#### A PHYLOGENETICS SURVEY

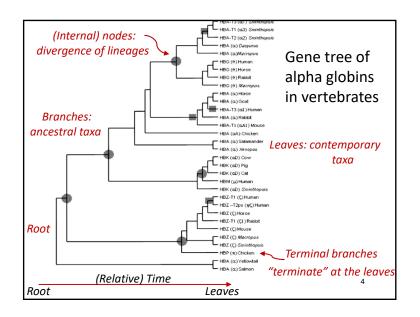
Dr. Maureen Stolzer October 2<sup>nd</sup>, 2025

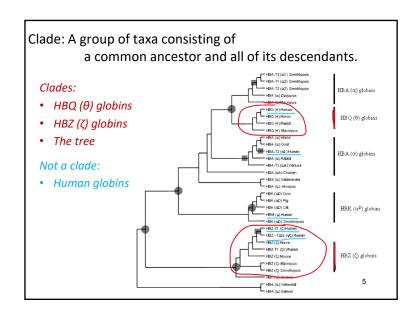
#### Information encoded in trees

- Tree terminology
- Branch lengths
- Similarity, common ancestry and relatedness
- Binary and Non-binary trees
- Rooted and unrooted trees

3

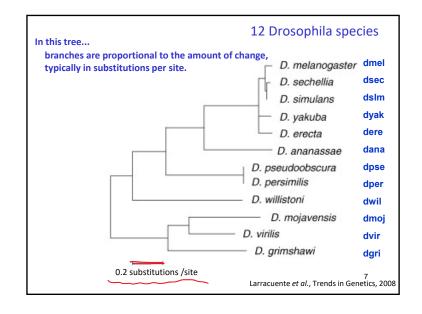
#### TREE TERMINOLOGY

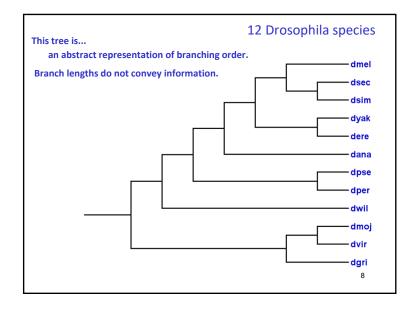


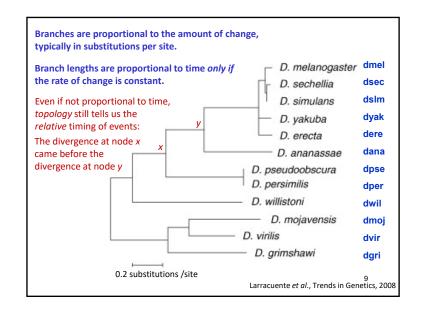


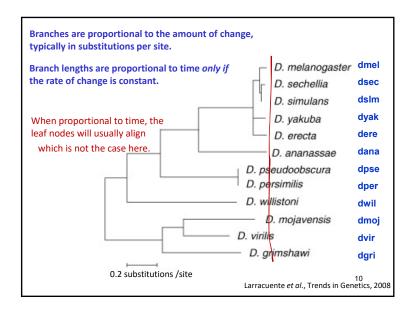
#### Information encoded in trees

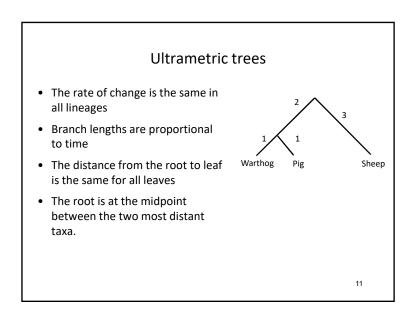
- Tree terminology
- ➤ Branch lengths vs topology
- Similarity, common ancestry and relatedness
- Binary and Non-binary trees
- Rooted and unrooted trees

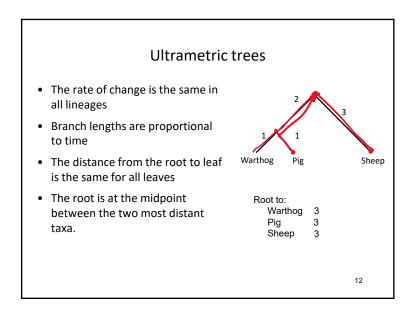


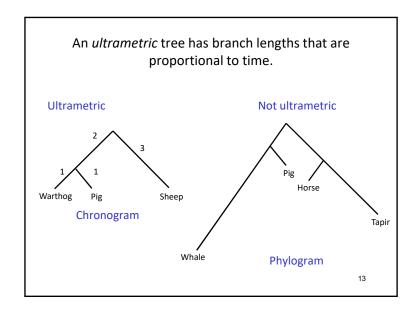












#### Information encoded in trees

- Tree terminology
- · Branch lengths
- Similarity, common ancestry and relatedness
- · Binary and Non-binary trees
- Rooted and unrooted trees

14

#### Similarity versus common ancestry

More closely related taxa are not guaranteed to be more similar except in ultrametric trees, where the rate of change is proportional to time.

