# Calibrated Species Trees
## Estimating species divergence times using StarBEAST2

*Huw A. Ogilvie*

## 1  Background

This tutorial assumes you have completed the first StarBEAST2 tutorial, "Species Trees with Relaxed Molecular Clocks." In that tutorial, the data suggested that zero shifts in substitution rates have occured within the pocket gopher genus *Thomomys*. No priors on node heights were used in the first tutorial, and the expected clock rate of each gene was 1.0, so the estimated speciation times were all in expected number of substituions.

We will improve on our previous analysis by calibrating the tree, so that we can estimate how many millions of years ago (ma) different *Thomomys* lineages diverged.

## 2  BEAUti

Run BEAUti by double clicking on its icon, or by launching the BEAUTi executable file from the command line in Linux.

### Beginning a strict clock analysis

Building on our previous results, we will use a strict clock for this tutorial. Compared to a relaxed clock, this will take less computational time and produce narrower credible ranges of divergence times. To set up a StarBEAST2 analysis using a strict clock, choose the **File/Template/StarBEAST2** menu item (Figure 1).
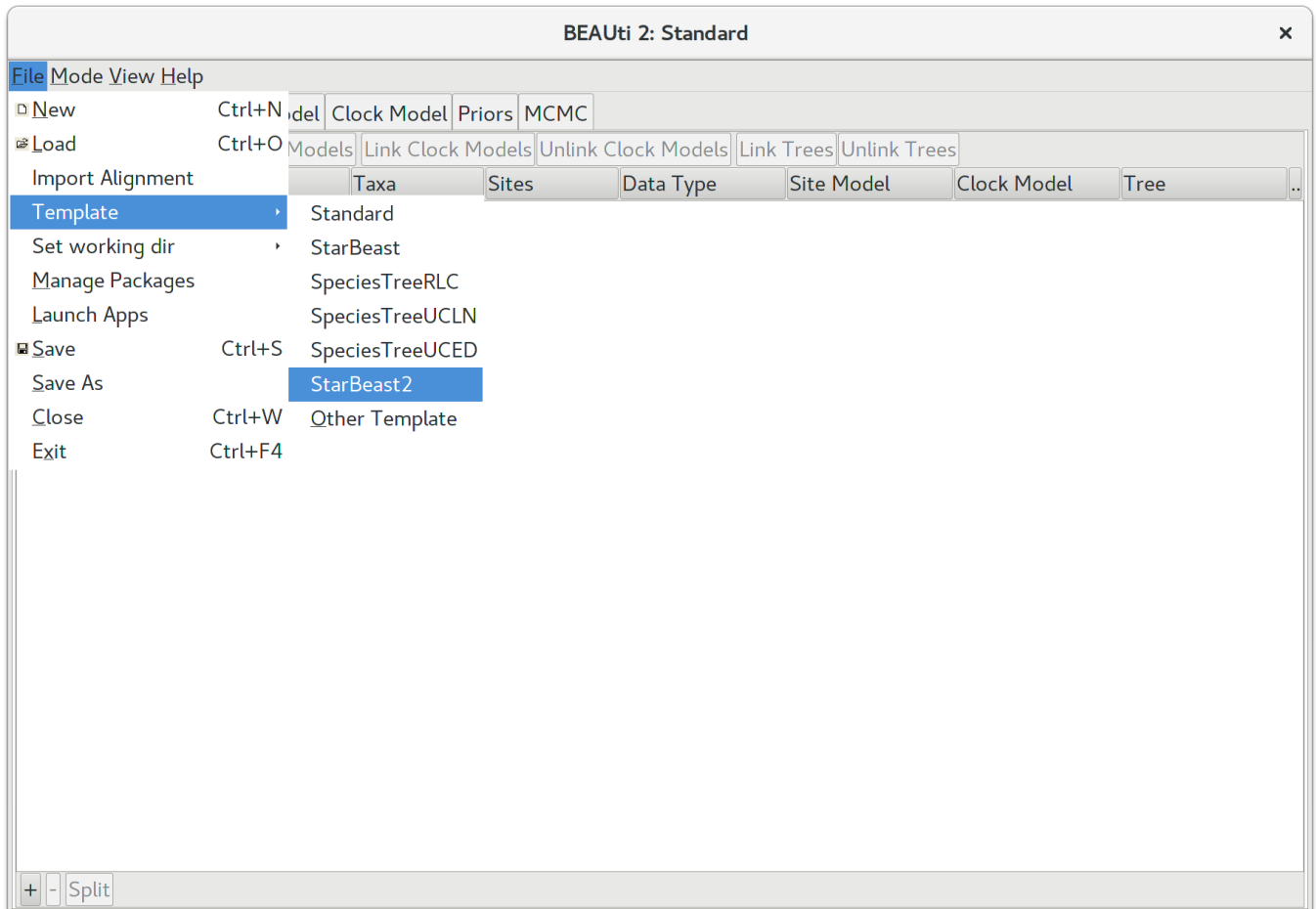
Figure 1: Select a species tree template in BEAUti.

