# FASTA: SIMILARITY SEARCHING AND ITS APPLICATIONS

# WHAT IS FASTA?

FASTA is a database similarity search tool which uses a standard format for sequence data of DNA and proteins. First developed by Lipman and Pearson, it was used to compare protein sequences against protein databases. But now it is used to compare both DNA and protein sequences against various databases.

FASTA uses a "hashing" strategy to find matches for a short stretch of identical residues with a length of k. Typically, a k-tuple is composed of two residues for protein sequences and six residues for DNA sequences.
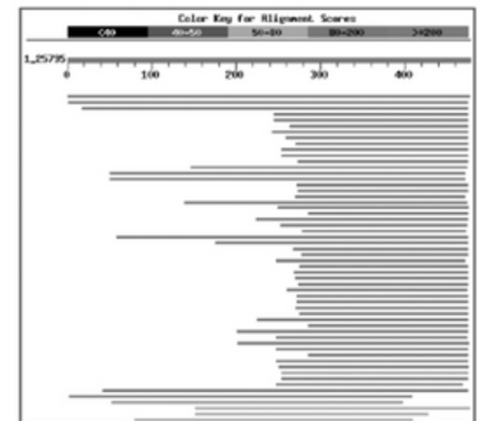
- **FASTA** – compares DNA/protein sequence against DNA/protein database respectively.
- **SSEARCH** – performs protein-protein or DNA-DNA comparisons using local alignment algorithm.
- **GGSEARCH/GLSEARC**H – works using global alignment (GGSEARCH) or a combination of global-local alignment (GLSEARCH) to compare protein and nucleotide sequence.
- **FASTX/FASTY** – compares a DNA sequences to protein database by translating the DNA sequence into 3 frames and allowing gaps and frameshifts.
- **TFASTX/TFASTY** – compares a protein sequence to a DNA database by translating the DNA sequence into 6 frames, 3 in the forward direction and 3 in the reverse direction.
- **FASTF/TFASTF** – compares a mixed peptide sequence against a protein database (FASTF) or translated DNA database (TFASTF).

FASTA works by comparing the query sequence to a database of sequences to identify similar matches. It uses heuristic algorithms to perform the searches and identify significant matches based on statistical parameters.

The mechanism involves 4 steps: identifying high similarity regions, re-scoring of the best aligned sequences, joining threshold to remove unlikely segments and final alignment of the new sequence.
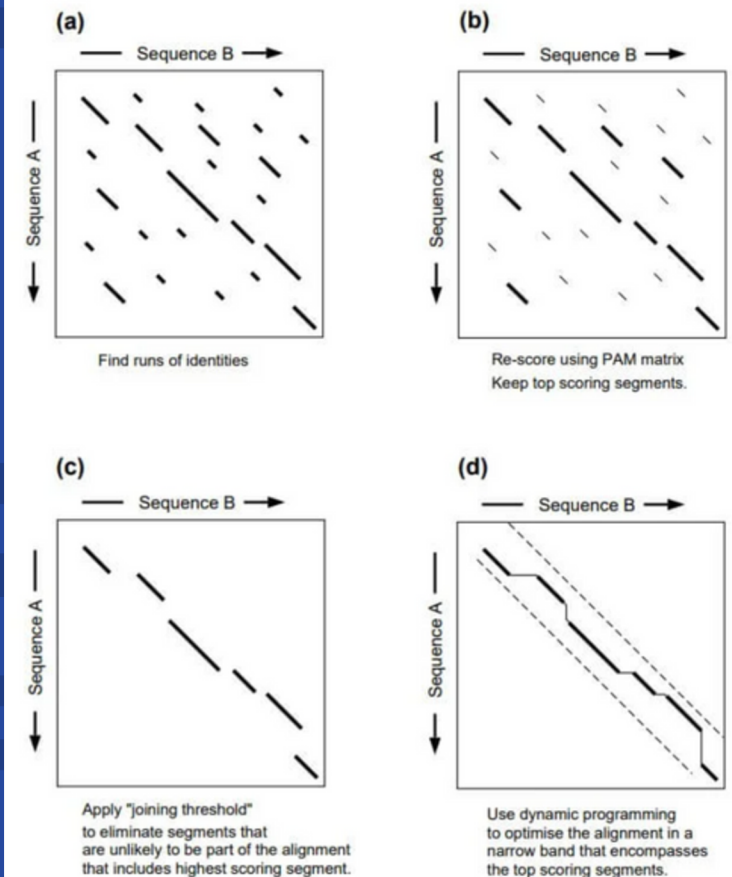
# HOW IT WORKS?



