

Phylogeographic inference in discrete space

A hands-on practical

This chapter provides a step-by-step tutorial on reconstructing the spatial dispersal and cross-species dynamics of rabies virus (RABV) in North American bat populations based on a set of 372 nucleoprotein gene sequences (nucleotide positions: 594–1353). The data set comprises a total of 17 bat species sampled between 1997 and 2006 across 14 states in the United States (Streicker *et al.*, Science, 2010, 329, 676-679). Following Faria *et al.* (Phil. Trans. Roy. Soc. B, 2013), two additional species that had been excluded from the original analysis owing to a limited amount of available sequences, *Myotis austroriparius* (Ma) and *Parastrellus hesperus* (Ph), are also included here. We also include a viral sequence with an unknown sampling date (accession no. TX5275, sampled in Texas from *Lasiurus borealis*), which will be adequately accommodated in our inference.

The aim of this tutorial is to estimate the ancestral locations of the virus using a Bayesian discrete phylogeographic approach and, at the same time, infer the history of host jumping using the same model approach. Using an extension of the discrete diffusion model, we will then test the factors that underly the host transition dynamics.

The first step will be to convert a NEXUS file with a DATA or CHARACTERS block into a BEAST XML input file. This is done using the program BEAUti (this stands for Bayesian Evolutionary Analysis Utility). This is a user-friendly program for setting the evolutionary model and options for the MCMC analysis. The second step is to actually run BEAST using the input file that contains the data, model and settings. The final step is to explore the output of BEAST in order to diagnose problems and to summarize the results.

To undertake this tutorial, you will need to download the following software packages in a format that is compatible with your computer system (all are available for Mac OS X, Windows and Linux/UNIX operating systems):

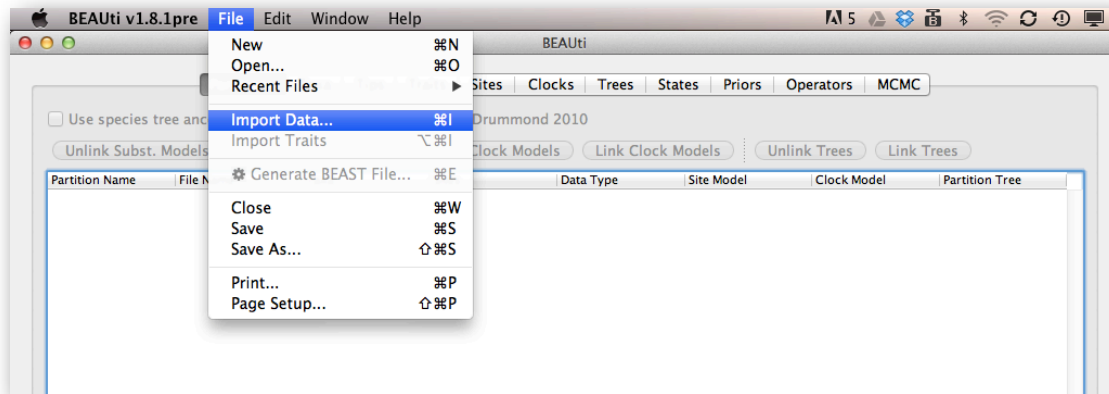
- **BEAST** - this package contains the *BEAST* program, *BEAUti* and a couple of utility programs. At the time of writing, the current version is v1.8.1. It is available for download from <http://beast.bio.ed.ac.uk> (or <http://beast-mcmc.googlecode.com>).
- **BEAGLE** - this is a high-performance library that can perform the core calculations at the heart of most Bayesian and Maximum Likelihood phylogenetics packages. It can make use of highly-parallel processors such as those in graphics cards (GPUs) found in many PCs. Binary installers and installation instructions can be found at: <http://beagle-lib.googlecode.com/>.
- **Tracer** - this program is used to explore the output of *BEAST* (and other Bayesian MCMC programs). It graphically and quantitatively summarizes the distributions of continuous parameters and provides diagnostic information. At the time of writing, the current version is v1.6. It is available for download from <http://tree.bio.ed.ac.uk/software/tracer/>.
- **FigTree** - this is an application for displaying and printing molecular phylogenies, in particular those obtained using *BEAST*. At the time of writing, the current version is v1.4.1. It is available for download from <http://tree.bio.ed.ac.uk/software/figtree>
- **SPREAD** - this is an application for the visualization of phylogeographic analyses performed with *BEAST*. At the time of writing, the current version is v1.0.6. It is available for download from <http://www.phylogeography.org/SPREAD>.
- **Google Earth** - this is a freely available virtual globe software that can be used to visualize KML output from SPREAD in an interactive fashion. **Google Earth** is available at <http://earth.google.com>.

Running BEAUti

The program **BEAUti** is a user-friendly program for setting the model parameters for BEAST. Run BEAUti by double clicking on its icon.

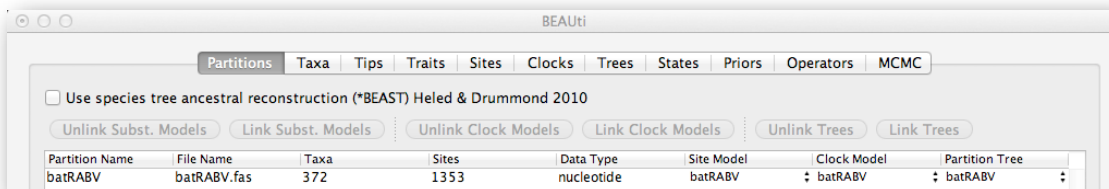
Loading the NEXUS file

To load a NEXUS or FASTA format alignment, simply select the **Import Data...** option from the **File** menu.



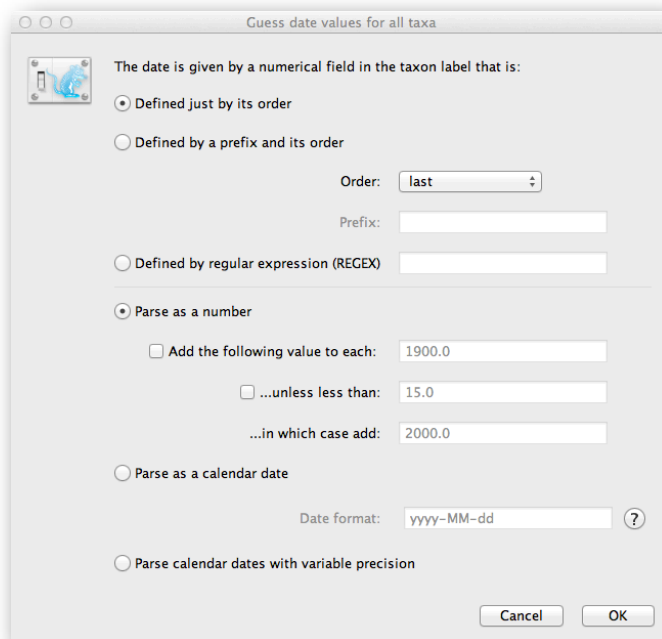
The NEXUS alignment

Select the file called **batRABV.fas**. This file contains an alignment of 372 nucleoprotein gene sequences of bat rabies viruses, 1353 nucleotides in length. Once loaded, the sequence data will be listed under **Partitions**:



Specifying the sampling date information

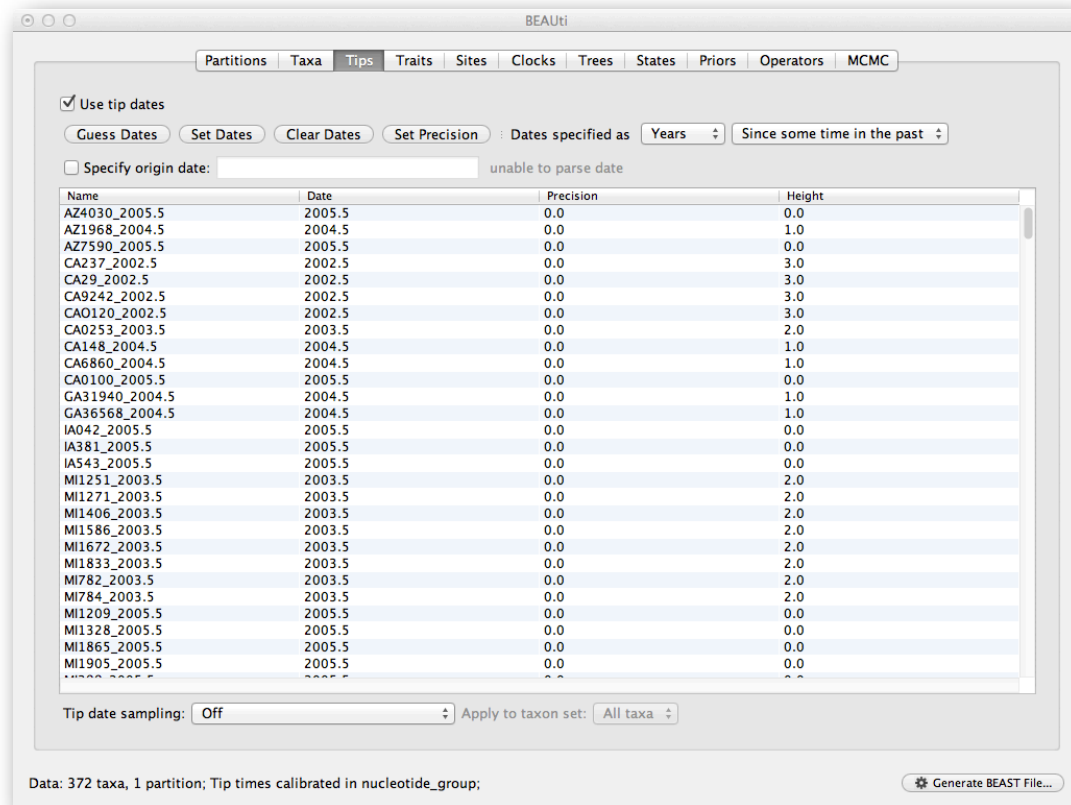
To inform BEAUti/BEAST about the sampling dates of the sequences, go to the **Tips** tab and select the **Use tip dates** option. By default all the taxa are assumed to have a date of zero (i.e. the sequences are assumed to be sampled at the same time). In this case, the RABV sequences have been sampled at various dates going back to 1997. The actual year of sampling is given in the name of each taxon and we could simply edit the value in the Date column of the table to reflect these. However, if the taxa names contain the calibration information, then a convenient way to specify the dates of the sequences in BEAUti is to use the **Guess Dates** button at the top of the Data tab. Clicking this will make a dialog box appear:



