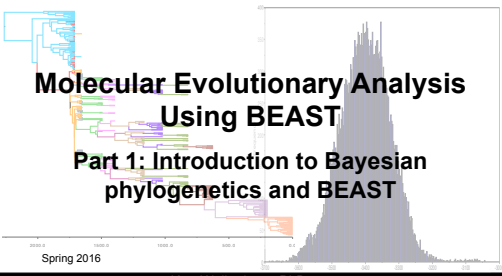


National Institute of Allergy and Infectious Diseases



**Molecular Evolutionary Analysis
Using BEAST**

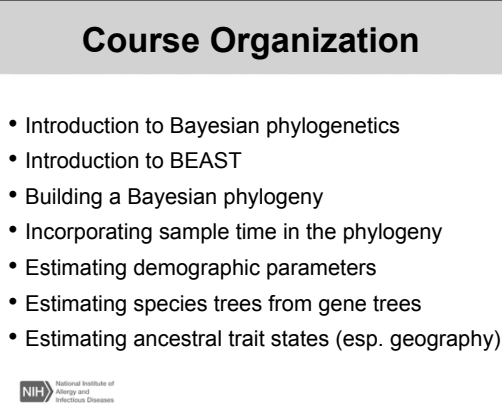
**Part 1: Introduction to Bayesian
phylogenetics and BEAST**

Spring 2016

Kurt Wollenberg, PhD
Phylogenetics Specialist
Bioinformatics and Computational Biosciences Branch
Office of Cyber Infrastructure and Computational Biology

Course Organization

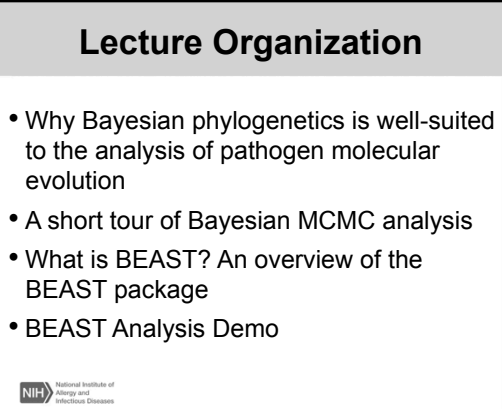
- Introduction to Bayesian phylogenetics
- Introduction to BEAST
- Building a Bayesian phylogeny
- Incorporating sample time in the phylogeny
- Estimating demographic parameters
- Estimating species trees from gene trees
- Estimating ancestral trait states (esp. geography)



NIAID

Lecture Organization

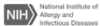
- Why Bayesian phylogenetics is well-suited to the analysis of pathogen molecular evolution
- A short tour of Bayesian MCMC analysis
- What is BEAST? An overview of the BEAST package
- BEAST Analysis Demo



NIAID

What's so special about pathogens?


- Short generation time
- Rapid evolution
- Genotypes - easy, phenotypes - hard
- Large populations
- Structured populations
- Rigorous temporal sampling of genotypes

 National Institute of Allergy and Infectious Diseases

NIAD

Why use Bayesian methods on pathogens?


- Coalescent approach is more appropriate
- Can incorporate temporal data
- Can incorporate geographical data
- Can incorporate host data

 National Institute of Allergy and Infectious Diseases

NIAD

What is Bayesian analysis?

- Calculation of the probability of parameters (tree, substitution model) given the data (sequence alignment)
- $p(\theta|D) = (\text{Likelihood} \times \text{prior}) / \text{probability of the data}$
- $p(\theta|D) = p(D|\theta)p(\theta)/p(D)$

 National Institute of Allergy and Infectious Diseases

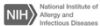

NIAD

Bayesian Analysis

Exploring the posterior probability distribution

Posterior probabilities of trees and parameters are approximated using Markov Chain Monte Carlo (MCMC) sampling

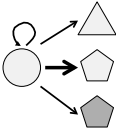
Markov Chain: A statement of the probability of moving from one state to another

What is MCMC?

