

## Lecture 2

# Sequence Searching: Pattern Matching and Blast

- Sequence format and notation
- Sequence comparison models
- Sequence searching phases
- BLAST: 3 step heuristic
- Hit statistics: P-value and E-value
- BLAST result interpretation
- BLAST flavours

# Sequence Searching

## Sequence

- a string of characters that represents the chain of building blocks in a biopolymer
- building blocks are amino acids (proteins) and nucleotides (DNA, RNA)
- the side chains define the sequence

## Sequence Searching

Sequence databank searching is the process of extracting homologues of one or several query sequence(s) from a sequence database.

# Sequence Format Conventions

- A sequence is composed of a name (often including an accession number) and the residue string.
- A sequence databank is a formatted (and often sorted) list of sequences (here FASTA format).

>Iepi epidermal growth factor (Mus musculus)

NSYPGCPSSYDGYCLNGGVCMHIESLDSYTCNCVIGYSGDRCQTRDLRWVWELR

>Iixa EGF-like module coagulation factor (Homo sapiens)

VDGDQCESNPCLNGGSCKDDINSYECWCPFGFEGKNCEL

- Proteins are written from the N-terminus to the C-terminus.
- Nucleotide sequences are defined within a 'reading frame'.

# Sequence Notation - Positions and Chains

## Mutations

- Y35G-BPTI (bovine pancreatic trypsin inhibitor)  
mutation from Tyr to Gly at position 35
- K(B29)P-insulin  
mutation from Lys to Pro at position 29 in chain B
- Des(B27-B30)-insulin-B26-carboxamide  
residues 27 to 30 deleted in chain B and C-terminus  
amidated

## Chain notation

Chains are denoted A, B, C, D,... in successive order.

# Sequence Searching and Sequence Alignment

- Sequence searching and sequence alignment are different techniques.
- Sequence searching uses sequence matching (like alignment), but sequence alignment works on a pre-defined sequence set.
- In sequence searching one tries to extract homologous sequences, whereas in sequence alignment one assumes homology and tries to correctly identify similarities between homologues.

# Sequence Searching Problem

given

required

# Sequence Searching Problem

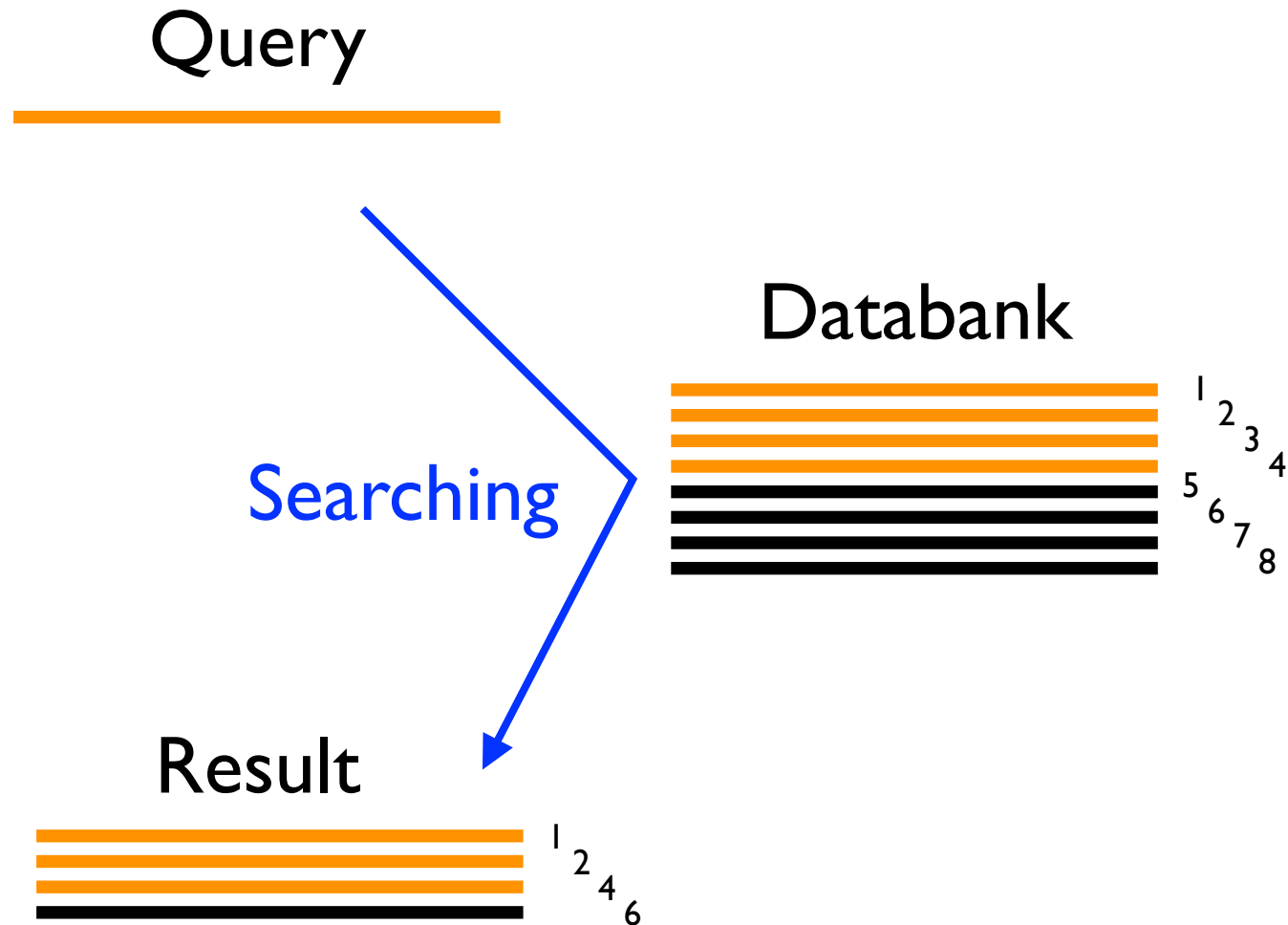
given

- Query: sequence to search with (assume 300 aa)
- Database (search space): very many sequences (assume  $> 100000$ )
- Goal: find sequences homologous to query

required

- a fast tool
- primarily a filter: most sequences will be unrelated to the query
- fine-tune the alignment later

