# BEAST

Bayesian Evolutionary Analysis Sampling Trees

## Revealing the evolutionary dynamics of influenza

### Introduction

This tutorial provides a step-by-step explanation on how to reconstruct the evolutionary dynamics of influenza based on a set of virus sequences which have been isolated at different points in time ('heterochronous' data) using *BEAST*. We will focus on influenza A virus evolution, in particular on the emergence of swine-origin influenza A (H1N1) virus in 2009 (H1N1/09) and on the epidemic dynamics of H3N2 in the New York State. The H1N1/09 data set is a subset of an analyzed set genomes in a study that provides insights into the origins and evolutionary genomics of this outbreak (Smith et al., 2009). The H3N2 data is a subset of a comprehensive data set spanning several epidemic seasons in the New York state, which has been used to unravel the genomic and epidemiological dynamics of this virus (Rambaut et al., 2008). In the first exercise, the aim is to obtain an estimate of the rate of molecular evolution, an estimate of the date of the most recent common ancestor, an estimate of the H1N1/09 epidemic growth or the H1N1/09 basic reproductive number. In the second exercise, we will examine how H3N2 diversity fluctuates through time.

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 three software packages in a format that is compatible with your computer system (all three 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 v1.8.4. *BEAST* releases are generally available for download from http://beast.bio.ed.ac.uk/, but the latest (pre-)releases can also be found at https://github.com/beast-dev/beast-mcmc/releases.

- **Tracer** - this program is used to explore the output of **BEAST** (and other Bayesian MCMC programs). It graphically and quantitively summarizes the empirical 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://beast.bio.ed.ac.uk/.

- **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 v1.4.3. It is available for download from http://tree.bio.ed.ac.uk/.

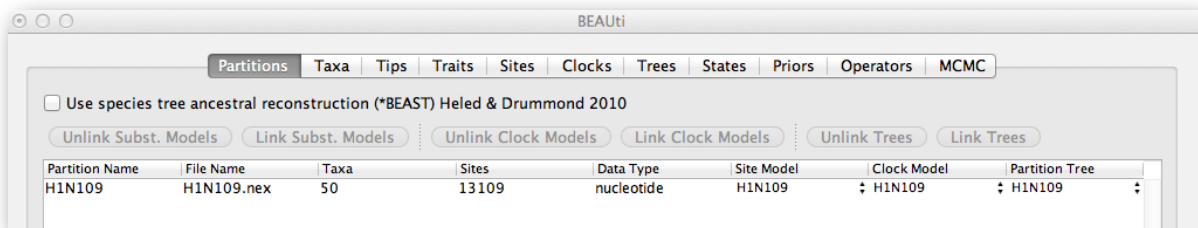# EXERCISE 1: The swine-origin influenza A outbreak

## 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.
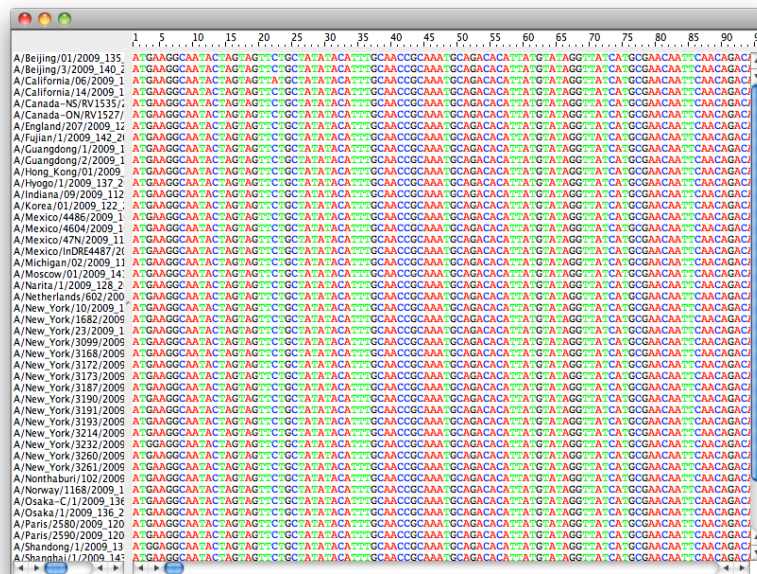
## Loading the NEXUS file

To load a NEXUS format alignment, simply select the Import NEXUS... option from the File menu.
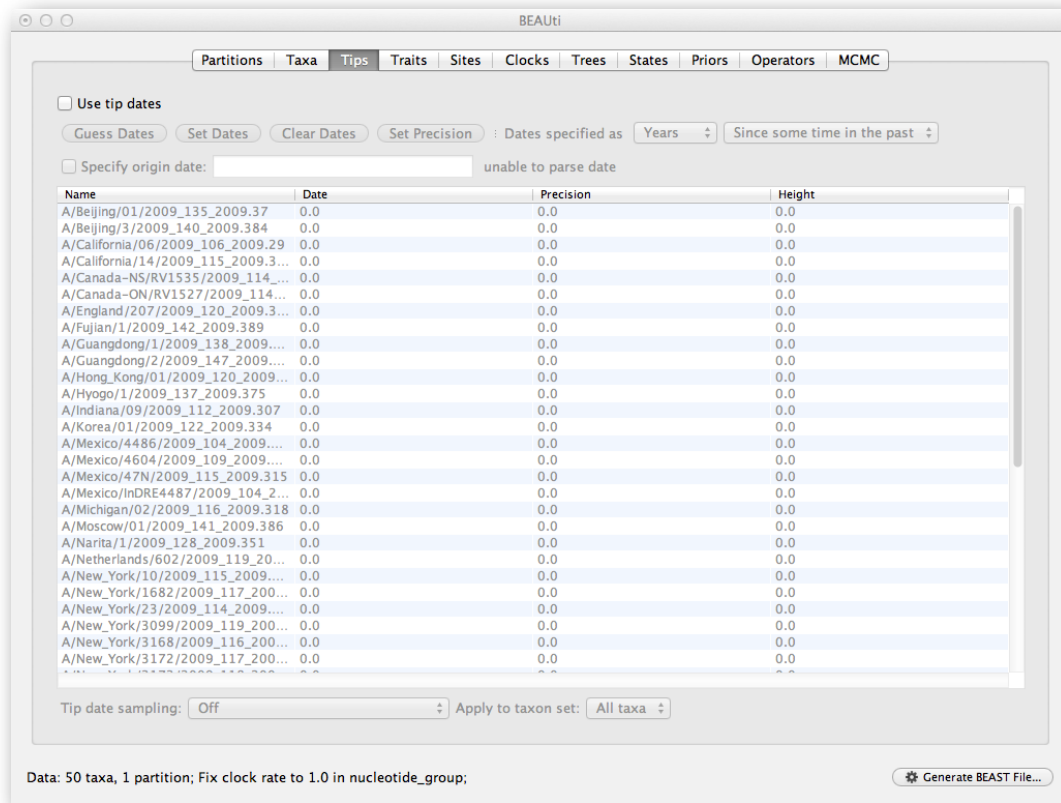
### The NEXUS alignment

