Creating Phylogenetic Trees with BEAST with Sequences Obtained from biomaRt in R

Roman Meneghini

Last Updated December 16, 2023

Video Tutorial: https://youtu.be/wH6m6RlYbqk

Phylogenetic trees are akin to the table of contents to the book of life. They can describe how different organisms connect to each other and when they may have split off from sister taxa. That said, when looking at organisms in nature, discerning what is and is not the same species can be difficult whether they are multiple species that appear to be one species, or one species that appears to be multiple (Bickford et al., 2007; Thompson et al., 2022). Using DNA from said organisms can make this task much easier as we can see on the genetic level how similar two individuals may be. Beyond that, the evolutionary relationships between them can be understood. Using tools like Bayesian Evolutionary Analysis Sampling Trees (BEAST) you can determine not only which organisms shared the most recent common ancestor but can also get an estimate on how long ago they shared that common ancestor.

This manual will demonstrate how to download BEAST and associated software, how to download and extract multiple sequences in R to export to a single FASTA file, how to align sequences using resources provided by EMBL-EBI, and finally how to use BEAST to produce phylogenetic trees using aligned sequence data and molecular clock calculations. This manual is aimed towards Windows users, so if using a different operating system, you may need to do some OS specific steps outside the manual. To demonstrate the effectiveness of BEAST, the tutorial will by default have you collect sequences from the reference genomes of several mammals. The tree produced by BEAST will then be compared to peer reviewed studies documenting the phylogenetic relationships of these animals to display the effectiveness of BEAST. Beyond BEAST using numerous iterations of possible evolutionary relationships to predict the most likely relationships, it is free software that has ongoing development support and troubleshooting available.
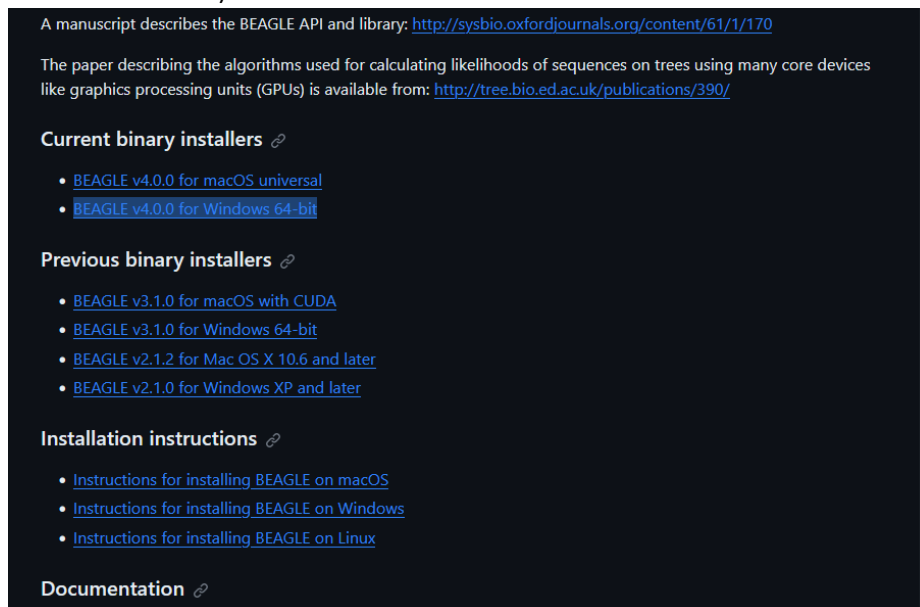
Video Link: https://youtu.be/wH6m6RlYbqk

Tools and Resources Used

- BEAST https://beast.community/#
  - Dependencies
    - BEAGLE https://beast.community/beagle
    - Fig Tree http://tree.bio.ed.ac.uk/software/figtree/ https://github.com/rambaut/figtree/releases
      - Technically a soft dependency as it is only required to view tree
    - Java Development Kit and Java Runtime Environment JDK and JRE
      - https://www.oracle.com/java/technologies/downloads/ JDK
      - https://www.java.com/en/ JRE
- EMBL-EBI multiple sequence alignment services: https://www.ebi.ac.uk/jdispatcher/msa
  - Requires single FASTA file

- RStudio https://www.r-project.org/
  - Rscript used https://github.com/Bearypterid/BEAST-Manual
- Basic Text Editor

