

# Sequence Alignment Algorithms

DEKM book

Notes from Dr. Bino John  
and Dr. Takis Benos

# To Do

- Global alignment
- Local alignment
- Gaps
  - Affine Gaps
  - Algorithm (blackboard)
- Statistical Significance
  - Notes (blackboard)
- Read up on database searches
  - BLAST
  - FASTA
  - CS tricks: suffix tree, ...
- PSSMs and Multiple Sequence Alignments

# Why compare sequences?

- Given a new sequence, infer its function based on similarity to another sequence
- Find important molecular regions – conserved across species

# Why compare sequences? Do more..

- Determine the evolutionary constraints at work
- Find mutations in a population or family of genes
- Find similar looking sequence in a database
- Find secondary/tertiary structure of a sequence of interest – molecular modeling using a template (homology modeling)

# Sequence alignment

- Are two sequences related?
  - Align sequences or parts of them
  - Decide if alignment is by chance or evolutionarily linked?
- Issues:
  - What sorts of alignments to consider?
  - How to score an alignment and hence rank?
  - Algorithm to find good alignments
  - Evaluate the significance of the alignment

# How do we use the matrices for sequence alignment?

AGGCTATCACCTGACCTCCAGGCCGATGCCC  
TAGCTATCACGACCGCGGTCGATTGCCCCGAC



-AGGCTATCACCTGACCTCCAGGCCGA--TGCCC---  
TAG-CTATCAC--GACCGC--GGTCGATTGCCCCGAC

## Definition

Given two strings  $x = x_1x_2\dots x_M$ ,  $y = y_1y_2\dots y_N$ ,

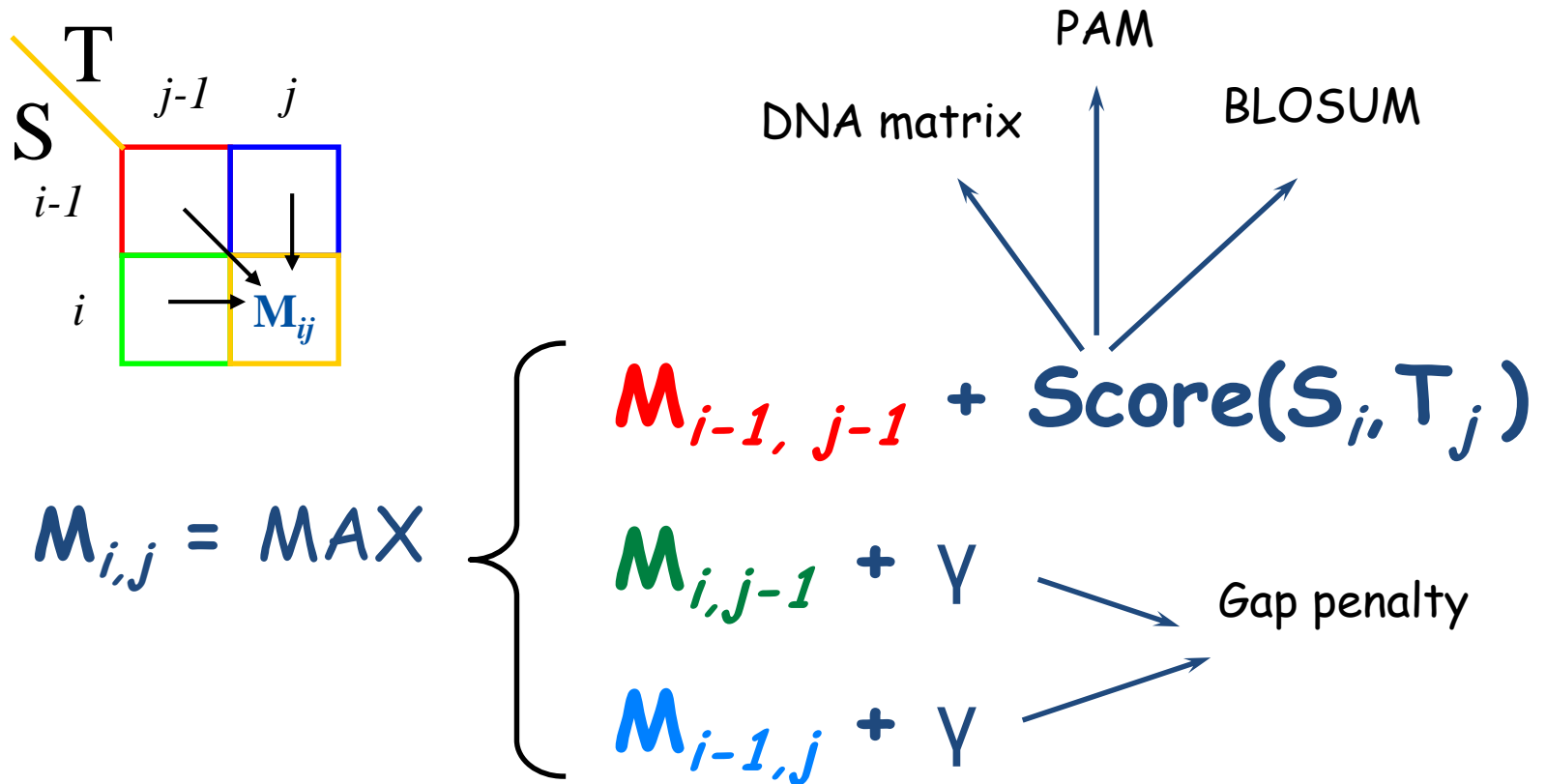
An **alignment** of two sequences  $x$  and  $y$  is an arrangement of  $x$  and  $y$  by position, where  $a$  and  $b$  can be padded with gap symbols to achieve the same length.

# Dynamic Programming

We apply **dynamic programming** when:

- There is only a polynomial number of subproblems
  - Align  $x_1 \dots x_i$  to  $y_1 \dots y_j$
- Original problem is one of the subproblems
  - Align  $x_1 \dots x_M$  to  $y_1 \dots y_N$
- Each subproblem is easily solved from smaller subproblems

# Global alignment



*Needleman & Wunsch, 1970*



# Dynamic programming for global alignment – simple case

We want the best alignment between two sequences  $x$  and  $y$

Consider two sequences:  $x_1 \dots x_M$ , and  $y_1 \dots y_M$

we have ONLY three choices to get the best score  $F(i,j)$

1.  $x_i$  aligns to  $y_j$

$x_1 \dots x_{i-1} \quad x_i$

$y_1 \dots y_{j-1} \quad y_j$

Align  $x[1 \dots i]$  with  $y[1 \dots j]$

2.  $x_i$  aligns to a gap

$x_1 \dots x_{i-1} \quad x_i$

$y_1 \dots y_j \quad -$

$x[1 \dots (i-1)]$  is already aligned with  $y[1 \dots (j)]$ , so align  $x[i]$  with a gap in  $y$

3.  $y_j$  aligns to a gap

$x_1 \dots x_i \quad -$

$y_1 \dots y_{j-1} \quad y_j$

$x[1 \dots i]$  is already aligned with  $y[1 \dots (j-1)]$ , so align a gap in  $x$  to  $y[j]$

1.  $x_i$  aligns to  $y_j$

$x_1 \dots x_{i-1} \quad x_i$

$y_1 \dots y_{j-1} \quad y_j$

$$F(i, j) = F(i-1, j-1) + s(i, j)$$

2.  $x_i$  aligns to a gap

$x_1 \dots x_{i-1} \quad x_i$

$y_1 \dots y_j \quad -$

$$F(i, j) = F(i-1, j) - \text{gap\_open\_penalty}(\text{go})$$

3.  $y_j$  aligns to a gap

$x_1 \dots x_i \quad -$

$y_1 \dots y_{j-1} \quad y_j$

$$F(i, j) = F(i, j-1) - \text{gap\_open\_penalty}(\text{go})$$

If we could make  $F(i, j-1)$ ,  $F(i-1, j)$ ,  $F(i-1, j-1)$  optimal, then we can make the next ones optimal as well

