

# BIOINFORMATICS 1

or why biologists need computers

<http://www.bioinformatics.uni-muenster.de/teaching/courses-2011/bioinf1/index.hbi>





# INTRODUCTION TO SEQUENCE ANALYSIS

dot plots, alignments, and similarity searches





# EVOLUTIONARY BASIS OF SEQUENCE ANALYSES

THIS IS AN ANCESTRAL SEQUENCE



# EVOLUTIONARY BASIS OF SEQUENCE ANALYSES

THIS IS AN ANCESTRAL SEQUENCE  
THIS IS AN **M** NCESTRAL SEQUENCE



# EVOLUTIONARY BASIS OF SEQUENCE ANALYSES

THIS IS AN ANCESTRAL SEQUENCE  
THIS IS AN **M** NCESTRAL SEQUENCE  
THIS IS AN **M** NCESTRAL **W** SEQUENCE



# EVOLUTIONARY BASIS OF SEQUENCE ANALYSES

THIS IS AN ANCESTRAL SEQUENCE  
THIS IS AN **M** NCESTRAL SEQUENCE  
THIS IS AN **M** NCESTRAL **W** SEQUENCE  
THIS IS AN **MP** CESTRAL **W** SEQUENCE



# EVOLUTIONARY BASIS OF SEQUENCE ANALYSES

THIS IS AN ANCESTRAL SEQUENCE  
THIS IS AN **M** NCESTRAL SEQUENCE  
THIS IS AN **M** NCESTRAL **W** SEQUENCE  
THIS IS AN **MP** CESTRAL **W** SEQUENCE  
THIS IS **CNMP** ESTRAL **W** SEQUENCE

Please note deletion of “C”





# EVOLUTIONARY BASIS OF SEQUENCE ANALYSES

THIS IS **CNMP** ESTRAW **W** SEQUENCE

Gene duplication or speciation!

THIS IS **CNMP** ESTRAW **W** SEQUENCE



# EVOLUTIONARY BASIS OF SEQUENCE ANALYSES

THIS IS **CNMP** ESTRA **W** SEQUENCE  
THIS IS **C** **OMP** **E** TRA **W** SEQUENCE

THIS IS **CNMP** ESTRA **W** SEQUENCE  
THIS IS **NMP** **ER** **SX** TRA SEQUENCE

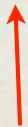
Please note deletion of “C” and “W”  
compensated by insertion of “R” and “X”



# EVOLUTIONARY BASIS OF SEQUENCE ANALYSES

THIS IS **COMP**ET**LA**W SEQUENCE

THIS IS **C****N****MP****E****X**TR A SEQUENCE



Please note insertion of "C"



# EVOLUTIONARY BASIS OF SEQUENCE ANALYSES

THIS IS **COMPLET** **NA** W **SEQUENCE**

THIS IS **CS** **MP** **E** **EX** **TR** **A** **SEQUENCE**



# EVOLUTIONARY BASIS OF SEQUENCE ANALYSES

THIS IS **COMPLET** **NA** W **SEQUENCE**

THIS IS **CSUP** **EX** **TR** **SEQUENCE**



# EVOLUTIONARY BASIS OF SEQUENCE ANALYSES

THIS IS **COMPLETELY NEW** SEQUENCE

THIS IS **CSUPEREXTRA** SEQUENCE



# EVOLUTIONARY BASIS OF SEQUENCE ANALYSES

THIS IS COMPLETELY NEW SEQUENCE

THIS IS SUPEREXTRA SEQUENCE

Please note another deletion of “C” and insertion of “R”



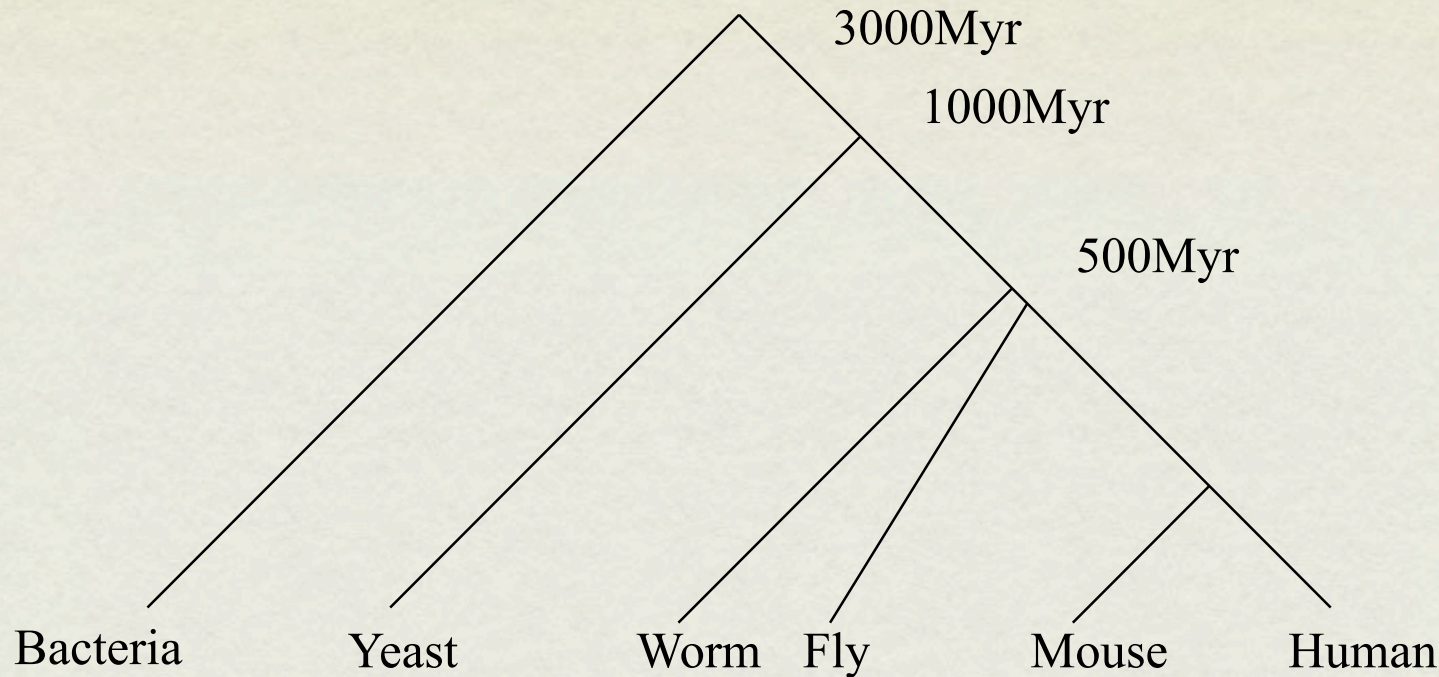


# EVOLUTIONARY BASIS OF SEQUENCE ANALYSES

THIS IS COMPLETELY NEW SEQUENCE  
THIS IS SUPEREXTRA SEQUENCE



# HUMAN COLON CANCER GENE AND BACTERIAL DNA REPAIR GENE

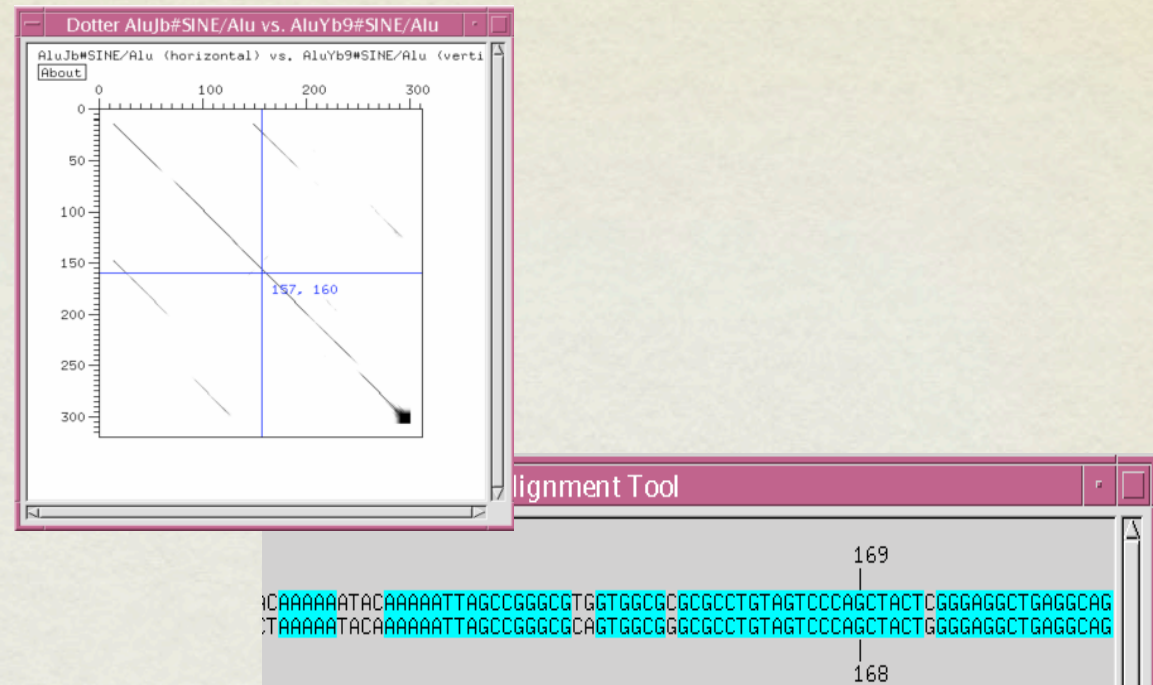


MSH2_Human	TGVIVLMAQIGCFVPCESAEVSI	VDCILARVGAGDSQLKGVSTFMAEMLETASILRSATK	
SPE1_DROME	VGTA	VLMAHIGAFVPCSLATISMVDSILGRVGASDNI	IKGLSTFMVEMIETSGIIRTATD
MSH2_Yeast	VGVISLMAQIGCFVPC	EEAEIAIVDAILCRVGAGDSQLKGVSTFMVEILETASILKNASK	
MUTS_ECOLI	TALIALMAYIGSYVPAQKVEIGPIDRIFTRVGAADDLASGRSTFMVEMTETANILRNATE		
	*** ** *	* * ****	***** * ** *



# MAJOR TECHNIQUES TO BE DISCUSSED

- Dot Matrix plots
- Sequence alignments
- Similarity searches





# HOW TO SOLVE THE PROBLEM - HUMAN OR COMPUTER?



- 🌀• very smart
- 🌀• slow
- 🌀• error prone
- 🌀• doesn't like repetitive tasks

not as cute looking as  
humans

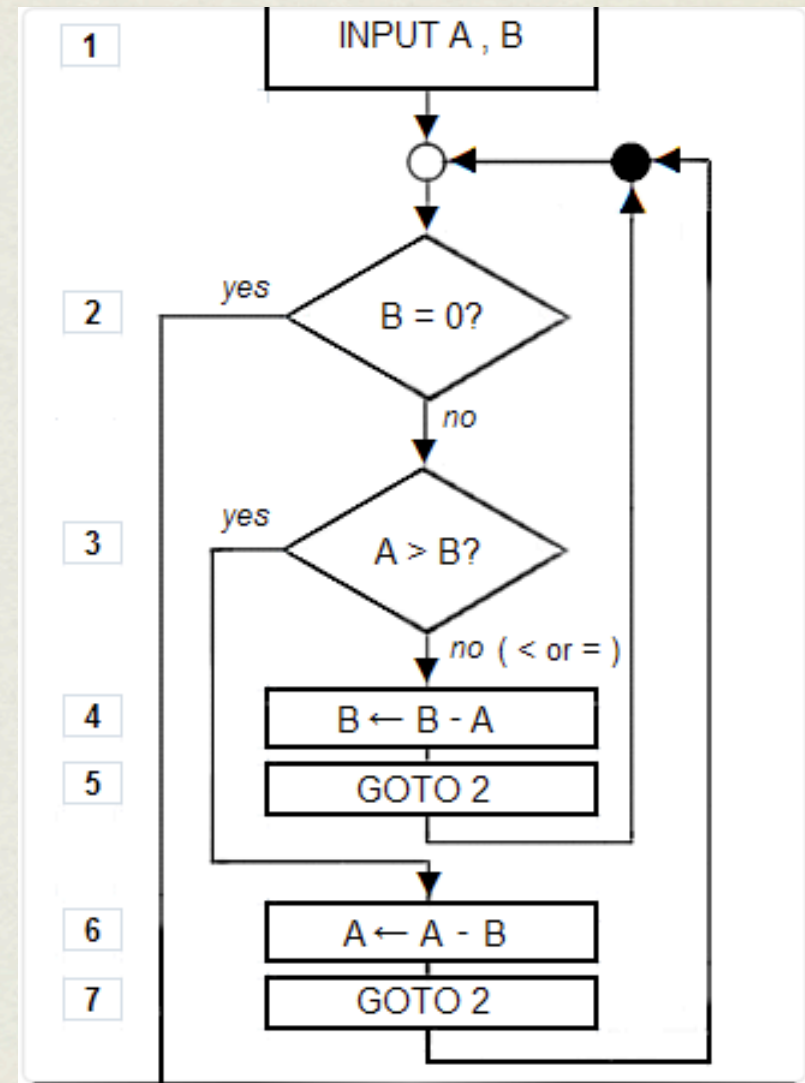
- 🌀• not so smart (stupid)
- 🌀• extremely fast
- 🌀• very accurate
- 🌀• doesn't understand human languages;  
needs instruction provided in a special way





# ALGORITHM

A step-by-step problem-solving procedure, especially an established, recursive computational procedure for solving a problem in a finite number of steps.



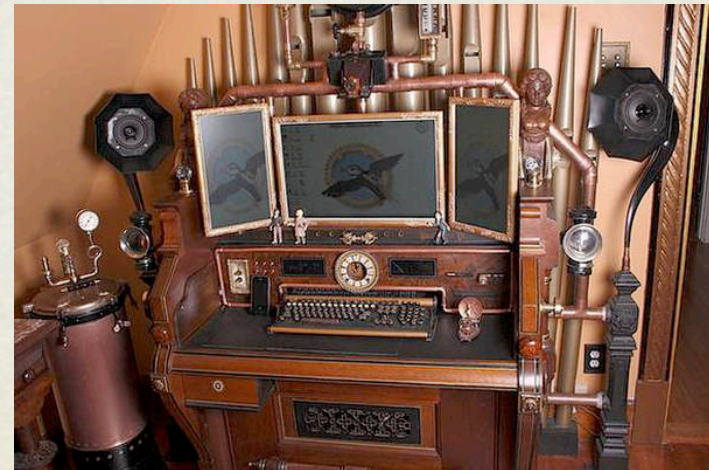


# EXAMPLE TASK: PUT SHOES ON!



A human just understands an order  
and often executes it automatically  
even without thinking

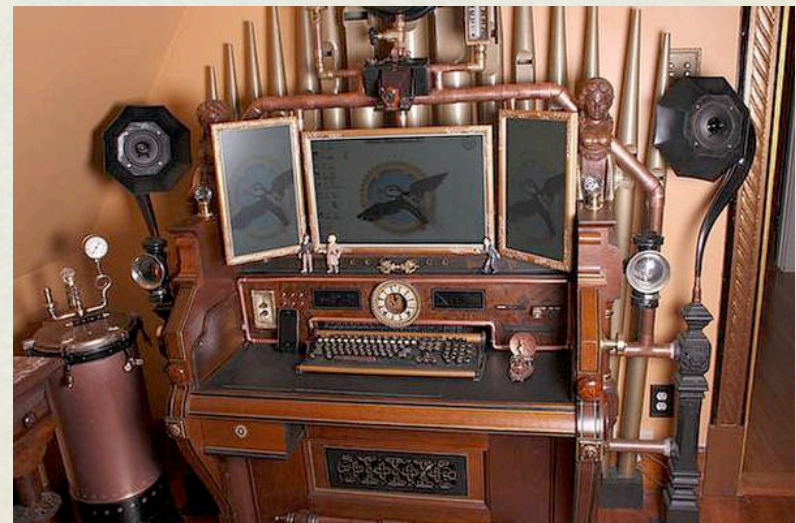
A computer needs detailed  
instruction (an algorithm)





# PUT SHOES ON! INSTRUCTION FOR A COMPUTER

1. Find two the same shoes
2. Check if you have left and right shoe
3. Check if their are of the same size
4. Check if this is the right size
5. Put the left shoe on
6. Put the right shoe on
7. Tie the laces







# Dot Matrix Plots



# DOT MATRIX PLOTS

- ✂• Sensitive qualitative indicators of similarity
- ✂• Better than alignments in some ways
  - ✂• rearrangements
  - ✂• repeated sequences
- ✂• Rely on visual perception (not quantitative)
- ✂• Useful for RNA structure



# DOT MATRIX PLOTS

- Simplest method - put a dot wherever sequences are identical
- A little better - use a scoring table, put a dot wherever the residues have better than a certain score (especially useful for amino acid sequence comparison)
- Or, put a dot wherever you get at least  $n$  matches in a row (identity matching, compare/word)
- Even better - filter the plot



# WINDOWED SCORES ALGORITHM

1. calculate a score within a window of a given size, for example six
2. plot a point if score is over a threshold (stringency), for example 70%
3. move the window over a given step, for example one
4. repeat step one to three till the end of sequence



# WINDOWED SCORES EXAMPLE

Let's compare two nucleotide sequences

ACCTTGTCCTCTTTGCCC

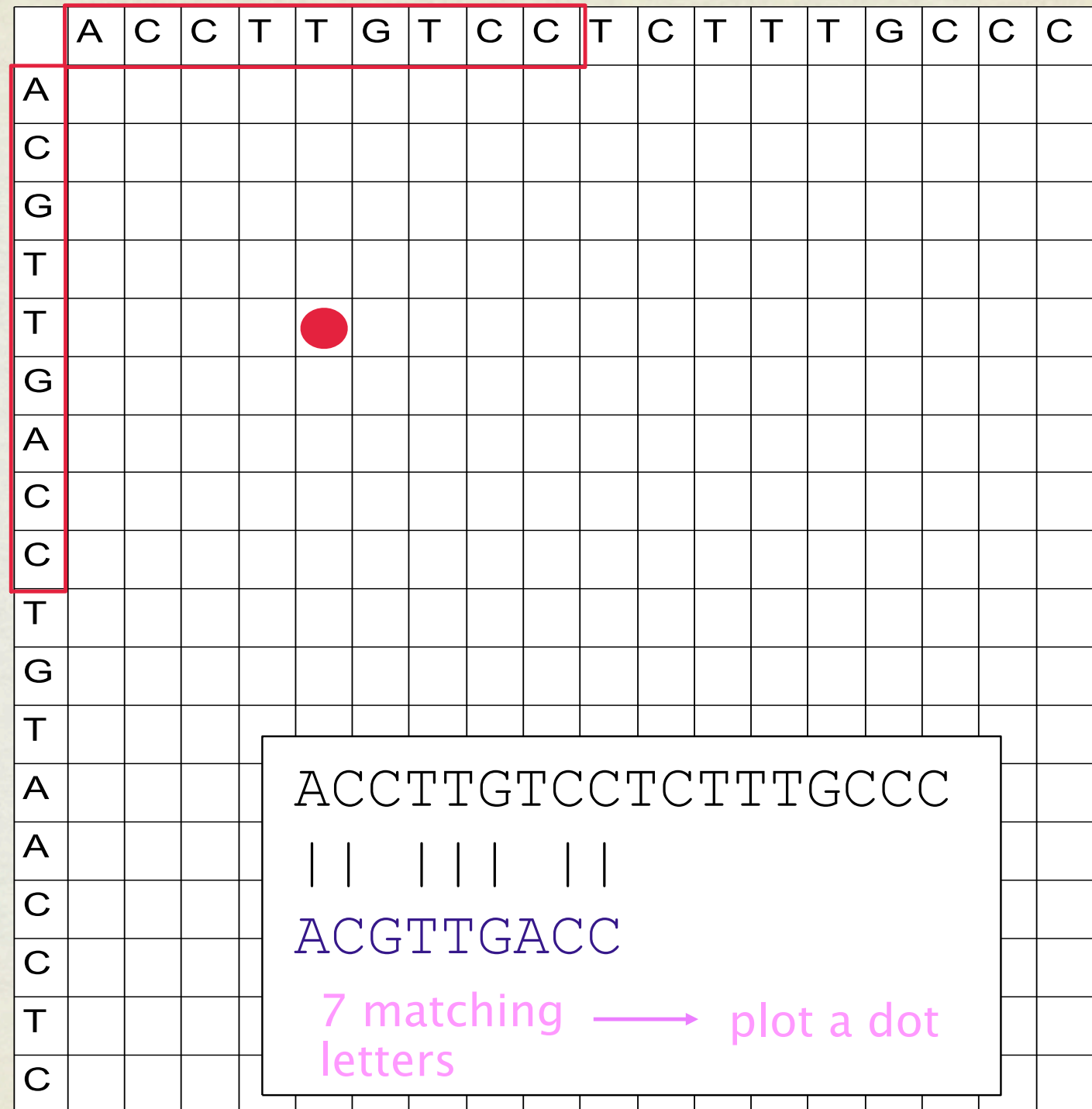
ACGTTGACCTGTAACTC

using following parameters:

