
Lecture 1.2

Evolutionary Models

Popular phylogenetic methods

1. Maximum parsimony
2. Distance-based methods
3. Maximum likelihood
4. Bayesian inference

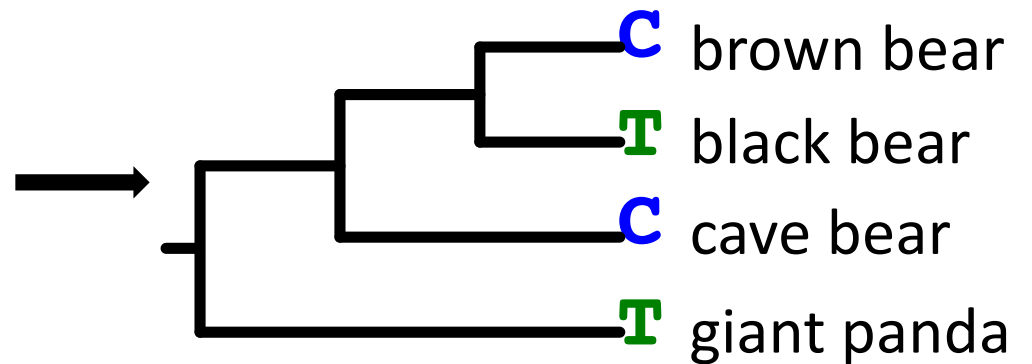
Model-based methods



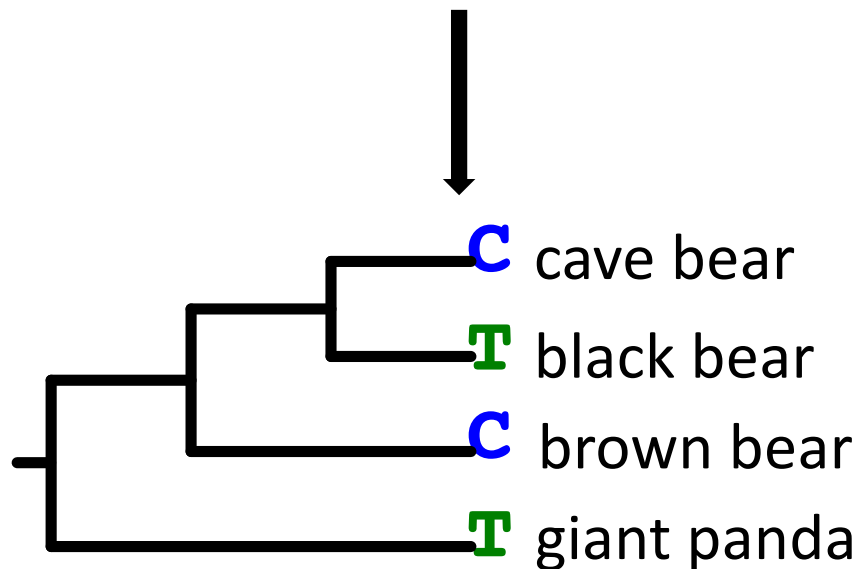
Maximum Parsimony

Maximum parsimony

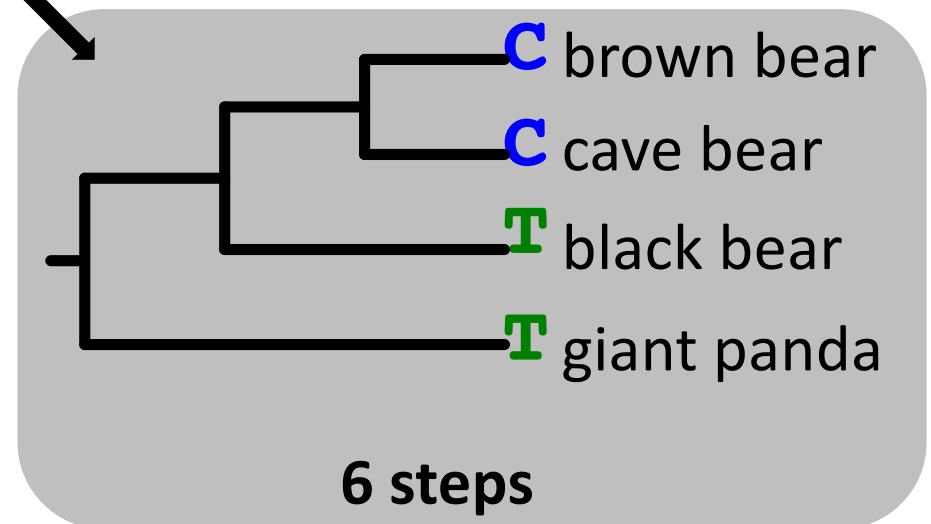
brown bear **C**G**T**T**A**G**T****A**C**A**C**T**
cave bear **C**G**A**T**A**G**T****T**C**A**C**T**
black bear **C**G**T**T**A**G**T****T**A**C**C
giant panda **C**A**T**T**G**G**T**T**A**C**T**



