

LinguaPhylo

A probabilistic model specification language for reproducible phylogenetic analyses

Alexei Drummond, Kylie Chen, Fabio K Mendes,
Dong Xie

BEAST: Bayesian evolutionary analysis by sampling trees

Alexei J Drummond^{*1,2} and Andrew Rambaut³

BMC Evol Biol, 2007

>12700 citations

Bayesian Phylogenetics with BEAUTi and the BEAST 1.7

Alexei J. Drummond,^{*1,2} Marc A. Suchard,^{*3,4} Dong Xie,^{1,2} and Andrew Rambaut^{*5}

Mol Biol Evol, 2012

>9700 citations

BEAST 2: A Software Platform for Bayesian Evolutionary Analysis

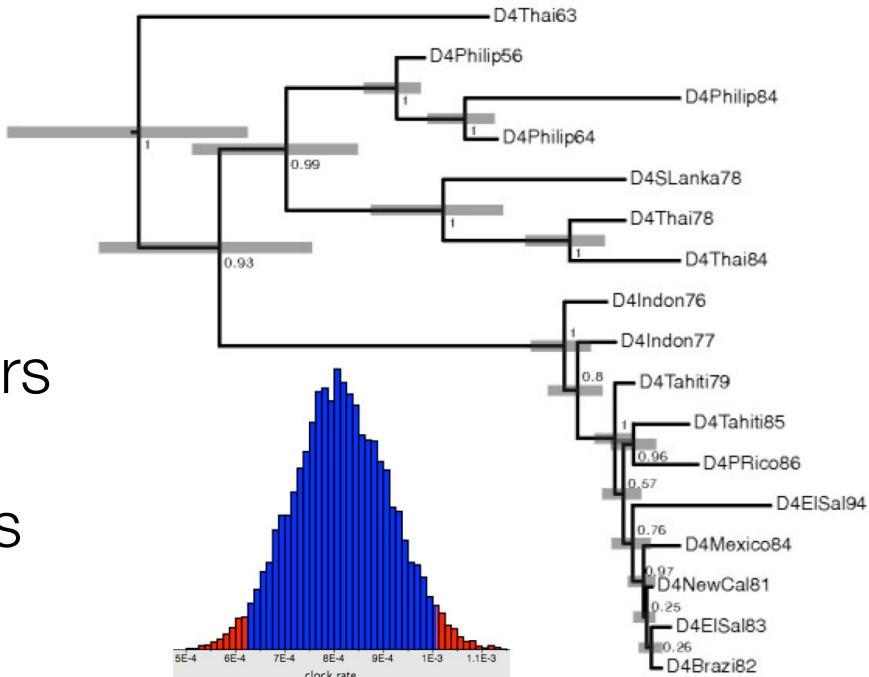
Remco Bouckaert^{1*}, Joseph Heled¹, Denise Kühnert^{1,2}, Tim Vaughan^{1,3}, Chieh-Hsi Wu¹, Dong Xie¹, Marc A. Suchard^{4,5}, Andrew Rambaut⁶, Alexei J. Drummond^{1,7*}

PLoS Comp Bio, 2014

>5100 citations

Estimation of

- evolutionary trees
- divergence times
- evolutionary rates
- evolutionary model parameters
- population history
- infectious disease parameters
- migration rates
- more...



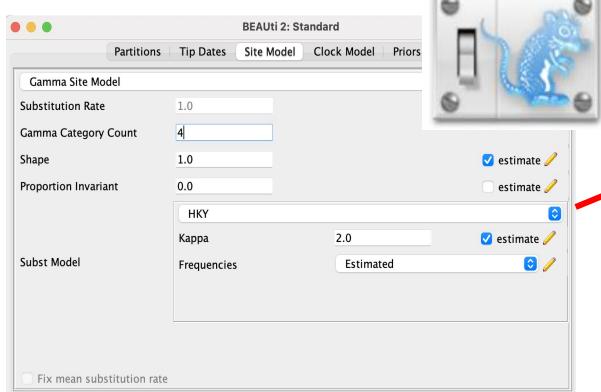
> 35,000 scientific studies have used this software since 2007;
> 2,000 studies last year alone (so far in 2023)

Bayesian phylogenetic inference software

- **MrBayes** (2003-): >79,000 citations
- **BEAST1 + 2** (2007-): >35,000 citations
- **PhyloBayes** (2007-): >1900 citations
- **BayesPhylogenies** (2006-): >1000 citations
- **BPP** (2015?-): >800 citations
- **RevBayes** (2016-): >500 citations
- **Phycas** (2005?-): >300 citations

Analysis Pipeline for BEAST 2

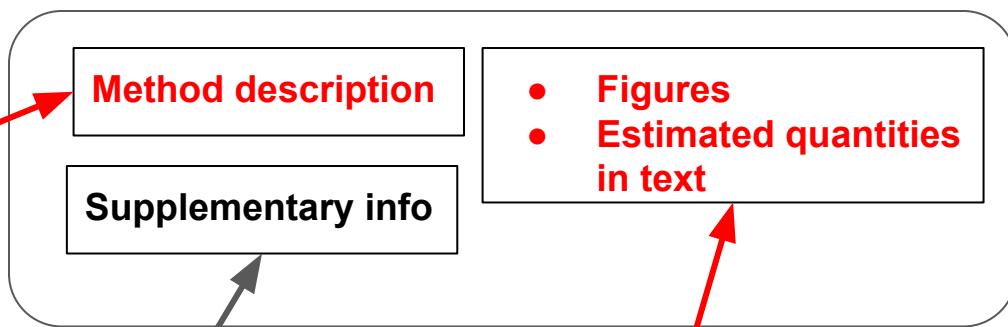
BEAUTi



BEAST2 XML

```
28 <state id="State" spec="State">
29   <parameter id="lambda" spec="parameter.RealParameter" lower="0.0" name="stateNode">9.591500353792183</parameter>
30   <stateNode id="psi" spec="beast.util.TreeParser" IsLabelledNewick="true" newick="((14:0.06962779760393568,(6:0.05
31     <taxonset id="TaxonSet" spec="TaxonSet"> ... </taxonset>
32   </stateNode>
33 </state>
34 <distribution id="posterior" spec="util.CompoundDistribution">
35   <distribution id="prior" spec="util.CompoundDistribution">
36     <distribution id="lambda_prior" spec="beast.math.distributions.Prior" x="@lambda">
37       <distr id="LogNormalDistributionModel" spec="beast.math.distributions.LogNormalDistributionModel">
38         <parameter id="RealParameter" spec="parameter.RealParameter" name="M">3.0</parameter>
39         <parameter id="RealParameter1" spec="parameter.RealParameter" name="S">1.0</parameter>
40     </distr>
41   </distribution>
42   <distribution id="YuleModel" spec="beast.evolution.speciation.YuleModel" birthDiffRate="@lambda" tree="@psi"/>
43 </distribution>
44 <distribution id="likelihood" spec="util.CompoundDistribution">
45   <distribution id="D.treeLikelihood" spec="ThreadedTreeLikelihood" tree="@psi">
46     <data idref="D"/>
47     <siteModel id="SiteModel" spec="SiteModel">
48       <subsiteModel id="JukesCantor" spec="JukesCantor"/>
49     </siteModel>
50     <branchRateModel id="StrictClockModel" spec="beast.evolution.branchratemodel.StrictClockModel">
51       <parameter id="RealParameter2" spec="parameter.RealParameter" name="clock_rate">1.0</parameter>
52     </branchRateModel>
53   </distribution>
54 </distribution>
55 </distribution>
```

Paper



Post-processing tools

- Tracer
- TreeAnnotator
- FigTree
- DensiTree
- IcyTree
- etc...

