

The FASTA program package

Introduction

This documentation describes the version 36 of the FASTA program package (see W. R. Pearson and D. J. Lipman (1988), "Improved Tools for Biological Sequence Analysis", PNAS 85:2444-2448, [16] W. R. Pearson (1996) "Effective protein sequence comparison" Meth. Enzymol. 266:227-258 [14]; and Pearson et. al. (1997) Genomics 46:24-36 [17]. Version 3 of the FASTA packages contains many programs for searching DNA and protein databases and for evaluating statistical significance from randomly shuffled sequences.

This document is divided into four sections: (1) A summary overview of the programs in the FASTA3 package; (2) A guide to using the FASTA programs; (3) A guide to installing the programs and databases. Section (4) provides answers to some Frequently Asked Questions (FAQs). In addition to this document, the `changes_v36.html`, `changes_v35.html` and `changes_v34.html` files list functional changes to the programs. The `readme.v30..v36` files provide a more complete revision history of the programs, including bug fixes.

The programs are easy to use; if you are using them on a machine that is administered by someone else, you can focus on sections (1) and (2) to learn how to use the programs. If you are installing the programs on your own machine, you will need to read section (3) carefully.

FASTA and BLAST – FASTA and BLAST have the same goal: to identify statistically significant sequence similarity that can be used to infer homology. The FASTA programs offer several advantages over BLAST:

1. Rigorous algorithms unavailable in BLAST (Table I). Smith-Waterman (`ssearch36`), global:global (`ggsearch36`), and global:local (`glsearch36`) programs are available, and these programs can be used with `psiblast` PSSM profiles.
2. Better translated alignments. `fastx36`, `fasty36`, `tfastx36`, and `tfasty36` allow frame-shifts in alignments; frame-shifts are treated like gap-penalties, alignments tend to be longer in error-prone reads.
3. Better statistics. BLAST calculates very accurate statistics for protein:protein alignments, but its model-based strategy is less robust for translated-DNA:protein and DNA:DNA scores. FASTA uses an empirical estimation strategy, and now provides both search-based, and high-scoring shuffle-based statistics (`-z 21`).
4. More flexible library sequence formats. The FASTA programs can read FASTA, NCBI/`formatdb`, and several other sequence formats, and can directly query MySQL and Postgres databases. The programs offer several strategies for specifying subsets of databases.
5. A very efficient threaded implementation. The FASTA programs are fully threaded; both similarity scores and alignments can be calculated in parallel on multi-core hardware. On multi-core machines, FASTA can be faster than BLAST while producing better alignments with more accurate statistical estimates.

In addition, the latest (fasta-36.3.4) version of the FASTA programs provides an option to produce very BLAST-like output (`-m BB`), so that analysis pipelines require minimal modification.

1 An overview of the FASTA programs

Although there are a large number of programs in this package, they belong to three groups: (1) Traditional similarity searching programs: `fasta36`, `fastx36`, `fasty36`, `tfastx36`, `tfasty36`, `ssearch36`, `ggsearch36`, and `glsearch36`; (2) Programs for searching with short fragments: `fasts36`, `fastf36`, `tfasts36`, `tfastf36`, and `fastm36`; (3) A program for finding non-overlapping local alignments: `lalign36`. Programs that start with `fast` search protein databases, while `tfast` programs search translated DNA databases. Table I gives a brief description of the programs.

In addition, there are several programs included. `map_db` is used to index FASTA format sequence databases for more efficient scanning. `lav2ps` and `lav2svg` plot the `.lav` files produced by `lalign -m 11` as postscript (`lav2ps`) or SVG (`lav2svg`) output.

2 Using the FASTA Package

2.1 Introduction/Overview

All the FASTA sequence comparison programs use similar command line options and arguments. The simplest arguments are: the name of a query sequence file, a library file, and (possibly) the *ktup* parameter. If command line options are provided, they *must* precede the standard query-file and library-file arguments. Thus:

```
fasta36 -s BP62 query.file library.file
```

will compare the sequences in `query.file` with those in `library.file` using the BLOSUM62 scoring matrix with BLASTP gap penalties (-11/-1).

Current versions of the FASTA programs expect a query file and library, if you simply type “`fasta36`”, you will see a short help message:

```
% ssearch36
USAGE
ssearch36 [-options] query_file library_file
ssearch36 -help for a complete option list

DESCRIPTION
SSEARCH performs a Smith-Waterman search
version: 36.3.4 Mar, 2011

COMMON OPTIONS (options must precede query_file library_file)
-s: [BL50] scoring matrix;
-f: [-10] gap-open penalty;
-g: [-2] gap-extension penalty;
-S filter lowercase (seg) residues;
-b: high scores reported (limited by -E by default);
-d: number of alignments shown (limited by -E by default);
-I interactive mode;
```

“`fasta36 -help`” (or any of the other program names in Table I) provides complete listing of the options available for the program and their default values.

Table 1: Comparison programs in the FASTA36 package

FASTA program	BLAST equiv.	Description
fasta36	blastp/ blastn	Compare a protein sequence to a protein sequence database or a DNA sequence to a DNA sequence database using the FASTA algorithm [14, 16]. Search speed and selectivity are controlled with the <i>ktup</i> (wordsize) parameter. For protein comparisons, <i>ktup</i> = 2 by default; <i>ktup</i> = 1 is more sensitive but slower. For DNA comparisons, <i>ktup</i> =6 by default; <i>ktup</i> =3 or <i>ktup</i> =4 provides higher sensitivity.
ssearch36		Compare a protein sequence to a protein sequence database or a DNA sequence to a DNA sequence database using the Smith-Waterman algorithm [20]. <i>ssearch36</i> uses SSE2 acceleration, and is only 2 - 5X slower than <i>fasta36</i> [5].
ggsearch36/ glsearch36		Compare a protein sequence to a protein sequence database or a DNA sequence to a DNA sequence database using an optimal global:global (<i>ggsearch36</i>) or global:local (<i>glsearch36</i>) algorithm.
fastx36/ fasty36	blastx	Compare a DNA sequence to a protein sequence database, by comparing the translated DNA sequence in three frames and allowing gaps and frameshifts. <i>fastx36</i> uses a simpler, faster algorithm for alignments that allows frameshifts only between codons; <i>fasty36</i> is slower but can produce better alignments because frameshifts are allowed within codons [24].
tfastx36/ tfasty36	tblastn	Compare a protein sequence to a DNA sequence database, calculating similarities with frameshifts to the forward and reverse orientations [24].
fastf36/ tfastf36		Compares an ordered peptide mixture, as would be obtained by Edman degradation of a CNBr cleavage of a protein, against a protein (<i>fastf</i>) or DNA (<i>tfastf</i>) database [10].
fasts36/ tfasts36		Compares set of short peptide fragments, as would be obtained from mass-spec. analysis of a protein, against a protein (<i>fasts</i>) or DNA (<i>tfasts</i>) database [10].
lalign36		Calculate multiple, non-intersecting alignments using the sim2 implementation of the Waterman-Eggert algorithm [21] developed by Xiaohui Huang and Web Miller [7]. Statistical estimates are calculated from Smith-Waterman scores of shuffled sequences.

The program can also be run by typing:

```
fasta36 -I
```

which presents the “classic” interactive mode (this was the default behavior before version 36.3.4). In interactive mode, you will be prompted for: (1) the name of the test sequence file; (2) the name of the library file; (3) whether you want ktup = 1 or 2. (1 – 6 for DNA sequences).

The package includes several test files. To check to make certain that everything is working, you can try:

```
fasta36 ../seq/musplfm.aa ../seq/prot_test.lib
or
tfastx36 ../seq/mgstml.aa ../seq/gst.nlib
```

2.2 Sequence files

The `fasta36` programs can read query and library files in many standard formats (see ??). The default file format for query and library files – the format that will be used if no additional file format information is provided – is FASTA format. Like BLAST, version 36 can compare a query file with multiple query sequences to a sequence database, performing an independent search with each sequence in the query file.

FASTA format files consist of a description line, beginning with a ‘>’ character, followed by the sequence itself:

```
>sequence name and description 1
A F A S Y T .... actual sequence.
F S S      .... second line of sequence.
>sequence name and description 2
PMILTYV ... sequence 2
```

All of the characters of the description line are read, and special characters can be used to indicate additional information about the sequence. In general, non-amino-acid/non-nucleotide sequences in the sequence lines are ignored.

FASTA format files from major sequence distributors, like the NCBI and EBI, have specially formatted description lines, e.g.:

```
>gi|54321|ref|np_12345| example NCBI refseq sequence
or
>sw:gstml_human P01234 glutathione transferase GSTM1 - human
```

Several sample test files are included with the FASTA distribution: `seq/*.aa` and `seq/*.seq`, as well as two small sequence libraries, `seq/prot_test.lib` and `seq/gst.nlib`.

