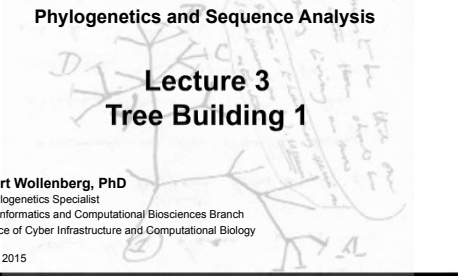National Institute of Allergy and Infectious Diseases

**Phylogenetics and Sequence Analysis**

# Lecture 3
# Tree Building 1

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

Fall 2015
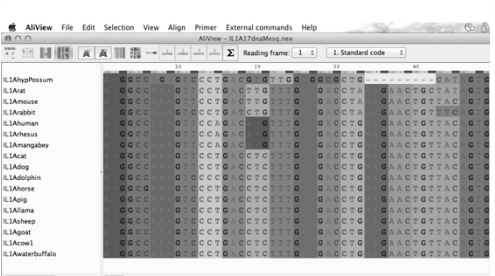
NIAID

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

---

# Multiple Sequence Alignment



NIAID

National Institute of Allergy and Infectious Diseases

## What's next?

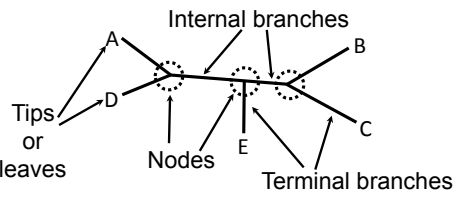Building trees with our MSA

NIAID

---

## What is a phylogenetic tree?

- Reconstruction of biological history
- Based on similarities and differences among homologous attributes (characters) of the entities under scrutiny
- Molecular characters (sequences, usually) are most often found only in extant organisms

NIAID

---

## What is a phylogenetic tree?



NIAID

## What is a phylogenetic tree?

Unrooted       Rooted

NIAID

NIH National Institute of Allergy and Infectious Diseases

## Two approaches to tree building

- Application of an algorithm to build the best tree from the data
- Evaluation of multiple possible best trees using an optimality criterion

NIAID

NIH National Institute of Allergy and Infectious Diseases

## The algorithm approach:
## Distance Methods

- Distance calculated based on a specific substitution model (J-C, Kimura, BLOSUM64, etc.)
- Distances from each sequence to all others are calculated and stored in a matrix
- Tree then calculated from the distance matrix using a specific tree-building algorithm

NIAID

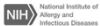NIH National Institute of Allergy and Infectious Diseases

## The algorithm approach: Distance Methods

Tree-Building Algorithms

- UPGMA
- Neighbor-Joining

NIAID

NIH National Institute of Allergy and Infectious Diseases

---

## The algorithm approach: Neighbor-joining Calculation

|        | A       | B       | C       | D       | E       | R      | R/(N−2) |
|--------|---------|---------|---------|---------|---------|--------|---------|
| A      | –       | 0.1715  | 0.2147  | 0.3091  | 0.2326  | 0.9279 | 0.3093  |
| B      | −0.4766 | –       | 0.2991  | 0.3399  | 0.2058  | 1.0163 | 0.3388  |
| C      | −0.4905 | −0.4356 | –       | 0.2795  | 0.3943  | 1.1876 | 0.3959  |
| D      | −0.4527 | −0.4514 | −0.5689 | –       | 0.4289  | 1.3574 | 0.4525  |
| E      | −0.4972 | −0.5535 | −0.4221 | −0.4441 | –       | 1.2616 | 0.4205  |

C to Node 1 distance = 0.2795/2 + (0.3959 − 0.4525)/2 = 0.1114
D to Node 1 distance = 0.2795 − 0.1114 = 0.1681

A to Node 1 distance = (0.2147 + 0.3091 − 0.2795)/2 = 0.1222
B to Node 1 distance = (0.2991 + 0.3399 − 0.2795)/2 = 0.1798
E to Node 1 distance = (0.3943 + 0.4298 − 0.2795)/2 = 0.2719