BEAST 2.7



Output files

- Log file
- Trees file

Analysis Pipeline for RevBayes

Rev Script

```
...  
### diversification = birth_rate - death_rate  
### assume an exponential prior distribution  
### and apply a scale proposal  
diversification ~ dnExponential(10.0)  
moves[mi++] = mvScale(diversification,lambda=1.0,tune=true,weight=3.0)  
  
### turnover = death_rate / birth_rate  
### this parameter can only take values between 0 and 1  
### use a Beta prior distribution  
### and a slide move  
turnover ~ dnBeta(2.0, 2.0)  
moves[mi++] = mvSlide(turnover,delta=1.0,tune=true,weight=3.0)  
  
### the parameters of the BDP include birth and death rates  
### these are deterministic variables of the diversification & turnover  
##### create a variable to ensure the rates are always positive (RealPos)  
denom := abs(1.0 - turnover)  
##### birth_rate = diversification / (1 - turnover)  
birth_rate := diversification / denom  
##### death_rate = (turnover * diversification) / (1 - turnover)  
death_rate := (turnover * diversification) / denom  
  
### rho is the probability of sampling species at the present  
### fix this to 0.068, since there are ~147 described species of caniforms (bears,  
dogs, mustelids, pinnipeds, etc.)  
### and we have sampled 10  
rho <- 0.068  
  
### the root age is an independent stochastic node with a lognormal prior  
### the lognormal distribution is off-set by the age of the canid fossil in the genus  
Hesperocyon  
### the observation time of Hesperocyon is ~ 38 Mya  
tHesperocyon <- 38.0  
### the mean of the lognormal distribution is set to 15 Mya older than the observed  
fossil  
### when offset by tHesperocyon, this gives a mean root_time of 49 Mya  
mean_ra <- 11.0  
stdv_ra <- 0.25  
### the lognormal distribution is parameterized by mu which is a function of the  
mean and standard deviation  
mu_ra <- ln(mean_ra) - ((stdv_ra*stdv_ra) * 0.5)  
root_time ~ dnLnorm(mu_ra, stdv_ra, offset=tHesperocyon)  
...
```

Paper

Method description

Supplementary info

- **Figures**
- **Estimated quantities in text**

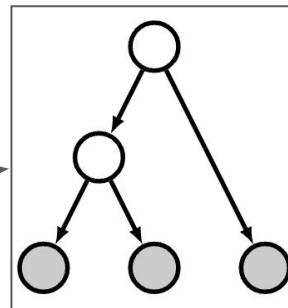
Post-processing tools

- **summarize()**
- **mapTree()**
- **Tracer**
- **etc**

Output files

- **Log file**
- **Trees file**
- **Map tree**
- **etc**

RevBayes



A phylogenetic model

To date the divergence time for *Y. pestis* and nodes within the *Y. pestis* clade we performed Bayesian Markov Chain Monte Carlo simulations using BEAST-2.3.0 (Bouckaert et al., 2014) and the BEAGLE library v2.1.2 (Ayres et al., 2012). We used the codon-partitioned supermatrix that included the two closest *Y. pseudotuberculosis* clades, with unlinked substitution models, GTR+G+I with eight gamma rate categories and unlinked clock models. Dates were set as years ago with the RISE509, RISE505, Justinian and Black Death samples set to 4,761, 3,701, 1,474, and 667 years ago, respectively. All unknown dates were set to 0 years ago. We followed previous work (Cui et al., 2013, Wagner et al., 2014) and applied a lognormal relaxed clock, assuming a constant population size.

Another phylogenetic model

RAxML was used to construct maximum likelihood phylogenies for each segment (33). The general-time reversible model of nucleotide substitution was assumed, and the rates were assumed to vary according to the discrete Gamma distribution with four rate categories. ... Bayesian phylogenetic trees were constructed for each segment using Bayesian Evolutionary Analysis Sampling Trees Software (BEAST) (35). The model of substitution and rate variation in each tree was the same as in the corresponding RAxML tree, with the additional assumption of a log-normally distributed molecular clock. BEAST trees were initiated with the RAxML majority rule consensus tree to reduce the computation time required for the MARcov Chain Monte Carlo (MCMC) sampler to reach the stationary distribution. Maximum clade credibility (MCC) trees were constructed using 2,000 equidistant samples from the final 30% of the posterior sample of trees ($n = 25 \times 106$ total samples) using a python script and tree annotator (35). The resulting MCC trees were partitioned into clades using a median branch length distance threshold of 0.20 (24). Briefly, clades with median root-to-tip distances equal to or less than 20% of the median root-to-tip distances of all clades in the tree were selected as clusters. Clades were then manually merged, if necessary, based on the branching posterior and reported classification (19, 20, 25). Each clade was assigned a unique clade ID.

A new time tree reveals Earth history's imprint on the evolution of modern birds

Santiago Claramunt* and Joel Cracraft*

2015 © The Authors, some rights reserved;
exclusive licensee American Association for
the Advancement of Science. Distributed
under a Creative Commons Attribution
NonCommercial License 4.0 (CC BY-NC).
10.1126/sciadv.1501005

The genomic data set was analyzed under a single GTR+ gamma substitution model. Bayesian inference was performed using BEAST 2 (44). Rate heterogeneity across lineages was modeled with a relaxed lognormal clock (39). Priors on rates of substitution and clocks were set to defaults (as in Beauti 2.1). We ran analyses using a birth-death tree prior with incomplete sampling that takes into account the fact that we sampled a small fraction of all avian diversity (Supplementary Materials).