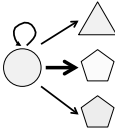
Markov Chain Monte Carlo

Markov chain





One link in the chain

Monte Carlo

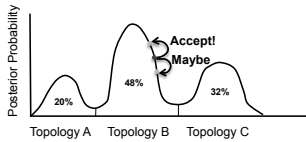




Choosing a link

What is MCMC?

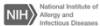
Markov Chain Monte Carlo: accept or reject? Metropolis-Hastings algorithm



$$\Pr(\text{accept}) = \min \left(1, \frac{\Pr(a) \Pr(X|a) \Pr(A|a)}{\Pr(A) \Pr(X|A) \Pr(a|A)} \right)$$



What is BEAST?

- **B**ayesian **E**volutionary **A**nalysis **S**ampling **T**rees
- A collection of programs for performing Bayesian MCMC analysis of molecular sequences
- Can incorporate sample time information
- Can perform a broad range of other evolutionary analyses using sequence data.


 National Institute of Allergy and Infectious Diseases

NIAD

What is BEAST?

The Programs:


- BEAUti - Creating XML input files
- BEAST - MCMC analysis of molecular sequences
- Tracer - Viewing MCMC output
- LogCombiner - Combining output files
- TreeAnnotator - Generate the consensus tree
- FigTree - Drawing a tree

 National Institute of Allergy and Infectious Diseases

NIAD

Different types of BEAST analyses

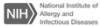

- Calculating a Bayesian coalescent phylogeny
- Calculating a Time-Stamped Bayesian coalescent
- Estimated population dynamics (Bayesian skyline/skyride/skygrid)
- Combined gene and species phylogeny estimate (*BEAST)
- Phylogeographic analysis (time and location data)

 National Institute of Allergy and Infectious Diseases

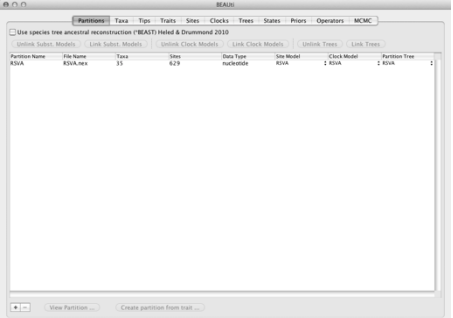


NIAD

Defining your analysis

- Prior knowledge of tree?
- Calibrating nodes?
- Substitution model?
- Effective population sizes?
- What priors to use?

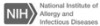




Setting up the analysis: BEAUTi

Setting up the analysis: BEAUTi

- Import data – Nexus or fasta format
- Incorporate known structure - taxa
- Substitution model parameters
- Strict or relaxed clock?
- Tree prior
- Substitution model priors
- Adjustments from previous runs (operators)
- Setting the chain






Setting up the analysis: BEAUTi

Import data: Nexus format

```
#NEXUS
[These are comments.
They are ignored by the program.]

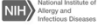

Begin data;
  dimensions ntax=5 nchar=15;
  format datatype=DNA gap=- missing=?;
  matrix
Bug1 ACCTGATTACGGGCA
Bug2 ACCCGAATACGGACA
Bug3 ACCTATTTACGCCCA
BugF ACTATATTACCGGCA
BugBX4W ACCAAA---CGGGCA
;
End;
```

Setting up the analysis: BEAUTi

Import data: Fasta format



```
>Bug1
ACCTGATTACGGGCA
>Bug2
ACCCGAATACGGACA
>Bug3
ACCTATTTACGCCCA
>BugF
ACTATATTACCGGCA
>BugBX4W
ACCAAA---CGGGCA
```

