

Current Topics in Genome Analysis 2012

Week 2: Biological Sequence Analysis I

Andy Baxevanis, Ph.D.



NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research



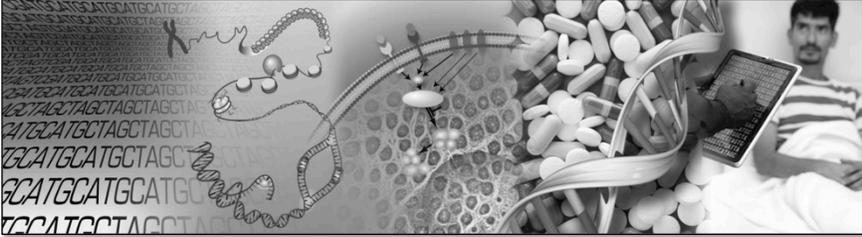
Current Topics in Genome Analysis 2012

Andy Baxevanis, Ph.D.

***No Relevant Financial Relationships with
Commercial Interests***



NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

Understanding the structure of genomes	Understanding the biology of genomes	Understanding the biology of disease	Advancing the science of medicine	Improving the effectiveness of healthcare
				
				
Bioinformatics and Computational Biology				

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

Overview

- **Week 2**
 - Similarity vs. Homology
 - Global vs. Local Alignments
 - Scoring Matrices
 - BLAST
 - BLAT
- **Week 4**
 - Profiles, Patterns, Motifs, and Domains
 - Structures: VAST, Cn3D, and *de novo* Prediction
 - Multiple Sequence Alignment

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

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, or 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

It is worth repeating here that homology, like pregnancy, is indivisible⁸. You either are homologous (pregnant) or you are not. Thus, if what one means to assert is that 80% of the character states are identical one should speak of 80% identity, and not 80% homology.

Fitch, Trends Genet. 16: 227-231, 2000



Defining the Terms

- 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*)

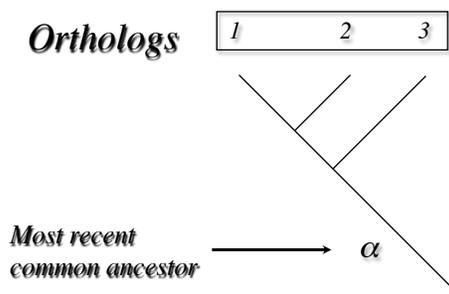


Defining the Terms

- **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



Defining the Terms

Orthologs **Paralogs**

Most recent common ancestor → α β

Gene duplication →

- 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)

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

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

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

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 alignment could be generated for any two sequences being compared
- Best for sequences that share some similarity, or for sequences of different lengths



Scoring Matrices

- Empirical weighting scheme representing physicochemical and biological characteristics of nucleotides and amino acids
 - Side chain structure and chemistry
 - Side chain function
- Amino acid-based examples:
 - 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

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



Matrix Structure: Nucleotides

- Simple match/mismatch scoring scheme:

Match +2
 Mismatch -3

	A	T	G	C
A	2	-3	-3	-3
T	-3	2	-3	-3
G	-3	-3	2	-3
C	-3	-3	-3	2

- 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	-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	2	-1	-1	-3	-4	-1	-3	-3	-1	0	-1	4	-3	-3	4	1	-1	-4	
C	0	-3	-3	-3	3	-3	-4	-3	-3	-1	-1	-3	-1	-2	-3	-1	-1	2	-2	-1	-3	-3	-2	-4
Q	-1	0	0	2	-4	2	5	-2	0	-3	-3	1	-2	-3	-1	0	-1	3	-2	-2	1	4	-1	-4
E	0	-2	0	-1	-3	-2	2	6	-2	-4	-4	-2	-3	-3	-2	0	-2	2	-3	-3	-1	-2	-1	-4
G	-2	0	1	-1	-3	0	0	-2	8	-3	-3	-1	-2	-1	-2	-1	-2	2	-3	0	0	-1	-4	
H	-1	-3	-3	-3	-1	-3	-3	-4	-3	4	2	-3	1	0	-3	-2	-1	3	-1	3	-3	-3	-1	-4
I	-1	-2	-3	-4	-1	-2	-3	-4	-3	2	4	-2	2	0	-3	-2	-1	2	-1	1	-4	-3	-1	-4
L	-1	2	0	-1	-3	1	1	-2	-1	-3	-2	5	-1	-3	-1	0	-1	3	-2	-2	0	1	-1	-4
K	-1	-1	-2	-3	-1	0	-2	-3	-2	1	2	-1	5	0	-2	-1	-1	1	-1	1	-3	-1	-1	-4
M	-2	-3	-3	-3	-2	-3	-3	-3	-1	0	0	-3	0	6	-4	-2	-2	1	3	-1	-3	-3	-1	-4
F	-1	-2	-2	-1	-3	-1	-1	-2	-2	-3	-3	-1	-2	-4	7	-1	-1	4	-3	-2	-2	-1	-2	-4
P	-1	1	0	-1	0	0	0	-1	-2	-2	0	-1	-2	-1	4	1	3	-2	-2	0	0	0	-4	
S	0	-1	0	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	1	5	2	-2	0	-1	-1	-1	0	-4
T	0	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
W	0	-3	-3	-3	-1	-2	-2	-3	-3	3	1	-2	1	-1	-2	-2	0	-3	-1	4	-3	-2	-1	-4
Y	0	-3	-3	-3	-1	-2	-2	-3	-3	3	1	-2	1	-1	-2	-2	0	-3	-1	4	-3	-2	-1	-4
V	-2	-1	3	4	-3	0	1	-1	0	-3	-4	0	-3	-3	-2	0	-1	-4	-3	-3	4	1	-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	0	-1	0	1	-3	3	4	-2	0	-3	-3	1	-1	-3	-1	0	-1	-3	-2	-2	1	4	-1	-4
X	-1	-1	-1	-2	-1	-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

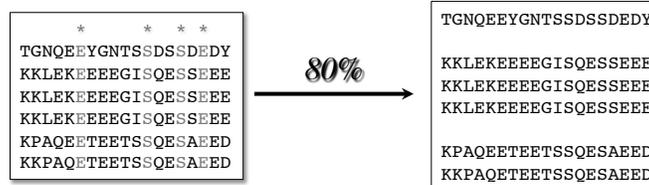
BLOSUM Matrices

- Henikoff and Henikoff, 1992
- **B**locks **S**ubstitution **M**atrix
 - 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**



Which one to choose?

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



So many matrices...

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

Further Reading

Unit 3.5 Current Protocols in Bioinformatics

- **PAM Matrices**
- **BLOSUM Matrices**
- **Specialized Scoring Matrices**

Selecting the Right Protein-Scoring Matrix

UNIT 3.5

OVERVIEW
Every program for searching protein sequences against a database includes a choice of "protein-scoring matrices," also called "weight matrices." Weight matrices add sensitivity to the search, while statistical significance adds selectivity (see *see it*). Virtually every user chooses the default, typically PAM250 or BLOSUM62. Despite the fact that the choice of matrix can strongly influence the outcome of the analysis, most users do not know why a particular matrix should be used. In general, scoring matrices implicitly represent a particular theory of protein sequence evolution. This unit provides guidance in the choice of a scoring matrix, in understanding the assumptions underlying the PAM and BLOSUM scoring matrices and in making the proper choice. The selection of PAM matrices is covered first, after which the selection of BLOSUM matrices is discussed, and finally a brief overview of the wide variety of specialized scoring matrices is provided.

PAM MATRICES
PAM, a two-letter matrix derived from Accepted Point Mutation (Dayhoff, 1978) is a probabilistic model for amino acid replacement defined by comparing the frequencies of replacements to closely related sequences to the frequency expected from the completely random replacement of amino acids. The basis of this scoring system is the observation that the evolution of protein sequences is a nonrandom process—i.e., some amino acid replacements occur much more frequently than others, especially in related sequences. Amino acid substitution tend to conserve charge, size, and hydrophobicity among other characteristics. One would expect that the substitution of glycine for alanine (CH₃ versus H) would have less of an effect on a protein's structure and function than the substitution of alanine for leucine (CH₃ versus substituted hydrocarbon). The ratio of the number of observed substitutions to the number of possible substitutions of those characteristics, the sequences are related. In a related discussion of the derivation and use of the PAM matrices is given in George et al. (1995).

PAM matrices are the result of comparing the probability of one substitution per 100 amino acids, called the PAM 1 matrix. Higher PAM matrices are obtained by multiplying the PAM 1 matrix by itself a defined number of times. Thus, a PAM 250 matrix is the result of performing 100 matrix multiplications of the PAM 1 matrix against itself. Similarly, the PAM 250 matrix is derived by multiplying the PAM 1 matrix against itself 250 times.

Biologically, the PAM50 matrix means that a 50 amino acids have been from 50 mutations, while the PAM 250 matrix means there have been 2.5 amino acid replacements at each site (see *see it*) regarding insertions and deletions. This is not unusual, but illustrates that over evolutionary time, it is possible that an amino acid changed to a glycine, then to a valine, and then back to an alanine. These direct substitutions are defined from observed amino acid frequency data in protein families and superfamilies.

Choosing a PAM Matrix
It is extremely important to note that PAM matrices are derived from protein sequence data available in the late 1960s and early 1970s. Most protein known at that time were small, globular, and hydrophilic. If the researcher believes their protein contains substantial hydrophobic regions, such as membrane-spanning helices or sheets, the PAM matrix is less useful than others described in this unit. Dayhoff et al. (1978) were the first to define the terms specific family and superfamily. A protein family is defined as sequences related from 30% identical or greater to each other. A protein superfamily is defined as sequences related from 20% identical or greater to each other. A protein superfamily may contain many protein families. The user should be aware that while the terms "family" and "superfamily" are widely used in biology, most of the time the original definition of Dayhoff and colleagues is not being used (see below).

