National Institute of Allergy and Infectious Diseases

# Virus Sequence Alignment and Phylogenetics

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

National Institute of
Allergy and
Infectious Diseases

21 June 2019

# We Are BCBB!



- Group of (so far) 52
  - Bioinformatics Software Developers
  - Computational Biologists
  - Project Management & Analysis Professionals
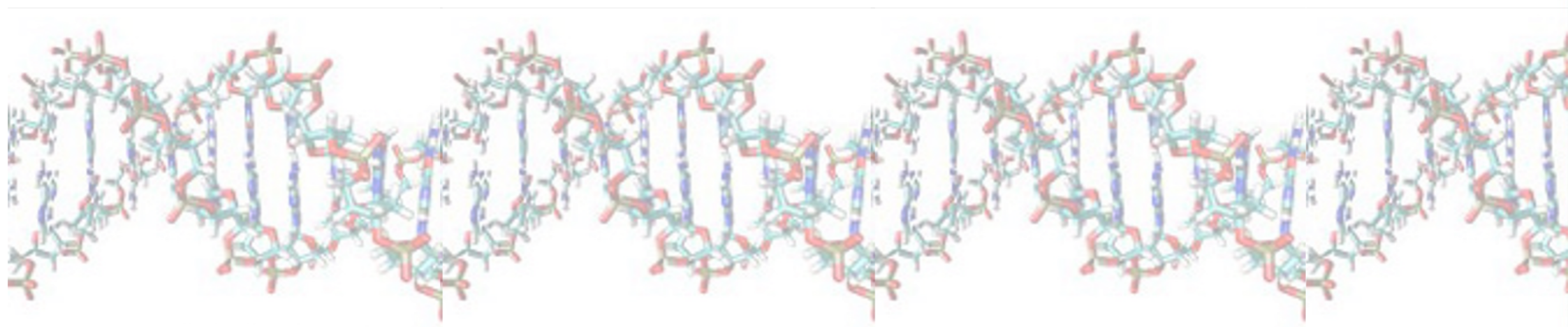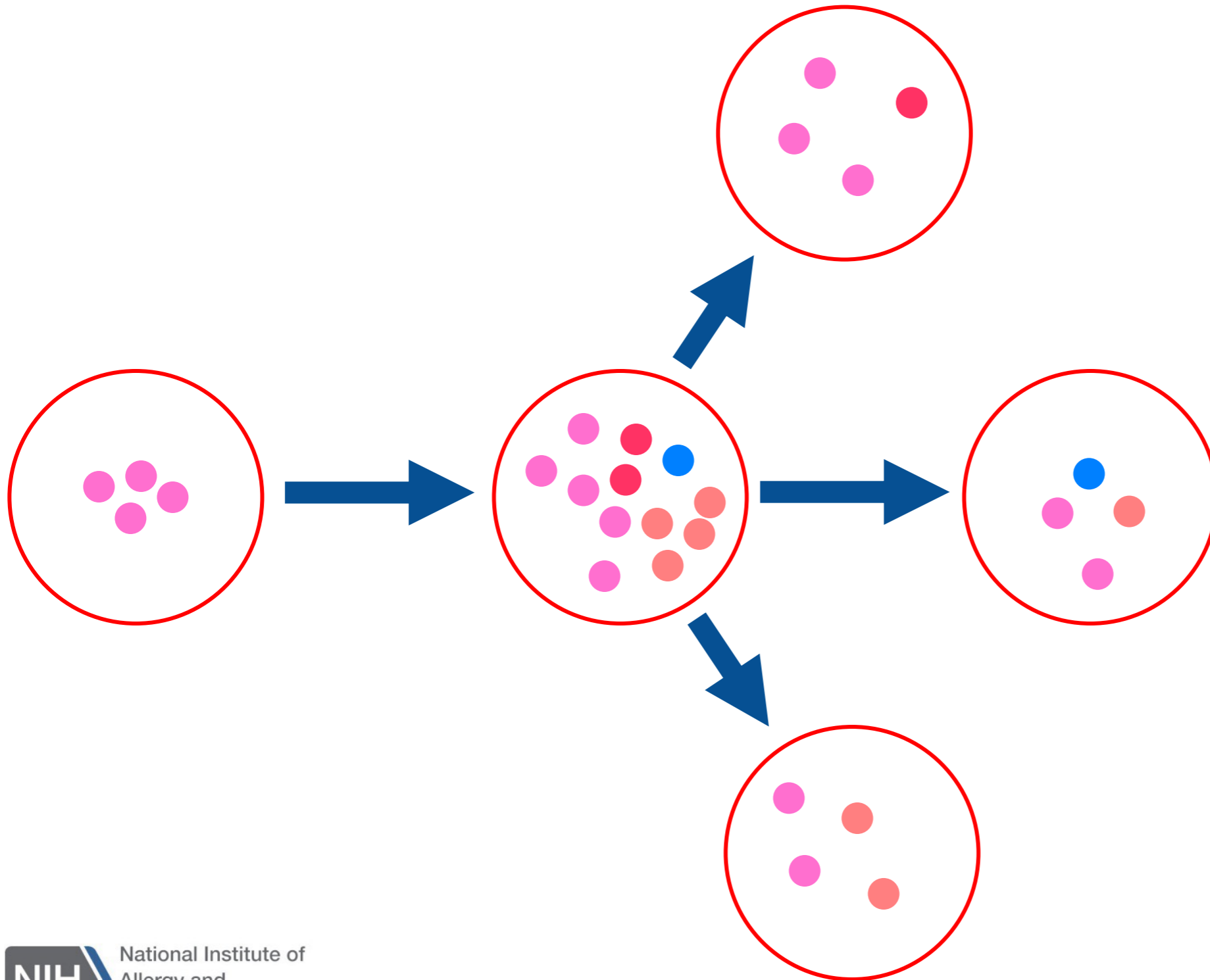
National Institute of Allergy and Infectious Diseases

# The Goal

Data → Phylogenetics Software → Tree

Data → Phylogenetics Software → Tree

National Institute of Allergy and Infectious Diseases

NIH

NIAID

3

# Biological sequences

# Why analyze biological sequences?

National Institute of Allergy and Infectious Diseases

# Biological sequences

- DNA contains the information basic to every process in a cell

- Proteins (and RNA) are the machines performing cellular processes

- Passed from one generation to the next

National Institute of Allergy and Infectious Diseases

# Sequence data are genealogical

# Comparative Methods

# Why analyze sequences using comparative methods?

# Comparative Methods

- Sequences related by common ancestry

- Analyzing samples with the trait against those without it

- The Grail: Finding nucleotide X at site Y in gene Z which correlates with the presence of the trait

- Correlation vs causation

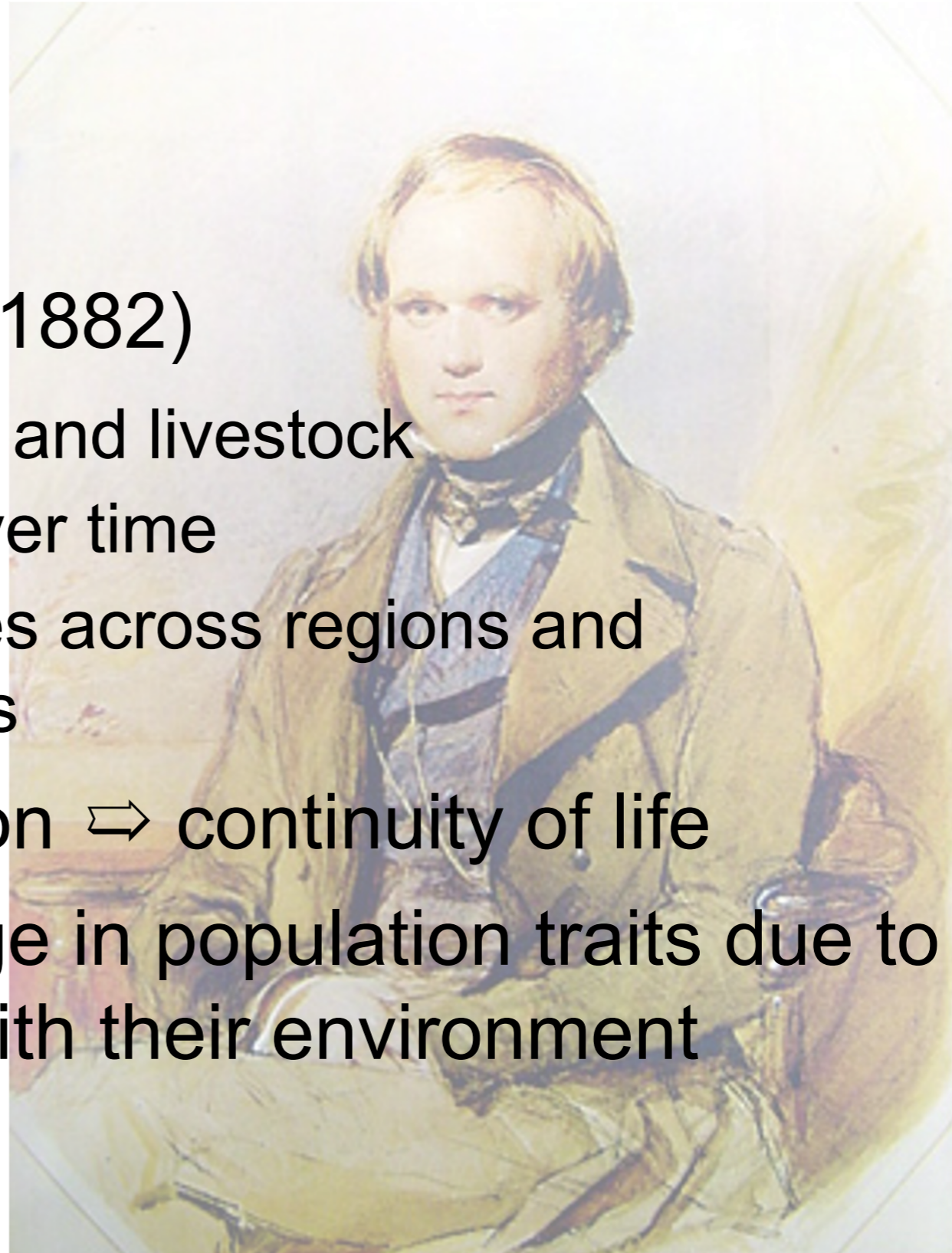National Institute of
Allergy and
Infectious Diseases

# Hierarchy of Life

- Carl Linnaeus (1707 - 1778)
  - Swedish physician/naturalist
  - Hierarchical organization of life
  - Binomial system of scientific names

