

For written notes on this lecture, please read chapter 11 of *The Practical Bioinformatician*, Chapters 7 & 8 of *Algorithms in Bioinformatics: A Practical Introduction*, and Chapter 17 of *Algorithms on Strings, Trees, and Sequences*.

CS2220 Introduction to Computational Biology Lecture 8: Phylogenetic Trees

Limsoon Wong



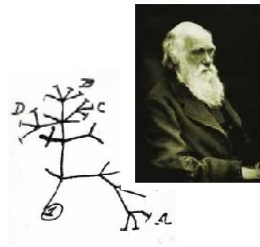
Evolution

- 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

Copyright 2011 © Limsoon Wong

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
- Edges of the tree represent the evolutionary relationships

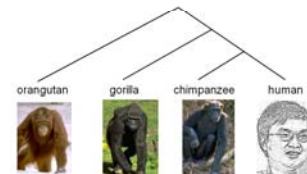


First Notebook on Transmutation of Species, 1837.

Copyright 2011 © Limsoon Wong

Phylogeny: An Example

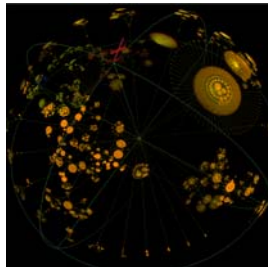
- By looking at extent of conserved positions in the “multiple seq alignment” of different groups of seqs, can infer when they last shared an ancestor
⇒ Construct “family tree” or phylogeny



Copyright 2011 © Limsoon Wong

Application of Phylogeny

- Understanding history of life
- Understanding rapidly mutating viruses (like HIV)
- Predict protein/RNA struct
- Do multiple seq alignment
- Explain and predict gene expression
- Explain and predict ligands
- Design enhanced organisms
- Design drug




Copyright 2011 © Limsoon Wong

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

Copyright 2011 © Limsoon Wong

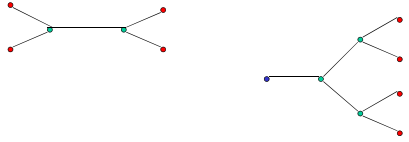
Phylogeny Reconstruction



Rooted and Unrooted Tree




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



Copyright 2011 © Limsoon Wong

How does outgroup work?



- More similar to outgroup ⇒ More "ancient"
- More diff from outgroup ⇒ More "recent", because more time to evolve

Species	1	2	3	4	Outgroup
Character state	a	a'	a	a'	a

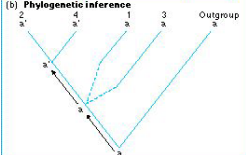



Image credit: Mark Ridley

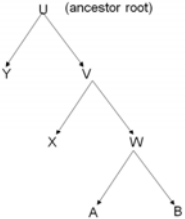
Copyright 2011 © Limsoon Wong

An Exercise



X	ACCTG-TACTTCGATAA
Y	ACCAG-TACTT-GATAA
A	ACCAGGTACTTCGATAT
B	ACCAGGTACTTCGATTT

1 2 3 4




- What is the most likely sequence for U?
- Hint: A phylogeny with fewer mutations is more likely than a phylogeny with fewer mutations

U = ACCAG-TACTT[C or -]GATAA

If position 1 is "T", then both Y and W has a mutation in this position. If position 1 is "A", then only X has a mutation in this position. By the parsimony assumption, position 1 must be "A".

Copyright 2011 © Limsoon Wong


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

Copyright 2011 © Limsoon Wong

Methods for Phylogenetic Reconstruction



- Maximum parsimony**
 - Exercise: What are the characteristics of max 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

Copyright 2011 © Limsoon Wong

13

When to Use Which Phylogenetic Prediction Method?

```

    graph TD
      A[Choose set of related sequences1] --> B[Obtain "mas"2 (Chapter 5)]
      B --> C{Is there strong sequence similarity?3}
      C -- yes --> D[Parimony or maximum likelihood methods]
      C -- No --> E{Is there clearly recognizable sequence similarity?4}
      E -- yes --> F[Distance methods]
      E -- no --> G[Analyze how well data support prediction5]
      G --> H[Try maximum likelihood methods, focus on regions of localized similarity or analysis may not be feasible6]
      H --> D
  
