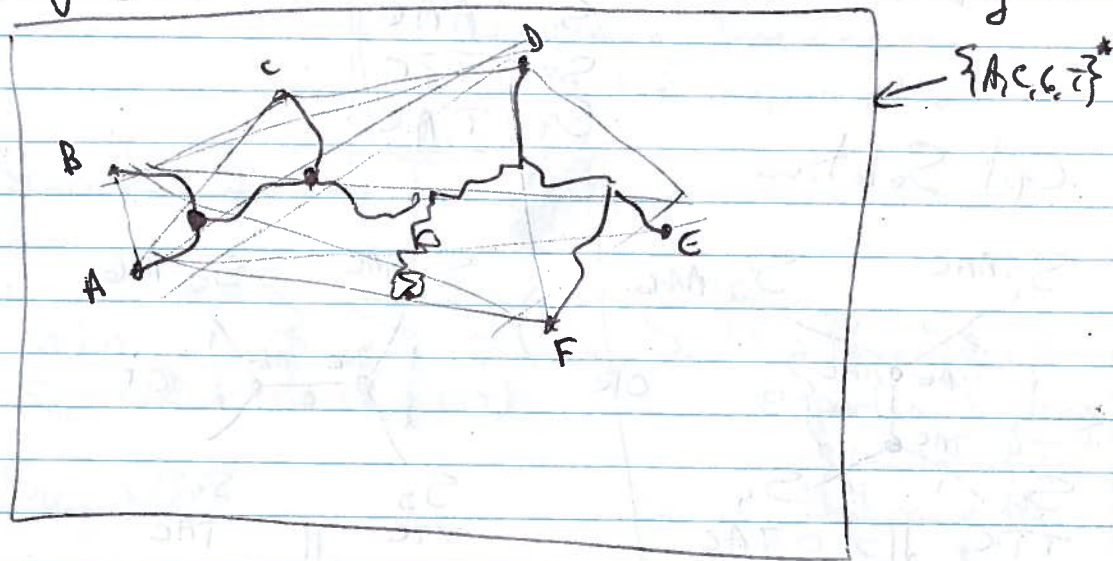


# Phylogenetic Inference - Maximum Parsimony



## 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 btw  $S_u, S_v$

Example: Input

MSA

$S_1: AAC$   
 $S_2: AAG$   
 $S_3: TTC$   
 $S_4: TAC$

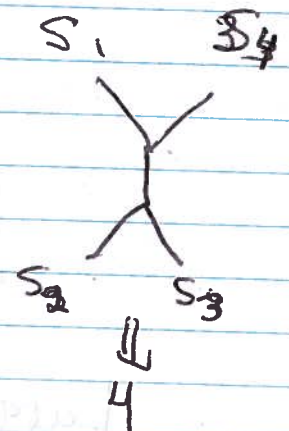
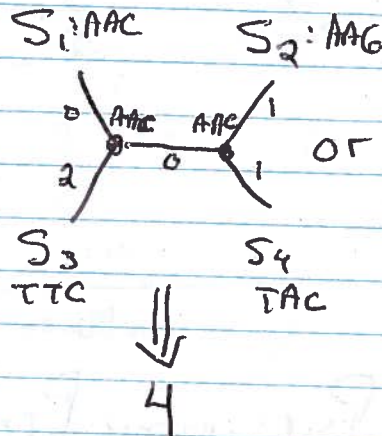
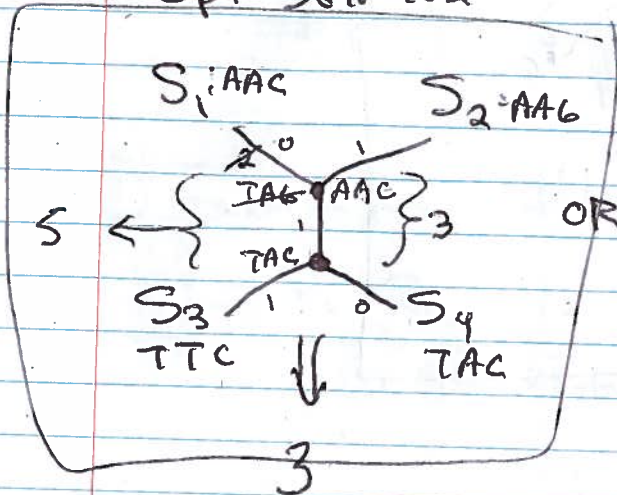
Unrooted tree



