



# FAST: Fast Analysis of Sequences Toolbox

Travis J. Lawrence<sup>1</sup>, Dana L. Carper<sup>1</sup>, Kyle T. Kauffman<sup>2</sup>, Katherine C.H. Amrine<sup>1,3</sup>, Raymond S. Lee<sup>4</sup>, Claudia J. Canales<sup>4</sup> and David H. Ardell<sup>1,2\*</sup>

<sup>1</sup>Quantitative and Systems Biology, University of California, Merced, CA, USA

<sup>2</sup>Molecular and Cell Biology Unit, University of California, Merced, CA, USA

<sup>3</sup>Dept. of Viticulture and Enology, University of California, Davis, CA, USA

<sup>4</sup>School of Engineering, University of California, Merced, CA, USA

Correspondence\*:

David H. Ardell

Molecular and Cell Biology, School of Natural Sciences, University of California, Merced, 5200 North Lake Road, Merced, CA, 95343, USA, dardell@ucmerced.edu

## ABSTRACT

FAST (Fast Analysis of Sequences Toolbox) provides simple command-line tools for rapid prototyping of powerful and reproducible bioinformatic workflows on flat-file biological sequence databases. Modeled after the GNU (GNU's Not UNIX) Textutils such as `grep`, `cut`, and `tr`, FAST tools such as `fasgrep`, `fascut`, and `fastr` are designed to be serially composed (linked together) in UNIX-style pipelines but process data per-record rather than per-line. This design enables efficient inline processing of biological sequence data in a familiar and time-tested idiom. Functionality highlights of FAST include feature-based slicing of sites from alignments, sequence filtering and transformation with Perl regular expressions, numerical and taxonomic selection and sorting of sequence records, annotation of sequences with molecular biological statistics including composition and codon usage, translation of gapped data, and statistics for molecular population genetic analysis. FAST's default data exchange format is MultiFastA (specifically, a restriction of BioPerl FastA format). Additionally, FAST has an option that allows the use of Sanger and Illumina 1.8+ FASTQ formatted files. Compared to software already available, FAST is distinguished by consistency in interface design across utilities, conformance with Unix philosophy, and the portability, ease of installation, and security inherent in its BioPerl and Perl foundations. FAST workflows reduce manual interventions in data processing and facilitate the interactive investigation and control of data by its users. Bringing all the advantages of open-source software consistent with scientific ideals, FAST has a shallow learning curve with a high return on investment of mastery because its interface design is consistent with the conventions of GNU utilities, Perl, BioPerl, and R. FAST brings the power of Perl and BioPerl to bioinformatic users at the command-line without requiring previous programming skills and experience.

**Keywords:** Unix philosophy, MultiFASTA, pipeline, bioinformatic workflow, open source, BioPerl, regular expression, NCBI Taxonomy

## 1 INTRODUCTION

The field of molecular biology has changed significantly with the advent of Next Generation Sequencing (NGS) technology. It is now commonplace to analyze gigabases of data per experiment. Commonly, bioinformatics programs are developed for visualization and basic sequence manipulation through a

**Table 1.** FAST 1.0 utilities

Tool	Function	Textutil analog	Default field processed
fasgrep	regex selection of records	grep	identifiers
fasfilter	numerical selection of records		identifiers
fastax	taxonomic selection of records		descriptions
fashead	order-based selection of records	head	
fastail	order-based selection of records	tail	
fascut	index-based selection and reordering of data	cut	sequences
fasuniq	record reduction by content and order	uniq	sequences
alncut	selection of sites by content		sequences
gbfalncut	selection of sites by features		sequences
fassort	numerical or text sorting of records	sort	identifiers
fastaxsort	taxonomic sorting of records		identifiers
faspaste	merging of records	paste	sequences
fastr	character transformations on records	tr	identifiers
fassub	regex substitutions on records		identifiers
faslen	annotate sequence lengths		descriptions
fascomp	annotate monomeric compositions		descriptions
fascodon	annotate codon usage		descriptions
fasxl	annotate biological translations		descriptions
fasrc	annotate reverse complements		descriptions
fasconvert	convert format of records		
gbfgrep	select feature neighborhoods by context	grep	features
gbf2fas	emit sequences by regex matching on features	grep	features
alnpi	molecular population genetic statistics		
faswc	tally sequences and characters	wc	sequences

28 monolithic application with a Graphical User Interface (GUI) (Smith et al., 1994; Rampp et al., 2006; ?).  
 29 Command-line utilities for bioinformatics such as the EMBOSS package (Rice et al., 2000) or the scripts  
 30 that come with BioPerl (Stajich et al., 2002) typically offer suites of tools with simple, well-defined  
 31 functions, but are not necessarily designed to interoperate with each other. Some monolithic applications  
 32 are not open-source, which invites the accumulation and persistence of incorrect code.

