

*Current Topics in Genome Analysis  
Spring 2005*

*Week 4  
Biological Sequence Analysis I*

*Andy Baxevanis, Ph.D.*



Overview

---

- Week 4: Comparative methods and concepts
  - **Similarity vs. Homology**
  - Global vs. Local Alignments
  - Scoring Matrices
  - BLAST
  - BLAT
- Week 5: Predictive methods and concepts
  - Profiles, patterns, motifs, and domains
  - Secondary structure prediction
  - Structures: VAST, Cn3D, and *de novo* prediction



## Why do sequence alignments?

---

- Provide a measure of relatedness between nucleotide or amino acid sequences
- Determining relatedness allows one to draw biological inferences regarding
  - structural relationships
  - functional relationships
  - evolutionary relationships

→ *importance of using correct terminology*



## Defining the Terms

---

- The quantitative measure: *Similarity*
  - Always based on an observable
  - Usually expressed as percent identity
  - Quantify changes that occur as two sequences diverge
    - substitutions
    - insertions
    - deletions
  - Identify residues crucial for maintaining a protein's structure or function
- High degrees of sequence similarity *might* imply
  - a common evolutionary history
  - possible commonality in biological function



## Defining the Terms

---

- The conclusion: **Homology**
  - Genes *are* or *are not* homologous (not measured in degrees)
  - Homology implies an evolutionary relationship
- The term “homolog” may apply to the relationship
  - between genes separated by the event of speciation (*orthology*)
  - between genes separated by the event of genetic duplication (*paralogy*)



## Defining the Terms

---

- Orthologs
  - Sequences are direct descendants of a sequence in a common ancestor
  - Most likely have similar domain structure, three-dimensional structure, and biological function
- Paralogs
  - Related through a gene duplication event
  - Provides insight into “evolutionary innovation” (adapting a pre-existing gene product for a new function)



## Defining the Terms

*Orthologs*



*Most recent common ancestor*

α



## Defining the Terms

*Paralogs*

*Orthologs*



*Most recent common ancestor*

α

*Gene duplication*

β

- Genes 1-3 are orthologous
- Genes 4-6 are orthologous
- Any pair of α and β genes are paralogous (genes related through a gene duplication event)



## Overview

---

- Week 4: Comparative methods and concepts
  - Similarity vs. Homology
  - Global vs. Local Alignments
  - Scoring Matrices
  - BLAST
  - BLAT
- Week 5: Predictive methods and concepts
  - Profiles, patterns, motifs, and domains
  - Secondary structure prediction
  - Structures: VAST, Cn3D, and *de novo* prediction



## Global Sequence Alignments

---

- Sequence comparison along the entire length of the two sequences being aligned
- Best for highly-similar sequences of similar length
- As the degree of sequence similarity declines, global alignment methods tend to miss important biological relationships



## Local Sequence Alignments

---

- Sequence comparison intended to find the most similar regions in the two sequences being aligned (“paired subsequences”)
- Regions outside the area of local alignment are excluded
- More than one local alignments could be generated for any two sequences being compared
- Best for sequences that share some similarity, or for sequences of different lengths



## Overview

---

- Week 4: Comparative methods and concepts
  - Similarity vs. Homology
  - Global vs. Local Alignments
  - Scoring Matrices
  - BLAST
  - BLAT
- Week 5: Predictive methods and concepts
  - Profiles, patterns, motifs, and domains
  - Secondary structure prediction
  - Structures: VAST, Cn3D, and *de novo* prediction



## Scoring Matrices

---

- Empirical weighting scheme to represent biology (side chain chemistry, structure, and function)
  - Cys/Pro important for structure and function
  - Trp has bulky side chain
  - Lys/Arg have positively-charged side chains



## Scoring Matrices

---

- **Conservation:** What residues can substitute for another residue and not adversely affect the function of the protein?
  - Ile/Val - both small and hydrophobic
  - Ser/Thr - both polar
  - *Conserve charge, size, hydrophobicity, other physicochemical factors*
- **Frequency:** How often does a particular residue occur amongst the entire constellation of proteins?



## Scoring Matrices

- Importance of understanding scoring matrices
  - Appear in all analyses involving sequence comparison
  - Implicitly represent particular evolutionary patterns
  - Choice of matrix can strongly influence outcomes



## Matrix Structure: Nucleotides

	A	T	G	C	S	W	R	Y	K	M	B	V	H	D	N
A	5	-4	-4	-4	-4	1	1	-4	-4	1	-4	-1	-1	-1	-2
T	-4	5	-4	-4	-4	1	-4	1	1	-4	-1	-4	-1	-1	-2
G	-4	-4	5	-4	1	-4	1	-4	1	-4	-1	-1	-4	-1	-2
C	-4	-4	-4	5	1	-4	-4	1	-4	1	-1	-1	-1	-4	-2
S	-4	-4	1	1	-1	-4	-2	-2	-2	-2	-1	-1	-3	-3	-1
W	1	1	-4	-4	-4	-1	-2	-2	-2	-2	-3	-3	-1	-1	-1
R	1	-4	1	-4	-2	-2	-1	-4	-2	-2	-3	-1	-3	-1	-1
Y	-4	1	-4	1	-2	-2	-4	-1	-2	-2	-1	-3	-1	-3	-1
K	-4	1	1	-4	-2	-2	-2	-2	-1	-4	-1	-3	-3	-1	-1
M	1	-4	-4	1	-2	-2	-2	-2	-4	-1	-3	-1	-1	-3	-1
B	-4	-1	-1	-1	-1	-3	-3	-1	-1	-3	-1	-2	-2	-2	-1
V	-1	-4	-1	-1	-1	-3	-1	-3	-3	-1	-2	-1	-2	-2	-1
H	-1	-1	-4	-1	-3	-1	-3	-1	-3	-1	-2	-2	-1	-2	-1
D	-1	-1	-1	-4	-3	-1	-1	-3	-1	-3	-2	-2	-2	-1	-1
N	-2	-2	-2	-2	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1

- Simple match/mismatch scoring scheme
- Assumes each nucleotide occurs 25% of the time





## Matrix Structure: Proteins

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

BLOSUM62

## PAM Matrices

- Margaret Dayhoff and colleagues, 1978
  - Look at patterns of substitutions in highly related proteins (> 85% similar) within multiple sequence alignments
  - Analysis documented 1572 changes in 71 groups of proteins examined
  - Substitution tables constructed based on results
  - Given high degree of similarity within original sequence set, results represent substitution pattern that would be expected over short evolutionary distances

## PAM Matrices

---

- Short evolutionary distance  
∴ change in function unlikely
- Point Accepted Mutation (PAM)
  - The new side chain must function the same way as the old one (“acceptance”)
  - On average, 1 PAM corresponds to 1 amino acid change per 100 residues
  - 1 PAM ~ 1% divergence
  - Extrapolate to predict patterns at longer evolutionary distances



## PAM Matrices: Assumptions

---

- All sites assumed to be equally mutable
- Replacement of amino acids is independent of previous mutations at the same position
- Replacement is independent of surrounding residues
- Forces responsible for sequence evolution over shorter time spans are the same as those over longer time spans



## PAM Matrices: Sources of Error

---

- Small, globular proteins of average composition used to derive matrices
- Errors in PAM 1 are magnified up to PAM 250 (only PAM 1 is based on direct observation)
- Does not account for conserved blocks or motifs



## BLOSUM Matrices

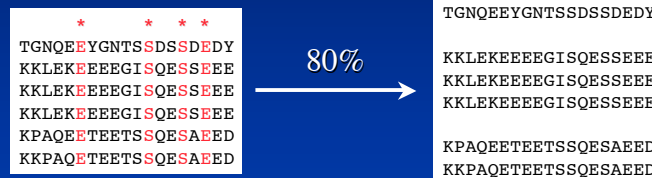
---

- Henikoff and Henikoff, 1992
- Blocks Substitution Matrix
  - Look only for differences in conserved, ungapped regions of a protein family (“blocks”)
  - Directly calculated, using no extrapolations
  - More sensitive to detecting structural or functional substitutions
  - Generally perform better than PAM matrices for local similarity searches (*Henikoff and Henikoff, 1993*)



## BLOSUM $n$

- Calculated from sequences sharing no more than  $n\%$  identity
- Contribution of sequences  $> n\%$  identical clustered and weighted to 1



*A+T Hook Domain (Block IPB000637B)*

2,000 blocks representing  $> 500$  groups of related proteins

## BLOSUM $n$

- Clustering reduces contribution of closely-related sequences (less bias towards substitutions that occur in the most closely related members of a family)
- Substitution frequencies are more heavily-influenced by sequences that are more divergent than this cutoff
- Reducing  $n$  yields more distantly-related sequences

## So many matrices...

---

### Triple-PAM Strategy (*Altschul, 1991*)

PAM 40	Short alignments, highly similar	70-90%
PAM 160	Detecting known members of a protein family	50-60%
PAM 250	Longer, weaker local alignments	~ 30%

### BLOSUM (*Henikoff, 1993*)

BLOSUM 90	Short alignments, highly similar	70-90%
BLOSUM 80	Detecting known members of a protein family	50-60%
<b>BLOSUM 62</b>	<b>Most effective in finding all potential similarities</b>	<b>30-40%</b>
BLOSUM 30	Longer, weaker local alignments	< 30%



## So many matrices...

---

- Matrix Equivalencies

PAM 250	~	BLOSUM 45
PAM 160	~	BLOSUM 62
PAM 120	~	BLOSUM 80

- Specialized matrices

- Transmembrane proteins
- Species-specific matrices



*Wheeler, 2003*

## So many matrices...

---

*No single matrix is  
the complete answer for  
all sequence comparisons*



## Gaps

---

- Compensate for insertions and deletions
- Used to improve alignments between two sequences
- Must be kept to a reasonable number, to not reflect a biological implausible scenario (~1 gap per 20 residues good rule-of-thumb)
- Cannot be scored simply as a “match” or a “mismatch”



## Affine Gap Penalty

Fixed deduction for introducing a gap *plus*  
an additional deduction proportional to the length of the gap

$$\text{Deduction for a gap} = G + Ln$$

		nuc	pro
where	$G$ = gap-opening penalty	5	11
	$L$ = gap-extension penalty	2	1
and	$n$ = length of the gap		

Can adjust scores to make gap insertion more or less permissive, but most programs will use values of  $G$  and  $L$  most appropriate for the scoring matrix selected

## Overview

- Week 4: Comparative methods and concepts
  - Similarity *vs.* Homology
  - Global *vs.* Local Alignments
  - Scoring Matrices
  - **BLAST**
  - BLAT
- Week 5: Predictive methods and concepts
  - Profiles, patterns, motifs, and domains
  - Secondary structure prediction
  - Structures: VAST, Cn3D, and *de novo* prediction

## BLAST

---

- Basic Local Alignment Search Tool
- Seeks high-scoring segment pairs (HSP)
  - pair of sequences that can be aligned without gaps
  - when aligned, have maximal aggregate score (score cannot be improved by extension or trimming)
  - score must be above score threshold  $S$
  - gapped or ungapped
- Results not limited to the “best HSP” for any given sequence pair



## BLAST Algorithms

---

<i>Program</i>	<i>Query Sequence</i>	<i>Target Sequence</i>
<b>BLASTN</b>	<b>Nucleotide</b>	<b>Nucleotide</b>
<b>BLASTP</b>	<b>Protein</b>	<b>Protein</b>
<b>BLASTX</b>	<b>Nucleotide, six-frame translation</b>	<b>Protein</b>
TBLASTN	Protein	Nucleotide, six-frame translation
TBLASTX	Nucleotide, six-frame translation	Nucleotide, six-frame translation





## Neighborhood Words

Query Word ( $W = 3$ )

Query: GSQSLAALLNKCKT **PQG** QRLVNQWIKQPLMDKNRIEERLNLVEAFVED

Neighborhood Words

PQG	18	= 7 + 5 + 6
PEG	15	
PRG	14	
PKG	14	
PNG	13	
PDG	13	
PHG	13	
PMG	13	
PSG	13	
PQA	12	
PQN	12	
etc.		

