

# **Bristol Bioinformatics Workshop 2018**

**Tom Williams, Gary Barker, Tom Batstone**

**5th-9th November, 2018**

# Course aims

- Introduction to **a general approach** to doing bioinformatics
- **Core concepts** and signposts to the literature
- **Practical experience** to get your hands dirty

# Course structure

- **Day 1 (today)**
  - Introduction
  - Comparative genomics (theory)
  - Comparative genomics (GUI/web practical)
- **Day 2:**
  - Introduction to UNIX
  - Comparative genomics (practical)
  - Phylogenetics (theory)
- **Day 3:**
  - Phylogenetics (practical)

# Today's schedule

<b>Time</b>	<b>Material</b>	<b>Lecturer</b>
10.00-10.30	Course introduction and overview	Tom W.
10.30-13.00	DNA sequencing, genome assembly and comparative genomics	Gary B.
14.00-17.00	Introduction to phylogenetics	Tom W.

# Protip

Participation is key to getting the most out of this kind of course.

Please ask (and answer) questions!

# The basic idea

Bioinformatics is an experimental science  
Bioinformatics is an experimental science



# Doing science: a workflow



Collect, collect, collect, collect, collect, collect, collect, analyze, Nature

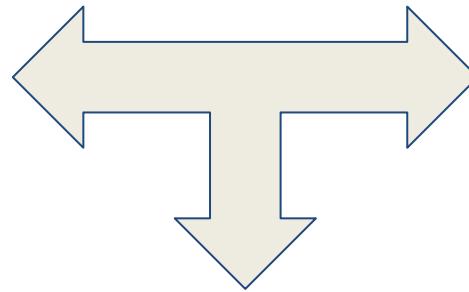
We use computers to explore data or test hypotheses



=



# Doing science in an integrated way



**“An hour of #bioinformatics can save a week in the lab. ALSO An hour in the lab can save a week of #bioinformatics” - Torsten Seeman**

# Two kinds of questions frequently encountered when doing bioinformatics

## Scientific questions:

- Is my hypothesis testable, and with what data?
- What assumptions can I make when analyzing the data?
- Which of the available analysis methods is most appropriate?

**Is what I'm doing sensible?**



## Technical questions:

- What software/methods are available for problem X?
- How do I install this software on my computer?
- How do I fix this error message?
- How do I make it give me the answer?

**Why won't it \$%\$£"!# work?**



# Two kinds of questions frequently encountered when doing bioinformatics

## Scientific questions:

- Is my hypothesis testable, and with what data?
- What assumptions can I make when analyzing the data?
- Which of the available analysis methods is most appropriate?

**Is what I'm doing sensible?**



## Technical questions:

- What software/methods are available for problem X?
- How do I install this software on my computer?
- How do I fix this error message?
- How do I make it give me the answer?

**Why won't it \$%\$£"!# work?**



# Bioinformatics: in sum

- Computers provide powerful tools for analyzing data
- They aren't black boxes, into which you feed data and out of which the result comes
- They can speed up (not replace) our thinking, and make many types of analysis possible
- Be as critical as you would be in any other area of science

# **Introduction to molecular phylogenetics**

**(a) Core concepts**

**Tom Williams**

# Schedule for the phylogenetics material

Today: Core concepts in phylogenetics.  
Core concepts: games and activities.

Friday morning: A workflow for phylogenetics.  
Some theory.  
The molecular clock.

Friday afternoon: **Phylogenetics practical.**  
Principles of phylogenomics.  
Games and activities.

# Core concepts in molecular phylogenetics

- Some justification for doing phylogenetics
- Homology and analogy
- Interpreting phylogenetic trees
- Naming groups on trees
- The root: a very special node on the tree

Some quizzes.

# What is phylogenetics?

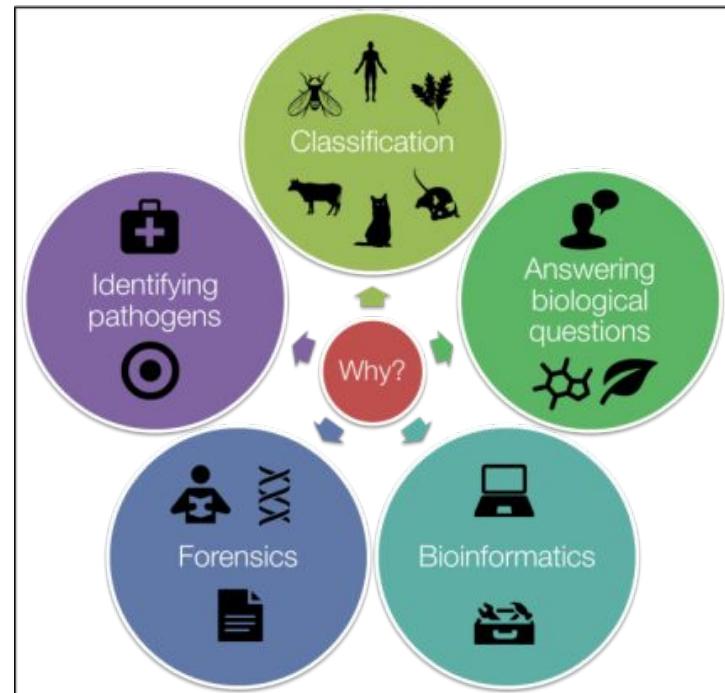
**The study of the evolutionary relationships among evolving entities (species, genomes, genes, individuals).**

“The time will come, though I shall not live to see it, when we shall have fairly true genealogical trees of each great kingdom of nature”

- Darwin (1857)

“One of the grand biological ideas is to be able to work out the complete detailed quantitative phylogenetic tree --- the history of the origins of all living species, back to the very beginning.”

- Dayhoff and Eck (1972)



(Various practical applications)

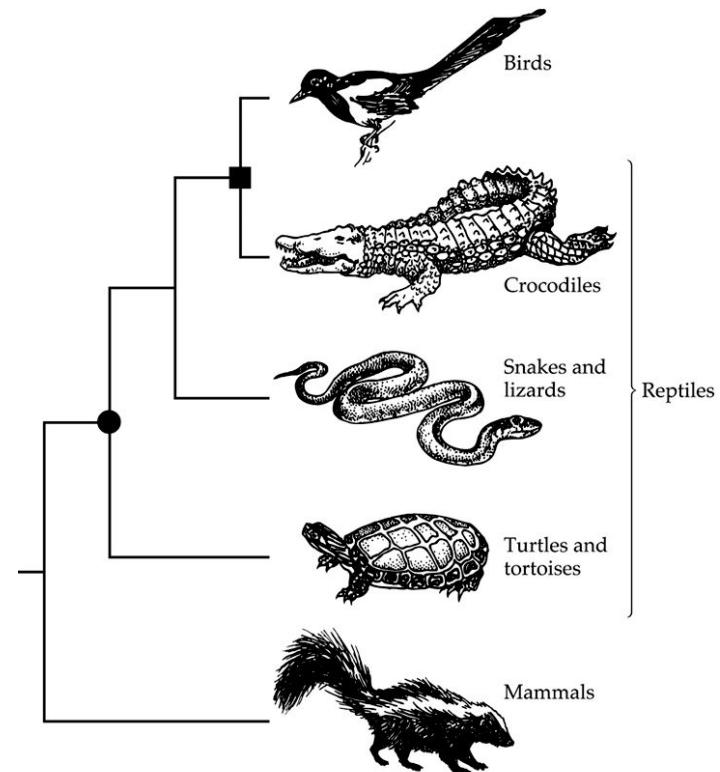
# What is a phylogeny (evolutionary tree)?

A branching diagram representing the genealogical relationships among species

Can be inferred using:

- Genetic (molecular) data
- Morphological data
- Any kind of character/trait data

Normally, we want to infer a tree (or trees) that **make sense of** the character data



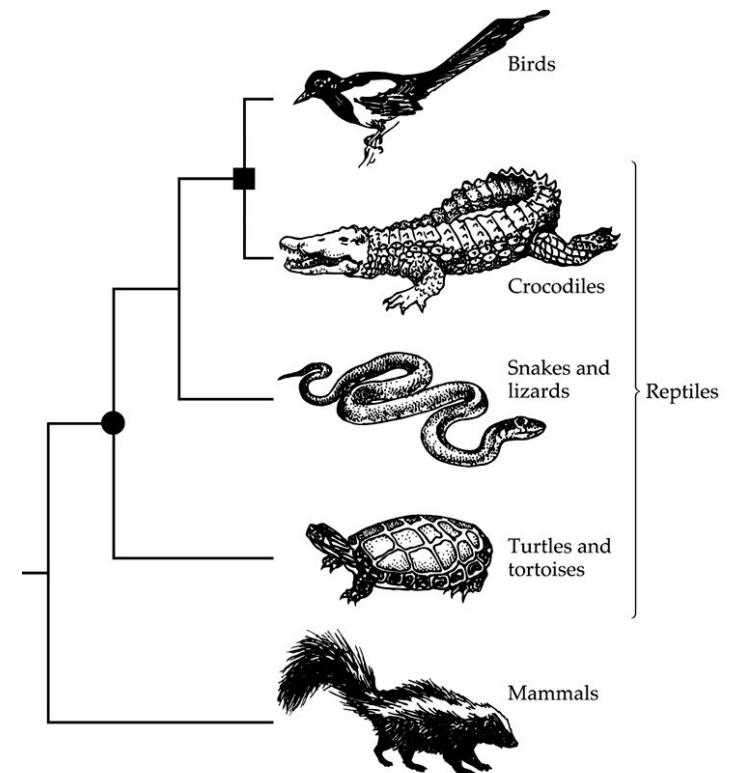
# 1. Homologous vs. analogous characters

- Homology is similarity due to common ancestry

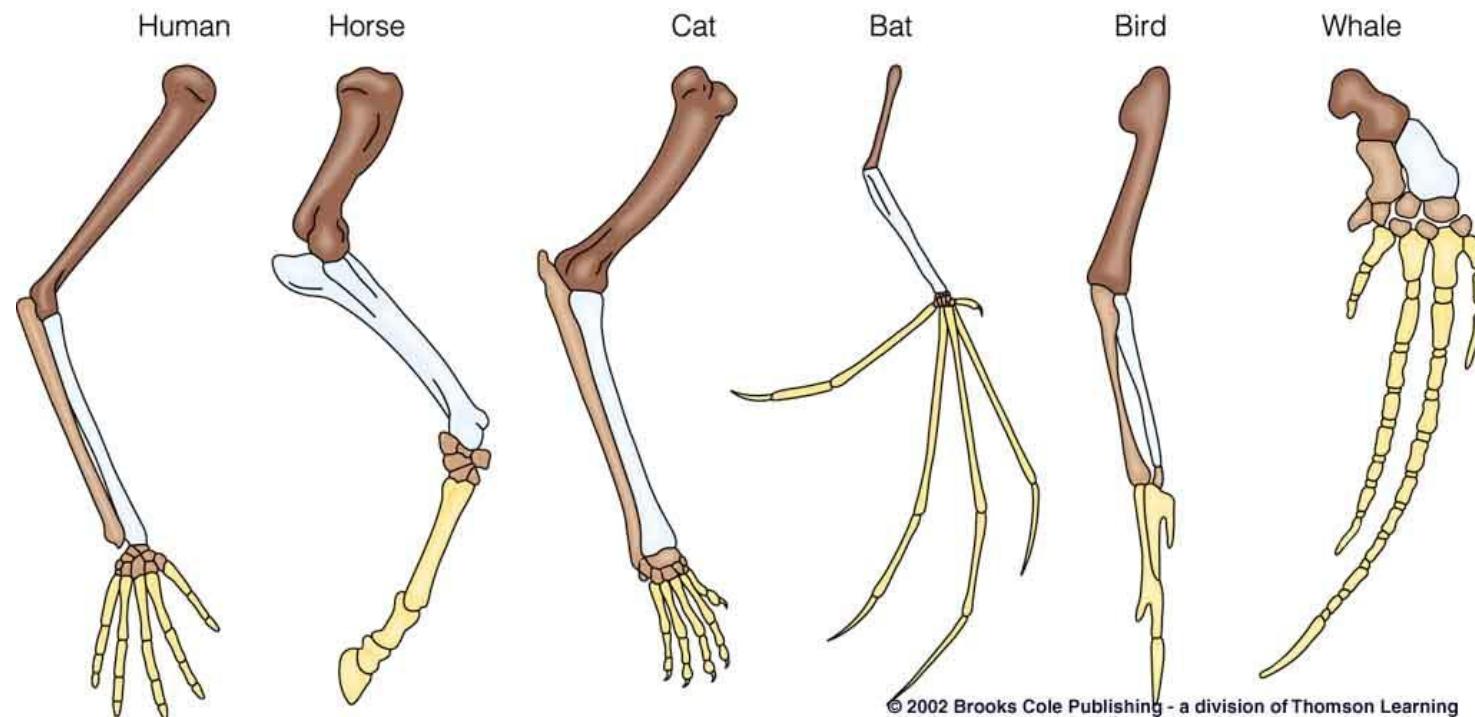
*“the same organ under every variety of form and function” (Owen)*

- Analogy is similarity due to chance or convergent evolution.

Superficial or misleading similarities.

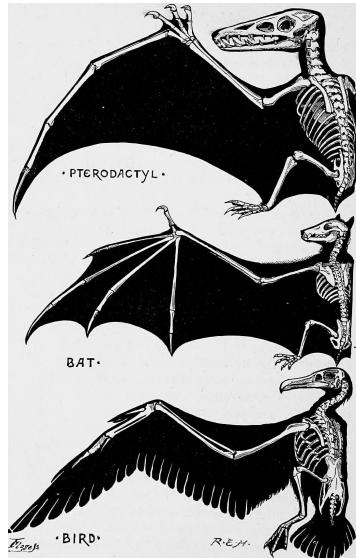


# Vertebrate limbs are homologous

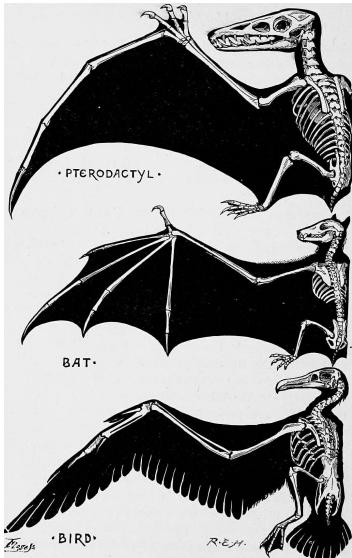


© 2002 Brooks Cole Publishing - a division of Thomson Learning

# Are these wings homologous?



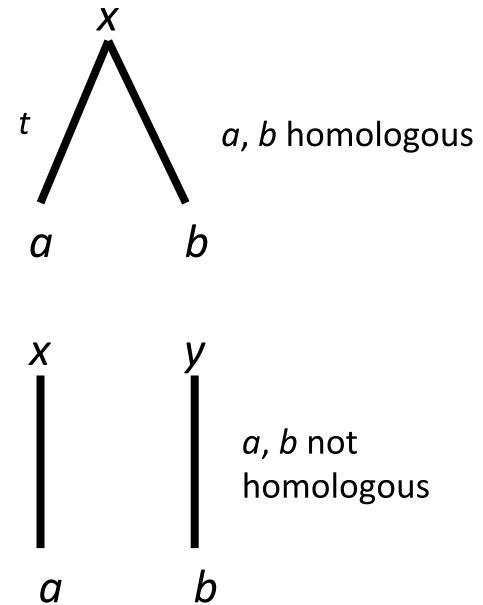
# Are these wings homologous?



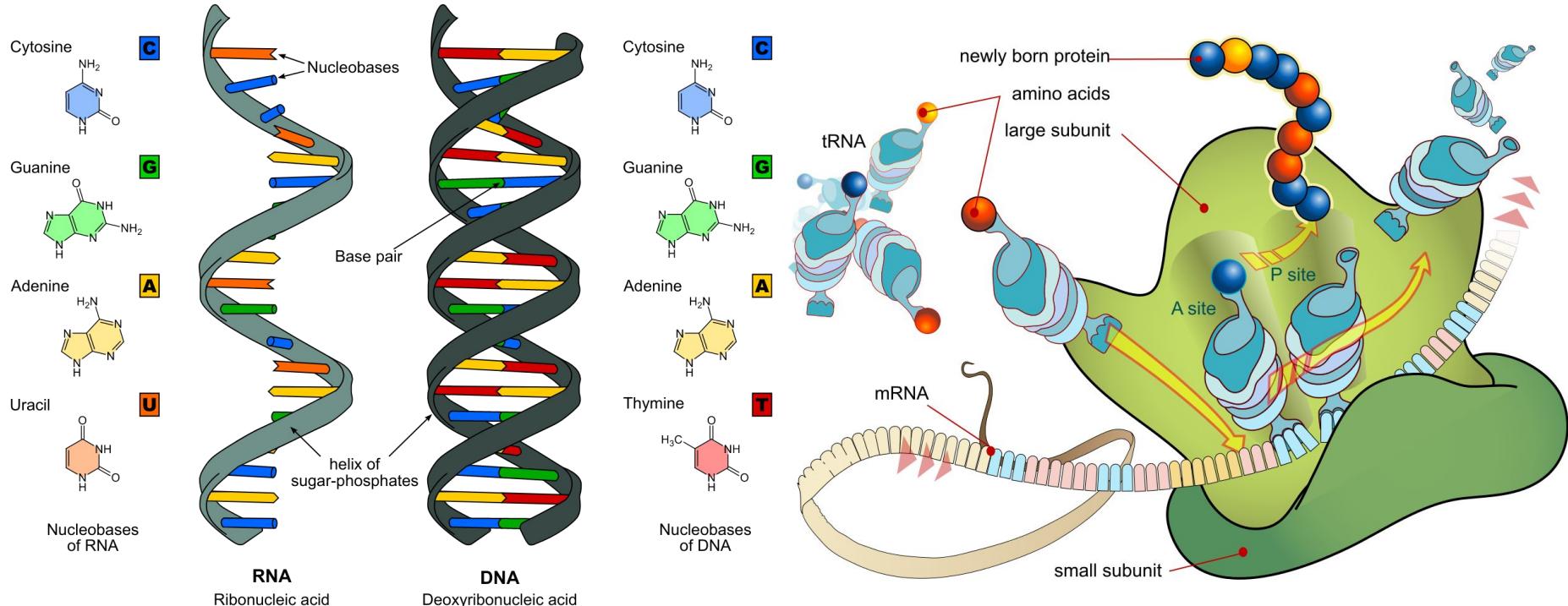
No: they do not descend from a common ancestral wing in the common ancestor of these animals.  
They evolved independently for the same reasons (flight).

# Homology

- Homologous characters can tell us about evolutionary relationships
- Homologous gene sequences: descend from a common ancestor
- A statement of homology is an **evolutionary hypothesis**.
- Homology is **either-or**: sequences either share a common ancestor or they don't.
- But, homology is often inferred from sequence similarity (e.g. BLAST searching).



# Molecular sequences provide digital information, many homologous characters



... — GTGCATCTGACTCCTGAGGGAGAAG ...    DNA  
... — CACGTAGACTGAGGAGTCCTCTTC ...



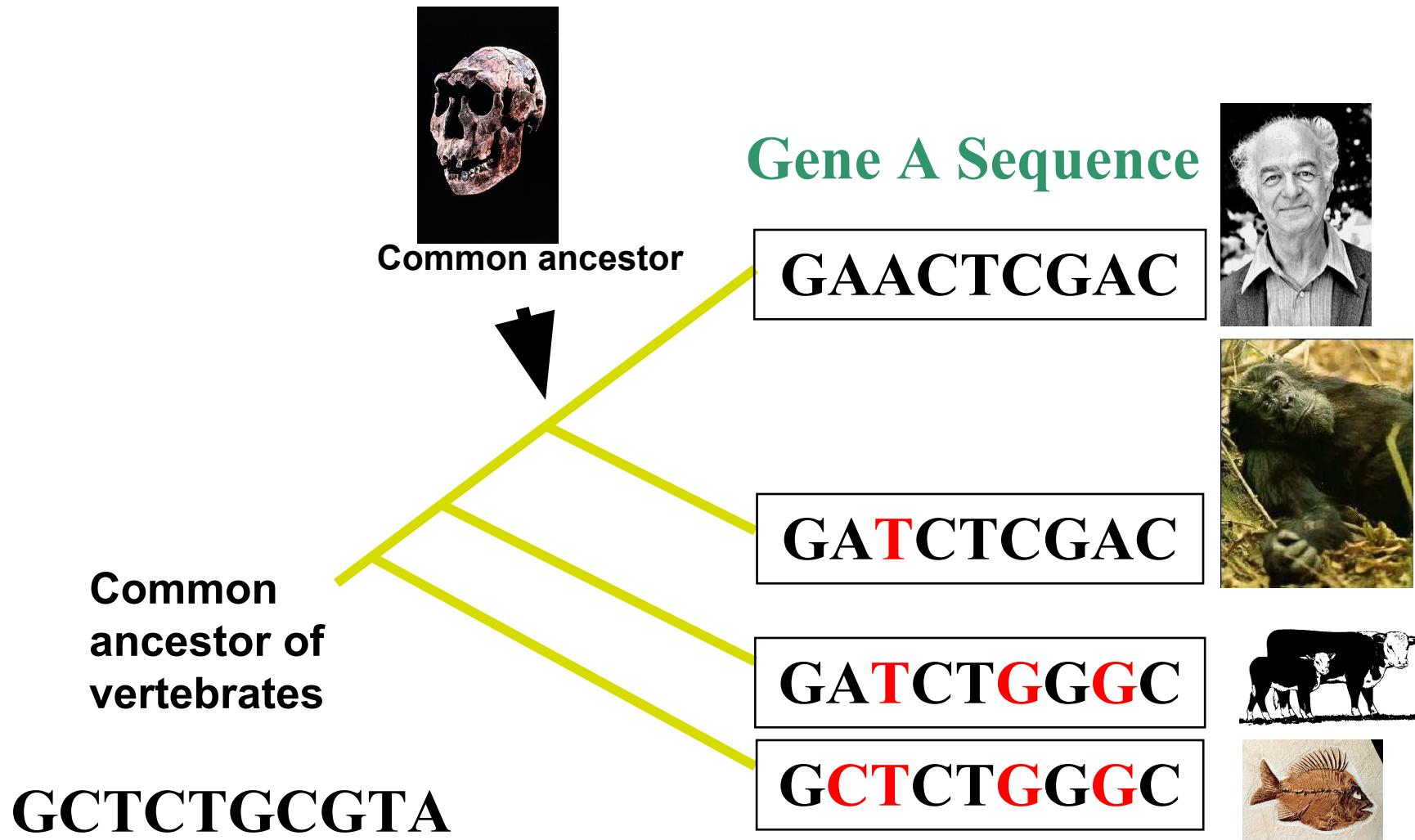
(transcription)

... — GUGCAUCUGACUCCUGAGGGAGAAG ...    RNA

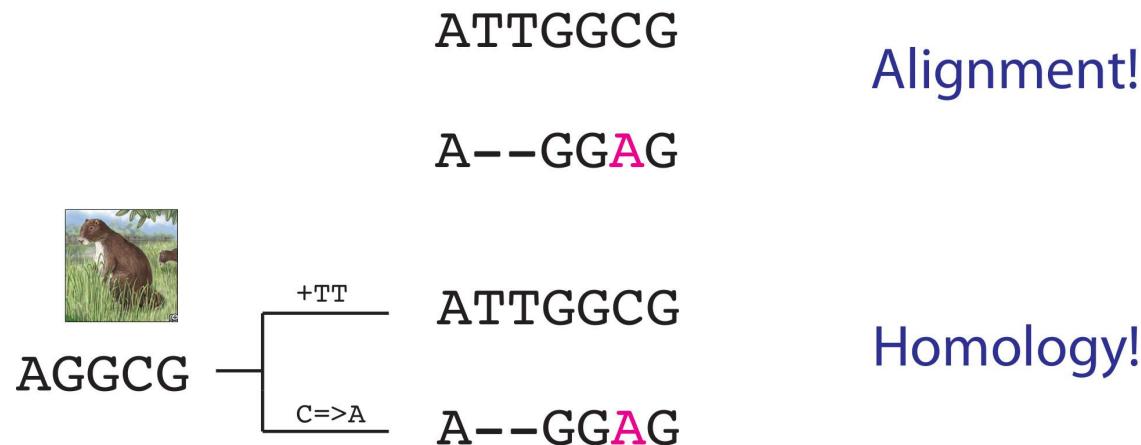
(translation)

... — V H L T P E E K ...    protein

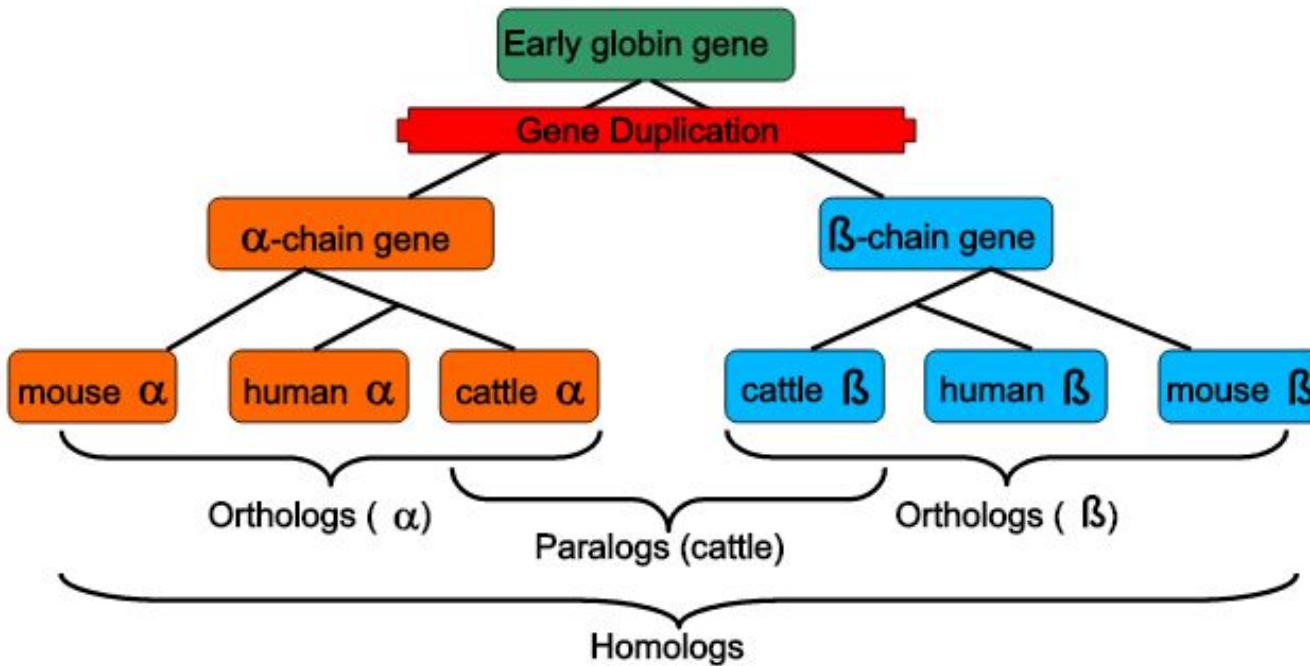
# DNA sequences can be used to make phylogenetic trees



# We use alignment to construct hypotheses of homology for sequence data



# Types of homology: orthology and paralogy



**Orthologues:** homologues related through a speciation event.

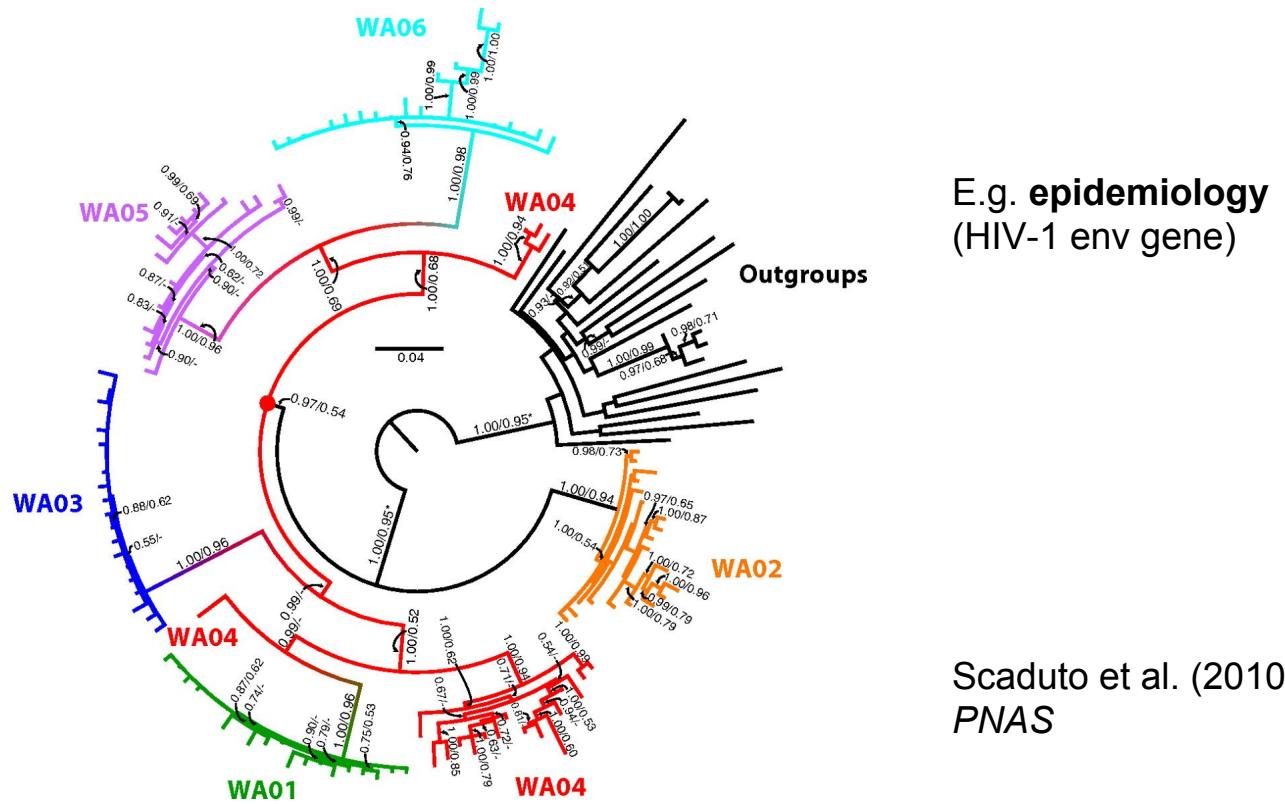
**Paralogues:** homologues related through a duplication event.

*(Don't mix them in your analyses!)*

## 2. Interpreting phylogenetic trees

What kind of information can we learn from phylogenies?

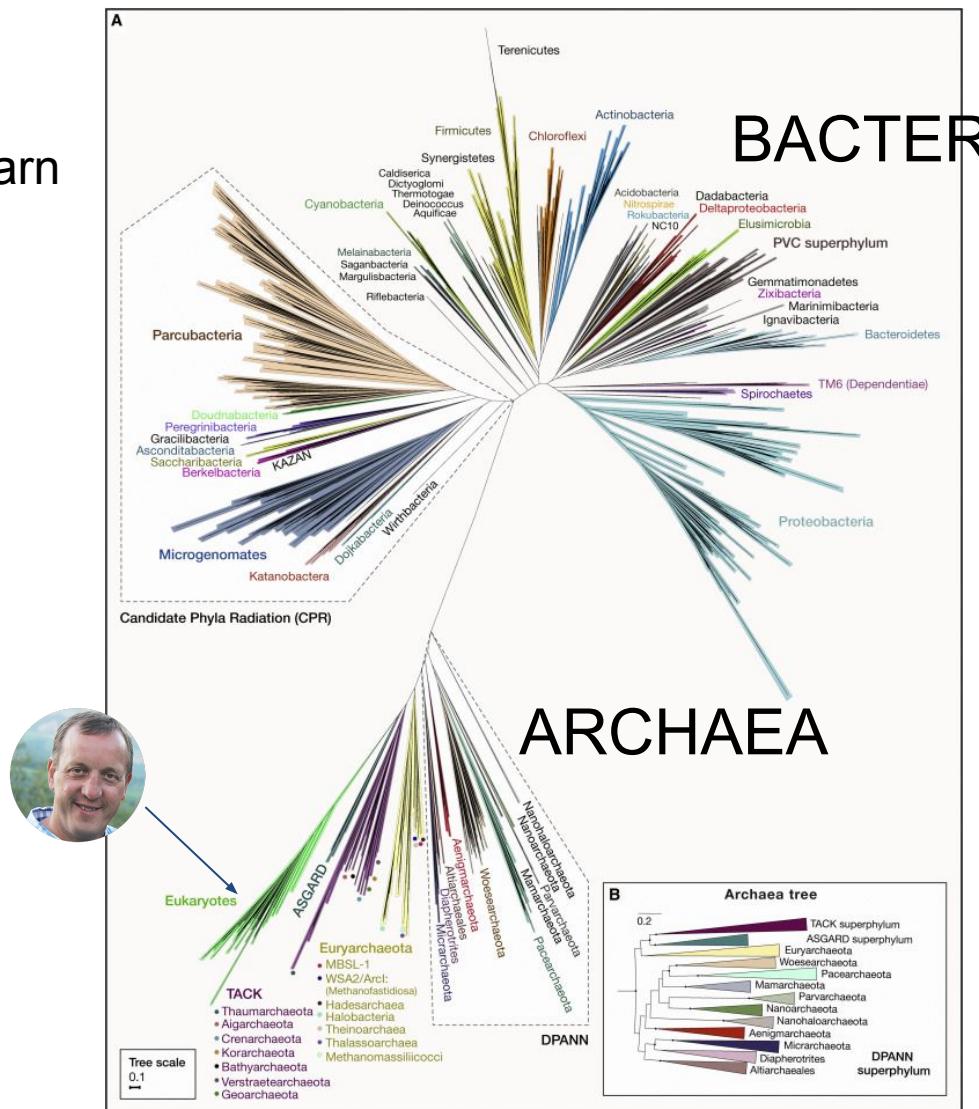
- Information about relationships (who is mostly closely related to whom?)