Graphical overview

Matching list

# 1. IDENTIFYING SIMILAR REGIONS

- FASTA identifies regions with high similarity by creating a lookup table by hashing method. The query is broken down to k-tuples.
- K-tuple values are increased to reduce the number of background hits, so it focuses more on the significant hits. K-tuple is 2 for proteins and 6 for nucleotides.
- The similar regions are plotted in a 2D matrix as diagonals and the top scoring diagonals are saved which are having the highest similarity.



**FASTA Algorithm**

(a) Sequence B → / Sequence A ↓
Find runs of identities

(b) Sequence B → / Sequence A ↓
Re-score using PAM matrix
Keep top scoring segments.

(c) Sequence B → / Sequence A ↓
Apply "joining threshold"
to eliminate segments that
are unlikely to be part of the alignment
that includes highest scoring segment.

(d) Sequence B → / Sequence A ↓
Use dynamic programming
to optimise the alignment in a
narrow band that encompasses
the top scoring segments.

# 2.RE-SCORING

- The 10 best diagonals are re-scored using scoring matrices – BLOSUM50 for proteins and identity matrices for DNA. A sub region with the highest score is identified for each of the diagonals which are called initial regions.

# 3.JOINING THRESHOLD

- A joining threshold is applied that excludes segments unlikely to be part of the final alignment. The selected regions with initial scores above the pre-set threshold are joined. This introduces gaps between diagonals while applying gap penalties. The score of the gapped alignment is calculated by subtracting a penalty for each gap, which is used to rank the database sequence by similarity.

Finally, the alignment is refined to produce the final alignment. This is done by using the banded Smith-Waterman algorithm, which is a dynamic programming algorithm that calculates the optimal score (opt) for alignment. This score is used for statistical calculations when measuring e-value, bit scores and finally, functional domains are studied and a phylogenetic analysis is performed for the given species.

1. Given two amino acid sequences for comparision:

sequence 1    **AMPSDGL**
sequence 2    **GPSDNAT**

2. Construct a hashing table:

| amino acid | sequence position | | offset |
| --- | --- | --- | --- |
| | seq 1 | seq 2 | |
| A | 1 | 6 | -5 |
| D | 5 | 4 | 1 |
| G | 6 | 1 | 5 |
| L | 7 | – | – |
| M | 2 | – | – |
| N | – | 5 | – |
| P | 3 | 2 | 1 |
| S | 4 | 3 | 1 |
| T | – | 7 | – |

3. Identify residues with the same offset values (highlighted in grey).

4. Find the matching word of three residues in the order of 3, 4 and 5 in one sequence and 2, 3,and 4 in the other.

5. This allows establishment of alignment between the two sequences.

sequence 1    **AMPSDGL-**
            **|||**
sequence 2    **-GPSDNAT**

# STATISTICAL SIGNIFICANCE OF FASTA

- FASTA provides an estimate of statistical significance of each alignment found, which is evaluated using E-value (the likelihood of obtaining a sequence alignment score by chance). Smaller the E-value, more significant is the alignment.
- FASTA also uses bit scores and similarity scores based on the scoring matrix and gap penalties, to evaluate the significance of sequence alignments.
- Z-score is another parameter that represents the number of standard deviations from the mean score of the database search. A higher z-score indicates a higher similarity match.

# APPLICATIONS OF FASTA

## SIMILARITY SEARCHING

Used to identify similar regions in protein and DNA sequences to understand conserved domains or motifs.

## FUNCTIONAL ANNONATION

Used to search database of sequences to identify homologous sequences to predict the function of a newly identified sequence.

## PHYLOGENETIC ANALYSIS

Multiple sequence alignment can be done to plot phylogenetic trees by identifying evolutionary relationships between species.