

“Google is arguably one of the greatest inventions. The search engine is one of the greatest inventions in human history.”

-- Franklin Foer

“Probabilistic models of transcriptomic (dys)regulation”

David Knowles

ECCR265

Thursday, February 22, 2018

Concepts of Sequence Similarity Searching

- The premise:

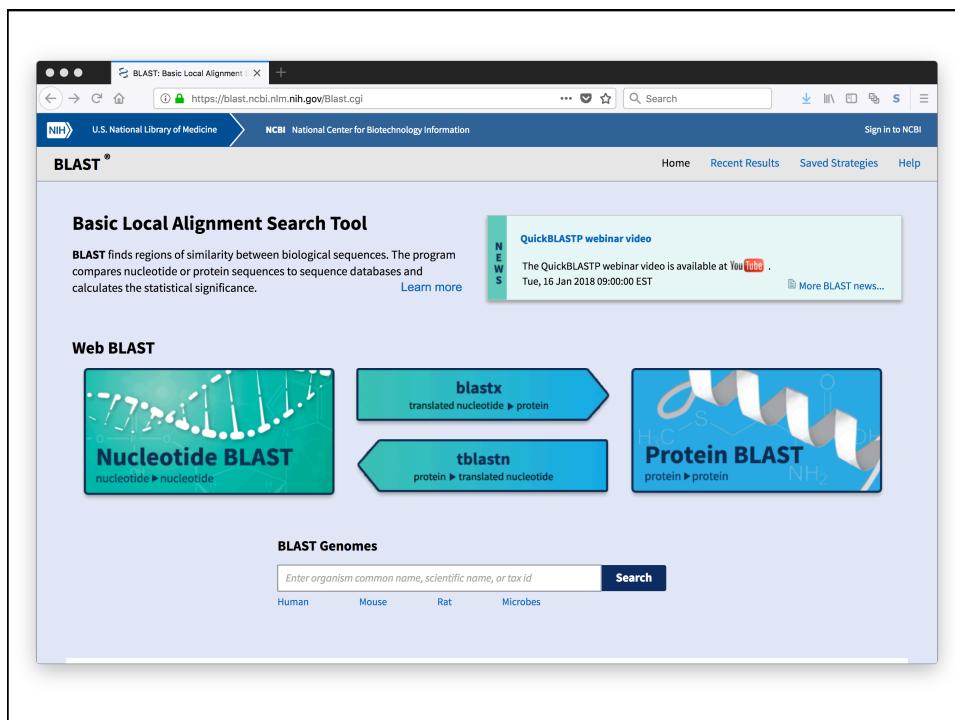
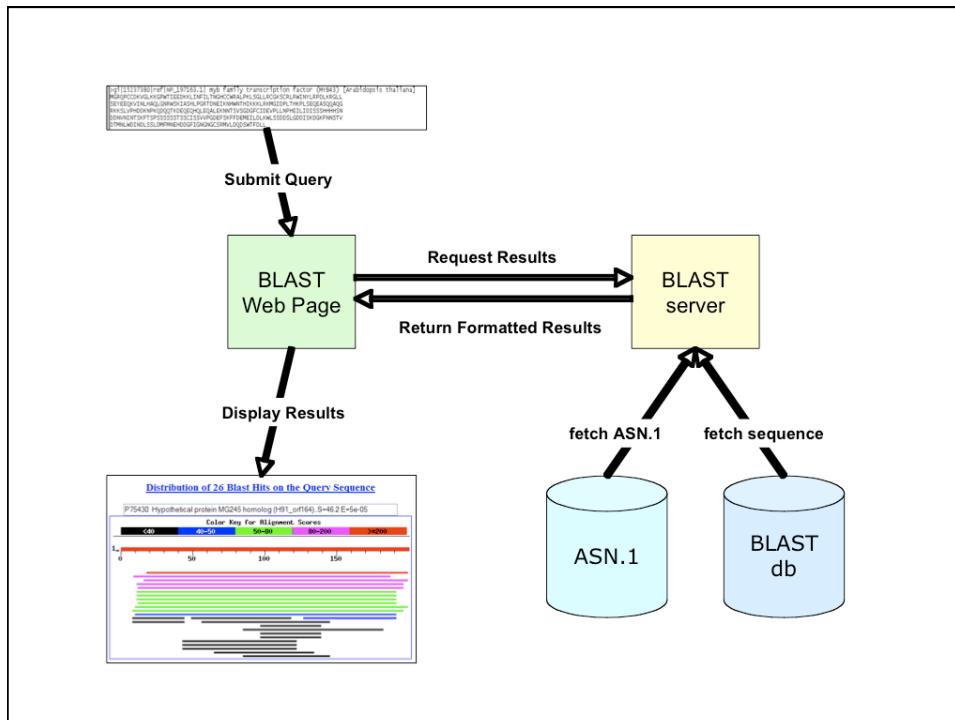
One sequence by itself is not informative; it must be analyzed by comparative methods against existing sequence databases to develop hypothesis concerning relatives and function.

3

The BLAST algorithm

The BLAST programs (Basic Local Alignment Search Tools) are a set of sequence comparison algorithms introduced in 1990 that are used to search sequence databases for optimal local alignments to a query.

- Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ (1990) "Basic local alignment search tool." *J. Mol. Biol.* 215:403-410.
- Altschul SF, Madden TL, Schaeffer AA, Zhang J, Zhang Z, Miller W, Lipman DJ (1997) "Gapped BLAST and PSI-BLAST: a new generation of protein database search programs." *NAR* 25:3389-3402.