Neighborhood Score Threshold ( $T = 13$ )



## High-Scoring Segment Pairs

PQG	18
PEG	15
PRG	14
PKG	14
PNG	13
PDG	13
PHG	13
PMG	13
PSG	13
PQA	12
PQN	12
etc.	

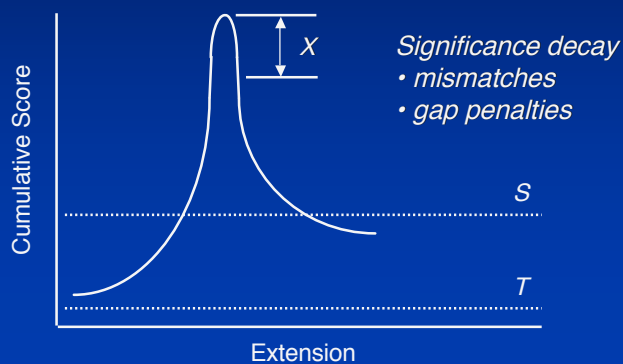
Query: 325 SLAALLNKCKT **PQG** QRLVNQWIKQPLMDKNRIEERLNLVEA 365  
 +LA++L TP+G R++ +W+ +P+ D + ER + A  
 Sbjct: 290 TLASVLDCTVT **PMG** SRMLKRWLHMPVRDTRVLLERQQTIGA 330



## Extension

←───────────────────▶

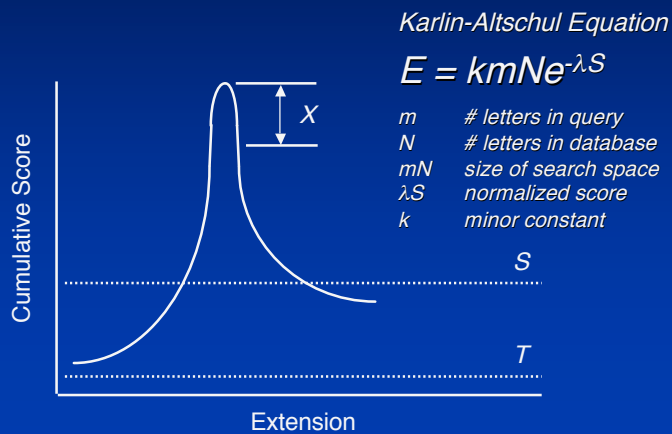
Query:	325	SLAALLNKCKT <b>PQG</b> QRLVNQWIKQPLMDKNRIEERLNLVEA	365
		+LA++L <b>TP+G</b> R++ +W+ +P+ D    + ER    + A	
Sbjct:	290	TLASVLDCTVT <b>PMG</b> SRMLKRWLHMPVRDTRVLLERQQTIGA	330



## Scores and Probabilities

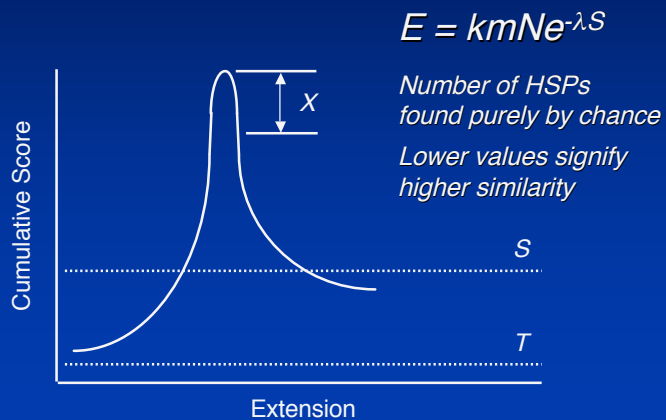
←───────────────────▶

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



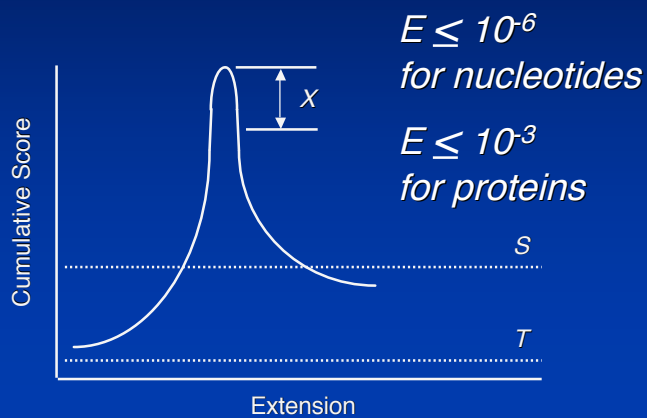
## Scores and Probabilities

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

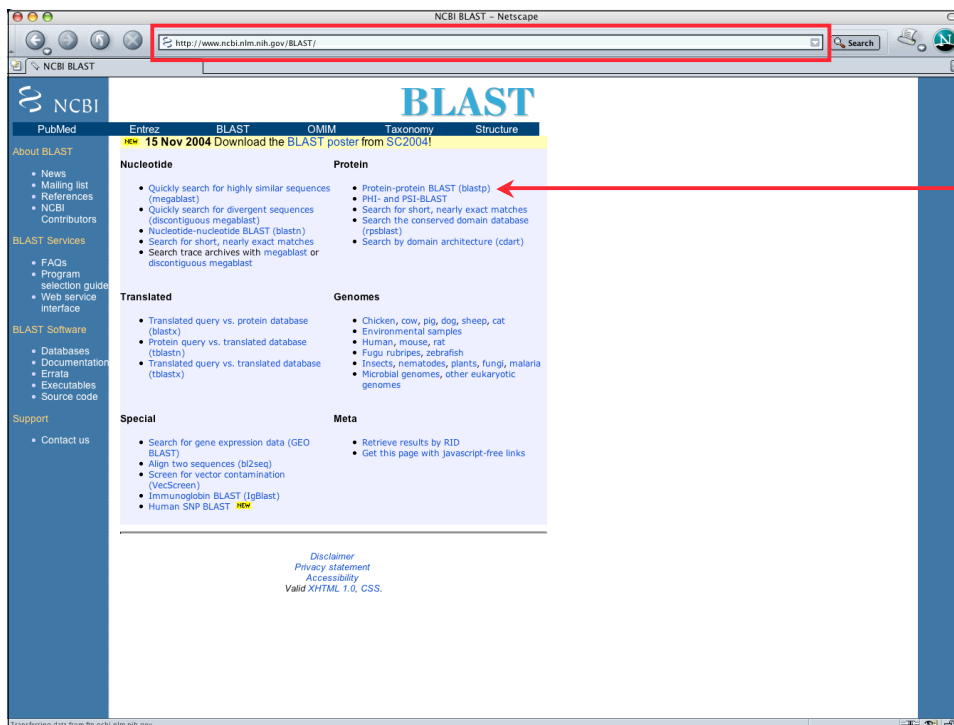
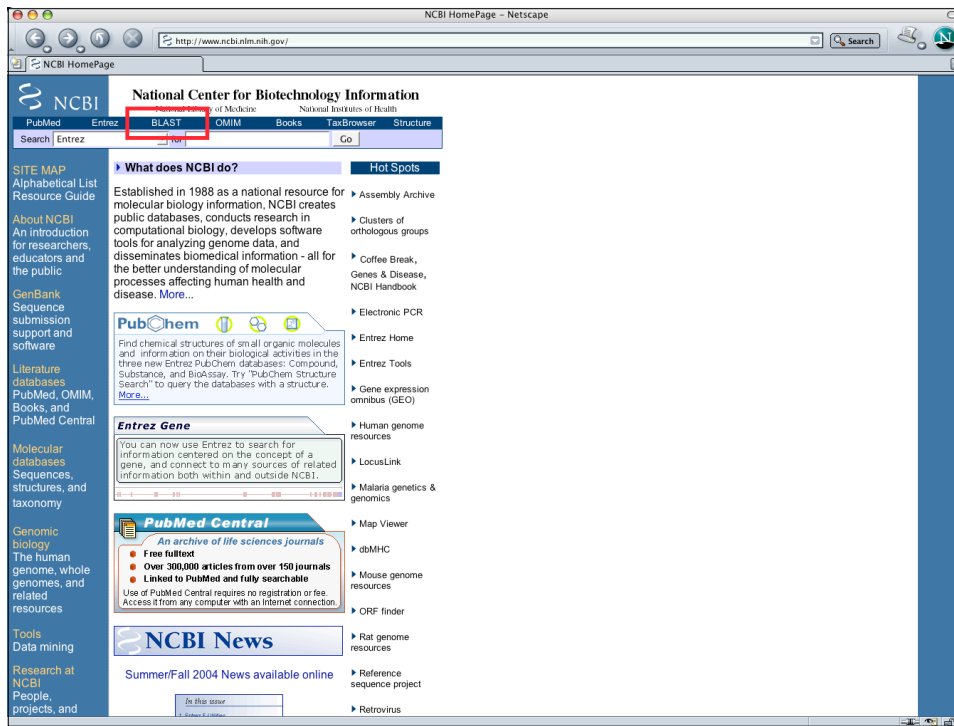


## Scores and Probabilities

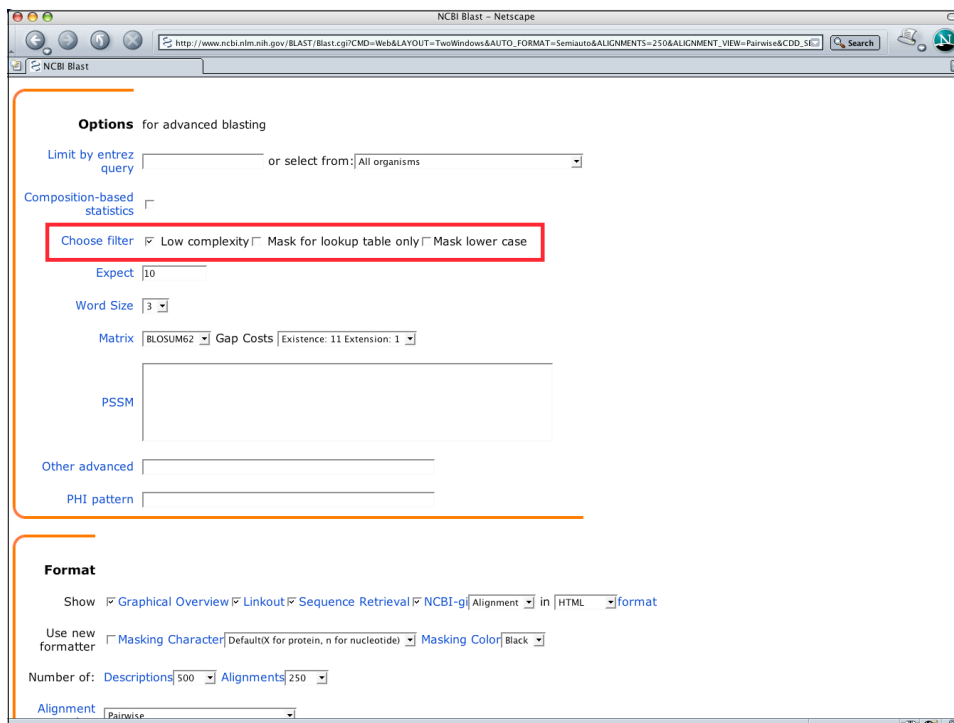
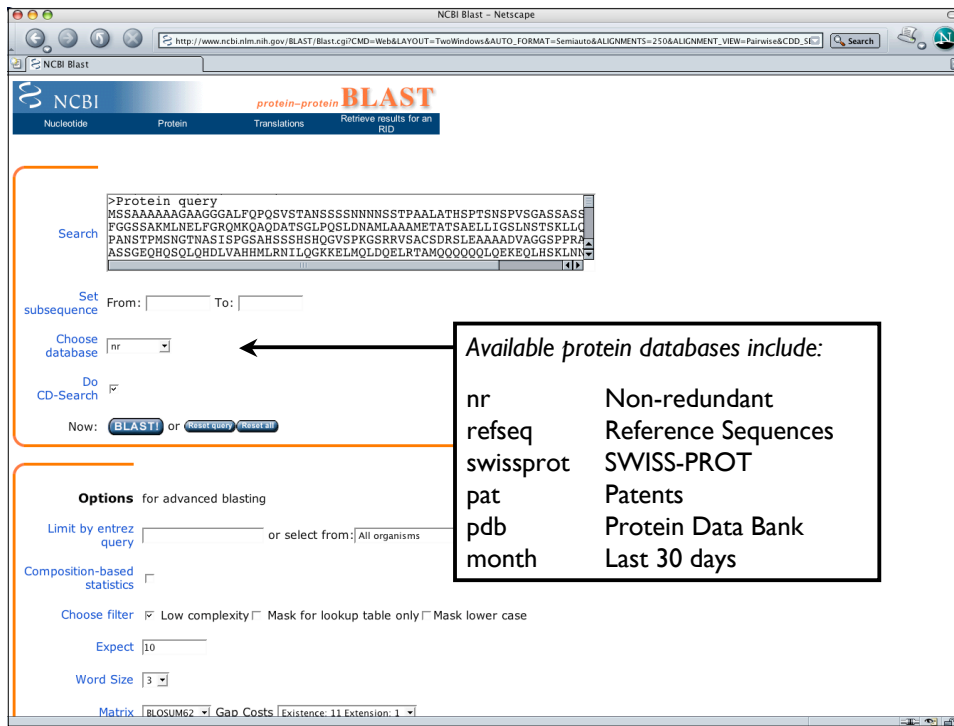
Query:	325	SLAALLNKCKT <b>PQG</b> QRLVNQWIKQPLMDKNRIEERLNLVEA	365
		+LA++L TP+G R++ +W+ +P+ D + ER + A	
Sbjct:	290	TLASVLDCTVT <b>PMG</b> SRMLKRWLHMPVRDTRVLLERQQTIGA	330