7 steps



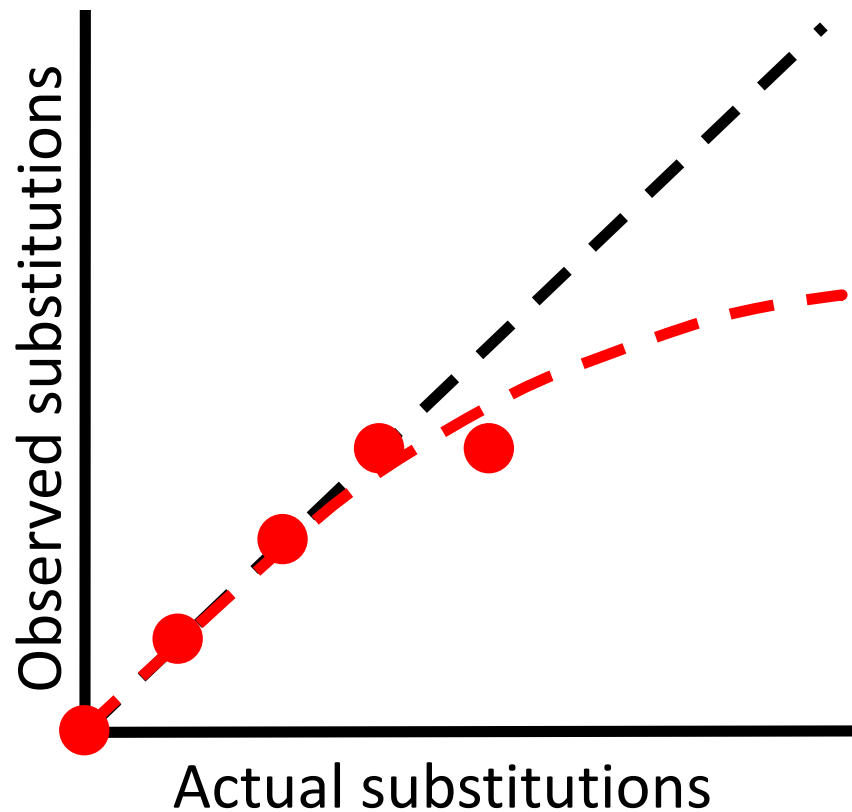
7 steps



6 steps

Maximum parsimony

- Identifies the tree topology that can explain the sequence data, using the smallest number of inferred substitution events
- Commonly used for morphological data
- Now *rarely used* for analysing genetic data
 - Cannot estimate evolutionary rates or timescales
 - Effects of multiple substitutions



A	A	A	A	A
A	T	T	T	T
C	C	G	G	G
A	A	A	A	A
T	T	T	T	T
T	T	T	T	T
A	A	A	A	A
G	G	G	G	G
T	T	T	A	C

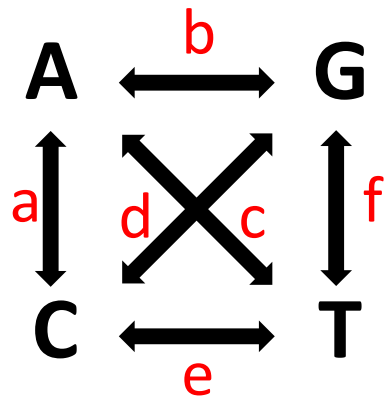
- Maximum parsimony does not correct for multiple substitutions at the same site
- This leads to a problem known as **long-branch attraction**
 - Long branch = many substitutions
 - Similarities arise by chance
 - Long branches cluster together