**Setting up**

1. Open the link to BEAST provided above
   a. Click on Install BEAST on Windows
   b. Download and extract BEAST in a directory you are comfortable working from
2. Open the links for Java provided above
   a. Download and install the latest Java Development Kit (JDK) using the Windows installer
   b. Download and install the Java Runtime Environment (JRE)
3. Open the link to BEAGLE
   a. Click the github link under downloading and installing BEAGLE
   b. Scroll all the way down and download the Binary installer for BEAGLE for Windows (see screenshot below).



   c. Use the installer
4. If you do not have RStudio installed, please download R and Rstudio using the link above
   a. When RStudio is installed, download the RScript linked above and place it in a new folder in your home directory. Make note of what the folder is called.
5. Download and extract FigTree https://github.com/rambaut/figtree/releases

**Walkthrough**

**Getting organism specific DNA sequences in RStudio**

1. Open up the Rscript and run it.  It will prompt you as needed with messages

```
 1  message("Create a directory in your home directory for this script to save files to, then change the working directory in the next line")
 2
 3  setwd("~/RProject")#instead of /RProject it should be whatever your folder of interest is called
 4  packages=c('biomaRt','BiocManager','plyr','dplyr','tibble','tidyselect')
 5  install.packages(setdiff(packages, rownames(installed.packages())))
 6  install.packages("devtools")
 7  # devtools::install_version("dplyr", version = "2.3.4") #This line is needed if the wrong dplyr version is installed and is interferring with biomart
 8  library("biomaRt")
 9  library("BiocManager")
10  library("plyr")
11  library("dplyr")
12  library("tibble")
13  library("tidyselect")
14  message("If biomart did not instal than uncomment and run the next two lines of code")
15  # BiocManager::install("biomaRt",force = TRUE)
16  # library("biomaRt")
```

2. These opening lines will install and load the packages needed.
3. The next set of lines allows you to browse the niche databases hosted by Ensembl as well as the more common and comprehensive ones.
    a. The default *datasetens* variable contains the database for vertebrates as well as *D. melanogaster* and *C. elegans*, while *datasetsanim* contains the database for all the remaining animals.

```
17  biomaRt::listMarts()# These 3 lines are available for you to use and browse the other more niche databases
18  listEnsembl()
19  listEnsemblGenomes()
20
21  ens=useMart("ensembl")
22  anim=useEnsemblGenomes("metazoa_mart")
23  singleeuk=useEnsemblGenomes("protists_mart")
24  fungi=useEnsemblGenomes("fungi_mart")
25  plants=useEnsemblGenomes("plants_mart")
26
27  datasetsens=listDatasets(ens)
28  datasetsanim=listDatasets(anim)
29  datasetsprotist=listDatasets(singleeuk)
30  datasetsfungi=listDatasets(fungi)
31  datasetsplants=listDatasets(plants)
32
```

4. Here you can select the database you would like to use. By default it will use the main Ensembl vertebrate database. Feel free to browse this database by clicking on it in your environment.

```
33  datasetused=datasetsens
34  datsetused=rownames(datasetused)=datasetused[,2]
35  message("Browse through yourdataset used to find your organisms of interest. Once you have found your ",
36          "organisms of interest, copy and past the text from the information cell into the list function below. ",
37          "The list function below will be used to search for and collect data on your species of interest, so make sure ",
38          "the names in the list are exactly as written in the database you browsed, otehrwise there will be an error. ",
39          "The first item in the list needs to be the database used, otherwise the script will not run.")
40
41  specieslist=list(ens,"Opossum","Wallaby","Platypus","Hyrax","Elephant genes","Koala","Tasmanian","Common wombat")
42
```

5. After browsing through the database, we will want to select our species of interest and make a list with them. The list needs to use the exact wording from the description column of the database for the keywords in the list to work correctly, however if you choose not to touch the list, the script will continue with the default animals.
    a. Notice that "common" is found in the names of multiple animals in the database, so we need "Common wombat". By contrast, Tasmanian is enough to get the Tasmanian Devil as that animal is the only one in the database at present with that adjective.
    b. Also note, the first word is always uppercase in the database, make sure your list uses the correct case as the code is case sensitive.

| | dataset | description | version |
|---|---|---|---|
| 25 | cccarpio_gene_ensembl | Common carp genes (Cypcar_WagV4.0) | Cypcar_WagV4.0 |
| 151 | pmuralis_gene_ensembl | Common wall lizard genes (PodMur_1.0) | PodMur_1.0 |
| 171 | scanaria_gene_ensembl | Common canary genes (SCA1) | SCA1 |
| 211 | vursinus_gene_ensembl | Common wombat genes (bare-nosed_wombat_geno... | bare-nosed_wombat_genome_assembly |

