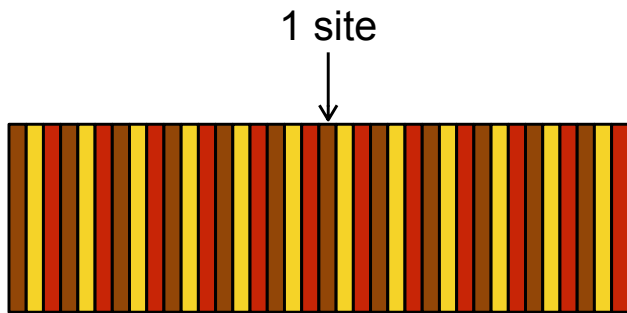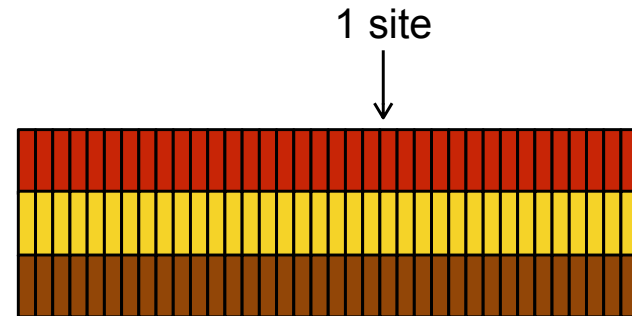# Two ways to model rate heterogeneity

1 site

1 site

## site-specific rates

each site assigned to 1
of 3 rate categories

## mixture model

each site has probability
1/3 of being in each of
the 3 rate categories

Dirichlet process priors provide a third option...

# Dirichlet Process Priors

A Bayesian Mixture Model for Across-Site Heterogeneities in the Amino-Acid Replacement Process

*Nicolas Lartillot and Hervé Philippe*

Lartillot and Philippe (2004)

PhyloBayes
http://www.atgc-montpellier.fr/phylobayes

A Dirichlet process model for detecting positive selection in protein-coding DNA sequences

John P. Huelsenbeck*[†], Sonia Jain[‡], Simon W. D. Frost[§], and Sergei L. Kosakovsky Pond[§]

Huelsenbeck et al. (2006)

Bayesian Estimation of Concordance among Gene Trees

*Cécile Ané,*[†] Bret Larget,*[†] David A. Baum,[†] Stacey D. Smith,[‡] and Antonis Rokas[§]

Ané et al (2007)

BUCKy
http://www.stat.wisc.edu/~ane/bucky/

A Nonparametric Method for Accommodating and Testing Across-Site Rate Variation

JOHN P. HUELSENBECK,[1] AND MARC A. SUCHARD[2,3,4]

Huelsenbeck and Suchard (2007)

# Dirichlet Process Prior

$$\frac{1}{\alpha+1}$$

(AB)

$$\frac{\alpha}{\alpha}$$

(A)

$$\frac{\alpha}{\alpha+1}$$

(A) (B)

Imagine you have a collection of objects (e.g. sites) labeled A, B, C, ...

B can either be added to A's group or form its own group

The parameter α determines the propensity for forming a new group

The third object C can either be added to an existing group...

...or form its own group

$$\frac{2}{\alpha + 2}$$ ABC

$$\frac{1}{\alpha + 1}$$ AB

$$\frac{\alpha}{\alpha + 2}$$ AB C

$$\frac{\alpha}{\alpha}$$ A

$$\frac{1}{\alpha + 2}$$ AC B

$$\frac{\alpha}{\alpha + 1}$$ A B

$$\frac{1}{\alpha + 2}$$ A BC

$$\frac{\alpha}{\alpha + 2}$$ A B C

After all objects have been considered, you can follow paths to determine the probability of different final configurations

Remember that this is a prior, so the data have a (usually big) say in how many clusters there are and what parameter values are assigned to each cluster.

5

$\frac{\alpha}{\alpha}$ (A)

$\frac{1}{\alpha+1}$

$\frac{\alpha}{\alpha+1}$

(AB)

$\frac{2}{\alpha+2}$

$\frac{\alpha}{\alpha+2}$

(A) (B)

$\frac{1}{\alpha+2}$

$\frac{1}{\alpha+2}$

$\frac{\alpha}{\alpha+2}$

(ABC)

$\frac{3}{\alpha+3}$ (ABCD) $\quad \left(\frac{\alpha}{\alpha}\right)\left(\frac{1}{\alpha+1}\right)\left(\frac{2}{\alpha+2}\right)\left(\frac{3}{\alpha+3}\right)$

