



# 3000788 Intro to Comp Molec Biol

## Lecture 7: Phylogenetics and evolutionary models

Fall 2025



**Sira Sriswasdi, PhD**

- Research Affairs
- Center of Excellence in Computational Molecular Biology (CMB)
- Center for Artificial Intelligence in Medicine (CU-AIM)

# Today's agenda



- Evolutionary perspective
- Phylogenetics



# Evolutionary perspective

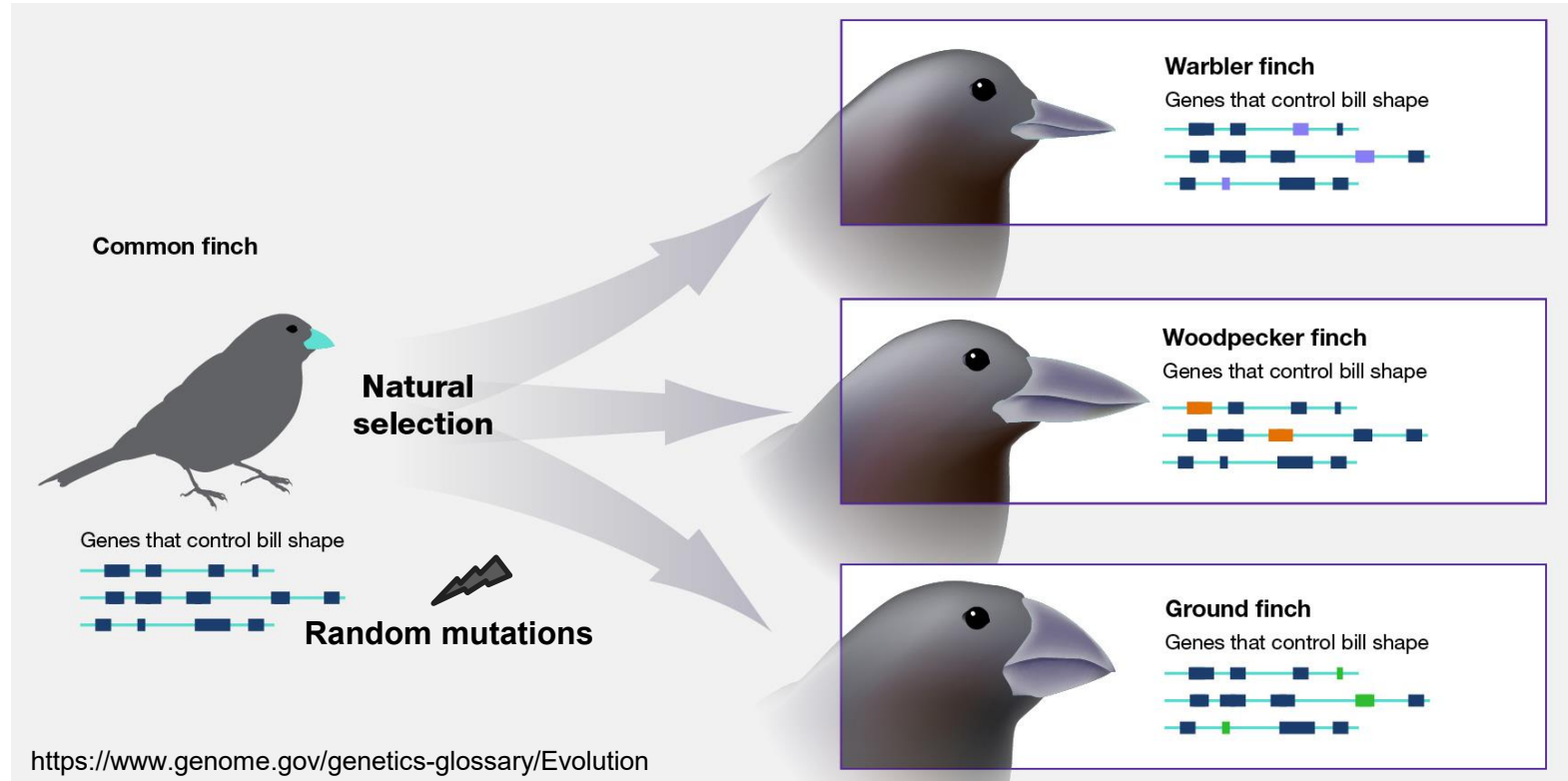
# Biology developed through and follows evolution



shutterstock.com · 1652460871

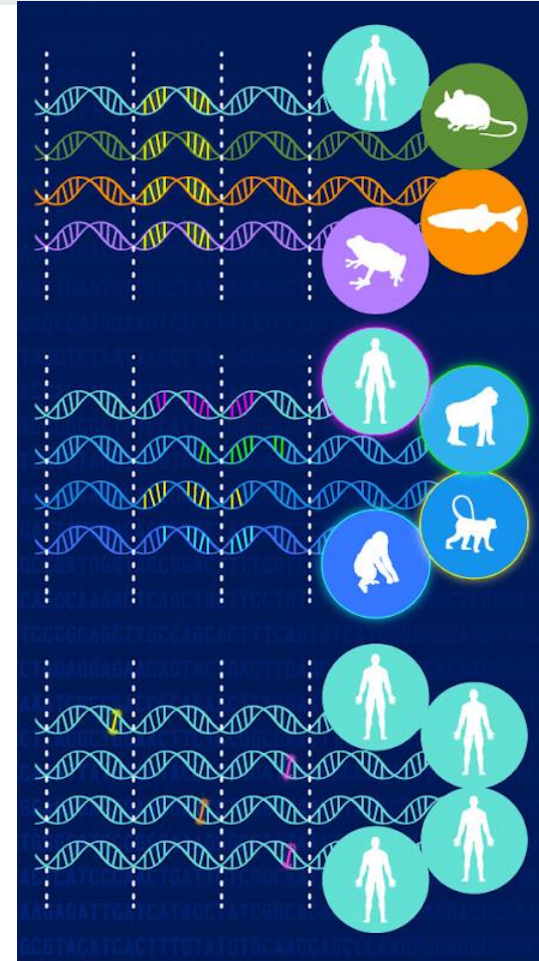
- Mutations occurred randomly without knowing their consequences
- Evolution is the process of competition and selection that weeds out unfit mutations and retains beneficial ones
- Nature is one large experiment

# Natural selection

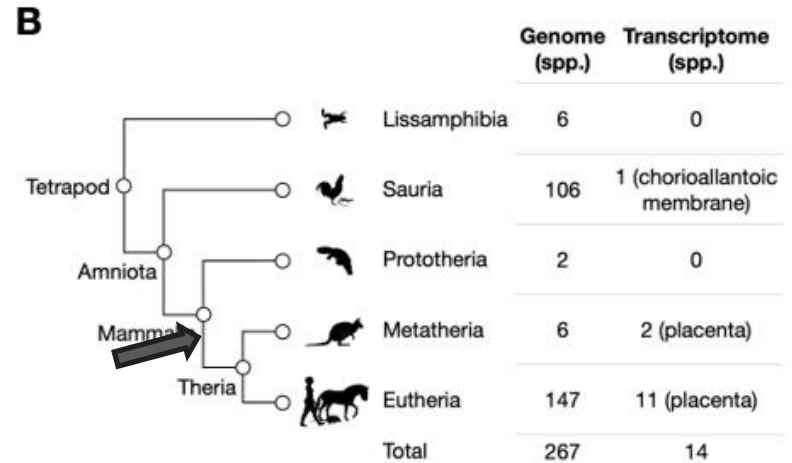
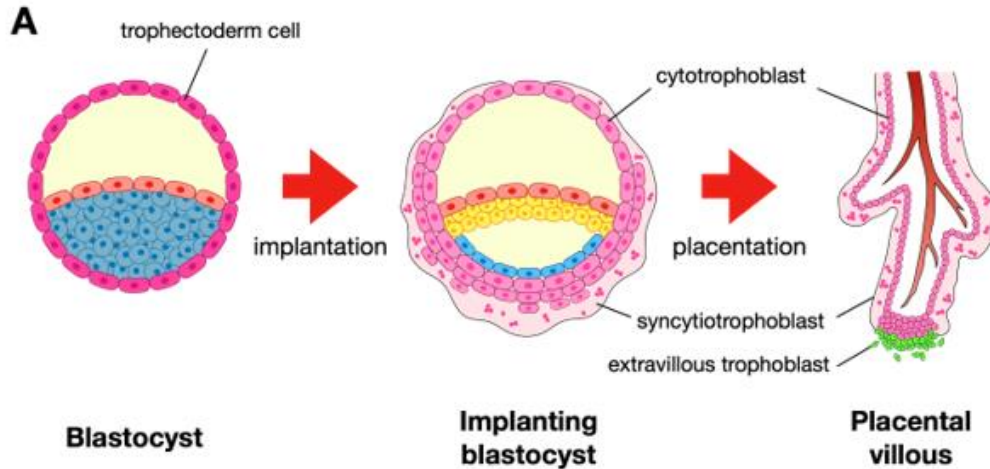


# Comparative genomics analysis

- Compare genomes of individuals/organisms with and without a phenotype
- Pinpoint causative genetic factors of the phenotype
- Many scopes
  - **Genetic disease:** Human populations
  - **Brain development:** Primates
  - **Novel enzymes:** Microorganisms



# The rise of placenta



Plianchaisuk, A. et al. Mol Biol Evol msac176 (2022)

- Placenta emerged as a new organ in mammalian ancestor
- Genes essential for placenta development should be acquired during evolution around the same period when placenta emerged

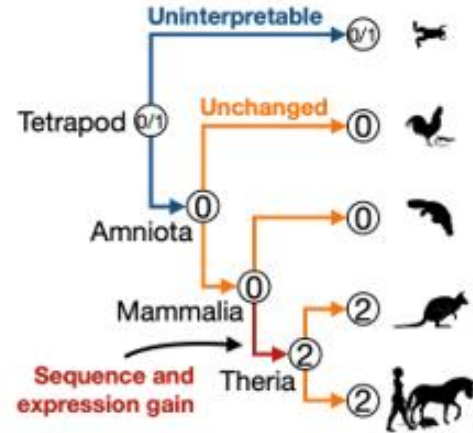
# The search for genes involved in placenta development



Plianchaisuk, A. et al. Mol Biol Evol msac176 (2022)

Without evolutionary perspective, we have to develop an animal model with defects in placenta development

