



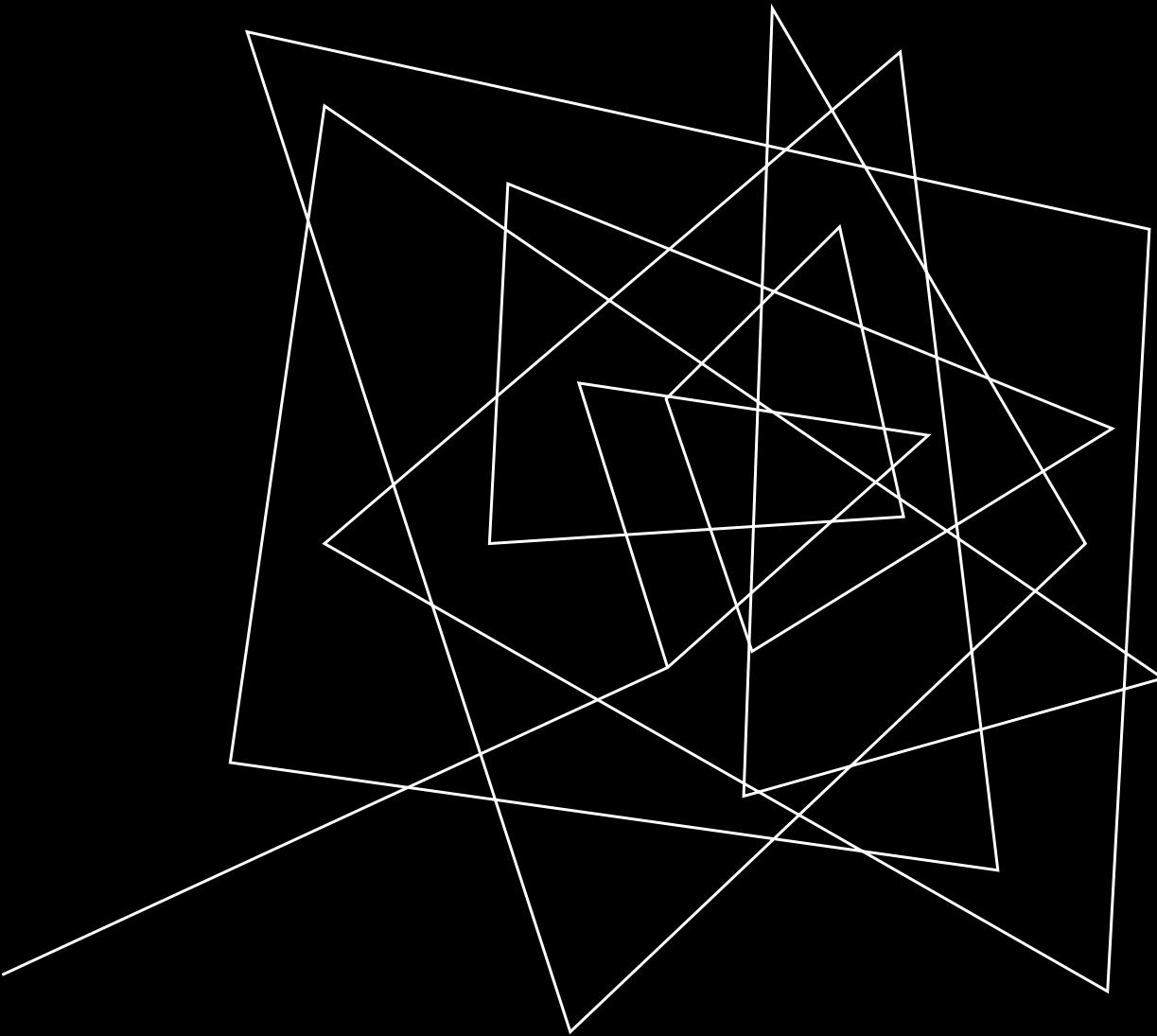
INTERNATIONAL  
BIOLOGY  
OLYMPIAD e.V.



# INTRODUCTION TO BIOINFORMATICS

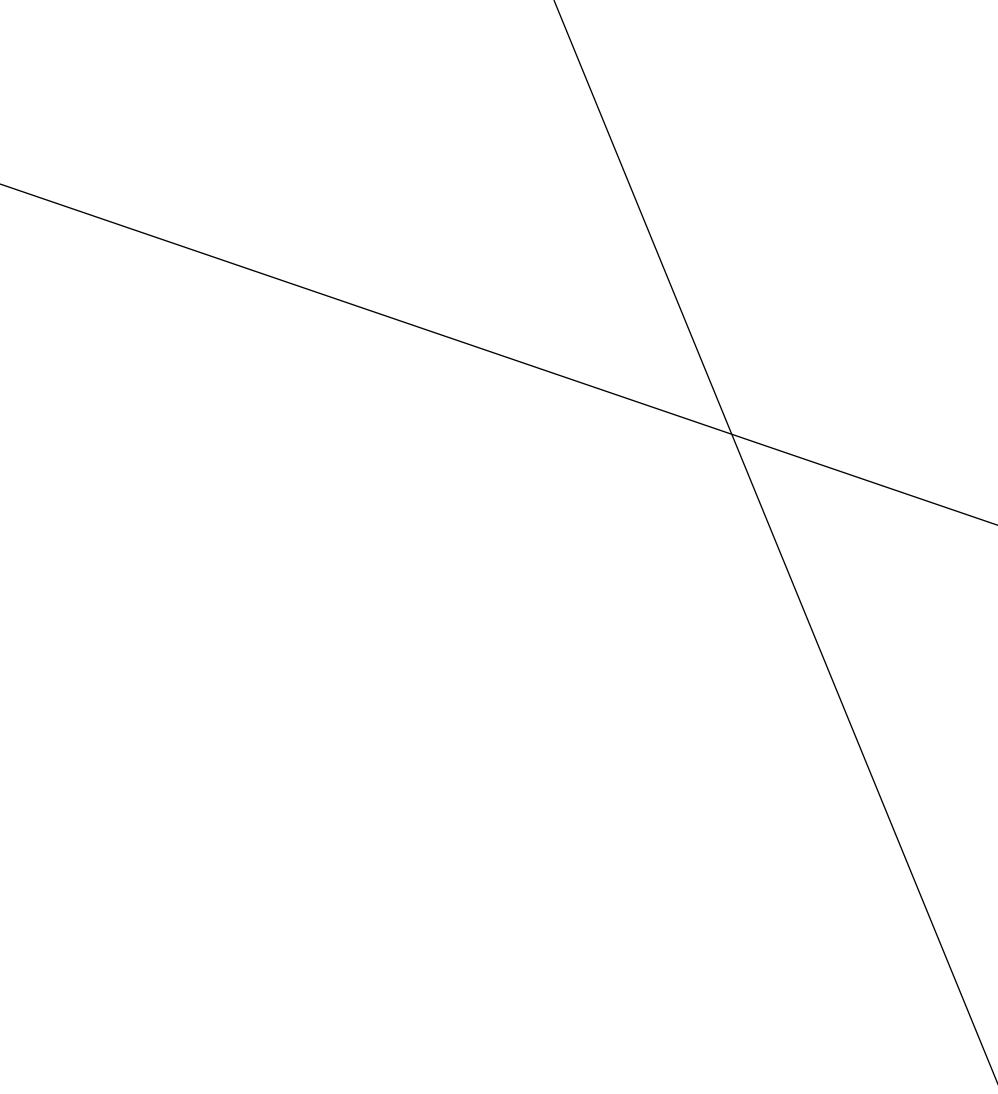
IrBO 25 | By Mohammad Ebrahim Katebi





# ALIGNMENT

THE FOUNDATION OF  
BIOINFORMATICS



# Why?

Conservation and Homology

- ❖ To obtain functional or mechanistic insight
- ❖ To find evolutionarily relation
- ❖ To find structurally or functionally similarity
- ❖ Similarity searching of databases
- ❖ Assembly of sequence reads

What?

### Pairwise Sequence Alignment

EEELTKPRLLWALYFNMRDALSSG  
VEKPRILYALYFNMRDSSDE



EEELT**KPR**~~L~~WA**L**YFNMRDAL**SSG**-  
---VE**KP**R**I**LYALYFNMRD--SSDE

Similarity and Identity

Mutations, Deletions and Insertions

What is the best Alignment?

- 4 matches
- 3 mismatches
- 0 gaps

CACTGTA  
||:;:||  
CATGTTA

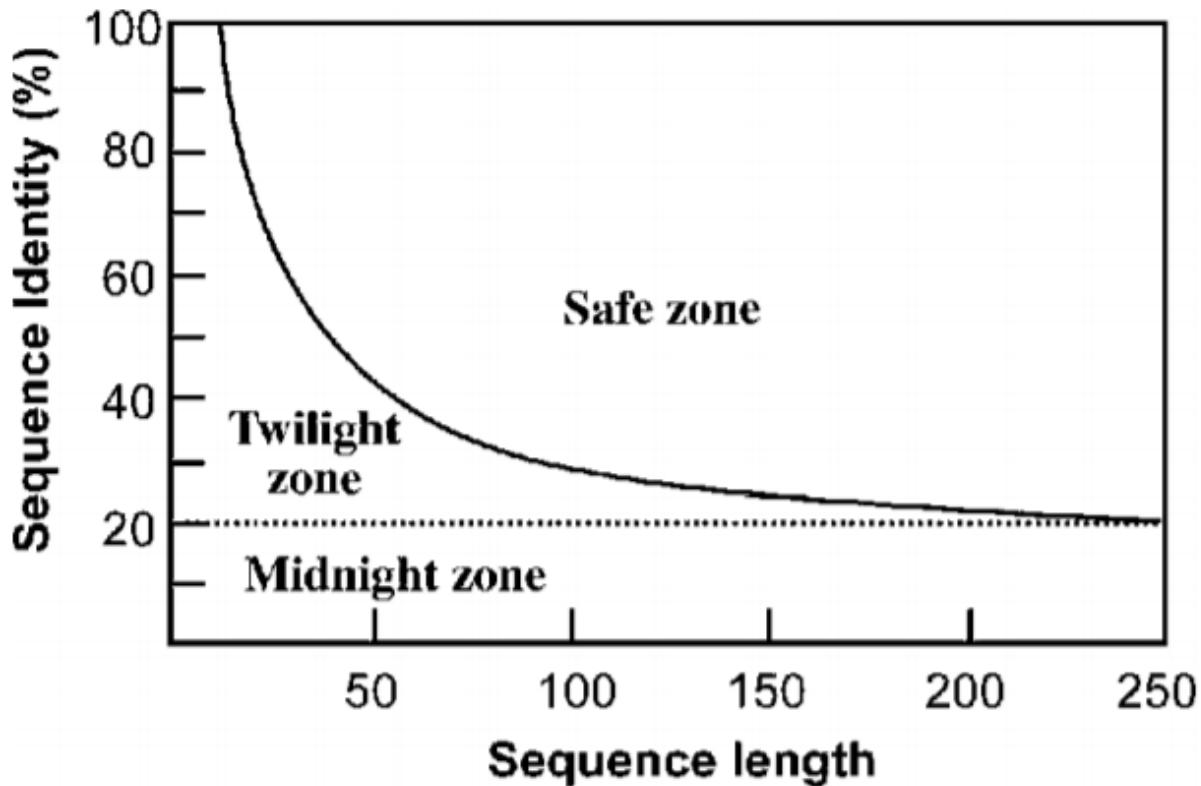
- 6 matches
- 0 mismatches
- 2 gaps

CACTGT-A  
|||:|||  
CA-TGTTA

- 5 matches
- 1 mismatch
- 2 gaps

CAC-TGTA  
||:  
CATGT-TA

What is the best Alignment?



What is the best Alignment?

- 4 matches
- 3 mismatches
- 0 gaps

CACTGTA  
||: : : ||  
CATGTTA

- 6 matches
- 0 mismatches
- 2 gaps

CACTGT-A  
|| | || |  
CA-TGTTA

- 5 matches
- 1 mismatch
- 2 gaps

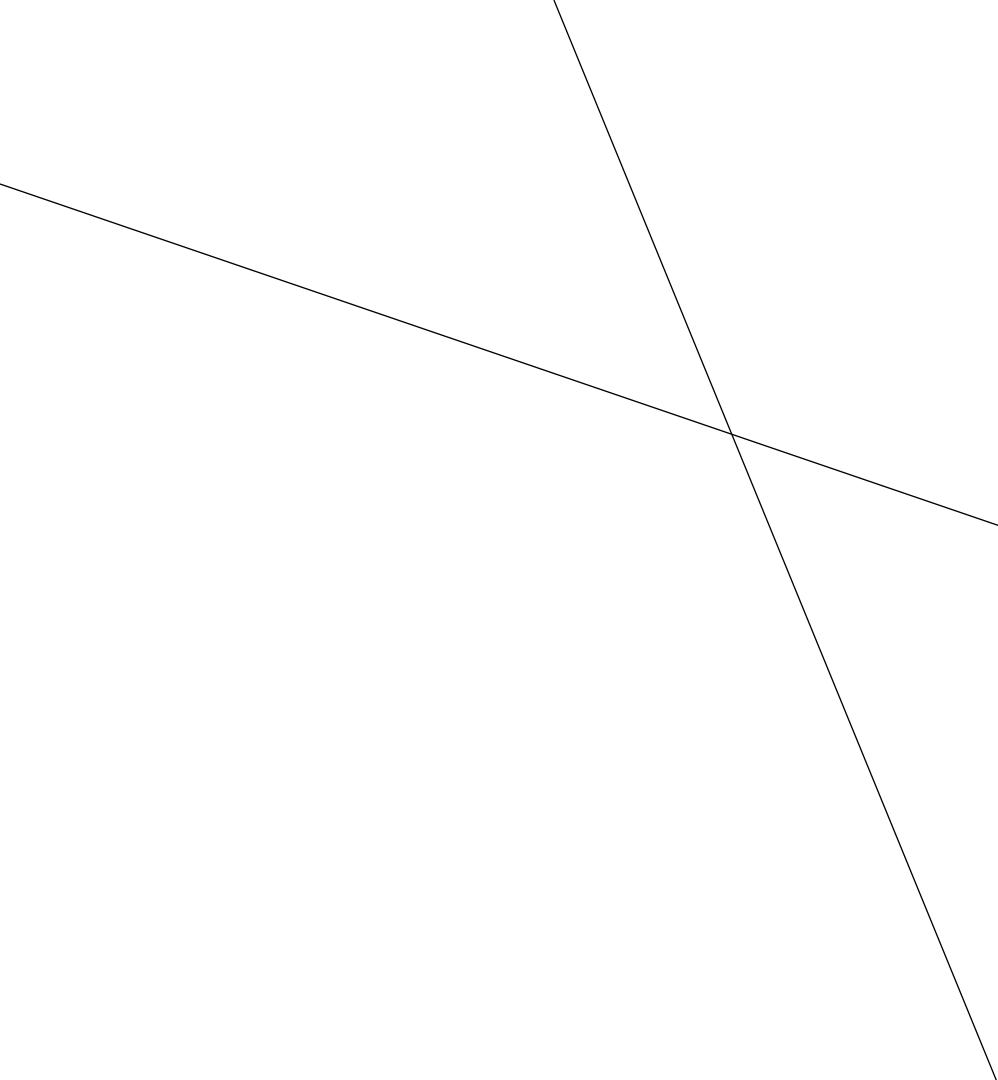
CAC-TGTA  
|| : | ||  
CATGT-TA

Scoring:

Match : +3

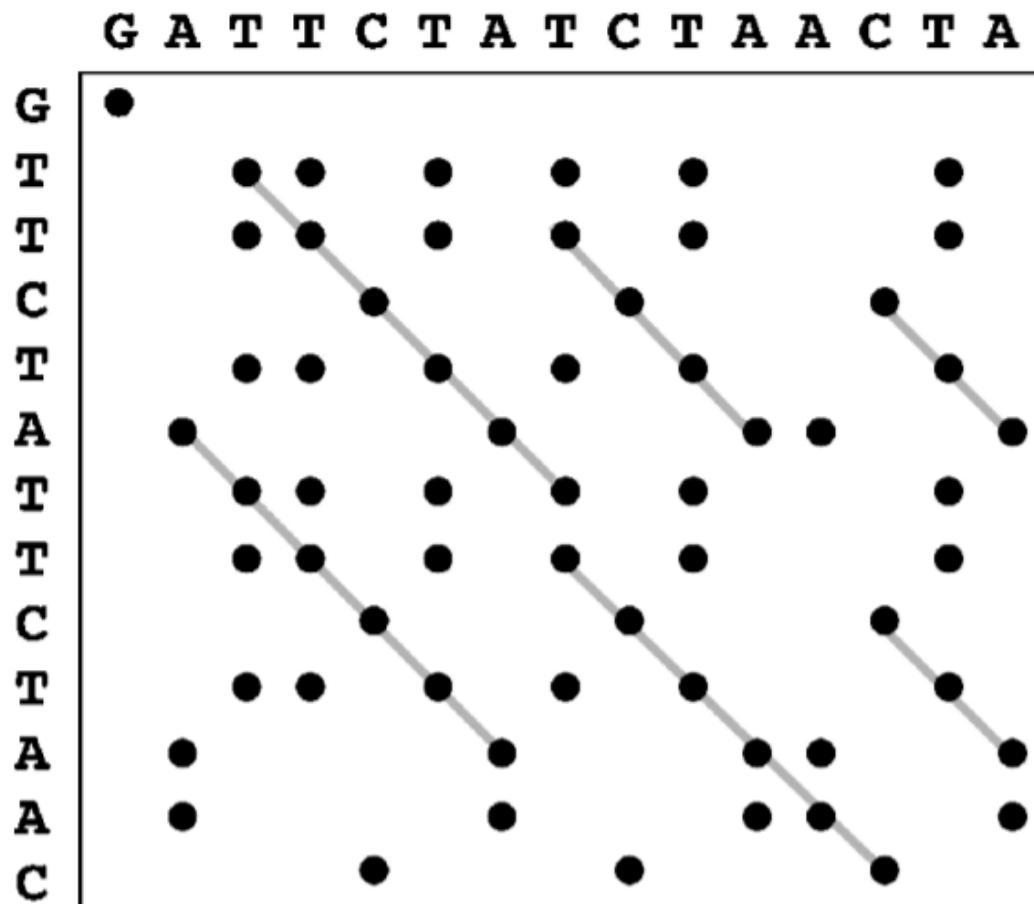
Mismatch : -1

Gap : -2



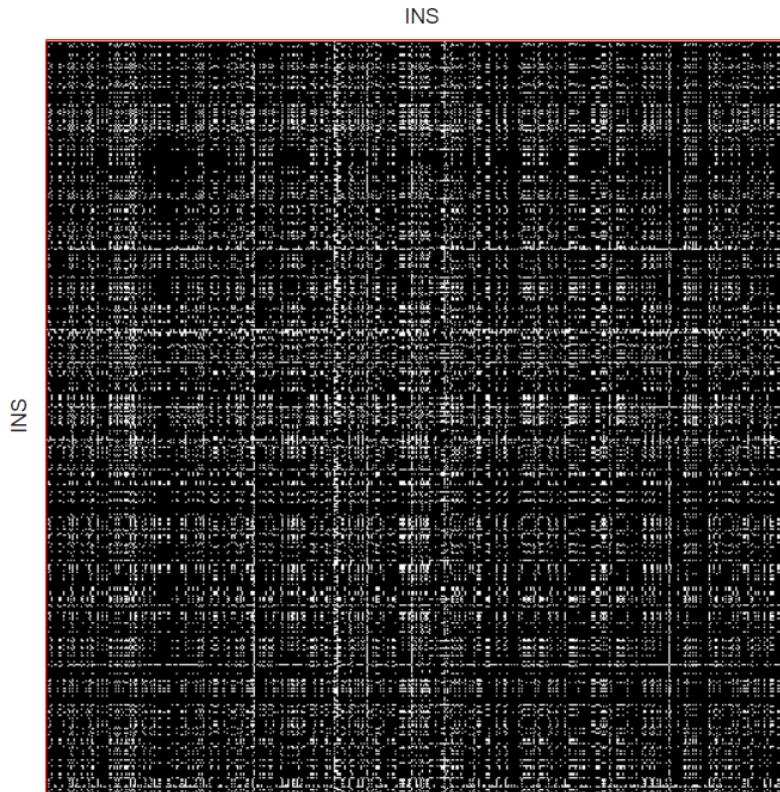
How?

Dot plots



How a Dotplot of Homo sapiens insulin (INS) against itself would look like?

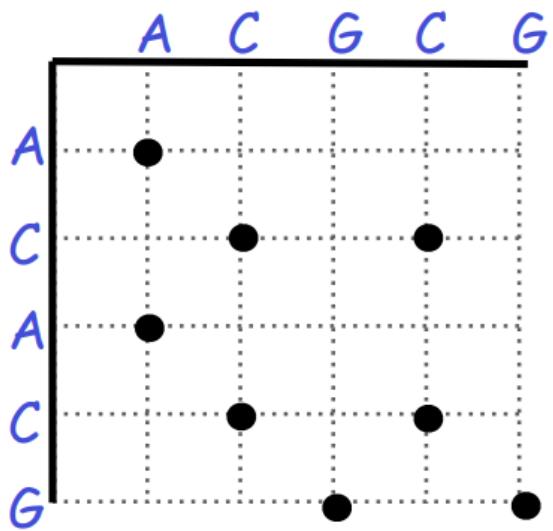
How a Dotplot of Homo sapiens insulin (INS) against itself would look like?



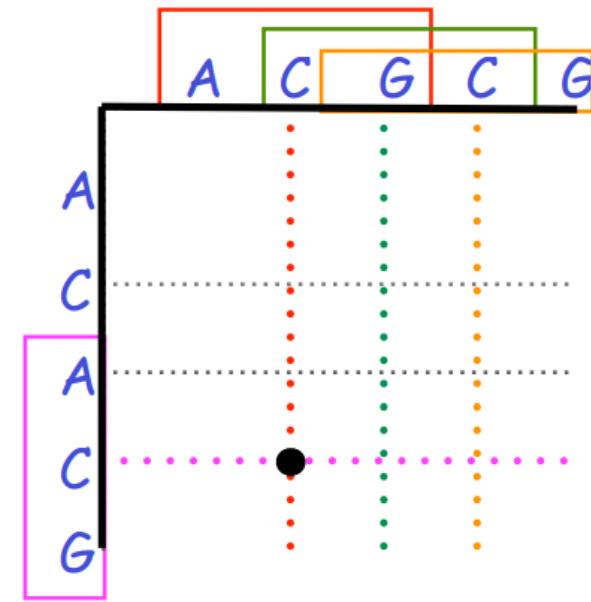
How a Dotplot of Homo sapiens insulin (INS) against itself would look like?

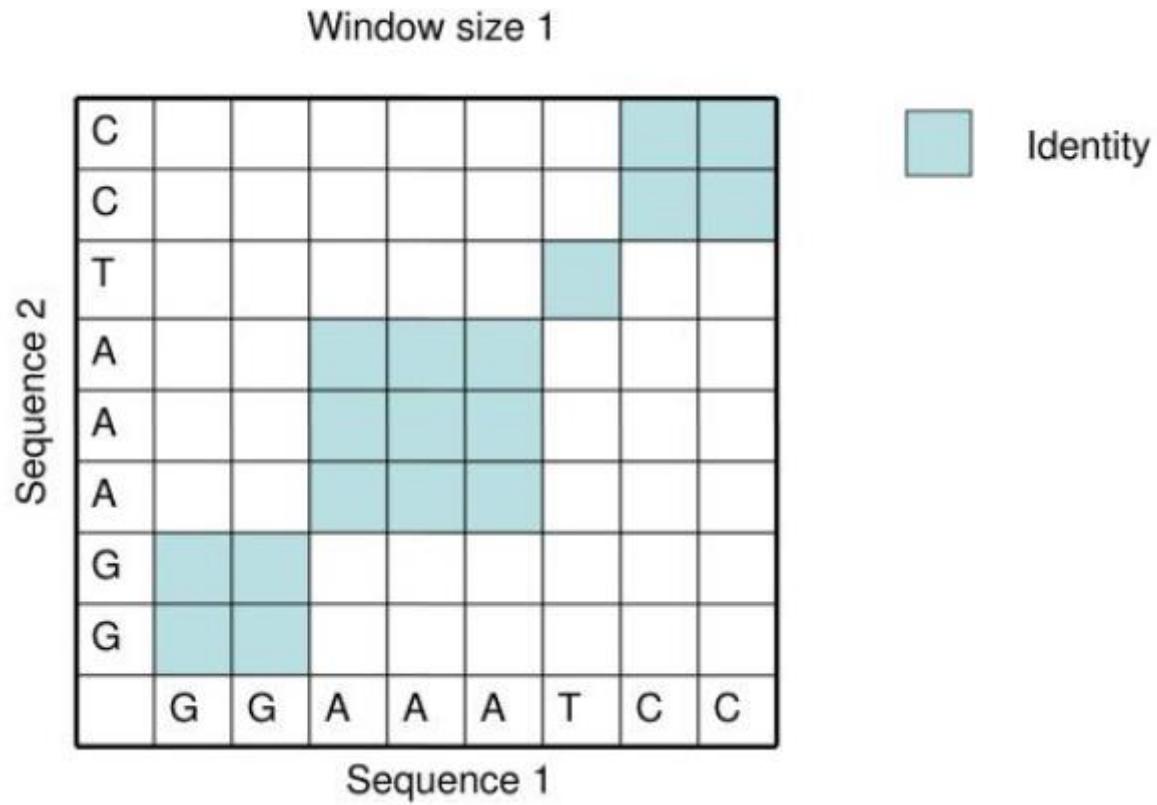


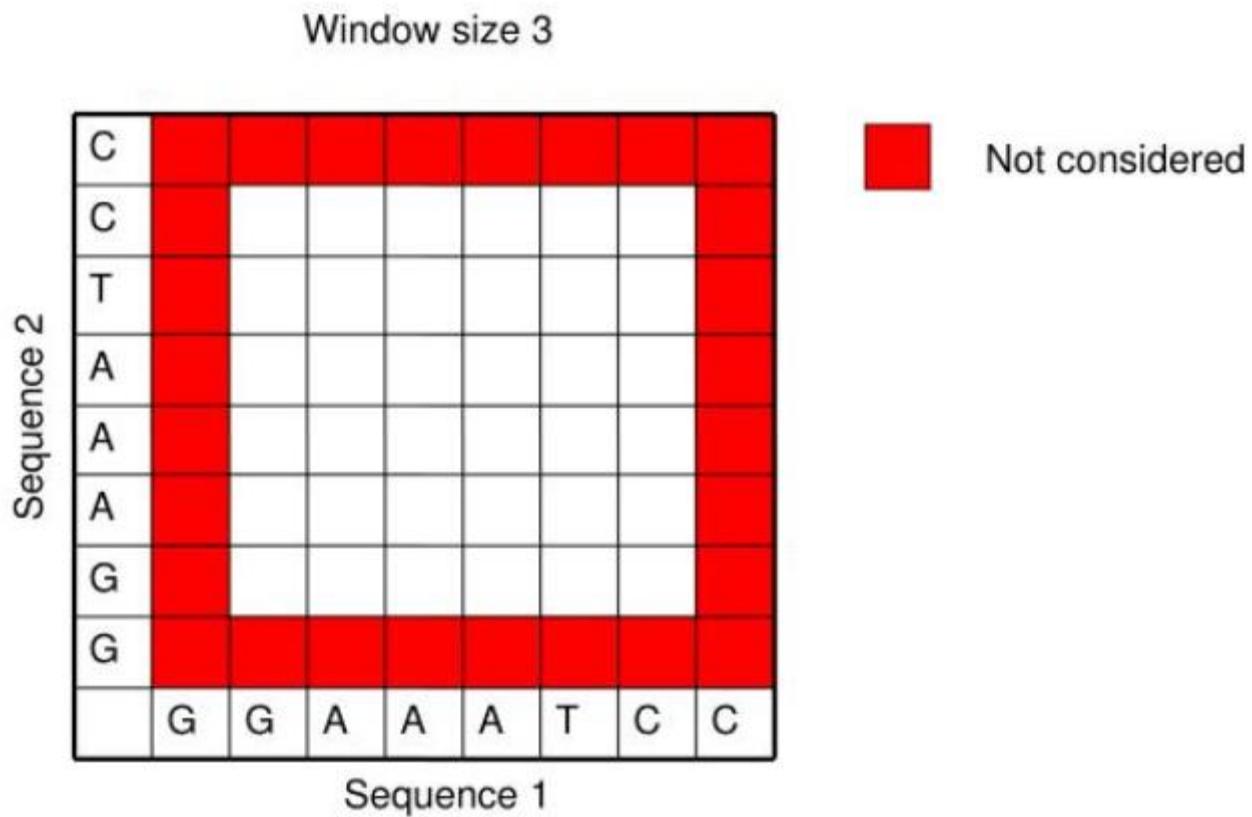
Noise reduction:  
Window and Threshold

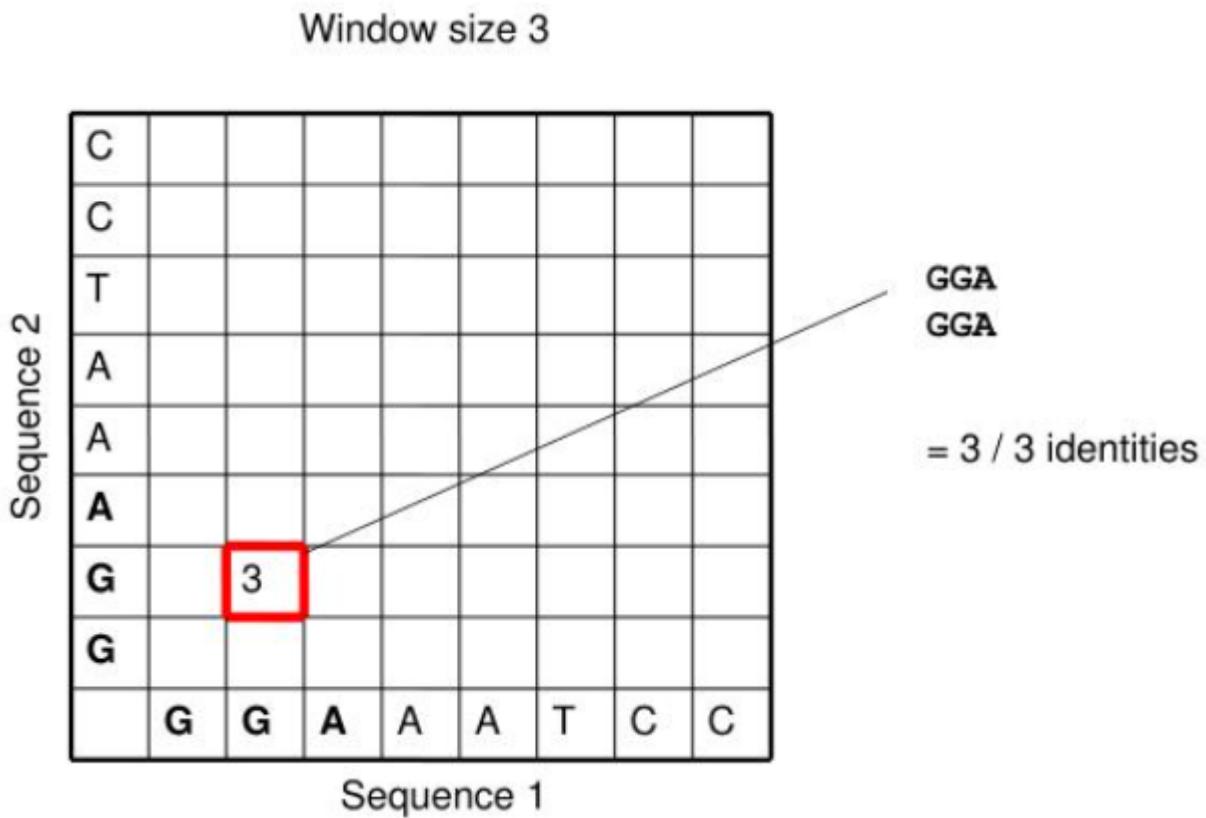


Filter  
→  
Window = 3  
Stringency = 3









Window size 3

	Sequence 1							
Sequence 2	G	G	A	A	A	T	C	C
C								
C								
T								
A								
<b>A</b>								
<b>A</b>			2					
<b>G</b>			3					
G								

**GGA**  
**GAA**

= 2 / 3 identities

Window size 3

C							
C	0	0	0	0	1	3	
T	0	0	1	1	3	1	
A	0	1	2	3	1	0	
A	1	2	3	2	1	0	
A	2	3	2	1	0	0	
G	3	2	1	0	0	0	
G							
	G	G	A	A	A	T	C
							C

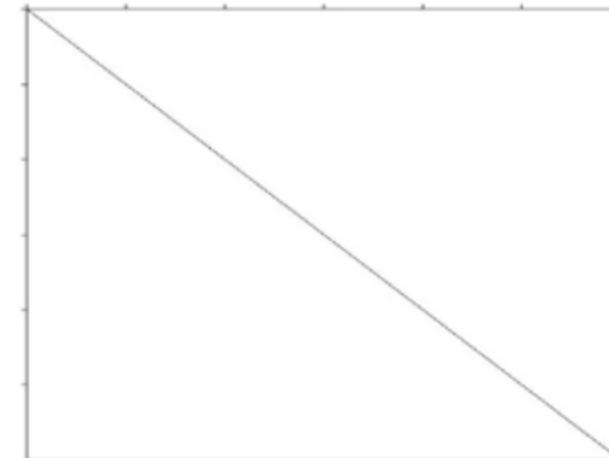
## Identical sequence

- These are the two identical sequences:
  - Seq1: MALWGRL
  - Seq2: MALWGRL

## Identical sequence

- These are the two identical sequences:
  - Seq1: MALWGRL
  - Seq2: MALWGRL

	M	A	L	W	G	R	L
M	*						
A		*					
L			*				
W				*			
G					*		
R						*	
L							*



## Direct repeat

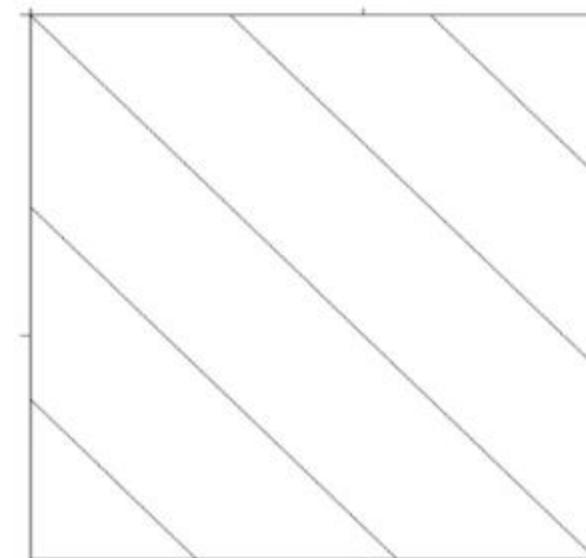
Direct repeats  
ACDEFGHIACDEFGHIACDEFGHI.

