National Institute of Allergy and Infectious Diseases

**Phylogenetics and Sequence Analysis**

# Lecture 4
# Tree Building 2

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

Fall 2015

NIAID — National Institute of Allergy and Infectious Diseases

---

# Course Organization

- Building a clean sequence
- Collecting homologs
- Aligning your sequences
- **Building trees**
- Further Analysis

NIAID — National Institute of Allergy and Infectious Diseases

---

# Tree building, so far

- Generate a distance tree
- Generate maximum parsimony tree(s)
- Calculate bootstrap support

NIAID — National Institute of Allergy and Infectious Diseases

## What's next?

Statistical Methods for Calculating Trees

- Maximum Likelihood
- Bayesian Phylogenetics

NIAID

NIH National Institute of Allergy and Infectious Diseases

---

## Optimality Criterion: Likelihood

Calculating likelihood

The Tree    The Data

S1: AAG

S2: AAA

S3: GGA

S4: AGA

$L(Tree) = Prob(Data|Tree) = \prod_i Prob(Data^{(i)}|Tree)$

NIAID

NIH National Institute of Allergy and Infectious Diseases

---

## Optimality Criterion: Likelihood

Calculating likelihood: Setting parameters

$L(Tree) = Prob(Data|Tree) = \prod_i Prob(Data^{(i)}|Tree)$

What values do you use for the substitution model?

Run jModelTest (or ProtTest for protein MSAs)

NIAID

NIH National Institute of Allergy and Infectious Diseases
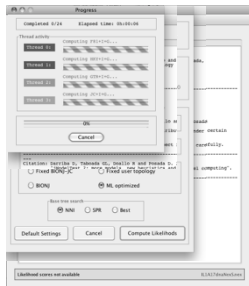
## Optimality Criterion: Likelihood

Calculating likelihood: jModelTest



## Optimality Criterion: Likelihood

Calculating likelihood: jModelTest



## Optimality Criterion: Likelihood

Calculating likelihood: jModelTest Results

## Optimality Criterion: Likelihood

Calculating likelihood: jModelTest Results

Substitution models
http://www.molecularevolution.org/resources/models/nucleotide

## Optimality Criterion: Likelihood

Calculating likelihood: jModelTest Results

Substitution models
http://www.molecularevolution.org/resources/models/nucleotide

## The Gamma Distribution

Mean = kθ   Shape parameter = θ
Coefficient of Variation = 1/√θ

## Calculating likelihood: Programs

PAUP* – Commercial, NIH Biowulf, or NIAID HPC
        DNA only
PHYLIP – Download, NIH Biowulf, or NIAID HPC
        dnaml and proml programs
PAML – Download, NIH Biowulf, or NIAID HPC
RaxML – Download or NIH BioWulf or webserver
PhyML – Download or NIAID HPC or webserver
GARLi – Download or NIAID HPC or webserver
MEGA6 - Download

Generally the user has more flexibility with a local program.
        But local programs can hog your computer.

NIH National Institute of Allergy and Infectious Diseases

NIAID

## Calculating likelihood: Programs

GARLi – Genetic Algorithm for Rapid Likelihood Inference

Input data format:
    PHYLIP or nexus **SEQUENTIAL**
Output files:
    *filename*.best.tree
    *filename*.best.all.tree
    *filename*.*.log – intermediate run output

NIH National Institute of Allergy and Infectious Diseases

NIAID

## Calculating likelihood: GARLi

Input data format:

PHYLIP or Nexus

SEQUENTIAL

not

INTERLEAVED

and now fasta

NIH National Institute of Allergy and Infectious Diseases

NIAID

## Calculating likelihood: GARLi

Specifying input parameters:
garli.conf

www.nescent.org/wg_garli/GARLI_Configuration_Settings

---

## Calculating likelihood: GARLi

### Running the program

---

## Calculating likelihood: GARLi

### Running the program

## Calculating likelihood: GARLi

### Running the program

```
Terminal — bash — 109×34

>>>Completed Search rep 4 (of 4)<<<

##################################################
Completed 4 replicate search(es) (of 4).

NOTE: Unless the following output indicates that search replicates found the
      same topology, you should assume that they found different topologies.
Results:
Replicate 1 : -4225.8261 (best)
Replicate 2 : -4225.8261    (same topology as 1)
Replicate 3 : -4225.8262    (same topology as 1)
Replicate 4 : -4225.8261    (same topology as 1)

Parameter estimates across search replicates:
        r(AC)  r(AG)  r(AT)  r(CG)  r(CT)  r(GT) pi(A) pi(C) pi(G) pi(T) alpha  pinv
rep 1:  1.462  9.208  2.05 0.5893  8.434      1 0.244 0.288 0.250 0.218 0.478 0.134
rep 2:  1.462  9.208 2.051 0.5894  8.434      1 0.244 0.288 0.250 0.218 0.478 0.134
rep 3:  1.462  9.209 2.051 0.5894  8.435      1 0.244 0.288 0.250 0.218 0.478 0.134
rep 4:  1.461  9.207  2.05 0.5893  8.434      1 0.244 0.288 0.250 0.218 0.478 0.134

Treelengths:
         TL
rep 1:  11.797
rep 2:  11.797
rep 3:  11.805
rep 4:  11.798

Saving final trees from all search reps to RBD_DNA_gtrg4i_1.best.all.tre

Saving final tree from best search rep (#1) to RBD_DNA_gtrg4i_1.best.tre
##################################################
w8107@hgagw:bin wollenbergk$
```

## Calculating likelihood: GARLi



The Tree

National Institute of Allergy and Infectious Diseases

## Calculating likelihood: GARLi

### Bootstrap results?

GARLi **does not** calculate a bootstrap consensus!
You must use another piece of software for this.

Dendropy:sumtrees – a very efficient python package for summarizing collections of phylogenetic trees
Dendropy – Freely distributed from:
        https://pythonhosted.org/DendroPy
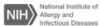sumtrees has also been implemented on the NIAID HPC

National Institute of Allergy and Infectious Diseases

## Bayesian Analysis

Calculating the posterior probability of the evolutionary parameters

$$\Pr(\tau, v, \theta | \text{Data}) = \frac{\Pr(D|\tau, v, \theta) \times \Pr(\tau, v, \theta)}{\Pr(D)}$$

where:
τ = tree topology
v = branch lengths
θ = substitution parameters

NIAID

NIH National Institute of Allergy and Infectious Diseases

## What is Bayesian Analysis?

- Calculation of the probability of parameters (tree, substitution model) given the data (sequence alignment)
- p(θ|D) = (Likelihood x Prior) / probability of the data
- p(θ|D) = p(D|θ)p(θ) / p(D)

NIAID

NIH National Institute of Allergy and Infectious Diseases

## 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

NIAID

NIH National Institute of Allergy and Infectious Diseases

## What is MCMC?

### Markov Chain Monte Carlo

Markov chain

Monte Carlo

One link in the chain

Choosing a link

NIAID

NIH National Institute of Allergy and Infectious Diseases

---

## Bayesian Analysis

Markov Chain example: Jukes-Cantor

$\mu t = 0.25$

|   | A | C | G | T |
|---|---|---|---|---|
| A | 0.5259 | 0.158 | 0.158 | 0.158 |
| C | 0.158 | 0.5259 | 0.158 | 0.158 |
| G | 0.158 | 0.158 | 0.5259 | 0.158 |
| T | 0.158 | 0.158 | 0.158 | 0.5259 |

NIAID

NIH National Institute of Allergy and Infectious Diseases

---

## Bayesian Analysis

Exploring the posterior probability distribution

The **posterior probability** of a specific tree is the number of times the Markov Chain visits that tree

**Posterior probability distribution** is summarized by the clade probabilities.

NIAID

NIH National Institute of Allergy and Infectious Diseases

## Bayesian Analysis

### Using MrBayes

- Input format = Nexus
- Choose a substitution model (jModelTest)
- Check for convergence

### Using BEAST

- Input format = XML (made using BEAUTi program)
- Choose a substitution model (jModelTest)
- Check for convergence (using Tracer program)

NIAID

---

## Bayesian Analysis

Running MrBayes: Model parameters

```
MrBayes> lset nst=6 rates=invgamma
MrBayes> showmodel
MrBayes> mcmc ngen=20000 samplefreq=100
printfreq=100 diagnfreq=100 burninfrac=0.25
```

NIAID

---

## Bayesian Analysis

Running MrBayes: Setting the Priors

- Generally, the default priors work well
- These are known as "uninformative " priors
- For implementing the Jukes-Cantor model, change statefreqpr to "fixed"

NIAID

**Bayesian Analysis**

Running MrBayes: Setting the Priors

Amino acid substitution models
- Poisson - equal rates, equal state frequencies
- Blosum62
- Dayhoff
- Mtrev, Mtmamm - mitochondrial models
- mixed - Let MrBayes choose among the many fixed-rate models

