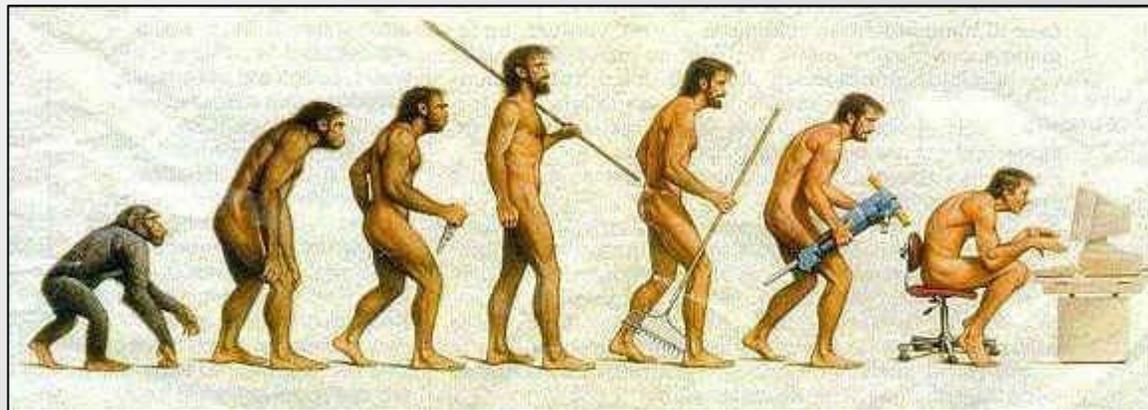


MIT OpenCourseWare
<http://ocw.mit.edu>

6.047 / 6.878 Computational Biology: Genomes, Networks, Evolution
Fall 2008

For information about citing these materials or our Terms of Use, visit: <http://ocw.mit.edu/terms>.

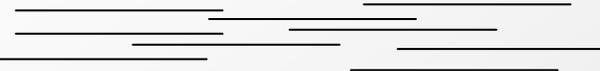
Molecular Evolution and Phylogenetics



Patrick Winston's 6.034

Somewhere, something went wrong...

Challenges in Computational Biology

④ Genome Assembly 

Regulatory motif discovery 

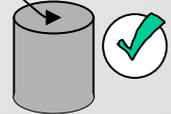
Gene Finding 

Sequence alignment 

Comparative Genomics 

⑦ Evolutionary Theory

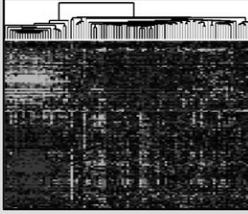
```
TCATGCTAT
TCGTGATAA
TGAGGATAT
TTATCATAT
TTATGATTT
```

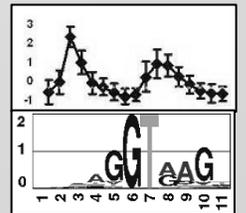
Database lookup 

RNA folding 

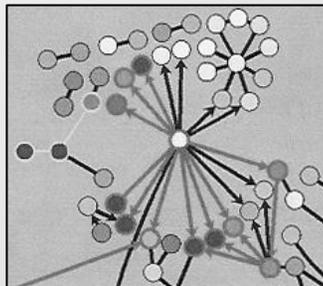
⑨ Gene expression analysis

RNA transcript 

⑩ Cluster discovery 



Gibbs sampling 

⑫ Protein network analysis 

⑬ Regulatory network inference

⑭ Emerging network properties

Goals for today

- Basics of phylogeny
 - Characters, traits, nodes, branches, lineages
 - Gene trees, species trees
- Modeling sequence evolution
 - Turning sequence data into distances
 - Probabilistic models of nucleotide divergence
 - Jukes-Cantor 1-parameter model, Kimura 2-parameter model
- From distances to trees
 - Ultrametric, Additive, General Distances
 - UPGMA, Neighbor Joining, guarantees and limitations
 - Least-squared error, minimum evolution
- From alignments to trees
 - Parsimony methods: set-based vs. dynamic programming
 - Maximum likelihood methods
 - MCMC and heuristic search

Open questions (?)

- Panda
 - Bear or raccoon?
- Out of Africa
 - mitochondrial evolution story?
- Human evolution
 - Did we ever meet Neanderthal?
- Primate evolution
 - Are we chimp-like or gorilla-like?
- Vertebrate evolution
 - How did complex body plans arise?
- Recent evolution
 - What genes are under selection?

Inferring Phylogenies: Traits and Characters

Trees can be inferred by several criteria:

- Morphology data

Image removed due to copyright restrictions.

- Molecular data



Traits – as many as we have letters in DNA

YAL042W
candida586
cdub17784
cgla72177
cgui48535
clus15345
ctro67868
klac20931

-MKRSTLLSLDAFAKTEEDVRVRTRAGGLITLSCILTTLFLLVNEWGQFNSVVTRPQLVV
MSSRPKLLSFDAFAKTVEDARIKTTSGGIITLICILITLVLIRNEYVDYTTIITRPELVV
MSSRPKLLSFDAFAKTVEDARIKTTSGGIITLICILITLVLIRNEYVDYTTIITRPELVV
-MKKSTLLSFDAFAKTEEDVRIRTRSGGFITLGLVVTLMMLLSEWRDFNSVVTRPELVI
-MPQPKLLSFDAFAKTVEDARVTRPAGGIITLICVIVVLYLIRNEYLEYTSIINRPELVV
MSSRPRLSLDAFAKTVEDARVKTASGGVITLVCVLIVLFLIRNEYSDYMLVVVRPELVV
MSSRPKLLSFDAFAKTVEDARIKTASGGIITLICVLITLILIRNEYIDYTTIITRPELVV
-MKKSPLLSIDAFGKTEEDVRVRTRTGGLITVSCIITMMLLVSEWKQFSTIVTRPDLVV
: . ***:***.** **.*:* :**.**: *:: .: *: .*: :: :: **:**:

YAL042W
candida586
cdub17784
cgla72177
cgui48535
clus15345
ctro67868
klac20931

DRDRHAKLELNMDVTFPSPMPCDLVNLDIMDDSGEMQLDILDAGFTMSRLNSEG-----R
DRDINKQLDINLDISFINLPCDLISIDLLDVTGDLNLNIDSGLKKIRLLKNKQGDVIVN
DRDINKQLDINLDISFINLPCDLISIDLLDVTGDLNLNIDSGLKKIRLLKNKQGDVIVN
DRDRSLRLDLNLDITFPSPMPCCELLTLDIMDDSGEVQLDIMNAGFEKTRLSKEG-----K
DRDINKKLEINLDISFPDIPCDVLTMDILDVSGDLQVDLLLSGFEEKFRLLKDG-----L
NRDVNRQLDINLDITFPDPVPCGVMSLDILDMTGDLHLDIVESGFEMFRVPLG-----E
DRDINKQLDINLDISFINLPCDLISVDLLDVTGDDQLDIIDSGLKKVRLLNKQGDVIIN
DRDRHLKLDLNLDTFPSPMPCNVNLNLDILDDSGEFQINLLDSGFTKIRISPEG-----K
:** :*:*:*:*:* :** :*:*:*:* :*: :*: :*: *:

