

Midterm Notes

Molecular Biology

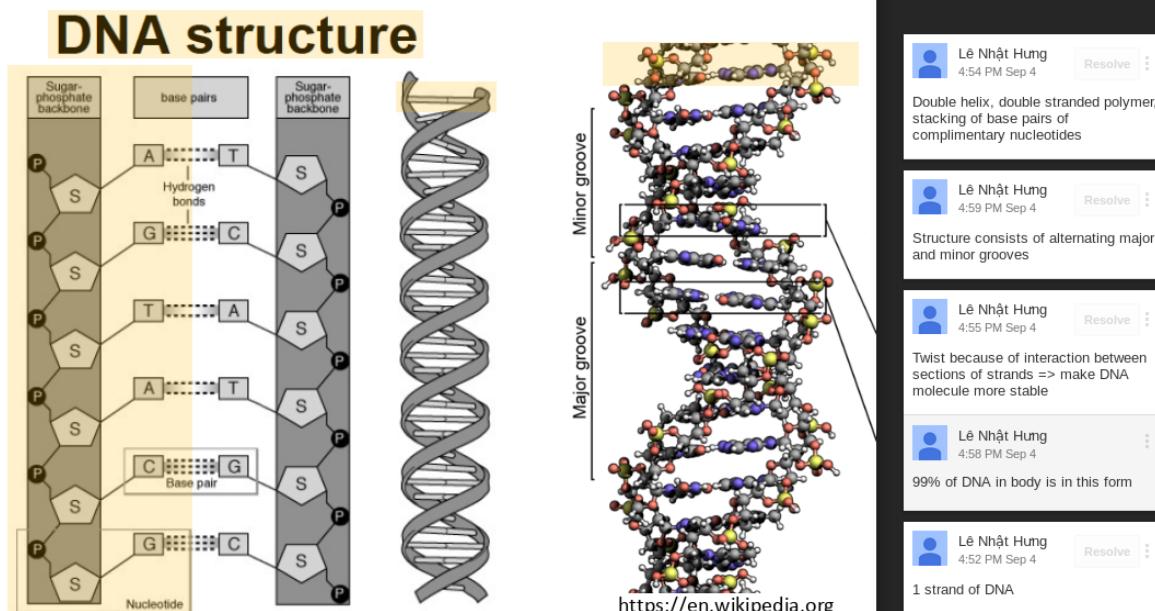
Cells:

- bags, fundamental working units of all living organisms
- each contains **genome** i.e. blueprint of all cellular structures and activities
- are of diff types (blood, skin, nerves) but all originate from **one cell: the fertilized egg**

Metazoa: mult cell organisms e.g. humans w trillions of cells

Protozoa: unicellular orgs with **no nucleus** e.g. bacteria

Eucaryote (e.g. humans) vs prokaryote: eukaryotic genes have introns, prokaryotic genes don't



DNA structure

- A **deoxyribonucleic acid** or **DNA** molecule is a double-stranded polymer composed of four basic molecular units called nucleotides.
- Each **nucleotide** comprises
 - a phosphate group;
 - a deoxyribose sugar;
 - one of four nitrogen bases:
 - purines: **adenine (A)** and **guanine (G)**,
 - pyrimidines: **cytosine (C)** and **thymine (T)**.

DNA structure

- Polynucleotide chains are **directional** molecules, with slightly different structures marking the two ends of the chains, the so-called **3' end** and **5' end**.
- The 3' and 5' notation refers to the numbering of carbon atoms in the sugar ring.
- The 3' end carries a sugar group and the 5' end carries a phosphate group.
- The two complementary strands of DNA are **antiparallel** (i.e., 5' end to 3' end directions for each strand are opposite)

Genomes

The *genome* of a cell is the entirety of its DNA content.

A genome is made of one or more *chromosomes*:
contiguous piece of double-stranded DNA

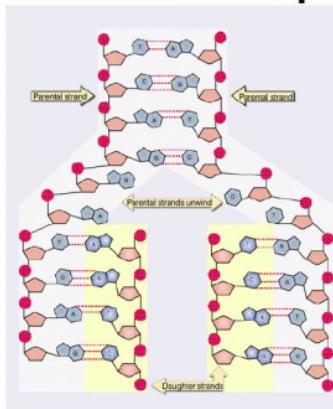
In bacteria (prokaryotes):

- One circular circular chromosome (2-10 Mb)
- Some small chromosomes called plasmids

In human (eukaryote):

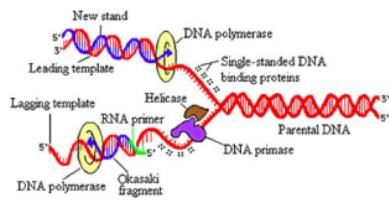
- 23 pairs of chromosomes = 22 autosomes pairs + 1 pair of sex chromosome (XX or XY)
- Each chromosome is 50 – 250 Mb
- Total genome size: 3,000,000,000 bp
- Total length of DNA in one nucleus: 2 meters!

DNA replication



Base pairing provides
the mechanism for
DNA replication.

