

Sequence Comparison: Local Alignment

Genome 373
Genomic Informatics
Elhanan Borenstein

A quick review: Global Alignment

Global Alignment Mission:

Find the best global alignment between two sequences.

An algorithm for finding the alignment with the best score

A method for scoring alignments

A quick review: Global Alignment

Global Alignment Mission:

Find the best global alignment between two sequences.

An algorithm for finding the alignment with the best score

A method for scoring alignments



- Substitution matrix:

	A	C	G	T
A	10	-5	0	-5
C	-5	10	-5	0
G	0	-5	10	-5
T	-5	0	-5	10

- Gap penalty:

- Linear gap penalty
- Affine gap penalty

Purine	A	G
Pyrimidine	C	T

GAAT-C **d=-4**
CA-TAC
-5 + 10 + -4 + 10 + -4 + 10 = 17

Review: Global Alignment

- Three Possible Moves:
 - A diagonal move aligns a character from each sequence.
 - A horizontal move aligns a gap in the seq along the left edge
 - A vertical move aligns a gap in the seq along the top edge.
- The move you keep is the best scoring of the three.

		G	A	A	T	C
	0	→ -4	→ -8	→ -12	→ -16	→ -20
C	↓ -4	↘ -5				
A	↓ -8	↘ -4	?			
T	↓ -12					
A	↓ -16					
C	↓ -20					

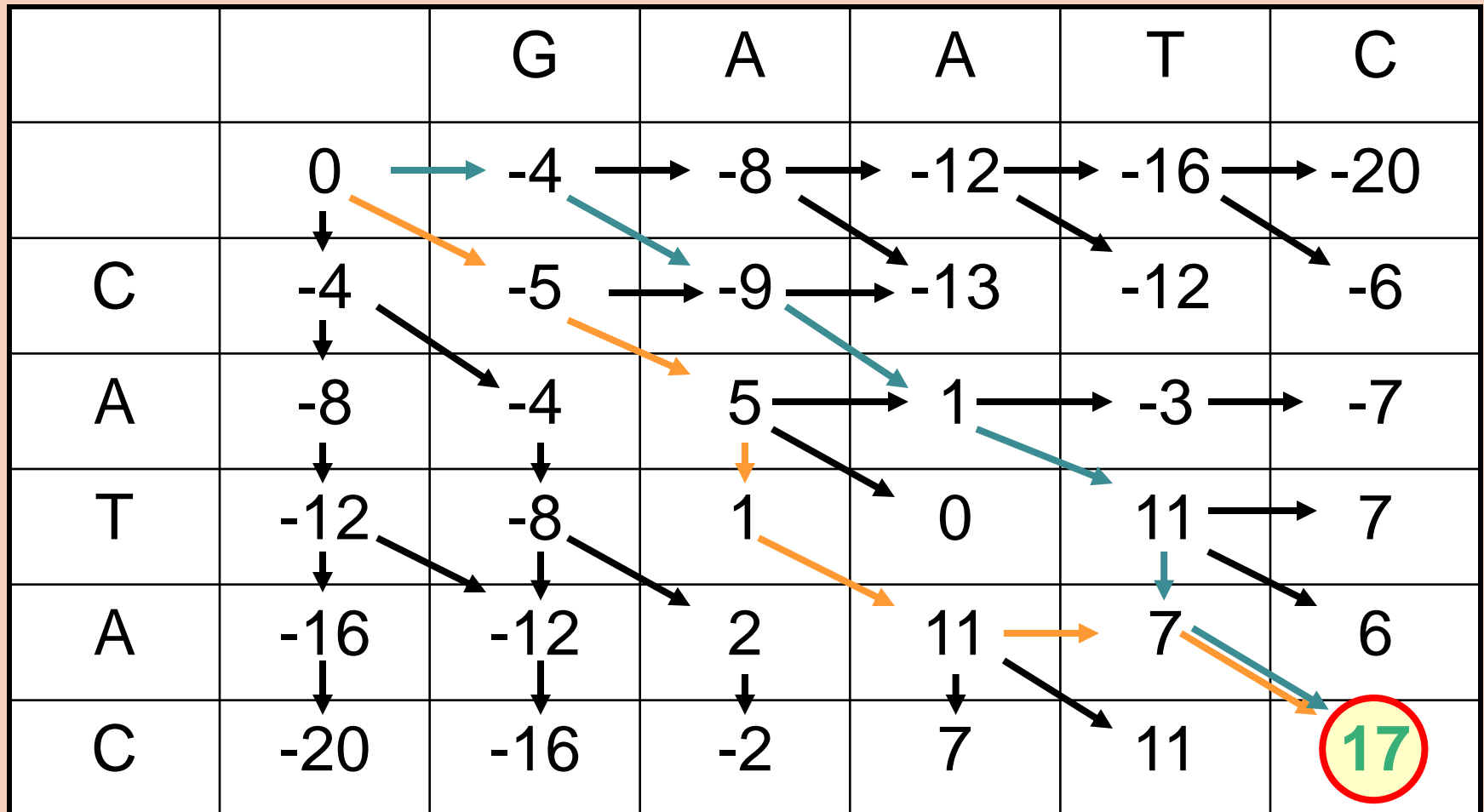
Review: Global Alignment

Fill DP matrix from upper left to lower right.
Traceback alignment from lower right corner.