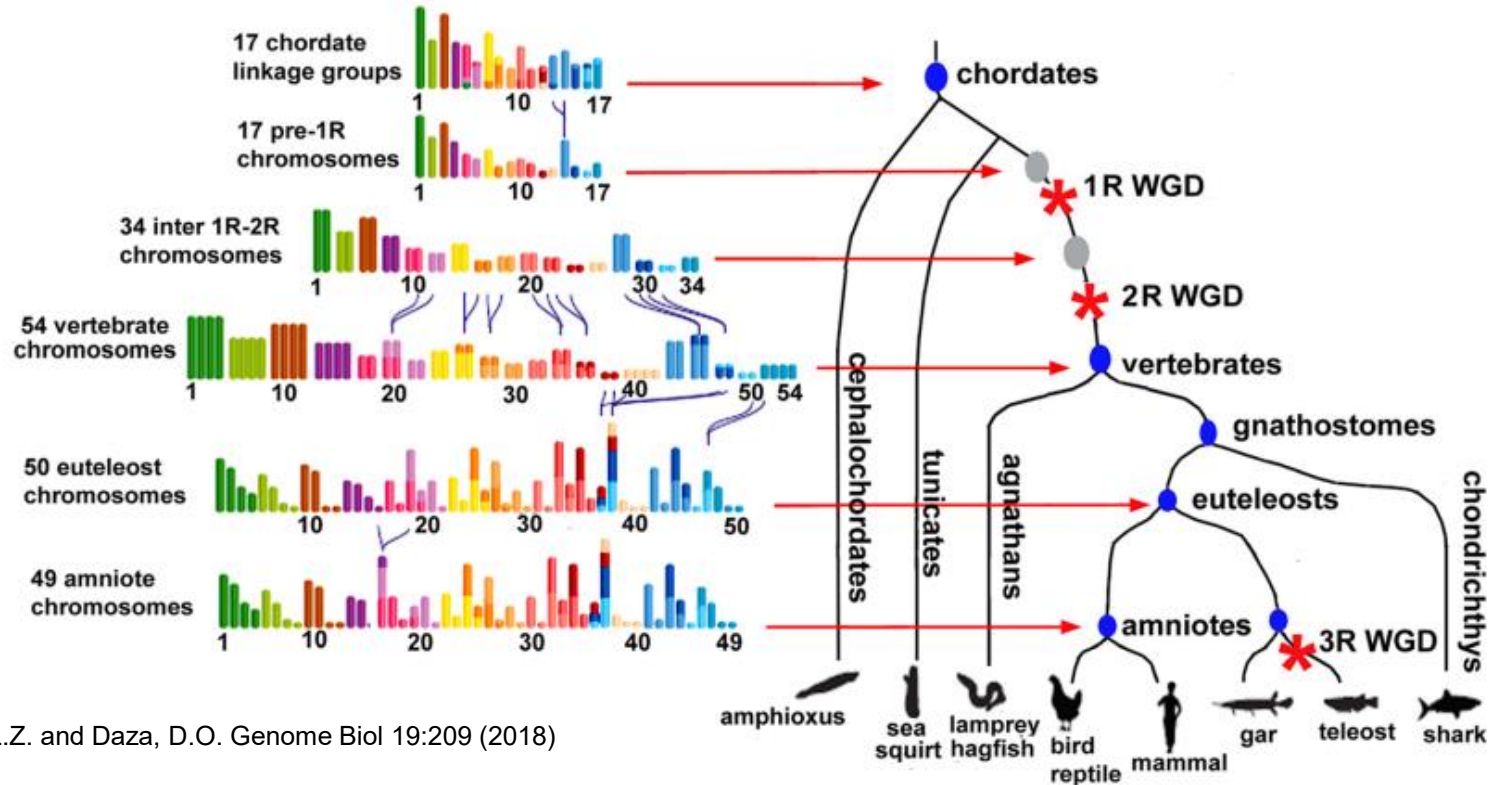
- Model function acquisition as two stages:
  - Gain of gene sequence (viral integration, gene duplication, etc.)
  - Activation of gene expression



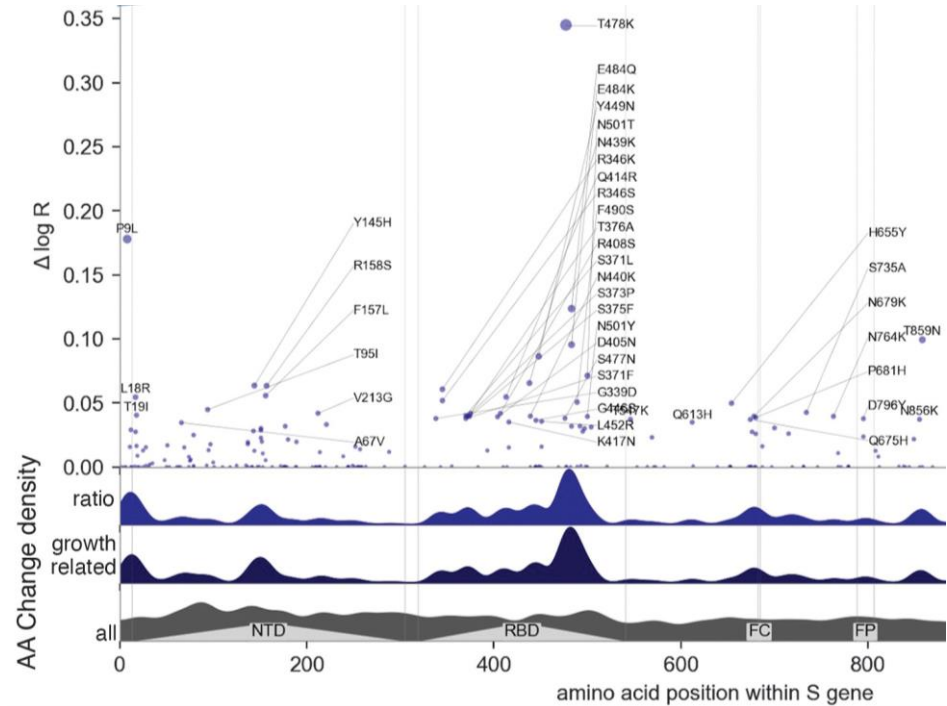
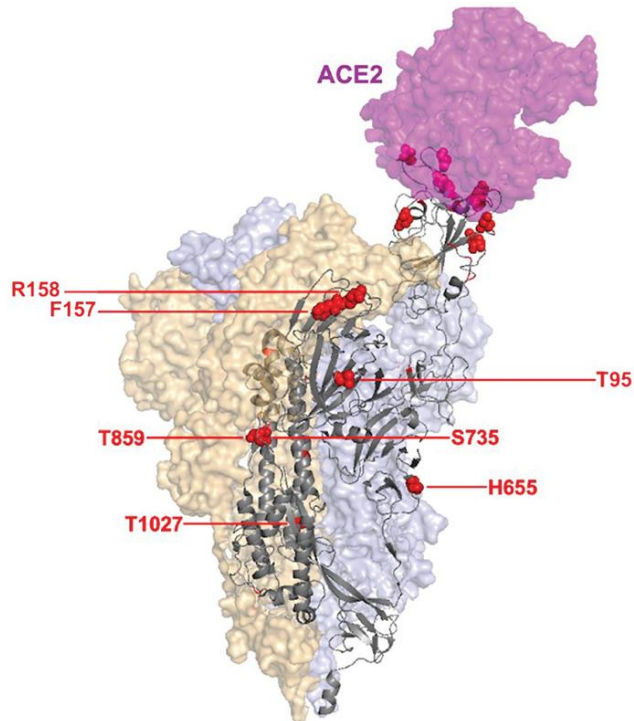
Interpreting character state changes from ancestors to descendants to transition events



# Genome expansions that gave rise to us



# Amino acid evolution in Coronavirus's Spike protein



# Applications of evolutionary perspective

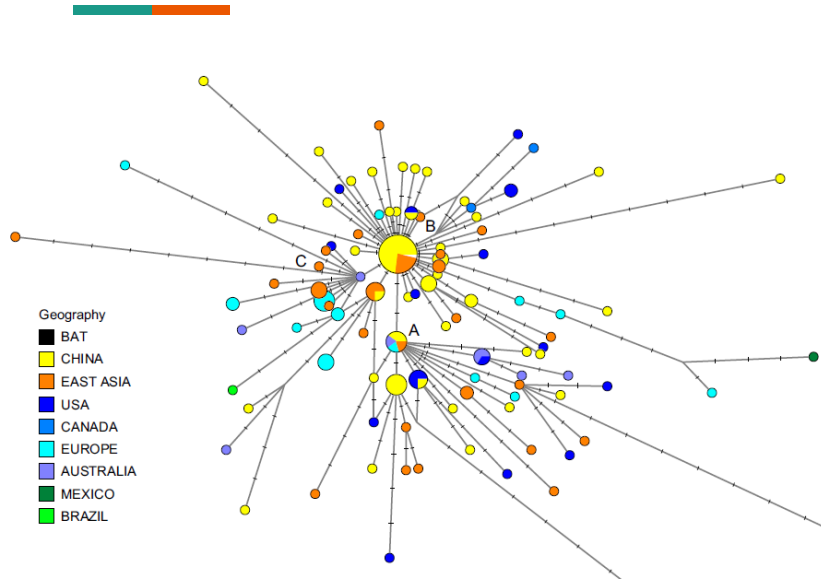


- Can be utilized in broad biology, not just evolutionary research
- Mutations in **conserved regions** probably have more chance to cause some phenotypes than those in **variable regions**
- Where to modify your enzyme to achieve new specificity?
- Where on a viral protein to target with antibody?
- Which group of genes are necessary for certain functions?

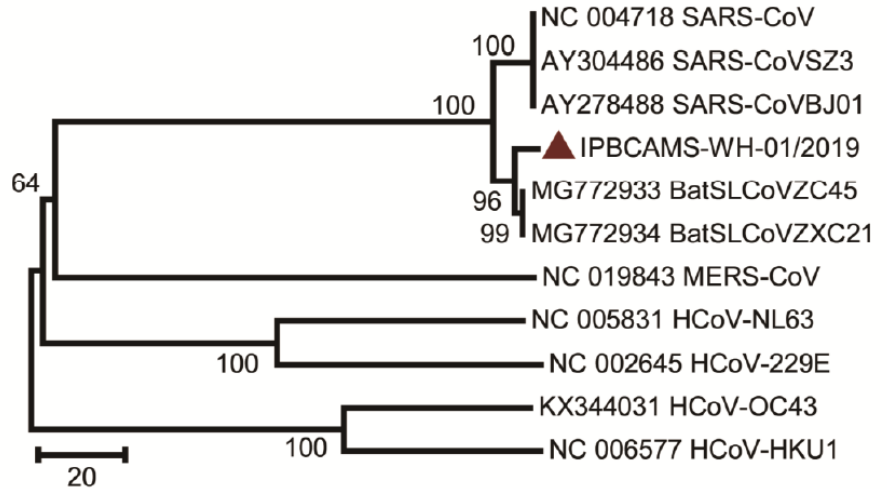


# Phylogenetics

# Phylogeny



Forster *et al.* PNAS (2020)



Guo *et al.* Clin Infect Dis (2020)

- Clustering of similar taxa with respect to evolutionary similarity
- Branch length reflect **time** and **mutation rate**

