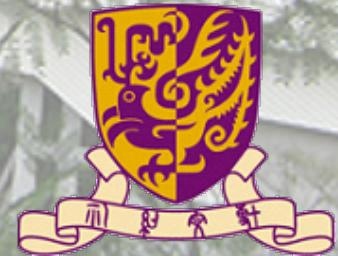


Tutorial 2. Optimal Sequence Alignment

The Chinese University of Hong Kong

BMEG3102 Bioinformatics

TA: Yizhen Chen



Agenda

- What is sequence alignment?
- Sequence data and databases
- Global alignment
- Local alignment
- Scoring matrix

What is sequence alignment?

- Biological problem:
 - Identify consistent patterns in DNA, RNA or protein sequences associated with genetic diseases or important biological functions (e.g. TF binding)
 - Evaluate the similarity of the sequences of different individuals or different species
- Computational problem:
 - Given 2 or more DNA, RNA or protein sequences, find the best alignment according to some scoring matrix, σ .

Sequence data and databases

- Sequence data
 - DNA: {A, C, G, T}
 - RNA: {A, C, G, U}
 - Protein: 20 amino acids
- Nucleotide databases
 - [GenBank](#)
 - [EMBL-EBI](#)
 - [DDBJ](#)
- Protein sequence database
 - [UniProt](#)

Dynamic Programming (DP)

- **Dynamic programming** is a method for simplifying a complicated problem by breaking it down into **simpler sub-problems** in a **recursive** manner, solving each of those subproblems just **once**, and **storing their solutions**.

Key Idea: Store the results of overlapping smaller sub-problems to avoid repetitive computations.

Global alignment

- Find the optimal alignment for the whole two sequences, r and s .
- Needleman-Wunsch algorithm
- Initialization: depend on the gap penalty
- DP Table Construction:

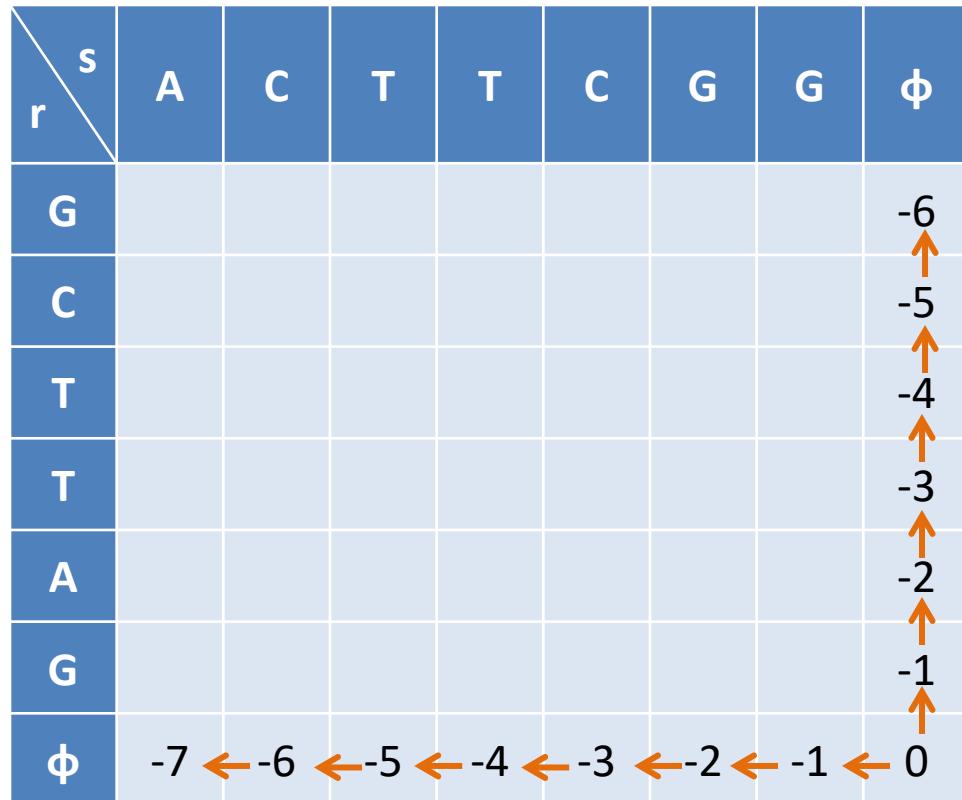
$$\max \begin{cases} V(i + 1, j + 1) + \sigma(r[i], s[j]) \\ V(i, j + 1) + \sigma(' ', s[j]) \\ V(i + 1, j) + \sigma(r[i], ' ') \end{cases}$$