YAL042W
candida586
cdub17784
cgla72177
cgui48535
clus15345
ctro67868
klac20931

PVGDATELHVGGNGDGTAPV--NND---PNY-CGPCYGAQDQSQN-ENLAQEEKVCCQDC
EIEDDEPAFNNDIELSDLAKGLPEGSDENAY-CGSCYALPQDK-----KQFCCNDC
EIEDDEPAFNNDIELTDLAKGLPEGSDENAY-CGSCYALPQDK-----KQFCCNDC
VLGTA-DMKIGEAAKKDKEA--QLAKLGANY-CGNCYGARDQGKNNDTTPRDQWVCCQTC
EIRDESPVMSSAGELEERAR---GRAPDGL-CGSCYALPQDEN-----LDYCCNDC
EISDDLPLLSGAKKFEDVCGPLTEDEISRGVPCGPCYGAVDQTD-----NKRCNDC
EIEDDKPALNSDVSLKELAKGLPEGSQONAY-CGPCYALPQDK-----KQFCCNDC
ELSKE-KFQVGDKS--SKQS--FNE--EGY-CGPCYALDQSKN-DELPQDQKVVCCQTC
: . ** **** * . **:

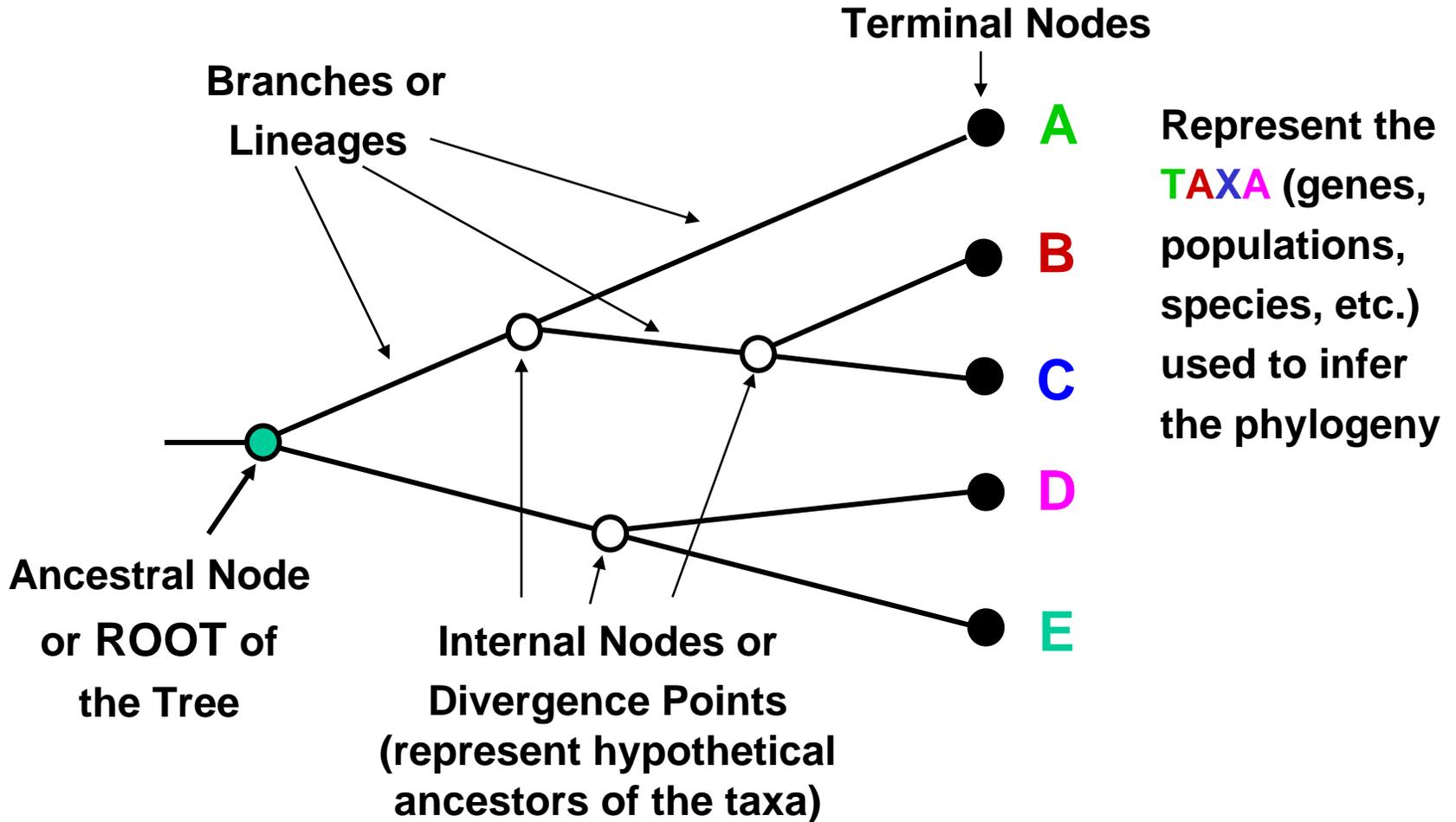
YAL042W
candida586
cdub17784
cgla72177
cgui48535
clus15345
ctro67868
klac20931

DAVRSAYLEAGWAFFDGKNIEQCEREGYVSKINEHLN--EGCRIKGSQINRIQGNLHFA
NTVRRAYAEEKHWSFYDGENIEQCEKEGYVGRRLRERINNNEGCRIGTKINRVSGTMDFA
NTVRRAYAEEKHWSFYDGENIEQCEKEGYVARLRERINNNEGCRIGTKINRVSGTMDFA
DDVRQAYFEKNWAFFDGKIDIEQCEREGYVQKIADQLQ--EGCRVSGSAQLNRIDGNLHFA
ETVRLAYAQAQAWGFFDGENIEQCEREGYVARLNEKINNFEGRIGTKINRISGNLHFA
EAVRMAYAVQEWGFFDGSNIEQCEREGYVEKMSRINNNEGCRIGKSAKINRISGNLHFA
NTVRRAYAEEKQWQFFDGENIEQCEKEGYVKRLRERINNNEGCRIGSTKINRVSGTMDFA
DDVRAAYGQKQWAFKDGKQVEQCEREGYVESINARIH--EGCRVQGRAQLNRIQGTIHFG
: ** ** * * **..:****:**** : :*: ***:.* :*:*:*:*:.

