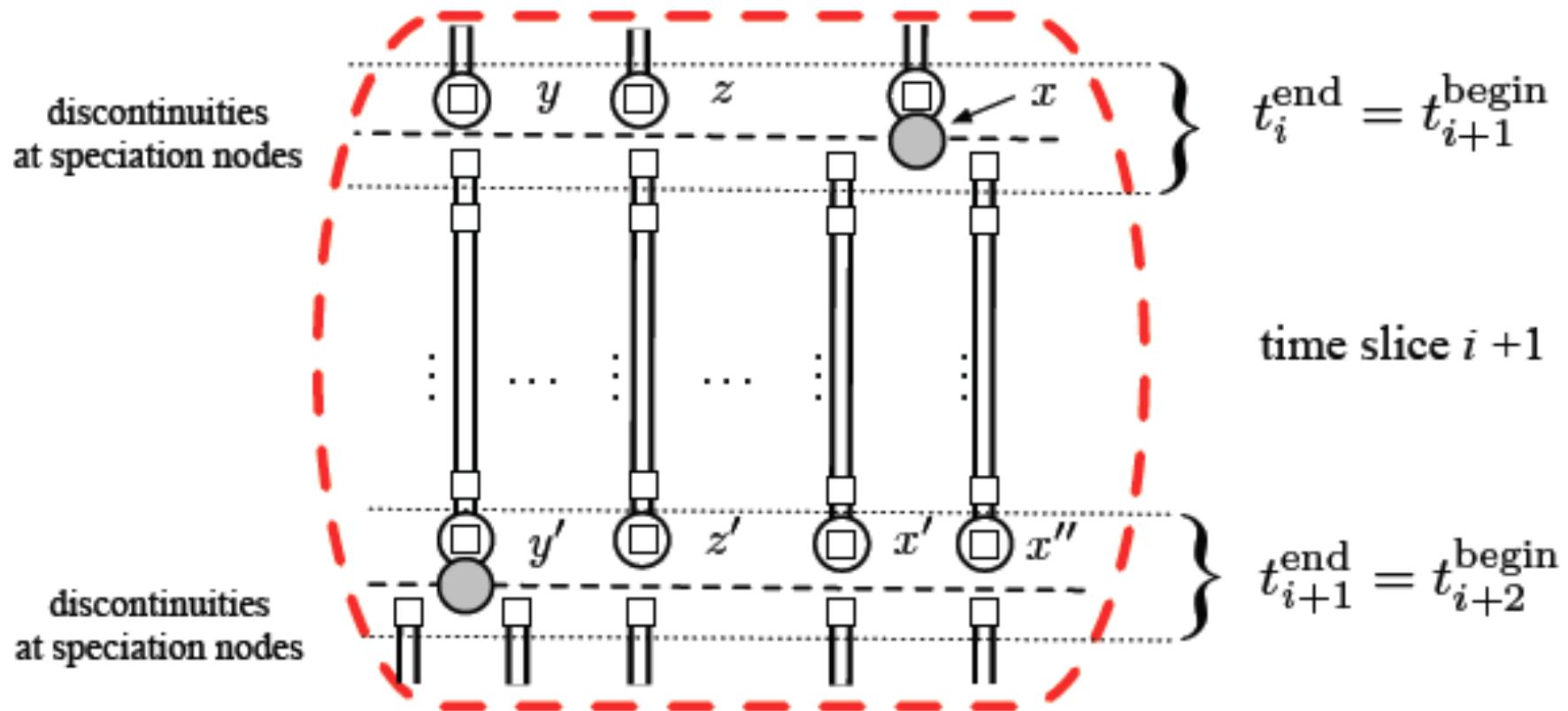
Use of the 6 events from Doyon et al. 2011



