

# BLAST

Getting the most from your cycles.

Tom Madden  
NCBI/NLM/NIH  
madden@ncbi.nlm.nih.gov



---

---

---

---

---

---

---

---

## The BLAST algorithm

---

---

---

---

---

---

---

---

## What is BLAST?

- **Basic Local Alignment Search Tool**
- Calculates similarity for biological sequences.
- Produces local alignments: only a portion of each sequence must be aligned.
- Uses statistical theory to determine if a match might have occurred by chance.

---

---

---

---

---

---

---

---

The BLAST family of programs allows all combinations of DNA or protein query sequences with searches against DNA or protein databases:

Protein-protein (blastp): compares an amino acid sequence against a protein sequence database.

Nucl.-nucl (blastn): compares a nucleotide query sequence against a nucleotide sequence database (in general optimized for speed, not sensitivity).

Translated nucl.-protein (blastx): compares the six-frame conceptual translation products of a nucleotide query against a protein sequence database.

Protein-translated nucl (tblastn): compares a protein query sequence against a sequence database dynamically translated in all six reading frames (useful for searching proteins against EST's).

Translated nucl-translated nucl. (tblastx): compares the six frame translation of a nucleotide query sequence against the six-frame translations of a nucleotide sequence database.

---

---

---

---

---

---

---

---

## BLAST is a heuristic.

- A lookup table is made of all the “words” (short subsequences) in the query sequence. In many types of searches “neighboring” words are included.
- The database is scanned for matching words (“hot spots”).
- Gapped and un-gapped extensions are initiated from these matches.

---

---

---

---

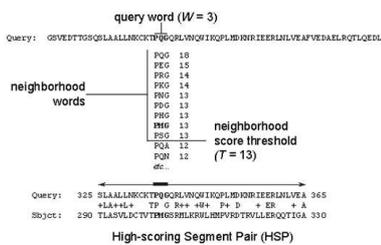
---

---

---

---

## The BLAST Search Algorithm




---

---

---

---

---

---

---

---





There are drawbacks to parsing the BLAST report and Hit-table.

- No way to automatically check for truncated output.
- No way to rigorously check for syntax changes in the output.

---

---

---

---

---

---

---

Structured output allows automatic and rigorous checks for syntax errors and changes.

---

---

---

---

---

---

---

### Abstract Syntax Notation 1 (ASN.1)

- Is an International Standards Organization (ISO) standard for describing structured data and reliably encoding it.
- Used extensively in the telecommunications industry.
- Both a binary and a text format.
- NCBI data model is written in ASN.1.
- Asntool can produce C object loaders from an ASN.1 specification.
- Datatool can produce C++ classes from an ASN.1 specification.

---

---

---

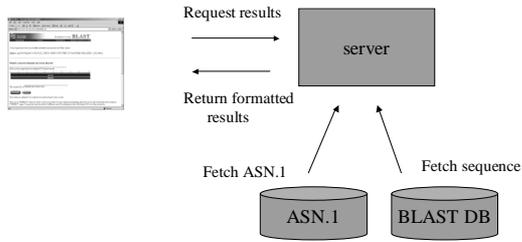
---

---

---

---

## ASN.1 is used for the NCBI BLAST Web page.




---

---

---

---

---

---

---

---

## Different reports can be produced from the ASN.1 of one search.

---

---

---

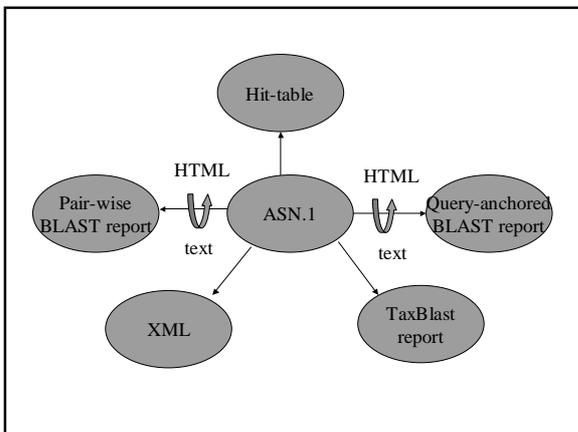
---

---

---

---

---




---

---

---

---

---

---

---

---

## The BLAST ASN.1 (“SeqAlign”) contains:

- Start, stop, and gap information (zero-offset).
- Score, bit-score, expect-value.
- Sequence identifiers.
- Strand information.

```

ASN1_SEQUENCE {
  SeqAlign {
    id Object-id OPTIONAL,
    value CHOICE {
      real REAL,
      int INTEGER }
  }
}

```

---

---

---

---

---

---

---

---

## Three flavors of Seq-Align, Score-block(s) plus one of:

- Dense-diag: series of unconnected diagonals. No coordinate “stretching” (e.g., cannot be used for protein-nucl. alignments). Used for ungapped BLASTN/BLASTP.
- Dense-seg: describes an alignment containing many segments. No coordinate “stretching”. Used for gapped BLASTN/BLASTP.
- Std-seg: a collection of locations. No restriction on stretching of coordinates. Used for gapped/ungapped translating searches. Generic.

