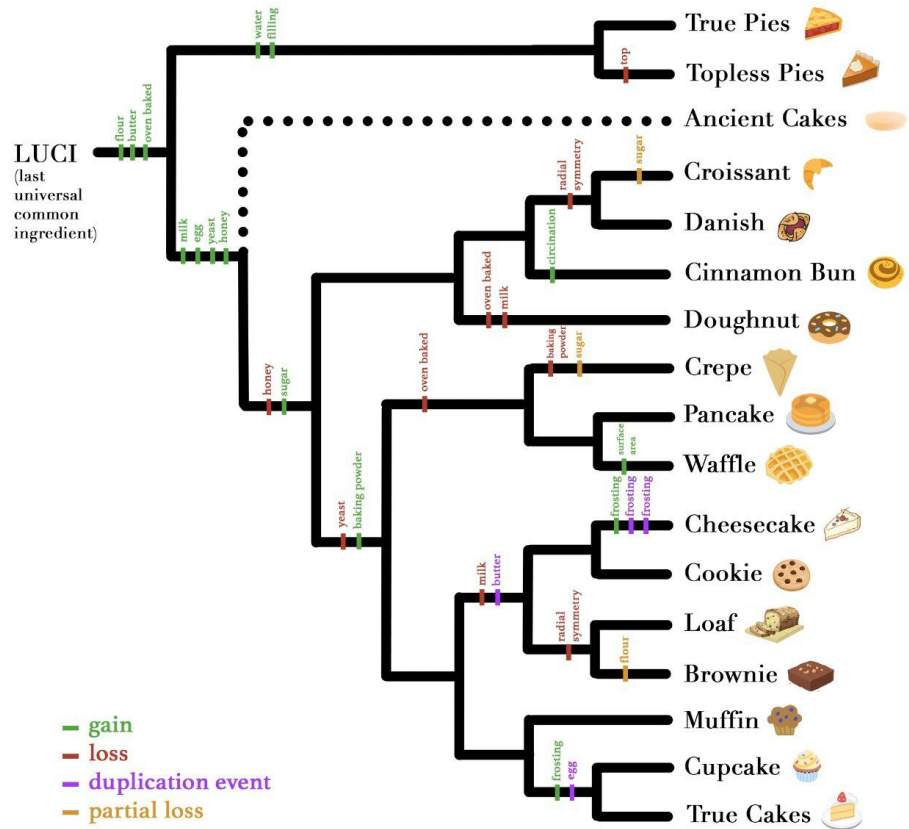


Comparative genomics

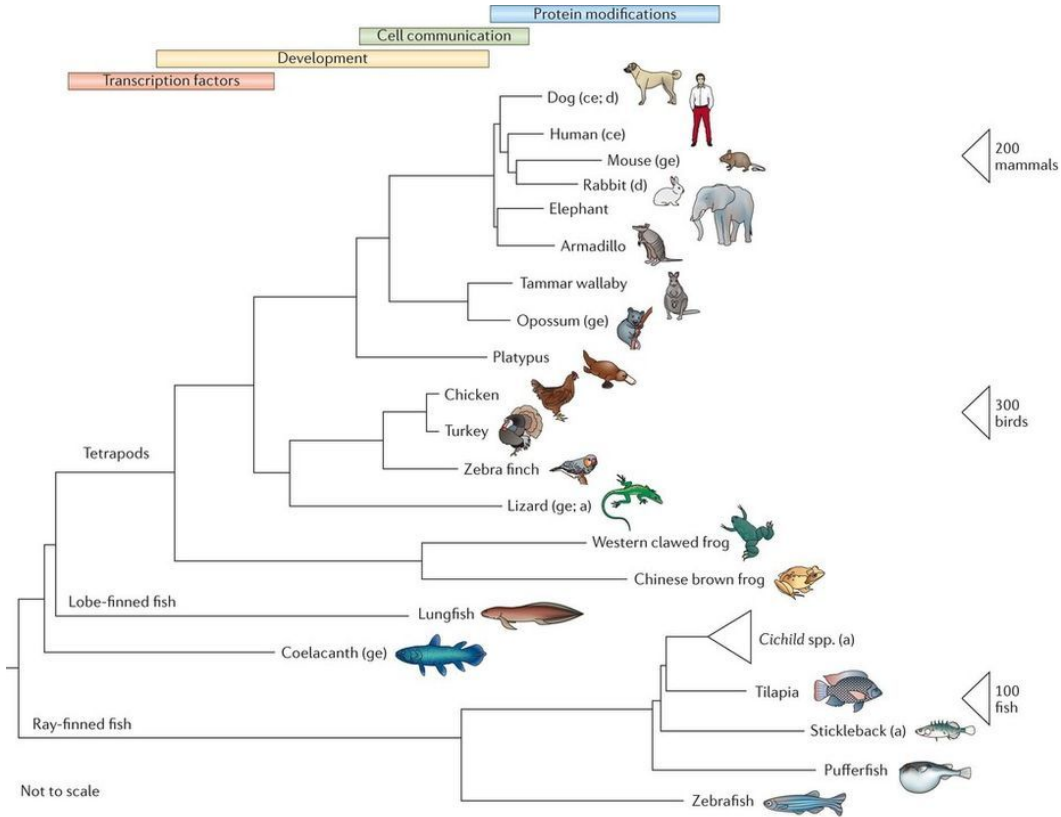
- Gene/species trees
 - Phylogenetic trees
 - Gene orthology & functional analysis