BLAST programs

Program	Description
blastp	Compares an amino acid query sequence against a protein sequence database.
blastn	Compares a nucleotide query sequence against a nucleotide sequence database.
blastx	Compares a nucleotide query sequence translated in all reading frames against a protein sequence database. You could use this option to find potential translation products of an unknown nucleotide sequence.
tblastn	Compares a protein query sequence against a nucleotide sequence database dynamically translated in all reading frames.
tblastx	Compares the six-frame translations of a nucleotide query sequence against the six-frame translations of a nucleotide sequence database.

BLAST is an experiment!

- BLAST is not Google
- BLAST is like doing an experiment: to get good, meaningful results, you need to optimize the experimental conditions
 1. What kind of BLAST?
 2. Pick an appropriate database
 3. Pick the right algorithm
 4. Choose parameters

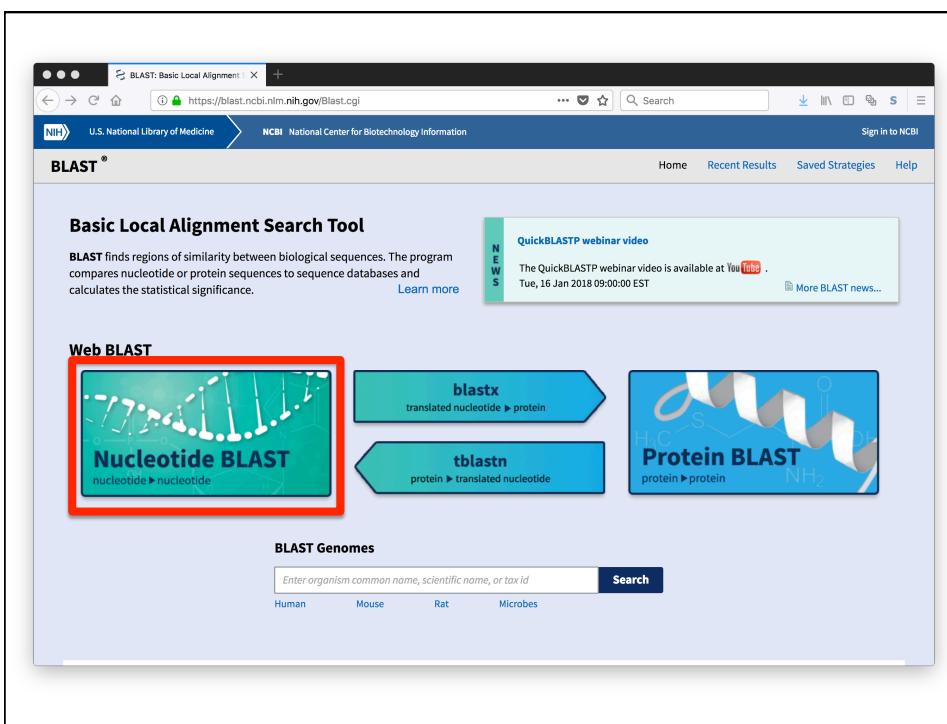
Step 1:

Nucleotide BLAST vs. protein BLAST

- Largely determined by your query sequence

BUT

- If your nucleotide sequence can be translated to a peptide sequence, you probably want to do it
- Protein blasts are more sensitive and biologically significant
- Sometimes it makes sense to use other blasts



Step 2: Choose a Database

- Too large:
 - Takes longer
 - Too many results
 - More random, meaningless matches
- Too small or wrong one:
 - Miss significant matches

The screenshot shows the NCBI Nucleotide BLAST search interface. The main header reads "Nucleotide BLAST: Search nuc...". Below it, the URL is https://blast.ncbi.nlm.nih.gov/Blast.cgi?PROGRAM=blastn&PAGE_TYPE=BlastSearch. The page title is "BLAST" followed by "blastn suite". The top navigation bar includes links for Home, Recent Results, Saved Strategies, and Help. A red arrow points to the "Choose Search Set" section.

Standard Nucleotide BLAST

Enter Query Sequence
Enter accession number(s), g(s), or FASTA sequence(s)

From:
To:

Or, upload file No file selected.

Job Title:

Enter a descriptive title for your BLAST search

Align two or more sequences

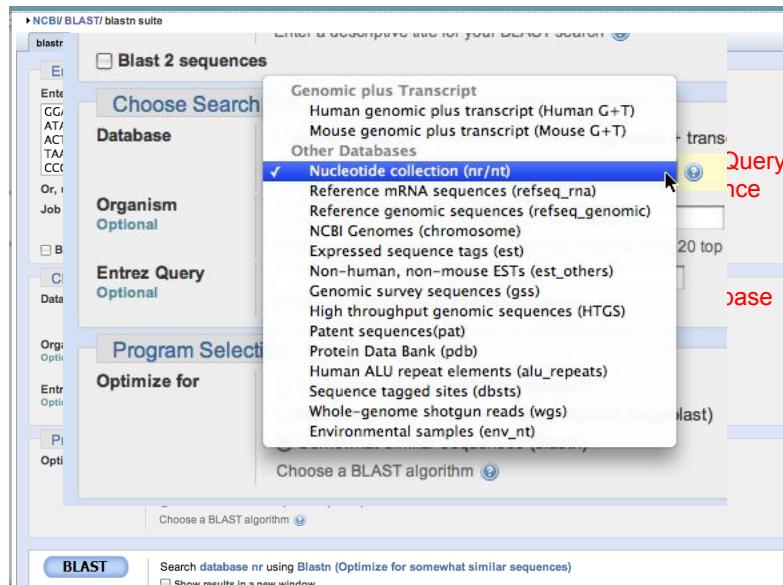
Choose Search Set

Database: Human genomic + transcript Mouse genomic + transcript Others (nr etc.) Nucleotide collection (nt/nt)
Organism: Enter organism name or ID—completions will be suggested Exclude
Exclude: Models (XM/XP) Uncultured/environmental sample sequences
Limit to: Sequences from type material
Entrez Query: You
Optional: Enter an Entrez query to limit search