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

Where

$s(x_i, y_j)$  = Score for a match, if  $x_i = y_j$ ;  
score for a mismatch, if  $x_i \neq y_j$ ;

# The Needleman-Wunsch Algorithm

– pioneering application of DP to biological sequences

## 1. Initialization.

a.  $F(0, 0) = 0$

b.  $F(0, j) = -j \times \text{go}$

c.  $F(i, 0) = -i \times \text{go}$

$O(NM)$

## 2. Main Iteration. Filling-in partial alignments

For each  $i = 1 \dots M$

For each  $j = 1 \dots N$

$$F(i, j) = \max \begin{cases} F(i-1, j-1) + s(x_i, y_j) & [\text{case 1}] \\ F(i-1, j) - \text{go} & [\text{case 2}] \\ F(i, j-1) - \text{go} & [\text{case 3}] \end{cases}$$

$$\text{Ptr}(i, j) = \begin{cases} \text{DIAG}, & \text{if } [\text{case 1}] \\ \text{UP}, & \text{if } [\text{case 2}] \\ \text{LEFT}, & \text{if } [\text{case 3}] \end{cases}$$

## 3. Termination. $F(M, N)$ is the optimal score, and from $\text{Ptr}(M, N)$ can trace back optimal alignment

# Alignment: adding scores (cntd)

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

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

Score(match) = 1

Score(mismatch) = 0

Score(gap) = 0

# Alignment: adding scores

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

Diagram illustrating the first step of sequence alignment. The table shows scores for aligning the sequence "G A A T T C A G T T A" (columns) with the sequence "G G A T C G A" (rows). The scores are calculated based on matches (1) and mismatches (0). The alignment path is highlighted with arrows: red arrows point from (row 1, col 3) to (row 2, col 4), (row 2, col 4) to (row 3, col 5), (row 3, col 5) to (row 4, col 6), (row 4, col 6) to (row 5, col 7), (row 5, col 7) to (row 6, col 8), (row 6, col 8) to (row 7, col 9), and (row 7, col 9) to (row 8, col 10). The final score is 6.

		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

Diagram illustrating the second step of sequence alignment. The table shows scores for aligning the sequence "G A A T T C A G T T A" (columns) with the sequence "G G A T C G A" (rows). The scores are calculated based on matches (1) and mismatches (0). The alignment path is highlighted with arrows: red arrows point from (row 1, col 3) to (row 2, col 4), (row 2, col 4) to (row 3, col 5), (row 3, col 5) to (row 4, col 6), (row 4, col 6) to (row 5, col 7), (row 5, col 7) to (row 6, col 8), (row 6, col 8) to (row 7, col 9), and (row 7, col 9) to (row 8, col 10). The final score is 6.

# Alignment: adding scores (cntd)

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

(Seq #1)

Alignment:

(Seq #2)

	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	
A	0	1	1	1	1	1	1	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	
A											6

A  
|  
A

# Alignment: adding scores (cntd)

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

(Seq #1)

Alignment:

(Seq #2)

T A  
|  
- A



# Alignment: adding scores (cntd)

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

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

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

Alignment:

*6 matches, 1 mism., 4 gaps*

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

$s(x_i, y_j) = 1$  for match,  
 $-1$  for mismatch

$go = 2$

		T	C	G	C	A
		0 ← -2 ← -4 ← -6 ← -8 ← -10				
T		-2 ↑ ↖ 1 ← -1 ← -3 ← -5 ← -7				
C		-4 ↑ -1 ↖ 2 ← 0 ← -2 ← -4				
C		-6 ↑ -3 ↖ 0 ← 1 1 ← -1				
A		-8 ↑ -5 ↖ -2 ↖ -1 ↖ 0 ↖ 2				

	T	C	G	C	A	
T	0	-2	-4	-6	-8	-10
C	-2	1	-1	-3	-5	-7
C	-4	-1	2	0	-2	-4
C	-6	-3	0	1	1	-1
A	-8	-5	-2	-1	0	2

T(1) C(2) G(3) C(4) A(5)

	T	C	G	C	A
T (1)	-	D	-	-	-
C (2)	-	-	D	L	-
C (3)	-	-	-	-	D
A (4)	-	-	-	-	-

$$F(i, j) = \max \begin{cases} F(i-1, j-1) + s(x_i, y_j) \text{ [case 1]} \\ F(i-1, j) - \text{go} \text{ [case 2]} \\ F(i, j-1) - \text{go} \text{ [case 3]} \end{cases}$$

$$\text{Ptr}(i, j) = \begin{cases} \text{DIAG OR D, if [case 1]} \\ \text{UP OR U, if [case 2]} \\ \text{LEFT OR L, if [case 3]} \end{cases}$$

D at 1,1 => 1,1 is paired & decrease i and j by 1 => T,T & go to 0,0 - **END**

D at 2,2 => 2,2 is paired & decrease i and j by 1 => C,C & go to 1,1

L at 2,3 => -,3 is paired & decrease j by 1 => -,G & go to 2,2

D at 3,4 => 3,4 is paired and decrease i and j by 1 => CC & go to 2,3

D at 4,5 => 4,5 are paired and decrease i and j by 1 => AA & go to 3,4

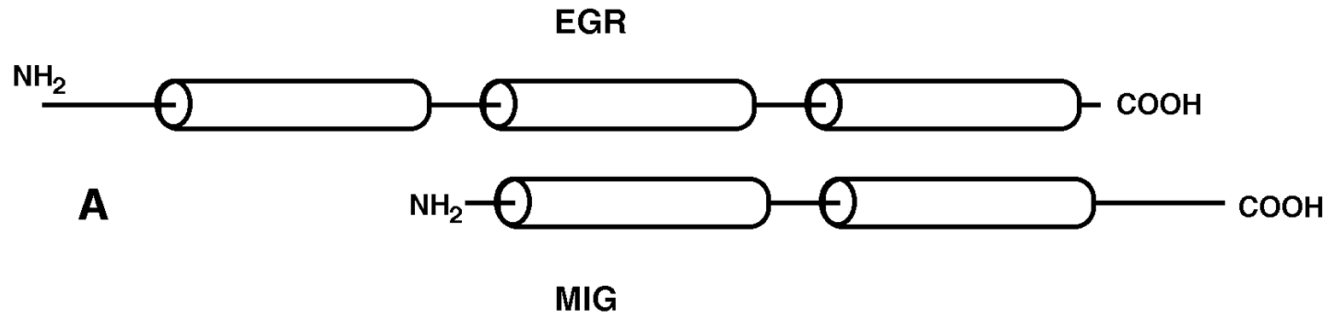
T	C	-	C	A
T	C	G	C	A

# Local alignment

Given two sequences,  $S$  and  $T$ , find two subsequences,  $s$  and  $t$ , whose alignment has the highest “score” amongst all subsequence pairs.

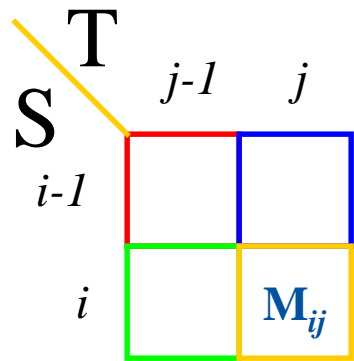
Question: Why do we need local alignment, if we have the global one?