## Allow clock rates to vary

By default BEAUTi fixes the clock rate of the first partition to 1, so that the rates of other loci are estimated relative to the first locus. This is generally inappropriate for StarBEAST2 analyses, so it is very important to disable this behaviour by deselecting the **Mode/Automatic set clock rate** menu item (Figure 2).
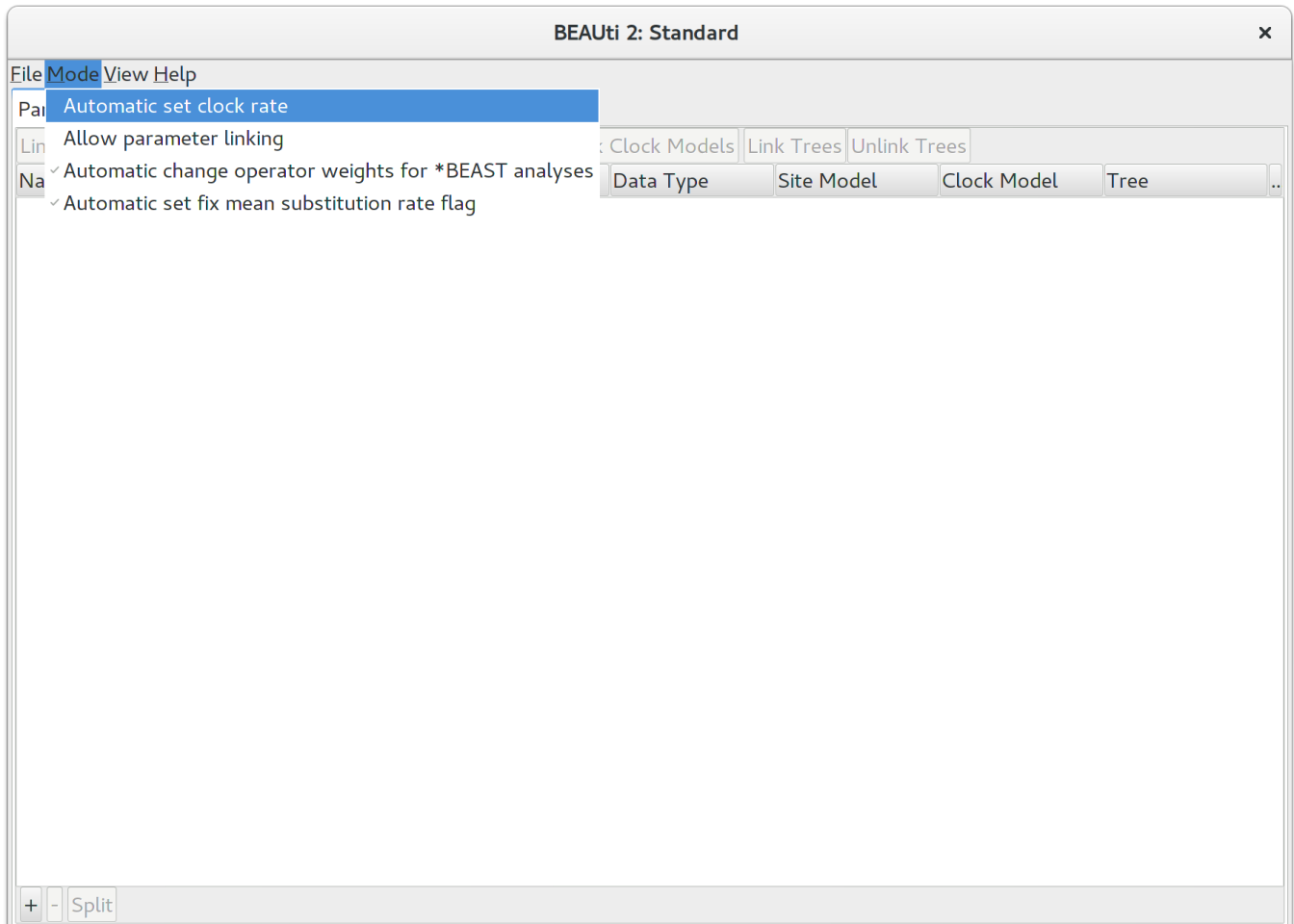


Figure 2: Disable automatic setting of clock rates.

## Loading the NEXUS files

To load a NEXUS format alignment, click the button with the plus symbol (+) in the lower left corner of the main **Partitions** tab. We will be using the same sequence files as in the previous tutorial; navigate to the **examples/nexus** subfolder inside the **beast** application folder, and select all of the first seven NEXUS files. They should be numbered 26, 29, 47, 53, 59, 64, and 72 (Figure 3).
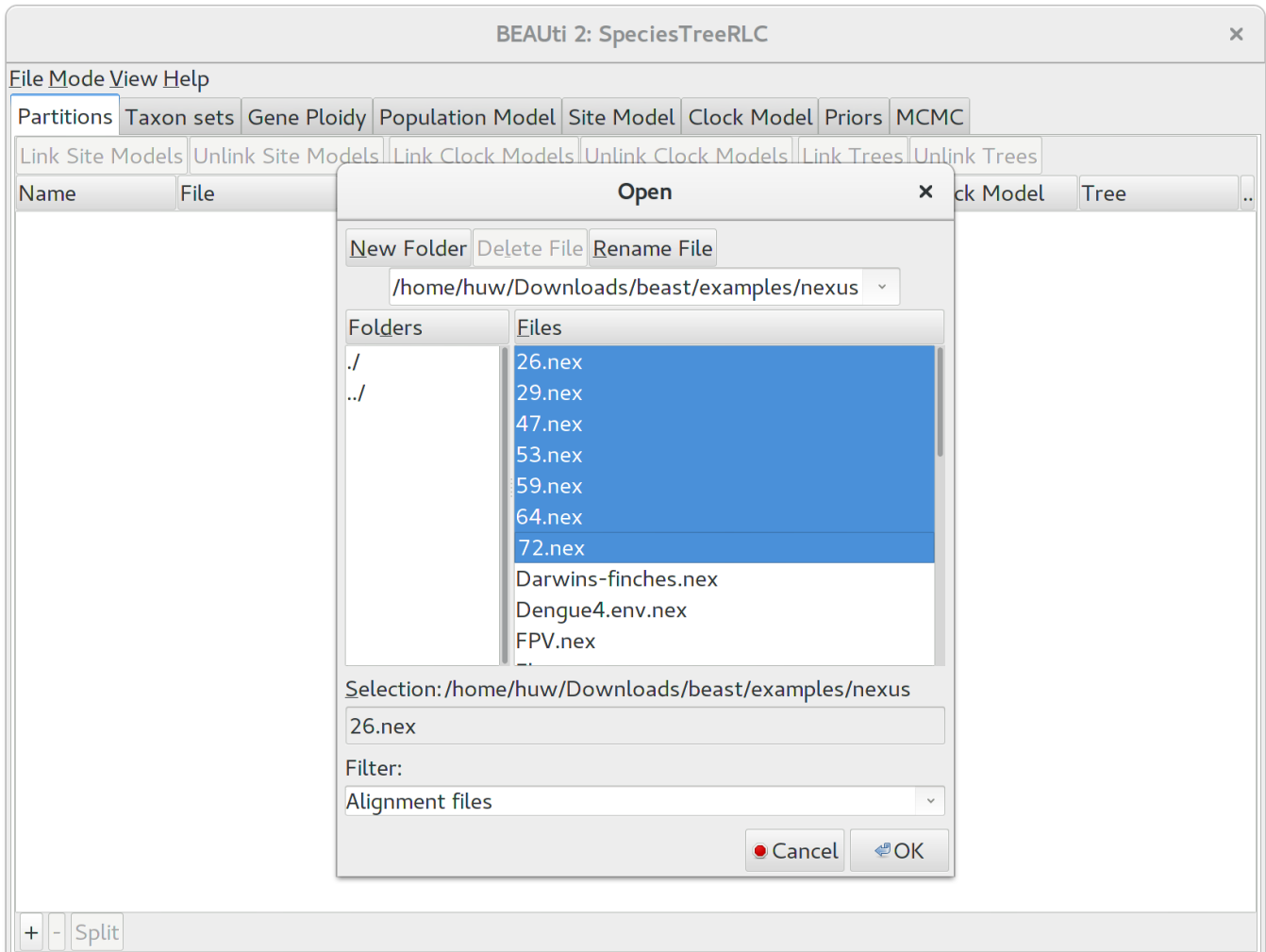


Figure 3: Selecting NEXUS alignment files to import.