Phylogeny & Phylogenetic tree



**On the Origin of Baked Goods
by Means of Natural Consumption**

Evolutionary relationships between species



a = adaptation

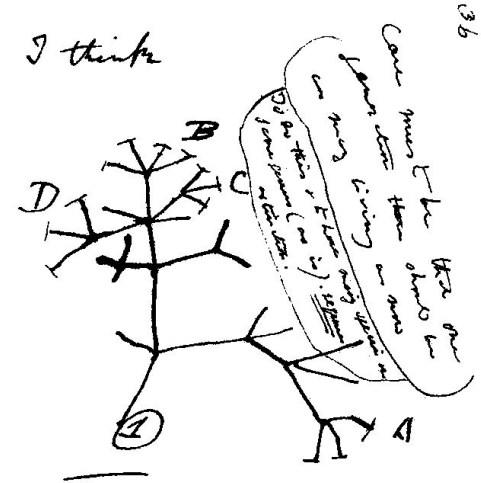
d = domestication

ge = genome evolution

ce = convergent evolution

Phylogeny & Phylogenetic tree

- Useful for:
 - a. organizing knowledge of biological diversity,
 - b. structuring classifications, and
 - c. providing insight into events that occurred during evolution.
- Diagram that depicts the lines of evolutionary descent of different species, organisms, or genes from a common ancestor.
- Trees show descent from a common ancestor.

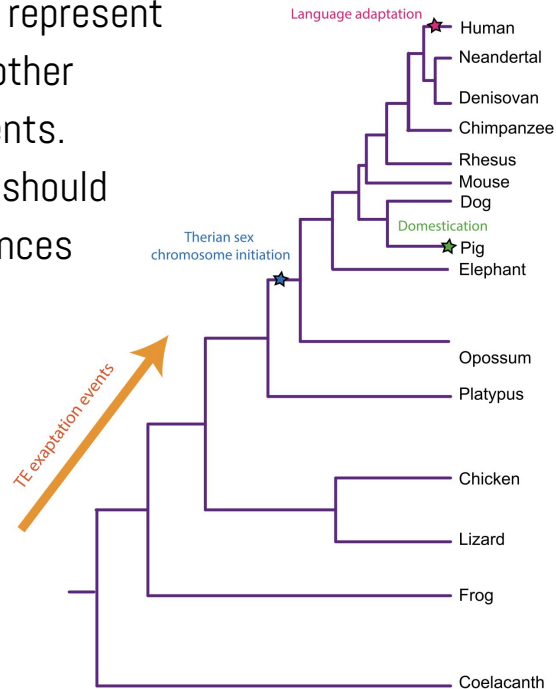


Then between A & B. various
size of relation. C & B. The
first predation, B & D
rather greater distinction
Then genus would be
formed. - binary relation

Species tree vs gene tree

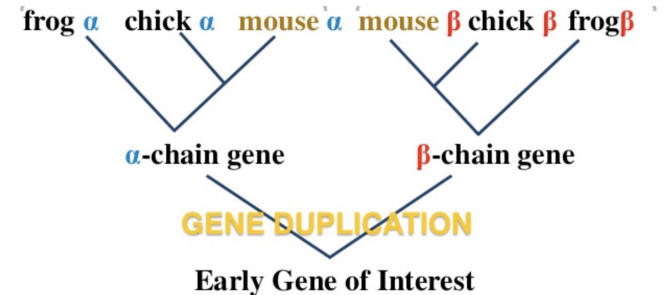
Species trees recover the genealogy of taxa, individuals of a population, etc.

- Internal nodes represent speciation or other taxonomic events.
- Species trees should contain sequences from only orthologous genes.



Gene trees represent the evolutionary history of the genes included in the study.

- Gene trees can provide evidence for gene duplication events, as well as speciation events.
- Sequences from different homologs can be included in a gene tree; the subsequent analyses should cluster orthologs