---

---

---

---

---

---

---

---

## Score Block

```

Score ::= SEQUENCE {
  id Object-id OPTIONAL, -- identifies Score type
  value CHOICE {
    real REAL, -- floating point value
    int INTEGER } -- integer
}

```

SEQUENCE is an ordered list of elements, each of which is an ASN.1 type. Required unless DEFAULT or OPTIONAL.

---

---

---

---

---

---

---

---







Specification (i.e., “data model”) issues should not be confused with the question about whether to use ASN.1 or XML.

---

---

---

---

---

---

---

Structured output is not a panacea.

- Design issues must still be addressed.
- Semantic issues still exist, e.g. is a start/stop value zero-offset or one-offset.
- Data issues still exist, e.g., is the correct sequence shown, are the offsets correct, was the DNA translated with the correct genetic code?

---

---

---

---

---

---

---

IMPROVING BLAST THROUGHPUT

---

---

---

---

---

---

---

Use megablast to align very similar sequences.

- Best if alignments between query and target will be 97-99% identical.
- Word-size: 28; an exact match of 28-31 bases required to initiate extensions.
- A greedy gapped alignment routine with non-affine gapping (constant cost per insertion/deletion) used to perform extensions
- Use for aligning sequences from the same organism.

---

---

---

---

---

---

---

---

Example: search u93237 (Human MEN1 gene) vs human EST's with filtering for low-complexity and human repeats and expect value 1.0e-6

MEGABLAST:

- word size: 28
- run time: 19 seconds
- 469 alignments found
- 345 alignments more than 98% identical

BLASTN

- word size: 11
- run time: 148 seconds
- 491 alignments found
- 359 alignments more than 98% identical

---

---

---

---

---

---

---

---

BLASTN alignment:

Score = 224 bits (112), Expect = 4e-55  
 Identities = 112/112 (100%)  
 Strand = Plus / Plus

Query: 6647 gggagctctacaaagggtctcttgaagtagccaaatgatgctatccccaaacctgctgaagga 6706  
 |||  
 Sbjct: 79 gggagctctacaaagggtctcttgaagtagccaaatgatgctatccccaaacctgctgaagga 138

Query: 6707 gggagctctacaaagggtctcttgaagtagccaaatgatgctatccccaaacctgctgaagga 6759  
 |||  
 Sbjct: 139 gggagctctacaaagggtctcttgaagtagccaaatgatgctatccccaaacctgctgaagga 191

MEGABLAST alignment:

Score = 226 bits (113), Expect = 4e-56  
 Identities = 113/113 (99%), Gaps = 6/142 (4%)  
 Strand = Plus / Plus

Query: 6623 ctacaa-ctac-tgcgggaa-gac-ga-ggagatctacaaagggtctcttgaagtagcc 6677  
 |||  
 Sbjct: 51 ctacaaagctacag-cgggaatgacggatgagatctacaaagggtctcttgaagtagcc 109

Query: 6678 aatgatctacccccaaacctgctgaagtagccaaatgatgctatccccaaacctgctgaagga 6737  
 |||  
 Sbjct: 130 aatgatctacccccaaacctgctgaagtagccaaatgatgctatccccaaacctgctgaagga 169

Query: 6738 gggagctctacaaagggtctcttgaagtagccaaatgatgctatccccaaacctgctgaagga 6759  
 |||  
 Sbjct: 170 gggagctctacaaagggtctcttgaagtagccaaatgatgctatccccaaacctgctgaagga 191

---

---

---

---

---

---

---

---

**Increasing the threshold for blast[px]/tblast[nx] speeds up search.**

- these programs use exact and “neighboring” (three-letter) words as initial hits.
- increasing threshold decreases the number of neighboring words.
- fewer neighboring words mean fewer extensions.
- if the threshold is a high value (e.g., 100) only exact matches are used.
- more subtle alignments will be missed.

