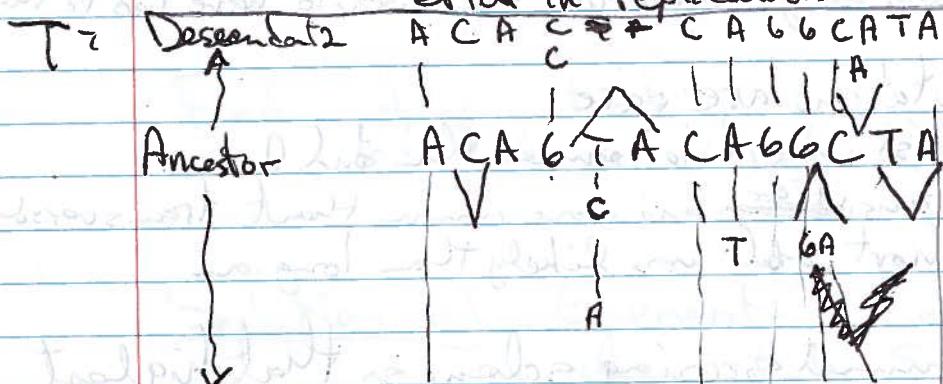


Sept 12 2016

Sequence alignment

Mutation: Change in sequence due to ~~error of~~ error in replication or unrepaird DNA damage



S: Descendant

Substitution: Change one nucleotide for another

A \leftrightarrow T transition

C \leftrightarrow G transition
(~3 times less common)

Insertion: Insert one more nucleotide

Deletion: Delete one or more nucleotides

Alignment of S, T is obtained by inserting gaps in S and T to obtain S' and T', so that in S' and T', nucleotides derived from same ancestor occur at same position.
Gaps reflect nucleotides missing due to del. or insertion.

T: A C A C - - C A G G - - C A T A

S: A - A G A A C T G G G A C - -

Match Mismatch Indel

Problem: Ancestor and mutation scenarios are generally not known.

Idealized Pairwise Global Alignment Problem

Given: S, T

Find: Alignment of S and T that is most likely to reflect the evolutionary scenario that would have led to them

- Key idea:
- ① Mutations are rare
 - ② Subst. are more common than indel
 - ③ Transitions are more common than transversions
 - ④ Short indels more likely than long ones

⇒ Define alignment scoring scheme so that highest scoring alignment is the one that is the most likely to be correct.

Scoring a given alignment

S: A C A G T C - - T A

T: A T A - T A A A T A

$$1 -1 1 -2 1 -1 -2 -2 1 = -3$$

Score: Subst. score + indel score

Subst. / indel cost matrix

$$M = \begin{pmatrix} A & C & G & T \\ C & +1 & -2 & -1 & +2 \\ G & -2 & +1 & -2 & -1 \\ T & -1 & -1 & -1 & -1 \end{pmatrix}$$

cost of gap of length l
Gap penalties: $g(l)$

- linear gap penalty: $b_l = -2l$

cost c per gap character

- affine: $a + bl = -5 - 1l$

Pairwise seq. Aln. Problem (linear gap penalty)

Given: $S = S_1 \dots S_m$
 $T = T_1 \dots T_n$ Sequences
 M subst. cost matrix
 c

Find: Alignment of S, T
s.t. that $\text{score}(\text{Aln}(S, T))$ is maximized

Solution #1: Enumerate and evaluate all possible alignments for S on T , report best

Problem: Way too slow $\mathcal{O}(2^{m+n})$

Solution #2: Needleman - Wunck algo. (1970)

$S = AGCT$
 $T = ACAG$

Dynamic Prog Algo

$s_1 s_2 s_3 s_4$
 S = A C A T
 T = A G T
 $t_1 t_2 t_3$

	A	G	T	
X =	C			
	A			
	T	O		

$X_{i,j}$ = Score of the best alignment for $s_1 \dots s_i$ against $t_1 \dots t_j$

Note: We want $X_{m,n}$, and we want the alignment that achieves this score.

How to calculate $X_{i,j}$?

Alignment

Score

Best Aln $(s_1 \dots s_i \mid t_1 \dots t_j)$ either looks like

Best Aln $(s_1 \dots s_{i-1} \mid t_1 \dots t_{j-1})$

$\rightarrow X_{i-1,j-1} + M(s_i, t_j)$

OR

Best Aln $(s_1 \dots s_{i-1} \mid t_1 \dots t_j)$ $s_i \rightarrow X_{i-1,j} + c$

OR

Best Aln $(s_1 \dots s_i \mid t_1 \dots t_{j-1})$ $t_j \rightarrow X_{i,j-1} + c$

So: $X_{i,j} = \max \begin{cases} X_{i-1,j-1} + M(s_i, t_j) \\ X_{i-1,j} + c \\ X_{i,j-1} + c \end{cases}$

$$X_{0,j} = j \cdot c$$

$$X_{i,0} = i \cdot c$$

$$X_{0,0} = 0$$

.	A	6	T	
I	0	-2	-4	-6
A	-2	+1	0	-3
C	-4	-1	1	-1
A	-6	-3	-2	-3
T	-8	-5	-4	-1

$$c = -2$$

Calculate X entries to right, top

For each entry, keep to entry that achieved

$X_{4,3} = -1$: Score of best alignment

We recover best alignment following backpointers from (m, n)

↖ : Match

← : Gap in S

↑ : Gap in T

Best = A C A T

A m - 6 T

Running time: $O(m \cdot n)$
space: $O(m \cdot n)$

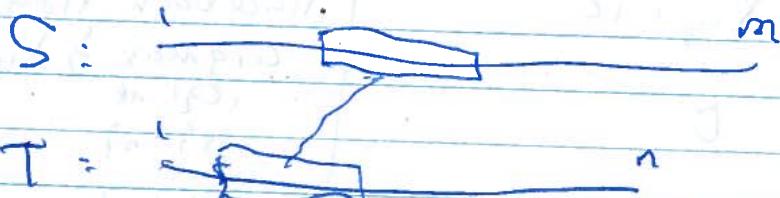
Note: NW algo can be modified with affine gap score
 → Why does it not align?
 → What is the modification?

To recover all optimal alignments, do Depth-First search on back pointers starting from (m, n) .

Pairwise local alignment problem.

Idea: Often, two sequences only share significant sim over a particular region.

Examples? Two



Given: $\left\{ \begin{array}{l} \text{Sequences } S = s_1 \dots s_m \\ \text{Subst. matrix } M \\ \text{gap penalty } c \end{array} \right.$

Find: ~~the best local alignment~~ Indices i, j, k, l , where $\left\{ \begin{array}{l} i \leq j \leq m \\ k \leq l \leq n \end{array} \right.$
such that Score(Best Aln $(s_i \dots s_j, t_k \dots t_l)$) is maximized

Can we do better than?

for $i=1 \dots m$

for $j=1 \dots m$

for $k=1 \dots n$

for $l=k \dots n$

NW($s_i \dots s_j, t_k \dots t_l$)

$O(m^3 n^3)$

Smith Waterman algo.

