

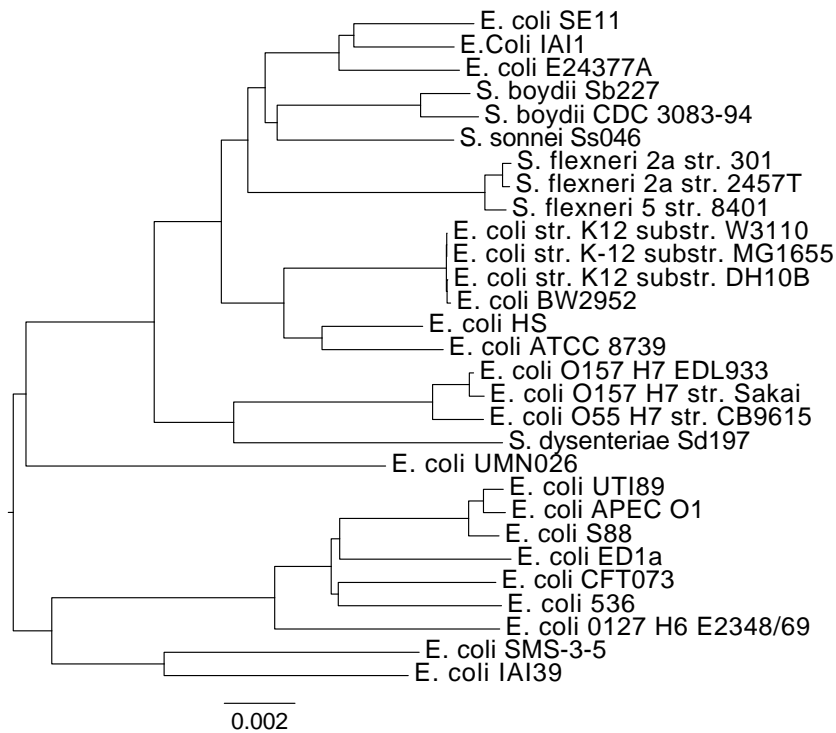
Documentation of ANDI

Rapid Estimation of Evolutionary Distances between Genomes

<https://github.com/EvolBioInf/andi>

Fabian Klötzl
kloetzl@evolbio.mpg.de

Version 0.9, 2015-06-02



Abstract

This is the documentation of the ANDI program for estimating the evolutionary distance between closely related genomes. These distances can be used to rapidly infer phylogenies for big sets of genomes. Because ANDI does not compute full alignments, it is so efficient that it scales well up to thousands of bacterial genomes.

This is scientific software. Please cite our paper [HKP15] if you use ANDI in your publication. Also refer to the paper for the internal of ANDI. Additionally, there is a Master's thesis with even more in depth analysis of ANDI [Klö15].

License

This document is released under the Creative Commons Attribution Share-Alike license. This means, you are free to copy and redistribute this document. You may even remix, tweak and build upon this document, as long as you credit me for the work I've done and release your document under the identical terms. The full legal code is available online: <https://creativecommons.org/licenses/by-sa/4.0/legalcode>.

Contents

1	Installation	5
1.1	Installing from Source Package	5
1.2	Installing without LIBDIVSUFSORT	5
1.3	Installing from Git Repository	6
2	Usage	7
2.1	Input	7
2.2	Output	7
2.3	Options	8
2.4	Example: ECO29	8
3	DevOps	11
3.1	Dependencies	11
3.2	Code Documentation	11
3.3	Unit Tests	11
3.4	Known Issues	12
3.5	Creating a Release	12

1 Installation

The easiest way to obtain your own copy of ANDI, is to download the latest release from GitHub. Please note, that ANDI requires `LIBDIVSUFSORT`¹ for optimal performance [Mor05]. If you wish to install ANDI without `LIBDIVSUFSORT`, consult Section 1.2.

ANDI runs best under a UNIX commandline such as `bash`. All following examples are also intended for that environment.

1.1 Installing from Source Package

Once you have downloaded the package, unzip it and change into the newly created directory.

```
~ % tar -xzf andi-0.9.tar.gz
~ % cd andi-0.9
```

Now build and install ANDI .

```
~/andi-0.9 % ./configure
~/andi-0.9 % make
~/andi-0.9 % sudo make install
```

This installs ANDI for all users on your system. If you do not have root privileges, you will find a working copy of ANDI in the `src` subdirectory. For the rest of this documentation, I will assume, that ANDI is in your `$PATH`.