window size = 9, step = 3, threshold = 4



window size = 9  
step = 3  
threshold = 4





window size = 9  
step = 3  
threshold = 4

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

ACCTTGTCCTCTTTGCC

|||

ACGTTGACC

3 matching → no action letters



window size = 9  
step = 3  
threshold = 4

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

ACCTTGTCCTCTTTGCC

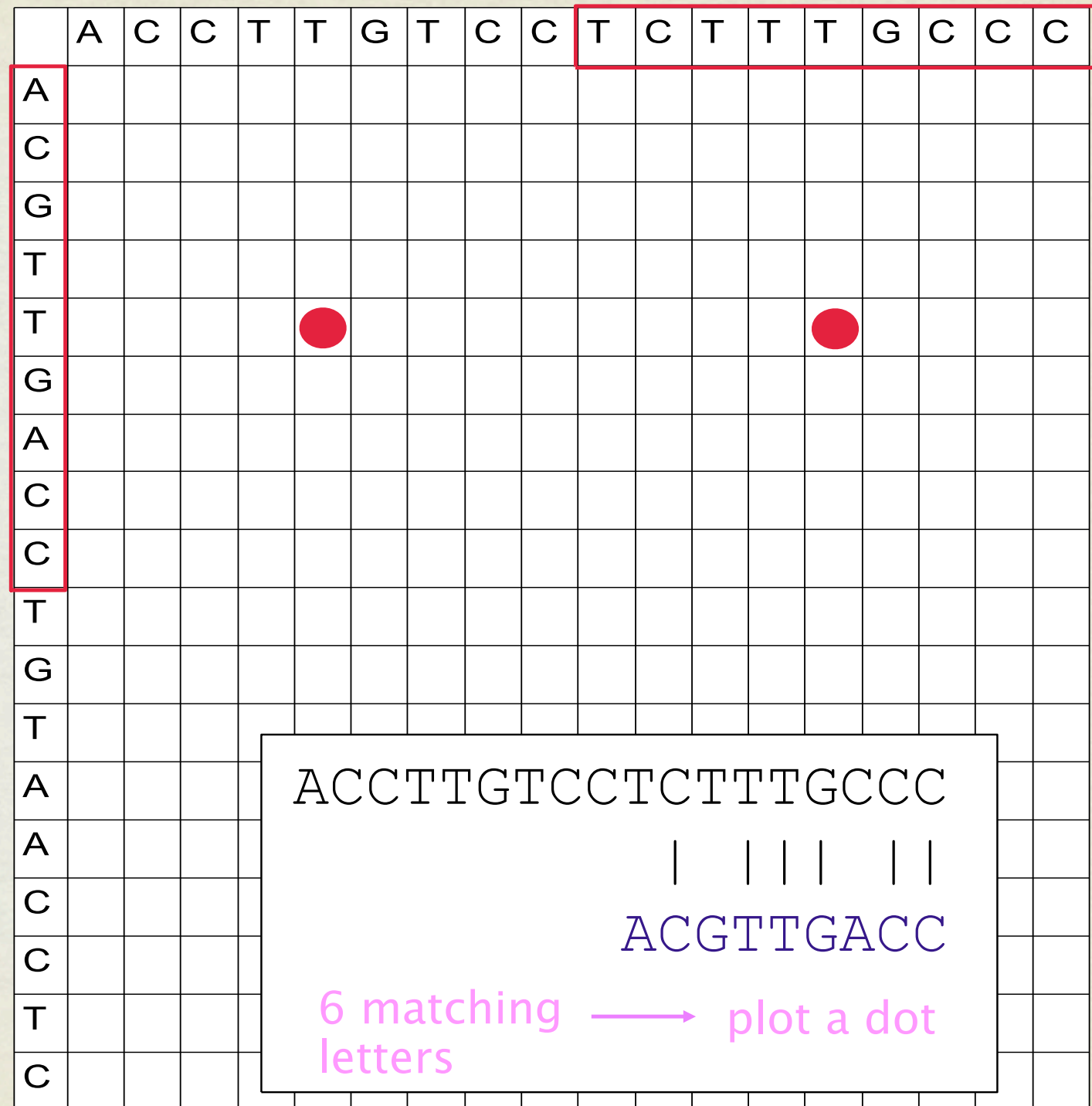
| |

ACGTTGACC

2 matching letters → no action



window size = 9  
step = 3  
threshold = 4





window size = 9  
step = 3  
threshold = 4

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

ACCTTGTCCTCTTTGCC

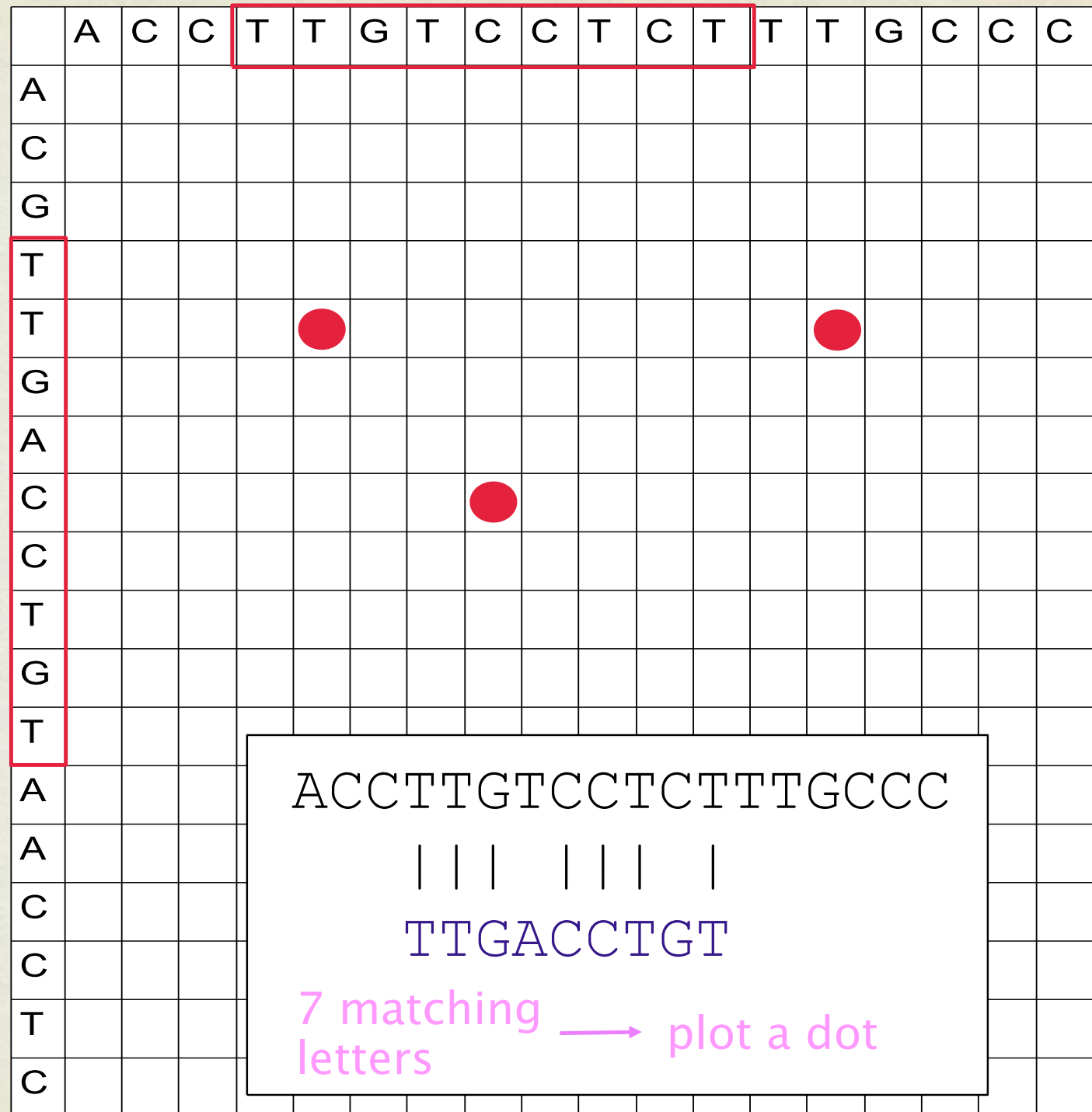
|

TTGACCTGT

1 matching letter → no action



window size = 9  
step = 3  
threshold = 4





window size = 9  
step = 3  
threshold = 4

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

ACCTTGTCCTCTTTGCC

| | |

TTGACCTGT

1 matching letters → no action



window size = 9  
step = 3  
threshold = 4

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

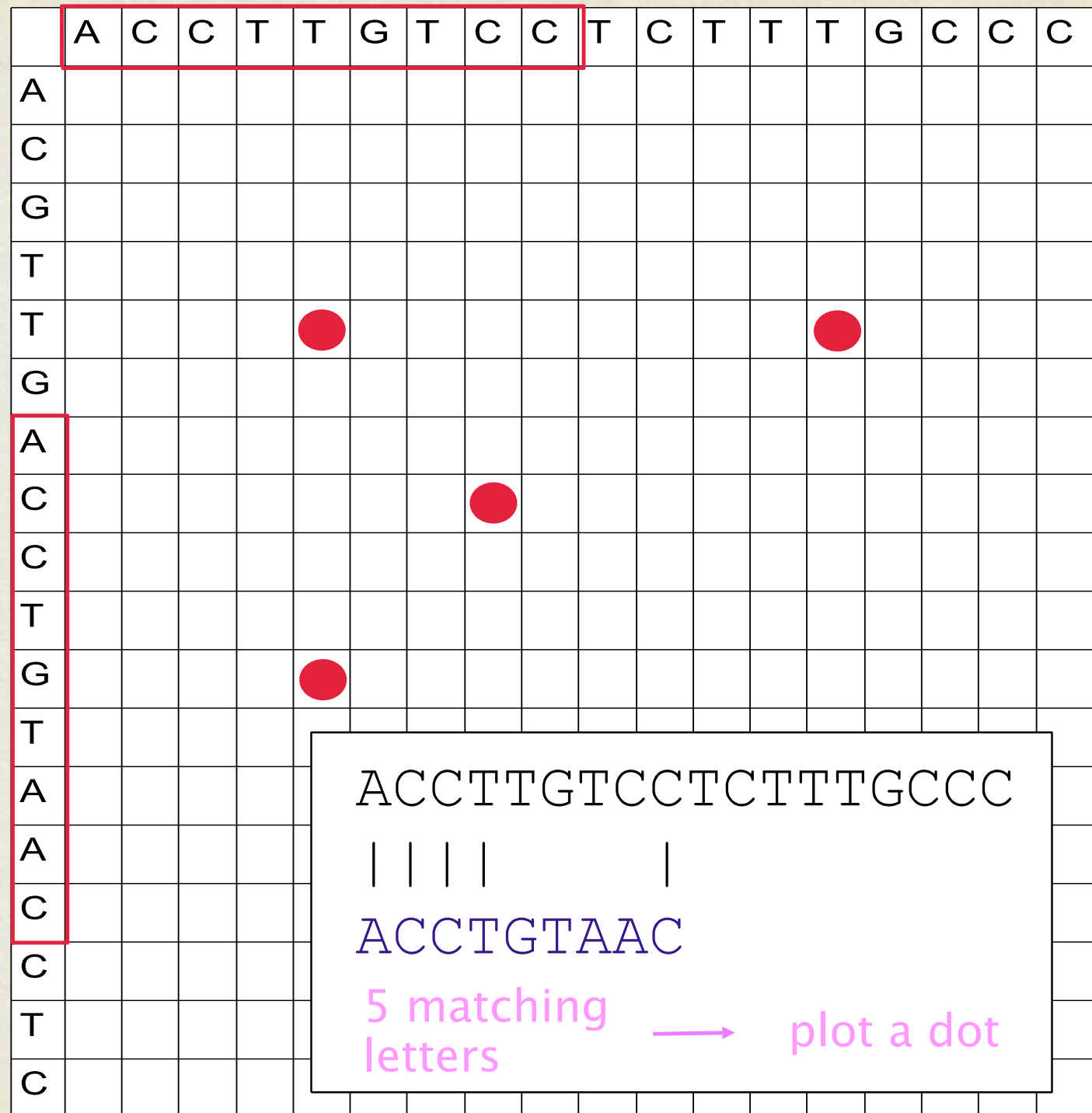
ACCTTGTCCTCTTTGCC

|

TTGACCTGT

1 matching letter → no action

window size = 9  
step = 3  
threshold = 4





window size = 9  
step = 3  
threshold = 4

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

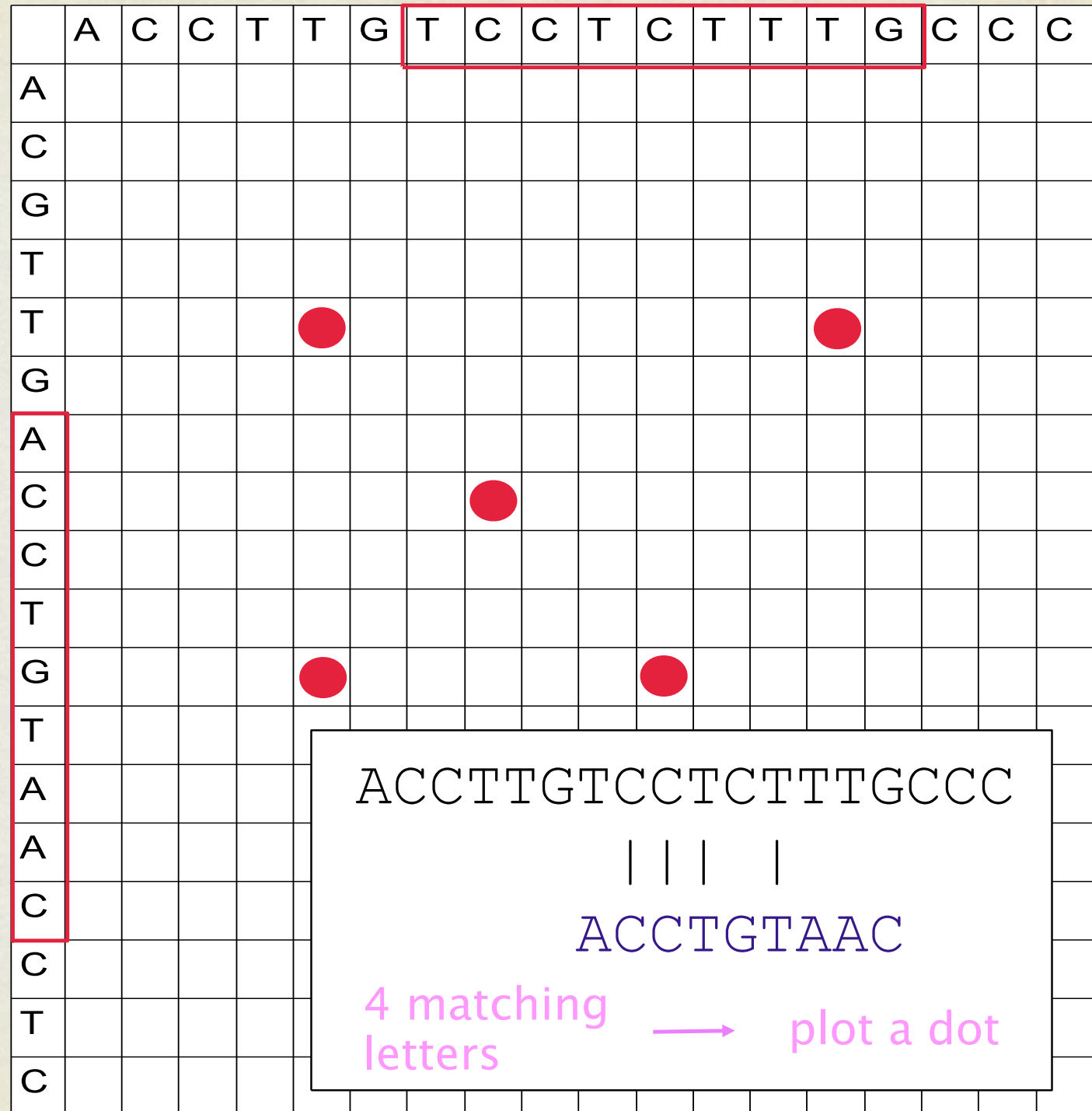
ACCTTGTCCTCTTTGCCC

|

ACCTGTAAC

1 matching letter → no action

window size = 9  
step = 3  
threshold = 4





window size = 9  
step = 3  
threshold = 4

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

ACCTTGTCCTCTTTGCCC

| | | |

ACCTGTAAC

3 matching letters → no action







Diagram illustrating a sequence alignment between two DNA sequences:

- Sequence 1 (Top): TCTTTTGCCCC
- Sequence 2 (Bottom): AACCTC

The alignment is shown on a grid. Red dots indicate matches (T-T, C-C, T-T, C-C, C-C). Green dots indicate mismatches (A-T, A-T, C-T, C-T, C-T). A red box highlights the first column (T vs A), and a green box highlights the last column (C vs C). A callout box points to the first column with the text "no action".

ACCTTGTCCTCTTGCCC

TGTAACCTC

2 matching letters → no action

```
window size = 9
      step = 3
      threshold = 4
```



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

ACCTTGTCCTCTTTGCCC

| | | | |

TGTAACCTC

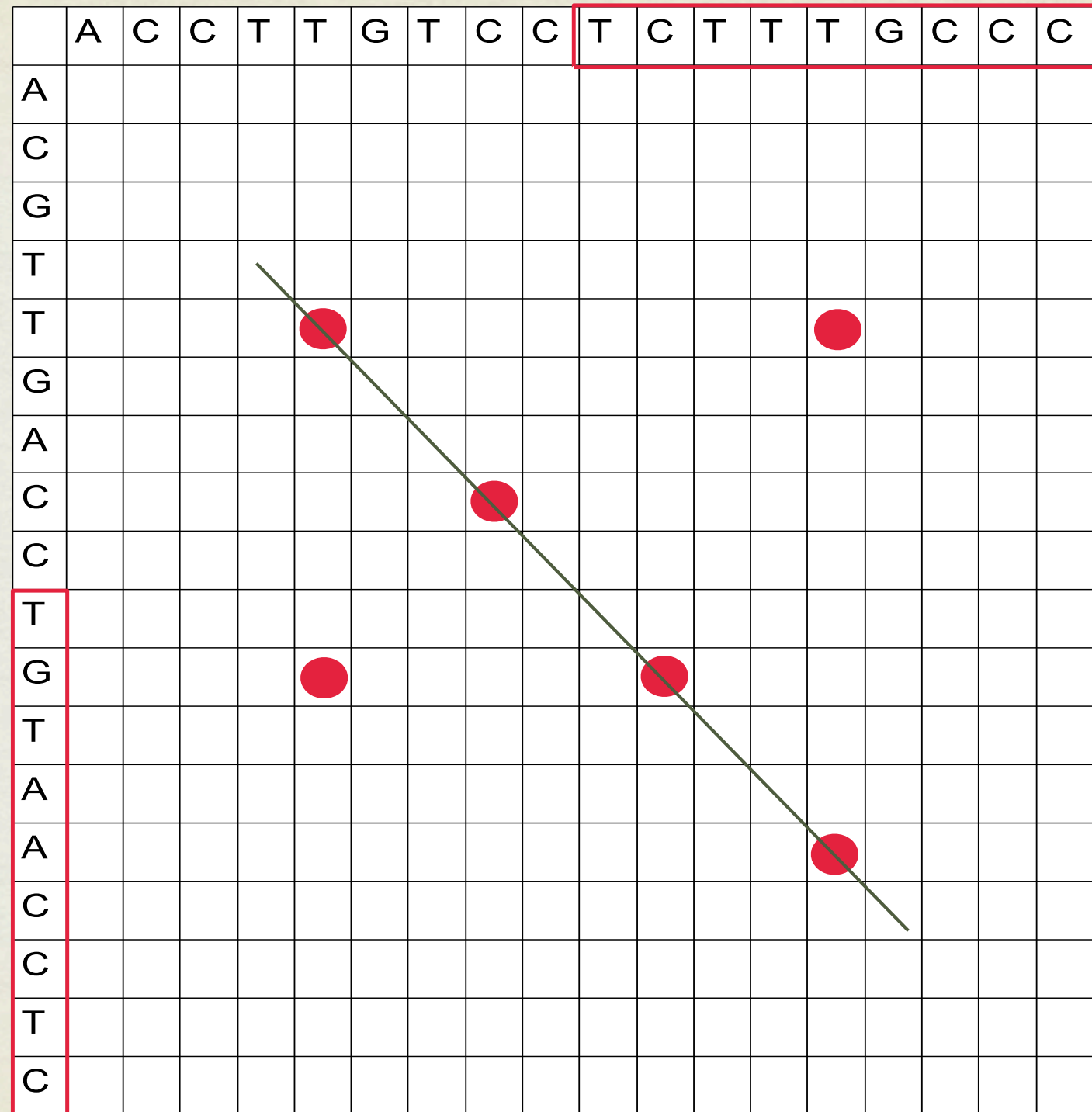
1 matching letters → plot a dot

window size = 9  
step = 3  
threshold = 4

[illegible]

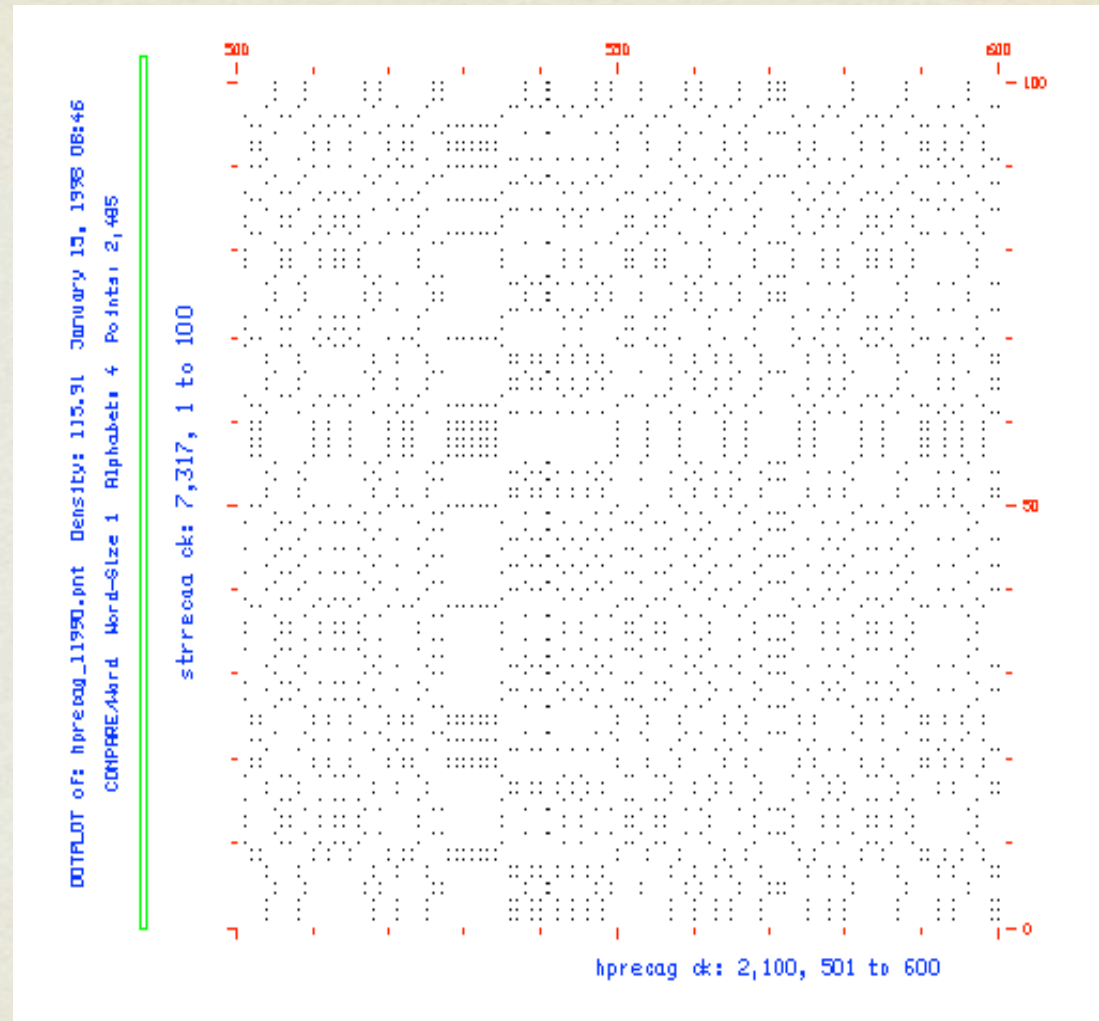
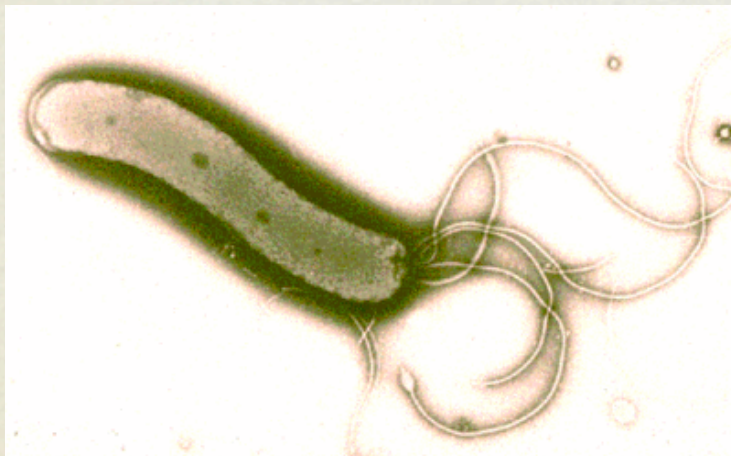


