How to use **PAML** on *BisonNet*

<u>Module:</u> molecular_evolution

Version: 4.9h

<u>Description on BisonNet:</u> Testing genes for selection

What it really does:

PAML is a group of programs that allows the phylogenetic analysis of DNA or protein sequences through use of maximum likelihood. **PAML** is an acronym which stands for *phylogenetic analysis by maximum likelihood*, the description is in the name! It is mainly a collection of the following programs when instructed to do so: *baseml, basemlg, codeml, evolver, pamp, ynOO, mcmctree, and chi2*. These programs offer a wide range of analyses from comparing and testing phylogenetic trees to estimating synonymous and nonsynonymous substitution rates. The full list of **PAML** capabilities is listed below:

- Comparison and tests of phylogenetic trees (baseml and codeml);
- Estimation of parameters in sophisticated substitution models, including models of variable rates among sites and models for combined analysis of multiple genes or site partitions (baseml and codeml);
- Likelihood ratio tests of hypotheses through comparison of implemented models (baseml, codeml, chi2);
- Estimation of divergence times under global and local clock models (baseml and codeml);
- Likelihood (Empirical Bayes) reconstruction of ancestral sequences using nucleotide, amino
 - acid and codon models (baseml and codeml);
- Generation of datasets of nucleotide, codon, and amino acid sequence by Monte Carlo simulation (evolver);
- Estimation of synonymous and nonsynonymous substitution rates and detection of positive selection in protein-coding DNA sequences (yn00 and codeml).
- Bayesian estimation of species divergence times incorporating uncertainties in fossil calibrations (mcmctree).

How to Use **PAML** on BisonNet:

In preparation for using **PAML**, you need to make sure your sequence is in <u>one</u> of the following formats:

Sequence data file ("PHYLIP" format)-

```
To Note:
                                              The "4" represents the # of species.
         4 60
                                              The "60" represents the # of nucleotides.
      sequence 1
                                                      - divide by 3 to get # of codons
      AAGCTTCACCGGCGCAGTCATTCTCATAAT
      CGCCCACGGACTTACATCCTCATTACTATT
                                              The sequence name should be limited to 10 characters.
      sequence 2
      AAGCTTCACCGGCGCAATTATCCTCATAAT
                                              Name and sequence should be separated by at least two spaces
      CGCCCACGGACTTACATCCTCATTATTATT
      sequence 3
      AAGCTTCACCGGCGCAGTTGTTCTTATAAT
      TGCCCACGGACTTACATCATCATTATTATT
      sequence 4
      AAGCTTCACCGGCGCAACCACCCTCATGAT
      TGCCCATGGACTCACATCCTCCCTACTGTT
Tree file-
((1,2),3,4);
                OR
((human:.1,chimpanzee:.2):.05,gorilla:.3,orangutan:.5);
Control file (____.ctl)
```

```
To Note:
```

```
- Make sure code name is
  seqfile = seqfile.txt
                           * sequence data filename
 outfile = results_0.txt * main result file name
                                                              codeml.ctl
               * 0,1,2,3,9: how much rubbish on the screen
              * 1:detailed output
* -2:pairwise
 verbose = 1
 runmode = -2
NSsites = 0
   icode = 0
               * 0:universal code
fix_kappa = 0 * 1:kappa fixed, 0:kappa to be estimated
                * initial or fixed kappa
   kappa = 2
                                                              - change omega values for
fix_omega = 0
              * 1:omega fixed, 0:omega to be estimated
   omega = 0.09
                * 1st fixed omega value [change this]
                                                              desired dN/dS
  * EXCERCISE 1
  *alternate fixed omega values
  *omega = 0.005 * 2nd fixed value
  *omega = 0.01 * 3rd fixed value
  *omega = 1.60 * 9th fixed value
*omega = 2.00 * 10th fixed value
```

Once the data is formatted, this code can be followed in order to run PAML.

Load the module

Katie Wendell November 2020

```
module load molecular evolution
```

Move the desired .ctl file file into your home directory and rename it codeml.ctl.

```
cp -rp /data/courses/biol325_evolgen/PAML_activity/ex1_codeml.ctl .
mv ex1 codeml.ctl codeml.ctl
```

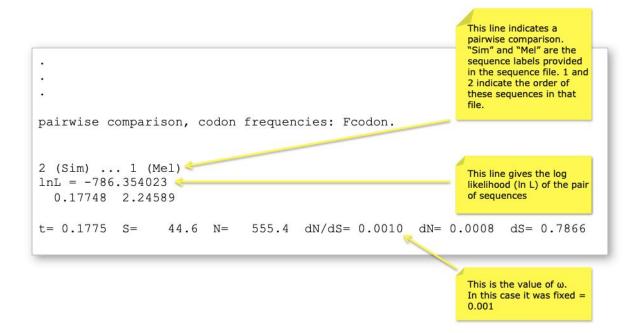
Create directories for each omega value results. Then move your desired sequence and the codeml.ctl file into the first directory.

```
mkdir results_for_omega_equals_0.001
mv codeml.ctl results_for_omega_equals_0.001
mv seqfile.txt results for omega equals 0.001
```

Run CODEML using the following command, changing the .ctl file for the corresponding omega values. Repeat for each value.

```
nano codeml.ctl
/software/apps/paml/current/bin/codeml codeml.ctl
```

Finally, observe and record the results utilizing this graphic.



Katie Wendell November 2020