# Pathogen sequences and phylogenetic analysis

Biotechnology Solutions for Infectious Disease

frederick.r.jaya@student.uts.edu.au
github.com/fredjaya
@fredjaya1

# Previous lecture



STEC and raw milk



Outbreak of STEC O157:H7 associated with contaminated salad leaves



WGS of MRSA

**Objective:** To understand how pathogen genetics is applied to epidemiological investigation and studies

# Assessment 2

- Journal article (1500 words)
- Presenting your results on ARGs in bacteria
- **Phylogenetic tree**
  - Construction
  - Presentation

## Exercises

- Finding/identifying gene sequences
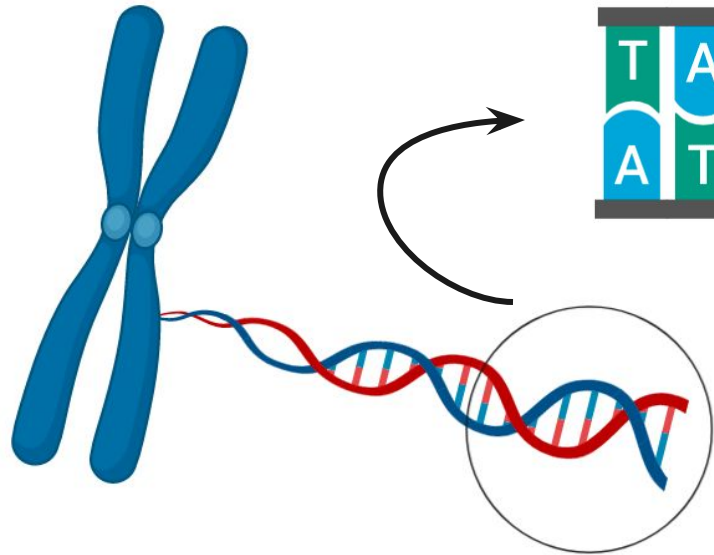- Applications of phylogeny

