

For written notes on this lecture, please read chapter 11 of *The Practical Bioinformatician*

CS2220 Introduction to Computational Biology

Lecture 9: Phylogenetic Trees

Limsoon Wong
24 March 2006



Evolution

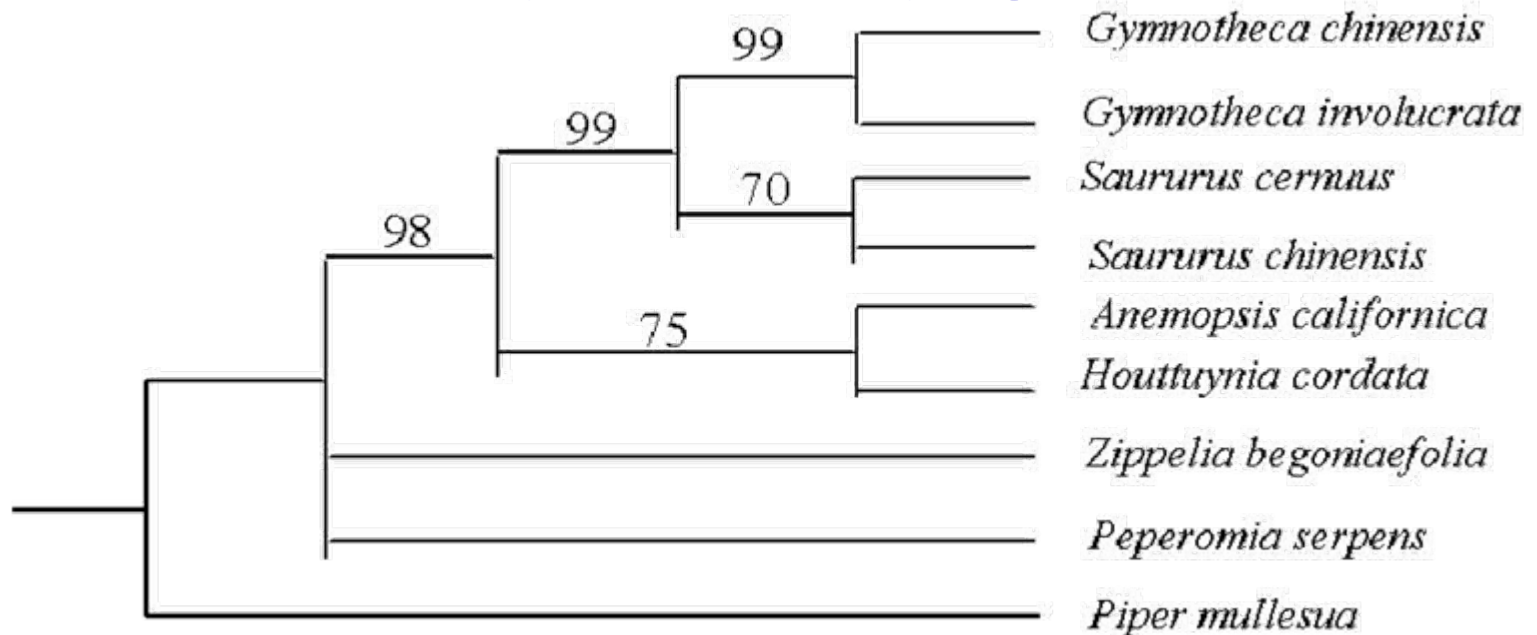
- Recall that DNA encodes blue print of life
- Living things pass DNA info to their children
- Due to mutations, DNA is changed a little bit
- After a long time, different species would evolve
- Phylogenetics studies genetic relationship between different species

Definition of Phylogeny

- **Phylogeny: Reconstruction of evolutionary history of a set of species**
- **Usually, it is a leaf-labeled tree where the internal nodes refer the hypothetical ancestors and the leaves are labeled by the species**
- **The edges of the tree represent the evolutionary relationships**

Phylogeny: An Example

- By looking at extent of conserved positions in the multiple seq alignment of different groups of seqs, can infer when they last shared an ancestor
 ⇒ Construct “family tree” or phylogeny



Application of Phylogeny

- **Understanding history of life**
- **Understanding rapidly mutating viruses (like HIV)**
- **Help to predict protein/RNA structure**
- **Help to do multiple sequence alignment**
- **Explaining and predicting gene expression**
- **Explaining and predicting ligands**
- **Help to design enhanced organisms**
- **Help to design drug**

Caution

- **Genomes of most organisms have complex origin**
 - Some parts of the genome are passed by vertical descent thru normal reproductive cycle
 - Some parts may have arisen by horizontal xfer of genetic material thru a virus, symbiosis, etc.
- ⇒ **When a particular gene is being subjected to phylogenetic analysis, the evolutionary history of that gene may not coincide with the evolutionary history of another gene**
- ⇒ **Try to use molecules that carry a great deal of evolutionary history, like mitochondrial DNA, and ribosomal RNA**

Phylogeny Reconstruction

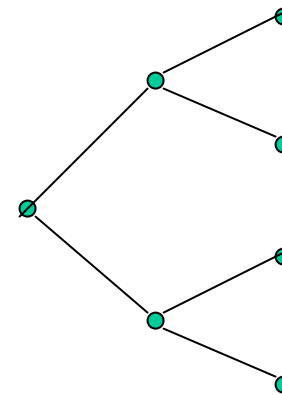


Rooted and Unrooted Tree

- Normally, the reconstructed tree is unrooted since estimating the root is scientifically difficult



- Rooted tree can be reconstructed by systematic biologists based on using outgroup
 - Outgroup is a species which is clearly less related with all other species in the phylogeny



Excercise: Why is a phylogenetic tree a binary tree?

Choosing Outgroup

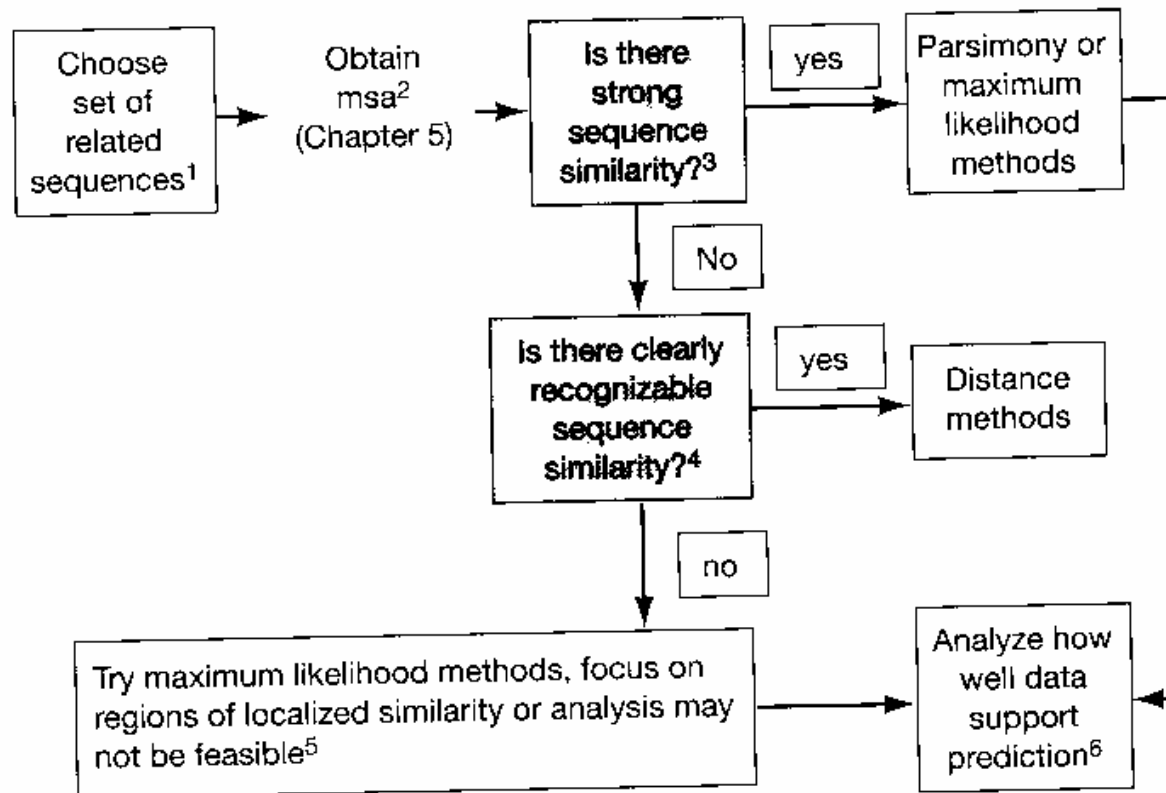
- **Outgroup seq should be closely related to rest of seqs, but there should also be significantly more diff betw outgroup and rest of seqs**
- **Outgroup that is too distant may lead to incorrect tree because of more random & complex nature of diff betw outgroup and rest of seqs**
- **In choosing outgroup, one assumes that the evolutionary history of the gene is same as rest of seqs. If this assumption is incorrect (e.g., horizontal gene xfer has occurred), an incorrect analysis could result**

Methods for Phylogenetic Reconstruction

- **Maximum parsimony**
- **Distance**
 - Straightforward
 - Applicable to large number of seqs
 - ⇒ Commonly used in mol biol labs
 - ⇒ We consider only this one here!
- **Maximum likelihood**
 - Require more understanding of evolutionary models on which they are based
 - Involve exponential number of steps
 - ⇒ Limited to small number of seqs

Exercise: What are the characteristics of max parsimony?