NHGRI Current Topics in Genome Analysis 2005  
 Biological Sequence Analysis I



NHGRI Current Topics in Genome Analysis 2005  
 Biological Sequence Analysis I



## Low-Complexity Regions

Defined as regions of biased composition

- Homopolymeric runs
- Short-period repeats
- Subtle over-representation of several residues

```
>gi|20455478|sp|P50553|ASC1_HUMAN_Achaete-scute_homolog_1_(HASH1)
MESSAKMESGGAGQQPQPQPQPFLPPAACFFAIAAAAAAAAAAAQSAQQQQQQQQQQQAPQLRPAA
DQQPSSGGGHSAPKQVKRQRSSSPELMRCKRRLLNFSGFCYSLPQQQIAAVARRNERERNRVKLVNLGFAT
LREHVPNGAANKKMSKVETLRSAVEYIRALQQLLDEHDVSAAFQACVLSPTISPNYSNDLNSMAGSPVS
SYSSDEGSYDPLSPPEEQELLDFTNWF
```

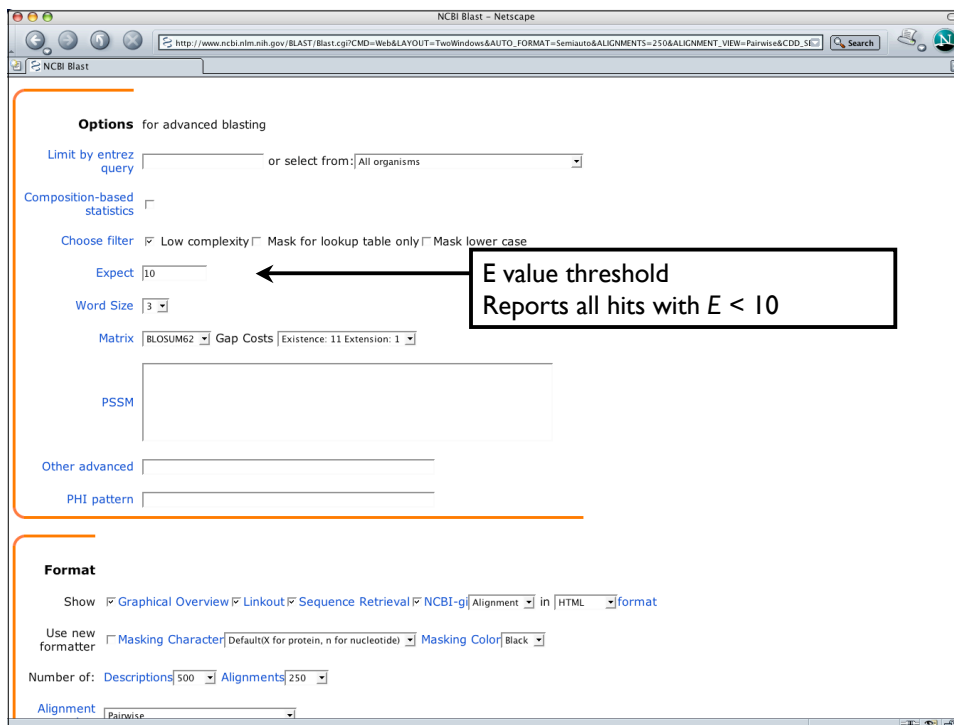
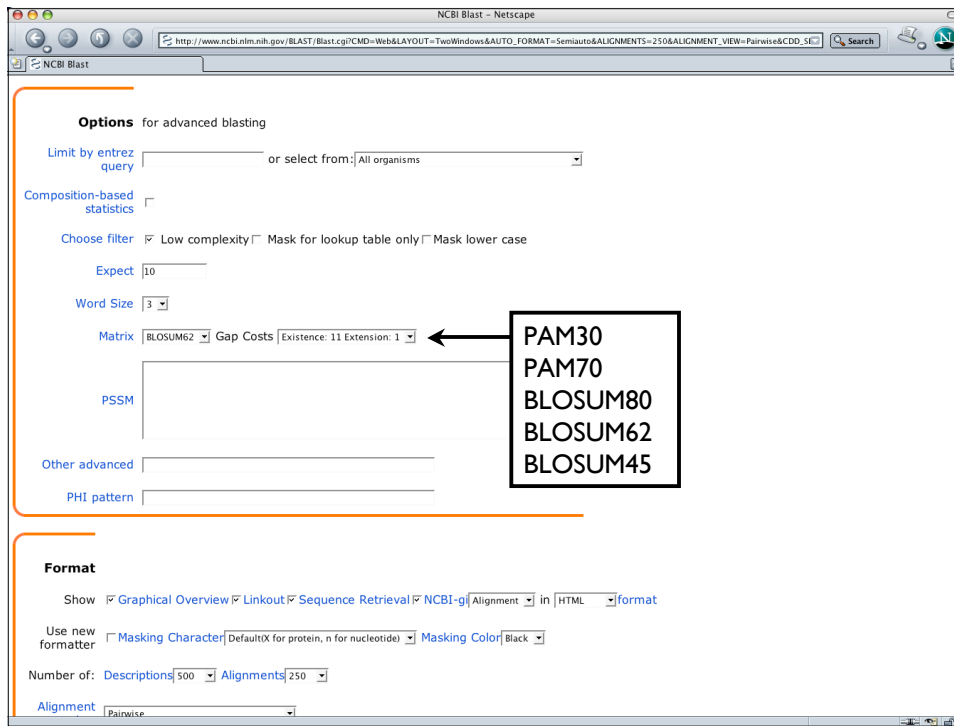
*Homopolymeric  
alanine-glutamine tract*



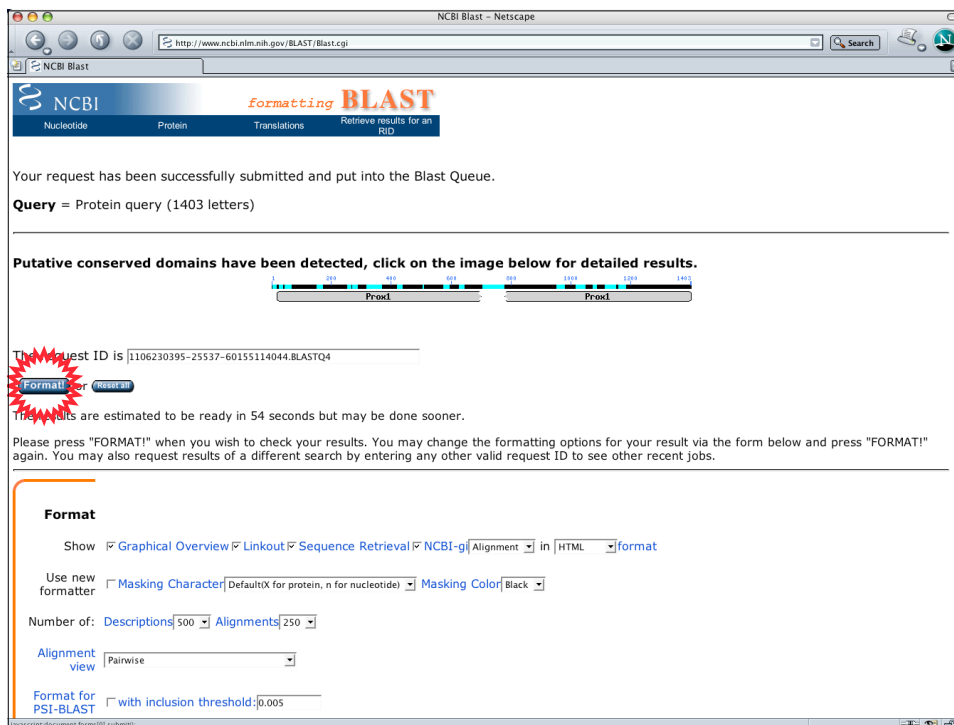
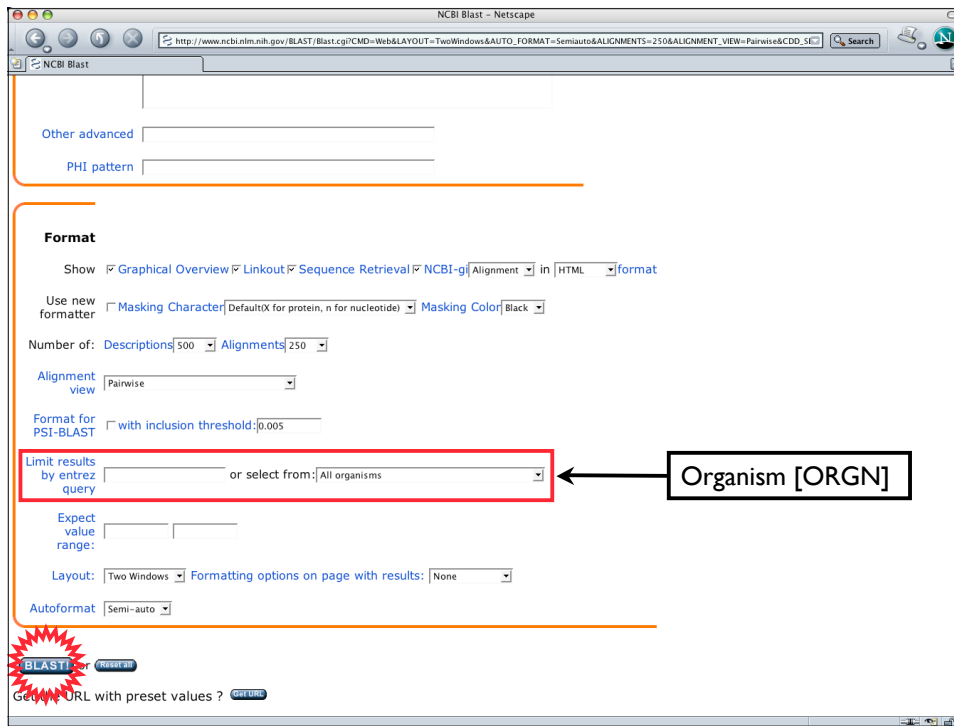
## Identifying Low-Complexity Regions

- Biological origins and role not well-understood
  - DNA replication errors (polymerase slippage)?
  - Unequal crossing-over?
- May confound sequence analysis
  - BLAST relies on uniformly-distributed amino acid frequencies
  - Often lead to false positives
  - Filtering is advised (and usually enabled by default)

