

# Comparison of sequences

---



INDRAPRASTHA INSTITUTE *of*  
INFORMATION TECHNOLOGY **DELHI**

**Dr. Jaspreet Kaur Dhanjal**

**Assistant Professor, Center for Computational Biology**

Email ID: [jaspreet@iiitd.ac.in](mailto:jaspreet@iiitd.ac.in)

*August 26, 2025*

---

# Multiple sequence alignment

---

# Multiple sequence alignment (MSA)

---

- Aligning more than 2 sequences together to find conserved residues or regions
- Input for phylogenetic reconstruction
- Given k sequences, find an alignment that maximizes the alignment score

```
VTISCTGSSSNIGAGNHVKWYQQLPG
VTISCTGTSSNIGSITVNWYQQLPG
LRLSCSSSGFIFSSYAMYWVRQAPG
LSLTCTVSGTSFDDYYSTWVRQPPG
PEVTCVVVDVSHEDPQVKFNWYVDG
ATLVCLISDFYPGAVTVAWKADS
AALGCLVKDYFPEPVTVSWNSG
VSLTCLVKGFYPSDIAVEWWSNG
```



```
VTISCTGSSSNIGAG-NHVKWYQQLPG
VTISCTGTSSNIGS--ITVNWYQQLPG
LRLSCSSSGFIFSS--YAMYWVRQAPG
LSLTCTVSGTSFDD--YYSTWVRQPPG
PEVTCVVVDVSHEDPQVKFNWYVDG--
ATLVCLISDFYPGA--VTVAWKADS--
AALGCLVKDYFPEP--VTVSWNSG---
VSLTCLVKGFYPSD--IAVEWWSNG--
```

# Scoring multiple sequence alignment

a: VSLSCTGSSSNIGAG-NHVKWYQQLPG  
 b: VTISCTGTSSNIG--SITVNWYQQLPG  
 c: ATLVCLISDFYPGA-SVTVAWKADS--  
 d: AALGCLVKDYFPEP--VTVSWNSG---  
 e: --LTCLVKGFYPSD--IAVEWESNG--

## Sum of pairs (SP) score

V
V
A
A
-

Score of column = sum of all pairs

$$\begin{aligned}
 &= S_{VV} + S_{VA} + S_{VA} + S_{V-} \\
 &\quad + S_{VA} + S_{VA} + S_{V-} \\
 &\quad + S_{AA} + S_{A-} \\
 &\quad + S_{A-}
 \end{aligned}$$

Multiple alignment score = sum of all column scores

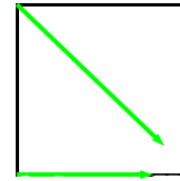
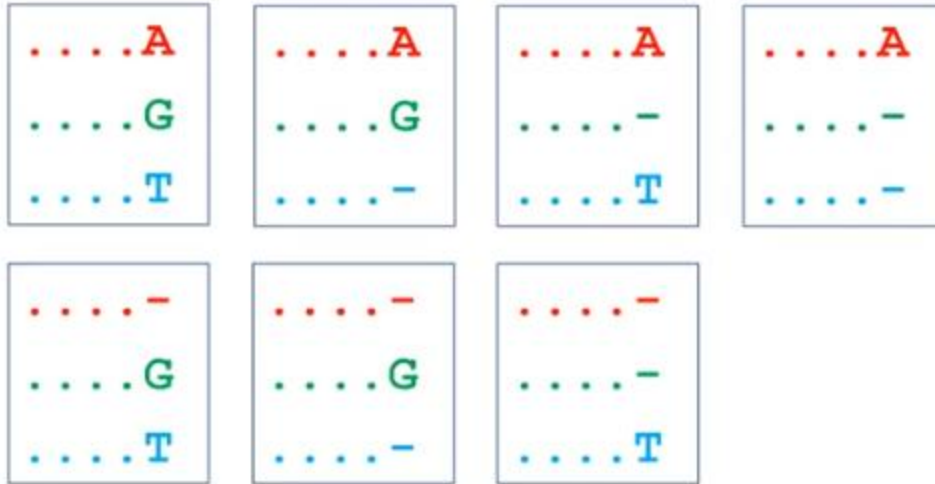
Pairwise sequence alignment score:  $S_{ab} = S_{VV} + S_{ST} + S_{LI} + S_{SS} + \dots$