Selecting all potential candidates: PAM 250
The most widely used PAM matrix is PAM 250 (Fig. 3.5.1). It has been chosen because it is capable of accurately detecting similarities in the 30% range (i.e., superfamilies), that is, when the two proteins are 30% different from each other (George et al., 1995). Another way to think about this is that the PAM 250

FIGURE 3.5.1
FIGURE 3.5.2
FIGURE 3.5.3
FIGURE 3.5.4
FIGURE 3.5.5
FIGURE 3.5.6
FIGURE 3.5.7
FIGURE 3.5.8
FIGURE 3.5.9
FIGURE 3.5.10
FIGURE 3.5.11
FIGURE 3.5.12
FIGURE 3.5.13
FIGURE 3.5.14
FIGURE 3.5.15
FIGURE 3.5.16
FIGURE 3.5.17
FIGURE 3.5.18
FIGURE 3.5.19
FIGURE 3.5.20
FIGURE 3.5.21
FIGURE 3.5.22
FIGURE 3.5.23
FIGURE 3.5.24
FIGURE 3.5.25
FIGURE 3.5.26
FIGURE 3.5.27
FIGURE 3.5.28
FIGURE 3.5.29
FIGURE 3.5.30
FIGURE 3.5.31
FIGURE 3.5.32
FIGURE 3.5.33
FIGURE 3.5.34
FIGURE 3.5.35
FIGURE 3.5.36
FIGURE 3.5.37
FIGURE 3.5.38
FIGURE 3.5.39
FIGURE 3.5.40
FIGURE 3.5.41
FIGURE 3.5.42
FIGURE 3.5.43
FIGURE 3.5.44
FIGURE 3.5.45
FIGURE 3.5.46
FIGURE 3.5.47
FIGURE 3.5.48
FIGURE 3.5.49
FIGURE 3.5.50
FIGURE 3.5.51
FIGURE 3.5.52
FIGURE 3.5.53
FIGURE 3.5.54
FIGURE 3.5.55
FIGURE 3.5.56
FIGURE 3.5.57
FIGURE 3.5.58
FIGURE 3.5.59
FIGURE 3.5.60
FIGURE 3.5.61
FIGURE 3.5.62
FIGURE 3.5.63
FIGURE 3.5.64
FIGURE 3.5.65
FIGURE 3.5.66
FIGURE 3.5.67
FIGURE 3.5.68
FIGURE 3.5.69
FIGURE 3.5.70
FIGURE 3.5.71
FIGURE 3.5.72
FIGURE 3.5.73
FIGURE 3.5.74
FIGURE 3.5.75
FIGURE 3.5.76
FIGURE 3.5.77
FIGURE 3.5.78
FIGURE 3.5.79
FIGURE 3.5.80
FIGURE 3.5.81
FIGURE 3.5.82
FIGURE 3.5.83
FIGURE 3.5.84
FIGURE 3.5.85
FIGURE 3.5.86
FIGURE 3.5.87
FIGURE 3.5.88
FIGURE 3.5.89
FIGURE 3.5.90
FIGURE 3.5.91
FIGURE 3.5.92
FIGURE 3.5.93
FIGURE 3.5.94
FIGURE 3.5.95
FIGURE 3.5.96
FIGURE 3.5.97
FIGURE 3.5.98
FIGURE 3.5.99
FIGURE 3.5.100

Contributed by David Wheeler
Current Protocols in Bioinformatics (2011) 3.5.1-3.5.6
Copyright © 2011 by John Wiley & Sons, Inc.

FIGURE 3.5.1
FIGURE 3.5.2
FIGURE 3.5.3
FIGURE 3.5.4
FIGURE 3.5.5
FIGURE 3.5.6
FIGURE 3.5.7
FIGURE 3.5.8
FIGURE 3.5.9
FIGURE 3.5.10
FIGURE 3.5.11
FIGURE 3.5.12
FIGURE 3.5.13
FIGURE 3.5.14
FIGURE 3.5.15
FIGURE 3.5.16
FIGURE 3.5.17
FIGURE 3.5.18
FIGURE 3.5.19
FIGURE 3.5.20
FIGURE 3.5.21
FIGURE 3.5.22
FIGURE 3.5.23
FIGURE 3.5.24
FIGURE 3.5.25
FIGURE 3.5.26
FIGURE 3.5.27
FIGURE 3.5.28
FIGURE 3.5.29
FIGURE 3.5.30
FIGURE 3.5.31
FIGURE 3.5.32
FIGURE 3.5.33
FIGURE 3.5.34
FIGURE 3.5.35
FIGURE 3.5.36
FIGURE 3.5.37
FIGURE 3.5.38
FIGURE 3.5.39
FIGURE 3.5.40
FIGURE 3.5.41
FIGURE 3.5.42
FIGURE 3.5.43
FIGURE 3.5.44
FIGURE 3.5.45
FIGURE 3.5.46
FIGURE 3.5.47
FIGURE 3.5.48
FIGURE 3.5.49
FIGURE 3.5.50
FIGURE 3.5.51
FIGURE 3.5.52
FIGURE 3.5.53
FIGURE 3.5.54
FIGURE 3.5.55
FIGURE 3.5.56
FIGURE 3.5.57
FIGURE 3.5.58
FIGURE 3.5.59
FIGURE 3.5.60
FIGURE 3.5.61
FIGURE 3.5.62
FIGURE 3.5.63
FIGURE 3.5.64
FIGURE 3.5.65
FIGURE 3.5.66
FIGURE 3.5.67
FIGURE 3.5.68
FIGURE 3.5.69
FIGURE 3.5.70
FIGURE 3.5.71
FIGURE 3.5.72
FIGURE 3.5.73
FIGURE 3.5.74
FIGURE 3.5.75
FIGURE 3.5.76
FIGURE 3.5.77
FIGURE 3.5.78
FIGURE 3.5.79
FIGURE 3.5.80
FIGURE 3.5.81
FIGURE 3.5.82
FIGURE 3.5.83
FIGURE 3.5.84
FIGURE 3.5.85
FIGURE 3.5.86
FIGURE 3.5.87
FIGURE 3.5.88
FIGURE 3.5.89
FIGURE 3.5.90
FIGURE 3.5.91
FIGURE 3.5.92
FIGURE 3.5.93
FIGURE 3.5.94
FIGURE 3.5.95
FIGURE 3.5.96
FIGURE 3.5.97
FIGURE 3.5.98
FIGURE 3.5.99
FIGURE 3.5.100

Gaps

- Used to improve alignments between two sequences
- Compensate for insertions and deletions →
gaps represent biological events
- 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

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

		nucleotide	protein
where	$G =$ gap-opening penalty	5	11
	$L =$ gap-extension penalty	2	1
	$n =$ length of the gap		
and	$G > L$		



BLAST

- **Basic Local Alignment Search Tool**
- **Seeks high-scoring segment pairs (HSP)**
 - Pair of sequences that can be aligned with one another
 - 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>
BLASTN	Nucleotide	Nucleotide
BLASTP	Protein	Protein
BLASTX	Nucleotide, six-frame translation	Protein
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 PEG 15 PRG 14 PKG 14 PNG 13 PDG 13 PHG 13 PMG 13 PSG 13 PQA 12 PQN 12 etc.	= 7 + 5 + 6 Neighborhood Score Threshold (T = 13)
---------------------------	--	---

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
 Division of Intramural Research

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.	

↓

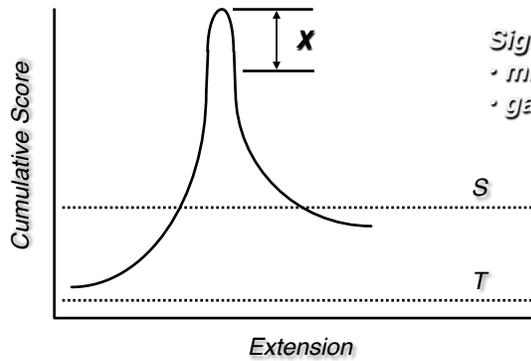
	←	□	→	
Query:	325	SLAALLNKCKTPQGQRLVNQWIKQPLMDKNRIEERLNLVEA	365	
		+LA++L TP+G R++ +W+ +P+ D + ER + A		
Sbjct:	290	TLASVLDCTVTPMGSRLKRWLHMPVRDTRVLLERQQTIGA	330	

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
 Division of Intramural Research

Extension

←—————|—————→

Query:	325	SLAALLNKCKTPQGQRLVNQWIKOPLMDKNRIEERLNLVEA	365
		+LA++L TP+G R++ +W+ +P+ D + ER + A	
Sbjct:	290	TLASVLDCTVTPMGSRMLKRWLHMPVRDTRVLLERQQTIGA	330



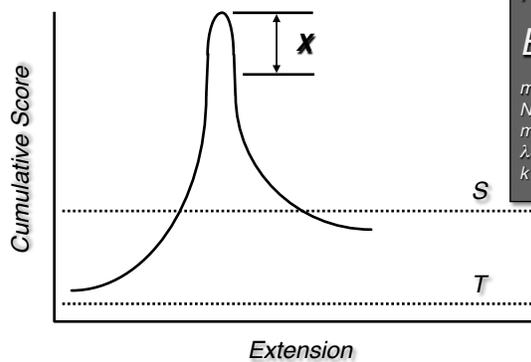
Significance decay
 · mismatches
 · gap penalties



Scores and Probabilities

←—————|—————→

Query:	325	SLAALLNKCKTPQGQRLVNQWIKOPLMDKNRIEERLNLVEA	365
		+LA++L TP+G R++ +W+ +P+ D + ER + A	
Sbjct:	290	TLASVLDCTVTPMGSRMLKRWLHMPVRDTRVLLERQQTIGA	330



Karlin-Altschul Equation

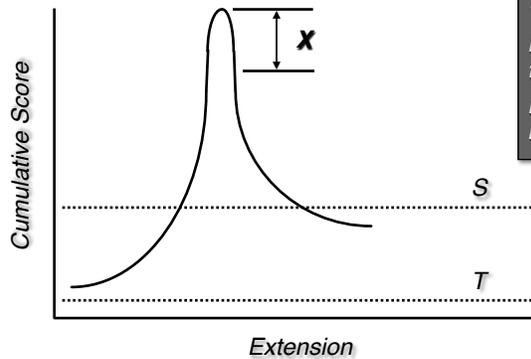
$$E = kmNe^{-\lambda S}$$

m # letters in query
N # letters in database
mN size of search space
 λS normalized score
k minor constant



Scores and Probabilities

Query:	325	SLAALLNKCKTPQGQRLVNQWIKQPLMDKNRIEERLNLVEA	365
		+LA++L TP+G R++ +W+ +P+ D + ER + A	
Sbjct:	290	TLASVLDCTVTPMGSRMLKRWLHMPVRDTRVLLERQQTIGA	330



