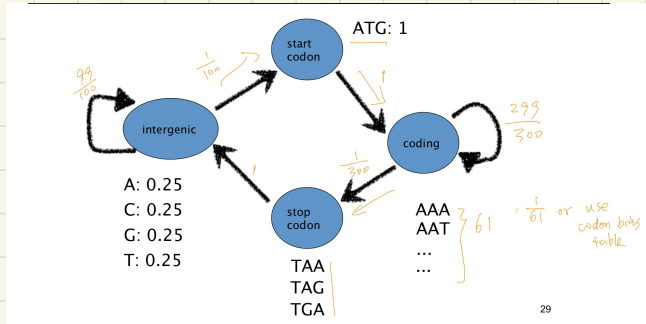


# Review

## ① HMM for prokaryotes gene prediction



- Dynamic programming

- recurrence relation different for intergenic & other states.

② promoters at upstream of the gene.

③ PSWM, PSFM, PSSM.

- motif.

④ phylogeny, virus, chain letter.

# A sample letter:

Trust in the Lord with all your heart

Trust in the Lord with all your heart and he will acknowledge and He will light the way. This Prayer has been sent to you for good luck. The original copy is from the Netherlands. It has been around the world nine times. The luck has been brought to you. You are to receive good luck within four days of receiving this letter. This is nojoke. You will receive it in the mail. Send copies of this letter to people you think need good luck. Do not send money. Do not keep this letter. It must leave your hands within ninety six hours after you receive it. An RAF officer received \$70,000. Don Elliott received \$50,000 and lost it because he broke the chain. While in the Phillipines, General Welch lost his life six days after he received this letter. He failed to circulate the Prayer. However, before his death, he received \$775,000. Please send twenty copies and see what happens to you on the fourth day. This chain comes from Venezuela and was written by Sol Anthony De Cadif, a missionary from South America. Since this chain must make a tour of the world. you must make twenty copies identical to this one and send it to your friends, parents, and acquaintances. After a few days you will get a surprise. This is true, even if you are not superstitious. Take note of the following. Constantine Diaz received the chain in 1953. He asked his secretary to make twenty copies and send them. A few days later he won a lottery for two million dollars in his country. Carlo Craduit, and office employee, received the chain. He forgot it and in a few days lost his job. He found the chain and sent it to twenty people. Five days later he got an even better job. Dolin Moirchild received the chain and not believing in it, threw it away. Nine days later he died. For no reason what so ever should this chain be broken

send money. Do not keep this letter. It must leave your hands within ninety six hours after you receive it. An RAF officer re

A few days later he won a lottery



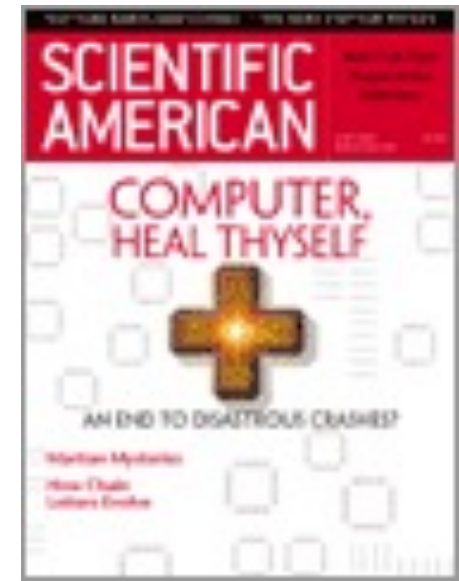
the chain. While in the Phillipines, General Welch lost his life six days after he received this letter. He failed to circulate the Prayer. However, before his death, he received \$775,000. Please send twenty copies and see what happens to you on the fourth day. This chain