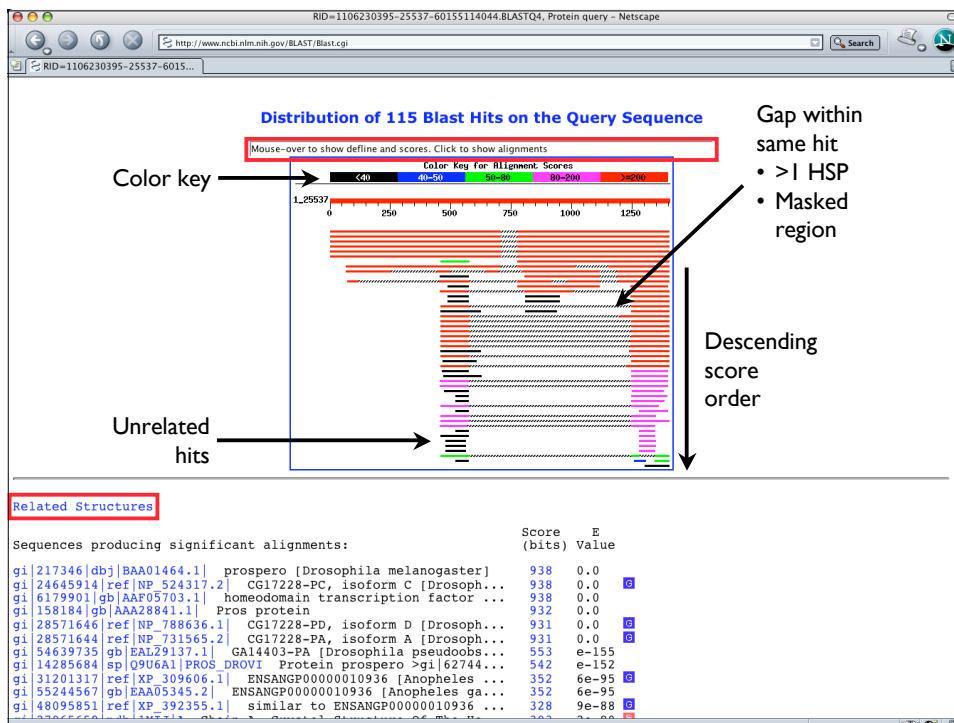
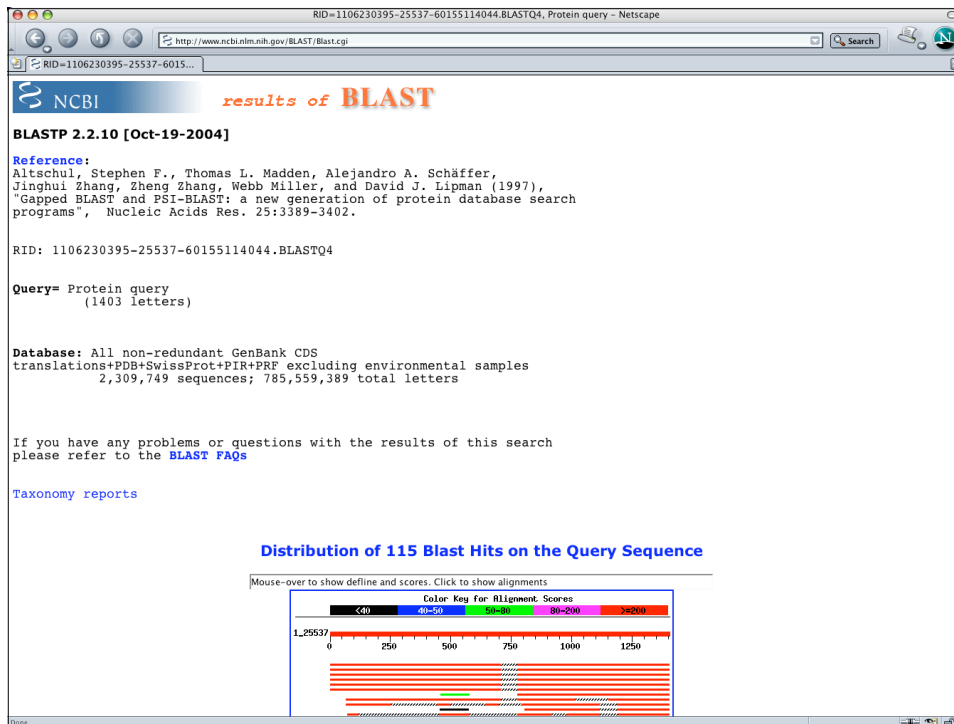




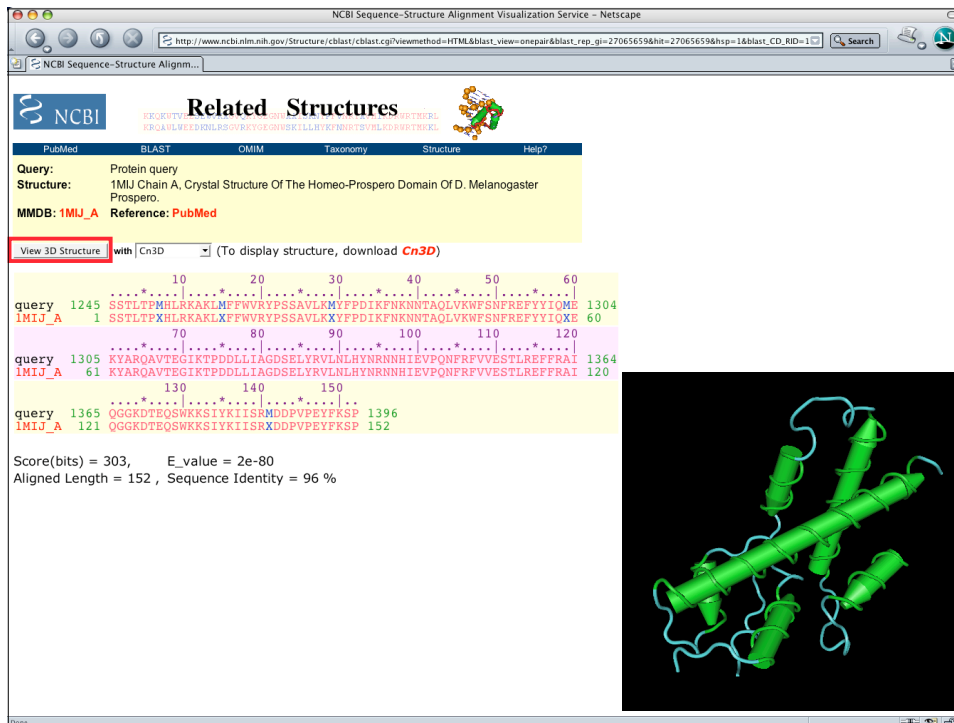
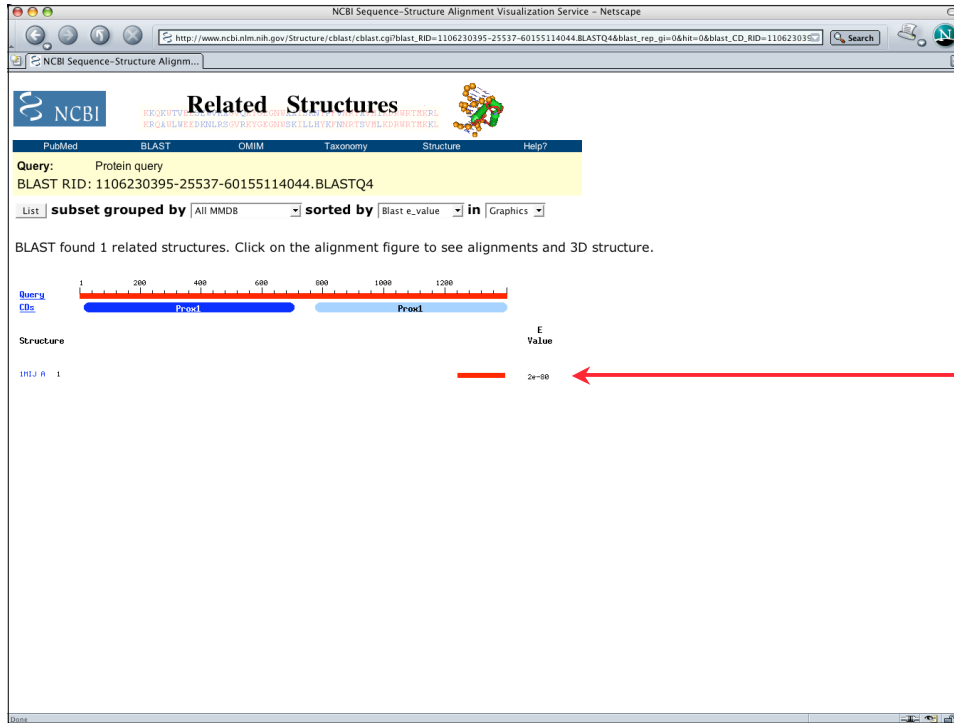
NHGRI Current Topics in Genome Analysis 2005  
Biological Sequence Analysis I







NHGRI Current Topics in Genome Analysis 2005  
 Biological Sequence Analysis I



NHGRI Current Topics in Genome Analysis 2005  
Biological Sequence Analysis I

Sequences producing significant alignments:

gi	Accession	Description	Score	E Value
gi 217346	dbj BAA01464.1	prospero [Drosophila melanogaster]	938	0.0
gi 24645914	ref NP_524317.2	CG17228-PC, isoform A [Drosophila melanogaster]	938	0.0
gi 158184	gb AAZ28841.1	Pros protein [Drosophila melanogaster]	932	0.0
gi 28571646	ref NP_788636.1	CG17228-PD, isoform A [Drosophila melanogaster]	931	0.0
gi 28571644	ref NP_731565.2	CG17228-PA, isoform A [Drosophila melanogaster]	931	0.0
gi 54639735	gb EAL29137.1	GAL4403-PA [Drosophila pseudoobscura]	553	e-155
gi 14285684	sp Q9U6A1	PROS_DROVI Protein prospero >gi 62744... [Drosophila melanogaster]	542	e-152
gi 31291317	ref XP_309606.1	ENSANGP0000010936 [Anopheles gambiae]	352	6e-95
gi 55244567	gb EAA05345.2	ENSANGP0000010936 [Anopheles gambiae]	352	6e-95
gi 48095851	ref XP_392355.1	similar to ENSANGP0000010936... [Anopheles gambiae]	328	9e-88
gi 27065659	pdb 1MIJ	Chain A, Crystal Structure Of The Homeobox Protein Prospero [Cupiennius saltator]	303	2e-80
gi 32261038	emb CAE00181.1	prospero protein [Cupiennius saltator]	279	4e-73
gi 1678018	gb AAL28228.1	GH11848p [Drosophila melanogaster]	273	3e-71
gi 39587414	emb CAE75068.1	hypothetical protein CB22984 [Cupiennius saltator]	243	4e-62
gi 17552742	ref NP_498760.1	C. Elegans Homeobox (ceh-26) [Caenorhabditis elegans]	239	4e-61
gi 3024449	sp Q92786	PRX1 HUMAN Homeobox prospero-like protein [Homo sapiens]	212	6e-53
gi 546374	gb AAB30541.1	Prox 1=homeobox gene prospero homolog [Drosophila melanogaster]	212	6e-53
gi 55589302	ref XP_514189.1	PREDICTED: similar to prospero... [Drosophila melanogaster]	212	8e-53
gi 21359846	ref NP_002754.2	prospero-related homeobox 1 [Homo sapiens]	211	2e-52
gi 6679483	ref NP_032963.1	prospero-related homeobox 1 [Mus musculus]	211	2e-52
gi 56785422	ref NP_001005616.1	PROX 1 protein [Gallus gallus]	211	2e-52
gi 7512233	pir JC5495	Prox 1 protein - chicken [Gallus gallus]	211	2e-52
gi 40254702	ref NP_571480.2	prospero-related homeobox gene... [Drosophila melanogaster]	208	8e-52
gi 3834411	gb AAC70926.1	homeodomain protein [Danio rerio]	208	8e-52
gi 57090743	ref XP_547908.1	PREDICTED: similar to RIKEN cDNA... [Xenopus laevis]	207	1e-51
gi 11071924	dbj BAB17310.1	Prox 1 [Xenopus laevis]	207	2e-51
gi 30424822	ref NP_780407.1	RIKEN cDNA 1700058C01 [Mus musculus]	205	7e-51
gi 27680210	ref NP_223067.1	similar to prospero-related homeobox... [Tetraodon lineatus]	200	2e-49
gi 47205868	emb CAF52934.1	unnamed protein product [Tetraodon lineatus]	191	2e-46
gi 47227457	emb CAG04605.1	unnamed protein product [Tetraodon lineatus]	188	1e-45
gi 47230216	emb CAG10630.1	unnamed protein product [Tetraodon lineatus]	182	8e-44
gi 47206446	emb CAF95276.1	unnamed protein product [Tetraodon lineatus]	182	8e-44
gi 3372869	gb AAC28353.1	Prox1 [Xenopus laevis]	178	9e-43
gi 47224292	emb CAG09138.1	unnamed protein product [Tetraodon lineatus]	172	7e-41
gi 1117962	gb AAC59781.1	prospero like protein [Tetraodon lineatus]	152	5e-35
gi 21753053	dbj BAC04278.1	unnamed protein product [Homo sapiens]	151	1e-34
gi 11071926	dbj BAB17311.1	Prox 1 [Cynops pyrrhogaster]	151	1e-34
gi 55961898	emb CAI15309.1	OTTHUMP0000061061 [Homo sapiens]	142	6e-32
gi 57089333	ref XP_514741.1	PREDICTED: similar to prospero... [Drosophila melanogaster]	140	2e-31
gi 47224321	emb CAG09167.1	unnamed protein product [Tetraodon lineatus]	139	8e-31
gi 47204095	emb CAG13403.1	unnamed protein product [Tetraodon lineatus]	96	8e-18
gi 34935368	ref XP_234418.2	similar to Homeobox prospero-like protein... [Drosophila melanogaster]	95	1e-17
gi 55641159	ref XP_522907.1	PREDICTED: similar to RIKEN cDNA... [Drosophila melanogaster]	94	2e-17

