

# Phylogenomics of closely related species and individuals



Matthew Rasmussen

Siepel lab, Cornell University

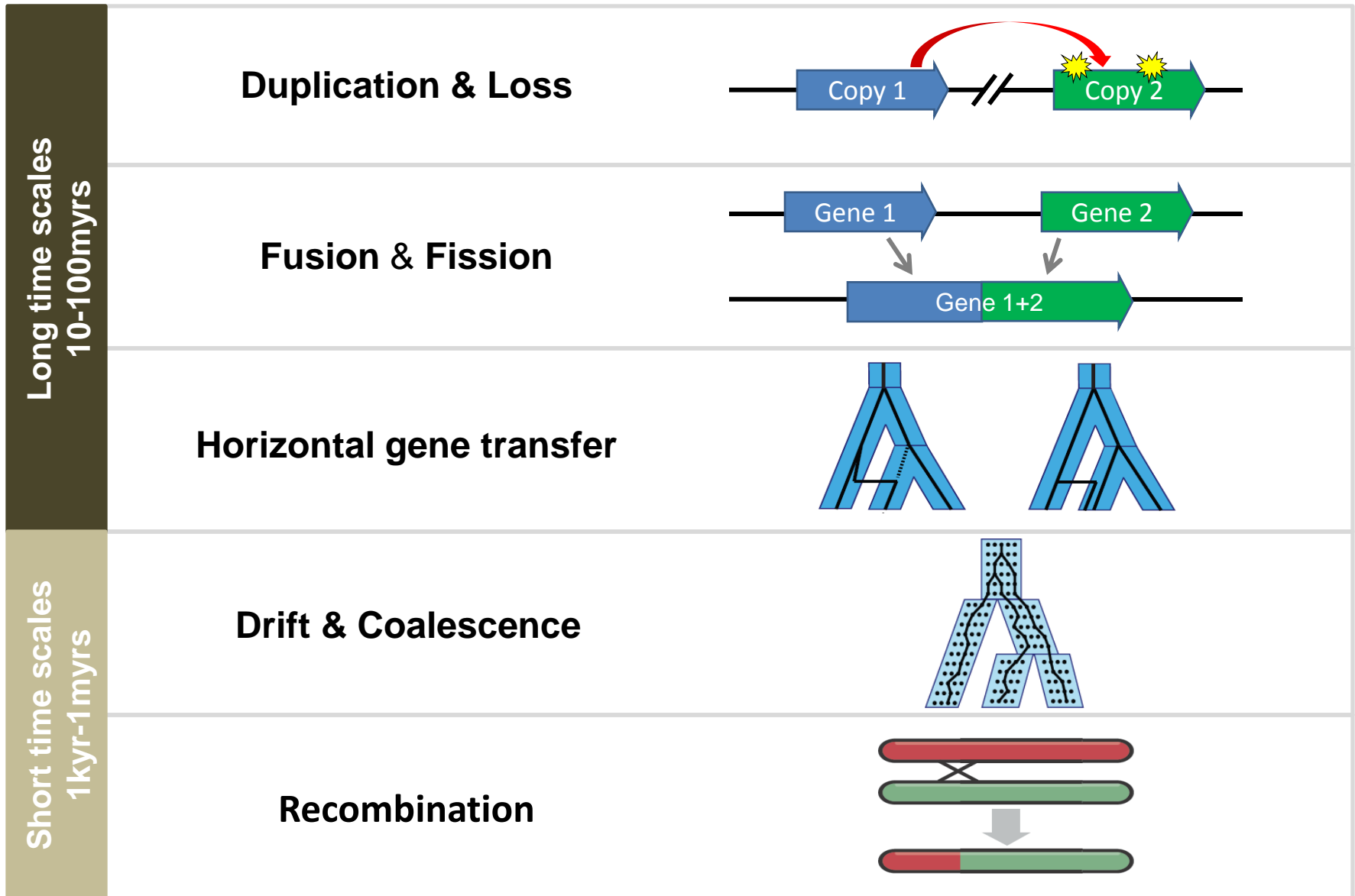
In collaboration with Manolis Kellis, MIT CSAIL

February, 2013



**BROAD**  
INSTITUTE

# How genomes evolve over short & long time scales

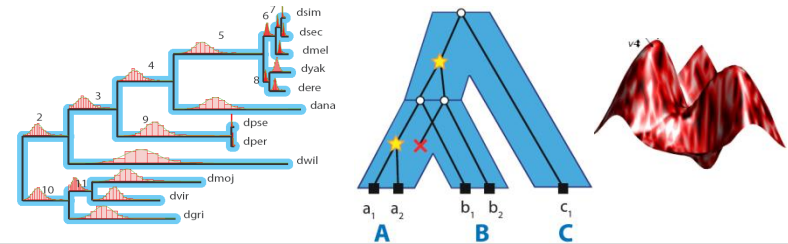


# Developing models of evolution

## Combining different evolutionary events & processes

### Duplication, Loss & Substitution

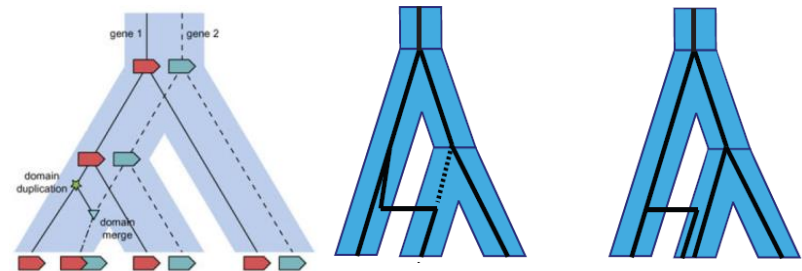
SPIDIR (GR 2007), SPIMAP (MBE 2010),  
TreeFix (SystBio 2012)



### Fusion & Fission

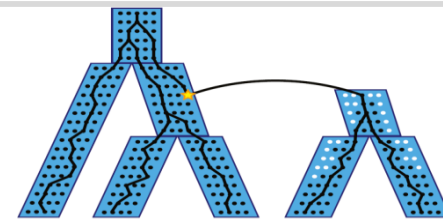
STAR-MP (MBE 2011)

Additive & Replacing transfer  
(MBE 2012)



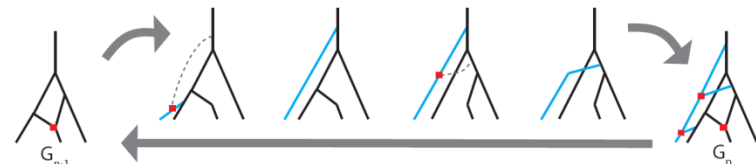
### Duplication, Loss & Coalescence

DLCoal (GR 2012)



### Coalescence & Recombination

ARGHMM

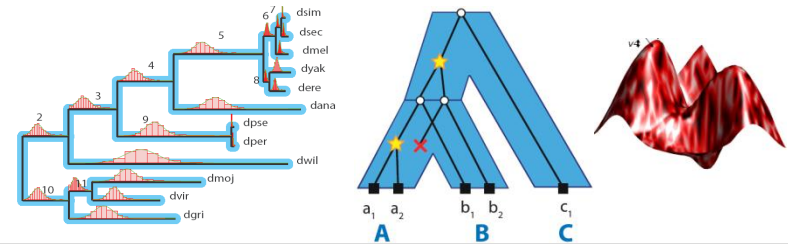


# Developing models of evolution

Combining different evolutionary events & processes

## Duplication, Loss & Substitution

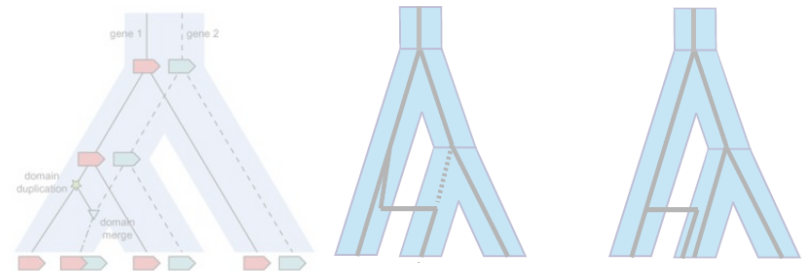
SPIDIR (GR 2007), SPIMAP (MBE 2010),  
TreeFix (SystBio 2012)



## Fusion & Fission

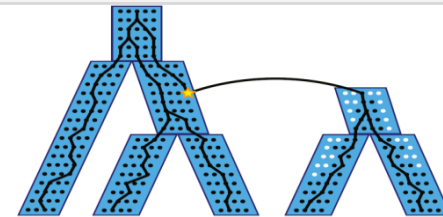
STAR-MP (MBE 2011)

Additive & Replacing transfer  
(MBE 2012)



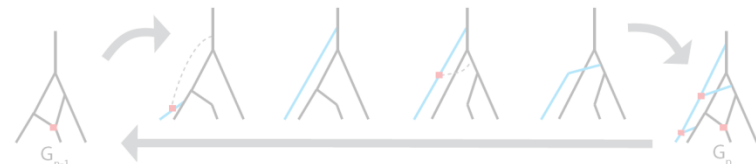
## Duplication, Loss & Coalescence

DLCoal (GR 2012)



## Coalescence & Recombination

ARGHMM

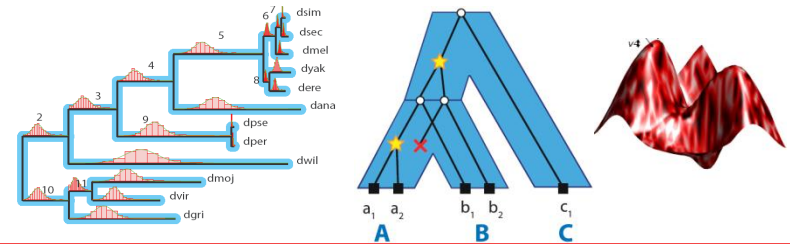


# Developing models of evolution

## Combining different evolutionary events & processes

### Duplication, Loss & Substitution

SPIDIR (GR 2007), SPIMAP (MBE 2010),  
TreeFix (SystBio 2012)

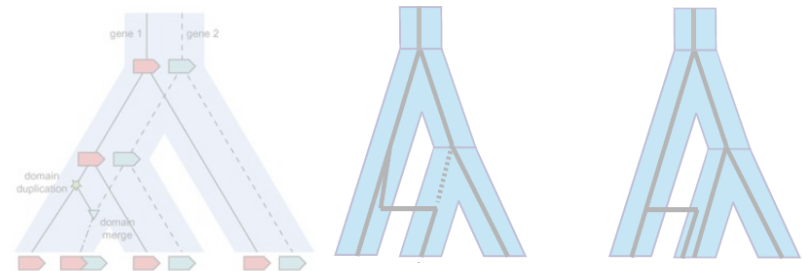


### Fusion & Fission

STAR-MP (MBE 2011)

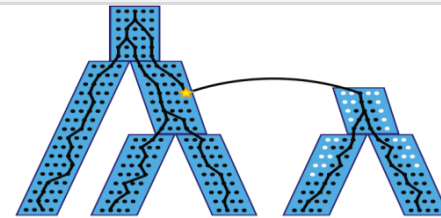
### Additive & Replacing transfer

(MBE 2012)



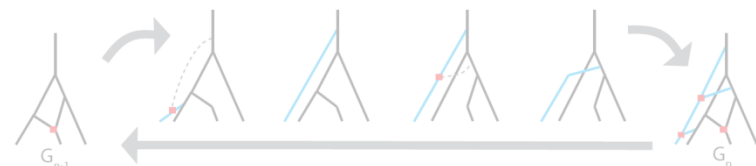
### Duplication, Loss & Coalescence

DLCoal (GR 2012)

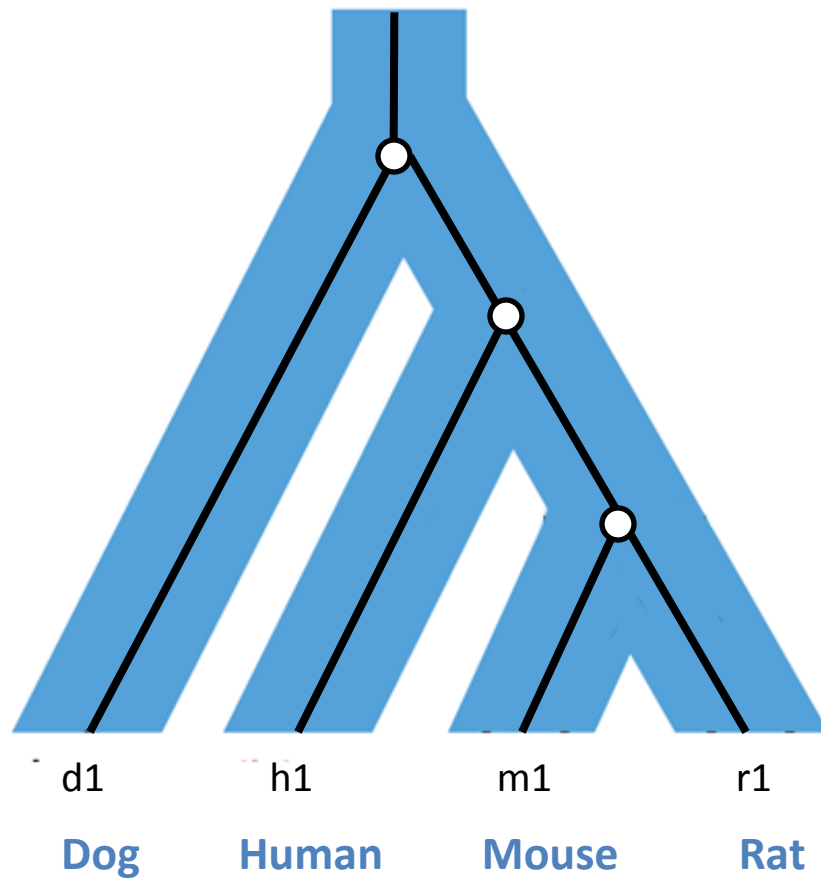


### Coalescence & Recombination

ARGHMM



# Gene family evolution

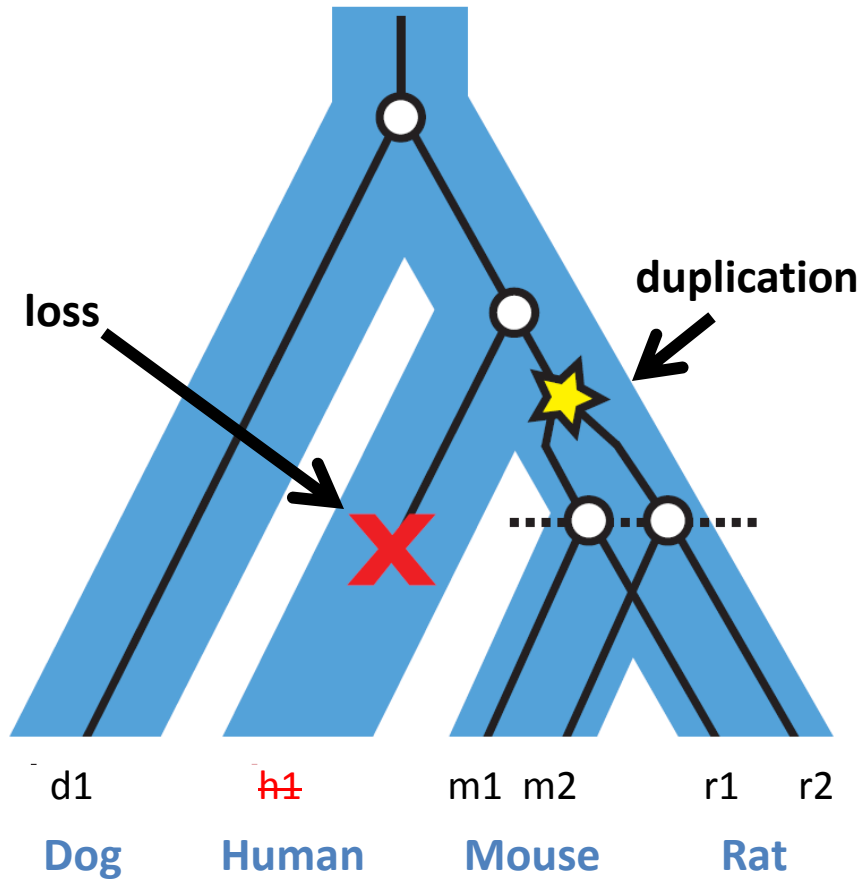


Species tree

Gene tree

d1	ATGCCT <b>GA</b> ACCCGTTCTC
h1	ATG <b>ACTGAT</b> CCAGTTCTC
m1	ATGCCTCCCC <b>AGTAGGC</b>
r1	ATGCCTCCCC <b>AGTAGGC</b>

# Gene family evolution



Species tree

Gene tree

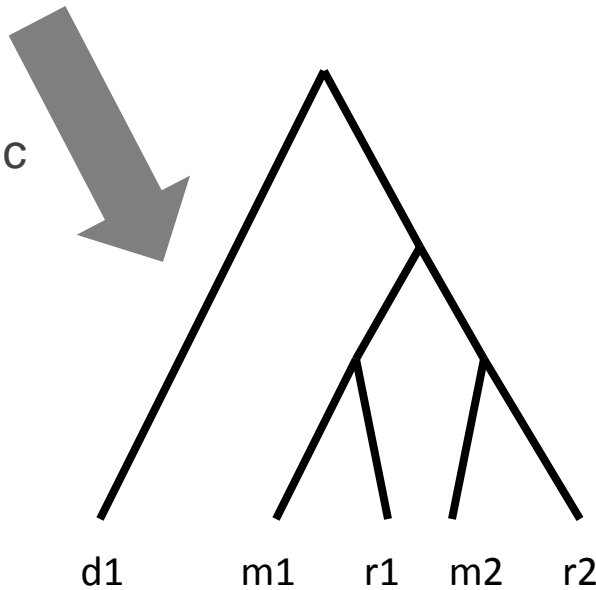
- Two genes are **orthologs** if their MRCA is a *speciation*: ○
- Two genes are **paralogs** if their MRCA is a *duplication*: ★