## Assigning the correct species to each sequence

As in the previous tutorial, splitting the name on the underscore character '_' and selecting the second group will give us the mapping that we need (Figure 4).
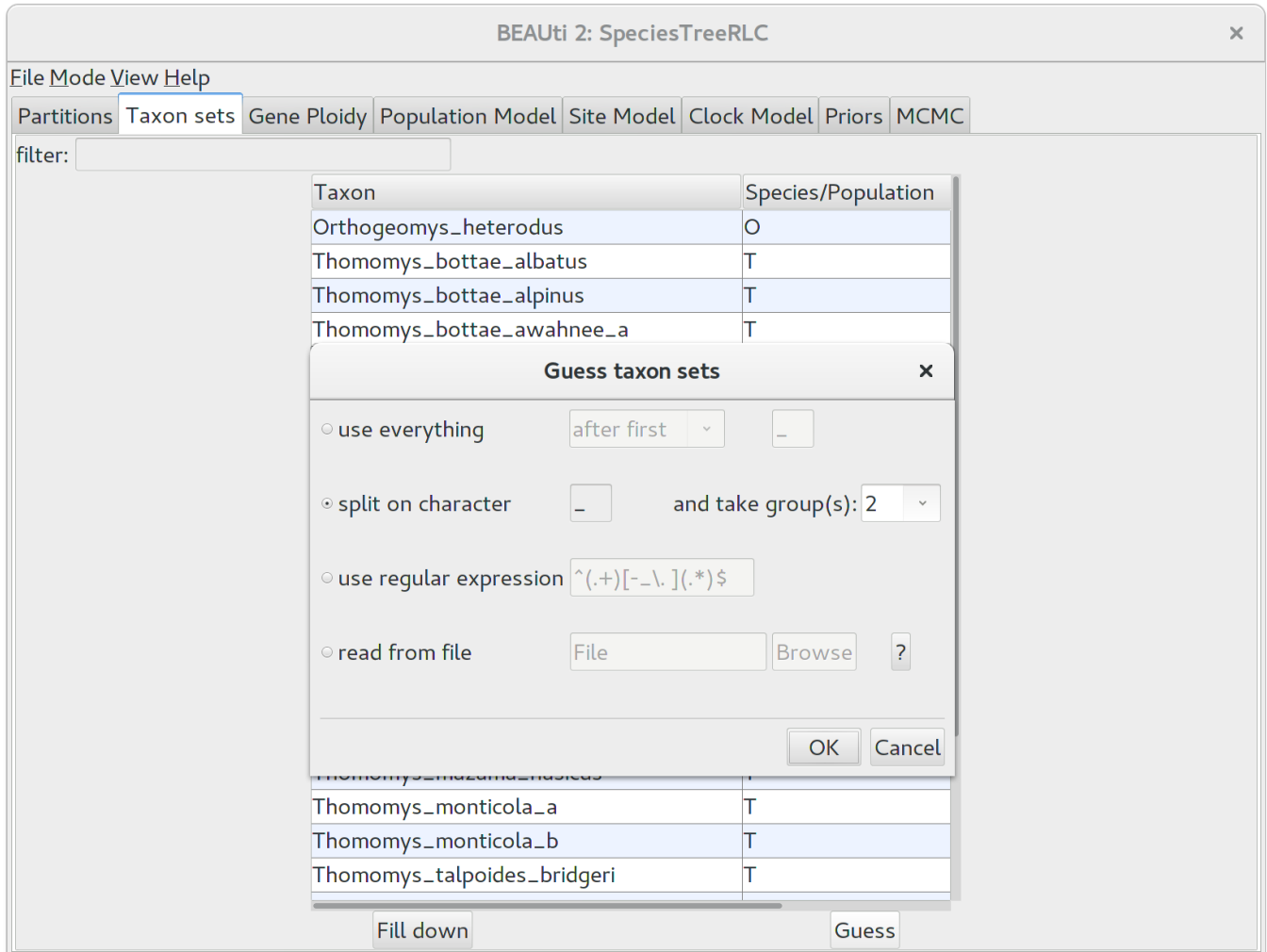


Figure 4: Assigning species to sequences in BEAUti using the guess dialog.

## Ploidy and population models

As in the previous tutorial, leave the ploidy of each gene and the population model at their default values.

## Setting the substitution model

For this tutorial we will also use the HKY substitution model, but will use a shortcut to set it for all genes at once. For the first gene (which should be "26") select "HKY" for substitution model (**Subst Model** in Figure 5) Now select the other six genes using the shift key. The right panel will now allow you to set the same substitution model for those six genes by cloning the parameters from the first gene. Click **OK** to proceed (Figure 6).
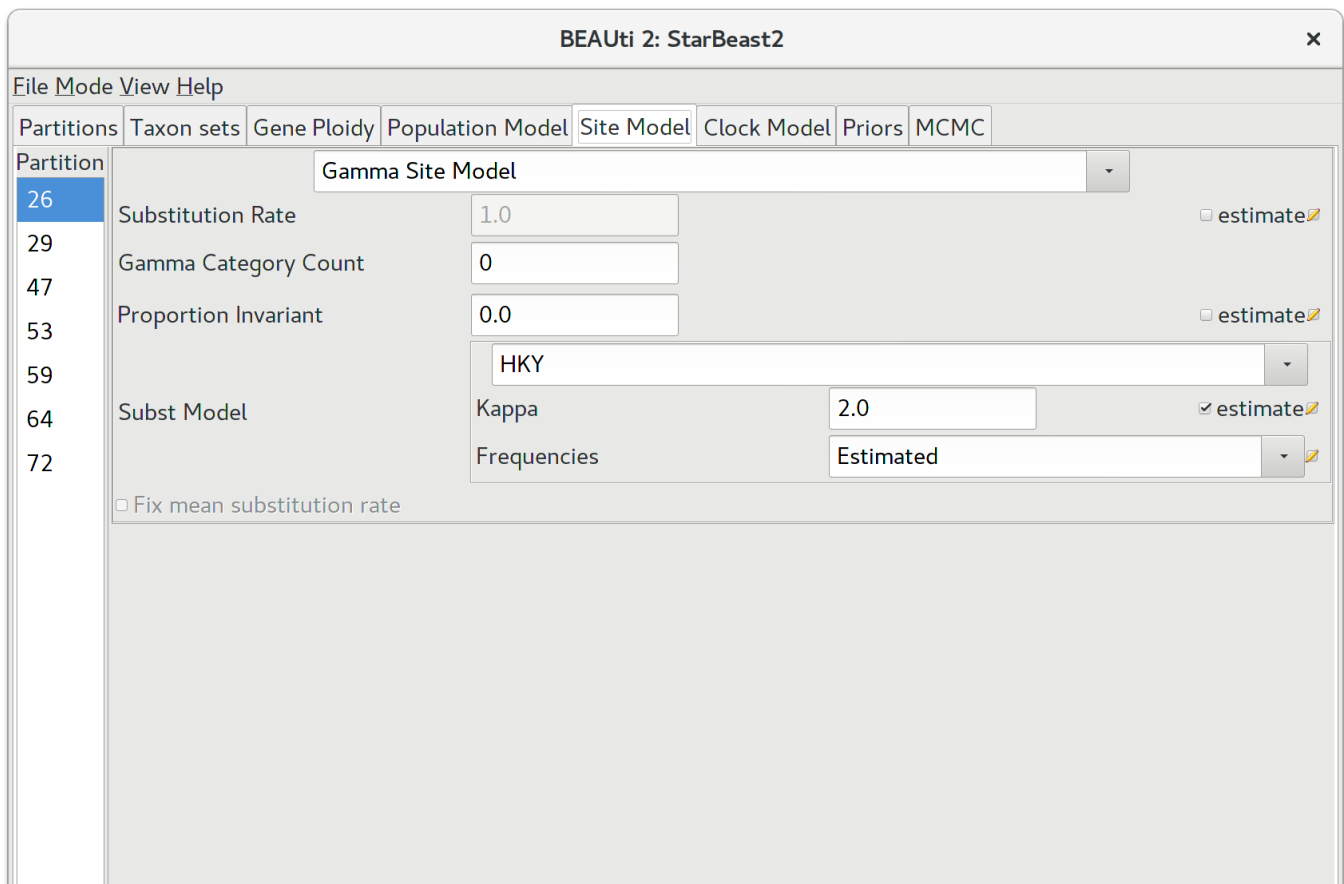


