Dr. Gregory A. Babbitt



**BabbittLab@RIT**

COMPUTATIONAL MOLECULAR EVOLUTION AND ECOLOGY

some of our students



Erin Coppola BME

Jamie Mortensen BME

Justin Liao BME

We are a mostly "dry-bench" or computer-based laboratory studying the evolution of the molecular components of cells. Biophysics and biochemistry underlies all molecular processes in the cell, but only some molecular structure and behavior, genetic and epigenetic, can contain information that is subject to heritable change over potentially deep timescales.

Mohammed Alawad
Bioinformatics

Katharina Schulze
Bioinformatics

Our long-term research goals are aimed at a more biophysically-grounded understanding of molecular evolution in the cell. We combine biophysical and molecular evolutionary modeling to investigate the boundary between evolvable and non-evolvable biophysics.
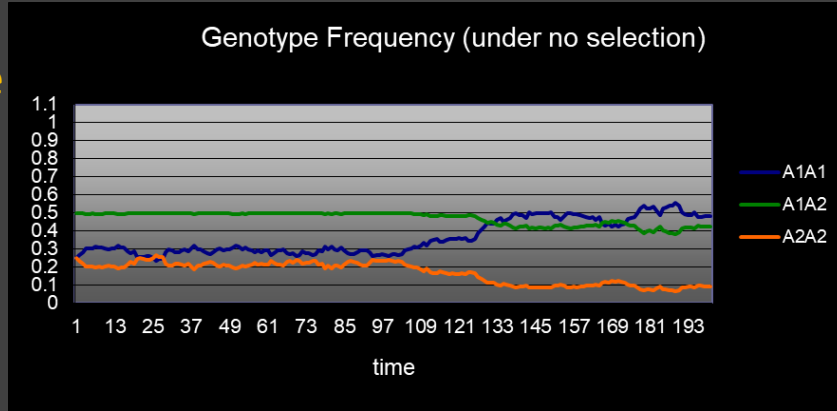
Realistic expectations about genetic variants (especially human) requires an appreciation of how evolution and migration affects molecular evolution

1. Almost all molecular variant evolution is neutral
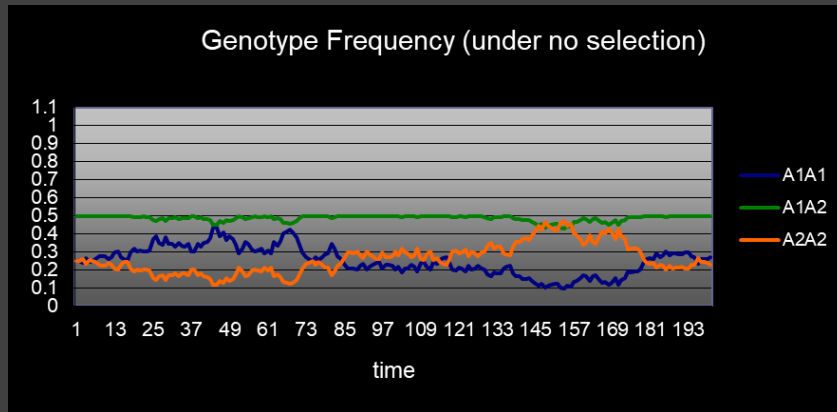
2. Small populations randomly fix traits in populations
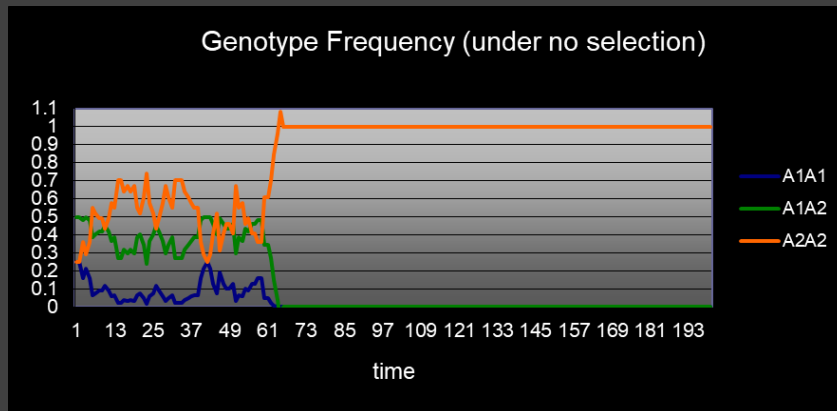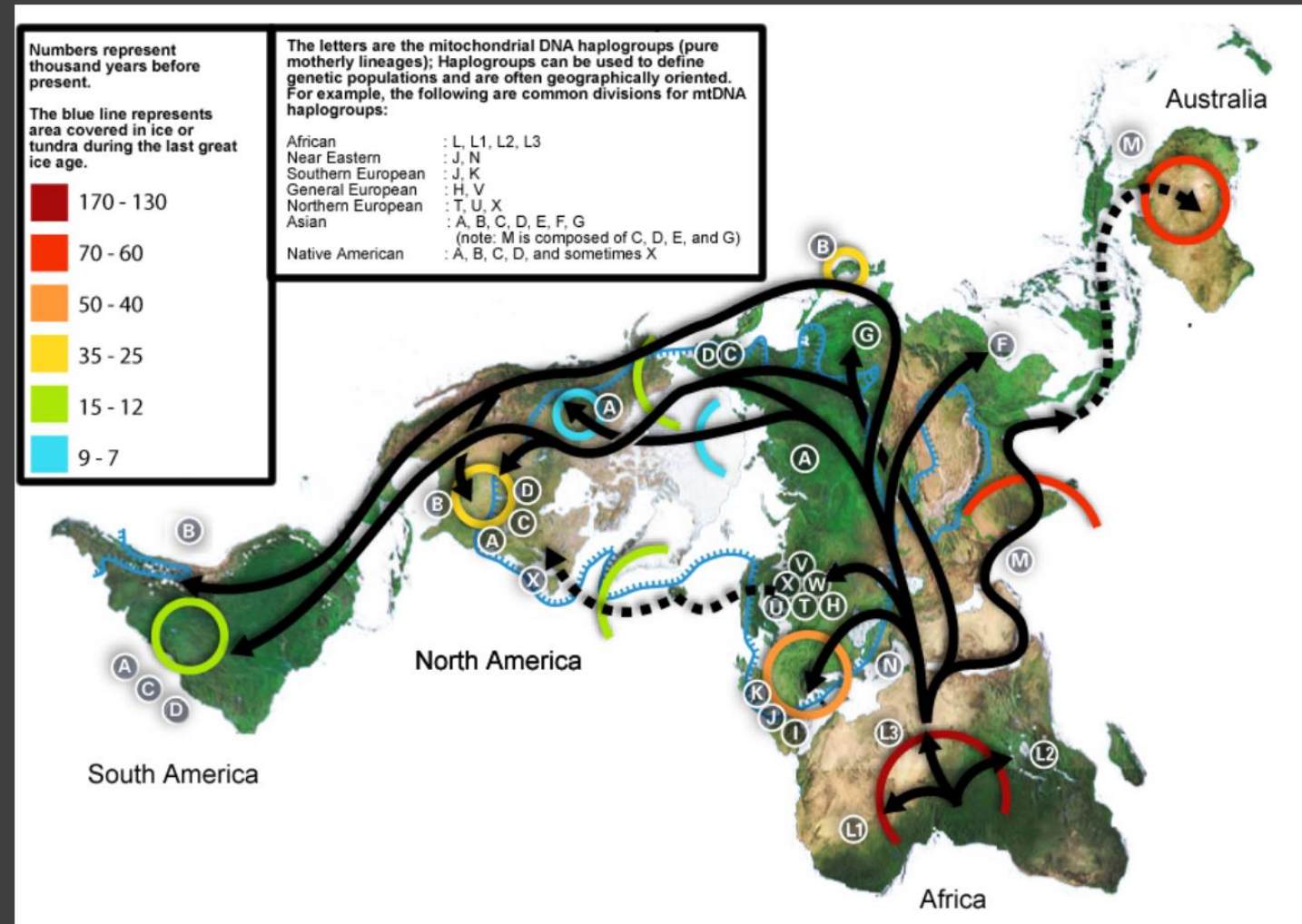
# Wright-Fisher two allele dynamic pop gen model