Sequences producing significant alignments:

gi	Accession	Description	Score	E Value
gi 55641159	ref XP_522907.1	PREDICTED: similar to RIKEN cDNA... [Drosophila melanogaster]	94	2e-17
gi 51493257	ref XP_372528.3	PREDICTED: hypothetical protein... [Drosophila melanogaster]	93	5e-17
gi 4809335	gb AAD30180.1	homeobox prospero-like protein [Homo sapiens]	91	2e-16
gi 7512234	pir JC5496	Prox 1 protein - chicken [Gallus gallus]	76	6e-12
gi 50749012	ref XP_426445.1	PREDICTED: similar to Homeobox... [Drosophila melanogaster]	62	1e-07
gi 6466795	gb AAFI3029.1	transcription factor Prox1 [Notopneustes	54	3e-05
gi 47202992	emb CAF94749.1	unnamed protein product [Tetraodon lineatus]	41	0.23
gi 40743593	gb EAA62783.1	hypothetical protein AN5690.2 [Acanthamoeba castellanii]	39	1.1
gi 30024062	ref NP_268084.2	hypothetical protein L188798 [Drosophila melanogaster]	39	1.1
gi 12724966	gb AAK06025.1	HYPOTHETICAL PROTEIN [Lactococcus lactis]	39	1.1
gi 47230218	emb CAG10632.1	unnamed protein product [Tetraodon lineatus]	39	1.5
gi 50365279	ref YP_053704.1	putative multidrug ABC transporter... [Staphylococcus aureus]	38	2.6
gi 56467270	gb EAL45239.1	calponin homology domain protein... [Staphylococcus aureus]	38	2.6
gi 56464504	gb EAL42983.1	calponin homology domain protein... [Staphylococcus aureus]	38	2.6
gi 383763	prf 1904201A	enn4 gene [Drosophila melanogaster]	37	3.3
gi 24662634	ref NP_648457.1	CG6175-PB [Drosophila melanogaster]	37	4.4
gi 37533942	ref NP_921273.1	Conserved hypothetical protein... [Drosophila melanogaster]	37	4.4
gi 21430574	gb AAM50965.1	RED0799p [Drosophila melanogaster]	37	4.4
gi 37676035	ref NP_936431.1	TPR repeat containing protein... [Drosophila melanogaster]	37	5.7
gi 10039425	dbj BAB13348.1	ALR protein [Equus caballus]	37	5.7
gi 57285852	gb AAW37946.1	ATP-dependent Clp protease, ATP-binding... [Staphylococcus aureus]	36	7.4
gi 49483136	ref YP_040360.1	putative ATPase subunit of an... [Staphylococcus aureus]	36	7.4
gi 21282586	ref NP_645674.1	hypothetical protein MW0857 [Staphylococcus aureus]	36	7.4
gi 15926564	ref NP_374097.1	hypothetical protein SA0835 [Staphylococcus aureus]	36	7.4
gi 51091708	dbj BAD36509.1	putative SMA-9 class B (Oryza sativa)	36	7.4
gi 27467592	ref NP_764229.1	clpB protein [Staphylococcus aureus]	36	9.7
gi 24653358	ref NP_610871.1	CG4744-PA [Drosophila melanogaster]	36	9.7
gi 14029840	gb AAK52834.1	cytoplasmic polyadenylation element... [Drosophila melanogaster]	36	9.7
gi 21428884	gb AAM50161.1	GH12467p [Drosophila melanogaster]	36	9.7
gi 25009677	gb AAN71015.1	AT02321p [Drosophila melanogaster]	36	9.7

Alignments

Get selected sequences | Select all | Deselect all

>gi|217346|dbj|BAA01464.1| prospero [Drosophila melanogaster]  
Length = 1403

Score = 938 bits (2424), Expect = 0.0  
Identities = 493/627 (78%), Positives = 493/627 (78%)

Query: 777 HVATAAPRQMHHPAPARLPTRMGAAGHTALKSELSEKFMRLRANNSSMMRMSGTDLE 836  
HVATAAPRQMHHPAPARLPTRMGAAGHTALKSELSEKFMRLRANNSSMMRMSGTDLE  
Sbjct: 777 HVATAAPRQMHHPAPARLPTRMGAAGHTALKSELSEKFMRLRANNSSMMRMSGTDLE 836

NHGRI Current Topics in Genome Analysis 2005  
 Biological Sequence Analysis I

RID=1106230395-25537-60155114044.BLASTQ4, Protein query - Netscape  
 http://www.ncbi.nlm.nih.gov/BLAST/Blast.cgi#6179901

>gi|6179901|gb|AAF05703.1| homeodomain transcription factor Prospero [Drosophila mel...]  
 Length = 1403

Score = 938 bits (2424), Expect = 0.0  
 Identities = 493/627 (78%), Positives = 493/627 (78%)

Query: 777 HVATAAPRPMQHHPAPARLPTRMGAAGHTALKSESEKQFQMLRANNSSMMRMSGTDLE 836  
 HVATAAPRPMQHHPAPARLPTRMGAAGHTALKSESEKQFQMLRANNSSMMRMSGTDLE 836  
 Sbjct: 777 HVATAAPRPMQHHPAPARLPTRMGAAGHTALKSESEKQFQMLRANNSSMMRMSGTDLE 836

Query: 837 GLADVLKSEITTSLSALVDTIVTRFVHORRLFSKQADSVTAAEQLNKDLLASQILDRK 896  
 GLADVLKSEITTSLSALVDTIVTRFVHORRLFSKQADSVTAAEQLNKDLLASQILDRK 896  
 Sbjct: 837 GLADVLKSEITTSLSALVDTIVTRFVHORRLFSKQADSVTAAEQLNKDLLASQILDRK 896

Query: 897 SPRTKVADRPQNGPTPATQSAAMFQAPKTPQGMNPAALYNSMTGPFCLPPDXXXXX 956  
 SPRTKVADRPQNGPTPATQSAAMFQAPKTPQGMNPAALYNSMTGPFCLPPD 956  
 Sbjct: 897 SPRTKVADRPQNGPTPATQSAAMFQAPKTPQGMNPAALYNSMTGPFCLPPDQOQQO 956

Query: 957 XXXXXXXXXXXXXXXXXXXXLEQNEALSUVTPKKRHKVTDTRITPRTVSRILAQD 1016  
 LEQNEALSUVTPKKRHKVTDTRITPRTVSRILAQD 1016  
 Sbjct: 957 QTAQQQSAQQQQSSQOQQLEQNEALSUVTPKKRHKVTDTRITPRTVSRILAQD 1016

Query: 1017 XXXXXXXXXXXXXXXXXXXXASNGGNSNATPAQSPTRSSGGAAYHXX 1076  
 ASNGGNSNATPAQSPTRSSGGAAYH 1076  
 Sbjct: 1017 VVPPTGGPSTPQQOQQOQQOQQOQQOQQOQQOQQOQQOQQOQQOQQOQQOQQO 1076

Query: 1077 XXXXXXXXXXXVSLPTSVAIPNPSLHESKVFSPYSPFFNPXXXXXXXXXXXXXXXX 1136  
 VSLPTSVAIPNPSLHESKVFSPYSPFFNP 1136  
 Sbjct: 1077 PPPPPMMPVSLPTSVAIPNPSLHESKVFSPYSPFFNPAAAAQATAAGLHHQHQQHHH 1136

Query: 1137 XXXXXXXXXXXXALMDSRDXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXKTC 1196  
 ALMDSRDX 1196  
 Sbjct: 1137 HQSMQLSSPFGSLGALMDSRDXSPLPPLPPSMLHPALLAAHHGSPDYKTCRAVMDAQ 1196

Query: 1197 DRQSECSADMQFDGMAPTISFYKQMLKTEHOESLMKHCESLTPHSSLTPMHLRKA 1256  
 DRQSECSADMQFDGMAPTISFYKQMLKTEHOESLMKHCESLTPHSSLTPMHLRKA 1256  
 Sbjct: 1197 DRQSECSADMQFDGMAPTISFYKQMLKTEHOESLMKHCESLTPHSSLTPMHLRKA 1256

Query: 1257 KLMFWRVPSSAVLKMYPFDIKFNKNNTAQLVKWFSNREFYIOMEKYARQAVTEG 1316  
 KLMFWRVPSSAVLKMYPFDIKFNKNNTAQLVKWFSNREFYIOMEKYARQAVTEG 1316  
 Sbjct: 1257 KLMFWRVPSSAVLKMYPFDIKFNKNNTAQLVKWFSNREFYIOMEKYARQAVTEG 1316

Query: 1317 TPDDLLIAGDSELYRVLNHYRNNHIEVPQNFVVESTLREFFRAIOGGKDEQSW 1376  
 TPDDLLIAGDSELYRVLNHYRNNHIEVPQNFVVESTLREFFRAIOGGKDEQSW 1376  
 Sbjct: 1317 TPDDLLIAGDSELYRVLNHYRNNHIEVPQNFVVESTLREFFRAIOGGKDEQSW 1376

Query: 1377 SIYKII SRMDDPVPEYFKSPNFLEQL 1403  
 SIYKII SRMDDPVPEYFKSPNFLEQL 1403  
 Sbjct: 1377 SIYKII SRMDDPVPEYFKSPNFLEQL 1403

Annotations:  
 - ≥ 25% for proteins  
 - ≥ 70% for nucleotides  
 - Gap  
 - x Low-Complexity

RID=1106230395-25537-60155114044.BLASTQ4, Protein query - Netscape  
 http://www.ncbi.nlm.nih.gov/BLAST/Blast.cgi#6179901

>gi|6179901|gb|AAF05703.1| homeodomain transcription factor Prospero [Drosophila mel...]  
 Length = 1403

Score = 826 bits (2134), Expect = 0.0  
 Identities = 454/704 (64%), Positives = 461/704 (65%)

Query: 1 MSSXXXXXXXXXXXXLFPQSVSTAXXXXXXXXXXTPAALATHXXXXXXXXXXXX 60  
 MSS LFPQSVSTA TPAALATH 60  
 Sbjct: 1 MSSAAAAAAGAAGGALFQPQSVSTANSSSSNNNSSTPAALATHSPTSNSPVSGASSAS 60

