**Dept. of Computer Science and Engineering**
**University of Rajshahi**
**www.ru.ac.bd**

**Dr. Shamim Ahmad**

# Local alignment

**Local alignment** methods f ind related regions *within* sequences - they can consist of a **subset** of the characters within each sequence.

For example, positions 20-40 of *sequence A* might be aligned with positions 50-70 of *sequence B*.

This is a more f lexible technique than *global alignment* and has the adv antage that *related regions* which appear in a dif f erent order in the two proteins (which is known as **domain shuffling**) can be identif ied as being related.

This is **not** possible with *global alignment* methods.

2

## Similarity Matrix

|   | A | G | C | T |
|---|---|---|---|---|
| **A** | 1 | -1 | -1 | -1 |
| **G** | -1 | 1 | -1 | -1 |
| **C** | -1 | -1 | 1 | -1 |
| **T** | −1 | −1 | −1 | 4 |

# The Smith Waterman algorithm

The **Smith-Waterman algorithm** (1981) is for determining similar regions between two nucleotide or protein sequences.

Smith-Waterman is also a dynamic programming algorithm and improves on Needleman-Wunsch. As such, it has the desirable property that it is guaranteed to find the **optimal local alignment** with respect to the scoring system being used (which includes the substitution matrix and the gap-scoring scheme).

However, the Smith-Waterman algorithm is **demanding of time and memory** resources: in order to align two sequences of lengths m and n, O(mn) time and space are required.

As a result, it has largely been replaced in practical use by the **BLAST algorithm**; although not guaranteed to find optimal alignments, BLAST is much more efficient.

4

## Smith–Waterman Algorithm

| | Smith–Waterman algorithm | Needleman–Wunsch algorithm |
|---|---|---|
| Initialization | First row and first column are set to 0 | First row and first column are subject to gap penalty |
| Scoring | Negative score is set to 0 | Score can be negative |
| Traceback | Begin with the highest score, end when 0 is encountered | Begin with the cell at the lower right of the matrix, end at top left cell |

## The Smith-Waterman algorithm

1. Create a table of size $(m+1) \times (n+1)$ for sequences **s** and **t** of lengths $m$ and $n$,

2. Fill table entries $(1,1:m+1)$ and $(1:n+1,1)$ with zeros.

3. Starting from the top left, compute each entry using the recursive relation:

$$M_{i,j} = \max \begin{Bmatrix} M_{i-1,j-1} + \sigma(\mathbf{s}_i, \mathbf{t}_j) \\ M_{i-1,j} + \sigma(\mathbf{s}_i, -) \\ M_{i,j-1} + \sigma(-, \mathbf{t}_j) \\ 0 \end{Bmatrix}$$

4. Perform the trace-back procedure from the maximum element in the table to the first zero element on the trace-back path.

3. Fill the scoring matrix using the equation below.

$$H_{ij} = \max \begin{cases} H_{i-1,j-1} + s(a_i, b_j), \\ \max_{k \geq 1}\{H_{i-k,j} - W_k\}, \\ \max_{l \geq 1}\{H_{i,j-l} - W_l\}, \\ 0 \end{cases} \quad (1 \leq i \leq n, 1 \leq j \leq m)$$

where

$H_{i-1,j-1} + s(a_i, b_j)$ is the score of aligning $a_i$ and $b_j$,
$H_{i-k,j} - W_k$ is the score if $a_i$ is at the end of a gap of length $k$,
$H_{i,j-l} - W_l$ is the score if $b_j$ is at the end of a gap of length $l$,
0 means there is no similarity up to $a_i$ and $b_i$.

7

## Step 3: Computing the length of a LCS



LCS-LENGTH(X, Y)
1  m = X.length
2  n = Y.length
3  let b[1...m, 1...n] and c[0...m, 0...n] be new tables
4  for i = 1 to m
5      c[i, 0] = 0
6  for j = 0 to n
7      c[0, j] = 0
8  for i = 1 to m
9      for j = 1 to n
10         if x_i == y_j
11             c[i, j] = c[i − 1, j − 1] + 1
12             b[i, j] = "↖"
13         elseif c[i − 1, j] ≥ c[i, j − 1]
14             c[i, j] = c[i − 1, j]
15             b[i, j] = "↑"
16         else c[i, j] = c[i, j − 1]
17             b[i, j] = "←"
18  return c and b

BCBA ⟸ AB C BDAB ⟸ BDCAB A

---

# Simplified Smith–Waterman algorithm

When linear gap penalty function is used
A linear gap penalty has the same scores for opening and extending a gap:

**Linear**   [ edit ]

A linear gap penalty has the same scores for opening and extending a gap:

$$W_k = kW_1,$$

where $W_1$ is the cost of a single gap.