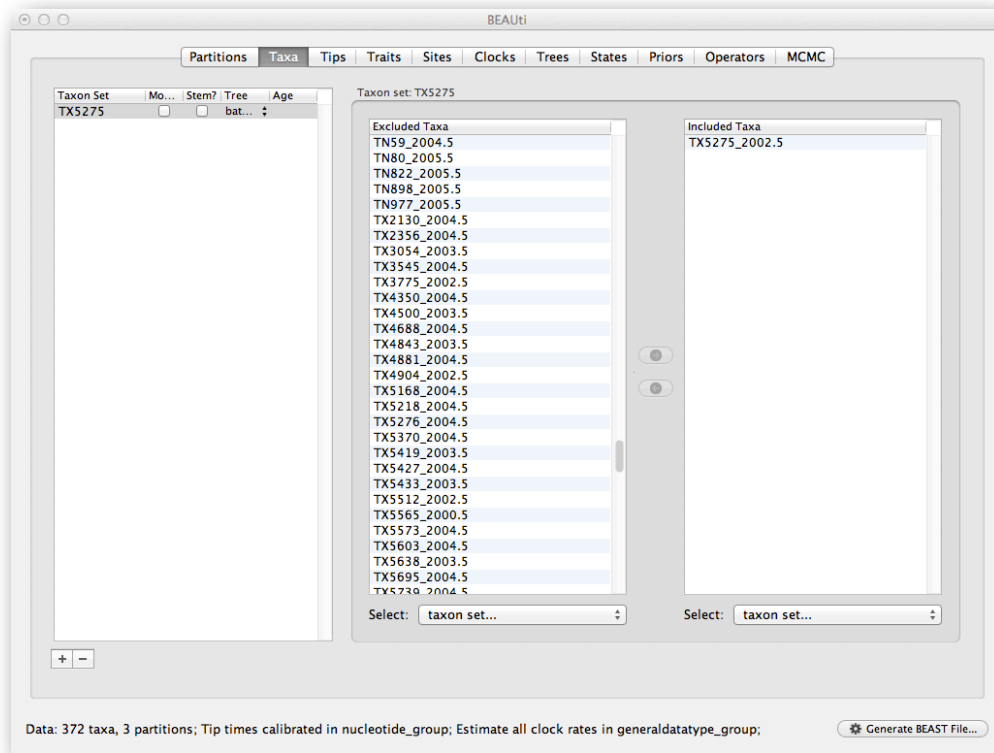
This operation attempts to guess what the dates are from information contained within the taxon names. It works by trying to find a numerical field within each name. If the taxon names contain more than one numerical field then you can specify how to find the one that corresponds to the date of sampling. You can (1) specify the order that the date field comes (e.g., first, last or various positions in between) or (2) specify a prefix (some characters that come immediately before the date field in each name) and the order of the field, or (3) define a regular expression (REGEX).

When parsing a number, you can ask BEAUTi to add a fixed value to each guessed date. For example, the value "1900" can be added to turn the dates from 2 digit years to 4 digit. Any dates in the taxon names given as "00" would thus become "1900". However, if these "00" or "01", etc. represent sequences sampled in 2000, 2001, etc., "2000" needs to be added to those. This can be achieved by selecting the "unless less than: .." and "...in which case add:.." option adding for example 2000 to any date less than 10. These operations are not necessary in our case since the dates are fully specified at the end of the sequence names. There is also an option to parse calendar dates and one for calendar dates with various precisions. For the H1N1/09 sequences you can keep the default **'Defined just by its order'** and select **'last'** from the drop-down menu for the order and press **'OK'**. The dates will appear in the appropriate column of the main window. You can then check these and edit them manually as required. At the top of the window you can set the units that the dates are given in (years, months, days) and whether they are specified relative to a point in the past (as would be the case for years such as 2005) or backwards in time from the present (as in the case of radiocarbon ages). In addition, the

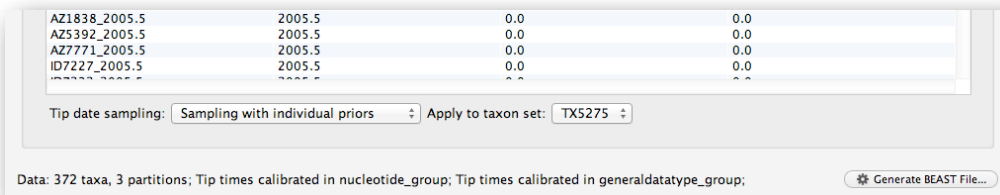
The "Height" column lists the ages of the tips relative to time 0 (in our case 2005.5). The Precision column allows specifying with what precision the sampling time is known. To include taxa only known up to the year of sampling for example (e.g., 2005), a precision of 1 year can be set and the age of those tips can be integrated over the time interval of 1 year using the Tip date sampling option at the bottom left of the Tips panel.



In our data set, the sampling date is unknown for one particular sequence (TX5275_2002.5, the '2002.5' is simply an arbitrary date that will be used as a starting value). To appropriately accommodate the uncertainty on the age of this tip, we will instruct *BEAST* to integrate over a particular sampling time interval for this tip. First, go back to the **Taxa** tab that we skipped, and make a taxon set for only that particular sequence. Press the small "plus" button at the bottom left of the panel ; this creates a new taxon set. Rename it by double-clicking on the entry that appears (it will initially be called **untitled1**). Call it TX5275 and keep the default settings. Move TX5275_2002.5 from the **Excluded Taxa** window to the **Included Taxa** window:

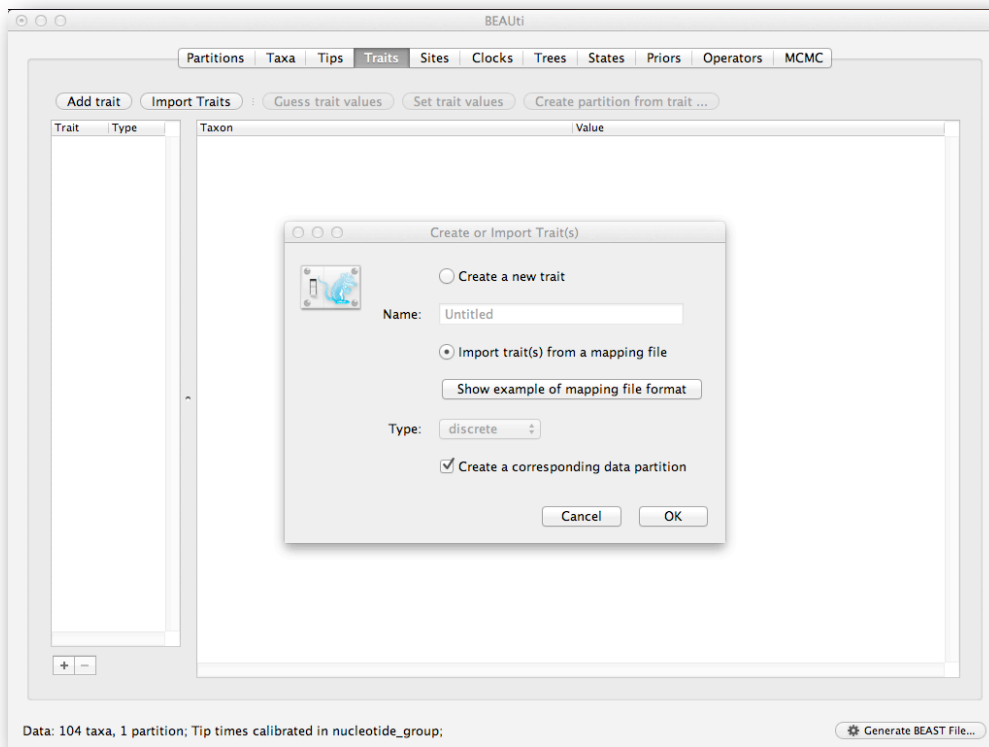


Go back to the **Tips** tab, and in the bottom left, select the **sampling with individual priors** as **Tip date sampling** option. Apply this to the **TX5275** taxa set instead of the default **All taxa** option. We will set a prior on its age when we get to the **Priors** tab.



Specifying the trait information for the sequences

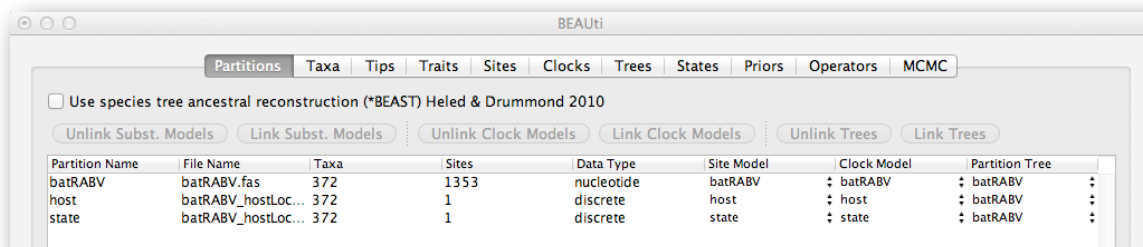
The next thing to do is to click on the **Traits** tab at the top of the main window. A trait can be any characteristic that is inherent to the specific taxon, for example, geographical location or species host. This step will assign a specific host and geographical location to each taxa. To associate the sequences with these traits, we need to add a new trait under the **Traits** tab (click **Add trait**). This will open a new window to **Create or Import Trait(s)**:



Select **Import trait(s) from a mapping file format** (the format of such a file can be shown). Browse to and load the **batRABV_hostLocation.txt** tab-delimited file which contains the discrete host and location for each sequence. Note that the host species is specified using a two-character abbreviation (e.g. Ef for *Eptesicus fuscus*, three characters for Lbl).

```
traits host state
AZ4030_2005.5Ap Arizona
AZ1968_2004.5Ef Arizona
AZ7590_2005.5Ef Arizona
CA237_2002.5 Ef California
CA29_2002.5 Ef California
CA9242_2002.5 Ef California
CAO120_2002.5 Ef California
CA0253_2003.5 Ef California
CA148_2004.5 Ef California
CA6860_2004.5 Ef California
CA0100_2005.5 Ef California
GA31940_2004.5 Ef Georgia
....
TX3545_2004.5Tb Texas
```

After clicking **OK**, select the host trait and click on **create partition from trait..!**. This new partition will be shown under the **Partitions** tab. Do the same for the location trait (state), resulting in three partitions in the **Partitions** tab:

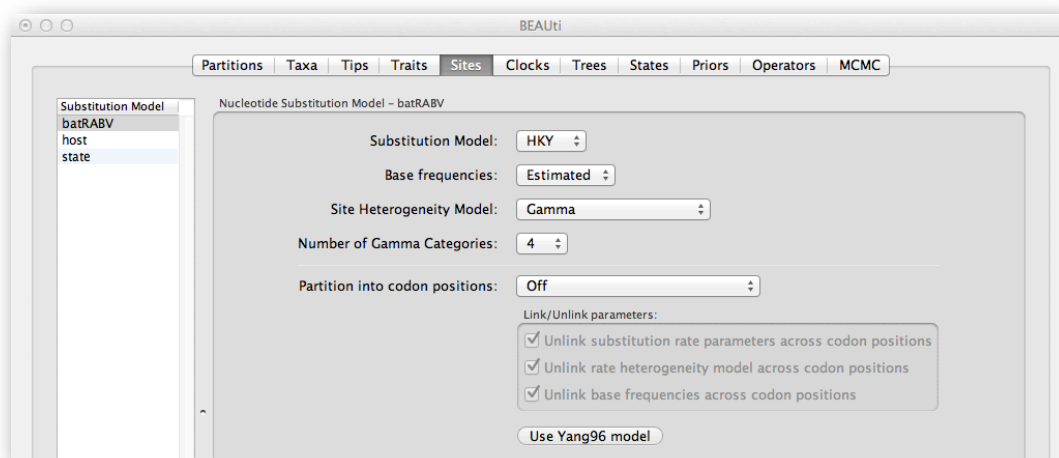


Setting the evolutionary and diffusion models

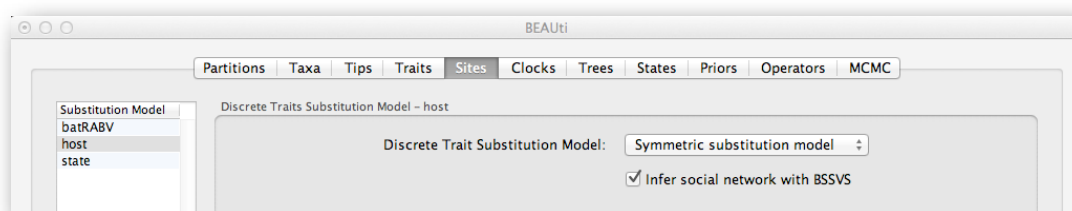
The next thing to do is to click on the **Sites** tab at the top of the main window. This will reveal the evolutionary model settings for BEAST. Exactly which options appear depend on whether the data are nucleotides, amino acids or traits. This tutorial assumes that you are familiar with the evolutionary models available; however there are a couple of points to note about selecting a model in **BEAUti**:

- Selecting the **Partition into codon positions** option assumes that the data are aligned as codons. This option will then estimate a separate rate of substitution for each codon position, or for 1+2 versus 3, depending on the setting.
- Selecting the **Unlink substitution model across codon positions** will specify that BEAST should estimate a separate transition-transversion ratio or general time reversible rate matrix for each codon position.
- Selecting the **Unlink rate heterogeneity model across codon positions** will specify that BEAST should estimate set of rate heterogeneity parameters (gamma shape parameter and/or proportion of invariant sites) for each codon position.

For the nucleotide model in this tutorial, keep the default **HKY** substitution model, set base frequencies to **Empirical**, and use **Gamma**-distributed rate variation among sites (with 4 discrete categories):



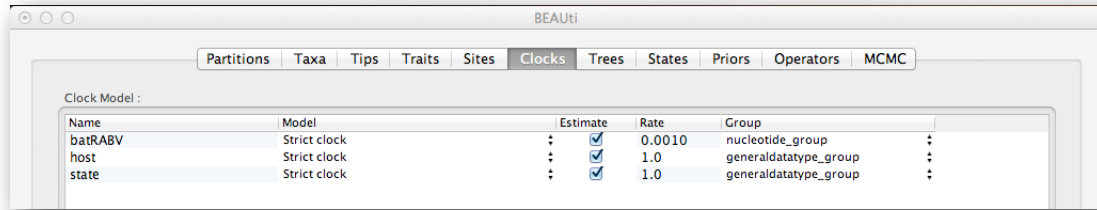
Click on **host** in the **Substitution model** window and keep the **Discrete Trait Substitution Model** to **Symmetric substitution model** and select the option to perform **BSSVS (Infer social network with BSSVS)**. The **Symmetric substitution model** specifies a discrete phylogeographic analysis using a standard continuous-time Markov chain (CTMC), in which the transition rates between locations are reversible. The alternative **Asymmetric substitution model** specifies a discrete phylogeographic analysis using a nonreversible CTMC. Selecting the **BSSVS** option enables the Bayesian Stochastic Search Variable Selection procedure. This procedure will attempt to invoke a limited number of rates (at least $k-1$) to adequately explain the phylogenetic diffusion process.



Apply the same discrete diffusion model settings to the spatial 'state' trait.

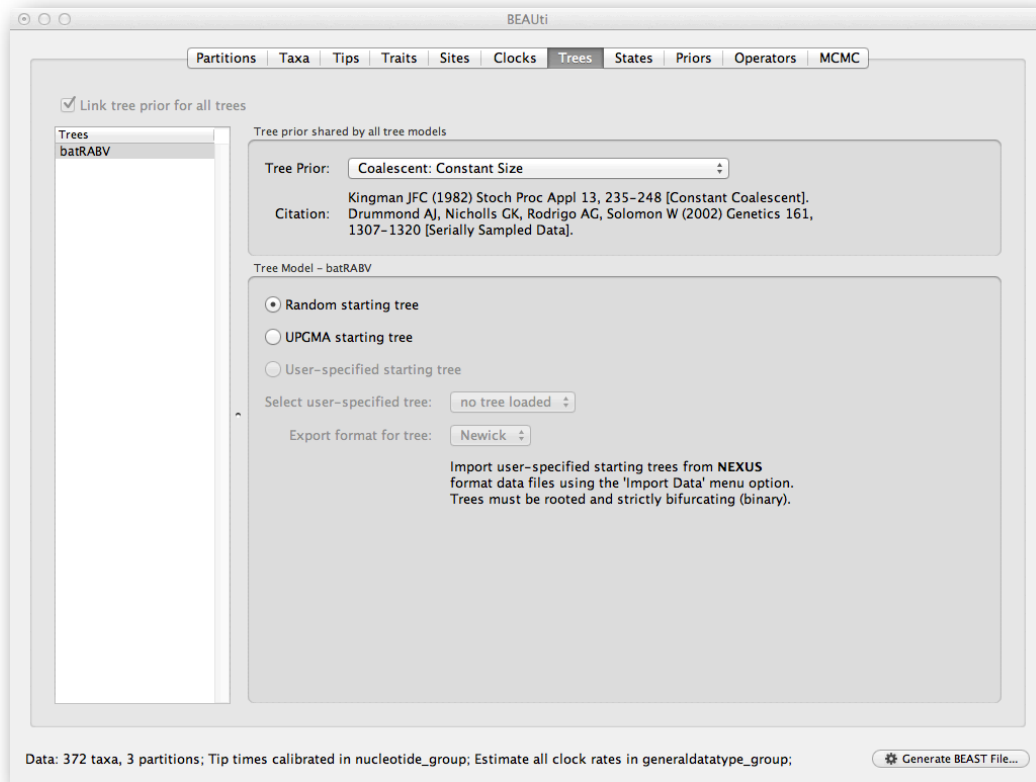
Setting the clock model

Click on the **Clocks** tab at the top of the main window. We will perform our run using the default **Strict clock** model and set the initial value for the **Rate** to 0.001. We can keep default settings for overall rate scalers in the host and location state transition processes.

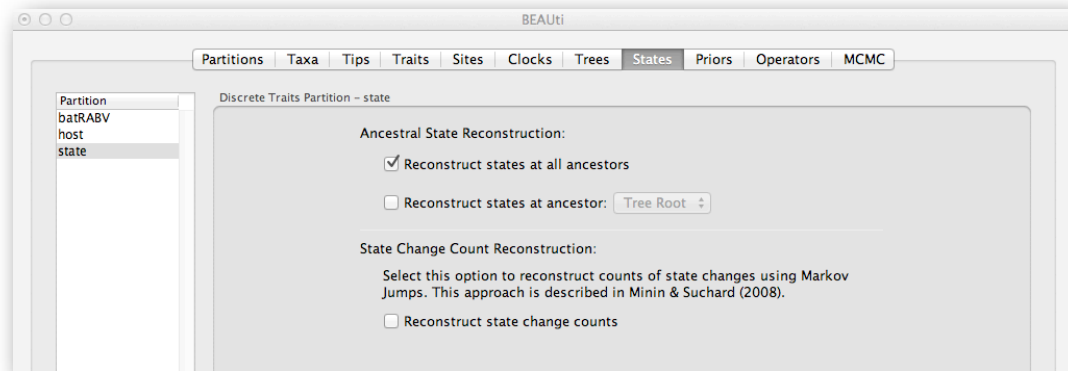


Setting the starting tree and tree prior

Click on the **Trees** tab at the top of the main window. We will select a simple constant size coalescent tree prior (**Coalescent: Constant Size**) and keep the default random starting tree.

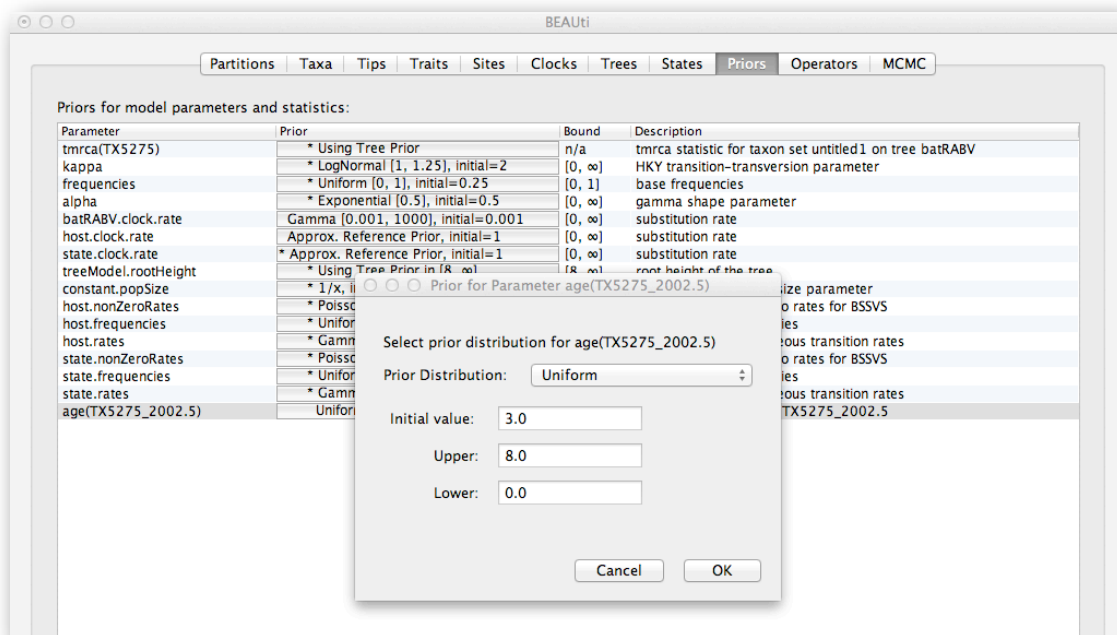


In the **States** tab, check that for the **host** and **state** partition the option to **Reconstruct states at all ancestors** is selected (by default).