# Local alignment: an example



EGR4_HUMAN	KA	[FACPVE	SCVRSF	ARSDEL	NRHLRIH]	TGHKP	[FQCRIC	LRNF	<b>RS</b>	<b>DH</b>	<b>LT</b>	<b>SH</b>	VRTH]	TGEKP	[FACDV--	CGRRF	<b>ARS</b>	<b>DE</b>	<b>KK</b>	<b>RH</b>	SKVH]		
EGR4_RAT	KA	[FACPVE	SCVRTF	ARSDEL	NRHLRIH]	TGHKP	[FQCRIC	LRNF	<b>RS</b>	<b>DH</b>	<b>LT</b>	<b>TH</b>	VRTH]	TGEKP	[FACDV--	CGRRF	<b>ARS</b>	<b>DE</b>	<b>KK</b>	<b>RH</b>	SKVH]		
EGR3_HUMAN	RP	[HACPAE	GCDRRF	SRSE	DELTRHLRIH]	TGHKP	[FQCRIC	MRSF	<b>RS</b>	<b>DH</b>	<b>LT</b>	<b>TH</b>	IRTH]	TGEKP	[FACEF--	CGRKF	<b>ARS</b>	<b>DE</b>	<b>RK</b>	<b>RH</b>	AKIH]		
EGR3_RAT	RP	[HACPAE	GCDRRF	SRSE	DELTRHLRIH]	TGHKP	[FQCRIC	MRSF	<b>RS</b>	<b>DH</b>	<b>LT</b>	<b>TH</b>	IRTH]	TGEKP	[FACEF--	CGRKF	<b>ARS</b>	<b>DE</b>	<b>RK</b>	<b>RH</b>	AKIH]		
EGR1_HUMAN	RP	[YACPVE	SCDRRF	SRSE	DELTRHIRIH]	TGQKP	[FQCRIC	MRNF	<b>RS</b>	<b>DH</b>	<b>LT</b>	<b>TH</b>	IRTH]	TGEKP	[FACDI--	CGRKF	<b>ARS</b>	<b>DE</b>	<b>RK</b>	<b>RH</b>	TKIH]		
EGR1_MOUSE	RP	[YACPVE	SCDRRF	SRSE	DELTRHIRIH]	TGQKP	[FQCRIC	MRNF	<b>RS</b>	<b>DH</b>	<b>LT</b>	<b>TH</b>	IRTH]	TGEKP	[FACDI--	CGRKF	<b>ARS</b>	<b>DE</b>	<b>RK</b>	<b>RH</b>	TKIH]		
EGR1_RAT	RP	[YACPVE	SCDRRF	SRSE	DELTRHIRIH]	TGQKP	[FQCRIC	MRNF	<b>RS</b>	<b>DH</b>	<b>LT</b>	<b>TH</b>	IRTH]	TGEKP	[FACDI--	CGRKF	<b>ARS</b>	<b>DE</b>	<b>RK</b>	<b>RH</b>	TKIH]		
EGR1_BRARE	RP	[YACPVE	TCDRRF	SRSE	DELTRHIRIH]	TGQKP	[FQCRIC	MRNF	<b>RS</b>	<b>DH</b>	<b>LT</b>	<b>TH</b>	IRTH]	TGEKP	[FACEI--	CGRKF	<b>ARS</b>	<b>DE</b>	<b>RK</b>	<b>RH</b>	TKIH]		
EGR2_RAT	RP	[YPCPAE	GCDRRF	SRSE	DELTRHIRIH]	TGHKP	[FQCRIC	MRNF	<b>RS</b>	<b>DH</b>	<b>LT</b>	<b>TH</b>	IRTH]	TGEKP	[FACDY--	CGRKF	<b>ARS</b>	<b>DE</b>	<b>RK</b>	<b>RH</b>	TKIH]		
EGR2_XENLA	RP	[YPCPAE	GCDRRF	SRSE	DELTRHIRIH]	TGHKP	[FQCRIC	MRNF	<b>RS</b>	<b>DH</b>	<b>LT</b>	<b>TH</b>	IRTH]	TGEKP	[FACDY--	CGRKF	<b>ARS</b>	<b>DE</b>	<b>RK</b>	<b>RH</b>	TKIH]		
EGR2_MOUSE	RP	[YPCPAE	GCDRRF	SRSE	DELTRHIRIH]	TGHKP	[FQCRIC	MRNF	<b>RS</b>	<b>DH</b>	<b>LT</b>	<b>TH</b>	IRTH]	TGEKP	[FACDY--	CGRKF	<b>ARS</b>	<b>DE</b>	<b>RK</b>	<b>RH</b>	TKIH]		
EGR2_HUMAN	RP	[YPCPAE	GCDRRF	SRSE	DELTRHIRIH]	TGHKP	[FQCRIC	MRNF	<b>RS</b>	<b>DH</b>	<b>LT</b>	<b>TH</b>	IRTH]	TGEKP	[FACDY--	CGRKF	<b>ARS</b>	<b>DE</b>	<b>RK</b>	<b>RH</b>	TKIH]		
EGR2_BRARE	RP	[YPCPAE	GCDRRF	SRSE	DELTRHIRIH]	TGHKP	[FQCRIC	MRNF	<b>RS</b>	<b>DH</b>	<b>LT</b>	<b>TH</b>	IRTH]	TGEKP	[FACDF--	CGRKF	<b>ARS</b>	<b>DE</b>	<b>RK</b>	<b>RH</b>	TKIH]		
MIG1_KLULA	--	[-----	-----	-----	-----]	---RP	[YVCPIC	QRGF	<b>HR</b>	<b>LE</b>	<b>HQ</b>	<b>TR</b>	IRTH]	TGERP	[HACDFP	GC	SK	<b>RS</b>	<b>DE</b>	<b>LT</b>	<b>RH</b>	RRIH]	
MIG1_KLUMA	--	[-----	-----	-----	-----]	---RP	[YMCPI	CHRG	<b>FH</b>	<b>RL</b>	<b>EH</b>	<b>Q</b>	<b>TR</b>	IRTH]	TGERP	[HACDFP	GC	AK	<b>RS</b>	<b>DE</b>	<b>LT</b>	<b>RH</b>	RRIH]
MIG1_YEAST	--	[-----	-----	-----	-----]	---RP	[HACPIC	HRA	<b>FH</b>	<b>RL</b>	<b>EH</b>	<b>Q</b>	<b>TR</b>	HMRIH]	TGEKP	[HACDFP	GC	VK	<b>RS</b>	<b>DE</b>	<b>LT</b>	<b>RH</b>	RRIH]
MIG2_YEAST	--	[-----	-----	-----	-----]	---RP	[FRCDT	CHRG	<b>FH</b>	<b>RL</b>	<b>EH</b>	<b>KK</b>	<b>RH</b>	LRTH]	TGEKP	[HHCAFP	GC	GK	<b>RS</b>	<b>DE</b>	<b>LT</b>	<b>RH</b>	MRTH]
		[			]	:	*		*	*	*	*	*	:	*	.	*	:	*	*	*	*	

# Local alignment (cntd)



$$M_{i,j} = \text{MAX}$$

$$\left\{ \begin{array}{l} 0 \\ M_{i-1,j-1} + \text{Score}(S_i, T_j) \\ M_{i,j-1} + Y \\ M_{i-1,j} + Y \end{array} \right.$$

Arrows from the three non-zero options point to a common point, with labels: DNA matrix (pointing to the first option), PAM (pointing to the second option), and BLOSUM (pointing to the third option). An arrow from the 'Y' in the third and fourth options points to the label 'Gap penalty'.

