

StarBEAST2 v15 tutorial

Huw A. Ogilvie^{*1}

¹Department of Computer Science, Rice University, Houston, Texas

December 2, 2020

^{*}huw.a.ogilvie@rice.edu

1 Introduction

StarBEAST2 is a method for multilocus, multispecies coalescent phylogenetic inference (Ogilvie et al. 2017). Multilocus methods like StarBEAST2 take as input one or more multiple sequence alignments derived from genomic loci. There is assumed to be no recombination within each locus and free recombination between loci. Therefore the sequence alignments should be relatively short, and relatively distantly spaced along the genome. Genomic data sets which are typically a good fit to the StarBEAST2 model include RAD-seq loci (Ogilvie et al. 2016) or exons (Scornavacca and Galtier 2017).

StarBEAST2 can be used to estimate species tree topologies and divergence times, the topologies and coalescence times of gene trees, the substitution model rates and base frequencies for those gene trees, per-species population sizes and per-species molecular clock rates.

Newer versions of StarBEAST2 supports the integration of molecular sequences and morphological characters for so-called “total evidence” analyses, and the use of fossil data for tip-dating. Fossil data can include ancient DNA, morphological characters, and the estimated ages of fossils. Tip-dating calibrates the trees in absolute time, typically millions of years.

This tutorial will cover using StarBEAST2 to estimate a species tree including per-species population sizes of the genus *Canis* from exon sequences (Section 3), to estimate per-species molecular clock rates (Section 4), and to conduct a total evidence tip-dating study (Section 5).

2 Preliminaries

Before doing anything, we need to install StarBEAST2, which is available as a package for the phylogenetic software platform BEAST2. If you have not yet installed BEAST2, or if you have an older version, you can download the latest version from <http://beast2.org>. Make sure you have the latest version before proceeding.

After downloading and installing BEAST2 on your computer, open BEAUTi — the graphical user interface for BEAST2. To install StarBEAST2 (or any BEAST2 package), open the File menu and select Manage Packages (Figure 1).

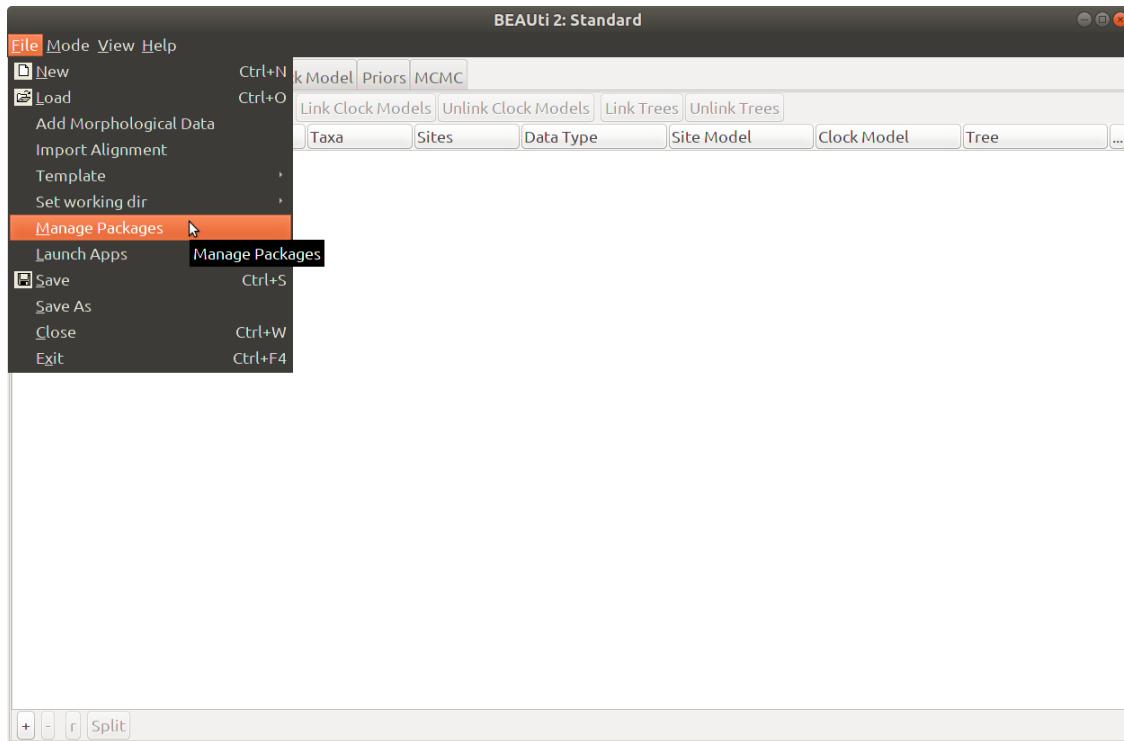


Figure 1: Opening the package manager

Select StarBEAST2 from the list of available packages and click the install button, which will take care of the installation for you (Figure 2). After the installation of StarBEAST2 is finished, just like for any BEAST2 package, you **must** quit and relaunch BEAUTi.

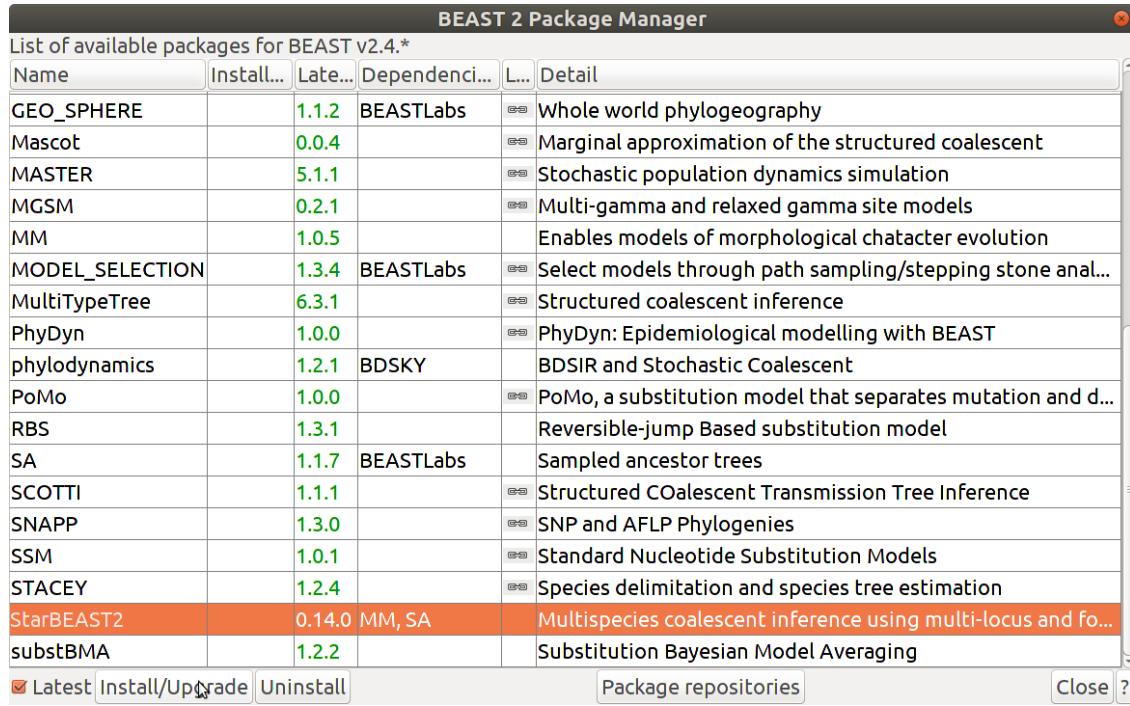


Figure 2: Installing StarBEAST2

After relaunching BEAUTi, you will notice four new templates (Figure 3). The first is “StarBeast2”, which is for a strict molecular clock or relaxed gene tree clocks uncorrelated with the species tree lineages. The others are “SpeciesTreeUCLN”, “SpeciesTreeRLC” and “SpeciesTreeUCED” which enable per-species clock rates, as described by Ogilvie et al. 2016. The relaxed clock models used by those templates are uncorrelated log-normal (UCLN), relaxed local clock (RLC), and uncorrelated exponential distribution (UCED) respectively.

To visualize posterior distributions of trees, we will be using DensiTree (Bouckaert 2010) which is included with BEAST2. We will also use FigTree to view summary trees, and Tracer to check on parameters and statistics. FigTree is available from <http://tree.bio.ed.ac.uk/software/figtree/> and Tracer from <http://tree.bio.ed.ac.uk/software/tracer>. Please make sure both are installed before proceeding.

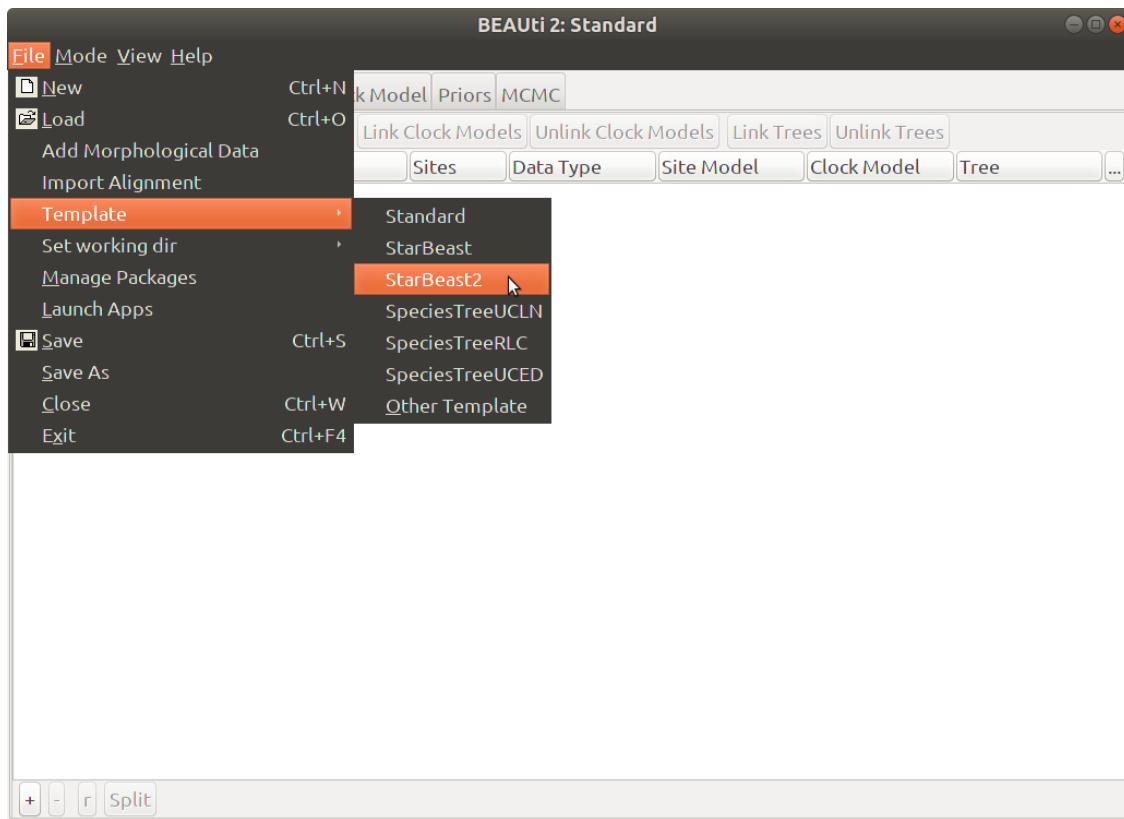


Figure 3: Selecting a StarBEAST2 template

3 Reconstructing the species tree of *Canis*

The genus *Canis* includes species such as *Canis lupus* (wolves), *Canis latrans* (coyotes), and *Canis dirus* (dire wolves, extinct). Because the taxonomy of this genus is quite messy, it also includes *Cuon alpinus* and *Lycaon pictus*. In this section we will reconstruct the species tree of *Canis* using StarBEAST2.

3.1 Importing alignments

First up, create a new folder somewhere with a sensible name like “CanisPhylogeny”. Relaunch BEAUTi and load up the strict clock template by selecting “StarBeast2” from the Templates submenu of the File menu (Figure 3). Now we will import multiple sequence alignments of the 16 nuclear loci sequenced by Lindblad-Toh et al. 2005. Select Import Alignments from the File menu, and navigate to the “data” subfolder of the tutorial. Select all 16 FASTA files, and click OK to import (Figure 4). Do not select “morphology-canis.nexus” which is for the total evidence study.

After you click OK, you will be asked what type of sequences are in the files. As all 16 files are MSAs of nucleotide sequences, select “all are nucleotide” and click OK again to continue. There should now be 16 partitions listed in BEAUTi.

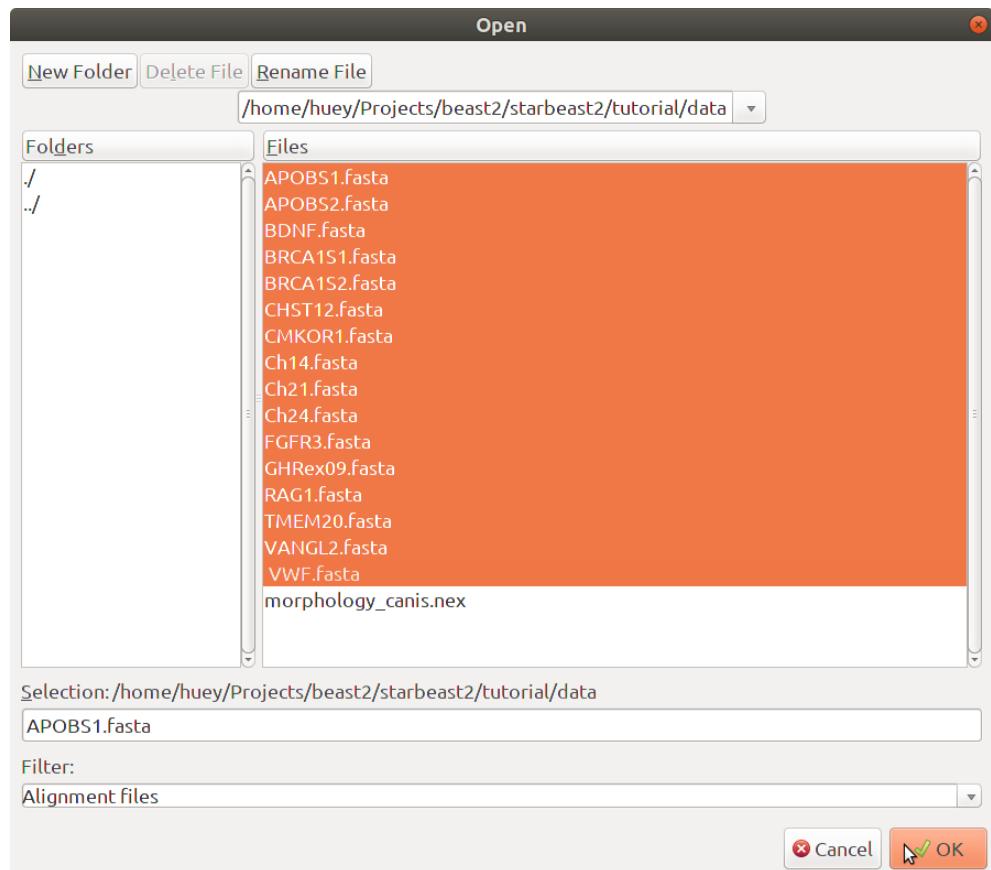


Figure 4: Importing *Canis* multiple sequence alignments

3.2 Linking models

As a general rule, clock models should be linked when using the strict clock StarBEAST2 template. Rates between loci are still allowed to vary and the mean rate among all sites is fixed at 1. Therefore by linking the clocks, the average clock rate can be fixed at an *a priori* value to calibrate the analysis, or estimated when node-dating or tip-dating is used. Select all 16 loci and click the “Link Clock Models” button. Now all loci should have “APOBS1” as the clock model name. The site model, clock model and tree names are text fields that you can click on to edit. Change the name of the clock model to “molecular” and it should look like Figure 5.