6. Run the next set of code, it requires no input. It simply saves the information needed to use the specific datasets for each organism into a variable for each.

```
42
43   numberspecieslist=length(specieslist)
44 ▾ for (i in 1:numberspecieslist) {#This for loop creates the marts needed for the get bm function
45     maht=paste("mart",i,sep = "")
46     # speciesgen=paste("speciesgenes",i,sep = "")
47 ▾   if (i==1) {
48       assign(maht,specieslist[[i]])
49 ▾   } else{
50       searchspecies=specieslist[[i]]
51       speciesdata=datasetused[c(paste(searchspecies)),]
52       speciesdata=speciesdata[1,1]
53       speciesgenes=useDataset(paste(speciesdata,sep = ""),mart=mart1)
54       assign(maht,speciesgenes)
55 ▴   }
56 ▴ }
```

7. The next set of code saves the whole dataset for our first organism into a dataframe you can view by clicking on it in the environment. It will look similar to the dataframe below. The first column will contain the ensembl gene id, the second will contain the go enrichment term attached to the gene, and the last column will provide the gene name. The gene name is a term that will be shared between species. Feel free to browse around looking at what genes are named, and what are not.

```
57   message("Use speciesgenes0 to browse for possible genes of interest")
58
59   speciesgenes0=getBM(attributes=c("ensembl_gene_id","name_1006","external_gene_name"),mart=mart2)#Use to find Genes
```

| | ensembl_gene_id | name_1006 | external_gene_name |
|---|---|---|---|
| 1 | ENSUAMG00000000001 | membrane | SLC25A42 |
| 2 | ENSUAMG00000000001 | transmembrane transport | SLC25A42 |
| 3 | ENSUAMG00000000001 | mitochondrial inner membrane | SLC25A42 |
| 4 | ENSUAMG00000000001 | ATP transmembrane transporter activity | SLC25A42 |
| 5 | ENSUAMG00000000001 | ADP transmembrane transporter activity | SLC25A42 |
| 6 | ENSUAMG00000000001 | coenzyme A transmembrane transporter activity | SLC25A42 |
| 7 | ENSUAMG00000000001 | ADP phosphatase activity | SLC25A42 |
| 8 | ENSUAMG00000000001 | AMP transmembrane transporter activity | SLC25A42 |
| 9 | ENSUAMG00000000001 | ADP transport | SLC25A42 |
| 10 | ENSUAMG00000000001 | ATP transport | SLC25A42 |
| 11 | ENSUAMG00000000001 | coenzyme A transmembrane transport | SLC25A42 |
| 12 | ENSUAMG00000000001 | AMP transport | SLC25A42 |
| 13 | ENSUAMG00000000001 | mitochondrion | SLC25A42 |
| 14 | ENSUAMG00000000002 | membrane | TMEM161A |
| 15 | ENSUAMG00000000002 | response to retinoic acid | TMEM161A |
| 16 | ENSUAMG00000000002 | cellular response to oxidative stress | TMEM161A |
| 17 | ENSUAMG00000000002 | cellular response to UV | TMEM161A |
| 18 | ENSUAMG00000000002 | positive regulation of DNA repair | TMEM161A |
| 19 | ENSUAMG00000000002 | negative regulation of intrinsic apoptotic signaling p... | TMEM161A |
| 20 | ENSUAMG00000000003 | RNA polymerase II transcription regulatory region s... | MEF2B |

8. The next two lines relate to our gene of interest. The variable should equal a gene name found in the above dataframe. The default gene used is PAOX as it was documented in all the organisms of interest, but you can and should try other genes as the script will check if all reference genomes used have it.

```
61    message("This is the gene of interest the rest of the code will use.")
62    filtervalue="PAOX" #Gene of interest
```

9. The next for loop checks for if our gene of interest is in each of our organisms of interest and then saves the database for that organism if it is. If the gene is not present in any of the organisms, the loop will end and you will be given a message prompting you to try a different gene.
   a. Note, this for loop will take a bit to run as scans through 100,000 of thousands of rows for each organism to determine whether or not the gene of interest is there.

```
64  for (i in 2:numberspecieslist) {
65    speciesgen=paste("speciesgenes",i,sep = "")
66    speciesgenepresent=getBM(attributes=c("ensembl_gene_id","name_1006","external_gene_name"),mart=get(paste("mart",i,sep = "")),filters = "external_gene_n
67    assign(speciesgen,speciesgenepresent)
68    if (is.na(speciesgenepresent[1,1])) {
69      message("The gene you have selected is not present in all species of interest.",
70              " Please select a different gene and try again")
71      {break}
72    }
73  }
```