Query: 61 XXXXFXGNLFGGSSAKMLNELFRGMKQADATSGLPQSLDNAMLAAMETATS AELL 120  
 FGNLFGGSS + + QSLDNAMLAAMETATS AELL 120  
 Sbjct: 61 SLLTAAFGNLFGGSSQDAERAVWPPDEAGPGRNEWPAQSLDNAMLAAMETATS AELL 120

Query: 121 GSLNSTSKLLQQQHNNNSIAPANSTPHNSGNTNXXXXXXXXXXXXXXXXXXXXXGSR 180  
 +L ++ P TPMSNGTN KGSRRVSA 180  
 Sbjct: 121 LALQFHVQAAAAAITTALLPPIGTPHNSGNTNASISPGSAHSSSHSQVSPKGSRRVSA 180

Query: 181 CSDRSLEAAAADVAGSPPRAASVSSLNGASSGEQHSQQLQHDLVAHHMLRNILQGR 240  
 CSDRSLEAAAADVAGSPPRAASVSSLNGASSGEQHSQQLQHDLVAHHMLRNILQGR 240  
 Sbjct: 181 CSDRSLEAAAADVAGSPPRAASVSSLNGASSGEQHSQQLQHDLVAHHMLRNILQGR 240

Query: 241 LMQLDQELRTAMXXXXXXXXXXXXHSLKXXXXXXXXXXXXXXXXXXMESINLID 300  
 LMQLDQELRTAM HSKL MESINLID 300  
 Sbjct: 241 LMQLDQELRTAMQQQQQLQEKEQLHSLKLNNNNNNIAATANNNNNTMESINLID 300

Query: 301 ADIKIKSEFPQAPQOXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXHG 360  
 ADIKIKSEFPQAPQO HG 360  
 Sbjct: 301 ADIKIKSEFPQAPQOQSPHGSSHSRSGSGSHSMASDGLRRKSSDSDSHGAQDD 360

Query: 361 XXXXXXXPTGQRESRAPEEPQLPTKESVDDMLDEVELLGLHSRGSMDSLASPS 420  
 PTGQRESRAPEEPQLPTKESVDDMLDEVELLGLHSRGSMDSLASPS 420  
 Sbjct: 361 AQDEEDAAPTGQRESRAPEEPQLPTKESVDDMLDEVELLGLHSRGSMDSLASPS 420

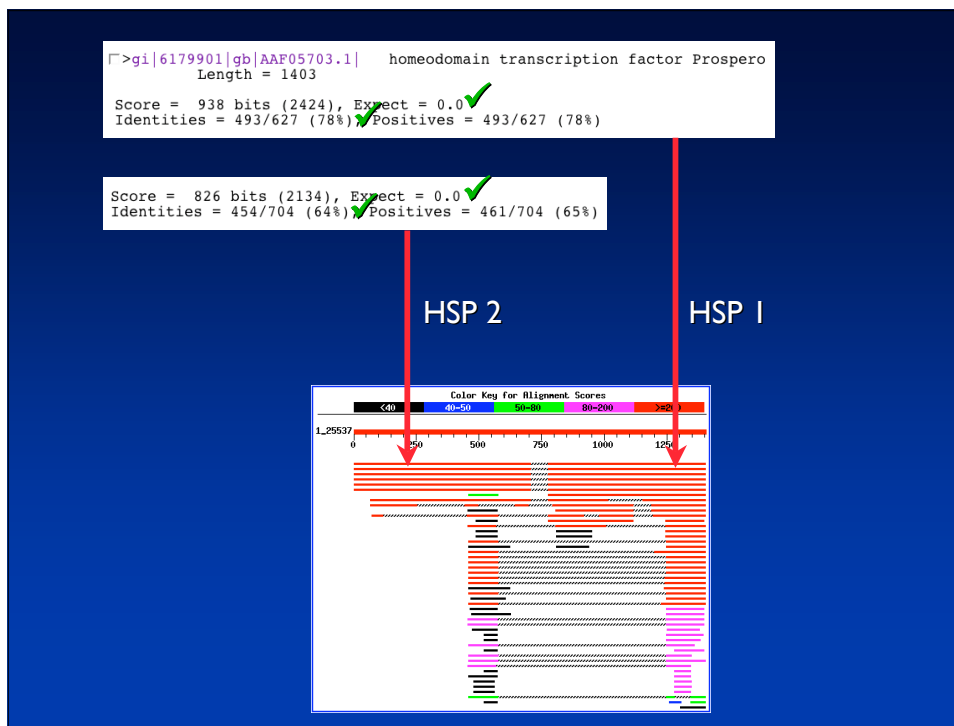
Query: 421 XXXXXXXXXXXXXXXXCVQKTSGSGCLKKPGMDLKRARVENIVSGMRCPSSGLA 480  
 CVQKTSGSGCLKKPGMDLKRARVENIVSGMRCPSSGLA 480  
 Sbjct: 421 MLLLDKDDVLEDDDDDCVQKTSGSGCLKKPGMDLKRARVENIVSGMRCPSSGLA 480

Query: 481 QLVNGCKRRKLYQPOHAMERYVXXXXGLNGLNLSMMLDQEDSENELESFQIQ 540  
 QLVNGCKRRKLYQPOHAMERYV GLNGLNLSMMLDQEDSENELESFQIQ 540  
 Sbjct: 481 QLVNGCKRRKLYQPOHAMERYVAAAAGLNFGLNLSMMLDQEDSENELESFQIQ 540

Query: 541 VEKNALKSQLSMQEQLAEMQOQYVQLCSRMEQESXXXXXXXXXXXXXXXXXNGSD 600  
 VEKNALKSQLSMQEQLAEMQOQYVQLCSRMEQES NGSD 600  
 Sbjct: 541 VEKNALKSQLSMQEQLAEMQOQYVQLCSRMEQESCEQLDQDQVEQEFPDNGSD 600

Query: 601 ELSPSPTLTGDGVSFNHKEETGQERXXXXXXXXXXXXXXXXXXXXSDGANMLSQ 660  
 ELSPSPTLTGDGVSFNHKEETGQER ESSDGANMLSQ 660  
 Sbjct: 601 ELSPSPTLTGDGVSFNHKEETGQERPGSSSPSPPLKPKTSLGESSDGANMLSQ 660

Annotation:  
 - No definition line ∴ second HSP identified



## Suggested BLAST Cutoffs

	<i>E</i> value	Sequence Identity
Nucleotide	$\leq 10^{-6}$	$\geq 70\%$
Protein	$\leq 10^{-3}$	$\geq 25\%$

- Do not use these cutoffs blindly!
- Pay attention to alignments on either side of the dividing line
- Do not ignore biology!



## Database Searching Artifacts

---

- Low-complexity regions
  - Nucleotide searches: removed with DUST (→ N)
  - Protein searches: removed with SEG (→ X)
- Repetitive elements
  - LINE, SINE, Alu
  - Automatic masking “experimental and still under development”
  - RepeatMasker  
*<http://www.repeatmasker.org>*



## Database Searching Artifacts

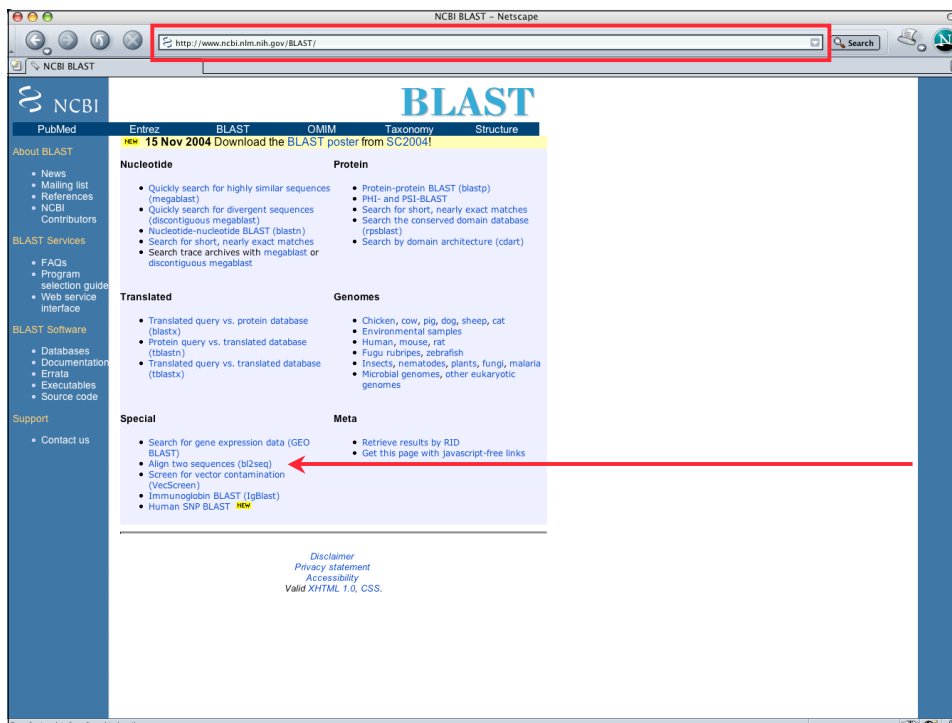
---

- Low-quality sequence hits
  - Expressed sequence tags (ESTs)
  - Single-pass sequence reads from large-scale sequencing (possibly with vector contaminants)

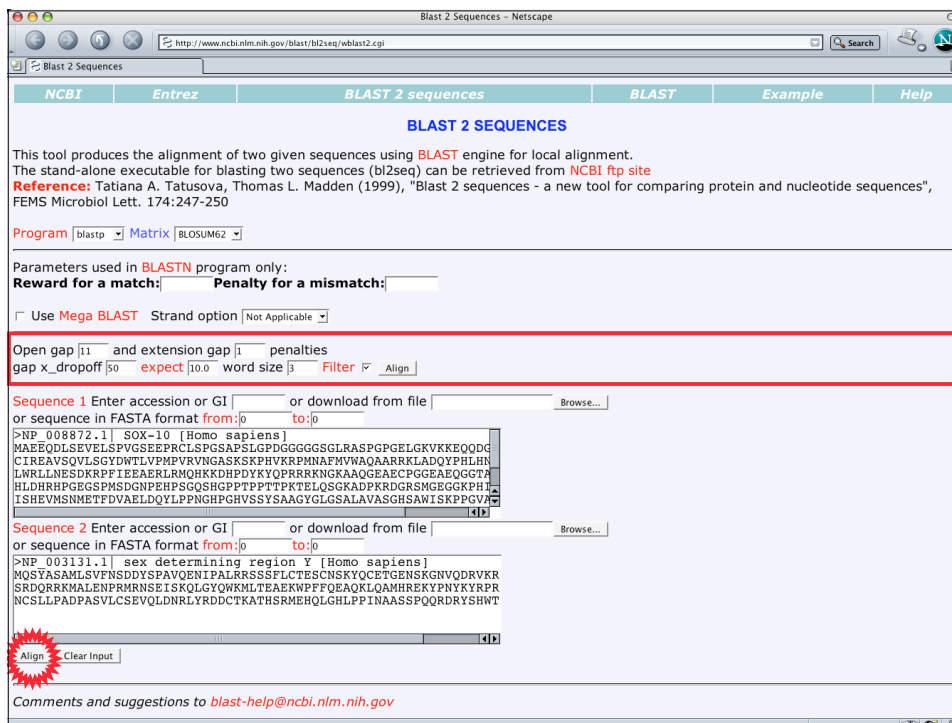
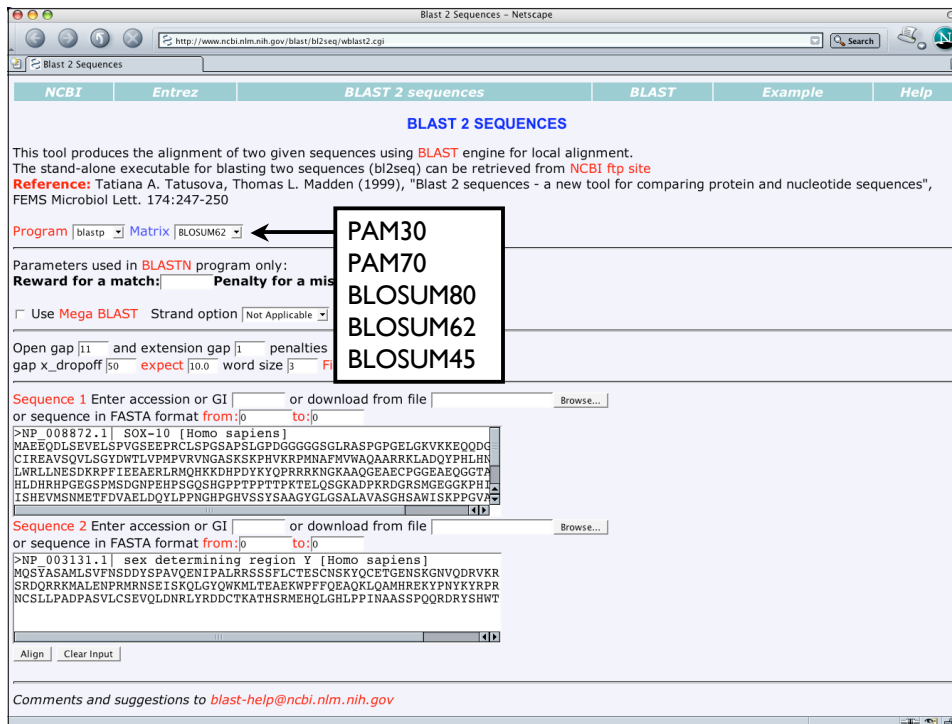