## Human haplotype variation

Pop size

2000

500
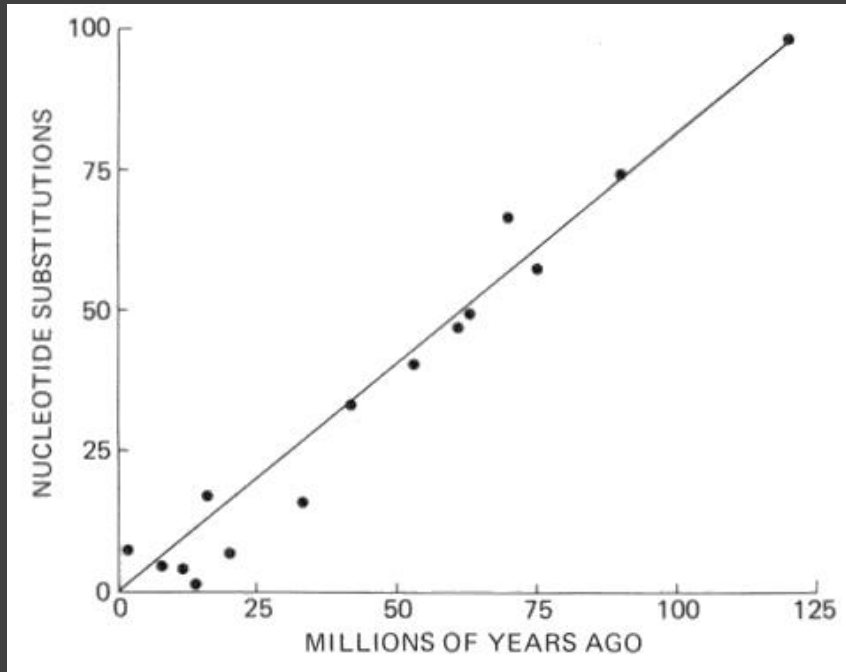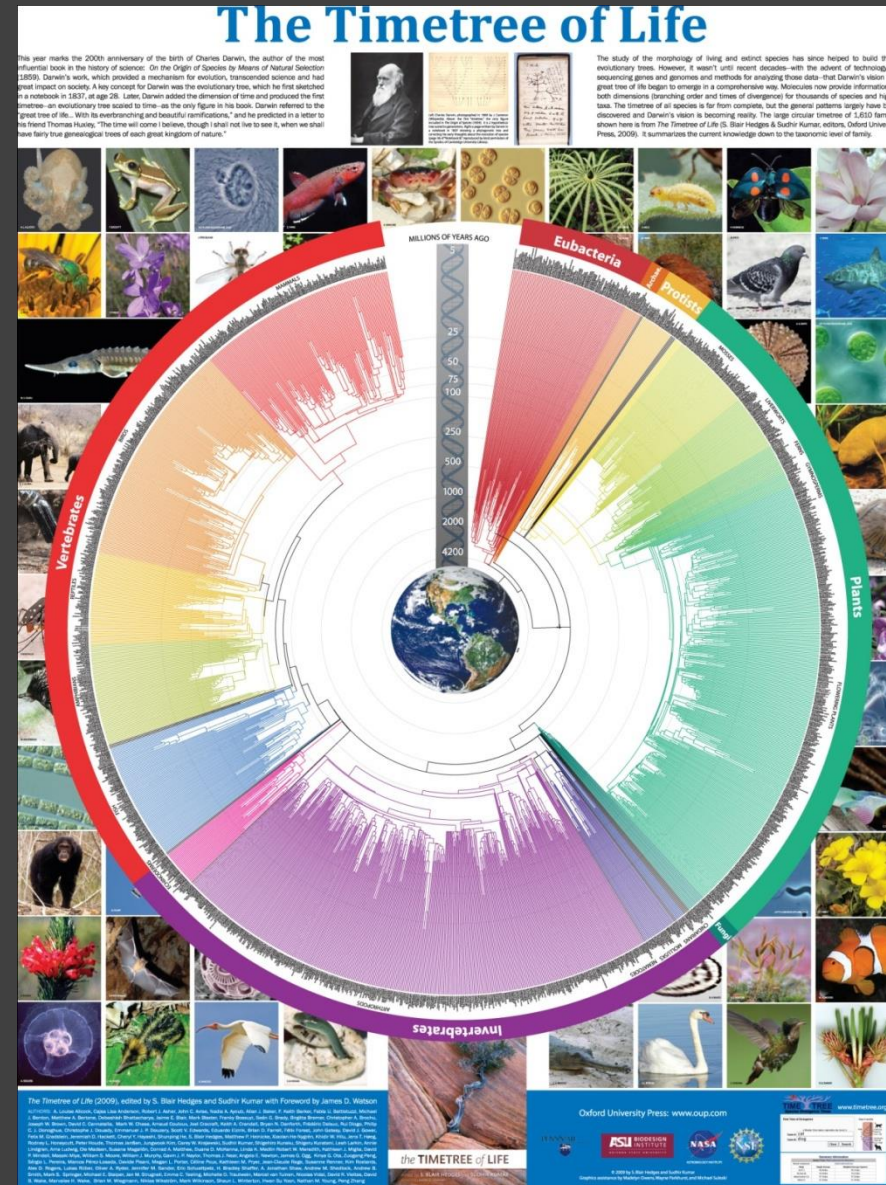
50