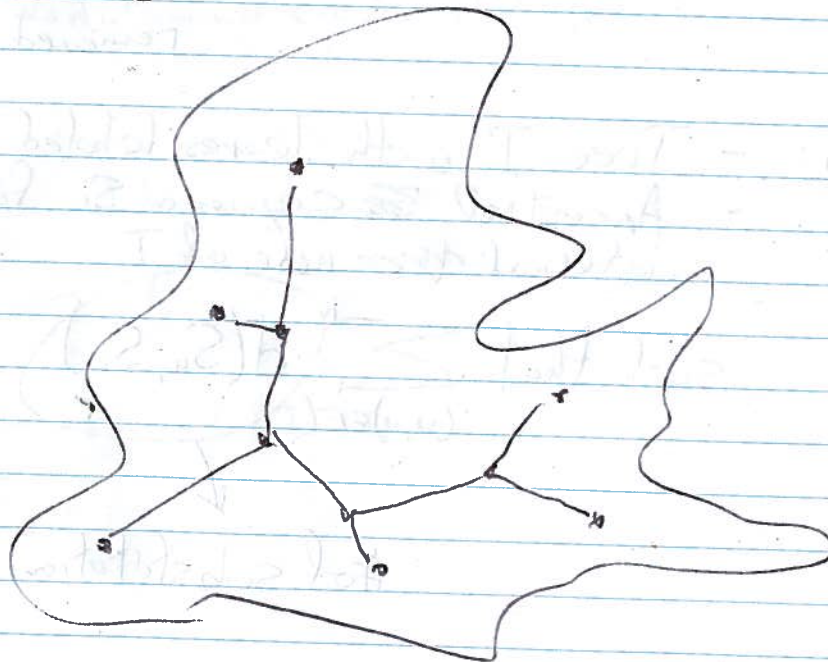
Rooted tree



Opt. Solution



Steiner tree





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

$$\text{ParseScore}(S_1, S_2, \dots, S_n, T) =$$

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

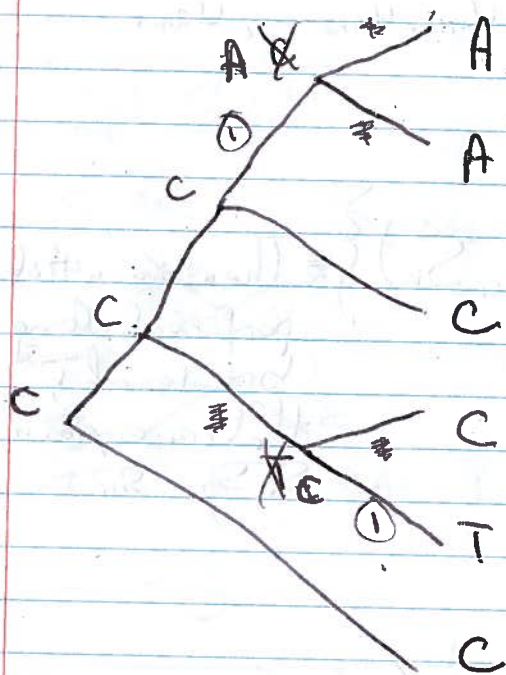
$$g = \sum_{i=1}^L \text{ParseScore}(S_1[i], \dots, S_n[i], T)$$

→ Small Parsimony Problem

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

Goal: Calculate Pars Score ( $S_1, \dots, S_n, T$ )

Example = Small Parsimony Problem



Score = 4

Score = 2 [Optimal sol.]

$$\Rightarrow \text{PairScore}(A, A, C, C, T, C, \text{Tree}) = 2$$

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