10. If you were not given a message prompting you to use a different gene of interest, simply proceed with the code. The next for loop extracts the DNA sequences for the gene of interest for each organisms and binds them into one FASTA file.

```
74  message("If your gene of interest was found in all species of interest, proceed.")
75  for (i in 2:numberspecieslist) {
76    searchspecies=specieslist[[i]]
77    speciesdata=datasetused[c(paste(searchspecies)),]
78    geneseq=getSequence(id = filtervalue,type ="external_gene_name",seqType = "cdna",mart = get(paste("mart",i,sep = "")))
79    geneid=get(paste("speciesgenes",i,sep = ""))
80    geneid=geneid[1,1]
81    Organism=speciesdata[1,2]
82    geneseq[geneseq == filtervalue]=paste(geneid,Organism,filtervalue,"cDNA",sep = ".")
83    if (i==2) {
84      Fastaseq=geneseq
85    }else{
86    Fastaseq=rbind.data.frame(Fastaseq,geneseq)
87    }
88  }
```

11. The last set of code will remind you to delete any files with the existing name of our FASTA file before exporting the file created by the script. The single file produced is needed for the alignment tools used further on.

```
89  message("Make sure to delete any duplicate fasta files before the next line is run.",
90          " Otherwise it will simply add your new sequences to the old one including any duplicates")
91  exportFASTA(Fastaseq,file = paste(filtervalue,"fasta",sep = "."))
```

12. This script requires minimal input. The majority of work required on your end is to set up the packages, choose the organisms and gene of interest, and then run it. It is designed to reduce the time and effort needed on your end when changing the genes and organisms used.

**Alignment**

1. Navigate to EMBL-EBI multiple sequence alignment services using the link https://www.ebi.ac.uk/jdispatcher/msa
2. There are many tools available, but for this walkthrough, we will use MAFFT
   a. The tools vary in speed and the size of regions they can align. MAFFT is just suitable for the default data, if you are using more organisms or larger segments of DNA, consider one of the other MSA tools available on the site.

   MAFFT
   MSA tool that uses Fast Fourier Transforms. Suitable for medium-large alignments.
   ↘Launch MAFFT

