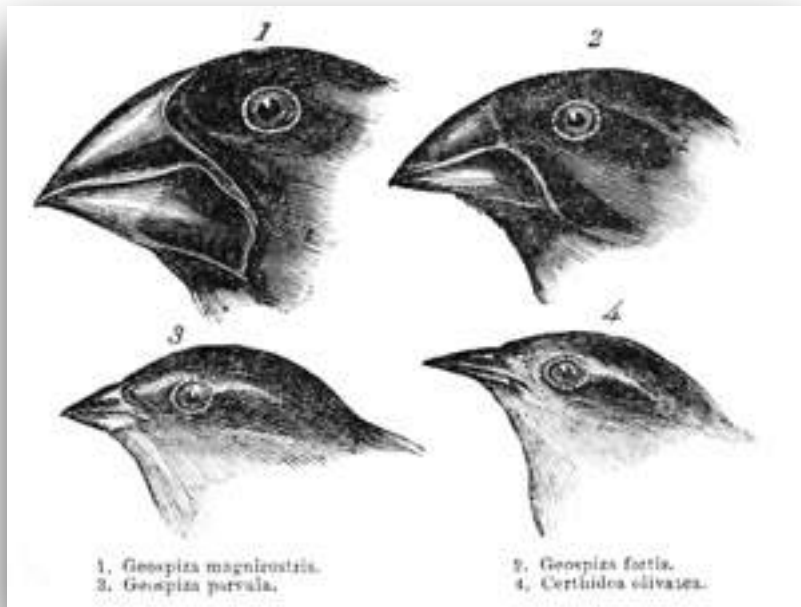# Mathematics of the Tree of Life
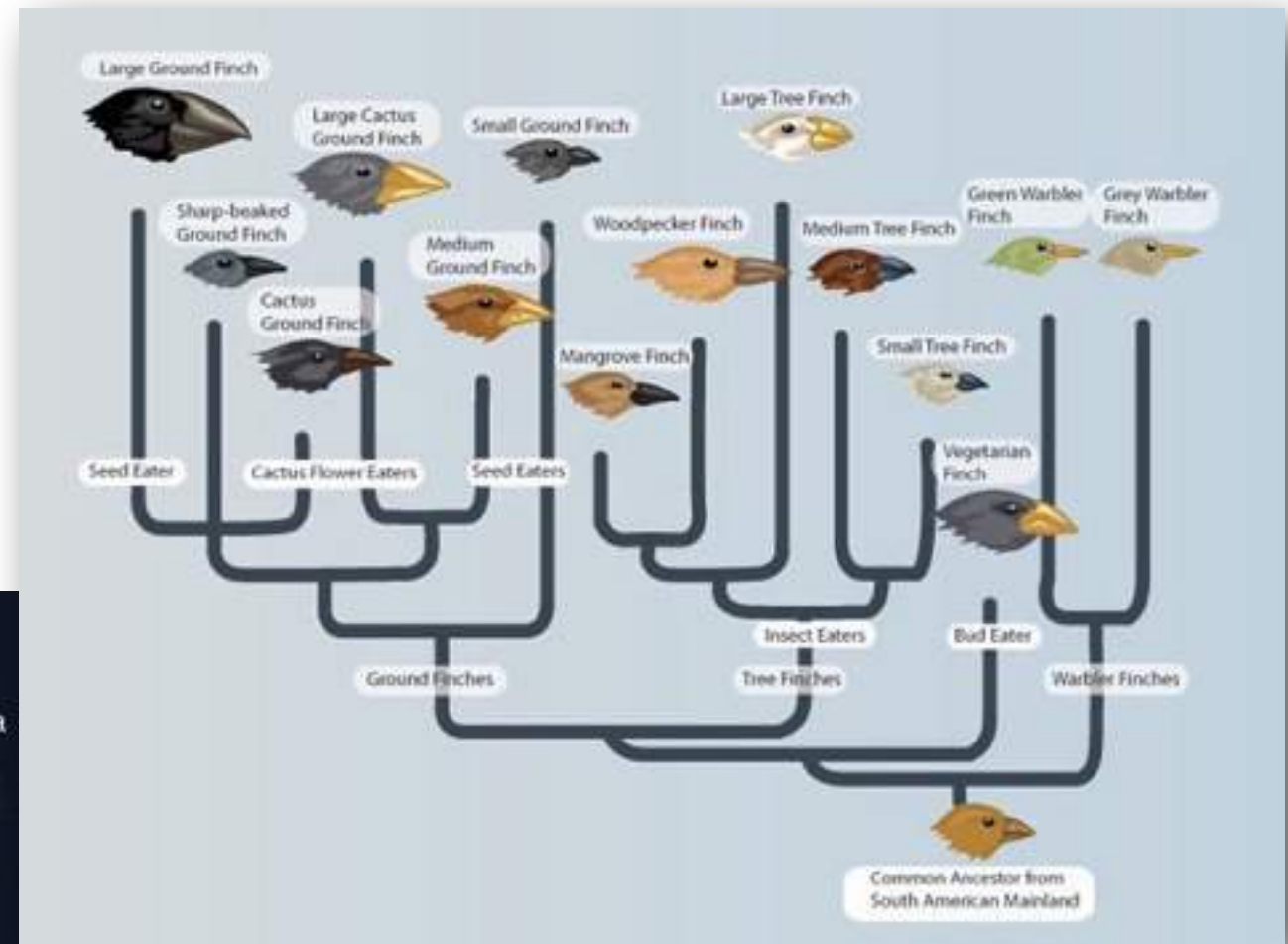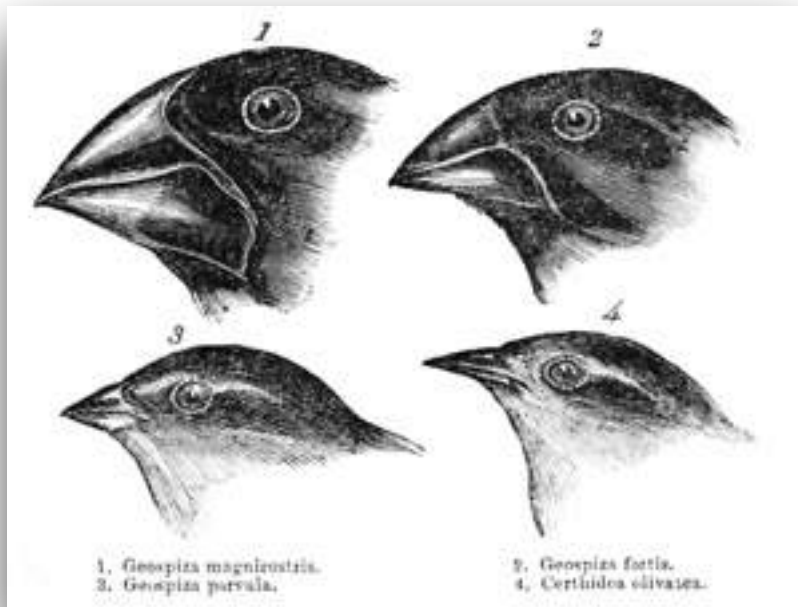## From Genomes to Phylogenetic Trees and Beyond

Sébastien Roch

Department of Mathematics

UW-Madison

# Darwin's finches
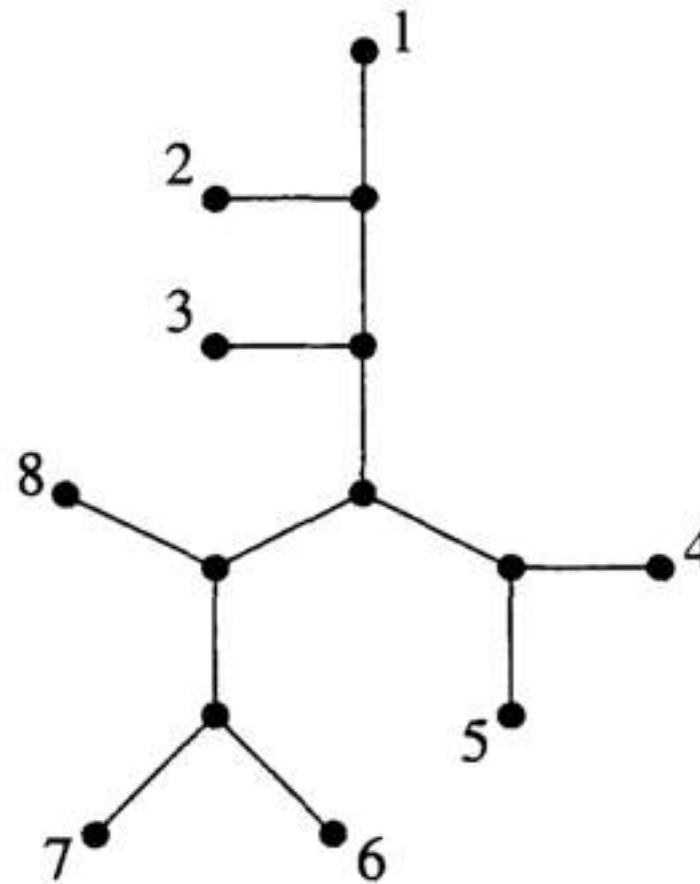
# Darwin's finches

# Phylogenetic *X*-trees

An *X-tree* is a pair $(T; \phi)$ where $T$ is a tree and $\phi : X \to V(T)$ is a labeling such that $\deg(v) \le 2 \implies v \in \phi(X)$. It is a *phylogenetic X-tree* if $\phi$ is a bijection into the leaves.

# Phylogenetic *X*-trees

## Definition

An *X-tree* is a pair $(T; \phi)$ where $T$ is a tree and $\phi : X \to V(T)$ is a labeling such that $\deg(v) \leq 2 \implies v \in \phi(X)$. It is a *phylogenetic X-tree* if $\phi$ is a bijection into the leaves.

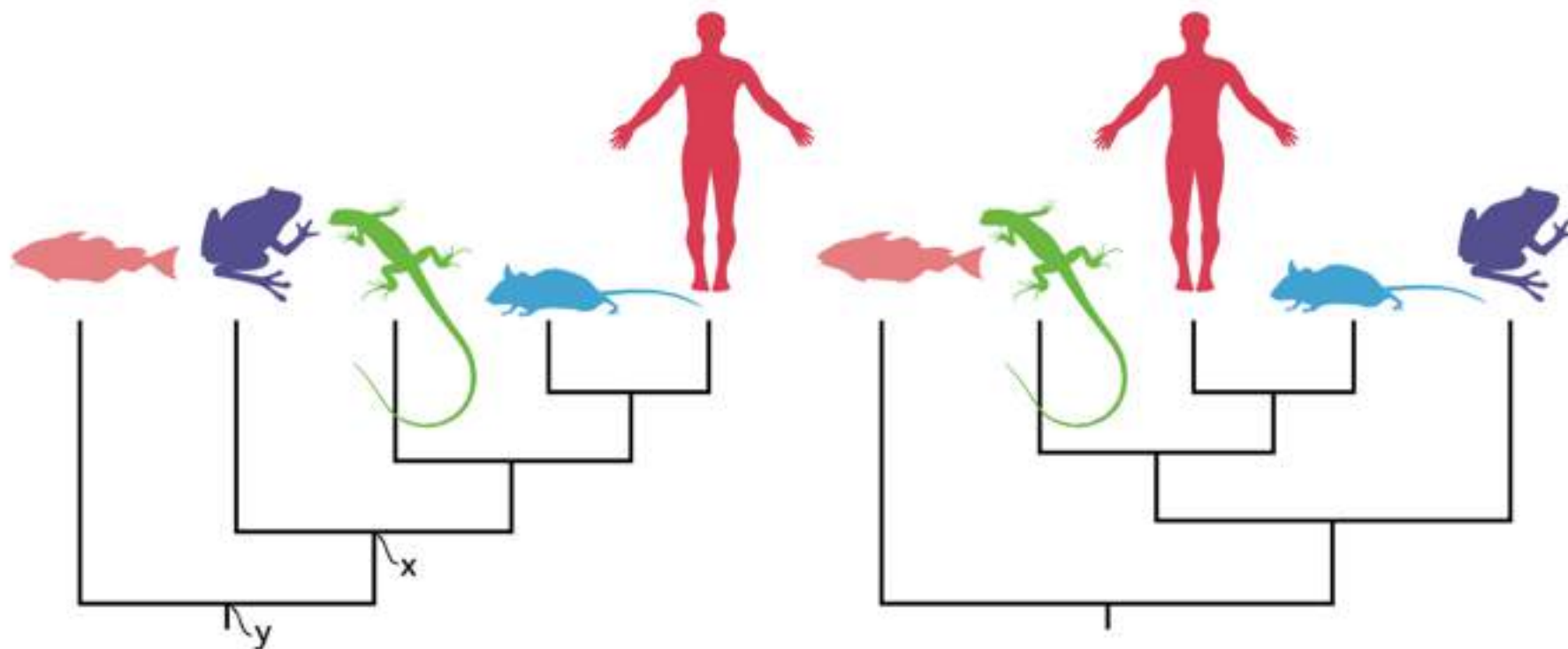## Definition

Two *X*-trees $(T_1; \phi_1)$ and $(T_2; \phi_2)$ are *isomorphic* if there is a graph isomorphism $\psi$ between $T_1$ and $T_2$ such that $\phi_2 = \psi \circ \phi_1$.
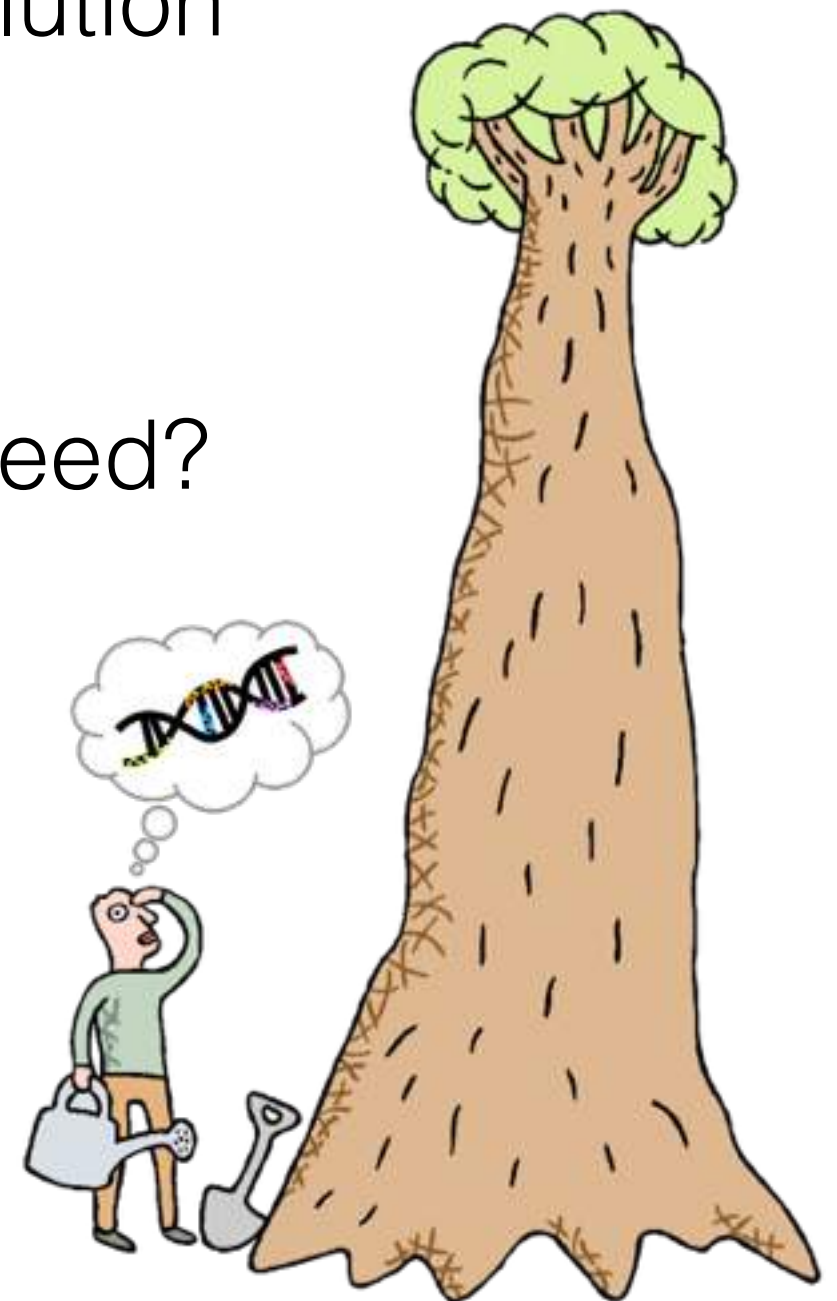
# So how is the Tree of Life inferred?

I. From Darwin's finches to HIV evolution

**II. Pre-genomics era**

III. Transition: How much data do I need?

**IV. More data, more problems**
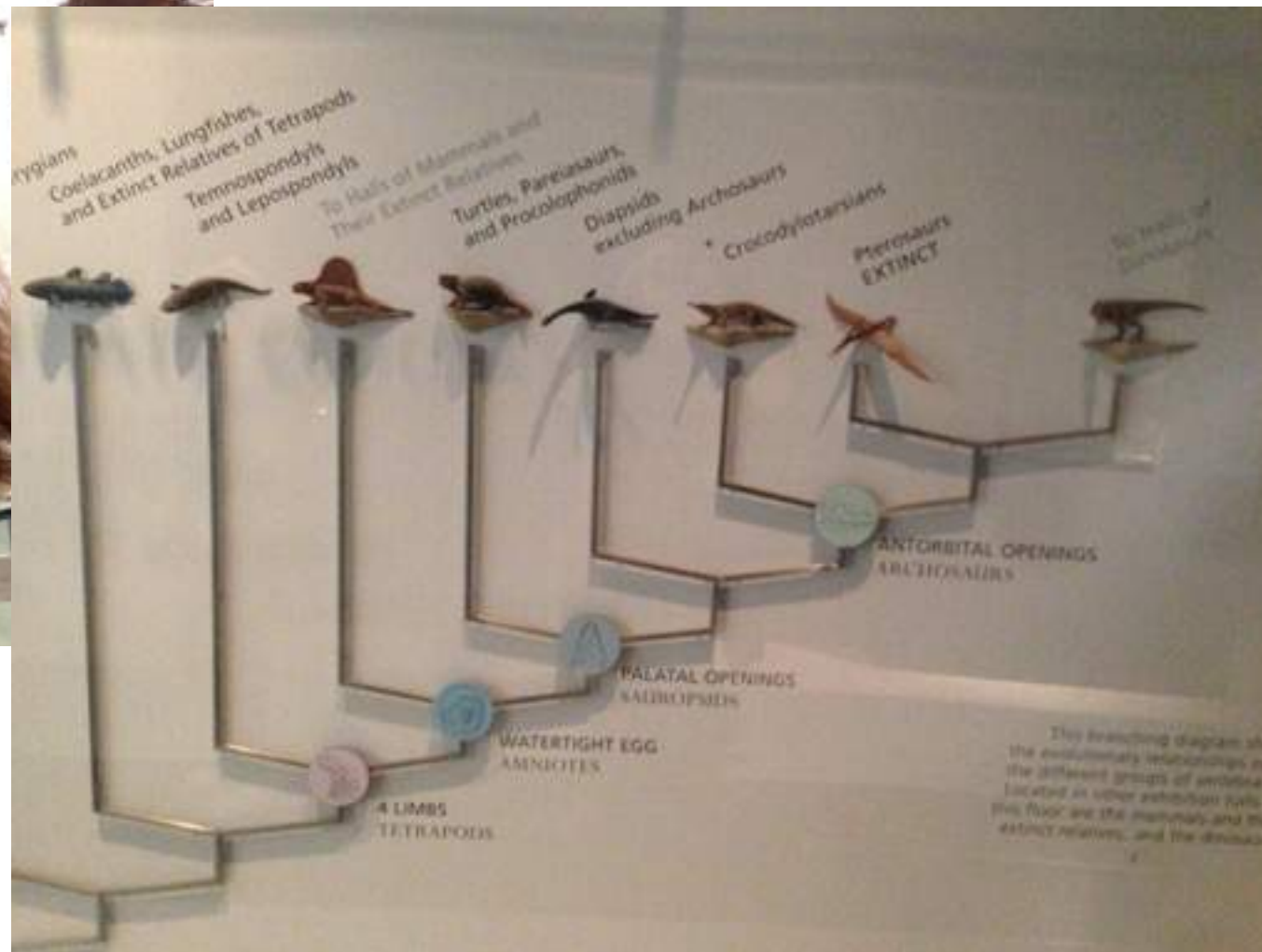
V. Is the Tree of Life even a tree?

Pre-genomics era

# Strolling down the Tree of Life

# Strolling down the Tree of Life

# Strolling down the Tree of Life
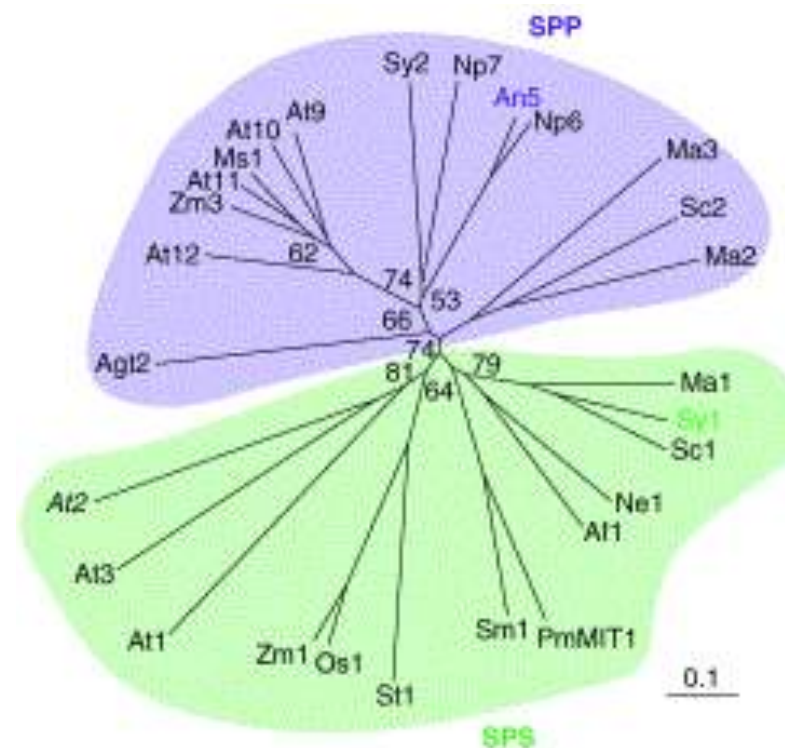
# Strolling down the Tree of Life

# Compatible splits

**Definition**

An *X-split* $A|B$ is a bipartition of $X$ into non-empty subsets $A$, $B$.