---

---

---

---

---

---

---

---

**Example: blastx search of NT\_078011 (41887 bases human contig) against nr.**

Threshold	Time (seconds)	Alignments with expect $\leq 1.0e-6$
12	4050	962
13	2452	955
14	1581	954
15	1139	947
17	770	916
100	658	879

---

---

---

---

---

---

---

---

```
>refINP_524599.11 Serotonin receptor 7 CG12073-PA [Drosophila melanogaster] ...
Length = 564
Score = 63.2 bits (152), Expect = 4e-07
Identities = 28/70 (40%), Positives = 37/70 (52%)
Frame = -3
Query: 41255 ARGKSFTFLVAVMGMVFLCHPFFFFSYSLYGICREACQVGPLFKFFFIWIGYCNSSLNP 41076
      A+ K + L ++M F +CN PFF + E VP L F W+GY NS LNP
Sbjct: 447 AKEKKASTLGLIMSFTVCLPFFLILALIRPF--ETMHWPRSLSSLFLMLGYANSLNP 504
Query: 41075 VITYVFNQDF 41046
      +IY N+DF
Sbjct: 505 IITYATLNRDF 514
```

This alignment is found with thresholds 12,13,14,15, and 17. It is not found if only exact matches are used (threshold = 100).

The default mode of blast[px] and tblast[nx] requires two hits on the same diagonal to initiate an extension.

---

---

---

---

---

---

---

---

**Do as little work (with the database)  
as possible.**

- Scanning the database and/or reading it from disk can be a significant portion of some searches.
- Use the “concatenation” feature of megablast to concatenate multiple queries and scan the database only once for all queries.
- The OS will cache recently used files. Make use of this by grouping together queries for one database before searching another one.
- Buy more memory if you cannot fit a database into memory.

---

---

---

---

---

---

---

---

**MEGABLAST concatenates queries**

- searching 50 human EST’s (total of 15,700 bases) against the human EST database took 23.5 seconds (668 bases/second).
- searching one human EST (gi|272208, 211 bases) took 13 seconds (16 bases/second).
- Concatenation minimizes time spent scanning the database.
- The more stringent the search, the more the savings.

---

---

---

---

---

---

---

---

**Three different strategies to search  
three EST’s against the nt and est  
databases.**

- 26 minutes (wall clock time) if searches are grouped by query.
- 10 minutes if searches are grouped by database.
- 7.5 minutes if megablast concatenates the queries.

---

---

---

---

---

---

---

---

# BLAST DATABASES

---

---

---

---

---

---

---

---

- ### BLAST databases:
- can be produced with stand-alone formatdb and a FASTA file.
  - are always (?) produced with the formatdb API (e.g., stand-alone formatdb).
  - are almost always read with the readdb API (recommended).
  - pack nucleotide sequences 4-to-1.
  - are architecture independent.

---

---

---

---

---

---

---

---

### The (physical) BLAST databases comprise files in binary format.

pin or nin	Index into sequence and header files
psq or nsq	Sequence data
phr or nhr	Sequence identifier, definition, taxonomic information, etc. stored as binary ASN.1
pni or nni *	ISAM index file for GI identifiers
pnd or nnd *	ISAM data file for GI identifiers
psi or nsi *	ISAM index file for other identifiers
psd or nsd *	ISAM data file for other identifiers

\* only created if "-o" option used with formatdb.

---

---

---

---

---

---

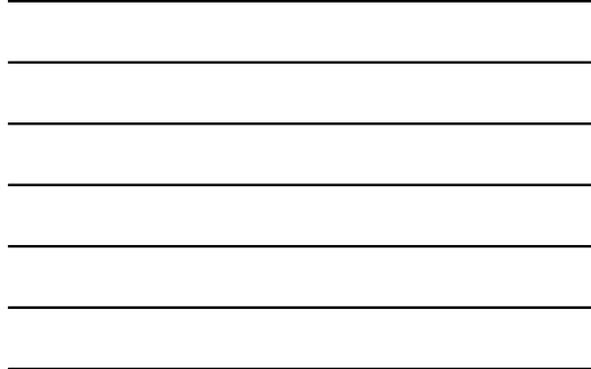
---

---

## ASN.1 spec. used for header files

-- one Blast-def-line-set for each entry  
 Blast-def-line-set ::= SEQUENCE OF Blast-def-line

```
Blast-def-line ::= SEQUENCE {
  title VisibleString OPTIONAL,      -- simple title
  seqid SEQUENCE OF Seq-id,         -- Regular NCBI Seq-id
  taxid INTEGER OPTIONAL,           -- taxonomy id
  memberships SEQUENCE OF INTEGER OPTIONAL, -- bit arrays
  links SEQUENCE OF INTEGER OPTIONAL, -- bit arrays
  other-info SEQUENCE OF INTEGER OPTIONAL -- future use
}
```



