

# Sequence alignments and scoring matrices

Arne Elofsson

To read: [http://perso.fundp.ac.be/~lambertc/DEA-bioinfo/CLambert\\_curr\\_gen\\_2003.pdf](http://perso.fundp.ac.be/~lambertc/DEA-bioinfo/CLambert_curr_gen_2003.pdf)

To read: Wikipedia about Sequence Alignment

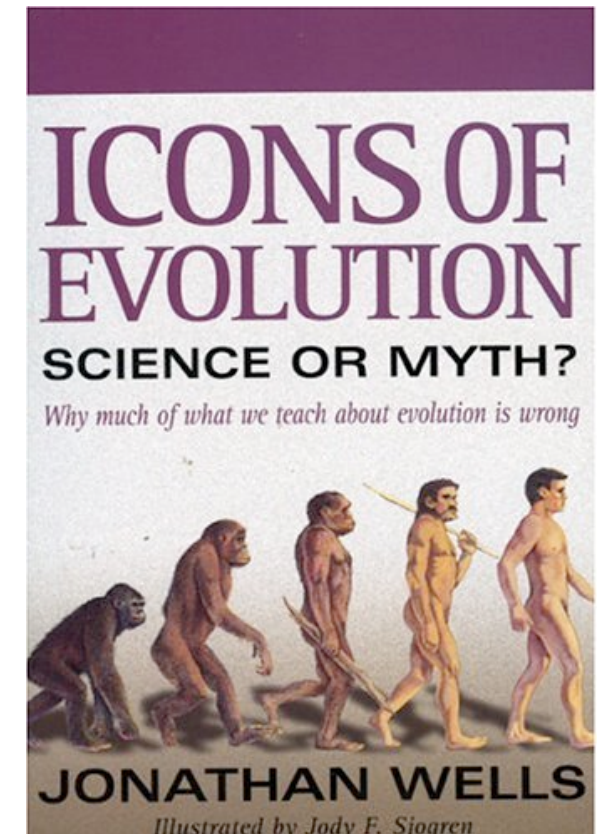
# Why alignments ?

```

AAB24882      TYHMCQFHCERYVNNHSGEKLIECNERSKAFSCPSHLQCHKRRQIGKTHEHNQCGKAFPT 60
AAB24881      -----YECNQCGKAFAQHSSLKCHYRTHIGKPYECNQCGKAFSK 40
                ****: .***: * *:*** * :****. :* *****,.

AAB24882      PSHLQYHERTHTGKPYECHQCGQAFKKCSLLQRHKRTHTGKPYE-CNQCGKAFAQ- 116
AAB24881      HSHLQCHKRTHTGKPYECNQCGKAFSQHGLLQRHKRTHTGKPYMNVINMVKPLHNS 98
                **** *:*****:****:*. : .*****: : *: :
    
```

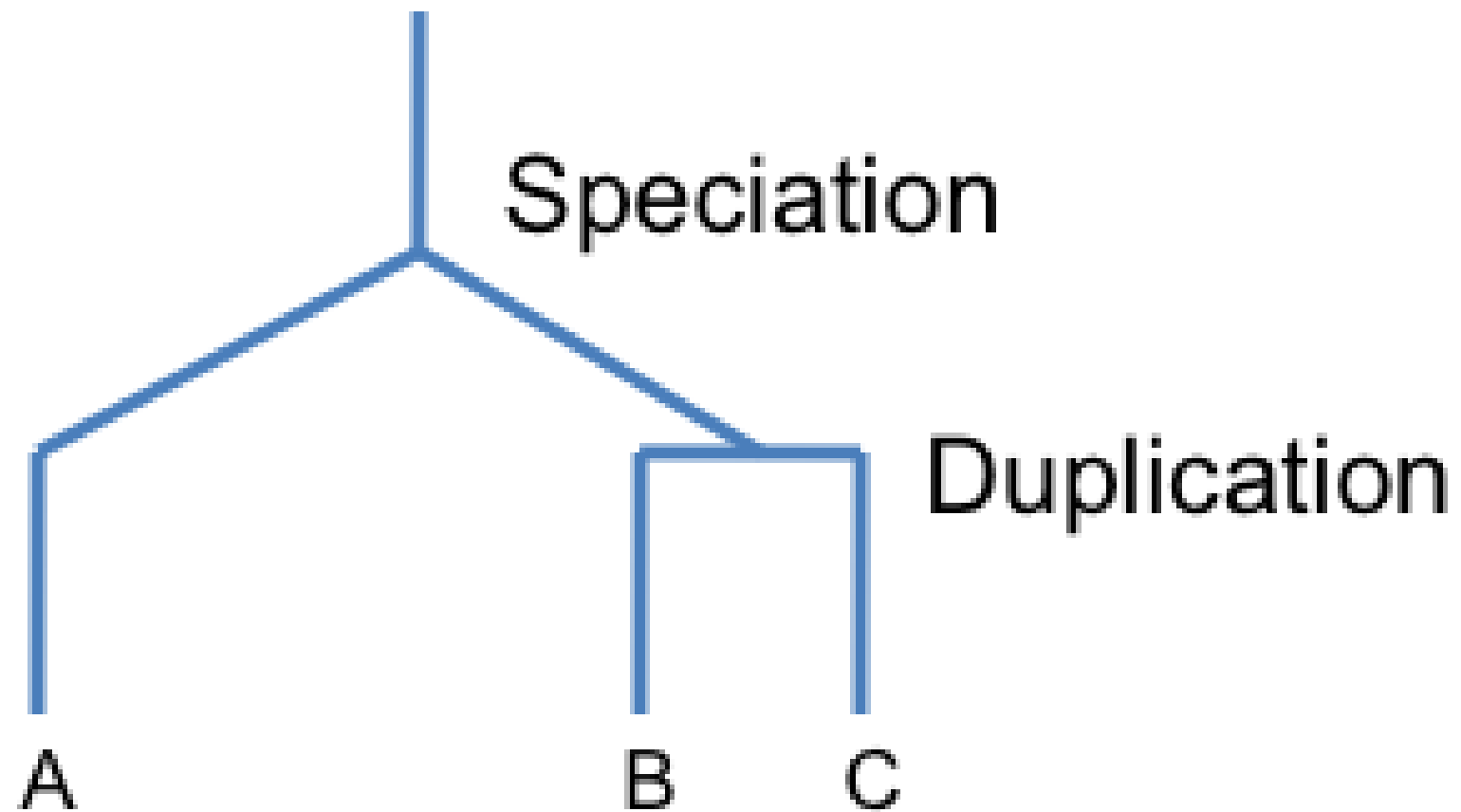
- Detect homology
- Study evolution
- Predict functions
- Model 3D-structure



# Sequence similarity

- Homologs have a common ancestor
- Gene duplication or speciation
- High sequence similarity indicates homology
- Homologs have similar 3D-structure