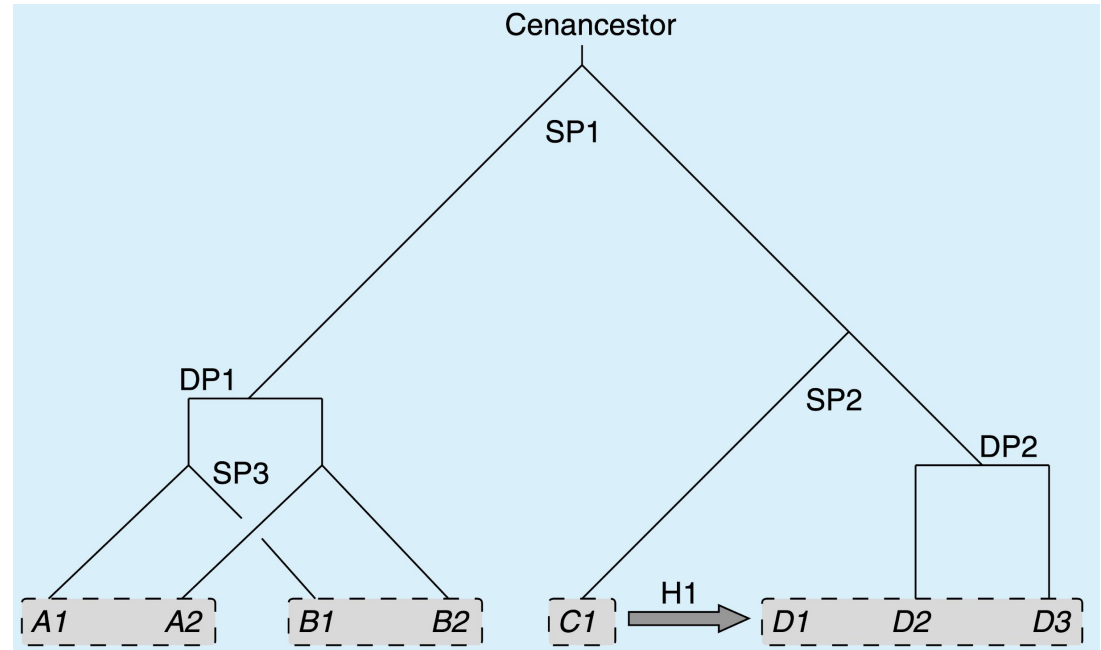
Evolutionary relationships between genes in different species

Evolutionary relationships:

- Orthologs
- Paralog
 - Subfunctionalization
 - Neofunctionalization

Complicated evolutionary processes:

- gene fusion and fission
- horizontal gene transfer
- whole gene deletion



Approaches for constructing phylogenetic trees

Distance-based methods

- UPGMA & Neighbor-Joining
- Calculate pairwise distances & then build tree

Character-based methods

- Maximum parsimony & Maximum likelihood
- Directly build tree by coupling tree proposal & scoring

Distance-based methods for constructing phylogenetic trees

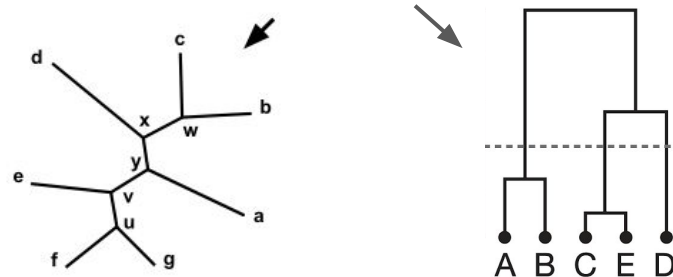
Multiple sequence alignment

```
B9SI54|B9SI54_RICCO_263_570      RILTNVYMGDGIRTIIISGSKHTM-DGLPAYRTATVAVLGDGFVCKSMTIQNSATSD-K
Q01I60|Q01I60_ORYSA_160_476      YEKTNILLVGDGIGATVITASRSVGIDGIGTYETATVAVIGDGFRAKDITFENGAGAGAH
C5Y8S2|C5Y8S2_SORBI_153_466      YEKTNILLMGEGMGATVITASRSVGIDGLGTHETATVAVIGDGFRAKDITFENSAGARAH
B4FRR6|B4FRR6_MAIZE_154_469      YEKANILLMGEGMGATVITASRSVGIDGLGTYETATVDVIGDGFRAKDITFENSAGAGAH
D7U4G4|D7U4G4_VITVI_82_394       LEKKNVVFLGDGMGKTVITGSLNVGQPGISTYNSATVGAGDGFMASGLTMENTAGPDEH
D7M270|D7M270_ARALY_263_574      FEKKNVVFLGDGMGKTVITGSLNAGMPGITTNTATVGVGVDGFMADLTFQNTAGPDAH
Q8L7Q7|PME64_ARATH_283_601       FEKKNVVFLGDGMGKTVITGSLNVGQPGMTTFESATVGVLGDGFMARDLTIENTAGADAH
D8QSM2|D8QSM2_SELML_242_541      DSKSMIMLVGAGARKTIIISGNVYVR-EGVTTMDTATVLVAGDGFVARDLTIRNTAGPELH
A9TZ89|A9TZ89_PHYPA_262_575      KQKTNLMFLGDGTDKTIITGSLSDSQPGMITWATATVAVSGSGFIARGITFQNTAGPAGR
D8SH72|D8SH72_SELML_209_529      LQKSNLMFVGDGMDKTIIRGSMVSKGGTTTFASATLAVNGKGFRLARDLTVENTAGPEGH
                                     1 1 * * * 2 * . .      * 1 1 2 2 * * . . 1 * . * *
```