# Types of phylogeny

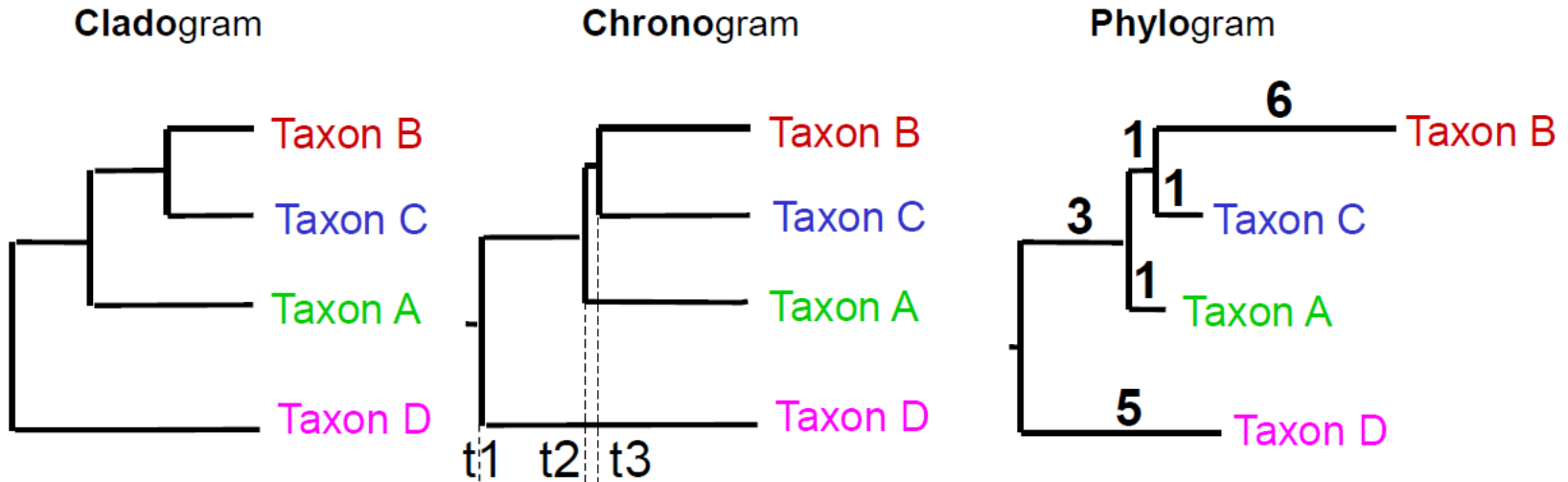


Image from Lecture 18 6.047 MIT OCW

- Depends research question: Is grouping or time more important?

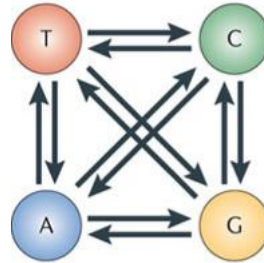
# Ingredients for phylogenetic reconstruction

## Multiple sequence alignment

Scarites	C	T	T	A	G	A	T	C	G	T	A	C	C	A	A	-	-	-	A	T	A	T	T	A	C	
Carenum	C	T	T	A	G	A	T	C	G	T	A	C	C	A	C	-	-	T	A	C	-	T	T	T	A	C
Pasimachus	A	T	T	A	G	A	T	C	G	T	A	C	C	A	C	T	A	T	A	A	G	T	T	T	A	C
Pheropsophus	C	T	T	A	G	A	T	C	G	T	T	C	C	A	C	-	-	-	A	C	A	T	A	T	A	C
Brachinus armiger	A	T	T	A	G	A	T	C	G	T	A	C	C	A	C	-	-	-	A	T	A	T	A	T	T	C
Brachinus hirsutus	A	T	T	A	G	A	T	C	G	T	A	C	C	A	C	-	-	-	A	T	A	T	A	T	A	C
Aptinus	C	T	T	A	G	A	T	C	G	T	A	C	C	A	C	-	-	-	A	C	A	T	A	T	A	C
Pseudomorpha	C	T	T	A	G	A	T	C	G	T	A	C	C	A	C	-	-	-	A	C	A	T	A	T	A	C

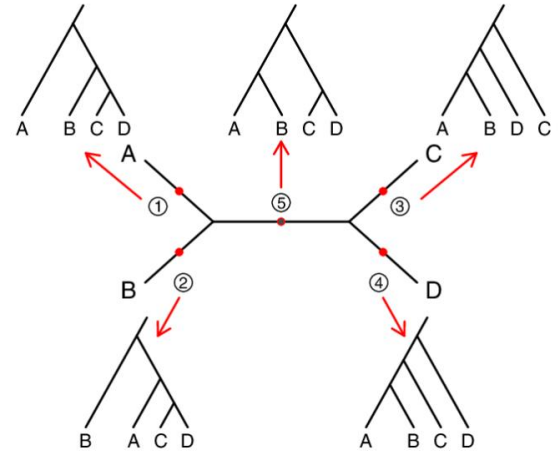
Image from [www.mcqbiology.com](http://www.mcqbiology.com)

## Evolutionary model



Yang & Rannala. Nat Rev Genetics (2012)

## Tree building algorithm



Tian, Y. and Kubatko, L.S. BMC Evol Biol 17(1) (2017)

- Sequence data + **evolutionary model** + tree construction
- Other constraints and non-molecular data can also aid the reconstruction
  - Sampling date, phenotype, lifestyle, geography, etc.

# Extra features for phylogenetic reconstruction

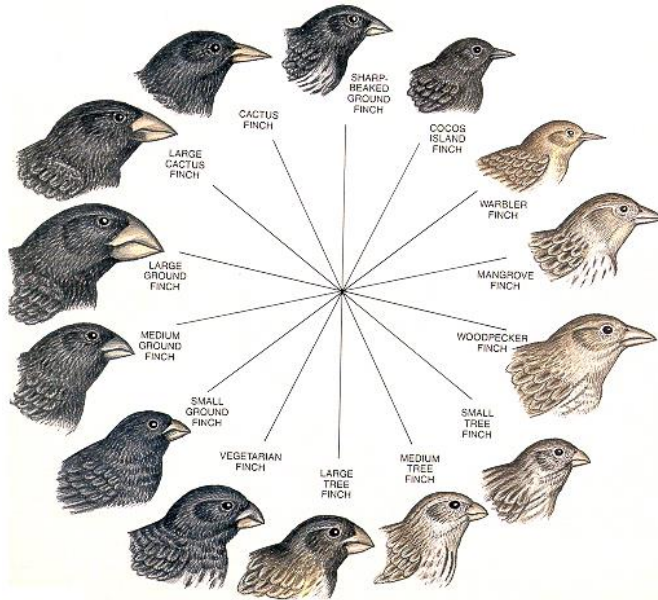
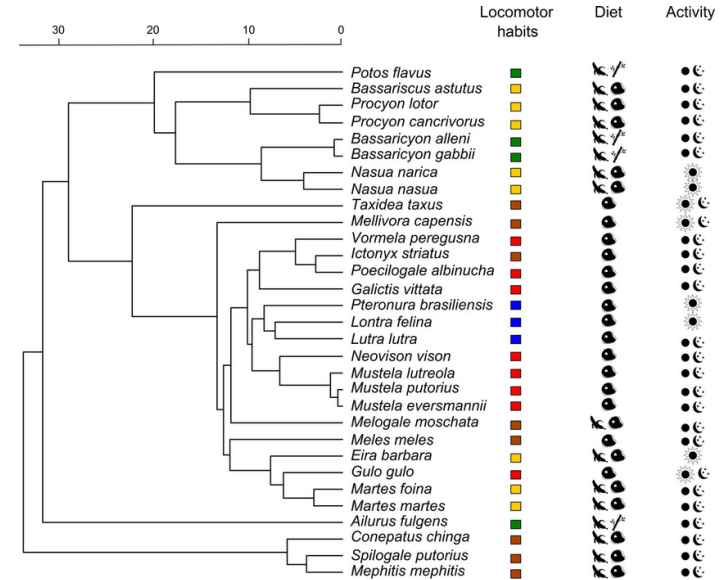


Image from pinterest



## Legends:

Locomotor habits

- Aquatic
- Terrestrial
- Arboreal
- Semi-fossorial
- Semi-arboreal

Diet

- Carnivorous
- Herbivorous
- Omnivorous

Activity

- Diurnal
- Nocturnal/crepuscular
- Arrhythmic





# **Evolutionary models:**

## **Nucleotide/amino acid substitution models**

# Recall nucleotide alignment scores

**Scoring Parameters**

Match/Mismatch Scores	1,-2 ?
Gap Costs	Linear ?

$$\text{Score} = +1+1-1-1-1+1+1+1+1 = +3$$

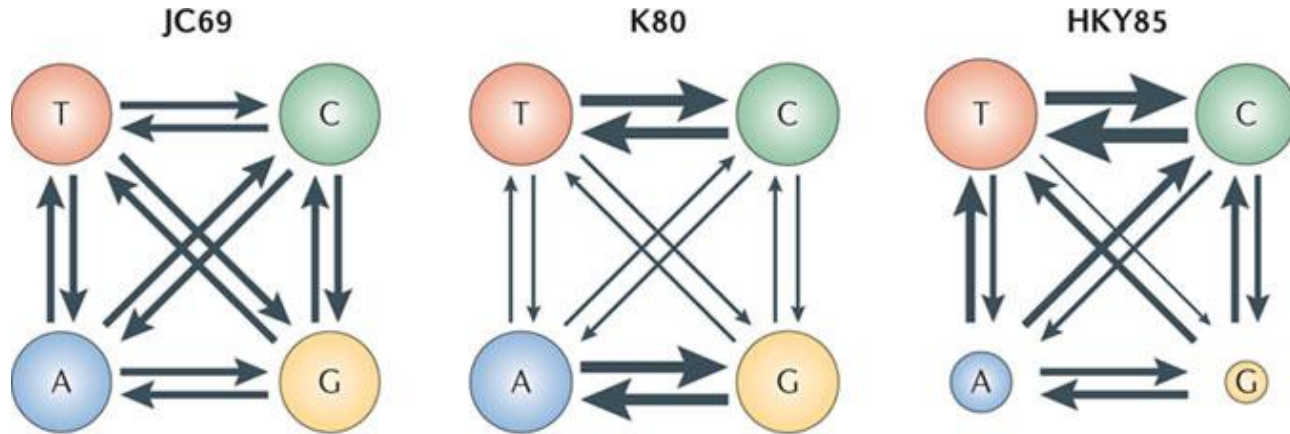
**Scoring Parameters**

Match/Mismatch Scores	2,-3 ?
Gap Costs	Existence: 5 Extension: 2 ?

$$\text{Score} = +2+2-5-2-2-2+2+2+2+2 = +1$$

- A special case of substitution model that considers only match or mismatch, each position independently
- That's not really how mutations occurred in nature


# Theoretical nucleotide substitution models



Yang & Rannala. Nat Rev Genetics, 13: 303-314 (2012)

- Juke-Cantor assumes equal base frequencies and equal substitution rates
- Kimura adds transition rate and transversion rate
- Hasegawa-Kishino-Yano adds different base frequencies

# How many parameters do these models used?



JC69	A	C	G	T
A	$1-3x_1$	$x_1$	$x_1$	$x_1$
C	$x_1$	$1-3x_1$	$x_1$	$x_1$
G	$x_1$	$x_1$	$1-3x_1$	$x_1$
T	$x_1$	$x_1$	$x_1$	$1-3x_1$