The screenshot shows the BEAUti 2: StarBeast2 software interface. The main window title is "BEAUti 2: StarBeast2". The menu bar includes "File", "Mode", "View", and "Help". Below the menu is a toolbar with buttons for "Partitions", "Taxon sets", "Tip Dates", "Gene Ploidy", "Population Model", "Site Model", "Clock Model", "Priors", "MCMC", "Link Site Models", "Unlink Site Models", "Link Clock Models", "Unlink Clock Models", "Link Trees", and "Unlink Trees". The main area is a table with columns: Name, File, Taxa, Sites, Data Type, Site Model, Clock Model, and Tree. The table lists 16 entries, each corresponding to a gene: APOBS1, APOBS2, BDNF, BRCA1S1, BRCA1S2, CHST12, CMKOR1, Ch14, Ch21, Ch24, FGFR3, GHRex09, RAG1, TMEM20, VANGL2, and VWF. All entries have "molecular" selected in the Site Model, Clock Model, and Tree dropdown menus. The "Clock Model" column for all entries shows "APOBS1". The "Site Model" column for all entries shows "molecular". The "Tree" column for all entries shows the gene name. At the bottom of the table is a toolbar with buttons for "+", "-", "r", "Split", and a search field.

Name	File	Taxa	Sites	Data Type	Site Model	Clock Model	Tree
APOBS1	APOBS1	16	702	nucleotide	APOBS1	molecular	APOBS1
APOBS2	APOBS2	16	633	nucleotide	APOBS2	molecular	APOBS2
BDNF	BDNF	16	489	nucleotide	BDNF	molecular	BDNF
BRCA1S1	BRCA1S1	16	705	nucleotide	BRCA1S1	molecular	BRCA1S1
BRCA1S2	BRCA1S2	16	741	nucleotide	BRCA1S2	molecular	BRCA1S2
CHST12	CHST12	16	705	nucleotide	CHST12	molecular	CHST12
CMKOR1	CMKOR1	16	735	nucleotide	CMKOR1	molecular	CMKOR1
Ch14	Ch14	16	921	nucleotide	Ch14	molecular	Ch14
Ch21	Ch21	16	601	nucleotide	Ch21	molecular	Ch21
Ch24	Ch24	16	730	nucleotide	Ch24	molecular	Ch24
FGFR3	FGFR3	16	503	nucleotide	FGFR3	molecular	FGFR3
GHRex09	GHRex09	16	736	nucleotide	GHRex09	molecular	GHRex09
RAG1	RAG1	16	741	nucleotide	RAG1	molecular	RAG1
TMEM20	TMEM20	16	615	nucleotide	TMEM20	molecular	TMEM20
VANGL2	VANGL2	16	546	nucleotide	VANGL2	molecular	VANGL2
VWF	VWF	16	732	nucleotide	VWF	molecular	VWF

Figure 5: Linking clock models

Mixing nuclear and mitochondrial loci

StarBEAST2 can be used with a combination of mitochondrial and nuclear loci, but this requires careful thought when specifying the model. One possibility is to link the nuclear loci as above, but not to link the mitochondrial loci. Then when editing the site models, untick “estimate” for the substitution rate of each mitochondrial locus. Now separate clock rates will be used for each mitochondrial locus, in addition to the overall clock rate which applies to nuclear loci.

If the species tree is calibrated using an *a priori* nuclear or mitochondrial molecular clock rate, you can set that rate in the Clock Model panel. If the mitochondrial rate is fixed, the nuclear rate should probably be estimated and *vice versa*. If the species tree is calibrated using node or tip dating, then all clock rates should probably be estimated.

Make sure appropriate priors are used for all estimated clock rates. This could be a $1/X$ prior for rates where you genuinely have no prior knowledge. However in most circumstances you will have a general idea of the clock rate, e.g. within an order of magnitude. In that case, you can set a broad prior with a mean based on your prior knowledge. Broad priors include Log Normal with a large standard deviation (up to about 2), or a Gamma distribution with a small α (set “alpha” to 1 or 2).

Linking site models

When clock models are linked as is recommended by this tutorial, linking site models will result in every gene having exactly the same clock rate. As this is unrealistic, generally site models should **not** be linked using BEAUTi when setting up a StarBEAST2 analysis.

3.3 Specifying species names

In the FASTA files, each sequence has a name like “Canis_anthus_a”. Here *Canis anthus* is the binomial name, and “a” is a haplotype. For the data supplied with this tutorial, each species has two haplotypes “a” and “b”, both from the same diploid individual. However in other studies multiple individuals may be sequenced, so an arbitrary number of haplotypes are available per species.

To assign haplotypes to species, select the Taxon Sets tab in BEAUTi, then click the Guess button. To assign the correct names, keep “use everything” selected, but change “after first” to “before last”. Leave the underscore in the text box and click OK (Figure 6). This way everything before the last underscore in the haplotype name will be used as the species name, so “Canis_anthus” will be the species name for “Canis_anthus_a”.

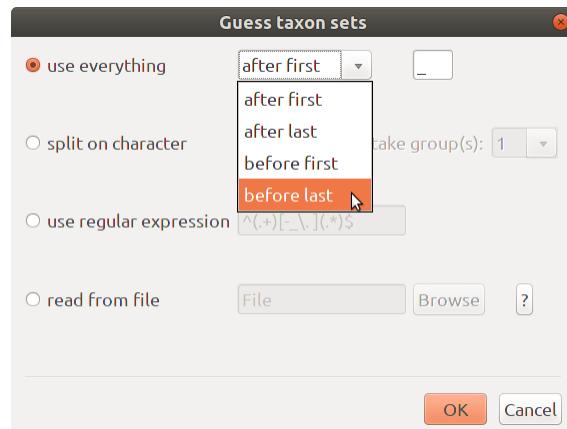
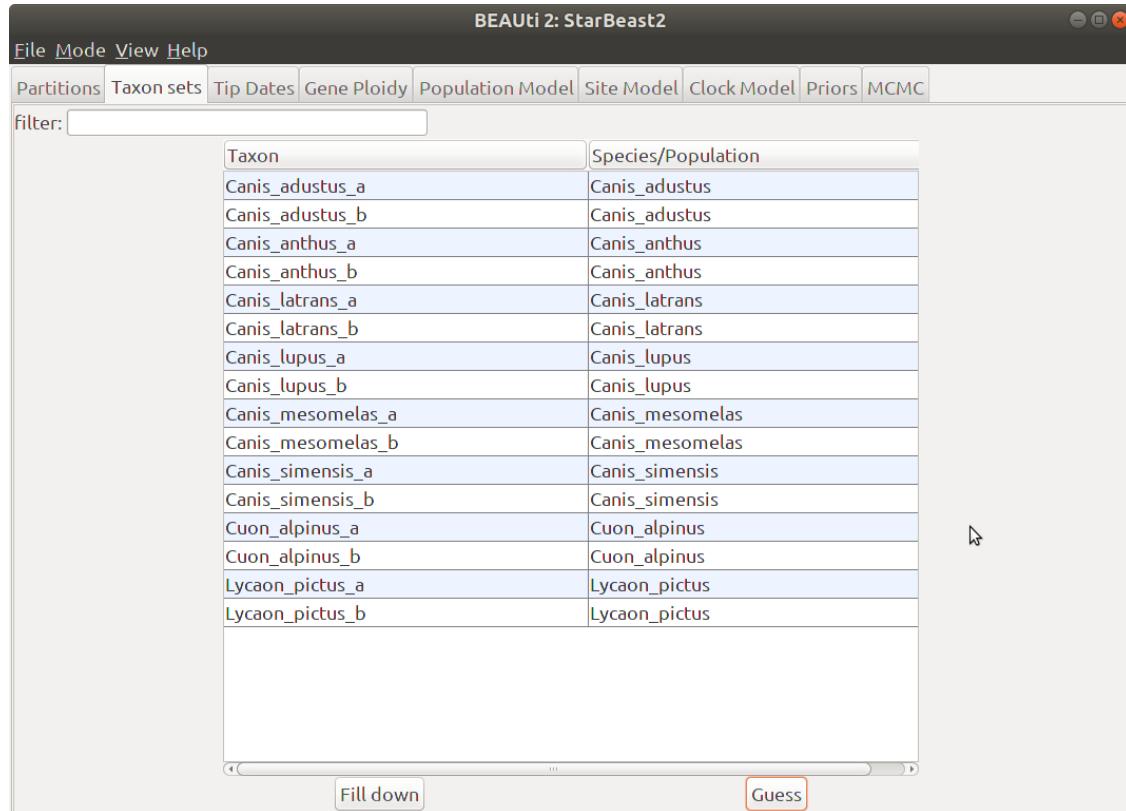


Figure 6: Guessing species names from haplotype names

Now each haplotype (Taxon) should have a corresponding species (Species/Population). Based on how we assigned the species names, there should be two haplotypes “a” and “b” for each species (Figure 7).



The screenshot shows the BEAUti 2 interface for StarBeast2. The window title is "BEAUti 2: StarBeast2". The menu bar includes "File", "Mode", "View", and "Help". Below the menu is a tab bar with "Partitions", "Taxon sets", "Tip Dates", "Gene Ploidy", "Population Model", "Site Model", "Clock Model", "Priors", and "MCMC". The "Taxon sets" tab is selected. A "filter:" input field is present. The main area contains a table with two columns: "Taxon" and "Species/Population". The table lists 18 entries, grouped by species: Canis_adustus_a, Canis_adustus_b; Canis_anthus_a, Canis_anthus_b; Canis_latrans_a, Canis_latrans_b; Canis_lupus_a, Canis_lupus_b; Canis_mesomelas_a, Canis_mesomelas_b; Canis_simensis_a, Canis_simensis_b; Cuon_alpinus_a, Cuon_alpinus_b; Lycaon_pictus_a, Lycaon_pictus_b. The "Species/Population" column shows the species name for each taxon. At the bottom of the table are "Fill down" and "Guess" buttons.

Taxon	Species/Population
Canis_adustus_a	Canis_adustus
Canis_adustus_b	Canis_adustus
Canis_anthus_a	Canis_anthus
Canis_anthus_b	Canis_anthus
Canis_latrans_a	Canis_latrans
Canis_latrans_b	Canis_latrans
Canis_lupus_a	Canis_lupus
Canis_lupus_b	Canis_lupus
Canis_mesomelas_a	Canis_mesomelas
Canis_mesomelas_b	Canis_mesomelas
Canis_simensis_a	Canis_simensis
Canis_simensis_b	Canis_simensis
Cuon_alpinus_a	Cuon_alpinus
Cuon_alpinus_b	Cuon_alpinus
Lycaon_pictus_a	Lycaon_pictus
Lycaon_pictus_b	Lycaon_pictus

Figure 7: The assignment of haplotypes to species

3.4 Gene Ploidy and Population Model

Ignore the Tip Dates tab, which is for tip-dating and will be covered in section 5 of the tutorial.

Open the Gene Ploidy tab, and observe that the default value for all loci is 2.0. This is because there are two copies of a locus in each individual for diploid populations, so we scale the effective population sizes by 2.0. If you use any mitochondrial or Y/W chromosomal loci, you should change their ploidy to 0.5, because there is on average only 0.5 copies per individual. The ploidy of X/Z chromosomal loci should be set to 1.5 for the same reason.

Select the Population Model tab. By default analytical integration is used, which is slightly faster but does not produce estimates of per-species population sizes. Change the model to “Constant Populations”, which will add effective population sizes to the species tree output (Figure 8).

