
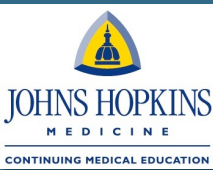



NATIONAL HUMAN GENOME RESEARCH INSTITUTE *Division of Intramural Research*




Current Topics in Genome Analysis 2016
Week 1: Biological Sequence Analysis I
Andy Baxevanis, Ph.D.

U.S. DEPARTMENT OF HEALTH AND HUMAN SERVICES | NATIONAL INSTITUTES OF HEALTH | genome.gov/DIR



Current Topics in Genome Analysis 2016
Andy Baxevanis, Ph.D.
***No Relevant Financial Relationships with
Commercial Interests***

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research



Sequence Alignments: Determining Similarity and Deducing Homology



Why construct 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
- Important to use correct terminology when describing phylogenetic relationships



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**

- **Homology:** Implies an evolutionary relationship
- **Homologs:** Genes that have arisen from a common ancestor
- Genes either *are* or *are not* homologous (not measured in degrees)

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

Orthologs: Genes that diverged as a result of a speciation event

- Sequences are direct descendants of a sequence in a common ancestor (share a common origin)
- Most likely have similar domain and three-dimensional structure
- Usually retain same biological function over evolutionary time
- Can be used to predict gene function in novel genomes



Defining the Terms

Paralogs: Genes that arose by the duplication of a single gene in a particular lineage

- Perhaps less likely to perform similar functions
- Can take on new functions over evolutionary time
- Provides insight into 'evolutionary innovation'



Defining the Terms

Paralogs

- 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

Orthology and Paralogy: Further Reading

Homology
 a personal view on some of the problems

Eugene Koonin
Annu. Rev. Genet.
39: 309-338, 2005

Orthologs, Paralogs, and Evolutionary Genomics¹

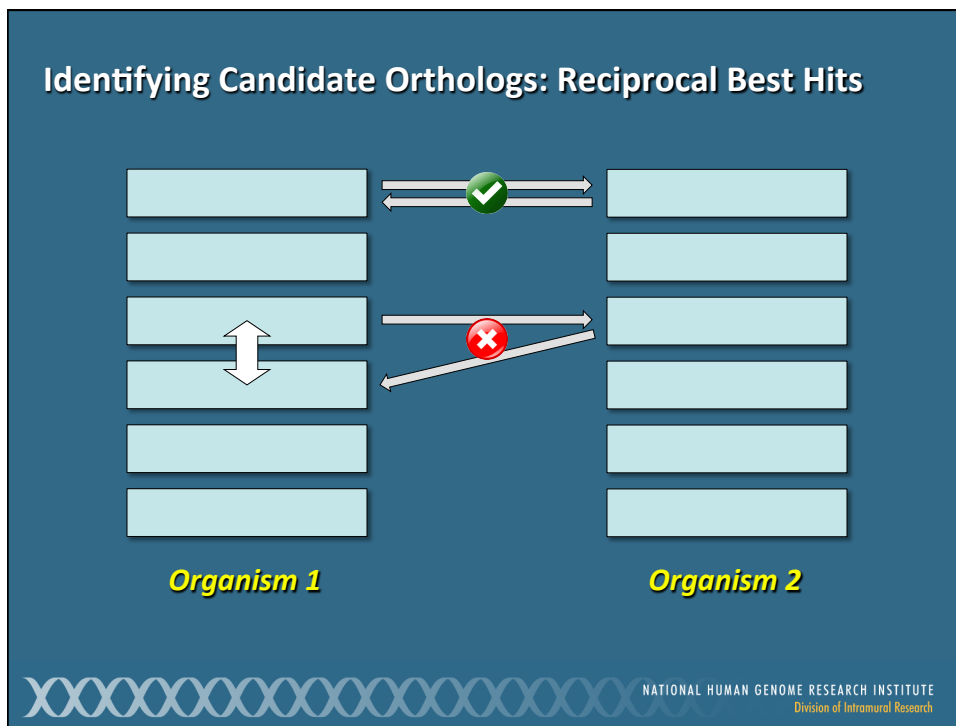
Walter Fitch
Trends Genet.
16: 227-231, 2000

Eugene V. Koonin
 National Center for Biotechnology Information, National Library of Medicine, Bethesda, Maryland 20894, USA
 e-mail: e.koonin@nih.gov

Key Words
 homologing, ortholog, paralog, pseudortholog, paralogizing, isolog

Abstract
 Orthologs and paralogs are two fundamentally different types of homologous genes that evolved, respectively, by vertical descent from a single ancestral gene and by duplication. Ortholog and paralog are key concepts of evolutionary genetics. A clear distinction between orthologs and paralogs is critical for the reconstruction of a robust evolutionary classification of genes and reliable functional annotation of newly sequenced genomes. Greater conceptual clarity in their nomenclature relationships with genes from taxonomically distinct species may be established for the majority of the genes from such sequenced genomes. This review examines in depth the definitions and subtleties of orthologs and paralogs, reviews the principal methodological approaches employed for identification of orthologs and paralogs, and reviews evolutionary and functional implications of these concepts.

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: Construction and Proper Selection



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 of considerations:
 - Cys/Pro are important for structure and function
 - Trp has a 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, additional 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	-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	-3	0	2	-1	-1	-3	-4	-1	-3	-3	-1	0	-1	-4	-3	-3	4	1	-1	-4
C	0	-3	-3	-3	6	-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	-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	-2	0	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	-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	-2	-1	-2	-4
S	1	-1	1	0	-1	0	0	0	-1	-2	-2	0	-1	-2	-1	4	1	-3	-2	-2	0	0	0	-4
T	0	-1	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	1	5	-2	-2	0	-1	-1	0	-4
W	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	11	2	-3	-4	-3	-2	-4
Y	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	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	-2	1	4	-1	-4
X	0	-1	-1	-1	-2	-1	-1	-1	-1	-1	-1	-1	-1	-2	0	0	-2	-1	-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

- Look only for differences in conserved, ungapped regions of a protein family ('blocks')
- Directly calculated based on local alignments
 - Substitution probabilities (*conservation*)
 - Overall *frequency* of amino acids
- Sensitive to detecting structural or functional substitutions
- Generally perform better than PAM matrices for local similarity searches (*Henikoff and Henikoff, 1993*)
- BLOSUM series can be used to identify both closely and distantly related sequences



BLOSUM n

- Built using sequences sharing no more than $n\%$ identity
- Contribution of sequences $> n\%$ identical clustered and replaced by a sequence that represents the cluster

Diagram illustrating the construction of BLOSUM80:

Initial sequences (with asterisks indicating high identity):

```
TGNQEYGTSSDSDSDY
KKLEKEEEEGISQESSEEE
KKLEKEEEEGISQESSEEE
KPAQEETEETSSQESAEED
```

After 80% identity threshold:

```
TGNQEYGTSSDSDSDY
KKLEKEEEEGISQESSEEE
KPAQEETEETSSQESAEED
```

After Clustering (replacing identical sequences with a representative):

```
TGNQEYGTSSDSDSDY
KKLEKEEEEGISQESSEEE
KPAQEETEETSSQESAEED
```

Final step: Calculate BLOSUM80

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

BLOSUM n

- Clustering reduces contribution of closely related sequences (less bias towards substitutions that occur in the most closely related members of a family)
- Reducing n yields more distantly related sequences
- Increasing n yields more closely related sequences

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

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

The takeaway...

No single matrix is the complete answer for all sequence comparisons

David Wheeler
 Curr. Protoc. Bioinformatics
 3.5.1 – 3.5.6, 2003

Selecting the Right Protein-Scoring Matrix

UNIT 3.5

OVERVIEW
 Every program for searching protein sequences against a database includes a choice of a "protein-scoring matrix," also called a "weight matrix." Weight matrices add sensitivity to the search, whose statistical significance is calculated by multiplying the PAM 250 matrix by the PAM 250 matrix.

Biologically, the PAM 250 matrix means that 100 amino acids, there have been 250 substitutions, while the PAM 100 matrix means that there have been 100 substitutions (at each site) over 100 generations (and mutations). This sounds unusual, but remember that our evolutionary time is a period that an alanine was changed to a glycine, then to a valine, and then back to alanine. These three substitutions are derived 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 proteins known at that time were small, globular, and hydrophilic. If the researcher believes that their protein contains substantial hydrophobic regions, such as membrane-spanning helices or sheets, the PAM matrices are less useful than others described in this unit. Dayhoff et al. (1978) were the first to define the scores for families and superfamilies. A protein family is defined as sequences related from 50% identical or greater to each other. A protein superfamily is defined as sequences related from 30% identical or greater to each other. While the terms "family" and "superfamily" are widely used in biology, most of the time the original definitions of Dayhoff and collaborators is not being used (see below).

Choosing all potential similarities: 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), but in which the two proteins are up to 300% different from each other (George et al., 1995). Another way to think about this is that the PAM 250

Contributed by David Wheeler
 Current Protocols in Bioinformatics (2003) 3.5.1-3.5.6
 Copyright © 2003 by John Wiley & Sons, Inc.