# Interpreting phylogenetic trees

What kind of information can we learn from phylogenies?

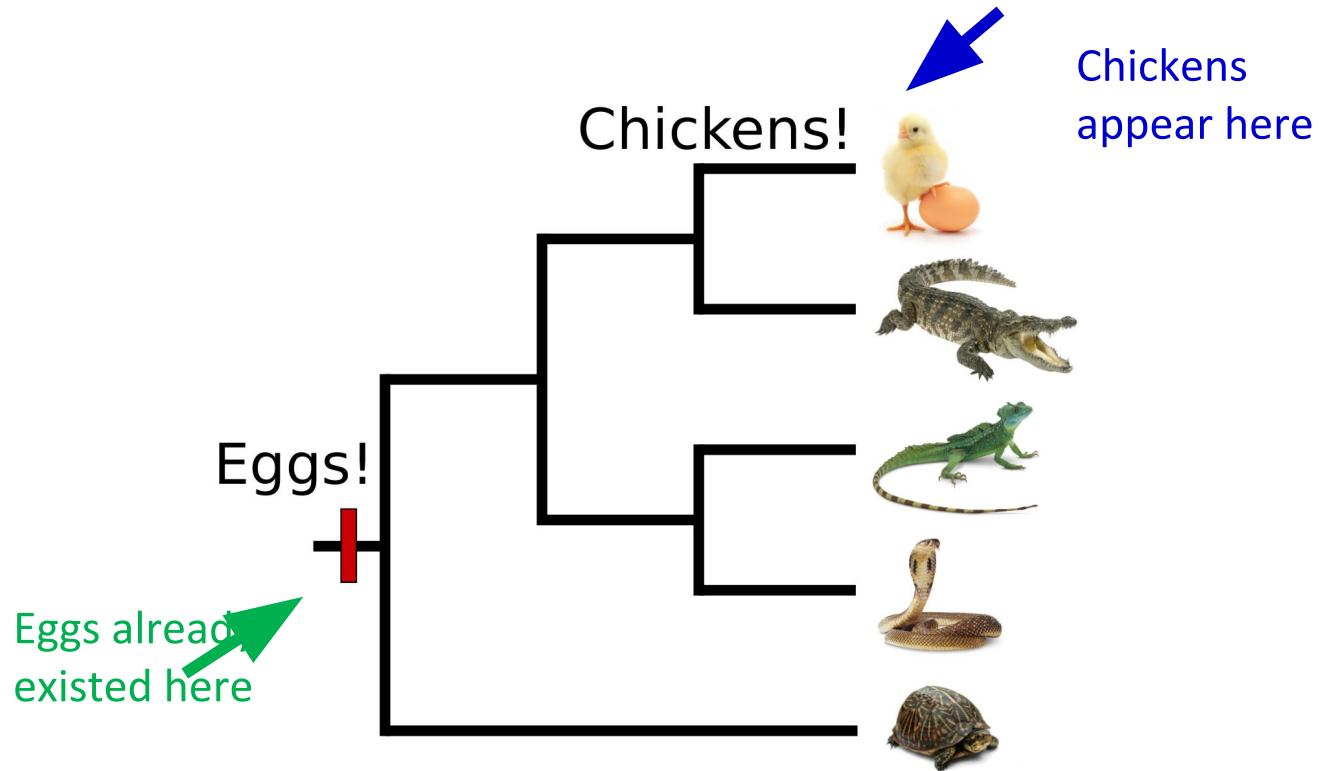
- Information about genetic (or other) diversity



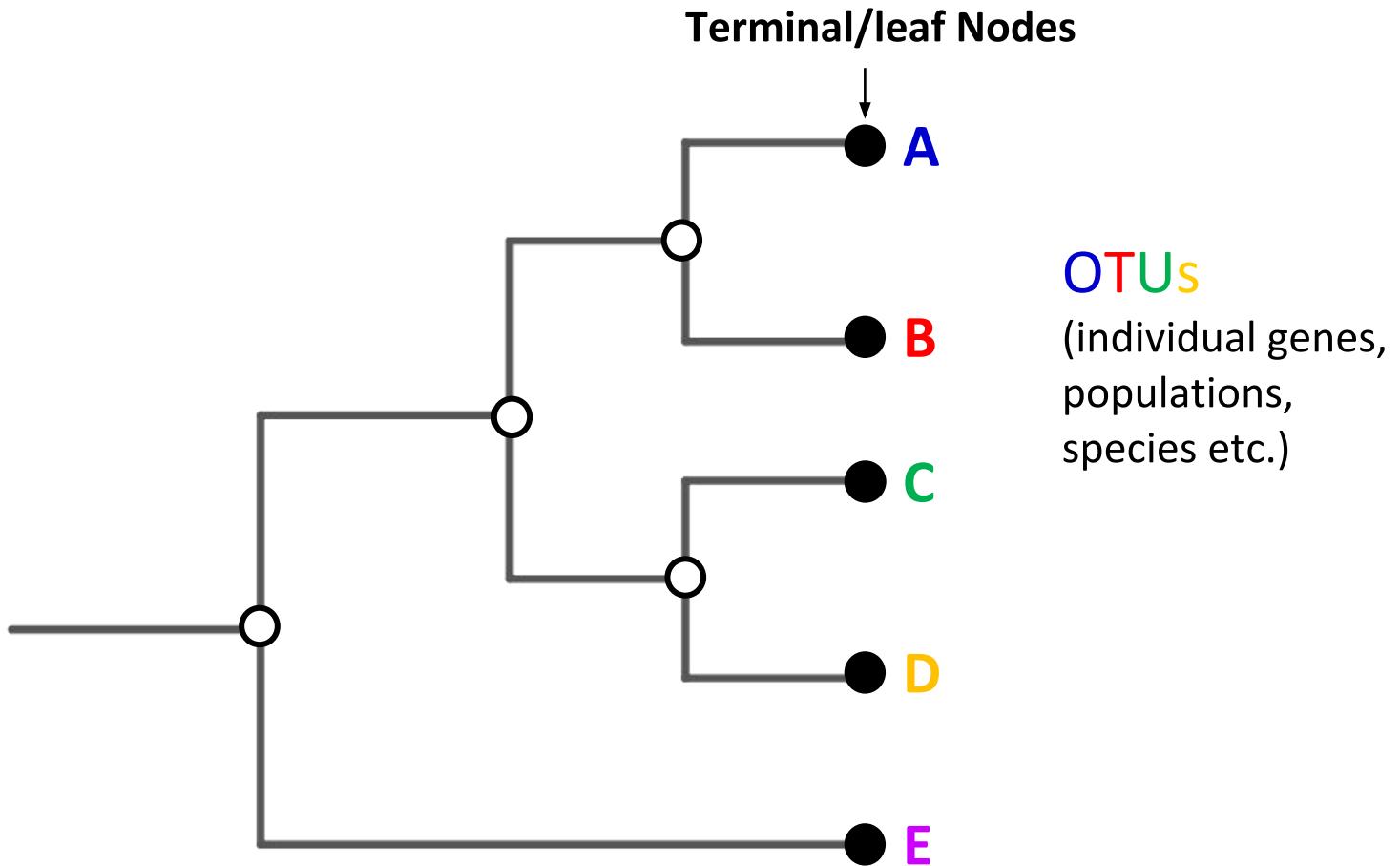
# Interpreting phylogenetic trees

What kind of information can we learn from phylogenies?

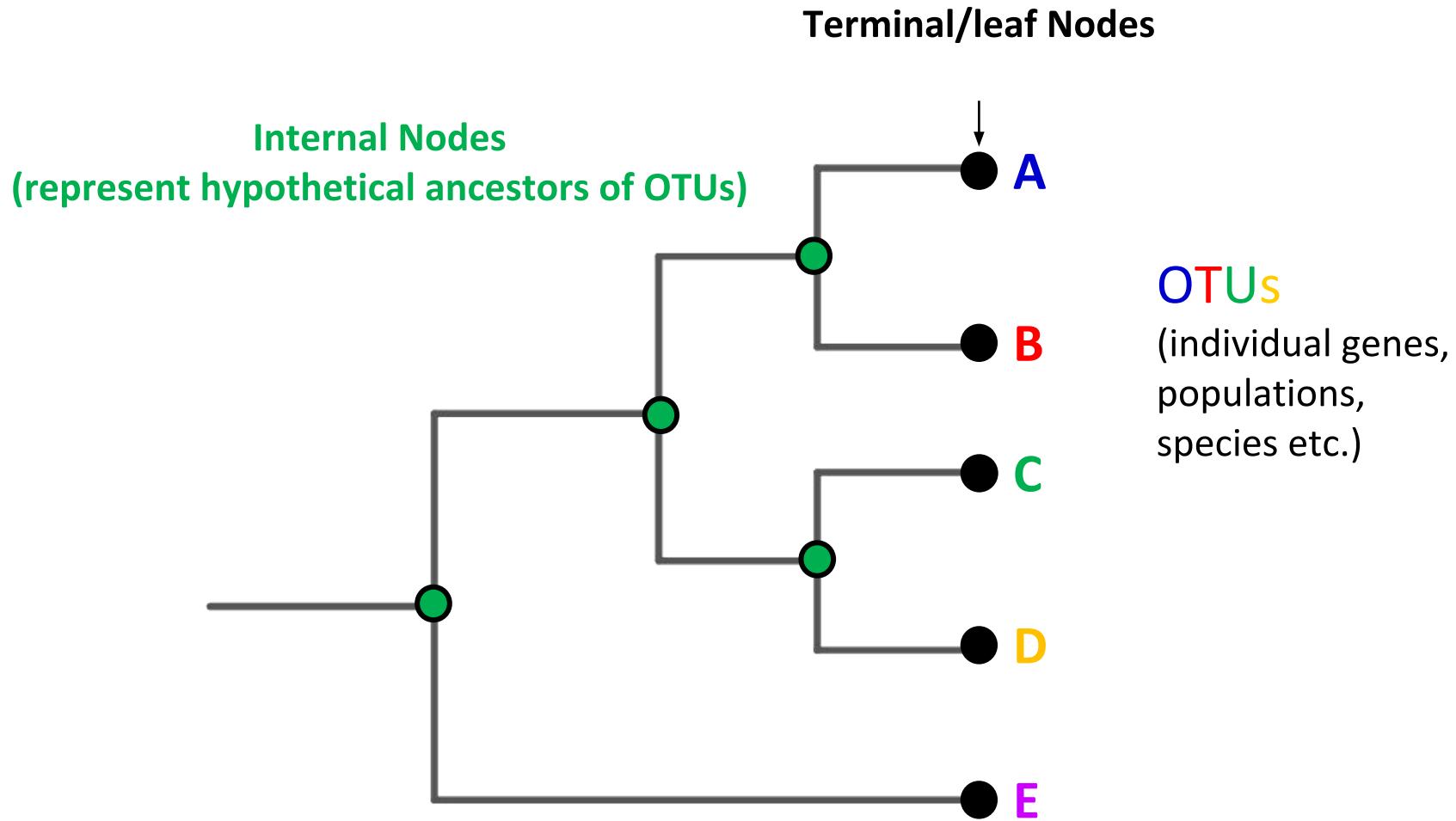
- The history of trait evolution



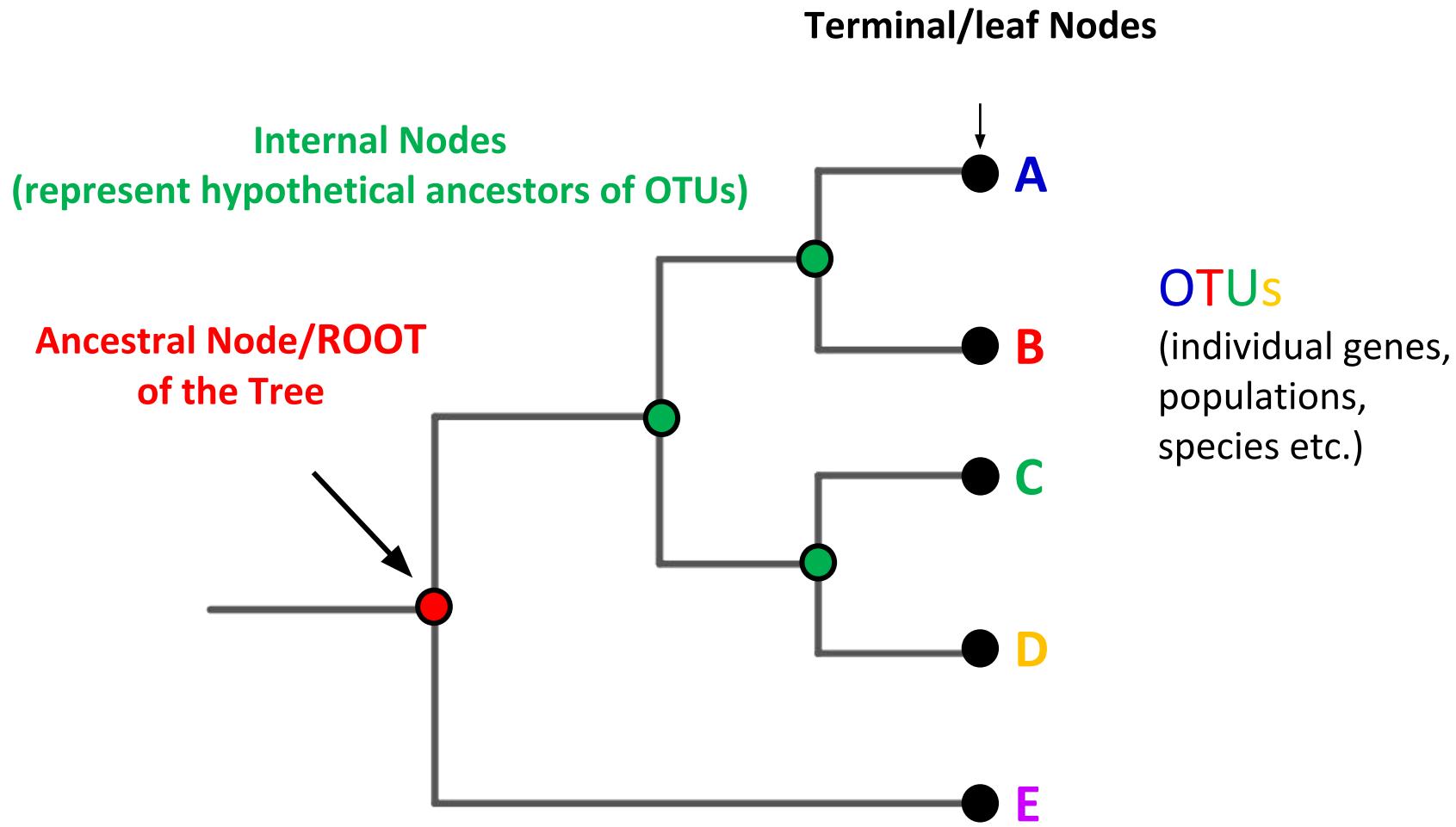
# Anatomy of a phylogenetic tree



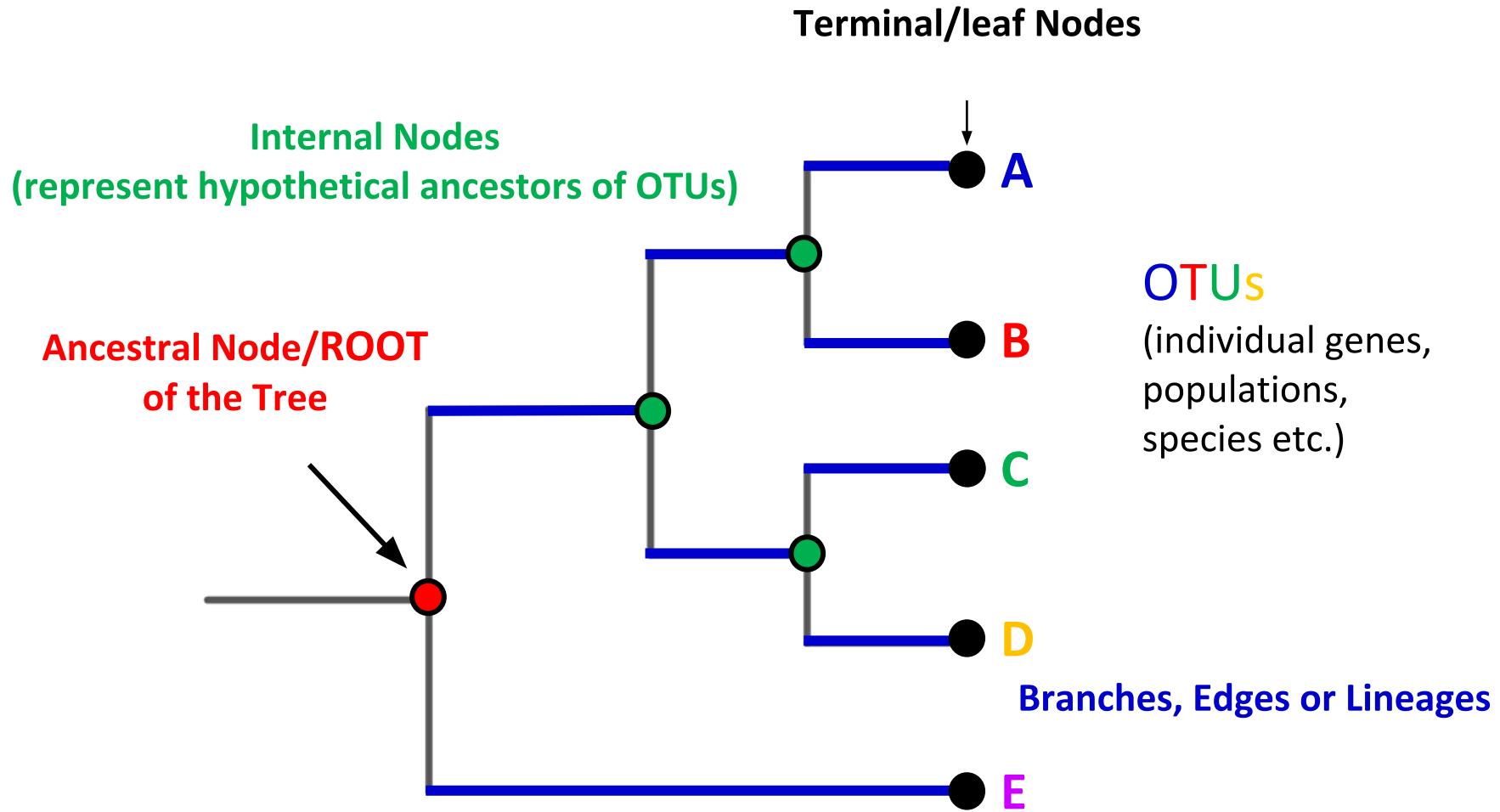
# Anatomy of a phylogenetic tree



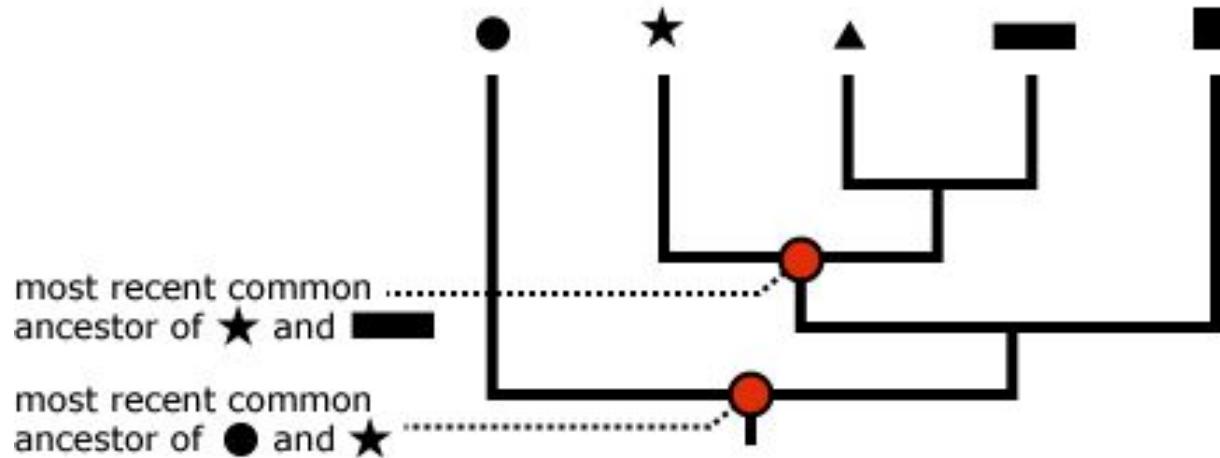
# Anatomy of a phylogenetic tree



# Anatomy of a phylogenetic tree

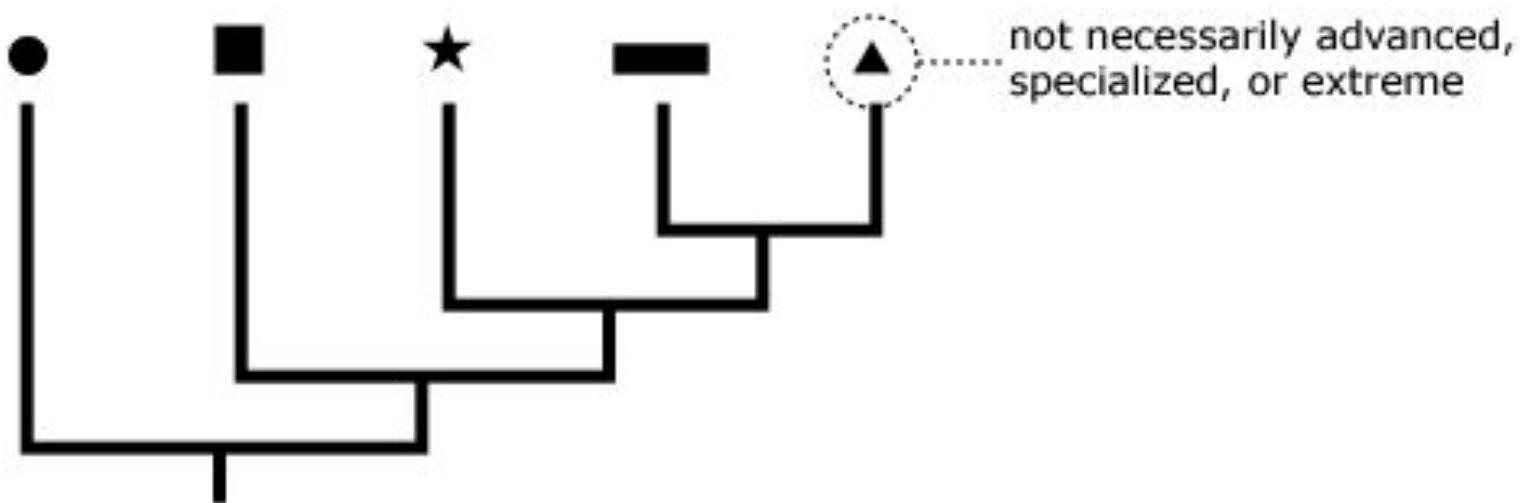


# Reading trees: patterns of relatedness



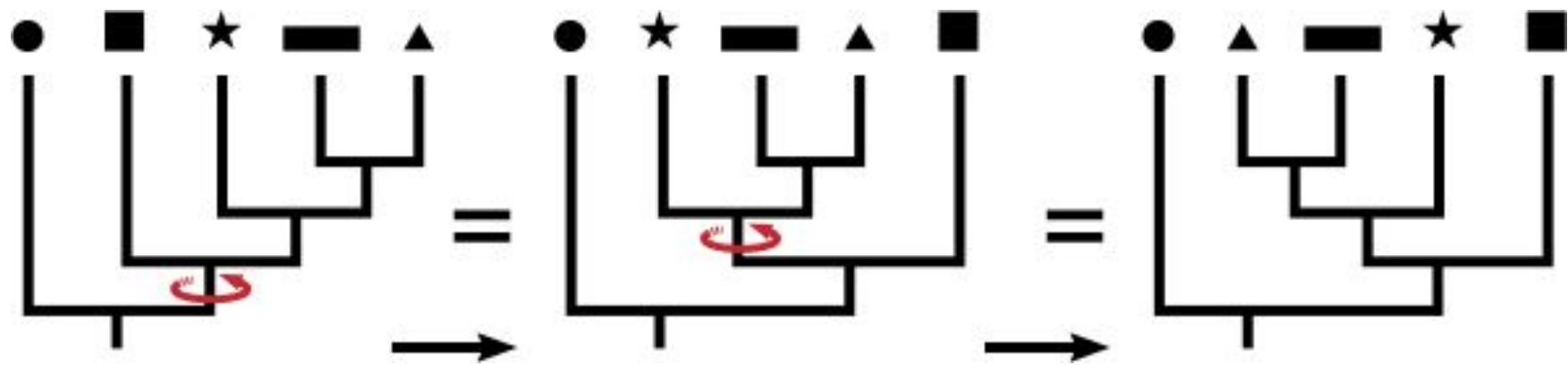
***The more recently two species share a common ancestor, the more closely related they are.***

# Reading trees: patterns of relatedness



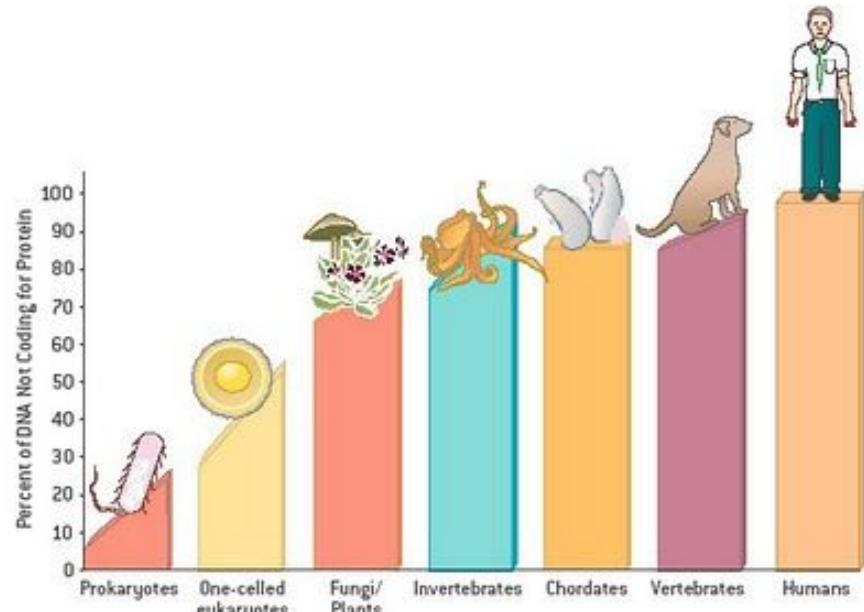
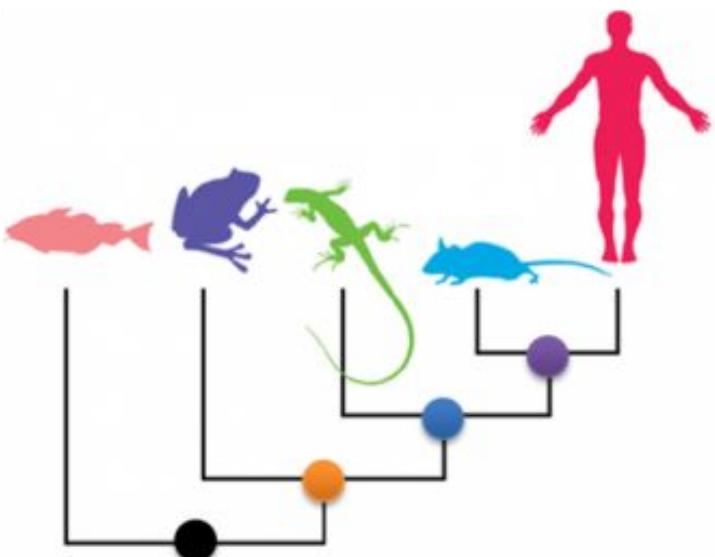
***Trees depict evolutionary relationships, not evolutionary progress (!)***

## Reading trees: patterns of relatedness



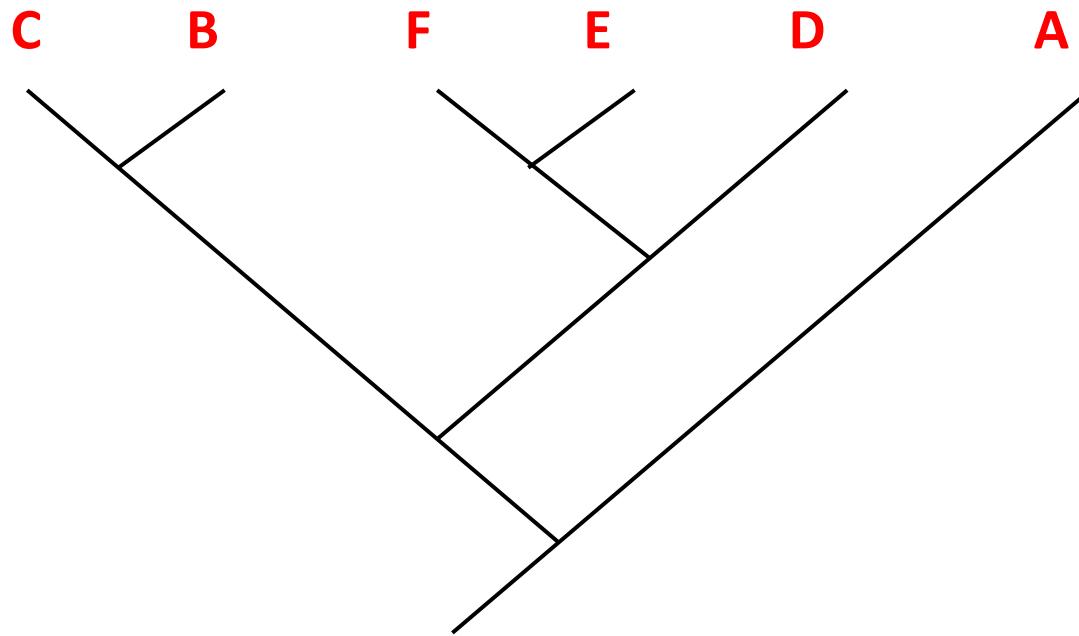
***Rotation around internal nodes produces equivalent trees.***

# Beware of “Dog’s Ass” plots



NONPROTEIN-CODING SEQUENCES make up only a small fraction of the DNA of prokaryotes. Among eukaryotes, as their complexity increases, generally so, too, does the proportion of their DNA that does not code for protein. The noncoding sequences have been considered junk, but perhaps it actually helps to explain organisms' complexity.

