

# CSE 566 Spring 2023

## **Multiple Sequence Alignment**

Instructor: Mingfu Shao

# What We've Learned

| Problem                              | Algorithm                          |
|--------------------------------------|------------------------------------|
| Querying P against text T            | Z-algo, KMP                        |
| Constructing suffix tree             | Ukkonen's algorithm                |
| Constructing suffix array            | Skew algorithm                     |
| Querying q against an suffix array   | $O( q  + \log S )$ -time algorithm |
| Constructing/Recovering/Querying BWT | $O( s )$ -time algorithms          |
| Constructing local/global alignments | $O(nm)$ -time algorithms           |
| Edit distance problem                | $O(nd)$ -time, $O(n)$ -space       |
| ...                                  | ...                                |

# Tackling NP-hard Problems

- NP-hard: there does not exist polynomial-time exact/optimal algorithm unless  $P = NP$ .
- Still need exact/optimal solution
  - exponential algorithms
- Don't require exact/optimal solution, but some theoretical guarantee
  - Approximation algorithms:  $algo(I)/OPT(I) \leq c$ , for any  $I$ .
  - Randomized algorithms
- Heuristics, often need experimental comparisons

# Multiple Sequence Alignment (MSA)

|                                   |                                 |
|-----------------------------------|---------------------------------|
| S 1 : A T C G A C T A C T C G T C | A T C - G A C T A C T C - G T C |
| S 2 : A T C G G A C T C G T T     | A T C G G A - - - C T C - G T T |
| S 3 : A C G G A C A C T C G T T   | A - C G G A C - A C T C - G T T |
| S 4 : A C G G C T A C T C A G T T | A - C G G - C T A C T C A G T T |

- Widely used to find common patterns, conserved regions, comparative analysis, etc.
- Input: sequences  $S_1, S_2, \dots, S_m$
- Output: an MSA  $M$ , such the cost of  $M$  is minimized.
- Different cost functions lead to different formulations.

# The First Formulation

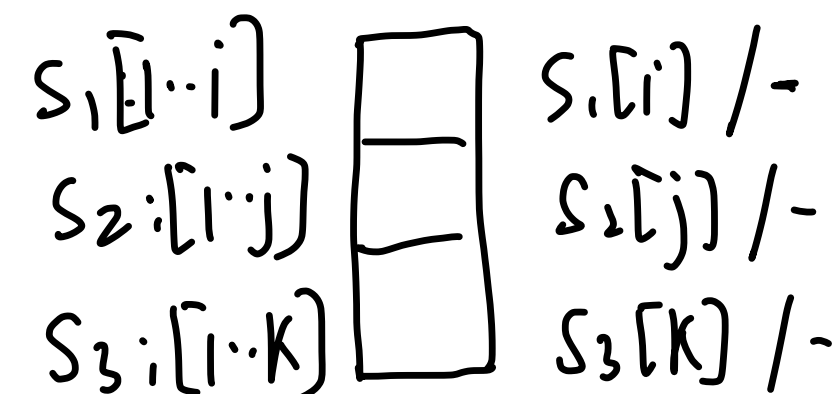
- Assume that columns of  $M$  are independent
- $cost(M)$  =  $\sum_c cost(M_c)$ , where  $M_c$  is the  $c$ -th column of  $M$
- This cost function generalizes the unit gap cost of 2 sequences to  $m$  sequences
- This formulation is NP-hard.

A T C - G A C T A C T C - G T C  
A T C G G A - - - C T C - G T T  
A - C G G A C - A C T C - G T T  
A - C G G - C T A C T C A G T T

# Dynamic Programming Algorithm

- Take  $m=3$  as example. Define  $OPT(i, j, k)$  as the minimized cost between  $S_1[1..i]$ ,  $S_2[1..j]$ , and  $S_3[1..k]$   $O(n^3)$
- Consider all possibilities of the last column

$$OPT(i, j, k) = \min_{\text{last column}} \begin{cases} cost(S_1[i], S_2[j], S_3[k]) + OPT(i-1, j-1, k-1) \\ cost(S_1[i], S_2[j], -) + OPT(i-1, j-1, k) \\ cost(S_1[i], -, S_3[k]) + OPT(i-1, j, k-1) \\ cost(S_1[i], -, -) + OPT(i-1, j, k) \\ cost(-, S_2[j], S_3[k]) + OPT(i, j-1, k-1) \\ cost(-, S_2[j], -) + OPT(i, j-1, k) \\ cost(-, -, S_3[k]) + OPT(i, j, k-1) \end{cases}$$