You can build your own library by concatenating several sequence files. Just be sure that each sequence is preceded by a line beginning with a ‘>’ followed by a sequence name/description. Sequences entered with word processors should use a “text” mode, e.g. “Save as text” with MS-WORD, with end of line characters and no special formatting characters in the file. The FASTA program cannot read Microsoft Word .DOC files, or rich text (.RTF) files; query and library sequence files should contain only sequence descriptions, sequences, and end-of-line characters.

2.3 Running the programs

As mentioned earlier, the FASTA programs can be run either interactively, by typing the name of a FASTA program (and possibly command line options), followed by `-I` (`fasta36 -I`) or from the command line, entering command line options, and the query and library file names. For searches of large databases that may take several minutes (or longer), it is more convenient to run searches from the command line, e.g.:

```
fasta36 query.file library.file > output.file
```

The command line shown above could be typed in a Unix or MacOSX terminal window, or from the MS-Windows command line interface (`command.exe`). The command line syntax shown above works for all the FASTA programs, e.g.:

```
lalign36 mchu.aa mchu.aa > mchu.laln
fastx36 mgstm1.seq prot_test.lseg > mgstm1.fx_out
sssearch36 mgstm1.aa xurtg.aa > mgstm1_xurtg.ss
```

Command line options – The FASTA programs provide a variety of command line options that modify the default scoring matrix (`-sBL62`) and gap penalties (`-f -11, -g -1`), other algorithm parameters, the output options (`-E 0.1, -d 20, -m 9i`), and statistical procedures (`z -2`). A complete list of command line options is shown near the end of this document. Unlike the BLAST programs, all FASTA command line options must precede the query file name and library file name (and there are no command line options available to specify the query and library file names). Thus, you should type:

```
sssearch36 -s BL62 -f -11 -g -1 query.file library.file > output.file
```

If you include `-I` as one of the options, you can provide command line options (e.g. to change the scoring matrix or gap penalties) without a query file or library file, and the program will use the options but prompt for the necessary files .

2.4 Interpreting the results

Fig. 1 shows the output from a typical FASTA program (`sssearch36`). The output file can be viewed as four parts: (a) the initial command line and description of the query sequence used (`mgstm1.aa`, 218 aa) and library (PIR1, 13,351 entries); (b) a description of the search statistics, algorithm (Smith-Waterman, SSE2 accelerated), and search parameters (BLOSUM50 matrix, gap penalties: -10 to open a gap, -2 for each residue in a gap); (c) a list of high scoring library sequences, descriptions, similarity scores, and statistical significance; (d) the alignments that produced the scores.

2.4.1 Identifying homologs

In the description section (which starts: `The best scores are:`), four numbers after the description of each library sequence are shown: (i) (in parentheses) the length of the library sequence; (ii) the raw Smith-Waterman score for the alignment (`s-w`; for the `fasta36`, `[t]fast[x,y]36` programs, this column would be labeled `opt`, for the *optimized* – banded Smith-Waterman – score), (iii) the *bit* score, and (iv) the expectation (`E()`), or statistical significance, of the alignment score.

Figure 1: ssearch36 results

6

The E()-value depends on the size of the database searched, in this case, 13,351 sequences, so the database size is given at the top of the list.

The bit score is equivalent to a BLAST bit score; together with the length of query and library sequences, it can be used to calculate the significance of the alignment.¹ Bit scores are convenient because they provide a matrix independent score that can be compared with other searches performed with other matrices and gap penalties against other databases. However, the E()-value, or expectation, provides the most direct measure of the statistical significance of the match.

In this example, the GSTP1_RAT, GSTA1_RAT, and GSTA4_RAT proteins share strong significant similarity (better than $E() < 6.2 \times 10^{-7}$), while the GSTF1_MAIZE, GSTF3_MAIZE, and GSTT1_DROME sequences do not share significant similarity (better than 0.001). However, GSTF1_MAIZE, GSTF3_MAIZE, and GSTT1_DROME are all glutathione transferase homologs, they simply do not share statistically significant similarity with this particular mGSTM1 query. Statistically significant sequence similarity scores *can* be used to infer *homology* (common ancestry), but non-significant scores *cannot* be used to infer *non-homology*.

While percent identity is often used to characterize the quality of an alignment and the likelihood that it reflects homology, the E()-value is a much more reliable value. Often sequences that share less than 30% identity will share very significant similarity (in the example above mgstm1.aa and GSTA4_RAT, with $E() < 6.2E-07$ are 25.6% identical). The expectation value captures information about conservative replacements, identities, and alignment length to provide a *single* value that captures the significance of the alignment.

For protein searches, library sequences with E()-values < 0.001 for searches of a 10,000 entry protein database are almost always homologous. Some sequences with E()-values from 1 - 10 may also be related, but unrelated sequences (1–10 per search) will have scores in this range as well.

E()-values < 0.001 can reliably be used to infer homology, assuming that the statistical estimates are accurate. The two most common causes of statistical problems are low-complexity regions and amino-acid composition bias. Low-complexity regions can be identified using the pseg program [22], and filtered out using the -S option. Composition bias rarely produces highly-significant E()-values, but can increase the number of sequences with E()-values between 0.01 and 0.001. The FASTA programs offer two shuffle-based strategy for evaluating composition bias; calculating similarity scores for random sequences with the same length and amino acid composition (-z 11 .. 16), this is done for pairwise alignments and lalign36 by default, or by showing statistical estimates derived from shuffles of the high-scoring sequences (-z 21, 22, 24, 25, 26).

The statistical routines assume that the library contains a large sample of unrelated sequences. If the library contains fewer than 500 sequences (MAX_RSTATS), then the library sequences are shuffled to produce 500 random scores, from which lambda and K statistical parameters are estimated. If the library contains a large number of *related* sequences, then the statistical parameters can be estimated by using the -z 11-15, options. -z options greater than 10 calculate a shuffled similarity score for each library sequence, in addition to the unshuffled score, and estimate the statistical parameters from the scores of the shuffled sequences.

¹ $E(D) = Dmn2^{-b}$, where D is the number of sequences in the database, m, n are the lengths of the two sequences, and b is the bit score.

```

>>GST26_SCHMA Glutathione S-transferase class-mu (218 aa)
  initn: 422 init1: 359 opt: 407 Z-score: 836.8 bits: 162.0 E(437847): 3.7e-39
Smith-Waterman score: 451; 42.4% identity (73.4% similar) in 203 aa overlap (6-208:6-203)

      10      20      30      40      50      60      70      80
mGSTM1 MPMILGYWNVRLTHPIRMLLEYTDSSYDEKRYTMGDAPDFDRSQWLNEKFKLGLDFPNLPYLIDGSHKITQSNAILRYL
      : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : :
GST26_ MAPKFGYWKVKGLVQPTRLllehleeetyeeRAY--DRNEIDA--WSNDKFKLGLGFNPYPYIDGDFKLTQSMAlIRYI
      10      20      30      40      50      60      70

      90     100     110     120     130     140     150     160
mGSTM1 ARKHHLDGETEEERIRADIVENQVMDTRMQLIMLCYNPDFEKQKPEFLKTIPEKMKLYSEFLGKRPWFAGDKVTVYVDFLA
      : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : :
GST26_ ADKHNMLGACPKERAIEISMLEGAVLDIRMGVLRlAYNKEYETLKVDFLNKLPGRLKMFEDRLSNKTYLNGNCVTHPDFML
      80      90      100      110      120      130      140      150

      170     180     190     200     210
mGSTM1 YDILDQYRMFEPKCLDAFPNLRDFLARFEGLKKISAYMKSSRYIATPIFSKMAHWSNK
      : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : :
GST26_ YDALDVVLYMDSQCLNEFPKLVSFKKCIEDLPQIKNYLNSSRYIKWPLQGWDATFGGGDTPPK
      160     170     180     190     200     210

```

Figure 2: Alignment with -s filtered sequence

2.4.2 Looking at alignments

The description section described above contains the critical information for inferring homology, the $E()$ -value. The alignment section shows the actual alignments that produced the similarity score and statistical estimates. In Fig. 1, the alignment display reports the percent identity, percent similarity (number of aligned residues with BLOSUM50 values ≥ 0), and the boundaries of the alignment. Note that for the `ssearch36` and `fasta36`, the alignment shown can include residues that are not part of the best local alignment (e.g. residues 1–5 and 207–218 in `mGSTM1` in Fig. 1). The amount of additional sequence context shown is the alignment line length (60 residues, set by `-w len`) divided by 2 by default, but can be adjusted with the `-W` context option.