33 To reduce human error and increase reproducibility in managing biological data, it is desirable to create  
 34 and document automated bioinformatic workflows governed by machine processing of biological data.  
 35 Web-based workflow suites such as Galaxy (?), Taverna (?) and BioExtract (?) are not quite as agile for  
 36 rapid prototyping of workflows as are UNIX pipelines.

37 The FAST utilities are modeled after the standard Unix toolkit (Peek, 2001) and follow the Unix  
 38 philosophy to “do one thing and do it well” (Stutz, 2000). They are written in Perl using BioPerl packages  
 39 (Stajich et al., 2002). This makes FAST utilities easy to adopt if you are familiar with the Unix toolbox  
 40 and allows fast sequence analysis even on large datasets. Extensive documentation has been developed  
 41 for each utility along with useful error messages following the recommendations of Seemann (2013) to  
 42 increase usability. Lastly, FAST is open source, which makes it available to anyone free of cost. This is  
 43 in line with the call to make science more assessable, open, and reproducible by other scientists and the  
 44 public (Groves and Godlee, 2012).

## 2 DESIGN

45 An overview of Version 1.0 of the FAST project is shown in Figure 1. Descriptions of the function of  
 46 utilities along with GNU Textutil analogs are given in Table 1.

FAST is split into three categories selection, transformation, and annotation and analysis. The selection category contains utilities designed to select sequences and sites from alignments based on several different criteria. For example `fasgrep` selects sequences by matching a regular expression to the ID, description, or sequence. The transformation utilities are used to modify the ID, description, sequence, or order of sequences using several criteria. For example, `fastaxsort` sorts sequences within a multifasta file based on NCBI taxonomy (Benson et al., 2009; Sayers et al., 2009). The annotation and analysis category contains utilities to calculate sequence composition, codon usage, sequence length, and basic population genetic statistics. Additionally these utilities can also append the results of the analysis to the sequence description, which then can be used as selection and sorting criteria by the utilities in the selection category.

Learnability of the FAST tools is helped by making interface components such as specific options, consistent with the standard UNIX tools and across the FAST suite. Learning one FAST tool generally helps the user anticipate how to use others. In addition, specification of numerical ranges, regular expressions and other useful parameters follows standard Perl and UNIX conventions, all with the intent of making the tools fast and easy to learn.

### 3 IMPLEMENTATION DETAILS AND BENCHMARKING

All FAST utilities can process files or input on what is called the “standard input” stream. They all by default out to the “standard output” stream which may be connected to the “standard input” of another utility by a UNIX “pipe.” It is this latter facility that eases serial processing of data. Since most FAST utilities process sequences inline, they should mostly linear runtime complexity in number of sequences. Two exceptions to inline processing in FAST utilities concern `fassort` and `fastail` which both require some paging of data into temporary files. However, some preliminary benchmarking suggests that `fassort` runtimes also scale linearly in sequence number (fig 2).

FAST is compatible with the zero-based indexing if the sequence identifier is thought as the zero<sup>th</sup> field of the identifier line. This field must exist in Data selection in FAST is one-based as is conventional BioPerl coordinates and bioinformatics generally.

FAST supports automated logging to ease the reproducibility of workflows. The BioPerl backend of FAST was version 1.6.901 downloaded in January, 2012. BioPerl dependencies of FAST scripts were analyzed with the Cava packager (<http://www.cavapackager.com>). To fix some problems with I/O, the SeqIO components were updated to version 1.6.923 on June 4, 2014 and the Root components were updated on July 10, 2014 (github commit 50f87e9a4d). Some customizations of the BioPerl code-base were introduced, primarily to enable the translation of sequences containing gaps.

### 4 INSTALLATION AND USAGE EXAMPLES

#### 4.1 INSTALLATION AND DEPENDENCIES

FAST requires a working Perl installation and is distributed through the Comprehensive Perl Archive Network (CPAN). In a manual install, after download, installation follows standard Perl install procedure: `perl Makefile.PL; make; make test; (sudo) make install`. A small footprint of BioPerl dependencies has been packaged together in the FAST namespace. Other CPAN dependencies can be detected and installed by the `cpan` package manager. A fully automated install may on many systems be initiated by executing `perl -MCPAN -e 'install FAST'`.