Figure 5: Setting up substitution and site models for the gopher alignments.
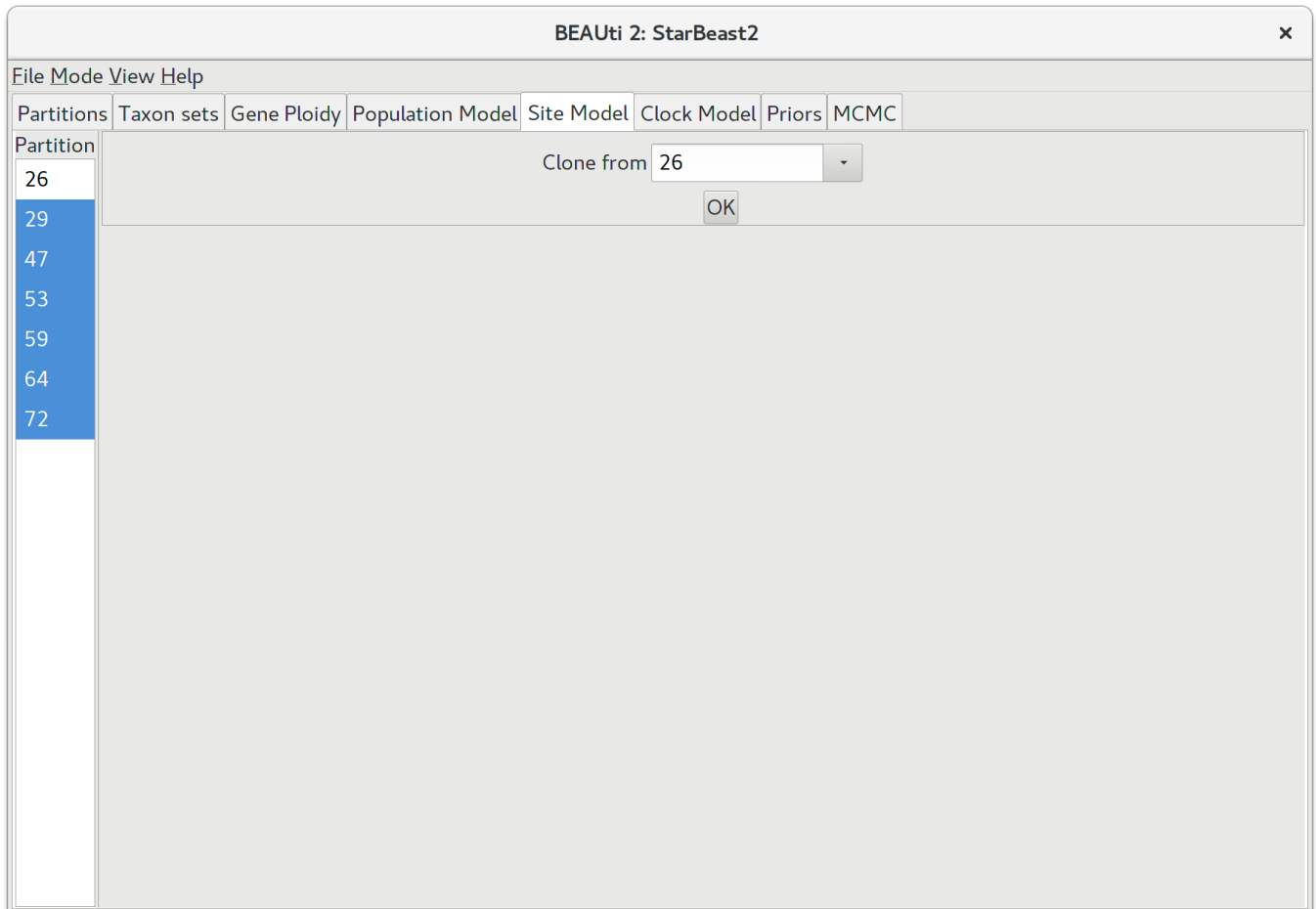
Figure 6: Setting multiple substitution and site models at once.

### Setting the clock model

Click on the **Clock Model** tab at the top of the main window. In this panel you can configure the mean clock rate for each locus. If you followed the earlier instructions to disable automatic setting of clock rates, the mean clock rate "Clock.rate" of all partitions should have the **estimate** box ticked (Figure 7).
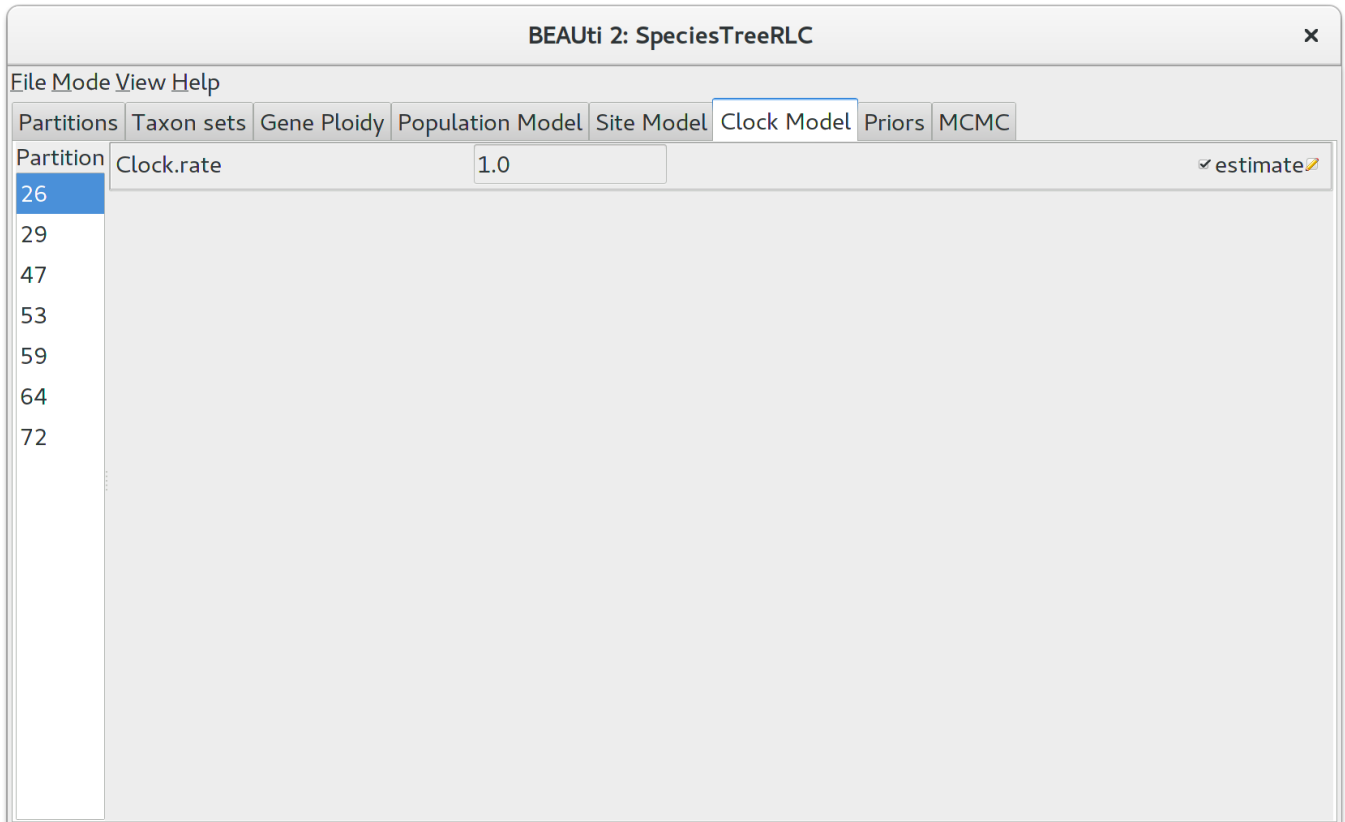