We can correct for multiple hits using substitution models

Substitution Models

Nucleotide substitution models

Rate Matrix



Base Frequencies

$$\pi_A + \pi_C + \pi_G + \pi_T = 1$$

JC

$$a=b=c=d=e=f$$

$$\pi_A = \pi_C = \pi_G = \pi_T$$

HKY

$$a=c=d=f, b=e$$

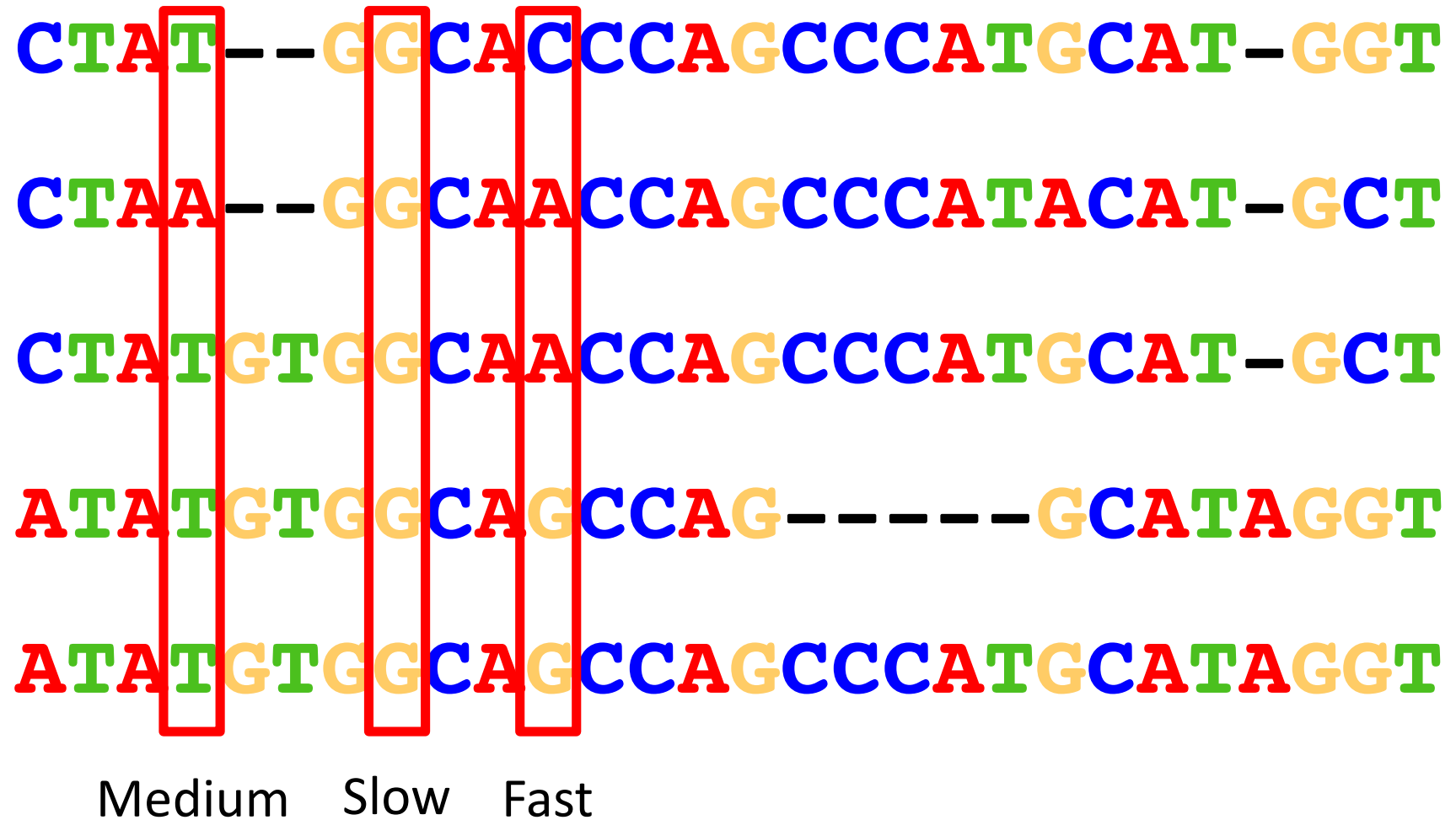
$$\pi_A, \pi_C, \pi_G, \pi_T$$

GTR