Select the file called **H1N109.nex**. This file contains an alignment of 50 genomes (concatenated segments), 13109 nucleotides in length. Once loaded, the new data will be listed under **Partitions** as shown in the figure:
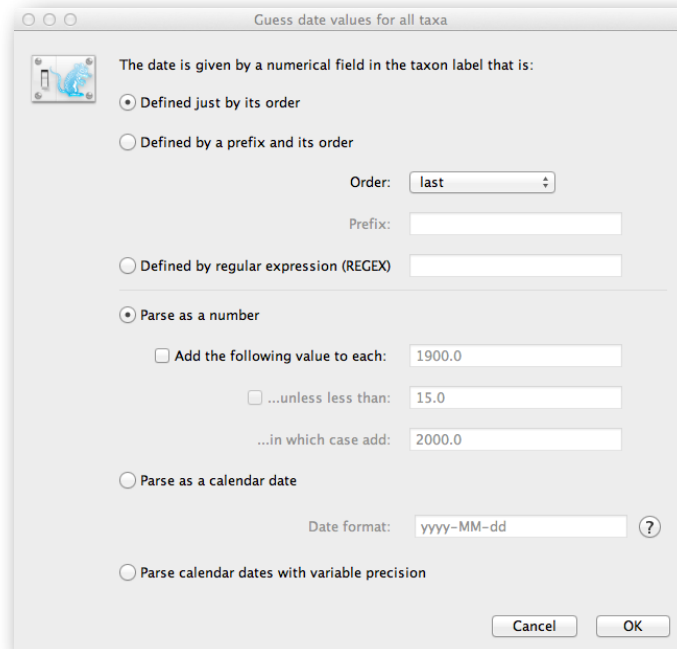


Double-click on the row of the table (but not on Partition Name) to display the actual sequence alignment:



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 sequences have been sampled from the H1N1/09 epidemic between March and May 2009. To set these dates switch to the 'Tips' panel using the tabs at the top of the window:
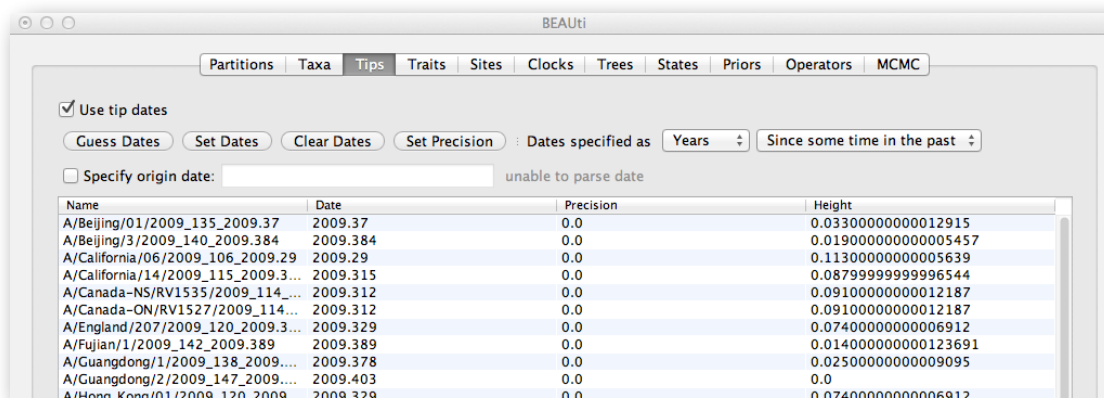
Select the box labelled 'Use tip dates'. The actual sampling in fractional years is encoded in the name of each taxon and we could simply edit the value in the 'Date' column of the table to reflect these. However, if the taxa names contain the calibration information, then a convenient way to specify the dates of the sequences in *BEAUti* is to use the 'Guess Dates' button at the top of the 'Data' panel. 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 2009) or backwards in time from the present (as in the case of radiocarbon ages).
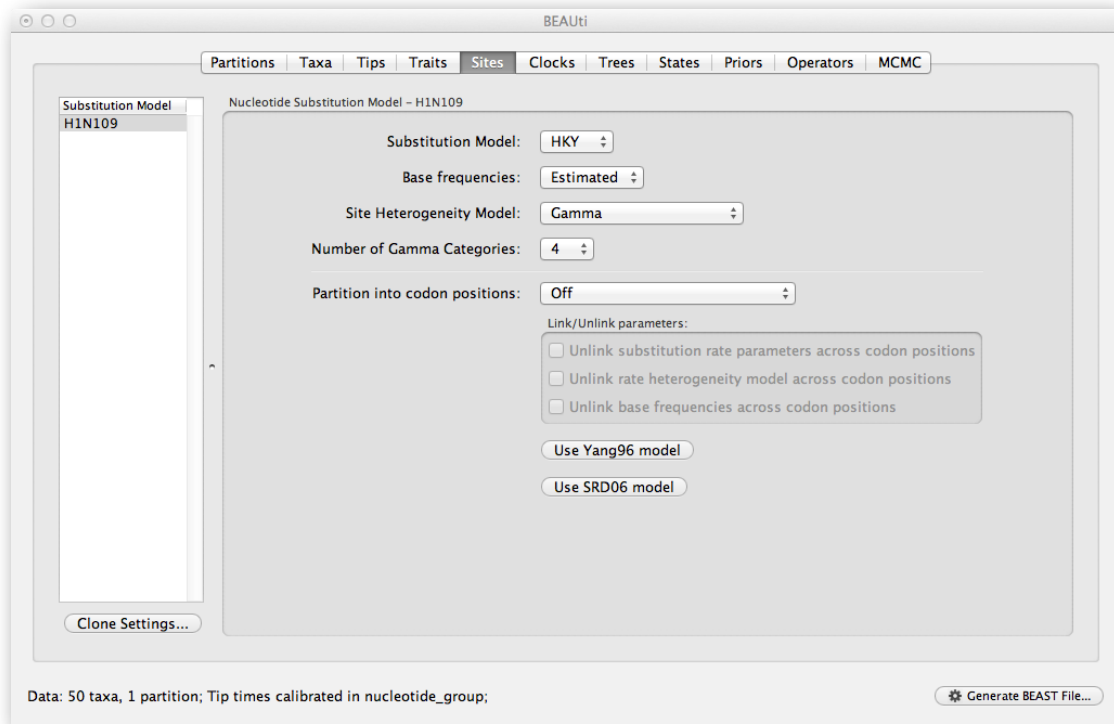


The "Height" column lists the ages of the tips relative to time 0 (in our case 2009.403). The Precision column allows specifying with what precision the sampling time is know. To include taxa only known up to the year of sampling for example (e.g., 2009), a precision of 1 year can be set and the age 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.

