

Introduction to molecular phylogenetic reconstruction with the ETE3 Toolkit

Frank Aylward

Abstract

This is a short tutorial on how to get started analyzing FASTA files via the command line.

Code is intended for use on an Ubuntu 16.04 LTS OS, but it may work on other Unix or Unix-like systems.

Here we will use the ETE3 toolkit, which is a very nice tool for phylogenetic reconstruction. The main website is here: <http://etetoolkit.org/>

We will also look at the EggNOG database, which is a very useful database of protein families: <http://eggnogdb.embl.de/#/app/home>

Citation: Frank Aylward Introduction to molecular phylogenetic reconstruction with the ETE3 Toolkit. **protocols.io** [dx.doi.org/10.17504/protocols.io.qhydt7w](https://doi.org/10.17504/protocols.io.qhydt7w)

Published: 01 Jun 2018

Protocol

Download a working directory

Step 1.

This tutorial is designed to provide experience creating molecular phylogenies of marker genes. You may find yourself in a position where you have a gene of interest and you wish to know it's phylogenetic relationships to other known genes of the same gene/protein family. Marker genes can be markers of particular metabolic processes (functional marker genes) or markers for phylogenetic diversity (phylogenetic marker genes). An example of the former is *nifH*, which is a gene which encodes for a core component of the nitrogenase enzyme that is responsible for nitrogen fixation. If an organism has a *nifH* gene, there is a pretty good chance that it has the capacity to fix dinitrogen gas into ammonia. A good example of a phylogenetic marker gene would be *rpoB*, which encodes for the beta subunit of RNA polymerase and is a highly-conserved protein found across Bacteria, Archaea, and Eukaryotes, and can therefore be used to analyze phylogenetic relationships between many different organisms. From the perspective of microbial ecology, marker gene surveys can be useful

methods to learn about what kinds of microbes are present in a given environment, and what metabolic activities they may be engaged in.

Here we will practice creating a reference phylogeny for the marker gene narG, which encodes a key subunit in the nitrate reductase complex.

First let's download some practice data for GitHub:

git clone <https://github.com/faylward/bioinfo-tutorials>

You should see a folder called 'bioinfo-tutorials', and inside that should be another folder called 'marker_gene_phylogenetics'. Let's navigate to that folder:

cd bioinfo-tutorials/marker_gene_phylogenetics

Inside this folder you should see a file called 'mystery_protein.faa' which contains the amino acid sequence of a single protein. Take a look in this file just to be sure:

head mystery_protein.faa

You can imagine that you found this protein encoded in a genome that you are analyzing and you think it may have been acquired through lateral gene transfer. Perhaps you did some gene annotation and you know it encodes for a NarG protein, but you don't know anything else about it. Creating a phylogeny with this protein and some reference sequences is a great way to investigate possible lateral gene transfers and see what other sequences are closely related. Or perhaps you found this gene in a particular environmental gene survey and you are interested in knowing what reference sequences are closely related. Either way, a phylogeny is informative.

```

frankayward@AYLWARD-9H6YHK2:~$
frankayward@AYLWARD-9H6YHK2:~$ git clone https://github.com/faylward/bioinfo-tutorials
Cloning into 'bioinfo-tutorials'...
remote: Counting objects: 36, done.
remote: Compressing objects: 100% (34/34), done.
remote: Total 36 (delta 8), reused 8 (delta 8), pack-reused 8
Unpacking objects: 100% (36/36), done.
Checking connectivity... done.
frankayward@AYLWARD-9H6YHK2:~$
frankayward@AYLWARD-9H6YHK2:~$ cd bioinfo-tutorials/marker_gene_phylogenetics/
frankayward@AYLWARD-9H6YHK2:~/bioinfo-tutorials/marker_gene_phylogenetics$
frankayward@AYLWARD-9H6YHK2:~/bioinfo-tutorials/marker_gene_phylogenetics$ head mystery_protein.faa
>Unccharacterized_Narg
MSNNKLSHWROEJAPENSKWEEFYRHRWQYDKVVRSTHGVNCTGSCCTQIHVKDGIVTMERQCLDYPKLEHGIPPEPRGCGRGISFSHYLYSPLRVKYPTIRCVLLDLHLKAKKEFDOPVEAWQSMVNNPVKRRWRARGKGGFRRTSHDTYKEIIAASMAHTIKQHGPDRIACFSPIPAKSKISY
ASGARLNLQIGGVSLSFYDNYCDLPMASPEAMGEQTDVQESADNYNAKMLAVMGANLNRTRTPDVHFAEAARYNGTKMNVVSPDFSQVSKFSDENIPINAGQDGAYWNAVNHVILKEFHHEKQESFLDYSKQYTDSPFLVELTKDGD SYKAGQLLRANRLNQYKDMENGQWQFLMWDEEENRTKVPK
GSGVNRWAKENKQHMNLKLEDNADGSKTHPKLTLKDSQHVSVELDDFGQGETITRYVPVKMIETADGQVPVATVYDLMLAQYQVGRGLEGYPENYEEDATYTPAWAEKYTGIAETLIRFAREMASTAETKGGKCTVIIAGGINHYHANLMYRAAIALHMF TGCICKNGGLAHYVQGKELAPGE
PWTATLAKDWFGPSRVQNPASWHYVHSDQHRYEKSF TDYHTVPEKQENTLAKGHTIDINVKAVKNGWLPFPYQVEANPLALAKEARENGATDEKSIDYVNVNKLKNKELKYSIEDPNSDNWPRVWFVWRGNAIGGSAGKHEFFLDHYLGTHNQVSDPQFAKGSTDEVVHDKVPQKDLVVDL
NFRMDTSALYSIDVLPAAATWYEKADLNSDMHSFIHPLAEAVPPAMESKSDMDIFGSAETFSKMSKHPKQMEDVVSVALAHDPAEIAQKEIRDQSQCEVDVPGKTHPGMKVVKRDYKNLYKRYISYCPNVPANGLGAGHTHYEVKDEYDLKIECPTETMGQQTYP SLKESKDV CNTILELAT
VTNGELAYRSYKMEERTGLVLADLAENKRGVRHSYDDLCSPRRLHNSPMWSGLENGRAYSPFTYNNVERLVPWRTLGRQHFLYDHPGYIQYGEHLPTYKPKPTPTQYAEATSESEHKTKIFNYLTPHGKMHISTYCDNHRMSTLSRCVEPFWINDKDAEELDIDVNDWVEVYNDNGVVVTRAS
VSARIPRGIVIIYHATERTLSVPKSPLRGNKRCNNLSLTRLRKPNLHVGGYGGFTYHFNYWPTONNRDTFIMVRKLPVLNW*
frankayward@AYLWARD-9H6YHK2:~/bioinfo-tutorials/marker_gene_phylogenetics$
frankayward@AYLWARD-9H6YHK2:~/bioinfo-tutorials/marker_gene_phylogenetics$

```

Download a FASTA file of protein sequences

Step 2.

So let's proceed with getting some reference sequences. The EggNOG database has a very nice compilation of marker proteins, so we can download sequences from there:

```
wget -O narg.faa.gz http://eggnogapi.embl.de/nog_data/text/fasta/COG5013
```

And because the file is gzipped, let's uncompress it:

```
gunzip narg.faa.gz
```

Go ahead and take a look at the narg.faa file and make sure it's in the format you think it should be in (FASTA amino acid).