Substitution matrix				
	A	C	G	T
A	10	-5	0	-5
C	-5	10	-5	0
G	0	-5	10	-5
T	-5	0	-5	10

Gap penalty: $d = -4$

		G	A	A	T	C
	0	-4	-8	-12	-16	-20
C	-4	-5	-9	-13	-12	-6
A	-8	-4	5	1	-3	-7
T	-12	-8	1	0	11	7
A	-16	-12	2	11	7	6
C	-20	-16	-2	7	11	17



DP in equation form

		G	A	A	T	C
	0	-4	-8	-12	-16	-20
C	-4	-5				
A	-8	-4	?			
T	-12					
A	-16					
C	-20					

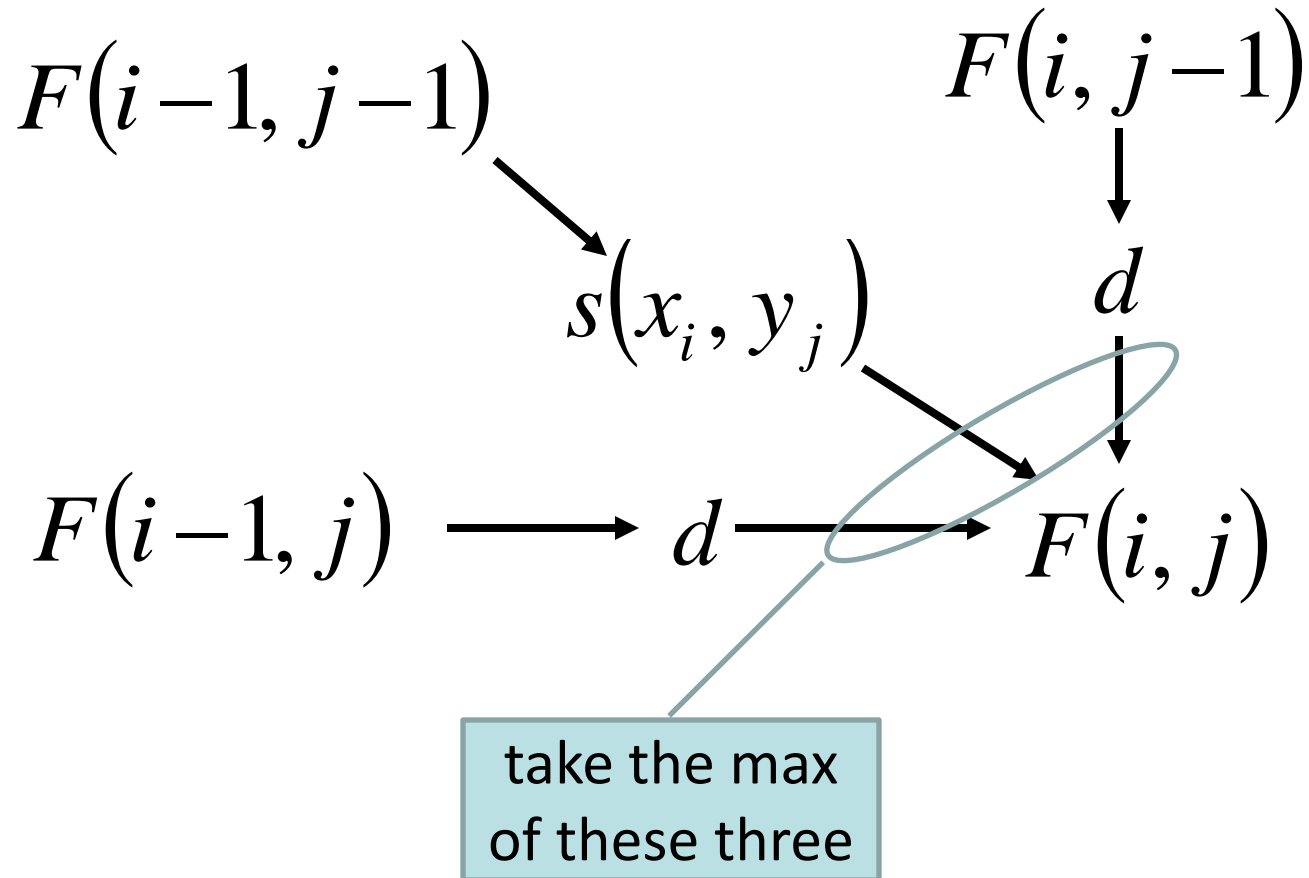
- Align sequence **x** and **y**.
- F** is the DP matrix; **s** is the substitution matrix;
d is the linear gap penalty.

$$F(0,0) = 0$$

$$F(i, j) = \max \begin{cases} F(i-1, j-1) + s(x_i, y_j) \\ F(i-1, j) + d \\ F(i, j-1) + d \end{cases}$$

DP equation graphically

		G	A	A	T	C
	0	→ -4	→ -8	→ -12	→ -16	→ -20
C	-4	↘ -5				
A	-8	↘ -4	?			
T	-12					
A	-16					
C	-20					



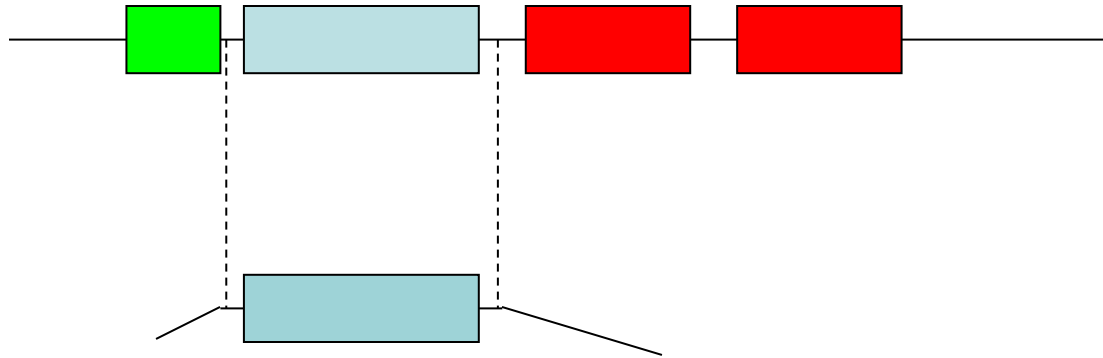
Local alignment

Mission:

Find best partial alignment
between two sequences.

Why?