# Inferring events in a gene family

## Gene sequences

d1	ATGCCTGAACCCGTTCTC
m1	ATGACTGATCCAGTTCTC
r1	ATGCCTCCCCCAGTAGGC
m2	ATGCCTCGGGCAGTAGGC
r2	ATGCCTCCCCCAGTAGGC

Phylogenetic  
method



Gene tree

Obtained previously  
using a species tree  
method



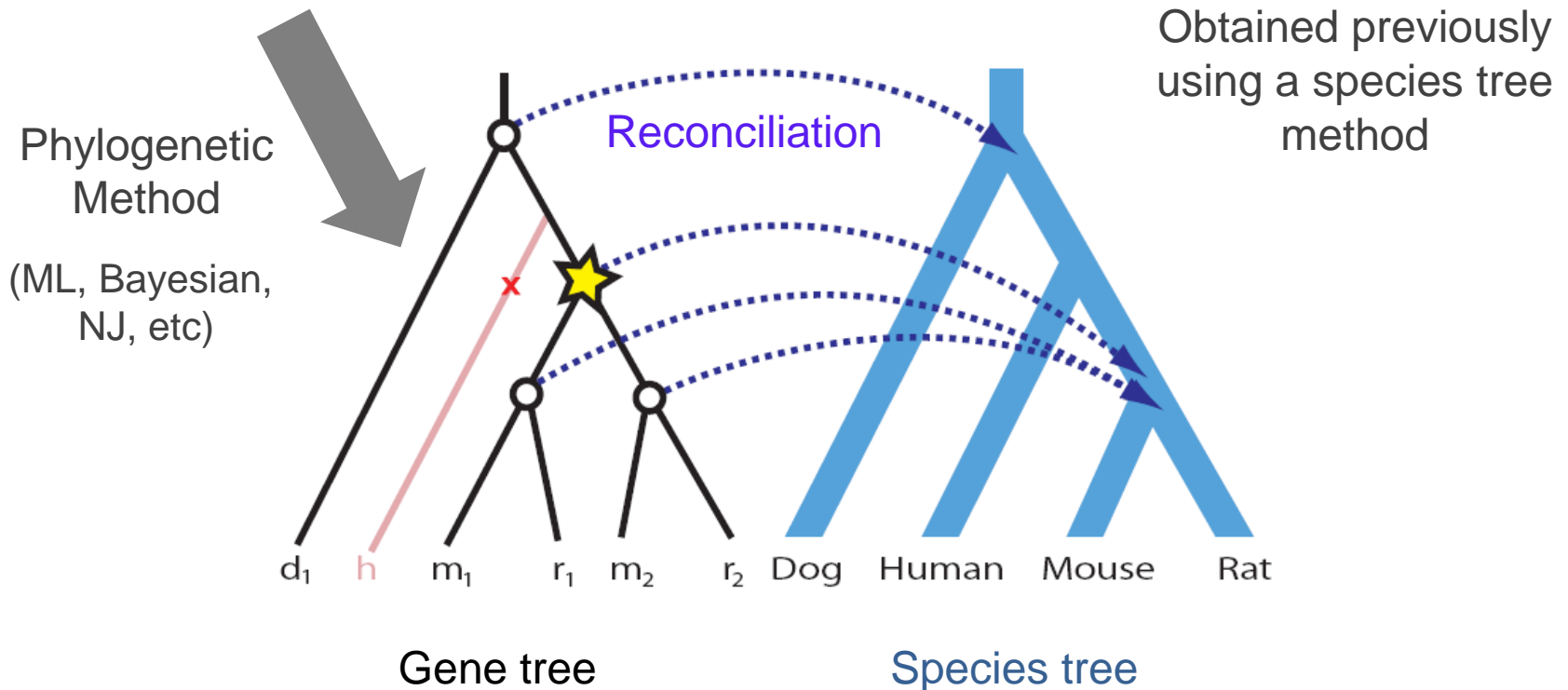
Species tree



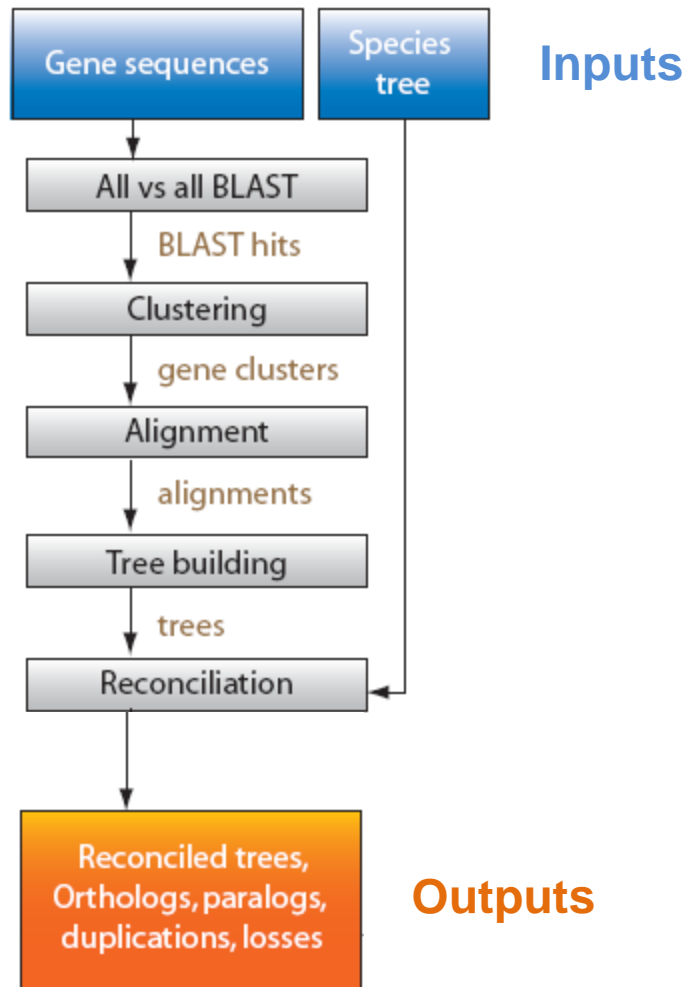
# Inferring events in a gene family

## Gene sequences

d1	ATGCCTGAACCCGTTCTC
m1	ATGACTGATCCAGTTCTC
r1	ATGCCTCCCCCAGTAGGC
m2	ATGCCTCGGGCAGTAGGC
r2	ATGCCTCCCCCAGTAGGC



# Phylogenomic Pipeline



## Discovery of functional elements in 12 *Drosophila* genomes using evolutionary signatures

Alexander Stark<sup>1,2\*</sup>, Michael F. Lin<sup>1,2\*</sup>, Pouya Kheradpour<sup>2\*</sup>, Jakob S. Pedersen<sup>3,4\*</sup>, Leopold Parts<sup>5,6</sup>, Joseph W. Carlson<sup>7</sup>, Madeline A. Crosby<sup>8</sup>, Matthew D. Rasmussen<sup>2</sup>, Sushmita Roy<sup>9</sup>, Ameya N. Deoras<sup>2</sup>, J. Graham Ruby<sup>10,11</sup>, Julius Brennecke<sup>12</sup>, Harvard FlyBase curators<sup>†</sup>, Berkeley *Drosophila* Genome Project<sup>†</sup>, Emily Hodges<sup>12</sup>, Angie S. Hinrichs<sup>4</sup>, Anat Caspi<sup>13</sup>, Benedict Paten<sup>4,5,14</sup>, Seung-Won Park<sup>15</sup>, Mira V. Han<sup>16</sup>, Morgan L. Maeder<sup>17</sup>, Benjamin J. Polansky<sup>17</sup>, Bryanne E. Robson<sup>17</sup>, Stein Aerts<sup>18,19</sup>, Jacques van Helden<sup>20</sup>, Bassem Hassan<sup>18,19</sup>, Donald G. Gilbert<sup>21</sup>, Deborah A. Eastman<sup>17</sup>, Michael Rice<sup>22</sup>, Michael Weir<sup>23</sup>, Matthew W. Hahn<sup>16</sup>, Yongkyu Park<sup>15</sup>, Colin N. Dewey<sup>24</sup>, Lior Pachter<sup>25,26</sup>, W. James Kent<sup>4</sup>, David Haussler<sup>4</sup>, Eric C. Lai<sup>27</sup>, David P. Bartel<sup>10,11</sup>, Gregory J. Hannon<sup>12</sup>, Thomas C. Kaufman<sup>21</sup>, Michael B. Eisen<sup>28,29</sup>, Andrew G. Clark<sup>30</sup>, Douglas Smith<sup>31</sup>, Susan E. Celniker<sup>7</sup>, William M. Gelbart<sup>8,32</sup> & Manolis Kellis<sup>1,2</sup>

Nature. 2007

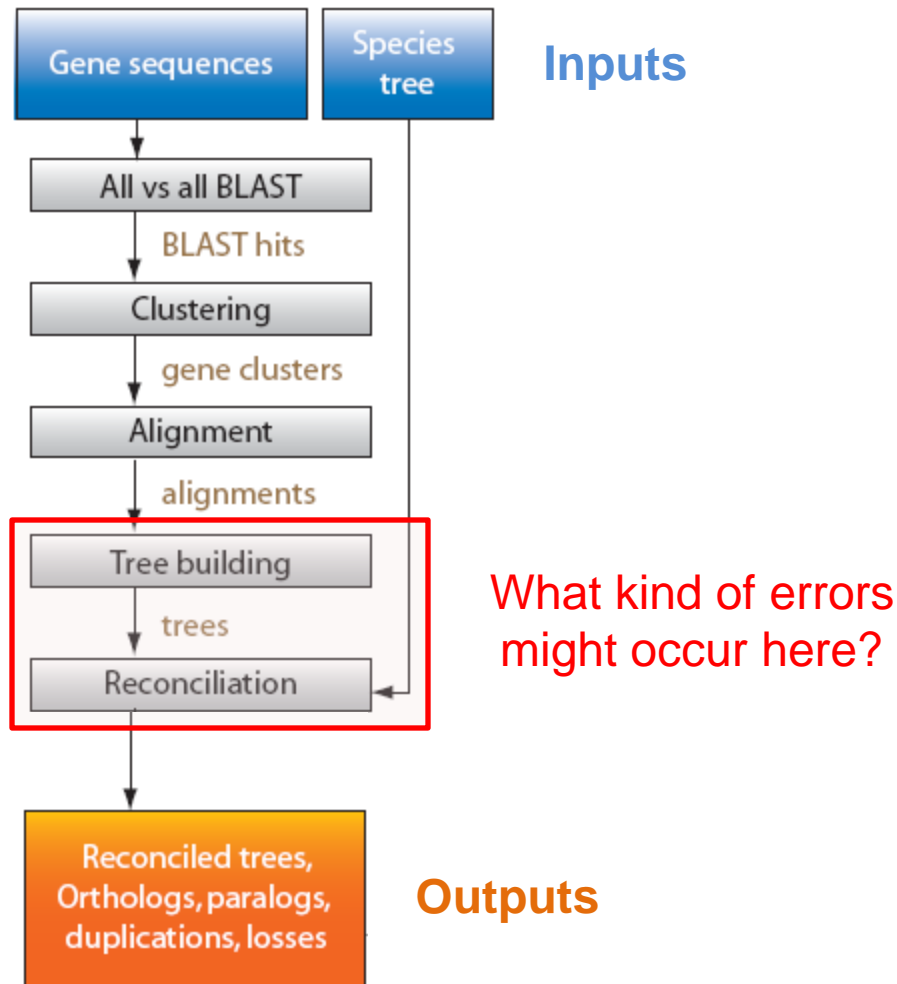
*Drosophila* Genomes/Methods

## Accurate gene-tree reconstruction by learning gene- and species-specific substitution rates across multiple complete genomes

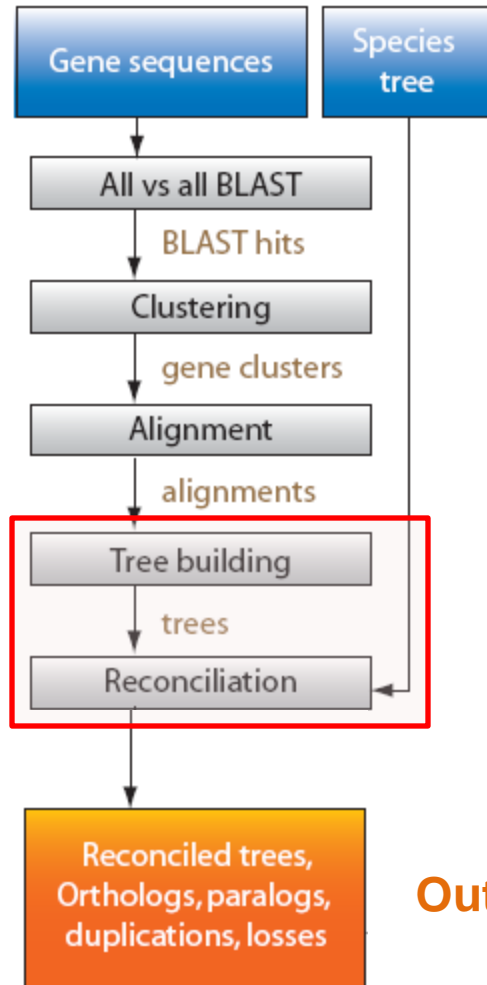
Matthew D. Rasmussen<sup>1</sup> and Manolis Kellis<sup>1,2,3</sup>

Genome Research, 2007

# Phylogenomic Pipeline



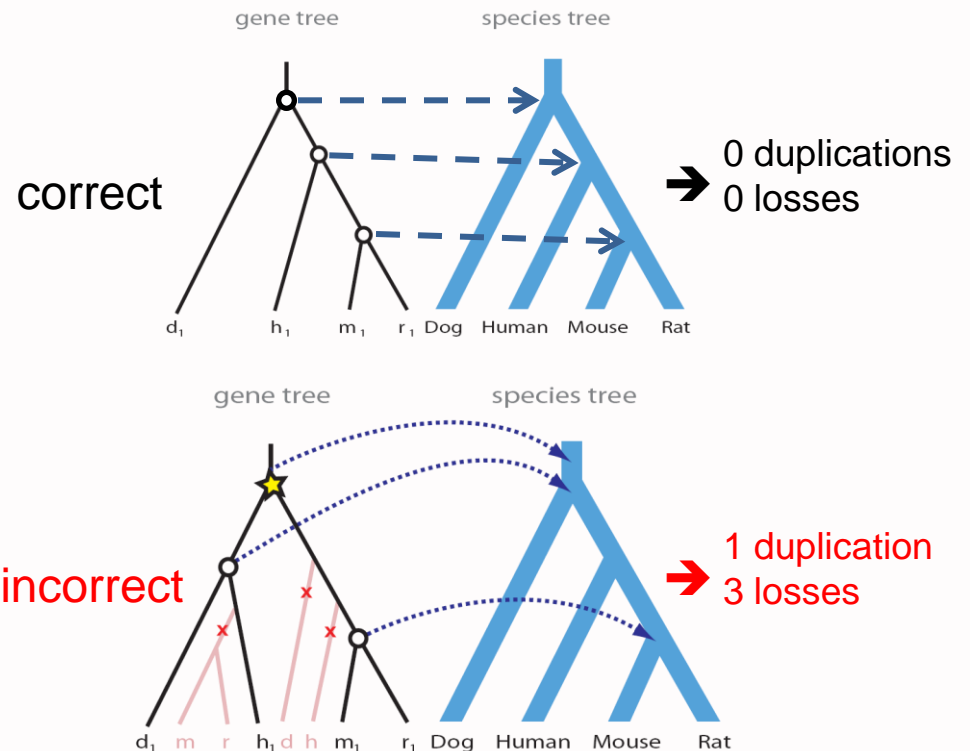
# Phylogenomic Pipeline



Inputs

Outputs

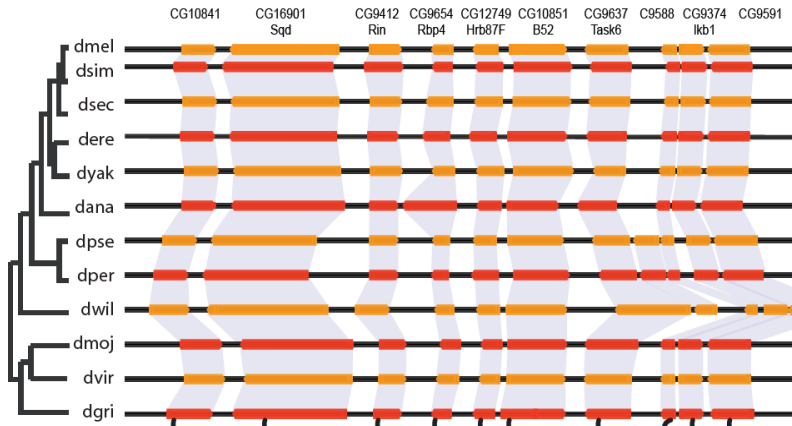
Tree error leads to dup/loss error



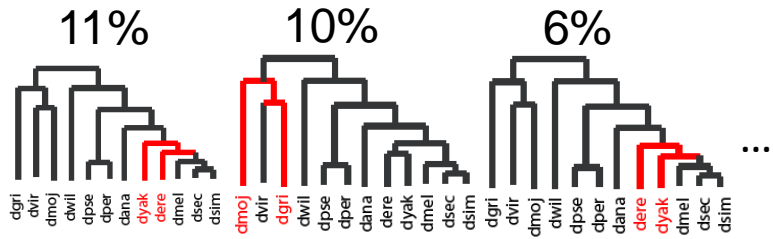
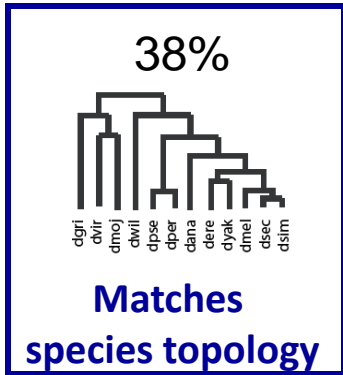
In general, gene tree errors lead to over estimating dup-loss by 2-3 fold [Hahn 2007]

# Existing phylogenetic methods are not accurate enough for phylogenomics

~5000 syntenic one-to-one orthologs



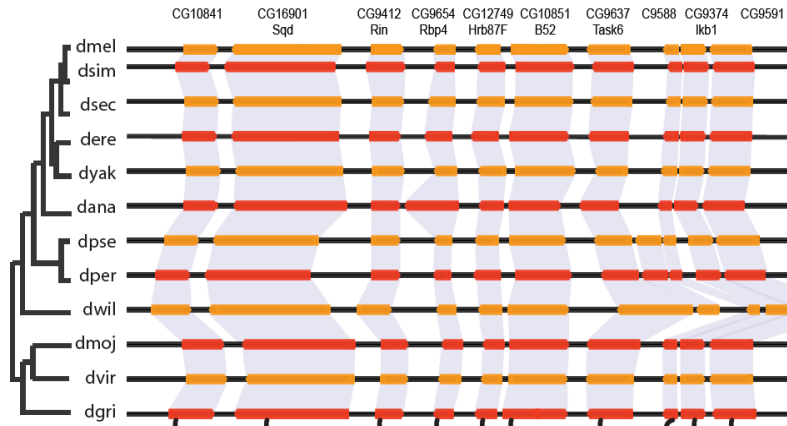
Build trees



316 incongruent topologies

# Existing phylogenetic methods are not accurate enough for phylogenomics

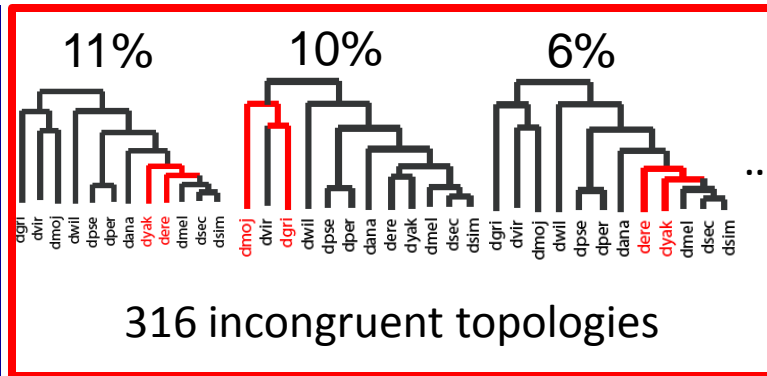
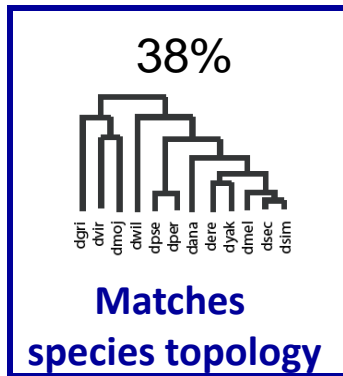
~5000 syntenic one-to-one orthologs



