




NATIONAL HUMAN GENOME RESEARCH INSTITUTE *Division of Intramural Research*




*Current Topics in Genome Analysis 2014*  
*Week 2: Biological Sequence Analysis I*  
*Andy Baxevanis, Ph.D.*

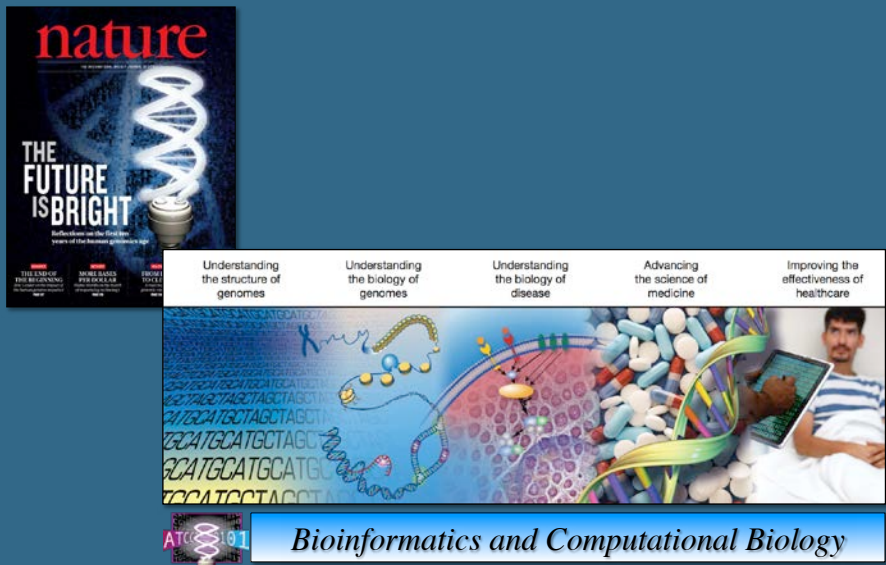
U.S. DEPARTMENT OF HEALTH AND HUMAN SERVICES | NATIONAL INSTITUTES OF HEALTH | genome.gov/DIR



*Current Topics in Genome Analysis 2014*  
*Andy Baxevanis, Ph.D.*  
*No Relevant Financial Relationships with  
Commercial Interests*

NATIONAL HUMAN GENOME RESEARCH INSTITUTE  
*Division of Intramural Research*






**nature**  
THE FUTURE IS BRIGHT  
Dedication on the 10th anniversary of the Human Genome Project

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

*Sequence Alignments:  
Determining Similarity and  
Deducing Homology*



NATIONAL HUMAN GENOME RESEARCH INSTITUTE  
Division of Intramural Research

## 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 indivisible<sup>8</sup>. 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
- 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



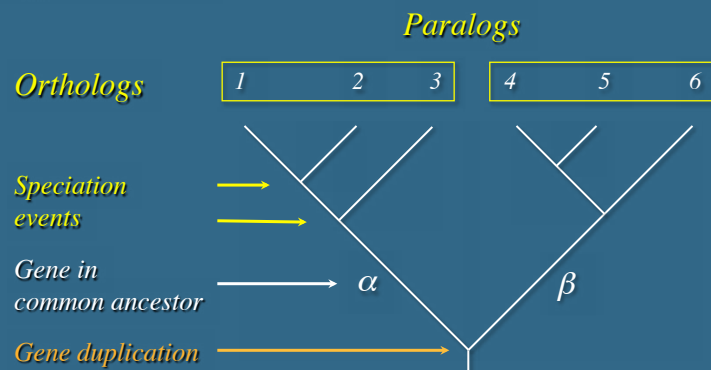
## Defining the Terms

*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”

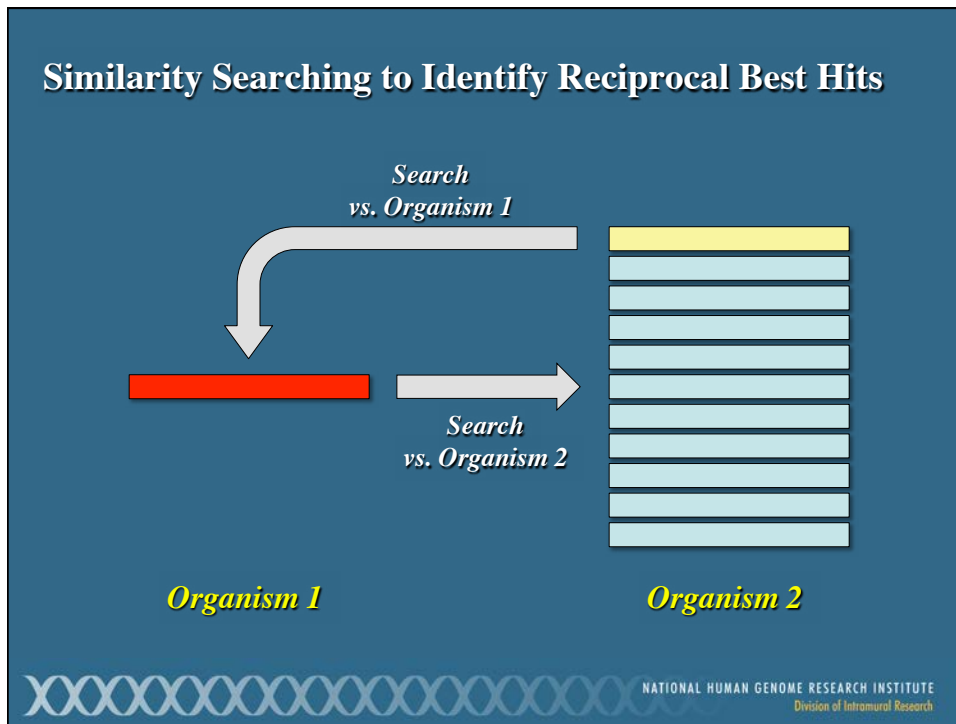


## Defining the Terms



- Genes 1-3 are orthologous
- Genes 4-6 are orthologous
- Any pair of  $\alpha$  and  $\beta$  genes are paralogous (genes related through a gene duplication event)





### Global Sequence Alignments

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

## Local Sequence Alignments

- Sequence comparison intended to find the most similar regions in the two sequences being aligned (“paired subsequences”)
- Regions outside the area of local alignment are excluded
- More than one local 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 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, 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	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	1	1	2	2	2	0	2	2	2	2	1	1	1	3	2	11	2	-3	-4	-3	-2	-4
Y	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	2	2	2	7	-1	-3	-2	-1	-4
V	0	-3	-3	-3	-1	-2	-2	-3	-3	3	1	-2	1	-1	-2	-2	0	-3	-1	4	-3	-2	-1	-4
B	-2	-1	3	4	-3	0	1	-1	0	-3	-4	0	-3	-3	-2	0	-1	-4	-3	-3	4	1	-1	-4
Z	-1	0	0	1	-3	3	4	-2	0	-3	-3	1	-1	-3	-1	0	-1	-3	-2	-2	1	4	-1	-4
X	0	-1	-1	-1	-2	-1	-1	-1	-1	-1	-1	-1	-1	-1	-2	0	0	-2	-1	-1	-1	-1	-1	-4
*	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	1

### BLOSUM62



## 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



## 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

## 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*



## Gaps

- Used to improve alignments between two sequences
  - Compensate for insertions and deletions
  - As such, *gaps represent biological events*
- Must be kept to a reasonable number, to not reflect a biologically 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

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

where	$G$ = gap-opening penalty	nucleotide	5	protein	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



## Neighborhood Words

Query Word ( $W = 3$ )

Query: GSQSLAALLNKCKT **PQG** ORLVNQWIKQPLMDKNRIEERLNLVEAFVED

Neighborhood Words


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

Neighborhood Score Threshold  
 ( $T = 13$ )



## High-Scoring Segment Pairs


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



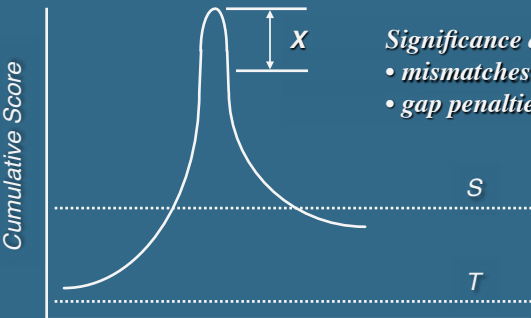
	←		→
Query:	325	SLAALLNKCKT <b>PQG</b> QRLVNQWIKQPLMDKNRIEERLNLVEA	365
		+LA++L TP G R++ +W+ +P+ D + ER + A	
Sbjct:	290	TLASVLDCTVT <b>PMG</b> SRMLKRWLHMPVRDTRVLLERQQTIGA	330