## Direct repeat

Direct repeats

ACDEFGHIACDEFGHIACDEFGHI.

	A	B	C	D	E	F	G	H	I	A	B	C	D	E	F	G	H	I	A	B	C	D	E	F	G	H	I
A	*									*									*								
B	*									*									*								
C	*									*									*								
D	*									*									*								
E	*									*									*								
F	*									*									*								
G	*									*									*								
H	*									*									*								
I	*									*									*								
A	*									*									*								
B	*									*									*								
C	*									*									*								
D	*									*									*								
E	*									*									*								
F	*									*									*								
G	*									*									*								
H	*									*									*								
I	*									*									*								



## **Inverted repeat**

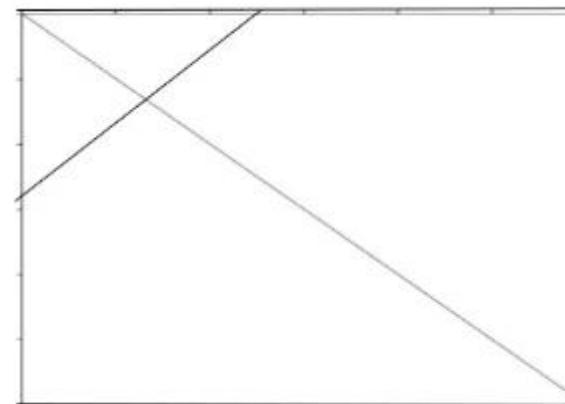
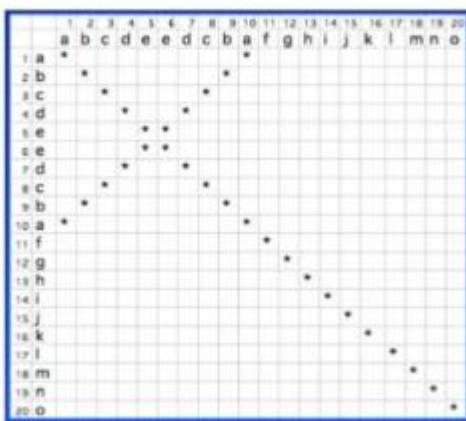
An inverted repeat is sequence of nucleotides followed downstream by its reverse complement.

Inverted repeat: abcdeedcba $\xrightarrow{\quad\quad}$ fghijklmno

## Inverted repeat

An inverted repeat is sequence of nucleotides followed downstream by its reverse complement.

Inverted repeat: abcdeedcbafghijklmn → ←



# Frame shifts

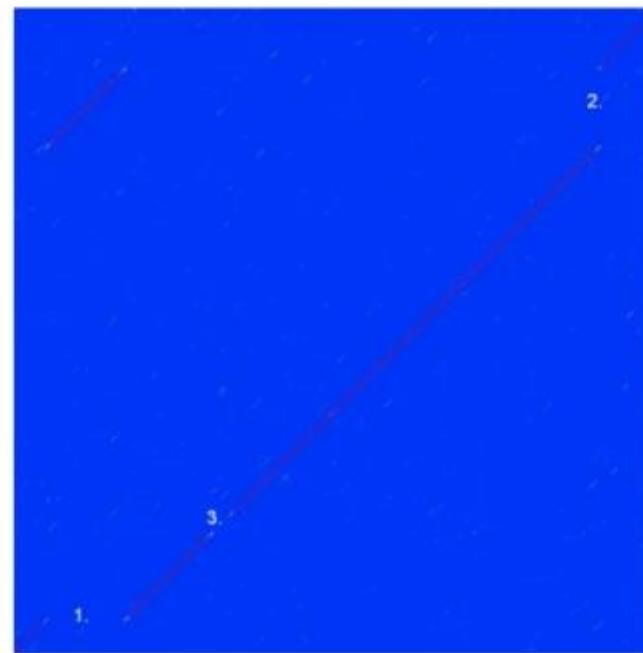
Frame shifts in a nucleotide sequence can occur due to insertions, deletions or mutations.

1. Deletion of nucleotides
2. Insertion of nucleotides
3. Mutation (out of frame)

## Frame shifts

Frame shifts in a nucleotide sequence can occur due to insertions, deletions or mutations.

1. Deletion of nucleotides
2. Insertion of nucleotides
3. Mutation (out of frame)



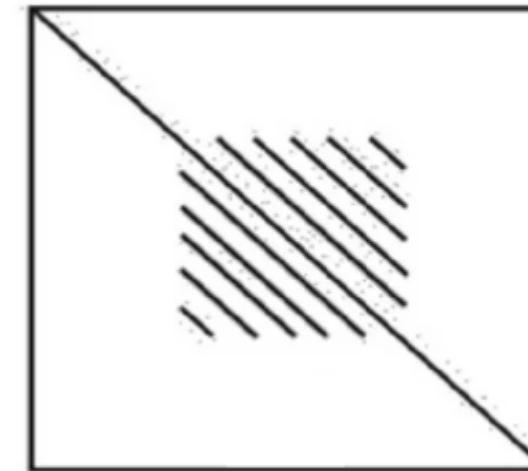
## Low complexity region

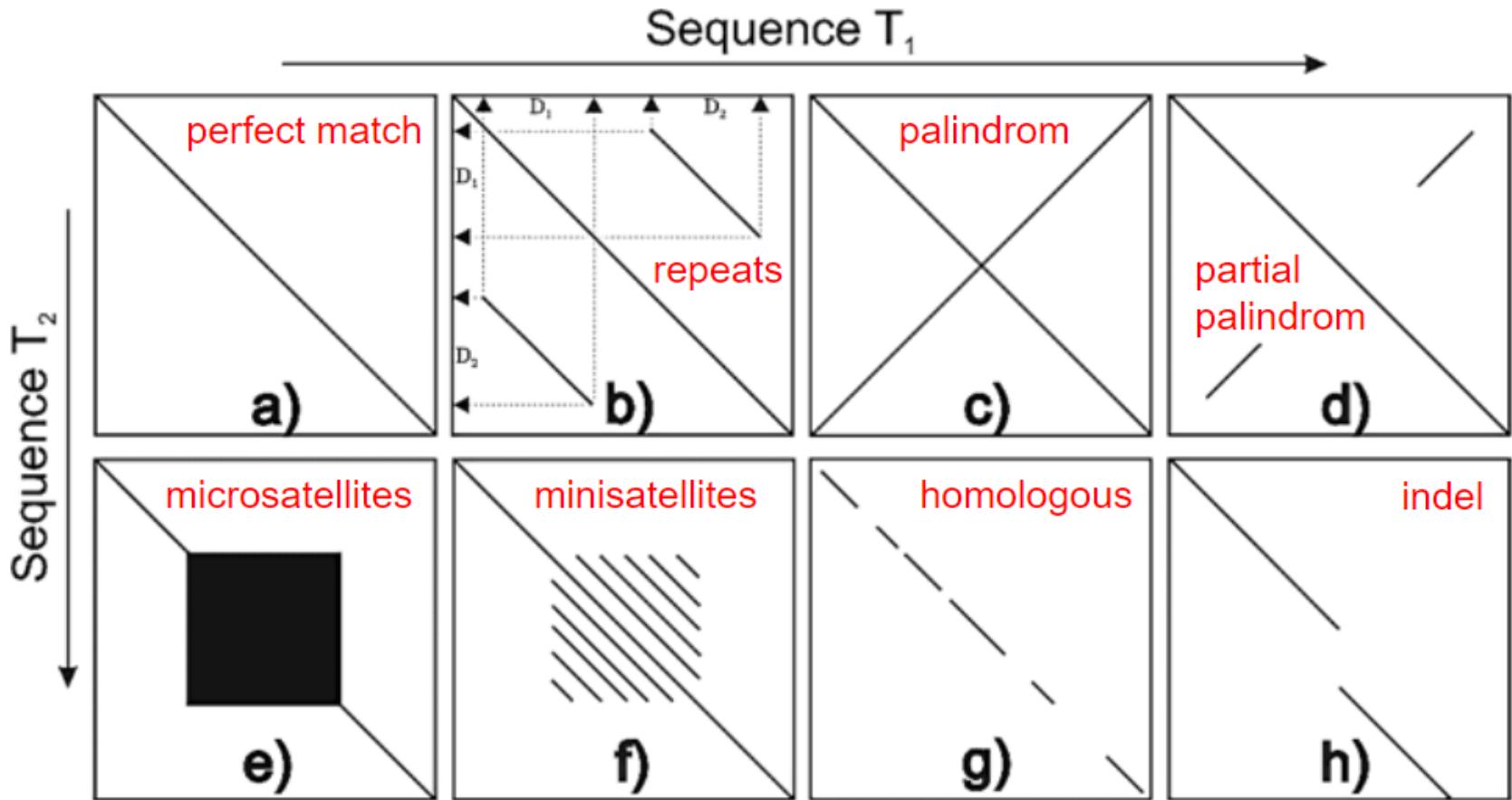
- Low-complexity regions in sequences can be found as regions around the diagonal all obtaining a high score. Low complexity regions are calculated from the redundancy of amino acids within a limited region [Wootton and Federhen,1993].

## Low complexity region

- Low-complexity regions in sequences can be found as regions around the diagonal all obtaining a high score. Low complexity regions are calculated from the redundancy of amino acids within a limited region [Wootton and Federhen,1993].

	C	D	E	E	E	E	F	G
C	0							
D		0						
E			0	0	0	0	0	
E			0	0	0	0	0	
E			0	0	0	0	0	
E			0	0	0	0	0	
E			0	0	0	0	0	
F						0		
G							0	

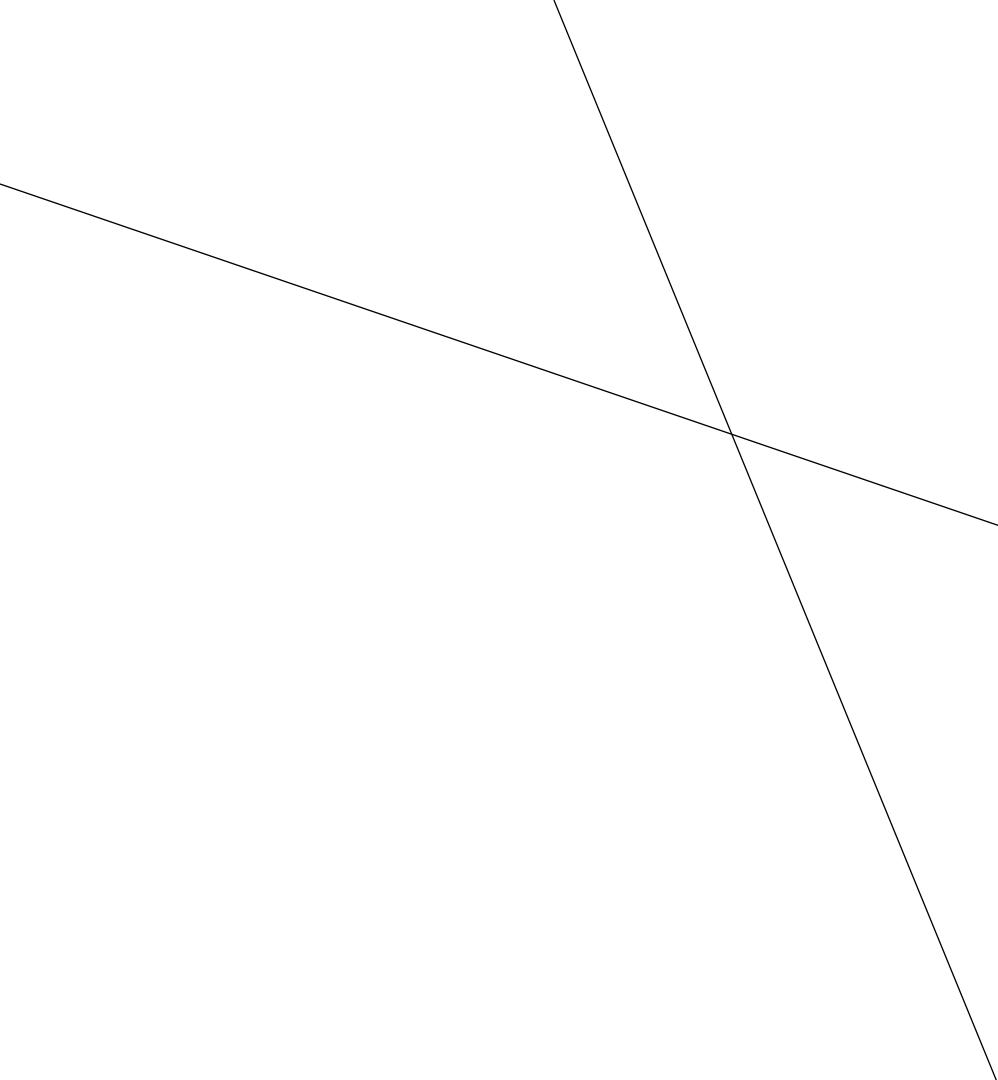




Your Turn! :)

<https://dotlet.vital-it.ch/>

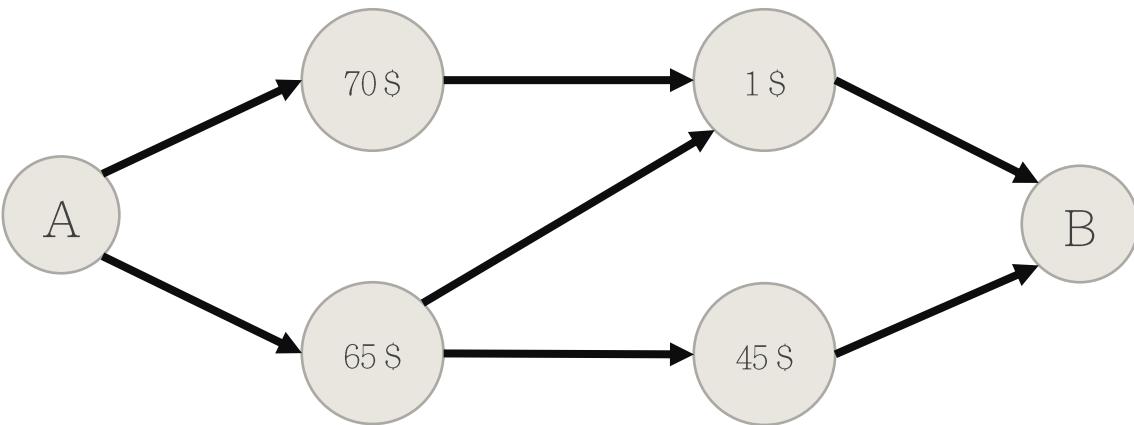
<https://www.ebi.ac.uk/Tools/seqstats/>



How?  
Dynamic Programming

Optimization:

- ❖ Exact algorithms
- ❖ Heuristic algorithms (Greedy Algorithms)



Dynamic Programming: simplifying a complicated problem by breaking it down into simpler sub-problems.

Dynamic programming may be used to solve problems with:

- ❖ Optimal Substructure: The optimal solution to an instance of the problem contains optimal solutions to subproblems.
- ❖ Overlapping Subproblems: There are a limited number of subproblems, many/most of which are repeated many times.
- ✓ Alignment scores are additive.

This allows us to create a **matrix** M indexed by i and j, which are positions in two sequences S and T to be aligned. The best alignment of S and T corresponds with the best path through the matrix M after it is filled in using a recursive formula.