For similar  
sequences  
dots form  
a diagonal  
line



# DOT PLOT - EXAMPLES

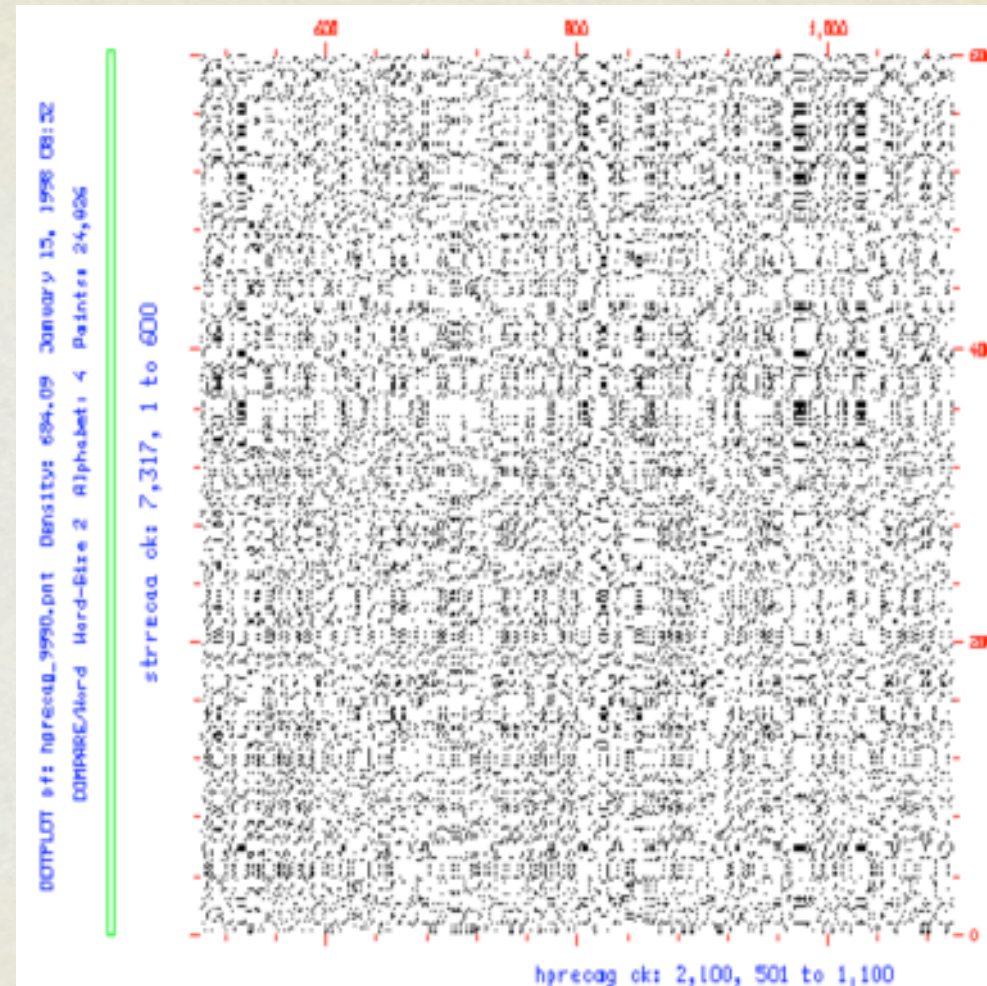
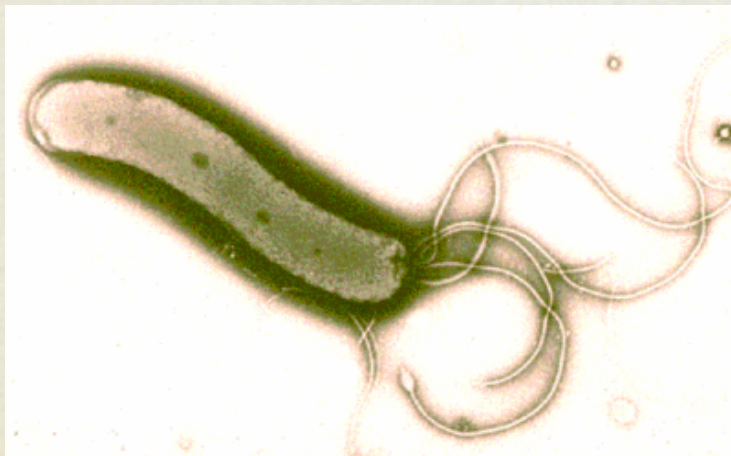
RecA DNA sequence  
from *Helicobacter pylori*  
and Streptococcus  
mutant, window=1  
match=1





# DOT PLOT - EXAMPLES

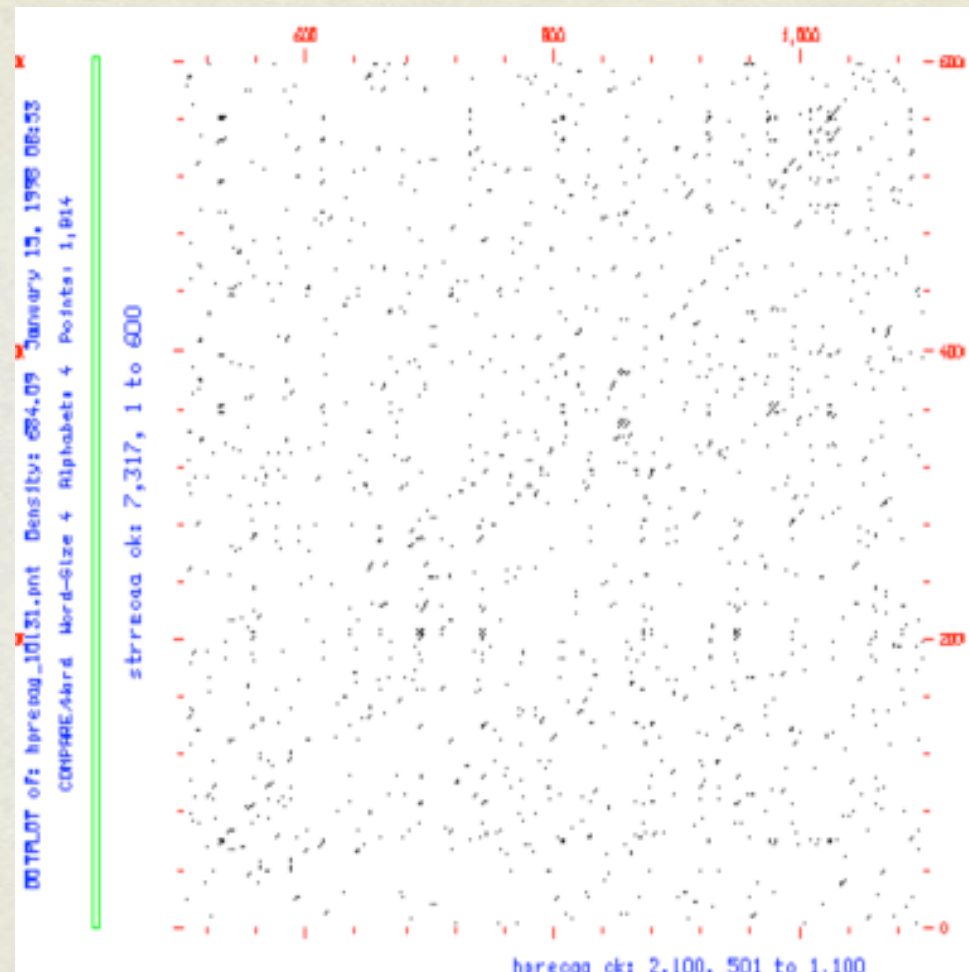
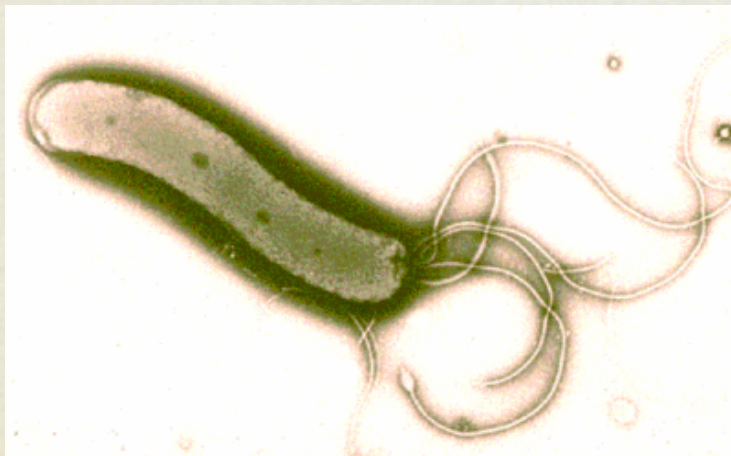
RecA DNA sequence  
from *Helicobacter pylori*  
and *Streptococcus*  
mutant, window=2  
match=2





# DOT PLOT - EXAMPLES

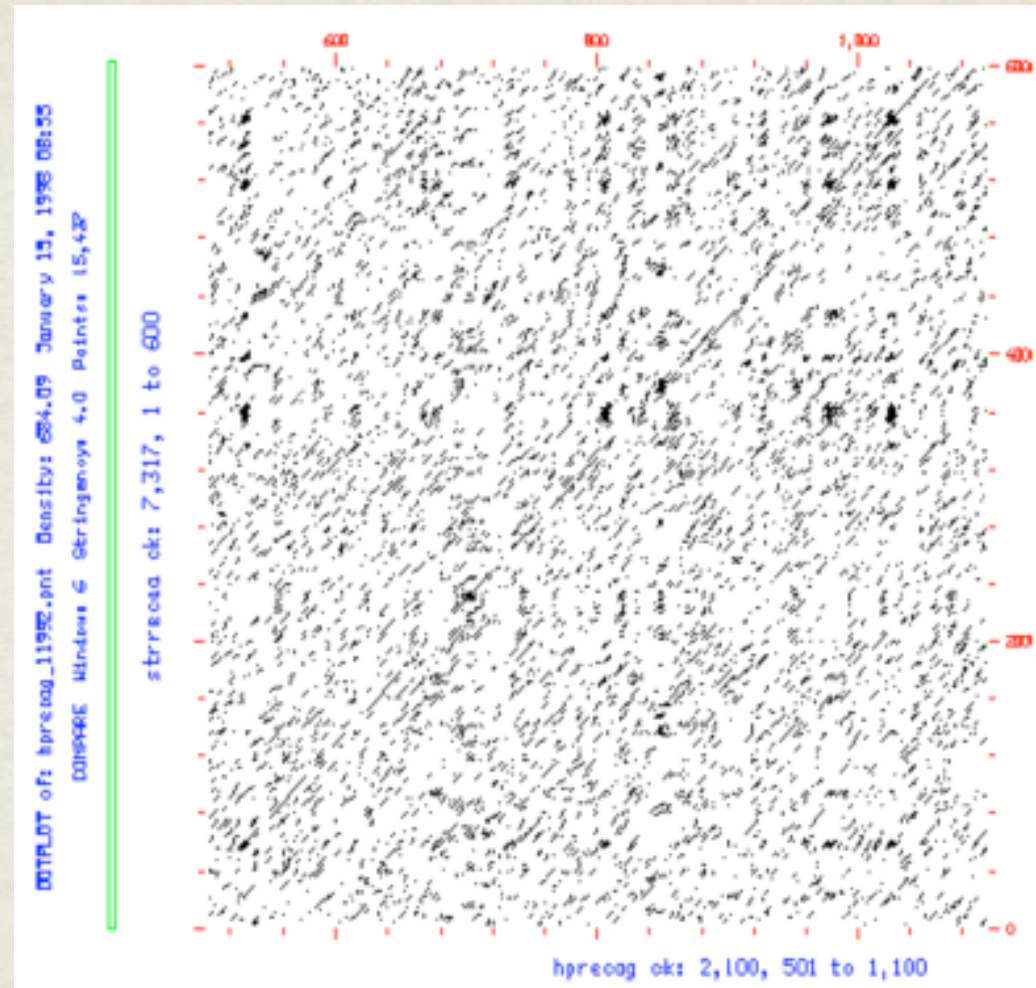
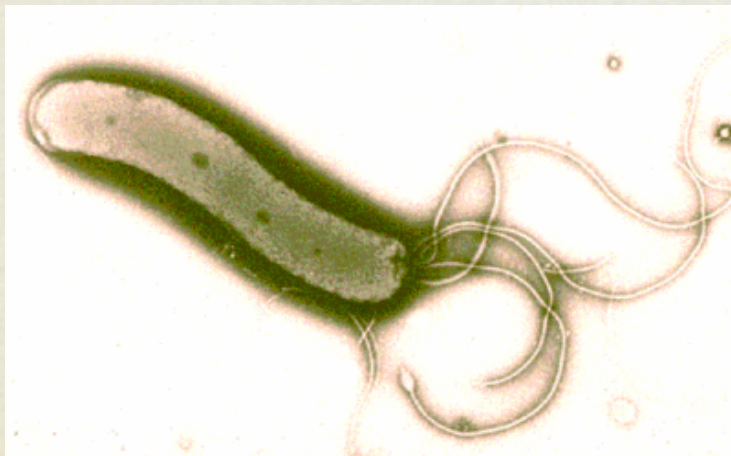
RecA DNA sequence  
from *Helicobacter pylori*  
and *Streptococcus*  
mutant, window=4  
match=4





# DOT PLOT - EXAMPLES

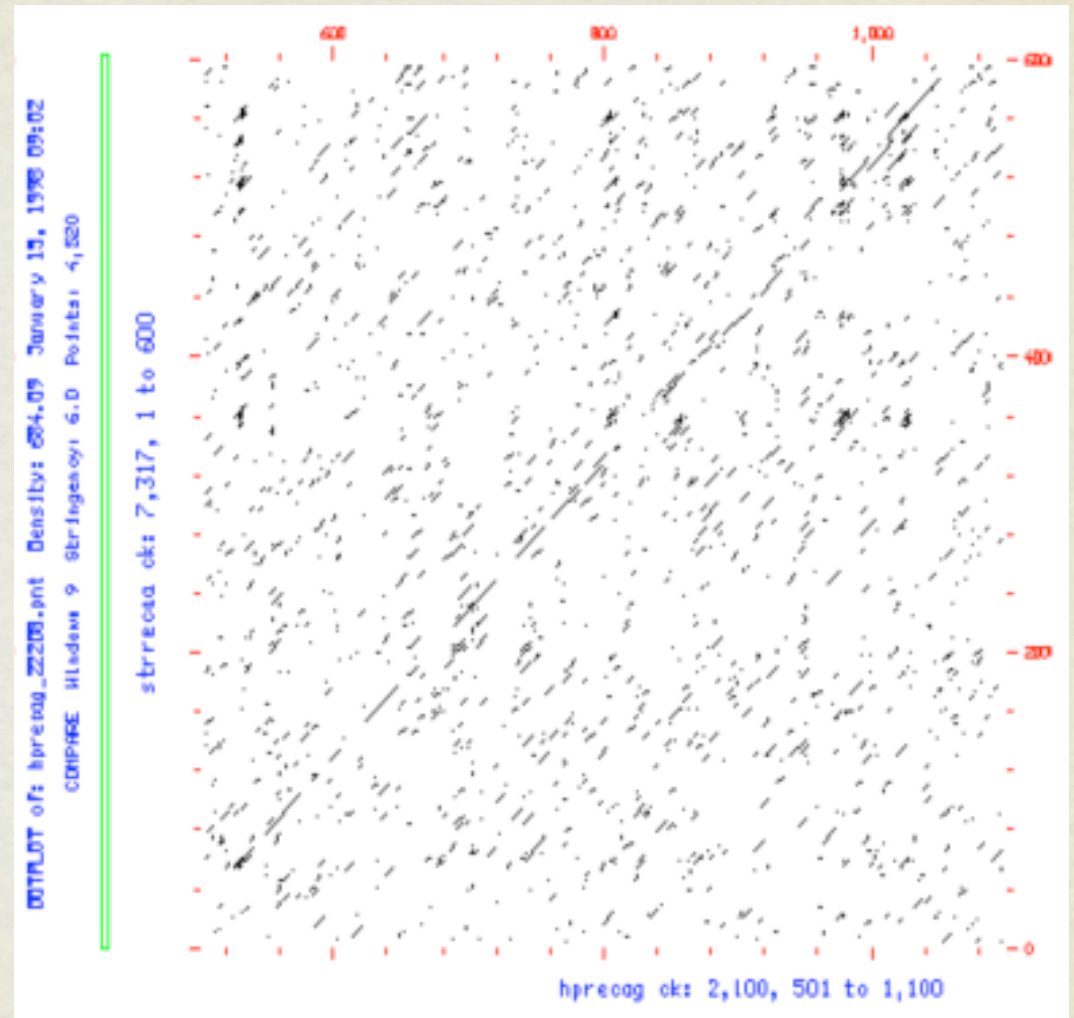
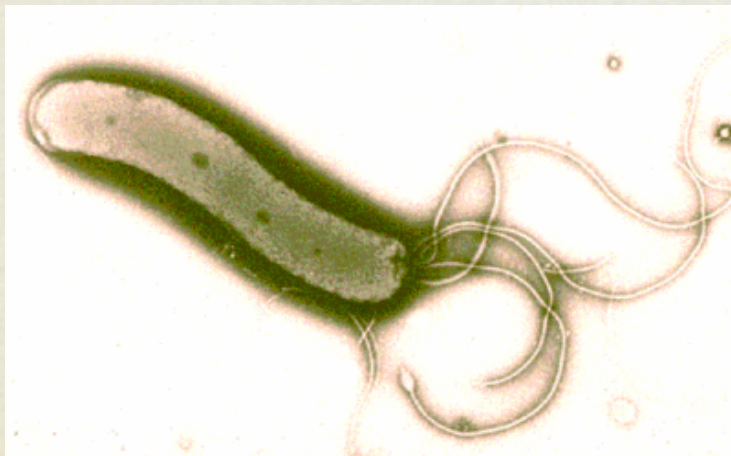
RecA DNA sequence  
from *Helicobacter pylori*  
and *Streptococcus*  
mutant, window=6  
match=4





# DOT PLOT - EXAMPLES

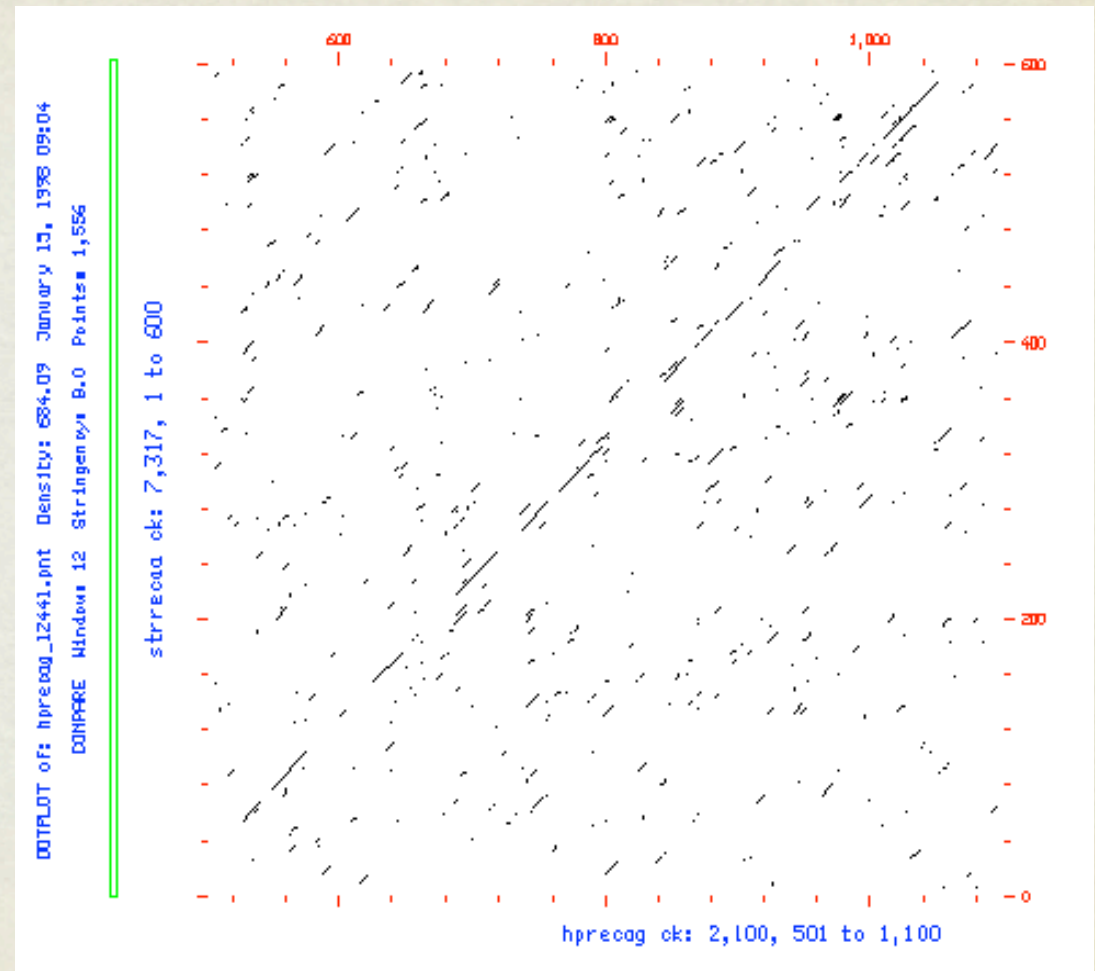
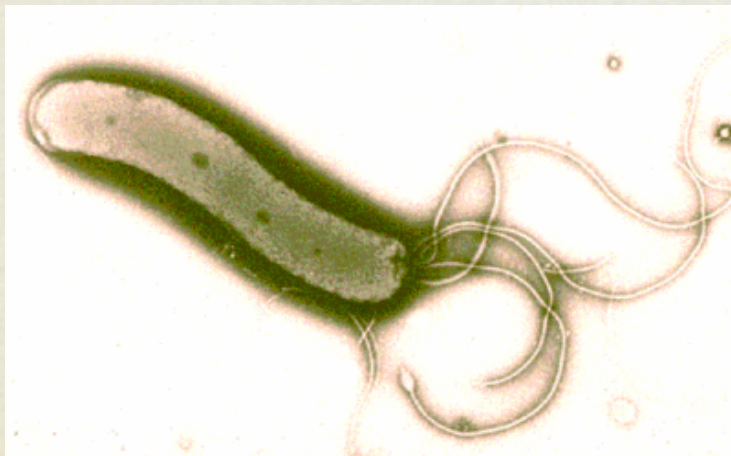
RecA DNA sequence  
from *Helicobacter pylori*  
and *Streptococcus*  
mutant, window=9  
match=6





# DOT PLOT - EXAMPLES

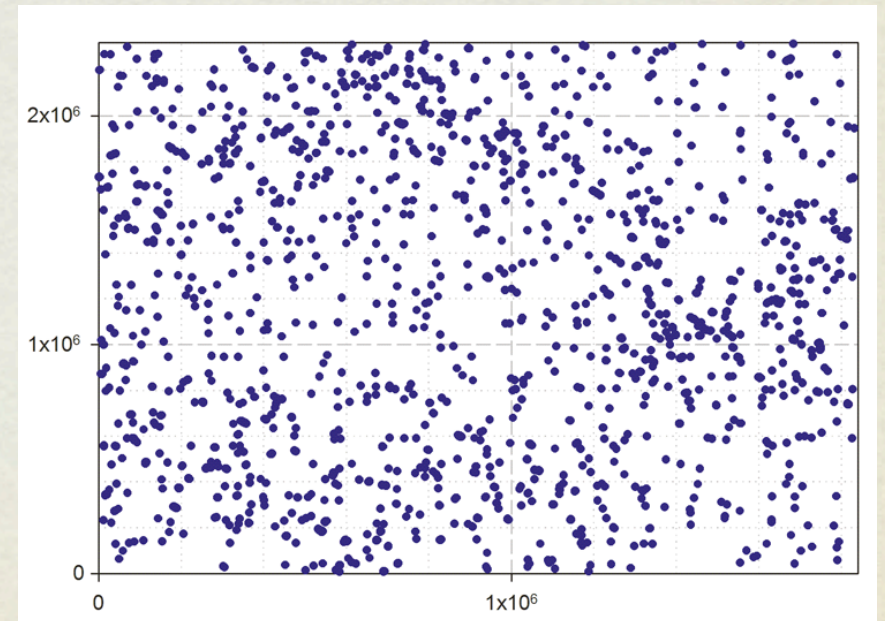
RecA DNA sequence  
from *Helicobacter pylori*  
and *Streptococcus*  
mutant, window=12  
match=8