$$a, b, c, d, e, f$$

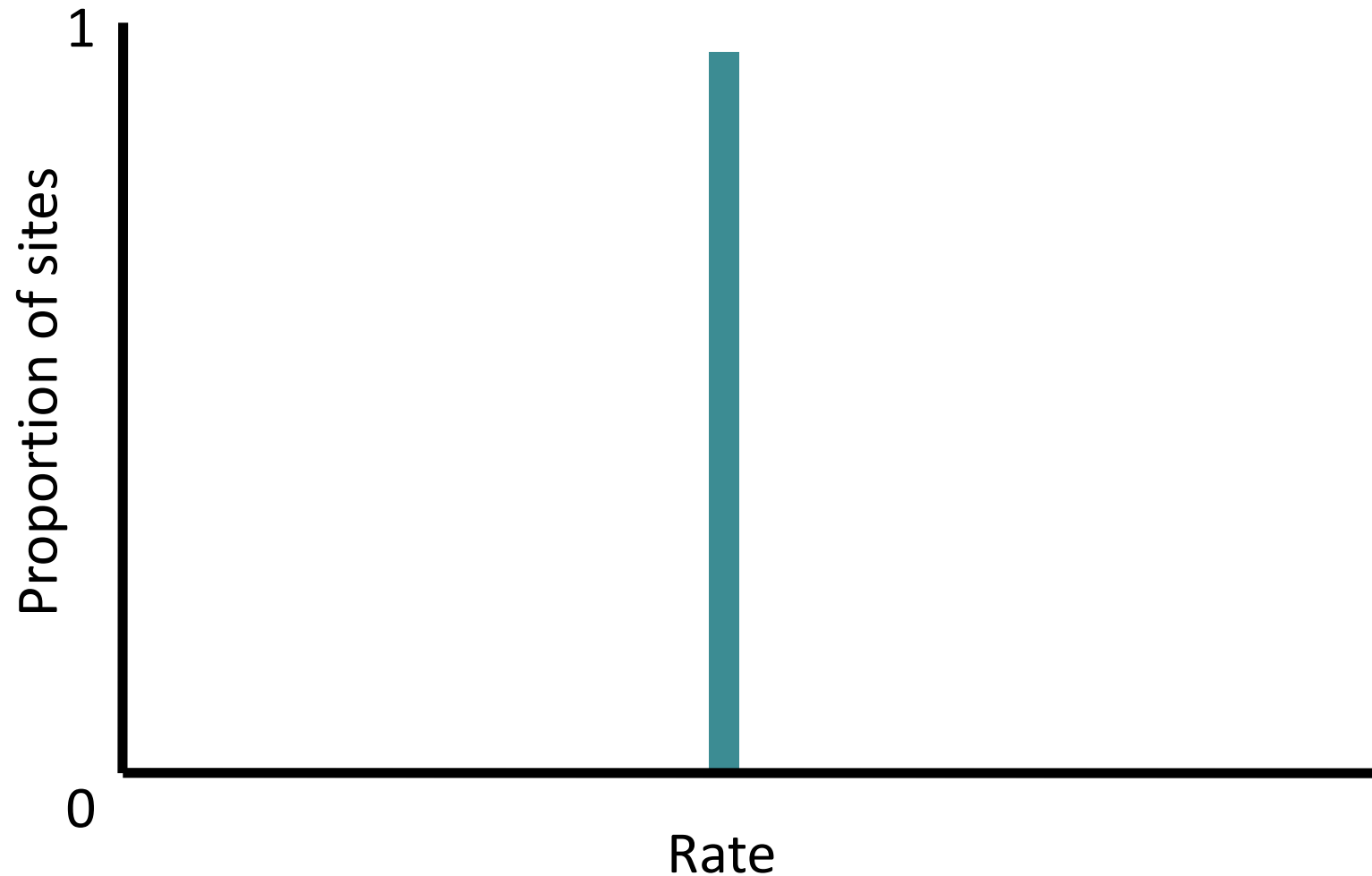
$$\pi_A, \pi_C, \pi_G, \pi_T$$

Rate variation across sites



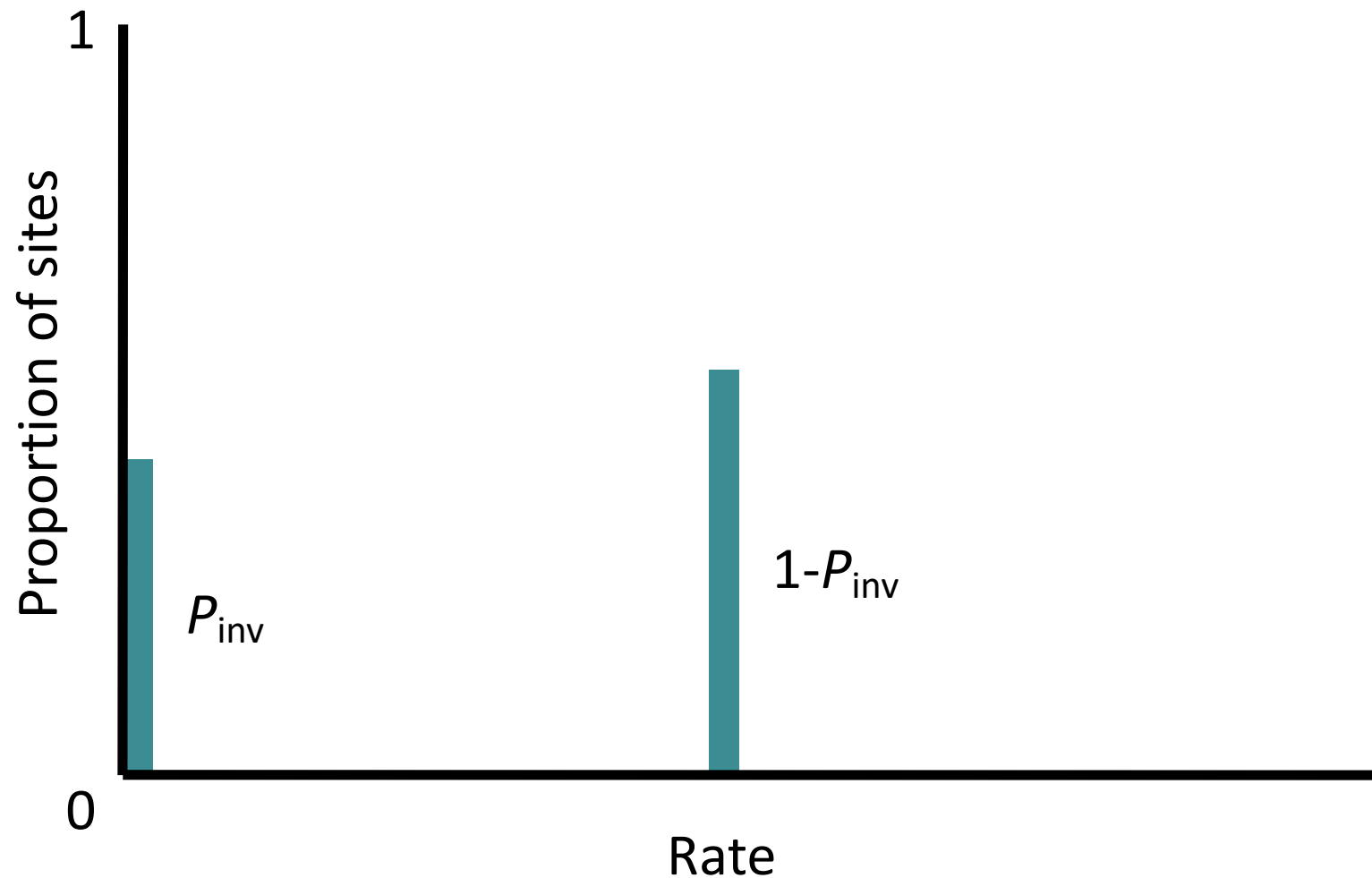
Rate variation across sites

- Equal rates among sites



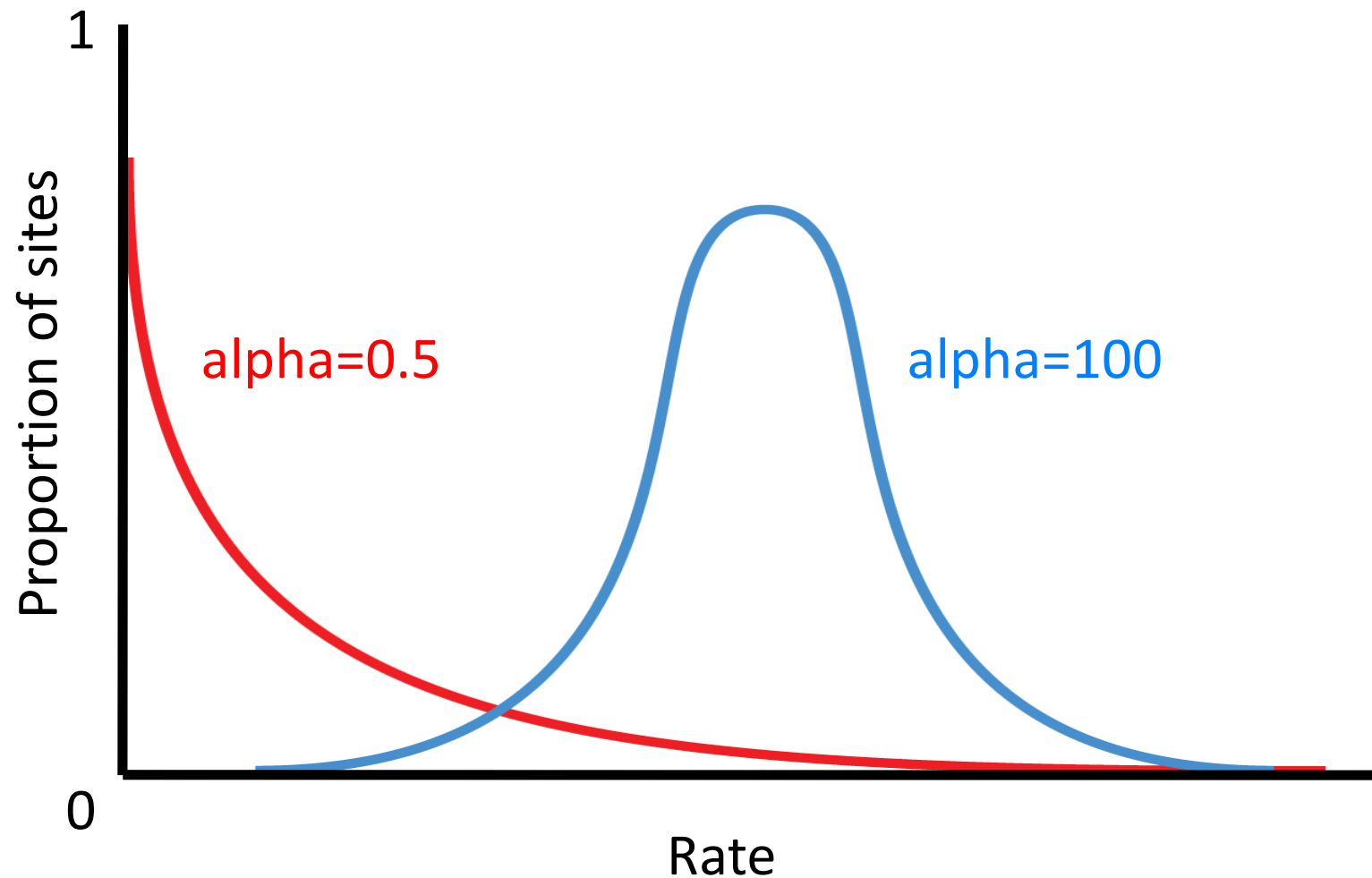
Rate variation across sites

- Proportion of invariable sites (+I models)