Fig. 2 shows an example of a `fasta36` alignment produced using the `-s` option to filter out lower-case (low complexity) residues. Here, additional scores (`initn`, `init1` are shown, in addition to the `opt` score which is used to rank the sequences and calculate statistical significance. The `init1` score is the highest scoring alignment without gaps; `initn` is a score that combines consistent (non-overlapping) runs without gaps, and `opt` is the score of a banded Smith-Waterman of width 16 for `ktup=2` that is applied to sequences with `initn` scores over the optimization threshold. For `fasta36` with proteins, the final alignment and score is calculated with the Smith-Waterman algorithm. For DNA sequences, a banded Smith-Waterman is used. (The `-A` option produces banded Smith-Waterman alignments for proteins, and full Smith-Waterman for DNA.) In Fig. 2, all four scores are different. The `opt` score and Smith-Waterman scores are calculated on exactly the same alignment, but the `opt` score excludes the contribution from the “low-complexity” region between 19–30 in `GST26_SCHMA`. The `init1` score is based on the long, un-gapped region from residues 46–208 in `mGSTM1`, while the `initn` and `opt` scores include the other regions joined by gaps. The `initn` score is higher than the `opt` score, because it uses a simpler, length-independent, gap penalty.

2.4.3 Results without alignments

While sequence alignments are very informative, it is often not practical to examine all the statistically significant alignments in large-scale searches. The `-m 9` and `-m 8` options present summaries of each alignment (alignment boundaries, percent identity, and other information) in a much more compact form. In addition, `-m 8c` or `-m 8C` (see options below) provide a detailed encoding of the alignment, that allows it to be reconstructed. For large-scale searches, we routinely use `-m 8` with the `-d 0` option, which sets the number of alignments shown to 0 (thus none are shown). Alternatively, the `-m 8` and `-m 8C` output options produce BLAST-format tabular results summaries (`-m 8C` provides commented tabular results). `-m BB` produces an output that mimics BLAST output (with alignments).

2.5 Program Options

Command line options are available to change the scoring parameters and output display. Unlike the NCBI BLAST programs, command line options *must* precede the query file name and library file name arguments. To see the command-line options for a program and their defaults, type `program_name -help`, e.g. `fasta36 -help` or `ssearch36 -help`. For a quick list of the most common options, just type the program name without any options (e.g. `fasta36<ret>`).

2.5.1 Command line options

- a (fasta36, ssearch36, glsearch36, fasts36) show both sequences in their entirety.
- A force Smith-Waterman alignments for fasta36 DNA sequences. By default, only fasta36 protein sequence comparisons use Smith-Waterman alignments. Likewise, for proteins, use band alignments (Smith-Waterman is used by default).
- B Show the z-score, rather than the bit-score in the list of best scores (rarely used, provided for backward compatibility).
- b # Number of sequence scores to be shown on output. In the absence of this option, fasta36 (and ssearch36) display all library sequences obtaining similarity scores with expectations less than the expectation (-E) threshold, 10.0 for proteins, and 2.0 for DNA:DNA and protein:translated DNA. The -b # option can limit the display further. In addition, -b =100 will force 100 high scores to be displayed, regardless of the expectation (-E) threshold.
- c #,# (fasta36, [t]fast[x,y]36 only) Fraction of alignments optimized (second value is fraction of sequences joined). FASTA36 uses a statistical threshold strategy that joins and optimizes only the fraction of the alignments with an `initn` score expected -c times. Thus, -c 0.05 should optimize about 5% of sequences. The actual number of sequences optimized (and joined) is displayed in the scoring parameters line. Thus:

```
Parameters: BL50 matrix (15:-5), open/ext: -10/-2
            ktup: 2, E-join: 1 (0.687), E-opt: 0.2 (0.294), width: 16
```

reports that 20% of the sequences in the database should have been band-optimized, and 29.4% were. Reducing the -c opt fraction improves performance, but dropping the fraction below 0.02 can reduce the accuracy of the statistical estimates.

- c 0 (letter 'O') sets the joining/optimization thresholds as they were prior to fasta-36.3.3 (original thresholds). Positive values set the thresholds to specific score values, as was the case in older versions of fasta.
- C length of the sequence name printed at the beginning of alignment lines (default 6 characters).
- d # Maximum number of alignments to be displayed (must be \leq to the number of descriptions, -b #)
- E e_cut [e_cut_r] Limit the number of scores and alignments shown based on the expected number of scores. Used to override the expectation value of 10.0 (protein:protein; 5.0 translated-DNA:protein; 2.0 DNA:DNA) used by default. -E 2.0 will show all library sequences with scores with an expectation value ≤ 2.0 . With fasta-36.3.4, a second value, e_cut_r is available to limit the E()-values of additional sequence alignments between the query and library sequences. If not given, the threshold is e_cut/10.0. If given with a value > 1.0 , e_cut_r = e_cut / value; for a value < 1.0 , e_cut_r = value; If e_cut_r ≤ 0.0 , then the additional alignment option is disabled.
- f # Gap open penalty (-10 by default for proteins, -12 for DNA, -12 for [t]fast[xy]).
- F # Limit the number of scores and alignments shown based on the expected number of scores. -E # sets the highest E()-value shown; -F # sets the lowest E()-value displayed. Thus, -F 0.0001 will not show any matches or alignments with $E() < 0.0001$. This allows one to skip over close relationships to search for more distant relationships.
- g # Penalty per residue in a gap (-2 by default for proteins, -4 for DNA, -2 for [t]fast[xy]). A single residue gap costs f + g.
- h Short help message. Help options with ':', e.g. -s:, require an argument (-s BP62). Defaults are shown in square brackets, e.g.: -s: [BL50].
- help Long help message
- H Show histogram.
- i DNA queries - search with reverse complement. For tfastx36/y36, search the reverse complement of the library sequence only (complement of -3 option).
- I Interactive mode (the default for versions older than fasta-36.3.4).
- j # Penalty for frameshift between codons ([t]fastx36, [t]fasty36) and within a codon (fasty36/tfasty36 only).
- J (lalign36 only) show the identity alignment (normally suppressed, -I in versions before fasta-36.3.4).
- k # number of shuffles for statistical estimates from shuffling.
- l file Location of library menu file (FASTLIBS).
- L Display longer library sequence description.
- M low-high Range of amino acid sequence lengths to be included in the search.

-m # Specify alignment type: 0, 1, 2, 3, 4, 5, 6, 8, 9, 10, 11, B, BB

```

      -m 0      -m 1      -m 2      -m 3      -m 4
MWRTCGPPYT  MWRTCGPPYT  MWRTCGPPYT  MWKSCGYPYT  MWRTCGPPYT
:::  :::  xx X  ..KS..Y...  MWKSCGYPYT  -----
MWKSCGYPYT  MWKSCGYPYT

```

-m 5: a combination of **-m 4** and **-m 0**. **-m 6** provides **-m 5** plus HTML formatting. In addition, independent **-m** options can be combined. Thus, one can use **-m 1 -m 6 -m 9**.

-m 8 provides BLAST tabular format output (a tab delimited line with the query name, library name, percent identity, and other alignment information). “**-m 8c**” provides some additional information provided by the BLAST tabular format with comment lines.

-m 9 display alignment coordinates and scores with the best score information. **-m 9i** provides alignment length, percent identity, and percent similarity only. **-m 9** extends the normal best score information:

```

The best scores are:                                opt bits E(14548)
XURTG4 glutathione transferase (EC 2.5.1.18) 4 -    ( 219) 1248 291.7 1.1e-79

```

to include the additional information (on the same line, separated by <tab> characters):

```

%_id %_gid  sw  alen  an0  ax0  pn0  px0  an1  ax1  pn1  px1  gapq  gapl  fs
0.771 0.771 1248 218   1  218   1  218   1  218   1  219   0   0   0

```

-m 9c provides additional information: an encoded alignment string. For example, the alignment:

```

      10      20      30      40      50      60      70
GT8.7  NVRGLTHPIRMLLEYTDSSYDEKRYTMGDAPDFDRSQWLNEKFKL--GLDFPNLPYL-IDGSHKITQ
      :::  . ::: .  :::  :  :  :  :  :  :  :  :  :  :  :  :  :  :  :  :  :  :  :
XURTG  NARGRMECIRWLLAAAGVEFDEK-----FIQSPEDLEKCLKKDGNLNMFQVPMVEIDG-MKLAQ
      20      30      40      50      60

```

would be encoded: =23+9=13-2=10-1=3+1=5 . The numbers in the alignment encoding is with respect to the beginning of the alignment, not the sequences. The beginning coordinate of the alignment is given earlier in the **-m 9c** line. **-m 9c** provides the alignment encoding in CIGAR format: 28M9D13M2I10M1I3M1D5M .