# DOT PLOT - WHAT CAN YOU SEE THERE?

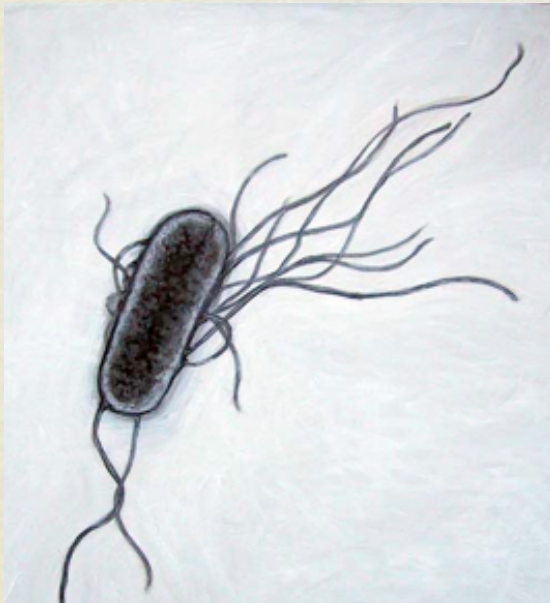
- Similar regions
- Repeated sequences
- Sequence rearrangements
- RNA structures
- Gene order





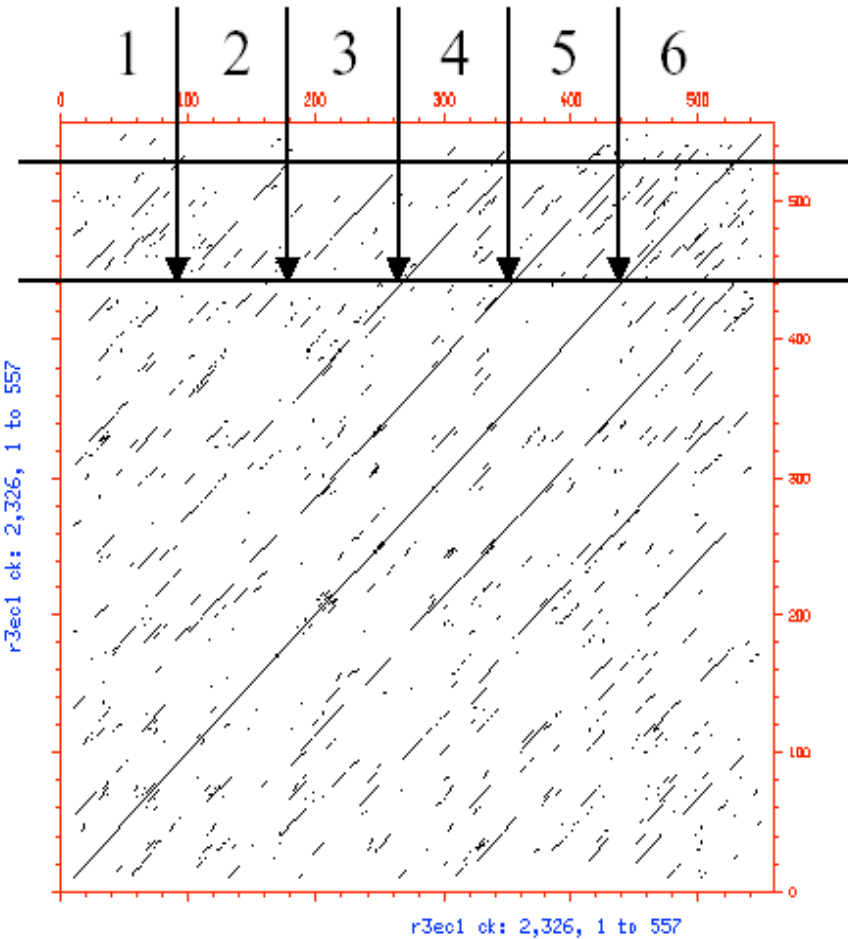
# DOT PLOT EXAMPLES - REPEATS

Repeated sequence  
in *Escherichia coli*  
ribosomal protein S1



DOTPLOT of: r3ec1\_12319.pnt Density: 633.23 January 13, 1998 09:14  
COMPARE Window 20 Stringency 10.0 Points 4,496

r3ec1 ck: 2,326, 1 to 557

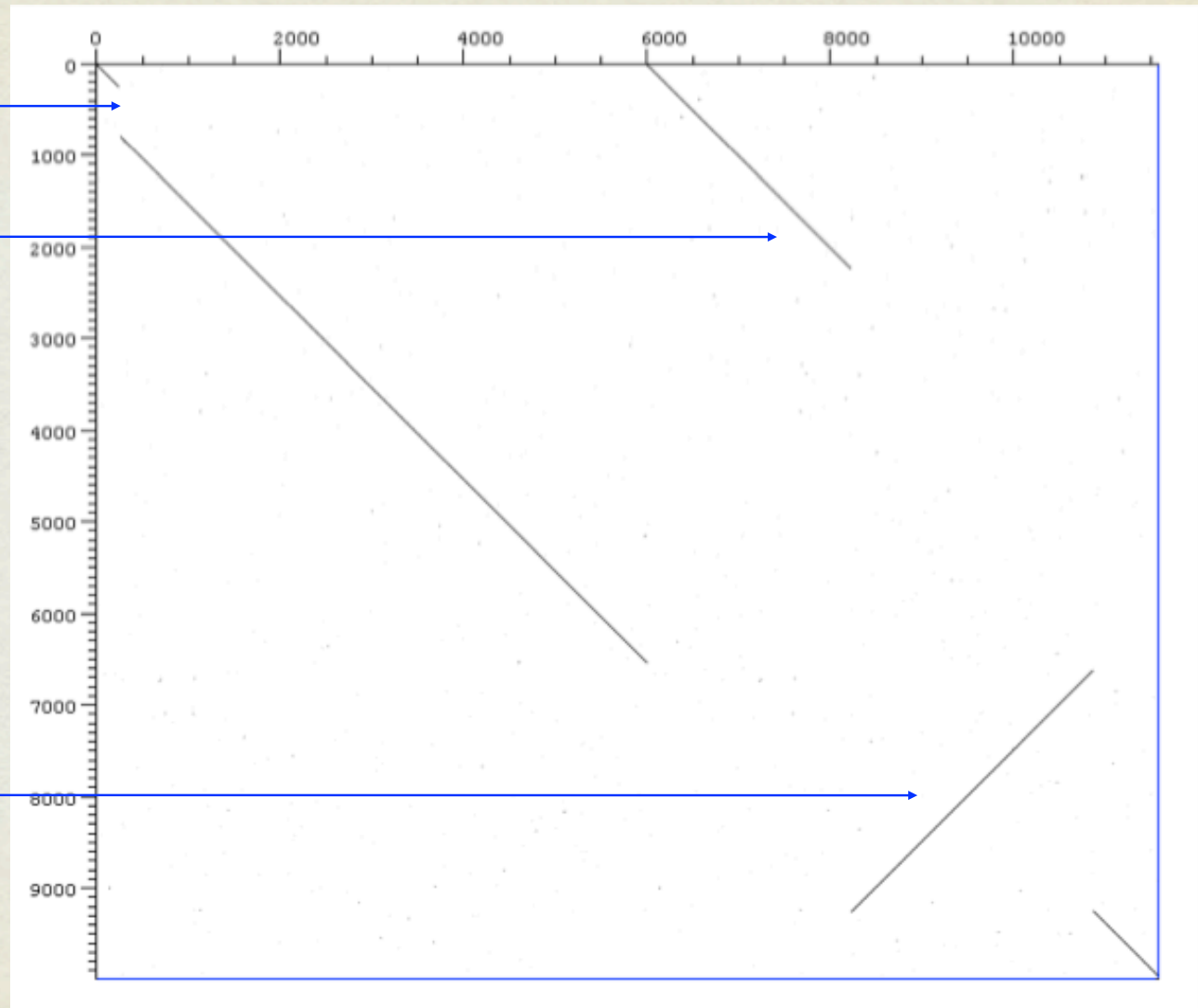


# DOT PLOT EXAMPLES - REARRANGEMENTS

deletion

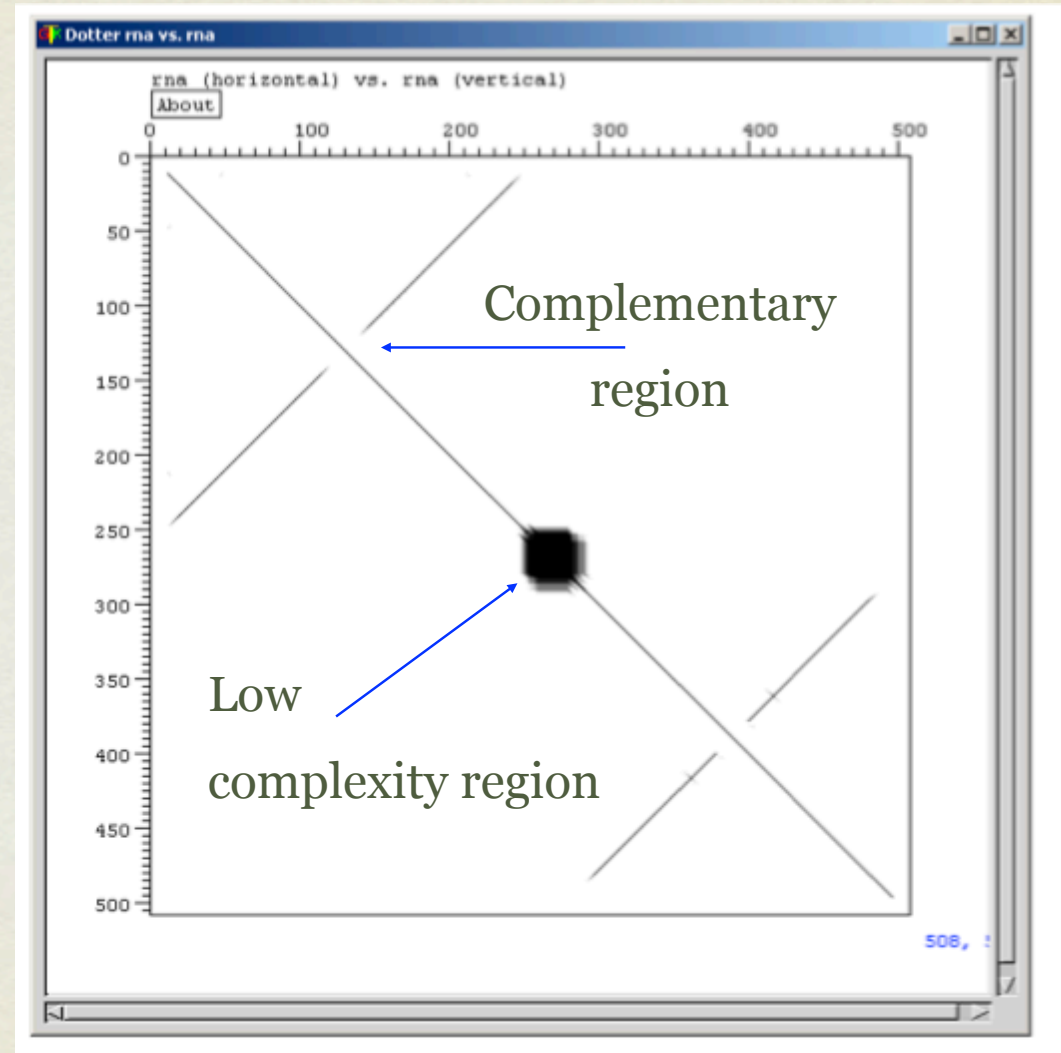
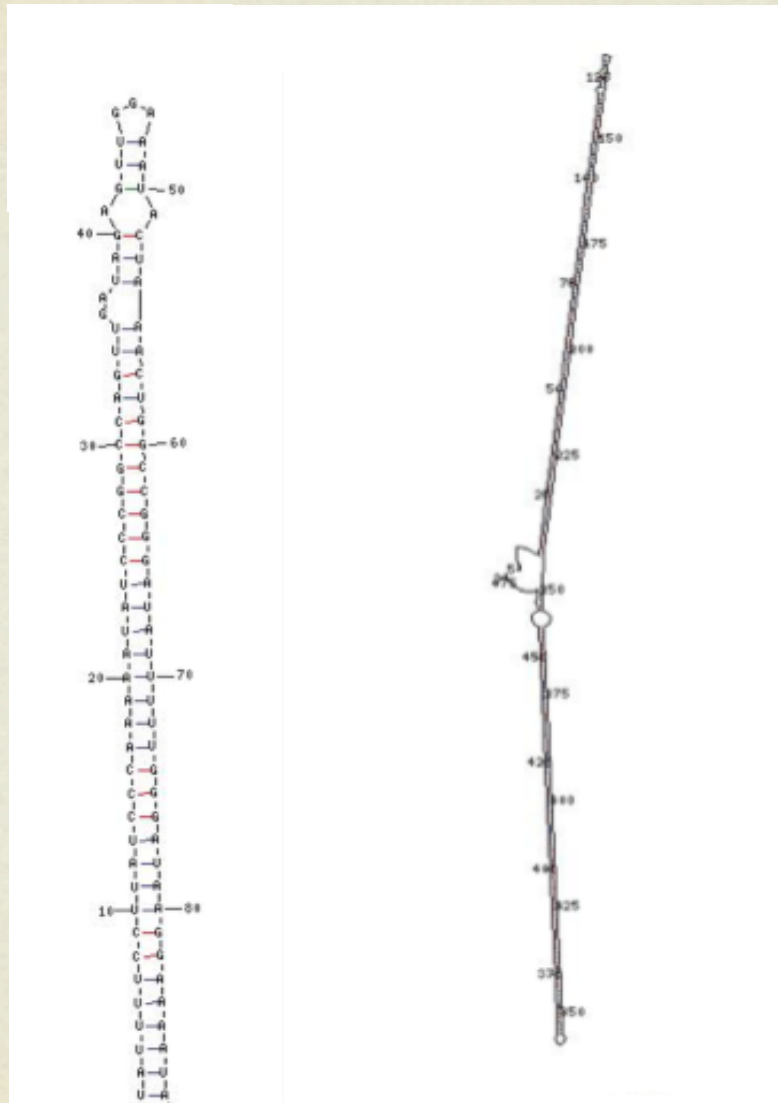
duplication

inversion





# DOT PLOT EXAMPLES - RNA STRUCTURE

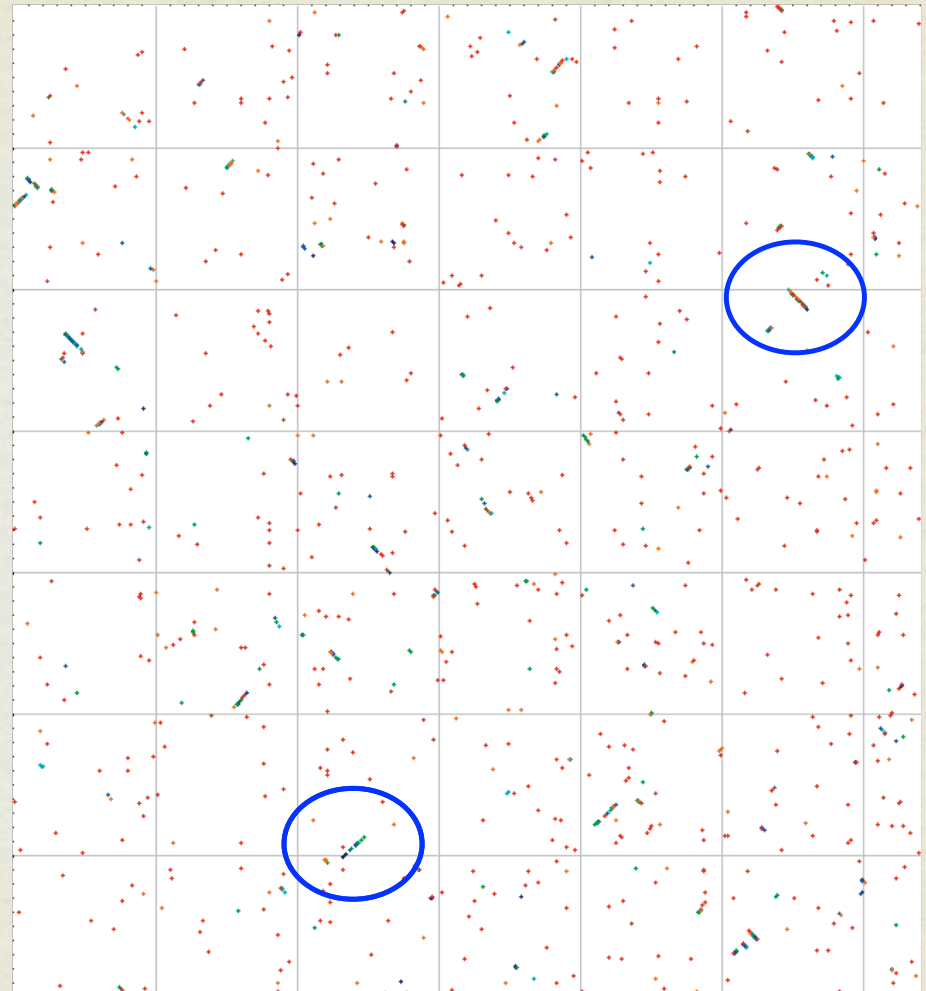


# DOT PLOT EXAMPLES - GENE ORDER

Whole genome  
comparison of  
*Buchnera* against  
*Wigglesworthia*

red dots - genes on the same  
strand

green dots - genes on  
opposite strand



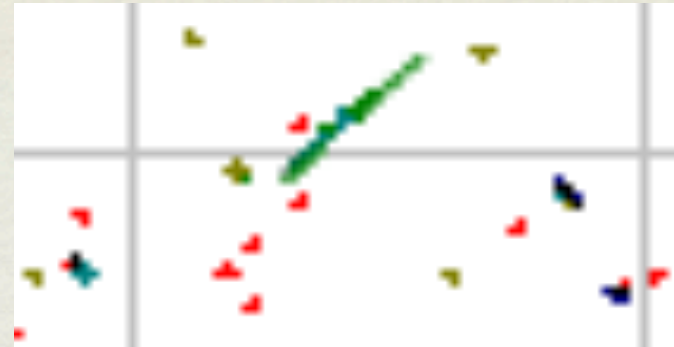


# DOT PLOT EXAMPLES - POTENTIAL OPERONS

Whole genome  
comparison of  
*Buchnera* against  
*Wigglesworthia*

red dots - genes on the same  
strand

green dots - genes on  
opposite strand



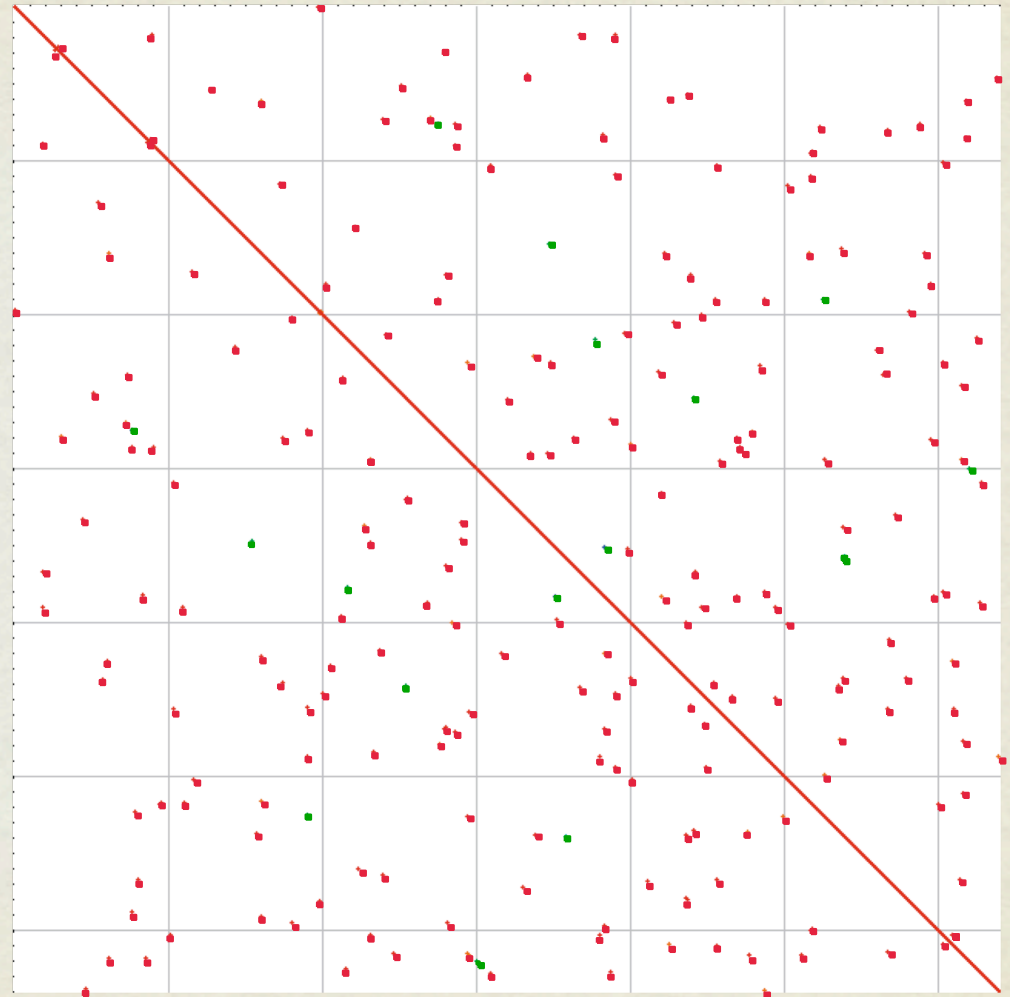
# DOT PLOT EXAMPLES - PARALOGOUS GENES

Whole genome  
comparison of  
*Wigglesworthia*

**red dots** - paralogs on the  
same strand

**green dots** - paralogs on  
opposite strand

Note: self-hits of all genes  
form **red** diagonal line





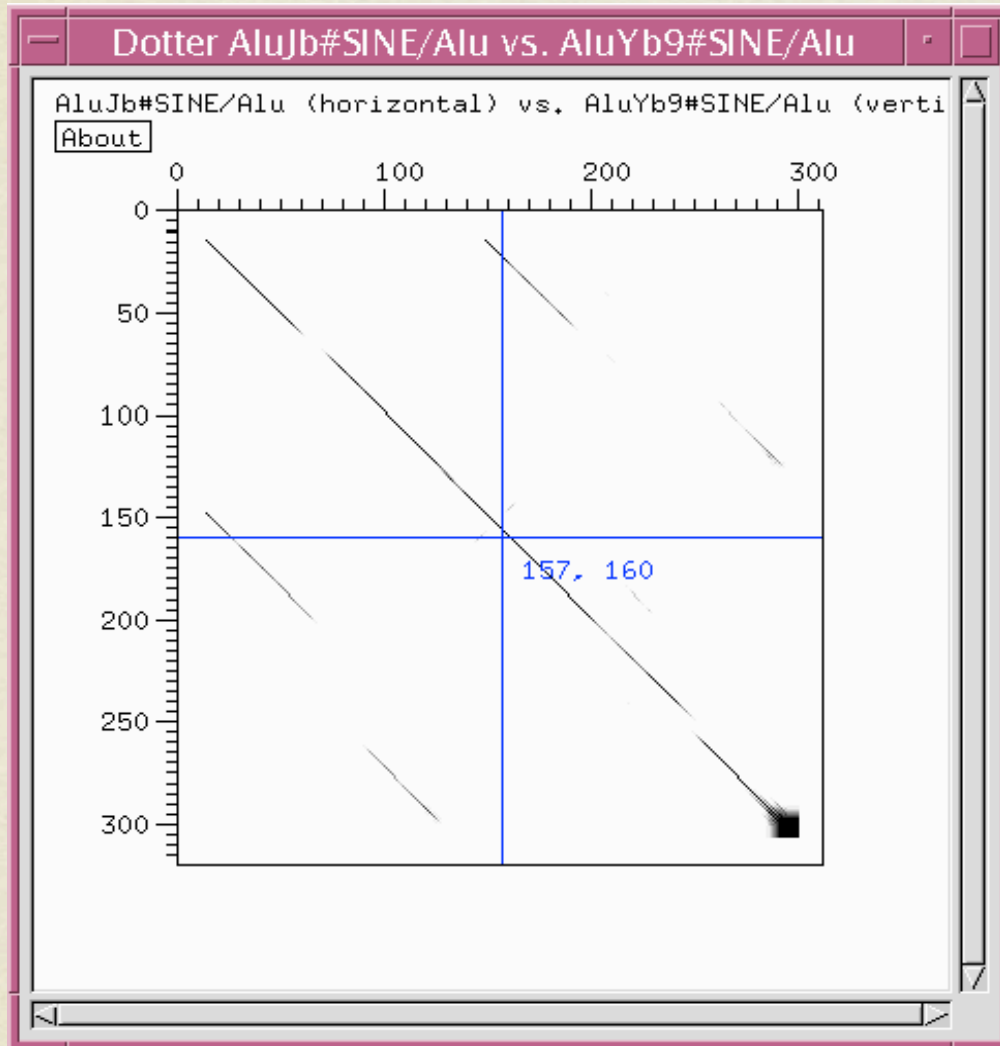
# DOT PLOTS

## RULES OF THUMB

- Don't get too many points, about 3-5 times the length of the sequence is about right (1-2%)
- Window size about 20 for distant proteins and 12 for nucleic acid (try stringency 50%)
- Check sequence against itself
  - Finds internal repeats
- Check sequence against another sequence
  - Finds repeats and rearrangements
- The best programs should have dynamic adjustment of parameters
  - dotlet: <http://myhits.isb-sib.ch/cgi-bin/dotlet>
  - dotter: <http://sonnhammer.sbc.su.se/Dotter.html>