Setting up the analysis: Models

Substitution Models

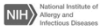

- HKY - Unequal base frequencies and transition/transversion rate ratio
 - Must specify prior and initial estimates for transition/transversion rate ratio
- GTR - Unequal base frequencies and each substitution has its own rate parameter
 - Must specify prior and initial estimates for each substitution rate (relative to C-T rate)

Setting up the analysis: Models

Site Models

- Site heterogeneity models
 - Gamma
 - Modeling rate of change using a discrete gamma distribution
 - Invariant
 - Percent of non-variable sites in the data

Setting up the analysis: Models

Estimating best-fit models and initial parameters: jModelTest



```

Model selected: TVM+I+G
-lnL = 1676.8109
k = 9
AIC = 3371.6218

Base frequencies:
freqA = 0.2259
freqC = 0.2199
freqG = 0.2465
freqT = 0.2137

Substitution model:
Rate matrix
R(a) [A-C] = 0.2494
R(b) [A-G] = 4.8655
R(c) [A-T] = 0.7435
R(d) [C-G] = 0.3907
R(e) [C-T] = 4.8655
R(f) [G-T] = 1.0000



Among-site rate variation
Proportion of invariable sites (I) = 0.6508
Variable sites (V)
Gamma distribution shape parameter = 0.5913
    
```

Setting up the analysis: Models

Site heterogeneity models: The Gamma Distribution


Mean = $k\theta$
 Shape parameter = θ
 Coefficient of Variation = $1/\sqrt{\theta}$

Setting up the analysis: Models

Clock Models

- Strict clock – same rate for all branches
- Relaxed clock – independent rate among branches
 - Exponential or Lognormal distribution of rates
- For contemporaneous data setting a fixed mean substitution rate of 1.0 (uncheck "Estimate") results in node ages as substitutions per site (MrBayes branch lengths)

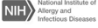
 National Institute of Allergy and Infectious Diseases

NIAD

Setting up the analysis: Models

Tree Prior

- Coalescent
 - constant size
 - exponential growth
 - GMRF Bayesian Skygrid
- Speciation
 - Yule process
 - Birth-Death
- Epidemiology

 National Institute of Allergy and Infectious Diseases


NIAD

Setting up the analysis: Models

Testing Models and Priors

Path Sampling/Stepping Stone analysis

- Estimation of marginal likelihoods under different analysis parameters.
- Invoke on MCMC tab in BEAUti.
- Separate runs necessary for each changed parameter.
- Runs a complete MCMC analysis, then the X PS/SS iterations.

 National Institute of Allergy and Infectious Diseases

NIAD

Setting up the analysis: Models

Testing Models and Priors

Path Sampling/Stepping Stone analysis

NIH National Institute of Allergy and Infectious Diseases

NIAID

Setting up the analysis: Models

Testing Models and Priors

Path Sampling/Stepping Stone analysis

log marginal likelihoods

	Path Sampling	Stepping Stone
HKY/strict clock	-4725.85	-4728.68
HKY+gi/strict	-4515.99	-4518.05
HKY+gi/LN relaxed	-4436.10	-4438.75
GTR/strict clock	-4746.62	-4749.14
GTR+gi/strict	-4526.87	-4529.05
GTR+gi/LN relaxed	-4548.39	-4551.22

NIH National Institute of Allergy and Infectious Diseases

NIAID

Setting up the analysis: Models

Testing Models and Priors

Does the relaxed clock fit the data?

NIH National Institute of Allergy and Infectious Diseases

NIAID

Setting up the analysis: Models

Testing Models and Priors

Does the relaxed clock fit the data?

NIH National Institute of Allergy and Infectious Diseases

NIAD

Setting up the analysis: Models

Testing Models and Priors

Does the relaxed clock fit the data?