Setting up the priors

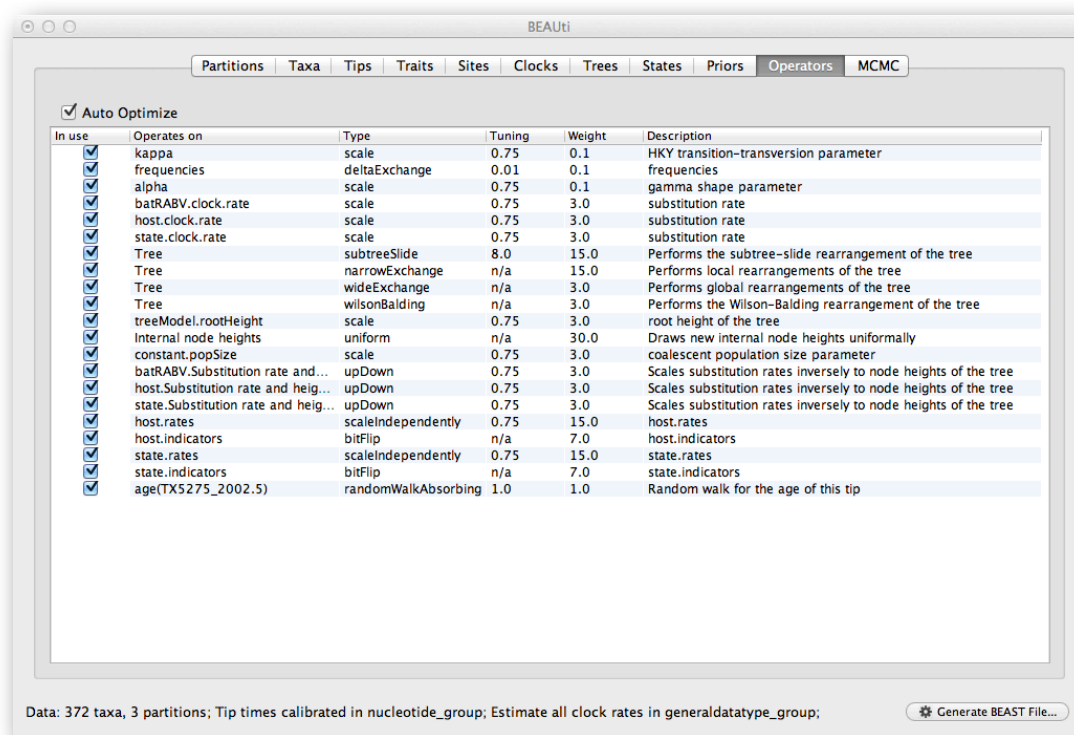
Review the prior settings under the **Priors** tab. Priors that have not been set to a proper prior distribution yet appear in red (**batRABV.clock.rate**). Click on the prior for the **batRABV.clock.rate** and a prior selection window will appear. Set the prior to a gamma distribution with shape = 0.001 and scale = 1000. The graphical representation of this prior distribution indicates that most prior mass is put on small values, but the density remains sufficiently diffuse. Notice that the prior setting turns black after confirming this setting by clicking **OK**. For the discrete host and location state rate, an approximation of a conditional reference prior (**Approx. Reference Prior**) (Ferreira and Suchard, 2008) is used. There is also a default uniform prior specification for the age of TX5275 (**age(TX5275_2002.5)**). We will assume that its sampling time for this tip is bounded by the sampling time distribution of the this data set, implying that it is sampled between 1997.5 and 2005.5. Click on the current uniform prior setting, set the **Upper** age to 8 years (reflecting the 1997.5 boundary) and click **OK**.



Setting up the operators

Each parameter in the model has one or more 'operators' (these are variously called *moves*, *proposals* or *transition kernels* by other MCMC software packages such as **MrBayes** and **LAMARC**). The operators specify how the parameters change as the MCMC runs. The operators tab in **BEAUti** has a table that lists the parameters, their operators and the tuning settings for these operators. In the first column are the parameter names while the next column has the type of operators that are acting on each parameter. For example, the scale operator scales the parameter up or down by a random

proportion and the uniform operator simply picks a new value uniformly within a range. Some parameters relate to the tree or to the divergence times of the nodes of the tree and these have special operators.



The next column, labelled **Tuning**, gives a tuning setting to the operator. Some operators don't have any tuning settings so have **n/a** under this column. The tuning parameter will determine how large a move each operator will make which will affect how often that change is accepted by the MCMC which will affect the efficiency of the analysis. For most operators (like the subtree slide operator) a larger tuning parameter means larger moves. However for the scale operator a tuning parameter value closer to 0.0 means bigger moves. At the top of the window is an option called **Auto Optimize** which, when selected, will automatically adjust the tuning setting as the MCMC runs to try to achieve maximum efficiency. At the end of the run a table of the operators, their performance and the final values of these tuning settings can be written to standard output.

The next column, labelled **Weight**, specifies how often each operator is applied relative to the others. Some parameters tend to be sampled very efficiently - an example is the kappa parameter - these parameters can have their operators down-weighted so that they are not changed as often.

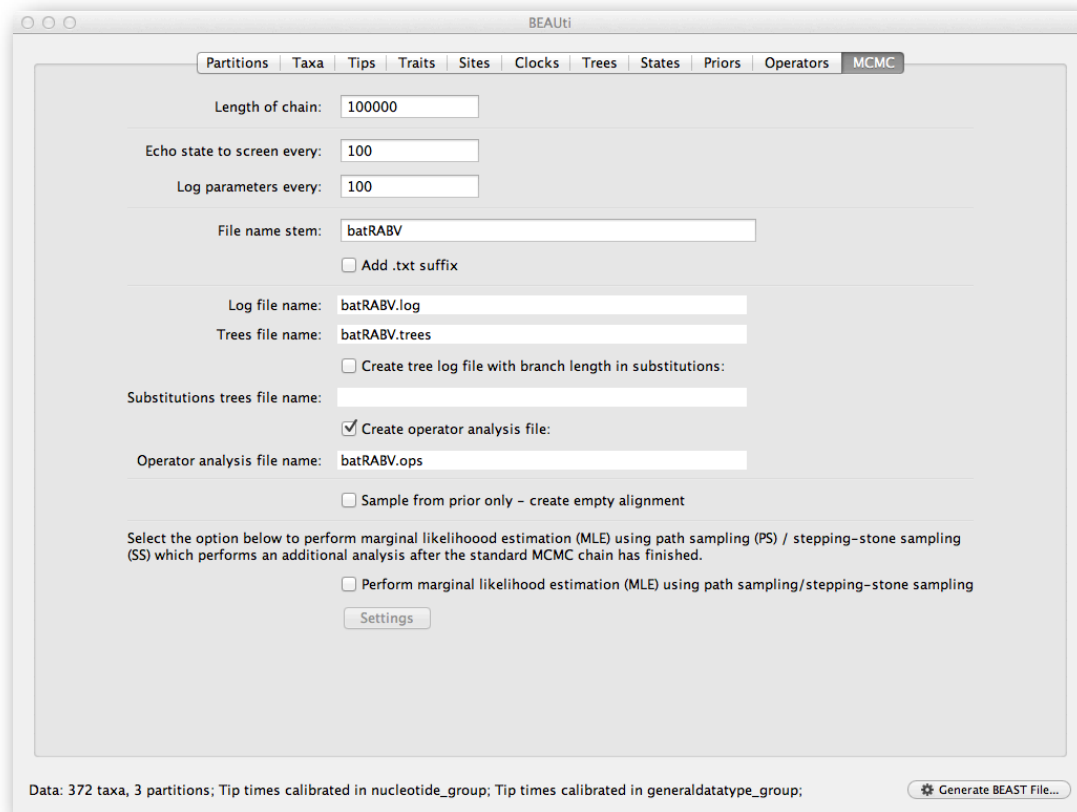
We can keep the default operator settings for the current analysis.

Setting the MCMC options

The **MCMC** tab in BEAUti provides settings to control the MCMC chain. Firstly we have the **Length of chain**. This is the number of steps the MCMC will make in the chain before finishing. How long this should depend on the size of the dataset, the complexity of the model and the precision of the answer required. The default value of 10,000,000 is entirely arbitrary and should be adjusted according to the size of your dataset. We will see later how the resulting log file can be analyzed using Tracer in order to examine whether a particular chain length is adequate.

The next couple of options specify how often the current parameter values should be displayed on the screen and recorded in the log file. The screen output is simply for monitoring the program's progress so can be set to any value (although if set too small, the sheer quantity of information being displayed on the screen will slow the program down). For the log file, the value should be set relative to the total length of the chain. Sampling too often will result in very large files with little extra

benefit in terms of the precision of the estimates. Sample too infrequently and the log file will not contain much information about the distributions of the parameters. You probably want to aim to store no more than 10,000 samples so this should be set to something \geq chain length / 10,000. For this dataset let's initially set the chain length to 100,000 as this will run quickly on most modern computers. Although the suggestion above would indicate a lower sampling frequency, in this case set both the sampling frequencies to 100.



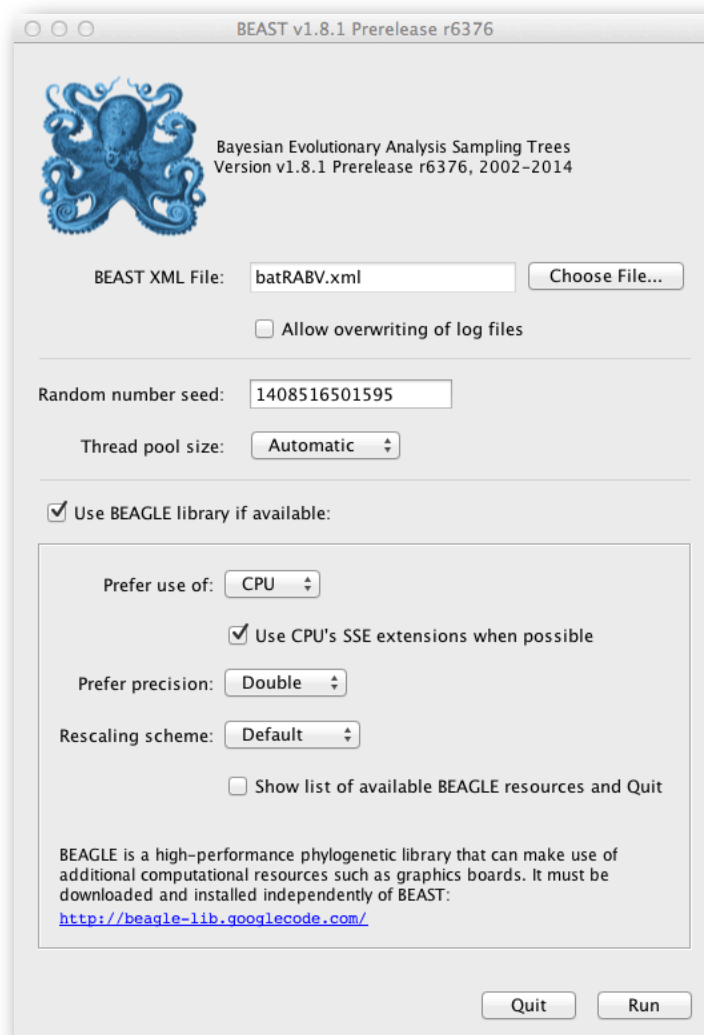
The next option allows the user to set the File stem name, which is set to **batRABV**. The next two options give the file names of the log files for the parameters and the trees. These will be set to a default based on the file stem name. By default, an operator analysis file is also created. Finally, an option is available to sample from the prior only, which can be useful to evaluate how divergent our posterior estimates are when information is drawn from the data. Here, we will not select this option, but analyze the actual data.

At this point we are ready to generate a BEAST XML file and to use this to run the Bayesian evolutionary analysis. To do this, either select the **Generate BEAST File...** option from the File menu or click the similarly labelled button at the bottom of the window. BEAUti will ask you to review the prior settings one more time before saving the file. Continue and choose a name for the file (for example, **batRABV.xml** by adding the xml extension to the file name stem) and save the file. For convenience, you can leave the *BEAUti* window open so that you can change the values and re-generate the *BEAST* file if necessary.