Low information → Incongruence

- Shorter sequences (<900bp)
- Too conserved (>90% identity)
- Too diverged (<30% identity)
- Rarely supported by LRT (<5.7%)

Build trees



# Phylogenomics needs a new approach

## Average gene alignment contains too little information for phylogenomic analysis

- Existing algorithms ignore species (species unaware, uses only sequence)

## Our approach:

- Use species tree to inform the gene tree reconstruction (species aware)

## New Methods:

*Drosophila* Genomes/Methods

Accurate gene-tree reconstruction by learning gene- and species-specific substitution rates across multiple complete genomes

Matthew D. Rasmussen<sup>1</sup> and Manolis Kellis<sup>1,2,3</sup> **Genome Research, 2007**

### A Bayesian Approach for Fast and Accurate Gene Tree Reconstruction

Matthew D. Rasmussen<sup>\*,1</sup> and Manolis Kellis<sup>\*,1,2</sup>

<sup>1</sup>Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology

<sup>2</sup>Broad Institute of MIT and Harvard

**Molecular Biology & Evolution, 2011**

TreeFix: Statistically Informed Gene Tree Error Correction Using Species Trees

YI-CHIEH WU<sup>1</sup>, MATTHEW D. RASMUSSEN<sup>2</sup>, MUKUL S. BANSAL<sup>1</sup>, AND MANOLIS KELLIS<sup>1,3,\*</sup>

**Systematic Biology, 2013**

## SPIDIR (Species Informed Distance-based Reconstruction)

- Maximum Likelihood

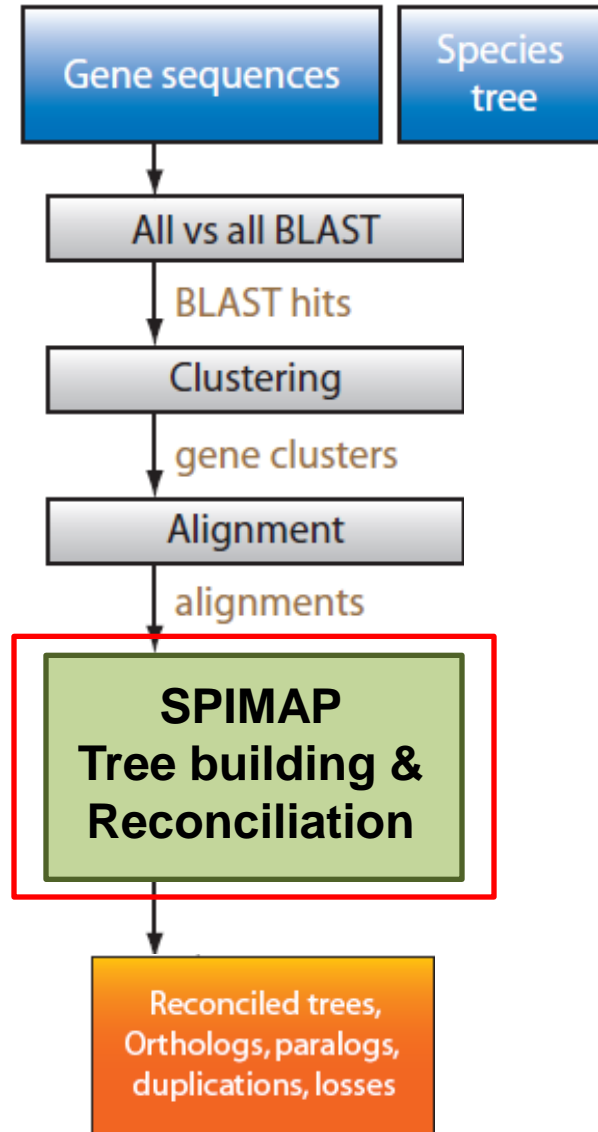
## SPIMAP (Species Informed Maximum A Posteriori)

- Bayesian

## TreeFix

- Fewer parameters, Likelihood Ratio

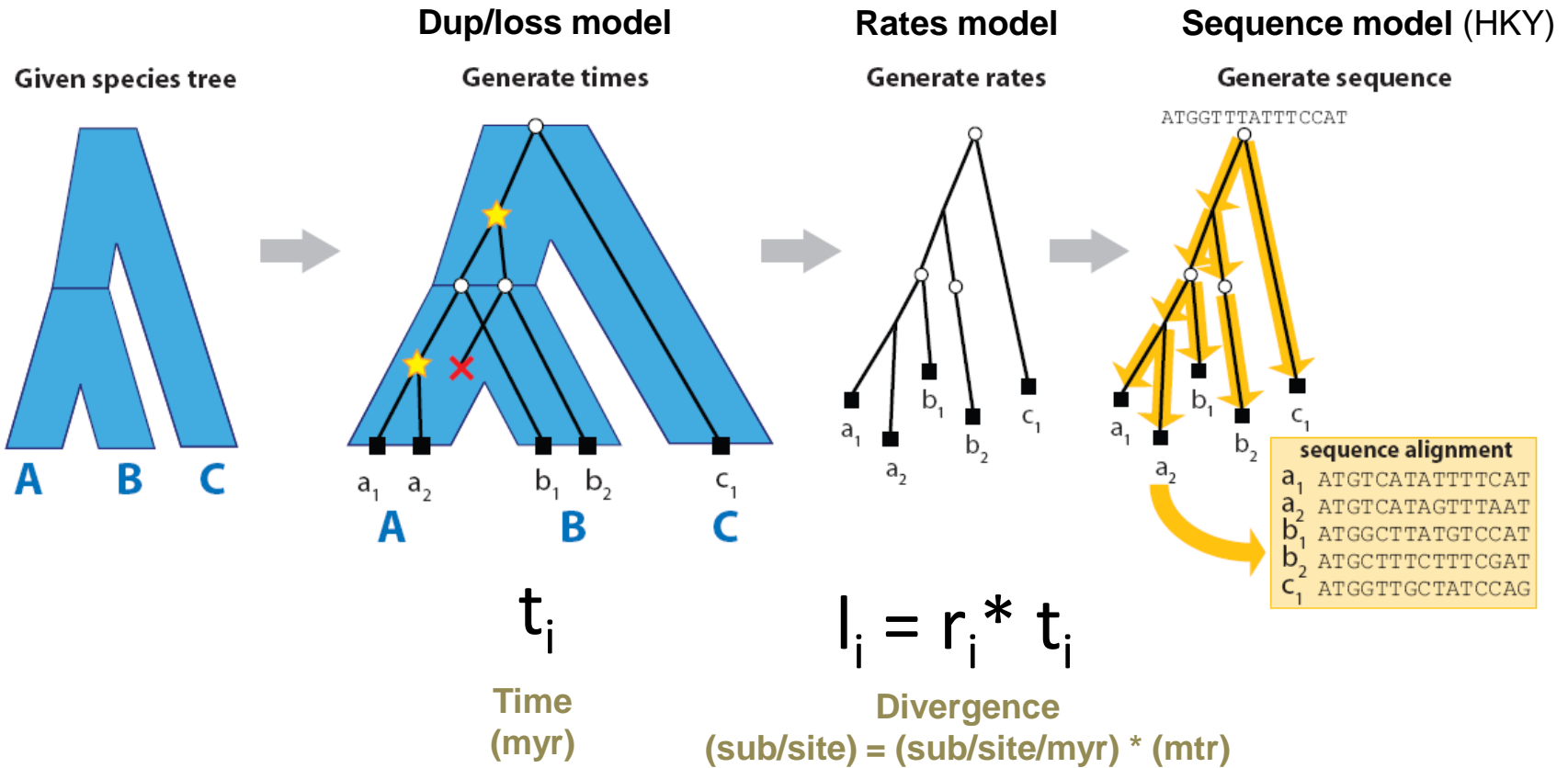
# SPIMAP's Phylogenomic Pipeline





# SPIMAP (Species Informed Maximum A Posteriori)

## Generative model



# Reconstruction using SPIMAP model

We find the maximum *a posteriori* tree

$$\hat{\mathbf{I}}, \hat{T}, \hat{R} = \underset{\mathbf{I}, T, R}{\operatorname{argmax}} P(\mathbf{I}, T, R | \mathbf{D}, \theta) = \underset{\mathbf{I}, T, R}{\operatorname{argmax}} \underbrace{P(\mathbf{D} | \mathbf{I}, T)}_{\text{likelihood}} \underbrace{P(\mathbf{I} | T, R, \theta) P(T, R | \theta)}_{\text{prior}}$$

**Max posterior**   **Model parameters**   **Likelihood**   **Branch prior**   **Topology prior**  
Sequence model   Rates model   Dup/loss model

$\mathbf{I}$  = vector of branch lengths

$T$  = gene tree topology

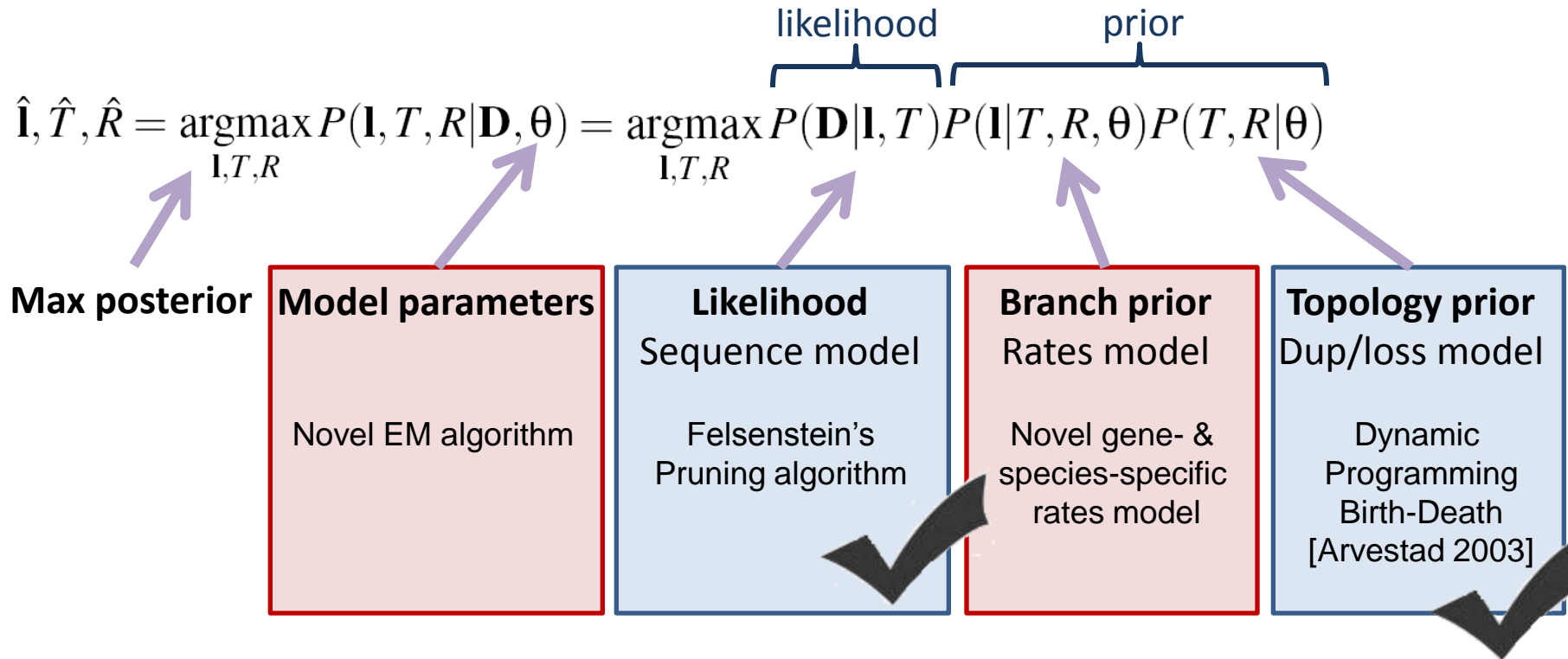
$R$  = reconciliation mapping

$\mathbf{D}$  = alignment data

$\theta$  = (species tree  $S$ , and other model parameters)

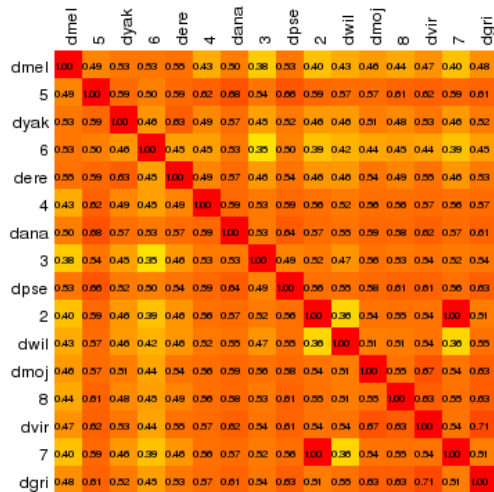
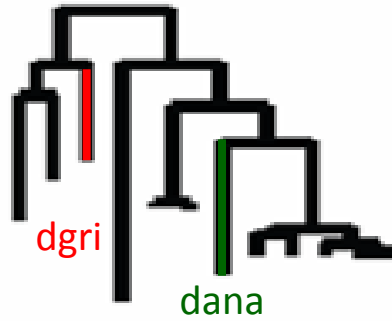
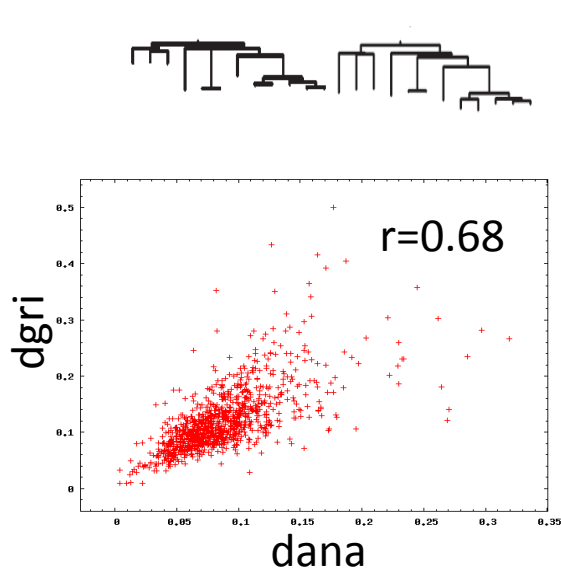
# Reconstruction using SPIMAP model

We find the maximum *a posteriori* tree

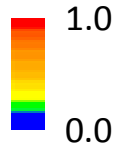


# Rates model: branch length correlations

## Absolute lengths

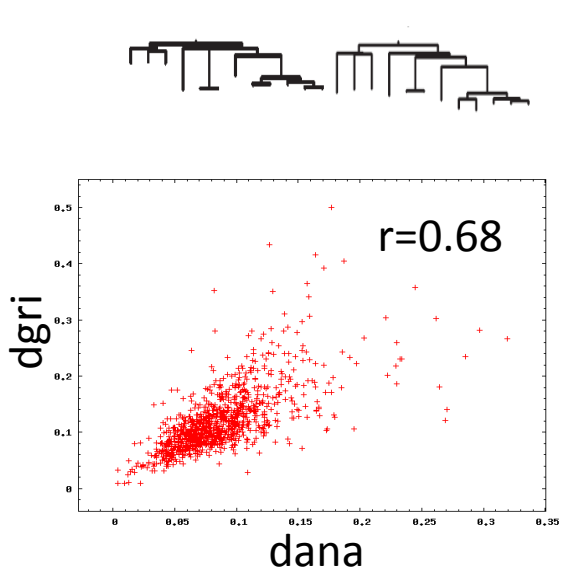


average  
 $r = 0.56$

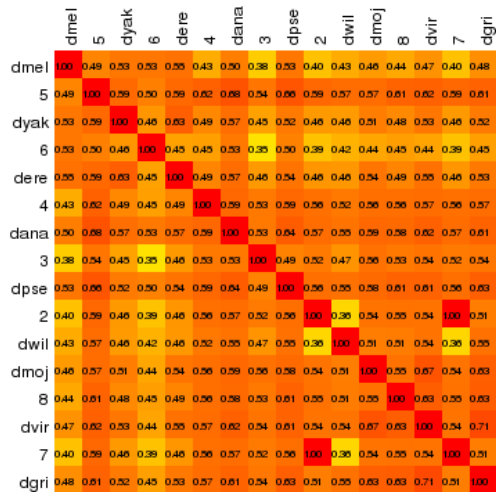
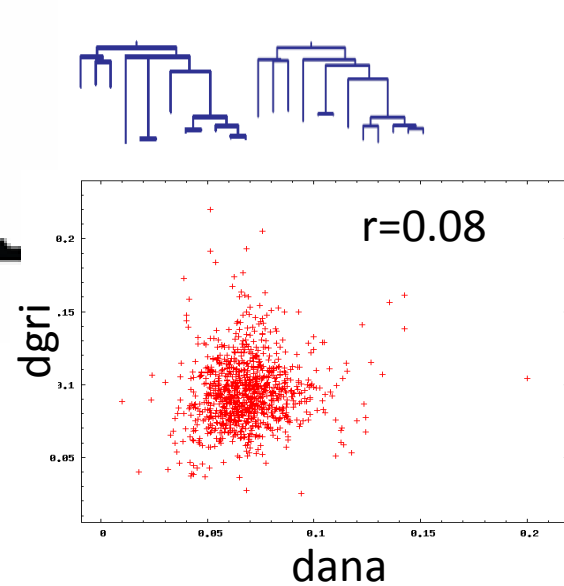