*Smith & Waterman, 1981*

Similarity Scoring Expected value:  
negative for random alignments  
positive for highly similar sequences

# The Smith-Waterman Algorithm

## 1. Initialization

$$F(0,0) = F(0,j) = F(i,0) = 0$$

## 2. Iteration

for  $i=1,\dots,M$

for  $j=1,\dots,N$

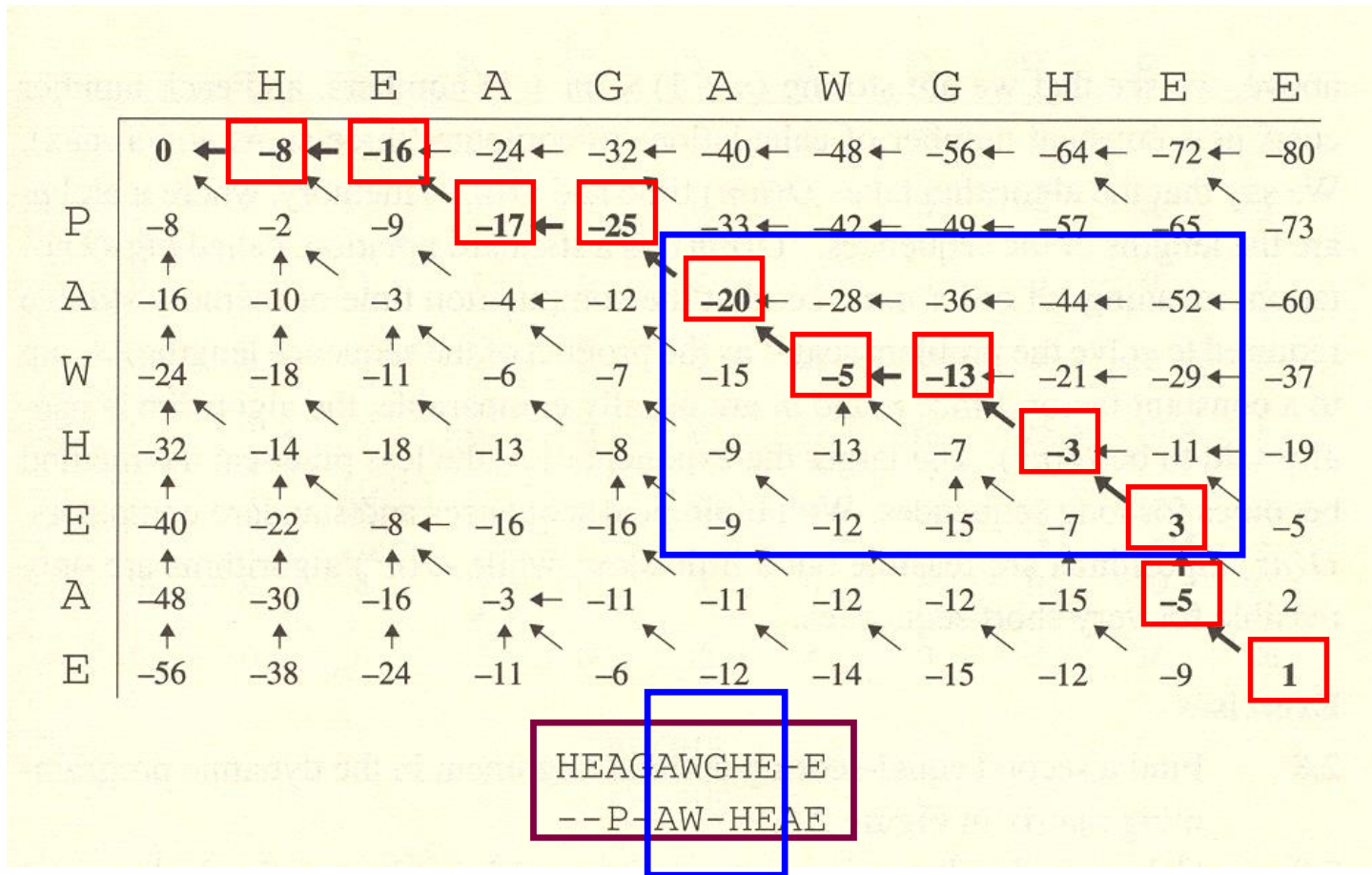
- calculate optimal  $F(i,j)$

- store  $\text{Ptr}(i,j)$

## 3. Termination

- Find the end of the best alignment with  $F_{\text{OPT}} = \max_{\{i,j\}} F(i,j)$  and trace back OR
- Find all alignments with  $F(i,j) > \text{threshold}$  and trace back

# Local vs. global alignment





# Local vs. global alignment (cntd)

		H	E	A	G	A	W	G	H	E	E
P	0	0	0	0	0	0	0	0	0	0	0
A	0	0	0	0	0	0	0	0	0	0	0
W	0	0	0	0	2	0	0	0	0	0	0
H	0	10	2	0	0	0	12	18	4	0	0
E	0	2	16	8	0	0	4	10	18	0	20
A	0	0	8	21	13	5	0	4	10	20	27
E	0	0	6	13	18	12	4	0	4	16	26

AWGHE  
 AW-HE

# Local alignment (cntd)

## Characteristics of local alignments:

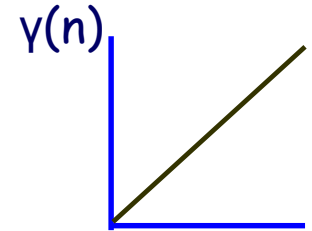
- The alignment can start/end at any point in the matrix.
- No negative scores in the alignment.
- The mean value of the scoring matrix (e.g. PAM, BLOSUM) should be negative, but there should be positive scores in the scoring matrix.

# Scoring the gaps more accurately

- **A naive model**

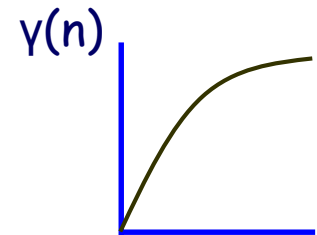
**Gap penalty is linear to the gap length**

Nature “prefers” to place gaps where other gaps exist



- **Convex gap penalty function**

$$\gamma(n+1) - \gamma(n) \leq \gamma(n) - \gamma(n-1)$$



**Time**  $O(N^2M)$       **Space**  $O(NM)$   
(assume  $N > M$ )

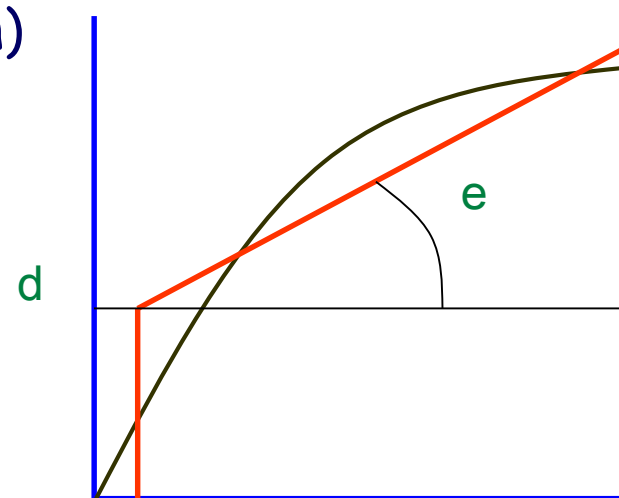
# Scoring gaps: affine gaps

- Affine gaps: a compromise between linear and convex gap penalties

$$\gamma(n) = -d - e * (n-1)$$

$d$ : gap initiation penalty

$e$ : gap extension penalty



# Convex gap dynamic programming

Initialization: same as before

Iteration:

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

Termination: same

Running Time:  $O(N^2M)$  (assume  $N > M$ )