Figure 7: The default when automatic clock rate setting is disabled.

## Priors

Before applying the date calibration first change the prior on the species tree, "Tree.t:Species", to **calibrated Yule** (Heled and Drummond 2012). Click the top-leftmost arrow to expand the options available for the calibrated Yule model, leave the Birth Rate at 1.0, and uncheck the **estimate** box to make this a fixed parameter. For real analyses you should almost certainly estimate this value, but a fixed value will help us complete the tutorial in a reasonable time frame.

The default prior for mean clock rates in StarBEAST2 is a lognormal distribution with a mean (in real space) of 1. This will not be appropriate when using fossil (or other external) calibrations. For this tutorial we will use the 1/X prior, which is uninformative and will allow the calibration(s) to guide the clock rates. For each clockRate.c:*gene* parameter, set the prior distribution to be 1/X (Figure 8)
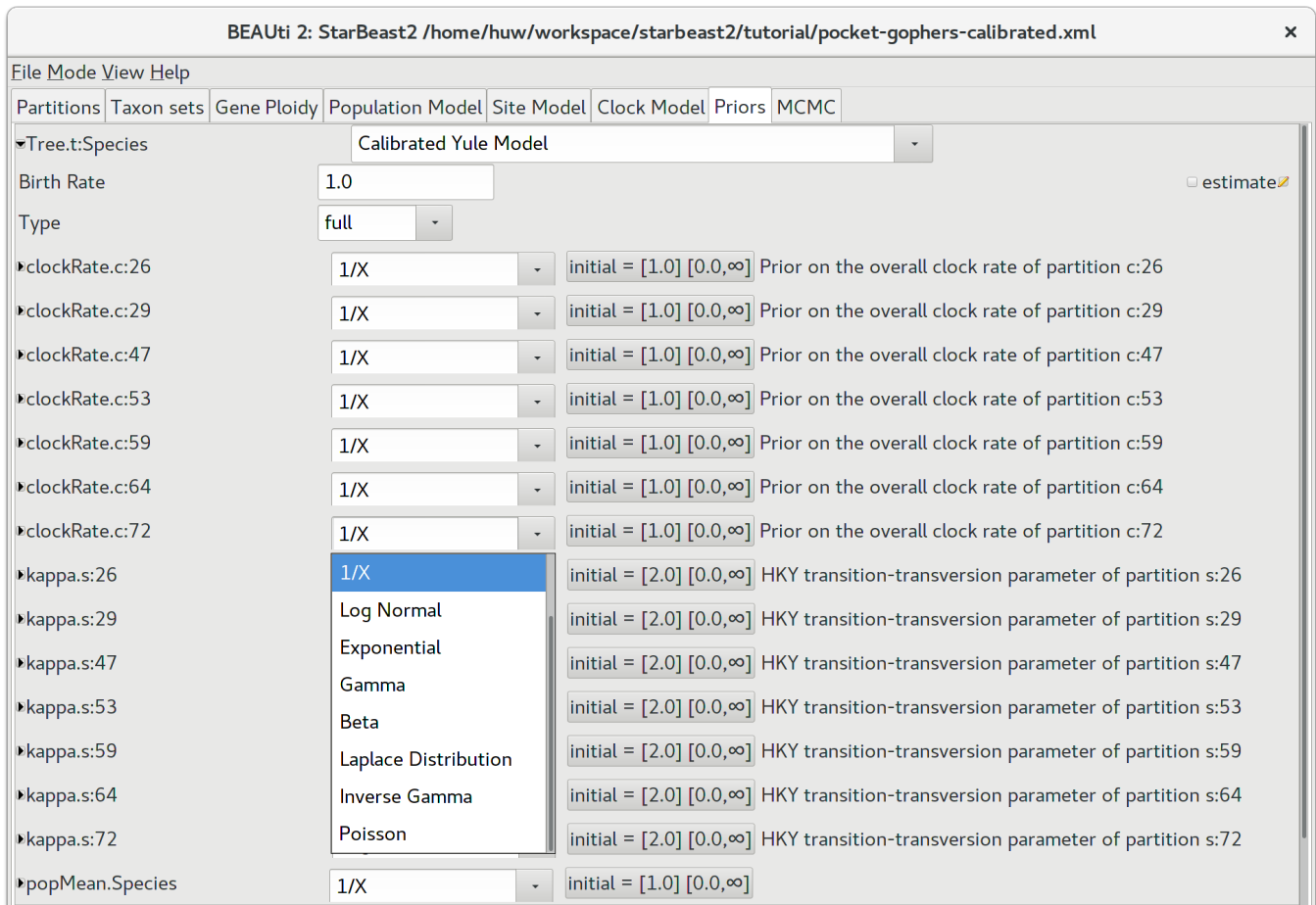


Figure 8: Setting an uninformative prior on clock rates.

We can calibrate this analysis by applying a date calibration to the root of the species tree. This is equivalent to the split between *Thomomys* and *Orthogeomys*, estimated to have occured between 6.5ma and 6.8ma (Belfiore et al. 2008).

To add the new date calibration, click the button with the plus symbol (+) at the very bottom of the list of priors, and click OK to add a new MRCA (most recent common ancestor) prior. Select "Tree.t:Species" to apply this prior to the species tree.

