

345680000000120036938643906042007957000086-V

7566619000087056340010023976010

# Bioinformatics

**Sequence Analysis**  
**Sequence Alignment**

**Part 10**

**Mahdi Vasighi**

700043061-M

F5-223000010023661160900

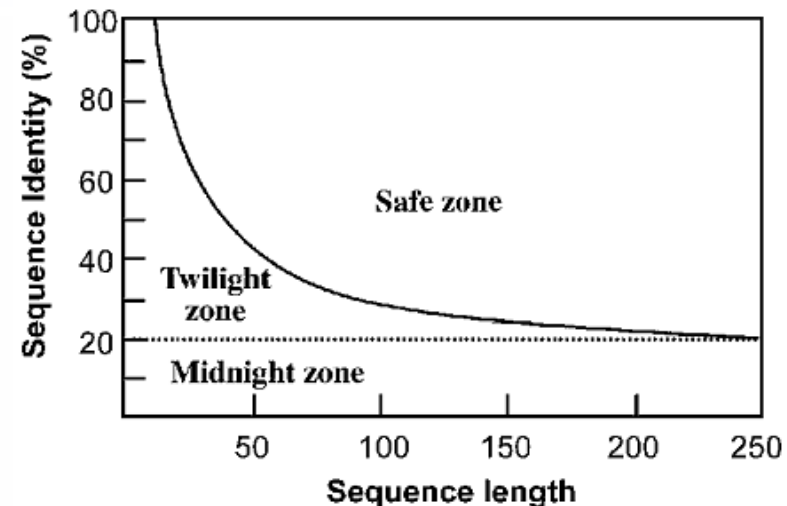
60619667700345-EH N1



# Pairwise Sequence Alignment

- An important concept in sequence analysis is **sequence homology**
  - A conclusion about a common ancestral relationship drawn from sequence similarity comparison
- A related but different term is **sequence similarity**, which is the percentage of aligned residues
  - A direct result of observation from the sequence alignment.
  - Can be quantified using percentages
  - similarity level → evolutionary relationship can be inferred

Sequence length is also a crucial factor. The shorter the sequence, the higher the chance that some alignment is attributable to random chance.





# Pairwise Sequence Alignment

- Sequence similarity vs. Sequence identity
  - Synonymous for nucleotide sequences.
  - In a protein sequence alignment, sequence identity refers to the percentage of matches of the same amino acid residues.
  - Sequence similarity refers to the percentage of aligned residues that have similar physicochemical characteristics.
- There are two ways to calculate the sequence similarity/identity:

$$S = [(L_s \times 2) / (L_a + L_b)] \times 100$$

$$I = [(L_i \times 2) / (L_a + L_b)] \times 100$$



# Pairwise Sequence Alignment

## Dynamic programming

Dynamic programming is a computational method that can be applied only to problems exhibiting the properties of overlapping subproblems. Examples include:

- Travelling salesman problem
- finding the best chess move

The method is very important for sequence analysis because it provides the **very best** or **optimal alignment** between sequences.

**Idea:** Build up an optimal alignment using previous solutions for optimal alignments of smaller subsequences.



# Pairwise Sequence Alignment

Dynamic programming

**Global Alignment : Needleman-Wunsch technique**

The following is an example of global sequence alignment using **Needleman-Wunsch** technique:

Seq.1:    G A A T T C A G T T A  
Seq.2:    G G A T C G A

A simple scoring scheme is assumed:

Status	Score
Match	+1
Mismatch	0
Indel (gap)	0



# Pairwise Sequence Alignment

## Dynamic programming

### Global Alignment : Needleman-Wunsch technique

Needleman-Wunsch technique has three steps:

1. Initialization
2. Matrix fill (scoring)
3. Traceback (alignment)

The first step in the global alignment dynamic programming approach is to create a matrix with  $M + 1$  columns and  $N + 1$  rows where  $M$  and  $N$  correspond to the size of the sequences to be aligned.



# Pairwise Sequence Alignment

Dynamic programming

**Global Alignment : Needleman-Wunsch technique**

## 1. Initialization

	G	A	A	T	T	C	A	G	T	T	A
	0	0	0	0	0	0	0	0	0	0	0
G	0										
G	0										
A	0										
T	0										
C	0										
G	0										
A	0										

This element should always be set to zero

These elements are set according to gap penalty  
Gap x 1  
Gap x 2  
Gap x 3  
...

$i = 4$   
 $j = 6$

$i = 9$   
 $j = 2$



# Pairwise Sequence Alignment

$$M_{i,j} = \text{MAXIMUM} [ M_{i-1,j-1} + S_{i,j}, M_{i,j-1} + w, M_{i-1,j} + w ]$$

Dynamic programming

## Global Alignment : Needleman-Wunsch technique

### 2. Matrix fill (scoring)

One possible solution of the matrix fill step finds the maximum global alignment score by starting in the upper left hand corner in the matrix and finding the maximal score  $M_{i,j}$  for each position in the matrix.

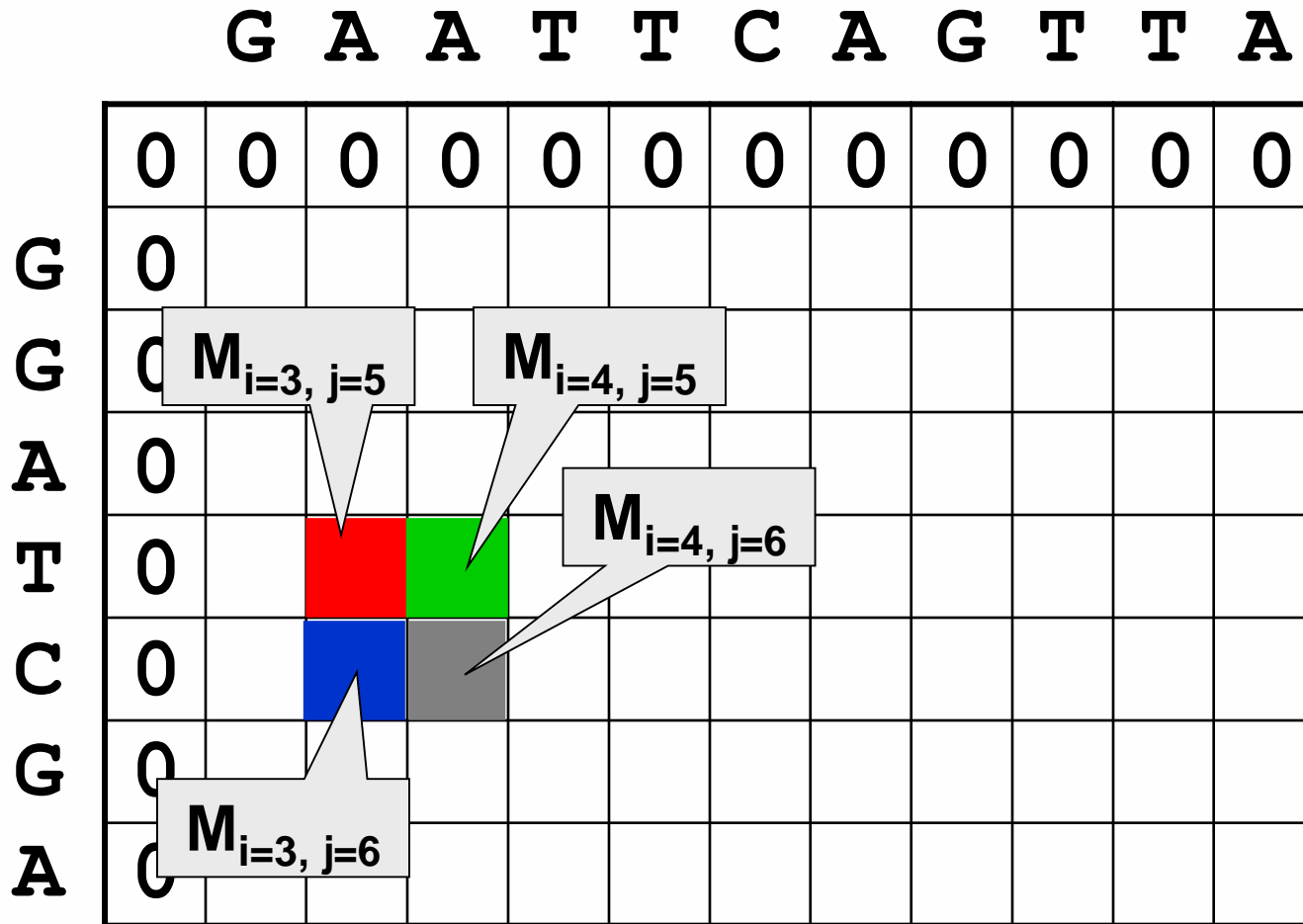
$$M_{i,j} = \text{MAXIMUM} [ M_{i-1,j-1} + S_{i,j}, M_{i,j-1} + w, M_{i-1,j} + w ]$$

match/mismatch in the diagonal      gap in sequence #1      gap in sequence #2



# Global Alignment : Needleman-Wunsch technique

$$M_{i,j} = \text{MAXIMUM} [ M_{i-1,j-1} + S_{i,j}, M_{i,j-1} + w, M_{i-1,j} + w ]$$