$X_{i,j}$ = Score of best local alignment for $s_i \dots s_j, t_k \dots t_l$
where s_i and t_j are included in alignment

$$X_{i,j} = \max \left(X_{i,j-1} + M(s_i, t_j), X_{i-1,j} + M(s_i, t_j) \right)$$

Max

$$X_{i-1,j} + c$$

$$X_{i,j-1} + c$$

0

$$X_{i,0} = 0 \quad \forall i=1 \dots m$$

$$X_{0,j} = 0 \quad \forall j=1 \dots n$$

Traceback from ~~max~~

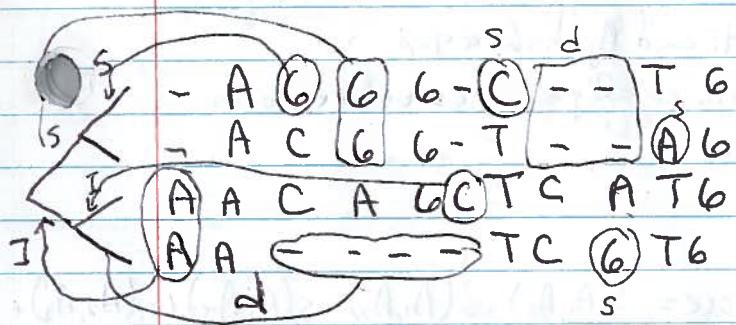
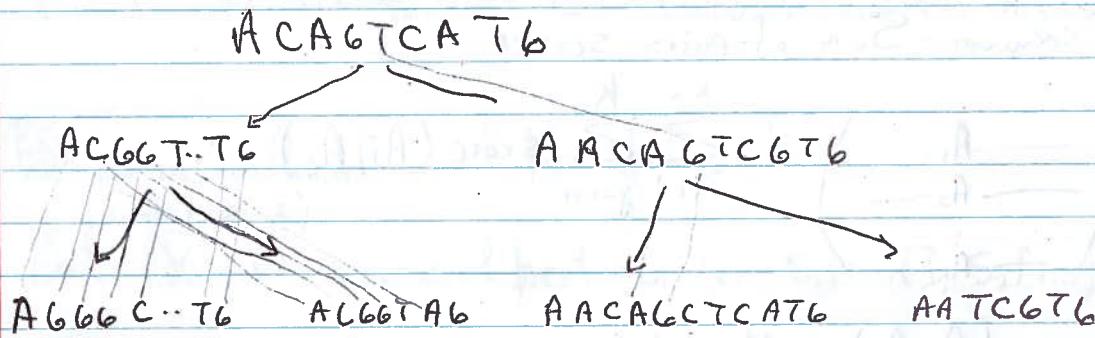
$$\arg \max_{\substack{i \in \{1 \dots m\} \\ j \in \{1 \dots n\}}} \{ X_{i,j} \}$$

until hitting zero

4520 Pontiac
→ 3452

Notes: Local alignment problem only makes sense
when ~~Σ M_{match}~~ < 0
 $\sum_{\text{Match}} - \sum_{\text{Mismatches}}$

Multiple seq. alignment.



Biological problem: Given n sequences
Phylogenetic tree T , where each sequence
is assigned to a leaf

Find: Multiple alignment where all
nucleotides derived from the same
common ancestor are aligned to each other.

Scoring a given MSA

Goal: Find a function $\text{score}: \text{MSA} \rightarrow \mathbb{R}$

such that $\text{score}(\text{MSA}) > \text{score}(\text{MSA}') \Rightarrow \text{MSA} \text{ is more likely to be biologically correct than MSA}'$

Many choices of scoring functions are possible

Most common: Sum-of-pairs scores

$$\text{Score} \left(\begin{array}{c} A_1 \\ A_2 \\ \vdots \\ A_K \end{array} \right) = \sum_{i=1}^{K-1} \sum_{j=i+1}^K \text{score}(A_i, A_j)$$

where $\text{score}(A_i, A_j)$ is obtained by

- ① Removing positions where both A_i and A_j have a gap
- ② Scoring the remaining pairwise alignment as before, with subst. Matrix M, gap penalty c.

Example:

① A C T
 ② A G -
 ③ - G T
 ④ - G T

$$\text{Score} = s(A_1, A_2) + s(A_1, A_3) + s(A_1, A_4) + s(A_2, A_3) + s(A_2, A_4) + s(A_3, A_4) + s(A_3, A_4)$$

$$= -3 + -6 + -6 + -3 + -3 + -3 + 2$$

$$\approx -19$$

Optimal MSA problem

Given: Sequences $S_1 \dots S_n$
 Subst. matrix M
 Gap penalty c

Find: MSA for $S_1 \dots S_n$ that achieves the highest possible score

~~generalization of N-W algorithm~~
 (For 3 seq.)

Let $X_{i,j,k}$ = score of best align for $S_1[i \dots i], S_2[j \dots j], S_3[k \dots k]$

$$X_{i,j,k} = \max \left\{ \begin{array}{l} X_{i-1,j-1,k-1} + M(S_1[i], S_2[j]) + M(S_1[i], S_3[k]) + M(S_2[j], S_3[k]) \\ X_{i-1,j-1,k} + M(S_1[i], S_2[j]) + 2c \\ \vdots \\ X_{i,j,k-1} + 2c \end{array} \right.$$

For 3 seq. of length l

Dyn Prog table has $l \times l \times l$ entries

For n seq of length l : $O(l^n \cdot 2^n)$

↓
 Exponential time algorithm

↗ Heuristic, with no guarantee of optimality

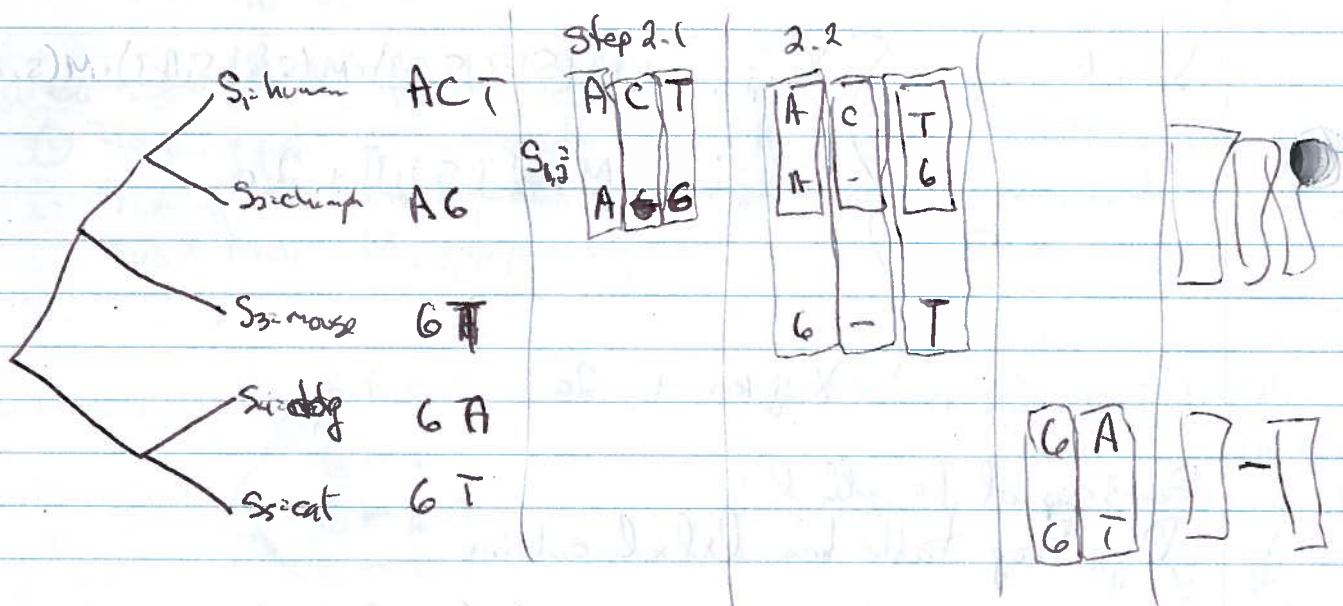
Progressive alignment Algo.

