

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 recent outbreak of swine-origin influenza A (H1N1) virus (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 of 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 v1.7.2. It is available for download from <http://beast.bio.ed.ac.uk/>.
- **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 v1.5. 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 v1.3.1. It is available for download from <http://tree.bio.ed.ac.uk/>.

EXERCISE 1: The swine-origin influenza A outbreak

Running *BEAUti*

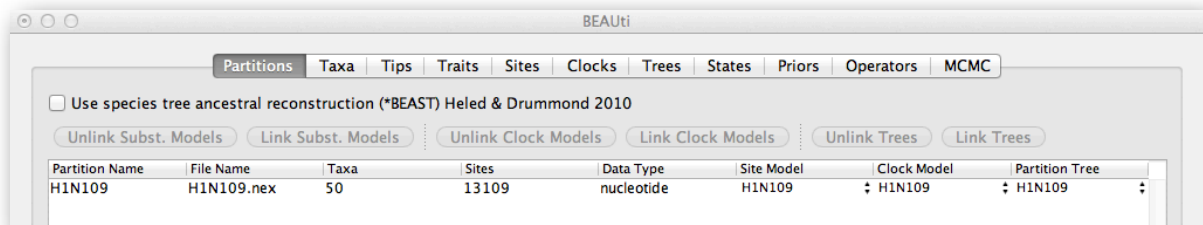
The program *BEAUti* is a user-friendly program for setting the model parameters for *BEAST*. Run *BEAUti* by double clicking on its icon.

Loading the NEXUS file

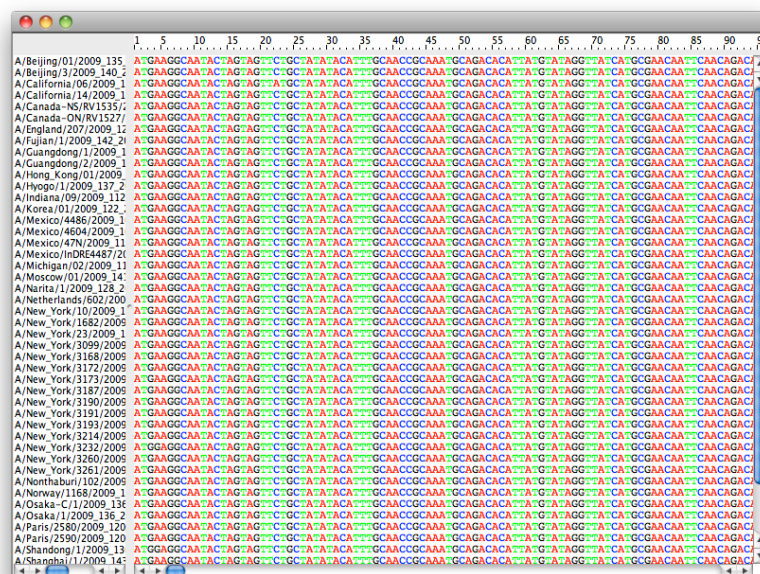
To load a NEXUS format alignment, simply select the Import NEXUS... option from the File menu.

The NEXUS alignment

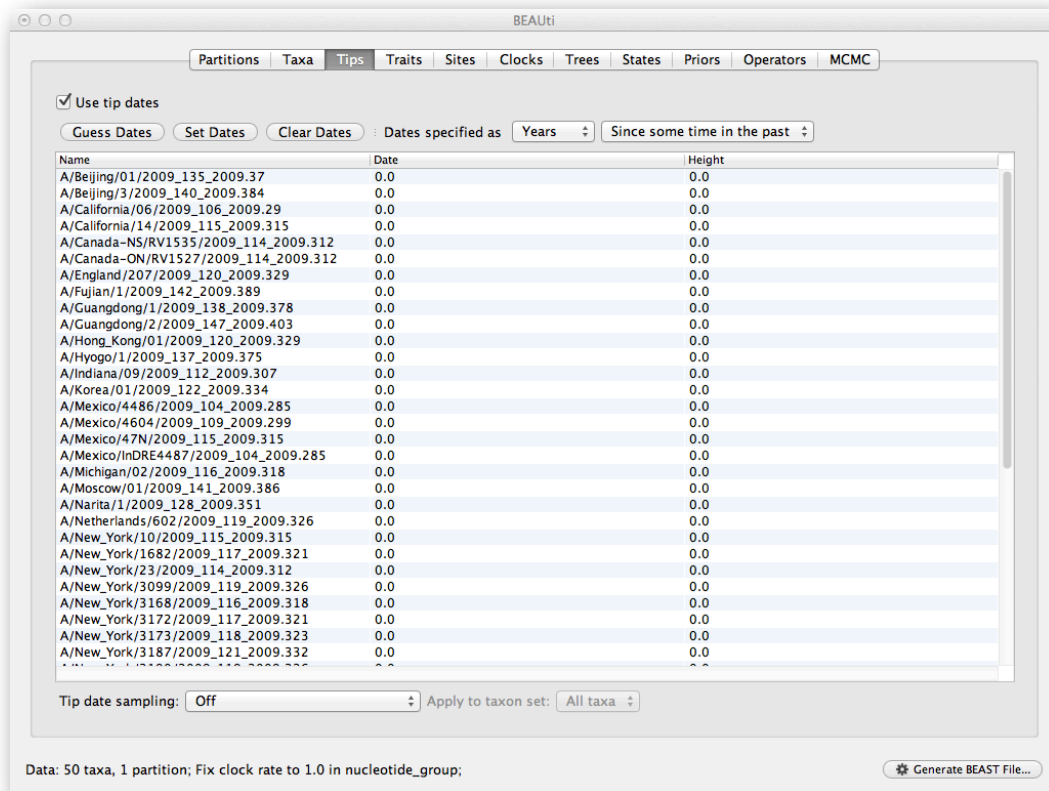
Select the file called H1N109.nex. This file contains an alignment of 50 genomes (concatenated segments), 13109 nucleotides in length. Once loaded, the new data will be listed in the table as shown in the figure:



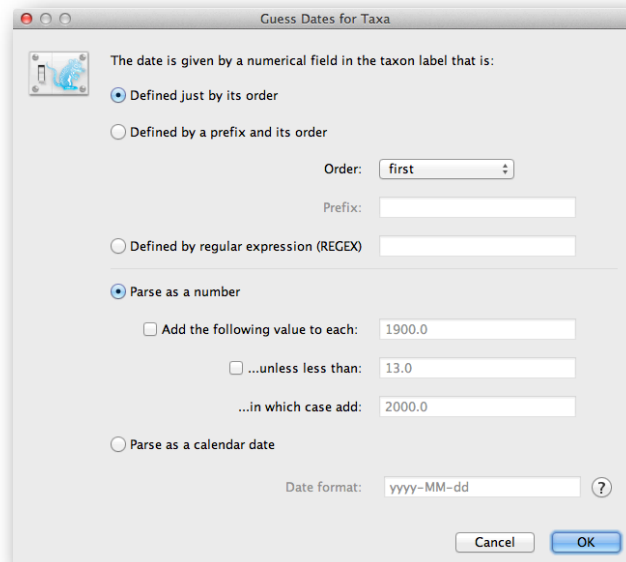
Double-click on the row of the table to display the actual sequence alignment:



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:

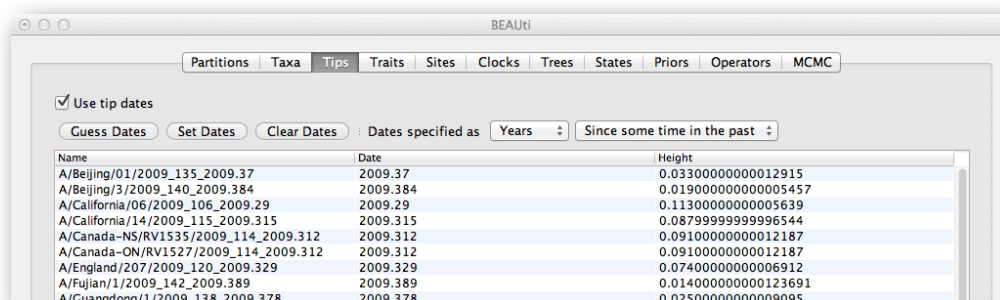


Select the box labelled 'Use tip dates'. The actual sampling in fractional years is encoded in the name of each taxon and we could simply edit the value in the 'Date' column of the table to reflect these. However, if the taxa names contain the calibration information, then a convenient way to specify 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 (such as the some YFV sequences, above) 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). For the YFV sequences you can keep the default 'Defined just by its order' and 'Order: first'.

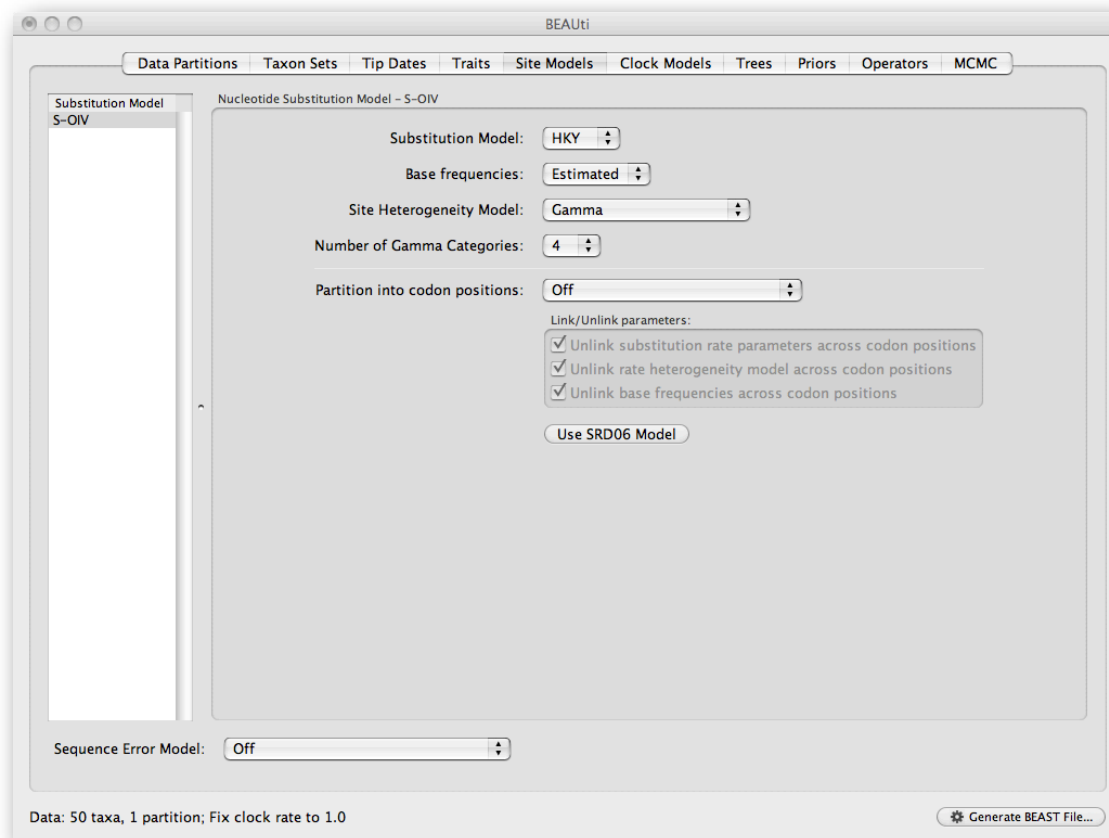
When parsing a number, you can ask BEAUti to add a fixed value to each guessed date. For example, the value "1900" can be added to turn the dates from 2 digit years to 4 digit. Any dates in the taxon names given as "00" would thus become "1900". However, if these "00" or "01", etc. represent sequences 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. 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).



Name	Date	Height
A/Beijing/01/2009_135_2009.37	2009.37	0.03300000000012915
A/Beijing/3/2009_140_2009.384	2009.384	0.019000000000005457
A/California/06/2009_106_2009.29	2009.29	0.113000000000005639
A/California/14/2009_115_2009.315	2009.315	0.087999999999996544
A/Canada-NS/RV1535/2009_114_2009.312	2009.312	0.09100000000012187
A/Canada-ON/RV1527/2009_114_2009.312	2009.312	0.09100000000012187
A/England/207/2009_120_2009.329	2009.329	0.074000000000006912
A/Fujian/1/2009_142_2009.389	2009.389	0.014000000000123691
A/Guangdong/1/2009_138_2009.378	2009.378	0.025000000000009095

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.