NATIONAL HUMAN GENOME RESEARCH INSTITUTE  
 Division of Intramural Research

## Extension



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

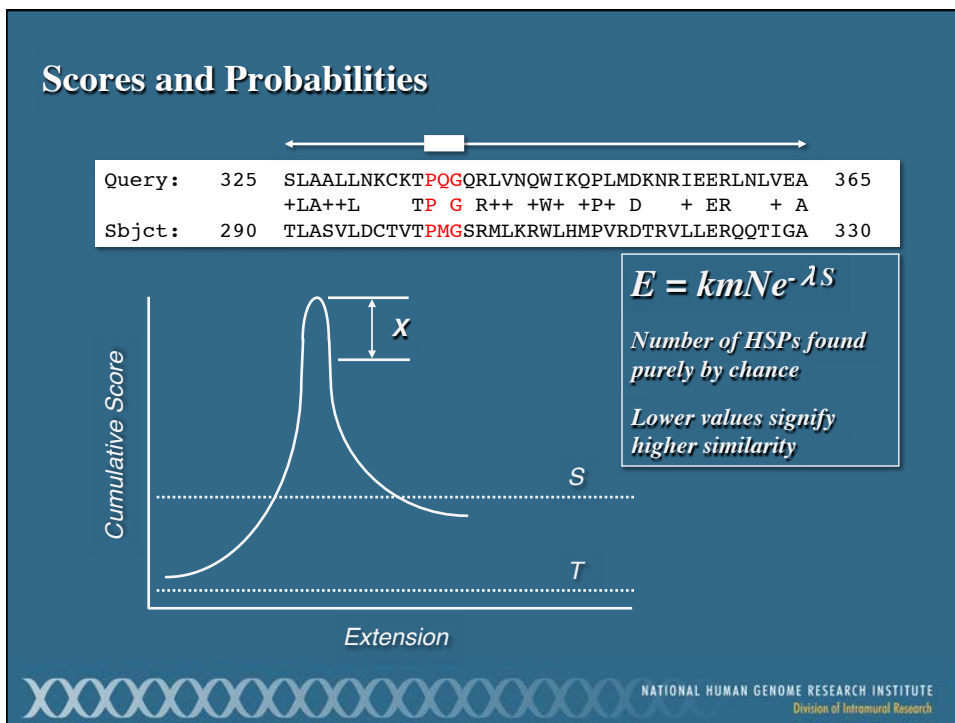
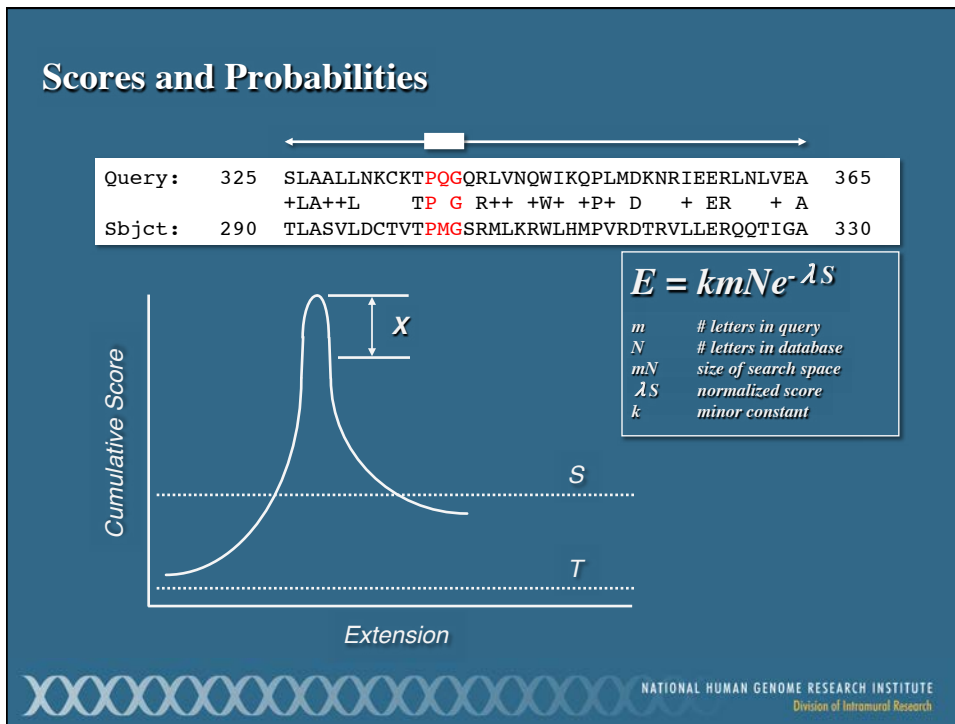


*Significance decay*

- mismatches
- gap penalties

NATIONAL HUMAN GENOME RESEARCH INSTITUTE  
 Division of Intramural Research





## Scores and Probabilities

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

Extension

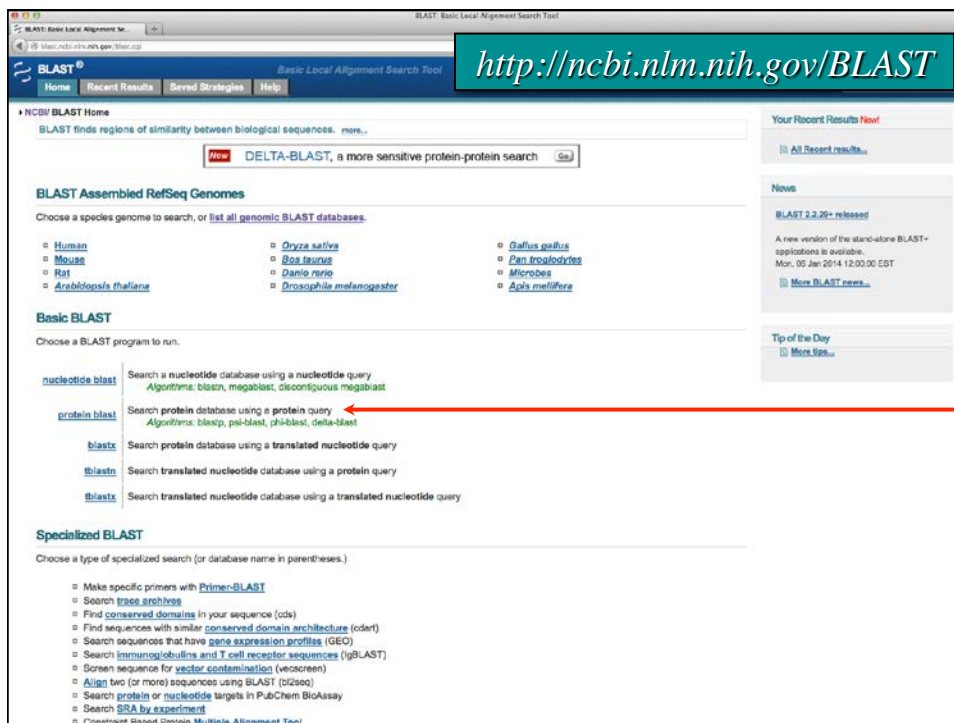
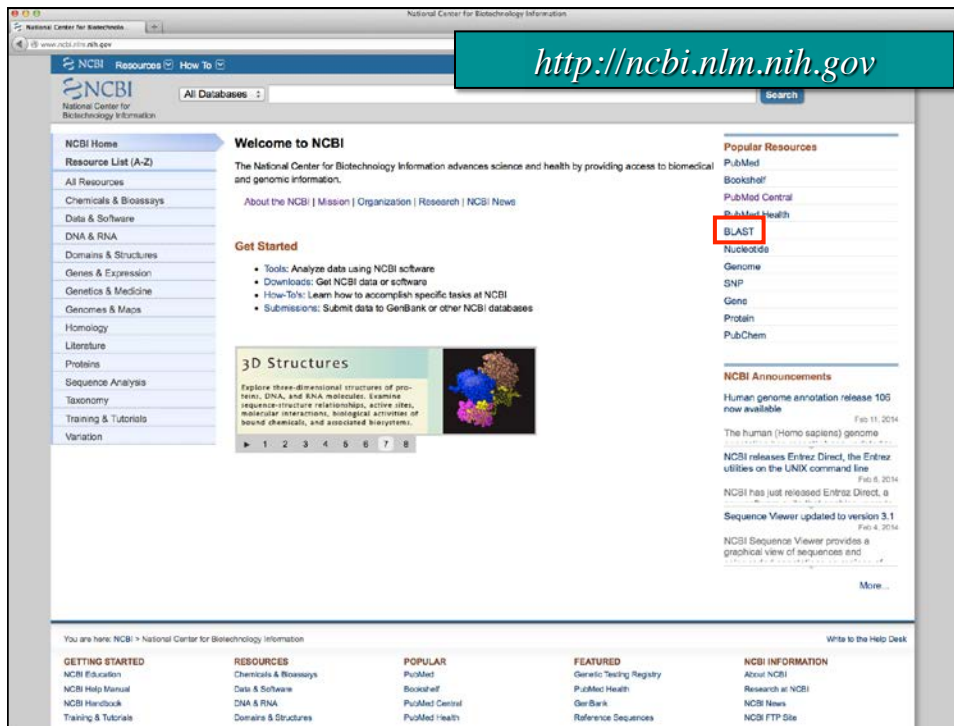
$E \leq 10^{-6}$   
for nucleotides