# Phylogenetic reconstruction from aligned sequence data





Substitutions accumulate in a "clock-like" linear fashion ….. and so divergence times can be estimated by "molecular clocks"
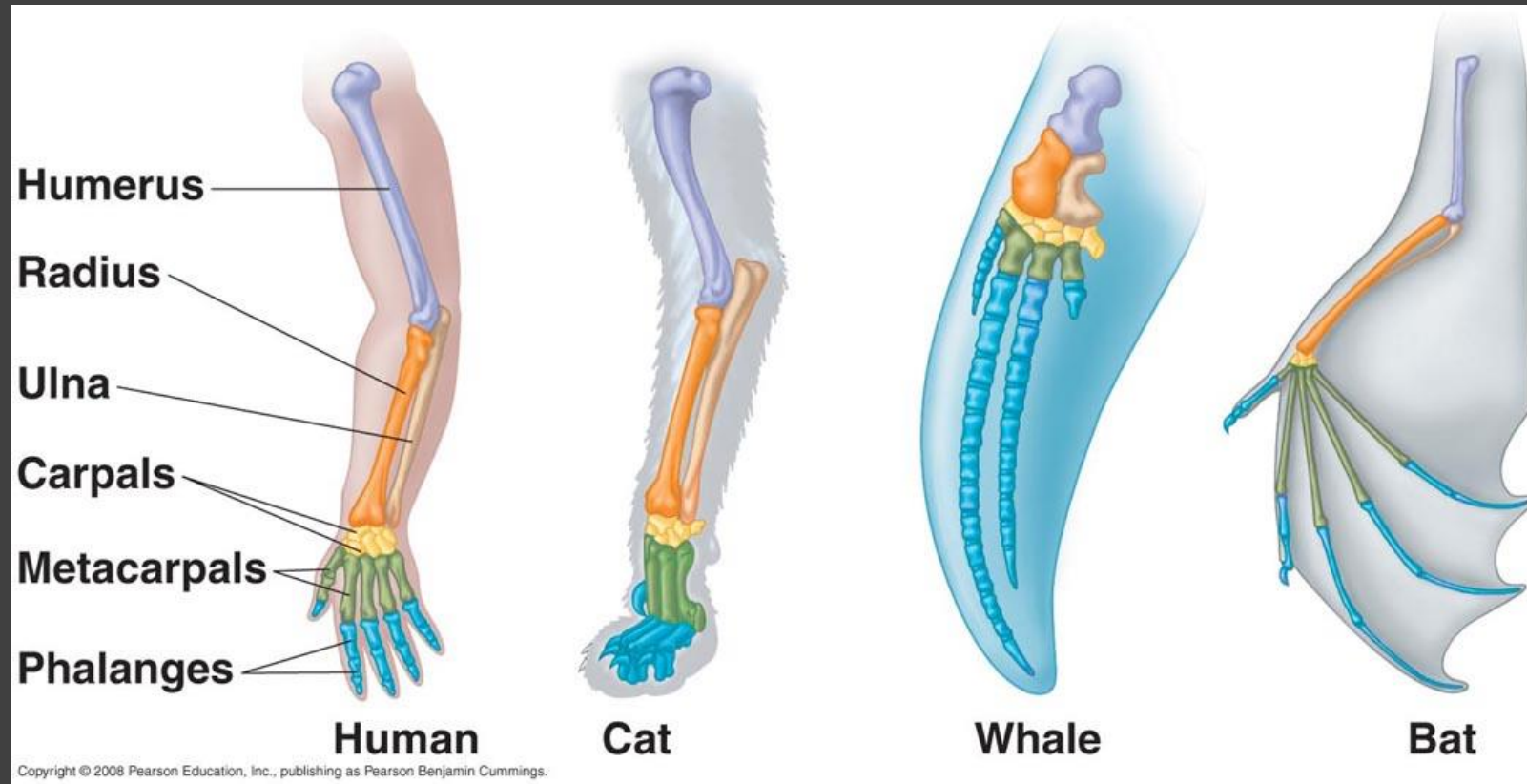
…most mutations are assumed neutral in their effects …"neutral theory of mol evol" Many are synonymous and/or fixed randomly

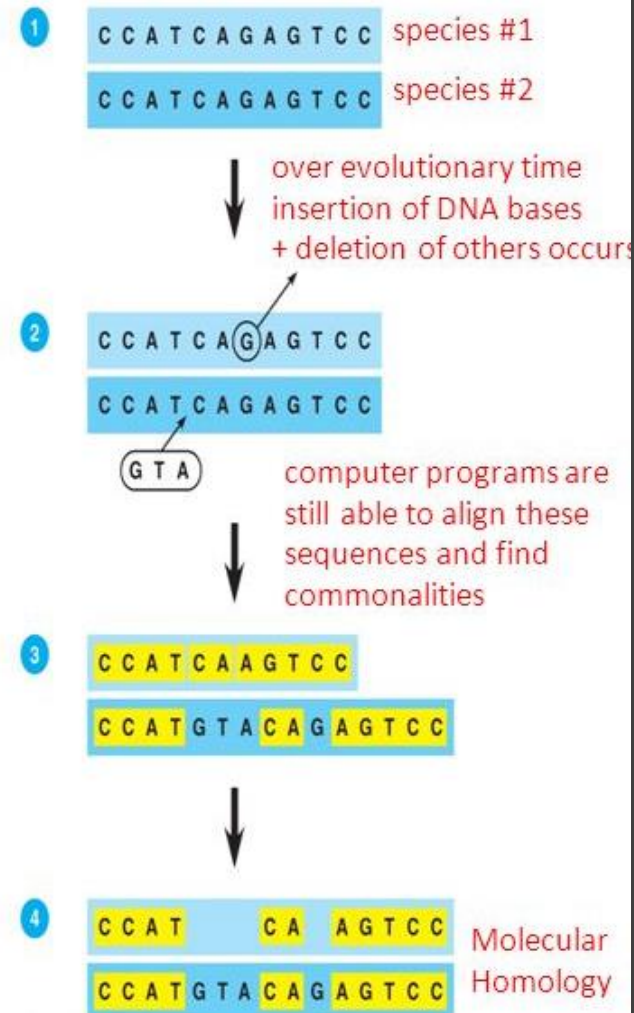Correct molecular homology is crucial for comparative genomics

1. Scrutinize all results of alignment algorithms every time you use one (visually or using mismatch filter to mask poorly aligned regions)

2. Do not mix mutation events derived via gene duplication with speciation into same phylogenetic inference (i.e. analyze orthologs and paralogs separately)