## Setting the substitution model

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 or amino acids.

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*:
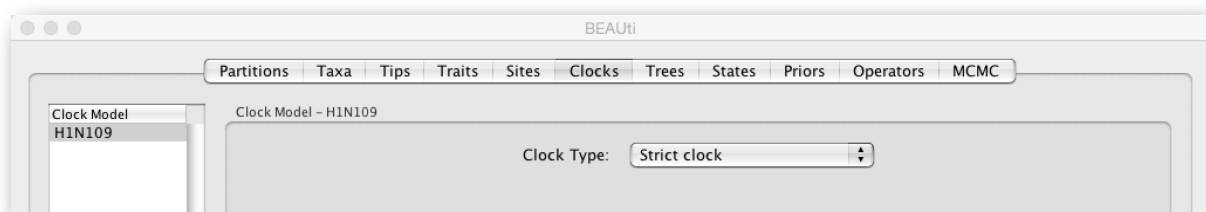
• Selecting the 'Partition into codon positions' option assumes that the data are aligned as codons. This option will then estimate a separate rate of substitution for each codon position, or for 1+2 versus 3, depending on the setting.

• 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 a set of rate heterogeneity parameters (gamma shape parameter and/or proportion of invariant sites) for each codon position.

For this tutorial, keep the default 'HKY' model, the default 'Estimated' base frequencies and select 'Gamma' as 'Site Heterogeneity Model' (with 4 discrete categories) before proceeding to the 'Clocks' tab.

## Setting the 'molecular clock' model

The 'Molecular Clock Model' options allows us to choose between a strict and a relaxed (uncorrelated lognormal or uncorrelated exponential) clock. Because of the low diversity data we analyze here, a relaxed clock would probably be over-parameterization. Hence, we keep a strict clock setting.
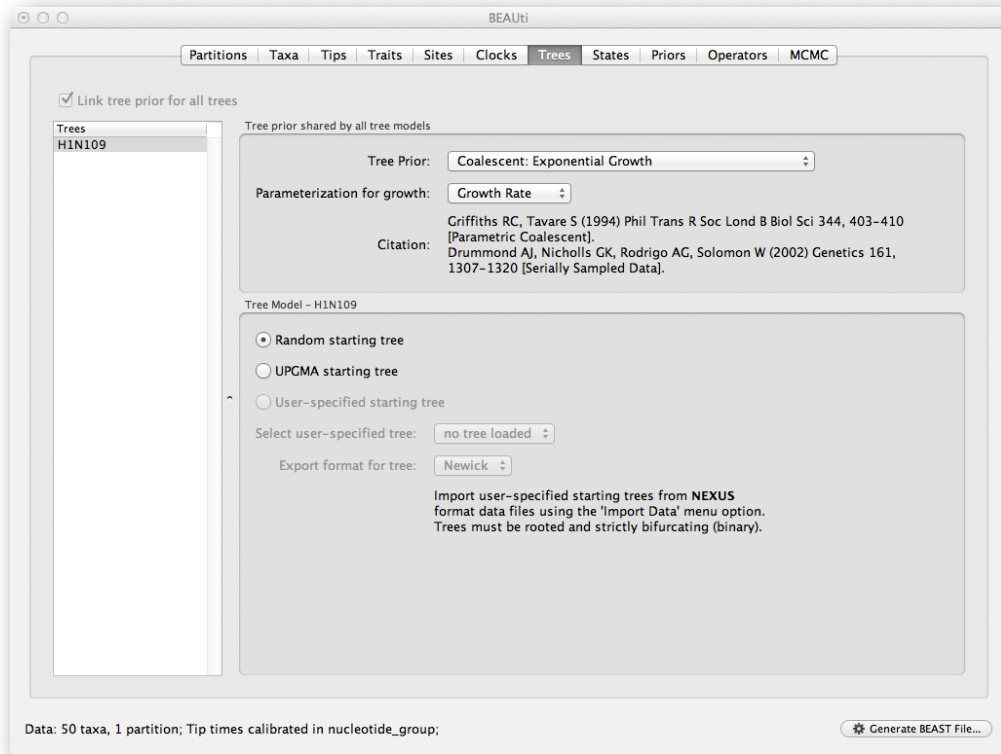


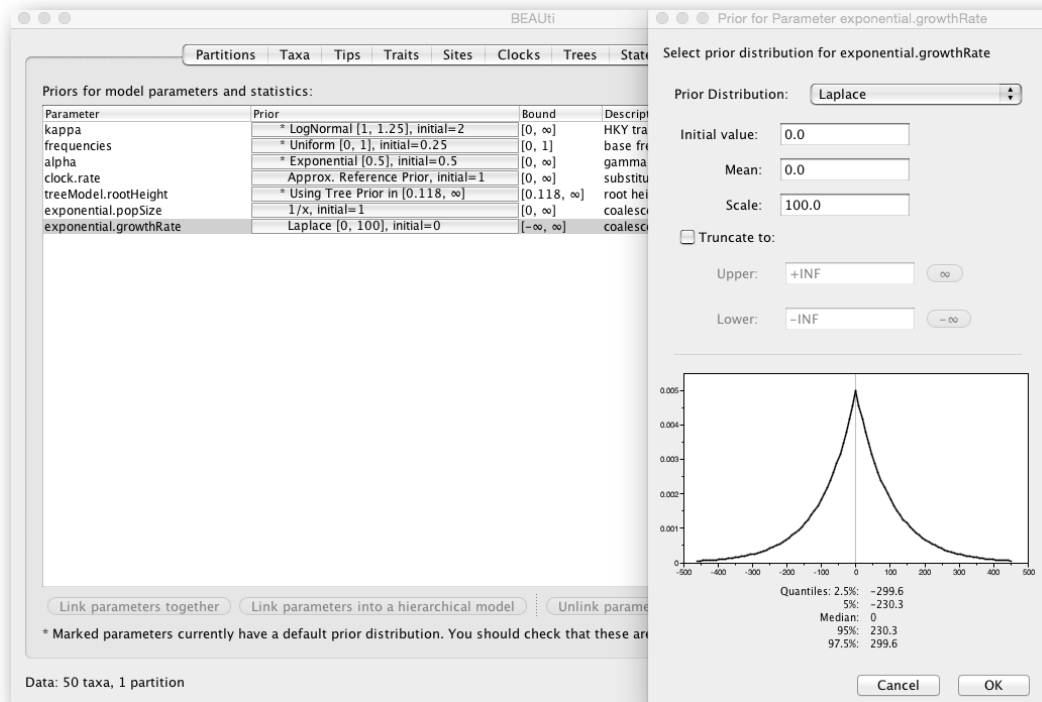Now move on to the 'Trees' panel.