Whale

 The pig is more closely related to the whale than to the horse (common ancestor is closer to the leaves)

• The pig is *more similar* to the horse than to the whale (path is shorter)

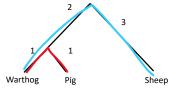
Pig Horse Tapir

#### Similarity versus common ancestry

More closely related taxa are not guaranteed to be more similar except in ultrametric trees, where the rate of change is proportional to time.

an ultrametric tree

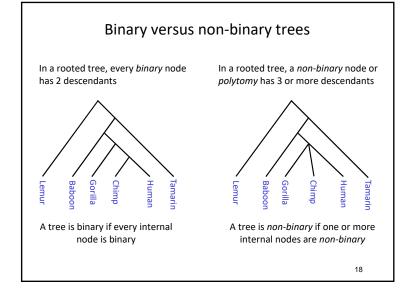
- The warthog is more closely related to the pig than to the sheep
- The warthog is more similar to the pig than to the sheep



#### Information encoded in trees

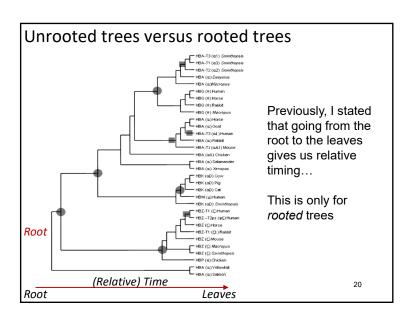
- Tree terminology
- Branch lengths
- Similarity, common ancestry and relatedness
- Binary and Non-binary trees
- Rooted and unrooted trees

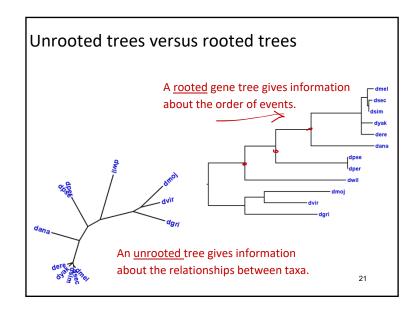
17

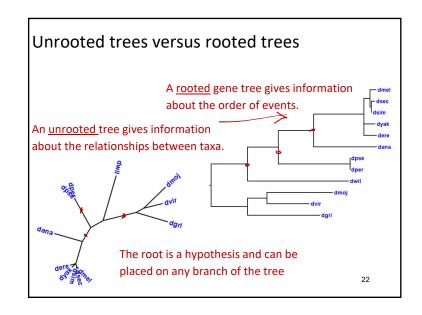


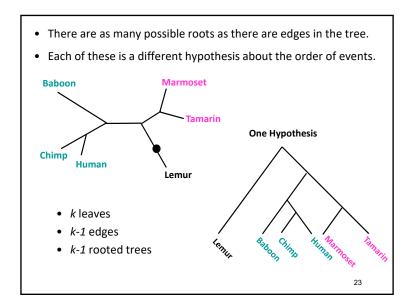
#### Information encoded in trees

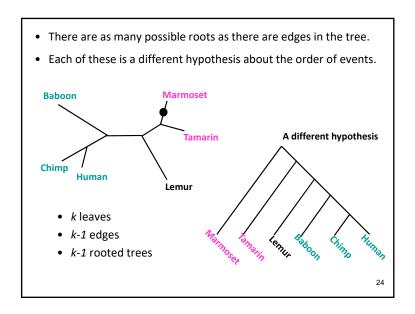
- Tree terminology
- · Branch lengths
- Similarity, common ancestry and relatedness
- · Binary and Non-binary trees
- > Rooted and unrooted trees