Distance-based methods

- UPGMA & Neighbor-Joining
- Calculate pairwise distances & then build tree

	A	B	C	D	E
A	0				
B	5	0			
C	10	3	0		
D	15	6	7	0	
E	20	8	2	11	0



Distance-based methods for constructing phylogenetic trees

UPGMA (Unweighted Pair Group Method with Arithmetic Mean)

- Rooted tree
- Assumes constant-rate

Distance b/w any two clusters
A and **B**, each of size = the
farthest distance between
elements of each cluster

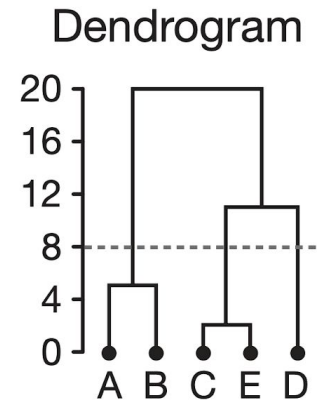
A	A	T	C	G	T	G	G	T	A	C	T	G
B	C	C	G	G	A	G	A	A	C	T	A	G
C	A	A	C	G	T	G	C	T	A	C	T	G
D	A	T	G	G	T	G	A	A	G	T	G	
E	C	C	G	G	A	A	A	A	C	T	T	G

			1		
	A	B	C	D	E
A	0				
B	5	0			
C	10	3	0		
D	15	6	7	0	
E	20	8	2	11	0

			2		
	A	B	CE	D	
A	0				
B	5	0			
CE	20	8	0		
D	15	6	11	0	

			3		
	AB	CE	D		
AB	0				
CE	20	0			
D	15	11	0		

			4		
	AB	CED			
AB	0				
CED	20	0			



Distance-based methods for constructing phylogenetic trees

Neighbor-Joining

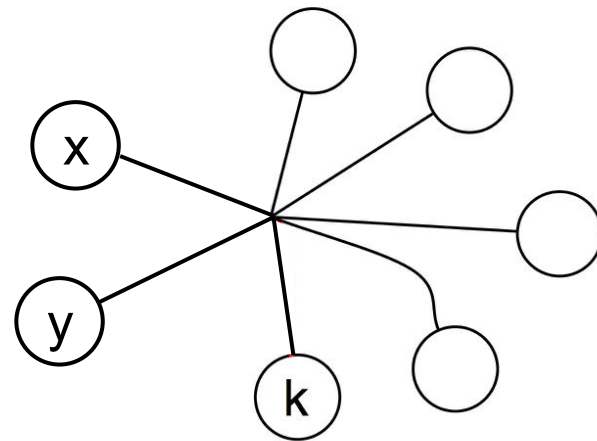
- Unrooted tree
- Does not assume constant-rate

Choose x, y to merge
that minimize:

$$Q(x, y) := (n - 2)D_{xy} - \left(\sum_{k=1}^n D_{xk} + \sum_{k=1}^n D_{yk} \right)$$

Update lengths:

	A	B	C	D	E
A	0				
B	5	0			
C	10	3	0		
D	15	6	7	0	
E	20	8	2	11	0



Distance-based methods for constructing phylogenetic trees

Neighbor-Joining

- Unrooted tree
- Does not assume constant-rate

Choose x, y to merge
that minimize:

$$Q(x, y) := (n - 2)D_{xy} - \left(\sum_{k=1}^n D_{xk} + \sum_{k=1}^n D_{yk} \right)$$

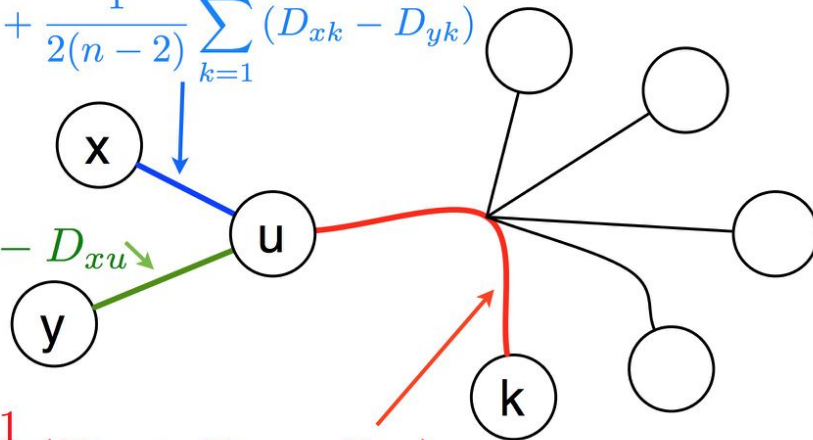
Update lengths:

$$D_{xu} := \frac{1}{2}D_{xy} + \frac{1}{2(n-2)} \sum_{k=1}^n (D_{xk} - D_{yk})$$