Local alignment



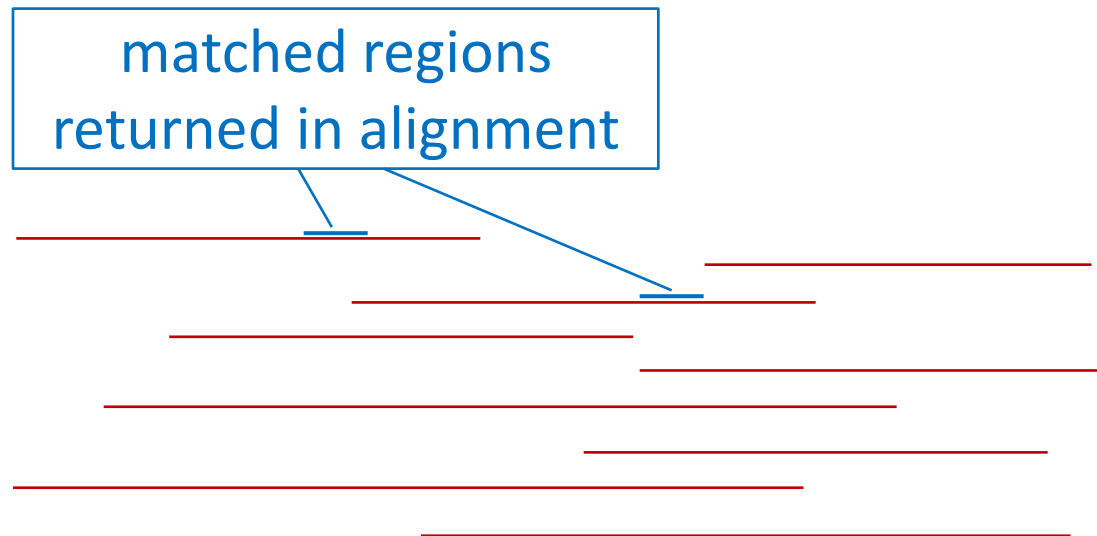
- A single-**domain** protein may be similar only to one region within a multi-domain protein.
- A DNA/RNA **query** may align to a small part of a genome/genomes/metagenomes.
- An alignment that spans the complete length of both sequences may be undesirable.

BLAST does local alignments

- Typical search has a short query against long targets.
- The alignments returned show only the well-aligned match region of both query and target.

Query: —————

Targets: (e.g. genome
contigs, full genomes,
metagenomes)



How can we modify the *Needleman-Wunsch* DP algorithm (for finding global alignment) such that it will find instead the best **local** alignment??

		G	A	A	T	C
	0	-4	-8	-12	-16	-20
C	-4	-5	-9	-13	-12	-6
A	-8	-4	5	1	-3	-7
T	-12	-8	1	0	11	7
A	-16	-12	2	11	7	6
C	-20	-16	-2	7	11	17

Substitution matrix

	A	C	G	T
A	10	-5	0	-5
C	-5	10	-5	0
G	0	-5	10	-5
T	-5	0	-5	10

Gap penalty: $d = -4$

↓

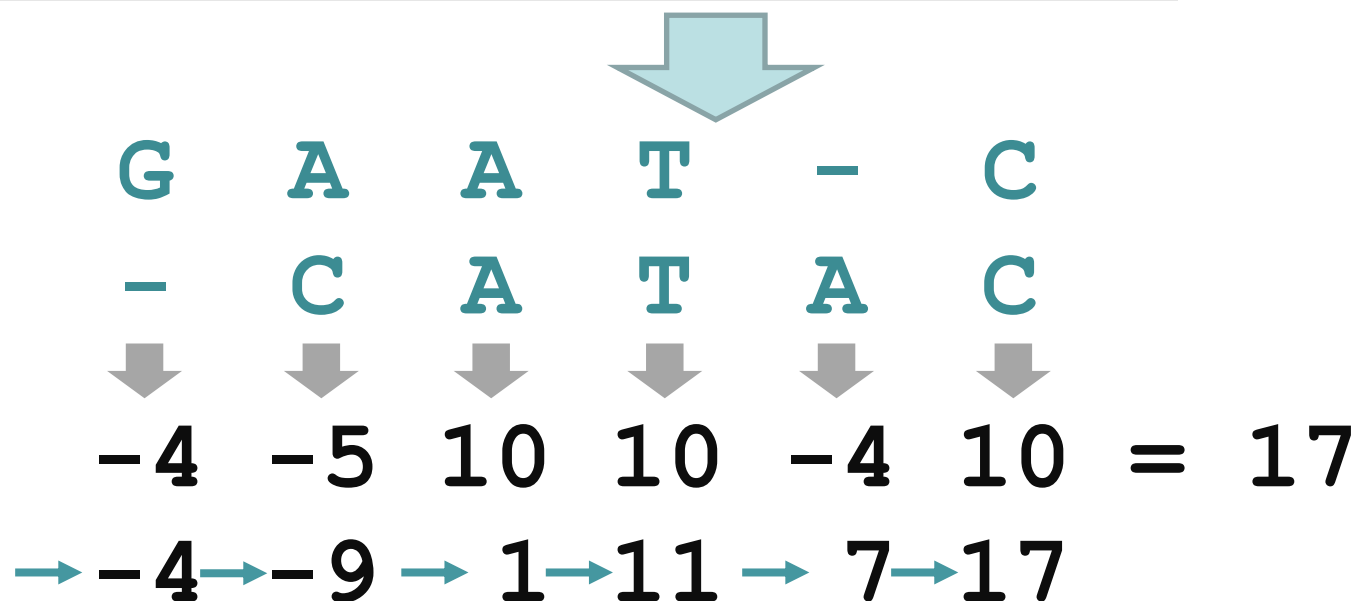
G	A	A	T	-	C	
-	C	A	T	A	C	
↓	↓	↓	↓	↓	↓	
-4	-5	10	10	-4	10	= 17