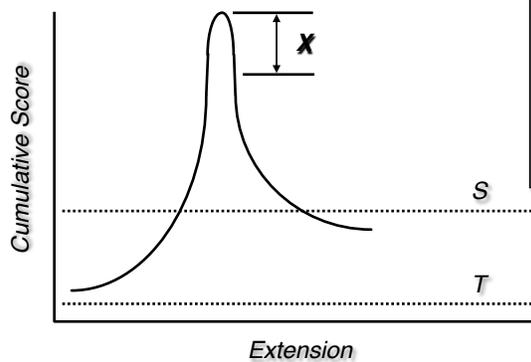
$$E = kmNe^{-\lambda S}$$

Number of HSPs
 found purely by chance
 Lower values signify
 higher similarity



Scores and Probabilities

Query:	325	SLAALLNKCKTPQGQRLVNQWIKQPLMDKNRIEERLNLVEA	365
		+LA++L TP+G R++ +W+ +P+ D + ER + A	
Sbjct:	290	TLASVLDCTVTPMGSRMLKRWLHMPVRDTRVLLERQQTIGA	330



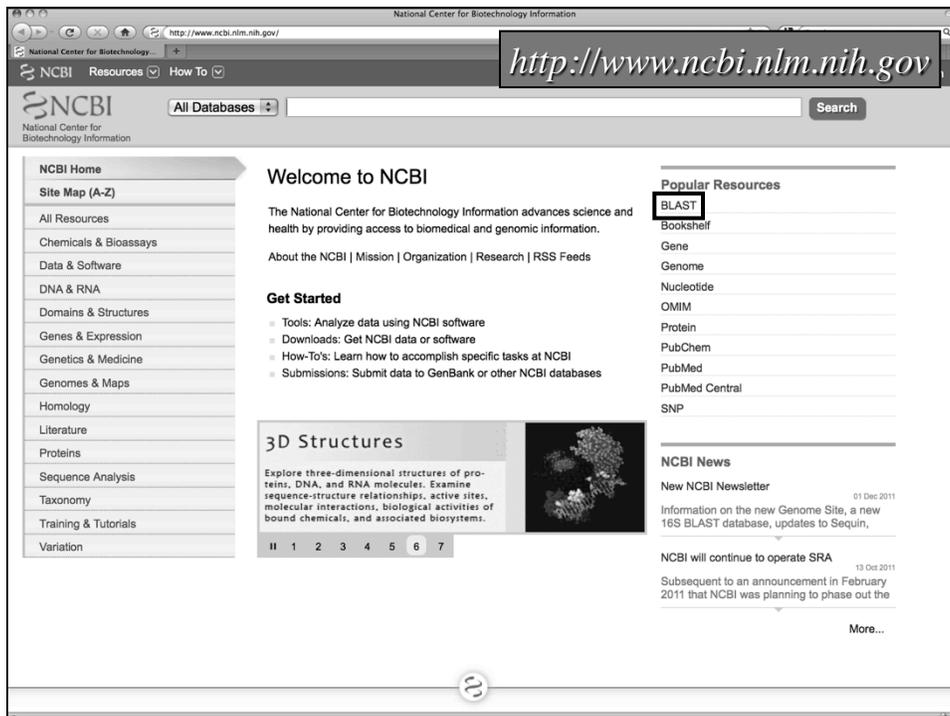
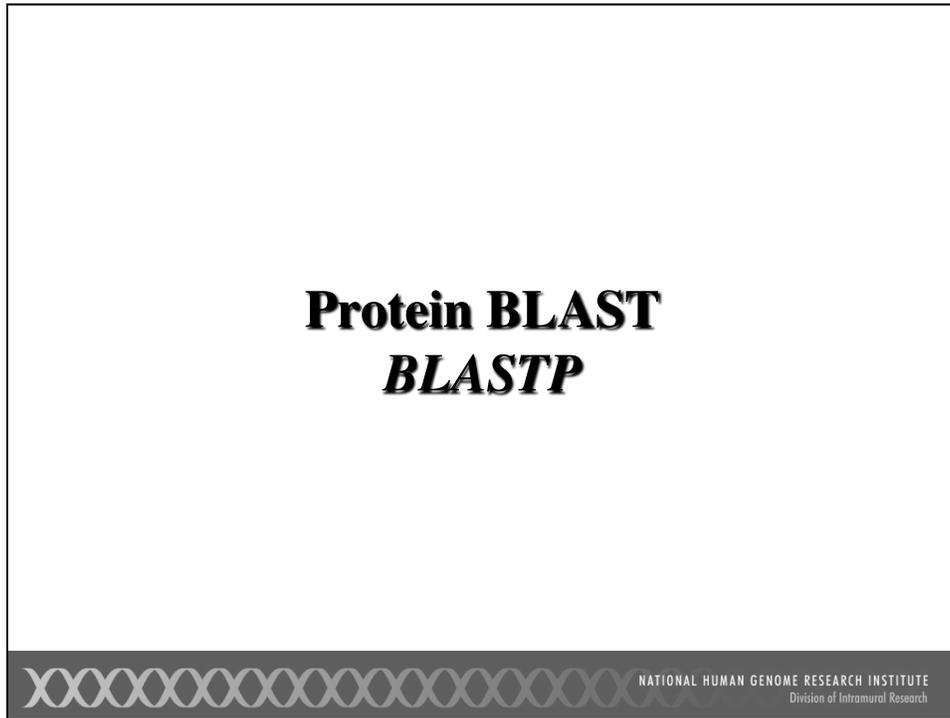
$$E \leq 10^{-6}$$

for nucleotides

$$E \leq 10^{-3}$$

for proteins





Available protein databases include:

<i>nr</i>	<i>Non-redundant</i>
<i>refseq</i>	<i>Reference Sequences</i>
<i>swissprot</i>	<i>SWISS-PROT</i>
<i>pat</i>	<i>Patents</i>
<i>pdb</i>	<i>Protein Data Bank</i>
<i>env_nr</i>	<i>Environmental samples</i>

RefSeq

- **Goal:** Provide a single reference sequence for each molecule of the central dogma (DNA, mRNA, protein)
- **Distinguishing Features**
 - Non-redundancy
 - Updates to reflect the current knowledge of sequence data and biology
 - Ongoing curation by NCBI staff and collaborators, with review status indicated on each record

RefSeq Accession Format

From curation of GenBank entries:

NT_123456	Genomic contigs
NM_123456	mRNAs
NP_123456	Proteins
NR_123456	Non-coding transcripts

From genome annotation:

XM_123456	Model mRNA
XP_123456	Model proteins

Complete key at

<http://www.ncbi.nlm.nih.gov/RefSeq/key.html>



Protein BLAST: search protein databases using a protein query

BLAST[®] Basic Local Alignment Search Tool

Standard Protein BLAST

Enter Query Sequence

Enter accession number(s), gi(s), or FASTA sequence(s)

Query sequence: MSAA... (truncated)

Job Title: []

Choose Search Set

Database: Non-redundant protein sequences (nr)

Organism: []

Exclude: Models (XM/XP) Uncultured/environmental sample sequences

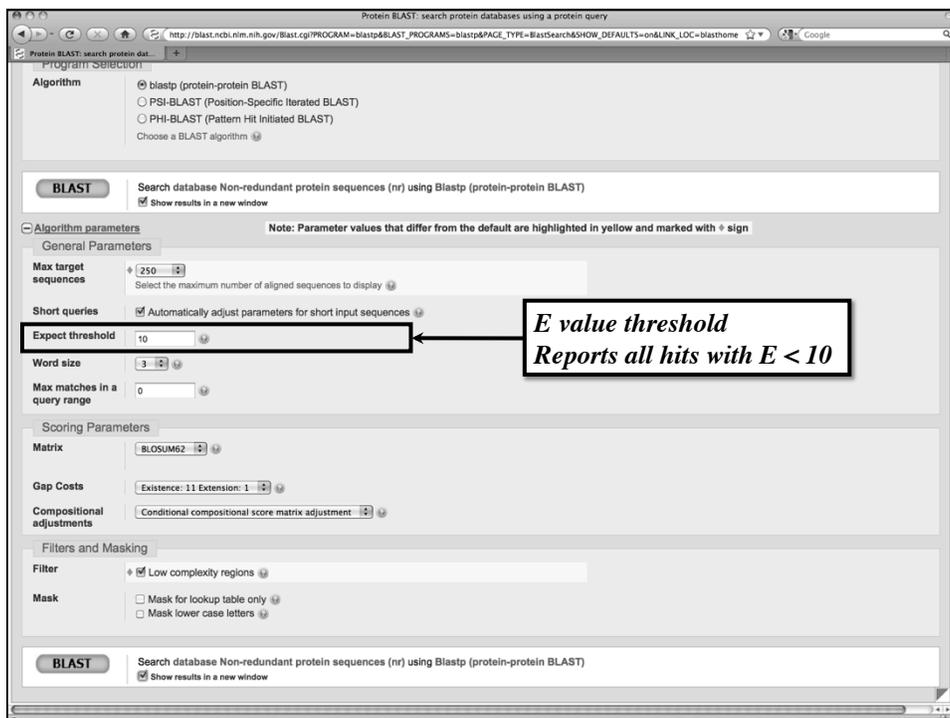
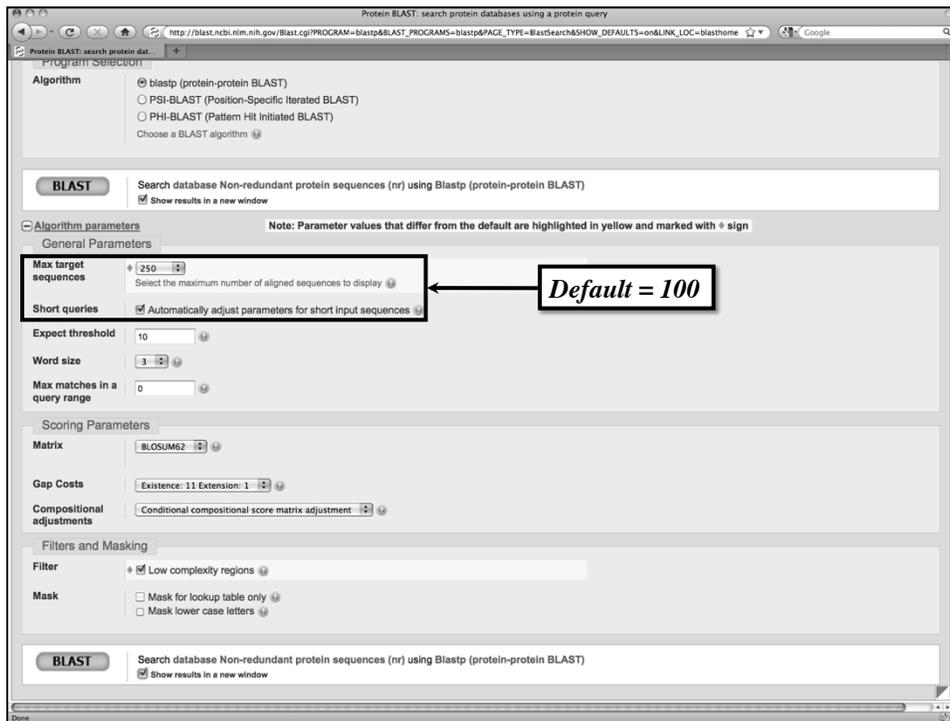
Program Selection