NIAID

NIH National Institute of Allergy and Infectious Diseases

Hillis, Moritz, and Mable 1996, p. 489

---

## The algorithm approach: Neighbor-joining Calculation

|        | A       | B       | E       | Node 1  | R      | R/(N−2) |
|--------|---------|---------|---------|---------|--------|---------|
| A      | –       | 0.1715  | 0.2326  | 0.1222  | 0.5263 | 0.2631  |
| B      | −0.3701 | –       | 0.2058  | 0.1798  | 0.5571 | 0.2785  |
| E      | −0.3856 | −0.4278 | –       | 0.2719  | 0.7103 | 0.3551  |
| Node 1 | −0.4278 | −0.3856 | −0.3701 | –       | 0.5739 | 0.2869  |

A to Node 2 distance = 0.1222/2 + (0.2631 − 0.2869)/2 = 0.0492
Node 1 to Node 2 distance = 0.1222 − 0.0492 = 0.0730

B to Node 2 distance = (0.1715 + 0.1798 − 0.1222)/2 = 0.1146
E to Node 2 distance = (0.2326 + 0.2719 − 0.1222)/2 = 0.1912

NIAID

NIH National Institute of Allergy and Infectious Diseases

Hillis, Moritz, and Mable 1996, p. 489

## Slide 1

**The algorithm approach:**
**Neighbor-joining Calculation**

| | B | E | Node 2 | R | R/(N−2) |
|---|---|---|---|---|---|
| B | – | 0.2058 | 0.1146 | 0.3204 | 0.3204 |
| E | −0.5116 | – | 0.1912 | 0.3970 | 0.3970 |
| Node 2 | −0.5116 | −0.5116 | – | 0.3058 | 0.3058 |

B to Node 3 distance = 0.1146/2 + (0.3204 − 0.3058)/2 = 0.0646
Node 2 to Node 3 distance = 0.1146 − 0.0646 = 0.0500

E to Node 3 distance = (0.2058 0.1912 − 0.1146)/2 = 0.1412

Hillis, Moritz, and Mable 1996, p. 489

## Slide 2

**The algorithm approach:**
**Neighbor-joining Calculation**

D
0.168
A 0.049 B
0.073 0.065
0.050
0.111 0.141
C E

## Slide 3

## The optimality criterion approach

- Build a tree or trees
- Evaluate the tree(s) using a specific numerical optimality criterion
- Most common optimality criteria
  - Maximum parsimony
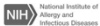  - Maximum likelihood
- Explore tree space to find the optimal tree

## Optimality Criterion: Parsimony

Occam's Razor: The simplest explanation is the preferred explanation.

The tree requiring the minimal number of changes is the optimal tree.

A step is any change in the data from one state to another

NIAID

NIH National Institute of Allergy and Infectious Diseases

---

## The optimality criterion approach

- Build the initial tree
  - Construct a neighbor-joining tree
  - Stepwise addition

- Calculate the tree score
  - Count steps (parsimony)
  - Calculate likelihood of the data given the tree

- Explore tree space
  - Branch swapping
    - Tree bisection and reconnection (TBR)

- Is this the best tree? (Stopping criteria)

NIAID

NIH National Institute of Allergy and Infectious Diseases

---

## The optimality criterion approach

### Building the initial tree

- Stepwise addition

- Choose three taxa and join

- Random, or closest

- Select a new taxon to add

- Calculate the optimal 4-taxa tree

- Repeat until all taxa are joined

NIAID

NIH National Institute of Allergy and Infectious Diseases

## The optimality criterion approach

**Building the initial tree**



---

## The optimality criterion approach

**Exploring tree space: Branch swapping**

- Nearest neighbor interchange
- Subtree pruning and regrafting
- Tree bisection and reconnection