From physiological traits to DNA characters

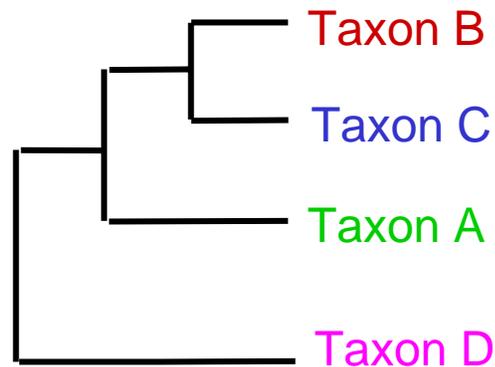
- Traditional phylogenetics
 - Building species trees
 - Small number of traits
 - Hoofs, nails, teeth, horns
 - Well-behaved traits, each arose once
 - Parsimony principle, Occam's razor
- Modern phylogenetics
 - Building gene trees and species trees
 - Very large number of traits
 - Every DNA base and every protein residue
 - Frequently ill-behaved traits
 - Back-mutations are frequent (convergent evolution)
 - Small number of letters, arise many times independently

Common Phylogenetic Tree Terminology



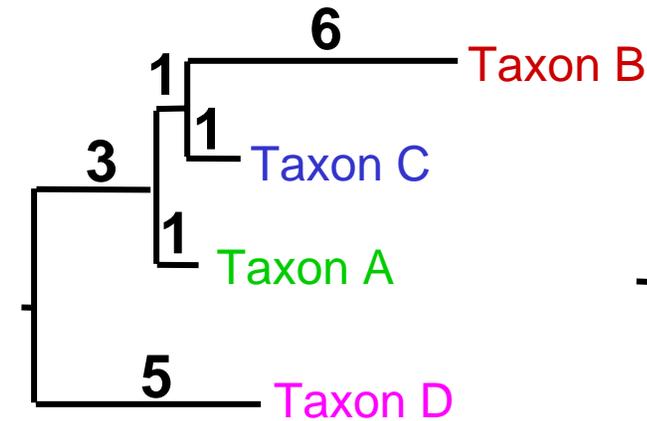
Three types of trees

Cladogram



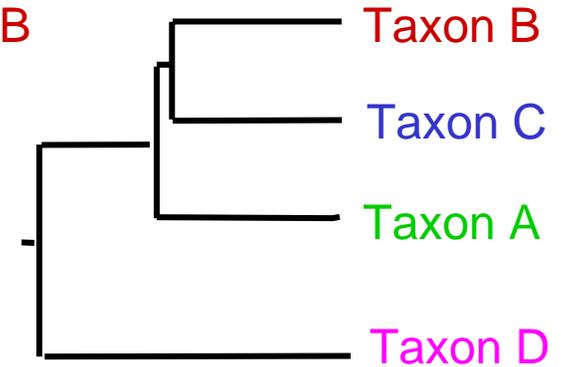
no
meaning

Phylogram



genetic
change

Ultrametric tree



time

All show the same evolutionary relationships, or branching orders, between the taxa.

Molecular phylogenetic tree building methods:

Are mathematical and/or statistical methods for inferring the divergence order of taxa, as well as the lengths of the branches that connect them. There are many phylogenetic methods available today, each having strengths and weaknesses. Most can be classified as follows:

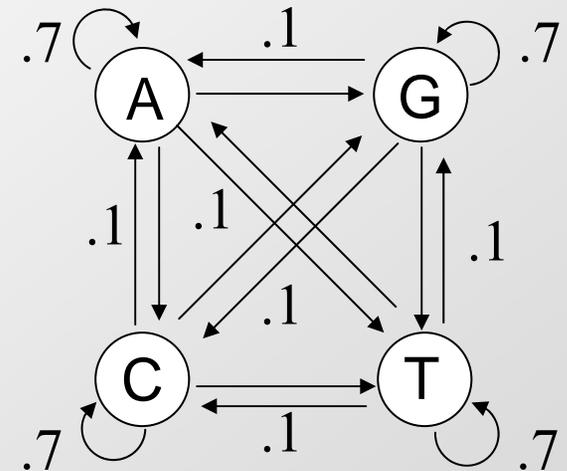
		COMPUTATIONAL METHOD	
		Optimality criterion	Clustering algorithm
DATA TYPE	Characters	PARSIMONY MAXIMUM LIKELIHOOD	
	Distances	MINIMUM EVOLUTION LEAST SQUARES	UPGMA NEIGHBOR-JOINING

2. Modeling evolution

Inferring evolutionary distance

'Evolving' a nucleotide under random model

- At time step 0, start with letter A
- At time step 1:
 - Remain A with probability 0.7
 - Change to C,G,T with prob. 0.1 each
- At time step 2:
 - In state A with probability 0.52
 - Remain A with probability $0.7 * 0.7$
 - Go back to A from C,G,T with $0.1*0.1$ each
 - In states C,G,T with prob. 0.16 each



	t=1	t=2	t=3	t=4	t=5
A	1	0.7	0.52	0.412	0.3472
C	0	0.1	0.16	0.196	0.2176
G	0	0.1	0.16	0.196	0.2176
T	0	0.1	0.16	0.196	0.2176

Modeling Nucleotide Evolution

During infinitesimal time Δt , there is not enough time for two substitutions to happen on the same nucleotide

So we can estimate $P(x | y, \Delta t)$, for $x, y \in \{A, C, G, T\}$

Then let

$$S(\Delta t) = \begin{pmatrix} P(A|A, \Delta t) & \dots & P(A|T, \Delta t) \\ \dots & & \dots \\ P(T|A, \Delta t) & \dots & P(T|T, \Delta t) \end{pmatrix}$$

Modeling Nucleotide Evolution

Reasonable assumption: multiplicative
(implying a stationary Markov process)

$$S(t+t') = S(t)S(t')$$

That is, $P(x | y, t+t') = \sum_z P(x | z, t) P(z | y, t')$

Jukes-Cantor: constant rate of evolution

$$\text{For short time } \varepsilon, S(\varepsilon) = \begin{pmatrix} 1 - 3\alpha\varepsilon & \alpha\varepsilon & \alpha\varepsilon & \alpha\varepsilon \\ \alpha\varepsilon & 1 - 3\alpha\varepsilon & \alpha\varepsilon & \alpha\varepsilon \\ \alpha\varepsilon & \alpha\varepsilon & 1 - 3\alpha\varepsilon & \alpha\varepsilon \\ \alpha\varepsilon & \alpha\varepsilon & \alpha\varepsilon & 1 - 3\alpha\varepsilon \end{pmatrix}$$

Modeling Nucleotide Evolution

Jukes-Cantor:

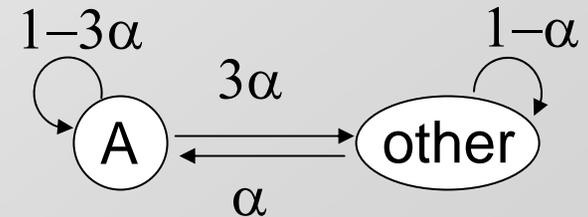
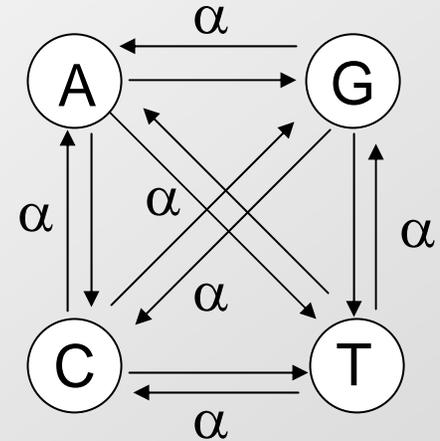
For longer times,

$$S(t) = \begin{pmatrix} r(t) & s(t) & s(t) & s(t) \\ s(t) & r(t) & s(t) & s(t) \\ s(t) & s(t) & r(t) & s(t) \\ s(t) & s(t) & s(t) & r(t) \end{pmatrix}$$

Where we can derive:

$$r(t) = \frac{1}{4} (1 + 3 e^{-4\alpha t})$$

$$s(t) = \frac{1}{4} (1 - e^{-4\alpha t})$$



Modeling Nucleotide Evolution

Kimura:

Transitions: A/G, C/T

Transversions: A/T, A/C, G/T, C/G

Transitions (rate α) are much more likely than transversions (rate β)

$$S(t) = \begin{matrix} & \begin{matrix} A & G & C & T \end{matrix} \\ \begin{matrix} A \\ G \\ C \\ T \end{matrix} & \left(\begin{array}{cccc} r(t) & s(t) & u(t) & u(t) \\ s(t) & r(t) & u(t) & u(t) \\ u(t) & u(t) & r(t) & s(t) \\ u(t) & u(t) & s(t) & r(t) \end{array} \right) \end{matrix}$$

Where

$$s(t) = \frac{1}{4} (1 - e^{-4\beta t})$$

$$u(t) = \frac{1}{4} (1 + e^{-4\beta t} - e^{-2(\alpha+\beta)t})$$

$$r(t) = 1 - 2s(t) - u(t)$$

Distance between two sequences

Given (well-aligned portion of) sequences x^i, x^j ,

Define

d_{ij} = distance between the two sequences

One possible definition:

d_{ij} = fraction f of sites u where $x^i[u] \neq x^j[u]$

Better model (Jukes-Cantor):

$$d_{ij} = -\frac{3}{4} \log(1 - 4f/3)$$

$$r(t) = \frac{1}{4} (1 + 3 e^{-4\alpha t})$$

$$s(t) = \frac{1}{4} (1 - e^{-4\alpha t})$$

Observed $F = [0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7]$

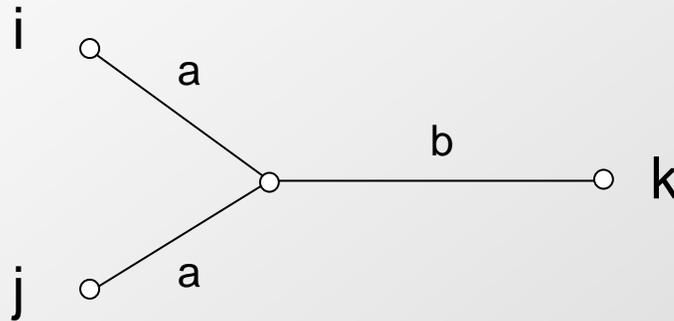
Actual $D = [0.11, 0.23, 0.38, 0.57, 0.82, 1.21, 2.03]$

3. From distances to trees

Ultrametric, additive, and general
distance matrices

3a. Ultrametric distances

- For all points i, j, k
 - two distances are equal and third is smaller
 $d(i,j) \leq d(i,k) = d(j,k)$
 $a+a \leq a+b = a+b$



where $a \leq b$

- Result:
 - All paths from labels are equidistant to the root
 - Rooted tree with uniform rates of evolution

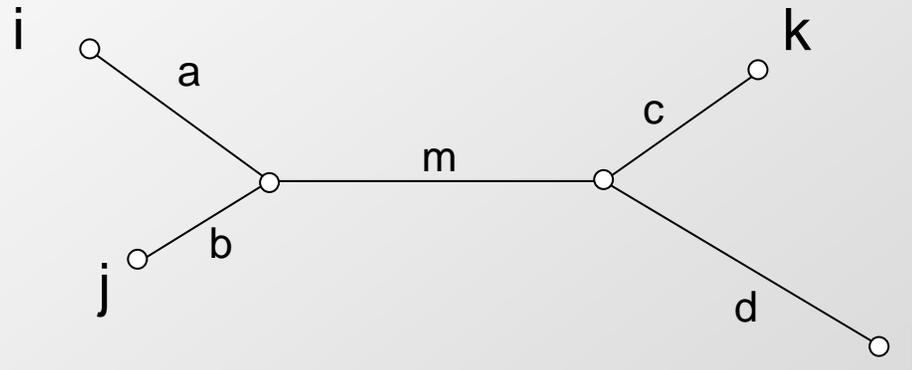
3b. Additive distances

- All distances satisfy the four-point condition

– For all i, j, k, l :

- $d(i, j) + d(k, l) \leq d(i, k) + d(j, l) = d(i, l) + d(j, k)$

- $(a+b) + (c+d) \leq (a+m+c) + (b+m+d) = (a+m+d) + (b+m+c)$



- Result:

– All pairwise distances obtained by traversing a tree

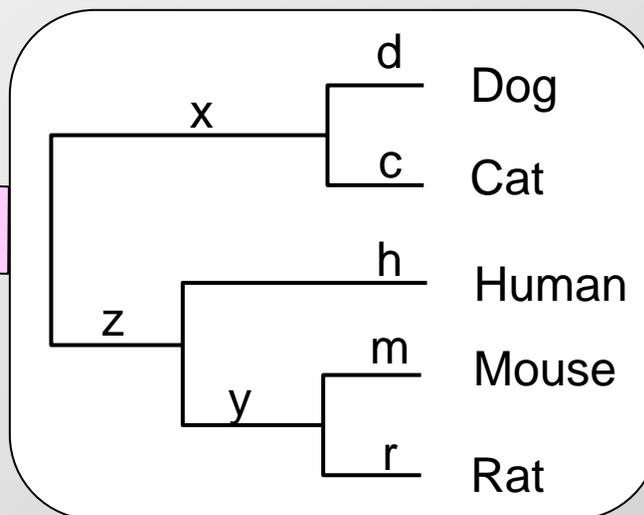
3c. General distances

- In practice, a distance matrix is neither ultrametric nor additive
 - Noise
 - Measured distances are not exact
 - Evolutionary model is not exact
 - Fluctuations
 - Regions used to measure distances not representative of the species tree
 - Gene replacement (gene conversion), lateral transfer
 - Varying rates of mutation can lead to discrepancies
- In the general case, tree-building algorithms generate an approximation to the distance matrix
 - Such a tree can be obtained by
 - Enumeration and scoring of all trees (too expensive)
 - Neighbor-Joining (typically gives a good tree)
 - UPGMA (typically gives a poor tree)

Distance matrix \Leftrightarrow Phylogenetic tree

	Hum	Mou	Rat	Dog	Cat
Human	0	4	5	7	6
Mouse	h.y.m	0	3	8	5
Rat	h.y.r	m.r	0	9	7
Dog	h.z.x.d	m.y.z.x.d	r.y.z.x.d	0	2
Cat	h.z.x.c	m.y.z.x.c	r.y.z.x.c	d.c	0

Tree implies
a distance matrix
 M_{ij}



Map distances D_{ij}
to a tree

$$\min \sum_{ij} (D_{ij} - M_{ij})^2$$

Goal:

Minimize discrepancy between **observed distances** and **tree-based distances**

4. Tree-building algorithms

Mapping a distance matrix to a tree

4a: UPGMA (aka. Hierarchical Clustering)

(Unweighted Pair Group Method with Arithmetic mean)

Initialization:

Assign each x_i into its own cluster C_i

Define one leaf per sequence, height 0

Iteration:

Find two clusters C_i, C_j s.t. d_{ij} is min

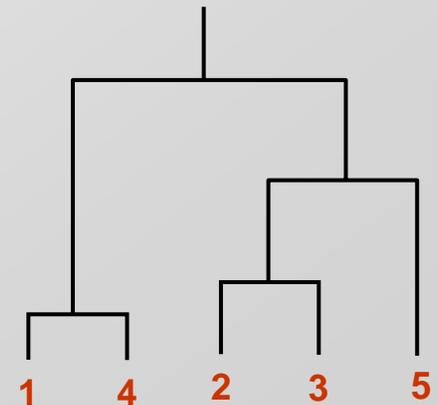
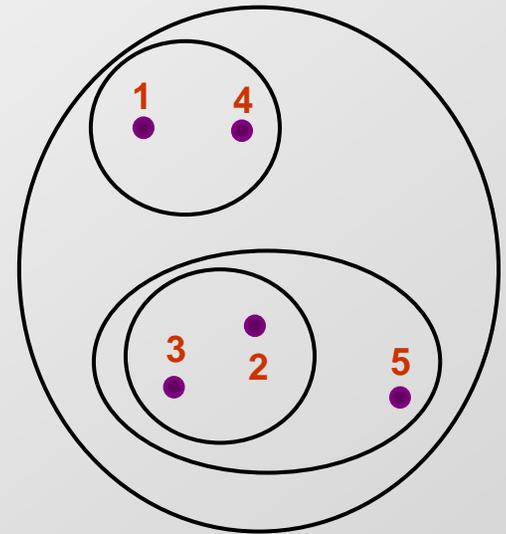
Let $C_k = C_i \cup C_j$

Define node connecting C_i, C_j ,
& place it at height $d_{ij}/2$

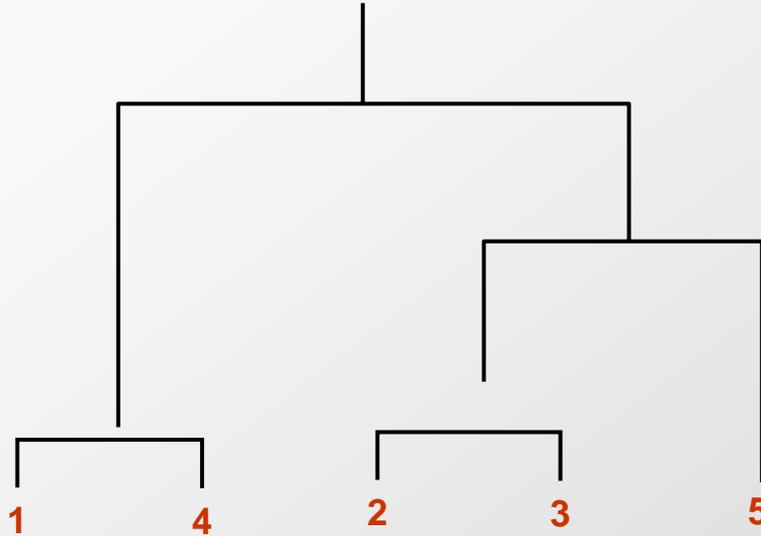
Delete C_i, C_j

Termination:

When two clusters i, j remain,
place root at height $d_{ij}/2$



Ultrametric Distances & UPGMA



UPGMA is guaranteed to build the correct tree if distance is ultrametric

Proof:

1. The tree topology is unique, given that the tree is binary
2. UPGMA constructs a tree obeying the pairwise distances

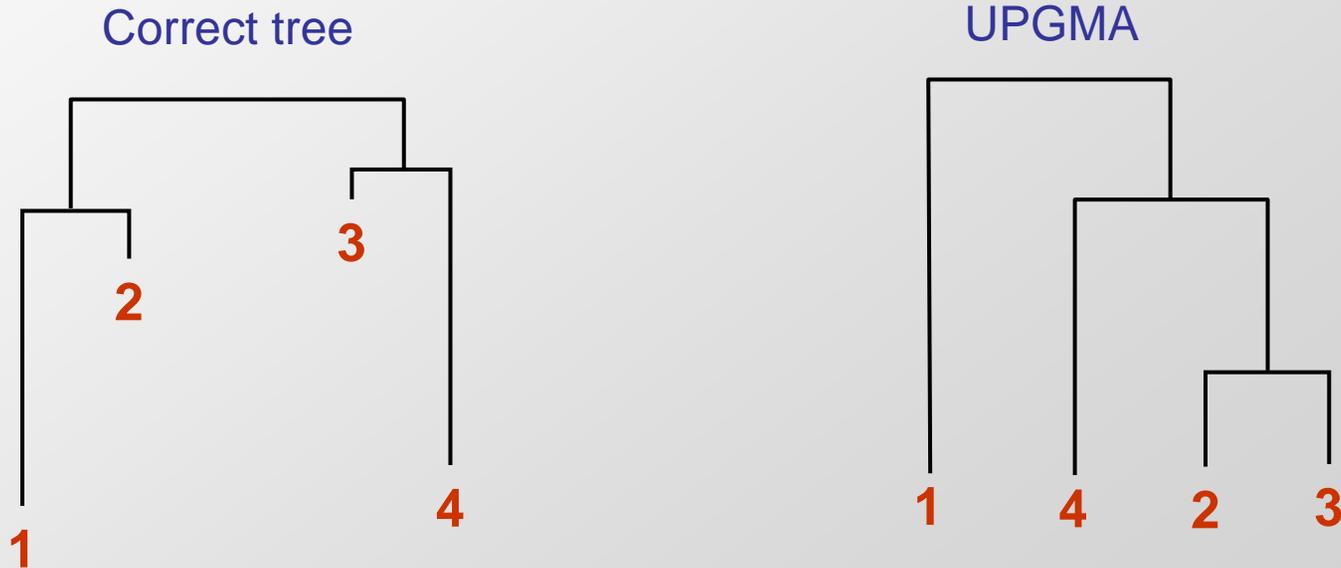
Weakness of UPGMA

Molecular clock assumption:

implies time is constant for all species

However, certain species (e.g., mouse, rat) evolve much faster

Example where UPGMA messes up:



4b. Neighbor-Joining

- Guaranteed to produce the correct tree if distance is additive
- May produce a good tree even when distance is not additive

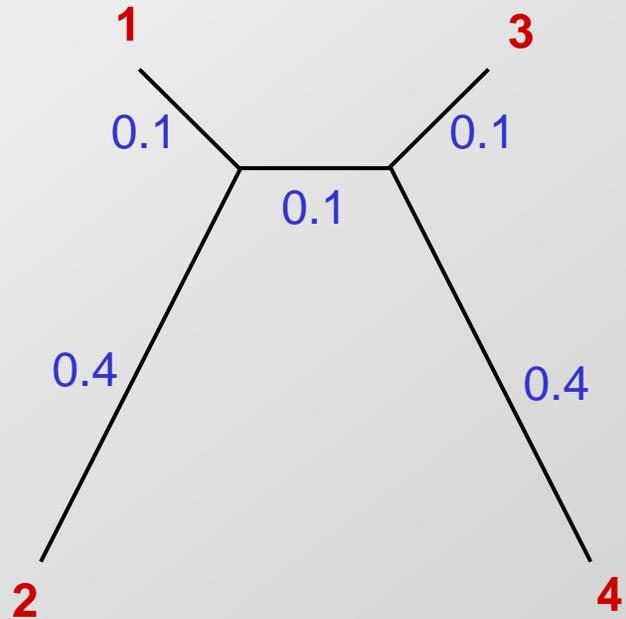
Step 1: Finding neighboring leaves

Define

$$D_{ij} = d_{ij} - (r_i + r_j)$$

Where

$$r_i = \frac{1}{|L| - 2} \sum_k d_{ik}$$



Claim: The above “magic trick” ensures that D_{ij} is minimal **iff** i, j are neighbors