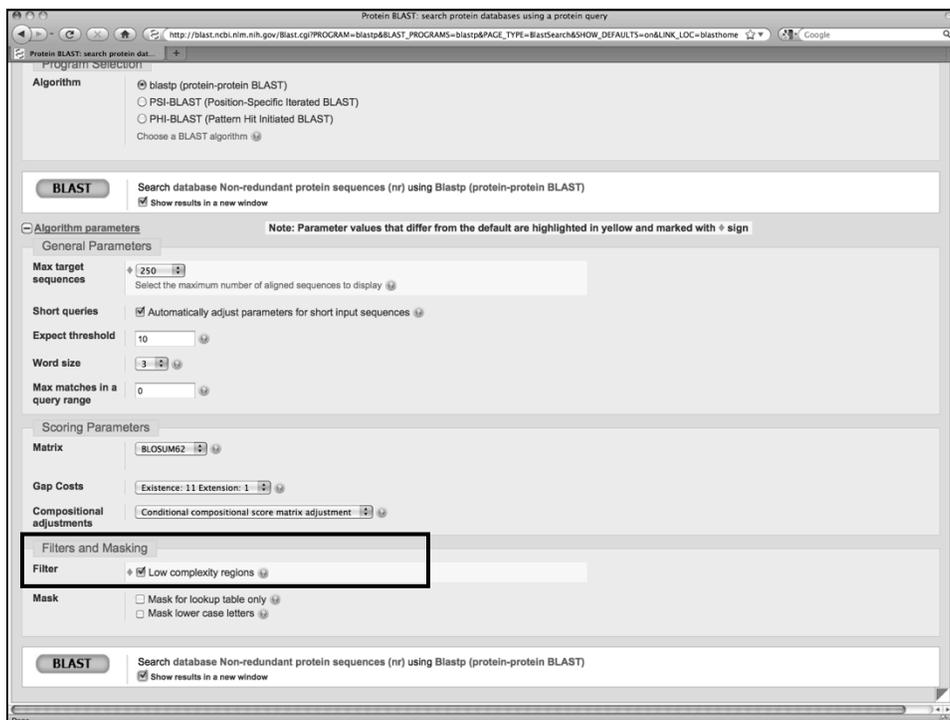
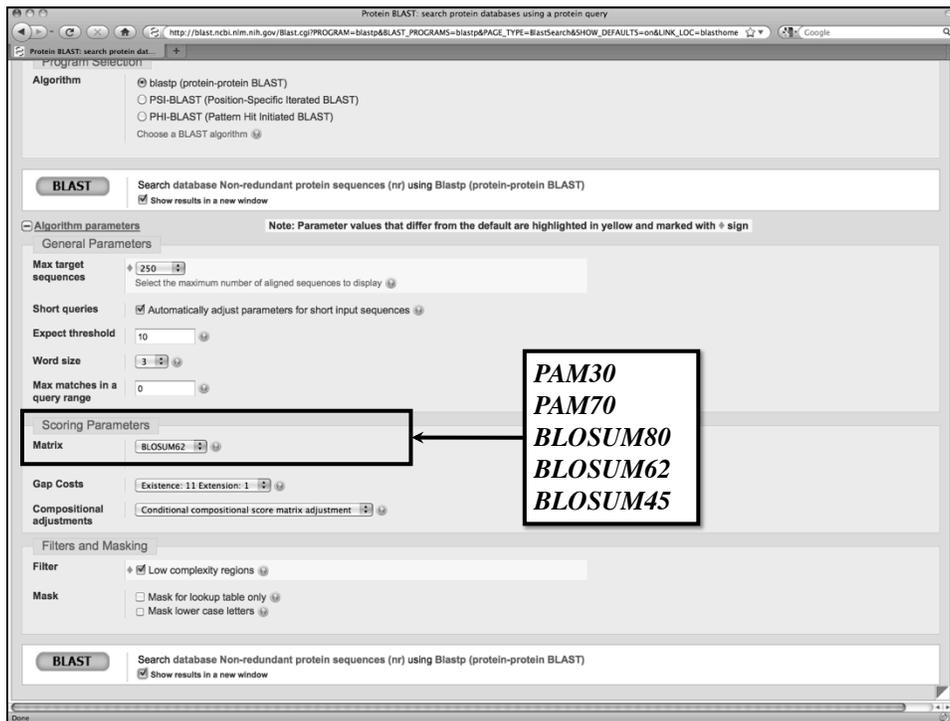
Algorithm: blastp (protein-protein BLAST)

BLAST Search database Non-redundant protein sequences (nr) using Blastp (protein-protein BLAST)