9

Because this prior applies to the root node, move all the species from the left hand column to the right hand column using the >> button. Give this taxon set a sensible name, for example "root" (Figure 9).
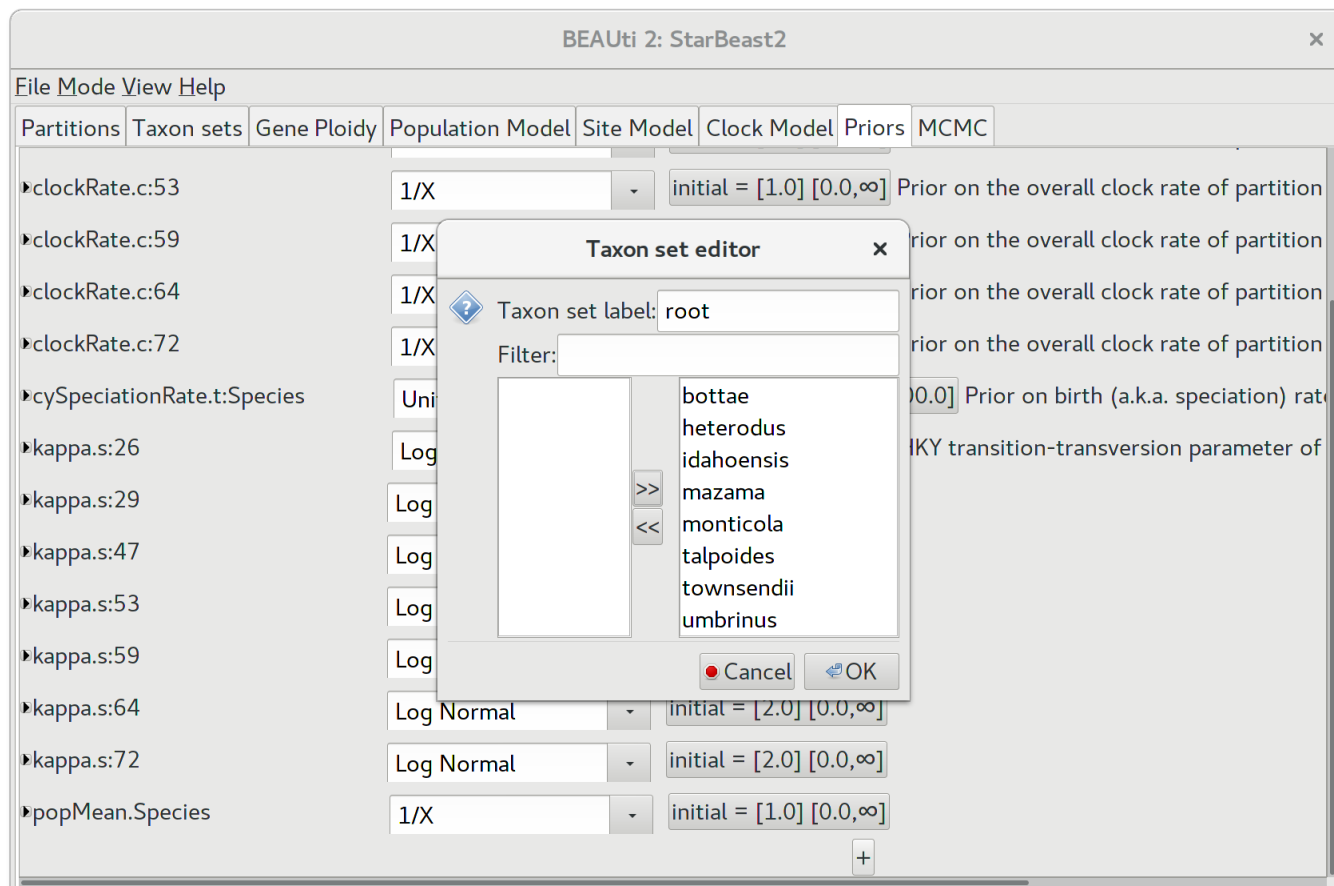


Figure 9: Setting an uninformative prior on clock rates.

Make this prior monophyletic by checking the "monophyletic" tick box. If 6.5ma is a lower bound, we can use an exponential prior with an offset of 6.5 to set a lower limit on the root node height. Set the mean of the exponential distribution to 0.15, so that the mean of the prior is equal to 6.65ma, the midpoint of the estimated range (Figure 10).

## Generating the BEAST XML file

As in the previous tutorial, stick to the default settings in the MCMC tab. We are now ready to create the BEAST XML file. To do this, either select the **File/Save** or **File/Save As** menu options. Save the file with an appropriate name (we usually end the filename with ".xml", e.g. "pocket-gophers-calibrated.xml"). We are now ready to run the file through BEAST.

## Running BEAST

Now run BEAST and when it asks for an input file, provide your newly created XML file as input by clicking **Choose File...**, and then click **Run**. In Linux BEAST will immediately launch a file opening dialog box, which is to select the BEAST XML to run. BEAST will then run until it has finished reporting information to the screen. The actual results files are saved to the disk in the same location as your input file.
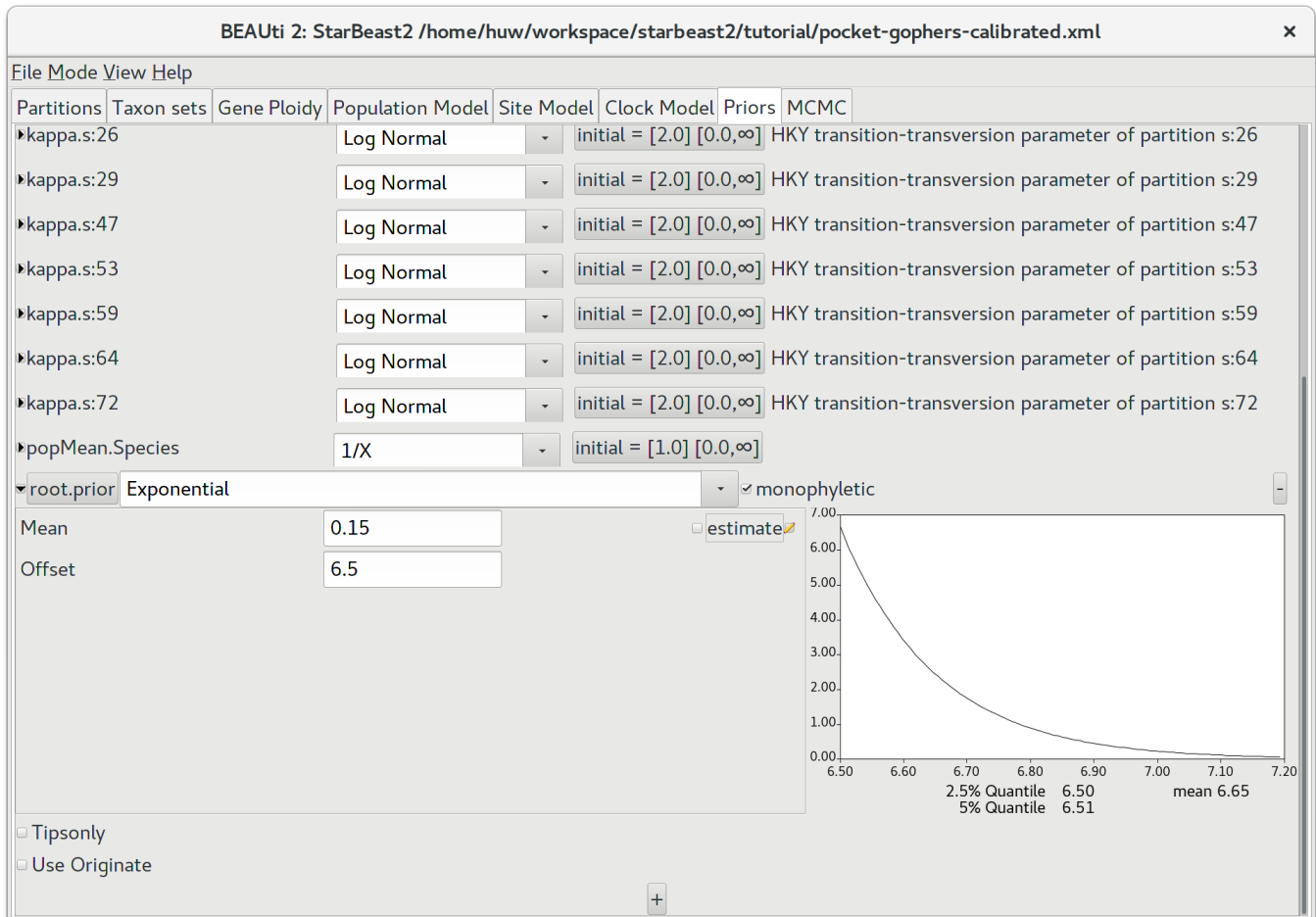
Figure 10: Setting an uninformative prior on clock rates.