Input: $\{S_1 \dots S_n\}$
M
c

① Guess phylogenetic tree for $S_1 \dots S_n$

② For each internal node in ~~in~~ order traversal
~~(bottom up)~~
(from leaves to root)

~ Find optimal alignment b/w pairs of seq or align two chil(dren)

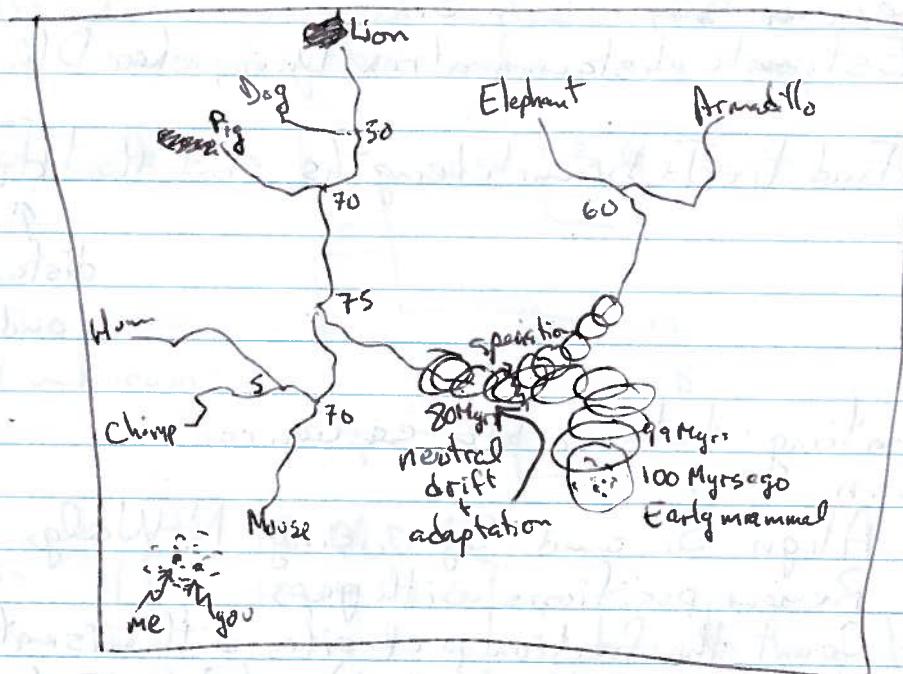


2.1. Align S_1 vs S_2 to obtain $S_{1,2}$

2.2 Align $S_{1,2}$ vs S_3 to obtain $S_{1,2,3}$

2.3 Align $S_{1,2,3}$ vs S_4 to obtain $S_{1,2,3,4}$

Sequence Evolution + Phylogenetics



$\{A, C, G, T\}^*$
= Space of all DNA sequences

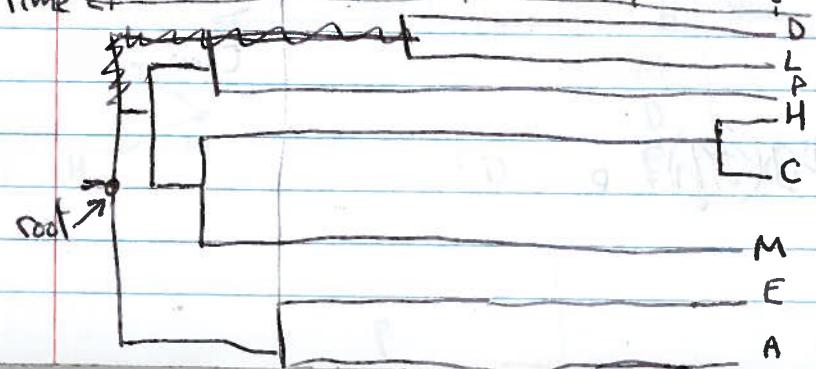
Phylogenetic Inference Problem (Biological version)

Given: DNA sequences from multiple species

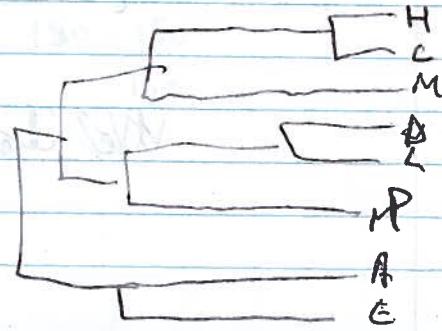
Find: Evolutionary tree that describes their evolution

Topology, Dating of speciation events

For our example, the solution would be:



Note: this is the same topology



Distance-based Phylo. Inference methods

Idea: Given seq. S_1, \dots, S_n

(A) Estimate distance matrix $D_{n \times n}$, where $D(i,j) = \text{distance}$
btw S_i, S_j

(B) Find tree + Branch lengths such that $d_T(i,j) \approx D(i,j)$

distance btwnodes
i and j in tree T

(A) Estimating distance btwn sequences

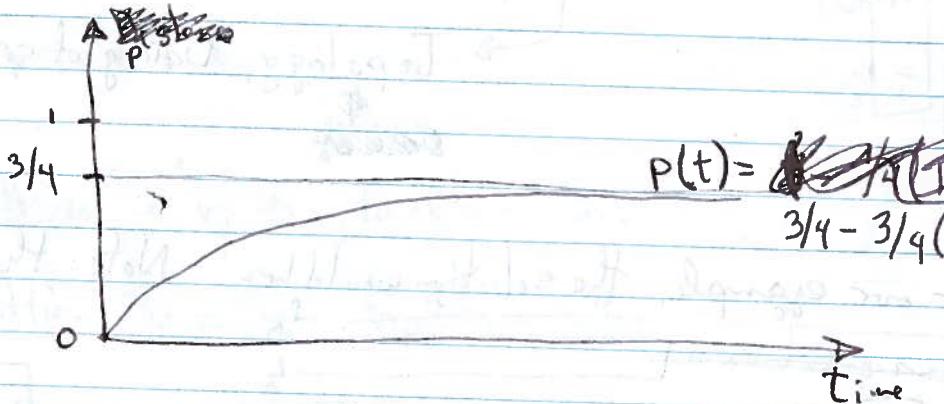
For $i, j = 1 \dots n$

1. Align S_i and S_j using N-Walg
2. Remove positions with gaps
3. Count the fraction p of sites with mismatches
4. $D(i,j) = -\frac{3}{4} \log(1 - \frac{4}{3}p)$ (Jukes-Cantor model)