K80	A	C	G	T
A	$1-x_1-x_2$	$x_2$	$x_1$	$x_2$
C	$x_2$	$1-x_1-x_2$	$x_2$	$x_1$
G	$x_1$	$x_2$	$1-x_1-x_2$	$x_2$
T	$x_2$	$x_1$	$x_2$	$1-x_1-x_2$

- Juke-Cantor = 1 (substitution rate)
- Kimura = 2 (transition rate, transversion rate)
- Hasegawa-Kishino-Yano = 5 (transition rate, transversion rate, 3 base frequencies)

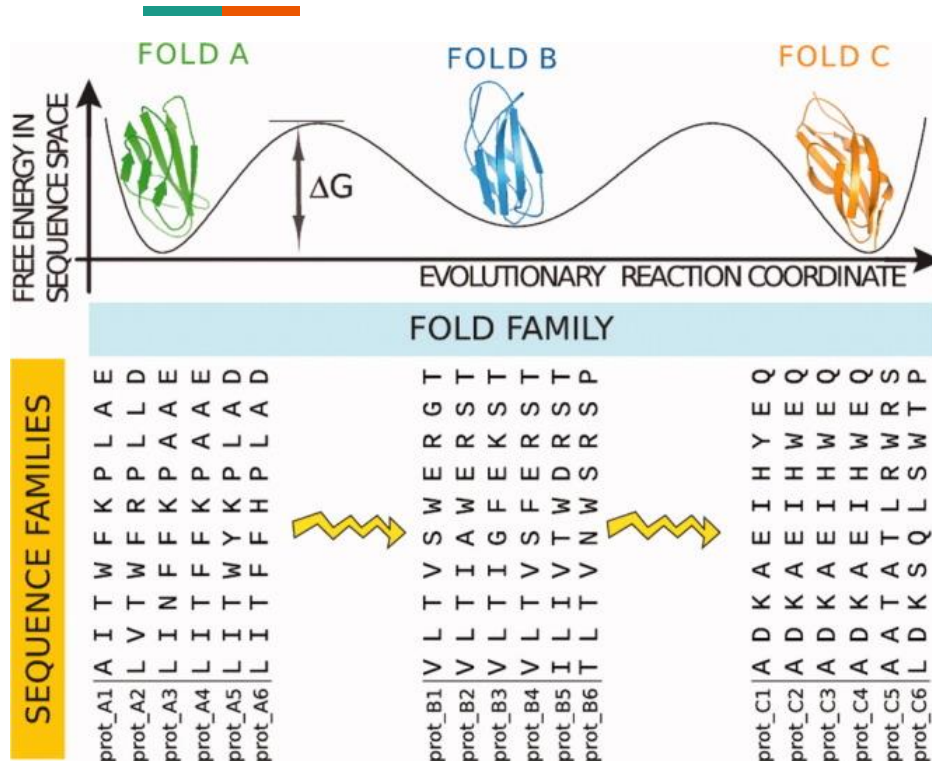
# General time-reversible model (GTR)



- **Symmetric** substitution rates:  $P(A \rightarrow G) = P(G \rightarrow A)$ 
  - 6 parameters (4 nucleotides choose 2)
  - **Time-reversible**: switching ancestral and descendant taxa does not change the calculation
  - Useful in practice because we often don't know which taxon came first
- Different nucleotide frequencies: 3 parameters
- This is the most generalized time-reversible model possible

	Ala	Arg	Asn	Asp	Cys	Gln	Glu	Gly	His	Ile	Leu	Lys	Met	Phe	Pro	Ser	Thr	Trp	Tyr	Val
Ala	4																			
Arg	-1	5																		
Asn	-2	0	6																	
Asp	-2	-2	1	6																
Cys	0	-3	-3	-3	9															
Gln	-1	1	0	0	-3	5														
Glu	-1	0	0	2	-4	2	5													
Gly	0	-2	0	-1	-3	-2	-2	6												
His	-2	0	1	-1	-3	0	0	-2	8											
Ile	-1	-3	-3	-3	-1	-3	-3	-4	-3	4										
Leu	-1	-2	-3	-4	-1	-2	-3	-4	-3	2	4									
Lys	-1	2	0	-1	-3	1	1	-2	-1	-3	-2	5								
Met	-1	-1	-2	-3	-1	0	-2	-3	-2	1	2	-1	5							
Phe	-2	-3	-3	-3	-2	-3	-3	-3	-1	0	0	-3	0	6						
Pro	-1	-2	-2	-1	-3	-1	-1	-2	-2	-3	-3	-1	-2	-4	7					
Ser	1	-1	1	0	-1	0	0	0	-1	-2	-2	0	-1	-2	-1	4				
Thr	0	-1	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	1	5			
Trp	-3	-3	-4	-4	-2	-2	-3	-2	-2	-3	-2	-3	-1	1	-4	-3	-2	11		
Tyr	-2	-2	-2	-3	-2	-1	-2	-3	2	-1	-1	-2	-1	3	-3	-2	-2	2	7	
Val	0	-3	-3	-3	-1	-2	-2	-3	-3	3	1	-2	1	-1	-2	-2	0	-3	-1	4

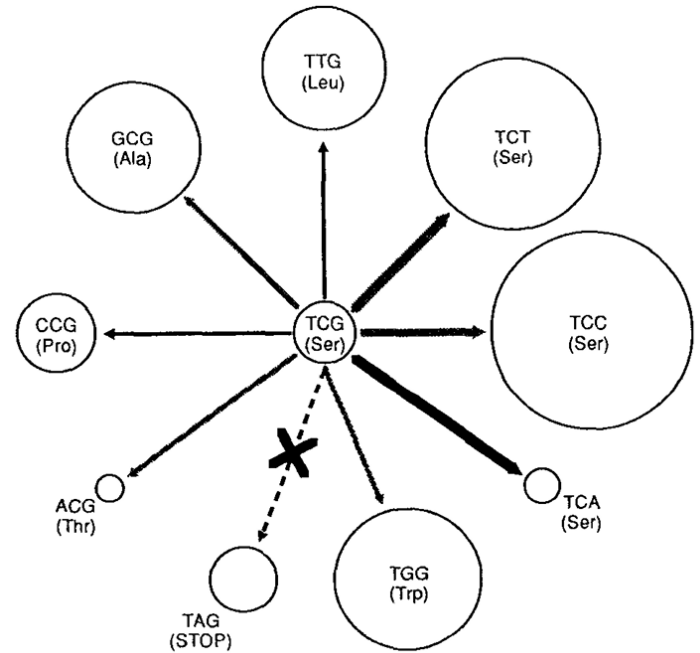
# Structural context for substitution model



- Different positions on different protein families have different evolutionary constraints
- Use energy from interactions between residues in 3D to develop substitution matrix
- Similar idea to **Position-Specific Substitution Matrix (PSSM)** used by PSI-BLAST

# Codon substitution model

- GY94 (Goldman-Yang)
- Learn base codon frequencies from nucleotide and amino acid data
- Codon neighborhood
  - Substitution rate depends on similarity between coded amino acids
- Non-synonymous / synonymous rate
- Similar idea to **BLASTX** and **tBLASTN**





# Codon alignment: amino acid → nucleotide

## A. DNA alignment

DNA alignment can disrupt codon frames

Q9FPK4	ATGGGTGTTTTCAGCTACGAGGATGAGGCCACC-TCCGTTATCCTCCCGCTA-GGCTGTTCAAGTCC-TTGTCCTAGATGCCGAC-AACCTCATT
Q9FPK3	ATGGGTGTTTTCAGCTACGAGGATGAGGCCACC-TCCGTTATCCTCCCGCTA-GGCTGTTCAAGTCC-TTGTCCTAGATGCCGAC-AACCTCATT
Q945E7	ATGGGTGTTGTGAGTTATGAGTTTGAGGTAACC-TCCCAATTGCTCCAGCCA-GGCTTTTCAAGGCTTTTGTTCTTGAGGCTGCC-AAGATTGG
Q6XC94	ATGGGTGTTGTGAGTTATGAGTTTGAGGTAACC-TCCCAATTGCTCCAGCCA-GGCTTTTCAAGGCTTTTGTTCTTGAGGCTGCC-AAGATTGG
Q6Q4B5	ATGGGTGTTGTGAGTTATGAGTTTGAGGTAACC-TCCCAATTGCTCCAGCCAAGGCTTTTCAAGGCTTTTGTTCTTGAGGCTGCCAAGGTTGG
Q43549	ATGGGTGTTTTCAATTACGAACTGAGTTTACC-TCCGTCATTCCCCCTGCTA-GGTTGTTCAATGCC-TTGTTCTTGATGCTGAC-AACCTCATC
Q4VPJ1	ATGGGTGTTTTACATACGAATCTGAGTCCACC-TCCGTCATCCCCCTGCTA-GGTTGTTCAATGCCACTGCTCTTGATGGTGAC-AAACTCATC
Q84LA7	ATGGGTGCTTTCACATACGAATCCGAA-TTACC-TCCGTCATCCCCCTGCTA-GGTTGTTCAATGCC-TTGTTCTTGATGCTGAC-AACCTCATC
Q4VPI3	ATGGGTGTTTTACATACGAATCCGAGTTTACC-TCTATCATCCCCCTGCTA-GGTTGTTCAATGCC-TTGTTCTTGATGCTGAC-AACCTCATC

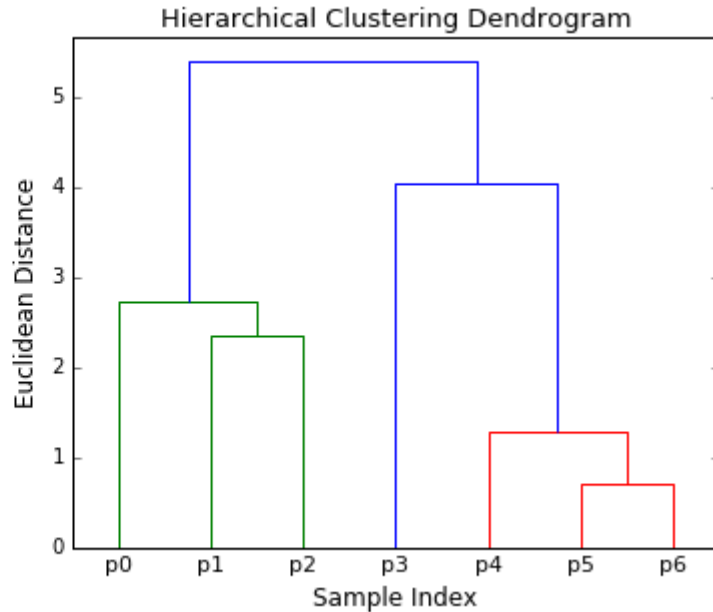
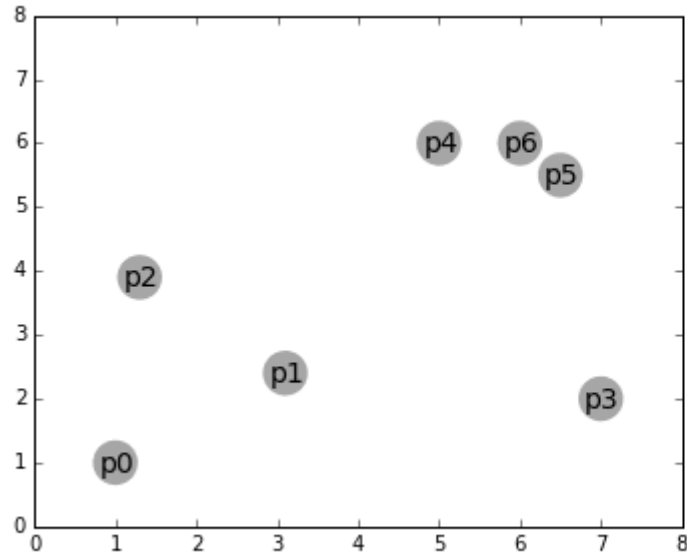
## B. Back-translation from protein alignment