Substitution Model: HKY

Base frequencies: Estimated

Site Heterogeneity Model: Gamma

Number of Gamma Categories: 4

Partition into codon positions: Off

Link/Unlink parameters:

- ☒ Unlink substitution rate parameters across codon positions
- ☒ Unlink rate heterogeneity model across codon positions
- ☒ Unlink base frequencies across codon positions

Sequence Error Model: Off

Data: 50 taxa, 1 partition; Fix clock rate to 1.0

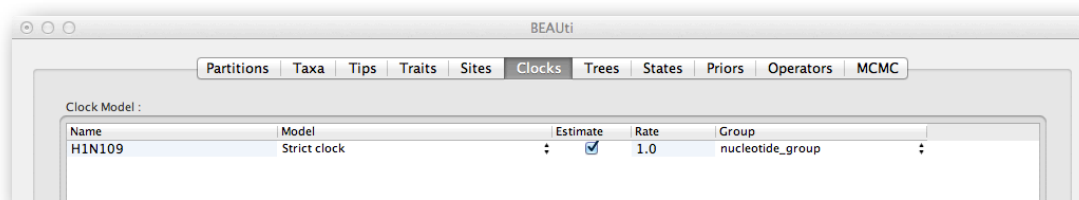
This tutorial assumes that you are familiar with the evolutionary models available, however there are a couple of points to note about selecting a model in *BEAUti*:

- 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**' 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.



Keep the default option to '**Estimate**' so that the rate of molecular evolution is estimated from the temporal sampling of the sampled viruses.

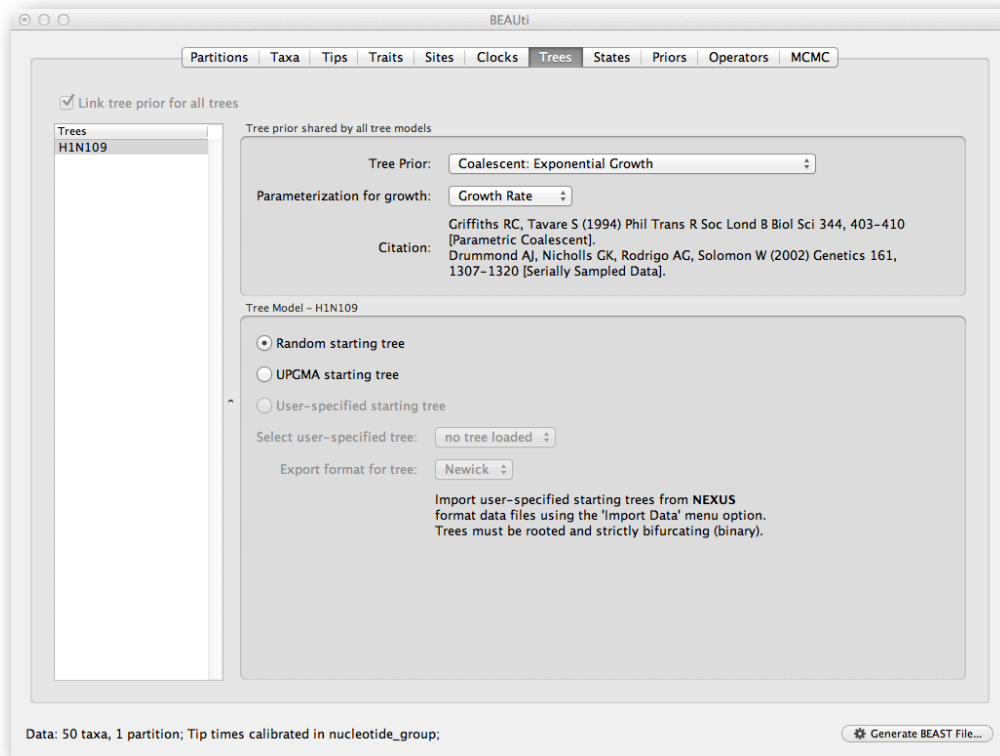
If there are no dates for the sequences (they are contemporaneous) then you can specify a fixed mean substitution rate obtained from another source. Setting this to 1.0 will result in the ages of the nodes of the tree being estimated in units of substitutions per site (i.e. the normal units of branch lengths in popular packages such as *MrBayes*).

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.

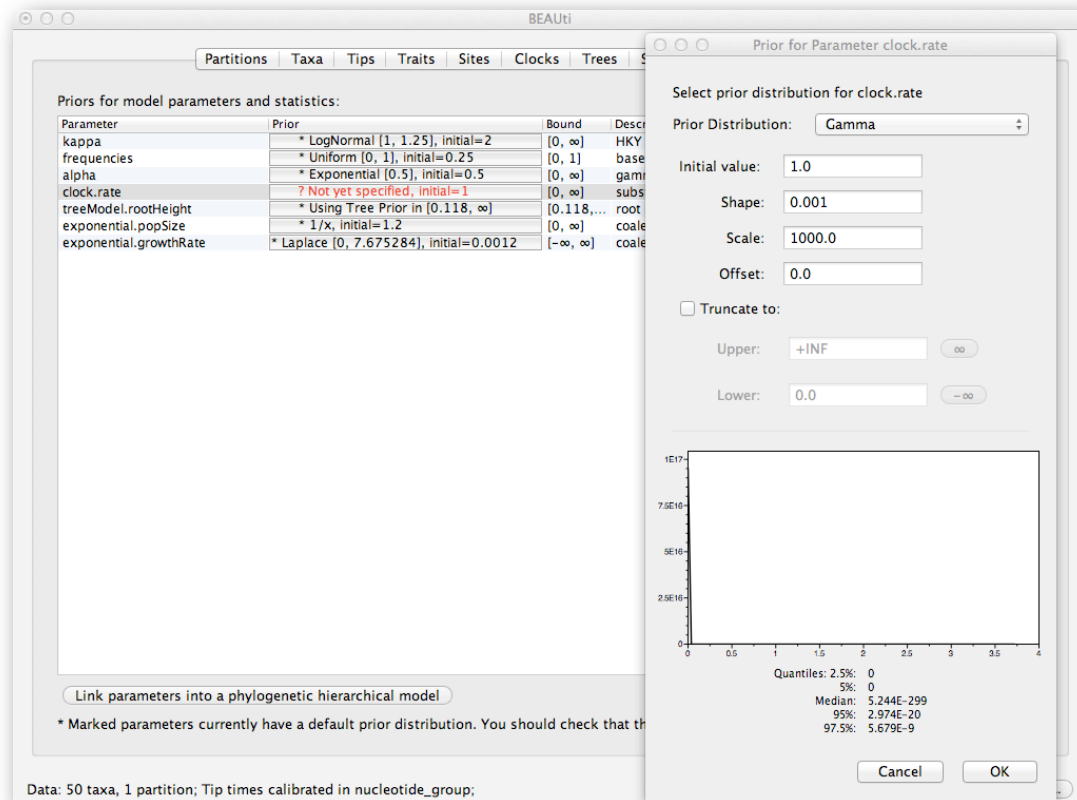
To estimate the epidemic growth rate, we will change this demographic model to an exponential growth coalescent prior, which is intuitively appealing for viral outbreaks. Switch the option for '**Tree Prior**' to '**Coalescent: Exponential Growth**'.