# Chain letters – old style

---

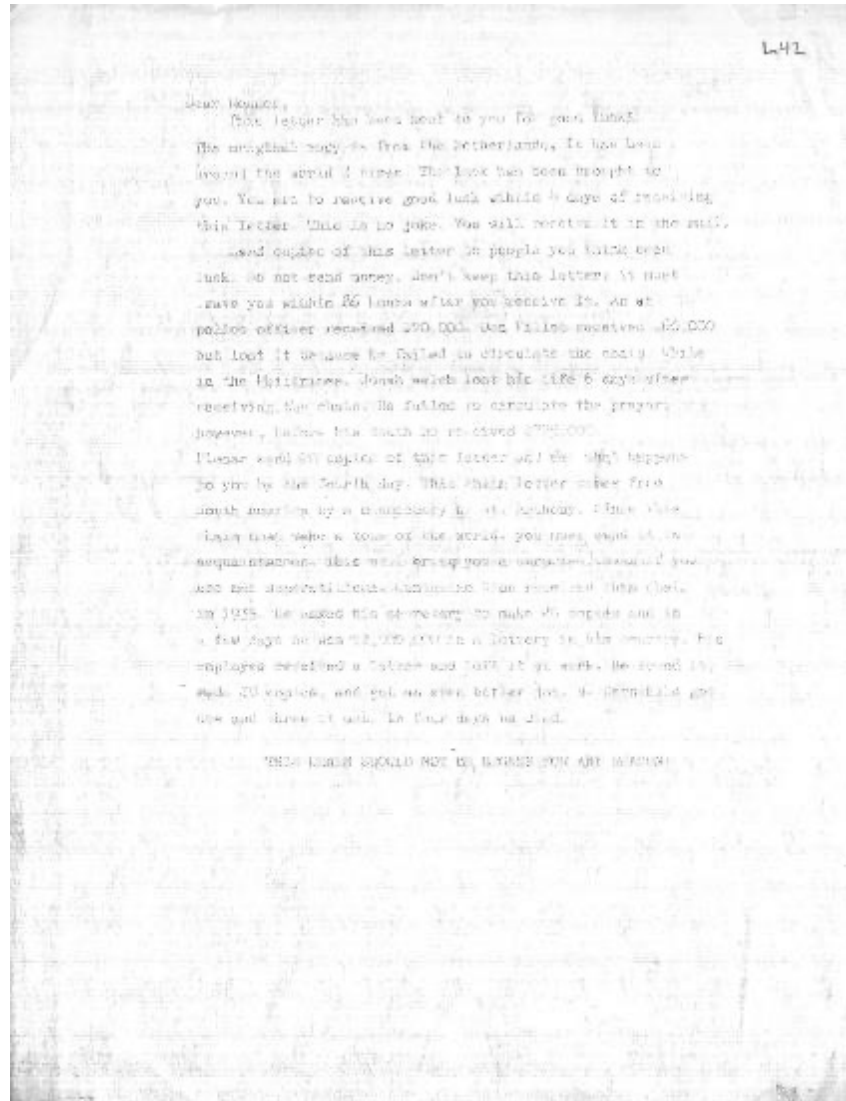
- These letters are different but appear to have the same origin.
- We were interested in reconstructing the evolutionary history of these chain letters.
- Because these chain letters are readable, they provide a perfect tool for classroom teaching of phylogeny methods and test for such methods.
- *Scientific American*: Jun. 2003

C. Bennett, M. Li, B. Ma: Chain Letters & Evolutionary Histories



# *An unclear letter reveals evolutionary path: ((copy)\*mutate)\**

---



# Why bother with chain letters?

---

<http://www.silcom.com/~barnowl/chain-letter/evolution.html>

- Like a virus, it has reached billions of people, literally.
- Like a gene, they are about 2000 characters;
- It even resembles some subtle phenomenon in biological evolution!

cause he broke the chain. While in the Philippines, Gene Walsh lost his wife six days after receiving the letter. He failed to circulate the letter. However, before her death he received \$7,755,000. Please

Life

his

WITH LOVE ALL THINGS ARE POSSIBLE

This paper has been sent to you for good luck. The original copy is in New England. It has been around the world nine times. The luck has now been sent to you. You will receive good luck within four days of receiving this letter, providing, you in turn send it on. This is no joke. You will receive it in the mail. Send copies to people you think need good luck. Don't send money as fate has no price. Do not keep this letter. It must leave your hands within 96 hours. An RAF officer received \$70,000. Joe Elliot received \$40,000 and lost it because he broke the chain. While in the Philippines, Gene Walsh lost his wife six days after receiving the letter. He failed to circulate the letter. However, before her death he received \$7,755,000. Please send 20 copies of the letter and see what happens in four days. The chain comes from Venezuela and was written by Saul Anthony Decroup, a missionary from South America. Since the copy must make a tour around the world, you must make 20 copies and send them to friends and associates. After a few days you will get a surprise. This is true even if you aren't superstitious. Do note the following: Constantion Dias received the chain in 1953. He asked his secretary to make 20 copies and send them out. A few days later he won the lottery of two million dollars. Carle Dadditt, an office employee, received the letter and forgot it had to leave his hands within 96 hours. He lost his job. Later, after finding the letter again, he mailed out the 20 copies. A few days later he got a better job. Dalan Fairchild received the letter and not believing, threw the letter away. Nine days later he died. Remember, send no money, and please don't ignore this.

IT WORKS

Coevolution

Life → wife

His → her

# Methods for constructing phylogeny

---

- Character Based Method
  - Parsimony Method: Find a tree topology so that the total number of mutations on the edges is the smallest.
  - Perfect Phylogeny.
  - Maximum Likelihood.
- Distance Based Method
  - UPGMA
  - Neighbour Joining
- Quartet Method
- Whole Genome Phylogeny.



# Character Based Method

---

- The first category of phylogeny methods do three things:
  - Define characters/features for each taxon.
  - Define a score function for each tree based on the characters.
  - Find the optimal tree



# Basics

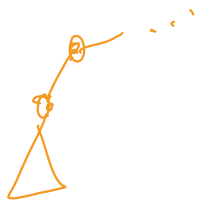
- A **character** is a “feature” in the species.
  - Vertebrate / invertebrate } *phenotype*
  - Has hooves / does not.
  - A letter in multiple sequence alignment. → *genotype*
  - The title is “Trust in lord ...” or “With love all things are possible”.
- An **evolutionary tree** is a rooted and leaf-labeled binary tree.



```
RLA0_METJA -----METKVKAHVAPWKIEEVKTLKGLIKSKPVAIVDMMDVPAPOLEIRDKIR-DKVKLRMSRNTLIIIRALKEAAEELN|
RLA0_PYRAB -----MAHVAEWKKKEVEELANLIKSYPPVIALVDVSSMPAYPLSQMRRLIRENGLLRVSRNLTIELAIKKAQELG|
RLA0_PYRHO -----MAHVAEWKKKEVEELAKLIKSYPPVIALVDVSSMPAYPLSQMRRLIRENGLLRVSRNLTIELAIKKAQELG|
RLA0_PYRFU -----MAHVAEWKKKEVEELANLIKSYPPVIALVDVSSMPAYPLSQMRRLIRENGLLRVSRNLTIELAIKKAQELG|
RLA0_PYRKO -----MAHVAEWKKKEVEELANLIKSYPPVIALVDVAGVPAYPLSKMRDKIR-GKAIIRVSRNTLIEIATKRAAQEIG|
```

# More details

- Input: An  $n \times m$  matrix of aligned characters. ( $n$  taxa,  $m$  characters)  
In the case of a multiple alignment, these are columns with no gaps.
- Output: A labeled tree with least number of mutations.