# Sequence Databank Searching



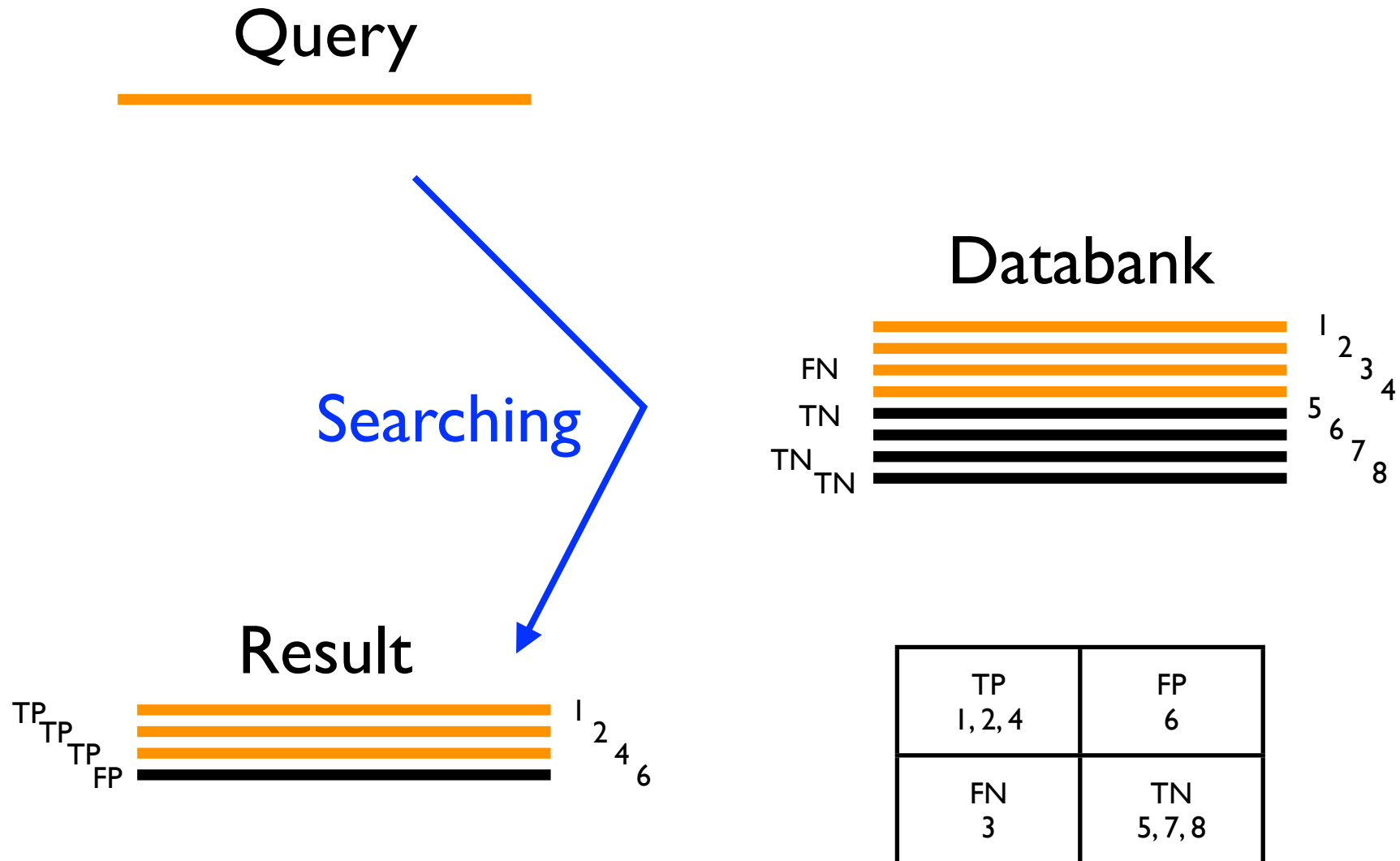


# Prediction : Contingency Table

Sequence searching is a prediction about the homology relation between query and DB sequence.

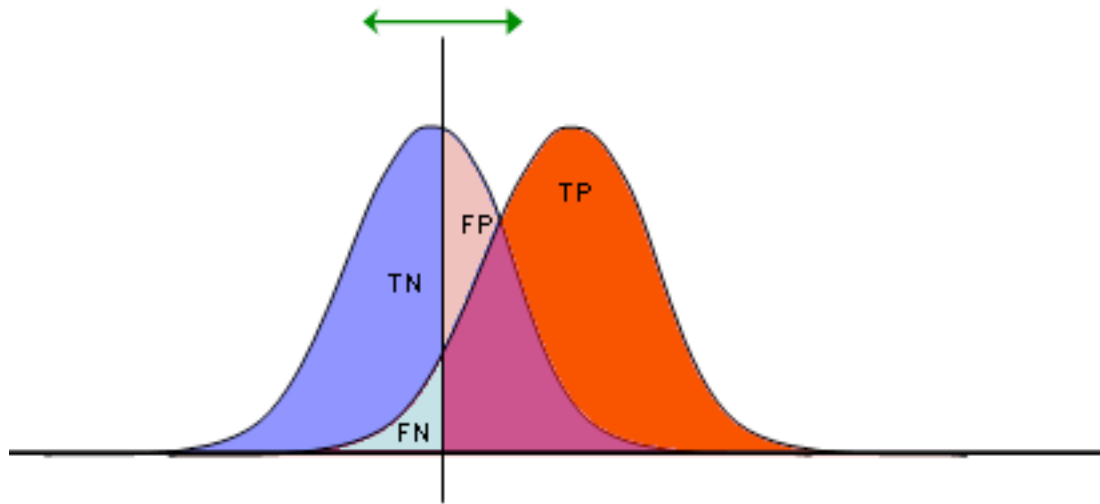
		reality	
		positive	negative
prediction	positive	true positive (TP)	false positive (FP)
	negative	false negative (FN)	true negative (TN)

