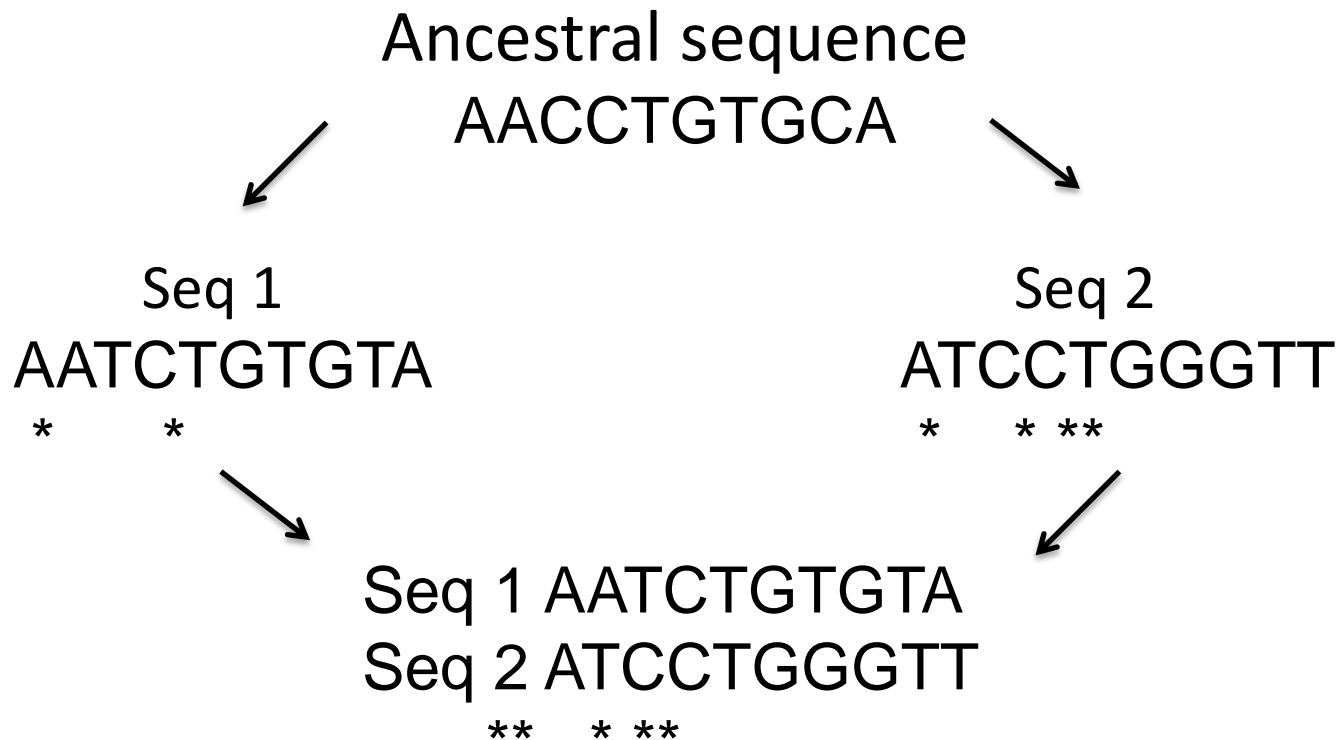# Introduction to Phylogenetics
# Week 4

# Phylogenetic Models

# I. Models

- Genetic distance
    - Used to determine divergence between sequences
    - Two identical sequences will diverge based on standard evolutionary rates
    - Rate depends a lot on how you model evolution
    - The evolutionary model you use is critical to obtaining a robust phylogenetic structure