Algorithm parameters

Limit by organism or taxonomic group





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 (ASH1)  
MESSAKMESGGAGQQPQPQQPFLPPAACFFAIAAAAAAAAAAAAAQSAQQQQQQQQQQQAPQLRPAA  
DGQPSGGGHSAPKQVKRQRSSPELNRCKRRLNFSGFGYSLPQQQMAAVARRNERERNRVRLVNLGFAT  
LREHVPNGAANKKMSKVTLSAVEYIRALQQLLDERDAVSAAFQAVLSPTTISPNYNDLNSMAGSPVS  
SYSSDEGSYDPLSPERQELLDFTNWF
```

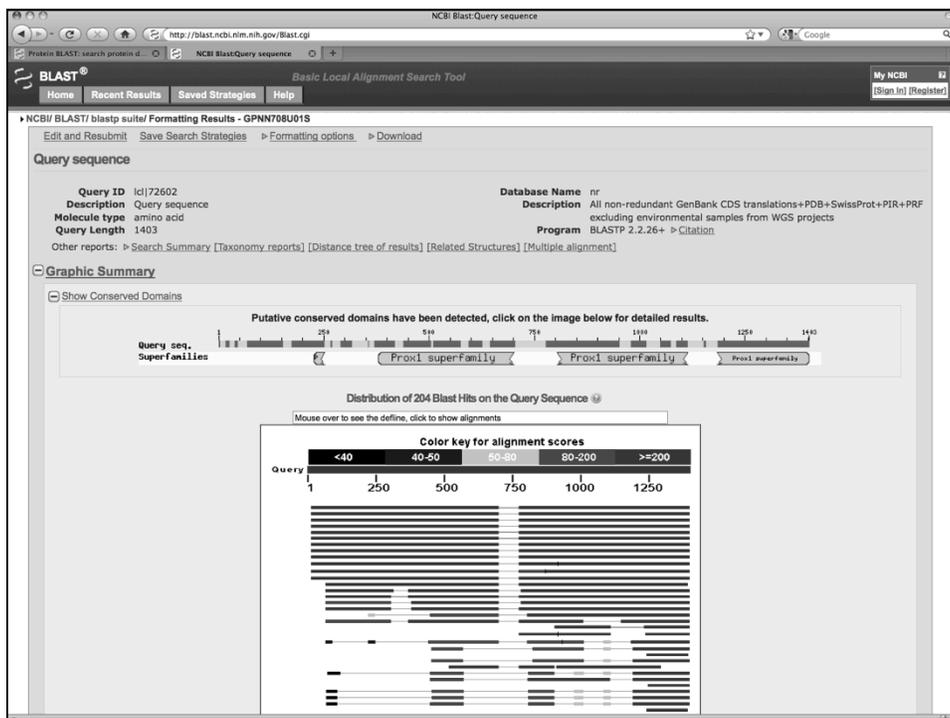
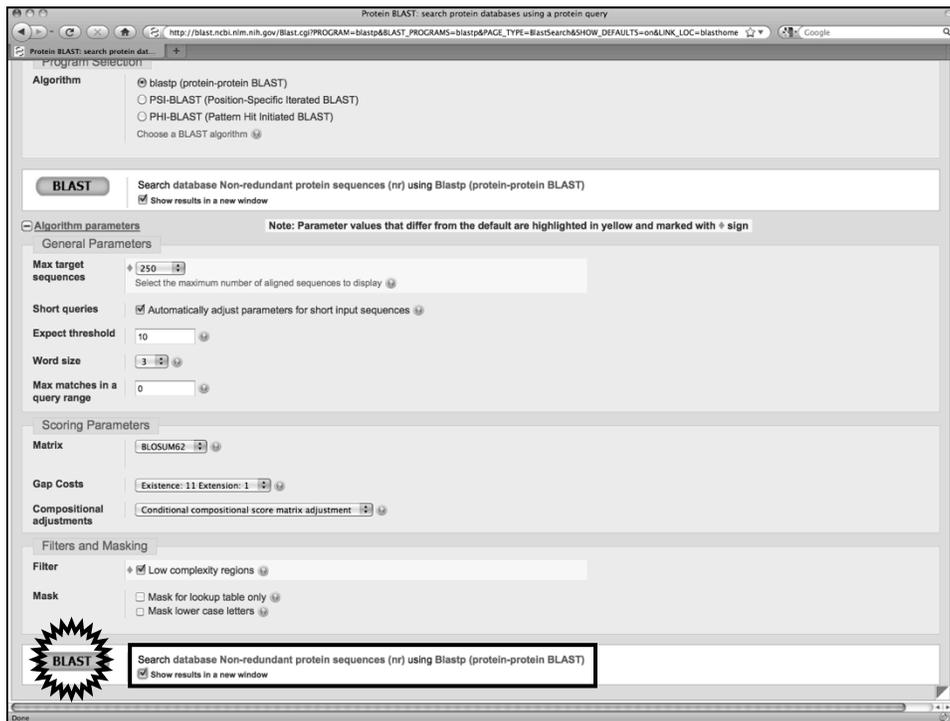
*Homopolymeric
alanine-glutamine tract*

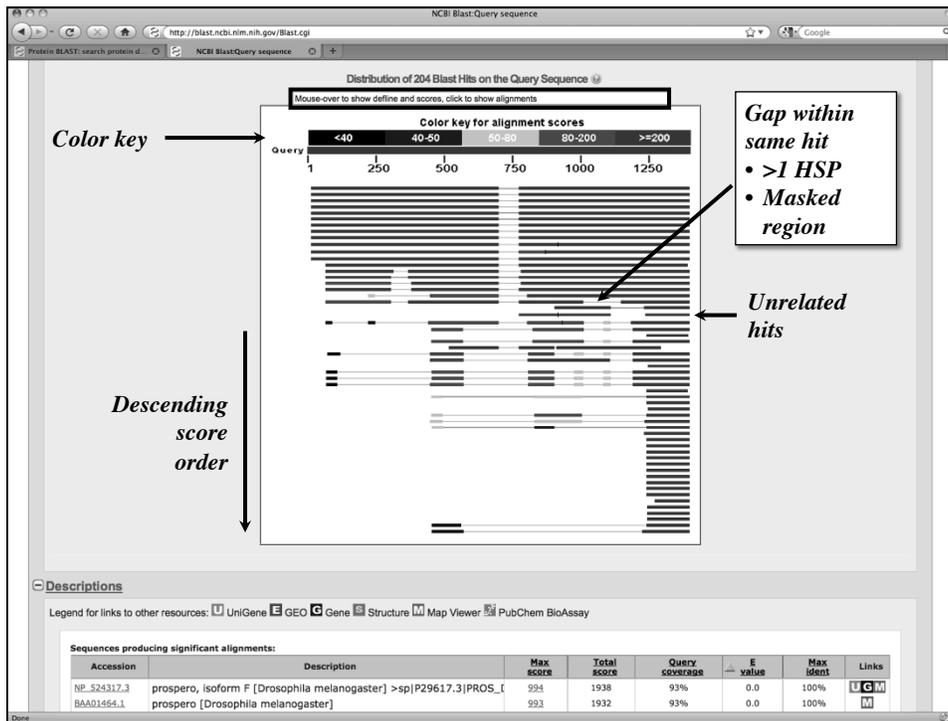


Identifying Low-Complexity Regions

- 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 (but *not* enabled by default)







Descriptions

Legend for links to other resources: UniGene GEO Gene Structure Map Viewer PubChem BioAssay

Sequences producing significant alignments:

Accession	Description	Max score	Total score	Query coverage	E value	Max ident	Links
NP_524317.3	prospero, isoform F [Drosophila melanogaster] >sp P29617.3 PROS_	994	1938	93%	0.0	100%	U G M
BAA01464.1	prospero [Drosophila melanogaster]	993	1932	93%	0.0	100%	M
CAA77802.1	prospero [Drosophila melanogaster]	993	1936	93%	0.0	100%	M
AAF05703.1	homeodomain transcription factor Prospero [Drosophila melanogaster]	990	1821	93%	0.0	100%	M
XP_001980573.1	GG18089 [Drosophila erecta] >gb EDV49531.1 GG18089 [Drosophila	989	1885	93%	0.0	99%	G M
AAA28841.1	Pros protein [Drosophila melanogaster]	982	1811	93%	0.0	97%	M
XP_002097201.1	GE26090 [Drosophila yakuba] >gb EDW96913.1 GE26090 [Drosophi	981	1885	93%	0.0	97%	G M
NP_788636.2	prospero, isoform G [Drosophila melanogaster] >gb AAN13500.3 pro	944	1862	93%	0.0	100%	U G M
AAT94492.1	LD37627p [Drosophila melanogaster]	943	1858	93%	0.0	100%	M
NP_731565.3	prospero, isoform E [Drosophila melanogaster] >gb AAN13501.3 pro	942	1864	93%	0.0	100%	U G M
XP_002031631.1	GM23939 [Drosophila sechellia] >gb EDW42617.1 GM23939 [Drosop	935	1987	93%	0.0	98%	G M
XP_001954214.1	GF16857 [Drosophila ananassae] >gb EDV42775.1 GF16857 [Drosop	904	1673	93%	0.0	92%	G M
XP_001359985.2	GA14403 [Drosophila pseudoobscura pseudoobscura] >gb EAL29137.	869	1499	89%	0.0	83%	G M
XP_002069959.1	GK11290 [Drosophila williston] >gb EDW80945.1 GK11290 [Drosopi	845	1532	85%	0.0	80%	G M
XP_002053284.1	pros [Drosophila virilis] >gb EDW66804.1 pros [Drosophila virilis]	821	1429	85%	0.0	79%	G M
XP_001994360.1	GH21437 [Drosophila grimshawi] >gb EDV95096.1 GH21437 [Droso	809	1374	84%	0.0	80%	G M
XP_002000130.1	G122696 [Drosophila mojavensis] >gb EDW15591.1 G122696 [Droso	804	1392	84%	0.0	80%	G M
XP_001655942.1	homeobox protein prospero/prox-1 [Aedes aegypti] >gb EAT46002.1	571	770	62%	0.0	85%	U G M
Q9J6A1.1	RecName: Full=Homeobox protein prospero >gb AAF06660.1 AF1904	430	1299	75%	4e-124	79%	M
XP_002103874.1	prospero [Drosophila simulans] >gb EDX13377.1 prospero [Drosophi	372	600	26%	1e-111	100%	U G M
IXPX_A	Chain A, Structural Basis Of Prospero-Dna Interaction; Implications Fo	347	347	11%			
XP_309606.5	AGAP004052-PA [Anopheles gambiae str. PEST] >gb EAA05345.5 AG	382	915	54%			
EFA07555.1	prospero [Tribolium castaneum]	365	706	38%			
1MJ1_A	Chain A, Crystal Structure Of The Homeo-Prospero Domain Of D. Mela	315	315	10%	4e-97	97%	
XP_971664.2	PREDICTED: similar to homeobox protein prospero/prox-1 [Tribolium c	342	702	38%	4e-97	89%	U G M
XP_002019831.1	GL11998 [Drosophila persimilis] >gb EDW38465.1 GL11998 [Drosop	341	783	49%	2e-96	93%	U G M
NP_001164363.1	homeobox protein prospero [Nasonia vitripennis]	346	700	38%	3e-96	94%	U G M
XP_002427668.1	homeobox protein prospero/prox-1/ceh-26, putative [Pediculus huma	345	705	45%	4e-96	74%	G
E1071784.1	hypothetical protein KGM_10139 [Danaus plexippus]	311	311	10%	1e-95	92%	
XP_003395765.1	PREDICTED: hypothetical protein LOC100645194 [Bombus terrestris]	335	698	37%	3e-93	97%	G M
XP_003489214.1	PREDICTED: hypothetical protein LOC100746817 [Bombus impatiens]	335	697	37%	3e-93	97%	G M
XP_392355.4	PREDICTED: hypothetical protein LOC406073 [Apis mellifera]	333	700	38%	1e-92	97%	U G M
BAH83641.1	prospero [Bombyx mori]	298	298	10%	4e-91	91%	U G M
XP_002410204.1	prospero protein, putative [Ixodes scapularis] >gb EEC12842.1 pros	301	425	26%	7e-86	83%	U G M

Annotations:

- 0.0 means $\leq 10^{-1000}$** : Points to the E value column.
- $4e-97 = 4 \times 10^{-97}$** : Points to a specific E value.

Accession	Description	Length	Score	E-value	Percent Identity
CAG09138.1	unnamed protein product [Tetraodon nigroviridis]	175	175	10%	5e-48
XP_003087294.1	hypothetical protein CRE_23850 [Caenorhabditis remanei] >gb EFO96	174	174	7%	3e-46
AAC59781.1	prospero_like protein [Takifugu rubripes]	156	156	9%	1e-41
XP_002475867.1	homeobox protein prospero/prox-1/ceh-26 [Schistosoma mansoni] >e	167	167	10%	4e-41
XP_003314467.1	PREDICTED: prospero homeobox protein 2 isoform 1 [Pan troglodytes	160	160	8%	1e-40
BAC04278.1	unnamed protein product [Homo sapiens]	157	157	8%	2e-40
EH015430.1	Prospero homeobox protein 2 [Heterocephalus glaber]	157	157	8%	2e-40
XP_002805213.1	PREDICTED: prospero homeobox protein 2 isoform 2 [Macaca mulatta	160	160	8%	3e-40
XP_003260805.1	PREDICTED: prospero homeobox protein 2-like [Nomascus leucogenys]	160	160	8%	3e-40
XP_002824990.1	PREDICTED: prospero homeobox protein 2-like isoform 2 [Pongo abeli	160	160	8%	4e-40
NP_001073877.2	prospero homeobox protein 2 isoform 2 [Homo sapiens] >gb EAWB11	158	158	8%	2e-39
BAB17311.1	Prox 1 [Cynops pyrrhogaster]	161	203	16%	3e-38