```

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

Copyright 2011 © Limsoon Wong

14

Allan Wilson

- **"Molecular clock": Dating by genetic mutations**
 - Deduced in 60s that proto-hominids evolved 5m yrs ago, contrary to the 25m yrs believed by anthropologists
 - In 80s, his findings became more widely accepted
- **Molecular approach to understand evolution**
 - Concluded in 80s that modern man evolved from "African Eve"
 - 20 yrs to convince palaeontologists, but when they did, it married their science with that of genetics

Copyright 2011 © Limsoon Wong

15

Distance Between Species

- In character-based methods, we try to minimize # of mutations
- Species which look similar should be evolutionary more related

⇒ Define distance betw two species to be # of mutations needed to change one species to another

- Try to construct a phylogeny based on distance info among species

Copyright 2011 © Limsoon Wong

16

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 of as 2 mutations
- However, this fails to capture "multiple" mutations on the same site
- In practice, need to apply some corrective distance transformation

Copyright 2011 © Limsoon Wong

17

Distance Based

- **Input: Distance matrix M satisfying constraints**
 - M should satisfy metric space properties
 - 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

Copyright 2011 © Limsoon Wong

18

Metric Space

- A distance metric M which satisfies
 - Symmetry

$$M_{ij} = M_{ji} \geq 0$$
 - Self identity

$$M_{ii} = 0$$
 - Triangular inequality

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

Copyright 2011 © Limsoon Wong

19

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

Copyright 2011 © Limsoon Wong

20

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

Copyright 2011 © Limsoon Wong

21

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

Copyright 2011 © Limsoon Wong

22

Properties of Additive Metric

- Buneman's 4-point condition

M is additive if and only if for every 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

Copyright 2011 © Limsoon Wong

23

Proof

Figure 8.3: Buneman's 4-Point Condition

$$\begin{aligned}
 & M_{ik} + M_{jl} \\
 &= (M_{ix} + M_{xy} + M_{yk}) + (M_{jx} + M_{xy} + M_{yl}) \\
 &= M_{ix} + M_{jx} + M_{yk} + M_{yl} + 2M_{xy} \\
 &= M_{jk} + M_{il} \\
 &= (M_{jx} + M_{xy} + M_{xk}) + (M_{il} + M_{xy} + M_{yl}) \\
 &= M_{jx} + M_{il} + M_{yk} + M_{yl} + 2M_{xy} \\
 &= M_{ij} + M_{kl} \\
 &= M_{ix} + M_{xj} + M_{yk} + M_{yl}
 \end{aligned}$$

So it can be easily verified that: $M_{ik} + M_{jl} = M_{il} + M_{jk} \geq M_{ij} + M_{kl}$.
 (⇒) Will not present here. ■

Copyright 2011 © Limsoon Wong

24

Peter Buneman

JOURNAL OF COMBINATORIAL THEORY (B) 17, 48-50 (1974)