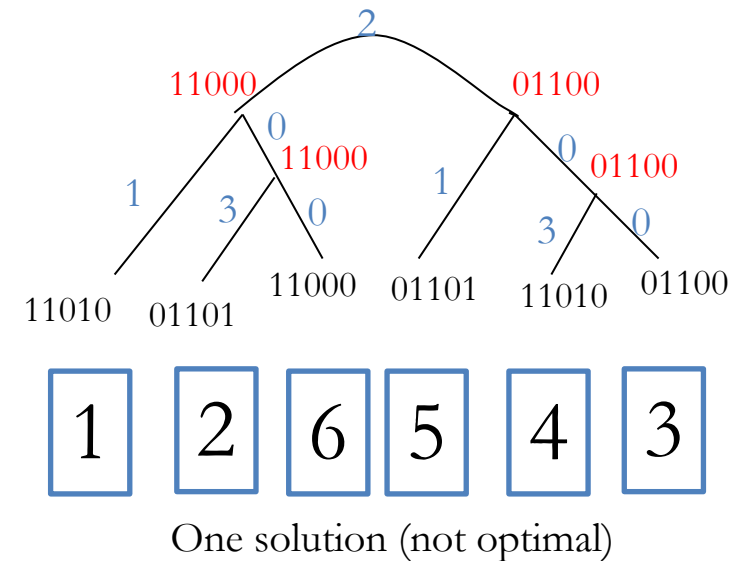


taxa

Characters (features)

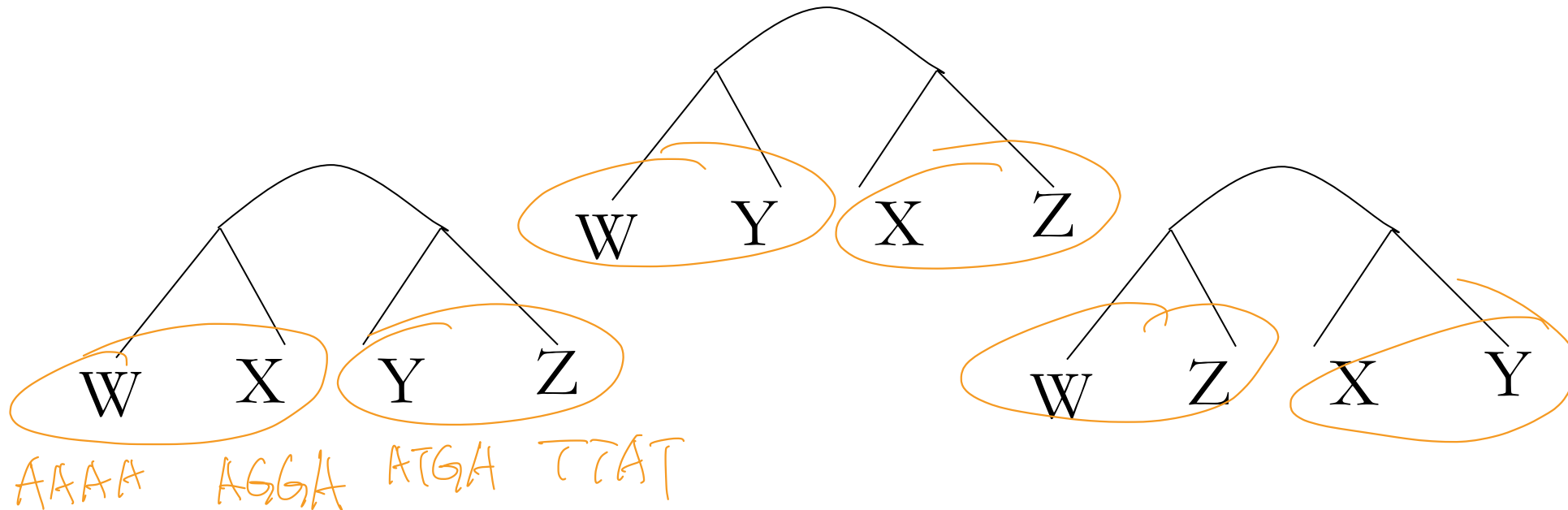
	a	b	c	d	e
1	1	1	0	1	0
2	0	1	1	0	1
3	0	1	1	0	0
4	1	1	0	1	0
5	0	1	1	0	1
6	1	1	0	0	0

Value of character e of taxon 2.