## Quiz: reading trees

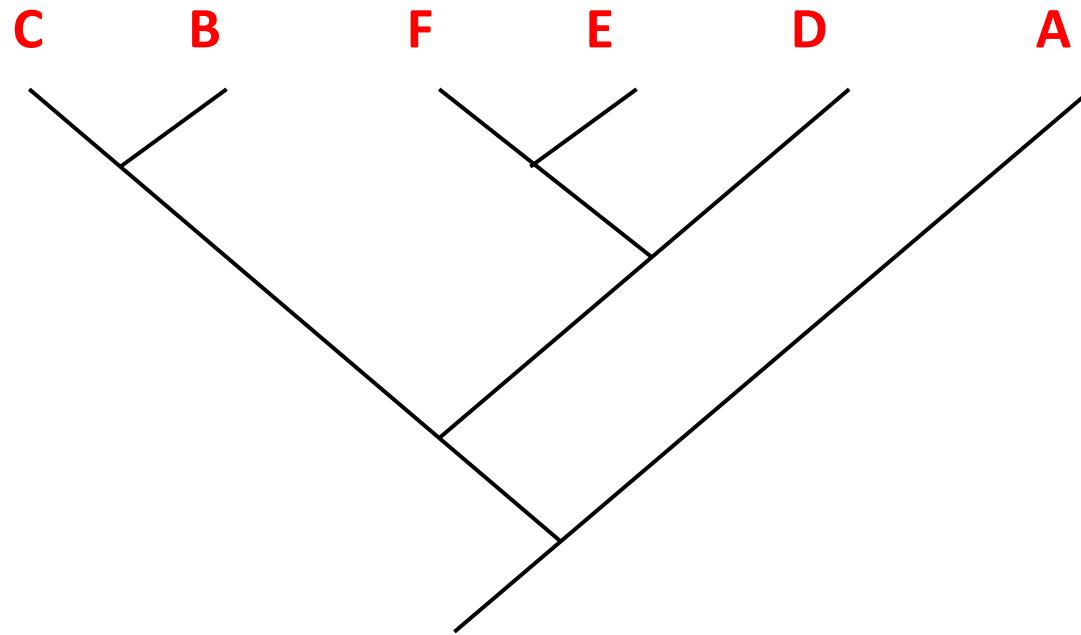


Is “E” more closely related to “D” or to “F”?

Is “E” more closely related to “B” or to “A”?

Is “E” more closely related to “B” or to “C”?

## Quiz: reading trees



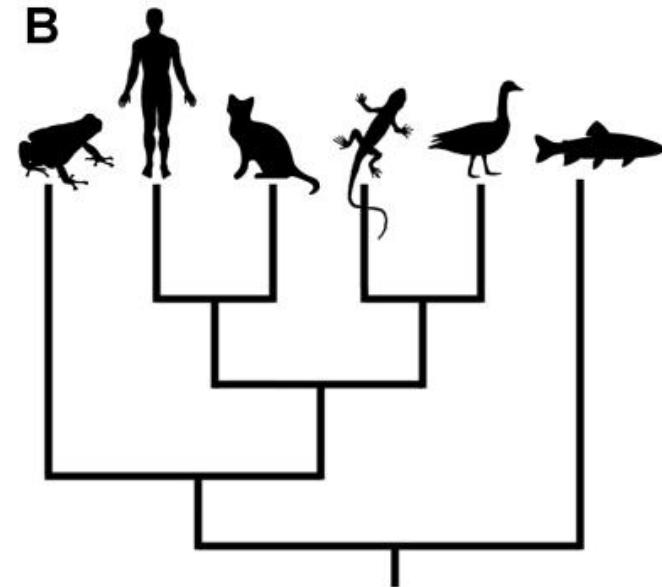
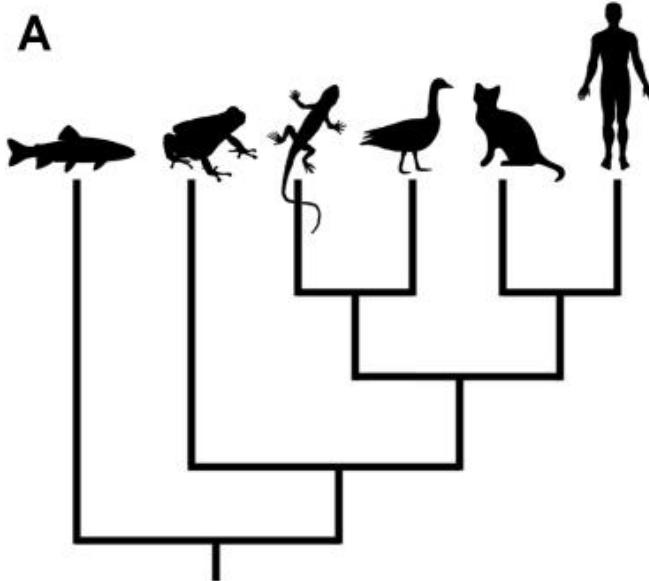
“E” more closely related to “F”

“E” more closely related to “B”

“E” more closely related to neither (equally to “B” & “C”)

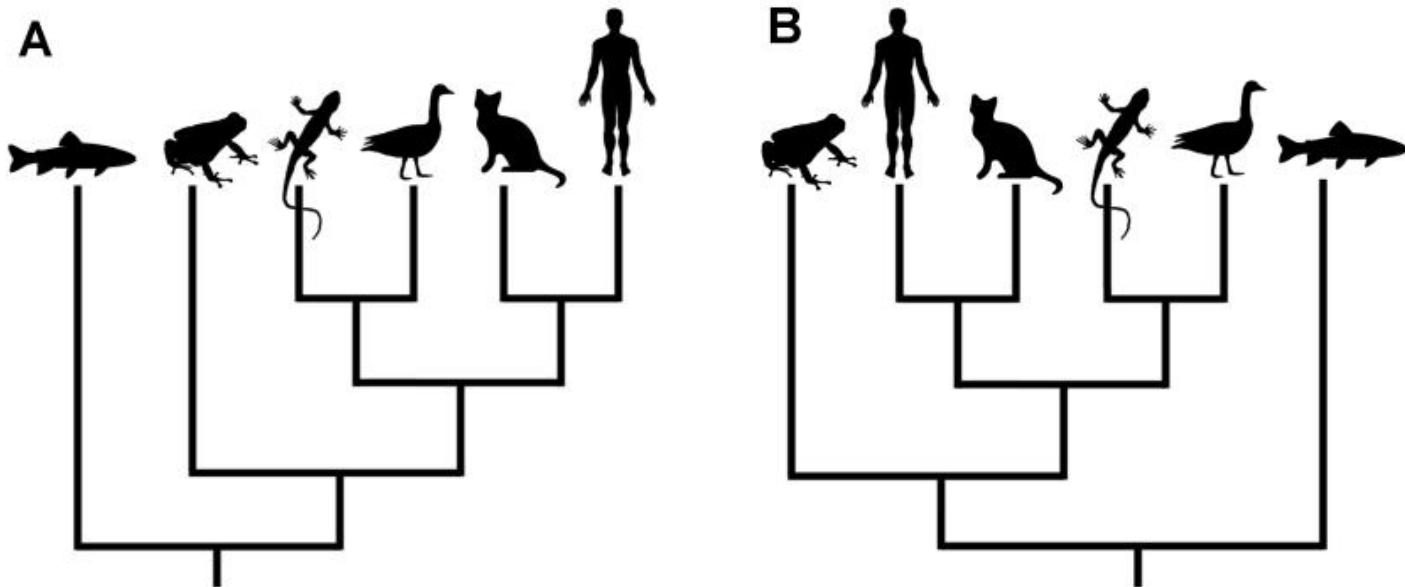
# Quiz: reading trees

Find the differences between the two trees!



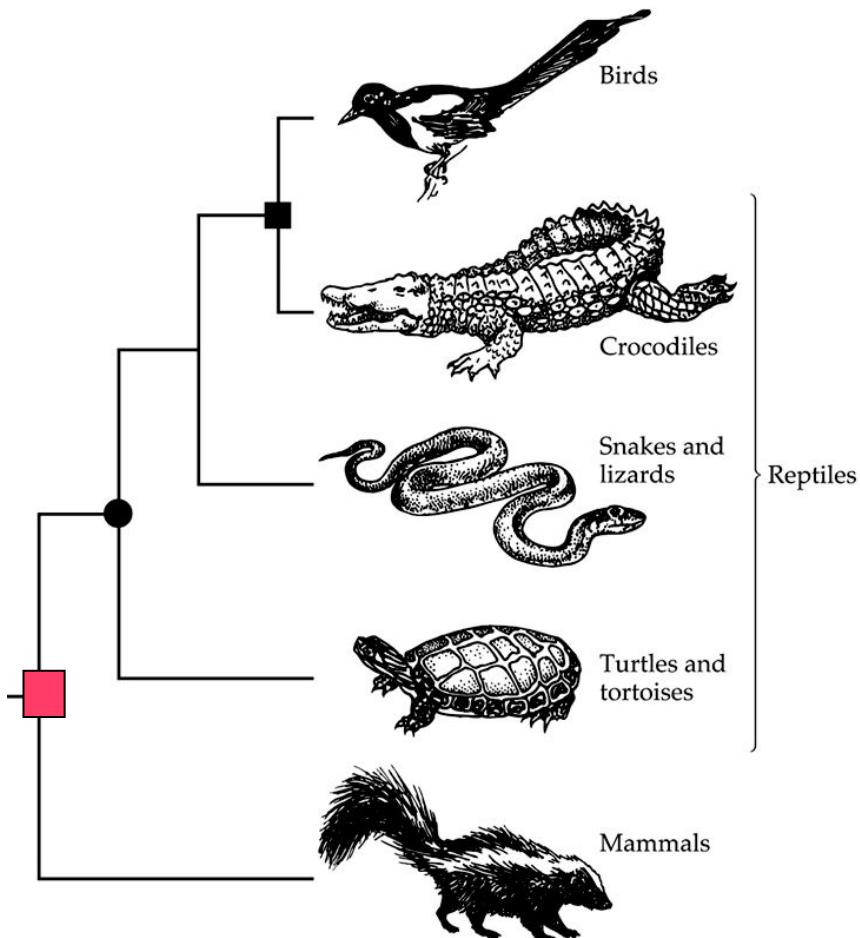
# Quiz: reading trees

Find the differences between the two trees!



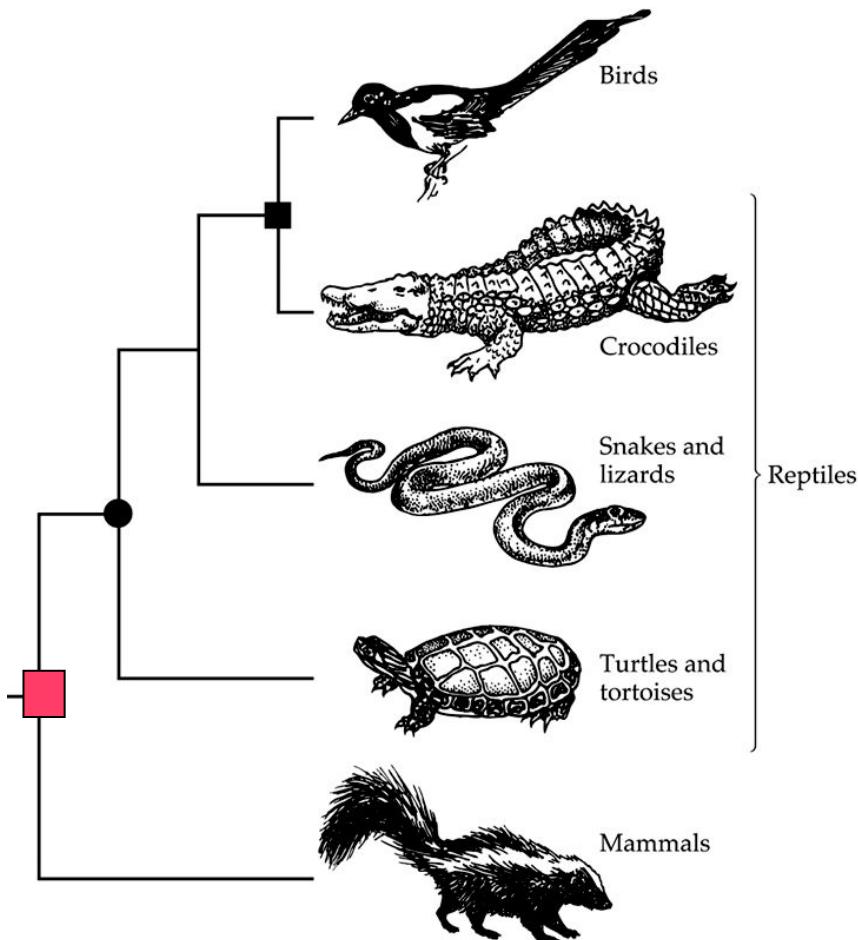
They are the same.

### 3. Naming groups on trees



Monophyletic group: a node and all its descendants (e.g. “amniotes”)

# Naming groups on trees

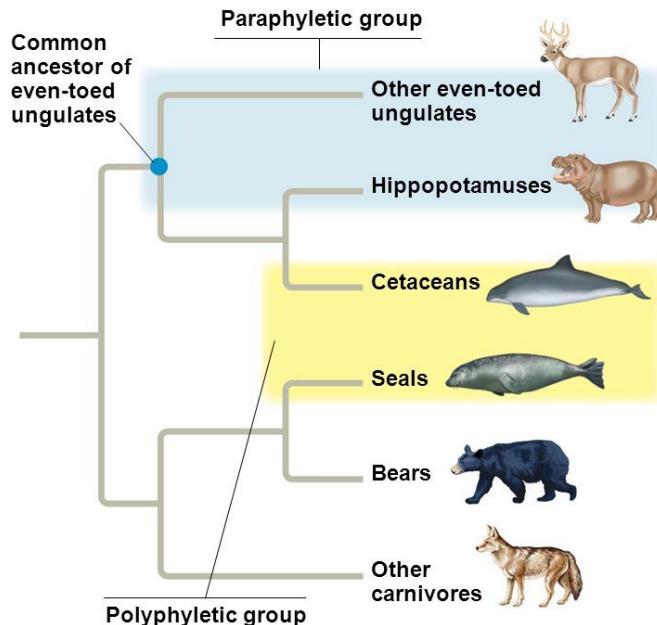


**Monophyletic** group: a node and all its descendants (e.g. “amniotes”)

**Paraphyletic** group: a node and some of its descendants (e.g. “reptiles”, “prokaryotes”)

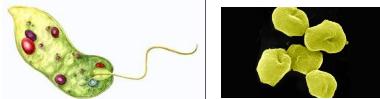
# Naming groups on trees

Figure 26.11



© 2014 Pearson Education, Inc.

(Don't feel bad if your favourite group is paraphyletic!  
Join the club!)



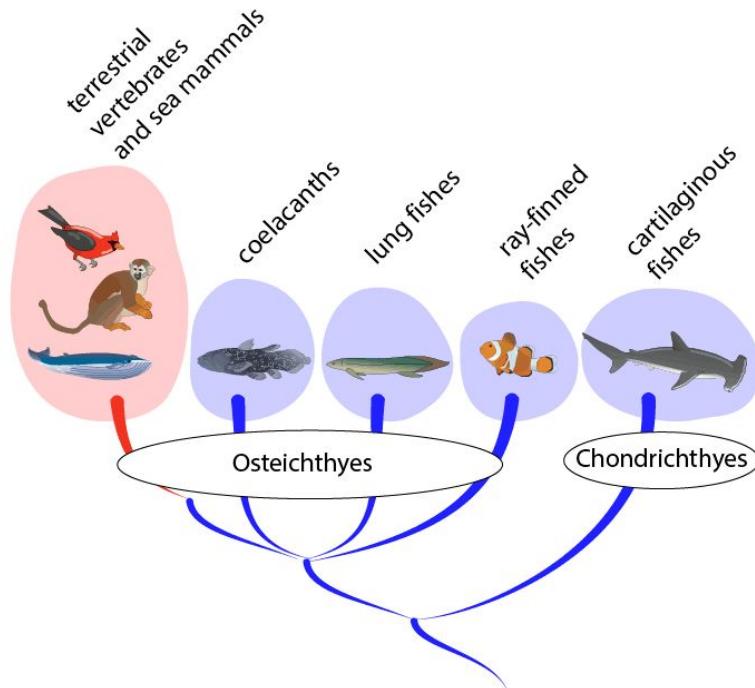
**Monophyletic** group: a node and all its descendants (e.g. “amniotes”)

**Paraphyletic** group: a node and some of its descendants (e.g. “reptiles”, “prokaryotes”)

**Polyphyletic** group: some taxa, but not their most recent common ancestor (e.g. “marine mammals”)

# Quiz: naming groups

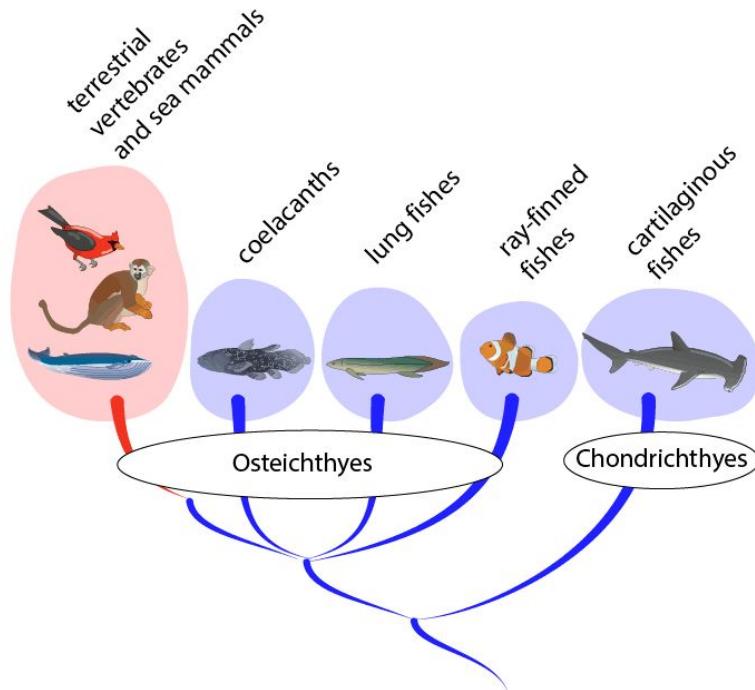
What kind of group are fishes?



- (a) Monophyletic
- (b) Paraphyletic
- (c) Polyphyletic

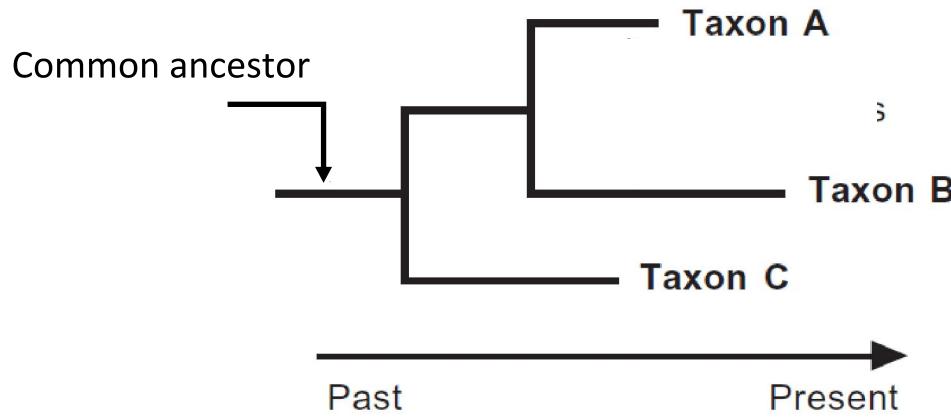
# Quiz: naming groups

What kind of group are fishes?



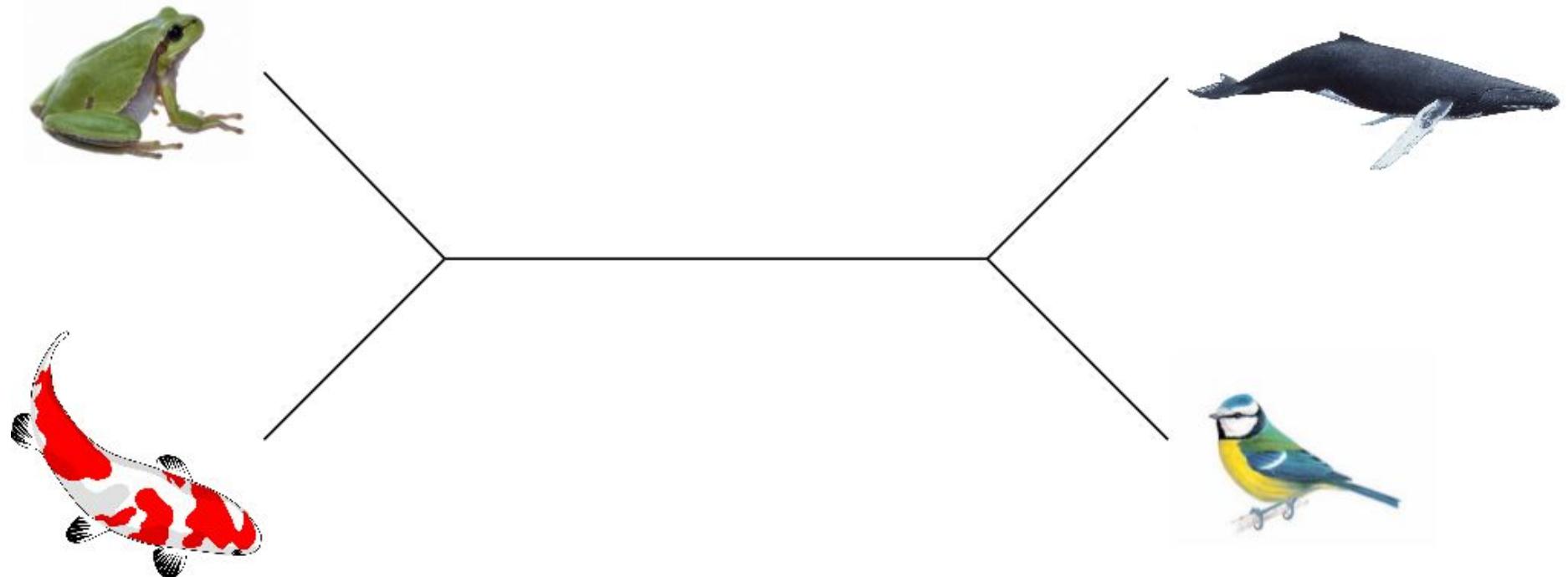
- (a) Monophyletic
- (b) Paraphyletic**
- (c) Polyphyletic

## 4. The root: a very special node on the tree

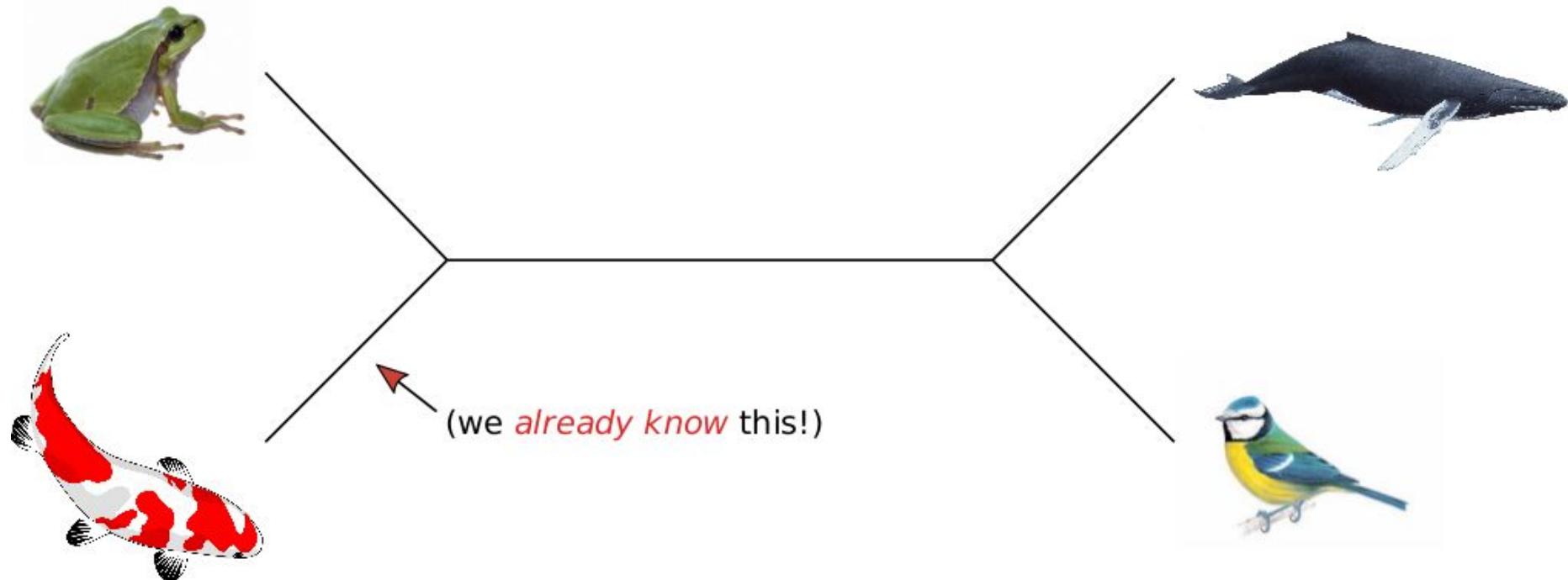


- The root is the oldest point on the tree (most recent common ancestor of all taxa). It orients the tree in time.
- All inferences about ancestor-descendant relationships **and relatedness** depend on the root position.
- Most phylogenetic methods do not estimate the root!

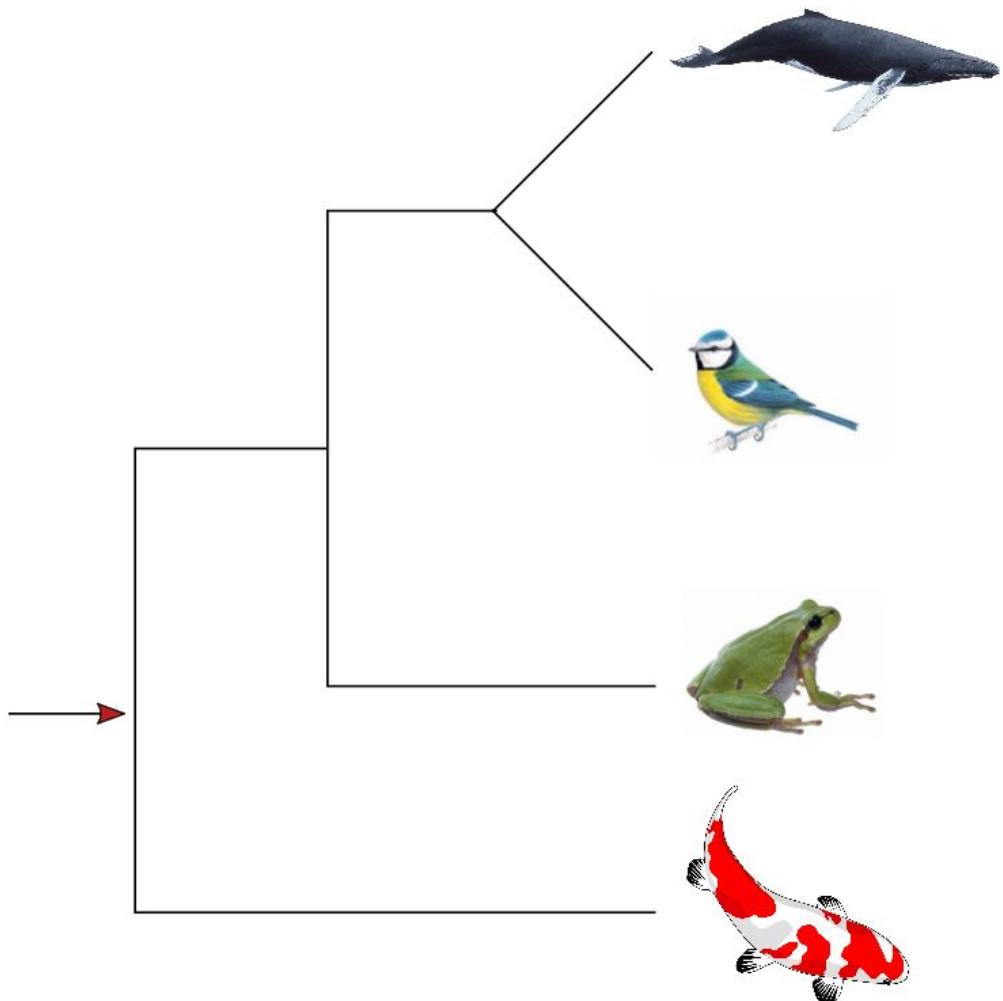
# Rooting an unrooted tree using *a priori* knowledge (outgroup)



# Rooting an unrooted tree using *a priori* knowledge (outgroup)



# Rooting an unrooted tree using *a priori* knowledge (outgroup)



- Most packages will arbitrarily root the tree on some branch, for no particular reason
- You need to re-root on the outgroup branch

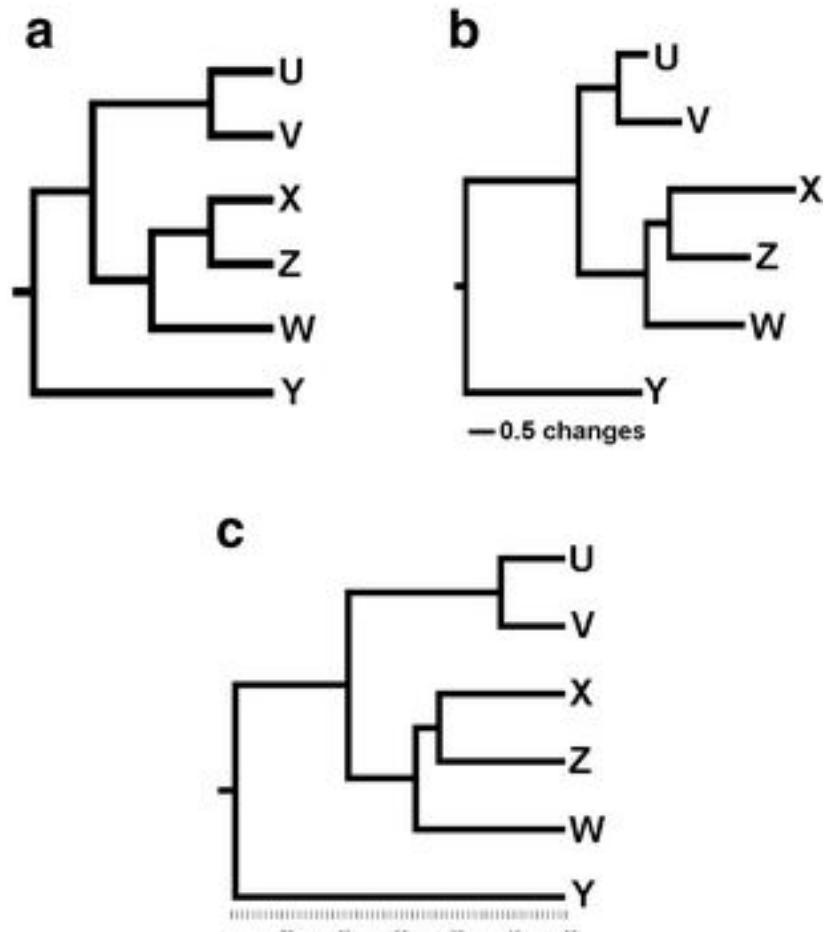
This doesn't apply if you use a method that can infer the root:

- Molecular clock
- Gene tree-species tree reconciliation
- Non-reversible/non-stationary phylogenetic model
- Some pop. gen. settings (if you know the ancestral/derived allele)

# Interpreting branch lengths

Depending on the method of analysis, branch lengths can have different interpretations.