Example $S_i: A C - T G A C T G \rightarrow p = 3/7$
 $S_j: A G T T A - C A G$
 $D(i,j) = 0.63$

Note:

If the two seqs.
are very long



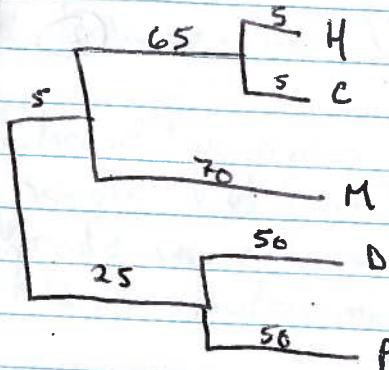
We choose $D(i,j) = p$

Under Kimura 3-parameter model with rates to transition

Notes More advanced distance measures consider separately transitions and transversions

Example:

True tree
(unknown)



$$D = \begin{matrix} & H & M & C & D & P \\ H & - & 10 & 140 & 150 & 150 \\ C & - & 140 & 150 & 150 & \\ M & - & 150 & 150 & & \\ D & - & 100 & & & \\ P & - & & & & \end{matrix}$$

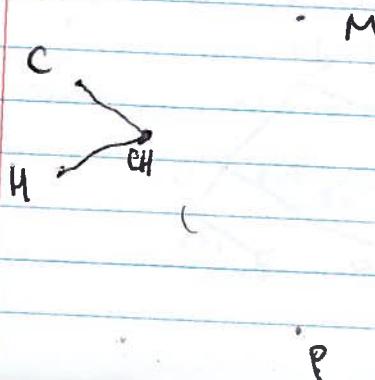
How to find tree?

UPGMA algorithm

Report

- ① Choose two closest nodes according to D_{ij}
- ② Merge two nodes into a new node w
- ③ Remove row/col u, v from matrix
- ④ Insert new row/col for w : $D(w,i) = D(u,i) - \frac{1}{2} D(u,v)$

	C	M	D	P
C		185	145	145
M		-	150	150
D				100
P				

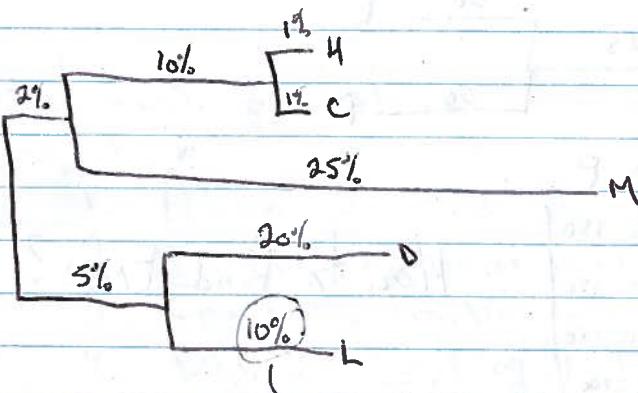


Running time: ?

Problem: Only works if

- ① Distances are estimated perfectly accurately
- ② Mutation rate is constant

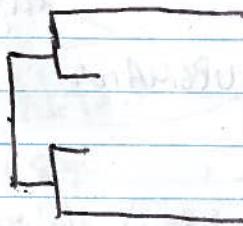
More realistic tree:



Mutation rate varies

⇒ Distances b/w seq
cannot be calculated in Myo
but instead in Expected
substitution per site

→ This means that ~~10%~~ positions will
have gotten mutated, the expected # of subst.
per site is 0.1



⇒ UPGMA w/|| produce the wrong tree

Unweighted Pair-Group Method with Arithmetic Mean

Reminder: UPGMA Algo

Input: $S_1 \dots S_n$

Output: Tree T with branch lengths

- ① Estimate distance matrix D from pairwise alignments.
- ② Repeat

2.1. Choose two nodes to pair up (u, v)

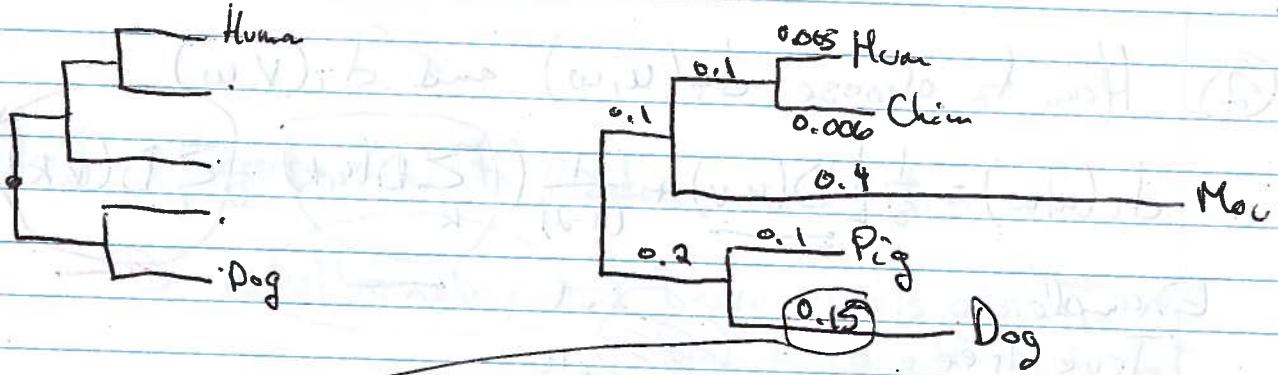
2.2. Remove (u, v) from D

2.3. Create new node w as ancestor of u, v

2.4. Add new row/column for w in D .

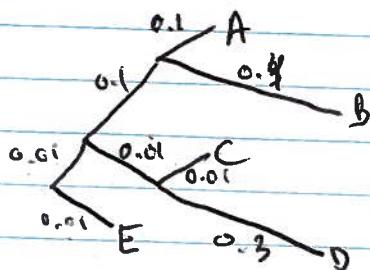
UPGMA produces correct output if

- ① Distances are estimated perfectly accurately
- ② Mutation rate is constant

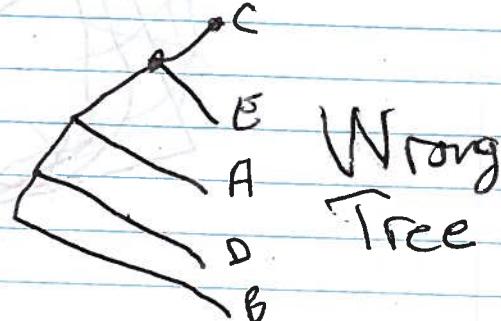


→ Expected number of subst. per site

Suppose true tree



UPGMA



Neighbor-joining Algo (Nei, Saitou)