Running BEAST

Once the **BEAST** XML file has been created the analysis itself can be performed using **BEAST**. The exact instructions for running **BEAST** depends on the computer you are using, but in most cases a standard file dialog box will appear in which you select the XML file: If the command line version is being used then the name of the XML file is given after the name of the **BEAST** executable. If not selected by default, select the option **Use BEAGLE library if available** (see <http://beagle->

lib.googlecode.com/ to install BEAGLE library), which is required to run the discrete diffusion models, and use the default settings. When pressing **Run**, the analysis will be performed with detailed information about the progress of the run being written to the screen. When it has finished, the log file and the trees file will have been created in the same location as your XML file.



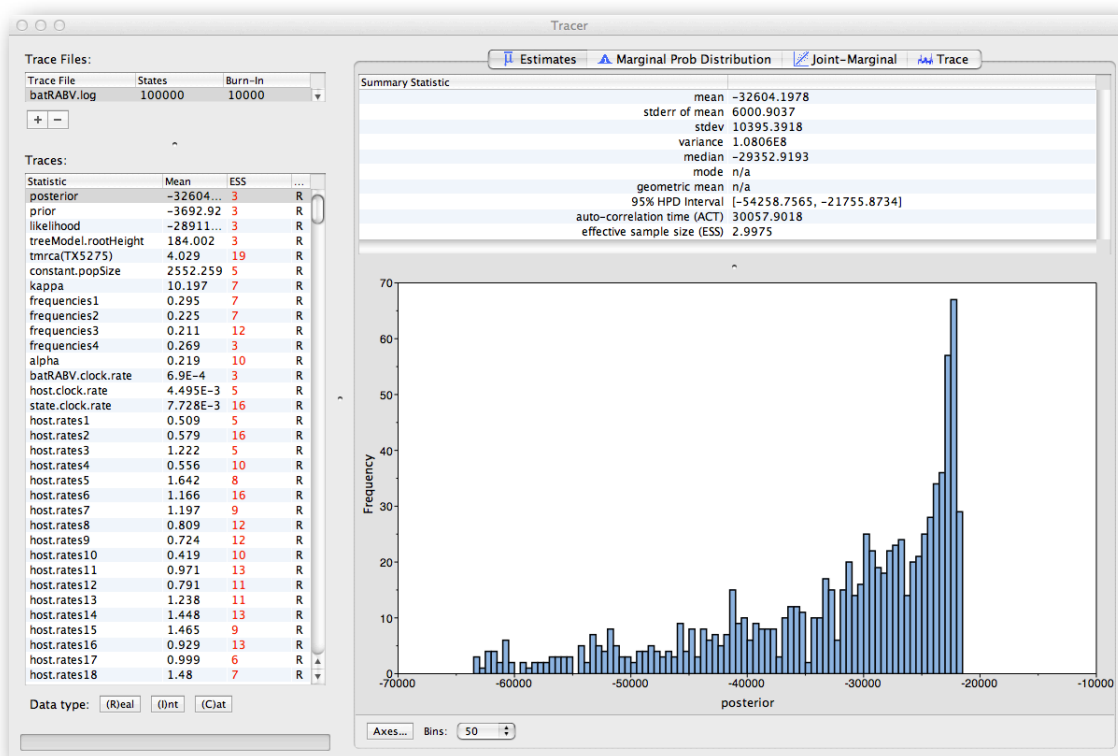
Analysing the BEAST output

To analyze the results of running BEAST we are going to use the program **Tracer**. The exact instructions for running Tracer differs depending on which computer you are using. Please see the README text file that was distributed with the version you downloaded. Once running, Tracer will look similar irrespective of which computer system it is running on.

Select the **Import Trace File...** option from the **File** menu. If you have it available, select the log file that you created in the previous section (**batRABV.log**). The file will load and you will be presented with a window similar to the one below. Remember that MCMC is a stochastic algorithm so the actual numbers will not be exactly the same.

On the left hand side is the name of the log file loaded and the traces that it contains. There are traces for a quantity proportional to posterior (this is the product of the data likelihood and the prior probabilities, on the log-scale), and the continuous parameters. Selecting a trace on the left brings up analyses for this trace on the right hand side depending on

tab that is selected. When first opened, the **posterior** trace is selected and various statistics of this trace are shown under the **Estimates** tab.



In the top right of the window is a table of calculated statistics for the selected trace. The statistics and their meaning are described in the table below.

Mean - The mean value of the samples (excluding the burn-in).

Stdev of mean - The standard error of the mean. This takes into account the effective sample size so a small ESS will give a large standard error.

Median - The median value of the samples (excluding the burn-in).

Geometric mean - The central tendency or typical value of the set of samples (excluding the burn-in).

95% HPD Lower - The lower bound of the highest posterior density (HPD) interval. The HPD is the shortest interval that contains 95% of the sampled values.

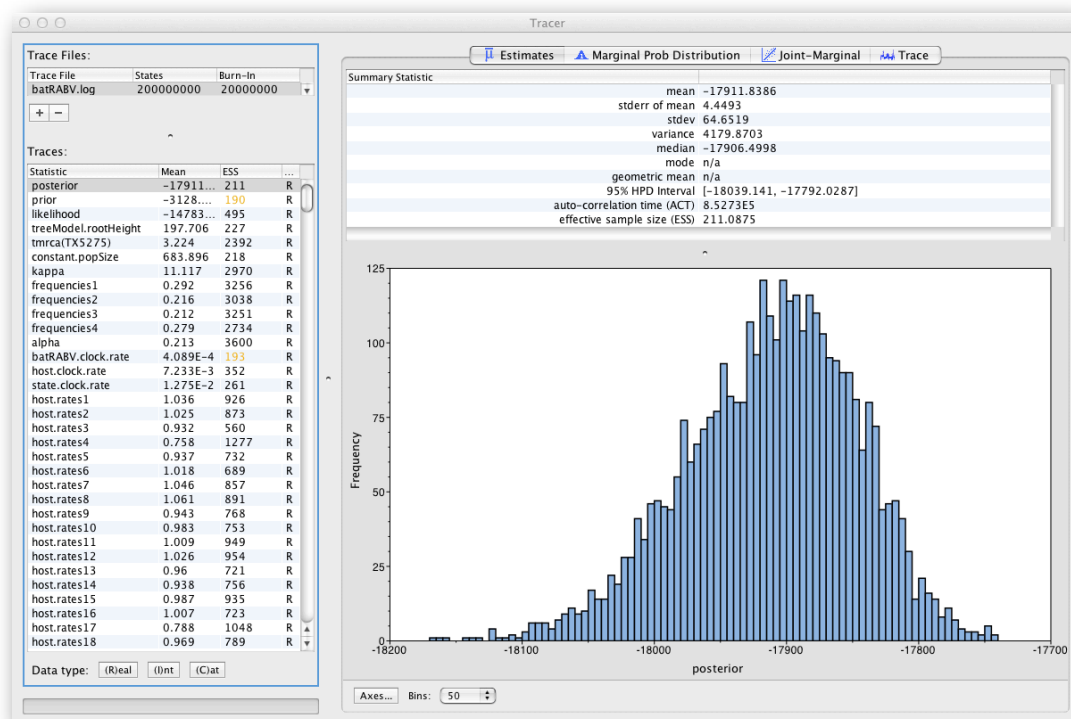
95% HPD Upper - The upper bound of the highest posterior density (HPD) interval.

Auto-Correlation Time (ACT) - The average number of states in the MCMC chain that two samples have to be separated by for them to be uncorrelated (i.e. independent samples from the posterior). The ACT is estimated from the samples in the trace (excluding the burn-in).

Effective Sample Size (ESS) - The effective sample size (ESS) is the number of independent samples that the trace is equivalent to. This is calculated as the chain length (excluding the burn-in) divided by the ACT.

Note that the effective sample sizes (ESSs) for all the traces are small (ESSs less than 100 are highlighted in red by Tracer and values > 100 but < 200 are in yellow). This is not good. A low ESS means that the trace contained a lot of correlated

samples and thus may not represent the posterior distribution well. In the bottom right of the window is a frequency plot of the samples which is expected given the low ESSs is extremely rough. Inspecting the **Trace** of many continuous parameters shows that the chain is still in the burn-in phase (the posterior values are still increasing over the entire chain), and this run does not allow us to summarize marginal posterior probability distributions for the parameters. The simple response to this situation is that we need to run the chain for longer. The example below was run for 200 million steps, sampling every 50,000th step, which means that 4,000 samples were stored in the log file. In this case, the MCMC run has reached stationarity, and almost all parameter traces still show satisfactory ESSs.

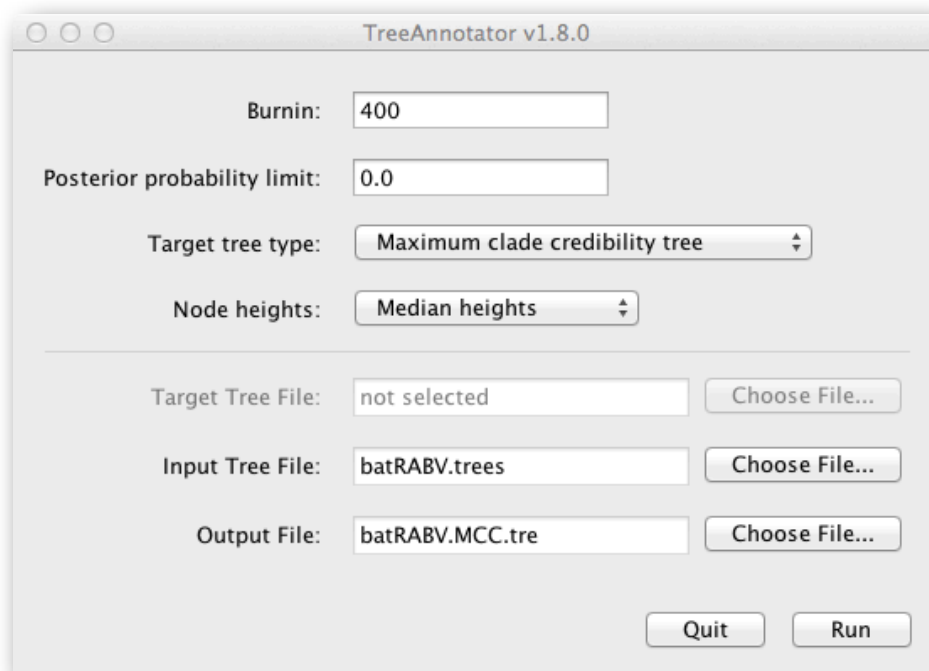


We can continue to summarize the annotated phylogeographic tree inferred with the BSSVS procedure and estimate the most significant rates of diffusion. If you are only interested in summarizing the Bayes Factor rates from the BSSVS analysis and not in summarizing the tree from your run, jump to the last section of this tutorial entitled **Identifying well-supported BF rates using Bayes factor test in SPREAD**. If you are also interested in summarizing the tree, continue to next section.