		G	A	A	T	C
	0	-4	-8	-12	-16	-20
C	-4	-5	-9	-13	-12	-6
A	-8	-4	5	1	-3	-7
T	-12	-8	1	0	11	7
A	-16	-12	2	11	7	6
C	-20	-16	-2	7	11	17

Substitution matrix

	A	C	G	T
A	10	-5	0	-5
C	-5	10	-5	0
G	0	-5	10	-5
T	-5	0	-5	10

Gap penalty: $d = -4$



Remember: Global alignment DP

- Align sequence x and y .
- F is the DP matrix; s is the substitution matrix; d is the linear gap penalty.

$$F(0,0) = 0$$


$$F(i, j) = \max \begin{cases} F(i-1, j-1) + s(x_i, y_j) \\ F(i-1, j) + d \\ F(i, j-1) + d \end{cases}$$

Local alignment DP

- Align sequence x and y.
- F is the DP matrix; s is the substitution matrix; d is the linear gap penalty.

$$F(0,0) = 0$$

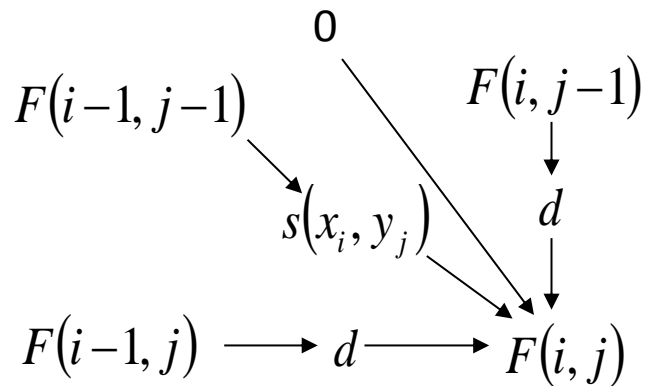
$$F(i, j) = \max \begin{cases} F(i-1, j-1) + s(x_i, y_j) \\ F(i-1, j) + d \\ F(i, j-1) + d \\ 0 \end{cases}$$

 (corresponds to start of alignment)

A simple example

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2

d = -5



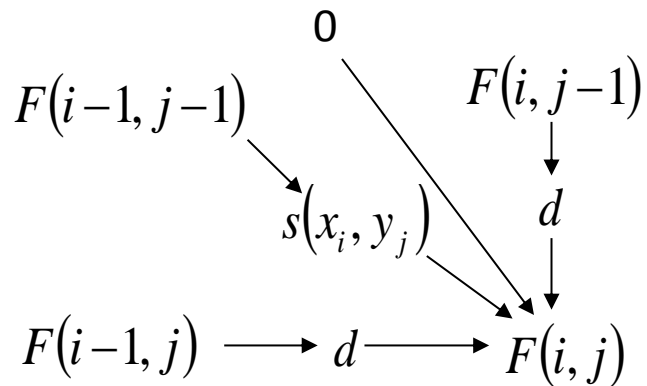
initialize the same way as
for global alignment

		A	A	G
	0			
A				
G				
C				

A simple example

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2

$d = -5$

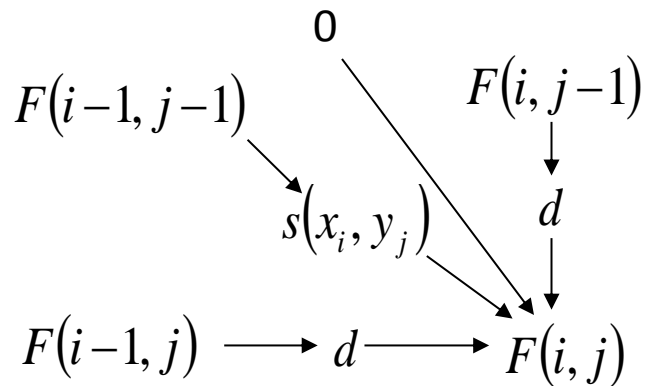


		A	A	G
	0	?	?	?
A	?			
G	?			
C	?			