Supplementary Materials for

A new time tree reveals Earth history's imprint on the evolution of modern birds

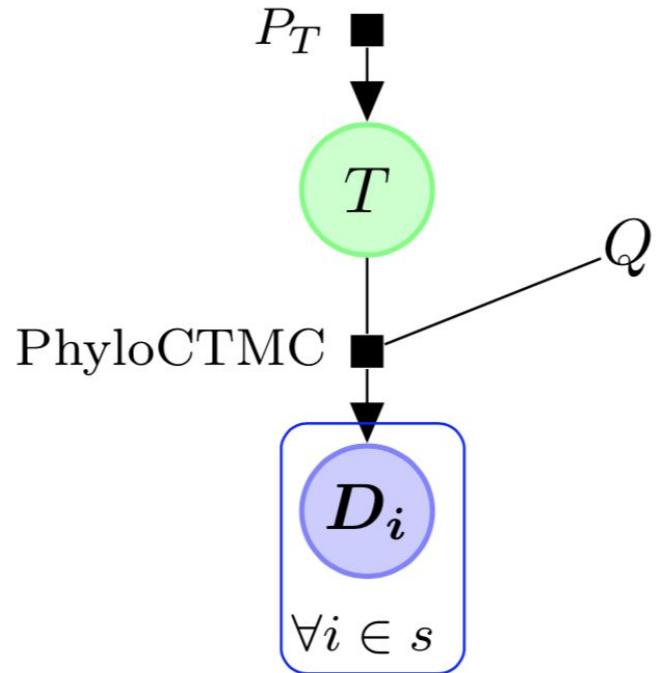
Santiago Claramunt and Joel Cracraft

Published 11 December 2015, *Sci. Adv.* 1, e1501005 (2015)

DOI: 10.1126/sciadv.1501005

We ran the main analyses using a birth-death tree prior with incomplete sampling (158) that takes into account that we sampled a small fraction of all avian species diversity; using the ‘birth-death skyline contemporary’ prior, part of the BEAST BDSKY 1.1.0 add-on (159), we set a uniform prior (0, infinity) for the birth/total death ratio (effective reproduction number R₀), a uniform prior (0, infinity) for the “total death” rate (“become uninfected”) prior, and a prior for the sampling probability (ρ) to take into account uncertainty regarding the total number of extant avian species: we set the parameters of a beta distribution so that it resulted in a bellshaped distribution with the 5% quantile at 0.023, the proportion of the currently assumed 10,000 avian species in the RAG dataset (0.0048 for the genome dataset), and the 95% quantile at 0.008 (0.0016 for the genome dataset), corresponding to assuming about 30,000 extant avian species.

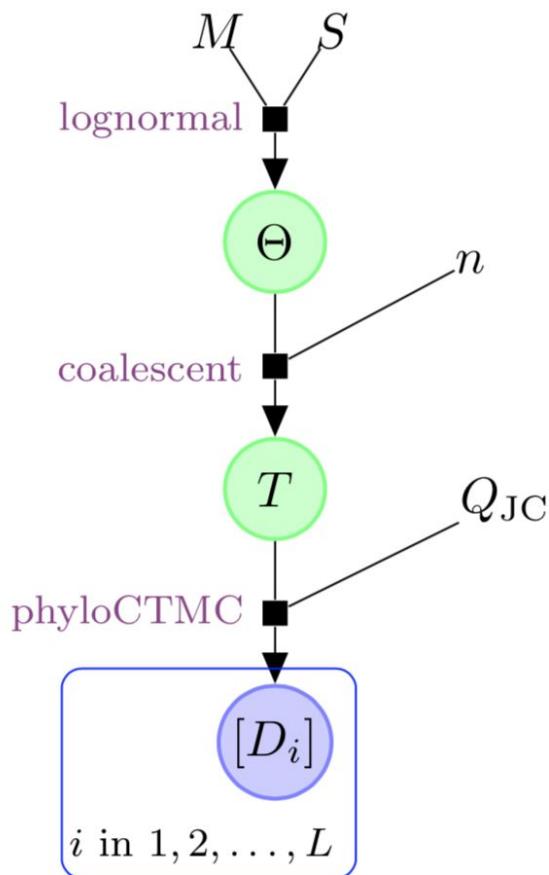
A phylogenetic model



$$P(T|\mathbf{D}) \propto \prod_i \Pr(\mathbf{D}_i|T)P(T)$$

A simple phylogenetic model

(a) Graphical model



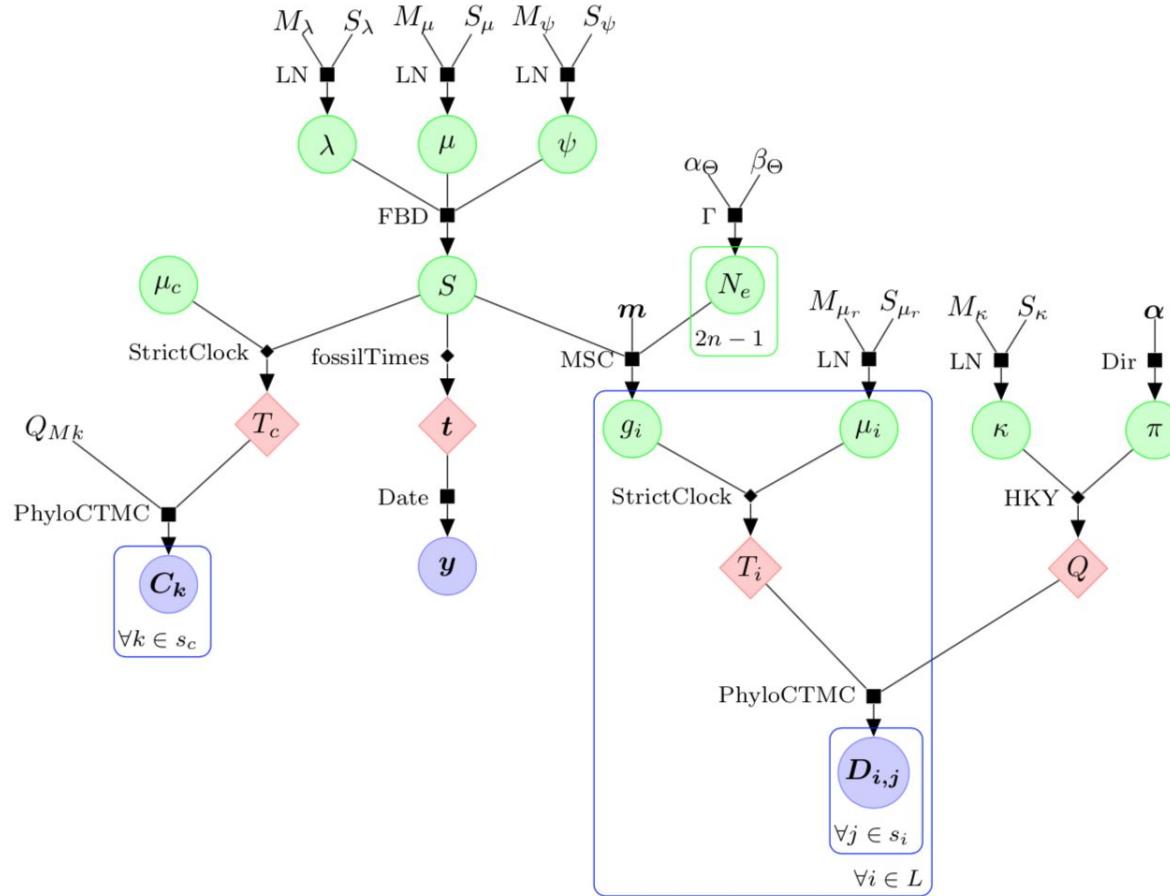
(b) Canonical phylogenetic posterior

$$P(T, \Theta | \mathbf{D}) = \frac{1}{P(\mathbf{D})} \prod_i [\Pr(D_i | T)] P(T | \Theta) P(\Theta)$$

