Learn more ...

| | | Learn more ... |
|---|---|---|
| | Zebrafish | Learn more ... |
| | Insects | Learn more ... |
| | Nematodes | Learn more ... |
| | Plants | Learn more ... |
| | Fungi | Learn more ... |
| | Malaria | Learn more ... |
| | Other eukaryotes genomes | Learn more ... |
| | Microbial genomes | Learn more ... |
| | Trace MEGABLAST | Learn more ... |
| | VecScreen | Learn more ... |

-F F -e 20000 -W 2                                           -F F -e 2000 -W 7

**NOTE:**
**GenBank®**      **BLAST®**

blast-help@ncbi.nlm.nih.gov
info@ncbi.nlm.nih.gov

Back to top

## 4. Explanation for the program choices given in Tables 2.1 to 2.2

### 4.1 MEGABLAST is the tool of choice to identify a nucleotide sequence.

Batch Search

Back to top

**4.2 Discontiguous MEGABLAST** NEW **is better at finding nucleotide sequences similar, but not identical, to your nucleotide query.**

The BLAST nucleotide algorithm finds similar sequences by breakin the query into short subsequences called words. The program identifies the exact matches to the query words first (word hits). BLAST program then extends these word hits in multiple steps to generate the final gapped alignments.

One of the important parameters governing the sensitivity of BLAST searches is the length of the initial words, or word size as it is called. The most important reason that blastn is more sensitive than MEGABLAST is that it uses a shorter default word size. Because of this, blastn is better than MEGABLAST at finding alignments to related nucleotide sequences from other organisms. The word size is adjustable in blastn and can be reduced from the default value of 11 to a minimum of 7 to increase search sensitivity.

The search sensitivity can further improved by using the newly introduced discontiguous megablast page. This page uses an algorithm with the same name, which is similar to that reported by *Ma et.al.* Rather than requiring exact word matches as seeds for alignment extension, discontiguous megablast uses non-contiguous work within a longer window of template. In coding mode, the third base wobbling is taken into consideration by focusing on finding matches at the first and second codon positions while ignoring the mismatches in the third position. Searching in discontiguous MEGABLAST using the same word size is more sensitive and efficient than standard blastn using the same word size. For this reason, it is now the recommended tool for this type of search. Alternative non-coding patterns can also be specified if desired. Additional details on discontiguous are available at:

http://www.ncbi.nlm.nih.gov/blast/discontiguous.html
http://www.ncbi.nlm.nih.gov/Web/Newsltr/FallWinter02/blastlab.html

It is important to point out that nucleotide-nucleotide searches are not the best method for finding homologous protein coding regions in other organisms. That task is better accomplished by performing searches at the protein level, by direct protein-protein BLAST searches or by translated BLAST searches. This is because of the codon degeneracy, the greater information available in amino acid sequence, and the more sophisticated algorithm in protein-protein BLAST.

[Back to top]

**4.3 "Search for short nearly exact matches" is useful for primer or short nucleotide motif searches.**

Short sequences (less than 20 bases) will often not find any significant matches to the database entries under the standard nucleotide-nucleotide BLAST settings. The usual reasons for this are that the significance threshold governed by the expect value parameter is set too stringently and the default word size parameter is set too high.

You can adjust both the word size and the expect value on the standard BLAST pages to work with short sequences. However, we do provide a BLAST page with these values preset to give optimum results with short sequences. This page ("Search for short nearly exact matches") is linked under the nucleotide BLAST section of the main BLAST page.

| Table 4.3.1 Parameter settings for standard blastn and<br>"Search for short and nearly exact matches" | | | |
|---|---|---|---|
| Program | Word Size | DUST Filter Setting | Expect Value |
| Standard blastn | 11 | On | 10 |
| Search for short nearly exact matches | 7 | Off | 1000 |

A common use of this page is to check the specificity of PCR or hybridization. A useful way to check a pair of PCR primers is to first concatenate them by inserting string of 20 or more N's in between the two primers, and then search the concatenated pair as one sequence. Since BLAST looks for local

alignments and automatically searches both strands, there is no need to reverse complement one of the primers before doing the concatenation or the search.

