

# Homology

Bio5488

Ting Wang

1/26/26, 1/28/26

.....ACGTTGCCACTTTCGGGGCCACCTGGCCACCTTATTTTCGGAAATATACCGGGCCTTTTTT.....

|||||x||||x|||||||  
CTTTCCCGGCCTCCTGGCCA

match: +1

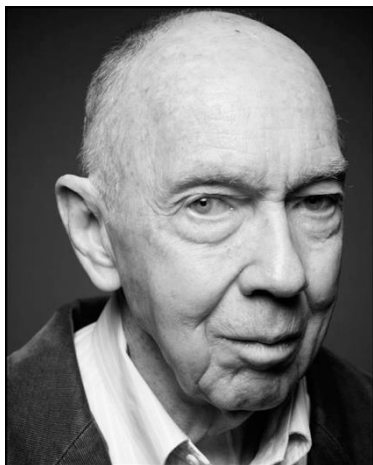
mismatch: -1

matching score = 16

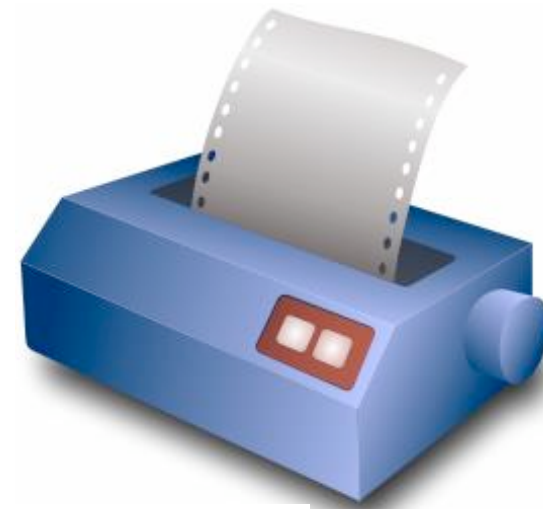
- **How to align them?**
- **Why can we align them?**
- **Why +1 for match, and -1 for mismatch?**
- **What does the score mean?**
- **Is 16 a good score?**

# Outline

- Nobel-prize-worthy work on homology
- What is homology?
- How to detect homology?
- How to quantify homology?
- How to use homology?
- Homology beyond sequence analysis
- Next-gen sequencing alignment



## Russell Doolittle (Bishop and Varmus)



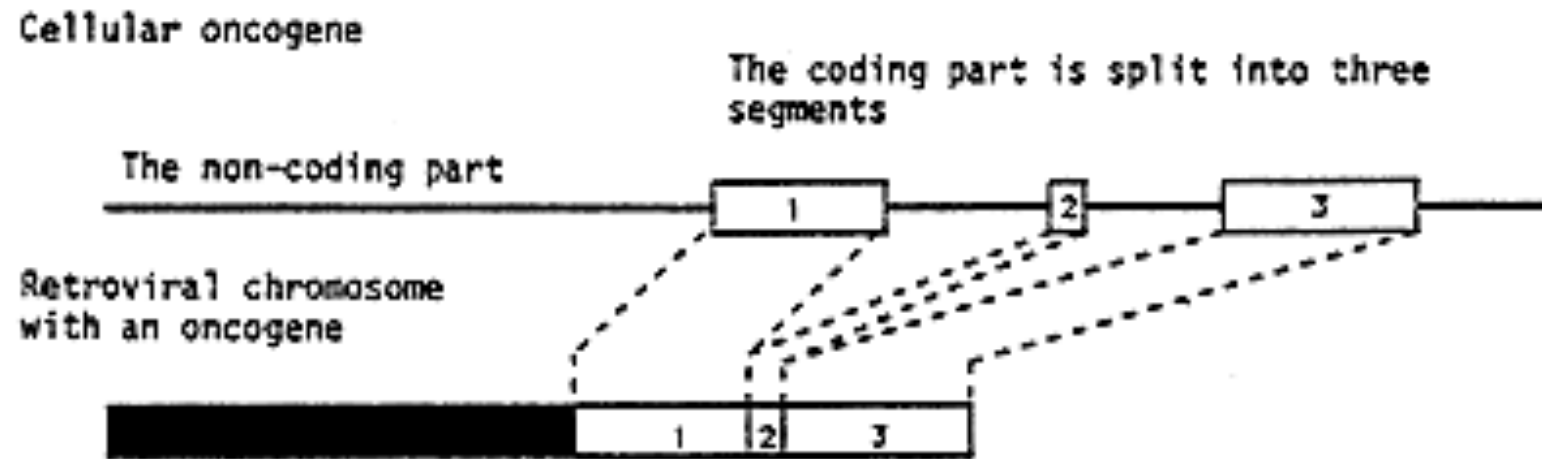
p28sis	1	MTLTWGGDPIPEELYKMLSGHSIARSFDDLQRLLOGDSGKEDGAELDNLMT	50
p28sis	51	RSHSGGELES LARGKRSLSLSVAEPAMIAECKTRTEVF EISARLIIDRTN	100
PDGF-2	1	SLGSLTIAEPAMIAECKTRTEVF CICAL?DR??	34
PDGF-1	1	SIEEAVPAVCKTRIVIVIEISARELD???	28
p28sis	101	ANFLVWPPCVEVGRCSGCCNNRNVQCAPTQVQLAPVQVAKIEIVAKKPIF	150
PDGF-2	35	?????PPCVEVKACTGCCNNRNVKCAPSQVQLAP?QVAKIEIVAK[	80
PDGF-1	29	ANFL [	32
p28sis	151	KKATVTLEDHLACKCEIVAAARAVTRSPGTSQEGRAKTTGSAVTIATVRV	200
PDGF-2			
PDGF-1			
p28sis	201	RRPPKGKHKCKHTHDKTALKETLGA	226
PDGF-2		]	
PDGF-1		]	

Simian Sarcoma Virus onc Gene, v-sis, is Derived from the Gene (or Genes) Encoding a Platelet-Derived Growth Factor

Author(s): Russell F. Doolittle, Michael W. Hunkapiller, Leroy E. Hood, Sushilkumar G. Devare, Keith C. Robbins, Stuart A. Aaronson, Harry N. Antoniades

Source: *Science*, New Series, Vol. 221, No. 4607 (Jul. 15, 1983), pp. 275-277

## Bishop and Varmus strategy (Nobel prize 1989)



## Doolittle strategy (could be the first Nobel prize for computational biology)

tween these proteins. This similarity was discovered by one of us (R.F.D.) during a search for sequence homology between the PDGF amino-terminal sequences and the other protein sequences in the Newat sequence data base at the University of California, San Diego (19). Subsequent-

base searched included 145,581 amino acid residues comprising 684 individual sequences in the Newat list and 121,098 residues from 1081 sequences in the 1978 Dayhoff collection [*Protein*

**What is the significance?**

# The Nucleotide Sequence of *Saccharomyces cerevisiae*

## 5.8 S Ribosomal Ribonucleic Acid

(Received for publication, November 20, 1972)

GERALD M. RUBIN\*

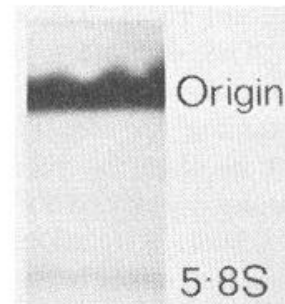
*From the Medical Research Council Laboratory of Molecular Biology, Cambridge, CB2 2QH, England*

### SUMMARY

The nucleotide sequence of *Saccharomyces cerevisiae* 5.8 S ribosomal RNA (also known as the 7 S or 18S RNA species) has been determined to be pApApApCpUpUpUpCpApApCpApApCpGpGpApUpCpUpCpUpUpGpGpUpUpCpUpCpGpCpApUpCpGpApUpGpApApGpApApCpGpCpApGpCpGpApApApUpGpCpGpApUpApCpGpUpApApUpGpUpGpApAΨpUpGpCpApGpApApUpUpCpCpGpUpGpApApUpCpApUpCpGpApApUpCpUpUpUpGpApApCpGpCpApCpApUpUpGpCpGpCpCpCpUpUpGpGpUpApUpUpCpCpApGpGpGpGpGpCpApUpGpCpCpUpGpUpUpUpGpApGpCpGpUpCpApUpUpU.

Ribosomes from the cytoplasm of eukaryotic cells contain two low molecular weight RNA species: the 5 S and the 5.8 S RNA (1-9). Both RNA species are structural components of the 18S subunit and each is found in equimolar amount to the 28S

*Low Phosphate Medium*—Inorganic phosphate was precipitated (as  $\text{MgNH}_4\text{PO}_4$ ) from 10% Bacto-yeast extract and 20% Bacto-peptone by the addition of 10 ml of 1 M  $\text{MgSO}_4$  and 10 ml of concentrated aqueous ammonia per liter. The phosphates were allowed to precipitate at room temperature for 30 min, and the precipitate was removed by filtration through Whatman No. 1 filter paper. The filtrate was adjusted to pH 5.8 with HCl and autoclaved. Sterile glucose was added to a final concentration of 2%.



# A few Definitions

**Homologs:** genes/sequences sharing a common origin

**Orthologs:** genes originating from a single ancestral gene in the last common ancestor of the compared genomes; genes related via speciation

**Paralogs:** genes related via duplication

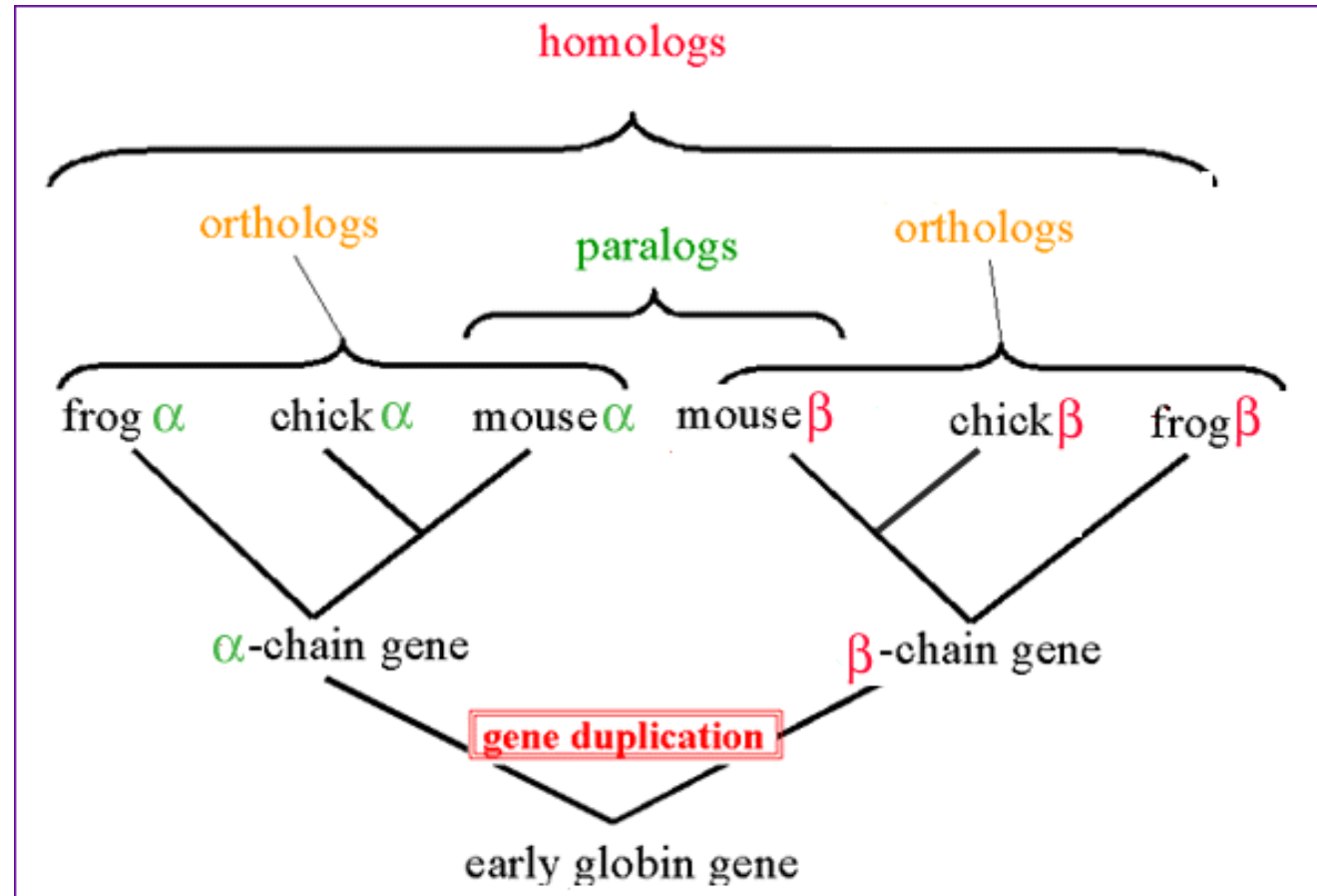
**Xenolog:** sequences that have arisen out of horizontal transfer events (symbiosis, viruses, etc)

**Co-orthologs:** two or more genes in one lineage that are, collectively, orthologous to one or more genes in another lineage due to a lineage-specific duplication(s)

**Outparalogs:** paralogous genes resulting from a duplication(s) preceding a given speciation event

**Inparalogs:** paralogous genes resulting from a lineage-specific duplication(s) subsequent to a given speciation event

# Relation of sequences



*Need ancestral sequences to distinguish orthologs and paralogs*



# Similarity versus Homology

- Similarity refers to the likeness or % identity between 2 sequences
- Similarity means sharing a statistically significant number of bases or amino acids
- **Similarity does not imply homology**
- Similarity can be quantified
- It is ok to say that two sequences are X% identical
- It is ok to say that two sequences have a similarity score of Z
- It is generally **incorrect** to say that two sequences are X% *similar*
- Homology refers to shared ancestry
- Two sequences are homologous if they are derived from a common ancestral sequence
- **Homology usually implies similarity**
- Low complexity regions can be highly similar without being homologous
- Homologous sequences are not always highly similar
- A sequence is either homologous or not.
- **Never say two things are X% homologous**

# Why Compare Sequences?

- Sequence comparisons lie at the heart of all bioinformatics
- Identify sequences
  - What is this thing I just found?
- Compare new genes to known ones
- Compare genes from different species
  - information about evolution
- Guess functions for entire genomes full of new gene sequences
  - Metagenomics
- What does it matter if two sequences are similar or not?
  - Globally similar sequences are likely to have the same biological function or role
  - Locally similar sequences are likely to have some physical shape or property with similar biochemical roles
  - If we can figure out what one does, we may be able to figure out what they all do

# Sequence alignment

- How to optimally align two sequences
  - Dot plots
  - Dynamic programming
    - Global alignment
    - Local alignment
- How to score an alignment
- Fast similar sequence search
  - BLAST
  - BLAT
  - More recent development: short read alignment
- Determine statistical significance
- Using information in multiple sequence alignment to improve sensitivity
- Alignment beyond sequences

## Visual Alignments (Dot Plots)

- Build a comparison matrix
  - Rows: Sequence #1
  - Columns: Sequence #2
- Filling
  - For each coordinate, if the character in the row matches the one in the column, fill in the cell
  - Continue until all coordinates have been examined