- (a) **Cladogram**: just the topology matters (not so useful).
- (b) **Phylogram**: branch lengths correspond to genetic distance (changes/alignment site).
- (c) **Chronogram**: branch lengths correspond to relative or absolute times (number of years)

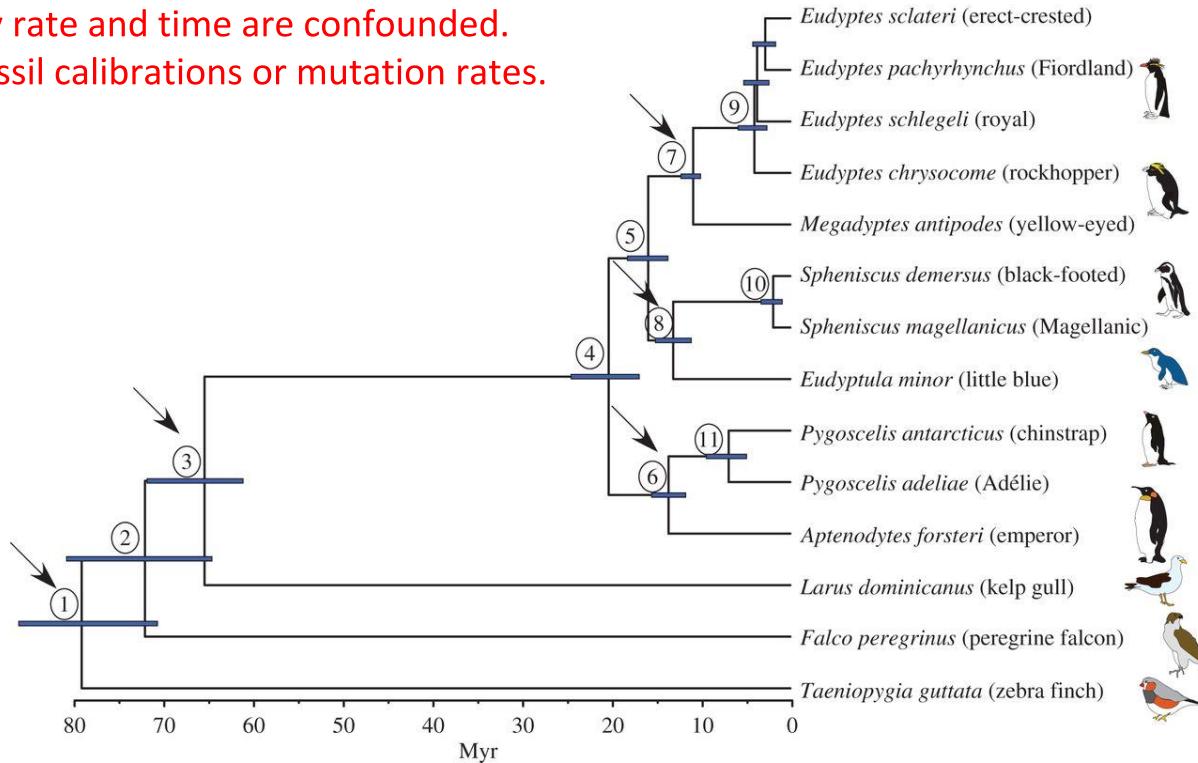


# Chronograms: branch lengths in time (estimated with a molecular clock)

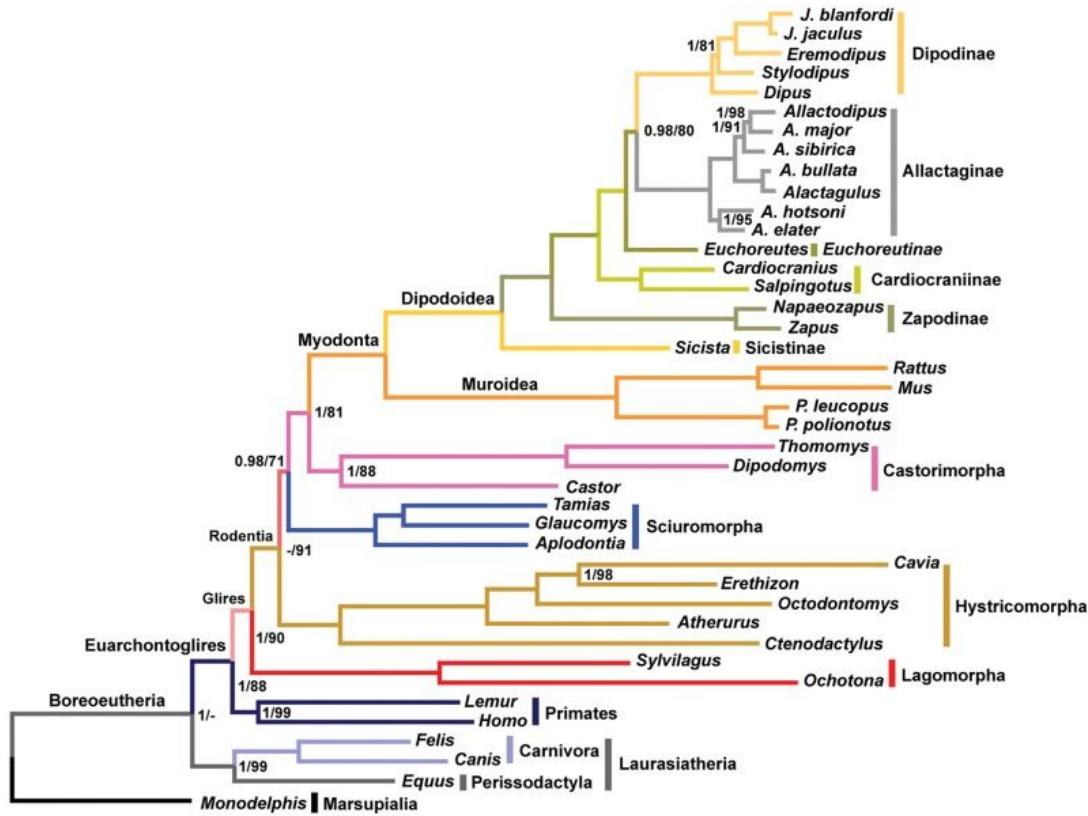
$$\text{Genetic change} = \text{Evolutionary rate} \times \text{Divergence time}$$

(substitutions/site)      (substitutions/site/year)      (years)

Evolutionary rate and time are confounded.  
Need e.g. fossil calibrations or mutation rates.



# Uncertainty in phylogenetic trees: support values



- Bootstraps (ML; 0-100) and posterior probabilities (Bayesian; 0-1)
- Different interpretations; higher is better
- Don't draw strong conclusions if the support values are low
- Apply to splits (branches) on the tree, not nodes
- Sometimes useful to summarise on same tree diagram
- Always think about supports, and include them in your published phylogenies!

# Bioinformatics Workshop Day 3: Schedule

Time	Event	Lecturer
10-12.30	Phylogenetics workflow & theory	Tom W.
1.30-3.30	Phylogenetics practical	Tom W.
3.45-4.30	Phylogenomics: theory and some games	Tom W.

# **Introduction to molecular phylogenetics**

**(b) A basic workflow**

**Tom Williams**

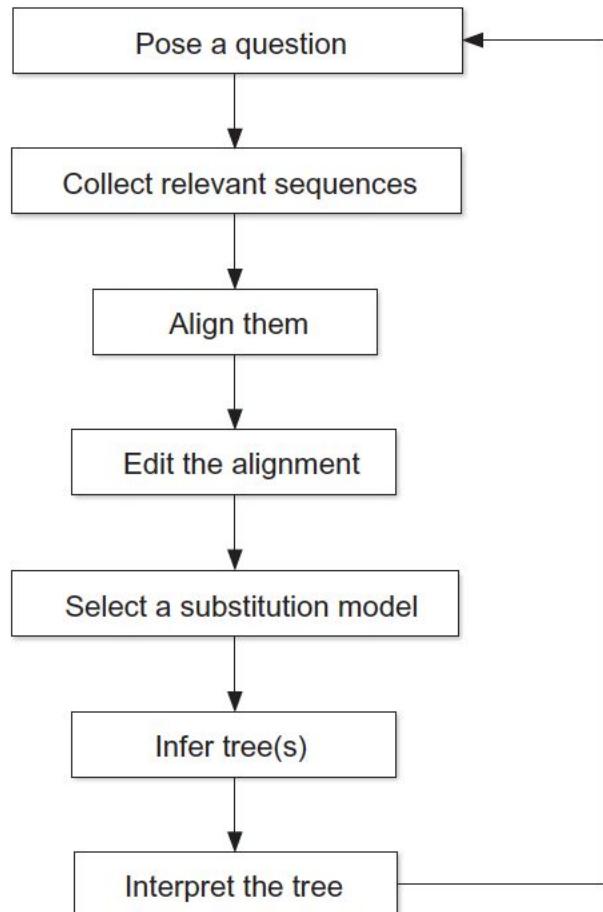
# A basic workflow for molecular phylogenetics

- Typical analysis steps
- Some discussion of the methodology at each step

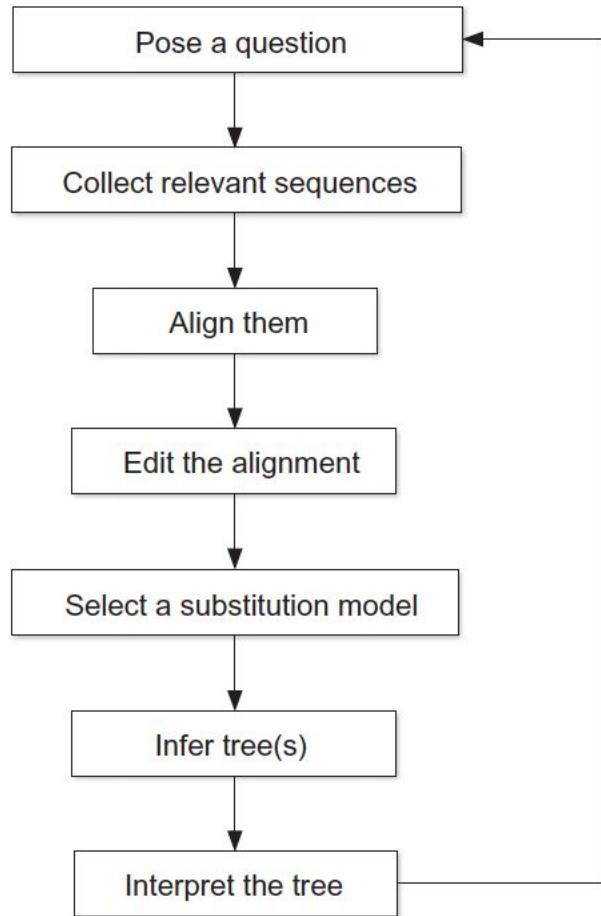
# A basic workflow for molecular phylogenetics

- Typical analysis steps
- Some discussion of the methodology at each step
- Core ideas/elements
- Somewhat more technical material

# Molecular phylogenetics: a possible flowchart

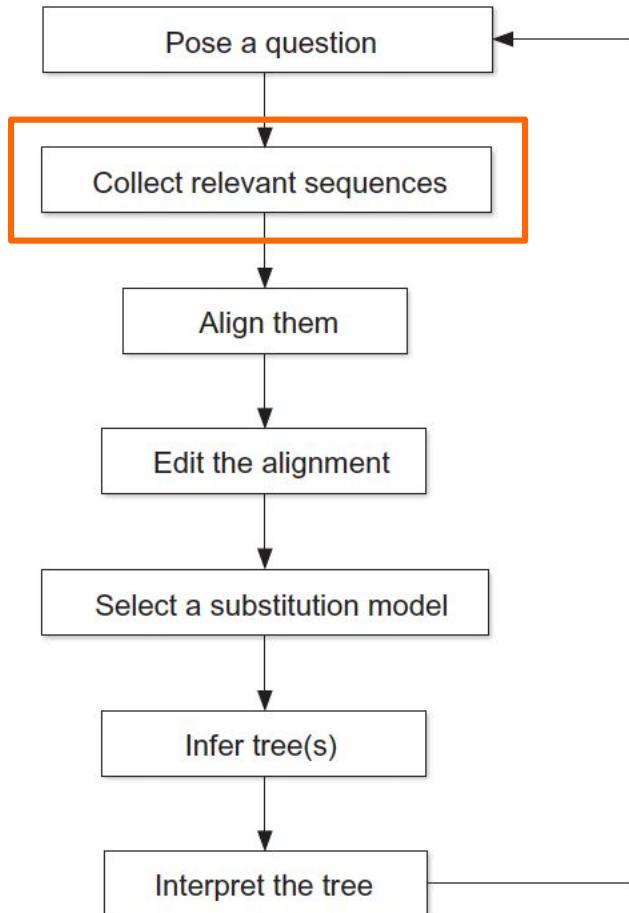


# Molecular phylogenetics: a possible flowchart



**A critical approach to analysis is key: don't fall into the black box!**

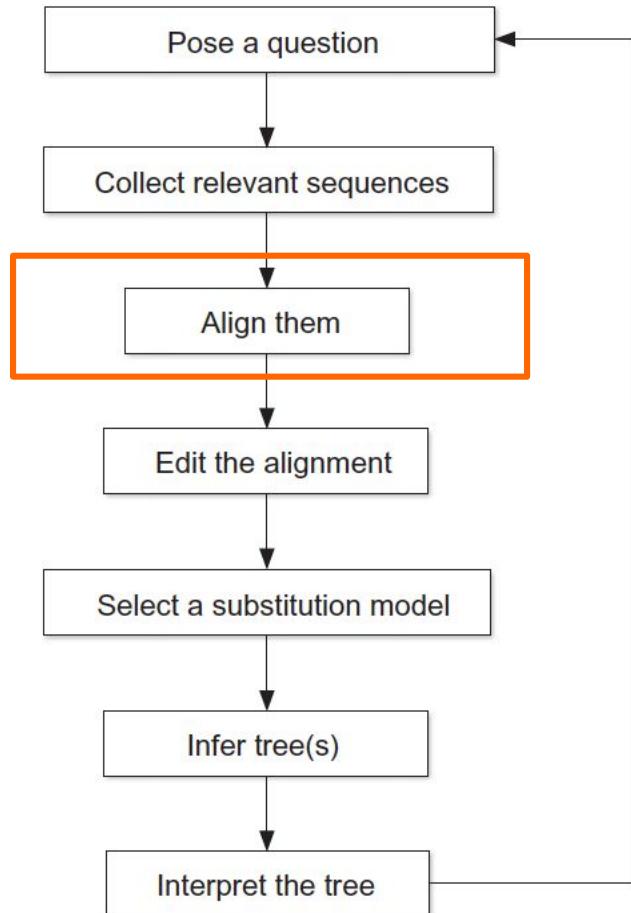
# Molecular phylogenetics: a possible flowchart



Did I miss any  
important ones?

**A critical approach to analysis is key: don't fall into the black box!**

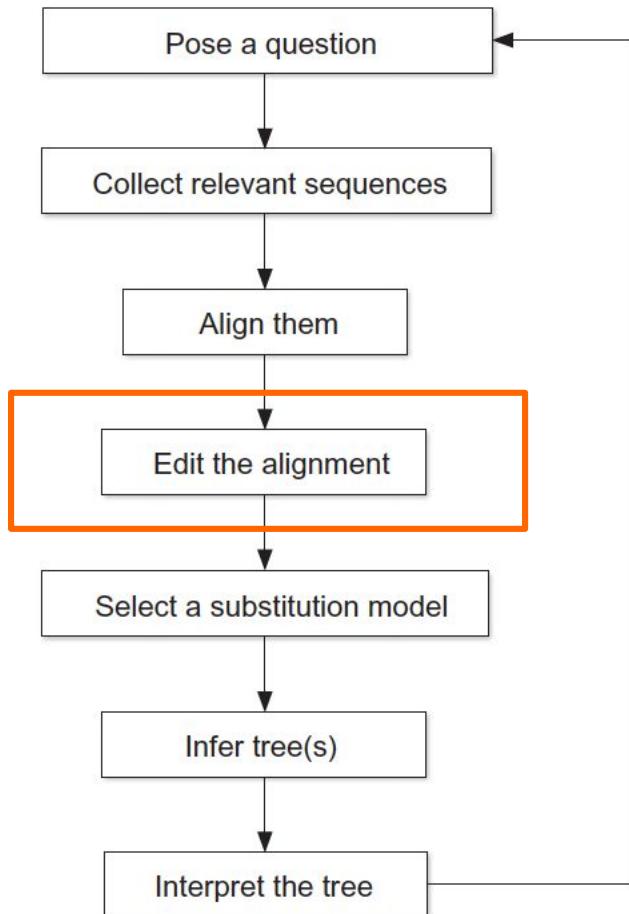
# Molecular phylogenetics: a possible flowchart



How certain is my alignment?  
Would other alignments give a different tree?

A critical approach to analysis is key: don't fall into the black box!

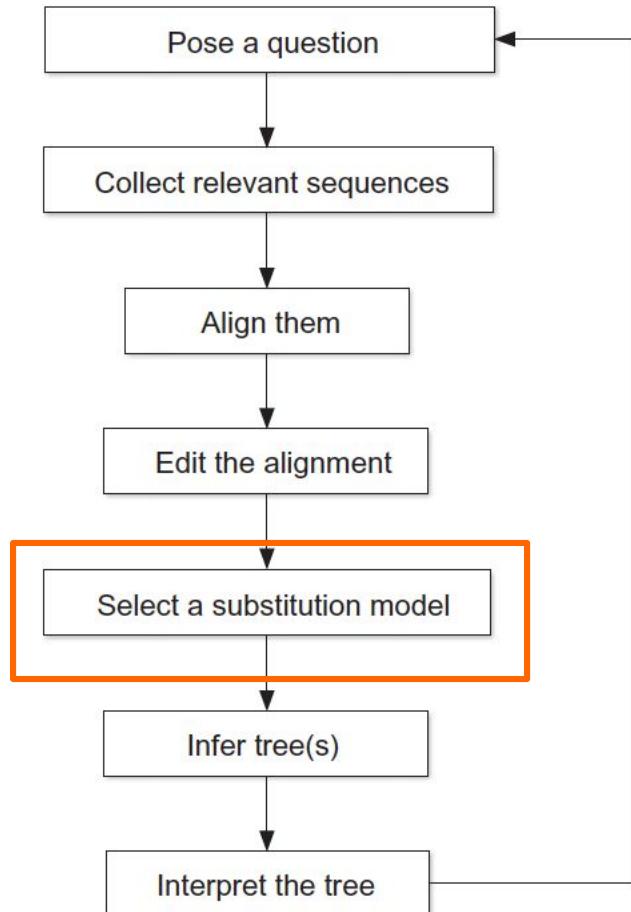
# Molecular phylogenetics: a possible flowchart



Which bits, if any,  
should I remove?

**A critical approach to analysis is key: don't fall into the black box!**

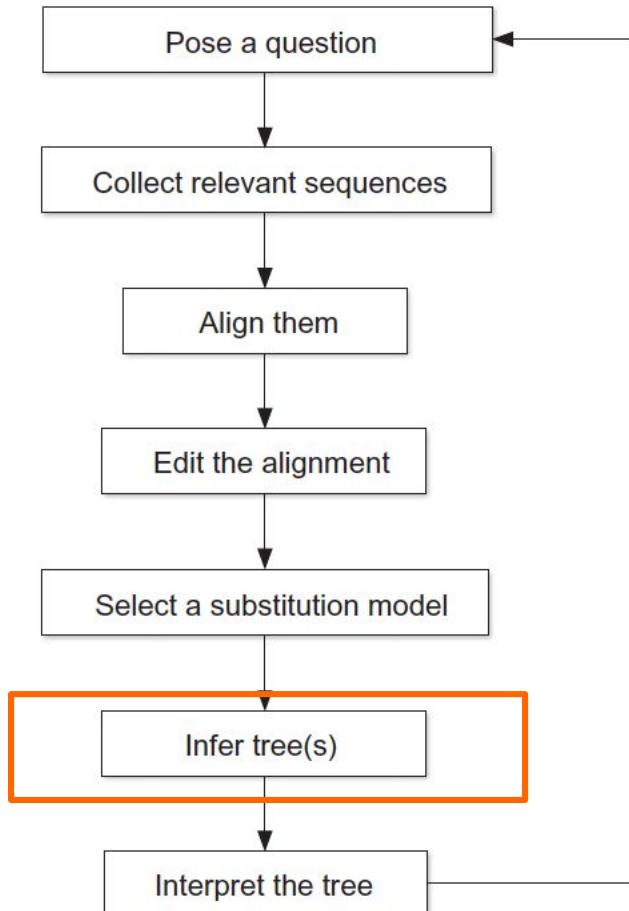
# Molecular phylogenetics: a possible flowchart



Is this the best  
available model?  
Is it good enough?

**A critical approach to analysis is key: don't fall into the black box!**

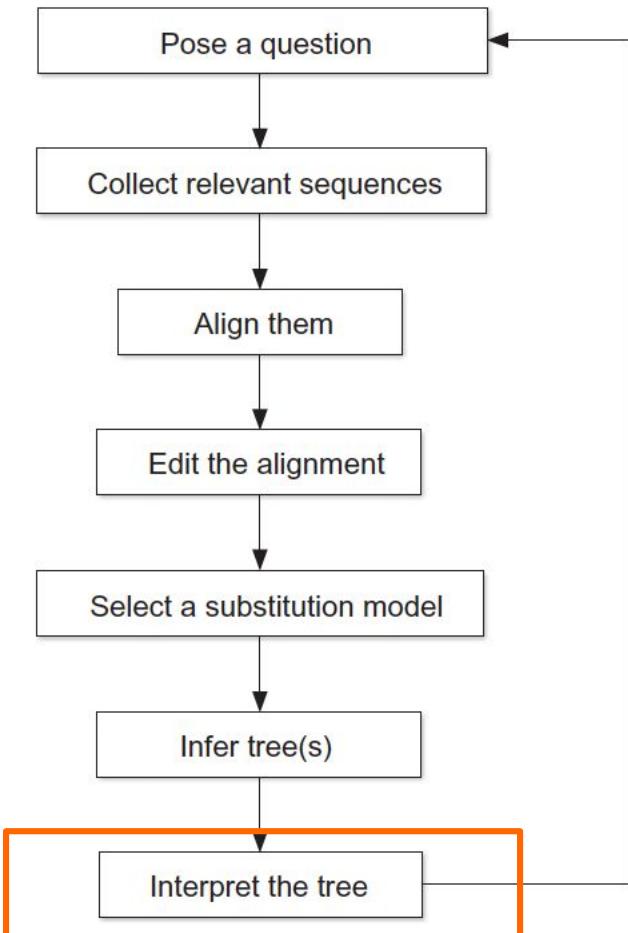
# Molecular phylogenetics: a possible flowchart



Are my ML or Bayesian analyses working properly? (Local maxima, convergence)

A critical approach to analysis is key: don't fall into the black box!

# Molecular phylogenetics: a possible flowchart



How well-supported  
are the key branches?

What does the tree  
mean for my  
hypothesis?

Can alternative  
hypotheses be  
rejected?

A critical approach to analysis is key: don't fall into the black box!

## 1. Pose a question

- Systematics, classification, character evolution in a clade of interest
- Use phylogenies as part of modelling a broader process

## 2. Assembling a dataset

- Where do I get the data?
- Use an existing dataset
- Use your own data
- Use published data to test a new hypothesis  
(GenBank, Ensembl, PDB, Dryad, FigShare...)

## 2. Assembling a dataset

- Starting with a query, search sequence databases.
- An **iterative process**: start broad and refine using initial trees
- Should I include it or not? **Do the experiment.**

## 2. Assembling a dataset

- Starting with a query, search sequence databases (BLAST, PSI-BLAST, HMMER).

Descriptions

Sequences producing significant alignments:							
Select: All None Selected:0							
<a href="#">Alignments</a> <a href="#">Download</a> <a href="#">GenPept</a> <a href="#">Graphics</a> <a href="#">Distance tree of results</a> <a href="#">Multiple alignment</a> <a href="#">Help</a>							
Description						Max score	
				Query cover	E value	Ident	Accession
<input type="checkbox"/> <a href="#">Translation elongation factor aEF2 [Lokiarchaeum sp. GC14_75]</a>				100%	0.0	100%	<a href="#">KKK44407.1</a>
<input type="checkbox"/> <a href="#">elongation factor aEF-2 [Candidatus Lokiarchaeota archaeon CR_4]</a>				94.6%	0.0	61%	<a href="#">QLS15932.1</a>
<input type="checkbox"/> <a href="#">hypothetical protein AM325_09125 [Candidatus Thorarchaeota archaeon SMTZ1-45]</a>				89.4%	0.0	59%	<a href="#">KXH72312.1</a>
<input type="checkbox"/> <a href="#">hypothetical protein AM324_08425 [Candidatus Thorarchaeota archaeon SMTZ1-83]</a>				88.5%	0.0	58%	<a href="#">KXH71555.1</a>
<input type="checkbox"/> <a href="#">Elongation factor 2 [Candidatus Odinarchaeota archaeon LCB_4]</a>				87.2%	0.0	57%	<a href="#">QLS17158.1</a>
<input type="checkbox"/> <a href="#">translation elongation factor aEF2 (CTD), partial [Lokiarchaeum sp. GC14_75]</a>				86.4%	0.0	94%	<a href="#">KKK44774.1</a>
<input type="checkbox"/> <a href="#">Elongation factor 2 [Candidatus Heimdallarchaeota archaeon AB_125]</a>				84.7%	0.0	56%	<a href="#">QLS3283.1</a>
<input type="checkbox"/> <a href="#">elongation factor EF-2 [archaeon ex4484_74]</a>				78.4%	0.0	52%	<a href="#">QYT35986.1</a>
<input type="checkbox"/> <a href="#">putative elongation factor 2 [uncultured organism]</a>				76.6%	0.0	52%	<a href="#">AKC94987.1</a>
<input type="checkbox"/> <a href="#">translation elongation factor aEF2 (NTD), partial [Lokiarchaeum sp. GC14_75]</a>				55.0%	0.0	92%	<a href="#">KKK43242.1</a>
<input type="checkbox"/> <a href="#">elongation factor EF-2 [Thermofilum pendens]</a>				53.2%	0.0	39%	<a href="#">WP_011752282.1</a>
<input type="checkbox"/> <a href="#">elongation factor EF-2 [Thermofilum uzonense]</a>				53.2%	0.0	39%	<a href="#">WP_052884495.1</a>

## 2. Assembling a dataset

sequence.fasta (~/Downloads) - Pluma

File Edit View Search Tools Documents Help

sequence.txt sequence.fasta

```
1 >M90923.1 Human immunodeficiency virus type 1, viral sample LC03D, V3 region
2 TCTGAAAATTTCACGGACAATACTAAACCATAATAGTACAGCTGAATACATCTGTAAACAATTAAATTGTA
3 CAAGACCTGGCAACAATAAAGAAAAAGTATAACTATGGACCGGGAAAGTATTTTATGCAGGAGAAAT
4 AATAGGAGATATAAGACAAGCACATTGTAACTTAGTAGAGCAGCATGGAATGACACTTTAACAGATA
5 GTTGGAAAATTACAAGAACATTGGGAATAAAACAATAGTCTTAATCACTCCTCAGGAGGGGACCCAG
6 AAATTGTAATGCACAGTTT
7
8 >M90927.1 Human immunodeficiency virus type 1, viral sample LC03.DA08, V3 region
9 CTAGCAGAAGGAGAGGTAGTAATTAGATCTGAAAATTTCACGAACAATGCTAAACCATAATAGTACAGC
10 TGAATACATCTGTAAACAATTAAATTGTACAAGACCTGGCAACAATAAGAGAAAAAGTATAACTATGGGACC
11 GGGGAAAGTATTTTATGCAGGAGAAATAATAGGAGATATAAGACAAGCACATTGTAACTTAGTAGAGCA
12 GCATGGAATGACACTTAAACAGATAGTTGAAAATTGCAAGAACATTGGGAATAAAACAATAGTCT
13 TTAATCACTCCTCAGGAGGGGACCCAGAAATTGTGATGCACAGTTT
14
15 >M90926.1 Human immunodeficiency virus type 1, viral sample LC03.DA07, V3 region
16 CTAGCAGAAAAGAGGTAGTAATTAGATCTGAAAATTTCACGGACAATACTAAACCATAATAACAGC
17 TAAATACATCTGTAAACAATTAAATTGTACAAGACCGGGCAACAATAACAAGAAAAAGTATAACTATGGGACC
18 GGGGAAAGTATTTTATGCAGGAGAAATAATAGGAGATATAAGACAAGCACATTGTAACTTAGTAGAGCA
19 GCATGGAATGACACTTAAACAGATAGTTGAAAATTACAAGAACATTGGGAATAAAACAATAGTCT
20 TTAATCACTCCTCAGGAGGGGACCCAGAAATTGTGATGCACAGTTT
21
22 >M90925.1 Human immunodeficiency virus type 1, viral sample LC03.DA04, V3 region
23 CTAGCAGAAAAGAGGTAAATTAGATCTGAAAATTTCACGGACAATACTAAACCATAATAACAGC
24 TGAATACATCTGTAAACAATTAAATTGTACAAGACCTGGCAACAATAACAAGAAAAAGTATAACTATGGGACC
25 GGGGAAAGTATTTTATGCAGGAGAAATAATAGGAGATATAAGACAAGCACATTGTAACTTAGTAGAGCA
26 GCATGGAATGACACTTAAACAGATAGTTGAAAATTACAAGAACATTGGGAATAAAACAATAATCT
27 TTAATCACTCCTCAGGAGGGGACCCAGAAATTGTGATGCACAGTTT
28
29 >M90924.1 Human immunodeficiency virus type 1, viral sample LC03.DA02, V3 region
30 CTAGCAGAAAAGAGGTAGTAATTAGATCTGAAAATTTCACGGACAATACTAAACCATAATAGTACAGC
31 TAAATACATCTGTAAACAATTAAATTGTACAAGACCTGGCAACAATAACAAGAAAAAGTATAACTATGGGACC
32 GGGGAAAGTATTTTATGCAGGAGAAATAATAGGAGATATAAGACAAGCACATTGTAACTTAGTAGAGCA
33 GCATGGAATGACACTTAAACAGATAGTTGAAAATTACAAGAACATTGGGAATAAAACAATAAGTCT
34 TTAATCACTCCTCAGGAGGGGACCCAGAAATTGTGATGCACAGTTT
35
36 >M90929.1 Human immunodeficiency virus type 1, viral sample LC03.DA15, V3 region
37 CTAGCAGAAGGAGAGGTAGTAATTAGATCTGAAAATTTCACGAACAATGCTAAACCATAATAGTACAGC
38 TGAATACATCTGTAAACAATTAAATTGTACAAGACCTGGCAACAATAACAAGAAAAAGTATAACTATGGGACC
39 GGGGAAAGTATTTTATGCAGGAGAAATAATAGGAGATATAAGACAAGCACATTGTAACTTAGTAGAGCA
40 GCATGGAATGACACTTAAACAGATAGTTGAAAATTACAAGAACATTGGGAATAAAACAATAAGTCT
```

Plain Text ▾ Tab Width: 4 ▾ Ln 1, Col 1 INS

### 3. Alignments: hypotheses of homology

ATTGGCG  
AGGAG



Sequences!

ATTGGCG  
A--GGAG



AGGCG → +TT  
C=>A

ATTGGCG  
A--GGAG

Alignment!  
Homology!

# Some bits are easier to align than others

The diagram illustrates a sequence alignment between two proteins. The top sequence is aligned with the bottom sequence. Regions of high conservation or similarity are highlighted with colored boxes. The colors used include blue, green, yellow, orange, red, purple, and pink. The alignment shows that some regions are highly conserved across both sequences, while others show more variation.

Sequence 1 (Top):

```
Y EPGDTVTIYPCNTD- EDVSRFLANQSHWL----- EIADKPLNFTSGVPNDLKDGGLVRPMTLRNLL  
FAAGDVVLIQPSNSA-AHVQR- FCQVLGLDP----- DQLFMLQPREDPVSSPTRLPQPCSMRHLV  
FA-GQTLLIYPKNYP-KDVQK- LIDLMGWSEV----- AEQRIEIDWVKGTRPRDYHFLKDATIRDVL  
FSPGDCLSFCPENYN--- YKEF-MKYNGM-E----- EDVDG----- ISSSLMM  
FEIGDCLAVYPENYN--- YEEF-VRYNNI-K----- DNT----- LVKYI  
FFPGDCLSLLPENYN--- YREF-MSYNGI-G----- DGDDLG----- VSSVWMLI  
FEPGDTIRIYPSNYN--- WREF-CDYIGN-V----- DDE----- DHV  
YCLGDSLALYQQNPV-NEAIK-AIEMFGYNPYSLLRLSINEENEANNTNKVNQRYSSLFGYDITVLQLF  
FVPGQYATMRYE-----  
YKAGQFITLGLPNPV-----
```

Sequence 2 (Bottom):

```
LATGTPNLTSIQESQQLSSL-NKDFTHTF----- DLEGNSRDIISSSKFTKKD-- SIFISTDRKI  
IAAGTKPNLVSIKENQQFRVLDQDFTHTF----- DLQGNNIEITTSPKPTKKD-- SIFISPDKKI  
YEAGDALGIWPTNSP-ELVAE-ILSALKLSP----- ATVVTVKDK-- GDMTIAEAL  
YEAGDALGVWPVNRP-DLVAE-WLAUTGLDP----- ADTVTVAD-- GPVPLGEAL  
YEPGDCIVVLPQNDP-SIVEL-LISTLGWDY----- EQVLIINEDG-- DTLNLLEEAL  
YEAGDALGVIPVNNEP-SLVSL-LLTQLNADY----- QTPVPGF-- DRSLGDLL  
YQPGDHLCVVPEENDD-AVVER-LLRRFNLD----- ADTYVRIESRSEMRCGPFPSCSTFSVYNLA  
YAPGDAALGLIAHNDD-GVVDA-LLKALSFDP----- AIQVPSNG-- RMISVASAL
```

Homologous!

???