**Definition**

A pair of $X$-splits $A_1|B_1$ and $A_2|B_2$ is *compatible* if at least one of the sets $A_1 \cap A_2$, $A_1 \cap B_2$, $B_1 \cap A_2$, or $B_1 \cap B_2$ is the empty set.

**Theorem (Splits-equivalence theorem; Buneman (1971))**

A set of $X$-splits is induced by an $X$-tree iff it is compatible.
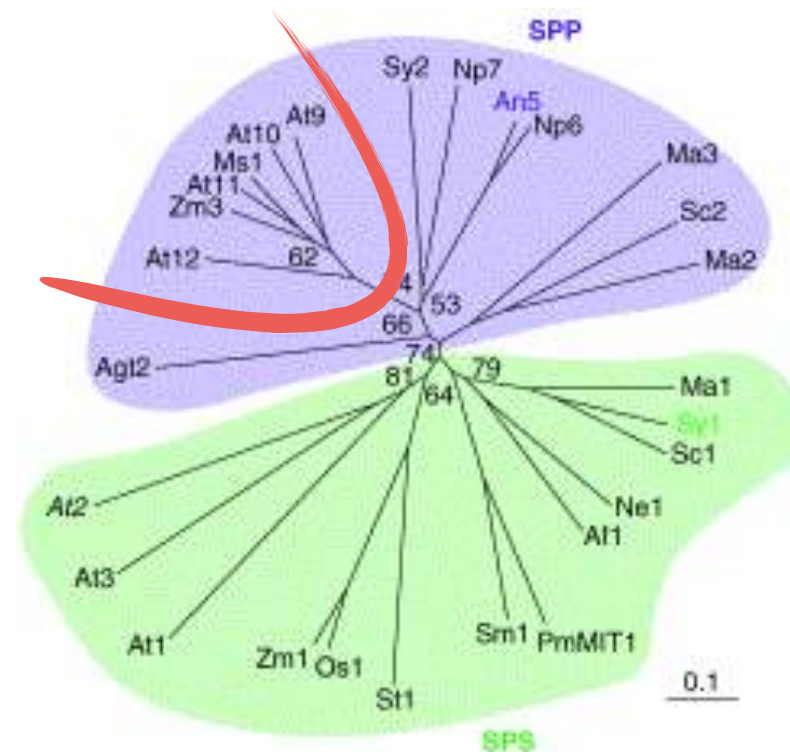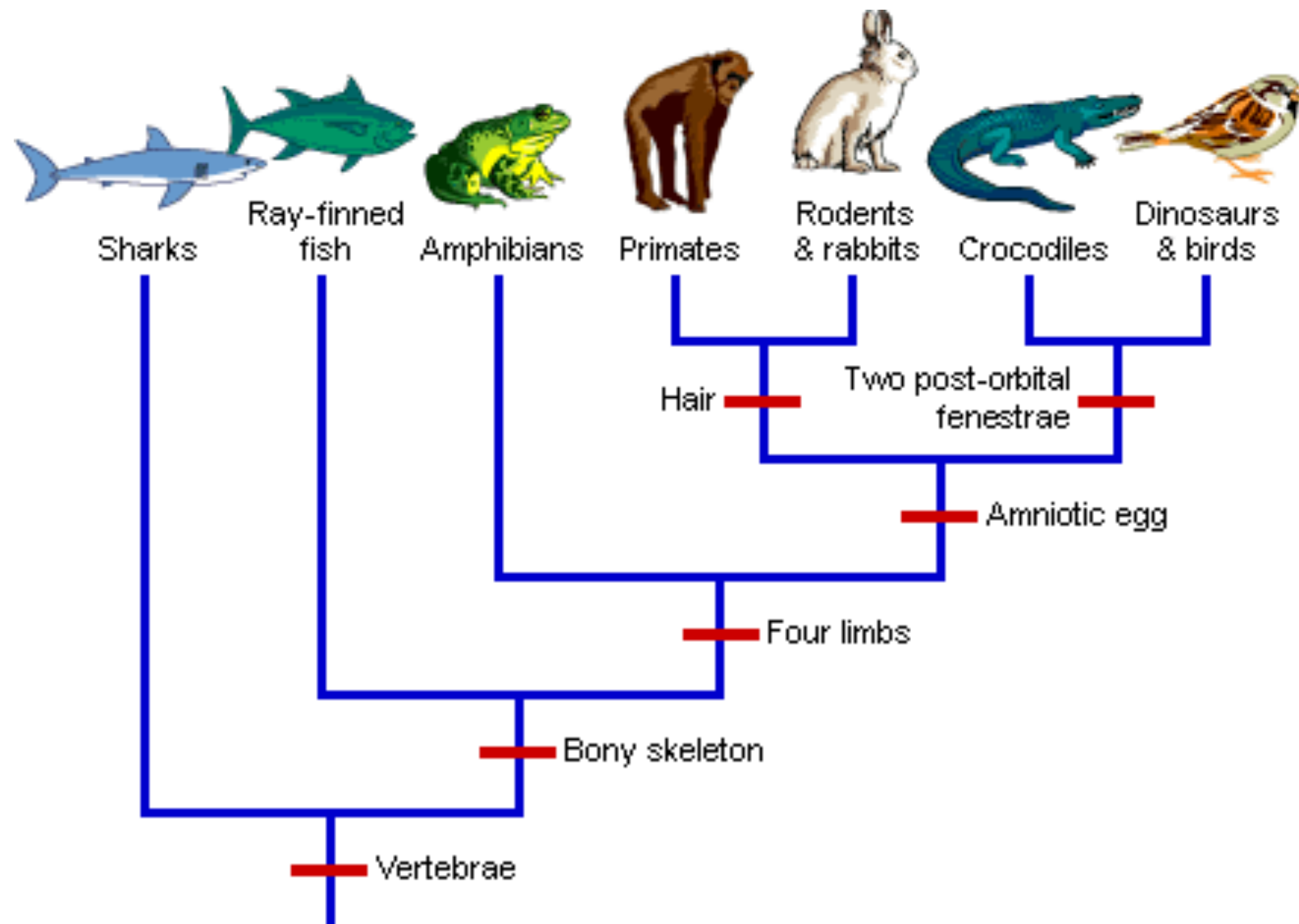
# Compatible splits

## Definition
An *X-split* $A|B$ is a bipartition of $X$ into non-empty subsets $A$, $B$.

## Definition
A pair of *X*-splits $A_1|B_1$ and $A_2|B_2$ is *compatible* if at least one of the sets $A_1 \cap A_2$, $A_1 \cap B_2$, $B_1 \cap A_2$, or $B_1 \cap B_2$ is the empty set.

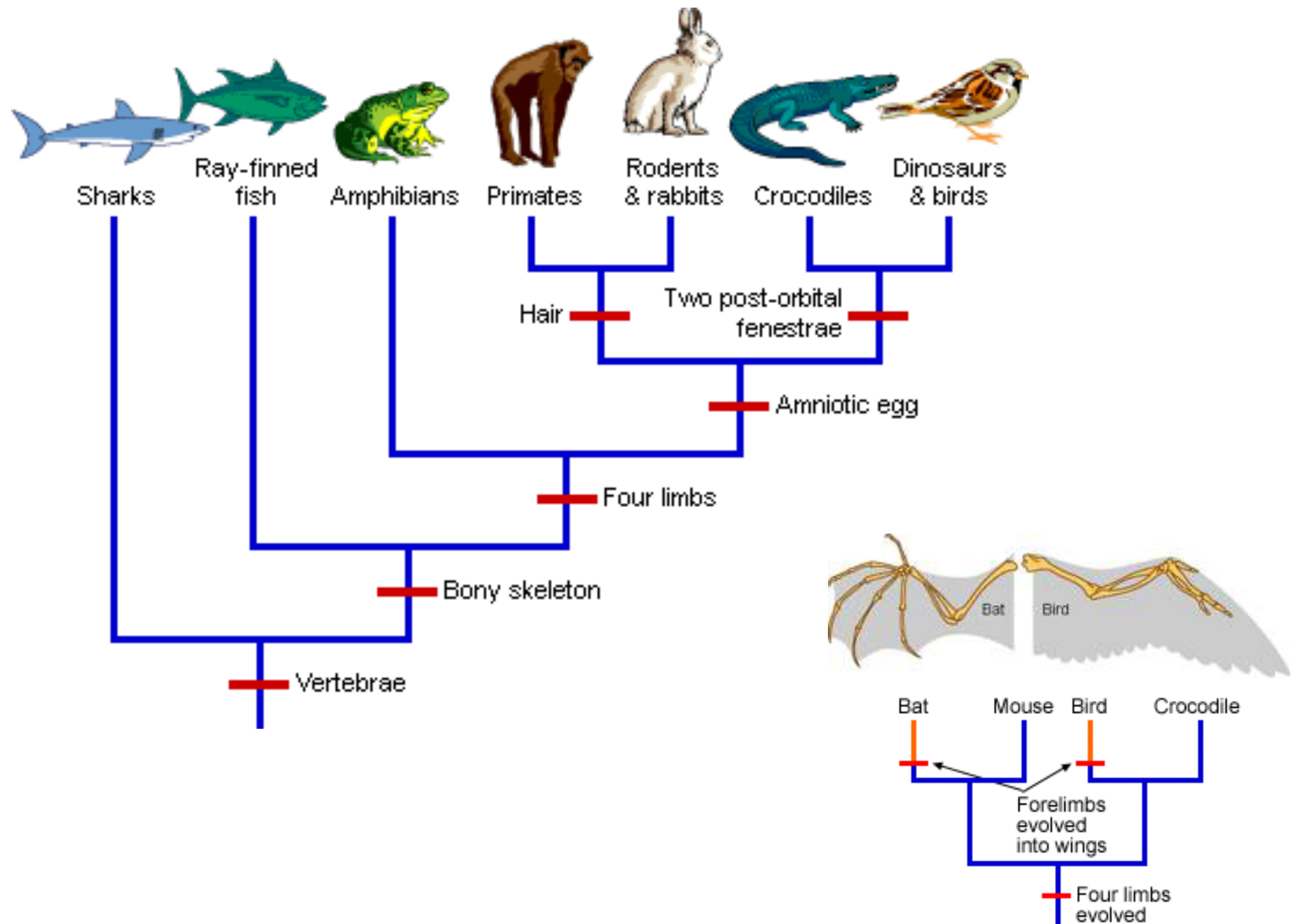## Theorem (Splits-equivalence theorem; Buneman (1971))
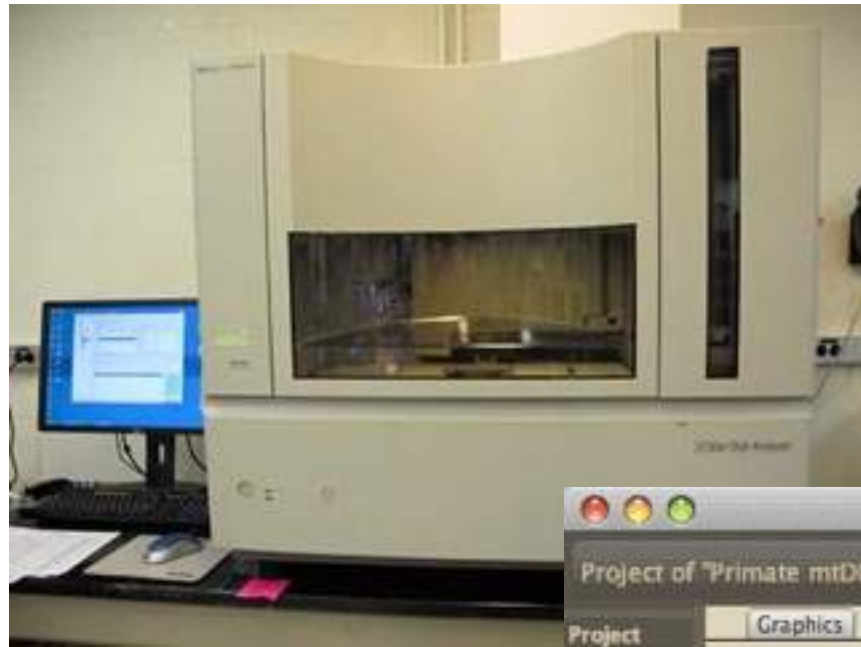A set of *X*-splits is induced by an *X*-tree iff it is compatible.

# Synapomorphies & homoplasies



Sharks  Ray-finned fish  Amphibians  Primates  Rodents & rabbits  Crocodiles  Dinosaurs & birds

Hair

Two post-orbital fenestrae

Amniotic egg

Four limbs

Bony skeleton

Vertebrae

# Synapomorphies & homoplasies



Figures by University of California Museum of Paleontology's Understanding Evolution

# Molecular systematics



*Snapshot from Mesquite*

# Molecular systematics



*Snapshot from Mesquite*

# Molecular systematics

# Tree metrics

**Definition**

A function $\delta : X \times X \to \mathbb{R}$ is a *tree metric* if there is an *X*-tree $\mathcal{T} = (T; \phi)$ and a weighting $w : E(T) \to \mathbb{R}_+$ such that for all $x, y$

$$\delta(x, y) = d_{(\mathcal{T};w)}(x, y) := \sum_{e \in P(\mathcal{T};x,y)} w(e),$$

where $P(\mathcal{T}; x, y)$ is the unique path between $\phi(x)$ and $\phi(y)$. The *tree metric representation* $(\mathcal{T}; w)$ of $\delta$ is unique (and efficiently computable).



|   | a | b | c | d | e | f |
|---|---|---|---|---|---|---|
| a | 0 | 3 | 8 | 1 | 10 | 8 |
| b |   | 0 | 9 | 2 | 11 | 9 |
| c |   |   | 0 | 7 | 8 | 6 |
| d |   |   |   | 0 | 9 | 7 |
| e |   |   |   |   | 0 | 4 |
| f |   |   |   |   |   | 0 |

# Tree metrics

**Definition**

A function $\delta : X \times X \to \mathbb{R}$ is a *tree metric* if there is an *X*-tree $\mathcal{T} = (T; \phi)$ and a weighting $w : E(T) \to \mathbb{R}_+$ such that for all $x, y$

$$\delta(x, y) = d_{(\mathcal{T};w)}(x, y) := \sum_{e \in P(\mathcal{T};x,y)} w(e),$$

where $P(\mathcal{T}; x, y)$ is the unique path between $\phi(x)$ and $\phi(y)$. The *tree metric representation* $(\mathcal{T}; w)$ of $\delta$ is unique (and efficiently computable).

# Tree metrics

**Definition**

A function $\delta : X \times X \to \mathbb{R}$ is a *tree metric* if there is an $X$-tree $\mathcal{T} = (T; \phi)$ and a weighting $w : E(T) \to \mathbb{R}_+$ such that for all $x, y$

$$\delta(x, y) = d_{(\mathcal{T};w)}(x, y) := \sum_{e \in P(\mathcal{T};x,y)} w(e),$$

where $P(\mathcal{T}; x, y)$ is the unique path between $\phi(x)$ and $\phi(y)$. The *tree metric representation* $(\mathcal{T}; w)$ of $\delta$ is unique (and efficiently computable).
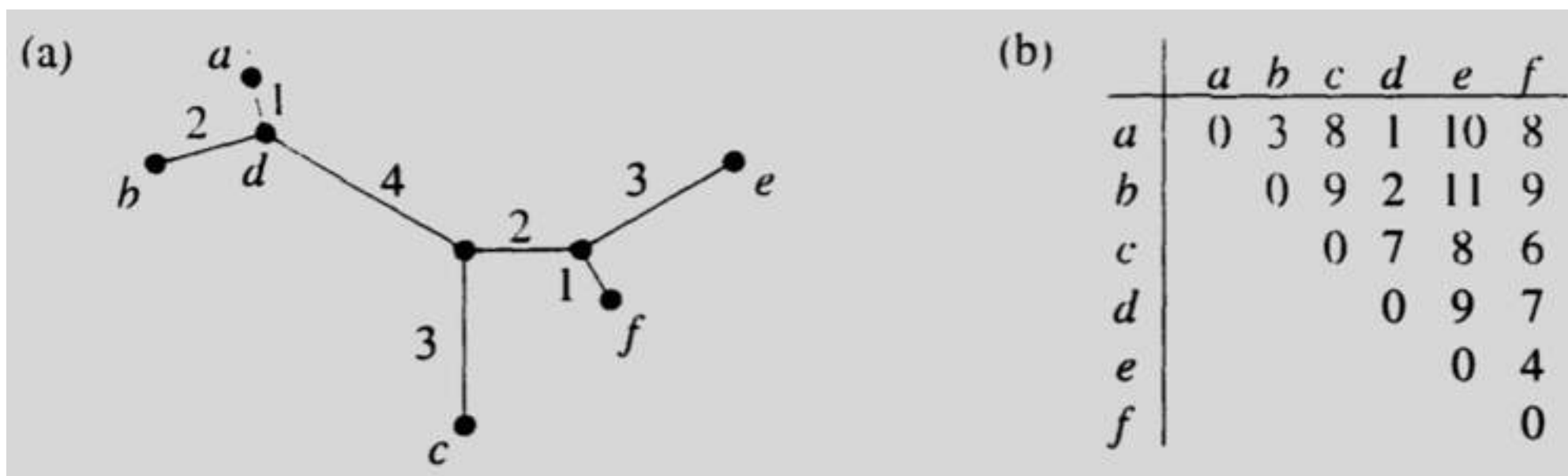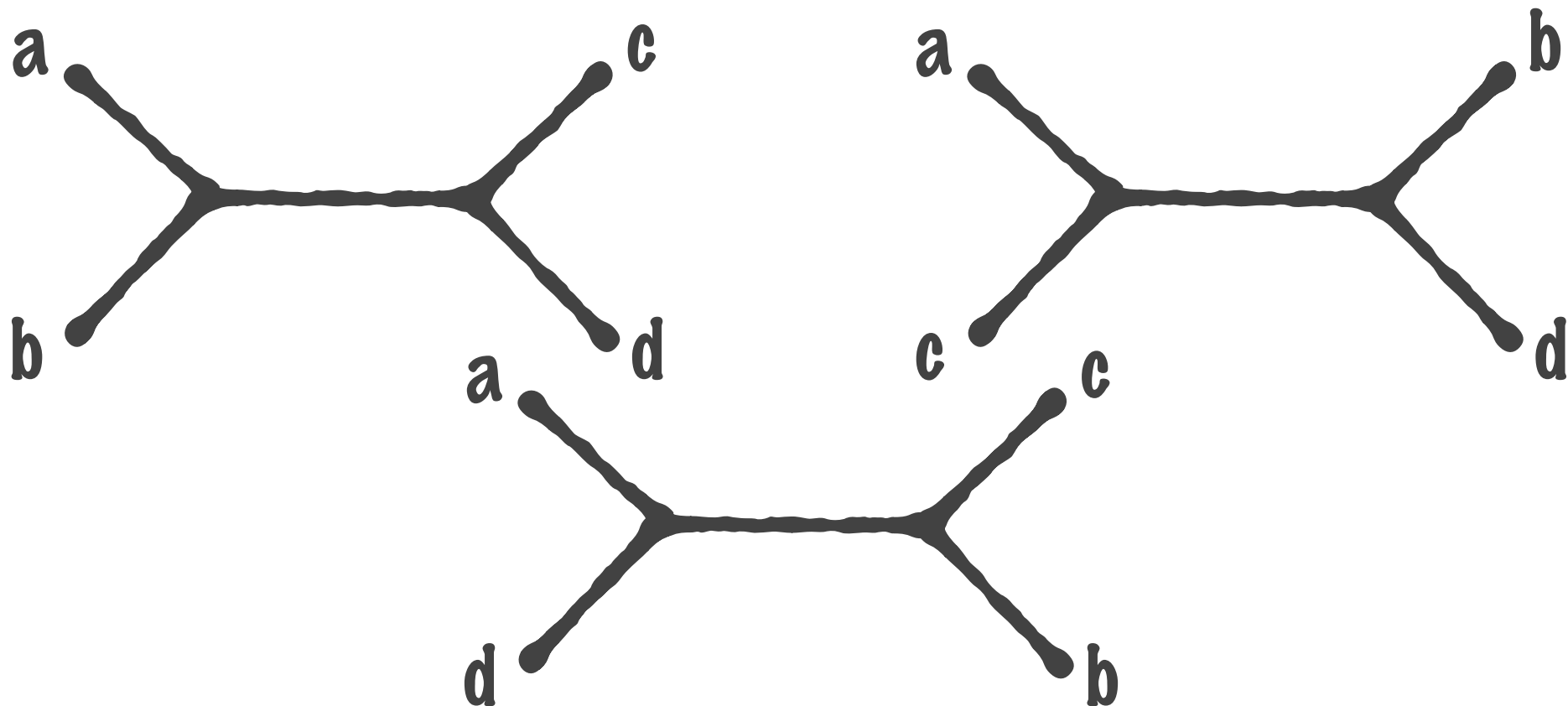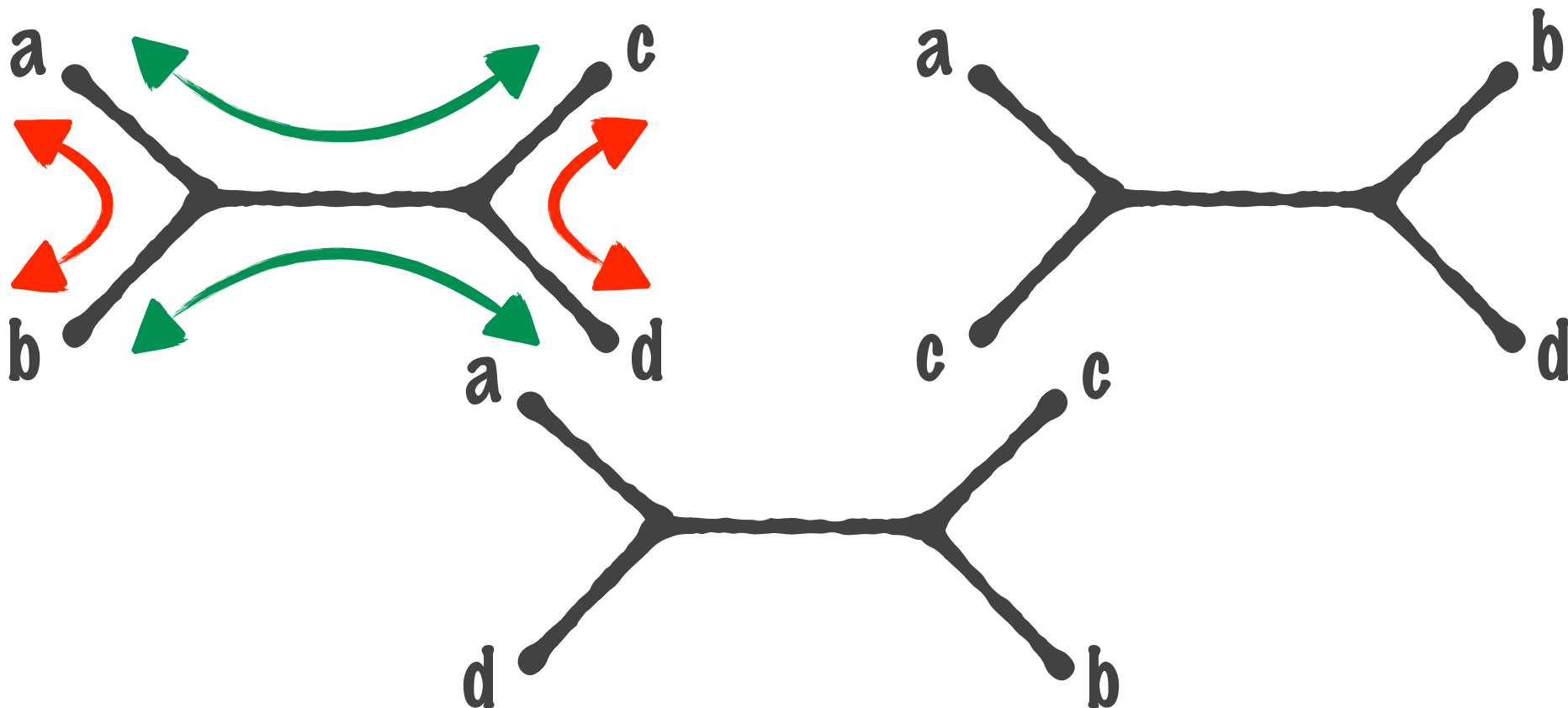
# Tree metrics

