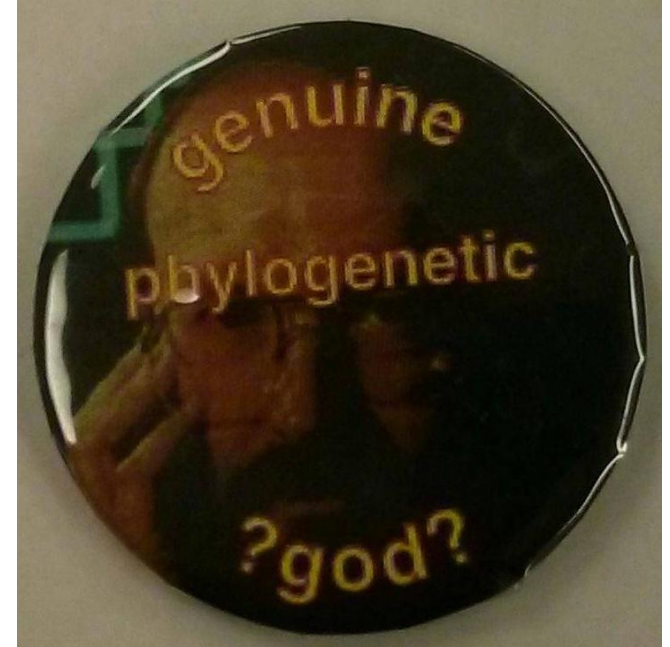


# Conflict in science

MATT SIMON SCIENCE 02.03.16 07:00 AM

## TWITTER NERD-FIGHT REVEALS A LONG, BIZARRE SCIENTIFIC FEUD

#ParsimonyGate



The epistemological paradigm of this journal is parsimony. There are strong philosophical arguments in support of parsimony versus other methods of phylogenetic inference (e.g. Farris, [1983](#)).

The high citation index of *Cladistics* shows that the journal is publishing some of the most ground-breaking empirical and theoretical research on the history of life, and we remain committed to the publication of outstanding systematics research. As a community of scientists, the Willi Hennig Society is always open to new methods and ideas, and to well-reasoned criticisms of old ones. However, we do not hold in special esteem any method solely because it is novel or purportedly sophisticated.

Phylogenetic data sets submitted to this journal should be analysed using parsimony. If alternative methods are also used and there is no difference among the results, the author should defer to the principles of the Society and present the tree obtained by parsimony. Unless there is a pertinent reason to include multiple trees from alternative methods, a tree based on parsimony is sufficient as an intelligible, informative and repeatable hypothesis of relationships, and articles should not be cluttered with multiple, often redundant, trees produced from other methods. If alternative methods give different results and the author prefers an unparsimonious topology, he or she is welcome to present that result, but should be prepared to defend it on philosophical grounds.

In keeping with numerous theoretical and empirical discussions of methodology published in this journal, we do not consider the hypothetical problem of statistical inconsistency to constitute a philosophical argument for the rejection of parsimony. All phylogenetic methods, including parsimony, may produce inconsistent or otherwise inaccurate results for a given data set. The absence of certain truth represents a philosophical limit of empirical science.

*Cladistics* will publish research based on methods that are repeatable, clearly articulated and philosophically sound. We believe these guidelines implement the vision of Willi Hennig ([1965](#), p. 97), who said, "(i)nvestigation of the phylogenetic relationship between all existing species and the expression of the results of this research in a form which cannot be misunderstood, is the task of phylogenetic systematics."

# Parsimony still commonly used for:

Morphological analyses (though probabilistic models exist here too)

Starting trees

Understanding history of trait change

# Practical considerations

How do you find the most parsimonious tree?

# Heuristic searches

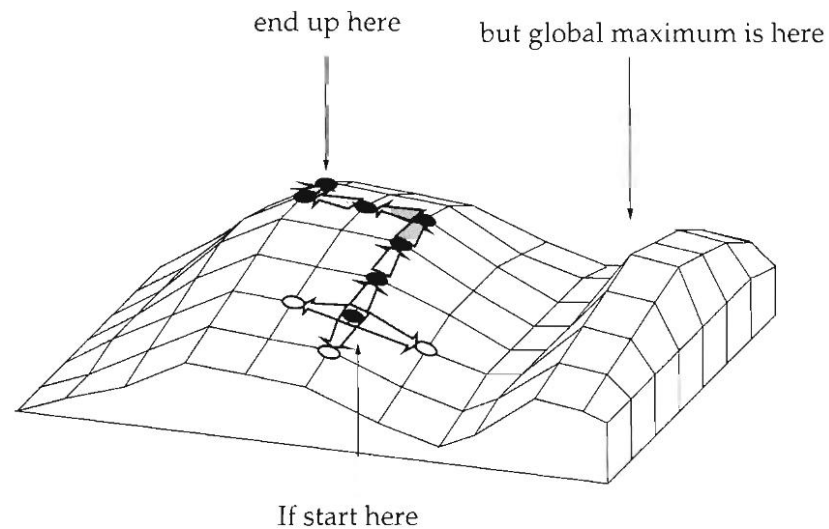


Figure 4.1: A surface rising above a two-dimensional plain (or plane). The process of climbing uphill on the surface is illustrated, as well as the failure to find a higher peak by this "greedy" method.

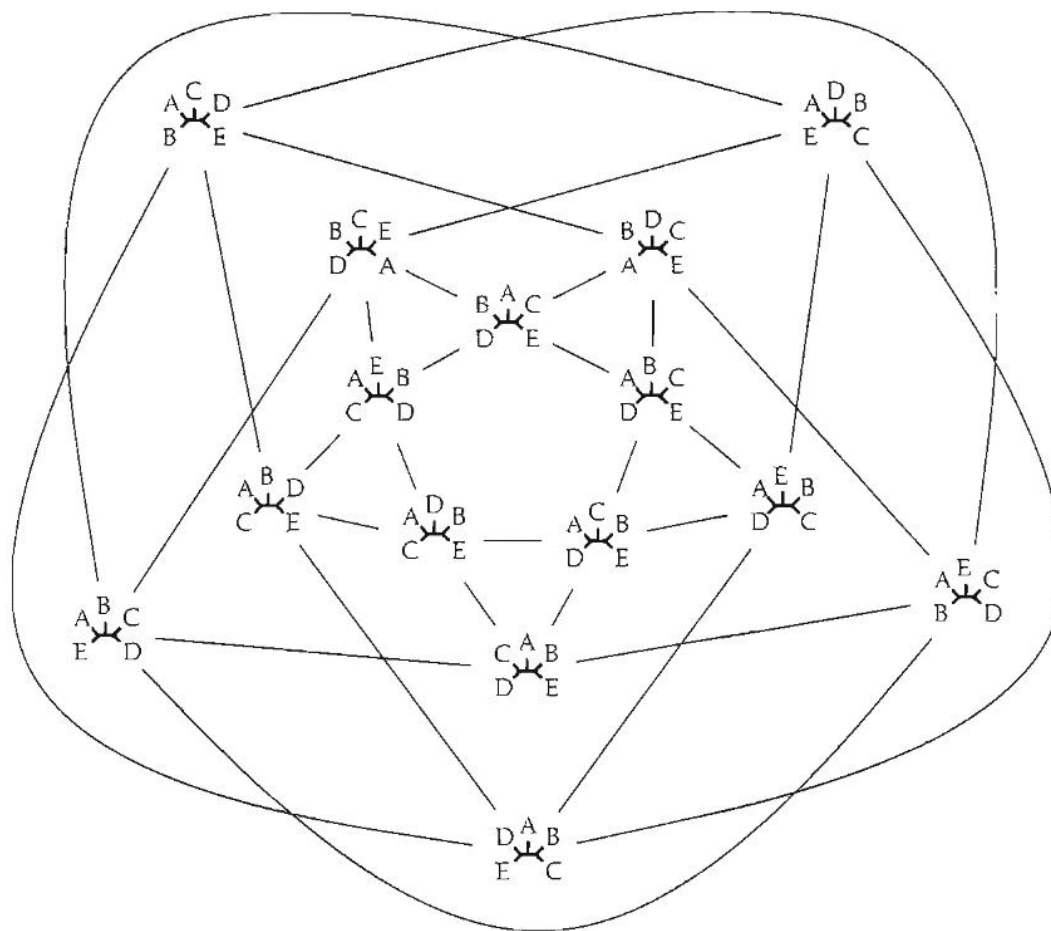


Figure 4.3: The space of all 15 possible unrooted trees with 5 tips. Neighbors are connected by lines when a nearest-neighbor interchange can convert one into the other. The labels A–E correspond to the species names Alpha through Epsilon in that data set. This symmetric arrangement of nodes was discovered by Ben Rudd Schoenberg (personal communication), and we thus denote this graph the Schoenberg graph.

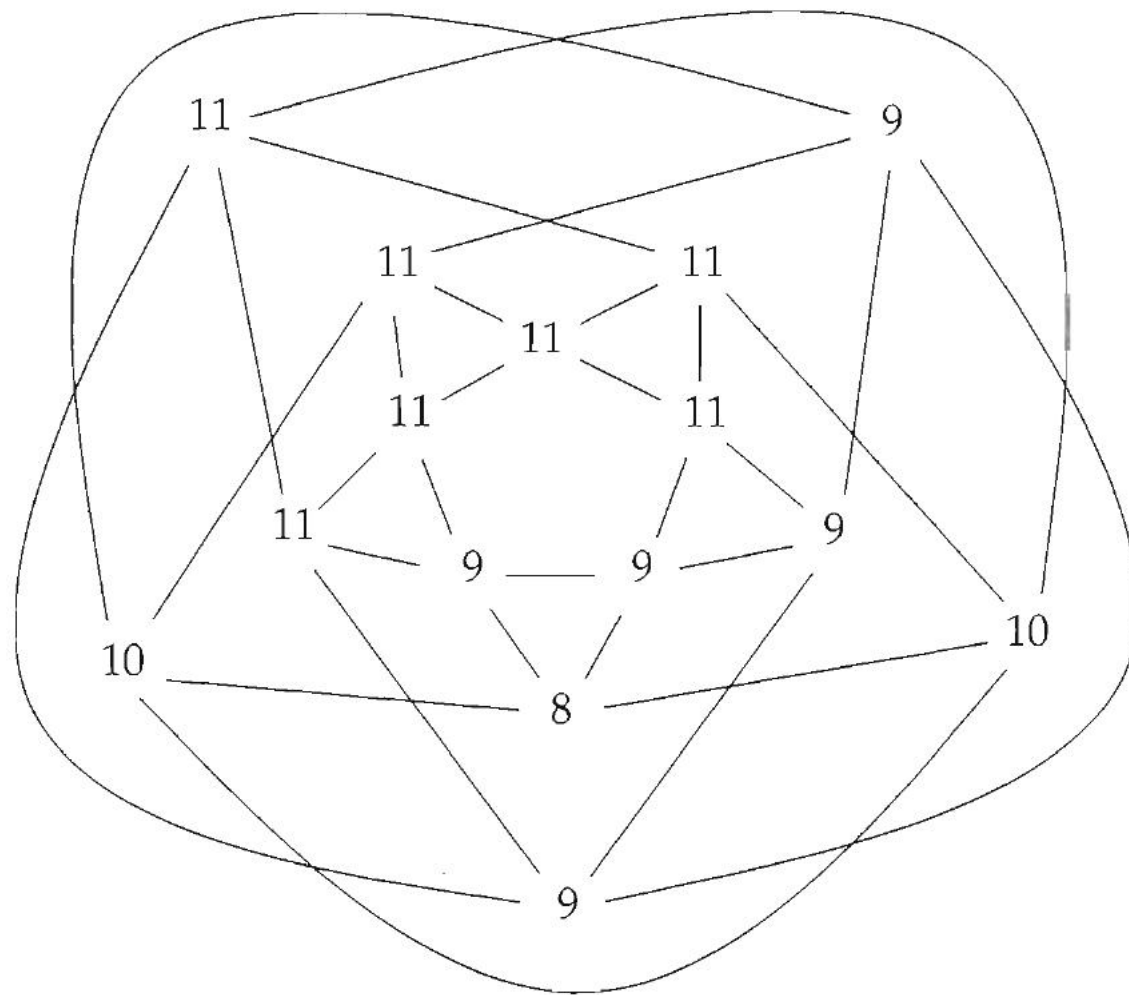
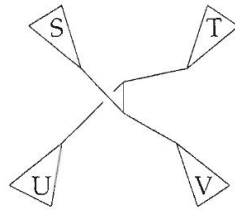
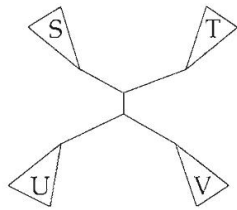


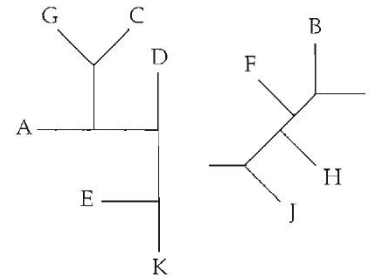
Figure 4.4: The space of all 15 possible trees, as in Figure 4.3, where the number of changes of state on the data set of Table 1.1 is shown. Nearest-neighbor interchanges search for the most parsimonious tree by moving in this graph.

# Tree space "moves"

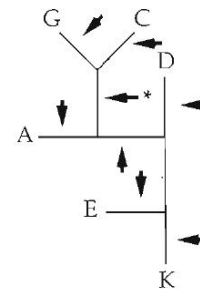


Nearest Neighbor Interchange  
(NNI)

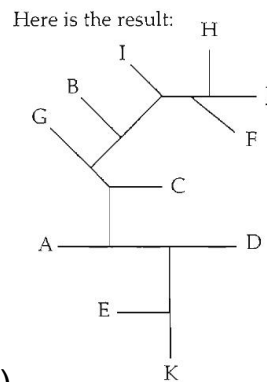
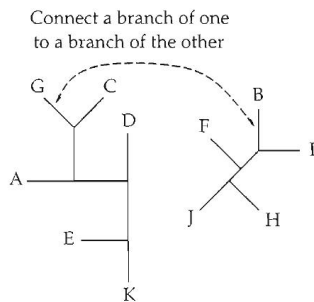
Break a branch, remove a subtree



Add it in, attaching it to one (\*)  
of the other branches

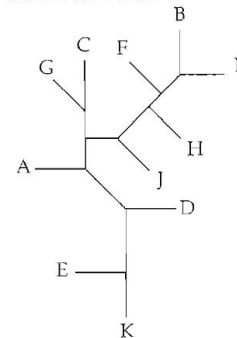


Subtree pruning and regrafting  
(SPR)



Tree bisection and  
reconnection (TBR)

Here is the result:





## Branch and bound

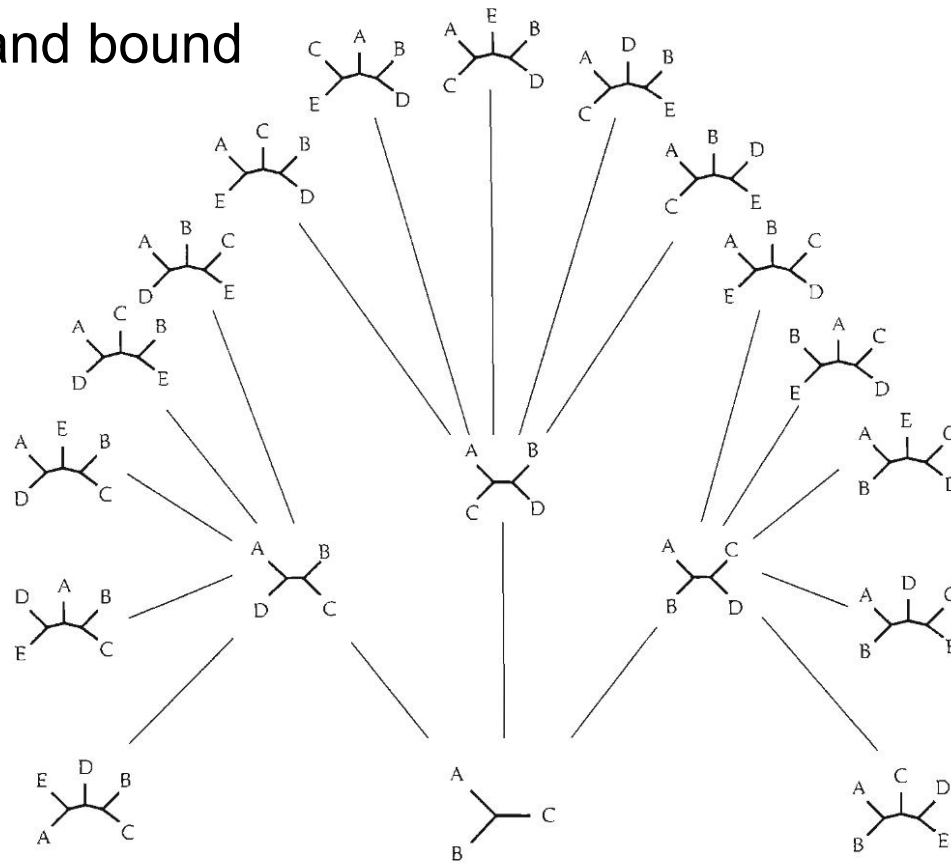


Figure 5.3: Search tree for most parsimonious tree in a five-species case.

## Branch and bound

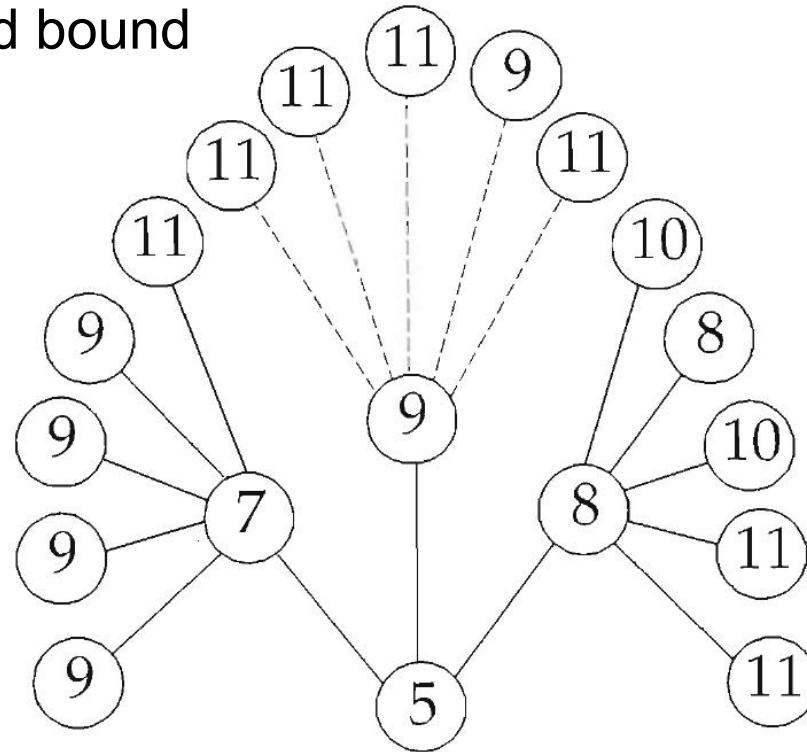


Figure 5.4: Search tree for most parsimonious tree for five species, using the data of Table 1.1. Trees are shown in Figure 5.3. Dashed lines are those not traversed by a branch and bound method. The species names in the data set correspond to labels A through E in Figure 5.3.

Software: PAUP\*, TNT, Mesquite,  
others...