# 3   Analyzing the results

Run the program called **Tracer** to analyze the output of BEAST. When the main window has opened, choose **Import Trace File...** from the **File** menu and select the file that BEAST has created called "starbeast.log". Select the clock rates for all genes using the shift key, then click on the "Marginal Prob Distribution" tab. You should now see a window like in Figure 11.

This display shows the marginal distributions of clock rates for all loci. Remember that MCMC is a stochastic algorithm so the actual distributions will not be exactly the same. The clock rates (in units of substitutions per million years) range from about 0.001 for genes 29 and 72 to 0.004 for locus 64, demonstrating that the loci sequenced by Belfiore et al. 2008 have slower and faster rates of molecular evolution respectively.
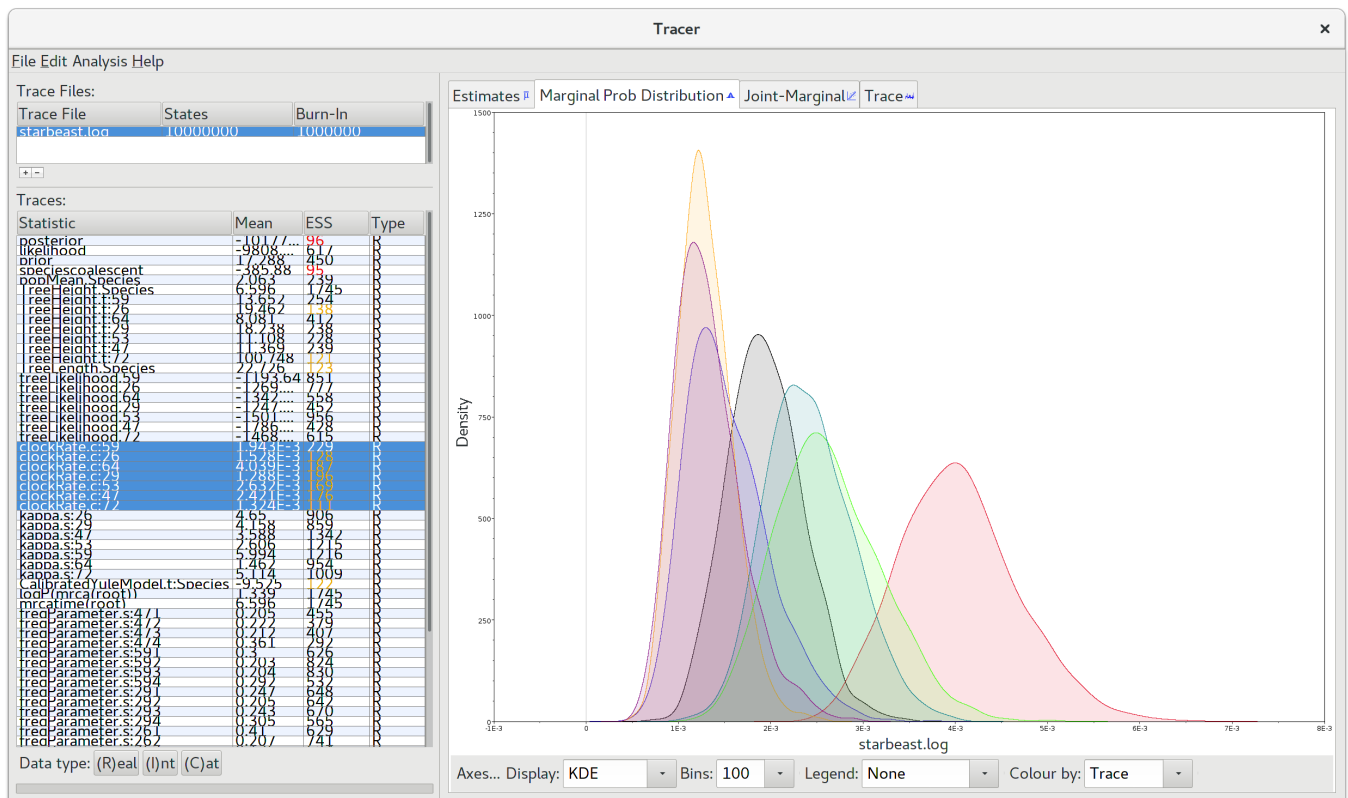
Figure 11: Tracer with the gopher data.