# Homology



# Convergent evolution



# What is an alignment

THISSEQUENCE

| |   | | | | | | |

10/12 Identical

THATSEQUENCE

THATSEQUENCE

| |     |   |

4/12 Identical

THISISASEQUENCE

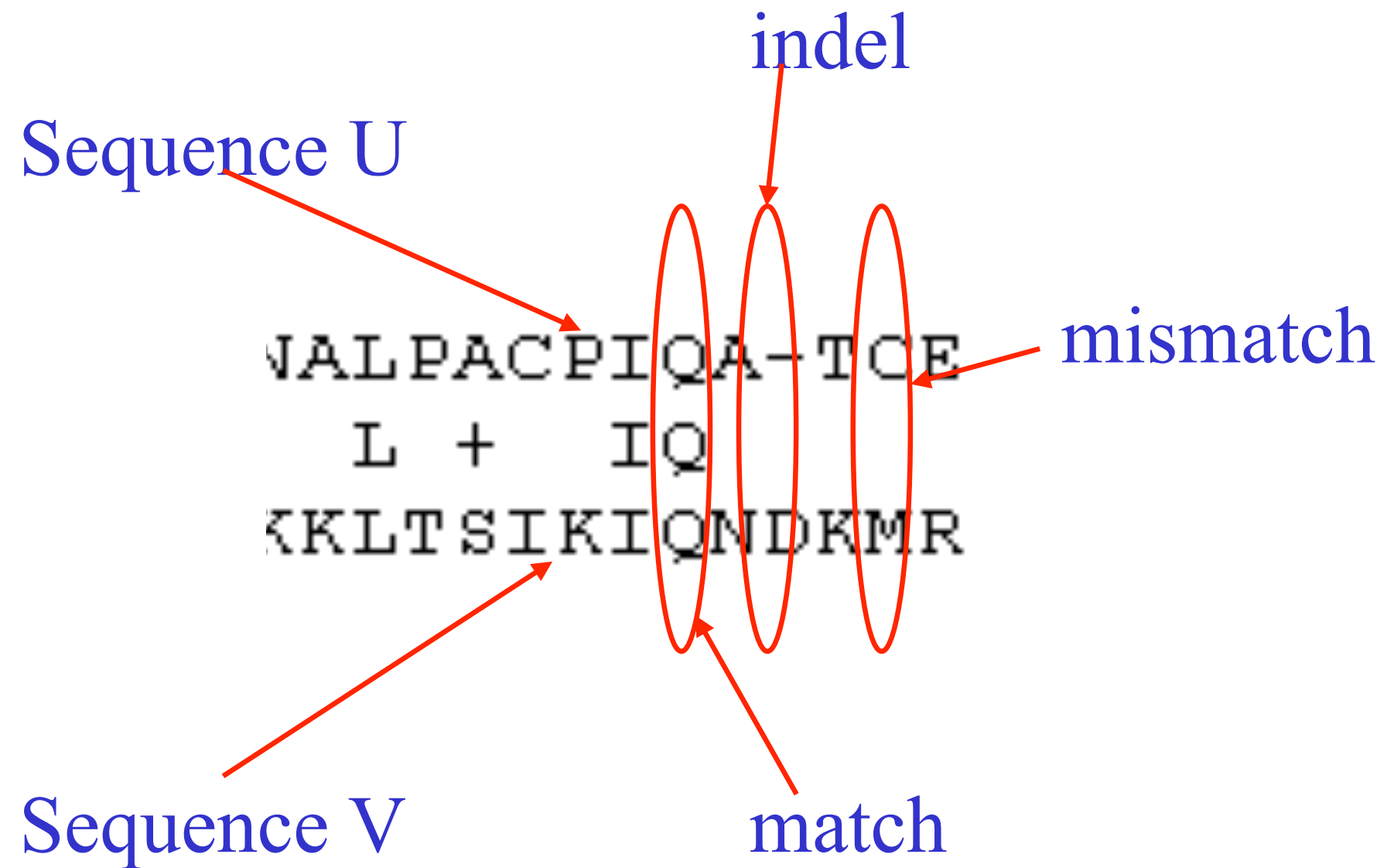
THISISA-SEQUENCE

| |   | | | | | | |

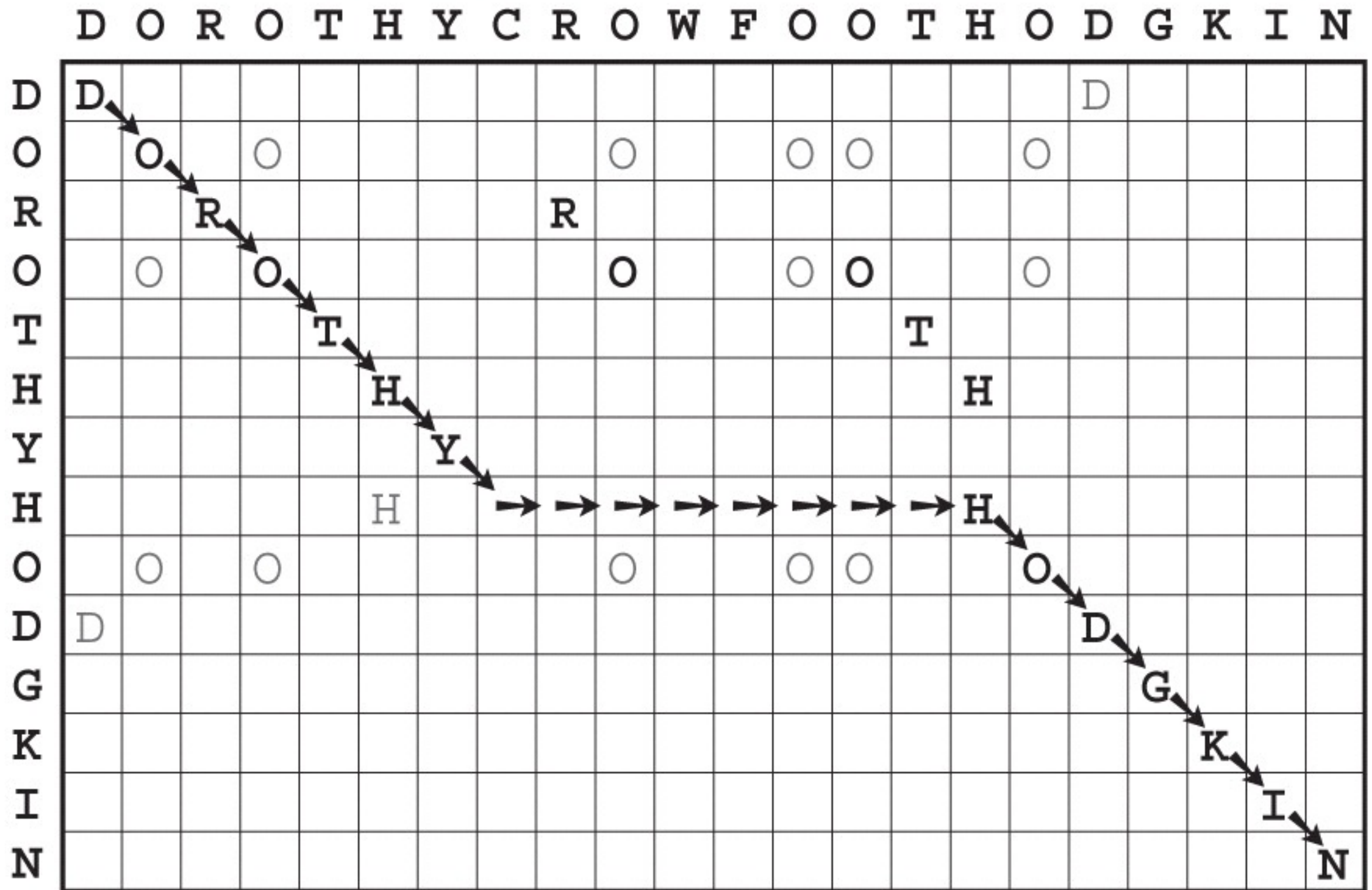
11/12 Identical

TH---ATSEQUENCE

# What can an alignment say ?



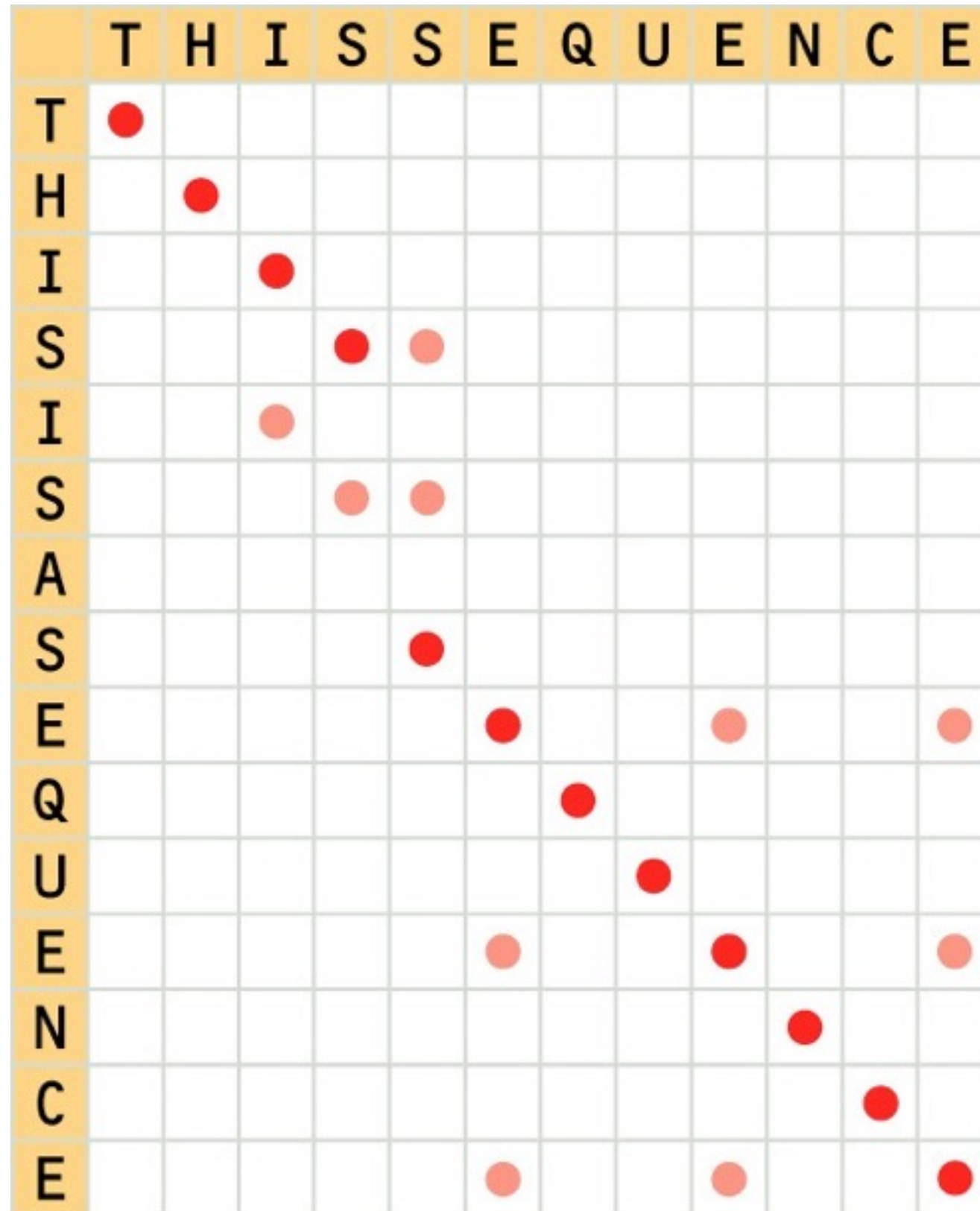
# An alignment matrix





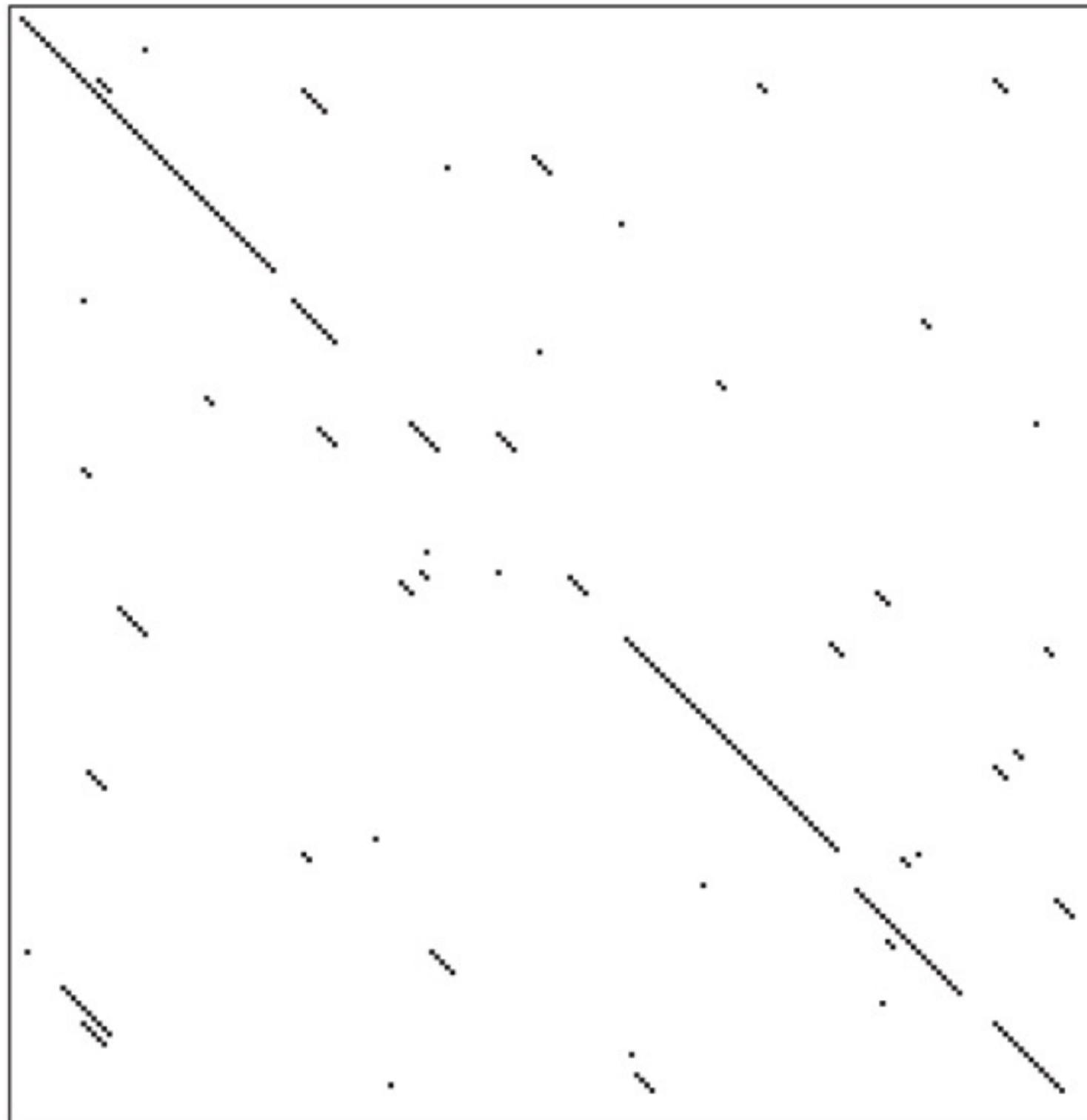
# Dotplots

# Dotplots



# Dotplots

PAPA\_CARPA / ACTN\_ACTCH



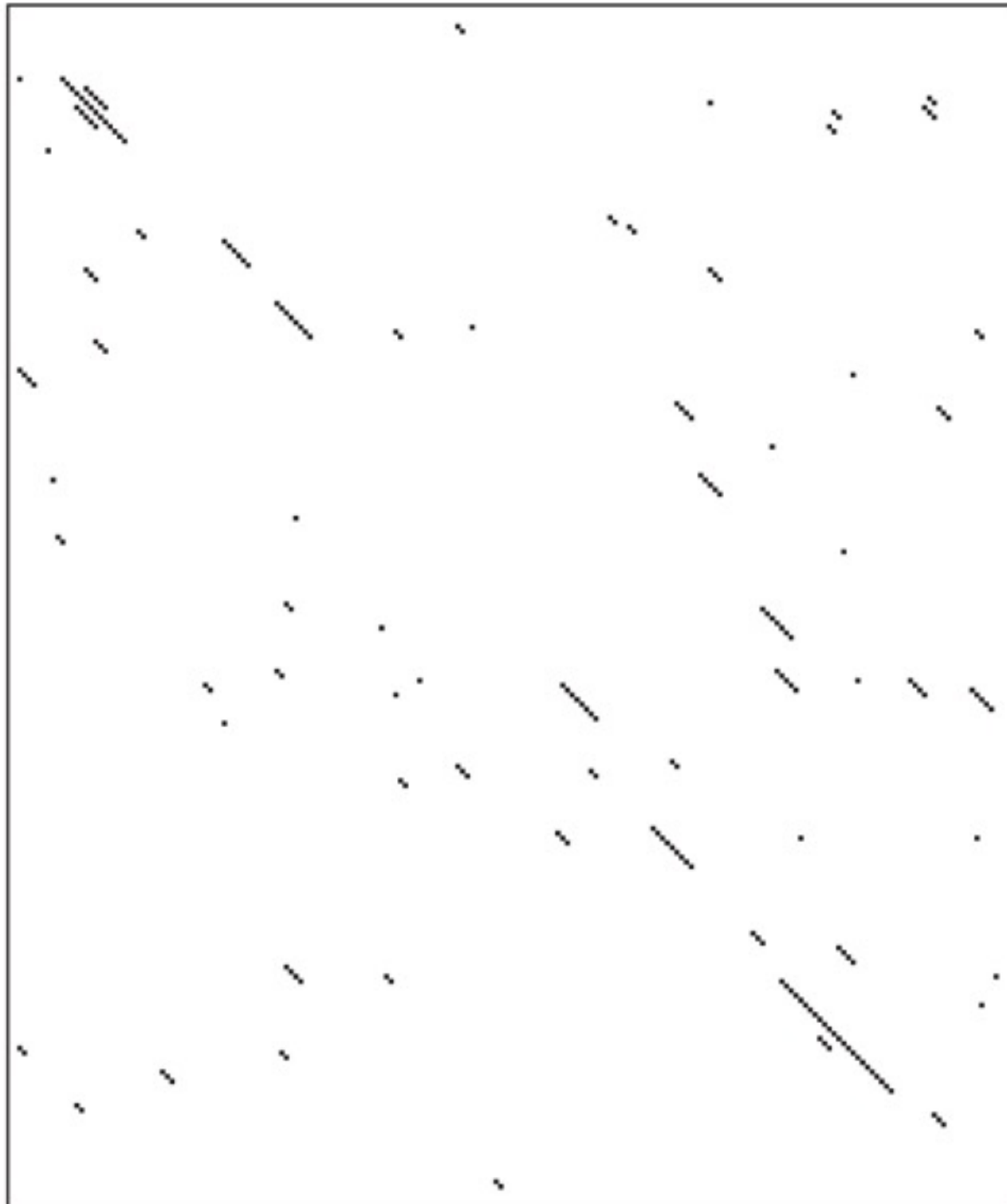
# Dotplots

PAPA\_CARPA / CATL\_HUMAN



# Dotplots

PAPA\_CARPA / CATB\_HUMAN

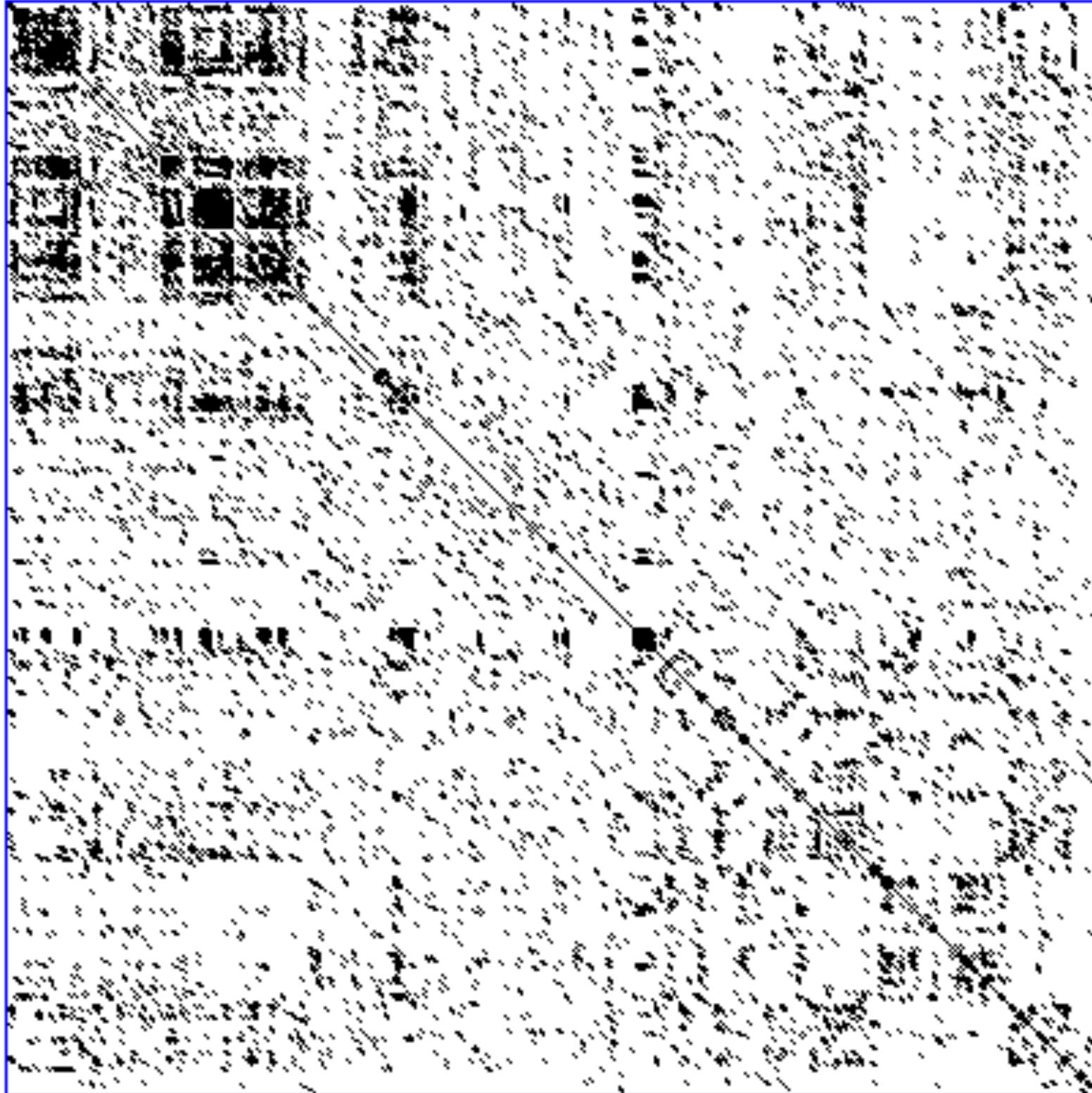


# Dotplots

PAPA\_CARPA / STPA\_STAAU



# Dotplots



# Types of alignment



# Types of alignment

## (A) local

PI3-kinase DRHNSNIMVKDDGQLFHI DFG  
cAMP PK DLKPENLLIDQQGYIQVT DFG