## BLAST 2 Sequences

- Finds local alignments between two protein or nucleotide sequences of interest
  - All BLAST programs available
  - Select BLOSUM and PAM matrices available for protein comparisons
  - Same affine gap costs (adjustable)
  - Input sequences can be masked
- Implementations
  - NCBI Web interface
  - bl2seq downloadable executable  
*ftp://ncbi.nlm.nih.gov/blast/executables/*



NHGRI Current Topics in Genome Analysis 2005  
 Biological Sequence Analysis I





**Blast 2 Sequences results**

BLAST 2 SEQUENCES RESULTS VERSION BLASTP 2.2.10 [Oct-19-2004]

Matrix: BLOSUM62 gap open: 11 gap extension: 1  
 x\_dropoff: 50 expect: 10.000 wordsize: 3 Filter:  Align:

**Sequence 1** |cl|tmpseq\_0 SOX-10 [Homo sapiens] **Length** 466 (1.. 466)  
**Sequence 2** |cl|tmpseq\_1 sex determining region Y [Homo sapiens] **Length** 204 (1.. 204)

**NOTE:** The statistics (bitscore and expect value) is calculated based on the size of nr database

**Score = 94.7 bits (234), Expect = 8e-18**  
**Identities = 39/84 (46%), Positives = 62/84 (73%)**

Query: 95 NGASKSPHVSRPMNFMVWQAARRKLDQYPHLHNAELSKTLGKLRLLNESDKRPF 154  
 N + VKRPMNAF+VW++ RRR+A + P + N+E+SK LG W+L E+K PF  
 Sbjct: 51 NSKGNVQDRVKRPMNAFVWVSRDQRRKMALENPRMRNSEISKOLGYQWKMLTEAEKWPFF 110

Query: 155 EEAERLRMQHKDHPDYKQPRRR 178  
 +EA++L+ H++ +P+YKY+PRR+  
 Sbjct: 111 QEAQKLQAMHREKYPNKYRPRRK 134

CPU time: 0.03 user secs. 0.01 sys. secs 0.04 total secs.

Lambda K H  
 0.311 0.130 0.399

Gapped

## MegaBLAST

- Optimized for aligning very long and/or highly-similar sequences
- Good for batch nucleotide searches
- Search targets include
  - Entire eukaryotic genomes
  - Complete chromosomes and contigs from RefSeq
- Run speeds approximately 10 times faster than BLASTN
  - Adjusted word size
  - Different gap scoring scheme

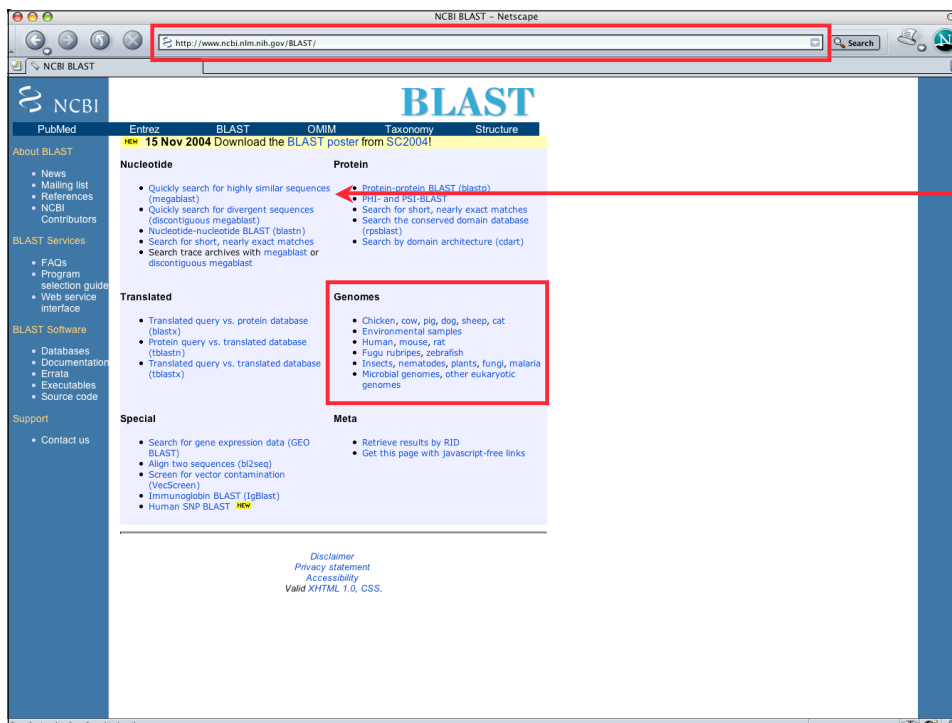


## BLASTN vs. MegaBLAST

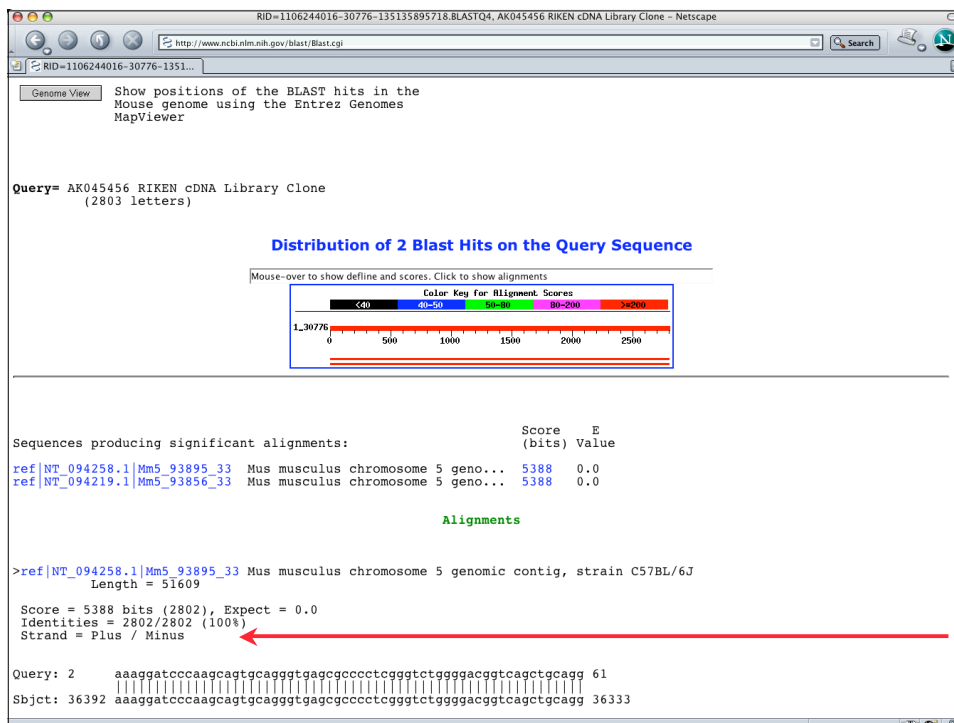
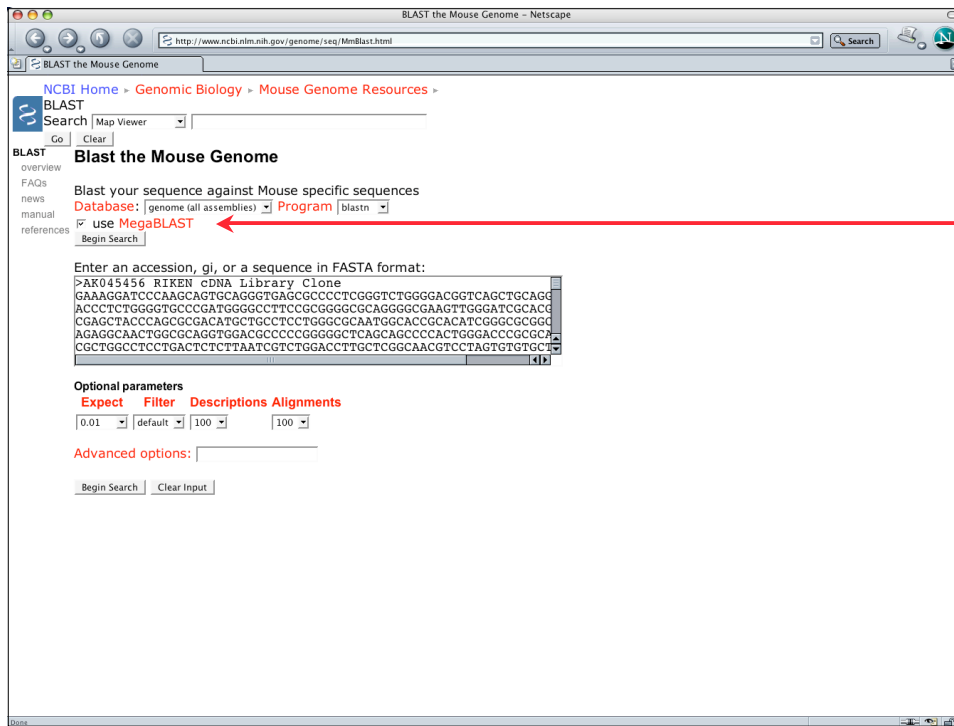
- Word size
  - BLASTN default = 11
  - MegaBLAST default = 28
- *Non-affine* gap penalties

$$\text{Deduction for a gap} = r/2 - q$$

where  $r$  = match reward (default 1)  
 $q$  = mismatch penalty (default -2)  
 and **no penalty for opening the gap**