---

## Step 4: Constructing a LCS (Backtracking)

BCBA ⟸ AB C BDAB ⟸ BDCAB A



PRINT-LCS(b, X, i, j)
1  if i == 0 or j == 0
2      return
3  if b[i, j] == "↖"
4      PRINT-LCS(b, X, i − 1, j − 1)
5      print x_i
6  elseif b[i, j] == "↑"
7      PRINT-LCS(b, X, i − 1, j)
8  else PRINT-LCS(b, X, i, j − 1)

---

# Similarity Matrix

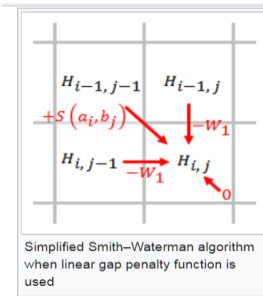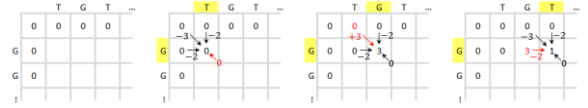|   | A | G | C | T |
|---|---|---|---|---|
| A | 1 | -1 | -1 | -1 |
| G | -1 | 1 | -1 | -1 |
| C | -1 | -1 | 1 | -1 |
| T | -1 | -1 | -1 | 1 |

This substitution matrix can be described as: $s(a_i, b_j) = \begin{cases} +1, & a_i = b_j \\ -1, & a_i \neq b_j \end{cases}$

## Initialize the scoring matrix

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

Substitution matrix:
$$S(a_i, b_j) = \begin{cases} +3, & a_i = b_j \\ -3, & a_i \neq b_j \end{cases}$$

Gap penalty:
$$W_k = kW_1$$
$$W_1 = 2$$





Simplified Smith–Waterman algorithm when linear gap penalty function is used

$$H_{ij} = \max \begin{cases} H_{i-1,j-1} + s(a_i, b_j), \\ H_{i-1,j} - W_1, \\ H_{i,j-1} - W_1, \\ 0 \end{cases}$$

|   | | T | G | T | T | A | C | G | G |
|---|---|---|---|---|---|---|---|---|---|
|   | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| G | 0 | 0 | 3 | 1 | 0 | 0 | 0 | 3 | 3 |
| G | 0 | 0 | 3 | 1 | 0 | 0 | 0 | 3 | 6 |
| T | 0 | 3 | 1 | 6 | 4 | 2 | 0 | 1 | 4 |
| T | 0 | 3 | 1 | 4 | 9 | 7 | 5 | 3 | 2 |
| G | 0 | 1 | 6 | 4 | 7 | 6 | 4 | 8 | 6 |
| A | 0 | 0 | 4 | 3 | 5 | 10 | 8 | 6 | 5 |
| C | 0 | 0 | 2 | 1 | 3 | 8 | 13 | 11 | 9 |
| T | 0 | 3 | 1 | 5 | 4 | 6 | 11 | 10 | 8 |
| A | 0 | 1 | 0 | 3 | 2 | 7 | 9 | 8 | 7 |

4

## Dynamic Programming

|     | GAP | M | N | A | L | S | D | R | T |
|-----|-----|---|---|---|---|---|---|---|---|
| GAP | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| M | 0 | 6 | 0 | 0 | 4 | 0 | 0 | 0 | 0 |
| G | 0 | 0 | 6 | 1 | 0 | 5 | 1 | 0 | 0 |
| S | 0 | 0 | 1 | 7 | 0 | 2 | 5 | 1 | 1 |
| D | 0 | 0 | 2 | 1 | 3 | 0 | 6 | 4 | 1 |
| R | 0 | 0 | 0 | 0 | 0 | 3 | 0 | 12 | 3 |
| T | 0 | 0 | 0 | 1 | 0 | 1 | 3 | 0 | 15 |
| T | 0 | 0 | 0 | 1 | 0 | 1 | 1 | 2 | 3 |
| E | 0 | 0 | 1 | 0 | 0 | 0 | 4 | 0 | 2 |
| T | 0 | 0 | 0 | 2 | 0 | 1 | 0 | 3 | 3 |

```
SDRT
SDRT
```

## Substitution Score

### Substitution matrix (BLOSUM 50 matrix)

Log odds score can be positive (identities, conservative replacements) and negative

## Aligning locally using BLOSUM 62

|     |   | A | A | E | E | K | K | L | A | A | A |
|-----|---|---|---|---|---|---|---|---|---|---|---|
|     | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| A | 0 | 4 | 4 | 0 | 0 | 0 | 0 | 0 | 4 | 4 | 4 |
| A | 0 | 4 | 8 | 3 | 0 | 0 | 0 | 0 | 4 | 8 | 8 |
| R | 0 | 0 | 3 | 8 | 3 | 2 | 2 | 0 | 0 | 3 | 7 |
| R | 0 | 0 | 0 | 3 | 8 | 5 | 4 | 0 | 0 | 0 | 2 |
| I | 0 | 0 | 0 | 0 | 0 | 5 | 2 | 6 | 0 | 0 | 0 |
| A | 0 | 4 | 4 | 0 | 0 | 0 | 4 | 1 | 10 | 4 | 4 |

```
KKLA
RRIA
```
Score: 10

5

## Algorithms

KATHOLIEKE UNIVERSITEIT
**LEUVEN**

**Pairwise Alignment**

**Dynamic programming**

**Needleman Wunsch**

**(global)**

Chapter 1

**Smith Waterman**

**(local)**

**Heuristic approaches**

FastA

**Blast**

**Database searches**

Chapter 1