Choosing
 Substitution and
 Scoring
 Matrices
 35.1

Gaps

- Used to improve alignments between two sequences
 - Compensate for insertions and deletions
 - As such, *gaps represent biological events*
- Gaps must be kept to a reasonable number, to not reflect a biologically implausible scenario. About one gap per 20 residues is a 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: The Basic Local Alignment Search Tool



BLAST

- Seeks high-scoring segment pairs (HSPs)
 - 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' high-scoring segment pair for the two sequences being aligned

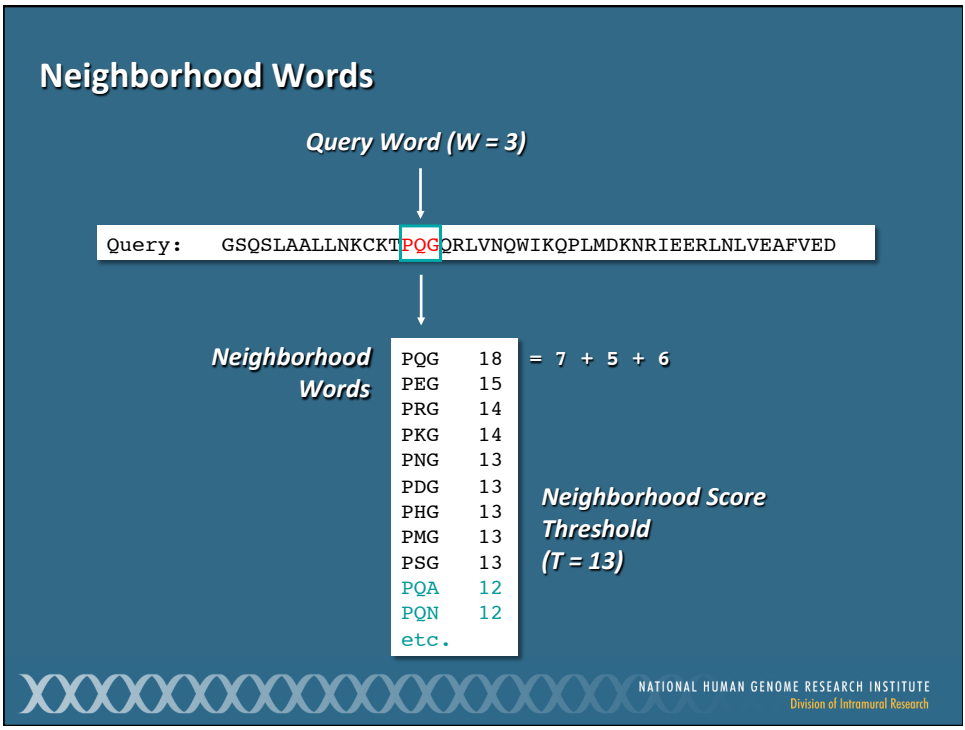
Altschul et al., J. Mol. Biol. 215: 403-410, 1990



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

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

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
 Division of Intramural Research

Extension

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

The graph shows a bell-shaped curve representing the cumulative score of an alignment as it is extended. The y-axis is labeled 'Cumulative Score' and the x-axis is 'Length of Alignment'. The curve starts at a baseline 'T', rises to a peak, and then decays towards a significance threshold 'S'. A vertical double-headed arrow labeled 'X' indicates the height of the peak above the threshold 'S'. To the right of the graph, the text 'Significance decay' is followed by a bulleted list: 'mismatches' and 'gap penalties'.

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
 Division of Intramural Research

Scores and Alignment Length Don't Tell the Whole Story

```
Query: 1 SGLKSLVGKTALLSGTSSKL 20
        SGLKSLVGKTALLSGTSSKL
Sbjct: 1 SGLKSLVGKTALLSGTSSKL 20
```

Score = 91

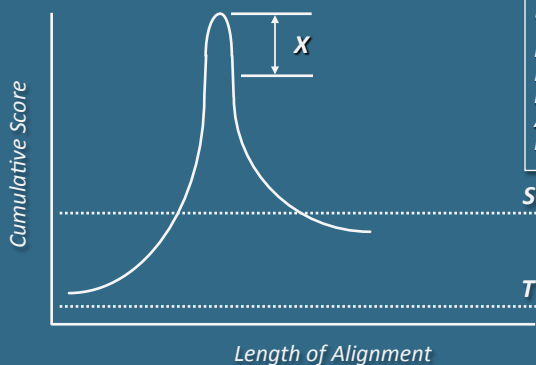
```
Query: 1 CQHMWYQWMIQCIWMYHCMQ 20
        CQHMWYQWMIQCIWMYHCMQ
Sbjct: 1 CQHMWYQWMIQCIWMYHCMQ 20
```

Score = 138



Scores and Probabilities

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



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

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



Scores and Probabilities

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

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

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

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
 Division of Intramural Research

Scores and Probabilities

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

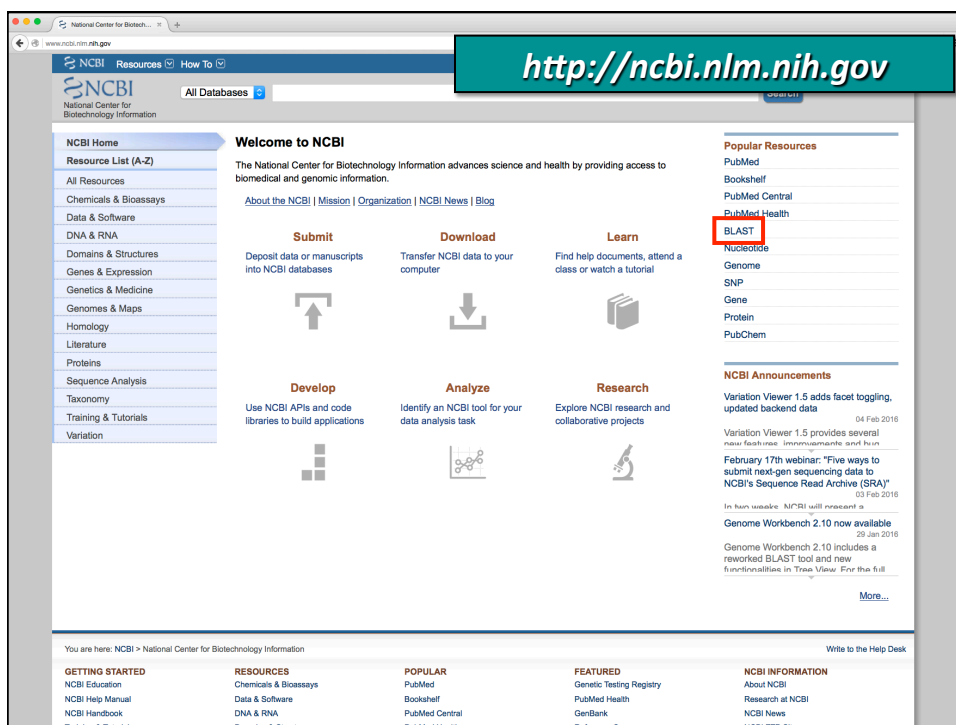
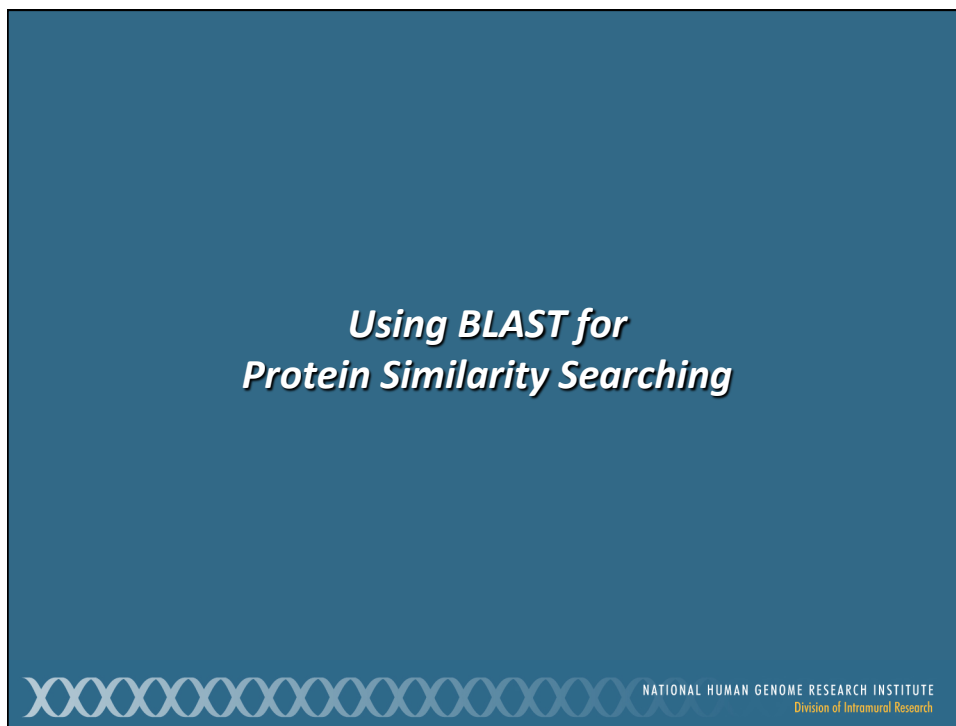
$$E \leq 10^{-6}$$

for nucleotides

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

for proteins

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
 Division of Intramural Research



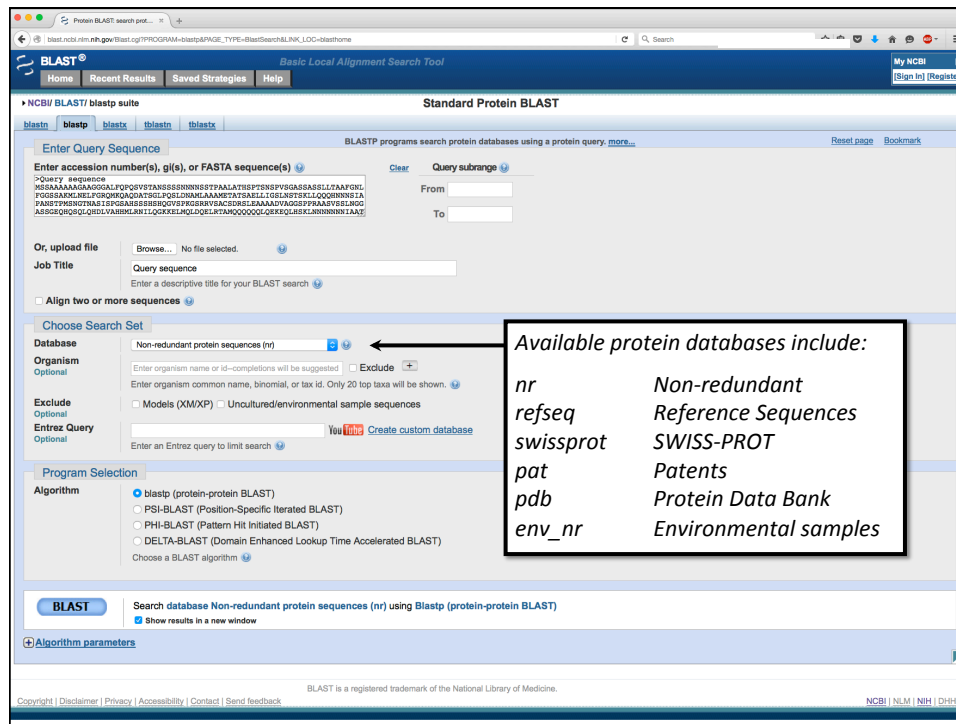
The screenshot shows the NCBI BLAST website. At the top, the URL **http://ncbi.nlm.nih.gov/BLAST** is highlighted in green. The page features a navigation bar with 'Home', 'Recent Results', 'Saved Strategies', and 'Help'. The main content is organized into sections: 'BLAST Assembled Genomes' with a list of organisms (Human, Mouse, Rat, Cow, Pig, Dog, Rabbit, Chimp, Guinea pig, Fruit fly, Honey bee, Chicken, Zebrafish, Clawed frog, Arabidopsis, Rice, Yeast, Microbes); 'Basic BLAST' with search options for nucleotide, protein, blastx, and tblastn; 'Specialized BLAST' with various advanced search tools; and a 'Your Recent Results' section on the right. A red arrow points to the 'protein blast' option.