- Optimal alignment score: $V(1, 1)$

Global alignment: An example

- r : GCTTAG
- s : ACTTCGG

Match: +1; otherwise: -1



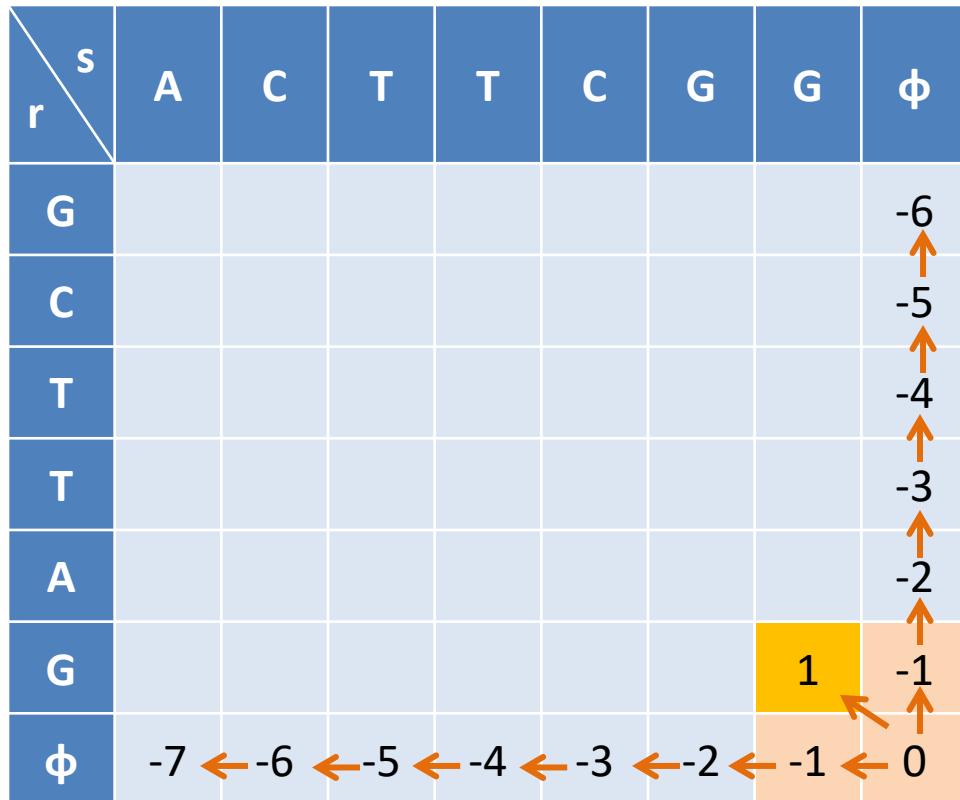
Recall: $V(i, j)$ equals the optimal (i.e., highest) alignment score of the suffixes $r[i..m]$ and $s[j..n]$

Initialization:
Depend on the gap penalty

Global alignment: An example

- r : GCTTAG
- s : ACTTCGG

$$V(i, j) = \max \begin{cases} V(i + 1, j + 1) + \sigma(r[i], s[j]) \\ V(i, j + 1) + \sigma(' ', s[j]) \\ V(i + 1, j) + \sigma(r[i], ' ') \end{cases}$$

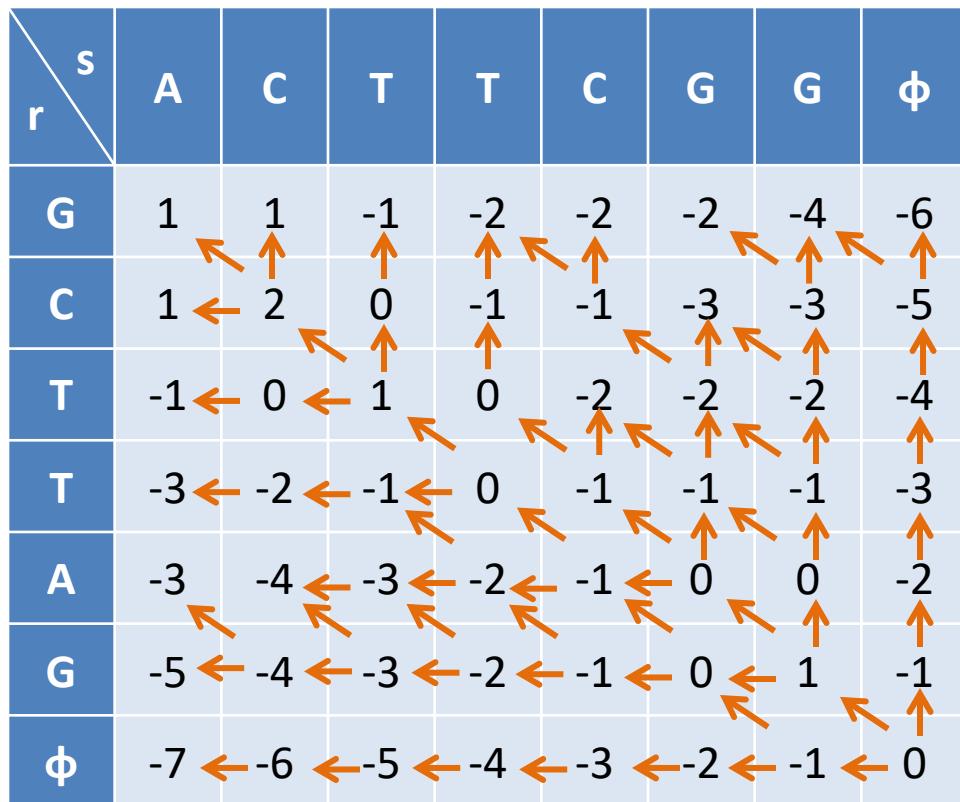


Align $r[6..6]$ with $s[7..7]$:
 $\max(0+1, -1-1, -1-1) =$
 $\max(1, -2, -2) = 1$