---

## The optimality criterion approach

**Branch swapping: Tree bisection and reconnection**

**The optimality criterion approach**

**Exploring tree space**

**Beware!** Hill climbing can often lead to local maxima rather than a global solution.

NIAID

NIH National Institute of Allergy and Infectious Diseases

---

**The optimality criterion approach**

Exploring tree space

NIAID

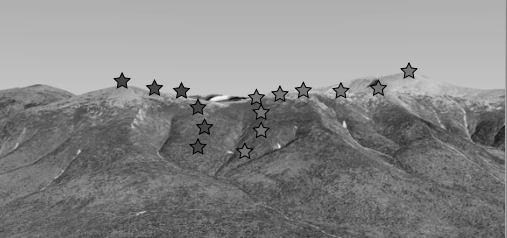NIH National Institute of Allergy and Infectious Diseases

---

**The optimality criterion approach**

**Is this tree optimal?**

Counting changes (Fitch parsimony)

The tree          The data

S1: AAG
S2: AAA
S3: GGA
S4: AGA

NIAID

NIH National Institute of Allergy and Infectious Diseases

**The optimality criterion approach**

**Is this tree optimal?**

Counting changes (Fitch parsimony)

Position 1   The tree          The data

{A} ∩ {A} = {A}                    S1: A
0 step                             S2: A

{A} ∩ {AG} = {A}
0 step

{A} ∩ {G} = Ø                      S3: G
{A} U {G} = {AG}                   S4: A
1 step

NIAID
NIH National Institute of Allergy and Infectious Diseases

---

**The optimality criterion approach**

**Is this tree optimal?**

Counting changes (Fitch parsimony)

Position 2   The tree          The data

{A} ∩ {A} = {A}                    S1: A
0 step                             S2: A

{A} U {G} = {AG}
1 step

{G} ∩ {G} = {G}                    S3: G
0 step                             S4: G

NIAID
NIH National Institute of Allergy and Infectious Diseases

---

**The optimality criterion approach**

**Is this tree optimal?**

Counting changes (Fitch parsimony)

Position 3   The tree          The data

{A} U {G} = {AG}                   S1: A
1 step                             S2: G

{AG} ∩ {A} = {A}
0 step

{A} ∩ {A} = {A}                    S3: A
0 step                             S4: A

NIAID
NIH National Institute of Allergy and Infectious Diseases

9

## The optimality criterion approach

**Is this tree optimal?**

Counting changes (Fitch parsimony)

The tree          The data

S1: AAG

S2: AAA          Total steps: 3
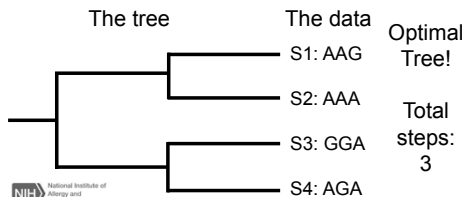
S3: GGA

S4: AGA

NIAID

NIH National Institute of Allergy and Infectious Diseases

---

## The optimality criterion approach

**Is this tree optimal?**

Counting changes (Fitch parsimony)

The tree          The data

S1: AAG

S3: GGA          Total steps: 4

S2: AAA

S4: AGA

NIAID

NIH National Institute of Allergy and Infectious Diseases

---

## The optimality criterion approach

**Is this tree optimal?**

Counting changes (Fitch parsimony)

The tree          The data

S1: AAG

S4: AGA          Total steps: 4

S2: AAA

S3: GGA

NIAID

NIH National Institute of Allergy and Infectious Diseases

## The optimality criterion approach

### Is this tree optimal?

Counting changes (Fitch parsimony)

The tree     The data

S1: AAG    Optimal Tree!