Same as UPGMA, but different pair selection rule

Calculate Q_{min} , where $Q(i,j) = \sum_{k=1}^n D(i,k) + \sum_{k=i+1}^n D(j,k) - (n-2)D_{ij}$

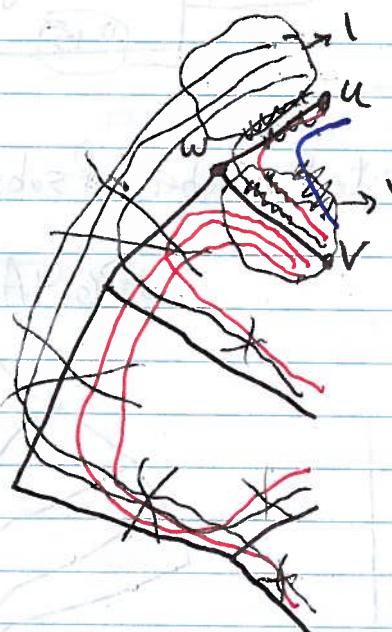
- ① Theorem: $\Leftrightarrow u, v \leftarrow \arg \max_{i, j} \{Q(i, j)\}$, then nodes u and v must be a cherry



- ② How to choose $d_T(u, w)$ and $d_T(v, w)$

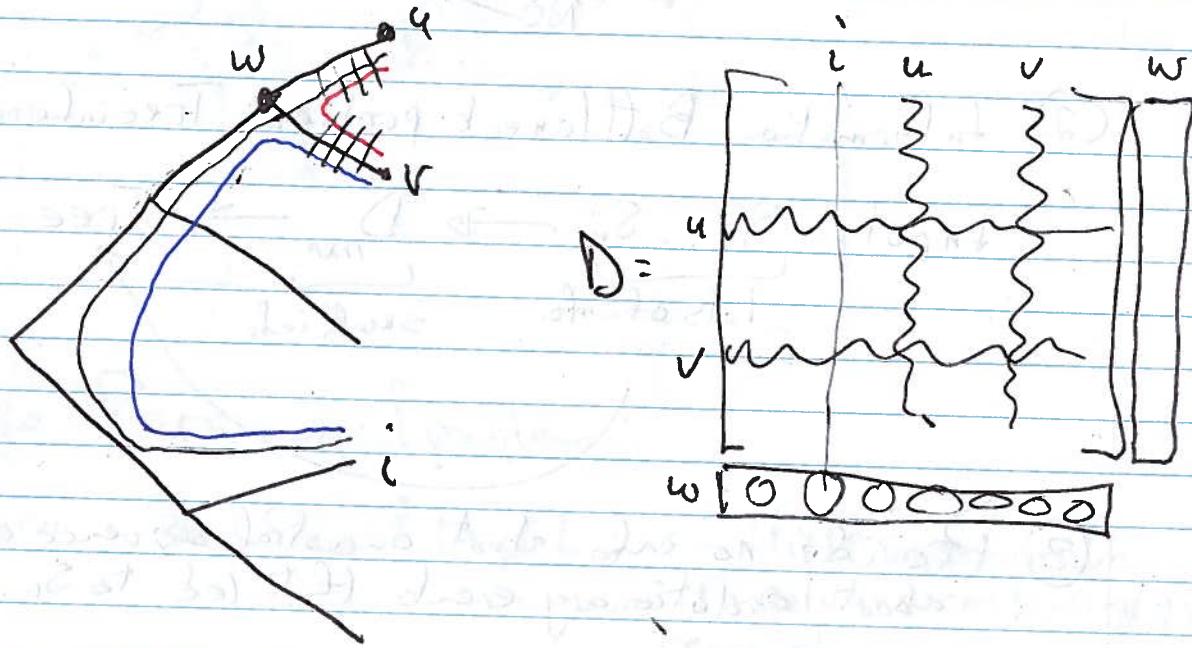
$$d_T(u, w) = \frac{1}{2} \left[\underbrace{D(u, v)}_{=} + \frac{1}{(n-2)} \left(\underbrace{\sum_k D(u, k)}_{} - \underbrace{\sum_k D(v, k)}_{} \right) \right]$$

Example:
True tree



Adding row/col for w in D

$$D(w, i) = \frac{1}{2} (D(u, i) + D(v, i) - D(u, v))$$



N-J is guaranteed to produce correct tree if

D is ultrametric \rightarrow There exists a tree with branch length such that

$$d_T(i, j) = D(i, j)$$

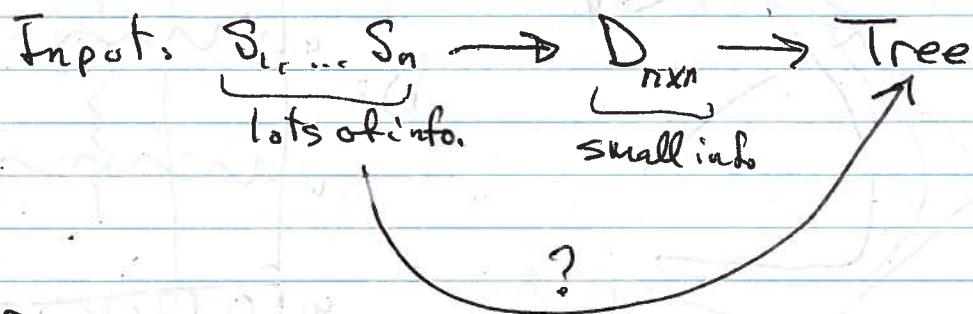
Running time: $O(n^4)$, but can be reduced to $O(n^3)$

Problems with N-J algo

① If D is not ultrametric (D is not estimated ~~perfectly~~)