3. Change the setting from automatic to DNA, upload the FASTA file produced above and submit the job. It will only take a minute or two at most for this example, but if you are using hundreds of sequences, it will take longer.
4. After the job finishes, click tool output. It will produce an aligned FASTA file that looks like this.
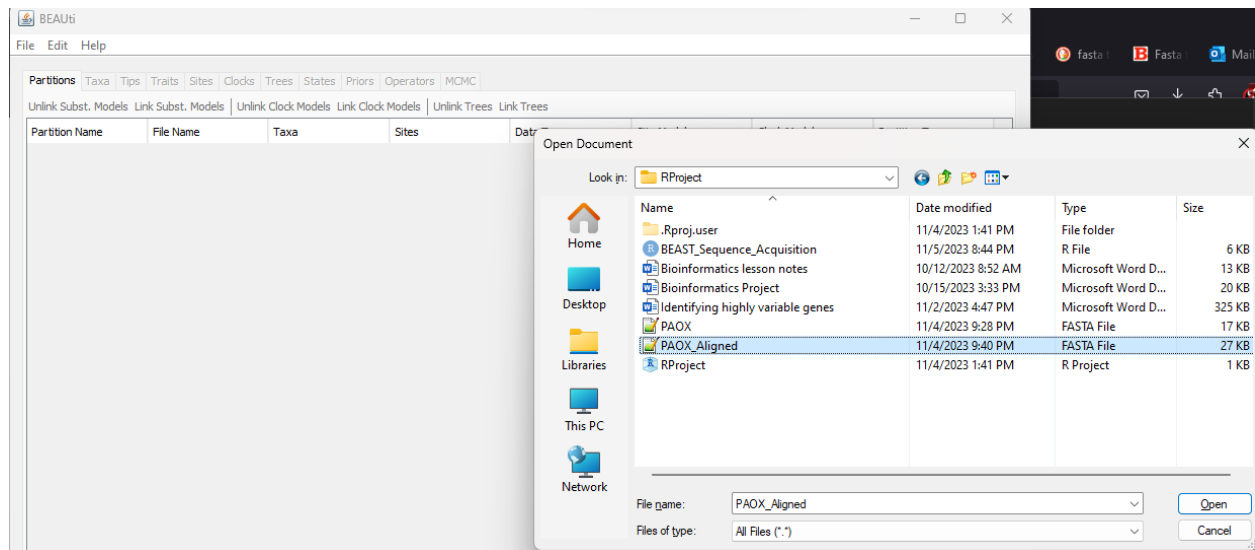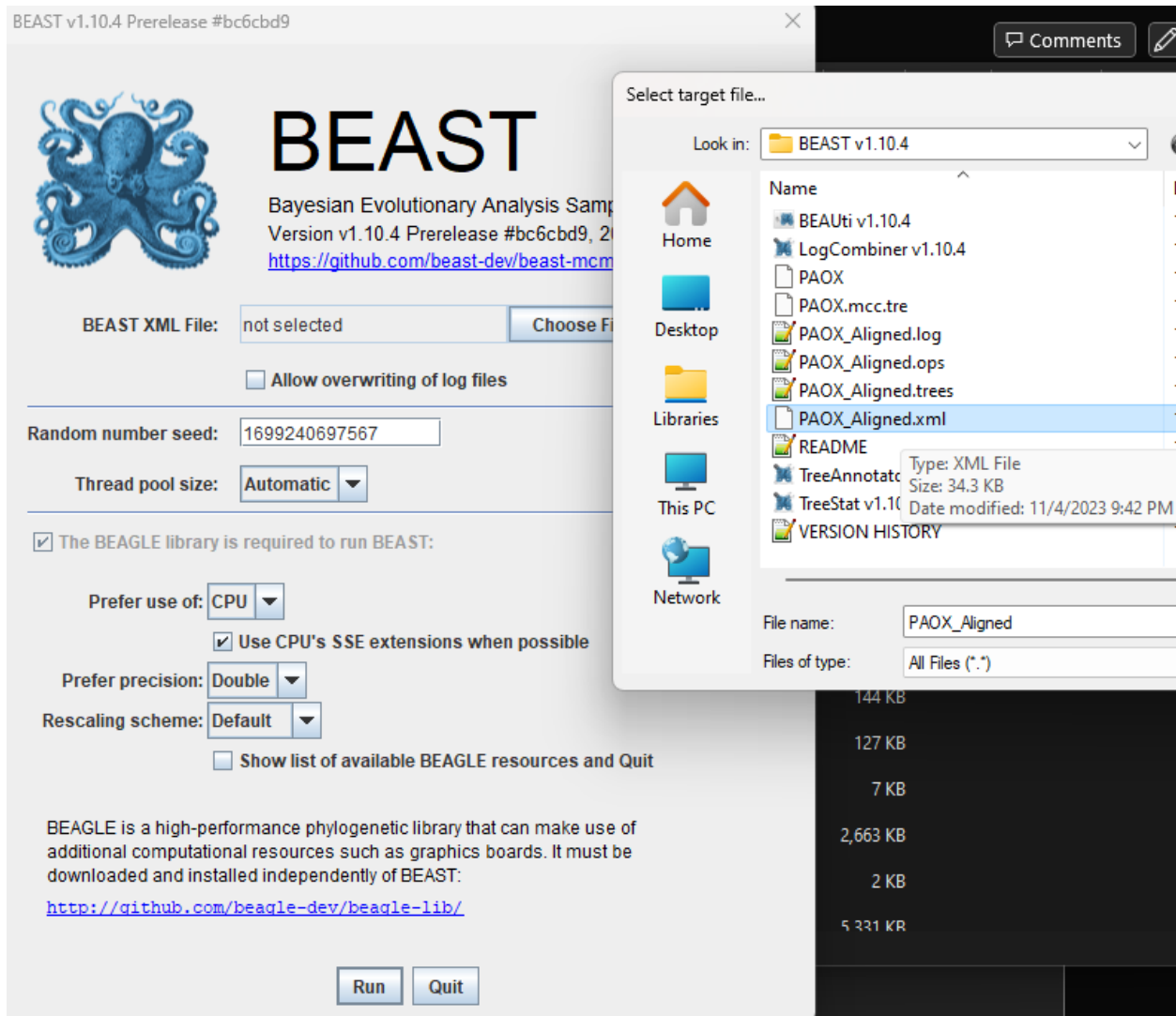