Sequences Used in Examples

[http://research.nhgri.nih.gov/
teaching/seq_analysis.shtml](http://research.nhgri.nih.gov/teaching/seq_analysis.shtml)

This screenshot displays the NHGRI website page for 'Current Topics in Genome Analysis 2016'. The main heading is 'Weeks 1 and 4: Biological Sequence Analysis Protein and Nucleotide Sequences for Analysis'. The page shows BLAST search results for a query sequence. It includes alignment details such as 'BLAST 2 Sequences', 'Query: [sequence]', and 'Subject: [sequence]'. The alignment shows a high-scoring match between the query and subject sequences, with alignment scores and E-values provided for the top hit.





NCBI RefSeq Database

- *Goal:* Provide a single reference sequence for each molecule of the central dogma (DNA, mRNA, and protein)
- Distinguishing features
 - Non-redundancy
 - Updates to reflect the current knowledge of sequence data and biology
 - Includes biological attributes of the gene, gene transcript, or protein
 - Encompasses a wide taxonomic range, with primary focus on mammalian and human species
 - Ongoing updates and curation (both automated and manual review), with review status indicated on each record

Pruitt et al., *Nucleic Acids Res.* 42: D756-D763, 2014



RefSeq Accession Number Prefixes

From curation of GenBank entries:

- NT_** Genomic contigs
- NM_** mRNAs
- NP_** Proteins
- NR_** Non-coding transcripts

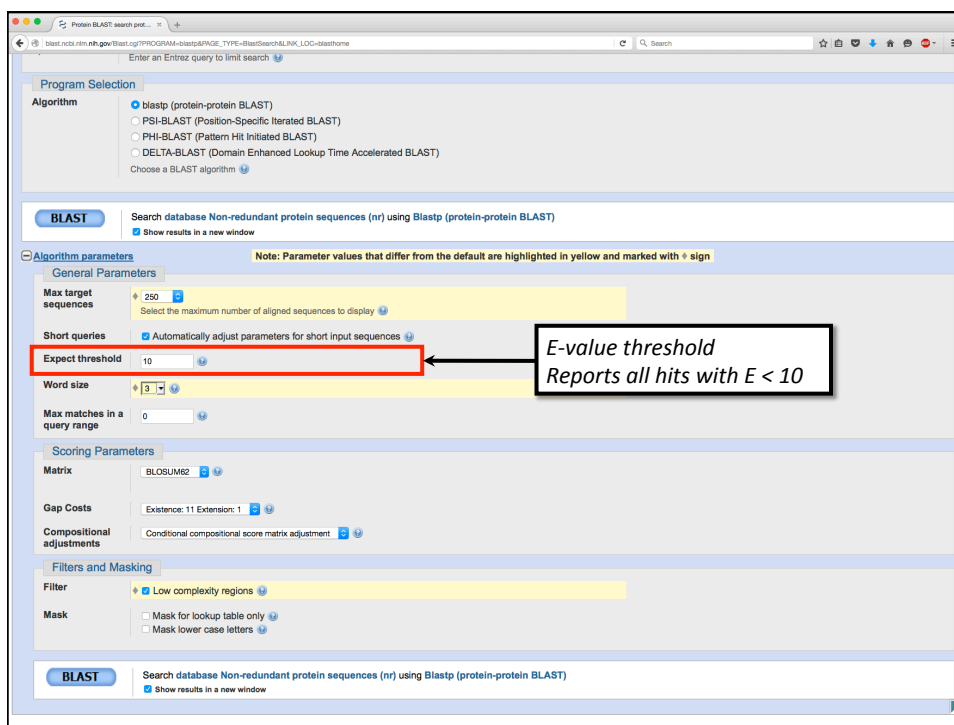
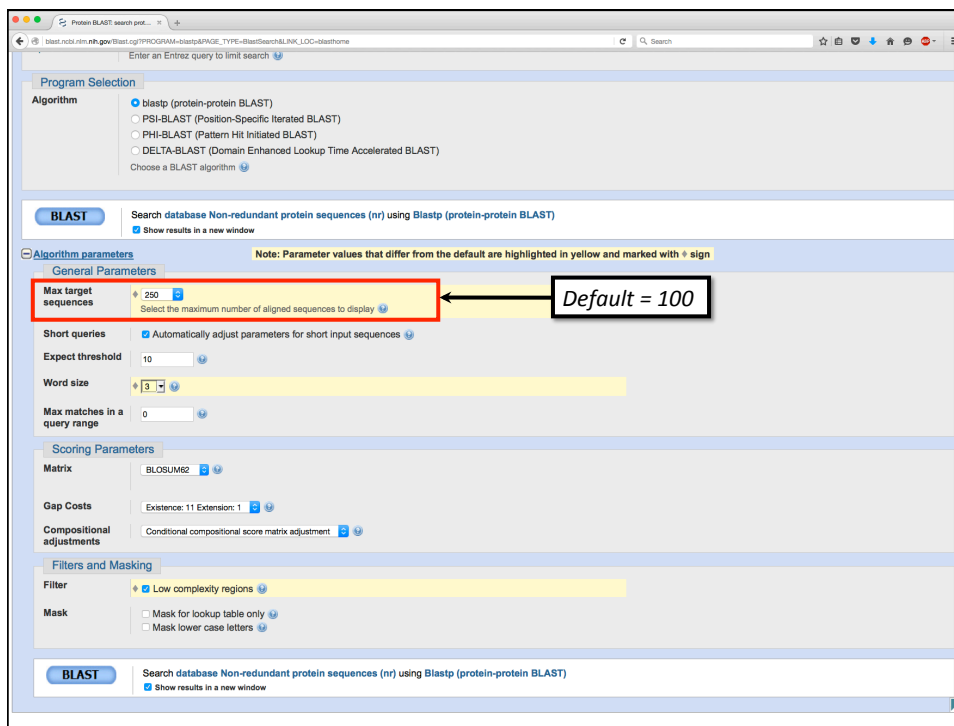
From genome annotation:

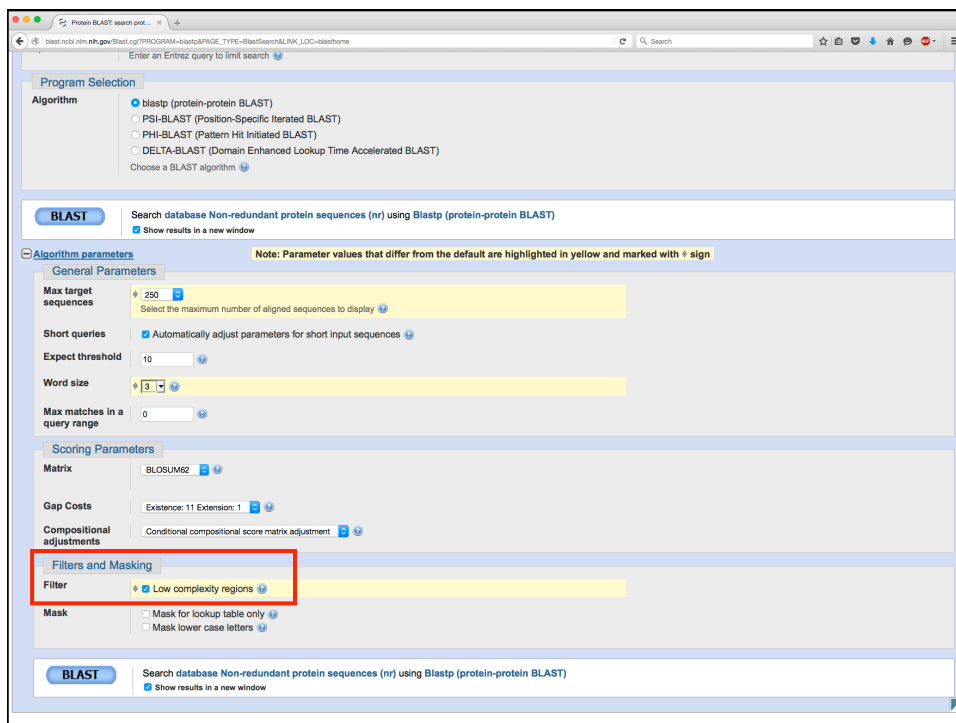
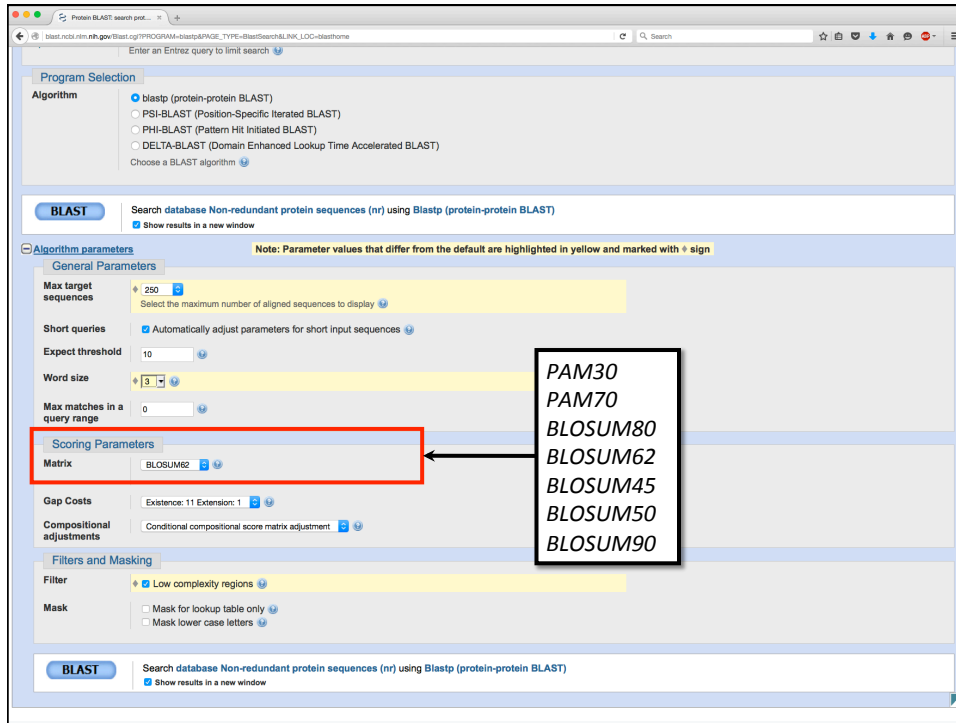
- XM_** Model mRNA
- XP_** Model proteins