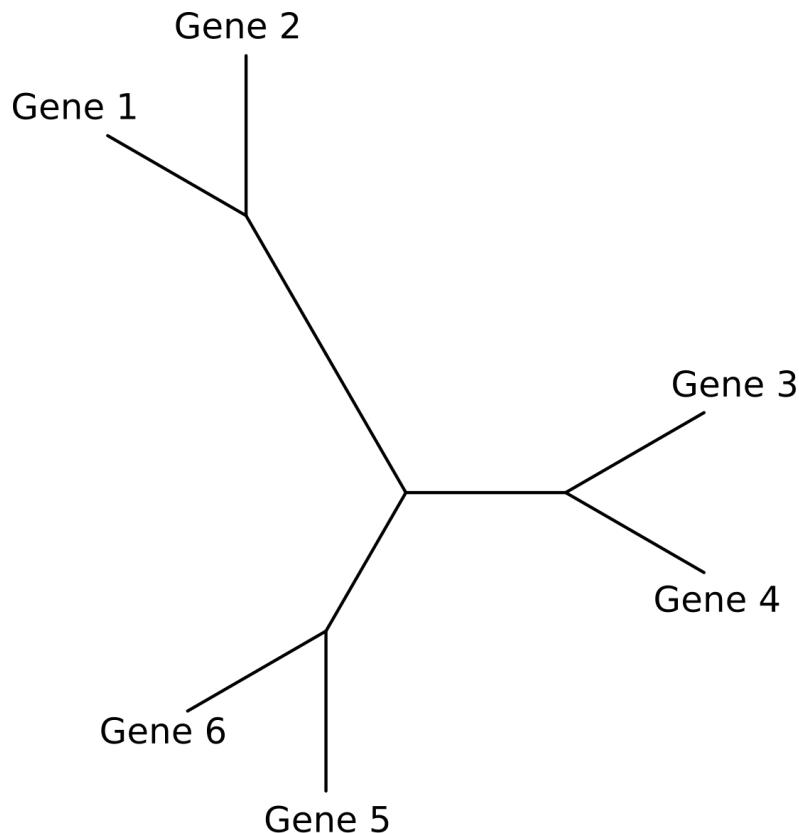
### 3. Alignments: some considerations

- Tree inference methods assume characters in each alignment column are homologous
- Automated alignment methods often disagree
- The alignment is uncertain, but often fixed in downstream analyses. **It matters!**

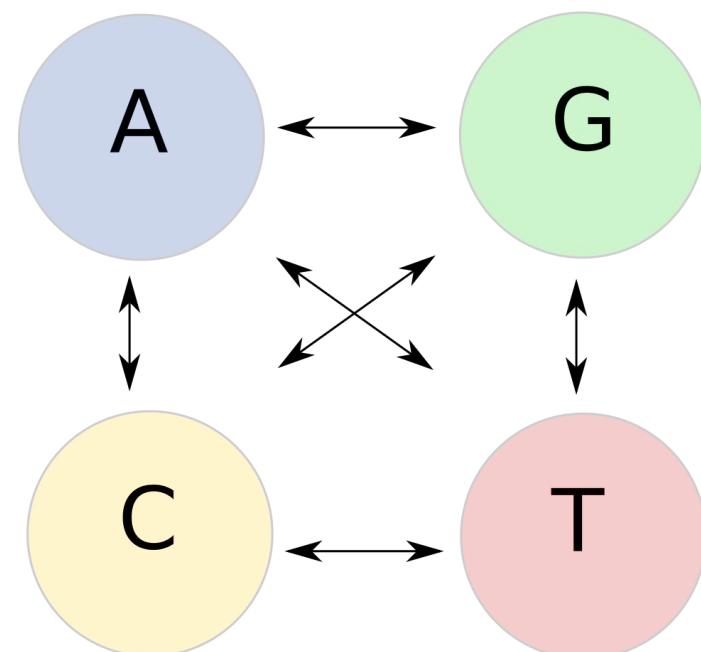
## 4. Modelling sequence evolution and inferring trees

- Modern phylogenetic analysis is **model-based**
- Fitting this model to the sequence alignment helps us learn about the parameters of the model (**including the tree topology**)
- Inference depends on **likelihood**: the probability of the data given the model.

# A basic model



(a) Phylogenetic tree  
(with branch lengths)

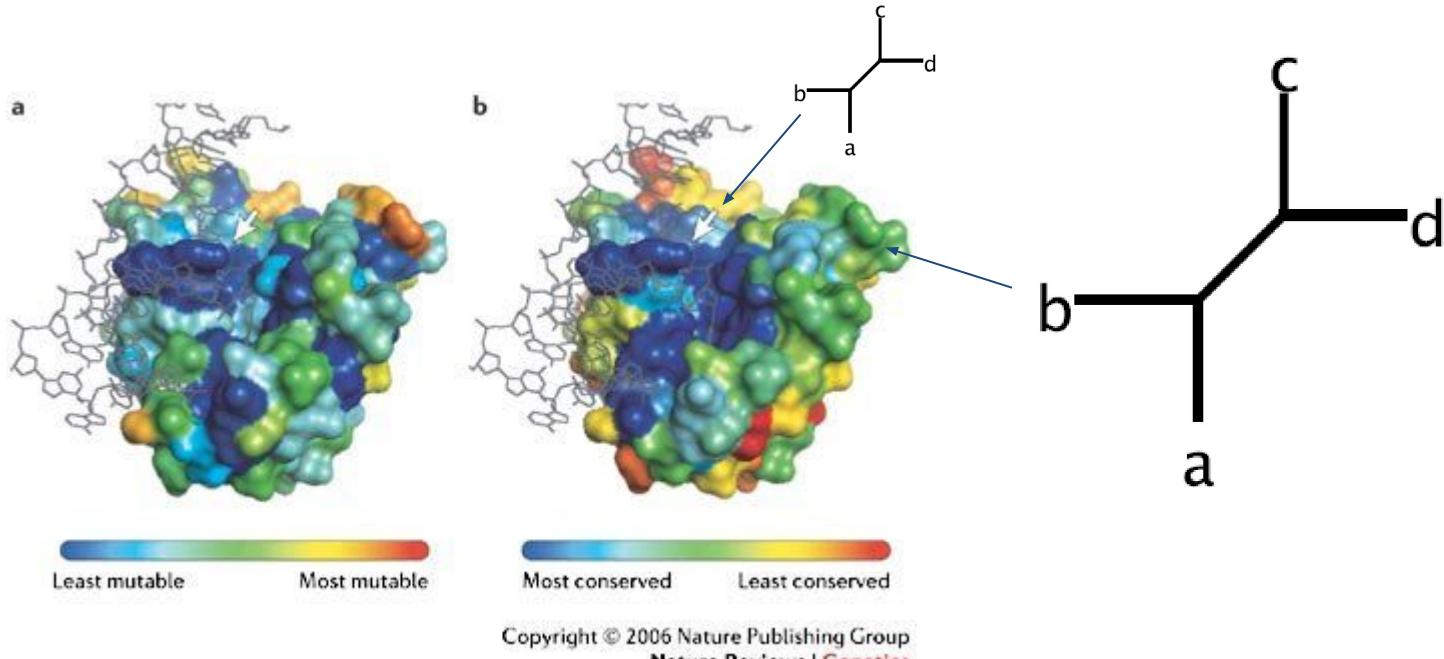


(b) Substitution process  
Rates, compositions

## Adding complexity (and parameters) to the basic model

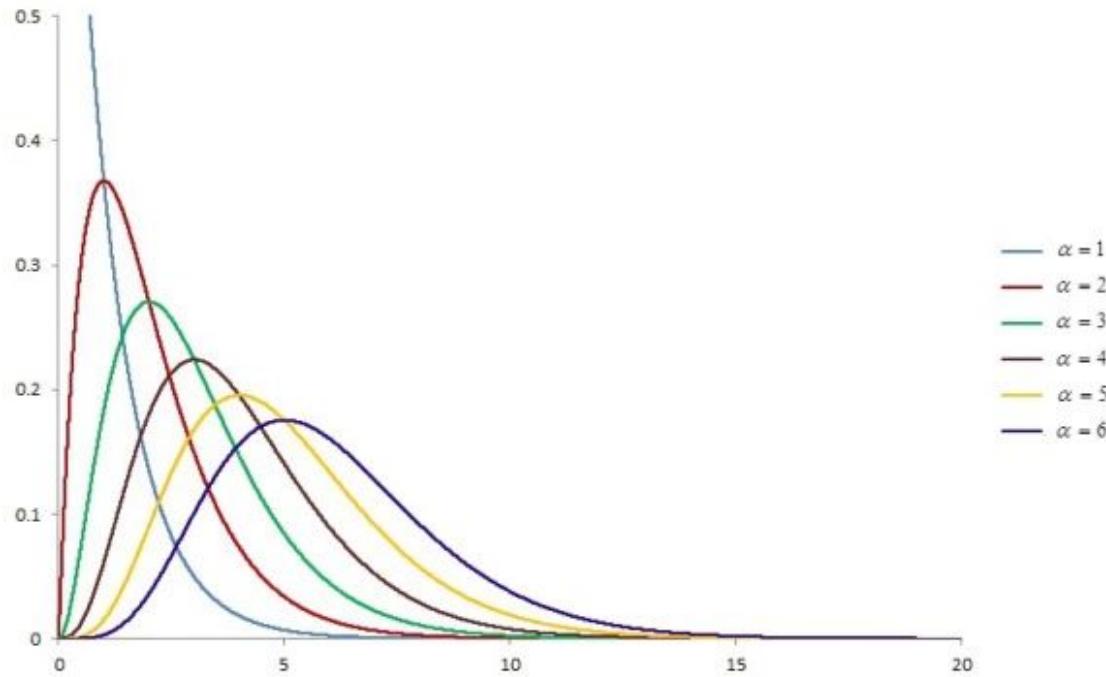
- Different **sites** evolve at different rates
- Different **genes** evolve at different rates
- The substitution process might **vary across the tree**  
(descent with modification)

# Adding complexity: site rates



- Some sites are more important for function than others; they tend to evolve more slowly.
- Model this by having several different rates (tree lengths)

## Adding complexity: site rates



- Gamma distribution can take many shapes, depending on a single parameter (alpha). So can model the distribution of rates-across-sites with just one more parameter!
- A **\*\*mixture model\*\***: average the likelihood per site over the possible rates!

# Choosing an appropriate model

Likelihood =  $P(\text{alignment} \mid \text{model})$

Tree  
Branch lengths  
Evolutionary rates  
Exchangeabilities  
etc.

- The evidence we extract from the data totally depends on the model.
- **All models are wrong!**  
....but some are useful (George Box).

How can we pick a useful model for our data?

# Choosing an appropriate model

$$\text{Likelihood} = P(\text{alignment} \mid \text{model})$$

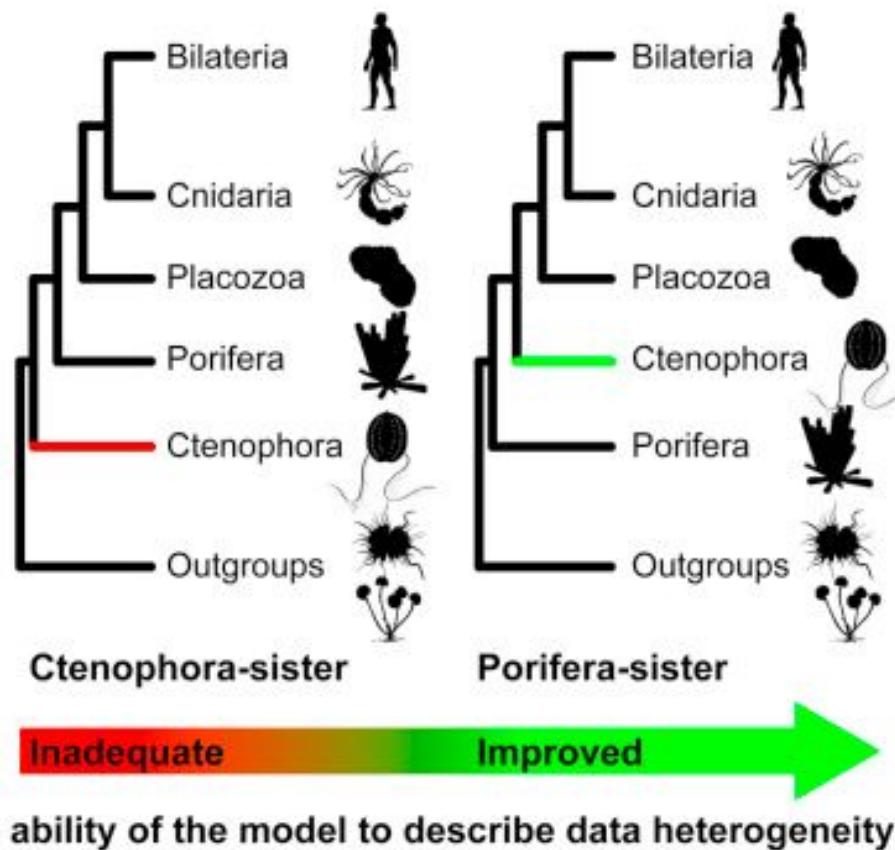
Tree  
Branch lengths  
Evolutionary rates  
Exchangeabilities  
etc.

- The likelihood provides a natural way of comparing models.
- Which model makes the data most likely?
- Subtle difficulties of maximum likelihood: model fit is not quite so easy. Use scores that are modifications of the maximum likelihood.

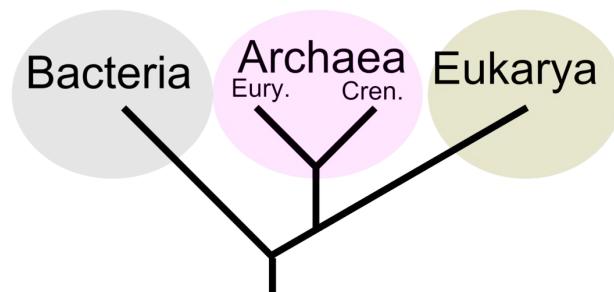
## Model selection

- Phylogeny packages implement a set of candidate models.
- Use statistical tests (AIC, BIC - based on the maximum likelihood) to pick **the best model from that set** for your data.
- It's possible that none of the implemented models is adequate for your data (use simulations?).
- \*\*\*Many published analyses are of very poor quality: **be critical and don't fall into this trap.**

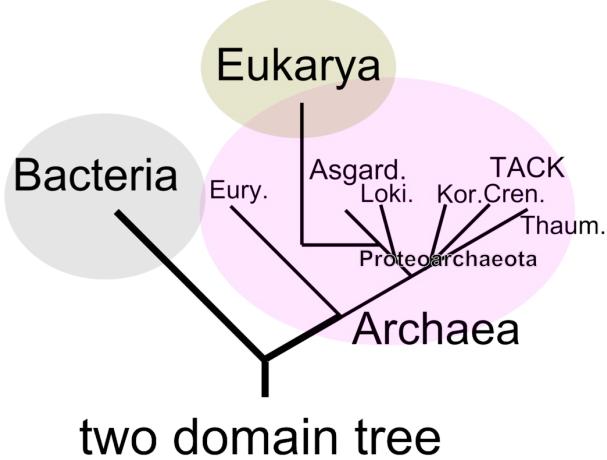
# Model selection matters: nervous system origins



# Model selection matters: The tree of life



Carl Woese's three domain tree



two domain tree

# Fitting models: maximum likelihood vs. Bayesian inference



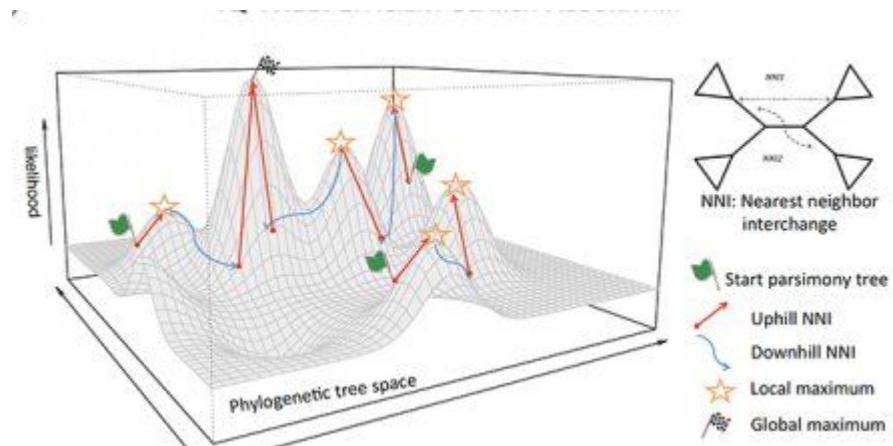
- . To a first approximation, the fit of the model is more important!

# Maximum likelihood

**Task:** find values for all of the parameters of the model that make the probability of the alignment (likelihood) as high as possible.

ML programs implement (hill climbing) algorithms for joint estimation of these values.

The ML tree is the best point estimate of the true tree, if the model is correct.



# Maximum likelihood: assessing support

- ML tree is best estimate of the true tree. But is it much better than any other estimate? (Hypothesis testing, e.g. **AU-test**)
- How certain are we of the relationships on the tree? Are some better supported than others? (**Bootstrapping**)

# Maximum likelihood: assessing support with the bootstrap

Bootstrapping is a technique for estimating the confidence interval around a parameter by **resampling and reanalysis** of the original data (alignment).

1. Sample, with replacement, the same number of characters as the original dataset.
2. Repeat tree inference on each of many bootstrap replicates.
3. Assess confidence as the proportion of samples in which each split on the tree appears.

Taxon	1	2	3	4
A	V	W	R	A
B	V	W	R	L
C	I	Y	K	M
D	M	F	K	A

Original dataset (4 characters)

Taxon	1	1	3	2
A	V	V	R	W
B	V	V	R	W
C	I	i	K	Y
D	M	M	K	F

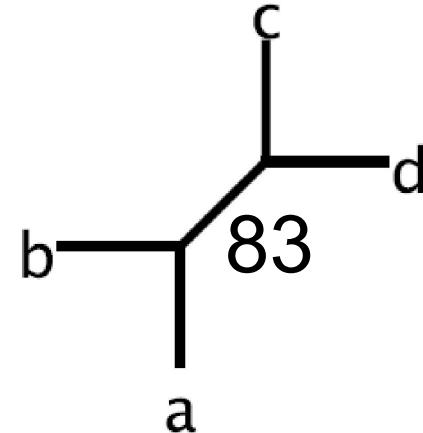
One bootstrap dataset

# Maximum likelihood: assessing support with the bootstrap

Bootstrapping is a technique for estimating the confidence interval around a parameter by **resampling and reanalysis** of the original data (alignment).

1. Sample, with replacement, the same number of characters as the original dataset.
2. Repeat tree inference on each of many bootstrap replicates.
3. Assess confidence as the proportion of samples in which each **split** on the tree appears.

Taxon	1	2	3	4
A	V	W	R	A
B	V	W	R	L
C	I	Y	K	M
D	M	F	K	A



Original dataset (4 characters)

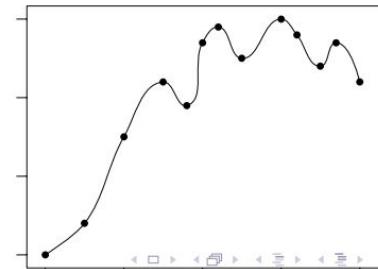
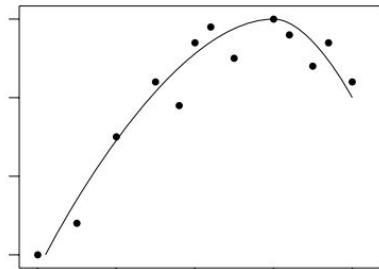
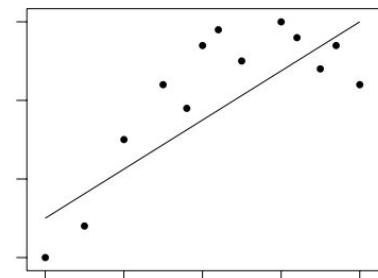
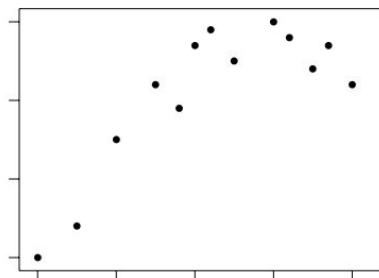
# Maximum likelihood: assessing support with the bootstrap

## Interpreting bootstrap values:

- **\*\*\*The bootstrap is not the probability of a split being correct!**
- Early on, some justifications by comparison to bootstrapping confidence intervals in other settings (**variance of bootstrap trees around ML tree ~ variance of ML trees around true tree**)
- These days, mostly thought of as a way of getting at the robustness of the signal for relationships in the alignment.
- Values of >70 usually taken as being somewhat reliable.

# Maximum likelihood: the problem of model complexity

- Since we always find the optimal values for all parameters in the model, adding parameters will not make the fit to the data worse! (So, can't directly compare the maximum likelihoods)
- We run the risk of overparameterisation (too many, or useless, parameters). This interferes with our ability to generalise/predict from the model!



# Maximum likelihood: the problem of model complexity

To address this problem, likelihood can be penalized by the number of model parameters when comparing models, e.g. with the Akaike Information Criterion:

$$\text{AIC} = 2k - 2\ln(ML)$$

[ $k$  = number of parameters;  $\ln(ML)$  = log of maximum likelihood.

**Choose smallest AIC!**

E.g.

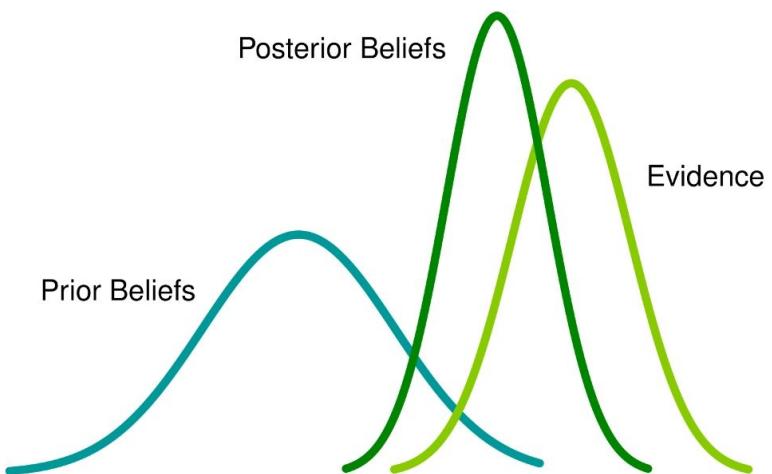
$$\ln(ML) = -17456$$

$$-2\ln(ML) = 34912$$

$$\text{AIC} = 34912 + 2(\text{number of parameters})$$

**Smallest AIC: best ML with fewest parameters.**

# Bayesian inference



Conceptualise the analysis as an **updating** of our prior beliefs, based on new **evidence**.

Our beliefs about parameters are expressed as **distributions**.

The **evidence = likelihood** of the data (as in ML).

This seems to follow human reasoning.

# Bayes theorem

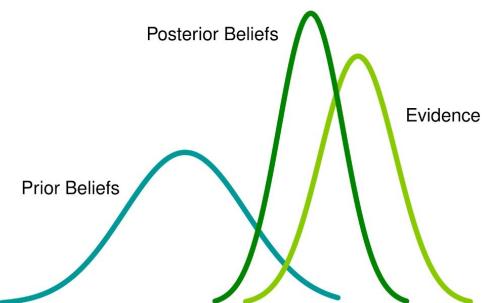
$$P(H|D) = \frac{P(H) P(D|H)}{P(D)}$$

**P(H|D)**: Probability of Hypothesis (tree, model...) given Data (alignment)

**P(H)**: Prior probability of Hypothesis

**P(D|H)**: Probability of Data given Hypothesis  
**(the likelihood!)**

**P(D)**: Prior probability of the Data



# Prior beliefs are a central part of logical reasoning

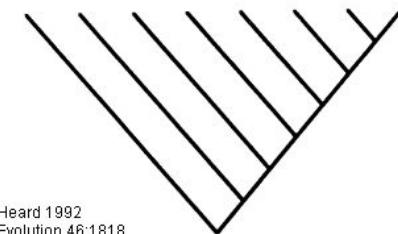
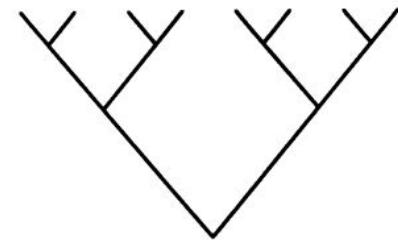
$$P(H|D) = \frac{P(H) P(D|H)}{P(D)}$$

- 1 out of 100 people in the population at large have disease X.
- A **diagnostic test** exists: it has a false negative rate of 0% and a false positive rate of 10%.
- As part of a routine checkup, Billy takes the test and it is positive.
- What is the probability that Billy has disease X?



# How can we have prior beliefs about e.g. the tree?

- We do: e.g. which of these trees is more probable, all else being equal?
- We usually want to use “weak” (flat) priors
- Some priors are fairly strong (e.g. on branch lengths)
- Some priors are extremely strong (e.g. node ages in molecular clock)



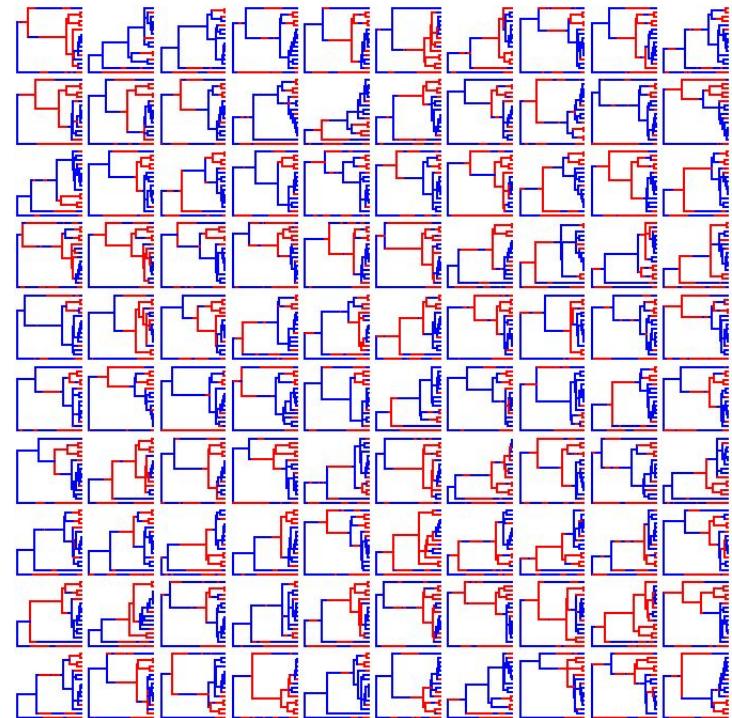
Heard 1992  
Evolution 46:1818

# Indisputable superiority of Bayesian approach (1): quantifying support

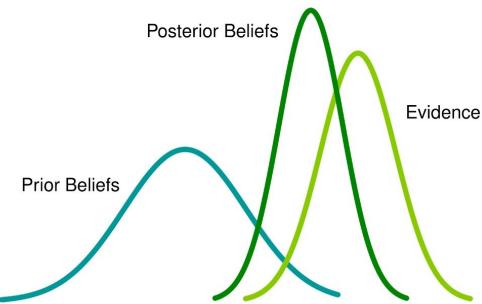
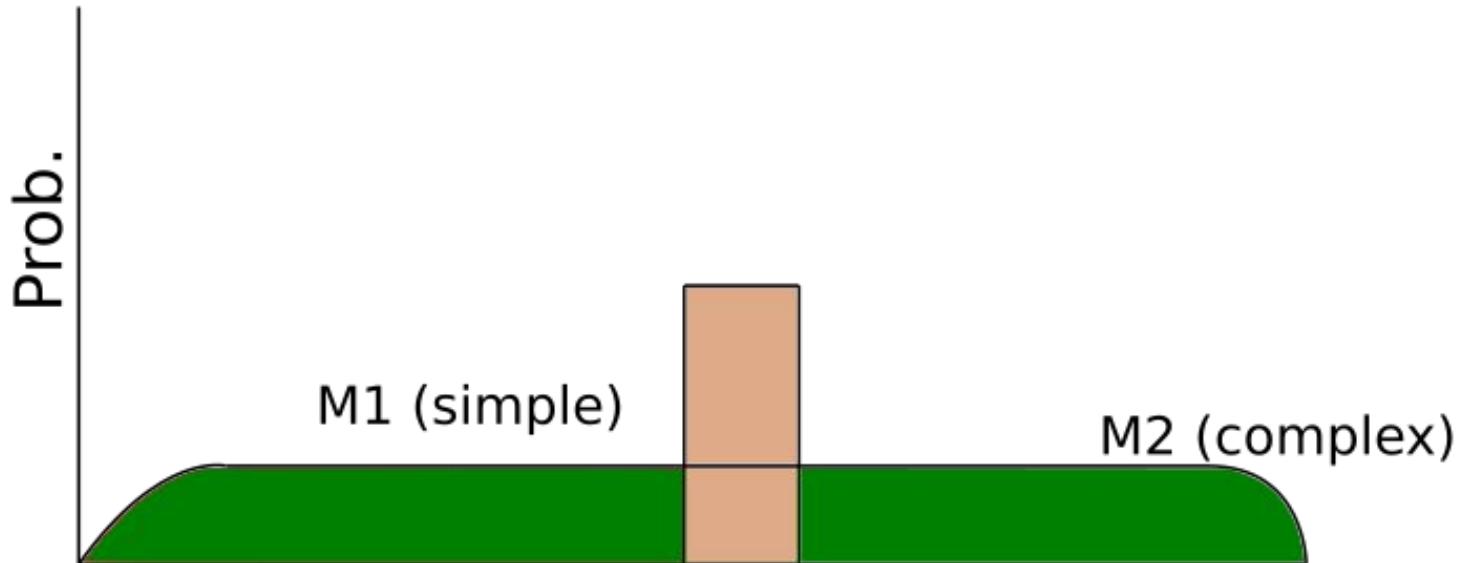
We have a prior distribution on trees, updated by the evidence, to give a posterior **distribution** of trees: a natural measure of uncertainty.

No need for bootstrapping.

Has a direct interpretation: the **frequency of a tree (or a split) in the posterior distribution is the probability that it is true**, under the model.