When to Use Which Phylogenetic Prediction Method?



Source: D.W.Mount, Bioinformatics: Sequence and Genome Analysis, Cold Spring Harbor Press, 2004

Distance Between Species

- In character-based methods, we try to minimize the number of mutations
- Species which look similar should be evolutionary more related
- ⇒ Define distance betw two species to be number of mutations need to change one species to another
- Try to construct a phylogeny based on distance info among species

Finding Distance Betw Two Species

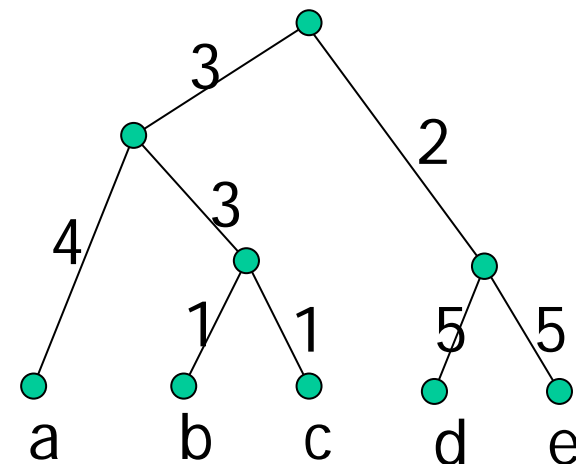
- **Consider two species with these DNA fragments:**
 - Species i: (A, C, G, C, T)
 - Species j: (C, C, A, C, T)
- **2 mismatches, so can estimate distance to be 2**
- **Looks reasonable, as 2 mismatches can be thought as 2 mutations**

- **However, this fails to capture “multiple” mutations on the same site**
- **In practice, need to apply some corrective distance transformation**

Distance Based

- **Input: Distance matrix M satisfying constraints**
 - M should satisfy the metric space
 - M is an additive metric
 - M is ultrametric (optional)
- **Output: Tree of degree 3 that is consistent with M**

M	a	b	c	d	e
a	0	8	8	14	14
b	8	0	2	14	14
c	8	2	0	14	14
d	14	14	14	0	10
e	14	14	14	10	0



Metric Space

- **A distance metric M which satisfies**

- Symmetry

$$M_{ij} = M_{ji} \geq 0$$

- Self identity

$$M_{ii} = 0$$

- Triangular inequality

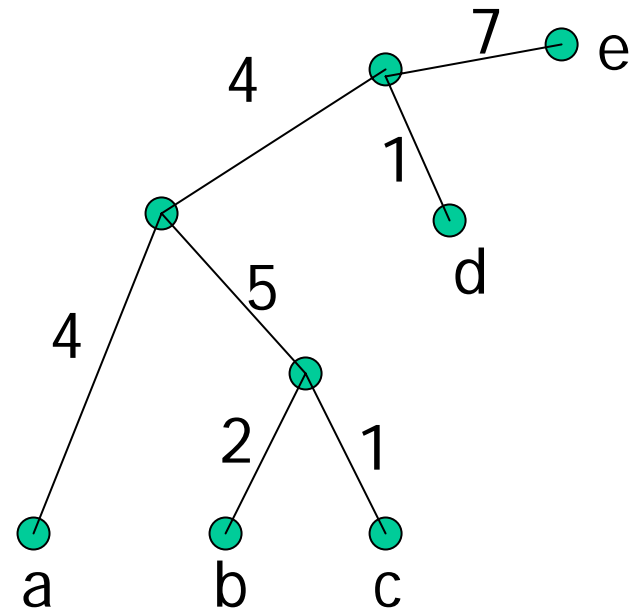
$$M_{ij} + M_{jk} \geq M_{ik}$$

Additive Metric

- **Let S be a set of species**
- **Let M be distance matrix for S**
- **If there is a rooted tree T where**
 - every edge has a positive weight and every leaf is labeled by a distinct species in S ; and
 - for every $i, j \in S$, M_{ij} = the sum of the edge weights along the path from i to j
- **Then M is called an additive metric**
- **The corresponding tree T is called additive tree**

Additive Metric Example

	a	b	c	d	e
a	0	11	10	9	15
b	11	0	3	12	18
c	10	3	0	11	17
d	9	12	11	0	8
e	15	18	17	8	0



- **Don't know the root! We can only build an unrooted phylogeny**

Why Additive Metric?

- **Distance captures actual number of mutations between a pair of species**
- **If (1) the correct tree for a set of species is known and (2) we get the exact number of mutations for each edge,**
 - The distance (the number of mutations) betw two species i and j should be the sum of the edge weights along the path from i to j

⇒ **Additive metric seems reasonable**

Properties of Additive Metric

- **Buneman's 4-point condition**

**M is additive if and only if
for any four species in S,
we can label them i, j, k, l such that**

$$M_{ik} + M_{jl} = M_{il} + M_{jk} \geq M_{ij} + M_{kl}$$

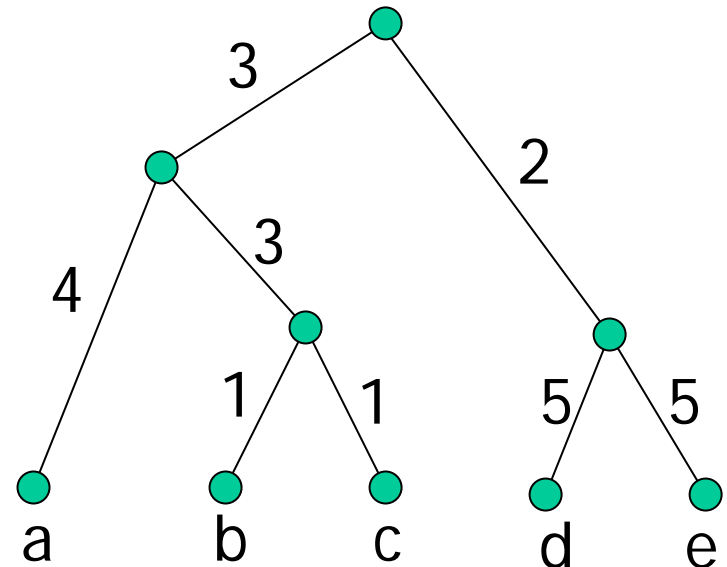
- **Based on the 4-point condition, we can check whether a matrix M is additive or not**

Ultrametric

- **Assume M is additive. That is, there exists a tree T such that**
 - the distance between any two species i and j equals the sum of the edge weights along the path from i to j .
- **If we can further identify a root such that the path length from the root of T to every leaf is identical, then M is called an ultrametric**
- **A tree T that satisfies ultrametric is an ultrametric tree**

Ultrametric Example

	a	b	c	d	e
a	0	8	8	14	14
b	8	0	2	14	14
c	8	2	0	14	14
d	14	14	14	0	10
e	14	14	14	10	0



- **Every path from root to leaf has the same length!**

Properties of Ultrametric

- **Ultrametric is an additive metric**
⇒ **It satisfies 4-point condition**
- **Additional property: 3-point condition**
M is ultrametric if and only if
for any three species in S,
we can label them i, j, k such that

$$M_{ik} = M_{jk} \geq M_{ij}$$

- **Based on the 3-point condition, we can check whether a matrix M is ultrametric or not**