DNA replication



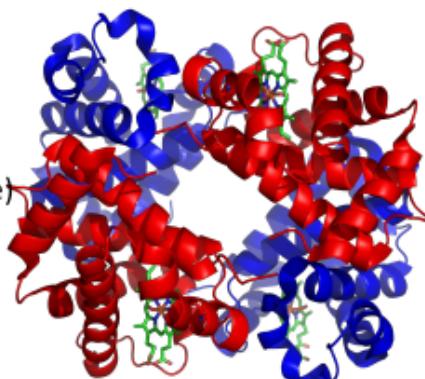
Collaboration of Proteins
at the Replication Fork

Useful video: <https://www.youtube.com/watch?v=TNK>

Proteins

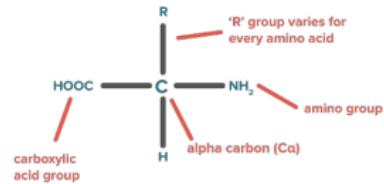
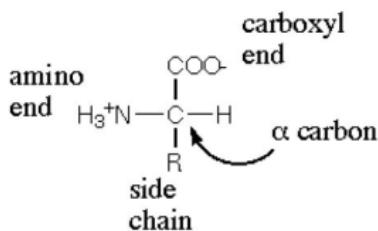
Proteins are molecules that perform a huge diversity of functions in the cell:

- Structure (actin, tubulin)
- DNA replication (DNA polymerase) + repairs
- DNA transcription (DNA transcriptase)
- Transport of small molecules (hemoglobin)
- Signaling (kinases)
- Regulation (transcription factors)
- Catalyze reactions (enzymes)
- Etc. etc.



Amino acids

Amino acids



- b) Name and describe briefly the three main biochemical steps that lead to the production of a protein from a eukaryotic gene.

1) *Transcription, which transcribes the DNA of a gene into pre-messenger RNA*

2) *Splicing, which removes introns, to obtain the mature messenger RNA.*

3) *Translation, which converts the mRNA to a protein sequence.*

- c) In protein-coding regions of genomes, why are insertions and deletions of 3 consecutive nucleotides more frequently observed than those of 1 or 2 nucleotides?

Because indels of size that is not a multiple of 3 result in frame shifts, so that all codons downstream are affected. This means that the amino acid encoded by the sequence downstream of the indel are likely to be completely changed. Indels of length 3 affect only one or two amino acids.

- d) If only 16 amino acids existed (instead of 20), what would be the minimal length of codons? Why? Watch out: that's a trick question!

3. You'd need 16 codons to encode the 16 amino acids, plus one for the stop codon. So you'd need ceiling ($\log_4(16+1)$) = 3 nucleotides per codon.

- e) Give one advantage and one disadvantage of using a Blast seed of size $w=11$ versus a blast seed of size $w=12$.

Advantage of w=11: Better sensitivity.

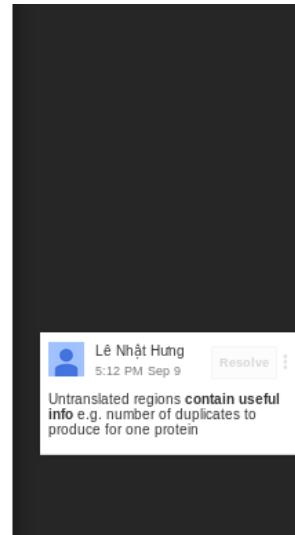
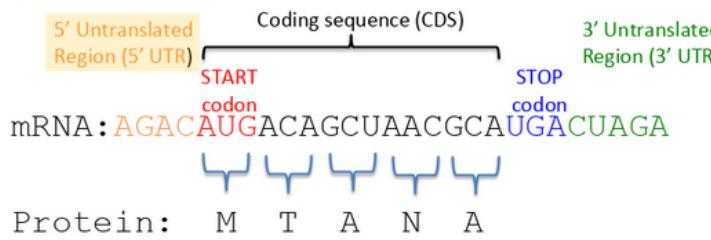
Disadvantage of w=11: Higher running time.

- f) What is the expected length of an open reading frame (ORF) in a random sequence where the nucleotides are chosen independently with probability $p_A = p_T = 1/3$, $p_C = p_G = 1/6$? Recall that the three stop codons are TAA, TAG, and TGA.

Probability of a stop codon = $1/3 \cdot 1/3 \cdot 1/3 + 1/3 \cdot 1/3 \cdot 1/6 + 1/3 \cdot 1/6 \cdot 1/3 = 2/27$. So the expected length of an ORF is $1 / (2/27) = 27/2 = 13.5$ codons = 40.5 nucleotides.

Translation

- 1) Ribosome searches for the first START codon (but there are many exceptions)
- 2) From there, non-overlapping triplets (codons) are translated to an amino acid
- 3) Until the ribosome encounters an in-frame STOP codon



Sequence Alignment

Goal: compare 2 biological seqs, to

- measure similarity
- highlight corresponding regions
- infer evolutionary history of 2 seqs derived from same ancestor

Mutations are random. Types of mutations:

1. Substitutions
2. Insertion
3. Deletion

Types of nucleotide **substitutions**

1. **Transitions:** happen more often

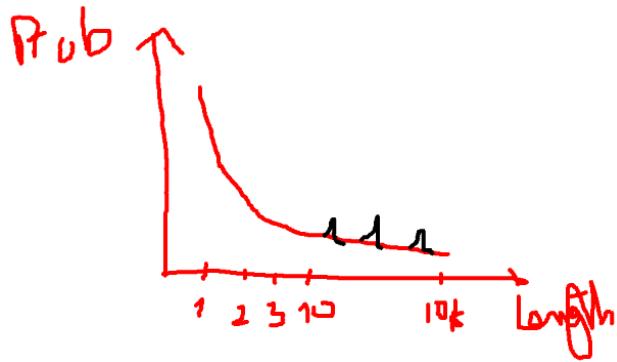
- A <-> G
- C <-> T

2. **Transversions**

- A → C,T
- C → A,G
- G → C,T
- T → A,G

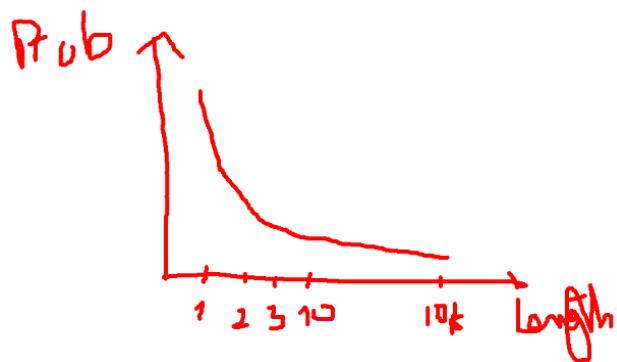
Insertion: insert n new nucleotides b/w 2 existing ones

Prob depends on insertion length



Spikes are due to **transposants**: organisms inserting nucleotides into seq

Deletion:



No peaks

Rates:

- transitions = $2.5 * \text{transversions} = 10$ del of size 1 ($N=1$) = 10 ins ($N=1$)

Pairwise Sequence Alignment

Global alignment: Needleman-Wunsch

$X(i, j)$ = score of optimal alignment of $s_1 \dots s_i$ vs $t_1 \dots t_j$

We want $X(m, n)$ = score of opt aln S vs T

Initialization: $X(i, 0) = b^* i$, $X(0, j) = b^* j$

Best aln ($s_1 \dots s_i, t_1 \dots t_j$) = $\max(\text{best}(s_1 \dots s_{i-1}, t_1 \dots t_{j-1}) + [s_i, t_j] = X(i-1, j-1) + \text{cost matrix } M(s_i, t_j),$

$$\text{best}(s_1 \dots s_{i-1}, t_1 \dots t_j) + [s_i, -] = X(i-1, j) + b,$$

$$\text{best}(s_1 \dots s_i, t_1 \dots t_{j-1}) + [-, t_j] = X(i, j-1) + b)$$

Filling out table

X_(m+1 x n+1) =	.	A	G	T	
.	0	-2	-4	-6	
A	-2	1	-1	-3	
C	-4	-1	0	-2	
A	-6	-3	-2	-1	
T	-8	-5	-2	-1	

Recover alignment **from right to left**

Arrow to left \Rightarrow col's nucleotide is aligned with gap

Arrow to top \Rightarrow row's nucleotide is aligned with gap

2 alignments:

```

Red: ACAT
      AG-T
Blue: ACAT
      A-GT
  
```

Running time: $O(n^*m)$ to fill table, $O(n+m)$ for traceback \Rightarrow total $O(n^*m)$

Space: $O(m^*n)$

Multiple Sequence Alignment (MSA)

Scoring a MSA: **sum-of-pairs score**

```

S1: AC-T
S2: AG-T
S3: CGAT
S4: CAAT

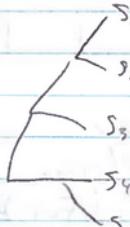
Ignore pairs of gaps
score(AC-T, AG-T) -> score(AC-T, AG-T)
  
```

Not making alignments, just scoring current alignments

$$\text{score} = \sum_{i,j} \text{score}(\text{Aln}(S_i, S_j))$$

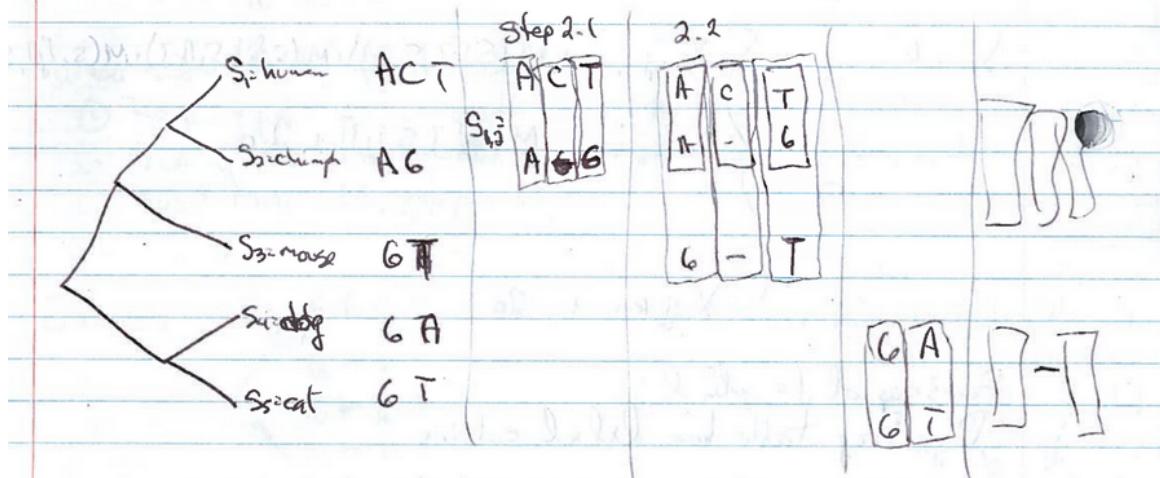
Progressive Sequence Alignment:

① Guess phylogenetic trees for $S_1 \dots S_n$



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

~ Find optimal alignment b/w pairs of seqs originating two children



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

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

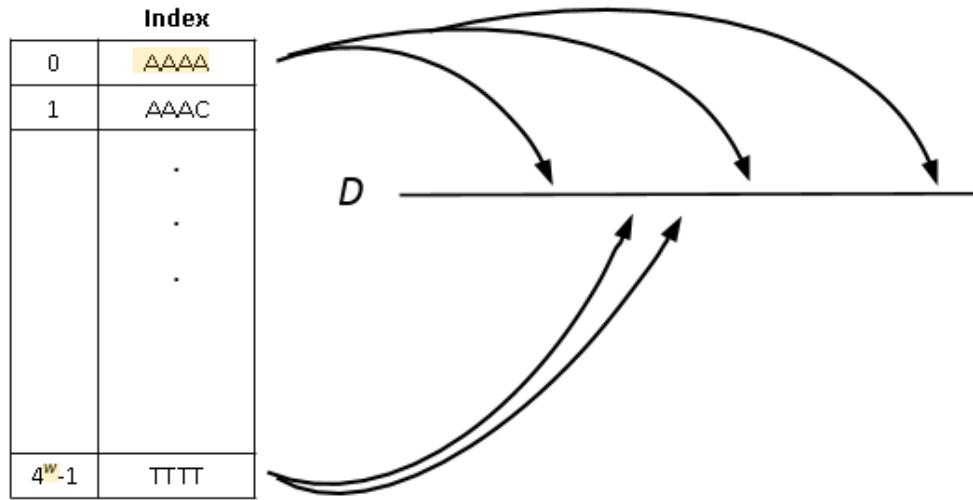
2.3 Align $S_{1,2,3}$ vs S_4 to obtain $S_{1,2,3,4}$

Fast Alignment Heuristics

Local alignment: BLAST

Query q and database D

1. Database indexing



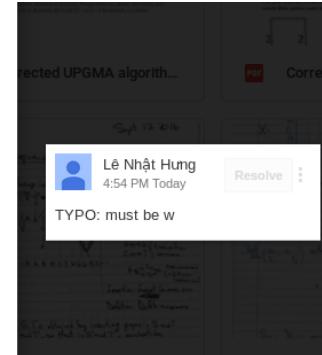
2. Scan for hit in D

- Given a query, q

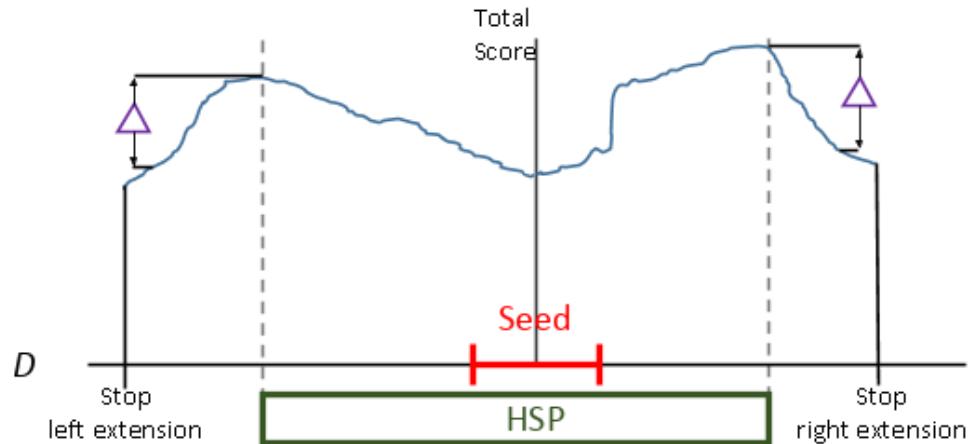
For each w -mer in q $O(|q|)$
 For each matches of q in D $O(w)$
 Investigate match further...

How many hits do we expect for a w -mer of size 11?

$$\frac{3 \times 10^9}{4^{11}} = 1000$$



3. Ungapped extension



Time? Linear in size of extension

4. If score of HSP in ungapped extension, feed HSP in gapped extension

Phylogenetic Inference

Distance based

Estimate distance matrix D:

Idea: Given seq. $S_1 \dots S_n$

- (A) Estimate distance matrix $D_{n \times n}$, where $D(i,j) = \text{distance b/w } S_i, S_j$
- (B) Find tree + Branch lengths such that $d_T(i,j) \approx D(i,j)$

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

Alg: UPGMA

Given: a distance matrix D with n species:

- 1) Initialize n clusters, C_1, \dots, C_n , each with a single species in it. Create a leaf node for each of the clusters.
- 2) Define the distance between two clusters as the average pairwise distance between members of the two clusters:

$$d(C_i, C_j) = \frac{\sum_{a \in C_i} \sum_{b \in C_j} D(a, b)}{|C_i| * |C_j|}$$

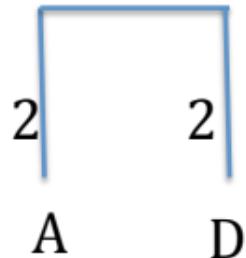
- 3) Repeat:
 - 3.1 Pick the two clusters C_i and C_j such that $d(C_i, C_j)$ is minimized.
 - 3.2 Create a new cluster $C_k = C_i \cup C_j$
 - 3.3 Create a new node in the tree, make it the parent of nodes i and j , at height $d(C_i, C_j)/2$.
 - 3.4 Add cluster C_k to the list of clusters, and remove clusters C_i and C_j .