Within-slice discretization for solving ordinary differential equations

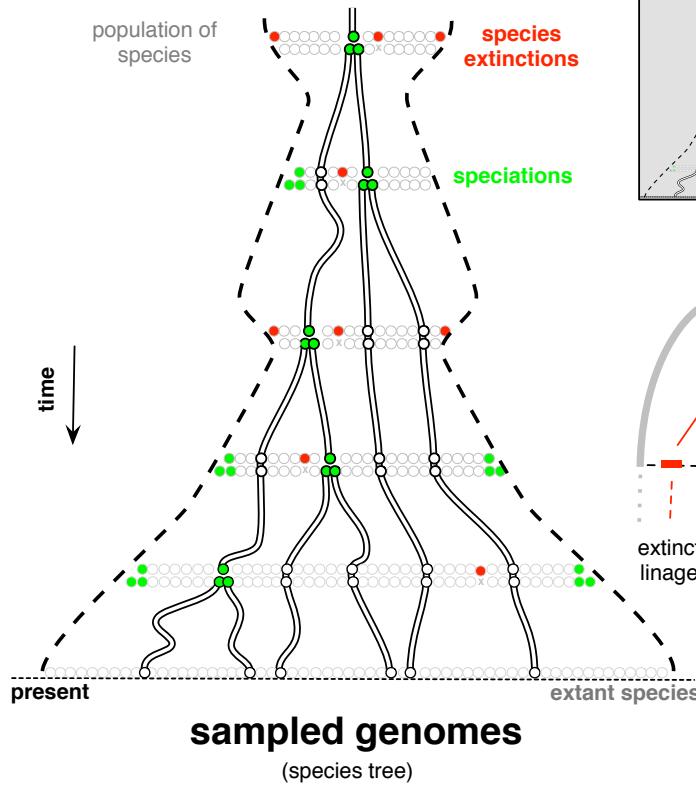
C

within slice discretization



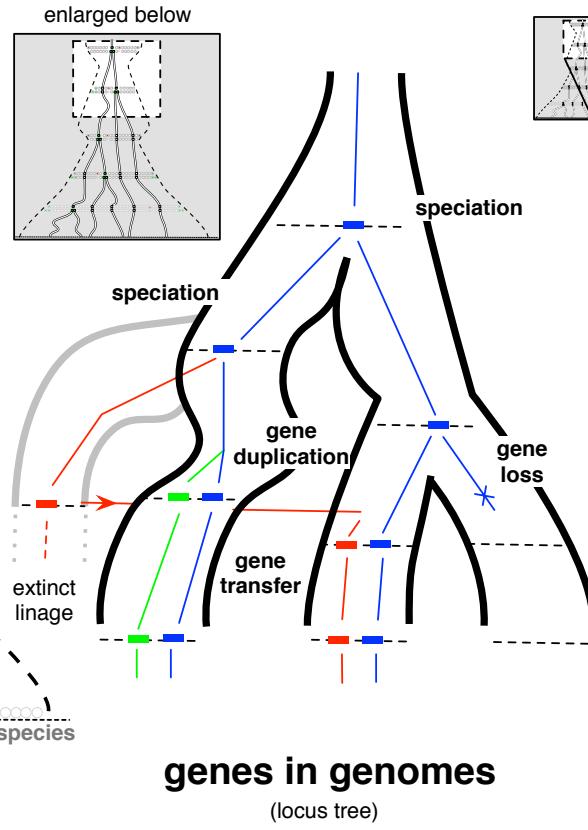
A hierarchy of processes

a) species diversification



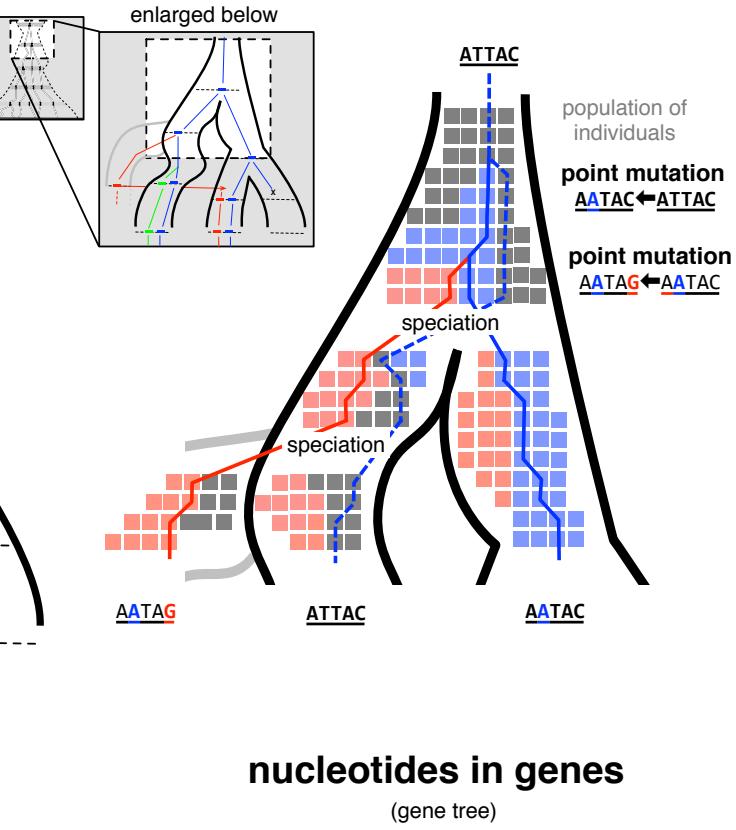
b)