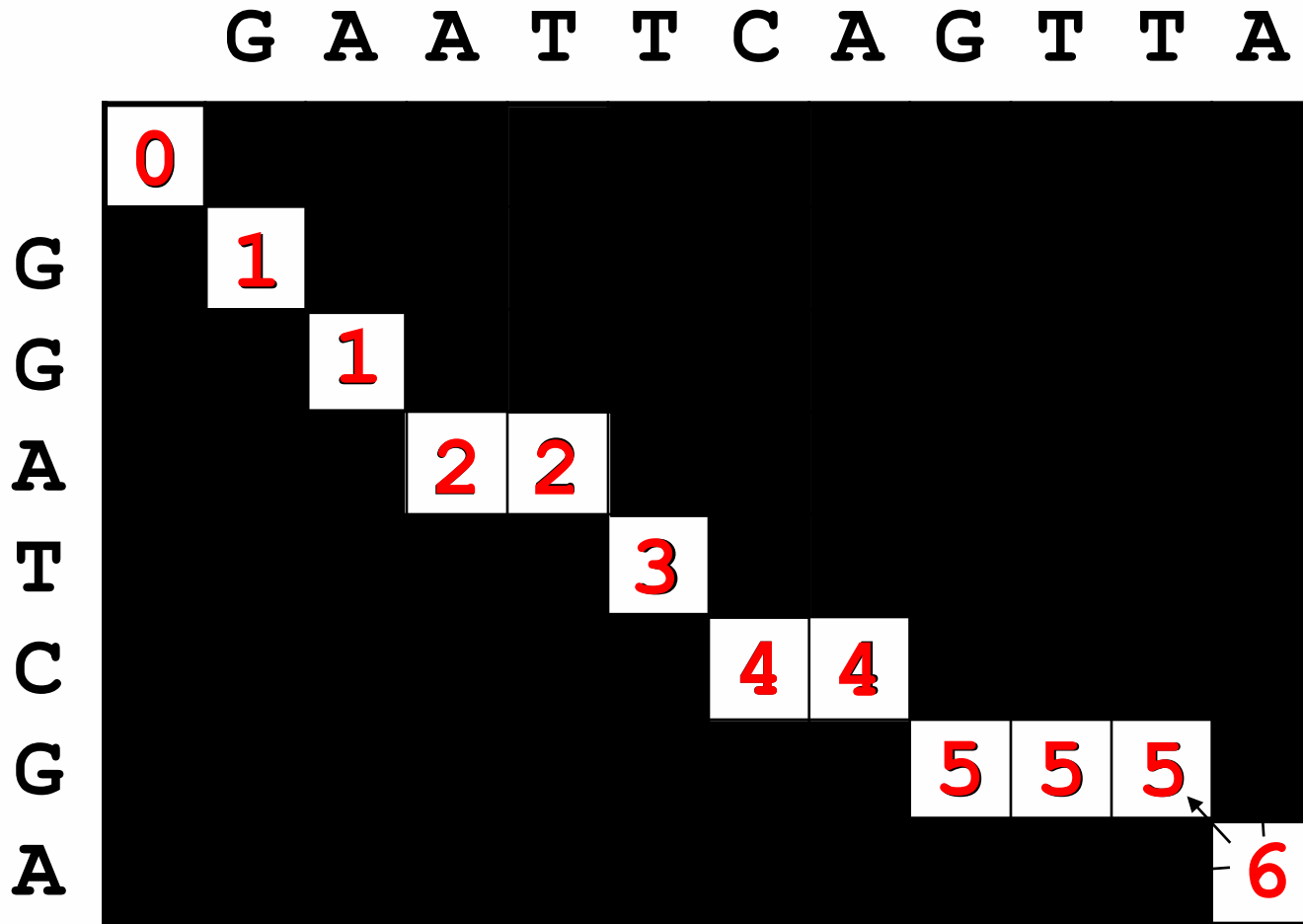


# Global Alignment : Needleman-Wunsch technique

the maximum alignment score for the two test sequences is 6. The traceback step determines the actual alignment(s) that result in the maximum score.

		G	A	A	T	T	C	A	G	T	T	A
	0	0	0	0	0	0	0	0	0	0	0	0
G	0	1	1	1	1	1	1	1	1	1	1	1
G	0	1	1	1	1	1	1	1	2	2	2	2
A	0	1	2	2	2	2	2	2	2	2	2	3
T	0	1	2	2	3	3	3	3	3	3	3	3
C	0	1	2	2	3	3	4	4	4	4	4	4
G	0	1	2	2	3	3	4	4	5	5	5	5
A	0	1	2	3	3	3	4	5	5	5	5	6

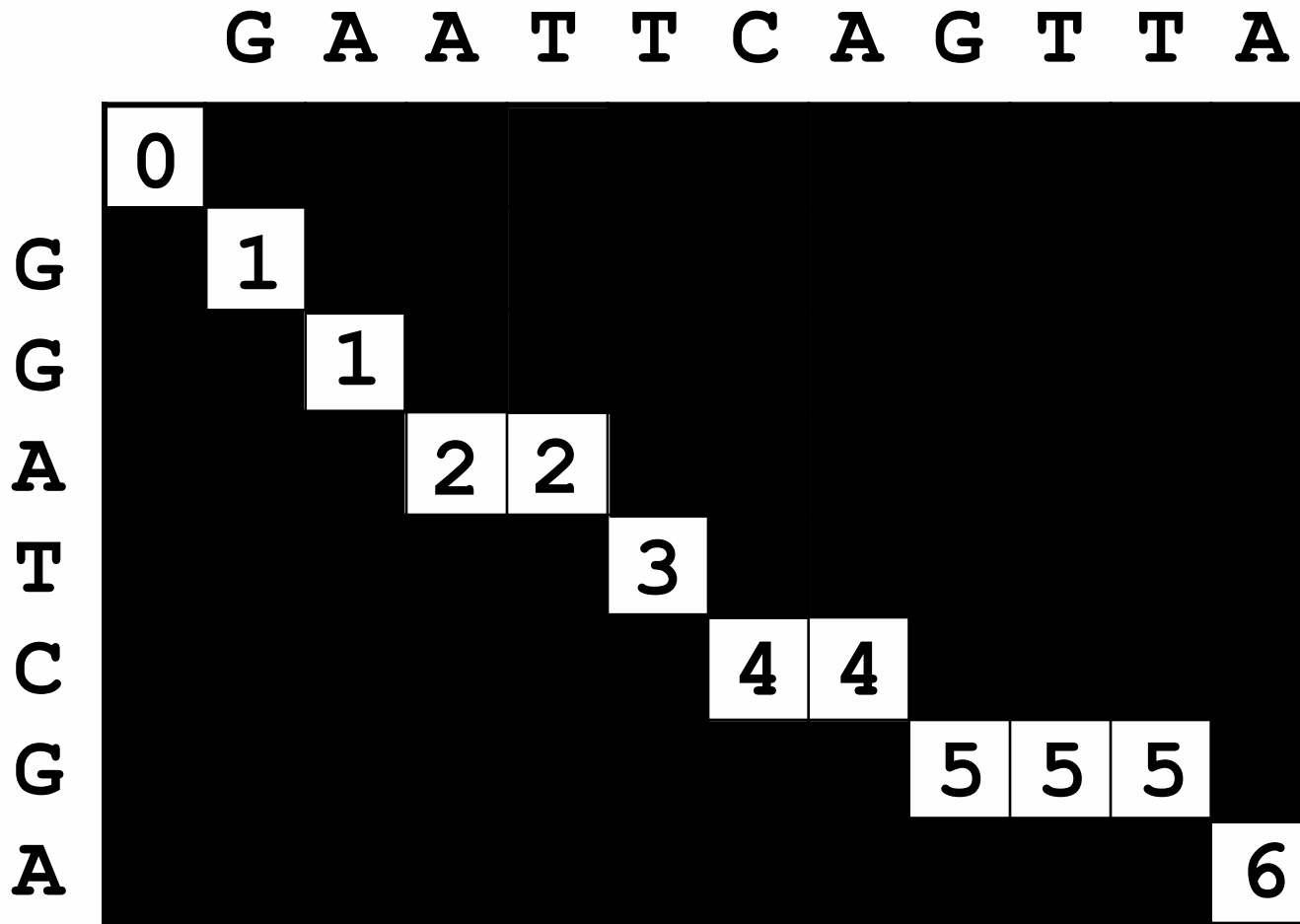
# Global Alignment : Needleman-Wunsch technique