(c) Probabilistic modelling language (EPML)

```
theta ~ lognormal(M=-5, S=1)
T ~ coalescent(theta, n=10)
D ~ phyloCTMC(T, Q=jukesCantor(), L=1000)
```

A phylogenetic model



$$\begin{aligned}
 P(S, \mathbf{N}_e, \mathbf{g}, \boldsymbol{\mu}, \lambda, Q | \mathbf{D}, \mathbf{C}, \mathbf{y}) \propto & \prod_i \left[P(g_i | S, \mathbf{N}_e) \prod_j \Pr(\mathbf{D}_{i,j} | g_i \mu_i, Q) \right] \times \\
 & \prod_k \Pr(C_k | \mu_c S, Q_{Mk}) \times \\
 & P(\mathbf{y} | S) P(S | \lambda) P(\lambda) P(\mathbf{N}_e) P(\boldsymbol{\mu}) P(\mu_c) P(Q)
 \end{aligned}$$

Reproducible and reusable research

Some good/necessary trends in science:

- The data science notebook paradigm (e.g. Jupyter / R notebooks)
- Everything open source, open access, open data and online
- Live science
- Not just reproducible but reusable
 - What happens if we change a key assumption in the model?
 - What happens if we add a little more data, or use different data?

The reality is that most research isn't reproducible

- Only a degree of reproducibility and reusability
 - If you have the XML file then you might be able to run it again
 - Need to find the version of the software that was used
 - No way of telling if the XML that was used to produce figures in the paper is the same as the XML that is provided in supplemental information
 - XML is very difficult for a normal person to read
 - Might be better for other formats (NEXUS?) but still rather complex and non-normalized
 - Users understanding of the model they ran is imperfect, leading to errors and omissions in model descriptions in methods section of paper
 - Can be hard to reuse XML on a new data set.

LinguaPhylo language

Language features

The LPhy language provides a simple syntax for specifying and simulating under different evolutionary models.

```
data {
    L = 200;
    taxa = taxa(names=1:10);
}
model {
    Θ ~ LogNormal(meanlog=3.0, sdlog=1.0);
    ψ ~ Coalescent(theta=Θ, taxa=taxa);
    D ~ PhyloCTMC(L=L, Q=jukesCantor(), tree=ψ);
}
```

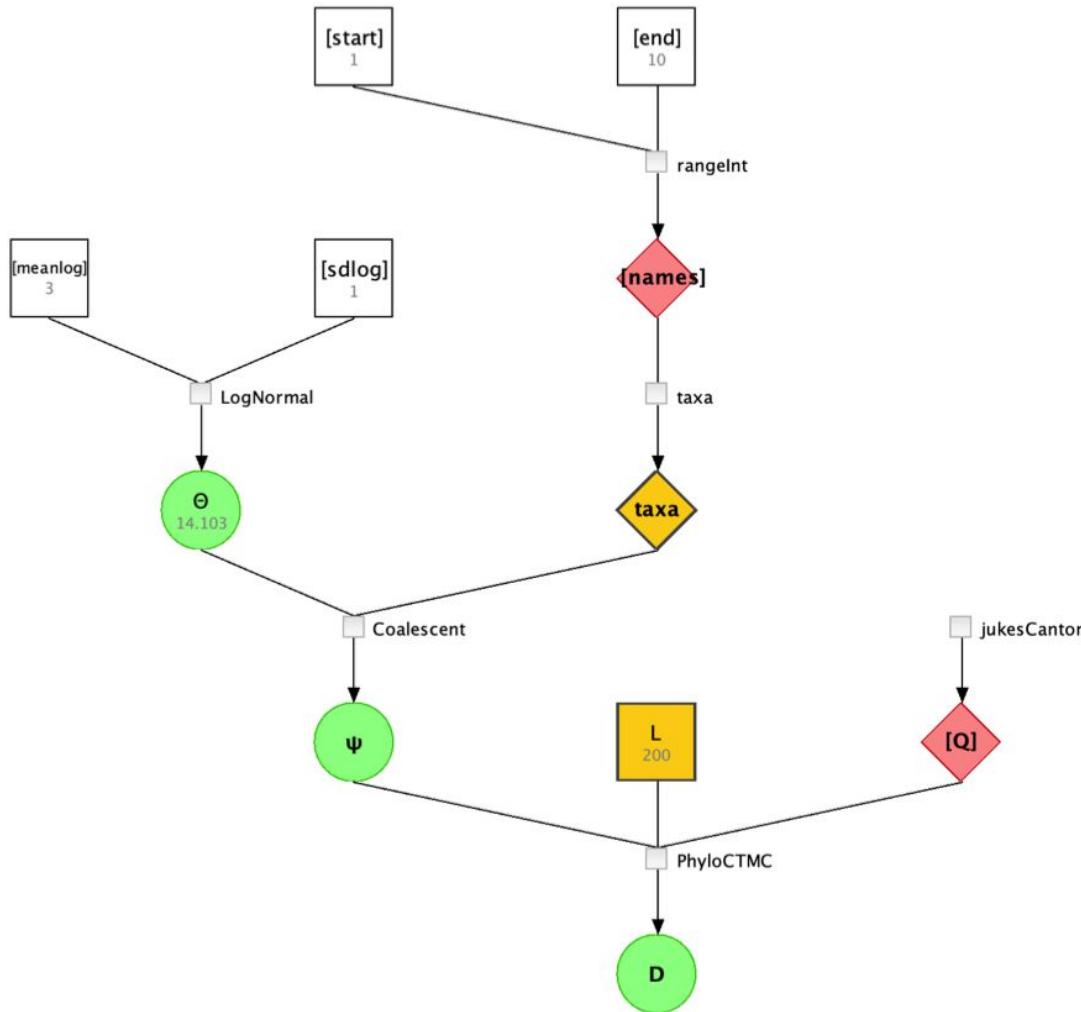
Unicode characters, such as greek letters, can be used for variable names to aid in consistency of model specification across code and mathematical descriptions.

Language features

```
data {
    L = 200;
    taxa = taxa(names=1:10);
}
model {
    Θ ~ LogNormal(meanlog=3.0, sdlog=1.0);
    ψ ~ Coalescent(theta=Θ, taxa=taxa);
    D ~ PhyloCTMC(L=L, Q=jukesCantor(), tree=ψ);
}
```

This listing above specifies a complete phylogenetic model using only five lines of code inside two blocks. Constant nodes with fixed values are shown in magenta. The first block specifies “200” nucleotide sites in “L” and ten different taxa named ‘1’ to ‘10’ in “taxa”. The second block declares the stochastic nodes or random variables in green, and their generative distributions in blue.