Constant Molecular Clock

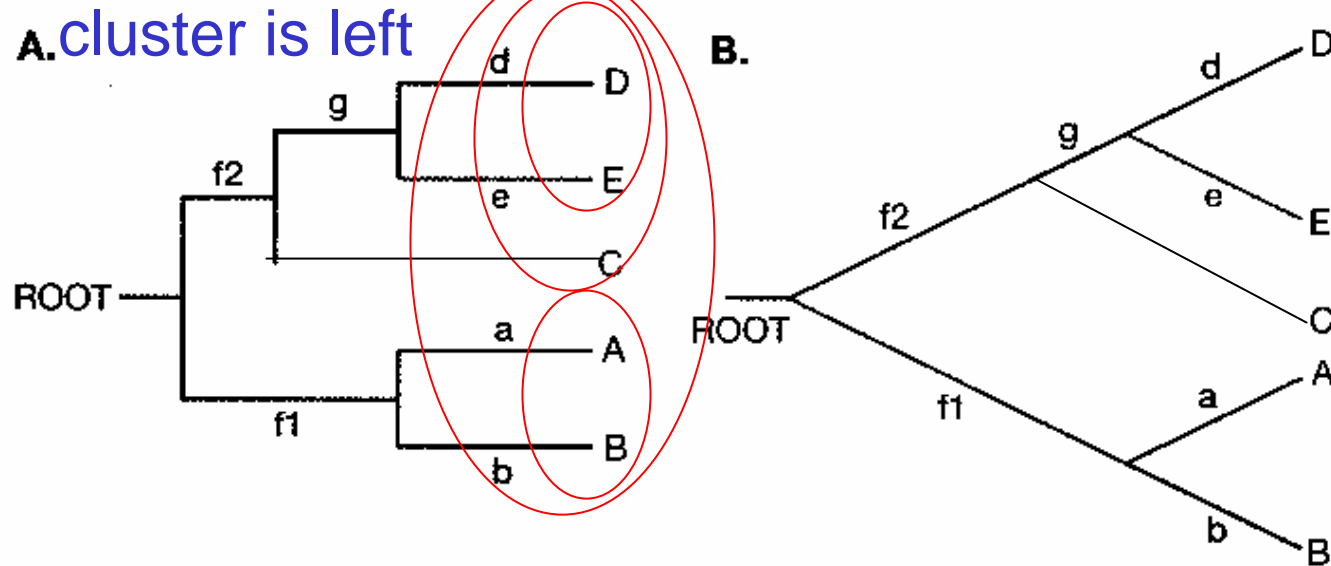
- **Constant molecular clock is an assumption in biology**
 - It states that the number of accepted mutations occurring in any time interval is proportional to the length of that interval
 - ⇒ All species evolved at equal rate from a common ancestor
- **Ultrametric tree states that distance from root to all species are the same. Thus, its correctness is based the constant molecular clock assumption, which is rarely correct!**

Some Computational Problems

- **Let M be a distance matrix for a set of species S**
 - If M is ultrametric, can we reconstruct the corresponding ultrametric tree T in polynomial time? (only consider this one!)
 - If M is additive, can we have an polynomial time algorithm to recover the corresponding additive tree T ?
 - If M is not exactly additive, can we find the nearest additive tree T ?

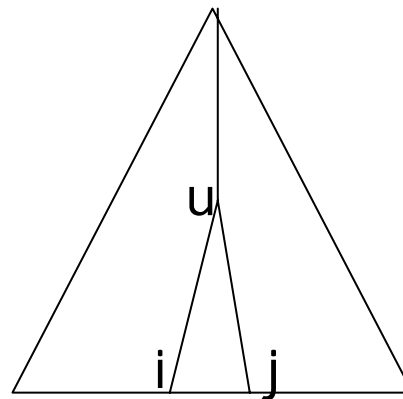
Unweighted Pair Group Method With Arithmetic Mean (UPGMA)

- Consider ultrametric tree T. If a subset of species S forms a subtree of T, we call it a cluster
- Idea:
 - Every species forms a cluster
 - Iteratively connect two nearest clusters, until one cluster is left



Definition - Height

- For a node u , define $\text{height}(u)$ be path length from u to any of its descendent leaf. (Since T is ultrametric, every path should have the same length!)
- Let i and j be descendent leaves of u in two different subtrees. To ensure that distance from the root to both i and j are the same, $\text{height}(u) = M_{ij}/2$



Distance Betw Two Clusters

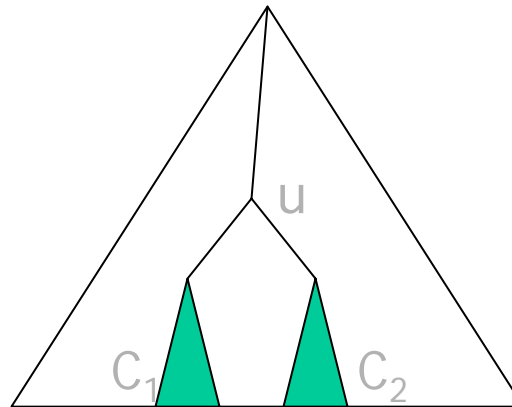
- For any two clusters C_1 and C_2 of T

- Define

$$\text{dist}(C_1, C_2) = \frac{\sum_{i \in C_1, j \in C_2} M_{ij}}{|C_1| \cdot |C_2|}$$

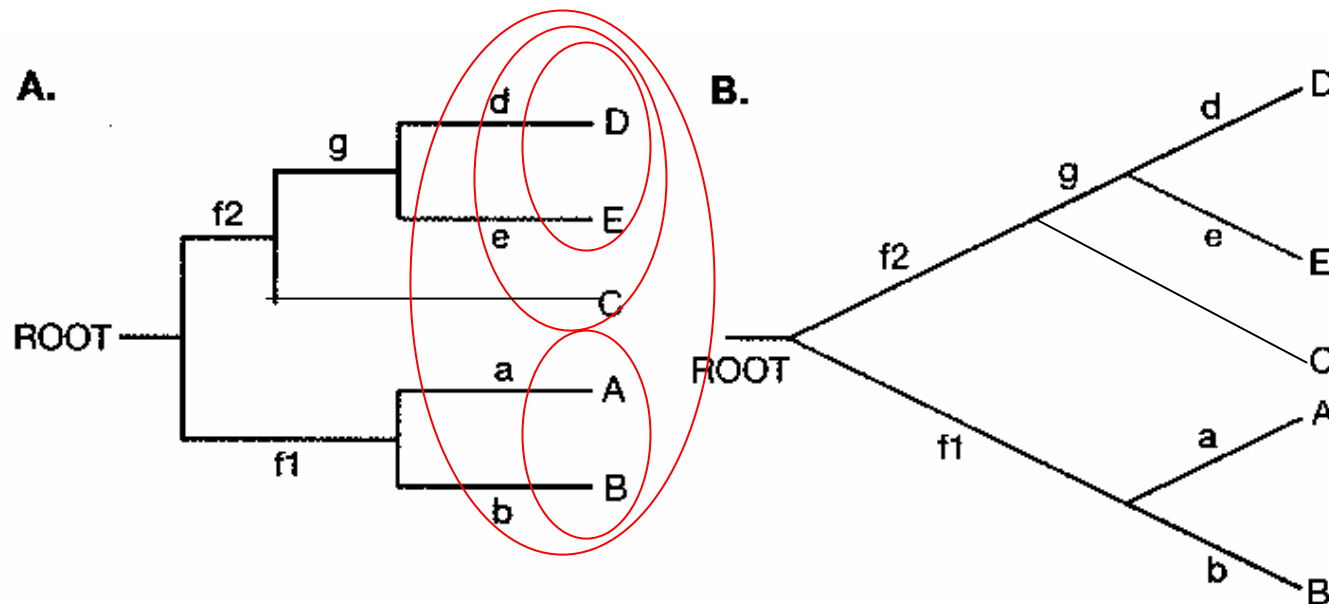
- Note that $\text{dist}(C_1, C_2) = M_{ij}$ for all $i \in C_1$ and $j \in C_2$ Why?

- Let u be lowest common ancestor of i and j .
 $\text{dist}(C_1, C_2) = 2 * \text{height}(u)$!



Idea of the UPGMA Algorithm

- Consider a set Z of clusters
- Let A, B be two clusters st $\text{dist}(A, B)$ is min
- Let C be tree formed by joining A and B w/ a root
- Repeat this until no more clusters to merge

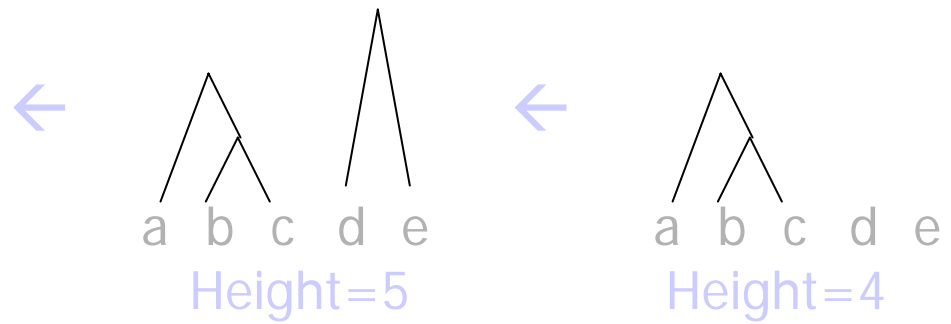


Algorithm

- Given $n \times n$ ultrametric distance matrix M
- Initialize set Z to consist of n initial singleton clusters $\{1\}, \{2\}, \dots, \{n\}$
- For all $\{i\}, \{j\} \in Z$, initialize $\text{dist}(\{i\}, \{j\}) = M_{ij}$
- Repeat $n-1$ times
 - Determine cluster $A, B \in Z$ where $\text{dist}(A, B)$ is min
 - Define a new cluster $C = A \cup B$
 - $Z := Z - \{A, B\} \cup \{C\}$
 - Define new node c and let c be parent of a and b . Also, define $\text{height}(c) = \text{dist}(A, B)/2$
 - For all $D \in Z - \{C\}$, define $\text{dist}(D, C) = \text{dist}(C, D) = (\text{dist}(A, D) + \text{dist}(B, D)) / 2$

Example

M	a	b	c	d	e
a	0	8	8	14	14
b	8	0	2	14	14
c	8	2	0	14	14
d	14	14	14	0	10
e	14	14	14	10	0



Time Complexity

- Initialization can be done in $O(n^2)$ time
- There are $n-1$ iterations, each iteration takes $O(n)$ time
- The total time complexity is $O(n^2)$

Phylogenetic Tree Comparison



Why Tree Comparison?