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 marginal prior distribution is for each. A marginal prior allows the user to 'inform' the analysis by selecting a particular distribution. Although some of the default marginal priors are improper (e.g. indicated in yellow), with sufficiently informative data the posterior becomes proper. Priors that have not been set yet appear in red (e.g. clock.rate). Click on the prior for this parameter 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".

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.



Setting up the operators

Each parameter in the model has one or more “operators” (these are variously called moves and proposals 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:

Partitions Taxa Tips Traits Sites Clocks Trees States Priors Operators MCMC					
<input checked="" type="checkbox"/> Auto Optimize					
In use	Operates on	Type	Tuning	Weight	Description
<input checked="" type="checkbox"/>	kappa	scale	0.75	0.1	HKY transition-transversion parameter
<input checked="" type="checkbox"/>	frequencies	deltaExchange	0.01	0.1	frequencies
<input checked="" type="checkbox"/>	alpha	scale	0.75	0.1	gamma shape parameter
<input checked="" type="checkbox"/>	clock.rate	scale	0.75	3.0	substitution rate
<input checked="" type="checkbox"/>	Tree	subtreeSlide	0.12	15.0	Performs the subtree-slide rearrangement of the tree
<input checked="" type="checkbox"/>	Tree	narrowExchange	n/a	15.0	Performs local rearrangements of the tree
<input checked="" type="checkbox"/>	Tree	wideExchange	n/a	3.0	Performs global rearrangements of the tree
<input checked="" type="checkbox"/>	Tree	wilsonBalding	n/a	3.0	Performs the Wilson-Balding rearrangement of the tree
<input checked="" type="checkbox"/>	treeModel.rootHeight	scale	0.75	3.0	root height of the tree
<input checked="" type="checkbox"/>	Internal node heights	uniform	n/a	30.0	Draws new internal node heights uniformly
<input checked="" type="checkbox"/>	exponential.popSize	scale	0.75	3.0	coalescent population size parameter
<input checked="" type="checkbox"/>	exponential.growthRate	randomWalk	1.0	3.0	exponential.growthRate
<input checked="" type="checkbox"/>	Substitution rate and heights	upDown	0.75	3.0	Scales substitution rates inversely to node heights of the tree

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.

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 screenshot shows the 'MCMC' tab in the BEAUti software interface. The settings are as follows:

- Length of chain:** 10000000
- Echo state to screen every:** 1000
- Log parameters every:** 1000
- File name stem:** H1N109
- ☐ Add .txt suffix
- Log file name:** H1N109.log
- Trees file name:** H1N109.trees
- ☐ Create tree log file with branch length in substitutions:
- Substitutions trees file name:** (empty)
- ☒ Create operator analysis file:
- Operator analysis file name:** H1N109.ops
- ☐ Sample from prior only - create empty alignment

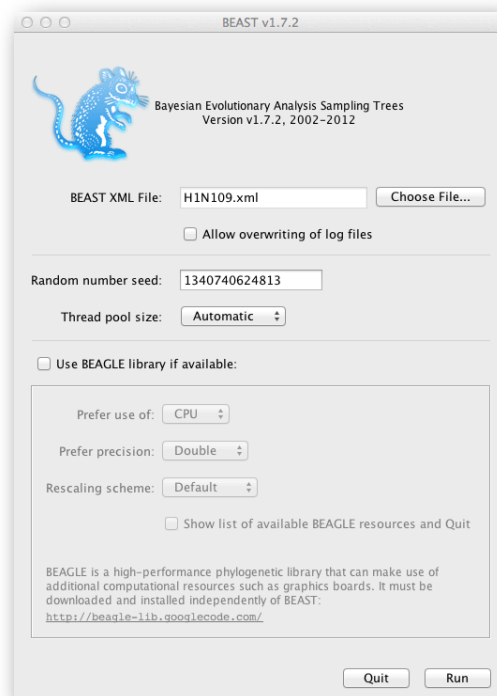
The next couple of options specify how often the current parameter values should be displayed on the screen and recorded in the log file. The screen output is simply for monitoring the program's progress so can be set to any value (although if set too small, the sheer quantity of information being displayed 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. We can keep the default settings for our analyses

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. 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 (H1N109.xml) and save the file.

For convenience, leave the BEAUi window open so that you can change settings 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 dialog box will appear in which you select the XML file:

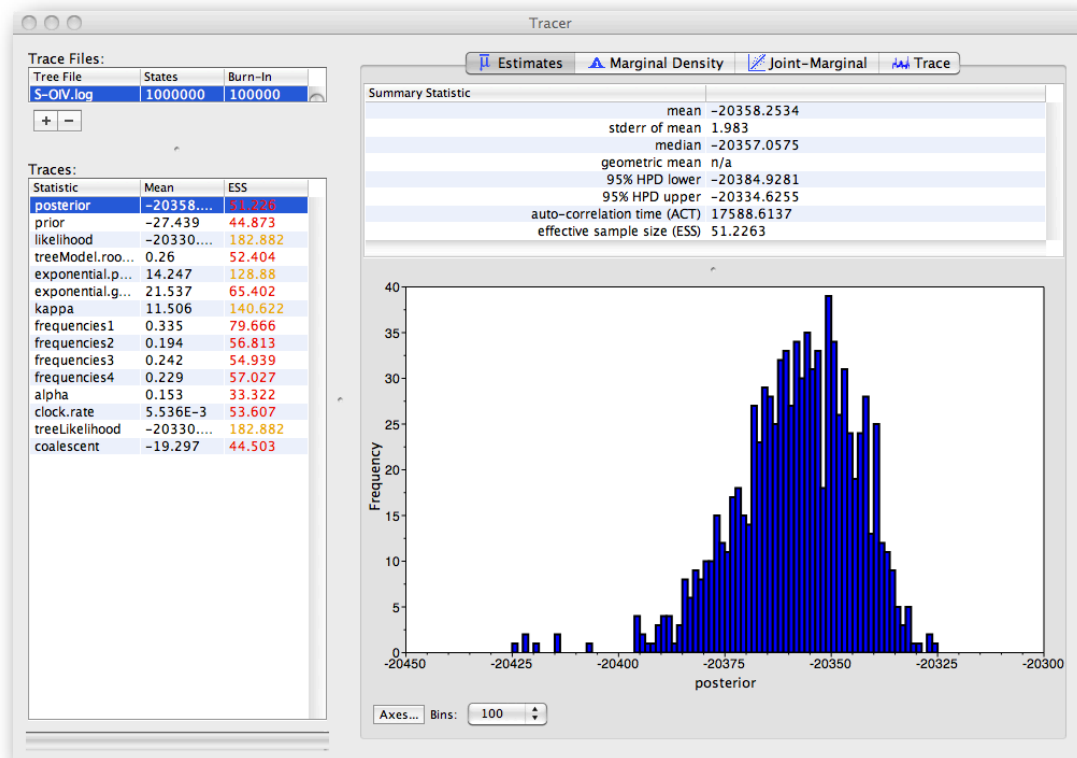


Press the 'Choose File' button and select the XML file you just created and press 'Run'. If the command line version is being used then the name of the XML file is given after the name of the *BEAST* executable. The analysis will then be performed with detailed information about the progress of the run 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.