A Note on the Metric Properties of Trees*

PETER BUNEMAN*

Communicated by Frank Harary
Received February 21, 1973

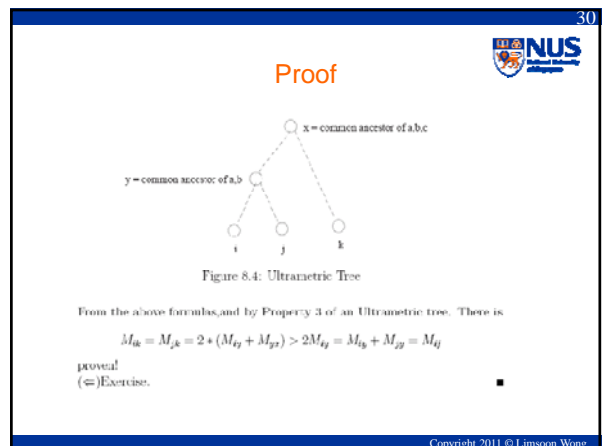
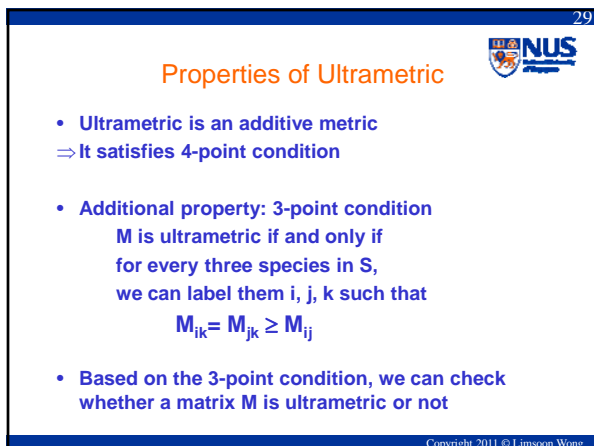
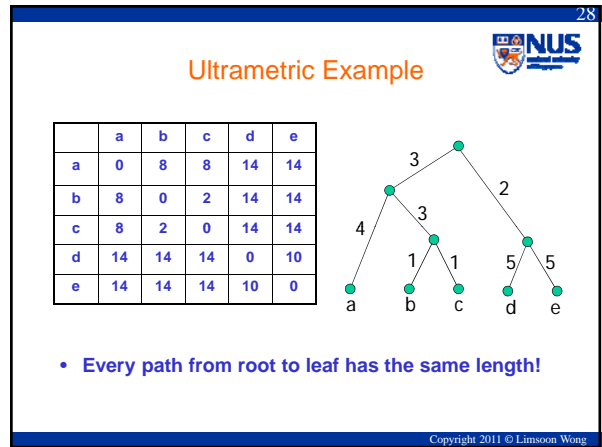
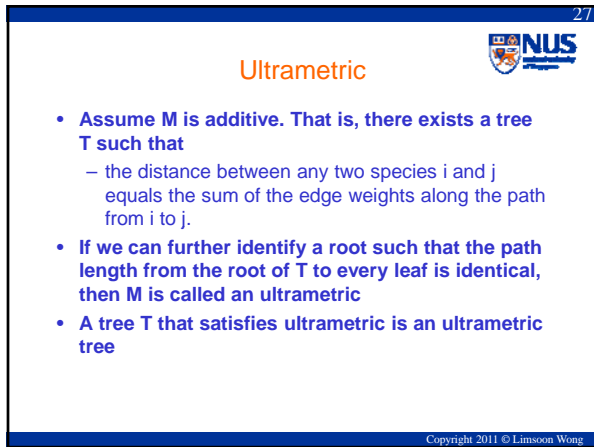
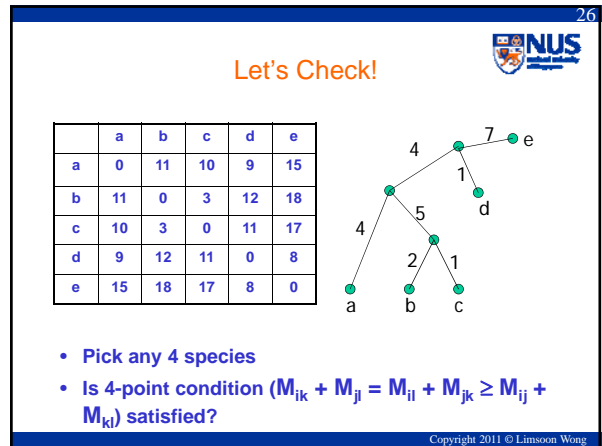
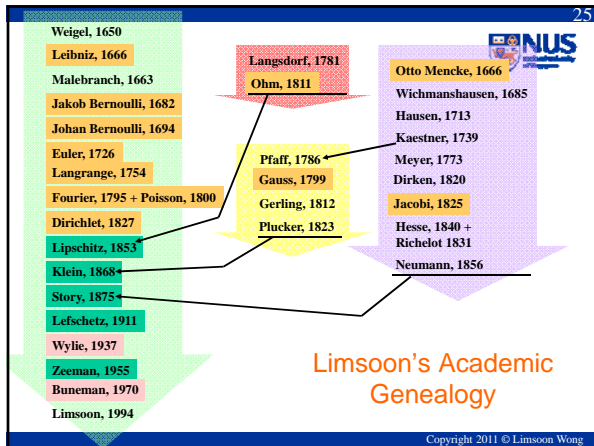
By checking the possible configurations of paths which can connect four points a, b, c, d in a tree, it can be seen that the graphical distance [1] must satisfy the inequality:

$$d(a, b) + d(c, d) \leq \max \{ d(a, c) + d(b, d), d(a, d) + d(b, c) \}$$

We shall refer to this condition as the four point condition: it is stronger than the triangle inequality (put $c = d$) and is equivalent to saying that of the three sums $d(a, b) + d(c, d)$, $d(a, c) + d(b, d)$, and $d(a, d) + d(b, c)$ two are equal and not less than the third. The four-point condition is also a sufficient condition for a graph to be a tree in the following sense.

THEOREM 1. A graph is a tree iff it is connected, contains no triangles, and has graphical distance satisfying the four-point condition.

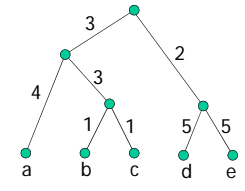
Copyright 2011 © Limsoon Wong



31

Let's Check!

	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



- Pick any 3 species
- Is 3-point condition ($M_{ik} = M_{jk} \geq M_{ij}$) satisfied?

Copyright 2011 © Limsoon Wong

32

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!**

Copyright 2011 © Limsoon Wong

33

Some Computational Problems

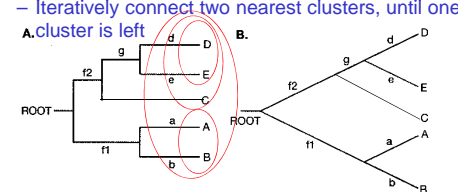
- Let M be a distance matrix for a set of species S
 - If M is ultrametric, can we reconstruct the corresponding ultrametric tree T in polynomial time? (only consider this one!)
 - If M is additive, can we have a polynomial time algorithm to recover the corresponding additive tree T ?
 - If M is not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

Copyright 2011 © Limsoon Wong

34

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

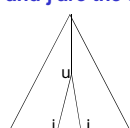


Copyright 2011 © Limsoon Wong

35

Definition - Height

- For a node u , define $height(u)$ 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, $height(u) = M_{ij}/2$



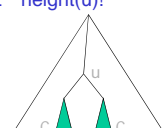
Copyright 2011 © Limsoon Wong

36

Distance Betw Two Clusters

- For any two clusters C_1 and C_2 of T
 - Define

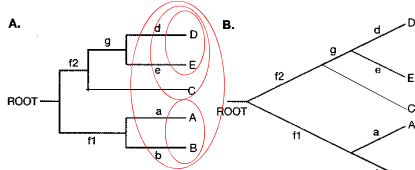
$$dist(C_1, C_2) = \frac{\sum_{i \in C_1, j \in C_2} M_{ij}}{|C_1| \cdot |C_2|}$$
 - Note that $dist(C_1, C_2) = M_{ij}$ for all $i \in C_1$ and $j \in C_2$. Why?
 - Let u be lowest common ancestor of i and j .
 $dist(C_1, C_2) = 2 * height(u)$



Copyright 2011 © Limsoon Wong

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

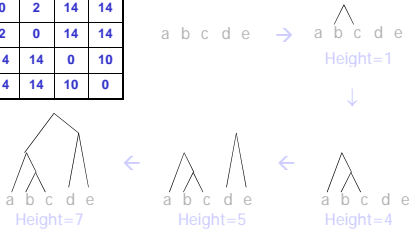


Algorithm

- 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 $O(n^2)$


The above is not obvious.
Can you identify the difficulty?
Can you solve it?