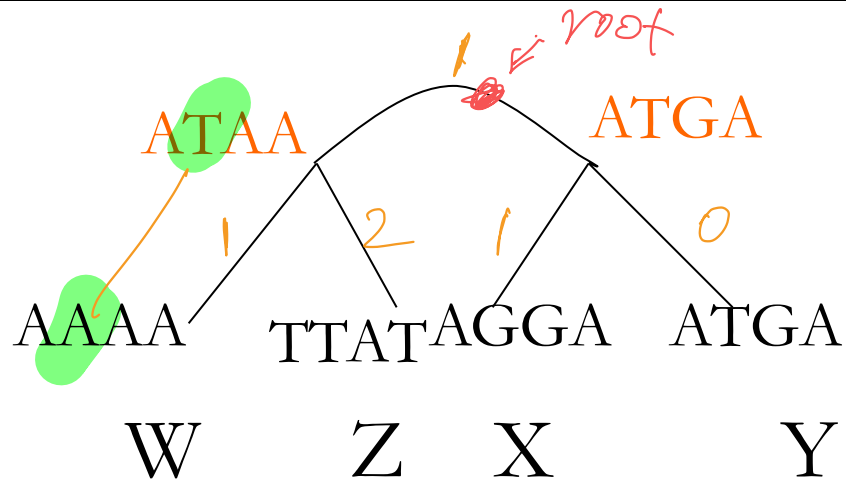
# Example:

- Suppose we have four taxa:
- W: AAAA
- X: AGGA
- Y: ATGA
- Z: TTAT
- There're only 3 possible (unrooted) trees on 4 taxa. Which has the least mutations?



# Parsimony example

---

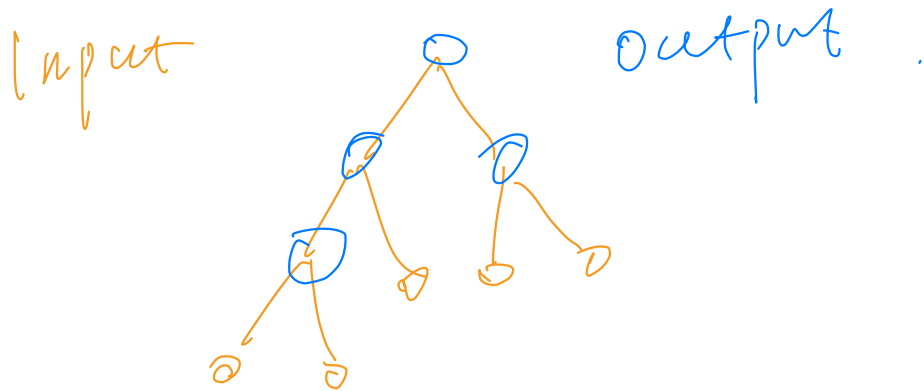


- In this case, we need 5 mutations. The other 2 require 6.
- So the “cheapest” tree joins W and Z on one side and X and Y on the other.
- Where is the root of the tree?

# Ancestor Reconstruction

---

- For a given topology, how to construct the ancestors? (in order to calculate the score of the tree)
- First observation: We can solve each column separately. So we can just solve for 1-character strings.
- Algorithm by Sankoff: tree-based dynamic programming.

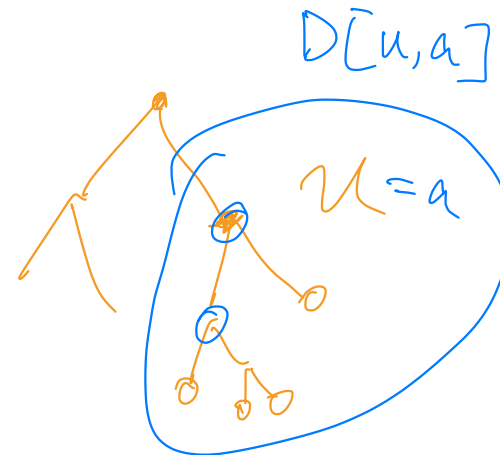


# Tree DP, details

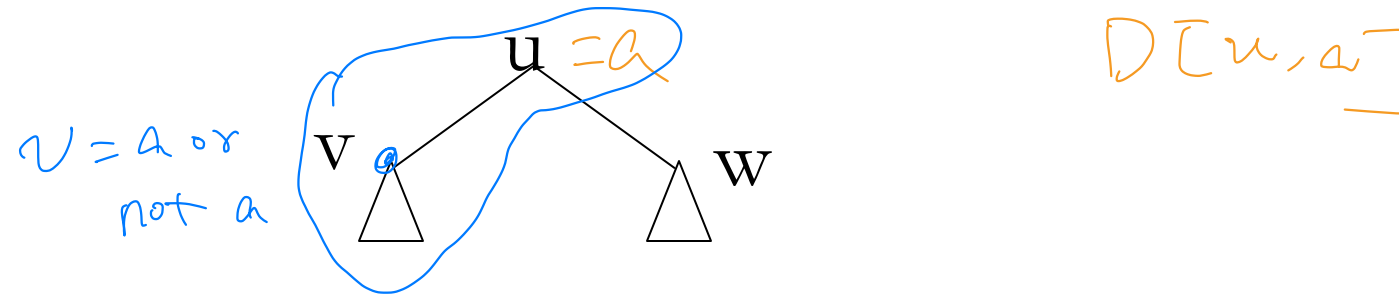
- For every node  $u$  of the tree and letter  $a$  of the alphabet  $\Sigma$ , let  $D[u, a] = \min \#$  of mutations in  $T_u$  if  $u$ 's label is  $a$ .
- Let  $r$  be the root. We want  $\min_x D[r, x]$ .
- For a leaf node  $v$ , if the character at leaf  $v$  is  $a$ , then  $D[v, a] = 0$ , and  $D[v, b] = 1$  for all other letters  $b$ .
- For an internal node  $u$ , with children  $v$  and  $w$ , suppose we know all of the values of  $D[v, *]$  and  $D[w, *]$ .
- How to compute  $D[u, *]$ ?