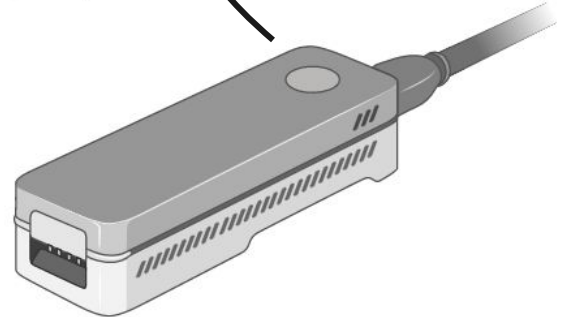1. GenBank

2. BLAST

3. Haiti outbreak

4. Vaccinating ebola

5. HIV transmission

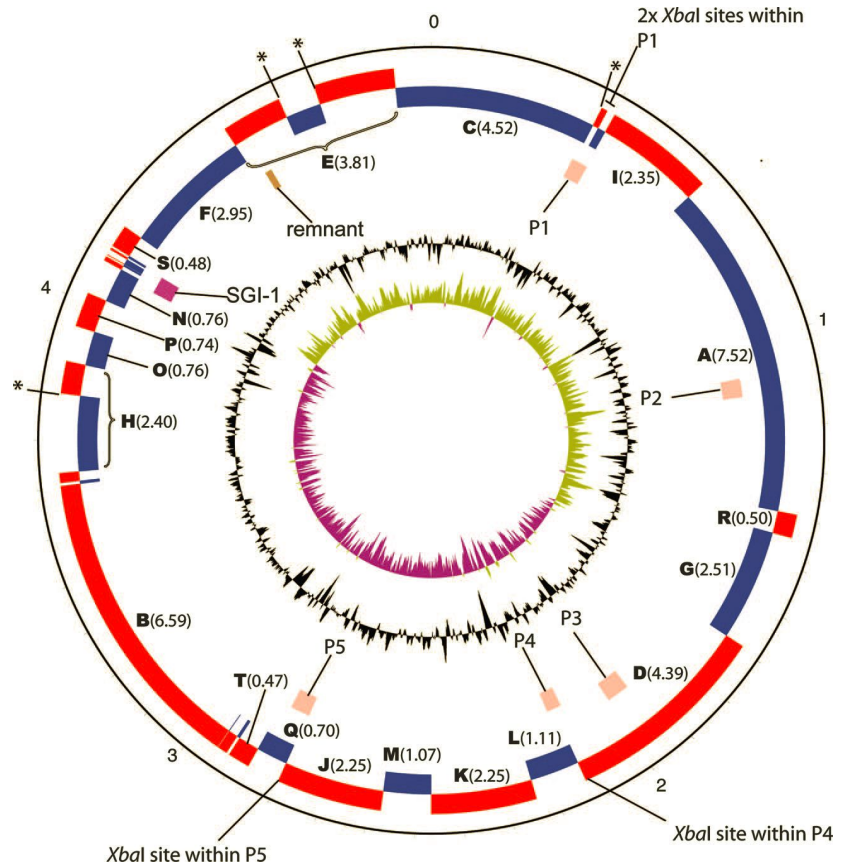# Pathogen sequences

# Genetic sequences



>seq1_F
TACGCTGA
>seq1_R
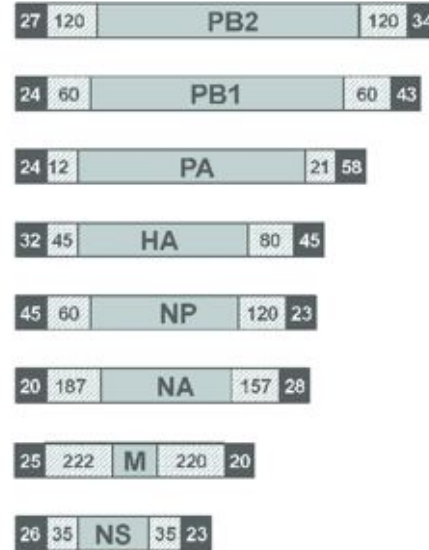ATGCGACT

# Bacterial whole genome

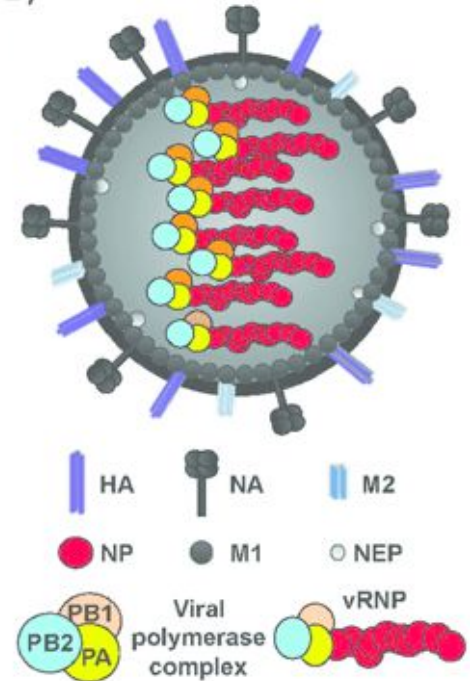- Physical map of the Salmonella serovar Typhimurium NCTC113348 genome
- circular

# Viral whole genome

- Influenza A virus genome organisation and virion structure
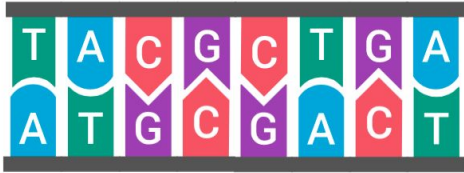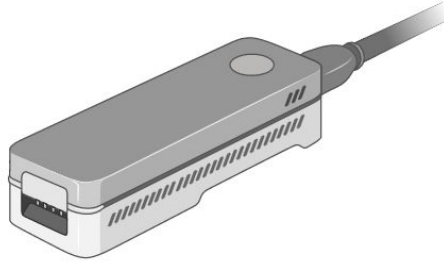- Linear, segmented

# Genetic sequences



**Biochemistry**

DNA/RNA, Protein



**Sequencing**

Whole genome, SNPs



**File format**

.FASTA, .FASTQ, .NEX

# Genetic sequences (FASTA)

**Salmonella enterica subsp. enterica serovar Typhimurium strain ABBSB1189-1 scaffold00001, whole genome shotgun sequence**
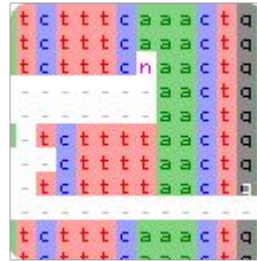
GenBank: LAPF01000001.1