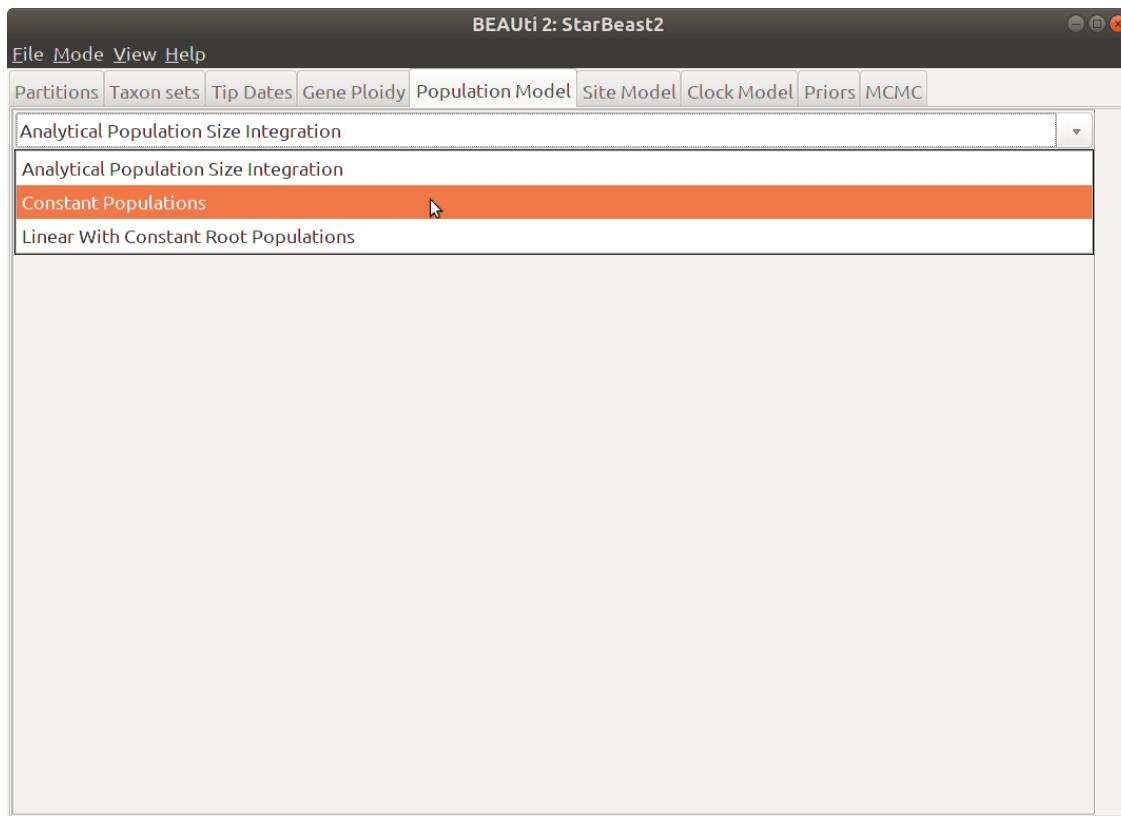


Figure 8: Changing the population model

3.5 Site Model

Select the Site Model tab, and you will see the site model for the first partition displayed. Change “JC69” to “HKY”, and then set the frequencies to empirical (Figure 9). HKY is more flexible because it allows nucleotide transitions to have a different rate relative to transversions (Hasegawa et al. 1985).

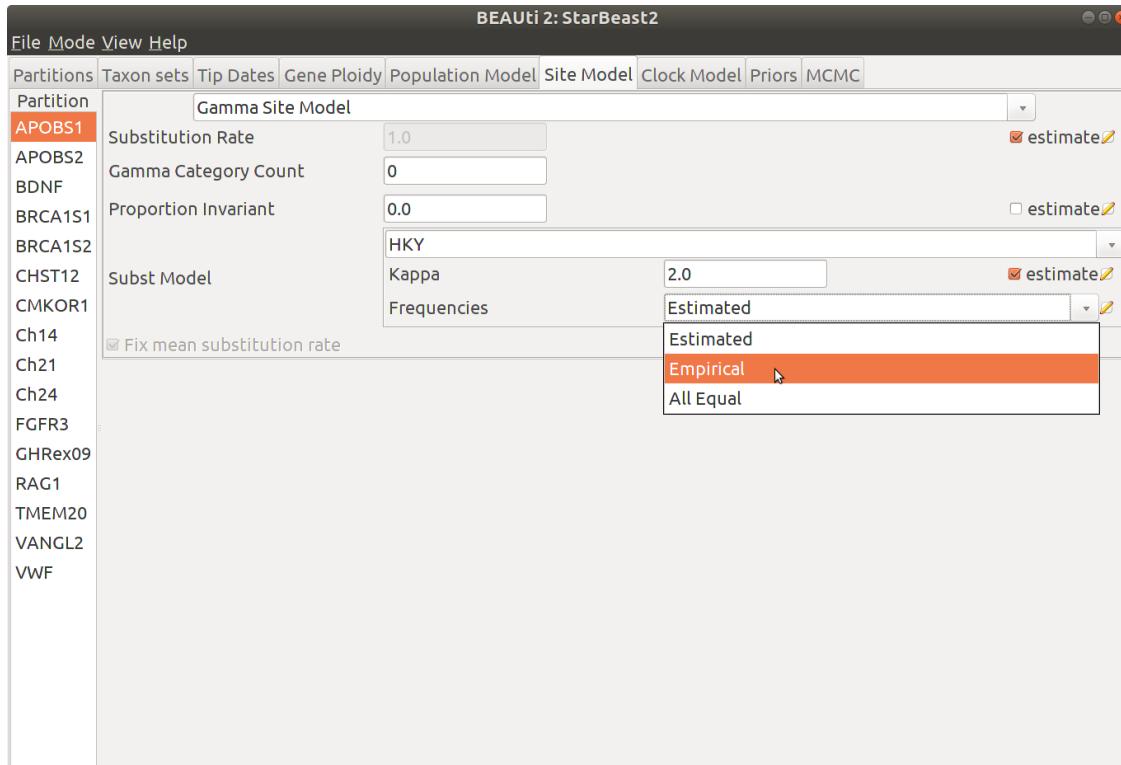


Figure 9: Setting the site model to HKY with empirical base frequencies

Now to change the site model for all loci to HKY with empirical base frequencies, select all of the partitions in the left hand column with the shift key. Then click “OK” to apply the same site model used for the first locus to everything else (Figure 10).

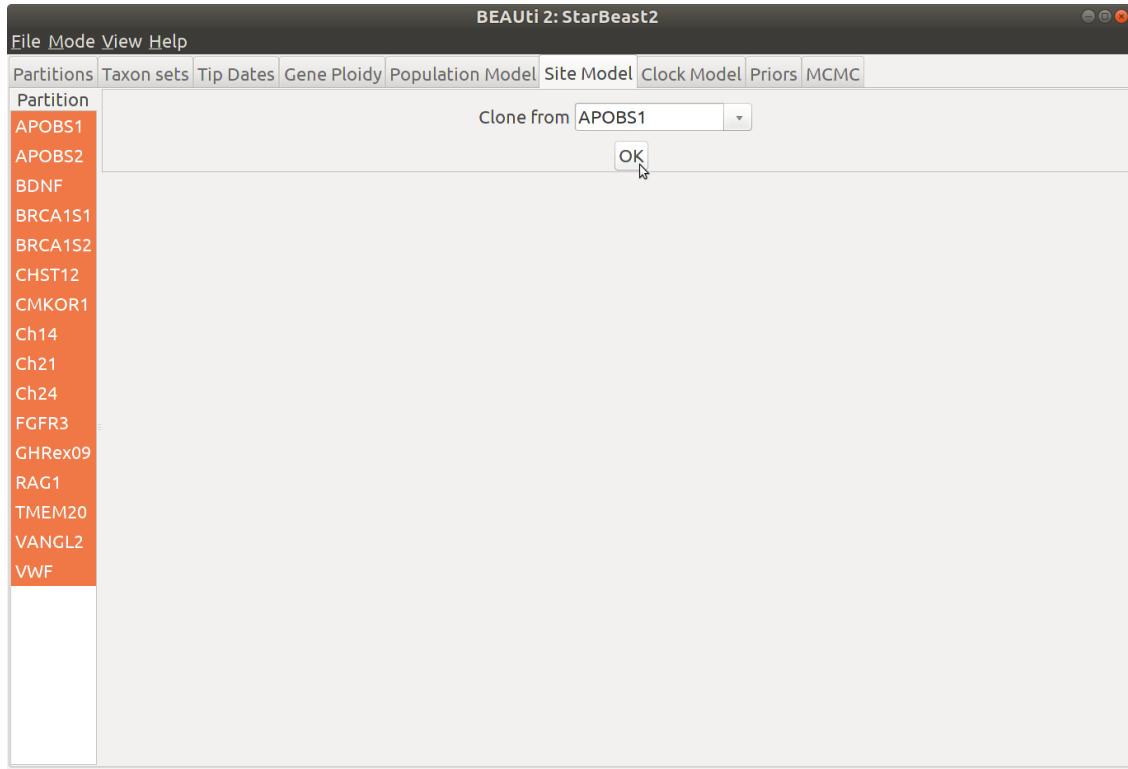


Figure 10: Cloning the site model

In a more serious analysis, you might consider setting the number of gamma categories to 4. This allows for different sites to evolve at different rates, an obviously more realistic model (Yang 1994). However the phylogenetic likelihood must be calculated once for each category, so 4 categories will be 4× slower. That likelihood calculation is a major part of the StarBEAST2 algorithm, so more gamma rate categories will require more computer time.

3.6 Clock Model

Hugall et al. 2007 estimated that the molecular clock rate for the RAG-1 nuclear coding gene in mammals is approximately 10^{-3} substitutions per site per million years. While we don't know the average rate across all loci in our data within *Canis*, we can use it as a **very approximate** calibration. Open the clock model panel and set the rate to “0.001”, to match the *a priori* estimate (Figure 11).

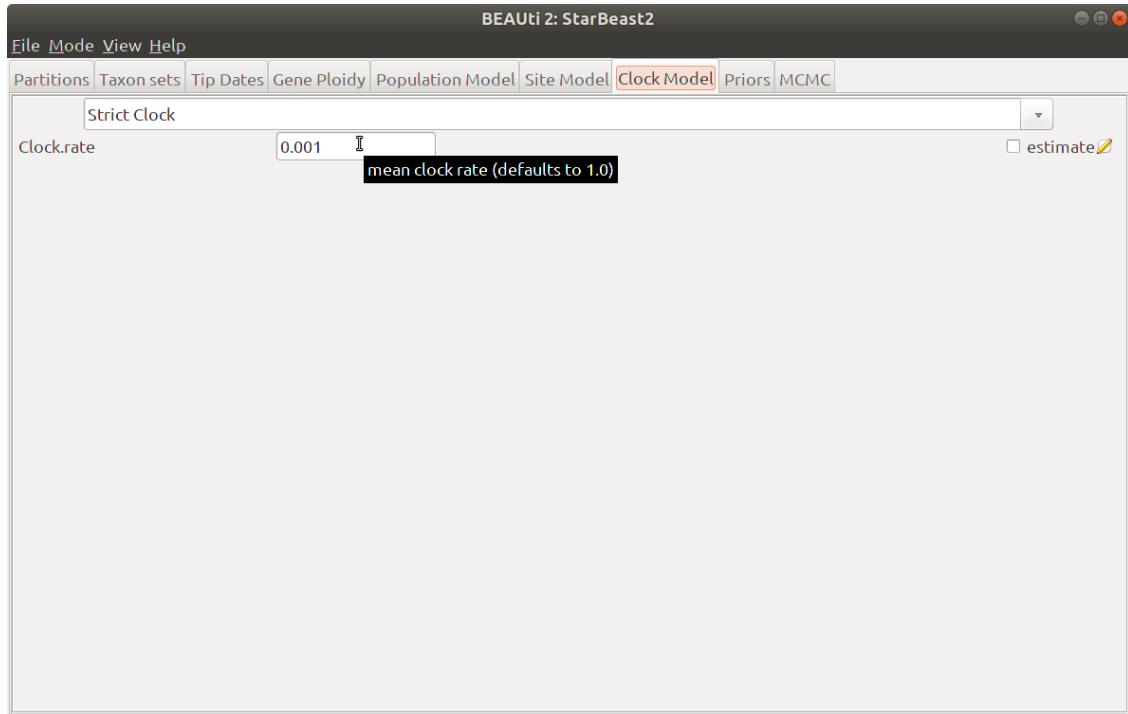


Figure 11: Using an *a priori* clock rate for calibration

You can also use the scientific notation shorthand, 1e-3, which is equivalent to 0.001.

3.7 Priors and MCMC

Open the Priors panel and change “Yule Model” to “Birth Death Model”. These models are identical, except the birth-death model allows for a non-zero rate of extinction.

StarBEAST2 is an MCMC method. These kind of methods do better at estimating the posterior distributions (of trees or other parameters) the longer they are run, although after a point there are diminishing returns. The default chain length in StarBEAST2 is 10 million states, but for this analysis we need a bit more for good estimates; change the chain length to 40 million (Figure 12)

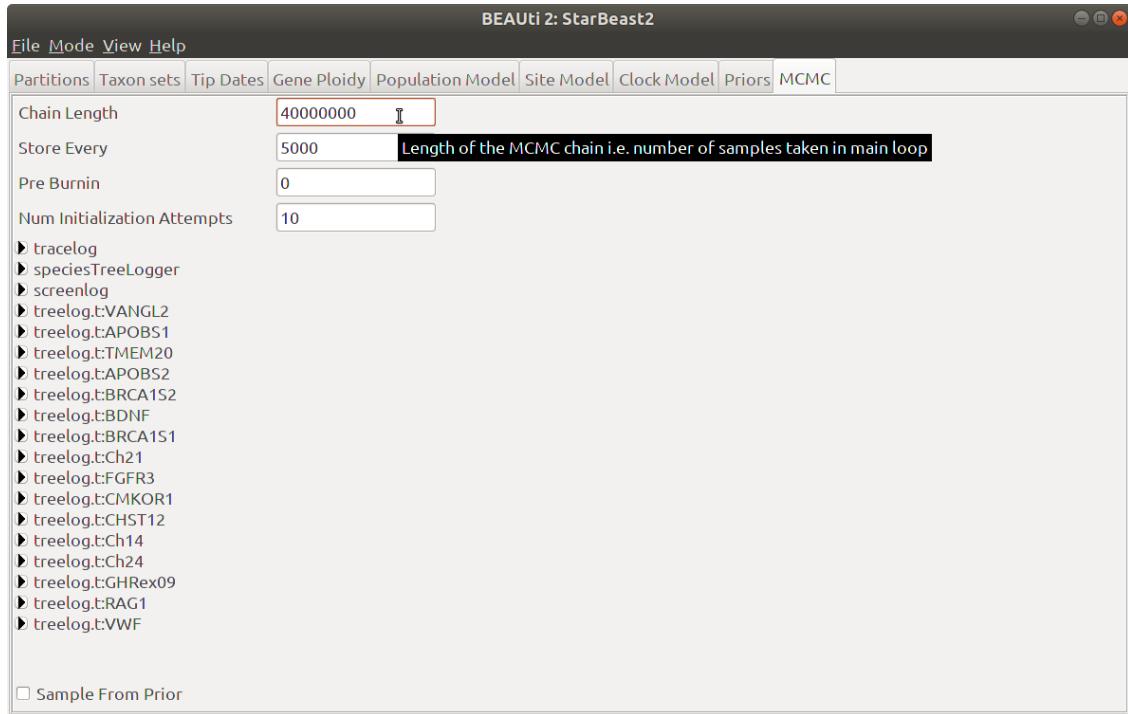


Figure 12: Setting a longer MCMC chain length

Save your model as an XML file – in the folder you created before importing alignments – by clicking Save As in the File menu. Navigate to the folder and give your XML file a name like “CanisPhylogeny.xml”.