## Setting the tree prior

This panel contains settings about the tree. Firstly the starting tree is specified to be 'randomly generated'. The other main setting here is to specify the 'Tree prior' which describes how the population size is expected to change over time for coalescent models. The default tree prior is set to a constant size coalescent prior.

To estimate the epidemic growth rate, we will change this demographic model to an exponential growth coalescent prior, which is intuitively appealing for viral outbreaks. Switch the option for 'Tree Prior' to 'Coalescent: Exponential Growth'.

## Setting up the priors

Now switch to the 'Priors' tab. This panel has a table showing every parameter of the currently selected model and what the prior distribution is for each. A prior allows the user to 'inform' the analysis by selecting a particular distribution. Although some of the default priors may be improper, with sufficiently informative data the posterior becomes proper. If priors or not set, they will appear in red.
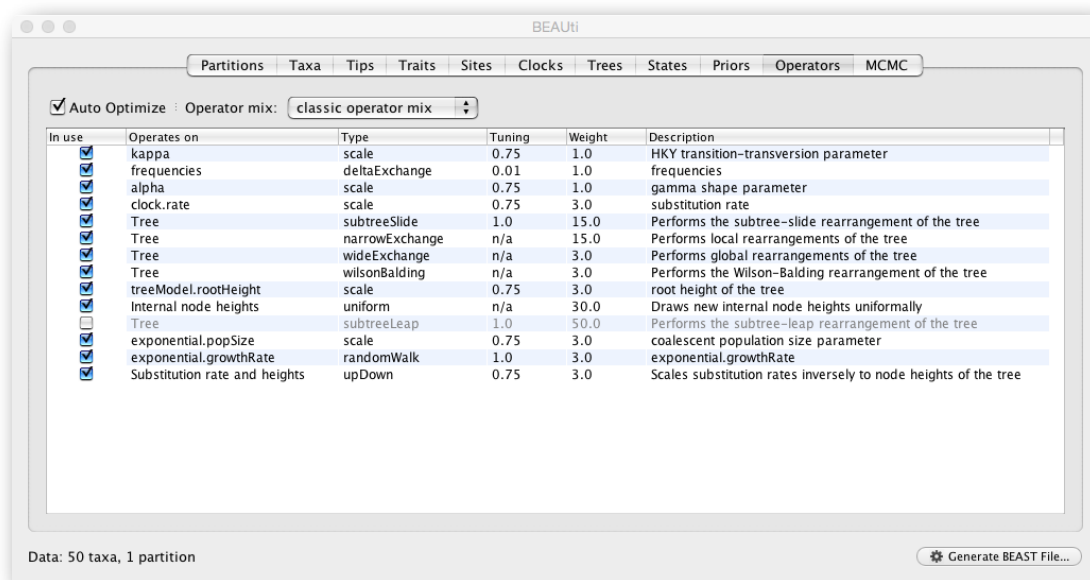
Note that a prior distribution must be specified for every parameter and whilst *BEAUti* provides default options these are not necessarily tailored to the problem and data being analyzed. In this case, the default Laplace prior prefers relatively small growth rates while this parameter take can take on relatively large values on this epidemic scale. Therefore, we will increase the variance of this prior distribution by setting the scale to 100.

The default prior on the rate of evolution (clock.rate) is an approximation of a conditional reference prior (**Approx. Reference Prior**) (Ferreira and Suchard, 2008). If the sequences are not associated with different sampling dates (they are contemporaneous), or when the sampling time range is trivial for the evolutionary scale of the taxa, the substitution rate can be fixed to a value based on another source, or better, a prior distribution can be specified to also incorporate the uncertainty of this 'external' rate. Fixing the rate to 1.0 will result in the ages of the nodes of the tree being estimated in units of substitutions per site (i.e. the normal units of branch lengths in popular packages such as *MrBayes*). Note that when selecting to fix the rate to a value, the transition kernel(s) on this parameter ('Operators' panel, see next section) will be automatically unselected.

## 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. These will be called things like **kappa** which means the HKY model's kappa parameter (the transition-transversion bias). 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 proportion, the random walk operator adds or subtracts an amount to the parameter 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 in turn affect the efficiency of the analysis. For most operators (like random walk and subtree slide operators) 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 will be written to standard output. These can then be used to set the starting tuning settings in order to minimize the amount of time taken to reach optimum performance in subsequent runs.

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 will start by using the default settings for this analysis. As of BEAST v1.8.4, different options are available w.r.t. exploring tree space. In this tutorial, we will use the 'classic operator mix', which consists of of set of tree transition kernels that propose changes to the tree. There is also an option to fix the tree topology as well as a 'new experimental mix', which is currently under development with the aim to improve mixing for large phylogenetic trees.

## 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 be depends 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 analysed using *Tracer* in order to examine whether a particular chain length is adequate. Change the chain length to 1,000,000 for our initial test run.



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 the chain length / 10,000. For this dataset let's initially set the chain length to 100,000 as this will run reasonably
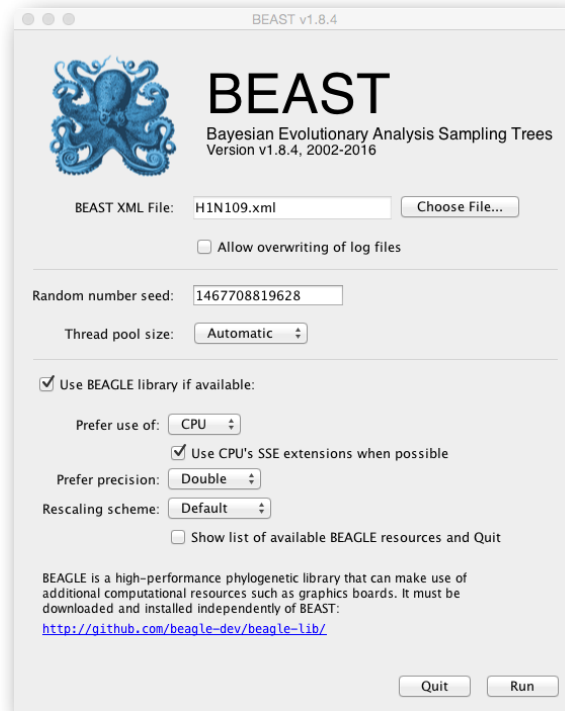
quickly on most modern computers. Although the suggestion above would indicate a lower sampling frequency, in this case set both the sampling frequencies to 100. A useful exercise could be to examine the sensitity of the growth rate estimates to different scale values for this prior distribution (e.g. scale = 1, 10, 100).