# Dynamic Programming Algorithm

- Generalizing to m sequence:
  - $\#(\text{subproblems}) = \underbrace{n^m}$ , where assuming  $|S_i| = n$
  - Solving each subproblem takes  $O(2^m)$  cases,
  - Assume cost of a column can be computed in  $O(m)$  time.
- Total running time:  $O(m \cdot n^m \cdot 2^m)$
- Only affordable for very small m.  $m=4 \Rightarrow O(\underline{4} \cdot n^4 \cdot 2^4)$

# The Second Formulation

- An MSA M implies all pairwise alignments.

$$\text{MSA } M \left\{ \begin{array}{l} S1 : A T C - G A C T A C T C - G T C \\ S2 : A T C G G A - - - C T C - G T T \\ S3 : A - C G G A C - A C T C - G T T \\ S4 : A - C G G - C T A C T C A G T T \end{array} \right.$$

$$\text{cost}(M) = M_{12} + M_{13} + M_{23} + M_{14} + M_{24} + M_{34}$$

$$\begin{array}{l} S2 : A T C G G A - - C T C G T T \\ S3 : A - C G G A C A C T C G T T \end{array} \quad M_{23}$$

$$\begin{array}{l} S2 : A T C G G A - - - C T C - G T T \\ S4 : A - C G G - C T A C T C A G T T \end{array} \quad M_{24}$$

- The implied pairwise alignment may not be optimal.



# The Second Formulation

- Assume the cost function for alignment between two sequences is given (say, unit gap cost, affine gap cost, etc)
- Assume that this cost function is a metric, i.e., satisfying the triangle inequality:  $cost(M_{ij}) \leq cost(M_{ik}) + cost(M_{kj}) \quad \forall k.$
- Sum-of-pairs cost:  $cost(M) = \sum_{i < j} cost(M_{i,j})$
- We are able to design a 2-approximation algorithm (STAR):

$$\frac{STAR(I)}{OPT(I)} = \frac{\sum_{i < j} cost(M_{i,j})}{\sum_{i < j} cost(M_{i,j}^*)} \leq 2$$

$\downarrow M$   
 $OPT: M^*$

# STAR Algorithm

$$O(\underline{m}^2 \cdot n^2)$$

- Step 1: build all optimal pairwise alignments; let  $M(S_i, S_j)$  be the alignment and let  $cost(S_i, S_j)$  be the cost.
- Step 2: find the sequence that is the closest to others, i.e.,  $S_x$  such that  $\sum_{i \neq x} cost(S_i, S_x)$  is minimized.
- Step 3: merge  $M(S_1, S_x)$ ,  $M(S_2, S_x)$ ,  $\dots$ ,  $M(S_m, S_x)$  into an MSA, called  $M$ .

# Merging $M(S_i, S_x), 1 \leq i \leq m$

- Keep all columns; add “-” as needed

$S_x =$  S2 : ATCGGA - - - CTCGTT ) MSA

S1 : ATC - GACTACTCGTC  
S2 : ATCGGA - - - CTCGTT  
S3 : A - CGGACACTCGTT

S2 : ATCGGA - - - CTC - GTT  
S4 : A - CGG - CTACTCAGTT

S1 : ATC - GACTACTCGTC  
S2 : ATCGGA - - - CTCGTT  
S3 : A - CGGAC - ACTCGTT

S1 : ATC - GACTACTC - GTC  
S2 : ATCGGA - - - CTC - GTT  
S3 : A - CGGAC - ACTC - GTT  
S4 : A - CGG - CTACTCAGTT

# Proof of Approximation Ratio

$$2 \cdot \text{cost}(M) = 2 \cdot \sum_{i < j} \text{cost}(M_{i,j}) = \sum_{i,j} \text{cost}(M_{i,j})$$

$$\leq \sum_{i,j} (\text{cost}(M_{i,x}) + \text{cost}(M_{x,j}))$$

triangle inequality

step 3

$$\leftarrow = \sum_{i,j} (\text{cost}(S_i, S_x) + \text{cost}(S_x, S_j))$$

$$\sum_j \left[ \sum_i \text{cost}(S_i, S_x) \right]$$

$$= \sum_{i,j} \text{cost}(S_i, S_x) + \sum_{i,j} \text{cost}(S_x, S_j)$$

$$= 2m \cdot \sum_i \text{cost}(S_i, S_x)$$

$$\underline{2 \cdot \text{cost}(M^*)} = 2 \cdot \sum_{i < j} \text{cost}(M_{i,j}^*) = \sum_{i,j} \text{cost}(M_{i,j}^*)$$