## 4.2 SELECTING SEQUENCES BY ENCODED MOTIFS

84 An advantage of the annotation approach in FAST is the ability to select and sort sequences by attributes  
85 computed and annotated into data by utilities upstream in the pipeline. For example, to select protein-  
86 coding genes from a file `cds.fas` whose translations contain the *N*-glycosylation amino acid motif `?`,  
87 one could execute:

```
88 fasx1 -a cds.fas | fasgrep -t x10 "N[^P][ST][^P]" | fascut -f 1..-2
```

89 The first command in the pipeline translates each sequence and appends the translation to the description  
90 with the tag “x10” (indicating translation in the zeroth reading frame). The second command in the pipeline  
91 does a regex match using the specified motif pattern on the value of a “name:value” pair in the description  
92 with tag “x10”, hence processing the annotations produced by `fasx1`. The regex argument to `fasgrep` is  
93 quoted to protect the argument from interpretation by the shell. The last command in the pipeline removes  
94 the last field in the description, restoring records as they were before they were annotated by `fasx1`.

## 4.3 SORTING SEQUENCES BY THIRD CODON POSITION COMPOSITION

95 Another example illustrates the powerful expression of ranges in `fascut`. An optional “by” parameter in  
96 ranges allows increments or decrements in steps larger than one. To extract the third codon position bases  
97 from a gene, annotate their compositions, and sort sequences by their third-position adenosine contents,  
98 do:

```
99 fascut 1:-1:3 cds.fas | fascomp | fassort -nt comp_A
```

## 5 CONCLUDING REMARKS

100 Planned additions in future versions of FAST include `fasrand` and `alnrand` for sampling,  
101 permutations and bootstrapping of sequences and sites, respectively.

## AVAILABILITY

102 FAST is available through the Comprehensive Perl Archive Network (CPAN) at <http://search.cpan.org/~dhard/FAST>  
103

## DISCLOSURE/CONFLICT-OF-INTEREST STATEMENT

104 The authors declare that the research was conducted in the absence of any commercial or financial  
105 relationships that could be construed as a potential conflict of interest.

## AUTHOR CONTRIBUTIONS

106 D.H.A. conceived, designed, and wrote much of FAST. T.J.L. contributed major code factorizations and  
107 reorganization and `fastail`. K.T.K. contributed code including `faspaste`, and `fashead`. R.S.L.  
108 contributed an analysis of code dependencies for the FAST installer. All authors, especially D.L.C.  
109 and C.J.C., contributed documentation, testing, and code fixes. K.C.H.A. and D.H.A. wrote the FAST  
110 Cookbook. D.H.A. wrote the paper with major contributions from D.L.C. and T.J.L. All authors made

111 minor contributions to the manuscript, reviewed the final version of the manuscript and agree to be  
 112 accountable for its contents.

## ACKNOWLEDGEMENT

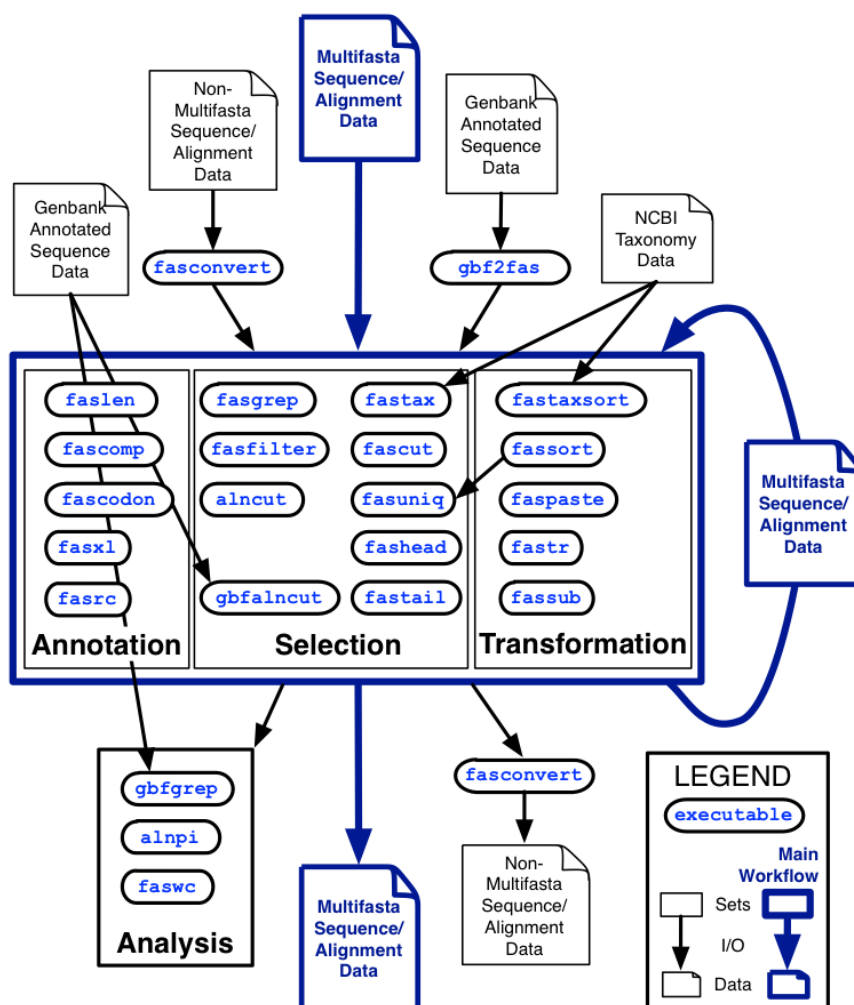
113 We acknowledge Christopher Clark for help in establishing a git repository for FAST. D.H.A. gratefully  
 114 acknowledges Professors Laura Landweber, Siv Andersson and Leif Kirsebom in whose laboratories the  
 115 FAST tools were first developed as well as the Linnaeus Centre for Bioinformatics at Uppsala University.