GenBank   Graphics

```
>LAPF01000001.1 Salmonella enterica subsp. enterica serovar Typhimurium strain
ABBSB1189-1 scaffold00001, whole genome shotgun sequence
CATTGTCATTGCGCATATTCAGGTTGATATGGCGTATAAACACCGGTCAGGCTTTCCAGGAAGGCGACGA
TATCGTCAATATCGTTTTGCGGCAGATCGGTGCCAACCTGATAACGCAGCATCAGTTTTACCGCTCCATC
CAGCGTCGGTACGTCGCCCCGATGGAAATAAGGCGCTGTTAACGCGACGTTGCGTAAGCCCGGGACTTTT
TGCCGTAATTTATCGCGAACCTCTTTGGTGACGTTCATACGACCAATATCCGCTGCGGTAATTTCGCCAA
AATTAAAGTCTCGCTTTAATCCCAATGGTTCAAAAGAGCGCCCGCCTAAAATGATACCGCCGTGACAGGT
TGCACATTTATTCTCTTTAAATAATTGATAACCGTGTTTCTGTTGCGCGGTCAGCGCATTTTCATCTCCA
CGTAGCCATTTATCAAAGGCGGAATCCGGCGTTATCAACGTTTTTTCGAATTCGGCGATCGCATCAGTAA
TATTTTCCCCGGTAAATCCTTGCGGATAAACCGCCTGGAAATCTTTTTTCAGGACAGGATCTTTATCAAG
CTTGCTAATAATTTCATCCCAGGATTTAGAGGCCATTTCAATAGGATTTAACGGTGGTCCTCCTGCTTGC
TCCTGCAGGGTTGCAGCACGACCATCCCAAAATTGTTCGATATTAAATACGGAGTTGAATACCGTCGGCG
CATTTATTGGTCCTACCGCACCGCCAACGCCAATTGAGGTTTTTCTGCCATCGACACCGCCCGCATTTAA
CGCGTGACAATGCGCACAGGATATTGTGCTGTCGCCGGATAAACGTTCATCATGATAAAGCCGGAAACCT
AAGTCGACTTTTTTCGCATCGACGGGAATATTGCGCGGAATAGGCTGAACGGGTTCATTCCGGTGCGCCG
```

Header

Sequence
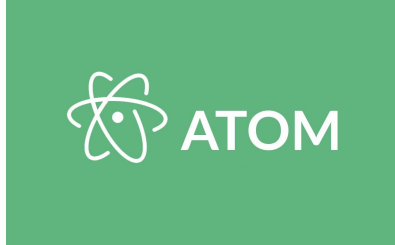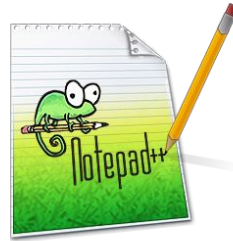
# Tools - Sequence alignment viewers


