

## Local similarity

- Up to this point we have focused on global comparison.
- That is, we have compared the complete sequence  $u$  with the complete sequence  $v$ .
- In DNA sequences we often have long non-coding regions and small coding regions.
- Thus if two coding regions in two large sequences are similar, this does not imply that the sequences have a small edit distance, see the Example 1.
- In an analogy, protein sequences often contain several domains which they share with related sequences.
- But if only the domains are similar, then a global comparison would not reveal the similarity of the domains.

## Example 1

The following global alignment shows two sequences with similarities over the entire range of positions. Hence a global alignment is an appropriate means to represent the similarities.

```
--T--CC-C-AGT--TATGT-CAGGGGACACG--A-GCATGCAGA-AC
  |  || |  ||  | | | |||      || |  | |  ||||  |
AATTGCCGCC-GTCGT-T-TTCAG----CA-GTTATG--TCAGAT--C
```

The middle line uses the symbol | to mark columns of identical characters. Sometimes only segments of the sequences show similarities.

```
CGGGGAGGCGCCCGCGCGGCCAGTTAT-TCAGAATTAAATCTTAGTAAACAT
      |||| | |||
AATTAAATCTTAGTAAACATCAGTT-TGTCAGCGGGGAGGCGCCCGCGCGGCC
```

Thus a global alignment does not make sense in this case and one should look for local similarities, displayed in a local alignment.

## Local similarity

- As a consequence, when comparing biological sequences it is sometimes important to perform local similarity comparisons: This means to find pairs of similar *substrings* of  $u$  and  $v$ .
- It does not make sense to look for pairs of substrings with some minimum distance, as  $\varepsilon$  is a substring of any sequence and the sequences  $\varepsilon$  and  $\varepsilon$  have distance 0.
- To clarify the notion of sequence similarity, we introduce score functions.

### Definition 1

A *score function*  $\sigma$  assigns to each edit operation  $\alpha \rightarrow \beta$  a *score*  $\sigma(\alpha \rightarrow \beta) \in \mathbb{R}$ . For each alignment  $A = (\alpha_1 \rightarrow \beta_1, \dots, \alpha_h \rightarrow \beta_h)$  we define the score  $\sigma(A) = \sum_{i=1}^h \sigma(\alpha_i \rightarrow \beta_i)$ . The similarity score of the sequences  $u'$  and  $v'$  is defined by

$$\text{score}_{\sigma}(u', v') = \max\{\sigma(A) \mid A \text{ is an alignment of } u' \text{ and } v'\}. \quad \square$$

**Table 1:** The BLOSUM62 similarity score matrix specifying a replacement score for each pair of amino acid. As the matrix is symmetric only the lower half of the matrix is shown. With some additional scores for insertions and deletions we would obtain a score function.

	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
A	4																			
R	-1	5																		
N	-2	0	6																	
D	-2	-2	1	6																
C	0	-3	-3	-3	9															
Q	-1	1	0	0	-3	5														
E	-1	0	0	2	-4	2	5													
G	0	-2	0	-1	-3	-2	-2	6												
H	-2	0	1	-1	-3	0	0	-2	8											
I	-1	-3	-3	-3	-1	-3	-3	-4	-3	4										
L	-1	-2	-3	-4	-1	-2	-3	-4	-3	2	4									
K	-1	2	0	-1	-3	1	1	-2	-1	-3	-2	5								
M	-1	-1	-2	-3	-1	0	-2	-3	-2	1	2	-1	5							
F	-2	-3	-3	-3	-2	-3	-3	-3	-1	0	0	-3	0	6						
P	-1	-2	-2	-1	-3	-1	-1	-2	-2	-3	-3	-1	-2	-4	7					
S	1	-1	1	0	-1	0	0	0	-1	-2	-2	0	-1	-2	-1	4				
T	0	-1	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	1	5			
W	-3	-3	-4	-4	-2	-2	-3	-2	-2	-3	-2	-3	-1	1	-4	-3	-2	11		
Y	-2	-2	-2	-3	-2	-1	-2	-3	2	-1	-1	-2	-1	3	-3	-2	-2	2	7	
V	0	-3	-3	-3	-1	-2	-2	-3	-3	3	1	-2	1	-1	-2	-2	0	-3	-1	4

# Local similarity

## Definition 2

Let  $\sigma$  be a score function. We define

- 1  $loc_{\sigma}(u, v) = \max\{score_{\sigma}(u', v') \mid \begin{array}{l} u' \text{ is substring of } u \text{ and} \\ v' \text{ is substring of } v \end{array}\}$

That is,  $loc_{\sigma}(u, v)$  is the maximum score over all pairs of substrings of  $u$  and  $v$ .