National Institute of
Allergy and
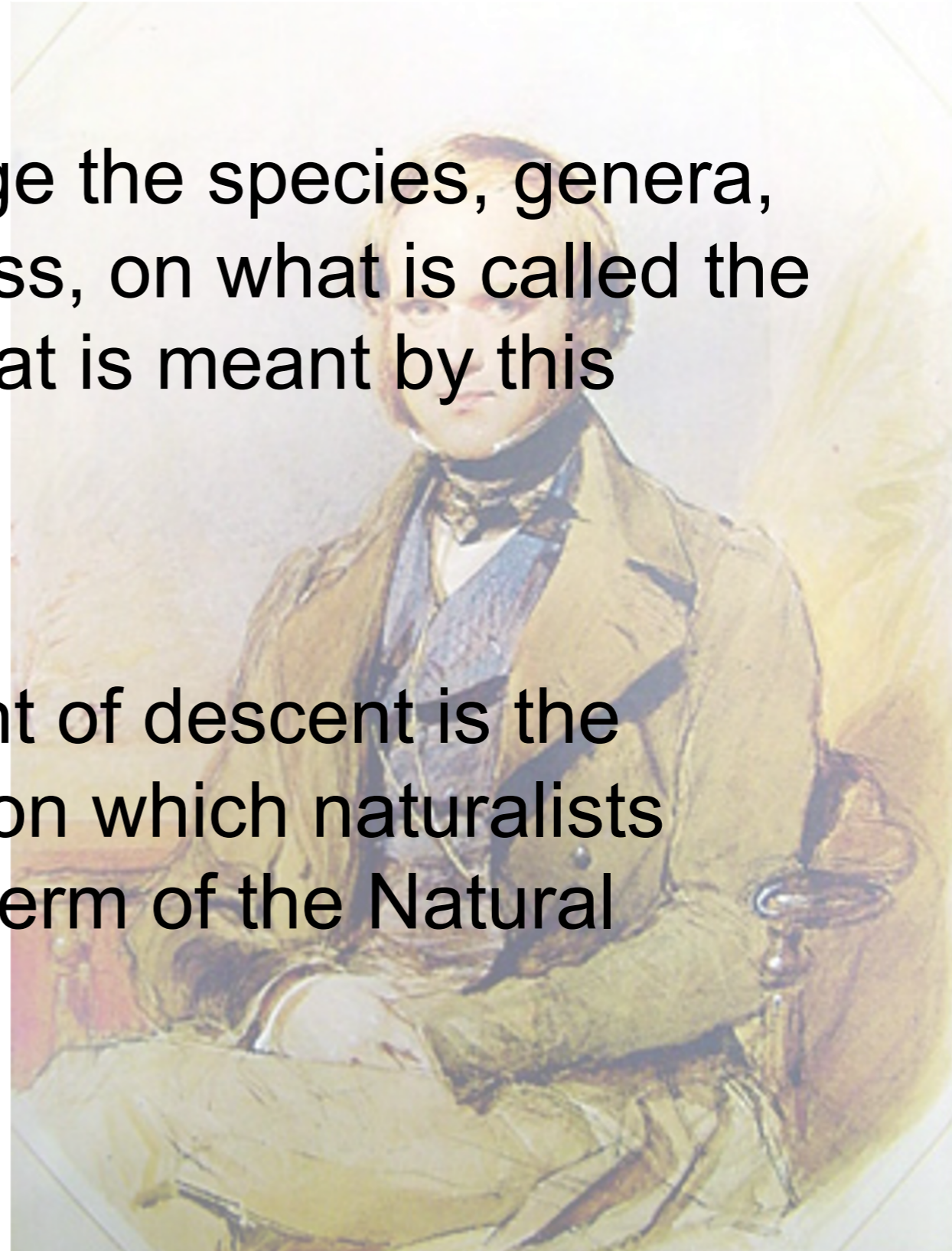Infectious Diseases

# Common Ancestry

- Charles Darwin (1809 - 1882)
  - Artificial selection: crops and livestock
  - Fossil record: change over time
  - Biogeography: similarities across regions and differences within locales
- Descent with modification ⇨ continuity of life
- Natural selection: change in population traits due to individual interactions with their environment

National Institute of Allergy and Infectious Diseases

NIAID

11

# Common Ancestry

"Naturalists try to arrange the species, genera, and families in each class, on what is called the Natural System. But what is meant by this system?" p.413

"… I believe this element of descent is the hidden bond of connexion which naturalists have sought under the term of the Natural System" p. 433

# What's so special about viruses?

- Short generation time
- Rapid evolution
- Genotypes - easy, phenotypes - hard
- Large populations
- Structured populations
- Rigorous temporal sampling of genotypes
- Shorter genomes → more WGS data

# PAIRWISE ALIGNMENT

- Sequence Alignment: Assigning homology to sites among a group of known sequences

- BLAST: Alignment of one sequence with many unknown

National Institute of
Allergy and
Infectious Diseases

# HOMOLOGY vs. ANALOGY

## common ancestry

## convergence

# PAIRWISE ALIGNMENT

- Sequence Alignment: Assigning homology to sites among a group of known sequences

  - Alignment of single loci
    - Clustal(W,X,Omega), MUSCLE, TCoffee, MAFFT

  - Alignment of overlapping contigs
    - Sequencher, Lasergene

  - Alignment of short reads
    - BWA, Bowtie, SOAP, MAQ

National Institute of
Allergy and
Infectious Diseases

# PAIRWISE ALIGNMENT

- Single locus

```
>GeneA_Human
ATGGGCCTTATATGCGTGATGCTGAAAG
>GeneA_Gorilla
ATGGGACTTATCTGCGTGATGCTGACAG
>GeneA_Macaque
ATGGGTCTCATATGTGTGATGCTTACAG
>GeneA_Mouse
ATGGCCCTGATATGCGTGATGCTGAACG
>GeneA_Sheep
ATGGCCCTAATATGC---AGGCTGAACG
```

# PAIRWISE ALIGNMENT

- Overlapping contigs

```
ATGGGCCTTATATGCGTGATGCTGAAAG
        TTATATGCGTGATGCTGAAAGGGCTTAG
          ATATGCGTGATGCTGAAAGGGCTTAGAAAT
            TGCGTGATGCTGAAAGGGCTTAGAAATT
                ATGCTGAAAGGGCTTAGAAATTCGG
                    AAAGGGCTTAGAAATTGCGGCTAGGCCTCC
                            CGGCTAGGCCTCCGAACGC
TACCCGGAATATACGCACTA
            CACTACGACTTTCCCGAATCTTTAAGCC
                CTTTCCCGAATCTTTAAGCCGATCCGGA
```

# PAIRWISE ALIGNMENT

- Short reads

# PAIRWISE ALIGNMENT

```
HBA_HUMAN     GSAQVKGHGKKVADALTNAVAHVDDMPNALSALSDLHAHKL
              G+ +VK+HGKKV   A+++++AH+D++ +++++LS+LH   KL
HBB_HUMAN     GNPKVKAHGKKVLGAFSDGLAHLDNLKGTFATLSELHCDKL


HBA_HUMAN     GSAQVKGHGKKVADALTNAVAHV---D--DMPNALSALSDLHAHKL
              ++ ++++H+ KV    + +A  ++              +L+ L+++H+ K
LGB2_LUPLU    NNPELQAHAGKVFKLVYEAAIQLQVTGVVVTDATLKNLGSVHVSKG


HBA_HUMAN     GSAQVKGHGKKVADALTNAVAHVDDMPNALSALSD----LHAHKL
              GS+ + G +    +D L  ++ H+ D+  A +AL D     ++AH+
F11G11.2      GSGYLVGDSLTFVDLL--VAQHTADLLAANAALLDEFPQFKAHQE
```

# PAIRWISE ALIGNMENT

Jukes-Cantor Substitution Probabilities

$\mu t = 0.25$

|   | A | C | G | T |
|---|---|---|---|---|
| A | 0.5259 | 0.1580 | 0.1580 | 0.1580 |
| C | 0.1580 | 0.5259 | 0.1580 | 0.1580 |
| G | 0.1580 | 0.1580 | 0.5259 | 0.1580 |
| T | 0.1580 | 0.1580 | 0.1580 | 0.5259 |

# PAIRWISE ALIGNMENT

Protein Score Matrices
Similarity of Amino Acids



From Esquivel RO, et al.. 2013. Advances in Quantum Mechanics, Chapter 27 InTech.

# PAIRWISE ALIGNMENT

Protein Score Matrices

- Derived from empirical data
- Account for depth of relationship among the data
- Expressed as log-odds ratio:
  - Logarithm of the ratio of the probabilities of two residues being aligned due to homology versus random chance

National Institute of Allergy and Infectious Diseases

# PAIRWISE ALIGNMENT

Protein Substitution Matrices

- PAM250: Based on phylogenies where all sequences differ by no more than 15%.

- BLOSUM62: Based on clusters of sequences with greater than 62% identical residues.

National Institute of Allergy and Infectious Diseases

# Protein Substitution Matrices

BLOSUM62

| | C | S | T | P | A | G | N | D | E | Q | H | R | K | M | I | L | V | F | Y | W |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| C | 9 | | | | | | | | | | | | | | | | | | | |
| S | −1 | 4 | | | | | | | | | | | | | | | | | | |
| T | −1 | 1 | 5 | | | | | | | | | | | | | | | | | |
| P | −3 | −1 | −1 | 7 | | | | | | | | | | | | | | | | |
| A | 0 | 1 | 0 | −1 | 4 | | | | | | | | | | | | | | | |
| G | −3 | 0 | −2 | −2 | 0 | 6 | | | | | | | | | | | | | | |
| N | −3 | 1 | 0 | −2 | −2 | 0 | 6 | | | | | | | | | | | | | |
| D | −3 | 0 | −1 | −1 | −2 | −1 | 1 | 6 | | | | | | | | | | | | |
| E | −4 | 0 | −1 | −1 | −1 | −2 | 0 | 2 | 5 | | | | | | | | | | | |
| Q | −3 | 0 | −1 | −1 | −1 | −2 | 0 | 0 | 2 | 5 | | | | | | | | | | |
| H | −3 | −1 | −2 | −2 | −2 | −2 | 1 | −1 | 0 | 0 | 8 | | | | | | | | | |
| R | −3 | −1 | −1 | −2 | −1 | −2 | 0 | −2 | 0 | 1 | 0 | 5 | | | | | | | | |
| K | −3 | 0 | −1 | −1 | −1 | −2 | 0 | −1 | 1 | 1 | −1 | 2 | 5 | | | | | | | |
| M | −1 | −2 | −1 | −2 | −1 | −3 | −2 | −3 | −2 | 0 | −2 | −1 | −1 | 5 | | | | | | |
| I | −1 | −2 | −1 | −3 | −1 | −4 | −3 | −3 | −3 | −3 | −3 | −3 | −3 | 1 | 4 | | | | | |
| L | −1 | −2 | −1 | −3 | −1 | −4 | −3 | −4 | −3 | −2 | −3 | −2 | −2 | 2 | 2 | 4 | | | | |
| V | −1 | −2 | 0 | −2 | 0 | −3 | −3 | −3 | −2 | −2 | −3 | −3 | −2 | 1 | 3 | 1 | 4 | | | |
| F | −2 | −2 | −2 | −4 | −2 | −3 | −3 | −3 | −3 | −3 | −1 | −3 | −3 | 0 | 0 | 0 | −1 | 6 | | |
| Y | −2 | −2 | −2 | −3 | −2 | −3 | −2 | −3 | −2 | −1 | 2 | −2 | −2 | −1 | −1 | −1 | −1 | 3 | 7 | |
| W | −2 | −3 | −2 | −4 | −3 | −2 | −4 | −4 | −3 | −2 | −2 | −3 | −3 | −1 | −3 | −2 | −3 | 1 | 2 | 11 |

# Multiple Sequence Alignment

- Global alignment (Needleman-Wunsch)
  - Assign homology across the entire sequence
  - Clustal

- Local alignment (Smith-Waterman)
  - Assign homology for subsequences
  - MUSCLE and BLAST
  - Good for aligning very divergent sequences
- **Inspect and edit your alignment!**

National Institute of
Allergy and
Infectious Diseases

# Multiple Sequence Alignment

## The Progressive Alignment Algorithm



Unaligned Sequences

Distance Matrix

UPGMA tree

Aligned Sequences MSA1

Kimura Distance Matrix

Subtree Alignment Profiles

Delete edge from UPGMA tree 2

Aligned Sequences MSA2

UPGMA tree 2

Better!