G	A	A	T	T	C	A	G	T	T	A	Seq (1)
G	G	A	-	T	C	-	G	-	-	A	Seq (2)

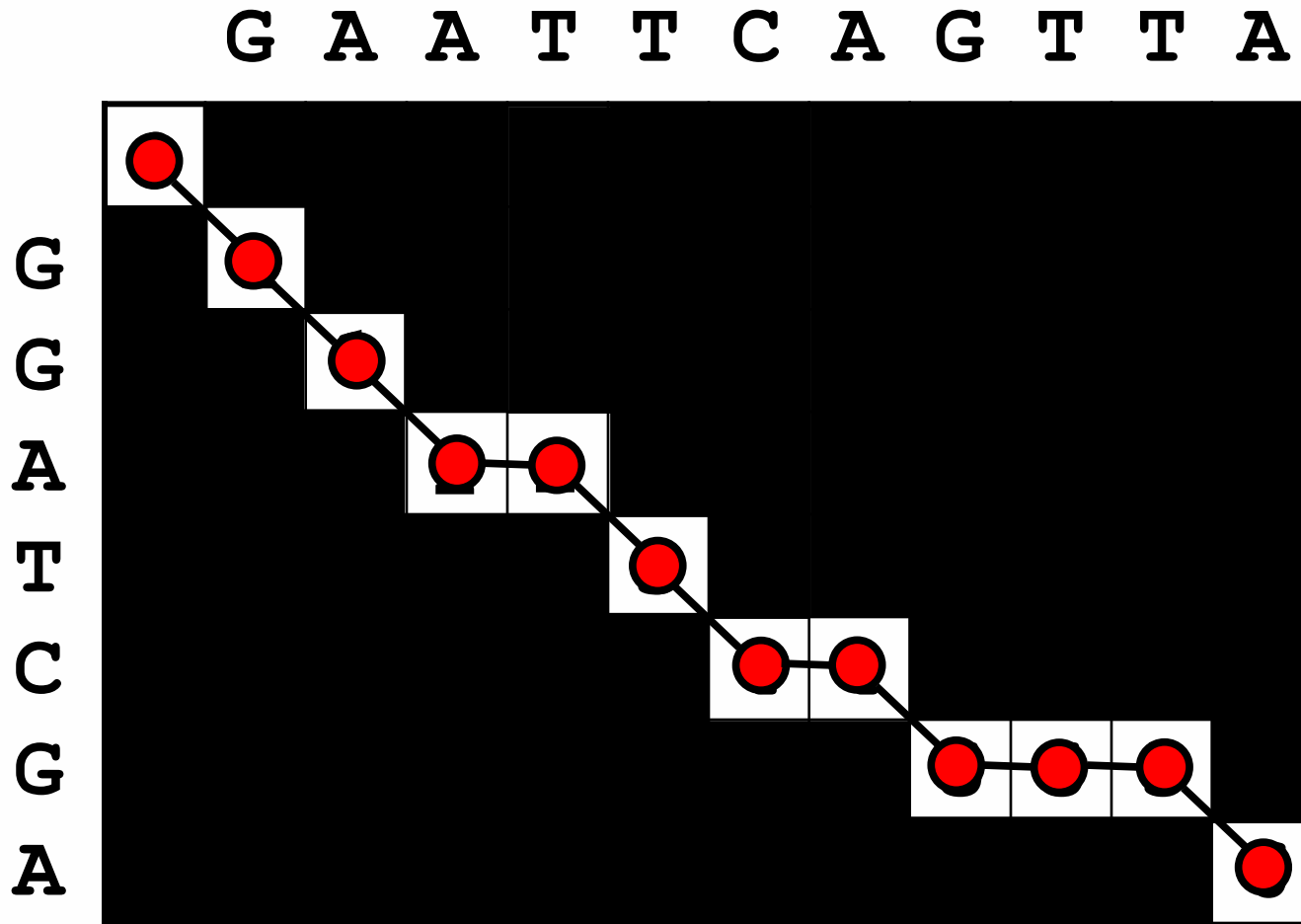
# Global Alignment : Needleman-Wunsch technique

There are more alternative solutions each resulting in a maximal global alignment score of 6.



G	-	A	A	T	T	C	A	G	T	T	A	Seq (1)
G	G	-	A	-	T	C	-	G	-	-	A	Seq (2)

# Global Alignment : Needleman-Wunsch technique



G	-	A	A	T	T	C	A	G	T	T	A	Seq (1)
G	G	-	A	-	T	C	-	G	-	-	A	Seq (2)

## Global Alignment : Needleman-Wunsch technique

Here is another example of global sequence alignment by Needleman-Wunsch technique using a different scoring scheme:

Seq. 1 :    G   A   A   T   T   C   A   G   T   T   A  
Seq. 2 :    G   G   A   T   C   G   A

The scoring scheme is:

Status	Score
Match	+2
Mismatch	-1
Indel (gap)	-2

# Global Alignment : Needleman-Wunsch technique

## 1. Initialization

		<b>G</b>	<b>A</b>	<b>A</b>	<b>T</b>	<b>T</b>	<b>C</b>	<b>A</b>	<b>G</b>	<b>T</b>	<b>T</b>	<b>A</b>
	0	-2	-4	-6	-8	-10	-12	-14	-16	-18	-20	-22
<b>G</b>	-2											
<b>G</b>	-4											
<b>A</b>	-6											
<b>T</b>	-8											
<b>C</b>	-10											
<b>G</b>	-12											
<b>A</b>	-14											

$$M_{i,j} = \text{MAXIMUM} [ M_{i-1,j-1} + S_{i,j}, M_{i,j-1} + w, M_{i-1,j} + w ]$$



# Global Alignment : Needleman-Wunsch technique

## 2. Matrix filling

**G A A T T C A G T T A**

	<b>0</b>	<b>-2</b>	<b>-4</b>	<b>-6</b>	<b>-8</b>	<b>-10</b>	<b>-12</b>	<b>-14</b>	<b>-16</b>	<b>-18</b>	<b>-20</b>	<b>-22</b>
<b>G</b>	<b>-2</b>	<b>2</b>										
<b>G</b>	<b>-4</b>											
<b>A</b>	<b>-6</b>											
<b>T</b>	<b>-8</b>											
<b>C</b>	<b>-10</b>											
<b>G</b>	<b>-12</b>											
<b>A</b>	<b>-14</b>											

$$M_{i-1, j-1} + S_{i,j}, M_{i,j-1} + w, M_{i-1, j} + w$$

$$M_{1,1} = \text{MAXIMUM} [ 0 + 2, -2 - 2, -2 - 2 ]$$

# Global Alignment : Needleman-Wunsch technique

## 2. Matrix filling

		<b>G</b>	<b>A</b>	<b>A</b>	<b>T</b>	<b>T</b>	<b>C</b>	<b>A</b>	<b>G</b>	<b>T</b>	<b>T</b>	<b>A</b>
<b>G</b>	0	-2	-4	-6	-8	-10	-12	-14	-16	-18	-20	-22
<b>G</b>	-2	2										
<b>A</b>	-4	0										
<b>A</b>	-6											
<b>T</b>	-8											
<b>C</b>	-10											
<b>G</b>	-12											
<b>A</b>	-14											

$$M_{i-1, j-1} + S_{i,j}, \quad M_{i,j-1} + w, \quad M_{i-1, j} + w$$

$$M_{1,2} = \text{MAXIMUM} [ \quad -2 \quad +2 \quad , \quad 2 \quad -2 \quad , \quad -4 \quad -2 \quad ]$$