# Felsenstein & the birth of statistical phylogenetics



## Joe Felsenstein

Professor of Genome Sciences, and Professor of Biology, [University of Washington, Seattle](#)

Verified email at gs.washington.edu - [Homepage](#)

[Evolutionary biology](#) [phylogenetic methods](#) [population genetics](#)

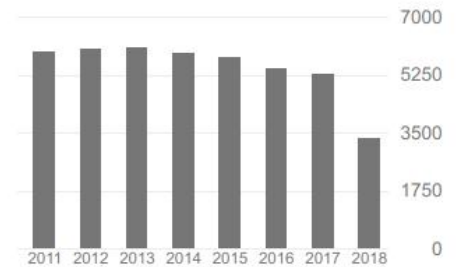
FOLLOW

TITLE	CITED BY	YEAR
<a href="#">Confidence limits on phylogenies: an approach using the bootstrap</a> J Felsenstein Evolution 39 (4), 783-791	36415	1985
<a href="#">PHYLIP (phylogeny inference package), version 3.5 c</a> J Felsenstein Joseph Felsenstein.	27086 *	1993
<a href="#">Evolutionary trees from DNA sequences: a maximum likelihood approach</a> J Felsenstein Journal of molecular evolution 17 (6), 368-376	10773	1981
<a href="#">Phylogenies and the comparative method</a> J Felsenstein The American Naturalist 125 (1), 1-15	7697	1985
<a href="#">Inferring phylogenies</a> J Felsenstein, J Felsenstein Sinauer associates	4570	2004
<a href="#">Cases in which parsimony or compatibility methods will be positively misleading</a> J Felsenstein Systematic zoology 27 (4), 401-410	3408	1978
<a href="#">Phylogenies from molecular sequences: inference and reliability</a> J Felsenstein Annual review of genetics 22 (1), 521-565	2471	1988
<a href="#">Maximum likelihood estimation of a migration matrix and effective population sizes in n subpopulations by using a coalescent approach</a> P Beerli, J Felsenstein Proceedings of the National Academy of Sciences 98 (8), 4563-4568	1583	2001
<a href="#">The evolutionary advantage of recombination</a> J Felsenstein Genetics 78 (2), 737-756	1312	1974

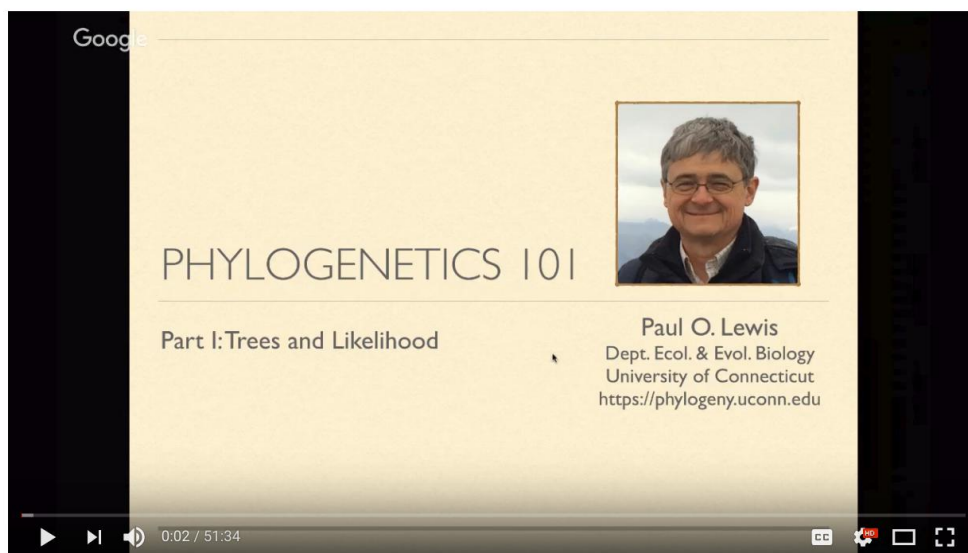
Cited by

[VIEW ALL](#)

	All	Since 2013
Citations	118614	31934
h-index	75	44
i10-index	148	81



Full disclaimer: I'm borrowing some of this material from Paul Lewis (Uconn)  
(Check out his teaching materials!)



Phyloseminar #76: Paul Lewis (UConn) Primer part 1

701 views

LIKE DISLIKE SHARE



phyloseminar.org

Streamed live on Apr 18, 2018

SUBSCRIBED 973

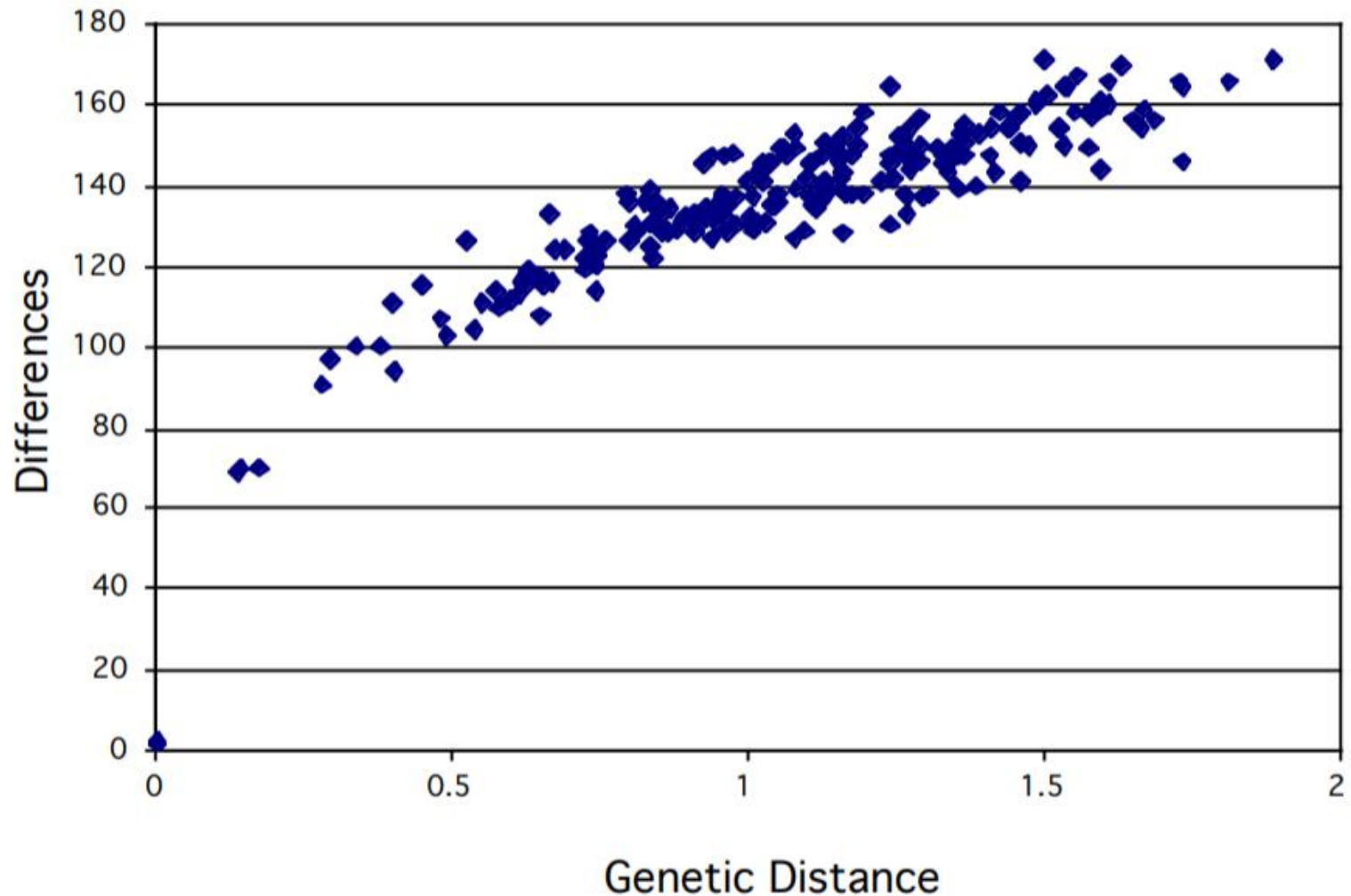


Primer part 1: tree terminology and substitution models

Slides: <https://git.io/vplW9>

SHOW MORE

# Why do we need statistics?



If two sequences are unrelated, what % of bases (aligned sites) do you expect to be identical?

A. 50%

B. 25%

C. 0%

D. I need more information



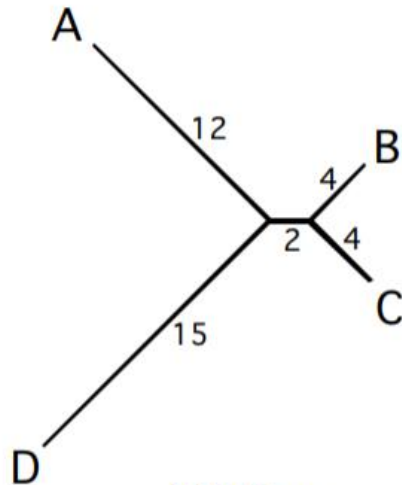
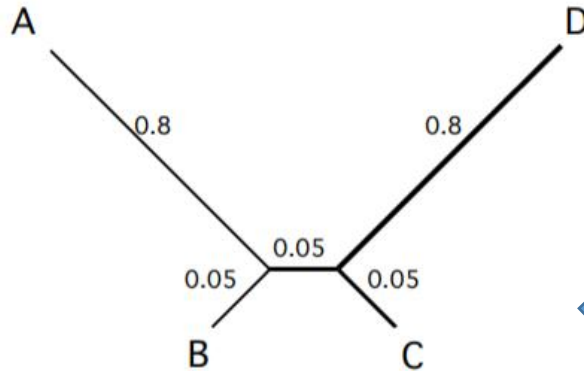
# "Long-branch attraction"

## CASES IN WHICH PARSIMONY OR COMPATIBILITY METHODS WILL BE POSITIVELY MISLEADING<sup>1</sup>

JOSEPH FELSENSTEIN

### Abstract

Felsenstein, J. (Department of Genetics, University of Washington, Seattle, WA 98195) 1978. Cases in which parsimony or compatibility methods will be positively misleading. *Syst. Zool.* 27:401-410.—For some simple three- and four-species cases involving a character with two states, it is determined under what conditions several methods of phylogenetic inference will fail to converge to the true phylogeny as more and more data are accumulated. The methods are the Camin-Sokal parsimony method, the compatibility method, and Farris's unrooted Wagner tree parsimony method. In all cases the conditions for this failure (which is the failure to be statistically consistent) are essentially that parallel changes exceed informative, nonparallel changes. It is possible for these methods to be inconsistent even when change is improbable a priori, provided that evolutionary rates in different lineages are sufficiently unequal. It is by extension of this approach that we may provide a sound methodology for evaluating methods of phylogenetic inference. [Numerical cladistics; phylogenetic inference; maximum likelihood estimation; parsimony; compatibility.]



MP Tree

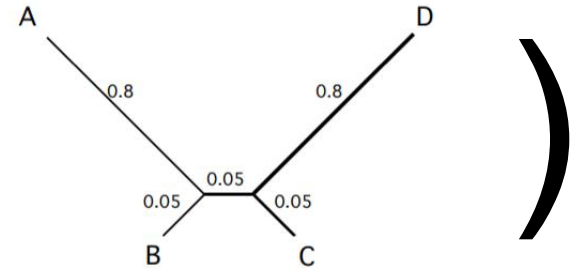
A  
B  
C  
D

ATCGAGCAGCCTGGGAGAGAGACTTATTTGACAAACGTAA  
ATTGGGGAGTAGCGTAAACACTCTTATTTGACGAAATTAT  
ATCGTGGGTTAGAGTAGAGACTCTCATTGACGAAATTAT  
AACGTGGCGAATAGTAGTCAAAAAATGTGTACCAGATTAC

P(

ATCGAGCAGCCTGGGAGAGAGACTTATTTGACAAACGTAA  
ATTGGGGAGTAGCGTAAACACTCTTATTTGACGAAATTAT  
ATCGTGGGTTAGAGTAGAGACTCTCATTTGACGAAATTAT  
AACGTGGCGAATAGTAGTCAAAAAATGTGTACCAGATTAC

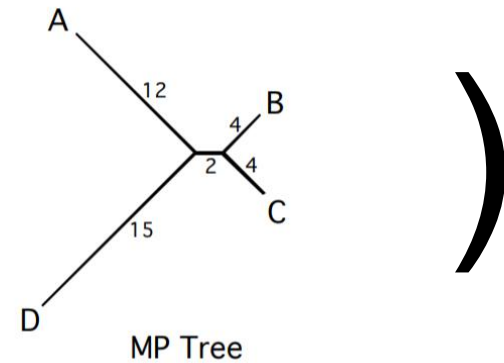
|



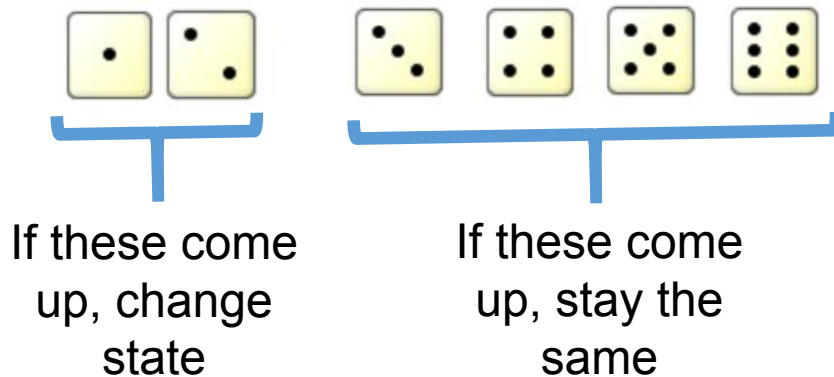
P(

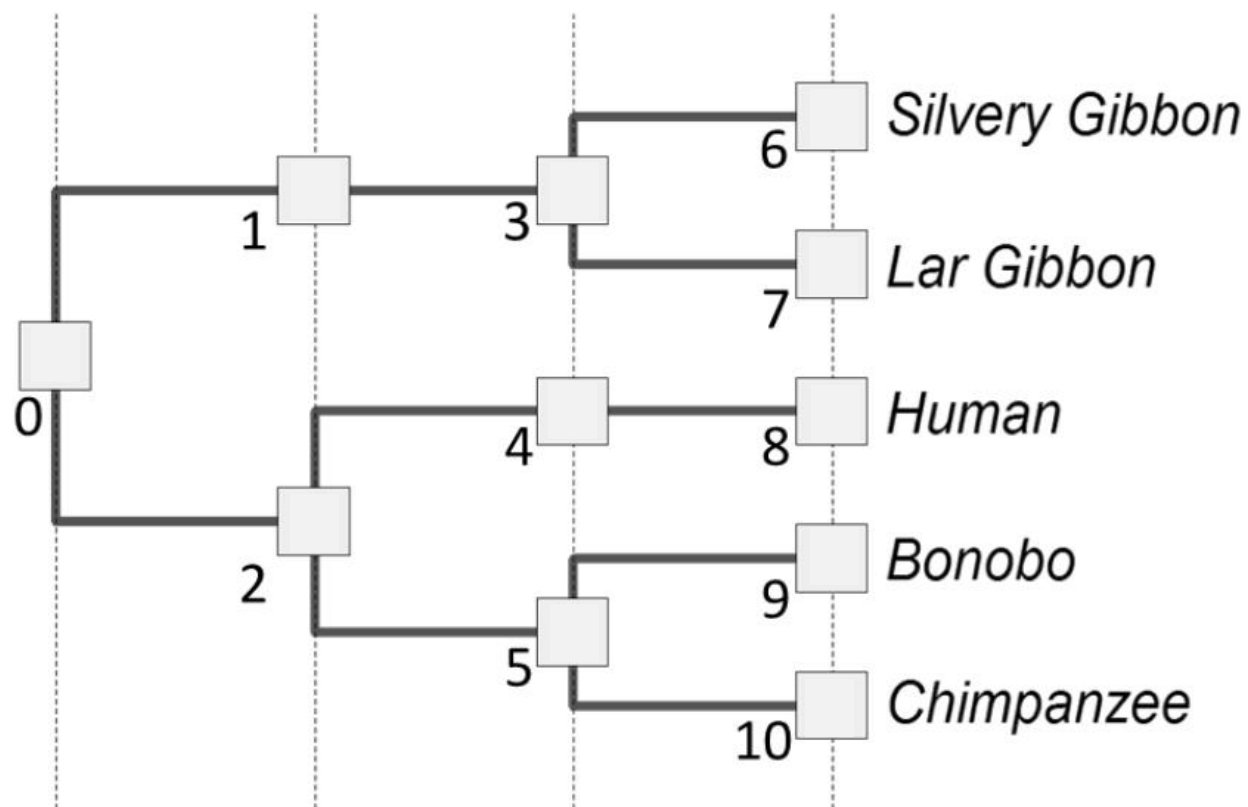
ATCGAGCAGCCTGGGAGAGAGACTTATTTGACAAACGTAA  
ATTGGGGAGTAGCGTAAACACTCTTATTTGACGAAATTAT  
ATCGTGGGTTAGAGTAGAGACTCTCATTTGACGAAATTAT  
AACGTGGCGAATAGTAGTCAAAAAATGTGTACCAGATTAC