A simple example

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2

d = -5

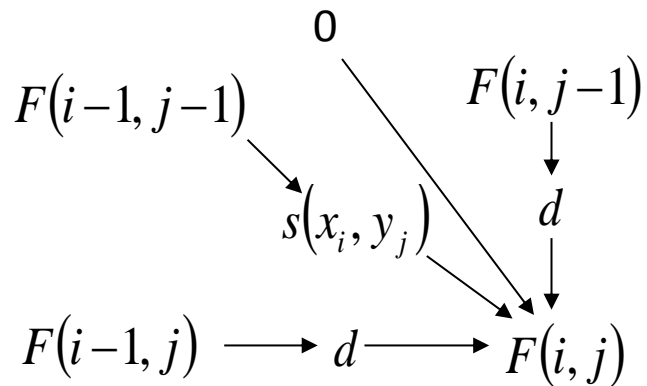


		A	A	G
		0	0	0
A		0		
G		0		
C		0		

A simple example

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2

$d = -5$

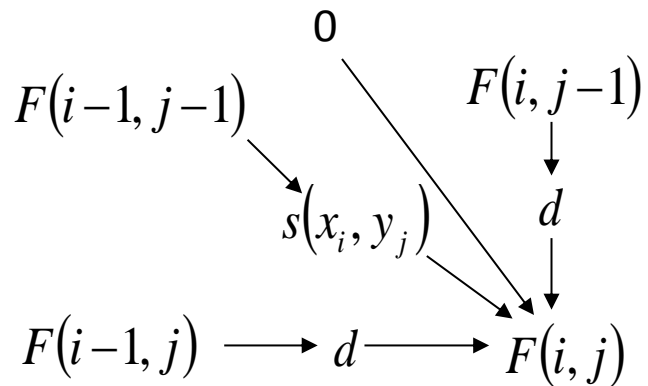


		A	A	G
		0	0	0
A		0	?	
G		0		
C		0		

A simple example

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2

d = -5



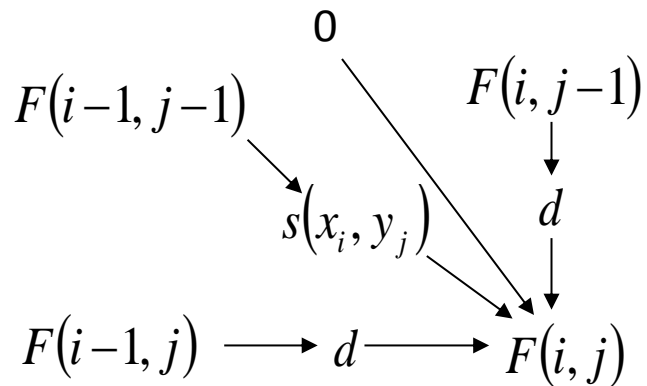
		A	A	G
		0	0	0
A	0	2	-5	
G	0	-5	0	
C	0			

A
A

A simple example

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2

d = -5

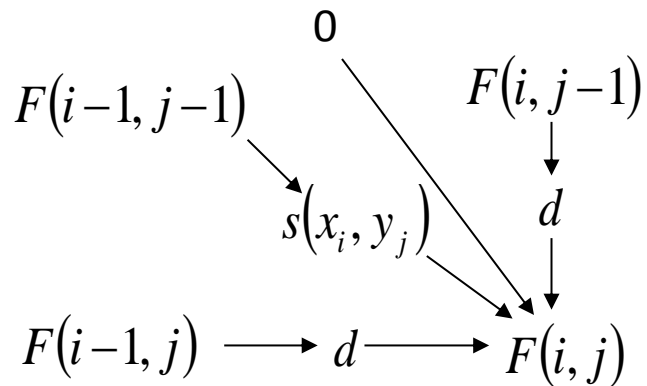


		A	A	G
		0	0	0
A	0	2		
G	0			
C	0			

A simple example

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2

d = -5

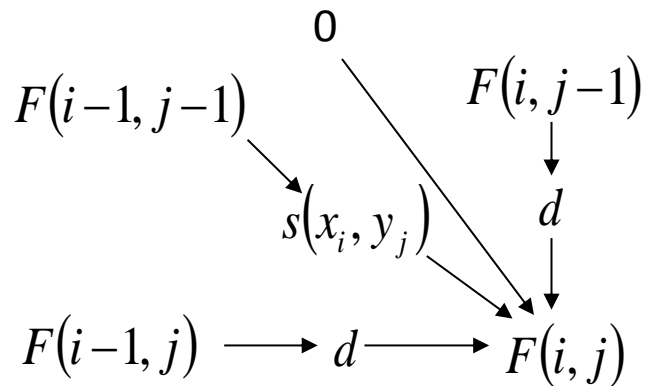


		A	A	G
		0	0	0
A	0	2		
G	0	?		
C	0	?		

A simple example

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2

d = -5

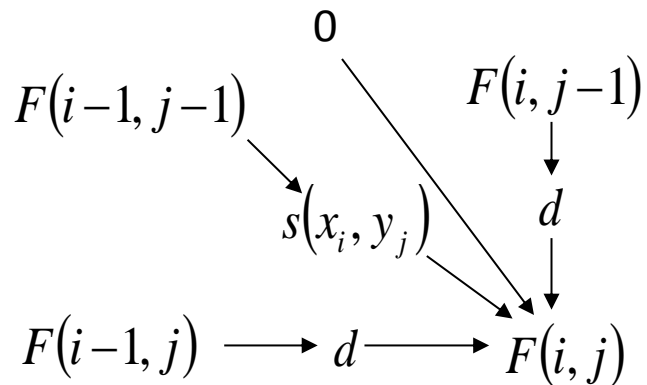


		A	A	G
		0	0	0
A	0	2		
G	0	-5	-3	
C	0	-5	0	

A simple example

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2

d = -5



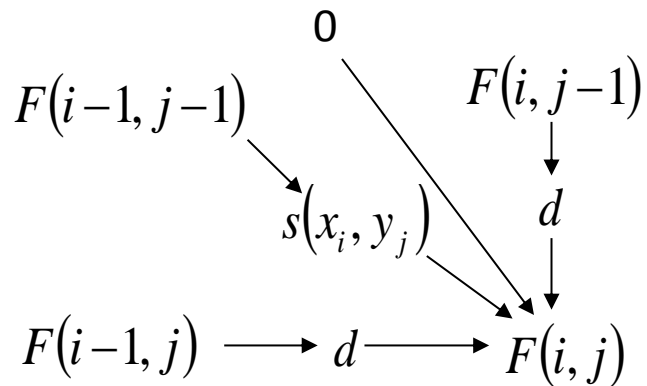
		A	A	G
		0	0	0
A		0	2	
G		0	0	
C		0	?	

(signify no preceding alignment with no arrow)

A simple example

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2

d = -5

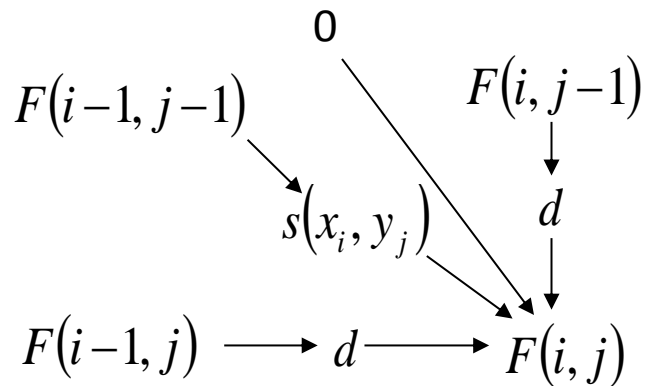


		A	A	G
		0	0	0
A		0	2	?
G		0	0	?
C		0	0	?