Graphical model representation



Factor nodes: nodes representing a generative distribution (e.g. Coalescent) or a deterministic function (e.g. rangelInt) with zero or more inputs and one output.

Value nodes: The result of a factor (random variables generated by a distribution, or functions of previous values), and constant nodes.



LPhy Studio version 1.3.1 - simpleYule.lphy

AutoLayout Interactive

Current Constants Variables Model Narrative Latex Variable Summary ▶

reps: 1 Sample Layering: From Sinks Show constants Edit values

D: 1 4 13 6 15 5 11 12 7 14 0 8 10 3 2 9
 ψ : 1 4 13 6 15 5 12 11 7 14 0 8 10 3 2 9
 λ : 5.171416194048878

data model

```

1 λ ~ LogNormal(meanlog=3.0, sdlog=1.0);
2 ψ ~ Yule(lambda=λ, n=16);
3 D ~ PhyloCTMC(L=200, Q=jukesCantor(), tree=ψ);
4
  
```

Read LPhy script /Users/dxie004/WorkSpace/linguaPhylo/lphy-studio/build/distributions/lphy-studio-1.3.1/examples/simpleYule.lphy from /Users/dxie004/WorkSpace/linguaPhylo/lphy-studio/build/distributions/lphy-studio-1.3.1/examples

Integrative total-evidence models

LPhy Studio version 1.3.1 - totalEvidence0.lphy

AutoLayout Interactive

Current Constants Variables Model Narrative Latex Variable Summary ▶

Dirichlet → π
LogNormal → κ
LogNormal → θ
 π , κ , θ → hky
Coalescent → ψ
LogNormal → σ^2
Normal → y_0
hky, ψ , σ^2 , y_0 → D
 ψ , σ^2 , y_0 → PhyloCTMC
 σ^2 , y_0 → PhyloBrownian
D → y

rep: 1 Sample Layering: From Sinks Show constants Edit values

D: 14
[Q]:
[T]: [0.257596, 0.339928, 0.264279, 0.138196]
[K]: 2.177835863795988
Ψ:
Φ: 8.52287283006234
y: Taxa trait_0
0 2.4363331908701467
1 -0.42259350740189583
10 -0.27892720287267325
11 4.731203763816379
12 -3.0401049675470886
13 -7.1757519356512525
14 -4.422286994108373
15 -1.5361726146487946
2 16.650631037049852
3 -7.555915480926887
4 -8.664948373041426

data model

```
1  $\kappa \sim \text{LogNormal}(\text{meanlog}=1.0, \text{sdlog}=0.5);$ 
2  $\pi \sim \text{Dirichlet}(\text{conc}=[2.0, 2.0, 2.0, 2.0]);$ 
3  $\theta \sim \text{LogNormal}(\text{meanlog}=3.0, \text{sdlog}=1.0);$ 
4  $\psi \sim \text{Coalescent}(n=16, \text{theta}=0);$ 
5  $\sigma^2 \sim \text{LogNormal}(\text{meanlog}=1.0, \text{sdlog}=0.5);$ 
6  $y_0 \sim \text{Normal}(\text{mean}=0.0, \text{sd}=1.0);$ 
7  $D \sim \text{PhyloCTMC}(L=200, Q=hky(\kappaappa=\pi, \text{freq}=\psi), \text{tree}=\psi);$ 
8  $y \sim \text{PhyloBrownian}(\text{diffRate}=\sigma^2, y_0=y_0, \text{tree}=\psi);$ 
9
```

Read LPhy script /Users/dxie004/WorkSpace/linguaPhylo/lphy-studio/build/distributions/lphy-studio-1.3.1/examples/totalEvidence0.lphy from /Users/dxie004/WorkSpace/linguaPhylo/lphy-studio/build/distributions/lphy-studio-1.3.1/examples

Integrative total-evidence models

LPhy Studio version 1.3.1 - totalEvidence0.lphy

AutoLayout Interactive Current Constants Variables Model Narrative Latex Variable Summary ▶

```

graph TD
    Dirichlet[Dirichlet] --> pi((π))
    LogNormalK[LogNormal] --> kappa((κ))
    LogNormalTheta[LogNormal] --> theta((Θ))
    LogNormalSigma2[LogNormal] --> sigma2((σ²))
    NormalY0[Normal] --> y0((y₀))
    hky[hky] --> Q[Q]
    Coalescent[Coalescent] --> psi((ψ))
    PhyloCTMC[PhyloCTMC] --> D((D))
    PhyloBrownian[PhyloBrownian] --> y((y))
    
```

totalEvidence0.lphy

Code

```

model {
    n ~ Dirichlet(conc=[2.0, 2.0, 2.0, 2.0]);
    κ ~ LogNormal(meanlog=1.0, sdlog=0.5);
    Θ ~ LogNormal(meanlog=3.0, sdlog=1.0);
    ψ ~ Coalescent(n=16, theta=Θ);
    D ~ PhyloCTMC(l=200, Q=hky(kappa=κ, freq=n), tree=ψ);
    σ² ~ LogNormal(meanlog=1.0, sdlog=0.5);
    y₀ ~ Normal(mean=0.0, sd=1.0);
    y ~ PhyloBrownian(diffRate=σ², tree=ψ, y₀=y₀);
}
    
```

Model

The alignment, D is assumed to have evolved under a phylogenetic continuous time Markov process ([Felsenstein: 1981](#)) on phylogenetic time tree, ψ , with an instantaneous rate matrix and a length of 200. The instantaneous rate matrix is the HKY model ([Hasegawa et al: 1985](#)) with transition bias parameter, κ and base frequency vector, π . The base frequency vector, π have a Dirichlet distribution prior with a concentration of [2.0, 2.0, 2.0, 2.0]. The transition bias parameter, κ has a log-normal prior with a mean in log space of 1.0 and a standard deviation in log space of 0.5. The time tree, ψ is assumed to come from a Kingman's coalescent tree prior ([Rodrigo and Felsenstein: 1999](#)) with coalescent parameter, Θ and an n of 16. The coalescent parameter, Θ has a log-normal prior with a mean in log space of 3.0 and a standard deviation in log space of 1.0. The continuous character data, y is assumed to have evolved under a phylogenetic Brownian motion process ([Felsenstein: 1973](#)) with tree, y , evolutionary rate, σ^2 and root value, y_0 . The evolutionary rate, σ^2 has a log-normal prior with a mean in log space of 1.0 and a standard deviation in log space of 0.5. The root value, y_0 has a normal prior with a mean of 0.0 and a standard deviation of 1.0.