5. It can also produce a simple phylogenetic tree, but this is typically not suitable for any serious analysis of relationships.
6. In the tool output tab, click download to download the aligned file. Open it in a program like notepad or notepad ++ and save it to either the work directory folder used by R or the folder BEAST is installed in. Save the file as your "Gene of interest_Aligned.FASTA". This will save our alignment as a new aligned FASTA file to use with BEAST.
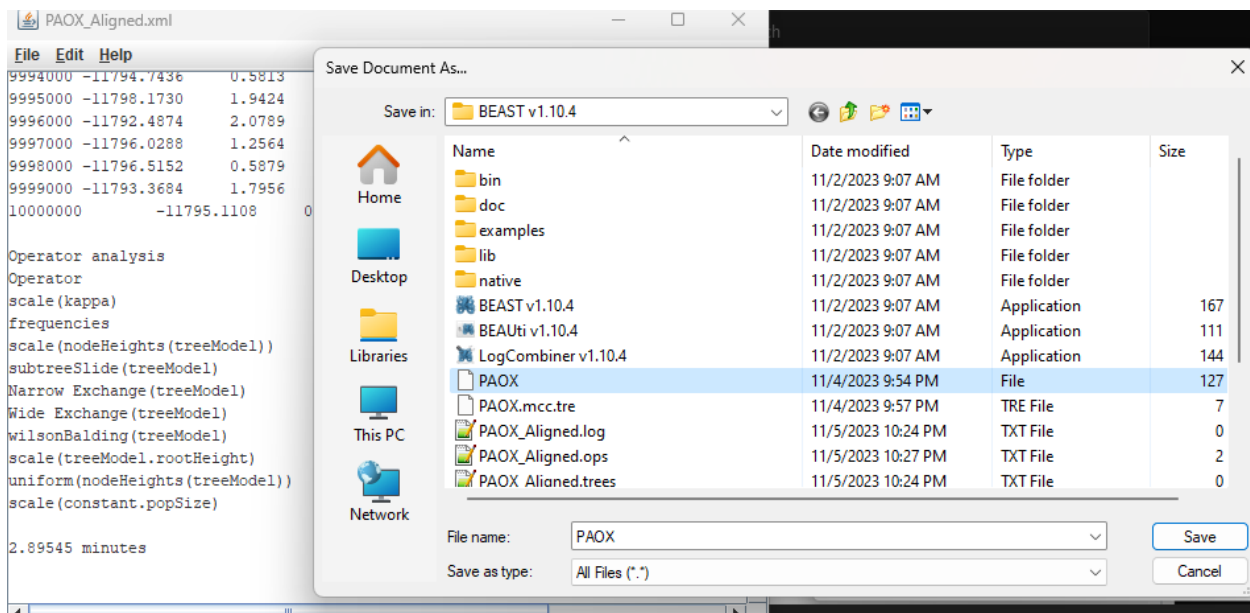