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

NATIONAL HUMAN GENOME RESEARCH INSTITUTE  
Division of Intramural Research

## Using BLAST for Protein Similarity Searching

NATIONAL HUMAN GENOME RESEARCH INSTITUTE  
Division of Intramural Research



## Sequences Used in Examples

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

The screenshot shows the NHGRI website interface. At the top, there's a navigation bar with links like 'Research Funding', 'Research at NHGRI', 'Health', 'Education', 'Issues in Genetics', 'Newsroom', and 'Careers & Training'. Below this, the main content area is titled 'Current Topics in Genome Analysis 2014' and 'Weeks 2 and 3: Biological Sequence Analysis'. It features sections for 'BLAST' and 'BLAST 2 Sequences', each containing a block of DNA sequence data.



NATIONAL HUMAN GENOME RESEARCH INSTITUTE  
 Division of Intramural Research

The screenshot shows the NCBI BLAST search interface. The 'Choose Search Set' section is highlighted, showing a dropdown menu for 'Database' set to 'Non-redundant protein sequences (nr)'. A text box on the right lists available protein databases: nr (Non-redundant), refseq (Reference Sequences), swissprot (SWISS-PROT), pat (Patents), pdb (Protein Data Bank), and env\_nr (Environmental samples). The 'Program Selection' section shows 'blastp (protein-protein BLAST)' selected.

Available protein databases include:

- nr* Non-redundant
- refseq* Reference Sequences
- swissprot* SWISS-PROT
- pat* Patents
- pdb* Protein Data Bank
- env\_nr* Environmental samples

## 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:*

<b>NT_</b>	Genomic contigs
<b>NM_</b>	mRNAs
<b>NP_</b>	Proteins
<b>NR_</b>	Non-coding transcripts

*From genome annotation:*

<b>XM_</b>	Model mRNA
<b>XP_</b>	Model proteins

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



BLAST® Basic Local Alignment Search Tool

Standard Protein BLAST

Enter Query Sequence

Enter accession number(s), g(i)s, or FASTA sequence(s)

Or, upload file

Job Title

Align two or more sequences

Choose Search Set

Database: Non-redundant protein sequences (nr)

Organism

Exclude

Entrez Query

Program Selection

Algorithm

blastp (protein-protein BLAST)

PSI-BLAST (Position-Specific Iterated BLAST)

PHI-BLAST (Pattern Hit Initiated BLAST)

DELTA-BLAST (Domain Enhanced Lookup Time Accelerated BLAST)

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

Algorithm parameters

Note: Parameter values that differ from the default are highlighted in yellow and marked with + sign

Algorithm parameters

Note: Parameter values that differ from the default are highlighted in yellow and marked with + sign

General Parameters

Max target sequences: 250

Short queries:  Automatically adjust parameters for short input sequences

Expect threshold: 10

Word size: 3

Max matches in a query range: 0

Scoring Parameters

Matrix: BLOSUM62

Gap Costs: Existence: 11 Extension: 1

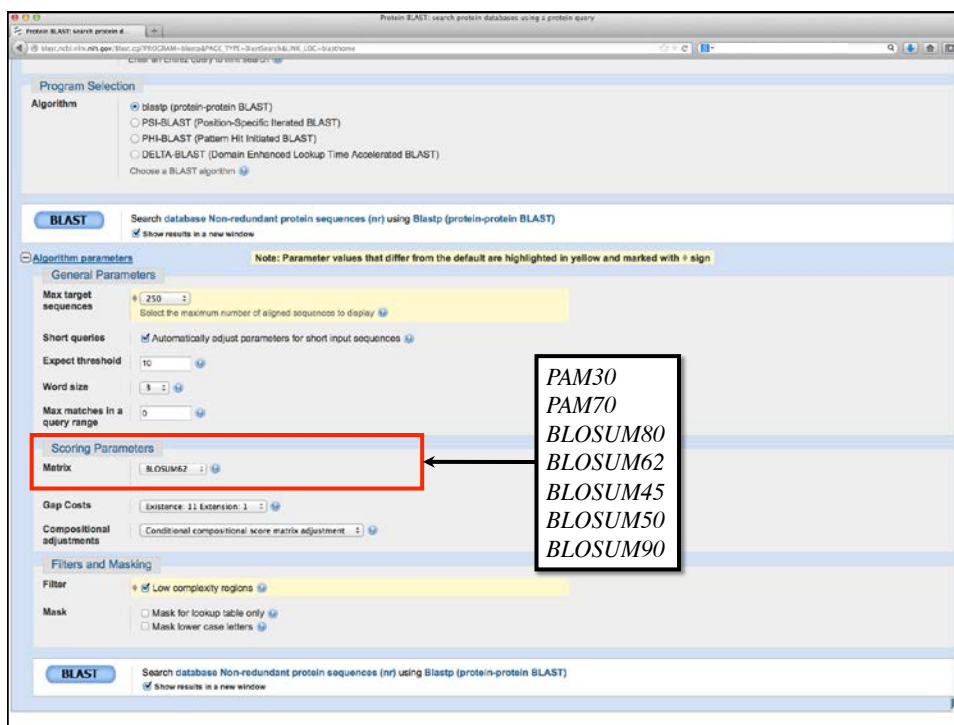
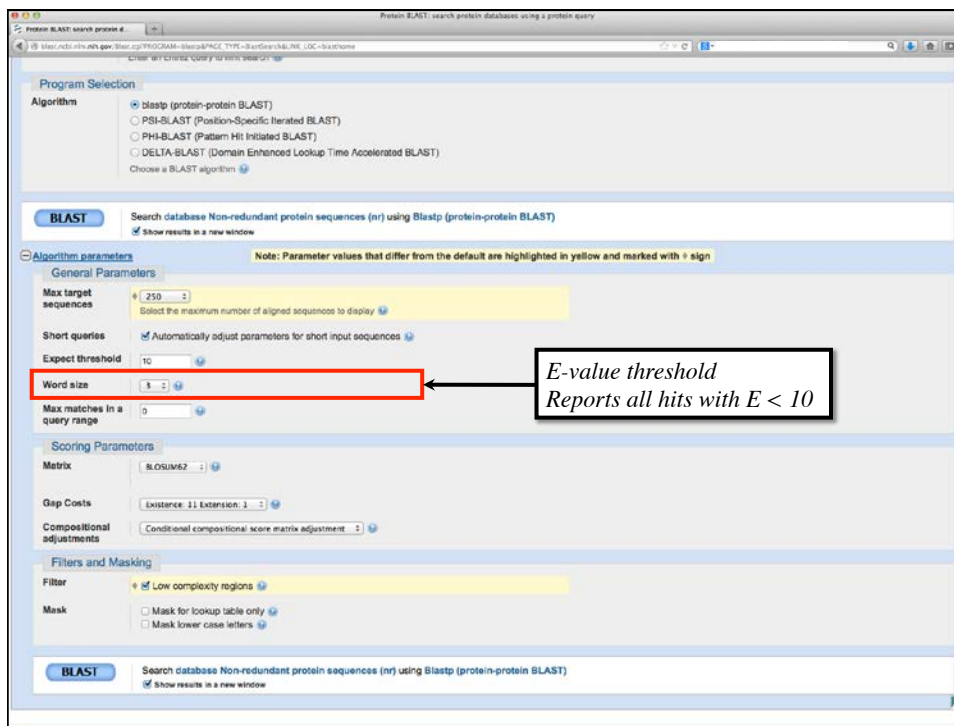
Compositional adjustments:  Conditional compositional score matrix adjustment