Some simple poking around and basic QC is always a good idea:

```
head narg.faa
```

```
grep -c '^>' narg.faa
```

I also like using seqtk and datamash on combination to get some basic stats about the proteins involved (see 'Introduction to analyzing FASTA files' for details here):

this:

```
seqtk comp narg.faa | awk '$2>755' | cut -f 1 > narg.long_proteins.list
```

This command does the following: 1) It takes the proteins in the narg.faa file and gets their general stats with seqtk, 2) uses an AWK command to filter through the seqtk comp output and provide only those lines where the second column value is > 755, and 3) then cuts out the first column and puts that in a file called narg.long_proteins.list.

Now we can use 'seqtk subseq' to get a FASTA file of only the proteins longer than 755 amino acids:

```
seqtk subseq narg.faa narg.long_proteins.list > narg.long_proteins.faa
```

How many proteins did we filter out with our length cutoff? We can check with:

```
grep -c '^>' narg.faa
```

```
grep -c '^>' narg.long_proteins.faa
```

I got 393 and 367, so we removed 26 of the shortest proteins with our length cutoff.

```
frankaylward@AYLWARD-9H6YHK2:~/bioinfo-tutorials/marker_gene_phylogenetics$
frankaylward@AYLWARD-9H6YHK2:~/bioinfo-tutorials/marker_gene_phylogenetics$ seqtk comp narg.faa | awk '$2>755' | cut -f 1 > narg.long_proteins.list
frankaylward@AYLWARD-9H6YHK2:~/bioinfo-tutorials/marker_gene_phylogenetics$ ls -l
total 472
-rw-rw-r-- 1 frankaylward frankaylward 1236 Jun  1 15:34 mystery_protein.faa
-rw-rw-r-- 1 frankaylward frankaylward 470958 Jun  1 15:39 narg.faa
-rw-rw-r-- 1 frankaylward frankaylward 6747 Jun  1 15:41 narg.long_proteins.list
frankaylward@AYLWARD-9H6YHK2:~/bioinfo-tutorials/marker_gene_phylogenetics$
frankaylward@AYLWARD-9H6YHK2:~/bioinfo-tutorials/marker_gene_phylogenetics$ seqtk subseq narg.faa narg.long_proteins.list > narg.long_proteins.faa
frankaylward@AYLWARD-9H6YHK2:~/bioinfo-tutorials/marker_gene_phylogenetics$
frankaylward@AYLWARD-9H6YHK2:~/bioinfo-tutorials/marker_gene_phylogenetics$ ls
mystery_protein.faa  narg.faa  narg.long_proteins.faa  narg.long_proteins.list
frankaylward@AYLWARD-9H6YHK2:~/bioinfo-tutorials/marker_gene_phylogenetics$ grep -c '^>' narg.faa
393
frankaylward@AYLWARD-9H6YHK2:~/bioinfo-tutorials/marker_gene_phylogenetics$
frankaylward@AYLWARD-9H6YHK2:~/bioinfo-tutorials/marker_gene_phylogenetics$ grep -c '^>' narg.long_proteins.faa
367
frankaylward@AYLWARD-9H6YHK2:~/bioinfo-tutorials/marker_gene_phylogenetics$
frankaylward@AYLWARD-9H6YHK2:~/bioinfo-tutorials/marker_gene_phylogenetics$
```

Remove redundant sequences from the reference file

Step 4.

Sometimes you will find that the reference sequences are highly redundant, meaning that a large number of sequences in the file are identical or nearly identical. Usually we are interested in getting a pretty broad view of phylogenetic relationships at this stage, so it may be worth while to de-replicate

the reference sequences rather than spend a lot of time later computing phylogenies of sequences that are 99% identical.

For sequence-based dereplication a handy tool is CD-HIT. You should be able to install this tool fairly easily using:

```
sudo apt install cd-hit
```

And if you just type cd-hit into the command line afterwards you should be able to view the many different options that this tool provides.

By default cd-hit will take a FASTA file of proteins and identify clusters that are 90% or more identical over 90% the length of the shorter protein. These sequences are then grouped together, and one (usually the longest protein) is chosen as a 'representative'. This way we can just use cluster representatives in our final analysis rather than using every single sequence. Note that the % identity and % overlap thresholds can be modified using various flags in the help menu, so you could opt to remove redundancy only at the 99% amino acid identity level if you were interested in removing only the very very similar sequences, for example.