**Definition**

A function $\delta : X \times X \to \mathbb{R}$ is a *tree metric* if there is an *X*-tree $\mathcal{T} = (T; \phi)$ and a weighting $w : E(T) \to \mathbb{R}_+$ such that for all $x, y$

$$\delta(x, y) = d_{(\mathcal{T};w)}(x, y) := \sum_{e \in P(\mathcal{T};x,y)} w(e),$$

where $P(\mathcal{T}; x, y)$ is the unique path between $\phi(x)$ and $\phi(y)$. The *tree metric representation* $(\mathcal{T}; w)$ of $\delta$ is unique (and efficiently computable).
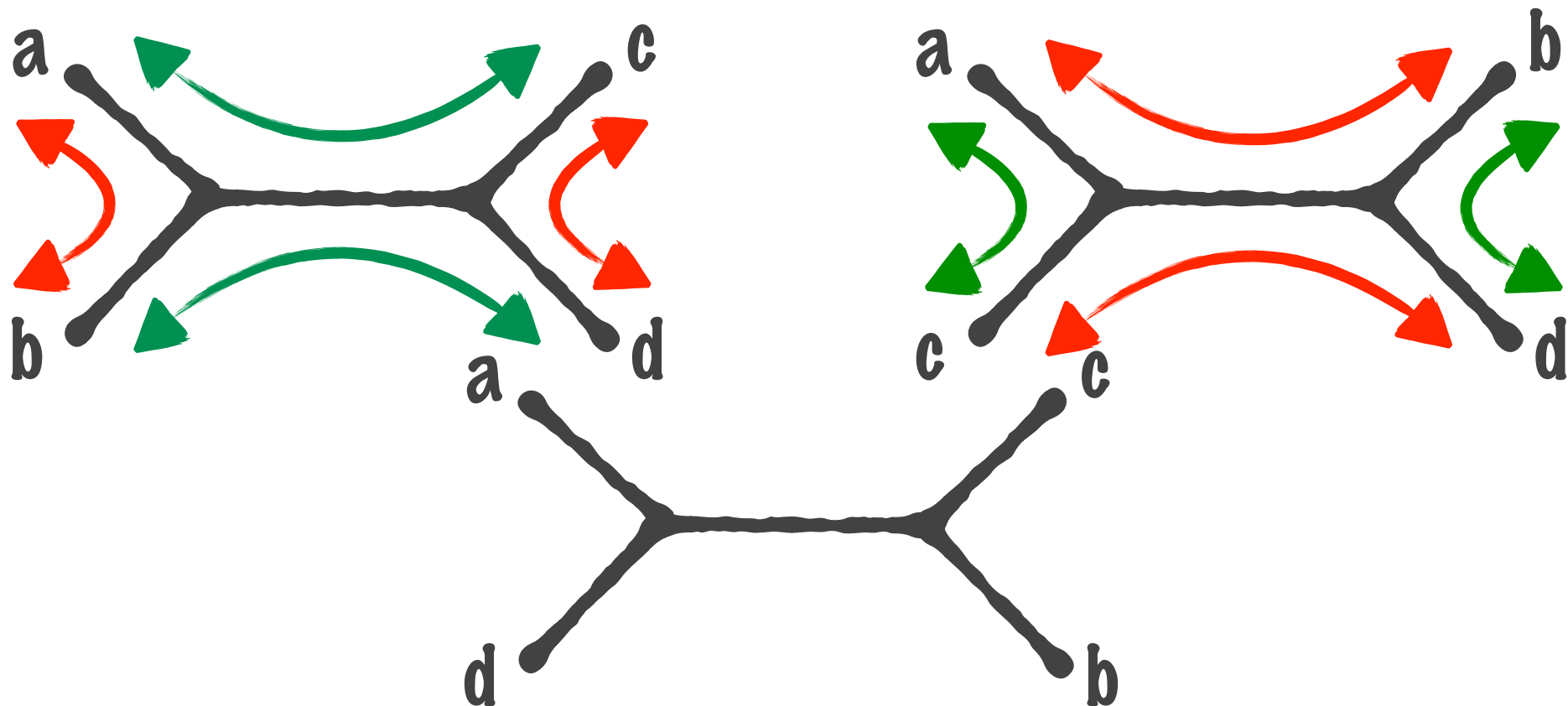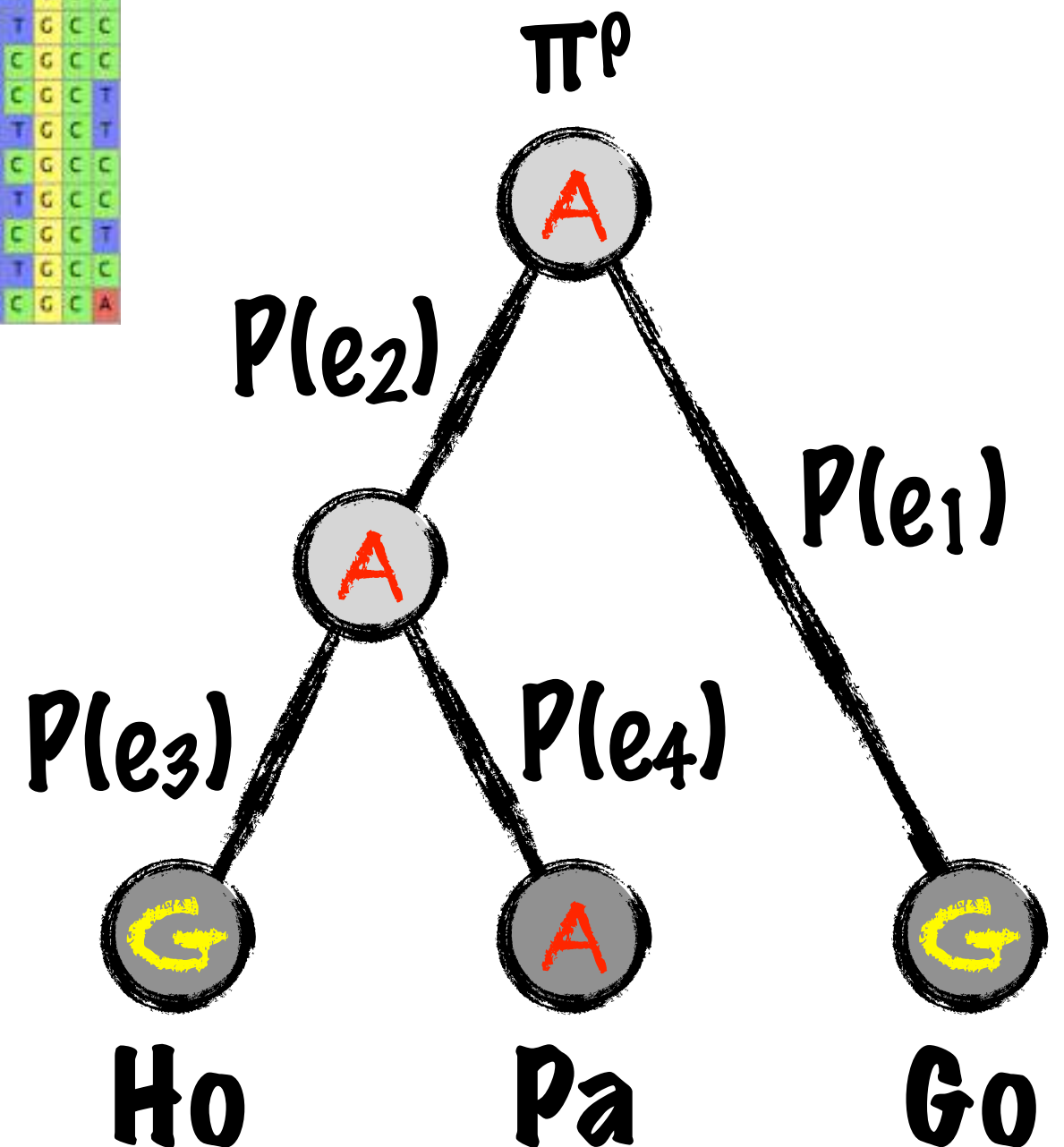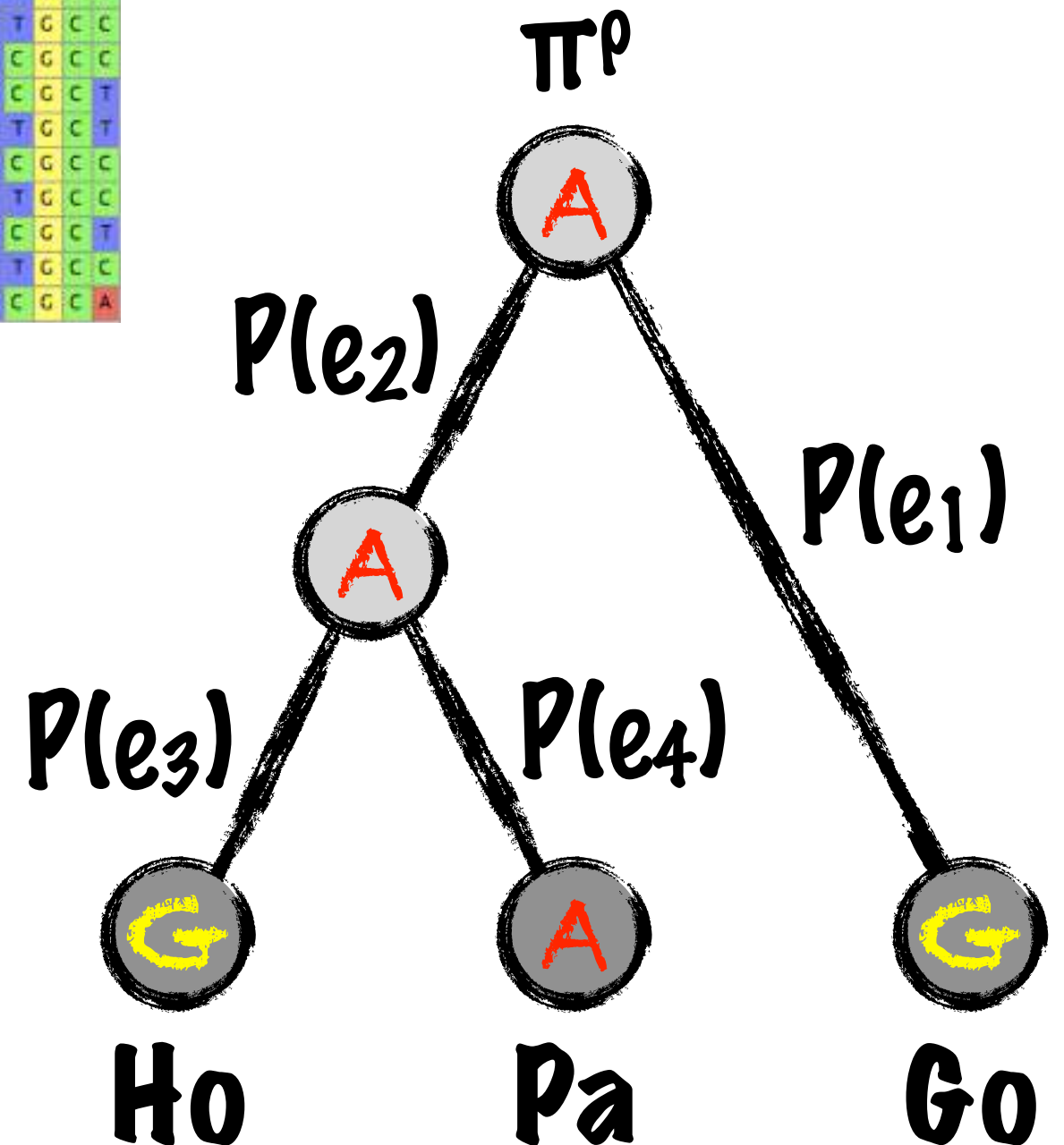
# Markov process on a tree

# Markov process on a tree

# Markov process on a tree

**Definition**

Let $\mathcal{T} = (T; \phi)$ be a phylogenetic $X$-tree with root $\rho$. Let $\pi^\rho$ be a distribution over $C = \{A, C, G, T\}$ and, for each $e \in E(T)$ (away from the root), let $P(e) = [P(e)_{\alpha,\beta}]$ be a Markov transition matrix over $C$. Let $\theta = (\pi^\rho; P(e), e \in E(T))$. The distribution of a state vector $\chi : X \to C$ at the leaves is defined as

$$p_\chi^{\mathcal{T}}(\theta) := \sum_{\substack{\bar{\chi}:V(T)\to C \\ \bar{\chi} \circ \phi = \chi}} \pi^\rho_{\bar{\chi}(\rho)} \prod_{e=(u,v)\in E(T)} P(e)_{\bar{\chi}(u),\bar{\chi}(v)}.$$

A natural choice is $P(e) = e^{\mu_e t_e Q}$ for a fixed rate matrix $Q$.

# Markov process on a tree

## Definition

Let $\mathcal{T} = (T; \phi)$ be a phylogenetic $X$-tree with root $\rho$. Let $\pi^\rho$ be a distribution over $C = \{A, C, G, T\}$ and, for each $e \in E(T)$ (away from the root), let $P(e) = [P(e)_{\alpha, \beta}]$ be a Markov transition matrix over $C$. Let $\theta = (\pi^\rho; P(e), e \in E(T))$. The distribution of a state vector $\chi : X \to C$ at the leaves is defined as

$$p_\chi^{\mathcal{T}}(\theta) := \sum_{\substack{\bar{\chi}:V(T)\to C \\ \bar{\chi} \circ \phi = \chi}} \pi^\rho_{\bar{\chi}(\rho)} \prod_{e=(u,v)\in E(T)} P(e)_{\bar{\chi}(u),\bar{\chi}(v)}.$$

A natural choice is $P(e) = e^{\mu_e t_e Q}$ for a fixed rate matrix $Q$.



k columns
=
k i.i.d. samples

# Back to tree metrics

## Definition

Let $F^{xy}$ be the matrix whose entries correspond to the joint distribution at the leaves $\phi(x)$ and $\phi(y)$. The *log-det distance* is

$$\delta(x, y) = -\log(\det(F^{xy})).$$

## Theorem (Steel, AML (1994))

*Assume $\pi^\rho > 0$ and $|\det P(e)| \neq 0, 1$ for all e. Then the log-det distance is a tree metric with corresponding X-tree $\mathcal{T}$.*

# Back to Darwin's finches



NJ tree of combined cytb and cr sequences.
(From: Akie Sato et al. PNAS 1999;96:5101-5106)

# Identifiability

Recall:

$$p_\chi^T(\theta) := \sum_{\substack{\bar{\chi}:V(T)\to C \\ \bar{\chi}\circ\phi=\chi}} \pi_{\bar{\chi}(\rho)}^\rho \prod_{e=(u,v)\in E(T)} P(e)_{\bar{\chi}(u),\bar{\chi}(v)}.$$

Let $n$ be the number of leaves.

## Definition
We say that the model is *identifiable* if, whenever $(\mathcal{T};\theta) \neq (\mathcal{T}';\theta')$, we have $p^T(\theta) \neq p^{T'}(\theta')$ as vectors in $\mathbb{R}^{4^n}$.

## Theorem (Steel, AML (1994); Chang, MB (1996))
*If $\pi^\rho > 0$ and $|\det P(e)| \neq 0, 1$ for all e, the model is* identifiable *(up to degeneracies).*

# Identifiability



$\mathbb{R}^{4^n}$

$p^T(\Theta)$

$p^{T'}(\Theta)$

## Definition

We say that the model is *identifiable* if, whenever $(\mathcal{T}; \theta) \neq (\mathcal{T}'; \theta')$, we have $p^{\mathcal{T}}(\theta) \neq p^{\mathcal{T}'}(\theta')$ as vectors in $\mathbb{R}^{4^n}$.

