

CS1810 9/18/25

# Computational Biology

## Local Alignment:

Given: 2 sequences  $X, Y$  over the same alphabet, and a scoring scheme  $S$ .

Compute: the maximum score local alignment of  $X$  with  $Y$ .

↳ If 2 random sequences are aligned of length  $N$ , the length of the global alignment on average is  $O(N)$ , while the length of the local alignment on average is  $O(\log N)$ .

- The optimal local alignment is the alignment of 2 substrings in  $X$  and  $Y$  that produce the highest possible score. In other words, choose a substring of  $X$  and a substring of  $Y$  such that when these two substrings are aligned, the maximum score is produced.

$X - \boxed{\text{---}}$

Ex.

$Y - \boxed{\text{---}}$

$\alpha^*$

$\beta^*$

$\beta^*$

is the optimal local alignment.

# Computational Biology

Suffixes, Prefixes, and Substrings:

Suppose  $X = ACBAAA$ .

$\hookrightarrow \text{Pref}(X) = \{A, AC, ACB, ACBA, ACBAA, ACBAAA\}$ .

$\hookrightarrow \text{Suff}(X) = \{A, AA, AAA, BAAA, CBAAA, ACBAAA\}$ .

$\hookrightarrow \text{Sub}(X)$ , there are  $O(n^2)$  substrings for a string of length  $n$ .

Theorem: The prefixes of suffixes or the suffixes of prefixes for a string is equivalent to the set of substrings.

$\hookrightarrow \text{Pref}(\text{Suff}(X)) = \text{Suff}(\text{Pref}(X)) = \text{Sub}(X)$ .

That begs the question: For sequences  $X$  and  $Y$  of length  $M$  and  $N$  respectively, there are  $O(N^2M^2)$  possible local alignments possible. How does local alignment do it in  $O(NM)$  time?

Suppose  $X_i$  is a prefix of  $X$  up to position  $i$  and  $Y_j$  is a prefix of  $Y$  up to position  $j$ .

$X_i : X - \underbrace{\dots}_{i}, Y_j : Y - \underbrace{\dots}_{j}$

Now, define  $V(i, j)$  as the value of the optimal global alignment between all pairs of suffixes of  $X_i$  and  $Y_j$ .

# Computational Biology

Local Alignment algorithm:

$\Rightarrow$  Initialize  $V(i, 0) = 0$  for all  $i, j$

↳ Initialize  $V(0,j) = 0$

~~FOR i>0 and j>0 DO~~

$$V(i, j) = \max \left\{ \begin{array}{l} V(i-1, j-1) + S(x_i, y_j) \\ V(i-1, j) + S(x_i, -) \\ V(i, j-1) + S(-, y_j) \end{array} \right\}$$