Not Better

Align Subtree Profiles MSAsp

# Multiple Sequence Alignment

## MSA with MEGA7

# Multiple Sequence Alignment

## MSA with MEGA7

# Multiple Sequence Alignment

## NEVER

directly input the output of a MSA program into an analysis program!

## ALWAYS

inspect the alignment to improve it.

National Institute of Allergy and Infectious Diseases

# Multiple Sequence Alignment

# Multiple Sequence Alignment

# Multiple Sequence Alignment

## Programs

- Clustal
  - Your own computer
  - Web Server
  - NIAID HPC cluster
- MUSCLE
  - Your own computer
  - Web Server
  - NIAID HPC cluster
- MAFFT
  - Web Server

# Multiple Sequence Alignment

## Multiple Sequence Alignment Editors

- Geneious

- MacVector

- MegAlign (Lasergene)

- AliView

- GeneDoc

- BioEdit

# Web Resources

**ClustalW**

http://www.clustal.org/

**Muscle**

http://www.drive5.com/muscle/download3.6.html

**MAFFT**

http://mafft.cbrc.jp/alignment/server/

**MEGA7**

https://www.megasoftware.net/

**AliView**

http://www.ormbunkar.se/aliview/

**GeneDoc**

http://genedoc.software.informer.com/

**BioEdit**

http://www.mbio.ncsu.edu/BioEdit/bioedit.html

# What's next?

After the break


Building trees with our MSA

# 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

National Institute of Allergy and Infectious Diseases

# What is a phylogenetic tree?

# What is a phylogenetic tree?

Unrooted

Rooted

# 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

# 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

National Institute of Allergy and Infectious Diseases

# SUBSTITUTION MODEL

Jukes-Cantor Substitution Probabilities

$\mu t = 0.25$

|   | A | C | G | T |
|---|---|---|---|---|
| A | 0.5259 | 0.1580 | 0.1580 | 0.1580 |
| C | 0.1580 | 0.5259 | 0.1580 | 0.1580 |
| G | 0.1580 | 0.1580 | 0.5259 | 0.1580 |
| T | 0.1580 | 0.1580 | 0.1580 | 0.5259 |

# SUBSTITUTION MODEL

BLOSUM62

```
C    9
S   -1    4
T   -1    1    5
P   -3   -1   -1    7
A    0    1    0   -1    4
G   -3    0   -2   -2    0    6
N   -3    1    0   -2   -2    0    6
D   -3    0   -1   -1   -2   -1    1    6
E   -4    0   -1   -1   -1   -2    0    2    5
Q   -3    0   -1   -1   -1   -2    0    0    2    5
H   -3   -1   -2   -2   -2   -2    1   -1    0    0    8
R   -3   -1   -1   -2   -1   -2    0   -2    0    1    0    5
K   -3    0   -1   -1   -1   -2    0   -1    1    1   -1    2    5
M   -1   -2   -1   -2   -1   -3   -2   -3   -2    0   -2   -1   -1    5
I   -1   -2   -1   -3   -1   -4   -3   -3   -3   -3   -3   -3   -3    1    4
L   -1   -2   -1   -3   -1   -4   -3   -4   -3   -2   -3   -2   -2    2    2    4
V   -1   -2    0   -2    0   -3   -3   -3   -2   -2   -3   -3   -2    1    3    1    4
F   -2   -2   -2   -4   -2   -3   -3   -3   -3   -3   -1   -3   -3    0    0    0   -1    6
Y   -2   -2   -2   -3   -2   -3   -2   -3   -2   -1    2   -2   -2   -1   -1   -1   -1    3    7
W   -2   -3   -2   -4   -3   -2   -4   -4   -3   -2   -2   -3   -3   -1   -3   -2   -3    1    2   11
     C    S    T    P    A    G    N    D    E    Q    H    R    K    M    I    L    V    F    Y    W
```

# The algorithm approach: Distance Methods

## Tree-Building Algorithms

- UPGMA

- Neighbor-Joining

# The algorithm approach: Neighbor-joining Calculation

|   | A | B | C | D | E | R |
|---|---|---|---|---|---|---|
| A | – | 0.1715 | 0.2147 | 0.3091 | 0.2326 | 0.9279 |
| B | –0.4766 | – | 0.2991 | 0.3399 | 0.2058 | 1.0163 |
| C | –0.4905 | –0.4356 | – | 0.2795 | 0.3943 | 1.1876 |
| D | –0.4527 | –0.4514 | –0.5689 | – | 0.4289 | 1.3574 |
| E | –0.4972 | –0.5535 | –0.4221 | –0.4441 | – | 1.2616 |

C to Node 1 distance = 0.2795/2 + (1.1876 – 1.3574)/6 = 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

# The algorithm approach: Neighbor-joining Calculation

|        | A       | B       | E       | Node 1  | R      |
|--------|---------|---------|---------|---------|--------|
| A      | –       | 0.1715  | 0.2326  | 0.1222  | 0.5263 |
| B      | –0.3701 | –       | 0.2058  | 0.1798  | 0.5571 |
| E      | –0.3856 | –0.4278 | –       | 0.2719  | 0.7103 |
| Node 1 | –0.4278 | –0.3856 | –0.3701 | –       | 0.5739 |

A to Node 2 distance = 0.1222/2 + (0.5263 – 0.5739)/4 = 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

National Institute of
Allergy and
Infectious Diseases

Hillis, Moritz, and Mable 1996, p. 489

# The algorithm approach: Neighbor-joining Calculation

|        | B       | E       | Node 2  | R      |
|--------|---------|---------|---------|--------|
| B      | –       | 0.2058  | 0.1146  | 0.3204 |
| E      | –0.5116 | –       | 0.1912  | 0.3970 |
| Node 2 | –0.5116 | –0.5116 | –       | 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

# MEGA7

Phylogeny Construction: Neighbor-Joining

# MEGA7

Phylogeny Construction: Neighbor-Joining

# MEGA7

## Phylogeny Construction: Neighbor-Joining

# MEGA7

Phylogeny Construction: Neighbor-Joining

# MEGA7

Phylogeny Construction: Neighbor-Joining

# 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

National Institute of
Allergy and
Infectious Diseases

# 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

# 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)

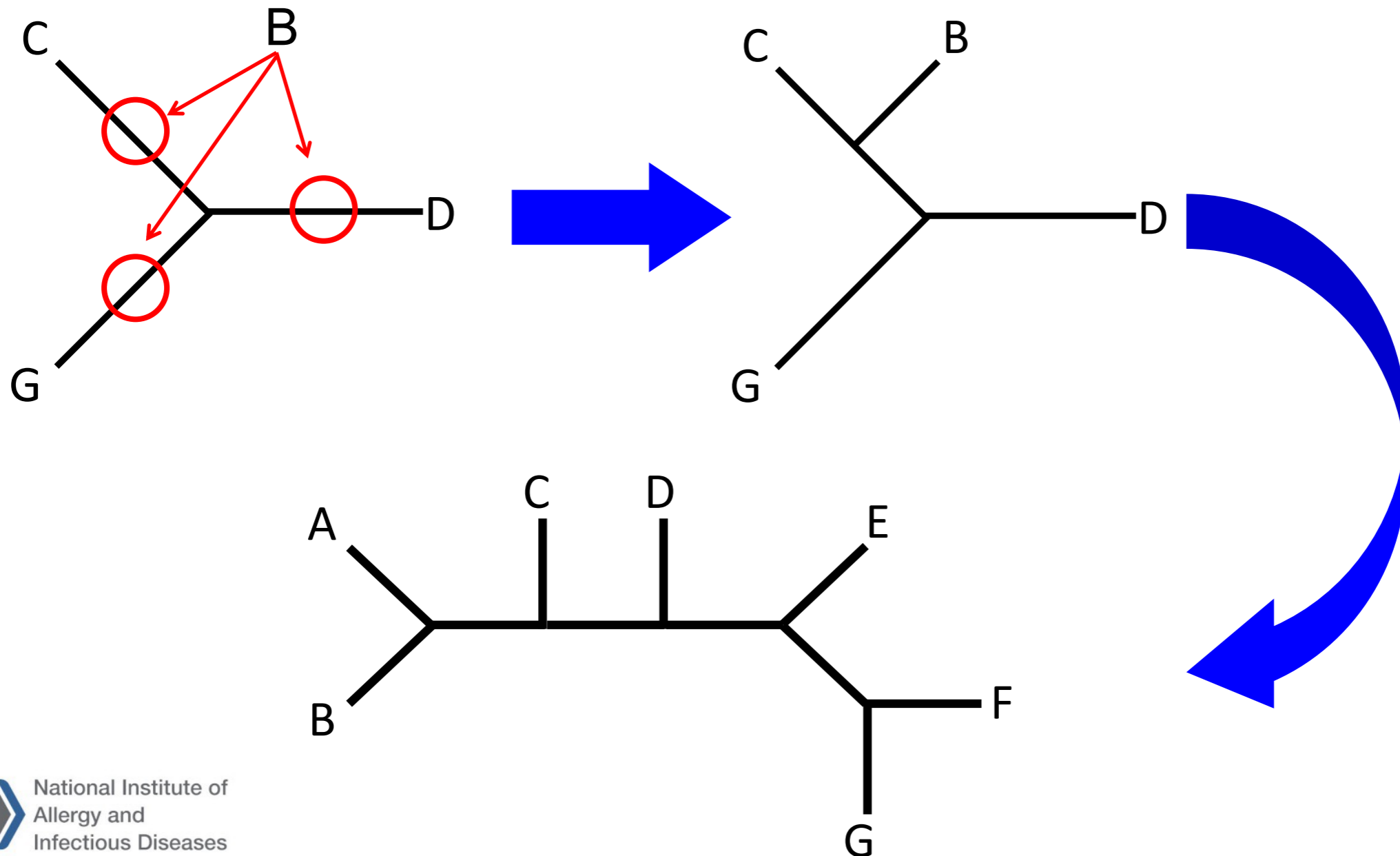- Is this the best tree? (Stopping criteria)