# DOT PLOTS VERSUS ALIGNMENTS



Dotter – Alignment Tool

169

CAAAAATTAGCCGGGCGTGGTGGCGCGCGCTGTAGTCCCAGCTACTGGGAGGCTGAGGCAG

AAAAAATTAGCCGGGCGCAGTGGCGGGCGCTGTAGTCCCAGCTACTGGGAGGCTGAGGCAG

168



# ALIGNMENT

- Linear representation of relation between sequences that shows one-to-one correspondence between amino acid or nucleotide residue
- How can we define a quantitative measure of sequence similarity?
  - match
  - mismatch
  - gap

**gctg-aa-cg**  
**-ctataa-tc**



# ALIGNMENT PROBLEM

THIS IS COMPLETELY NEW SEQUENCE

THIS IS SUPEREXTRA SEQUENCE



# ALIGNMENT PROBLEM

THIS IS AN ANCESTRAL SEQUENCE  
THIS IS COMPLETELY NEW SEQUENCE

THIS IS AN ANCESTRAL SEQUENCE  
THIS IS SUPEREXTRA SEQUENCE



# ALIGNMENT PROBLEM

THISISANANCEST-R--ALSEQUENCE  
THISISCOMP-LETELYNEWSEQUENCE

THISISANANCES-TRALSEQUENCE  
THISISSU-PEREXTRA-SEQUENCE



# ALIGNMENT PROBLEM

THIS IS COMP-LETELY NEW SEQUENCE  
THIS IS ANANCEST-R--AL SEQUENCE  
THIS IS ANANCES-TRAL SEQUENCE  
THIS IS SU-PEREXTRA-SEQUENCE



# ALIGNMENT PROBLEM

THIS IS COMP-LE-TELY NEW SEQUENCE  
THIS IS ANANCES-T-R--AL SEQUENCE  
THIS IS ANANCES-T-R--AL SEQUENCE  
THIS IS SU-PEREXT-R--A-SEQUENCE



# ALIGNMENT PROBLEM

THIS IS COMP-LE-TELY NEW SEQUENCE  
THIS IS SUPEREXT-R--A-SEQUENCE

The problem is that we need to model evolutionary events based on extant sequences, without knowing an ancestral sequence!



# ALIGNMENT

- ✂ Any assignment of correspondences that preserves the order of residues within the sequence is an alignment
- ✂ It is the basic tool of bioinformatics
- ✂ Computational challenge - introduction of insertions and deletions (gaps) that correspond to evolutionary events
- ✂ We must define criteria so that an algorithm can choose the best alignment



# ALIGNMENT AN EXAMPLE

Let's compare two strings **gctgaacg** and **ctataatc**

an uninformative alignment

```
-----gctgaacg  
ctataatc-----
```

an alignment without gaps

```
gctgaacg  
ctataatc
```

an alignment with gaps

```
gctga-a--cg  
--ct-ataatc
```

another alignment with gaps

```
gctg-aa-cg  
-ctataa-tc
```

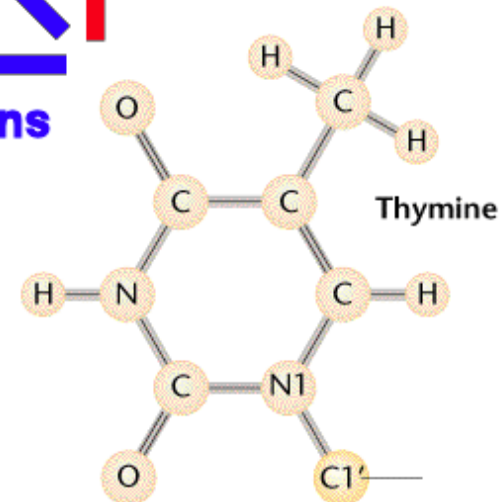
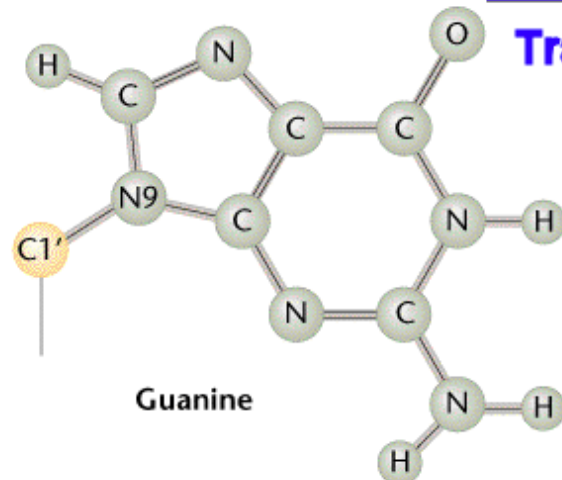
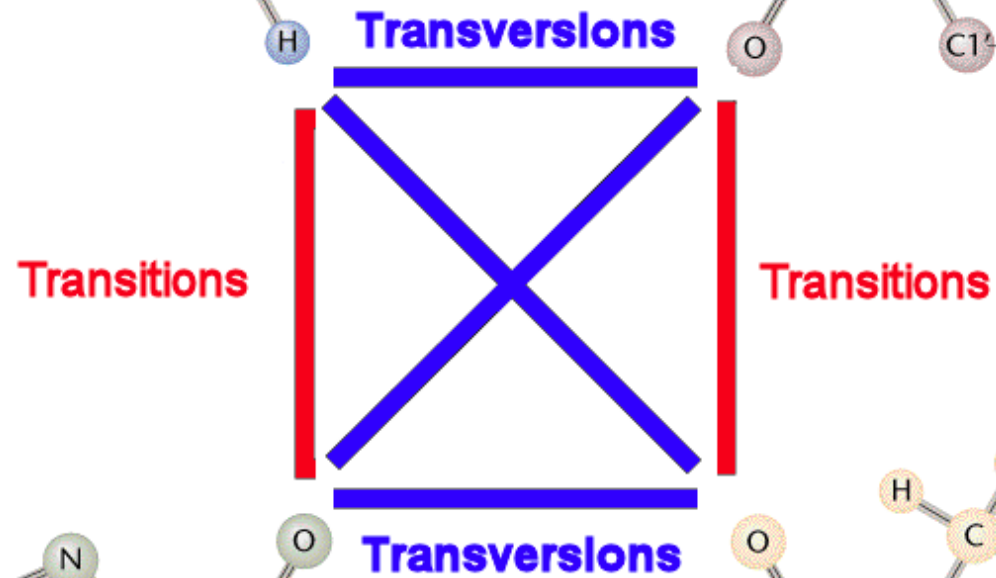
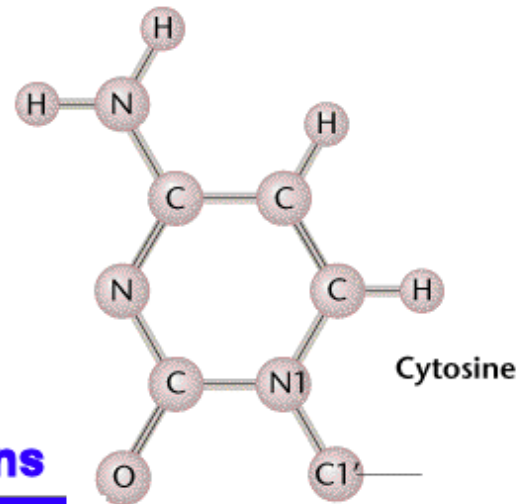
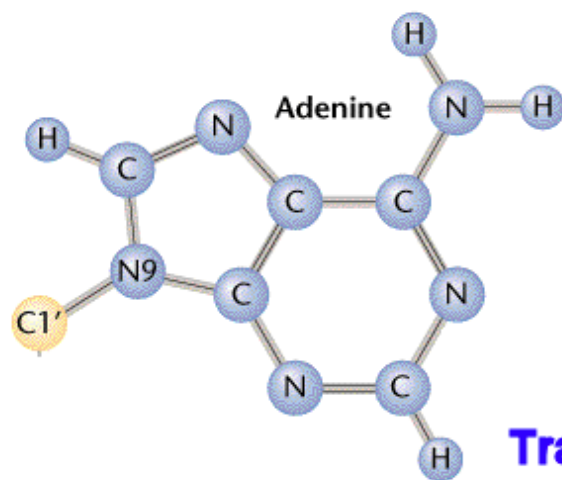




# SCORING SCHEMES

- ✂• A scoring system must account for residue substitution, and insertions or deletions (indels)
- ✂• Indels (gaps) will have scores that depend on their length
- ✂• For nucleic acid sequences, it is common to use a simple scheme for substitutions, e.g. +1 for a match, -1 for a mismatch
- ✂• More realistic would be to take into account nucleotide frequencies (sequence composition) and fact that transitions are more frequent than transversions





# GAP SCORING SYSTEMS

- non-affine model - each gap position treated the same, e.g. match = 4, mismatch = -3, gap -4
- affine model - first gap position penalized more than others, e.g. match = 4, mismatch = -3, gap opening = -8, gap = -4



# GAP SCORING

## AN EXAMPLE

non-affine gapping score - the second alignment is "better"

GGTGCCAC-TCCAC-----CTG  
AGTGCCACCCCCCAATGCCGCTG  
-3 4 4 4 4 4 4 4 -4 -3 4 4 4 -3 -4 -4 -4 -4 -4 4 4 4 = 23

GGTGCCAC-TCCA---C---CTG  
AGTGCCACCCCCCAATGCCGCTG  
-3 4 4 4 4 4 4 4 -4 -3 4 4 4 -4 -4 -4 4 -4 -4 4 4 4 = 26



# GAP SCORING

## AN EXAMPLE

affine gapping score - the first alignment is "better"

GGTGCCAC-TCCAC-----CTG  
AGTGCCACCCCCCAATGCCGCTG  
-3 4 4 4 4 4 4 4 -12 -3 4 4 4 -3 -12 -4 -4 -4 -4 4 4 4 = 7

GGTGCCAC-TCCA---C---CTG  
AGTGCCACCCCCCAATGCCGCTG  
-3 4 4 4 4 4 4 4 -12 -3 4 4 4 -12 -4 -4 4 -12 -4 4 4 4 = 2



# GAP SCORING

## AN EXAMPLE

### Equivalent alignments

GGTGCCAC-TCCA---C--CTG  
AGTGCCACCCCCCAATGCCGCTG  
-3 4 4 4 4 4 4 4 -12 -3 4 4 4 -12 -4 -4 4 -12 -4 4 4 4 = 2

GGTGCCACT-CCA---C--CTG  
AGTGCCACCCCCCAATGCCGCTG  
-3 4 4 4 4 4 4 4 -3 -12 4 4 4 -12 -4 -4 4 -12 -4 4 4 4 = 2



# AMINO ACID SCORING SYSTEMS

- more complicated than nucleotide matrices
- first, we can align two homologous protein sequences and count the number of any particular substitution, for instance Serine to Threonine
- a likely change should score higher than a rare one
- we have to take into account that several the same position mutated several times after sequence divergence - this could bias statistics



# AMINO ACID SCORING SYSTEMS

- to avoid this problem one can compare very similar sequences so one can assume that no position has changed more than once
- Margret Dayhoff introduced the PAM system (Percent of Accepted Mutations)



- 1 PAM - two sequence have 99% identical residues
- 10 PAM - two sequence have 90% identical residues



# APPROXIMATE RELATION BETWEEN PAM AND SEQUENCE IDENTITY

PAM	0	30	80	110	200	250
AA sequence identity (%)	100	75	50	60	25	20

PAM matrix is expressed as log-odds values multiplied by 10 simply to avoid decimal points



# PAM MATRIX CALCULATION

$$\text{score of substitution } i \leftrightarrow j = \log \frac{\text{observed } i \leftrightarrow j \text{ mutation rate}}{\text{mutation rate expected from amino acids frequencies}}$$

For instance, a value 2 implies that in related sequences the mutation would be expected to occur 1.6 times more frequently than random.

The calculation: The matrix entry 2 corresponds to the actual value 0.2 because of the scaling. The value 0.2 is  $\log_{10}$  of the relative expectation value of the mutation. Therefore, the expectation value is  $10^{0.2} = 1.6$



# AMINO ACID MATRICES

- ✂ Problem with PAM schema lies in that the high number matrices are extrapolated from closely related sequences
- ✂ Henikoffs developed the family of BLOSUM matrices based on the BLOCKS database of aligned protein sequences, hence the name BLOcks SUBstitution Matrix
- ✂ observed substitution frequencies taken from conserved regions of proteins (blocks), not the whole proteins as in case of Dayhoff's work
- ✂ to avoid overweighting closely related sequences, the Henikoffs replaced groups of proteins that have sequence identities higher than a threshold by either a single representative or a weighted average, e.g. for the commonly used BLOSUM62 matrix the threshold is 62%
- ✂ NOTE reversed numbering of PAM and BLOSUM matrices



# BLOSUM 62 SCORING MATRIX

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

some replacement are more frequent than others

score system based on comparison of homologous domains



# BLOSUM 62 SCORING MATRIX

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

identical amino acids  
positive score but  
- frequent amino acids, e.g. alanine, get lower scores  
rare amino acids, e.g. tryptophan

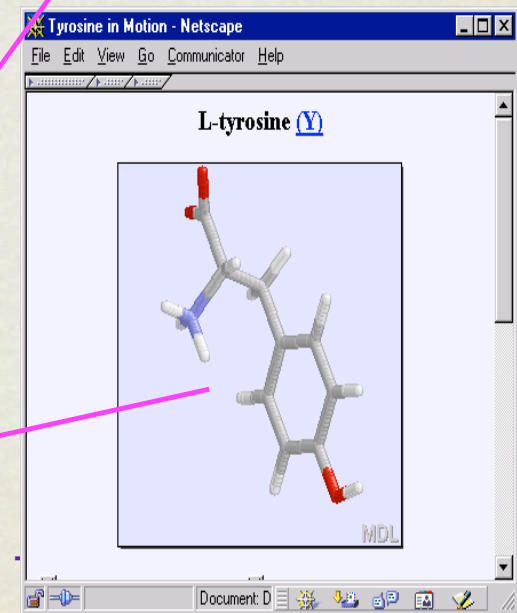
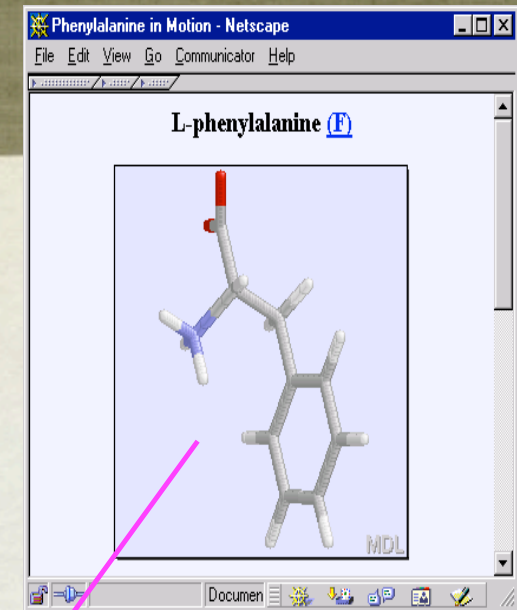
identical amino acids get positive score but not the same  
 – frequent amino acids, e.g. alanine, get lower score than rare amino acids, e.g. tryptophan



# BLOSUM 62 SCORING MATRIX

substitutions to amino acids of similar properties give a positive score

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

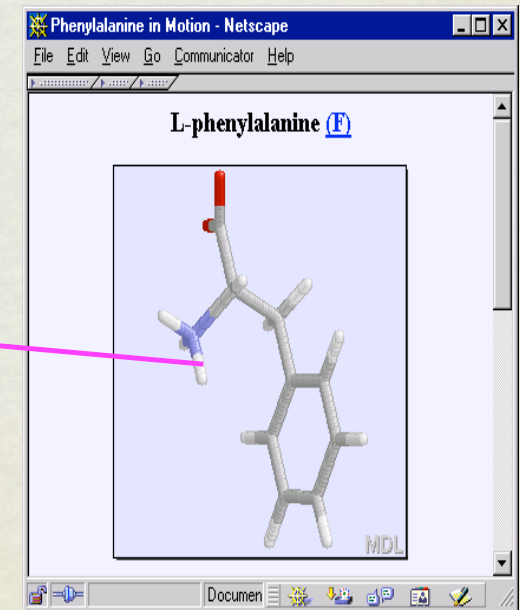
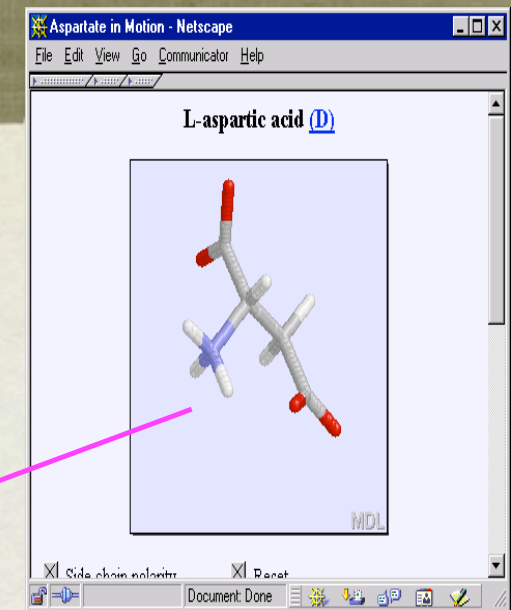




# BLOSUM 62 SCORING MATRIX

substitutions to amino acids of different properties give a negative score

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





# POSITION SPECIFIC SUBSTITUTION MATRIX

	210	220	230	240	250	260	270
	.....*	.....*	.....*	.....*	.....*	.....*	.....*
154	AT-CRR-AYSGG-----	PAITDNMLCAGGLE--	GGKDACQGD	SGGPLVC	NDP----	RWVLVGI	
707	KV-CNR-YEYLG-----	GKVSPNELCAGHLA--	GGIDSCQGD	SGGPLVC	fEKD----	KYILQGV	
187	KGsCER-DAQYApgydkvkdis	EVVTPRFLCTGGV	SpYADPNTCRGD	SGGPLIV	-HKr----	sRFIQVGV	
166	SS-CKS-AYPG-----	QITSNMFCAGYLE--	GGKDSQGD	SGGPVVC	-S-----	GKLQGI	
166	TN-CK--KYWG-----	TKIKDAMICAGA---	SGVSSCMGD	SGGPLVC	-KKn----	gAWTLVGI	
160	SD-CNN-SYPG-----	MITNAMFCAGYLE--	GGKDSQGD	SGGPVVC	-N-----	GELQGV	
167	KV-CNRyEFLNG-----	RVQSTELCAGHI--	GGTDSQGD	SGGPLVC	fEk----	dKYILQGV	
156	AT-CSKpGWWG-----	STVKTNMICAGGI--	GIISSCNGD	SGGPLM	QGan----	gQWQVHGI	
156	AI-CSSsSYWG-----	SVRSQCQGD	SGGPL			HGV	
169	EH-CSQyDWWG-----	RSQCDGD	SGG			HGV	
156	HI-CDAkYHLGAytg----	ddv	DSQGD			AGV	
168	AT-CSQrDWWG-----	ISACNGD	SGG			RGI	
165	EE-CS--QTWGN-----	AGATSCMGD	SGGP			VGI	
185	NQ-CR--QYWG-----	STTDSMICAGS	AGASSCQGD	SGGPLV	qKq----	ntWVLIGI	
518	ER-CSS-PEVHGd-----	AFLSG-MLCAGFI	g--gTDACQGD	SGGPLVC	Edea-aehRLILRGI		
320	DV-CNGaDFYGN-----	QIKPKMFCAGYFe--	GGIDACQGD	SGGPVVC	EDsisrtpRWRLCGI		
183	DE-CEK-AHVQ-----	KVTDFMLCVGHLE--	GGKDTCVGD	SGGPLMC	-D-----	GVLQGV	
160	KN-CDD-AHIA-----	NVTGTMLCAGDLA--	GGKDTCVGD	SGGPLIC	-D-----	GVLQGL	
182	DM-CAR-AYSE-----	KVTEFMLCAGLWT--	GGKDTCCGD	SGGPLVC	-N-----	GVLQGI	
678	RF-CKE-RYKG-----	LFTGRMLCAGNLQ	edNRVDSQGD	SGGPLMC	-EKpd---	eSWVVG	
408	QR-CNSrYVYDN-----	LITPAMICAGFLq--	GNVDSQGD	SGGPLVT	-SNn----	niWWLIGD	
475	TL-CNSrQLYDH-----	MIDDSMICAGNLQ	k-pGQDTCQGD	SGGPLTC	-EKd----	gTYYVYGI	
181	AH-CSRyDWWG-----	SLVTTSMVCAGGD--	GVLABSCNGD	SGGPLNC	-QNad---	gSWDVHGV	
176	AE-CAA-ALVNv-----	vPVTEQMICAGYAA	g--GKDSQGD	SGGPLVS	-GD-----	KLGV	
181	RE-CNChYQTILeq-----	ddEVIKQDMLCAGSE--	G-HDSQQMD	SGGPLVC	-RWk----	cTWIQVGV	

weakly  
conserved  
serine

active site  
serine



# POSITION SPECIFIC SUBSTITUTION MATRIX

		A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
206	D	0	-2	0	2	-4	2	4	-4	-3	-5	-4	0	-2	-6	1	0	-1	-6	-4	-1
207	G	-2	-1	0	-2	-4	-3	-3	6	-4	-5	-5	0	-2	-3	-2	-2	-1	0	-6	-5
208	V	-1	1	-3	-3	-5	-1	-2	6	-1	-4	-5	1	-5	-6	-4	0	-2	-6	-4	-2
209	I	-3	3	-3	-4	-6	0	-1	-4	-1	2	-4	6	-2	-5	-5	-3	0	-1	-4	0
210	S	-2	-5	0	8	-5	-3	-2	-1	-4	-7	-6	-4	-6	-7	-5	1	-3	-7	-5	-6
211	S	4	-4	-4	-4	-4	-1	-4	-2	-3	-3	-5	-4	-4	-5	-1	4	3	-6	-5	-3
212	C	-4	-7	-6	-7	12	-7	-7	-5	-6	-5	-5	-7	-5	0	-7	-4	-4	-5	0	-4
213	N	-2	0	2	-1	-6	7	0	-2	0	-6	-4	2	0	-2	-5	-1	-3	-3	-4	-3
214	G	-2	-3	-3	-4	-4	-4	-5	7	-4	-7	-7	-5	-4	-4	-6	-3	-5	-6	-6	-6
215	D	-5	-5	-2	9	-7	-4	-1	-5	-5	-7	-7	-4	-7	-7	-5	-4	-4	-8	-7	-7
216	S	-2	-4	-2	-4	-4	-3	-3	-3	-4	-6	-6	-3	-5	-6	-4	7	-2	-6	-5	-5
217	G	-3	-6	-4	-5	-6	-5	-6	8	-6	-8	-7	-5	-6	-7	-6	-4	-5	-6	-7	-7
218	G	-4	-5	-4	-5	-6	-5	-6	8	-6	-7	-7	-5	-6	-7	-6	-2	-4	-6	-7	-7
219					-5	-6	-5	-5	-6	-6	-6	-7	-4	-6	-7	9	-4	-4	-7	-7	-6
220					-7	-5	-5	-6	-7	0	-1	6	-6	1	0	-6	-6	-5	-5	-4	0
221	N			0	-6	-4	-4	-6	-6	-1	3	0	-5	4	-3	-6	-2	-1	-6	-1	6
222	C	0	-4	-5	-5	10	-2	-5	-5	1	-1	-1	-5	0	-1	-4	-1	0	-5	0	0
223	Q	0	1	4	2	-5	2	0	0	0	-4	-2	1	0	0	0	-1	-1	-3	-3	-4
224	A	-1	-1	1	3	-4	-1	1	4	-3	-4	-3	-1	-2	-2	-3	0	-2	-2	-2	-3

active  
center



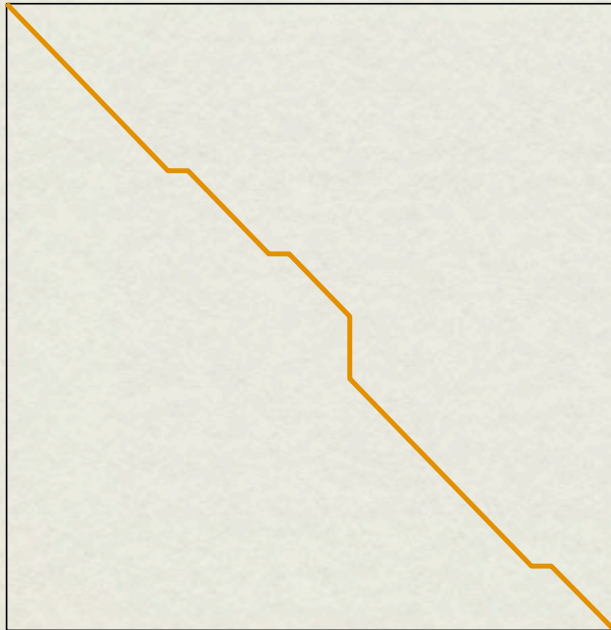
# SCORING RECOMMENDATIONS

- ✂• nucleotide sequence comparison
  - ✂• match +10, mismatch -3, gap opening -50, gap extension -5
- ✂• amino acid sequence comparison
  - ✂• for general use (e.g. unknown sequence similarity) - BLOSUM62
  - ✂• for diverged proteins - PAM250 or BLOSUM30
  - ✂• for similar sequences - PAM15 or BLOSUM80