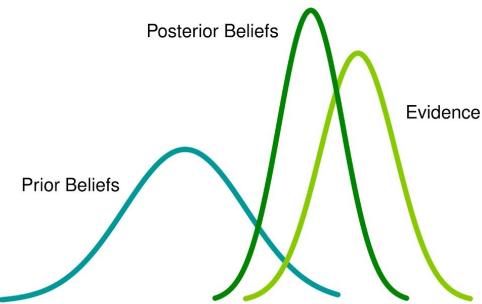
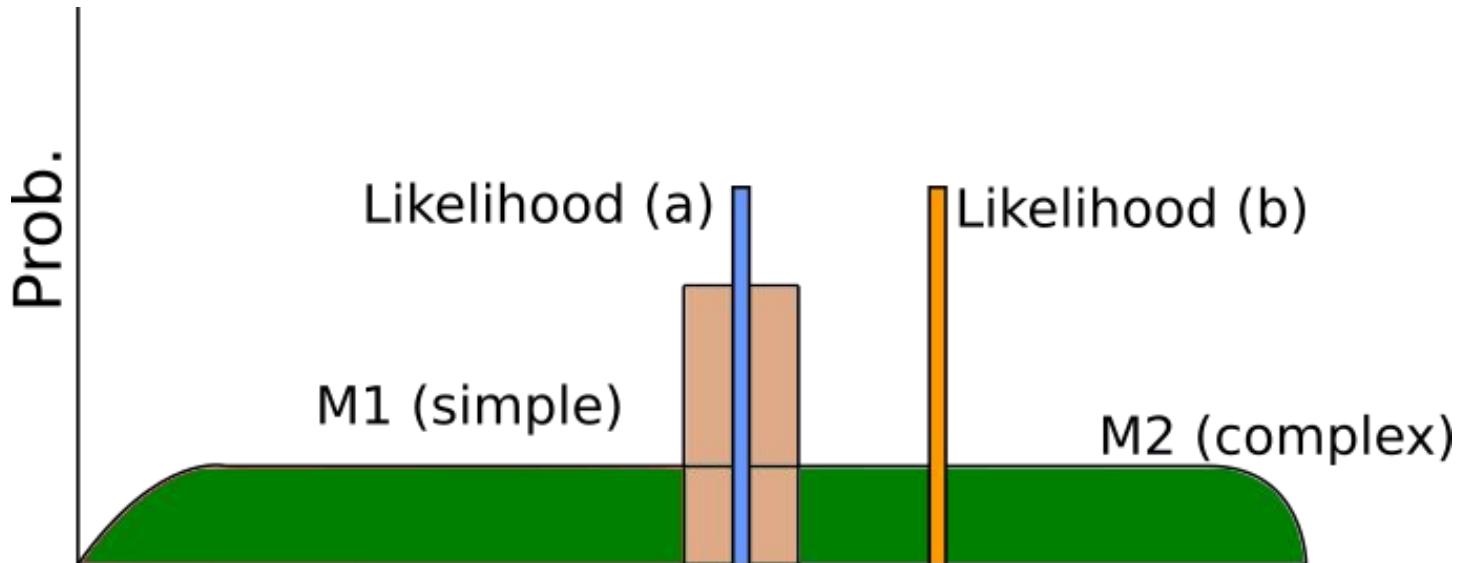


## Indisputable superiority of Bayesian approach (2): natural incorporation of model complexity



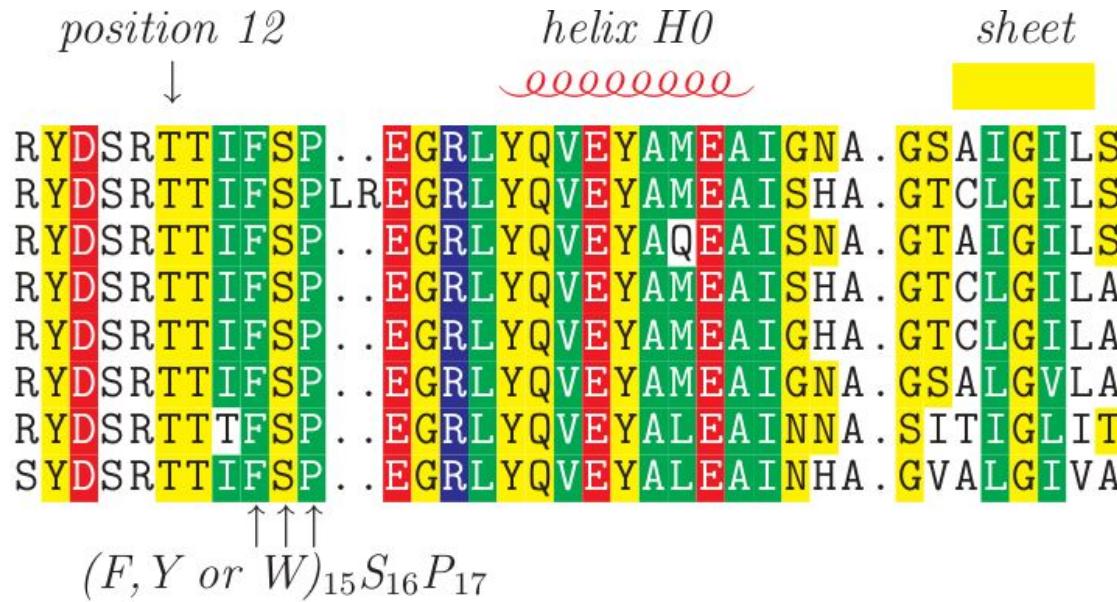
**Complex models predict a wider range of outcomes *a priori*, but  $P(\text{outcomes})$  must sum to 1.**

## Indisputable superiority of Bayesian approach (2): natural incorporation of model complexity



**Complex models predict a wider range of outcomes *a priori*, but  $P(\text{outcomes})$  must sum to 1!**  
**Priors provide a natural penalty for model complexity.**

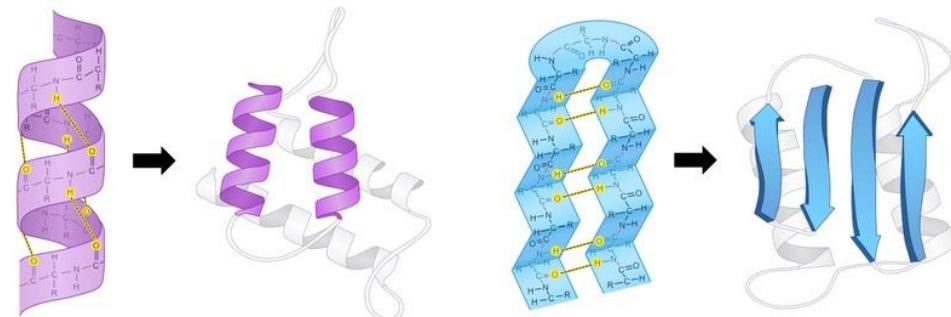
# A complex model: site-specific sequence compositions



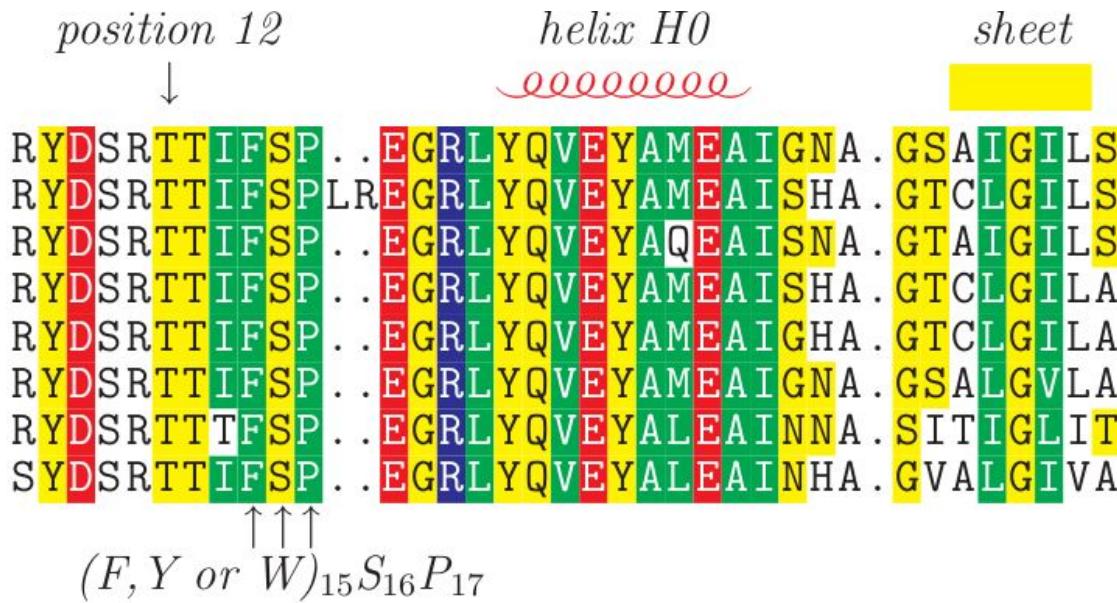
**How the hell to model this?**

How many different compositions are there?

How are sites allocated to compositions?



# The Dirichlet process: a prior distribution on discrete distributions



We can draw sets of composition “centroids” (vectors of e.g. 20 proportions) from a Dirichlet process prior.

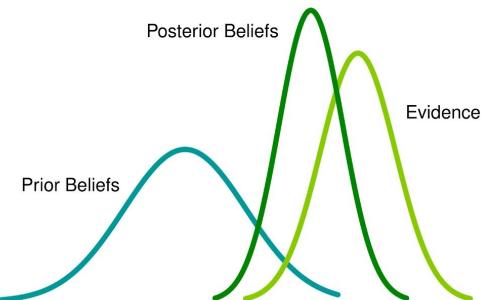
- The result is an **\*\*\*infinite mixture model\*\*\***: we average over an unknown, but possibly very large, number of different compositions!
- This will usually have uncountably many parameters: can't compute BIC or similar (but prior puts constraints on the draws via a concentration parameter).
- **Usually fits the real data very well.** Can't be done with ML.

# Sampling from $P(H|D)$ using MCMC

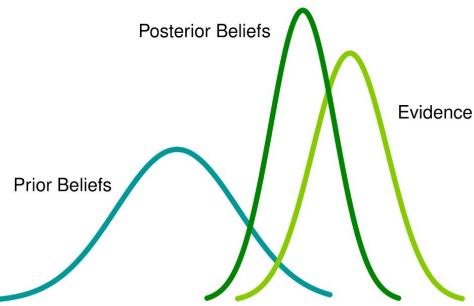
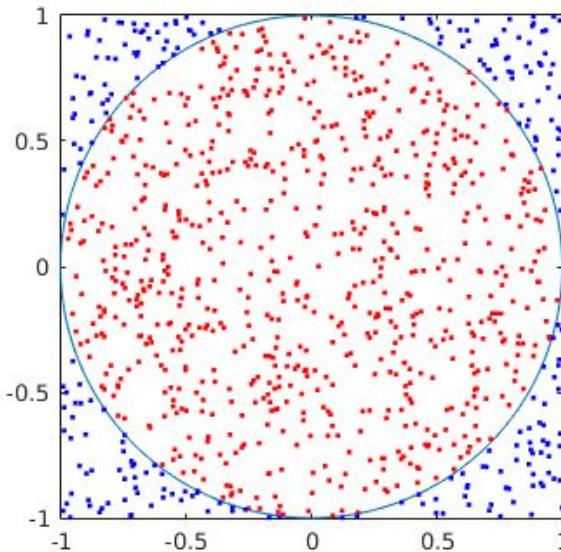
$$P(H|D) = \frac{P(H) P(D|H)}{P(D)}$$

Uh oh!

- We can't analytically find  $P(H|D)$  in the phylogenetic case because  $P(D)$  is rather difficult to calculate (average over all trees, branch lengths, other model parameters...)
- We use numerical methods (Markov Chain Monte Carlo) to sample from  $P(H|D)$ ; with enough samples, we obtain a reasonable approximation of  $P(H|D)$ .

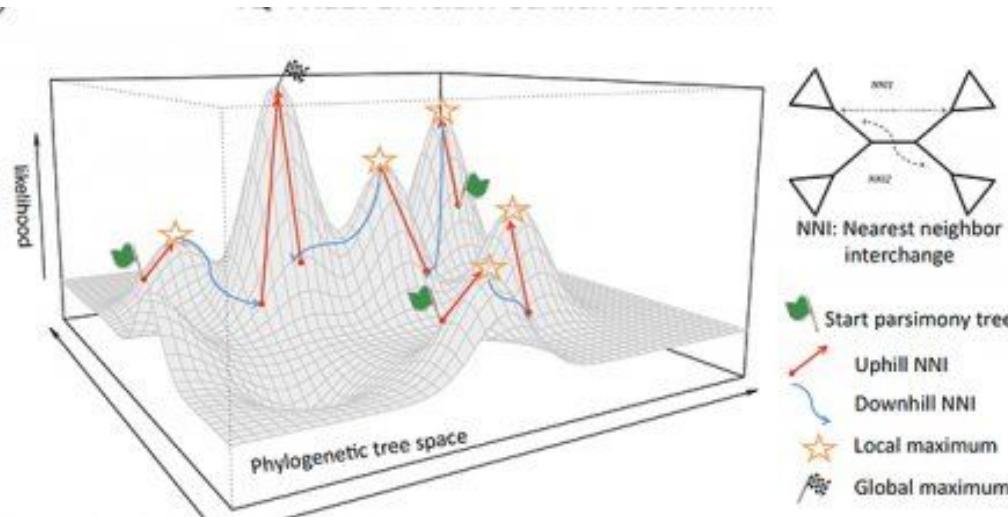


# Sampling from $P(H|D)$ using MCMC

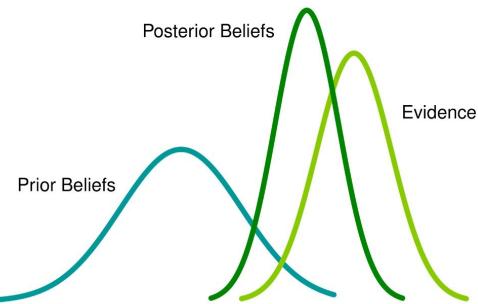


How could we estimate the area of this circle, without knowing much about geometry?

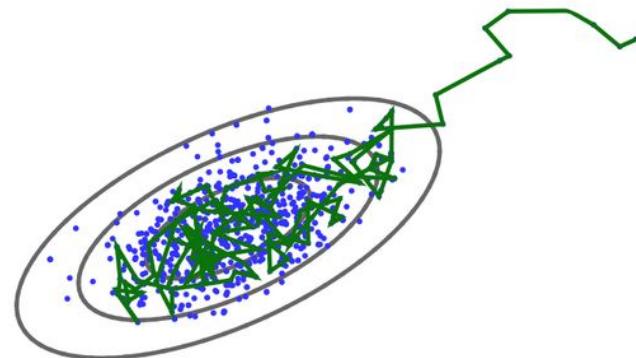
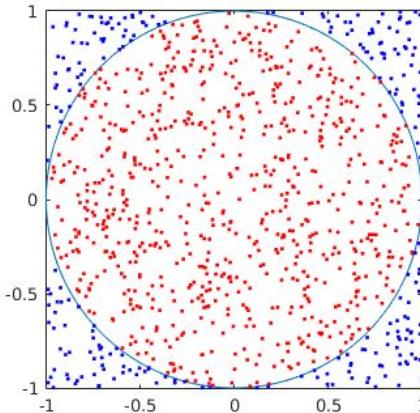
# Sampling from $P(H|D)$ using MCMC



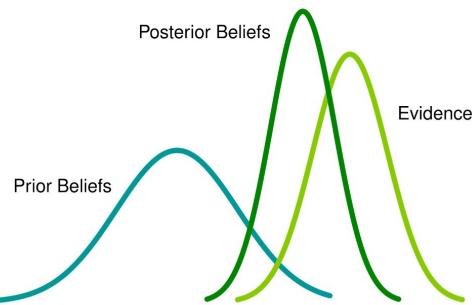
What about a more complex, multidimensional shape (like a probability distribution!)?



# Sampling from $P(H|D)$ using MCMC



- We draw samples from  $P(H|D)$  using a biased random walk, such that histograms of the samples/set of trees drawn will approximate the distribution.
- We need to sample in a biased way because  $P(H|D)$  has many dimensions, and the region(s) of credible values with  $P > 0$  is quite small.
- Can mostly sample from high-probability regions by biasing the “walk” using the ratio of  $P(D|H)P(H)$  at the current and proposed region to decide whether to jump or not.
- **NB: We sample to estimate the whole density: NOT just the maximum likelihood!**



# Analysis with ML and Bayesian methods: practical considerations

(Methods have a lot in common; either is often OK, and better than alternatives).

## Maximum likelihood:

- + Very efficient implementations (scale to 100s or 1000s of sequences, 10,000s of alignment sites)
- + Fitting the models usually a bit less involved (but: local maxima)
- Support values hard to interpret
- Limited to simple models

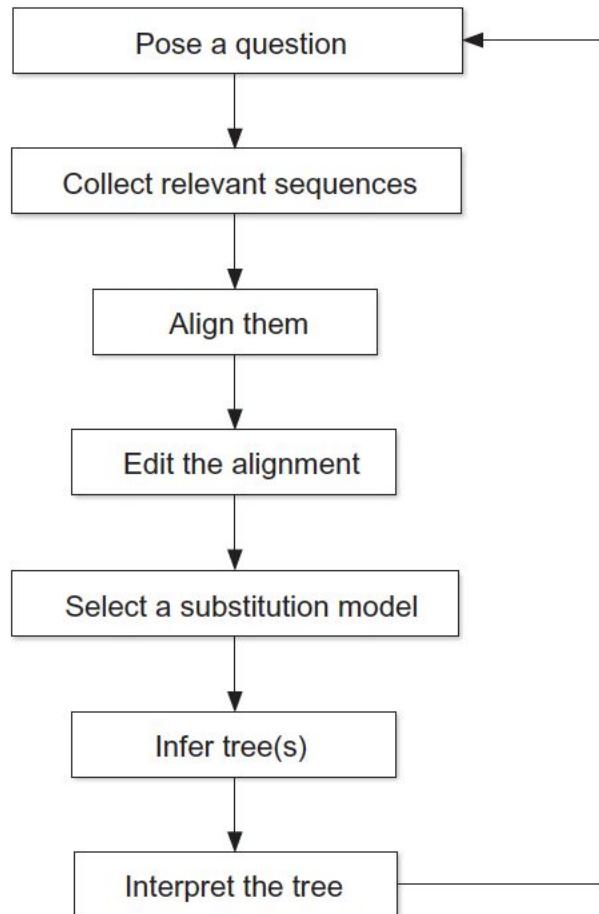
## Bayesian:

- + A principled framework for inference (...)
- + Natural way to fit more complex models (best-fitting models are all Bayesian-only)
- + Natural support values
- Getting MCMC to work can be difficult
- Don't scale to very large datasets

# Fitting phylogenetic models: summary

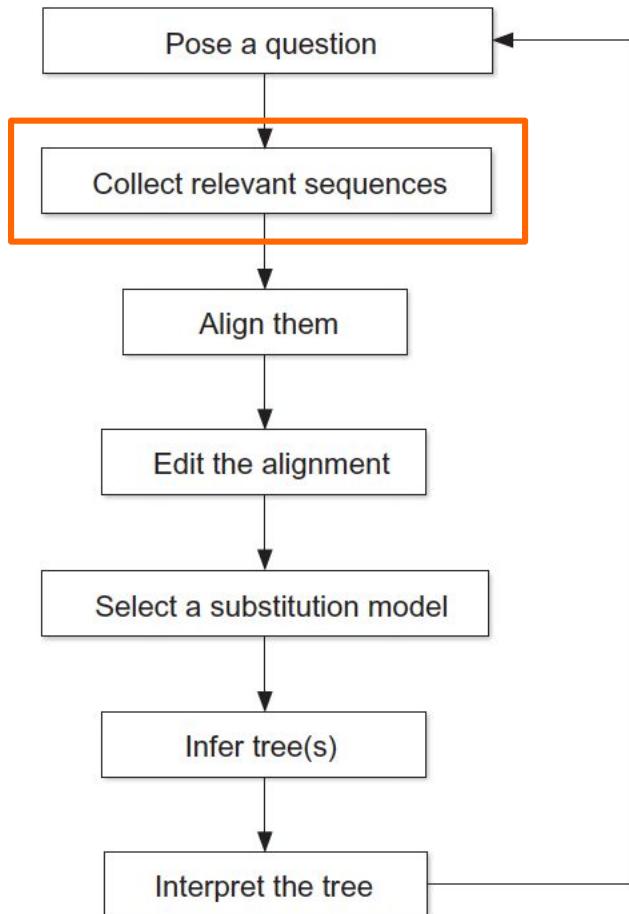
- Phylogenetic trees are statistical inferences
- The inferences are based on a model: it's important, so do experiments!
- Quantifying the uncertainty in our inferences is critically important for interpretation
- ML & Bayesian methods both provide a framework for assessing support

# Molecular phylogenetics: a possible flowchart



**A critical approach to analysis is key: don't fall into the black box!**

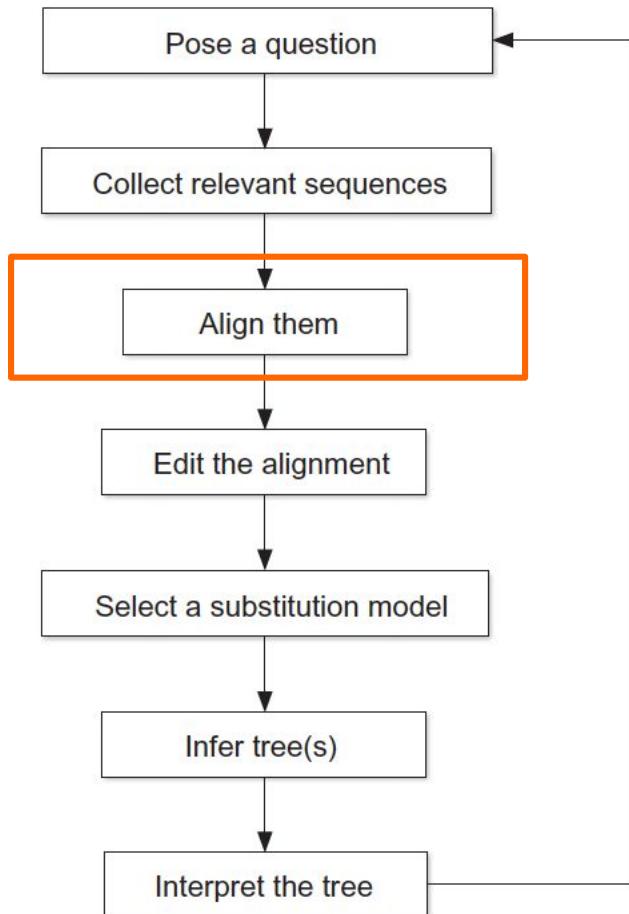
# Molecular phylogenetics: a possible flowchart



Did I miss any  
important ones?

A critical approach to analysis is key: don't fall into the black box!

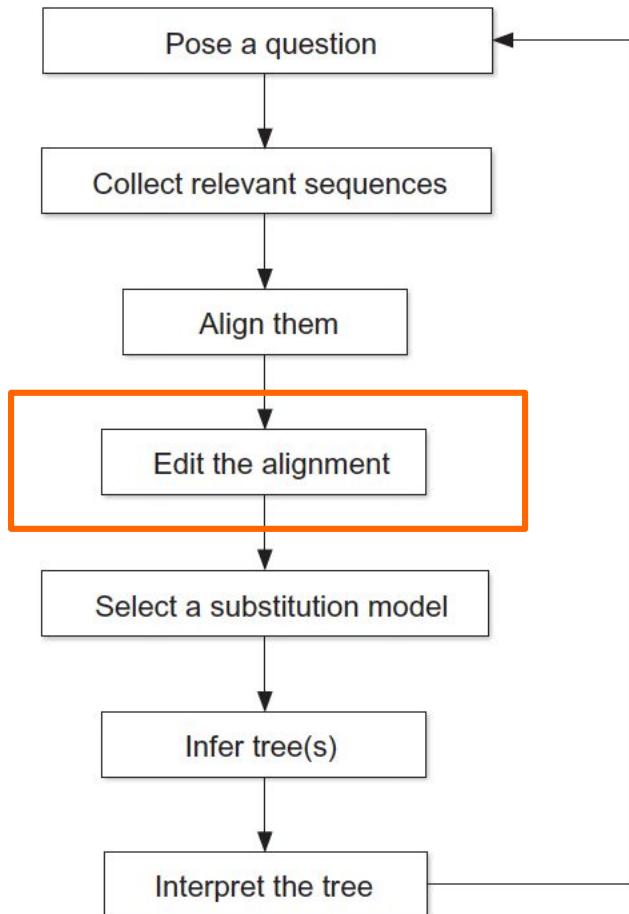
# Molecular phylogenetics: a possible flowchart



How certain is my alignment?  
Would other alignments give a different tree?

A critical approach to analysis is key: don't fall into the black box!

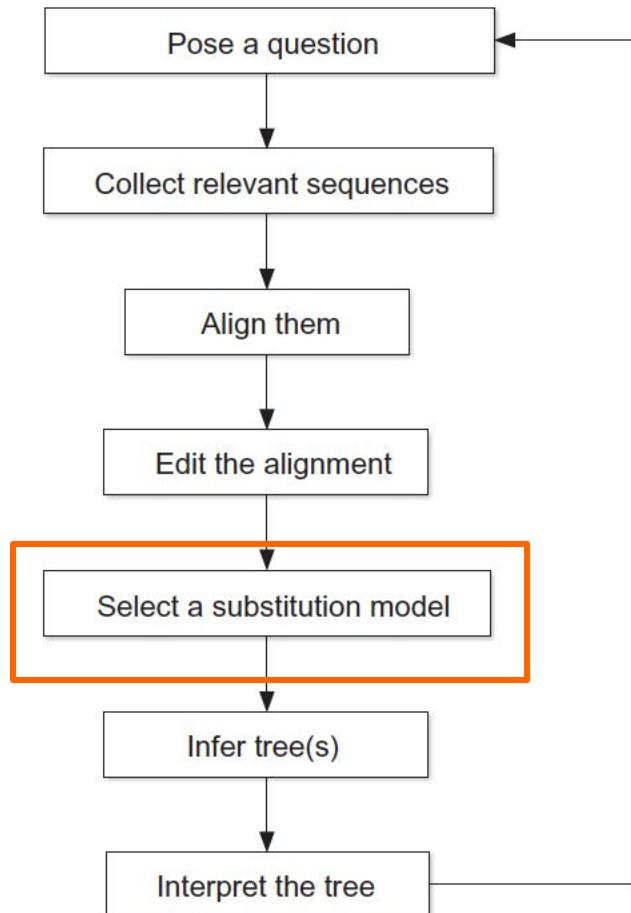
# Molecular phylogenetics: a possible flowchart



Which bits, if any,  
should I remove?

A critical approach to analysis is key: don't fall into the black box!

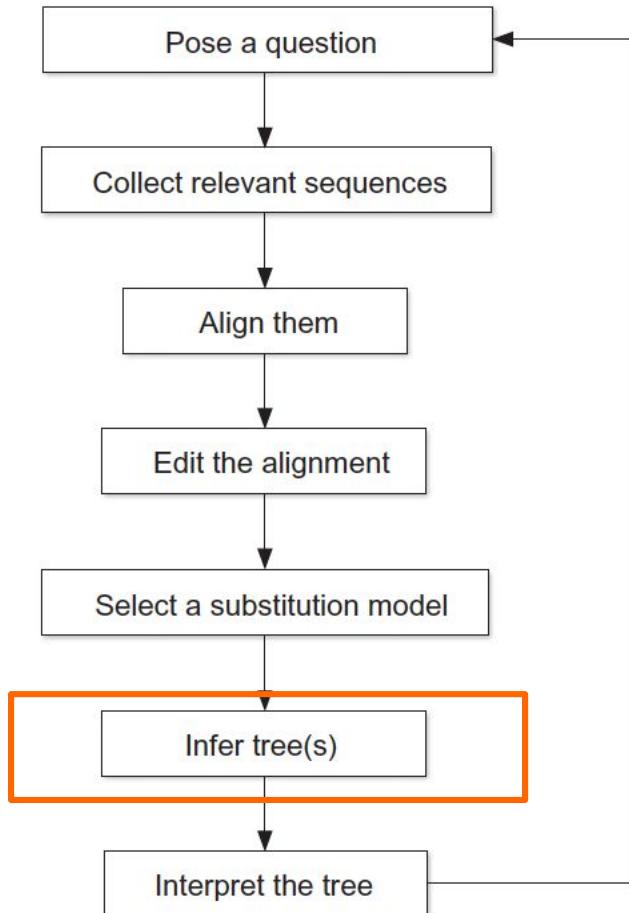
# Molecular phylogenetics: a possible flowchart



Is this the best  
available model?  
Is it good enough?

A critical approach to analysis is key: don't fall into the black box!

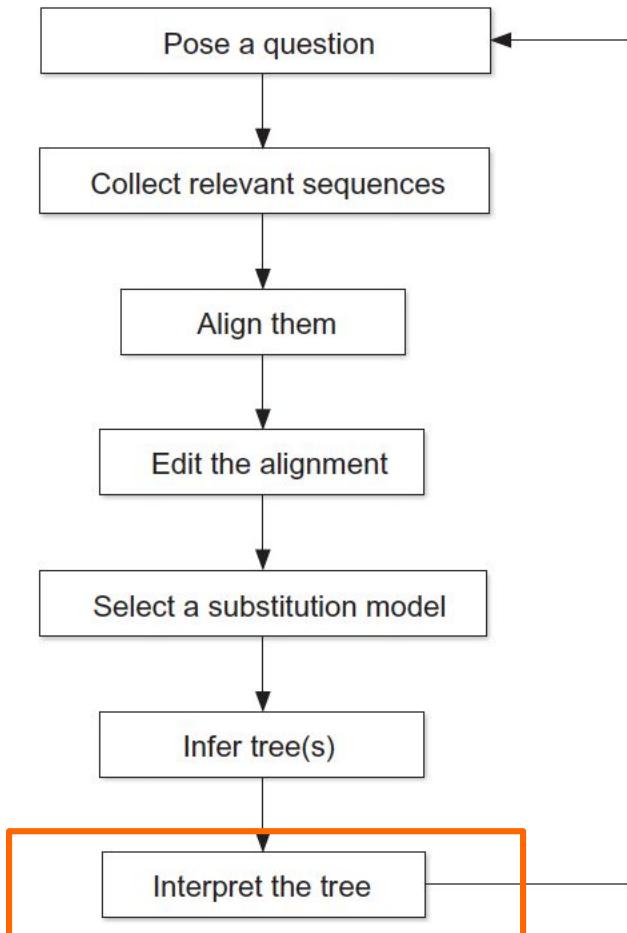
# Molecular phylogenetics: a possible flowchart



Are my ML or Bayesian analyses working properly? (Local maxima, convergence)

A critical approach to analysis is key: don't fall into the black box!

# Molecular phylogenetics: a possible flowchart



How well-supported  
are the key branches?

What does the tree  
mean for my  
hypothesis?

Can alternative  
hypotheses be  
rejected?

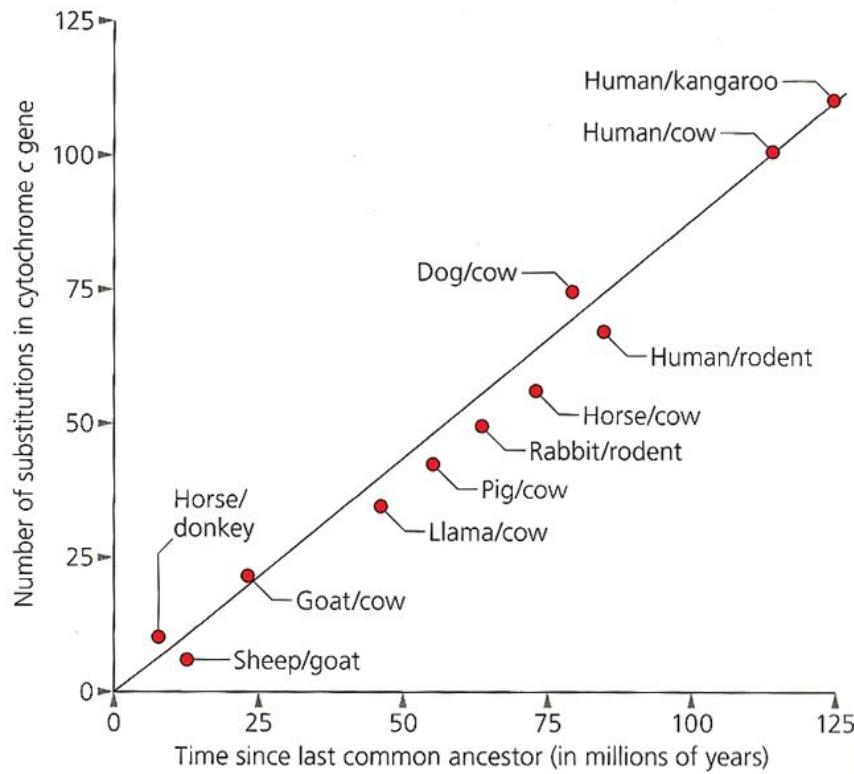
A critical approach to analysis is key: don't fall into the black box!