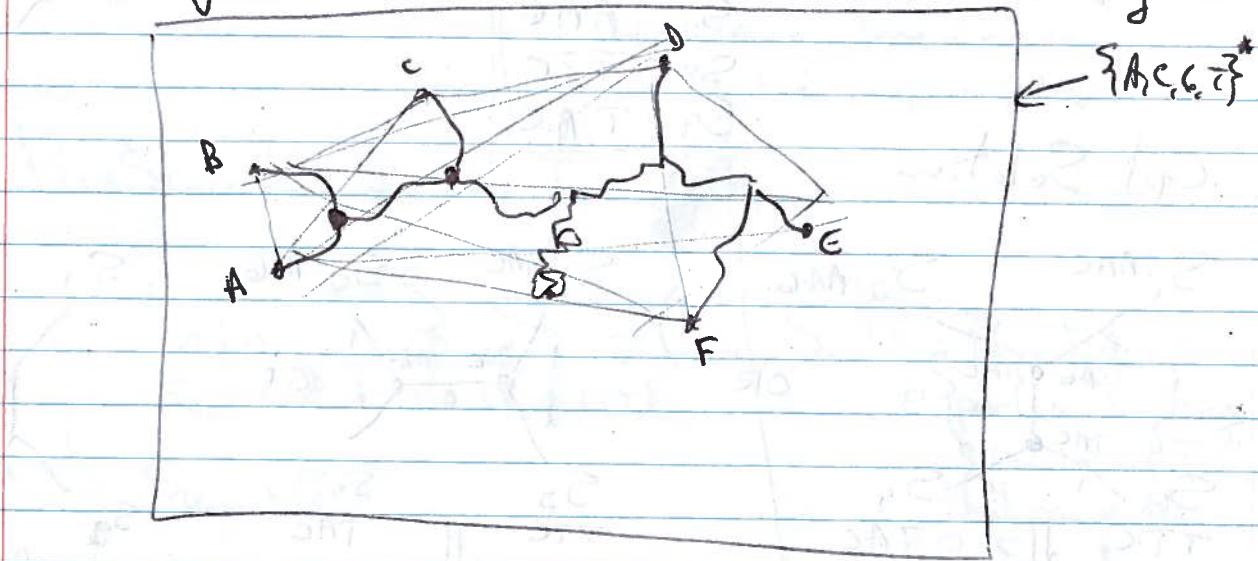
\Rightarrow N-J comes with guarantees about optimality
~~No~~

② Information Bottleneck problem: Tree inferred may be inaccurate



③ Provide no info about ancestral sequence or about evolutionary events that led to S_1, \dots, S_n

Phylogenetic Inference - Maximum Parsimony



Large Parsimony Problem:

Given: $S_1 \dots S_n$ of length L , in multiple sequence alignment, with columns containing gaps removed

- Find:
- Tree T with leaves labeled with $S_1 \dots S_n$
 - Ancestral ~~seq~~ sequence S_u for each internal ~~tree~~ node of T

such that $\sum_{(u,v) \in E(T)} d(S_u, S_v)$ is minimized

of substitutions b/w S_u, S_v

Example Input

MSA

$S_1: A\text{AC}$
 $S_2: A\text{AG}$
 $S_3: \text{TTC}$
 $S_4: \text{TAC}$

Unrooted tree



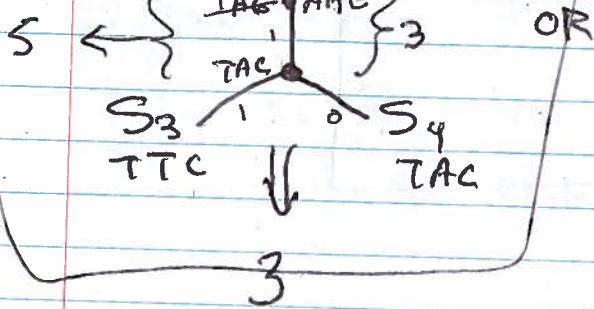
Rooted tree



Opt. Solution

$S_1: \text{AAC}$

$S_2: \text{AAG}$



OR

$S_1: \text{AAC}$

$S_2: \text{AAG}$

S_1

S_4

S_3

S_4

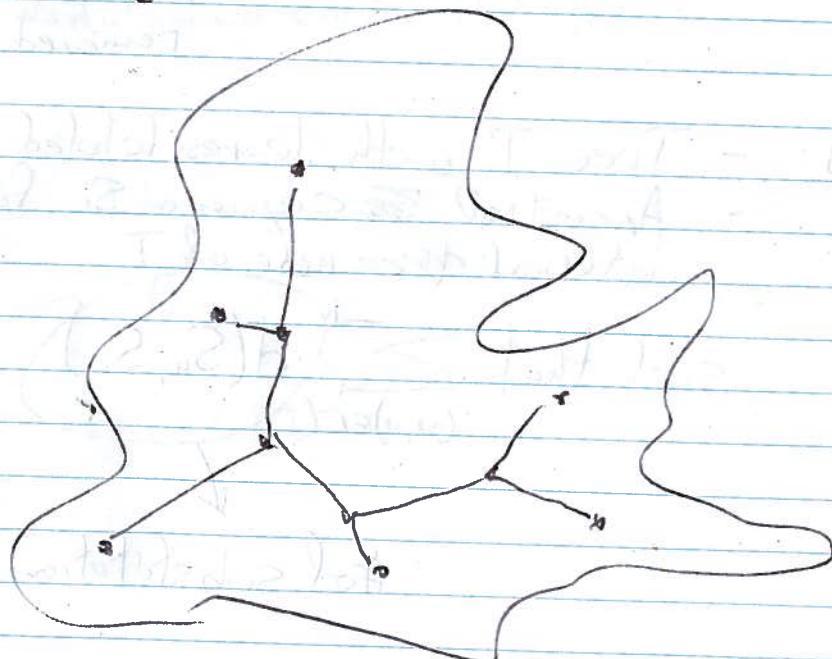
S_2

S_3

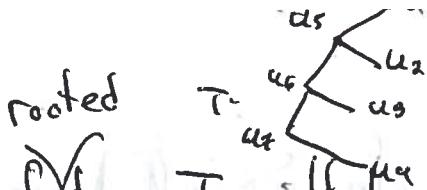
↓ 4

↓ 4

Steiner tree



Define Parsimony Score of tree T with leaves u_1, u_2, \dots, u_n
and internal nodes
 $u_{n+1}, u_{n+2}, \dots, u_{2n-1}$



$\text{ParsScore}(S_1, S_2, \dots, S_n, T) = \cancel{\text{min}}$

$$= \min_{\substack{S_{u_{n+1}} \in (\Sigma)^n \\ S_{u_{n+2}} \in (\Sigma)^n \\ \vdots \\ S_{u_{2n-1}} \in \Sigma^n}} \left\{ \sum_{(u, v) \in E(T)} d(S_u, S_v) \right\} \quad \begin{array}{l} \text{The minimum # of steps} \\ \text{performed along} \\ \text{branches of } T \\ \text{that can explain} \\ S_1, S_2, \dots, S_n \end{array}$$