A simple example

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2

d = -5

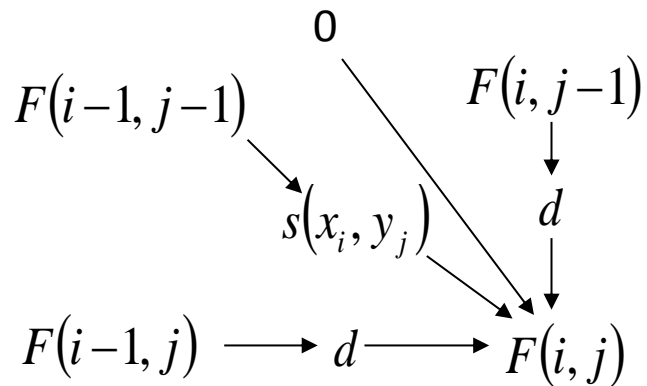


		A	A	G
		0	0	0
A		0	2	2
G		0	0	0
C		0	0	0

A simple example

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2

d = -5

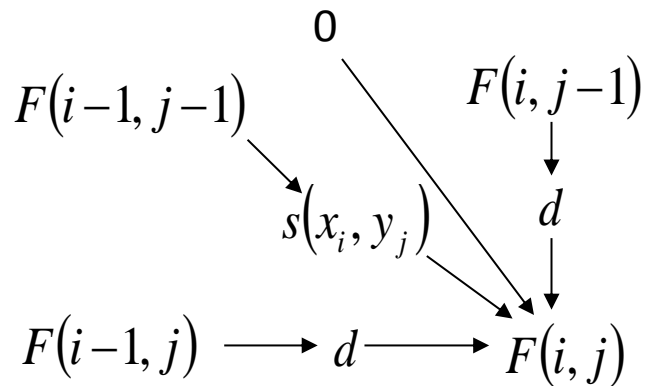


		A	A	G
		0	0	0
A		0	2	?
G		0	0	?
C		0	0	?

A simple example

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2

d = -5



		A	A	G
	0	0	0	0
A	0	2	2	0
G	0	0	0	4
C	0	0	0	0

What's different about the DP matrix

Global Alignment DP Matrix

		G	A	A	T	C
	0	-4	-8	-12	-16	-20
C	-4	-5	-9	-13	-12	-6
A	-8	-4	5	1	-3	-7
T	-12	-8	1	0	11	7
A	-16	-12	2	11	7	6
C	-20	-16	-2	7	11	17

Global Alignment DP Matrix showing sequence alignment between CAGATA and GATA. The matrix includes gap penalties (blue arrows) and match/mismatch scores (orange arrows). The final score 17 is circled in red.

Local Alignment DP Matrix

		A	A	G
	0	0	0	0
A	0	2	2	0
G	0	0	0	4
C	0	0	0	0

Local Alignment DP Matrix showing sequence alignment between AAG and AA. The matrix includes gap penalties (blue arrows) and match/mismatch scores (orange arrows). The final score 4 is circled in red.