Global alignment: An example

- r : GCTTAG
- s : ACTTCGGG

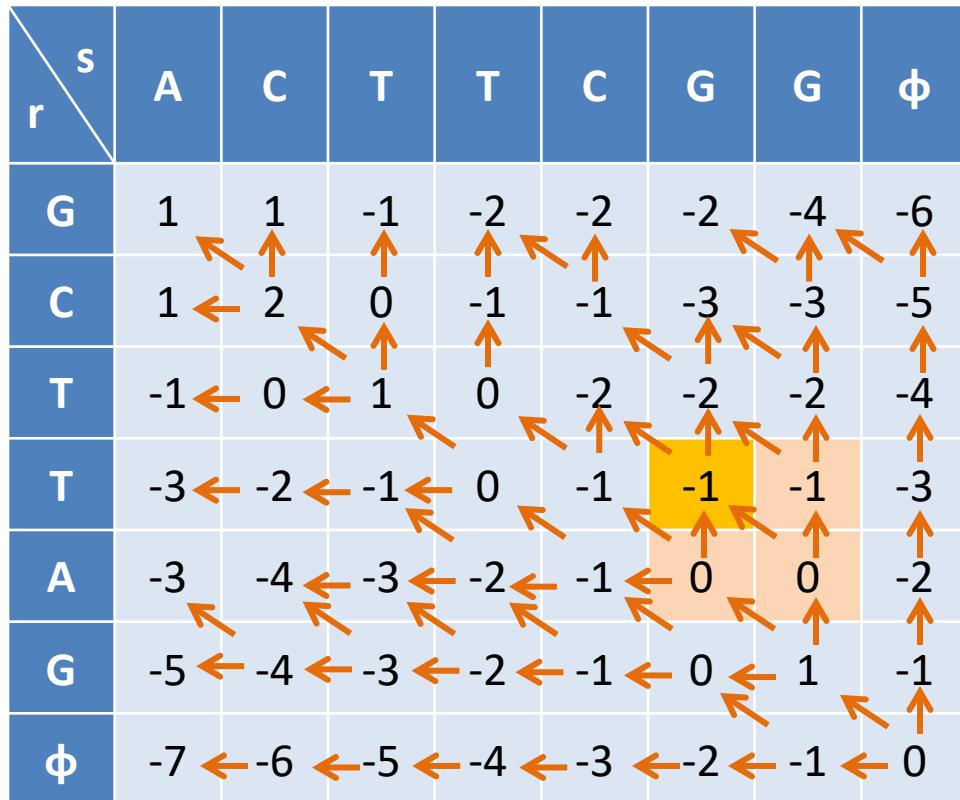
$$V(i, j) = \max \begin{cases} V(i + 1, j + 1) + \sigma(r[i], s[j]) \\ V(i, j + 1) + \sigma(' ', s[j]) \\ V(i + 1, j) + \sigma(r[i], ' ') \end{cases}$$



Global alignment: An example

- r : GCTTAG
- s : ACTTCGGG

$$V(i, j) = \max \begin{cases} V(i + 1, j + 1) + \sigma(r[i], s[j]) \\ V(i, j + 1) + \sigma(' ', s[j]) \\ V(i + 1, j) + \sigma(r[i], ' ') \end{cases}$$