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

National Institute of
Allergy and
Infectious Diseases

NIH

# 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

National Institute of Allergy and Infectious Diseases

# 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.
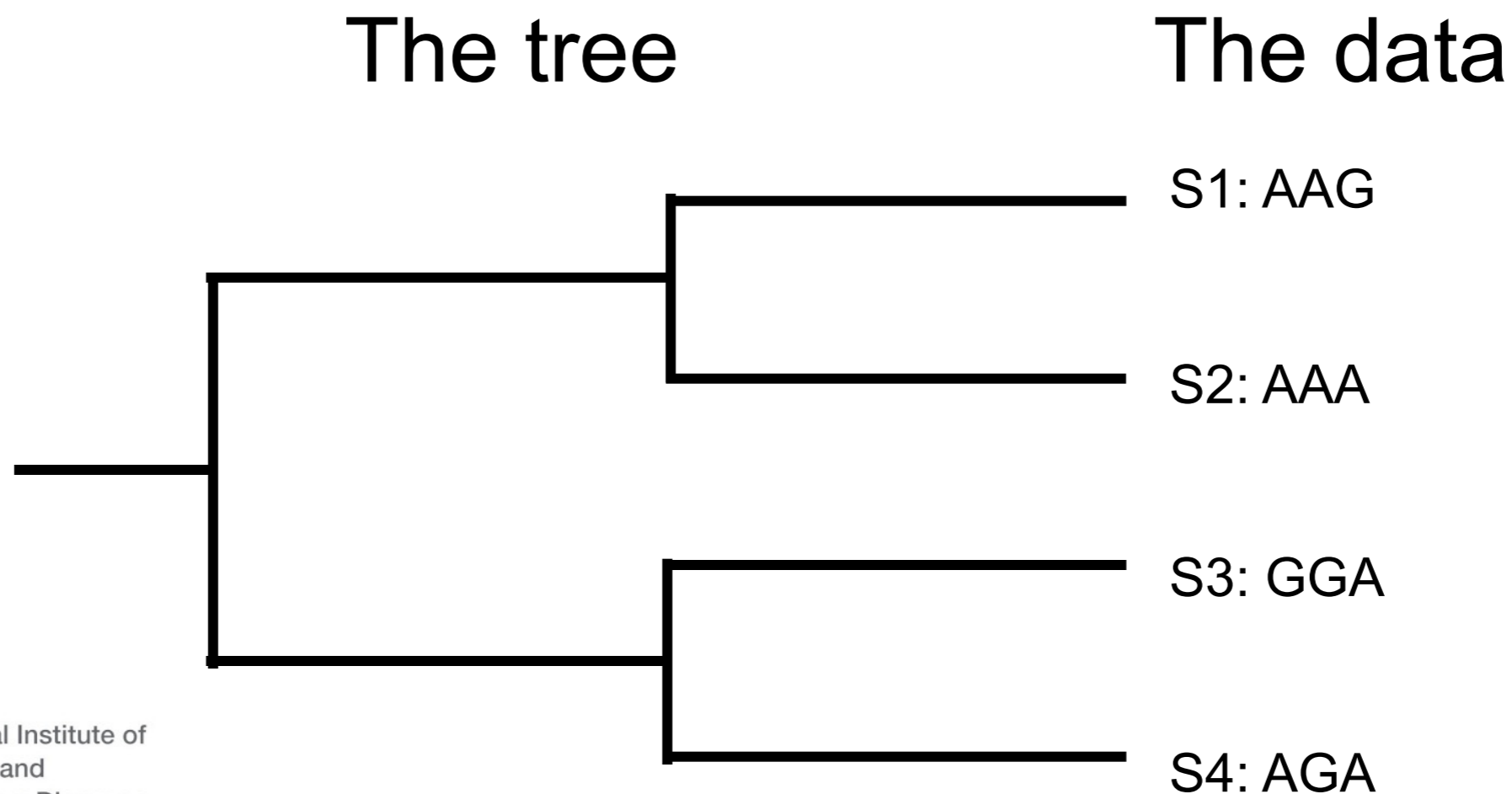
# The optimality criterion approach

## Exploring tree space

# The optimality criterion approach

## Is this tree optimal?
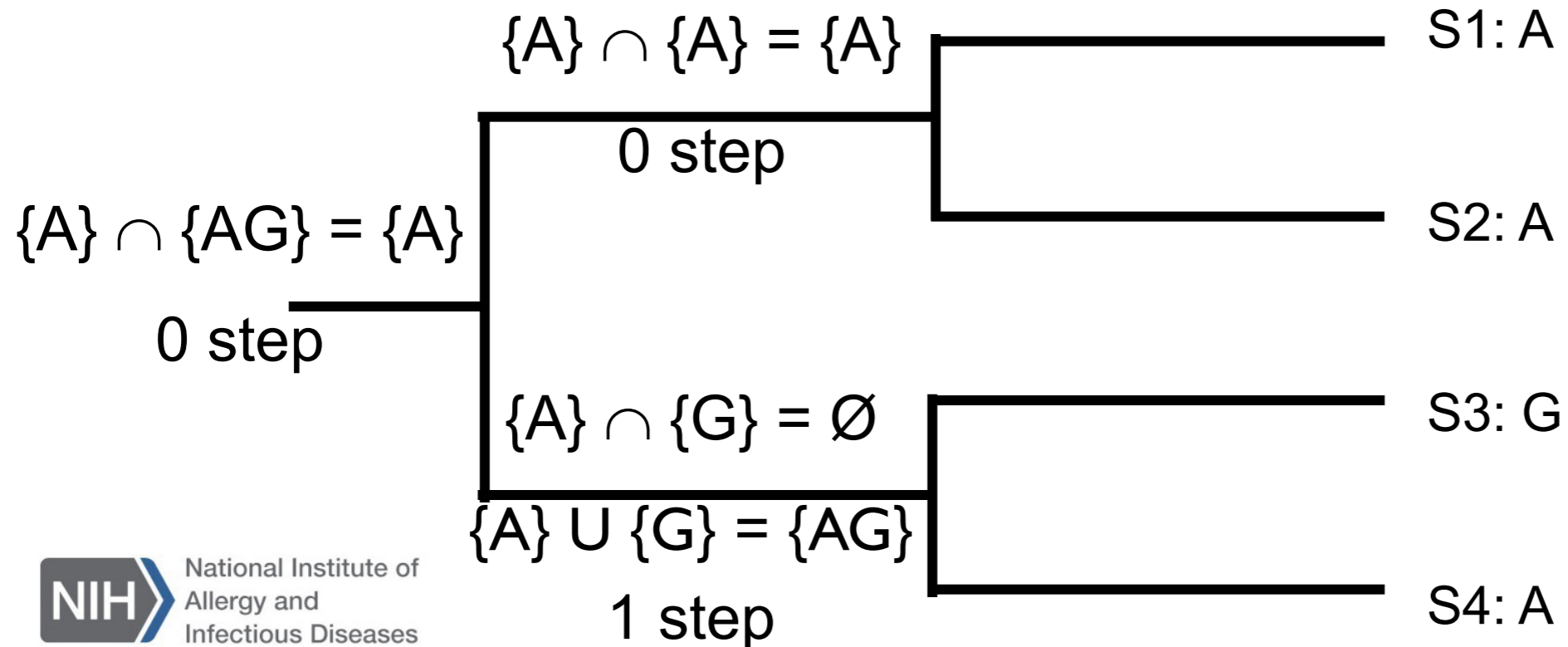
Counting changes (Fitch parsimony)



The tree

The data

S1: AAG

S2: AAA

S3: GGA

S4: AGA

# The optimality criterion approach

## Is this tree optimal?

Counting changes (Fitch parsimony)
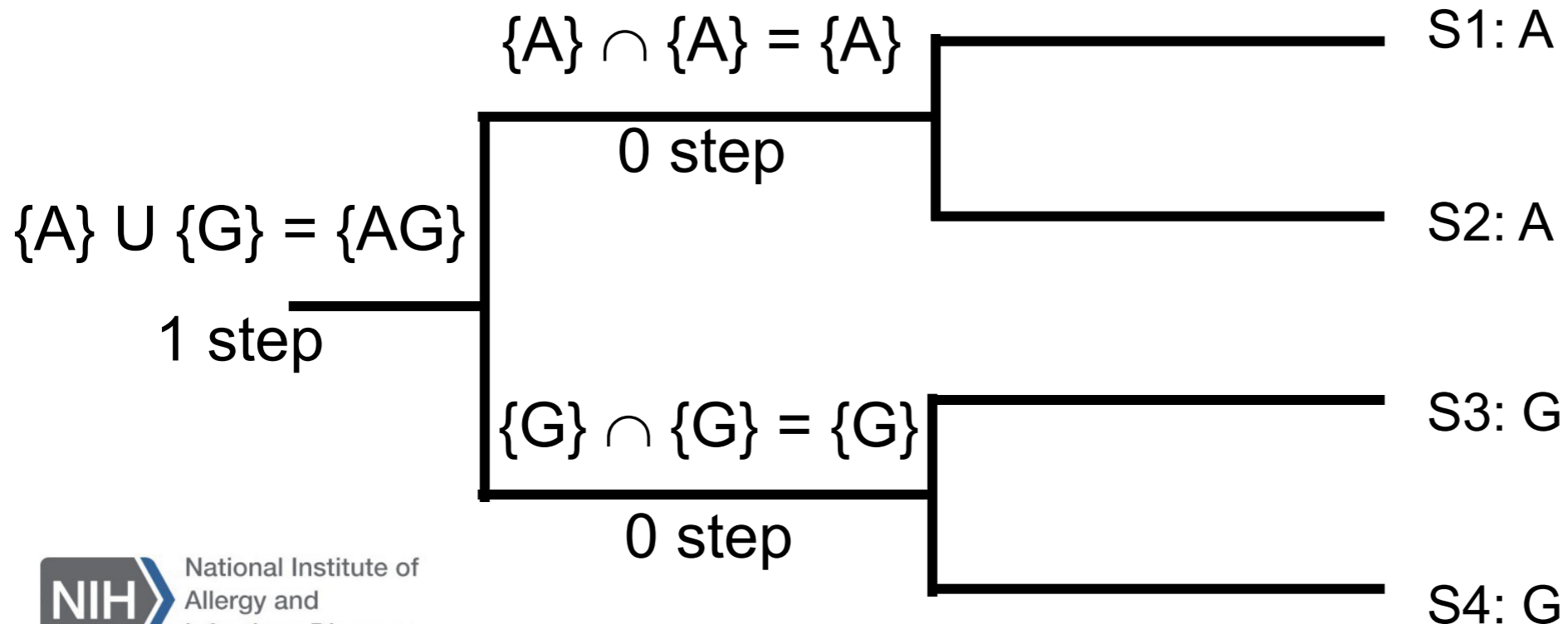
Position 1     The tree     The data



{A} ∩ {A} = {A}

0 step

{A} ∩ {AG} = {A}

0 step

{A} ∩ {G} = Ø

{A} U {G} = {AG}

1 step

S1: A

S2: A

S3: G

S4: A

# The optimality criterion approach

## Is this tree optimal?

Counting changes (Fitch parsimony)