# Global Alignment : Needleman-Wunsch technique

## 2. Matrix filling

		<b>G</b>	<b>A</b>	<b>A</b>	<b>T</b>	<b>T</b>	<b>C</b>	<b>A</b>	<b>G</b>	<b>T</b>	<b>T</b>	<b>A</b>
<b>G</b>	0	-2	-4	-6	-8	-10	-12	-14	-16	-18	-20	-22
<b>G</b>	-2	<b>2</b>										
<b>G</b>	-4	<b>0</b>										
<b>A</b>	-6	<b>-2</b>										
<b>T</b>	-8											
<b>C</b>	-10											
<b>G</b>	-12											
<b>A</b>	-14											

$$M_{i-1, j-1} + S_{i,j}, \quad M_{i,j-1} + w, \quad M_{i-1, j} + w$$

$$M_{1,3} = \text{MAXIMUM} [ \text{-4 -1}, \quad \text{0 -2}, \quad \text{-6-2} ]$$

# Global Alignment : Needleman-Wunsch technique

## 2. Matrix filling

		<b>G</b>	<b>A</b>	<b>A</b>	<b>T</b>	<b>T</b>	<b>C</b>	<b>A</b>	<b>G</b>	<b>T</b>	<b>T</b>	<b>A</b>
<b>G</b>	0	-2	-4	-6	-8	-10	-12	-14	-16	-18	-20	-22
<b>G</b>	-2	<b>2</b>										
<b>G</b>	-4	<b>0</b>										
<b>A</b>	-6	<b>-2</b>										
<b>T</b>	-8	<b>-4</b>										
<b>C</b>	-10											
<b>G</b>	-12											
<b>A</b>	-14											

$$M_{i-1, j-1} + S_{i,j} , M_{i,j-1} + w , M_{i-1, j} + w$$

$$M_{1,4} = \text{MAXIMUM} [ -6 - 1 , -2 - 2 , -8 - 2 ]$$

# Global Alignment : Needleman-Wunsch technique

## 2. Matrix filling

	G	A	A	T	T	C	A	G	T	T	A
0	-2	-4	-6	-8	-10	-12	-14	-16	-18	-20	-22
G	-2	2	0	-2							
G	-4	0	1	-1							
A	-6	-2	2								
T	-8	-4	0								
C	-10	-6	-2								
G	-12	-8	-4								
A	-14	-10	-6								

$$M_{i-1, j-1} + S_{i,j} , M_{i, j-1} + w , M_{i-1, j} + w$$

$$M_{3,2} = \text{MAXIMUM} [ 0 \quad -1 , -2 - 2 , 1 - 2 ]$$

# Global Alignment : Needleman-Wunsch technique

## 2. Matrix filling

		G	A	A	T	T	C	A	G	T	T	A
	0	-2	-4	-6	-8	-10	-12	-14	-16	-18	-20	-22
G	-2	2	0	-2	-4	-6	-8	-10	-12	-14	-16	-18
G	-4	0	1	-1	-3	-5	-7	-9	-8	-10	-12	-14
A	-6	-2	2	3	1	-1	-3	-5	-7	-9	-11	-10
T	-8	-4	0	1	5	3	1	-1	-3	-5	-7	-9
C	-10	-6	-2	-1	3	4	5	3	1	-1	-3	-5
G	-12	-8	-4	-3	1	2	3	4	5	3	1	-1
A	-14	-10	-6	-2	1	0	1	5	3	4	2	3

After the matrix fill step, the maximum global alignment score for the two sequences is 3.

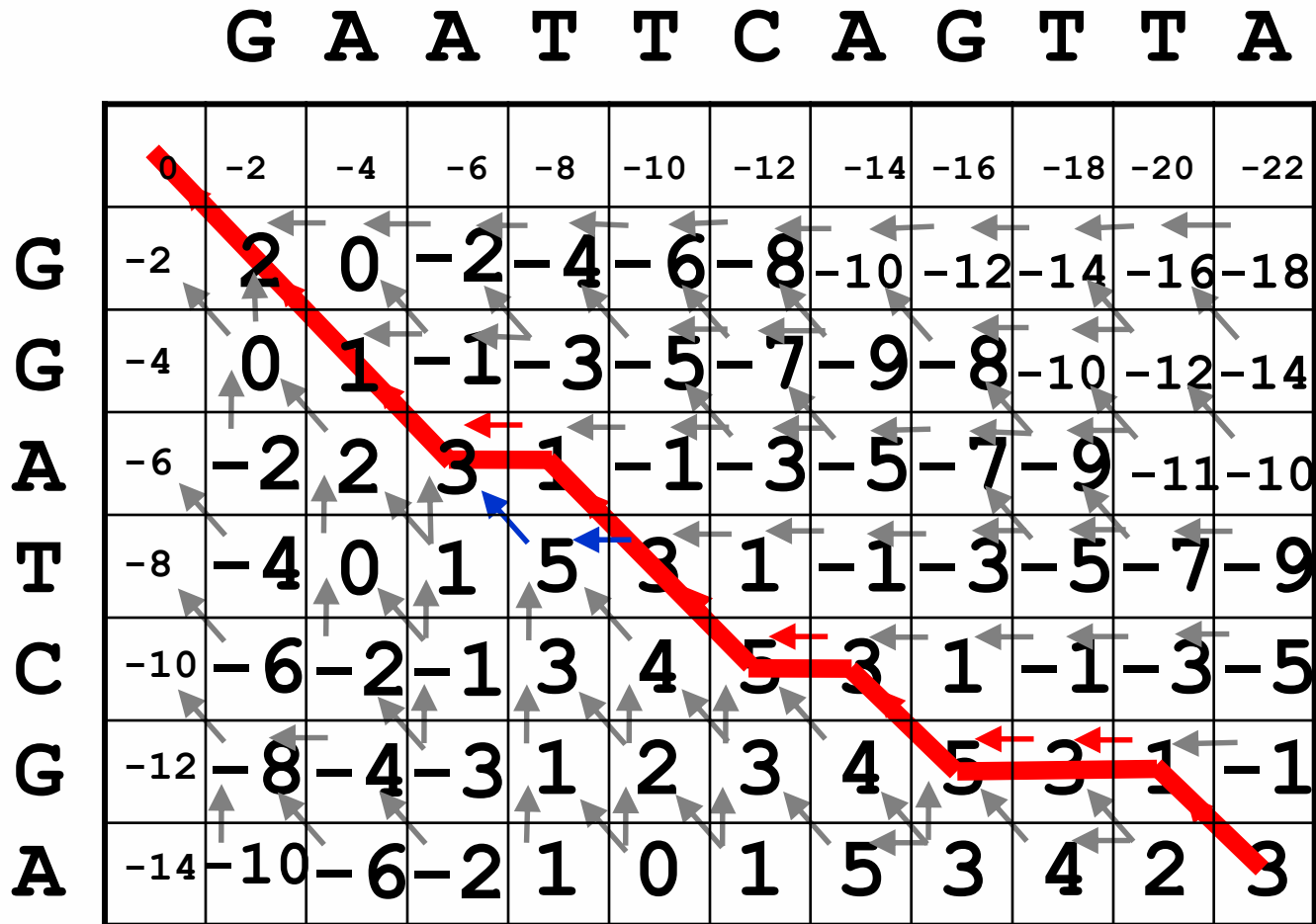
# Global Alignment : Needleman-Wunsch technique