# Noise in Dot Plots

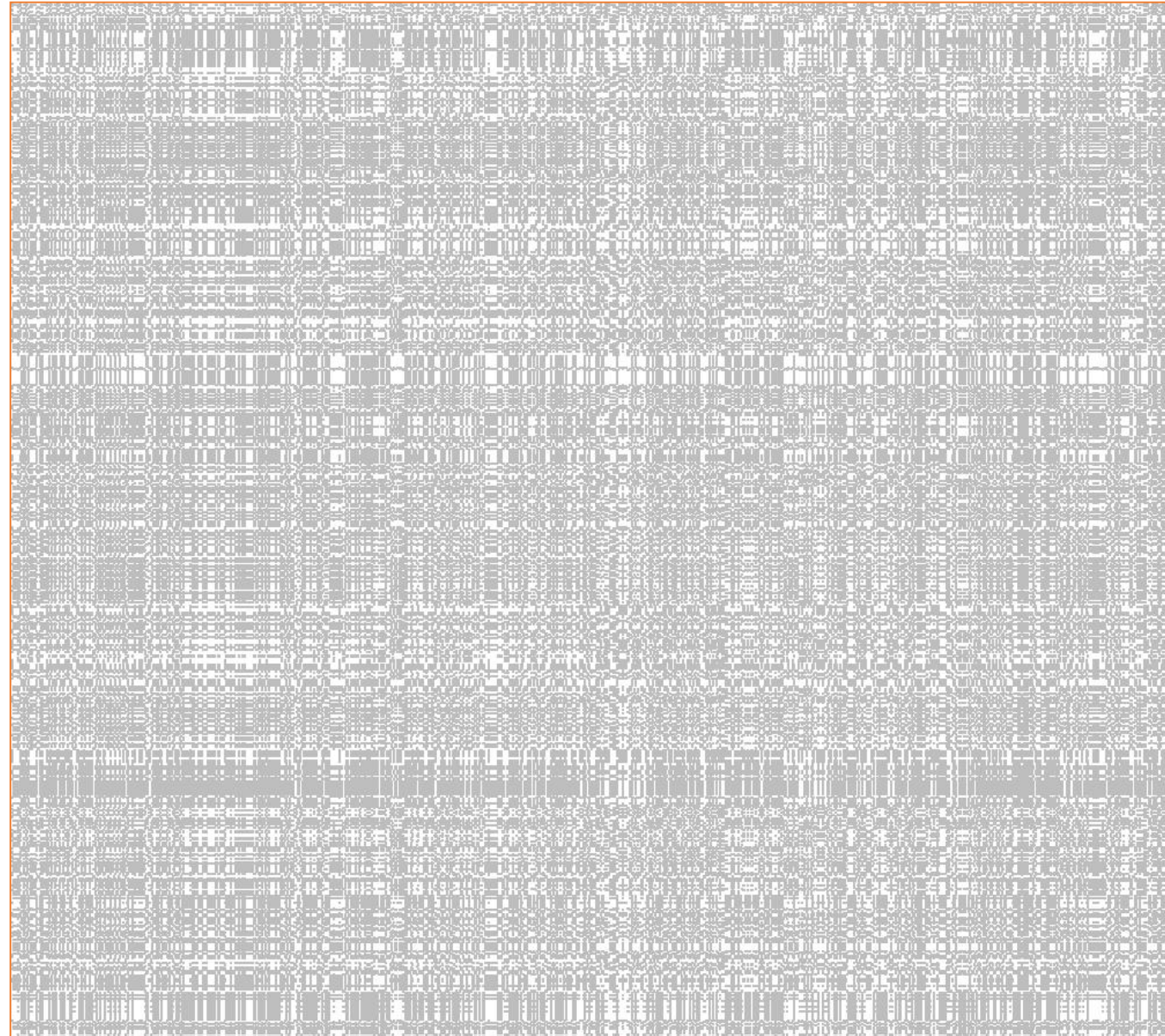
- Nucleic Acids (DNA, RNA)
  - 1 out of 4 bases matches at random
- Windowing helps reduce noise
  - Can require  $>X$  bp match before plotting
  - Percentage of bases matching in the window is set as threshold

Met14 vs  
Met2  
“DotPlot”

MET14 (1000nt)

MET2(895nt)