Now ANDI should be ready for use. Try invoking the help.

```
~/andi-0.9 % andi -h
Usage: andi [-jrv] [-p FLOAT] FILES...
      FILES... can be any sequence of FASTA files. If no files are
      supplied, stdin is used instead.
Options:
  -j, --join    Treat all sequences from one file as a single genome
  -m, --low-memory Use less memory at the cost of speed
  -p <FLOAT>    Significance of an anchor pair; default: 0.05
  -r, --raw      Calculates raw distances; default: Jukes-Cantor
                  corrected
  -v, --verbose  Prints additional information
  -t <INT>       The number of threads to be used; default: 1
  -h, --help    Display this help and exit
  --version     Output version information and acknowledgments
```

ANDI also comes with a man page, which can be accessed via **man andi**.

1.2 Installing without LIBDIVSUFSORT

In case you cannot or do not want to install `LIBDIVSUFSORT`, ANDI comes with its own implementation of a *suffix array construction algorithm*. It is called `PSUFSORT` and is also available, separately.² To activate it for ANDI, proceed as described in Section 1.1, but with the following twist:

```
~/andi % ./configure --without-libdivsufsort
```

¹<https://github.com/y-256/libdivsufsort>

²<https://github.com/kloetzel/psufsort>

1 Installation

Please note that this requires a C++11 compiler. Also, PSUFSORT may be slower than LIBDIVSUFSORT and thus lead to inferior runtimes.

1.3 Installing from Git Repository

To build ANDI from the GIT repo, you will also need the AUTOTOOLS. Refer to your OS documentation for installation instructions. Once done, execute the following steps.

```
~ % git clone git@github.com:EvolBioInf/andi.git
~ % cd andi
~/andi % autoreconf -i
```

For old versions of AUTOCONF you may instead have to use `autoreconf -i -Im4`. Continue with the GNU trinity as described in Section 1.1.

2 Usage

The input sequences for ANDI should be in FASTA format. Any number of files can be passed. Each file may contain more than one sequence.

```
~ % andi S1.fasta S2.fasta
2
S1      0.0000 0.0979
S2      0.0979 0.0000
```

If no file argument is given, ANDI reads the input from STDIN. This makes it convenient to use in UNIX pipelines.

```
~ % cat S1.fasta S2.fasta | andi
2
S1      0.0000 0.0979
S2      0.0979 0.0000
```

The output of ANDI is a matrix in PHYLIP style: On the first line the number of compared sequences is given, 2 in our example. Then the matrix is printed, where each line is preceded by the name of the i th sequence. Note that the matrix is symmetric and the main diagonal contains only zeros. The numbers themselves are evolutionary distances, estimates as substitution rates.

2.1 Input

As mentioned before, ANDI reads in FASTA files. It recognizes only the four standard bases and is case insensitive (Regex: `[acgtACGT]`). All other residue symbols are excluded from the analysis and ANDI prints a warning, when this happens.

If instead of distinct sequences, a FASTA file contains contigs belonging to a single taxon, ANDI will treat them as a unit when switched into JOIN mode. This can be achieved by using the `-j` or `--join` command line switch.

```
~ % andi --join E_coli.fasta Shigella.fasta
[Output]
```

When the JOIN mode is active, the file names are used to label the individual sequences. Thus, in JOIN mode, each genome has to be in its own file, and furthermore, at least one filename has to be given via the command line.

If ANDI seems to take unusually long, or requires huge amounts of memory, then you might have forgotten the JOIN switch. This makes ANDI compare each contig instead of each genome, resulting in many more comparisons! To make ANDI output the number of genome it about to compare, use the `--verbose` switch.

2.2 Output

The output of ANDI is written to `stdout`. This makes it easy to use on the command line and within shell scripts. As seen before, the matrix, computed by ANDI, is given in PHYLIP format [Fel05].

```
~ % cat S1.fasta S2.fasta | andi
2
S1      0.0000 0.0979
S2      0.0979 0.0000
```

If the computation completed successfully, ANDI exits with the status code 0. Otherwise, the value of `errno` is used as the exit code. ANDI can also produce warnings and error messages for the user's convenience. These messages are printed to `stderr` and thus do not interfere with the normal output.

2.3 Options

ANDI takes a small number of commandline options, of which even fewer are of interest on a day-to-day basis. If **andi** -h displays a -t option, then ANDI was compiled with multi-threading support (implemented using OPENMP). By default ANDI uses all available processors. However, to restrict the number of threads, use -t.

```
~ % time andi ../test/1M.1.fasta -t 1
2
S1      0.0000 0.0995
S2      0.0995 0.0000
./andi ../test/1M.1.fasta 0,60s user 0,01s system 99% cpu 0,613 total
~ % time andi ../test/1M.1.fasta -t 2
2
S1      0.0000 0.0995
S2      0.0995 0.0000
./andi ../test/1M.1.fasta -t 2 0,67s user 0,03s system 195% cpu 0,362
total
```

In the above examples the runtime dropped from 0.613 s, to 0.362 s using two threads. Giving ANDI more threads than input genomes leads to no further speed improvement. The other important option is `--join` (see Section 2.1).