## Theorem (Steel, AML (1994); Chang, MB (1996))

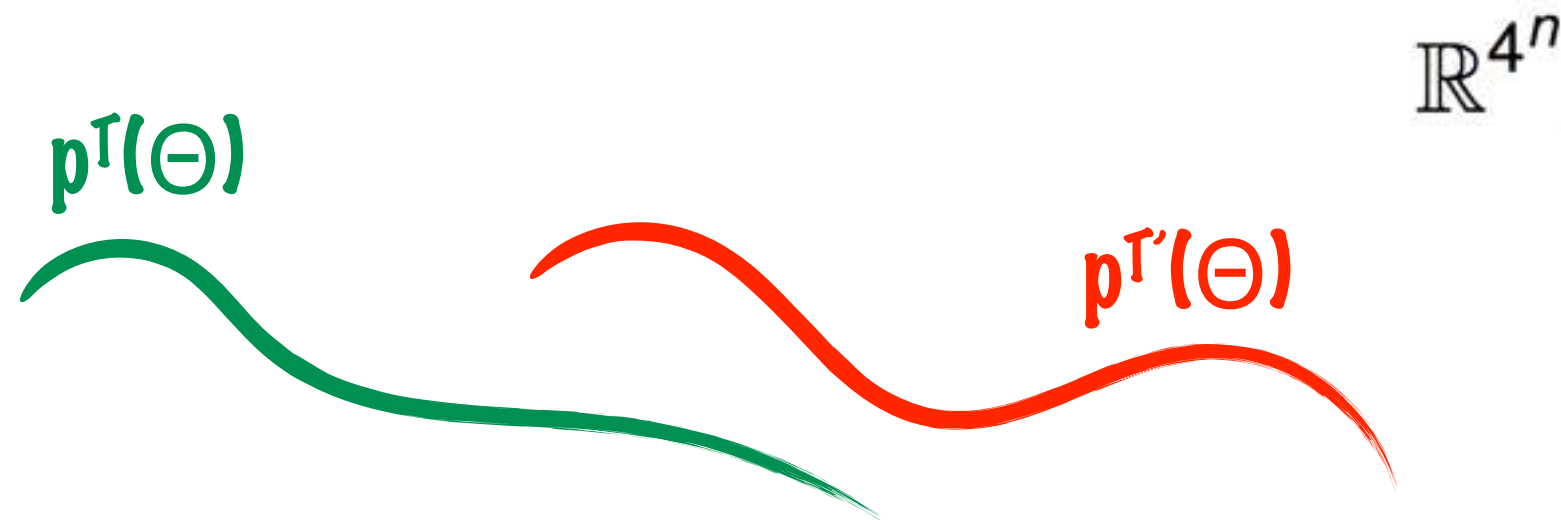If $\pi^\rho > 0$ and $|\det P(e)| \neq 0, 1$ for all $e$, the model is identifiable (up to degeneracies).

# Likelihood-based inference

## Definition

Given sequences of length $k$, i.e., $(\chi^i)_{i=1}^k$, the maximum likelihood estimator (MLE) is

$$\hat{\mathcal{T}} \in \arg\max \left\{ \prod_{i=1}^k p_{\chi^i}^{\mathcal{T}}(\theta) \; : \; \mathcal{T}, \theta \in \Theta \right\}.$$

## Theorem (Chang, MB (1996))

*The MLE is* consistent, *i.e.,* $\hat{\mathcal{T}} \to \mathcal{T}$ *as* $k \to +\infty$.

## Theorem (Chor-Tuller, JACM (2006); Roch, TCBB (2006))

*Computing the MLE is NP-hard.*

How much data do I need?

# Adaptive radiation



Genome-scale phylogeny of birds. (From: Erich D. Jarvis et al. Science 2014;346:1320-1331)

# Short branches

Theorem (Steel & Székely, SIDMA (2002))

*Under the symmetric 2-state Markov model on four taxa with internal branch of weight $f$, reconstructing the phylogeny with high probability requires $k = \Omega(f^{-2})$ as $f \to 0$.*

# Short branches

## Theorem (Steel & Székely, SIDMA (2002))

*Under the symmetric 2-state Markov model on four taxa with internal branch of weight $f$, reconstructing the phylogeny with high probability requires $k = \Omega(f^{-2})$ as $f \to 0$.*

# Short branches

## Theorem (Steel & Székely, SIDMA (2002))

*Under the symmetric 2-state Markov model on four taxa with internal branch of weight f, reconstructing the phylogeny with high probability requires $k = \Omega(f^{-2})$ as $f \to 0$.*

# Short branches
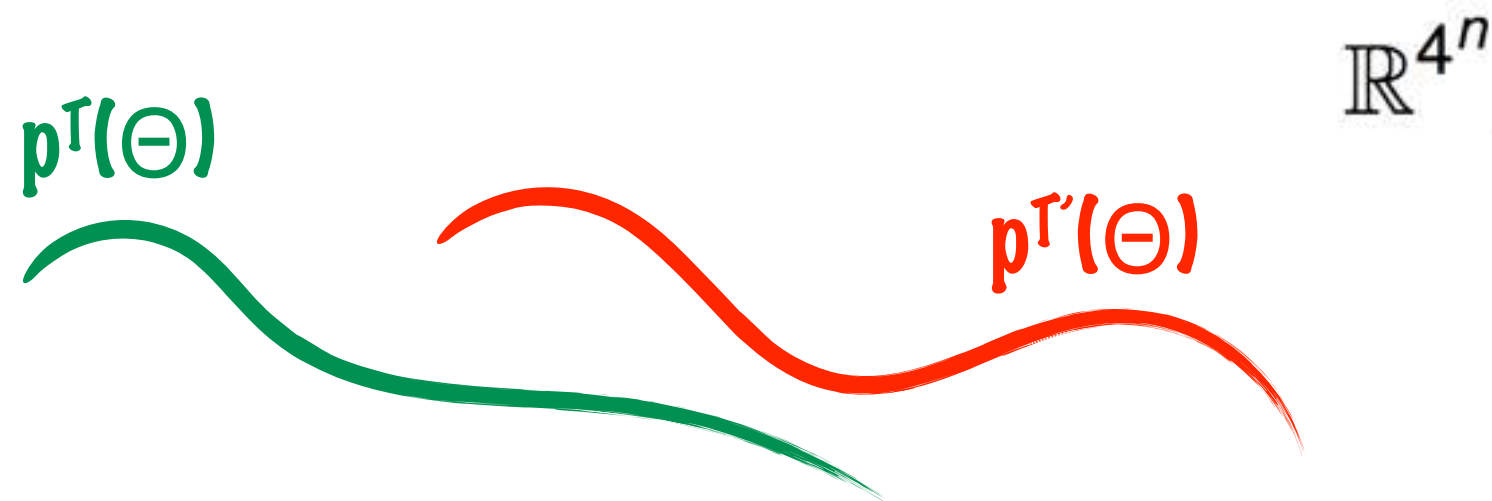
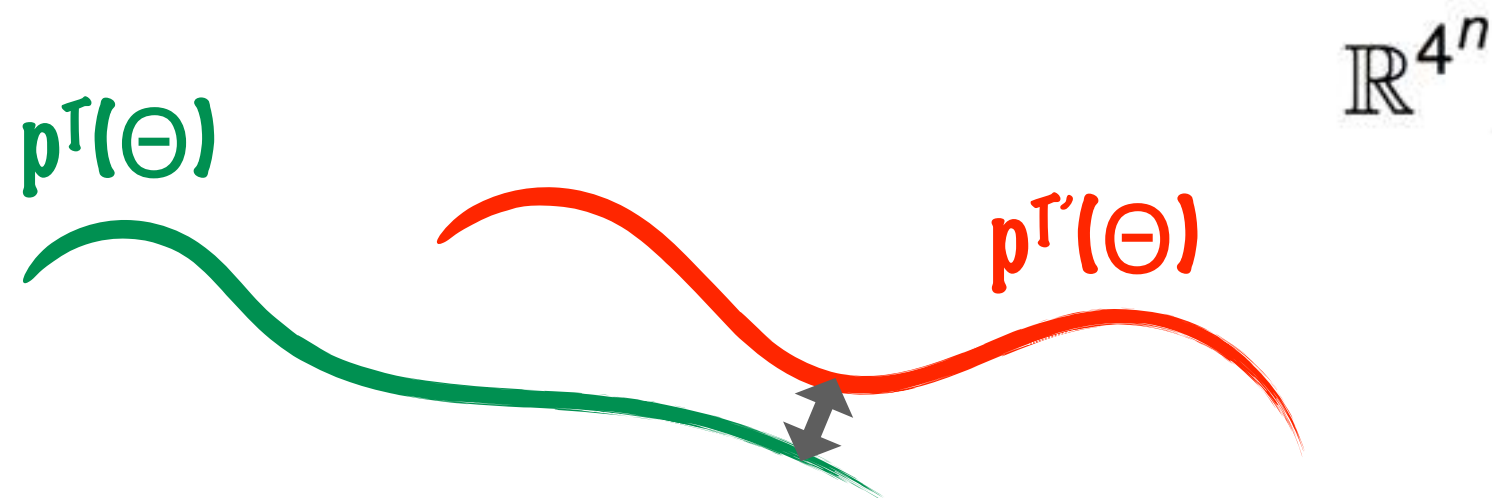## Theorem (Steel & Székely, SIDMA (2002))

*Under the symmetric 2-state Markov model on four taxa with internal branch of weight $f$, reconstructing the phylogeny with high probability requires $k = \Omega(f^{-2})$ as $f \to 0$.*

# Short branches

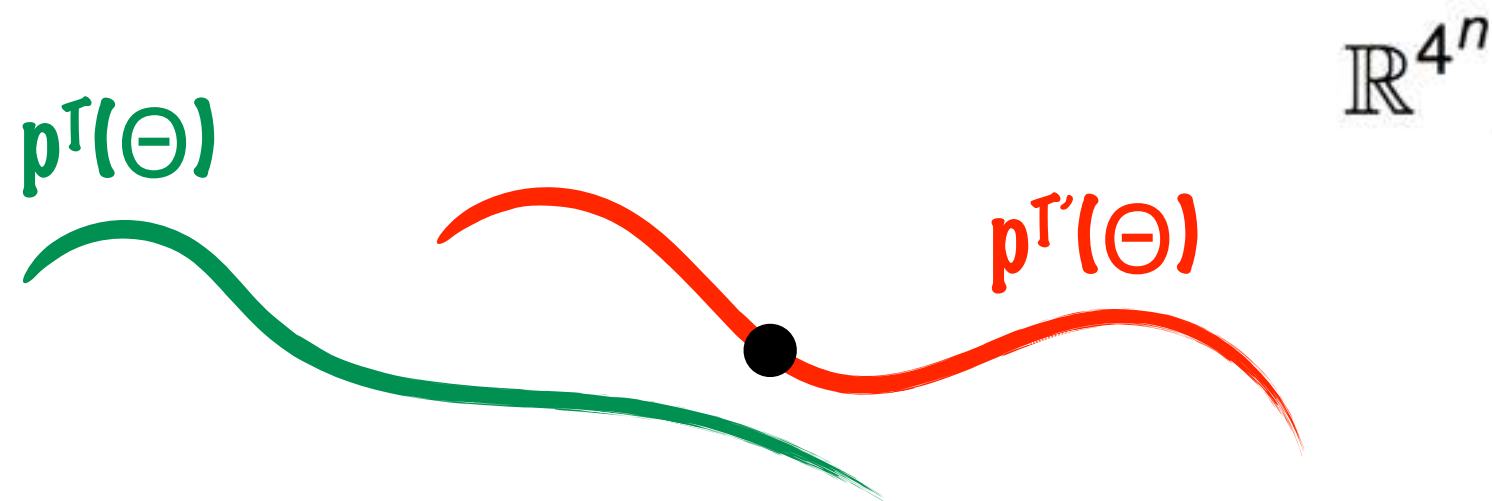## Theorem (Steel & Székely, SIDMA (2002))

*Under the symmetric 2-state Markov model on four taxa with internal branch of weight $f$, reconstructing the phylogeny with high probability requires $k = \Omega(f^{-2})$ as $f \to 0$.*

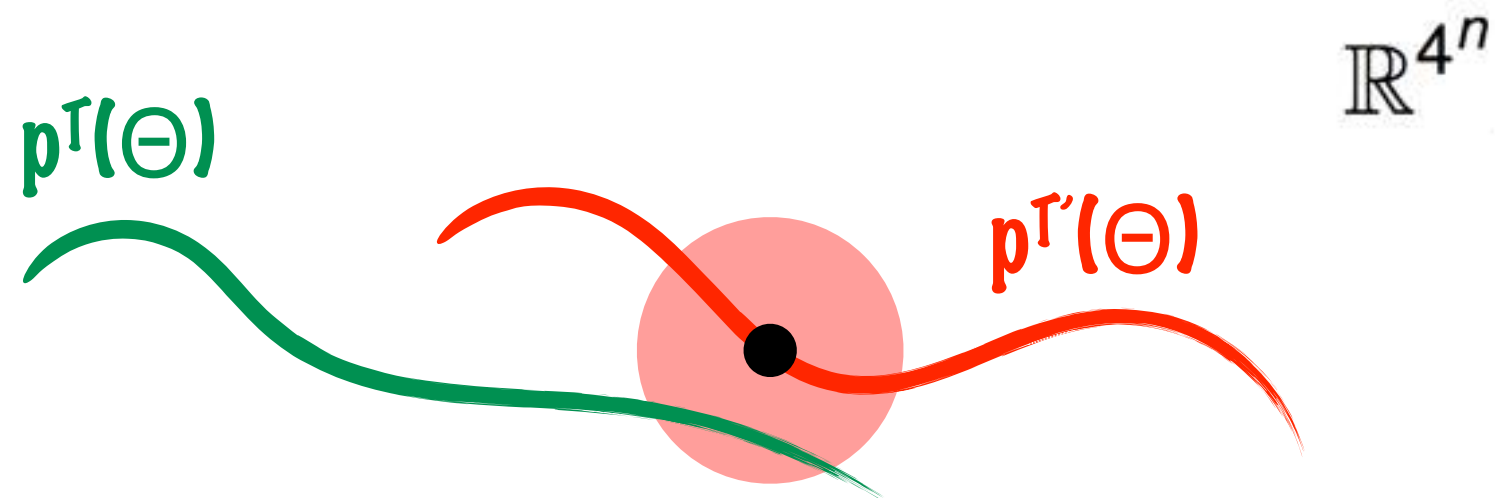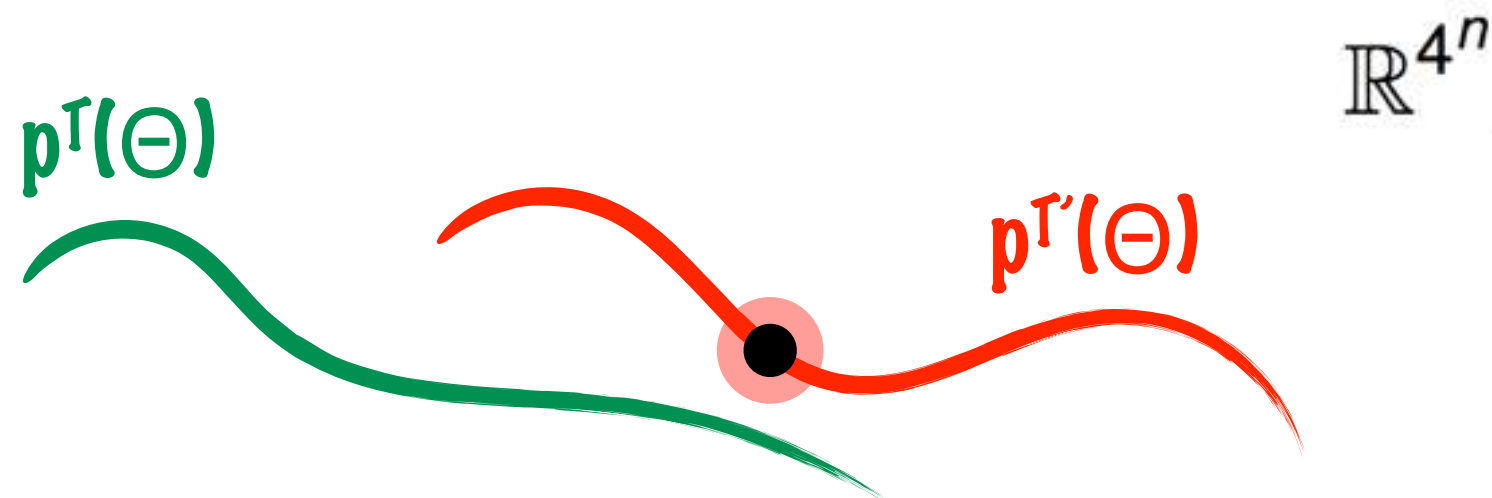# Short branches

**Theorem (Steel & Székely, SIDMA (2002))**

*Under the symmetric 2-state Markov model on four taxa with internal branch of weight f, reconstructing the phylogeny with high probability requires $k = \Omega(f^{-2})$ as $f \to 0$.*

$\mathbb{R}^{4^n}$

$\mathbf{p}^T(\Theta)$

$\mathbf{p}^{T'}(\Theta)$

# Depth

A special case of a more general phenomenon:

Theorem (Mossel, TAMS (2004))

*Under the symmetric 2-state Markov model on n taxa with branches of weight f, reconstructing the phylogeny with high probability requires in general*

$$k = \begin{cases} \Theta(f^{-2} \log n), & \text{if } f < f^*, \\ n^{\Theta(f)}, & \text{if } f \geq f^*. \end{cases}$$