$\frac{\alpha}{\alpha+3}$ (ABC) (D) $\quad \left(\frac{\alpha}{\alpha}\right)\left(\frac{1}{\alpha+1}\right)\left(\frac{2}{\alpha+2}\right)\left(\frac{\alpha}{\alpha+3}\right)$

(AB) (C)

$\frac{2}{\alpha+3}$ (ABD) (C) $\quad \left(\frac{\alpha}{\alpha}\right)\left(\frac{1}{\alpha+1}\right)\left(\frac{\alpha}{\alpha+2}\right)\left(\frac{2}{\alpha+3}\right)$
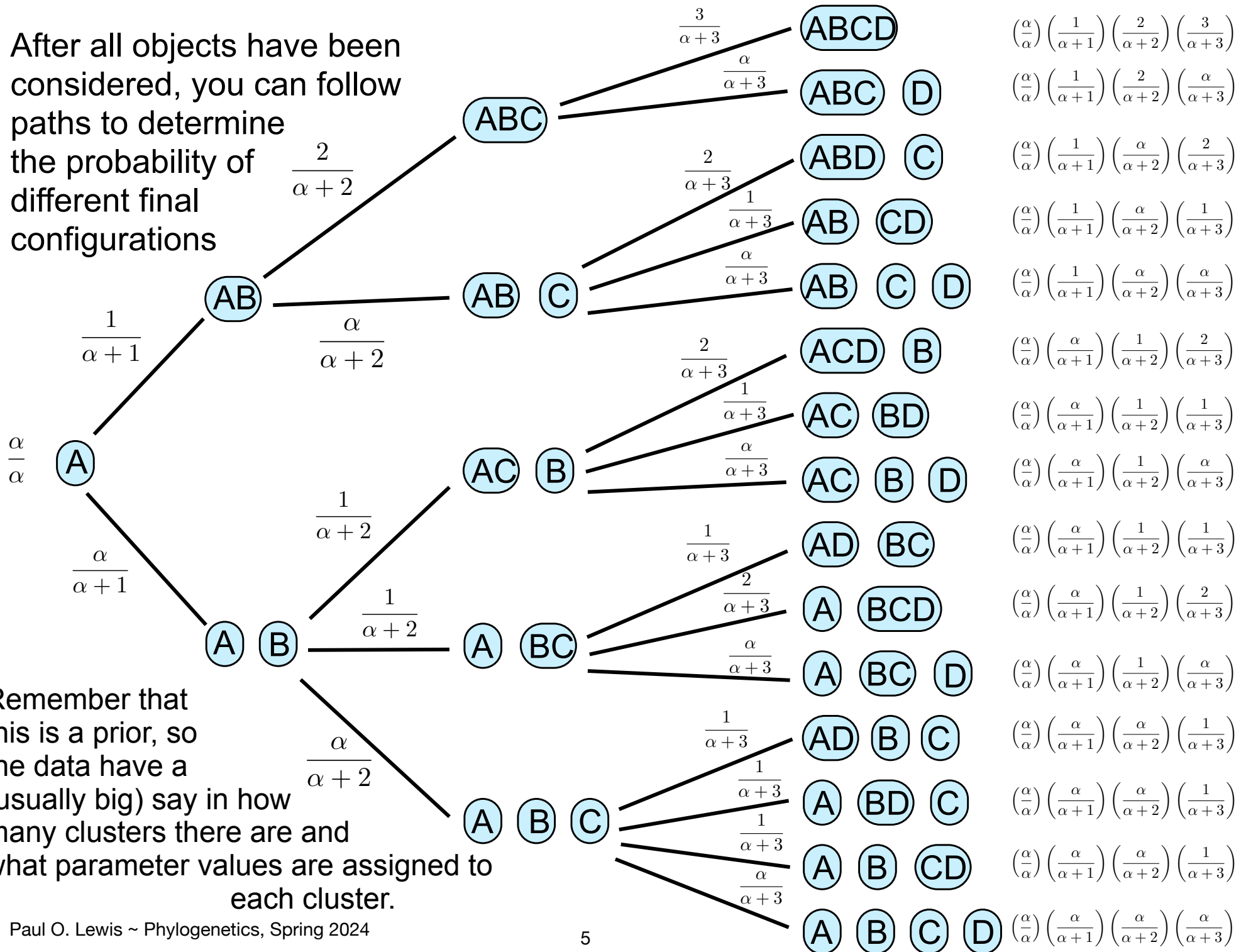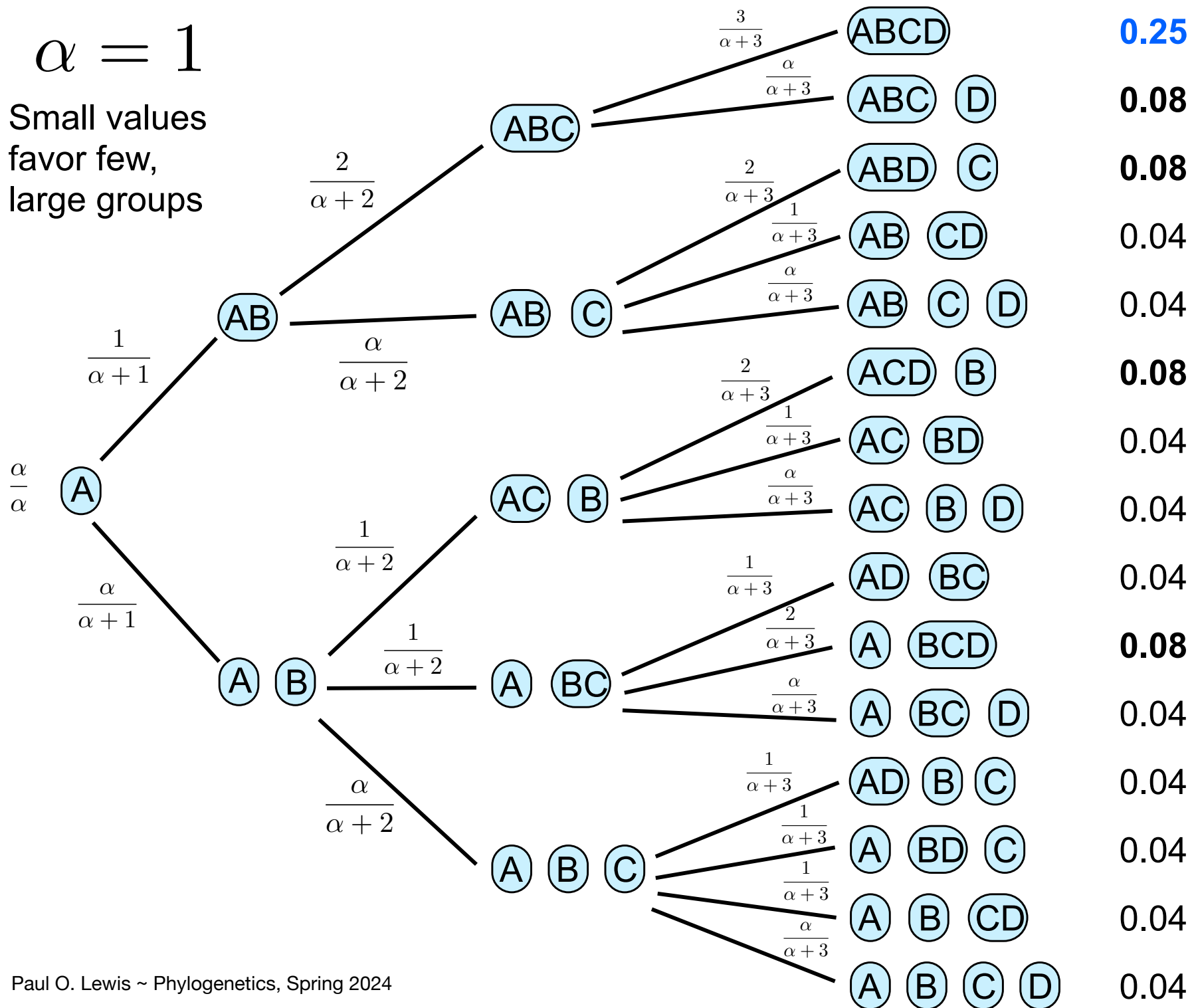