XP_003149047.1	hypothetical protein LOAG_13494 [Loa loa] >gb EFO15022.1 hypothe	145	145	7%	5e-38
EFN67531.1	Homeobox protein prospero [Camponotus floridanus]	157	422	23%	2e-36
ED02840.1	RIKEN cDNA 1700058C01, isoform CRA_b [Mus musculus]	154	154	8%	2e-36
GA051489.1	prospero homeobox protein 2 [Clonorchis sinensis]	147	147	14%	1e-35
CA115309.1	prospero homeobox 1 [Homo sapiens]	154	198	15%	2e-35
EFB18550.1	hypothetical protein PANDA_009835 [Alluorhoda melanoleuca]	152	197	16%	6e-35
CAG09167.1	unnamed protein product [Tetraodon nigroviridis]	150	190	18%	4e-34
EFZ18533.1	hypothetical protein SINV_16510 [Solenopsis invicta]	126	126	4%	1e-31
EH071783.1	prospero [Danaus plexippus]	141	383	36%	3e-31
EGW02786.1	Prospero homeobox protein 2 [Cricetus griseus]	101	101	5%	4e-23
CAG13403.1	unnamed protein product [Tetraodon nigroviridis]	100	100	4%	8e-23
XP_003150996.1	hypothetical protein LOAG_14553 [Loa loa] >gb EFO13973.1 hypothe	105	105	3%	8e-22
EG67129.1	Homeobox protein prospero [Acromyrmex echinator]	104	385	21%	1e-19
EFN07231.1	Homeobox protein prospero [Harpagathos saltator]	104	387	21%	1e-19
EH063774.1	hypothetical protein EGM_16808 [Macaca fascicularis]	99.8	99.8	4%	2e-19
EH028047.1	hypothetical protein EGK_18383 [Macaca mulatta]	99.8	99.8	4%	2e-19
AA030180.1	homeobox prospero-like protein [Homo sapiens]	97.4	97.4	4%	1e-18
XC5499	Prox 1 protein 671 - chicken	80.1	183	19%	3e-12
XP_003366161.1	homeobox protein ceh-26 [Trichinella spiralis] >gb EFV48171.1 home	57.0	57.0	3%	1e-07
CAG94249.1	unnamed protein product [Tetraodon nigroviridis]	43.5	43.5	3%	0.006
NP_001106671.1	prospero homeobox protein 1 [Rattus norvegicus] >gb EDL94973.1 p	42.0	42.0	8%	0.19
CAP58229.1	Prox1 protein [Xenopus (Silurana) tropicalis]	42.0	42.0	8%	1.2
AA110209.1	transcription factor Prox1 [Notophthalmus viridescens]	40.4	40.4	7%	3.0
YP_004342610.1	hypothetical protein Arcve_1902 [Archaeoglobus venificus SNP6] >gb	38.1	38.1	2%	5.1
ABG29070.1	transcription factor Prox1 [Pleurodeles waltii]	38.9	38.9	7%	6.7

>ref|NP_731565.3| UGM prospero, isoform E [Drosophila melanogaster]
 gb|AA013501.3| prospero, isoform E [Drosophila melanogaster]
 Length=1835

GENE ID: 41363 pros | prospero [Drosophila melanogaster]
 (Over 100 PubMed links)

Sort alignments for this subject sequence by:
 E value Score Percent identity
 Query start position Subject start position
 Identities = 688/688 (100%), Positives = 688/688 (100%), Gaps = 0/688 (0%)

Query	Subject	Score	Percent Identity
Query 17	Sbjct 317	76	376
Query 77	Sbjct 136	136	136
Query 137	Sbjct 437	496	496
Query 197	Sbjct 256	256	256
Query 257	Sbjct 616	616	616
Query 317	Sbjct 376	376	376
Query 617	Sbjct 736	736	736
Query 437	Sbjct 796	796	796
Query 797	Sbjct 796	796	796
Query 497	Sbjct 856	856	856
Query 557	Sbjct 616	616	616

NCBI Blast Query sequence

Query 557 LAEMQKQYVQLCSRMEQSEccgldqddvegeqepdNGSSDHIELSPSPITLGDGDVSP 616
 Sbjct 857 LAEMQKQYVQLCSRMEQSECCQLDQDDVEQEPEPNDGSSDHIELSPSPITLGDGDVSP 916
 LAEMQKQYVQLCSRMEQSECCQLDQDDVEQEPEPNDGSSDHIELSPSPITLGDGDVSP

Query 617 NHKEETGQERPGSSSPSPSPKPKTSLGESSDSGANMLSQMMSKMSGKLNHPLVGVGHP 676
 Sbjct 917 NHKEETGQERPGSSSPSPSPKPKTSLGESSDSGANMLSQMMSKMSGKLNHPLVGVGHP 976

Query 677 ALPQGFPPLLQHMGMDSHAAMYYQFFF 704
 Sbjct 977 ALPQGFPPLLQHMGMDSHAAMYYQFFF 1004

Score = 636 bits (1640), Expect = 0.0, Method: Compositional matrix adjust.
 Identities = 461/498 (93%), Positives = 463/498 (93%), Gaps = 32/498 (6%)

Query 906 PQNGPTPATQSAAMFQAPKTPQGMNFVAAAALYNSMTGPFCLPPDqqqqqtaqqqssa 965
 P P+P +AAAMFQAPKTPQGMNFVAAAALYNSMTGPFCLPPDQQQQQTAQQQSA 1426
 PHIRFSP---TAAAMFQAPKTPQGMNFVAAAALYNSMTGPFCLPPDQQQQQTAQQQSA

Query 966 qqqqqssggtqqqLEQNEALSLVVTFKKRHKVTDTRITPRTVSRILAQDvvpptggpp 1025
 Sbjct 1427 QQQQQSSQQTQQQLEQNEALSLVVTFKKRHKVTDTRITPRTVSRILAQDGVVPTGGPP 1486
 QQQQQSSQQTQQQLEQNEALSLVVTFKKRHKVTDTRITPRTVSRILAQDGVVPTGGPP

Query 1026 stpqggggggggggggggggggggASNGNSNATPAQSPTRSSGGAAYHppppppppppmp 1085
 Sbjct 1487 STPQQQQQQQQQQQQQQQQQQQQASNGNSNATPAQSPTRSSGGAAYHPPPPPPPPMP 1546
 STPQQQQQQQQQQQQQQQQQQQQASNGNSNATPAQSPTRSSGGAAYHPPPPPPPPMP

Query 1086 VSLPTSVAIPNPSLHESKVFSPYSPFFNPhaaaggataaqlhghqghhphhgsmglsss 1145
 Sbjct 1547 VSLPTSVAIPNPSLHESKVFSPYSPFFNPhaaaggataaqlhghqghhphhgsmglsss 1606
 VSLPTSVAIPNPSLHESKVFSPYSPFFNPhaaaggataaqlhghqghhphhgsmglsss

Query 1146 ppgslGALMDSRDSppLpHppsmLhpallaahhggsPDYKTCRAVMDAQDRGSECNSA 1205
 Sbjct 1607 PPGSLGALMDSRDSFPLPpHppsmLhpallaahhggsPDYKTCRAVMDAQDRGSECNSA 1666
 PPGSLGALMDSRDSFPLPpHppsmLhpallaahhggsPDYKTCRAVMDAQDRGSECNSA

Query 1206 DMQFDGMAPTISFYKQMLKTEHQESLMAKHCESLTPHSSLTPMHLRKAALMFVWRY 1265
 Sbjct 1667 DMQFDGMAPT-----SSTLTPMHLRKAALMFVWRY 1697
 DMQFDGMAPT-----SSTLTPMHLRKAALMFVWRY

Query 1266 PSSAVLKMYPDIKFNKNNTAQLVKWFSNFRFYIOMEKYARQAVTEGIKTPDDLIIAG 1325
 Sbjct 1698 PSSAVLKMYPDIKFNKNNTAQLVKWFSNFRFYIOMEKYARQAVTEGIKTPDDLIIAG 1757
 PSSAVLKMYPDIKFNKNNTAQLVKWFSNFRFYIOMEKYARQAVTEGIKTPDDLIIAG

Query 1326 DSELYRVLNLYNRNHHIEVQNFVVESTLREFFRAIQGGKDEQSWKSIYKILSRM 1385
 Sbjct 1758 DSELYRVLNLYNRNHHIEVQNFVVESTLREFFRAIQGGKDEQSWKSIYKILSRM 1817
 DSELYRVLNLYNRNHHIEVQNFVVESTLREFFRAIQGGKDEQSWKSIYKILSRM

Query 1386 DDPVPEYFKSPNFLEQLE 1403
 Sbjct 1818 DDPVPEYFKSPNFLEQLE 1835

No definition line -> Second HSP identified

- Gap a Low-Complexity

Score = 942 bits (2435), Expect = 0.0 ✓ Method: Compositional matrix adjust.
 Identities = 688/688 (100%) ✓ Positives = 688/688 (100%), Gaps = 0/688 (0%)