Proof: Beyond the scope of this lecture (Durbin book, p. 189)

Algorithm: Neighbor-joining

Initialization:

Define T to be the set of leaf nodes, one per sequence

Let $L = T$

Iteration:

Pick i, j s.t. D_{ij} is minimal

Define a new node k , and set $d_{km} = \frac{1}{2} (d_{im} + d_{jm} - d_{ij})$ for all $m \in L$

Add k to T , with edges of lengths $d_{ik} = \frac{1}{2} (d_{ij} + r_i - r_j)$

Remove i, j from L ;

Add k to L

Termination:

When L consists of two nodes, i, j , and the edge between them of length d_{ij}

5. Alignment-based algorithms

Parsimony (set-based)

Parsimony (Dynamic Programming)

Maximum Likelihood

5a. Parsimony

- One of the most popular methods

Idea:

Find the tree that explains the observed sequences with a minimal number of substitutions

Two computational sub-problems:

1. Find the parsimony cost of a given tree (easy)
2. Search through all tree topologies (hard)

Parsimony Scoring

Given a tree, and an alignment column

Label internal nodes to minimize the number of required substitutions

Initialization:

Set cost $C = 0$; $k = 2N - 1$

Iteration:

If k is a leaf, set $R_k = \{ x^k[u] \}$

If k is not a leaf,

Let i, j be the daughter nodes;

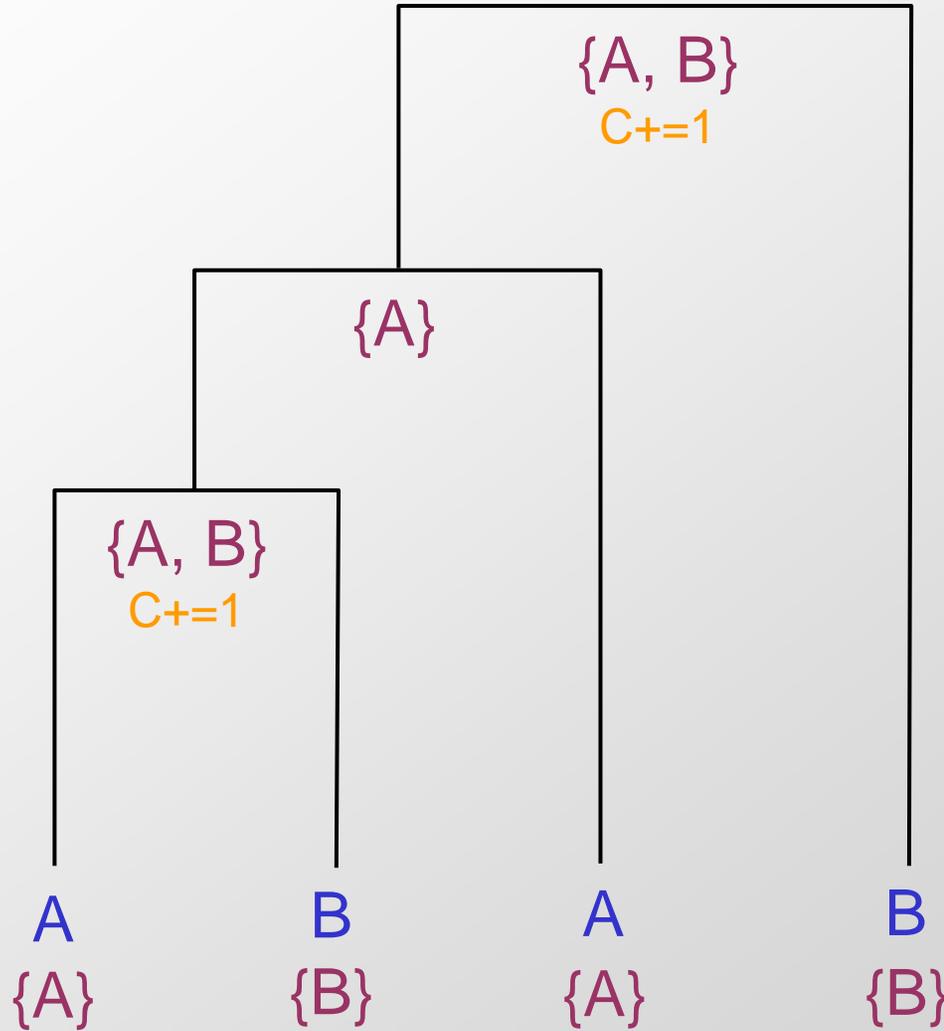
Set $R_k = R_i \cap R_j$ if intersection is nonempty

Set $R_k = R_i \cup R_j$, and $C += 1$, if intersection is empty

Termination:

Minimal cost of tree for column u , = C

Example



Traceback to find ancestral nucleotides

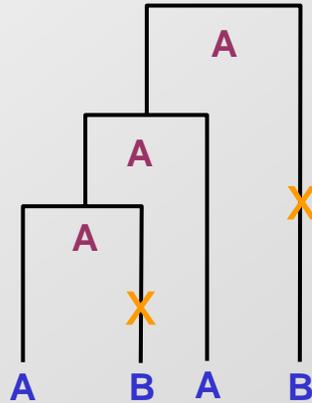
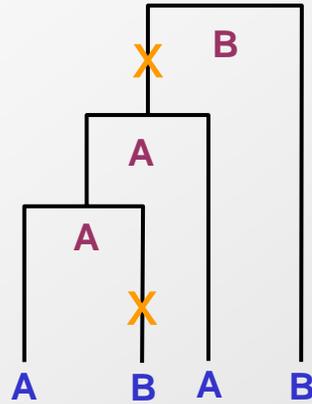
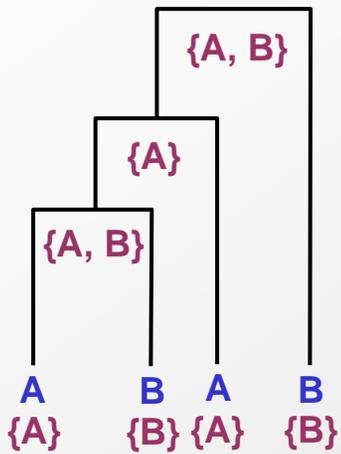
Traceback:

1. Choose an arbitrary nucleotide from R_{2N-1} for the root
2. Having chosen nucleotide r for parent k ,
If $r \in R_i$ choose r for daughter i
Else, choose arbitrary nucleotide from R_i

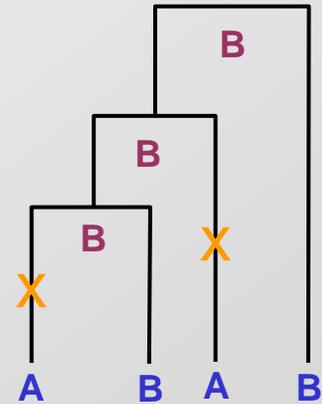
Easy to see that this traceback produces some assignment of cost C