Q9FPK4	ATGGGTGTTTTCAGCTACGAGGATGAGGCCACC-TCCGTTATCCTCCCGCTAGGCTGTTT---	AAGTCC-TTGTCCTAGATGCCGACAACCTCATT
Q9FPK3	ATGGGTGTTTTCAGCTACGAGGATGAGGCCACC-TCCGTTATCCTCCCGCTAGGCTGTTT---	AAGTCC-TTGTCCTAGATGCCGACAACCTCATT
Q945E7	ATGGGTGTTGTGAGTTATGAGTTTGAGGTAACC-TCCCAATTGCTCCAGCCAAGGCTTTT---	AAGGCTTTTGTTCTTGAGGCTGCCAAGATTGG
Q6XC94	ATGGGTGTTGTGAGTTATGAGTTTGAGGTAACC-TCCCAATTGCTCCAGCCAAGGCTTTT---	AAGGCTTTTGTTCTTGAGGCTGCCAAGATTGG
Q6Q4B5	ATGGGTGTTGTGAGTTATGAGTTTGAGGTAACC-TCCCAATTGCTCCAGCCAAGGCTTTTCAAGGCTTTTGTTCTTGAGGCTGCCAAGGTTGG	
Q43549	ATGGGTGTTTTCAATTACGAACTGAGTTTACC-TCCGTCATTCCCCCTGCTAGGTTGTTT---	AATGCC-TTGTTCTTGATGCTGACAACCTCATC
Q4VPJ1	ATGGGTGTTTTACATACGAATCTGAGTCCACC-TCCGTCATCCCCCTGCTAGGTTGTTT---	AATGCCACTGCTCTTGATGGTGACAAACTCATC
Q84LA7	ATGGGTGCTTTCACATACGAATCCGAA-TTACC-TCCGTCATCCCCCTGCTAGGTTGTTT---	AATGCC-TTGTTCTTGATGCTGACAACCTCATC
Q4VPI3	ATGGGTGTTTTACATACGAATCCGAGTTTACC-TCTATCATCCCCCTGCTAGGTTGTTT---	AATGCC-TTGTTCTTGATGCTGACAACCTCATC



# Phylogenetic tree construction

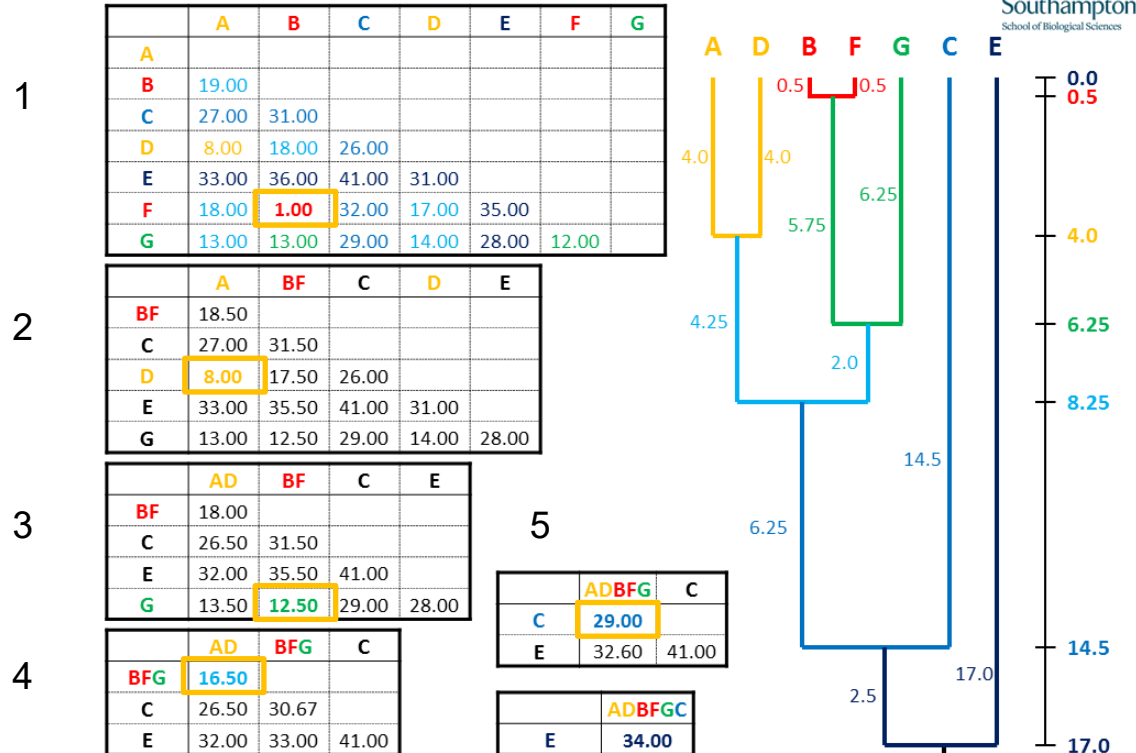
# Hierarchical/agglomerative clustering



<https://dashee87.github.io/data%20science/general/Clustering-with-Scikit-with-GIFs/>

- **Iterative** grouping of the **most similar** samples/groups until all are connected

# UPGMA: A hierarchical clustering approach



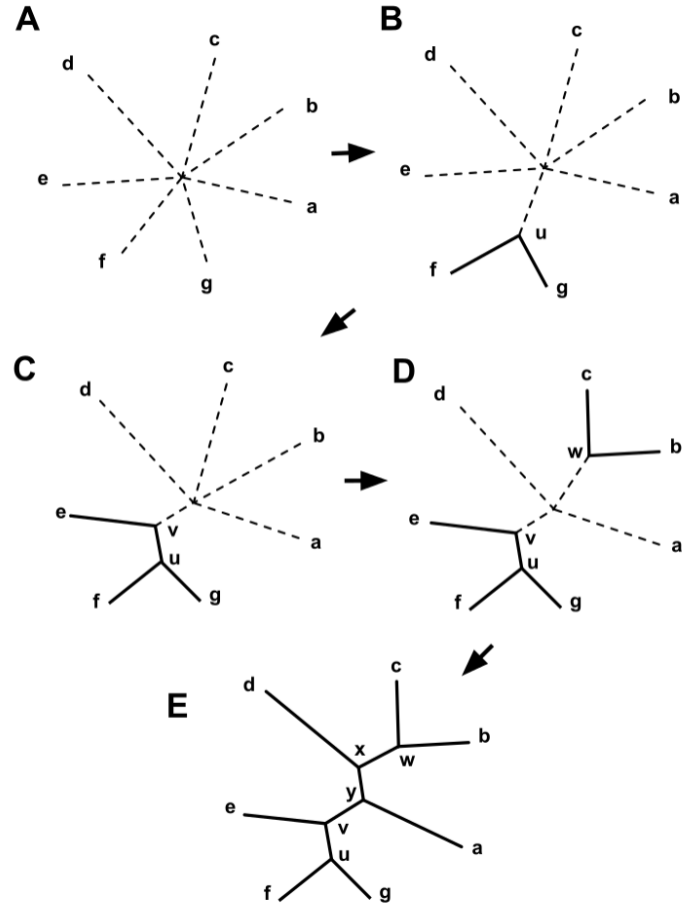
- Distance between groups of taxa to average distance between all pairs
- Calculated based on a selected **substitution model**

# NJ: Neighbor joining

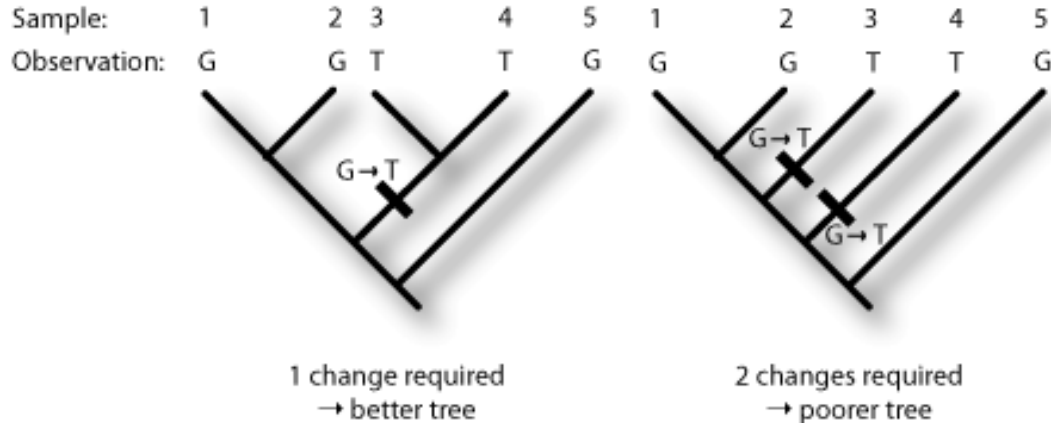
- Start with a star tree
- Join taxa  $i, j$  with smallest Q score

$$Q(i, j) = (n - 2)d(i, j) - \sum_{k=1}^n [d(i, k) + d(j, k)]$$

- Join taxa that are similar to each other but dissimilar from other taxa
- $d(i, j)$  is based on a substitution model



# Maximum parsimony / Minimum evolution



<https://biology.stackexchange.com/questions/60161/which-nucleotide-as-start-point-for-maximum-parsimony>

- Explanation requiring **minimum number of changes** is the most likely
- Fail under **long evolutionary time**:  $A \rightarrow T \rightarrow A$
- Cannot distinguish **convergent evolution**



# Maximum likelihood principle

## Let's review some terminology



- $P(A)$  = probability that A occurs

**Joint probability:**  $P(A, B)$  = probability that A and B occurs

**Conditional probability:**

- $P(A | B)$  = probability that A occurs given that B already occurred
- $P(A | B) = P(A, B) / P(B)$
- $P(A, B) = P(A | B) P(B) = P(B | A) P(A)$
- $P(\text{Shopping} | \text{Sunny}) = ?$  How about  $P(\text{Sunny} | \text{Shopping})$ ?
  - Which one is more sensible?