Program Selection

Optimize for: Highly similar sequences (megablast) More dissimilar sequences (discontiguous megablast) Somewhat similar sequences (blastn)

blastn Options



Protein Databases

- Non-redundant protein sequences (nr)
 - Kitchen-sink:
 - Translations of GenBank coding sequences (CDS)
 - RefSeq Proteins
 - PDB (RCSB Protein Data Bank - 3d-structure)
 - SwissProt
 - Protein Information Resource (PIR)
 - Protein Research Foundation (Japanese DB)
- Reference proteins (refseq_protein)
 - NCBI Reference Sequences: Comprehensive, integrated, non-redundant, well-annotated set of sequences
- Swissprot protein sequences (swissprot)
 - Swiss-Prot: European protein database (no incremental updates)

The screenshot shows the 'Standard Nucleotide BLAST' search page. The 'Program Selection' section is highlighted with a large red arrow pointing to it. This section contains the following options:

- Optimize for:**
 - Highly similar sequences (megablast)
 - More dissimilar sequences (discontiguous megablast)
 - Somewhat similar sequences (blastn)
- Choose a BLAST algorithm:** A link to select a different algorithm.

Step 3: Choose an Algorithm

- How close a match are you looking for?
- Determines how similarities are “scored”
- Affects speed of search and chance of missing match
- Again, what is the goal of the search?

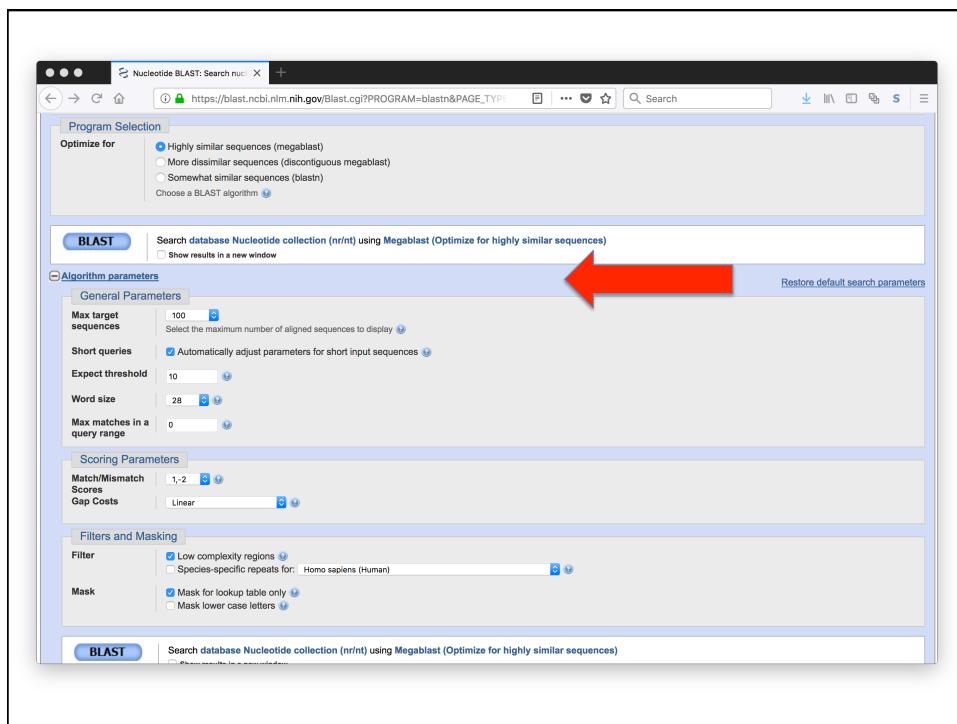
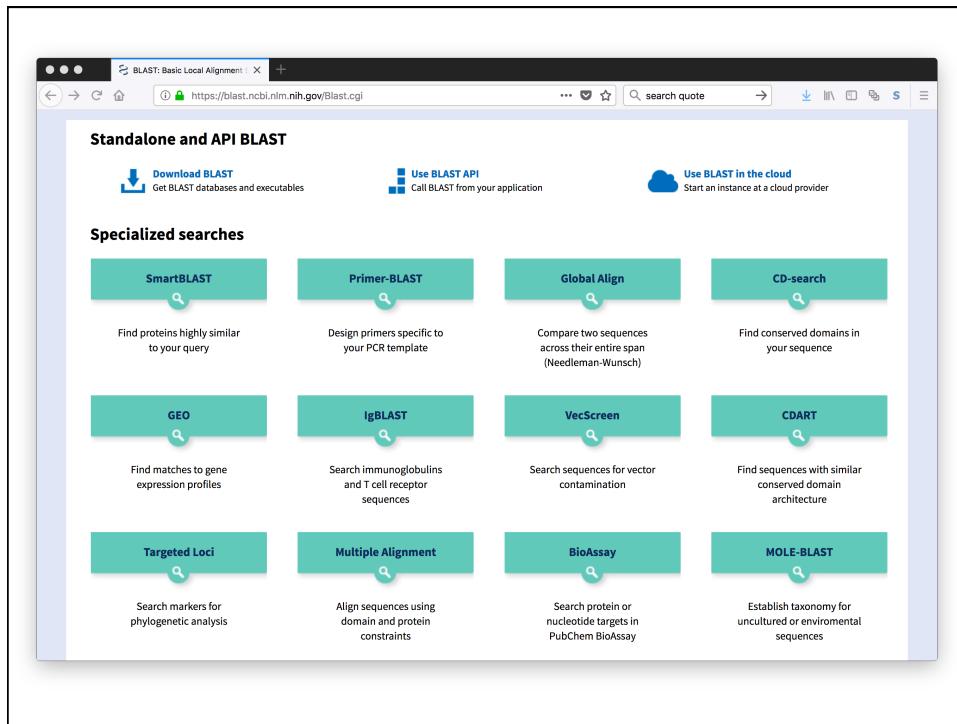
more BLAST programs