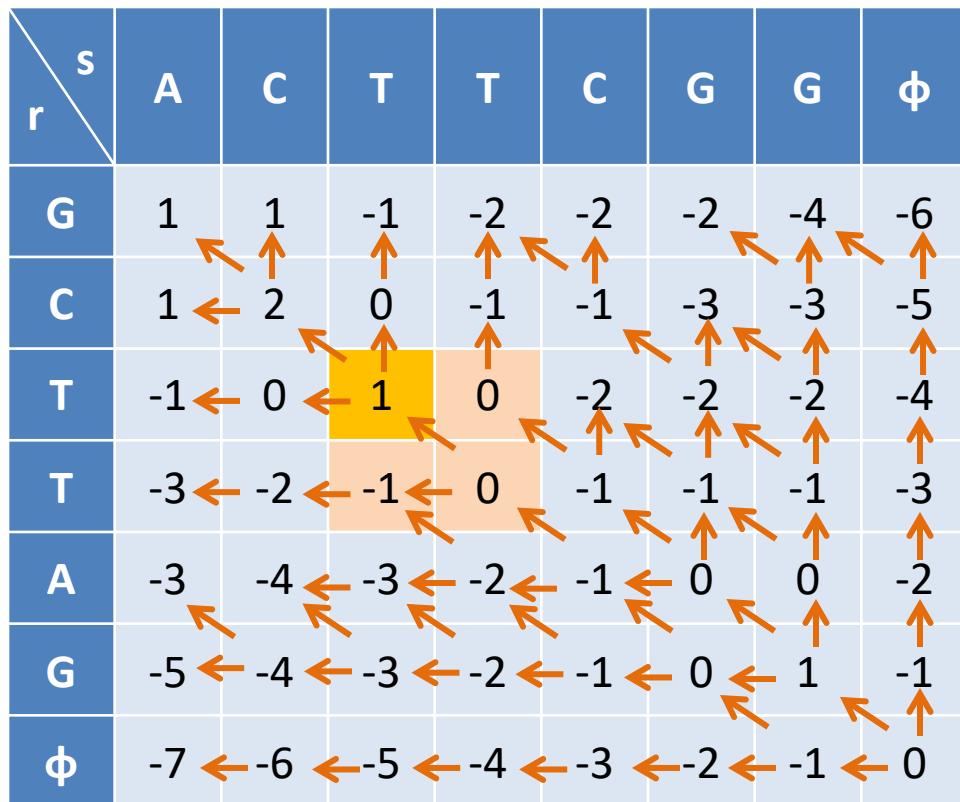


Align $r[4..6]$ with $s[6..7]$:
 $\max(0-1, -1-1, 0-1) =$
 $\max(-1, -2, -1) = -1$

Global alignment: An example

- r : GCTTAG
- s : ACTTCGGG

$$V(i, j) = \max \begin{cases} V(i + 1, j + 1) + \sigma(r[i], s[j]) \\ V(i, j + 1) + \sigma(' ', s[j]) \\ V(i + 1, j) + \sigma(r[i], ' ') \end{cases}$$



Align $r[3..6]$ with $s[3..7]$:
 $\max(0+1, 0-1, -1-1) =$
 $\max(1, -1, -2) = 1$

Local alignment

- Find the optimal alignment for the subsequences of two sequences, r and s .
- Smith-Waterman Algorithm
- Initialization: 0 (usually)
- DP Table Construction:

$$\max \left\{ \begin{array}{l} V(i + 1, j + 1) + \sigma(r[i], s[j]) \\ V(i, j + 1) + \sigma(' ', s[j]) \\ V(i + 1, j) + \sigma(r[i], ' ') \\ 0 \end{array} \right.$$

- Optimal alignment score: $\max_{i,j} V(i, j)$

Local alignment: An example

- r : GCTTAG
- s : ACTTCGG

Match: +1; otherwise: -1

$r \backslash s$	A	C	T	T	C	G	G	ϕ
G								0
C								0
T								0
T								0
A								0
G								0
ϕ	0	0	0	0	0	0	0	0

Recall: $V(i, j)$ equals the optimal (i.e., highest) alignment score of all the prefixes of suffixes $r[i..m]$ and $s[j..n]$

Initialization:
zeros

Local alignment: An example

- r : GCTTAG
- s : ACTTCGG

$$\max \begin{cases} V(i + 1, j + 1) + \sigma(r[i], s[j]) \\ V(i, j + 1) + \sigma(' ', s[j]) \\ V(i + 1, j) + \sigma(r[i], ' ') \\ 0 \end{cases}$$

r \ s	A	C	T	T	C	G	G	ϕ
G								0
C								0
T								0
T								0
A								0
G						1	0	
ϕ	0	0	0	0	0	0	0	

Align $r[6..6]$ with $s[7..7]$:
 $\max(0+1, 0-1, 0-1, 0) =$
 $\max(1, -1, -1, 0) = 1$

Local alignment: An example

- r : GCTTAG
- s : ACTTCGGG

$$\max \begin{cases} V(i + 1, j + 1) + \sigma(r[i], s[j]) \\ V(i, j + 1) + \sigma(' ', s[j]) \\ V(i + 1, j) + \sigma(r[i], ' ') \\ 0 \end{cases}$$

r \ s	A	C	T	T	C	G	G	ϕ
G	2	2	0	0	0	1	1	0
C	2	3	1	0	1	0	0	0
T	0	1	2	1	0	0	0	0
T	0	0	1	1	0	0	0	0
A	1	0	0	0	0	0	0	0
G	0	0	0	0	0	1	1	0
ϕ	0	0	0	0	0	0	0	0

Local alignment: An example

- r : GCTTAG
- s : ACTTCGGG

$$\max \begin{cases} V(i + 1, j + 1) + \sigma(r[i], s[j]) \\ V(i, j + 1) + \sigma(' ', s[j]) \\ V(i + 1, j) + \sigma(r[i], ' ') \\ 0 \end{cases}$$

r \ s	A	C	T	T	C	G	G	ϕ
G	2	2	0	0	0	1	1	0
C	2	3	1	0	1	0	0	0
T	0	1	2	1	0	0	0	0
T	0	0	1	1	0	0	0	0
A	1	0	0	0	0	0	0	0
G	0	0	0	0	0	1	1	0
ϕ	0	0	0	0	0	0	0	0

Align $r[4..6]$ with $s[6..7]$:
 $\max(0-1, 0-1, 0-1, 0) =$
 $\max(-1, -1, -1, 0) = 0$