## 3. Traceback

		G	A	A	T	T	C	A	G	T	T	A
	0	-2	-4	-6	-8	-10	-12	-14	-16	-18	-20	-22
G	-2	2	0	-2	-4	-6	-8	-10	-12	-14	-16	-18
G	-4	0	1	-1	-3	-5	-7	-9	-8	-10	-12	-14
A	-6	-2	2	3	1	-1	-3	-5	-7	-9	-11	-10
T	-8	-4	0	1	5	3	1	-1	-3	-5	-7	-9
C	-10	-6	-2	-1	3	4	5	3	1	-1	-3	-5
G	-12	-8	-4	-3	1	2	3	4	5	3	1	-1
A	-14	-10	-6	-2	1	0	1	5	3	4	2	3

# Global Alignment : Needleman-Wunsch technique

## 3. Traceback

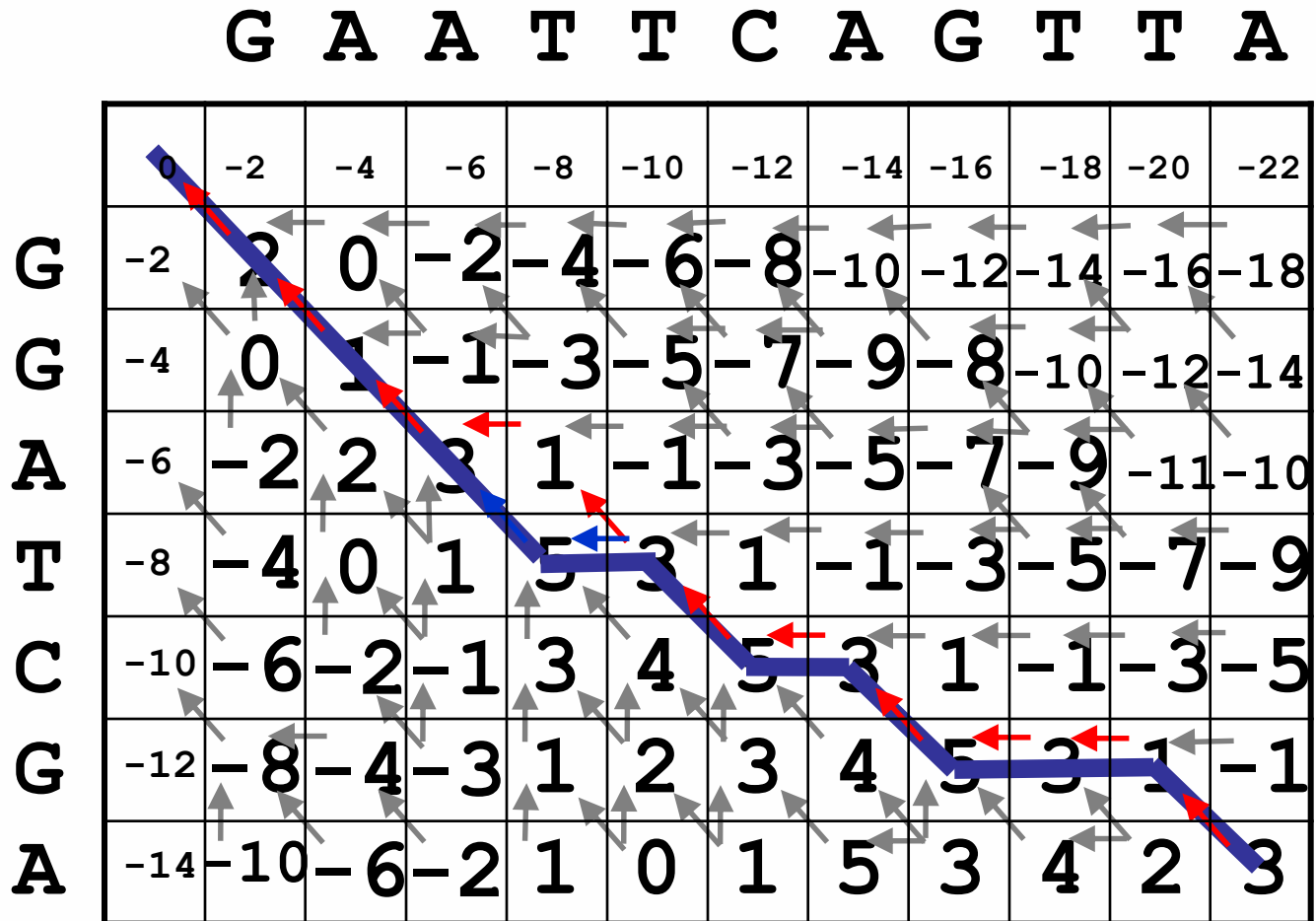


G	A	A	T	T	C	A	G	T	T	A
G	G	A	-	T	C	-	G	-	-	A



# Global Alignment : Needleman-Wunsch technique

## 3. Traceback



<b>G</b>	<b>A</b>	<b>A</b>	<b>T</b>	<b>T</b>	<b>C</b>	<b>A</b>	<b>G</b>	<b>T</b>	<b>T</b>	<b>A</b>
<b>G</b>	<b>G</b>	<b>A</b>	<b>T</b>	<b>-</b>	<b>C</b>	<b>-</b>	<b>G</b>	<b>-</b>	<b>-</b>	<b>A</b>

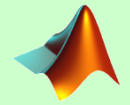
# Global Alignment : Needleman-Wunsch technique

**G A A T T C A G T T A**

	0	-2	-4	-6	-8	-10	-12	-14	-16	-18	-20	-22
<b>G</b>	-2	<b>2</b>	<b>0</b>	<b>-2</b>	<b>-4</b>	<b>-6</b>	<b>-8</b>	-10	-12	-14	-16	-18
<b>G</b>	-4	<b>0</b>	<b>1</b>	<b>-1</b>	<b>-3</b>	<b>-5</b>	<b>-7</b>	<b>-9</b>	<b>-8</b>	-10	-12	-14
<b>A</b>	-6	<b>-2</b>	<b>2</b>	<b>3</b>	<b>1</b>	<b>-1</b>	<b>-3</b>	<b>-5</b>	<b>-7</b>	<b>-9</b>	-11	-10
<b>T</b>	-8	<b>-4</b>	<b>0</b>	<b>1</b>	<b>5</b>	<b>3</b>	<b>1</b>	<b>-1</b>	<b>-3</b>	<b>-5</b>	<b>-7</b>	<b>-9</b>
<b>C</b>	-10	<b>-6</b>	<b>-2</b>	<b>-1</b>	<b>3</b>	<b>4</b>	<b>5</b>	<b>3</b>	<b>1</b>	<b>-1</b>	<b>-3</b>	<b>-5</b>
<b>G</b>	-12	<b>-8</b>	<b>-4</b>	<b>-3</b>	<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>	<b>5</b>	<b>3</b>	<b>1</b>	<b>-1</b>
<b>A</b>	-14	<b>-10</b>	<b>-6</b>	<b>-2</b>	<b>1</b>	<b>0</b>	<b>1</b>	<b>5</b>	<b>3</b>	<b>4</b>	<b>2</b>	<b>3</b>

# Global Alignment : Needleman-Wunsch technique

## nwalign



**Purpose:** Globally align two sequences using the Needleman-Wunsch algorithm

### Syntax:

```
[Score, Alignment] = nwalign(Seq1, Seq2, 'PropertyName',  
PropertyValue)
```

Values are 'PAM40', 'PAM250',  
'BLOSUM30' 'BLOSUM62', or  
'BLOSUM100'.