Score = 636 bits (1640), Expect = 0.0 ✓ Method: Compositional matrix adjust.
 Identities = 461/498 (93%) ✓ Positives = 463/498 (93%), Gaps = 32/498 (6%)

HSP 1
 Q: 17- 704
 S: 317-1004

HSP 2
 Q: 906-1403
 S: 1370-1835

Color key for alignment score

Query	<40	40-50	50-60	60-70	70-80	80-90	>=200
1	200	500	750	1000	1250		

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
- Repetitive elements
 - LINEs, SINEs, retroviral repeats
 - Choose “Filter: Species-Specific Repeats” with BLASTN
 - RepeatMasker
<http://www.repeatmasker.org>
- 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



BLAST: Basic Local Alignment Search Tool

<http://www.ncbi.nlm.nih.gov/BLAST>

Human
Mouse
Rat
Arabidopsis thaliana
Oryza sativa
Bos taurus
Danio rerio
Drosophila melanogaster
Pan troglodytes
Microbes
Apis mellifera

Basic BLAST

Choose a BLAST program to run.

nucleotide blast Search a nucleotide database using a nucleotide query
Algorithms: blastn, megablast, discontinuous megablast

protein blast Search protein database using a protein query
Algorithms: blastp, psi-blast, phi-blast

blastx Search protein database using a translated nucleotide query

tblastn Search translated nucleotide database using a protein query

tblastx Search translated nucleotide database using a translated nucleotide query

Specialized BLAST

Choose a type of specialized search (or database name in parentheses.)

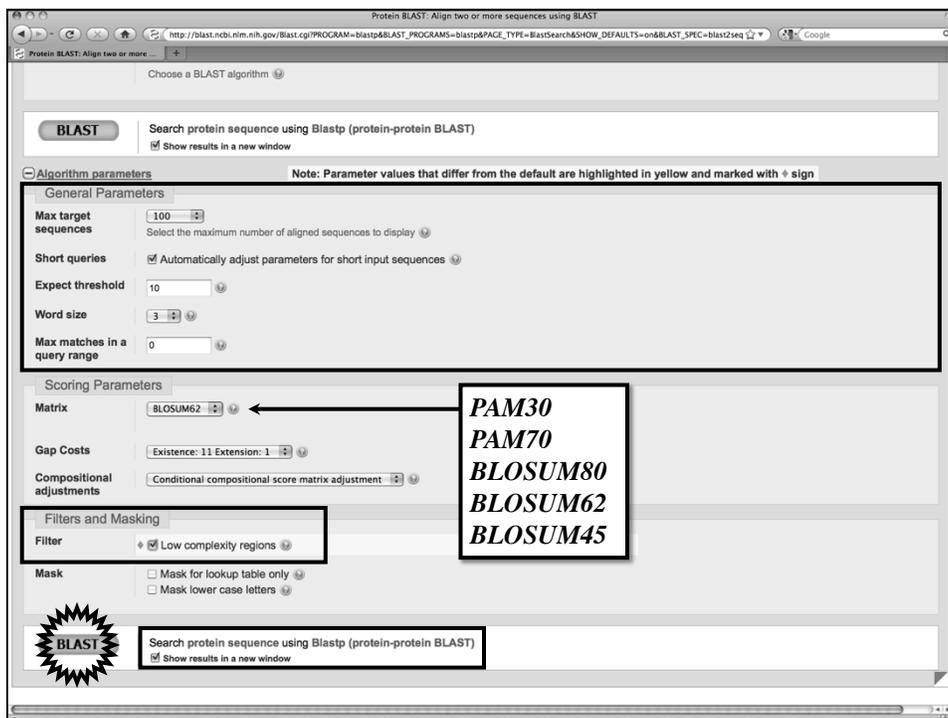
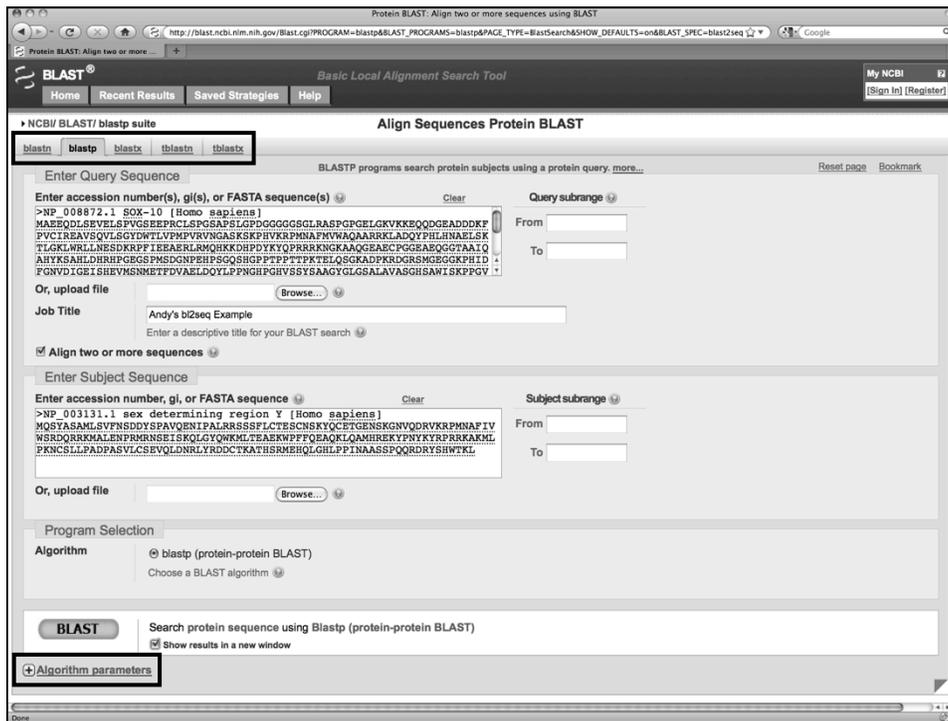
- Make specific primers with **Primer-BLAST**
- Search **trace archives**
- Find **conserved domains** in your sequence (cde)
- Find sequences with similar **conserved domain architecture** (cdart)
- Search sequences that have **gene expression profiles** (GEO)
- Search **immunoglobulins** (IgBLAST)
- Search using **SNP flanks**
- Screen sequence for **vector contamination** (vecscreen)
- Align two (or more) sequences** using BLAST (bl2seq) ←
- Search **protein or nucleotide targets** in PubChem BioAssay
- Search SRA **transcript and genomic libraries**
- Constraint Based Protein **Multiple Alignment Tool**
- Needleman-Wunsch **Global Sequence Alignment Tool**
- Search **RefSeqGene**
- Search **WGS sequences** grouped by organism

Tip of the Day

Use Genomic BLAST to see the genomic context

If you are interested in the evolution of a particular gene or gene family it is often interesting to examine the intro-exon structure even across species.

[More tips...](#)



NCBI BLAST: Andy's bl2seq Example

Protein BLAST: Align two or more sequences against a protein database

Basic Local Alignment Search Tool

Home Recent Results Saved Strategies Help

My NCBI [Sign In] [Register]

NCBI/BLAST/blastp suite-2sequences/ Formatting Results - GR3U63S2112

Edit and Resubmit Save Search Strategies Formatting options Download

Andy's bl2seq Example

Query ID |cl|57577
Description NP_008872.1 SOX-10 [Homo sapiens]
Molecule type amino acid
Query Length 466

Subject ID 57579
Description NP_003131.1 sex determining region Y [Homo sapiens]
Molecule type amino acid
Subject Length 204
Program BLASTP 2.2.26+ Citation

Other reports: Search Summary Taxonomy reports Multiple alignment

Graphic Summary

Distribution of 2 Blast Hits on the Query Sequence

Mouse over to see the define, click to show alignments

Color key for alignment scores

Score Range	Count
<40	90
40-50	180
50-80	270
80-200	360
>=200	450

Dot Matrix View

Descriptions

Legend for links to other resources: UniGene GEO Gene Structure Map Viewer PubChem BioAssay

Sequences producing significant alignments:

Accession	Description	Max score	Total score	Query coverage	E value	Max ident	Links
57579	NP_003131.1 sex determining region Y [Homo sapiens]	94.0	109	19%	4e-24	46%	

NCBI BLAST: Andy's bl2seq Example

Protein BLAST: Align two or more sequences against a protein database

Basic Local Alignment Search Tool

Home Recent Results Saved Strategies Help

My NCBI [Sign In] [Register]

NCBI/BLAST/blastp suite-2sequences/ Formatting Results - GR3U63S2112

Edit and Resubmit Save Search Strategies Formatting options Download

Andy's bl2seq Example

Query ID |cl|57577
Description NP_008872.1 SOX-10 [Homo sapiens]
Molecule type amino acid
Query Length 466

Subject ID 57579
Description NP_003131.1 sex determining region Y [Homo sapiens]
Molecule type amino acid
Subject Length 204
Program BLASTP 2.2.26+ Citation

Other reports: Search Summary Taxonomy reports Multiple alignment

Graphic Summary

Distribution of 2 Blast Hits on the Query Sequence

Mouse over to see the define, click to show alignments

Color key for alignment scores

Score Range	Count
<40	90
40-50	180
50-80	270
80-200	360
>=200	450

Dot Matrix View

Descriptions

Legend for links to other resources: UniGene GEO Gene Structure Map Viewer PubChem BioAssay

Sequences producing significant alignments:

Accession	Description	Max score	Total score	Query coverage	E value	Max ident	Links
57579	NP_003131.1 sex determining region Y [Homo sapiens]	94.0	109	19%	4e-24	46%	

Alignments

>|cl|57579 NP_003131.1 sex determining region Y [Homo sapiens]
 Length=204

Sort alignments for this subject sequence by:
 E value Score Percent identity
 Query start position Subject start position

Score = 94.0 bits (232), Expect = 4e-24, Method: Compositional matrix adjust.
 Identities = 39/84 (46%), Positives = 62/84 (74%), Gaps = 0/84 (0%)

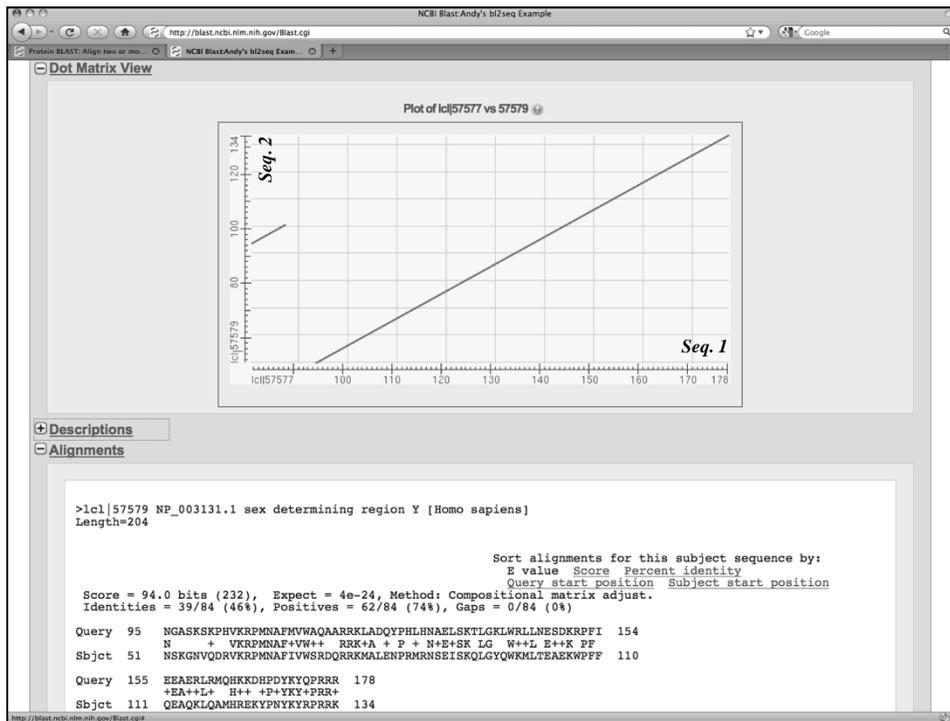
```