NHGRI Current Topics in Genome Analysis 2005  
 Biological Sequence Analysis I



NHGRI Current Topics in Genome Analysis 2005  
 Biological Sequence Analysis I

```

RID=1106244016-30776-135135895718.BLASTQ4, AK045456 RIKEN cDNA Library Clone - Netscape
http://www.ncbi.nlm.nih.gov/blast/Blast.cgi
RID=1106244016-30776-1351...
>ref|NT_094219.1|Mm5_93856_33 Mus musculus chromosome 5 genomic contig, strain C57BL/6J
Length = 194175
Score = 5388 bits (2802), Expect = 0.0
Identities = 2802/2802 (100%)
Strand = Plus / Plus
Query: 2 aaaggatcccaagcagtgagggtgagcgcacctcggtctggggacggtcagctgcagg 61
Sbjct: 107035 aaaggatcccaagcagtgagggtgagcgcacctcggtctggggacggtcagctgcagg 107094
Query: 62 gccgggagcacacctctgggggtgcccgatggggcctcccgggggcgcaggggcgaattg 121
Sbjct: 107095 gccgggagcacacctctgggggtgcccgatggggcctcccgggggcgcaggggcgaattg 107154
Query: 122 ggatgcacgcagtgagcccgagctaccacgcgcacatgctgcctctgggccaatgg 181
Sbjct: 107155 ggatgcacgcagtgagcccgagctaccacgcgcacatgctgcctctgggccaatgg 107214
Query: 182 caccgcacatcgggcggctgggggtgcagaggcaactggcgcaggtggacgccccgg 241
Sbjct: 107215 caccgcacatcgggcggctgggggtgcagaggcaactggcgcaggtggacgccccgg 107274
Query: 242 gggctcagcagccccactgggaccgcgcagggtgaccgcctgctcctgactctct 301
Sbjct: 107275 gggctcagcagccccactgggaccgcgcagggtgaccgcctgctcctgactctct 107334
Query: 302 aatcgtctggaccttgcctggcaacgctcctagtgtgtgctctatcgtccgcagccca 361
Sbjct: 107335 aatcgtctggaccttgcctggcaacgctcctagtgtgtgctctatcgtccgcagccca 107394
Query: 362 cctgcgcgcaaatgaccaacatcttcctgctctctggcctctcagacctctctg 421
Sbjct: 107395 cctgcgcgcaaatgaccaacatcttcctgctctctggcctctcagacctctctg 107454
Query: 422 ggcattgctggtcatgctctggaagccgtggctgaggtggcgggtactgccccttgg 481
Sbjct: 107455 ggcattgctggtcatgctctggaagccgtggctgaggtggcgggtactgccccttgg 107514
Query: 482 ggcattctcgcacatctgggtggccttgacatcatgctcaccactgcttcacctgaa 541
    
```

Genome View Show positions of the BLAST hits in the Mouse genome using the Entrez Genomes MapViewer

Query= AK045456 RIKEN cDNA Library Clone (2803 letters)

**Distribution of 2 Blast Hits on the Query Sequence**

Mouse-over to show define and scores. Click to show alignments

Color Key for Alignment Scores

<40	40-50	50-60	60-70	70-80	80-200	>200
-----	-------	-------	-------	-------	--------	------

Sequences producing significant alignments:

ref	Score	E
	(bits)	Value
ref NT_094258.1 Mm5_93895_33	5388	0.0
ref NT_094219.1 Mm5_93856_33	5388	0.0

**Alignments**

```

RID=1106244016-30776-135135895718.BLASTQ4, AK045456 RIKEN cDNA Library Clone - Netscape
http://www.ncbi.nlm.nih.gov/blast/Blast.cgi
RID=1106244016-30776-1351...
>ref|NT_094258.1|Mm5_93895_33 Mus musculus chromosome 5 genomic contig, strain C57BL/6J
Length = 51609
Score = 5388 bits (2802), Expect = 0.0
Identities = 2802/2802 (100%)
Strand = Plus / Minus
Query: 2 aaaggatcccaagcagtgagggtgagcgcacctcggtctggggacggtcagctgcagg 61
Sbjct: 36392 aaaggatcccaagcagtgagggtgagcgcacctcggtctggggacggtcagctgcagg 36333
    
```

The screenshot shows the Entrez Genome view interface. The search query is "AK045456 RIKEN cDNA Library Clone". The results table shows two hits, both "not mapped" to any specific region on chromosome 5. A callout box points to these two entries with the text: "Overlapping clones? Two separate regions of chromosome 5? Finished sequence needed Check subsequent builds of mouse genome".

Chr	Map element	Type	BLAST results
			Hits Score E value
not mapped	NT_094258	CONTIG 1	5388 0.0
not mapped	NT_094219	CONTIG 1	5388 0.0

## Overview

- Week 4: Comparative methods and concepts
  - Similarity vs. Homology
  - Global vs. Local Alignments
  - Scoring Matrices
  - BLAST
  - BLAT
- Week 5: Predictive methods and concepts
  - Profiles, patterns, motifs, and domains
  - Secondary structure prediction
  - Structures: VAST, Cn3D, and *de novo* prediction



## BLAT

---

- “BLAST-Like Alignment Tool”
- Designed to rapidly-align longer nucleotide sequences ( $L \geq 40$ ) having > 95% sequence similarity
- Can find exact matches reliably down to  $L = 33$
- Method of choice when looking for exact matches in nucleotide databases
- 500 times faster for mRNA/DNA searches
- May miss divergent or shorter sequence alignments
- Can be used on protein sequences



## When to Use BLAT

---

- To characterize an unknown gene or sequence fragment
  - Find its genomic coordinates
  - Determine gene structure (the presence and position of exons)
  - Identify markers of interest in the vicinity of a sequence
- To find highly-similar sequences
  - Identify gene family members
  - Identify putative homologs
- To display a specific sequence as a separate track





Rat BLAT Results - Netscape

http://genome.ucsc.edu/cgi-bin/hgBlat

Rat BLAT Results

Home Genomes Gene Sorter Blat Tables FAQ Help

### Rat BLAT Results

#### BLAT Search Results

ACTIONS	QUERY	SCORE	START	END	QSIZE	IDENTITY	CHRO	STRAND	START	END	SPAN
<a href="#">browser details</a>	CB312815	710	1	733	768	98.1%	5	+	101460825	101461549	725

Two red arrows point to the 'browser details' link in the first row of the table.

Rat chr5:101,460,643-101,461,730 - UCSC Genome Browser v95 - Netscape

http://genome.ucsc.edu/cgi-bin/hgTracks

Rat chr5:101,460,643-101,461,730

Home Genomes BLAT DNA Tables Gene Sorter Convert Ensembl PDF/PS Help

### UCSC Genome Browser on Rat Jun 2003 Assembly

move <<< << < > >> >>> zoom in 1.5x 3x 10x base zoom out 1.5x 3x 10x  
 position chr5:101,460,643-101,461,730 size 1,088 bp. image width: 1250 jump

Base Position (bp) 101460750 101460800 101460850 101460900 101460950 101461000 101461050 101461100 101461150 101461200 101461250 101461300 101461350 101461400 101461450 101461500 101461550 101461600 101461650 101461700

BLAT

RefSeq Genes

Ensembl Genes

UCSC ESTs

Non-RAT mRNAs

House Chai in

House Het

Neocathoxen

Gap Locations

Your Sequence from BLAT Search

Known Genes Based on S115-PROT, TRINL, MESH, and RefSeq

Rat Gene

Rat Gene Database Curated Genes

Ensembl Gene Predictions

Rat mRNAs from GenBank

Rat ESTs That Have Been Sliced

Non-RAT mRNAs from GenBank

House mRNAs from GenBank

House (Only Genes with At Least One)

Repeat Map Elements by RepeatMasker

move start < 2.0 > move end < 2.0 >

Click on a feature for details. Click on base position to zoom in around cursor. Click on left mini-buttons for track-specific options.

reset all hide all Chromosome  Guidelines  Labels: left  center  refresh

**Chromosome Color Key:**

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 X Y M Un

Use drop down controls below and press refresh to alter tracks displayed. Tracks with lots of items will automatically be displayed in more compact modes.

#### Mapping and Sequencing Tracks

Base Position	Chromosome Band	RGD QTL	STS Markers	Recomb Rate
full	full	hide	hide	hide
Assembly	Gap	Bactigs	BAC End Pairs	GC Percent
full	full	full	hide	hide
Short Match	BLAT Sequence			
hide	pack			

#### Genes and Gene Prediction Tracks

Known Genes	RefSeq Genes	RGD Genes	MGC Genes	Ensembl Genes
pack	pack	dense	hide	dense





## FASTA

---

- Identifies regions of local alignment
- Employs an approximation of the Smith-Waterman algorithm to determine the best alignment between two sequences
- Method is significantly different from that used by BLAST
- Online implementations at  
<http://fasta.bioch.virginia.edu>  
<http://www.ebi.ac.uk/fasta33>



## Overview

---

- Week 4: Comparative methods and concepts
  - Similarity *vs.* Homology
  - Global *vs.* Local Alignments
  - Scoring Matrices
  - BLAST
  - BLAT
- Week 5: Predictive methods and concepts
  - Profiles, patterns, motifs, and domains
  - Secondary structure prediction
  - Structures: VAST, Cn3D, and *de novo* prediction



## Further Reading

Altschul, S.F., Boguski, M.S., Gish, W., and Wootton, J.C. 1994. Issues in searching molecular sequence databases. *Nat. Genet.* 6: 119-129. A review of the issues that are of importance in using sequence similarity search programs, including potential pitfalls.

Baxevanis, A.D. Assessing pairwise sequence similarity: BLAST and FASTA. In *Bioinformatics: A Practical Guide to the Analysis of Genes and Proteins*, third edition (Baxevanis, A.D. and Ouellette, B.F.F., eds.), John Wiley and Sons, 2005. An overview of the methods used to generate pairwise sequence alignments and assess the biological significance of results.

Henikoff, S. and Henikoff, J.G. 2000. Amino acid substitution matrices. *Adv. Protein Chem.* 54: 73-97. A comprehensive review covering the factors critical to the construction of protein scoring matrices.

Korf, I., Yandell, M., and Bedell, J. BLAST. O'Reilly and Associates, 2003. An in-depth treatment of the BLAST algorithm, its applications, as well as installation, hardware, and software considerations. The book provides "documentation" that is not easily found elsewhere.

Pearson, W.R. Finding protein and nucleotide similarities with FASTA. 2003. *Current Protocols in Bioinformatics* 3.9.1-3.9.23. An in-depth discussion of the FASTA algorithm, including worked examples and additional information regarding run options and use scenarios.

Wheeler, D.G. Selecting the right protein scoring matrix. 2003. *Current Protocols in Bioinformatics* 3.5.1-3.5.6. A discussion of PAM, BLOSUM, and specialized scoring matrices, with guidance regarding the proper choice of matrices for particular types of protein-based analyses.

## References

Altschul, S.F., Gish, W., Miller, W., Myers, E.W., and Lipman, D.J. 1991. 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.

Dayhoff, M.O., Schwartz, R.M., Orcutt, B.C. 1978. A model of evolutionary change in proteins. *In Atlas of Protein Sequence and Structure*, M.O. Dayhoff, ed., National Biomedical Research Foundation, Washington, 5: 345-352.

Henikoff, S. and Henikoff, J.G. 1991. Automated assembly of protein blocks for database searching. *Nucleic Acids Res.* 19: 6565-6572.

Henikoff, S. and Henikoff, J.G. 1992. Amino acid substitution matrices from protein blocks. *Proc. Natl. Acad. Sci. USA* 89: 10915-10919.

Henikoff, S. and Henikoff, J.G. 1993. Performance evaluation of amino acid substitution matrices. *Proteins Struct. Funct. Genet.* 17: 49-61.

Henikoff, S. and Henikoff, J.G. 2000. Amino acid substitution matrices. *Adv. Protein Chem.* 54: 73-97.

Karlin, S. and Altschul, S.F. 1990. Methods for assessing the statistical significance of molecular sequence features by using general scoring schemes. *Proc. Natl. Acad. Sci. USA* 87: 2264-2268.

Kent, W.J. 2002. BLAT: the BLAST-like alignment tool. *Genome Res.* 12: 656-664.

Pearson, W.R. 1995. Comparison of methods for searching protein sequence databases. *Protein Sci.* 4: 1145-1160.

Pearson, W.R. 2000. Flexible sequence similarity searching with the FASTA3 program package. *Methods Mol. Biol.* 132: 185-219.

Pearson, W.R. Finding protein and nucleotide similarities with FASTA. 2003. *Current Protocols in Bioinformatics* 3.9.1-3.9.23.

Pearson, W.R. and Lipman, D.J. 1988. Improved tools for biological sequence comparison. *Proc. Natl. Acad. Sci. USA* 85: 2444-2448.

Smith, T.F. and Waterman, M.S. 1981. Identification of common molecular subsequences. *J. Mol. Biol.* 147: 195-197.

Tatusova, T.A. and Madden, T.L. 1999. BLAST2Sequences, a new tool for comparing protein and nucleotide sequences. *FEMS Microbio. Lett.* 174: 247-250.

Wootton, J.C. and Federhen, S. 1993. Statistics of local complexity in amino acid sequences and sequence databases. *Comput. Chem.* 17: 149-163.