Posterior

$$P(\pi, \kappa, \psi, \Theta, \sigma^2, y_0 | D, y) \propto P(D|\psi, Q)P(\pi)P(\kappa)P(\psi|\Theta)P(\Theta)P(y|\psi, \sigma^2, y_0)P(\sigma^2)P(y_0)$$

Include
Code
Data
Model
Posterior
Graphical Model
References

Exclude

reps: 1 Sample Layering: From Sinks Show constants Edit values

data model

```

1 κ ~ LogNormal(meanlog=1.0, sdlog=0.5);
2 π ~ Dirichlet(conc=[2.0, 2.0, 2.0, 2.0]);
3 Θ ~ LogNormal(meanlog=3.0, sdlog=1.0);
4 ψ ~ Coalescent(n=16, theta=Θ);
5 σ² ~ LogNormal(meanlog=1.0, sdlog=0.5);
6 y₀ ~ Normal(mean=0.0, sd=1.0);
7 D ~ PhyloCTMC(l=200, Q=hky(kappa=κ, freq=n), tree=ψ);
8 y ~ PhyloBrownian(diffRate=σ², tree=ψ, y₀=y₀);
    
```

Read LPhy script /Users/dxie004/WorkSpace/linguPhylo/lphy-studio/build/distributions/lphy-studio-1.3.1/examples/totalEvidence0.lphy from /Users/dxie004/WorkSpace/linguPhylo/lphy-studio/build/distributions/lphy-studio-1.3.1/examples

Vectorization / array programming

Named variables in LPhy can be scalars or vectors. Any generator can be vectorized to produce a vector of i.i.d. values by using the `replicates` keyword:

```
 $\kappa$  ~ LogNormal(meanlog=0.5, sdlog=1.0, replicates=3);
```

Vectorization can also be applied to a distribution that already produces vectors. In which case, the output will be a matrix as in the following example, where the major dimension has size 3, with each element of the π random variable being a vector of base frequencies.

```
 $\pi$  ~ Dirichlet(conc=[2.0, 2.0, 2.0, 2.0], replicates=3);
```

Finally, vectorization can be coerced simply by passing a vector of elements instead of a single element to one or more of the inputs of a generator:

```
Q = hky(kappa= $\kappa$ , freq= $\pi$ );
```

Parametric Distributions

LPhy core has the familiar parametric distributions available as generative distributions:

- Beta
- Cauchy
- Dirichlet
- Exp
- Gamma
- InverseGamma
- LogNormal
- MVN
- Normal
- Uniform
- Weibull
- Bernoulli
- Geometric
- Poisson

Version 1.4 has some core phylogenetic distributions

Tree models

- BirthDeathSerialSampling
- Coalescent
- MultispeciesCoalescent
- SkylineCoalescent
- StructuredCoalescent
- SimFossilsPoisson
- Yule
- ...

Data sampling distributions

- ErrorModel
- PhyloBrownian
- PhyloCTMC
- ...

Version 1.4 has some core functions

Rate matrices

- BinaryRateMatrix
- GTR
- TN93
- HKY
- F81
- K80
- JukesCantor

Phylogenetic functions

- Newick

Standard functions

- Exp
- Floor
- Log
- Rep
- ...

The LPhy 1.3.1 Doc is available at

<https://github.com/LinguaPhylo/linguaPhylo/blob/1.3.1/lphy/doc/index.md>

Model Guide version 0.1.0

Models Latex

Model category : Tree functions

Name	Category	Description
countMigrations	Tree functions	The number of single-child nodes in the tree where the 'deme' attribute changes.
extantTree	Tree functions	A tree pruned from a larger tree by retaining only the tips at time zero.
localBranchRates	Tree functions	A function that returns branch rates for the given tree, indicator mask and raw ...
newick	Tree functions	A function that parses a tree from a newick formatted string.
nodecount	Tree functions	The number of nodes in the tree
pruneTree	Tree functions	A tree pruned from a larger tree by retaining only nodes subtending nodes with n ...
treeLength	Tree functions	The sum of all the branch lengths in the tree.
TimeTree	Tree functions	Method calls

countMigrations(*tree*)

The number of single-child nodes in the tree where the 'deme' attribute changes.

Parameters:

- TimeTree **tree** - the tree.

Return type:

- Integer

Examples

simpleStructuredCoalescent.lphy

Coalescent model

The simplest coalescent model LPhy implements is the constant-population size coalescent, which can be extended to generate serially sampled (heterochronous) data

```
 $\psi$  ~ Coalescent(theta=Θ, taxa=taxa(names=["a", "b", "c", "d"],  
ages=[0.0, 1.0, 2.0, 3.0]));
```

The taxa argument is optional, and only needed if you want to give the taxa specific names, as in the case when doing inference against specific input data.

Structured coalescent model

The structured coalescent generalizes the constant-population size coalescent by allowing multiple demes, each of which are characterized by a distinct population size (in the simplest case this population size does not change through time).

```
M = migrationMatrix(theta=[0.1, 0.1], m=[1.0, 1.0]);
g ~ StructuredCoalescent(M=M, n=[15, 15]);
```

Skyline coalescent model

The following script specifies a Skyline coalescent model with 10 coalescent intervals (hence 11 taxa), governed by four distinct population sizes.

```
g ~ SkylineCoalescent(theta=[0.1, 0.2, 0.3, 0.4], groupSizes=[4,3,2,1]);
```

Here, “g” is a stochastic node in the PGM, with its value sampled from the 160 SkylineCoalescent generative distribution. Ten coalescent intervals are defined 161 through the “groupSizes” argument: the first four coalescent intervals will be drawn 162 assuming a “theta” of 0.1, the next three intervals with “theta” equal to 0.2, and so on.

Birth-death-serial-sampling model

Birth-death models are commonly used in macroevolution as sampling distributions for species trees.

```
ages = [0.0, 1.0, 2.0, 3.0, 4.0];
tree ~ BirthDeathSerialSampling(lambda=1, mu=0.5, rho=0.1,
psi=1, rootAge=5, ages=ages);
```

Other tree models including different conditions on birth-death models, the fossilized birth-death model and the Yule process are also implemented.

Substitution models

Substitution models consist of continuous-time Markov chains (CTMC) used to model the evolution of discrete characters, such as nucleotides and amino acid residues.

LPhy implements a general formulation of a phylogenetic CTMC, known as the GTR model, under which several nested models can be specified. The first line below constructs the instantaneous rate matrix (Q) for an HKY model, which is then used to in PhyloCTMC, the generative distribution over sequence alignment data (D):

```
Q = hky(kappa=2.0, freq=[0.2, 0.25, 0.3, 0.25]);  
Θ ~ LogNormal(meanlog=3.0, sdlog=1.0);  
ψ ~ Coalescent(theta=Θ, taxa=taxa);  
D ~ PhyloCTMC(L=200, Q=Q, tree=ψ);
```

Evolutionary clock models

The simplest clock model is the strict clock, under which the evolutionary rate remains constant over the entire tree. Specifying a strict molecular clock can be done by specifying the “mu” parameter in the PhyloCTMC distribution (default is 1.0)

```
 $\lambda$  ~ LogNormal(meanlog=3.0, sdlog=1.0);  
 $\psi$  ~ Yule(lambda= $\lambda$ , n=16);  
D ~ PhyloCTMC(L=200, Q=jukesCantor(), tree= $\psi$ , mu=0.5);
```

Relaxed phylogenetic clock models

More realistic clock models like the uncorrelated relaxed clock model assume the rate for each branch is drawn according to a parametric distribution. For example, a relaxed clock with rates drawn from a log-normal distribution can be constructed as follows:

```
 $\lambda$  ~ LogNormal(meanlog=3.0, sdlog=1.0);
 $\psi$  ~ Yule(lambda= $\lambda$ , n=16);
branchRates ~ LogNormal(sdlog=0.5, meanlog=-0.25, replicates= $\psi$ .branchCount());
D ~ PhyloCTMC(L=200, Q=jukesCantor(), branchRates=branchRates, tree= $\psi$ );
```

Documentation

The art of constructing generative distribution specifications is to do so in the most canonical and intuitive way possible. However there may still be ambiguity or uncertainty, so good documentation is essential (<https://linguaphylo.github.io>).