Let's try a simple command with the default parameters:

```
cdhit -i narg.long_proteins.faa -o narg.nr.faa
```

CD-HIT will print out some general log of what it's doing before finishing. You can check out the files that were created with 'ls -l':

You should see the output file we specified 'all_narg.nr.faa', which will have the cluster representatives, as well as a .clstr file which contains information about what sequences were clustered together. We will work mainly with the .nr.faa file.

Now let's check to see how many sequences were in the original file, and how many fewer are in new .nr.faa file (i.e, how many cluster representatives there are compared to raw sequences).

```
grep -c '^>' narg.long_proteins.faa
```

```
grep -c '^>' narg.nr.faa
```

I got 367 and 226, which indicates we removed 141 sequences. This will save us some time when we generate the phylogeny in the next steps.

Before we proceed, let's combined our non-redundant reference sequences with our mystery sequence so that we have one consolidated FASTA file to use for phylogenetic reconstruction:

```
cat mystery_protein.faa narg.nr.faa > all_narg.faa
```

Overall the commands should look something like this:

```
frankaylward@AYLWARD-9H6YHK2:~/bioinfo-tutorials/marker_gene_phylogenetics$
frankaylward@AYLWARD-9H6YHK2:~/bioinfo-tutorials/marker_gene_phylogenetics$ cdhit -i narg.long_proteins.faa -o narg.nr.faa
=====
Program: CD-HIT, V4.6 (+OpenMP), Jan 23 2016, 05:09:49
Command: cdhit -i narg.long_proteins.faa -o narg.nr.faa

Started: Fri Jun 1 15:27:52 2018
=====
Output
-----
total seq: 367
longest and shortest : 1296 and 848
Total letters: 451933
Sequences have been sorted

Approximated minimal memory consumption:
Sequence      : 0M
Buffer        : 1 X 10M = 10M
Table         : 1 X 65M = 65M
Miscellaneous  : 0M
Total         : 76M

Table limit with the given memory limit:
Max number of representatives: 537308
Max number of word counting entries: 90422116

comparing sequences from      0 to      367
367 finished      226 clusters

Apprximated maximum memory consumption: 78M
writing new database
writing clustering information
program completed !

Total CPU time 0.15
frankaylward@AYLWARD-9H6YHK2:~/bioinfo-tutorials/marker_gene_phylogenetics$ grep -c '^>' narg.long_proteins.faa
367
frankaylward@AYLWARD-9H6YHK2:~/bioinfo-tutorials/marker_gene_phylogenetics$ grep -c '^>' narg.nr.faa
226
frankaylward@AYLWARD-9H6YHK2:~/bioinfo-tutorials/marker_gene_phylogenetics$
frankaylward@AYLWARD-9H6YHK2:~/bioinfo-tutorials/marker_gene_phylogenetics$
frankaylward@AYLWARD-9H6YHK2:~/bioinfo-tutorials/marker_gene_phylogenetics$ cat mystery_protein.faa narg.nr.faa > all_narg.faa
frankaylward@AYLWARD-9H6YHK2:~/bioinfo-tutorials/marker_gene_phylogenetics$
```

Generate the phylogeny

Step 5.

Now that we have our final NarG amino acid file, we can use it as input for ete3, which is a very useful tool that integrates different programs for sequence alignment, alignment trimming, and phylogenetic reconstruction into one interface.

Let's try the following command:

ete3 build -a all_narg.faa -o narg_phylogeny -w standard_trimmed_fasttree --sname-delimiter -C 16

This will take a bit to run (5-10 minutes), and it will provide a running log of the processes as they run. A few notes on the flags:

-a is the input. If we had wanted to use nucleic acids we would have used -d.