## (B) global

PI3-kinase HQLGNLR--LEEERI---MSSAKRPLWLNWENPDIMSELLFQNNETIFKNGDDLRRQDMLT  
cAMP PK GNAAAAGKKGXEQESVKEFLAKAKEDFLKKWENPAQNTAHLDDQFERIKTLGTGSFGRVML-

PI3-kinase LQIIRIME--NIWQNGGLDLRMLPYGCLSIGDCVGLIEVVRNSHTIMQ-IQCKGGLKGAAL  
cAMP PK ---VKHMETGNHYAMKILDKQKVVK-----LKQIEHTLNEKRILQAVNFPFLVKLEF

PI3-kinase QFNSHT-LHQWLKDKNKGEIYDAA--IDLFTTRSCAGYCVATFILGIGDRHNSNIMVKD-D  
cAMP PK SFKDNSNLYMVMEYVPGGEMFSLRRIIGRFSEPHARFYAAQIVLTFFEYLHSLDLIYRDLK

PI3-kinase GQLFHI DFGHFLDHKKKKFGYKRERVP-----FVLTQDFL---IVISKGAQECTKTREFE  
cAMP PK PENLLIDQQGYI--QVT DFGFAK-RVKGRTWXLCTPEYLAPEIILSKGYNKAVDWWALG

PI3-kinase RF-QEMC--YKAYLAIRQHANLFINLFSMMLGSGMPQLQSFDDIAYIRKTLALDKTEQEA  
cAMP PK VLIYEMAAGYPPFFA-DQPIQIYEKIVSGKVR--FPSHFSSDLKDLLRNLLQVDLTKR--

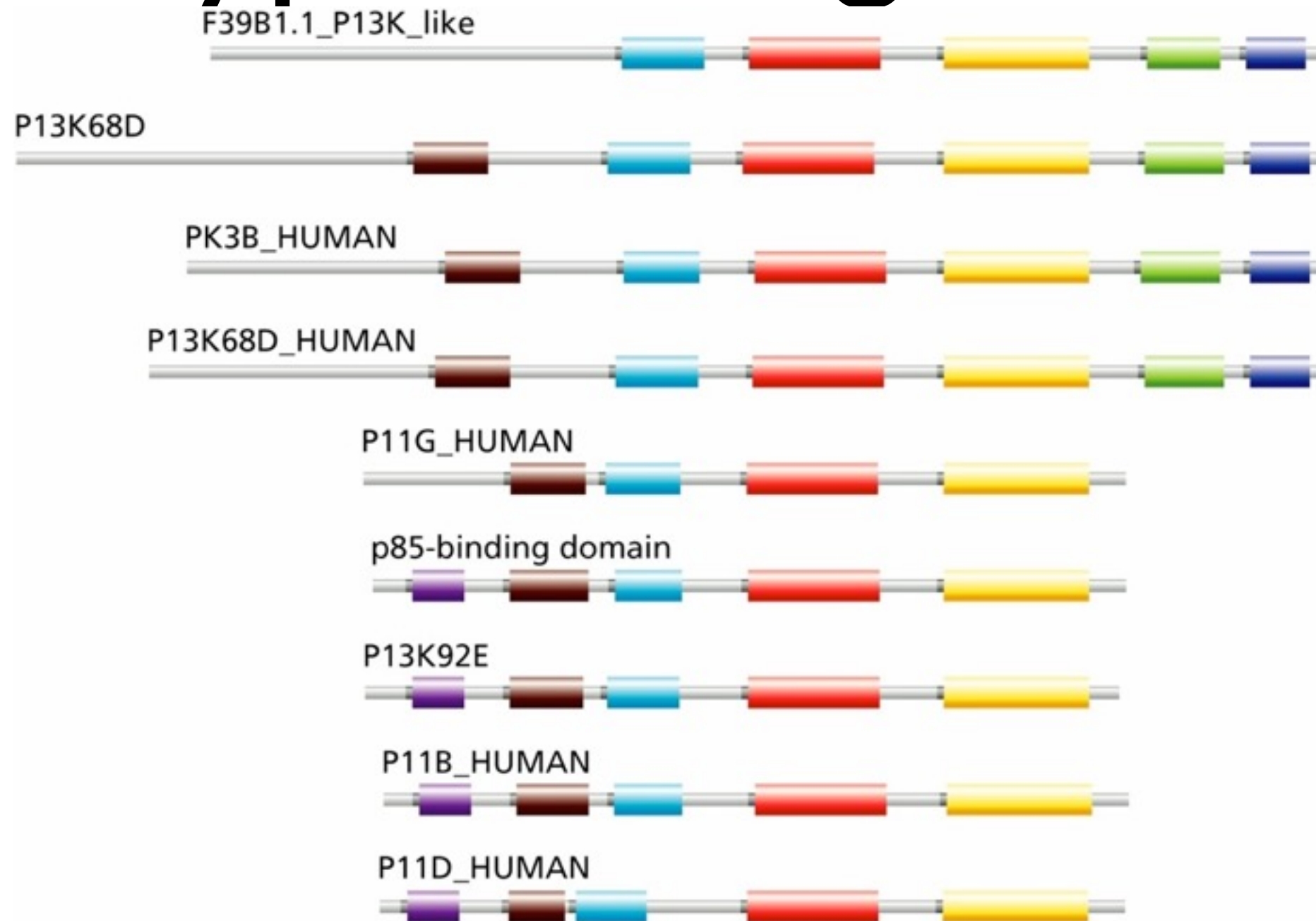
PI3-kinase LEYFMKQMNDAAHHGGWTTKMDWI-----FHTIKQHALN-----  
cAMP PK FGNLKNQVNDIKNHKWFATTDWIAIYQRKVEAPFIPKFKGPGDTSNFDDEEEEEIRVXIN