Summarizing the trees

We have seen how we can diagnose our MCMC run using Tracer and produce estimates of the marginal posterior distributions of parameters of our model. However, BEAST also samples trees (either phylogenies or genealogies) at the same time as the other parameters of the model. These are written to a separate file called the 'trees' file. This file is a standard NEXUS format file. As such it can easily be loaded into other software in order to examine the trees it contains. One possibility is to load the trees into a program such as MrBayes or PAUP* and construct a consensus tree in a similar manner to summarizing a set of bootstrap trees. In this case, the support values reported for the resolved nodes in the consensus tree will be the posterior probability of those clades.

In this tutorial, however, we are going to use a tool that is provided as part of the BEAST package to summarize the information contained within our sampled trees. The tool is called **TreeAnnotator** and once running, you will be presented with a window like the one below.



TreeAnnotator takes a single 'target' tree and annotates it with the summarized information from the entire sample of trees. The summarized information includes the average node ages (along with the HPD intervals), the posterior support and the average rate of evolution on each branch (for models where this can vary). The program calculates these values for each node or clade observed in the specified 'target' tree.

- **Burnin** - This is the number of trees in the input file that should be excluded from the summarization. This value is given as the number of trees rather than the number of steps in the MCMC chain. Thus for the example above, with a chain of 1,000,000 steps, sampling every 500 steps, there are 10,000 trees in the file. To obtain a 10% burn-in, set this value to 1,000.
- **Posterior probability limit** - This is the minimum posterior probability for a node in order for TreeAnnotator to store the annotated information. The default is 0.0 so every node, no matter what its support, will have information summarized. Make sure this value remains 0.0 as every node will require location annotation for further visualization.
- **Target tree type** - This has two options "Maximum clade credibility" or "User target tree". For the latter option, a NEXUS tree file can be specified as the Target Tree File, below. For the former option, TreeAnnotator will examine every tree in the Input Tree File and select the tree that has the highest sum of the posterior probabilities of all its nodes.
- **Node heights** - This option specifies what node heights (times) should be used for the output tree. If the "Keep target heights" is selected, then the node heights will be the same as the target tree. The other two options give node heights as an average (Mean or Median) over the sample of trees. Keep the default median node heights for the time being.
- **Target Tree File** - If the "User target tree" option is selected then you can use "Choose File..." to select a NEXUS file containing the target tree.
- **Input Tree File** - Use the "Choose File..." button to select an input trees file. This will be the trees file produced by BEAST.
- **Output File** - Select a name for the output tree file (e.g., batRABV.MCC.tre).

Once you have selected all the options above, press the "Run" button. TreeAnnotator will analyze the input tree file and write the summary tree to the file you specified. This tree is in standard NEXUS tree file format so may be loaded into any tree drawing package that supports this. However, it also contains additional information that can only be displayed using the FigTree program.

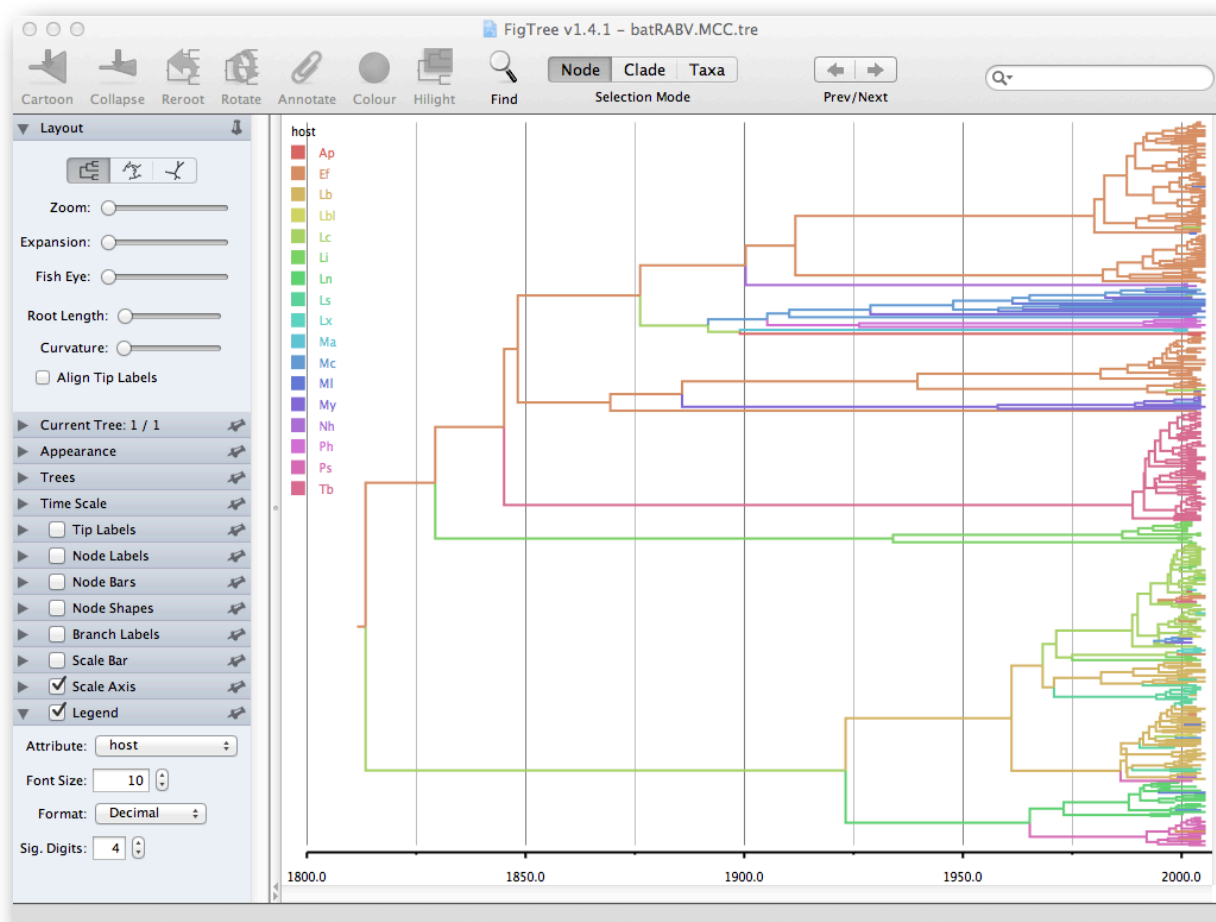
Viewing the annotated tree

Run FigTree and select the **Open...** command from the **File** menu. Select the tree file you created using TreeAnnotator in the previous section. The tree will be displayed in the FigTree window. On the left hand side of the window are the options and settings which control how the tree is displayed. In this case we want to display the posterior probabilities of each of the clades present in the tree and estimates of the age of each node. In order to do this you need to change some of the settings.

First open the **Branch Labels** section of the control panel on the left. Now select **posterior** from the **Display popup** menu. The posterior probabilities won't actually be displayed until you tick the check-box next to the **Branch Labels** title.

We now want to display bars on the tree to represent the estimated uncertainty in the date for each node. TreeAnnotator will have placed this information in the tree file in the shape of the 95% highest posterior density (HPD) intervals (see the description of HPDs, above). Open the **Node Bars** section of the control panel and you will notice that it is already set to display the 95% HPDs of the node heights so all you need to do is to select the check-box in order to turn the node bars on. We can also plot a time scale axis for this evolutionary history (select '**Scale Axis**' and deselect '**Scale bar**'). For appropriate scaling, open the '**Time Scale**' section of the control panel, set the '**Offset**' to 2005.5 (date of the most recent sample), the scale factor to -1.0. and '**Reverse Axis**' under '**Scale Axis**'.

Open the **Appearance** panel and alter the **Line Weight** to 2 in order to draw the tree with thicker lines. Under the same panel, alter **Colour by** and select **state**. Alternatively, color the branches by **host**. Unselect the **Tip Labels** and the **Scale Bar** option. Finally, in the **Legend** panel select the **state or host** attribute depending on what you have used to color the branches with. None of the options actually alter the tree's topology or branch lengths in anyway so feel free to explore the options and settings. You can also save the tree and this will save all your settings so that when you load it into FigTree again it will be displayed exactly as you selected.



Identifying well-supported BF rates using Bayes factors test in SPREAD

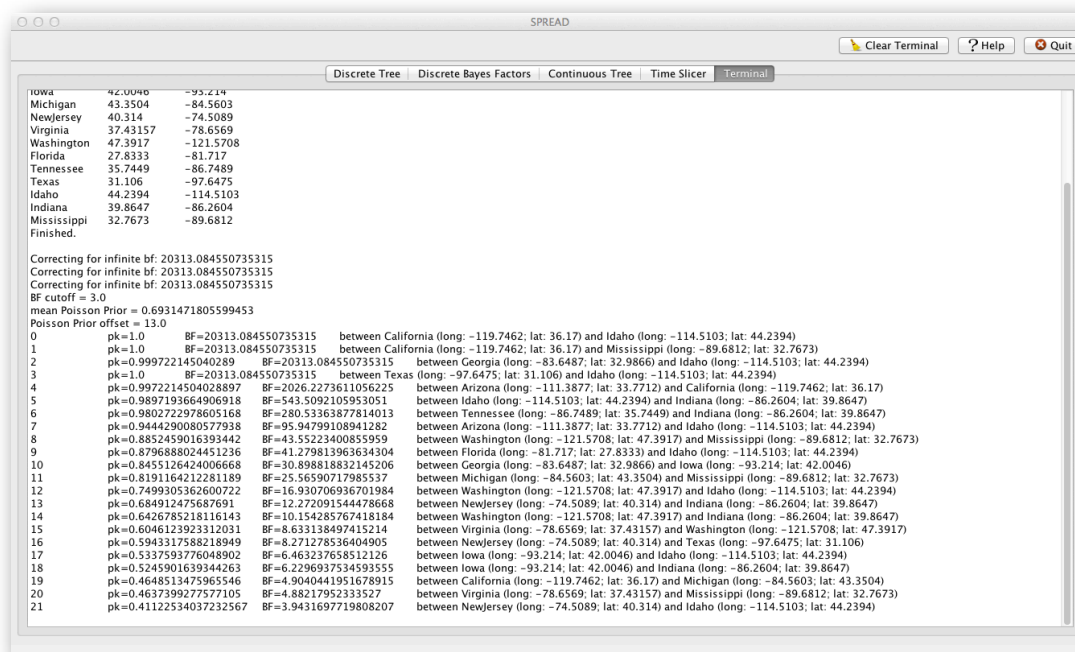
SPREAD (Spatial Phylogenetic Reconstruction of Evolutionary Dynamics) is a software to visualize the output from Bayesian phylogeographic analysis. SPREAD comes with its own map and virtual globe software, and it is able to generate KML files to visualize the output in GoogleEarth. Some of the functionalities of **SPREAD** that relate to the discrete phylogeographic analysis performed previously include visualizing location-annotated MCC trees, generation of KML output files for **Google Earth** and identification of well-supported rates using Bayes Factor test. The later option takes as input the rate matrix file (`batRABV.state.rates.log` for location states and `batRABV.host.rates.log` for host states) generated under the analysis using the Bayesian Stochastic Search Variable Selection (**BSSVS**) procedure. This test aims at identifying frequently invoked rates to explain the diffusion process and, in case of locations, visualize them using SPREAD's own map and in virtual globe software. A detailed tutorial for this particular step is available at http://www.phylogeography.org/tutorial/spread_tutorial.html#toc-Section-3.