$$\geq \sum_{i,j} \text{cost}(S_i, S_j)$$

$$= \sum_i \text{cost}(S_i, S_1) + \sum_i \text{cost}(S_i, S_2) + \dots + \sum_i \text{cost}(S_i, S_m)$$

$$\geq m \cdot \sum_i \text{cost}(S_i, S_x)$$

step 2

# Transform MSA to Profile

$n$  columns

$m$  seqs

|   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|
| A | T | C | - | G | A | C | T | A |
| A | T | C | G | G | A | - | - | - |
| A | - | C | G | G | A | C | - | A |
| A | - | C | G | G | - | C | T | A |
| A | T | - | G | T | - | C | G | - |

MSA

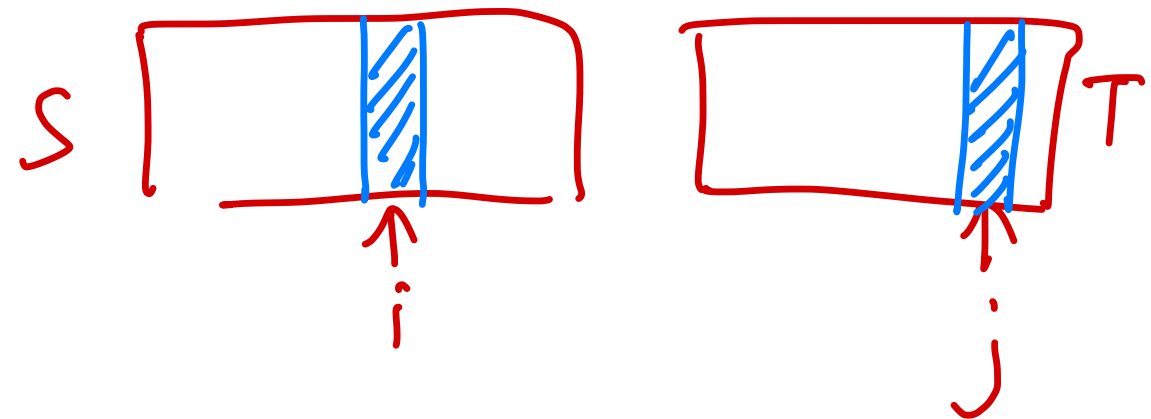
$m \times n$

|   |   |     |     |     |     |     |     |     |     |
|---|---|-----|-----|-----|-----|-----|-----|-----|-----|
| A | 1 | 0   | 0   | 0   | 0   | 0.6 | 0   | 0   | 0.6 |
| C | 0 | 0   | 0.8 | 0   | 0   | 0   | 0.8 | 0   | 0   |
| G | 0 | 0   | 0   | 0.8 | 0.8 | 0   | 0   | 0.2 | 0   |
| T | 0 | 0.6 | 0   | 0   | 0.2 | 0   | 0   | 0.4 | 0   |
| - | 0 | 0.4 | 0.2 | 0.2 | 0   | 0.4 | 0.2 | 0.4 | 0.4 |

$5 \times n$

# Aligning Two Profiles

- Input: two profiles  $S$  and  $T$
- DP algorithm:



$$OPT(i, j) = \max \begin{cases} OPT(i-1, j-1) + score(S[i], T[j]) \\ OPT(i-1, j) + score(S[i], -) \\ OPT(i, j-1) + score(-, T[j]) \end{cases}$$

$$score(S[i], T[j]) = \sum_{c \in \Sigma} S[i][c] \cdot T[j][c]$$

|   |   |
|---|---|
| 0 | A |
| 0 | C |
| 0 | G |
| 0 | T |
| 1 | - |

- More sophisticated approaches exist, such as profile HMM; see more in “Biological sequence analysis”, by R. Durbin, et al.

# Progressive Approach for MSA

- Step 1: build all pairwise alignments to get the cost matrix.
- Step 2: build a guide tree (clustering, phylogeny, etc).
- Step 3: progressively align two profiles following the tree.

|   | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| 1 |   |   |   |   |   |   |   |
| 2 |   |   |   |   |   |   |   |
| 3 |   |   |   |   |   |   |   |
| 4 |   |   |   |   |   |   |   |
| 5 |   |   |   |   |   |   |   |
| 6 |   |   |   |   |   |   |   |
| 7 |   |   |   |   |   |   |   |