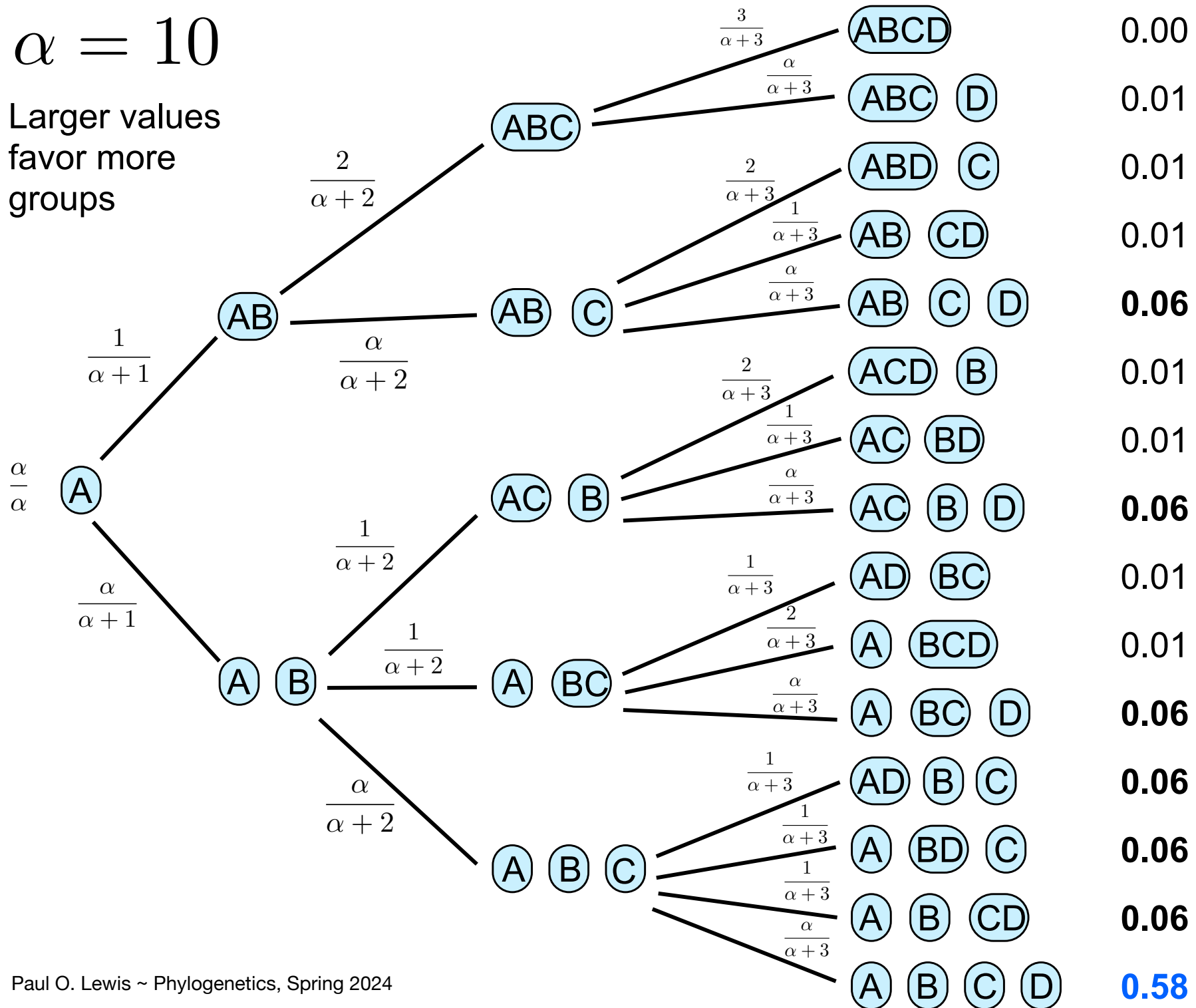
$\frac{1}{\alpha+3}$ (AB) (CD) $\quad \left(\frac{\alpha}{\alpha}\right)\left(\frac{1}{\alpha+1}\right)\left(\frac{\alpha}{\alpha+2}\right)\left(\frac{1}{\alpha+3}\right)$