Briefly, go to the **Discrete Bayes Factor** menu and using *Load log file* upload the output BEAST file containing the spatial rates and rate indicators (`batRABV.state.rates.log`). To visualize the results in the log file, a tab delimited file with location names and corresponding latitude and longitude coordinates needs to be uploaded using the *Load locations file*. The locations file should look like this:

Arizona	33.7712	-111.3877
California	36.17	-119.7462
Georgia	32.9866	-83.6487
Iowa	42.0046	-93.214
Michigan	43.3504	-84.5603
New Jersey	40.314	-74.5089
Virginia	18.0001	-64.8199
Washington	47.3917	-121.5708
Florida	27.8333	-81.717
Tennessee	35.7449	-86.7489
Texas	31.106	-97.6475
Idaho	44.2394	-114.5103
Indiana	39.8647	-86.2604
Mississippi	32.7673	-89.6812

Click done to upload your location file (**locationStates.txt**). The next step is to set up visualization attributes. Here you can specify the burn-in (default is 10%), the Poisson prior mean and offset, the Bayes factor cut-off, the color of the mapped rates, the KML name and more technical attributes. Once this is done, go to *Generate KML / Plot map* and click on *Plot* to visualize the Bayes factor rates in a map and click on *Generate* to save a KML output file that can be further inspected using Google Earth.

Finally, go to the menu **Terminal** to visualize the values of the Bayes factors for the rates that achieve a support beyond the specified cut-off and the respective locations involved for these rates. Note that for the reversible model the order of the locations, between 'X' and 'Y', is arbitrary as there is no directionality in this case.



Which rates receive the highest Bayes factor support? Try to compute similar support for the host transition rates.

Exercise 2: identifying predictors for the host transition process

This exercise builds on the previous analyses and aims at testing the factors that drive the host transition process for bat rabies viruses in North America. The original analyses resorted a population genetic approach and post hoc statistical procedures to test such predictors (Streicker et al., 2010), but here we adopt an extension of the discrete diffusion model as applied by Faria et al., 2013. This extension parameterizes the CTMC matrix as generalized linear model (GLM), in which log CTMC rates are a log function linear function of several potential predictors (most of the detail on the model can be found in Lemey et al., 2014). We use the predictors originally proposed by Streicker et al. (2010): host phylogenetic distance (based on host mitochondrial DNA), geographic range overlap, roost site overlap, and foraging niche overlap as approximated using three morphological measurements: wing aspect ratio, wing loading and body length, which are associated with foraging behavior in bats.

The GLM model allows us to estimate the support for an arbitrary number of predictors and quantify their contribution to the viral diffusion intensities while estimating the bat rabies evolutionary history. So, although they are also considered in the analysis, we will not focus anymore on the sequence evolutionary process or spatial diffusion across states in this part of the tutorial. At present, there is no BEAuti functionality to set up a GLM-discrete diffusion model, so we would have to manually edit our previous XML using a text editor (like TextWrangler). An xml ([batRABV_glm.xml](#)) is provided in which the model and analysis has been set up (with GLM specific edits indicated with comments) as well as output files resulting from this analysis.

The host state order

Because the GLM model requires us to specify predictors for the discrete diffusion process, we will need to get acquainted with the way the CTMC rates are associated with the discrete (host) states. Note that the order of the discrete states specified in the `generalDataType` will be important in this respect (which is alphabetical in our case):

```

<!-- START Discrete Traits Model -->

<!-- general data type for discrete trait model, 'host' -->
<generalDataType id="host.dataType">

  <!-- Number Of States = 17 -->
  <state code="Ap"/>
  <state code="Ef"/>
  <state code="Lb"/>
  <state code="Lb1"/>
  <state code="Lc"/>
  <state code="Li"/>
  <state code="Ln"/>
  <state code="Ls"/>
  <state code="Lx"/>
  <state code="Ma"/>
  <state code="Mc"/>
  <state code="Ml"/>
  <state code="My"/>
  <state code="Nh"/>
  <state code="Ph"/>
  <state code="Ps"/>
  <state code="Tb"/>
</generalDataType>

```

The CTMC matrix we will need to consider for the predictors follows this order:

	Ap	Ef	Lb	Lbl	Lc	Li	Ln	Ls	Lx	Ma	Mc	MI	My	Nh	Ph	Ps	Tb
Ap																	
Ef																	
Lb																	
Lbl																	
Lc																	
Li																	
Ln																	
Ls																	
Lx																	
Ma																	
Mc																	
MI																	
My																	
Nh																	
Ph																	
Ps																	
Tb																	

The GLM substitution model specification

The GLM model is specified as a `glmSubstitutionModel` element, which replaces the `generalSubstitutionModel` in the XML. As in a `generalSubstitutionModel` model, this specification contains a `frequencyModel`, while the `glmModel` part contains the six predictors (in linearized format in the `designMatrix`) and coefficients and indicators associated with these predictors. The critical part involves the predictor specification which takes the format of a linearized vector of values corresponding to the pair-wise host entries in the CTMC matrix. For each pair of hosts (i and j), the rate of transition (λ) between them is parameterized as:

$$\log \lambda_{ij} = \beta_1 \delta_1 \log(p_{1\{ij\}}) + \beta_2 \delta_2 \log(p_{2\{ij\}}) + \dots + \beta_6 \delta_6 \log(p_{6\{ij\}}),$$

where the β 's represent the coefficients in log space and the δ 's represent the indicators that determine the inclusion or exclusion of the predictors from the model. $p_{1\{ij\}}$ is the ij^{th} element of p_1 (= host distance); p_2 = range overlap; ... ; p_6 = body size differences. The predictor specification involves ordering the values according to the CTMC matrix entries, transforming them to log space, and standardizing them (to grant the predictors equal variance *a priori*).

```

<glmSubstitutionModel id="host.model">
  <dataType idref="host.dataType"/>
  <rootFrequencies>
    <frequencyModel normalize="true">
      <dataType idref="host.dataType"/>
      <frequencies>
        <parameter dimension="17"/>
      </frequencies>
    </frequencyModel>
  </rootFrequencies>
  <glmModel id="glmModel" family="logLinear" checkIdentifiability="true">
    <independentVariables>
      <parameter id="glmCoefficients" value="1 1 1 1 1 1"/>
      <indicator>
        <parameter id="coefIndicator" value="1 1 1 1 1 1"/>
      </indicator>
      <designMatrix id="designMatrix">
        <parameter id="hostDistance" value="0.132827181 0.585371903 0.329059784 0.00025
        <parameter id="rangeOverlap" value="0.476943868 -0.06457814 0.718750304 0.47893
        <parameter id="roostStructures" value="1 0 0 0 0 1 0 0 1 1 1 1 1 1 1 0 0 0
        <parameter id="wingAspectRatio" value="-0.06405182 0.05647927 0.05647927 0.8524
        <parameter id="wingLoading" value="-0.559133226 0.513303039 0.513303039 0.76378
        <parameter id="bodySize" value="-1.236980974 0.335824861 0.243898091 -0.0668675
      </designMatrix>
    </independentVariables>
  </glmModel>
</glmSubstitutionModel>

```

Take for example the host phylogenetic distances (first predictor), which are taken from the study by Streicker et al. (2010) (also provided in the `hostDistances.csv` file) and specified in our matrix format as follows (rounded to two decimals for representation here):

	Ap	Ef	Lb	Lbl	Lc	Li	Ln	Ls	Lx	Ma	Mc	MI	My	Nh	Ph	Ps	Tb
Ap		0.75	0.88	0.8	0.71	0.86	0.71	0.92	0.79	0.79	0.72	0.74	0.78	0.74	0.6	0.67	0.92
Ef	0.75		0.97	0.89	0.8	0.95	0.59	1.01	0.88	0.89	0.81	0.83	0.87	0.62	0.69	0.76	1.02
Lb	0.88	0.97		0.32	0.55	0.7	0.94	0.12	0.62	0.94	0.87	0.89	0.93	0.97	0.75	0.82	0.84
Lbl	0.8	0.89	0.32		0.47	0.62	0.86	0.35	0.55	0.87	0.79	0.81	0.85	0.89	0.67	0.74	0.76
Lc	0.71	0.8	0.55	0.47		0.65	0.77	0.58	0.58	0.78	0.7	0.72	0.76	0.8	0.58	0.65	0.67
Li	0.86	0.95	0.7	0.62	0.65		0.92	0.73	0.36	0.93	0.85	0.87	0.91	0.95	0.73	0.8	0.82
Ln	0.71	0.59	0.94	0.86	0.77	0.92		0.98	0.85	0.85	0.78	0.8	0.84	0.53	0.66	0.73	0.98
Ls	0.92	1.01	0.12	0.35	0.58	0.73	0.98		0.66	0.98	0.91	0.92	0.97	1	0.78	0.86	0.88
Lx	0.79	0.88	0.62	0.55	0.58	0.36	0.85	0.66		0.85	0.78	0.8	0.84	0.88	0.66	0.73	0.75
Ma	0.79	0.89	0.94	0.87	0.78	0.93	0.85	0.98	0.85		0.27	0.29	0.2	0.88	0.58	0.65	0.99
Mc	0.72	0.81	0.87	0.79	0.7	0.85	0.78	0.91	0.78	0.27		0.14	0.26	0.81	0.5	0.57	0.91
MI	0.74	0.83	0.89	0.81	0.72	0.87	0.8	0.92	0.8	0.29	0.14		0.28	0.82	0.52	0.59	0.93
My	0.78	0.87	0.93	0.85	0.76	0.91	0.84	0.97	0.84	0.2	0.26	0.28		0.87	0.56	0.64	0.98
Nh	0.74	0.62	0.97	0.89	0.8	0.95	0.53	1	0.88	0.88	0.81	0.82	0.87		0.68	0.76	1.01
Ph	0.6	0.69	0.75	0.67	0.58	0.73	0.66	0.78	0.66	0.58	0.5	0.52	0.56	0.68		0.43	0.79
Ps	0.67	0.76	0.82	0.74	0.65	0.8	0.73	0.86	0.73	0.65	0.57	0.59	0.64	0.76	0.43		0.87

Tb	0.92	1.02	0.84	0.76	0.67	0.82	0.98	0.88	0.75	0.99	0.91	0.93	0.98	1.01	0.79	0.87	
-----------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	--

After log-transformation and standardization, we get the following matrix (again rounded to two decimals for the sake of representation here):