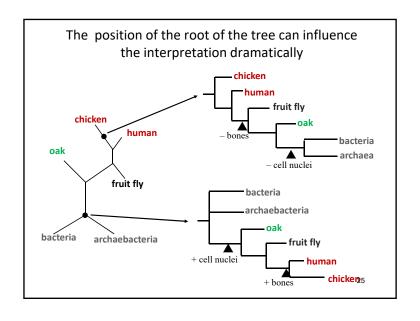


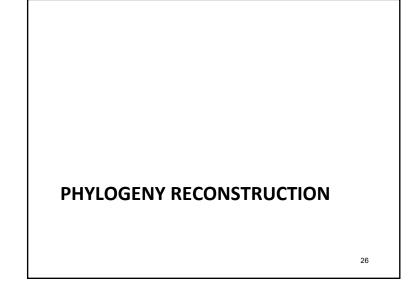


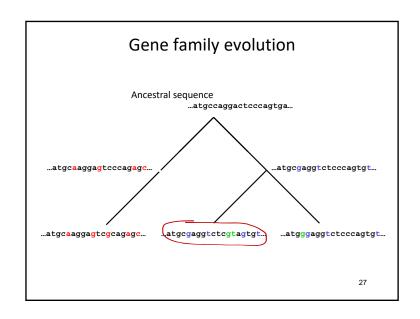


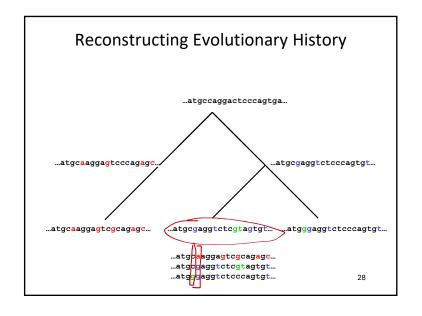


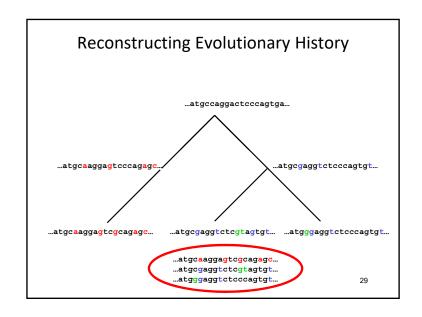


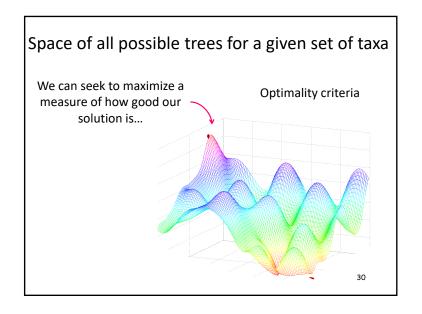


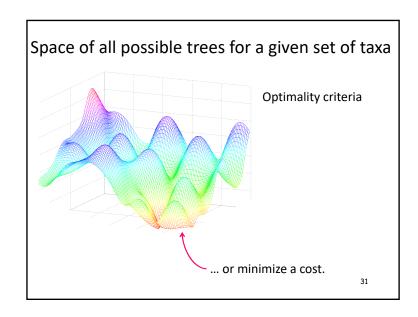




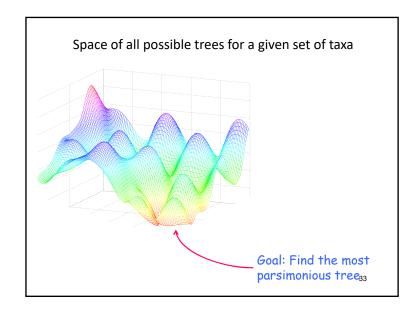












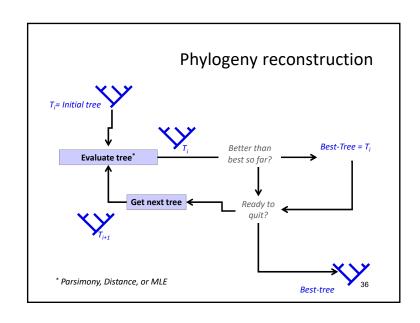
#### Tree reconstruction

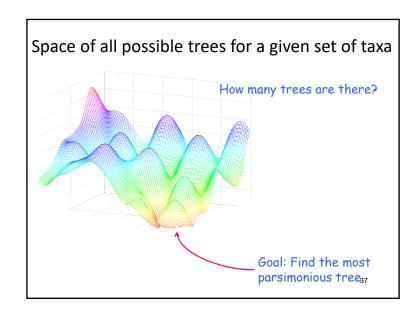
- Construct a tree that fits the data
   This is only possible when the data satisfies some very restrictive conditions.
- 2. Score each possible tree with *k* leaves and select the best one

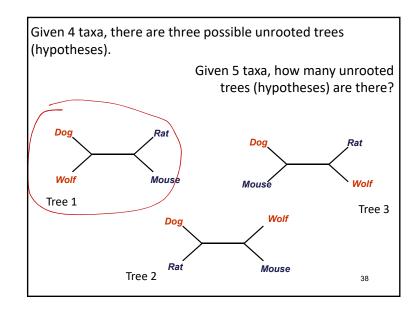
34

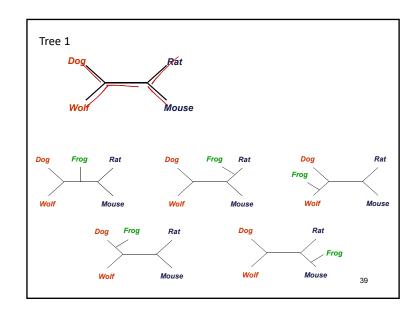
#### Tree reconstruction

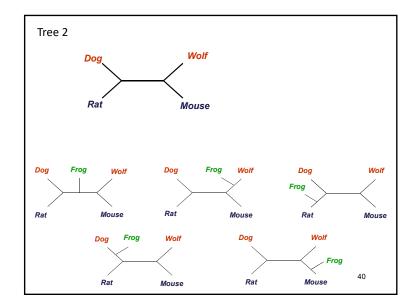
- Construct a tree that fits the data
   This is only possible when the data satisfies some very restrictive conditions
- 2. Score each possible tree with *k* leaves and select the best one

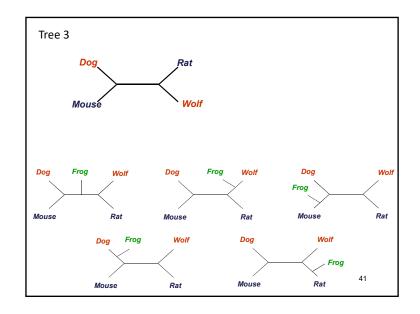


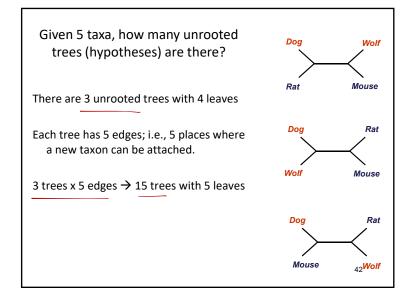












Given 5 taxa, how many unrooted trees (hypotheses) are there?

There are 3 unrooted trees with 4 leaves

Each tree has 5 edges; i.e., 5 places where a new taxon can be attached.

3 trees x 5 edges → 15 trees with 5 leaves

15 t x 7e

Wolf

Rat

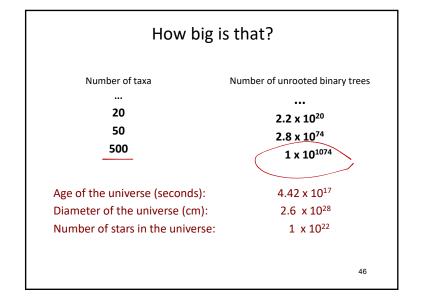
Mouse

In general, how many hypotheses must we consider for *k* leaf taxa?