SkylineCoalescent distribution

SkylineCoalescent(Double[] theta, Integer[] groupSizes, Integer n, Taxa taxa, Double[] ages)

The skyline coalescent distribution over tip-labelled time trees. If no group sizes are specified, then there is one population parameter per coalescent event (as per classic skyline coalescent of Pybus, Rambaut and Harvey 2000)

Parameters

- Double[] **theta** - effective population size, one value for each group of coalescent intervals, ordered from present to past. Possibly scaled to mutations or calendar units. If no groupSizes are specified, then the number of coalescent intervals will be equal to the number of population size parameters.
- Integer[] **groupSizes** - A tuple of group sizes. The sum of this tuple determines the number of coalescent events in the tree and thus the number of taxa. By default all group sizes are 1 which is equivalent to the classic skyline coalescent.
- Integer **n** - number of taxa.
- Taxa **taxa** - Taxa object, (e.g. Taxa or Object[])
- Double[] **ages** - an array of leaf node ages.

Return type

- TimeTree

Reference

Drummond, A. J., Rambaut, A., Shapiro, B., & Pybus, O. G. (2005). Bayesian coalescent inference of past population dynamics from molecular sequences. *Molecular biology and evolution*, 22(5), 1185-1192. <http://doi.org/10.1093/molbev/msi103>

Inference and Data clamping

```
data {
    options = {ageDirection="forward", ageRegex="s(\d+)"};
    nexusFilePath = "tutorials/data/RSV2.nex";
    D = readNexus(file=nexusFilePath, options=options);
    codon = D.charset(["3-629\3", "1-629\3", "2-629\3"]);
    n = 3;
    L = [209, 210, 210];
    taxa = D.taxa();
}
model {
     $\pi$  ~ Dirichlet(replicates=n, conc=[2.0, 2.0, 2.0, 2.0]);
     $\kappa$  ~ LogNormal(sdlog=0.5, meanlog=1.0, replicates=n);
     $r$  ~ WeightedDirichlet(conc=rep(element=1.0, times=n), weights=L);
     $\mu$  ~ LogNormal(meanlog=-5.0, sdlog=1.25);
     $\Theta$  ~ LogNormal(meanlog=3.0, sdlog=2.0);
     $\psi$  ~ Coalescent(taxa=taxa, theta= $\Theta$ );
    Q = hky(kappa= $\kappa$ , freq= $\pi$ , meanRate= $r$ );
    codon ~ PhyloCTMC(L=L, Q=Q, mu= $\mu$ , tree= $\psi$ );
}
```

LPhyStudio

LPhy Studio version 1.3.0 - RSV2.lphy

AutoLayout Interactive

RSV2.lphy

Model

The vector of integers, L are the number of characters in vector of alignments, codon . The vector of alignments, codon are the character sets ["3-6293", "1-6293", "2-6293"] of metadatalignment, D . The metadatalignment, D is read from the Nexus file with a file name of "data/RSV2.nex" and options. The options are ageDirection="forward" and ageRegex="s (d+)". The integer, n is the length of vector of alignments, codon . The taxa is the list of taxa of metadatalignment, D .

Posterior

$$P(\pi, \kappa, r, \mu, \psi | \text{codon}) \propto \prod_{i=0}^{n-1} P(\text{codon}_i | \psi, \mu, Q_i) \prod_{i=0}^{n-1} P(\pi_i) \prod_{i=0}^{n-1} P(\kappa_i) P(r) P(\mu) P(\psi | \Theta) P(\Theta)$$

References

- Felsenstein, J. (1981). Evolutionary trees from DNA sequences: a maximum likelihood approach. Journal of molecular evolution, 17(6), 368-376. <https://doi.org/10.1007/BF01734359>
- Hasegawa, M., Kishino, H., & Yano, T. (1985) Dating of the human-ape splitting by a molecular clock of mitochondrial DNA. J Mol Evol 22, 160-174 <https://doi.org/10.1007/BF02101694>

Include	Exclude
Data	Code
Model	Graphical Model
Posterior	
References	

reps: 1 Sample Layering: From Sinks Show constants Edit values

data model

```

1 | K ~ LogNormal(meanlog=1.0, sdlog=0.5, replicates=n);
2 | π ~ Dirichlet(conc=[2.0, 2.0, 2.0, 2.0], replicates=n);
3 | r ~ WeightedDirichlet(conc=rep(element=1.0, times=n), weights=L);
4 | μ ~ LogNormal(meanlog=-5.0, sdlog=1.25);
5 | ψ ~ LogNormal(meanlog=3.0, sdlog=2.0);
6 | Θ ~ Coalescent(taxa=taxa, theta=8);
7 | codon ~ PhyloCTMC(L=L, Q=hky(kappa=K, freq=π, meanRate=r), mu=μ, tree=ψ);
8 |
  
```

Load nucleotide alignment 129 by 629, age direction is forward

Simulation example

```
model {
    π ~ Dirichlet(conc=[3.0, 3.0, 3.0, 3.0,
                         3.0, 3.0, 3.0, 3.0,
                         3.0, 3.0, 3.0, 3.0,
                         3.0, 3.0, 3.0, 3.0]);
    rates ~ Dirichlet(conc=[1.0, 2.0, 1.0, 1.0, 2.0, 1.0]);
    Q = gt16(rates=rates, freq=π);
    Θ ~ LogNormal(meanlog=-2.0, sdlog=1.0);
    ψ ~ Coalescent(n=16, theta=Θ);
    A ~ PhyloCTMC(L=200, Q=Q, dataType=phasedGenotype(), tree=ψ);
    delta ~ Beta(alpha=1.5, beta=4.5);
    epsilon ~ Beta(alpha=2, beta=18);
    E ~ GT16ErrorModel(alignment=A, delta=delta, epsilon=epsilon);
}
```

Model in LPhyStudio

LPhy Studio version 1.2.0 - gt16CoalErrModel.lphy

AutoLayout Interactive Current Constants Variables Model Narrative Latex Variable Summary ▶

Diagram description: The model consists of several components. At the top, four observed variables ([conc], [conc], [sdlog], [meanlog]) provide priors for parameters pi (~Dirichlet), rates (~Dirichlet), theta (~LogNormal), and n (~Coalescent). The parameter gt16 is influenced by pi, rates, and theta. The parameter A is influenced by rates, theta, and a Coalescent process. The parameter Psi is influenced by theta and a PhylloCTMC process. The parameter epsilon is influenced by two Beta priors. The parameter delta is influenced by two Beta priors. The error node E is influenced by the GT16ErrorModel, which takes epsilon, delta, and alignment A as inputs.

