



FAST: Fast Analysis of Sequences Toolbox

Travis J. Lawrence¹, Dana L. Carper¹, Katherine C.H. Amrine^{1,2}, Kyle Kauffman³, Claudia Canales³ and David H. Ardell^{1,*}

¹Quantitative and Systems Biology, University of California, Merced, CA, USA

²Present address: Department of Viticulture and Enology, University of California, Davis, CA, USA

³School of Natural Sciences, University of California, Merced, CA, USA

Correspondence*:

David Ardell

Quantitative and Systems Biology, University of California, Merced, 5200 North Lake Road, Merced, CA , 95343, USA, dardell@ucmerced.edu

ABSTRACT

Keywords: Text Text Text Text Text Text Text Text

1 INTRODUCTION

The field of molecular biology has changed significantly with the advent of next generation sequencing technology. It is now commonplace to analyze gigabases worth of data per experiment. Traditionally programs were developed for visualization and for basic sequence manipulation by a GUI interface (Smith et al., 1994; Rampp et al., 2006). Most available bioinformatic toolkits are designed for specific types of data or analysis requiring several toolkits to be installed. Moreover, each toolkit often requires a different file format making data analysis difficult.

The FAST utilities are modeled after the standard Unix toolkit (Peek, 2001), follow the Unix philosophy of “do one thing and do it well” (Stutz, 2000), and are written in PERL using bioperl packages (Stajich et al., 2002). This makes FAST utilities easy to adopt if you are familiar with the Unix toolkit and allows fast sequence analysis even on large datasets. FAST utilities have a uniform interface requiring FASTA formatted files and are capable of reading data from STDIN. This allows quick prototyping of sequence analysis problems by piping data between several utilities. Additionally, fasconvert can convert to/from fasta from/to several formats increasing the flexibility and usability of FAST. Extensive documentation has been developed for each utility along with useful error messages following the recommendations of Seemann (2013) to increase usability. Lastly, FAST is open source, which makes it available to anyone free of cost. This is in line with the call to make science more assessable, open, and reproducible by other scientists and the public (Groves and Godlee, 2012).

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 faspick 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

29 population genetic statistics. Additionally these utilities can also append the results of the analysis to the
30 sequence description, which then can be used as selection criteria by the utilities in the selection category.

31 Some utilities within FAST have overlapping function with those found within other toolkits. For
32 example sequence composition, sequence translation, and codon usage are available in the EMBOSS
33 package (Rice et al., 2000). Another example is the Bioinformatics Toolbox (White et al., 2014) that
34 has utilities to select only unique sequences and extract sequences from Genbank files based on gene
35 name, which have overlapping function with fasuniq and gbfeat2fas respectively. However, the utilities in
36 EMBOSS (Rice et al., 2000) and Bioinformatics Toolbox (White et al., 2014) lack a uniform interface,
37 are not modeled after the Unix toolkit, and do not have the ability to use regular expressions to select and
38 manipulate sequences. However, FAST also contains several utilities that have unique functionality. For
39 example gbalncut takes a multiple sequence alignment annotated with a genomic feature and a genbank
40 file and allows you to select certain regions of the alignment such as all the exons or the coding sequence.
41 Another example is fastaxsort that allows sorting of a multifasta file based on NCBI taxonomy (Benson
42 et al., 2009; Sayers et al., 2009).

2 DESCRIPTION

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

48 FAST is compatible with the zero-based indexing if the sequence identifier is thought as the zeroth
49 field of the identifier line. This field must exist in Data selection in FAST is one-based as is conventional
50 BioPerl coordinates and bioinformatics generally.

2.1 SELECTION UTILITIES

2.2 TRANSFORMATION UTILITIES

2.3 ANNOTATION AND ANALYSIS UTILITIES

2.4 USABILITY AND SCALABILITY

3 DISCUSSION

51 Text
52 Text Text Text Text. Additional Requirements:

3.1 DATA SHARING

53 Frontiers supports the policy of data sharing, and authors are advised to make freely available any
54 materials and information described in their article, and any data relevant to the article (while not
55 compromising confidentiality in the context of human-subject research) that may be reasonably requested
56 by others for the purpose of academic and non-commercial research. In regards to deposition of data and
57 data sharing through databases, Frontiers urges authors to comply with the current best practices within
58 their discipline.

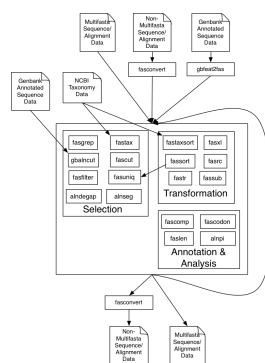


Figure 1. Enter the caption for your figure here. Repeat as necessary for each of your figures

92 Stutz, M. (2000), Linux and the Tools Philosophy, <http://www.linuxdevcenter.com/pub/a/linux/2000/07/25/LivingLinux.html>
93
94 White, B. P., Pilgrim, E. M., Boykin, L. M., Stein, E. D., and Mazor, R. D. (2014), Comparison of
95 four species-delimitation methods applied to a DNA barcode data set of insect larvae for use in routine
96 bioassessment, *Freshwater Science*, 33, 1, 338–348, doi:10.1086/674982

FIGURES