# How to run BLAST locally

Jarek Bryk

Huddersfield, 29th March 2017

# Basic Local Alignment Search Tool

Stephen F. Altschul[1], Warren Gish[1], Webb Miller[2]
Eugene W. Myers[3] and David J. Lipman[1]

[1]*National Center for Biotechnology Information*
*National Library of Medicine, National Institutes of Health*
*Bethesda, MD 20894, U.S.A.*

[2]*Department of Computer Science*
*The Pennsylvania State University, University Park, PA 16802, U.S.A.*

[3]*Department of Computer Science*
*University of Arizona, Tucson, AZ 85721, U.S.A.*

A new approach to rapid sequence comparison, basic local alignment search tool (BLAST), directly approximates alignments that optimize a measure of local similarity, the maximal segment pair (MSP) score. Recent mathematical results on the stochastic properties of MSP scores allow an analysis of the performance of this method as well as the statistical significance of alignments it generates. The basic algorithm is simple and robust; it can be implemented in a number of ways and applied in a variety of contexts including straight-forward DNA and protein sequence database searches, motif searches, gene identification searches, and in the analysis of multiple regions of similarity in long DNA sequences. In addition to its flexibility and tractability to mathematical analysis, BLAST is an order of magnitude faster than existing sequence comparison tools of comparable sensitivity.

# blog.thegrandlocus.com/2014/06/once-upon-a-blast

(the origin story of BLAST)

Support Center Home    Write to the Help Desk

# NCBI Support Center

| | Search |
|---|---|

All Help Topics › BLAST

# How many sequences can I submit for a BLAST search at one time?

- Print
- Email this page
- + Share

There is no hard limit on the number of sequences that you can enter into the **Enter Query Sequence** input window or upload in a file to the BLAST® web interface. The total number of sequences that you can search with will change depending on the nature of the query sequences and the database being searched. Long sequences and large numbers of sequences in a search increase the possibility that the search will not complete due to a SIGXCPU error. If you are planning a large BLAST project you should download BLAST+ software and databases or run BLAST searches at a cloud provider.

Need More Help?    **Contact Us**

## Rate the article

Rating: ★★★★☆ 4 Votes

Was this answer helpful? 👍 👎    **Submit**

NIH  U.S. National Library of Medicine    NCBI    National Center for Biotechnology Information

Support Center Home    Write to the Help Desk

## NCBI Support Center

Search

All Help Topics > BLAST

# How many sequences can I submit for a BLAST search at one time?

There is no hard limit on the number of sequences that you can enter into the **Enter Query Sequence** input window or upload in a file to the BLAST® web interface.

you are planning a large BLAST project you should download BLAST+ software and databases or run BLAST searches at a cloud provider.

Need More Help?    **Contact Us**
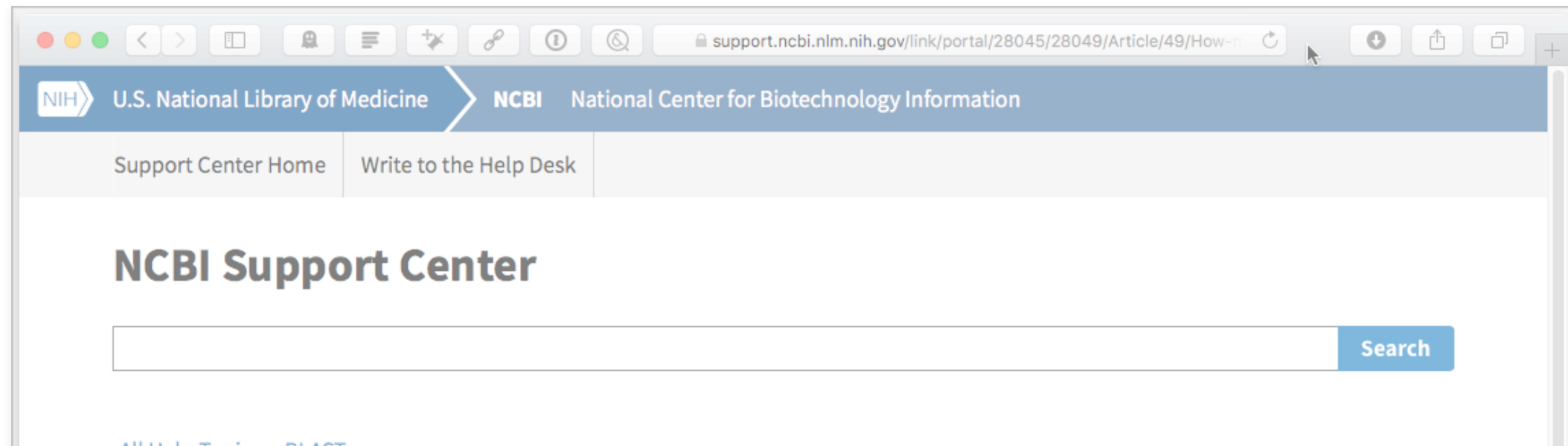
**Rate the article**