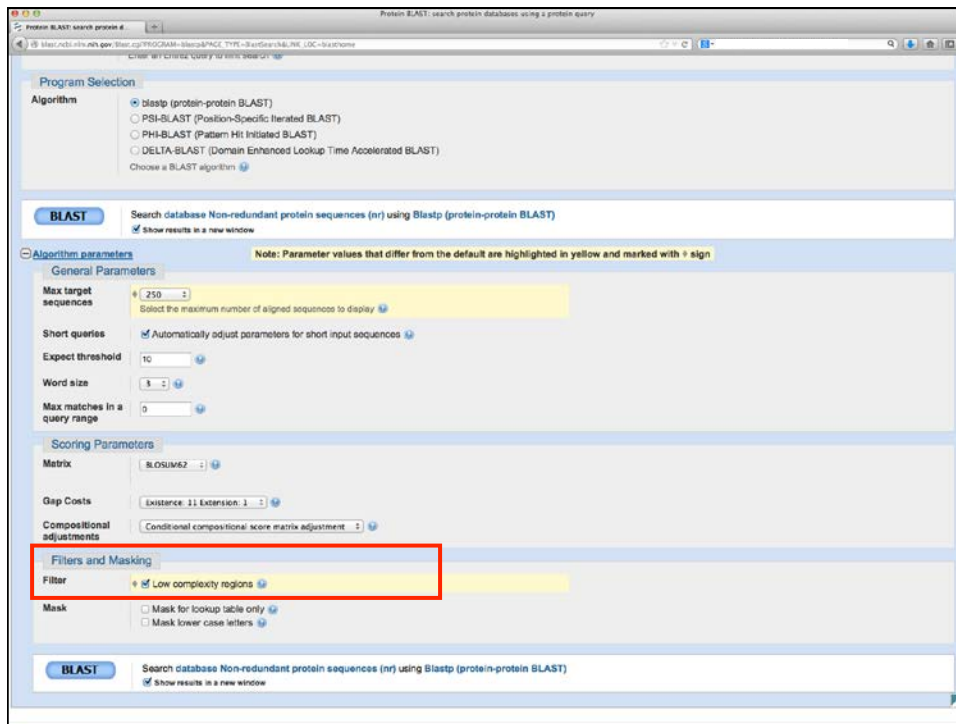
Filters and Masking

Filter:  Low complexity regions

Mask:  Mask for lookup table only  Mask lower case letters

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

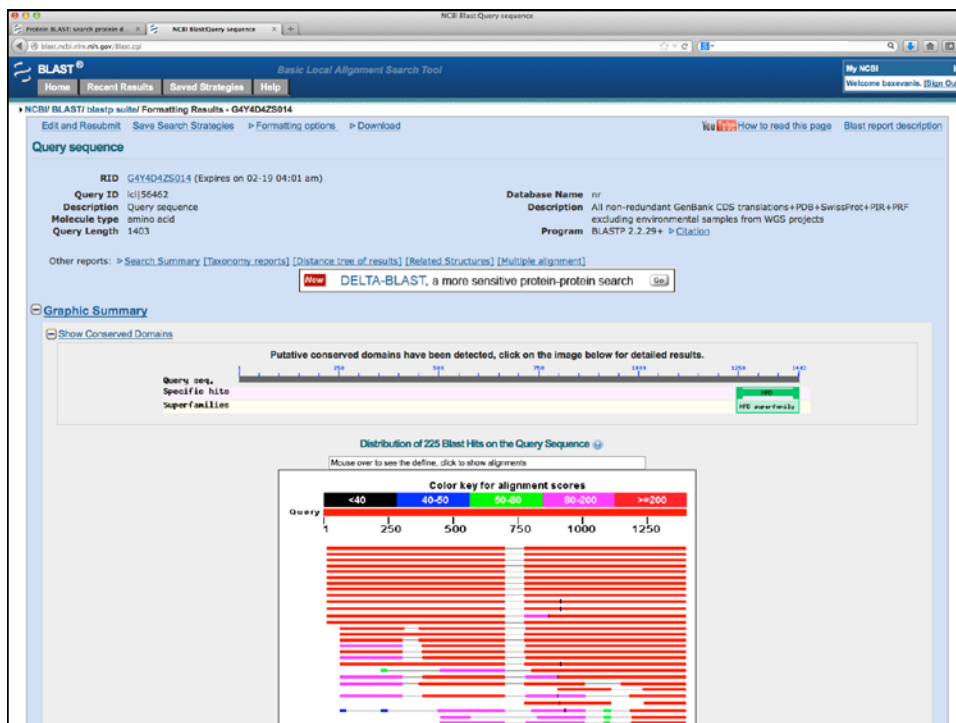
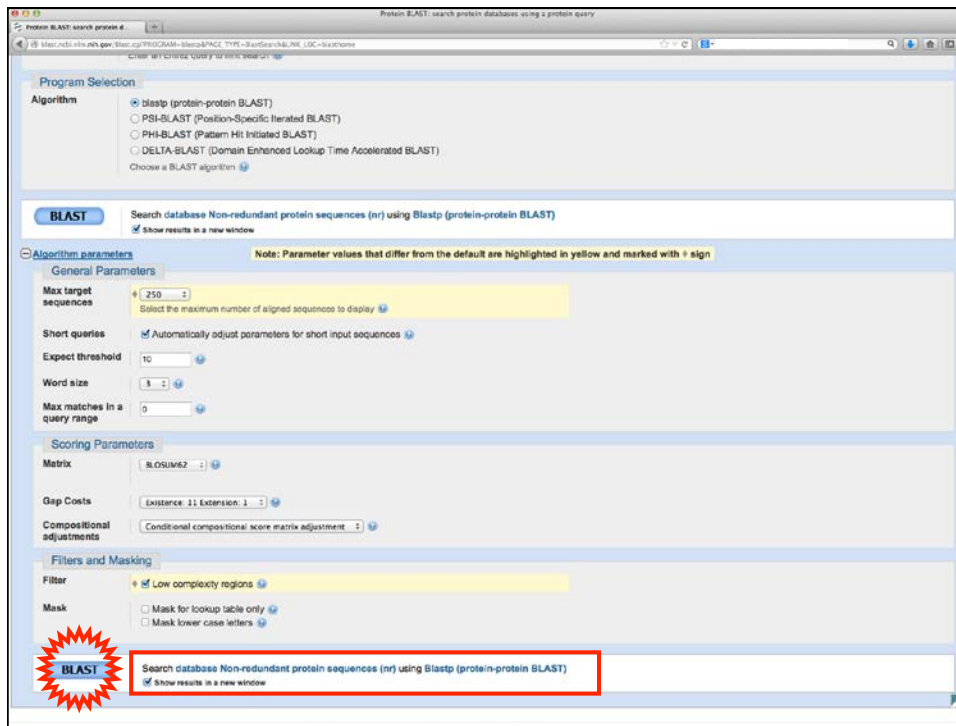