The next option allows the user to set the File stem name; if not set to 'H1N109' by default, you can type this in here. The next two options give the file names of the log files for the parameters and the trees. These will be set to based on the file stem name. You can also log the operator analysis to a file. An option is also 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. Finally, one can select to perform marginal likelihood estimation to assess model fit, which is not needed in this exercise. So, 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. BEAST will ask you to review the prior settings one more time before saving the file. Continue and keep the default name for the file, add the xml extension (H1N109.xml) and save the file.

## 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 dialog box will appear in which you select the XML file:



Press the 'Choose File' button and select the XML file you just created and press 'Run'. When you have installed the *BEAGLE* library (https://github.com/beagle-dev/beagle-lib), you can use this in conjunction with *BEAST* to speed up the calculations. If not installed, unselect the use of the *BEAGLE* library. If the command line version of *BEAST* is being used then the name of the XML file is given after the name of the *BEAST* executable. The analysis will then 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.

## Analyzing 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. Double click on the *Tracer* icon; 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. 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 the posterior (this is the log of the product of the tree likelihood and the prior probabilities), 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** - 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).

- **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 200 and 100 are highlighted in yellow and red respectively by *Tracer*). 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  - as expected given the low ESSs - is extremely rough.
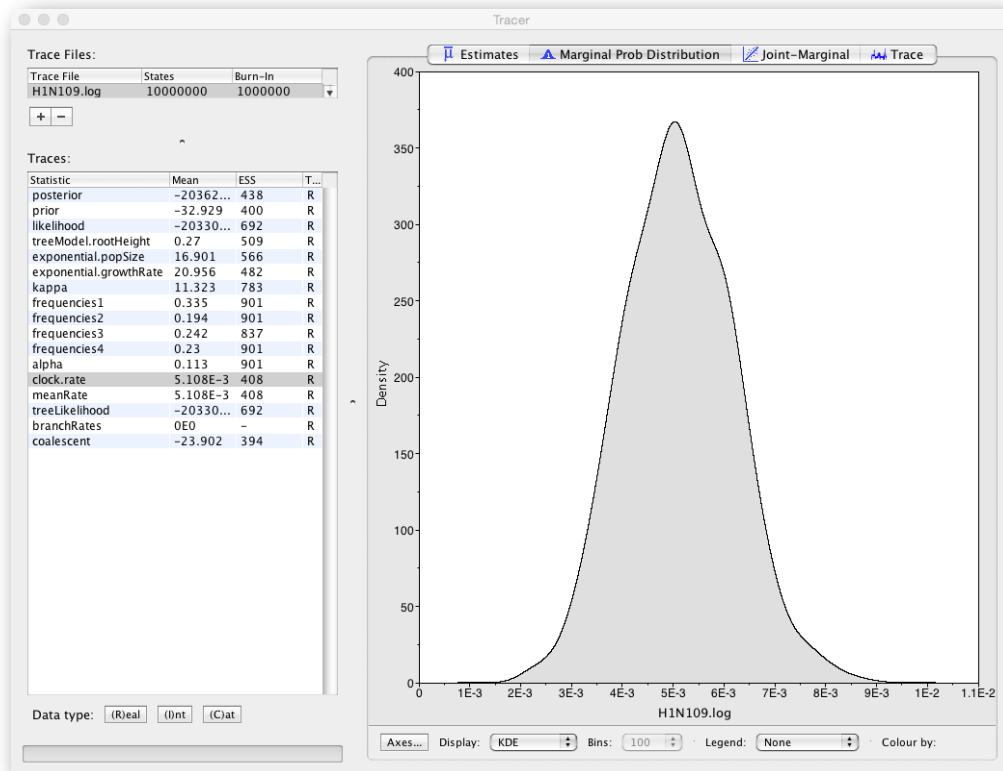
If we select the tab on the right-hand-side labelled `Trace' we can view the raw trace (e.g for treeModel.rootHeight), that is, the sampled values against the  step in the MCMC chain.



Here you can see how the samples are correlated. There are 1000 samples in the trace (we ran the MCMC for 1,000,000 steps sampling every 1000) but it is clear that adjacent samples often tend to have similar values. The ESS for the clock.rate is about 17, so we are only getting 1 independent sample to every 60 actual samples). It also seems that the default burn-in of 10% of the chain length is inadequate (the posterior values are still increasing over the first part of the chain). Not excluding enough of the start of the chain as burn-in will bias the results and render estimates of ESS unreliable.

The analysis needs to be run longer. The lowest ESS of about 3 suggests that we have to run it at least 35 times longer to get ESSs that are >100. However, it would be better to aim higher (e.g. a chain length of 10,000,000 and sampling every 10,000 generations). If the previous analysis ran reasonably fast and if time permits, you can can go back to *BEAUti* and set up and run this longer analysis, but it is probably advisable to proceed with summarizing the longer runs that are provided with this tutorial. Load the new log file into *Tracer* (you can leave the old one loaded for comparison). Click on the Trace tab and look at the raw trace plot.