$\left(4^L\right)^{n-1}$
 $= 4^{L(n-1)}$ cases

Observation: $\text{ParsScore}(S_1, \dots, S_n, T) = \cancel{\text{ParsScore}}$

$$\Rightarrow \sum_{i=1}^L \text{ParsScore}(S_i[i], \dots, S_n[i], T)$$

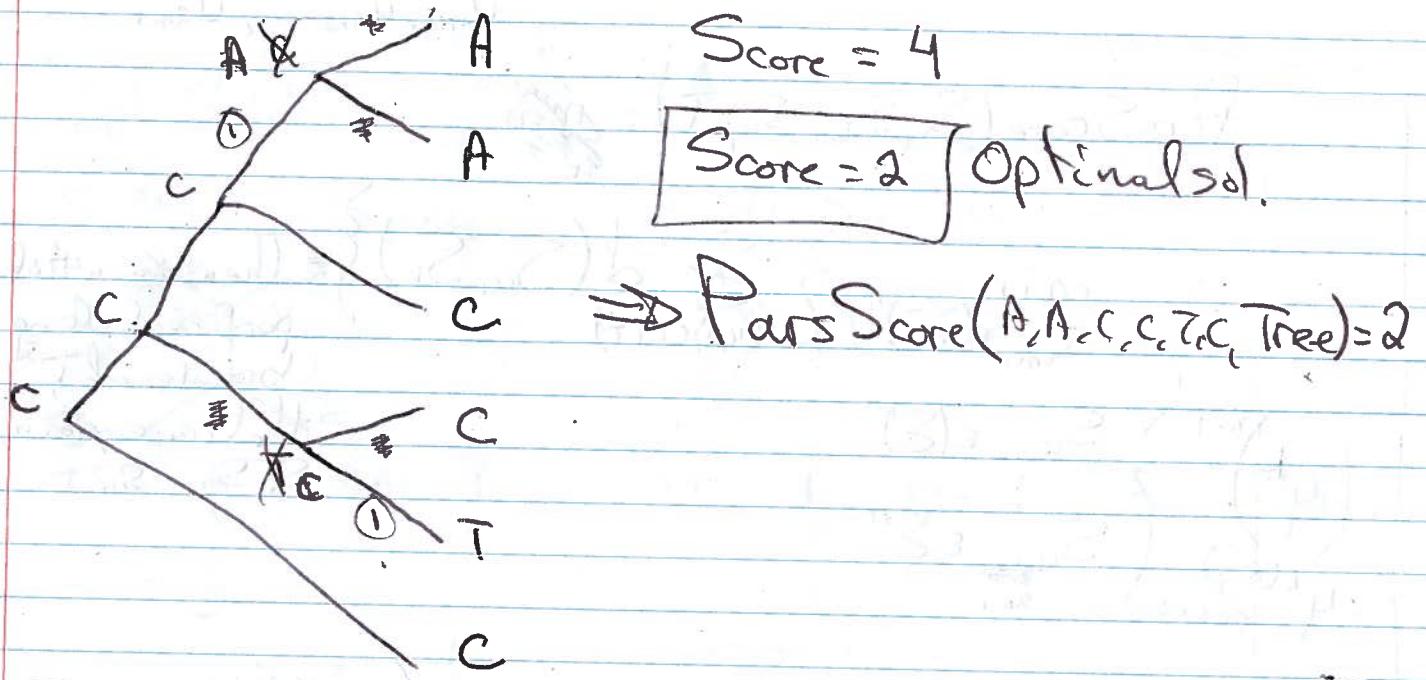
Small Parsimony Problem

Given:

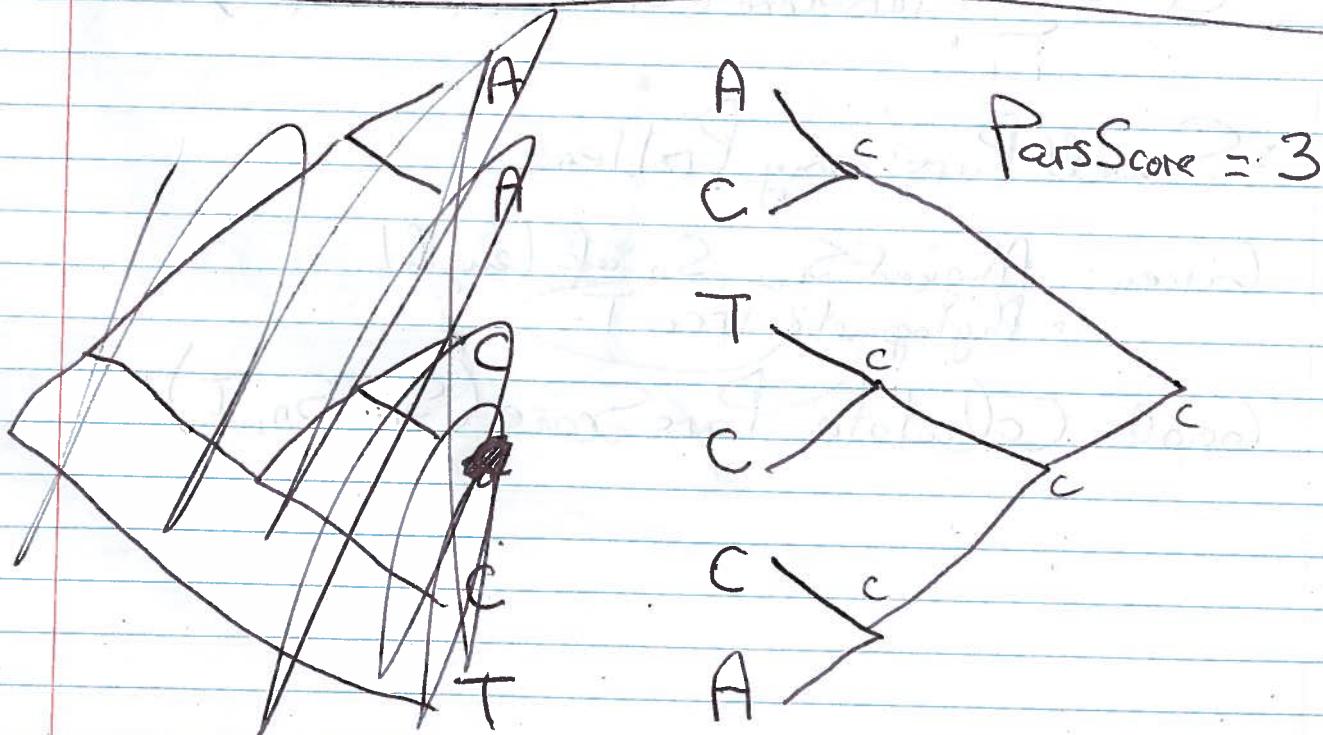
- Aligned S_1, \dots, S_n of length L
- Phylogenetic tree T

Goal: Calculate $\text{ParsScore}(S_1, \dots, S_n, T)$

Example = Small Parsimony Problem



How to calculate $\text{ParsScore}[S_1[i], S_2[i], \dots, S_n[i], T]$?