# character homology – paleontology



Human — Cat — Whale — Bat

Humerus
Radius
Ulna
Carpals
Metacarpals
Phalanges

Copyright © 2008 Pearson Education, Inc., publishing as Pearson Benjamin Cummings.
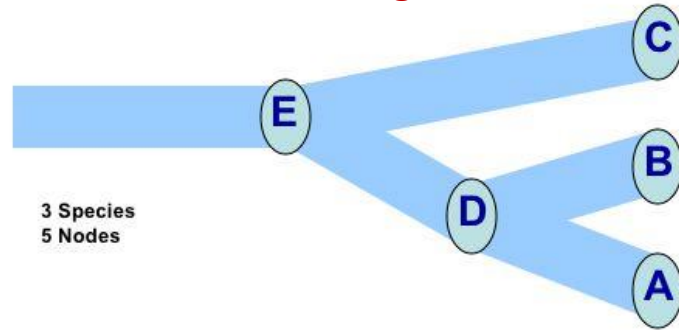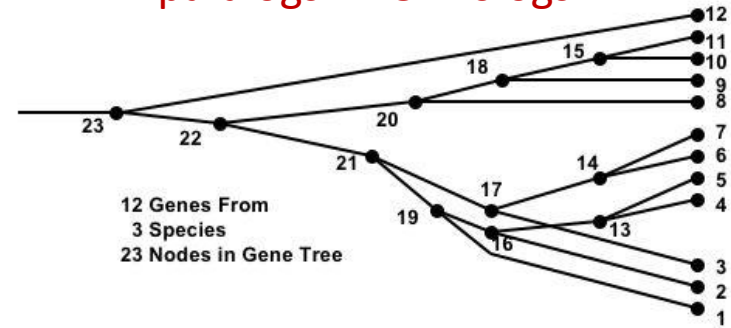
# molecular homology – sequence level

- to evaluate molecular homology requires analysis of DNA sequences
  - extract the DNA, sequence the DNA and align them in terms of similar sequences
  - alignment done by powerful computer programs that take into account deletions of bases or additions of bases that can "shift" the coding and non-coding sequences back or forward
  - also determine if the similarities are just a coincidence (molecular homoplasy or analogy)
- so looking at the DNA sequences of the Australian and N.A. moles identifies numerous differences in DNA sequences that can't be aligned
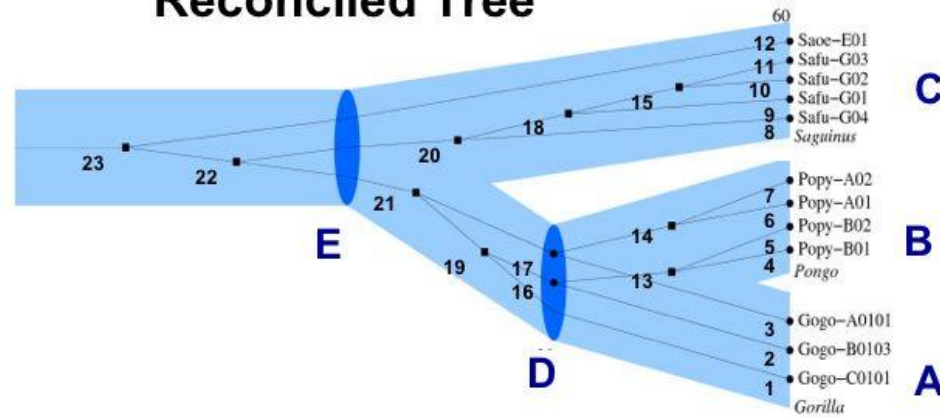  - do not share a common ancestor and their phylogenetic trees will differ



1. CCATCAGAGTCC species #1
   CCATCAGAGTCC species #2

over evolutionary time insertion of DNA bases + deletion of others occurs

2. CCATCA(G)AGTCC
   CCATCAGAGTCC
   (GTA)

computer programs are still able to align these sequences and find commonalities

3. CCATCAAGTCC
   CCATGTACAGAGTCC

4. CCAT   CA AGTCC   Molecular
   CCATGTACAGAGTCC   Homology

**Species Tree**
orthologs

C
E
B
D
A

3 Species
5 Nodes

**Gene Tree**
paralogs = "Ohnologs"

23
22
20
21
18
15
17
14
19
16
13

12
11
10
9
8
7
6
5
4
3
2
1

12 Genes From
3 Species
23 Nodes in Gene Tree

**Reconciled Tree**

60

23
22
20
21
18
15
19
17
16
14
13
E
D

12 • Saoe–E01
11 • Safu–G03
10 • Safu–G02
9 • Safu–G01
8 • Safu–G04
  *Saguinus*
7 • Popy–A02
6 • Popy–A01
5 • Popy–B02
4 • Popy–B01
  *Pongo*
3 • Gogo–A0101
2 • Gogo–B0103
1 • Gogo–C0101
  *Gorilla*

C

B

A

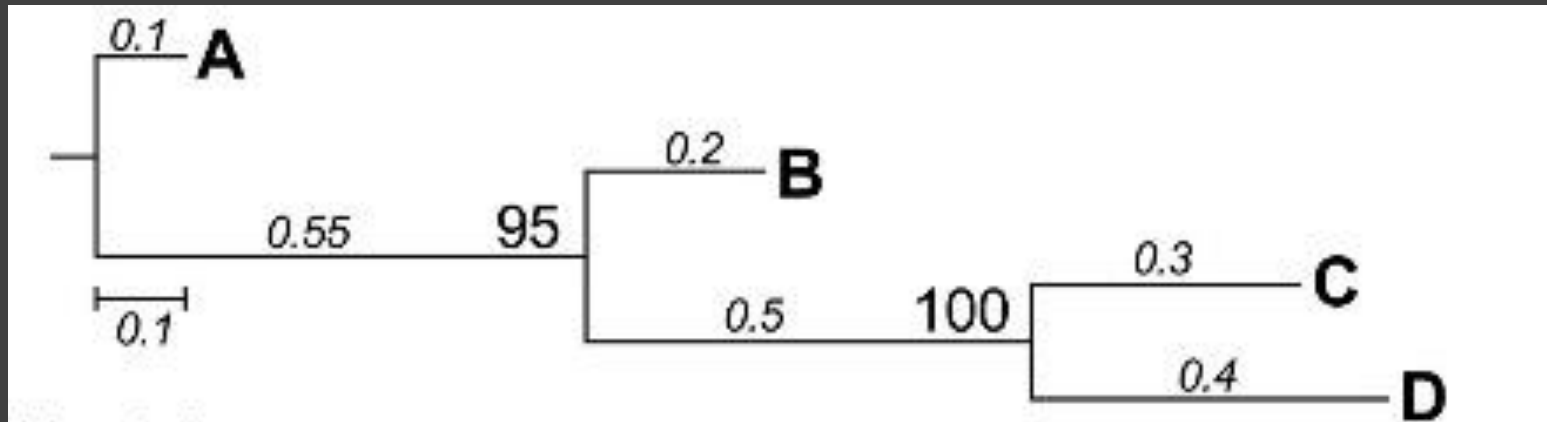# Inferring the most probable tree topology