$v = a$

$$D[v, a] = 0, \quad D[v, b] = 1$$

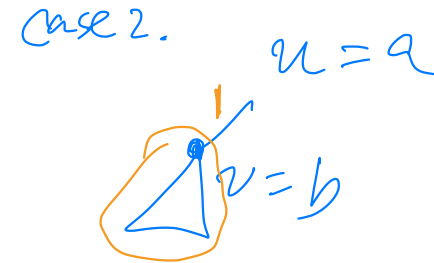


# Tree DP details (end)



- If we put letter “a” at node  $u$ , the cost of the left branch of the tree is the minimum of

- Case 1.  $D[v, a]$ , if  $v = a$
- Case 2.  $1 + \min_{b \neq a} D[v, b]$ , if  $v \neq a$ ,



- The same argument holds for the right branch. So

$$D[u, a] = \min(D[v, a], 1 + \min_{b \neq a} D[v, b]) + \min(D[w, a], 1 + \min_{b \neq a} D[w, b])$$

$\rightarrow = \min(D[v, a], 1 + \min_{b \in \Sigma} D[v, b])$

- Order of computation? Depth-first

- Time complexity?  $\left. \begin{array}{l} \text{for all } u \\ \text{for all } a \end{array} \right\} n \times |\Sigma|^2 = n \cdot \sigma^2$   $|\Sigma| = \sigma$



# Total runtime

- We can ignore  $b \neq a$  and minimize on all  $b$  without changing the value.
- Note:  $\min_b D[v,b]$  only needs to be computed once, not once every letter  $a$  for  $\min_{b \neq a} D[v,b]$ .
- If the tree is binary, and the size of the alphabet is  $\sigma$ , this algorithm takes  $O(n\sigma)$  time, since it's just  $O(\sigma)$  time at each node

$$D[u,a] = \frac{\min(D[v,a], 1 + \min_{b \neq a} D[v,b])}{\min(D[w,a], 1 + \min_{b \neq a} D[w,b])} = \frac{\min(D[v,a], 1 + \min_{b \in \Sigma} D[v,b])}{\dots}$$

for every  $u$   $\rightarrow$  let  $x = \min_{b \in \Sigma} D[v,b] : O(n \cdot \sigma)$   
 for every  $a$