S2: AAA

S3: GGA    Total steps: 3

S4: AGA

NIAID

National Institute of Allergy and Infectious Diseases

---

## Multiple Sequence Alignment

Export MSA as FASTA and PHYLIP formats

FASTA            PHYLIP

NIAID

National Institute of Allergy and Infectious Diseases

---

## Multiple Sequence Alignment

Phylogenetics program input file formats

PHYLIP
1st line: Number of sequences(space)Number of sites
2nd line: Sequence ID (10 characters max) Sequence

```
      9  1823
HCVT050   GGTCTTGGTCTACTGTGAGCGAGGAGGCCGGTGAGGACGTCGTCTGCTGC
HCVT142   GGTCTTGGTCTACCGTGAGTGAGGAGGCCACTGAGGACGTCGTCTGCTGC
HCVT169   GGTCTTGGTCTACCGTGAGCGAGGAGGCTAGTGAGGACGTCGTCTGCTGC
SE0307168 GGTCGTGGTCCACCGTGAACGAGGAGGCTGGTGAGGACGTCGTCTGCTGC
HCVT221   GGTCTTGGTCTACCGTGAGCGAGGAGGCCAGTGAAGACGTTGTCTGCTGC
MD2_2     GGTCTTGGTCTACTGTAAGCGAGGAGGCTAGTGAAGACGTCGTCTGCTGC
HCV1b     GGTCTTGGTCTACCGTGAGCGAGGAGGCTGGTGAGGATGTCGTCTGCTGC
Cont1g0001GGTCTTGGTCTACCGTGAGCGAGGAGGCTAGTGAGGACGTCGTCTGCTGC
HCVT140   GGTCTTGGTCTACTGTGAGCGAGGAGGCTAGTGAGGATGTCGTCTGTTGC
```

NIAID

National Institute of Allergy and Infectious Diseases

## Multiple Sequence Alignment

Phylogenetics program input file formats

NEXUS

```
#NEXUS
[Name: HCVT050        Len: 1823  Check: 5A341084]
[Name: HCVT142        Len: 1823  Check: AB5C0B76]
[Name: HCVT169        Len: 1823  Check: 7EAF66DA]
[Name: SE03071689     Len: 1823  Check: 1EFF8405]
[Name: HCVT221        Len: 1823  Check: 3D0C96F0]
[Name: MD2_2          Len: 1823  Check: 1E2A0948]
[Name: HCV1b          Len: 1823  Check: BC29D7FB]
[Name: Contig0001     Len: 1823  Check: CD240524]
[Name: HCVT140        Len: 1823  Check: 2A5C0D4E]

begin data;
    dimensions ntax=9 nchar=1823;
    format datatype=dna interleave missing=-;
matrix
HCVT050   GGTCTTGGTCTACTGTGAGC GAGGAGGCCGGTGAGGACGT
HCVT142   GGTCTTGGTCTACCGTGAGT GAGGAGGCCACTGAGGACGT
HCVT169   GGTCTTGGTCTACCGTGAGC GAGGAGGCTAGTGAGGACGT
SE03071689 GGTCGTGGTCCACCGTGAAC GAGGAGGCTGGTGAGGACGT
HCVT221   GGTCTTGGTCTACCGTGAGC GAGGAGGCCAGTGAAGACGT
MD2_2     GGTCTTGGTCTACTGTAAGC GAGGAGGCTAGTGAGGACGT
HCV1b     GGTCTTGGTCTACCGTGAGC GAAGAGGCTGGTGAGGACGT
Contig000 GGTCTTGGTCTACCGTGAGC GAGGAGGCTAGTGAGGACGT
HCVT140   GGTCTTGGTCTACTGTGAGC GAGGAGGCTAGTGAGGATGT
;
end;
```

## Multiple Sequence Alignment

Phylogenetics program input data guidelines

• Make sequence IDs different in the first ten characters
• Only letters, numbers, and "_" in sequence IDs
• Make sure all sequences overlap each other

## How reliable are my trees?

Bootstrapping (nonparametric)

**Bootstrapping**
the ideal world

Build replicates by resampling from unlimited data



**Bootstrapping**
the real world

Build pseudoreplicates of unlimited data by sampling with replacement from limited data



**Let's build some trees!**

Phylogenetic programs (distance and parsimony)

PAUP* – Phylogenetic Analysis Using Parsimony

　　　*and other methods

PHYLIP – a suite of phylogenetic programs

MEGA – An integrated phylogenetic analysis package

**Let's build some trees!**



**MEGA6**

Importing a FASTA-formatted alignment
File > Convert to MEGA format
Navigate to FASTA file
Data Format > .fasta (FASTA format)