# Sequence Databank Searching



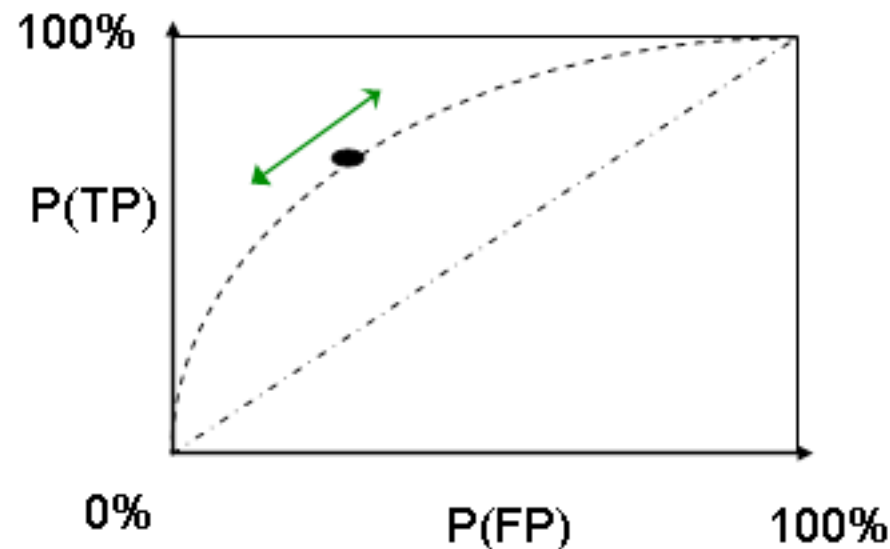
# Inference / Prediction Statistics

## Random and Target Distributions



## Contingency Table

TP	FP
FN	TN
1	1



ROC or  
benchmark  
curve

# Models Underlying Sequence Comparison

Sequence comparison (searching and alignment) is predominantly based on three models:

## 1. Scoring Model PAM or BLOSUM matrix + gap penalties

Information content of matched amino acid pair is defined as

$$\text{Score}(XZ) = \log [P(XZ) / P(X)P(Z)]$$

with  $P(XZ)$ : observed probability of matched XY in trusted alignments  
and  $P(X), P(Z)$ : probabilities of X and Z in a random sequence

## 2. Alignment Model

local or global alignment of homologues

=> alignment scores of biologically meaningful alignments

## 3. Random Model

local or global alignment of un-related (=random) sequence pair

=> alignment scores of random alignments

# Scoring Model: Substitution Matrix

<b>A</b>	2																					
<b>R</b>	-2	6																				
<b>N</b>	0	0	2																			
<b>D</b>	0	-1	2	4																		
<b>C</b>	-2	-4	-4	-5	12																	
<b>Q</b>	0	1	1	2	-5	4																
<b>E</b>	0	-1	1	3	-5	2	4															
<b>G</b>	1	-3	0	1	-3	-1	0	5														
<b>H</b>	-1	2	2	1	-3	3	1	-2	6													
<b>I</b>	-1	-2	-2	-2	-2	-2	-2	-3	-2	5												
<b>L</b>	-2	-3	-3	-4	-6	-2	-3	-4	-2	2	6											
<b>K</b>	-1	3	1	0	-5	1	0	-2	0	-2	-3	5										
<b>M</b>	-1	0	-2	-3	-5	-1	-2	-3	-2	2	4	0	6									
<b>F</b>	-4	-4	-4	-6	-4	-5	-5	-5	-2	1	2	-5	0	9								
<b>P</b>	1	0	-1	-1	-3	0	-1	-1	0	-2	-3	-1	-2	-5	6							
<b>S</b>	1	0	1	0	0	-1	0	1	-1	-1	-3	0	-2	-3	1	2						
<b>T</b>	1	-1	0	0	-2	-1	0	0	-1	0	-2	0	-1	-3	0	1	3					
<b>W</b>	-6	2	-4	-7	-8	-5	-7	-7	-3	-5	-2	-3	-4	0	-6	-2	-5	17				
<b>Y</b>	-3	-4	-2	-4	0	-4	-4	-5	0	-1	-1	-4	-2	7	-5	-3	-3	0	10			
<b>V</b>	0	-2	-2	-2	-2	-2	-2	-1	-2	4	2	-2	2	-1	-1	-1	0	-6	-2	4		
<b>B</b>	0	-1	2	3	-4	1	2	0	1	-2	-3	1	-2	-5	-1	0	0	-5	-3	-2	2	
<b>Z</b>	0	0	1	3	-5	3	3	-1	2	-2	-3	0	-2	-5	0	0	-1	-6	-4	-2	2	3
<b>A</b>		<b>R</b>	<b>N</b>	<b>D</b>	<b>C</b>	<b>Q</b>	<b>E</b>	<b>G</b>	<b>H</b>	<b>I</b>	<b>L</b>	<b>K</b>	<b>M</b>	<b>F</b>	<b>P</b>	<b>S</b>	<b>T</b>	<b>W</b>	<b>Y</b>	<b>V</b>	<b>B</b>	<b>Z</b>

# Scoring Model: PAM Matrices

- Derive permissible mutations from trusted alignments.
- Step 1: Construct phylogenetic tree from protein alignments.
- Step 2: Fill substitution matrix with observed transition probabilities.
- A 'PAM' is an evolutionary distance:
  - 1 PAM = 1 accepted mutation per 100 amino acids
  - 250 PAM = 2.5 accepted mutations per amino acid
- There are modern versions of the PAM matrix: Gonnet matrices or Jones-Taylor matrices.

# Scoring Model:

## BLOSUM Matrices

- Mostly used family of amino acid substitution matrices is the BLOSUM series (BLOSUM50, BLOSUM62).
- The BLOSUM matrices are derived from the BLOCKS database of multiple alignments (Henikoff & Henikoff, 1992).
- BLOSUM50 is derived from BLOCKS (core) alignment regions with  $\geq 50\%$  sequence identity, Blosum62 from those  $\geq 62\%$ , etc.
- Note that there is NO underlying Markov model of amino acid substitution; BLOSUM is based on observed pairs of aligned residues. Higher numbers mean shorter evolutionary distance (opposite to PAM!).