Matched for MLE (Roch & Sly (2015)) and some tree metric methods (Roch, Science (2010)). In contrast, NJ requires an exponential in *n* amount of data.
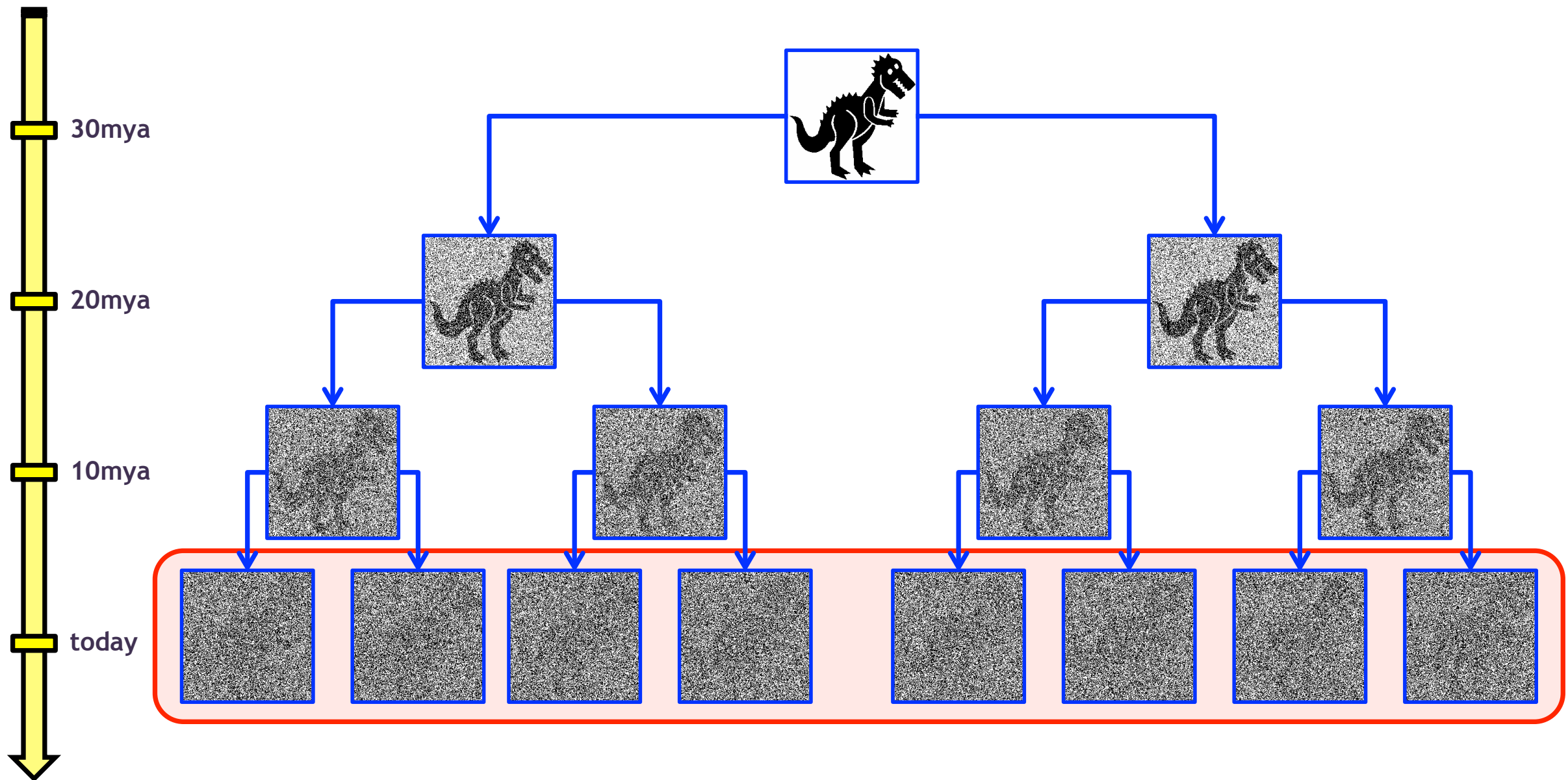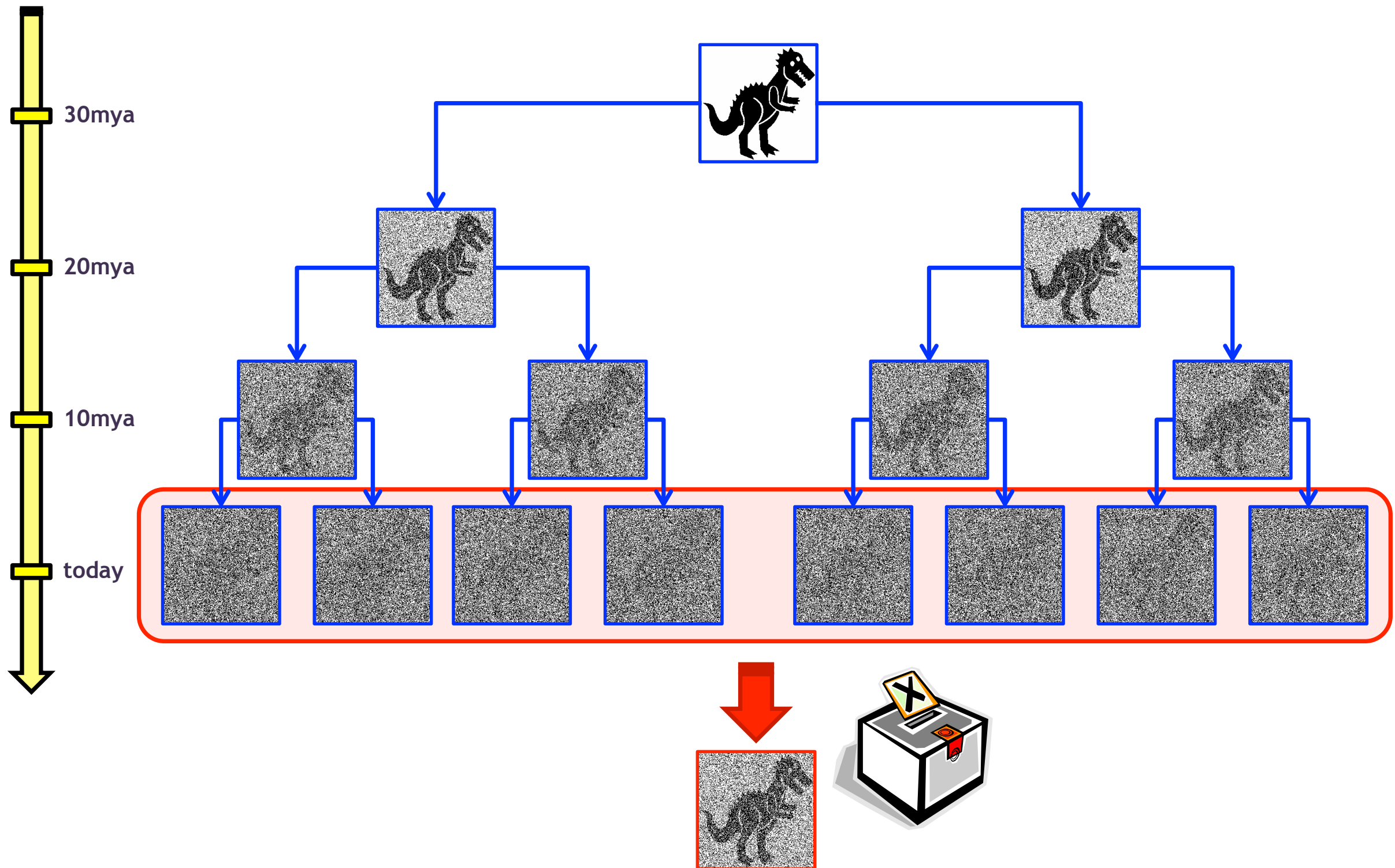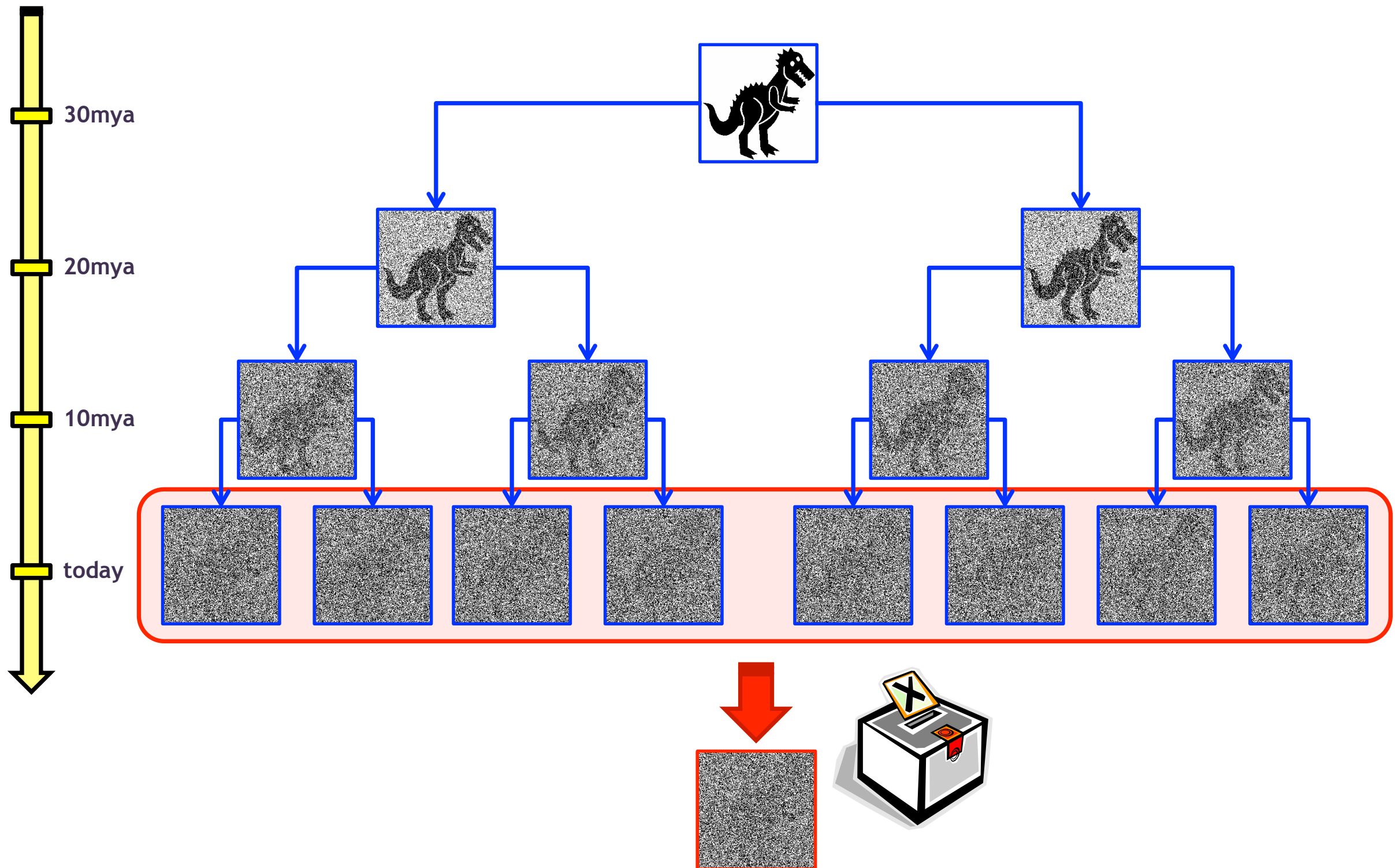
# Correlation decay

# Correlation decay
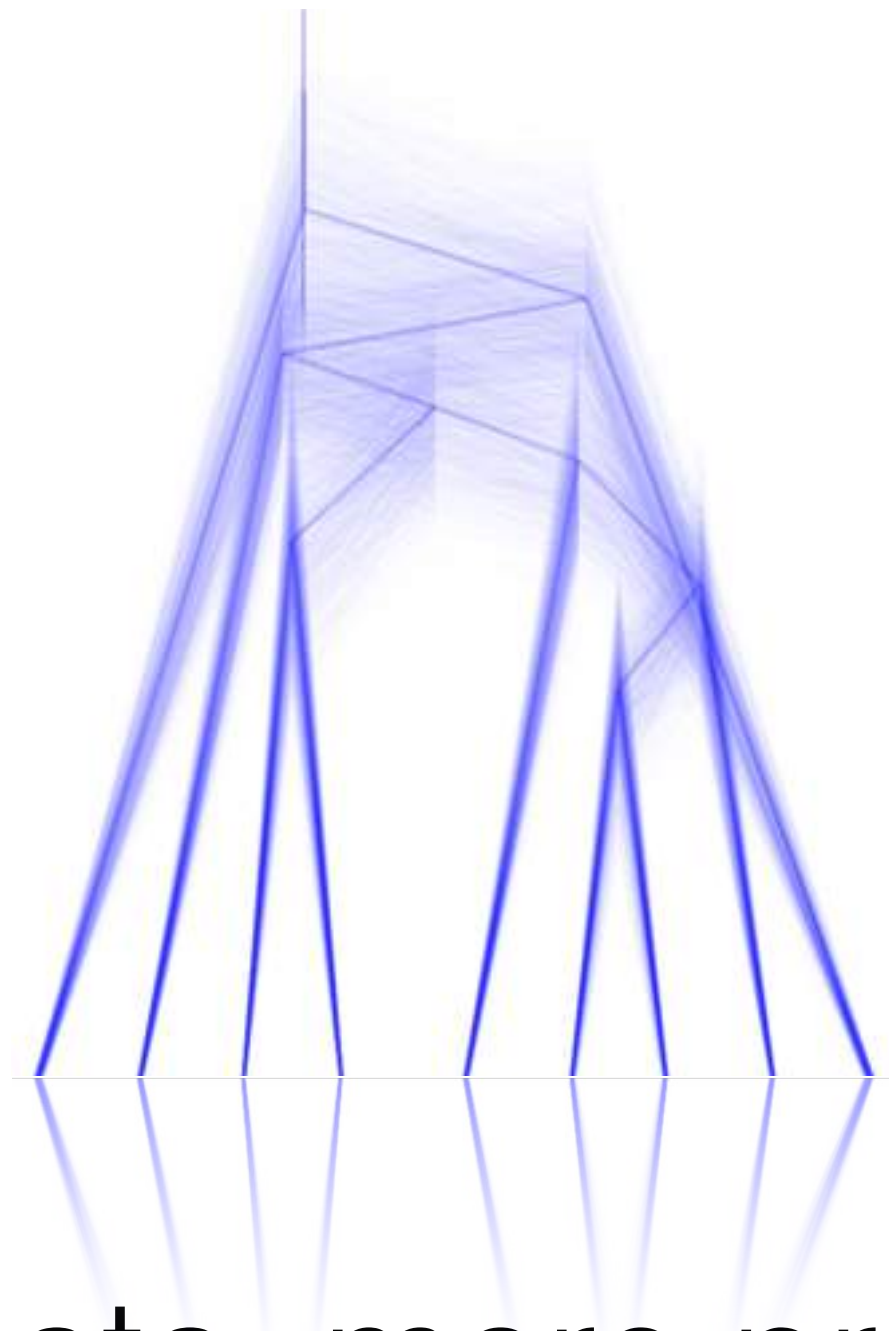
# Correlation decay

# Correlation decay

# Correlation decay

More data, more problems

# Next-generation sequencing

# Concatenating genes

# Concatenating genes

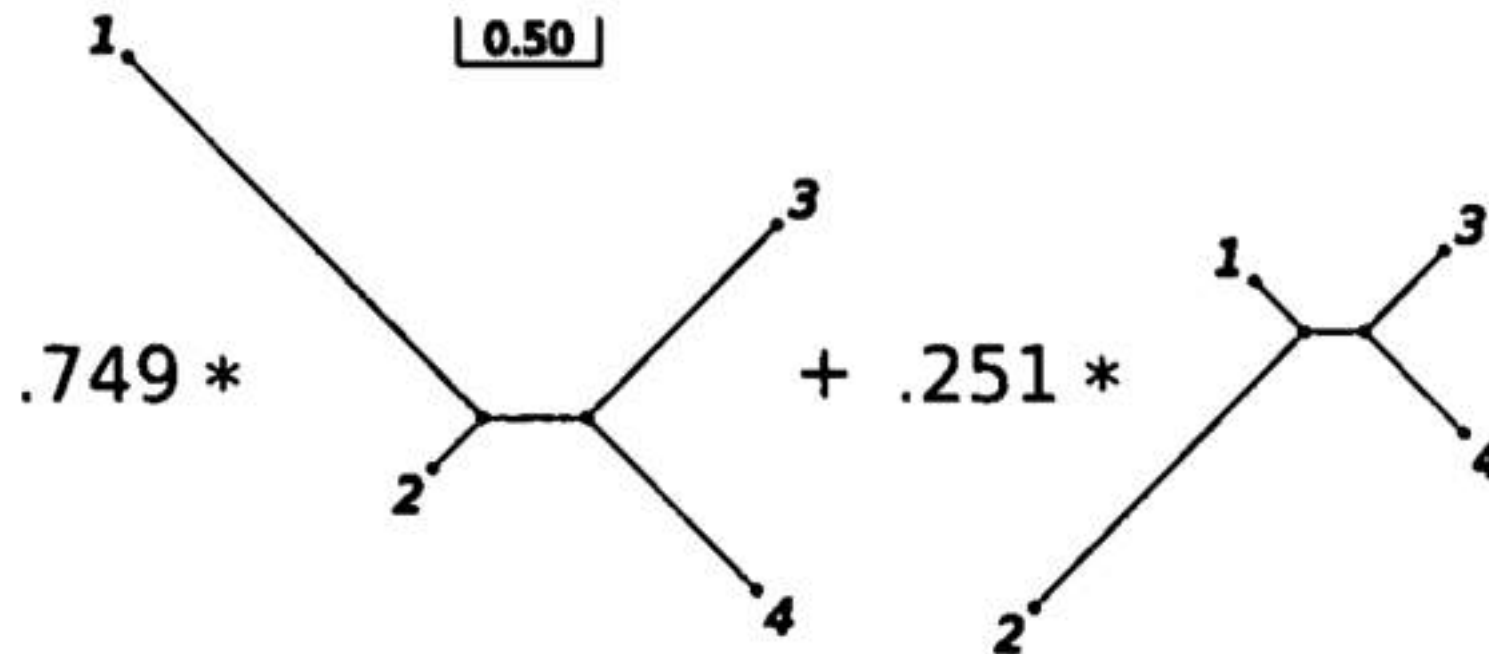Supergene of length mk

# Concatenating genes



ML

Supergene of length mk

# Mixed-up trees

Using algebraic geometry (Sturmfels & Sullivant, JCB (2005)):

Theorem (Matsen & Steel, SB (2007))

*Phylogenetic mixtures on a single tree can mimic a tree of another topology.*

# Mixed-up trees

Using algebraic geometry (Sturmfels & Sullivant, JCB (2005)):

Theorem (Matsen & Steel, SB (2007))

*Phylogenetic mixtures on a single tree can mimic a tree of another topology.*
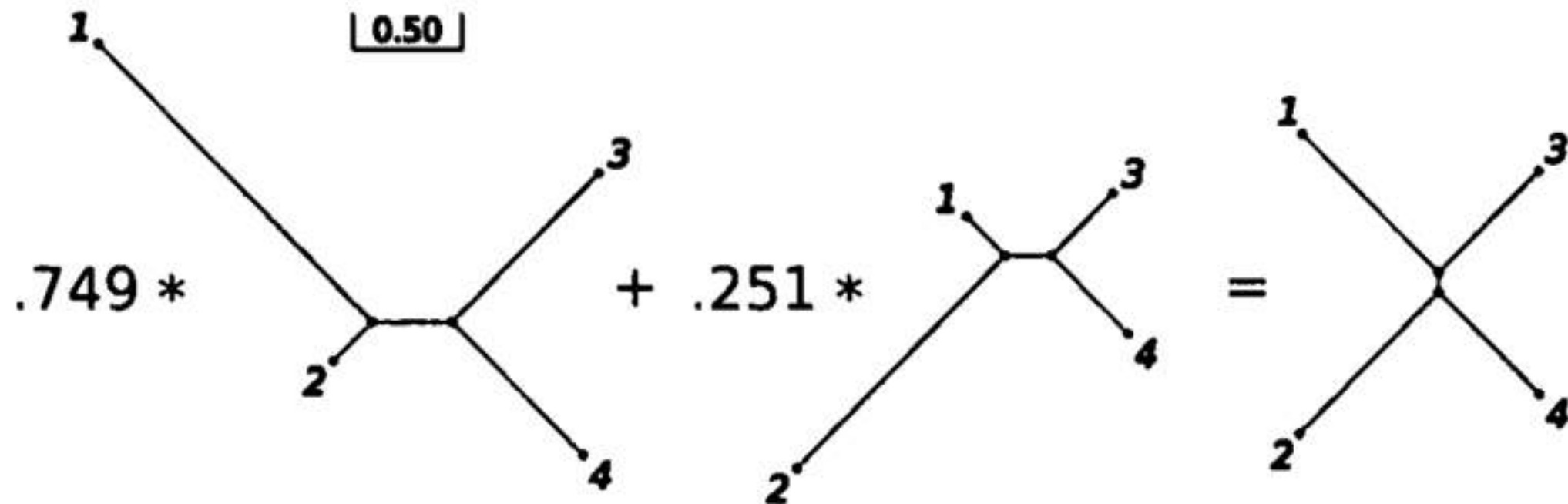


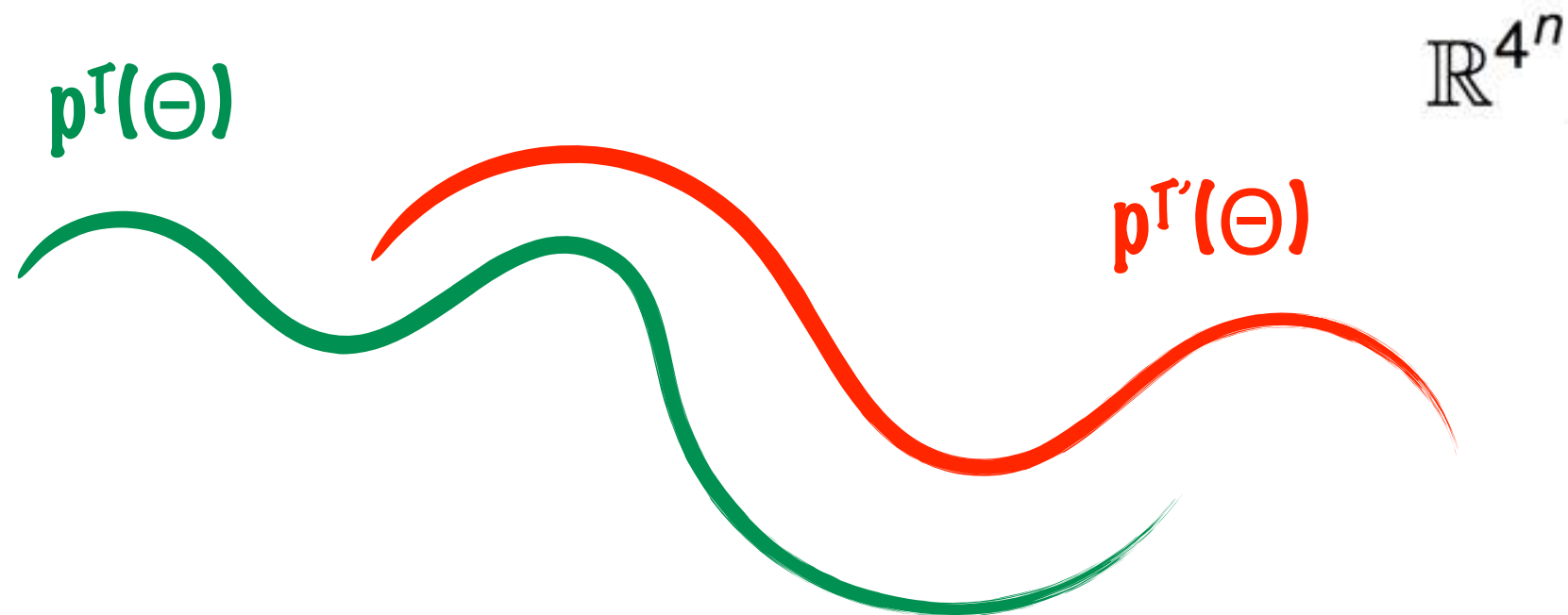.749 * [tree] + .251 * [tree] = [tree]

# Mixed-up trees

Using algebraic geometry (Sturmfels & Sullivant, JCB (2005)):