apply the recurrence relation,

$$O(n \cdot \sigma) \rightarrow D[u,a] = \min(D[v,a], 1 + x) \dots$$

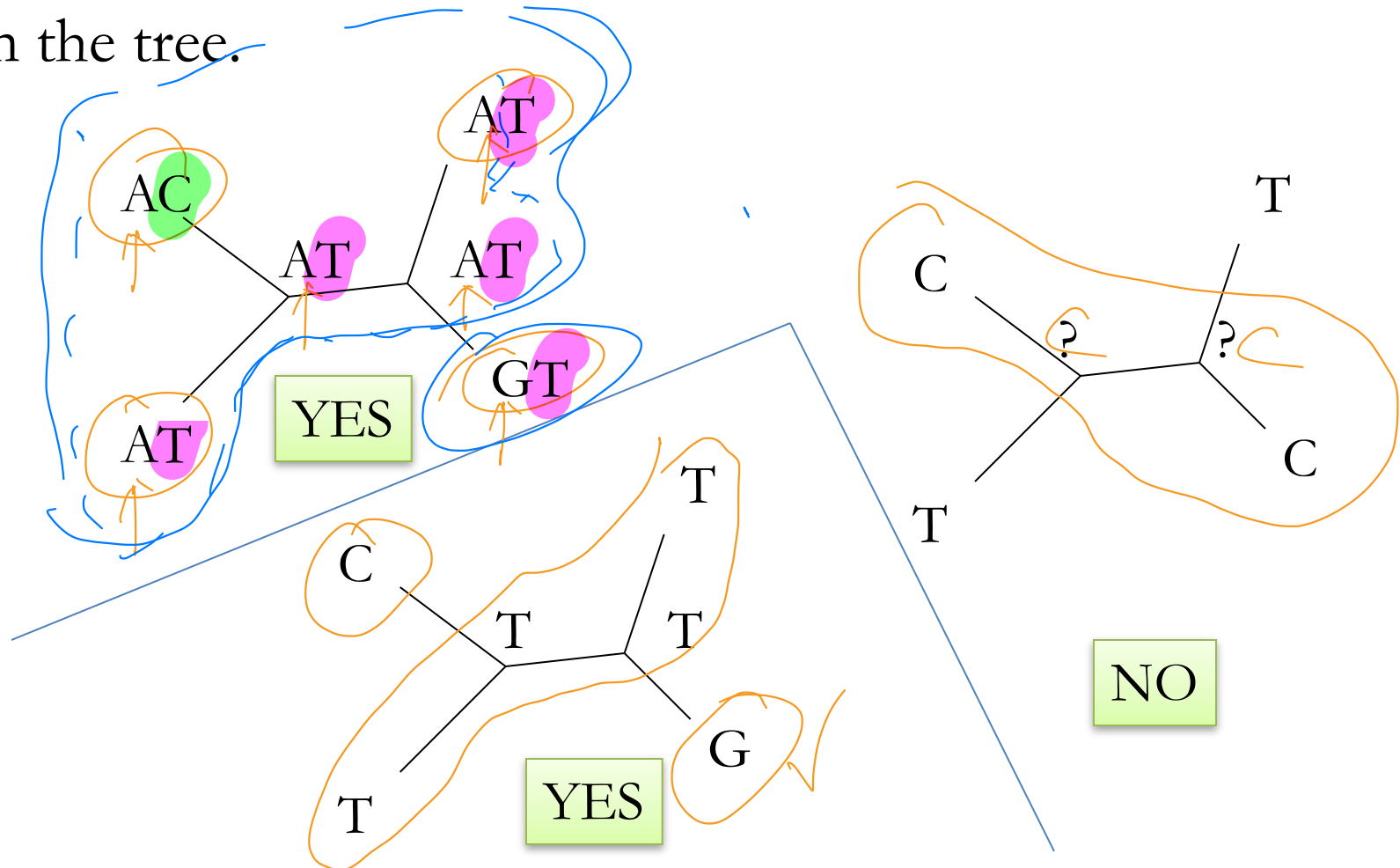
# Parsimony

---

- Parsimony Method: Find a tree topology so that the total number of mutations on the edges is the smallest.
  - NP-hard.
  - Algorithm: For each possible tree topology, uses DP to compute cost. Output the best tree.
- Suppose there are  $f(n)$  trees on  $n$  taxa.
- Total runtime:  $O(nm\sigma f(n))$ .
- Unfortunately,  $f(n) = 1*3*\dots*(2n-5)$ . (Roughly  $n!2^n$  or so).

# Perfect Phylogeny

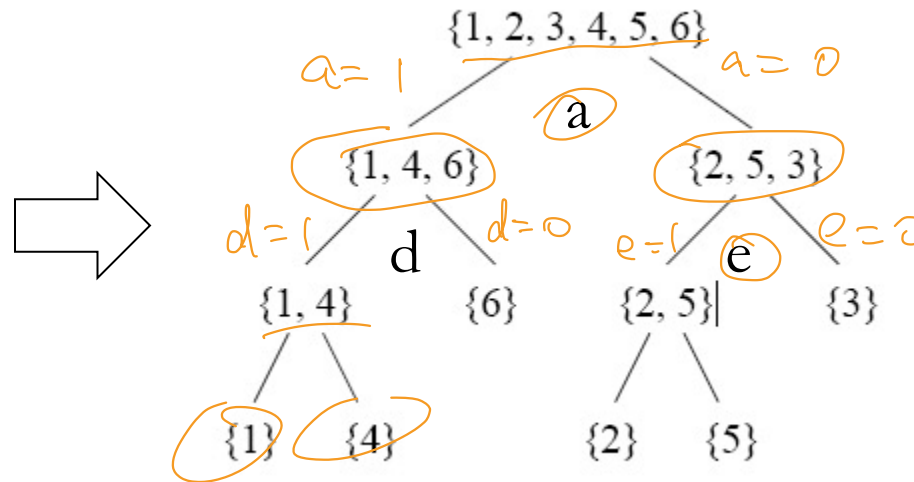
- A perfect phylogeny is such that for every character (every column), all species with the same state of that character is a connected component on the tree.



# Algorithm for Binary Case

- Algorithm 1: Start with a set of all taxa. Find a character and split the set into two. Recursion until each set has only one taxon.

	a	b	c	d	e
1	1	1	0	1	0
2	0	1	1	0	1
3	0	1	1	0	0
4	1	1	0	1	0
5	0	1	1	0	1
6	1	1	0	0	0



# Perfect Phylogeny

---

- For binary characters, Algorithm 1 is a polynomial time algorithm. If there is a perfect phylogeny, it outputs the perfect phylogeny.
  - Equivalently, if the output is not a perfect phylogeny, then there is no perfect phylogeny for the input.
- Theorem: If there is a perfect phylogeny for the input, and there are constant number of states for the characters, then a perfect phylogeny can be computed in polynomial time.
- r states, n taxa, m characters:  $O(2^{2r} nm^2)$ .

# Important Fact

---

- If a column has  $k$  different states, then
  - any phylogeny requires at least  $k-1$  mutations for the column.
  - a perfect phylogeny only has  $k-1$  mutations for the column.
- Conclusion: a perfect phylogeny is the best you can get for parsimony.
- Not all input matrix can cause a perfect phylogeny.

# Maximum Likelihood

---

- The score function used in parsimony (and perfect phylogeny) is too simple, especially when sequencing data become available.
- For example: the multiple alignment of proteins can be used as the input
  - Thousands of columns (characters).
  - Mutations between different pairs of amino acids have different rates.
  - Different substitution matrices on different columns.
- The maximum likelihood method aims to provide a better scoring function.