geneIOUS


AliView


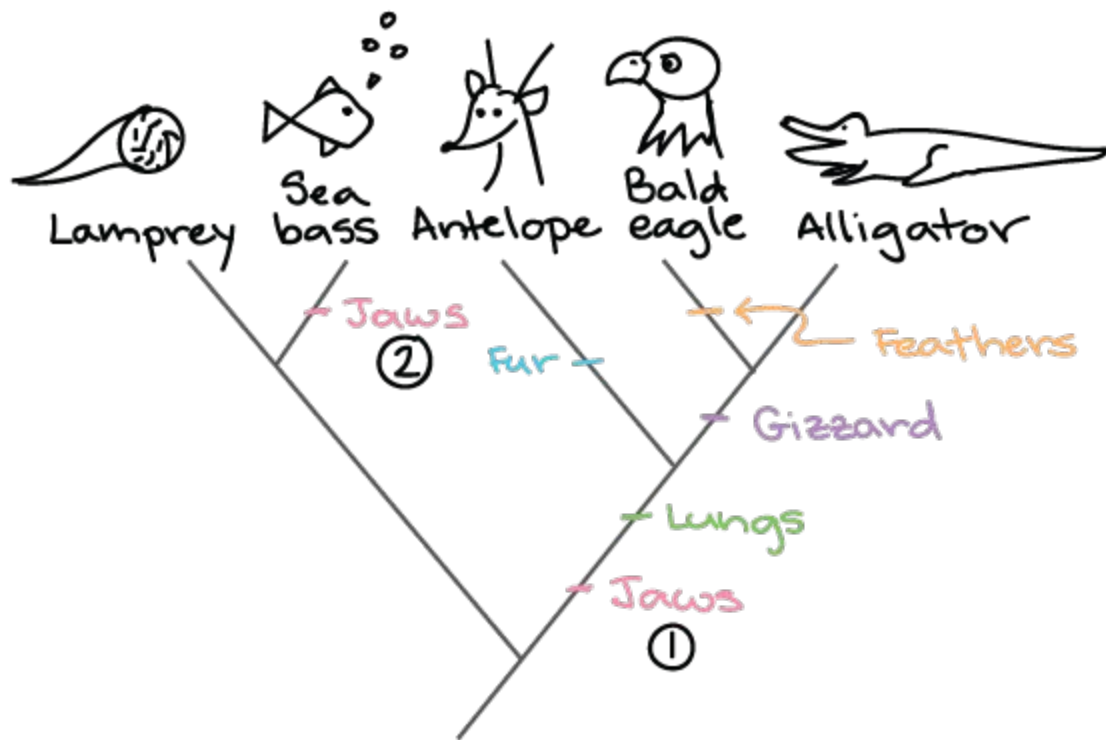Seqotron

# Tools - text editors



Sublime
Text 3

Notepad++

Notepad

# Exercise 1 - NCBI GenBank

# Exercise 2 - BLAST

# Phylogenetics

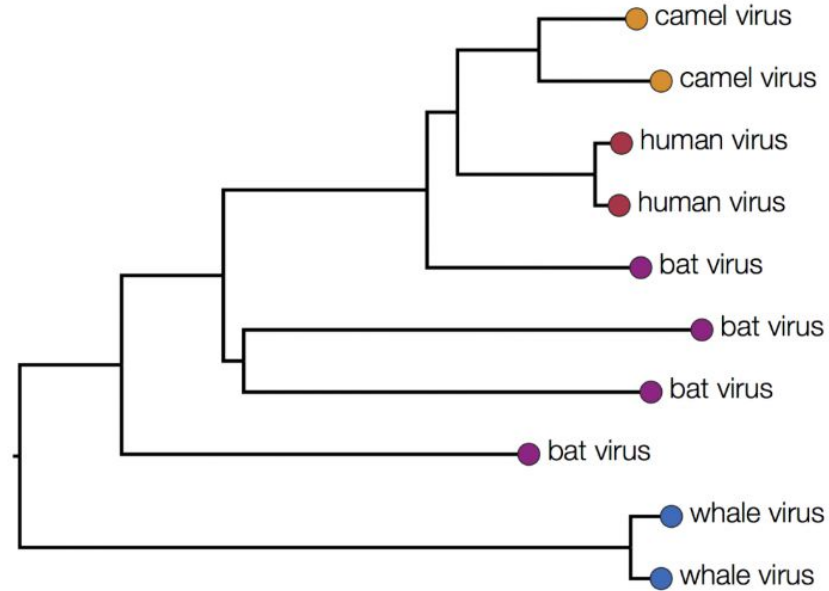# 1. Relationships between species

# 2. Evolutionary changes and history

# Phylogenetic tree



Darwin's notes (1837)



"Modern" phylogeny

(Cladogram)