Local alignment: An example

- r : GCTTAG
- s : ACTTCGGG

$$\max \begin{cases} V(i + 1, j + 1) + \sigma(r[i], s[j]) \\ V(i, j + 1) + \sigma(' ', s[j]) \\ V(i + 1, j) + \sigma(r[i], ' ') \\ 0 \end{cases}$$

r \ s	A	C	T	T	C	G	G	ϕ
G	2	2	0	0	0	1	1	0
C	2	3	1	0	1	0	0	0
T	0	1	2	1	0	0	0	0
T	0	0	1	1	0	0	0	0
A	1	0	0	0	0	0	0	0
G	0	0	0	0	0	1	1	0
ϕ	0	0	0	0	0	0	0	0

Align $r[3..6]$ with $s[3..7]$:
 $\max(1+1, 1-1, 1-1, 0) =$
 $\max(2, 0, 0, 0) = 2$

Comparison between global and local alignment

	Global Alignment	Local Alignment
Description	Alignment of 2 or more sequences	Alignment of subsequences of 2 or more sequences
Application	Whole genome sequence alignment	Find motifs, domains
Initialization	According to the gap penalty	Zeros (usually)
DP table, $V(i, j) =$	$\max \begin{cases} V(i + 1, j + 1) + \sigma(r[i], s[j]) \\ V(i, j + 1) + \sigma(' ', s[j]) \\ V(i + 1, j) + \sigma(r[i], ' ') \end{cases}$	$\max \begin{cases} V(i + 1, j + 1) + \sigma(r[i], s[j]) \\ V(i, j + 1) + \sigma(' ', s[j]) \\ V(i + 1, j) + \sigma(r[i], ' ') \\ \mathbf{0} \end{cases}$
Optimal alignment score	Top left corner, i.e., $V(1, 1)$	The largest value

Demonstration: Web tools for alignment

- Find a Web tool to perform Needleman-Wunsch optimal global alignment on the sequences below:
 - GCTTAG
 - ACTTCGG

Step 1: Find tools on Google

needle-wunsch web tool

All Images Shopping News Videos More Settings Tools

About 1,310,000 results (0.42 seconds)

Pairwise Sequence Alignment Tools < EMBL-EBI

Stretcher uses a modification of the **Needleman-Wunsch** algorithm ...

[EMBOSS Needle](#) · [EMBOSS Water](#) · [EMBOSS Stretcher](#) · [EMBOSS Matcher](#)

[www.ebi.ac.uk](#) › training › online › glossary › needleman-wunsch-alg... ▾

Needleman-Wunsch algorithm | EMBL-EBI Train online

Needleman-Wunsch algorithm. This global sequence alignment method – the first to apply dynamic programming techniques to biological sequence analysis ...

[rna.informatik.uni-freiburg.de](#) › Teaching › toolName=Needleman-Wun...

Teaching - Needleman-Wunsch - Freiburg RNA Tools

Teaching - Needleman-Wunsch : global, linear gap cost source at ... Saul B. Needleman and Christian D. Wunsch introduced 1970 an approach to ... Freiburg RNA tools: a central online resource for RNA-focused research and teaching

[blast.ncbi.nlm.nih.gov](#) › Blast ▾

Needleman-Wunsch alignment of two nucleotide sequences

Needleman-Wunsch Global Align Nucleotide Sequences. Nucleotide · Protein. Needleman-

<https://www.google.com/preferences> alignment of two nucleotide sequences [?]. Reset page ...

Step 2: Select the tool

The screenshot shows a web browser window for the "Pairwise Sequence Alignment" tool at ebi.ac.uk/Tools/psa/. The page has a teal header with the title "Pairwise Sequence Alignment". Below the header, there's a navigation bar with links for "EMBL-EBI", "Services", "Research", "Training", "Industry", "About us", and a search icon. On the right side of the header is the "EMBL-EBI Hinxton" logo. The main content area has a large heading "Pairwise Sequence Alignment". Below it, a sub-header "Tools > Pairwise Sequence Alignment" is followed by a paragraph explaining what pairwise sequence alignment is. Another paragraph contrasts it with multiple sequence alignment (MSA). A section titled "Global Alignment" is shown, with a sub-section for "Needle (EMBOSS)". It describes the EMBOSS Needle algorithm and provides a "Launch Needle" button, which is circled in red. A second section for "Stretcher (EMBOSS)" is also present.

Pairwise Sequence Alignment is used to identify regions of similarity that may indicate functional, structural and/or evolutionary relationships between two biological sequences (protein or nucleic acid).

By contrast, **Multiple Sequence Alignment (MSA)** is the alignment of three or more biological sequences of similar length. From the output of MSA applications, homology can be inferred and the evolutionary relationship between the sequences studied.

Global Alignment

Global alignment tools create an end-to-end alignment of the sequences to be aligned.

Needle (EMBOSS)

EMBOSS Needle creates an optimal global alignment of two sequences using the Needleman-Wunsch algorithm.

[Launch Needle](#)

Stretcher (EMBOSS)

EMBOSS Stretcher uses a modification of the Needleman-Wunsch algorithm that allows larger sequences to be globally aligned.