Complete list of molecule types in Chapter 18 of the NCBI Handbook
<http://ncbi.nlm.nih.gov/books/NBK21091>

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

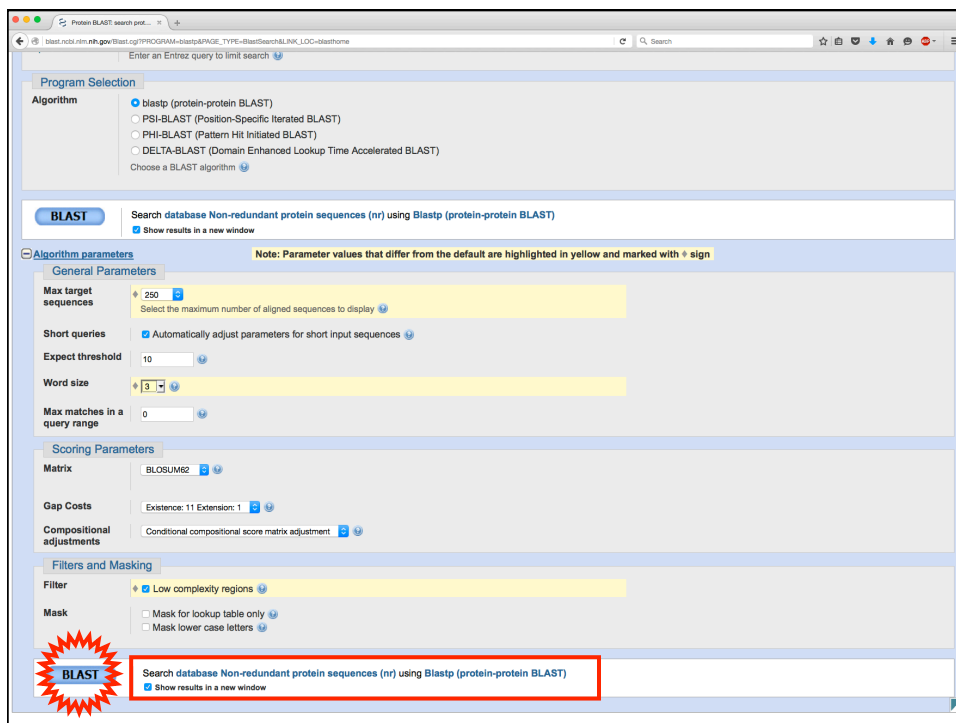
The screenshot shows the NCBI BLAST web interface. The main heading is "Standard Protein BLAST". Under "Choose Search Set", the "Database" is set to "Non-redundant protein sequences (nr)". The "Organism" field is highlighted with a callout box containing the text "Limit by organism or taxonomic group". Other options include "Exclude" (Models (XM/XP) and Uncultured/environmental sample sequences) and "Entrez Query". Under "Program Selection", "blastp (protein-protein BLAST)" is selected. At the bottom, there is a "BLAST" button and a link to "Algorithm parameters".

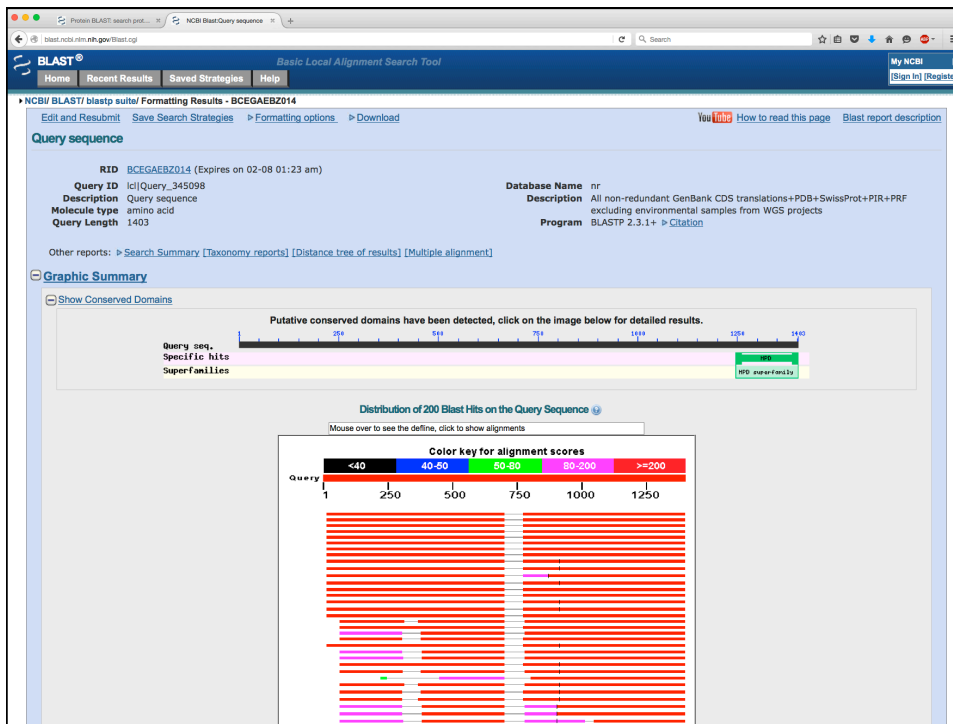




Low-Complexity Regions

- Defined as regions of 'biased composition'
 - Homopolymeric runs
 - Short-period repeats
 - Subtle over-representation of several residues
- 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)





Sequences producing significant alignments:

Select: All None Selected: 0

Alignments

Description	Max score	Total score	Query cover	E value	Ident	Accession
<input type="checkbox"/> prospero, isoform M [Drosophila melanogaster]	994	1938	93%	0.0	100%	NP_001247046.1
<input type="checkbox"/> prospero, isoform J [Drosophila melanogaster]	993	1936	93%	0.0	100%	NP_624317.4
<input type="checkbox"/> prospero [Drosophila melanogaster]	993	1932	93%	0.0	100%	BAAD1464.1
<input type="checkbox"/> homeodomain transcription factor Prospero [Drosophila melanogaster]	990	1821	93%	0.0	100%	AAF0703.1
<input type="checkbox"/> uncharacterized protein Dere_GG18089, isoform A [Drosophila erecta]	989	1885	93%	0.0	99%	XP_001980573.2
<input type="checkbox"/> Pros protein [Drosophila melanogaster]	982	1811	93%	0.0	97%	AAA28841.1
<input type="checkbox"/> prospero, isoform H [Drosophila melanogaster]	944	1862	93%	0.0	100%	NP_001247044.1
<input type="checkbox"/> prospero, isoform L [Drosophila melanogaster]	943	1858	93%	0.0	100%	NP_788636.3
<input type="checkbox"/> prospero, isoform I [Drosophila melanogaster]	942	1864	93%	0.0	100%	NP_001247045.1
<input type="checkbox"/> prospero, isoform K [Drosophila melanogaster]	942	1863	93%	0.0	100%	NP_731565.4
<input type="checkbox"/> GM23939 [Drosophila sechellia]	935	1987	93%	0.0	98%	XP_002031631.1
<input type="checkbox"/> LOW QUALITY PROTEIN: prospero [Drosophila simulans]	932	1827	93%	0.0	98%	KMZ04286.1
<input type="checkbox"/> uncharacterized protein Dere_GG18089, isoform B [Drosophila erecta]	915	1810	93%	0.0	95%	XP_015010089.1
<input type="checkbox"/> uncharacterized protein Dana_GF16857, isoform A [Drosophila ananassae]	904	1673	93%	0.0	92%	XP_001954214.2
<input type="checkbox"/> uncharacterized protein Dyak_GE26090 [Drosophila yakuba]	903	1816	93%	0.0	96%	XP_002097201.2
<input type="checkbox"/> uncharacterized protein Dere_GG18089, isoform C [Drosophila erecta]	894	1814	93%	0.0	97%	XP_015010070.1
<input type="checkbox"/> uncharacterized protein Dana_GF16857, isoform C [Drosophila ananassae]	855	1623	93%	0.0	90%	XP_014786172.1
<input type="checkbox"/> uncharacterized protein Owl_GK11290, isoform A [Drosophila willistonii]	845	1532	85%	0.0	83%	XP_002089989.2
<input type="checkbox"/> uncharacterized protein Dose_GA14403, isoform I [Drosophila pseudoobscura pseudoobscura]	825	1456	90%	0.0	82%	XP_001359985.4
<input type="checkbox"/> GH21437 [Drosophila grimshawi]	809	1374	84%	0.0	80%	XP_001994360.1
<input type="checkbox"/> uncharacterized protein Dmsl_GI22895, isoform B [Drosophila mojavensis]	799	1386	84%	0.0	78%	XP_002000130.2
<input type="checkbox"/> uncharacterized protein Dana_GF16857, isoform B [Drosophila ananassae]	787	1627	93%	0.0	83%	XP_014786171.1
<input type="checkbox"/> PREDICTED: homeobox protein prospero isoform X3 [Carattis capitata]	692	1111	84%	0.0	66%	XP_004529243.2
<input type="checkbox"/> PREDICTED: homeobox protein prospero [Bactrocera oleae]	690	1115	84%	0.0	70%	XP_014096508.1
<input type="checkbox"/> uncharacterized protein Dose_GA14403, isoform D [Drosophila pseudoobscura pseudoobscura]	612	14				
<input type="checkbox"/> pros [Drosophila busckii]	611	14				
<input type="checkbox"/> AAEL002769-PA [Aedes aegypti]	571	770	62%	8e-179	59%	XP_001655942.1
<input type="checkbox"/> uncharacterized protein Owl_GK11290, isoform B [Drosophila willistonii]	571	1501	85%	2e-171	77%	XP_015032827.1