**Theorem (Matsen & Steel, SB (2007))**

*Phylogenetic mixtures on a single tree can mimic a tree of another topology.*

$$\mathbb{R}^{4^n}$$

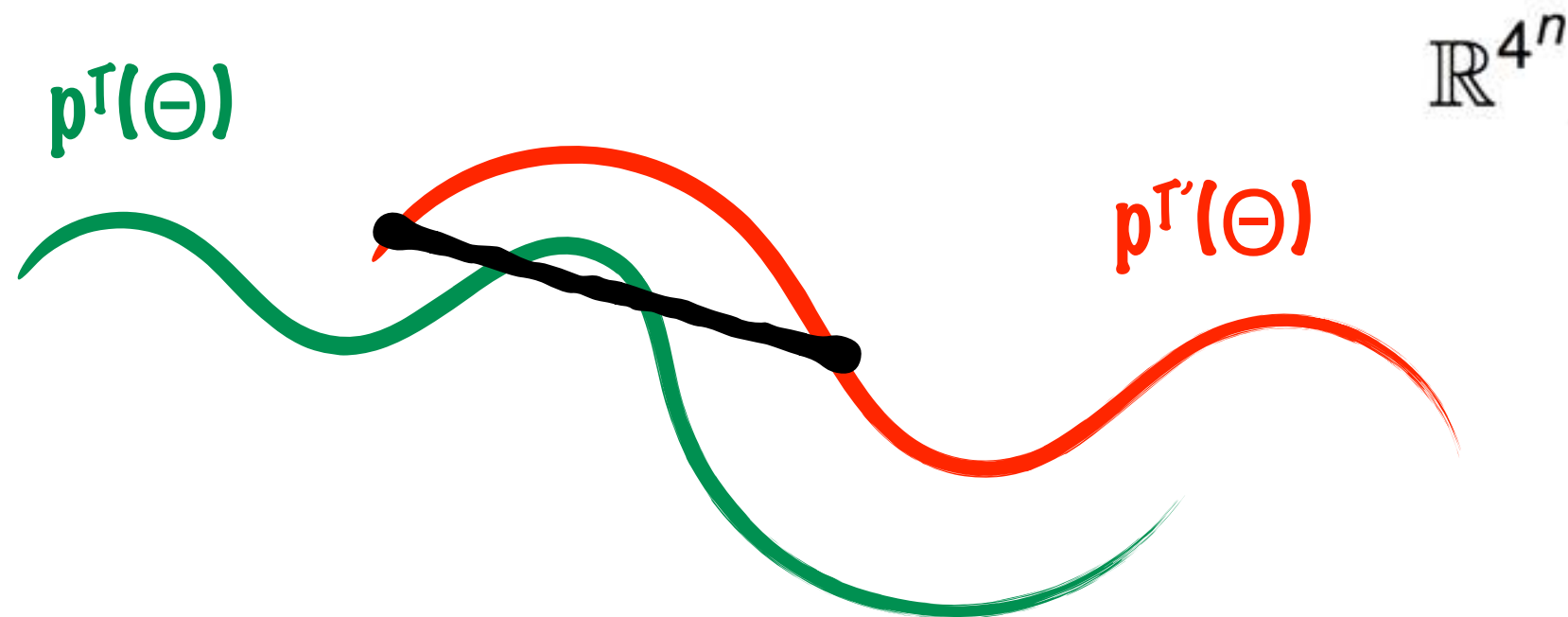$$\mathbf{p}^T(\Theta)$$

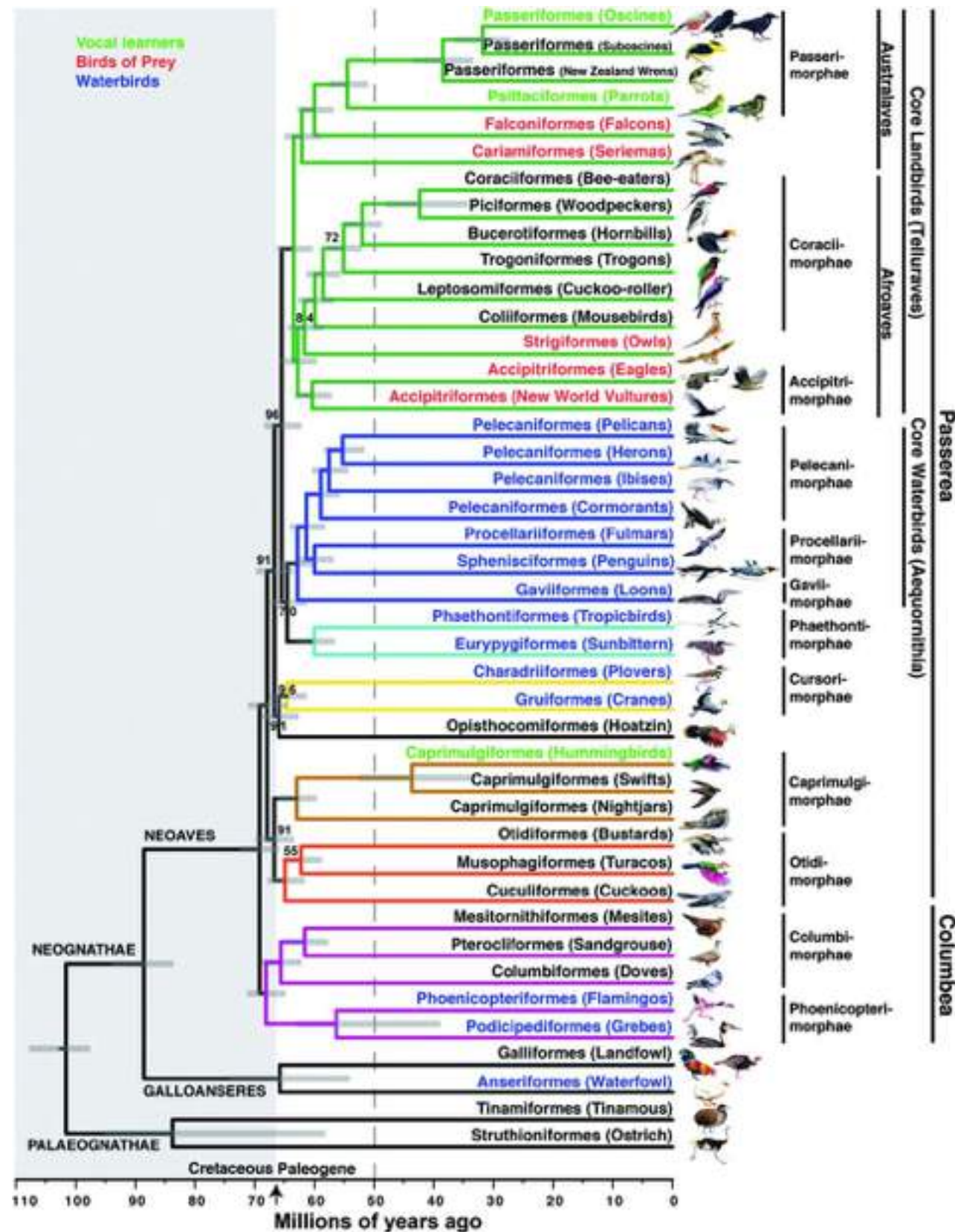$$\mathbf{p}^{T'}(\Theta)$$

# Mixed-up trees

Using algebraic geometry (Sturmfels & Sullivant, JCB (2005)):

Theorem (Matsen & Steel, SB (2007))

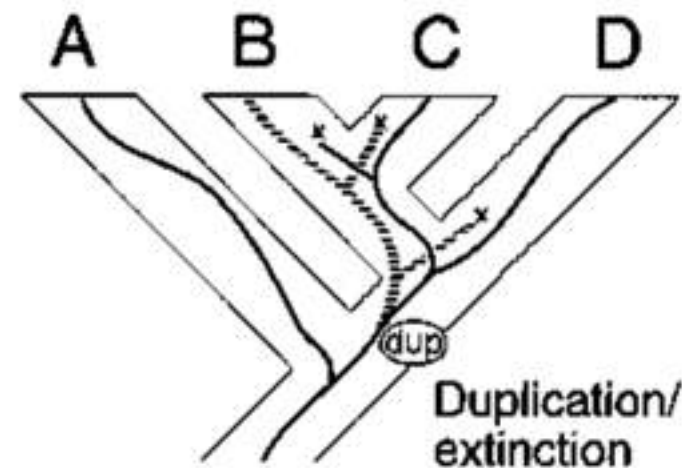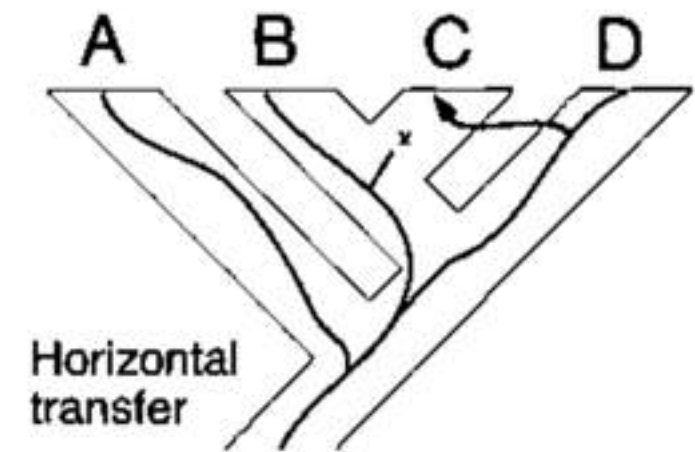*Phylogenetic mixtures on a single tree can mimic a tree of another topology.*
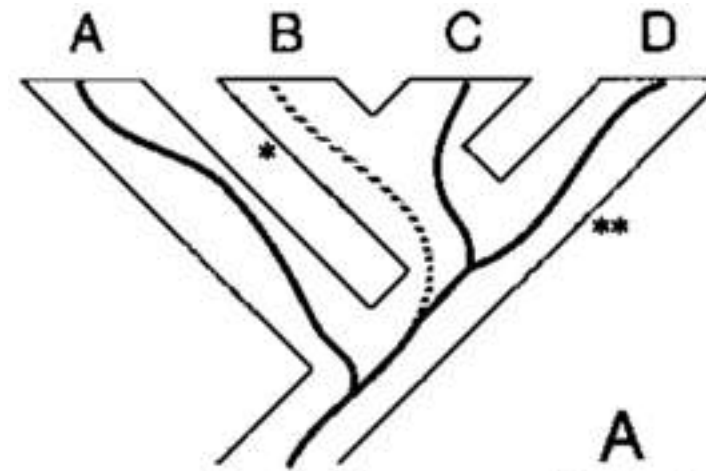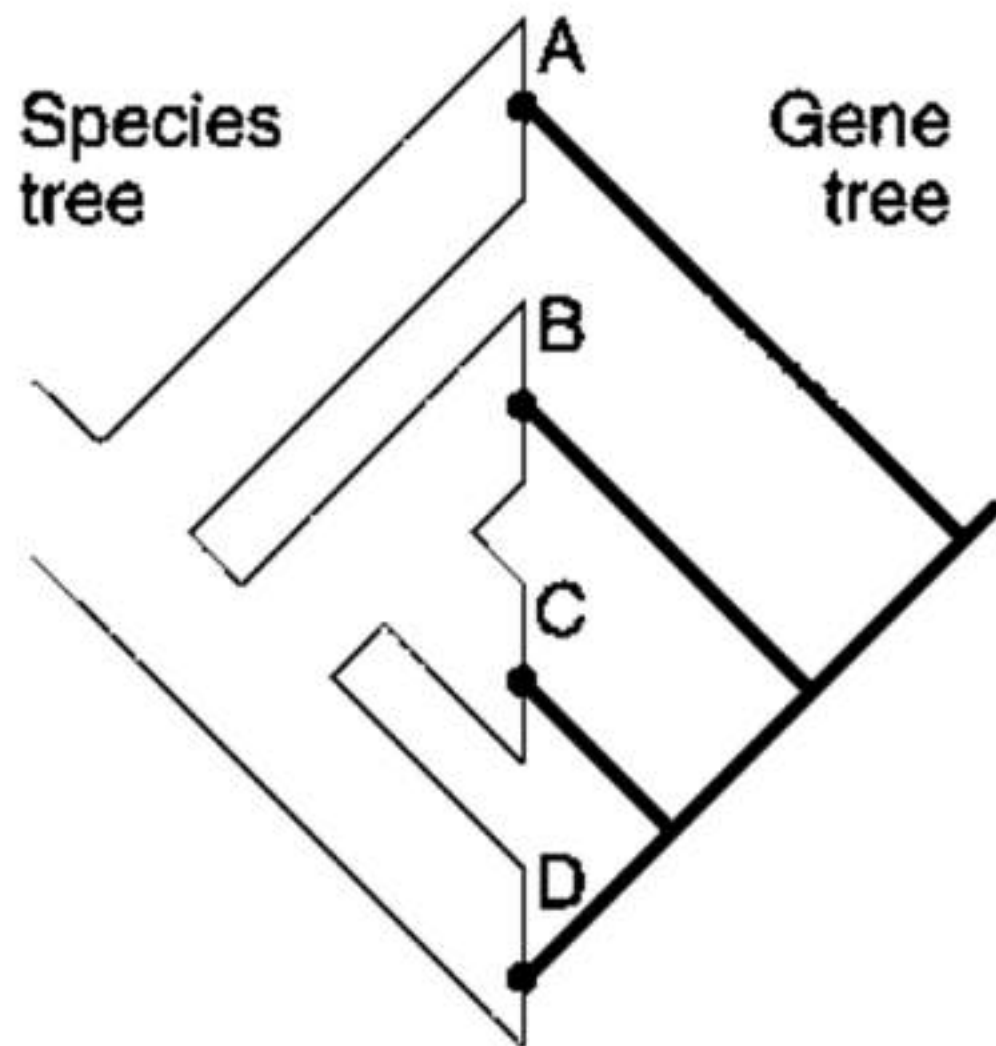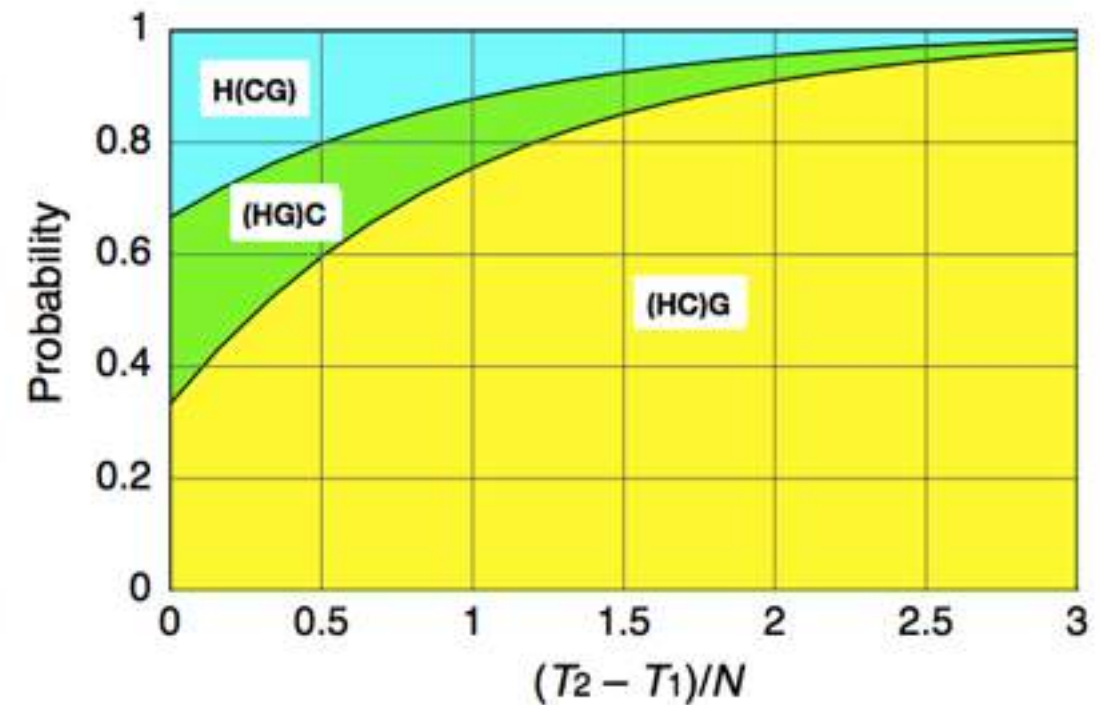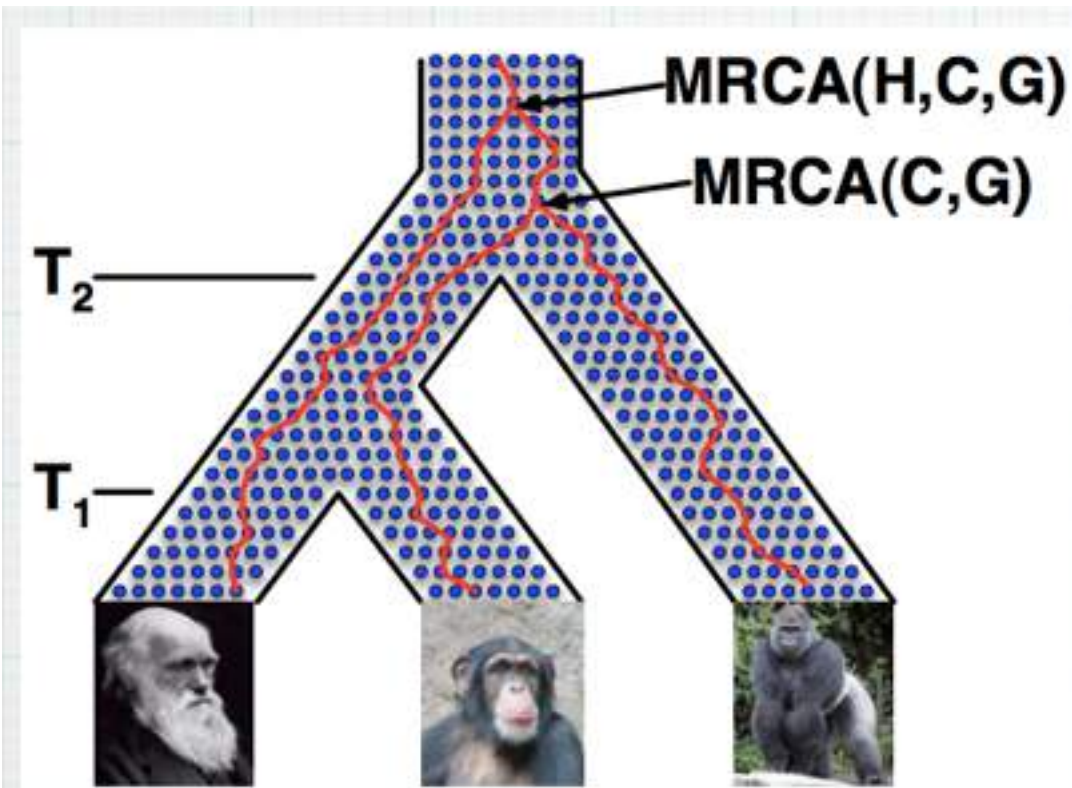
# Back to the birds



Genome-scale phylogeny of birds. (From: Erich D. Jarvis et al. Science 2014;346:1320-1331)

# Species tree v. "gene" trees



Species tree

Gene tree

Horizontal transfer

Duplication/extinction

# A source of discordance:
# Deep coalescence



$$\mathbf{P}[((H,C),G)] = 1 - \frac{2}{3}e^{-(T_2-T_1)/N}$$

$$\mathbf{P}[((H,G),C)] = \frac{1}{3}e^{-(T_2-T_1)/N}$$

$$\mathbf{P}[(H,(C,G))] = \frac{1}{3}e^{-(T_2-T_1)/N}$$
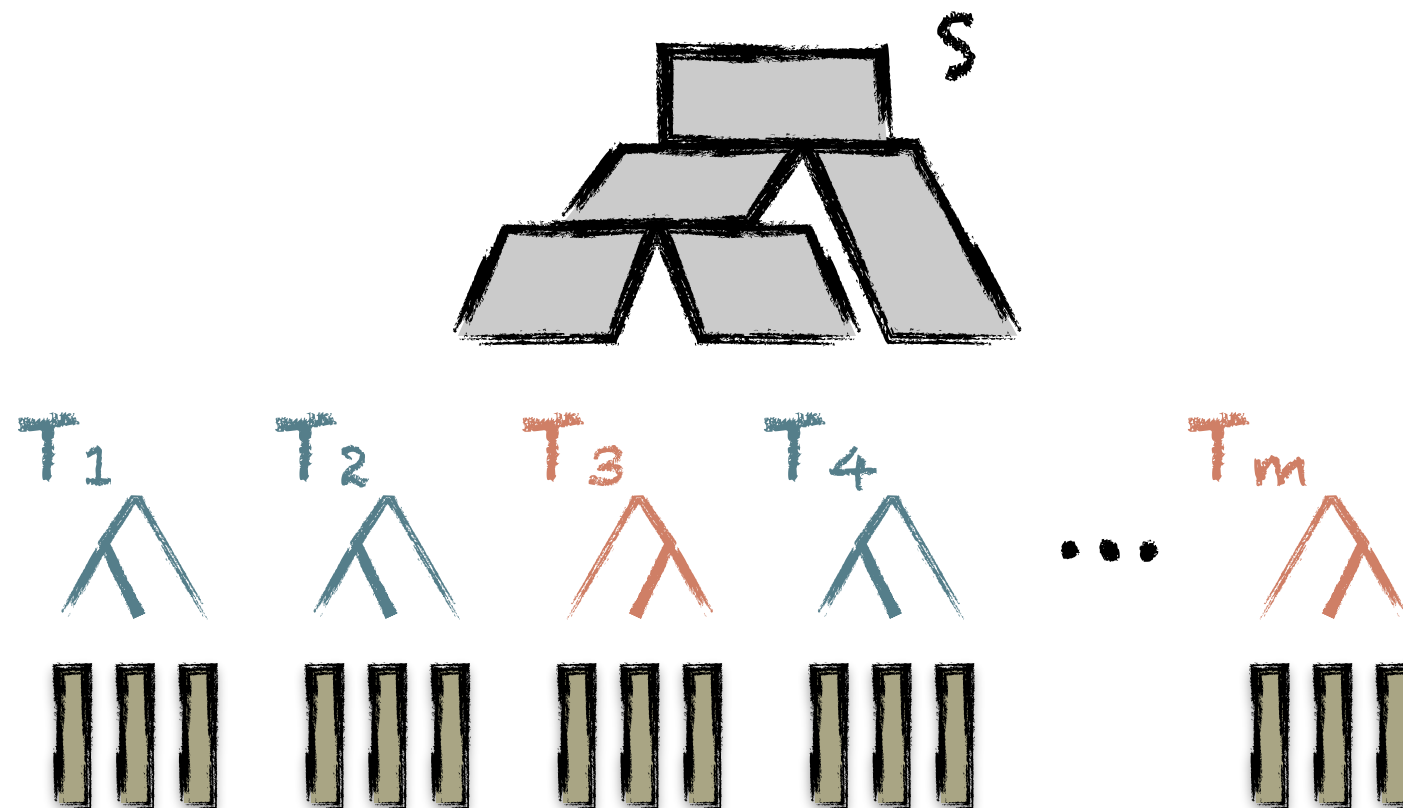
# Anomaly zone