NIAID

NIH National Institute of Allergy and Infectious Diseases

---

**Bayesian Analysis**

Running MrBayes: General

- burnin - initial portion of the run to discard
  - Generally, 25% of the samples
- samplefreq - how often to sample the Markov chain
  - More frequently for small analyses
  - Less frequently for low-complexity data
- printfreq - how often output is sent to the log file(s)

NIAID

NIH National Institute of Allergy and Infectious Diseases

---

**Bayesian Analysis**

Running MrBayes: General

```
#NEXUS
begin mrbayes;
    set autoclose=yes nowarn=yes;
    execute /pathtodata/InputData.nex;
    lset nst=6 rates=invgamma;
    mcmc stoprule=yes stopval=0.009;
end;
```

NIAID

NIH National Institute of Allergy and Infectious Diseases

## Bayesian Analysis

### Running MrBayes: General

Burn in



## Bayesian Analysis

### Running MrBayes: Summarizing results

```
MrBayes> sump (burninfrac=0.25)
MrBayes> sumt (burninfrac=0.25)
```

## Bayesian Analysis

### Using MrBayes: Convergence

```
Chain results:

      1 -- [-5762.003] (-5753.828) [...6 remote chains...]
   1000 -- (-4832.654) (-4844.806) [...6 remote chains...] -- 0:16:39

   Average standard deviation of split frequencies: 0.143471

   2000 -- (-4748.109) (-4762.679) [...6 remote chains...] -- 0:24:57

            *************** [SNIP] ***************

 999000 -- (-4886.847) [-4876.966] [...6 remote chains...] -- 0:00:06

Average standard deviation of split frequencies: 0.002371
1000000 -- (-4885.621) [-4889.536] [...6 remote chains...] -- 0:00:00

Average standard deviation of split frequencies: 0.002413
```

## Bayesian Analysis

### Using MrBayes: Convergence

#### Log-probability plot appears stochastic

```
Overlay plot for both runs:
(1 = Run number 1; 2 = Run number 2; * = Both runs)
+------------------------------------------------------------+ -4879.06
|     2                                        1             |
|     1          1  11              2                    1   |
|        2    1 2                        1  1 1              |
|     1     2 1                        2       11 1          |
|  2   1 2          *   1        1    2    1   1 * 2  1   22  |
|          1  2       12 * 1222     2     11 22  1 21     1 2 |
|  2 *21 1 2              221         2   2     211 21 1  1  1|
|     1 2    12  2  2 1  1  12   2 2           2   2    1     |
|           122      2  2 1     21    1  22           1  1 2  |
|  1              1 2              1        2    2           |
|                 1              11         2 2              |
|  1                               1                  22     |
|                                                           |
|                  2                                        |
|                                                           |
|           2                                               |
* +-----+-----+-----+-----+-----+-----+-----+-----+-----+---+ -4882.76
^
100000                                           1000000
```
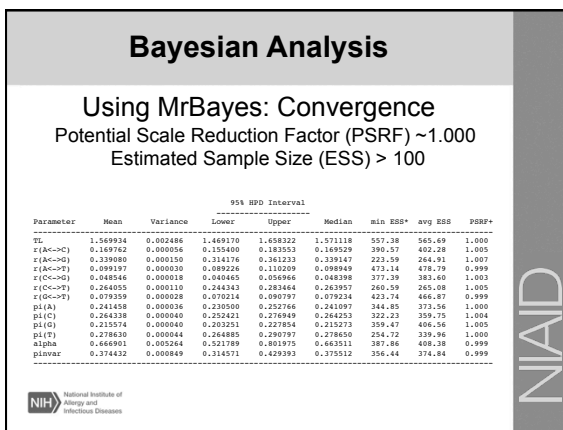
## Bayesian Analysis