Multiple alignment score = sum of all pairwise alignments =  $S_{ab} + S_{ac} + S_{ad} + S_{ae} + S_{bc} + S_{bd} + S_{be} + S_{cd} + S_{ce} + S_{de}$

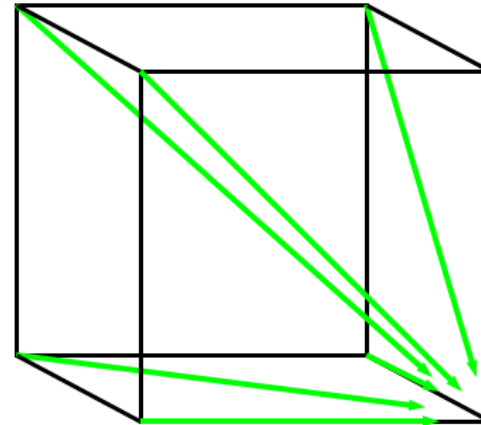
# Dynamic programming for MSA

Aligning three sequences:

.....A  
.....G  
.....T

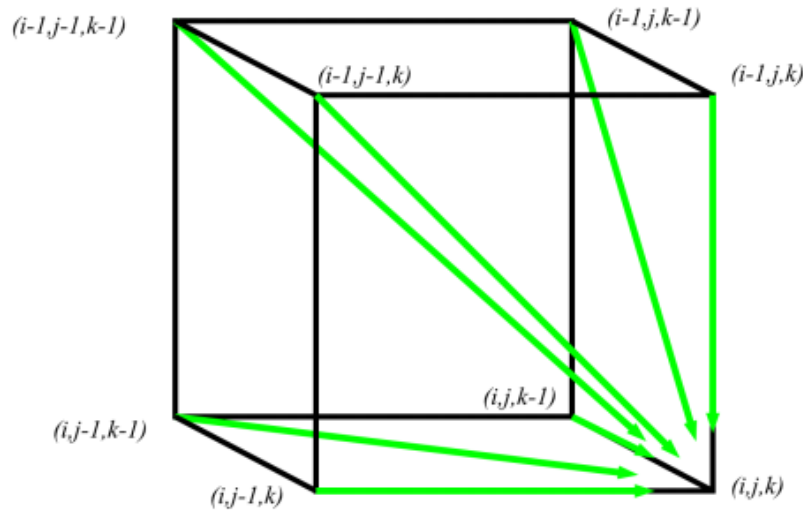


Pairwise: 3 possible paths  
(match/mismatch, insertion, and  
deletion)



In 3-D,  
7 edges in each  
unit cube

# Dynamic programming for MSA



$$S_{i,j,k} = \max$$

$$\left\{ \begin{array}{l} S_{i-1,j-1,k-1} + \delta(a_i, b_j, c_k) \\ S_{i-1,j-1,k} + \delta(a_i, b_j, \_) \\ S_{i-1,j,k-1} + \delta(a_i, \_, c_k) \\ S_{i,j-1,k-1} + \delta(\_, b_j, c_k) \\ S_{i-1,j,k} + \delta(a_i, \_, \_) \\ S_{i,j-1,k} + \delta(\_, b_j, \_) \\ S_{i,j,k-1} + \delta(\_, \_, c_k) \end{array} \right\}$$