# Alignment Model

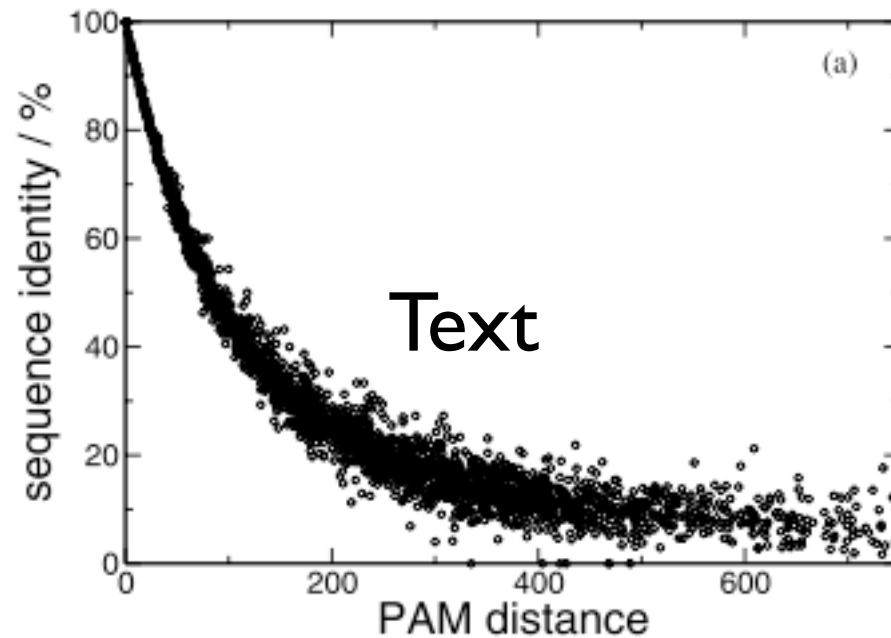
- The alignment model is the result of a meaningful sequence alignment, i.e. a correct alignment of two homologous sequences. In the ideal case this is identical to the alignment derived from structure superpositioning.
- Such meaningful alignments are called 'reference alignments' or 'trusted alignments'.
- Alignment models include local, global, probabilistic or other forms matching residue pairs.
- Generally we want to maximise  $P(XZ)$ .



# Random Model

- The random model is the result of an alignment of random sequences.
- The probability of finding a matched pair is just a chance event, which is given by the product of the probability to find the individual amino acids in the sequences (joint probability of independent events):  $P(X) * P(Z)$ .
- The random model is important, because amino acids have different frequencies of occurrence.
- Is that the only random model you can think of?

# Evolutionary Distance and Sequence Identity



Distant sequences are difficult to match, because the amount of identical or similar residues is low (signal / noise ratio is low).

# Sequence Searching Phases

## Matching Phase

The query sequence is matched (partially) with the databank sequence.

## Scoring Phase

The matched residues pairs are scored the scores are summed up.

## Selection phase

Based on a statistical criterion, database sequences with score above a user-defined cutoff score are returned as hits. Each sequence in the hitlist is a potential homologue.

# **BLAST:**

## **Basic Local Alignment Search Tool**

BLAST is a program designed to compare a query sequence with every sequence in a database and to report the most similar sequences.

Basic idea:

- High scoring segments have well conserved (almost identical) parts.
- After well conserved parts have been identified, extend them to the real alignment.

# BLAST: History

Smith and Waterman, 1981

Exact Local Dynamic Programming

Lipman and Pearson, 1985

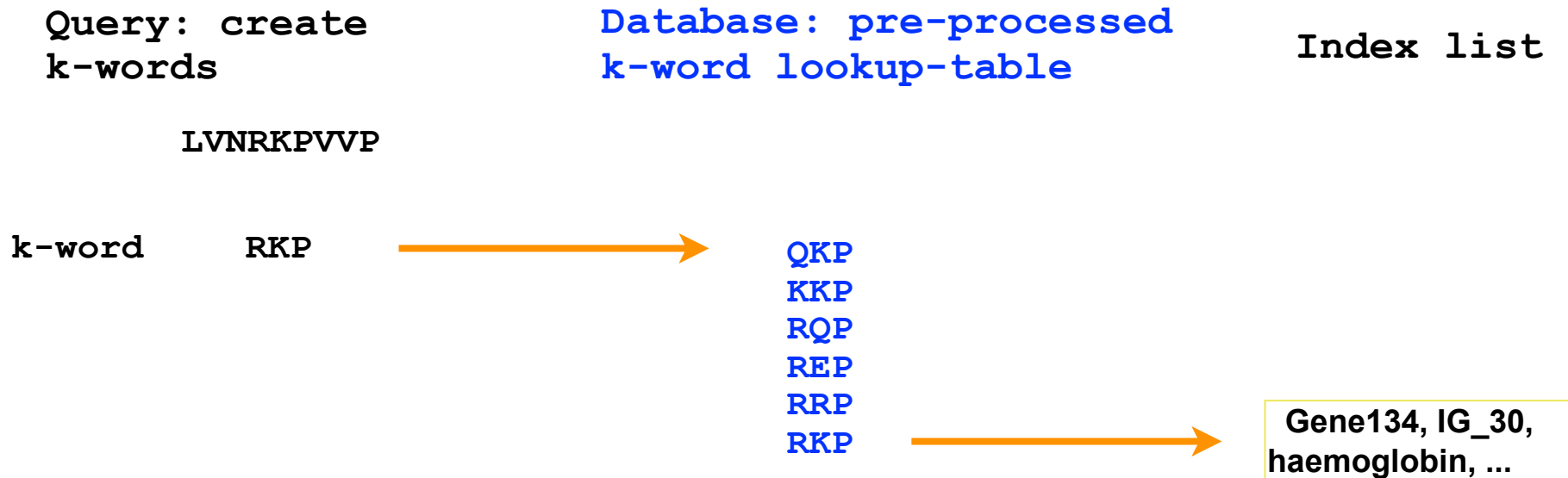
FASTA finds similar words (k-tup) on the same diagonal.

Altschul et al., 1990

BLAST: The most widely cited tool in Biology

# BLAST: k-word search

Considers only sequences with at least two k-word matches above a certain score.