3.8 Running BEAST

You can run BEAST from the command line or using the GUI. To run your XML from the command line, first navigate to the folder you saved the XML file into. Then run “/path/to/beast/bin/beast CanisPhylogeny.xml”. Make sure to change /path/to/beast to match the folder where BEAST is installed on your computer.

BEAST should take about 10 to 15 minutes to finish the chain. Sampled statistics and various parameters will be saved to “starbeast.log”, species trees to “species.trees”, and separate gene tree files

will be created for each locus. The command line output should start off looking something like what follows:

```

Checking out /home/huey/.beast/2.4/BEAST/lib
Loaded URL file:/home/huey/.beast/2.4/BEAST/lib/beast.jar
jardir = /home/huey/beast/lib/launcher.jar
Loading package BEAST v2.4.7
Loading package MM v1.0.5
Loading package SA v1.1.7
Loading package BEAST v2.4.7
Loading package starbeast2 v0.14.0
Loading package BEASTLabs v1.7.1

        BEAST v2.4.7 Prerelease, 2002-2017
        Bayesian Evolutionary Analysis Sampling Trees
        Designed and developed by
        Remco Bouckaert, Alexei J. Drummond, Andrew Rambaut & Marc A. Suchard

        Department of Computer Science
        University of Auckland
        remco@cs.auckland.ac.nz
        alexei@cs.auckland.ac.nz

        Institute of Evolutionary Biology
        University of Edinburgh
        a.rambaut@ed.ac.uk

        David Geffen School of Medicine
        University of California, Los Angeles
        msuchard@ucla.edu

        Downloads, Help & Resources:
        http://beast2.org/

Source code distributed under the GNU Lesser General Public License:
        http://github.com/CompEvol/beast2

        BEAST developers:
        Alex Alekseyenko, Trevor Bedford, Erik Bloomquist, Joseph Heled,
        Sebastian Hoehna, Denise Kuehnert, Philippe Lemey, Wai Lok Sibon Li,
        Gerton Lunter, Sidney Markowitz, Vladimir Minin, Michael Defoin Platel,
        Oliver Pybus, Chieh-Hsi Wu, Walter Xie

        Thanks to:
        Roald Forsberg, Beth Shapiro and Korbinian Strimmer

Random number seed: 1512645447217
File: CanisPhylogeny.xml seed: 1512645447217 threads: 1

```

It will end with a bunch of statistics describing the performance of the MCMC operators:

Operator	Tuning	#accept	#reject	Pr(m)	Pr(acc m)
UpDownOperator(clockUpDownOperator.c:BRCA1S1)	0.8522	18169	50825	0.0019	0.2633
ScaleOperator(TreeScaler.t:BRCA1S1)	0.8682	19465	49496	0.0019	0.2823
ScaleOperator(TreeRootScaler.t:BRCA1S1)	0.4248	14918	54203	0.0019	0.2158
Uniform(UniformOperator.t:BRCA1S1)	-	146565	198613	0.0094	0.4246
SubtreeSlide(SubtreeSlide.t:BRCA1S1)	0.7795	120013	225596	0.0094	0.3473
Exchange(Narrow.t:BRCA1S1)	-	109470	235681	0.0094	0.3172
Exchange(Wide.t:BRCA1S1)	-	8117	336186	0.0094	0.0236
WilsonBalding(WilsonBalding.t:BRCA1S1)	-	10344	335125	0.0094	0.0299
[...]					
DeltaExchangeOperator(FixMeanMutationRatesOperator)	2.0294	88056	648890	0.0201	0.1195
ScaleOperator(KappaScaler.s:AP0BS1)	0.1987	7059	15949	0.0006	0.3068
ScaleOperator(KappaScaler.s:AP0BS2)	0.2034	7243	15463	0.0006	0.3190
ScaleOperator(KappaScaler.s:BDNF)	0.1660	6783	16531	0.0006	0.2909
ScaleOperator(KappaScaler.s:BRCA1S1)	0.1857	6903	15965	0.0006	0.3019
ScaleOperator(KappaScaler.s:BRCA1S2)	0.1976	6794	16280	0.0006	0.2944
ScaleOperator(KappaScaler.s:CHST12)	0.1751	6930	16559	0.0006	0.2950
ScaleOperator(KappaScaler.s:CMK0R1)	0.1976	7162	15697	0.0006	0.3133
ScaleOperator(KappaScaler.s:Ch14)	0.2973	6333	16470	0.0006	0.2777
ScaleOperator(KappaScaler.s:Ch21)	0.2553	7159	15909	0.0006	0.3103
ScaleOperator(KappaScaler.s:Ch24)	0.2272	6810	16334	0.0006	0.2942
ScaleOperator(KappaScaler.s:FGFR3)	0.2009	7028	15842	0.0006	0.3073
ScaleOperator(KappaScaler.s:GHREx09)	0.1926	7244	15595	0.0006	0.3172
ScaleOperator(KappaScaler.s:RAG1)	0.2242	7336	15792	0.0006	0.3172
ScaleOperator(KappaScaler.s:TMEM20)	0.2046	7077	15773	0.0006	0.3097
ScaleOperator(KappaScaler.s:VANGL2)	0.1901	6838	15992	0.0006	0.2995
ScaleOperator(KappaScaler.s:VWF)	0.2174	7618	15565	0.0006	0.3286
starbeast2.NodeReheight2(Reheight.t:Species)	-	190932	2793110	0.0471	0.0640
starbeast2.CoordinatedUniform(coordinatedUniform.t:Species)	-	277295	318958	0.0094	0.4651
starbeast2.CoordinatedExponential(coordinatedExponential.t:Species)	0.0724	357000	240475	0.0094	0.5975
UpDownOperator(updownAll.Species)	0.6391	50943	186734	0.0038	0.2143

starbeast2.RealCycle(constPopSizesSwap.Species)	2.0000	23521	95986	0.0019	0.1968	k = 2
ScaleOperator(constPopSizesScale.Species)	0.2694	28048	91487	0.0019	0.2346	
ScaleOperator(constPopMeanScale.Species)	0.4184	11724	27838	0.0006	0.2963	
ScaleOperator(netDiversificationRateScale.t:Species)	0.2074	10390	29138	0.0006	0.2629	
ScaleOperator(ExtinctionFractionScale.t:Species)	0.1502	5748	14261	0.0003	0.2873	
UniformOperator(ExtinctionFractionUniform.t:Species)	-	10374	9609	0.0003	0.5191	
SubtreeSlide(bdSubtreeSlide.t:Species)	0.0697	177124	419681	0.0094	0.2968	
WilsonBalding(bdWilsonBalding.t:Species)	-	371	597396	0.0094	0.0006	
Exchange(bdWide.t:Species)	-	3255	594405	0.0094	0.0054	
Exchange(bdNarrow.t:Species)	-	29181	567726	0.0094	0.0489	
Uniform(bdUniformOperator.t:Species)	-	45371	551315	0.0094	0.0760	
ScaleOperator(bdTreeRootScaler.t:Species)	0.9259	9555	109831	0.0019	0.0800	
ScaleOperator(bdTreeScaler.t:Species)	0.9876	21203	98760	0.0019	0.1767	

Tuning: The value of the operator's tuning parameter, or '-' if the operator can't be optimized.

#accept: The total number of times a proposal by this operator has been accepted.
#reject: The total number of times a proposal by this operator has been rejected.
Pr(m): The probability this operator is chosen in a step of the MCMC (i.e. the normalized weight).
Pr(acc|m): The acceptance probability (#accept as a fraction of the total proposals for this operator).

Total calculation time: 756.696 seconds
End likelihood: -16953.439185360676

3.9 Checking the log file

To check parameters other than species tree topologies and branch values, or to verify that the chain has been run long enough to reliably represent the posterior distribution, we will use Tracer. Start the Tracer app and then open the “starbeast.log” file. The first statistic that will be displayed is a histogram of the log posterior probability (Figure 13).

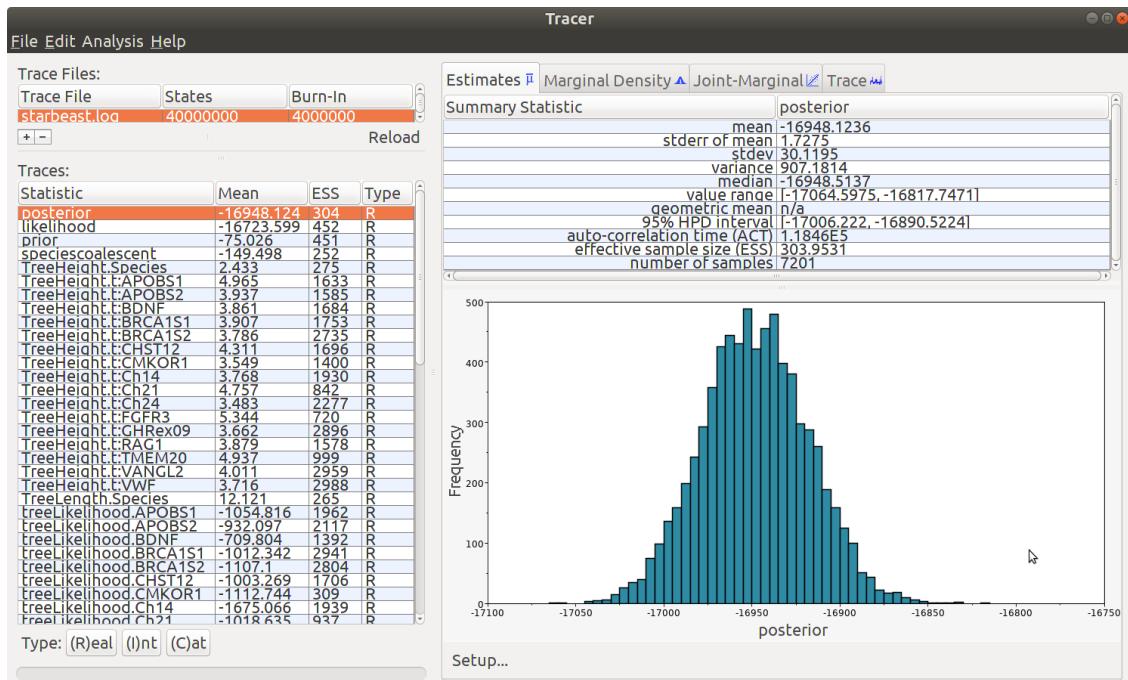


Figure 13: Opening a log file in Tracer

The posterior probability is the sum of the likelihood (which is the sum of log phylogenetic likelihoods for all sites for all loci), the prior probability (which is the sum of log prior probabilities for all parameters), and the speciescoalescent (which is the sum of log coalescent probabilities for all gene trees). TreeHeight.Species is the height of the root node of the species tree, and TreeLength.Species is the sum of all branch lengths in the species tree.

Tracer computes effective sample sizes (ESS) for each logged statistic and parameter. As a rule, ESS values should be at least 200, particularly for the important statistics just noted and for parameters of interest. Because of the stochastic nature of MCMC algorithms, your values **will be different** to those in Figure 13.

3.10 Checking the species trees

Start the DensiTree app (included with BEAST2), and then open the “species.trees” file. Under the Show panel, enable the Root Canal tree to get an idea of the most plausible species tree. Then open the Grid panel, and enable the full grid (Figure 14). If you want to check the clade posterior probabilities, select “View clade toolbar” from the Window menu.

You can see that using a fixed clock rate of 10^{-3} substitutions per site per year, the split between *Canis latrans* (coyotes) and *Canis lupus* (wolves) is probably less than 500,000 years ago. The age of the most recent common ancestor (MRCA) of extant *Canis* taxa is approximately 2.5 million years ago, more recent than the split between humans and chimpanzees (Prado-Martinez et al. 2013).

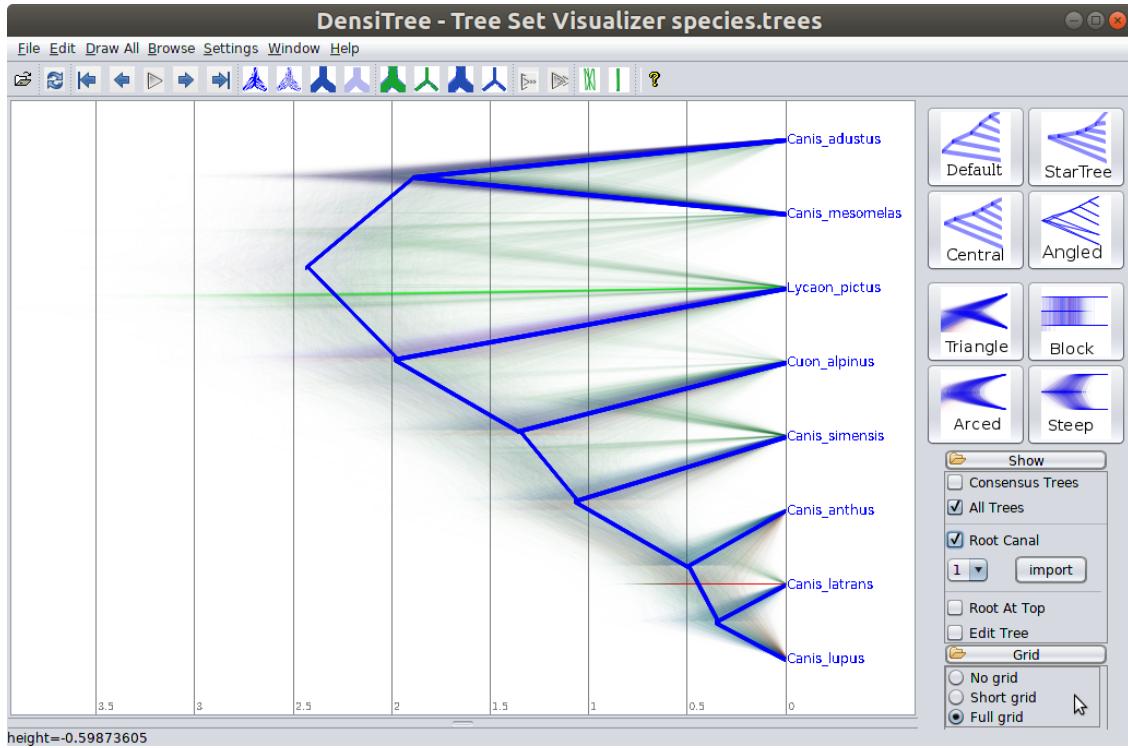


Figure 14: Viewing the species trees in DensiTree

3.11 Generating a summary tree

Summary trees are a way of reducing a posterior distribution of tree topologies and times to a single tree. This is often more readable than the cloud of trees that DensiTree displays, but researchers should be cautious not to give it too much emphasis because for most clades and data sets there is substantial uncertainty in the posterior distribution of topologies and times.