**MEGA6**

MEGA data file format

**MEGA6**

Imported sequence alignment: Alignment Explorer



**MEGA6**

Imported sequence alignment: Alignment Explorer



**MEGA6**

Imported sequence alignment: Alignment Explorer

## MEGA6

Phylogeny Construction: Neighbor-Joining



## MEGA6

Phylogeny Construction: Neighbor-Joining



## MEGA6

Phylogeny Construction: Neighbor-Joining

**MEGA6**

Phylogeny Construction: Neighbor-Joining



**MEGA6**

Phylogeny Construction: Neighbor-Joining



**MEGA6**

Phylogeny Construction: The Caption Option

## Recapitulation

Where have we been? What have we done?

- What is a phylogenetic tree?
- How to build trees
  - Distance
  - Parsimony
- Calculated bootstrap support



## 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/

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



## Retrieving Slides/Handouts

## Retrieving Slides/Handouts

**BCBB Seminar Follow-Up Site**

Email BCBB For Help    Menu of BCBB Services    Upcoming BCBB Seminars

| 1. Select a Subject Matter | 2. Select a Seminar Title |
| --- | --- |
| 2013 Bioinformatics Festival | **Building Trees – Phylogenetics I** |
| Bioinformatics Development | Building Trees – Phylogenetics II |
| Bioinformatics Festival | Homology Searching and Sequence Alignment |
| Biostatistics | Introduction to Phylogenetics and Sequence Assembly |
| Data Presentation | Making Presentation Quality Phylogenetic Trees |
| General Bioinformatics | Making Publication-Quality Trees Figures |
| Genomics | More Tools for Adding Genomic and Functional Context to Your Data |
| Literature Management | Pathogen Analysis Using BEAST |
| Microarray Analysis | Phylogenetics: Molecular Evolution I |
| Next-Gen Sequencing | Phylogenetics: Molecular Evolution II |
| NIAID BRCs | Selection Analysis Using HyPhy |
| Pathway Analysis | Selection Analysis Using PAML |
| **Phylogenetics** | |
| XML Boot Camp | |
| Sequence Analysis | |
| Structural Biology | |

This lecture

These slides

**Seminar Details**

This lecture is part of a six-part series on Phylogenetics presented by the NIAID Bioinformatics and Computational Biosciences Branch (BCBB).

This course will cover:
What is a tree/phylogeny?
Distance methods for tree building
Optimality criterion methods for tree building: Parsimony
Input file formats
Measuring tree reliability: Bootstrapping

**Seminar Handouts and Reference Documents**

| File | Link to file (right click and select "Save As") | Created | Modified |
| --- | --- | --- | --- |
| txt | Contig0001_ClwPhY1.txt | 9/12/2011 1:57 PM | 9/12/2011 1:58 PM |
| zip | PhylogeneticsLecture1_F2011.zip | 9/12/2011 1:30 PM | 9/12/2011 1:32 PM |
| zip | PhylogeneticsLecture2_F2011.zip | 11/11/2012 3:09 PM | 11/11/2012 3:11 PM |
| pdf | PhyloSeqAnFall2010_2.pdf | 11/22/2010 3:24 PM | 11/22/2010 3:36 PM |
| pdf | PhyloSeqAnApr2010_2.pdf | 1/28/2010 11:43 AM | 1/28/2010 11:43 AM |

**Links Relevant to this Seminar**

There are no items to show in this view.

**Seminar Recording Links**

| Seminar Date | Web Seminar Recording URL |
| --- | --- |
| 1/28/2010 9:30 AM | https://webmeeting.nih.gov/p402447561/ |
| 2/29/2013 9:30 AM | https://webmeeting.nih.gov/p70391841/ |

NIAID

NIH National Institute of Allergy and Infectious Diseases

58

## 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

NIAID

NIH National Institute of Allergy and Infectious Diseases

59

## Next

### Tuesday, 17 November at 0930

jModeltest: Fitting analysis parameters to your data

GARLi: Genetic Algorithm for Rapid Likelihood inference

MrBayes: Robust statistical phylogeny inference

NIAID

NIH National Institute of Allergy and Infectious Diseases

60