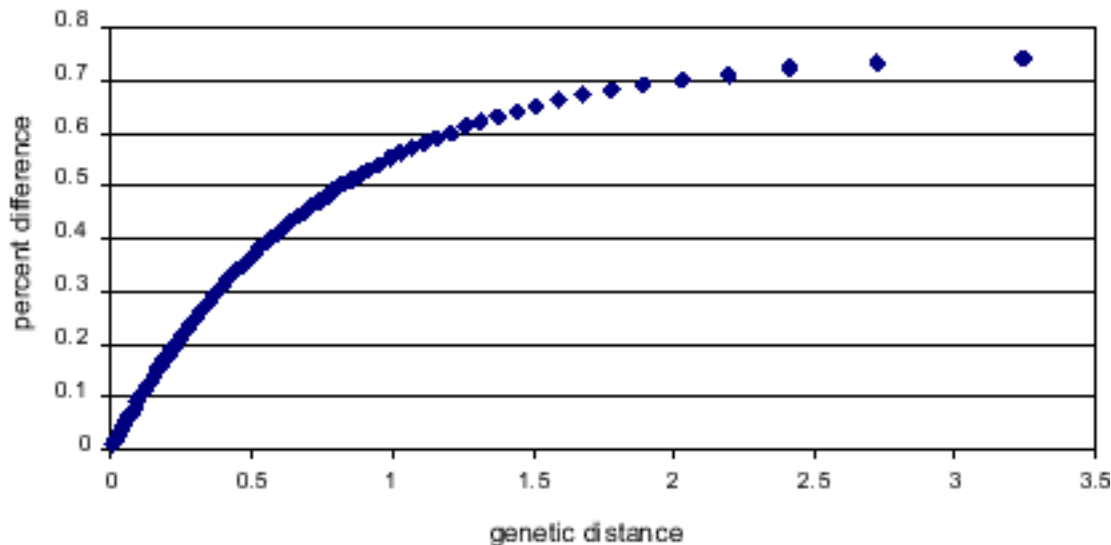# II. Observed and expected

- Simplest approach – count differences

Ancestral sequence
AACCTGTGCA

Seq 1
AATCTGTGTA
\*    \*

Seq 2
ATCCTGGGTT
\*    \* \*\*

Seq 1 AATCTGTGTA
Seq 2 ATCCTGGGTT
\*\*    \* \*\*

- Proportion is observed distance or $p$-distance

# II. Observed and expected

- Relatedness based on $p$-distance (Hamming distance) easy to imagine
- Only based on identity between sequences
- Cannot account for type of change/back-mutations
- No account of evolutionary processes (i.e. transition vs. transversion)
- Observed distance ($p$) underestimates true genetic distance ($d$)

# II. Observed and expected

- Overtime substitutions at each site accumulate and sequences are saturated



- Can use this to determine genetic distance for tree – not very robust.

# III. Mutations and time

- Substitutions assumed to be random event
- Substitution model provides statistical description of stochastic process
- Number of mutations X($t$) over time $t$.
- Use Poisson (P) distribution
  - Discrete probability distribution
  - Most basic model assumes mutation equally likely at each site
  - Occurs at rate $\mu$
  - $P_n$(t) probability n mutations in $t$

# III. Mutations and time

$$P_n(t) = [(\mu t)^n \exp(-\mu t)] / n!$$

- Number of substitutions up to time $t$ is distributed with factor $\mu t$ with variance $\mu t$
- Nucleotide substitution rate therefore is tied to $t$

# IV. Calculating Nucleotide Subs.

- Rates of substitution a Markov process
  - Stochastic model
  - Random system changes state according to transition rule
  - Can make predictions on future based on present state of system
- Basically – you know what the nucleotide is now, you can determine likelihood of assuming future state

# IV. Calculating Nucleotide Subs.

- Q matrix specifies rate of change for each nucleotide
- Way of describing all possible changes between states
- Probability departing from state $i$, arriving at state $j$
- Assumes state prior to $i$ has no impact on probability of $j$
- Rate of change modeled differently by different evolutionary models

$$\begin{array}{cccc}
\quad\quad A & \quad\quad C & \quad\quad G & \quad\quad T
\end{array}$$

$$Q = \begin{bmatrix}
-m(ap_C + bp_G + cp_T) & amp_C & bmp_G & cmp_T \\
gmp_A & -m(gp_A + dp_G + ep_T) & dmp_G & emp_T \\
hmp_A & imp_C & -m(hp_A + jp_C + fp_T) & fmp_T \\
jmp_A & kmp_C & lmp_G & -m(ip_A + kp_C + lp_T)
\end{bmatrix}$$