-m 10 a parseable format for use with other programs.

-m 11 Provide lav-like output (used by **lalign**) for graphical output.

```
lalign36 -m 11 mchu.aa mchu.aa | lav2ps > mchu_laln.ps
```

Produces a postscript plot of the local alignments. Likewise, **lav2svg** produces SVG output.

-m BB Format output to mimic BLAST format. **-m B** formats alignments to look like BLAST alignments (Query/Sbjct), but is FASTA output otherwise. **-mBB** imitates BLAST as much as possible, and cannot be used with other **-m** options.

-M low-high Include library sequences with lengths between low and high.

- n Force the query sequence to be treated as a DNA sequence. Useful when query sequences contain a large number of ambiguous residues, e.g. transcription factor binding sites.
- N # break long library sequences into blocks of # residues. Useful for bacterial genomes, which have only one sequence entry. -N 2000 works well for well for bacterial genomes. (This option was required when FASTA only provided one alignment between the query and library sequence. It is not as useful, now that multiple alignments are available.)
- O Send a copy of results to filename. Helpful for environments without STDOUT, but should be avoided (use > filename instead).
- o (fasta36, [t]fast[x,y] only) Turn off the default opt score calculation and sort results by initn scores (reduces sensitivity and statistical accuracy, obsolete).
- p Force query to be treated as protein sequence.
- P PSSM_file Specify a PSI-BLAST format PSSM (Position Specific Scoring Matrix) file. ssearch36, ggsearch36, and glsearch36 can use a PSSM file to improve the sensitivity of a search. The FASTA programs accept two PSSM file formats:

format	blastpgp option
0	blastpgp -C PSSM_file -u 0 byte-encoded
1	blastpgp -C PSSM_file -u 1 text ASN.1

which can be specified after the file name, e.g.:

```
ssearch36 -P 'gstt1_drome.pssm 1' gstt1_drome.aa +sp+
```

Searches with a PSI-BLAST PSSM must still include a query sequence file, and the query sequence file must match the PSSM seed sequence. The format 0 byte-encoded PSSM is machine dependent; it must be created by blastpgp on the same architecture as ssearch36.

- Q, -q Quiet - does not prompt for any input. Writes scores and alignments to the terminal or standard output file (on by default, turned off with -I).
- r +n/-m Specify match/mismatch scores for DNA comparisons. The default is +5/-4. +3/-2 can perform better in some cases.
- R file Save a results summary line for every sequence in the sequence library. The summary line includes the sequence identifier, superfamily number (if available) position in the library, and the similarity scores calculated. This option can be used to evaluate the sensitivity and selectivity of different search strategies [13,15].
- s file Specify the scoring matrix file. fasta36 uses the same scoring matrix format as Blast. Several scoring matrix files are included in the standard distribution in the data/ directory. For protein sequences: codaa.mat - based on minimum mutation matrix; idnaa.mat - identity matrix; pam250.mat - the PAM250 matrix; [4], (-s P250), and pam120.mat - a PAM120 matrix (-s P120). The default scoring matrix is BLOSUM50 (-s BL50). Other matrices include a series of modern PAM-based matrices [8]: MDM40/-s MD40, MDM20/-s MD20, and MDM10/-s MD10, and a selection from the BLOSUM series [6] BLOSUM50, 62, and 80/-s

BL50, -s BL62, -s BL80. -s BP62 sets the scoring matrix to BLOSUM62 and the gap penalties to -11/-1, identical to BLASTP. In addition, the VTML160 matrix (-s VT160) [12] and OPTIMA_5 (-s OPT5) [9] are available.

If the scoring matrix is prefaced by a question mark, e.g. ?BP62, then the scoring matrix is adjusted for each query to ensure that a 100% identical match can produce a score of at least 40 bits. This is designed for *fastx36* searches with potentially short DNA queries; A 120 nt DNA query can only produce a 40 amino-acid alignment, which, with BLOSUM62 -11/-1, cannot produce more than 23 bits of score. A scoring matrix with a higher information content is required; in the set available by default, MD40, with 2.22 bits/position, would be used. For more information about alignment length and information content, see [1].

- S Filter out lower-case characters in the query or library sequences for the initial score calculation (used to filter low-complexity – *seg-ed* – residues). The *pseg* program [22] can be used to lower-case mask low complexity regions in protein sequences. With the -S option, lower case characters in the query or database sequences are treated as X's during the initial scan, but are treated as normal residues during the final alignment display. Since statistical significance is calculated from the similarity score calculated during the library search, the lower case residues do not contribute to the score. However, if a significant alignment contains low complexity regions, the residues are shown (as lower case characters, Fig. 2). The *pseg* program can be used to produce databases (or query sequences) with lower case residues indicating low complexity regions using the command:

```
pseg ./swissprot.fasta -z 1 -q > swissprot.lseg
```

The -S option should always be used with FASTX/Y and TFASTX/Y because out-of-frame translations often generate low-complexity protein sequences. However, only lower case characters in the protein sequence (or protein database) are masked; lower case DNA sequences are translated into upper case protein sequences, and not treated as low complexity by the translated alignment programs. (There is an option in the Makefile, -DDNALIB_LC, to enable preserving case in DNA sequences.)

- t # Translation table - *fastx36*, *tfastx36*, *fasty36*, and *tfasty3* now support the BLAST translation tables. See <http://www.ncbi.nih.gov/Taxonomy/Utils/wprintgc.cgi>.
-t t or -t t# enables the addition of an implicit termination codon to a protein:translated DNA match. That is, each protein sequence implicitly ends with *, which matches the termination codes for the appropriate genetic code. -t t# sets implicit termination and a different genetic code.
- T # set number of threads/workers. Normally on a multi-core machine, the maximum number of processors/cores is used.
- U Treat the query sequence an RNA sequence. In addition to selecting a DNA/RNA alphabet, this option causes changes to the scoring matrix so that G:A , T:C or U:C are scored as G:G -3.
- v # Do window shuffles with the window size specified.
- V str Specify annotation characters that can be included (and will be ignored), in the query sequence file, but are displayed in the alignments. If a query file contains "ACVS*ITRLFT?",

where "*" and "?" are used to indicate phosphorylation, giving the option `-V '*?'`, the annotated characters in the query will (S*, F?) will be highlighted in the alignment (on the number line). A `fasts36` alignment of `seq/ngts.aa` compared to `seq/mgstm1.aa` with `-V '*?'` produces:

```

          *               10??
GT8.7      ILGYWN-----EYTDSSSYDEKR-----
          :::::          :::::::::::
GT8.7      MPMILGYWNVRGLTHPIRMLLEYTDSSSYDEKRYTMGDAPDFDRSQWLNEKFKLGLDFPNL
          10          20          30          40          50          60

```

In addition to showing the alignments of post-translationally modified sites, the `-V` option can be used to highlight active sites in library sequences. In the `-m 9c` output, the state of the annotated sites is summarized when `-V` is used.

- `-w #` Line length (width) = number (<200)
- `-W #` context length (default is 1/2 of line width `-w`) for alignment, for programs like `fasta36` and `ssearch36`, that provide additional sequence context.
- `-x #match,#mismatch` Specify the penalty for a match to an X, and mismatch to X, independently of the PAM matrix. Particularly useful for `fastx3/fasty36`, where termination codons are encoded as X.
- `-X off1,off2` Specifies offsets for the beginning of the query and library sequence. For example, if you are comparing upstream regions for two genes, and the first sequence contains 500 nt of upstream sequence while the second contains 300 nt of upstream sequence, you might try:

```
fasta -X "-500 -300" seq1.nt seq2.nt
```

If the `-X` option is not used, FASTA assumes numbering starts with 1. (You should double check to be certain the negative numbering works properly.)

- `-y` Set the width of the band used for calculating "optimized" scores. For proteins and `ktup=2`, the width is 16. For proteins with `ktup=1`, the width is 32 by default. For DNA the width is 16.
- `-z -1,0,1,2,3,4,5,6`
 - `-z -1` turns off statistical calculations. `z 0` estimates the significance of the match from the mean and standard deviation of the library scores, without correcting for library sequence length. `-z 1` (the default) uses a weighted regression of average score vs library sequence length; `-z 2` uses maximum likelihood estimates of λ and K ; `-z 3` uses Altschul-Gish parameters [2]; `-z 4 - 5` uses two variations on the `-z 1` strategy. `-z 1` and `-z 2` are the best methods, in general.
- `-z 11,12,14,15,16`
 - estimate the statistical parameters from shuffled copies of each library sequence. This allows accurate statistics to be estimated for libraries comprised of a single protein family.
- `-z 21,22,24,25,26`
 - estimate the statistical parameters from shuffled copies of the highest scoring sequences reported in the search. library sequence. This shuffling strategy is much more like `prss`, since the sequences shuffled share compositional similarity to the query.