Step 3: Inputs

The screenshot shows a web browser window for the EMBOSS Needle tool at ebi.ac.uk/Tools/psa/emboss_needle/. The page title is "Pairwise Sequence Alignment". A red circle highlights the "DNA" selection in the dropdown menu under "Enter a pair of". Another red circle highlights the sequence "GCTTAG" entered in the first sequence input field. A third red circle highlights the sequence "ACTTCGG" entered in the second sequence input field.

EMBOSS Needle reads two input sequences and writes their optimal global sequence alignment to file.

STEP 1 - Enter your nucleotide sequences

Enter a pair of

DNA

sequences. Enter or paste your first **nucleotide** sequence in any supported format:

GCTTAG

Or, upload a file: Choose File No file chosen

Use a example sequence | Clear sequence | See more example inputs

AND

Enter or paste your second **nucleotide** sequence in any supported format:

ACTTCGG

Step 4: Parameters

- End gap penalty: true
- Others: default settings
- Note: match: +5; mismatch: -4

STEP 2 - Set your pairwise alignment options

OUTPUT FORMAT

pair

MATRIX	GAP OPEN	GAP EXTEND	END GAP PENALTY	END GAP OPEN	END GAP EXTEND
DNAfull	10	0.5	true	10	0.5

STEP 3 - Submit your job

Be notified by email (*Tick this box if you want to be notified by email when the results are available*)

Submit

Step 5: Show the results

The screenshot shows a web browser window for the EMBoss Needle - Alignment tool. The URL is ebi.ac.uk/Tools/services/web/toolresult.ebi?jobId=emboss_needle-I20200118-164618-0248-44475... The page has a green header bar with links for Input form, Web services, Help & Documentation, Bioinformatics Tools FAQ, Feedback, and Share.

```
#####
# Program: needle
# Rundate: Sat 18 Jan 2020 16:46:20
# Commandline: needle
#   -auto
#   -stdout
#   -asequence emboss_needle-I20200118-164618-0248-444753-p2m.asequence
#   -bsequence emboss_needle-I20200118-164618-0248-444753-p2m.bsequence
#   -datafile EDNAFULL
#   -gapopen 10.0
#   -gapextend 0.5
#   -endweight
#   -endopen 10.0
#   -endextend 0.5
#   -aformat3 pair
#   -snucleotide1
#   -snucleotide2
# Align_format: pair
# Report_file: stdout
#####
=====
#
# Aligned sequences: 2
# 1: EMBOS_001
# 2: EMBOS_001
# Matrix: EDNAFULL
# Gap_penalty: 10.0
# Extend_penalty: 0.5
#
# Length: 7
# Identity:      4/7 (57.1%)
# Similarity:    4/7 (57.1%)
# Gaps:          1/7 (14.3%)
# Score:         2.0
#
# =====
EMBOSS_001      1 GCTT-AG      6
                  .||| .|
EMBOSS_001      1 ACTTCGG     7
```

Exercise 1: Pairwise sequence alignment

- Find the optimal **global** alignment(s) and **local** alignment(s) for the following sequences and find the alignment score for each optimal alignment.
 - r : GCTTAG
 - s : ACTTCGG
- Match: +1; otherwise: -1

Answer: optimal global alignment

$r \backslash s$	A	C	T	T	C	G	G	'''
G	1	1	-1	-2	-2	-2	-4	-6
C	1	2	0	-1	-1	-3	-3	-5
T	-1	0	1	0	-2	-2	-2	-4
T	-3	-2	-1	0	-1	-1	-1	-3
A	-3	-4	-3	-2	-1	0	0	-2
G	-5	-4	-3	-2	-1	0	1	-1
''	-7	-6	-5	-4	-3	-2	-1	0

Answer: optimal global alignment

$r \backslash s$	A	C	T	T	C	G	G	'''
G	1	1	-1	-2	-2	-2	-4	-6
C	1	2	0	-1	-1	-3	-3	-5
T	-1	0	1	0	-2	-2	-2	-4
T	-3	-2	-1	0	-1	-1	-1	-3
A	-3	-4	-3	-2	-1	0	0	-2
G	-5	-4	-3	-2	-1	0	1	-1
''	-7	-6	-5	-4	-3	-2	-1	0

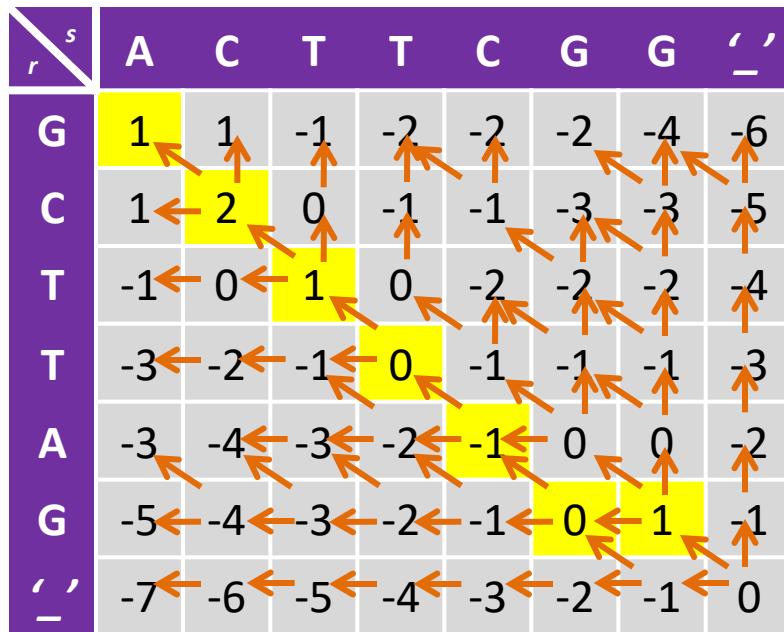
Best alignment score: 1

Best alignment 1:

r : GCTT_AG

s : ACTTCGG

Answer: optimal global alignment



Best alignment score: 1

Best alignment 2:

r : GCTTA_G

s : ACTTCGG

Answer: optimal global alignment

$r \backslash s$	A	C	T	T	C	G	G	'''
G	1	1	-1	-2	-2	-2	-4	-6
C	1	2	0	-1	-1	-3	-3	-5
T	-1	0	1	0	-2	-2	-2	-4
T	-3	-2	-1	0	-1	-1	-1	-3
A	-3	-4	-3	-2	-1	0	0	-2
G	-5	-4	-3	-2	-1	0	1	-1
''	-7	-6	-5	-4	-3	-2	-1	0

Best alignment score: 1

Best alignment 3:

r : GCTTAG_

s : ACTTCGG

Answer: optimal local alignment

$r \backslash s$	A	C	T	T	C	G	G	'''
G	2	2	0	0	0	1	1	0
C	2	3	1	0	1	0	0	0
T	0	1	2	1	0	0	0	0
T	0	0	1	1	0	0	0	0
A	1	0	0	0	0	0	0	0
G	0	0	0	0	0	1	1	0
'''	0	0	0	0	0	0	0	0

Best alignment score: 3

Best alignment 1:

r : CTT

s : CTT

Best alignment 2:

r : CTTAG

s : CTTCG

Scoring matrix

- Give higher score for more similar characters
- Handle gaps more properly
- Always be **symmetric**

	A	C	G	T	φ
A	1	-1	-1	-1	-1
C	-1	1	-1	-1	-1
G	-1	-1	1	-1	-1
T	-1	-1	-1	1	-1
φ	-1	-1	-1	-1	NA

	A	C	G	T	φ
A	5	-3	-2	-3	-7
C	-3	6	-3	-2	-5
G	-2	-3	6	-3	-7
T	-3	-2	-3	5	-5
φ	-7	-5	-7	-5	NA

More about scoring matrix

- Different matrices for DNA sequence alignment and protein sequence alignment
- For similarity of DNA, consider *transition* and *transversion*
- For similarity of amino acid, consider *charge*, *hydrophobicity* and *size* of amino acid
- Affine gap penalty (separate costs for gap opening and gap size)

$$y = -a - bx$$

$-a$: Gap opening penalty
 $-b$: Gap size penalty

Exercise 2: Scoring matrix and alignment

- A) Suppose the mismatch penalties for transition and transversion are -1 and -2 respectively, the indel penalty is -3, and the score of a match is 2. Write down the scoring matrix.