# Types of alignment

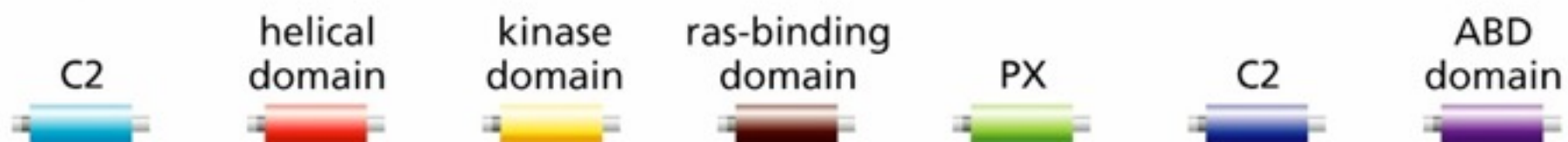
Global    FTFTALILLAVAV  
          F--TAL--LLA--AV

Local     FTFTALILL--AVAV  
          --FTAL--LLAAV--

# Types of alignment



KEY:



# Inserting gaps

(A)

```
Bovine PI-3Kinase p110a      LNWENPDIMSELLFQNNELIFKNGDDLRRQDMLTLQIIRIMENIWQNGGLDLRMLPYGCLSIGDCVGLIEVVRNSHTIMQIQCKGGLKGAL
cAMP-dependent protein kinase --WENPAQNTAHLDDQFERIKTLGTGSFGRVMLVKHMETGNHYAMKILDKQKVVKLKQIEHTLNEKRILQAVNFPFLVKLEFSFKDNSNLY

Bovine PI-3Kinase p110a      QFNSHTLHQWLKDKNKGEIYDAAIDLFTTRSCAGYCVATFILGIGDRHNSNIMVKDDGQLFHIDFGHFLDHKKKKFGYKRERVPFVLTQDF
cAMP-dependent protein kinase MVMEYVPGGEMFSLRRIGRFSEPHARFYAAQIVLTFEYLHSLDLIYRDLKPENLLIDQQGYIQVTDFGFAKRVKGRTWXLCGTPEYLAP

Bovine PI-3Kinase p110a      LIVISKGAQECTKTREFERFQEMCYKAYLAIRQHANLFINLFSMMLGSGMPELQSFDDIAYIRKTLALDKTEQEAELEYFMKQMNDAAHHGG
cAMP-dependent protein kinase EIILSKGYNKAVDWWALGVLIYEMAAGYPPFFADQPIQIYEKIVSGKVRFP SHFSSDLKDLLRNLLQVDLTKRFGNLKNGVNDIKNHKWF

Bovine PI-3Kinase p110a      WTTKMDWIFHTIKQHALN-----
cAMP-dependent protein kinase ATTDWIAIYQRKVEAPFIPKFKGPGDTSNFDDYEEEEIRVXINEKCGKEFSEF
```

(B)

```
Bovine PI-3Kinase p110a      LNWENPDIMSELLFQNNELIFKNGDDLRRQDMLTLQIIRIMENIWQNGGLDLRMLPYGCLSIGDCVGLIEVVRNSHTIMQIQCKGGLKGAL
cAMP-dependent protein kinase ?-WENPAQNTAHLDDQFERIKTLGTGSFGRVMLVKHM--ETGNHYAMKILDKQKV-VKLLQIEHTLNEKRILQAVNFPFLVKLEFSFKDN-

Bovine PI-3Kinase p110a      QFNSHTLHQWLKDKNKGEIYDAAIDLFTTRSCAGYCVATFILGIGDRHNSNIMVKD-DGQLFHIDFGHFLDHKKKKFGYKRERVPFVL--T
cAMP-dependent protein kinase -SNLYMVMEYVPGGEMFSLRR-IGRFSEPHARFYAAQIVLTFEYLHSLDLIYRDLKPENLLIDQQGYIQVTDFGFAKRVKGRTWXLCGT

Bovine PI-3Kinase p110a      QDFL---IVISKGAQECTKTREFERF-QEMC--YKAYLAIRQHANLFINLFSMMLGSGMPELQSFDDIAYIRKTLALDKTEQEAELEYFMK
cAMP-dependent protein kinase PEYLAPEIILSKGYNKAVDWWALGVLIYEMAAGYPPFFA-DQPIQIYEKIVSGKVRFP--PSHFSSDLKDLLRNLLQVDLTKR--FGNLKN

Bovine PI-3Kinase p110a      QMNDAAHHGGWTTKMDWI-----FHTIKQHAL---N-----
cAMP-dependent protein kinase GVNDIKNHKWFATTDWIAIYQRKVEAPFIPKFKGPGDTSNFDDYEEEEIRVXINEKCGKEFSEF
```