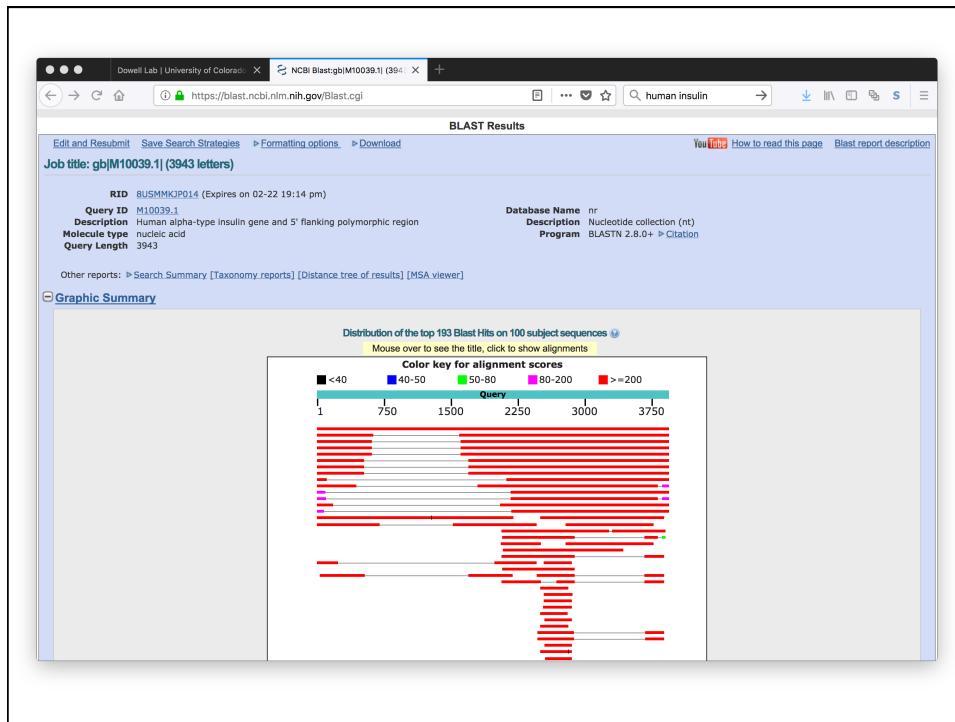
Program		Notes
Megablast	Contiguous	Nearly identical sequences
	Discontiguous	Cross-species comparison
Position Specific	PSI-BLAST	Automatically generates a position specific score matrix (PSSM)
	RPS-BLAST	Searches a database of PSI-BLAST PSSMs

 nucleotide only
 protein only

Specialized BLAST versions

- **Megablast** – (nucleotides) very fast for highly similar sequences
- **PSI-Blast** – (proteins) Position Specific Iterated (profile searching)
- Whole slew of even more specialized BLAST versions for specific reoccurring cases ...

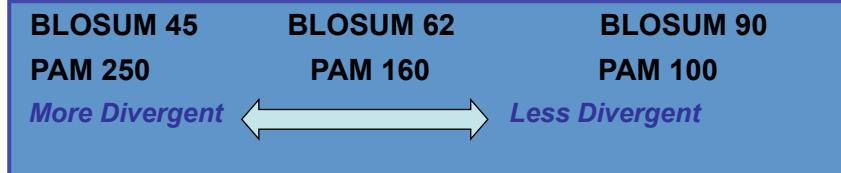




Where does the score (S) come from?

- The quality of each pair-wise alignment is represented as a score and the scores are ranked.
- **Scoring matrices** are used to calculate the score of the alignment base by base (DNA) or amino acid by amino acid (protein).
- The alignment score will be the sum of the scores for each position.

BLOSUM vs PAM



- BLOSUM 62 is the default matrix in BLASTp 2.0. Though it is tailored for comparisons of moderately distant proteins, it performs well in detecting closer relationships. A search for distant relatives may be more sensitive with a different matrix.

Database searching: What is a good hit?