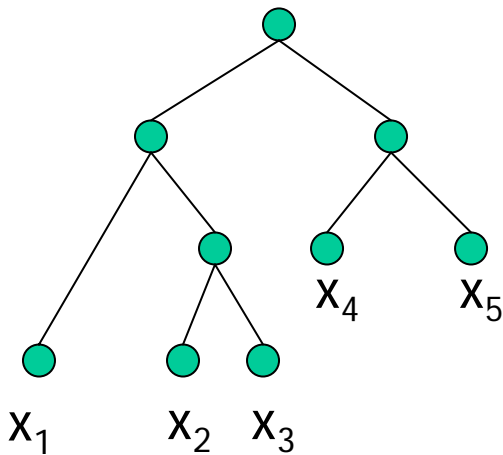
- **We learn a number of methods to reconstruct phylogeny for the same set of species**
- **Different phylogenies are resulted using**
 - Different data (different segments of genomes)
 - Different model (CF model, Jukes-Cantor Model)
 - Different reconstruction algorithms
- **Tree comparison helps us to gain information from multiple trees**

Two Types of Comparisons

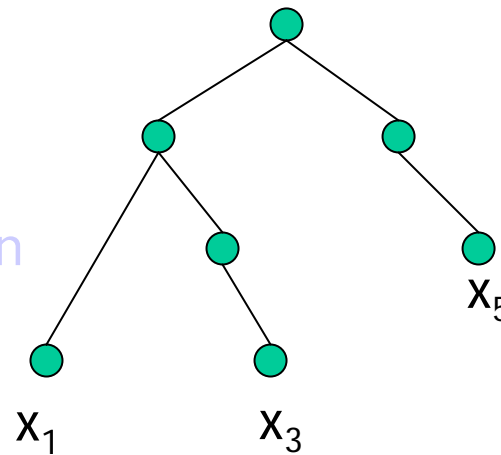
- **Similarity measurement**
 - Find common structure among given trees
 - **Maximum Agreement Subtree**
- **Dissimilarity measurement**
 - Determine differences among given trees
 - **Robinson-Foulds distance**
 - **Nearest-neighbor interchange**
 - **Subtree transfer distance**
- **In this lecture, we will discuss the first method**

Restricted Subtree

- Consider tree T



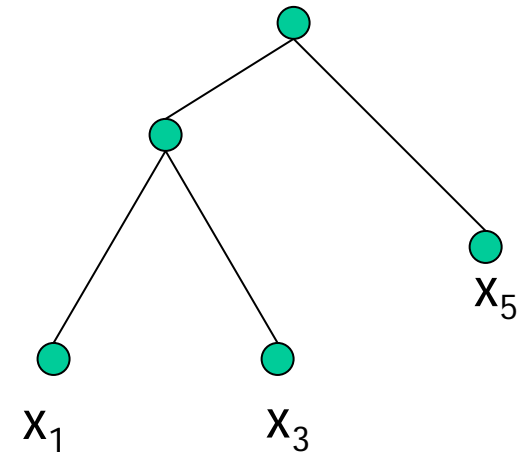
→
Restricted on
 X_1, X_3, X_5



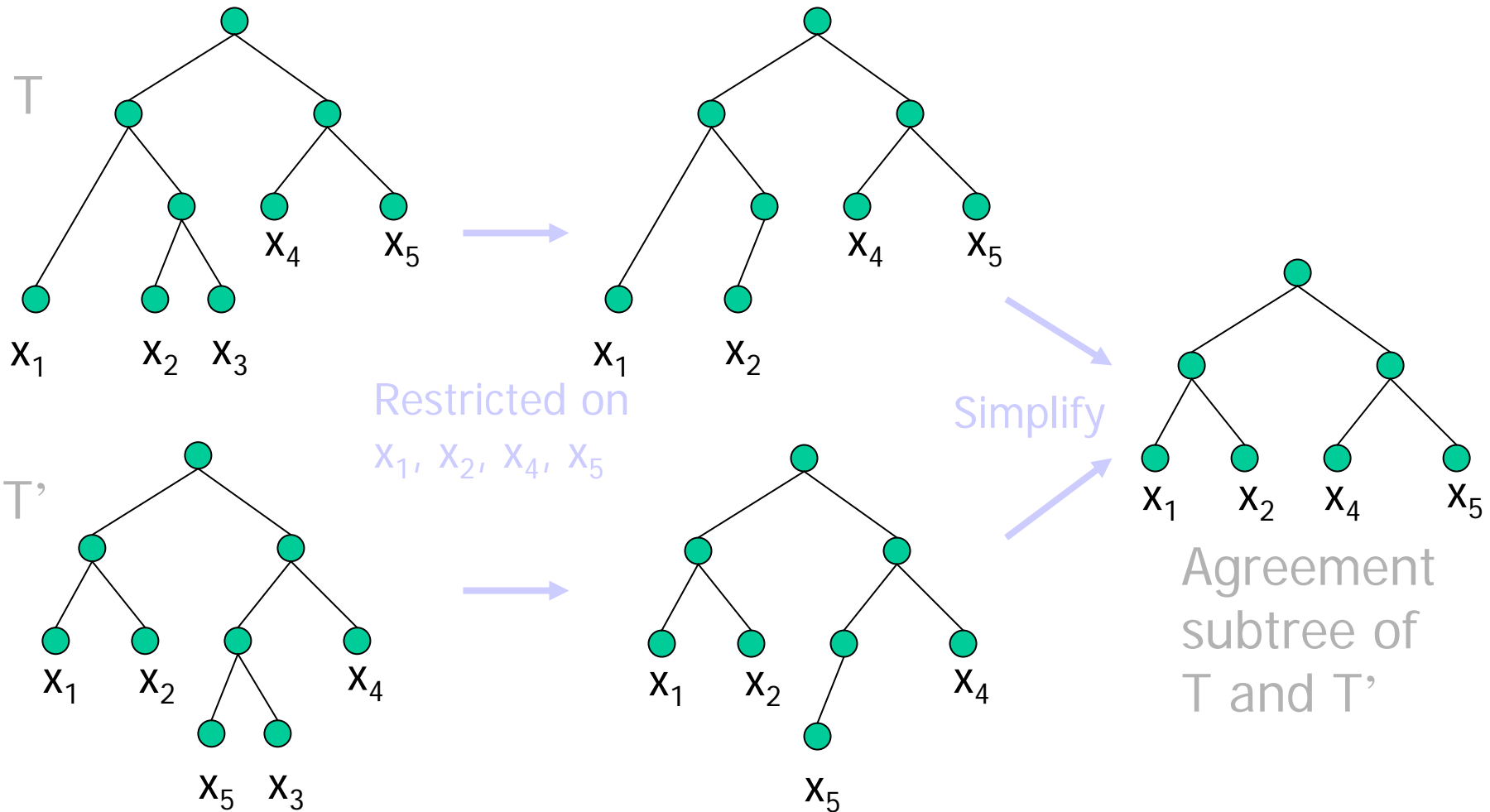
Evolution
information
of X_1, X_3, X_5

Evolution
information of $X_1,$
 X_2, X_3, X_4, X_5

↓
Simplify



Agreement Subtree



Maximum Agreement Subtree (MAST)

- **Given two trees T_1 and T_2**
- **Agreement subtree of T_1 and T_2 is the common info agreed by both trees**
 - Since it is agreed by both trees, the evolution of the agreement subtree is more reliable
- **Maximum agreement subtree problem**
 - Find the agreement subtree with largest possible number of leaves
 - Such agreement subtree is called the maximum agreement subtree

MAST for Rooted Trees

- **MAST of two degree- d rooted trees T_1 and T_2 with n leaves can be computed in**

$$O(\sqrt{d}n \log(\frac{n}{d})) \text{ time}$$

- **But the algo for the above is complicated**
- **So here we show you a $O(n^2)$ -time algorithm which computes the maximum agreement subtree of two binary trees with n leaves**