NIH National Institute of Allergy and Infectious Diseases

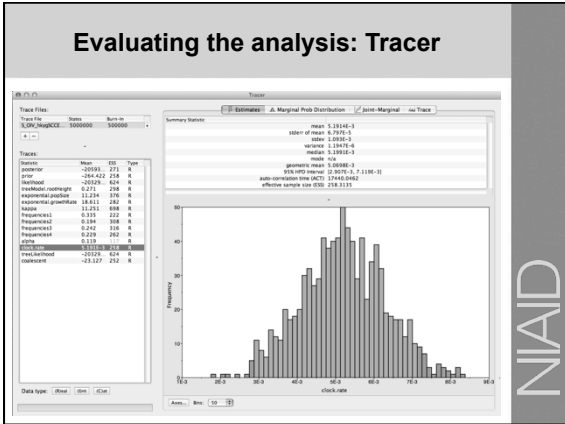
NIAD



Running the analysis: BEAST



- Load your input file
- That's it

NIH National Institute of Allergy and Infectious Diseases

NIAD

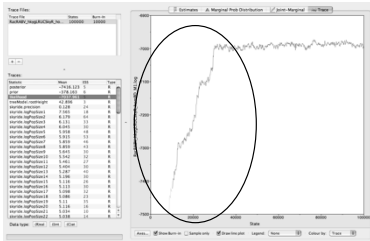


- ### Evaluating the analysis: Tracer
- Check for convergence
 - Evaluating ESS values
 - Viewing behavior of parameter estimates
 - Examining traces
 - Extracting parameter estimates and statistics
- 
- 

- ### Evaluating the analysis: Tracer
- What if my analysis didn't converge?
 - Can I make multiple simultaneous runs?
 - Swarm on Biowulf
- 
- 

Evaluating the analysis: Tracer

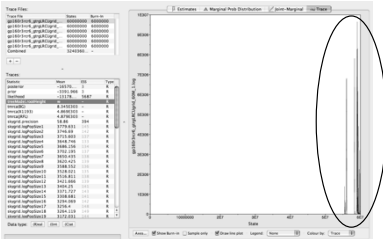
What if my analysis didn't converge?



NIAD

Evaluating the analysis: Tracer

What if my analysis didn't converge?



NIAD

What if my analysis didn't converge?

Partitions Taxa Tips Traits Sites Clocks Trees States Priors Operators MCMC

Priors for model parameters and statistics:

Name	Bound	Description
invariant	0.5	invariant sites
invariantI	0.5	invariant sites for tree set I
invariantII	0.5	invariant sites for tree set II
invariantIII	0.5	invariant sites for tree set III
invariantIV	0.5	invariant sites for tree set IV
invariantV	0.5	invariant sites for tree set V
invariantVI	0.5	invariant sites for tree set VI
invariantVII	0.5	invariant sites for tree set VII
invariantVIII	0.5	invariant sites for tree set VIII
invariantIX	0.5	invariant sites for tree set IX
invariantX	0.5	invariant sites for tree set X
invariantXI	0.5	invariant sites for tree set XI
invariantXII	0.5	invariant sites for tree set XII
invariantXIII	0.5	invariant sites for tree set XIII
invariantXIV	0.5	invariant sites for tree set XIV
invariantXV	0.5	invariant sites for tree set XV
invariantXVI	0.5	invariant sites for tree set XVI
invariantXVII	0.5	invariant sites for tree set XVII
invariantXVIII	0.5	invariant sites for tree set XVIII
invariantXIX	0.5	invariant sites for tree set XIX
invariantXX	0.5	invariant sites for tree set XX
invariantXXI	0.5	invariant sites for tree set XXI
invariantXXII	0.5	invariant sites for tree set XXII
invariantXXIII	0.5	invariant sites for tree set XXIII
invariantXXIV	0.5	invariant sites for tree set XXIV
invariantXXV	0.5	invariant sites for tree set XXV
invariantXXVI	0.5	invariant sites for tree set XXVI
invariantXXVII	0.5	invariant sites for tree set XXVII
invariantXXVIII	0.5	invariant sites for tree set XXVIII
invariantXXIX	0.5	invariant sites for tree set XXIX
invariantXXX	0.5	invariant sites for tree set XXX
invariantXXXI	0.5	invariant sites for tree set XXXI
invariantXXXII	0.5	invariant sites for tree set XXXII
invariantXXXIII	0.5	invariant sites for tree set XXXIII
invariantXXXIV	0.5	invariant sites for tree set XXXIV
invariantXXXV	0.5	invariant sites for tree set XXXV
invariantXXXVI	0.5	invariant sites for tree set XXXVI
invariantXXXVII	0.5	invariant sites for tree set XXXVII
invariantXXXVIII	0.5	invariant sites for tree set XXXVIII
invariantXXXIX	0.5	invariant sites for tree set XXXIX
invariantXXXX	0.5	invariant sites for tree set XXXIV