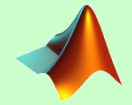
```
nwalign(..., 'ScoringMatrix', ScoringMatrixValue)  
nwalign(..., 'Alphabet', AlphabetVlaue)
```

Property to select the type of sequence. Value is  
either 'AA' or 'NT'. The default value is 'AA'.

```
nwalign(..., 'Showscore', ShowscoreValue)
```

# Global Alignment : Needleman-Wunsch technique

## nwalign



### Example

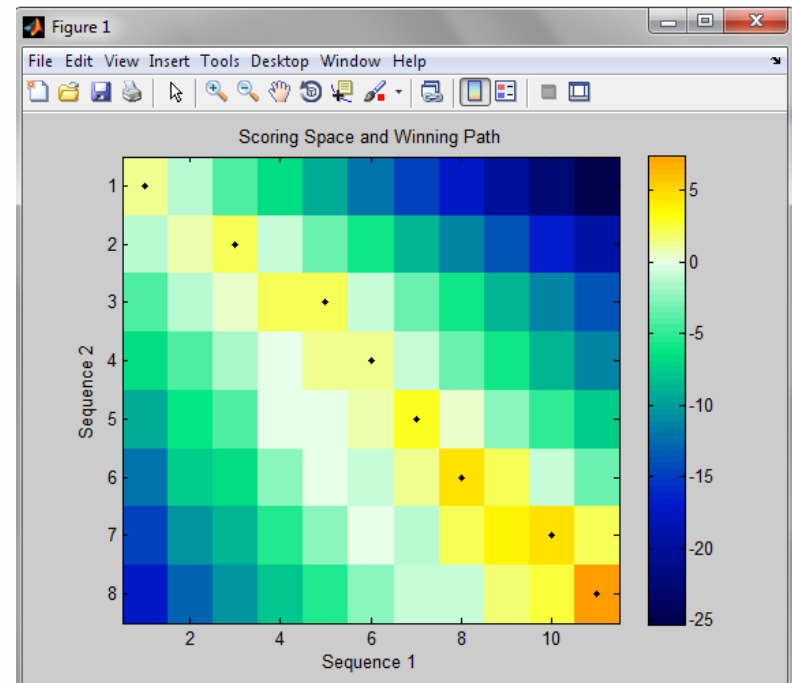
```
[Score, Alignment] = nwalign('IGRHRYHIGG', 'SRYIGRG', ...  
                             'scoringmatrix', 'pam250')
```

```
[Score, Alignment] = nwalign('VSPAGMASGYD', 'IPGKASYD', 'Showscore', 1)
```

```
Score =  
  
7.3333
```

```
Alignment =
```

```
VSPAGMASGYD  
: | | | | |  
I-P-GKAS-YD
```





# Pairwise Sequence Alignment

## Dynamic programming

### Local Alignment : Smith-Waterman technique

The following is an example of local sequence alignment using **Smith-Waterman** technique:

Seq.1:    A T C A G A G T C  
Seq.2:    G T C A G T C A

A simple scoring scheme is assumed:

Status	Score
Match	+2
Mismatch	-1
Indel (gap)	-2



# Local Alignment : Smith-Waterman technique

## 2. Matrix fill (scoring)

The main difference to the Needleman-Wunsch algorithm is that negative scoring matrix cells are set to zero, which renders the (thus positively scoring) local alignments visible.

$$M_{i,j} = \text{MAXIMUM} [ M_{i-1,j-1} + S_{i,j}, M_{i,j-1} + w, M_{i-1,j} + w, 0 ]$$

match/mismatch  
in the diagonal

gap in  
sequence #1

gap in  
sequence #2





$$M_{i-1, j-1} + S_{i,j}, M_{i,j-1} + w, M_{i-1, j} + w$$

$$H_{i,j} = \text{MAXIMUM} [ 0 + 2, 0 - 2, 0 - 2, 0 ]$$

	A	T	C	A	G	A	G	T	C
G	0	0	0	0	0	0	0	0	0
T	0	0	2						
C	0	0							
A	0	2							
G	0	0							
T	0	0							
C	0	0							
A	0	2							

$$M_{i-1, j-1} + S_{i,j}, M_{i,j-1} + w, M_{i-1,j} + w$$

$$H_{i,j} = \text{MAXIMUM} [ 0 \quad -1, 2 \quad -2, 0 \quad -2, 0 ]$$

	A	T	C	A	G	A	G	T	C
	0	0	0	0	0	0	0	0	0
G	0	0	0						
T	0	0	2						
C	0	0	0						
A	0	2	0						
G	0	0							
T	0	0							
C	0	0							
A	0	2							

$$H_{i,j} = \text{MAXIMUM} [ \overset{\text{red}}{M_{i-1,j-1} + S_{i,j}}, \overset{\text{green}}{M_{i,j-1} + w}, \overset{\text{blue}}{M_{i-1,j} + w}, 0 ]$$

	A	T	C	A	G	A	G	T	C
G	0	0	0	0	0	0	0	0	0
T	0	0	2	0	0	0	1	0	4
C	0	0	0	4	2	0	0	0	2
A	0	2	0	2	6	4	2	0	0
G	0	0	1	0	4	8	6	4	2
T	0	0	2	0	2	6	7	5	6
C	0	0	0	4	2	4	5	6	4
A	0	2	0	2	6	4	6	4	5

To obtain the optimum local alignment, we start with the highest value in the matrix (i,j). Then, we go backwards to one of positions (i-1,j), (i,j-1), and (i-1,j-1) depending on the direction of movement used to construct the matrix. We keep the process until we reach a matrix cell with zero value

	A	T	C	A	G	A	G	T	C
G	0	0	0	0	0	0	0	0	0
T	0	0	2	0	0	0	1	0	4
C	0	0	0	4	2	0	0	0	2
A	0	2	0	2	6	4	2	0	0
G	0	0	1	0	4	8	6	4	2
T	0	0	2	0	2	6	7	5	6
C	0	0	0	4	2	4	5	6	4
A	0	2	0	2	6	4	6	4	5

A T C A G A G T C

G T C A G T C A

A T C A G A G T C

	0	0	0	0	0	0	0	0	0
G	0	0	0	0	0	2	0	2	0
T	0	0	2	0	0	0	1	0	4
C	0	0	0	4	2	0	0	0	2
A	0	2	0	2	6	4	2	0	0
G	0	0	1	0	4	8	6	4	2
T	0	0	2	0	2	6	7	5	6
C	0	0	0	4	2	4	5	6	4
A	0	2	0	2	6	4	6	4	5

**A T C A G A G T C**

**G T C A G T C A**

**A T C A G A G T C**

	0	0	0	0	0	0	0	0	0
<b>G</b>	0	0	0	0	0	2	0	2	0
<b>T</b>	0	0	2	0	0	0	1	0	4
<b>C</b>	0	0	0	4	2	0	0	0	2
<b>A</b>	0	2	0	2	6	4	2	0	0
<b>G</b>	0	0	1	0	4	8	6	4	2
<b>T</b>	0	0	2	0	2	6	7	5	6
<b>C</b>	0	0	0	4	2	4	5	6	4
<b>A</b>	0	2	0	2	6	4	6	4	5

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

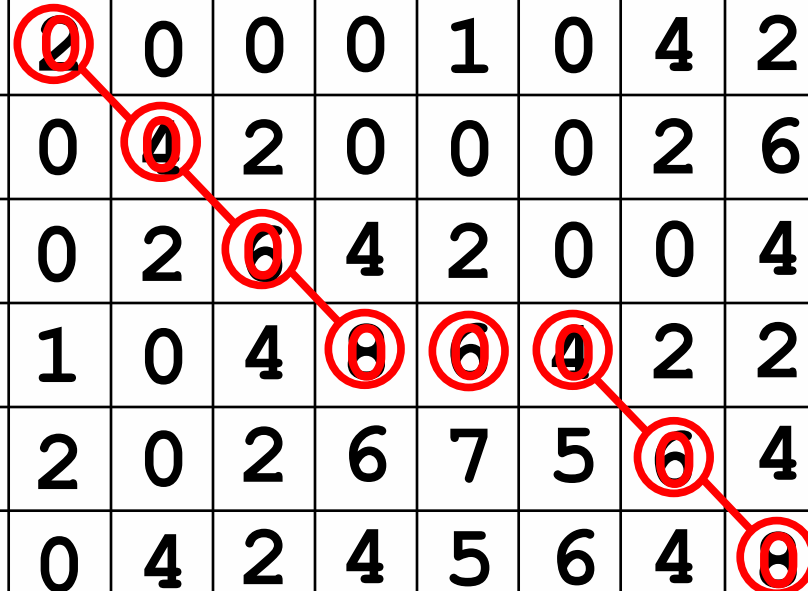
**A T C A G A G T C**

	0	0	0	0	0	0	0	0	0
<b>G</b>	0	0	0	0	0	2	0	2	0
<b>T</b>	0	0	2	0	0	0	1	0	4
<b>C</b>	0	0	0	4	2	0	0	0	2
<b>A</b>	0	2	0	2	6	4	2	0	0
<b>G</b>	0	0	1	0	4	8	6	4	2
<b>T</b>	0	0	2	0	2	6	7	5	6
<b>C</b>	0	0	0	4	2	4	5	6	4
<b>A</b>	0	2	0	2	6	4	6	4	5