To generate a summary tree, open the Tree Annotator app included with BEAST2. Set the burnin percentage to 10 and choose the “species.trees” file generated by StarBEAST2 as the input file. Specify the output file to be something sensible like “summary.tree” and click Run to generate the summary tree (Figure 15).

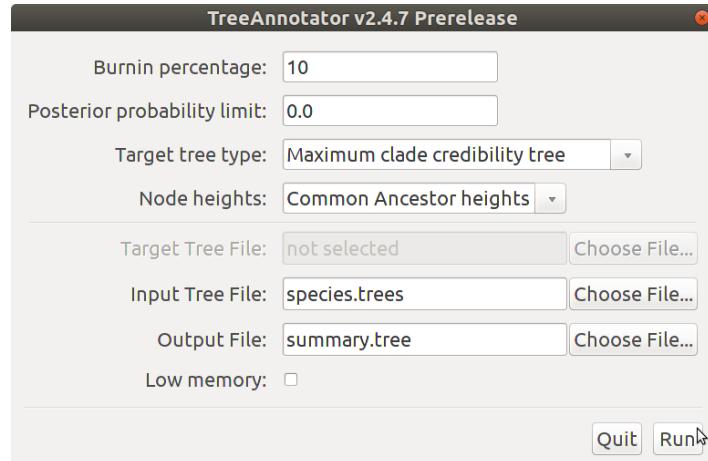


Figure 15: Running Tree Annotator

Now start FigTree and open the summary tree file. Enable Branch Labels, and expand the Branch Labels panel. Choose “dmv1_95%_HPD”, and set the number of significant digits to 2. Now the 95% highest posterior density (HPD) intervals of effective population sizes will be displayed on each branch (Figure 16).

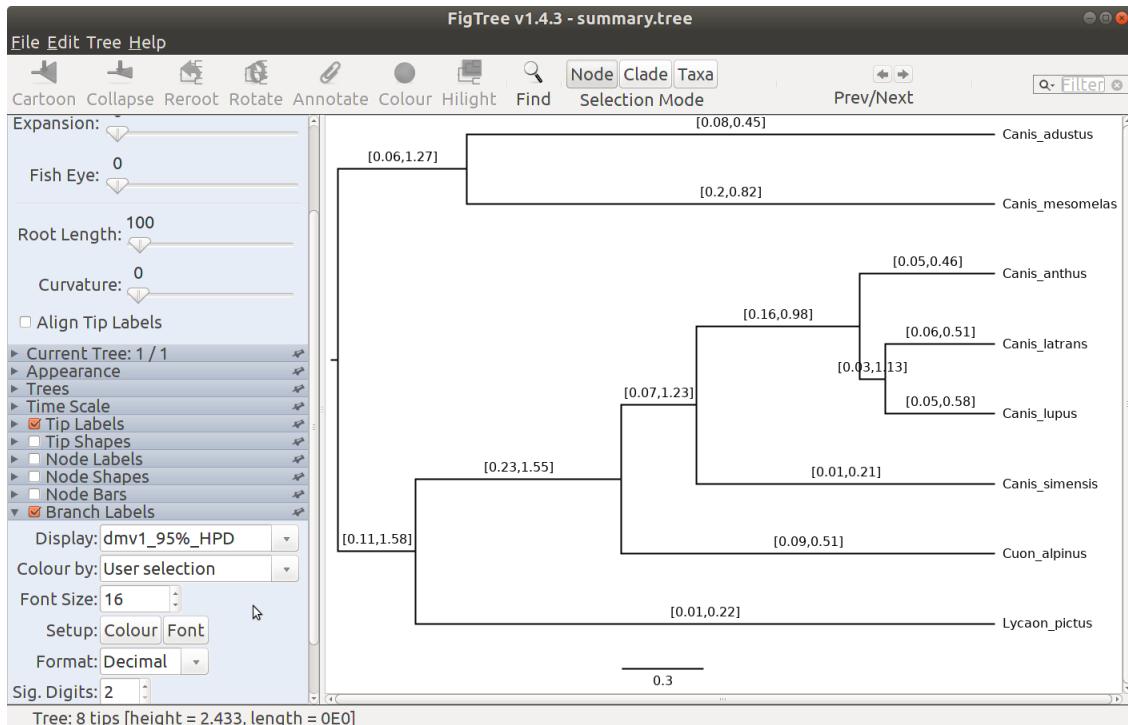


Figure 16: Running FigTree

In StarBEAST2 and many other BEAST packages, effective population sizes are scaled by generation time. That means if generation times (the average number of years from zygote to zygote) vary between species in your analysis, the effective population sizes of different branches cannot be directly compared.

If the average generation time is 5 years, and our tree is scaled in millions of years, then the generation time is 5×10^{-6} . To get effective population sizes in numbers of individuals, divide each value by the generation time. So if the scaled effective population size is 1, that will correspond to 200,000 individuals.

The effective population sizes inferred in this tutorial appear to have roughly an order of magnitude of uncertainty, and more precise estimates would require a more informative data set. This could be achieved by sampling more individuals or more loci.

4 Estimating per-species clock rates

We can estimate molecular clock rates separately for each species using a relaxed clock model, for example the uncorrelated log-normal (UCLN) model. First create a new folder for this analysis with a sensible name, something like “CanisUCLN”. Relaunch BEAUti, and select “SpeciesTreeUCLN” from the Templates submenu of the File menu (Figure 17).

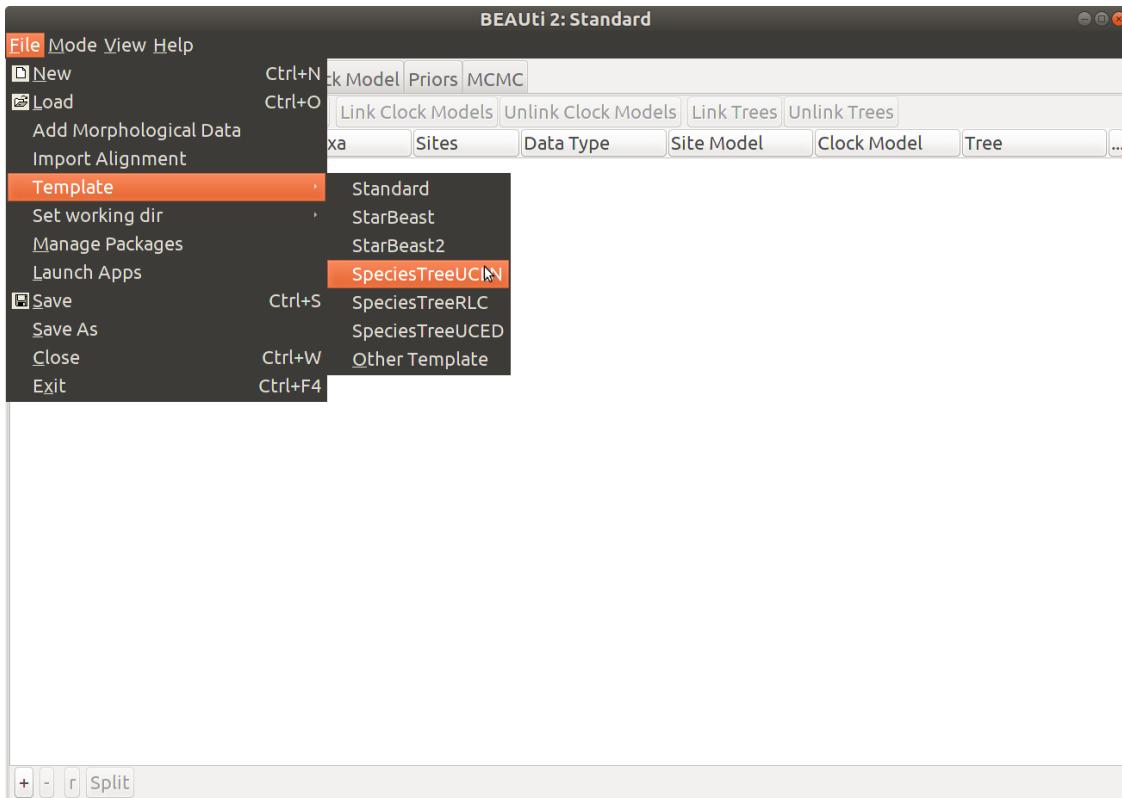


Figure 17: Starting a UCLN relaxed clock analysis

Import the same FASTA files as in subsection 3.1 (Figure 4). For any of the species tree relaxed clock templates, do **NOT** link the clock models. This will break the inference of per-species clock rates.

Assign the same species names as in subsection 3.3 (Figure 6, 7). This time we will keep the default setting for the population model (Analytical Integration), so that our analysis will run slightly faster. Set all the site models to HKY with empirical frequencies as in subsection 3.5 (Figure 9, 10).

Open the Clock Model tab, and manually set the Clock.rate to 0.001 (or equivalently 1e-3) for **every** locus (Figure 18). This is necessary because for technical reasons we cannot link the clock models for the species tree relaxed clock templates.

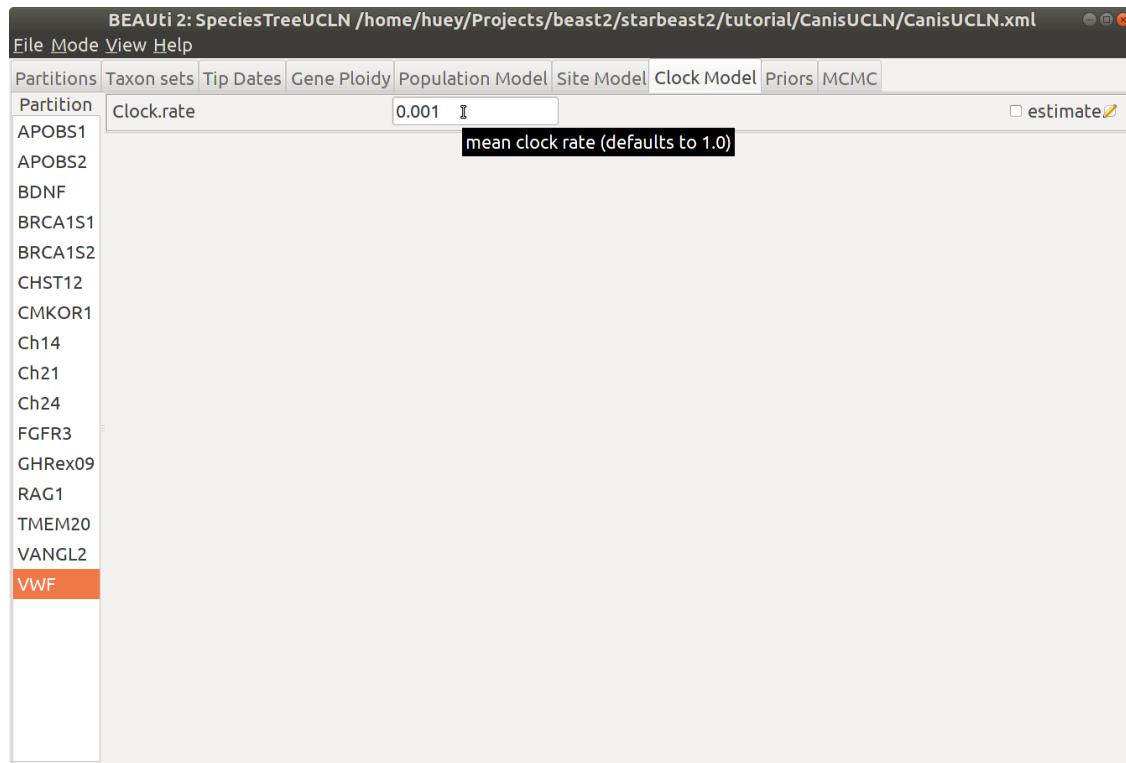


Figure 18: Manually setting every clock rate

Open the Priors tab and change the species tree prior from Yule to Birth-Death to allow for extinction. Finally, open the MCMC tab and change the length of the chain to 40 million, as in subsection 3.7 (Figure 12). Save the file in the folder you created previously for this analysis, and give it a sensible name like “CanisUCLN.xml”.

Run the MCMC chain by opening the XML file in BEAST, or running it from the command line as in subsection 3.8. This will take about twice as long (20 to 30 minutes) as the fixed clock analysis, because relaxed clocks are more computationally intensive.

After BEAST has finished, open the “starbeast.log” file in Tracer to check that the important statistics have ESS values of at least 200. Select the branchRatesStdev.Species parameter (Figure 19).

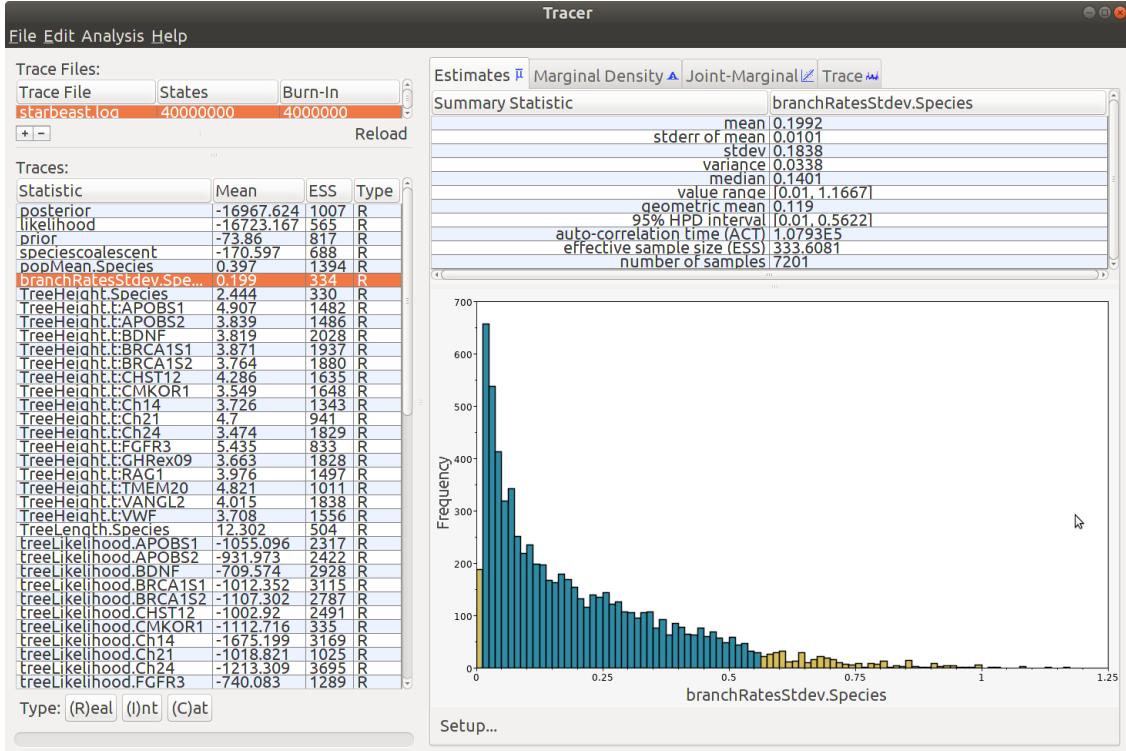


Figure 19: Checking the species tree relaxed clock analysis in Tracer

This parameter models the spread of molecular clock rates among branches in the species tree. The default prior for this parameter in StarBEAST2 is a log-normal distribution with a mean of 1, but the posterior distribution of this parameter has a mean of about 0.2. This means that the data is pulling this parameter lower, probably because there is very little variation in clock rates between species. Indeed the mode of the posterior distribution is about zero, suggesting that a strict clock may be more appropriate (Figure 19).