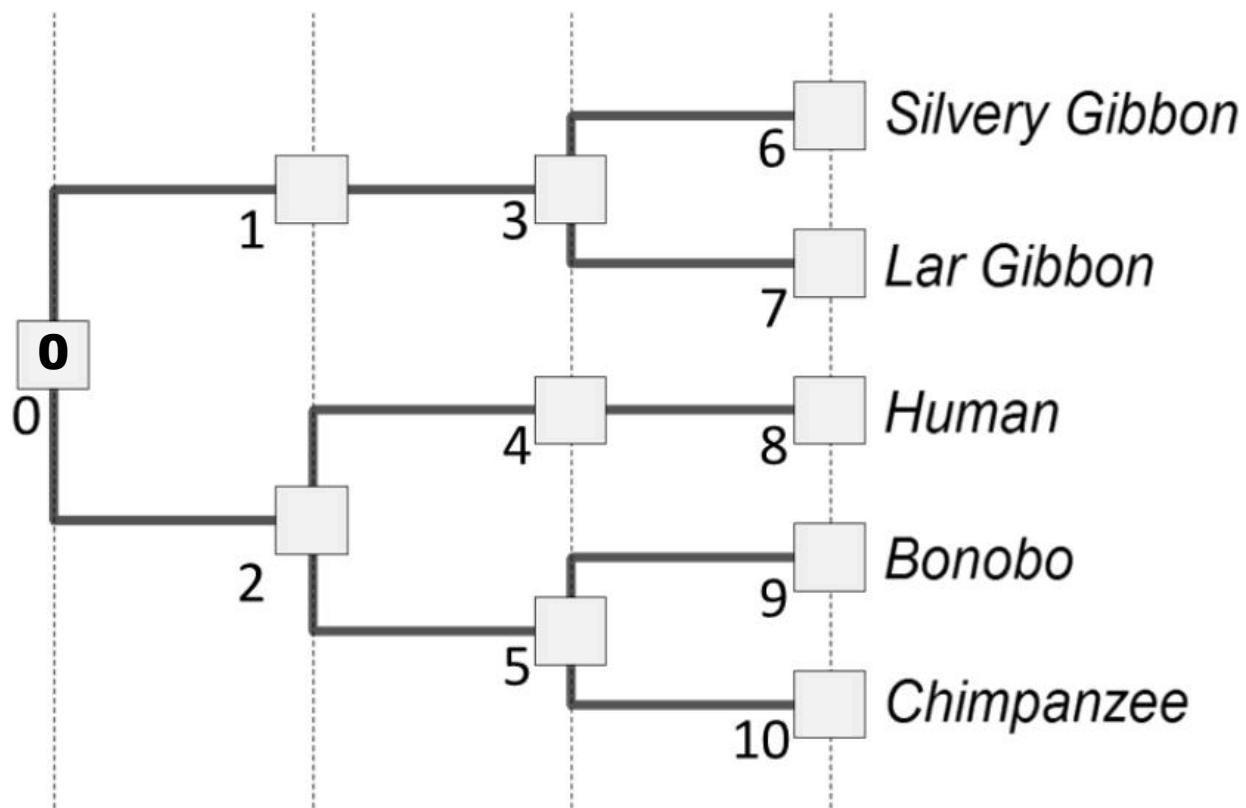
|

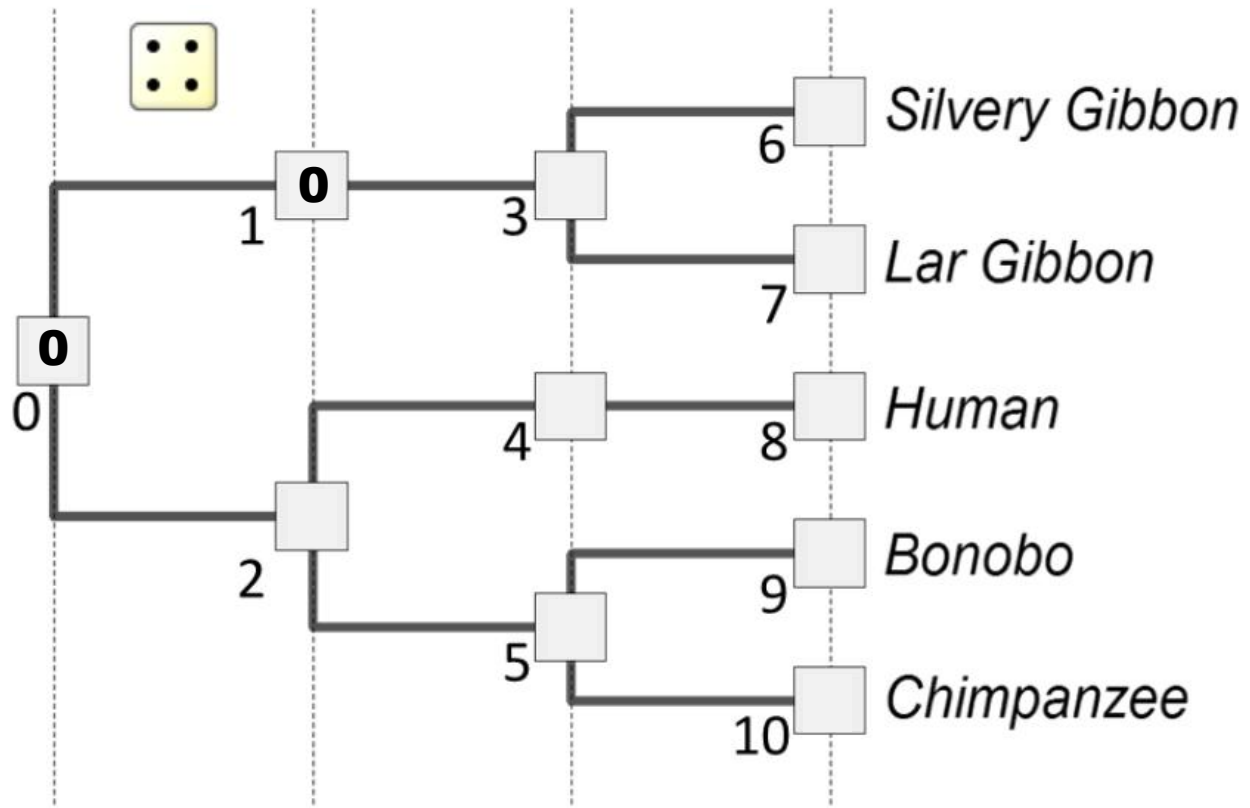


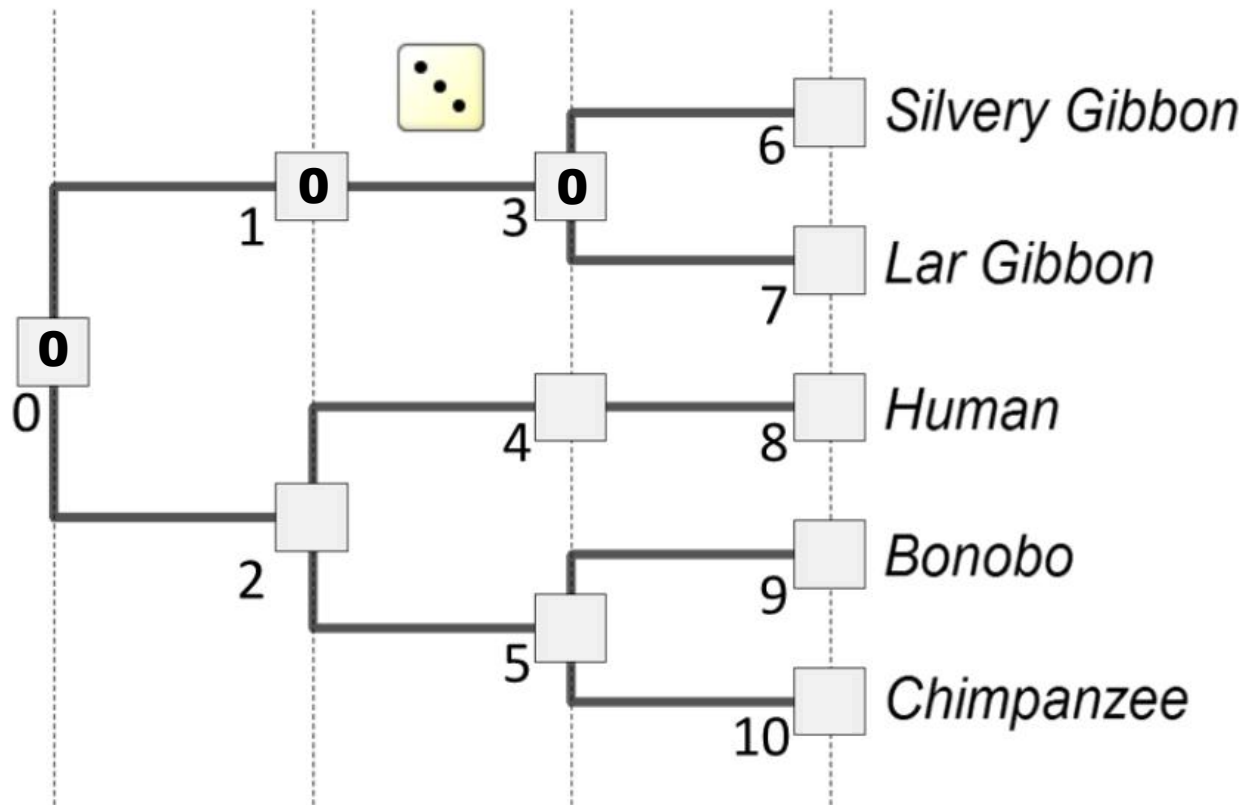
# How do we treat trees probabilistically?