# Rates model: branch length correlations

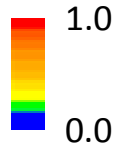
Absolute lengths



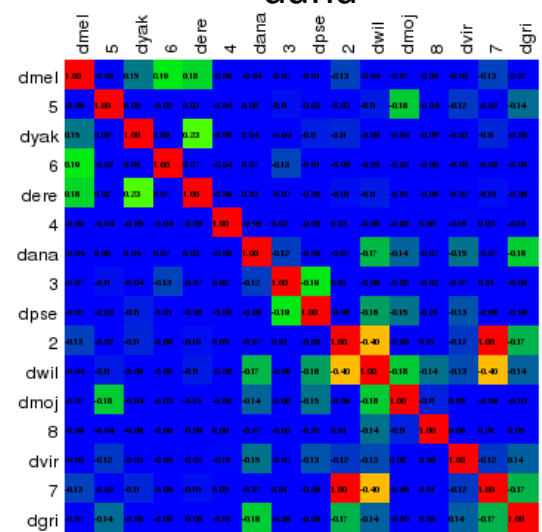
Relative lengths



average  
 $r = 0.56$

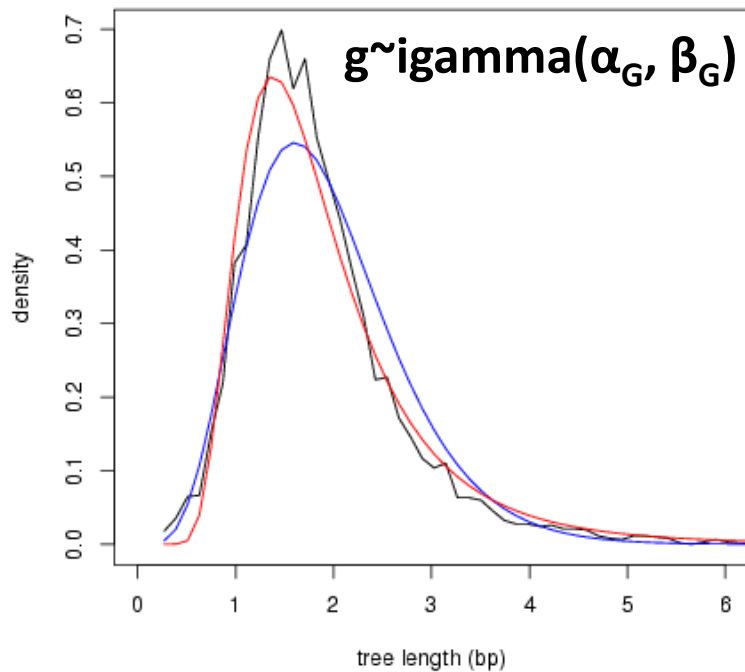


average  
 $r = 0.09$

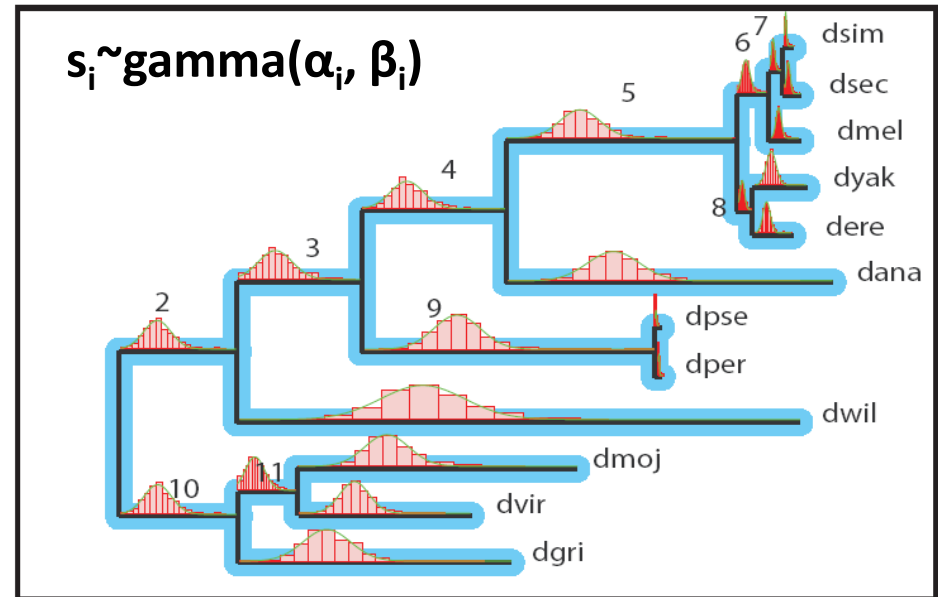


# Rates model: rate distributions

Gene rate distribution  
(inverse-gamma)



Species rate distributions  
(gamma)



$$l_i = r_i * t_i$$

$$r_i = g * s_i$$

$l_i$  (sub/site)

$r_i$  (sub/site/myr)

$t_i$  (myr)

# Method Evaluation

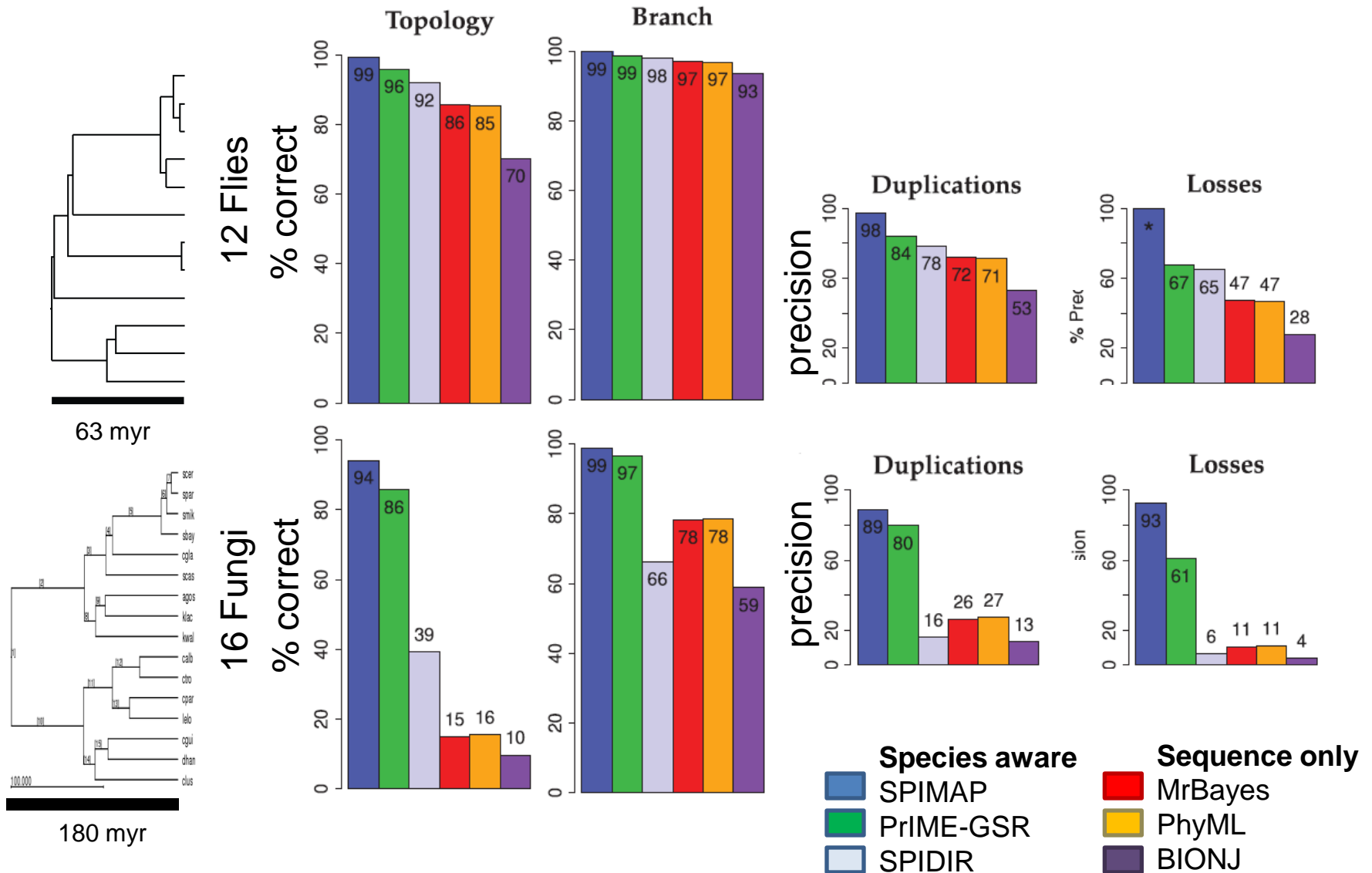
## “Sequence-only” methods

- PHYML (ML) [Guidon 2003](#)
- RAxML (ML) [Stamatakis 2006](#)
- MrBayes (Bayesian) [Ronquist 2003](#)
- BIONJ (NJ) [Gascuel 1997](#)

## Species-aware methods:

- SPIMAP (Bayesian) [Rasmussen, Kellis 2010](#)
- SPIDIR (dist. ML) [Rasmussen, Kellis 2007](#)
- SYNERGY (dist/parsimony) [Wapinski 2007](#)
- PrIME-GSR (Bayesian, i.i.d.) [Akerborg 2009](#)

# Evaluation on fly & fungi simulated data





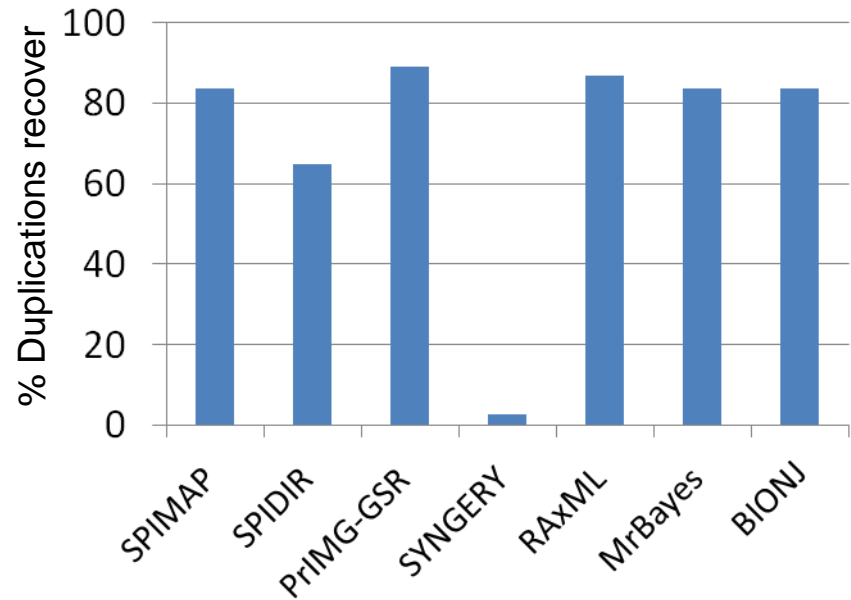
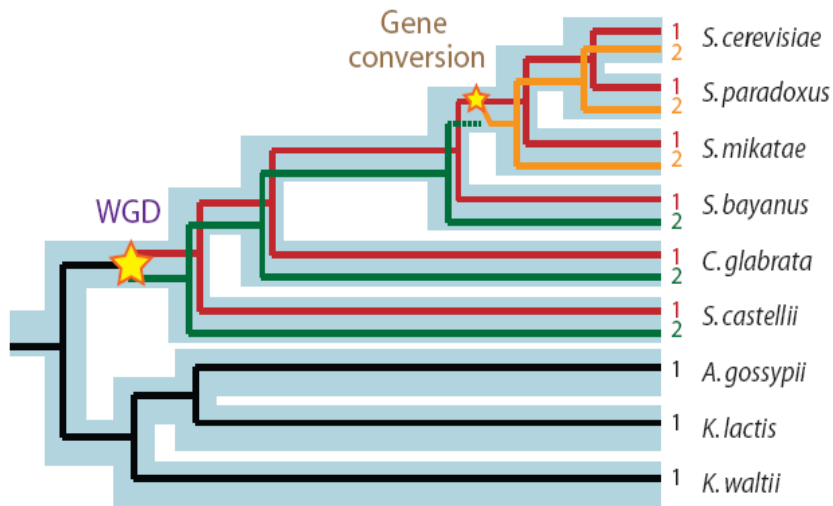
# Evaluation on 16 fungal species

	Program	Orthologs	# Dup	# Loss	Run time
Species aware	SPIMAP	96.2%	5,541	10,884	1.0 min
	SPIDIR	83.3%	10,177	33,550	2.2 min
	PrIMG-GSR	90.7%	7,951	21,099	53.1 min
	SYNGERY	99.2%	4,604	8,179	---
Sequence only	RAxML	63.3%	21,485	65,392	18.4 s
	MrBayes	64.2%	21,307	65,238	43.2 s
	BIONJ	60.4%	22,396	71,231	0.5 s

(1) Species-aware methods recover far more syntenic orthologs

(2) Sequence-only methods infer far more (likely spurious) duplications-losses

### (3) High recovery of duplications due to gene conversion



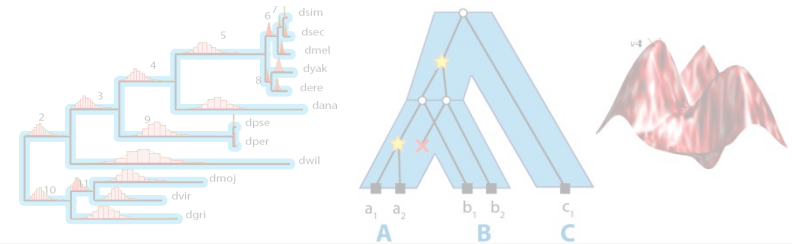
Methods that **balance** sequence information vs dup-loss and synteny information perform best

# Developing models of evolution

## Combining different evolutionary events & processes

### Duplication, Loss & Substitution

SPIDIR (GR 2007), SPIMAP (MBE 2010),  
TreeFix (SystBio 2012)

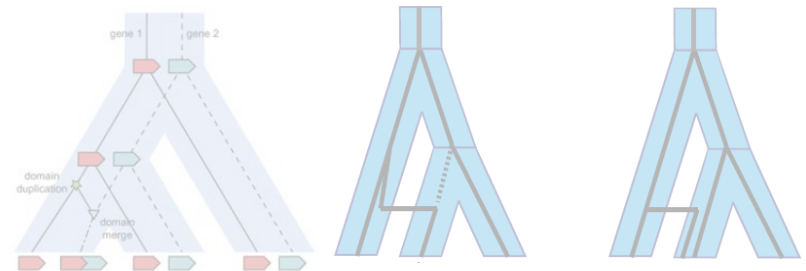


### Fusion & Fission

STAR-MP (MBE 2011)

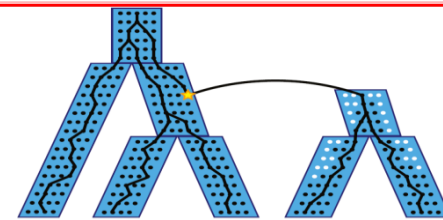
### Additive & Replacing transfer

(MBE 2012)



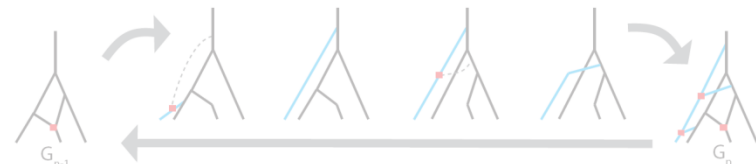
### Duplication, Loss & Coalescence

DLCoal (GR 2012)