Rating: ★★★☆☆ 4 Votes

Was this answer helpful? 👍 👎    Submit
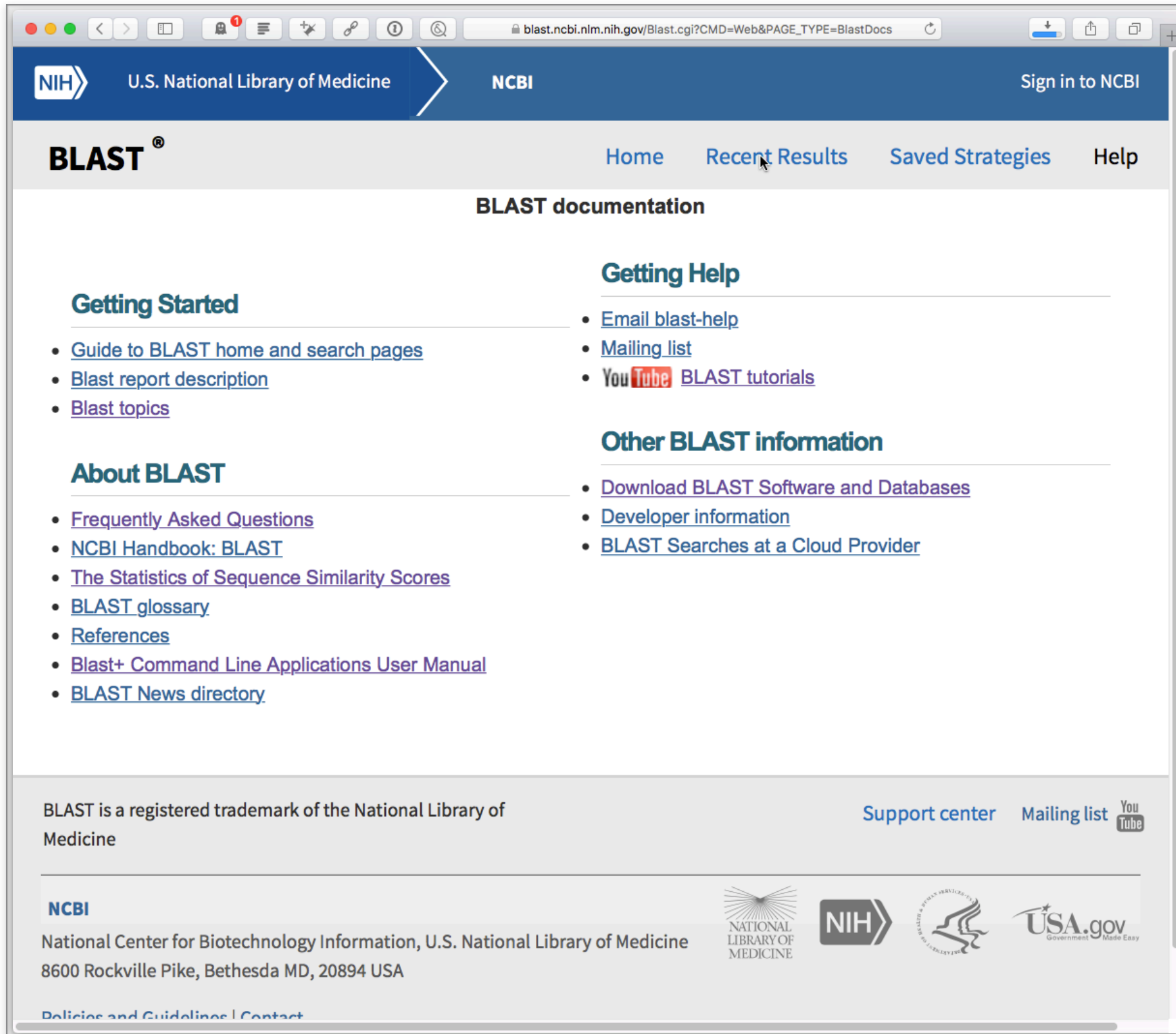
blast.ncbi.nlm.nih.gov/Blast.cgi?CMD=Web&PAGE_TYPE=BlastDocs

Install BLAST+

Download or make the database

Prepare your query

Run BLAST with desired options

NIH ❯ **U.S. National Library of Medicine**    **NCBI**    Sign in to NCBI

## BLAST ®

Home        Recent Results        Saved Strategies        Help

## Download BLAST Software and Databases

### BLAST+ executables

BLAST+ is a suite of command-line tools to run BLAST. For details, please see the BLAST+ user manual, the BLAST Help manual, the BLAST releases notes, and the article in BMC Bioinformatics (PubMed link). BLAST+ is the most current version of BLAST and is the only supported version.