Space:  $O(NM)$

# Smith, waterman, and Beyer Algorithm for affine gap

## 1. Initialization.

a.  $F(0, 0) = 0$

b.  $F(0, j) = -j \times g_o$

c.  $F(i, 0) = -i \times g_o$

$O(N^3)$

## 2. Main Iteration. Filling-in partial alignments

For each  $i = 1 \dots M$

For each  $j = 1 \dots N$

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

# Fast implementation of affine gap penalty

We need three matrices for tracking scores

Matrix-1:  $a[i,j]$  = to store maximum score of an alignment that

ends in  $x[i]$  matched to  $y[j]$

... $x_i$   
... $y_j$

Matrix-2:  $b[i,j]$  = to store maximum score of an alignment that

ends in gap matched to  $y[j]$

...-  
... $y_j$

Matrix-3:  $c[i,j]$  = to store maximum score of an alignment that

ends in gap matched to  $x[i]$

... $x_i$   
... -

## Implementation – Cont'd

$$a[i, j] = \max \begin{cases} a[i-1, j-1] \\ b[i-1, j-1] + s(i, j) \\ c[i-1, j-1] \end{cases} \quad b[i, j] = \max \begin{cases} a[i, j-1] + go \\ b[i, j-1] + ge \\ c[i, j-1] + go \end{cases}$$

$$c[i, j] = \max \begin{cases} a[i-1, j] + go \\ b[i-1, j] + go \\ c[i-1, j] + ge \end{cases}$$

**Pointer-matrices:** Three matrices to figure out which state within each score matrix maximization was used to obtain the optimal alignment of position  $i, j$ . Of course you need to know which matrix yielded the best score in the end  $(m, n)$  as well



## The twist with affine gap penalty

$$\begin{array}{ccc} \text{score} \left( \begin{array}{c} WA--- \\ -AGC \end{array} \right) & \text{vs} & \text{score} \left( \begin{array}{c} WA- \\ -AG \end{array} \right) + \text{score} \left( \begin{array}{c} - \\ C \end{array} \right) \\ g(1) + s(A, A) + g(2) & \neq & g(1) + s(A, A) + g(1) + g(1) \end{array}$$

# Fast implementation of affine gap penalty

We need three matrices for tracking scores

Matrix-1:  $a[i,j]$  = to store maximum score of an alignment that ends in  $x[i]$  matched to  $y[j]$

... $x_i$
... $y_j$

Matrix-2:  $b[i,j]$  = to store maximum score of an alignment that ends in gap matched to  $y[j]$

...-
... $y_j$

Matrix-3:  $c[i,j]$  = to store maximum score of an alignment that ends in gap matched to  $x[i]$

... $x_i$
... -

## Implementation – Cont'd

$$a[i, j] = \max \begin{cases} a[i-1, j-1] \\ b[i-1, j-1] + s(i, j) \\ c[i-1, j-1] \end{cases} \quad b[i, j] = \max \begin{cases} a[i, j-1] + go \\ b[i, j-1] + ge \\ c[i, j-1] + go \end{cases}$$

$$c[i, j] = \max \begin{cases} a[i-1, j] + go \\ b[i-1, j] + go \\ c[i-1, j] + ge \end{cases}$$

Pointer-matrices: Three matrices to figure out which state within each score matrix maximization was used to obtain the optimal alignment of position  $i, j$ . Of course you need to know which matrix yielded the best score in the end  $(m, n)$  as well

<http://www.ebi.ac.uk/Tools/emboss/align/index.html>

### EMBOSS Pairwise Alignment Algorithms

This tool is used to compare 2 sequences. When you want an alignment that covers the whole length of both sequences, use [needle](#). When you are trying to find the best region of similarity between two sequences, use [water](#).

For checking results of your code

Method: **EMBOSS::water (local)**

Gap Open: **10.0**

Gap Extend: **0.5**

Molecule: **DNA**

Matrix: **DNAfull**

Sequence 1: paste Sequence in any format OR upload a file:

Seq. 1 Upload a file:

Choose the alignment method : ☐ local (default) ☐ global ☒ global without end-gap penalty

Number of reported sub-alignments : **3**

Scoring matrix : **BLOSUM62**

Opening gap penalty : **0** (default -14)

Extending gap penalty : **0** (default -4)

First sequence title (optional):

Input sequence format: **Plain Text**

1st Query sequence: **GGATCGA**

or ID or AC or GI

Use a checkbox for valid

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

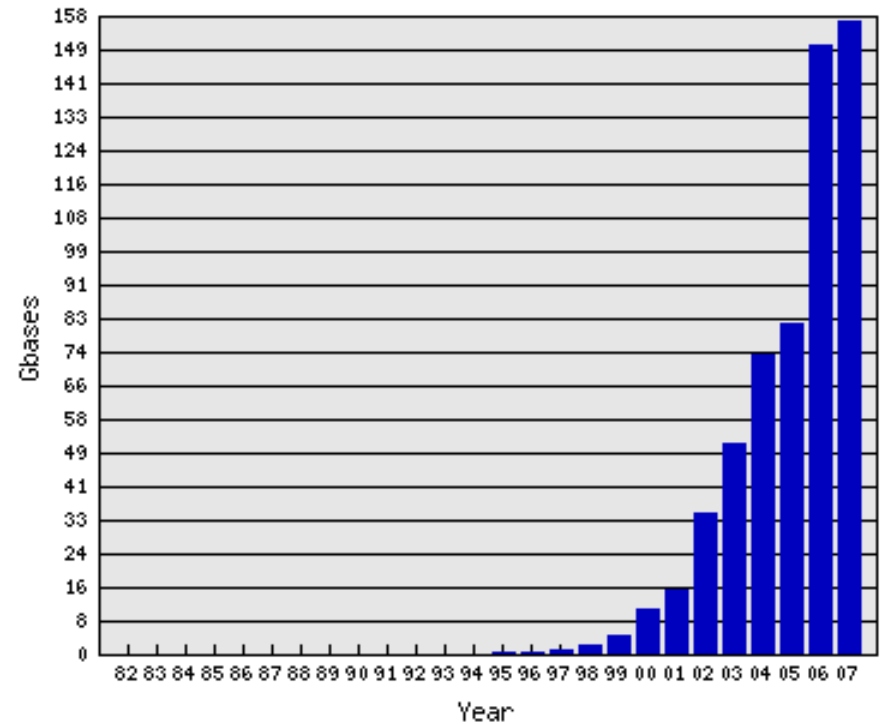
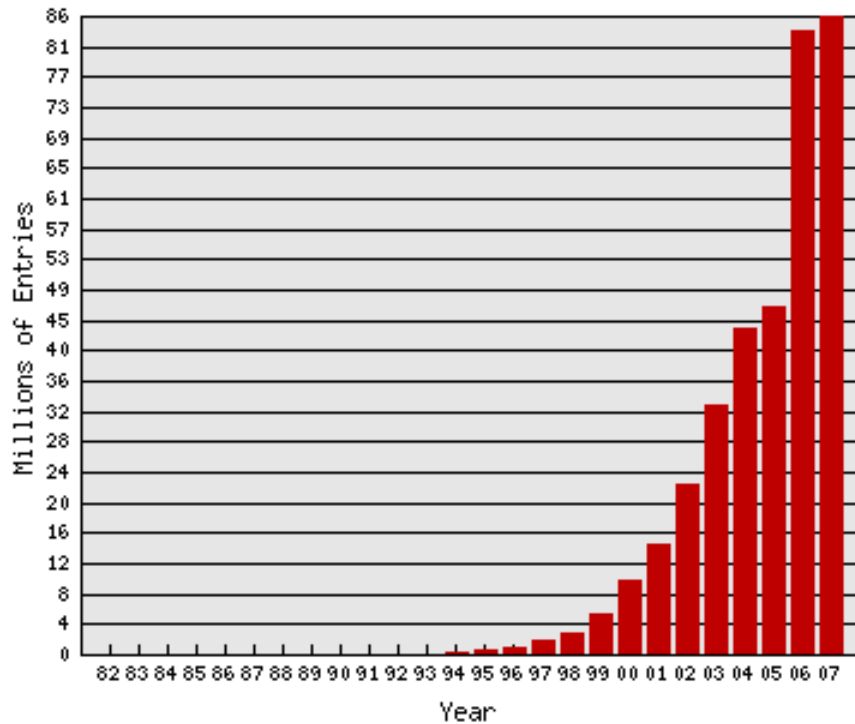
## Alignment Programs