Example

Accessible to traceback



Still optimal, but not found by traceback

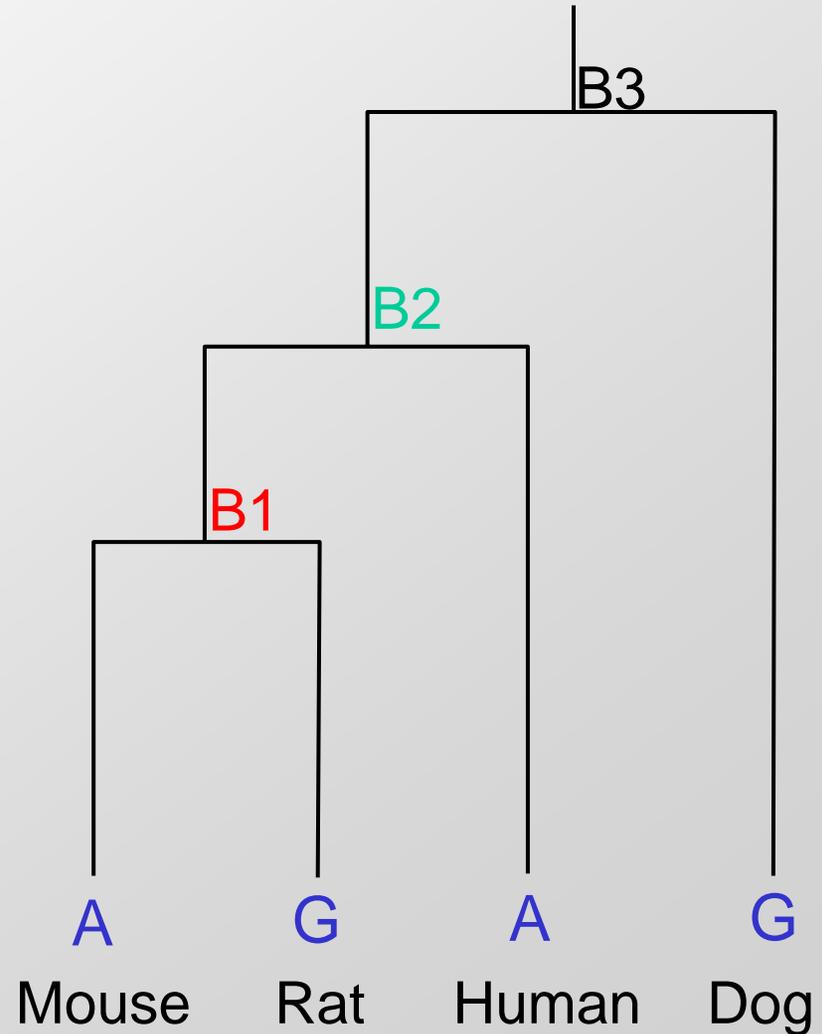


5b. Parsimony with dynamic programming

	M	R	B1	H	B2	D	B3
A	0	1	1	0	1	1	2
C	1	1	2	1	3	1	4
G	1	0	1	1	2	0	2
T	1	1	2	1	3	1	4



- Each cell (N,C) represents the min cost of the subtree rooted at N, if the label at N is C.
- Update table by walking up the tree from the leaves to the root, remembering max choices.
- Traceback from root to leaves to construct a min cost assignment



5c. Maximum Likelihood Methods

Input: Proposed topology T

Output: Prob. that proposed tree gave rise to observed data

Search: Heuristic MCMC search for max likelihood tree.

$$\begin{aligned} B^{\wedge}, T^{\wedge} &= \operatorname{argmax}_{B, T} P(D, B, T) \\ &= \operatorname{argmax}_{B, T} P(D|B, T) P(B, T) \end{aligned}$$

Likelihood $P(\text{Data}|\text{BranchLengths}, \text{Topology})$

Prior $P(B, T)$: typically uniform/can use to guide search

Iterate: Iterate over proposed topologies.

- Given current topology T , branch lengths B :
 - Propose many alternative (T', B') , by modifying existing $T \rightarrow T'$, and inferring branch lengths B' that maximize $P(D|B', T')$
 - Evaluate $P(D|B, T)$ and $P(D|B', T')$
 - Select one T' at random based on increase in likelihood
- Heuristics for proposing new topology T'
 - Nearest-neighbor interchange, subtree cut-and-paste, rotations

Advantages/disadvantages of ML methods

- Advantages:

- Are inherently statistical and evolutionary model-based.
- Usually the most ‘consistent’ of the methods available.
- Can be used for character (can infer the exact substitutions) and rate analysis.
- Can be used to infer the sequences of the extinct (hypothetical) ancestors.
- Can help account for branch-length effects in unbalanced trees.
- Can be applied to nucleotide or amino acid sequences, and other types of data.

- Disadvantages:

- Are not as simple and intuitive as many other methods.
- Are computationally very intense (limits number of taxa and length of sequence).
- Like parsimony, can be fooled by high levels of homoplasy.
- Violations of the assumed model can lead to incorrect trees.

Bootstrapping to get the best trees

Main outline of algorithm

1. Select random columns from a multiple alignment – one column can then appear several times
2. Build a phylogenetic tree based on the random sample from (1)
3. Repeat (1), (2) many (say, 1000) times
4. Output the tree that is constructed most frequently

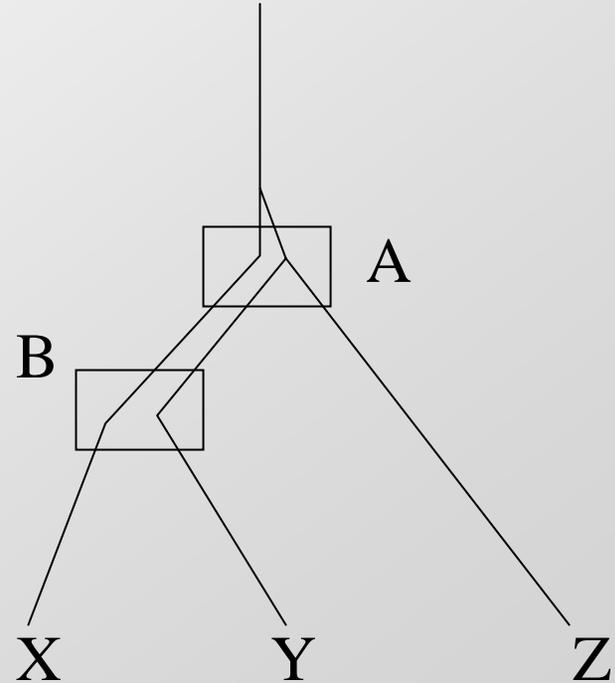
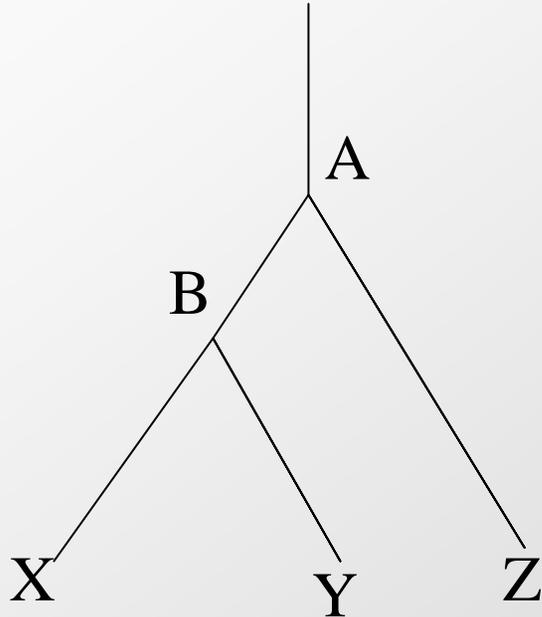
Summary

- Basics of phylogeny
 - Characters, traits, nodes, branches, lineages
 - Gene trees, species trees
- Modeling sequence evolution
 - Turning sequence data into distances
 - Probabilistic models of nucleotide divergence
 - Jukes-Cantor 1-parameter model, Kimura 2-parameter model
- From distances to trees
 - Ultrametric, Additive, General Distances
 - UPGMA, Neighbor Joining, guarantees and limitations
 - Least-squared error, minimum evolution
- From alignments to trees
 - Parsimony methods: set-based vs. dynamic programming
 - Maximum likelihood methods
 - MCMC and heuristic search

Extra Time?