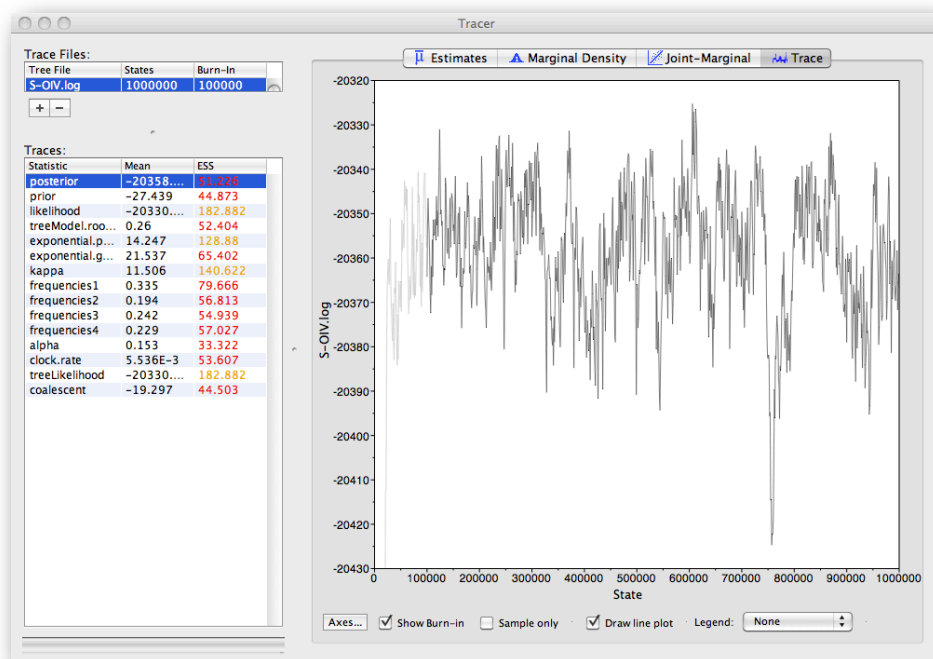
On the left hand side is the name of the log file loaded and the traces that it contains. There are traces for the posterior (this is the log of the product of the tree likelihood and the prior probabilities), and the continuous parameters. Selecting a trace on the left brings up analyses for this trace 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).
- **Stddev** - 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.

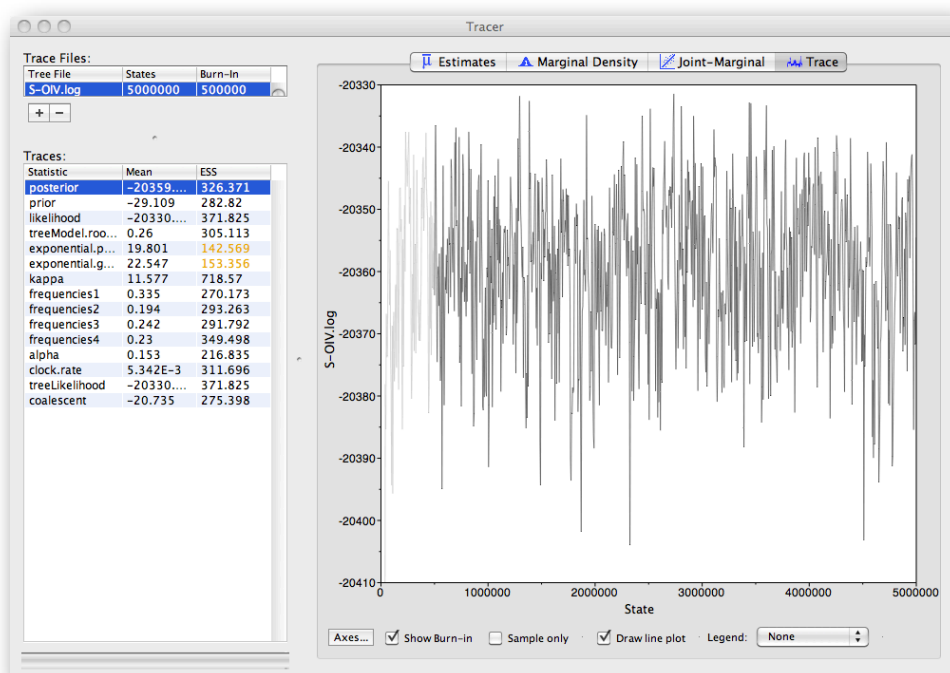
Note that the effective sample sizes (ESSs) for all the traces are small (ESSs less than 100 are highlighted in red by *Tracer*). This is not good. A low ESS means that the trace contained a lot of correlated samples and thus may not represent the posterior distribution well. In the bottom right of the 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, that is, the sampled values against the step in the MCMC chain.