## The Needleman—Wunsch algorithm

Scoring:

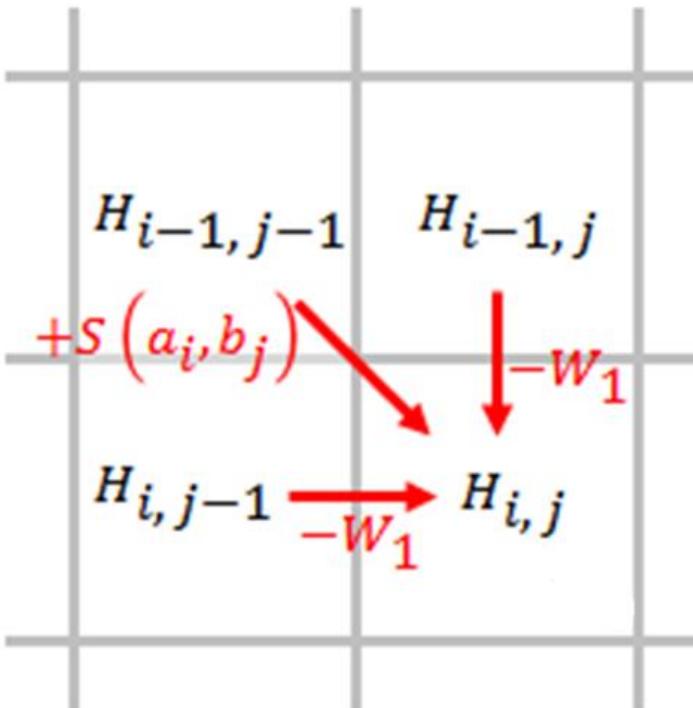
Match : +3

Mismatch : -1

Gap : -2

		A	G	T	T	C
		0				
A						
	T					
	T					
	G					
	C					

## The Needleman-Wunsch algorithm



$$F_{ij} = \max(F_{i-1,j-1} + S(A_i, B_j), F_{i,j-1} + d, F_{i-1,j} + d)$$

## The Needleman—Wunsch algorithm

Scoring:

Match : +3

Mismatch : -1

Gap : -2

		A	G	T	T	C
		0				
A						
	T					
	T					
	G					
	C					

## The Needleman—Wunsch algorithm

Scoring:

Match : +3

Mismatch : -1

Gap : -2

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

## The Needleman—Wunsch algorithm

Scoring:

Match : +3

Mismatch : -1

Gap : -2

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

## The Needleman—Wunsch algorithm

Scoring:

Match : +3

Mismatch : -1

Gap : -2

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

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

Score = 8

Your Turn! :)

Seq 1: TTAGCG  
Seq 2: AAGCTG

The Needleman—Wunsch algorithm

Scoring

Match : +3  
Mismatch : -1  
Gap : -1

Your Turn! :)

Seq 1: TTAGCG  
Seq 2: AAGCTG

The Needleman—Wunsch algorithm

<http://rna.informatik.uni-freiburg.de/Teaching/>

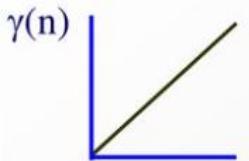
Scoring

Match : +3

Mismatch : -1

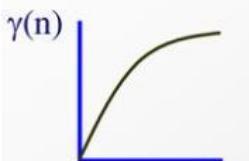
Gap : -1

## Generalized gap penalty



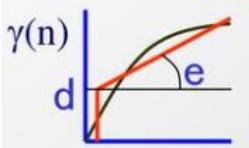
**Linear gap penalty:**  $w(k) = k * p$

- State: Current index tells if in a gap or not
- Achievable using quadratic algorithm (even w/ linear space)



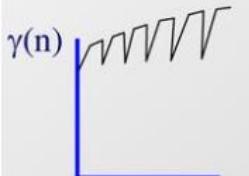
**Quadratic:**  $w(k) = p + q * k + r * k^2$ .

- State: needs to encode the length of the gap, which can be  $O(n)$
- To encode it we need  $O(\log n)$  bits of information. Not feasible



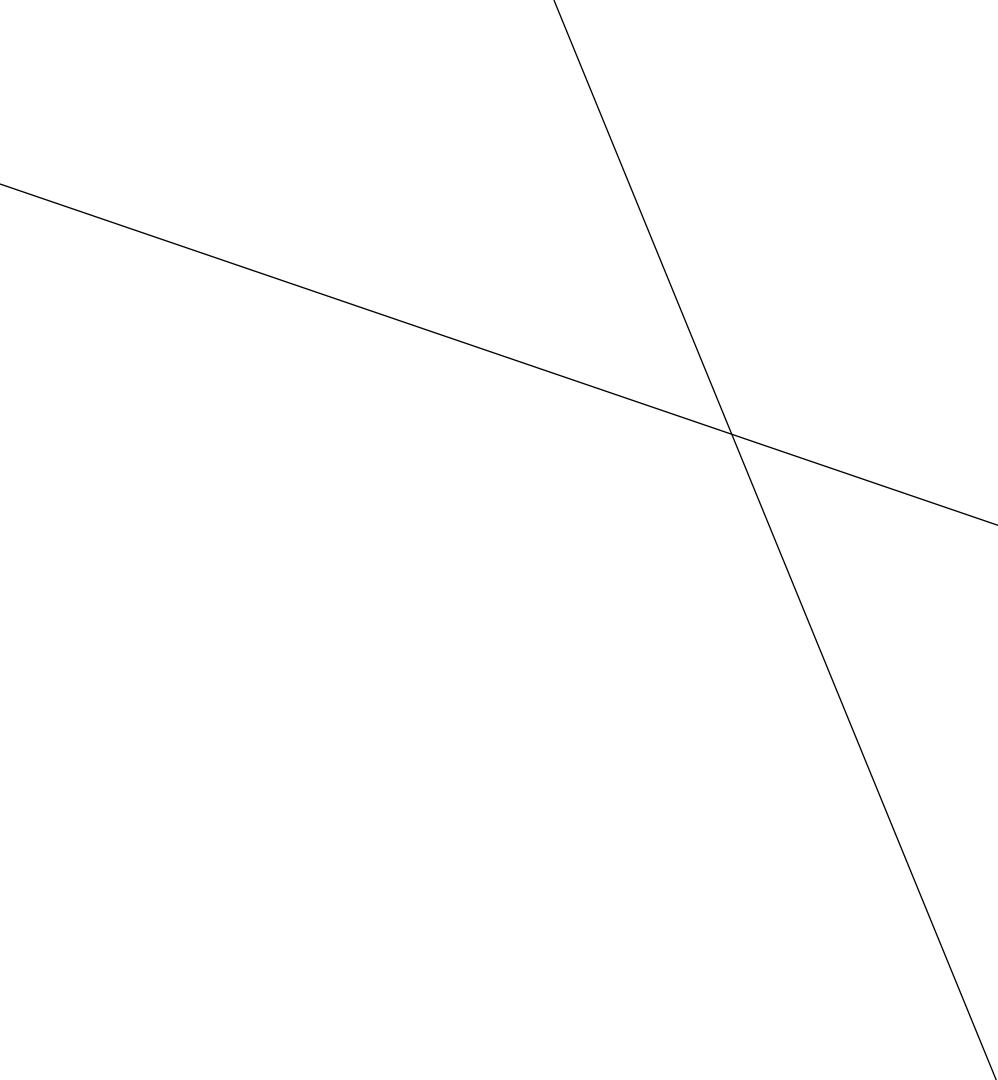
**Affine gap penalty:**  $w(k) = p + q * k$ , where  $q < p$

- State: add binary value for each sequence: starting a gap or not
- Implementation: add second matrix for already-in-gap (recitation)



**Length (mod 3) gap penalty for protein-coding regions**

- Gaps of length divisible by 3 are penalized less: conserve frame
- This is feasible, but requires more possible states
- Possible states are: starting, mod 3=1, mod 3=2, mod 3=0



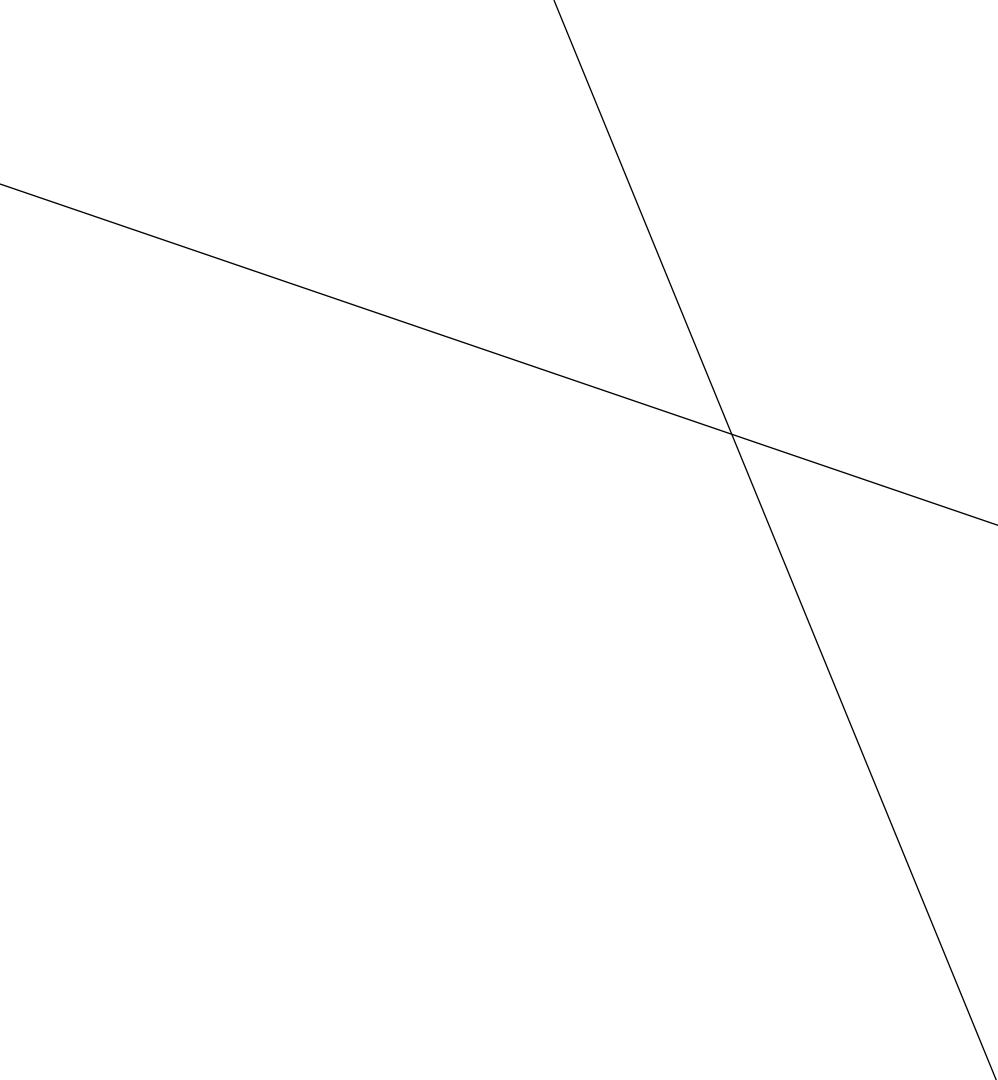
## Semi-global Alignment

## Semi-Global Alignment

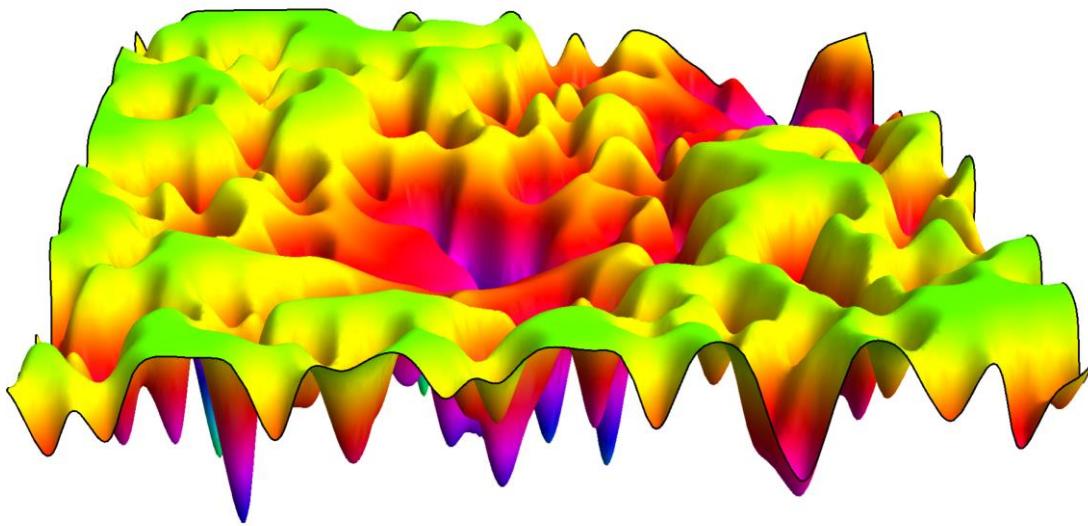
	0	S	T	R	I	N	G	1
S	□	0	0	0	0	0	0	0
T	□							□
R	□							□
I	□							□
N	□							□
G	□							□
2	2	-7	-8	-3	-4	-5	-4	-5
								-6

*For free initial gaps in string 2, initialize this row to all "0"s*

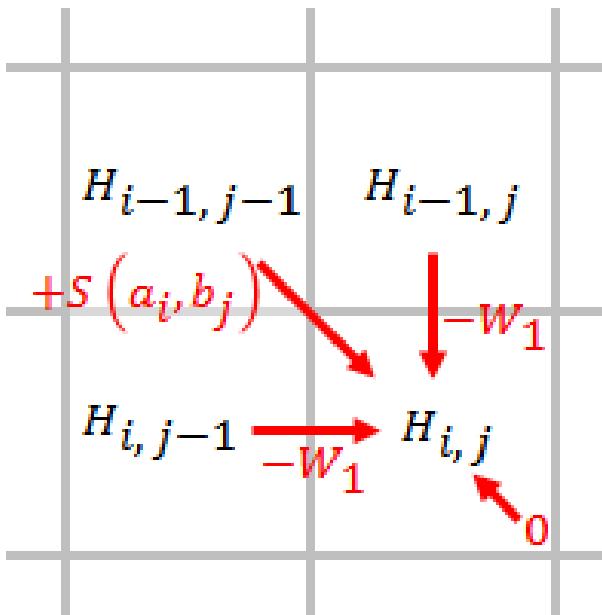
*For free end gaps in string 2, select the greatest element in the last row, and align accordingly*



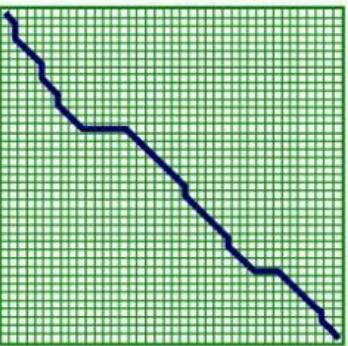
# Local Alignment



## The Smith—Waterman algorithm



$$H_{ij} = \max \begin{cases} H_{i-1,j-1} + s(a_i, b_j), \\ H_{i-1,j} - W_1, \\ H_{i,j-1} - W_1, \\ 0 \end{cases}$$



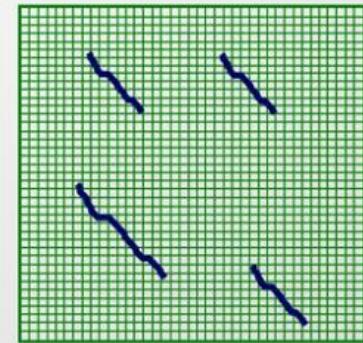
### Needleman-Wunsch algorithm

Initialization:  $F(0, 0) = 0$

Iteration:

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

Termination: Bottom right



### Smith-Waterman algorithm

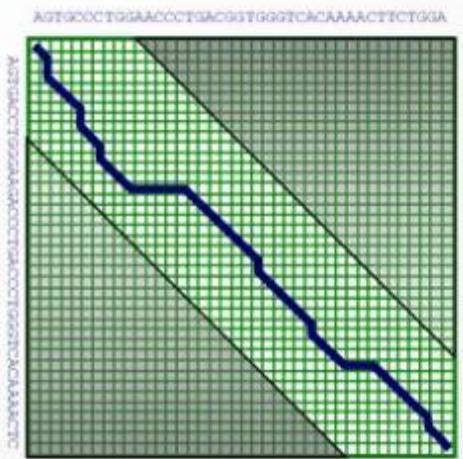
Initialization:  $F(0, j) = F(i, 0) = 0$

Iteration:

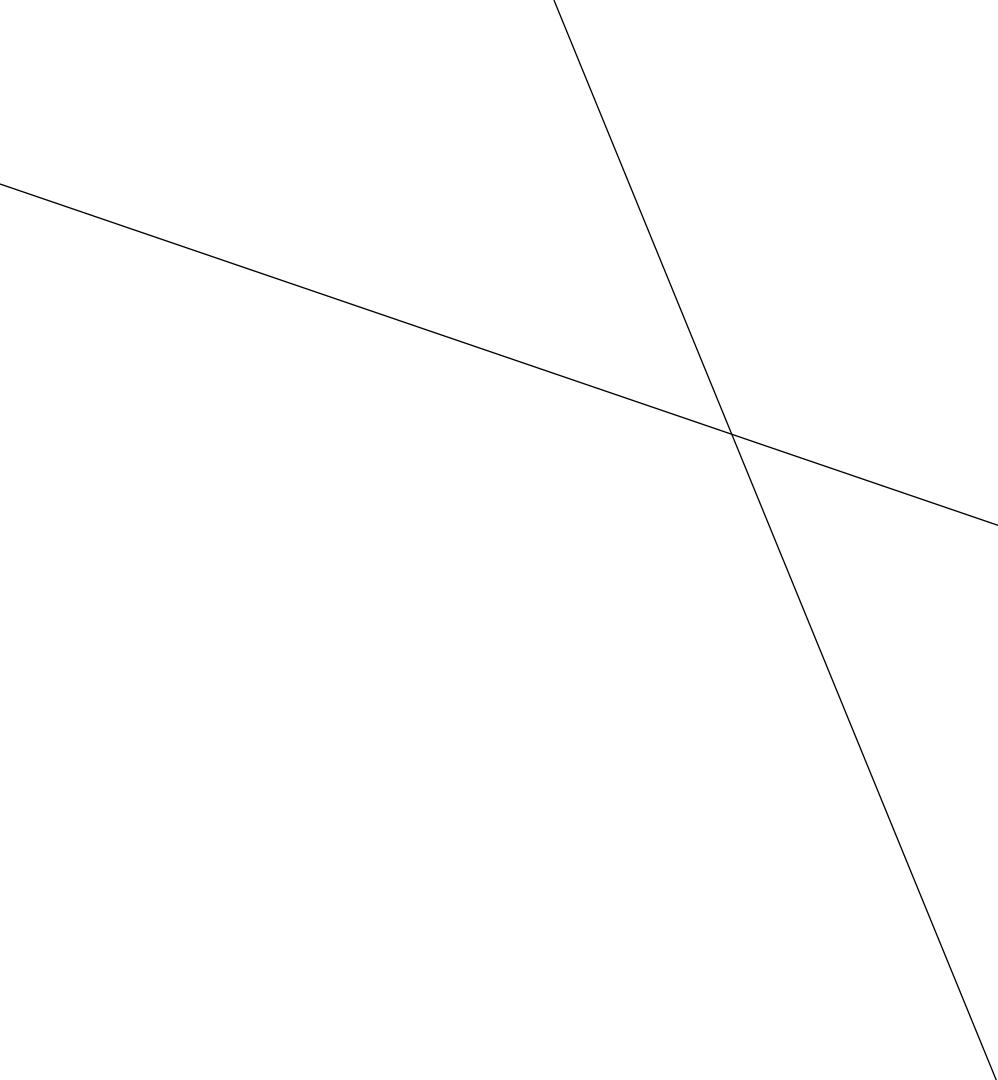
$$F(i, j) = \max \begin{cases} 0 \\ F(i - 1, j) - d \\ F(i, j - 1) - d \\ F(i - 1, j - 1) + s(x_i, y_j) \end{cases}$$