# Finding the sub-optimal alignments

1. fix optimal alignment positions to zero

		A	T	C	A	G	A	G	T	C
	0	0	0	0	0	0	0	0	0	0
G	0	0	0	0	0	2	0	2	0	0
T	0	0	0	0	0	0	1	0	4	2
C	0	0	0	0	2	0	0	0	2	6
A	0	2	0	2	0	4	2	0	0	4
G	0	0	1	0	4	0	0	0	2	2
T	0	0	2	0	2	6	7	5	0	4
C	0	0	0	4	2	4	5	6	4	0
A	0	2	0	2	6	4	6	4	5	6





# Finding the sub-optimal alignments

## 2. Recalculate the matrix elements again

	A	T	C	A	G	A	G	T	C
G	0	0	0	0	0	0	0	0	0
T	0	0	0	0	0	2	0	2	0
C	0	0	0	0	0	0	1	0	4
A	0	0	0	0	2	0	0	0	2
G	0	2	0	2	0	4	2	0	0
T	0	0	1	0	4	0	0	0	4
C	0	0	2	0	2	6	7	5	0
A	0	0	0	4	2	4	5	6	4
G	0	2	0	2	6	4	6	4	5
T	0	2	0	2	6	4	6	4	5
C	0	2	0	2	6	4	6	4	5
A	0	2	0	2	6	4	6	4	5

# Finding the sub-optimal alignments

## 2. Recalculate the matrix elements again

		A	T	C	A	G	A	G	T	C
	G	0	0	0	0	0	0	0	0	0
	T	0	0	0	0	0	2	0	2	0
	C	0	0	0	0	0	0	0	2	6
	A	0	2	0	0	0	2	0	0	4
	G	0	0	1	0	0	0	0	0	2
	T	0	0	2	0	0	0	0	0	0
	C	0	0	0	4	2	0	0	0	0
	A	0	2	0	2	6	4	2	0	0

# Finding the sub-optimal alignments

## 2. Recalculate the matrix elements again

	A	T	C	A	G	A	G	T	C
G	0	0	0	0	0	0	0	0	0
T	0	0	0	0	0	2	0	2	0
C	0	0	0	0	0	0	0	0	0
A	0	0	0	0	0	1	0	4	2
G	0	0	0	0	0	0	0	2	6
T	0	2	0	0	0	0	0	0	4
C	0	0	1	0	0	0	0	0	2
A	0	0	2	0	0	0	0	0	0
G	0	0	0	4	2	0	0	0	0
T	0	0	0	4	2	0	0	0	0
C	0	0	0	2	0	0	0	0	0
A	0	2	0	2	6	4	2	0	0

→ New values

A T C A G A G T C

G T C A G T C A

A T C A G A G T C

	0	0	0	0	0	0	0	0	0
G	0	0	0	0	0	2	0	2	0
T	0	0	0	0	0	0	1	0	4
C	0	0	0	0	0	0	0	0	2
A	0	2	0	0	0	0	2	0	0
G	0	0	1	0	0	0	0	0	2
T	0	0	2	0	0	0	0	0	0
C	0	0	0	4	2	0	0	0	0
A	0	2	0	2	6	4	2	0	0

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

	0	0	0	0	0	0	0	0	0	0
G	0	0	0	0	0	2	0	<span style="border: 1px solid green; border-radius: 50%; padding: 2px;">2</span>	0	0
T	0	0	0	0	0	0	1	0	<span style="border: 1px solid green; border-radius: 50%; padding: 2px;">4</span>	2
C	0	0	0	0	0	0	0	0	2	<span style="border: 1px solid green; border-radius: 50%; padding: 2px;">6</span>
A	0	2	0	0	0	0	2	0	0	4
G	0	0	1	0	0	0	0	0	0	2
T	0	0	<span style="border: 1px solid yellow; border-radius: 50%; padding: 2px;">2</span>	0	0	0	0	0	0	0
C	0	0	0	<span style="border: 1px solid yellow; border-radius: 50%; padding: 2px;">4</span>	2	0	0	0	0	0
A	0	2	0	2	<span style="border: 1px solid yellow; border-radius: 50%; padding: 2px;">6</span>	4	2	0	0	0

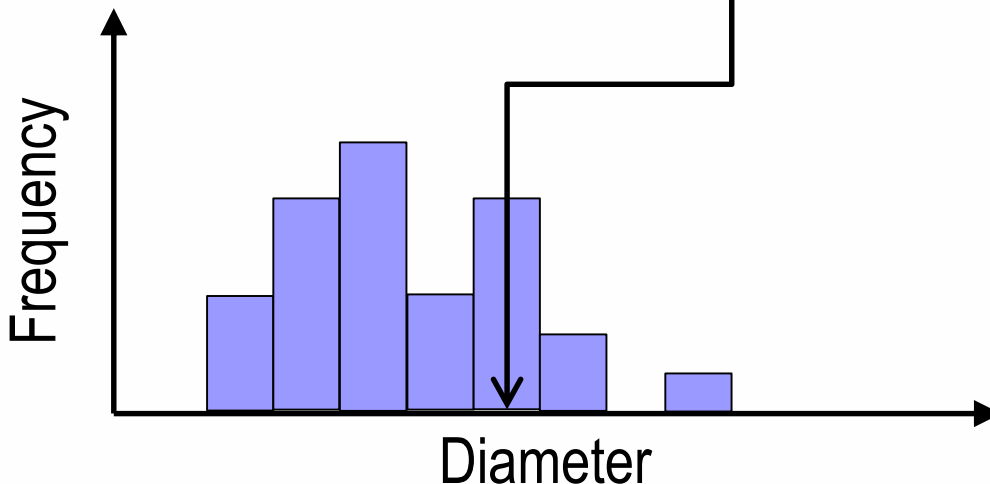


# Significance of Sequence Alignments

Assessing significance requires a *distribution*



I have an apple of diameter 10 cm. Is that unusual?





# Significance of Sequence Alignments

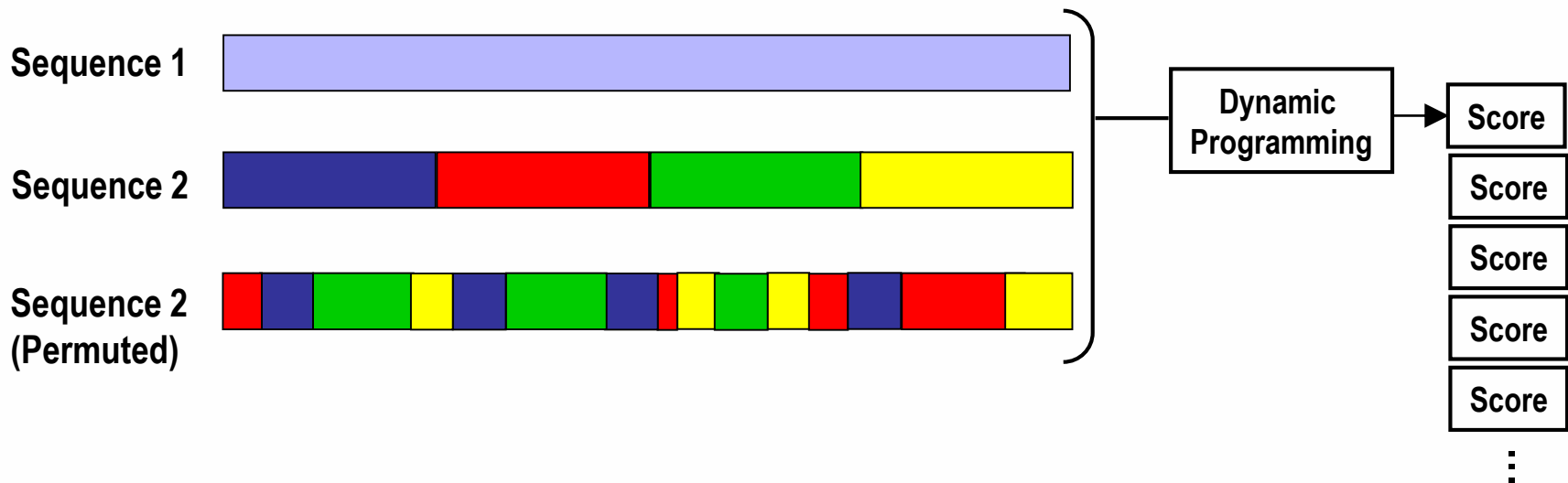
- When given a sequence alignment showing a certain degree of similarity, it is often important to ask whether the observed sequence alignment can occur by random chance or the alignment is indeed statistically sound.
- The truly statistically significant sequence alignment will be able to provide evidence of homology between the sequences involved.
- Solving this problem requires a distribution of the alignment scores of two unrelated sequences of the same length.