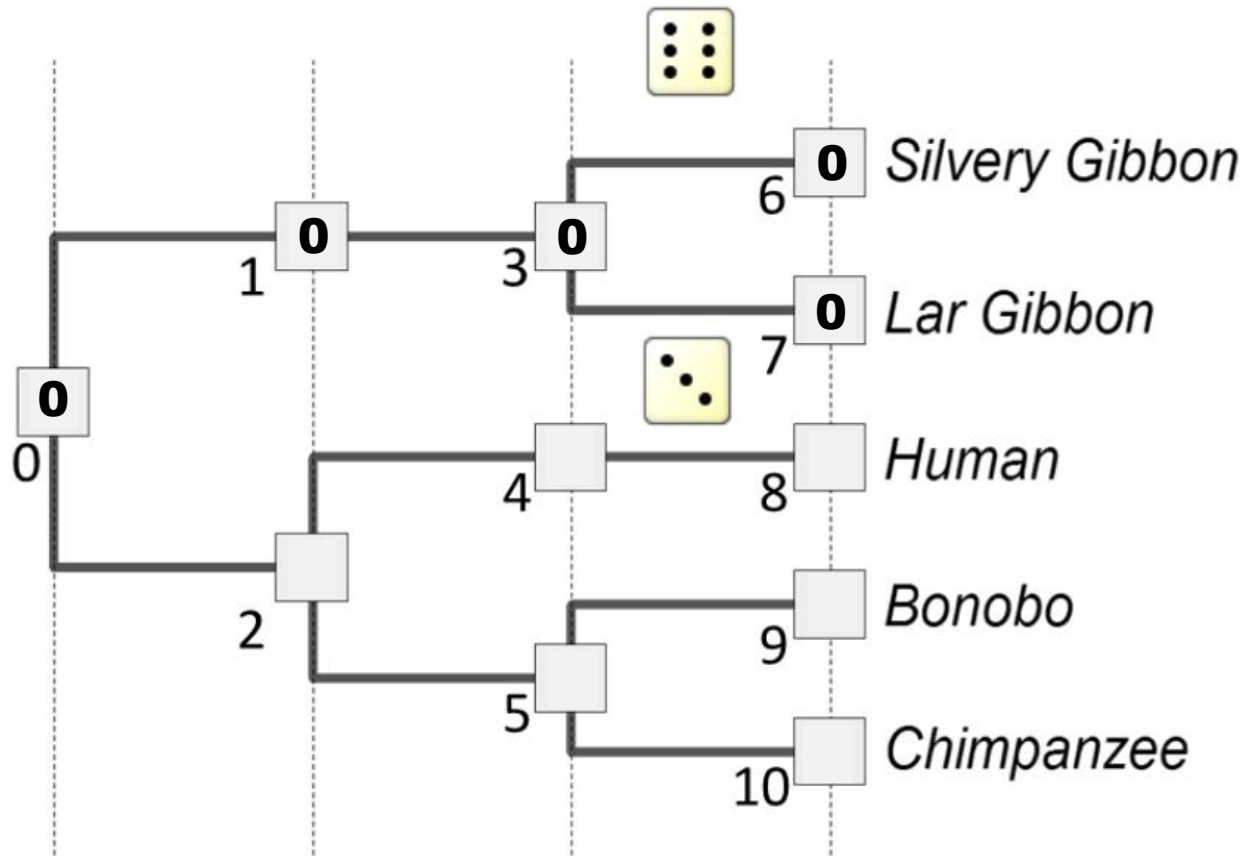




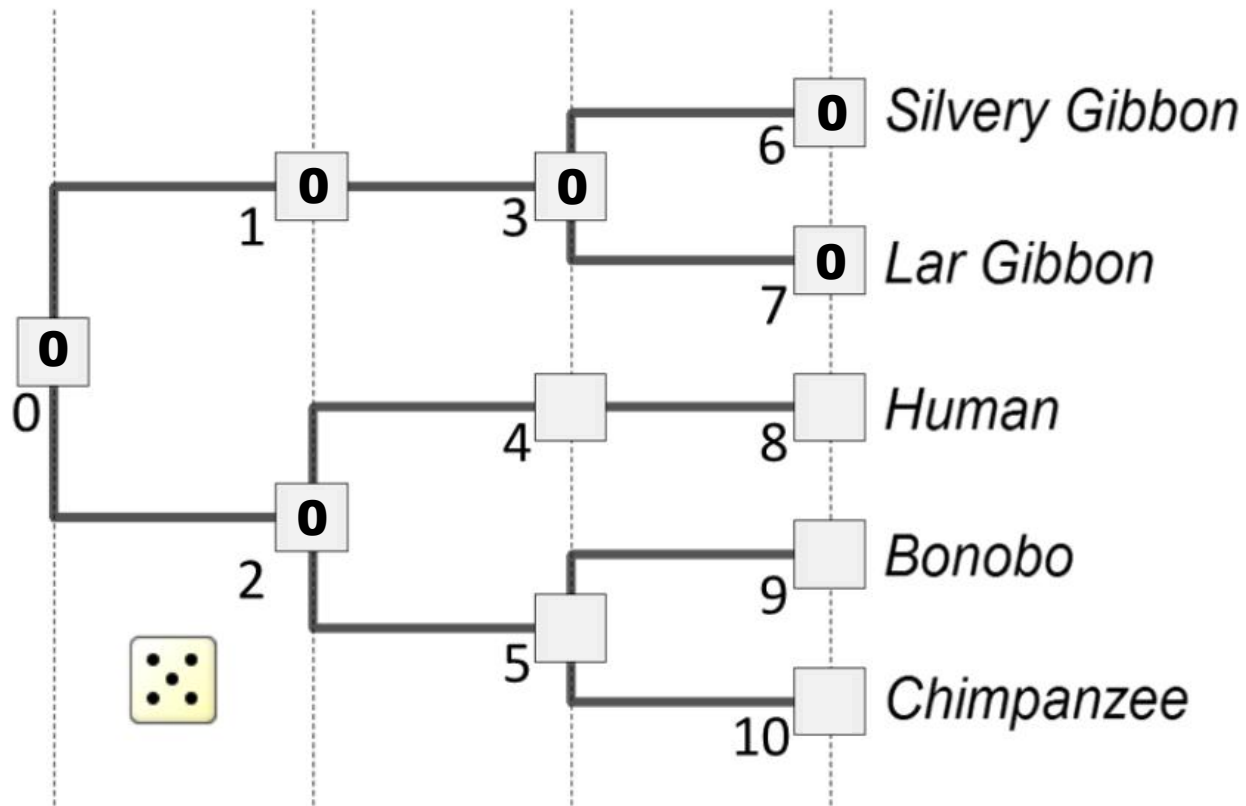


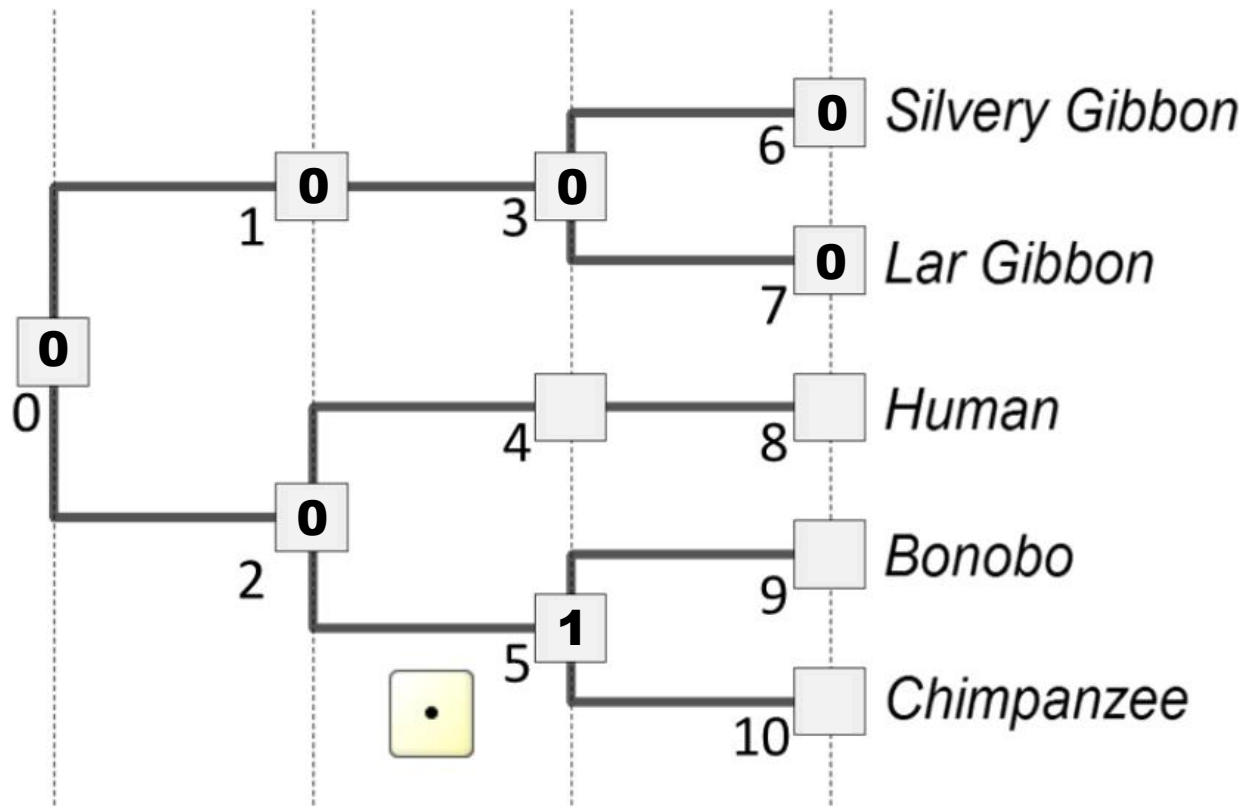


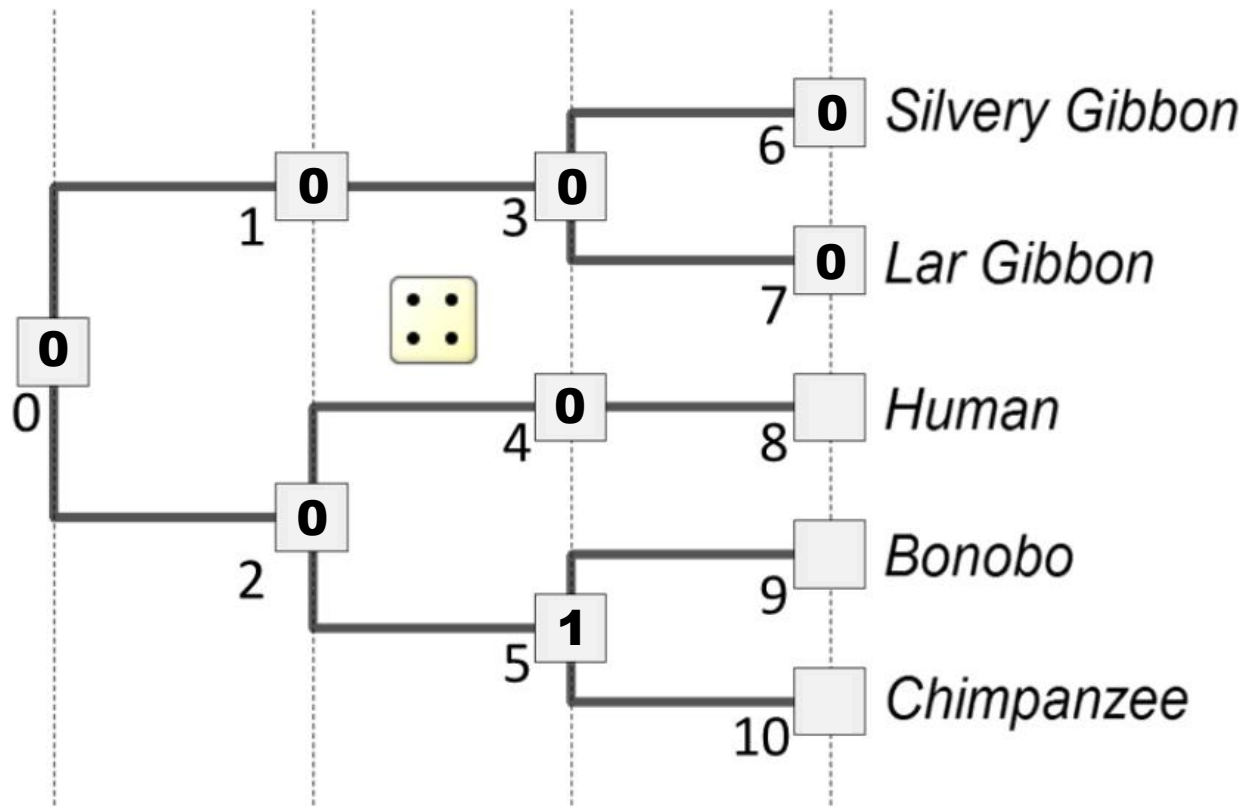


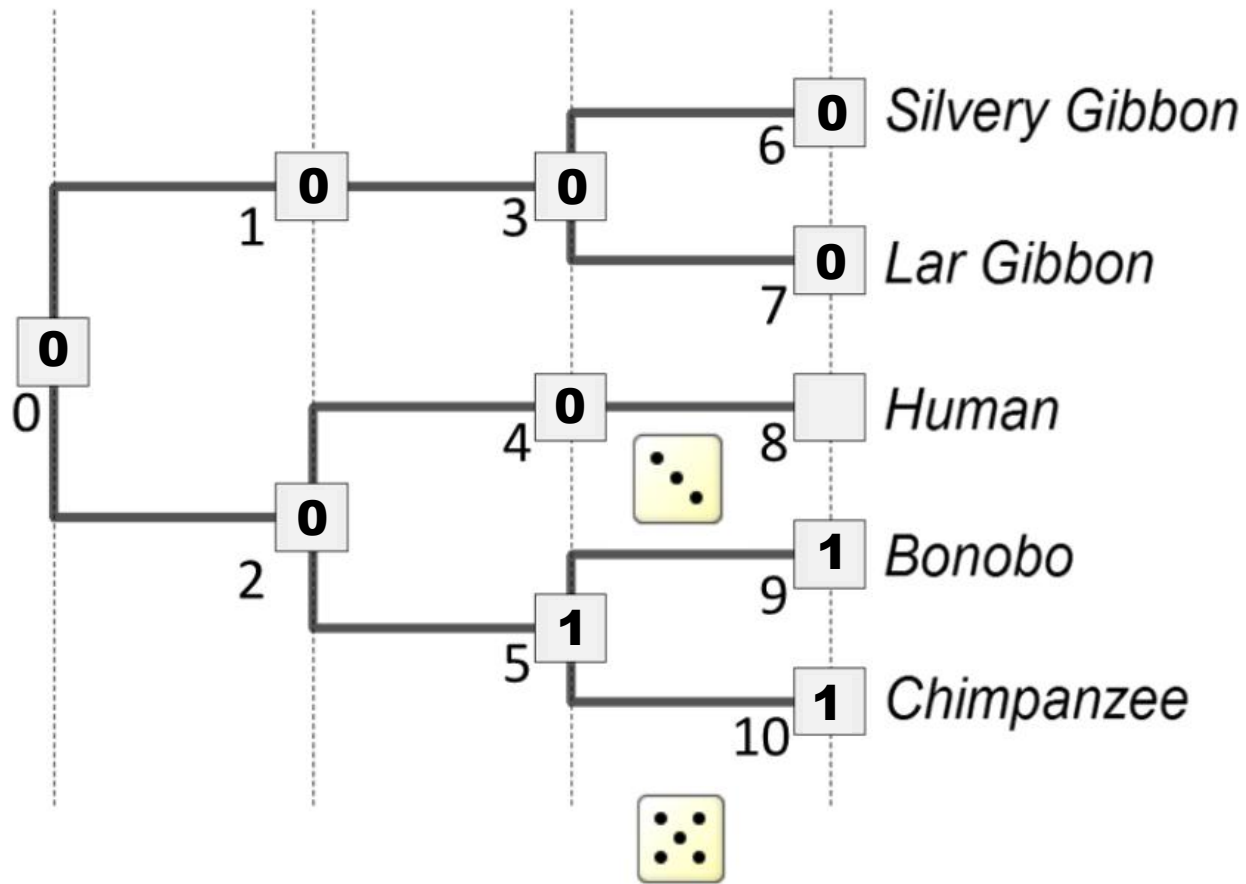


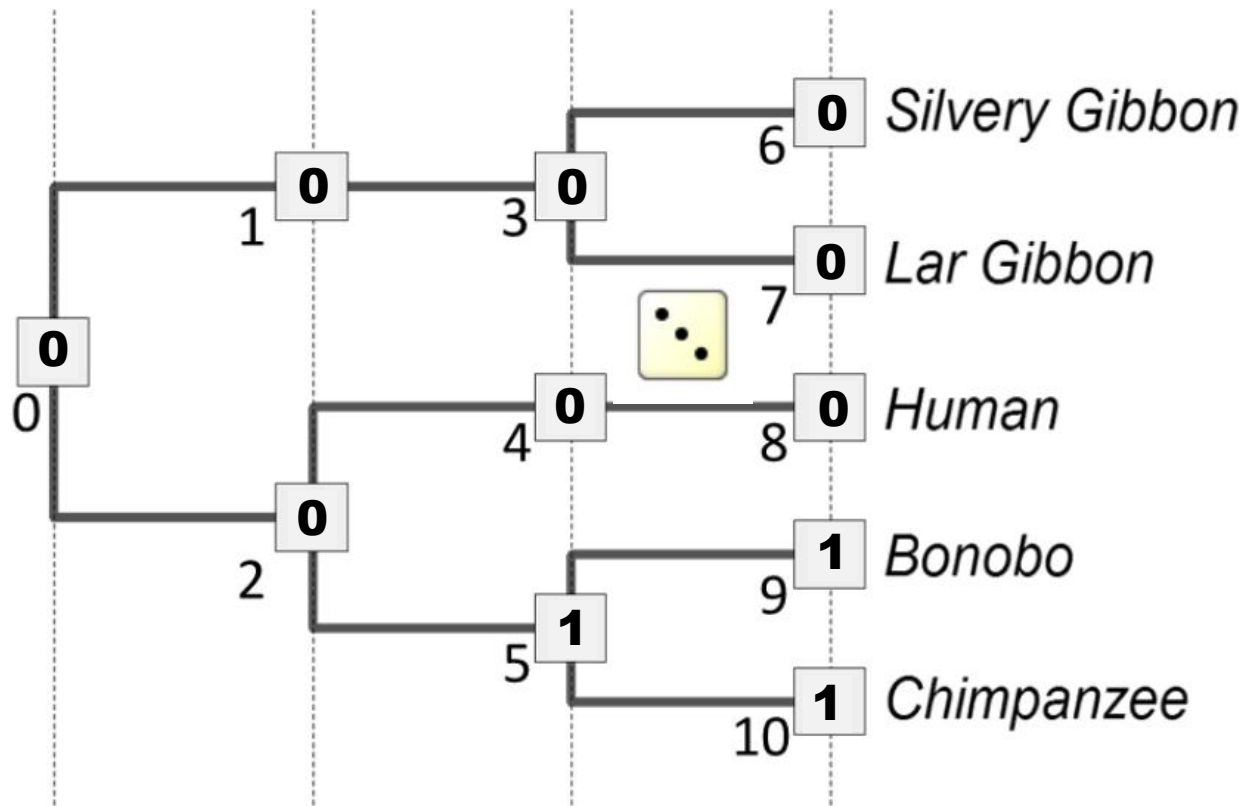












# Combining probabilities: The AND rule

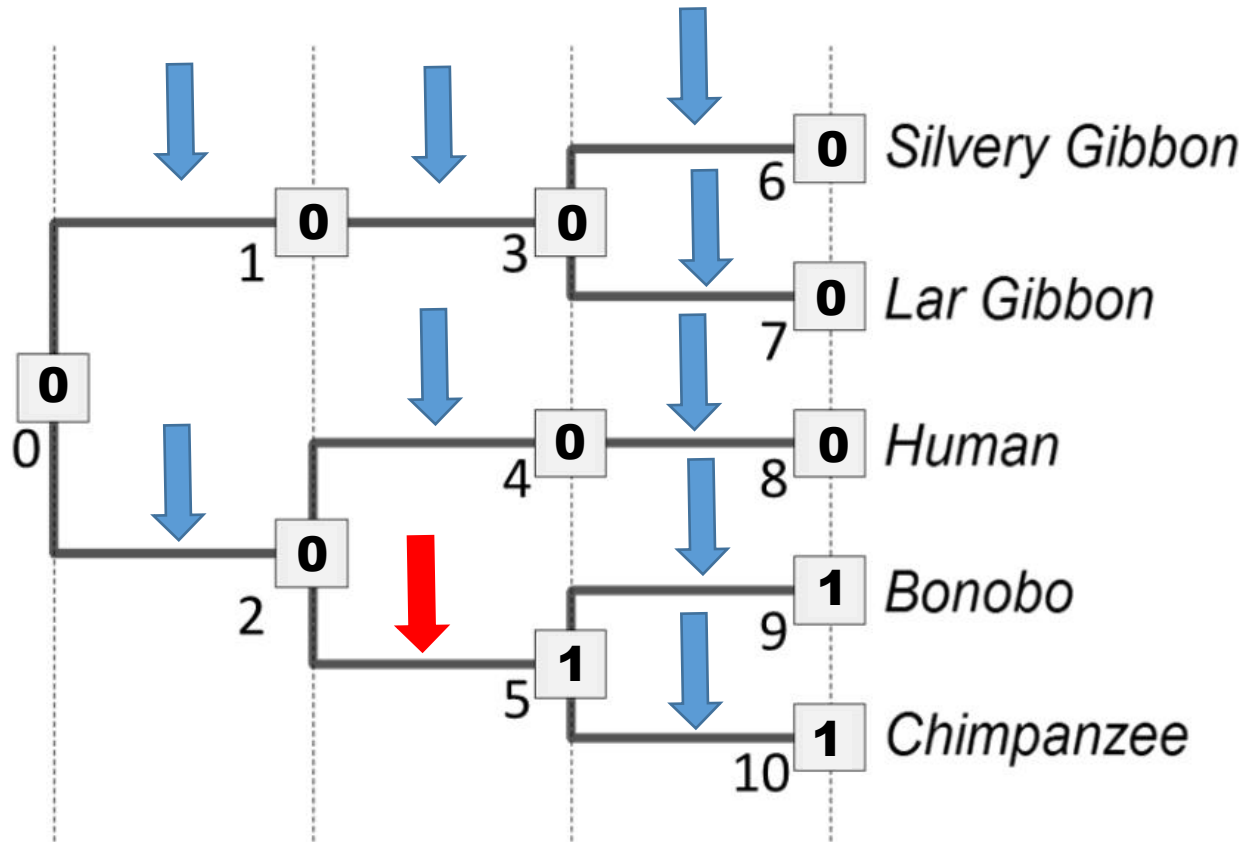
If two independent events occur, multiply their individual probabilities to get the full probability of an event

Using 2 dice, what is the probability of



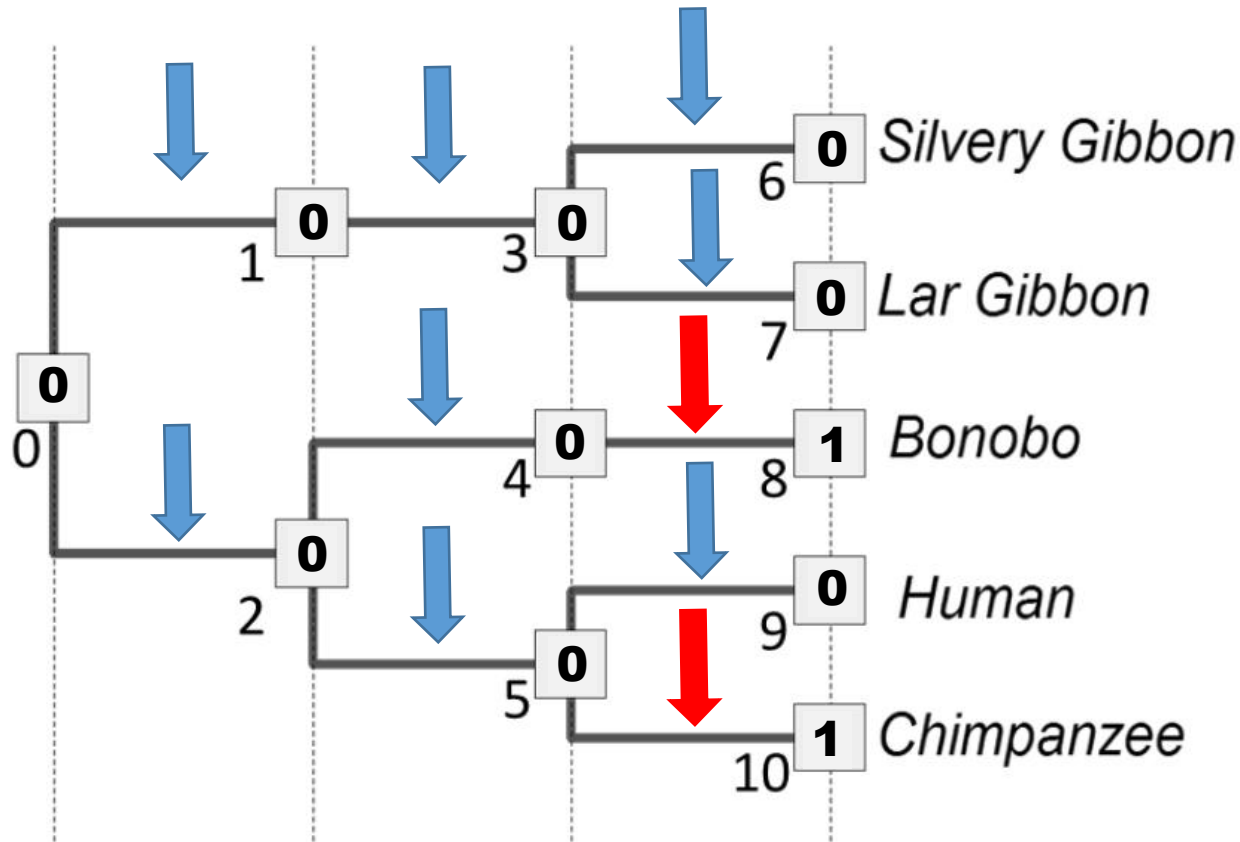
$$(1/6) \times (1/6) = 1/36$$

# Combining probabilities: The AND rule in phylogenetics



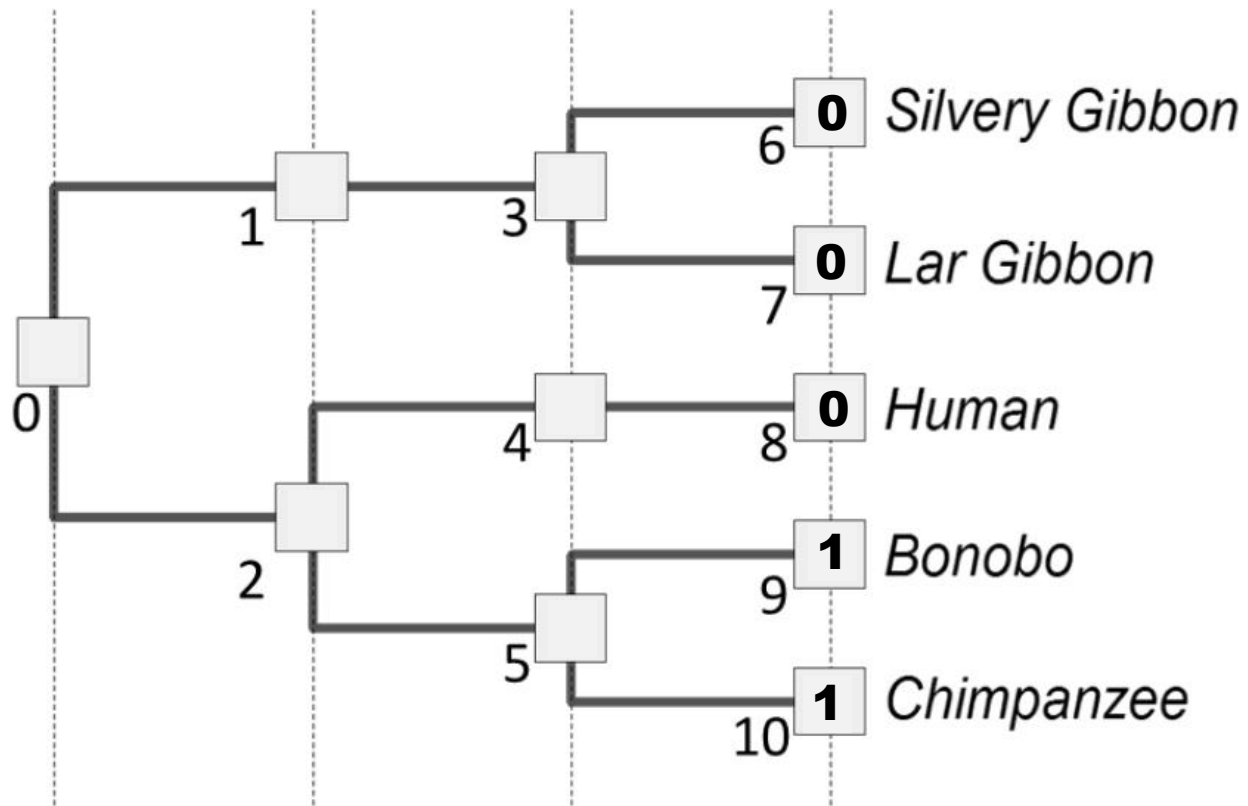
$$\text{Likelihood}(D, N_i | T_1)^* = (4/6)^9 \times (2/6)^1 = 0.0087$$

# Combining probabilities: The AND rule in phylogenetics



$$\text{Likelihood}(D, N_j|T_2)^* = (4/6)^8 \times (2/6)^2 = 0.0043$$





# Combining probabilities: The OR rule

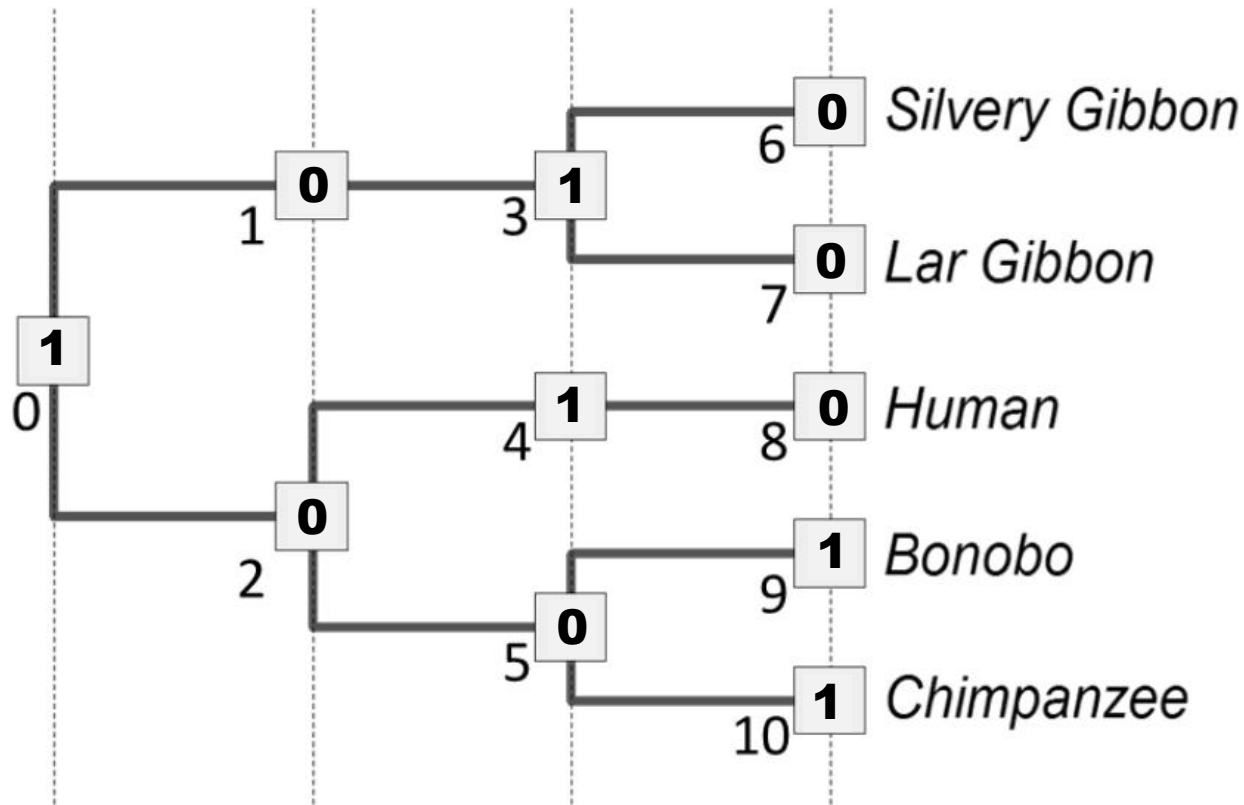
- Two mutually exclusive probabilities should be ADDED together to get the total probability of the two events