# BLAST: match extension

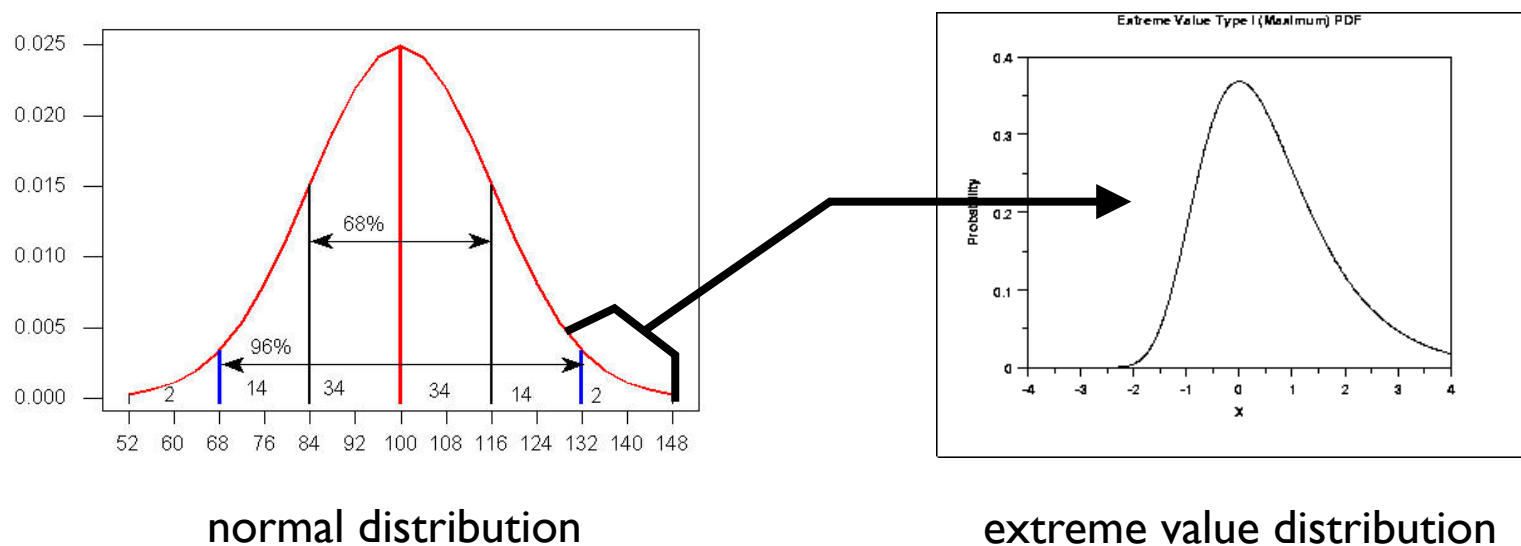
Elongate k-word matches in both directions as far as possible.

Query:	L	V	N	<b>R</b>	<b>K</b>	<b>P</b>	V	V	P
Database:	G	V	C	<b>R</b>	<b>K</b>	<b>P</b>	L	K	C
Score:	-3	4	-3	5	2	7	1	-2	-3

Extension terminates if score drops below value X (ungapped local alignments).

# BLAST: score evaluation

The BLAST search yields a match (alignment) score for each query/hit sequence pair.



Sampling from the extreme end of a normal distribution yields an extreme value distribution. BLAST scores are extreme value distributed.



# BLAST: Score Evaluation

P-value:

Probability of finding at least a match score  $x$ .

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

**K** is calibrated with the database composition.

**Lambda** is calibrated with the substitution matrix.

**mn** is the search space, i.e. length 'm' of query sequence and length 'n' of the database sequence.

# BLAST: Score Threshold

The likelihood of random sequences to yield a score greater than  $T$  increases linearly with the logarithm of the search space  $mn$ .

This gives a formula for accepting hits:

$$S > T + \log(mn) - \lambda$$

# BLAST: E-value

E-value: P-value adjusted by database size.

Number of database hits expected by chance

E-values are easier to compare than P-values, because they represent a sequence count rather than a probability.



# P31383 Vs YEAST

<input type="checkbox"/>	sp	<a href="#">P31383</a>	2AAA_YEAST Protein phosphatase PP2A regulatory subunit...	<a href="#">1170</a>	0.0
<input type="checkbox"/>	sp	<a href="#">P33892</a>	GCN1_YEAST Translational activator GCN1 [GCN1] [Saccha...	<a href="#">47</a>	1e-05
<input type="checkbox"/>	sp	<a href="#">P53978</a>	EF3B_YEAST Elongation factor 3B (EF-3B) [YEF3B] [Sacch...	<a href="#">46</a>	2e-05
<input type="checkbox"/>	sp	<a href="#">P22219</a>	VP15_YEAST Protein kinase VPS15 (EC 2.7.1.-) [VPS15] [...	<a href="#">40</a>	0.001
<input type="checkbox"/>	sp	<a href="#">P32337</a>	IMB3_YEAST Importin beta-3 subunit (Karyopherin beta-3...	<a href="#">39</a>	0.002
<input type="checkbox"/>	sp	<a href="#">P49955</a>	S3B1_YEAST U2 snRNP component HSH155 [HSH155] [Sacchar...	<a href="#">38</a>	0.005
<input type="checkbox"/>	sp	<a href="#">Q06708</a>	YL86_YEAST Hypothetical 99.8 kDa protein in IKI3-RPS29...	<a href="#">38</a>	0.006
<input type="checkbox"/>	sp	<a href="#">P35194</a>	YBA4_YEAST Hypothetical 287.5 kDa protein in PDR3-HTA2...	<a href="#">37</a>	0.010
<input type="checkbox"/>	sp	<a href="#">P32917</a>	STE5_YEAST STE5 protein [STE5] [Saccharomyces cerevisi...	<a href="#">33</a>	0.15
<input type="checkbox"/>	sp	<a href="#">P32074</a>	COPG_YEAST Coatomer gamma subunit (Gamma-coat protein)...	<a href="#">32</a>	0.25
<input type="checkbox"/>	tr	<a href="#">Q12150</a>	Chromosome XII COSMID 9449 [CSF1] [Saccharomyces cerev...	<a href="#">32</a>	0.32

# P31383 Vs YEAST

<input type="checkbox"/>	sp	<a href="#">P31383</a>	2AAA_YEAST Protein phosphatase PP2A regulatory subunit...	<a href="#">1170</a>	0.0
<input type="checkbox"/>	sp	<a href="#">P33892</a>	GCN1_YEAST Translational activator GCN1 [GCN1] [Saccha...	<a href="#">47</a>	1e-05
<input type="checkbox"/>	sp	<a href="#">P53978</a>	EF3B_YEAST Elongation factor 3B (EF-3B) [YEF3B] [Sacch...	<a href="#">46</a>	2e-05
<input type="checkbox"/>	sp	<a href="#">P22219</a>	VP15_YEAST Protein kinase VPS15 (EC 2.7.1.-) [VPS15] [...	<a href="#">40</a>	0.001
<input type="checkbox"/>	sp	<a href="#">P32337</a>	IMB3_YEAST Importin beta-3 subunit (Karyopherin beta-3...	<a href="#">39</a>	0.002
<input type="checkbox"/>	sp	<a href="#">P49955</a>	S3B1_YEAST U2 snRNP component HSH155 [HSH155] [Sacchar...	<a href="#">38</a>	0.005
<input type="checkbox"/>	sp	<a href="#">Q06708</a>	YL86_YEAST Hypothetical 99.8 kDa protein in IKI3-RPS29...	<a href="#">38</a>	0.006
<input type="checkbox"/>	sp	<a href="#">P35194</a>	YBA4_YEAST Hypothetical 287.5 kDa protein in PDR3-HTA2...	<a href="#">37</a>	0.010
<input type="checkbox"/>	sp	<a href="#">P32917</a>	STE5_YEAST STE5 protein [STE5] [Saccharomyces cerevisi...	<a href="#">33</a>	0.15
<input type="checkbox"/>	sp	<a href="#">P32074</a>	COPG_YEAST Coatomer gamma subunit (Gamma-coat protein)...	<a href="#">32</a>	0.25
<input type="checkbox"/>	tr	<a href="#">Q12150</a>	Chromosome XII COSMID 9449 [CSF1] [Saccharomyces cerev...	<a href="#">32</a>	0.32