# What is an optimal alignment ?

T H I S S E Q U E N C E

| |       | | | | | | |

10/12 Identical

T H A T S E Q U E N C E

T H A T S E Q U E N C E

| |                       |       |

4/12 Identical

T H I S I S A S E Q U E N C E

T H I S I S A - S E Q U E N C E

| |                       |       | | | | | | |

11/12 Identical

T H - - - - A T S E Q U E N C E



# Different scoring

T H I S S E Q U E N C E

5 8 -1 1 4 5 6 0 5 6 9 5

Score = 52

T H A T S E Q U E N C E

T H A T S E Q U E N C E

5 8 -1 -1 -2 0 -1 0 5 0 0 5

Score = 18

T H I S I S A S E Q U E N C E

T H I S I S A - S E Q U E N C E

5 8 0 0 0 0 4 0 4 5 6 0 5 6 9 5

Score = 56

T H - - - A T S E Q U E N C E

# With Gap cost

T H I S S E Q U E N C E

5 8-1 1 4 5 6 0 5 6 9 5

Score = 52

T H A T S E Q U E N C E

T H A T S E Q U E N C E

5 8-1-1-2 0-1 0 5 0 0 5

Score = 18

T H I S I S A S E Q U E N C E

T H I S I S A - S E Q U E N C E

5 8-1-1-1-1 4-1 4 5 6 0 5 6 9 5

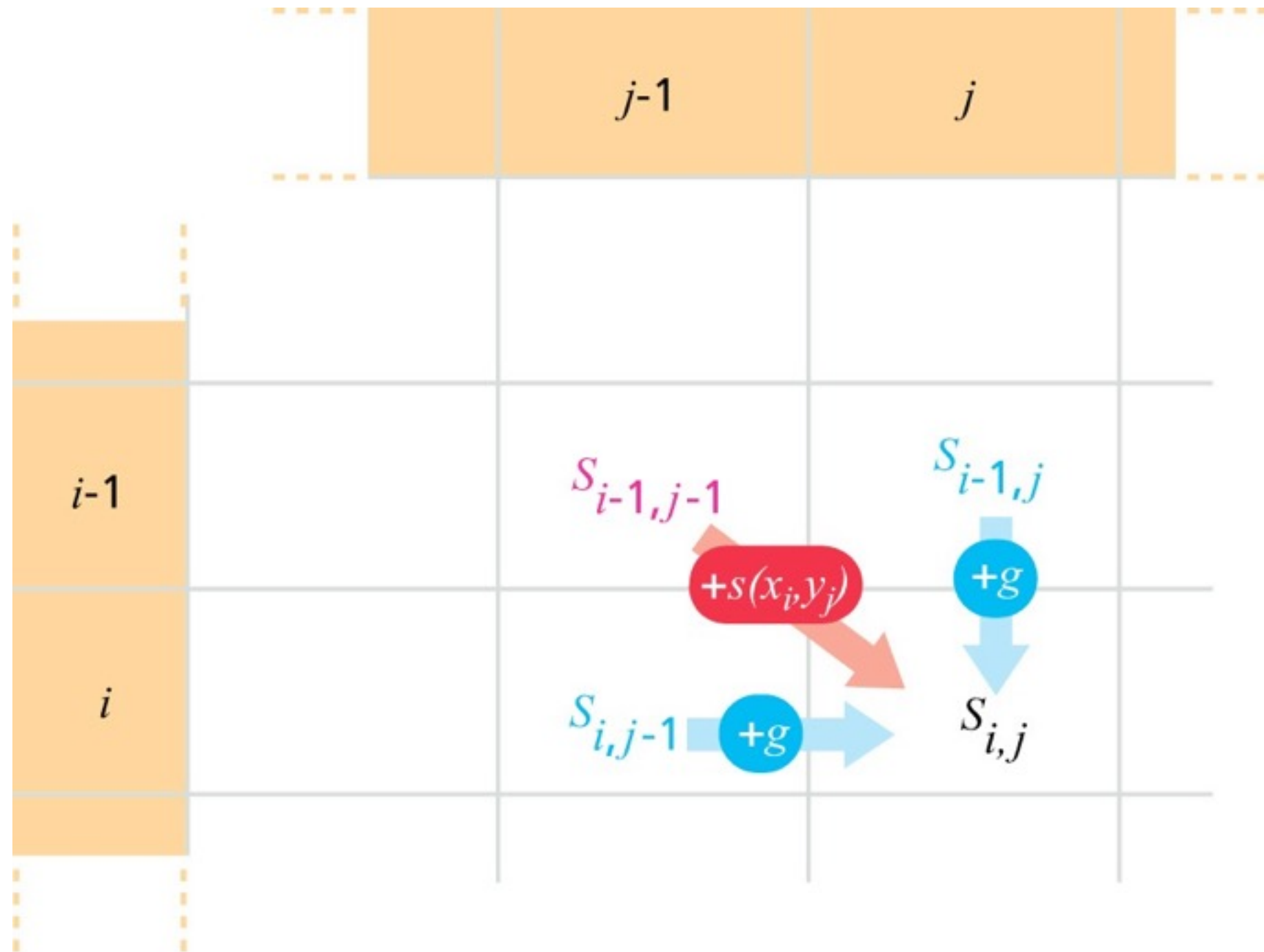
Score = 51

T H - - - A T S E Q U E N C E

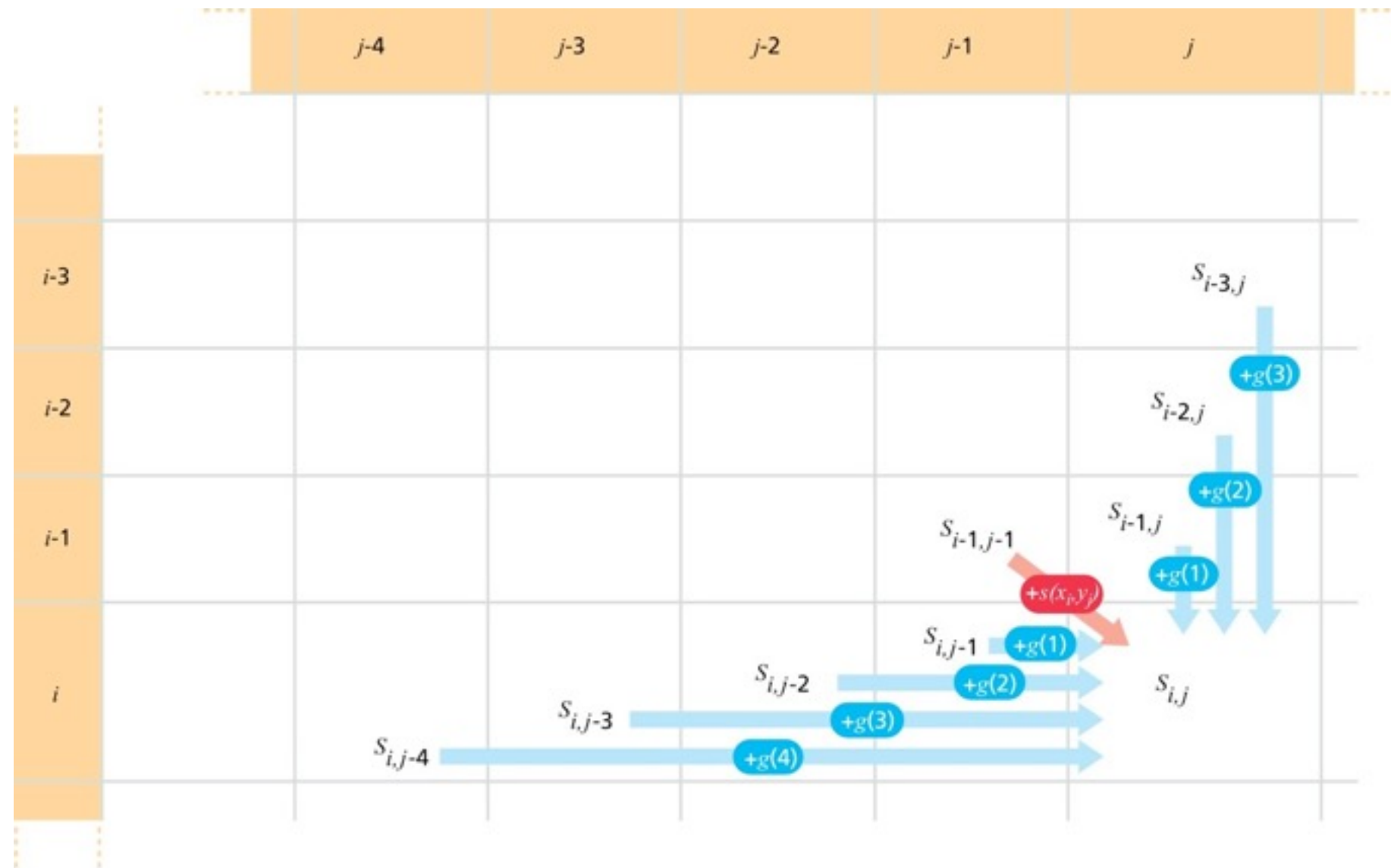
# Dynamic programming



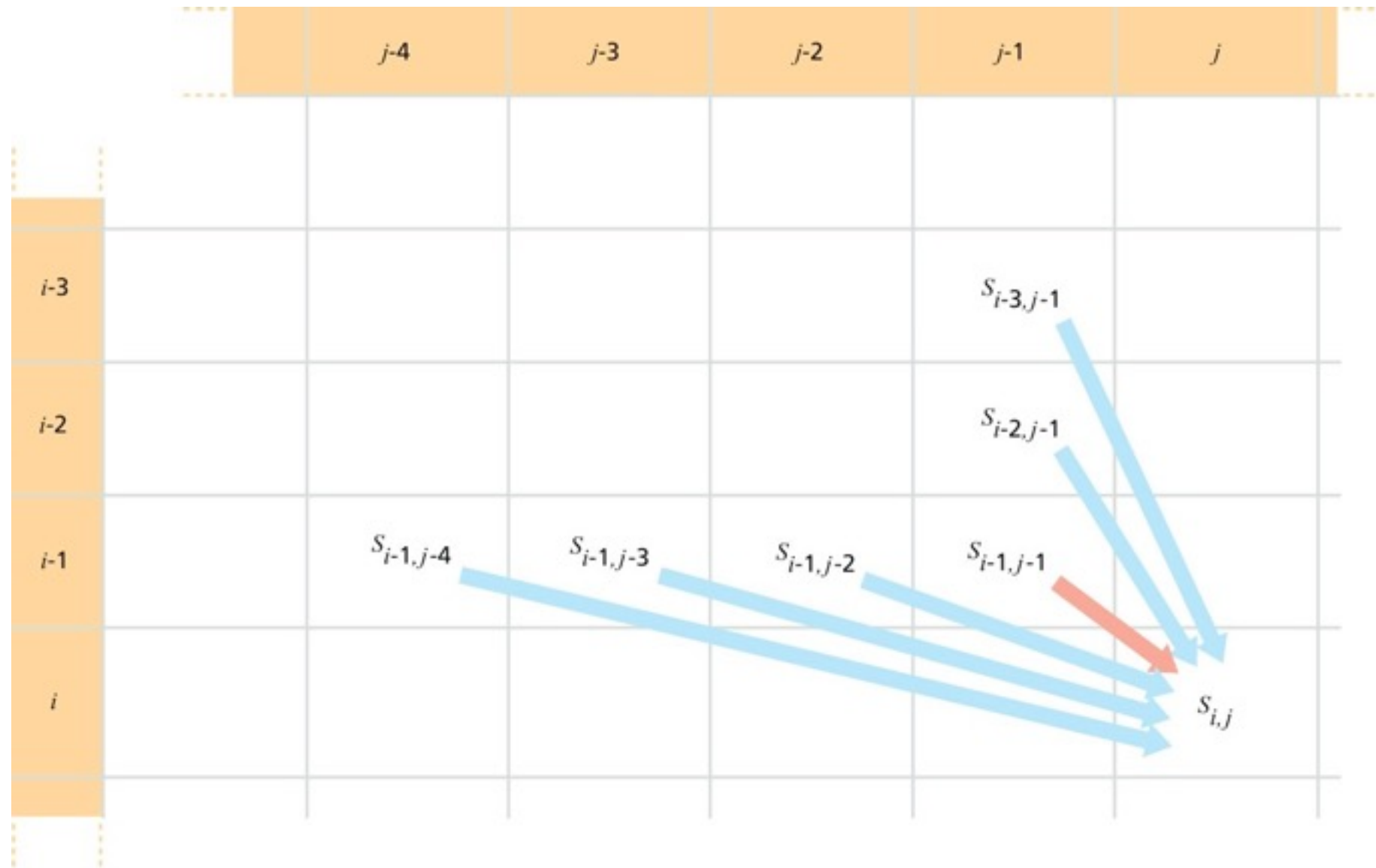
# Dynamic programming



# Dynamic programming



# Dynamic programming

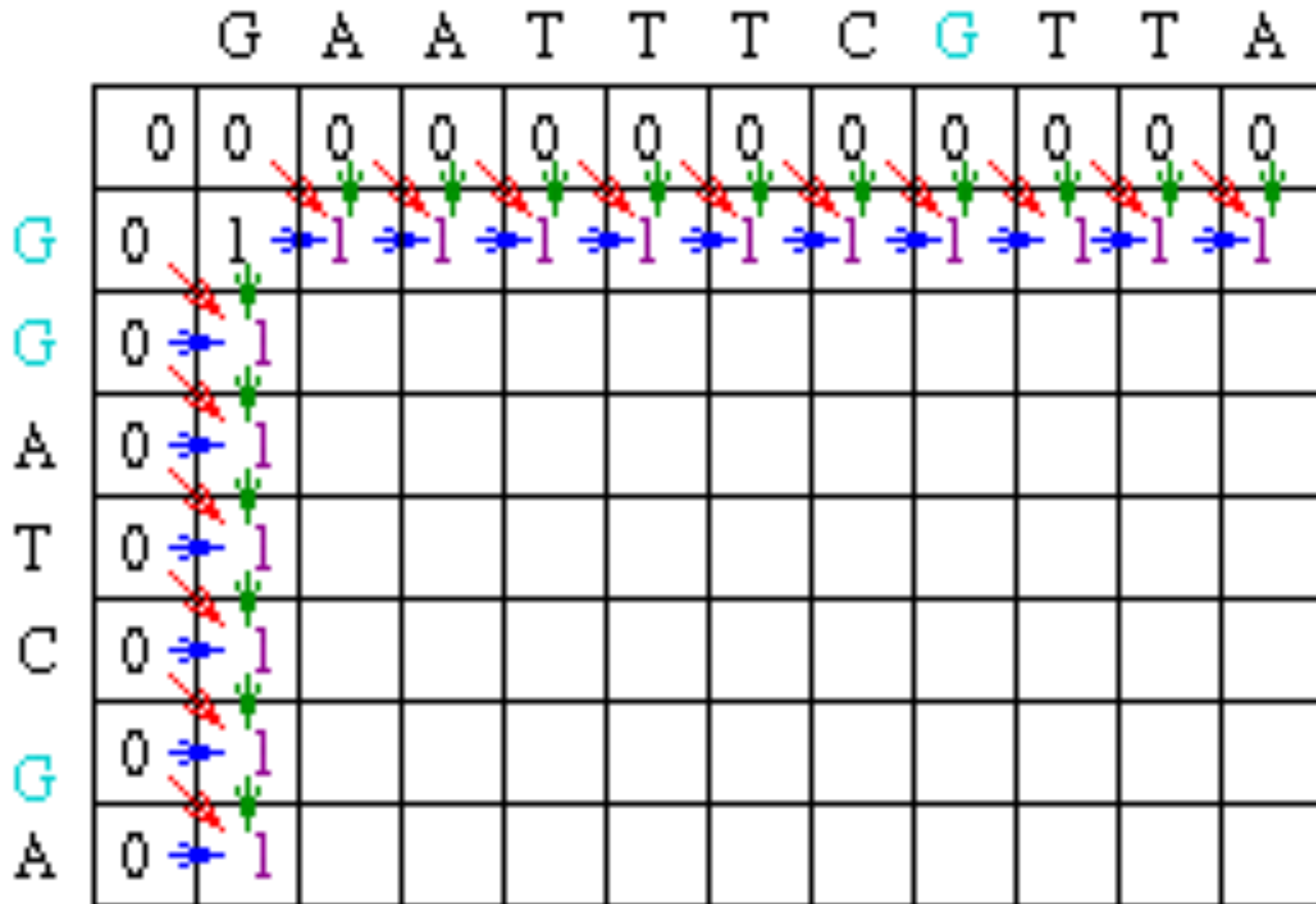


Initialisation step: Create Matrix with  $M + 1$  columns and  $N + 1$  rows. First row and column filled with 0.

[illegible]

$$M_{i,j} = \text{MAXIMUM [}$$
$$M_{i,j-1} + w \text{ (gap in sequence \#1)}$$
$$M_{i-1,j} + w(\text{gap in sequence \#2})]$$
[illegible]

## Fill in rest of row 1 and column 1





Fill in column 3

		G	A	A	T	T	C	A	G	T	T	A
		0	0	0	0	0	0	0	0	0	0	0
G		0	1	1	1	1	1	1	1	1	1	1
G		0	1	1								
A		0	1	2								
T		0	1	2								
C		0	1	2								
G		0	1	2								
A		0	1	2								



Column 3 with answers

		G	A	A	T	T	C	A	G	T	T	A
		0	0	0	0	0	0	0	0	0	0	0
G		0	1	1	1	1	1	1	1	1	1	1
G		0	1	1								
A		0	1	2								
T		0	1	2								
C		0	1	2								
G		0	1	2								
A		0	1	2								


Fill in rest of matrix with answers

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

Traceback step:

Position at current cell and look at direct predecessors

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



Traceback step:

Position at current cell and look at direct predecessors

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

Seq#1 A

|

Seq#2 A

## Position at current cell and look at direct predecessors

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

## Position at current cell and look at direct predecessors

[illegible]

## Position at current cell and look at direct predecessors

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

## Position at current cell and look at direct predecessors

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



## Position at current cell and look at direct predecessors

[illegible]

## Position at current cell and look at direct predecessors

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

## Position at current cell and look at direct predecessors

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

## Position at current cell and look at direct predecessors

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

## Position at current cell and look at direct predecessors

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

## Position at current cell and look at direct predecessors

[illegible]



Traceback step:

Position at current cell and look at direct predecessors

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

Seq#1	G	A	A	T	T	C	A	G	T	T	A
Seq#2	G	G	A	T	-	C	-	G	-	-	A



# Pseudocode

```
for i=0 to length(A)
    F(i,0) ← d*i
for j=0 to length(B)
    F(0,j) ← d*j
for i=1 to length(A)
    for j=1 to length(B)
    {
        Match ← F(i-1,j-1) + S(Ai, Bj)
        Delete ← F(i-1, j) + d
        Insert ← F(i, j-1) + d
        F(i,j) ← max(Match, Insert,
Delete)
    }
```

# Traceback

```
AlignmentA ← ""
AlignmentB ← ""
i ← length(A)
j ← length(B)
while (i > 0 or j > 0)
{
if (i > 0 and j > 0 and F(i,j) == F(i-1,j-1) + S(Ai, Bj))
{
AlignmentA ← Ai + AlignmentA
AlignmentB ← Bj + AlignmentB
i ← i - 1
j ← j - 1
}
else if (i > 0 and F(i,j) == F(i-1,j) + d)
{
AlignmentA ← Ai + AlignmentA
AlignmentB ← "-" + AlignmentB
i ← i - 1
}
else (j > 0 and F(i,j) == F(i,j-1) + d)
{
AlignmentA ← "-" + AlignmentA
AlignmentB ← Bj + AlignmentB
j ← j - 1
}
}
```

# Substitution matrices

# Substitution matrices

## Scoring Matrices

$S = [s_{ij}]$  gives score of aligning character  $i$  with character  $j$  for every pair  $i, j$ .