0.0 means $\leq 10^{-1000}$

$8e-179 = 8 \times 10^{-179}$

<input type="checkbox"/> Prospero homeobox protein 1 [Chlamydomonas reinhardtii]	226	270	19%	6e-58	62%	KFP45850.1
<input type="checkbox"/> Prospero homeobox protein 1 [Cuculus canorus]	226	270	19%	6e-58	62%	KFO75119.1
<input type="checkbox"/> PREDICTED: prospero homeobox protein 1-like [Poecilia formosa]	228	228	12%	6e-58	57%	XP_007567659.1
<input type="checkbox"/> Prospero homeobox protein 1 [Pterodroma gutturalis]	225	269	19%	6e-58	63%	KFV13067.1
<input type="checkbox"/> homeobox protein prospero/brx-1 [Culex quinquefasciatus]	209	209	9%	6e-58	76%	XP_001845683.1
<input type="checkbox"/> PREDICTED: prospero homeobox protein 1 [Octodon degus]	226	270	19%	7e-58	63%	XP_004626924.1
<input type="checkbox"/> Prospero homeobox protein 1 [Charadrius vociferus]	226	270	19%	7e-58	62%	KGL88768.1
<input type="checkbox"/> PREDICTED: prospero homeobox protein 1 isoform X2 [Chinchilla lanigera]	226	270	19%	8e-58	63%	XP_005374780.1
<input type="checkbox"/> PREDICTED: prospero homeobox protein 1 isoform X1 [Fukuyma diamantensis]	226	270	19%	8e-58	63%	XP_010640836.1
<input type="checkbox"/> PREDICTED: prospero homeobox protein 1 isoform X2 [Cavia porcellus]	226	270	19%	8e-58	63%	XP_003474644.1
<input type="checkbox"/> PREDICTED: prospero homeobox protein 1 isoform X1 [Saimiri boliviensis boliviensis]	226	270	19%	8e-58	63%	XP_010339250.1
<input type="checkbox"/> PREDICTED: prospero homeobox protein 1 [Peromyscus maniculatus bairdii]	226	270	19%	8e-58	63%	XP_006972145.1
<input type="checkbox"/> PREDICTED: prospero homeobox protein 1 [Chaetura pelagica]	225	270	19%	8e-58	63%	XP_009693032.1
<input type="checkbox"/> PREDICTED: prospero homeobox protein 1 isoform X2 [Callithrix jacchus]	225	270	19%	8e-58	63%	
<input type="checkbox"/> PREDICTED: prospero homeobox protein 1 isoform X1 [Heteroscephalus glaber]	225	270	19%	8e-58	63%	
<input type="checkbox"/> PREDICTED: prospero homeobox protein 1 isoform X2 [Osteimur garnettii]	225	269	19%	8e-58	63%	
<input type="checkbox"/> PREDICTED: prospero homeobox protein 1 [Cuculus canorus]	225	270	19%	8e-58	63%	
<input type="checkbox"/> PREDICTED: prospero homeobox protein 1 isoform X2 [Equus asinus]	225	269	19%	8e-58	63%	XP_014680416.1
<input type="checkbox"/> PREDICTED: prospero homeobox protein 1 isoform X2 [Propithecus coquerellei]	225	270	19%	8e-58	63%	XP_012483821.1
<input type="checkbox"/> PREDICTED: prospero homeobox protein 1 [Colobus angolensis palliatus]	225	269	19%	8e-58	63%	XP_011784800.1
<input type="checkbox"/> PREDICTED: prospero homeobox protein 1 [Mandrillus leucophaeus]	225	270	19%	8e-58	63%	XP_011851547.1
<input type="checkbox"/> prospero homeobox protein 1 [Homo sapiens]	225	270	19%	8e-58	63%	NP_002754.2
<input type="checkbox"/> PREDICTED: LOW QUALITY PROTEIN: prospero homeobox protein 1-like [Colius striatus]	225	270	19%	8e-58	63%	XP_010205580.1
<input type="checkbox"/> PREDICTED: prospero homeobox protein 1 isoform X2 [Columba livia]	225	271	19%	8e-58	63%	XP_005502819.1
<input type="checkbox"/> PREDICTED: prospero homeobox protein 1 [Falco cherrug]	225	270	19%	9e-58	63%	XP_005441959.1
<input type="checkbox"/> PREDICTED: prospero homeobox protein 1 [Marmota marmota marmota]	225	270	19%	9e-58	63%	XP_015339134.1
<input type="checkbox"/> hypothetical protein EGM_01399 [Macaca fascicularis]	225	270	19%	9e-58	63%	EHH50546.1
<input type="checkbox"/> PREDICTED: prospero homeobox protein 1 isoform X2 [Nannospalax galii]	225	270	19%	9e-58	63%	XP_008623079.1
<input type="checkbox"/> PREDICTED: prospero homeobox protein 1 [Ochotona princeps]	225	269	19%	9e-58	63%	XP_004578703.1

Accept (for now)

Reject above desired threshold ($E \leq 10^{-3}$)

prospero, isoform L [Drosophila melanogaster]
 Sequence ID: [ref|NP_788636.3](#) Length: 1374 Number of Matches: 2
[See 2 more title\(s\)](#)

Range 1: 17 to 704 GenPept Graphics

Score	Expect	Method	Identities	Positives	Gaps
943 bits(2437)	0.0	Compositional matrix adjust.	688/688(100%)	688/688(100%)	0/688(0%)

Query 17 LFQPQSVSTANSSSSNNNSSTPAALATHSPTNSPVGASASSLLTAAFGNLFGGSSA 76
 Sbjct 17 LFQPQSVSTANSSSSNNNSSTPAALATHSPTNSPVGASASSLLTAAFGNLFGGSSA 76

Query 77 KMLNELFGROMKQADATSLPQSLDNAMLAAMETATSSELLIGSLNSTSKLLOQOHNN 136
 Sbjct 77 KMLNELFGROMKQADATSLPQSLDNAMLAAMETATSSELLIGSLNSTSKLLOQOHNN 136

Query 137 NSIAPANSTPMSNGTNASISFGSAHSSSHHGQVSPKGSRRVSACDRSLEAAADVAG 196
 Sbjct 137 NSIAPANSTPMSNGTNASISFGSAHSSSHHGQVSPKGSRRVSACDRSLEAAADVAG 196

Query 197 SPPRAASVSLGGASGEHQSQQLQHLVAHMLRNILQCKKELMLQDLQELRTAMQQQQ 256
 Sbjct 197 SPPRAASVSLGGASGEHQSQQLQHLVAHMLRNILQCKKELMLQDLQELRTAMQQQQ 256

Query 257 qqlqekelHSKlnnnnnnlaatannnnntMESINLIDSEMDIKIKSEPQTAPOPO 316
 Sbjct 257 qqlqekelHSKlnnnnnnlaatannnnntMESINLIDSEMDIKIKSEPQTAPOPO 316

Query 317 QspghshssrsgsgsgshsmasdgslrtksdsldaHgqddagdeedaPTQRSP 376
 Sbjct 317 QSPHSGSHSSRSGSGSGSHSMASDGLRKRSSDLSHGAQDDAQDEEADPTQRSP 376

Query 377 RAPEEPQLPTKESVDMLDEVLLGLHRSQSDMSLASPSHSDMLDKDDVLEDDDD 436
 Sbjct 377 RAPEEPQLPTKESVDMLDEVLLGLHRSQSDMSLASPSHSDMLDKDDVLEDDDD 436

Query 437 DCVEQKTSGGCLKPKMDLKRARVENIVSGMRCSPSSGLAAGQLVNGCKRKLQYPO 496
 Sbjct 437 DCVEQKTSGGCLKPKMDLKRARVENIVSGMRCSPSSGLAAGQLVNGCKRKLQYPO 496

Query 497 QHAMERYVAAAAGLNFGLNLSMMLDQEDSESNLESPQIQKRVEKNALKSQRSMQEO 556
 Sbjct 497 QHAMERYVAAAAGLNFGLNLSMMLDQEDSESNLESPQIQKRVEKNALKSQRSMQEO 556

Query 557 LAEMQQRVQLCSRMEQSECEGLDQDQVQEQEPDNGSDHIELSPSPITLGDGVDP 616
 Sbjct 557 LAEMQQRVQLCSRMEQSECEGLDQDQVQEQEPDNGSDHIELSPSPITLGDGVDP 616

Query 617 NHKEETGQERogsspsplkpktaLgESSDGSANLMSQMSKMSGKLNHPLVGVGHP 676
 Sbjct 617 NHKEETGQERogsspsplkpktaLgESSDGSANLMSQMSKMSGKLNHPLVGVGHP 676

Query 677 ALPQGFPLLQHMGMDSHAAMVQFFF 704
 Sbjct 677 ALPQGFPLLQHMGMDSHAAMVQFFF 704

Range 2: 777 to 1374 GenPept Graphics

Score	Expect	Method	Identities	Positives	Gaps
915 bits(2365)	0.0	Compositional matrix adjust.	598/627(95%)	598/627(95%)	29/627(4%)