# P31383 Vs UniProt

<input type="checkbox"/>	sp	<a href="#">Q10178</a>	S3B1 SCHPO U2 snRNP component prp10 [PRP10] [Schizosac...	50	2e-04
<input type="checkbox"/>	tr	<a href="#">Q9FMF9</a>	Nuclear protein-like [Arabidopsis thaliana (Mouse-ear ...	49	2e-04
<input type="checkbox"/>	tr	<a href="#">Q17873</a>	Hypothetical protein F46C5.6 [F46C5.6] [Caenorhabditis...	49	3e-04
<input type="checkbox"/>	tr	<a href="#">Q7PLL6</a>	CG17514-PA.3 [CG17514] [Drosophila melanogaster (Fruit...	49	3e-04
<input type="checkbox"/>	tr	<a href="#">Q86JC5</a>	Hypothetical protein [Dictyostelium discoideum (Slime ...	49	3e-04
<input type="checkbox"/>	tr	<a href="#">Q7RM47</a>	Hypothetical protein [PYO2341] [Plasmodium yoelii yoelii]	48	4e-04
<input type="checkbox"/>	tr	<a href="#">Q42900</a>	Protein kinase with calcium binding domain [SPBC119.07...	48	6e-04
<input type="checkbox"/>	sp	<a href="#">P33892</a>	GCN1_YEAST Translational activator GCN1 [GCN1] [Saccha...	47	0.001
<input type="checkbox"/>	sp	<a href="#">P53978</a>	EF3B_YEAST Elongation factor 3B (EF-3B) [YEF3B] [Sacch...	46	0.002
<input type="checkbox"/>	tr	<a href="#">Q8C0Y0</a>	Hypothetical heat repeat containing protein [8430415E0...	46	0.002
<input type="checkbox"/>	tr	<a href="#">Q9CRR0</a>	8430415E04Rik protein (Fragment) [8430415E04RIK] [Mus ...	46	0.002
<input type="checkbox"/>	tr	<a href="#">Q7SDG6</a>	Hypothetical protein [NCU03042.1] [Neurospora crassa]	45	0.003
<input type="checkbox"/>	tr	<a href="#">Q7QEW2</a>	AgCP13260 (Fragment) [AGCG49048] [Anopheles gambiae st...	45	0.004
<input type="checkbox"/>	tr	<a href="#">Q9ECF0</a>	Hypothetical protein KIAA1622 (Fragment) [KIAA1622] [H...	45	0.004
<input type="checkbox"/>	tr	<a href="#">Q8VD65</a>	Hypothetical protein [PIK3R4] [Mus musculus (Mouse)]	45	0.005
<input type="checkbox"/>	tr	<a href="#">Q8C948</a>	Phosphatidylinositol 3 kinase [C730038E05RIK] [Mus mus...	45	0.005
<input type="checkbox"/>	tr	<a href="#">Q99570</a>	Adaptor protein [P150] [Homo sapiens (Human)]	45	0.005
<input type="checkbox"/>	sp	<a href="#">Q10105</a>	YAO5 SCHPO Putative translational activator C18G6.05c ...	44	0.006
<input type="checkbox"/>	tr	<a href="#">Q7XJN7</a>	At2g40730 protein [AT2G40730] [Arabidopsis thaliana (M...	44	0.006
<input type="checkbox"/>	tr	<a href="#">Q8LQE7</a>	Kinase-like protein [OJ1529_G03.14] [Oryza sativa (jap...	44	0.006
<input type="checkbox"/>	tr	<a href="#">Q21909</a>	T08A11.2 protein [T08A11.2] [Caenorhabditis elegans]	44	0.006
<input type="checkbox"/>	tr	<a href="#">Q9CXV4</a>	13 days embryo head cDNA, RIKEN full-length enriched l...	44	0.008
<input type="checkbox"/>	tr	<a href="#">Q7QF20</a>	AgCP13790 (Fragment) [AGCG51792] [Anopheles gambiae st...	44	0.008
<input type="checkbox"/>	sp	<a href="#">Q57683</a>	S3B1_XENLA Splicing factor 3B subunit 1 (Spliceosome a...	44	0.011
<input type="checkbox"/>	sp	<a href="#">Q99NB9</a>	S3B1_MOUSE Splicing factor 3B subunit 1 (Spliceosome a...	44	0.011
<input type="checkbox"/>	sp	<a href="#">Q75533</a>	S3B1_HUMAN Splicing factor 3B subunit 1 (Spliceosome a...	44	0.011
<input type="checkbox"/>	tr	<a href="#">Q9VPR5</a>	CG2807 protein [CG2807] [Drosophila melanogaster (Fru...	44	0.011
<input type="checkbox"/>	tr	<a href="#">Q8SQL2</a>	Protein phosphatase PP2-A regulatory subunit A [ECU09_...	43	0.014
<input type="checkbox"/>	tr	<a href="#">Q59191</a>	Hypothetical protein PH1522 [PH1522] [Pyrococcus horik...	43	0.014
<input type="checkbox"/>	tr	<a href="#">Q8II56</a>	PF16 protein, putative [PF11_0318] [Plasmodium falcipa...	42	0.031
<input type="checkbox"/>	tr	<a href="#">Q9V021</a>	Hypothetical protein PYRAB06470 [PYRAB06470] [Pyrococc...	42	0.031
<input type="checkbox"/>	tr	<a href="#">Q7T160</a>	SI:d2146N9.1 (Novel protein similar to human general c...	42	0.041
<input type="checkbox"/>	tr	<a href="#">Q7PY15</a>	EbiP8519 (Fragment) [EBI8519] [Anopheles gambiae str....	42	0.041
<input type="checkbox"/>	tr	<a href="#">Q86Y56</a>	Hypothetical protein FLJ20397 [Homo sapiens (Human)]	42	0.041
<input type="checkbox"/>	sp	<a href="#">P22219</a>	VP15_YEAST Protein kinase VPS15 (EC 2.7.1.-) [VPS15] [...	40	0.091