- Z `db_size` sets the apparent size of the database to be used when calculating expectation E()-values. If you searched a database with 1,000 sequences, but would like to have the E()-values calculated in the context of a 100,000 sequence database, use `-Z 100000`.
- 1 sort output by `init1` score (for compatibility with FASTP; obsolete).
- 3 translate only three forward frames or search with only the forward strand (complement of `-i`).

Thus, to tell `fasta36` to align `seq1.aa` with `seq2.aa` showing the entirety of both sequences, with 80 characters per line, one would type:

```
fasta36 -w 80 -s BP62 -a seq1.aa seq2.aa
```

The `-w 80` and `-a` options must precede the file names. If you just enter the options on the command line followed by `-I`, the program will prompt for the file names.

In addition, the FASTA programs can accept query sequence data from STDIN. To specify that stdin be used as the query or library file, the file name should be specified as `@`. Thus:

```
cat query.aa | fasta36 @:25-75 /slib/swissprot
```

would take residues 25-75 from `query.aa` and search the `/slib/swissprot`.

2.5.2 Environment variables

FASTA allows virtually every option to be set on the command line (except the *ktup*, which must be set as the third command line argument), but it is often convenient to set the FASTLIBS environment variable to specify the location of the `fastlibs` database description file.

FASTLIBS – FASTLIBS specifies the location of the file that contains the list of library descriptions, locations, and library types (see section on finding library files).

REF_URL, **SRCH_URL** and **SRCH_URL1** – These environment variables are used in HTML mode (`-m 6`) to provide links from the sequence alignment (see the links at http://fasta.bioch.virginia.edu/fasta_www2/). **REF_URL** is associated with the Entrez Lookup link; **SRCH_URL** with the Re-search database link, and **SRCH_URL1** with the General re-search link. In each case, the text corresponds to a HTML URL, but with positions containing the `%s` or `%ld` (for numbers) part of a 'C' `sprintf()` call for specific variables. **REF_URL** uses the database (protein or nucleotide), together with a query term (typically the `gi` number). **SRCH_URL** and **SRCH_URL1** use `db`, `query` (`gi`, `pgm` (`fa`, `ss`, `fx`, etc.)), and `start`, `stop`, and `n1` (library sequence length), where `start` and `stop` are the boundaries of the alignment, for sub-sequence searches. The values of these environment variables are used with `sprintf` to build a new URL that is linked in the output.

In addition, environment variables can be used inside both the `fastlibs` file and in the `@db.nam` files of file names. The `fasta36/conf/fast_libs_e.www` file, included with the distribution, shows an example, as do the descriptions of file of file names files shown below. Whenever a word of the form `${WORD}` is found in `fastlibs` or a file of file names, the `${WORD}` environment variable is expanded and inserted in the string. Thus, if `<${SLIB}/blast_dbs/` describes where a list of files will be found and `${SLIB}` is `"/seqdata"`, then the resulting substitution yields: `</seqdata/blast_dbs/`.

3 Installing FASTA and the sequence databases

3.1 Obtaining/preparing the sequence libraries

The FASTA program package does not include any protein or DNA sequence libraries. Protein and DNA sequence databases are available via anonymous FTP from the NCBI (<ftp://ftp.ncbi.nih.gov/blast/db>), UniProt (<ftp://ftp.uniprot.org/pub/databases/uniprot>), and the EBI (<ftp://ftp.ebi.ac.uk/pub/databases>).

Protein Sequence Databases – Protein sequence databases are available from the NCBI, UniProt, and the EBI. The NCBI provides a “raw” database, `nr`, and a well-curated, less redundant database, `refseq_protein`, and a copy of the very well annotated `swissprot` database. Protein sequence databases can also be downloaded from UniProt and the EBI; both sites provide the same UniProt [3] database.

Protein libraries, particularly those used for translated-DNA:protein comparisons with `fastx36` or `fasty36`, should be scanned to remove low-complexity regions. Matches between low complexity regions can violate the composition assumptions used by the FASTA statistical estimates. The `pseg` program ([22], <ftp://ftp.ncbi.nih.gov/pub/seg/pseg>) can be used to lower-case low complexity regions, which then can be ignored during the initial database search by using the `-S` option. To lower-case low complexity regions, run the `pseg` program against the protein sequence database:

```
pseg /seqdata/swissprot.fa -z 1 -q > /seqdata/swissprot.lseg
```

And then you can run most FASTA programs with `-S`:

```
ssearch36 -S mgstm1.aa /seqdata/swissprot.lseg
```

Fig. 2 shows the effect of including the `-S` option with lower-cased low-complexity sequences. The `opt` score (407), which is used to sort the results and calculate statistics, is lower than the Smith-Waterman score (451), even though exactly the same residues are aligned for each score. The `opt` score excludes residues 19-30, because they were marked as low-complexity by `pseg`; thus they are shown as lower-case. The Smith-Waterman score includes the contribution from that part of the alignment.

Out-of-frame translated DNA sequences often produce low-complexity regions [17], so it is particularly important to avoid low-complexity alignments when using `fastx36` and `fasty36`.

3.2 Searching taxonomic subsets

Because increasing database size reduces search sensitivity (an alignment with an $E()$ -value of 0.001 in a search of a 100,000 entry database will have an $E()$ -value of 0.1, not significant, if found in a database of 10,000,000 sequences), it is much more effective to search smaller, less redundant databases (you can always search the larger database later). Thus, the `refseq_protein` database from the NCBI is preferred over `nr`; even better are databases that reflect a limited phylogenetic range (e.g. `refseq_human` for vertebrate sequences).

While the NCBI provides organism-specific `refseq` subsets on their FTP site, they can be difficult to find. Alternatively, you can use the NCBI Entrez web site to download a list of `gi` numbers specific to a particular organism or taxonomic range. The FASTA programs can search a subset of a large sequence database that is specified by a list of `gi` numbers by using library

format 10. For example, given a list of `gi` numbers for the human proteins in `swissprot.lseg`, the file `sp_human.db`, with the content:

```
<${SLIB}/swissprot.lseg 0:2 4|
3121763
51701705
7404340
205831112
74735515
...
```

could be used to search the human subset of `swissprot.lseg`. The `gi` numbers for the Swiss-Prot entries begin with the second line. The first line specifies the location of the file where the sequences containing the `gi` numbers can be found (`${SLIB}/swissprot.lseg`, the `libtype` of that file (`0:fasta`), the character offset to the beginning of the sequence identifier in that file (`2`), the identifier type (`4`), and the character that separates the fields in the FASTA descriptor (`|`). The identifier type can take four formats:

-
- 1 ordered accession strings (letters or numbers)
 - 2 ordered numbers (digits only)
 - 3 un-ordered accession strings
 - 4 un-ordered numbers
-

(Ordered accession strings/numbers are ordered in both the library and the subset file.)

Thus, given the `0:2 4|` specification above, the line:

```
>gi|3121763|sp|015143.3|ARC1B_HUMAN Actin-related protein 2/3 ...
```

would be parsed, looking for a number starting at column 4 (the first column is numbered 0), and ending with `|`. The order of sequences in the library do not have to correspond to the order in the `sp_human.db` file (un-ordered). Given a the `sp_human.db` file, a file `swissprot.lseg` in the directory specified by the environment variable `${SLIB}`, and a command of the form:

```
fasta36 -S mgstm1.aa 'sp_human.db 10'
```

Would use the `sp_human.db` file to search the subset of `swissprot.lseg` that contained the specified `gi` numbers.

3.3 DNA sequence libraries

Because of the large size of DNA databases, you will probably want to keep DNA databases in only one formats. The FASTA3 programs that search DNA databases - `fasta36`, `fastm36`, and `tfastx/y36` can read DNA databases in Genbank flatfile (not ASN.1), FASTA, and BLAST2.0 (`formatdb`) formats, as well as EMBL format. BLAST2.0 format is preferred for DNA sequence libraries, because the files are considerably more compact than GenBank format. The NCBI does not provide software for converting from Genbank flat files to Blast2.0 DNA databases, but you can use the `Blast formatdb` program to convert ASN.1 formatted Genbank files, which are available from the NCBI `ftp` site.