**Definition (Degnan & Rosenberg (2006))**

The *anomaly zone* is the region of the parameter space in the multispecies coalescent where the most likely gene tree topology does not coincide with the species tree.
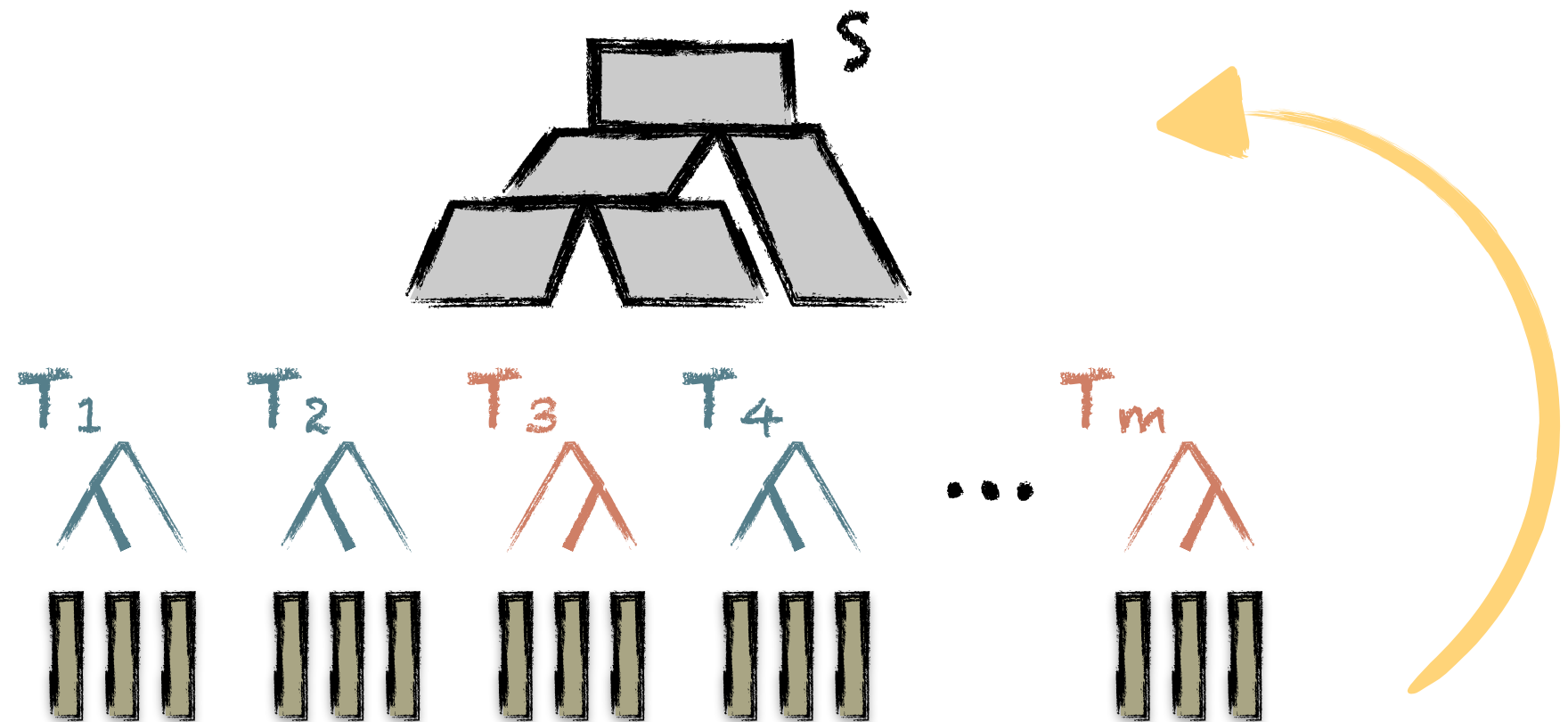
# An extra layer

- Species tree: S

- Two-stage hierarchical model: for each gene g (independently and identically),

  - Generate a gene tree $T_g$ for g using the multispecies coalescent on S

  - Generate sequence data of length k on $T_g$ using a Markov model
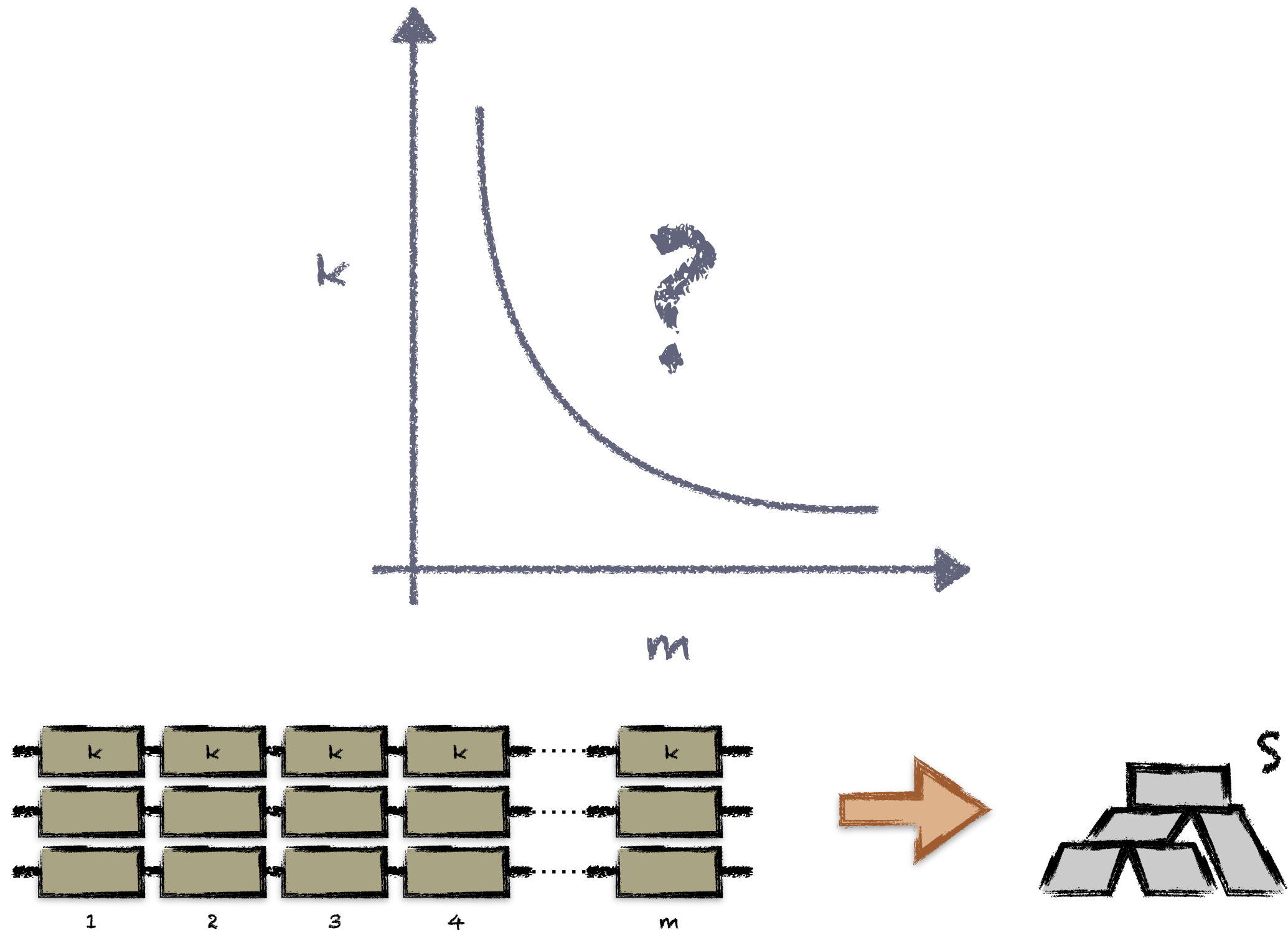
- Goal: recover S from sequences

# An extra layer

- Species tree: S

- Two-stage hierarchical model: for each gene g (independently and identically),

  - Generate a gene tree $T_g$ for g using the multispecies coalescent on S

  - Generate sequence data of length k on $T_g$ using a Markov model
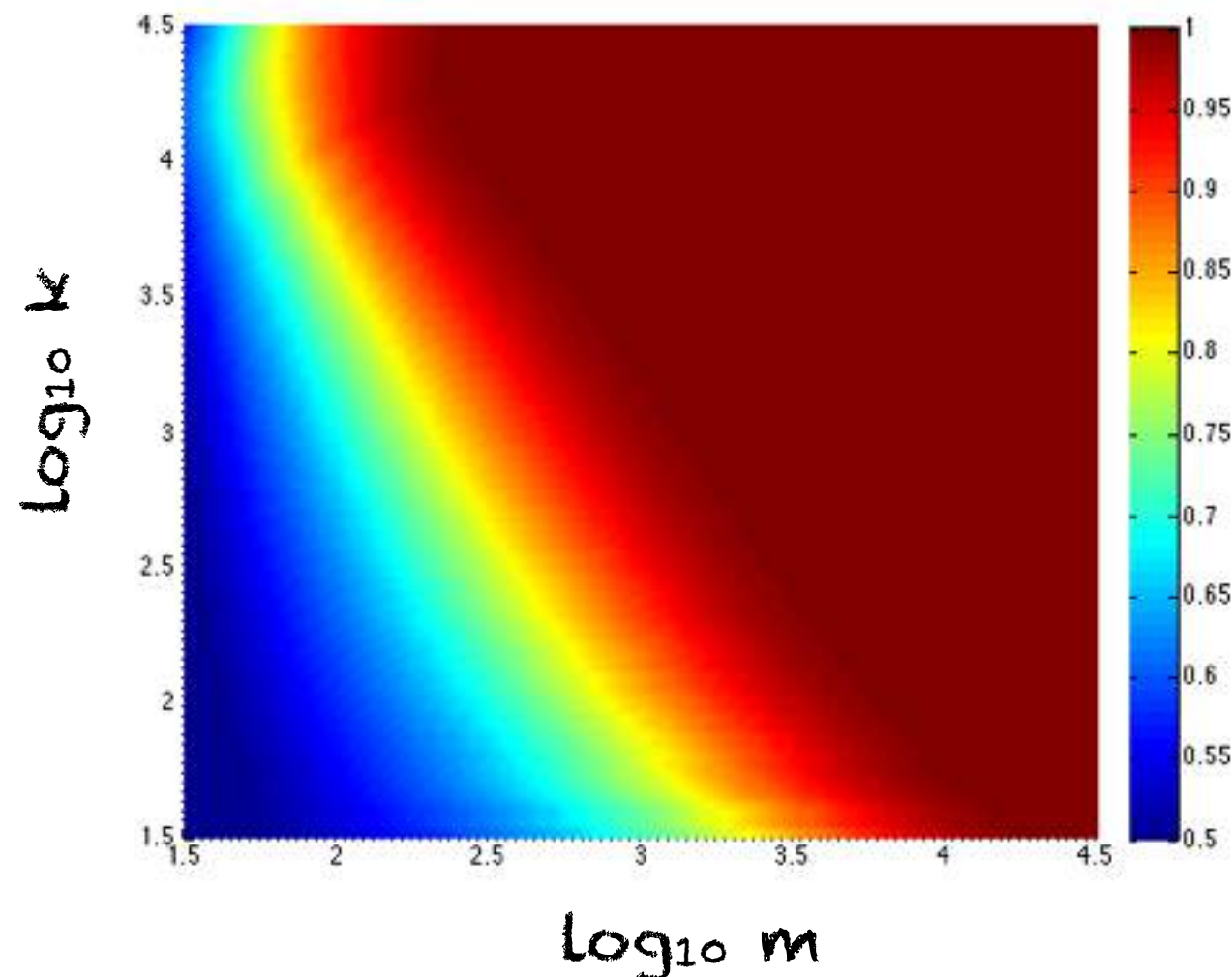
- Goal: recover S from sequences

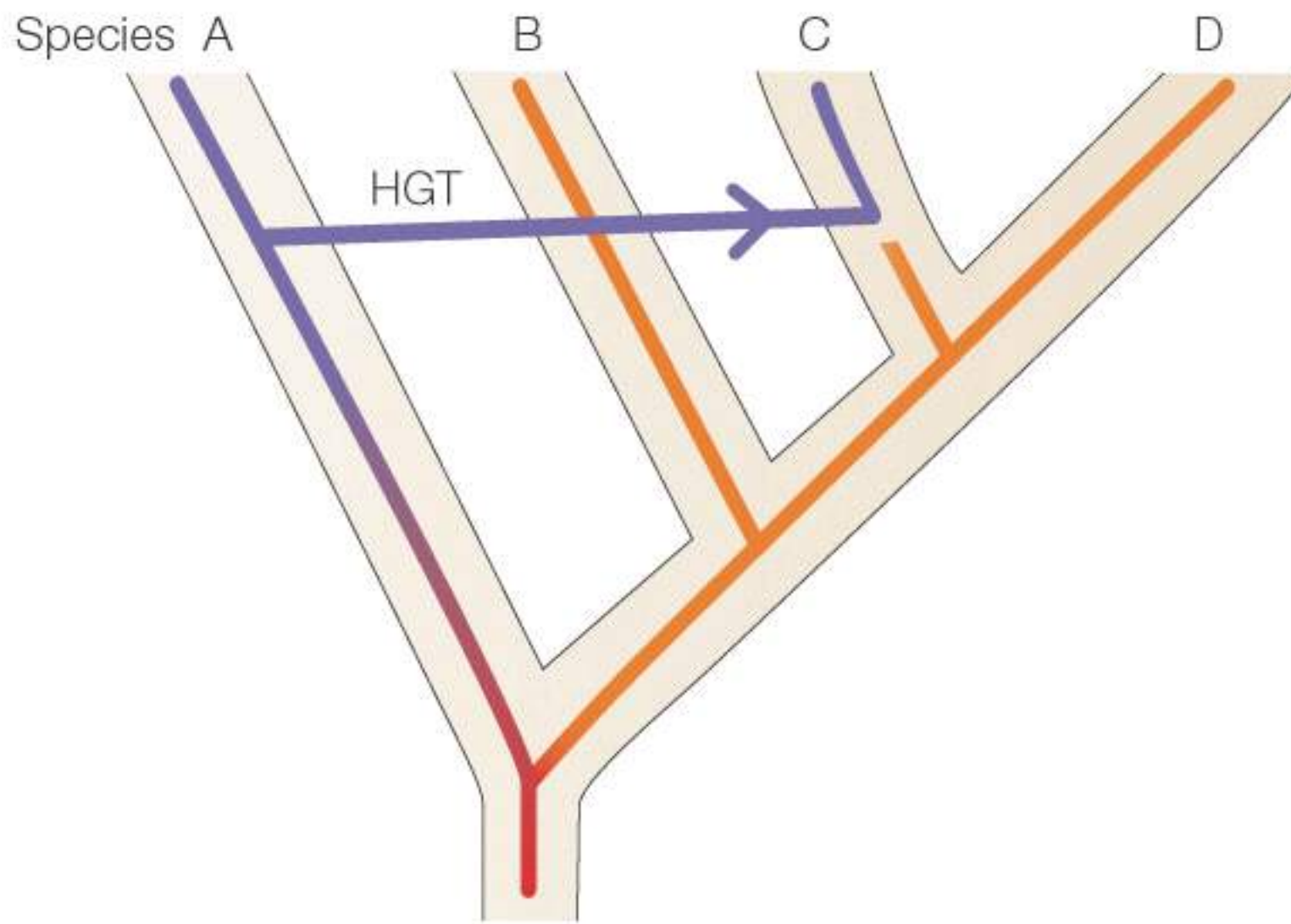# Question:
# How much data is needed?

# An unexpected trade-off
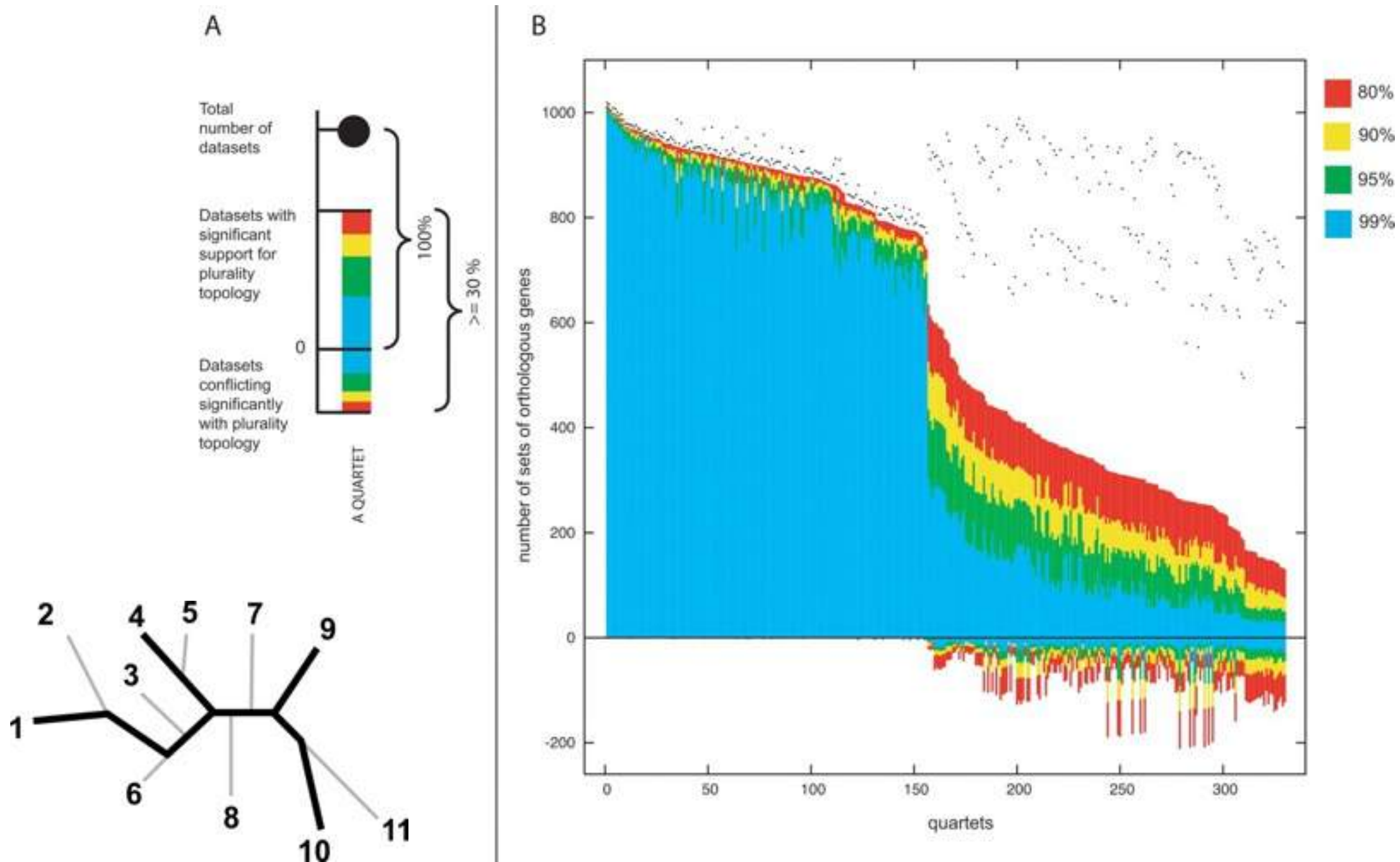
**Theorem (Mossel & R. (2015))**

*Under the 3-taxon multispecies coalescent with 4-state symmetric Markov model, reconstructing the species tree requires $m = \Theta\left(f^{-2}/\sqrt{k}\right)$ when $k = O(f^{-2})$ as $f \to 0$.*
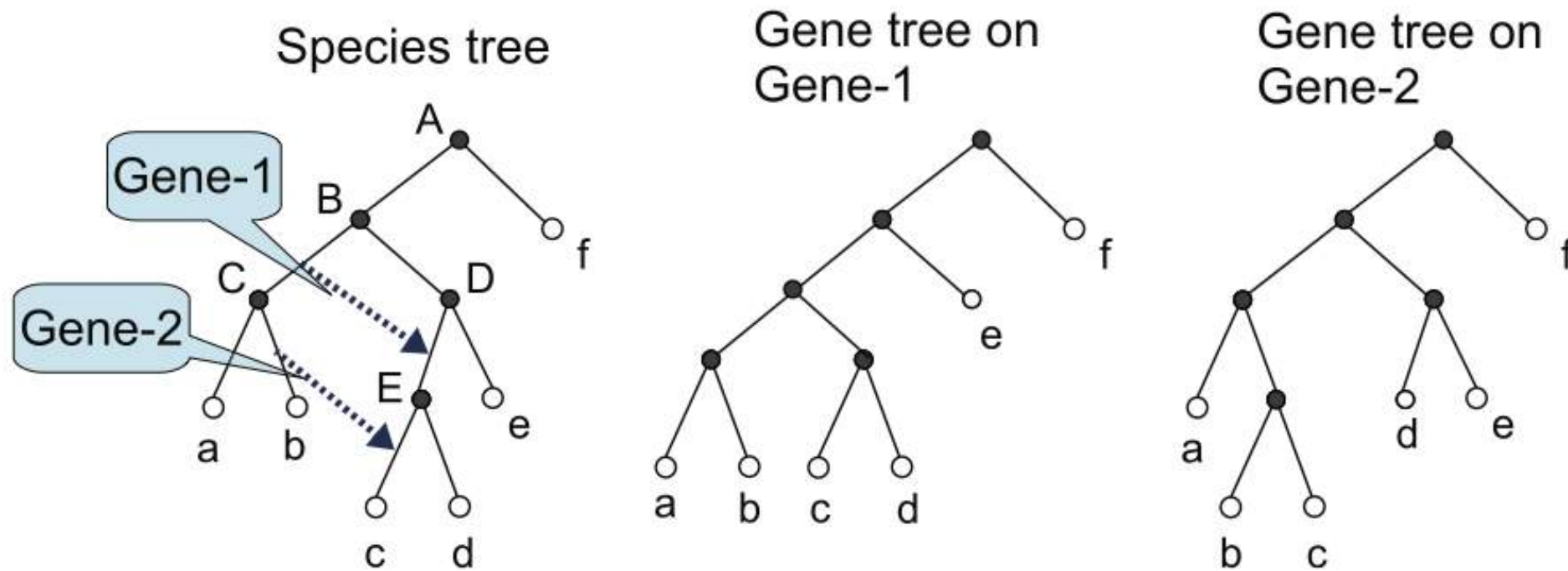
# A source of discordance: Horizontal gene transfer (HGT)