- The Fasta program **Align** –Global.
- **EMBOSS Needle/ Stretcher**–Global
- Fasta program **LALIGN** – Local
- **BLAST 2 Sequences @ NCBI** – Local
  - <http://www.ncbi.nlm.nih.gov/blast/bl2seq/bl2.html>
- **AVID** –Genome Scale alignment (LONG Seqs)

# Database searches

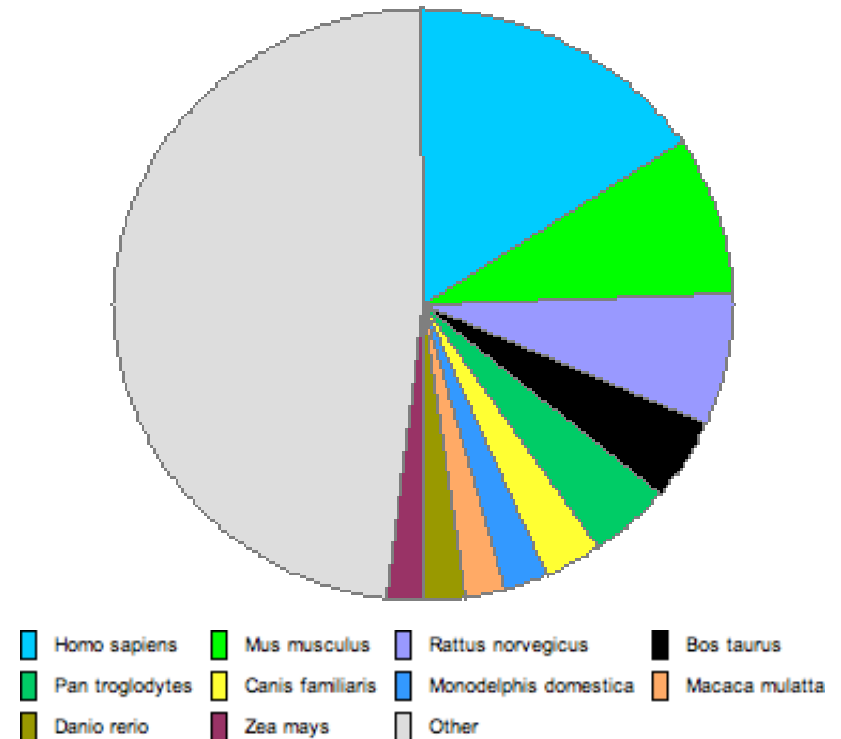
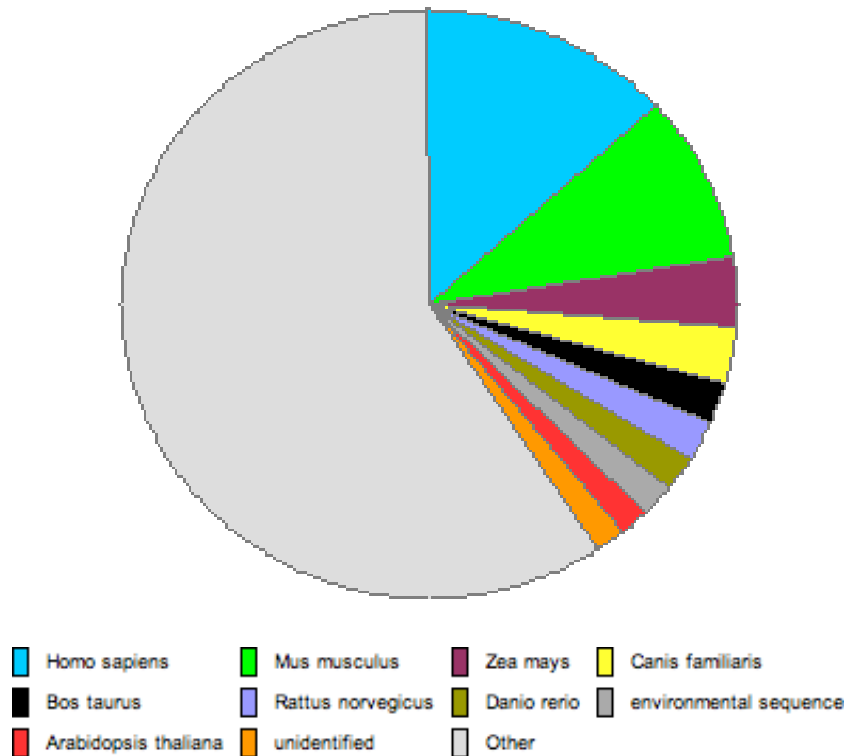
# DNA and protein databases

- EMBL/GenBank/DDBJ database of nucleic acids



# DNA and protein databases

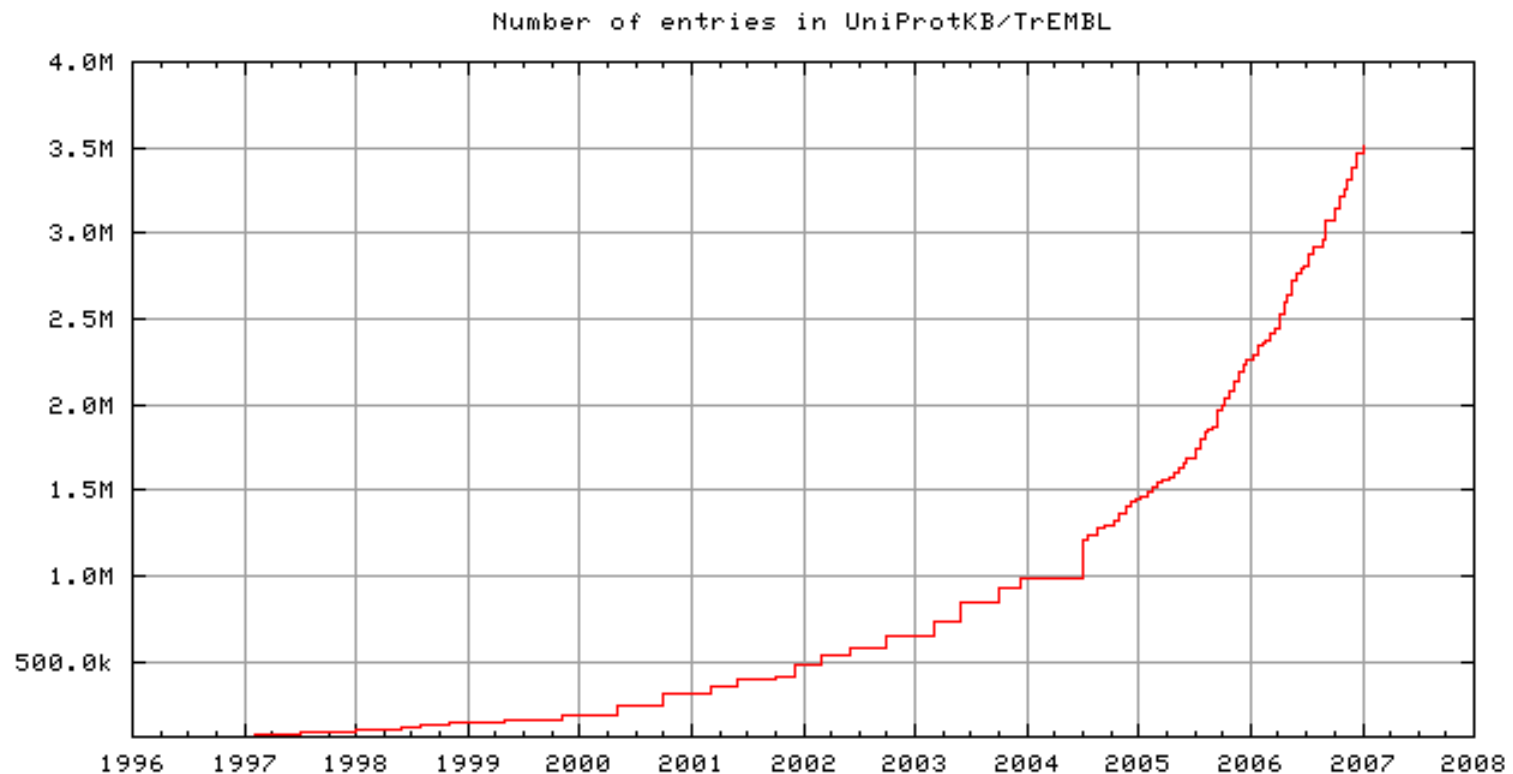
- EMBL/GenBank/DDBJ database of nucleic acids (cntd)