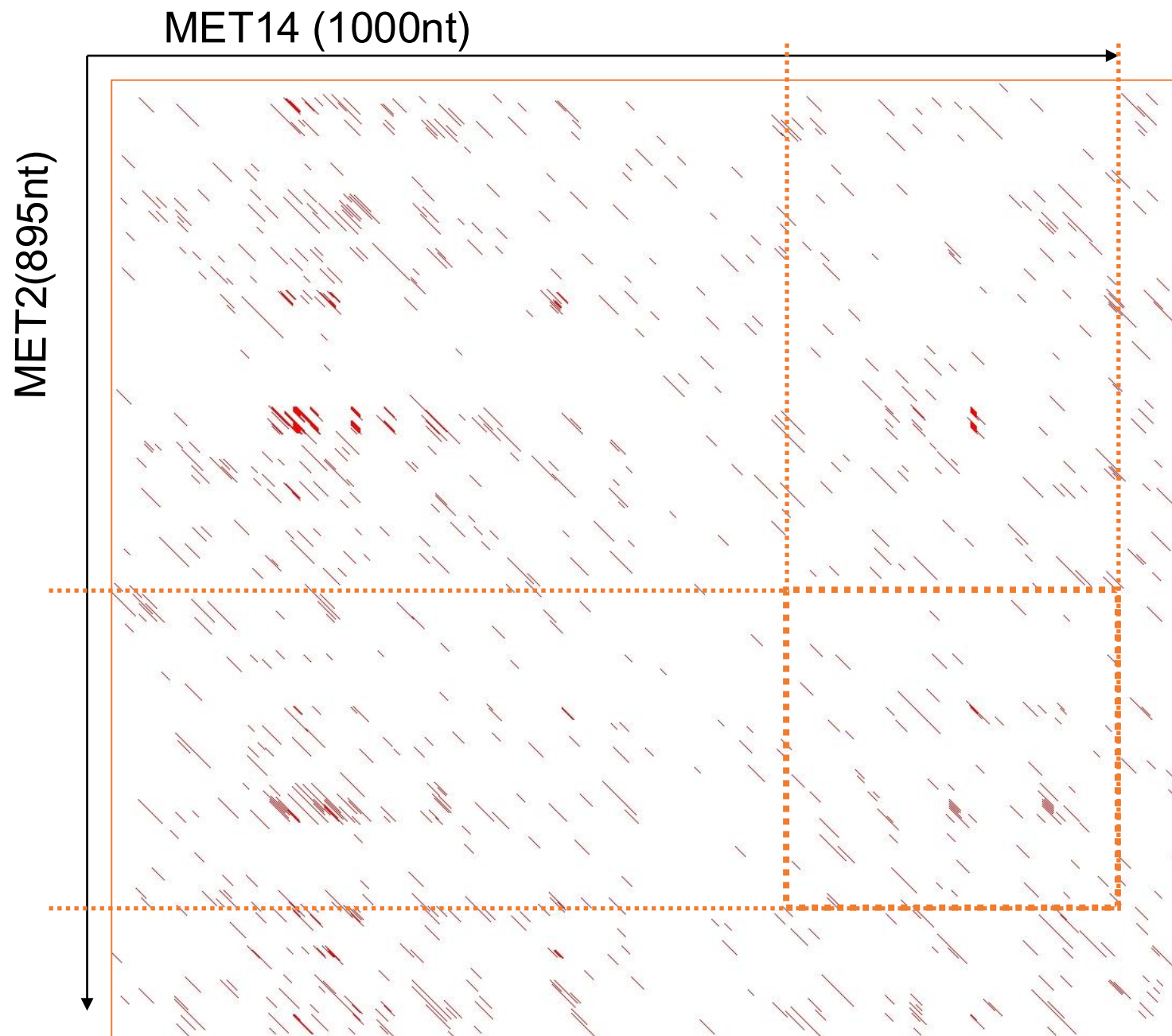
Match = 1  
Mismatch = -1  
Gray: 1





# Met14 vs Met2

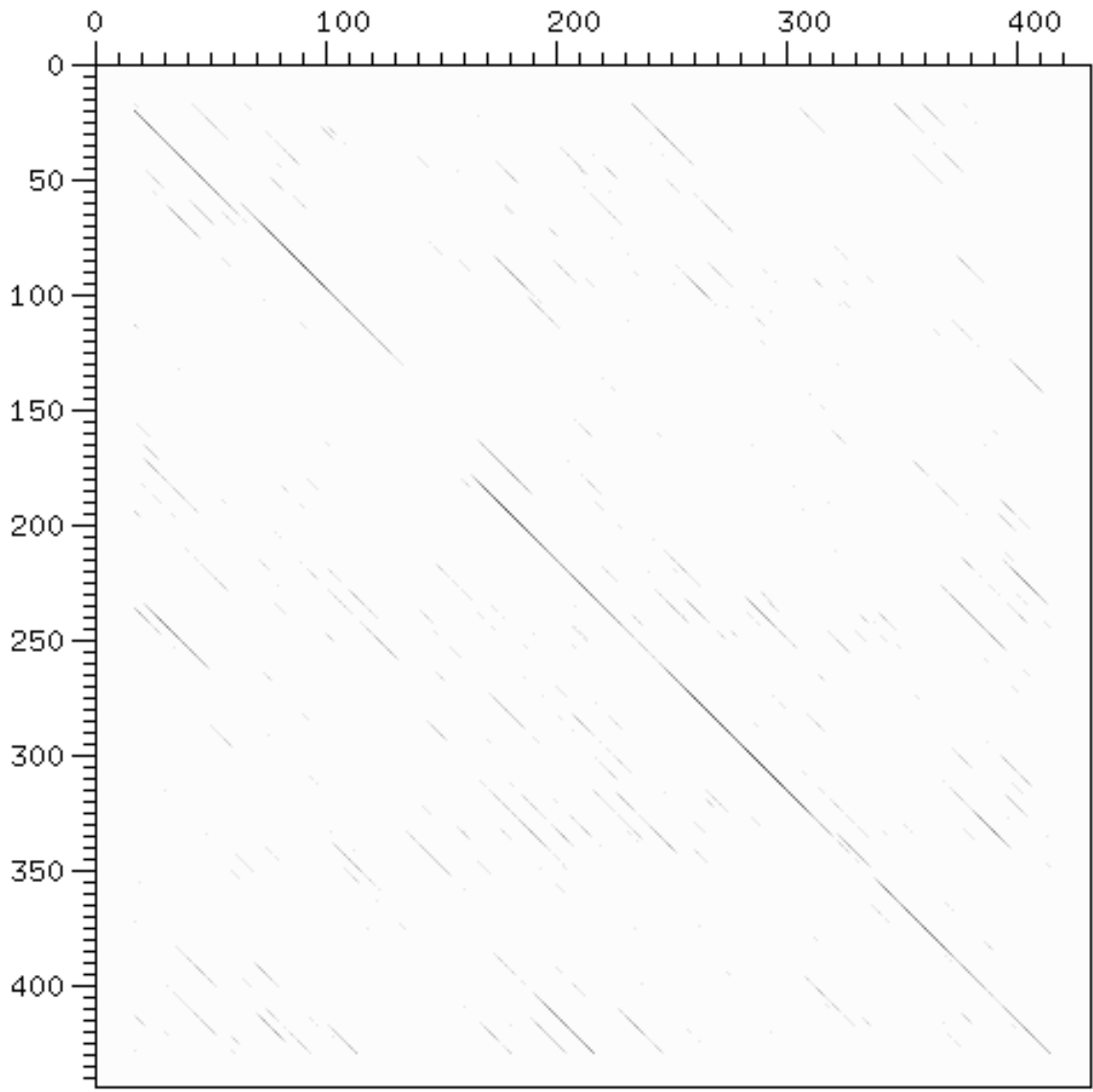
Red: >5



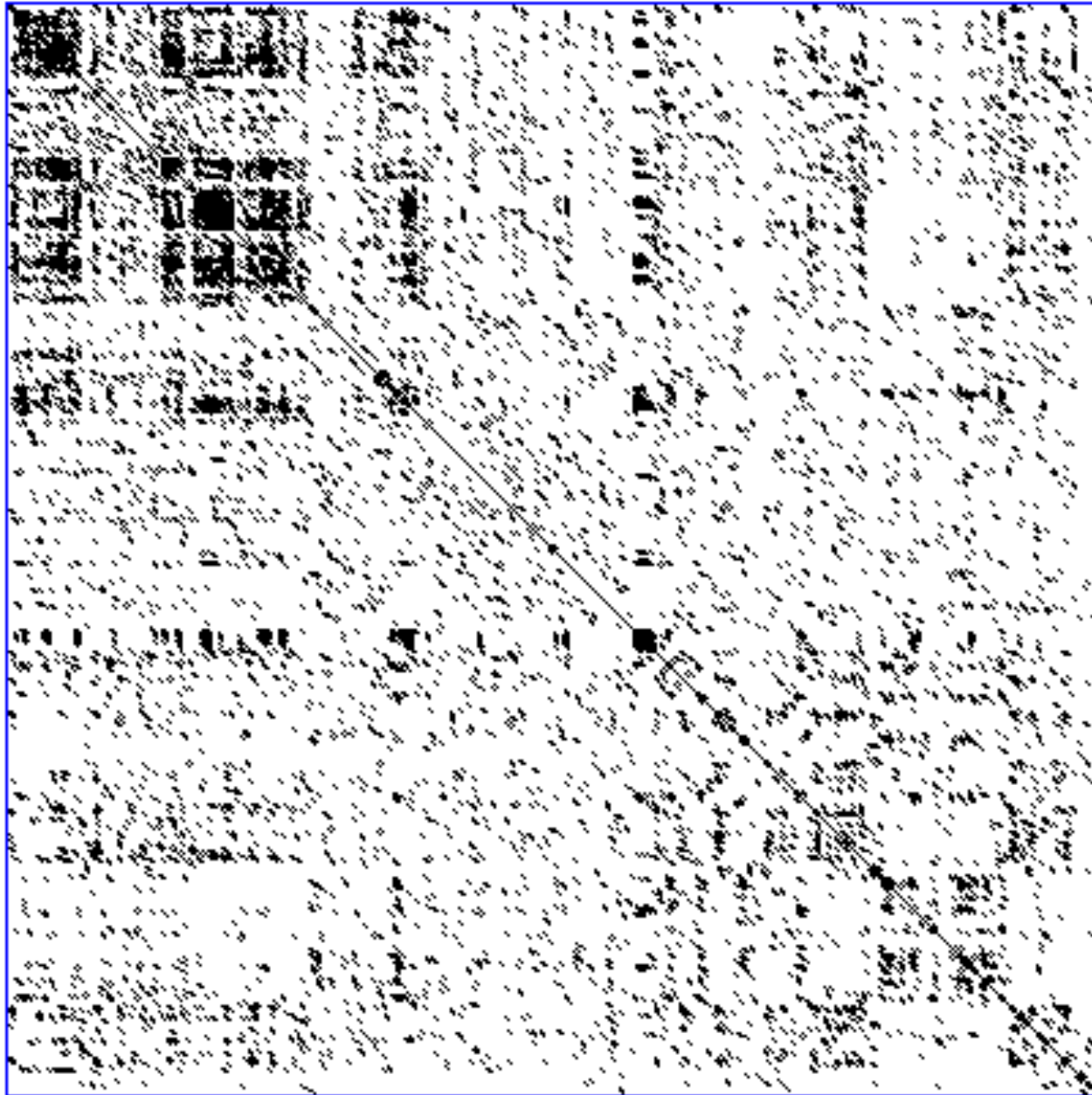


**$\alpha$  chain of human hemoglobin**

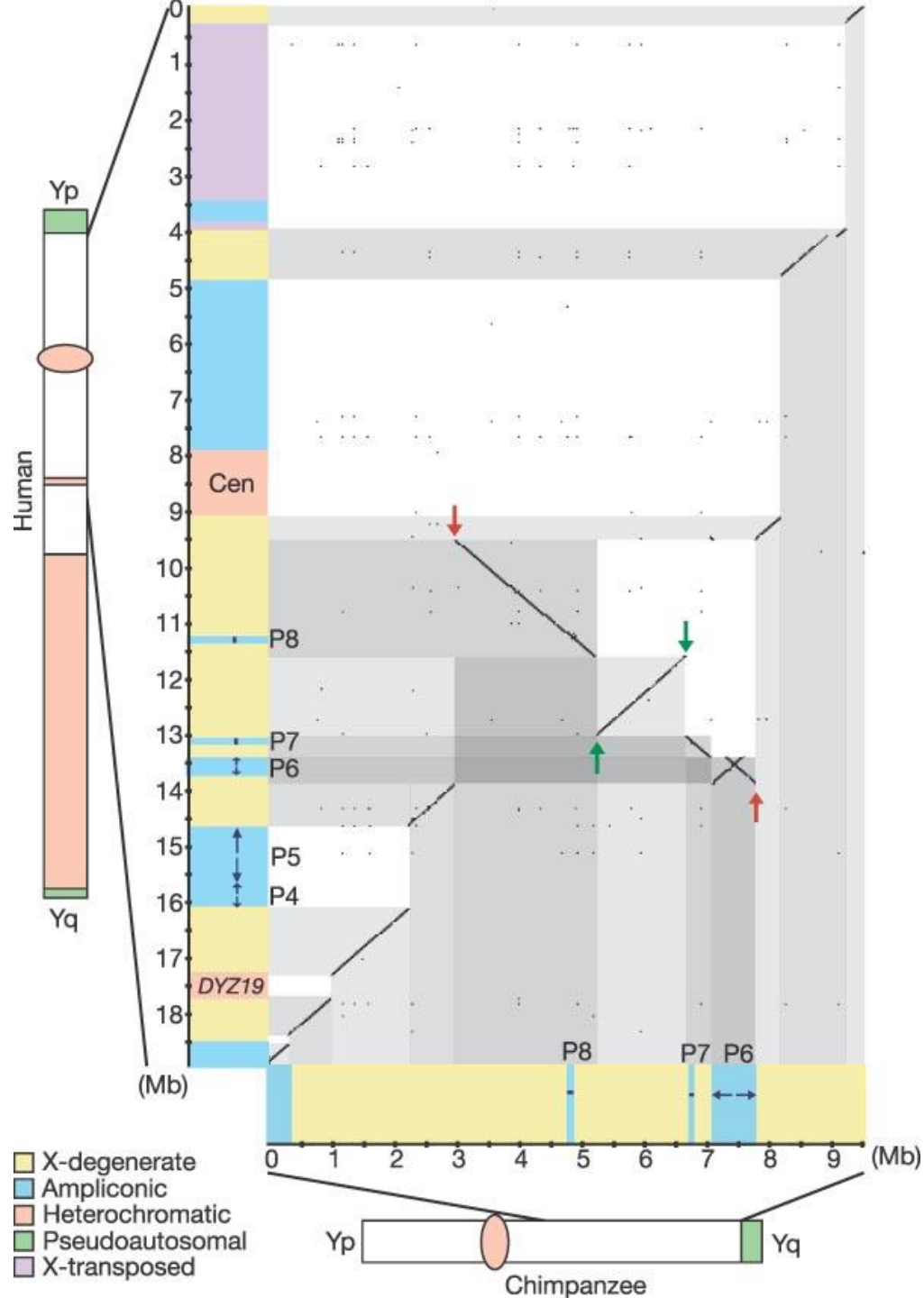
**$\beta$  chain of human hemoglobin**



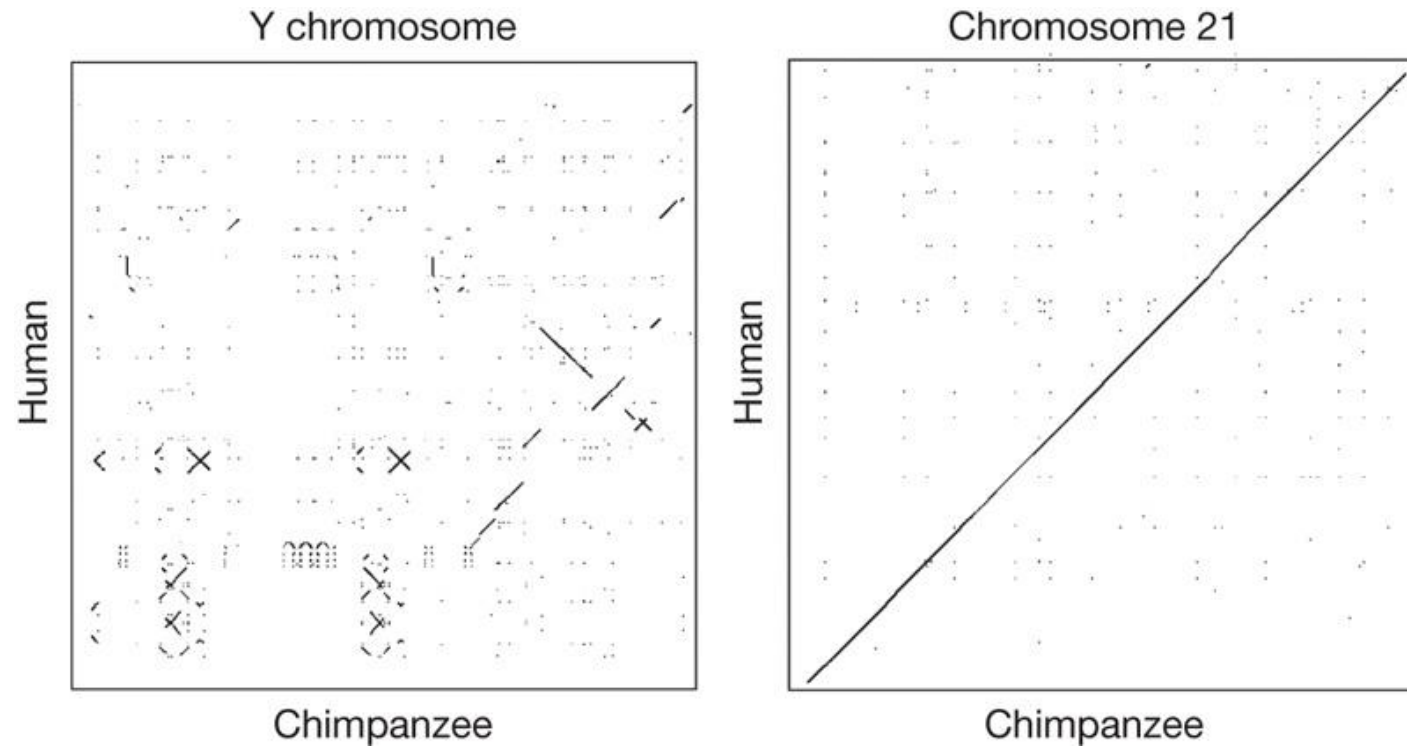
## MAZ: Myc associated zinc finger isoform 1 self alignment



# Human vs Chimp Y chromosome comparison



# Dot plots of DNA sequence identity between chimpanzee and human Y chromosomes and chromosomes 21



JF Hughes *et al. Nature* **000**, 1-4 (2010) doi:10.1038/nature08700

nature

# Aligning sequences by residue

- Match: **award**
- Mismatch (substitution or mutation): **penalize**
- Insertion/Deletion (INDELS – gaps): **penalize** (gap open, gap extension)

A	L	I	G	N	M	E	N	T
-	L	I	G	A	M	E	N	T

## More than one solution is possible

- Which alignment is best?

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

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

# Alignment Scoring Scheme

- Possible scoring scheme:

match: +2

mismatch: -1

indel -2

- Alignment 1:  $5*2 + 1*-1 + 4*-2 = 10 - 1 - 8 = 1$
- Alignment 2:  $6*2 + 1*-1 + 2*-2 = 12 - 1 - 4 = 7$

# Dynamic Programming

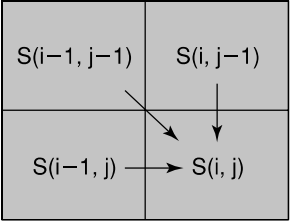
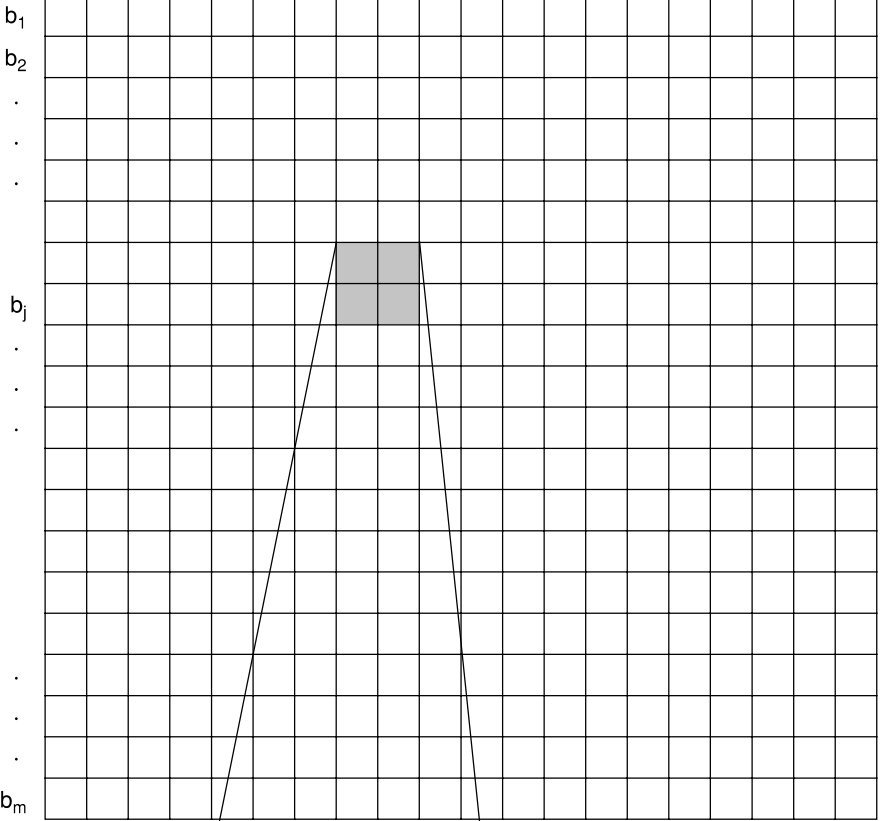
- Global Alignments:
  - Needleman S.B. and Wunsch C.D. (1970) *J. Mol. Biol.* 48, 443-453
- Local Alignments:
  - Smith T.F. and Waterman M.S. (1981) *J. Mol. Biol.* 147, 195-197
  - One simple modification of Needleman/Wunsch: when a value in the score matrix becomes negative, reset it to zero (begin of new alignment)
- Guaranteed to be mathematically optimal:
  - Given two sequences (and a scoring system) these algorithms are guaranteed to find the very best alignment between the two sequences!
- Slow  $N^2$  algorithm
- Performed in 2 stages
  - Prepare a scoring matrix using recursive function
  - Scan matrix diagonally using traceback protocol



Seq A:

$a_1$   $a_2$     $\cdot$     $\cdot$     $\cdot$     $a_i$     $\cdot$     $\cdot$     $\cdot$     $\cdot$     $\cdot$     $\cdot$     $a_n$

Seq B:



$$S(i,j) = \max \begin{cases} S(i-1, j-1) + \text{score}(a_i, b_j) \\ S(i, j-1) + \delta \\ S(i-1, j) + \delta \end{cases}$$

# Dynamic Programming

	G	E	N	E	T	I	C	S
G	10	0	0	0	0	0	0	0
E	0	10	0	10	0	0	0	0
N	0	0	10	0	0	0	0	0
E	0	0	0	10	0	0	0	0
S	0	0	0	0	0	0	0	10
I	0	0	0	0	0	10	0	0
S	0	0	0	0	0	0	0	10

	G	E	N	E	T	I	C	S
G	60	40	30	20	20	0	10	0
E	40	50	30	30	20	0	10	0
N	30	30	40	20	20	0	10	0
E	20	20	20	30	20	10	10	0
S	20	20	20	20	20	0	10	10
I	10	10	10	10	10	20	10	0
S	0	0	0	0	0	0	0	10

G	E	N	E	T	I	C	S
				*			
G	E	N	E	S	I		S

# DP (demo)

- Match=5, mismatch=-3, gap=-4

-	G	A	T	C	A	G	T	T	A
-	0	0	0	0	0	0	0	0	0
G	0								
G	0								
A	0								
T	0								
C	0								
G	0								
A	0								

# DP (demo)

$$S_{1,1} = \text{MAX}\{S_{0,0} + 5, S_{1,0} - 4, S_{0,1} - 4, 0\} = \text{MAX}\{5, -4, -4, 0\} = 5$$

[illegible]

# DP (demo)

$$S_{1,2} = \text{MAX}\{S_{0,1} - 3, S_{1,1} - 4, S_{0,2} - 4, 0\} = \text{MAX}\{0 - 3, 5 - 4, 0 - 4, 0\} = \text{MAX}\{-3, 1, -4, 0\} = 1$$

[illegible]

# DP (demo)

$$S_{1,3} = \text{MAX}\{S_{0,2} - 3, S_{1,2} - 4, S_{0,3} - 4, 0\} = \text{MAX}\{0 - 3, 1 - 4, 0 - 4, 0\} = \text{MAX}\{-3, -3, -4, 0\} = 0$$

[illegible]

## Trace Back (Local Alignment)

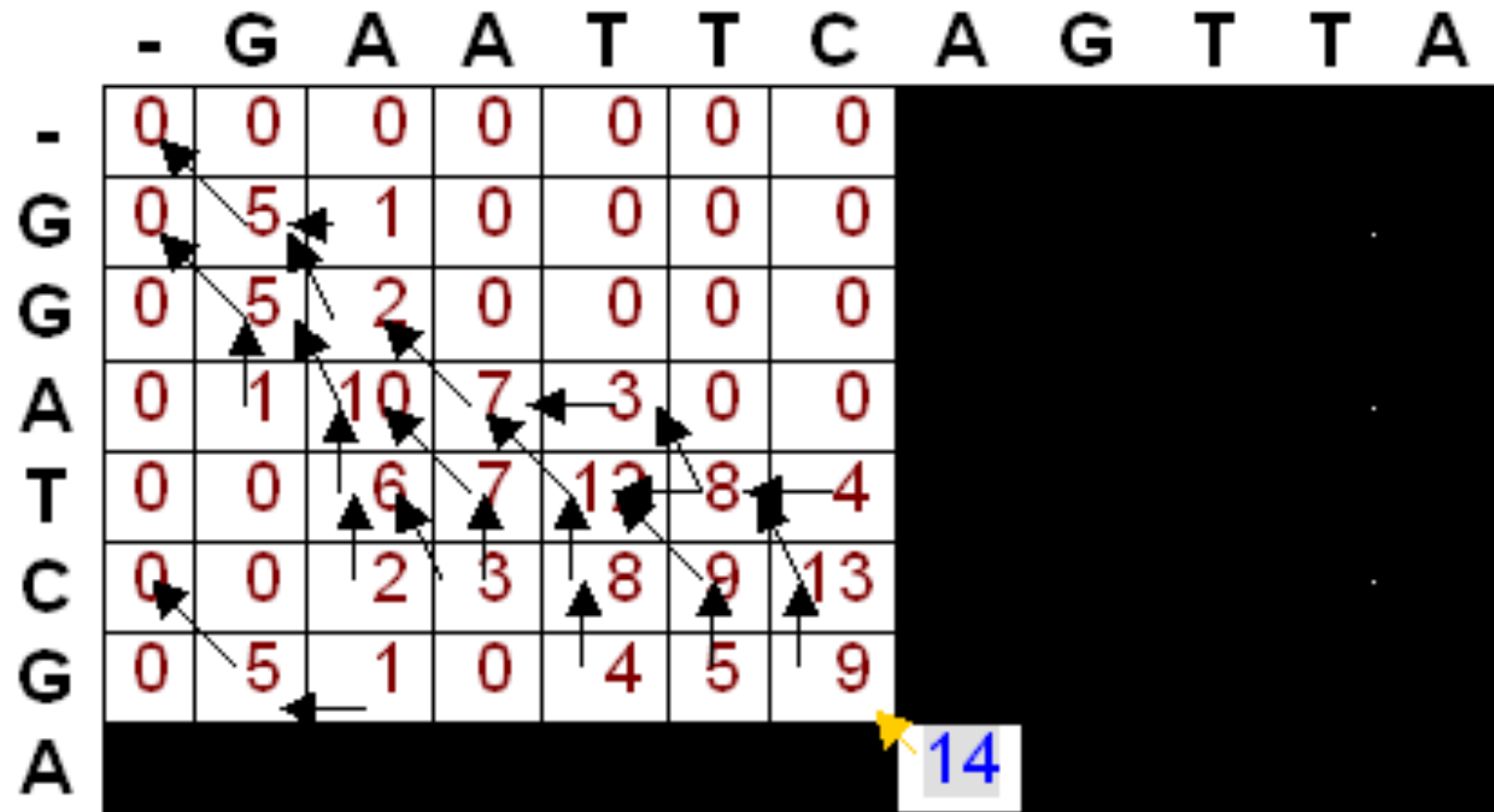
- Maximum local alignment score is the highest score anywhere in the matrix (14 in this example)
- 14 is found in two separate cells, indicating two possible multiple alignments producing the maximal local alignment score