Prior Distribution: None (Tree Prior Only)

Cancel OK

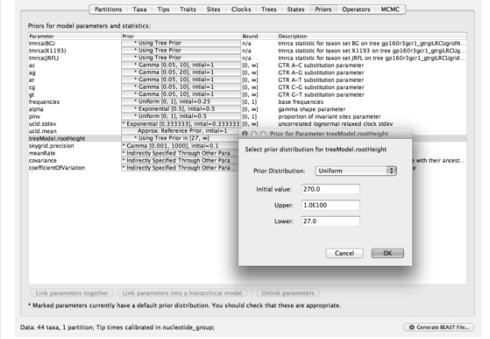
Link parameters together. Link parameters into a hierarchical model. Unlink parameters.

* Marked parameters currently have a default prior distribution. You should check that these are appropriate.

Data: 44 taxa, 1 partition. Tip times calibrated in nucleotide_group. © Cytosine MAAT (P...)

NIAD

What if my analysis didn't converge?



Running BEAST: swarm on Biowulf

- Requires a .swarm file
 - A text file containing

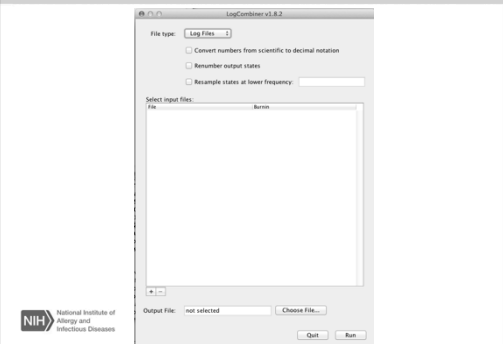
```
beast beastJob_1.xml > beastJob_1out.txt
sleep 2; beast beastJob_2.xml > beastJob_2out.txt
sleep 4; beast beastJob_3.xml > beastJob_3out.txt
sleep 6; beast beastJob_4.xml > beastJob_4out.txt
sleep 8; beast beastJob_5.xml > beastJob_5out.txt
```

- Run in command line

```
[username]$ swarm -f beastInput.swarm --module BEAST
```




Merging output files: LogCombiner





Merging output files: LogCombiner



- Log files vs Tree files
- Selecting files
- Specifying burn-in (number of steps or trees)
- Specifying subsampling
- Specifying output file



NIAD


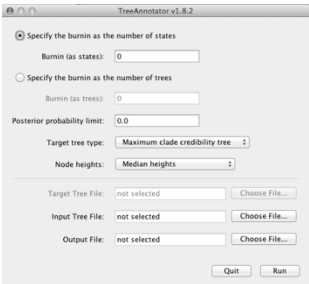
Merging output files: LogCombiner

- Burn in?