-o is the output directory. This is where all of the output files will go.

-w is the workflow name. There are many different workflows that are possible, depending on what programs you wish to use for alignment, trimming, and phylogenetic reconstruction. This particular workflow uses Clustal Omega for alignment, Trimal for alignment trimming, and FastTree for phylogenetic inference. A full list of options is available on the ete3 toolkit website:

At the end you should get some command line output that looks like this:

```
INFO - Launched 0 jobs. 1(R), 0(W). Cores usage: 2/8
INFO - Updating tasks status: (Fri Jun 1 14:50:53 2018)
INFO - Thread clustalo default-trinal01-none-fasttree_full: pending tasks: 1 of sizes: 227
INFO - (R) TreeTask (227 aa seqs, FastTree, /clustalo d...tree_full)
INFO - Waiting 2 seconds
INFO - Launched 0 jobs. 1(R), 0(W). Cores usage: 2/8
INFO - Updating tasks status: (Fri Jun 1 14:50:54 2018)
INFO - Thread clustalo default-trinal01-none-fasttree_full: pending tasks: 1 of sizes: 227
INFO - (R) TreeTask (227 aa seqs, FastTree, /clustalo d...tree_full)
INFO - Waiting 2 seconds
INFO - Updating tasks status: (Fri Jun 1 14:50:56 2018)
INFO - Thread clustalo default-trinal01-none-fasttree_full: pending tasks: 1 of sizes: 227
INFO - (R) TreeTask (227 aa seqs, FastTree, /clustalo d...tree_full)
INFO - Waiting 2 seconds
INFO - Launched 0 jobs. 1(R), 0(W). Cores usage: 2/8
INFO - Updating tasks status: (Fri Jun 1 14:50:58 2018)
INFO - Thread clustalo default-trinal01-none-fasttree_full: pending tasks: 1 of sizes: 227
INFO - (R) TreeTask (227 aa seqs, FastTree, /clustalo d...tree_full)
INFO - Waiting 2 seconds
INFO - Launched 0 jobs. 1(R), 0(W). Cores usage: 2/8
INFO - Updating tasks status: (Fri Jun 1 14:51:00 2018)
INFO - Thread clustalo default-trinal01-none-fasttree_full: pending tasks: 1 of sizes: 227
INFO - (R) TreeTask (227 aa seqs, FastTree, /clustalo d...tree_full)
INFO - Waiting 2 seconds
INFO - Updating tasks status: (Fri Jun 1 14:51:02 2018)
INFO - Thread clustalo default-trinal01-none-fasttree_full: pending tasks: 1 of sizes: 227
INFO - (R) TreeTask (227 aa seqs, FastTree, /clustalo d...tree_full)
INFO - Waiting 2 seconds
INFO - Launched 0 jobs. 0(R), 0(W). Cores usage: 0/8
INFO - Updating tasks status: (Fri Jun 1 14:51:04 2018)
INFO - Thread clustalo default-trinal01-none-fasttree_full: pending tasks: 1 of sizes: 227
INFO - (W) TreeMergeTask (227 aa seqs, TreeMerger, /clustalo d...tree_full)
INFO - (D) TreeMergeTask (227 aa seqs, TreeMerger, /clustalo d...tree_full)
INFO - Waiting 2 seconds
INFO - Assembling final tree...
INFO - Done thread clustalo default-trinal01-none-fasttree_full in 1 iteration(s)
INFO - Writing final tree for clustalo default-trinal01-none-fasttree_full
INFO - /home/frankayward/bioinfo-tutorials/marker_gene_phylogenetics/narg_phylogeny/clustalo_default-trinal01-none-fasttree_full/all_narg.faa.final_tree.nw
INFO - /home/frankayward/bioinfo-tutorials/marker_gene_phylogenetics/narg_phylogeny/clustalo_default-trinal01-none-fasttree_full/all_narg.faa.final_tree.nwx (newick extended)
INFO - Writing final tree alignment clustalo default-trinal01-none-fasttree_full