## Trace Back (Local Alignment)

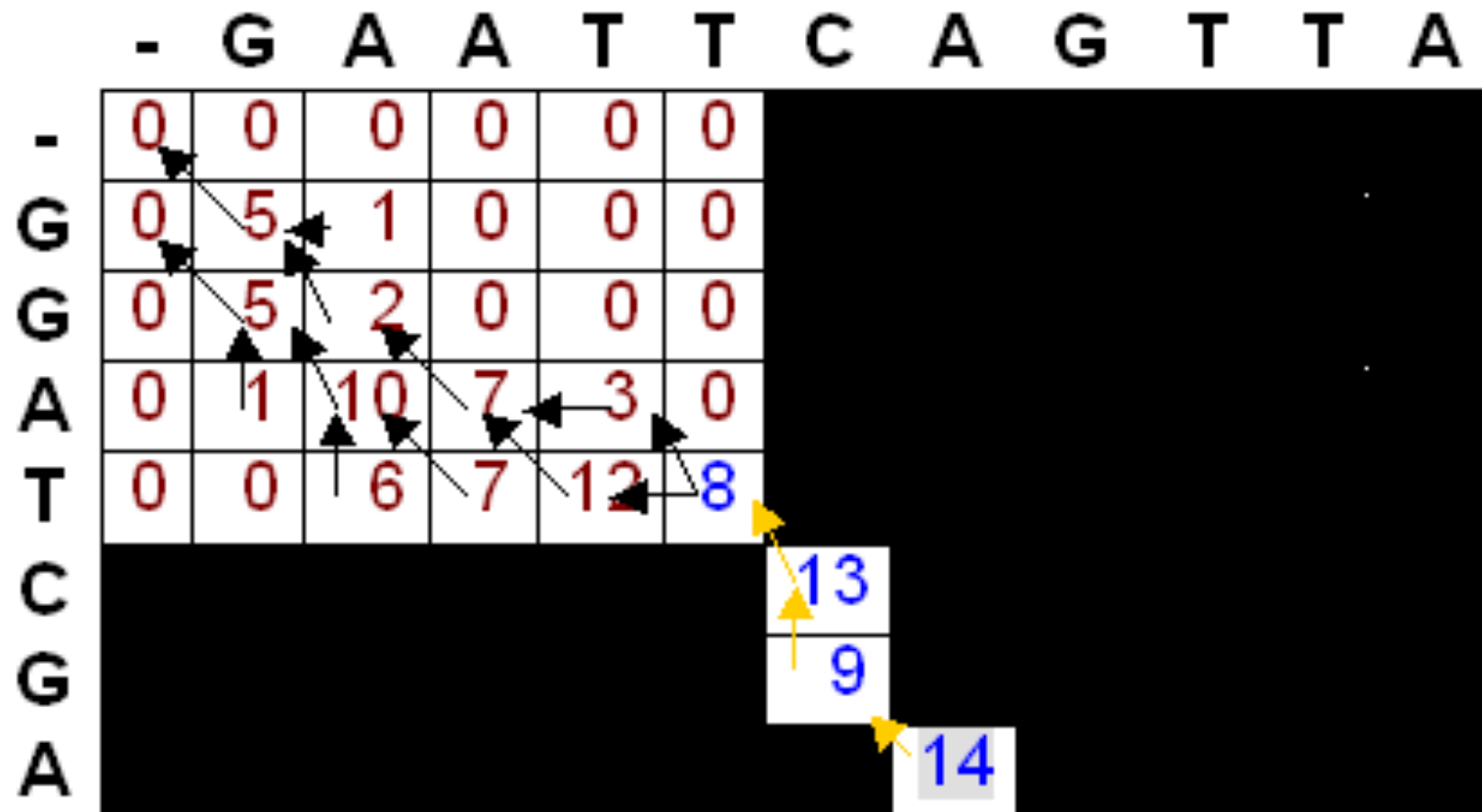
- Trace-back begins in the position with the highest value.
- At each cell, we look to see where we move next according to the pointers
- When a cell is reached where there is not a pointer to a previous cell, we have reached the beginning of the alignment



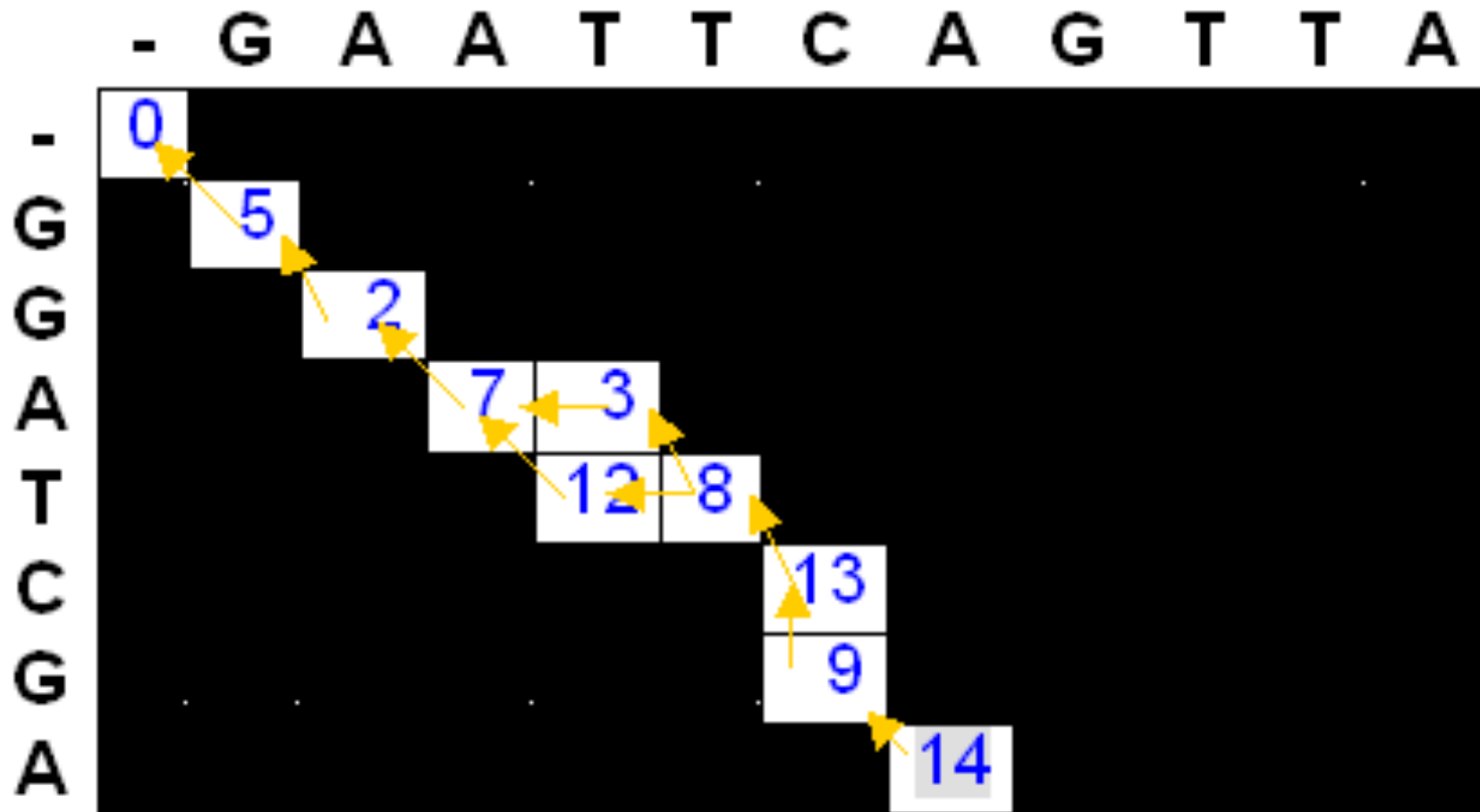
## Trace Back Demo



## Trace Back Demo



## Trace Back Demo



## Maximum Local Alignment

G A A T T C - A

|       | |       |       |

G G A T - C G A

+ - + + - + - +

5 3 5 5 4 5 4 5

**=14**

G A A T T C - A

|       |       | |       |

G G A - T C G A

+ - + - + + - +

5 3 5 4 5 5 4 5

**=14**

## Linear vs. Affine Gaps

- So far, gaps have been modeled as linear
- More likely contiguous block of residues inserted or deleted
  - 1 gap of length  $k$  rather than  $k$  gaps of length 1
- Can create scoring scheme to penalize big gaps relatively less
  - Biggest cost is to open new gap, but extending is not so costly

## Affine Gap Penalty

$$w_x = g + r(x-1)$$

- $w_x$  : total gap penalty
  - $g$ : gap open penalty
  - $r$ : gap extend penalty
  - $x$ : gap length
- 
- gap penalty chosen relative to score matrix

# Scoring Alignments

- Pick a scoring matrix
  - BLOSUM62
  - PAM250
  - Match=5, mismatch=-4
- Decide on gap penalties
  - -gap opening penalty (-8)
  - -gap extension penalty (-1)
- Assume every position is independent
- Sum scores at each position
  - $[\log(x*y)=\log x+\log y]$

## Scoring Matrices

$$S_{ij} = \frac{\log\left(\frac{q_{ij}}{p_i p_j}\right)}{\lambda}$$

- An empirical model of evolution, biology and chemistry all wrapped up in a 20 X 20 (or 4 X 4) table of numbers
- Structurally or chemically similar residues should ideally have high diagonal or off-diagonal numbers
- Structurally or chemically dissimilar residues should ideally have low diagonal or off-diagonal numbers
- What does the score mean: The likelihood of seeing two residues align (preserved) than random expected.



# Scoring Alignments

Blosum62 Scoring Matrix

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

# BLOSUM substitution matrices

Developed for distantly related proteins

Substitutions only from multiple alignments of conserved regions of protein families, hand curated, constitute the known homologous blocks

Identity threshold to define conserved blocks can be varied, e.g. 62% identity gives BLOSUM62

Scores calculated from frequency of amino acids in aligned pairs compared to what would be expected due to abundance alone, given all sequences

## Blosum Matrices

What score should we give to a ser residue aligned with a thr residue?

$$\text{score}(S : T) \propto \log_2 \frac{P(S : T \mid \text{homology})}{P(S : T \mid \text{random})}$$

# Example of deriving Blossum scores for S:S, S:T, and T:T

Database of known alignments

<b>S</b> DH I P	HK <b>S</b> A	WMFET <b>T</b>	R <b>T</b> Q C
<b>S</b> DH L P	HR <b>T</b> A	WMF D <b>T</b>	R <b>T</b> N C
<b>S</b> DH I P	HK <b>S</b> G	WLF D <b>T</b>	K <b>T</b> Q C
<b>S</b> EHL P			K <b>S</b> Q C
<b>S</b> EHL P			K <b>T</b> Q C

Homology Model (consider each pair of sequences separately)

S:S pairs in alignments = 11

S:T pairs in alignments = 6

T:T pairs in alignments = 9

$P(S:S|homology) = 11/117 = .094$

$P(S:T|homology) = 6/117 = .051$

$P(T:T|homology) = 9/117 = .078$

Total pairs in alignments = 117

# Example of deriving Blossum scores for S:S, S:T, and T:T

Database of known alignments

<b>S</b> DH I P	HK <b>S</b> A	WMFET <b>T</b>	R <b>T</b> QC
<b>S</b> DHLP	HR <b>T</b> A	WMFD <b>T</b>	R <b>T</b> NC
<b>S</b> DH I P	HK <b>S</b> G	WLFD <b>T</b>	K <b>T</b> QC
<b>S</b> EHLP			K <b>S</b> QC
<b>S</b> EHLP			K <b>T</b> QC

## Random Model

Number of S residues = 8       $P(S:S|random) = P(S)P(S) = (8/72)^2 = .012$   
Number of T residues = 8       $P(S:T|random) = 2 * P(S)P(T) = 2 * (8/72)^2 = .024$   
Total residues = 72               $P(T:T|random) = P(T)P(T) = (8/72)^2 = .012$

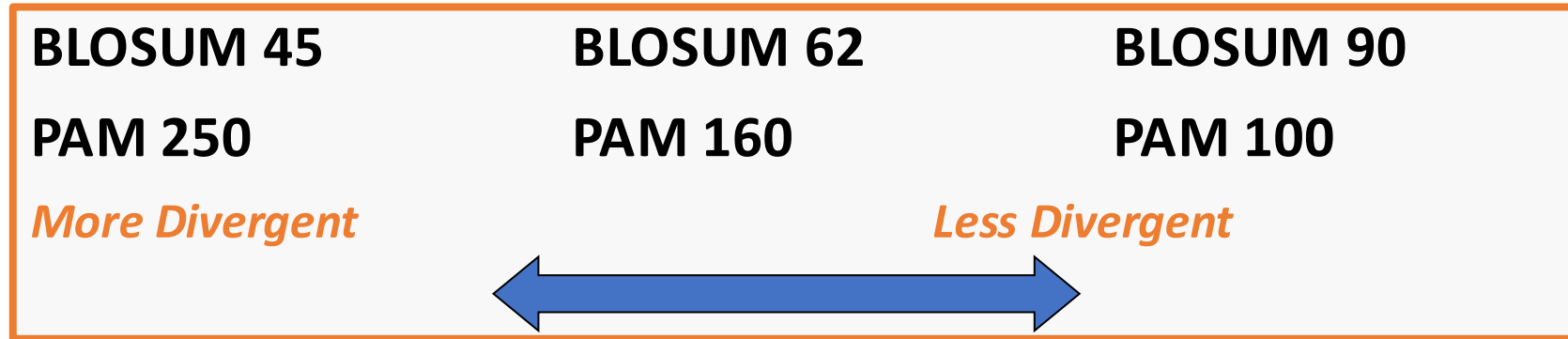
## Example of deriving Blossum scores for S:S, S:T, and T:T

$$\text{score}(S : S) = \log_2 \frac{P(S : S \mid \text{homology})}{P(S : S \mid \text{random})} = \log_2 \frac{.094}{.012} = 2.96$$

$$\text{score}(S : T) = \log_2 \frac{P(S : T \mid \text{homology})}{P(S : T \mid \text{random})} = \log_2 \frac{.051}{.024} = 1.09$$

$$\text{score}(T : T) = \log_2 \frac{P(T : T \mid \text{homology})}{P(T : T \mid \text{random})} = \log_2 \frac{.078}{.012} = 2.70$$