# phylogenetic inference

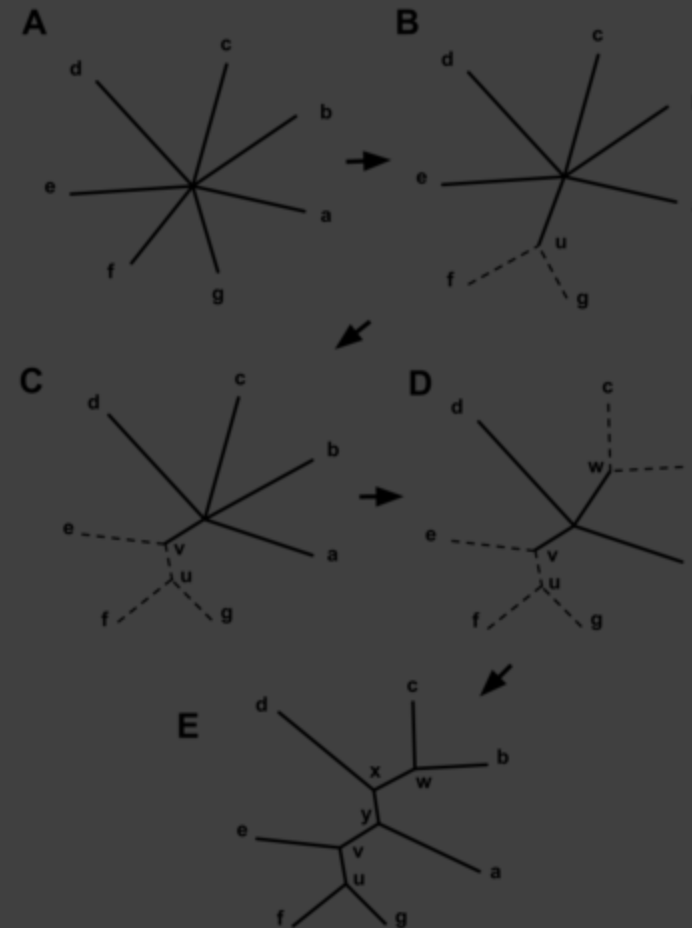# phylogenetic tree viewing

# phylogenetic tree representation – computer format



Newick:
(A:0.1,(B:0.2,(C:0.3,D:0.4)100:0.5)95:0.55);

Extended Newick (eNewick):
(A:0.1,(B:0.2,(C:0.3,D:0.4)0.5[100])0.55[95]);

Starting with a star tree (A), the Q matrix is calculated and used to choose a pair of nodes for joining, in this case f and g. These are joined to a newly created node, u, as shown in (B). The part of the tree shown as dotted lines is now fixed and will not be changed in subsequent joining steps. The distances from node u to the nodes a-e are computed from the formula given in the text. This process is then repeated, using a matrix of just the distances between the nodes, a,b,c,d,e, and u, and a Q matrix derived from it. In this case u and e are joined to the newly created v, as shown in (C). Two more iterations lead first to (D), and then to (E), at which point the algorithm is done, as the tree is fully resolved



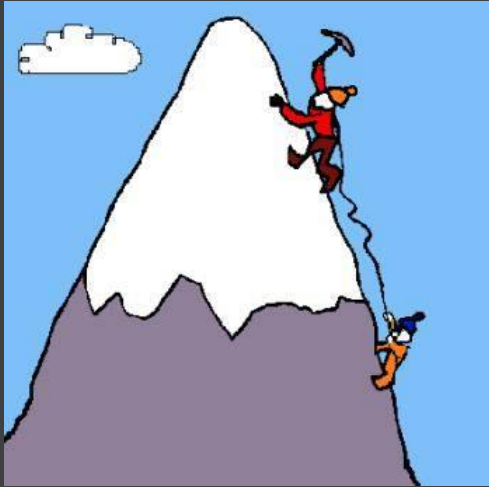Based on a distance matrix relating the r taxa, calculate Q as follows:

$$Q(i,j) = (r-2)d(i,j) - \sum_{k=1}^{r} d(i,k) - \sum_{k=1}^{r} d(j,k)$$

where   d(i,j)   is the distance between taxa   i and j   .

Note: evolutionary distances are just mutation counts adjusted for the probability of multiple hits at the same sites
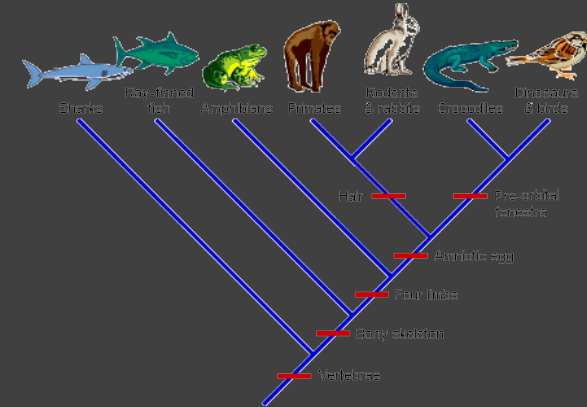
# Frequentist paradigm

**Maximum likelihood estimate**

**Maximum likelihood tree**





$$P(x) = freq\ (x \mid \theta)$$

$$P(E) = freq\ (E \mid H)$$

Likelihoods are much more useful

$$L(H) = prob\ (H \mid E)$$

$$L(model) = joint\ prob\ (model \mid \text{[x x x x x x x x x x x]})$$



Ronald Alymer Fisher

# Probability versus Likelihood

P(x) = P(occurrence of observation given particular model of known parameters)

More simply ….**P(data|model parameters)**     Is this often useful?

Do we always 'know' what model fits?

likelihood flips this relation          **L(model parameters | data)**
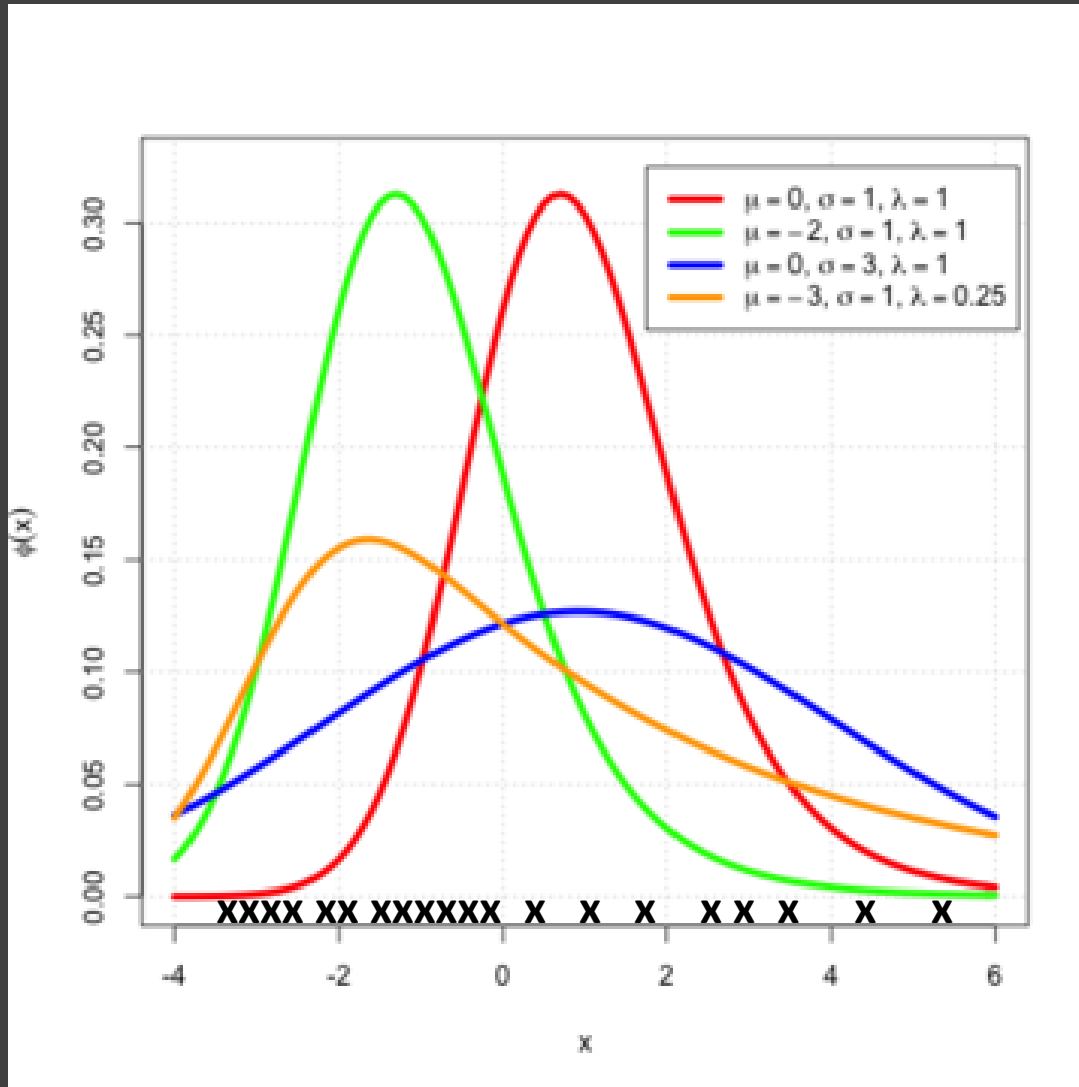
Thus L(θ/x)  = P(x/θ)