- 2 Let  $u'$  be a substring of  $u$  and  $v'$  be a substring of  $v$  such that

$$score_{\sigma}(u', v') = loc_{\sigma}(u, v)$$

An alignment  $A$  of  $u'$  and  $v'$  satisfying  $score_{\sigma}(u', v') = \sigma(A)$  is an *optimal local alignment* of  $u$  and  $v$ .

- 3 The *local alignment problem* is to compute  $loc_{\sigma}(u, v)$  and an optimal local alignment of  $u$  and  $v$ .  $\square$

## Example 2

Consider the sequences  $u = \text{PQRAFADCSTVQ}$  and  $v = \text{FYAFDACSL}$  and the following similarity score function:

$$\sigma(\alpha \rightarrow \beta) = \begin{cases} -1 & \text{if } \alpha = \varepsilon \text{ or } \beta = \varepsilon \\ -2 & \text{if } \alpha, \beta \in \mathcal{A} \text{ and } \alpha \neq \beta \\ +2 & \text{if } \alpha, \beta \in \mathcal{A} \text{ and } \alpha = \beta \end{cases}$$

Then  $\text{loc}_\sigma(u, v) = 8 = \text{score}_\sigma(u', v')$ , where  $u' = \text{AFADCS}$ ,  $v' = \text{AFDACS}$  and

A	F	-	A	D	C	S
A	F	D	A	-	C	S
2	2	-1	2	-1	2	2

of score 8 is an optimal local alignment (with the score of each column in the last line).

## Local similarity

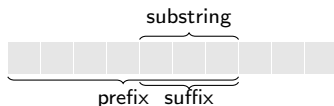
A brute force solution to the local alignment problem would be as follows:

*for each pair  $(u', v')$  of substrings  $u'$  of  $u$  and  $v'$  of  $v$  compute  $\text{score}_\sigma(u', v')$  and determine the maximum over all such values*

- Computing  $\text{score}_\sigma(u', v')$  is a global alignment problem where, instead of minimizing distances, similarities are maximized.
- Therefore, we can modify the DP Algorithm for computing the edit distance in such a way that a score function is used and maxima are computed instead of minima.
- Thus  $\text{score}_\sigma(u', v')$  can be computed in  $O(|u'| |v'|) = O(mn)$  time.
- There are  $O(m^2 n^2)$  pairs  $(u', v')$  of substrings of  $u$  and  $v$ .
- Thus, this method would require  $O(m^2 n^2 mn) = O(m^3 n^3)$  time, which is, of course, too much.

## Local similarity

- Note that a substring of a string is a suffix of a prefix of that string:



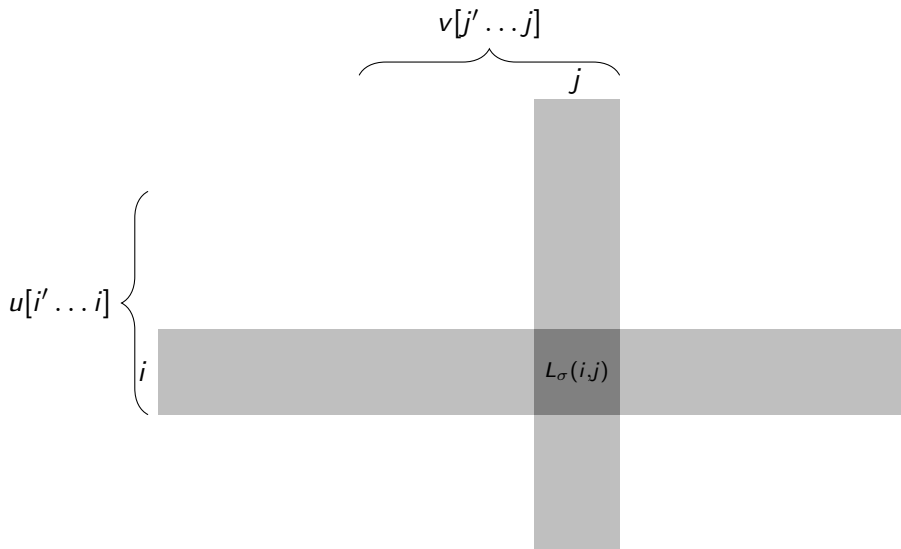
- So substring  $u'$  of  $u$  is a suffix of a prefix of  $u$  and substring  $v'$  of  $v$  is a suffix of a prefix of  $v$ .
- For example, if  $u' = u[p \dots q]$  (i.e.  $u'$  is a substring of  $u$ ), then  $u[1 \dots q]$  is a prefix of  $u$  and  $u'$  is a suffix of this prefix.
- The idea is to compute a matrix where each entry  $(i, j)$  contains the maximal score for all pairs of suffixes of prefix  $u[1 \dots i]$  of  $u$  and prefix  $v[1 \dots j]$  of  $v$ .
- That is, we compute an  $(m + 1) \times (n + 1)$ -matrix  $L_\sigma$  defined by

