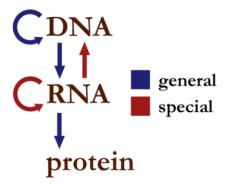
BMEG 3105: Data Analytics for Personalized Genomics and Precision Medicine

10/09/2025

Scriber: Geoffrey Li

Sequence and Dynamic Programming (L3)

A) Sequence Data



DNA store genetic information

RNA helps express information and turns to protein.

Phenotype = Genotype + Environment

Phenotype: How we look, for example skin colours.

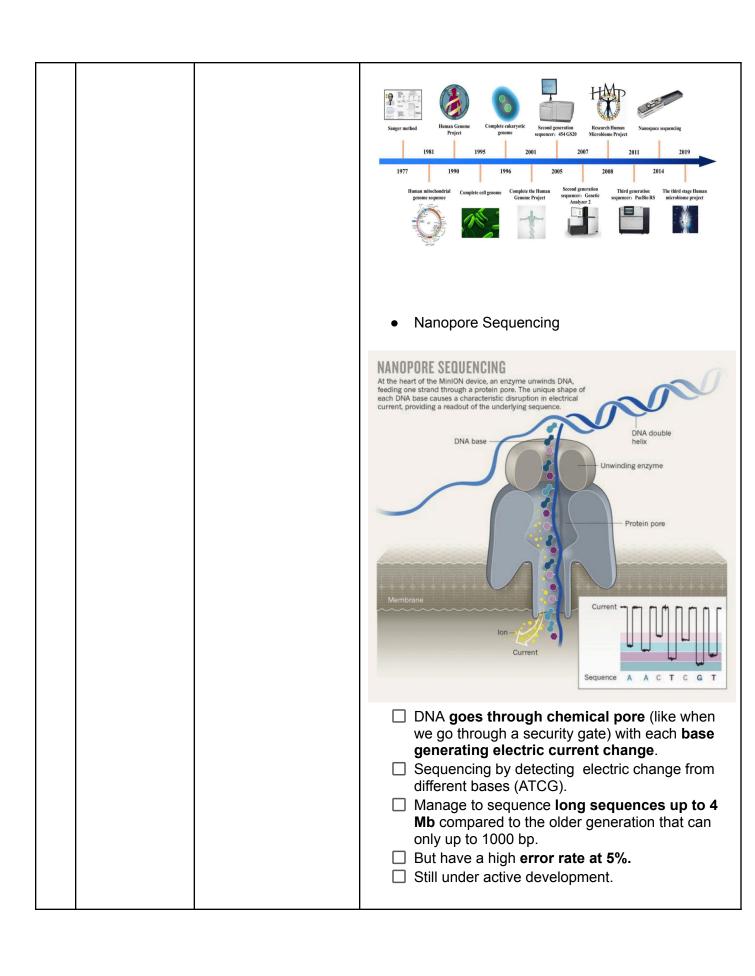
Genotype: DNA sequence / gene sequence. 99% dna sequence of each individuals is the

same

Environment: That affect the gene expression

1) Sequence Type

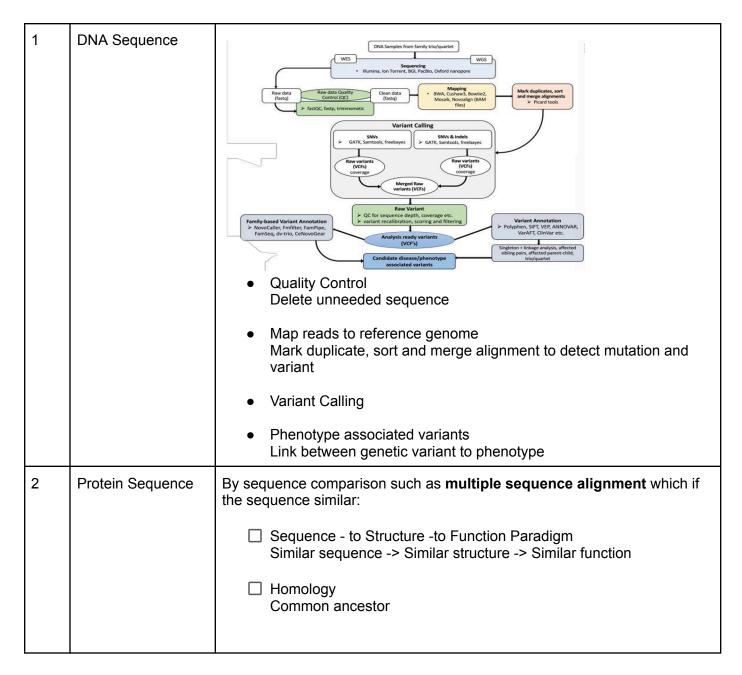
No.	Sequence		How to get the sequence
	DNA Sequence	 ATCG Double Strand There are 3 billions of this pair 	As years go, there are more advanced ways to achieve sequence and it still under development From short read to long reads
	RNA Sequence	AUCGSingle Strand	



Protein Sequence	20 amino acids Multiple sequence alignment	Protein Sequencing Step 1: Digestion Step 2: HPLC Step 3: Mass spectrometry Step 4: Computational sequence assembly Trends in Blochemical Sciences Based on Mass Spectrometry (MS) The process starts by breaking it down into short pieces then being determined by MS based on weight before being assembled into raw sequence (similar like a jigsaw puzzle).