**1. Relationships between species**
**TOPOLOGY**

**2. Evolutionary changes and history**
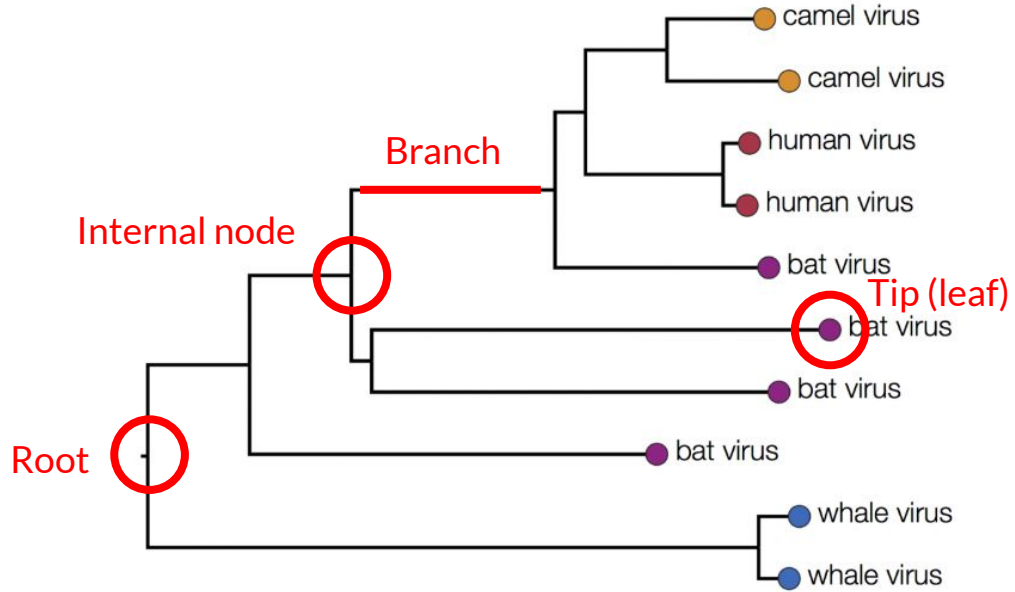**BRANCH LENGTH**

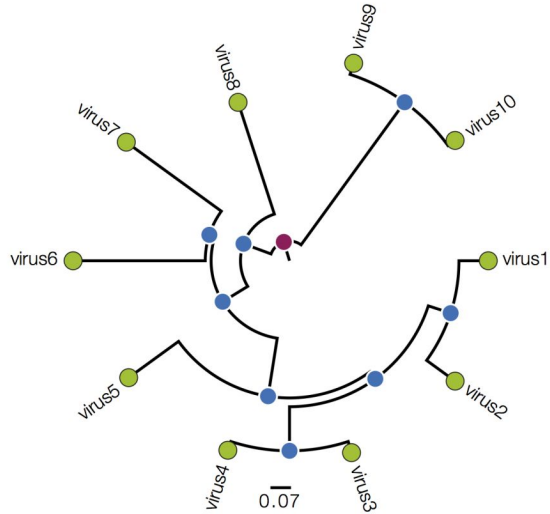# BRANCH LENGTH



TOPOLOGY

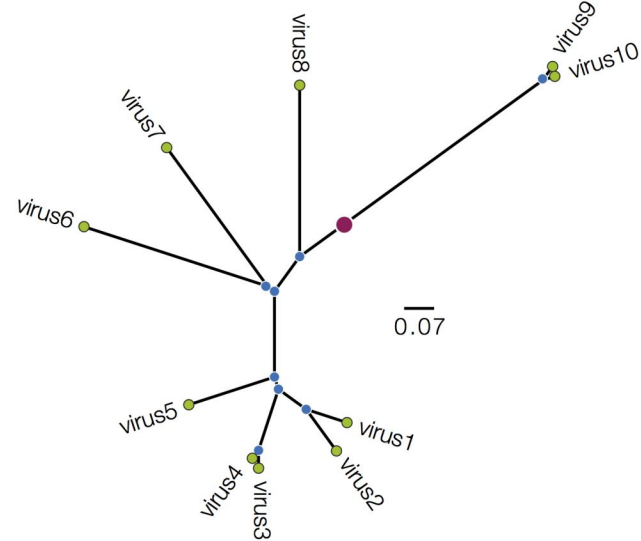Components of a phylogenetic tree

Components of a phylogenetic tree

# Variety of tree formats...



Circular/polar



Unrooted

# Applications

- Wide range of applications
- Epidemiology - origin, transmission, monitor outbreaks
- Clinical - Drug vaccine, design
- Conservation - identify diversity hotspots
- Cultural - Evolution of art, music, linguistics
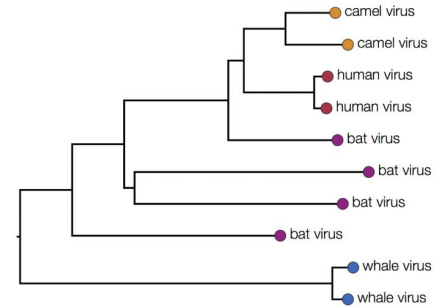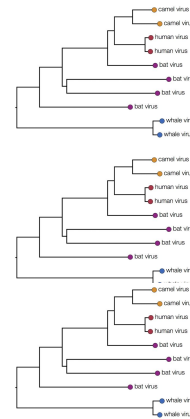
PHYLOGENETICS PROGRAMS