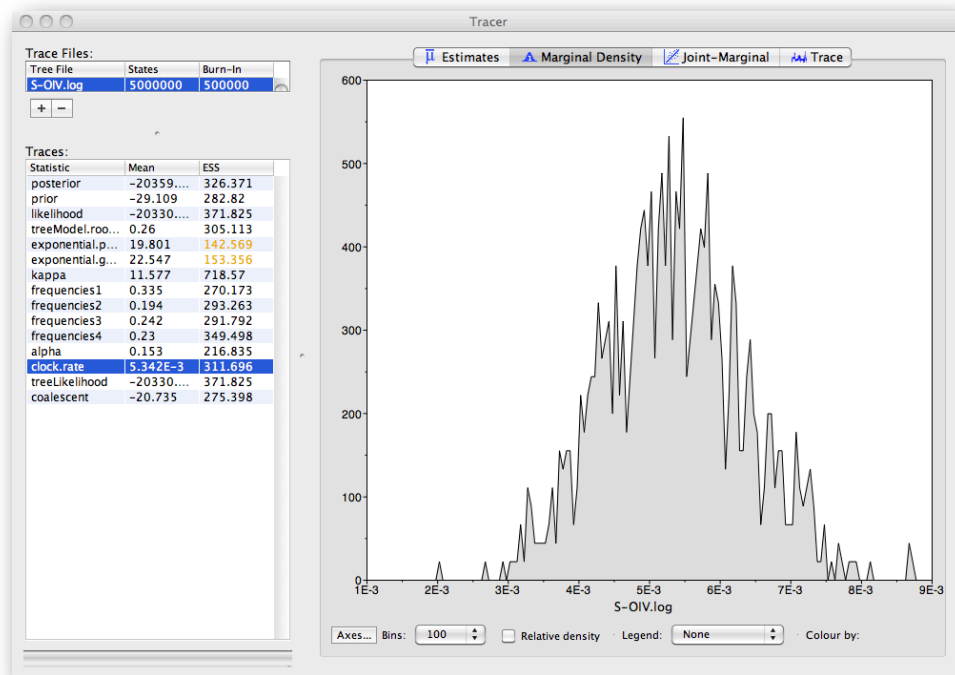
Here you can see how the samples are correlated. There are 1000 samples in the trace (we ran the MCMC for 1,000,000 steps sampling every 1000) but it is clear that adjacent samples often tend to have similar values. The ESS for the clock.rate is about 53 so we are only getting 1 independent sample to every 18 actual samples). Despite this behavior, it seems like there is no particular problem with the burn-in.

The analysis needs to be run longer. The lowest ESS of about 30 suggests that we have to run it at least 3 times longer to get ESSs that are >100. However, it would be better to aim higher (e.g. a chain length of 5,000,000 and sampling every 5000 generations). If the previous analysis ran reasonably fast and if time permits, you 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 about 300 there is still auto-correlation between the samples but 300 effectively independent samples will now provide a reasonable estimate of the posterior distribution. *Tracer* still highlights the ESSs < 200 in yellow, so these can still be improved. Poor prior choices for these parameters are likely the cause of poor performance. Fortunately, 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 'Density'. This shows a plot of the posterior probability density of this parameter. You should see a plot similar to this:

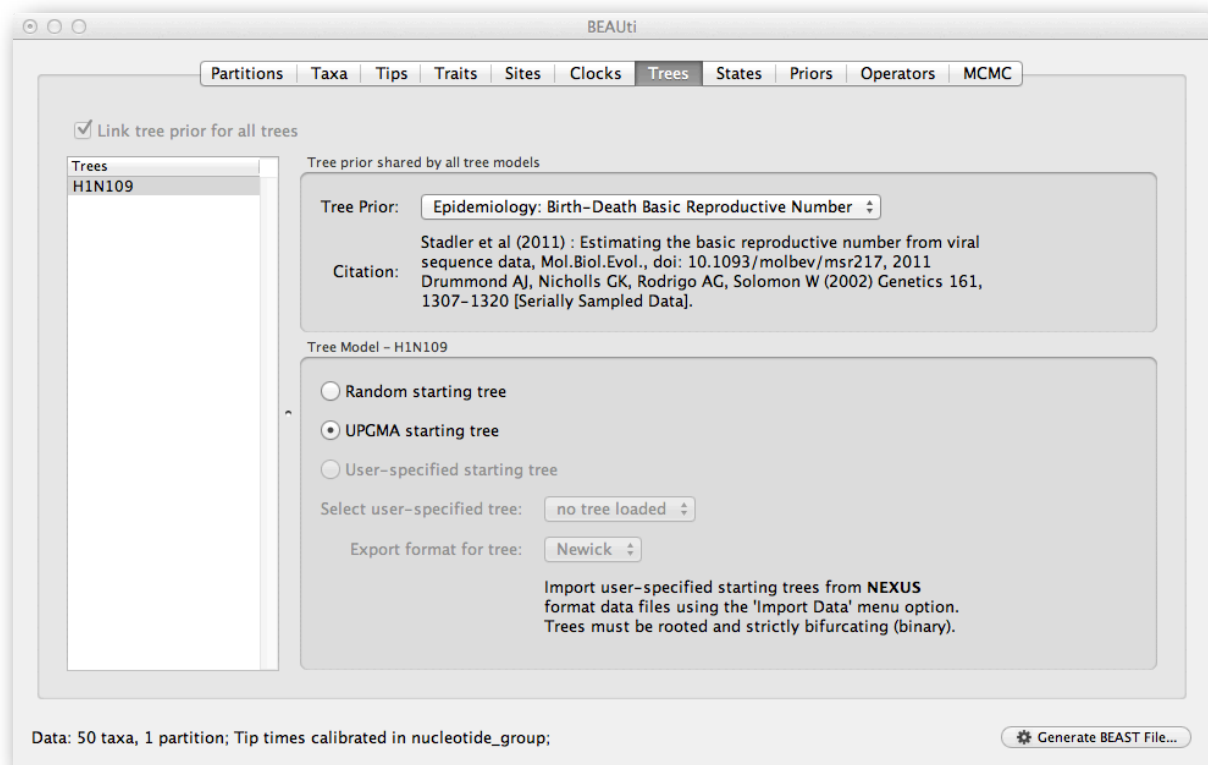


As you can see the posterior probability density is roughly bell-shaped. There is some sampling noise which would be reduced if we ran the chain for longer but we already have a reasonable estimate of the mean and HPD interval. The **treeModel.rootHeight** parameter provides an estimate of the time to 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 $r = 22.5$ is about 0.03 years or 11 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/groups/influenza/wiki/d6504/>). 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 $R_0 = (1 + r/b)^{a/b}$, 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 ?

Because there are some limitations to applying the coalescent framework to R_0 estimation, alternatively approaches are currently being explored, such as birth-death models (BDM). A BDM for estimating R_0 has recently been implemented in BEAST (Stadler, et al., MBE, 2012) and this parametrizes the infection process using the time of the origin of a new epidemic (x_0), R_0 , a recovery rate (μ) and a sampling probability (ψ), where R_0 is $\lambda/(\mu+\psi)$ and λ is the transmission rate. The BDM is a forward in time description of the epidemiological process and ψ allows to specify the sampling intensity. To set up this model, go back to **Trees** panel in *BEAUti* and select “Epidemiology: Birth-Death Basic Reproductive Number” instead

of a coalescent prior. Also select a UPGMA starting tree to ensure a reasonable starting value for the most recent common ancestor (MRCA), on which we can rely to set an appropriate starting value for x_0 .