INFO - /home/frankayward/bioinfo-tutorials/marker_gene_phylogenetics/narg_phylogeny/clustalo_default-trinal01-none-fasttree_full/all_narg.faa.final_tree.used_alg.faa
INFO - Writing root node alignment clustalo default-trinal01-none-fasttree_full
INFO - /home/frankayward/bioinfo-tutorials/marker_gene_phylogenetics/narg_phylogeny/clustalo_default-trinal01-none-fasttree_full/all_narg.faa.final_tree.faa
INFO - Writing root node trimmed alignment clustalo default-trinal01-none-fasttree_full
INFO - /home/frankayward/bioinfo-tutorials/marker_gene_phylogenetics/narg_phylogeny/clustalo_default-trinal01-none-fasttree_full/all_narg.faa.final_tree.trimmed.faa
INFO - Generating tree image for clustalo default-trinal01-none-fasttree_full
INFO - /home/frankayward/bioinfo-tutorials/marker_gene_phylogenetics/narg_phylogeny/clustalo_default-trinal01-none-fasttree_full/all_narg.faa.final_tree.png
INFO - Launched 0 jobs. 0(R), 0(W). Cores usage: 0/8
INFO - Launched 0 jobs. 0(R), 0(W). Cores usage: 0/8
INFO - Done
INFO - Deleting temporal data...

=====
The following published software and/or methods were used.
*** Please, do not forget to cite them! ***
=====
Stievers F, Wlin A, Dineen D, Gibson TJ, Karplus K, Li W, Lopez R,
McWilliam H, Remmert M, Söding J, Thompson JD, Higgins DG. Fast,
scalable generation of high-quality protein multiple sequence
alignments using Clustal Omega. Mol Syst Biol. 2011 Oct 11;7:539.
doi: 10.1038/msb.2011.75.
Capella-Gutierrez S, Silla-Martinez JM, Gabaldón T. trimAl: a tool for
automated alignment trimming in large-scale phylogenetic analyses.
Bioinformatics. 2009 Aug 1;25(15):1972-3.
Huerta-Cepas J, Serra F, Bork P. ETE 3: Reconstruction, analysis and
visualization of phylogenomic data. Mol Biol Evol (2016) doi:
10.1093/molbev/msw046
Price MN, Dehal PS, Arkin AP. FastTree 2 - approximately maximum-
likelihood trees for large alignments. PLoS One. 2010 Mar
10;5(3):e9490.
frankayward@AYLWARD-9HG6YHK2:~/bioinfo-tutorials/marker_gene_phylogenetics$
frankayward@AYLWARD-9HG6YHK2:~/bioinfo-tutorials/marker_gene_phylogenetics$
frankayward@AYLWARD-9HG6YHK2:~/bioinfo-tutorials/marker_gene_phylogenetics$
frankayward@AYLWARD-9HG6YHK2:~/bioinfo-tutorials/marker_gene_phylogenetics$
```


At the end you should have a directory called `narg_phylogeny`, and inside that should be another directory called something like `'clustalo_default-trimal01-none-fasttree_full'`. Inside that directory will be the `.nwk` and `.fa` files with the trees and alignments, respectively, as well as several visualizations of the trees.

If you look at the figures it may take you a bit to find the mystery NarG protein we started with- it should be called 'uncharacterized narg' since that is the name that it had in the `mystery_protien.faa` file. Once you locate it you can see what other proteins are similar. Here, since we are using sequences from the EggNOG database, you will want to look up what some of the species codes stand for. For example, two of the proteins that is similar to our mystery protein have names that start with the prefix 314278.NB231, which, if we look up on EggNOG, is the species code for *Nitrococcus mobilis*, which is a marine nitrate-reducing bacterium. So the presence of NarG in this organism makes sense given what is known about its physiology.