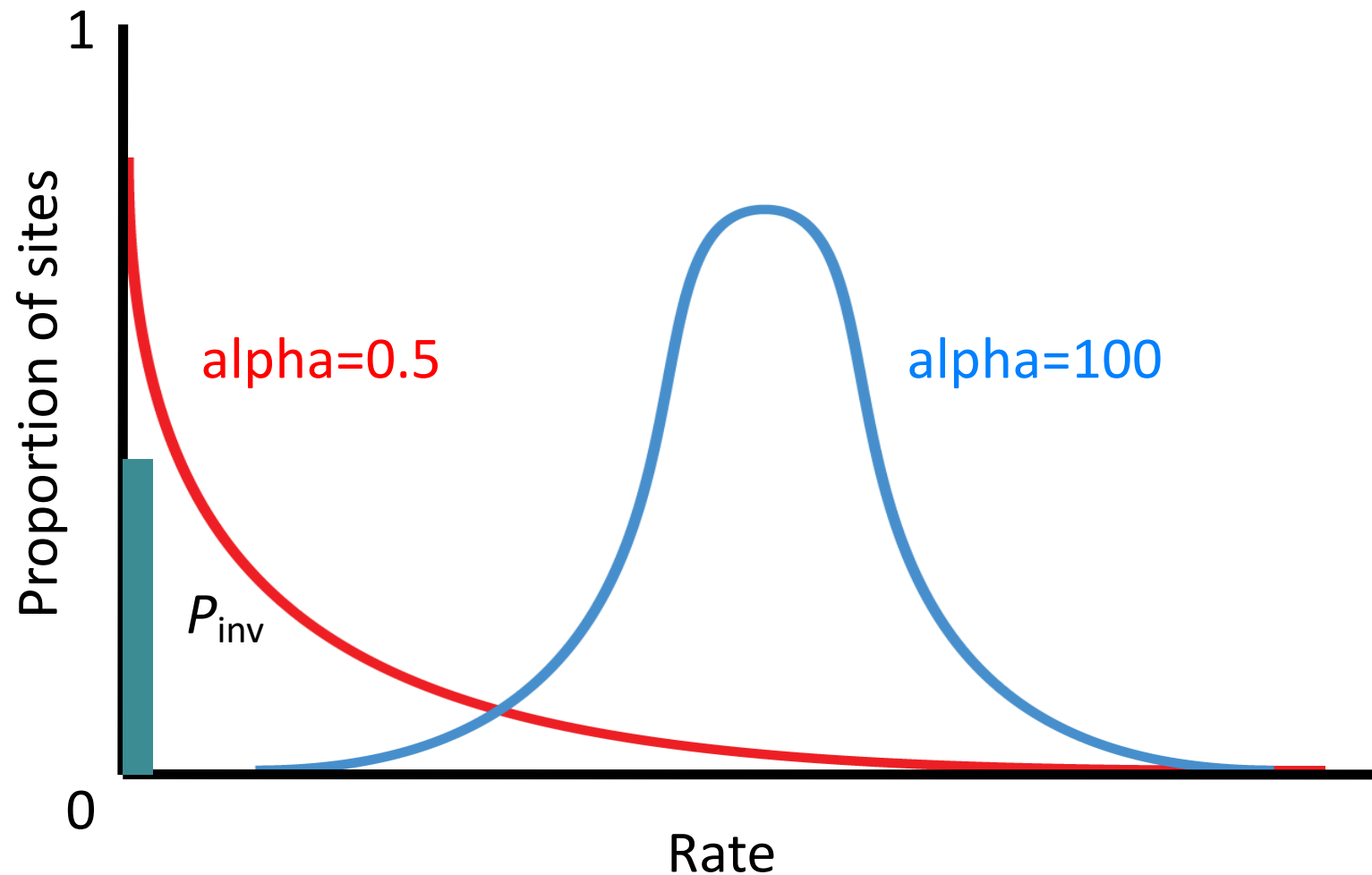
Rate variation across sites

- Gamma-distributed rate variation across sites (**+G** models)



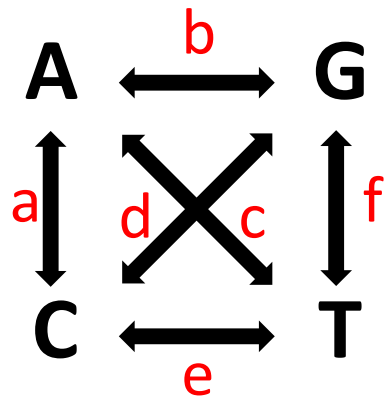
Rate variation among sites

- Gamma-distributed rate variation across sites and a proportion of invariable sites (**+G+I** models)



Nucleotide substitution models

Rate Matrix



Base Frequencies

$$\pi_A + \pi_C + \pi_G + \pi_T = 1$$

Site Rates

$$+ I + G$$

JC

$$a=b=c=d=e=f$$

$$\pi_A = \pi_C = \pi_G = \pi_T$$

HKY

$$a=c=d=f, b=e$$

$$\pi_A, \pi_C, \pi_G, \pi_T$$

GTR

$$a, b, c, d, e, f$$

$$\pi_A, \pi_C, \pi_G, \pi_T$$

GTR+I+G

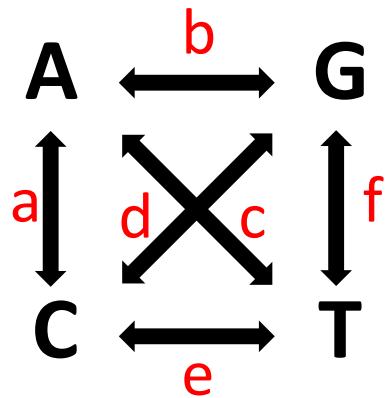
$$a, b, c, d, e, f$$

$$\pi_A, \pi_C, \pi_G, \pi_T$$

$$I, G$$

Nucleotide substitution models

Rate Matrix



Base Frequencies

$$\pi_A + \pi_C + \pi_G + \pi_T = 1$$

Site Rates

$$+ I + G$$

#Models

$$203 \quad \times \quad 15 \quad \times \quad 4 \quad = \quad 12,180$$

In phylogenetics, we typically consider a small subset of these

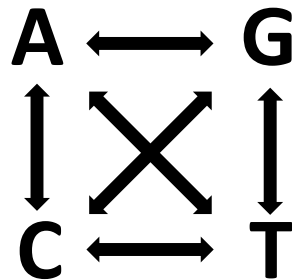
Proportion of invariable sites

- Often overestimated in analyses of intraspecific data
- Unable to distinguish between:
 - Sites that are **invariable** and unable to change
 - Sites that are **constant** and by chance have not mutated
- Not biologically meaningful
- Slowly evolving sites taken into account by **+G**