2) Raw Data Handling

No. Sequence Process	
----------------------	--



3) Sequence Alignment

To determine the similarity between sequences and identify regions of similarity

Why is it important?:

- For biomolecular function and property prediction Similar sequence -> Similar structure -> Similar function
- For evolution, identifying conservative region, investing mechanism Common ancestor

Pairwise sequence alignment

Maximise the similarity of the two sequences by inserting the gap in the sequences and score the alignments.

__ ATCG

ATCG__ -> Even though its has the same formation sequence but the position of the gap in the arrangement makes them very different if we calculated it from the alignment score directly which ended up being wrong.

How to do the scoring?

Scoring Type	Example	
Match (identical bases)	A-AT-TC-CG-G	
Mismatch (Substitution)	A-TT-GC-Aetc	
Gap (Insertion / deletion)	• A • T • C • G	

	Α	С	G	Т
Α	2	-7	-5	-7
С	-7	2	-7	-5
G	-5	-7	2	-7
Т	-7	-5	-7	2

AGGCCG ->
$$2 + (-7) + 2 + 2 + (-10) + 2$$

ATGC_G = -9
AGGCCG -> $2 + (-7) + 2 + 2 + (-7) + (-10)$
ATGCG_ = -18

Higher the alignment score means the more similarity between the sequences.

How can we find the best alignment?

enumeration

The straight forward solution is to **enumerate down all** the possibilities then **combine the score alignment** then find the **highest**.

But the problem is there are too many alignments

$${\binom{2n}{n}} = \frac{(2n)!}{(n!)^2}$$
 If n = 300
There are 7 x 10^88 possibilities

Dynamic Programming
 The solution to enumeration problem

B) Dynamic Programming

- Break main problems into sub-problems
- Solve the sub-problems optimally and recursively
- Utilize these optimal solutions to build the best overall solution for the initial problem

Step:

Scoring matrix:

	Α	С	G	Т
Α	2	-7	-5	-7
С	-7	2	-7	-5
G	-5	-7	2	-7
Т	-7	-5	-7	2

	Α	С	С	G
Α				
С				
G				

1) Fill the **first row and and first column** which is the **gap penalties** and in this example is (-10)

		Α	С	С	G
	0	-10	-20	-30	-40
Α	-10				
С	-20				
G	-30				

As you go toward the left or down side of the gap penalties , add the gap penalties from the previous box number.

2) Next we start to fill the box on the diagonal of the box (0)

0	-10
-10	

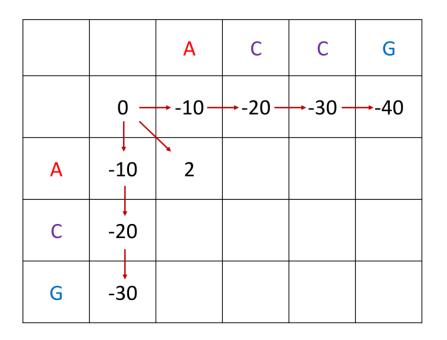
To fill the box, there are 3 options to be considered, **top**, **left**, **and diagonal**.

Top Left

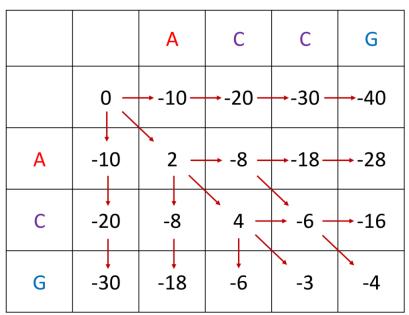
Diagonal

Previous number +match/mismatch (Since in this case is A-A so is a match) 0 + 2 2

From the results, we got -20, -20 and 2. Here we **choose the biggest number** to be put in the box which is 2. Also **put the arrow** to indicate which side the number comes from.



3) Do it to the other empty box until it fills the whole table.



4) Translate the arrow into alignment starting from the bottom right. It may have more than 1 alignment. Horizontal arrow translates to a gap while diagonal arrow means both sequences proceeding by one letter.

In this case there are 2 alignment