MAST by Dynamic Programming

Notations

- For any two binary rooted trees T_1 and T_2 , let $\text{MAST}(T_1, T_2)$ be number of leaves in the maximum agreement subtree
- For a tree T and a node u , T^u is the subtree of T rooted at u

Base Cases

- For any leaf x in T_1 and y in T_2 ,

$$MAST(x, y) = \max \begin{cases} 1 & \text{if } x = y \\ 0 & \text{otherwise} \end{cases}$$

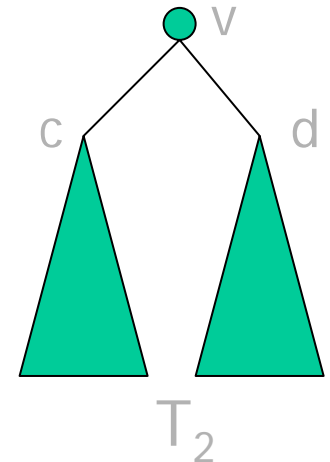
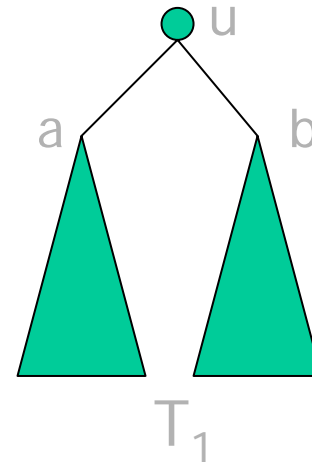
- For any node u in T_1 and v in T_2 ,

$$MAST(T_1^u, \Lambda) = 0, MAST(\Lambda, T_2^v) = 0$$

Recurrence (I)

$$MAST(T_1^u, T_2^v) =$$

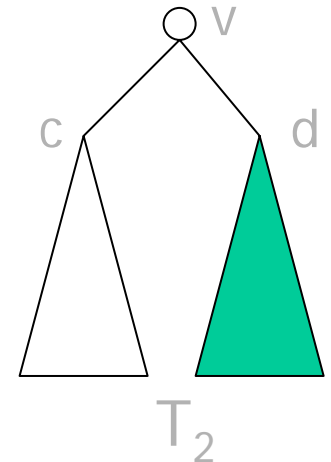
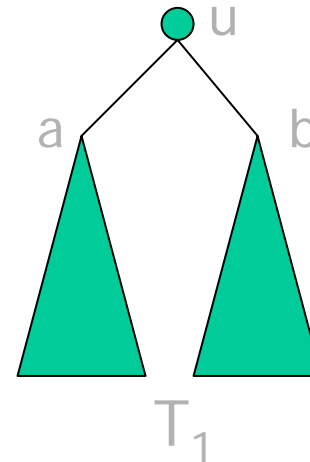
$$\max \left\{ \begin{array}{l} MAST(T_1^a, T_2^c) + MAST(T_1^b, T_2^d) \\ MAST(T_1^a, T_2^d) + MAST(T_1^b, T_2^c) \\ MAST(T_1^a, T_2^v) \\ MAST(T_1^b, T_2^v) \\ MAST(T_1^u, T_2^c) \\ MAST(T_1^u, T_2^d) \end{array} \right.$$



Recurrence (II)

$$MAST(T_1^u, T_2^v) =$$

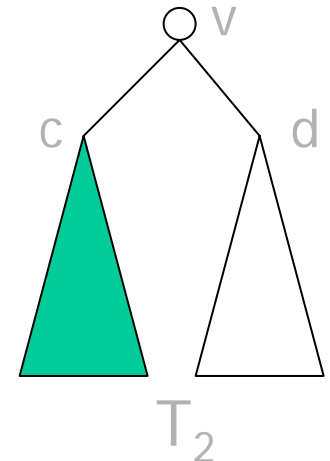
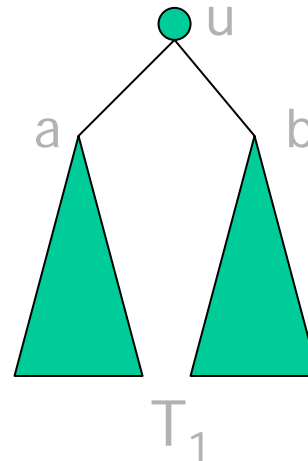
$$\max \left\{ \begin{array}{l} MAST(T_1^a, T_2^c) + MAST(T_1^b, T_2^d) \\ MAST(T_1^a, T_2^d) + MAST(T_1^b, T_2^c) \\ MAST(T_1^a, T_2^v) \\ MAST(T_1^b, T_2^v) \\ MAST(T_1^u, T_2^c) \\ MAST(T_1^u, T_2^d) \end{array} \right.$$



Recurrence (III)

$$MAST(T_1^u, T_2^v) =$$

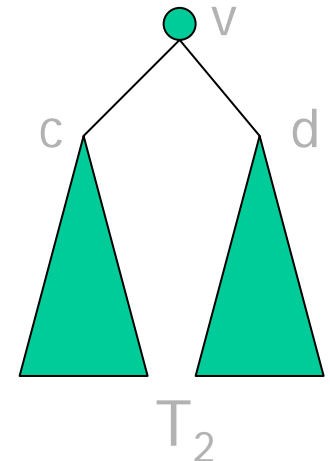
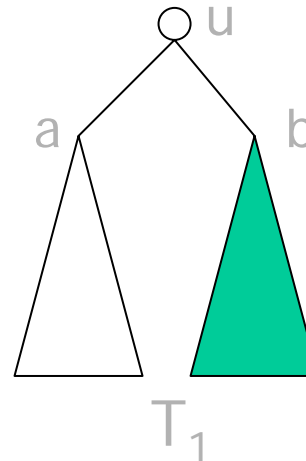
$$\max \left\{ \begin{array}{l} MAST(T_1^a, T_2^c) + MAST(T_1^b, T_2^d) \\ MAST(T_1^a, T_2^d) + MAST(T_1^b, T_2^c) \\ MAST(T_1^a, T_2^v) \\ MAST(T_1^b, T_2^v) \\ MAST(T_1^u, T_2^c) \\ MAST(T_1^u, T_2^d) \end{array} \right.$$



Recurrence (IV)

$$MAST(T_1^u, T_2^v) =$$

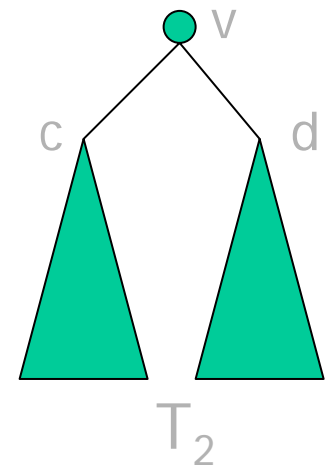
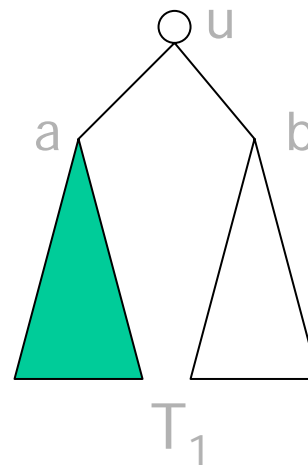
$$\max \left\{ \begin{array}{l} MAST(T_1^a, T_2^c) + MAST(T_1^b, T_2^d) \\ MAST(T_1^a, T_2^d) + MAST(T_1^b, T_2^c) \\ MAST(T_1^a, T_2^v) \\ MAST(T_1^b, T_2^v) \\ MAST(T_1^u, T_2^c) \\ MAST(T_1^u, T_2^d) \end{array} \right.$$



Recurrence (V)

$$MAST(T_1^u, T_2^v) =$$

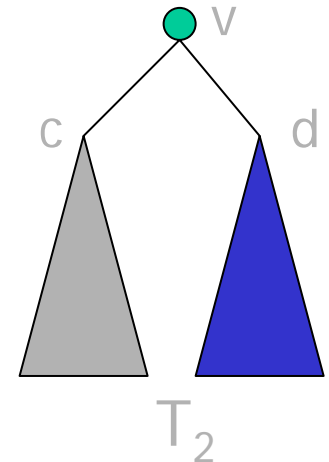
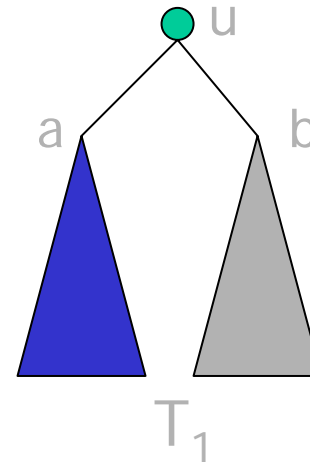
$$\max \left\{ \begin{array}{l} MAST(T_1^a, T_2^c) + MAST(T_1^b, T_2^d) \\ MAST(T_1^a, T_2^d) + MAST(T_1^b, T_2^c) \\ MAST(T_1^a, T_2^v) \leftarrow \\ MAST(T_1^b, T_2^v) \\ MAST(T_1^u, T_2^c) \\ MAST(T_1^u, T_2^d) \end{array} \right.$$