Using one die, what is the probability of



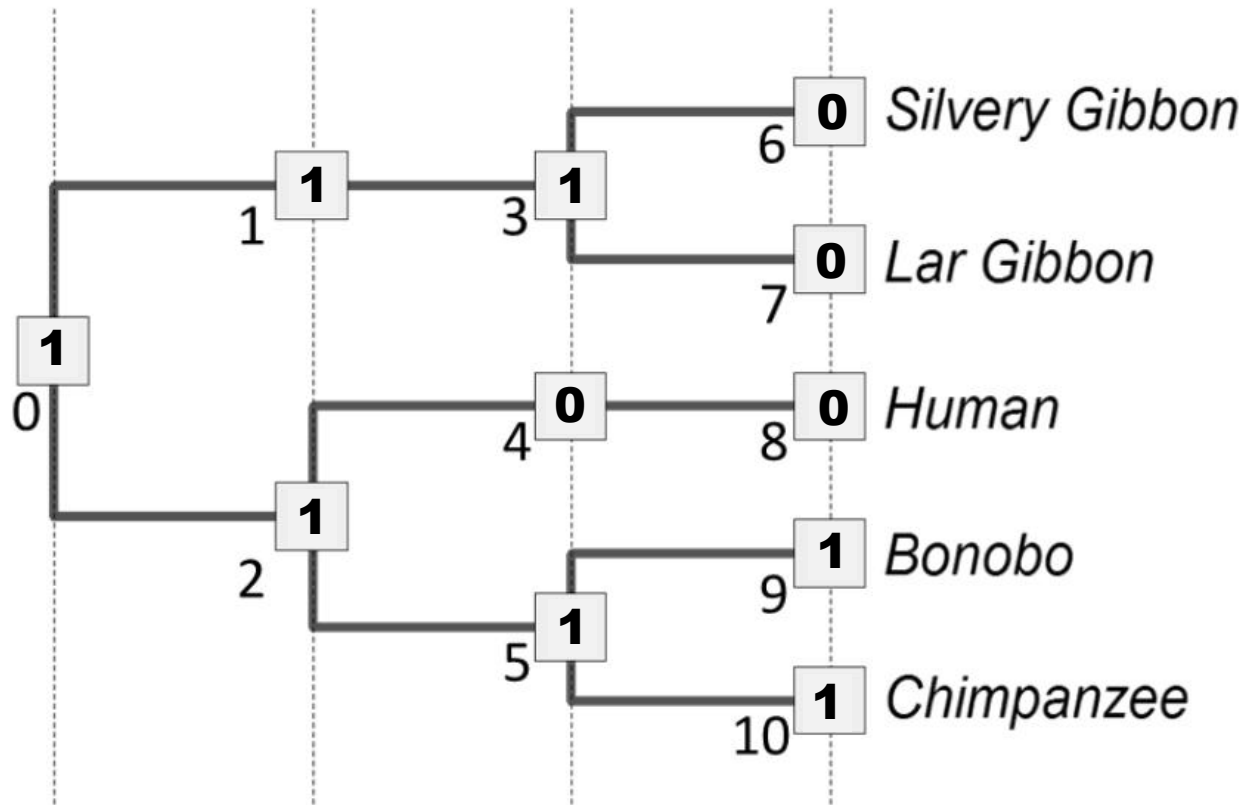
$$(1/6) + (1/6) = 1/3$$

# Combining probabilities: The OR rule in phylogenetics

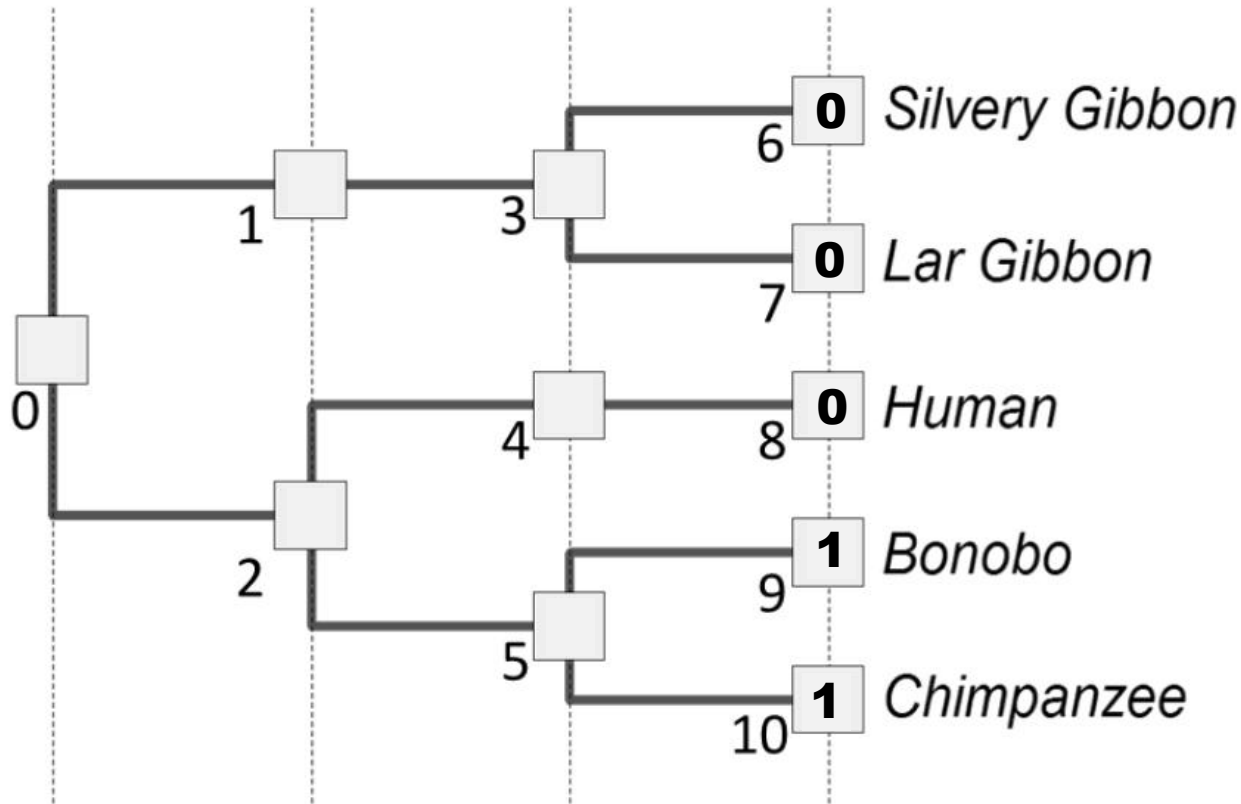


$$\text{Likelihood } (D, N_k | T_1) = (4/6)^1 \times (2/6)^9 = 0.000033$$

# Combining probabilities: The OR rule in phylogenetics



$$\text{Likelihood } (D, N_k | T_1) = (4/6)^7 \times (2/6)^3 = 0.00217$$



There are  $2^6 = 64$  possible node assignments. We could calculate the likelihood of each one, then (?) them together, to get the total  $L(\text{Data} \mid \text{Tree}_1)$ .

We used discrete time units.

Biology will want *continuous time*

# Continuous-Time Markov Models

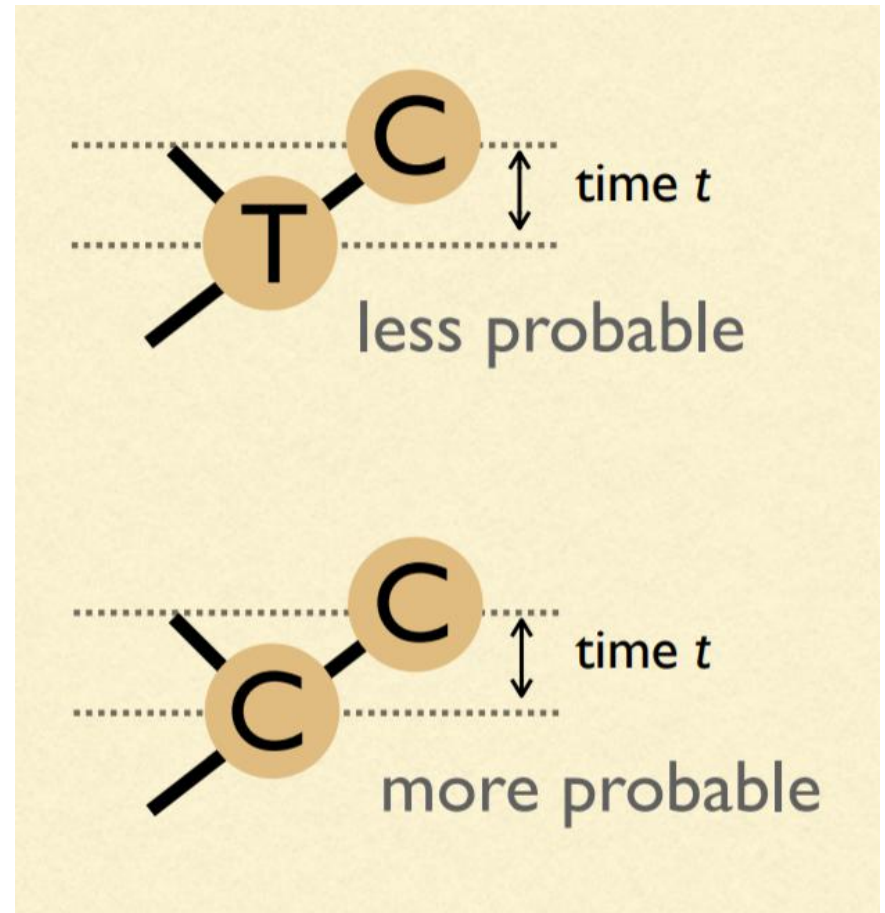
*Markov assumption* = probability of change depends only on current state, not how long it has been in that state

Our model of change depends on time:

We must estimate branch lengths

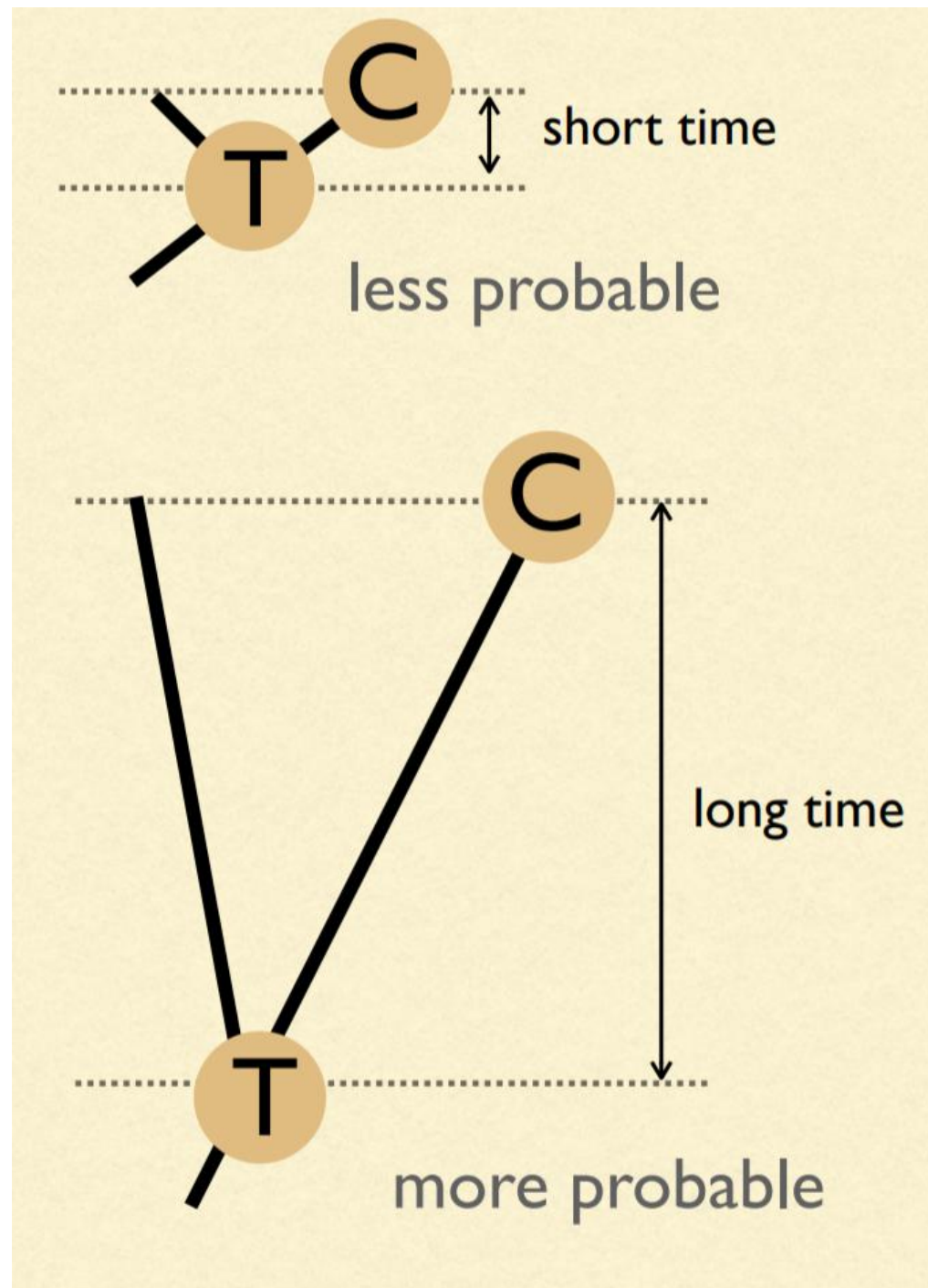
Units of branch length  
will be expected  
number of  
substitutions per site

(= rate of substitutions x  
time)

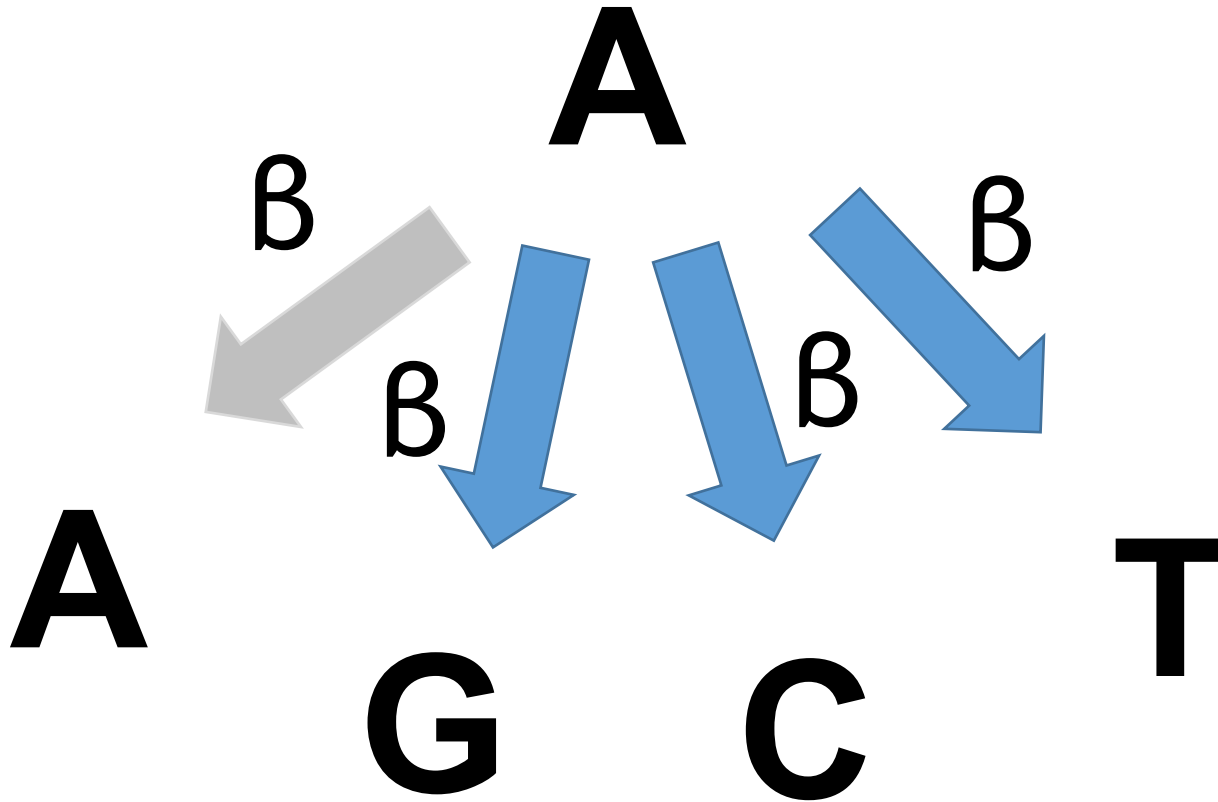




Probabilities  
are dependent  
on time



$$\mu = 4\beta$$



# A (very) simple phylogeny...



$$P_{AA} =$$

Probability nothing happened +

Probability something happened, but that the last thing that happened ended in an A

# A (very) simple phylogeny...



$$P_{AA} = (e^{-\mu t}) + (1 - e^{-\mu t})(1/4)$$

↑

Probability something doesn't happen

↑      ↑

Probability at least one thing happens      Probability that the last thing that happened ends in an A

# A (very) simple phylogeny...



$$P_{AG} =$$

$$(1 - e^{-\mu t})(1/4)$$



Probability  
at least  
one thing  
happens



Probability  
that the last  
thing that  
happened  
ends in an G

# A (very) simple phylogeny...



$$P_{AC} =$$

$$(1 - e^{-\mu t})(1/4)$$



Probability  
at least  
one thing  
happens



Probability  
that the last  
thing that  
happened  
ends in an C

# A (very) simple phylogeny...



$$P_{AT} =$$

$$(1 - e^{-\mu t})(1/4)$$



Probability  
at least  
one thing  
happens



Probability  
that the last  
thing that  
happened  
ends in an T

One last bit...substitutions vs.  
"events"

$$v = (3/4)\mu t = 3\beta t$$

$$4v/3 = \mu t$$

Only 3 out of 4 events results in a substitution. Thus, we can define the substitution rate  $v$ .