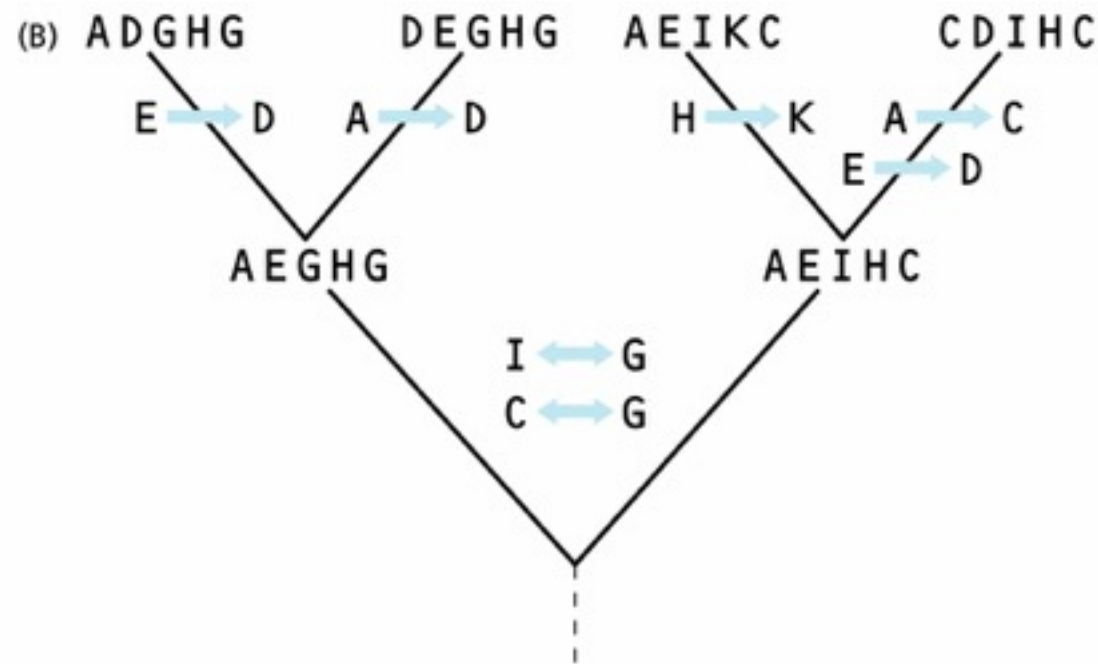
C	12				
S	0	2			
T	-2	1	3		
P	-3	1	0	6	
A	-2	1	1	1	2
	C	S	T	P	A

STPP  
CTCA

$$0 + 3 + (-3) + 1 = 1$$

# Substitution matrices

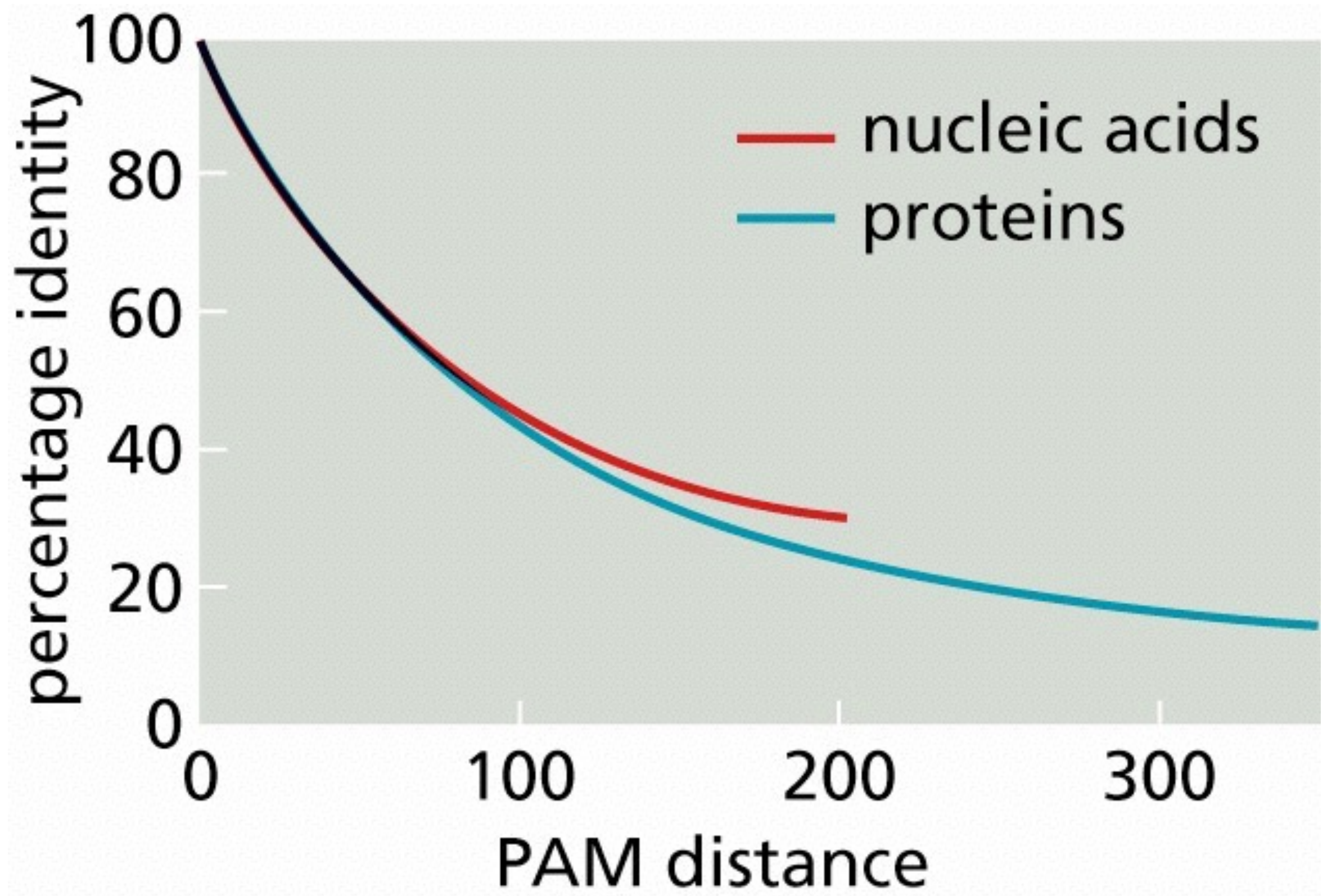
(A) DEGHG  
ADGHG  
CDIHC  
AEIKC



(C)

	A	C	D	E	G	H	I	K
A								
C		1	1					
D	1				1			
E				2				
G		1					1	
H								1
I					1			
K								1

# Substitution matrices



# Log Odds Ratios

$$S_{i,j} = \log \frac{p_i \cdot M_{i,j}}{p_i \cdot p_j} = \log \frac{M_{i,j}}{p_j} = \log \frac{\textit{observed frequency}}{\textit{expected frequency}}$$

# Point Accepted Mutations (PAM)

$$\text{PAM}_n(i, j) = \log \frac{f(i)M^n(i, j)}{f(i)f(j)} = \log \frac{M^n(i, j)}{f(j)}$$



# Point Accepted Mutations (PAM)

$$\text{PAM}_n(i, j) = \log \frac{f(i)M^n(i, j)}{f(i)f(j)} = \log \frac{M^n(i, j)}{f(j)}$$

## The PAM Family

Define a *family* of substitution matrices — PAM 1, PAM 2, etc. — where PAM  $n$  is used to compare sequences at distance  $n$  PAM.

$$\text{PAM } n = (\text{PAM } 1)^n$$

**Do not confuse with scoring matrices!**

Scoring matrices are derived from PAM matrices to yield log-odds scores.

# Point Accepted Mutations (PAM)

$$\text{PAM}_n(i, j) = \log \frac{f(i)M^n(i, j)}{f(i)f(j)} = \log \frac{M^n(i, j)}{f(j)}$$

## PAM matrices

- Let  $M$  be a PAM 1 matrix. Then,

$$\sum_i p_i (1 - M_{ii}) = 0.01$$

- **Reason:**  $M_{ii}$ s are the probabilities that a given amino acid does not change, so  $(1 - M_{ii})$  is the probability of mutating away from  $i$ .

# Point Accepted Mutations (PAM)

$$\text{PAM}_n(i, j) = \log \frac{f(i)M^n(i, j)}{f(i)f(j)} = \log \frac{M^n(i, j)}{f(j)}$$

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

# Blosum

# Blosum

(A)

	1	2	3	4	5
1	A	T	C	K	Q
2	A	T	C	R	N
3	A	S	C	K	N
4	S	S	C	R	N
5	S	D	C	E	Q
6	S	E	C	E	N
7	T	E	C	R	Q

(B)

	$q_{QN}$	$q_{NN}$	$q_{QQ}$	$p_N$	$p_Q$
$C=62\%$	0.114	0.057	0.029	0.114	0.086
$C=50\%$	0.117	0.025	0.058	0.084	0.117
$C=40\%$	—	—	—	—	—

# Blosum

**Equivalent PAM and Blossum matrices (according to *H*)**

- PAM100 ==> Blosum90
- PAM120 ==> Blosum80
- PAM160 ==> Blosum60
- PAM200 ==> Blosum52
- PAM250 ==> Blosum45

# Blosum

C	11																			
S	1	2																		
T	-1	1	2																	
P	-2	1	1	6																
A	-1	1	2	1	2															
G	-1	1	-1	-1	1	5														
N	-1	1	1	-1	0	0	3													
D	-3	0	-1	-2	0	1	2	5												
E	-4	-1	-1	-2	-1	0	1	4	5											
Q	-3	-1	-1	0	-1	-1	0	1	2	5										
H	0	-1	-1	0	-2	-2	1	0	0	2	6									
R	-1	-1	-1	-1	-1	0	0	-1	0	2	2	5								
K	-3	-1	-1	-2	-1	-1	1	0	1	2	1	4	5							
M	-2	-1	0	-2	-1	-3	-2	-3	-3	-2	-2	-2	-2	6						
I	-2	-1	1	-2	0	-3	-2	-3	-3	-3	-3	-3	-3	3	4					
L	-3	-2	-1	0	-1	-4	-3	-4	-4	-2	-2	-3	-3	3	2	5				
V	-2	-1	0	-1	1	-2	-2	-2	-2	-3	-3	-3	-3	2	4	2	4			
F	0	-2	-2	-3	-3	-5	-3	-5	-5	-4	0	-4	-5	0	0	2	0	8		
Y	2	-1	-3	-3	-3	-4	-1	-2	-4	-2	4	-2	-3	-2	-2	-1	-3	5	9	
W	1	-3	-4	-4	-4	-2	-5	-5	-5	-3	-3	0	-3	-3	-4	-2	-3	-1	0	15
	C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W

# Difference between Pam and Blosum

- PAM matrices are based on an explicit evolutionary model (i.e. replacements are counted on the branches of a phylogenetic tree), whereas the BLOSUM matrices are based on an implicit model of evolution.
- The PAM matrices are based on mutations observed throughout a global alignment, this includes both highly conserved and highly mutable regions. The BLOSUM matrices are based only on highly conserved regions in series of alignments forbidden to contain gaps.
- The method used to count the replacements is different: unlike the PAM matrix, the BLOSUM procedure uses groups of sequences within which not all mutations are counted the same.
- Higher numbers in the PAM matrix naming scheme denote larger evolutionary distance, while larger numbers in the BLOSUM matrix naming scheme denote higher sequence similarity and therefore smaller evolutionary distance. Example: PAM150 is used for more distant sequences than PAM100; BLOSUM62 is used for closer sequences than BLOSUM50.