Query 777 HVATAAPRQMHHFAPARLPTRMGGAAGHTALKSELSEKFQMLRANNSSMRMSGTDLE 836
 Sbjct 777 HVATAAPRQMHHFAPARLPTRMGGAAGHTALKSELSEKFQMLRANNSSMRMSGTDLE 836

Query 837 GLADVLKSEITTSLSALVDITVTRFHQRRLFSQADSVTAAEQLNKOLLASQILDRK 896
 Sbjct 837 GLADVLKSEITTSLSALVDITVTRFHQRRLFSQADSVTAAEQLNKOLLASQILDRK 896

Query 897 SFRKTVADEQNGPTTATQSAAMFQAKFTPGQMNVAALYNSMTGPFCLPDDQQQQ 956
 Sbjct 897 SFRKTVADEQNGPTTATQSAAMFQAKFTPGQMNVAALYNSMTGPFCLPDDQQQQ 956

Query 957 qt4gqgagqgqgagqgqgqLQNEALSIVTTPKKRKHVDTTRITRIVSRILAQDG 1016
 Sbjct 957 QTAQQQSAQQQQSQQTQQQLQNEALSIVTTPKKRKHVDTTRITRIVSRILAQDG 1016

Query 1017 VVPTGPPSTPQQQQQQQQQQQQQQQQQQASNGGNSNATPAQSPTRS SGGAAHPQP 1076
 Sbjct 1017 VVPTGPPSTPQQQQQQQQQQQQQQQQQQASNGGNSNATPAQSPTRS SGGAAHPQP 1076

Query 1077 pppppmmpVSLPTSVAIPNPSLHESKVSFSPYFFFNHAAAGQATAAQLHQHQHHPH 1136
 Sbjct 1077 pppppmmpVSLPTSVAIPNPSLHESKVSFSPYFFFNHAAAGQATAAQLHQHQHHPH 1136

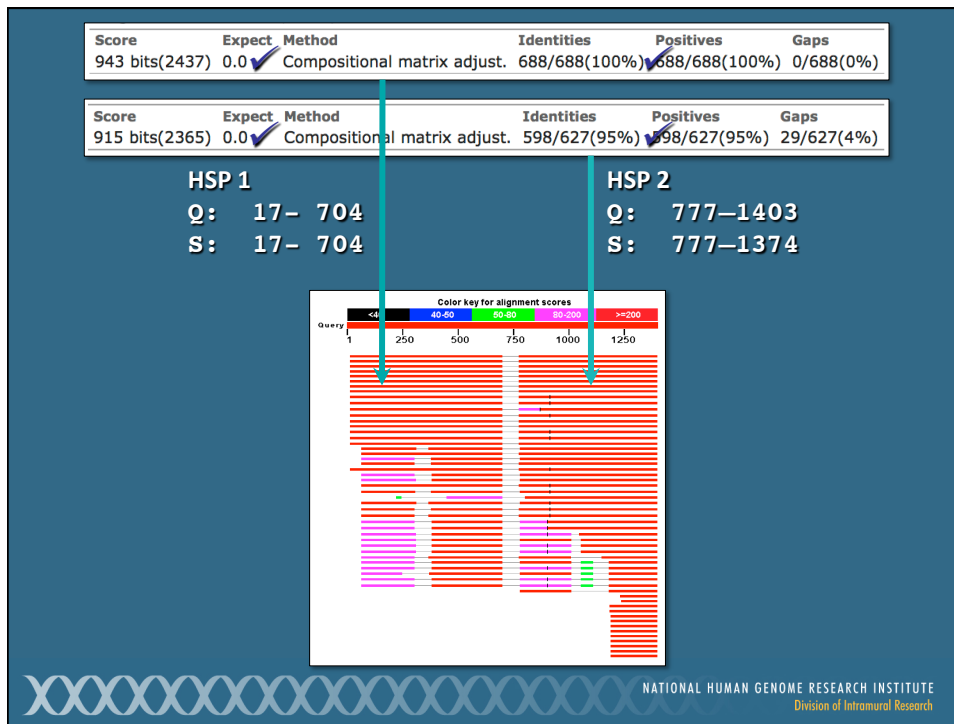
Query 1137 hqamllssppgslALMDSRDEppllpppsmlhpallaaahggsDYKTCRLRMDAQ 1196
 Sbjct 1137 HQMQLSSFPGLGALMDSRDEppllpppsmlhpallaaahggsDYKTCRLRMDAQ 1196

Query 1197 DRQSECNADQFDGMAPTISFYKQMLKTEHQESLMAKICSEPTLHSSLTTPMLRKA 1256
 Sbjct 1197 DRQSECNADQFDGMAPT-----SSTLTTPMLRKA 1227

Query 1257 KLMFVWVYRPSAVLKYFFDIKFNKNNTAQLKVMFSNFRFYIOMEKARQAVTEGIK 1316
 Sbjct 1257 KLMFVWVYRPSAVLKYFFDIKFNKNNTAQLKVMFSNFRFYIOMEKARQAVTEGIK 1287

Query 1317 TPDDLIAQDSELYRVLNLYRNNHIEVPQNFVVESTLREFFRAIQGGKDTQSWKK 1376
 Sbjct 1288 TPDDLIAQDSELYRVLNLYRNNHIEVPQNFVVESTLREFFRAIQGGKDTQSWKK 1347

Query 1377 SIYKISRMDDPVEYFKSPNLEOLE 1403
 Sbjct 1348 SIYKISRMDDPVEYFKSPNLEOLE 1374



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!

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



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

BLAST finds regions of similarity between biological sequences. [more...](#)

BLAST Assembled Genomes

Find Genomic BLAST pages:

Enter organism name or id--completions will be suggested

<input type="checkbox"/> Human	<input type="checkbox"/> Rabbit	<input type="checkbox"/> Zebrafish
<input type="checkbox"/> Mouse	<input type="checkbox"/> Chimp	<input type="checkbox"/> Clawed frog
<input type="checkbox"/> Rat	<input type="checkbox"/> Guinea pig	<input type="checkbox"/> Arabidopsis
<input type="checkbox"/> Cow	<input type="checkbox"/> Fruit fly	<input type="checkbox"/> Rice
<input type="checkbox"/> Pig	<input type="checkbox"/> Honey bee	<input type="checkbox"/> Yeast
<input type="checkbox"/> Dog	<input type="checkbox"/> Chicken	<input type="checkbox"/> Microbes

Basic BLAST

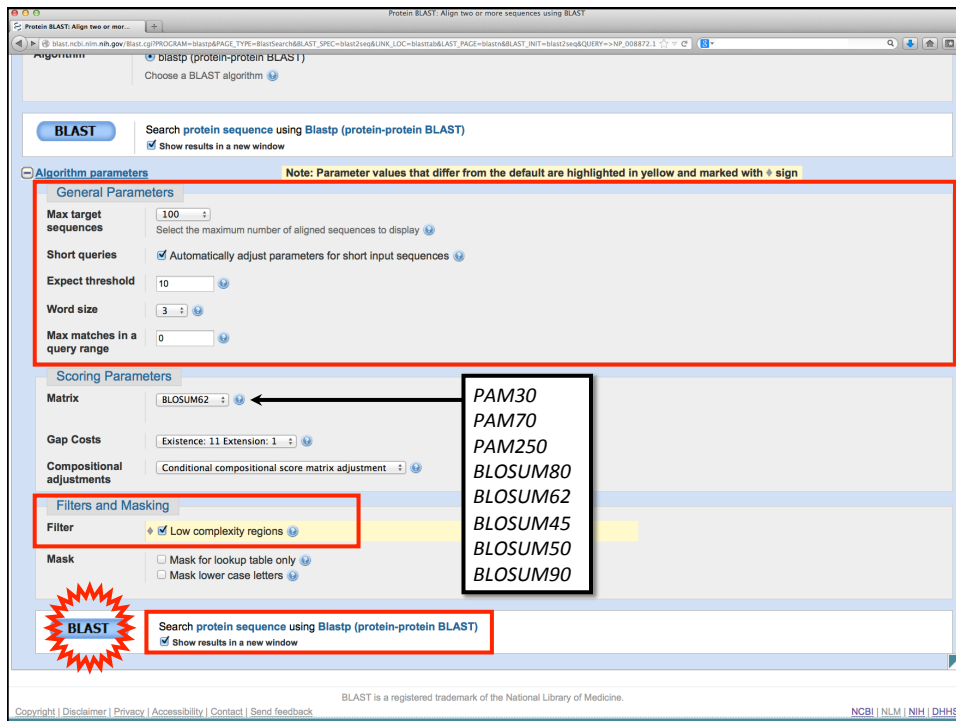
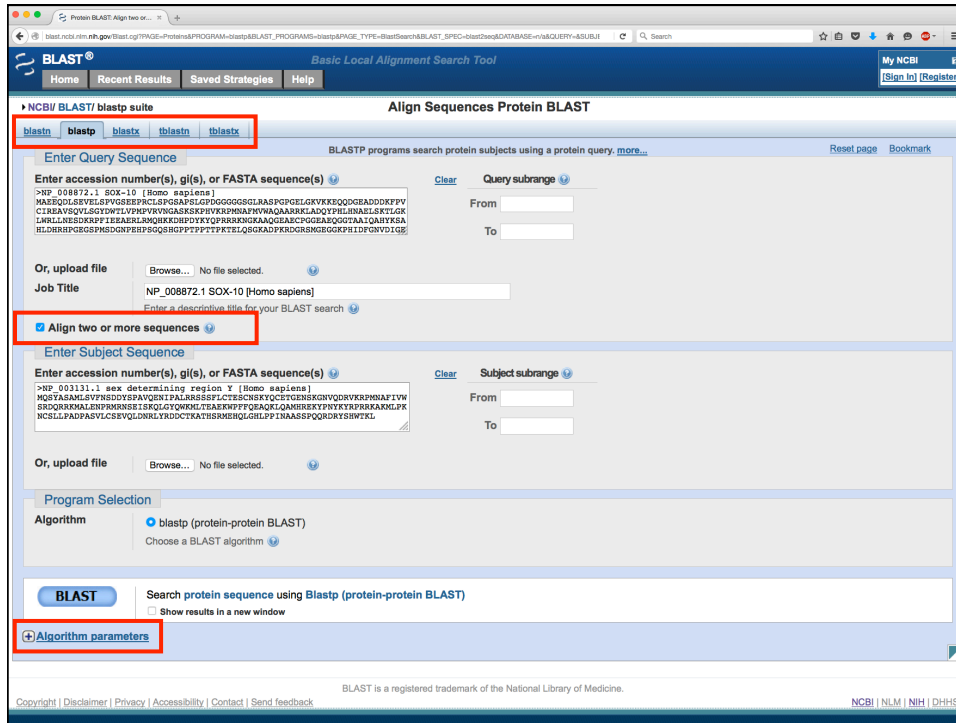
Choose a BLAST program to run.

