

Odds and Ends

Last day of class is usually a mix of trying to find gaps or fill in holes.

I discussed phylogenetic tree building.

```
git clone https://github.com/biodataprogram/GEN220_2020_examples.git
cd GEN220_2020_examples/Trees
module load muscle
module load fasttree
module load hmmer/3
module load IQ-TREE/2.1.1
module load trimal
```

```
# build an alignment of sequences already identified as homologs
# previously I had started with MET12 (S. cerevisiae) enzyme
more MET12.fa # single sequence
ls -l MET12.hit_seqs.fasta # the collection of homologs for MET12 in a few yeast fungi
# denovo multiple alignment - writes in multi-fasta format
muscle -in MET12.hit_seqs.fasta -out MET12.hit_seqs.fasaln
# denovo multiple alignment - writes in multi-fasta format - writes in Clustal format
muscle -in MET12.hit_seqs.fasta -out MET12.hit_seqs.fasaln.clw -clw
# trim sequences - using automated parameters - see http://trimal.cgenomics.org/trimal for more
trimal -automated1 -in MET12.hit_seqs.fasaln -out MET12.hit_seqs.mfa.trim

# build a tree w fasttree (FastTreeMP uses multiple processors, FastTree uses 1 processor only)
FastTreeMP < MET12.hit_seqs.fasaln > MET12.hit_seqs.tre
# build a tree with IQ-TREE2 - ultrafast bootstrap and first determine optimal number of processors
iqtree2 -s MET12.hit_seqs.fasaln -nt AUTO -bb 1000 -alrt 1000
```

Some links * [Muscle](#) - Multiple alignment tool * [TrimAl](#) - alignment trimming tool * [HMMER](#) - HMMER - Hidden Markov Model for biosequence analyses. * [FastTree](#) - Fast Phylogenetic Tree construction * [IQ-TREE](#) - Phylogenetic Tree construction * [RAxML; a tutorial](#) * [iTOL](#) - Tree visualization (web-based) tool * [FigTree](#) - Tree visualization (can run on HPCC if you have X11 enabled: `module load figtree; figtree`) * [ggtree](#) - R package for Tree rendering

I also showed how to use HMMER and hmmbuild

```
module load hmmer/3
```

```
# build an HMM from a multiple alignment
hmmbuild MET12.hmm MET12.hit_seqs.fasaln
```

This is a little circular I am searching the HMM back against the original sequences, but if you wanted to instead search this HMM against a database of proteins (eg swissprot or your collection of proteins from species)

```
module load hmmer/3
```

```
# domtbl is the result file which has columns of data that are parseable instead of more complex  
hmmsearch -E 1e-3 --domtblout MET12.search.domtbl MET12.hmm DATABASE > MET12.search.hmmsearch  
  
# to align a set of proteins back to an HMM (which is instead of doing a denovo multiple alignment)  
hmmalign MET12.hmm MET12.hit_seqs.fasta > MET12.hit_seqs.stk  
# convert the stockholm format to multifasta  
esl-reformat afa MET12.hit_seqs.stk > MET12.hit_seqs.hmmalign.fasaln  
# convert the stockholm format to clustal  
esl-reformat clustal MET12.hit_seqs.stk > MET12.hit_seqs.hmmalign.fasaln
```