Again we have chosen options that produce 1000 samples and with an ESS of > 400 for the coalescent model parameters there is little auto-correlation between the samples. There are no obvious trends in the plot which would suggest that the MCMC has not yet converged, and there are no large-scale fluctuations in the trace which would suggest poor mixing. As we are satisfied with the behavior of the MCMC we can now move on to one of the parameters of interest: substitution rate. Select clock.rate in the left-hand table. This is the average substitution rate across all sites in the alignment. Now choose the density plot by selecting the tab labeled 'Marginal Prob Distribution'. This shows a plot of the posterior probability density of this parameter. You should see a plot similar to this:
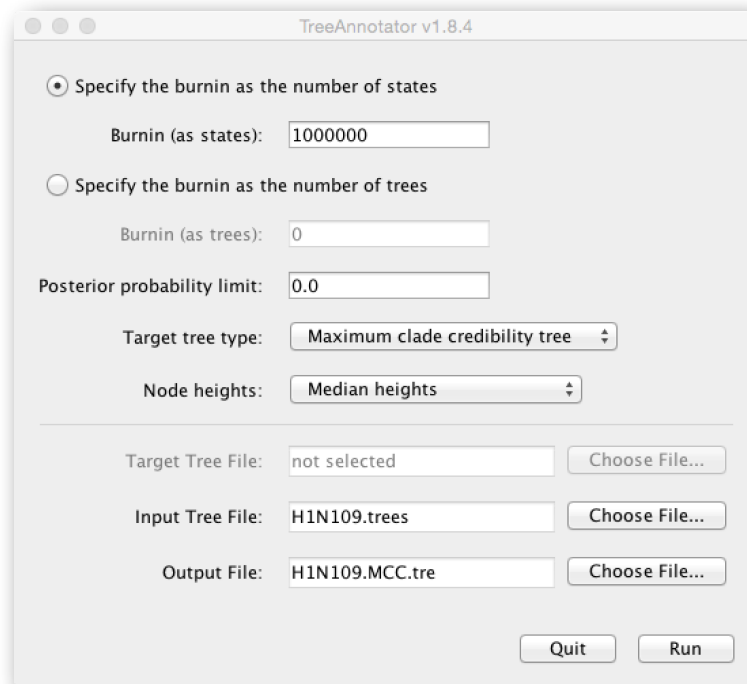
As you can see the posterior probability density is roughly bell-shaped. When looking at the equivalent histogram in the Estimates panel, there is some sampling noise which is smoothened by the KDE; this would be reduced if we ran the chain for longer but we already have a reasonable estimate of the mean and HPD interval. The **treeModel.rootHeight** parameter provides an estimate of the time to the most recent common ancestor since the most recent sampling data (in our case: 2009.403). What would be the mean estimate for the date of the MRCA?

The **exponential.growthRate** (r) provides an estimate of the epidemic growth of H1N1/09. Given that $N_t = N_0\, e^{-rt}$ (with $N_0$ being the population size at present), the doubling time for $r = 21$ is about 0.03 years or 12 days. Interestingly, it has been shown that the basic reproductive ratio ($R_0$) is related to the growth rate (see http://tree.bio.ed.ac.uk/wiki/pages/t769F5D1/Relationship_between_R0_and_the_epidemic_growth_rate.html). However, the basic reproductive number is dependent not just on an estimate of *r*, but also a good estimate of the generation time distribution, which reflects the time between successive infections in a chain of transmission. If we assume a generation time distribution that follows the gamma distribution, then $R_0 = (1 + r / b)^a$, where *a* and *b* are the parameters of the gamma distribution (and $a = \mu^2 / \sigma^2$, $b = \mu / \sigma^2$). Taking $\mu = 3$ days and $\sigma = 2$ days, what would be the mean estimate for the H1N1/09 $R_0$?

## 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 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.
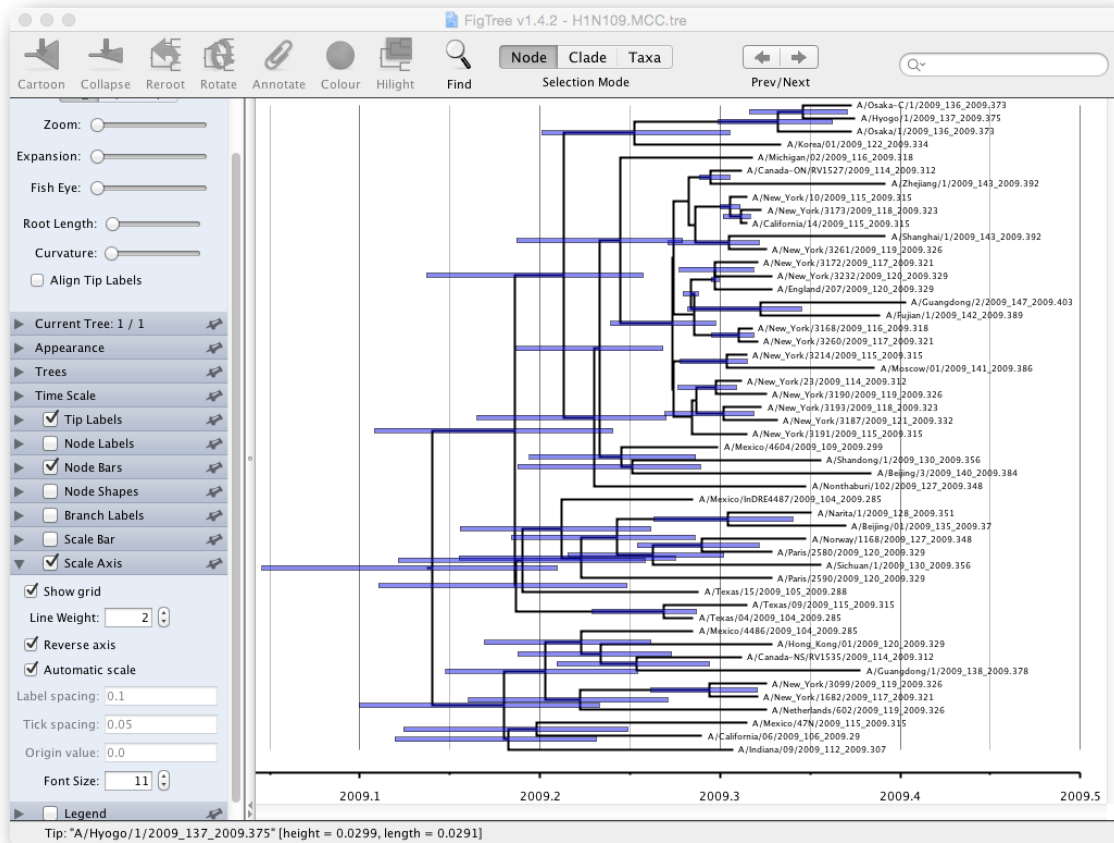
*TreeAnnotator* takes a single `target' tree and annotates it with the summarized information from the entire sample of trees. The summarized information includes the average node ages (along with the HPD intervals), the posterior support and the average rate of evolution on each branch (for relaxed clock 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 steps in the MCMC chain, Burnin (as states), or the number of trees, Burnin (as trees), that should be excluded from the summarization. For the example above, with a chain of 10,000,000 steps, a 10% burnin corresponds to 1,000,000 steps. Alternatively, sampling every 10,000 steps results in 1000 trees in the file, and to obtain at the same 10% burnin, the number of trees needs to be set to 100.

- **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 all nodes will have information summarized independent of their posterior probability.