The NCBI also provides the comprehensive `nt` DNA database, and several EST databases in BLAST2.0/`formatdb` format from `ftp://ncbi.nih.gov/blast/db`.

3.4 Finding the library files

All the FASTA programs comparison programs have the command line syntax:

```
fasta36 query.file /seqdata/library
```

However, in addition to simply specifying the location of the database to be searched (`/seqdata/library`), the FASTA programs provide several methods for referring to sequence databases without specifying a specific file. These methods can be used to provide abbreviations for sequence libraries, e.g.:

```
fasta36 query.file s Or fasta36 query.file +sp+
```

To use abbreviations like `'s'` or `'+sp+'` to reference a sequence database, a FASTLIBS file must be used, see section 3.5.

Large DNA and protein databases are often distributed across several files. For example, the NCBI `nr` protein database is found in 5 files, `nr.00 ... nr.04`. To search databases in multiple files, the names of the files are specified in a file of filenames, `nr.nam`:

```
<${SLIB}/blast_dbs/  
nr.00 12  
nr.01 12  
nr.02 12  
nr.03 12  
nr.04 12
```

In this file, the first line `<${SLIB2}/blast_dbs/`, beginning with `<`, specifies the location and format (Blast2.0 `formatdb`) the data files. Text of the form `${SLIB}` refers to Unix/MacOSX/Windows environment variables; the value of `${SLIB}` is set by a Unix/MacOSX shell environment command. Thus, if the value of `${SLIB}` is `/seqdata`, then the first sequence library file to be read will be `/seqdata/blast_dbs/nr.00`, in format 12 (Blast2.0 `formatdb`).

To refer to the `nr.nam` file as a file of file names, it must be prefixed by a `@` character, e.g.

```
fasta36 query.file @nr.nam
```

Files of file names can contain references to other files of file names:

```
<${SLIB}/fasta_dbs/  
@pdb.nam  
@swissprot.nam
```

The FASTA file of file names is similar to the NCBI `prot_db.pal` and `dna_db.nal`, files, but unfortunately they are different, and currently FASTA cannot read NCBI `.pal` or `.nal` files that contain a DBLIST line. FASTA can read NCBI `.pal` or `.nal` files that do not contain a DBLIST line.

3.5 FASTLIBS

All the search programs in the FASTA3 package can use the environment variable `FASTLIBS` to find the protein and DNA sequence libraries. (Alternatively, you can specify the `FASTLIBS` file with the `-l fastlibs.file` option.) The `FASTLIBS` variable contains the name of a file that has the actual filenames of the libraries. The `fastlibs` file included with the distribution on is an example of a file that can be referred to by `FASTLIBS`. To use the `fastlibs` file, type:

```
setenv FASTLIBS /seqdata/info/fastgbs (csh/tcsh)
or
export FASTLIBS=/seqdata/info/fastgbs (bash/ksh)
```

Then edit the `fastlibs` file to indicate the location of the protein and DNA sequence libraries. If protein sequence library is kept in the file `/seqdata/aa/swissprot.lseg` and your Genbank DNA sequence library is kept in the directory: `/seqdata/genbank`, then the `fastlibs` file might contain:

```
SwissProt$0P/seqdata/aa/swissprot.lseg 0
UniProt$0U+uniprot+@/seqdata/aa/uniprot.nam
GB Primate$1P@/seqdata/genbank/gpri.nam
GB Rodent$1R@/seqdata/genbank/grod.nam
GB Mammal$1M@/seqdata/genbank/gmammal.nam
^ 1 ^^^^ 4 ^ ^
      23 (5)
```

The first line of this file says that there is a copy of the SwissProt sequence database (a protein database) that can be selected by typing "P" on the command line or when the database menu is presented in interactive mode.

Note that there are 4 (or 5) fields in the lines in the `fastlibs` file. The first field describes library and is displayed by FASTA program; it ends with the '\$'. The second field (1 character), is a 0 if the library is a protein library and 1 if it is a DNA library. The third field can either be a single character (P) or a word surrounded by the + symbol (+uniprot+), and can be used to specify the library on the command line or in interactive mode.

The fourth field is the name of the library file. In the example above, the `/seqdata/aa/swissprot.lseg` file contains the entire protein sequence library. Alternatively, `/seqdata/aa/uniprot.nam` is a file of file names, which contains a list of one or more library files. Likewise, the DNA library files are files of file names.

In addition, an optional fifth field can be used to specify the format of the library file. Alternatively, you can specify the library format in a file of file names. This field must be separated from the file name by a space character (' ') from the filename. FASTA can read the libraries in the following formats:

- 0 FASTA (>SEQID - comment/sequence)
- 1 Uncompressed Genbank (LOCUS/DEFINITION/ORIGIN)
- 2 NBRF CODATA (ENTRY/SEQUENCE) (obsolete)
- 3 EMBL/SWISS-PROT (ID/DE/SQ)
- 4 Intelligenetics (;comment/SEQID/sequence) (obsolete)
- 5 NBRF/PIR VMS (>P1;SEQID/comment/sequence) (obsolete)
- 6 GCG (version 8.0) Unix Protein and DNA (compressed)
- 11 NCBI Blast1.3.2 format (unix only) (obsolete)
- 12 NCBI Blast2.0 format
- 15 MySQL (requires special compilation)
- 16 Postgres (requires special compilation)

Today, the most popular formats are FASTA, type '0', the default, and the NCBI Blast2.0 formatdb formats (type '12'). The FASTA programs cannot read NCBI ASN.1 formatted databases. If a library format is not specified, for example, because you are just comparing two sequences, FASTA (format 0) is used by default. To specify a library type on the command line, add it to the

library filename and surround the filename and library type in quotes:

```
fasta36 query.file "/seqdb/genbank/gbmam 12"
```

NCBI formatdb databases are built from multiple files, e.g. `gbmam.nsq`, `gbmam.nhr`, `gbmam.nin`; to refer to the complete set of files, simply use name before the suffixes, e.g. `gbmam`. When NCBI databases distributed across several files, e.g. `gbtct.00`, `gbtct.01`, etc, those files must be included in a `gbtct.nam` file of file names.

The major problem that most new users of the FASTA package have is in setting up the program to find the databases and their library type. In general, if you cannot get `fasta36` to read a sequence database, it is likely that something is wrong with the `FASTLIBS` file. A common problem is that the database file is found, but either no sequences are read, or an incorrect number of entries is read. This is almost always because the library format (`libtype`) is incorrect.

Test the setup by running FASTA. Enter the sequence file `'mgstm1.aa'` when the program requests it (this file is included with the programs). The program should then ask you to select a protein sequence library. Alternatively, if you run the `tfastx36 -I` program and use the `mgstm1.aa` query sequence, the program should show you a selection of DNA sequence libraries. Once the `fastlibs` file has been set up correctly, you can set `FASTLIBS=fastgbs` in your `AUTOEXEC.BAT` file, and you will not need to remember where the libraries are kept or how they are named.

4 Frequently Asked Questions (FAQs)

Where can I get FASTA? – <http://faculty.virginia.edu/wrpearson/fasta> has the latest versions of the FASTA programs. This document describes fasta-36.3.4, which is available from <http://faculty.virginia.edu/wrpearson/fasta/fasta36/fasta-36.3.4.tar.gz>. In addition, pre-compiled versions of the programs are available for MacOSX and Windows.

Which program should I use? – See Table I, also:

Query	Library	FASTA pgm.	BLAST pgm.	
Prot.	Prot.	fasta36 ssearch36 ggsearch36 ggsearch36	blastp	heuristic local similarity optimal local sim. global:global sim. global:local sim.
DNA	DNA	fasta36*	blastn	
Prot. DNA	Prot. DNA	lalign36		multiple non-intersecting alignments
DNA	Prot.	fastx36 fasty36	blastx	trans. DNA:protein sim.
Prot.	DNA	tfastx36 tfasty36	blastn	protein:trans. DNA
Prot.	Prot.	fasts36		Unordered peptides
Prot.	DNA	tfasts36		Unordered peptides
DNA	DNA	fasts36		Unordered oligonucleotides
Prot.	Prot.	fastm36		Ordered peptides
DNA	DNA	fastm36		Ordered oligos

*ssearch36 can also be used for DNA:DNA, but is much slower and no more sensitive.

How do I make FASTA act/look like BLAST? –

```
fasta36 -s BP62 -m BB query.file library.file
```

-s BP62 sets the same scoring matrix (BLOSUM62) and gap-penalties (-11/-1) as BLAST (FASTA uses BLOSUM50 by default). -m BB produces very BLAST-like output.