# Psi-BLAST

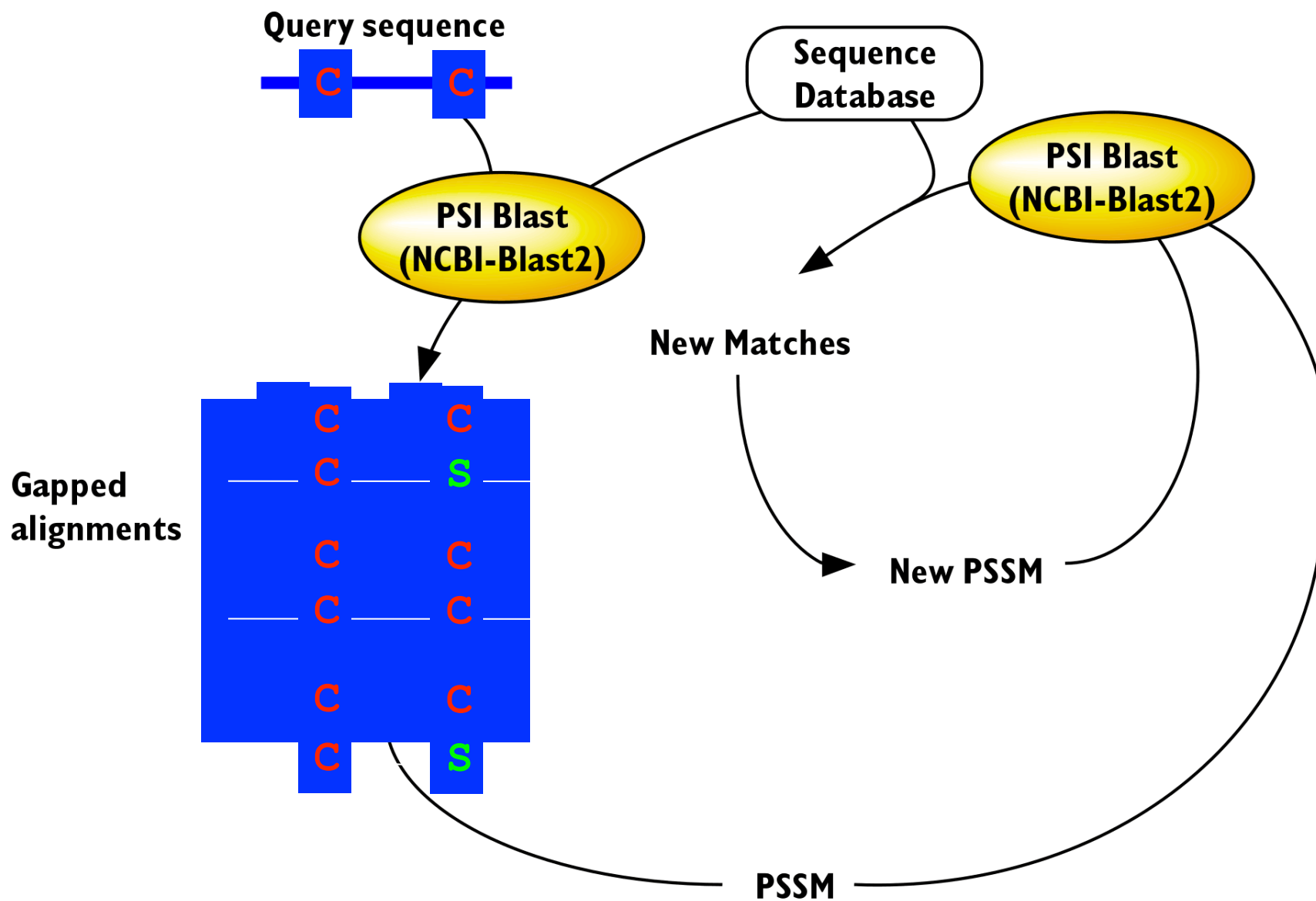
Position Specific Iterated version of BLAST.

Psi-BLAST performs a gapped BLAST database search.

Psi-BLAST program uses the information from any significant alignments returned to construct a position-specific score matrix, which replaces the query sequence for the next round of database searching.

Psi-BLAST may be iterated until no new significant alignments are found.

# Psi-BLAST Iteration





# BLAST: E-value

	M	Y	<b>C</b>	E	Q	U	E	N	<b>C</b>	E	S	.	.
A	0	2	-1	0	0	0	0	-1	0	-1	3		
S	-1	-1	-1	0	-1	0	0	0	<b>5</b>	-1	-1		
C	-1	-1	<b>10</b>	1	-1	0	0	5	<b>5</b>	4	-1		
.													
.													
Y	-1	6	-1	-1	-1	0	-1	-1	-1	-1	-1		
V	-1	1	-1	-1	-1	0	-1	-1	-1	1	-1		

[PubMed](#)
[Entrez](#)
[BLAST](#)
[OMIM](#)
[Taxonomy](#)
[Structure](#)

## Info

- [FAQs](#)
- [News](#)
- [References](#)
- [NCBI Contributors](#)

## Education

- [Program selection guide](#)
- [Tutorial](#)
- [URL API guide](#)

## Download

- [Databases](#)
- [Documentation](#)
- [Executables](#)
- [Source code](#)

## Support

- [Helpdesk](#)
- [Mailing list](#)

**NEW 12 May 2004** BLAST 2.2.9 has been released. [Read more...](#)

### Nucleotide

- [Discontiguous megablast](#)
- [Megablast](#)
- [Nucleotide-nucleotide BLAST \(blastn\)](#)
- [Search for short, nearly exact matches](#)
- [Search trace archives with megablast or discontiguous megablast](#)

### Protein

- [Protein-protein BLAST \(blastp\)](#)
- [PHI- and PSI-BLAST](#)
- [Search for short, nearly exact matches](#)
- [Search the conserved domain database \(rpsblast\)](#)
- [Search by domain architecture \(cdart\)](#)

### Translated

- [Translated query vs. protein database \(blastx\)](#)
- [Protein query vs. translated database \(tblastn\)](#)
- [Translated query vs. translated database \(tblastx\)](#)

### Genomes

- [Chicken, cow, pig, dog, sheep, cat](#) **NEW**
- [Environmental samples](#)
- [Human, mouse, rat](#)
- [Fugu rubripes, zebrafish](#)
- [Insects, nematodes, plants, fungi, malaria](#)
- [Microbial genomes, other eukaryotic genomes](#)

### Special

- [Search for gene expression data \(GEO BLAST\)](#)
- [Align two sequences \(bl2seq\)](#)
- [Screen for vector contamination \(VecScreen\)](#)
- [Immunoglobulin BLAST \(IgBlast\)](#)

### Meta

- [Retrieve results by RID](#)
- [Get this page with javascript-free links](#)

[Disclaimer](#)

[Privacy statement](#)

[Accessibility](#)

Valid [XHTML 1.0](#), [CSS](#).