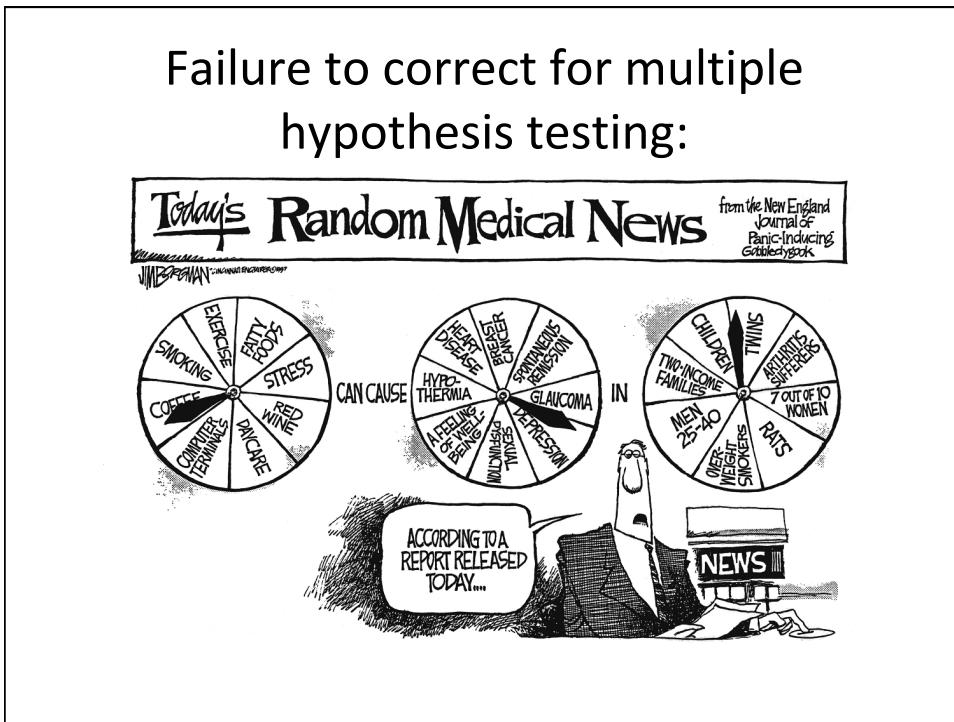
- Say that you search the non-redundant protein database at NCBI, containing roughly one million sequences. What p-value threshold should you use?

Multiple testing

- Say that you perform a statistical test with a 0.05 p-value threshold, but you repeat the test on twenty different observations.
- Assume that all of the observations are explainable by the null hypothesis.
- What is the chance that at least one of the observations will receive a p-value less than 0.05?

Multiple testing

- Say that you perform a statistical test with a 0.05 threshold, but you repeat the test on twenty different observations. Assuming that all of the observations are explainable by the null hypothesis, what is the chance that at least one of the observations will receive a p-value less than 0.05?
- $\Pr(\text{making a mistake}) = 0.05$
- $\Pr(\text{not making a mistake}) = 0.95$
- $\Pr(\text{not making any mistake}) = 0.95^{20} = 0.358$
- $\Pr(\text{making at least one mistake}) = 1 - 0.358 = 0.642$
- **There is a 64.2% chance of making at least one mistake.**



Bonferroni correction

- Assume that individual tests are *independent*.
(Is this a reasonable assumption?)
- Divide the desired p-value threshold by the number of tests performed.
- For the previous example, $0.05 / 20 = 0.0025$.

$\Pr(\text{making a mistake}) = 0.0025$
 $\Pr(\text{not making a mistake}) = 0.9975$
 $\Pr(\text{not making any mistake}) = 0.9975^{20} = 0.9512$
 $\Pr(\text{making at least one mistake}) = 1 - 0.9512 = 0.0488$

Database searching

- Say that you search the non-redundant protein database at NCBI, containing roughly one million sequences. What p-value threshold should you use?
- Say that you want to use a conservative p-value of 0.001.
- Recall that you would observe such a p-value by chance approximately every 1000 times in a random database.
- A Bonferroni correction would suggest using a p-value threshold of $0.001 / 1,000,000 = 0.000000001 = 10^{-9}$.

The screenshot shows a web browser window with the URL <https://blast.ncbi.nlm.nih.gov/Blast.cgi>. The search term 'human insulin' is entered in the search bar. The results table has columns for Description, Max score, Total score, Query cover, E value, Ident, and Accession. The table lists various insulin-related genes from different species, including Human, Homo sapiens, Gorilla, Pan troglodytes, Pongo abelii, and C. elegans. The top result is 'Human alpha-type insulin gene and 5' flanking polymorphic region' with a Max score of 7282 and 100% query cover.