Position 2　　　The tree　　　The data

{A} ∩ {A} = {A}

S1: A
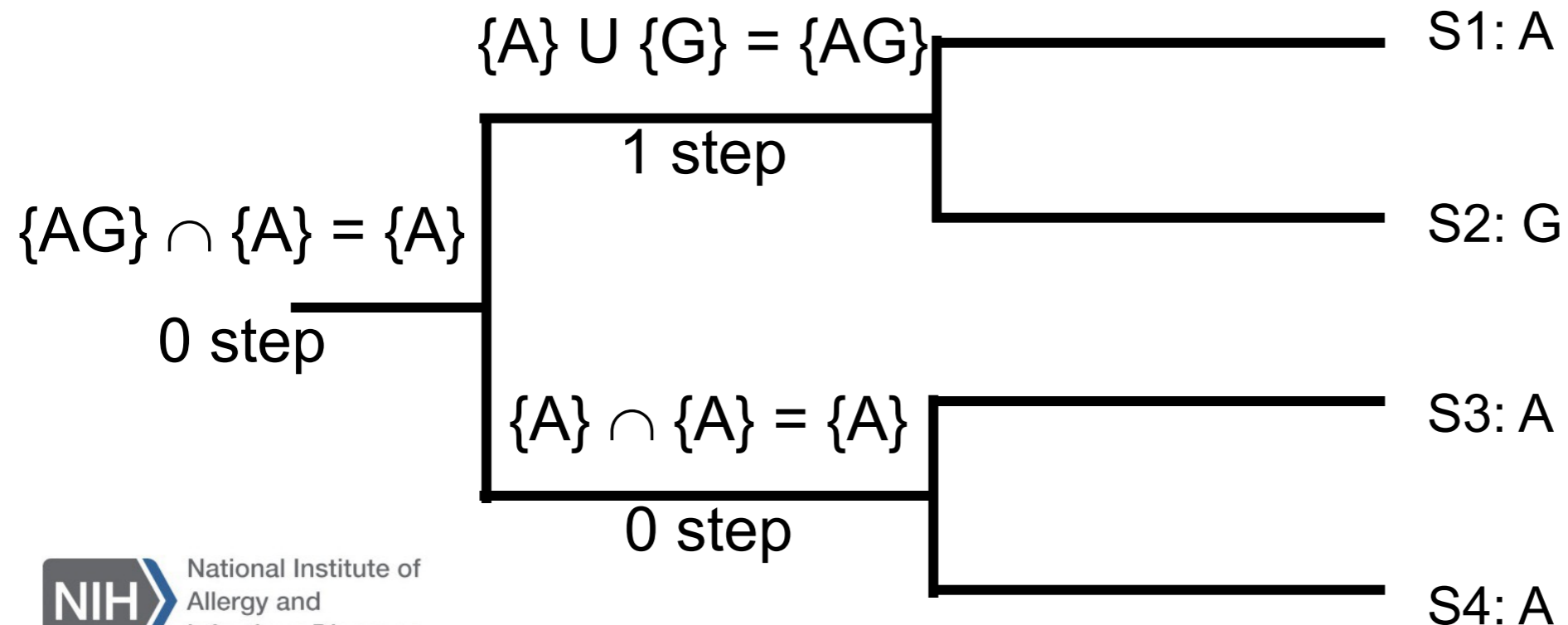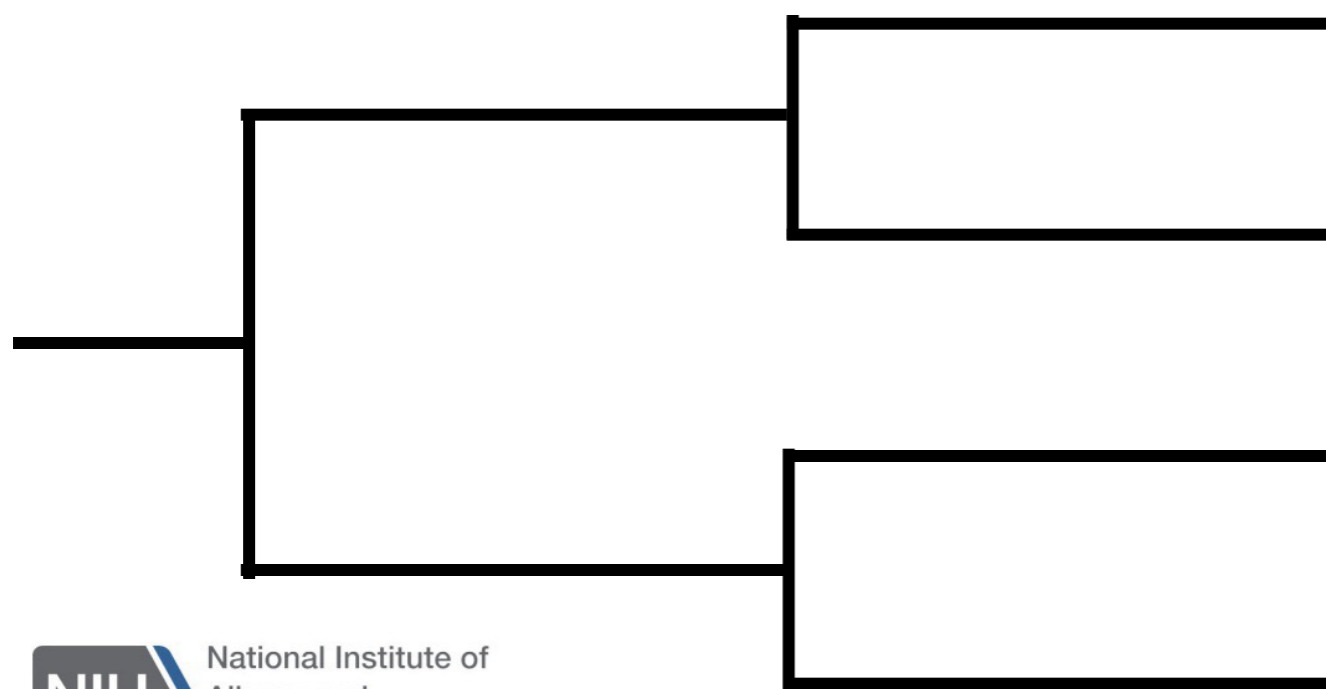
0 step

S2: A

{A} U {G} = {AG}

1 step

{G} ∩ {G} = {G}

S3: G

0 step

S4: G

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

{AG} ∩ {A} = {A}                                          S2: G

0 step

{A} ∩ {A} = {A}                                           S3: A

0 step

                                                          S4: A

# The optimality criterion approach

## Is this tree optimal?

Counting changes (Fitch parsimony)

The tree                    The data
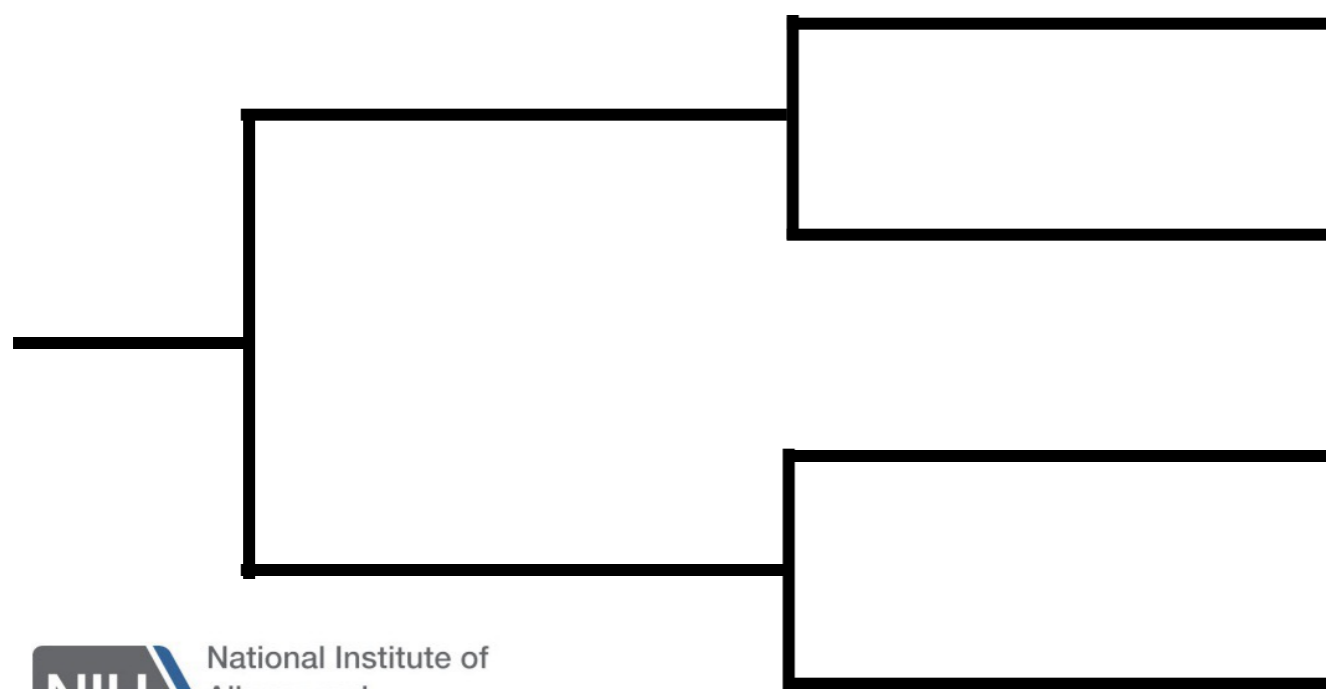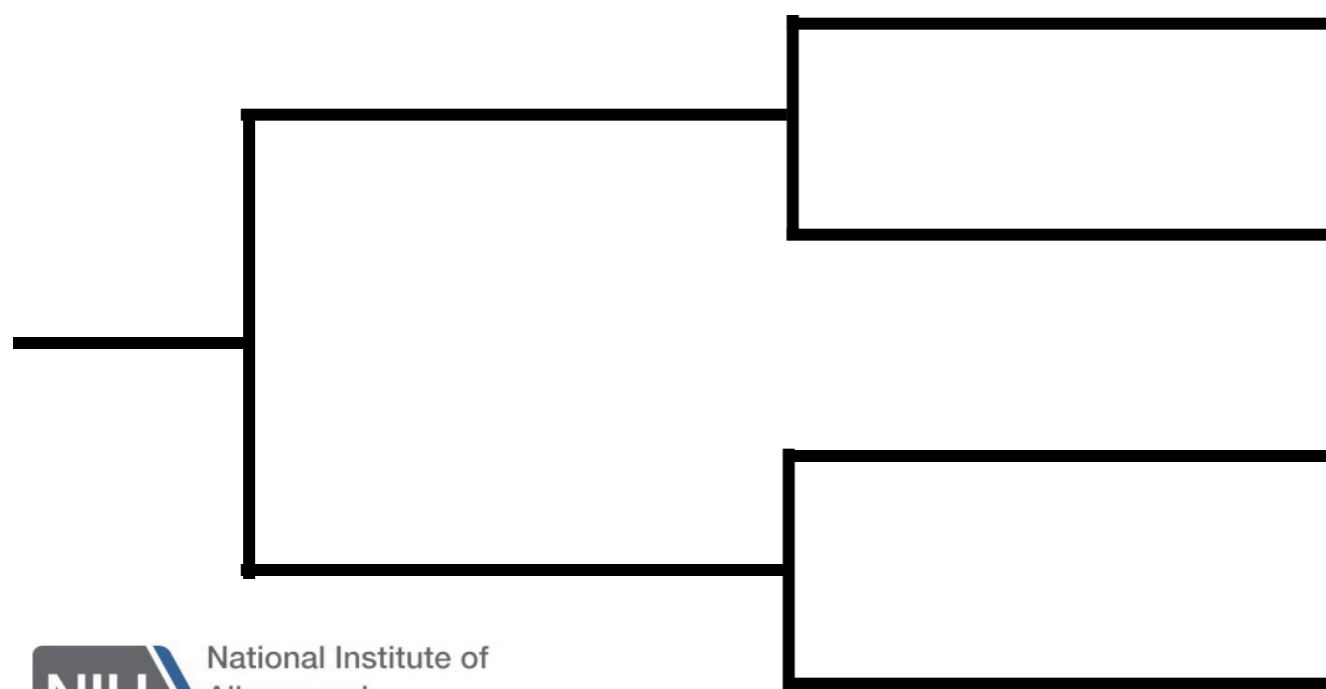