# A Multiple Alignment

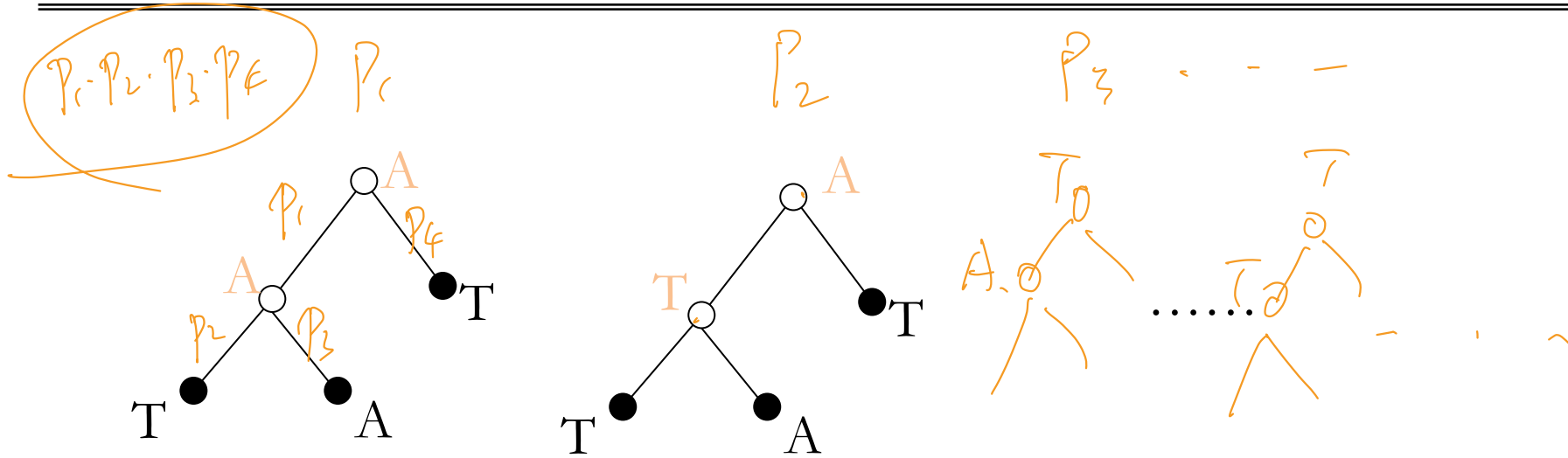
```

* . : * : : :
Q5E940_BOVIN -----MPREDRATWKSNYFLKIIQLLDDYPKCFIVGADNVGSKOMQQIRMSLRGK-AVVLGMGKNTMMRKAIRGHLENN--PALE 76
RLA0_HUMAN -----MPREDRATWKSNYFLKIIQLLDDYPKCFIVGADNVGSKOMQQIRMSLRGK-AVVLGMGKNTMMRKAIRGHLENN--PALE 76
RLA0_MOUSE -----MPREDRATWKSNYFLKIIQLLDDYPKCFIVGADNVGSKOMQQIRMSLRGK-AVVLGMGKNTMMRKAIRGHLENN--PALE 76
RLA0_RAT -----MPREDRATWKSNYFLKIIQLLDDYPKCFIVGADNVGSKOMQQIRMSLRGK-AVVLGMGKNTMMRKAIRGHLENN--PALE 76
RLA0_CHICK -----MPREDRATWKSNYFMKIIQLLDDYPKCFVVGADNVGSKOMQQIRMSLRGK-AVVLGMGKNTMMRKAIRGHLENN--PALE 76
RLA0_RANSY -----MPREDRATWKSNYFLKIIQLLDDYPKCFIVGADNVGSKOMQQIRMSLRGK-AVVLGMGKNTMMRKAIRGHLENN--SALE 76
Q7ZUG3_BRARE -----MPREDRATWKSNYFLKIIQLLDDYPKCFIVGADNVGSKOMQTIIRLSLRGK-AVVLGMGKNTMMRKAIRGHLENN--PALE 76
RLA0_ICTPU -----MPREDRATWKSNYFLKIIQLLNDYPKCFIVGADNVGSKOMQTIIRLSLRGK-AIVLMGKNTMMRKAIRGHLENN--PALE 76
RLA0_DROME -----MVENKAAWKAQYFIKVVELDFEFPKCFIVGADNVGSKOMQNIIRTSLRGL-AVVLGMGKNTMMRKAIRGHLENN--PQLE 76
RLA0_DICDI -----MSGAG-SKRKKLFIKATKLFTTYDKMIVAEADFGVSSQLQKIRKSIRGI-GAVLMGKNTMIRKVIIRDLADSK--PELD 75
Q54LP0_DICDI -----MSGAG-SKRKNVFIKATKLFTTYDKMIVAEADFGVSSQLQKIRKSIRGI-GAVLMGKNTMIRKVIIRDLADSK--PELD 75
RLA0_PLAF8 -----MAKLSKQKKQMYIEKLSSLIQQYSKILIVHVDNVGSNOMASVRKSLRGK-ATILMGKNTIRIRTALKKNLQAV--PQIE 76
RLA0_SULAC -----MIGLAVTTTTKIAKWKVDEVAELTEKCLKTHKTIIANIEGFPADKLHEIRKKLRGK-ADIKVTKNNLFIKALKNAG-----YDTK 79
RLA0_SULTO -----MRIMAVITQERKIAKWKIEEVKELEQKLRKYHTIIIANIEGFPADKLHDIRKKMRGM-AEIKVTKNTLFGIAAKNAG-----LDVS 80
RLA0_SULSO -----MKRLALALKQRKVASWKELEIKNSNTILIGNLEGFPADKLHEIRKKLRGK-ATIKVTKNTLFGIAAKNAG-----IDIE 80
RLA0_AERPE MSVVSIVGQMYKREKPIPEWKTLMLELEELFKSHRVVLFADLTGTPTFVVQVRVKKLWKK-YPMVAKKRIILRAMKAAGLE---LDDN 86
RLA0_PYRAE -MMLAIGKRRYVRTROYIPARKVKIVSEATELLQKYPYVFLFDLHGLSSRILHEYRYRLRRY-GVIKIIPKTLFKIAFTKVYGG---IPAE 85
RLA0_METAC -----MAEERHHTEHIPQWKKDEIENIKELIQSHKVFQGMVIEGILATKMKIRRDLDKDV-AVLKVSRTNLTERRALNQLG-----ETIP 78
RLA0_METMA -----MAEERHHTEHIPQWKKDEIENIKELIQSHKVFQGMVRIEILATKIQKIRRDLDKDV-AVLKVSRTNLTERRALNQLG-----ESIP 78
RLA0_ARCFU -----MAAVRGS---PPEYKVRAVEEIKRMISSKPVVAIVSFRNVPAGQMQKIRREFRGK-AEIKVVKNTLLEALDALG-----GDYL 75
RLA0_METKA MAVKAKGQPPSGYEPKVAEWRKREVKELMDEYENVGLVDLEGIPAPQLQEIIRAKLRERDTIIRMSRNTLMRIAEEKLDER--PELE 88
RLA0_METTH -----MAHVAEWKKKEVEQELHDLIKGYEVVGIANLADIPARQLQKMRQTLRDS-ALIRMSKKTLLISLALAKAGREL--ENVD 74