Description	Max score	Total score	Query cover	E value	Ident	Accession
Human alpha-type insulin gene and 5' flanking polymorphic region	7282	7282	100%	0.0	100%	M10039_1
Homo sapiens chromosome 11, clone RP4-539G11, complete sequence	4132	5137	75%	0.0	98%	AC130303.8
Homo sapiens insulin (INS) gene, complete cds	4065	5032	74%	0.0	98%	AH02844.2
Homo sapiens tyrosine hydroxylase (TH) gene, 3' end; insulin (INS) gene, complete cds; insulin-like growth factor 2 (IGF2) gene, 5' end	4065	5032	74%	0.0	98%	L15440_1
Human gene for proinsulin, from chromosome 11. Includes a highly polymorphic region upstream from the insulin gene containing tandemly repeated ss	4065	5032	74%	0.0	98%	Y00565_1
Home sapiens INS-IGF2 readthrough (INS-IGF2) RefSeqGene on chromosome 11	3982	4792	70%	0.0	99%	NG_050578_1
Home sapiens insulin (INS). RefSeqGene on chromosome 11	3982	4792	70%	0.0	99%	NG_007114_1
Home sapiens chromosome 11, clone RP11-889I17, complete sequence	3982	4792	70%	0.0	99%	AC132217_15
Home sapiens haplotype I (Ca tyrosine hydroxylase (TH) gene, partial sequence; insulin (INS) gene, complete cds; and insulin-like growth factor 2 (IGF2) ge	3341	3544	48%	0.0	99%	AH012037_2
P. troglodytes gene for proproinsulin	3291	4005	64%	0.0	96%	X61089_1
Gorilla gorilla tyrosine hydroxylase (TH) gene, partial cds; tyrosine hydroxylase/insulin intergenic spacer, partial sequence; and insulin precursor (INS) gene	2900	3065	47%	0.0	96%	AH011915_2
Pan troglodytes tyrosine hydroxylase (TH) gene, partial cds; and insulin precursor (INS) gene, complete cds	2870	3179	46%	0.0	98%	AH011814_2
Pongo abelii BAC clone CH276-476G11 from chromosome unknown, complete sequence	2863	3091	52%	0.0	94%	AC199962_4
Pongo pygmaeus tyrosine hydroxylase (TH) gene, partial cds; tyrosine hydroxylase/insulin intergenic spacer, partial sequence; and insulin precursor (INS)	2665	2788	46%	0.0	94%	AH011916_2
Home sapiens insr gene, partial	2545	2545	35%	0.0	99%	AJ038655_1
Human hypervariable DNA 6' to the insulin (14) gene	1888	3645	55%	0.0	93%	M28868_1
Gorilla gorilla insulin gene, partial cds	1589	1589	24%	0.0	96%	AY092023_1
C. elegans gene for proproinsulin	1517	2290	45%	0.0	90%	X61092_1
Human insulin gene; alpha allele S7flank (ulrich)	1465	2574	41%	0.0	95%	J00286_1
PREDICTED: Pan troglodytes insulin (INS) transcript variant X1, mRNA	1437	1768	25%	0.0	99%	XM_016919751_1
Pongo pygmaeus insulin gene, partial cds	1434	1434	24%	0.0	93%	AY092024_1

Sequence Similarity Searching – The statistics are important

- Discriminating between real and artifactual matches is done using an estimate of probability that the match might occur by chance.
- The Score (S) reflects the desired scoring scheme. Karlin-Altshul statistics gives us a sense for whether that's a good score (p-value), but BLAST reports E-values!

E-values

- A p-value is the probability of making a mistake ... but properly adjusting it requires *knowing the size of the database!*
- The E-value is the expected number of times that the given score would appear in a random database of the given size.
- One simple way to compute the E-value is to multiply the p-value times the size of the database.

BLAST actually calculates E-values in a much more complex way.

E-values

- Thus, for a p-value of 0.001 and a database of 1,000,000 sequences, the corresponding E-value is $0.001 \times 1,000,000 = 1,000$.
- When $E < 0.01$, e-value and p-value are essentially the same.
- Hence smaller numbers are more significant!

Short sequences can't get good e-values!!

- What is the probability of finding a 12 base fragment in a "random" genome?
 $4^{12} = 16,777,216$ (once per 16 million bases)
- What length DNA fragment is needed to define a unique location in the human genome?
 $4^{16} = 4,294,967,296$ (4 billion bases)

What do the Score and the e-value really mean?

- The quality of the alignment is represented by the **Score (S)**.
- The score of an alignment is calculated as the sum of substitution and gap scores. Substitution scores are given by a look-up table (PAM, BLOSUM) whereas gap scores are assigned empirically.
- The significance of each alignment is computed as an **E value (E)**.
- Expectation value. The number of different alignments with scores equivalent to or better than S that are expected to occur in a database search by chance. The lower the E value, the more significant the score.

The E-value is not a probability; it's an expected value!

The screenshot shows a web browser displaying the NCBI Blast search results for the query 'human insulin'. The URL is <https://blast.ncbi.nlm.nih.gov/Blast.cgi>. The search interface includes a 'Descriptions' section and a table of results. The table has columns for Description, Max score, Total score, Query cover, E value, Ident, and Accession. The results list various homologous sequences from different species, including Human, Gorilla, Pan, and Pongo. The table highlights several entries with blue borders, such as Human alpha-type insulin gene and Human gene for proinsulin.