When I search Genbank - the program reports: 0 residues in 0 sequences? This typically happens because the program does not know that you are searching a Genbank flatfile database and is looking for a FASTA format database. Be certain to specify the library type ("1" for Genbank flatfile) with the database name.

What is the difference between fastx3 and fasty3? (or tfastx3 and tfasty3)? – [t]fastx3 uses a simpler codon based model for alignments that does not allow frameshifts in some codon positions (see ref. [24]). fastx3 is about 30% faster, but fasty3 can produce higher quality alignments in some cases.

What is ktup? – All of the programs with `fast` in their name use a computer science method called a lookup table to speed the search. For proteins with `ktup=2`, this means that the program does not look at any sequence alignment that does not involve matching two identical residues in both sequences. Likewise with DNA and `ktup = 6`, the initial alignment of the sequences looks for 6 identical adjacent nucleotides in both sequences. Because it is less likely that two identical amino-acids will line up by chance in two unrelated proteins, this speeds up the comparison. But very distantly related sequences may never have two identical residues in a row but will have single aligned identities. In this case, `ktup = 1` may find alignments that `ktup=2` misses.

Where are prss and prfx? – Earlier FASTA3 releases included `prss3` and `prfx3`. With FASTA version 35 and 36, these programs have been incorporated into `ssearch36` and `fastx36`. FASTA version 35 and 36 programs now automatically estimate statistical parameters by shuffling - the function of `prss` and `prfx`, when searching for libraries with fewer than 500 members.

Where is tfasta? – Although it is possible to make `tfasta36`, it is not compiled by default. `tfastx36` and `tfasty36` allow frame-shifts to be joined into a single alignment; `tfasta` did not. `tfastx36` produces better alignments with better statistics.

Can I run the FASTA programs on a cluster? – With version 36.3.4, almost all of the FASTA programs can be run on clusters of computers using MPI (Message Passing Interface). The programs can be compiled using `make -f ../make/Makefile.mpi_sse2` from the `fasta36/src` directory. Except for `lalign36`, all the programs in Table I are available as `fasta36_mpi`, `ssearch36_mpi`, etc.

Unfortunately, the current MPI implementation involves substantially more communications overhead than the threaded versions. The FASTA programs are very efficient on threaded machines; if the preload option is used (edit `make/Makefile36m.common` to use `comp_lib7.c`), the FASTA programs can obtain more than 40-fold speedup on a 48-core machine (the largest I have tested).

Sometimes, in the list of best scores, the same sequence is shown twice with exactly the same score. Sometimes, the sequence is there twice, but the scores are slightly different? – When any of the FASTA programs searches a long sequence, it breaks the sequence up into *overlapping* pieces. If the highest scoring alignment is at the end of one piece, it will be scored again at the beginning of the next piece. If the alignment is not be completely included in the overlap region, one of the pieces will give a higher score than the other. These duplications can be detected by looking at the coordinates of the alignment. If either the beginning or end coordinate is identical in two alignments, the alignments are at least partially duplicates.

As always, please inform me of bugs as soon as possible.

William R. Pearson
Department of Biochemistry
Jordan Hall Box 800733
U. of Virginia
Charlottesville, VA
wrp@virginia.EDU

References

- [1] S. F. Altschul. Amino acid substitution matrices from an information theoretic perspective. *J. Mol. Biol.*, 219:555–65, 1991.
- [2] S. F. Altschul and W. Gish. Local alignment statistics. *Methods Enzymol.*, 266:460–480, 1996.
- [3] UniProt Consortium. Ongoing and future developments at the universal protein resource. *Nucleic Acids Res*, 39:D214–D219, 2011.
- [4] M. Dayhoff, R. M. Schwartz, and B. C. Orcutt. A model of evolutionary change in proteins. In M. Dayhoff, editor, *Atlas of Protein Sequence and Structure*, volume 5, supplement 3, pages 345–352. National Biomedical Research Foundation, Silver Spring, MD, 1978.
- [5] M. Farrar. Striped Smith-Waterman speeds database searches six times over other simd implementations. *Bioinformatics*, 23:156–161, 2007.
- [6] S. Henikoff and J. G. Henikoff. Amino acid substitutions matrices from protein blocks. *Proc. Natl. Acad. Sci. USA*, 89:10915–10919, 1992.
- [7] X. Huang, R. C. Hardison, and W. Miller. A space-efficient algorithm for local similarities. *Comp. Appl. Biosci.*, 6:373–381, 1990.
- [8] D. T. Jones, W. R. Taylor, and J. M. Thornton. The rapid generation of mutation data matrices from protein sequences. *Comp. Appl. Biosci.*, 8:275–282, 1992.
- [9] Maricel G Kann and Richard A Goldstein. Performance evaluation of a new algorithm for the detection of remote homologs with sequence comparison. *Proteins*, 48:367–76, Aug 2002.
- [10] A. J. Mackey, T. A. J. Haystead, and W. R. Pearson. Getting more from less: Algorithms for rapid protein identification with multiple short peptide sequences. *Mol. Cell. Proteomics*, 1:139–147, 2002.
- [11] R. Mott. Maximum-likelihood estimation of the statistical distribution of smith-waterman local sequence similarity scores. *Bull. Math. Biol.*, 54:59–75, 1992.
- [12] Tobias Muller, Rainer Spang, and Martin Vingron. Estimating amino acid substitution models: a comparison of dayhoff’s estimator, the resolvent approach and a maximum likelihood method. *Mol Biol Evol*, 19:8–13, 2002.
- [13] W. R. Pearson. Comparison of methods for searching protein sequence databases. *Prot. Sci.*, 4:1145–1160, 1995.
- [14] W. R. Pearson. Effective protein sequence comparison. *Methods Enzymol.*, 266:227–258, 1996.
- [15] W. R. Pearson. Empirical statistical estimates for sequence similarity searches. *J. Mol. Biol.*, 276:71–84, 1998.
- [16] W. R. Pearson and D. J. Lipman. Improved tools for biological sequence comparison. *Proc. Natl. Acad. Sci. USA*, 85:2444–2448, 1988.

- [17] W. R. Pearson, T. C. Wood, Z. Zhang, and W. Miller. Comparison of DNA sequences with protein sequences. *Genomics*, 46:24–36, 1997.
- [18] J. T. Reese and W. R. Pearson. Empirical determination of effective gap penalties for sequence comparison. *Bioinformatics*, 18:1500–1507, 2002.
- [19] T. Rognes and E. Seeberg. Six-fold speed-up of smith-waterman sequence database searches using parallel processing on common microprocessors. *Bioinformatics*, 16:699–706, 2000.
- [20] T. F. Smith and M. S. Waterman. Identification of common molecular subsequences. *J. Mol. Biol.*, 147:195–197, 1981.
- [21] M. S. Waterman and M. Eggert. A new algorithm for best subsequences alignment with application to tRNA-rRNA comparisons. *J. Mol. Biol.*, 197:723–728, 1987.
- [22] J. C. Wootton and S. Federhen. Statistics of local complexity in amino acid sequences and sequence databases. *Comput. Chem.*, 17:149–163, 1993.
- [23] A. Wozniak. Using video-oriented instructions to speed up sequence comparison. *Comput Appl Biosci*, 13:145–150, 1997.
- [24] Z. Zhang, W. R. Pearson, and W. Miller. Aligning a DNA sequence with a protein sequence. *J. Computational Biology*, 4:339–349, 1997.

Appendix

A FASTA Makefile compile time options

The `fasta-36/make` directory includes Makefiles appropriate for a broad range of environments, including Linux/Unix, BSD, MacOSX, and Windows. Makefiles are regularly tested against MacOSX, Linux, and Windows. Table 2 summarizes the Makefile options that can be modified.

As distributed, the Makefiles in `fasta36/make`, build a version of the FASTA programs that is optimized for single searches against arbitrary sized databases, using bit scores, efficient sampled statistics, and gap-open/extend penalties. The default compilation configuration can be changed either by changing the compile time defines (Table 2) in the main Makefile, e.g. `make/Makefile.linux64_sse2`, or by editing `make/Makefile36m.common`.

High-performance searches with many queries – By default, the `comp_lib5.c` program specified in `Makefile36m.common` builds FASTA programs that re-read the library sequence database for every query sequence. This has the advantage that sequence comparison begins almost immediately, but if thousands of searches are being performed, the database is re-read thousands of times. `Makefile36m.common` can be edited to use `comp_lib7.c` in place of `comp_lib5.c`, and the database is read only once, and held in memory for additional searches. Of course, if `comp_lib7.c` is used, the computer must have enough memory to store the complete database. Keeping the database in memory allows the FASTA programs to very efficiently used large, multicore computers.