Termination: Anywhere

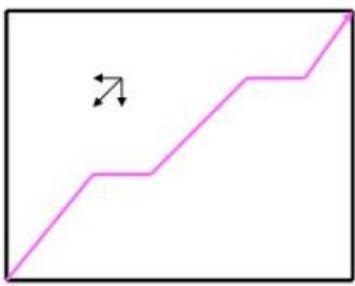
	Global	Local	Semi-global
Initialization	Top left	Top row/left col.	Top row
Iteration: max	$F(i - 1, j) - d$ $F(i, j - 1) - d$ $F(i - 1, j - 1) + s(x_i, y_j)$	0 $F(i - 1, j) - d$ $F(i, j - 1) - d$ $F(i - 1, j - 1) + s(x_i, y_j)$	$F(i - 1, j) - d$ $F(i, j - 1) - d$ $F(i - 1, j - 1) + s(x_i, y_j)$
Termination	Bottom right	Anywhere	Right column



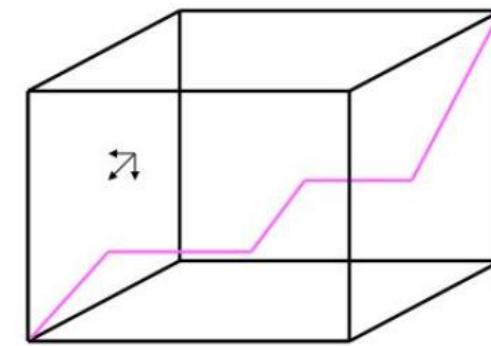
- Save time: Bounded-space computation
  - Space:  $O(k^*m)$
  - Time:  $O(k^*m)$ , where  $k$  = radius explored
  - Heuristic
    - Not guaranteed optimal answer
    - Works very well in practice
  - Practical interest



# Multiple Sequence Alignment (MSA)



Finding a PSA = Finding a path through a 2-Dim matrix.



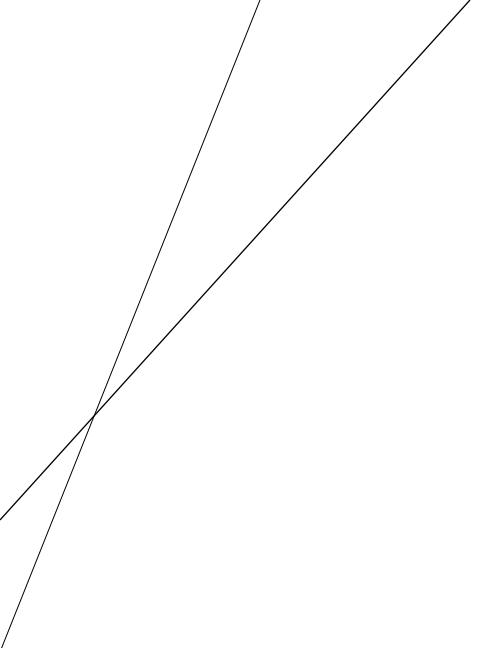
Finding an MSA = Finding a path through an N-Dim matrix. It's  $O(L^N)$ , where N is the number of sequences and L is the sequence length.

# Scoring

Construction of Substitution Matrices

The idea behind the scoring matrix is that the score of alignment should reflect the probability that two similar sequences are homologous.

- ❖ Hypothesis 1: That the alignment between the two sequence is due to chance and the sequences are, in fact, unrelated. (U)
- ❖ Hypothesis 2: That the alignment is due to common ancestry and the sequences are actually related. (R)



For **unrelated** sequences, the probability of having an n-residue alignment between x and y is a simple product of the probabilities of the individual sequences since the residue pairings are independent.

For **unrelated** sequences, the probability of having an n-residue alignment between x and y is a simple product of the probabilities of the individual sequences since the residue pairings are independent.

$$\mathbf{x} = \{x_1 \dots x_n\}$$

$$\mathbf{y} = \{y_1 \dots y_n\}$$

$$q_a = P(\text{ amino acid } a)$$

$$P(\mathbf{x}, \mathbf{y} \mid U) = \prod_{i=1}^n q_{x_i} \prod_{i=1}^n q_{y_i}$$

For related sequences, the residue pairings are no longer independent so we must use a different joint probability, assuming that each pair of aligned amino acids evolved from a common ancestor:

$$p_{ab} = P(\text{ evolution gave rise to } a \text{ in } \mathbf{x} \text{ and } b \text{ in } \mathbf{y})$$

$$P(\mathbf{x}, \mathbf{y} \mid R) = \prod_{i=1}^n p_{x_i y_i}$$

A principled way of comparing R and U:  
The Likelihood Ratio between the two

$$\begin{aligned}\frac{P(\mathbf{x}, \mathbf{y} | R)}{P(\mathbf{x}, \mathbf{y} | U)} &= \frac{\prod_{i=1}^n p_{x_i y_i}}{\prod_{i=1}^n q_{x_i} \prod_{i=1}^n q_{y_i}} \\ &= \frac{\prod_{i=1}^n p_{x_i y_i}}{\prod_{i=1}^n q_{x_i} q_{y_i}}\end{aligned}$$

Since we eventually want to compute a sum of scores and probabilities require add products, we take the log of the products.

$$\log \frac{\Pr(x, y | R)}{\Pr(x, y | U)} = \sum_i \log \left( \frac{p_{x_i y_i}}{q_{x_i} q_{y_i}} \right)$$

We define the alignment score as the log of the Likelihood Ratio:

$$S = \sum_i s(x_i, y_i) = \log \frac{\Pr(x, y | R)}{\Pr(x, y | U)} = \sum_i \log \left( \frac{p_{x_i y_i}}{q_{x_i} q_{y_i}} \right)$$

The substitution matrix score for a given pair a,b:

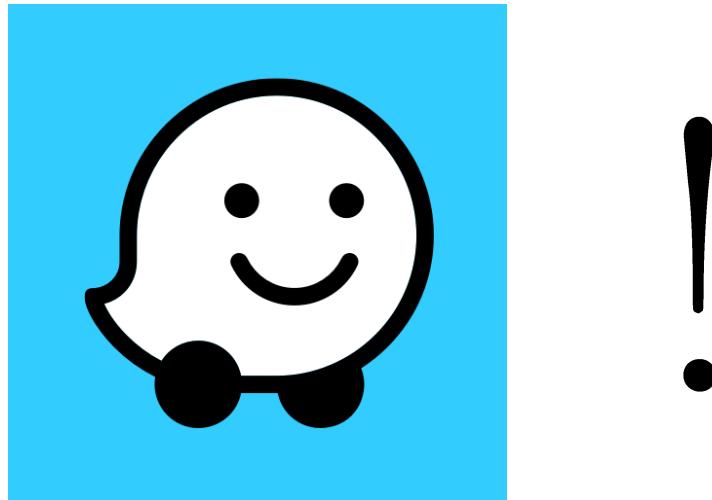
$$s(a, b) = \log\left(\frac{p_{ab}}{q_a q_b}\right)$$

Hypothesis we wish to test; two amino acids are correlated because they are homologous.

$$\text{score} = \log \text{ odds ratio} = s_{AB} \propto \log\left(\frac{\text{observed}}{\text{expected}}\right)$$

Null hypothesis; two amino acids occur independently (and are uncorrelated and unrelated).

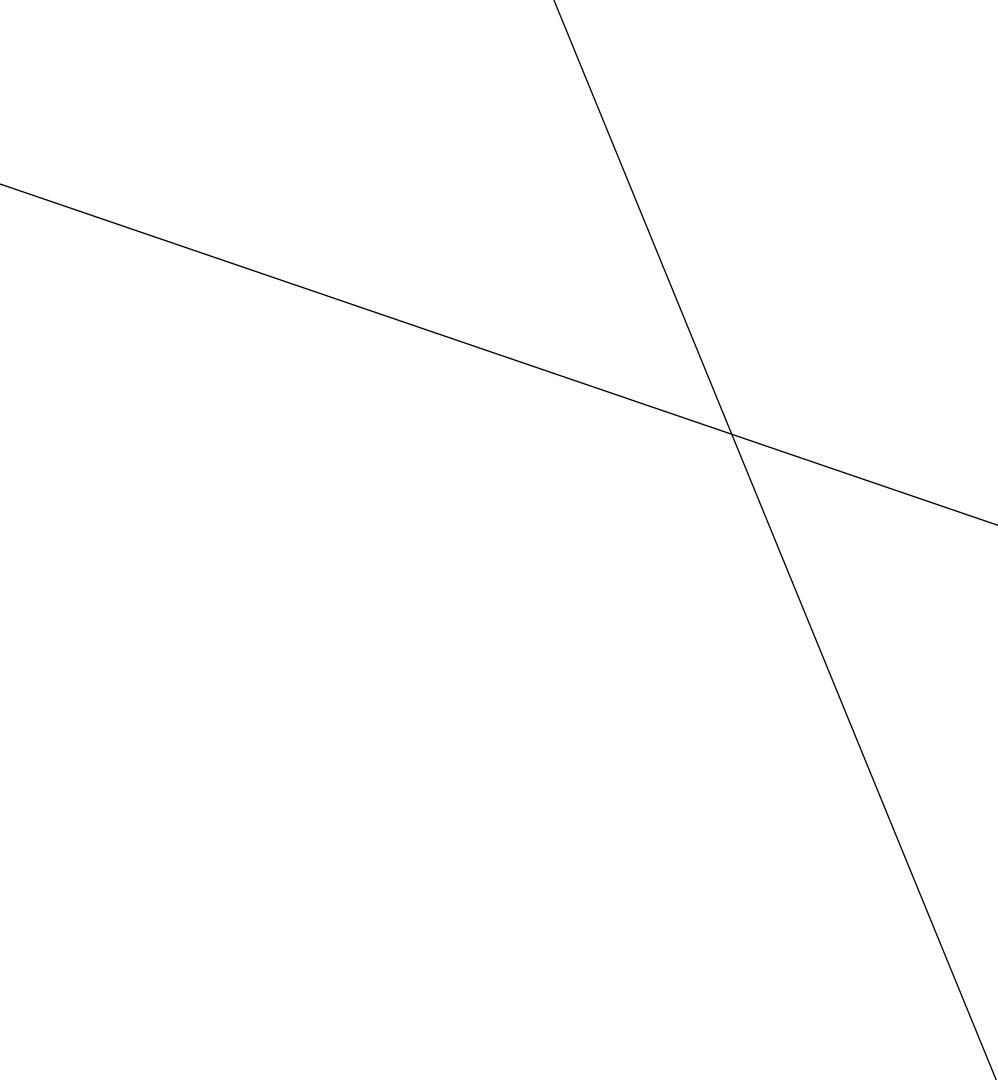
But how do we get  $P_{ab}$ ?



## Key Idea:

Empirical observations

Research first done by Margaret Dayhoff in 1970s and colleagues and later by Henikoff and Henikoff early 1990s resulted in PAM and BLOSUM Substitution matrices and are the most widely used today.



# BLOSUM Matrices

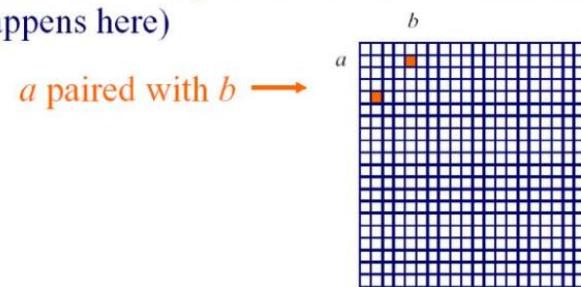
BLOcks SUbstitution Matrix

## BLOSUM Matrices

- [Henikoff & Henikoff, *PNAS* 1992]
- probabilities estimated from “blocks” of sequence fragments that represent structurally conserved regions in proteins
- transition frequencies observed directly by identifying blocks that are at least
  - 45% identical (BLOSUM-45)
  - 50% identical (BLOSUM-50)
  - 62% identical (BLOSUM-62)
  - etc.

## BLOSUM Matrices

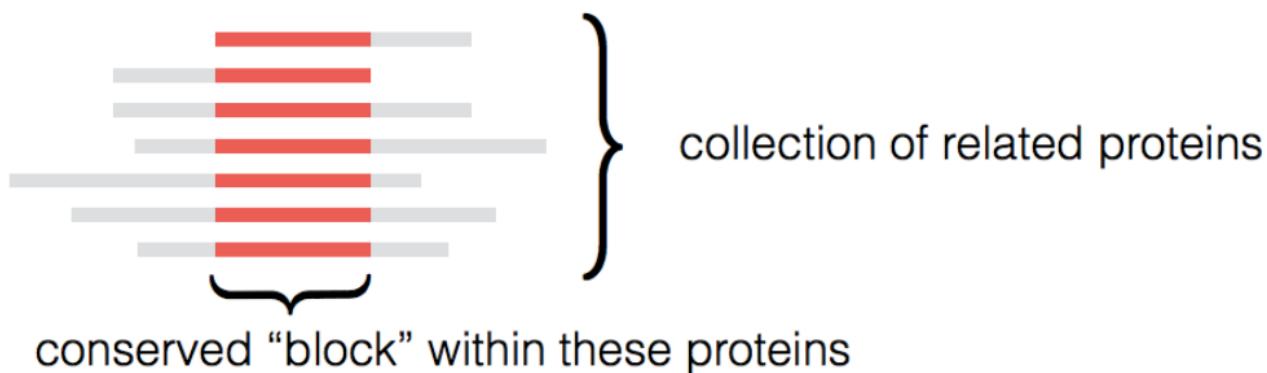
- given: a set of sequences in a block
- fill in matrix  $A$  with number of observed substitutions (we won’t worry about details of some normalization that happens here)



$$p_{ab} = \frac{A_{ab}}{\sum_{c,d} A_{cd}} \quad q_a = \frac{\sum_b A_{ab}}{\sum_{c,d} A_{cd}}$$

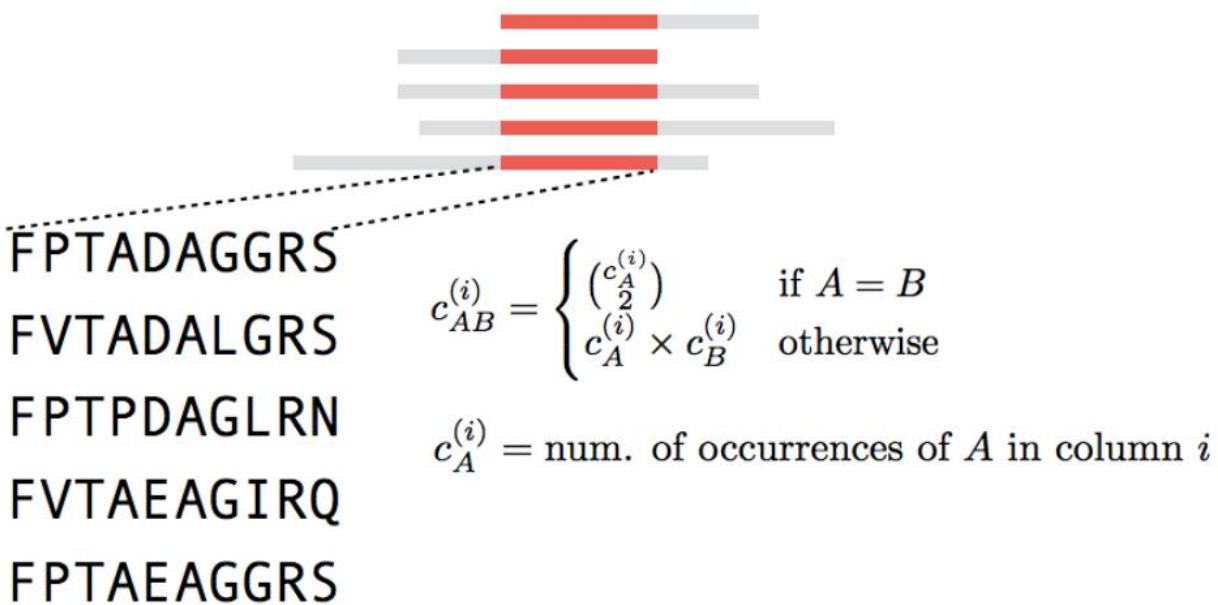
## The Blocks Database

1. Look for conserved (gapless) regions in alignments.

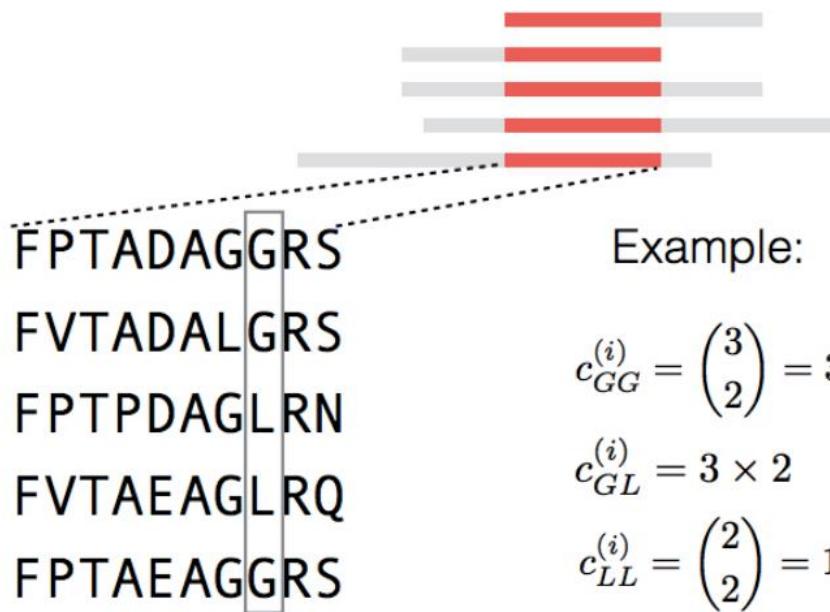


Henikoff and Henikoff analyzed conserved regions of related protein sequences they obtained from BLOCKS database. In total, they examined 2,000 blocks without gaps and 500 groups of related proteins by counting the number of matches and mismatches of each type of the 20 different amino acids.