S1: AAG

S2: AAA

S3: GGA

S4: AGA
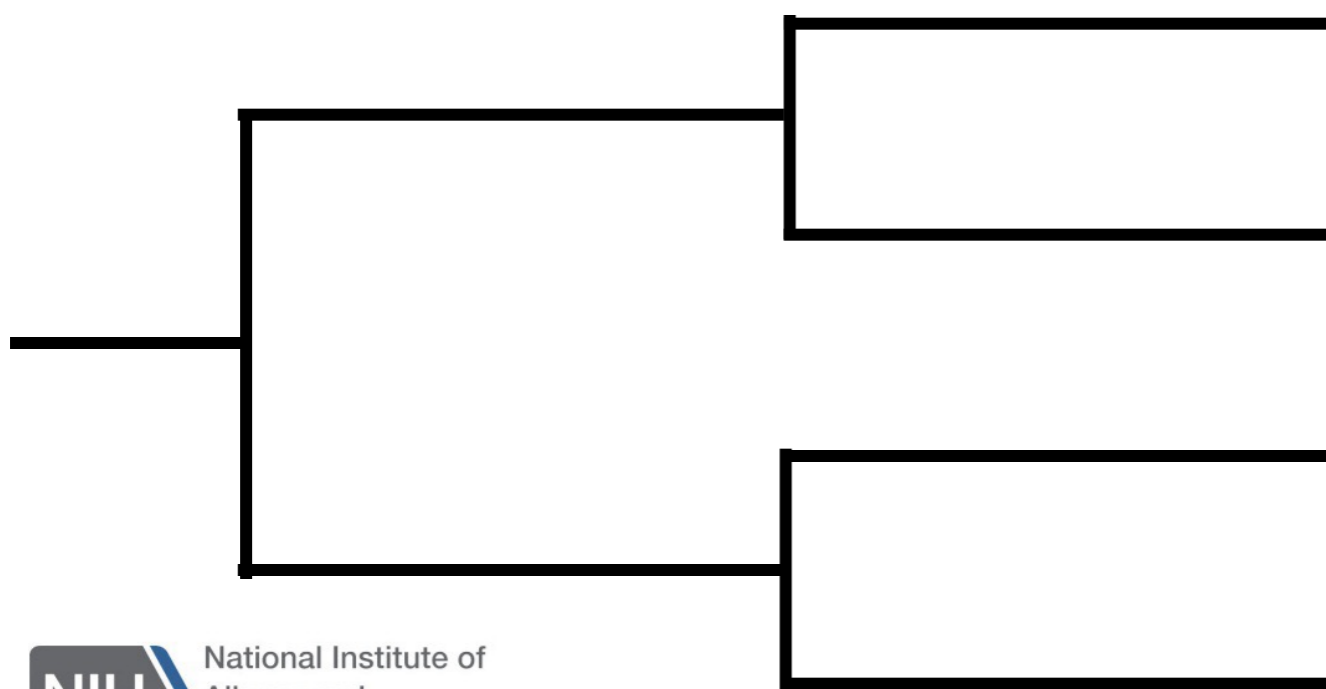
Total steps: 3

# 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

# The optimality criterion approach

## Is this tree optimal?

Counting changes (Fitch parsimony)

The tree                    The data



S1: AAG

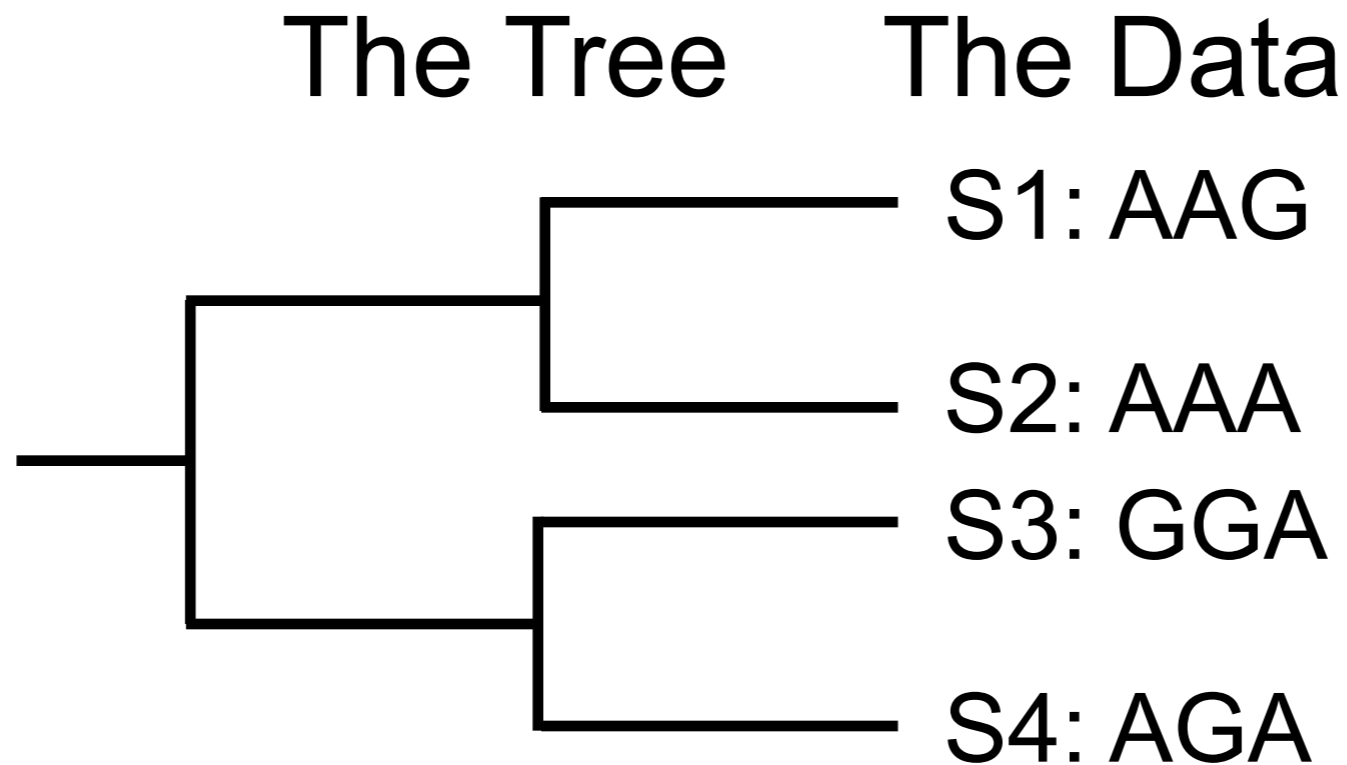S4: AGA                     Total
                            steps:
                            4
S2: AAA

S3: GGA

# Optimality Criterion: Likelihood

Calculating likelihood

The Tree       The Data

S1: AAG

S2: AAA

S3: GGA

S4: AGA

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

National Institute of
Allergy and
Infectious Diseases

# Optimality Criterion: Likelihood

Calculating likelihood: Setting parameters

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

What values do you use for the substitution model?

Run jModelTest (or ProtTest for protein MSAs)

# 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

# Optimality Criterion: Likelihood

## Calculating likelihood: MEGA7 Options

# Optimality Criterion: Likelihood

## Calculating likelihood: MEGA7 Results

Table. Maximum Likelihood fits of 24 different nucleotide substitution models

| Model | Parameters | BIC | AICc | lnL | (+I) | (+G) | R | f(A) | f(T) | f(C) | f(G) | r(AT) | r(AC) | r(AG) | r(TA) | r(TC) | r(TG) | r(CA) | r(CT) | r(CG) | r(GA) | r(GT) | r(GC) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| HKY+G | 36 | 12438.206 | 12167.533 | -6047.669 | n/a | 1.09 | 1.68 | 0.320 | 0.241 | 0.232 | 0.208 | 0.044 | 0.043 | 0.131 | 0.059 | 0.146 | 0.038 | 0.059 | 0.152 | 0.038 | 0.202 | 0.044 | 0.043 |
| TN93+G | 37 | 12445.207 | 12167.020 | -6046.407 | n/a | 1.09 | 1.68 | 0.320 | 0.241 | 0.232 | 0.208 | 0.044 | 0.043 | 0.121 | 0.059 | 0.160 | 0.038 | 0.059 | 0.166 | 0.038 | 0.186 | 0.044 | 0.043 |
| HKY+G+I | 37 | 12447.534 | 12169.347 | -6047.571 | 0.04 | 1.24 | 1.68 | 0.320 | 0.241 | 0.232 | 0.208 | 0.044 | 0.043 | 0.131 | 0.059 | 0.147 | 0.038 | 0.059 | 0.152 | 0.038 | 0.202 | 0.044 | 0.043 |
| GTR+G | 40 | 12450.726 | 12150.001 | -6034.880 | n/a | 1.07 | 1.67 | 0.320 | 0.241 | 0.232 | 0.208 | 0.034 | 0.058 | 0.120 | 0.045 | 0.160 | 0.025 | 0.080 | 0.166 | 0.046 | 0.185 | 0.029 | 0.052 |
| T92+G | 34 | 12453.816 | 12198.170 | -6064.998 | n/a | 1.09 | 1.68 | 0.280 | 0.280 | 0.220 | 0.220 | 0.052 | 0.041 | 0.138 | 0.052 | 0.138 | 0.041 | 0.052 | 0.176 | 0.041 | 0.176 | 0.052 | 0.041 |
| TN93+G+I | 38 | 12454.467 | 12168.767 | -6046.275 | 0.04 | 1.26 | 1.68 | 0.320 | 0.241 | 0.232 | 0.208 | 0.044 | 0.043 | 0.121 | 0.059 | 0.160 | 0.038 | 0.059 | 0.166 | 0.038 | 0.186 | 0.044 | 0.043 |
| K2+G | 33 | 12456.913 | 12208.781 | -6071.309 | n/a | 1.10 | 1.66 | 0.250 | 0.250 | 0.250 | 0.250 | 0.047 | 0.047 | 0.156 | 0.047 | 0.156 | 0.047 | 0.047 | 0.156 | 0.047 | 0.156 | 0.047 | 0.047 |
| GTR+G+I | 41 | 12460.097 | 12151.860 | -6034.804 | 0.03 | 1.19 | 1.67 | 0.320 | 0.241 | 0.232 | 0.208 | 0.034 | 0.058 | 0.120 | 0.045 | 0.161 | 0.025 | 0.080 | 0.167 | 0.046 | 0.185 | 0.029 | 0.052 |
| T92+G+I | 35 | 12463.297 | 12200.137 | -6064.976 | 0.02 | 1.16 | 1.68 | 0.280 | 0.280 | 0.220 | 0.220 | 0.052 | 0.041 | 0.138 | 0.052 | 0.138 | 0.041 | 0.052 | 0.177 | 0.041 | 0.177 | 0.052 | 0.041 |
| K2+G+I | 34 | 12466.397 | 12210.751 | -6071.288 | 0.02 | 1.17 | 1.66 | 0.250 | 0.250 | 0.250 | 0.250 | 0.047 | 0.047 | 0.156 | 0.047 | 0.156 | 0.047 | 0.047 | 0.156 | 0.047 | 0.156 | 0.047 | 0.047 |
| HKY+I | 36 | 12502.024 | 12231.351 | -6079.578 | 0.21 | n/a | 1.77 | 0.320 | 0.241 | 0.232 | 0.208 | 0.043 | 0.041 | 0.133 | 0.057 | 0.149 | 0.037 | 0.057 | 0.155 | 0.037 | 0.206 | 0.043 | 0.041 |
| TN93+I | 37 | 12507.897 | 12229.710 | -6077.752 | 0.21 | n/a | 1.77 | 0.320 | 0.241 | 0.232 | 0.208 | 0.043 | 0.041 | 0.122 | 0.057 | 0.164 | 0.037 | 0.057 | 0.170 | 0.037 | 0.188 | 0.043 | 0.041 |
| GTR+I | 40 | 12516.143 | 12215.418 | -6067.589 | 0.21 | n/a | 1.76 | 0.320 | 0.241 | 0.232 | 0.208 | 0.034 | 0.055 | 0.122 | 0.045 | 0.164 | 0.025 | 0.076 | 0.170 | 0.044 | 0.188 | 0.029 | 0.050 |

Substitution models

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

# The Gamma Distribution



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

# How reliable are my trees?

Bootstrapping (nonparametric)

# Calculating likelihood: Programs

PAUP* – Commercial, NIH Biowulf, or NIAID HPC

     DNA only

PHYLIP – Download, NIH Biowulf, or NIAID HPC

     dnaml and proml programs

MEGA – Download for free from www.megasoftware.net

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

Generally the user has more flexibility with a local program.