# **(Brief, crappy) introduction to the molecular clock**

# The molecular clock

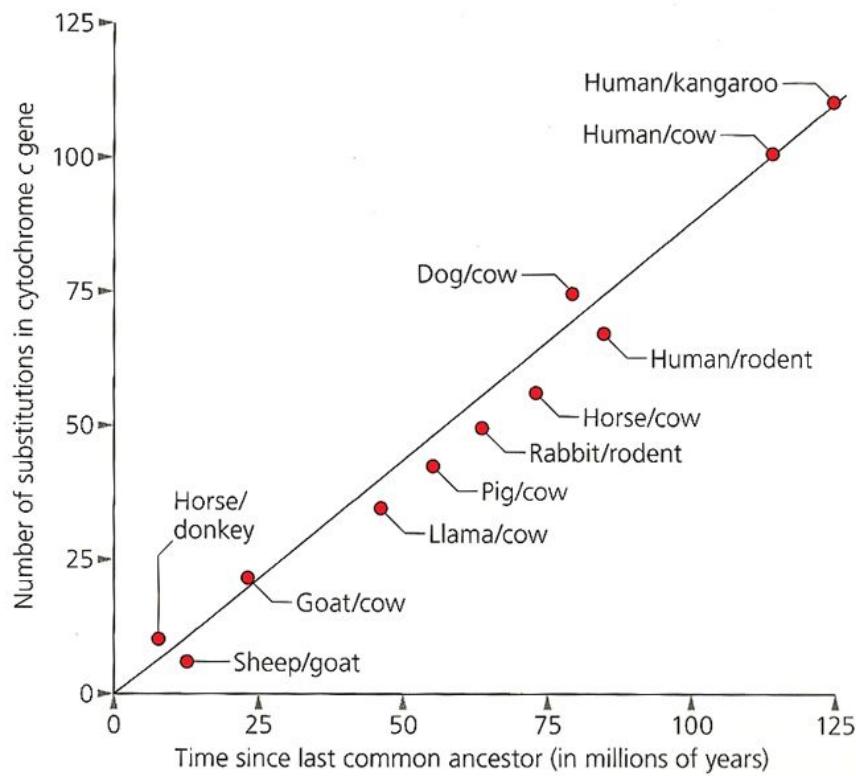
- **Genetic differences accumulate over time**
- **The amount of divergence between species is roughly proportional to the age of their last common ancestor**

# The molecular clock



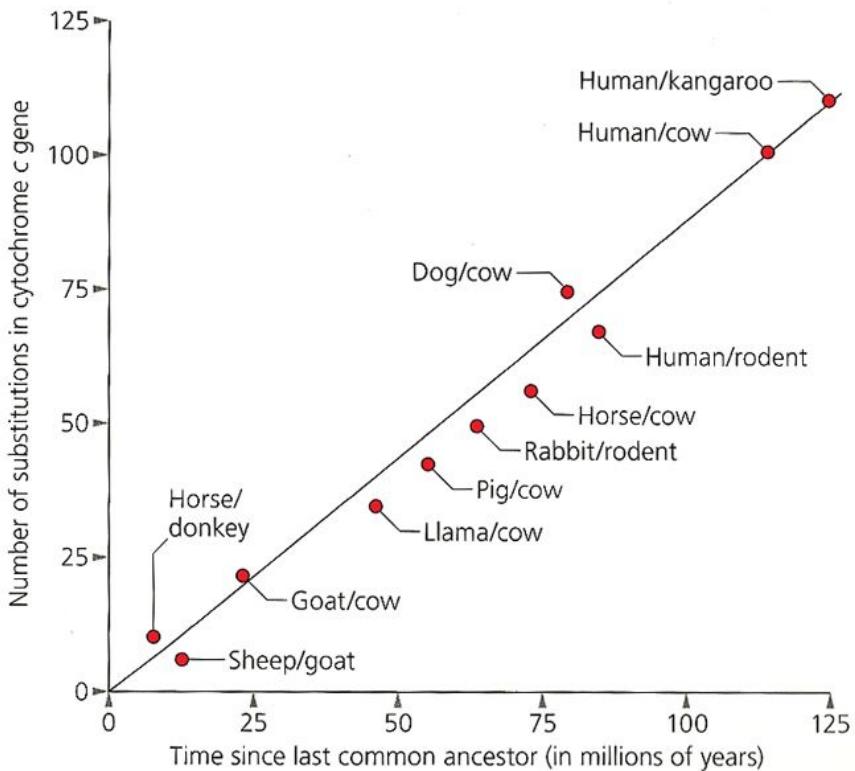
**Figure 7.8** DNA mutates at a roughly clock-like rate. This graph shows how distantly related pairs of species have a large number of different substitutions in the cytochrome c gene. (Adapted from Moore and Moore, 2004)

# Quiz: why is there a molecular clock?



**Figure 7.8** DNA mutates at a roughly clock-like rate. This graph shows how distantly related pairs of species have a large number of different substitutions in the cytochrome c gene. (Adapted from Moore and Moore, 2004)

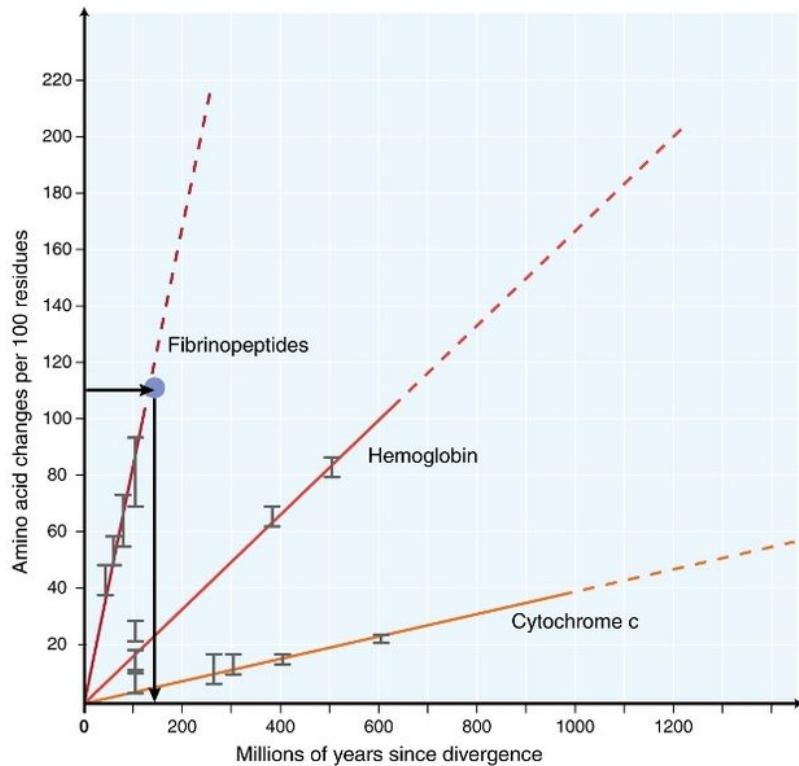
# Quiz: why is there a molecular clock?



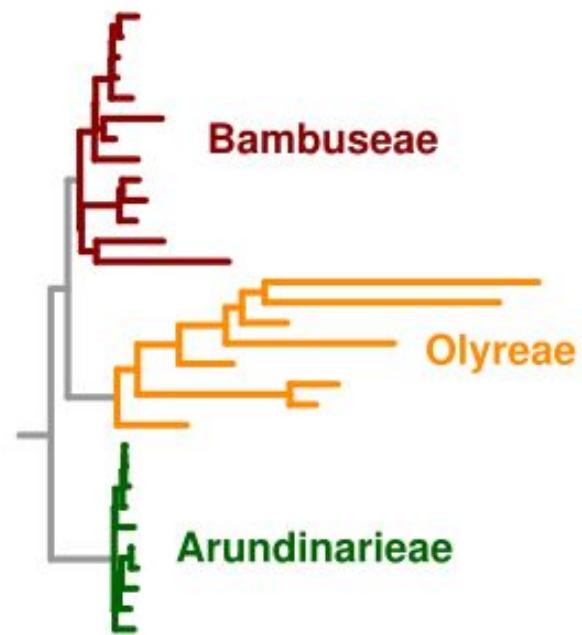
**Figure 7.8** DNA mutates at a roughly clock-like rate. This graph shows how distantly related pairs of species have a large number of different substitutions in the cytochrome c gene. (Adapted from Moore and Moore, 2004)

- **Selective neutrality (substitution rate = mutation rate)**
- **Pragmatic, long-term average?**

# Variation in the molecular clock



Across genes

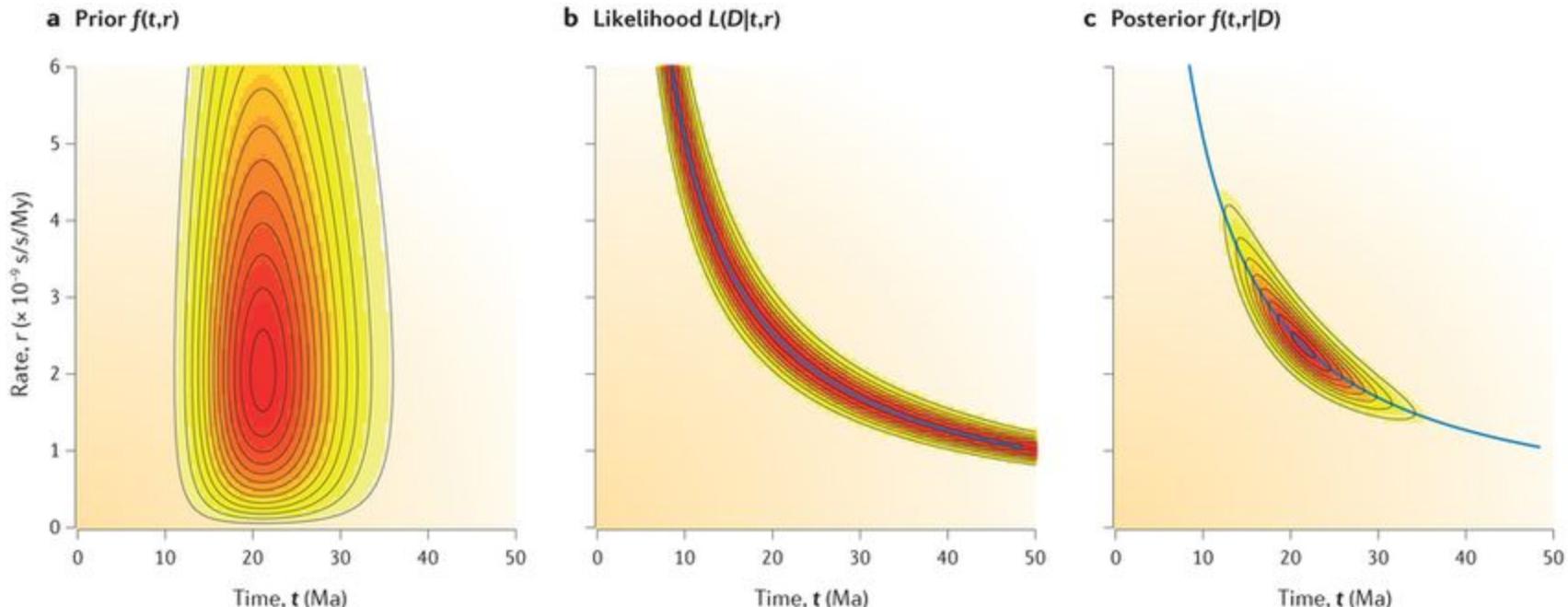


Across species

# No strict clock, but just relax...

- Clock ticks at different rates among:
  - sites in a molecule
  - genes
  - regions of genomes
  - genomes in the same cell
  - taxonomic groups
- Bayesian methods allow these assumptions to be relaxed, some dating information to be extracted.

# Inference under a relaxed molecular clock



Nature Reviews | Genetics

A very natural way to combine fossil calibrations and sequence data.

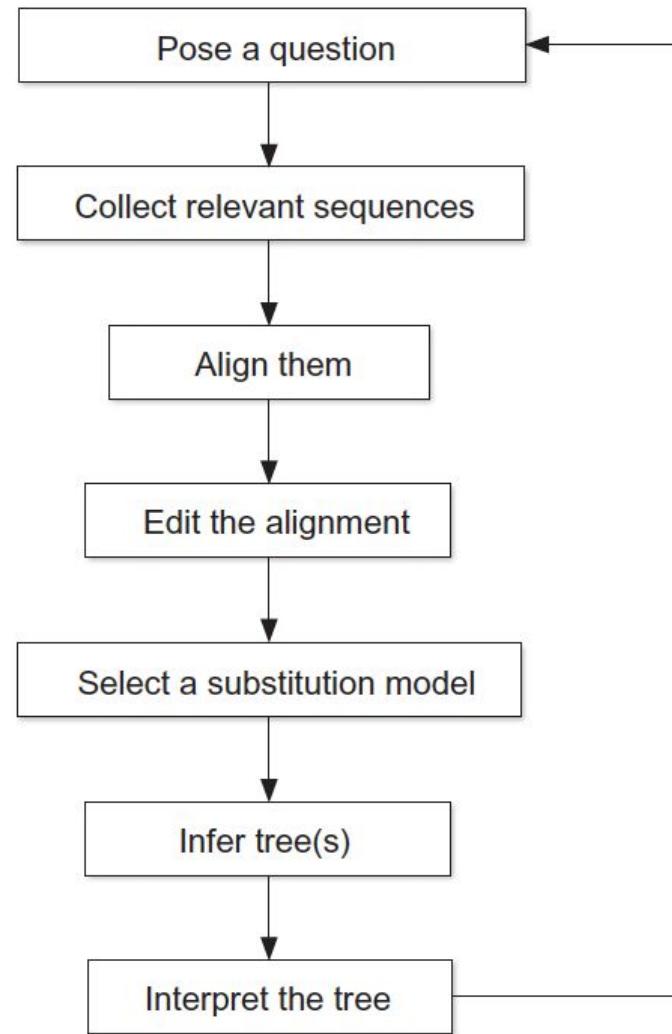
But, note the different role of the prior!

# Phylogenetics: some key questions

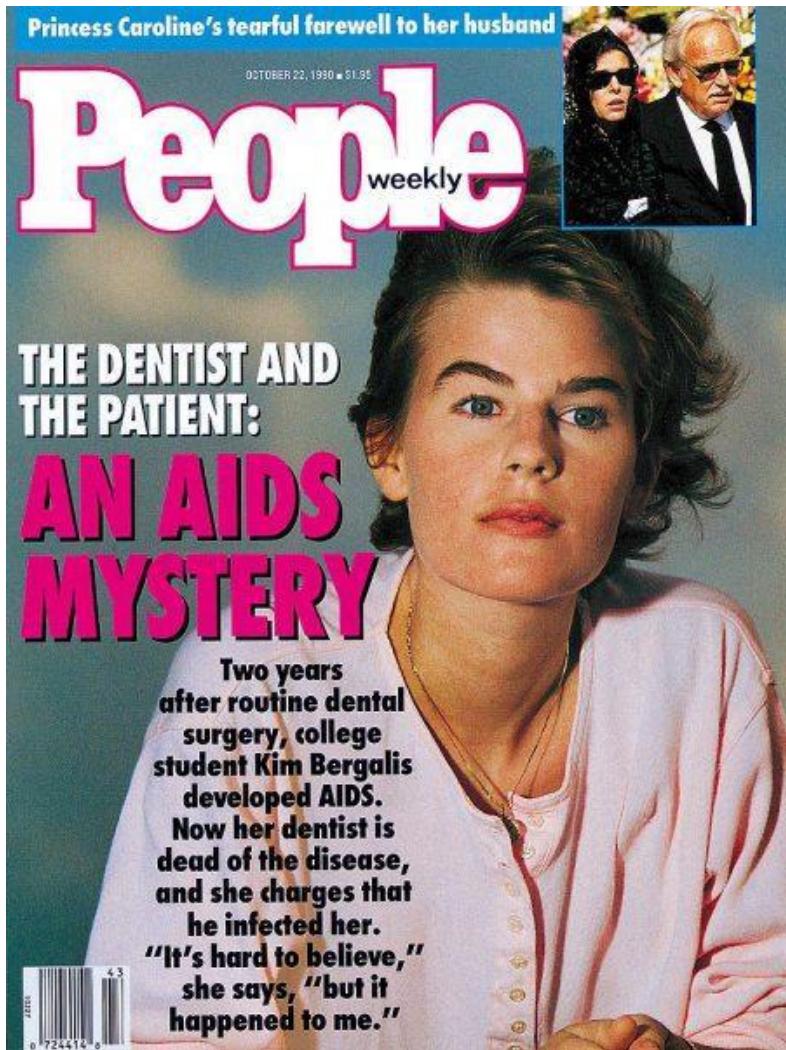
- Did I analyse all the relevant data?
- Does the model fit the data?
- How well-resolved is the tree? Can alternative hypotheses be rejected?
- What do these analyses mean for my question?

# Molecular phylogenetics practical

# Practical: A first phylogenetic analysis



# Dentist-patient transmission of HIV?



In the early 90s, a dentist was accused of infecting several of his patients with HIV during surgical procedures.

After a “low-risk” patient was diagnosed with HIV, other patients were screened. 10 had HIV.

**Did the dentist infect them? We will do a phylogenetic analysis to evaluate these claims.**

# The data

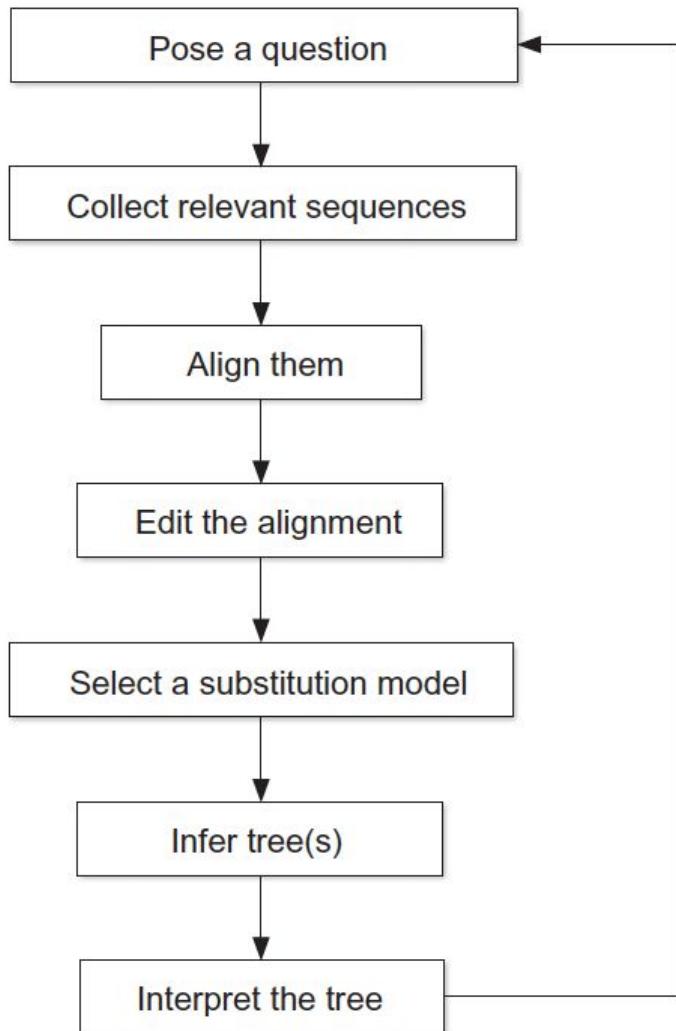
## **Molecular Epidemiology of HIV Transmission in a Dental Practice**

Chin-Yih Ou, Carol A. Ciesielski, Gerald Myers,  
Claudiu I. Bandea, Chi-Cheng Luo, Bette T. M. Korber,  
James I. Mullins, Gerald Schochetman, Ruth L. Berkelman,  
A. Nikki Economou, John J. Witte, Lawrence J. Furman,  
Glen A. Satten, Kersti A. MacInnes, James W. Curran,  
Harold W. Jaffe, Laboratory Investigation Group,\*  
Epidemiologic Investigation Group†

HIV *env* sequences from:

- The dentist
- His patients
- Local controls

# The workflow



**Q: Are patient sequences  
*descended* from the dentist's  
sequences?**

Let's use NCBI...

We'll use mafft.

\*\*\*

We'll use the best-fitting model in  
IQ-Tree...

Does the phylogeny support the claims?

# How to retrieve the data?

NCBI Resources How To

Nucleotide Nucleotide ou[au] ciesielski[au] V3 Create alert Advanced

Species Summary ▾ 20 per page ▾ Sort by Default order ▾  
Viruses (134)  
Customize ...

Molecule types Items: 1 to 20 of 134  
genomic DNA/RNA (134) << First < Prev Page 1 of 7 Next  
Customize ...

Source databases 1. 300 bp linear RNA  
INSDC (GenBank) (134) Accession: M90923.1 GI: 327315  
Customize ... GenBank FASTA Graphics

Sequence length 2. 327 bp linear RNA  
Custom range... Accession: M90927.1 GI: 327313  
GenBank FASTA Graphics

Release date 3. 327 bp linear RNA  
Custom range... Accession: M90926.1 GI: 327311  
GenBank FASTA Graphics

Revision date 4. 327 bp linear RNA  
Custom range... Accession: M90925.1 GI: 327309  
GenBank FASTA Graphics

[Clear all](#)

[Show additional filters](#)

# How to retrieve the data?

- **Dentist sequences: FLD1,2,4,5,7,8**
- **Patient sequences: FLP[A-H]: take 3 isolates/patient**
- **Local controls (LCXX, FLQ): take at least 10.**

# The workflow

- Download the sequences in FASTA format, save in a text file.
- Rename with sensible tags
- FASTA format:



A screenshot of a text editor window titled "sequence.fasta". The window has an "Open" button and a "+" button. The file content is a FASTA sequence with the following header and sequence:  
>FLD2\_M90849.1  
CTAGCAGAAGAAGAGATAGTAATTAGATCTGCCAATT  
>FLPB9\_M90870.1  
GGAGATATAAGACAAGCACATTGTAACATTAGTAGAG/  
>FLPBR3A\_M92115.1  
TGTACAAGACCCAACAACAATACAAGAAAAGGTATAC/  
>FLPH1D\_M90907.1  
CTAGCAGAAGGAGAGGTAATAATTAGATCTGAAAATT  
>FLPG6\_M90905.1  
GAAGAGGTTAGTAATTAGATCTGCCAATTACAGACAC/  
>FLDD\_M90847.1  
GAGGTAGTAATTAGATCTGCCAATTACAGACAAATGC

# Some useful commands

**Alignment:**

```
mafft --auto mySequences >  
myAlignment
```

**View the alignment in belvu/Jalview.**

**Tree inference (black box):**

```
iqtree -s myAlignment -m MFP -bb  
1000
```

# Questions

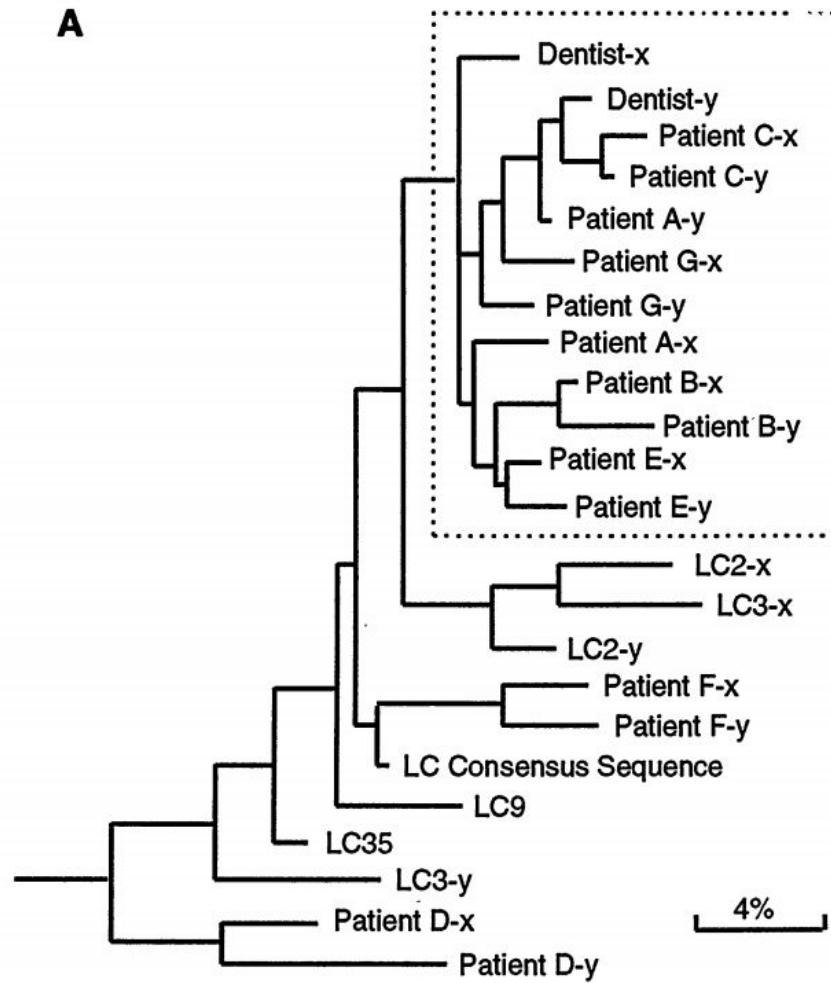
- 1. Did the patients get their viruses from the dentist?**
- 2. How much would you bet on it?**
- 3. Why include local controls in the analysis?**

**For a bonus practical, see:**

[https://github.com/Tancata/bioinf\\_workshop](https://github.com/Tancata/bioinf_workshop)

# Ou et al. (1992)

A



# Some properties of the patients

**Table 1.** Dental cohort clinical information and HIV nucleotide variation in the C2-V3 domain of the envelope gene.

Person	Known risk factor			M13 clones (no.)	Intraperson variation† (%)	Interperson variation† (%)	
	Sex	Clinical status*				To dentist	To 30 LCs‡
Dentist	M	Yes	AIDS	6	3.3 (0.8–5.4)		11.0 (5.8–16.0)
Patient A	F	No	AIDS	6	2.0 (0.0–4.5)	3.4 (0.8–6.2)	10.9 (5.4–14.8)
Patient B	F	No	Asymptomatic (CD4 = 222/ $\mu$ l)	12	1.9 (0.4–3.7)	4.4 (2.1–7.0)	11.2 (6.2–16.5)
Patient C	M	No§	Asymptomatic (CD4 = <50/ $\mu$ l)	5	1.2 (0.4–1.6)	3.4 (2.1–4.9)	11.1 (7.0–15.6)
Patient E	F	No	Asymptomatic (CD4 = 567/ $\mu$ l)	6	2.1 (0.4–3.7)	3.4 (1.2–6.6)	10.8 (5.8–14.8)
Patient G	M	No	Asymptomatic (CD4 = 400/ $\mu$ l)	5	2.8 (1.6–3.7)	4.9 (2.9–7.0)	11.8 (6.2–16.9)
Patient D	M	Yes	AIDS	5	7.5 (0.0–9.9)	13.6 (11.5–15.6)	13.1 (7.8–17.3)
Patient F	M	Yes	Asymptomatic (CD4 = 253/ $\mu$ l)	6	3.0 (0.8–5.8)	10.7 (8.2–13.6)	11.9 (7.0–17.3)

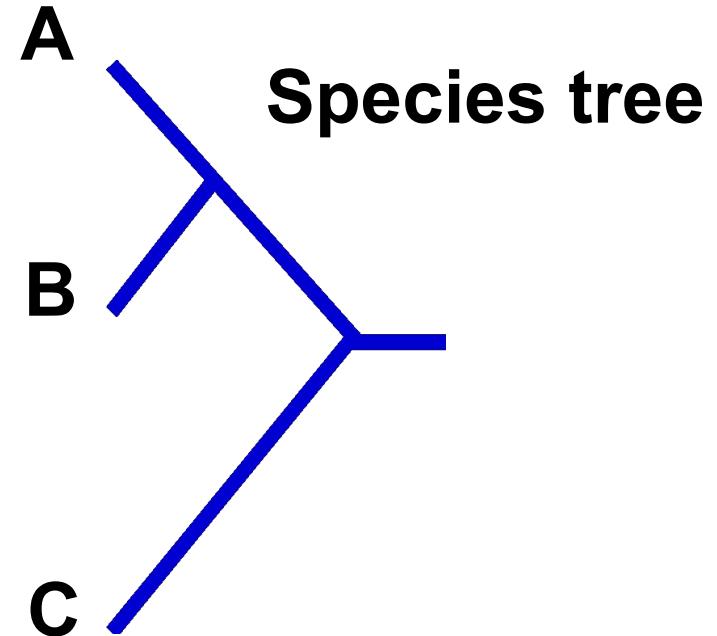
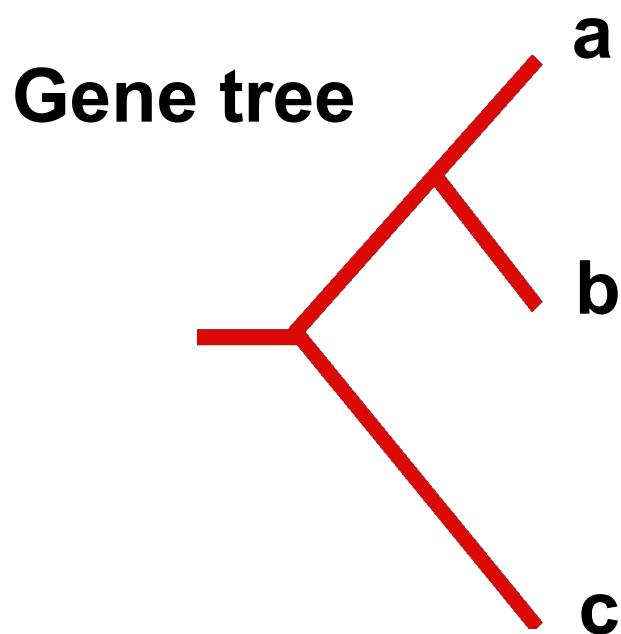
# Introduction to phylogenomics

# Phylogenetic analysis of genome data

- These days, data collection is at the genome scale (genomes, transcriptomes, genome-wide variation): **not just one “marker” gene, but all of them**
- Far more evolutionary information than from single genes, but we need to consider additional evolutionary processes (gene duplications, etc.)
- More data is better, but only if it is analysed properly!



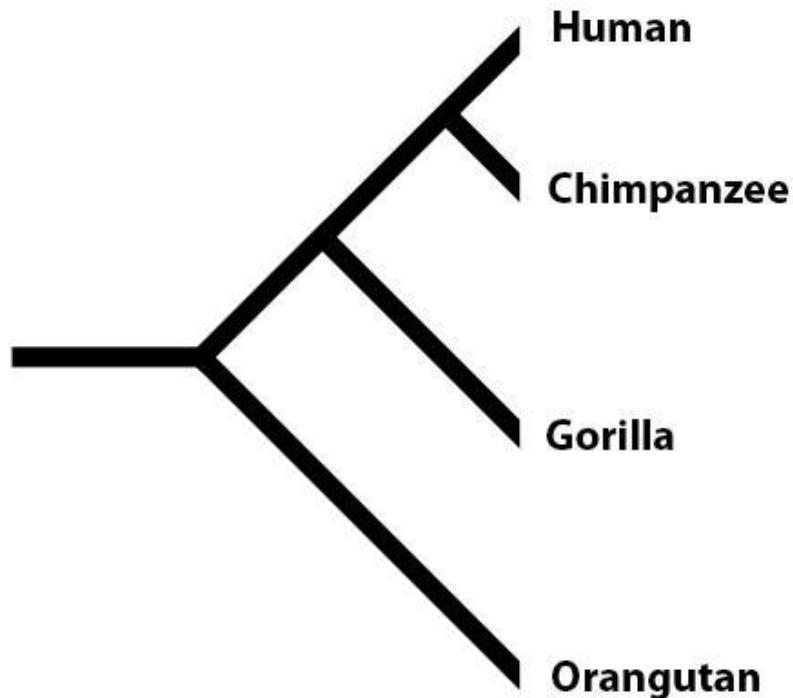
# Gene trees and species trees



We often assume that gene trees give us species trees:  
but different gene trees often disagree!

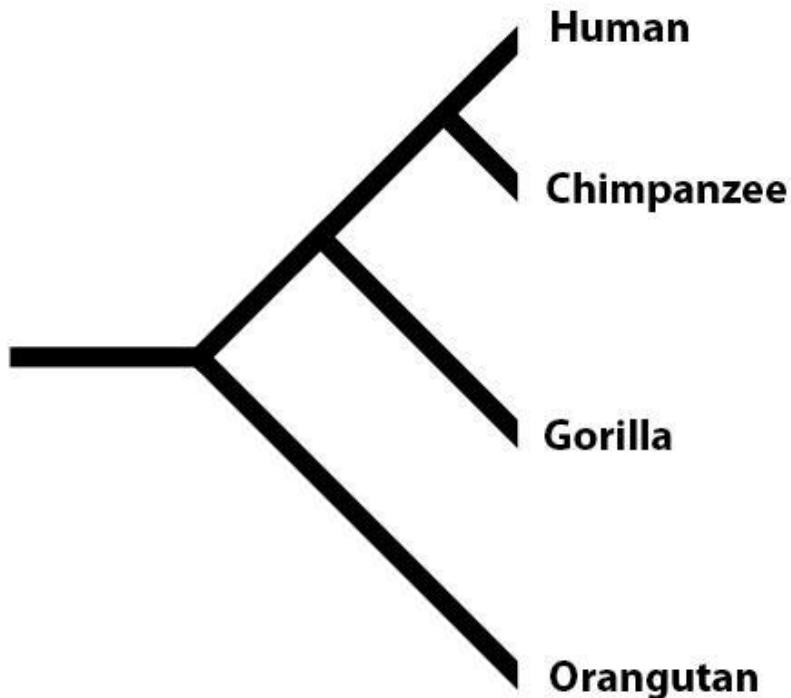
(and that's OK...)