The query sequence should contain no ambiguous bases. Consensus motifs with degenerate bases, such as KMKGSMGYYGSNNNNNNGCTYRGCWCSYTC or CNNGAANNTCCNNG will not work for this type of search.

[Back to top]

### 4.4 Use the Trace Archive BLAST page to search raw primary sequence trace files.

Trace data files are not official entries of the GenBank database and have no associated feature annotations. Despite this limitation, they are still a rich source of information, especially for organisms lacking a significant amount of regular mRNA or assembled genomic sequences. The sequences come from a variety of projects and sequencing strategies, including Whole Genome Shotgun (WGS), BAC end sequencing, and EST sequencing. The trace data are single pass sequencing reads not trimmed for quality or vector contamination. Their average lengths are between 500 to 700 bp.

A search against the Trace Archive can use MEGABLAST or discontiguous MEGABLAST. The former is better for indentifying exact matches in intra-species searches, such as looking for extra mRNA sequences or the genomic counterparts for a given gene, while the later is better for indentifying similar coding sequences from different species. Information on the Trace Archive is available from the Trace documentation page.

[Back to top]

### 4.5 Standard protein BLAST is designed for protein searches.

Standard protein-protein BLAST (blastp) is used for both identifying a query amino acid sequence and for finding similar sequences in protein databases. Like other BLAST programs, blastp is designed to find local regions of similarity. When sequence similarity spans the whole sequence, blastp will also report a global alignment, which is the preferred result for protein identification purposes.

For clear result in identification search, try taking off both "low complexity filter" and "Composition based statistics" function. Unlike nucleotide BLAST, there is no comparable MEGABLAST for protein searches, so batch search via the web is not possible.

### 4.6 PSI-BLAST is designed for more sensitive protein-protein similarity searches.

Position-Specific Iterated (PSI)-BLAST is the most sensitive BLAST program, making it useful for finding very distantly related proteins. Use PSI-BLAST when your standard protein-protein BLAST search either failed to find significant hits, or returned hits with descriptions such as "hypothetical protein" or "similar to..."

The first round of PSI-BLAST is a standard protein-protein BLAST search. The program builds a position-specific scoring matrix (PSSM or profile) from an alignment of the sequences returned with Expect values better (lower) than the inclusion threshold (default=0.005). In the second iteration the PSSM becomes the query in the search. Any new database hits below the inclusion threshold are included in the construction of the new PSSM. A PSI-BLAST search is said to have converged when no more new database sequences are added in subsequent iterations. You can add database hits that fall outside the inclusion threshold to your PSSM for the next round by checking the box next to the hit.

You can also save a PSSM created during a PSI-BLAST search of one database and use it to search a different database. To do this, change "Alignment" to "PSSM" in a pull-down menu in the Format section of a "Formatting BLAST" page (at any iteration after the first). Then format the search, copy the resulting PSSM and paste it into the PSSM window of a new PSI-BLAST search page.

[Back to top]

### 4.7 PHI-BLAST can do a restricted protein pattern search.

Pattern-Hit Initiated (PHI)-BLAST is designed to search for proteins that contain a pattern specified by the user, AND are similar to the query sequence in the vicinity of the pattern. This dual requirement is intended to reduce the number of database hits that contain the pattern, but are likely to have no true homology to the query.

To run PHI-BLAST, enter your query (which contains one or more instances of the pattern) into the "Search" box, and enter your pattern into the "PHI pattern" box in the "Options" section of the page. Patterns must follow the syntax conventions of PROSITE. Only one pattern can be used in a given search. The documentation on pattern syntax is at:

http://www.ncbi.nlm.nih.gov/blast/html/PHIsyntax.html

Sample query sequence, with modified defline and highlighted pattern occurance, and a sample pattern in ProSite format are given below:

```
>gi|4758958|ref|NP_004148.1| Human cAMP-dependent protein kinase
MSHIQIPPGLTELLQGYTVEVLRQQPPDLVEFAVEYFTRLREARAPASVLPAATPRQSLGHPPPEPGPDR
VADAKGDSESEEDEDLEVPVPSRFNRRVSVCAETYNPDEEEDTDPRVIHPKTDEQRCRLQEACKDILLF
KNLDQEQLSQVLDAMFERIVKADEHVIDQGDDGDNFYVIERGTYDILVTKDNQTRSVGQYDNRGSFGELA
LMYNTPRAATIVATSEGSLWGLDRVTFRRIIVKNNAKKRKMFESFIESVPLLKSLEVSERMKIVDVIGEK
IYKDGERIITQGEKADSFYIIESGEVSILIRSRTKSNKDGGNQEVEIARCHKGQYFGELALVTNKPRAAS
AYAVGDVKCLVMDVQAFERLLGPCMDIMKRNISHYEEQLVKMFGSSVDLGNLGQ

[LIVMF]-G-E-x-[GAS]-[LIVM]-x(5,11)-R-[STAQ]-A-x-[LIVMA]-x-[STACV].
```

[Back to top]

## 4.8 The protein "Search for short nearly exact matches" is optimized to find matches to a short peptide.

A short peptide (10-15mer or less) often will not find any significant matches to the database under the standard protein-protein BLAST settings. The usual reasons for this are that the significance threshold governed by the expect value parameter is set too stringently and the default word size parameter is set too high.

You could adjust both the word size and the expect value on the standard BLAST pages to make it work with short query sequences. NCBI provides a separate BLAST page with these values preset to optimize blastp searches with short query sequences. This page, "Search for short nearly exact matches", is available via a link under the Protein BLAST section of the BLAST home page. In addition, the more stringent PAM30 is used in lieu of BLOSUM62, and the composition-based statistics which takes the amino acid composition of the query sequence into account when calculating the score and significance of the alignments.

Composition based statistics takes the amino acid composition of the query and subject sequence into account when calculating the score and significance of the alignments. It can have a large effect on searches using queries with a biased amino acid composition. By definition, short peptides will have a biased compositions and should not be used with composition based statistics.

Due to the requirement that the query needs to be at least twice the word size, a query shorter than 5 residues is not recommended even though it can be as short as 4 residues when the word size is set to 2. In addition, since ambiguous residues break the query sequence, there should be no ambiguities in the query to ensure that the entire sequence can be used as seeds for the initial search.

| Table 4.8.1 Parameter settings for standard blastp and "Search for short and nearly exact matches" | | | | | |
|---|---|---|---|---|---|
| Program | Word Size | SEG Filter | Expect Value | Composition based Statistics | Score Matrix |
| Standard Protein Blast | 3 | On | 10 | On | BLOSUM62 |
| Search for short and nearly exact matches | 2 | Off | 20000 | Off | PAM30 |

[Back to top]

### 4.9 The "Nucleotide query - Protein db [blastx]" is useful for finding similar proteins to those encoded by a nucleotide query.

Translated BLAST services are useful when trying to find homologous proteins to a nucleotide coding region. Blastx compares the translation of the nucleotide query sequence to a protein database. Because blastx translates the query sequence in all six reading frames and provides combined significance statistics for hits to different frames, it is particularly useful when the reading frame of the query sequence is unknown or it contains errors that may lead to frame shifts or other coding errors. Thus blastx is often the first analysis performed with a newly determined nucleotide sequence and is used extensively in analyzing EST sequences.

### 4.10 The "Protein query - Translated db [tblastn]" search is useful for finding protein homologs in unannotated nucleotide data.

A tblastn search allows you to compare a protein sequence to the six-frame translations of a nucleotide database. It can be a very productive way of finding homologous protein coding regions in unannotated nucleotide sequences such as expressed sequence tags (ESTs) and draft genome records (HTG), located in the BLAST databases est and htgs, respectively.

ESTs are short, single-read cDNA sequences. They comprise the largest pool of sequence data for many organisms and contain portions of transcripts from many uncharacterized genes. Since ESTs have no annotated coding sequences, there are no corresponding protein translations in the BLAST protein databases. Hence a tblastn search is the only way to search for these potential coding regions at the protein level. The HTG sequences, draft sequences from various genome projects or large genomic clones, are another large source of unannotated coding regions.

Like all translating searches, the tblastn search is especially suited to working with error prone data like ESTs and draft genomic sequences from HTG because it combines BLAST statistics for hits to multiple reading frames and thus is robust to frame shifts introduced by sequencing error.