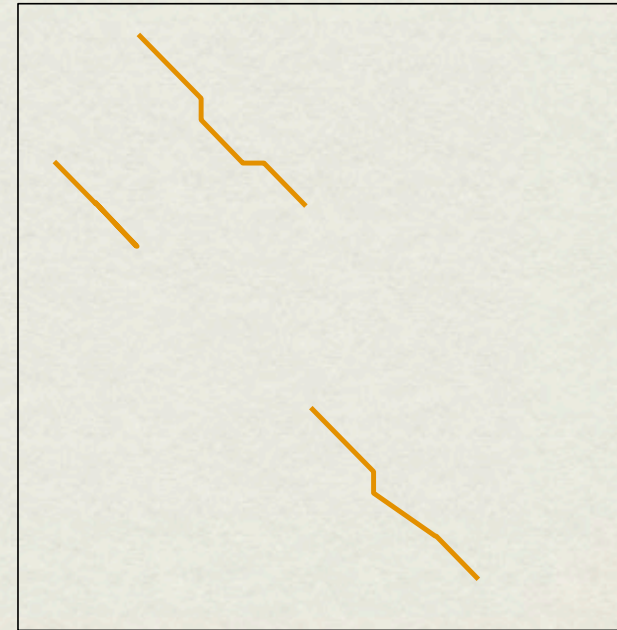
# GLOBAL VERSUS LOCAL ALIGNMENT

Optimal global  
alignment



Sequences align essentially  
from end to end. Needleman &  
Wunsch (1970)

Optimal local  
alignment



Sequences align only in small,  
isolated regions. Smith & Waterman  
(1981)



# Sequence alignment using dynamic programming



# Dynamic programming

Construct an optimal alignment of these two sequences:

**G A T A C T A**  
**G A T T A C C A**

Using these scoring rules:

Match: +1

Mismatch: -1

Gap: -1





# Dynamic programming

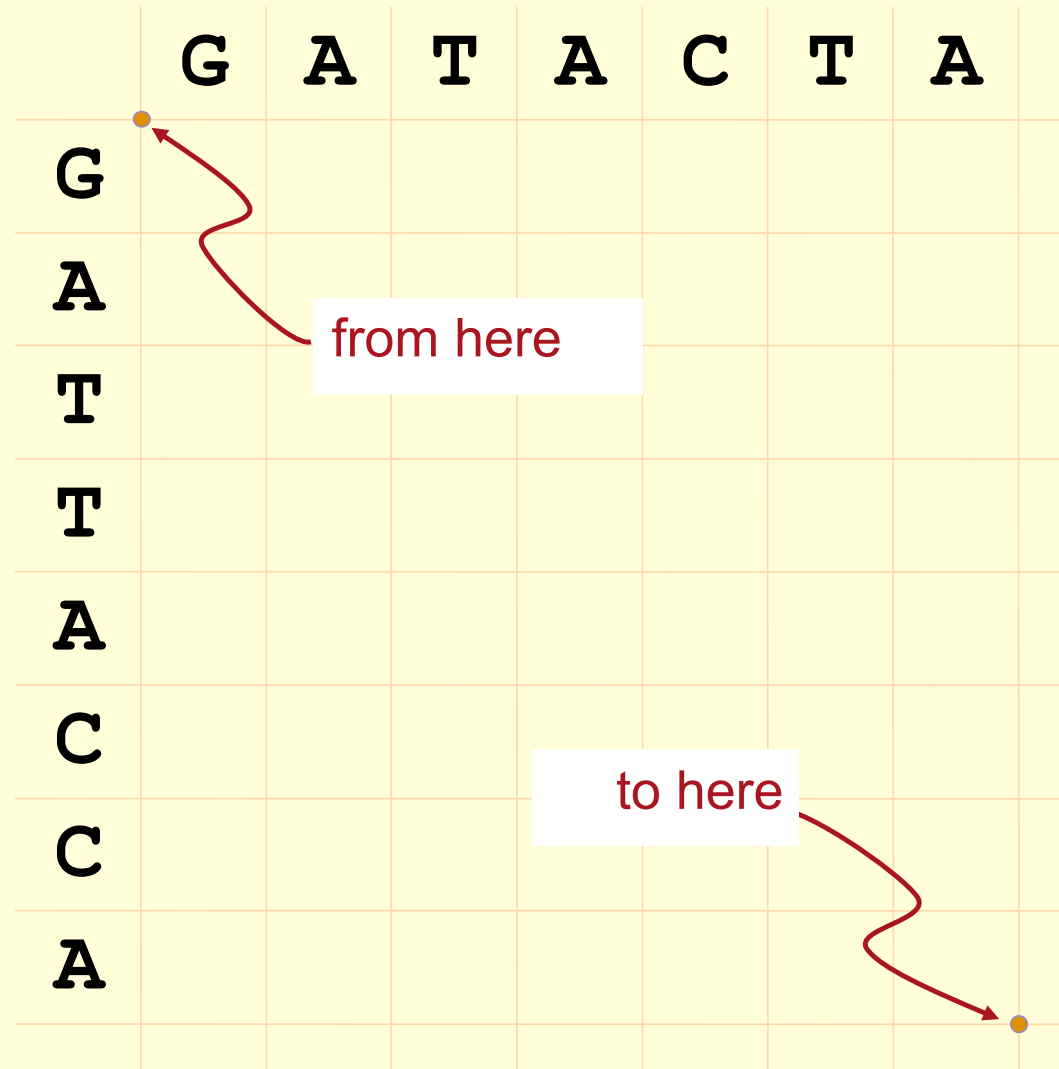
Arrange the  
sequence residues  
along a two-  
dimensional lattice