$$P_{AA} = (e^{-\mu t}) + (1 - e^{-\mu t})(1/4)$$

$$P_{AG} = (1 - e^{-\mu t})(1/4)$$

$$P_{AC} = (1 - e^{-\mu t})(1/4)$$

$$P_{AT} = (1 - e^{-\mu t})(1/4)$$

$$P_{AA} = (1/4) + (3/4)(e^{-4v/3})$$

$$P_{AG} = (1/4) - (1/4)(e^{-4v/3})$$

$$P_{AC} = (1/4) - (1/4)(e^{-4v/3})$$

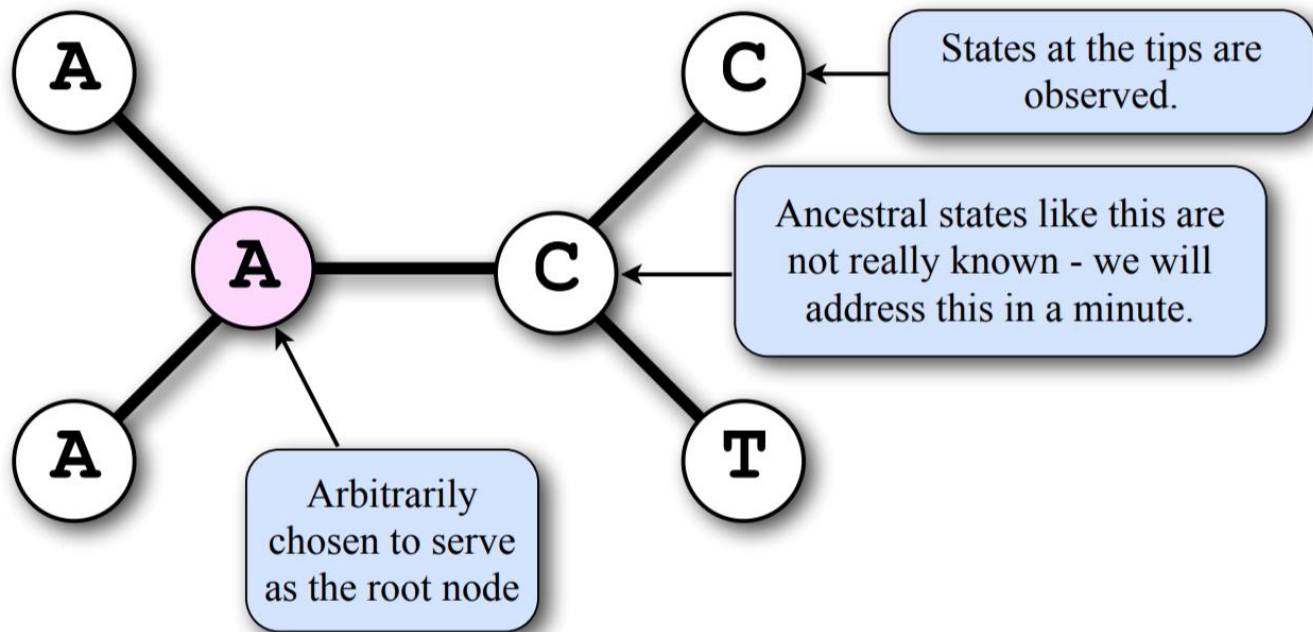
$$P_{AT} = (1/4) - (1/4)(e^{-4v/3})$$

Sanity check:

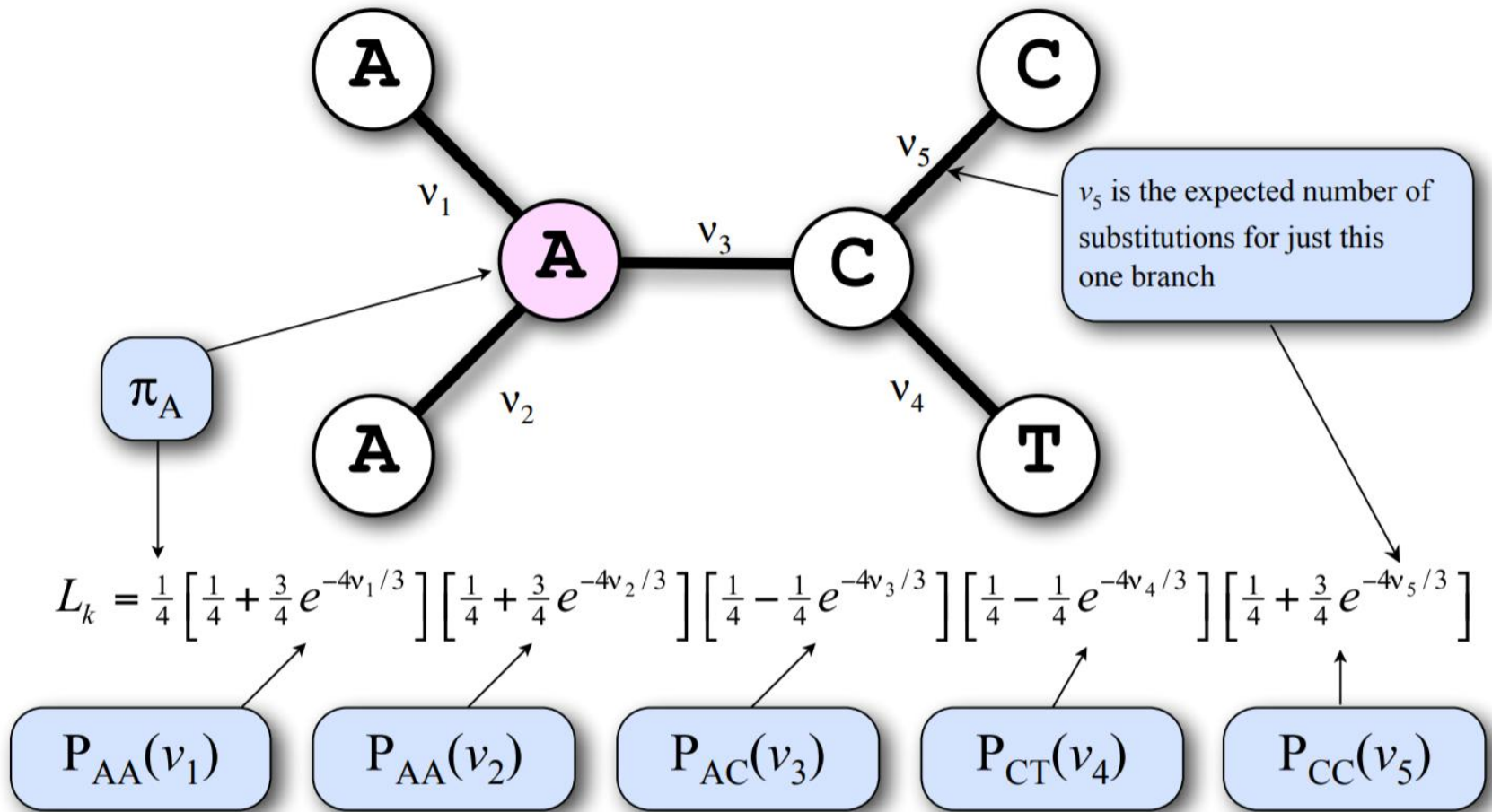
Do they all add to 1?

# Likelihood of an unrooted tree

(data shown for only one site)



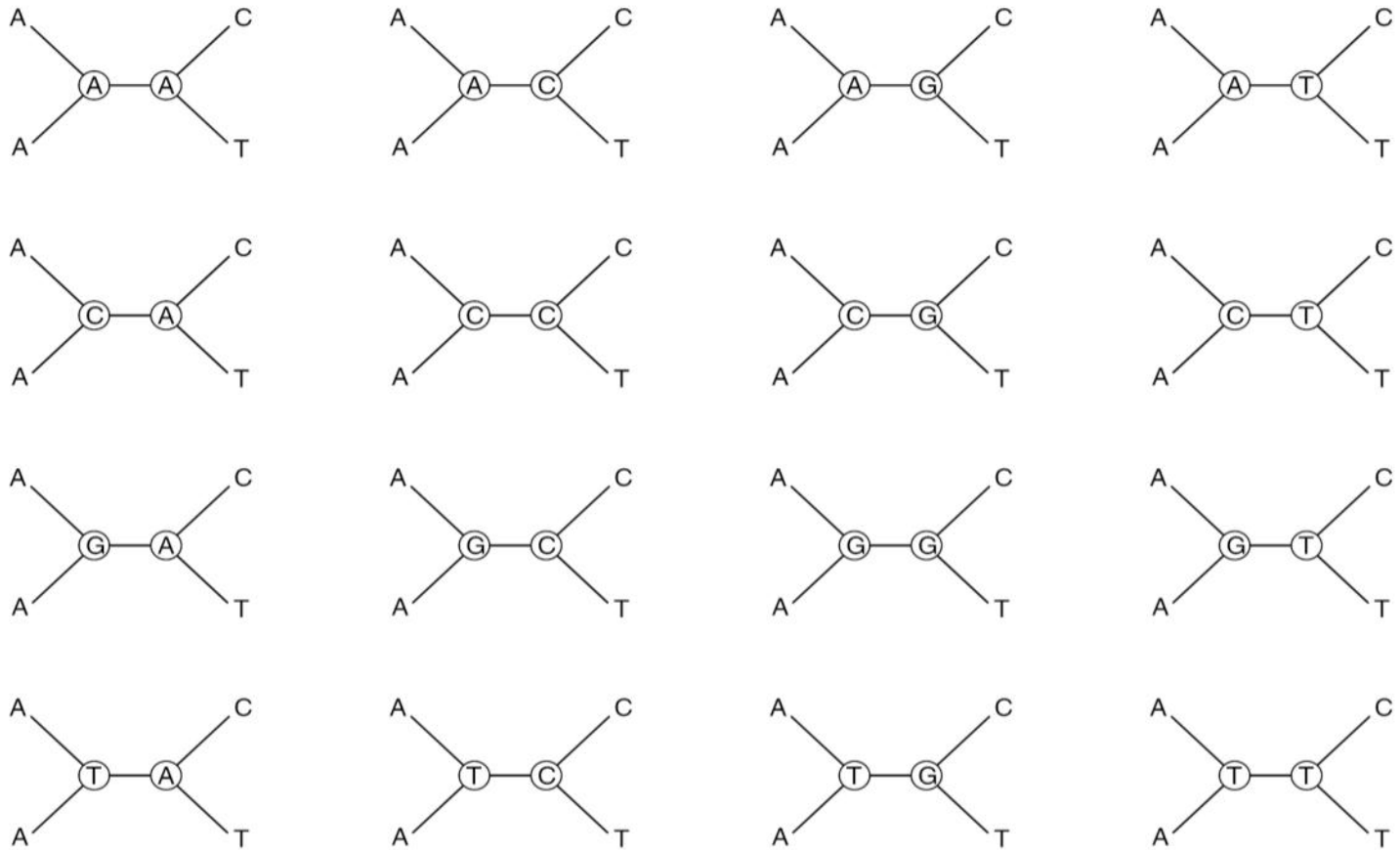
## Likelihood for site $k$



Paul O. Lewis (2014 Woods Hole Workshop in Molecular Evolution)

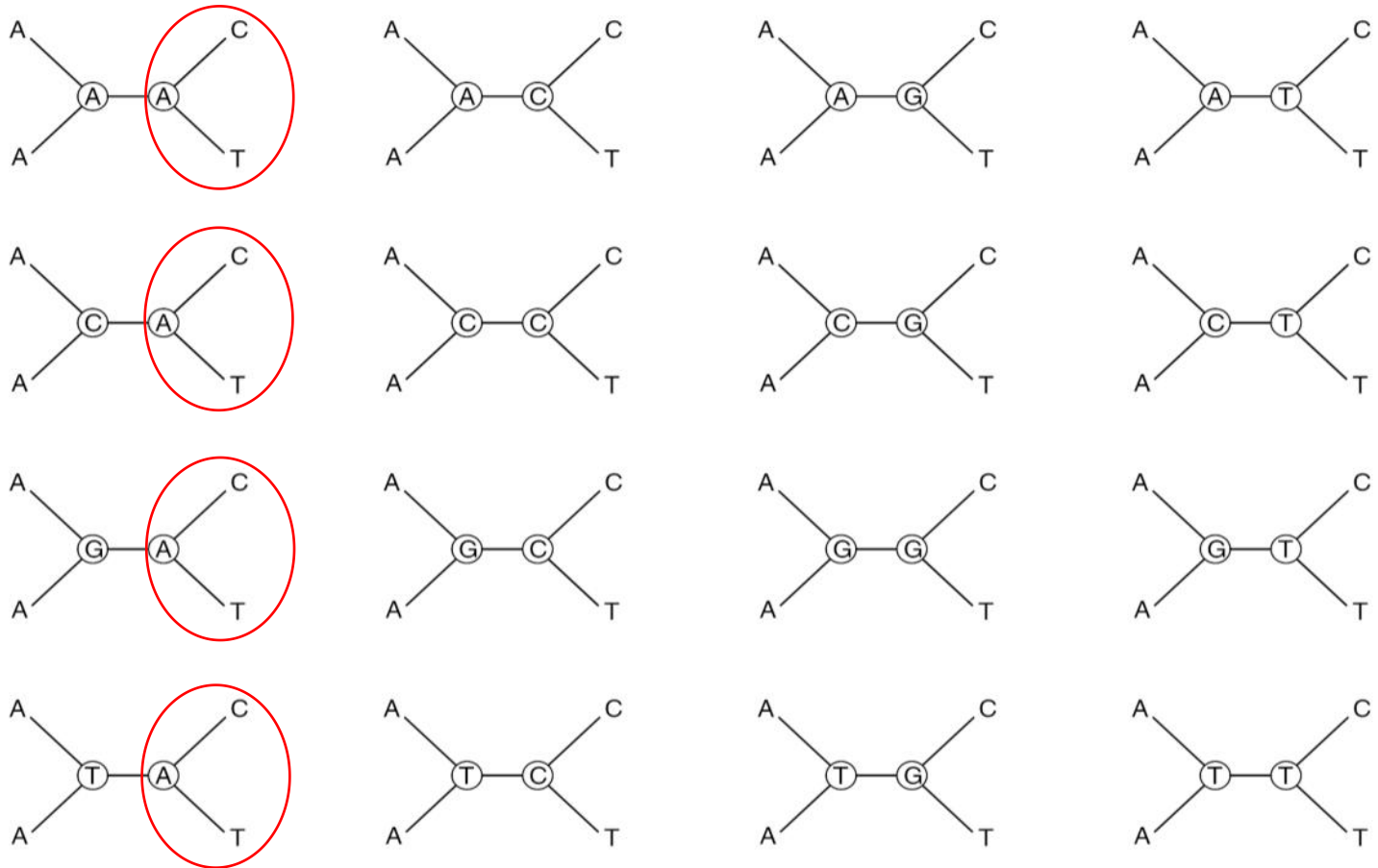
### Note use of the AND probability rule

Brute force approach would be to calculate  $L_k$  for all 16 combinations of ancestral states and sum them