Recurrence (VI)

$$MAST(T_1^u, T_2^v) =$$

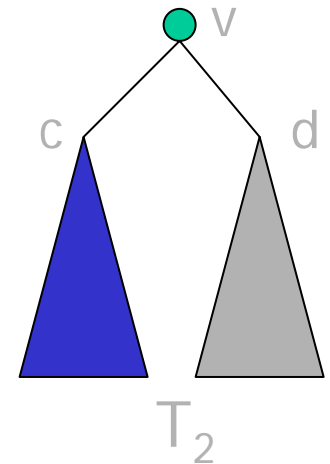
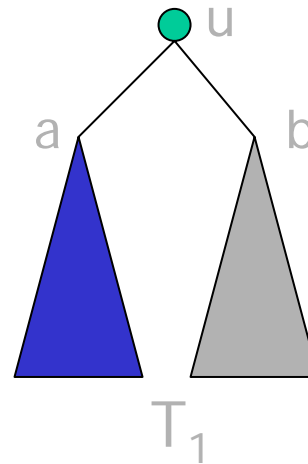
$$\max \left\{ \begin{array}{l} MAST(T_1^a, T_2^c) + MAST(T_1^b, T_2^d) \\ MAST(T_1^a, T_2^d) + MAST(T_1^b, T_2^c) \leftarrow \\ MAST(T_1^a, T_2^v) \\ MAST(T_1^b, T_2^v) \\ MAST(T_1^u, T_2^c) \\ MAST(T_1^u, T_2^d) \end{array} \right.$$



Recurrence (VII)

$$MAST(T_1^u, T_2^v) =$$

$$\max \left\{ \begin{array}{l} MAST(T_1^a, T_2^c) + MAST(T_1^b, T_2^d) \leftarrow \\ MAST(T_1^a, T_2^d) + MAST(T_1^b, T_2^c) \\ MAST(T_1^a, T_2^v) \\ MAST(T_1^b, T_2^v) \\ MAST(T_1^u, T_2^c) \\ MAST(T_1^u, T_2^d) \end{array} \right.$$



Time Complexity

- Suppose T_1 and T_2 are rooted phylogenies for n species
- We have to compute $\text{MAST}(T_1^u, T_2^v)$ for every u in T_1 and v in T_2
- Thus, we need to fill in n^2 entries
- Each entry can be computed in $O(1)$ time
- In total, the time complexity is $O(n^2)$

SNP: From Looking for Similarities To Looking for Differences



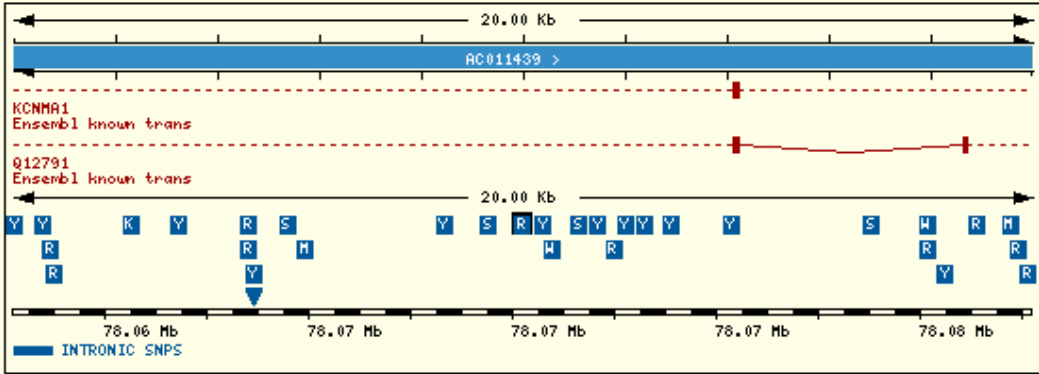
Single Nucleotide Polymorphism

- SNP occurs when a single nucleotide replaces one of the other three nucleotide letters
- E.g., the alteration of the DNA segment
 AAGGTTA to
 ATGGTTA
- SNPs occur in human population > 1% of the time
- Most SNPs are found outside of "coding seqs"
- SNPs found in a coding seq are of great interest as they are more likely to alter function of a protein

Exercise: Why are most SNPs found outside of coding seqs?

Example SNP Report

Ensembl SNP Report

SNP	1907745
Source	dbSNP
Synonyms	dbSNP: 1907745 TSC: TSC0953388 HGbase: SNP001275703
Score	1
Validation Status	proven by cluster (SNP tested and validated by a non-computational method)
Alleles	A G (ambiguity code: R)
Sequence Region	AGGCATCCAGTCTCGGTAAACCTAG R CAAGTAATATTATTAGTTGAGCATT (SNP highlighted)
SNP neighbourhood	 <p>The figure is a genomic map of a 20.00 Kb region. At the top, a scale bar shows the 20.00 Kb distance. Below it, a blue bar represents the AC011439 contig. Two red dashed lines indicate the Ensembl known transcripts for KCNMA1 and Q12791. The bottom part shows a sequence alignment with various amino acid residues (Y, K, R, S, H, W, R, Y, S, Y, Y, Y, Y, S, H, R, H, R) and a legend for INTRONIC SNPS.</p>

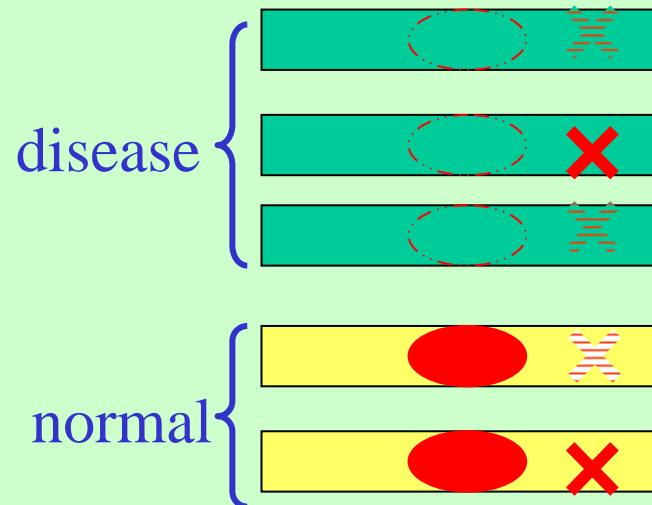
SNP Uses

- **Association studies**
 - Analyze DNA of group affected by disease for their SNP patterns
 - Compare to patterns obtained from group unaffected by disease
 - Detect diff betw SNP patterns of the two
 - Find pattern most likely associated with disease-causing gene

strong assoc ∞ weak assoc $\frac{1}{2}$

	Allele A	Allele B
Disease	π_{00}	π_{01}
Normal	π_{10}	π_{11}

odds ratio = $\frac{\pi_{00}}{\pi_{10}} / \frac{\pi_{01}}{\pi_{11}} = \frac{\pi_{00} * \pi_{11}}{\pi_{01} * \pi_{10}}$



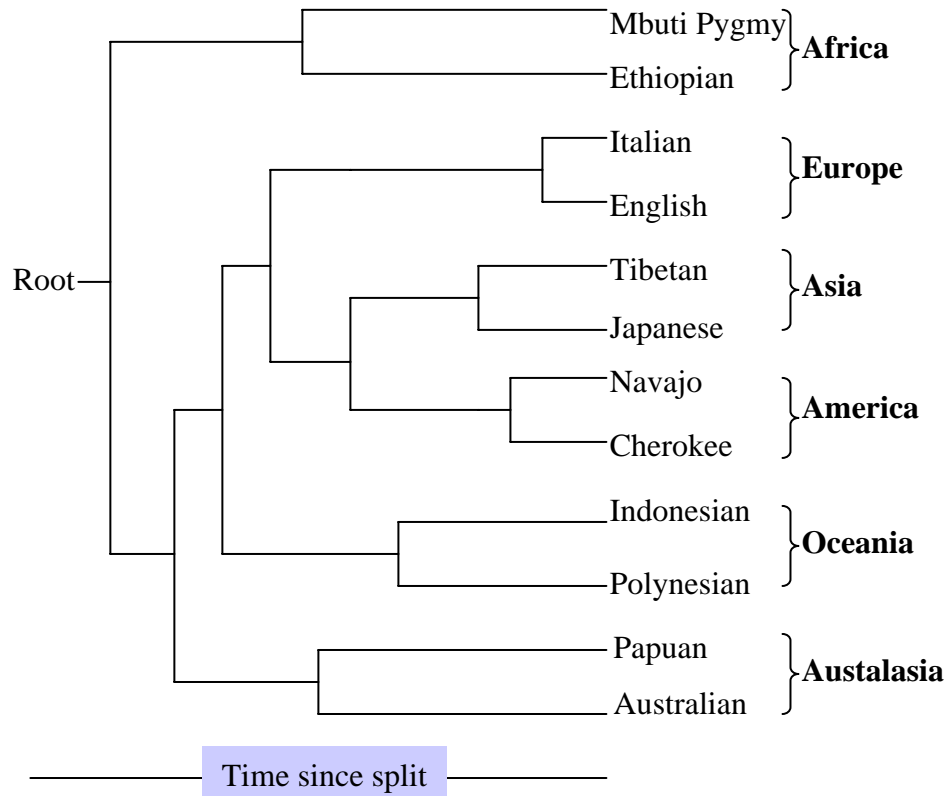
SNP Uses

- **Better evaluate role of non-genetic factors (e.g., behavior, diet, lifestyle)**
- **Determine why people differ in abilities to absorb or clear a drug**
- **Determine why an individual experiences side effect of a drug**