Vertices of the  
lattice fall between  
letters

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

# Dynamic programming

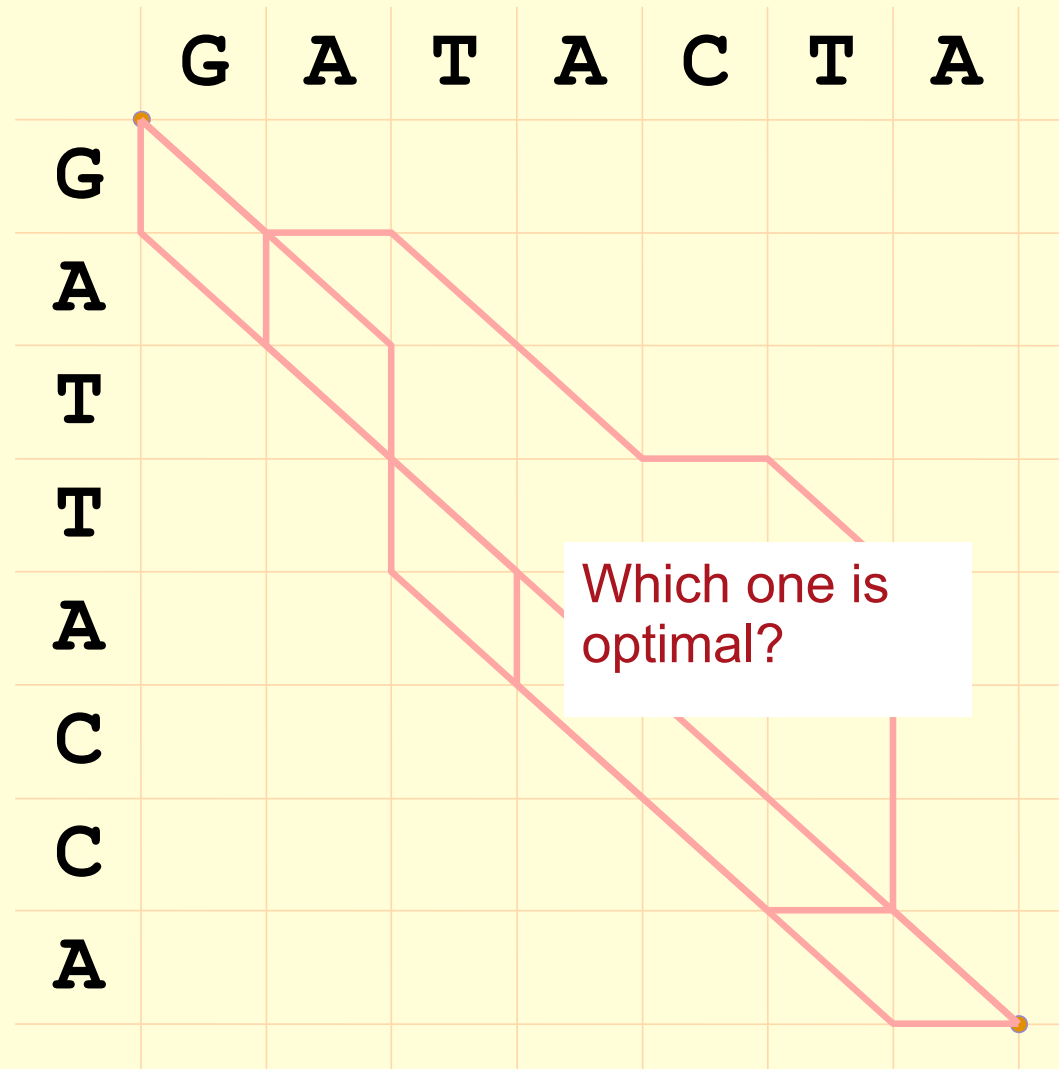
The goal is to find the optimal path





# Dynamic programming

Each path  
corresponds to a  
unique alignment



# Dynamic programming


The score for a path is the sum of its incremental edges scores

Match: +1

Mismatch: -1

Gap: -1

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



**A** aligned with **A**  
Match = +1




# Dynamic programming

The score for a path  
is the sum of its  
incremental edges  
scores

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

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



**A** aligned with **T**  
Mismatch = -1

# Dynamic programming

The score for a path is the sum of its incremental edges scores

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

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

**T** aligned with *NULL*

Gap = -1

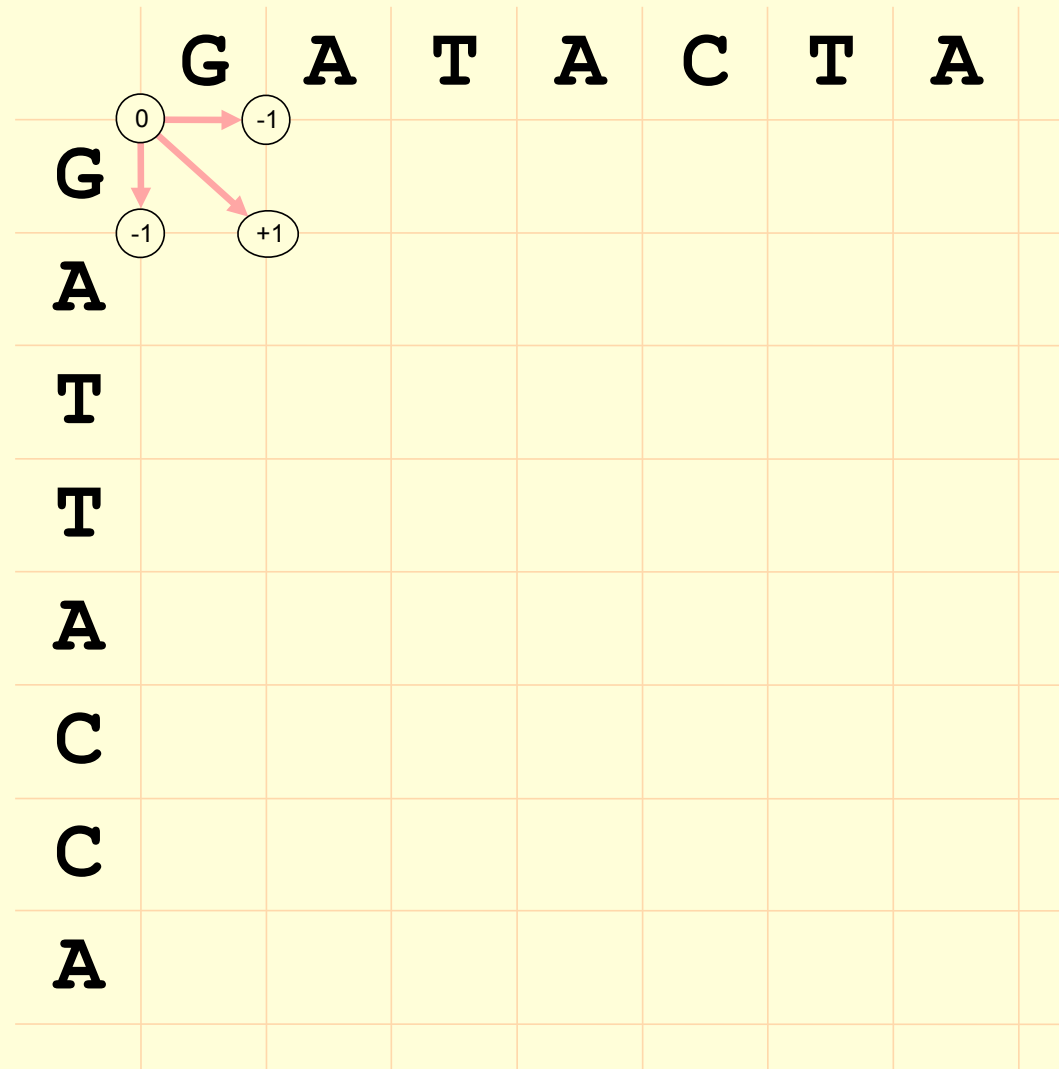
*NULL* aligned with **Tz**



# Dynamic programming

Incrementally extend  
the path

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



# Dynamic programming

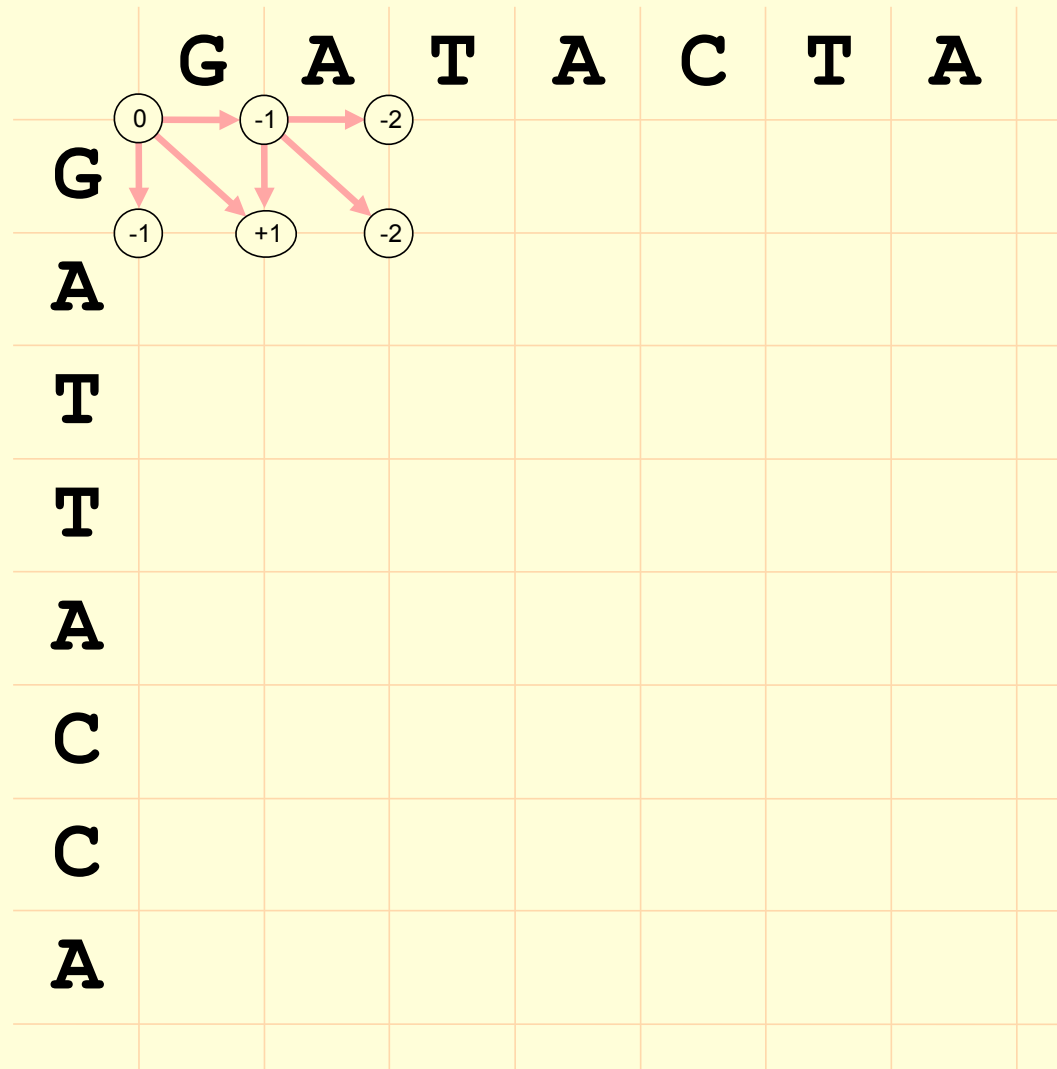
Incrementally extend  
the path

Remember the best  
sub-path leading to  
each point on the  
lattice

Match: +1

Mismatch: -1

Gap: -1





# Dynamic programming

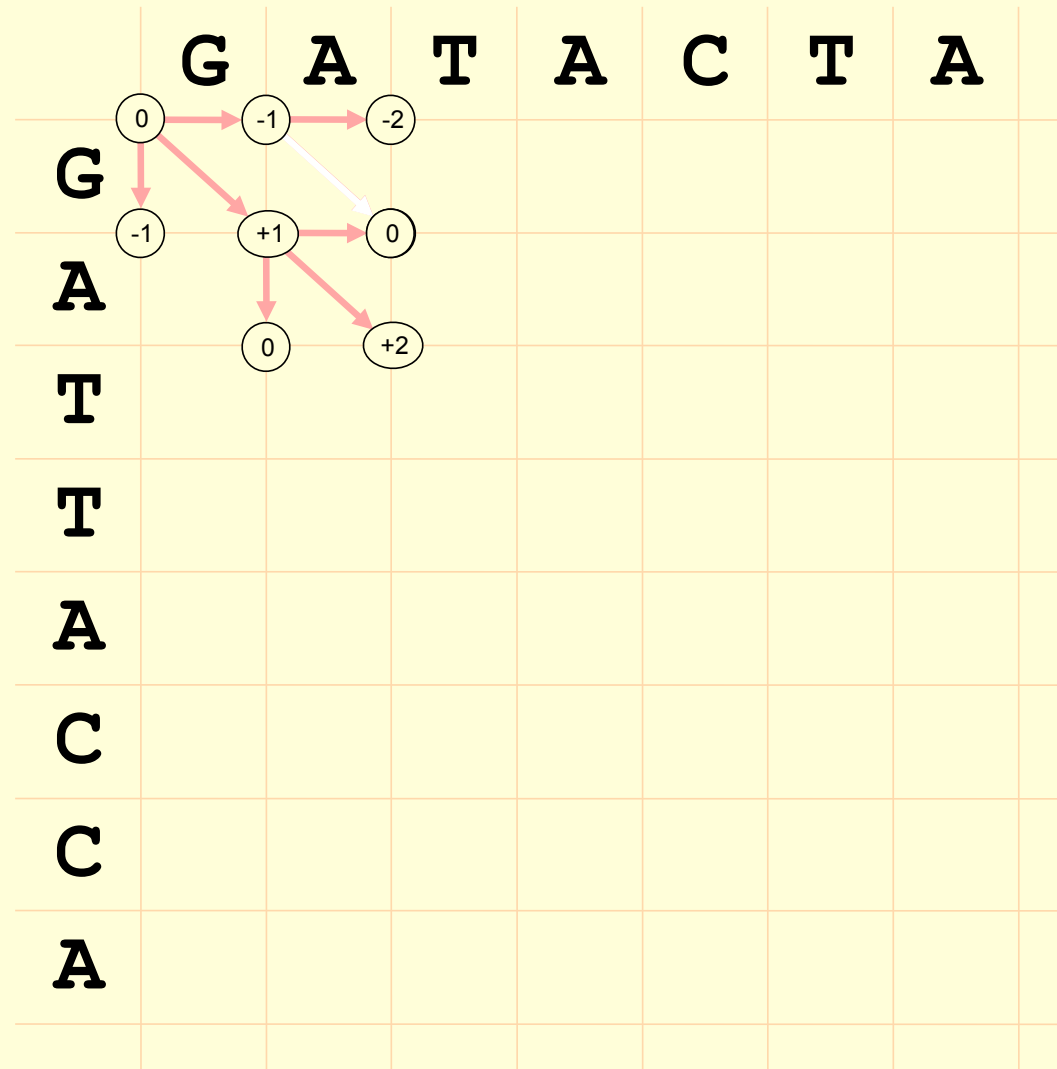
Incrementally extend  
the path

Remember the best  
sub-path leading to  
each point on the  
lattice

Match: +1

Mismatch: -1

Gap: -1



# Dynamic programming

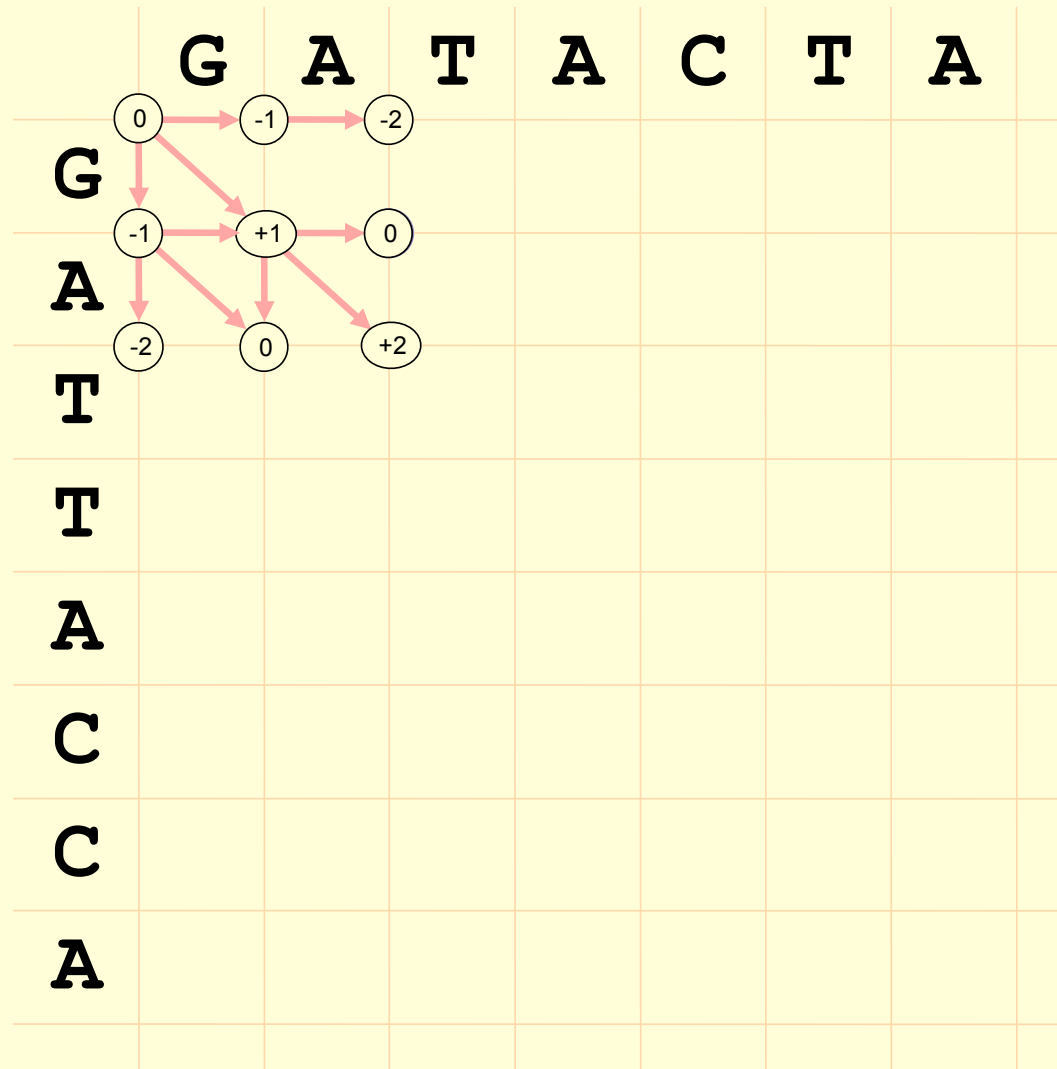
Incrementally extend  
the path

Remember the best  
sub-path leading to  
each point on the  
lattice

Match: +1

Mismatch: -1

Gap: -1





# Dynamic programming

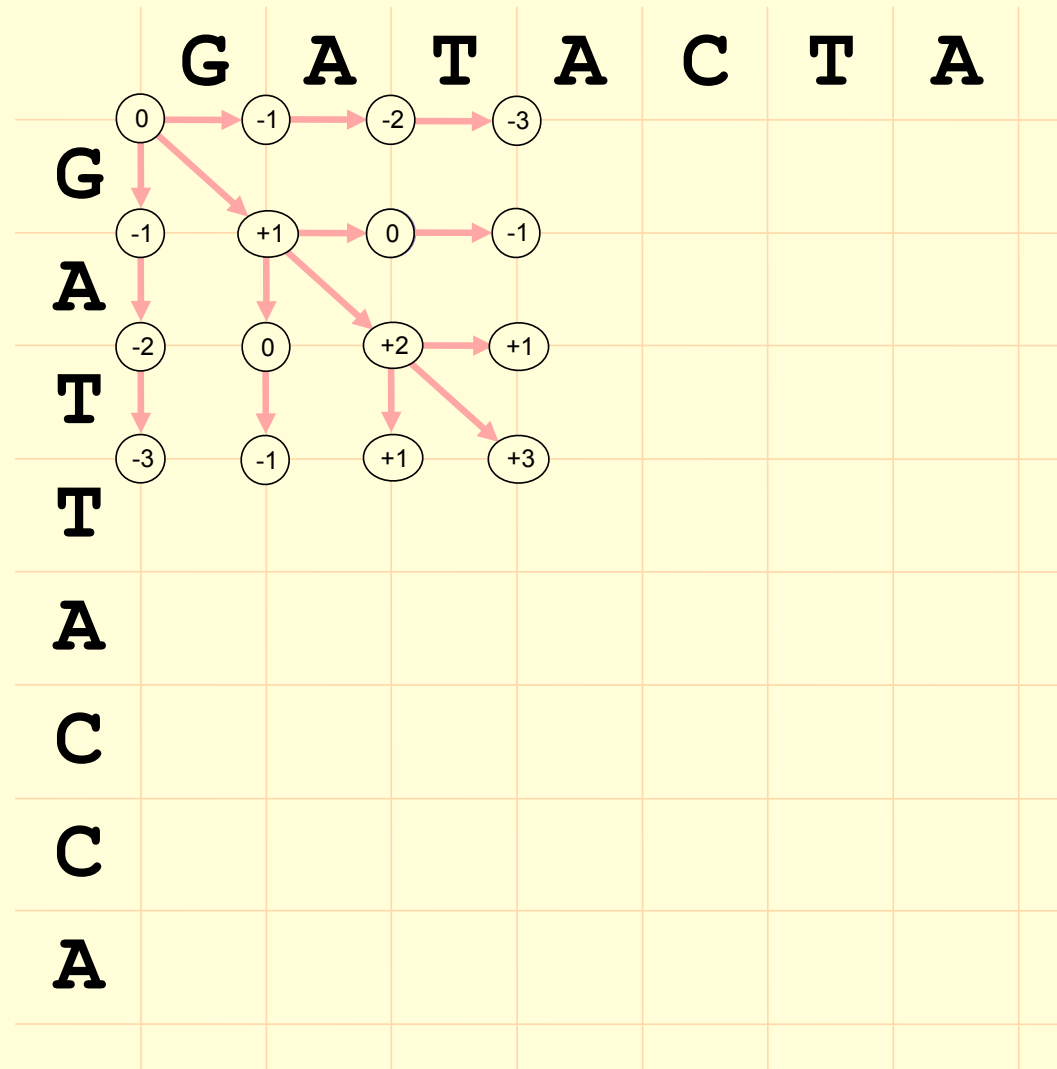
Incrementally extend  
the path

Remember the best  
sub-path leading to  
each point on the  
lattice

Match: +1

Mismatch: -1

Gap: -1



# Dynamic programming

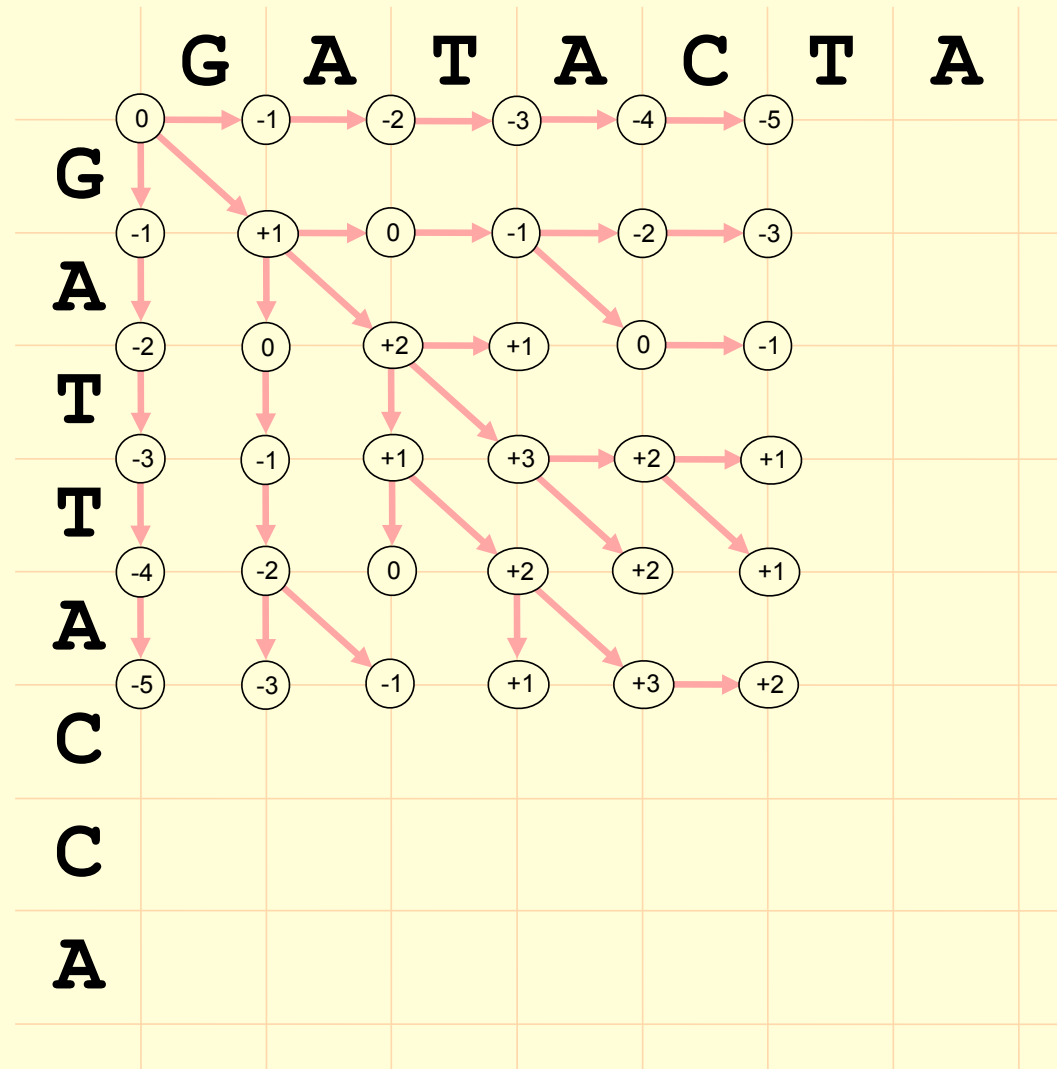
Incrementally extend  
the path

Remember the best  
sub-path leading to  
each point on the  
lattice

Match: +1

Mismatch: -1

Gap: -1





# Dynamic programming

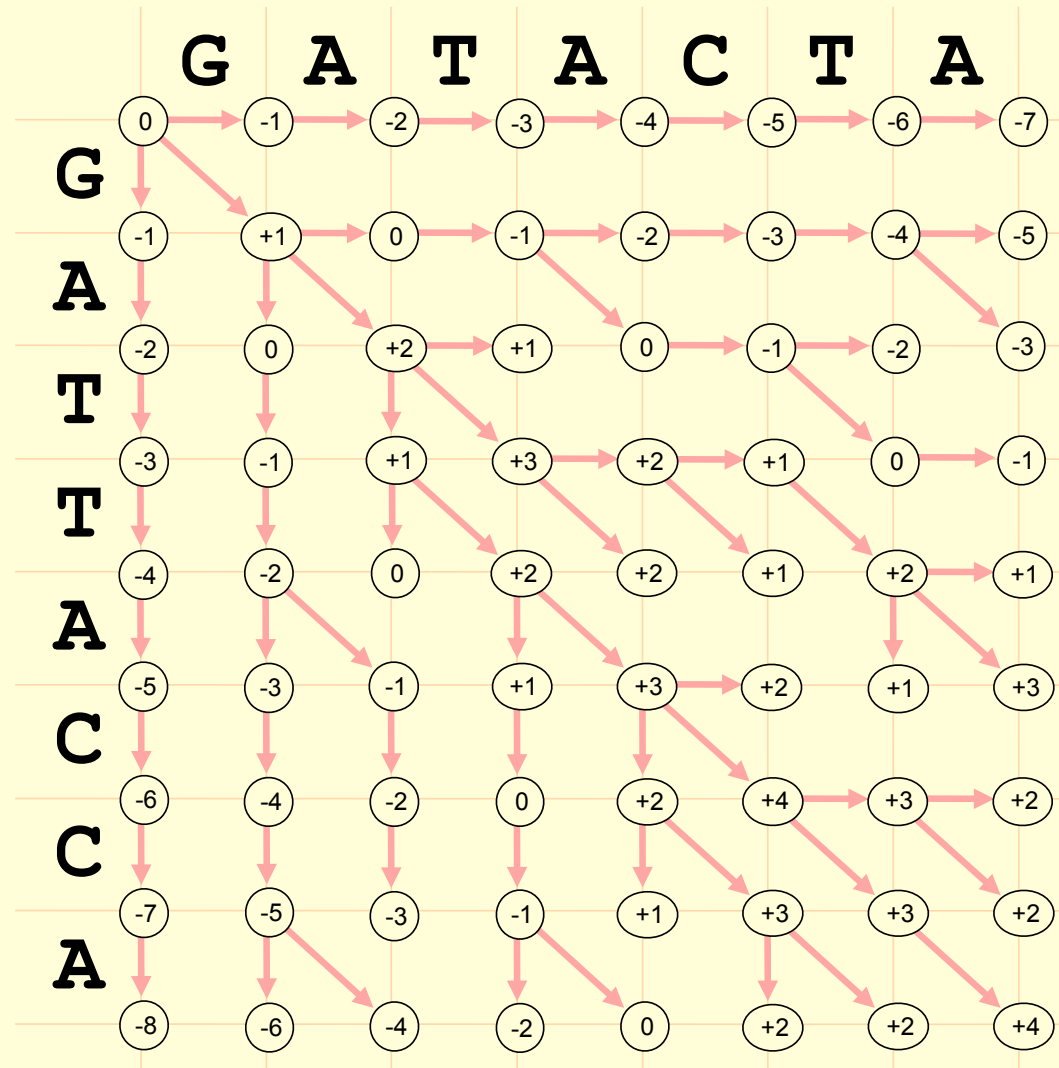
Incrementally extend  
the path

Remember the best  
sub-path leading to  
each point on the  
lattice

Match: +1

Mismatch: -1

Gap: -1



# Dynamic programming

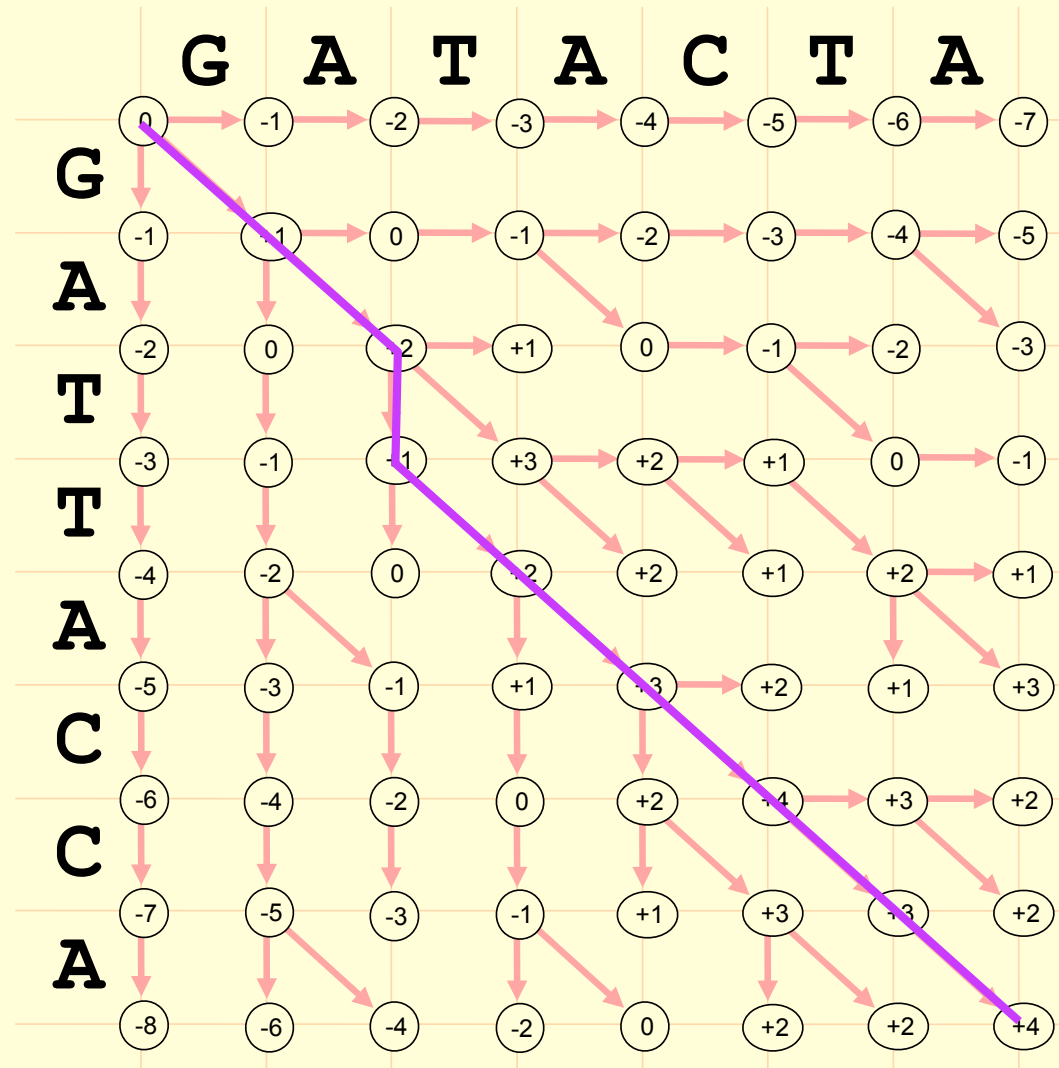
Incrementally extend  
the path

Remember the best  
sub-path leading to  
each point on the  
lattice

Match: +1

Mismatch: -1

Gap: -1

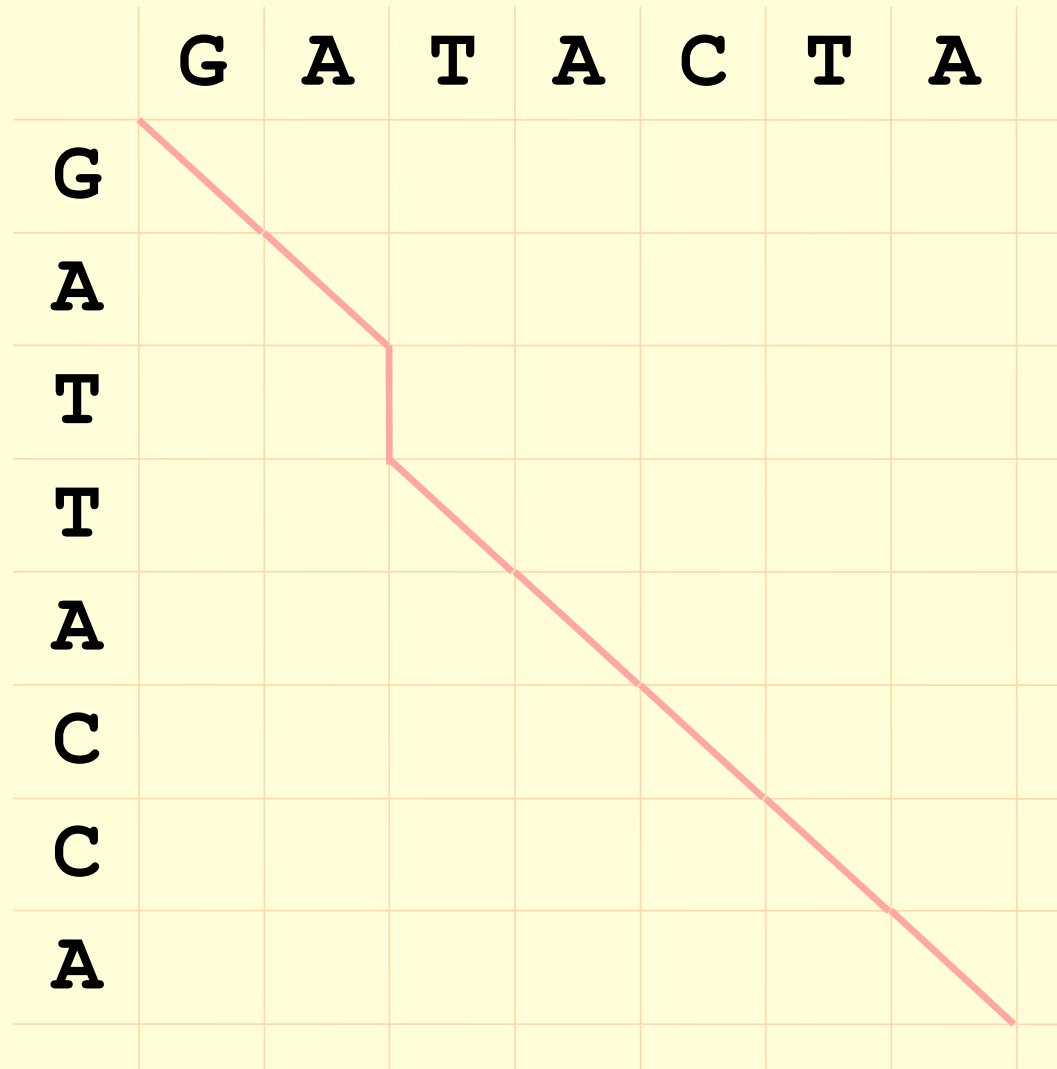




# Dynamic programming

Print out the  
alignment

**GA-TACTA**  
**GATTACCA**



# Dynamic programming

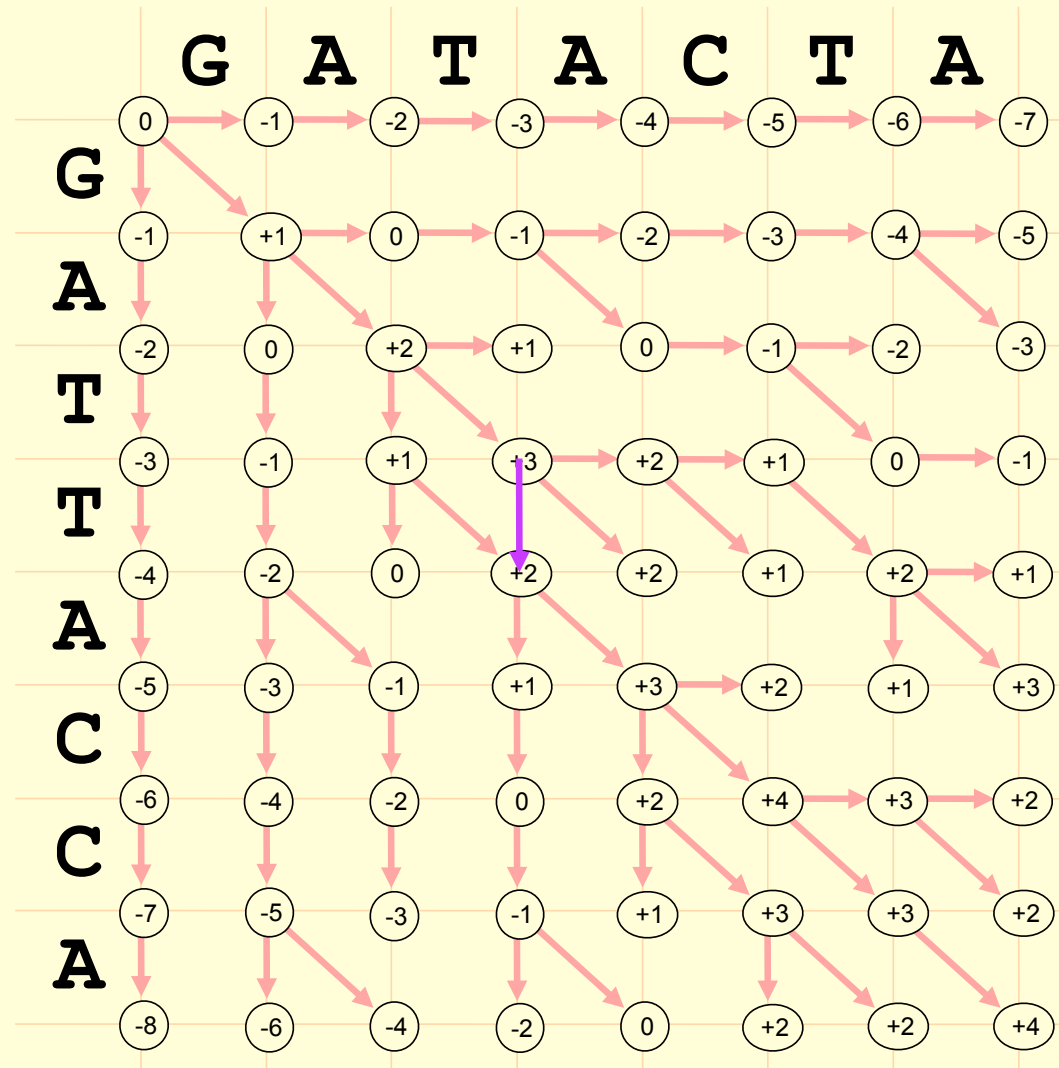
Incrementally extend  
the path

Remember the best  
sub-path leading to  
each point on the  
lattice

Match: +1

Mismatch: -1

Gap: -1





# Dynamic programming

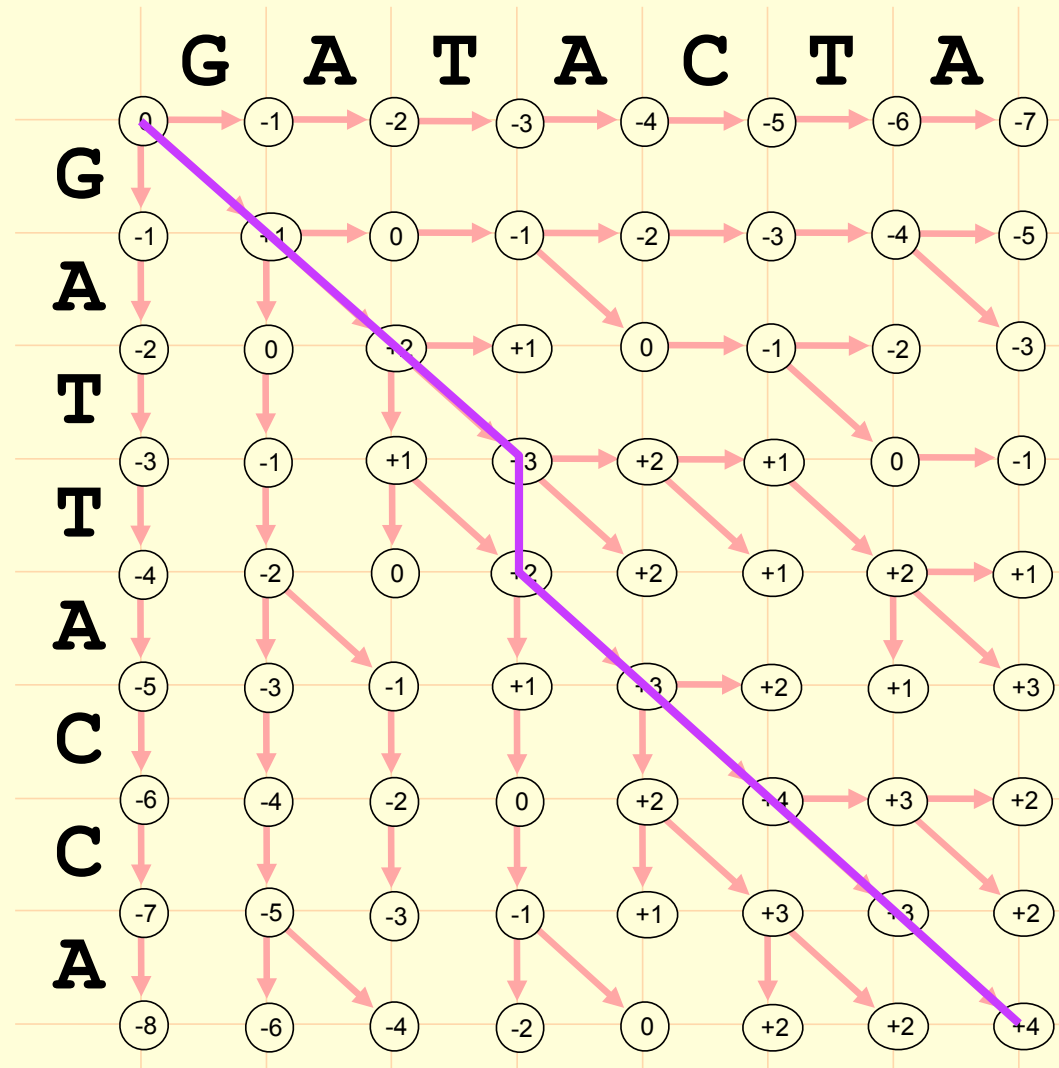
Incrementally extend  
the path

Remember the best  
sub-path leading to  
each point on the  
lattice

Match: +1

Mismatch: -1

Gap: -1



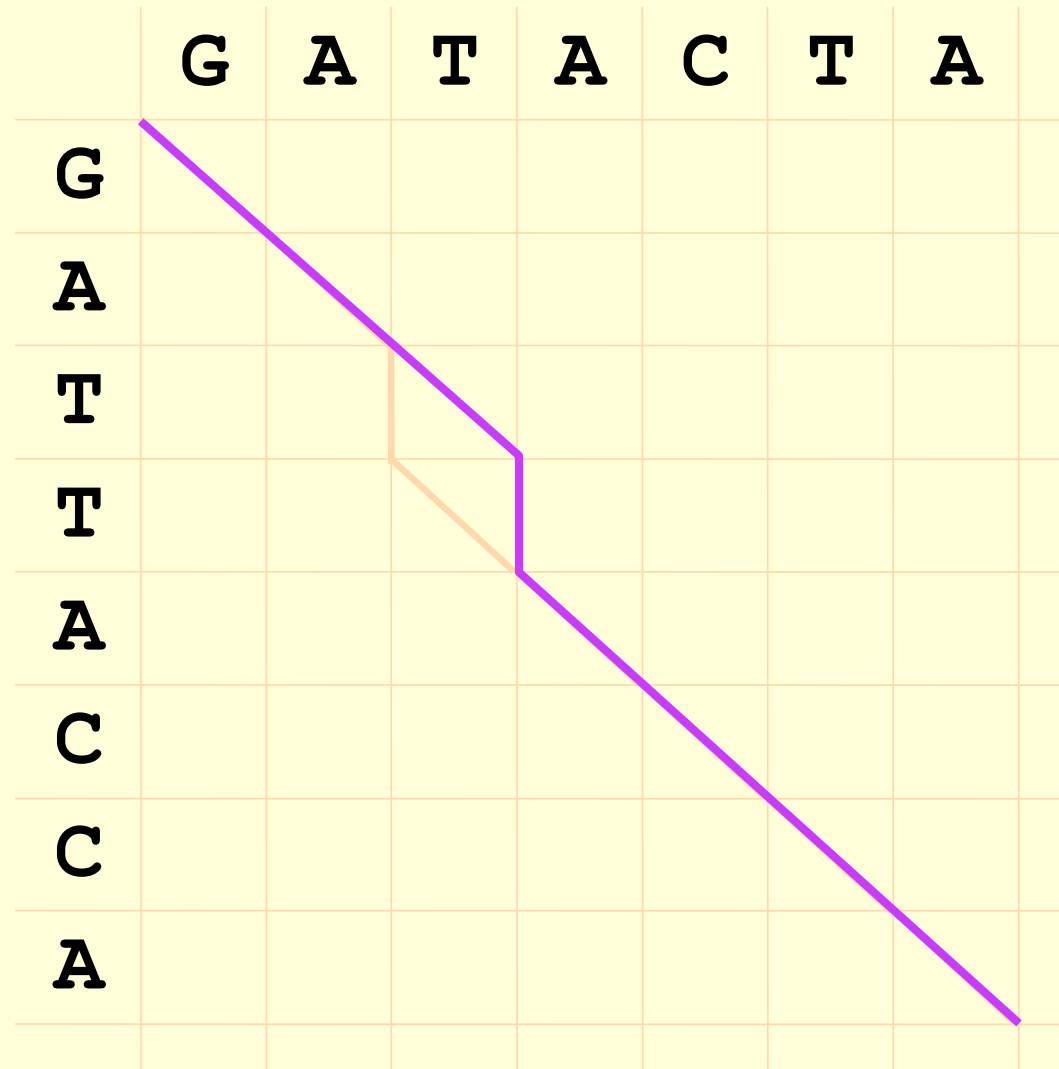
# Dynamic programming

Print out the  
alignment

**GA-TACTA**  
**GATTACCA**

**GAT-ACTA**  
**GATTACCA**

Both alignments are  
optimal - give the same  
max. score





# SEQUENCE SIMILARITY SEARCH



# BASICS OF DATABASE SEARCH

- Database searching is fundamentally different from alignment
- The goal is to find homologous sequences (often more than one), not to establish the correct one-to-one mapping of particular residues
- Usually, this is a necessary first step to making an information map between two sequences
- Database searching programs were originally thought of as approximations to dynamic programming alignments
- Assumption: the best database search conditions are those that would produce the “correct” alignment
- Key idea - most sequences don't match. If one can find a fast way to eliminate sequences that don't match, the search will go much faster



# BASICS OF DATABASE SEARCH

basic terminology:

**query** - sequence to be used for the database search

**subject** - sequence found in the database that meets some similarity criteria

**hit** - local alignment between query and subject



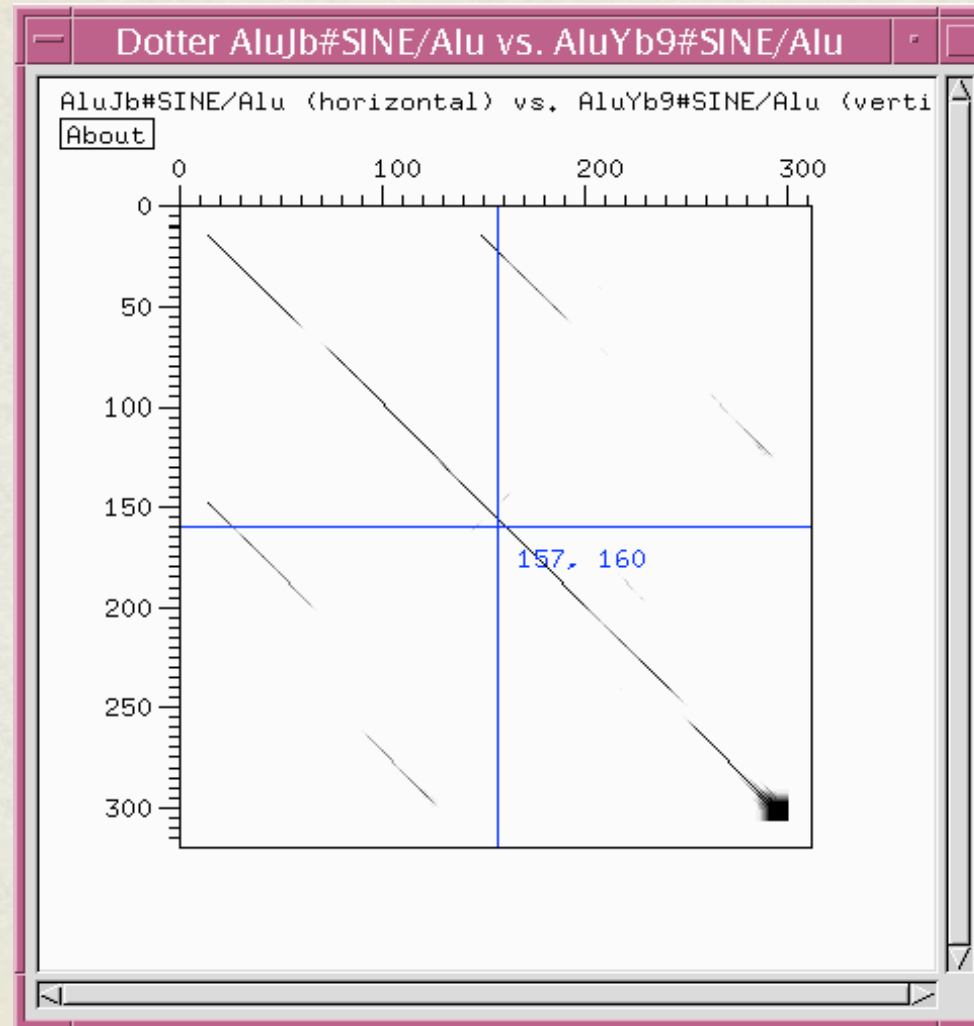
# BASICS OF DATABASE SEARCH

Through the influence of BLAST and FASTA, database searching programs have converged to a basic format

- a. a graphical depiction of the results
- b. a list of top scoring sequences from the databases
- c. a series of alignments for some of the top scoring sequences



# Related sequences have "diagonals" with high similarity





# BLAST

## Basic Local Alignment Search Tool

### References:

Altschul, S.F., Gish, W., Miller, W., Myers, E.W. & Lipman, D.J. (1990) "Basic local alignment search tool." J. Mol. Biol. 215:403-410.

Altschul, S.F., Madden, T.L., Schäffer, A.A., Zhang, J., Zhang, Z., Miller, W. & Lipman, D.J. (1997) "Gapped BLAST and PSI-BLAST: a new generation of protein database search programs." Nucleic Acids Res. 25:3389-3402



# NUCLEOTIDE BLAST ALGORITHM

1. Break down query sequence into overlapping words.
2. Scan databases for exact matches of size  $W$  (BLASTn) or 110110 pattern (MegaBlast).
3. Try to extend the word matches into the complete maximal scoring pair (MSP).  
Significance is easily calculated from Karlin-Altschul equation.
4. Perform local dynamic programming alignment around MSP regions