gene birth and death

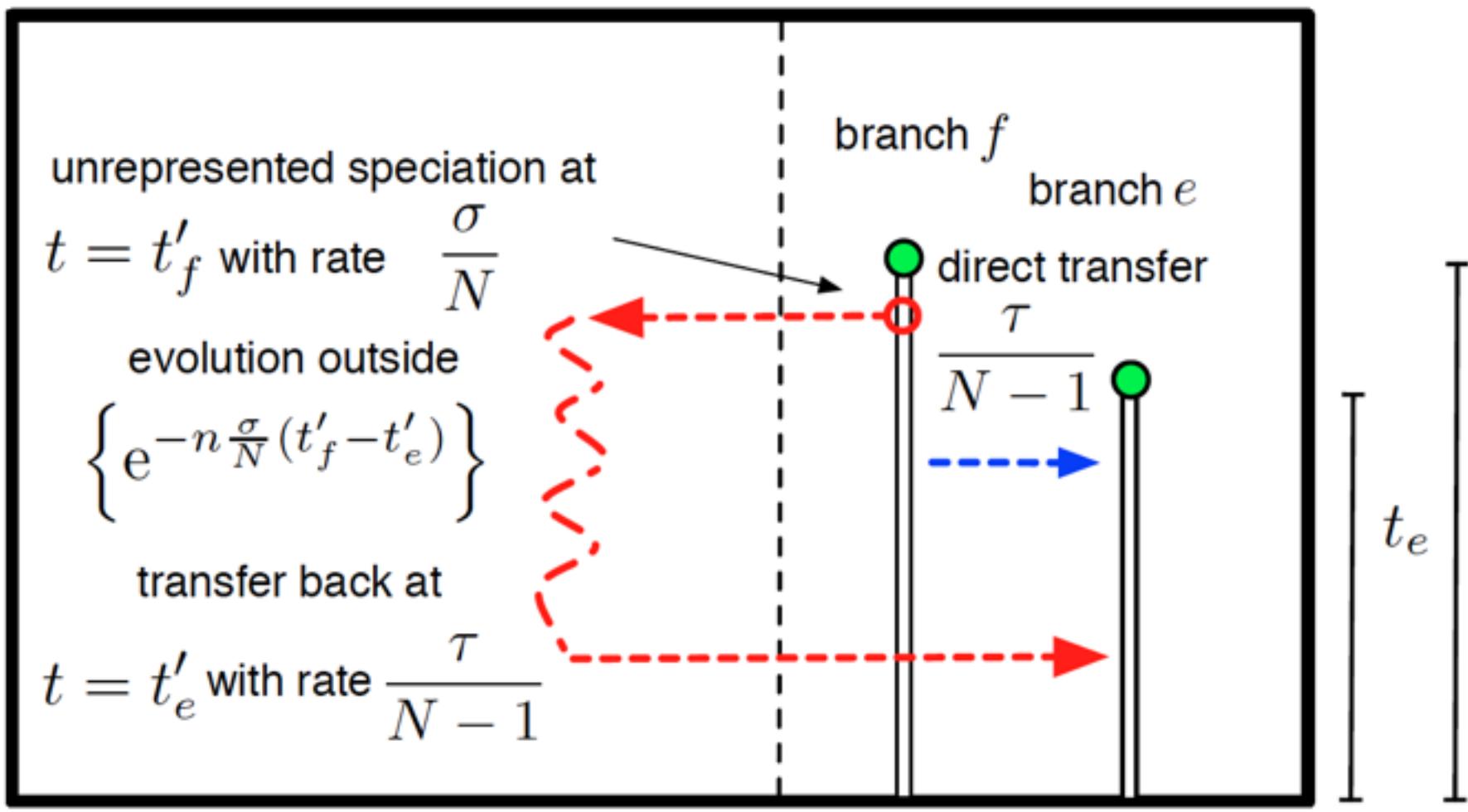


c)