# Cyanobacteria



Quartet decomposition analysis of cyanobacteria.
(From: Olga Zhaxybayeva et al. Genome Res. 2006;16:1099-1108)
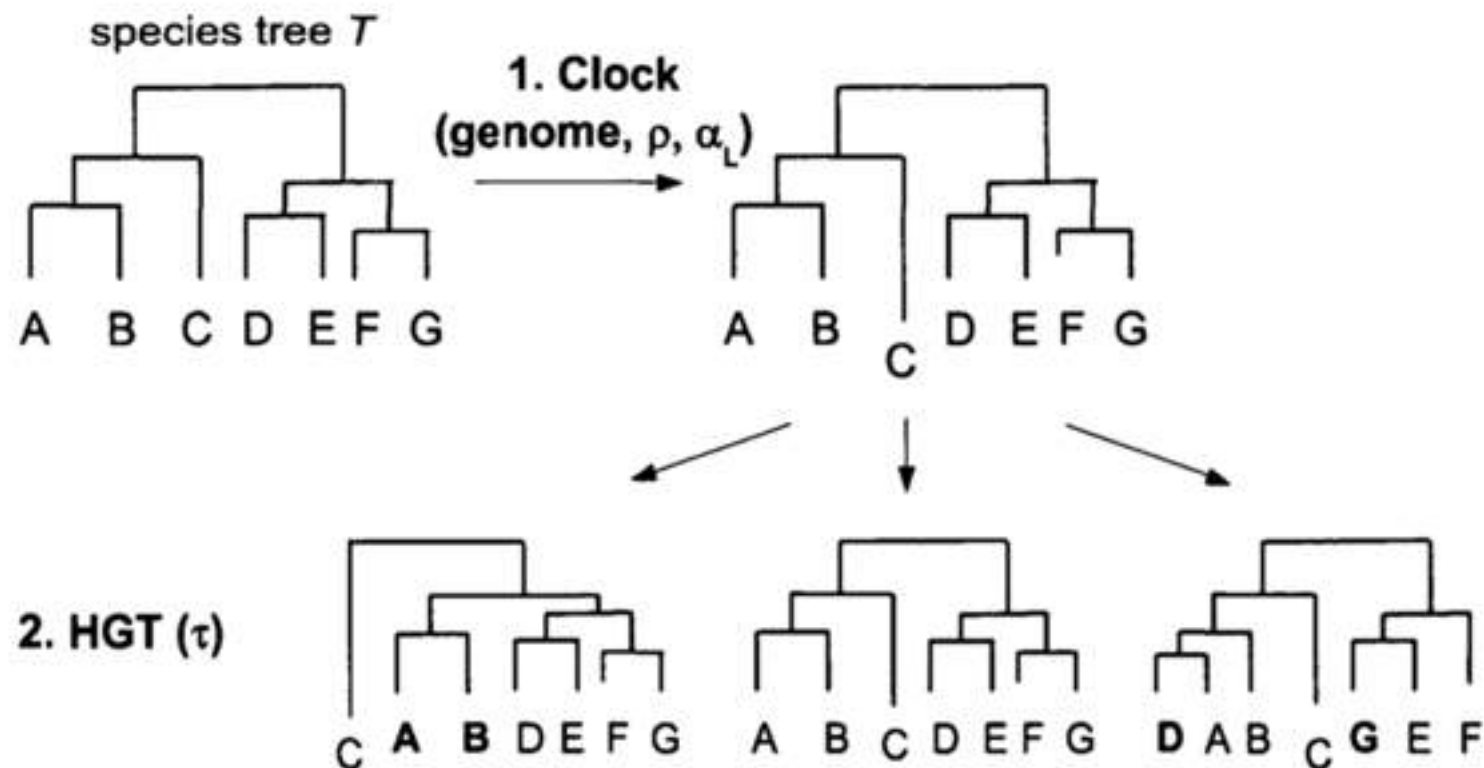
# Subtree-prune-regraft

# HGT as combinatorial noise

- Species tree: T

- Galtier's model: for each gene g (independently and identically),

    - HGTs occur at random positions with average number ρ of HGTs per gene

    - Receivers are chosen at random among contemporaneous positions
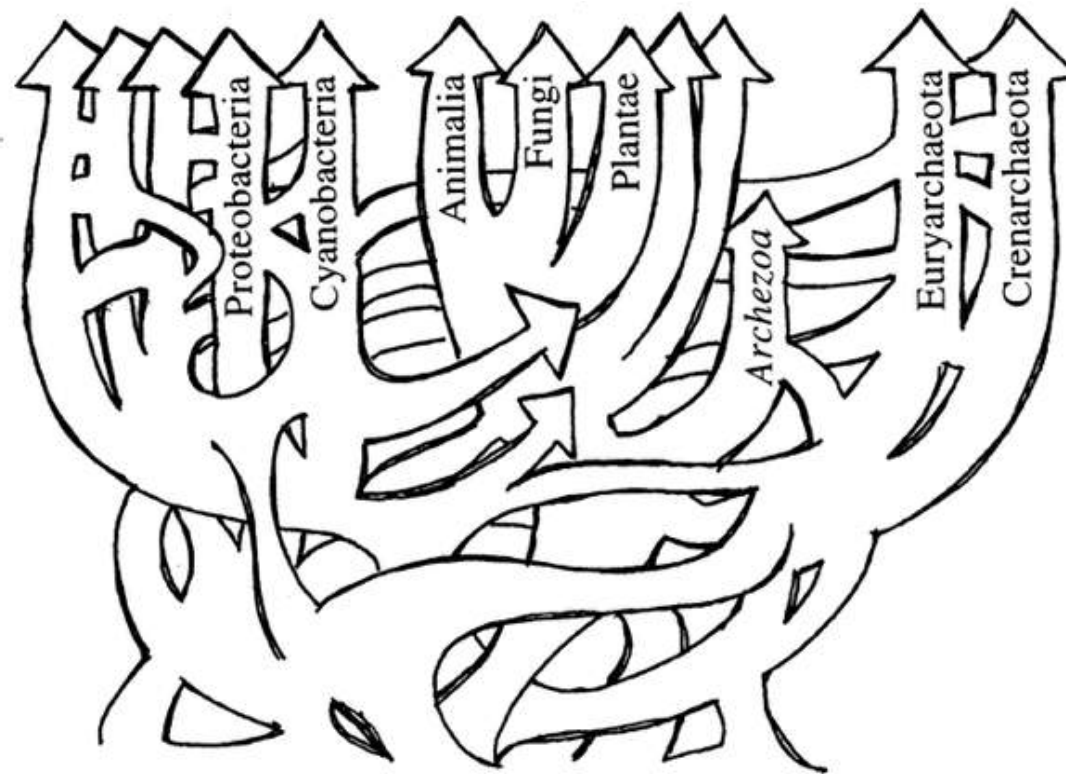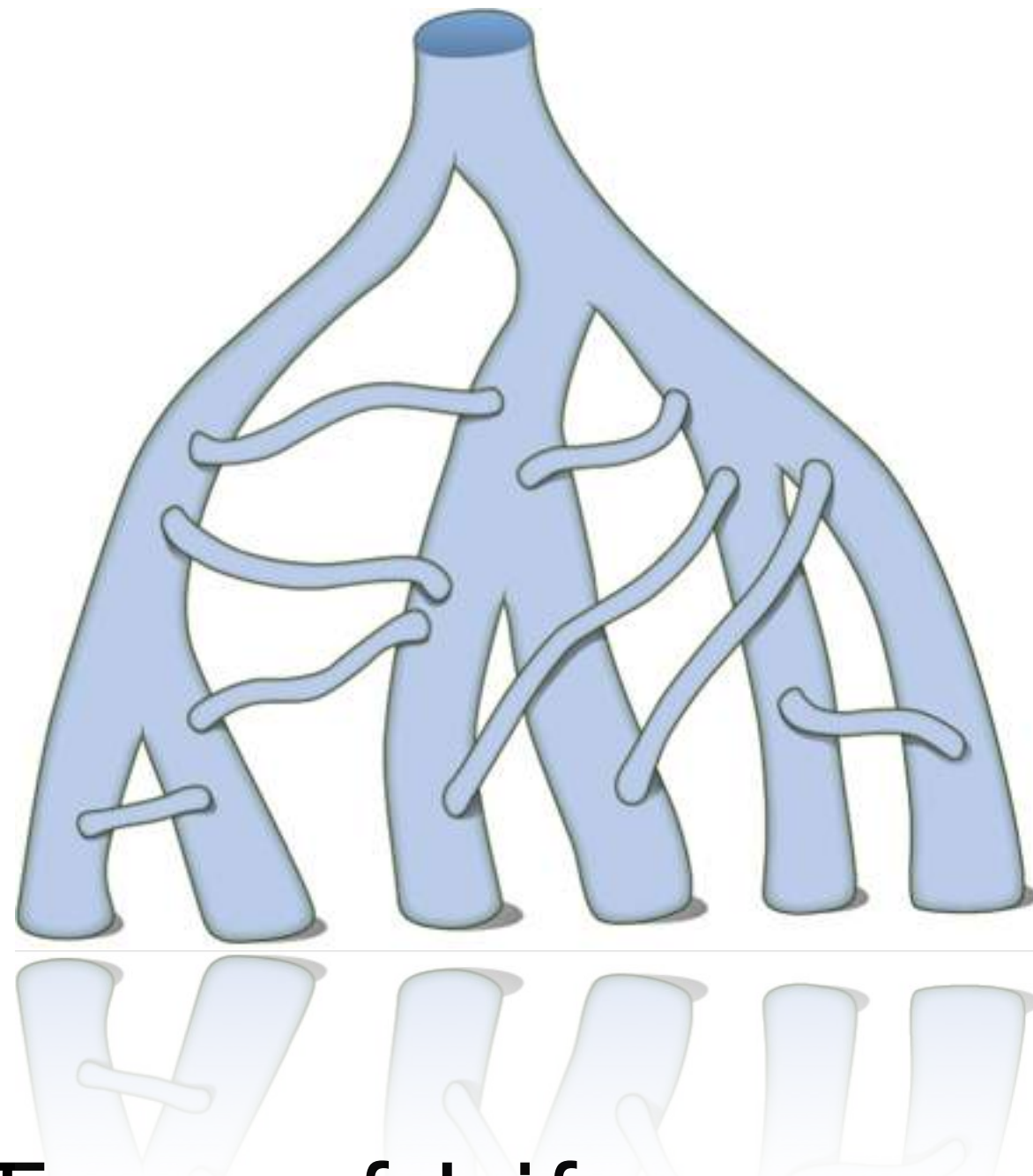
- Goal: recover T from gene trees

# Question:
# How much HGT is too much?

**Theorem (Daskalakis & R. SODA (2016))**

*Under Galtier's model with bounded branch lengths and a molecular clock, reconstructing the species tree from $\Omega(\log n)$ genes is possible as long as the HGT rate is constant.*
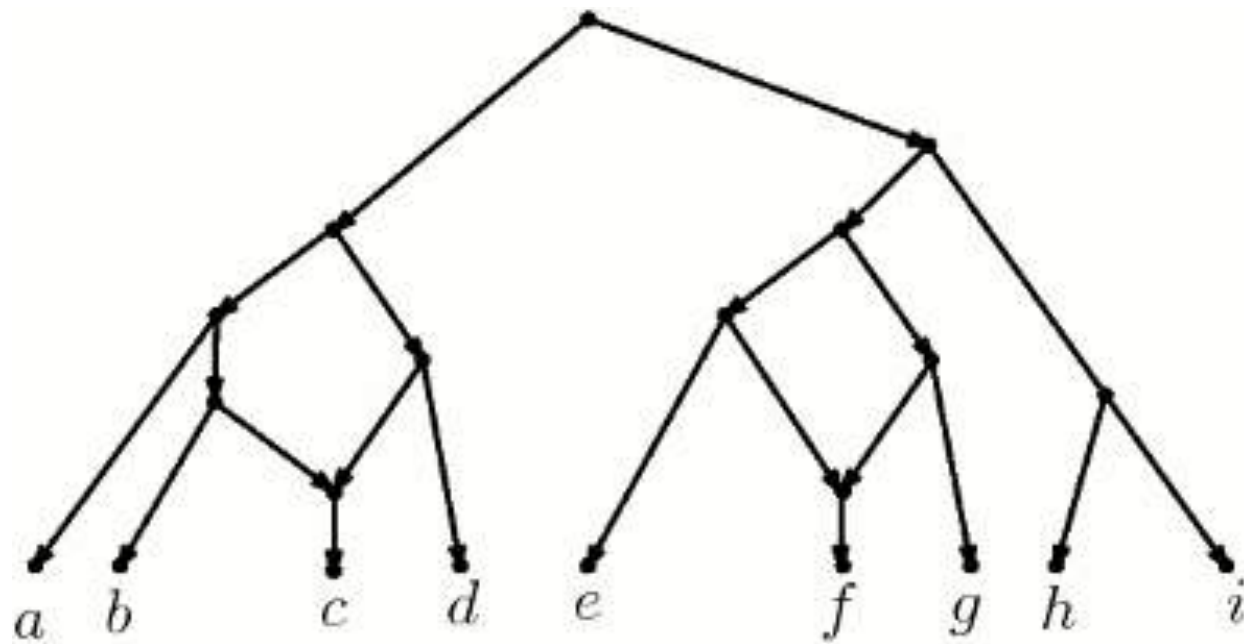
Is the Tree of Life even a tree?
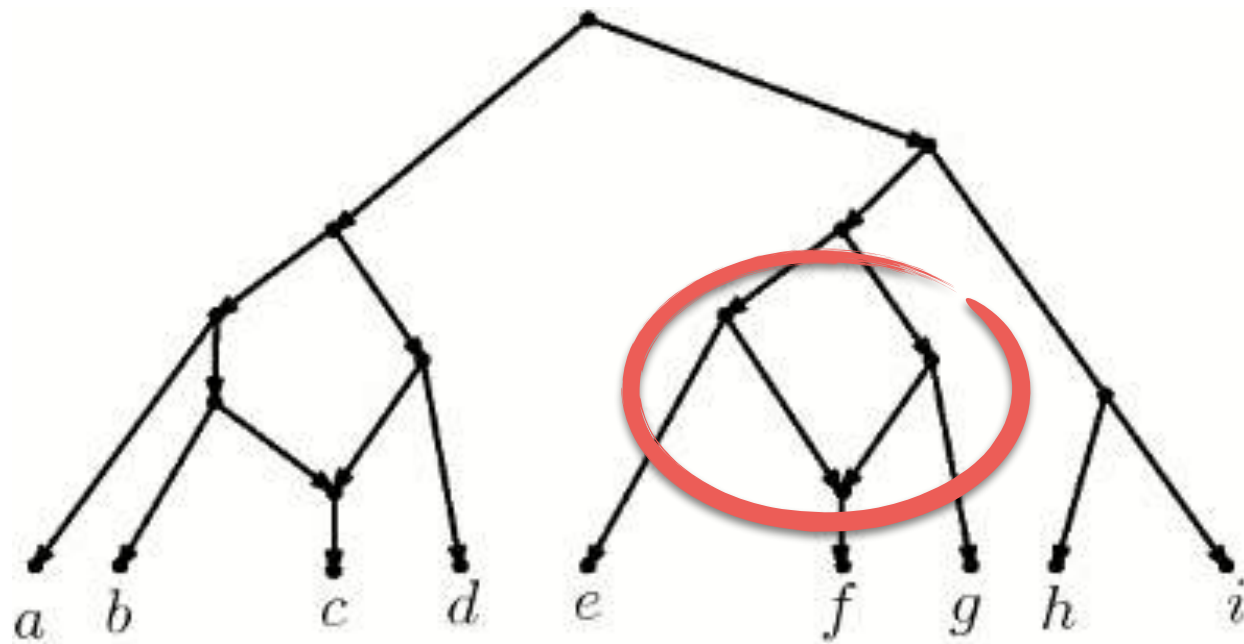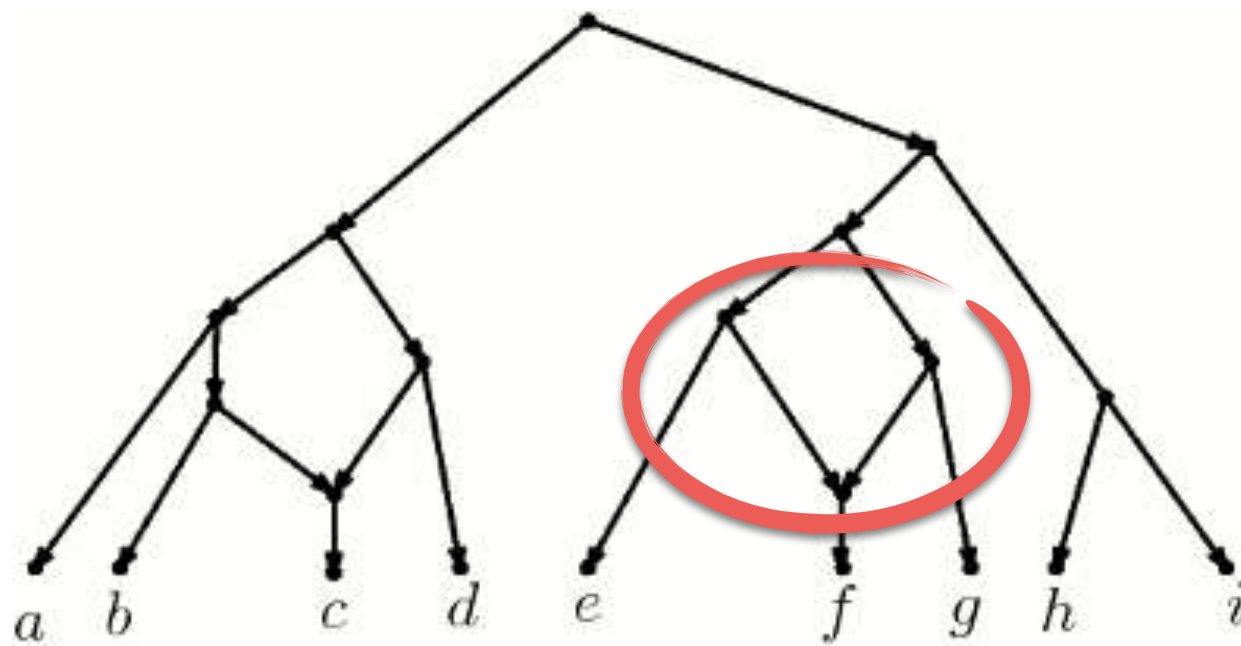
# Hybridization

# Hybridization

# Beyond trees



**Phylogenetic network**

# Beyond trees



**Phylogenetic network**

# Beyond trees



**Phylogenetic network**

**Split network**

For more: http://www.math.wisc.edu/~roch/evol-gen/

# Thanks

For more: http://www.math.wisc.edu/~roch/evol-gen/