     But local programs can hog your computer.

National Institute of
Allergy and
Infectious Diseases

# Input File Formats

Phylogenetics program input file formats

FASTA

```
>MC1_01B4fs
TGCACTAATAATCTGATT---------AATATCACTGAGAATACTAATAATACCATTACT
>MC1_01A10
TGCACT---AATCTGACAAAGGCTATTAAGACCAATGGGAATGCTAATAATACCAGTACT
>MC1_01C1
TGCACTAATAATCTGATT--------AATATCACTGAGAATACTAATAATACCATTACT
>MC1_01A20
TGCACTAATAATCTGACAAAGGCTAGTAATGCCACTGAGAAGGCTAATAATACCATTACT
>MC1_01TA1
TGCACTAATAATCTGATT--------AATATCACTGAGAATACTAATAATACCATTACT
```

# Input File Formats

Phylogenetics program input file formats

## PHYLIP

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

```
5 60
MC1_01B4fsTGCACTAATAATCTGATT---------AATATCACTGAGAATACTAATAATACCATTACT
MC1_01A10 TGCACT---AATCTGACAAAGGCTATTAAGACCAATGGGAATGCTAATAATACCAGTACT
MC1_01C1  TGCACTAATAATCTGATT--------AATATCACTGAGAATACTAATAATACCATTACT
MC1_01A20 TGCACTAATAATCTGACAAAGGCTAGTAATGCCACTGAGAAGGCTAATAATACCATTACT
MC1_01TA1 TGCACTAATAATCTGATT--------AATATCACTGAGAATACTAATAATACCATTACT
```

National Institute of
Allergy and
Infectious Diseases

# Input File Formats

Phylogenetics program input file formats

NEXUS

```
#NEXUS
begin data;
    dimensions ntax=9 nchar=1823;
    format datatype=dna interleave missing=-;
matrix
HCVT050   GGTCTTGGTCTACTGTGAGC GAGGAGGCCGGTGAGGACGT
HCVT142   GGTCTTGGTCTACCGTGAGT GAGGAGGCCACTGAGGACGT
HCVT169   GGTCTTGGTCTACCGTGAGC GAGGAGGCTAGTGAGGACGT
SE0307168 GGTCGTGGTCCACCGTGAAC GAGGAGGCTGGTGAGGACGT
HCVT221   GGTCTTGGTCTACCGTGAGC GAGGAGGCCAGTGAAGACGT
MD2_2     GGTCTTGGTCTACTGTAAGC GAGGAGGCTAGTGAGGACGT
HCV1b     GGTCTTGGTCTACCGTGAGC GAAGAGGCTGGTGAGGATGT
Contig000 GGTCTTGGTCTACCGTGAGC GAGGAGGCTAGTGAGGACGT
HCVT140   GGTCTTGGTCTACTGTGAGC GAGGAGGCTAGTGAGGATGT
;
end;
```

National Institute of
Allergy and
Infectious Diseases

NIH

# Input File Formats

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

National Institute of
Allergy and
Infectious Diseases

# What's next?

After the break


Building Bayesian trees with our MSA

# Bayesian Analysis

Calculating the posterior probability of the evolutionary parameters

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

where:

$\tau$ = tree topology

$v$ = branch lengths

$\theta$ = substitution parameters

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)

- $p(\theta|D)$ = (Likelihood x Prior) / probability of the data

- $p(\theta|D) = p(D|\theta)p(\theta) / p(D)$

National Institute of Allergy and Infectious Diseases

# What is Bayesian analysis?

Likelihood that this die is unbiased?

# 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

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

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 |

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.

# 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)

# 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
```

National Institute of
Allergy and
Infectious Diseases

# 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"

National Institute of
Allergy and
Infectious Diseases

# 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

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)

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;
```