Use +G models to account for rate variation across sites

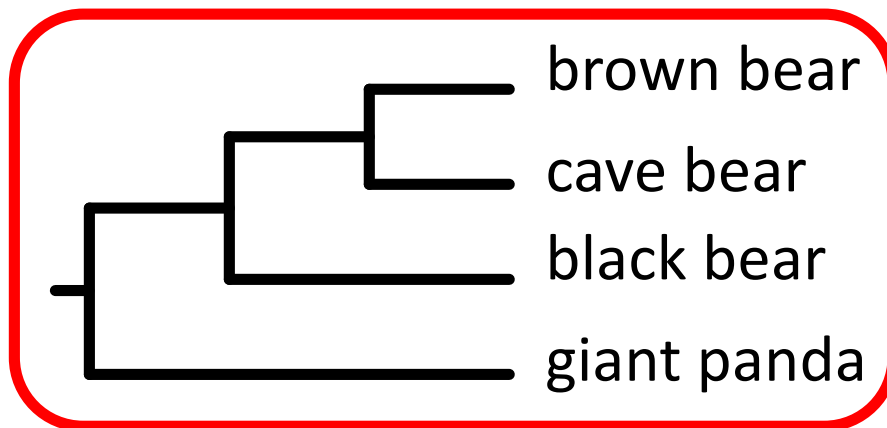
Fundamental assumptions

Reversible



Stationary

$$\pi_A + \pi_C + \pi_G + \pi_T = 1$$



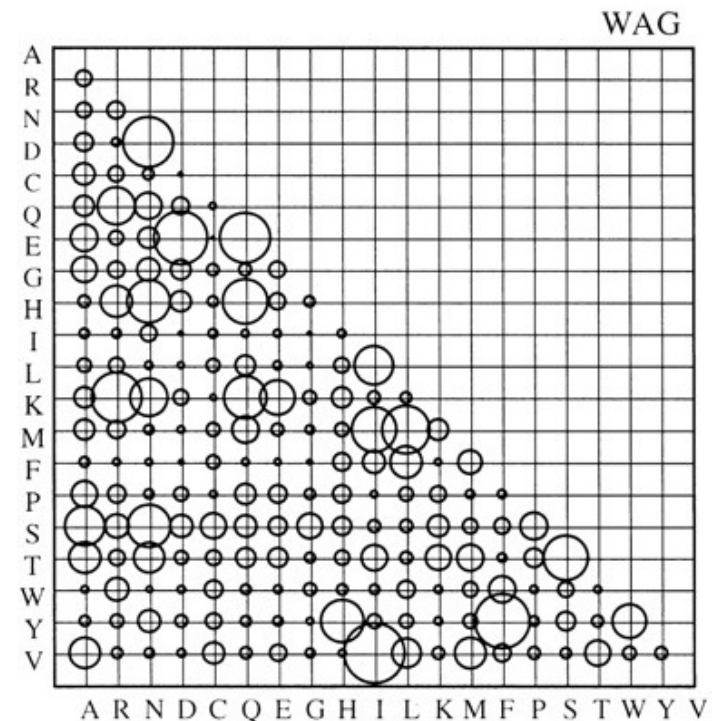
CGTTAGTACACT
CGATAGTTCACT
CGTTAGTTTACC
CATTGGTTTACT

Homogeneous

Independent across sites

Amino acid substitution matrices

- 20x20 matrix of substitution probabilities
- Too many parameters to estimate
 - GTR model for DNA: 6 parameters
 - GTR model for proteins: 190 parameters
- Estimate substitution probabilities using large data set
 - PAM
 - BLOSUM
 - JTT
 - WAG



Model Selection

Model selection

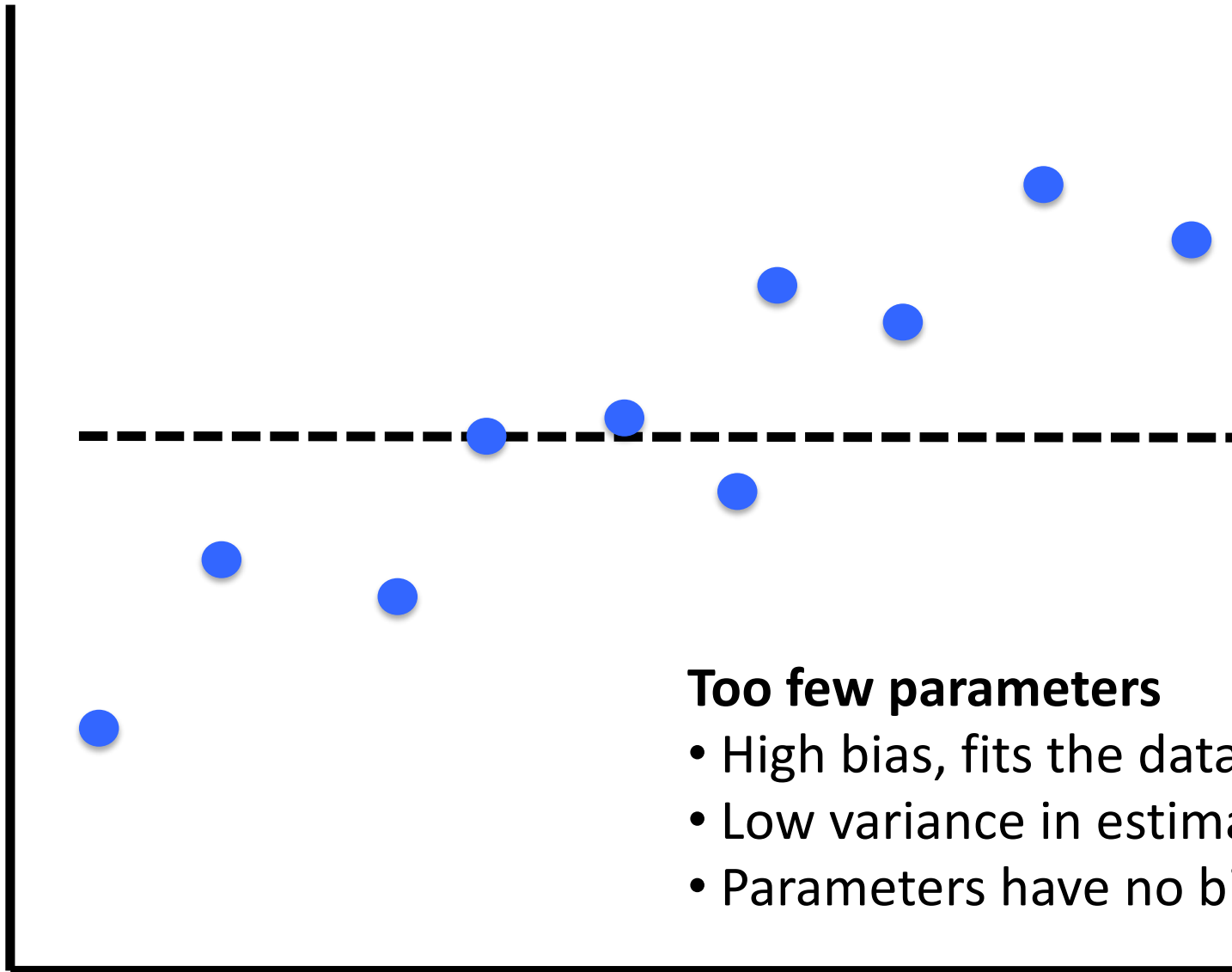
1. **Subjective model selection**

- Pick a model that seems sensible
- Balance the number of parameters against the amount of data
- Biological motivation