**BEAST**

1. Navigate to the file where BEAST was extracted and launch BEAUTi
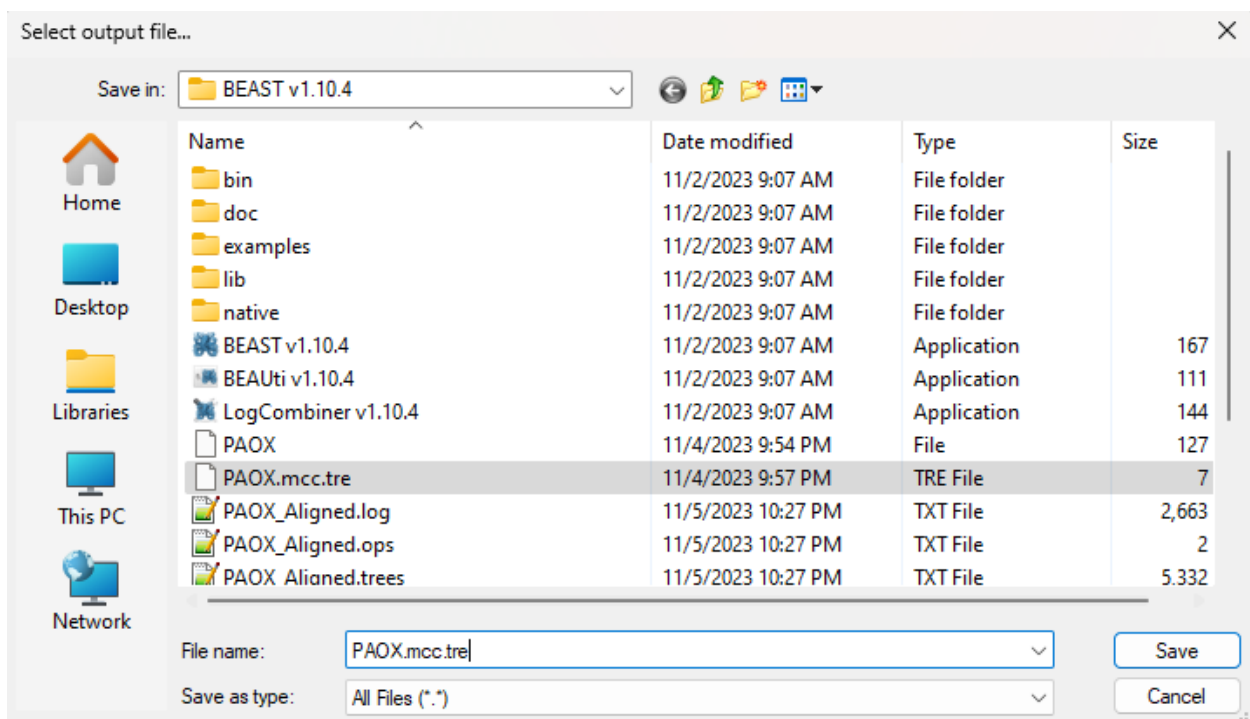   a. BEAST requires a specific xml file to run, and you obtain that xml file using BEAUTi.

2. Click file, then import data then select the aligned FASTA file made above and then click "generate BEAST file" once it has loaded the file. If you would like to alter the evolutionary model or the clock model used you may do so now, but we will stick with the default for now. Once you have generated the BEAST file, you may close BEAUTi

3. Launch BEAST and choose the newly generated XML file from part 2. Toggle "Allow overwriting of log files", the default is off and can produce an error. Click run. This will take some time and can be intensive on your computer. As it performs many repeated calculations over and over to produce numerous sets of possible trees.

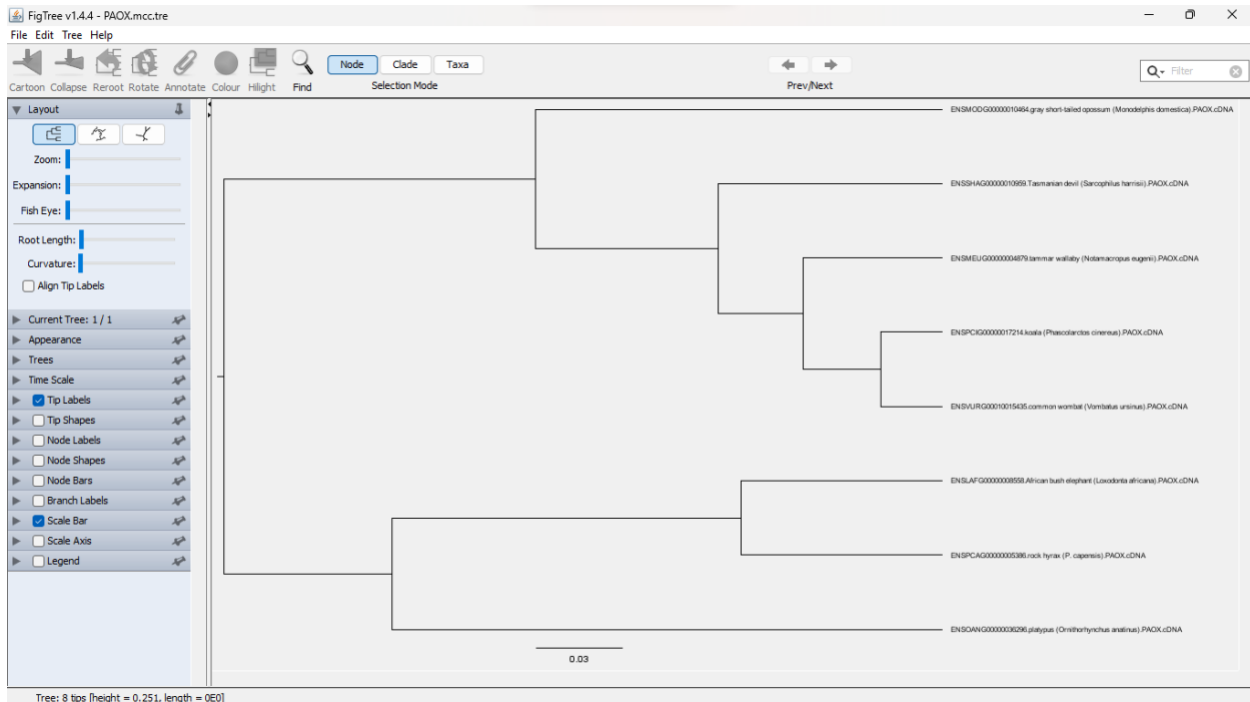4. Save the File. I used the gene of interest as the name for this example