$$L_\sigma(i, j) = \max\{\text{score}_\sigma(x, y) \mid x \text{ is suffix of } u[1 \dots i] \text{ and } y \text{ is suffix of } v[1 \dots j]\}$$

See also Figure 1 for an illustration.



**Figure 1:** The curly brackets mark the substrings of  $u$  and  $v$  an entry in  $L_\sigma$  refers to. In particular,  $L_\sigma(i, j)$  refers to pairs of suffixes  $u[i' \dots i]$  of  $u[1 \dots i]$  and suffixes  $v[j' \dots j]$  of  $v[1 \dots j]$ .



## Local similarity

$$\begin{aligned}loc_{\sigma}(u, v) &= \max\{score_{\sigma}(u', v') \mid u' \text{ is substring of } u, v' \text{ is substring of } v\} \\&= \max\{score_{\sigma}(x, y) \mid 0 \leq i \leq m, 0 \leq j \leq n, \\&\quad x \text{ is suffix of } u[1 \dots i], \\&\quad y \text{ is suffix of } v[1 \dots j]\} \\&= \max\{\max\{score_{\sigma}(x, y) \mid x \text{ is suffix of } u[1 \dots i], \\&\quad y \text{ is suffix of } v[1 \dots j]\} \mid 0 \leq i \leq m, \\&\quad 0 \leq j \leq n\} \\&= \max\{L_{\sigma}(i, j) \mid 0 \leq i \leq m, 0 \leq j \leq n\}\end{aligned}$$

- So  $loc_{\sigma}(u, v)$  can be computed by maximizing over all entries in  $L_{\sigma}$ .
- Consider an edit graph representing all local alignments.
- Since we are interested in alignments of all pairs of substrings of  $u$  and  $v$ , we are interested in each path.
- The paths do not necessarily have to start at  $(0, 0)$  or end at  $(m, n)$ .
- Since a path can begin at any node, we have to allow the score 0 in any entry of the matrix.

## Local similarity

These considerations lead to the following result:

### Theorem 3

Let  $\sigma$  be a score function satisfying

$$\text{loc}_{\sigma}(s, \varepsilon) = \text{loc}_{\sigma}(\varepsilon, s) = 0 \quad (1)$$

for any sequence  $s \in \mathcal{A}^*$ . Then the following holds:

- If  $i = 0$  or  $j = 0$ , then  $L_{\sigma}(i, j) = 0$ .
- Otherwise,

$$L_{\sigma}(i, j) = \max \left\{ \begin{array}{l} 0 \\ L_{\sigma}(i-1, j) + \sigma(u[i] \rightarrow \varepsilon) \\ L_{\sigma}(i, j-1) + \sigma(\varepsilon \rightarrow v[j]) \\ L_{\sigma}(i-1, j-1) + \sigma(u[i] \rightarrow v[j]) \end{array} \right\}$$

## Local similarity

- In general, it is not easy to verify condition (1).
- However, one usually requires that  $\sigma$  is biased towards negative indel scores, that is, it holds:  $\sigma(\alpha \rightarrow \beta) < 0$  for all insertions and deletions  $\alpha \rightarrow \beta$ .
- This condition then implies (1):

$$\begin{aligned} loc_{\sigma}(s, \varepsilon) &= \max\{score_{\sigma}(x, \varepsilon) \mid x \text{ is a substring of } s\} \\ &= \max\{score_{\sigma}(\varepsilon, \varepsilon)\} \cup \{score_{\sigma}(x, \varepsilon) \mid x \text{ is a substring of } s, x \neq \varepsilon\} \\ &= 0 \end{aligned}$$

One can similarly show  $loc_{\sigma}(\varepsilon, s) = 0$ .