# Nucleotide Matrices

Dayhoff's PAM matrix

	<i>A</i>	<i>R</i>	<i>N</i>	<i>D</i>	<i>C</i>
<i>A</i>	9867	2	9	10	3
<i>R</i>	1	9913	1	0	1
<i>N</i>	4	1	9822	36	0
<i>D</i>	6	0	42	9859	0
<i>C</i>	1	1	0	0	9973

All entries  $\times 10^4$

(A)

	<b>A</b>	<b>C</b>	<b>G</b>	<b>T</b>
<b>A</b>	67	-96	-20	-117
<b>C</b>	-96	100	-79	-20
<b>G</b>	-20	-79	100	-96
<b>T</b>	-117	-20	-96	67

(B)

	<b>A</b>	<b>C</b>	<b>G</b>	<b>T</b>
<b>A</b>	91	-114	-31	-123
<b>C</b>	-114	100	-125	-31
<b>G</b>	-31	-125	100	-114
<b>T</b>	-123	-31	-114	91

(C)

	<b>A</b>	<b>C</b>	<b>G</b>	<b>T</b>
<b>A</b>	100	-123	-28	-109
<b>C</b>	-123	91	-140	-28
<b>G</b>	-28	-140	91	-123
<b>T</b>	-109	-28	-123	100

<sup>16</sup>

# Gap models

- Gap-extension
- Gap opening cost

# Local and global

## (A) local

PI3-kinase **DRHNSN**IMVKDDGQLFHI**DFG**  
 cAMP PK **DLKPEN**LLIDQQGYIQVT**DFG**

## (B) global

PI3-kinase    10                      20                      30                      40                      50  
 HQLGNLR--L**EE**CRI---MSSAKRPLWLNWENPDIMSELL**FQ**NNEIIFKNGDDL**RQ**DMLT  
 cAMP PK    10                      20                      30                      40                      50  
 GNAAA**AK**KGX**EQ**ESVKEFLAKAKEDFLKKWENPAQNTAHL**DQ**FERIKTLGTGSFGRV**ML**-

PI3-kinase    60                      70                      80                      90                      100                      110  
 LQIIRIME--NIWQN**QGL**DLRMLPYGCLSIGDCVGL**IE**VVRNSHT**IM**Q-IQCKGGL**K**GAL  
 cAMP PK    60                      70                      80                      90                      100  
 ---V**KH**METGNHYAMKI**LD**KQKVVK-----L**KQ**IEHTLNEKRILQAVNFPFLV**K**LEF

PI3-kinase    120                      130                      140                      150                      160  
 QFNSHT-LHQWLKDKN**KGE**IYDAA--IDLFTRSCAGYCVATFILGIG**DRHNSN**IMVKD-D  
 cAMP PK    110                      120                      130                      140                      150                      160  
 SFKD**NS**NLYMVMEYVPG**GEM**FSHLRR**IGR**FSEPHARFYAAQIVLTFEY**LH**SLDLIYR**DLK**

PI3-kinase    170                      180                      190                      200                      210                      220  
 GQLFHI**DFG**HFLDHKKKK**FGY**KRERVP-----FVLT**QD**FL---IVISKGAQECTKTREFE  
 cAMP PK    170                      180                      190                      200                      210                      220  
**PEN**LLIDQQGYI--QVT**DFG**FAK-RVKGRTWXL**CGT**PEYLAPEIILSKGYNKAVDWWALG

PI3-kinase    230                      240                      250                      260                      270  
 RF-Q**EM**C--YKAYLAIRQHANLFINLFS**MML**SGSGMP**ELQS**FDDIAYIRKT**LAL**DKTEQEA  
 cAMP PK    230                      240                      250                      260                      270  
 VLIY**EMA**AGYPPFFA-DQPIQIYEKIVSGKVR--FPSHFSSDLKD**LLRN**LLQVDLT**KR**--

# Global alignment Needleman-Wunch

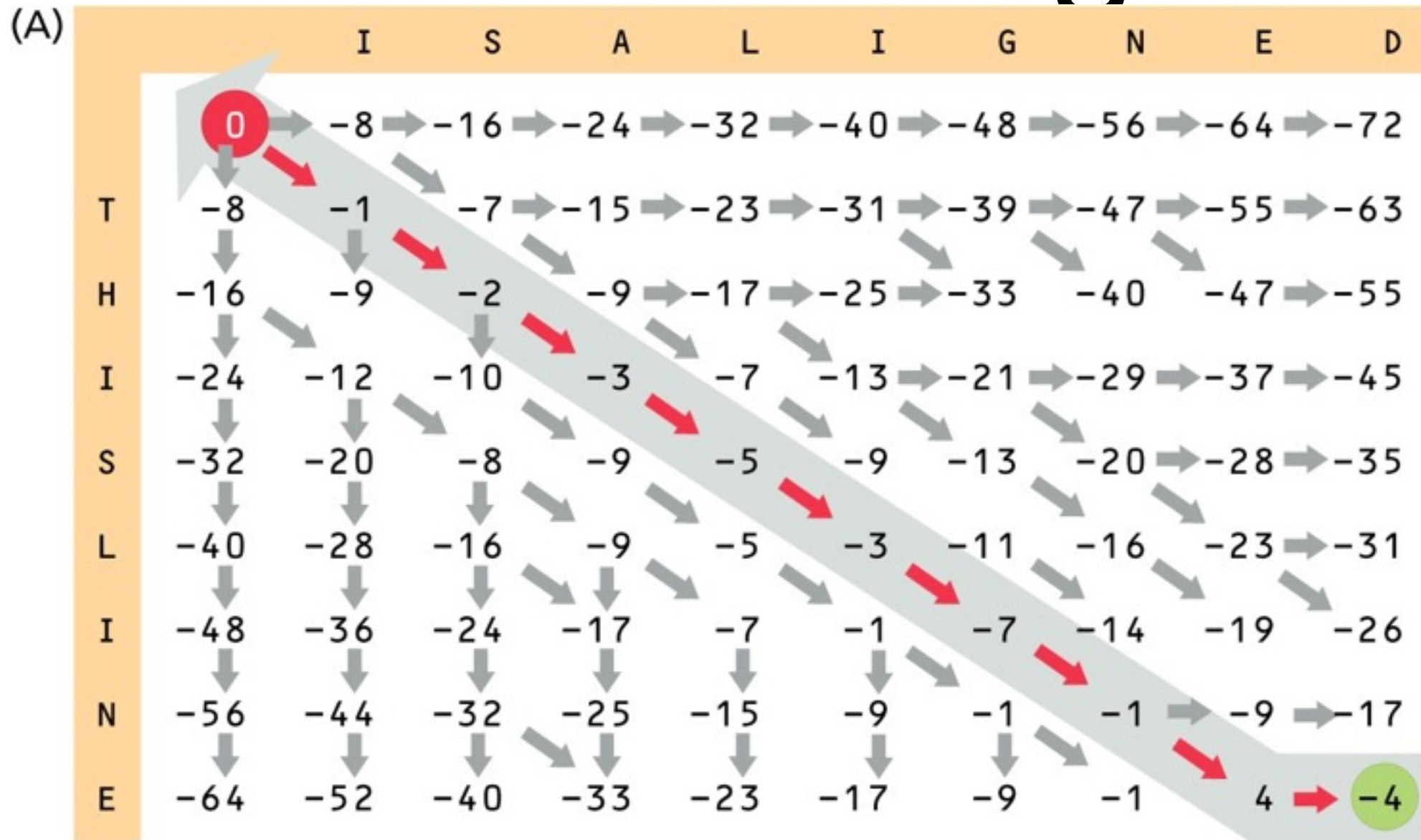
	GAP	M	N	A	L	S	D	R	T
GAP	0	-12	-16	-20	-24	-28	-32	-36	-40
M	-12	6 <sup>(6)</sup>	-6 <sup>(-2)</sup>	-10	-14	-18	-22	-26	-30
G	-16	-6 <sup>(-3)</sup>	6 <sup>(0)</sup>	-5	-10	-13	-17	-22	-26
S	-20	-10	-5	7	-5	-8	-13	-17	-21
D	-24	-14	-8	-5	3	-5	-4	-14	-17
R	-28	-18	-14	-9	-8	3	-6	2	-10
T	-32	-22	-18	-13	-11	-7	3	-7	5
T	-36	-26	-22	-17	-15	-10	-7	2	-4
E	-40	-30	-25	-21	-20	-15	-7	-8	2
T	-44	-34	-30	-24	-23	-19	-15	-8	-5