**Now we can think of the likelihood of a model given a data point**
**….or better many data points**

Thus L(θ; $x_1 \dots x_n$) = $\prod_{i=1}^{n} f(x_i|\theta) = \sum_{i=1}^{n} ln f(x_1|\theta)$



Ronald Alymer Fisher

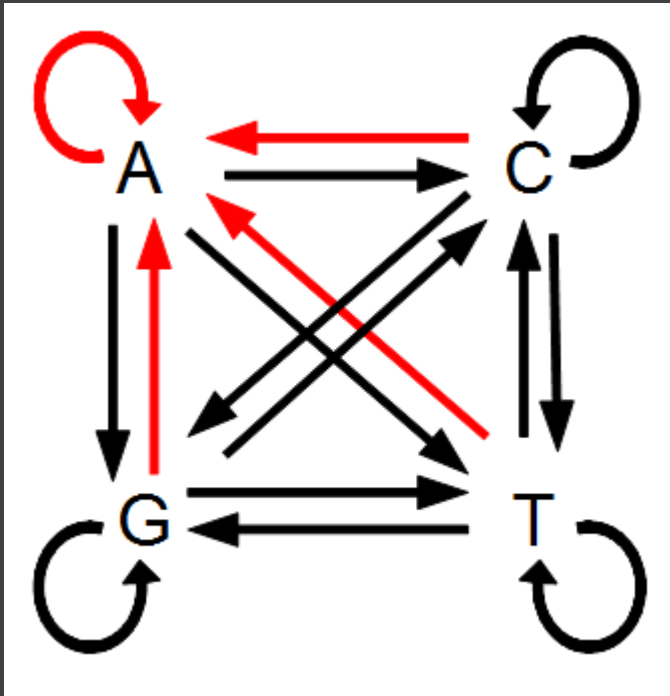# maximum likelihood estimation and multimodel inference



Ronald Alymer Fisher

What is conditional probability is generally better to know?

P(model|data) or P(data|model)

**Which density curve is most likely given the data below?**

# Underlying models of evolution (DNA substitution vs codon transitions)

Simplest model (Jukes Cantor JC69) all substitution
rates are assumed of equal potential probability
and related to overall mutation rate (i.e. one parameter model)



$$Q = \begin{pmatrix} -\mu_A & \mu_{GA} & \mu_{CA} & \mu_{TA} \\ \mu_{AG} & -\mu_G & \mu_{CG} & \mu_{TG} \\ \mu_{AC} & \mu_{GC} & -\mu_C & \mu_{TC} \\ \mu_{AT} & \mu_{GT} & \mu_{CT} & -\mu_T \end{pmatrix}$$
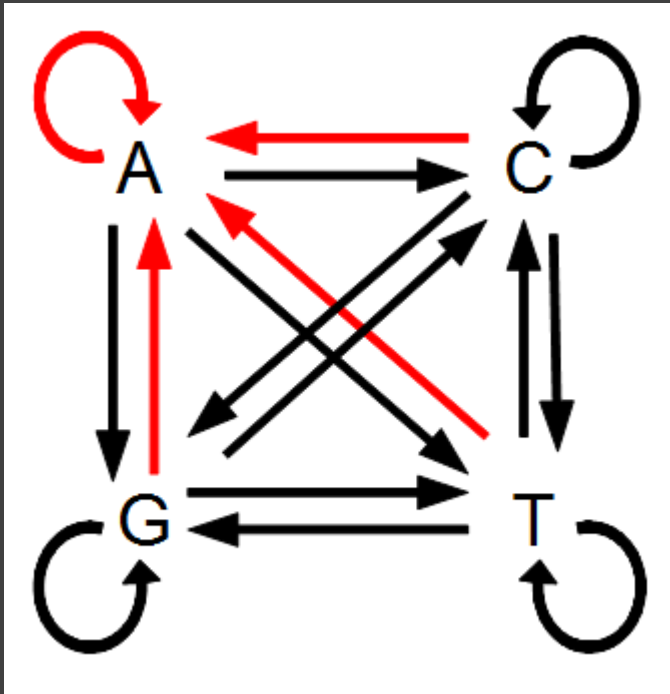
Subs rate matrix

$$P(t) = \begin{pmatrix} p_{AA}(t) & p_{GA}(t) & p_{CA}(t) & p_{TA}(t) \\ p_{AG}(t) & p_{GG}(t) & p_{CG}(t) & p_{TG}(t) \\ p_{AC}(t) & p_{GC}(t) & p_{CC}(t) & p_{TC}(t) \\ p_{AT}(t) & p_{GT}(t) & p_{CT}(t) & p_{TT}(t) \end{pmatrix}$$

probability matrix of observing subs – recall adjust for multiple hits

# Underlying models of evolution (DNA substitution vs codon transitions)

Simplest model (Jukes Cantor JC69) all substitution
rates are assumed of equal potential probability
and related to overall mutation rate (i.e. one parameter model)

$$Q = \begin{pmatrix} * & \frac{\mu}{4} & \frac{\mu}{4} & \frac{\mu}{4} \\ \frac{\mu}{4} & * & \frac{\mu}{4} & \frac{\mu}{4} \\ \frac{\mu}{4} & \frac{\mu}{4} & * & \frac{\mu}{4} \\ \frac{\mu}{4} & \frac{\mu}{4} & \frac{\mu}{4} & * \end{pmatrix}$$