$$D_{yu} := D_{xy} - D_{xu}$$

$$D_{uk} := \frac{1}{2} (D_{xk} + D_{yk} - D_{xy})$$

	A	B	C	D	E
A	0				
B	5	0			
C	10	3	0		
D	15	6	7	0	
E	20	8	2	11	0



Distance-based methods for constructing phylogenetic trees

Neighbor-Joining

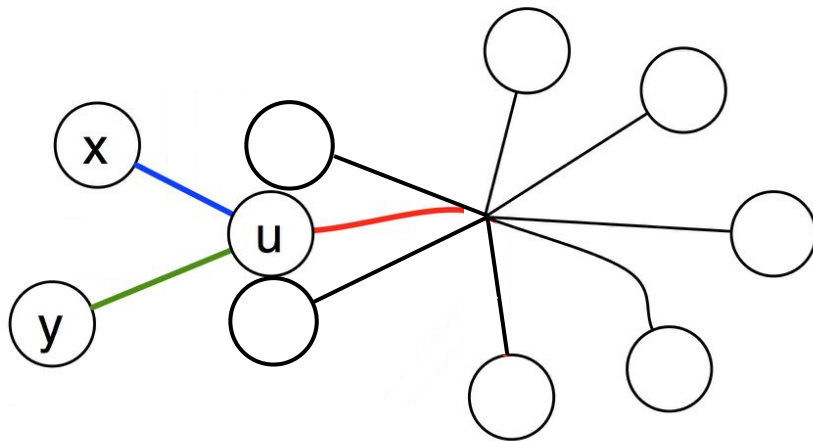
- Unrooted tree
- Does not assume constant-rate

Choose x, y to merge
that minimize:

$$Q(x, y) := (n - 2)D_{xy} - \left(\sum_{k=1}^n D_{xk} + \sum_{k=1}^n D_{yk} \right)$$

Update lengths:

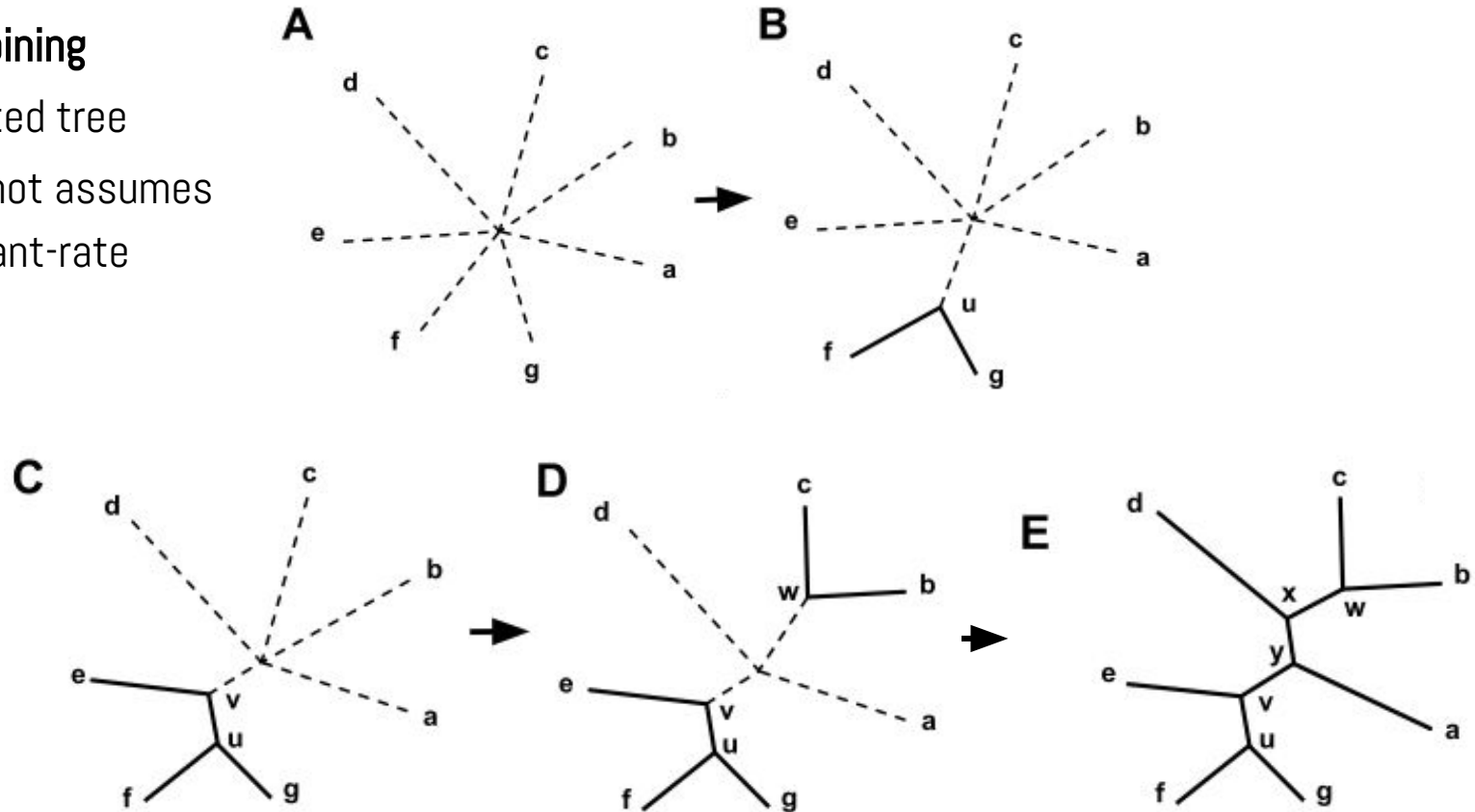
	A	B	C	D	E
A	0				
B	5	0			
C	10	3	0		
D	15	6	7	0	
E	20	8	2	11	0



Distance-based methods for constructing phylogenetic trees

Neighbor-Joining

- Unrooted tree
- Does not assume constant-rate



Distance-based methods for constructing phylogenetic trees

High computational efficiency (esp. NJ).