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 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                      |
   |         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
-------------------------------------------------------------------------------------------------
```
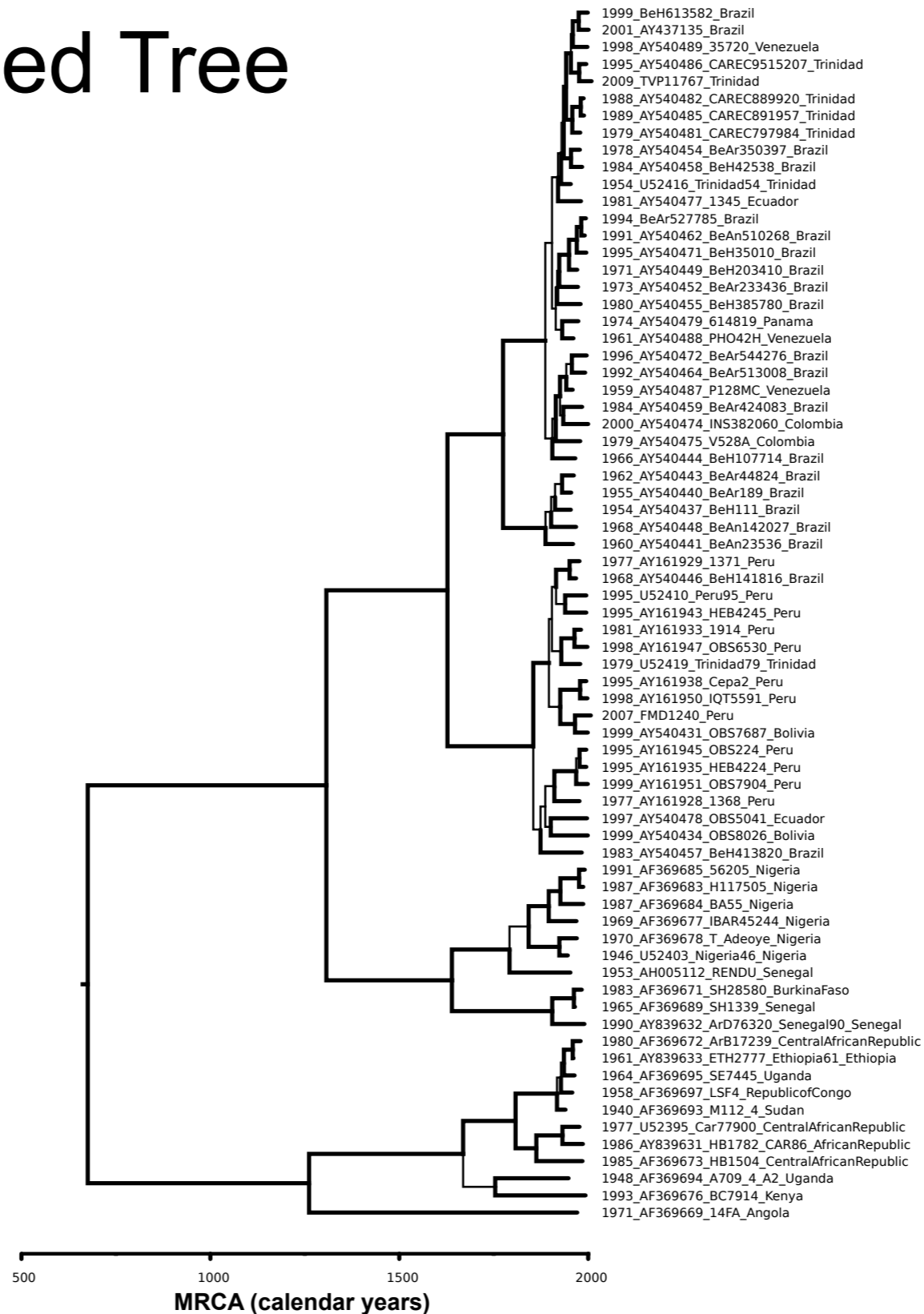
# Bayesian Analysis
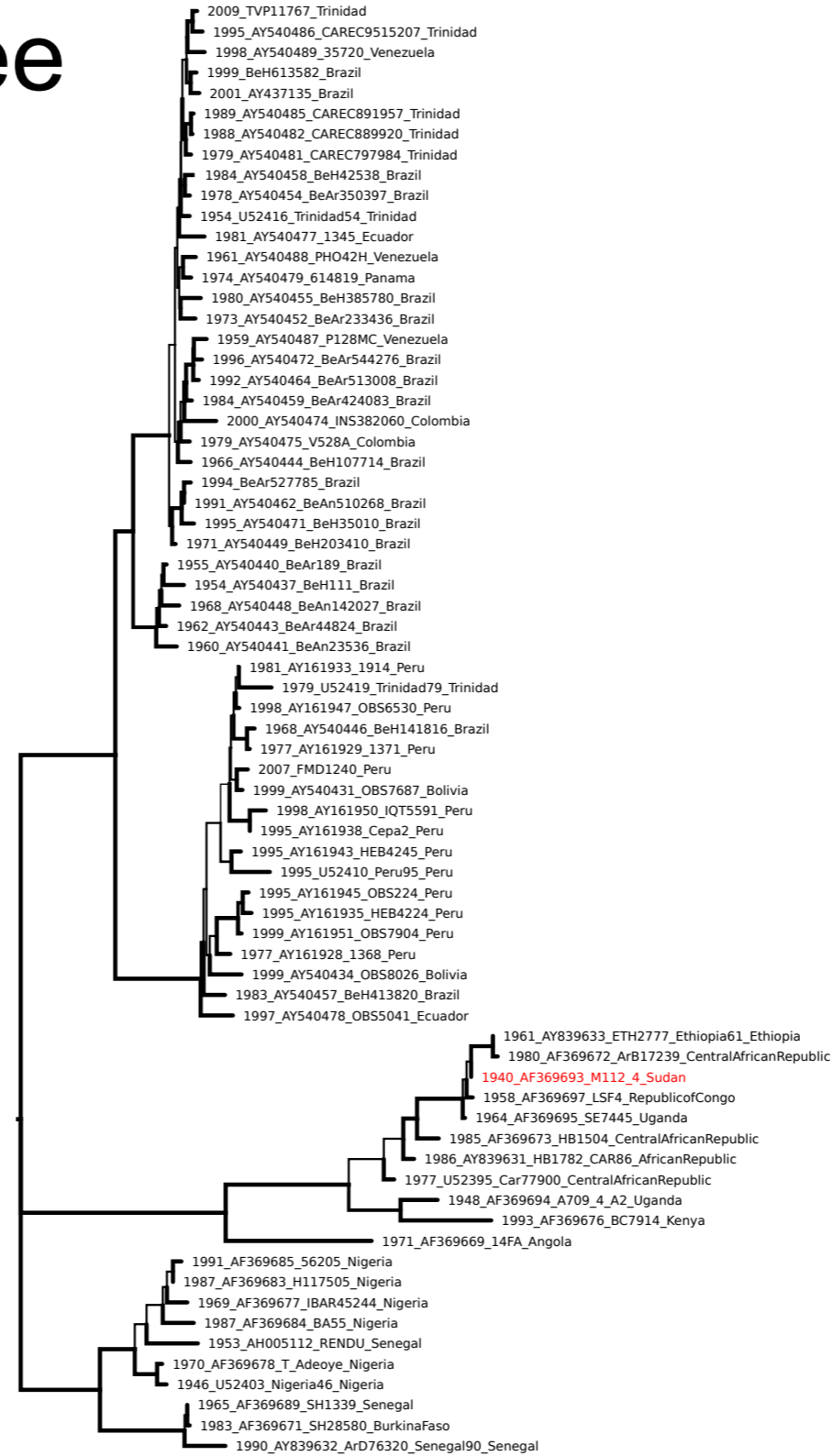
## The Consensus Tree

# Bayesian Analysis

The Time-Stamped Tree
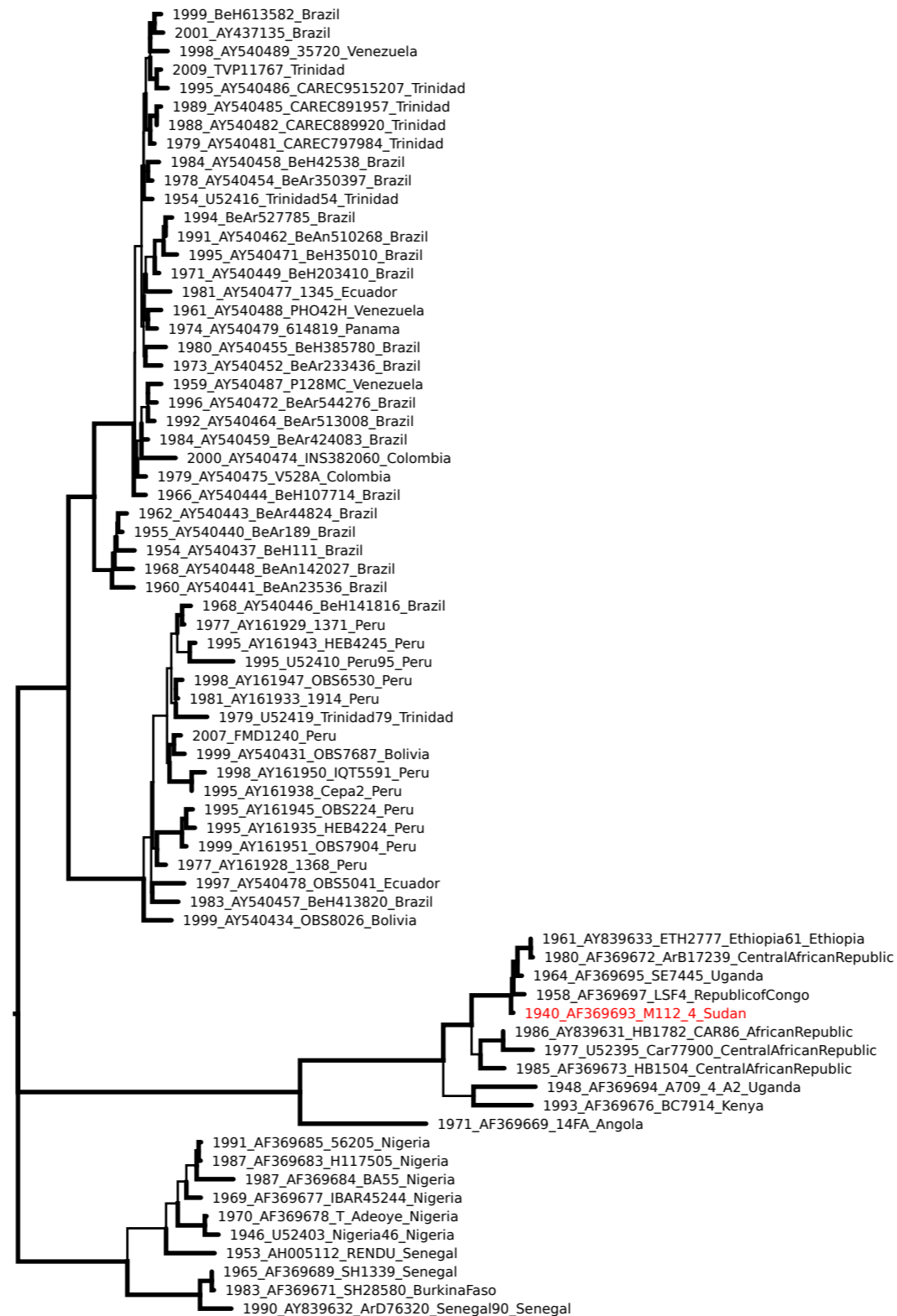
# Distance Analysis

## The Neighbor-Joining Tree

# Likelihood Analysis

## The Best Tree
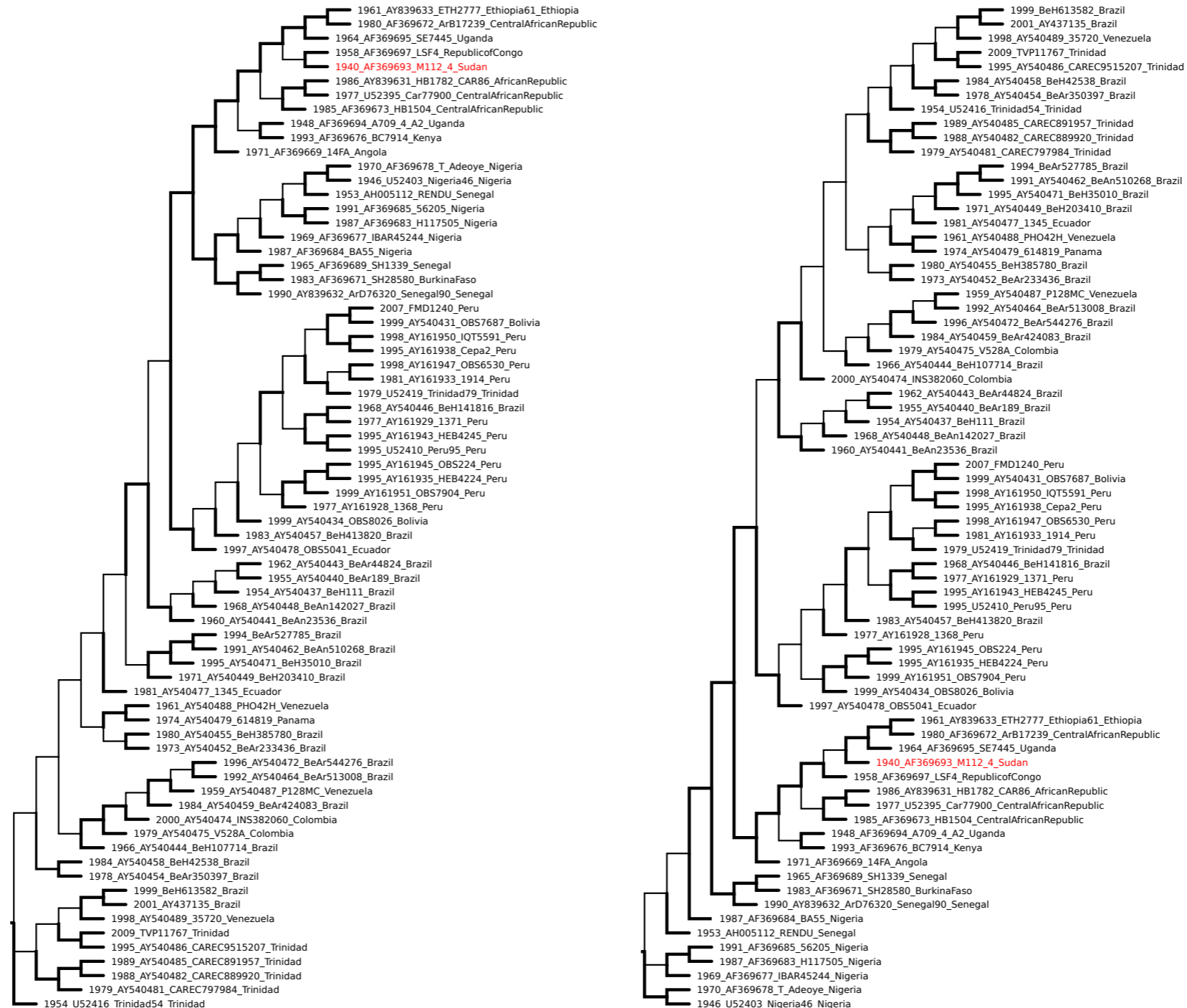


0.05

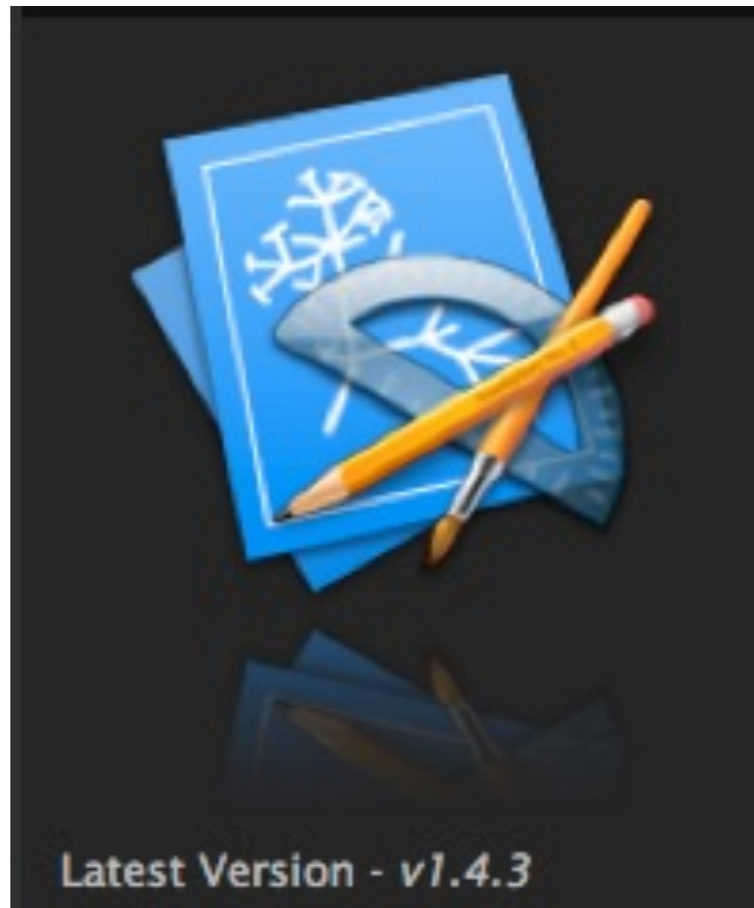# Parsimony Analysis

## Two Equally Parsimonious Trees

# Making a Tree Figure: FigTree



Latest Version - v1.4.3

http://tree.bio.ed.ac.uk/software/figtree/

National Institute of
Allergy and
Infectious Diseases

# In Conclusion

Where have we been? What have we done?

- Why are virus' biological sequences special?
- Calculating a multiple sequence alignment.
- Calculating trees using distance, parsimony, and likelihood.
  - How to calculate bootstrap support.
- Bayesian exploration of phylogeny posterior distribution.
- <span style="color:red">Always</span> use more than one tree generation algorithm
- Look for <span style="color:red">consensus</span> and investigate <span style="color:red">disagreement</span>

National Institute of
Allergy and
Infectious Diseases

NIH

# Questions?

Email us!

**bioinformatics@niaid.nih.gov**



National Institute of Allergy and Infectious Diseases