- $\mu$ is mean instantaneous substitution rate
- a, b, c… relative substitution rate (i.e. A to C)
- $\pi_G$, $\pi_A$, $\pi_T$… nucleotide frequencies
- Diagonal values so each row = 0 (no change)
- How you parameterize matrix determines how you model evolution

# V. Time Reversible Models

- Basic substitution models are probably not biologically relevant
- Do allow us to model stochastic events
- Time-reversible models assume rate of change $i$ to $j$ is the same as $j$ to $i$ (a = g, etc)
- Probability of nucleotide change at any site during evolutionary time ($t$)

$$P(t) = \exp(Qt)$$

# V. Jukes Cantor Model (JC69)

- When the probabilities of change $P(t)$ are known, can determine evolutionary distance between two sequences
- JC69 equilibrium frequency nucleotide = 25%

$$\pi_G = \pi_A = \pi_T = \pi_C = 1/4$$

- JC69 any nucleotide replaced by any other

$$a = b = c = d = e = f = g\ldots = 1$$

- Probability of nucleotide not changing $P_{ii}(t)$
- Probability of nucleotide replacement $P_{ij}(t)$

# V. Jukes Cantor Model (JC69)

Using JC69 Q matrix and given $P(t) = \exp(Qt)$

$$P_{ii}(t) = \frac{1}{4} + \frac{3}{4}\exp(-mt) \quad \text{and} \quad P_{ij}(t) = \frac{1}{4} - \frac{1}{4}\exp(-mt)$$

Comparing two sequences:   $p = \frac{3}{4}\left[1 - \exp(-2mt)\right]$

Solving for $\mu t$:
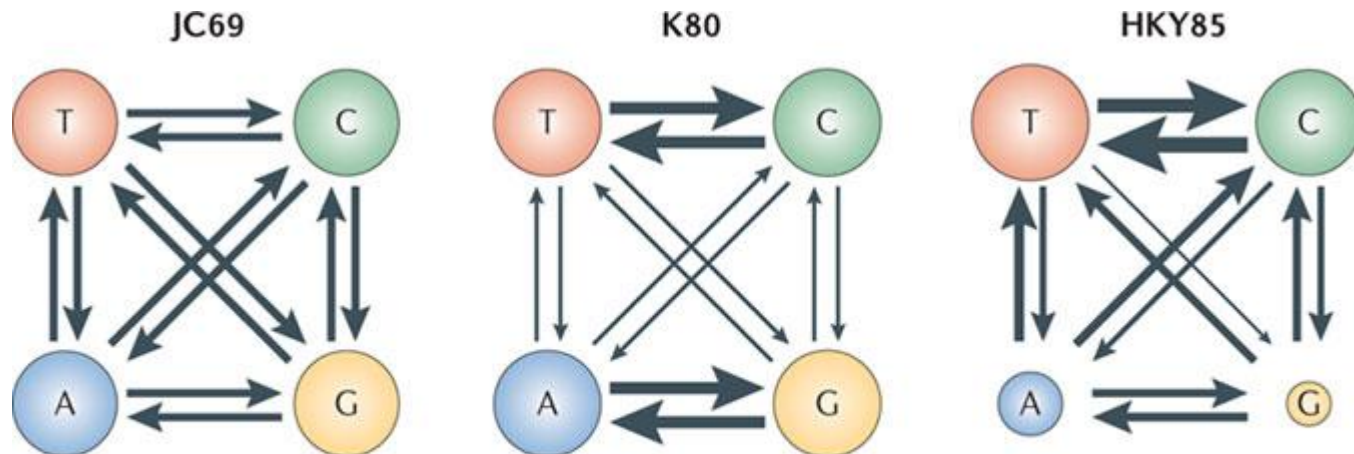
$$mt = -\frac{1}{2}\log(1 - \frac{4}{3}p)$$