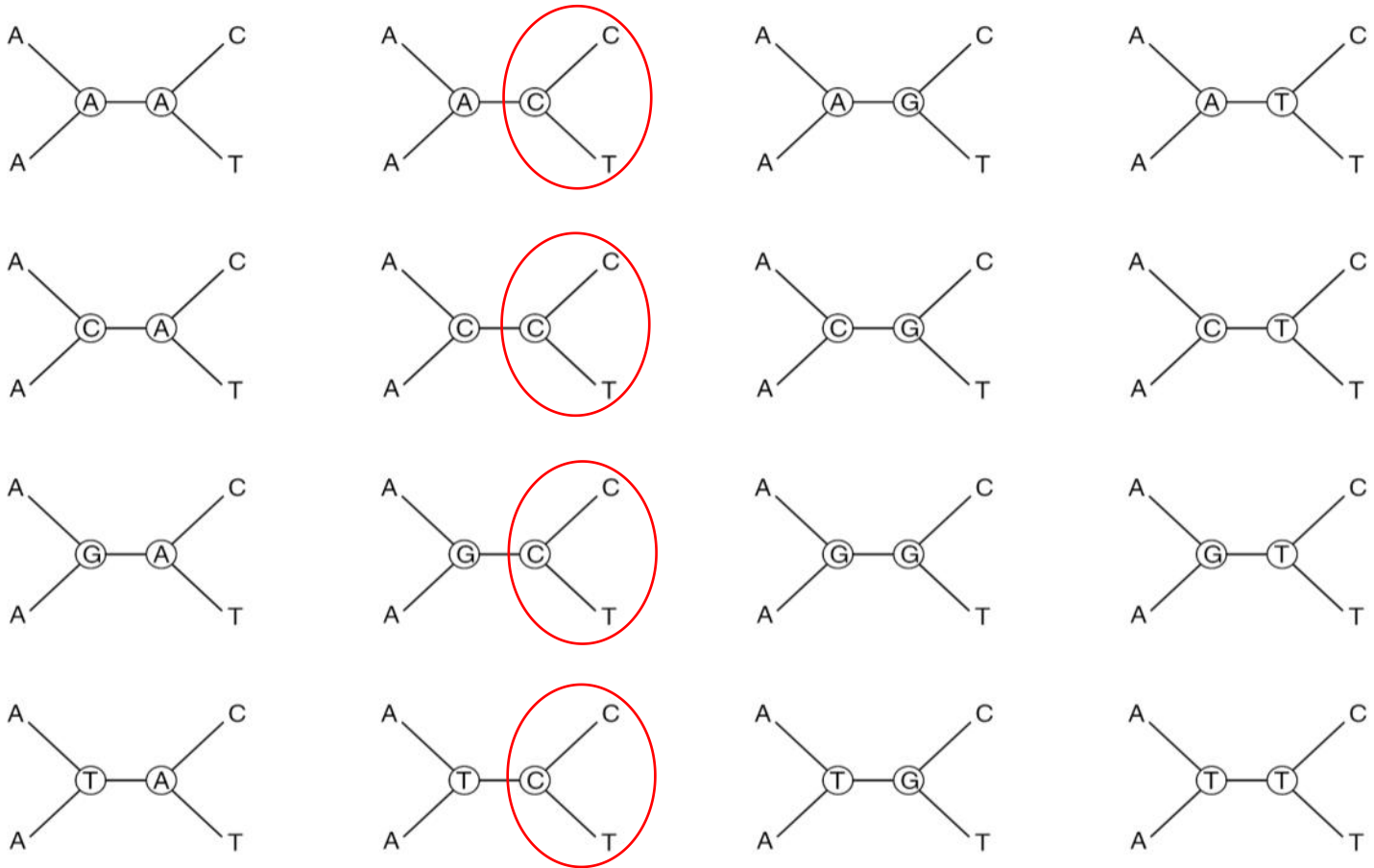
Note use of the OR probability rule

# Pruning algorithm



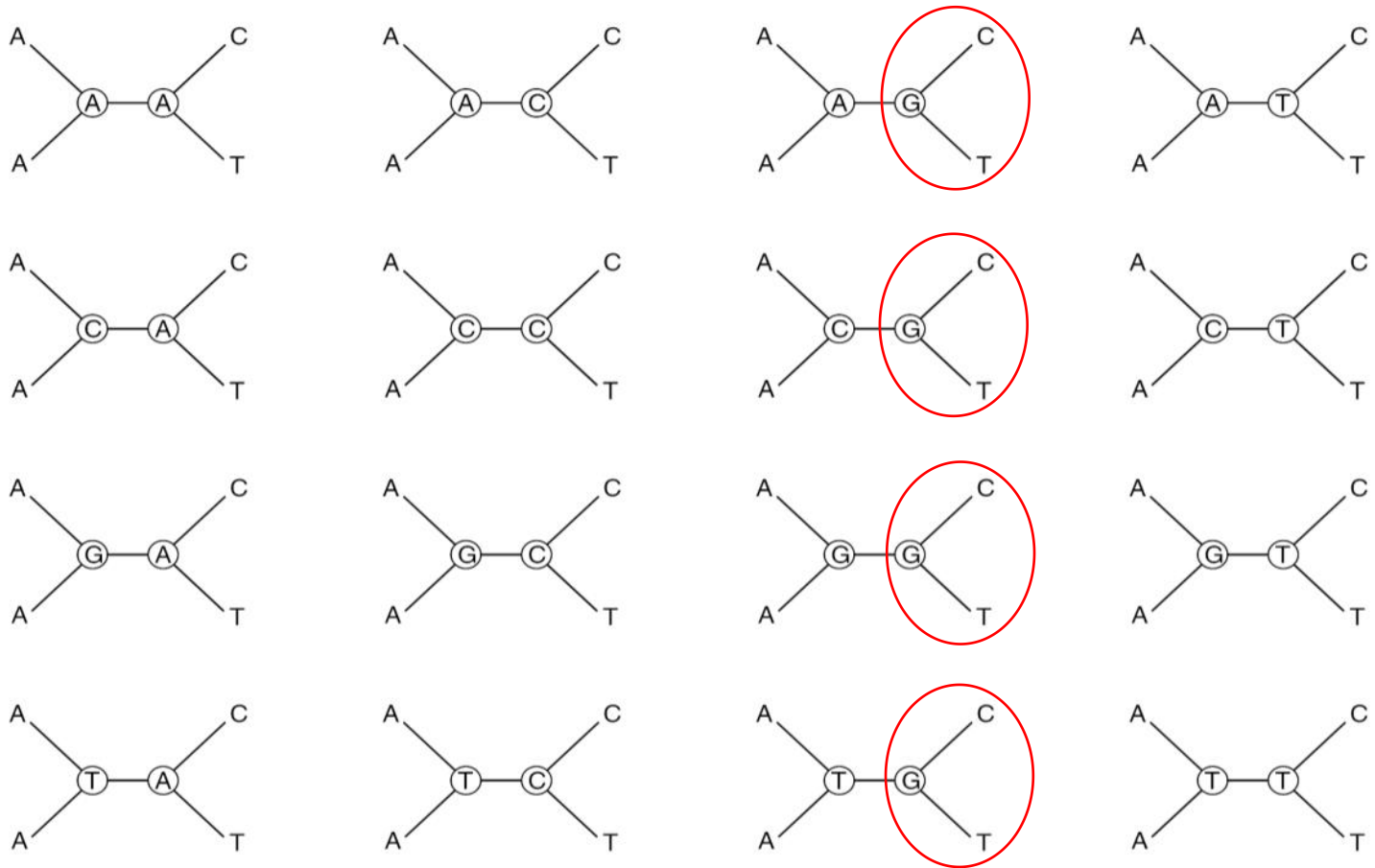
Note use of the OR probability rule

# Pruning algorithm



Note use of the OR probability rule

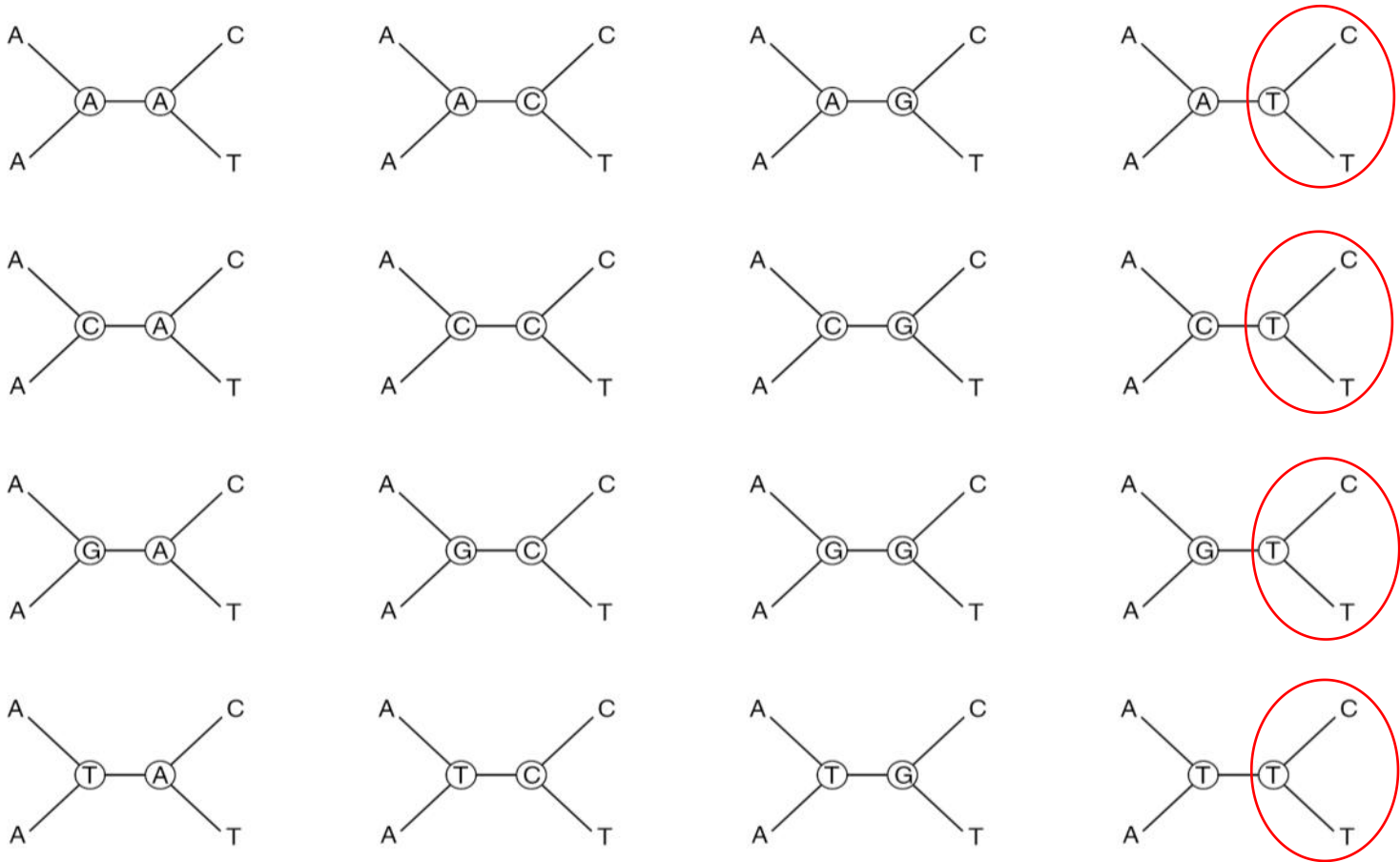
# Pruning algorithm



Note use of the OR probability rule

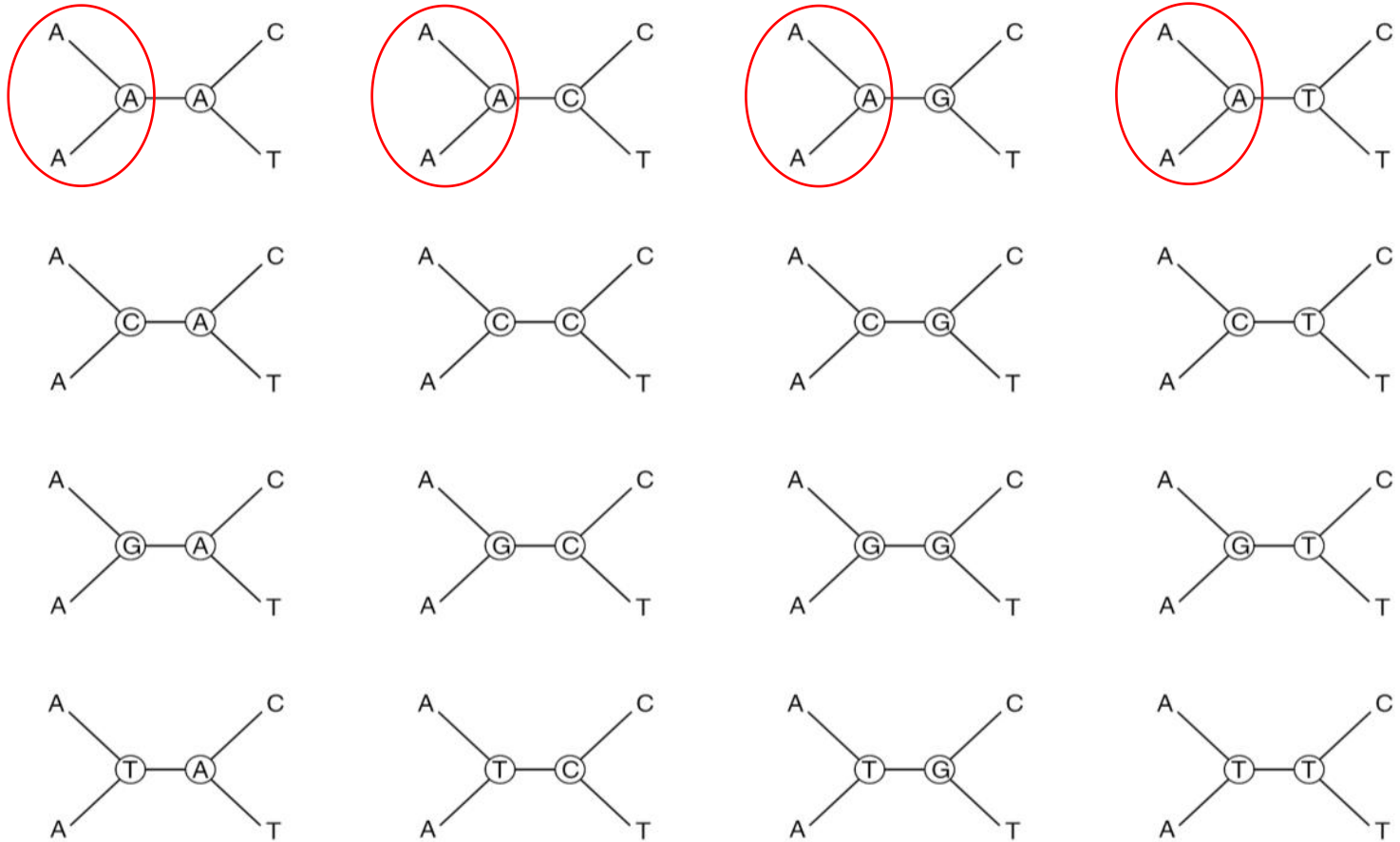


# Pruning algorithm



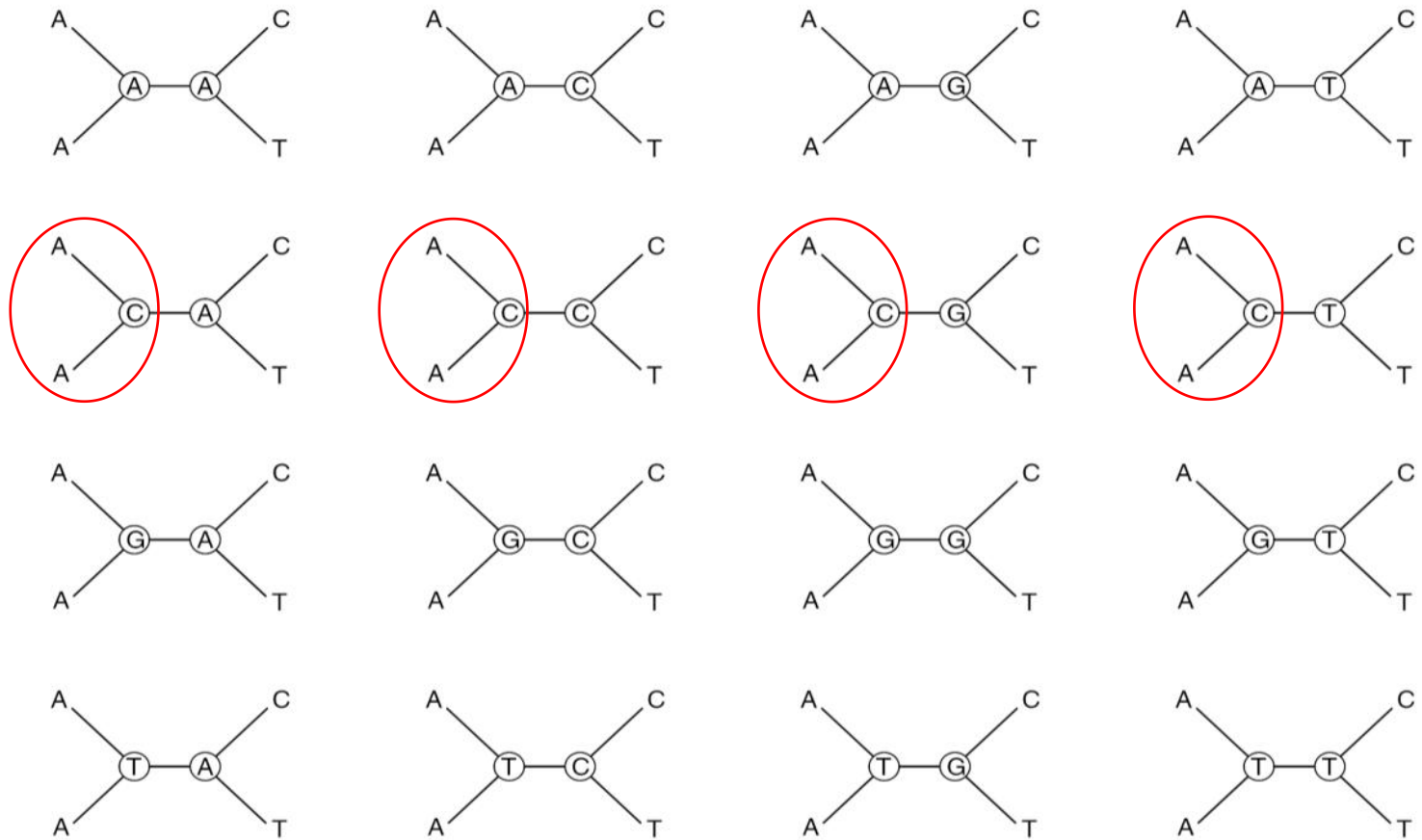
Note use of the OR probability rule

# Pruning algorithm



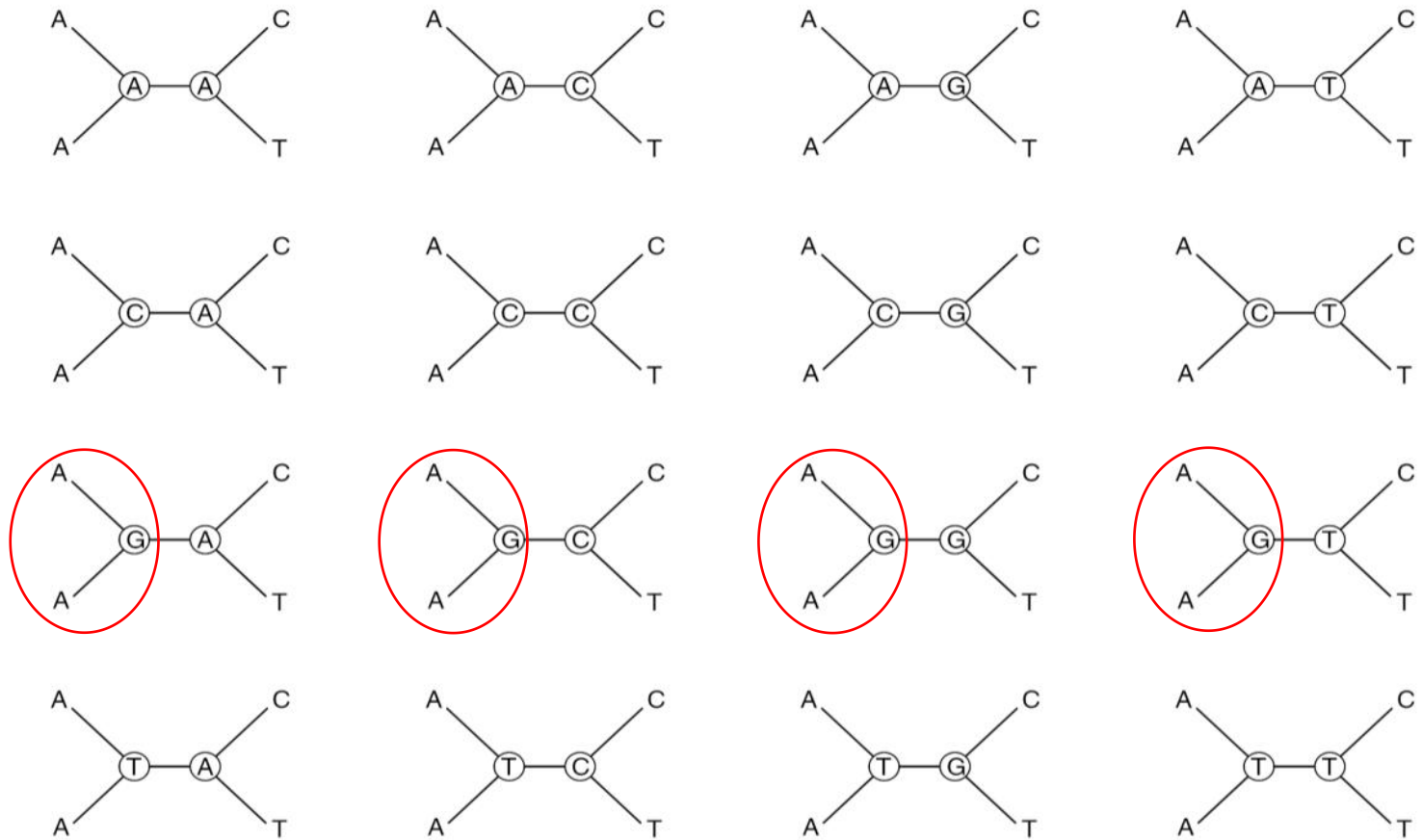
Note use of the OR probability rule

# Pruning algorithm



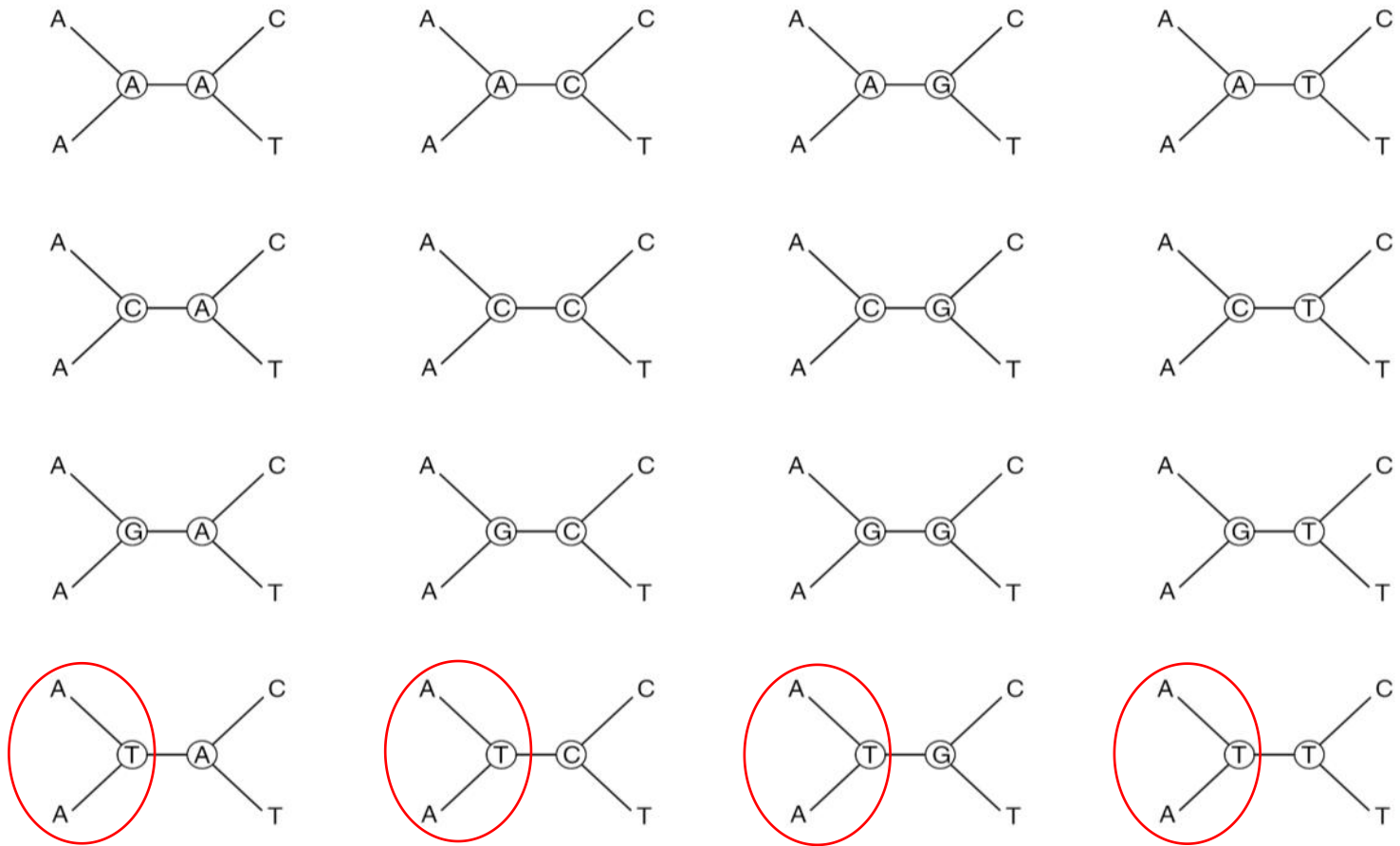
Note use of the OR probability rule

# Pruning algorithm



Note use of the OR probability rule

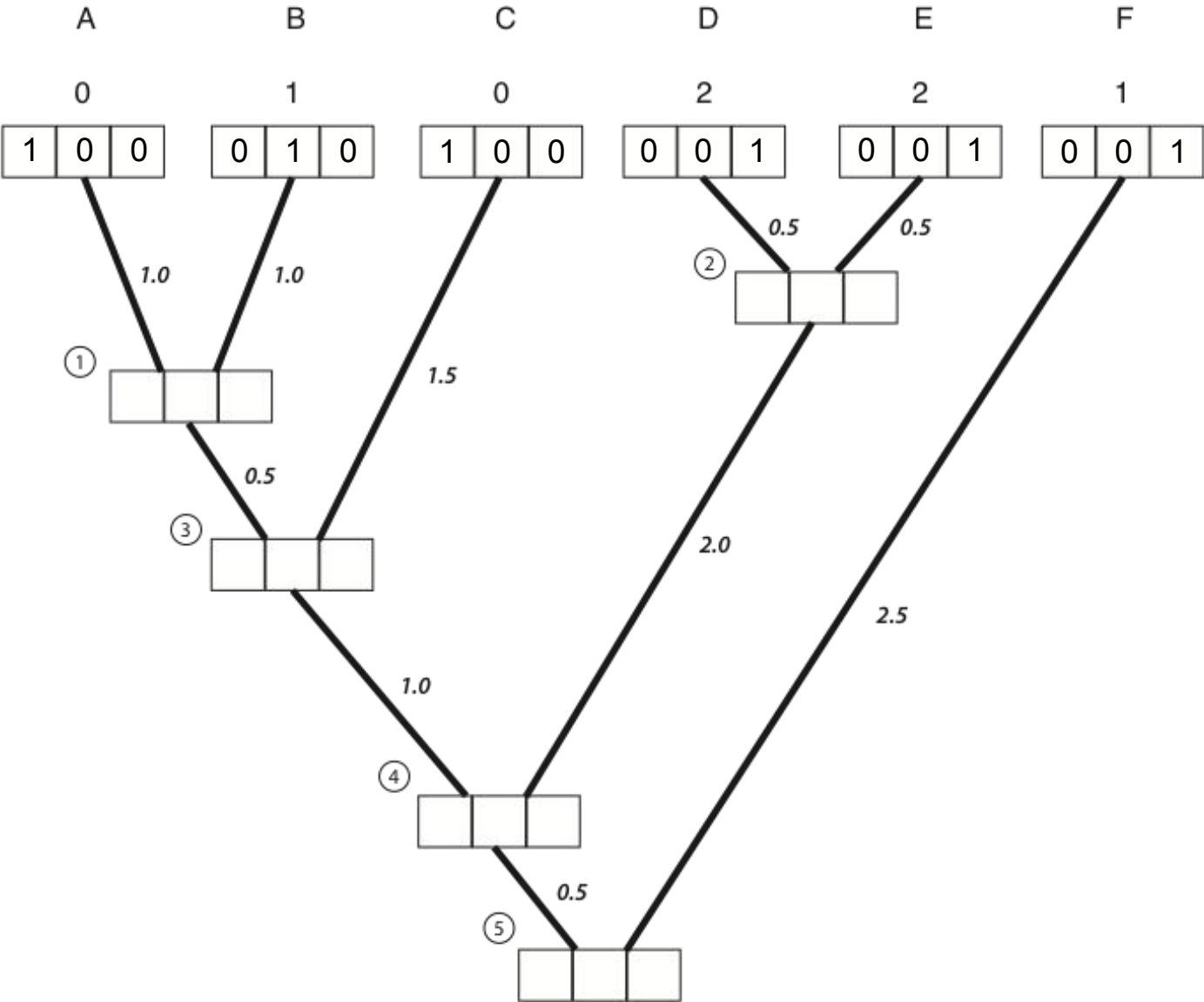
# Pruning algorithm



Note use of the OR probability rule

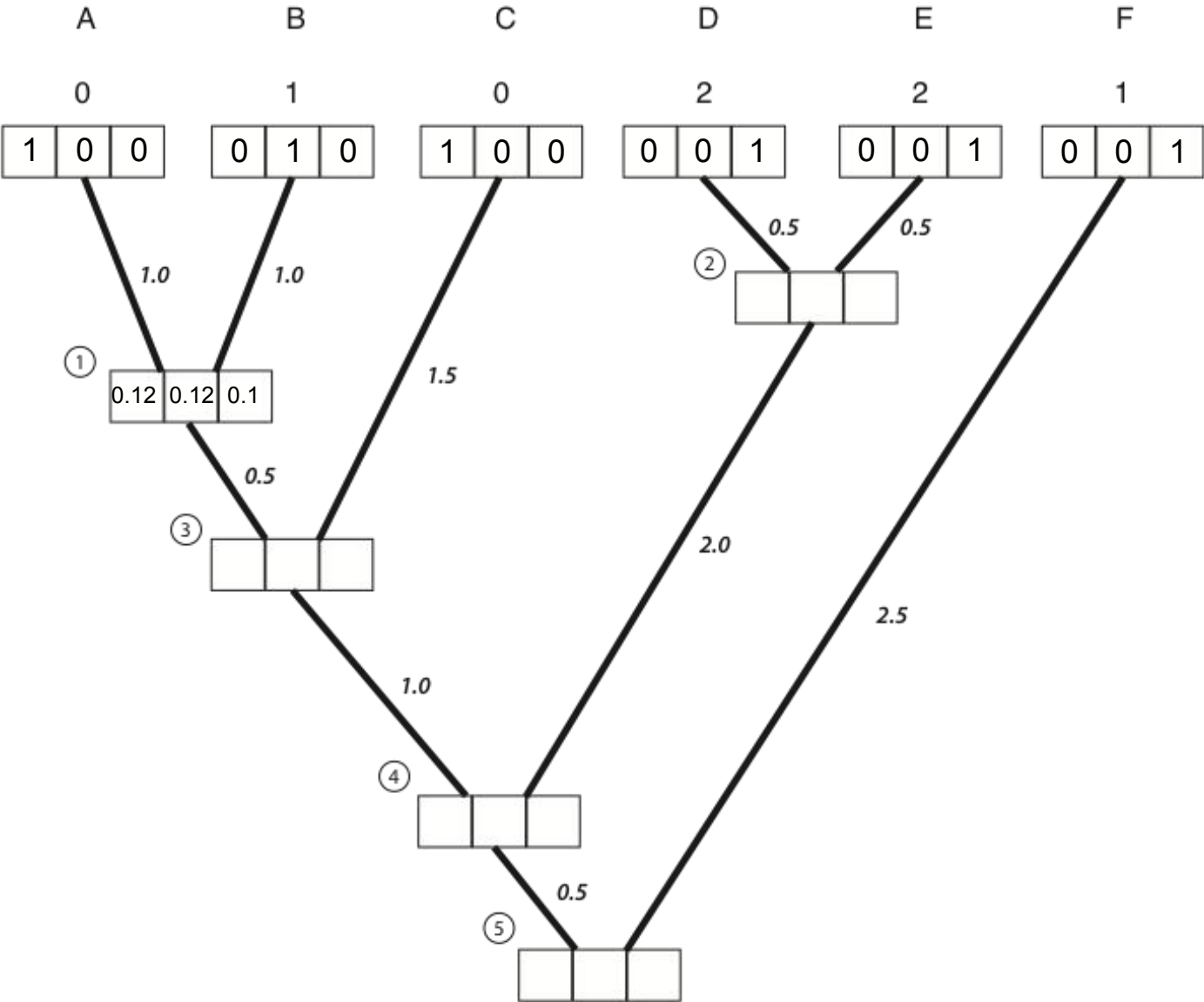
Species

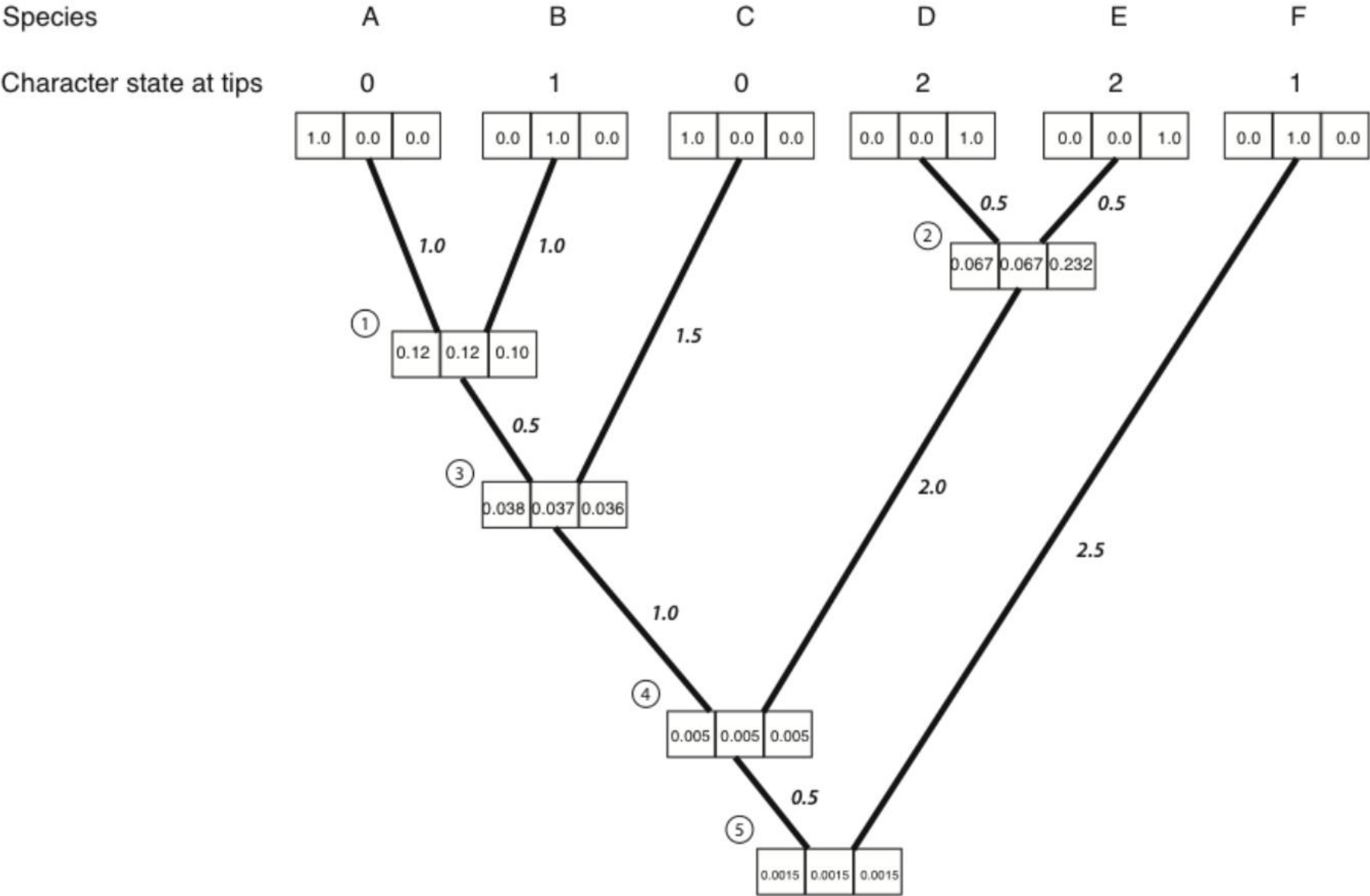
Character state at tips



Species

Character state at tips







# Jukes-Cantor Model (JC69)

Q matrix (*instantaneous rates*)

$$Q = \begin{array}{c} \begin{array}{c} A & C & G & T \end{array} \\ \begin{array}{c} A \\ C \\ G \\ T \end{array} \begin{bmatrix} -3\beta & \beta & \beta & \beta \\ \beta & -3\beta & \beta & \beta \\ \beta & \beta & -3\beta & \beta \\ \beta & \beta & \beta & -3\beta \end{bmatrix} \end{array}$$

# Jukes-Cantor Model (JC69)

*Transition probabilities:*

$$P = e^{Qt}$$