2. Count all pairs of amino acids in each column of the alignments.



- Count all pairs of amino acids in each column of the alignments.



Example:

$$c_{GG}^{(i)} = \binom{3}{2} = 3$$

$$c_{GL}^{(i)} = 3 \times 2$$

$$c_{LL}^{(i)} = \binom{2}{2} = 1$$

In this column, there are 3 ways to pair G with G, 6 potential ways to pair G with L and 1 potential way to pair L with L.

3. Use amino acid pair frequencies to derive “score” for a mutation/replacement

$$\text{Total # of potential align. between A \& B: } c_{AB} = \sum_i c_{AB}^{(i)}$$

$$\text{Total number of pairwise char. alignments: } T = \sum_{A \geq B} c_{AB}$$

$$\text{Normalized frequency of aligning A \& B: } q_{AB} = \frac{c_{AB}}{T}$$

3. Use amino acid pair frequencies to derive “score” for a mutation/replacement

Probability of occurrence of amino acid A in any {A,B} pair:

$$p_A = q_{AA} + \sum_{A \neq B} q_{AB}$$

Expected likelihood of each {A,B} pair, assuming independence:

$$e_{AB} = \begin{cases} (p_A)(p_B) = (p_A)^2 & \text{if } A = B \\ (p_A)(p_B) + (p_B)(p_A) = 2(p_A)(p_B) & \text{otherwise} \end{cases}$$

To get the final scores in the matrix, Henikoff and Henikoff further converted the Log–Odds Ratios into bit units and multiplied each bit score by a scaling factor of two and rounded to the nearest integer, producing the final scores in BLOSUM matrix.

$$\text{score} = \log \text{ odds ratio} = s_{AB} \propto \log \left( \frac{\text{observed}}{\text{expected}} \right)$$

$$\text{score} = \log \text{ odds ratio} = s_{AB} = \text{Round} \left( \left( \frac{1}{\lambda} \right) \log_2 \left( \frac{q_{AB}}{e_{AB}} \right) \right)$$

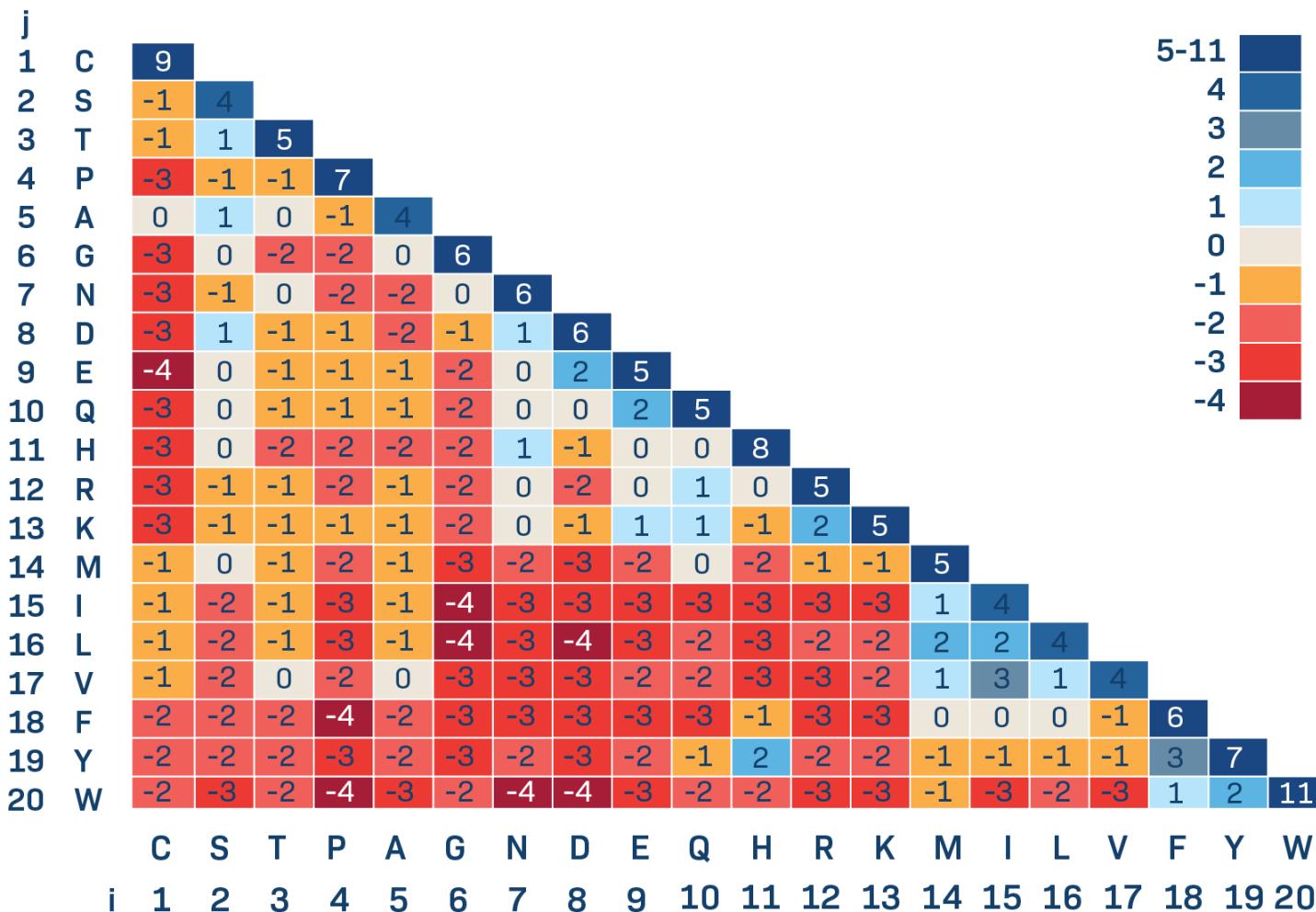
Scaling factor used to produce scores that can be rounded to integers; set to 0.5 in H&H '92.

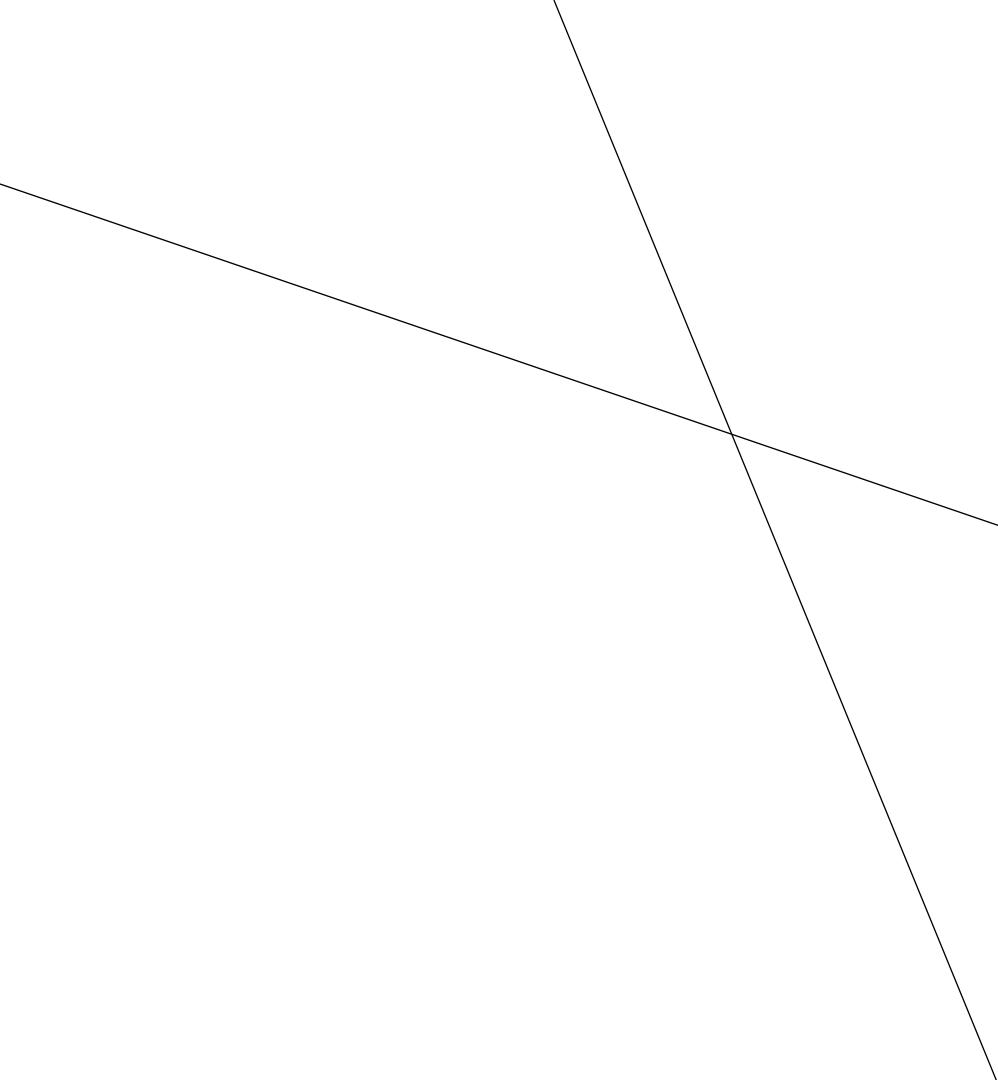
## A Family of Matrices

Henikoff and Henikoff divided the family clusters into sub-clusters by their percentage of similarity.

This division resulted in BLOSUM family of matrices where the associated number, e.g., BLOSUM65 means scores are from a cluster of sequences where sequences are at least 65% similar, in BLOSUM80 matrix scores are from clusters with at least 80% similarity and so on.