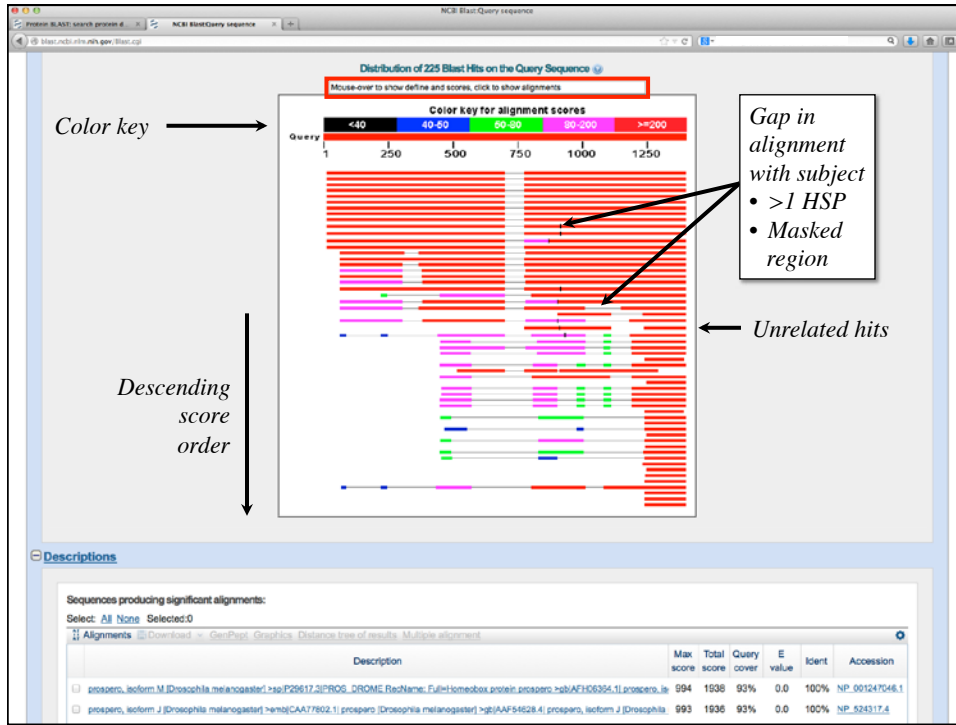


## Low-Complexity Regions

- 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)







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

3e-179 =  $3 \times 10^{-179}$

Descriptions

Sequences producing significant alignments:

Select: All None Selected:0

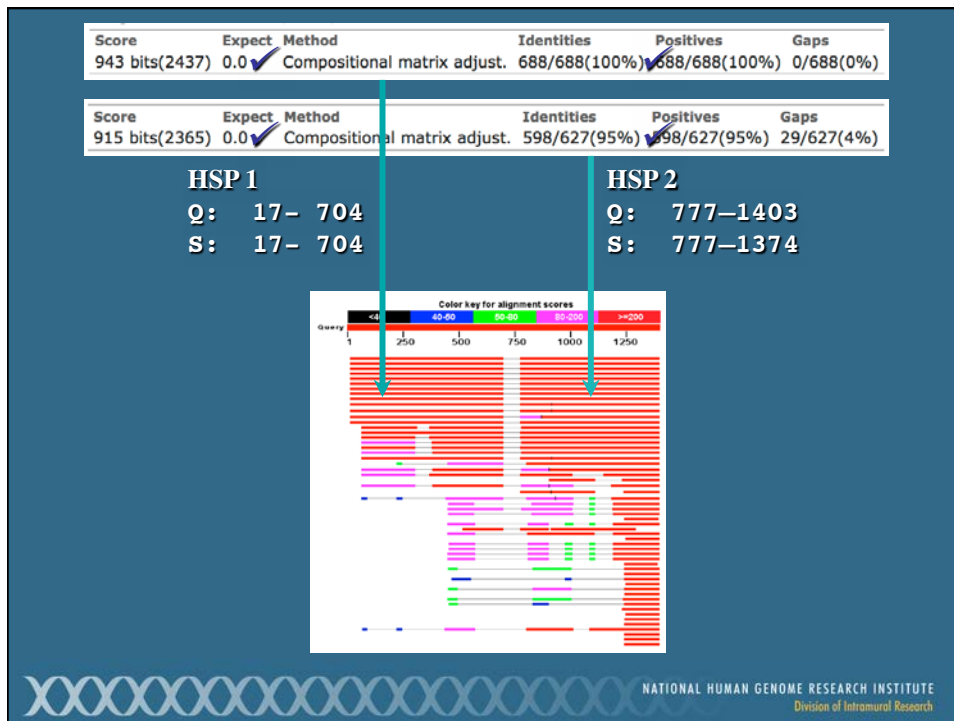
Description	Max score	Total score	Query cover	E value	Ident	Accession
prospero, isoform M [Drosophila melanogaster] :expP26617.3 PROG_DROME_RefName: Full-Homobox protein:prospero :gb AF156354.1 prospero_iso	594	1936	93%	0.0	100%	NP_001247046.1
prospero, isoform J [Drosophila melanogaster] :hmhCAA7802.1 prospero [Drosophila melanogaster] :gb AA15428.4 prospero, isoform J [Drosophila	593	1936	93%	0.0	100%	NP_524317.4
prospero [Drosophila melanogaster]	593	1932	93%	0.0	100%	AA01464.1
homeodomain transcription factor:Prospero [Drosophila melanogaster]	590	1621	93%	0.0	100%	AA05703.1
GG18089 [Drosophila erecta] :gb EDV4931.1 GG18089 [Drosophila erecta]	569	1865	93%	0.0	99%	XP_001580573.1
Proa protein [Drosophila melanogaster]	562	1611	93%	0.0	97%	AA42841.1
GE26090 [Drosophila yakuba] :gb E0Y999.1.3.1 GE26090 [Drosophila yakuba]	561	1865	93%	0.0	97%	XP_002097261.1
prospero, isoform H [Drosophila melanogaster] :gb AF156352.1 prospero, isoform H [Drosophila melanogaster]	544	1662	93%	0.0	100%	NP_001247044.1
prospero, isoform L [Drosophila melanogaster] :gb AA194462.1 LD37627p [Drosophila melanogaster] :gb AA113400.4 prospero, isoform L [Drosophila]	543	1858	93%	0.0	100%	NP_789636.3
prospero, isoform I [Drosophila melanogaster] :gb AF156353.1 prospero, isoform I [Drosophila melanogaster]	542	1864	93%	0.0	100%	NP_001247045.1
prospero, isoform K [Drosophila melanogaster] :gb AA113401.4 prospero, isoform K [Drosophila melanogaster]	542	1863	93%	0.0	100%	NP_731665.4
GM23939 [Drosophila sechellia] :gb EDV42817.1 GM23939 [Drosophila sechellia]	535	1967	93%	0.0	98%	XP_002031601.1
GF19657 [Drosophila aransaeae] :gb EDV42775.1 GF19657 [Drosophila aransaeae]	504	1673	93%	0.0	92%	XP_001554214.1
GK11290 [Drosophila williston] :gb EDV80845.1 GK11290 [Drosophila williston]	545	1532	85%	0.0	83%	XP_002069959.1
GA14403, isoform B [Drosophila pseudoobscura pseudoobscura] :gb FAI29137.3 GA14403, isoform B [Drosophila pseudoobscura pseudoobscura]	545	1470	90%	0.0	81%	XP_001559965.3
proa [Drosophila virilis] :gb EDV68804.1 proa [Drosophila virilis]	821	1429	85%	0.0	79%	XP_002032084.1
GH21437 [Drosophila grimshawi] :gb EDV69096.1 GH21437 [Drosophila grimshawi]	809	1374	84%	0.0	80%	XP_001559430.1
GI22095 [Drosophila melanogaster] :gb EDV15591.1 GI22095 [Drosophila melanogaster]	804	1392	84%	0.0	80%	XP_002000100.1
PSBEDCTED: homeobox protein:prospero-like isoform X2 [Conostictus capitata]	691	11				
GA14403, isoform C [Drosophila pseudoobscura pseudoobscura] :gb EM52398.1 GA14403, isoform C [Drosophila pseudoobscura pseudoobscura]	611	14				
homeobox protein:prospero-like 1 [Aedes aegypti] :gb EA146002.1 AAEL002769-PA, partial [Aedes aegypti]	571	770	62%	3e-179	59%	XP_001659942.1
PSBEDCTED: homeobox protein:prospero-like isoform X1 [Conostictus capitata]	491	1093	84%	5e-144	82%	XP_004529242.1
RefName: Full-Homobox protein:prospero :gb AF156354.1 AF156354.1, homeodomain transcription factor:Prospero [Drosophila virilis]	440	1309	75%	3e-126	80%	Q91461.1
prospero [Drosophila simulans] :gb EDX13377.1 prospero [Drosophila simulans]	378	606	26%	3e-112	100%	XP_002103974.1
PSBEDCTED: LOW QUALITY PROTEIN: homeobox protein:prospero-like [Musca domestica]	397	1080	72%	1e-110	82%	XP_005188268.1
Chain A, Structural Basis of Prospero-Dna Interaction; Implications For Transcription Regulation In Developing Cells	347	347	11%	2e-107	100%	1XPX_A

Accept (for now)