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, pti-blast, delta-blast
tblastx	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.)

- Get faster protein results with a graphical view using [SmartBLAST](#)
- Make specific primers with [Primer-BLAST](#)
- Cluster multiple sequences together with their database neighbors using [MOLE-BLAST](#)
- Find [conserved domains](#) in your sequence (cds)
- Find sequences with similar [conserved domain architecture](#) (cdart)
- Search sequences that have [gene expression profiles](#) (GEO)
- Search [immunoglobulins and T cell receptor sequences](#) (IqBLAST)
- Screen sequences for [vector contamination](#) (vecscreen)
- Align two (or more) sequences using BLAST (bl2seq)**
- Search [protein or nucleotide targets](#) in PubChem BioAssay
- Search [SRA by experiment](#)
- Constraint Based Protein [Multiple Alignment Tool](#)
- Needleman-Wunsch [Global Sequence Alignment Tool](#)
- Search [RefSeqGene](#)



NCBI BLAST/ blastp suite-2sequences/ Formatting Results - BCJA4YBV114

Blast 2 sequences

NP_008872.1 SOX-10 [Homo sapiens]

RID: BCJA4YBV114 (Expires on 02-08 02:28 am)

Query ID: lcl|Query_213409
Description: NP_008872.1 SOX-10 [Homo sapiens]
Molecule type: amino acid
Query Length: 466

Subject ID: lcl|Query_213411
Description: NP_003131.1 sex determining region Y [Homo sapiens]
Molecule type: amino acid
Subject Length: 204
Program: BLASTP 2.3.1+

Other reports: [Search Summary](#) [Multiple alignment](#)

Graphic Summary

Distribution of 2 Blast Hits on the Query Sequence

Color key for alignment scores

Dot Matrix View

Descriptions

Sequences producing significant alignments:

Select: All None Selected:0

Description	Max score	Total score	Query cover	E value	Ident	Accession
<input type="checkbox"/> NP_003131.1 sex determining region Y [Homo sapiens]	94.0	109	19%	1e-26	46%	Query_213411

Dot Matrix View

Descriptions

Sequences producing significant alignments:

Select: All None Selected:0

Description	Max score	Total score	Query cover	E value	Ident	Accession
<input type="checkbox"/> NP_003131.1 sex determining region Y [Homo sapiens]	94.0	109	19%	1e-26	46%	Query_213411

Alignments

Download Graphics Sort by: E value

NP_003131.1 sex determining region Y [Homo sapiens]
 Sequence ID: lcl|Query_213411 Length: 204 Number of Matches: 2

Range 1: 91 to 134

Score	Expect	Method	Identities	Positives	Gaps
94.0 bits(232)	1e-26	Compositional matrix adjust.	39/84(46%)	62/84(73%)	0/84(0%)

Query 95 NGASKSPKPKRPMNAPFYWAQAARRKLADYTPHLENABLSEKTLGKLRLLNESDKRPFT 154
 N + VRSPNKA+VH+ SRPA + P + H+S+K LG H+L B+K PF
 Sbjct 51 NSKGNVDQRYKRPNAP IVWSRQRRKMALENFRMRNSEISQLYQWRMLTEAEKWPF 110

Query 155 EEAERLRMQHKKDHPDYKYQPRRR 178
 *EA++L+ H++ +P+RY+PRR+
 Sbjct 111 QEAQELQAHREKYPWYKTRPRK 134

Range 2: 95 to 101

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

Query 82 GYDWTLY 88
 GE M ++
 Sbjct 95 GYQWRML 101

**Nucleotide Similarity Searching:
MegaBLAST, BLASTN, and BLAT**

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

The slide features a blue background with a white DNA double helix graphic at the bottom. The title is centered in a large, bold, white font.

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

BLAST finds regions of similarity between biological sequences. [more...](#)

BLAST Assembled Genomes

Find Genomic BLAST pages:

Enter organism name or id--completions will be suggested

- Human
- Mouse
- Rat
- Cow
- Pig
- Dog
- Rabbit
- Chimp
- Guinea pig
- Fruit fly
- Honey bee
- Chicken
- Zebrafish
- Clawed frog
- Arabidopsis
- Rice
- Yeast
- Microbes

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, pti-blast, delta-blast
- [tblastx](#) 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.)

- Get faster protein results with a graphical view using [SmartBLAST](#)
- Make specific primers with [Primer-BLAST](#)
- Cluster multiple sequences together with their database neighbors using [MOLE-BLAST](#)
- Find [conserved domains](#) in your sequence (cds)
- Find sequences with similar [conserved domain architecture](#) (cdart)
- Search sequences that have [gene expression profiles](#) (GEO)
- Search [immunoglobulins and T cell receptor sequences](#) (IgBLAST)
- Screen sequences for [vector contamination](#) (vecscreen)
- [Align](#) two (or more) sequences using BLAST (bl2seq)
- Search [protein or nucleotide targets](#) in PubChem BioAssay
- Search [SRA by experiment](#)
- Constraint Based Protein [Multiple Alignment Tool](#)
- Needleman-Wunsch [Global Sequence Alignment Tool](#)
- Search [RefSeqGene](#)

Your Recent Results [New!](#)

[All Recent results...](#)

News

[Searching Whole Genome Shotgun sequences](#)

It is now much easier to search WGS (Whole Genome Shotgun) with stand-alone BLAST on your own computer.

Wed, 20 Jan 2016 10:00:00 EST

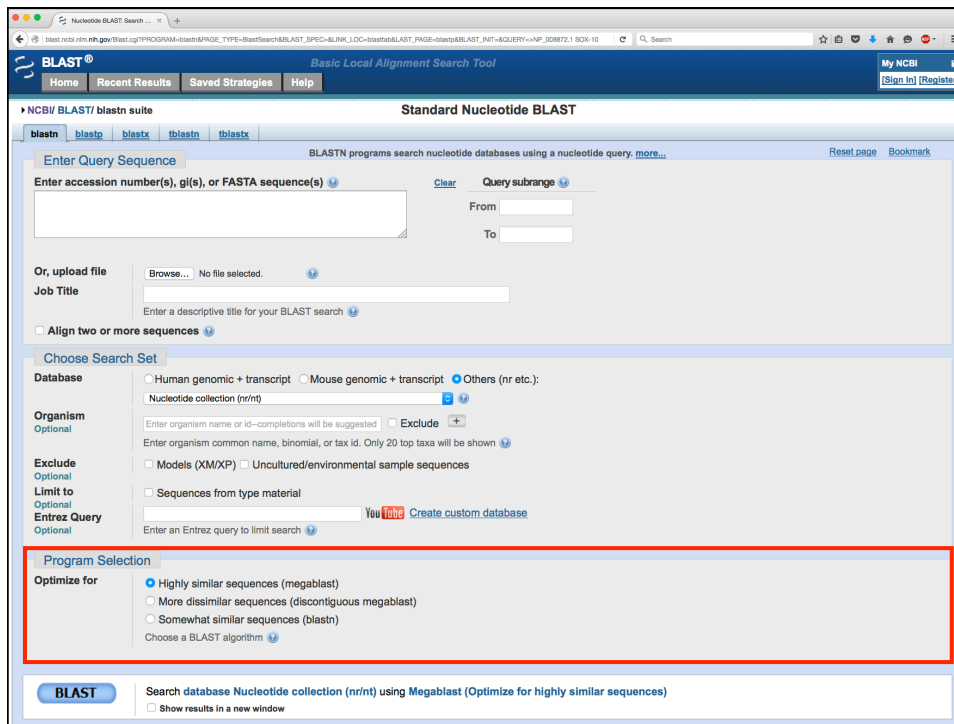
[More BLAST news...](#)

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...](#)



Nucleotide-Based BLAST Algorithms

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

BLAT

- “BLAST-Like Alignment Tool”
- Designed to rapidly align longer nucleotide sequences ($L \geq 40$) having $\geq 95\%$ 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 than BLAST for mRNA/DNA searches
- May miss divergent or shorter sequence alignments
- Can be used on protein sequences, but BLASTP is more efficient



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 (or identical) sequences
 - Alignment of mRNA sequences onto a genome assembly
 - Identification of gene family members
 - Cross-species alignment to identify putative homologs
- To display a specific sequence as a separate track within the UCSC Genome Browser



UCSC Genome Bioinformatics

<http://genome.ucsc.edu>

Genomes Genome Browser Tools Mirrors Downloads My Data Help About Us

About the UCSC Genome Bioinformatics Site

Welcome to the UCSC Genome Browser website. This site contains the reference sequence and working draft assemblies for a large collection of genomes. It also provides portals to [ENCODE](#) data at UCSC (2003 to 2012) and to the [Neanderthal](#) project. Download or purchase the Genome Browser source code, or the Genome Browser in a Box ([GBIB](#)) at our [online store](#).

We encourage you to explore these sequences with our tools. The [Genome Browser](#) zooms and scrolls over chromosomes, showing the work of annotators worldwide. The [Gene Sorter](#) shows expression, homology and other information on groups of genes that can be related in many ways. [Blat](#) quickly maps your sequence to the genome. The [Table Browser](#) provides convenient access to the underlying database. [VisiGene](#) lets you browse through a large collection of *in situ* mouse and frog images to examine expression patterns. [Genome Graphs](#) allows you to upload and display genome-wide data sets.

The UCSC Genome Browser is developed and maintained by the Genome Bioinformatics Group, a cross-departmental team within the [UC Santa Cruz Genomics Institute](#) at the University of California Santa Cruz ([UCSC](#)). If you have feedback or questions concerning the tools or data on this website, feel free to contact us on our [public mailing list](#).