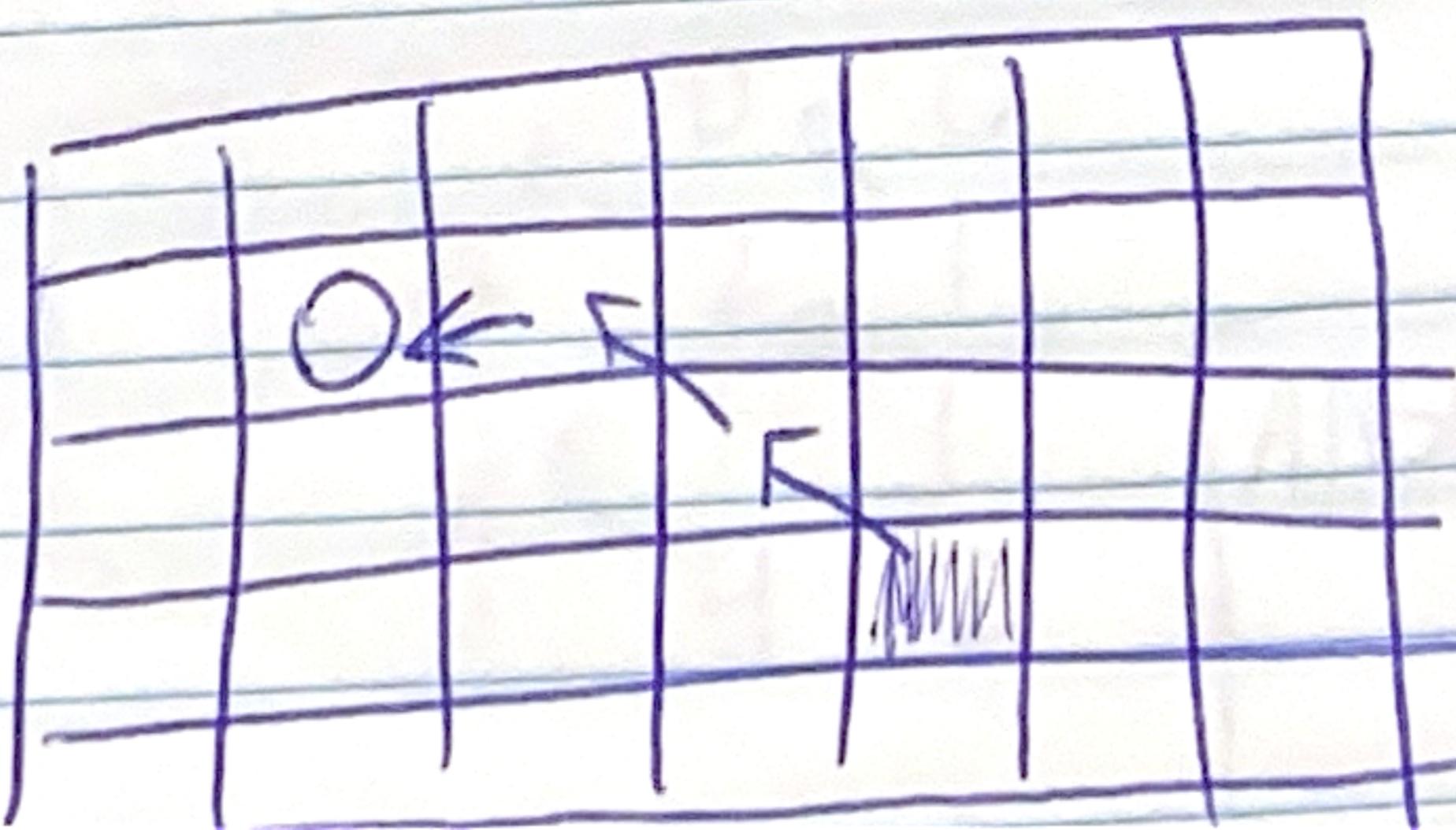
Finally, return  $V(X, Y) = \max_{\substack{0 \leq i \leq N \\ 0 \leq j \leq M}} [V(i, j)]$   
 (OLA = optimal local alignment)

Note: In global alignment we choose the bottom right corner of the DP table as the alignment score. However, in local alignment, we choose the maximum value among all the cells of the DP table.

To find the local alignment that is optimal, backtrack from the max score back to a 0 value. AKA follows the backpointers until you encounter a 0.

Ex.

1



# Computational Biology

To understand why this algorithm works, here is a useful theorem:

Thus:  $V^{OLD} = \text{MAX}[V(i,j), 1 \leq i \leq N, 1 \leq j \leq M]$ ,

where  $V^{OLD}$  is the alignment between substrings of X and Y resulting in the maximum score.

Proof: We know that  $V^{OLD} \geq \text{MAX}[V(i,j), 1 \leq i \leq N, 1 \leq j \leq M]$  because  $V(i,j)$  considers suffixes of prefixes, which are therefore also substrings.

We also know that  $\text{MAX}[V(i,j), 1 \leq i \leq N, 1 \leq j \leq M] \geq V^{OLD}$  because the left hand side considers all possible suffixes of prefixes and therefore all possible substring pairs. This means that the MAX score over that will be at least the score resulting from one specific pair of substrings (aka  $V^{OLD}$ ).

$\hookrightarrow V^{OLD} = \text{MAX}[V(i,j), 1 \leq i \leq N, 1 \leq j \leq M]$ .

Ex. A C T C C

Q < 0 < 0 < 0 < 0 < 0

A Q 1 0 0 0 0 0

MATCH : +1

C Q 0 2 0 1 1

MISMATCH : -1

A Q 1 0 1 0 0

INDEL : -2

C Q 0 2 0 2 1

T Q 0 0 3 1 1

C Q 0 1 1 4 2

ANSWER: ACTC

||||  
ACTC