RLA0_METTL -----MITAESEHKIAPWKIEEVNKLKELKNGQIVALVDMMEVPAVQLQEIIRDKIR-GTMTLKMSRNTLIERAIKEVAEETGNPEFA 82
RLA0_METVA -----MIDAKSEHKIAPWKIEEVNALKKELKLSANVIALIDMMEVPAVQLQEIIRDKIR-DQMTLKMSRNTLIKRAVEEVAEETGNPEFA 82
RLA0_METJA -----METKVAHVAPWKIEEVKTLKGLIKSKPVVAIVDMMDVPAVQLQEIIRDKIR-DKVKLRMSRNTLIIIRALKEAAEELNPKLA 81
RLA0_PYRAB -----MAHVAEWKKKEVEELANLIKSYVPIALVDVSSMPAYPLSQMRRLIRENGLLRVSRTNLTIELAIKKAQELGKPELE 77
RLA0_PYRHO -----MAHVAEWKKKEVEELAKLIKSYVPIALVDVSSMPAYPLSQMRRLIRENGLLRVSRTNLTIELAIKKAQELGKPELE 77
RLA0_PYRFU -----MAHVAEWKKKEVEELANLIKSYVPIALVDVSSMPAYPLSQMRRLIRENGLLRVSRTNLTIELAIKKAQELGKPELE 77
RLA0_PYRKO -----MAHVAEWKKKEVEELANLIKSYVPIALVDVAGVPAVPLSKMRDKLR-GKALLRVSRTNLTIELAIKRAQELGQPELE 76
RLA0_HALMA -----MSAESERKTETIPEWKQEEVDAIVEMIESYESVGVVNIAGIPSRQLQDMRRDLHGT-AELRVSRTNLTERRALDDVD-----DGLE 79
RLA0_HALVO -----MSESEVRQTEVIPQWKREEVDLVDLIESYESVGVVGVAGIPSRQLQSMRRE LHGS-AAVRMSRNTLVNRRALDEVN-----DGFE 79
RLA0_HALSA -----MSAEEQRTTEEVPEWKQEQEVAELVDLLETYSVGVVNVGTGIPSKQLQDMRRGLHGT-AALRMSRNTLLVRALEEAG-----DGLD 79
RLA0_THEAC -----MKEVSQKKELVNEITRIKASRSVAIVDTAGIRTRQIQDIRGKNRGK-INLKVIKKTLLFKALENLGD-----EKLS 72
RLA0_THEVO -----MRKINPKKKEIVSELAQDITKSKAVAVDIKGVRTROMQDIRAKNRDK-VKIKVVKKTLLFKALDSIND-----EKLT 72
RLA0_PICTO -----MTEPAQWKIDFVKNLENEINSRKVAIVSISKGLRNNFQKIRNSIRDK-ARIKVSRRARLLRLAIENTGK-----NNIV 72
ruler 1.....10.....20.....30.....40.....50.....60.....70.....80.....90

```

Max likelihood method starts with a multiple alignment. Different columns may have different substitution frequency matrix.

# Maximum Likelihood



- For each possible tree topology  $T$ , for each possible internal node assignment, and calculate the probability based on the substitution matrix.
- For each tree  $T$ , add up all probabilities of all possible internal nodes. This is the likelihood of the input tree  $T$ . Figure shows a single column of the multiple alignment.
- Find the tree that maximizes the likelihood.

# More Notes about Maximum Likelihood

---

- Often more accurate than other methods.
- Very time consuming. Usually heuristic algorithms and dynamic programming algorithms are used to assist the search and estimation of the likelihood.
- If desired, one can also allow the change of the edge length (the mutation rate at each edge).
- Software available: e.g. PhyML.