# Significance of Sequence Alignments

A practical approach to the problem is as follows:

We may randomize one of the sequences, many times, re-align each result to the second sequence (held fixed), and collect the distribution of resulting scores.



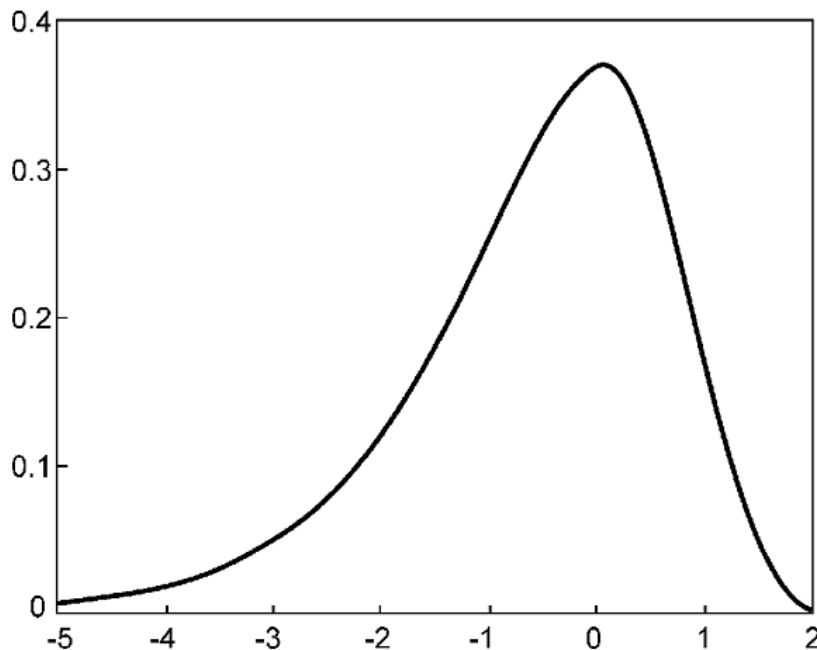




# Significance of Sequence Alignments

Many studies have demonstrated that the distribution of similarity scores assumes a peculiar shape that resembles a highly skewed normal distribution with a long tail on one side

## Gumble extreme value distribution



The distribution can be expressed as:

$$P = 1 - e^{-Kmn} e^{-\lambda x}$$

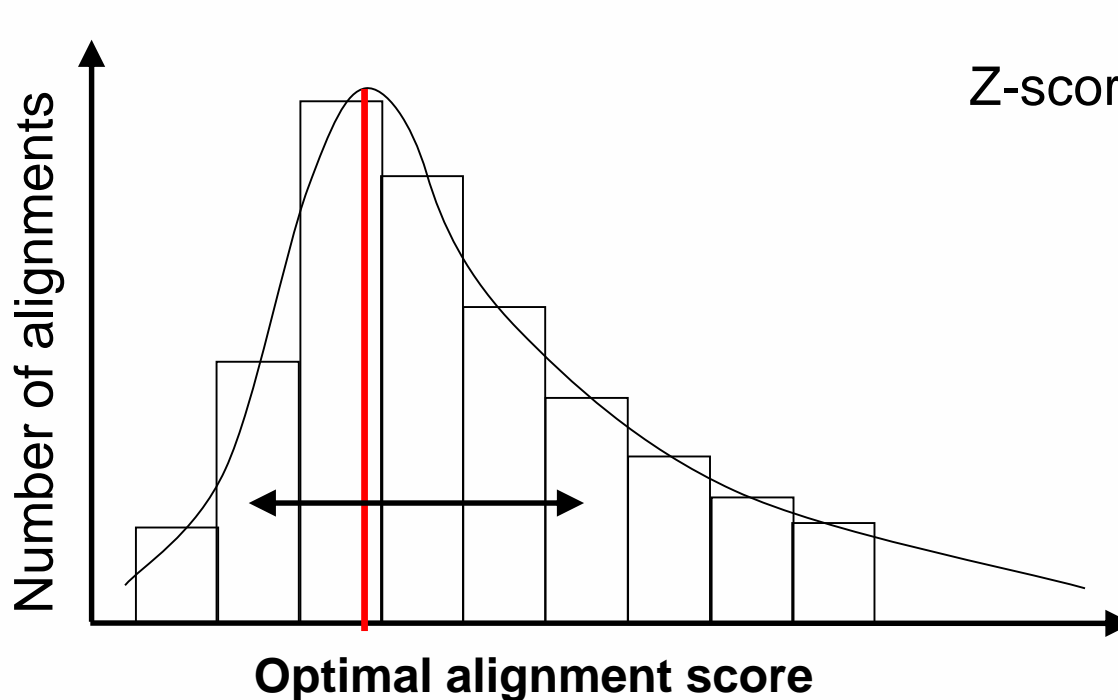
Where **m** and **n** are the sequence lengths,  **$\lambda$**  is a scaling factor for the scoring matrix used, and **K** is a constant that depends on the scoring matrix and gap penalty combination that is used.

A **P-value** is given to indicate the probability that the original alignment is due to random chance.



# Significance of Sequence Alignments

A population of alignment score will be obtained and z-score can be calculated as follow:



$$Z\text{-score} = \frac{\text{Score} - \text{mean}}{\text{Standard deviation}}$$

- The Z-score reflects the extent to which the original result is an outlier from the population.
- Experience suggests that Z-scores  $\geq 5$  are significant.



# Significance of Sequence Alignments

[http://www.ch.embnet.org/software/PRSS\\_form.html](http://www.ch.embnet.org/software/PRSS_form.html)



PRSS

A frequently used software program for assessing statistical significance of a pairwise alignment is the **PRSS3** (Probability of Random Shuffles) program.

## PRSS3 - evaluates the significance of a protein sequence alignment

prss3 is used to evaluate the significance of a protein or DNA sequence similarity score by comparing two sequences and calculating optimal similarity scores, and then repeatedly shuffling the second sequence, and calculating optimal similarity scores using the Smith-Waterman algorithm. An extreme value distribution is then fit to the shuffled-sequence scores. The characteristic parameters of the extreme value distribution are then used to estimate the probability that each of the unshuffled sequence scores would be obtained by chance in one sequence, or in a number of sequences equal to the number of shuffles.

This program is derived from rdf2, described by Pearson and Lipman, PNAS (1988) 85:2444-2448, and Pearson (Meth. Enz. 183:63-98). Use of the extreme value distribution for estimating the probabilities of similarity scores was described by Altshul and Karlin, PNAS (1990) 87:2264-2268. The 'z-values' calculated by rdf2 are not as informative as the P-values and expectations calculated by prdf. prss3 uses calculates optimal scores using the same rigorous Smith-Waterman algorithm (Smith and Waterman, J. Mol. Biol. (1983) 147:195-197) used by the ssearch3 program.

prss3 also allows a more sophisticated shuffling method: residues can be shuffled within a local window, so that the order of residues 1-10, 11-20, etc, is destroyed but a residue in the first 10 is never swapped with a residue outside the first ten, and so on for each local window.

This program is part of the FASTA package of sequence analysis program.

- **Usage:** Paste your two sequences in one of the supported **formats** into the sequence fields below and press the "Run PRSS" button.  
Make sure that both format buttons (next to the sequence fields) shows the correct formats

Number of shuffles : 200 ▾

Scoring matrix : default ▾

gap opening penalty: 12 gap extension penalty: 2

First sequence title (optional):

1. Globally align the **AK2H\_ECOLI** and **AK3\_ECOLI** using the toolbox and calculate similarity and identity..
2. How could the **divide and conquer** method be used to solve the global sequence alignment more efficiently?