# BLAST

PubMed

Entrez

BLAST

OMIM

Taxonomy

Structure

## Info

- FAQs
- News
- References
- NCBI Contributors

## Education

- Program selection guide
- Tutorial
- URL API guide

## Download

- Databases
- Documentation
- Executables
- Source code

## Support

- Helpdesk
- Mailing list

**NEW** 12 May 2004 BLAST 2.2.9 has been released. [Read more...](#)

### Nucleotide

- Discontiguous megablast
- Megablast
- Nucleotide-nucleotide BLAST (blastn)
- Search for short, nearly exact matches
- Search trace archives with [megablast](#) or [discontiguous megablast](#)

### Protein

- Protein-protein BLAST (blastp)
- PHI- and PSI-BLAST
- Search for short, nearly exact matches
- Search the conserved domain database ([rpsblast](#))
- Search by domain architecture ([cdart](#))

### Translated

- Translated query vs. protein database ([blastx](#))
- Protein query vs. translated database ([tblastn](#))
- Translated query vs. translated database ([tblastx](#))

### Genomes

- Chicken, cow, pig, dog, sheep, cat **NEW**
- Environmental samples
- Human, mouse, rat
- Fugu rubripes, zebrafish
- Insects, nematodes, plants, fungi, malaria
- Microbial genomes, other eukaryotic genomes

### Special

- Search for gene expression data (GEO BLAST)
- Align two sequences ([bl2seq](#))
- Screen for vector contamination ([VecScreen](#))
- Immunoglobulin BLAST ([IgBlast](#))

### Meta

- Retrieve results by RID
- Get this page with javascript-free links

[Disclaimer](#)

[Privacy statement](#)

[Accessibility](#)

Valid [XHTML 1.0](#), [CSS](#).

# Low Complexity Regions

Some genome sequences contain low-complexity regions. These can give false-positive hits.

## Example:

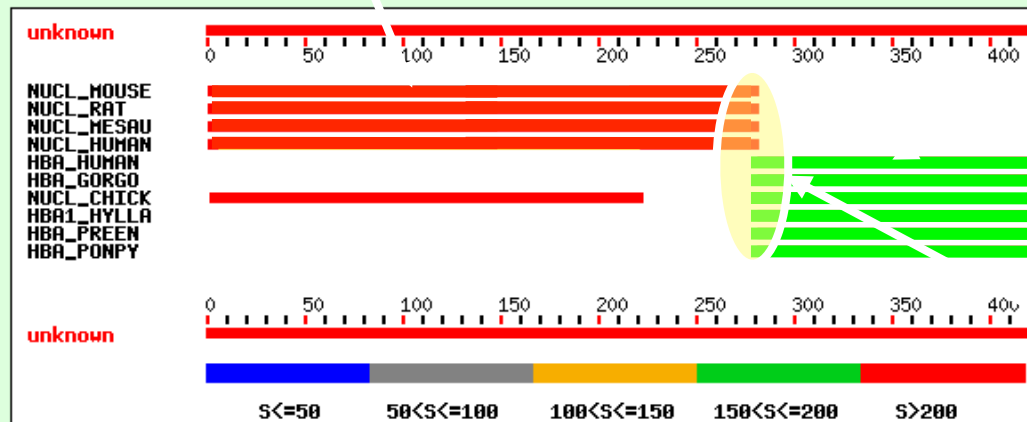
```
HSGDLPERTCPPCPPPCPPCPPPPCPPPCPCPPCPPPPPLWQPSSERTD
      |- low-complexity region -|
```

Most sequence searching programs use filters to recognise and skip such low-complexity regions. If such regions are by chance included in the hit, the output looks like:

```
HSGDLPERTXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXLWQPSSERTD
      |- low-complexity region -|
```

# Domain 1

# Domain 2



No Overlap

Sequences producing significant alignments:

Score E  
(bits) Value

sp P09405	NUCL_MOUSE 89505EE39C89F832	(NCL..)Nucleolin (Protein ...	425	e-119
sp P13383	NUCL_RAT 68774A214E550F90	(NCL..)Nucleolin (Protein C2...	407	e-113
sp P08199	NUCL_MESAU 79DDCF724CED7DB4	(NCL)Nucleolin (Protein C2...	397	e-110
sp P19338	NUCL_HUMAN 85A2F2CA22EA03DB	(NCL)Nucleolin (Protein C2...	371	e-102
sp P01922	HBA_HUMAN 34D13618E62A33C1	(HBA1..)Hemoglobin alpha ch...	285	2e-76
sp P01923	HBA_GORGO 25D13618E72A3306	(HBA)Hemoglobin alpha chain...	283	4e-76
sp P15771	NUCL_CHICK 7996C504BE9459A1	Nucleolin (Protein C23).[G...	283	5e-76
sp Q9TS35	HBA1_HYLLA 25D13618E36A7706	(HBA1)Hemoglobin alpha-1 c...	281	2e-75
sp P01924	HBA_PREEN 3771361D402A35C7	(HBA)Hemoglobin alpha chain...	280	5e-75
sp P06635	HBA_PONPY 37DE74049545CE88	(HBA)Hemoglobin alpha chain...	279	7e-75

# Tricky Problems

- Repeats
- Multi-domain proteins
- Low-complexity regions
- Redundancy
- Very short queries
- Very distant sequences
- Un-annotated sequences

# Post Processing of Search Results

- Match identification numbers and key words
- Grep sequences out of the database
- Compare results of different searches  
(on different databases)
- Filter sequences on specified criteria
- Align results using a multiple alignment
- Recognise family patterns or generate a family profile

# Typical Method Combinations

- Sequence search -> multiple alignment -> phylogenetic tree
- Sequence search -> sequence alignment -> homology modelling
- Sequence search -> sequence alignment -> mutation analysis
- Sequence search -> multiple alignment -> functional annotation



# Learning Outcomes

- Principles of Sequence Searching
- Substitution Matrix and Scoring
- BLAST heuristics: k-word matching and local alignment
- Extreme value distribution
- True/false positive/negative hit
- P-value and E-value
- Interpretation of hit list
- BLAST for nucleotides and proteins
- PHI- and PSI-BLAST