### Using MrBayes: Convergence
Potential Scale Reduction Factor (PSRF) ~1.000
Estimated Sample Size (ESS) > 100

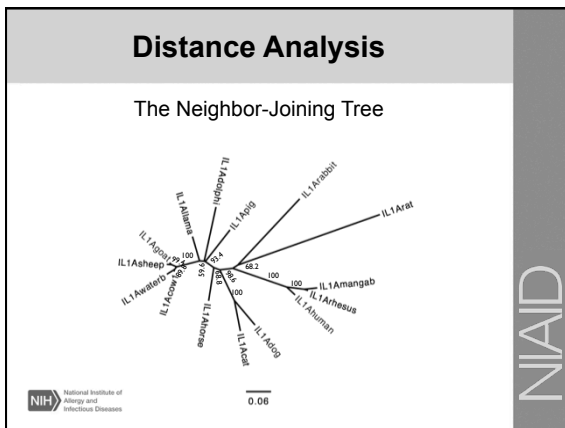|  |  |  | 95% HPD Interval | | | | | |
|  |  |  | --- | --- | | | | |
| Parameter | Mean | Variance | Lower | Upper | Median | min ESS* | avg ESS | PSRF+ |
| TL | 1.569934 | 0.002486 | 1.469170 | 1.658322 | 1.571118 | 557.38 | 565.69 | 1.000 |
| r(A<->C) | 0.169762 | 0.000056 | 0.155400 | 0.183553 | 0.169529 | 390.57 | 402.28 | 1.005 |
| r(A<->G) | 0.339080 | 0.000150 | 0.314176 | 0.361233 | 0.339147 | 223.59 | 264.91 | 1.007 |
| r(A<->T) | 0.099197 | 0.000030 | 0.089226 | 0.110209 | 0.098949 | 473.14 | 478.79 | 0.999 |
| r(C<->G) | 0.048546 | 0.000018 | 0.040465 | 0.056966 | 0.048398 | 377.39 | 383.60 | 1.003 |
| r(C<->T) | 0.264055 | 0.000110 | 0.244343 | 0.283464 | 0.263957 | 260.59 | 265.08 | 1.005 |
| r(G<->T) | 0.079359 | 0.000028 | 0.070214 | 0.090797 | 0.079234 | 423.74 | 466.87 | 0.999 |
| pi(A) | 0.241458 | 0.000036 | 0.230500 | 0.252766 | 0.241097 | 344.85 | 373.56 | 1.000 |
| pi(C) | 0.264338 | 0.000040 | 0.252421 | 0.276949 | 0.264253 | 322.23 | 359.75 | 1.004 |
| pi(G) | 0.215574 | 0.000040 | 0.203251 | 0.227854 | 0.215273 | 359.47 | 406.56 | 1.005 |
| pi(T) | 0.278630 | 0.000044 | 0.264885 | 0.290797 | 0.278650 | 254.72 | 339.96 | 1.000 |
| alpha | 0.666901 | 0.005264 | 0.521789 | 0.801975 | 0.663511 | 387.86 | 408.38 | 0.999 |
| pinvar | 0.374432 | 0.000849 | 0.314571 | 0.429393 | 0.375512 | 356.44 | 374.84 | 0.999 |

## Bayesian Analysis

### The Consensus Tree

## Likelihood Analysis

The Consensus Tree



## Distance Analysis

The Neighbor-Joining Tree



## Parsimony Analysis
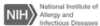
The Bootstrap Consensus Tree

## Tree Building, in conclusion

Where have we been? What have we done?

- How to generate trees using distance, parsimony, and likelihood
  - How to calculate bootstrap support
- Bayesian exploration of phylogeny posterior distribution
- Always use more than one tree generation algorithm
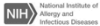- Look for consensus and investigate disagreement

NIH National Institute of Allergy and Infectious Diseases

NIAID

---

## Visualizing Trees

- **FigTree**
- **Dendroscope**

NIH National Institute of Allergy and Infectious Diseases

NIAID

---

## REFERENCES

✱ Inferring Phylogenies. J. Felsenstein. 2004
A good general reference written in Professor Felsenstein's unique style.

✱ The Phylogenetic Handbook. Edited by P. Lemey, et al. 2009
A very thorough exploration of theory and practice.

✱ Biological Sequence Analysis. R. Durbin, et al. 1998.
A good introduction to maximum likelihood and hidden Markov models (HMM).

✱ Phylogenetic Trees Made Easy. B. Hall
1st, 2nd Edition - A PAUP and PHYLIP manual, cookbook style.

3rd Edition - A MEGA4 manual, among other things, cookbook style.

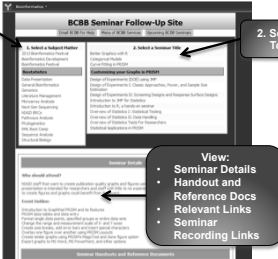4th Edition - MEGA5 and MrBayes 3.2 manuals, cookbook style.

NIH National Institute of Allergy and Infectious Diseases

NIAID

**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/



**Retrieving Slides/Handouts**



**Retrieving Slides/Handouts**

## Questions?

Consultation & Advice | Software Development | Biocomputing Resources

**ScienceApps@niaid.nih.gov**

BCBB Bioinformatics and Computational Biosciences Branch
NIAID Office of Cyber Infrastructure and Computational Biology

NIH National Institute of Allergy and Infectious Diseases

NIAID

49

## Next

**Tuesday, 10 November at 1300**

*Making Publication-Quality Phylogenetic Tree Figures*

NIH National Institute of Allergy and Infectious Diseases

NIAID

50