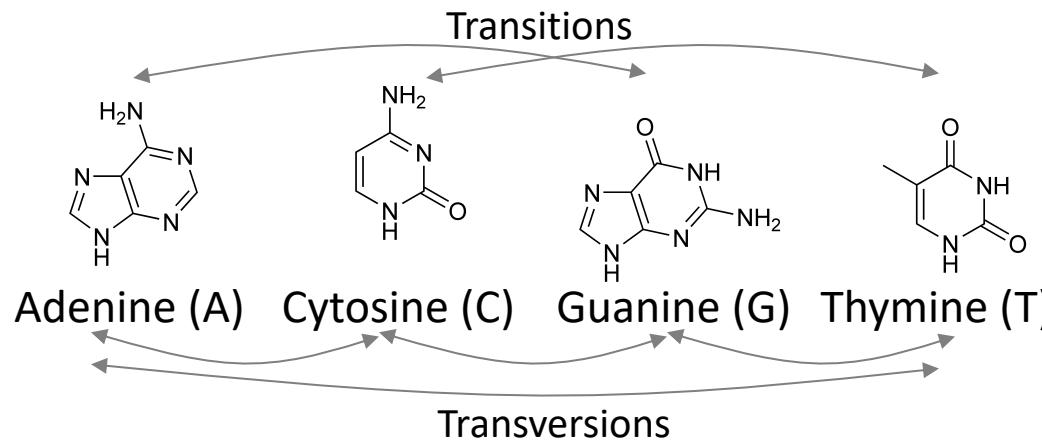


Image credit: Wikipedia

Answer: Scoring matrix

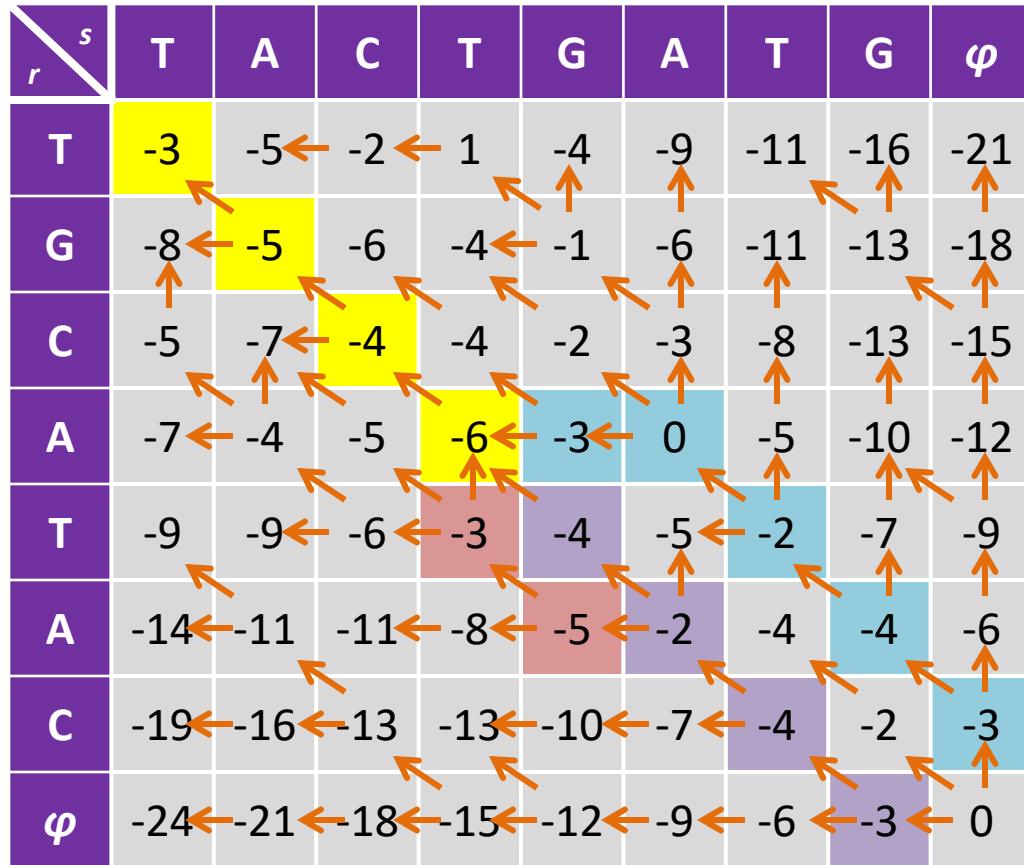
	A	C	G	T	φ
A	2	-2	-1	-2	-3
C	-2	2	-2	-1	-3
G	-1	-2	2	-2	-3
T	-2	-1	-2	2	-3
φ	-3	-3	-3	-3	NA

Exercise 2: Scoring matrix and alignment

- B) Using the previous scoring matrix, find the best global alignment of the following sequences:
 - TGCATAC
 - TACTGATG

	A	C	G	T	φ
A	2	-2	-1	-2	-3
C	-2	2	-2	-1	-3
G	-1	-2	2	-2	-3
T	-2	-1	-2	2	-3
φ	-3	-3	-3	-3	NA

Exercise 2: Scoring matrix and alignment



	A	C	G	T	φ
A	2	-2	-1	-2	-3
C	-2	2	-2	-1	-3
G	-1	-2	2	-2	-3
T	-2	-1	-2	2	-3
φ	-3	-3	-3	-3	NA

Best alignment score: 3

r : TGCATAC
 s : TACTGATG

r : TGC _ ATAC
 s : TACTGATG _

r : TGCAT _ AC
 s : TAC _ TGATG

Exercise 3: Affine gap penalty

- You are provided with a scoring matrix. Given that gap opening penalty = -5 and gap size penalty = -1, calculate the alignment score for the following alignments:

AACCTT
AAG _ TT

$$\begin{aligned} & 4*(5)+1*(-3)+1*(-5)+1*(-1) \\ & =11 \end{aligned}$$

ACGGCTTCGGCA
AC __ CTTC _ ACA

$$\begin{aligned} & 8*(5)+1*(-2)+2*(-5)+3*(-1) \\ & =25 \end{aligned}$$

CTCTCGGGGGGAACC
CCCTC _____ AACC

$$\begin{aligned} & 8*(5)+1*(-2)+1*(-5)+6*(-1) \\ & =27 \end{aligned}$$

	A	C	G	T
A	5	-3	-2	-3
C	-3	5	-3	-2
G	-2	-3	5	-3
T	-3	-2	-3	5

Check list

- What are two main ideas behind dynamic programming?
- What are the considerations in proposing the scoring rules in sequence alignment problem?
- What is affine gap penalty in solving sequence alignment problem?
- Are there any negative numbers in the D.P. table of the local alignment? Why?
- (Optional) Deduce the number of possible alignments of a pair of sequences with length m and n .
- (Optional) How can you perform optimal alignment for three sequences?