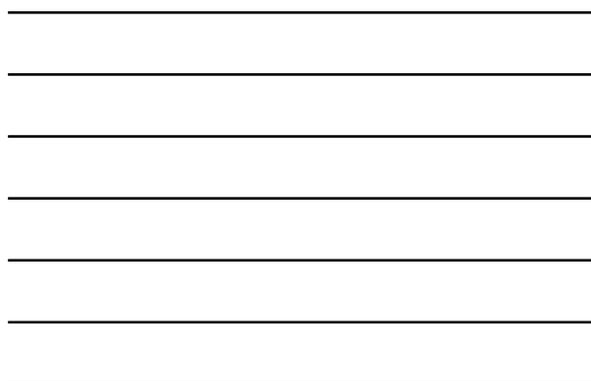
## Use asntool to view header file content.

```
Blast-def-line-set ::= {
  title "INA mismatch repair protein Mh1 (MutL protein homolog 1)" .
  seqid {
    gi 1576997 .
    swissprot {
      name "MLH_MOUSE" .
      accession "Q9J931" } } .
  taxid 10090 .
  memberships {
    1 } .
  links {
    1 } .
  other-info {
    157699 } } .
  {
    title "MutL homolog 1 protein (Mus musculus)" .
    seqid {
      gi 7595954 .
      swissprot {
        name "MF250844.1" .
        accession "M9F64514" .
        version 1 } } .
    taxid 10090 .
    links {
      1 } .
    other-info {
      157699 } } }
}
```

Annotations:

- GI number assigned by NCBI (points to gi 1576997)
- Swissprot locus name and accession (points to name "MLH\_MOUSE" and accession "Q9J931")
- Tax ID for "house mouse" (points to taxid 10090)
- Is in swissprot subset of nr (points to swissprot block)
- Has information in LocusLink (points to links { 1 })
- Used for Protein-Identifier-Group (points to other-info { 157699 })

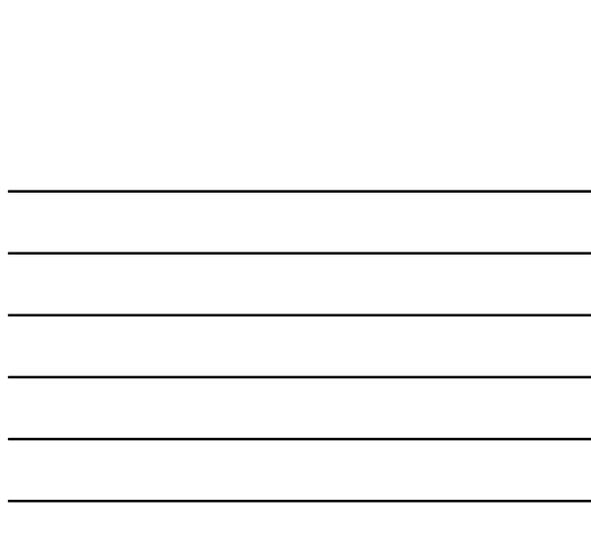
asntool -m fastadl.asn -M asn.all -d nr.phr -t Blast-def-line-set -p stdout



## Asntool can also produce header files in XML.

```
<Blast-def-line-set>
  <Blast-def-line>
    <Blast-def-line_title>INA mismatch repair protein Mh1 (MutL protein homolog 1)</Blast-def-line_title>
    <Blast-def-line_seqid>
      <Blast-def-line_seqid_gi>1576997</Blast-def-line_seqid_gi>
      <Blast-def-line_seqid_swissprot>
        <Blast-def-line_seqid_swissprot_name>MLH_MOUSE</Blast-def-line_seqid_swissprot_name>
        <Blast-def-line_seqid_swissprot_accession>Q9J931</Blast-def-line_seqid_swissprot_accession>
      </Blast-def-line_seqid_swissprot>
    </Blast-def-line_seqid>
    <Blast-def-line_taxid>10090</Blast-def-line_taxid>
    <Blast-def-line_memberships>
      <Blast-def-line_memberships_membership>1</Blast-def-line_memberships_membership>
    </Blast-def-line_memberships>
    <Blast-def-line_links>
      <Blast-def-line_links_link>1</Blast-def-line_links_link>
    </Blast-def-line_links>
    <Blast-def-line_other-info>
      <Blast-def-line_other-info_other-info_1>157699</Blast-def-line_other-info_other-info_1>
    </Blast-def-line_other-info>
  </Blast-def-line>
  <Blast-def-line>
    <Blast-def-line_title>MutL homolog 1 protein (Mus musculus)</Blast-def-line_title>
    <Blast-def-line_seqid>
      <Blast-def-line_seqid_gi>7595954</Blast-def-line_seqid_gi>
      <Blast-def-line_seqid_swissprot>
        <Blast-def-line_seqid_swissprot_name>MF250844.1</Blast-def-line_seqid_swissprot_name>
        <Blast-def-line_seqid_swissprot_accession>M9F64514</Blast-def-line_seqid_swissprot_accession>
        <Blast-def-line_seqid_swissprot_version>1</Blast-def-line_seqid_swissprot_version>
      </Blast-def-line_seqid_swissprot>
    </Blast-def-line_seqid>
    <Blast-def-line_taxid>10090</Blast-def-line_taxid>
    <Blast-def-line_memberships>
      <Blast-def-line_memberships_membership>1</Blast-def-line_memberships_membership>
    </Blast-def-line_memberships>
    <Blast-def-line_links>
      <Blast-def-line_links_link>1</Blast-def-line_links_link>
    </Blast-def-line_links>
    <Blast-def-line_other-info>
      <Blast-def-line_other-info_other-info_1>157699</Blast-def-line_other-info_other-info_1>
    </Blast-def-line_other-info>
  </Blast-def-line>
</Blast-def-line-set>
```