Example: Consider the following distance matrix D:

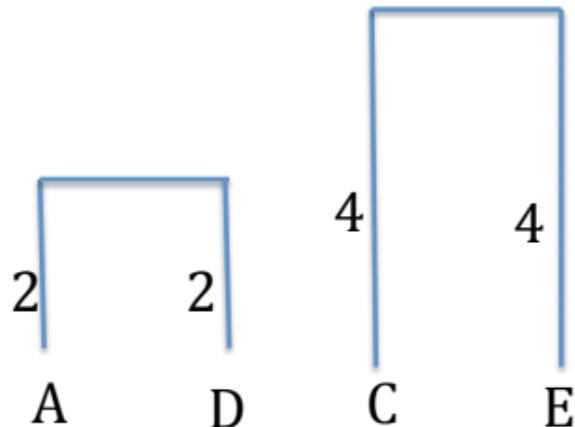
	A	B	C	D	E
A	-	16	16	4	16
B		-	10	16	10
C			-	16	8
D				-	16
E					-

First set $C_1 = \{A\}$, $C_2 = \{B\}$, $C_3 = \{C\}$, $C_4 = \{D\}$, $C_5 = \{E\}$.

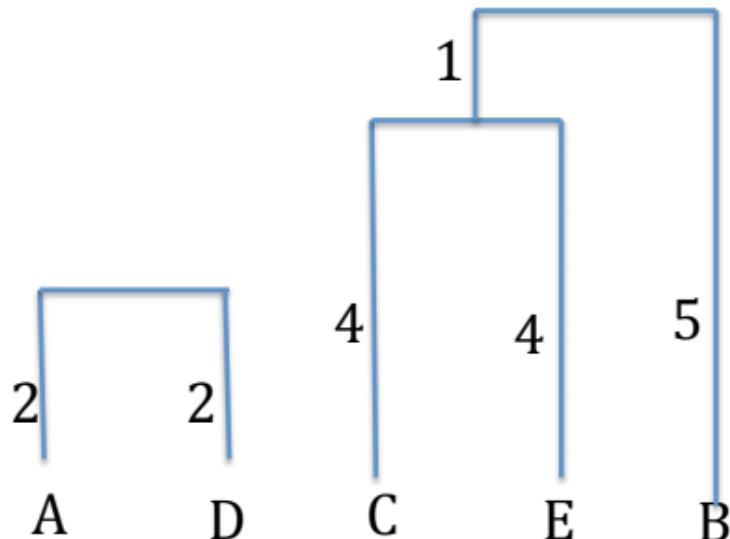
The pair with the smallest distance is (C_1, C_4) . Merge them to obtain $C_6 = \{A, D\}$ and create their parent node at distance $d(C_1, C_4)/2 = 4/2 = 2$ from each, to obtain:



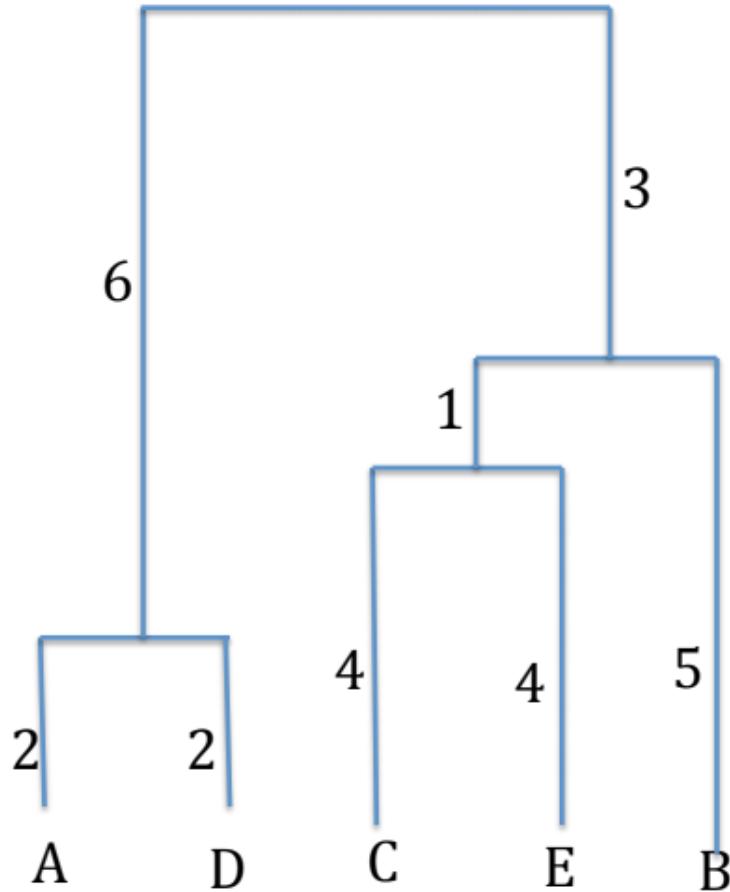
The next pair of clusters that is the closest is C_3 and C_5 , with $d(C_3, C_5) = 8$. Merge them to obtain $C_7 = \{C, E\}$ and create their parent node at distance $d(C_3, C_5)/2 = 8/2 = 4$ from each, to obtain:



The next pair of clusters that is the closest is $C_7 = \{C, E\}$ and $C_2 = \{B\}$, with $d(C_7, C_2) = (10 + 10) / 2 = 10$. Merge them to obtain $C_8 = \{C, E, B\}$ and create their parent node at height $d(C_7, C_2) / 2 = 10 / 2 = 5$, to obtain:



There are only two clusters C_6 and C_8 . Merge them and place their parent node at height $d(C_6, C_8) / 2 = ((16 + 16 + 16 + 16 + 16 + 16 + 16) / 6) / 2 = 8$, to obtain:



Parsimony based

Large Parsimony Problem

Given: n seqs S1, ..., Sn aligned to a MSA with columns with gaps removed

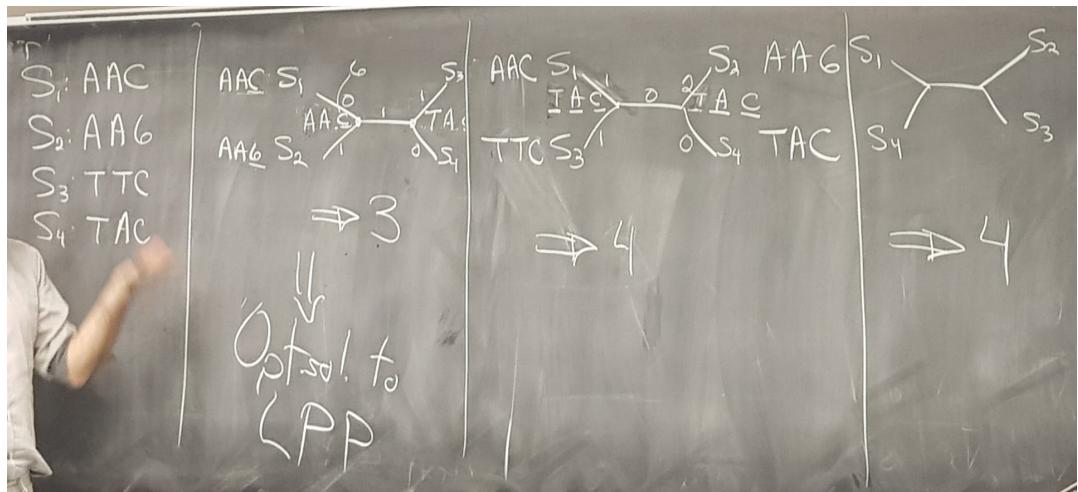
```

      x  x
S1: AC-TA      CTA
S2: -CCTA -> CTA
S3: ACCAA      CAA
  
```

Find:

- Unrooted tree T, whose leaves are S1, ..., Sn
- Sequences Sn+1, Sn+2, ..., Sn+(n-2) of all internal nodes of T, s.t.

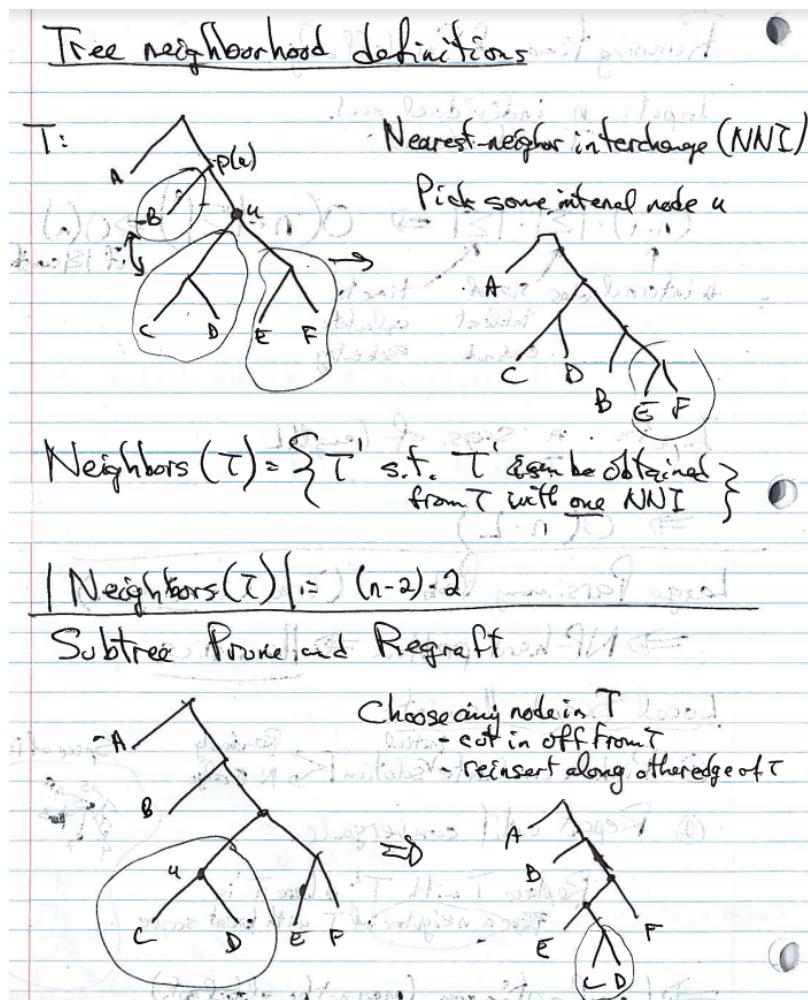
$$\begin{aligned}
 & \sum_{(i,j) \in E(T)} d(S_i, S_j) \text{ is minimized} \\
 d(S_i, S_j) &= \# \text{ subsets b/w } S_i, S_j \\
 & \quad E(T) \text{ edges of } T
 \end{aligned}$$



Greedy search using nearest neighbor interchange (NNI):

1. Choose tree T_0 randomly
2. Repeat...

$$T_{n+1} = \operatorname{argmin}\{\operatorname{parsScore}(T) : T \in \operatorname{NNI}(T_n)\}$$



Small Parsimony Problem

Sankoff Algorithm:

Sankoff Algorithm: Dynamic Programming

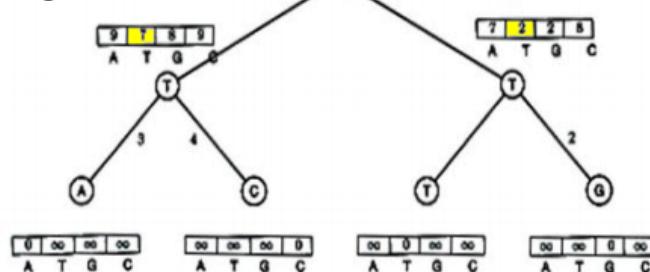
- Calculate and keep track of a score for every possible label at each vertex
 - $s_t(v)$ = minimum parsimony score of the **subtree** rooted at vertex v if v has character t
- The score at each vertex is based on scores of its children:
 - $s_t(\text{parent}) = \min_i \{s_i(\text{left child}) + \delta_{i,t}\} + \min_j \{s_j(\text{right child}) + \delta_{j,t}\}$

9 is derived from 7 + 2

So left child is T,

14	9	10	15
A	T	G	C

And right child is T



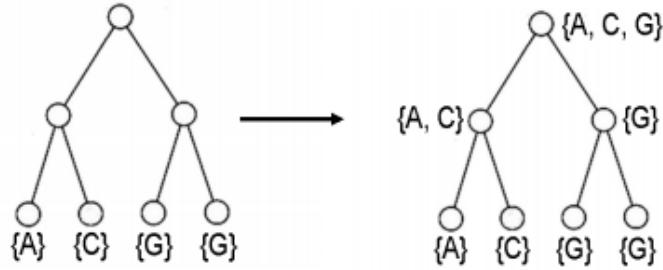
Pseudocode:

```

for each node u in tree (from leaves back to root):
    for each a in {A, C, G, T}:
        if u is leaf:
            Fu[a] = 0 if su == a else inf
        else:
            Fu[a] = min( Fv[a], min_{a != b}(Fv[b] + 1) )
                + min( Fw[a], min_{a != b}(Fw[b] + 1) )
    
```

Fitch's Algorithm:

As seen previously:



- Sankoff gives **same output** if have match=0 mismatch=1 scoring matrix

Proof:

For Sankoff, we have for each node u

$s_x(u)$ = minimum score of subtree rooted at node u if u had nucleotide x

We see that nucleotide x is optimal for node u if

$$s_x(u) = \min_i \{s_i(u)\}$$

Let

$$S_u = \text{set of optimal nucleotides for node } u$$

Then:

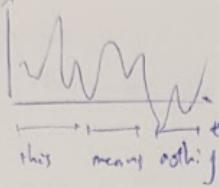
$$\begin{aligned} S_u &= S_v \cup S_w && \text{if } S_v \cap S_w = \emptyset \\ &= S_v \cap S_w && \text{otherwise} \\ &&& v, w \text{ child nodes of } u \end{aligned}$$

This is identical to Fitch's recurrence.

HHMs

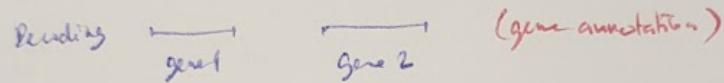
Goal: denote "seq of observations"

e.g.: Voice recog: x



Bioinformatics

Obs: $X = ACT \dots$



Obs: $X = EVLAZ \dots$

domain $\xrightarrow{\text{domain}}$ domain (protein domain annotation)

Prob: "Random walk in MFL"

Hear greeting every min $\{ b, h, v, q \}$
 ↑ ↑ ↑ ↑
 bonjour hello video namaste

$$= \sum (\text{alphabet of greetings})$$

$$= \{ \sigma_1, \sigma_2, \dots, \sigma_k \}$$

Obs: $X = x_1, \dots, x_n$ where $x_i \in \Sigma$

~~def~~

Know: \exists 3 types of neighborhood = states $= \{ F, E, C \}$
French lang claimed

Goal:

~~Given~~: X

Given: X

Find: most likely path $P = p_1, \dots, p_n$, where $p_i \in S$

Need to know:

1. Emission prob: $\Pr(x_i = \alpha | p_i = \beta)$

$$\mathbf{E} = \begin{bmatrix} b & h & n & a \\ \alpha \in \Sigma & \beta \in S \\ \begin{matrix} 0.6 & 0.3 & 0 & 0.1 \\ 0 & 0.9 & 0.1 & 0 \\ 0.3 & 0.3 & 0.3 & 0.1 \end{matrix} \end{matrix}$$

2. Transition probabilities: $\Pr(p_{i+1} = \alpha | p_i = \beta)$

$$\mathbf{T} = \begin{bmatrix} F & E & C \\ F & 0.8 & 0.2 \\ E & 0 & 0.9 & 0.1 \\ C & 0.2 & 0.3 & 0.5 \end{matrix} \quad \text{ES}$$

Assumptions:

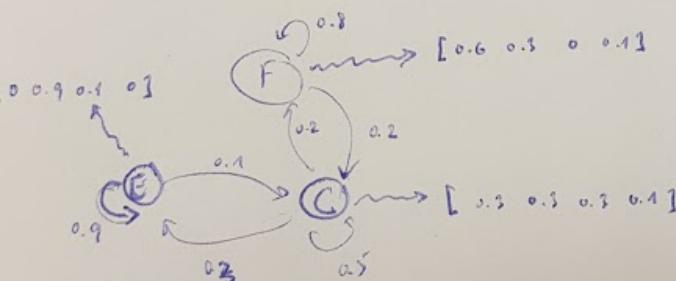
1. Markovian assumption:

Prob of state α from β doesn't depend on path to β

2. Obs are indep from each other, given state.

3. Initial state prob:

$$\Pr(p_1 = \alpha) \quad I = \begin{bmatrix} C & E & F \\ 0.1 & 0.5 & 0.4 \end{bmatrix} \quad \text{es}$$



HMM as generative model:

- generates path + seq of obs
- pick initial state p_1 randomly according to I
- Repeat
 - emit x_i from E given state p_i
 - transition to p_{i+1} given p_i

not a
finite
state
automata

$\hookrightarrow P: C C G E C F$
 $X: a n b h h n \dots$

Questions:

1. Maximum Likelihood path:
 Given: $\{x_i\}$ obs X
 $\{E, T, I, S, \text{ and } \Sigma\}$ of HMM
 Find: Path $P = p_1 \dots p_n$ s.t. $P(p_1 \dots p_n | X=x_1 \dots x_n)$ is max
 Alg: Viterbi alg.
 Viterbi alg.
2. Posterior decoding
 Given: $\{x_i\}$
 $\{H\}$
 $\{t_i\}$ of interest
 $\{s_i\}$ of interest
 Find: $\Pr(p_i = s_i | X=x_1 \dots x_n)$
3. Estimation problem
 Given: X
 HMM: S, Σ but not E, T, I
 Find: $\{E, T, I\}$ st $P(X=x_1 \dots x_n | E, T, I)$ is max

Maximum Likelihood Path

Given: $X = x_1 \dots x_n$

Find: $P = p_1 \dots p_n$ st $P(P|X)$ is max

Recall:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Therefore:

$$P(P|X) = \frac{P(X|P)P(P)}{P(X)}$$

1. From the independence assumptions:

$$P(X|P) = \prod_{i=0}^n P(x_i|p_i) = E(p_i, x_i)$$

2. Using independence and transition probabilities:

$$P(P) = P(p_1) \prod P(p_{i+1}|p_i) = I(P_1)T(P_i, P_{i+1})$$

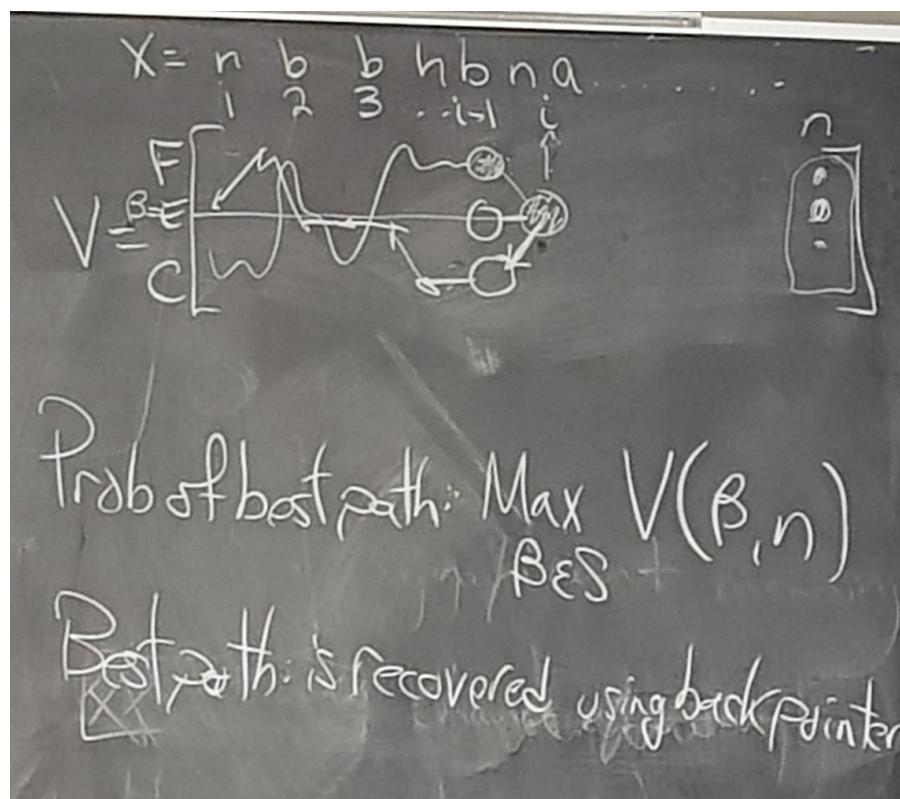
3. Denominator isn't needed: there's no P in it

Viterbi Algorithm

Define

$$V(\beta, i) = \text{prob of most likely path of len } i, \text{ given } X = x_1 \dots x_n, \text{ assuming } p_i = \beta \\ \beta \in S, i \in \{1 \dots n\}$$

$$V(\beta, i) = \max_{p_1 \dots p_i, p_i = \beta} \{P(P = p_1 \dots p_i, X = x_1 \dots x_i)\}$$



Fill out V:

For $i = 1$:

$$V(\beta, 1) = I(\beta)E(x_1, \beta) \quad \forall \beta \in S$$

For $i > 1$:

$$V(\beta, i) = E(\beta, x_i) \max_{\text{prev state } \delta} \{V(\delta, i-1)T(\delta, \beta)\}$$

Runtime:

$O(n|S|^2)$ time
 $n|S|$ cells, $|S|$ prev cells to find max for each cell