The Genome Browser project team relies on public funding to support our work. Donations are welcome -- we have many more ideas than our funding supports! If you have ideas, drop a comment in our [suggestion box](#). [DONATE NOW](#)

News [News Archives](#)

To receive announcements of new genome assembly releases, new software features, updates and training seminars by email, subscribe to the [genome-announce](#) mailing list. Please see our [blog](#) for posts about Genome Browser tools, features, projects and more.

20 Jan 2016 - dbSNP 142 Available for mm10

Data from dbSNP build 142 is now available for the most recent mouse assembly (mm10/GRCm38). As was the case for previous annotations based on dbSNP data, there are three tracks in this release. One is a track containing all mappings of reference SNPs to the mouse assembly, labeled "All SNPs (142)". The other two tracks are subsets of this track and show different interesting and easily defined subsets of dbSNP:

- Common SNPs (142): uniquely mapped variants that appear in at least 1% of the population
- Mult. SNPs (142): variants that have been mapped to more than one genomic location

By default, only the Common SNPs (142) are visible. The other tracks can be made visible using the track controls. These three SNPs (142) tracks can be found on the Mouse Dec. 2011 (mm10/GRCm38) browser in the "Variation and Repeats" group.

Thank you to the [dbSNP](#) group at NCBI for making these data publicly available. The tracks were produced at UCSC by Brian Raney, Angie Hinrichs and Matthew Speir.

08 January 2016 - dbSNP 144 Available for hg19 and hg38

We are pleased to announce the release of four tracks derived from NCBI [dbSNP](#) Build 144 data, available on the two most recent human assemblies GRCCh37/hg19 and GRCCh38/hg38.

Rhesus BLAT Search

BLAT Search Genome

Genome: Assembly: Query type: Sort output: Output type:

```
>CB312814 NICHDRh_Ov1 Macaca mulatta cDNA clone
GGGGGTGGAGCTGCCAGTAAAGCAAAGAGCAAGGAGCAGGCTGTTGGAAGGGGTTGTGACAGCCCC
AGCAATGTGGAGAGTCTGGGGCTTGCCCTGGCTCTCTGCTCTTCATCGGAGGAAACAGAGAGCCAG
GACAAAGCTCCTTGTGTAAGCAACCCAGCCTGGAGCATAAGAGATCAAGATCCAAATGCTAGACTCA
ATGGTTCAGTACGTGTGCTGCTCTTCTTCAAGCCAGCTGATAGCTTGCAATACGAGCAATCTAAATT
GGAAAGACTGCGGTAAGCTGGGAAAGAGGATATCTTAAATATCTTATATGGTGTATCATCA
GGATCTCTTCTCGATTAATAACACACATCTTTAGAAAAAGGTTTCAGAGCATATCTCTATATTC
CCAGAAGAAAACCAACCGATCTCTGGACTCTTTAATGGAAACCAAGAACCTCTCATATATGACGG
ATGTGGCTCTCTGAAAAACACCTGGTGGCCCTTTTCTTCCCACTTGGCGAATGGTAATAAAAAACC
CCTTAAATGGTTTTCCGGAAAAAAGTGGGAAATTTGCTCCCTCCCAATCTCAAAAAGAAAA
TTTTTGTAAAAAGGATCTTTTGGCACCGGGGAAAAAAAATTTGAAACCTCCCCACCCCTT
TTTCCCTTTGGGACTCTTCCCAATTCCGGGACATCCCCCT
```

Paste in a query sequence to find its location in the genome. Multiple sequences may be searched if separated by lines starting with '>' followed by the sequence name.

File Upload: Rather than pasting a sequence, you can choose to upload a text file containing the sequence.
 Upload sequence: No file selected.

Only DNA sequences of 25,000 or fewer bases and protein or translated sequence of 10000 or fewer letters will be processed. Up to 25 sequences can be submitted at the same time. The total limit for multiple sequence submissions is 50,000 bases or 25,000 letters.

For locating PCR primers, use [In-Silico PCR](#) for best results instead of BLAT.

I'm feeling lucky returns only the highest scoring alignment (direct path to genome browser)

Rhesus BLAT Results

BLAT Search Results

Go back to [chr6:43159698-43164683](#) on the Genome Browser.

ACTIONS	QUERY	SCORE	START	END	QSIZE	IDENTITY	CHRO	STRAND	START	END	SPAN
browser details	CB312814	380	1	418	677	96.2%	6	-	43159698	43161152	1455
browser details	CB312814	23	591	613	677	100.0%	4	-	148338464	148338466	23
browser details	CB312814	22	546	567	677	100.0%	12	-	39379930	39379951	22
browser details	CB312814	21	628	648	677	100.0%	16	+	20696166	20696186	21
browser details	CB312814	21	629	651	677	95.7%	1	+	134928210	134928232	23
browser details	CB312814	20	553	574	677	95.5%	11	-	4332856	4332877	22
browser details	CB312814	20	627	646	677	100.0%	1	-	187748214	187748233	20
browser details	CB312814	20	511	530	677	100.0%	1	-	90178654	90178673	20

[Missing a match?](#)

UCSC Genome Browser on Rhesus Oct. 2010 (BGI CR_1.0/rheMac3) Assembly

chr6:43,157,205-43,167,176 9,972 bp

Click on a feature for details. Click or drag in the base position track to zoom.

- 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

Rhesus BLAT Results

BLAT Search Results

Go back to [chr6:43159698-43164683](#) on the Genome Browser.

ACTIONS	QUERY	SCORE	START	END	QSIZE	IDENTITY	CHRO	STRAND	START	END	SPAN
browser details	CB312814	380	1	418	677	96.2%	6	-	43159698	43161152	1455
browser details	CB312814	23	591	613	677	100.0%	4	-	148338464	148339486	23
browser details	CB312814	22	546	567	677	100.0%	12	-	39379930	39379951	22
browser details	CB312814	21	628	648	677	100.0%	16	+	20696166	20696186	21
browser details	CB312814	21	629	651	677	95.7%	1	+	134928210	134928232	23
browser details	CB312814	20	553	574	677	95.5%	11	-	4332856	4332877	22
browser details	CB312814	20	627	646	677	100.0%	1	-	187748214	187748233	20
browser details	CB312814	20	511	530	677	100.0%	1	-	90178654	90178673	20

[Missing a match?](#)

Alignment of CB312814 and chr6:43159698-43161152

Click on links in the frame to the left to navigate through the alignment. Matching bases in cDNA and genomic sequences are colored blue and capitalized. Light blue bases mark the boundaries of gaps in either sequence (often splice sites).

[CB312814](#)
[Rhesus.chr6](#)
[block1](#)
[block2](#)
[together](#)

cDNA CB312814

```

AGCAATGTGG AGAAGTCTGG GCGTTGCCCT GCGTCTCTGT CTCTTCCAT 50
CGGGAGAAC AGAGAGCCAG GACCAAAGCT CCTTCTGTAA GCACCCCCA 100
GCGTGGAGCA TAAGAGATCA AGATCCAATG CTAGACTCCA ATGGTTCAGT 150
GACTGTGGTC GCTCTTCTTC AAGCCAGCTG ATACCTGTGC ATACTGCANg 200
CATCTAAATT GGAGAACTC CGAGTAAAC TGGAGAAAGA AGGATATTCF 250
AAATTTCTC ATATTGtgg TAATCATCAA GGGATCTCTT CTCGATTTAA 300
ATACACACAT CTTLAGAaaA AGGTTTCAG AGCATATTC TGTATATcA 350
CoAGAAGAAA ACCoAACCGA TGCTGGACT CTTTAAATGG AAoCAAGAA 400
GACCTCCTCA TATATGAagg atgtggcctt cctggaaaaa accctgggtg 450
gccttttcc tcccaactt tgggaatgg taaaaaaac cctttaaattg 500
gttttccggg aaaaaaaag tgggaatgg gtctctctcc aaactcAAA 550
aaagaaaaa tttttgtaa aagggatctt ttggggcacc ggggggaaaa 600
aaaaattga aaactctcc caccctctt tttccctctt tggggactcc 650
ttcccaaat cgggggacat cccctct
    
```

Genomic chr6 (reverse strand):

```

agtgaatta tgtctgcagg atttatagaa attcatagtt aggactgtga 43161203
agttaactat gaagaagagt gacaggtttt ctcttttaca ggaacgcccc 43161153
AGCAATGTGG AGAAGTCTGG GCGTTGCCCT GCGTCTCTGT CTCTTCCAT 43161103
CGGGAGAAC AGAGAGCCAG GACCAAAGCT CCTTCTGTAA GCACCCCCA 43161053
GCGTGGAGCA TAAGAGATCA AGATCCAATG CTAGACTCCA ATGGTTCAGT 43161003
GACTGTGGTC GCTCTTCTTC AAGCCAGCTG ATACCTGTGC ATACTGCANg 43160953
CATCTAgtA agacagtctt tctgtggett aaaaactctt aaagggaaag 43160903
ttattagata cacacatgca tatagacata aaagtgtaaa caataattta 43160853
agtcacactt ttgaaaaaac tatgtgtttg cacagaaact attagnaagag 43160803
agaaactagg atgatacaca ccaaaatgtt tacagtgggt ttccatagggt 43160753
tatggaaatt ttctctctt ttgtggacc tatattttat aacttctcac 43160703
taaatatgta ttacttgtgt aatgaaaatg ttctaaaatg tacttctgca 43160653
aatagaacag ttacttcaat acaagaagca gagaagactt ttgtcagagt 43160603
agaagaaca ctaggcttgc tatcaaaagt ttgggttat taaacaataa 43160553
aatatacaat atatttgtg agtatctcac cagatattga ttgggtgaa 43160503
ttctcttacc agtcaatcaa ttgtattctg geaggattgg ttccaggatc 43160453
ccctctcac accaaaatcc atggagctca aatcccttat ataaaatgac 43160403
atataaac tatgacata ttctacatg tttaaatcat ccttagatta 43160353
ottataaac ctaacagaat gtaaatgcta tgtaagtaa ttttatactg 43160303
tattgttag tgaataatga atgacatttt aaaaagtcta catgttcaat 43160253
    
```