Achieving Quadratic Complexity

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

- Use a vector $V[i]$ to record the column id j such that $M[i, j]$ is min of row i
- When searching for clusters to merge, look for $x = \text{argmin}_i V[i]$
- Then merge cluster x with $V[x]$

Phylogenetic Tree Comparison



43

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 (Cavender-Farris-Neyman model, Jukes-Cantor Model)
 - Different reconstruction algorithms
- Tree comparison helps us to gain information from multiple trees

Copyright 2011 © Limsoon Wong

44

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

Copyright 2011 © Limsoon Wong

45

Restricted Subtree

- Consider tree T

Copyright 2011 © Limsoon Wong

46

Agreement Subtree

Copyright 2011 © Limsoon Wong

47

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

Copyright 2011 © Limsoon Wong

48

MAST for Rooted Trees

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

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

Copyright 2011 © Limsoon Wong

MAST by Dynamic Programming



Notations

- For any two binary rooted trees T_1 and T_2 , let $MAST(T_1, T_2)$ 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 ,

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

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

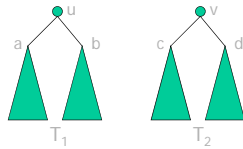
$$MAST(T_1^u, \Lambda) = 0, MAST(\Lambda, T_2^v) = 0$$

Recurrence (I)



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

$$\max \begin{cases} MAST(T_1^a, T_2^c) + MAST(T_1^b, T_2^d) \\ MAST(T_1^a, T_2^d) + MAST(T_1^b, T_2^c) \\ MAST(T_1^a, T_2^v) \\ MAST(T_1^b, T_2^v) \\ MAST(T_1^u, T_2^c) \\ MAST(T_1^u, T_2^d) \end{cases}$$

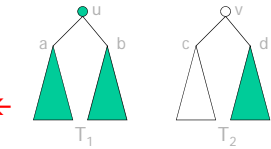


Recurrence (II)



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

$$\max \begin{cases} MAST(T_1^a, T_2^c) + MAST(T_1^b, T_2^d) \\ MAST(T_1^a, T_2^d) + MAST(T_1^b, T_2^c) \\ MAST(T_1^a, T_2^v) \\ MAST(T_1^b, T_2^v) \\ MAST(T_1^u, T_2^c) \\ MAST(T_1^u, T_2^d) \end{cases}$$



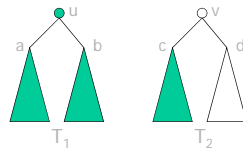
All the species in "agreement" are in right subtree of v.

Recurrence (III)



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

$$\max \begin{cases} MAST(T_1^a, T_2^c) + MAST(T_1^b, T_2^d) \\ MAST(T_1^a, T_2^d) + MAST(T_1^b, T_2^c) \\ MAST(T_1^a, T_2^v) \\ MAST(T_1^b, T_2^v) \\ MAST(T_1^u, T_2^c) \\ MAST(T_1^u, T_2^d) \end{cases}$$



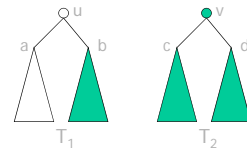
All the species in "agreement" are in left subtree of u.

Recurrence (IV)



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

$$\max \begin{cases} MAST(T_1^a, T_2^c) + MAST(T_1^b, T_2^d) \\ MAST(T_1^a, T_2^d) + MAST(T_1^b, T_2^c) \\ MAST(T_1^a, T_2^v) \\ MAST(T_1^b, T_2^v) \\ MAST(T_1^u, T_2^c) \\ MAST(T_1^u, T_2^d) \end{cases}$$



All the species in "agreement" are in right subtree of v.

61

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?

Copyright 2011 © Limsoon Wong

62

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

150000 years ago 100000 years ago 50000 years ago present

● African ○ Asian ■ Papuan □ European

Copyright 2011 © Limsoon Wong

63

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

Copyright 2011 © Limsoon Wong

64

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

Copyright 2011 © Limsoon Wong

65

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

Copyright 2011 © Limsoon Wong

66

Origin of Polynesians

- Do they come from Asia or America?

Copyright 2011 © Limsoon Wong

67

In the course of evolution...

Copyright 2011 © Limsoon Wong

68

Origin of Polynesians

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

Copyright 2011 © Limsoon Wong

69

Neanderthal vs Cro Magnon

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

Neanderthal

Cro Magnon

Copyright 2011 © Limsoon Wong

70

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

Copyright 2011 © Limsoon Wong

71

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

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

Copyright 2011 © Limsoon Wong

72

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?

Copyright 2011 © Limsoon Wong