- Useful for analysing large data sets with low levels of sequence divergence.

Can perform poorly for very divergent sequences.

- Large distances involve large sampling errors, and most distance methods (such as NJ) do not account for the high variances of large distance estimates.

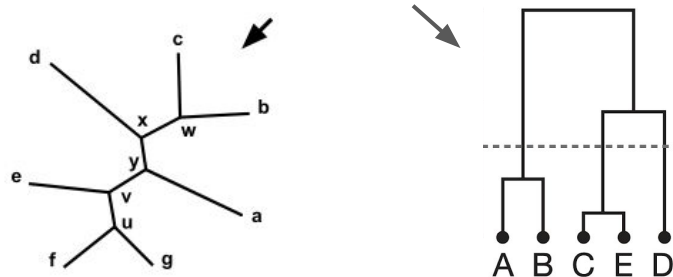
Need a realistic substitution model to calculate the pairwise distances. Also sensitive to gaps in the sequence alignment.

Multiple sequence alignment

B9SI54	B9SI54_RICCO	263	570
Q01160	Q01160_ORYSA	160	476
S5Y852	S5Y852_SORBI	153	466
B4FRR6	B4FRR6_MAIZE	154	469
D7U4G4	D7U4G4_VITVI	82	394
D7M270	D7M270_ARALY	263	574
Q8L7QE	Q8L7QE_ARATH	283	601
D8QSM2	D8QSM2_SELML	242	541
A9TZ89	A9TZ89_PHYPA	262	525
D8SH72	D8SH72_SELML	209	529

RILTVNLLVYMGDGDIIISGSKHMT-DGLPAYRTATVAVLGDGFVCKSMTIQNSATSD-K
YEKTNILLVGGDIGATVITASRSGVIGDIGYETAATVAVIGDGFRAKDIITFENGAGAGA
YEKTNILLMGVGDGATVITASRSGVIGDLGLETATVAVIGDGFRAKDIITFENGAGARH
YEKANILLMGGMGATVITASRSGVIGDLGLEYETAATVVDIGDGFRAKDIITFENGAGARH
LEKNNVFLGDGKGVITIGSLNVSQPGITSTYNSATVGVAGDGFMSGLTMENTAGPDEH
FEKNNVVFIDGMGKVTITGSLNMGPGITTYNTATVGVVDGDFMAHDLTQNTAGPDH
FEKNNVVFIDGMGKVTITGSLNVSQPGMTTFESATVGVGDGFMAHDLTIENTAGADAH
DSKSMIMLVGGAGARKTIISGNVYR-EGVTTMDATVLVAGDGFVARDLIRNTAGPELH
LQKTNLMFLGDGDKTITIGSLSDSQPGMTTWTATVAVVSGSGFLARGITFQNTAGPAGK
LQKSMLMVFGDGDGDKTIRGSMVSQGGITTFASATLVANGFLARDLIENTAGPELH

	A	B	C	D	E
A	0				
B	5	0			
C	10	3	0		
D	15	6	7	0	
E	20	8	2	11	0



Approaches for constructing phylogenetic trees

Distance-based methods

- UPGMA & Neighbor-Joining
- Calculate pairwise distances & then build tree

Character-based methods

- Maximum parsimony & Maximum likelihood
- Directly build tree by coupling tree proposal & scoring

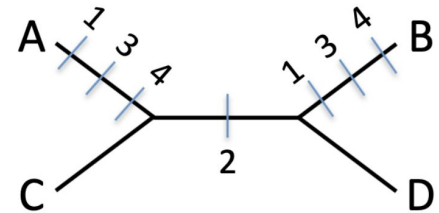
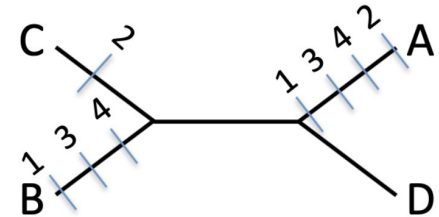
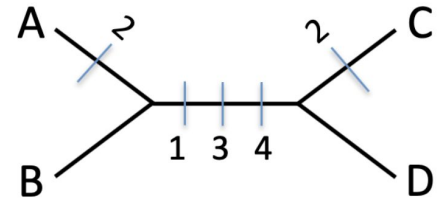
Maximum parsimony

MP **minimizes the number of changes on a phylogenetic tree** by assigning character states to interior nodes on the tree.