Query 95  NGASKSPHVKRPMPNMFVWQAARRKLDQYPHLHNAELSKTLGKLRLLNESDKRPFI 154
          N      + VKRPMNAF+VW++ RRK+A + P + N+E+SK LG W++L E++K PF
Sbjct 51  NSKGNVQDRVKRPMNAFVWGRDQRRKMALENFRMRNSEISKQLGYQWKMLTEAEKPPFF 110
    
```

Score = 15.4 bits (28), Expect = 1.9, Method: Compositional matrix adjust.
 Identities = 3/7 (43%), Positives = 5/7 (71%), Gaps = 0/7 (0%)

```

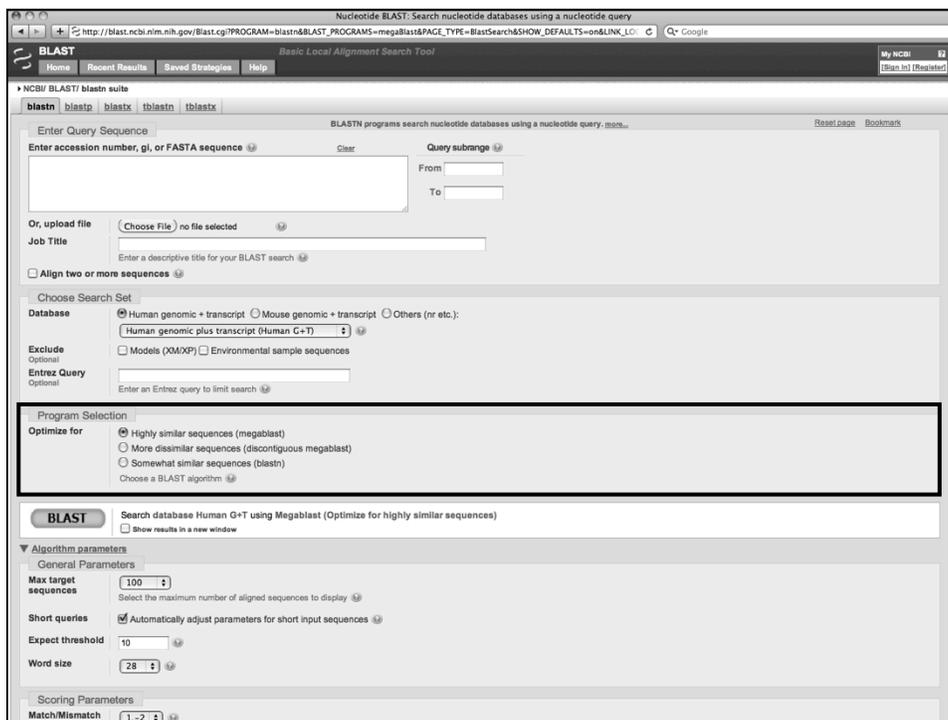
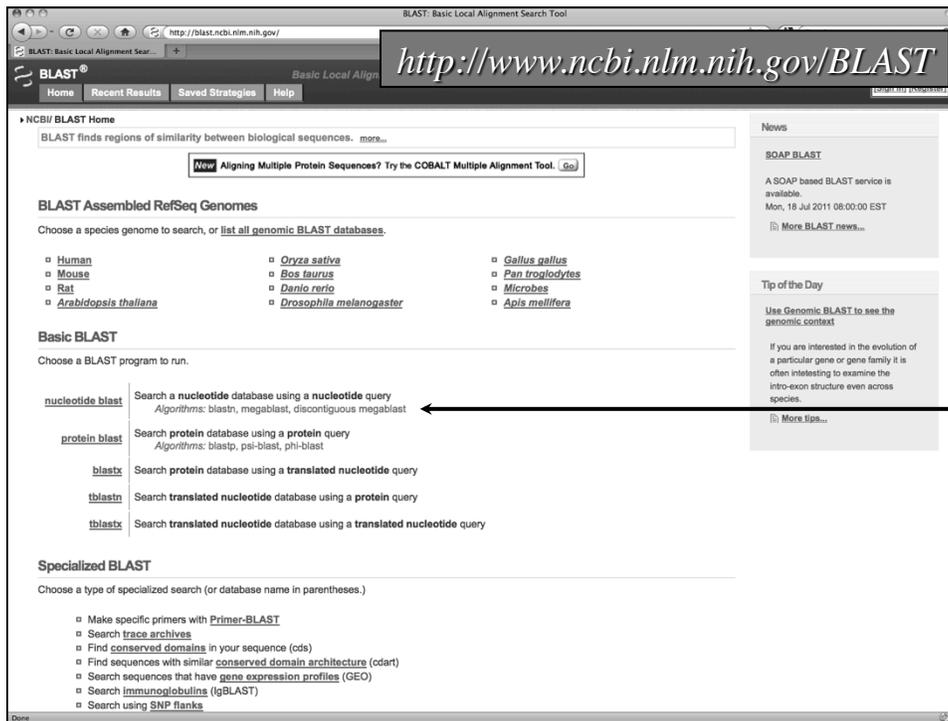
Query 82  GYDWLV 88
          GY W ++
Sbjct 95  GYQWKML 101
    
```



Nucleotide BLAST

MegaBLAST and BLASTN





Nucleotide-Based BLAST Algorithms

	<u>W</u>	<u>+/-</u>	<u>Gaps</u>
Optimized for aligning very long and/or highly similar sequences ($\geq 95\%$)			
MegaBLAST (default)	28	1, -2	Linear
Better for diverged sequences and/or cross-species comparisons ($< 80\%$)			
Discontiguous MegaBLAST	11	2, -3	Affine
BLASTN	11	2, -3	Affine
Finding short, nearly exact matches (< 20 bases)			
BLASTN <i>E = 1000, all filtering off</i>	7	2, -3	Affine



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



The screenshot shows the UCSC Genome Browser homepage. At the top, there is a navigation bar with links for Genomes, Blat, Tables, Gene Sorter, PCR, VisiGene, Proteome, Session, FAQ, and Help. A sidebar on the left lists various tools and resources like Genome Browser, ENCODE, Neandertal, Blat, Table Browser, Gene Sorter, In Silico PCR, Genome Graphs, Galaxy, VisiGene, Proteome Browser, Utilities, Downloads, Release Log, Custom Tracks, Microbial Genomes, Mirrors, Archives, and Training. The main content area features a section titled "About the UCSC Genome Bioinformatics Site" with a welcome message. Below that is a "News" section with two recent updates: "3 January 2012 - Roadmap Epigenomics Now Available through Data Hub at Washington University" and "19 December 2011 - Variant Call Format (VCF) Now Supported in Genome Browser".

ACTIONS	QUERY	SCORE	START	END	QSIZE	IDENTITY	CHRO	STRAND	START	END	SPAN
browser details	CB312815	710	1	733	768	98.1%	5	+	101455599	101456323	725
browser details	CB312815	29	501	537	768	89.2%	2	+	38736251	38736287	37
browser details	CB312815	25	501	529	768	93.2%	3	+	22960346	22960374	29
browser details	CB312815	22	341	363	768	100.0%	1	+	122930956	122930979	24
browser details	CB312815	21	202	222	768	100.0%	17	-	33248146	33248166	21
browser details	CB312815	21	706	727	768	100.0%	3	+	46857920	46857942	23
browser details	CB312815	21	552	574	768	95.7%	1	+	157973111	157973133	23
browser details	CB312815	20	277	298	768	95.5%	2	-	240446870	240446891	22
browser details	CB312815	20	442	461	768	100.0%	1	-	216323127	216323146	20
browser details	CB312815	20	508	527	768	100.0%	1	-	56102029	56102048	20
browser details	CB312815	20	453	474	768	95.5%	2	+	186587336	186587357	22

UCSC Genome Browser on Rat Nov. 2004 (Baylor 3.4/rn4) Assembly

position/search chr5:101,455,400-101,456,480 gene Jump clear size 1,081 bp. (configure)

chr5 (q31) 5011 5012 5013 502 5022 5024 5032 5033 5034 5035

Scale chr5: 101455898 101456000 101456098 101456196 101456294 101456392 101456490

STS Markers: STS Markers on Genetic and Radiation Hybrid Maps

Gap: Gap Locations

Other RefSeq: CB312815 Your Sequence from Blat Search

Rat mRNAs: Rat mRNAs From GenBank

Spliced ESTs: RAT ESTs That Have Been Spliced

Conservation: Vertebrate Multiz Alignment & Conservation

RepeatMasker: Repeat Masker

Repeat Elements by RepeatMasker

- red: Genome and query sequence have different bases at this position.
- orange: The query sequence has an insertion (or genome has a deletion / alignment gap) at this point.
- purple: The query sequence extends beyond the end of the alignment.
- green: The query sequence appears to have a polyA tail which is not aligned to the genome.

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
 Division of Intramural Research

Rat BLAT Results

BLAT Search Results

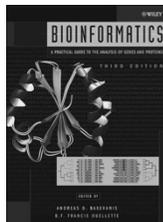
ACTIONS	QUERY	SCORE	START	END	QSIZE	IDENTITY	CHRO	STRAND	START	END	SPAN
browser details	CB312815	710	1	733	768	98.1%	5	+	101455599	101456323	725
browser details	CB312815	29	501	537	768	89.2%	2	+	38736251	38736287	37
browser details	CB312815	25	501	529	768	93.2%	3	+	22960346	22960374	29
browser details	CB312815	22	341	363	768	100.0%	1	+	122930956	122930979	24
browser details	CB312815	21	202	222	768	100.0%	17	-	33248146	33248166	21
browser details	CB312815	21	706	727	768	100.0%	3	+	46857920	46857942	23
browser details	CB312815	21	552	574	768	95.7%	1	+	157973111	157973133	23
browser details	CB312815	20	277	298	768	95.5%	2	-	240446870	240446891	22
browser details	CB312815	20	442	461	768	100.0%	1	-	216323127	216323146	20
browser details	CB312815	20	508	527	768	100.0%	1	-	56102029	56102048	20
browser details	CB312815	20	453	474	768	95.5%	2	+	186587336	186587357	22

FASTA

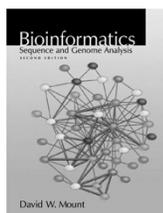
- Identifies regions of local alignment
- Employs an approximation of the Smith-Waterman algorithm to determine the best alignment between two sequences
- Method is significantly different from that used by BLAST
- Online implementations at
<http://fasta.bioch.virginia.edu>
<http://www.ebi.ac.uk/fasta33>



Further Reading



Chapter 11
Assessing Pairwise Sequence Similarity:
BLAST and FASTA



Chapter 6
Sequence Database Searching for
Similar Sequences