Now create a summary tree from the “species.trees” file using Tree Annotator, as in subsection 3.11 (Figure 15). Again give it a sensible name like “summary.tree”. Open up the summary tree file in FigTree. Enable Branch Labels, and expand the Branch Labels panel. Choose “rate_95%_HPD”, and set the number of significant digits to 2. Now the 95% HPD intervals of relative clock rates will be displayed on each branch (Figure 20).

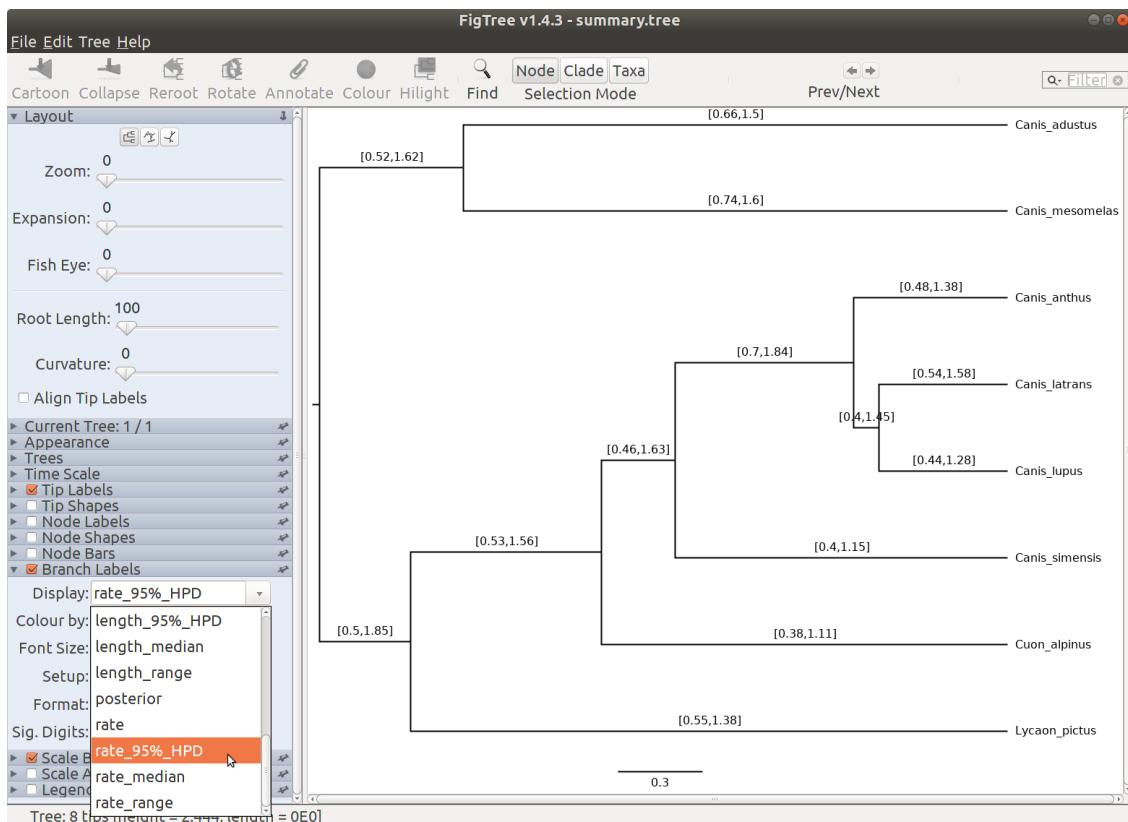


Figure 20: Checking the species tree relaxed clock analysis in FigTree

All of the species tree branches have clock rate HPD intervals that include 1, further suggesting that a strict clock may be more appropriate for this clade. These rates are obviously relative, and to get absolute rates in substitutions per site per million years, they must be scaled by the *a priori* clock rate of 0.001.

5 Total evidence tip-dating

The versions of StarBEAST2 from 14 onwards include the ability to combine tip-dating using the fossilized birth death (FBD) process with multispecies coalescent (MSC) inference of species and gene trees. We call this integrative model “FBD-MSC”. To begin your analysis, creat a new folder for this purpose with a sensible name, something like “CanisFBD”.

When setting up an FBD-MSC analysis in BEAUTi, it is necessary to import morphological data after specifying the taxon sets and the tree model, but before specifying the tip dates. The following steps will be consistent with that ordering. First, repeat section 3 from subsection 3.1 to and including 3.3. Leave the population model set at the default (Analytical Integration) to save time, then set all site models to HKY with empirical frequencies as in subsection 3.5.

Before adding the morphological data, go to the Priors panel and change the prior on the species tree from Yule to FBD (Figure 21). This is necessary because the inclusion of fossil taxa means our species tree will be *serially sampled* rather than *ultrametric*. StarBEAST2 gives you the option of various birth-death models – Yule, Calibrated Yule, Birth Death and FBD – but only FBD is valid for serially sampled trees.

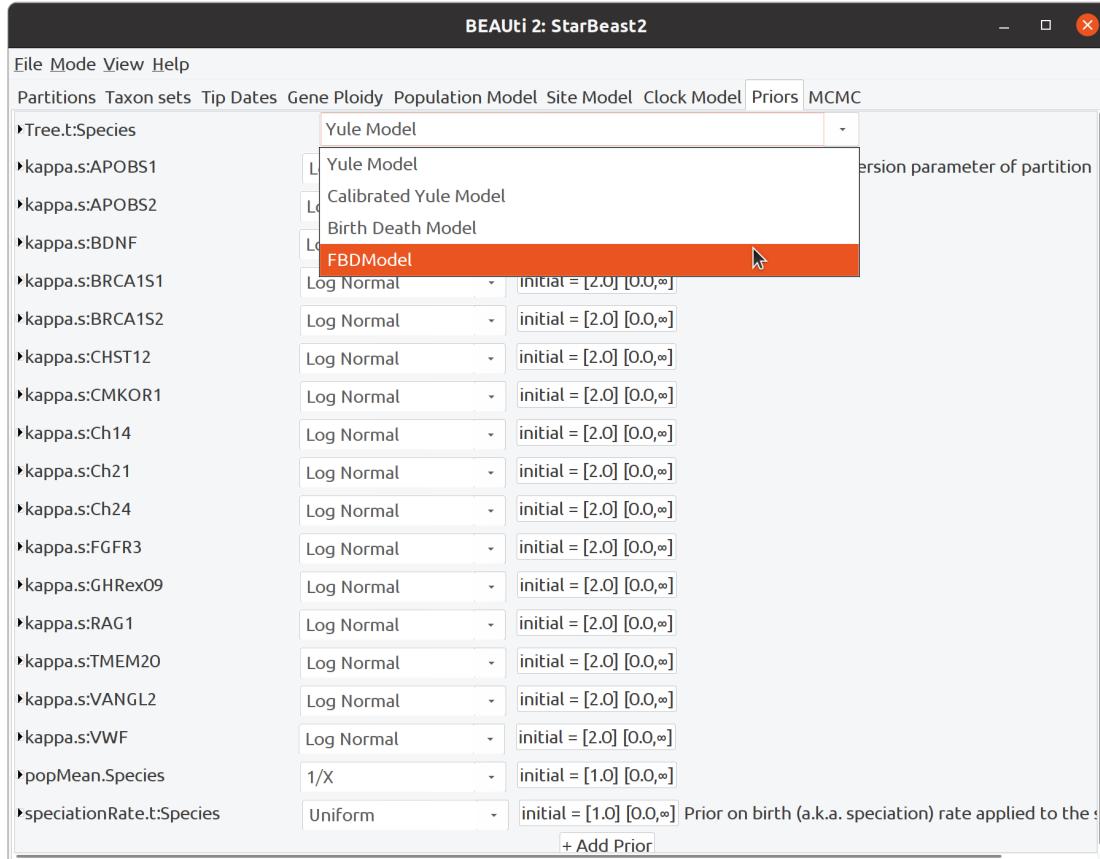


Figure 21: Change the tree model from Yule to FBD.

Now go back to the Partitions tab, and from the File menu select “Add Morphology to Species Tree” (Figure 22). Navigate to the data folder of the tutorial and import the “morphology_canis.nex” file.¹

Select Yes to condition on recording variable characters only (Mkv). Mkv is used when only characters that vary between species have been included in the character matrix, as is the case for this tutorial’s data set. You should now see four morphology partitions (Figure 23), one for each number of character states. For example, morphology_canis2 is for binary characters.

¹This file contains recorded states for 50 morphological characters of *Canis* species and some related species, taken from Slater 2015.

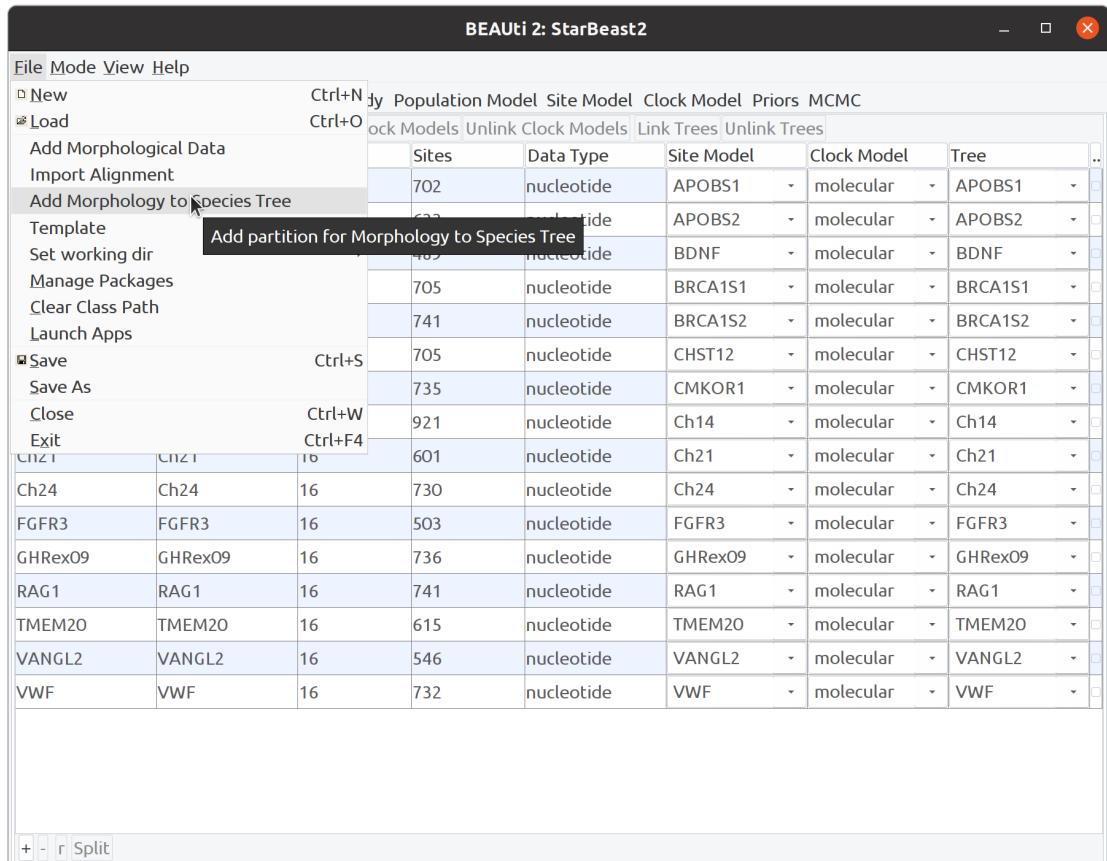


Figure 22: The option to add morphological data to the species tree

BEAUTi 2: StarBeast2

The screenshot shows the BEAUTi 2 software interface with the title "BEAUTi 2: StarBeast2" at the top. Below the title is a menu bar with "File", "Mode", "View", and "Help". The main area contains a table with the following columns: Name, File, Taxa, Sites, Data Type, Site Model, Clock Model, Tree, and a column for ellipsis (...). The table lists various partitions, each with its name, file name, number of taxa, number of sites, data type (nucleotide), site model, clock model, and tree model. The rows include APOBS1, APOBS2, BDNF, BRCA1S1, BRCA1S2, CHST12, CMKOR1, Ch14, Ch21, Ch24, FGFR3, GHREX09, RAG1, TMEM20, VANGL2, VWF, and morphology_canis1 through morphology_canis5. Most partitions have "molecular" as their clock model and tree model. Some partitions like APOBS1, APOBS2, and BDNF have "molecular" as their site model. The morphology partitions have "standard" as their data type. The tree column shows the tree models used for each partition.