The **character (or site) length** is the minimum number of changes required for that site, whereas the **tree score** is the sum of character lengths over all sites.

The **maximum parsimony tree** is the tree that minimizes the tree score.

	1	2	3	4
A:	A	C	G	T
B:	C	C	G	A
C:	G	C	G	C
D:	T	C	C	C



Maximum likelihood

Maximum Likelihood is a:

- general statistical method
- for **estimating unknown parameters** of a probabilistic model
- by maximizing a function, so that
- **under the assumed model,**
- **the observed data is most probable.**

█ Likelihood of hypothesis =
Probability of data given hypothesis

- Fair or unfair coin?

$$P_{\text{head}} = 0.5$$

Fair

$$P_{\text{head}} = 0.67$$

Unfair



- Flip coin 4 times, get:

3 heads, 1 tail

	Fair	Unfair
H x H x H x T	$1/2 \times 1/2 \times 1/2 \times 1/2 = 1/16$	$2/3 \times 2/3 \times 2/3 \times 1/3 = 8/81$
H x H x T x H	$1/2 \times 1/2 \times 1/2 \times 1/2 = 1/16$	$2/3 \times 2/3 \times 1/3 \times 2/3 = 8/81$
H x T x H x H	$1/2 \times 1/2 \times 1/2 \times 1/2 = 1/16$	$2/3 \times 1/3 \times 2/3 \times 2/3 = 8/81$
T x H x H x H	$1/2 \times 1/2 \times 1/2 \times 1/2 = 1/16$	$1/3 \times 2/3 \times 2/3 \times 2/3 = 8/81$
Total	$1/4$ (0.25)	$32/81$ (0.40)

Maximum likelihood

1. Given **data**, assume it **comes from a model** (e.g., normal/binomial distribution).
2. **Likelihood** ~ the probability of observing the data given the model: **$P(\text{Data} \mid \text{Model})$** .
3. Examine this likelihood function to see **where it is greatest** (meaning, different values of the parameters of the model: e.g. μ & σ).
4. The values of the parameters at that point is the **maximum likelihood estimate** of the parameters (found numerically by some iterative optimization procedure).

MLEs have desirable asymptotic properties:

1. Unbiased (expected value = true value of the parameter),
2. Consistent (approach true values), &
3. Efficient (have the smallest variance among unbiased estimates).

- Flip coin 4 times, get:

3 heads, 1 tail				
	Fair		Unfair	
H x H x H x T	$1/2 \times 1/2 \times 1/2 \times 1/2 = 1/16$		$2/3 \times 2/3 \times 2/3 \times 1/3 = 8/81$	
H x H x T x H	$1/2 \times 1/2 \times 1/2 \times 1/2 = 1/16$		$2/3 \times 2/3 \times 1/3 \times 2/3 = 8/81$	
H x T x H x H	$1/2 \times 1/2 \times 1/2 \times 1/2 = 1/16$		$2/3 \times 1/3 \times 2/3 \times 2/3 = 8/81$	
T x H x H x H	$1/2 \times 1/2 \times 1/2 \times 1/2 = 1/16$		$1/3 \times 2/3 \times 2/3 \times 2/3 = 8/81$	
Total	1/4 (0.25)		32/81 (0.40)	

Maximum likelihood for tree estimation

Model: The tree; **Parameters:** The tree's branch lengths.

ML for tree inference is equivalent to comparing many statistical models, each with the same number of parameters.

Use a **specific substitution model**:

- Assume independent evolution of sites in the sequence → likelihood = product of the probabilities for different sites.
- Probability at any particular site = average over the unobserved character states at the ancestral nodes.

Two optimization steps:

1. Optimization of branch lengths to calculate the tree score for each candidate tree.
2. A search in the tree space for the maximum likelihood tree.

Maximum likelihood for tree estimation

Maximum likelihood is **used exclusively these days** for inferring deep phylogenies using conserved proteins.

- All model **assumptions** are explicit, so that they can be evaluated and improved.
- Availability of a rich repertoire of sophisticated evolutionary **models**.
 - Including models that accommodate variable amino acid substitution rates among sites or different amino acid frequencies among sites.
- Great for **understanding the process** of sequence evolution.
 - The likelihood ratio test can be used to:
 - Examine the fit of evolutionary models
 - Test interesting biological hypotheses (e.g. molecular clock) and selection affecting protein evolution.

Maximum likelihood for tree estimation

There are some drawbacks!

- The attractive asymptotic properties of MLEs apply to parameter estimation when the true tree is given but not to the maximum likelihood tree.
- The likelihood calculation, particularly tree search under the likelihood criterion, is **computationally demanding**.
- The method has potentially **poor statistical properties if the model is misspecified**.