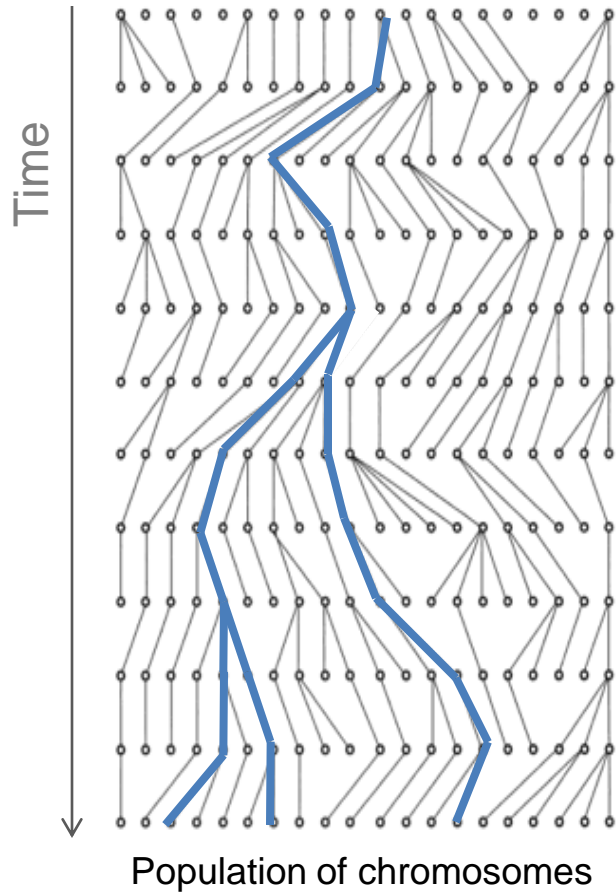
### Coalescence & Recombination

ARGHMM

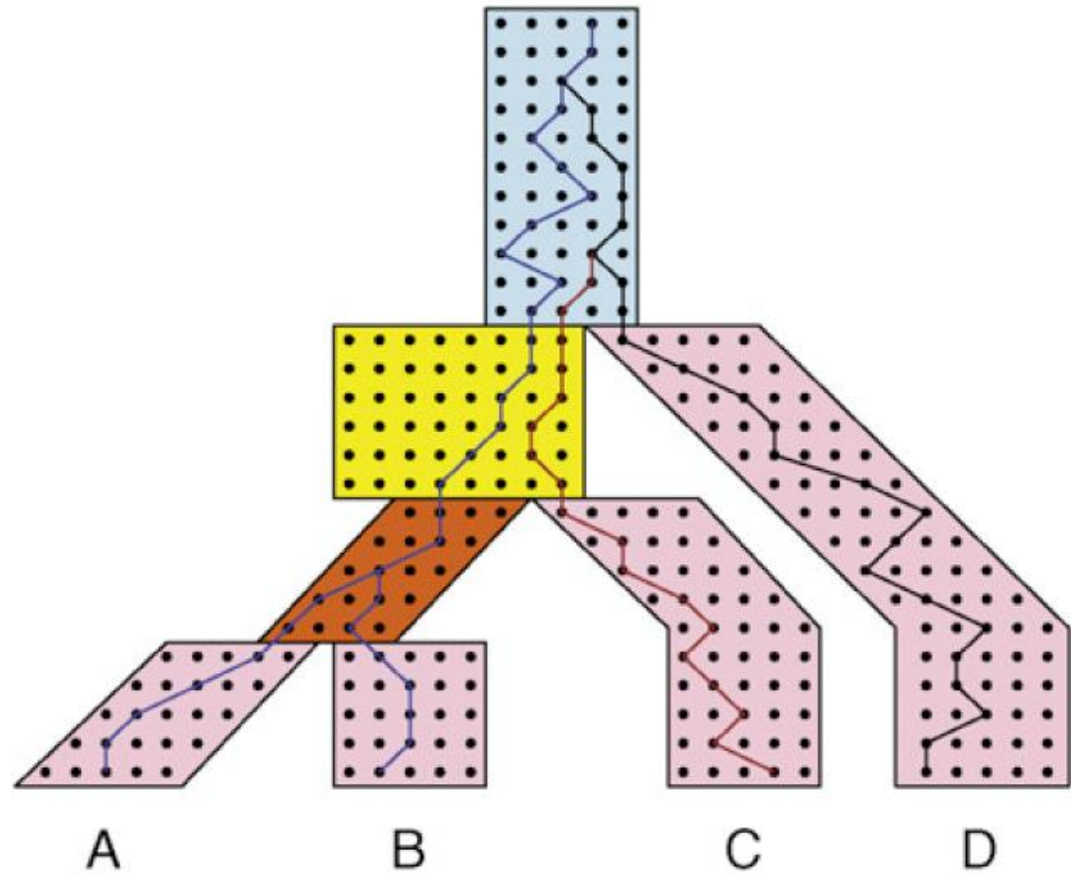


# Modeling drift with the coalescent

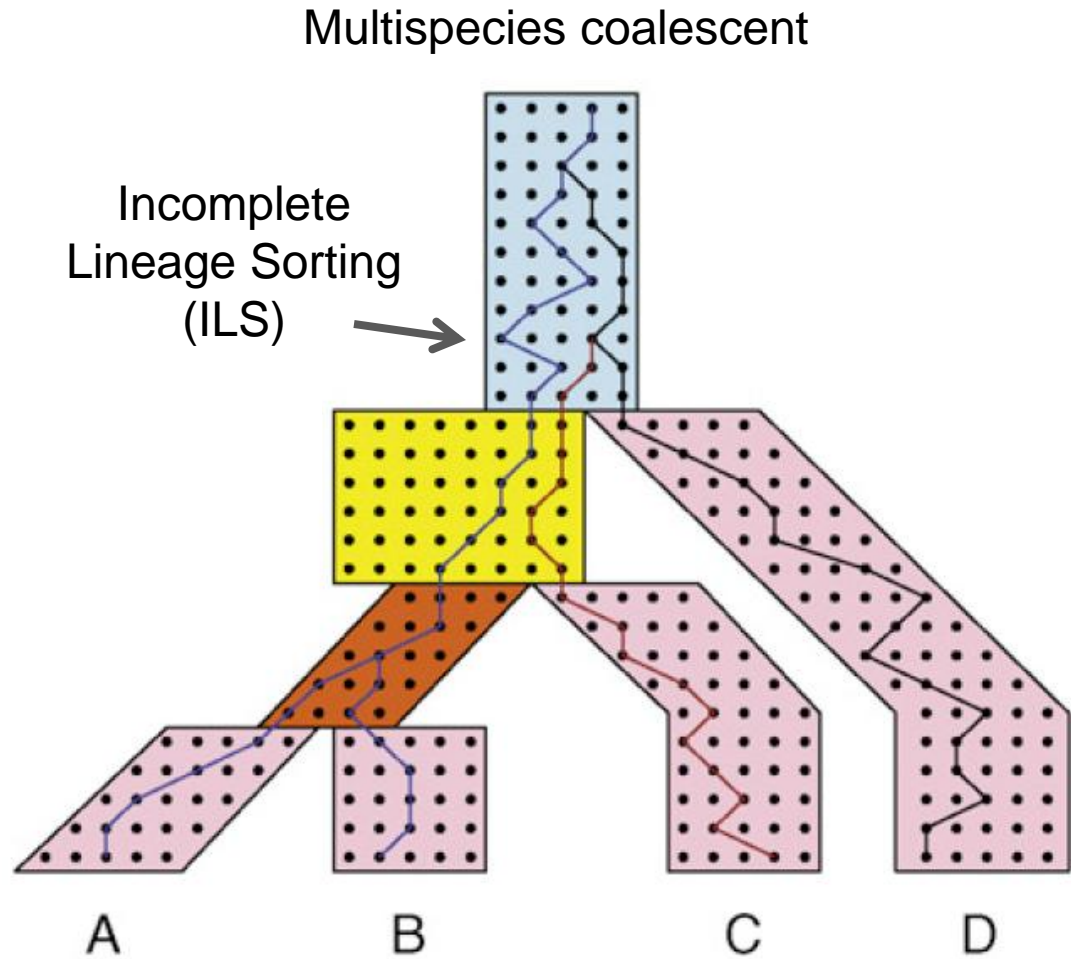
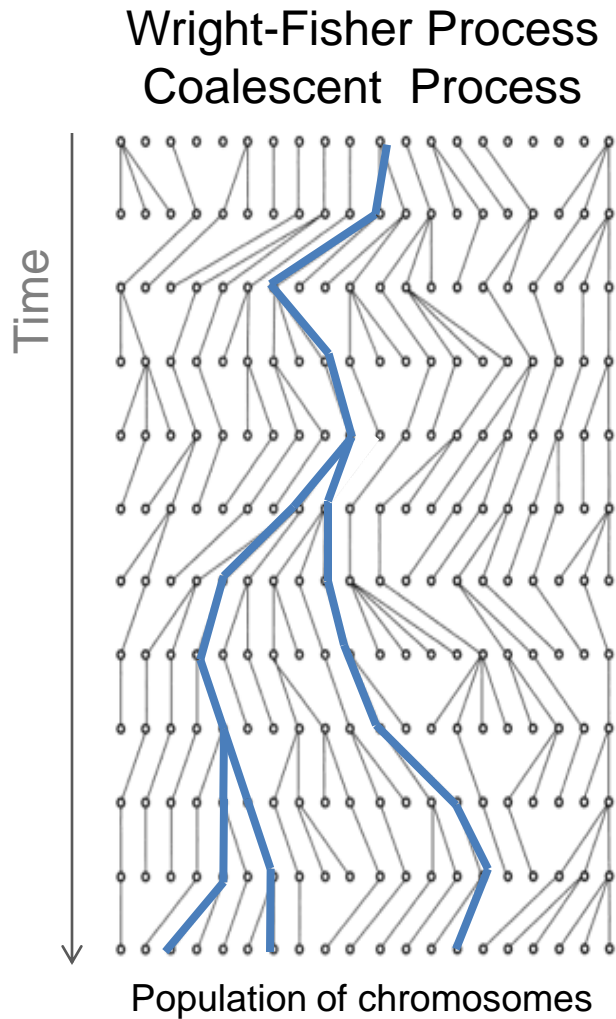
Wright-Fisher Process  
Coalescent Process