- **Target tree type** - This has three options 'Maximum clade credibility tree' or 'User target tree' For the latter option, a NEXUS tree file can be specified as the Target Tree File, below. Select the first option, *TreeAnnotator* will examine every tree in the Input Tree File and select the tree that has the highest product 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. Node heights can also be summarised as a Mean or a Median over the sample of trees. Sometimes a mean or median height for a node may actually be higher than the mean or median height of its parental node (because particular ancestral-descendent relationships in the MCC tree may still be different compared to a large number of other tree sampled). This will result in artifactual negative branch lengths, but can be avoided by the 'Common Ancestor heights' option. Let's use the default Median heights for our summary tree.

- **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., H1N109.MCC.tre).

Once you have selected all the options, above, press the 'Run' button. *TreeAnnotator* will analyse 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* now 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, re-order the node order by Increasing Node Order under the Tree Menu. Click on **Branch Labels** in the control panel on the left and open its section by clicking on the arrow on the left. Now select **posterior** under the **Display option**.

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). Select **Node Bars** in the control panel and open this section; select 'height_95%_HPD' to display the 95% HPDs of the node heights. 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 2009.403, the scale factor to -1.0. and '**Reverse Axis**' under '**Scale Axis**'.

Finally, open the **Appearance** panel and alter the **Line Weight** to draw the tree with thicker lines. 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 most of your settings so that when you load it into FigTree again it will be displayed almost exactly as you selected. The tree can also be exported to a graphics file (pdf, eps, etc.).
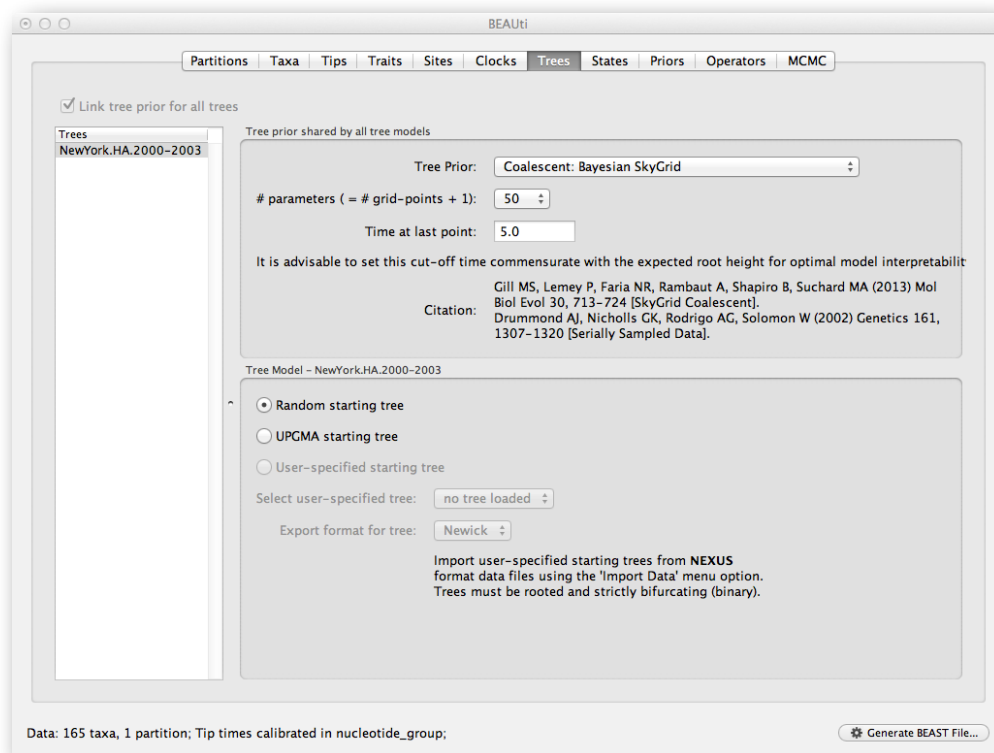
# A quick how-to summary for Exercise 1

- Run **BEAUti**.

    - Load a NEXUS format alignment by selecting the Import Data... option from the File menu. Select the file called H1N109.nex.

    - In the Tips tab, select the box labelled Use tip dates and click the Guess Dates button. Keep the default Defined just by its order and select last from the drop-down menu for the order and press OK.

    - In the Sites tab, keep the default HKY substitution model for the nucleotide data and base frequencies as Estimated. Select 'Gamma' as 'Site Heterogeneity Model' (with 4 discrete categories) before proceeding to the 'Clocks' tab.

    - In the Clocks tab, keep a strict clock model.

    - In the Trees tab, set the option for Tree Prior to Coalescent: Exponential Growth..

    - In the MCMC tab, set the chain length to 1,00,000 and both the sampling frequencies to 100. Set the File name stem to H1N109 and generate the beast file (H1N109.xml).

- Run **BEAST** and load the xml file.

- Analyze the output using **Tracer**. Analyze the output file for the longer runs.

- Calculate the growth rate for pandemic influenza H1N1 (see page 12 of this tutorial).

- Summarize the trees of the longer run using **treeAnnotator** (burn-in = 500,000 states or 100 trees).

- Visualize the tree in **FigTree**.

# EXERCISE 2: reconstructing H3N2 epidemic dynamics in the New York state.

In this exercise, we will reconstruct a *Bayesian skygrid* of H3N2 spread during three epidemic seasons. The data set, **NewYork.HA.2000-2003.nex**, contains 165 Hemagglutinin genes and takes more time to run in *BEAST* than available during a practical session. Therefore, this tutorial will discuss how to set up this analysis and how to summarize the results based on runs that have already been performed.
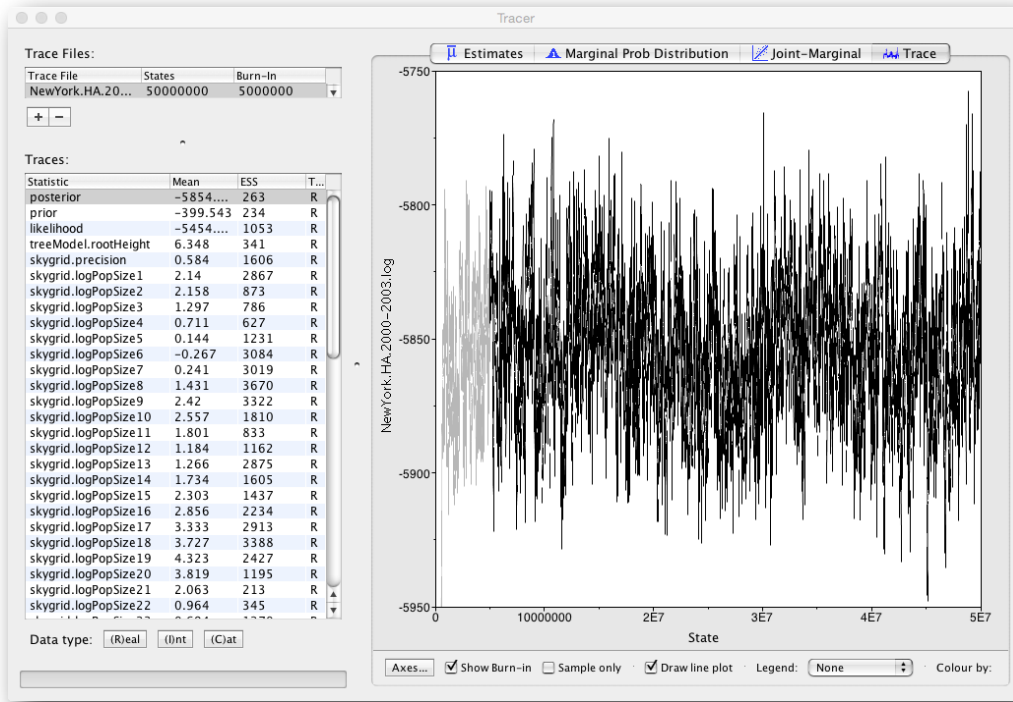
## Running BEAUti