cube diagonal:  
no indels

face diagonal:  
one indel

edge diagonal:  
two indels

# Computational complexity

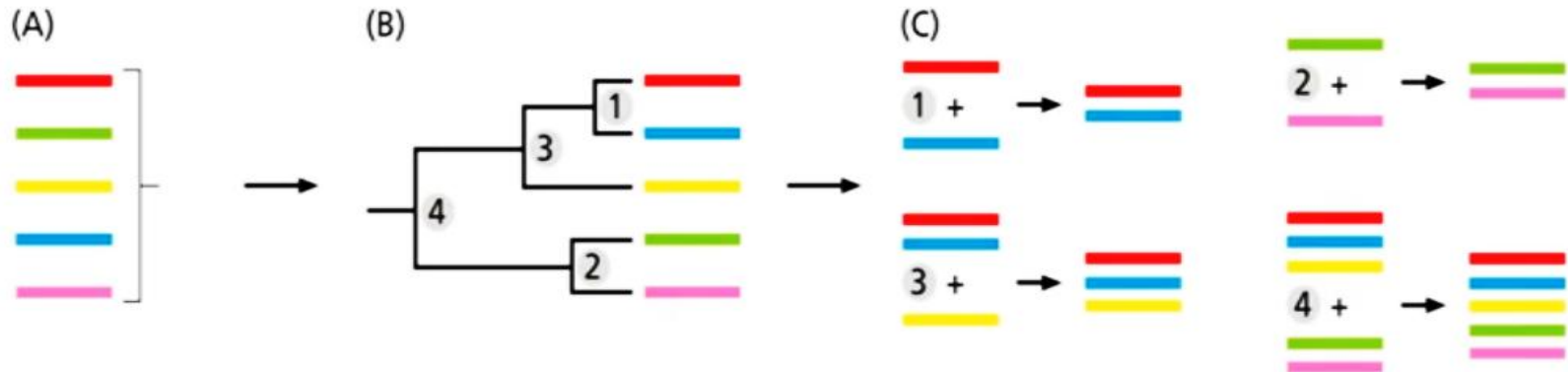
---

- For 3 sequences of length  $n$ , the run time:  $O(7n^3)$
- For 4 sequences of length  $n$ :  $O(15n^4)$
- For  $k$  sequences of length  $n$ :  $O((2^k-1)n^k) = O(2^k n^k)$
- Dynamic programming approach for alignment between two sequences is easily extended to  $k$  sequences but it is impractical due to exponential running time
- Computing exact MSA is computationally almost impossible, and in practice heuristics are used (progressive alignment)

# Progressive alignment for MSA

Three-step process:

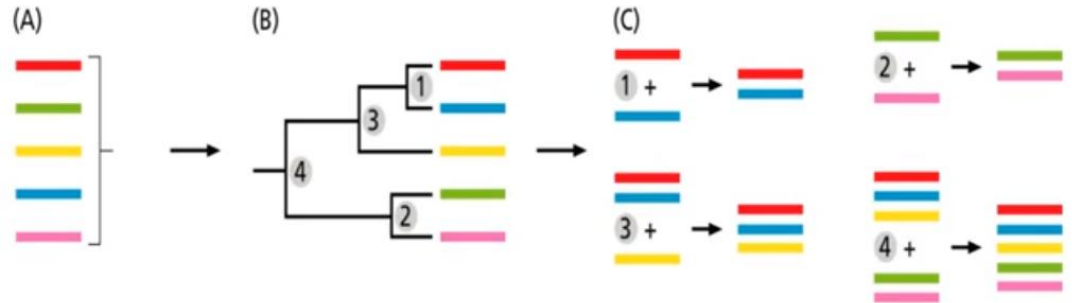
- 1) Construct pairwise alignments
- 2) Build guide tree
- 3) Progressive alignment guided by the tree