A simple example

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2

d = -5

		A	A	G
	0	0	0	0
A	0	2	2	0
	0	0	0	4
	0	0	0	0

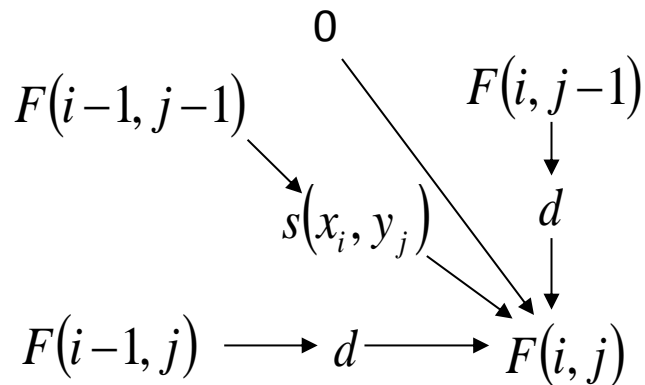
**But ...
how do we
traceback?**

Traceback

AG
AG

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2

d = -5

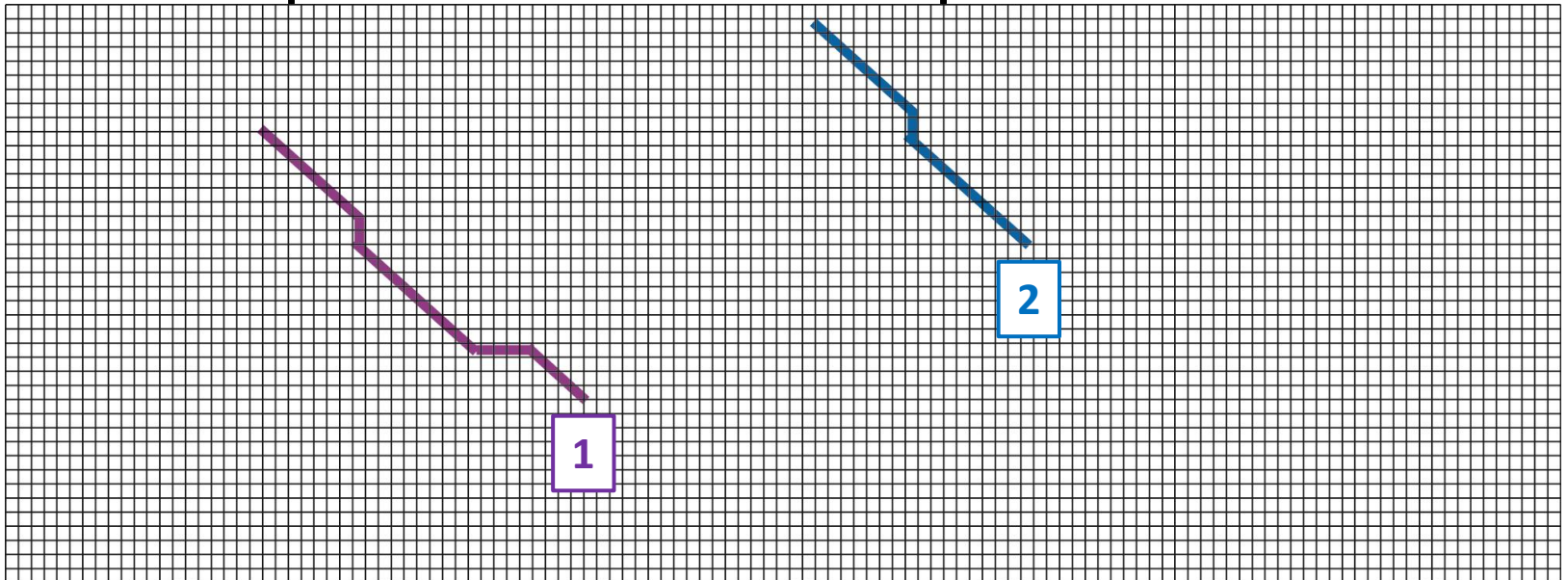


		A	A	G
		0	0	0
A	0	2	2	0
G	0	0	0	4
C	0	0	0	0

Start traceback at highest score anywhere in matrix, follow arrows back until you reach 0

Multiple local alignments

- Traceback from highest score, marking each DP matrix score along traceback.
- Now traceback from the remaining highest score, etc.
- The alignments may or may not include the same parts of the two sequences.



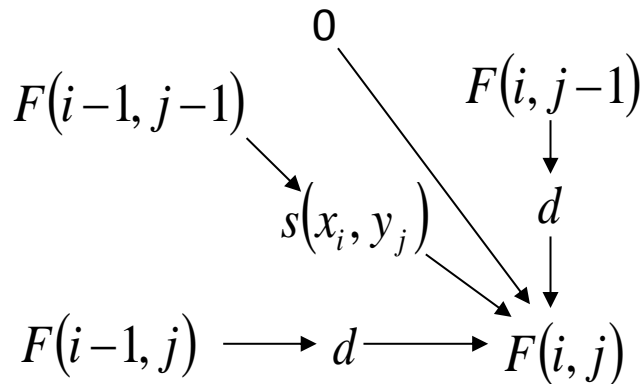
Local alignment

- Two differences from global alignment:
 - If a DP score is negative, replace with 0.
 - Traceback from the highest score in the matrix and continue until you reach 0.
- Global alignment algorithm: *Needleman-Wunsch*.
- Local alignment algorithm: *Smith-Waterman*.

Another example

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2

$d = -5$



Find the optimal **local** alignment of **AAG** and **GAAGGC**.

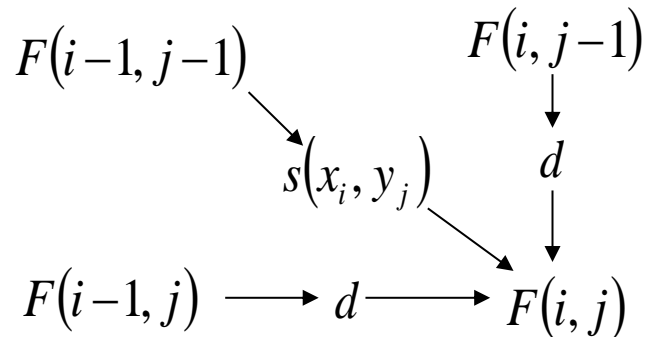
Use a gap penalty of $d = -5$.

		A	A	G
	0	0	0	0
G	0	0	0	2
A	0	2	2	0
A	0	2	4	0
G	0	0	0	6
G	0	0	0	2
C	0	0	0	0

Compare with the Best GLOBAL Alignment

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2

$d = -5$



(contrast with the best
local alignment)

Find the optimal **Global** alignment of
AAG and **GAAGGC**.

Use a gap penalty of $d = -5$.

		A	A	G
	0	-5	-10	-15
G	-5			
A	-10			
A	-15			
G	-20			
G	-25			
C	-30			

Summary

Global alignment algorithm:

Needleman-Wunsch.

		G	A	A	T	C
	0	-4	-8	-12	-16	-20
C	-4	-5	-9	-13	-12	-6
A	-8	-4	5	1	-3	-7
T	-12	-8	1	0	11	7
A	-16	-12	2	11	7	6
C	-20	-16	-2	7	11	17

Local alignment algorithm:

Smith-Waterman.

		A	A	G
	0	0	0	0
G	0	0	0	2
A	0	2	2	0
A	0	2	4	0
G	0	0	0	6
G	0	0	0	2
C	0	0	0	0

Using sequence alignment to study evolution

Are these proteins related?

The intuitive answer:

SEQ 1: RVVNLVPS--FWVL DATYKNYA INYNCDV TYKLY

L P L Y N Y C L

score = -1 → NO?

SEQ 2: QFFPLMPPAPYFILATDYENLPLVYSC TTFFWLF

SEQ 1: RVVNLVPS--FWVL DATYKNYA INYNCDV TYKLY

L P W L DATYKN A Y C L

score = 15 → PROBABLY?