Partitions								
	Taxon sets	Tip Dates	Gene Ploidy	Population Model	Site Model	Clock Model	Priors	MCMC
	Link Site Models	Unlink Site Models	Link Clock Models	Unlink Clock Models	Link Trees	Unlink Trees		
Name	File	Taxa	Sites	Data Type	Site Model	Clock Model	Tree	..
APOBS1	APOBS1	16	702	nucleotide	APOBS1	▼ molecular	APOBS1	▼
APOBS2	APOBS2	16	633	nucleotide	APOBS2	▼ molecular	APOBS2	▼
BDNF	BDNF	16	489	nucleotide	BDNF	▼ molecular	BDNF	▼
BRCA1S1	BRCA1S1	16	705	nucleotide	BRCA1S1	▼ molecular	BRCA1S1	▼
BRCA1S2	BRCA1S2	16	741	nucleotide	BRCA1S2	▼ molecular	BRCA1S2	▼
CHST12	CHST12	16	705	nucleotide	CHST12	▼ molecular	CHST12	▼
CMKOR1	CMKOR1	16	735	nucleotide	CMKOR1	▼ molecular	CMKOR1	▼
Ch14	Ch14	16	921	nucleotide	Ch14	▼ molecular	Ch14	▼
Ch21	Ch21	16	601	nucleotide	Ch21	▼ molecular	Ch21	▼
Ch24	Ch24	16	730	nucleotide	Ch24	▼ molecular	Ch24	▼
FGFR3	FGFR3	16	503	nucleotide	FGFR3	▼ molecular	FGFR3	▼
GHREX09	GHREX09	16	736	nucleotide	GHREX09	▼ molecular	GHREX09	▼
RAG1	RAG1	16	741	nucleotide	RAG1	▼ molecular	RAG1	▼
TMEM20	TMEM20	16	615	nucleotide	TMEM20	▼ molecular	TMEM20	▼
VANGL2	VANGL2	16	546	nucleotide	VANGL2	▼ molecular	VANGL2	▼
VWF	VWF	16	732	nucleotide	VWF	▼ molecular	VWF	▼
morphology_canis2	morphology_canis	24	29	standard	morphology_canis2	▼ morphology_canis	Species	▼
morphology_canis3	morphology_canis	24	16	standard	morphology_canis3	▼ morphology_canis	Species	▼
morphology_canis4	morphology_canis	24	11	standard	morphology_canis4	▼ morphology_canis	Species	▼
morphology_canis5	morphology_canis	24	8	standard	morphology_canis5	▼ morphology_canis	Species	▼

+ - r Split

Figure 23: After morphological data has been added to the species tree

Open the Tip Dates panel and enable “Use tip dates”. Change “Since some time in the past” to “before the present”. Click on Auto-configure and select “read from file”. Choose the “tip_dates_canis.txt” file in the data folder of the tutorial. Click OK, and your tip dates should look like Figure 24.

Name	Date	Height
Canis_mesomelas	0.000000	0.0
Canis_anthus	0.000000	0.0
Canis_lupus	0.000000	0.0
Canis_latrans	0.000000	0.0
Lycaon_pictus	0.000000	0.0
Cuon_alpinus	0.000000	0.0
Canis_simensis	0.000000	0.0
Canis_adustus	0.000000	0.0
Canis_antonii	1.200000	1.2
Canis_armbrusteri	0.300000	0.3
Canis_arnensis	1.500000	1.5
Canis_chihliensis	1.200000	1.2
Canis_dirus	0.012000	0.012
Canis_edwardii	0.300000	0.3
Canis_etruscus	1.700000	1.7
Canis_falconeri	1.500000	1.5
Canis_ferox	3.500000	3.5
Canis_lepophagus	1.800000	1.8
Canis_mosbachensis	0.500000	0.5
Canis_palmidens	1.500000	1.5
Canis_thoooides	1.800000	1.8
Canis_variabilis	0.500000	0.5
Eucyon_davisi	4.900000	4.9
Xenocyon_texanus	0.240000	0.24

Figure 24: Configured tip dates

Go to the Clock Model panel, and enable “estimate” for the molecular clock (by default called APOBS1 after the first locus, although this can be changed in the Partitions tab). Then select the morphological data partition and enable “estimate” so that the morphological clock rate will also be estimated (Figure 25). Make sure you have enabled “estimate” for **both** clocks, since we are using fossil data to calibrate our tree.

Next, open the Priors panel, scroll down to the strictClockRate.c:molecular parameter, and expand

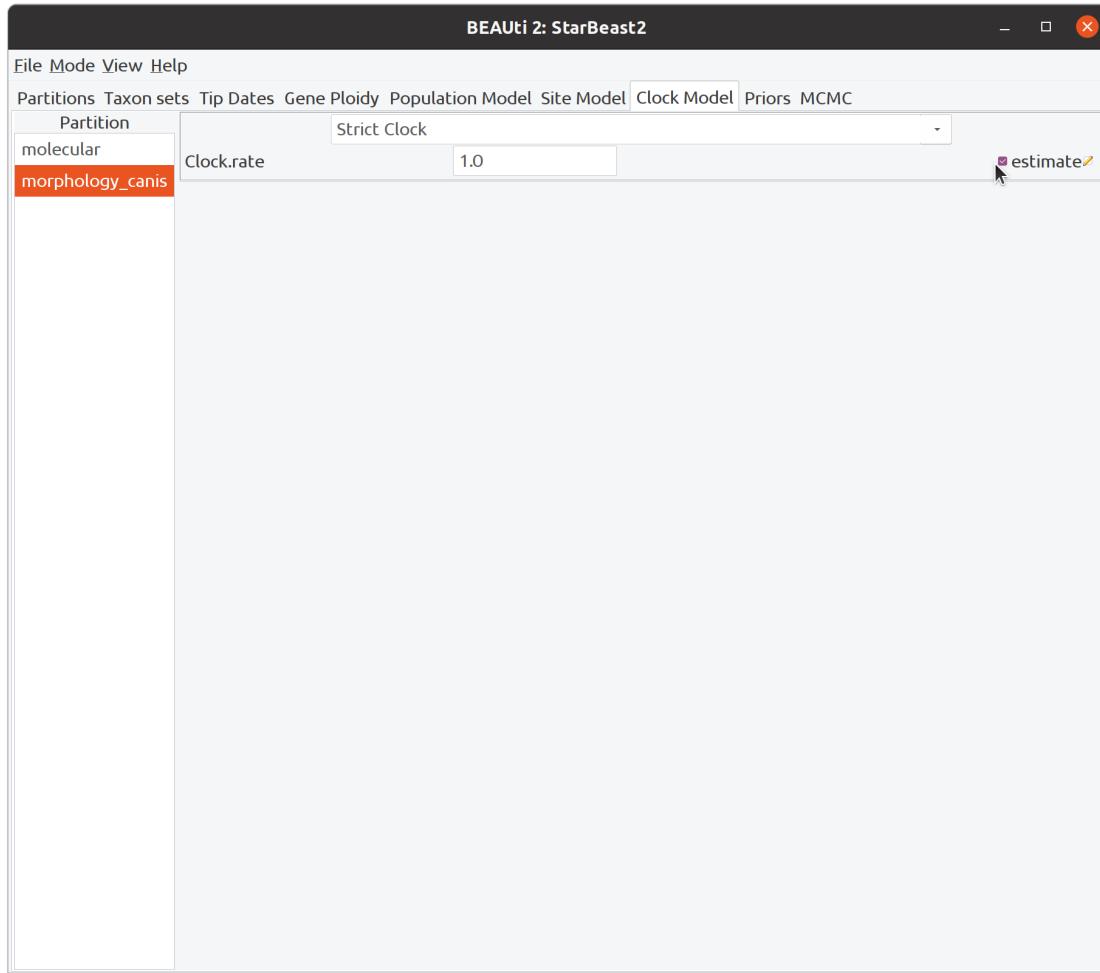


Figure 25: Estimating the clock rates

that prior. Change the mean to 0.001, the *a priori* molecular clock rate estimate from section 3. Leave the standard deviation set at 1, a moderately informative prior that will still allow the rate to be informed by the fossil data. Then change the prior distribution for the strictClockRate.c:morphology_canis prior to $1/X$ (Figure 26). This is an improper, uninformative prior so the inferred morphological clock rate will be estimated solely from the data.

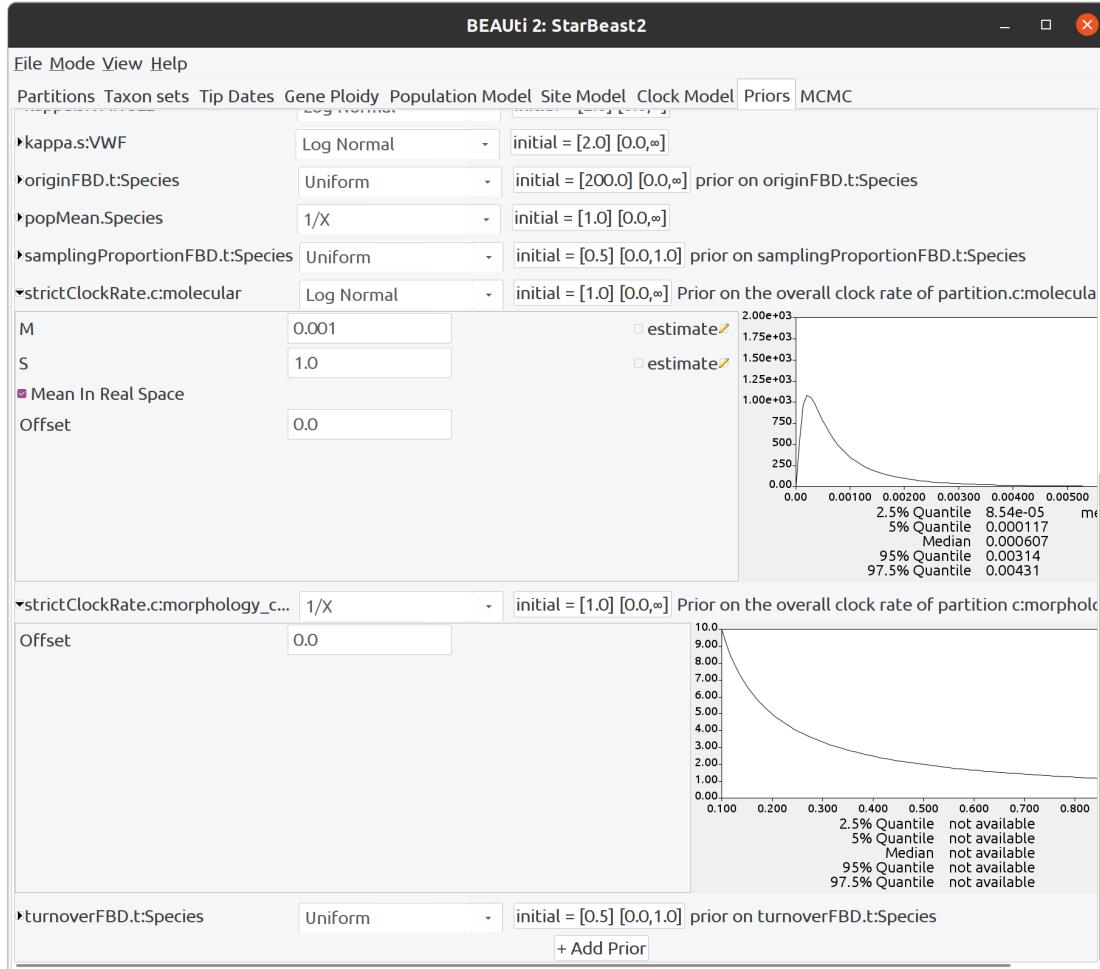


Figure 26: Configuring the clock prior distributions

The FBD-MSC model is quite complex and requires very long chains for reliable sampling. Open the MCMC tab and change the chain length to 200 million. We will need to reduce the sampling rate so that the output files remain managably small, so change the Store Every value to 50,000. Expand the tracelog, speciesTreeLogger, screenlog and every gene tree log, and set the Log Every value to 50,000 for each of them. This will limit the number of samples to 4,000 in each log file (Figure 27).

Save the XML file in the folder you created for this analysis, and name it something like “*CansFBD.xml*”. Run the MCMC chain by opening the XML file in BEAST, or running it from the command line as in subsection 3.8. The chain will take about 2 hours to finish, depending on the speed of your computer.

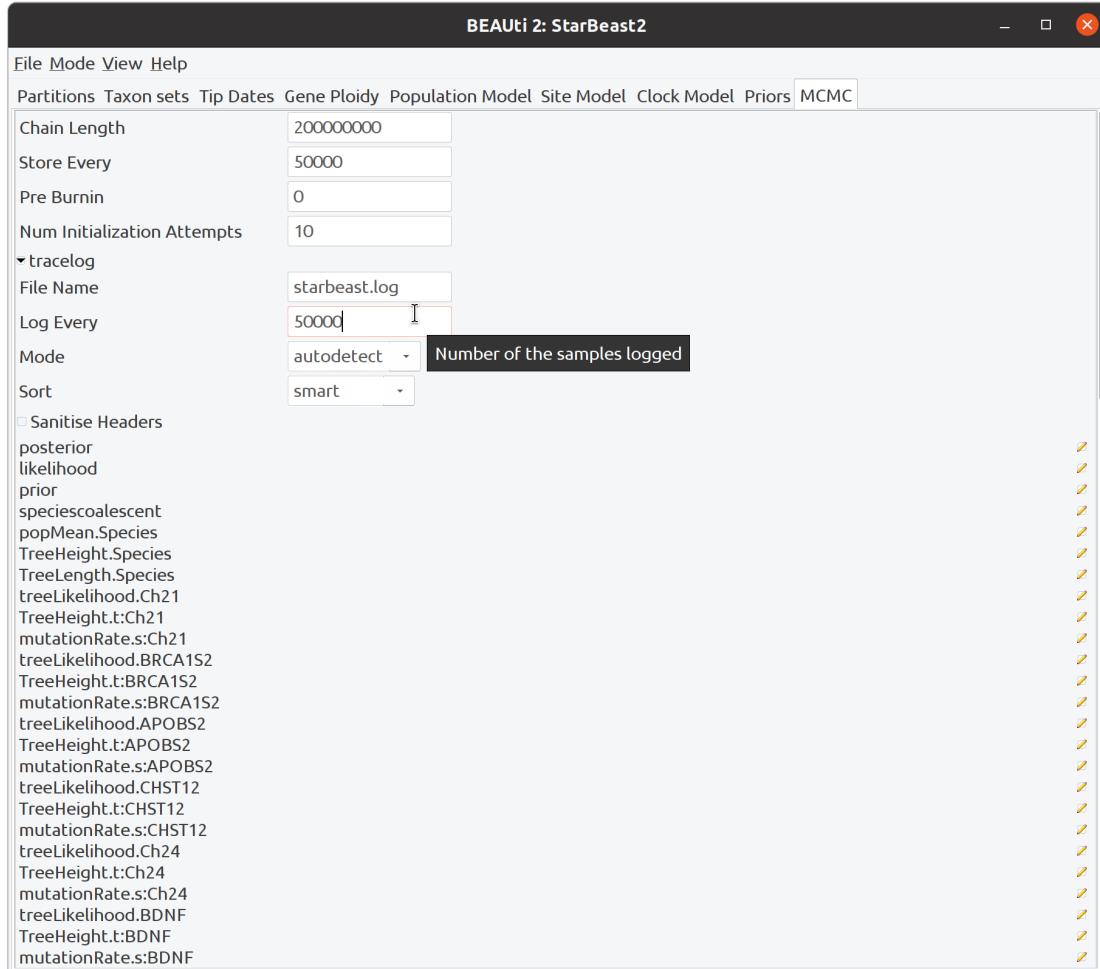


Figure 27: Setting a very long chain length

After BEAST has finished, open the “starbeast.log” file in Tracer to check on the ESS values. Select the strictClockRate.c:molecular parameter (Figure 28). The 95% HPD interval for this parameter spans approximately 0.0005 to 0.0013 substitutions per site per million years, in agreement with the *a priori* rate of 0.001.