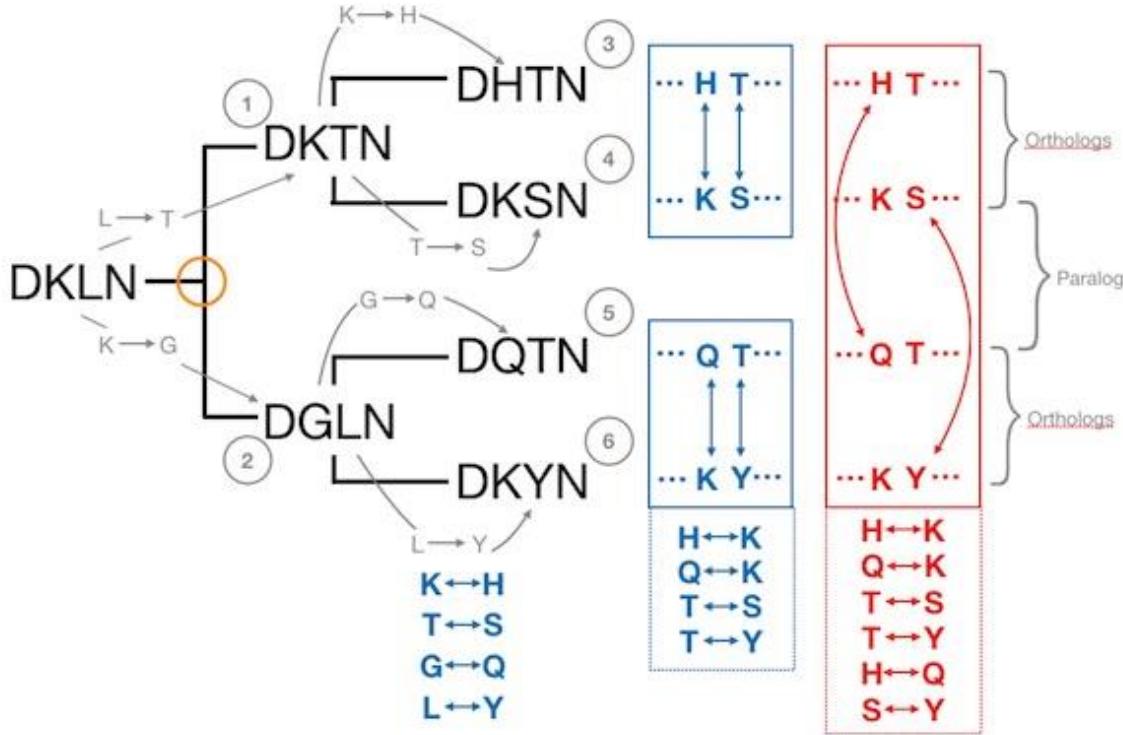
# BLOSUM62

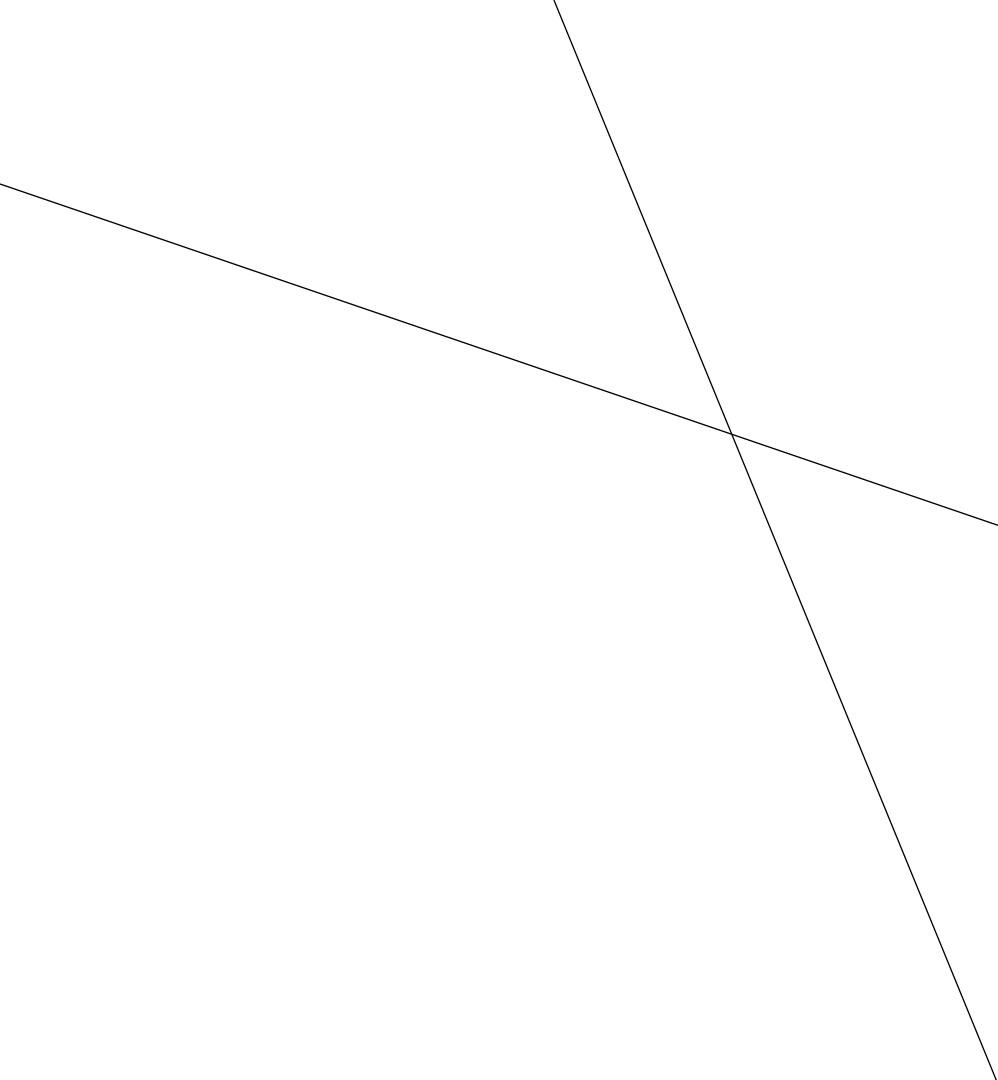




# PAM Matrices

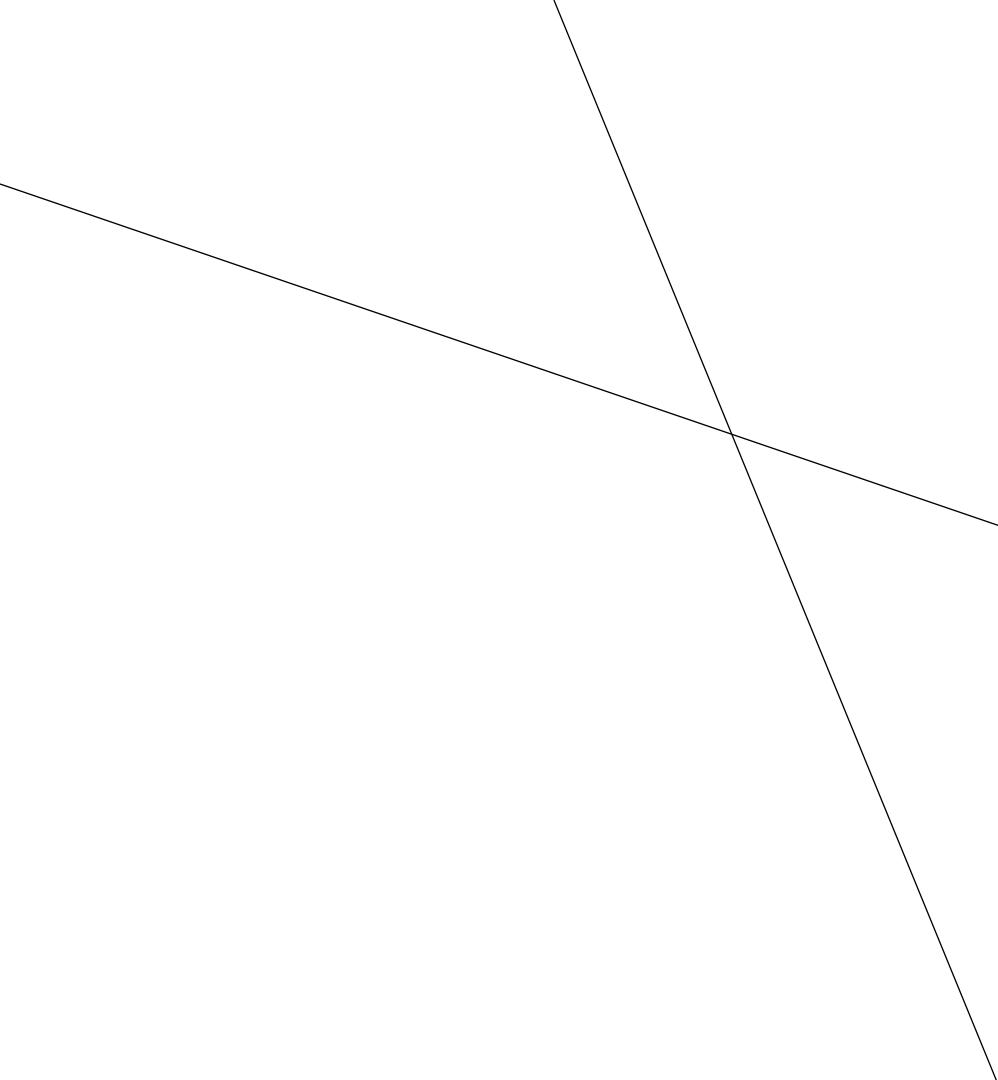
Accepted Point Mutation



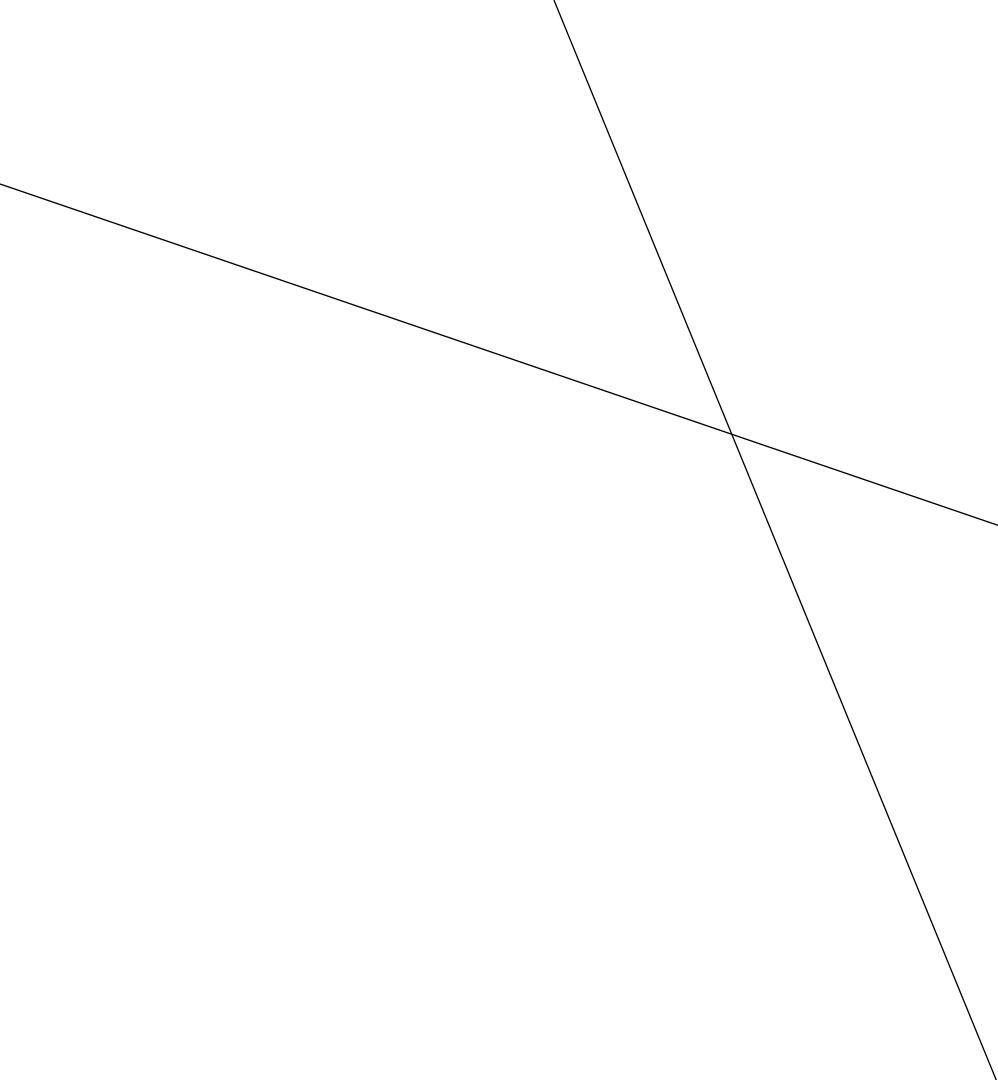


More information:

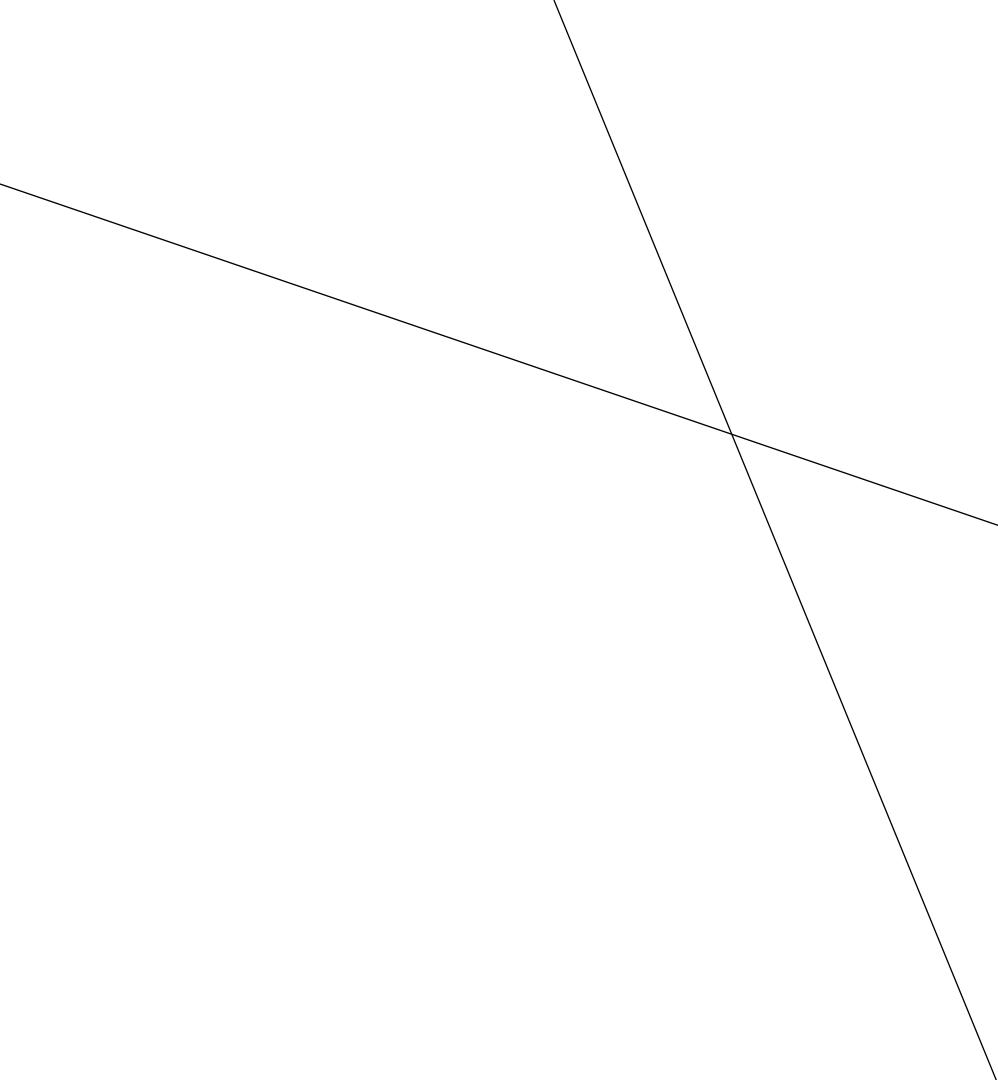
[https://bioinformaticshome.com/bioinformatics\\_tutorials/sequence\\_alignment/substitution\\_matrices.html](https://bioinformaticshome.com/bioinformatics_tutorials/sequence_alignment/substitution_matrices.html)



How to select the right  
substitution matrix?



## Statistical Significance of Alignment

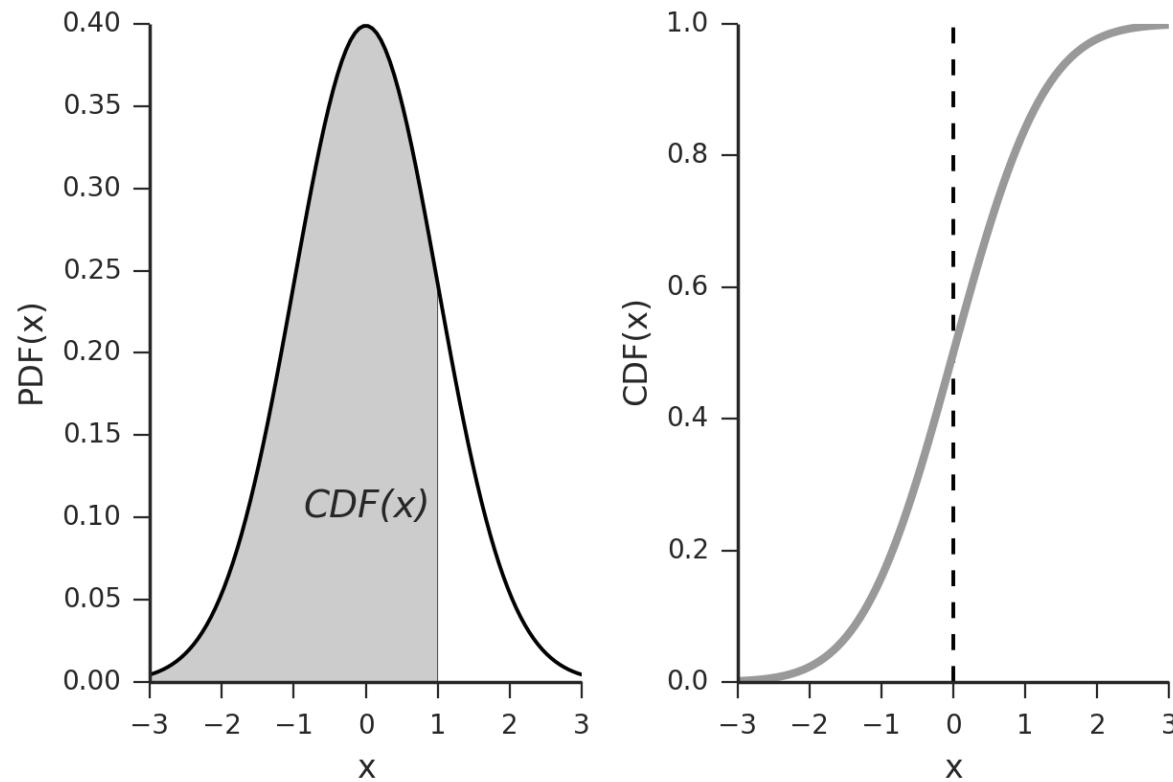


# Database Searching

## BLAST

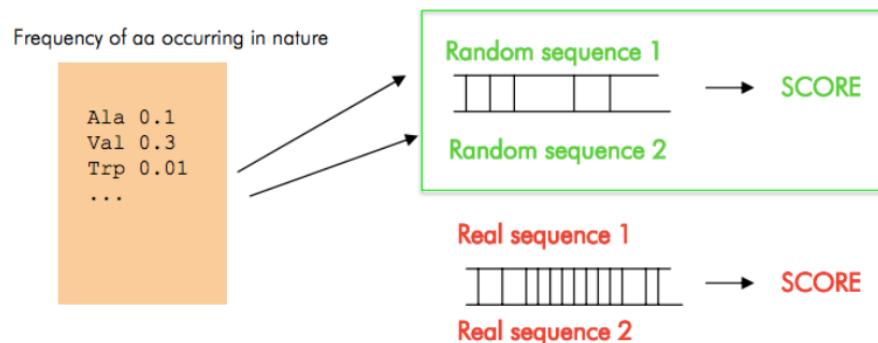
Given two random sequences of lengths m and n:  
What is the probability that they will produce an  
MSP score of  $\geq S$  ?

Given two random sequences of lengths m and n:  
What is the probability that they will produce an  
MSP score of  $\geq S$  ?

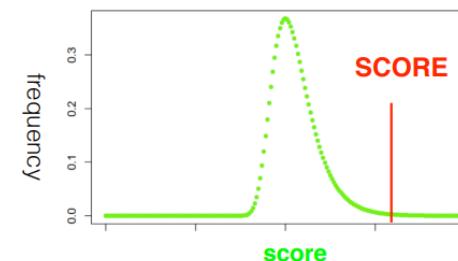
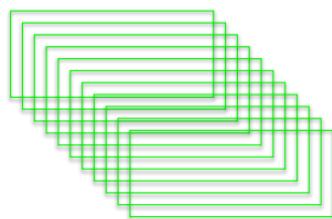


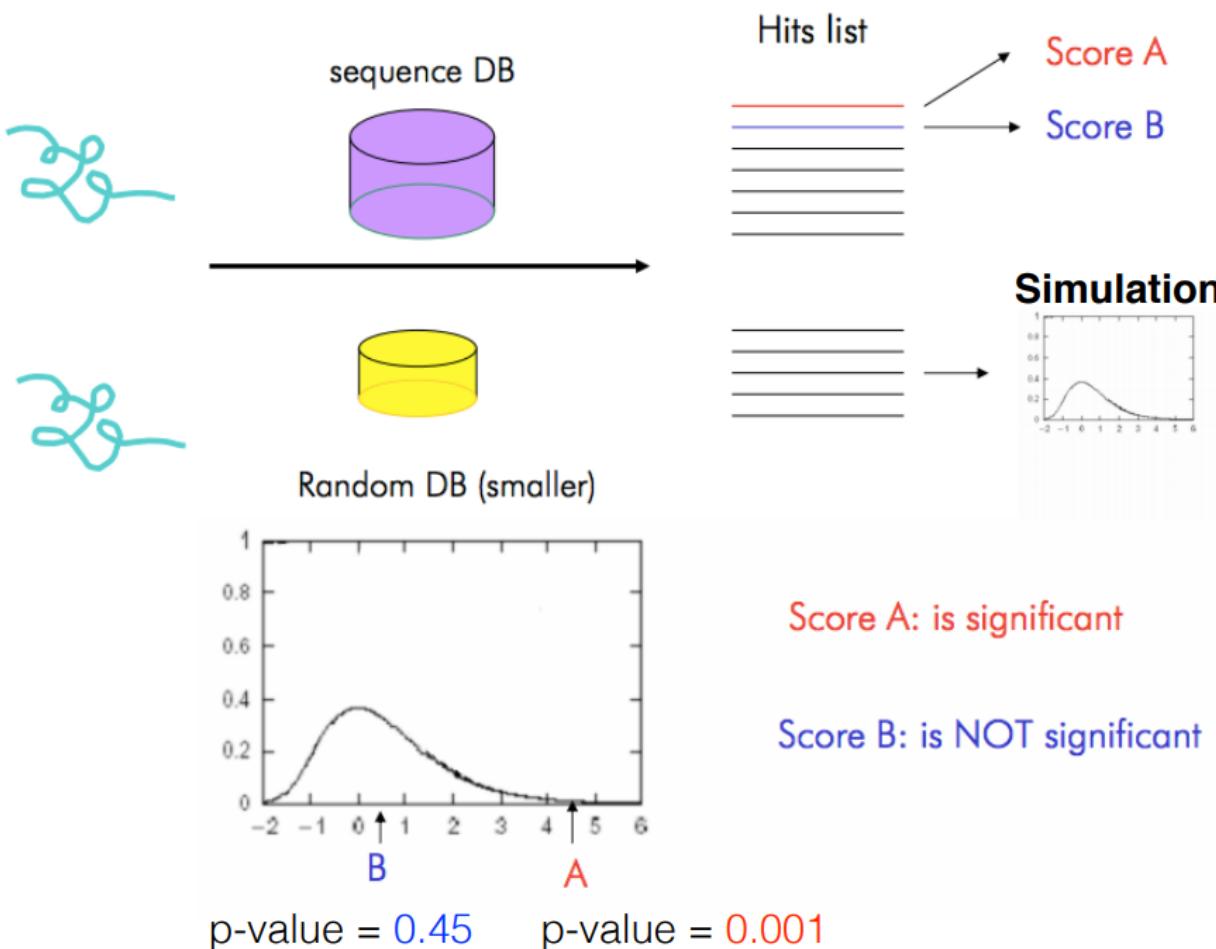
Given two random sequences of lengths m and n:  
What is the probability that they will produce an  
MSP score of  $\geq S$  ?