Run *BEAUti*, load the nexus file (**NewYork.HA.2000-2003.nex**) and set the dates to the last numerical field in the taxa names as previously. Set the same evolutionary model (including gamma distributed rate variation) and clock model as in the previous exercise. In the 'Trees' tab, select a 'Coalescent: Bayesian SkyGrid' as the 'Tree Prior'. We will construct a grid of 50 intervals over 5 years (**Time at point**) prior to the most recent sampling dates (2003.98 in our case, so going back to about 1999) requiring us to estimate 10 population sizes by year.



## Analyzing the BEAST output

Using *Tracer*, we can analyze the run based on the output files provided (load the file called 'NewYork.HA. 2000-2003.log'):

To reconstruct the Bayesian skygrid plot, select 'SkyGrid reconstruction...' under the Analysis window. The following window should appear:



Set the manual bin range from 1999 to 2004 and specify '2003.98' as the 'Age of the youngest tip' at the bottom. Press 'OK' and after some time, the following Bayesian skyGrid reconstruction should appear (with solid interval selected):

Output files for a Bayesian skyline plot analysis are also provided for comparison. To reconstruct a Bayesian skyline plot based on these, select 'Bayesian Skyline reconstruction' under the Analysis window.

## Some Questions

• What type of dynamics does the H3N2 skyride plot suggest? Would you expect to see the similar dynamics for H3N2 sampled in a southern hemisphere location?

• Is the H1N1/09 evolutionary rate similar to the H3N2 evolutionary rate? If not, what could explain their differences.

• Based on the H1N1/09 tree inferred from a limited sampling, how many H1N1/09 introductions in New York would you conclude for this sample?

## A quick how-to summary

• Run **BEAUti**.

  - Load a NEXUS format alignment by selecting the Import Data... option from the File menu. Select the file called NewYork.HA.2000-2003.nex.

  - In the Tips tab, select the box labelled Use tip dates and click the Guess Dates button. Keep the default Defined just by its order and select last from the drop-down menu for the order and press OK.

  - In the Sites tab, keep the default HKY substitution model for the nucleotide data and base frequencies as Estimated. Select 'Gamma' as 'Site Heterogeneity Model' (with 4 discrete categories) before proceeding to the 'Clocks' tab.

  - In the Clocks tab, keep a strict clock model.

  - In the Trees tab, set the option for Tree Prior to Coalescent: Bayesian SkyGrid with 50 parameters and a 5 year cut-off (Time at last point).

- In the **Priors** tab, keep the default prior settings.

- In the **MCMC** tab, set the chain length to 25,000,000 and both the sampling frequencies to 5000. Set the **File name stem** to NewYork.HA.2000-2003 and generate the beast file (**NewYork.HA.2000-2003.xml**).

• Run **BEAST** and load the xml file.

• Analyze the output using **Tracer**. Analyze the output file for the longer runs.

• Reconstruct a Bayesian skyride plot by selecting 'SkyGrid reconstruction...' under the **Analysis** window (bin range = 1999-2004 and age of the youngest tip = 2003.98).

## References

• Drummond AJ, Rambaut A (2007) *BEAST*: Bayesian evolutionary analysis by sampling trees. *BMC Evolutionary Biology* **7**: 214.

• Drummond AJ, Ho SYW, Phillips MJ & Rambaut A (2006) *PLoS Biology* 4, e88.

• Drummond AJ, Rambaut A & Shapiro B and Pybus OG (2005) *Mol Biol Evol* **22**, 1185-1192.

• Drummond AJ, Nicholls GK, Rodrigo AG & Solomon W (2002) *Genetics* **161**, 1307-1320.

• Ferreira, M. A. R. and M. A. Suchard. 2008. Bayesian analysis of elapsed times in continuous-time Markov chains. Can J Statistics, 36: 355–368. doi: 10.1002/cjs.5550360302

• Gill MS, Lemey P, Faria NR, Rambaut A, Shapiro B, and Suchard MA (2013) Improving Bayesian population dynamics inference: a coalescent-based model for multiple loci. *Mol Biol Evol* **30**, 713-724.

• Minin VN, Bloomquist EW and Suchard MA (2008) Smooth Skyride through a Rough Skyline: Bayesian Coalescent-Based Inference of Population Dynamics. *Molecular Biology and Evolution* **25**:1459-1471; doi:10.1093/molbev/msn090.

• Rambaut A, Pybus OG, Nelson MI, Viboud C, Taubenberger JK, Holmes EC (2008) The genomic and epidemiological dynamics of human influenza A virus. *Nature*, **453**: 615-9.

• Smith GJD, Vijaykrishna D, Bahl J, Lycett SJ, Worobey M, Pybus OG, Ma SK, Cheung CL, Raghwani J, Bhatt S, Peiris JSM, Guan Y & Rambaut A (2009) Origins and evolutionary genomics of the 2009 swine-origin H1N1 influenza A epidemic. *Nature* **459**, 1122-1125.

## Help and documentation

• The BEAST software download: http://beast-mcmc.googlecode.com/

• The BEAST website: http://beast.bio.ed.ac.uk/

• Tutorials: http://beast.bio.ed.ac.uk/Tutorials/

• Frequently asked questions: http://beast.bio.ed.ac.uk/FAQ/

• H1N1/09: http://tree.bio.ed.ac.uk/wiki/projects/influenza/