# V. Jukes Cantor Model (JC69)

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

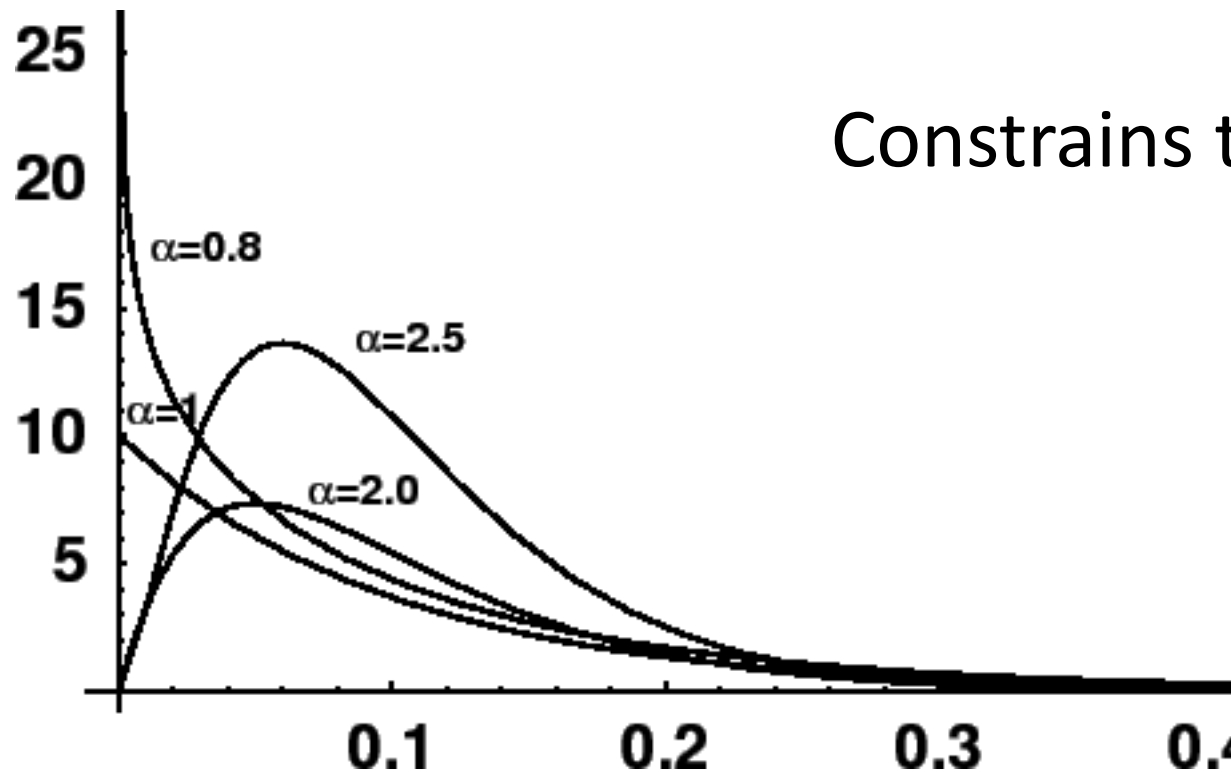$$d = -\frac{3}{4}\ln(1 - \frac{4}{3}p)$$

# VI. Nucleotide substitution models

- If all parameters of Q matrix determined – considered a general time reversible model (**GTR**)
- Parameterization reflects more of biological processes – rate heterogeneity between sites

# VI. Nucleotide substitution models

- Consider distribution of nucleotide changes
- Standard probability distributions - Gamma

Constrains the amount of site variability