Under the Prior panel, we can change the prior distribution on the origin (x_0) and recovery rate (μ) to a diffuse gamma(0.001,1000) prior and to ensure that the starting value for x_0 is larger than the TMRCA, set this to 10. In the MCMC panel, we can set the chain length to 5,000,000 and sampling frequency to 5000 generations, similar to the exponential growth analysis. Also for this analysis, log and tree files are available, both for the current settings as well as for a run using a beta(1,10) prior on Ψ instead of the default beta(1,1).

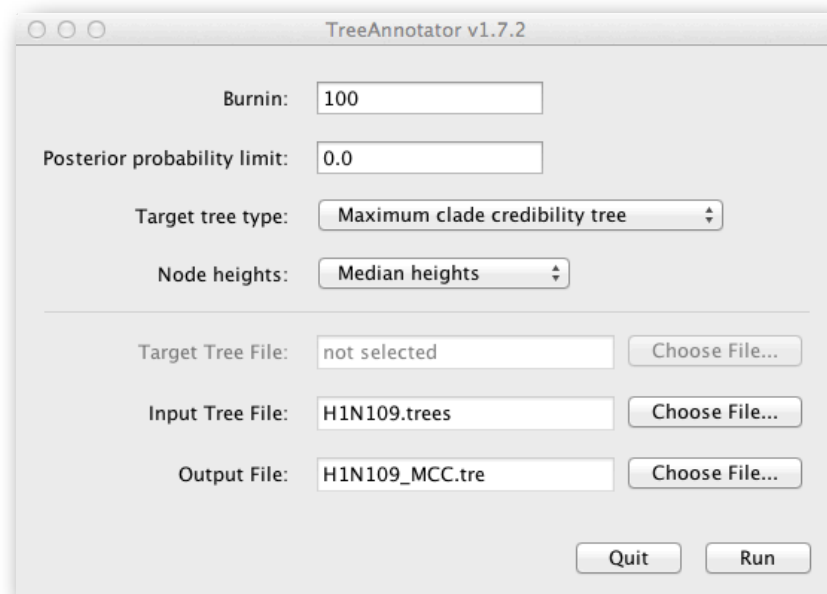
Some Questions

- Are the R_0 estimates affected by the beta prior on Ψ ? How do they relate to the estimate obtained using the growth rate?
- Do we get a reasonable posterior expectation for Ψ for the run with the beta(1,1) prior.
- Do the tree priors affect the MRCA and rate estimates?

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.

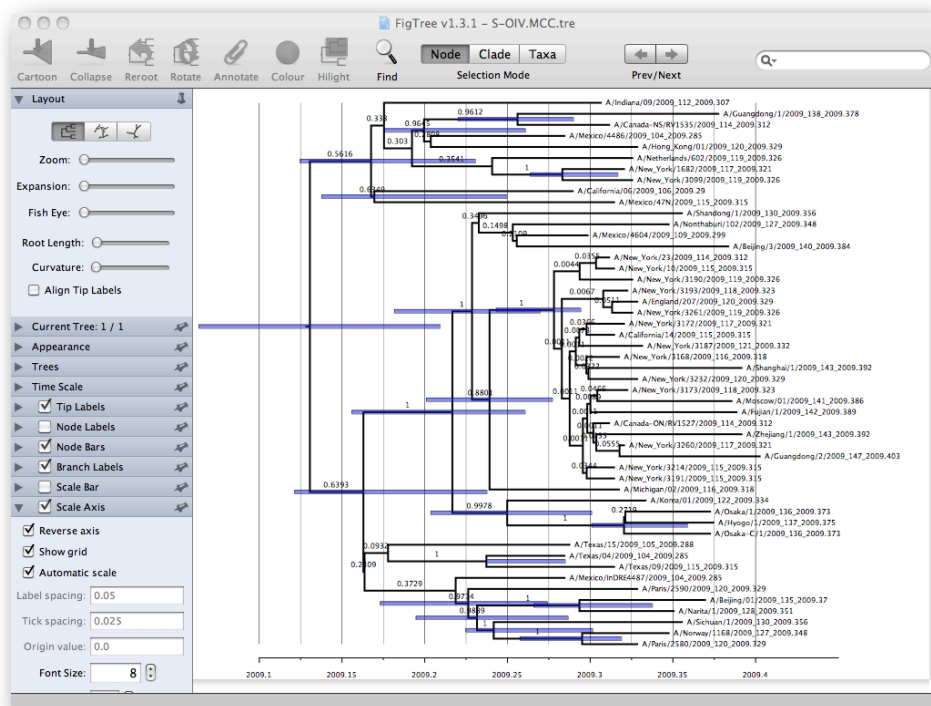
- **Burnin** - This is the number of trees in the input file that should be excluded from the summarization. This value is given as the number of trees rather than the number of steps in the MCMC chain. Thus for the example above, with a chain of 5,000,000 steps, sampling every 5000 steps, there are 1000 trees in the file. To obtain a 10% burn-in, set this value 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 0.5 so only nodes with this posterior probability or greater will have information summarized (the equivalent to the nodes in a majority-rule consensus tree).
- **Target tree type** - This has three options 'Maximum clade credibility tree', 'Maximum sum of clade credibilities' 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. The other two options give node heights as an average (Mean or Median) over the sample of trees. Select 'Mean heights' for our analysis.
- **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 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 2009.403, the scale factor to -1.0. and '**Reverse Axis**' under '**Scale Axis**'.