116 *Funding:* D.H.A. gratefully acknowledges an NSF-DBI Postdoctoral Fellowship in Biological  
 117 Informatics and awards to D.H.A. from UC Merced's Graduate Research Council and a Chancellor's  
 118 Award from UC Merced's second Chancellor Sung-Mo Kang.

## REFERENCES

- 119 Benson, D. A., Karsch-Mizrachi, I., Lipman, D. J., Ostell, J., and Sayers, E. W. (2009), GenBank., *Nucleic*  
 120 *acids research*, 37, Database issue, D26–31, doi:10.1093/nar/gkn723
- 121 Groves, T. and Godlee, F. (2012), Open science and reproducible research., *BMJ (Clinical research ed.)*,  
 122 344, jun26\_1, e4383, doi:10.1136/bmj.e4383
- 123 Peek, J. (2001), Why Use a Command Line Instead of Windows?, <http://www.linuxdevcenter.com/pub/a/linux/2001/11/15/learnunixos.html>
- 124 Rampp, M., Soddemann, T., and Lederer, H. (2006), The MIGenAS integrated bioinformatics toolkit for  
 125 web-based sequence analysis., *Nucleic acids research*, 34, Web Server issue, W15–9, doi:10.1093/nar/  
 126 gkl254
- 127 Rice, P., Longden, I., and Bleasby, A. (2000), EMBOS: The European Molecular Biology Open Software  
 128 Suite, *Trends in Genetics*, 16, 6, 276–277, doi:10.1016/S0168-9525(00)00204-2
- 129 Sayers, E. W., Barrett, T., Benson, D. A., Bryant, S. H., Canese, K., Chetvernin, V., et al. (2009), Database  
 130 resources of the National Center for Biotechnology Information., *Nucleic acids research*, 37, Database  
 131 issue, D5–15, doi:10.1093/nar/gkn741
- 132 Seemann, T. (2013), Ten recommendations for creating usable bioinformatics command line software.,  
 133 *GigaScience*, 2, 1, 15, doi:10.1186/2047-217X-2-15
- 134 Smith, S. W., Overbeek, R., Woese, C. R., Gilbert, W., and Gillevet, P. M. (1994), The genetic  
 135 data environment an expandable GUI for multiple sequence analysis., *Computer applications in the*  
 136 *biosciences : CABIOS*, 10, 6, 671–5
- 137 Stajich, J. E., Block, D., Boulez, K., Brenner, S. E., Chervitz, S. A., Dagdigian, C., et al. (2002), The  
 138 Bioperl toolkit: Perl modules for the life sciences., *Genome research*, 12, 10, 1611–8, doi:10.1101/gr.  
 139 361602
- 140 Stutz, M. (2000), Linux and the Tools Philosophy, [http://www.linuxdevcenter.com/pub/a/  
 141 linux/2000/07/25/LivingLinux.html](http://www.linuxdevcenter.com/pub/a/linux/2000/07/25/LivingLinux.html)

## FIGURES



**Figure 1.** FAST version 1.0 with data and workflow dependencies indicated.

**Figure 2.** Runtime benchmarks for `fassort`. Each point is the average of CPU runtime until completion over 10 runs with the Benchmark package. Left panel represents runs with a sequence length of 1000 bp. Right panel are runs with a sequence length of 100 bp.