# BLAST - Maximal Segment Pairs (MSP)

Highest scoring pair of identical length segments from two sequences

Local alignment without gaps

Expected distribution is known!

0121000123456567656543210  
TGCAATCGATCGTCGTCCGTATACA  
:: ::::: :: :  
AGCTCGTGATCGTGGTGGGATCGGT

← running sum  
match = +1  
mism. = -1

potential MSP



# BLAST - extend word matches

Most expensive step in BLAST algorithm

Extend to end of high scoring segment pair, or HSP. HSPs approximate maximal segment pairs or MSPs. They are only approximate because extension does not continue until running score reaches zero - drop off value concept.

After initial hit was found BLAST tries so called extension - an alignment is extended until the maximum value of the score drops by  $x$ , hence name  $x$  dropoff value



# PROTEIN BLAST ALGORITHM

- Break down query sequence into overlapping words and create a lookup table.
- For each word, determine a neighborhood of words that, if found in another sequence, would likely to be part of a significant maximum scoring pair (MSP).
- Scan databases for neighborhood words.
- If two words are found on the same diagonal within a specified distance, try to extend the word matches into the complete MSP. Significance is (relatively) easy calculated from Karlin-Altschul equation.
- Perform local dynamic programming alignment around MSP regions
- first step of BLASTp is controlled by three parameters and a score matrix
- $w$  - word length (k-tuple in FASTA terminology); default value is 3 (lowest possible is 2); two words on the same diagonal are required
- $f$  - score threshold; unlike FASTA BLAST allows mismatches at this step but overall score of the "mini-alignment" has to be above the threshold - the concept of "neighborhood words"



# BLASTp - neighborhood words

## Example - ITV triplet

	BLOSUM62	PAM230
ITV - ITV	$4+5+4 = 13$	$5+3+5 = 13$
ITV - MTV	$1+5+4 = 10$	$2+3+5 = 10$
ITV - ISV	$4+1+4 = 9$	$2+3+5 = 10$
ITV - LTV	$2+5+4 = 11$	$2+3+5 = 10$
ITV - LSV	$2+1+4 = 7$	$2+3+5 = 10$
ITV - MSV	$1+1+4 = 6$	$2+3+5 = 10$
ITV - IAV	$4+0+4 = 8$	$5+1+5 = 11$
ITV - MAV	$1+0+4 = 5$	$2+1+5 = 8$
ITV - ITL	$4+5+1 = 10$	$5+3+2 = 10$
ITV - LAV	$2+0+4 = 6$	$2+1+5 = 8$

# BLASTp - neighborhood words

Threshold  $f = 11$  (default for BLASTp)

$f=10$

	BLOSUM62	PAM230
ITV - ITV	4+5+4 = 13	5+3+5 = 13
ITV - MTV	1+5+4 = 10	2+3+5 = 10
ITV - ISV	4+1+4 = 9	2+3+5 = 10
ITV - LTV	2+5+4 = 11	2+3+5 = 10
ITV - LSV	2+1+4 = 7	2+3+5 = 10
ITV - MSV	1+1+4 = 6	2+3+5 = 10
ITV - IAV	4+0+4 = 8	5+1+5 = 11
ITV - MAV	1+0+4 = 5	2+1+5 = 8
ITV - ITL	4+5+1 = 10	5+3+2 = 10
ITV - LAV	2+0+4 = 6	2+1+5 = 8

	BLOSUM62	PAM230
ITV - ITV	4+5+4 = 13	5+3+5 = 13
ITV - MTV	1+5+4 = 10	2+3+5 = 10
ITV - ISV	4+1+4 = 9	2+3+5 = 10
ITV - LTV	2+5+4 = 11	2+3+5 = 10
ITV - LSV	2+1+4 = 7	2+3+5 = 10
ITV - MSV	1+1+4 = 6	2+3+5 = 10
ITV - IAV	4+0+4 = 8	5+1+5 = 11
ITV - MAV	1+0+4 = 5	2+1+5 = 8
ITV - ITL	4+5+1 = 10	5+3+2 = 10
ITV - LAV	2+0+4 = 6	2+1+5 = 8

Pairs marked in blue would initiate an alignment extension



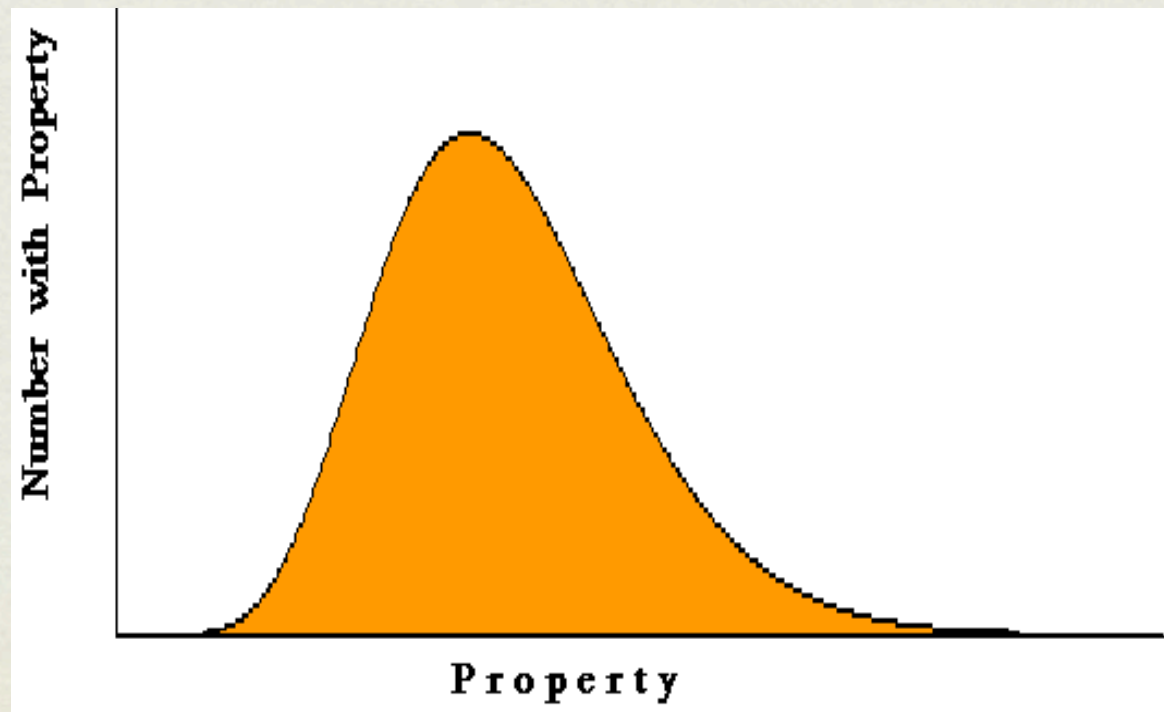
# BLAST - FINAL STEP

- ✂ Smith-Waterman algorithm (local dynamic programming), discussed before but limited to regions that include the HSPs
- ✂ Significance of alignment with gaps can be evaluated using  $K$  and  $\lambda$  estimated from alignments of random sequences with same gap penalty and scoring parameters
- ✂ In spite of claims of being “mathematically rigorous” these parameters can only be estimated empirically



# KARLIN-ALTCHUL STATISTICS

High scores of local alignments between two random sequences follow Extreme Value Distribution





# KARLIN-ALTCHUL STATISTICS

For ungapped alignments their expected number with score  $S$  or greater equals

$$E = Kmne^{-\lambda S}$$

$K$  i  $\lambda$ , are parameters related to a search space and scoring system, and  $m$ ,  $n$  represent a query and database length, respectively.

Score can be transformed to a bit-score according to formula  $S' = \text{bitscore} = (\lambda S - \ln K) / \ln 2$ , then

$$E = mn2^{-S'}$$



# KARLIN-ALTSCHUL STATISTICS

- ⌘• for ungapped alignments parameters  $K$  and  $\lambda$  are calculated algebraically but for gapped alignment a solid theory doesn't exist and these parameters are calculated by simulation which has to be run for every combination of scoring system including gap penalties
- ⌘• therefore not all gap opening and extension score combinations are available
- ⌘• more at <http://www.ncbi.nlm.nih.gov/BLAST/tutorial/Altschul-1.html>



# BLAST - KNOWN PROBLEMS

- ✂ Significance is calculated versus theoretic distribution using Karlin-Altschul equation not real sequences.
- ✂ Assumes sequences are random
- ✂ Assume database is one long sequence – length effects are not corrected for
- ✂ Statistics are very inaccurate for short queries (ca. 20 characters).
- ✂ Be careful when you change BLAST parameters, some of them should be coordinated, e.g. match/mismatch penalty and X-drop off value
- ✂ nucleotide BLAST - default parameters tuned up for speed not sensitivity [Gotea, Veeramachaneni, and Makalowski (2003) Mastering seeds for genomic size nucleotide BLAST searches. Nucleic Acids Res. 31(23):6935-41]



# BLAST ALGORITHM IMPLEMENTATION

Program	Query	Database type
blastn	nt	nt
megablast	nt	nt
blastp	aa	aa
blastx	nt	aa
tblastn	aa	nt
tblastx	nt	aa
blast2seq	nt, aa	nt, aa



# BIOINFORMATICS CREED

- Remember about biology
- Do not trust the data
- Use comparative approach
- Use statistics
- Know the limits
- Remember about biology!!!