NIAD

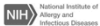
Calculating the tree: TreeAnnotator



NIAD

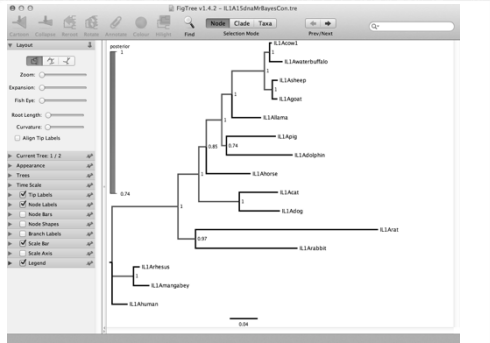
Calculating the tree: TreeAnnotator

- Burn in? Number of trees or the number of steps.
- Tree Type: MCC, Max sum of CC, or target
- Node heights: target, mean, or median
- Specify input and output files



NIAID


Drawing trees: FigTree



NIAID

Drawing trees: FigTree


- Specifying additional values (esp. posterior probabilities)
- Tree appearance
- Ordering branches
- Re-rooting
- Exporting graphics



NIAID

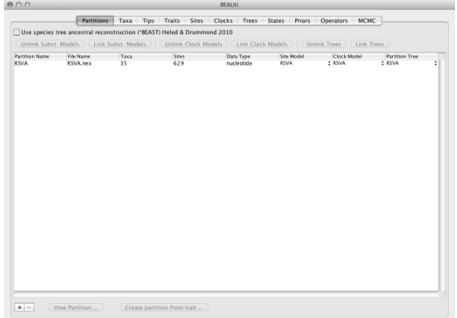
Running BEAST


DEMO



NIAD


Setting up the analysis: BEAUTi






NIAD

Running the analysis: BEAST





NIAD

Evaluating the analysis: Tracer

Summary Statistics:

mean	0.211415
sd	0.120115
mode	0.201015
median	0.211415
range	0.001015 - 0.411415
skewness	0.001015
kurtosis	0.001015
effective sample size	211.1115

NIAD

Merging output files: LogCombiner

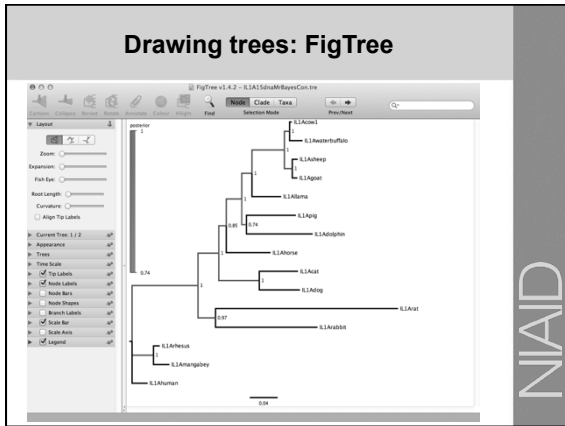
NIH National Institute of Allergy and Infectious Diseases

NIAD

Calculating the tree: TreeAnnotator

NIH National Institute of Allergy and Infectious Diseases

NIAD



BEAST2

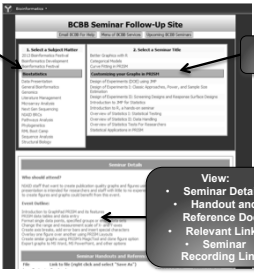
- Still ... **B**ayesian **E**volutionary **A**nalysis **S**ampling **T**rees
- Modular rewrite of the BEAST software
- Various evolutionary analyses performed through a system of independent software packages.
- Access software, documentation, etc., through the website beast2.org
- Still a few bugs in the system...

NIAD

Seminar Follow-Up Site

• For access to past recordings, handouts, slides visit this site from the NIH network: <http://collab.niaid.nih.gov/sites/research/SIG/Bioinformatics/>

1. Select a Subject Matter



2. Select a Topic

Recommended Browsers:

- IE for Windows,
- Safari for Mac (Firefox on a Mac is incompatible with NIH Authentication technology)



Login

- If prompted to log in use "NIH" in front of your username

NIAD

Next?

- Time-structured phylogenies
- Estimating demographic parameters
 - GMRF skyride analysis



58

Thank you



59