The 7 Daughters of Eve

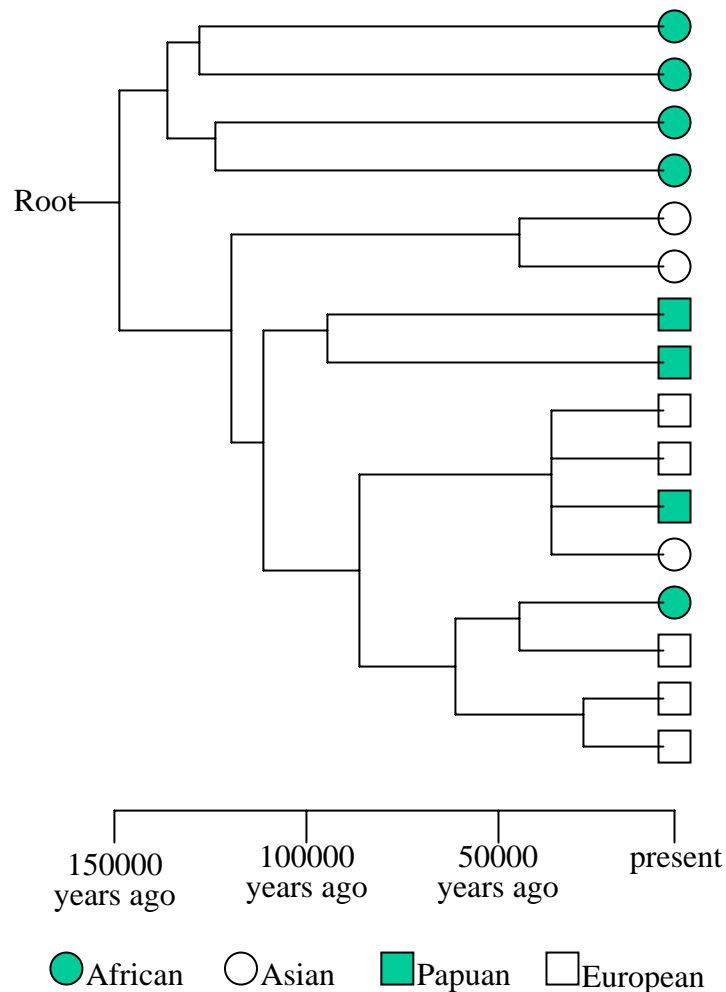


Population Tree



- Estimate order in which “populations” evolved
- Based on assimilated freq of many different genes
- But ...
 - is human evolution a succession of population fissions?
 - Is there such thing as a proto-Anglo-Italian population which split, never to meet again, and became inhabitants of England and Italy?

Evolution Tree



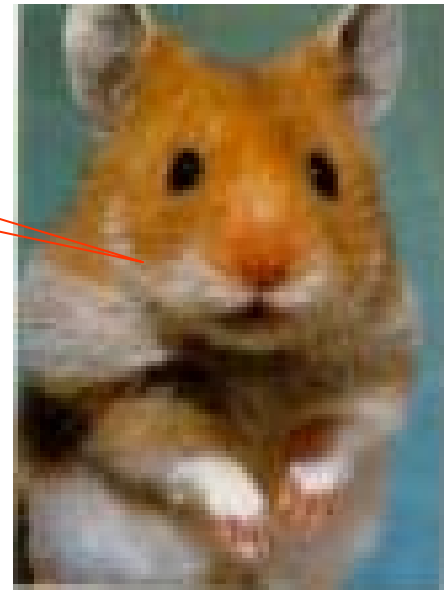
- Leaves and nodes are individual persons---real people, not hypothetical concept like “proto-population”
- Lines drawn to reflect genetic differences between them in one special gene called mitochondrial DNA

Why Mitochondrial DNA

- **Present in abundance in bone fossils**
- **Inherited only from mother**
- **Sufficient to look at the 500bp control region**
- **Accumulate more neutral mutations than nuclear DNA**
- **Accumulate mutations at the “right” rate, about 1 every 10,000 years**
- **No recombination, not shuffled at each generation**

Mutation Rates

- All pet golden hamsters in the world descend from a single female caught in 1930 in Syria
 - Golden hamsters “manage” ~4 generations a year :-)
 - So >250 hamster generations since 1930
 - Mitochondrial control regions of 35 (independent) golden hamsters were sequenced and compared
 - No mutation was found
- ⇒ Mitochondrial control region mutates at the “right” rate

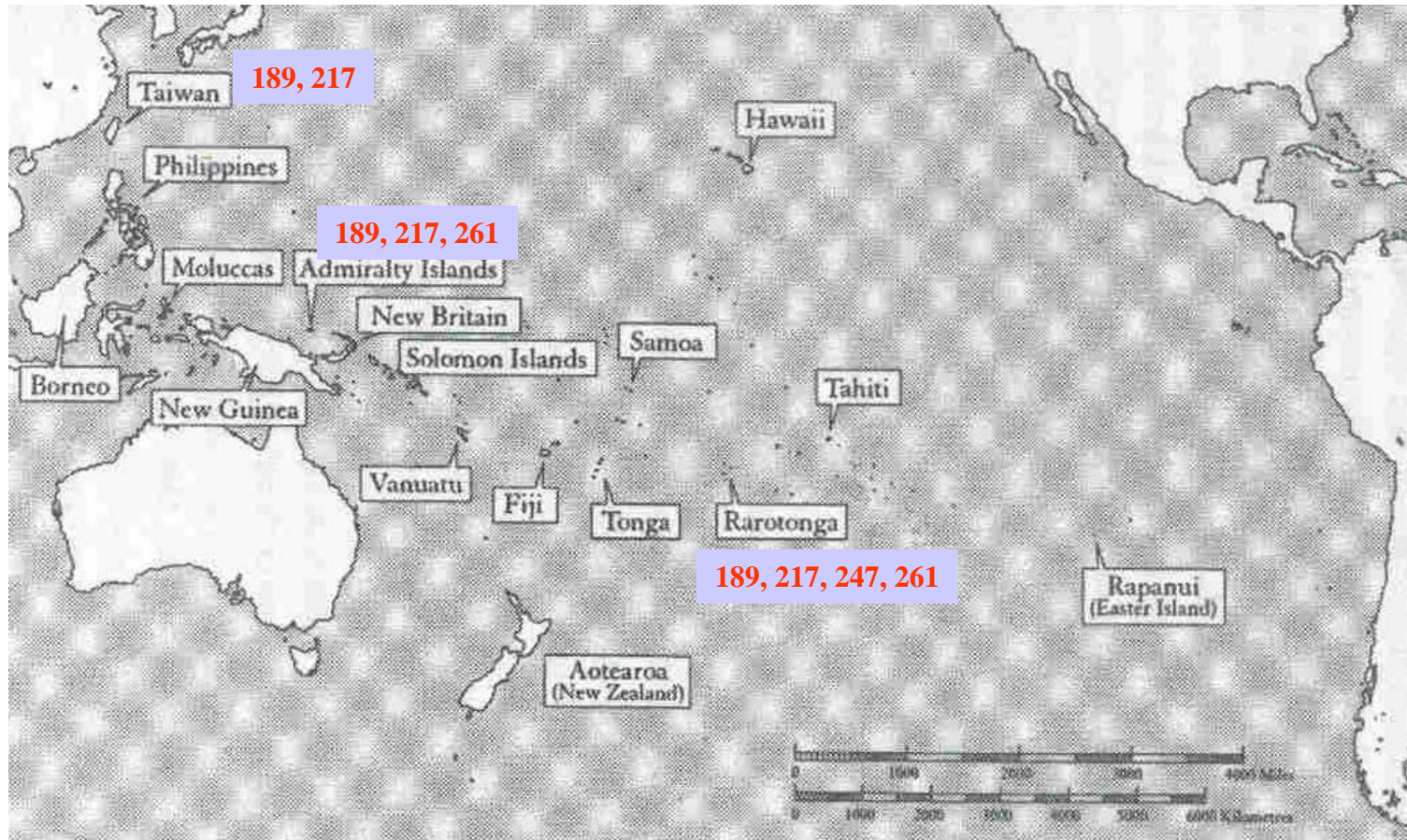


Contamination

- **Need to know if DNA extracted from old bones really from those bones, and not contaminated with modern human DNA**
- **Apply same procedure to old bones from animals, check if you see modern human DNA.**
- **If none, then procedure is OK**

Origin of Polynesians

- Do they come from Asia or America?