$\frac{\alpha}{\alpha+3}$ (AB) (C) (D) $\quad \left(\frac{\alpha}{\alpha}\right)\left(\frac{1}{\alpha+1}\right)\left(\frac{\alpha}{\alpha+2}\right)\left(\frac{\alpha}{\alpha+3}\right)$

(AC) (B)

$\frac{2}{\alpha+3}$ (ACD) (B) $\quad \left(\frac{\alpha}{\alpha}\right)\left(\frac{\alpha}{\alpha+1}\right)\left(\frac{1}{\alpha+2}\right)\left(\frac{2}{\alpha+3}\right)$

$\frac{1}{\alpha+3}$ (AC) (BD) $\quad \left(\frac{\alpha}{\alpha}\right)\left(\frac{\alpha}{\alpha+1}\right)\left(\frac{1}{\alpha+2}\right)\left(\frac{1}{\alpha+3}\right)$

$\frac{\alpha}{\alpha+3}$ (AC) (B) (D) $\quad \left(\frac{\alpha}{\alpha}\right)\left(\frac{\alpha}{\alpha+1}\right)\left(\frac{1}{\alpha+2}\right)\left(\frac{\alpha}{\alpha+3}\right)$

(A) (BC)

$\frac{1}{\alpha+3}$ (AD) (BC) $\quad \left(\frac{\alpha}{\alpha}\right)\left(\frac{\alpha}{\alpha+1}\right)\left(\frac{1}{\alpha+2}\right)\left(\frac{1}{\alpha+3}\right)$

$\frac{2}{\alpha+3}$ (A) (BCD) $\quad \left(\frac{\alpha}{\alpha}\right)\left(\frac{\alpha}{\alpha+1}\right)\left(\frac{1}{\alpha+2}\right)\left(\frac{2}{\alpha+3}\right)$

$\frac{\alpha}{\alpha+3}$ (A) (BC) (D) $\quad \left(\frac{\alpha}{\alpha}\right)\left(\frac{\alpha}{\alpha+1}\right)\left(\frac{1}{\alpha+2}\right)\left(\frac{\alpha}{\alpha+3}\right)$

(A) (B) (C)

$\frac{1}{\alpha+3}$ (AD) (B) (C) $\quad \left(\frac{\alpha}{\alpha}\right)\left(\frac{\alpha}{\alpha+1}\right)\left(\frac{\alpha}{\alpha+2}\right)\left(\frac{1}{\alpha+3}\right)$

$\frac{1}{\alpha+3}$ (A) (BD) (C) $\quad \left(\frac{\alpha}{\alpha}\right)\left(\frac{\alpha}{\alpha+1}\right)\left(\frac{\alpha}{\alpha+2}\right)\left(\frac{1}{\alpha+3}\right)$

$\frac{1}{\alpha+3}$ (A) (B) (CD) $\quad \left(\frac{\alpha}{\alpha}\right)\left(\frac{\alpha}{\alpha+1}\right)\left(\frac{\alpha}{\alpha+2}\right)\left(\frac{1}{\alpha+3}\right)$

$\frac{\alpha}{\alpha+3}$ (A) (B) (C) (D) $\quad \left(\frac{\alpha}{\alpha}\right)\left(\frac{\alpha}{\alpha+1}\right)\left(\frac{\alpha}{\alpha+2}\right)\left(\frac{\alpha}{\alpha+3}\right)$

$\alpha = 1$

Small values favor few, large groups

$\frac{\alpha}{\alpha}$ (A)

$\frac{1}{\alpha+1}$    (AB)

$\frac{\alpha}{\alpha+1}$    (A)(B)

$\frac{2}{\alpha+2}$    (ABC)

$\frac{\alpha}{\alpha+2}$    (AB)(C)

$\frac{1}{\alpha+2}$    (AC)(B)

$\frac{1}{\alpha+2}$    (A)(BC)

$\frac{\alpha}{\alpha+2}$    (A)(B)(C)

$\frac{3}{\alpha+3}$   (ABCD)    **0.25**

$\frac{\alpha}{\alpha+3}$   (ABC)(D)    **0.08**

$\frac{2}{\alpha+3}$   (ABD)(C)    **0.08**

$\frac{1}{\alpha+3}$   (AB)(CD)    0.04

$\frac{\alpha}{\alpha+3}$   (AB)(C)(D)    0.04

$\frac{2}{\alpha+3}$   (ACD)(B)    **0.08**

$\frac{1}{\alpha+3}$   (AC)(BD)    0.04

$\frac{\alpha}{\alpha+3}$   (AC)(B)(D)    0.04

$\frac{1}{\alpha+3}$   (AD)(BC)    0.04

$\frac{2}{\alpha+3}$   (A)(BCD)    **0.08**

$\frac{\alpha}{\alpha+3}$   (A)(BC)(D)    0.04

$\frac{1}{\alpha+3}$   (AD)(B)(C)    0.04

$\frac{1}{\alpha+3}$   (A)(BD)(C)    0.04

$\frac{1}{\alpha+3}$   (A)(B)(CD)    0.04

$\frac{\alpha}{\alpha+3}$   (A)(B)(C)(D)    0.04

$\alpha = 10$

Larger values favor more groups

| Tree | Probability |
|------|-------------|
| ABCD | 0.00 |
| ABC D | 0.01 |
| ABD C | 0.01 |
| AB CD | 0.01 |
| AB C D | **0.06** |
| ACD B | 0.01 |
| AC BD | 0.01 |
| AC B D | **0.06** |
| AD BC | 0.01 |
| A BCD | 0.01 |
| A BC D | **0.06** |
| AD B C | **0.06** |
| A BD C | **0.06** |
| A B CD | **0.06** |
| A B C D | **0.58** |

Branch probabilities shown on tree:

From A: $\frac{\alpha}{\alpha}$

$\frac{1}{\alpha+1}$ to AB, $\frac{\alpha}{\alpha+1}$ to A B

From AB: $\frac{2}{\alpha+2}$ to ABC, $\frac{\alpha}{\alpha+2}$ to AB C

From A B: $\frac{1}{\alpha+2}$ to AC B, $\frac{1}{\alpha+2}$ to A BC, $\frac{\alpha}{\alpha+2}$ to A B C

From ABC: $\frac{3}{\alpha+3}$ to ABCD, $\frac{\alpha}{\alpha+3}$ to ABC D

From AB C: $\frac{2}{\alpha+3}$ to ABD C, $\frac{1}{\alpha+3}$ to AB CD, $\frac{\alpha}{\alpha+3}$ to AB C D

From AC B: $\frac{2}{\alpha+3}$ to ACD B, $\frac{1}{\alpha+3}$ to AC BD, $\frac{\alpha}{\alpha+3}$ to AC B D

From A BC: $\frac{1}{\alpha+3}$ to AD BC, $\frac{2}{\alpha+3}$ to A BCD, $\frac{\alpha}{\alpha+3}$ to A BC D

From A B C: $\frac{1}{\alpha+3}$ to AD B C, $\frac{1}{\alpha+3}$ to A BD C, $\frac{1}{\alpha+3}$ to A B CD, $\frac{\alpha}{\alpha+3}$ to A B C D

# Expected Number of Groups

$$\sum_{i=0}^{n-1} \frac{\alpha}{\alpha + i}$$

For example, if *n* = 3 and **α** = 1:

$$\frac{\alpha}{\alpha + 0} + \frac{\alpha}{\alpha + 1} + \frac{\alpha}{\alpha + 2} = 1 + \frac{1}{2} + \frac{1}{3} = 1.83$$