Variables tab content:

- E: A heatmap showing posterior probabilities for each site across the tree.
- A: A dendrogram showing the phylogenetic tree.
- Q: A matrix of values corresponding to the Q matrix.
- PhylloCTMC: A matrix of values corresponding to the transition matrix.
- GT16ErrorModel: A matrix of values corresponding to the error model parameters.
- pi: [0.049543, 0.139006, 0.039881, 0.056556, 0.037495, 0.017662, 0.118413, 0.046075, 0.075668, 0.069751, 0.051081, 0.017082, 0.029874, 0.08551, 0.088075, 0.07833]
- rates: [0.0631, 0.623498, 0.125172, 0.02717, 0.117999, 0.043062]
- Psi: [0.158256, 0.147676, 0.41864, 0.11918, 0.03983, 0, 0, 0, 0.79428, 0.05263, -1.11437, 0.01824, 0.11235, 0, 0.01876, 0, 0, 0, 0.7321, 0.52005, 0.06359, -1.47224, 0.041, 0, 0, 0.12579, 0, 0, 0, 0.1044, 0.27615, 0.02891, -0.80279, 0, 0, 0, 0.04895, 0, 0.05263, 0, 0, -1.50543, 0.01876, 1.24298, 0.0971, 0.03461, 0, 0.14767, 0, 0, 0.03983, -0.53498, 0.05417, 0.09153, 0, 0.0319, 0, 0, 0.04237, 0, 0.39358, 0.00808, -0.67577, 0.0334, 0, 0, 0, 0.06008, 0.07901, 0.03509, 0.08585, -0.42345, 0, 0.52005, 0, 0, 0, 0.01715, 0, 0, 0, 0, -1.20516, 0.074, 0, 1.45915, 0, 0, 0, 0.00808, 0, 0, 0, 0.08038, -1.666, 0, 0, 0.41864, 0, 0, 0, 0.05417, 0, 0.79428, 0.0319, 0, 0, 0, 0.59367, 0, 0, 0, 0, 0.02108, 0.15946, 0.1385, 0.1044, 0, 0, 0, 0.07449, 0, 0, 0, 0, 0, 0.05486, 0, 0.29293, 0, 0, 0, 0.03509, 0, 0, 0, 0.23524, 0, 0, 0, 0, 0.08404, 0, 0, 0, 0, 0, 0.0505, 0, 0, 0, 0.11918, 0, 0, 0, 0, 0, 0.09153, 0]

Code tab content:

```

1 | θ ~ LogNormal(meanlog=-2.0, sdlog=1.0);
2 | ψ ~ Coalescent(n=16, theta=θ);
3 | π ~ Dirichlet(conc=[3.0, 3.0, 3.0, 3.0, 3.0, 3.0, 3.0, 3.0, 3.0, 3.0, 3.0, 3.0, 3.0, 3.0, 3.0, 3.0]);
4 | rates ~ Dirichlet(conc=[1.0, 2.0, 1.0, 1.0, 2.0, 1.0]);
5 | Q = gt16(freq=n, rates=rates);
6 | A ~ PhylloCTMC(L=200, Q=Q, tree=tree, dataType=phasedGenotype());
7 | epsilon ~ Beta(alpha=2, beta=18);
8 | delta ~ Beta(alpha=1.5, beta=4.5);
9 | E ~ GT16ErrorModel(delta=delta, epsilon=epsilon, alignment=A);
10 |
  
```

Message bar: sample(1) took 103 ms.

Method section generation

The alignment, \mathbf{E} is assumed to come from a GT16ErrorModel. The alignment, \mathbf{A} is assumed to have evolved under a phylogenetic continuous time Markov process [37] on phylogenetic time tree, ψ , with instantaneous rate matrix, \mathbf{Q} , a length of 200 and a dataType. The instantaneous rate matrix, \mathbf{Q} is the general time-reversible rate matrix on phased genotypes [35] with relative rates, **rates** and base frequencies, π . The base frequencies, π have a Dirichlet distribution prior with a concentration of [3.0, 3.0, 3.0, 3.0, 3.0, 3.0, 3.0, 3.0, 3.0, 3.0, 3.0, 3.0, 3.0, 3.0, 3.0, 3.0]. The relative rates, **rates** have a Dirichlet distribution prior with a concentration of [1.0, 2.0, 1.0, 1.0, 2.0, 1.0]. The dataType is the phased genotype data type. The phylogenetic time tree, ψ is assumed to come from a Kingman's coalescent tree prior [21] with coalescent parameter, Θ and an n of 16. The coalescent parameter, Θ has a log-normal prior with a mean in log space of -2.0 and a standard deviation in log space of 1.0. The delta has a Beta distribution prior with an alpha of 1.5 and a beta of 4.5. The epsilon has a Beta distribution prior with an alpha of 2 and a beta of 18.

LinguaPhylo is fully extensible by a Java API

```
public interface LPhyExtension {  
  
    /**  
     * @return the list of new {@link GenerativeDistribution} implemented  
     *         in the LPhy extension.  
     */  
    List<Class<? extends GenerativeDistribution>> getDistributions();  
  
    /**  
     * @return the list of new {@link Func} implemented in the LPhy extension.  
     */  
    List<Class<? extends Func>> getFunctions();  
  
    /**  
     * @return the map of new {@link SequenceType} implemented in the LPhy extension.  
     *         The string key is a keyword to represent this SequenceType.  
     *         The keyword can be used to identify and initialise the  
     *         corresponding sequence type.  
     */  
    Map<String, ? extends SequenceType> getSequenceTypes();  
}
```

GenerativeDistribution interface

```
12  public interface GenerativeDistribution<T> extends Generator<T> {  
13  
14      /**  
15      * The {@link RandomVariable} must be re-wrapped  
16      * to ensure correct behaviour downstream.  
17      * @return {@link RandomVariable} to connect to this {@link GenerativeDistribution}  
18      */  
19      RandomVariable<T> sample();
```

Func abstract class

```
/**  
 * @return a value generated by this generator.  
 */  
Value<T> generate();
```

LPhyBEAST

- Takes an LPhy script as input
- Generates BEAST 2.6 XML as output
- Can produce many model combinations that BEAUti can't
- Not all valid LPhy scripts are supported by BEAST 2.6

Future directions

- Supporting BEAST module developers to provide LPhy/LPhyBEAST support
- LPhyRevBayes
- LPhyMrBayes
- Integration of LPhy into Data Science notebook frameworks like Jupyter Notebooks and/or Stencila
- A formal language for defining the analytical outputs and figures of a particular analysis using a Bayesian model?