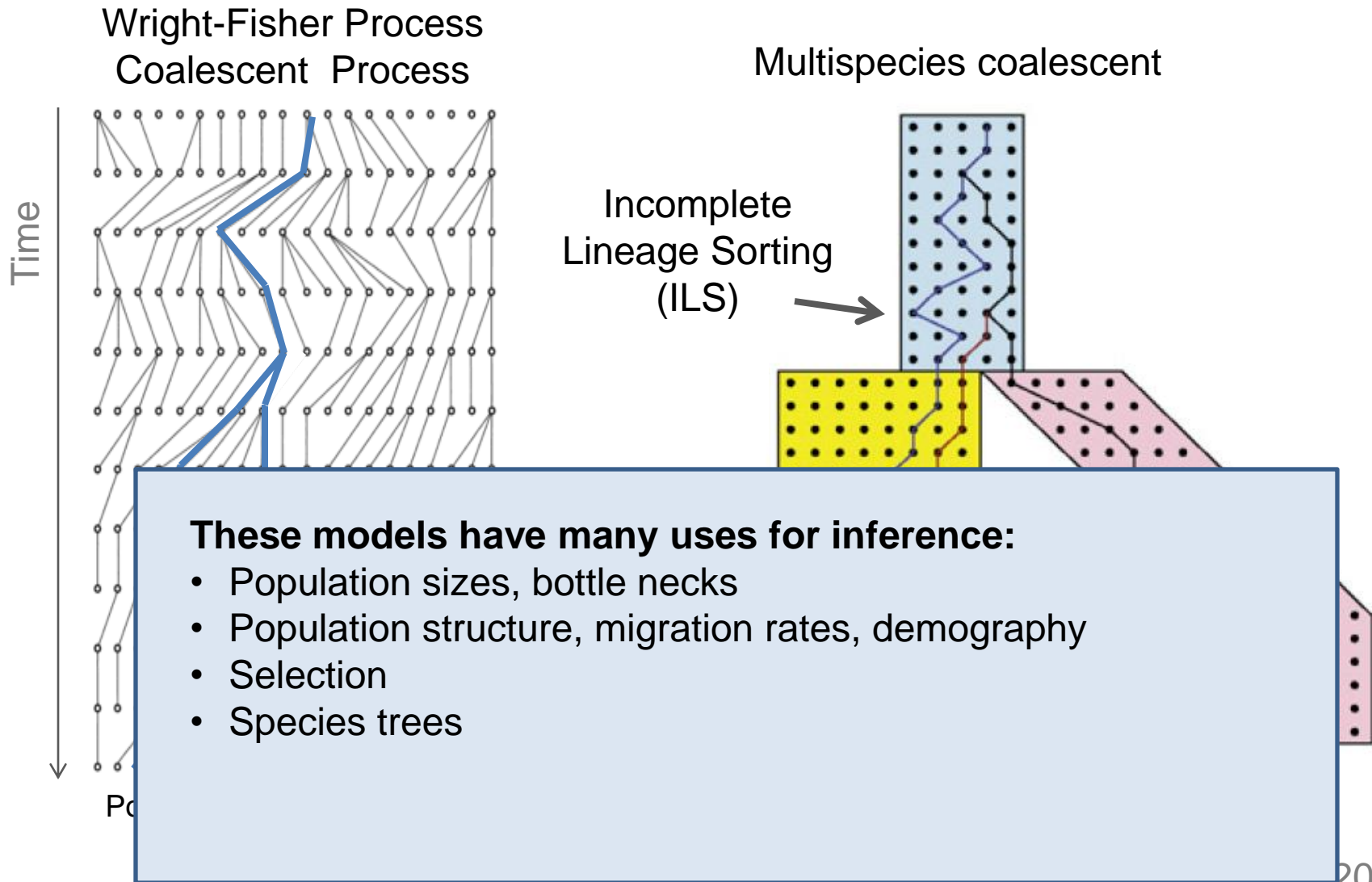
Multispecies coalescent



# Modeling drift with the coalescent

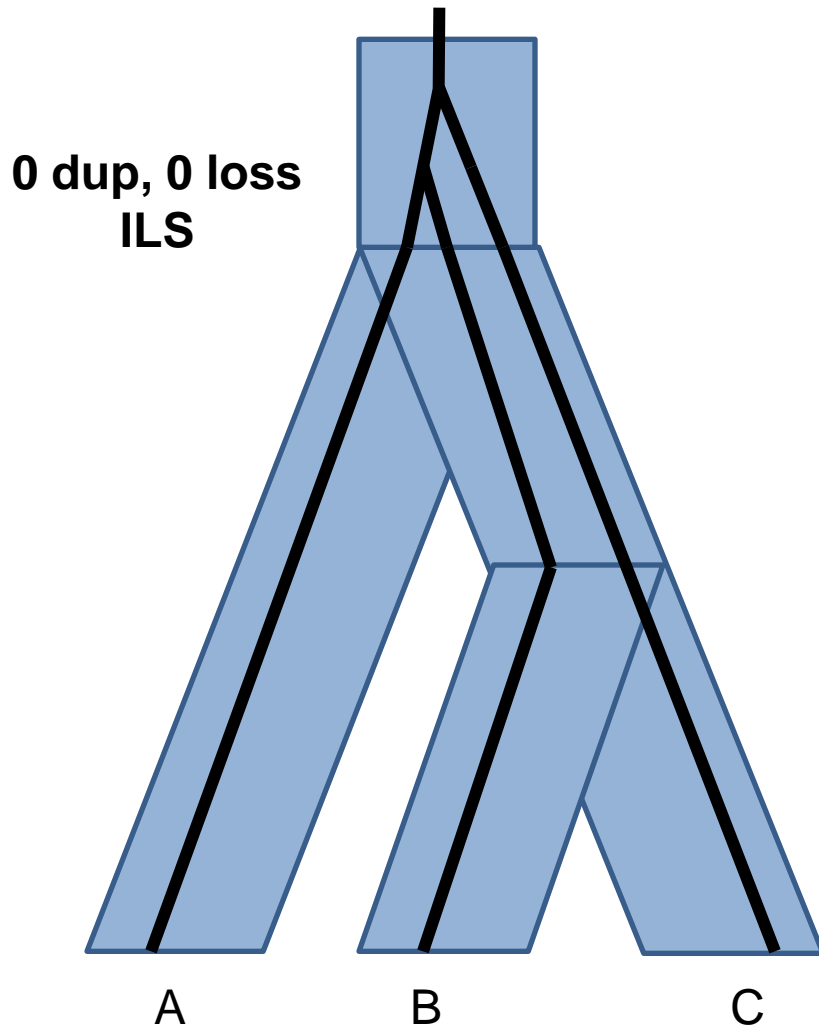


# Modeling drift with the coalescent

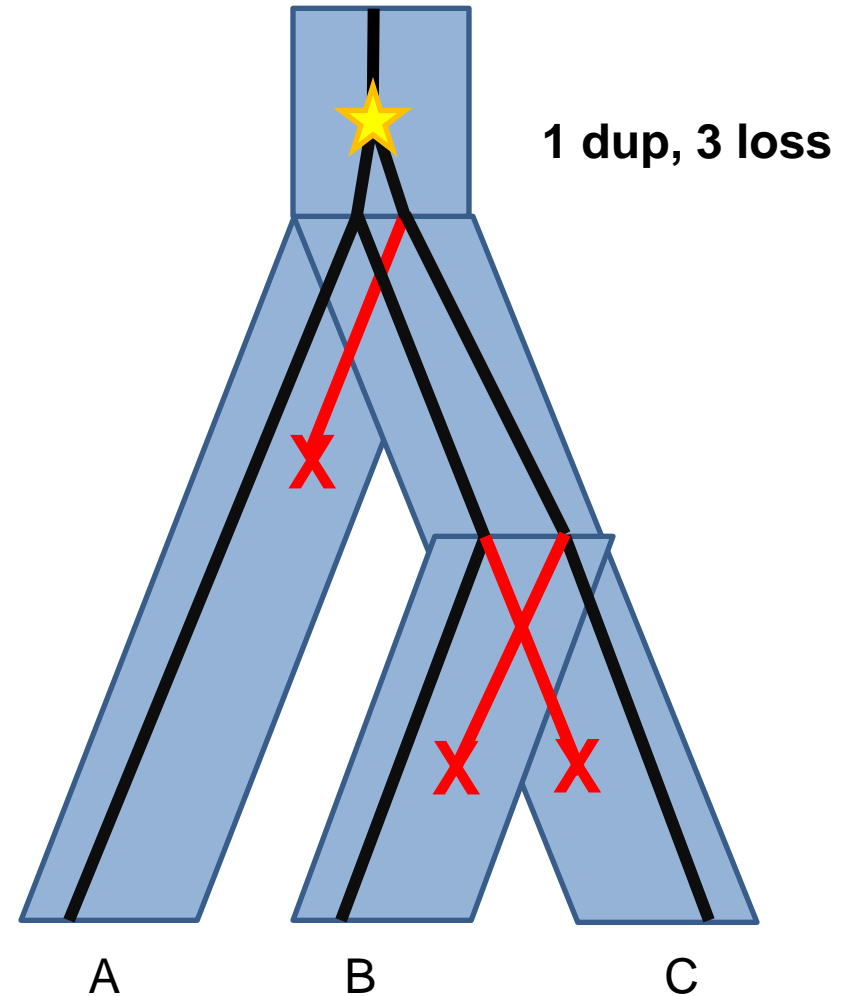


# The interpretation of a gene tree depends on your model

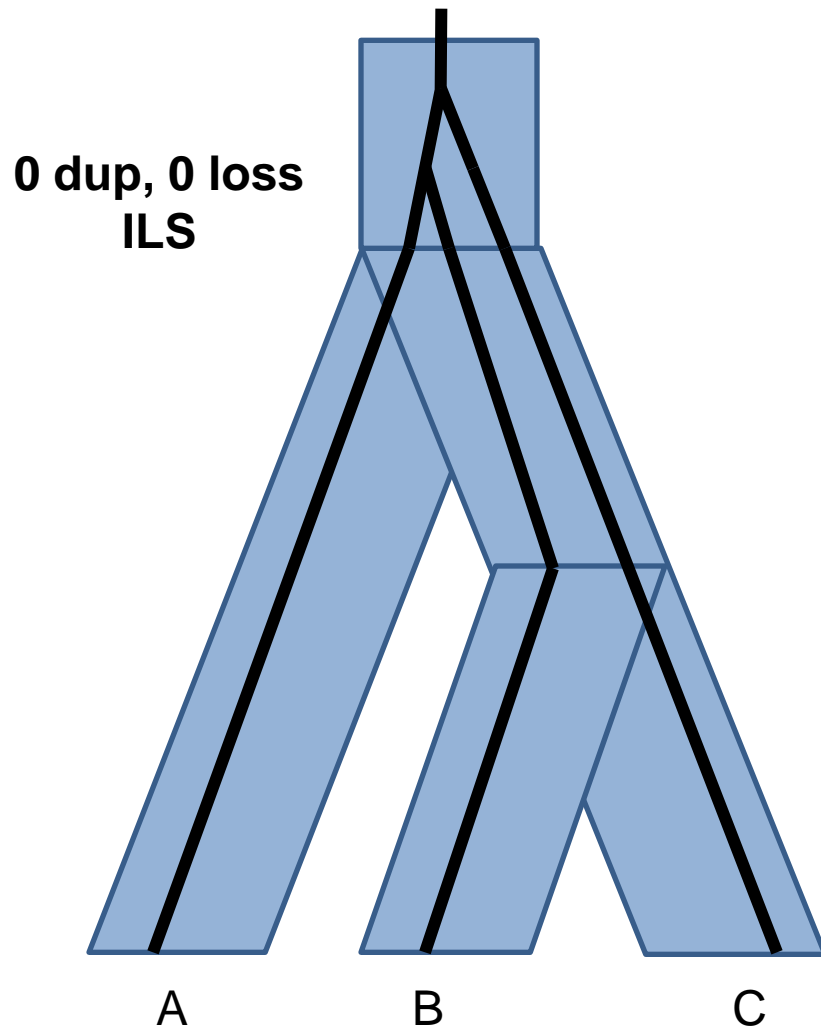
## Coalescence



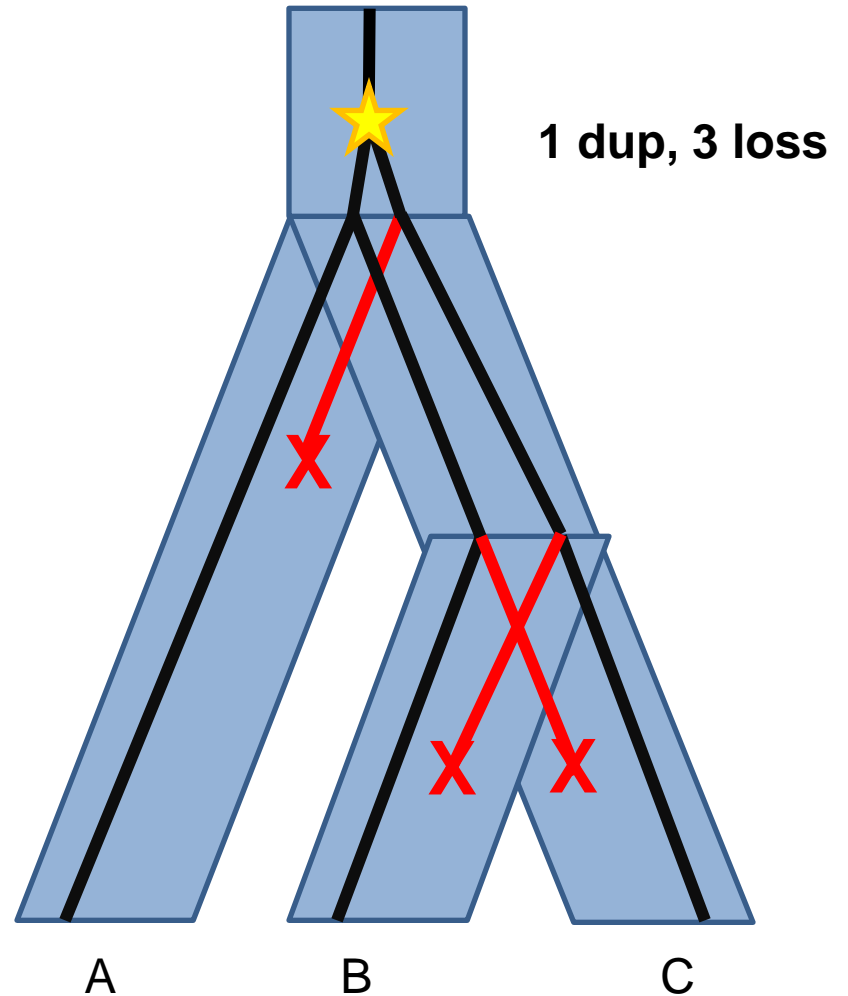
## Duplication & Loss



# Unified model could capture both

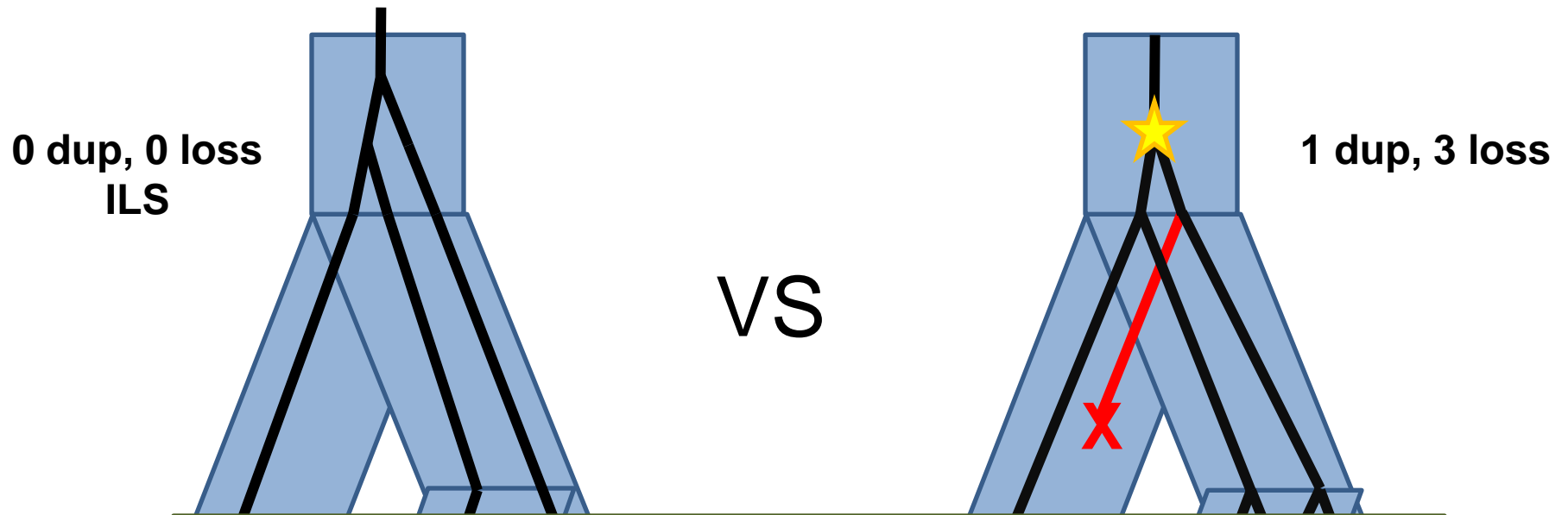


VS





# Unified model could capture both

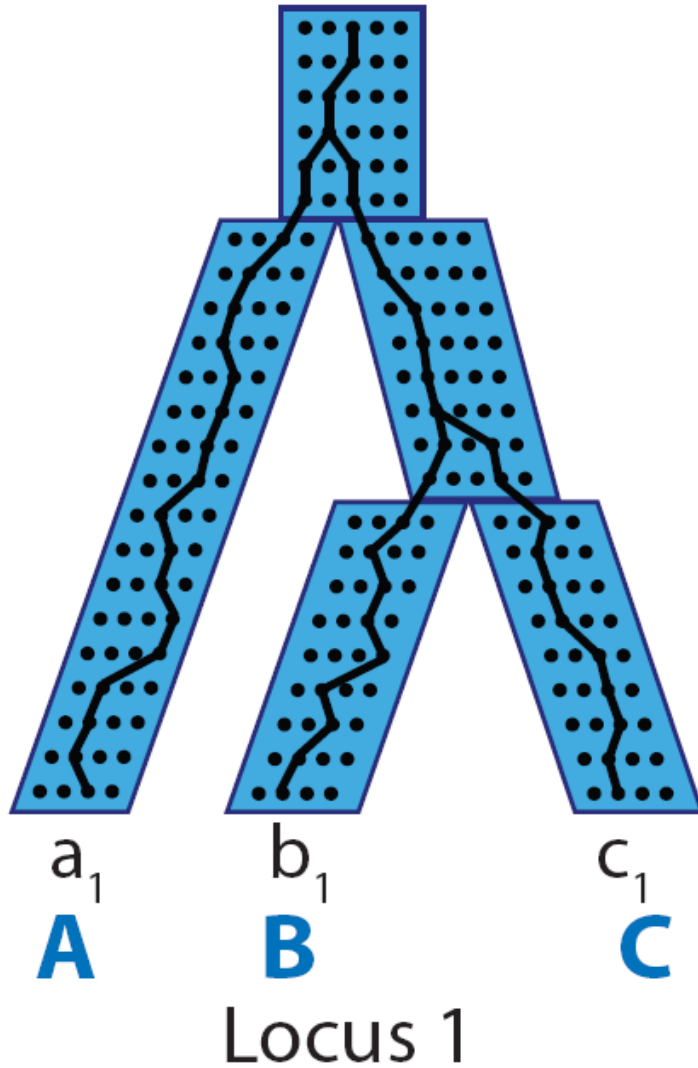


**Unified model would lead to several algorithms**

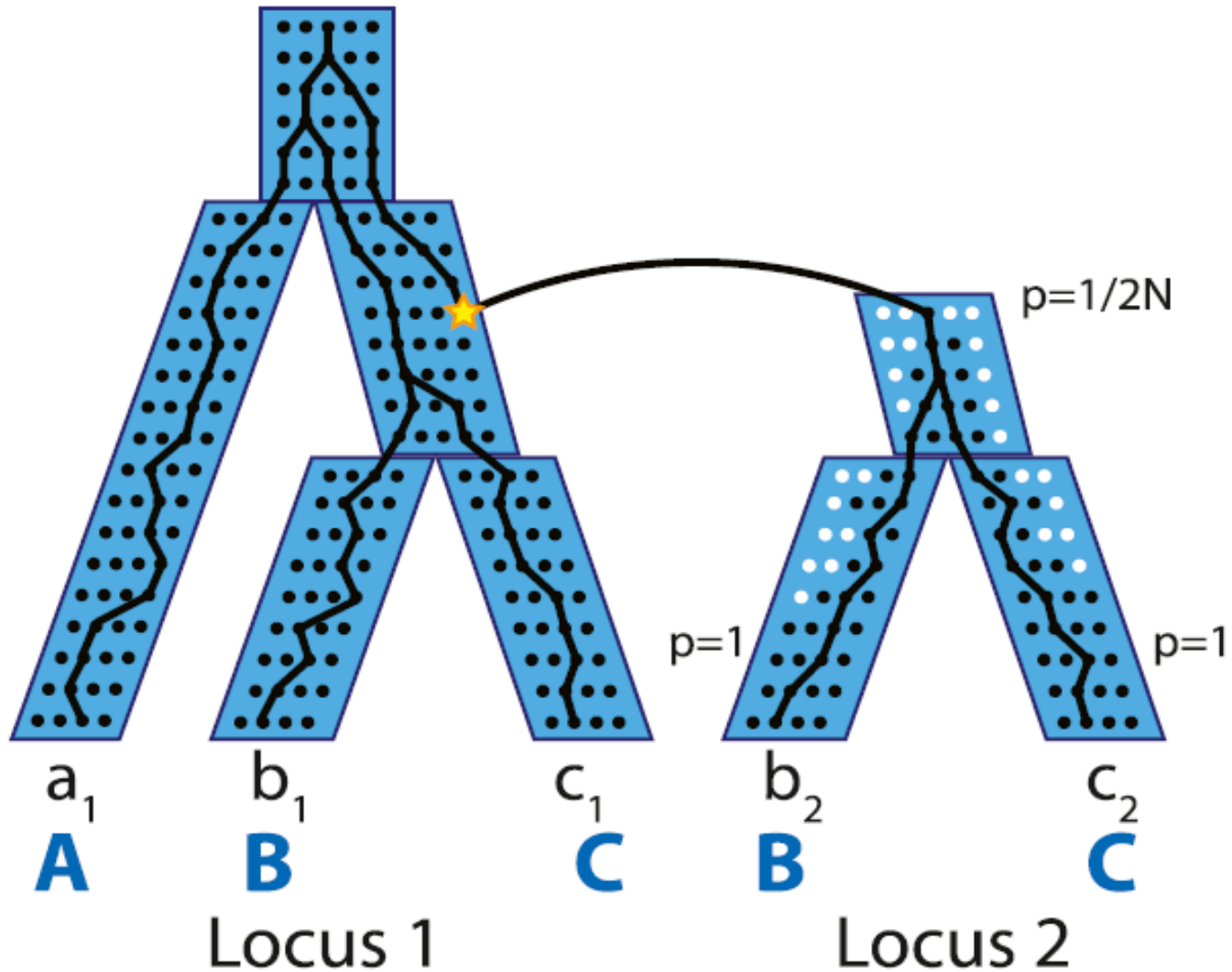
- Reconciliation
- Gene tree reconstruction
- Species tree reconstruction
- Population genetics applications

A

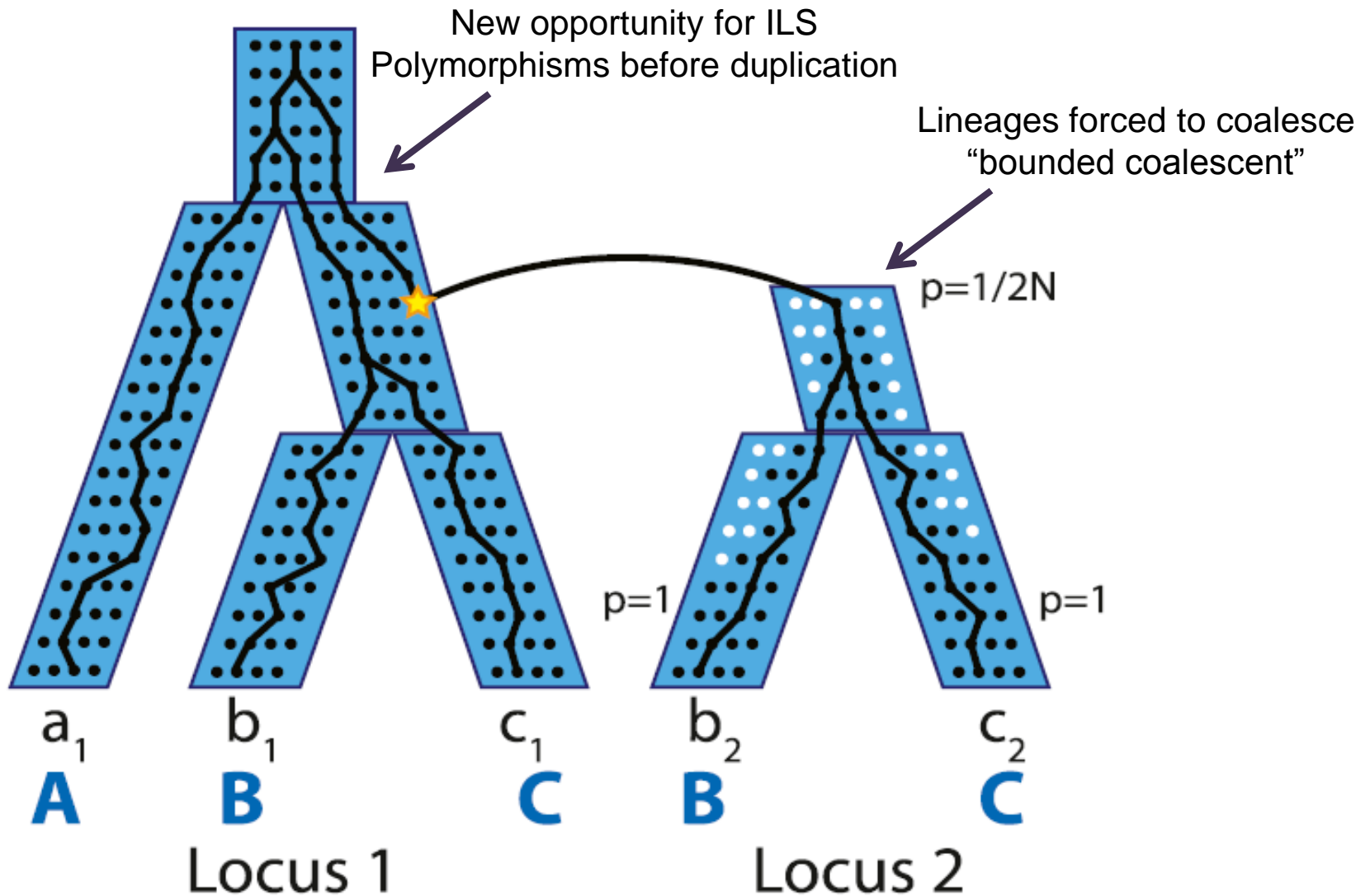
# Example of duplication in a population



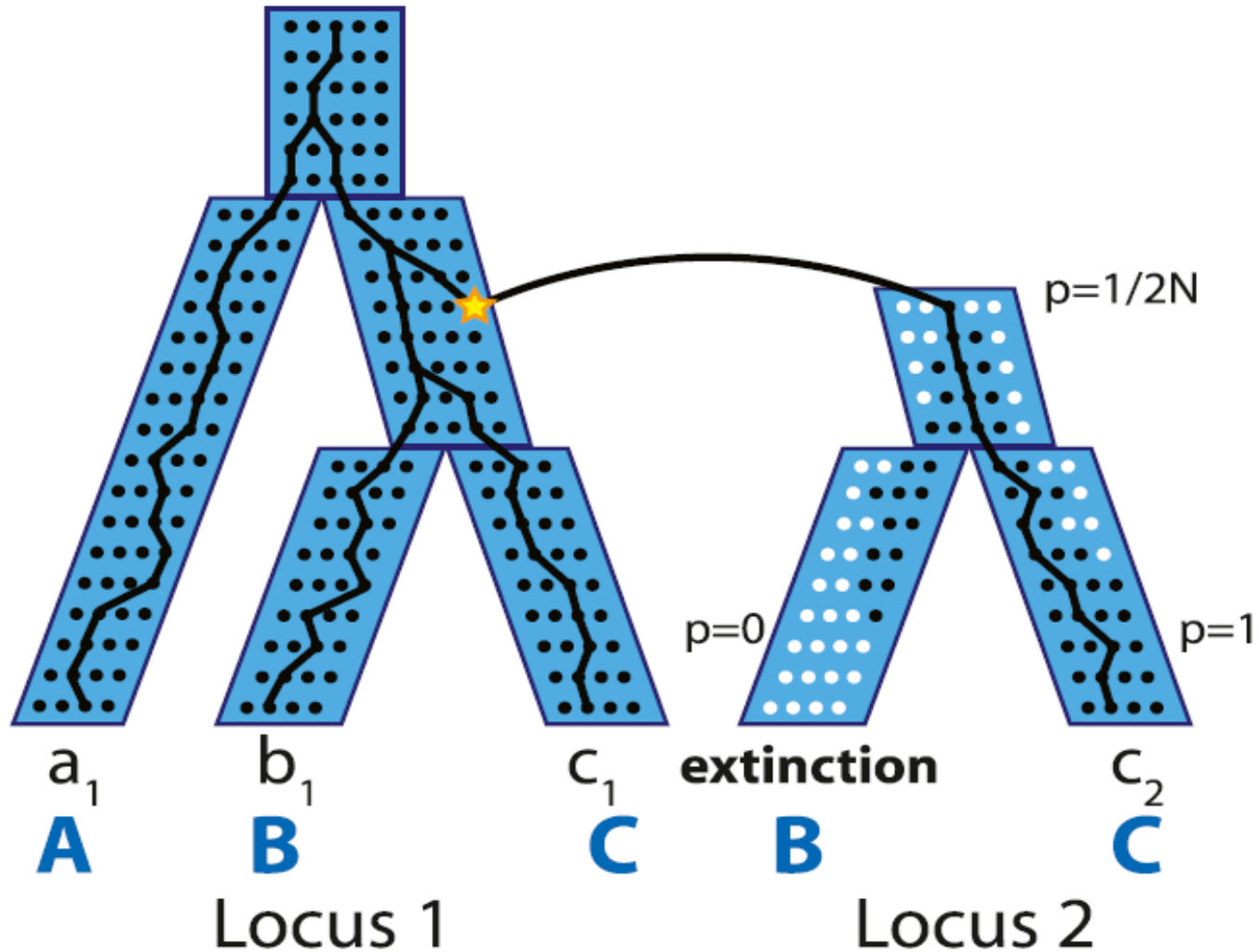
# Example of duplication in a population