	Ap	Ef	Lb	Lbl	Lc	Li	Ln	Ls	Lx	Ma	Mc	MI	My	Nh	Ph	Ps	Tb
Ap		0.13	0.59	0.33	0	0.53	0	0.7	0.29	0.31	0.03	0.1	0.26	0.11	-0.5	-0.2	0.72
Ef	0.13		0.86	0.63	0.34	0.81	-0.5	0.97	0.59	0.61	0.36	0.43	0.57	-0.4	-0.1	0.2	0.99
Lb	0.59	0.86		-2.2	-0.7	-0.1	0.76	-4.9	-0.4	0.79	0.55	0.61	0.74	0.85	0.14	0.4	0.46
Lbl	0.33	0.63	-2.2		-1.1	-0.4	0.53	-1.9	-0.7	0.55	0.29	0.36	0.5	0.62	-0.2	0.12	0.19
Lc	0	0.34	-0.7	-1.1		-0.2	0.22	-0.5	-0.6	0.25	-0	0.03	0.19	0.32	-0.6	-0.2	-0.2
Li	0.53	0.81	-0.1	-0.4	-0.2		0.71	0.09	-1.9	0.73	0.49	0.55	0.69	0.8	0.07	0.34	0.4
Ln	0	-0.5	0.76	0.53	0.22	0.71		0.87	0.49	0.51	0.25	0.31	0.46	-0.8	-0.2	0.07	0.89
Ls	0.7	0.97	-4.9	-1.9	-0.5	0.09	0.87		-0.2	0.89	0.67	0.73	0.85	0.95	0.27	0.52	0.58
Lx	0.29	0.59	-0.4	-0.7	-0.6	-1.9	0.49	-0.2		0.51	0.25	0.32	0.46	0.58	-0.2	0.08	0.15
Ma	0.31	0.61	0.79	0.55	0.25	0.73	0.51	0.89	0.51		-2.6	-2.4	-3.5	0.6	-0.6	-0.2	0.91
Mc	0.03	0.36	0.55	0.29	-0	0.49	0.25	0.67	0.25	-2.6		-4.5	-2.8	0.35	-1	-0.6	0.69
MI	0.1	0.43	0.61	0.36	0.03	0.55	0.31	0.73	0.32	-2.4	-4.5		-2.6	0.41	-0.9	-0.5	0.75
My	0.26	0.57	0.74	0.5	0.19	0.69	0.46	0.85	0.46	-3.5	-2.8	-2.6		0.55	-0.6	-0.3	0.87
Nh	0.11	-0.4	0.85	0.62	0.32	0.8	-0.8	0.95	0.58	0.6	0.35	0.41	0.55		-0.1	0.18	0.97
Ph	-0.5	-0.1	0.14	-0.2	-0.6	0.07	-0.2	0.27	-0.2	-0.6	-1	-0.9	-0.6	-0.1		-1.4	0.3
Ps	-0.2	0.2	0.4	0.12	-0.2	0.34	0.07	0.52	0.08	-0.2	-0.6	-0.5	-0.3	0.18	-1.4		0.54
Tb	0.72	0.99	0.46	0.19	-0.2	0.4	0.89	0.58	0.15	0.91	0.69	0.75	0.87	0.97	0.3	0.54	

The order in which these predictor values need to be linearized for xml vector representation is as follows (the entries now represent the ordering):

	Ap	Ef	Lb	Lbl	Lc	Li	Ln	Ls	Lx	Ma	Mc	MI	My	Nh	Ph	Ps	Tb
Ap		1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
Ef	137		17	18	19	20	21	22	23	24	25	26	27	28	29	30	31
Lb	138	153		32	33	34	35	36	37	38	39	40	41	42	43	44	45
Lbl	139	154	168		46	47	48	49	50	51	52	53	54	55	56	57	58
Lc	140	155	169	182		59	60	61	62	63	64	65	66	67	68	69	70
Li	141	156	170	183	195		71	72	73	74	75	76	77	78	79	80	81
Ln	142	157	171	184	196	207		82	83	84	85	86	87	88	89	90	91
Ls	143	158	172	185	197	208	218		92	93	94	95	96	97	98	99	100
Lx	144	159	173	186	198	209	219	228		101	102	103	104	105	106	107	108
Ma	145	160	174	187	199	210	220	229	237		109	110	111	112	113	114	115
Mc	146	161	175	188	200	211	221	230	238	245		116	117	118	119	120	121
MI	147	162	176	189	201	212	222	231	239	246	252		122	123	124	125	126
My	148	163	177	190	202	213	223	232	240	247	253	258		127	128	129	130
Nh	149	164	178	191	203	214	224	233	241	248	254	259	263		131	132	133
Ph	150	165	179	192	204	215	225	234	242	249	255	260	264	267		134	135
Ps	151	166	180	193	205	216	226	235	243	250	256	261	265	268	270		136
Tb	152	167	181	194	206	217	227	236	244	251	257	262	266	269	271	272	

Note that order in the linearized vector thus follows the upper matrix *row-by-row* and then the lower matrix *column-by-column*. The values are specified in this order under **value** in each predictor **parameter** in the **designMatrix**.

Additional XML edits

A statistic is added that returns the product of the coefficients and the respective indicators for the predictors (below the **glmSubstitutionModel**). This is not essential for our analysis.

```

<!-- GLM edit: statistic that returns the product of the coefficients and the respective indicators for the predictors -->
<productStatistic id="coefficientsTimesIndicators" elementwise="false">
  <parameter idref="glmCoefficients"/>
  <parameter idref="coefIndicator"/>
</productStatistic>

```

In the operators, the standard discrete model operators on the rates and rate indicators are substituted by the GLM model operators on the predictor indicators (also **bitFlipOperator**) and predictor coefficients (both a **randomWalkOperator** and **mvnOperator**).

```

<!-- GLM edit: standard discrete model operators on the rates and rate indicators substituted by glm model operators
<!--
<scaleOperator scaleFactor="0.75" weight="1" scaleAllIndependently="true">
  <parameter idref="host.rates"/>
</scaleOperator>
<bitFlipOperator weight="1">
  <parameter idref="host.indicators"/>
</bitFlipOperator>
-->
<bitFlipOperator weight="3">
  <parameter idref="coefIndicator"/>
</bitFlipOperator>
<randomWalkOperator windowSize="0.5" weight="1">
  <parameter idref="glmCoefficients"/>
</randomWalkOperator>
<mvnOperator scaleFactor="1" weight="5" formXtXInverse="true">
  <parameter idref="glmCoefficients"/>
  <varMatrix>
    <parameter idref="designMatrix"/>
  </varMatrix>
</mvnOperator>

```

In the priors, three modifications have been made: (i) the Poisson prior on the total number of non-zero rates in the standard discrete model with BSSVS is substituted by binomial prior on the predictor indicators, (ii) the uniform prior on the host frequencies is commented out as no `id` is specified for those in the GLM substitution model and they are not estimated but fixed to equal frequencies, (iii) the standard gamma prior on the rates in the discrete model is substituted by a normal prior on the GLM coefficients.

```

<!-- GLM edit: poissonPrior prior on the total number of non-zero rates in the standard discrete model with BSSVS
<!--
<poissonPrior mean="0.6931471805599453" offset="16.0">
  <statistic idref="host.nonZeroRates"/>
</poissonPrior>
-->
<!-- using the binomialLikelihood we specify a 50% prior mass on no predictors being included; for 6 trials (predi
<binomialLikelihood>
  <proportion>
    <parameter value="0.019"/>
  </proportion>
  <trials>
    <parameter value="1 1 1 1 1 1"/>
  </trials>
  <counts>
    <parameter idref="coefIndicator"/>
  </counts>
</binomialLikelihood>

<!-- GLM edit: uniform prior host frequencies removed as they are not defined in the GLM substitution model; they
<!--
<uniformPrior lower="0.0" upper="1.0">
  <parameter idref="host.frequencies"/>
</uniformPrior>
-->

<!-- GLM edit: standard gammaPrior prior on the rates in the discrete model substituted by normal prior on the glm
<!--
<cachedPrior>
  <gammaPrior shape="1.0" scale="1.0" offset="0.0">
    <parameter idref="host.rates"/>
  </gammaPrior>
  <parameter idref="host.rates"/>
</cachedPrior>
-->
<normalPrior mean="0" stdev="2">
  <parameter idref="glmCoefficients"/>
</normalPrior>

```


Using the **binomialLikelihood** we specify a 50% prior mass on no predictors being included. For none of the 6 predictors being included, the binomial distribution probability is 0.50 if the success probability for each predictor inclusion is set to 0.11.

In the screen log, the logging of the total number of non-zero rates in the standard discrete model with BSSVS is substituted by logging of the new product statistic (coefficients times indicators).

```

<!-- START Discrete Traits Model
<!-- GLM edit: logging of the total number of non-zero rates in the standard discrete model with BSSVS
<!--
<column label="host.nonZeroRates" sf="6" width="12">
  <sumStatistic idref="host.nonZeroRates"/>
</column>
-->
<column label="coefficientsTimesIndicators" sf="6" width="12">
  <productStatistic idref="coefficientsTimesIndicators"/>
</column>

```

In the file log, the logging of the logging of the standard discrete model (with BSSVS) parameters and statistic is commented out.

```

<!-- START Discrete Traits Model
<!-- GLM edit: logging of standard discrete model (with BSSVS) parameters and statistics removed;
<!--
<parameter idref="host.rates"/>
<parameter idref="host.indicators"/>
<sumStatistic idref="host.nonZeroRates"/>
-->

```

Finally, logging of standard discrete model (with BSSVS) parameters and statistic to a separate file is substituted by logging of the GLM model parameters and statistics to a separate log file (**batRABV_glm.host.model.log**).

```

<!-- START Discrete Traits Model
<!-- GLM edit: logging of standard discrete model (with BSSVS) parameters and statistics to a separate file substituted
<!--
<log id="batRABV.hostrateMatrixLog" logEvery="10000" fileName="batRABV_glm.host.rates.log">
  <parameter idref="host.rates"/>
  <parameter idref="host.indicators"/>
  <sumStatistic idref="host.nonZeroRates"/>
</log>
-->
<log id="glmFileLog" logEvery="10000" fileName="batRABV_glm.host.model.log">
  <parameter idref="coefIndicator"/>
  <parameter idref="glmCoefficients"/>
  <productStatistic idref="coefficientsTimesIndicators"/>
  <glmModel idref="glmModel"/>
  <parameter idref="host.clock.rate"/>
</log>

```

Evaluating the support for the predictors of rabies host transitioning

The XML file (`batRABV_glm.xml`) can be run using BEASTv1.8, but will take a long time to complete. Output files for this analysis are therefore provided. Here, we focus on the GLM-diffusion parameters (in the `batRABV_glm.host.model.log` file), which are simultaneously estimated the sequence and location diffusion (other output files). Examine the `batRABV_glm.host.model.log` file in TRACER and summarize the predictor inclusion probabilities, their Bayes factors (based on posterior odds over prior odds for predictor inclusion) and their effect sizes.

Conclusion and Resources

This tutorial only scratches the surface of the analyses that are possible to undertake using BEAST. It has hopefully provided a relatively gentle introduction to the fundamental steps that will be common to all BEAST analyses and provide a basis for more challenging investigations. BEAST is an ongoing development project with new models and techniques being added on a regular basis. The BEAST website provides details of the mailing list that is used to announce new features and to discuss the use of the package. The website also contains a list of tutorials and recipes to answer particular evolutionary questions using BEAST as well as a description of the XML input format, common questions and error messages.

- The BEAST website: <http://beast.bio.ed.ac.uk/> (or <http://mcmc.googlecode.com>)
- Tutorials: <http://beast.bio.ed.ac.uk/Tutorials/>
- Phylogeography: <http://www.phylogeography.org> (includes **SPREAD** and tutorial)
- Frequently asked questions: <http://beast.bio.ed.ac.uk/FAQ/>