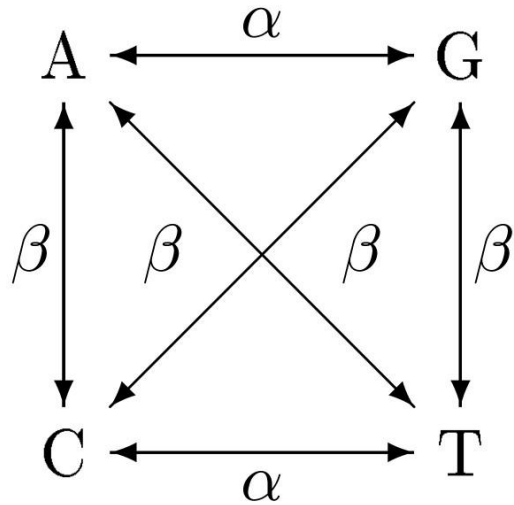


Matrix exponentiation

JC69 is our most basic model. We will be able to do amazing things with generalizations of this one model!

$$Q = \begin{matrix} & \begin{matrix} A & C & G & T \end{matrix} \\ \begin{matrix} A \\ C \\ G \\ T \end{matrix} & \begin{bmatrix} - & \beta & \beta & \beta \\ \beta & - & \beta & \beta \\ \beta & \beta & - & \beta \\ \beta & \beta & \beta & - \end{bmatrix} \end{matrix}$$

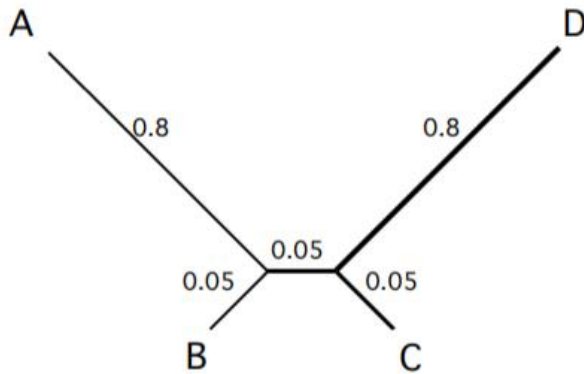
# Kimura 2 Parameter model: K2P



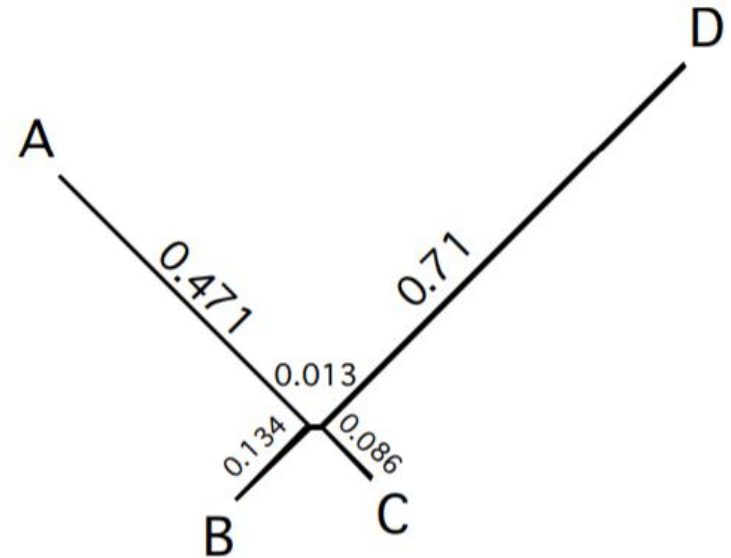
$Q =$

	A	C	G	T
A	-	$\beta$	$\alpha$	$\beta$
C	$\beta$	-	$\beta$	$\alpha$
G	$\alpha$	$\beta$	-	$\beta$
T	$\beta$	$\alpha$	$\beta$	-

# Returning to our original problem:



True Tree



ML Tree

# Differences between statistical phylogenetics and parsimony

We get branch lengths in expected number of changes rather than minimum # of changes

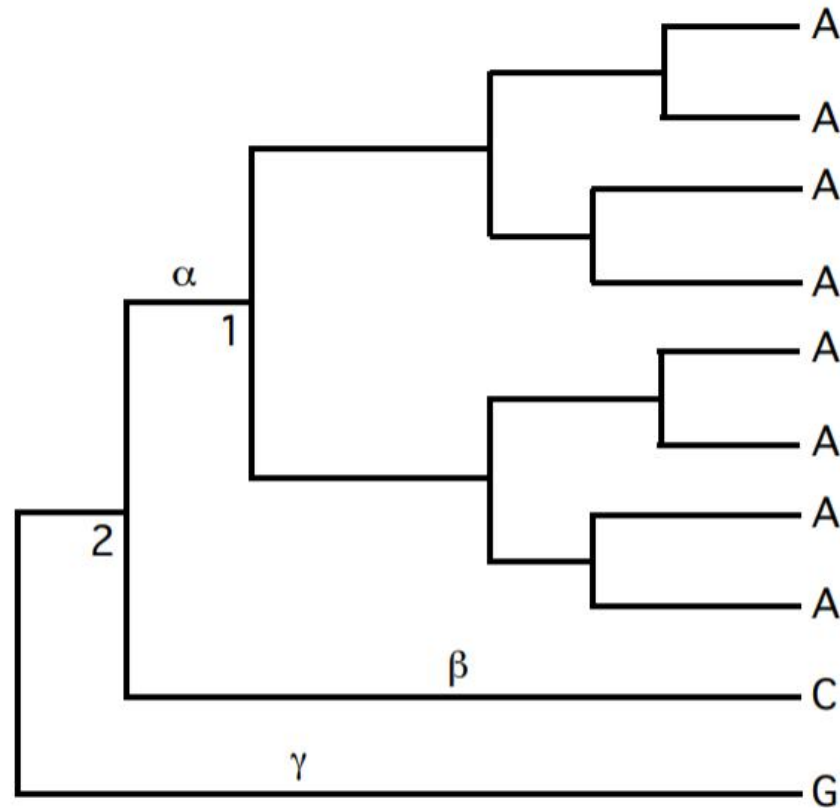
We expect and probabilistically incorporate all possible paths to the data, not just the shortest path

We have the flexibility to modify and compare models

We have a straight-forward way to convert branch lengths to time (with fossils or other constraints)

We use ALL OF THE DATA, not just parsimony-informative sites

Example: Parsimony says node 2 can be either A, C, or G with equal number of steps.



Example: Likelihood will say that node 2 is most likely A. Why?