# Computational complexity

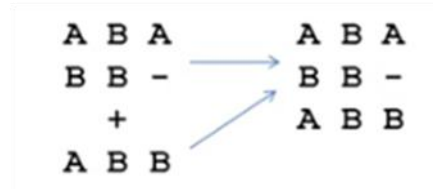
- $k$  sequences, each with length  $n$
- Each pairwise alignment:  $O(n^2)$
- Building the phylogenetic tree:
  - $O(k^2)$  pairwise comparisons
  - $O(k^2n^2)$  time
  - Done once to construct a phylogenetic tree
- Number of merge steps:
  - $O(k)$  steps
  - $O(kn^2)$  time
- Overall time:
  - $O(k^2n^2 + kn^2) = O(k^2n^2)$



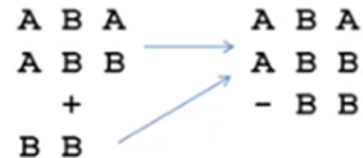
# Limitation of progressive alignment

The order of alignment matters

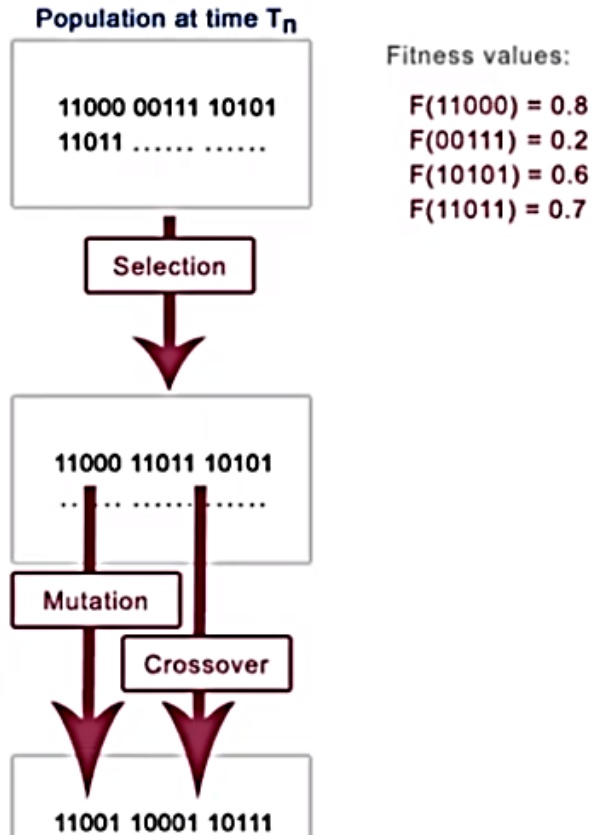
If we align ABA and BB first:



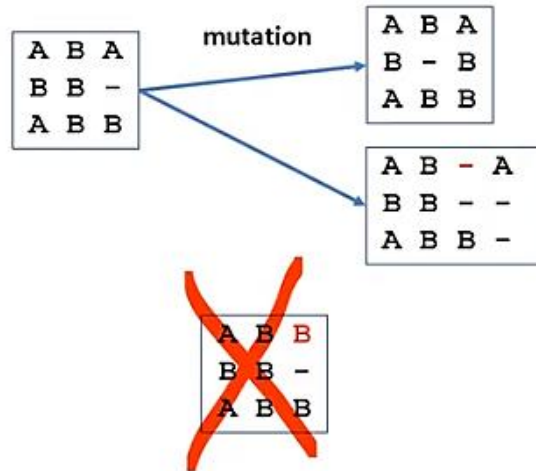
But if we align ABA and ABB first:



# MSA using Genetic algorithm



# MSA using Genetic algorithm



Mutations

WGKVN---VDEVGGEAL-  
WDKVNEEE---VGGEAL-  
WGKVG--AHAGEYGAEAL  
WSKVGGHA--GEYGAEAL

WGKV--NVDEVG-GEAL  
WDKV--NEEEVG-GEAL  
WGKVGA-HAGEYGAEAL  
WSKVGGHAGE-YGAEAL

+

--WGKVNDEVG-GEAL  
WD--KVNEEEVG-GEAL  
WGKVGA-HAGEYGAEAL  
WSKVGGHAGE-YGAEAL

one-point  
crossover

new offspring

--WGKVN---VDEVGGEAL-  
WD--KVNEEE---VGGEAL-  
WGKV--G--AHAGEYGAEAL  
WSKV--GGHA--GEYGAEAL

selection

WGKV--NVDEVG-GEAL  
WDKV--NEEEVG-GEAL  
WGKVGA-HAGEYGAEAL  
WSKVGGHAGE-YGAEAL

Crossover

---

# Database searching

---

# Database searching

---

- Search query sequence against all sequence in database
- Calculate score and select top sequences
- Dynamic programming is best
- Approximation Algorithms : FASTA BLAST

## FASTA

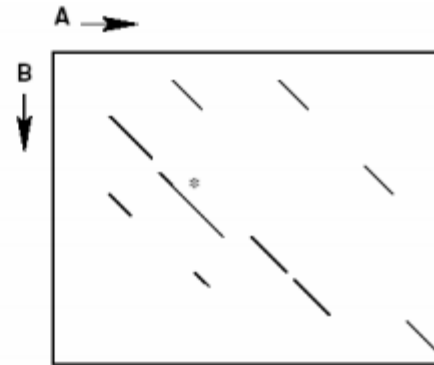
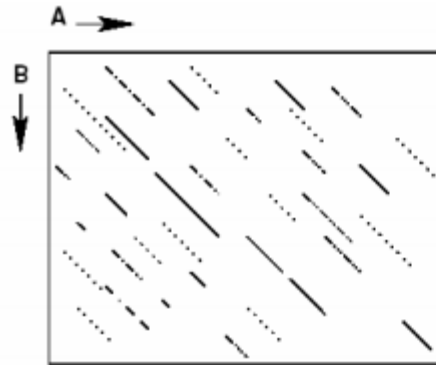
- Fast sequence search
- Based on dot plot
- Identify identical words (k-tuples)
- Search significant diagonals
- Use PAM 250 for further refinement
- Dynamic programming for narrow region

## BLAST

- Heuristic method to find the highest scoring based on local optimal alignments
- Based on ungapped sequence alignments
- Allow multiple hits to the same sequence

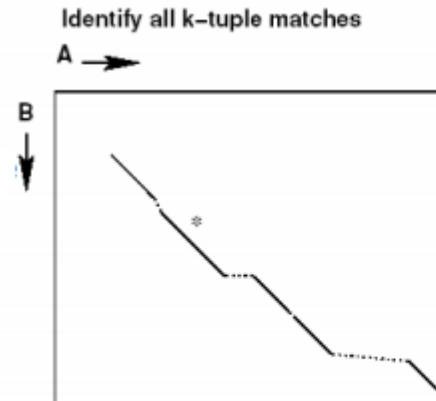
# FASTA

Localize the 10 best regions of similarity between the two seq. Each identity between two "word" is represented by a dot

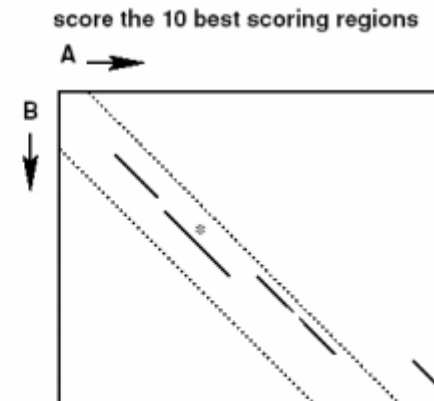


Each diagonal: ungapped alignment.  
Smaller the  $k$ , sensitive the method but slower

Find the best combination of the diagonals -> compute a score.  
Only those sequences with a score higher than a threshold will go to the fourth step