[Back to top]

### 4.11 The "Nucleotide query - Translated db [tblastx]" is useful for identifying novel genes in error prone query sequences.

tblastx takes a nucleotide query sequence, translates it in all six frames, and compares those translations to the database sequences dynamically translated in all six frames. This effectively performs a more sensitive blastp search without doing the manual translation.

tblastx gets around the potential frame-shift and ambiguities that may prevent certain open reading frames from being detected. This is very useful in identifying potential proteins encoded by single pass read ESTs. In addition, it can be a good tool for identifying novel genes.

This type of search is computationally intensive and searches with large genomic queries are not recommended. The best way to do this is to install standalone blast and perform the search locally. For more information on standalone blast, please read the documents for formatdb and standalone BLAST at:

ftp://ftp.ncbi.nih.gov/blast/documents/formatdb.txt
ftp://ftp.ncbi.nih.gov/blast/documents/netblast.txt

[Back to top]

### 4.12 The Conserved Domain Database (CDD) search service uses RPS-BLAST to identify conserved protein domains.

Reverse Position Specific BLAST (RPS-BLAST) is a more sensitive way of identifying conserved domains in proteins than standard BLAST searching. It compares a protein sequence against a

database of position specific scoring matrices (PSSMs). The PSSMs used in CDD search capture the substitution frequencies at each position in the multiple sequence alignments of recognized conserved domains. The conserved domain alignments are from the NCBI's CDD, which contains alignments from protein domain databases: Smart, Pfam, COG, KOG, and LOAD. For additional information, refer to CDD help document at: http://www.ncbi.nlm.nih.gov/Structure/cdd/cdd_help.shtml

[Back to top]

### 4.13 The Conserved Domain Architecture Retrieval Tool (CDART) explores the domain architectures of proteins.

CDART allows you to examine the domain structure of all proteins in the default BLAST protein database. The CDART tool first searches a query sequence for the presence of conserved domains using RPS-BLAST. It then allows you to retrieve proteins that share one or more protein domains in common with your query. Because CDART relies on RPS-BLAST, these searches are more sensitive than ordinary BLAST searches.

If the query does not contain any conserved domains, CDART will not report any result.

[Back to top]

### 5. Explanation for Program Choices Given in Table 3.3

### 5.1 "BLAST 2 Sequences" is designed for direct comparison of two sequences.

This program takes two input sequences and compares them directly. "Aligning Two Sequences" regards the second sequence as the database. Unlike the other BLAST programs, there is no need to format the database sequence in any special way. Since translated BLAST programs are incorporated in this program, the second sequence can be of different type so long as an appropriate BLAST program is selected. Appropriate query/program combination is listed in the table below.

| Table 5.1.1 Appropriate Query/Program Combinations for "BLAST 2 Sequences" | | |
|---|---|---|
| First Query | Second Query | Program to Use |
| Nucleotide | Nucleotide | blastn, megablast, or tblastx |
| Nucleotide | Protein | blastx |
| Protein | Nucleotide | tblastn |
| Protein | Protein | blastp |

If the database sequence or second query is present in an NCBI database, using the GI/Accession instead of the FASTA sequence allows the program to incorporate the translation and other sequence features, found in that record, into the final alignment making it more informative.

[Back to top]

### 5.2 The Human Genome BLAST page is for comparing a query against the NCBI's assembly of human genome, plus its derivative and related databases.

Like other BLAST search pages in this Genomes section, this page provides a centralized page to access specialized databases. In this case, the databases are the current NCBI human genome build and those derived from or related to it. All flavors of BLAST, except tblastx, are available with MEGABLAST set as default. Default filters are DUST and human repeats. The BLAST output links directly to the Human Genome MapViewer, where database hits can be visualized/analyzed in a genomic context, such as their relationship to other map elements like Transcript, SNPs, and Gene_seq. Names for the databases are being standardized, refer to Table 5.3.1 for details on database content.

[Back to top]

Gene Trap

Single pass sequence reads from numerous Zebrafish cDNA libraries.