Start the DensiTree app again, and then open the “species.trees” file. Under the Show panel, enable the Root Canal tree to get an idea of the most plausible species tree. Then open the Grid panel, and enable the full grid. Compared with the fixed clock analysis (Figure 14), the divergence time between *Canis lupus* and *Canis latrans* is a little older, and the age of the MRCA of extant species is a little younger (Figure 29).

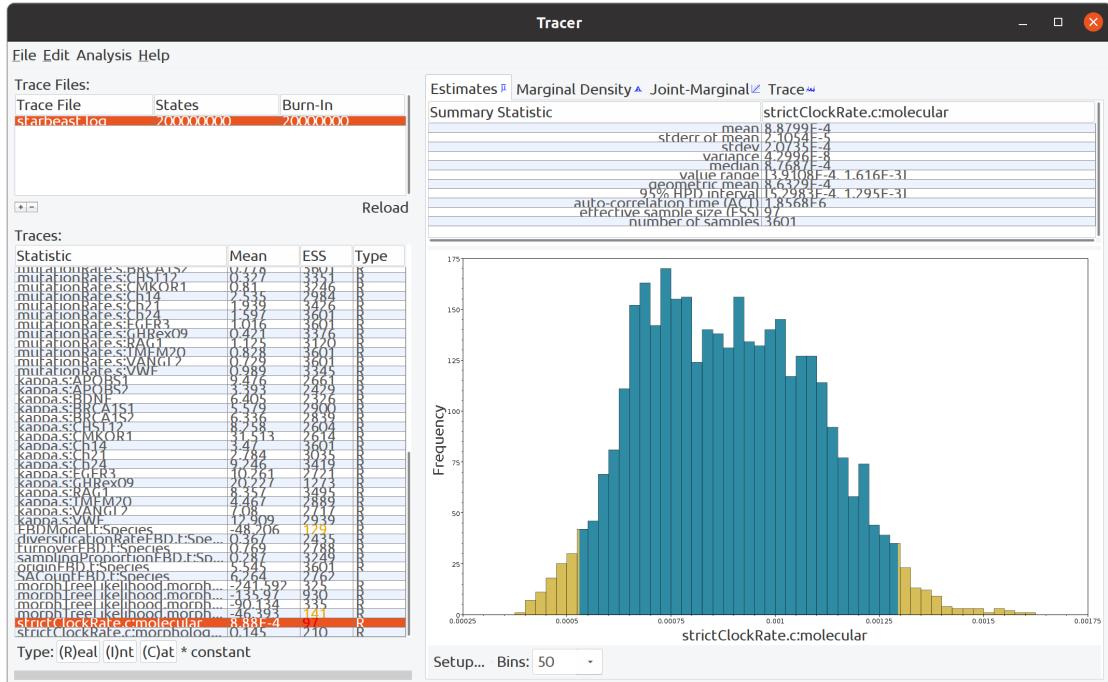


Figure 28: The posterior estimate of the molecular clock rate

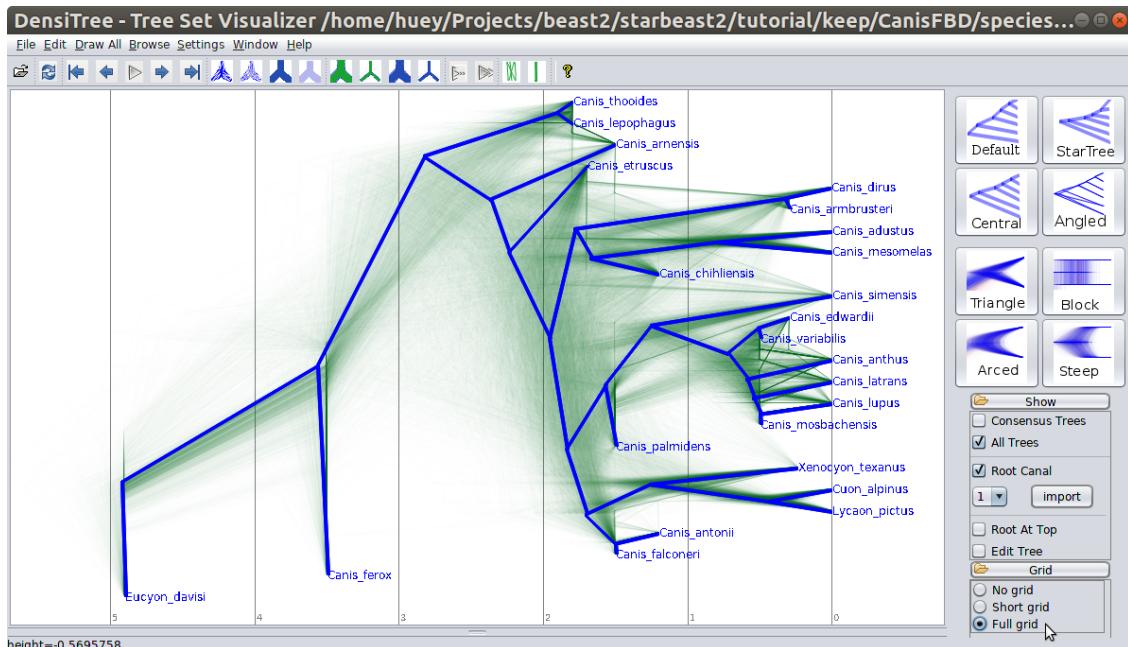


Figure 29: Cloud of FBD-MSC species trees

Open the Tree Annotator app included with BEAST2, and set the burnin percentage to 10. Trees inferred by the FBD-MSM model as implemented in StarBEAST2 can include sampled ancestors (Gavryushkina et al. 2014), which are incompatible with Common Ancestor heights (Heled and Bouckaert 2013) as implemented in Tree Annotator. So in order to generate a summary tree, change “Node heights” to “Mean heights” (Figure 30). Choose the “species.trees” file as the input file, and specify the output file to be something sensible like “summary.tree”, then click Run.

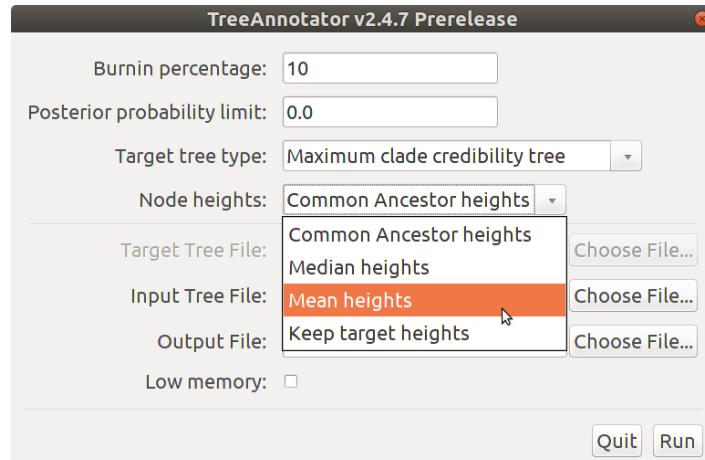


Figure 30: Running Tree Annotator for sampled ancestor trees

Start FigTree and open the summary tree file you have just created. Enable Node Labels, and expand the Node Labels panel. Choose “posterior”, and set the number of significant digits to 2. Now the posterior support for each clade will be displayed next to each node (Figure 31).

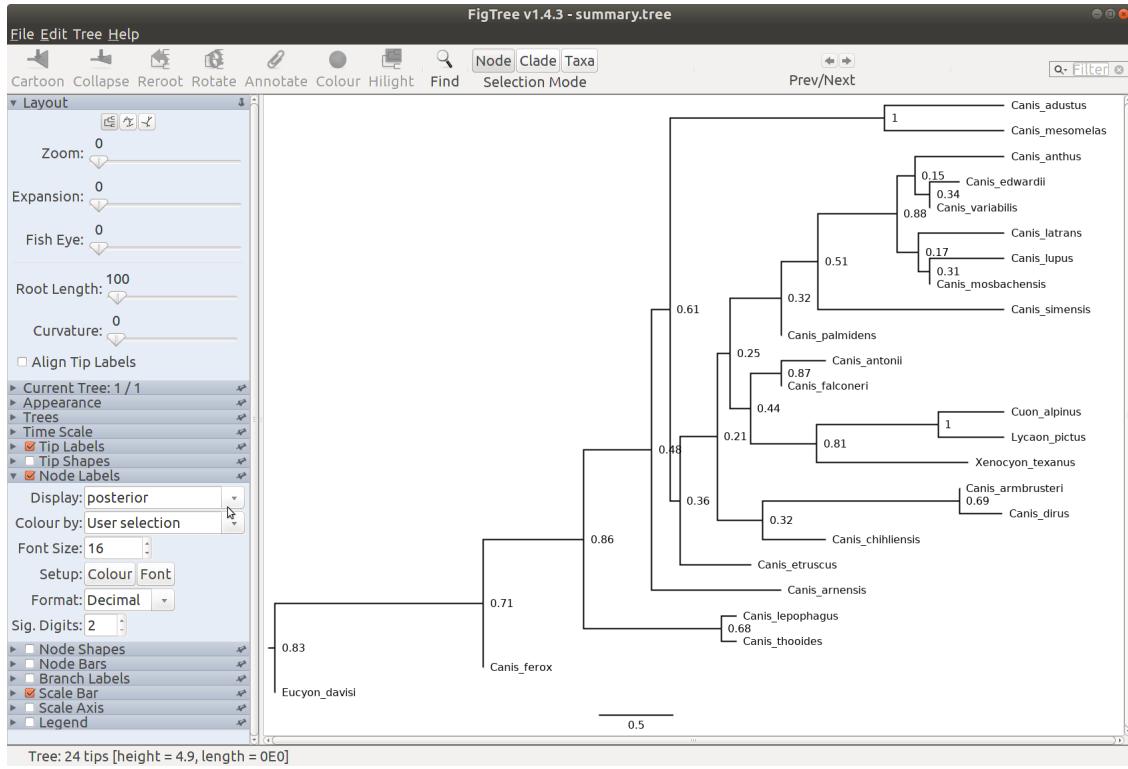


Figure 31: Showing posterior support for clades in FigTree

The tree topology is similar to the fixed clock analysis (Figure 16), except that *Cuon alpinus* and *Lycaon pictus* are now a clade with 100% support. A disagreement between molecular data which contradicts this clade and morphological data which supports it has been reported previously (Zrzavý and Řičánková 2004). This shows how morphological characters and fossil data can influence species tree estimates.

6 Wrapping up

That concludes the tutorial. One final note; when running any MCMC method, it is possible than the chain has gotten stuck in a region of parameter space with high local probability. So for any kind of analysis intended for publication, it is a good idea to run at least two MCMC chains. Open both chains in Tracer, and check that the traces for the logged statistics and parameters have the same distributions.

References

- Bouckaert, Remco R. (2010). “DensiTree: making sense of sets of phylogenetic trees”. In: *Bioinformatics* 26.10, pp. 1372–1373.
- Gavryushkina, Alexandra, David Welch, Tanja Stadler, and Alexei J. Drummond (2014). “Bayesian inference of sampled ancestor trees for epidemiology and fossil calibration”. In: *PLoS Computational Biology* 10.12, e1003919.
- Hasegawa, Masami, Hirohisa Kishino, and Takaaki Yano (Oct. 1985). “Dating of the human-ape splitting by a molecular clock of mitochondrial DNA”. In: *Journal of Molecular Evolution* 22.2, pp. 160–174. ISSN: 1432-1432.
- Heled, Joseph and Remco R. Bouckaert (Oct. 2013). “Looking for trees in the forest: summary tree from posterior samples”. In: *BMC Evolutionary Biology* 13.1, p. 221. ISSN: 1471-2148.
- Hugall, Andrew F., Ralph Foster, Michael S. Y. Lee, and Marshal Hedin (2007). “Calibration Choice, Rate Smoothing, and the Pattern of Tetrapod Diversification According to the Long Nuclear Gene RAG-1”. In: *Systematic Biology* 56.4, pp. 543–563.
- Lindblad-Toh, Kerstin et al. (2005). “Genome sequence, comparative analysis and haplotype structure of the domestic dog”. In: *Nature* 438.7069, pp. 803–819.

- Ogilvie, Huw A., Remco R. Bouckaert, and Alexei J. Drummond (2017). “StarBEAST2 Brings Faster Species Tree Inference and Accurate Estimates of Substitution Rates”. In: *Molecular Biology and Evolution* 34.8, pp. 2101–2114.
- Ogilvie, Huw A., Joseph Heled, Dong Xie, and Alexei J. Drummond (2016). “Computational Performance and Statistical Accuracy of *BEAST and Comparisons with Other Methods”. In: *Systematic Biology* 65.3, pp. 381–396.
- Prado-Martinez, Javier et al. (July 2013). “Great ape genetic diversity and population history”. In: *Nature* 499, pp. 471–475.
- Scornavacca, Celine and Nicolas Galtier (2017). “Incomplete Lineage Sorting in Mammalian Phylogenomics”. In: *Systematic Biology* 66.1, pp. 112–120.
- Slater, Graham J. (2015). “Iterative adaptive radiations of fossil canids show no evidence for diversity-dependent trait evolution”. In: *Proceedings of the National Academy of Sciences* 112.16, pp. 4897–4902.
- Yang, Ziheng (Sept. 1994). “Maximum likelihood phylogenetic estimation from DNA sequences with variable rates over sites: Approximate methods”. In: *Journal of Molecular Evolution* 39.3, pp. 306–314. ISSN: 1432-1432.
- Zrzavý, Jan and Věra Řičáková (2004). “Phylogeny of Recent Canidae (Mammalia, Carnivora): relative reliability and utility of morphological and molecular datasets”. In: *Zoologica Scripta* 33.4, pp. 311–333. ISSN: 1463-6409.