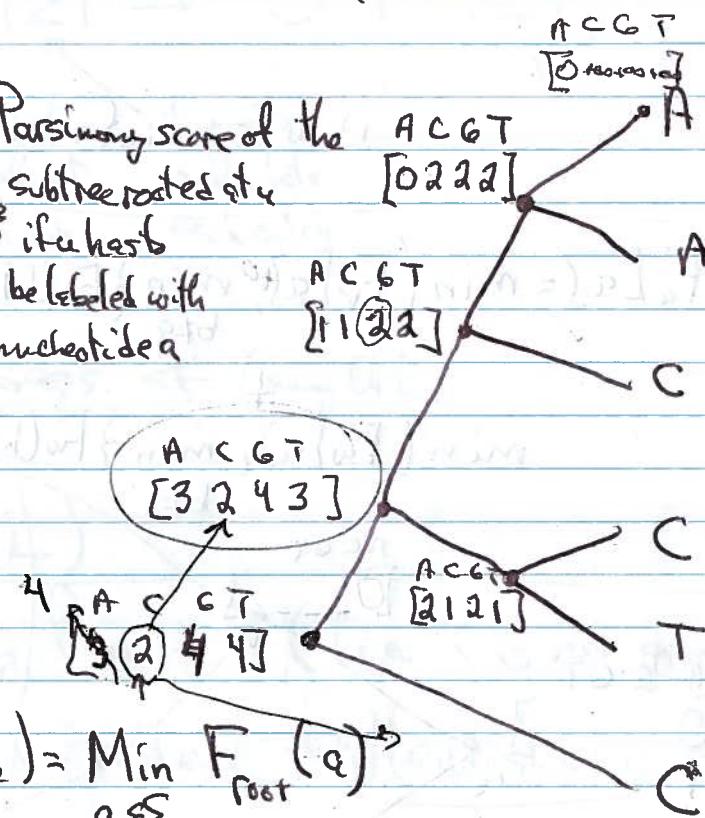


Sankoff Algorithm (1975)

Input: $\{S_1, S_2, \dots, S_n\}$ each of length l
 Tree T

Goal: Calculate ParsScore(S_1, \dots, S_n, T)

Define $F_u[a] =$ Parsimony score of the subtree rooted at node $u \in \{A, C, G, T\}$ if it has been labeled with nucleotide a



$$\text{ParsScore}(T) = \min_{a \in \Sigma} F_{\text{root}}(a)$$

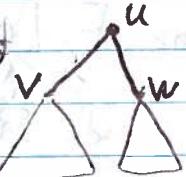
Algo: For each node u in tree (from leaves back to the root)

For each $a \in \{A, C, G, T\}$

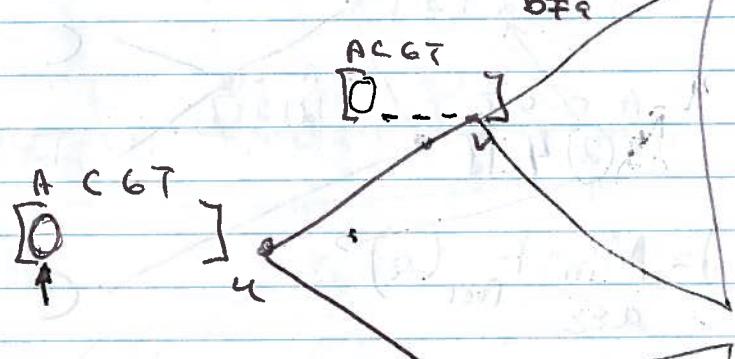
$$\text{if } u \text{ is a leaf } F_u[a] = \begin{cases} 0 & \text{if } S_u = a \\ +\infty & \text{if } S_u \neq a \end{cases}$$

if u is an internal node,

$$\Rightarrow F_u[a] = \min \left(F_v[a] + \min_{b \neq a} \{ F_v(b) \} + 1 \right)$$



$$\min \left(F_w[a] + \min_{b \neq a} \{ F_w(b) \} + 1 \right)$$



$v [3243]$

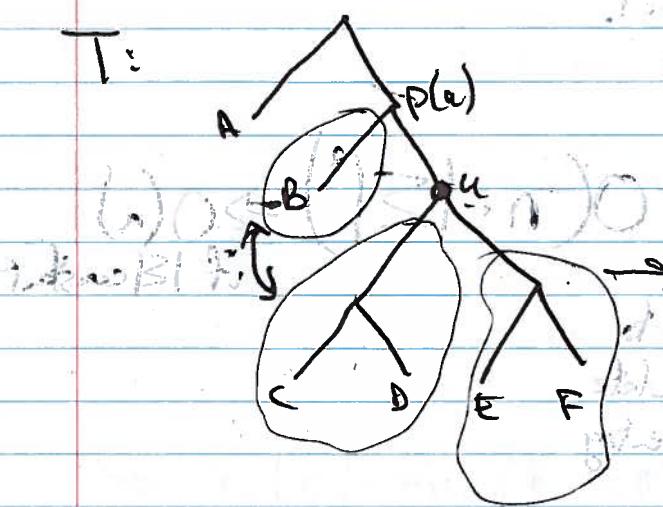
$$\min(+\infty, \min(0, \infty, +\infty) + 1)$$

$$= \min(30, \min(2043) + 1)$$

$$= \min(3, 2+1) = 3$$

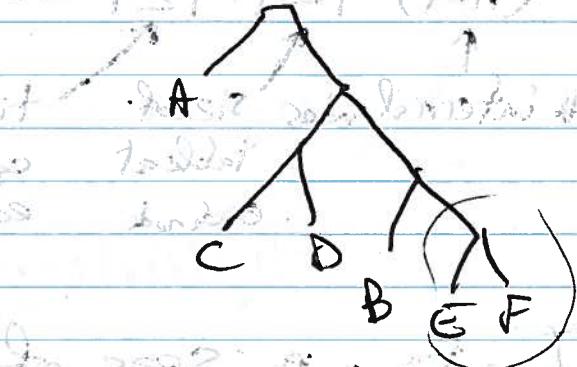
$w [+\infty 0 +\infty +\infty]$

Tree neighborhood definitions



Nearest-neighbor interchange (NNI)

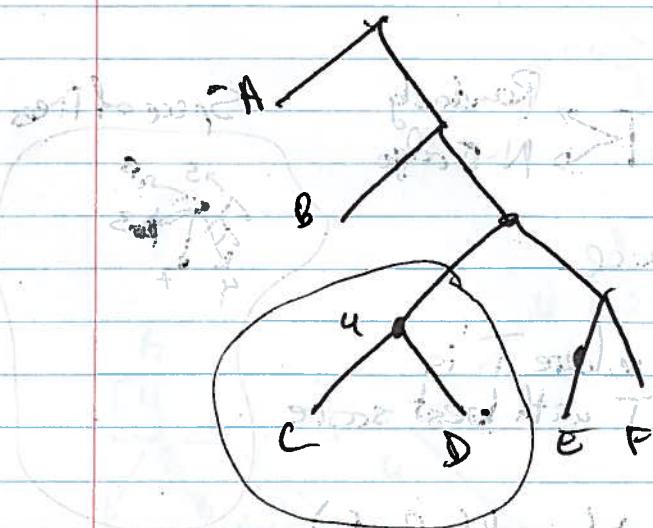
Pick some internal node u



$\text{Neighbors}(T) = \{ T' \text{ s.t. } T' \text{ can be obtained from } T \text{ with one NNI} \}$

$$|\text{Neighbors}(T)| = (n-2) \cdot 2$$

Subtree Pruned and Regrafted



Choose any node in T (root).

Cut - cut it off front

Reinsert - reinsert along other edge of T