# Example of duplication in a population



# Duplications can fail to fix

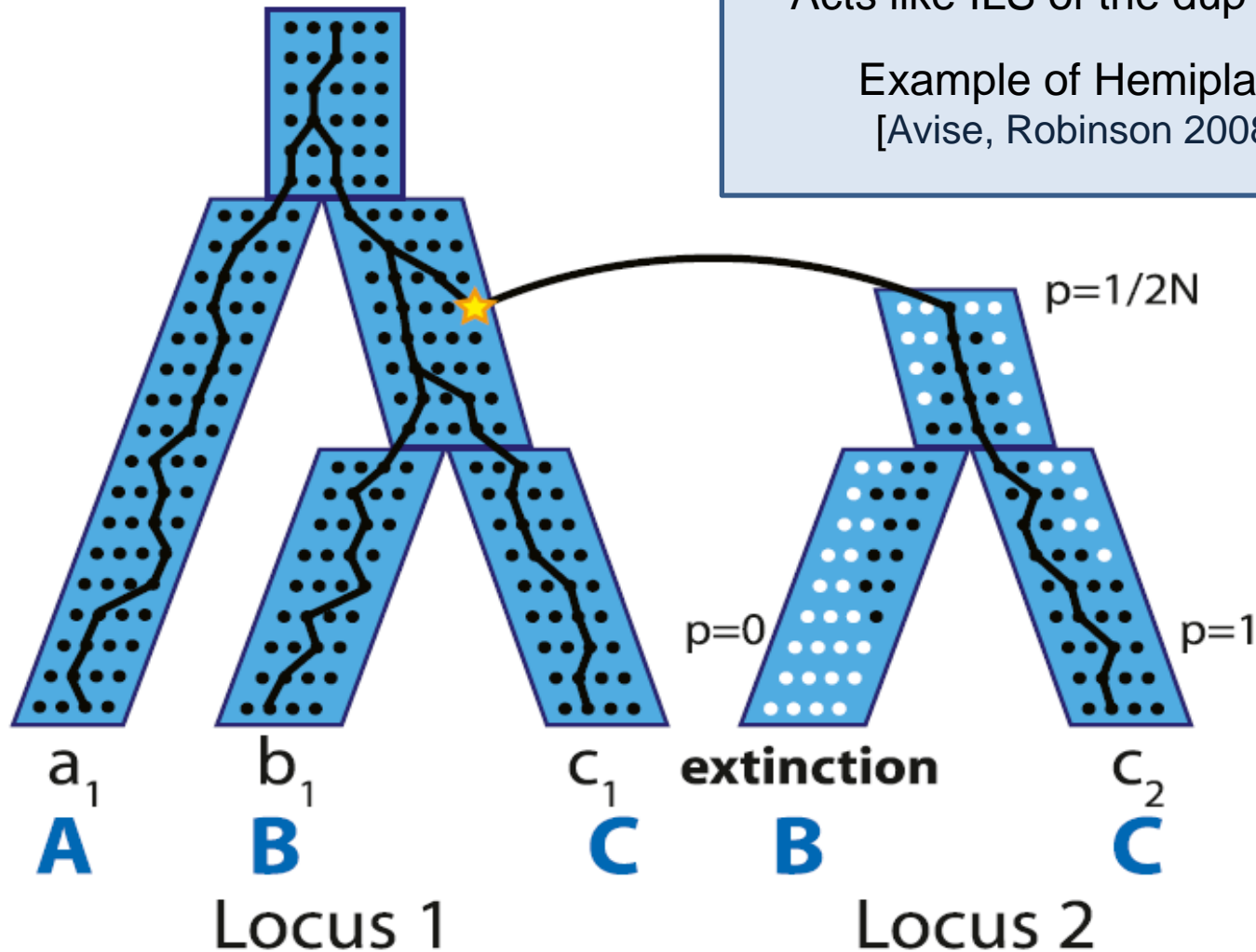


# Duplications can fail to fix

Acts like ILS of the dup event

Example of Hemiplasy

[Avice, Robinson 2008]

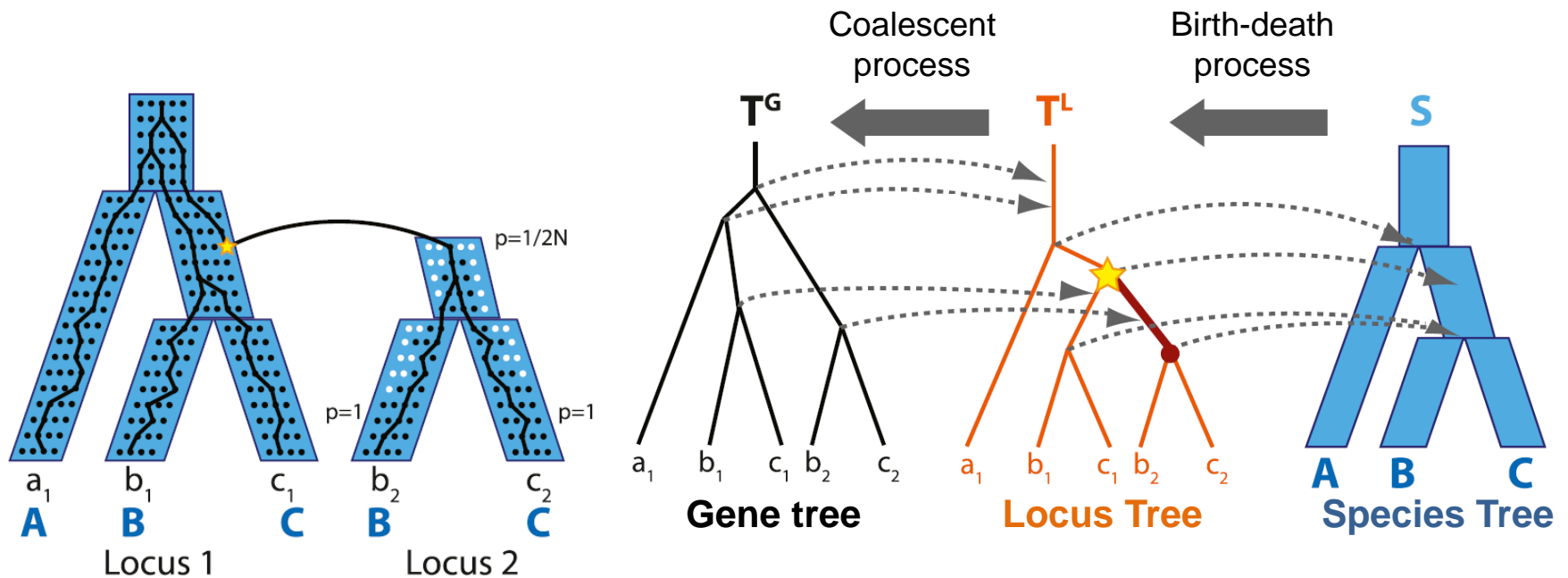


# How to unify dup-loss and coalescent models?

- **Realization:** the “gene trees” in dup-loss and coalescent models are distinct objects
  - In coalescent: describes the history of gene sequences
  - In dup-loss: describes history of loci (i.e. changes in copy number)
- **Resolution: “three tree model”**

# How to unify dup-loss and coalescent models?

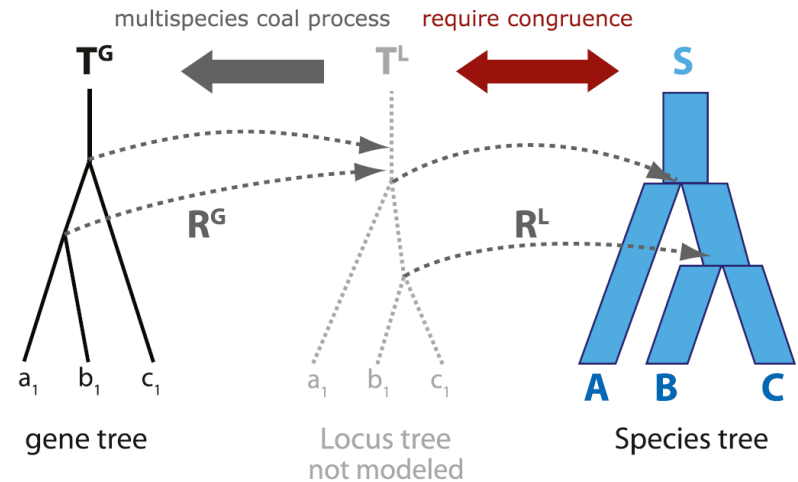
- **Realization:** the “gene trees” in dup-loss and coalescent models are distinct objects
  - In **coalescent**: describes the history of gene sequences
  - In **dup-loss**: describes history of loci (i.e. changes in copy number)
- **Resolution: “three tree model”**
  - Can track lineages across: individuals, loci, and species



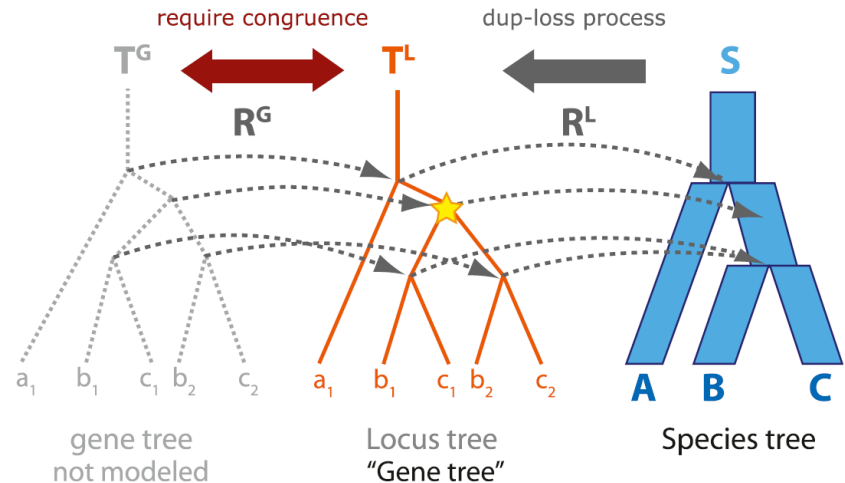


# DLCoal model generalizes dup-loss and multispecies coalescent models

- **Multispecies coalescent**  
assumes locus tree congruent to species tree (i.e. no paralogs)

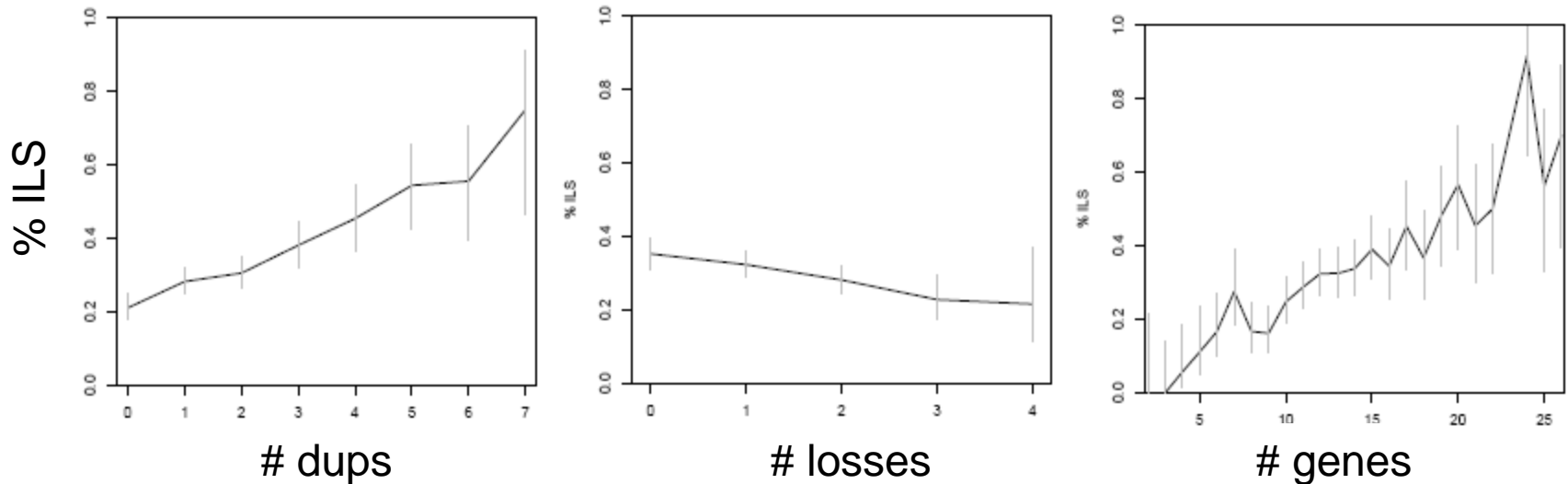


- Previous **dup-loss models**  
assume gene tree congruent to locus tree (i.e. no ILS)



# Simulation with dup,loss,coal:

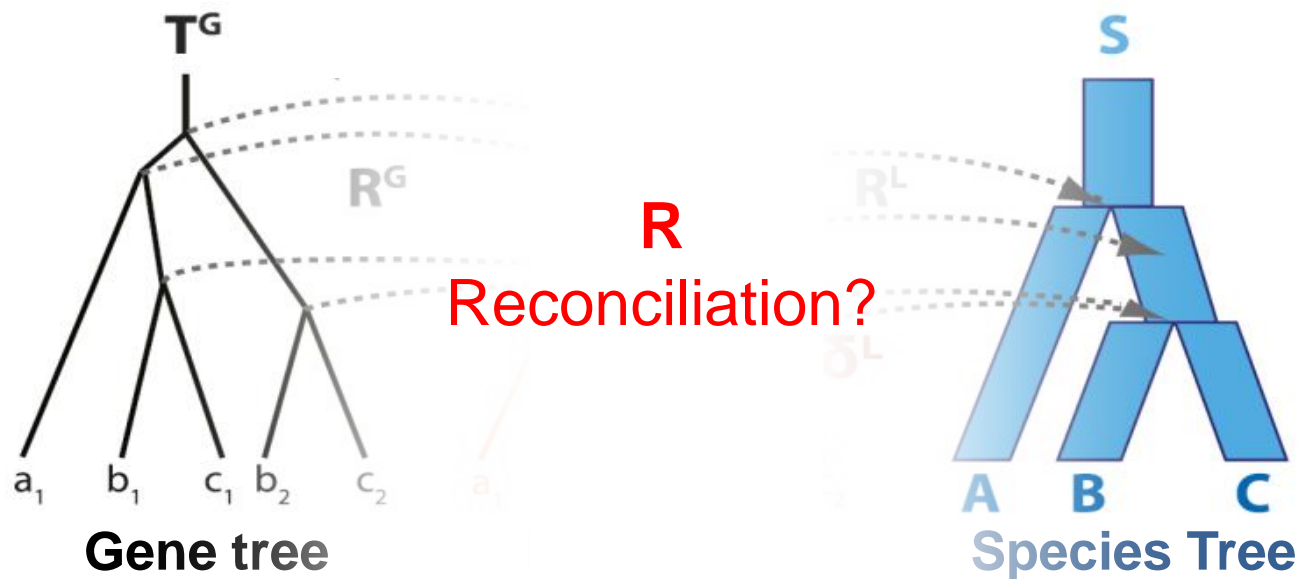
Large families (more dups) have higher ILS rate



- Duplications break up long branches in locus tree → ILS more likely
- Losses do the reverse, joining branches in the locus tree → ILS less likely

# A new reconciliation method: DLCoalRecon

**Problem:** Given gene tree and species tree, find most likely duplications and losses in the presence of ILS.

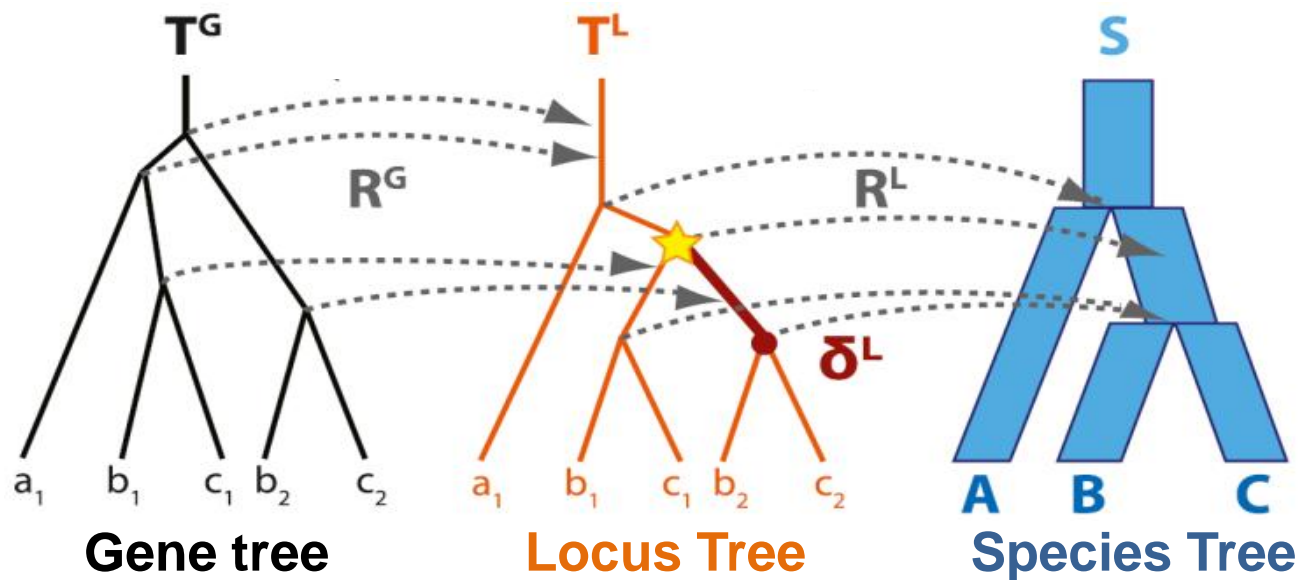


**Input:** Gene tree  $T^G$ , Species tree  $S$ , model parameters  $\theta$  (pop size, dup/loss rates)

**Output:**  $R = ?$

# A new reconciliation method: DLCoalRecon

**Problem:** Given gene tree and species tree, find most likely duplications and losses in the presence of ILS.



**Input:** Gene tree  $T^G$ , Species tree  $S$ , model parameters  $\theta$  (pop size, dup/loss rates)

**Output:**  $R = (\text{Locus tree } T^L, \text{ daughters } \delta^L, \text{ mappings } R^G, \text{ and } R^L)$

# DLCoalRecon: Maximum A Posterior reconciliation

**Input:** Gene tree  $T^G$ , Species tree  $S$ , model parameters  $\theta$  (pop size, dup/loss rates)

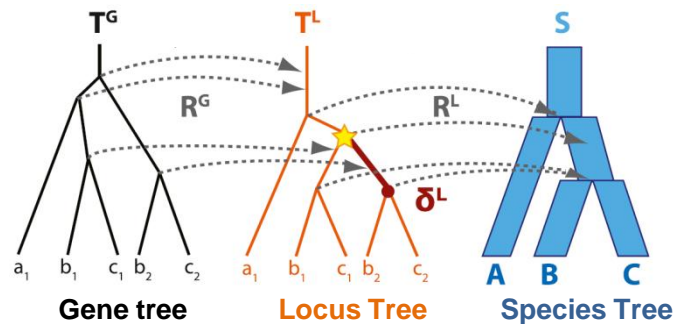
**Output:**  $R = (\text{Locus tree } T^L, \text{daughters } \delta^L, \text{mappings } R^G, \text{ and } R^L)$

$$\text{Goal: } \hat{\mathbb{R}} = \underset{\mathbb{R}}{\operatorname{argmax}} P(\mathbb{R} | T^G, S, \theta) = \underset{\mathbb{R}}{\operatorname{argmax}} P(\mathbb{R}, T^G | S, \theta)$$