### 1. Generate many random sequence pairs

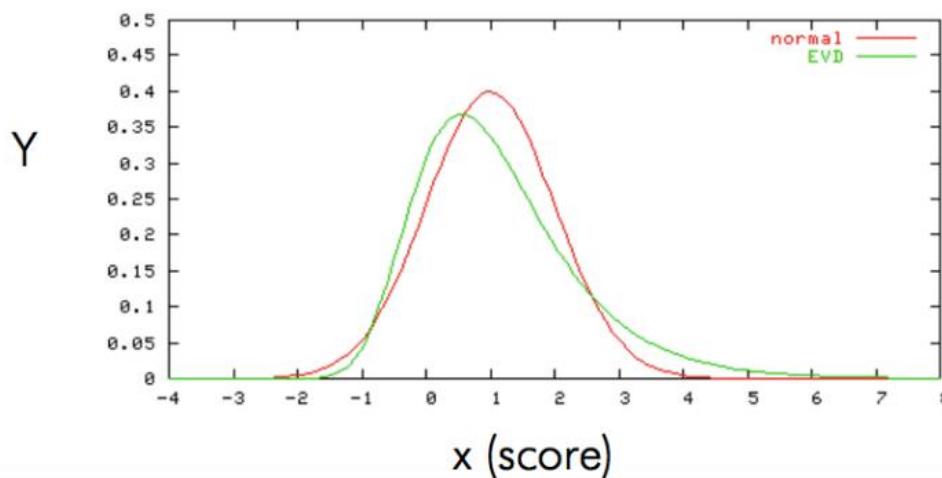


### 2. Compute the distribution of the SCOREs

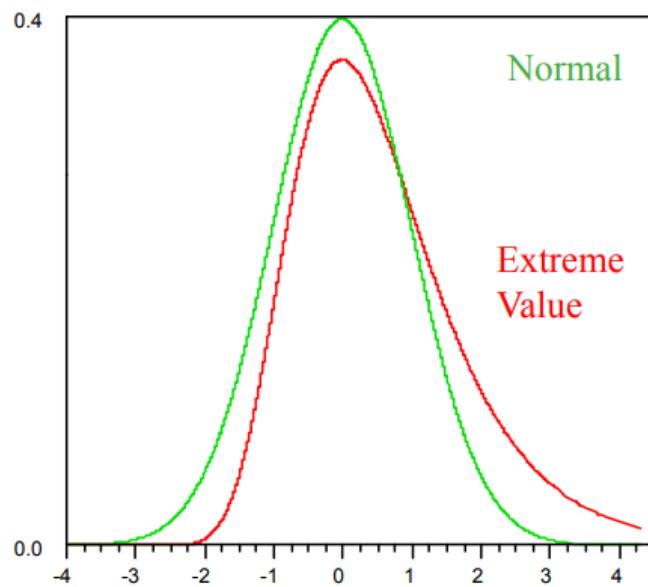




Karlin and Altschul observed that in the framework of local alignments without gaps: the distribution of random sequence alignment scores follow an EVD.



## Normal vs. Extreme Value Distribution



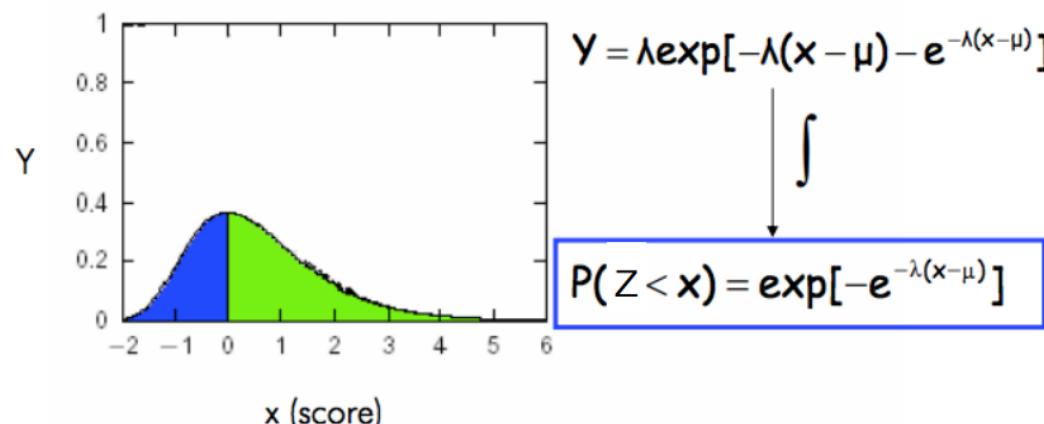
Normal distribution:

$$y = (1/\sqrt{2\pi})e^{-x^2/2}$$

Extreme value distribution:

$$y = e^{-x} - e^{-x}$$

# Extreme value distribution



$$P(Z \geq x) = 1 - \exp[-e^{-\lambda(x-\mu)}]$$

**P-value** = the probability of obtaining a score equal or greater than  $x$  by chance

In a database of size N:  $P \times N = E$

- **P-value:**

Probability that an alignment with this score occurs by chance in a database of size N.

The closer the P-value is towards 0, the better the alignment

- **E-value:**

Number of matches with this score one can expect to find by chance in a database of size N.

The closer the E-value is towards 0, the better the alignment

→ Smaller E-value, more significant in statistics

Bigger E-value , by chance

# Parameters

$$P(Z \geq x) = 1 - \exp[-e^{-\lambda(x-\mu)}] \quad (1)$$

$\mu, \lambda$  : parameters depend on the length and composition of the sequences and on the scoring system:  $\mu$  is the mode (highest point) of the distribution and  $\lambda$  is the decay parameter  
-They can be estimated by making many alignments of random or shuffled sequences.  
- For alignments without gaps they can be calculated from the scoring matrix and then :

$$P(Z \geq x) = 1 - \exp[-Kmne^{-\lambda x}] \quad (2)$$

$K$ : is a constant that depends on the scoring matrix values and the frequencies of the different residues in the sequences.

$m, n$  : sequence lengths

Approximation:

if  $x$  is very small, then  $1 - \exp(-x)$  can be approximated by  $x$

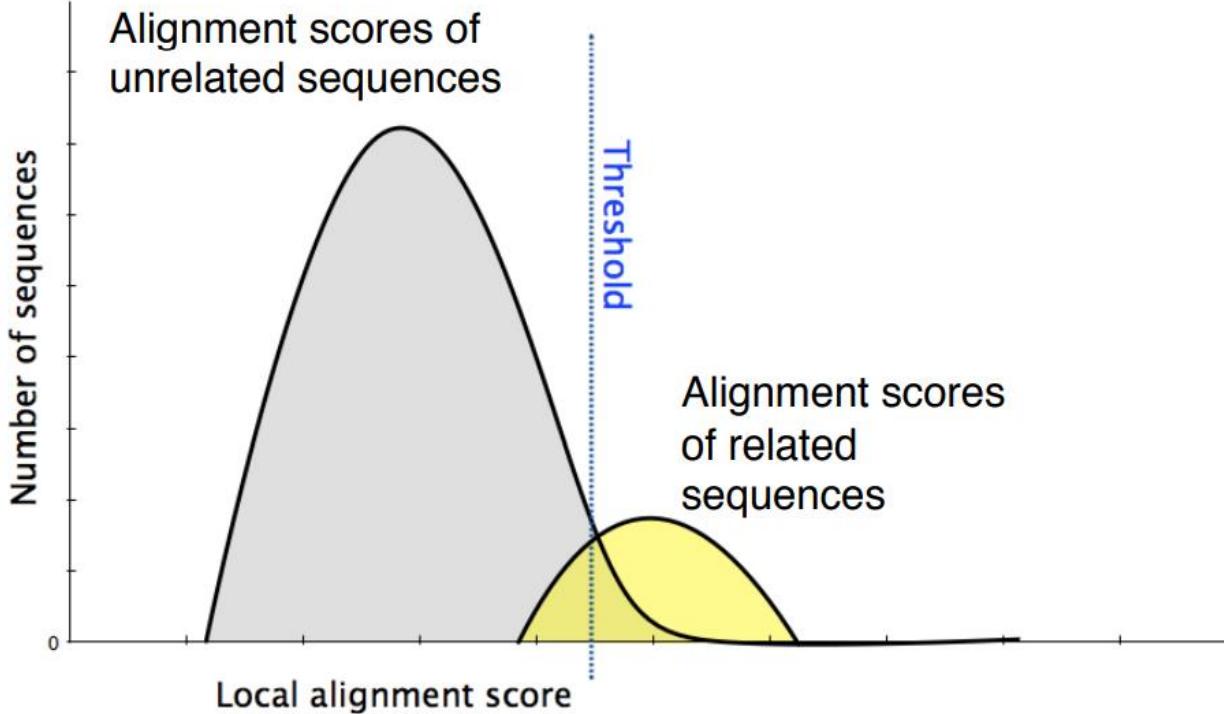
Therefore,

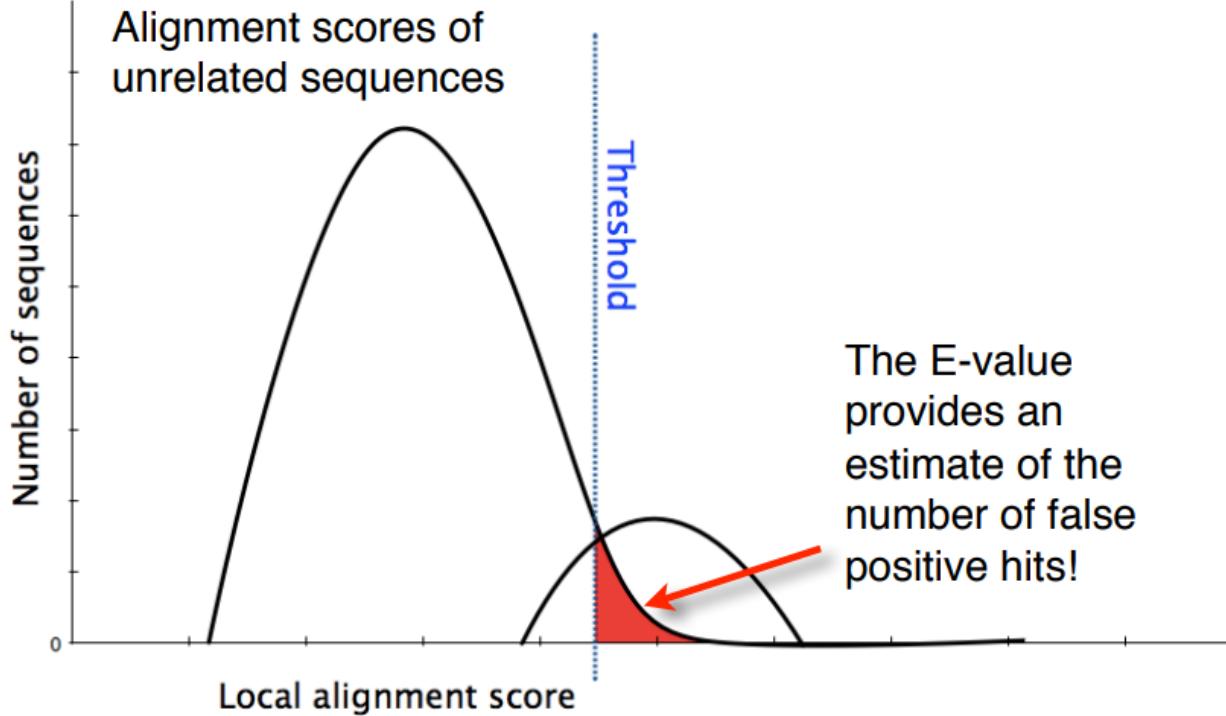
$$P(Z \geq x) \sim e^{-\lambda(x-\mu)} = K_{mn} e^{-\lambda x}$$

So E-value = DatabaseLength \* p-value

$$\text{E-value} = KNm e^{-\lambda x}$$

where  $N$  is the database size (not the aligned length  $n$ )

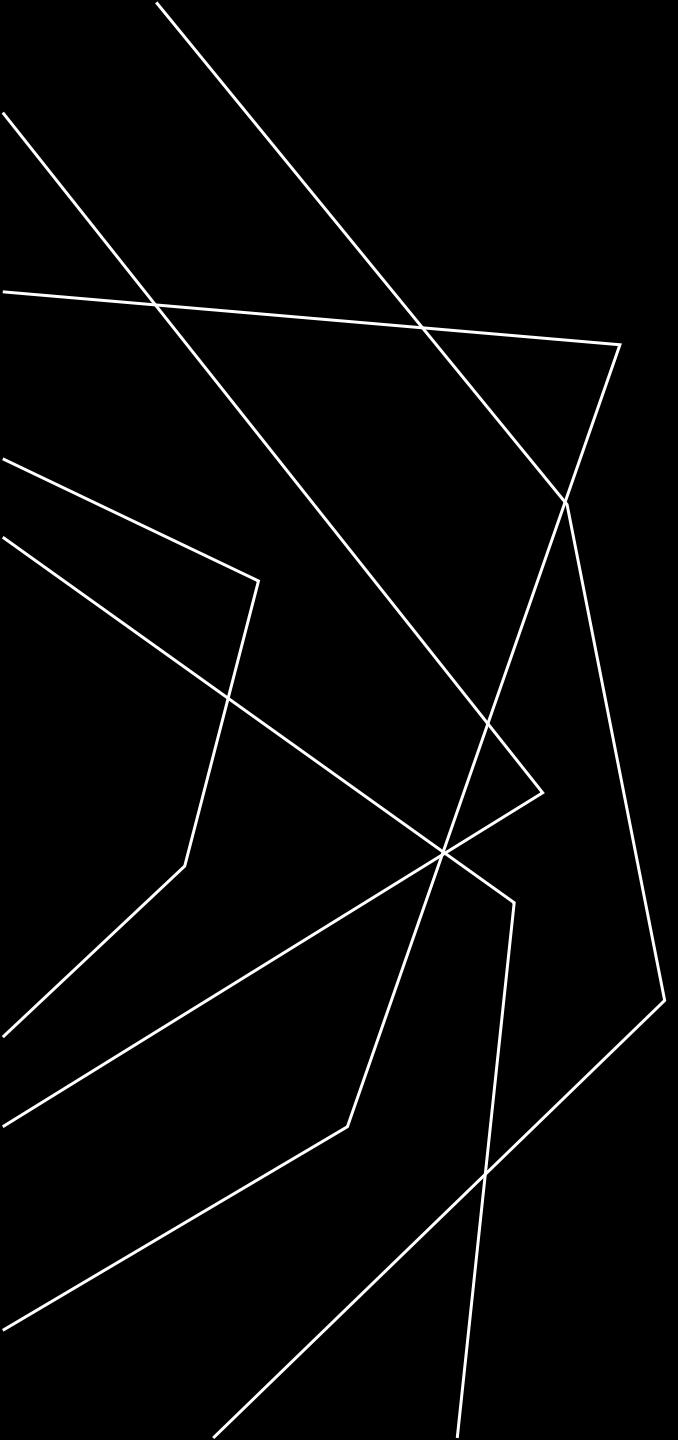




The E value is the expected number of hits that are as good or better than the observed local alignment score (with this score or better) if the query and database are random with respect to each other.

$$E(S \geq x) = Kmne^{-\lambda x}$$

- $K < 1$  is a proportionality constant that corrects the  $mn$  “space factor” for the fact that there are not really  $mn$  independent places that could have produced score  $S \geq x$ .
- $K$  has little effect on the statistical significance of a similarity score
- $\lambda$  is closely related to the scoring matrix used and it takes into account that the scoring matrices do not contain actual probabilities of co-occurrence, but instead a scaled version of those values. To understand how  $\lambda$  is computed, we have to look at the construction of scoring matrices.



THANK YOU  
Q&A