**Data**

**Model Selection**

**Tree inference**

**Estimation**

Inference workflow

# Sequence alignment



Camel ACTCGAT

Human ACTCT

Bat ACCCGT

Whale ACCCTTT

→

ACTCGAT
ACTC‑‑T
ACCCG‑T
ACCCTTT

Homologous sites aligned

# Inference methods

- Parsimony
- Distance
- Maximum likelihood
- Bayesian

+

$P(D|H)$

"Which phylogeny and evolutionary model best explains our data?"

Maximum likelihood estimation

# Evolutionary models

| Base frequencies | Substitution rates | Rate variation |
|:---:|:---:|:---:|
| $\pi_A + \pi_C + \pi_G + \pi_T = 1$ | $\rho, \alpha, \beta, \mu$ | +I +G |

Equal base frequencies (3 df)

| | JC | F81 | K80 | | | HKY | SYM | GTR |
|---|---|---|---|---|---|---|---|---|
| Base frequencies | $\pi$ | $\pi_A\pi_C\pi_G\pi_T$ | $\pi$ | | | $\pi_A\pi_C\pi_G\pi_T$ | $\pi$ | $\pi_A\pi_C\pi_G\pi_T$ |
| Substitution rates | $\rho$ | $\rho$ | $\alpha\beta$ | | | $\alpha\beta$ | $\mu_1\mu_2\mu_3\mu_4\mu_5\mu_6$ | $\mu_1\mu_2\mu_3\mu_4\mu_5\mu_6$ |

JC vs F81

A        R

Transition rate equals
Transversion rate (1 df)

JC vs K80        F81 vs HKY

A    R        A    R

Equal transition rates and
Equal transversion rates (4 df)

K80 vs SYM        HKY vs GTR

A    R        A    R

Rates equal among sites (1 df)

| JC vs JC+Γ | K80 vs K80+Γ | SYM vs SYM+Γ | F81 vs F81+Γ | HKY vs HKY+Γ | GTR vs GTR+Γ |
|---|---|---|---|---|---|

A  R   A  R   A  R   A  R   A  R   A  R   A  R

No invariable sites (1 df)

| JC vs JC+I | JC+Γ vs JC+I+Γ | K80 vs K80+I | K80+Γ vs K80+I+Γ | SYM vs SYM+I | SYM+Γ vs SYM+I+Γ | F81 vs F81+I | F81+Γ vs F81+I+Γ | HKY vs HKY+I | HKY+Γ vs HKY+I+Γ | GTR vs GTR+I | GTR+Γ vs GTR+I+Γ |
|---|---|---|---|---|---|---|---|---|---|---|---|