5. Close BEAST and open TreeAnnotator. You will find it in the same folder as BEAUTi and BEAST. Set the "Burnin (as states) to 100000 as this was the default used for BEAST. Load the file produced by BEAST as the input tree file. For the output file, click choose file, then type in the new window the name you want for the file, with the file extension of .mcc.tre as seen in the example below
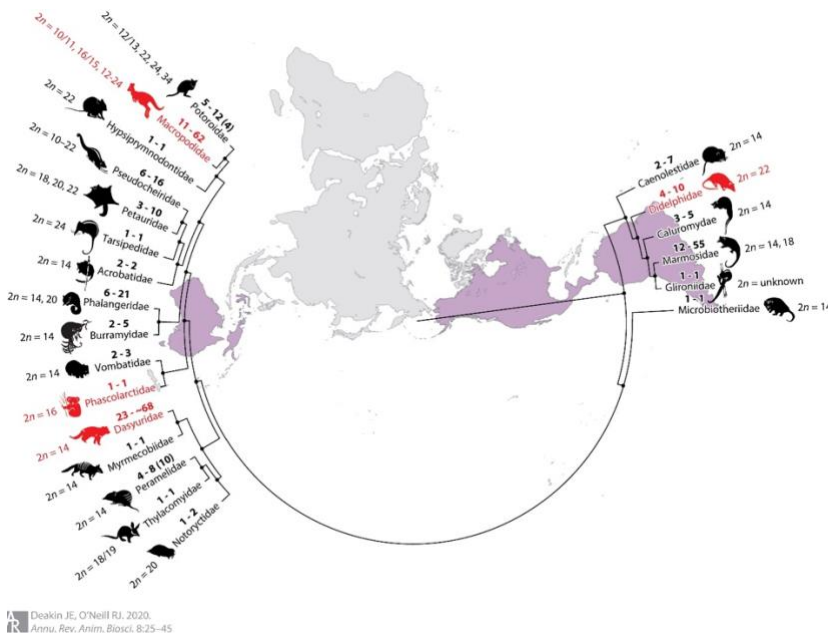


6. Click run. The program will now generate the data for trees based off of the median outcome of possibilities produced by BEAST.

7. To actually visualize the tree, you need to navigate to where FigTree was installed and open FigTree. Once FigTree is open, click file, then open the newly generated file from TreeAnnotator to see the newly created tree.



8. When compared to the phylogenetic tree found in Deakin & O'Neil (2020), we can see our small example phylogeny lines up fairly well to theirs when looking at the marsupials. Opossums are placed on their own branch, then Tasmanian Devils, then Wallabies, and finally Wombats and Koalas make up their own clade.



Figure 1 (Deakin & O'Neill, 2020)

**References**

Bickford, D., Lohman, D. J., Sodhi, N. S., Ng, P. K. L., Meier, R., Winker, K., Ingram, K. K., & Das, I. (2007).

    Cryptic species as a window on diversity and conservation. *Trends in Ecology & Evolution*, *22*(3),

    148–155. https://doi.org/10.1016/j.tree.2006.11.004

Deakin, J. E., & O'Neill, R. J. (2020). Evolution of Marsupial Genomes. *Annual Review of Animal*

    *Biosciences*, *8*(1), 25–45. https://doi.org/10.1146/annurev-animal-021419-083555

Thompson, M. J., Capilla-Lasheras, P., Dominoni, D. M., Réale, D., & Charmantier, A. (2022). Phenotypic

    variation in urban environments: Mechanisms and implications. *Trends in Ecology & Evolution*,

    *37*(2), 171–182. https://doi.org/10.1016/j.tree.2021.09.009