Installers and source code are available from ftp://ftp.ncbi.nlm.nih.gov/blast/executables/blast+/LATEST/ .

See our versioning policy.

The BLAST+ suite is the currently supported package. The older C toolkit executables are no longer supported.

### Databases

BLAST databases are updated daily and may be downloaded via FTP from ftp://ftp.ncbi.nlm.nih.gov/blast/db/. Database sets may be retrieved automatically with update_blastdb.pl, which is part of the BLAST+ suite. Please refer to the BLAST database documentation for more details.

Support center        Mailing list        YouTube

**NCBI**

National Center for Biotechnology Information, U.S. National Library of Medicine

8600 Rockville Pike, Bethesda MD, 20894 USA

Policies and Guidelines | Contact

# www.ncbi.nlm.nih.gov/taxonomy

# www.ncbi.nlm.nih.gov/sites/batchentrez

blastn -task megablast -db db/refseq_rna
-query test_query.txt -dust no -max_target_seqs 1
-outfmt "6 qseqid sseqid evalue pident stitle"

```
blastn -task megablast -db db/refseq_rna
-query test_query.txt -dust no -max_target_seqs 1
-outfmt "6 qseqid sseqid evalue pident stitle"
```

blastn -task megablast -db db/refseq_rna
-query test_query.txt -dust no -max_target_seqs 1
-outfmt "6 qseqid sseqid evalue pident stitle"

blastn -task megablast -db db/refseq_rna

location of the query (multi FASTA)

-query test_query.txt -dust no -max_target_seqs 1

-outfmt "6 qseqid sseqid evalue pident stitle"

blastn -task megablast -db db/refseq_rna

mask low complexity regions in query

-query test_query.txt -dust no -max_target_seqs 1

-outfmt "6 qseqid sseqid evalue pident stitle"

blastn -task megablast -db db/refseq_rna

-query test_query.txt -dust no -max_target_seqs 1

-outfmt "6 qseqid sseqid evalue pident stitle"

blastn -task megablast -db db/refseq_rna

-query test_query.txt -dust no -max_target_seqs 1

formatting options for hits

-outfmt "6 qseqid sseqid evalue pident stitle"

www.ncbi.nlm.nih.gov/books/NBK279675

S NCBI   Resources ⊙   How To ⊙

**Bookshelf**   [Books ▾]   [                              ]   [Search]

Browse Titles    Limits    Advanced                                    Help

**BLAST® Command Line Applications User Manual [Internet].**      [< Prev]  [Next >]

▸ Show details

**Contents** ⊙

[                    ]  [Search this book]

# Options for the command-line applications.

This appendix consists of several tables that list option names, types, default values, and a short description of the option. These tables were first published as an appendix to an article in BMC Bioinformatics (BLAST+: architecture and applications). They have been updated for this manual.

## Table C1:

Options common to all BLAST+ search applications. An option of type "flag" takes no argument, but if present is true. Some options are valid only for a local search ("remote" option not used), others are valid only for a remote search ("remote" option used).

| option | type | default value | description and notes |
|--------|------|---------------|-----------------------|
| db | string | none | BLAST database name. |
| query | string | stdin | Query file name. |
| query_loc | string | none | Location on the query sequence (Format: start-stop) |

**Views**

PubReader

Print View

Cite this Page

PDF version of this page (169K)

PDF version of this title (1.3M)

**Other titles in this collection**

NCBI Help Manual

**Recent Activity**

Turn Off   Clear

Options for the command-line applications. - BLAST® Command Line Applications Us...

Extracting data from BLAST databases with blastdbcmd - BLAST® Command Line Appli...

BLAST+ features - BLAST® Command Line Applications User Manual

Building Customized Data Pipelines Using the Entrez Programming Utilities (eUtil...

mus musculus (1)

blastn -task megablast -db db/refseq_rna

-query test_query.txt -dust no -max_target_seqs 1

-outfmt "6 qseqid sseqid evalue pident stitle"

name and desired location of the output file

-out outputfile.txt