By default, the distances computed by ANDI are *Jukes-Cantor* corrected [JC69]. This means, the output is substitution rates, which is greater than the mismatch rate. To obtain the pure mismatch rate, use the `--raw` switch.

2.4 Example: ECO29

Here follows a real-world example of how to use ANDI. It makes heavy use of the commandline and tools like PHYLIP. If you prefer R, check out this excellent blog post by Kathryn Holt.

As a data set we use ECO29; 29 genomes of *E. Coli* and *Shigella*. You can download the data here: <http://guanine.evolbio.mpg.de/andi/eco29.fasta.gz>. The genomes have an average length of 4.9 million nucleotides amounting to a total 138 MB.

ECO29 comes a single FASTA file, where each sequence is a genome. To calculate their pairwise distances, enter

```
~ % andi eco29.fasta > eco29.mat
andi: The input sequences contained characters other than acgtACGT.
These were automatically stripped to ensure correct results.
```

The ECO29 data set includes non-canonical nucleotides, such as Y, N, and P, which get stripped from the input sequences. The resulting matrix is stored in `eco29.mat`; Here is a small excerpt:

```
~ % head -n 5 eco29.mat | cut -d ' ' -f 1-5
29
gi|563845 0.0000e+00 1.8388e-02 1.8439e-02 2.6398e-02
gi|342360 1.8388e-02 0.0000e+00 4.4029e-04 2.6166e-02
gi|300439 1.8439e-02 4.4029e-04 0.0000e+00 2.6123e-02
gi|261117 2.6398e-02 2.6166e-02 2.6123e-02 0.0000e+00
```

From this we compute a tree via neighbor-joining using a PHYLIP wrapper called EMBASSY.¹

¹<http://emboss.sourceforge.net/embassy/#PHYLIP>


```
~ % fneighbor -datafile eco29.mat -outfile eco29.phylipdump
```

To make this tree easier to read, we can midpoint-root it.

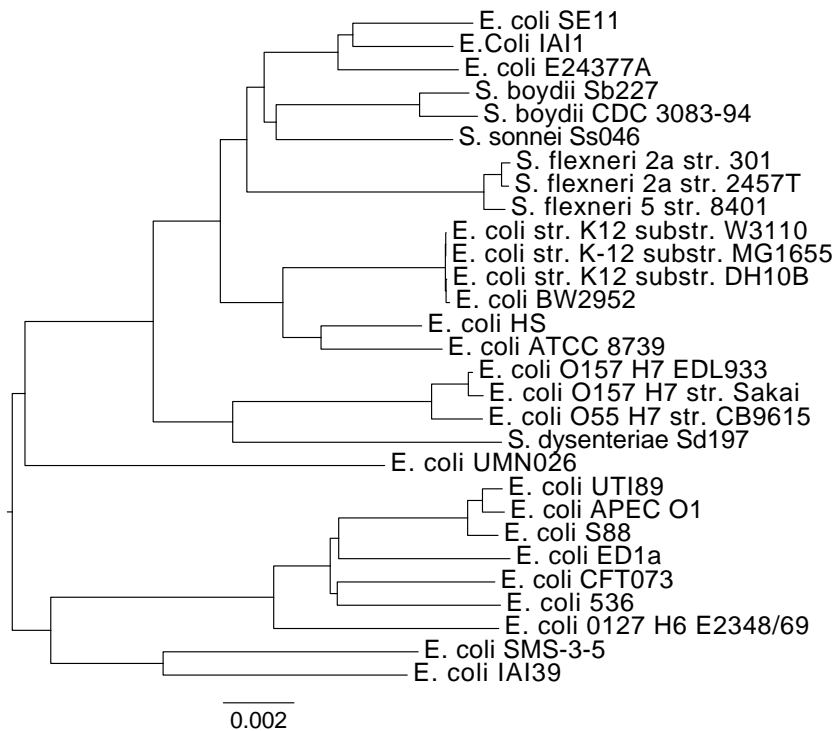
```
~ % fretree -spp 29 -intreefile eco29.treefile -outtreefile eco29.  
tree <<EOF
```

```
M  
X  
Y  
R  
EOF
```

The file `eco29.tree` now contains the tree in Newick format. This can be plotted using [Ram15]

```
~ % figtree eco29.tree &
```

to yield



3 DevOps

ANDI is written in C/C++; mostly C99 with some parts in C++11. The sources are released on GITHUB as *free software* under the GNU GENERAL PUBLIC LICENSE VERSION 3 [Fre07]. Prebundled packages using AUTOCONF are also available, with the latest release being v0.9 at the time of writing.