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

$\geq 25\%$  for proteins  
 $\geq 70\%$  for nucleotides

Gap  
 a Low complexity



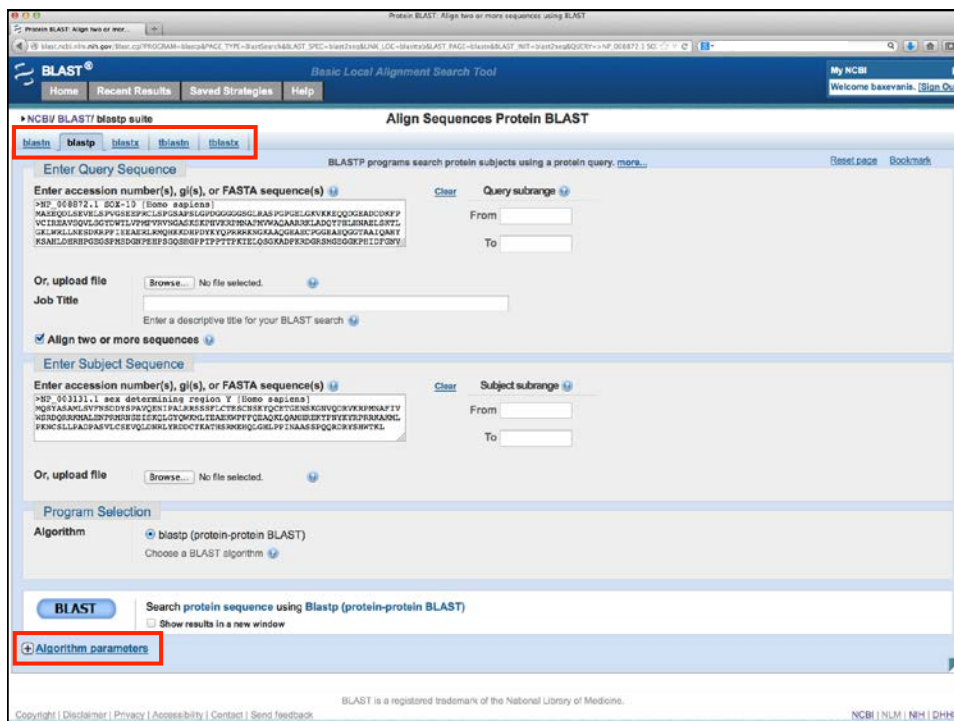
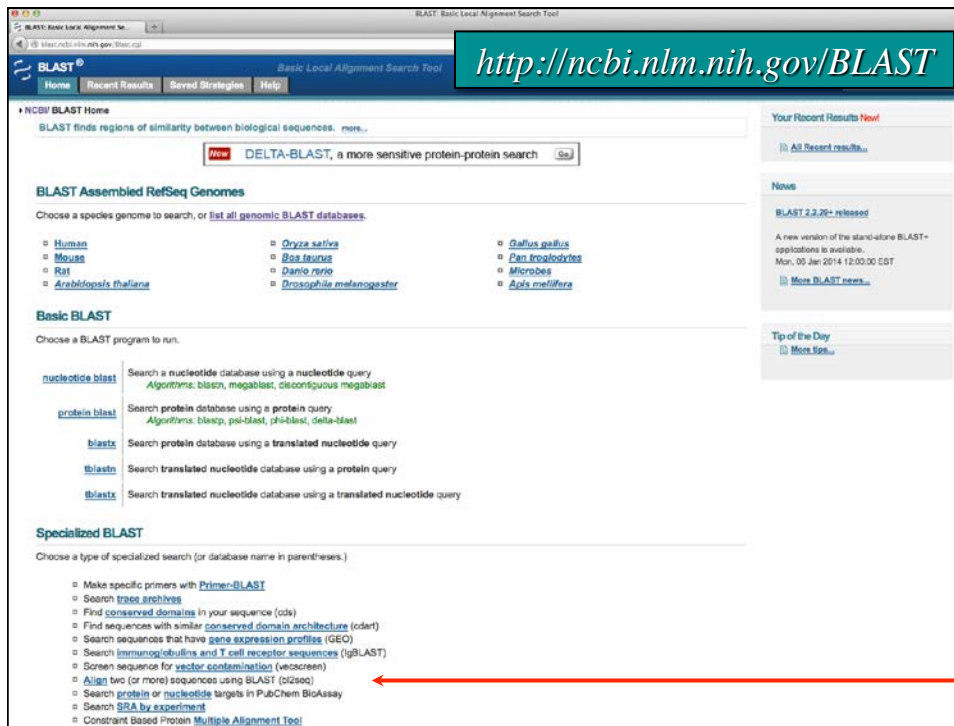
## Suggested BLAST Cutoffs

	<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



Protein BLAST: Align two or more sequences using BLAST

Search protein sequence using **Blastp** (protein-protein BLAST)

**Algorithm parameters** Note: Parameter values that differ from the default are highlighted in yellow and marked with + sign

**General Parameters**

- Max target sequences: 100
- Short queries:  Automatically adjust parameters for short input sequences
- Expect threshold: 10
- Word size: 3
- Max matches in a query range: 0

**Scoring Parameters**

- Matrix: BLOSUM62
- Gap Costs: Existence: 11 Extension: 1
- Compositional adjustments: Conditional compositional score matrix adjustment

**Filters and Masking**

- Filter:  Low complexity regions
- Mask:  Mask for lookup table only  Mask lower case letters

**BLAST** Search protein sequence using **Blastp** (protein-protein BLAST)  Show results in a new window

BLAST is a registered trademark of the National Library of Medicine.

NCBI BLAST: blastp suite 2-sequences/ Formatting Results - GAU650SP114

**Blast 2 sequences**

**NP\_008872.1 SOX-10 [Homo sapiens] (466)**

RID: GAU650SP114 (Expires on 02-19 02:53 am)

Query ID: lc|19239  
 Description: NP\_008872.1 SOX-10 [Homo sapiens]  
 Molecule type: amino acid  
 Query Length: 466

Subject ID: lc|19241  
 Description: NP\_003131.1 sex determining region Y [Homo sapiens]  
 Molecule type: amino acid  
 Subject Length: 204  
 Program: BLASTP 2.2.29+ Citation

Other reports: Search Summary Taxonomy reports Multiple alignment

**Graphic Summary**

Distribution of 2 Blast Hits on the Query Sequence

Color key for alignment scores

Score Range	Color
<40	Black
40-50	Blue
50-80	Green
80-200	Red
>=200	Magenta

Query: 1 90 180 270 360 450

**Dot Matrix View**

**Descriptions**

Sequences producing significant alignments:

Select: All None Selected: 0

Description	Max score	Total score	Query cover	E value	Ident	Accession
NP_003131.1 sex determining region Y [Homo sapiens]	94.0	109	19%	1e-26	46%	19241

NCBI Blast NP\_003131.1 S01-10 [Homo sapiens] 1165

Dot Matrix View

Descriptions

Sequences producing significant alignments:

Select: All None Selected:0

Alignments Download Graphics Multiple alignment

Description	Max score	Total score	Query cover	E value	Ident	Accession
NP_003131.1 sex determining region Y [Homo sapiens]	94.0	109	19%	1e-26	46%	19241

Alignments

Download Graphics Sort by: E value

NP\_003131.1 sex determining region Y [Homo sapiens]  
 Sequence ID: lc|19241 Length: 204 Number of Matches: 2

Related Information

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 NGASKKPHVRRPMNAFVVAQAARRKLADQYPLHNAELKTLGKLMWLLNESDKRPF 154  
 N + VERFMNAF+VN++ RRK+A = P + N+E+SK LG W++L E++K PF  
 Sbjct 51 NSKGRVQDRVRRPMNAFIVNSRQRRREALENFRMNSLISKOLGYQWFKLZAEKWPFF 110

Query 155 EEAERLRMQHRKDHFDYKQPRR 178  
 +EA++L+ H++ +P+YK+PRR+  
 Sbjct 111 QEAGKLGMMHREKYPYKYPFRK 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 GYDFILV 88  
 GY W ++  
 Sbjct 95 GYQNKGL 101