Back to top

**5.9** **Use the Plants genome BLAST pages to search against green plant genomes.**

| **Table 5.9.1 Plants Genome BLAST Database Content** | |
|---|---|
| | |
| | |
| | |
| | |
| | |
| | |
| | |
| | |
| | |
| | |
| | |

Back to top

**5.10** **The Nematode BLAST page.**

**5.11 Yeasts Genome BLAST page provides access to multiple yeast genomes.**

This page provides access to different yeast genomes and their protein translations. Sequences of other yeast strains are also available in addition to that for *Saccharomyces cerevisiae* and *Schizosaccharomyces pombe*. The databases for the two well known strains can be searched individually or together. Hits are linked to MapView. All flavors of BLAST, with the exception of tblastx, are available.

| Table 5.11.1 Database list for Yeasts page | |
|---|---|
| **Organism** | **Sequence available** |
| Schizosaccharomyces pombe | Genome, mRNAs, and Proteins |
| Saccharomyces cerevisiae | Genome, mRNAs, and Proteins |
| Saccharomyces paradoxus NRRL Y-17217 | Nucleotide only |
| Saccharomyces mikatae IFO1815(MIT) | Nucleotide only |
| Saccharomyces mikatae IFO1815(WashU) | Nucleotide only |
| Saccharomyces bayanus MCYC623(MIT) | Nucleotide only |
| Saccharomyces bayanus MCYC623(WashU) | Nucleotide only |
| Saccharomyces castellii NRRL Y-12630 | Nucleotide only |
| Saccharomyces kluyveri NRRL Y-12651 | Nucleotide only |
| Saccharomyces kudriavzevii IFO 1802 | Nucleotide only |
| All YEAST Genomes | Genomes, mRNAs, and Proteins |
| Neurospora crassa | Genome, mRNAs, and proteins |
| Magnaporthe grisea | Genome, mRNAs, and proteins |
| Aspergillus nidulans | Genome, mRNAs, and proteins |
| All Species | Genomes, mRNAs, and Proteins |

[Back to top]

**5.12 Use the Flies BLAST page to search the *Anopheles gambiae*, *Drosophila melanogaster*, and *Apis mellifera* genomes.**

This page provides access to the genome scaffold of Anopheles gambiae (mosquito) and Drosophila melanogaster (fruitfly) chromosomes. The proteins translated from the genome annotation are also available. The data available for Anopheles gambiae is from a NIAID publicly funded project with the sequencing and assembly performed by Celera Corporation. The data for Drosophila melanogaster come from FlyBase. Hits are linked to corresponding MapViewer pages, providing additional information. Apis mellifera sequences were added recently, which has no MapViewer link yet.

[Back to top]

**5.13 The VecScreen page is for identifying vector sequence contamination in a query sequence.**

VecScreen, under special section, is a rapid screening tool that checks the query sequence against a non-redundant vector database, UniVec, which contains one copy of every unique sequence segment from a large number of cloning vectors. In addition, UniVec contains sequences for adapters, linkers, stuffers, and primers that are commonly used in the cloning and manipulation of cDNA or genomic DNA. Detailed information on UniVec is at: http://www.ncbi.nlm.nih.gov/VecScreen/UniVec.html.

This page is generally used to screen for vector contamination in sequences before their submission to GenBank. The color-coded graphics in the result page makes the result easy to understand.

**6. Appendix**

| | | | |
|---|---|---|---|
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |

| **Table 6.2.2 Single Letter Amino Acid Code** | | | |
|---|---|---|---|
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |

Back to top

**6.3 Other alternative means for batch BLAST searches.**

| Table 6.3 Alternatives Means for Batch BLAST Searches | | | |
|---|---|---|---|
| Alternatives | Pros | Cons | Links |
| blastcl3 | • No database maintainance<br>• Simple to set up | • server/network fluctuation<br>• Relative low throughput<br>• No graphic output | document<br>program |
| URL-API | • Versatility<br>• No database maintainance | • Custom scripts needed<br>• Load restrictions | document |
| Standalone BLAST | • No server fluctuation<br>• Custom databases<br>• High throughput | • Needs database update<br>• No graphic | document<br>program |

[Back to top]

Updated on Sat Feb 14 14:12:00 2004