# Local alignment Smith Waterman

	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

# Different alignments

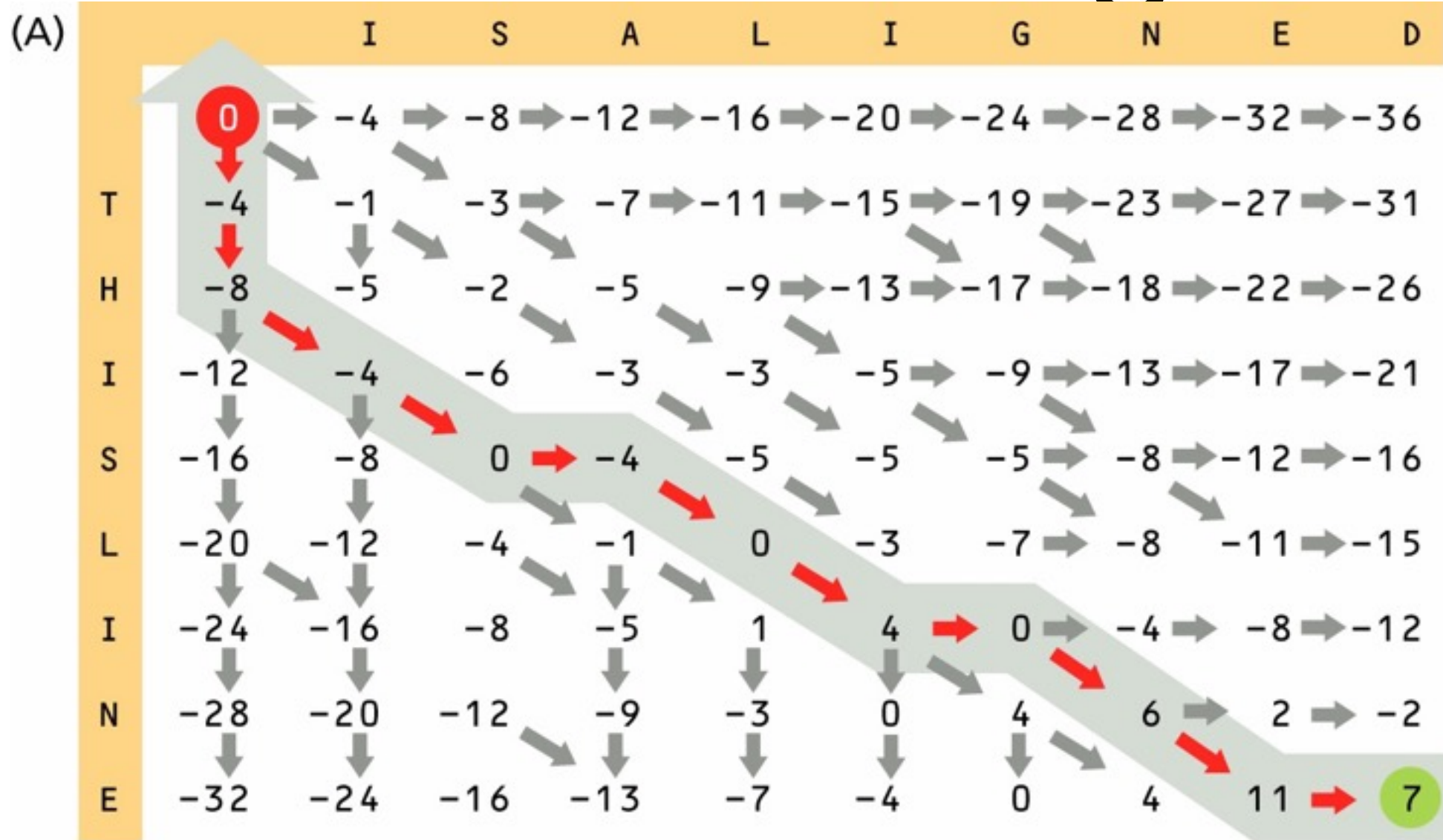
# Different alignments



(B) THISLINE-  
ISALIGNED



# Different alignments

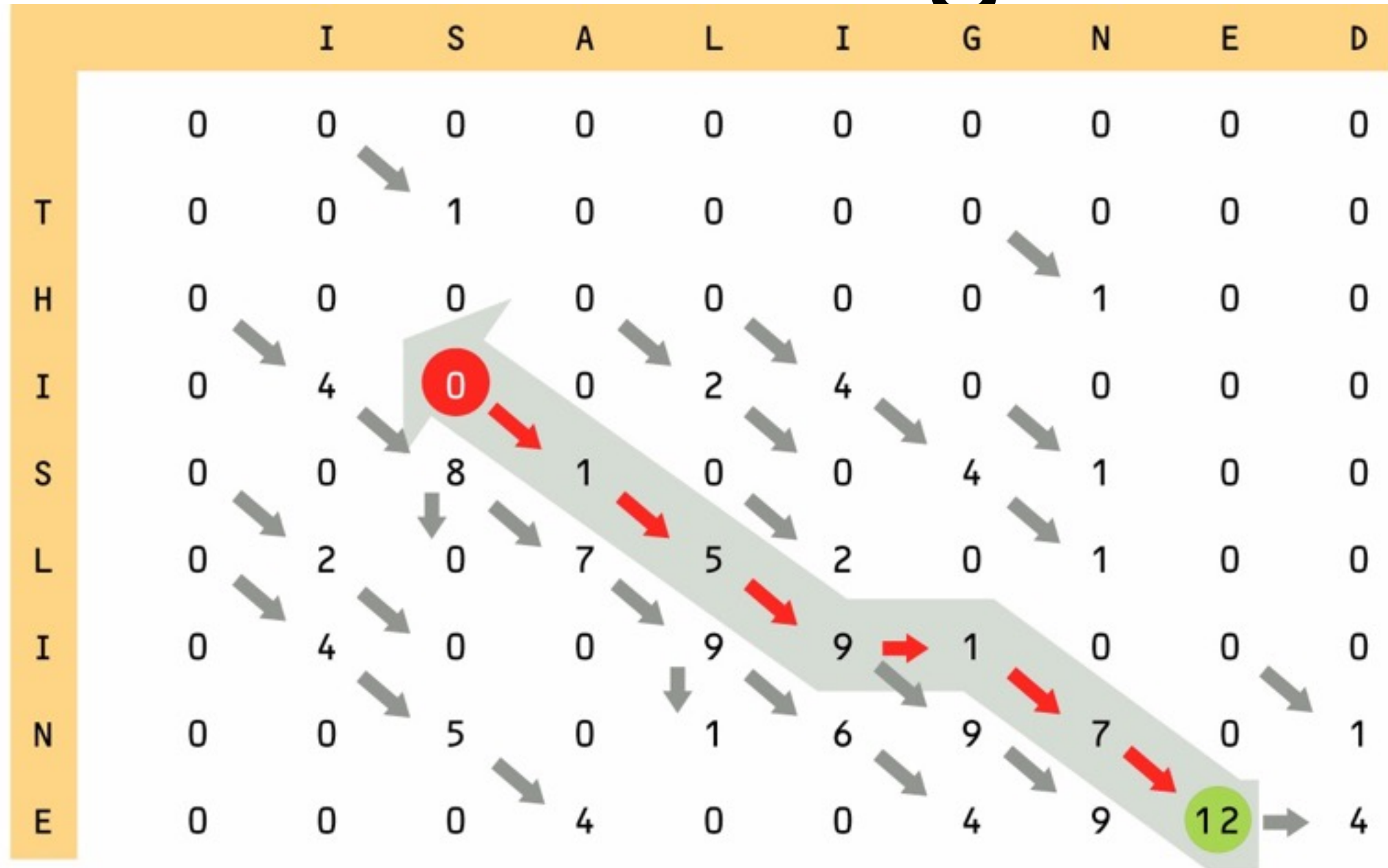


(B)

```
THIS-LI-NE-  
--ISALIGNED
```

# Different alignments

(A)



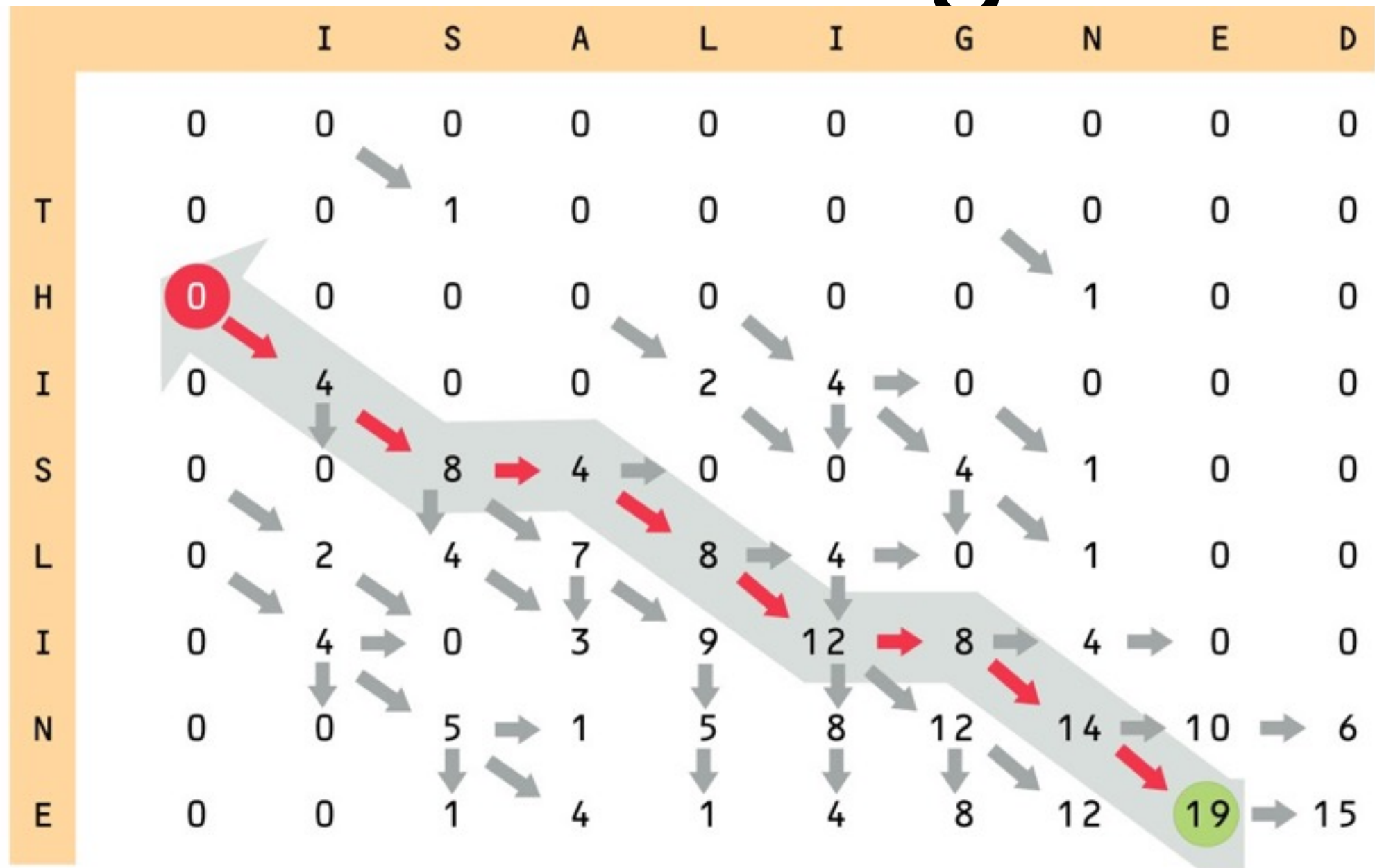
(B)

S L I - N E  
A L I G N E



# Different alignments

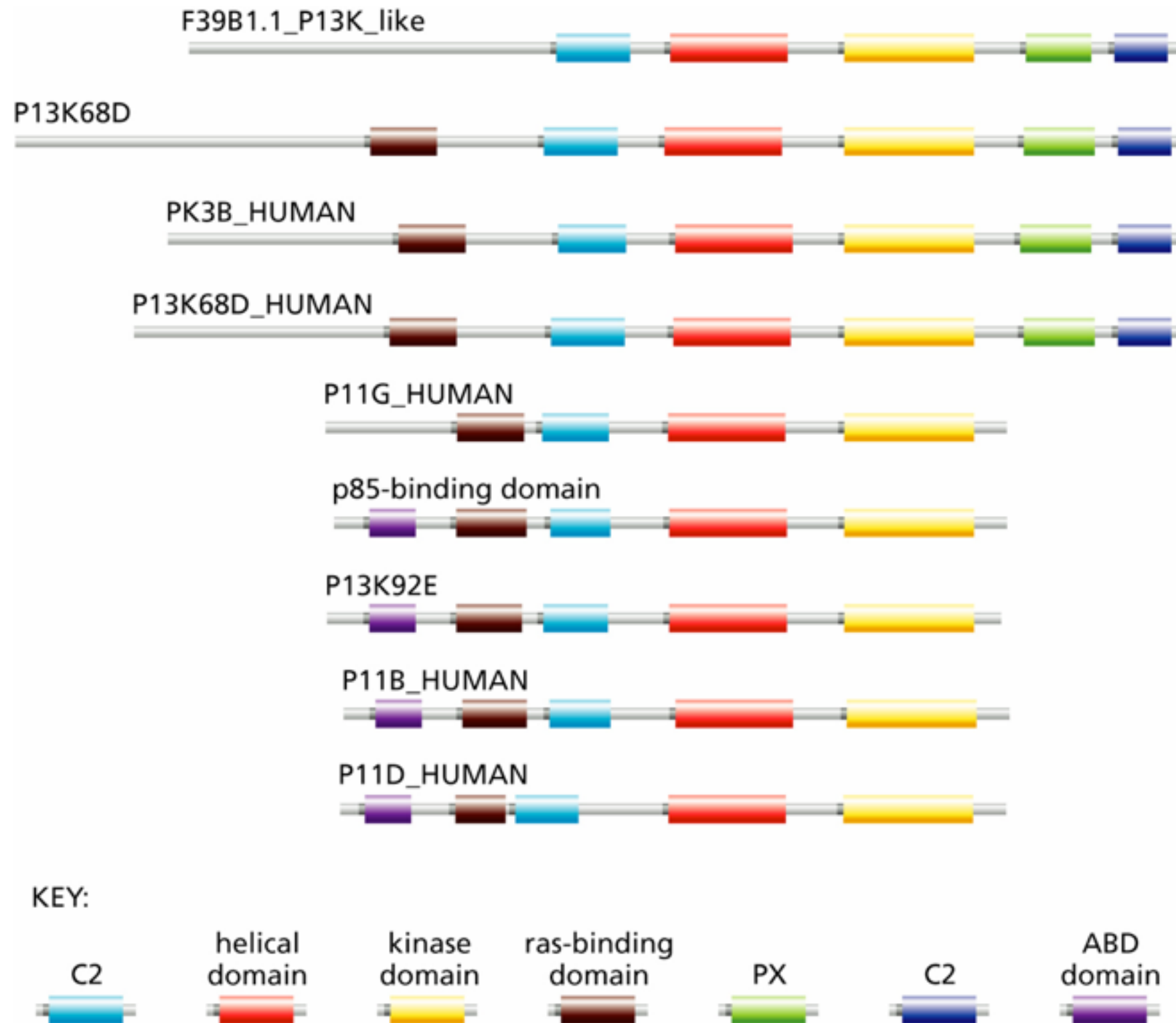
(A)



(B)

IS-LI-NE  
|||  
ISALIGNE

# Multidomain proteins



# Next lecture

- $O(nm)$  is too slow. How to speed up
- When is a “score” significant.