SEQ 2: QFFPLMPPAPYWIL DATYKNLALVYSC TTFFWLF

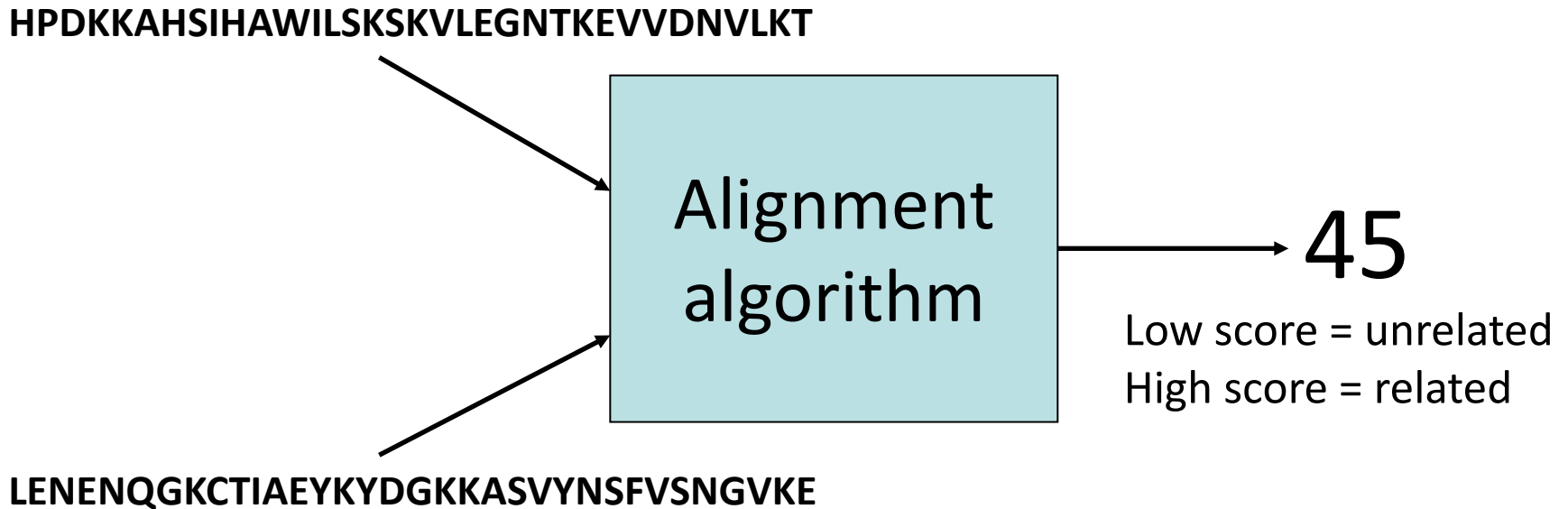
SEQ 1: RVVNLVPS--FWVL DATYKNYA INYNCDV TYKLY

RVV L PS W L DATYKNYA Y CDV TYKL

score = 37 → YES?

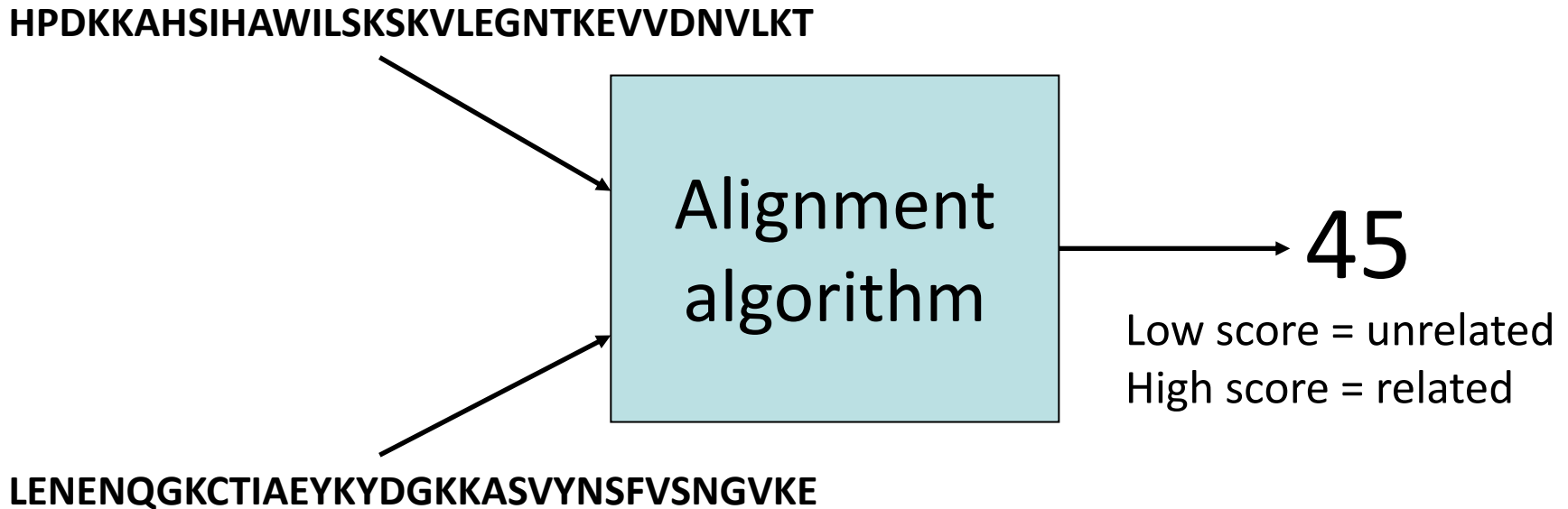
SEQ 2: RVVPLMPSAPYWIL DATYKNYALVYSC DV TYKLF

Significance of scores



But ... how high is high enough?

Significance of scores



But ... how high is high enough?

Subjective

Problem specific

Parameter specific

The null hypothesis

- We want to know how **surprising** a given score is, ...
assuming that the two sequences are not related.
- This assumption is called the **null hypothesis**.
- The purpose of most statistical tests is to determine whether the observed result provides a reason to reject the null hypothesis.
- We want to characterize the distribution of scores from pairwise sequence alignments.

Sequence similarity score distribution



- Search a database of **unrelated** sequences using a given query sequence.
- What will be the form of the resulting distribution of pairwise alignment scores?