Approaches for constructing phylogenetic trees

Distance-based methods

- UPGMA & Neighbor-Joining
- Calculate pairwise distances & then build tree

Character-based methods

- Maximum parsimony & Maximum likelihood
- Directly build tree by coupling tree proposal & scoring

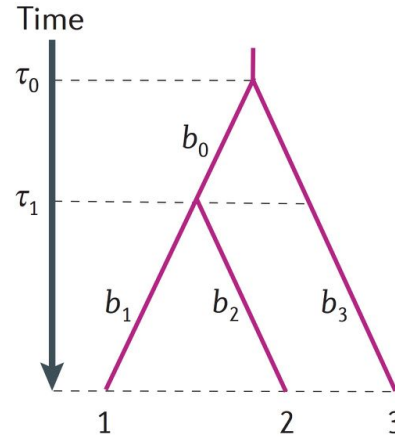
Rooted vs. Unrooted trees

Substitution rate is constant over time or among lineages → the molecular clock holds.

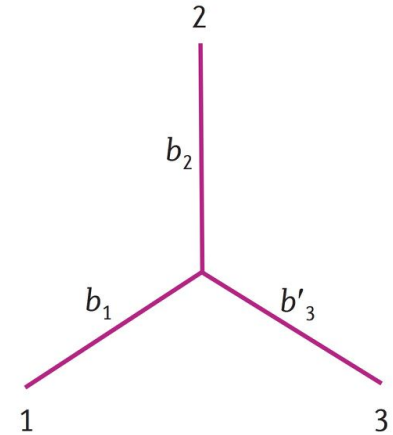
For distantly related species, the clock hypothesis should not be assumed.

If every branch on the tree is allowed to have an independent evolutionary rate → unrooted trees.

a Rooted tree



b Unrooted tree



Rooted vs. Unrooted trees

Substitution rate is constant over time or among lineages \rightarrow the molecular clock holds.

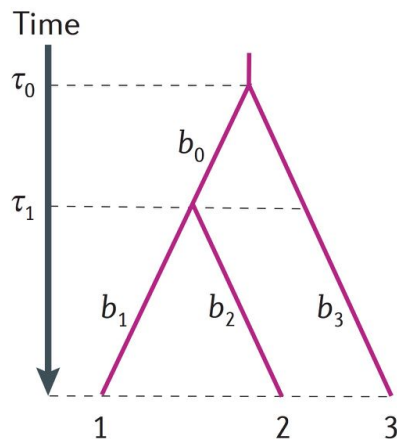
The tree will then have a root (inferring rooted tree is called molecular clock rooting).

- The tree will be ultrametric: distances from the tips of the tree to the root are all equal ($b_0 + b_1 = b_0 + b_2 = b_3$).

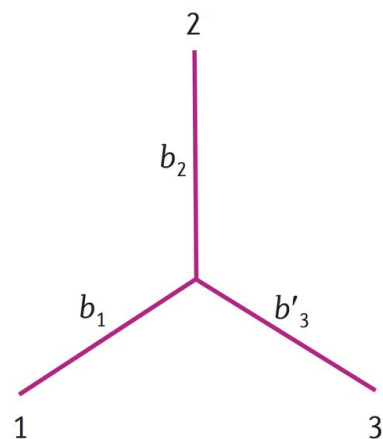
A rooted tree for s species:

- Can then be represented by the ages of the $s - 1$ ancestral nodes.
- Involves $s - 1$ branch-length parameters.

a Rooted tree



b Unrooted tree

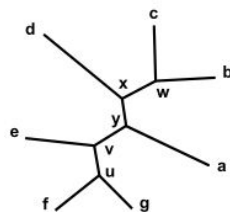


For distantly related species, the clock hypothesis should not be assumed.

Rooted vs. Unrooted trees

If every branch on the tree is allowed to have an independent evolutionary rate \rightarrow unrooted trees.

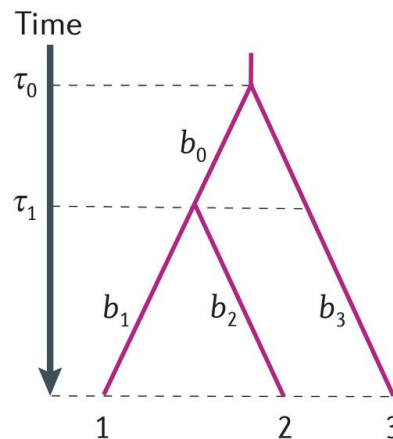
An unrooted tree for s species has $2s - 3$ branch length parameters.



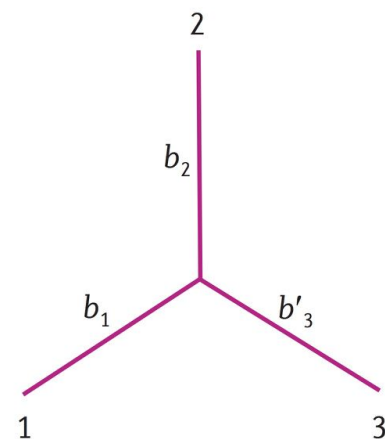
Rooting a tree using outgroup rooting:

- Include outgroup species (a species/genes known to be more distantly related than the species/genes of interest).
- Root is located along the branch that leads to the outgroup so that the tree for the ingroup species is rooted.

a Rooted tree



b Unrooted tree



Interpreting a tree