# BLOSUM and PAM



- BLOSUM 62 is the default matrix in BLAST 2.0. Though it is tailored for comparisons of moderately distant proteins, it performs well in detecting closer relationships. A search for distant relatives may be more sensitive with a different matrix.
- PAM matrices: point accepted mutation

# Scoring Matrices Take Home Points

- Based on log odds scores
  - Ratios  $>1$  give positive scores, ratios  $<1$  give negative scores
  - Because  $\log(x*y) = \log x + \log y$  the score of an alignment is the sum of the scores for each pair of aligned residues
- Assume independence of adjacent residues when scoring
- Introduced the concept that the frequency of a residue in a multiple alignment is informative



# Fast Similar Sequence Search

- Can we run Smith-Waterman between query and every DB sequence?
- Yes, but too slow!
- General approach
  - Break query and DB sequence to match subsequences
  - Extend the matched subsequences, filter hopeless sequences
  - Use dynamic programming to get optimal alignment

# BLAST

- Basic Local Alignment Search Tool
- Altschul et al. *J Mol Biol.* 1990
- One of the most widely used bioinformatics applications
  - Alignment quality not as good as Smith-Waterman
  - But much faster, supported at NCBI with big computer cluster

```
❑ >gi|2498170|sp|Q27974|AUXI BOVIN Auxilin
    Length = 910

Score = 107 bits (268), Expect = 4e-23
Identities = 76/275 (27%), Positives = 131/275 (47%), Gaps = 21/275 (7%)

Query: 22  DLDLTYYIPNIIAMGFPAERLEGVYRNNIDDVVRFLDSKHKHNYKIYNLCAERHYDTAKF 81
          DLD TY+   II M FP + ++ +RN +DD+  FLDS+H +HY +YNL + + Y TAKF
Sbjct: 60  DLDFTYVTSRIIVMSFPLDSVDIGFRNQVDDIRSFLDSRHLDDHYTVYNL-SPKSYRTAKF 118
```

## BLAST Algorithm Steps

- Query and DB sequences are optionally filtered to remove low-complexity regions
  - E.g. ACACACACA, TTTTTTTTTT

# BLAST Algorithm Steps

- Query and DB sequences are optionally filtered to remove low-complexity regions
- Break DB sequences into k-mer words and hash their locations to speed later searches
  - k is usually 11 for DNA/RNA and 3 for protein

LPPQGLL

LPP

PPQ

PQG

QGL

GLL

## BLAST Algorithm Steps

- Query and DB sequences are optionally filtered to remove low-complexity regions
- Break DB sequences into k-mer words and hash their locations to speed later searches
- Each k-mer in query find possible k-mers that matches well with it
  - “well” is evaluated by substitution matrices

# BLAST Algorithm Steps

- Only words with  $\geq T$  cutoff score is kept
  - T is usually 11-13, ~ 50 words make T cutoff
  - Note: this is 50 words at every query position
- For each DB sequence with a high scoring word, try to extend it in both ends

Query:	<b>LP PQG LL</b>
DB seq:	<b>MP PEG LL</b>
HSP score	$9 + 15 + 8 = 32$

- Form HSP (High-scoring Segment Pairs)
- Use BLOSUM to score the extended alignment
- No gaps allowed

# The BLAST Search Algorithm

## Query Word

Query: GSVEDTTGSQSLAALLNKCKT **PQG** QRLVNQWIKQPLMDKNRIEERLNLVEAFVEDAELRQTLQEDL

## Neighbourhood Words

PQG 18

PEG 15

PRG 14

PKG 14

PNG 13

PDG 13

PHG 13

**PMG** 13

PSG 13

PQA 12

PQN 12

...

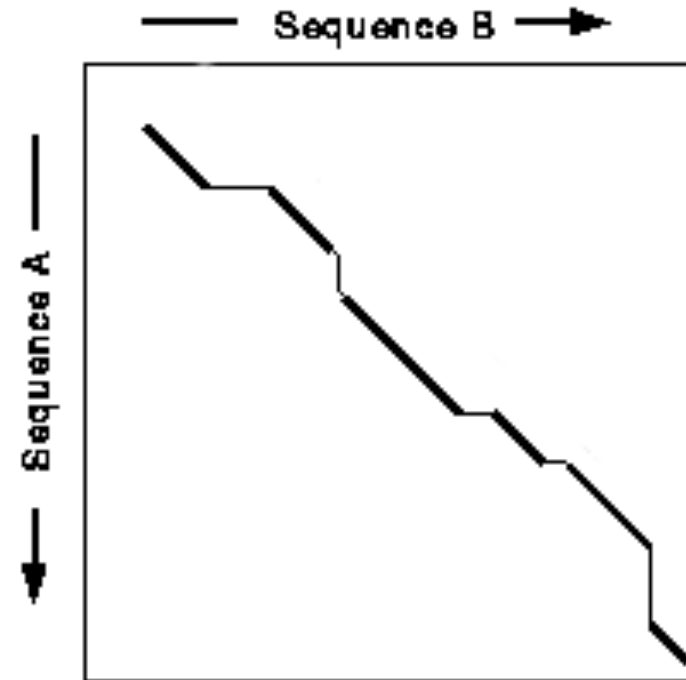
**Score Threshold (13)**

Query: 325 SLAALLNKCKT **PQG** QRLVNQWIKQPLMDKNRIEERLNLVEA 365  
+LA++L+ TP G R++ +W+ P+ D + ER + A  
Sbjct: 290 TLASVLDCTVT **PMG** SRMLKRWLHMPVRDTRVLLERQQTIGA 330

## High-scoring Segment Pair

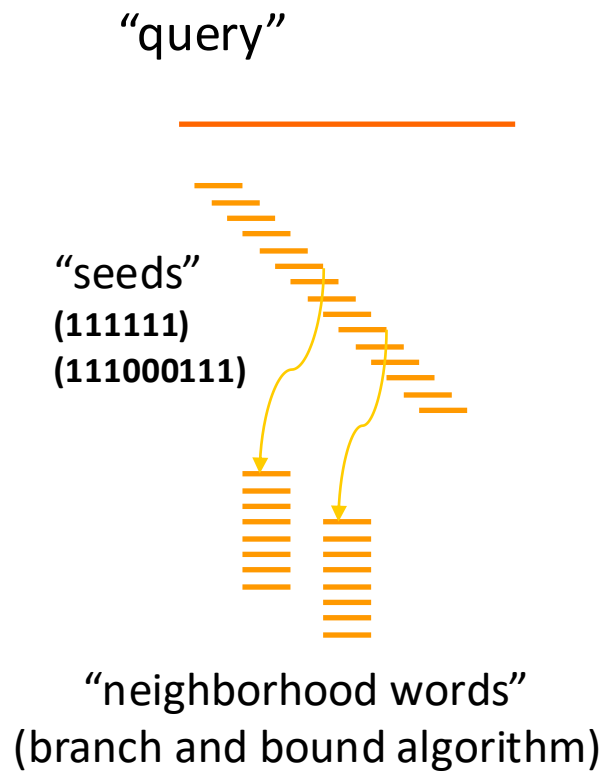
# BLAST Algorithm Steps

- Keep only statistically significant HSPs
  - Based on the scores of aligning 2 random seqs
- Use Smith-Waterman algorithm to join the HSPs and get optimal alignment
  - Gaps are allowed  
default (-11, -1)





# BLAST algorithm summary



Indexing all seeds

“subjects” (database)

Scan the index and find all word hits

DP extension to recover the high scoring pairs

Extending *high scoring pairs*

Evaluate Significance of HSPs by  
Karlin-Altschul Statistic:  $E = KMN \exp(-\lambda * S)$

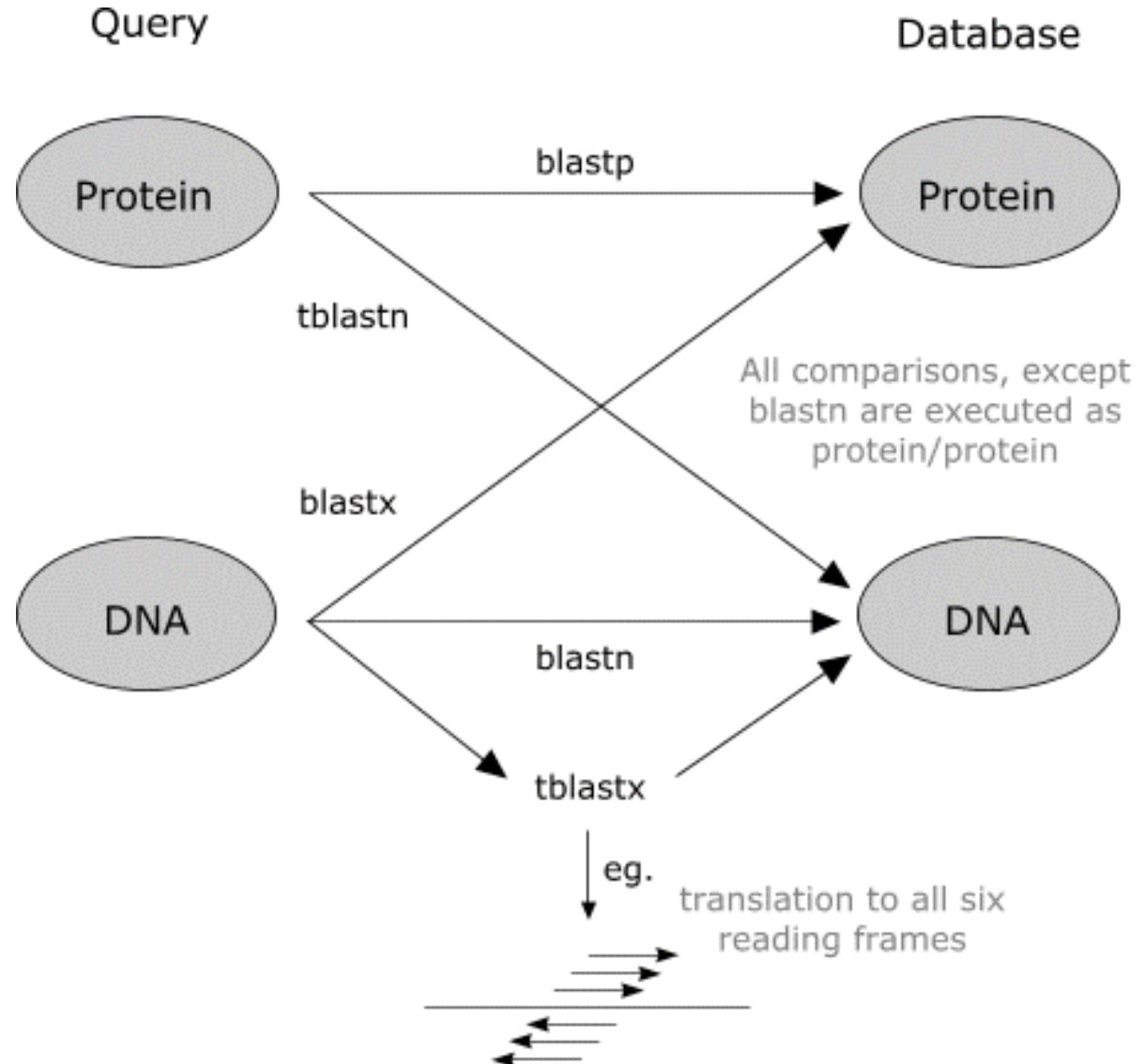
# Different BLAST Programs

BLAST DB:

- nr (non-redundant):
  - GenBank, RefSeq, EMBL...
- est:
  - expressed sequences (cDNA), redundant
- Swissprot and pdb:
  - protein databases

If query is DNA, but known to be coding (e.g. cDNA)

- Translate cDNA into protein
- Zero gap-extension penalty

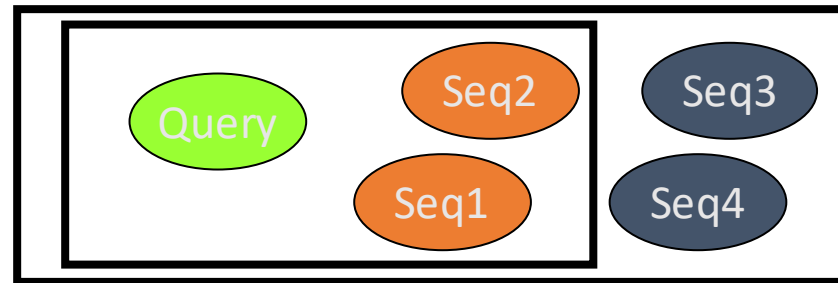


# Different BLAST Programs

Program	Description
blastp	Compares an amino acid query sequence against a protein sequence database.
blastn	Compares a nucleotide query sequence against a nucleotide sequence database.
blastx	Compares a nucleotide query sequence translated in all reading frames against a protein sequence database. You could use this option to find potential translation products of an unknown nucleotide sequence.
tblastn	Compares a protein query sequence against a nucleotide sequence database dynamically translated in all reading frames.
tblastx	Compares the six-frame translations of a nucleotide query sequence against the six-frame translations of a nucleotide sequence database. Please note that the tblastx program cannot be used with the nr database on the BLAST Web page because it is too computationally intensive.

