**<u>Lecture agenda</u>**:

- Recap of last lecture
- Discussion of further detail in dynamic programming (DP)
- Understand the function and application of gene expression matrix
- Introduction to sequence assembly and sequence mapping (Notes are included in slides. but due to limited time, it did not mention in this lecture)

**<u>Expected outcomes</u>**:

- Understand the principles of DP and able to apply it to sequence alignment
- Use DP to trace back the optimal alignment
- Able to transition from sequence data to data matrix

**<u>Feedback and comments from last lecture</u>**:

- Positive feedback:
  - Clear and easy to follow
  - Liked the calculation part
- Request and suggestion:
  - Slower explanation of calculations.
  - More examples needed
  - Use of animations for complex graphs

**<u>Recap</u>**:

***Fundamental understanding of Dynamic Programming***:

- Purpose:
  - it solves alignment problems by breaking them into smaller problems
    *[e.g. when we calculate F(4,3), we can break it into F(4,2) + F(3,3) + F(3,2). ]*
  - Find the optimal alignment score to determine the optimal alignment
- Matches of each base:
  - Finite choice for each base
    a. Align to another base
    b. Align to a gap

**<u>Lecture</u>**:

*Dynamic Programming (DP) in Sequence Alignment*:

- Further analysis/application of DP
  - Merge the result of small questions to fix the final problem
  - Arrows show the alignment pathway/arrangements
    a. <u>Trace back</u> from the right bottom conner to left upper conner
    b. Determine the optimal alignment <u>by reversing the alignment pathway</u> (the optimal alignments can be various)
  - Calculate the alignment score
    a. Directly <u>observe the base pair</u>
    b. According to the scoring matrix, add the <u>score together</u> corresponding to base pair (scoring matrix can be various)
    c. Optimal alignment score should be equal to the score on the right lower conner on DP
  - Sequence alignment can be used to identify sequence similarity
- Invent DP and DP process
  - Fill in the table according to the scoring matrix
  - <u>Preserve the arrows</u>
  - The value in <u>the last cell</u> is the <u>best alignment score</u>
  - <u>Trace back</u> the arrows to get the alignment.
- Further information provided by DP table
  - DP table stores answer of <u>sub-problems</u> and the <u>construction path</u>
    *[e.g. from DP which solve the problem of F(5,4), it contains the answer of F(3,3), F(4,2), etc..]*
- Concept of local alignment
  - Similar components, motifs and domains, in dissimilar sequences
  - Only care the local information between two sequencing; care the most important and the number in the cell

*Computational Analysis*:

- Using two sequences and scoring matrix
- Provide straight forward solution
- If using DP will be too much calculation
- Webserver for sequence alignment is provided at supporting link section

*Scoring Matrix*:

- Mismatch causes by mutation
- Insertion/deletion, or gene duplication due to additional insertion during transcription may cause a gap

- Scoring matrices are various
  - Different databases can build <u>different scoring matrices</u>
  - Different scoring matrices can <u>aim different needs</u>
  - There are <u>different types of matrices</u> including specific for DNA, RNA or protein
    *[e.g. Blocks Substitution Matrix (BLOSUM) is a protein scoring matrix]*
- It depends on how we define the similarity between two sequences


*Data Sequencing*:

- Purpose:
  - Reveal the genetic information which hidden in DNA sequences
- Since human genome is mostly the same, sequencing alignment can help find the differences
  - Gene expression difference is important for studying the phenotype difference


*Gene Expression Matrix*:

- Purpose:
  - Visualize the difference between different gene expressions across sample or environment, etc..
- Principle:
  - The amount of protein that is translated by a specific gene can reveal the gene expression level
  - Since the protein is hard to count, check the <u>counts of RNA copies</u> can also determine the <u>gene expression</u> due to central dogma
- Processing of building gene expression matrix:
  - <u>Map</u> the short read to the genome
  - <u>Count the number of reads</u>, which is content of gene expression matrix


*Potential Project – 1*

- A pipeline to get the gene expression matrix form reaeds
  - Find the genome
  - Find the reads
  - Map reads to reference genome
  - Count reads for each gene
  - Use Google to find the software and the data
  - Explain each step in the report to let us know you understand what you are doing

**Next lecture topic**:

- Where to find/ how to get reference genome
- How do we do genome assembly and mapping
- Data exploration and data cleaning

**Supporting Links**:

- Webserver for sequence alignment: https://www.ebi.ac.uk/Tools/psa/emboss_needle/
- Biopython: https://biopython.org/
- Post-lecture survey: https://forms.gle/4AyB35ztD7QWPdDv8

**Resource and related uncovered topics**:

- Bioinformatics: Sequence and Genome Analysis---Chapter 2&3
- Time complexity and space complexity analysis
- Local alignment
- Multiple sequence alignment
- Affine gap penalty
- Sequence database search: BLAST