- A tree with *k* leaves, has 2*k* 3 edges
- For a given tree with k leaves, we can add a new taxon to any edge to get a different tree with k + 1 leaves.
- That gives us 2k 3 new trees for each tree with k leaves

*Ergo,* the number of hypotheses we need to consider gets big <u>very</u>, <u>very quickly</u>, as the number of taxa increases

The number of trees gets big fast				
Number of taxa	Number of unrooted binary trees			
3	1			
4	3			
5	15			
6	105			
10	2,027,025			
20	2.2 x 10 <sup>20</sup>			
50	2.8 x 10 <sup>74</sup>			
500	1 x 10 <sup>1074</sup>			
	45			



In general, how many hypotheses must we consider for *k* leaf taxa?

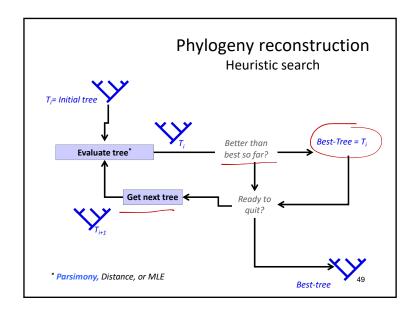
The number of hypotheses we need to consider gets big <u>very</u>, <u>very quickly</u>, as the number of taxa increases.

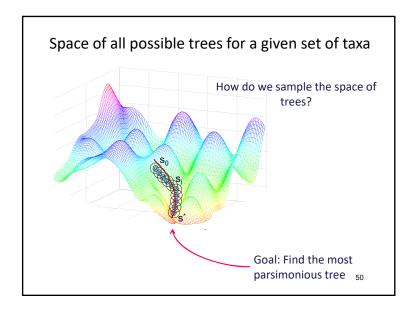
- If k is small (e.g.,  $k \le 10$ ), we can consider all hypotheses.
- $\triangleright$  Otherwise (k > 10), sample a subset of all possible hypotheses.

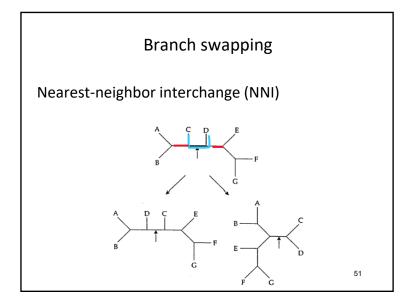
"Heuristic search"

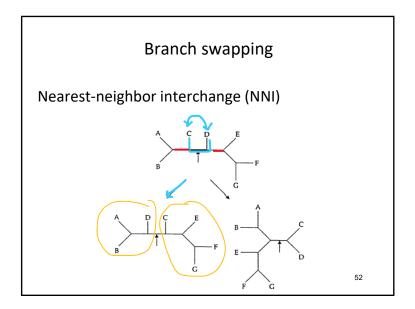
47

#### **SEARCHING TREE SPACE**



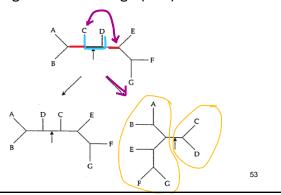






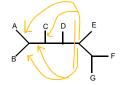
#### Branch swapping

Nearest-neighbor interchange (NNI)



#### **Branch swapping**

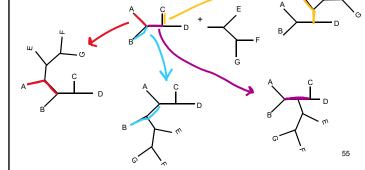
Subtree prune and regraft (SPR)



54

#### Branch swapping

Subtree prune and regraft (SPR)

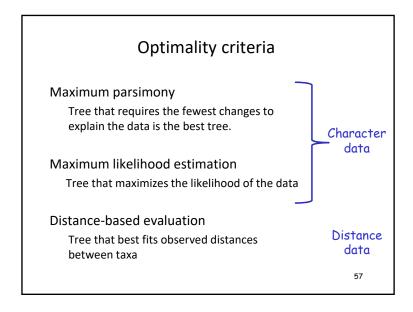


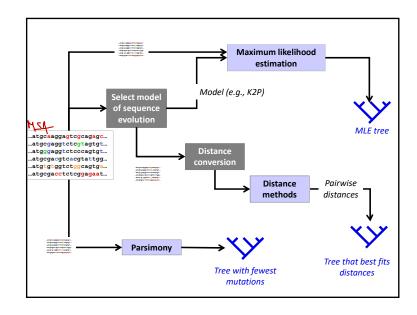
#### Phylogeny reconstruction

- Design a mathematical function for scoring evolutionary trees such that
  - Good trees get good scores

Optimality criterion

- Bad trees get bad scores
- Criterion
- Find the *optimal* tree; the tree with the best score
- Potential problems:
  - The optimal tree\_can be hard to find NP Hard
  - A good numerical score does not always translate to a good evolutionary hypothesis.





#### **PARSIMONY**

59

#### The parsimony criterion

- In regular speech, "parsimonious" means excessively thrifty or sparing, stingy, or frugal
- Phylogenetics: the most parsimonious tree is the tree that requires the fewest substitutions to explain the data.

62

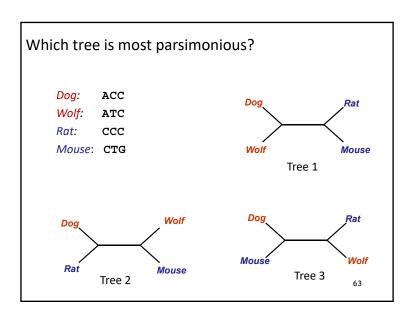
#### Evolutionary change on a tree

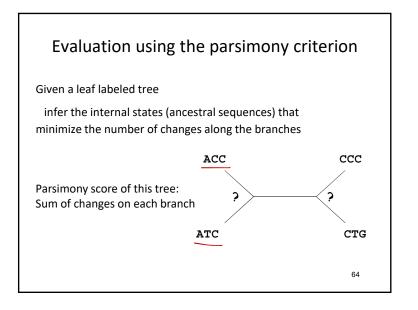
- Given
  - a tree
  - a set of characters that are variable for these taxa,
  - a character state matrix for the leaf taxa
- infer
  - the character states of each ancestral node and
  - the state changes along each branch
- such the *number of changes required is minimal*