NCBI Blast NP\_003131.1 S01-10 [Homo sapiens] 1165

Graphic Summary

Distribution of 2 Blast Hits on the Query Sequence

Color key for alignment scores

Score Range	Color
<40	Black
40-60	Blue
60-80	Green
80-200	Pink
>=200	Red

Query 1 90 180 270 360 450

Dot Matrix View


Plot of lc|19239 vs lc|19241

Descriptions

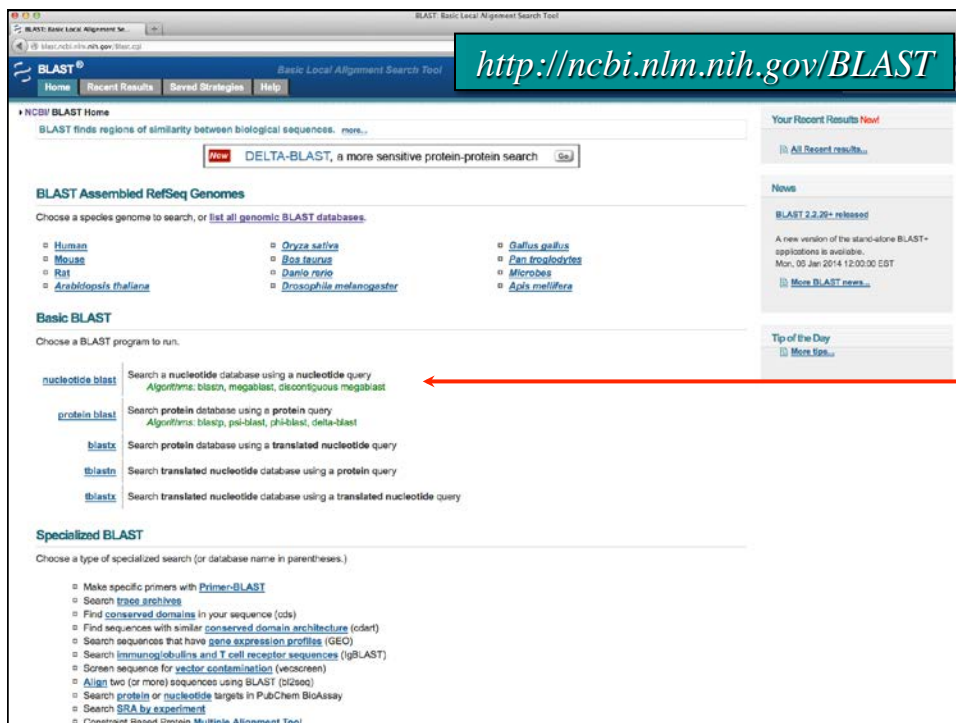
Alignments



# Nucleotide Similarity Searching MegaBLAST, BLASTN, and BLAT



NATIONAL HUMAN GENOME RESEARCH INSTITUTE  
Division of Intramural Research



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

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

**DELTA-BLAST**, a more sensitive protein-protein search

### BLAST Assembled RefSeq Genomes

Choose a species genome to search, or [list all genomic BLAST databases.](#)

<input type="checkbox"/> Human	<input type="checkbox"/> <i>Oryza sativa</i>	<input type="checkbox"/> <i>Gallus gallus</i>
<input type="checkbox"/> Mouse	<input type="checkbox"/> <i>Bos taurus</i>	<input type="checkbox"/> <i>Pan troglodytes</i>
<input type="checkbox"/> Rat	<input type="checkbox"/> <i>Danio rerio</i>	<input type="checkbox"/> <i>Morone saxatilis</i>
<input type="checkbox"/> <i>Arabidopsis thaliana</i>	<input type="checkbox"/> <i>Drosophila melanogaster</i>	<input type="checkbox"/> <i>Apis mellifera</i>

### Basic BLAST

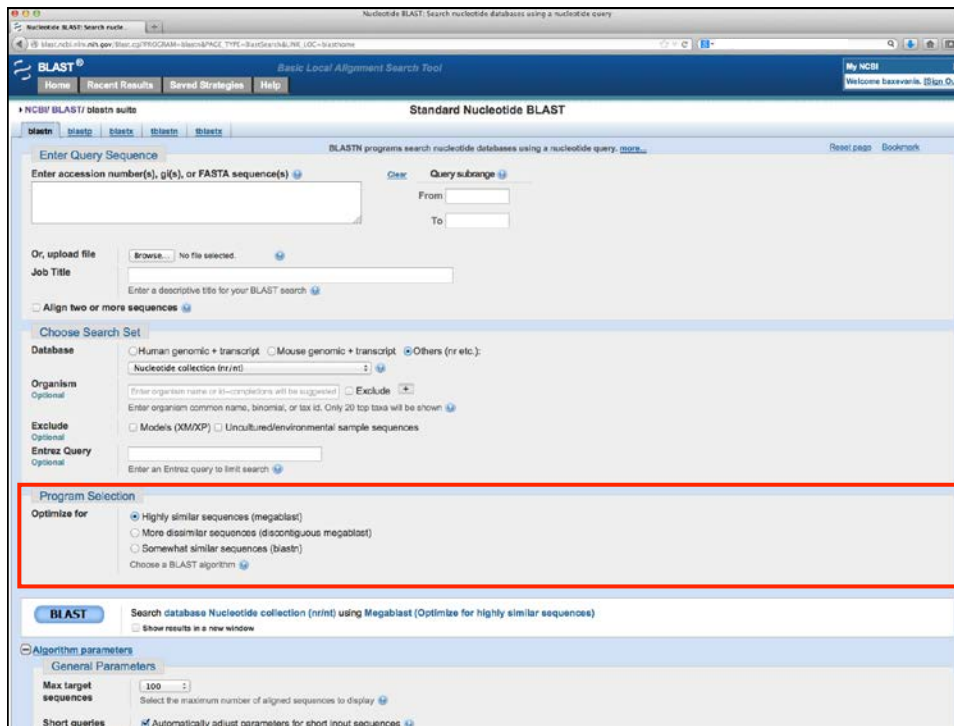
Choose a BLAST program to run:

<a href="#">nucleotide blast</a>	Search a nucleotide database using a nucleotide query <i>Algorithms: blastn, megablast, discontinuous megablast</i>
<a href="#">protein blast</a>	Search protein database using a protein query <i>Algorithms: blastp, psi-blast, phi-blast, delta-blast</i>
<a href="#">tblastx</a>	Search protein database using a translated nucleotide query
<a href="#">tblastn</a>	Search translated nucleotide database using a protein query
<a href="#">tblastx</a>	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 [tree archives](#)
- Find [conserved domains](#) in your sequence (cds)
- Find sequences with similar [conserved domain architecture](#) (cdart)
- Search sequences that have [gene expression profiles](#) (GEO)
- Search [immunoglobulin and T cell receptor sequences](#) (igBLAST)
- Screen sequence for [vector contamination](#) (vecscreen)
- [Align](#) two (or more) sequences using BLAST (tblast2)
- Search [protein or nucleotide targets](#) in PubChem BioAssay
- Search [SRA by experiment](#)
- [Constraint Based Protein Multiple Alignment Tool](#)



## Nucleotide-Based BLAST Algorithms

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



The screenshot shows the UCSC Genome Bioinformatics website homepage. The URL bar displays <http://genome.ucsc.edu>. The page features a navigation menu with links to Genomes, Blat, Tables, Gene Sorter, PCR, VialGene, Session, FAQ, and Help. A sidebar on the left lists various tools and resources such as Genome Browser, ENCODE, Neanderthal, Blat, Table Browser, Gene Sorter, In Silico PCR, Genome Graphs, Galaxy, VialGene, Utilities, Downloads, Release Log, Custom Tracks, Cancer Browser, Microbial Genomes, Mirrors, Training, Credits, Publications, and Cite Us. The main content area includes an "About the UCSC Genome Bioinformatics Site" section, a "News" section with recent announcements (e.g., "22 January 2014 - All the DNA on the internet now at your fingertips!"), and a "Publications" section.