# Restaurant Analogy

# Dirichlet process prior applet
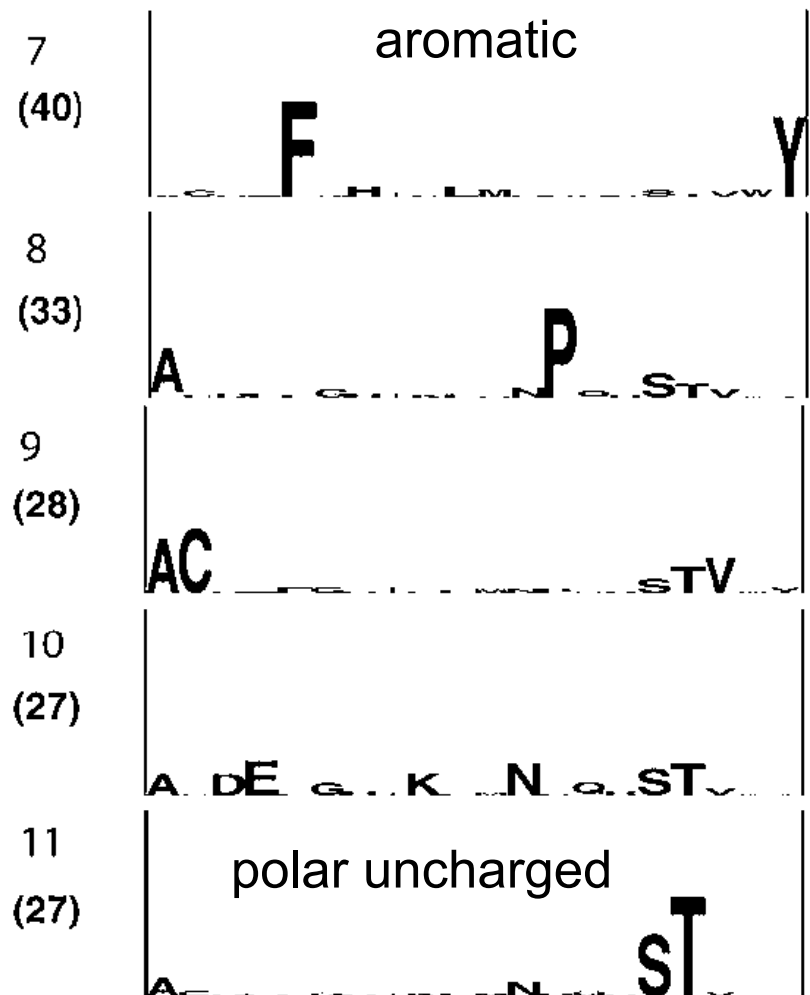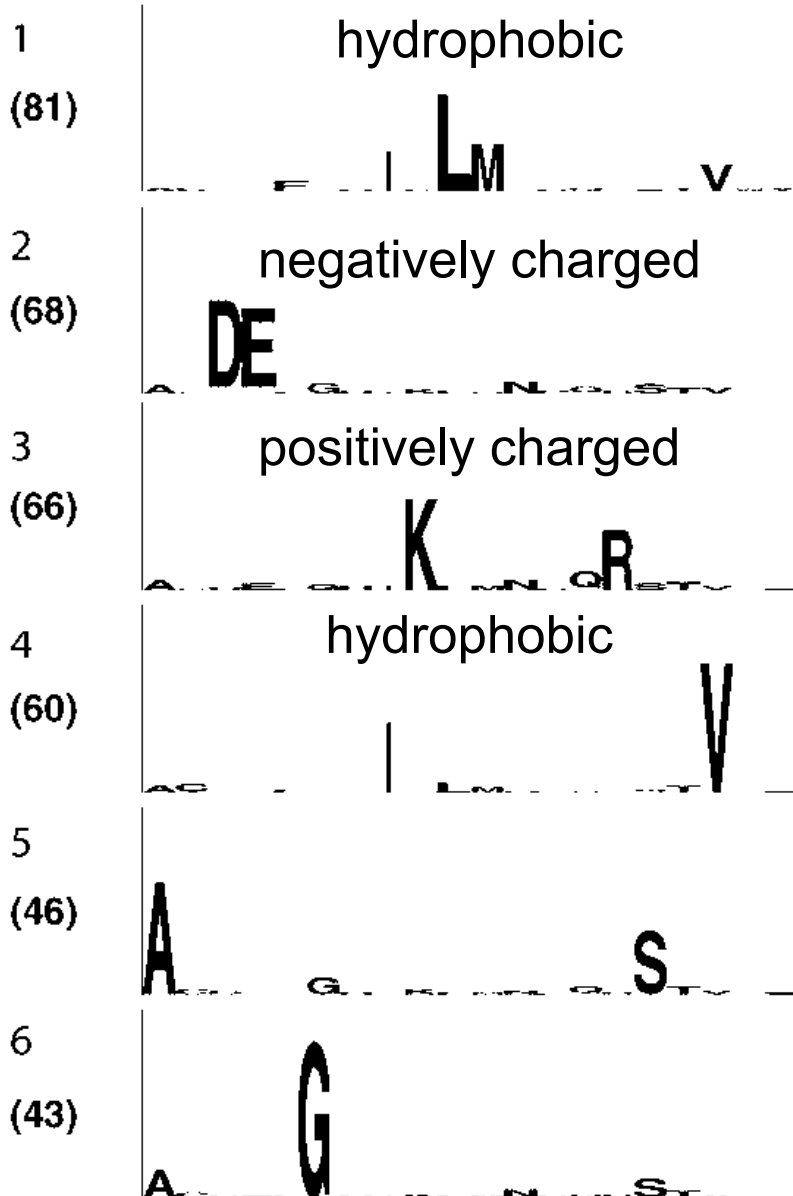
https://plewis.github.io/applets/dpp/

# Example 1:Elongation factor 2
# (software: PhyloBayes)



This model is implemented in PhyloBayes software: http://www.atgc-montpellier.fr/phylobayes/

Lartillot and Philippe (2004)

D

| | |
|---|---|
| 1 (81) | hydrophobic |
| 2 (68) | negatively charged |
| 3 (66) | positively charged |
| 4 (60) | hydrophobic |
| 5 (46) | |
| 6 (43) | |

| | |
|---|---|
| 7 (40) | aromatic |
| 8 (33) | |
| 9 (28) | |
| 10 (27) | |
| 11 (27) | polar uncharged |

L=leucine, M=methionine
D=aspartic acid, E=glutamic acid
K=lysine, R=arginine, G=glycine,
V=valine, I=isoleucine, A=alanine,
F=phenylalanine, Y=tyrosine
S=serine, T=threonine