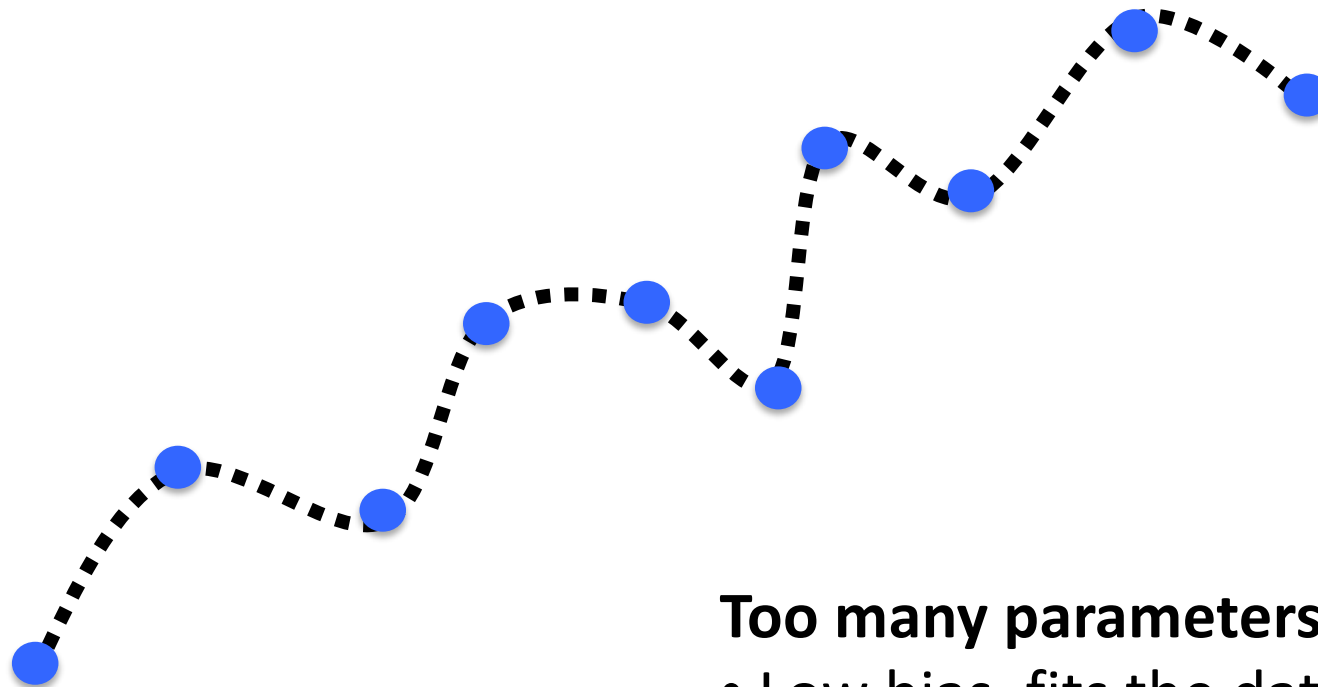
2. **Objective model selection**

- Use information theory and let a computer do it for you
- Statistical motivation