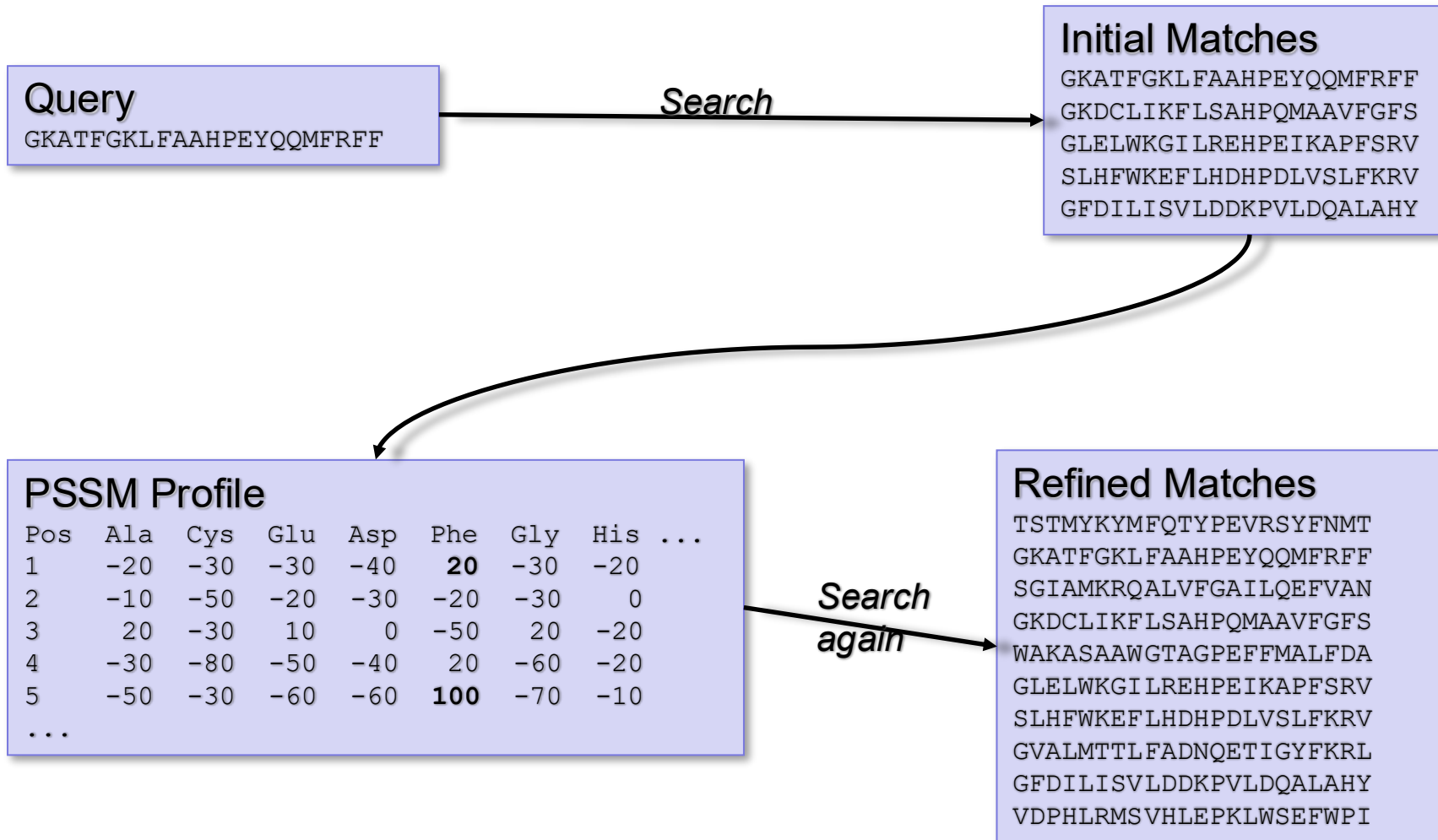
# PSI-BLAST

- Position Specific Iterative BLAST
  - Align high scoring hits in initial BLAST to construct a profile for the hits
  - Use profile (PSSM) for next iteration BLAST



- Find remote homologs or protein families
- FP sequences can degrade search quickly

# PSI-BLAST



<http://www.ncbi.nlm.nih.gov/blast/Blast.cgi>

<http://www.ebi.ac.uk/blastpgp/>

# Reciprocal Blast

- Search for orthologous sequences between two species
  - GeneA in Species1 BLAST Species2 → GeneB
  - GeneB in Species2 BLAST Species1 → GeneA
- GeneA      ←————→ GeneB  
                  orthologous
- Also called bi-directional best hit

# BLAT

- BLAST-Like Alignment Tool
  - Compare to BLAST, BLAT can align much longer regions (MB) really fast with little resources
  - E.g. can map a sequence to the genome in seconds on one Linux computer
  - Allow big gaps (mRNA to genome)
  - Need higher similarity (> 95% for DNA and 80% for proteins) for aligned sequences
- Basic approach
  - Break long sequence into blocks
  - Index k-mers, typically 8-13
  - Stitch blocks together for final alignment

# BLAT: Indexing

**Genome:** cacaattatcacgaccgc

**3-mers:** cac aat tat cac gac cgc

**Index:**

aat	3	gac	12
cac	0, 9	tat	6
cgc	15		

**cDNA (mRNA -> DNA):** aattctcac

**3-mers:**

aat	att	ttc	tct	ctc	tca	cac
0	1	2	3	4	5	6

**hits:**

aat	0, 3	-3
cac	6, 0	6
cac	6, 9	-3

**clump:** cac**AAT**tat**CAC**gaccgc  
          | | |          | | |  
          aattctcac



# Summary of Fast Search

- Fast sequence similarity search
  - Break seq, hash DB sub-seq, match sub-seq and extend, use DP for optimal alignment
  - \*BLAST, most widely used, many applications with sound statistical foundations
  - \*BLAT, align sequence to genome, fast yet need higher similarity

## BLAST score and significance

- Report DB sequences above a threshold
  - E value: Number (instead of probability  $\rightarrow$  pvalue) of matches expected merely by chance

$$E = Kmn e^{-\lambda S}$$

$$p(s \geq x) \approx 1 - \exp[-e^{-x}]$$

- m, n are query and DB length
- K,  $\lambda$  are constants
- Smaller E, more stringent

# Are these proteins homologs?

**SEQ 1:** RVVNLVPS--FWVL DATYKNYA INYNCDV TYKLY

L P      W L      Y N      Y C      L

Probably not (score = 9)

**SEQ 2:** QFFPLMPPAPYWILATDYENLPLVYSCTTFFWLF

**SEQ 1:** RVVNLVPS--FWVL DATYKNYA INYNCDV TYKLY

L P      W L DATYKNYA    Y C      L

MAYBE (score = 15)

**SEQ 2:** QFFPLMPPAPYWIL DATYKNYA LVYSCTTFFWLF

**SEQ 1:** RVVNLVPS--FWVL DATYKNYA INYNCDV TYKLY

RVV L PS      W L DATYKNYA    Y CDV TYKL

Most likely (score = 24)

**SEQ 2:** RVVPLMPSAPYWIL DATYKNYA LVYSCDV TYKLF

## Significance of scores

HPDKKAHSIHAWILSKSKVLEGNTKEVVDNVLKT

Homology  
detection  
algorithm

45

LENENQGKCTIAEYKYDGKKASVYNSFVSNGVKE

Low score = unrelated  
High score = homologs

*How high is high enough?*

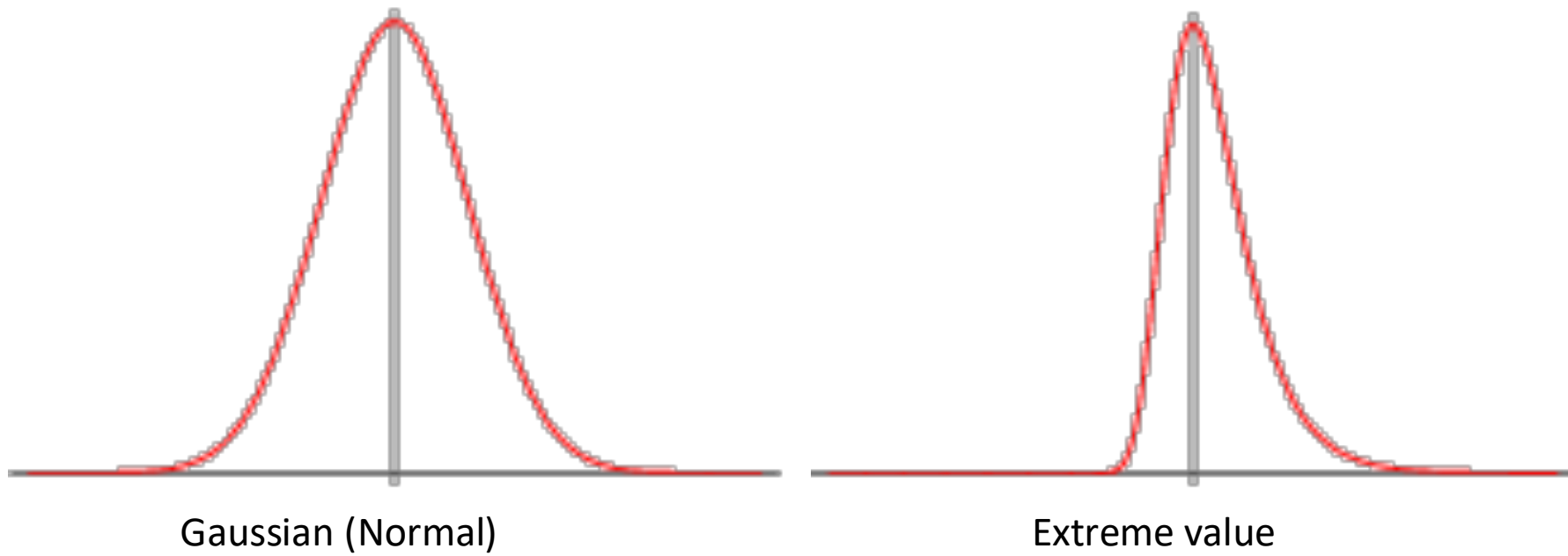
## Other significance questions

- Pairwise sequence comparison scores
- Microarray expression measurements
- Sequence motif scores
- Functional assignments of genes
- Call peaks from ChIP-seq data

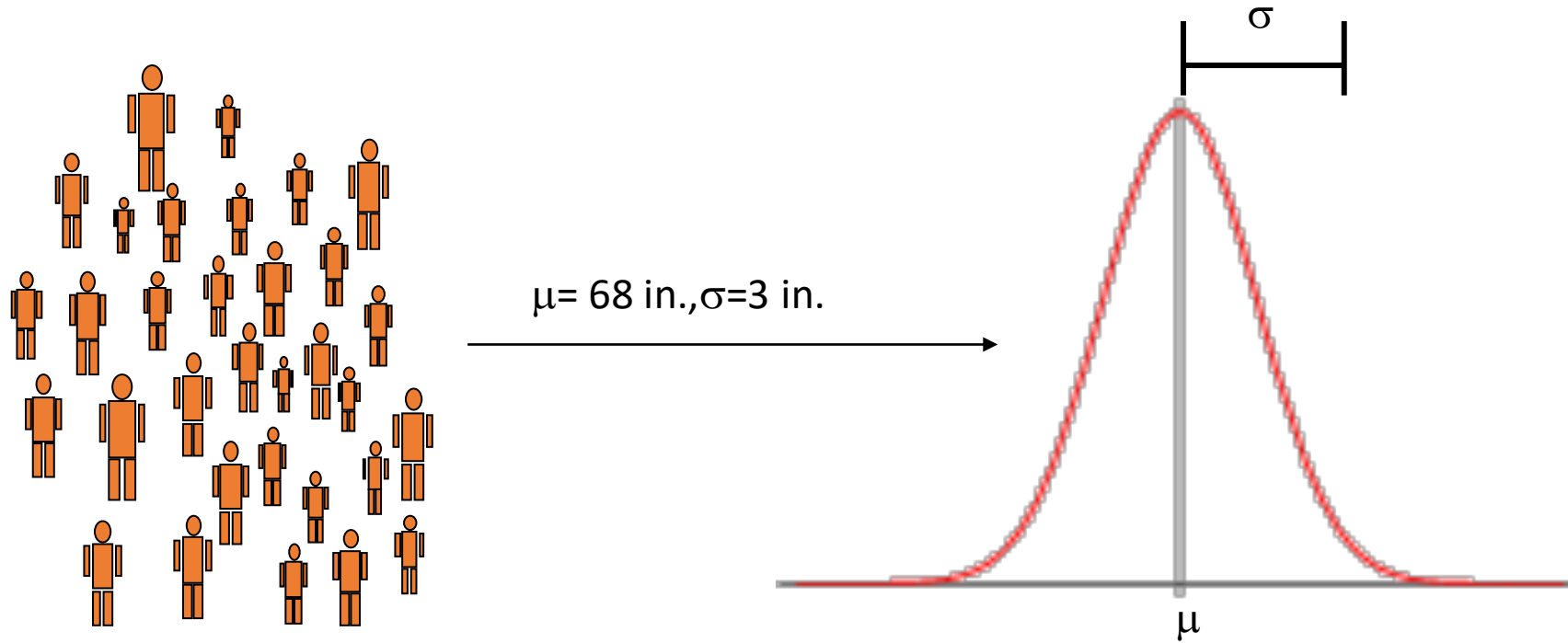
# The null hypothesis

- We are interested in characterizing the distribution of scores from sequence comparison algorithms.
- We would like to measure how surprising a given score is, *assuming that the two sequences are not related*.
- The assumption is called the **null hypothesis**.
- The purpose of most statistical tests is to determine whether the observed results provide a reason to reject the hypothesis that they are merely a product of chance factors.

# Gaussian vs. Extreme Value Distribution (EVD)



# Gaussian



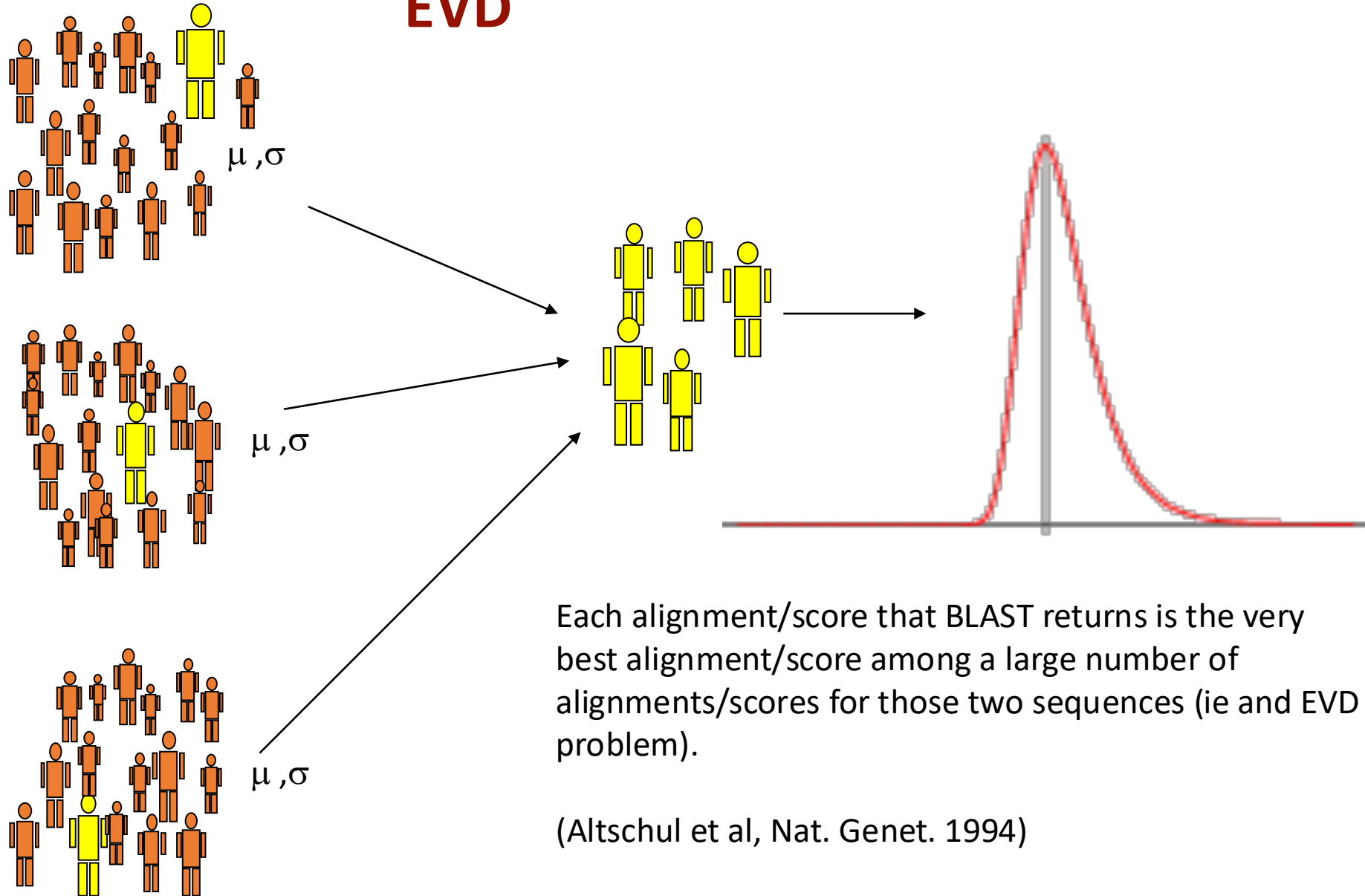
What is the chance of picking a person at least 75 in. tall  $P(X \geq 75)$ ?