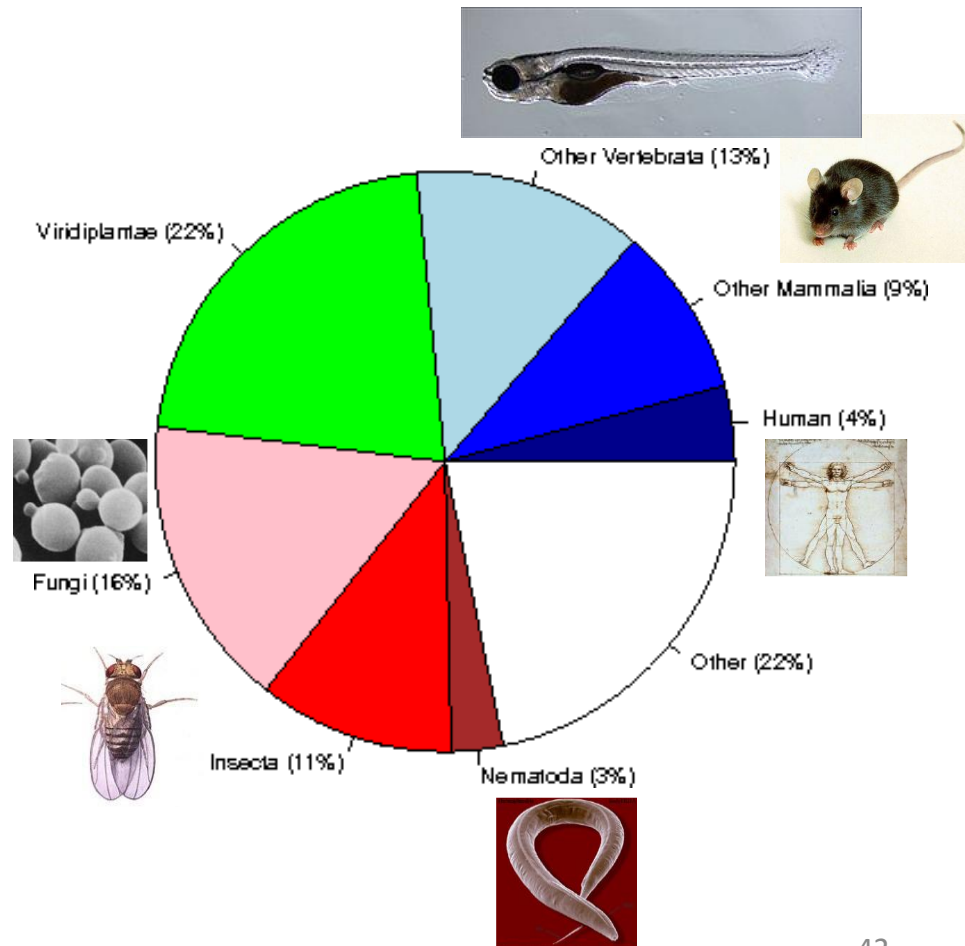
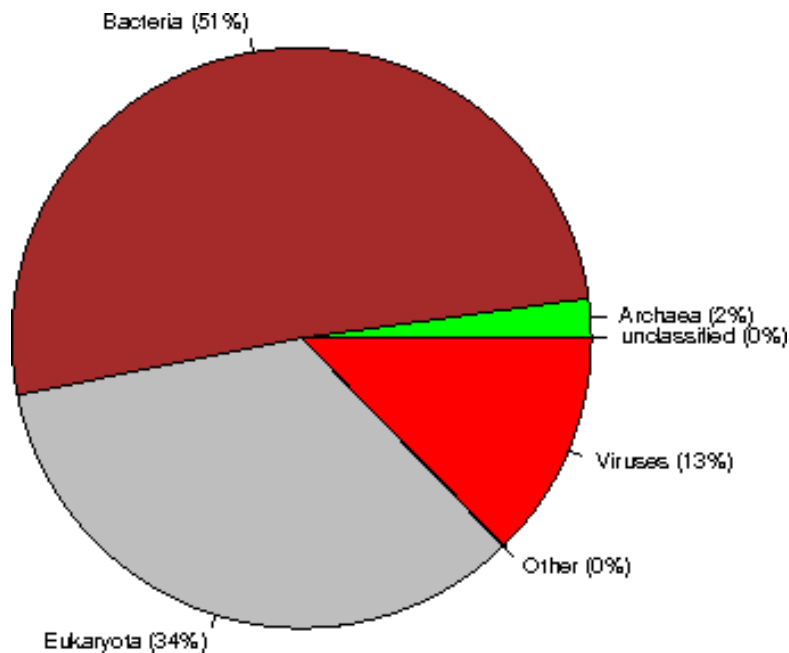
# DNA and protein databases

- SWISS-PROT & TrEMBL database of proteins



# DNA and protein databases

- SWISS-PROT & TrEMBL database of proteins



# Database searches

- Database searching consists of many pairwise alignments combined in one search.
- It helps determining the *function* and the evolutionary relationships
- Heuristic algorithms are used instead of DP. *Why?*
  - Size of SWISS-PROT + TrEMBL (Rel. 9.5):  
3.9M entries or 1,276M residues.
  - Exact algorithms are  $O(NM)$  fast.
- Heuristic methods can look at a small fraction of the searching space that will include all (or most) of the high scoring pairs.