Description	Max score	Total score	Query cover	E value	Ident	Accession
Human alpha-type insulin gene and 5' flanking polymorphic region	7282	7282	100%	0.0	100%	M10039.1
Human chromosome 11, clone RP4-539011, complete sequence	4132	5137	75%	0.0	98%	AC130303.8
Human sapiens insulin (INS) gene, complete cds	4065	6032	74%	0.0	98%	AH02844.2
Human sapiens tyrosine hydroxylase (TH) gene, 3' end; insulin (INS) gene, complete cds; insulin-like growth factor 2 (IGF2) gene, 5' end	4065	6032	74%	0.0	98%	L15440.1
Human gene for proproinsulin, from chromosome 11. Includes a highly polymorphic region upstream from the insulin gene containing tandemly repeated	4065	6032	74%	0.0	98%	V00565.1
Human sapiens INS-IGF2 readthrough (INS-IGF2), RefSeqGene on chromosome 11	3982	4792	70%	0.0	99%	NG_005078.1
Human sapiens insulin (INS), RefSeqGene on chromosome 11	3982	4792	70%	0.0	99%	NG_007114.1
Human sapiens chromosome 11, clone RP11-889I17, complete sequence	3982	4792	70%	0.0	99%	AC132217.15
Human sapiens haplotype (Ca tyrosine hydroxylase (TH) gene, partial sequence; insulin (INS) gene, complete cds; and insulin-like growth factor 2 (IGF2) ge	3341	3544	48%	0.0	99%	AH012037.2
Pirogotes gene for proproinsulin	3291	4005	64%	0.0	96%	X61089.1
Gorilla gorilla tyrosine hydroxylase (TH) gene, partial cds; tyrosine hydroxylase/insulin intergenic spacer, partial sequence; and insulin precursor (INS) gene	2900	3065	47%	0.0	96%	AH011815.2
Pan troglodytes tyrosine hydroxylase (TH) gene, partial cds; and insulin precursor (INS) gene, complete cds	2870	3179	46%	0.0	98%	AH011814.2
Pongo abelii BAC clone CH278-476G11 from chromosome unknown, complete sequence	2863	3091	52%	0.0	94%	AC199982.4
Pongo pygmaeus tyrosine hydroxylase (TH) gene, partial cds; tyrosine hydroxylase/insulin intergenic spacer, partial sequence; and insulin precursor (INS)	2665	2788	46%	0.0	94%	AH011816.2
Human sapiens ins gene, partial	2545	2545	35%	0.0	99%	AJ009655.1
Human hypervariable DNA 5' to the insulin (INS) gene	1888	3645	55%	0.0	93%	M26868.1
Gorilla gorilla insulin gene, partial cds	1589	1589	24%	0.0	96%	AY092023.1
Caethiops gene for proproinsulin	1517	2230	45%	0.0	90%	X61092.1
Human insulin gene, alpha allele 5'thank (Ulrich)	1465	2574	41%	0.0	95%	J00266.1
PREDICTED: Pan troglodytes insulin (INS), transcript variant X1, mRNA	1437	1768	25%	0.0	99%	XM_016919751.1
Pongo pygmaeus insulin gene, partial cds	1434	1434	24%	0.0	93%	AY092024.1

It is important to keep in mind that the E() value does not represent a measure of similarity between the two sequences.

E-values

- Statistical significance depends on both the size of the alignments and the size of the sequence database
 - ▶ Important consideration for comparing results across different searches
 - ▶ E-value increases as database gets bigger
 - ▶ E-value decreases as alignments get longer

Interpretation of E-values

- very low E() values are identical genes or possibly closely related genes
- moderate E() values are almost certainly related genes
- long list of gradually declining of E() values indicates a large gene family
- long regions of moderate similarity are more significant than short regions of high identity

Homology: Some Guidelines

- Similarity can be indicative of homology
- Generally, if two sequences are significantly similar over entire length they are likely homologous
- BUT … Low complexity regions can be highly similar without being homologous
- AND … Homologous sequences not always highly similar