## Bayes' rule



$$P(A \mid B) / P(A) = P(A, B) = P(B \mid A) / P(B)$$

- A = Evolutionary data is observed
- B = Proposed phylogenetics is true
- We want to identify phylogenetics that maximize  $P(B \mid A)$
- Using Bayes' rule:
  - $P(B \mid A) = P(A \mid B) \times P(B) / P(A)$
  - $P(A \mid B)$  is much easier (and more sensible) to compute!
  - Prior belief,  $P(B)$ , can be integrated with likelihood to get a better estimate

# Maximum likelihood principle



- Likelihood score =  $P(\text{data} \mid \text{model})$ , or  $P(\text{data} \mid \text{hypothesis})$
- Find the **model** or **hypothesis** that maximize the likelihood score
- Become an **optimization problem!**
  - How to identify such **model** or **hypothesis** if there are many of them?
  - **Calculus**
  - **Brute-force search**

## Maximum likelihood answer is often intuitive



- Gene A has two alleles, A and *a*
- A study of 1,000 Thai individuals found 700 with genotype AA, 200 with genotype Aa, and 100 with genotype aa
- What is the estimated allele frequency of A?
- **Total allele counts**
  - A:  $2 \times 700 + 200 = 1,600$
  - *a*:  $2 \times 100 + 200 = 400$
  - Frequency of A =  $1,600 / 2,000 = 0.8$

## Maximum likelihood answer is often intuitive



- In a study of 5 pancreatic cancer patients, they passed away after 1, 5, 3, 4, and 5 years, respectively
- What is the estimated yearly survival rate?
- **Observation:** SD, SSSSSD, SSDD, SSSSD, SSSSSD
  - 18 *S*'s and 5 *D*'s
  - 23 years total
  - Probability of *S* =  $18/23$

# MLE: Maximum likelihood estimator

- Gene A has two alleles, A and  $a$
- A study of 1,000 Thai individuals found 700 with genotype AA, 200 with genotype Aa, and 100 with genotype aa
- What is the estimated allele frequency of A?
  - Parametrize the allele frequencies:  $f_A = p$  and  $f_a = 1 - p$
  - $P(AA) = p^2$ ,  $P(Aa) = 2p(1 - p)$ , and  $P(aa) = (1 - p)^2$
  - **Likelihood** =  $P(AA)^{700} P(Aa)^{200} P(aa)^{100} = p^{1400} 2^{200} p^{200} (1 - p)^{200} (1 - p)^{200}$   
 $= 2^{200} p^{1600} (1 - p)^{400}$
  - Which  $p$  maximize the likelihood?
    - Solve the equation  $\frac{d\text{Likelihood}}{dp} = 0 \rightarrow p_{\text{MLE}} = 0.8$

# MLE: Maximum likelihood estimator



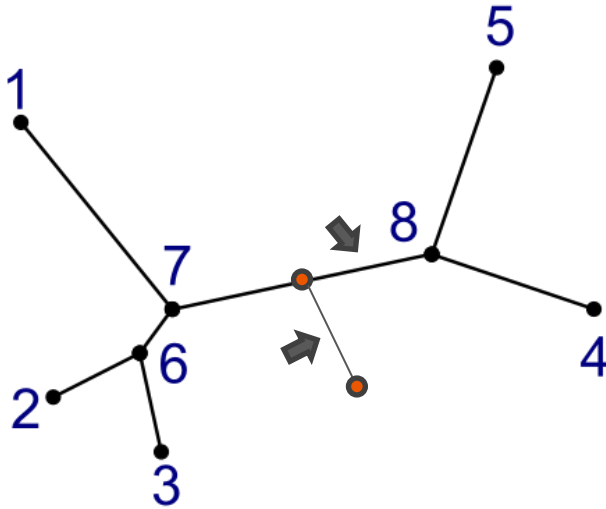
- In a study of 5 pancreatic cancer patients, they passed away after 1, 5, 3, 4, and 5 years, respectively
- What is the estimated yearly survival rate?
  - Let's set the yearly survival rate =  $r$
  - $P(\text{survive exactly } k \text{ years}) = r^k(1 - r)$
  - $P(\text{data} \mid r) = r^1(1 - r) r^5(1 - r) r^3(1 - r) r^4(1 - r) r^5(1 - r) = r^{18}(1 - r)^5$
  - Which  $r$  maximize the likelihood?
    - Solve the equation  $\frac{d\text{Likelihood}}{dr} = 0 \rightarrow r_{\text{MLE}} = 18/23$

# MLE for phylogenetics



- **Likelihood** =  $P(\text{sequence data} \mid \text{substitution model, phylogenetics tree})$
- We have a limited number of **substitution models**
- Given a fixed phylogenetic tree topology, we can alter the lengths of the branches to find the best answer:
- The problem is how to find the best **tree topology**

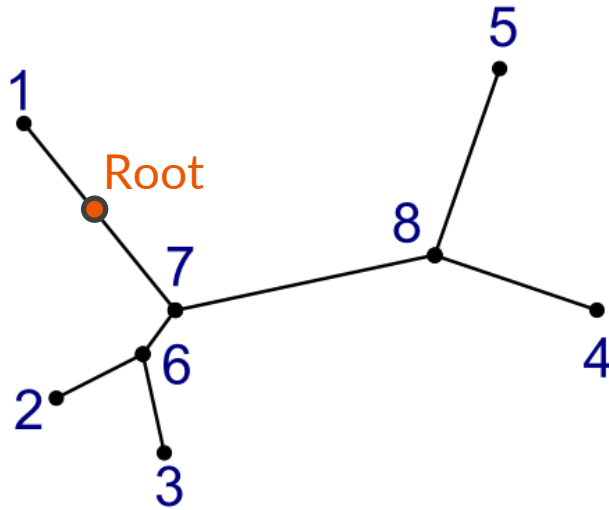
# Number of possible trees with distinct leaves



- Unrooted tree with  $n$  leaves has  $2n-2$  nodes
  - Adding a leaf will also create one internal node
- There are  $2n-3$  branches
  - Adding a leaf will create two branches
  - One new branch and split an existing branch
- From each unrooted tree with  $n$  leaves, there are  $2n-3$  locations to attach a new leaf
- Number of unrooted trees with  $n$  leaves, for  $n > 2$ , is  $(2n-5) \times (2n-7) \times (2n-9) \times \dots \times 1$

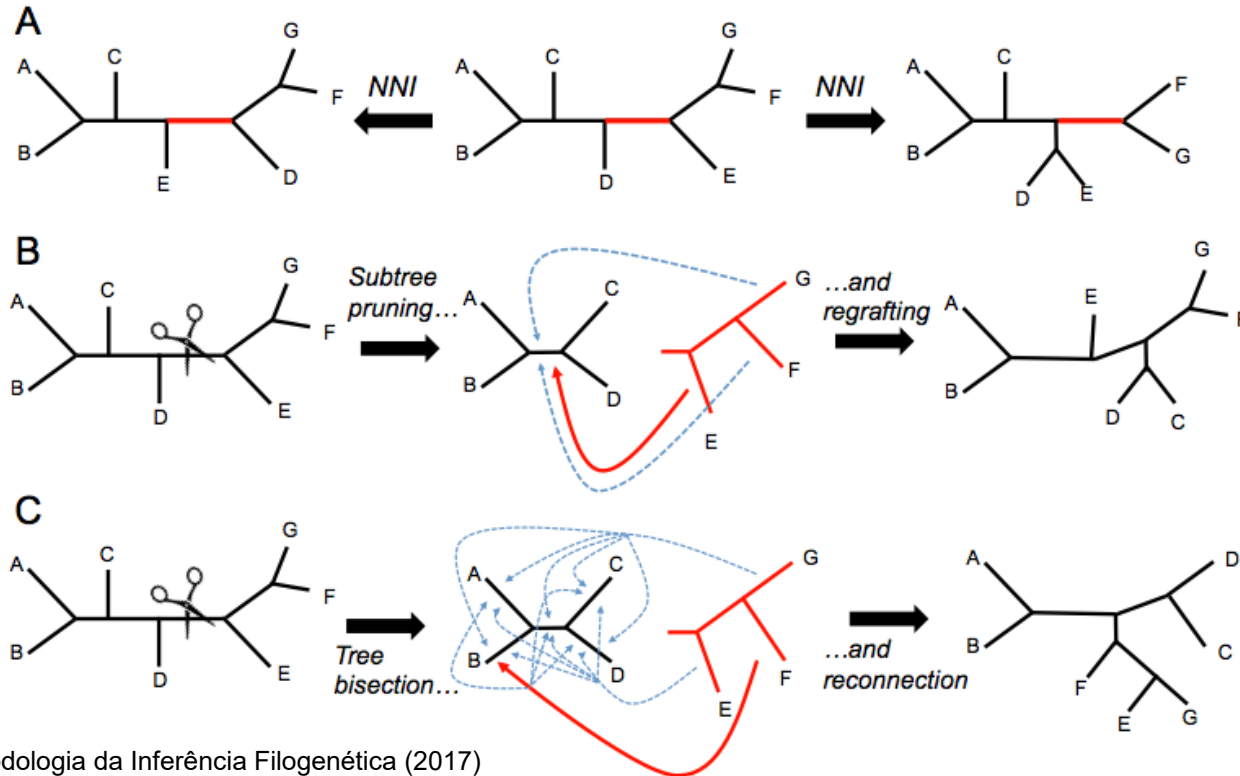


# Number of possible trees with distinct leaves



- For each unrooted tree with  $n$  leaves, there are  $2n-3$  locations to designate as the root (common ancestor)
- Number of rooted trees with  $n$  leaves, for  $n > 2$ , is  $(2n-3) \times (2n-5) \times (2n-7) \times (2n-9) \times \dots \times 1$
- Number of rooted trees with 10 leaves =  $17 \times 15 \times 13 \times 11 \times 9 \times 7 \times 5 \times 3 = 34,459,425$