$$z_{\text{score}}(x) = \frac{x - \mu}{\sigma} = \frac{75 - 68}{3} = 2.33$$

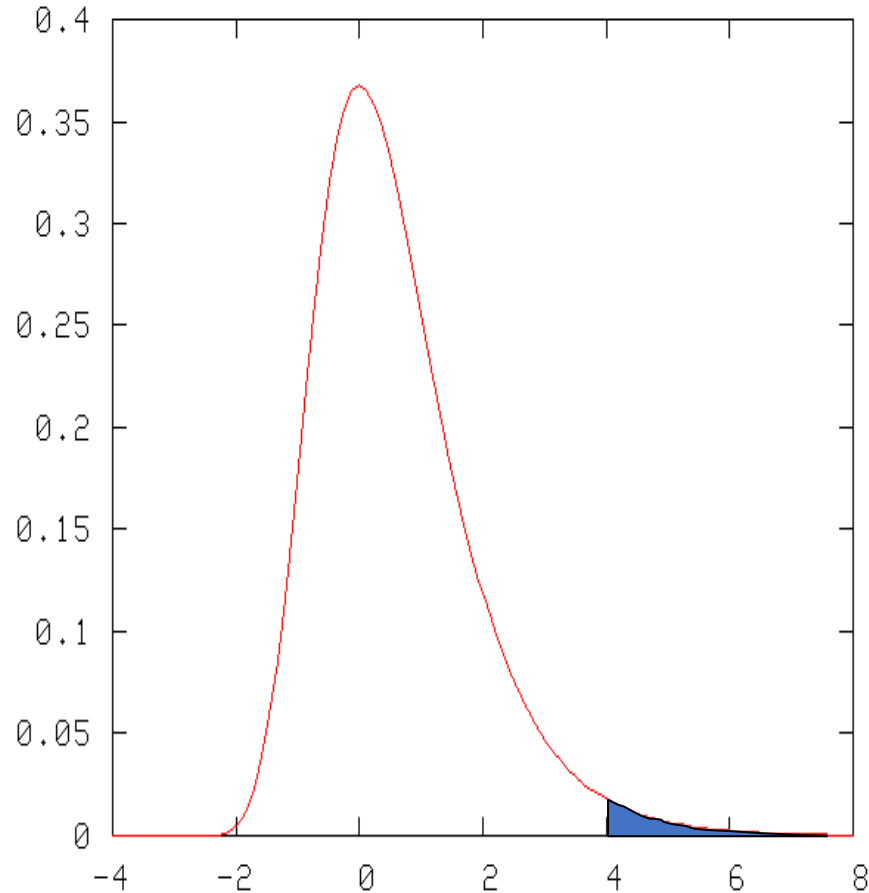
From Table:  
 $z=2.33 \rightarrow P=0.01$



# EVD

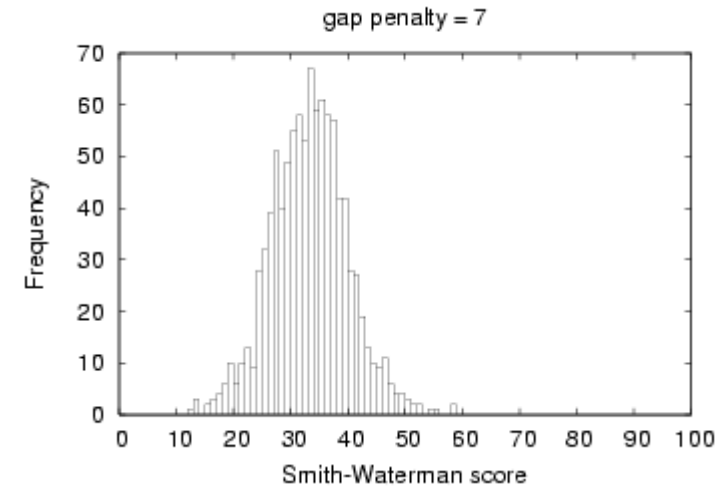
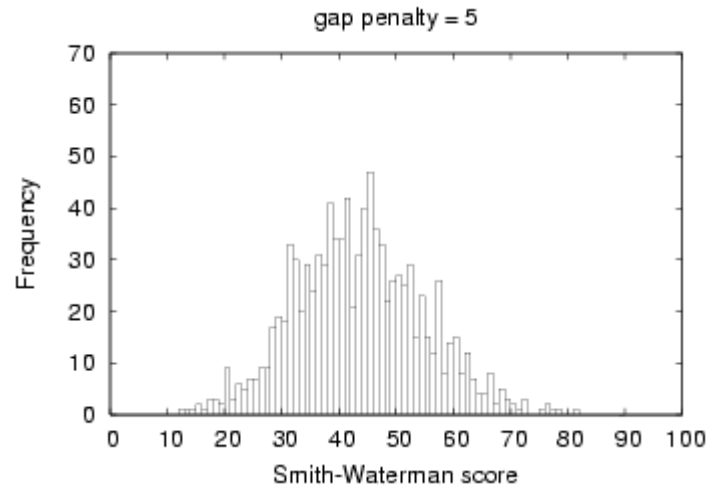


# Computing a p-value



- The probability of observing a score  $>4$  is the area under the curve to the right of 4.
- This probability is called a p-value.
- $\text{p-value} = \Pr(\text{data} | \text{null})$

# Scaling the EVD



- An extreme value distribution derived from, e.g., the Smith-Waterman algorithm will have a characteristic mode  $\mu$  and scale parameter  $\lambda$ .

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

- These parameters depend upon the size of the query, the size of the target database, the substitution matrix and the gap penalties.

## An example

You run BLAST and get a score of 45. You then run BLAST on a shuffled version of the database, and fit an extreme value distribution to the resulting empirical distribution. The parameters of the EVD are  $\mu = 25$  and  $\lambda = 0.693$ . What is the p-value associated with 45?

$$\begin{aligned}P(S \geq x) &= 1 - \exp\left[-e^{-\lambda(x-\mu)}\right] \\P(S \geq 45) &= 1 - \exp\left[-e^{-0.693(45-25)}\right] \\&= 1 - \exp\left[-e^{-13.86}\right] \\&= 1 - \exp\left[-9.565 \times 10^{-7}\right] \\&= 1 - 0.999999043 \\&= 9.565 \times 10^{-7}\end{aligned}$$

# Summary of statistical significance

- A distribution plots the frequency of a given type of observation.
- The area under the distribution is 1.
- Most statistical tests compare observed data to the expected result according to the null hypothesis.
- Sequence similarity scores follow an extreme value distribution, which is characterized by a larger tail.
- The p-value associated with a score is the area under the curve to the right of that score.

```
.....ACGTTGCCACTTTCCGGGCCACCTGGCCACCTTATTTTCGGAAATATACCGGGCCTTTTTT.....  
      |||||x||||x|||||||  
      CTTTCCCGGCCTCCTGGCCA
```

match: +1

mismatch: -1

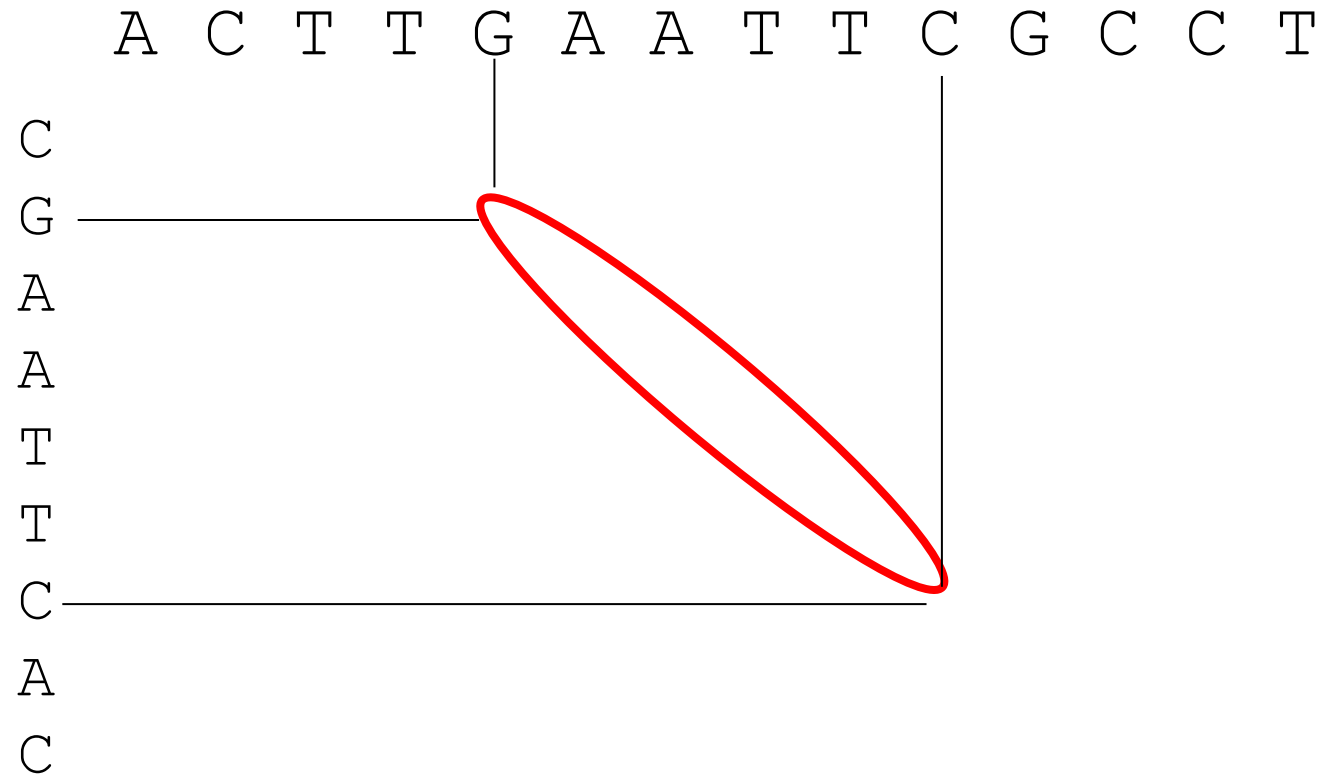
matching score = 16

- **How to align them?**
- **Why we can align them?**
- **Why +1 for match, and -1 for mismatch?**
- **What does the score mean?**
- **Is 16 a good score?**

# Applying homology: concept and technology

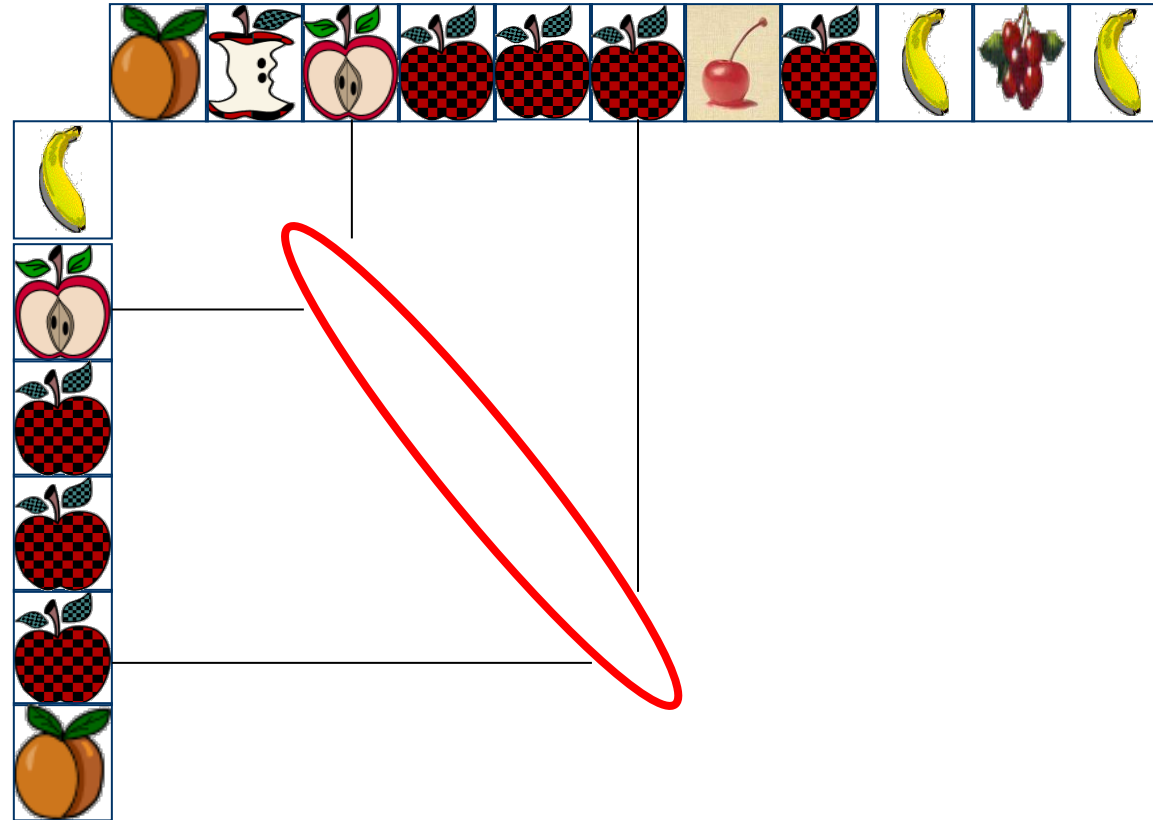
- Genome evolution
  - Mammalian genome evolution
  - Human genome variation
  - Cancer genome evolution
- Gene finding
  - Comparative approaches
  - Ab initio approaches
    - Hidden Markov Model
- Protein structure
  - Threading
- Regulatory motif finding
  - Profile comparison
- Pathway/Network comparison
  - PathBLAST
- Conservation
  - Ultra conserved elements
  - Human accelerated regions