asntool -m fastadl.asn -M asn.all -d nr.phr -t Blast-def-line-set -x stdout



## Alias files

- Virtual databases uses alias files to instruct BLAST which physical database(s) to search.
- Alias files have extensions “.nal” or “.pal”.
- Alias files can specify multiple databases.
- The search can be limited by a list of GI's or ordinal id's (sequences in BLAST database) specified in the alias file.
- Alias files hide physical databases with the same (root) name.

---

---

---

---

---

---

---

---

## Alias file using GI list

```
Alias file created Mon Nov 25 13:25:47 2002
*
*
*
TITLE Mus musculus RIN sequences
DBLIST ../newest_blast/blast/rin
GILIST rin.rn.gil
HEED 46392
LENGTH 4726730
```

Physical database to search

Limits search to this list of GI's.

Virtual database statistics

The GI list can be either text or binary; formatdb can produce a binary GI list from a text one.

Formatdb can be used to produce this alias file with statistics.

---

---

---

---

---

---

---

---

## Alias file using ordinal ID list.

```
Alias file generated by formatdb
* Date created: Fri May 10 03:26:48 2003
*
*
*
TITLE NCBI Truncated Reference Sequences
DBLIST rin
GILIST seqs_rna_00.gil
LENGTH 39102645
HEED 393762
MWORDID 1761901
* end of file 111
```

Physical database to search

Specifies subset to search.

Virtual database statistics

---

---

---

---

---

---

---

---



## ASN.1 RESOURCES

- The Open Book : A Practical Perspective on OSI by Marshall T. Rose (Prentice Hall).
- OSS Nokalva Web site:  
<http://www.oss.com/asn1/overview.html>
- NCBI toolkit documentation on ASN.1:  
<http://www.ncbi.nlm.nih.gov/IEB/ToolBox/SDKDOCS/ASNLIB.HTML>

---

---

---

---

---

---

---

---

## Email addresses

- General questions about running BLAST:  
[blast-help@ncbi.nlm.nih.gov](mailto:blast-help@ncbi.nlm.nih.gov)
- Questions about compiling the toolkit and requests for hard-copy of documentation:  
[toolbox@ncbi.nlm.nih.gov](mailto:toolbox@ncbi.nlm.nih.gov)

---

---

---

---

---

---

---

---

## Selected BLAST References

- Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool. *J Mol Biol.* 1990 Oct 5;215(3):403-10.
- Altschul SF, Madden TL, Schaffer AA, Zhang J, Zhang Z, Miller W, Lipman DJ. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res.* 1997 Sep 1;25(17):3389-402.
- Altschul SF. Amino acid substitution matrices from an information theoretic perspective. *J Mol Biol.* 1991 Jun 5;219(3):555-65.
- Altschul SF, Boguski MS, Gish W, Wootton JC. Issues in searching molecular sequence databases. *Nat Genet.* 1994 Feb;6(2):119-29.

---

---

---

---

---

---

---

---

(Some of the) People (currently)  
working on BLAST

- Kevin Bealer
- Christian Camacho
- George Coulouris
- Ilya Dondoshansky
- Tom Madden
- Yuri Merezhuk
- Yan Raytselis
- Jian Ye
- Richa Agarwala
- Stephen Altschul
- Peter Cooper
- Susan Dombrowski
- David Lipman
- Wayne Matten
- Scott McGinnis
- Alexander Morgulis
- Alejandro Schaffer
- Tao Tao
- David Wheeler

---

---

---

---

---

---

---

---