Biological Relevance

- It is up to you, the biologist to scrutinize these alignments and determine if they are significant.
- Were you looking for a short region of nearly identical sequence or a larger region of general similarity?
- Are the mismatches conservative ones?
- Are the matching regions important structural components of the genes or just introns and flanking regions?

```

>gb|AAL08419.1| PTEN [Takifugu rubripes]
Length=412

Score = 197 bits (501), Expect = 2e-49, Method: Composition-based stats.
Identities = 95/100 (95%), Positives = 98/100 (98%), Gaps = 0/100 (0%)

Query   2     IVSRNKKRRYQEDGFDLDTIYIPNIIAMGFFPAERLEGRGVYRNNNIDVVVRFLDSKHKNHYKI  61
+VSRNKKRRYQEDGFDLDTIYIPNIIAMGFFPAERLEGRGVYRNNNIDVVVRFLDSKHKNHYKI
Sbjct   8     MVSRNKKRRYQEDGFDLDTIYIPNIIAMGFFPAERLEGRGVYRNNNIDVVVRFLDSKHKNHYKI  67

Query   62    YNLCAERHYD TAKFNCRVAQYPFEDHNPPQLELIKPFQN  101
YNLCAERHYD AFKNCRVAQYPFEDHNPPQLELIKPF ++
Sbjct   68    YNLCAERHYD AAKFNCRVAQYPFEDHNPPQLELIKPFCED  107

Score = 83.6 bits (205), Expect = 4e-15, Method: Composition-based stats.
Identities = 60/103 (58%), Positives = 68/103 (66%), Gaps = 32/103 (31%)

Query   99    KQNKMLKKDKMPHFVWNNTFFIPGPBEEV-----D  126
KQNKMLKKDKMPHFVWNNTFFIPGPBEEV
Sbjct  260    KQNKMMKKDKMPHFVWNNTFFIPGPBEESRDKLENGAVNNADSQQGVPAPGQGQPQSACRE  319

Query  127    NDKEYLVLTkndldkankdkanRYFSPNPKVKLYFTKTVEE  169
+D++YL+KND DKANRDYKANRYFSPNPKVKL F+KTVEE
Sbjct  320    SDRDYILTLKNSKNDRKANKDKANRYFSPNPKVKLCFKSFTVEE  362

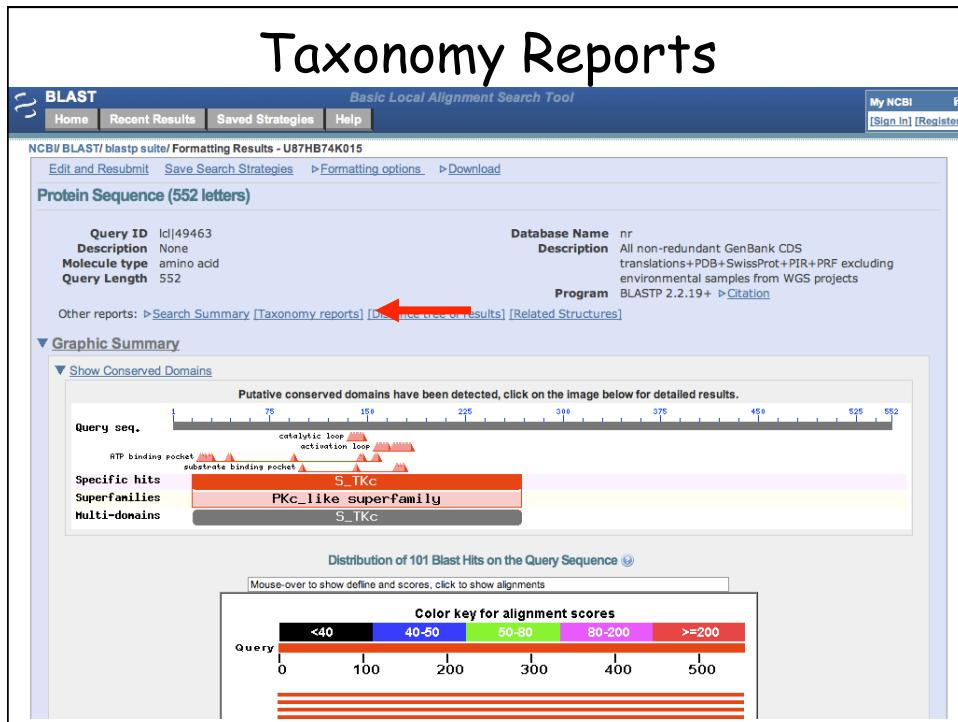
>gb|AAH93110.1| UG Ptenb protein [Danio rerio]
Length=289

Score = 197 bits (500), Expect = 2e-49, Method: Composition-based stats.
Identities = 95/99 (95%), Positives = 98/99 (98%), Gaps = 0/99 (0%)

Query   3     VSRNKRRYQEDGFDLDTIYIPNIIAMGFFPAERLEGRGVYRNNNIDVVVRFLDSKHKNHYKI  62
VSRNKRRYQEDGFDLDTIYIPNIIAMGFFPAERLEGRGVYRNNNIDVVVRFLDSKHKNHYKI
Sbjct   9     VSRNKRRYQEDGFDLDTIYIPNIIAMGFFPAERLEGRGVYRNNNIDVVVRFLDSKHKNHYKI  68

Query  63    NLCAERHYD TAKFNCRVAQYPFEDHNPPQLELIKPFQN  101
NLCAERHYD TAKFNCRVAQYPFEDHNPPQLELIKPF ++
Sbjct  69    NLCAERHYD TAKFNCRVAQYPFEDHNPPQLELIKPFCED  107

```



Lineage Report

Organism Report Taxonomy Report

Organism	Blast Name	Score	Number of Hits	Description
root		108		
· Simiiformes	primates	92		
· · Catarrini	primates	90		
· · · Hominoidea	primates	77		
· · · · Hominidae	primates	72		
· · · · · Homininae	primates	68		
· · · · · · Homo sapiens	primates	7282	53	Homo sapiens hits
· · · · · · Pan troglodytes	primates	3291	5	Pan troglodytes hits
· · · · · · Gorilla gorilla	primates	2900	2	Gorilla gorilla hits
· · · · · · Gorilla gorilla gorilla	primates	375	4	Gorilla gorilla gorilla hits
· · · · · · Pan paniscus	primates	370	4	Pan paniscus hits
· · · · · Pongo abelii	primates	2863	2	Pongo abelii hits
· · · · · Pongo pygmaeus	primates	2665	2	Pongo pygmaeus hits
· · · · · Nomascus leucogenys	primates	363	5	Nomascus leucogenys hits
· · · · · Chirocebus aethiops	primates	1517	1	Chirocebus aethiops hits
· · · · · Chirocebus sabaeus	primates	1088	2	Chirocebus sabaeus hits
· · · · · Papio anubis	primates	556	3	Papio anubis hits
· · · · · Macaca mulatta	primates	556	1	Macaca mulatta hits
· · · · · Pithecioides iheringi	primates	536	1	Pithecioides iheringi hits
· · · · · Rhinopithecus bieti	primates	492	2	Rhinopithecus bieti hits
· · · · · Macaca fascicularis	primates	327	1	Macaca fascicularis hits
· · · · · Macaca nemestrina	primates	322	2	Macaca nemestrina hits
· · · · · Aotus trivirgatus	primates	1371	1	Aotus trivirgatus hits
· · · · · Callithrix jacchus	primates	802	1	Callithrix jacchus hits
· · synthetic construct	other sequences	351	15	synthetic construct hits
· Human ORFeome Gateway entry vector	other sequences	348	1	Human ORFeome Gateway entry vector hits

Organism Report Taxonomy Report