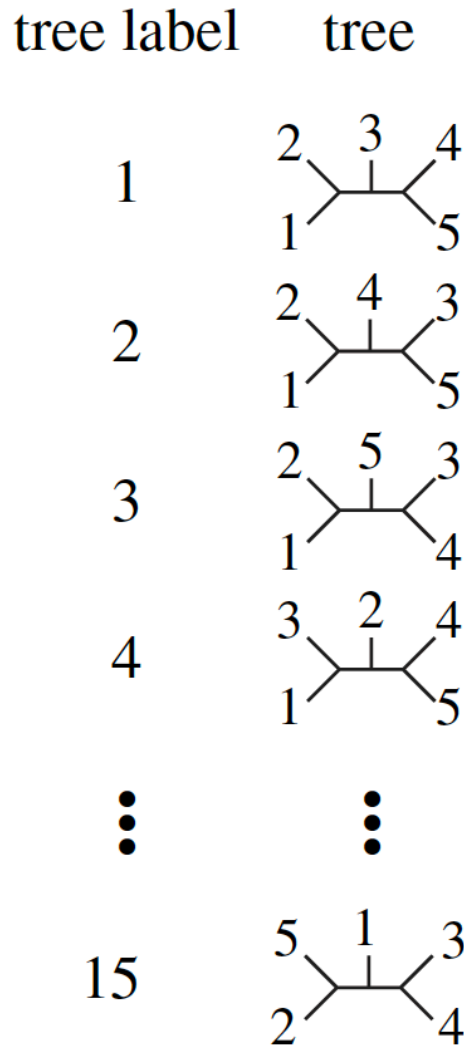# Example 2: Bayesian concordance analyses (software: BUCKy)

| tree label | tree | $m_1$ | | | $m_2$ | | |
|---|---|---|---|---|---|---|---|
| | | $g_1$ | $g_2$ | $g_3$ | $g_1$ | $g_2$ | $g_3$ |
| 1 | 2 3 4 / 1 5 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | 2 4 3 / 1 5 | **1** | **1** | **1** | **1** | **1** | 0 |
| 3 | 2 5 3 / 1 4 | 0 | 0 | 0 | 0 | 0 | **1** |
| 4 | 3 2 4 / 1 5 | 0 | 0 | 0 | 0 | 0 | 0 |
| ⋮ | ⋮ | | ⋮ | | | ⋮ | |
| 15 | 5 1 3 / 2 4 | 0 | 0 | 0 | 0 | 0 | 0 |
| | | (2,2,2) | | | (2,2,3) | | |

Marginal posterior distributions for each gene separately

| | gene | | |
|---|---|---|---|
| tree | 1 | 2 | 3 |
| 2 | 1 | 0 | 0 |
| 3 | 0 | 0.9 | 0.2 |
| 4 | 0 | 0.1 | 0.2 |
| 15 | 0 | 0 | 0.6 |

Gene-to-Tree Mapping (GTM)

Ané et al. (2007)

# Concordance Factors (CF)

tree label     tree



| GTM | Posterior | K | CF 12\|345 |
|---|---|---|---|
| (2,3,3) | 0.6600 | | |
| (2,3,4) | 0.0600 | | |
| (2,3,15) | 0.1800 | | |
| (2,4,3) | 0.0067 | | |
| (2,4,4) | 0.0733 | | |
| (2,4,15) | 0.0200 | | |

1.0000

# Concordance Factors (CF)

tree label     tree



| GTM | Posterior | K | CF 12\|345 |
|---|---|---|---|
| (2,3,3) | 0.6600 | 2 | |
| (2,3,4) | 0.0600 | 3 | |
| (2,3,15) | 0.1800 | 3 | |
| (2,4,3) | 0.0067 | 3 | |
| (2,4,4) | 0.0733 | 2 | |
| (2,4,15) | 0.0200 | 3 | |

           1.0000      2.27

# Concordance Factors (CF)

tree label    tree

1    (tree with tips 2, 3, 4 / 1, 5)

2    (tree with tips 2, 4, 3 / 1, 5)

3    (tree with tips 2, 5, 3 / 1, 4)

4    (tree with tips 3, 2, 4 / 1, 5)

⋮    ⋮

15    (tree with tips 5, 1, 3 / 2, 4)

| GTM | Posterior | K | CF 12\|345 |
|---|---|---|---|
| (2,3,3) | 0.6600 | 2 | 1 |
| (2,3,4) | 0.0600 | 3 | 2/3 |
| (2,3,15) | 0.1800 | 3 | 2/3 |
| (2,4,3) | 0.0067 | 3 | 2/3 |
| (2,4,4) | 0.0733 | 2 | 1/3 |
| (2,4,15) | 0.0200 | 3 | 1/3 |

        1.0000     2.27     0.86

# Dirichlet Process Priors

- To encourage **few, large** groups, use a **small** alpha value

- To encourage **lots of small** groups, use a **large** alpha value

- In practice, **hierarchical models** are used (i.e. alpha is a hyperparameter that can be estimated, so you need not worry about choosing the appropriate value for alpha)

- Bottom line: DP models are very nice for automatically grouping sites into clusters that have some property in common