Subs rate matrix



Given the proportion $p$ of sites that differ between the two sequences the Jukes-Cantor estimate of the evolutionary distance (in terms of the expected number of changes) between two sequences is given by
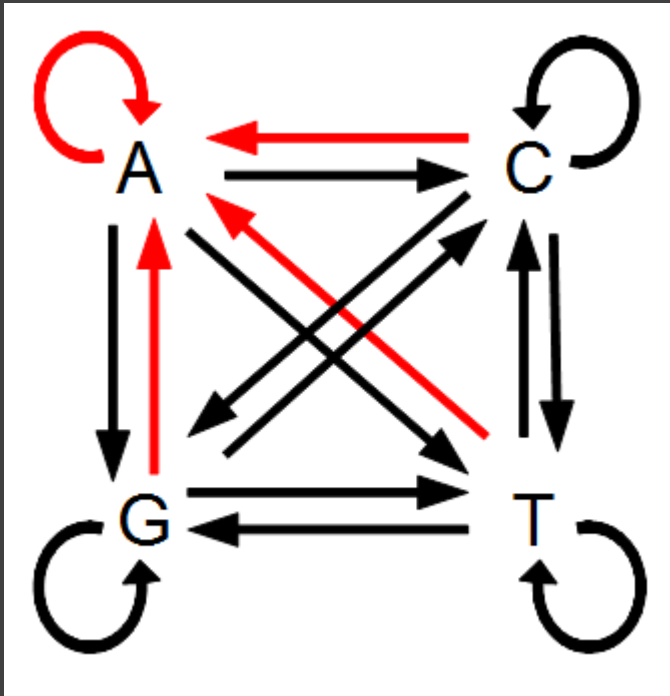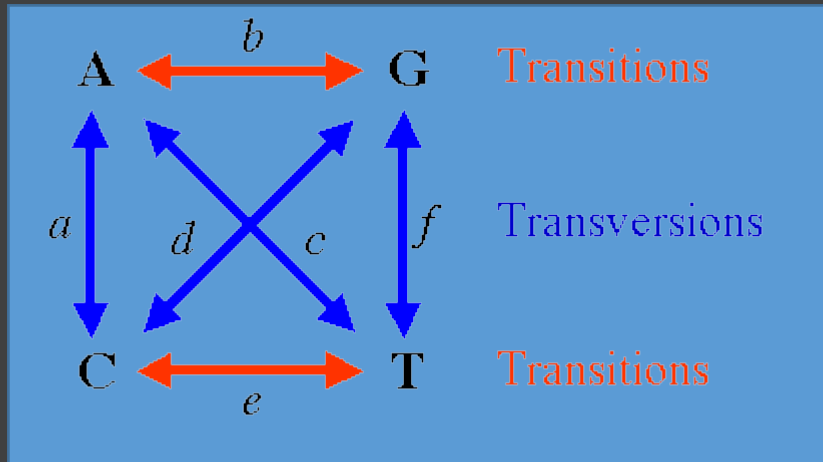
$$\hat{d} = -\frac{3}{4}\ln(1 - \frac{4}{3}p) = \hat{\nu}$$

d = est. subs per site
p = obs. proportion sites
that differ

probability matrix of observing subs – recall adjust for multiple hits

# Underlying models of evolution (DNA substitution models)

Simplest model (Jukes Cantor JC69) all substitution
rates are assumed of equal potential probability
and related to overall mutation rate (i.e. one parameter model)



$$Q = \begin{pmatrix} * & \frac{\mu}{4} & \frac{\mu}{4} & \frac{\mu}{4} \\ \frac{\mu}{4} & * & \frac{\mu}{4} & \frac{\mu}{4} \\ \frac{\mu}{4} & \frac{\mu}{4} & * & \frac{\mu}{4} \\ \frac{\mu}{4} & \frac{\mu}{4} & \frac{\mu}{4} & * \end{pmatrix}$$

Subs rate matrix

$$P = \begin{pmatrix} \frac{1}{4} + \frac{3}{4}e^{-t\mu} & \frac{1}{4} - \frac{1}{4}e^{-t\mu} & \frac{1}{4} - \frac{1}{4}e^{-t\mu} & \frac{1}{4} - \frac{1}{4}e^{-t\mu} \\ \frac{1}{4} - \frac{1}{4}e^{-t\mu} & \frac{1}{4} + \frac{3}{4}e^{-t\mu} & \frac{1}{4} - \frac{1}{4}e^{-t\mu} & \frac{1}{4} - \frac{1}{4}e^{-t\mu} \\ \frac{1}{4} - \frac{1}{4}e^{-t\mu} & \frac{1}{4} - \frac{1}{4}e^{-t\mu} & \frac{1}{4} + \frac{3}{4}e^{-t\mu} & \frac{1}{4} - \frac{1}{4}e^{-t\mu} \\ \frac{1}{4} - \frac{1}{4}e^{-t\mu} & \frac{1}{4} - \frac{1}{4}e^{-t\mu} & \frac{1}{4} - \frac{1}{4}e^{-t\mu} & \frac{1}{4} + \frac{3}{4}e^{-t\mu} \end{pmatrix}$$

probability matrix of observing subs – recall adjust for multiple hits

# Underlying models of evolution (DNA substitution models)

## K80 model – substitution rates can differ as two parameter model



A ⟷ G   Transitions
Transversions
C ⟷ T   Transitions

below). The K80 model assumes that all of the bases are equally frequent ($\pi T = \pi C = \pi A = \pi G = 0.25$).

Rate matrix $Q = \begin{pmatrix} * & \kappa & 1 & 1 \\ \kappa & * & 1 & 1 \\ 1 & 1 & * & \kappa \\ 1 & 1 & \kappa & * \end{pmatrix}$

The Kimura two-parameter distance is given by:

$$K = -\frac{1}{2}\ln((1-2p-q)\sqrt{1-2q})$$

p = obs Ts
q = obs Tv

Ts = transitions – do NOT change purine-pyrimidine state of sequences
Tv = transversions – do change purine-pyrimidine state of sequences

Ts are more common than Tv despite having half as many ways to happen
…this is largely due to spontaneous deamination of C -> T over time

In truth…mutational spectrum observed in MA lines in model organisms never seem to match those inferred in comparative genomics.   We (my lab) thinks this is due To universal weak selection on physical aspects of the genome (more to come…)

# Underlying models of evolution (DNA substitution models)

HKY85 model – accomodates base comp and Ts:Tv ratio

F81 model – base composition can vary

Felsenstein's 1981 model is an extension of the JC69 model in which base frequencies are allowed to vary from 0.25 ( $\pi_T \neq \pi_C \neq \pi_A \neq \pi_G$ )