## Obtaining an estimate of the phylogenetic tree

Summarize the posterior sample produced by BEAST by running the **TreeAnnotator** program and setting it up to look like in Figure 12.
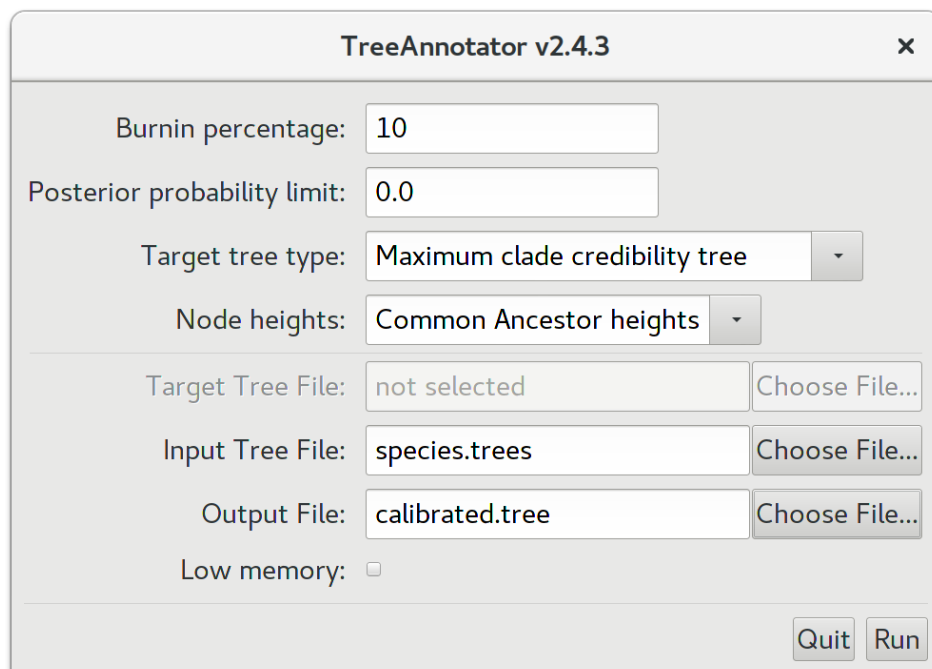


Figure 12: Using TreeAnnotator to summarise the tree set.

The **Burnin percentage** is the proportion of trees to remove from the start of the sample; for this tutorial, set a 10% burnin as shown in Figure 12.

The **Posterior probability limit** option specifies a limit such that if a node is found at less than this frequency in the sample of trees (i.e., has a posterior probability less than this limit), it will not be annotated.

For **Target tree type** you can either choose a specific tree from a file or ask TreeAnnotator to find a tree in your sample. The default option, **Maximum clade credibility tree**, finds the tree with the highest product of the posterior probability of all its nodes.

Keep "Common Ancestor heights" for **Node heights**. This sets the heights (ages) of each node in the tree to the mean height of the most recent common ancestor across the entire set of trees in the posterior.

For the input file, select the trees file that BEAST created (by default this will be called "species.trees") and select a file for the output (here we have called it "calibrated.tree"). Now press **Run** and wait for the program to finish.

## Viewing the species tree(s)

Run **FigTree** and open the "calibrated.tree" file by using the Open command in the File menu. The tree should appear. Select **Node Bars** to get node age error bars. Disable **Scale Bar** and select **Scale Axis** to indicate the ages when speciation events occured. You should end up with something like Figure 13.
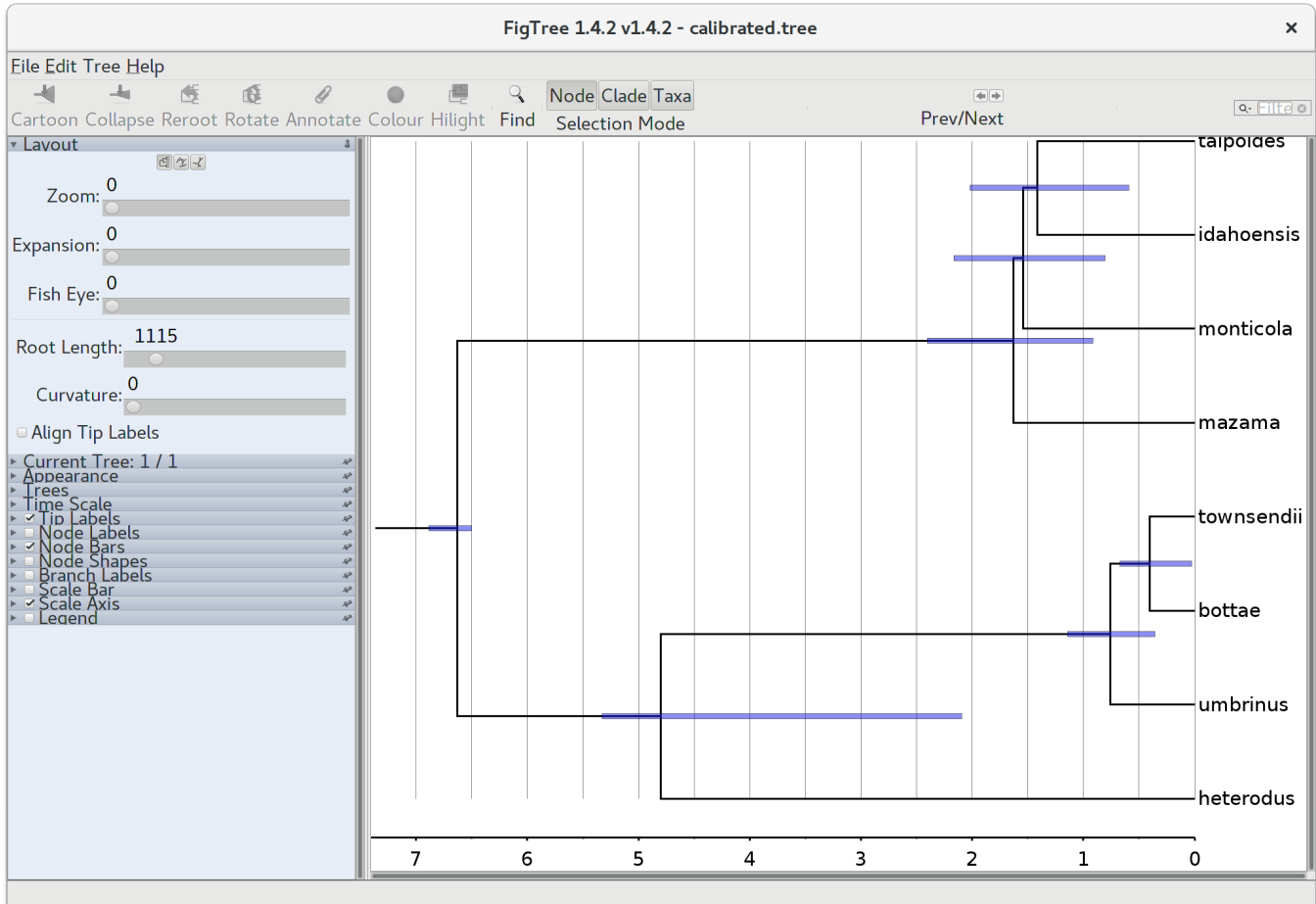


Figure 13: Figtree representation of the species tree.

According to our reanalysis, the radiation of *T. townsendii*, *T. bottae*, and *T. umbrinus* occured quite recently, within the last million years. On the other hand timing of the *T. heterodus* split is uncertain, and could have occured more than 5ma or as recently as 2ma.

Version dated: November 2, 2023

# Relevant References

Belfiore, NM, L Liu, and C Moritz. 2008. Multilocus phylogenetics of a rapid radiation in the genus thomomys (rodentia: geomyidae). *Systematic Biology* 57: 294–310.

Heled, J and AJ Drummond. 2012. Calibrated tree priors for relaxed phylogenetics and divergence time estimation. *Systematic Biology* 61: 138–149.