If you are interested in the internals of ANDI, consult the paper [HKP15] or my master thesis [Klö15]. Both explain the used approach in detail. The latter emphasizes the used algorithms, data structures and their efficient implementation.

3.1 Dependencies

Here is a complete list of dependencies required for developing ANDI.

- A C and a C++11 compiler,
- the AUTOTOOLS,
- PDLATEX with various packages for the manual,
- GIT,
- GLIB2 for the unit tests,
- DOXYGEN,
- and LIBDIVSUFSORT.

3.2 Code Documentation

Every function in ANDI is documented using DOXYGEN style comments. To create the documentation run **make** code-docs in the main directory. You will then find the documentation under ./docs.

3.3 Unit Tests

The unit tests are located in the ANDI repository under the ./test directory. Because they require GLIB2, and a C++11 compiler, they are deactivated by default. To enable them, execute

```
~/andi % ./configure --enable-unit-tests
```

during the installation process. You can then verify the build via

```
~/andi % make check
```

The unit tests are also checked each time a commit is sent to the repository. This is done via TRAVISCI.¹ Thus, a warning is produced, when the builds fail, or the unit tests do not run successfully. Currently, the unit tests cover more than 80% of the code. This is computed via the TRAVIS builds and a service called COVERALLS.²

¹<https://travis-ci.org/EvolBioInf/andi>

²<https://coveralls.io/r/EvolBioInf/andi>

3.4 Known Issues

These minor issues are known. I intend to fix them, when I have time.

1. Stripping the filename from its path in join mode will not work properly under Windows because of the different path separator.
2. I currently do not check for integer multiplication overflows in mallocs. To fix this, I might be using `reallocarray` from BSD.

3.5 Creating a Release

A release should be a stable version of ANDI with significant improvements over the last version. For bug fixes, `dotdot` release may be used.

Once ANDI is matured, the new features implemented, and all tests were run, a new release can be created. First, increase the version number in `configure.ac`. Commit that change in git, and tag this commit with `vX.y`. Create another commit, where you set the version number to the next release (e.g., `vX.z-beta`). This assures that there is only one commit and build with that specific version. Now return to the previous commit `git checkout vX.y`.

Ensure that ANDI is ready for packaging with AUTOCONF.

```
~ % make distcheck
make dist-gzip am__post_remove_distdir='@:'
make[1]: Entering directory '/home/kloetzl/Projects/andi'
if test -d "andi-0.9.1-beta"; then find "andi-0.9.1-beta" -type d !
  -perm -200 -exec chmod u+w {} ';' && rm -rf "andi-0.9.1-beta" || {
  sleep 5 && rm -rf "andi-0.9.1-beta"; }; else ;; fi
test -d "andi-0.9.1-beta" || mkdir "andi-0.9.1-beta"
(cd src && make top_distdir=../andi-0.9.1-beta distdir=../andi
  -0.9.1-beta/src \
  am__remove_distdir=: am__skip_length_check=: am__skip_mode_fix=:
  distdir)

... Loads of output ...

=====
andi-0.9.1-beta archives ready for distribution:
andi-0.9.1-beta.tar.gz
=====
```

If the command does not build successfully, no tarballs will be created. This may necessitate further study of AUTOCONF and AUTOMAKE.

Bibliography

- [Fel05] J. Felsenstein. Phylip (phylogeny inference package). Distributed by the author, 2005. Department of Genome Sciences, University of Washington.
- [Fre07] Free Software Foundation. Gnu general public license, 2007. <https://gnu.org/licenses/gpl.html>.
- [HKP15] Bernhard Haubold, Fabian Klötzl, and Peter Pfaffelhuber. andi: Fast and accurate estimation of evolutionary distances between closely related genomes. *Bioinformatics*, 31(8):1169–1175, 2015.
- [JC69] T. H. Jukes and C. R. Cantor. Evolution of protein molecules. *Mammalian protein metabolism*, 3:21–132, 1969.
- [Klö15] Fabian Klötzl. Efficient estimation of evolutionary distances. Master’s thesis, University of Lübeck, 2015.
- [Mor05] Yuta Mori. Short description of improved two-stage suffix sorting algorithm, 2005. <http://homepage3.nifty.com/wpage/software/itssort.txt>.
- [Ram15] Andrew Rambaut. Figtree, accessed 2015. <http://tree.bio.ed.ac.uk/software/figtree/>.