The screenshot shows the Rat BLAT Search interface. The page title is "Rat BLAT Search". The main heading is "BLAT Search Genome". The interface includes a search form with the following fields: Genome (set to Rat), Assembly (set to Nov. 2004 (Baylor 3.4/m4)), Query type (set to DNA), Sort output (set to query.score), and Output type (set to hyperlink). Below the form is a text area containing a DNA sequence: 

```
>CB312815 NICED_Rr Piti1 Rattus norvegicus cDNA clone
GGGCTCTGGCTGGCTGTGTCTCAGAAGCTGCTTCTCCACCTCTCCCTGTGAATTCCTAAACTCTC
TAGCTCTGGCTCACTGTCTCTCTGGATAGTCTGTGTGGCAAGAGCCCTTAAGAGAAATATGCAATGA
GCTATAAGAGTGTGGCTGTGGCTGTGGCTGTGGCTGTGGCTGTGGCTGTGGCTGTGGCTGTGGCTGTGG
GCTCTTAAGTACAGGTTATCCAGAGCCCACTTACCCCAAGAGACAGCCTCTCCCCATCCCTAGGAAA
CAGTAGAGCTTAGGAAAAGAAAGAAAGAAAGAAAGAAAGAAAGAAAGAAAGAAAGAAAGAAAGAAAGAA
GATTTTCAGTTCTGAAATCTGTCTTGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTGT
AAAAACCAATGAATTAAGGAAAATGAAAATCAAGCCCACTAAACAATGAAGAAAGTGTCTATGTCTT
GATCAGGGGGGGATATAGCTCAGTCAAGAGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTGT
CCGACACGAAAAGAAAAGAAAAGAAAAGAAAAGAAAAGAAAAGAAAAGAAAAGAAAAGAAAAGAAAAG
GGGAAAGGTTTAAACACCTACTGAGGAAATGGTAAAGCGGAGTGCCTTGTCTATACTCGGGATGGCT
AGTCAACAGTAAAGAAAATTTGAAAATCTAAATACCAATGGGATGGATCCCTTAAACCAATCC
```

 Below the sequence is a "submit" button and a "feeling lucky" link. The page also includes instructions for pasting a query sequence, a "File Upload" section, and an "About BLAT" section.

Rat BLAT Results

BLAT Search Results

ACTIONS	QUERY	SCORE	START	END	QSIZE	IDENTITY	CHRO	STRAND	START	END	SPAN
<a href="#">browser details</a>	CB312815	710	1	733	768	98.14	5	+	101455599	101456223	725
<a href="#">browser details</a>	CB312815	29	501	537	768	89.24	2	+	38736251	38736287	37
<a href="#">browser details</a>	CB312815	25	501	529	768	93.24	3	+	22960346	22960374	29
<a href="#">browser details</a>	CB312815	22	341	363	768	100.04	1	+	122930956	122930979	24
<a href="#">browser details</a>	CB312815	21	202	222	768	100.04	17	-	33246146	33246166	21
<a href="#">browser details</a>	CB312815	21	706	727	768	100.04	3	+	46857920	46857942	23
<a href="#">browser details</a>	CB312815	21	552	574	768	95.74	1	+	157973111	157973133	23
<a href="#">browser details</a>	CB312815	20	277	298	768	95.54	2	-	240446870	240446891	22
<a href="#">browser details</a>	CB312815	20	452	461	768	100.04	1	-	216323127	216323146	20
<a href="#">browser details</a>	CB312815	20	508	527	768	100.04	1	-	56102029	56102048	20
<a href="#">browser details</a>	CB312815	20	453	474	768	95.54	2	+	186587336	186587357	22

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

chr5:101,455,417-101,456,504 1,088 bp.

- 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

Genomes Genome Browser Tools Mirrors Downloads My Data About Us Help

### BLAT Search Results

ACTIONS	QUERY	SCORE	START	END	QSIZE	IDENTITY	CHRO	STRAND	START	END	SPAN
<a href="#">browser details</a>	CB312815	710	1	733	768	98.14	5	+	101455599	101456323	725
<a href="#">browser details</a>	CB312815	25	501	537	768	89.24	2	+	38736251	38736287	37
<a href="#">browser details</a>	CB312815	25	501	529	768	93.28	3	+	22960346	22960374	29
<a href="#">browser details</a>	CB312815	22	341	363	768	100.04	1	+	122930956	122930979	24
<a href="#">browser details</a>	CB312815	21	202	222	768	100.04	17	-	33246146	33248166	21
<a href="#">browser details</a>	CB312815	21	706	727	768	100.04	3	+	46857920	46857942	23
<a href="#">browser details</a>	CB312815	21	552	574	768	95.74	1	+	157973111	157973133	23
<a href="#">browser details</a>	CB312815	20	277	298	768	95.54	2	-	240446870	240446891	22
<a href="#">browser details</a>	CB312815	20	452	461	768	100.04	1	-	216323127	216323146	20
<a href="#">browser details</a>	CB312815	20	508	527	768	100.04	1	-	56102029	56102048	20
<a href="#">browser details</a>	CB312815	20	453	474	768	95.54	2	+	186587336	186587357	22

User Sequence vs Genome

Alignment of CB312815 and chr5:101455599-101456323

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).

CB312815  
 Rat.chr5  
 block1  
 together

**cDNA CB312815**

```

GgGGCTCTCG CTGGCCCTGG TCTCAGAAGC TGCFFTCCTC ACCTCTTCCT 50
TGTGAATTC CTAAACTCTC TACCTCTGGT TCATGFTCCG TCTCTGGAT 100
ACTCTGTCTG CAATGAGCCC TTAAGGAAAT ATFGCAATGA GCTATAGAG 150
TTGTGAGCCT GGGTAGGCA AGGCTTGAC TGGGACAGCA AAGGAAATTT 200
CATTGCATCT GCTCCTAAGT CACAGGTTAT CCAGAGCCCA CTTTACCCCA 250
AGAGCAGCC TCTCCCCCA? CCTAGGAAA CAGTAGAGCT TAGGAAATG 300
AATGACTCCA CCACATTCAA GAGGCTTCAA ATFGTAACT TGGCATTTC? 350
GATTCACGCT CTGAAATCT? GTCCCTTAG? CDTGGGAAA ATAGAAATG 400
GAGTTACACC TTGTCATTTA AAAAACCA? GAATTAAGG AAATGGAAA 450
TCATGCCCC ATAAACATG TATGGAAGT TTCATGTTT GATCATGGC 500
GGGATATAG CTCAGTCAT? GAGTCTTTC ATAGCAATG GCATAATCC 550
AGTTTCAMC CCGAGCACG AAAAAGAAA GCGGAGGAG TGGAGGCATT 600
CACAGCACG TTTTCAGT? AGGCGCAAG GGGAGGAG? TTAACACCT 650
ACTGAGGgAA TGGATAAGC GAGTGCCT? GTCTACTC GGGgatgGCT 700
AGTCATClog taGAAAAG? TTGgaAAATG ATAAaatacc aatgggatgg 750
atccccctta aaccatcc
    
```

**Genomic chr5 :**

```

cttggaaag gtaactata cattaatata gagccctott tttcttgcg 101455548
ggcccaggac acacaggacg gatgtttccc agtcaatcca gggacagcat 101455598
GgGGCTCTCG CTGGCCCTGG TCTCAGAAGC TGCFFTCCTC ACCTCTTCCT 101455648
TGTGAATTC CTAAACTCTC TACCTCTGGT TCATGFTCCG TCTCTGGAT 101455698
AGTCTGTCTG CAATGAGCCC TTAAGGAAAT ATFGCAATGA GCTATAGAG 101455748
TTGTGAGCCT GGGTAGGCA AGGCTTGAC TGGGACAGCA AAGGAAATTT 101455798
CATTGCATCT GCTCCTAAGT CACAGGTTAT CCAGAGCCCA CTTTACCCCA 101455848
AGAGCAGCC TCTCCCCCA? CCTAGGAAA CAGTAGAGCT TAGGAAATG 101455898
AATGACTCCA CCACATTCAA GAGGCTTCAA ATFGTAACT TGGCATTTC? 101455948
GATTCACGCT CTGAAATCT? GTCCCTTAG? CDTGGGAAA ATAGAAATG 101455998
GATTCACGCT TTGTCATTTA AAAAACCA? GAATTAAGG AAATGGAAA 101456048
TCATGCCCC ATAAACATG TATGGAAGT TTCATGTTT GATCATGGC 101456098
GGGATATAG CTCAGTCAT? GAGTCTTTC ATAGCAATG GCATAATCC 101456148
AGTTTCAMC CCGAGCACG AAAAAGAAA GCGGAGGAG TGGAGGCATT 101456198
CACAGCACG TTTTCAGT? AGGCGCAAG GGGAGGAG? TTAACACCT 101456248
ACTGAGGAA? GGATAAGCG AGTGCCT?GT CTACTCTCG gatctactc 101456298
ATCagTAGAA AAGTTTAAA TGTAgatacc gatggatgat cccctaaaca 101456348
tctccctaac taacagatcc agcctagaaa gggatattt ccatatttc 101456398
    
```

## 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>*