A  R  A  R  A  R  A  R  A  R  A  R  A  R  A  R  A  R  A  R  A  R  A  R

JC JC+I JC+Γ JC+I+Γ  K80 K80+I K80+Γ K80+I+Γ  SYM SYM+I SYM+Γ SYM+I+Γ  F81 F81+I F81+Γ F81+I+Γ  HKY HKY+I HKY+Γ HKY+I+Γ  GTR GTR+I GTR+Γ GTR+I+Γ
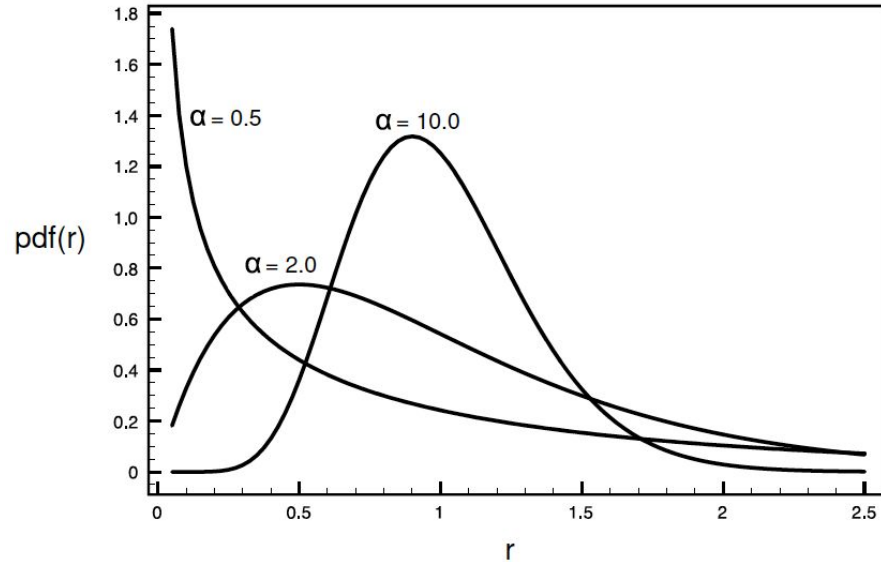
Hierarchical model test (Posada and Crandall, 1998)
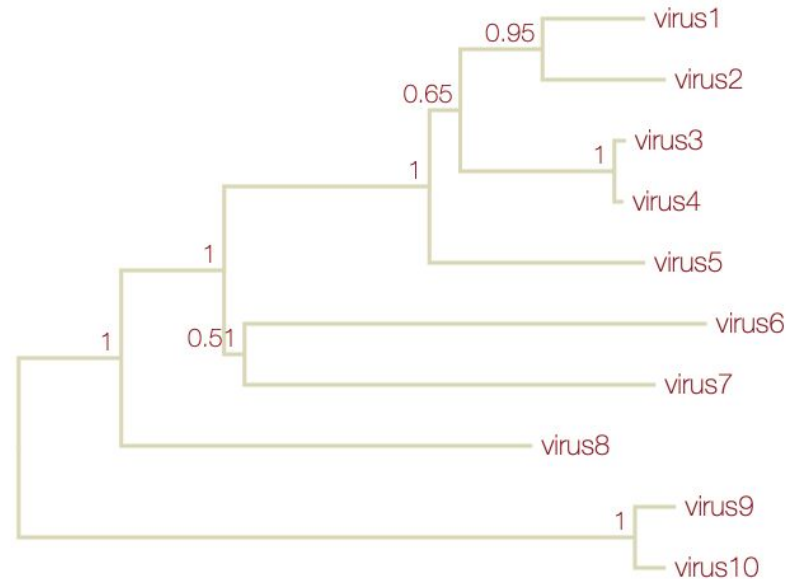
# Rate variation



Codon positions mutate at different rates (3rd > 2nd > 1st)

(**+G**) The gamma distribution models rate heterogeneity

(**+I**) Invariant sites assumes "sites do not vary"

# Statistical tests

- Bootstrap - measures the certainty of a tree estimate
- Model selection tests:
  - Likelihood ratio tests
  - Akaike Information Criteria (AIC)
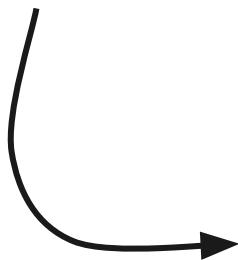  - BIC

# Not all sites evolve at the same rate

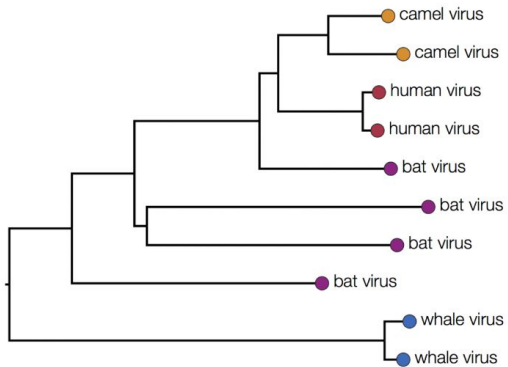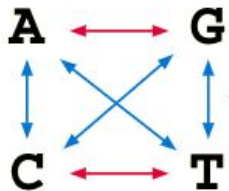$P(D|H)$

"Which mode of coin toss and type of coin best explains our data?"

A conceptual example

+

$P(D|H)$

"Which phylogeny and evolutionary model best explains our data?"

Maximum likelihood estimation

# Exercise 3. Haiti outbreak

# Exercise 4. Vaccinating ebola

# Exercise 5. HIV transmission