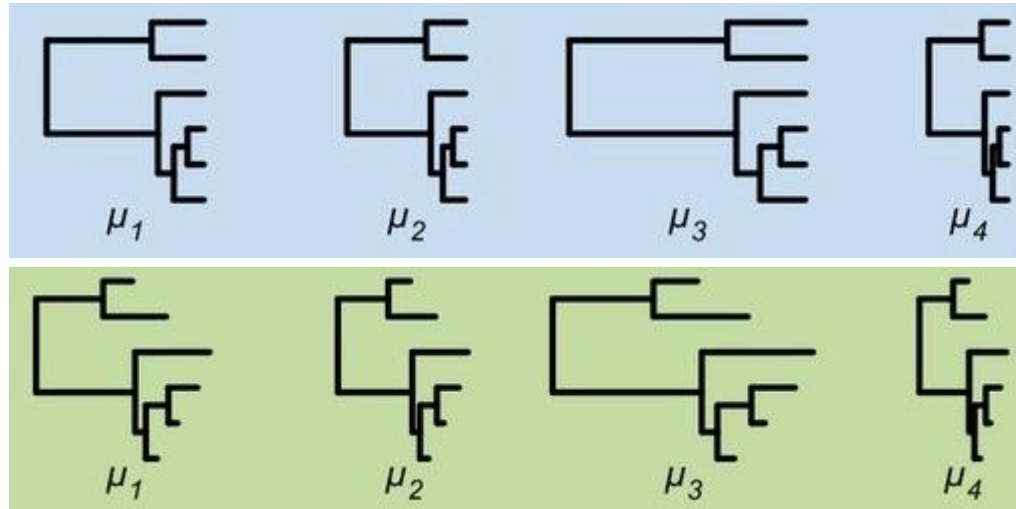
# Heuristic tree search algorithms





# **Additional evolutionary parameters**

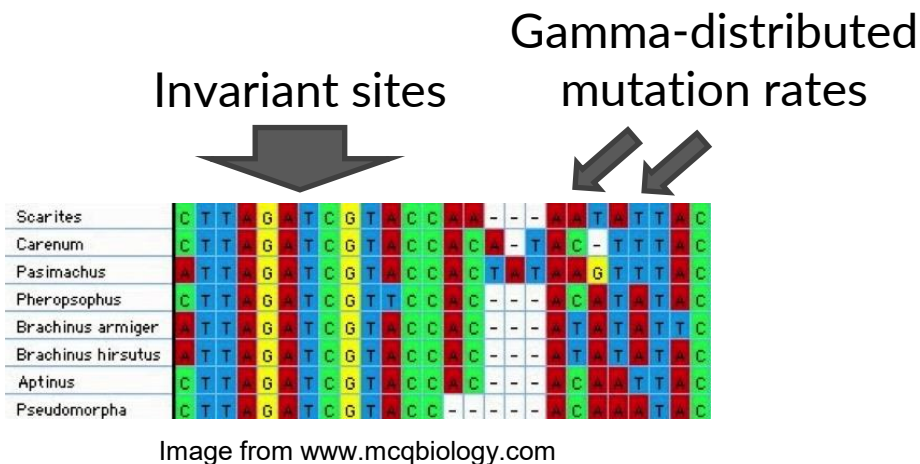
# Molecular clock assumption



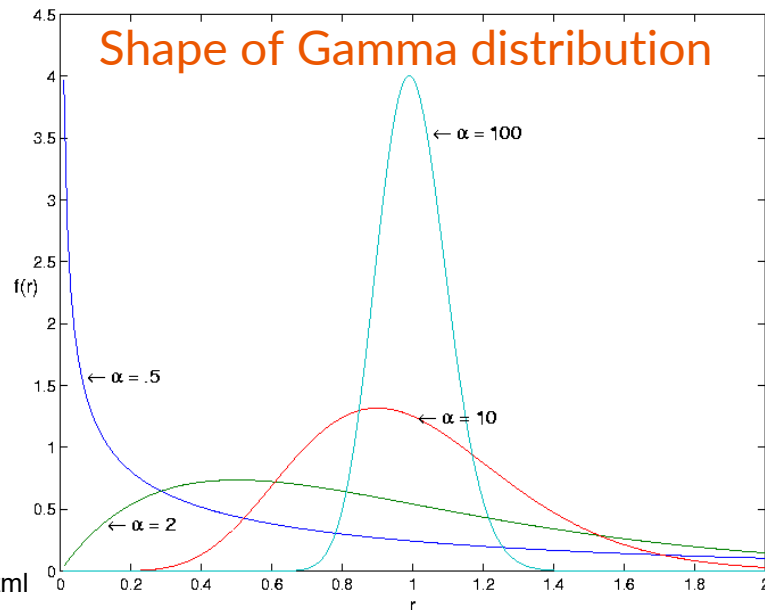
Ho, S.Y.W. and Duchene, S. Molecular Ecology 23:5947-65 (2014)

- Molecular clock assumes constant evolutionary rate throughout the tree
- Same root-to-tip distance → allow dating of evolutionary events

# Site-specific evolutionary models



<http://www.bioinf.man.ac.uk/resources/phase/manual/node81.html>



- Gamma distribution  $\Gamma(\alpha, \alpha)$  can mimic diverse shapes with mean = 1
- High  $\alpha$  results in bell shape while low  $\alpha$  yields higher variance =  $1/\alpha$

# Natural selection on codon model

Nonsynonymous / Synonymous substitution

<u>TCC</u>	GAT	<u>ATAT</u>	TGG	<u>CAAC</u>	CCC	<u>GAC</u>	AAA
S	D	I	W	Q	P	D	K

<u>TCA</u>	GAT	<u>CTAT</u>	TGG	<u>CAG</u>	CCC	<u>CAC</u>	AAA
S	D	L	W	Q	P	R	K

Luo, H. Frontiers in Microbiology 6:191 (2015)

- **Null hypothesis:** synonymous and non-synonymous occurs proportional to the number of corresponding codon positions ( $dN/dS = 1$ )
- **Alternative hypothesis:**  $dN/dS$  can differ from 1

# Nested model testing (chi-squared & likelihood)

Model complexity ↓

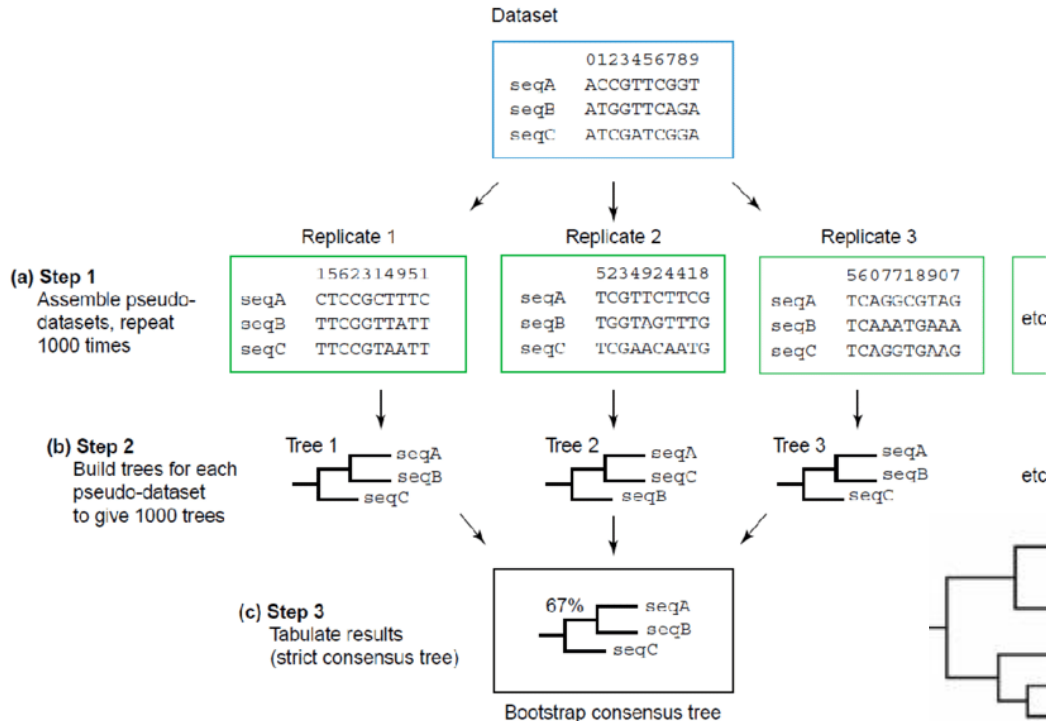
<u>Model</u>	<u>-lnL</u>	<u>models</u>	<u>diff. DF = q</u>	<u><math>\chi^2</math></u>	<u>P</u>	
JC69	3585.54820	JC-F81	3 - 0 = 3	155	0	F81 model fits the data significantly better than JC
F81	3508.04085					
HKY85	3233.34395	F81-HKY85	4 - 3 = 1	549.4	0	HKY85 model fits the data significantly better than F81
TrN93	3232.29439					
		KHY-TrN	5 - 4 = 1	2.1	0.15	TrN model does not fit the data significantly better than HKY85

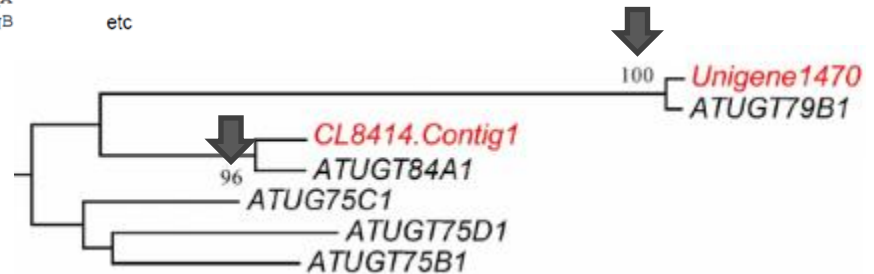
<u>Model</u>	<u>-lnL</u>	<u>models</u>	<u>diff. DF = q</u>	<u><math>\chi^2</math></u>	<u>P</u>	
HKY85	3233.34395	HKY85-vs. +G	1	176	0	Adding site-specific rate fits the data significantly better
HKY85 +G	3145.29031	HKY85+G vs. I+G	1	5.85	0.015	Adding invariant sites does not fit the data significantly better
HKY85 +I+G	3142.36439					

Model complexity ↓

# Bootstrap support for taxa group



- Bootstrapping = sampling with repeats
- Allow multiple trees to be constructed from a single sequence alignment





**AVPR1a tree**

**Consensus tree**


**B**

- Ren, D. et al. PLoS ONE 9:e222638 (2014)




# **Example of how to setup phylogenetic reconstruction**

# MEGA: A GUI tool for phylogenetic analysis



Molecular Evolutionary  
Genetics Analysis



## TUTORIALS

Below are links to online video lectures and tutorials for multiple versions of MEGA. The first section of videos were created by members of Dr. Sudhir Kumar's lab at the Institute for Genomics and Evolutionary Medicine (iGEM) at Temple University. The rest of the videos were produced by users of MEGA. To assemble this collection of videos, the MEGA team performed a search of YouTube for instructional MEGA videos and assembled this collection of the most popular videos found. If you would like to suggest additions to this collection, please contact us by using the [feedback page](#).