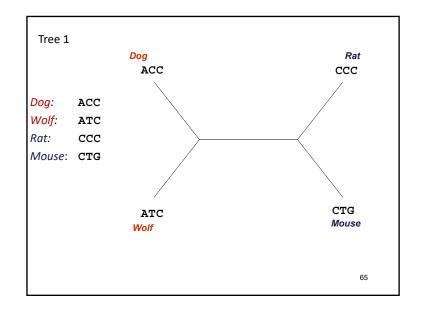
Parsimony

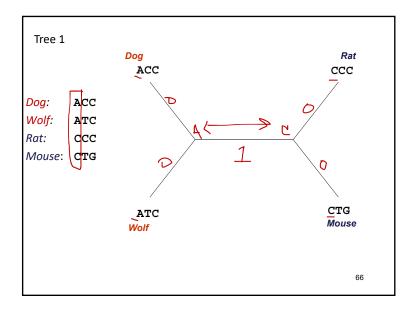
61

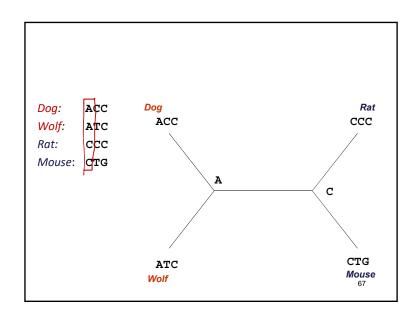
# An example: Given - 4 taxa - data describing each taxon find the most parsimonious tree with 4 leaves.

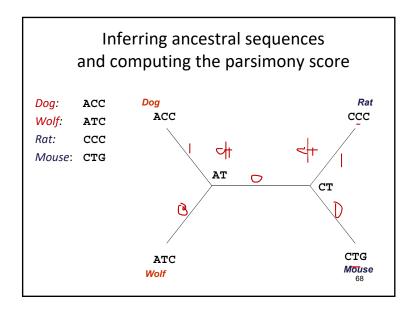


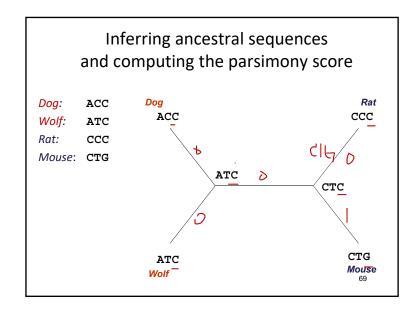


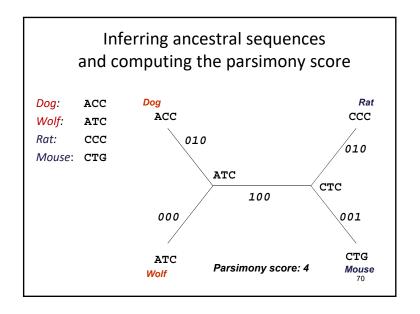


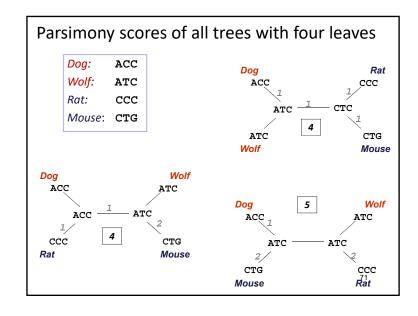


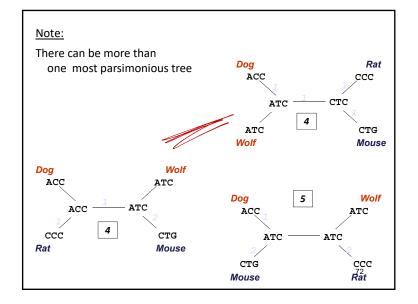












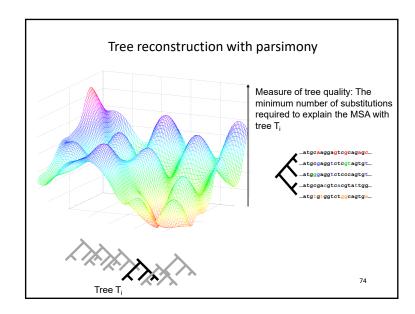
#### Phylogeny reconstruction

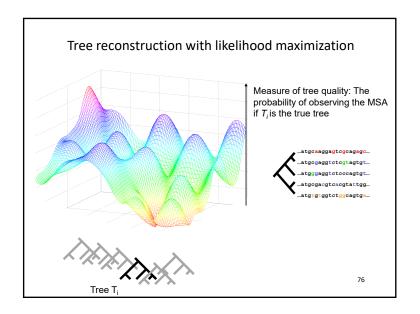
Potential problems with Maximum Parsimony:

- There can be more than one most parsimonious tree.
- · Not all sites are informative
- The assumption that mutations are rare may be wrong.