Recitation tomorrow: Gene vs. Species evolution

- Genes can start diverging before species separate
 - Genetic polymorphism within population could exist
 - After divergence, forms evolve differently in each species
 - Gene divergence could predate species divergence
 - Gene tree topology could be misleading



- Solution: Use multiple genes to infer a species tree

Phylogenomics

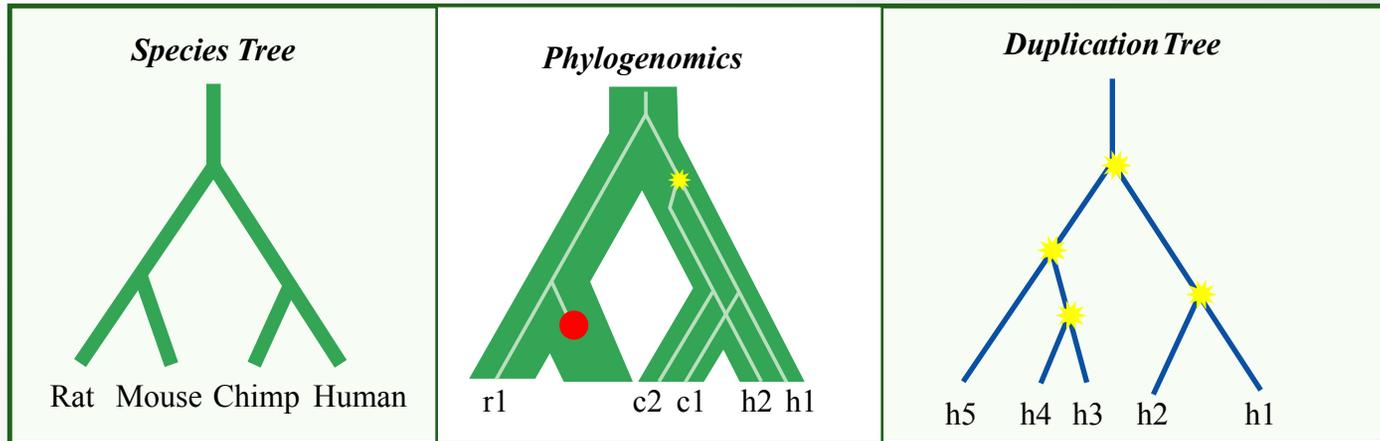


Figure by MIT OpenCourseWare.

Many species Many species One species
One gene Many genes Many genes

- Traditional phylogenetics focused on uniform trees
 - Any topology makes a good story
- Phylogenomics imposes additional constraints
 - Gene trees evolve inside species trees
 - Errors imply large-scale duplications and losses

Extending traditional max likelihood methods

- Traditional max likelihood (phylogenetics)

$$B^{\wedge}, T^{\wedge} = \operatorname{argmax}_{B, T} P(D, B, T)$$
$$\operatorname{argmax}_{B, T} P(D|B, T) \cancel{P(B, T)}$$

- Extended likelihood function (phylogenomics)

$$B^{\wedge}, T^{\wedge} = \operatorname{argmax}_{B, T} P(D, B, T, R | \mathbf{E})$$
$$\operatorname{argmax}_{B, T} P(D|B, T) P(B|T, R, \mathbf{E}) \cancel{P(R|T, \mathbf{E})} \cancel{P(T|\mathbf{E})}$$

Likelihood of data given
proposed branch lengths

Likelihood of proposed branch
lengths (given species evolution)

Evaluation: Large increase in accuracy

program	accuracy
SPIML	71.7%
SPIDIR	43.2%
PHYML	23.5%
BIONJ	29.6%
MrBayes	31.4%
DNAPARS	22.2%

Syntenic regions:

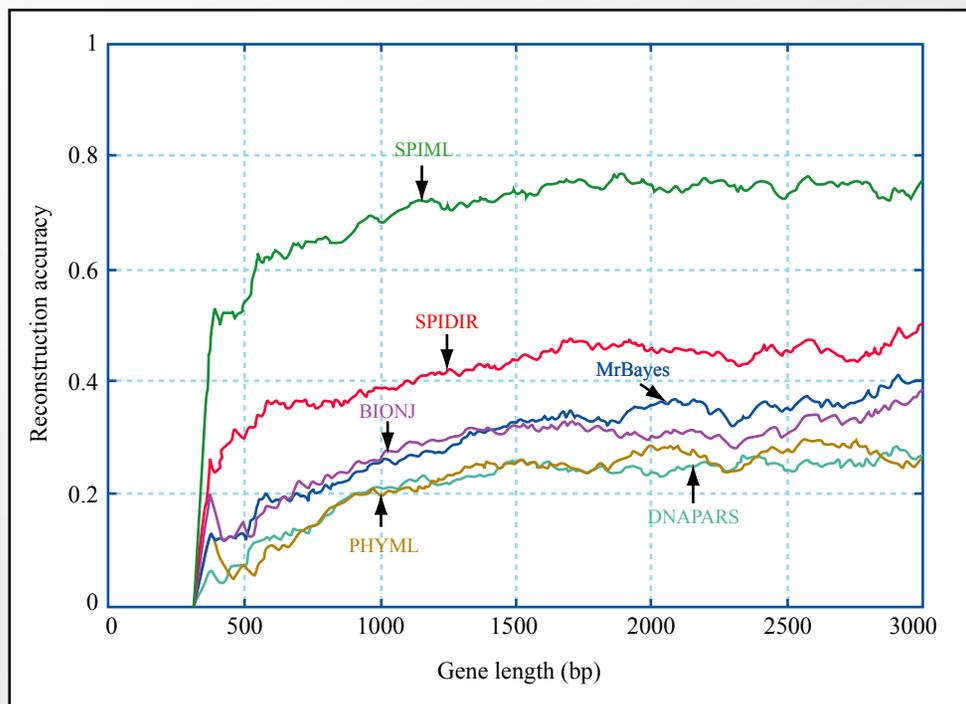
Increasing number of species

Diverse lineages

→ Great increase in accuracy

- Simulated data:

- Run (generative) model
- 1 dup event ⇔ many dup genes
- Method robust to dup/loss



Length correlation:

Maximally use data

More data → closer to truth

Still depends on gene length

But much higher than other methods