### KUMAR LAB VIDEOS

Molecular Dating with MEGA

Choosing and Acquiring Sequences Part 1

Choosing and Acquiring Sequences Part 2

Reconstructing Ancestral

Relative Rate Framework for

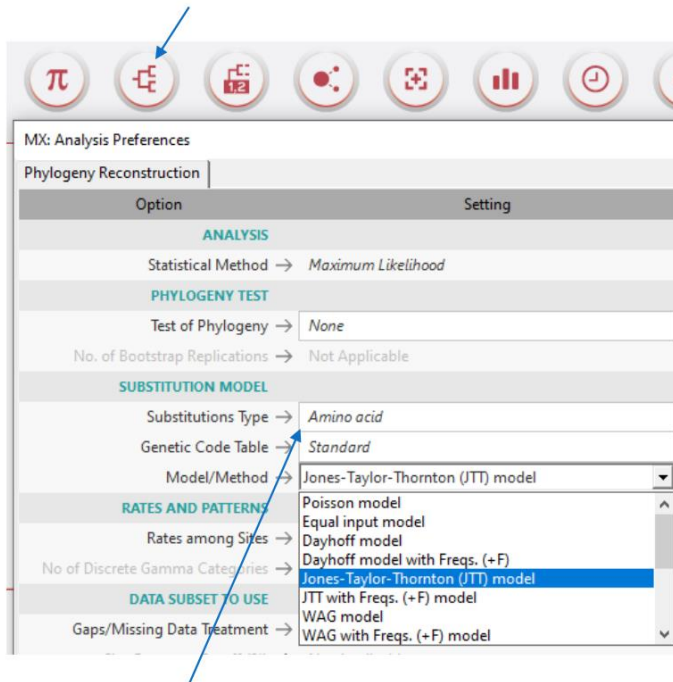
Inferring Selection with MEGA

Filter tutorials by topic:

- Align Sequences
- BLAST Search
- Bootstrap Tree
- Calibrate
- Distance Estimation
- Edit Sequences
- Edit Tree
- Evolutionary Probability
- Gene Duplication
- Grouping Taxa
- Installation
- Model Selection
- Pairwise Distances
- Phylogeny Construction
- Substitution Matrix
- Trace Files

# Substitution model choices

## PHYLOGENY mode

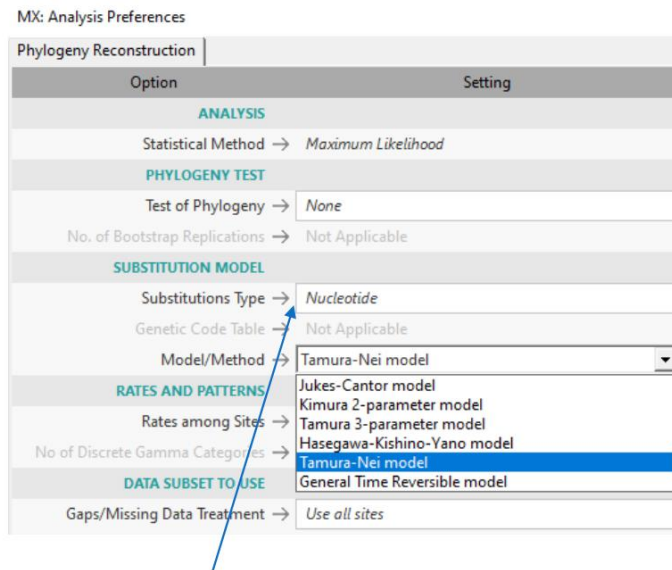


MX: Analysis Preferences

Phylogeny Reconstruction

Option	Setting
<b>ANALYSIS</b>	
Statistical Method →	Maximum Likelihood
<b>PHYLOGENY TEST</b>	
Test of Phylogeny →	None
No. of Bootstrap Replications →	Not Applicable
<b>SUBSTITUTION MODEL</b>	
Substitutions Type →	Amino acid
Genetic Code Table →	Standard
Model/Method →	Jones-Taylor-Thornton (JTT) model
<b>RATES AND PATTERNS</b>	
Rates among Sites →	Dayhoff model
No of Discrete Gamma Categories →	Jones-Taylor-Thornton (JTT) model
<b>DATA SUBSET TO USE</b>	
Gaps/Missing Data Treatment →	WAG model

Amino acid-based models



MX: Analysis Preferences

Phylogeny Reconstruction

Option	Setting
<b>ANALYSIS</b>	
Statistical Method →	Maximum Likelihood
<b>PHYLOGENY TEST</b>	
Test of Phylogeny →	None
No. of Bootstrap Replications →	Not Applicable
<b>SUBSTITUTION MODEL</b>	
Substitutions Type →	Nucleotide
Genetic Code Table →	Not Applicable
Model/Method →	Tamura-Nei model
<b>RATES AND PATTERNS</b>	
Rates among Sites →	Jukes-Cantor model
No of Discrete Gamma Categories →	Tamura-Nei model
<b>DATA SUBSET TO USE</b>	
Gaps/Missing Data Treatment →	Use all sites

Nucleotide-based models

# Maximum likelihood parameters

Phylogeny Reconstruction	
Option	Setting
ANALYSIS	
Statistical Method →	Maximum Likelihood
PHYLOGENY TEST	
Test of Phylogeny →	None
No. of Bootstrap Replications →	None
	Bootstrap method
SUBSTITUTION MODEL	
Substitutions Type →	Nucleotide
Model/Method →	Tamura-Nei model
RATES AND PATTERNS	
Rates among Sites →	Uniform Rates
No of Discrete Gamma Categories →	Not Applicable
DATA SUBSET TO USE	
Gaps/Missing Data Treatment →	Use all sites
Site Coverage Cutoff (%) →	Not Applicable
TREE INFERENCE OPTIONS	
ML Heuristic Method →	Nearest-Neighbor-Interchange (NNI)
Initial Tree for ML →	Make initial tree automatically (Default - NJ/BioNJ)
Initial Tree File →	Not Applicable
Branch Swap Filter →	None
SYSTEM RESOURCE USAGE	
Number of Threads →	7

Bootstrapping

Evolutionary / substitution model

Site-specific evolutionary rate

How to handle gap (indel) and missing data

Detailed tree building method

Parallel processing

# Maximum likelihood parameters

RATES AND PATTERNS	
Rates among Sites →	Uniform Rates
No of Discrete Gamma Categories →	Uniform Rates
	Gamma Distributed (G)
	Has Invariant Sites (I)
	Gamma Distributed With Invariant Sites (G+I)
Gaps/Missing Data Treatment →	Use all sites

Invariant site and Gamma-distributed site-specific substitution rates

TREE INFERENCE OPTIONS	
ML Heuristic Method →	Nearest-Neighbor-Interchange (NNI)
	Nearest-Neighbor-Interchange (NNI)
	Subtree-Pruning-Regrafting - Fast (SPR level 3)
	Subtree-Pruning-Regrafting - Extensive (SPR level 5)
Initial Tree for ML →	Make initial tree automatically (Default - NJ/BioNJ)
Initial Tree File →	Make initial tree automatically (Default - NJ/BioNJ)
	Make initial tree automatically (Maximum Parsimony)
Branch Swap Filter →	Make initial tree automatically (Neighbor Joining)
	Make initial tree automatically (BioNJ)
	Use tree from file
	Use Topology Editor

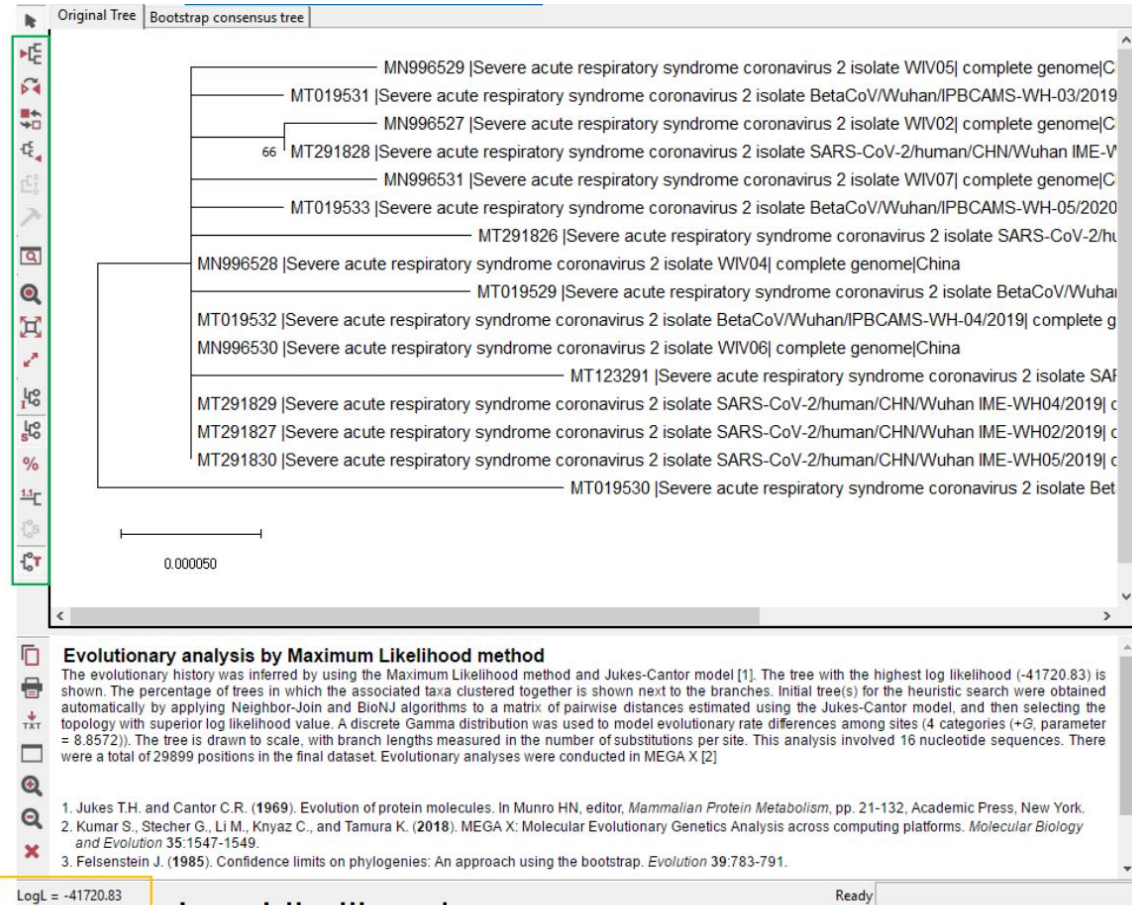
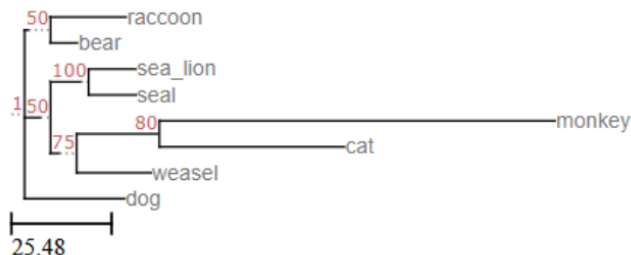
Tree search method

Initial tree can be built using quick and simple method like neighbor joining

# Output

## NEWICK format

```
((racoon:19.19959,bear:6.80041)50:0.84  
600,((sea_lion:11.99700,seal:12.00300)  
100:7.52973,((monkey:100.85930,cat:47.  
14069)80:20.59201,weasel:18.87953)75:  
2.09460)50:3.87382,dog:25.46154);
```



Log-Likelihood score

## Other useful tools



- RAxML, PAUP are standard phylogenetic reconstruction tools
- PAML specializes in natural selection analysis (dN/dS)
- FastTree specializes in speed for large dataset
- MrBayes, BEAST use Bayesian approach
  - $P(\text{tree topology, branch lengths} \mid \text{substitution model, sequence data})$
- More: <https://evolution.genetics.washington.edu/phylip/software.html#methods>



# Any question?



- See you next time