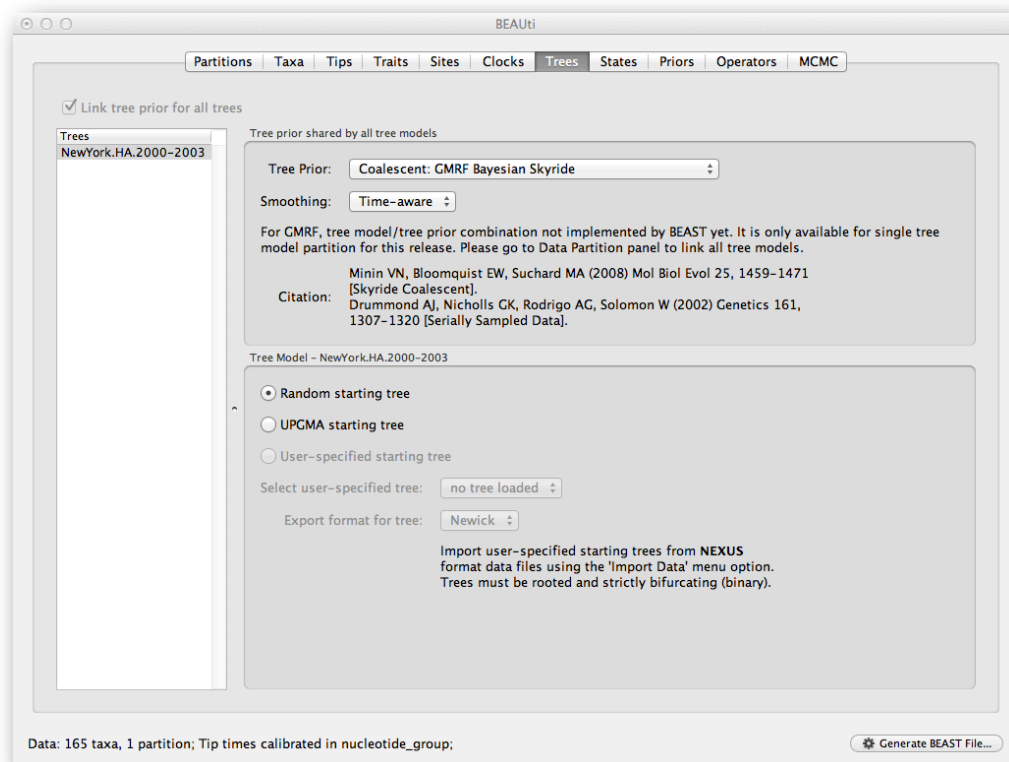
Finally, open the '**Appearance**' panel and alter the '**Line Weight**' to draw the tree with thicker lines. None of the options actually alter the tree's topology or branch lengths in anyway so feel free to explore the options and settings. You can also save the tree and this will save all your settings so that when you load it into *FigTree* again it should be displayed exactly as you selected. The tree can also be exported to a graphics file (pdf, eps, etc.).

EXERCISE 2: reconstructing H3N2 demographics in the New York state.

In this exercise, we will reconstruct a *Bayesian skyride* of H3N2 spread during three epidemic seasons. The data set, **NewYork.HA.2000-2003.nex**, contains 165 Hemagglutinin genes and takes considerable time to run in *BEAST*. 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.

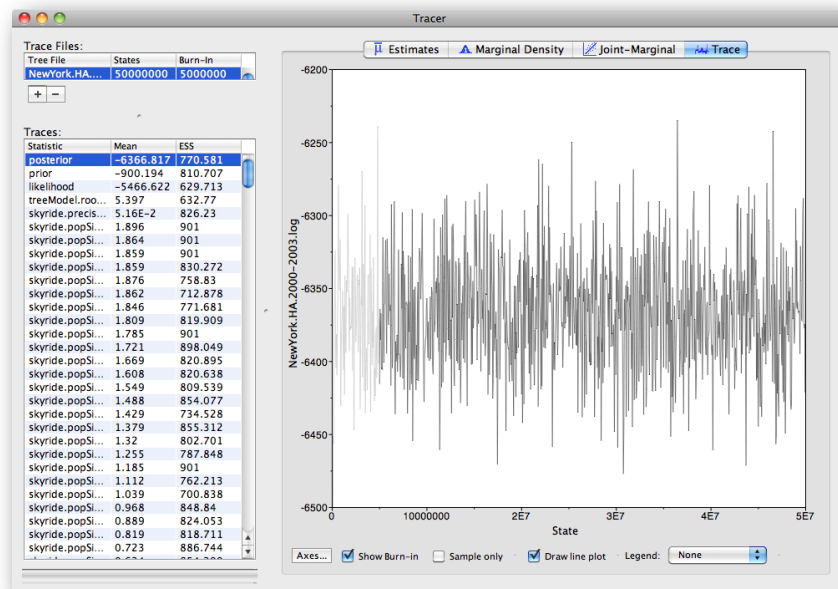
Running BEAUti

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 as in the previous exercise (including gamma distributed rate variation). In the 'Trees' tab, select a 'Coalescent: GMRF Bayesian Skyride' as the 'Tree Prior' with a default 'Time-aware' smoothing.



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 skyride plot, select 'GMRF Skyride reconstruction' under the Analysis window. The following window should appear:

GMRF Skyride Analysis

Warning! This analysis should only be run on traces where the GMRF Skyride model was specified as the demographic in BEAST. Any other model will produce meaningless results.

Trees Log File: NewYork.HA.2000-2003.trees [Choose File...]

Select the traces to use for the arguments:

Population Size: skyride.popSize

Maximum time is the root height's: Lower 95% HPD

Select the trace of the root height: treeModel.rootHeight

Number of bins: 100

☐ Use manual range for bins:

Minimum time: []

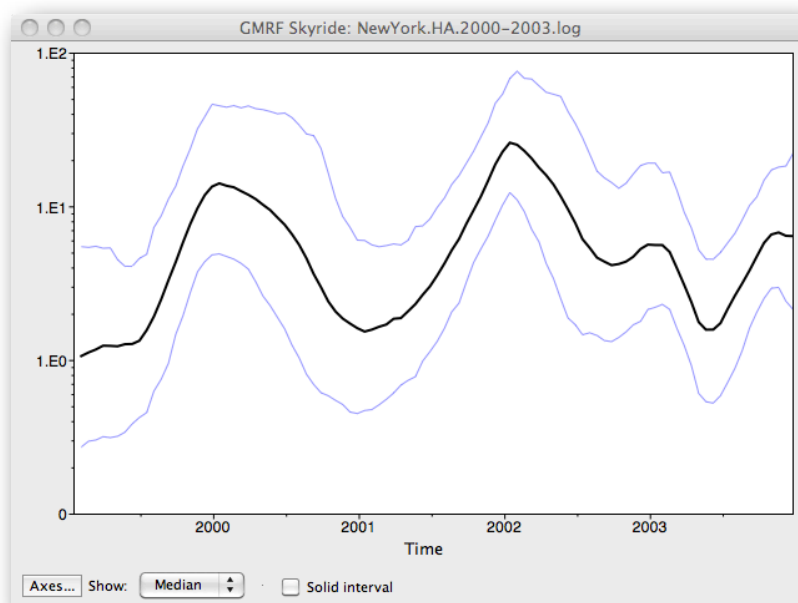
Maximum time: []

Age of youngest tip: 2003.98

You can set the age of sampling of the most recent tip in the tree. If this is set to zero then the plot is shown going backwards in time, otherwise forwards in time.

Buttons: Cancel, OK

Choose the 'NewYork.HA.2000-2003.trees' file as Tree Log File, enter 2003.98 as Age of the youngest tip at the bottom, and press 'OK'. After some time, the following Bayesian skyride reconstruction should appear:



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?

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.
- 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, Raghwan 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 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/groups/influenza/>