Model selection



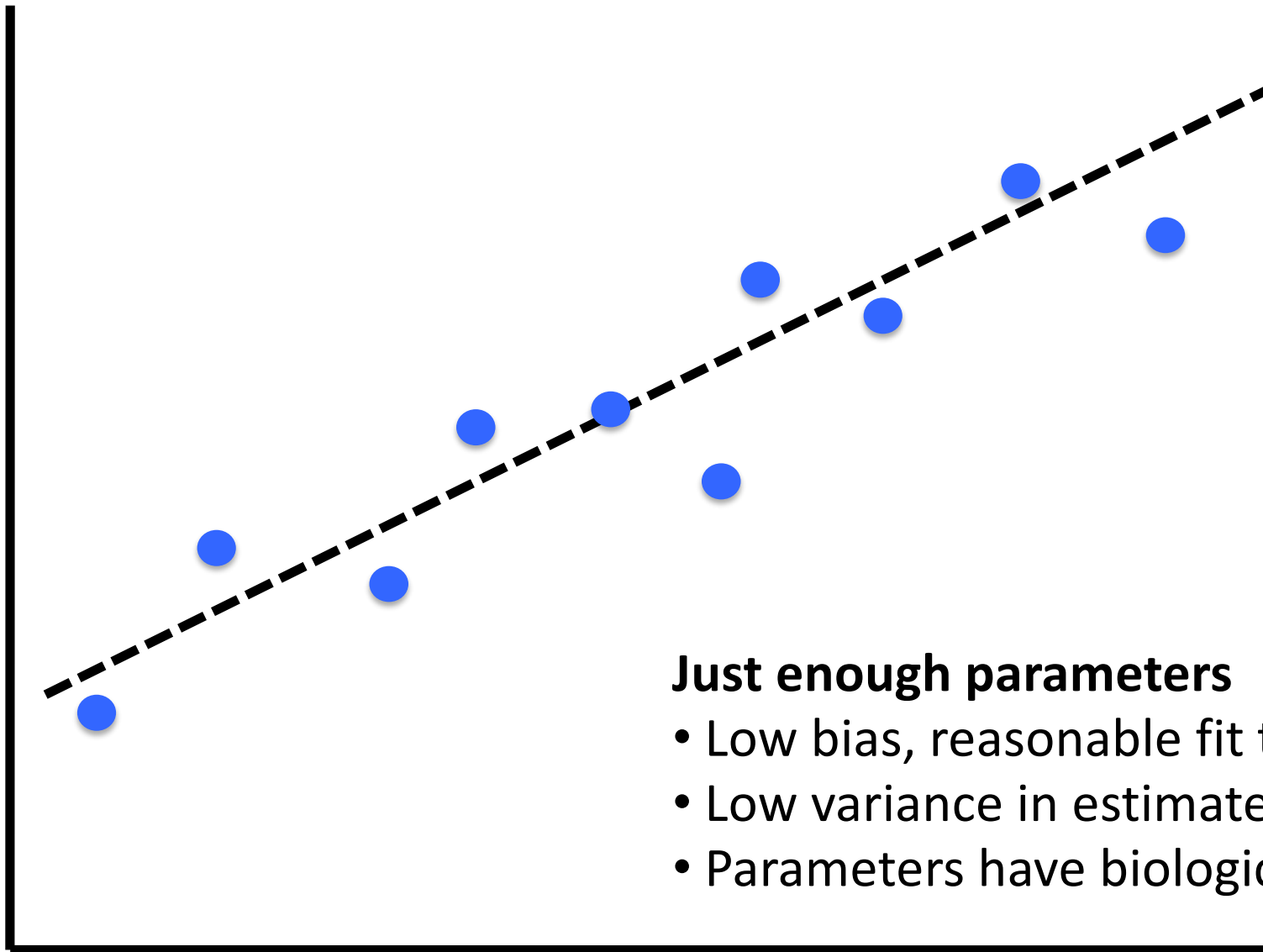
Model selection



Too many parameters

- Low bias, fits the data very well
- High variance in estimates
- Parameters have little biological meaning

Model selection



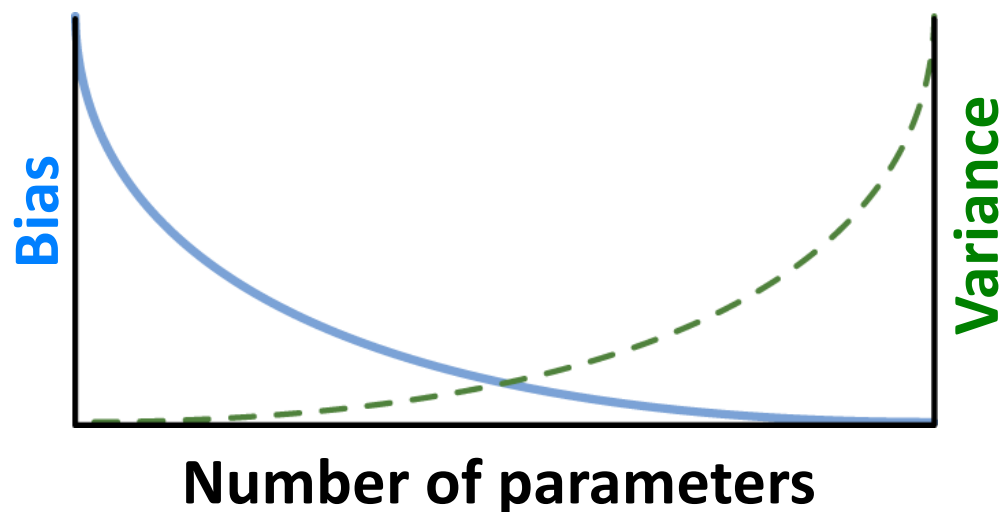
Just enough parameters

- Low bias, reasonable fit to data
- Low variance in estimates
- Parameters have biological meaning

Model selection

- Adding more parameters *always* improves the fit of the model to the observed data
- But more parameters leads to greater variance in the estimates of those parameters

Is the improvement in likelihood worth the cost of adding a parameter?

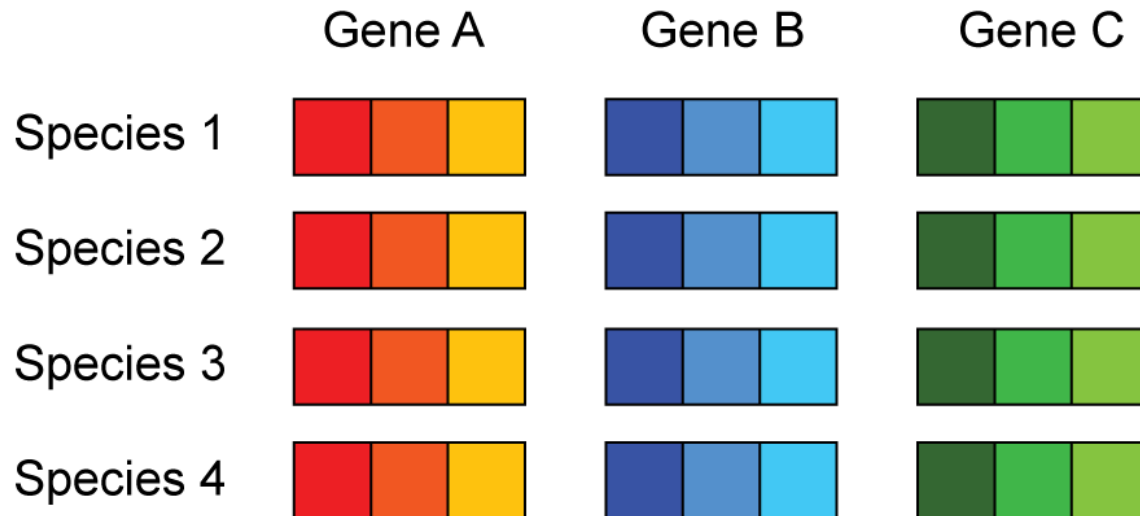


Model selection

- **Likelihood-ratio test (LRT)**
Used to compare nested models
- **Akaike information criterion (AIC)**
 $AIC = -2\ln(\text{likelihood}) + 2k$
- **Bayesian information criterion (BIC)**
 $BIC = -2\ln(\text{likelihood}) + k\ln(n)$

Data partitioning

- Separate substitution model for each gene and codon position?



- **Biological**

- Genome
- Genes
- Codon positions
- RNA stems vs loops
- Hydrophobic vs hydrophilic

- **Statistical**

PartitionFinder

- Too many possible partitioning schemes
 - 15 schemes for 4 genes
 - 52 schemes for 5 genes
 - 203 schemes for 6 genes

PartitionFinder: Combined Selection of Partitioning Schemes and Substitution Models for Phylogenetic Analyses

Robert Lanfear,^{*,1} Brett Calcott,^{1,2} Simon Y. W. Ho,³ and Stephane Guindon⁴

2012 – *Molecular Biology and Evolution*, 29: 1695–1701.

Substitution models in practice

- Phylogenetic estimates are usually robust to choice of model
- **GTR+G** is fine for most data sets
- Sensible data partitioning (*e.g.*, by codon position)

Useful references

- **Model selection in phylogenetics**
Sullivan & Joyce (2005) *Annual Review of Ecology, Evolution, and Systematics*, 36: 445–466.
- **The effects of partitioning on phylogenetic inference**
Kainer & Lanfear (2015) *Molecular Biology and Evolution*, 32: 1611–1627.