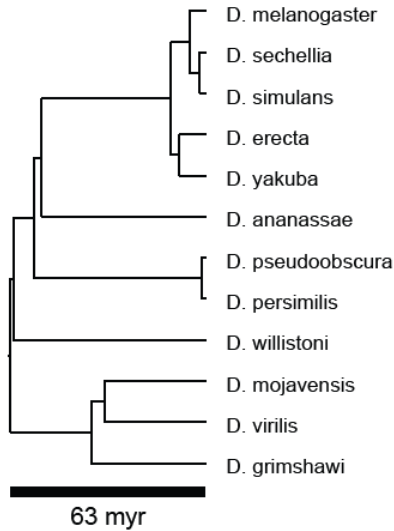
Factor probability into previously derived terms:

$$P(T^G, \mathbb{R} | S, \theta) = P(\delta^L | T^L, R^L, S) P(T^L, R^L | S, \theta) \int P(T^G, R^G | T^L, \mathbf{t}^L, \delta^L, \mathbf{N}^L) P(\mathbf{t}^L | T^L, R^L, S, \theta) d\mathbf{t}^L$$

↑ Daughters      ↑ Topology prior Dup/loss model      ↑ Topology prior Coal model      ↑ Branch length prior Dup/loss model



# DLCoalRecon outperforms on 500 simulated fly gene families

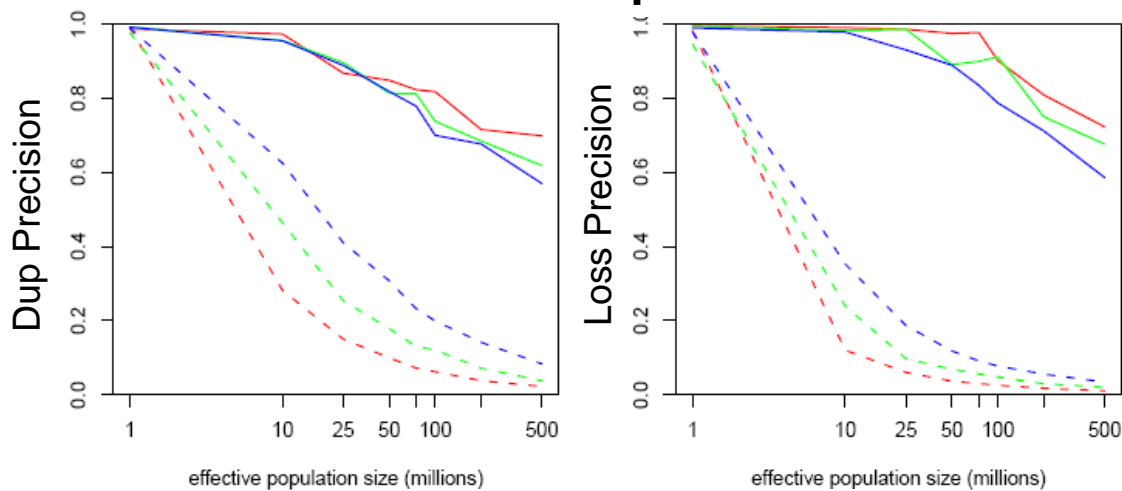


	Actual	MPR	DLCoal Recon
# dup	218	1562	231
dup sensitivity		71.6%	89.9%
dup precision		10.0%	84.8%
# loss	202	4873	200
loss sensitivity		88.6%	96.5%
loss precision		3.7%	97.5%
# orths	32,625	17,680	32,708
orth sensitivity		54.2%	99.6%
orth precision		100%	99.4%

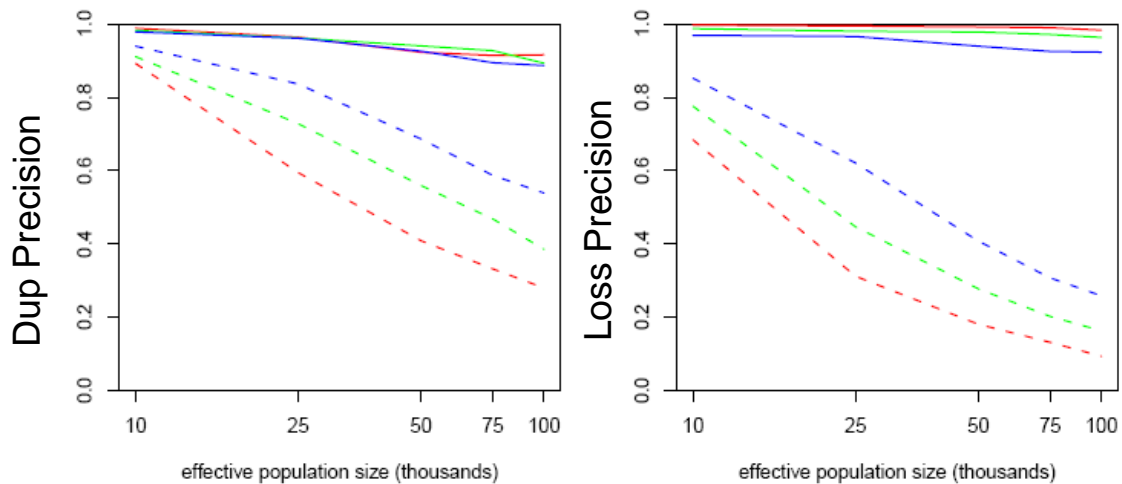
- 0.1 years/generation
- $N_e=50 \times 10^6$
- Dup-loss = 0.0012 event/gene/myr
  - Hahn et al 2007

# Performance holds over a variety of populations sizes and dup/loss rates

## Drosophila

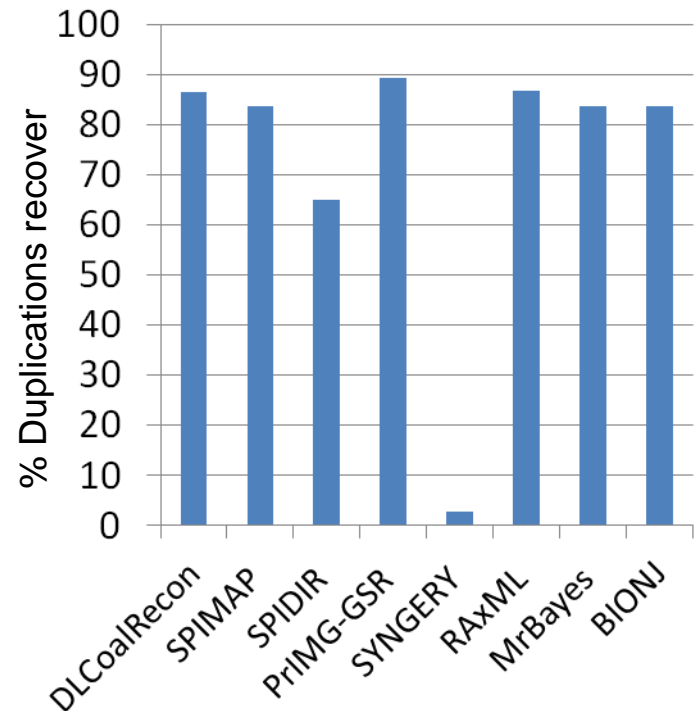


## Primates



# DLCoalRecon outperforms on 16 fungi genomes

Phylo Program	Recon program	% Orth	No. Orth	No. Dup	No. Loss
PhyML	DLCoalRecon	97.8%	575,374	4533	6398
PhyML	MPR	64.2%	464,479	21,264	64,391
RAXML	MPR	63.8%	463,020	21,485	65,392
MrBayes	MPR	63.9%	460,510	21,307	65,238
BIONJ	MPR	60.4%	439,193	22,396	71,231
SPIMAP	—	96.5%	557,981	5407	10,384
SYNERGY	—	99.2%	595,289	4604	8179
PrIME-GSR	—	88.9%	527,153	7951	21,099



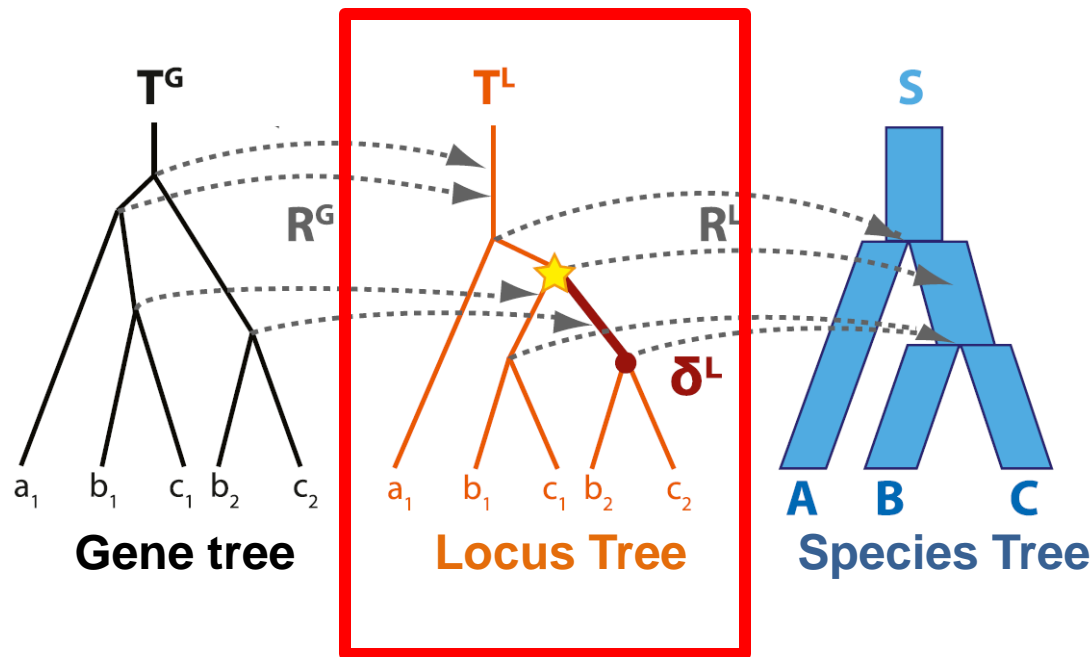
(1) Further improved syntenic ortholog recovery

(2) Even fewer dup-losses inferred

(3) Improved gene conversion recovery



# Importance of the locus tree



- Dup & losses can only meaningfully be annotated on the locus tree.
- Locus tree is a gene tree with the ILS removed and is often more meaningful
- DLCoal has been critical for accurate orthology determination in the modENCODE project

# Future directions

- **Continue to model more sources of incongruence**
  - Duplication-loss-horizontal transfer
  - Incomplete Lineage sorting (ILS)
  - Uncertainty of trees and alignments
  - Domain rearrangement
  - Gene conversion
  - Recombination
- **Apply models to new applications**
  - Species tree inference
  - Demography, migration rates
  - Selection

# Acknowledgements



**Adam Siepel**  
Cornel  
Postdoc advisor



**Manolis Kellis**  
MIT  
Postdoc & PhD  
advisor

## Candida sequencing and analysis

Geraldine Butler (UC Dublin)

Mike Lorenz (Texas)

Neil Gow (Aberdeen)

Christina Cuomo (Broad Institute)

Bruce Birren (Broad Institute)

Mike Lin (MIT)

## STAR-MP: domain evolution

Yi-chieh Wu

Manolis Kellis

## TreeFix: gene tree correction

Yi-chieh Wu

Mukul Bansal

## Horizontal transfer in Strep

Sang Chul Choi

Melissa Hubisz

Ilan Gronau

Michael Stanhope

## Funding sources

Ruth L. Kirschstein National Research

Service Award

Cornell Center for Comparative and

Population Genomics Fellowship

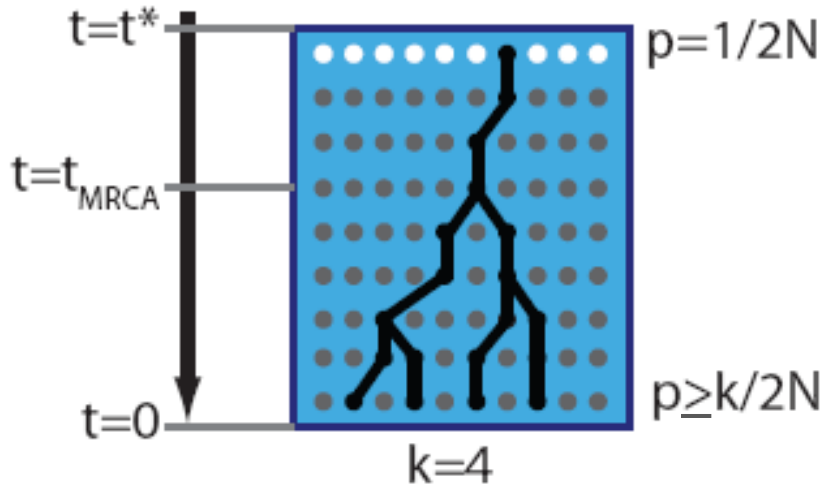


# First version of a dup,loss,coal model

For a first version, we make the following assumptions

- New duplicates begin at unlinked loci
- We model no gene conversion
- **Hemiplasy assumption:**
  - no event (duplications or losses) under goes hemiplasy
  - Namely, full extinction or no extinctions

# Building up the model: Bounded coalescent



- **Given:**

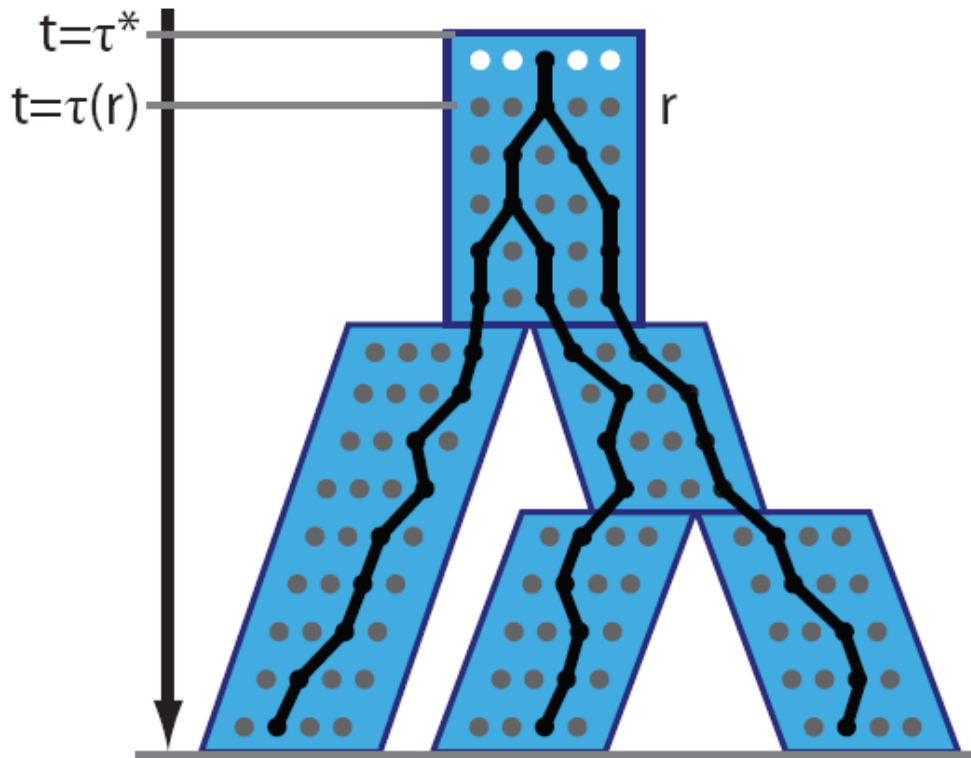
- new mutation (black) at  $t^*$
- $k$  lineages at  $t=0$ , all of them have mutation

- Equivalent to conditioning  $t_{MRCA} < t^*$

$$P(t | t_{MRCA} < t^*, k, N) = \frac{P(t, t_{MRCA} < t^* | k, N)}{P(t_{MRCA} < t^* | k, N)}$$

$$= \begin{cases} \frac{P(t | k, N)}{P(t_{MRCA} < t^* | k, N)}, & \text{if } t_{MRCA} < t^* \\ 0, & \text{otherwise} \end{cases}$$

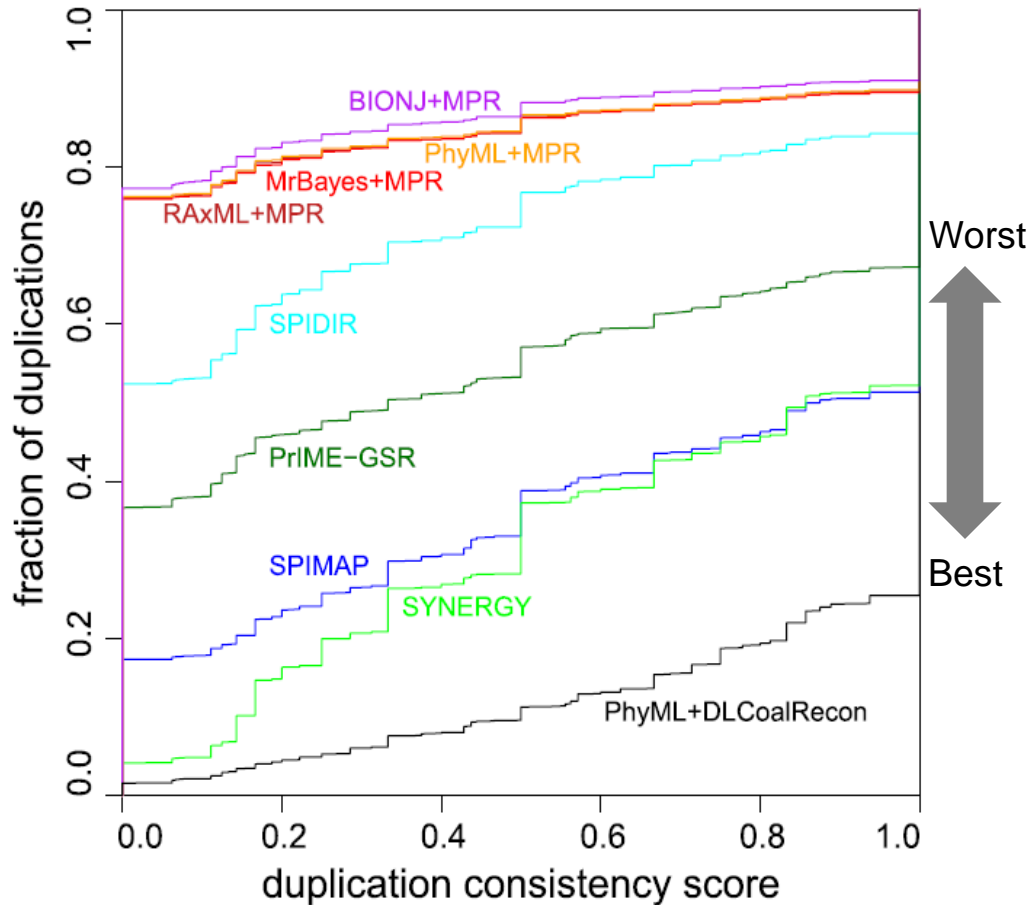
# Building up the model: Bounded multispecies coalescent



- Condition  $t(r) < t^*$
- Use time of MRCA of MC [Efromovich & Kubatko 2009]

# DLCoalRecon outperforms on 16 fungi genomes

(3) Improved duplication consistency



(4) Improved gene conversion recovery