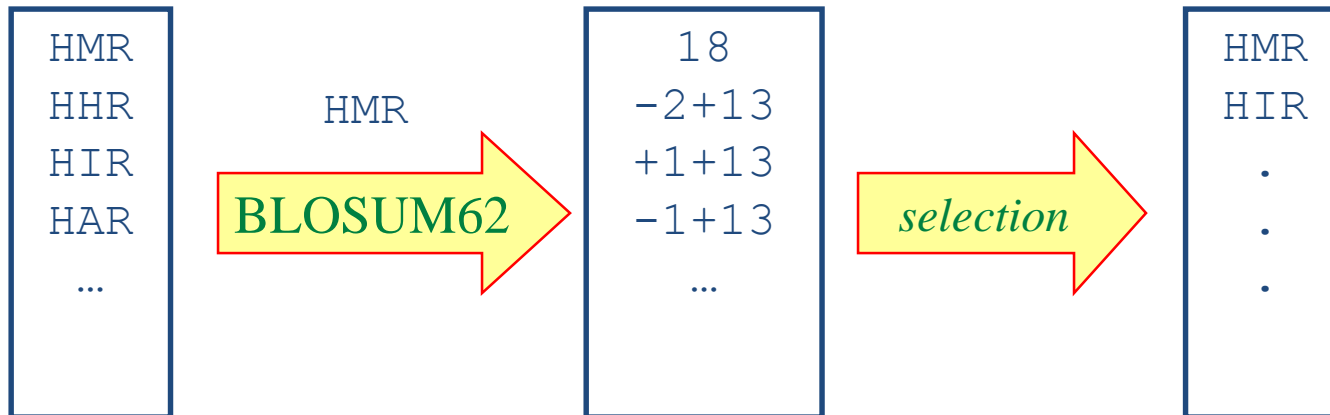
# BLAST algorithm

- Basic Local Alignment Search Tool - The method:
  - For each “word” (of fixed-length) in the query sequence, make a list of all neighbouring “words” that score above some threshold.
  - Scan the database for these words.
  - Perform (ungapped) “hit extension” until score < threshold.
  - Stop at maximum scoring extension.

# BLAST algorithm (cntd)

- An example:

Query: CPICHRAFHRLEHQTRHMR**I**HTGEKPHAC



# BLAST algorithm (cntd)

- An example:

Query: CPICHRAFHRLHQTRHMR IHTGEKPHAC

H+R

Sbjct: CPLCDKAFHRLHQTRHIR THTGEKPHAC

# BLAST algorithm (cntd)

- An example:

Query: CPICHRAFHRLHQTRHMR IHTGEKPHAC  
CP+C +AFHRLHQTR H+R HTGEKPHAC

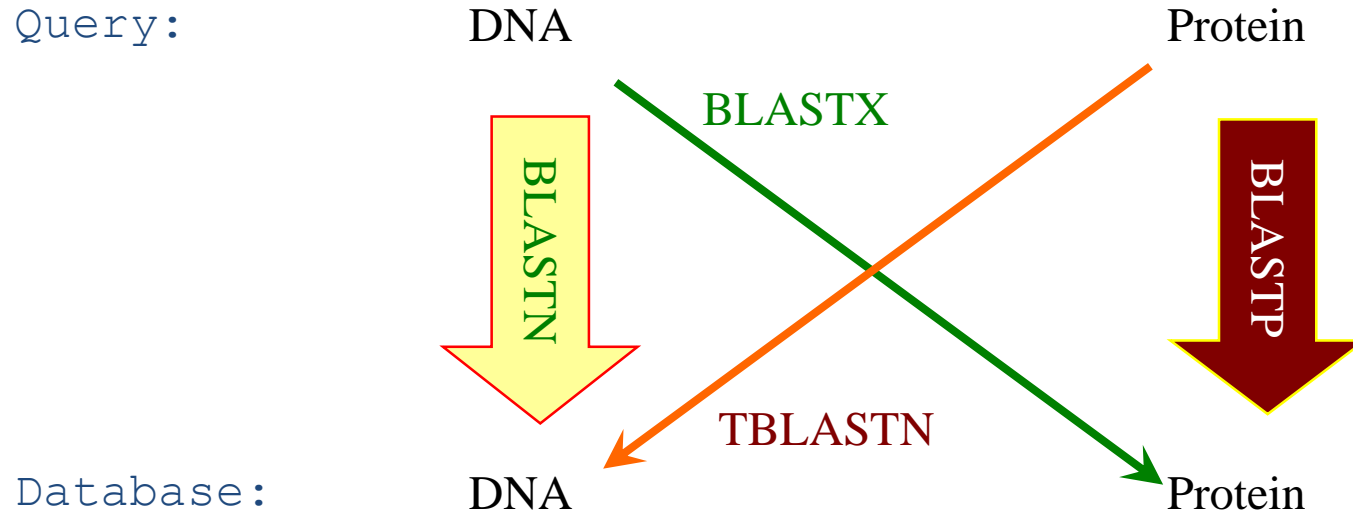
Sbjct: CPLCDKAFHRLHQTRHIRTHTGEKPHAC

# BLAST algorithm (cntd)

- **The idea:** a high scoring match alignment is very likely to contain a short stretch of very high scoring matches.
- **Word length:** 3 (proteins) and 11 (DNA).
- **HSP:** multiple HSPs can be reported for each database entry.
- **Gapped alignments:** more recent BLAST versions perform gapped alignments.



# BLAST flavours

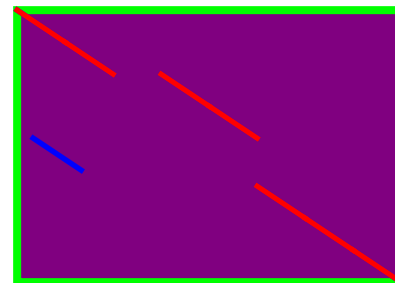
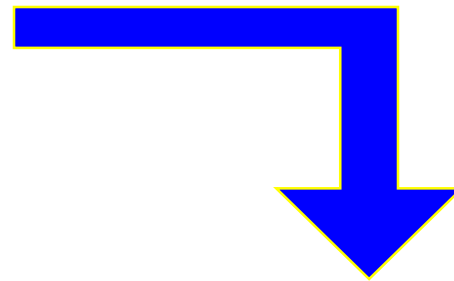
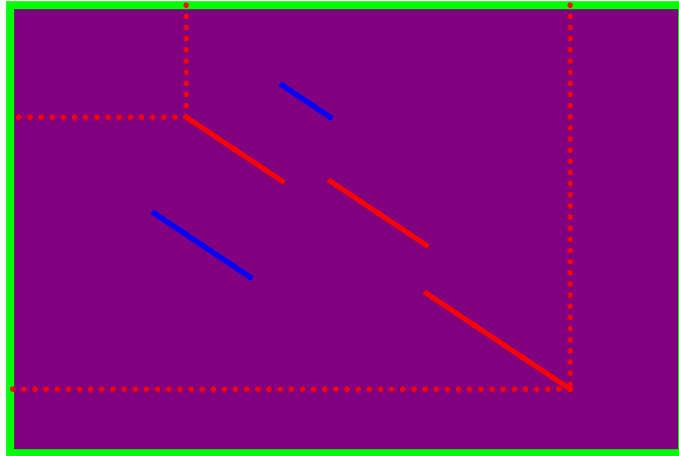


**TBLASTX:** DNA Query to DNA Database *via* translation

# FASTA algorithm

- The method:
  - For each pair of sequences (query, subject), identify all identical “word” matches of (fixed) length.
  - Look for diagonals with many mutually supporting “word” matches.
  - The best diagonals are used to extend the word matches to find the maximal scoring (ungapped) regions.
  - Join ungapped regions, using gap costs.
  - Align the two (sub)regions using full dynamic programming techniques.

# FASTA algorithm (cntd)



# FASTA algorithm (cntd)

- **The idea:** a high scoring match alignment is very likely to contain a short stretch of identities.
- **Word length:** 2 (proteins) and 4-6 (DNA).
- **HSP:** usually one (extended) gapped alignment is presented.

# FASTA flavours

Query:

DNA

Protein

FASTX3

FASTA3

FASTA3

TFASTA3

Database:

DNA

Protein