Table 2: FASTA Makefile compile time #defines

#define	Status*	Target file(s)	Function
ALLOCNO	obs	dropnfa.c, dropfx.c, dropfz2.c	allows FASTA algorithm to use memory ~ query length (n0), not query + library (n0+n1).
DNALIB_LC	undef	initfa.c	enable lower case masking for DNA libraries
HTML_HEAD	undef	comp_lib5.c, comp_lib7.c	wrap -m 6 HTML output with <html> <body> </body> </html>
M10_CONS	def	c_disp.c	show consensus line (: .) with -m 10 out- put.
OLD_FASTA_GAP	undef	drop*.c	use first-residue/additional residue penalties, not open/extend.
PGM_DOC	def	comp_lib5.c, comp_lib7.c	provide #pgm_name -opt1 -opt2 query file copy of command line
PROGRESS	def	comp_lib5.c, comp_lib7.c	provide progress symbols in interactive mode
SAMP_STATS	def	comp_lib5.c, comp_lib7.c	scores are sampled for statistical estimates
SAMP_STATS_LESS	def	compacc.c	a slower sampling strategy is used
SHOW_ALIGN_SCORE	undef	wm_align.c	print score, cumulative score, during align- ment (for teaching)
SHOW_HELP	def	comp_lib5.c, comp_lib7.c, initfa.c, doinit.c	print out help information with '-help', or no arguments given. Undef SHOW_HELP reverts to pre-fasta-35.4.4.
SHOW_HIST	undef	doinit.c	inverts current meaning of -H (shows by de- fault for non-PCOMPLIB (MPI) programs).
SHOWSIM	def	mshowbest.c mshowalign2.c	display percent similarity
USE_LNSTATS	obs	scaleswn.c	use $\ln()$ -scaling for scores, removed in fasta2.0.

*Status: def: #defined in standard Makefiles; undef: undefined; obs: obsolete, provided backwards compatibility with FASTA2.0 or earlier.

Parallel searches with MPI – By default under Unix/Linux/MacOSX, the FASTA programs are threaded; they will spawn as many threads as CPU cores are available (this can be limited with the `-t n-threads` option). Using `comp_lib7.c`, we see almost 48-fold speedup on a 48-core machine. The FASTA programs can also be run in parallel in the MPI environment on clusters of computers. To build the MPI versions of the programs, use `make ../make/Makefile.mpi_sse2 ssearch36_mpi, fastx36_mpi, etc.` The MPI programs currently substantially more communications overhead than the threaded versions, so they may not scale as well to large clusters.

B FASTA version history

FASTA v33, Oct, 1999 – Dec, 2000	
Oct 1999	<p>Add support for NCBI Blast2.0 formatted libraries, and memory mapped databases. FASTA now reads both BLAST1.4 and BLAST2.0 formatted databases. (version 3.2t08)</p> <p>Include Maximum Likelihood Estimates for Lambda and K (-z 2)</p> <p>Include a new strategy for searching with low complexity regions. The pseg program can produce libraries with low complexity regions as lower case characters, which can be ignored during the initial FASTA/SSEARCH scan, but are considered when producing the final alignments. (3.3t01)</p> <p>Change output to report bit scores, which are also used by BLAST.</p>
Mar 2000	<p>Another new statistics option, -z 6, uses Mott's approach [11] for calculating a composition dependent Lambda for each sequence. (3.3t05)</p>
Dec 2000	<p>Automatically change the gap penalties when alternate (known) scoring matrices are used using Reese and Pearson gap penalties [18]. First implementation to read from MySQL databases.</p>
May 2001	<p>change all FASTA gap penalties from first-residue, additional residue to the gap-open, gap-extend values used by BLAST.</p>
FASTA v34, Jan, 2001 – Jan, 2007	
Jun 2002	<p>Modify statistical estimation strategy to sample all the sequences in the database, not just the first 60,000. (3.4t11)</p>
Jan 2003	<p>Implementation of vector-accelerated (Altivec) code for Smith-Waterman (SSEARCH) and banded Smith-Waterman (FASTA) using the Rognes and Seeburg [19] algorithm. This code was removed in Sept, 2003, because of possible conflict with a patent application, but was restored using a different algorithm in Nov. 2004.</p>
Jun 2003	<p>Provide PSI-SEARCH — an implementation of SSEARCH that can search with PSI-BLAST PSSM profile files. PSI-SEARCH estimates statistical significance from the distribution of actual alignment scores; thus the estimates are much more reliable than PSI-BLAST estimates. Also, change the similarity display to work with profiles. (3.4t22)</p>
July 2003	<p>Provide ASN.1 definition line parsing for BLAST formatdb v.4 libraries. Restructure the programs to use a table-driven approach to parameter setting. Two tables now define the algorithm, query sequence type, library type, scoring matrix, and gap penalties for all programs.</p>
Sept 2003	<p>A new option -v for annotating alignments provided. Designed for highlighting post-translational modifications with <i>fasts</i>, it can also be used to highlight active sites and other conserved residues. (3.4t23)</p>
Dec 2003	<p>Addition of -U option for RNA sequence comparison. G:A matches score like G:G matches to account for G:U basepairs. Change default <i>ktup</i> for short query sequences. Increase band-width for DNA banded final alignments.</p>

FASTA version history (cont.)

July 2004	Allow searching of Postgres, as well as MySQL database queries.
Nov 2004	(fa34t24) Incorporation of Erik Lindahl "anti-diagonal" Altivec implementation of [23] for Smith-Waterman only. Altivec ssearch34 is now faster than fasta34 for query sequences < 250 amino acids.
Jan 2005	Change FASTS to accommodate very large numbers of peptides (>100) for full coverage on long proteins
Jun. 2006	(fa34t26) Incorporation of Smith-Waterman algorithm for the SSE2 vector instructions written by Michael Farrar [5]. The SSE code speeds up Smith-Waterman 8 – 16-fold.

FASTA v35, March, 2007 – March, 2010

Mar. 2007	<p>fasta v35 – Accurate shuffle-based $E()$-values for all searches and alignments; statistics from searches against small libraries are supplemented with shuffled alignments.</p> <p>More efficient threading strategies on multi-core computers, for 12X speedup on 16-core machines.</p> <p>Inclusion of lalign (SIM) local domain alignments. lalign alignments now have accurate shuffle-based $E()$-values.</p>
Apr. 2007	Introduction of ggsearch, for global alignment searches, and glsearch, for searches with scores that are global in the query and local in the library. ggsearch and glsearch calculate $E()$ -values using the normal distribution. Both programs can search with PSI-BLAST PSSMs.
Dec. 2007	Efficient strategy for searching subsets of databases (lists of GI or accession numbers)
Feb. 2008	Annotations in either query or library sequences can be highlighted in the alignment, and the state of annotated residues is compactly summarized with -m 9c.
Oct. 2008	Modification lsim4.c (lalign35) provided by Xiaoqui Huang to ensure that self-alignments do not cross the identity diagonal.

FASTA v36, March, 2010 –

Mar. 2010	<p>FASTA v36 displays all significant alignments between query and library sequence. BLAST has always displayed multiple high-scoring alignments (HSPs) between the query and library sequence; previous versions of the FASTA programs displayed only the best alignment, even when other high-scoring alignments were present.</p> <p>New statistical options, -z 21, 22, 26, provide a second $E2()$-value estimate based on shuffles of the highest scoring sequences.</p>
-----------	---

FASTA version history (cont.)

	Improved performance using statistics-based thresholds for gap-joining and band-optimization in the heuristic FASTA local alignment programs, increasing speed 2 - 3X.
	Greater flexibility in specifying combinations of library files and subsets of libraries. FASTA v36 programs can include indirect files of library names inside of indirect files of library names.
	FASTA 36 programs are fully threaded, both for searches, and for alignments. The programs routinely run 12 - 15X faster on 8-core machines with "hyperthreading" (effectively 16 cores).
	-z 21 .. 26 E2() statistical estimates from shuffled best scores.
Sep. 2010	-m 8, -m 8C BLAST tabular output.
Nov, 2010	Variable scoring matrices (-m ?BP62).
Dec, 2010	(fasta-36.3.1) SSE2 vectorized ggsearch36, glsearch36 (Michael Farrar).
Jan, 2011	(fasta-36.3.2) MPI versions implemented and tested.
Feb, 2011	Introduce -m B, -m BB BLAST-like output.
Mar, 2011	(fasta-36.3.4) Program is no longer interactive by default. fasta36 -h and fasta36 -help provide common/complete options, with many defaults. doc/fasta_guide.pdf available.