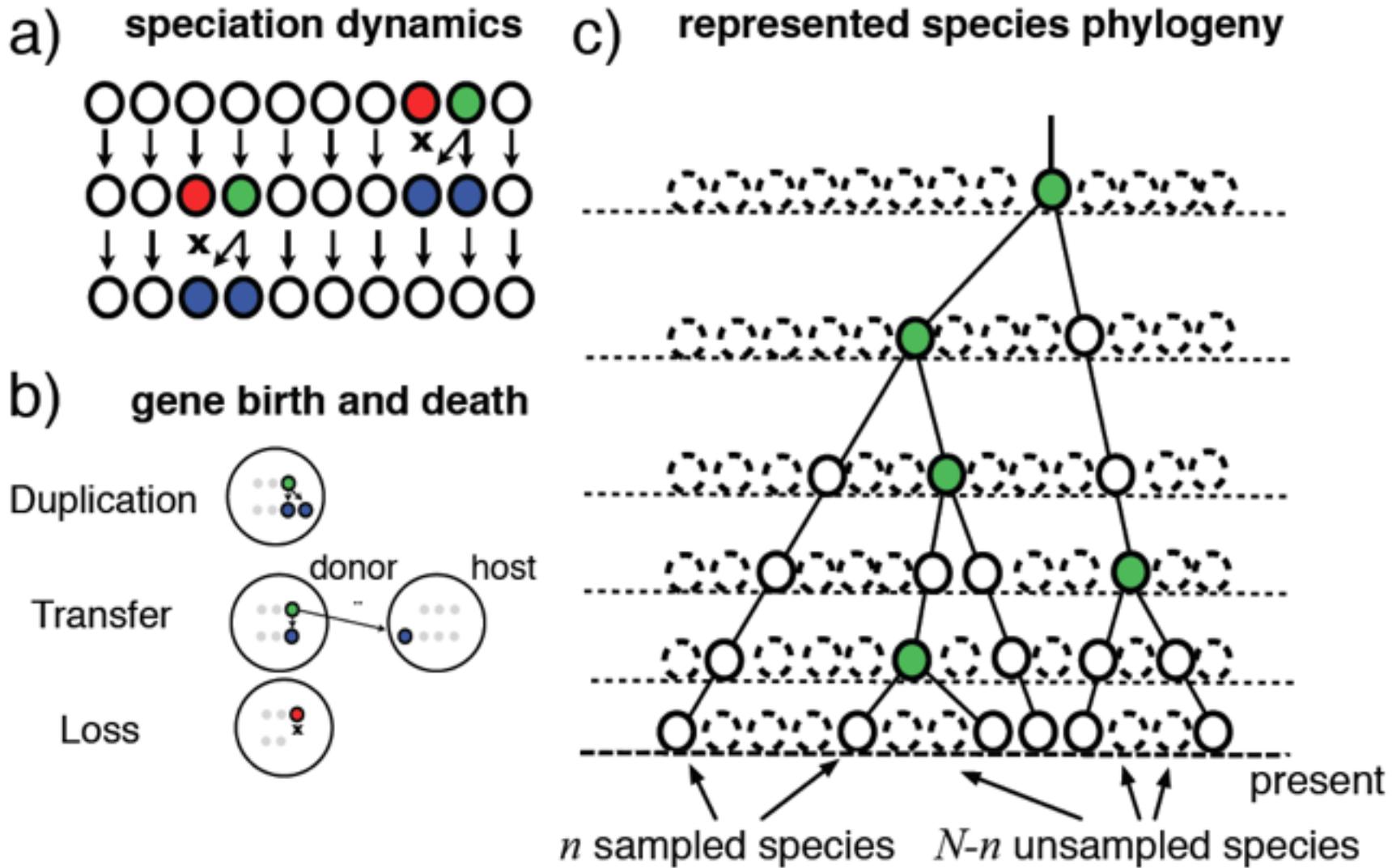
sequence substitution



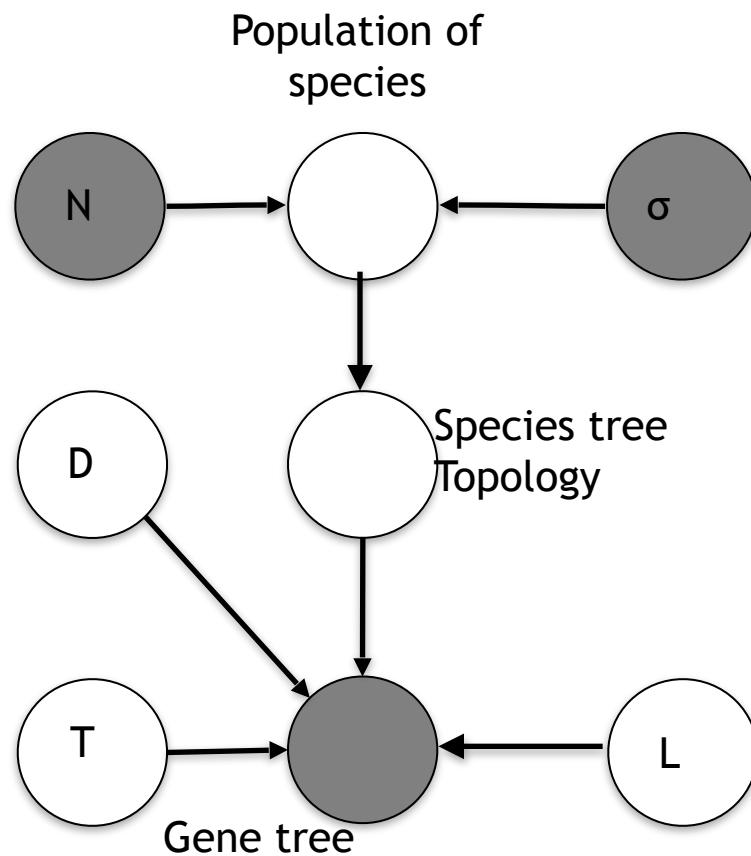
Further refinement: seeing dead lineages



Further refinement: seeing dead lineages



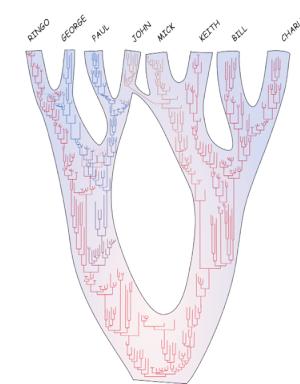
The exODT graphical model



Plan

Modeling the relationship between species tree and gene tree

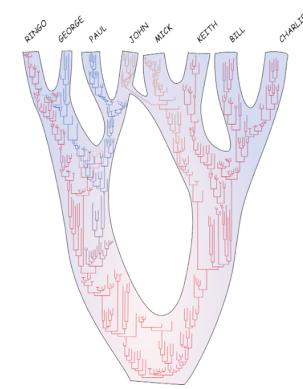
- coalescent models
- models of gene duplication and loss
- models of gene transfer
- models that combine the above



Gene tree species tree models in RevBayes, and tutorial

Gene tree species tree models in RevBayes

- Multispecies coalescent
- Pomo model
- In progress:
 - exODT (DTL model)
- Soon:
 - Akerborg et al. 2009 (DL model)
 - MP-EST



The tutorial

- A walk around the multispecies coalescent:
 1. Simulation with the multispecies coalescent (play with the parameters!)
 2. Then:
 - Inference with the multispecies coalescent from gene trees
 - Inference with the multispecies coalescent from gene alignments
 - Inference with the POMO model from gene alignments
 - Inference by concatenating the gene alignments

The tutorial

Files are in:

tutorials/RB_GeneTreeSpeciesTree_Tutorial/RBTutorial_files/RevBayes_scripts/
in five folders. **It is highly recommended to start RevBayes from each
of these folders when doing the exercises to avoid path problems.**

There is one simulation-only folder, used for better understanding the
multispecies coalescent:

UnderstandingMultiSpeciesCoalescent

Then there are 4 inference folders.

In each inference folder, there are 3 scripts. The file XXXMCMC.Rev is the
one that can launch the entire analysis, which consists of three parts:

Simulating the data

Constructing the model for inference

Preparing and running the MCMC