# Human-Chimp-Gorilla



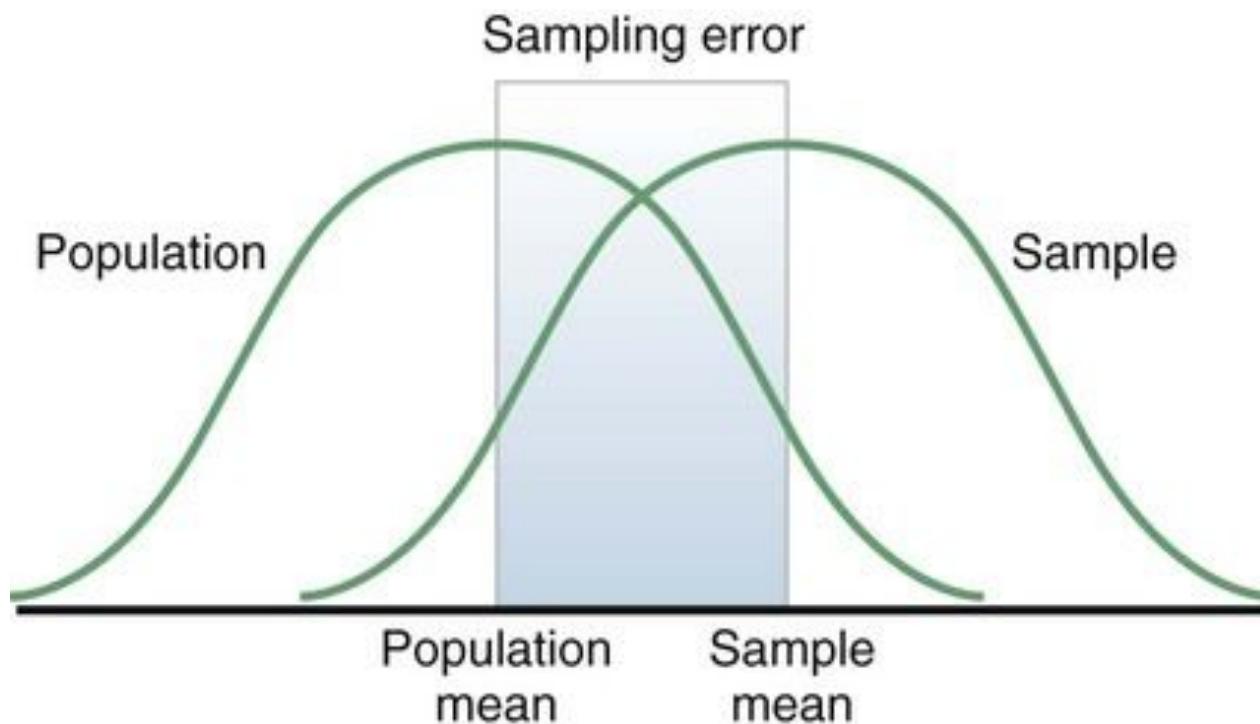
- ~15% Human genome more similar to Gorilla than Chimp
- ~15% Chimp genome more similar to Gorilla than Human

# Quiz: why?

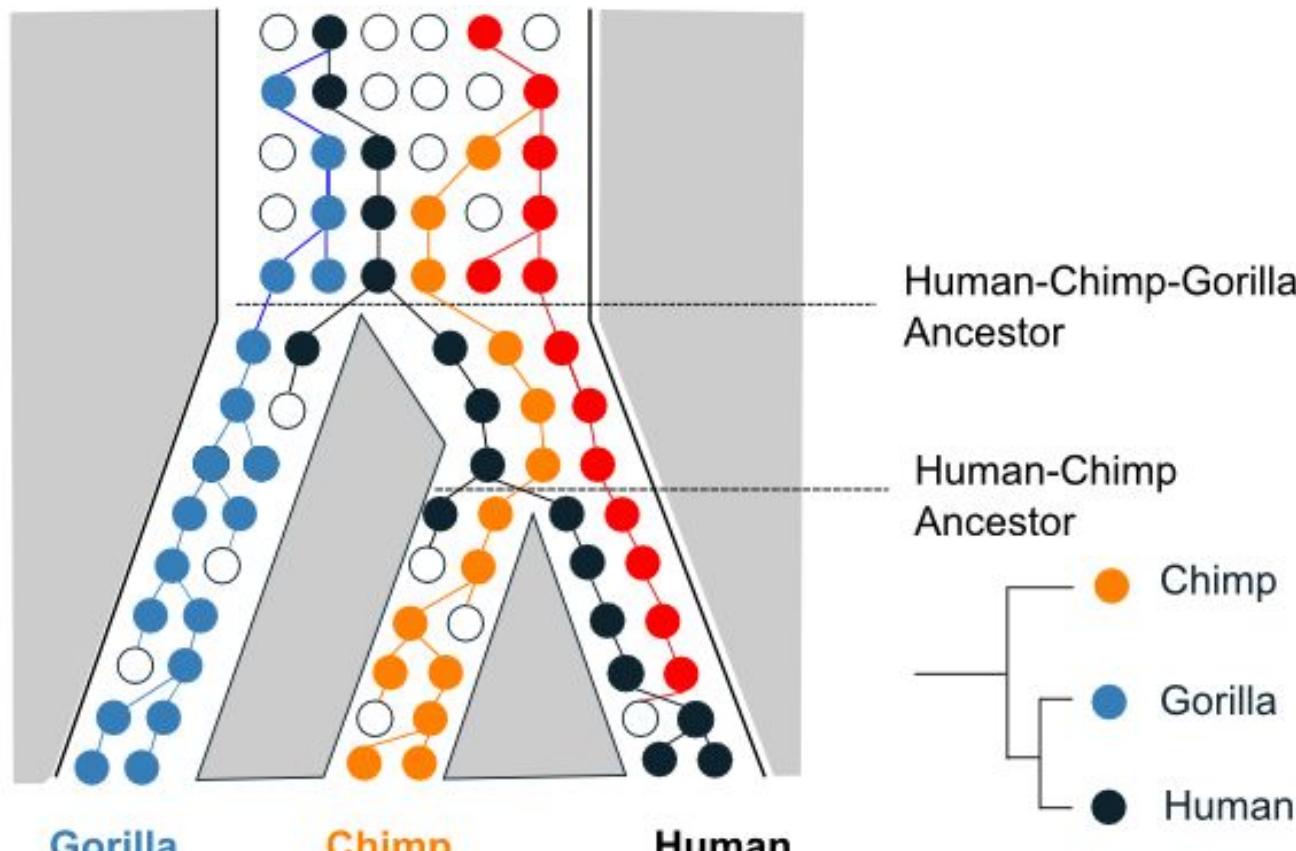


- ~15% Human genome more similar to Gorilla than Chimp
- ~15% Chimp genome more similar to Gorilla than Human

# Causes of disagreement (1): stochastic/sampling error



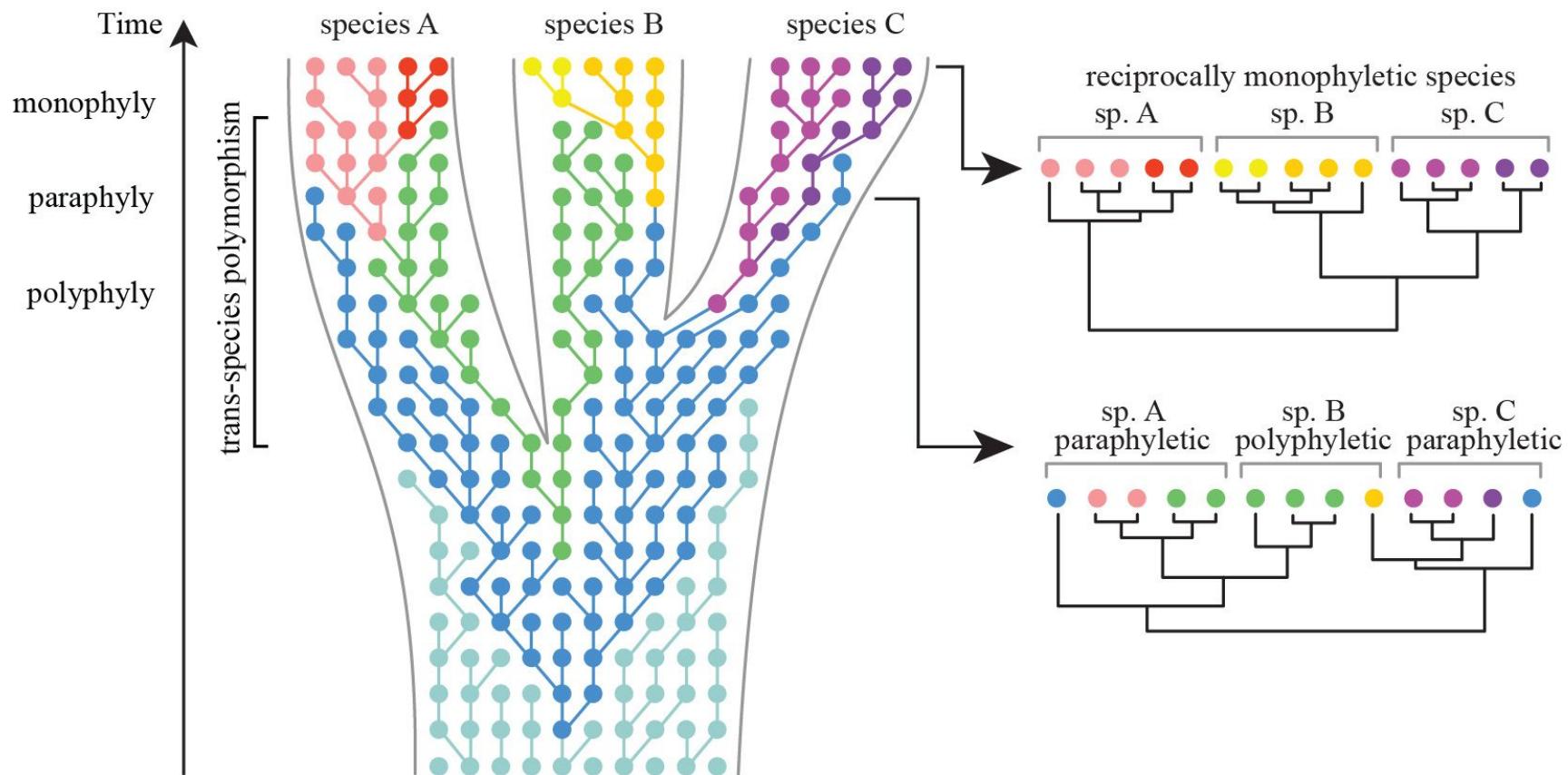
# Causes of disagreement (2): genuinely different gene histories



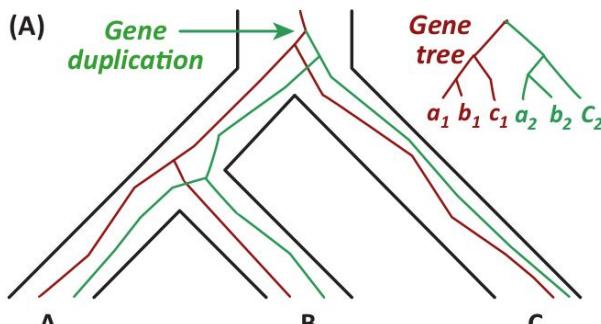
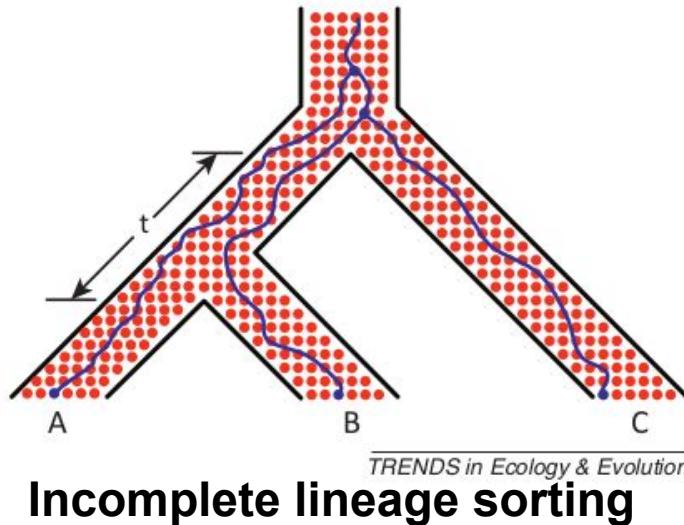
Incomplete lineage sorting

# Genome-level evolution

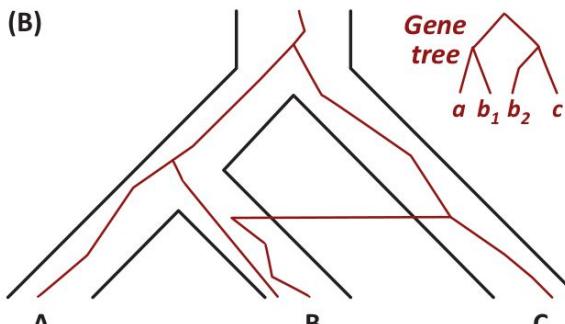
- An extra layer in our models of molecular evolution: genes (and alleles) evolve within the species tree



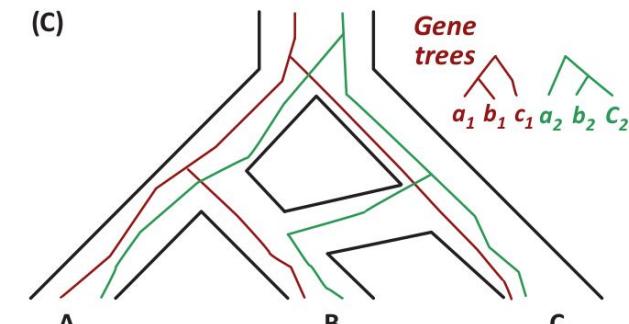
# Genome-level evolution on short and long timescales



Gene duplication



Horizontal transfer



Hybridization

# Genome evolution on short and long timescales

- Current methods for modelling genome evolution lie at the interface between population genetics and phylogenetics:
- Short timescales: incomplete lineage sorting, hybridization
- Long timescales: gene duplications, losses, transfers
- This is due to model limitations, not a disconnect in the underlying biology.
-

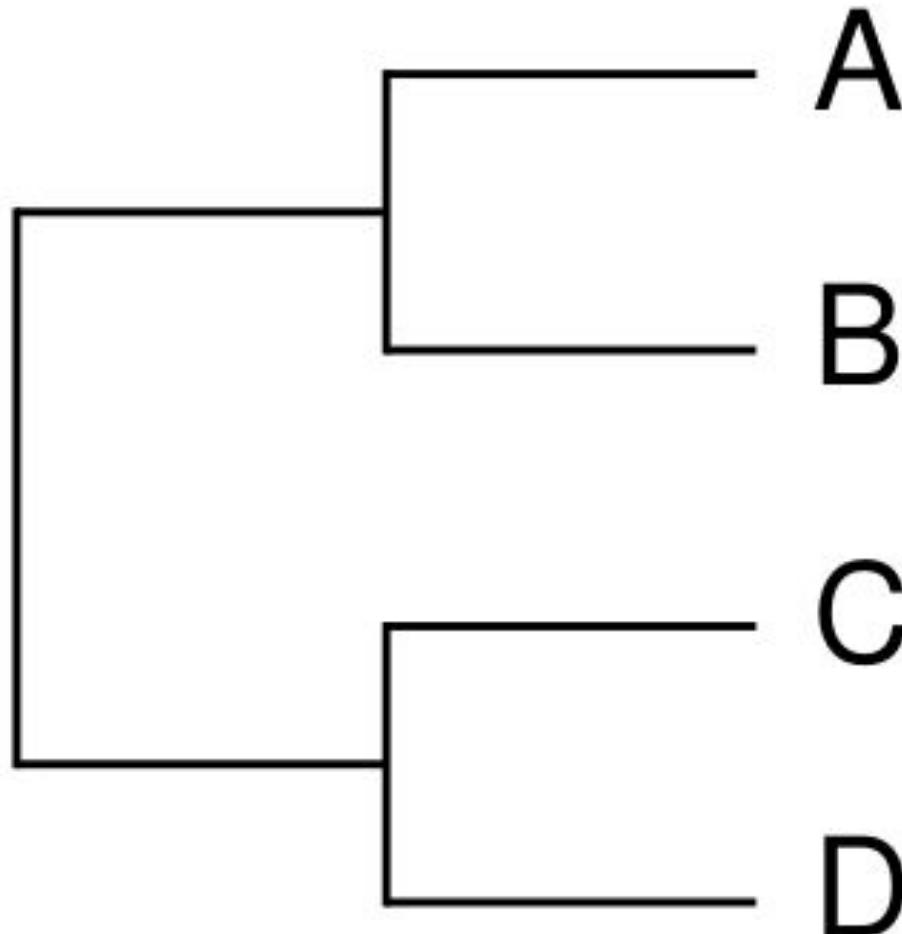
# Disagreements between gene trees and species trees

- **Whole-genome data are great**
  - More information to resolve evolutionary relationships
- **But, trees inferred from different genes don't always agree**
  - Stochastic error
  - Various biological processes

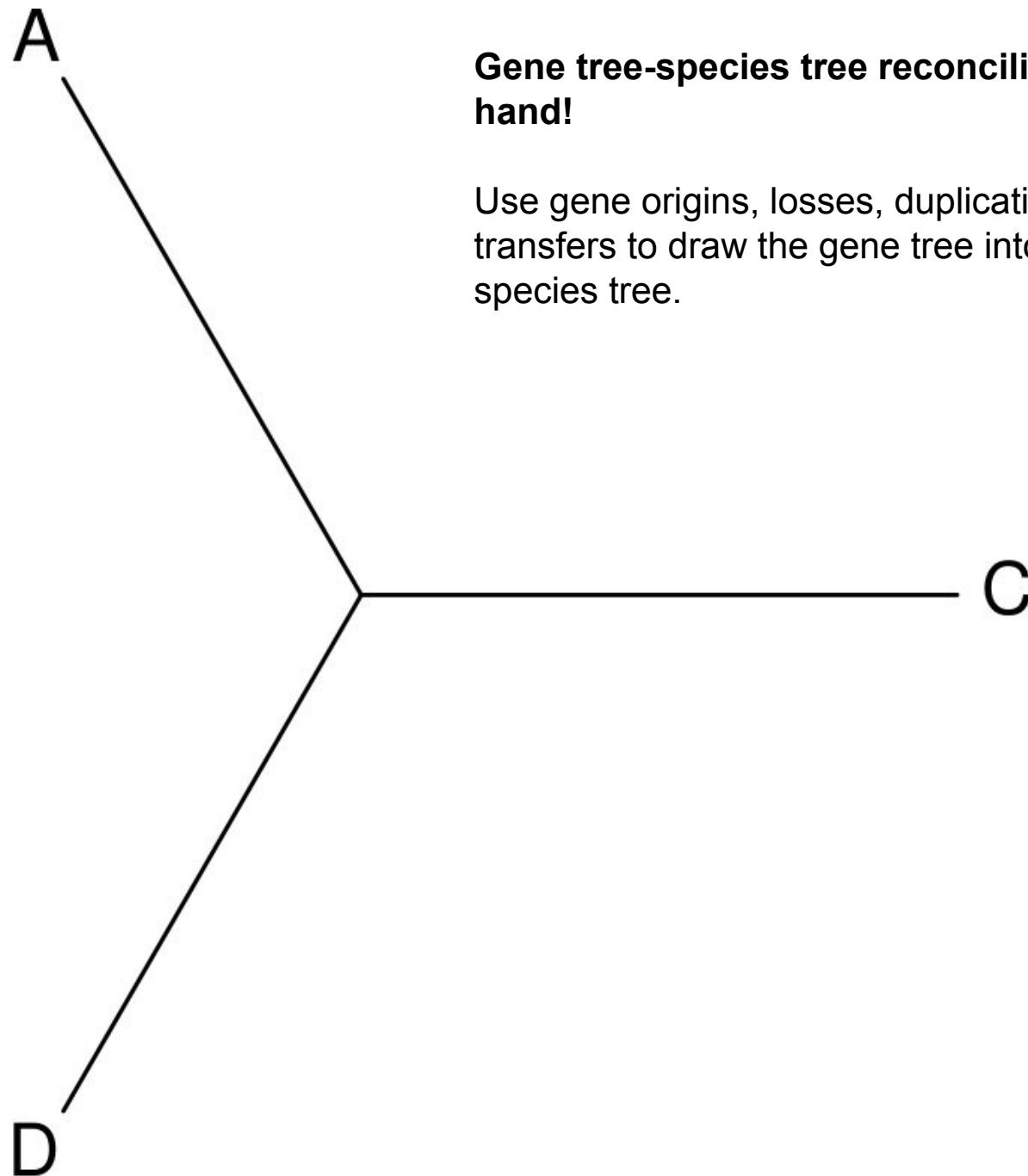
# Is a single tree appropriate?

- **Disagreement between gene trees and species trees**
  - Stochastic error
  - Recombination, sex, incomplete lineage sorting
  - Horizontal gene transfer
  - Mistaken orthology, paralogy
- **Alternative methods in these cases:**
  - Multispecies coalescent (Beast, ASTRAL, RevBayes...)
  - Gene tree-species tree reconciliation (guenomu, ALE, ...)
  - Similarity networks

# Games with species and gene trees

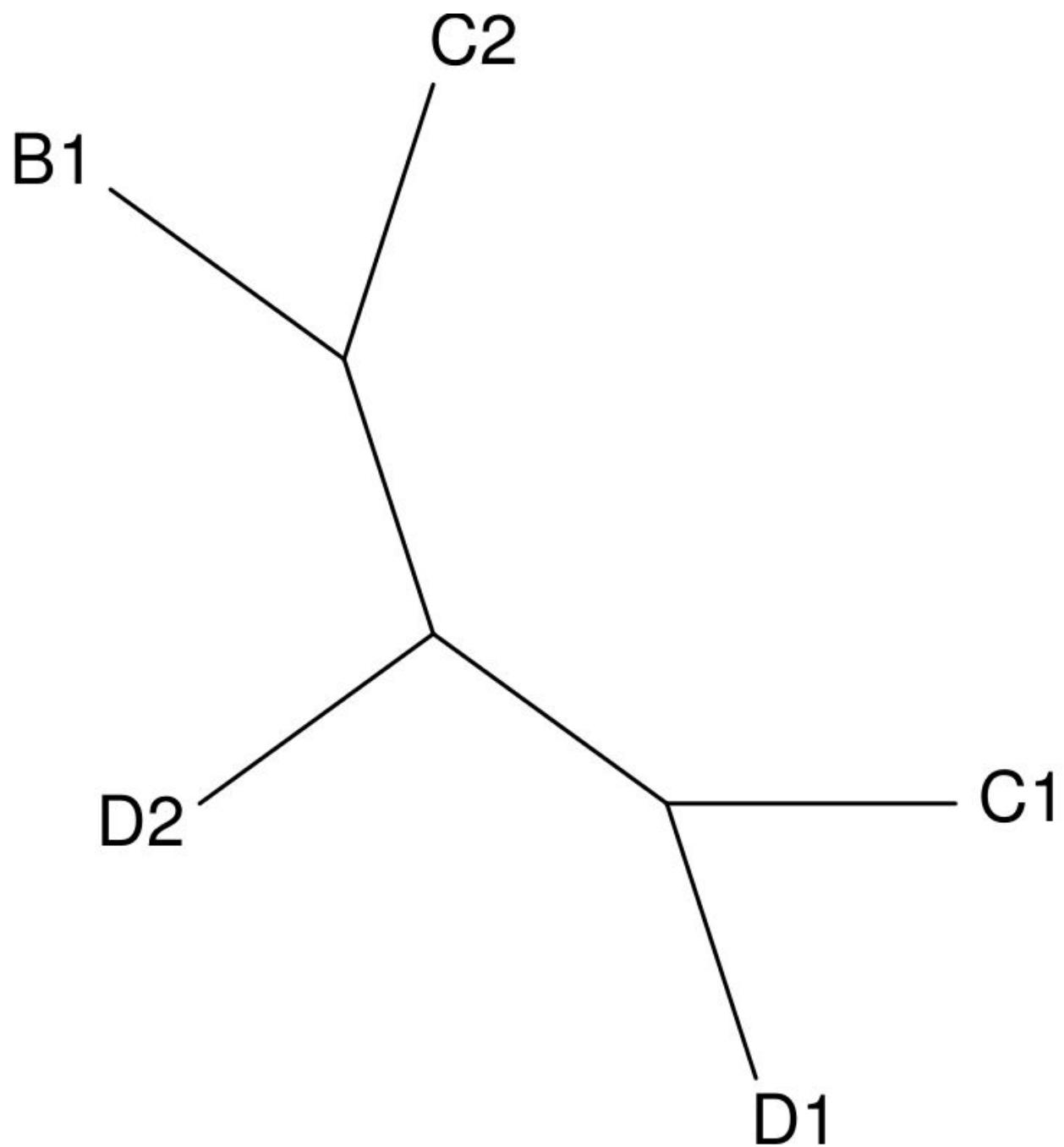


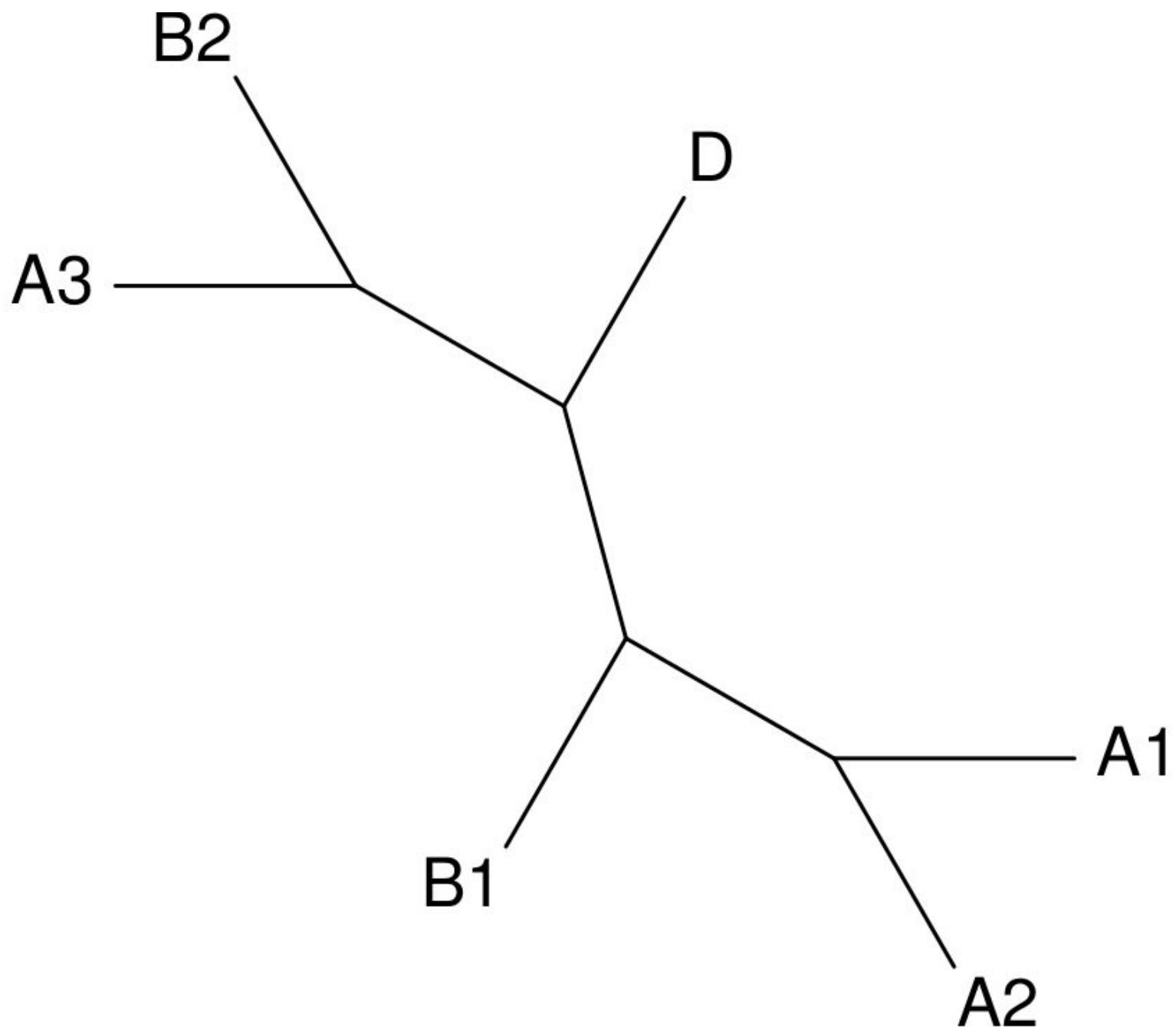
A rooted species tree, assumed known.

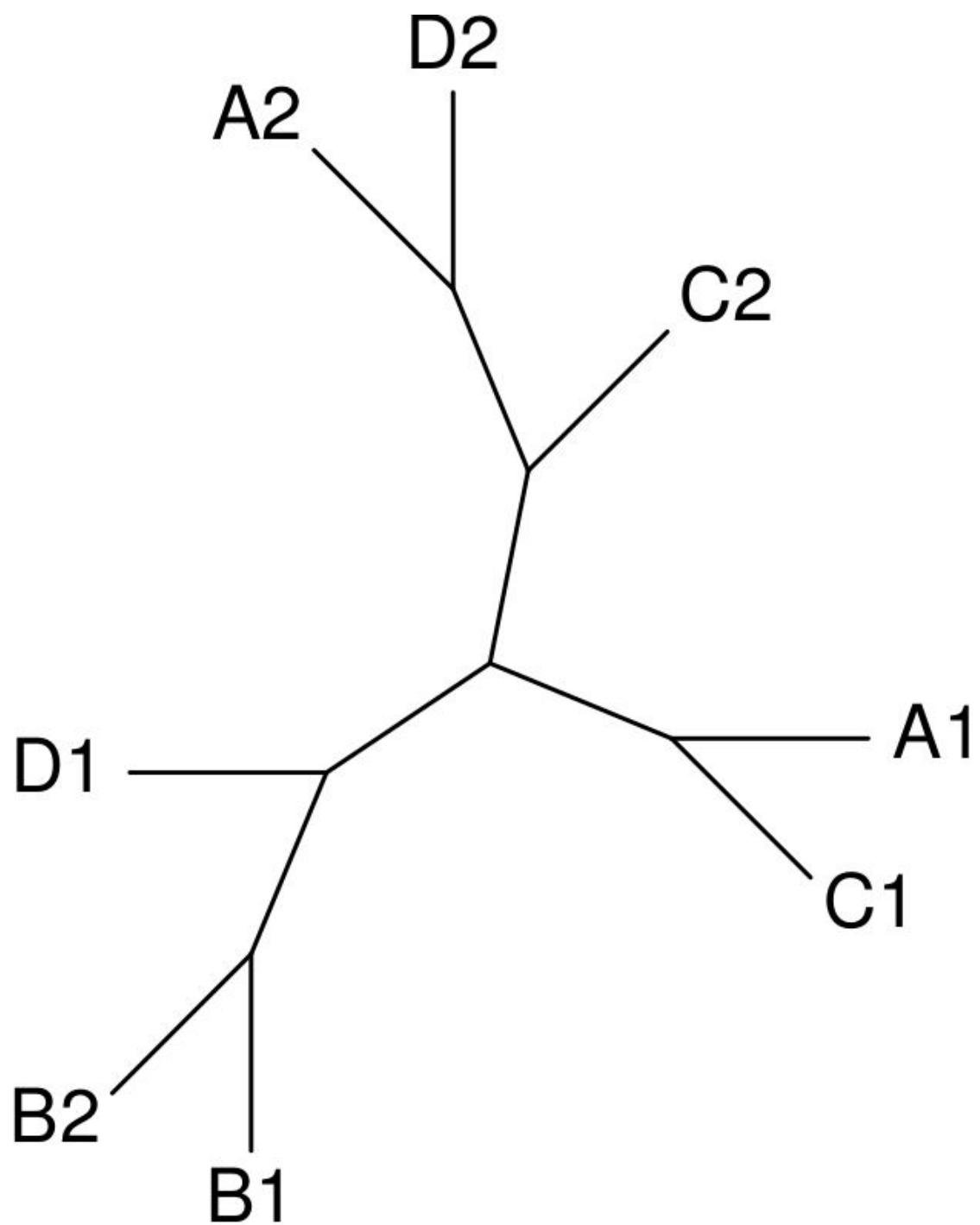


## Gene tree-species tree reconciliation: by hand!

Use gene origins, losses, duplications and transfers to draw the gene tree into the species tree.







# Disagreements among trees express information

- Several reconciliations may seem plausible
- Rates of gene duplications, losses, transfers can be inferred
- Would a different species tree make reconciliations easier?

# Some further reading

## Papers:

- Yang and Rannala (2008) “Molecular phylogenetics: principles and practice”. *Nat Rev Genet*
- Holder and Lewis (2003) “Phylogeny estimation: traditional and Bayesian approaches.” *Nat Rev Genet*

## Books:

- Bromham L. “An introduction to molecular evolution and phylogenetics.”
- Yang Z. “Computational molecular evolution: a statistical approach.”