Rate matrix:

$$Q = \begin{pmatrix} * & \pi_C & \pi_A & \pi_G \\ \pi_T & * & \pi_A & \pi_G \\ \pi_T & \pi_C & * & \pi_G \\ \pi_T & \pi_C & \pi_A & * \end{pmatrix}$$

The HKY85 model can be thought of as combining the extensions made in the Kimura80 and Felsenstein81 models. Namely, it distinguishes between the rate of transitions and transversions (using the $\kappa$ parameter), and it allows unequal base frequencies ( $\pi_T \neq \pi_C \neq \pi_A \neq \pi_G$ ). [ Felsenstein described a similar (but not equivalent) model in 1984 using a different parameterization;[5] that latter model is referred to as the F84 model.[6] ]

$$\text{Rate matrix } Q = \begin{pmatrix} * & \kappa\pi_C & \pi_A & \pi_G \\ \kappa\pi_T & * & \pi_A & \pi_G \\ \pi_T & \pi_C & * & \kappa\pi_G \\ \pi_T & \pi_C & \kappa\pi_A & * \end{pmatrix}$$

Other models

T92 – extends K80 to accommodate GC content change

TN93 – distinguishes between two types of Ts

GTR – general time reversible – 6 subs rate parameters and 4 equilibrium base
freq parameters   (206 parameters for 20 aa's)
(633 parameters for all 64 codons)

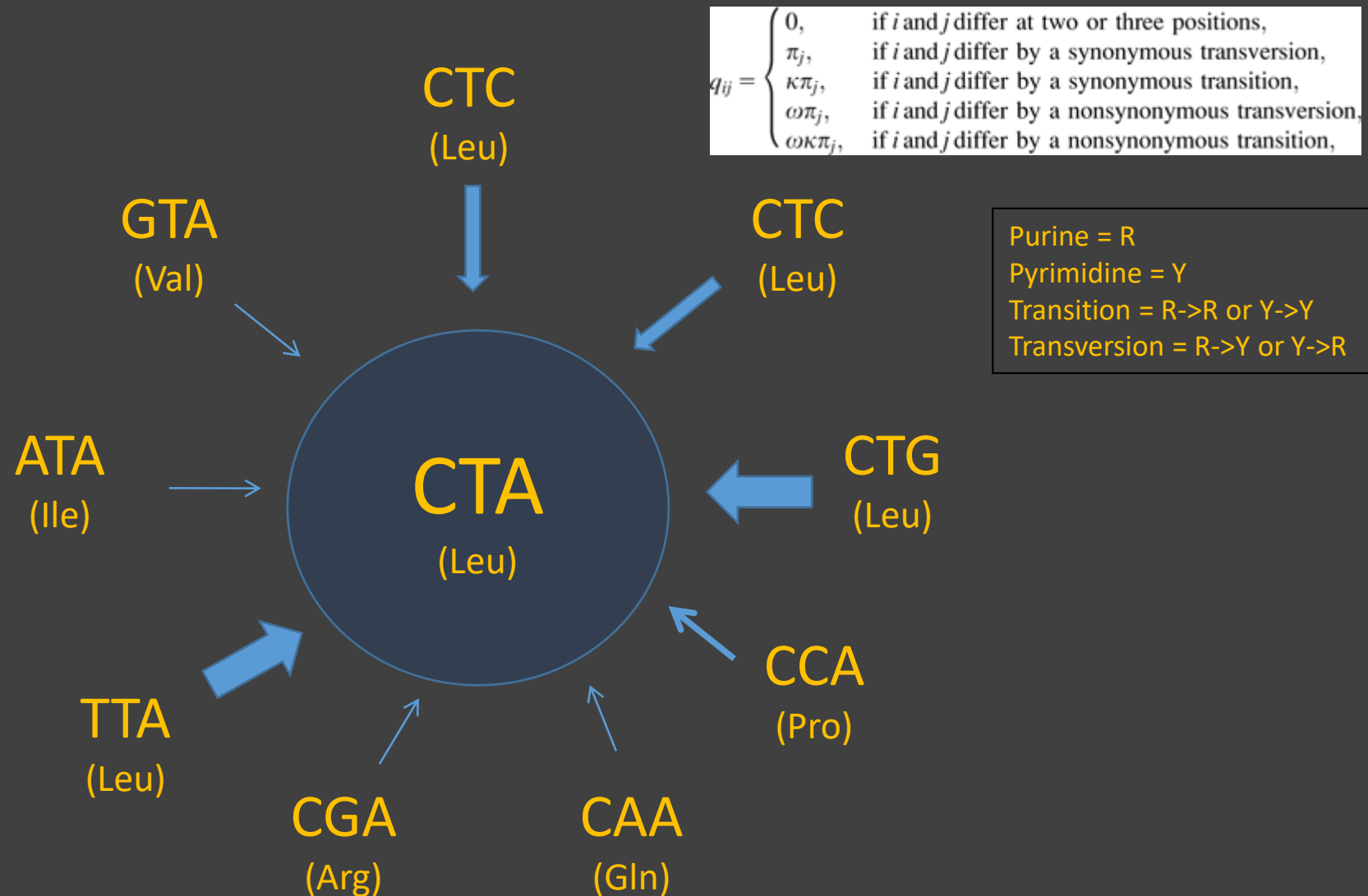# Protein evolution models - Dayhoff, JTT, WAG and LG models
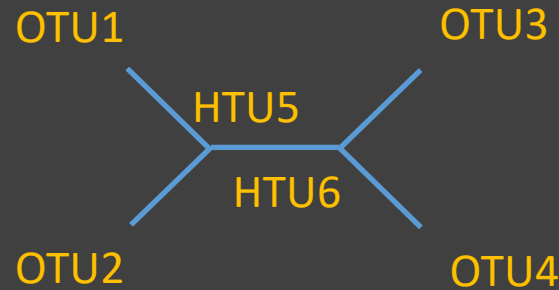


JTT – like PAM but based on more data

WAG – Whelan and Goldman
uses ML fitting very massive database to create subs matrix

LG – Le and Goldman
uses ML fitting of your data to create subs matrix

# Codon evolution models - 64 X 64 matrix

# Maximum likelihood- based phylogenetic reconstruction from aligned sequence data

OTU1    OTU3

HTU5

HTU6

OTU2    OTU4

See Figure 5.29 for application of ML to tree topology

1. starts with multiple sequence alignment
2. to compute L for given site (Ls), proceed over each site and sum probabilities of all possible sequence reconstructions for the each of the given sites for a given tree
3. to compute L for a given tree (Lt) take the product of all Ls for the given tree
4. or in practice take the sums of all ln(Ls)
5. search all tree topologies for maximum Lt

|      | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|------|---|---|---|---|---|---|---|---|
| OTU1 | A | G | G | A | C | T | G | G |
| OTU2 | A | A | C | T | C | T | G | C |
| OTU3 | A | A | A | C | A | T | G | A |
| OTU4 | A | A | A | G | G | T | G | A |



Total L for tree is sum of log L's for all sites (1-8)