## So far, only linear sequence comparison





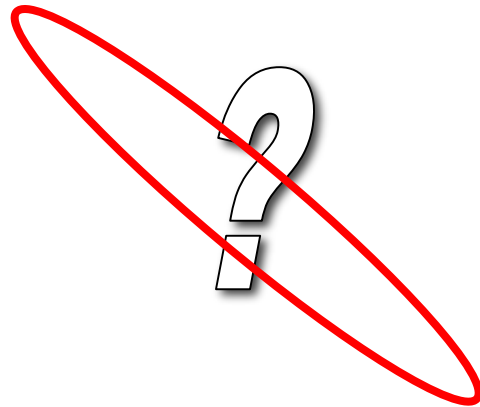
# Expanding the idea of a sequence



# Central theme of the new algorithm – compare profiles

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

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

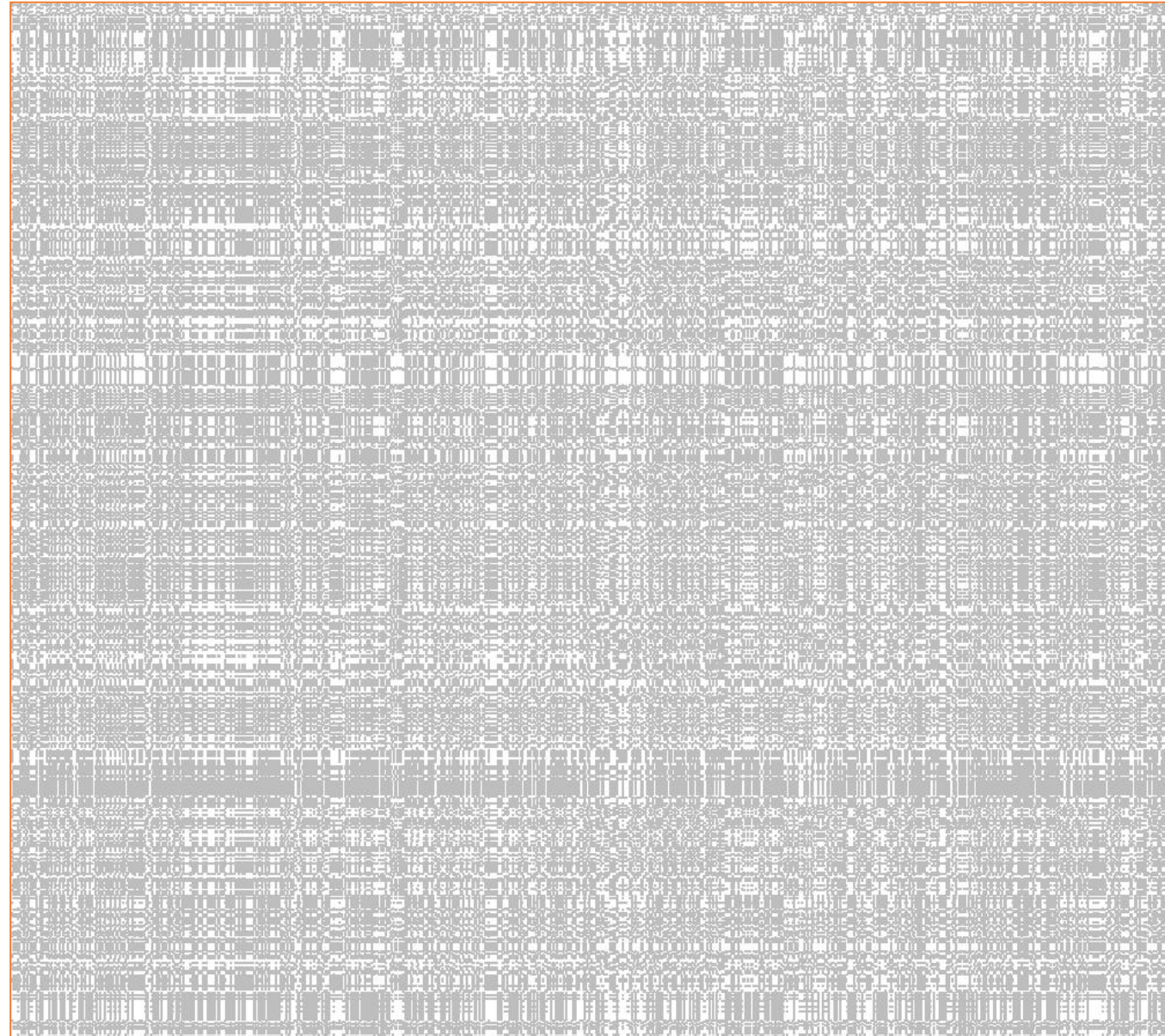


Met14 vs  
Met2  
“DotPlot”

MET14 (1000nt)

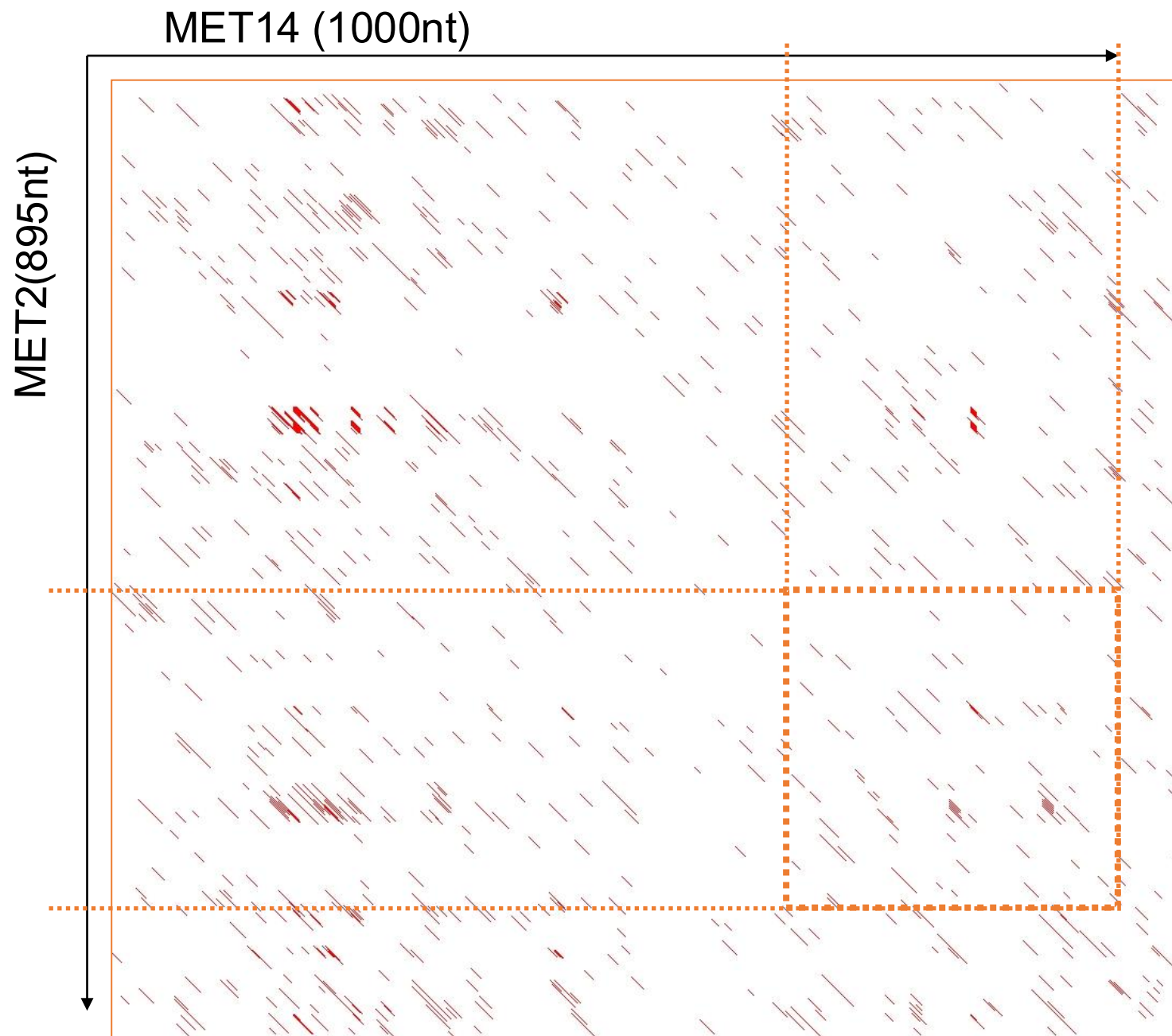
MET2(895nt)

Match = 1  
Mismatch = -1  
Gray: 1



# Met14 vs Met2

Red: >5

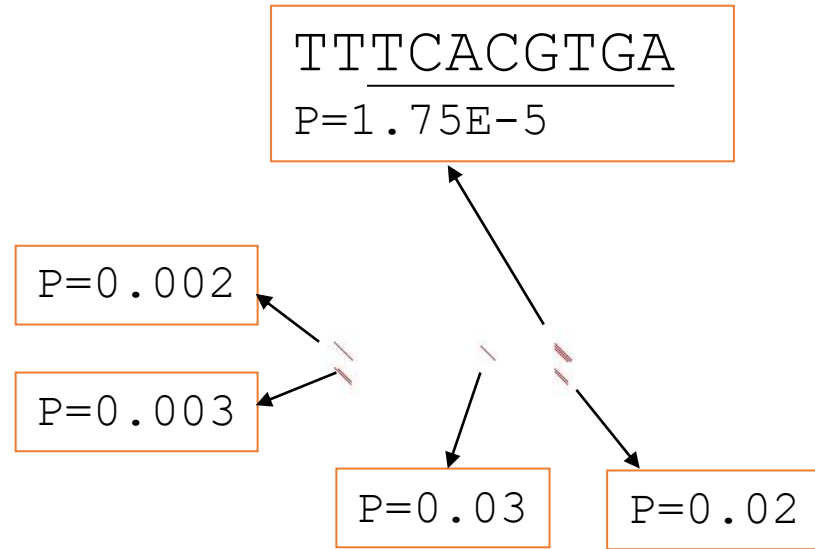


# Met14 vs Met2 PhyloNet

MET14 (1000nt)

MET2(895nt)

HSPs:  
 $E < 0.1$

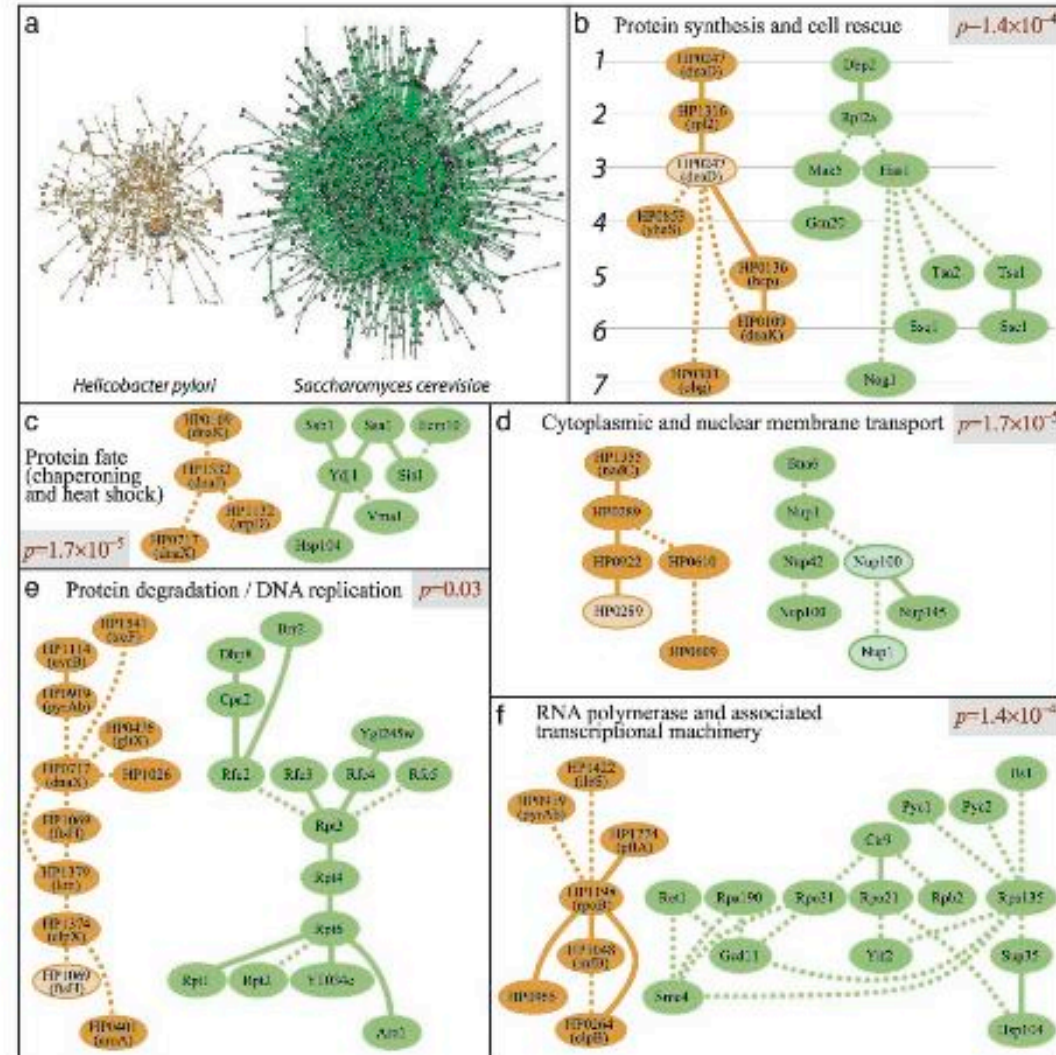
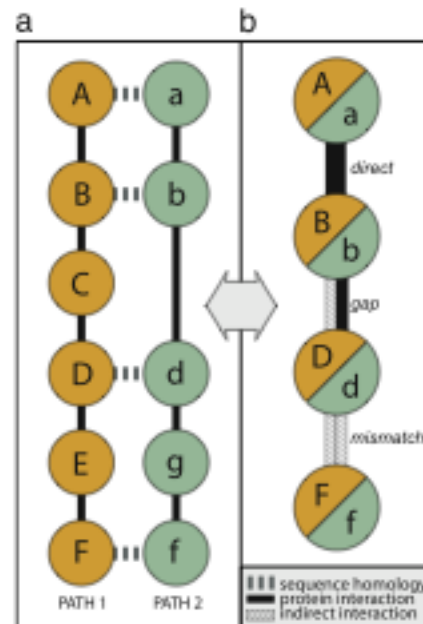




# PathBlast, NetworkBlast



Trey Ideker



# Aligning Short Reads

# NGS: Sequence alignment

- Map the **large** numbers of **short** reads to a reference genome
  - In a broader sense: Identify similar sequences (DNA, RNA, or protein) in consequence of functional, structural, or evolutionary relationships between the them
  - Applications: Genome assembly, SNP detection, homology search, etc
- **large**  $\Rightarrow$  faster search speed
- **short**  $\Rightarrow$  greater search sensitivity.



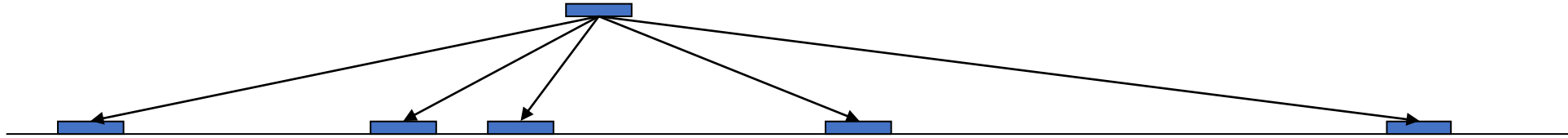
# Many short read aligners

- Bfast
- BioScope
- Bowtie
- BWA
- CLC bio
- CloudBurst
- Eland/Eland2
- GenomeMapper
- GnuMap
- Karma
- MAQ
- MOM
- Mosaik
- MrFAST/MrsFAST
- NovoAlign
- PASS
- PerM
- RazerS
- RMAP
- SSAHA2
- Segemehl
- SeqMap
- SHRiMP
- Slider/SliderII
- SOAP/SOAP2
- Srprism
- Stampy
- vmatch
- ZOOM
- .....

# Short read mapping

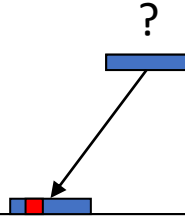
- Input:
  - A reference genome
  - A collection of many 25-100bp tags (reads)
  - User-specified parameters
- Output:
  - One or more genomic coordinates for each tag
- In practice, only 70-75% of tags successfully map to the reference genome. Why?

# Multiple mapping



- A single tag may occur more than once in the reference genome.
- The user may choose to ignore tags that appear more than  $n$  times.
- As  $n$  gets large, you get more data, but also more noise in the data.

# Inexact matching



- An observed tag may not exactly match any position in the reference genome.
- Sometimes, the tag *almost* matches one or more positions.
- Such mismatches may represent a SNP or a bad read-out.
- The user can specify the maximum number of mismatches, or a phred-style quality score threshold.
- As the number of allowed mismatches goes up, the number of mapped tags increases, but so does the number of incorrectly mapped tags.

# Mapping Reads Back

- Hash Table (Lookup table)
  - FAST, but requires perfect matches
- Array Scanning
  - Can handle mismatches, but not gaps
- Dynamic Programming (Smith Waterman, Forward, Viterbi)
  - Indels
  - Mathematically optimal solution
  - Slow (most programs use Hash Mapping as a prefilter)
- Burrows-Wheeler Transform (BW Transform)
  - FAST (memory efficient)
  - But for gaps/mismatches, it lacks sensitivity