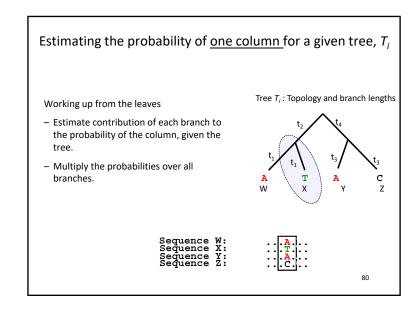
73

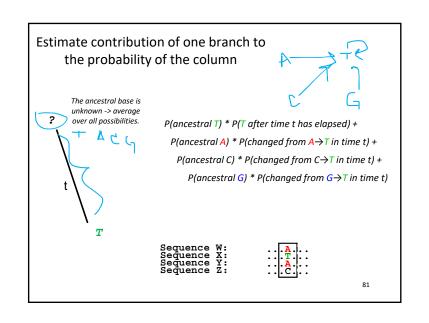
#### **MAXIMUM LIKELIHOOD**

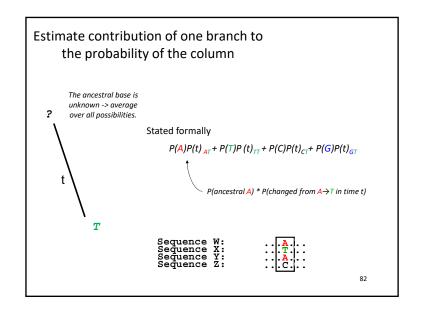


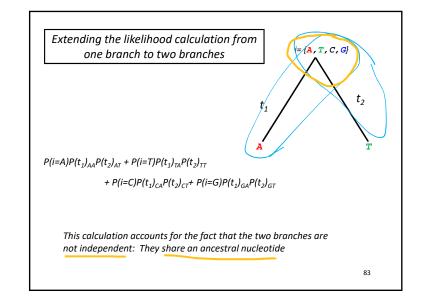


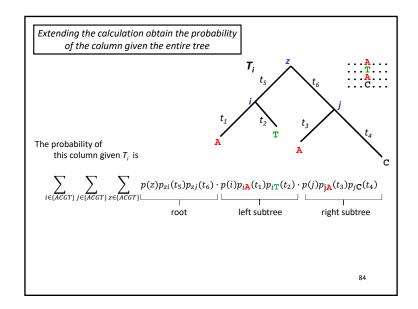
# Given - contemporary taxa: k genes, strains, species... - data describing each taxon ...atgcaaggagtcgcagagc... ...atgcaaggagtccgtagtgt... ...atgcgaggtctcgtagtgt... atgcgaggtctcccagtgt... find the tree with k leaves that explains the data. ➤ Find the best tree with respect to an optimality criterion

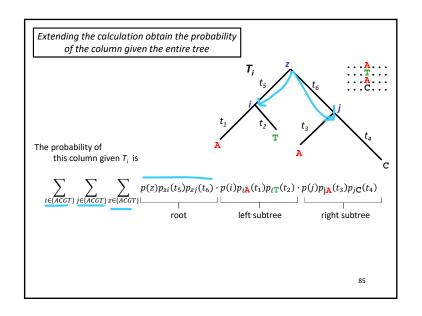


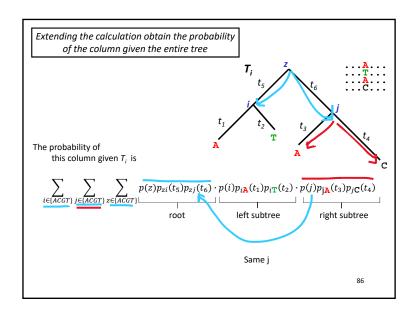


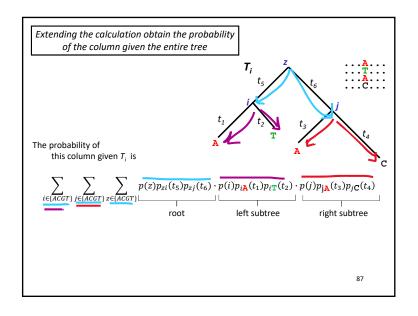


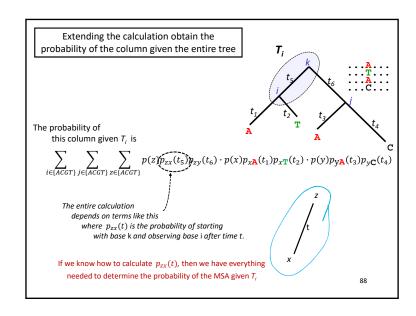


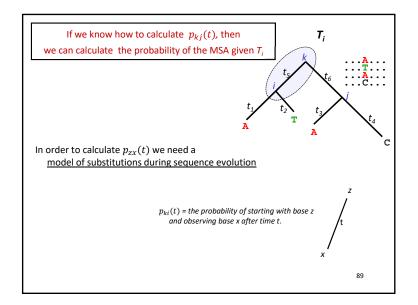








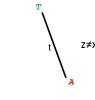




We discussed how to solve this problem with the Jukes Cantor model:

Given a site evolving according to Jukes Cantor with parameter *a*, what is the probability of observing *z* at time 0 and *x* at time *t*?



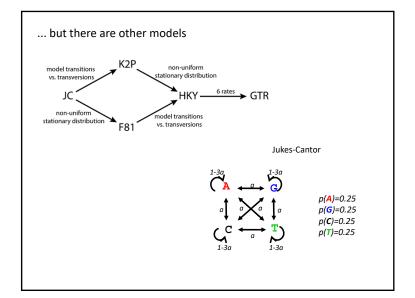


$$p_{xx}(\alpha, t) = \frac{1}{4} + \frac{3}{4}e^{-4\alpha t}$$

$$p_{xy}(\alpha, t) = \frac{1}{4} - \frac{1}{4}e^{-4\alpha t}$$

From last class

90



To estimate the probability of an MSA given a tree requires a model of substitutions during sequence evolution

- · Which evolutionary model?
  - Do all substitutions have the same rate?
- · What rates?
- Do all sites have the same rate?
- Are there invariant sites that never change?