Based on Theorem 3, we can derive Algorithm 1 which solves the local similarity search problem.

### Algorithm 1 (Smith-Waterman Algorithm)

**Input:** sequences  $u = u[1 \dots m]$  and  $v = v[1 \dots n]$   
score function  $\sigma$  satisfying (1)

**Output:**  $loc_{\sigma}(u, v)$  and an optimal local alignment of  $u$  and  $v$ .

- 1 Compute matrix  $L_{\sigma}$  according to Theorem 3.
- 2 Compute a maximal entry, say  $L_{\sigma}(i, j)$ , in  $L_{\sigma}$ .
- 3 Compute optimal local alignments by a traceback from  $(i, j)$  on a maximizing path until some entry  $(i', j')$  satisfying  $L_{\sigma}(i', j') = 0$  is found.

Algorithm 1 requires  $O(mn)$  time and space.

### Example 3

Consider the similarity score  $\sigma(\alpha \rightarrow \beta) = \begin{cases} -1 & \text{if } \alpha = \varepsilon \text{ or } \beta = \varepsilon \\ -2 & \text{if } \alpha, \beta \in \mathcal{A} \text{ and } \alpha \neq \beta \\ +2 & \text{if } \alpha, \beta \in \mathcal{A} \text{ and } \alpha = \beta \end{cases}$

and sequences  $u = \text{PQRAFADCSTVQ}$  and  $v = \text{FYAFDACSLL}$ .

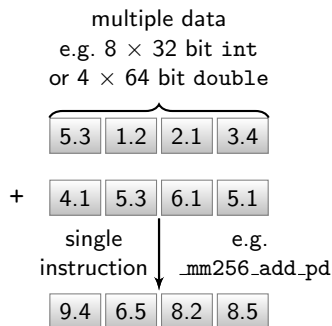
$L_\sigma$		F	Y	A	F	D	A	C	S	L	L
	0	0	0	0	0	0	0	0	0	0	0
P	0	0	0	0	0	0	0	0	0	0	0
Q	0	0	0	0	0	0	0	0	0	0	0
R	0	0	<b>0</b>	0	0	0	0	0	0	0	0
A	0	0	0	<b>2</b>	1	0	2	1	0	0	0
F	0	2	1	1	<b>4</b>	<b>3</b>	2	1	0	0	0
A	0	1	0	3	3	2	<b>5</b>	4	3	2	1
D	0	0	0	2	2	5	<b>4</b>	3	2	1	0
C	0	0	0	1	1	4	3	<b>6</b>	5	4	3
S	0	0	0	0	0	3	2	5	<b>8</b>	7	6
T	0	0	0	0	0	2	1	4	7	6	5
V	0	0	0	0	0	1	0	3	6	5	4
Q	0	0	0	0	0	0	0	2	5	4	3

optimal  
local  
alignment:

A F - A D C S  
A F D A - C S

## Local similarity

- as the Smith-Waterman Algorithm is so important, there are many implementations available
- the most efficient implementations make use of SIMD-instructions for vectorized computations (originally developed for image processing)
- instructions performed in parallel on short vectors
- widely available standards: AVX2 on X86-CPU's (vectors have 256 bits); Neon on M1-silicon (Mac): vectors have 128 bits



SIMD instructions provide potential speedup of 8/4 for arithmetic ops on int/double-values

# Local similarity

## Example application

- Given:  $k = 2\,550$  protein sequences; total length  $n = 997\,571$
- perform all-against-all comparison: compare sequence  $i$  and  $j$  for all  $0 \leq i \leq k - 2$  and  $i + 1 \leq j \leq k - 1$
- $\Rightarrow \frac{k \cdot (k-1)}{2} = 3\,186\,250$  pairwise sequence comparisons (using SW-alg.)
- total size of all matrices to be computed  $\approx 5 \cdot 10^{11}$
- computation time of software developed in the GI-group (and based on SSW<sup>a</sup>): 23 005 ms using 4 threads of an Intel i5 (3.2 GHz)
- $\frac{5 \cdot 10^{11}}{23\,005} = 21\,734\,405$  matrix entries per ms
- $\frac{3\,186\,250}{23\,005} \approx 138$  matrices per ms

---

<sup>a</sup>Zhao et. al. 2013