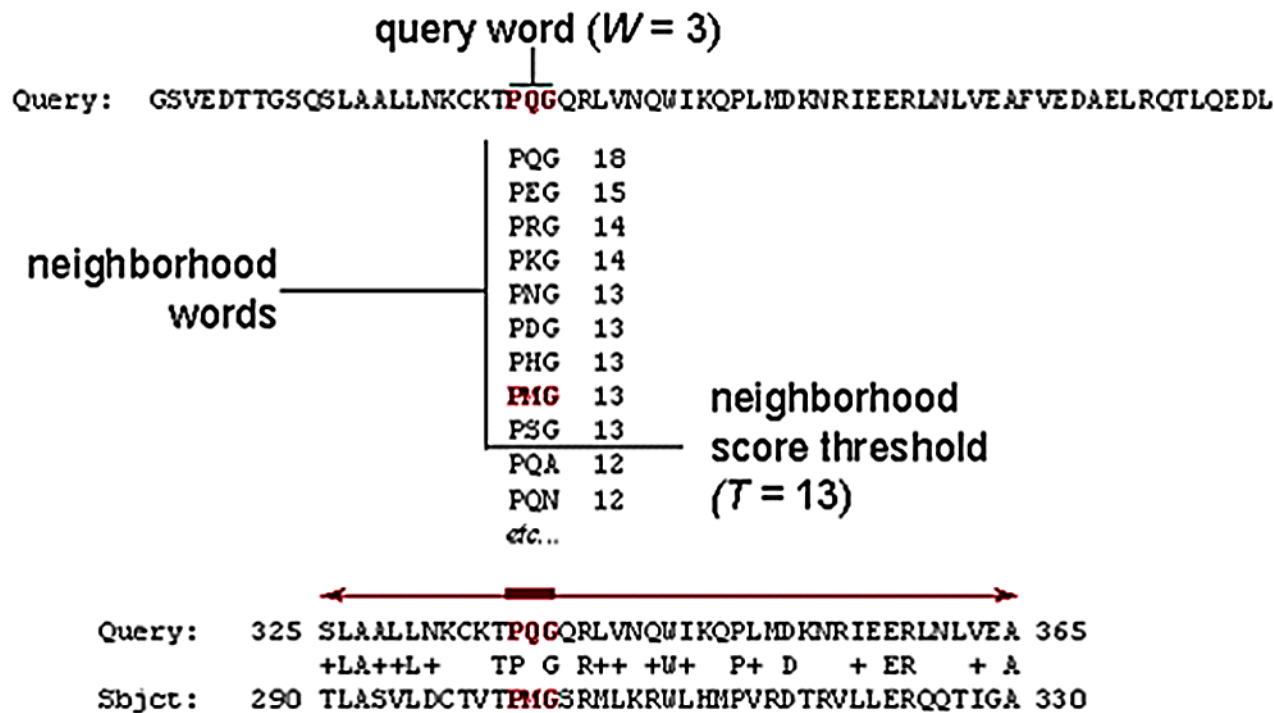
Apply joining procedure



Apply limited DP

DP applied around the best scoring diagonal.

# BLAST



High-scoring Segment Pair (HSP)



# BLAST

---

Scoring parameters: Score, E-value, Query coverage, Identity

Formula for the E-value (likely that a hit of a given score could occur just by chance in the database) in BLAST:

$$E = K \cdot m \cdot n \cdot e^{-\lambda S}$$

Where,

E = expected number of alignments with score  $\geq S$

m = length of the query sequence

n = total length of the database (sum of all sequences)

S = raw alignment score

K = a statistical constant, depends on the scoring system and sequence composition

$\lambda$  (lambda) = scaling factor (related to the scoring matrix and gap penalties) that scales the raw score into a probability scale.

Interpreting the E-value:

E = 10 → you'd expect ~10 hits with this score or better to occur just by chance.

E = 0.01 → only 1 in 100 searches of this size would produce a hit this good by chance.

E  $\approx$  0 (e.g., 1e-100) → essentially certain that the hit is real.

# BLAST variants