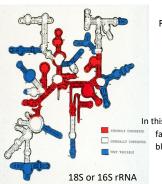
To estimate the probability of an MSA given a tree requires a model of substitutions during sequence evolution

- Which evolutionary model?
  - Do all substitutions have the same rate?
- · What rates?
- > Do all sites have the same rate?
- Are there invariant sites that never change?

93

#### Site heterogeneity:

#### Different sites may be changing at different rates



Rates may be site specific, depending on functional and structural constraints (e.g., codon position, or alpha helix etc.)

In this example, white sites are changing
faster than red sites and more slowly than
blue sites.

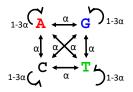
94

#### Modeling rate variation

It is not possible to assign a different rate to each position.

Instead, partition sites into a small number of rate categories, e.g., "slow, medium, fast, supersonic"

Use one model for the entire alignment, but different rates at different sites.

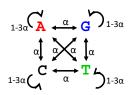


For example, all sites might be modeled with the Jukes Cantor model, but with different parameters:  $\alpha_1 > \alpha_2 > \alpha_3 > \alpha_4$ .

95

#### Modeling rate variation

In addition, some sites may not be changing at all.



<u>Invariant sites:</u> Sites that do not change...

- Examples: very recent divergence, purifying selection, parallel substitutions
- Incorporate invariable sites in model as an extra category. For every site, calculate the likelihood that it is an invariable site.

96

#### Discrete rates model

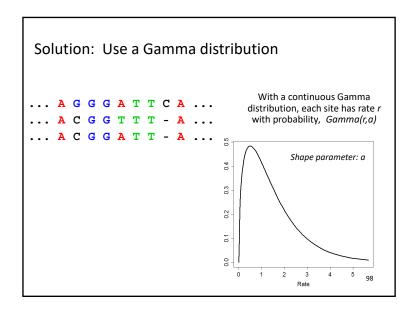
Given k rate classes,

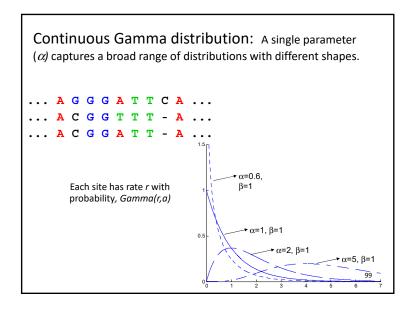
- Determine the rate,  $r_k$ , for each class
- Determine the probability that site *i* is in class *k*.

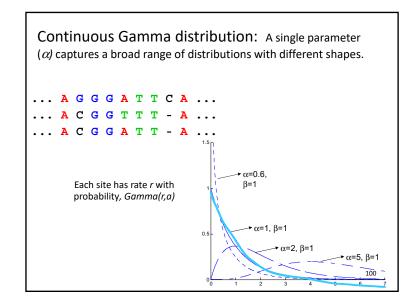
Site class	1	2	3	 К
Probability	$p_1$			 $p_{K}$
Rate	$r_1$			 $r_k$

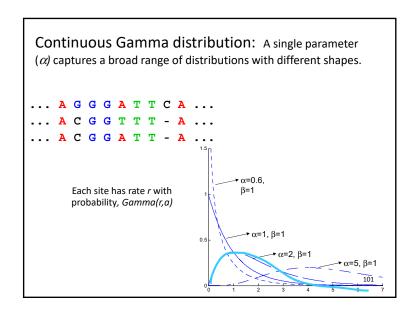
```
r, r, ... r,
... A G G G A T T C A ...
A C G G T T T - A ...
... A C G G A T T - A ...
```

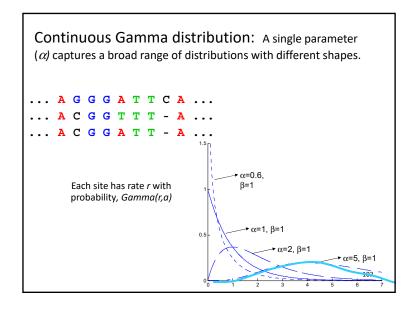
Problem: Too many parameters

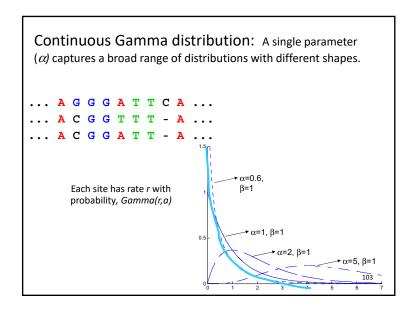


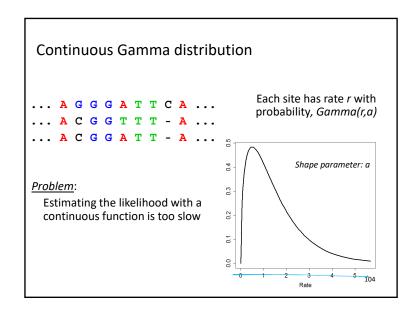


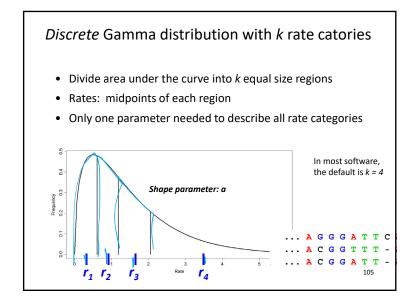










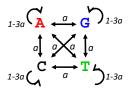


In summary, use models of sequence evolution to capture the following features:

Frequency of substitutions between

bases/residues:

Base/residue frequencies: p(A), p(C), p(G), p(T)