Origin of Polynesians

- **Common mitochondrial control seq from Rarotonga have variants at positions 189, 217, 247, 261. Less common ones have 189, 217, 261**
- **Seq from Taiwan natives have variants 189, 217**
- **Seq from regions in betw have variants 189, 217, 261.**
- **More 189, 217 closer to Taiwan. More 189, 217, 261 closer to Rarotonga**
- **247 not found in America**
⇒ **Polynesians came from Taiwan!**
- **Taiwan seq sometimes have extra mutations not found in other parts**
⇒ **These are mutations that happened since Polynesians left Taiwan!**

Neanderthal vs Cro Magnon

- Are Europeans descended purely from Cro Magnons? Pure Neanderthals? Or mixed?



Neanderthal



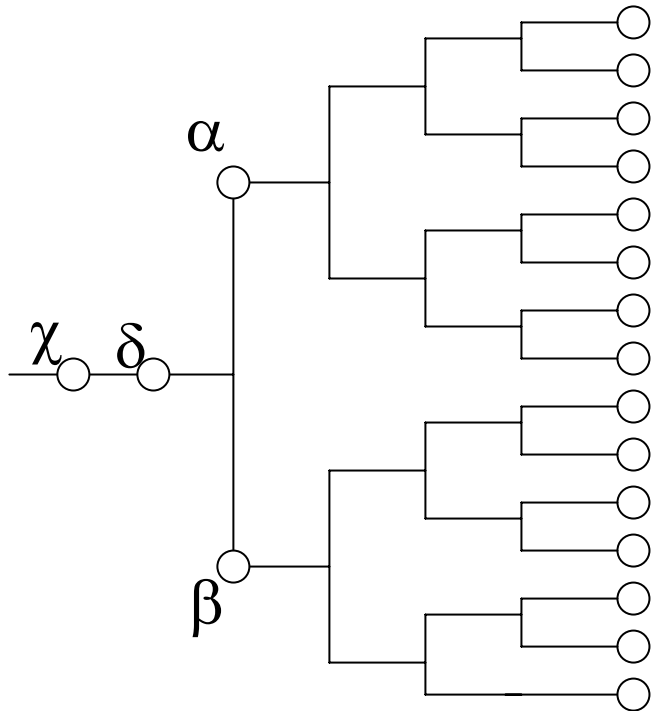
Cro Magnon



Neanderthal vs Cro Magnon

- Based on palaeontology, Neanderthal & Cro Magnon last shared an ancestor 250000 yrs ago
- Mitochondrial control regions accumulate 1 mutation per 10000 yrs
 - ⇒ If Europeans have mixed ancestry, the mitochondrial control regions betw 2 Europeans should have ~25 diff w/ high probability
- The number of diff betw Welsh is ~3, & at most 8.
- When compared w/ other Europeans, 14 diff at most
 - ⇒ Ancestor either 100% Neanderthal or 100% Cro Magnon
- Mitochondrial control seq from Neanderthal have 26 diff from Europeans
 - ⇒ Ancestor must be 100% Cro Magnon

Clan Mother



Exercise: Which of α , β , χ , δ is the clan mother?

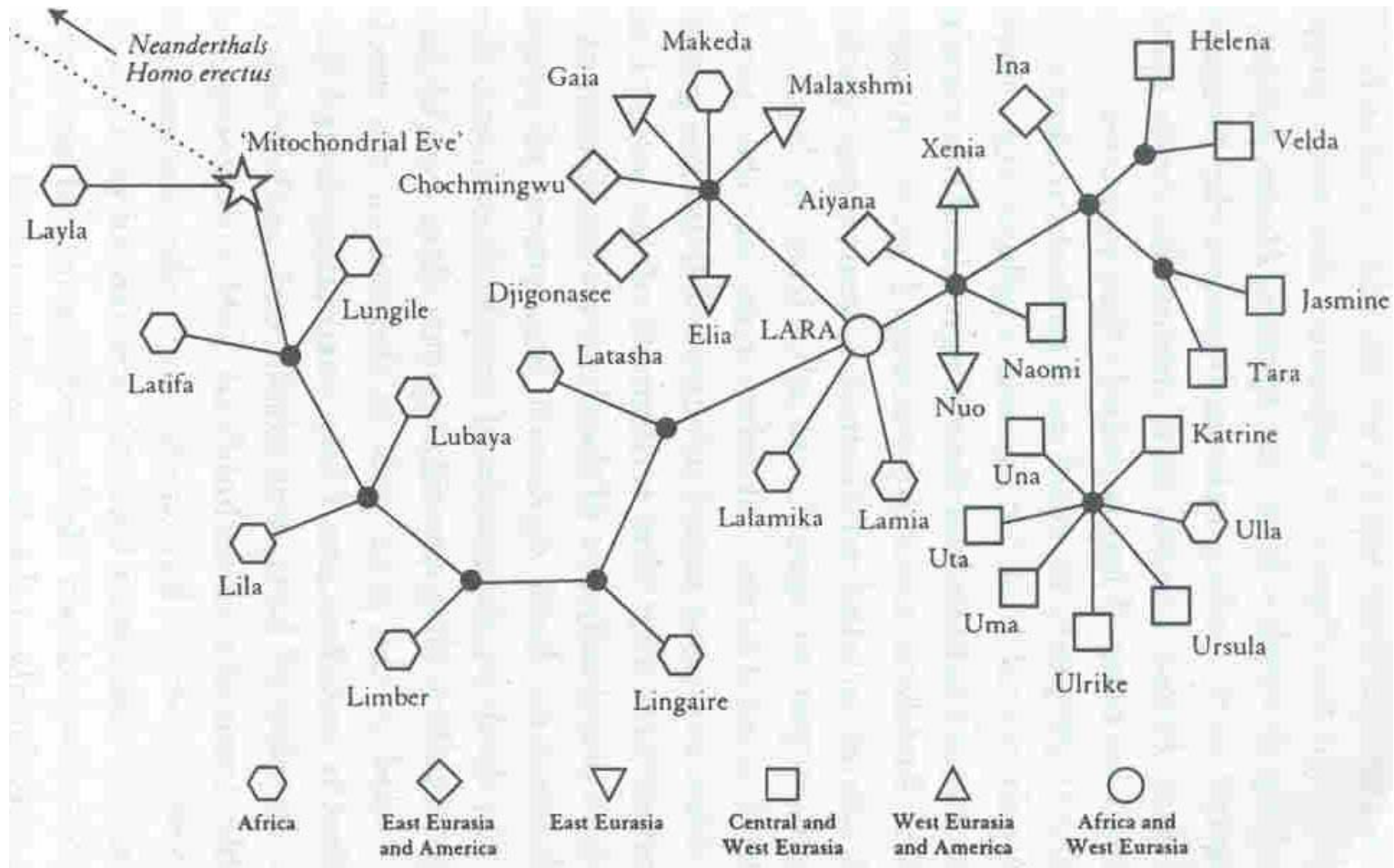
- **Clan mother is the most recent maternal ancestor common to all members of the clan**
- **A woman with only sons cant be clan mother---her mitochondrial DNA cant be passed on**
- **A woman cant be clan mother if she has only 1 daughter---she is not most recent maternal ancestor**

How many clans in Europe?

- Cluster seq according to mutations
- Each cluster thus represents a major clan
- European seq cluster into 7 major clans
- The 7 clusters age betw 45000 and 10000 years (length of time taken for all mutations in a cluster to arise from a single founder seq)
- The founder seq carried by just 1 woman in each case--the clan mother
- Note that the clan mother did not need to be alone. There could be other women, it was just that their descendants eventually died out

Exercise: How about clan father?

World Clans



Any Question?



Acknowledgements

- **A lot of the slides from this lecture are given to me by Ken Sung**

References

- B. Sykes. *The seven daughters of Eve*, Gorgi Books, 2002
- S.-W. Meng. Analysis of Phylogeny: A Case Study on Saururacea, *The Practical Bioinformatician*, chapter 11, pages 245—268, WSPC, 2004