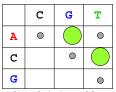
Invariant sites: sites that do not change

#### Multiple rate categories:

- Same model, but with different values of the parameter(s)

106

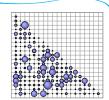
## DNA substitution models



DNA substitution model (e.g., JC, K2P, GTr)

- Four states (A, C, G, T)
- Model specifies the relative rates of substitution for all possible pairs of nucleotides

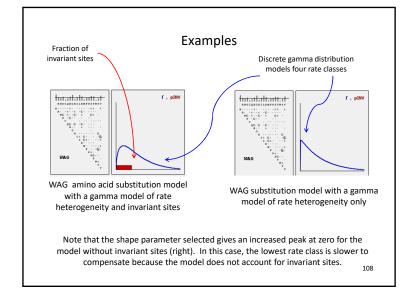
### Amino acid substitution models



Amino acid substitution matrix (e.g., PAM, WAG, JTT, MtREV etc)

- Twenty states (A, C, ... Y)
- Model specifies the relative rates of substitution for all possible pairs of amino acids

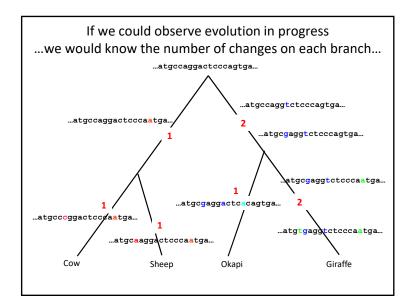
107

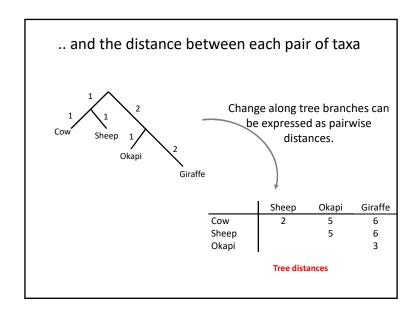


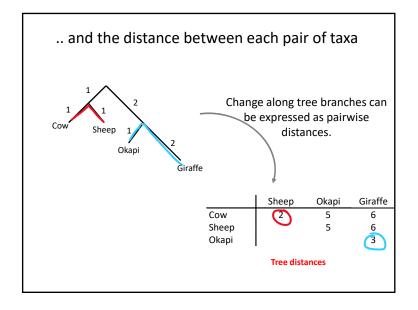
# DISTANCE-BASED TREE RECONSTRUCTION

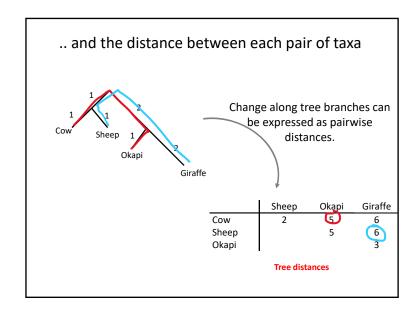
# Phylogeny reconstruction outline

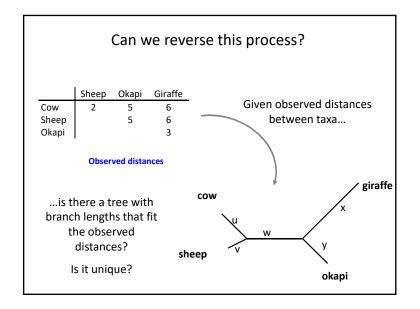
- > Optimality criteria for scoring evolutionary trees
  - Maximum parsimony
  - Maximum likelihood estimation
  - **Distance-based evaluation**
- Find the *optimal* tree; the tree with the best score
  - How many hypotheses (trees with k leaves) are there?
  - Sampling hypotheses: how to search through the space of trees with k leaves?











#### Outline

Distance-based phylogeny reconstruction Convert multiple alignments into distances

Properties of pairwise distances between taxa

- Additive distances
- Ultrametric distances

Using pairwise distances to infer a tree: constructive methods

UPGMA

If you've seen hierarchical clustering, this may seem similar

- Neighbor Joining

116

#### Calculating distances from MSAs

- GCTTGTCCGTTACGAT

- ACTTGACCGTTTCCTT - ACTTGTCCGAAACGAT

- ACTTGTCTGTTACGAT

	Sheep	Okapi	Giraffe
Cow	2	5	6
Sheep Okapi		5	6
Okapi			3
	•	1	1
	_		

#### For each pair of taxa

- Use the pairwise alignment induced by the MSA.
- · Count substitutions to obtain pairwise distances
- · Correct for multiple substitutions

#### Correcting for multiple substitutions with Jukes-Cantor

Given an alignment of n nucleotides that differs at m positions, the expected number of substitutions since the divergence of the two sequences is given by

*n* nucleotides with  $\mu$  mismatches

$$D = \frac{-3}{4} \ln \left( 1 - \frac{4}{3} \frac{\mu}{n} \right) n$$

From last lecture

For example, if we observe 200 mismatches in an alignment of 1000 nucleotides, then the number of actual substitutions is

 $\frac{-3}{4}\ln\left(1 - \frac{800}{3000}\frac{\mu}{n}\right)3000 = 233$ 

#### **Distance Takeaways**

- Several algorithms based on various assumptions about the data.
- If the data satisfy the assumptions, then distancebased algorithms will return the correct tree.
- If the data do not deviate too far from the assumptions, these methods give reasonable approximations
- They run in polynomial time.

#### Outline

Distance-based phylogeny reconstruction

Convert multiple alignments into distances

Properties of pairwise distances between taxa

- Additive distances
- Ultrametric distances

Using pairwise distances to infer a tree: constructive methods

- Neighbor Joining
- UPGMA