

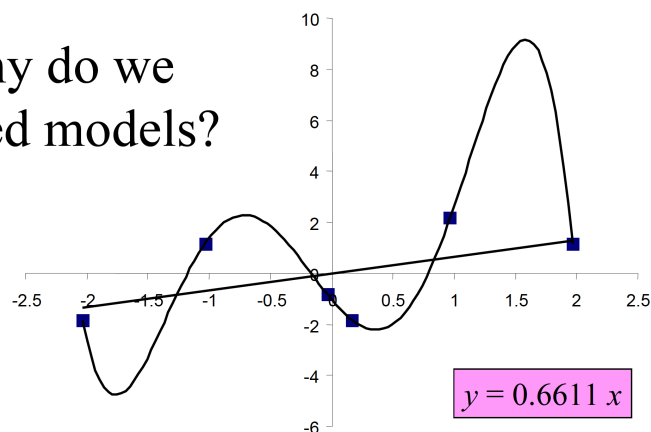
# Introduction to model-based methods

Niklas Wahlberg

Slides by Jadranka Rota, Paul Lewis and Chris Simon (University of Connecticut, USA)

$$y = -1.5972 x^5 + 23.167 x^4 - 126.18 x^3 + 319.17 x^2 - 369.22 x + 155.67$$

Why do we  
need models?



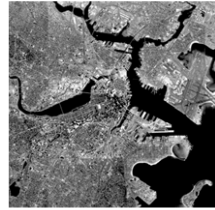
$$y = 0.6611 x$$

A very *practical* MBTA subway map



17

A very *realistic* MBTA subway map



18

► Which is more useful?

© 2005 by Paul O. Lewis

## Models

- Models help us intelligently **interpolate between our observations** for purposes of predicting future observations
- **Adding parameters** to a model generally increases its fit to the data
- **Underparameterized** models lead to poor fit to observed data points
- **Overparameterized** models lead to poor prediction of future observations
- Criteria for choosing models include likelihood ratio tests, AIC, BIC, Bayes Factors, etc.
  - all provide a way to choose a model that is neither underparameterized nor overparameterized

© 2005 by Paul O. Lewis

## Modelling nucleotide substitution

- With thousands of genomes sequenced
  - Good understanding of how DNA sequences evolve
  - Different **regions** of the genome have their own substitution dynamics
  - Different **lineages** may have their own substitution dynamics

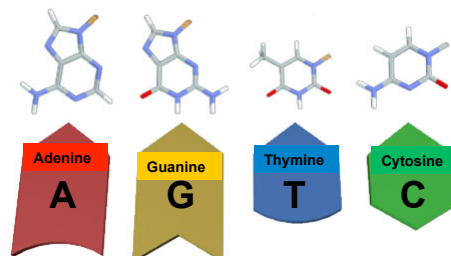
## Main Challenge

- ▶ DNA has only four characters

Purines

Pyrimidines

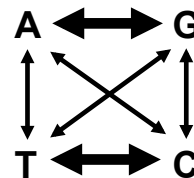
Figure B-3: The Four Nitrogenous Bases



Each base has a distinct shape that can be used to distinguish it from the others. 3D representations of the four bases are shown, with the corresponding chemical structures drawn above.

## Substitution types

- ▶ Purines: A, G
- ▶ Pyrimidines: C, T
- ▶ Transversions

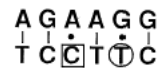


- Pu → Pyr

- Pyr → Pu

Pur - Pyr mispairs lead to transitions

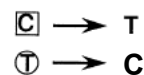
- ▶ Transitions – more common



In next round of replication

- Pu → Pu

- Pyr → Pyr



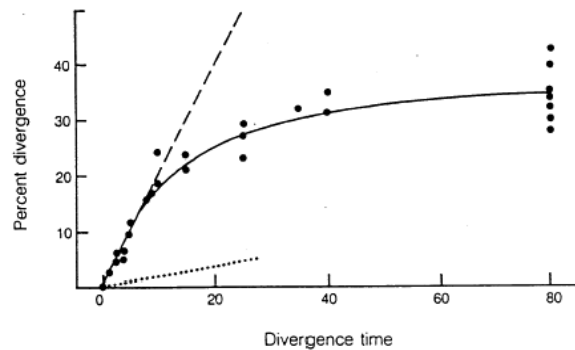
Slide by Chris Simon 2005

## Saturation in sequence data:

- Saturation is due to **multiple substitutions at the same site** subsequent to lineage splitting
- Models of evolution attempt to infer the missing information through correcting for “**multiple hits**”
- Most data will contain some fast evolving sites which are potentially saturated
  - e.g. in protein-coding genes codon position 3
- In severe cases the data become essentially random and all information about relationships can be lost
- **Probabilistic models of sequence evolution** are used to calculate expected changes

## Misleading DNA evolution

Multiple substitutions hide previous changes



Slide by Chris Simon 2005

Brown et al. 1979. PNAS 76:1967

## Difference between mutation and substitution

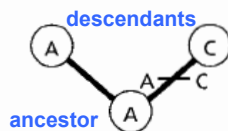
- **Substitutions** = mutational changes observed in populations
- **Mutations** = not all observed in populations, randomly distributed
  - 1) removed by proof reading enzymes
  - 2) cause death of cell, gamete, embryo

Slide by Chris Simon 2005

## Types of Substitutions

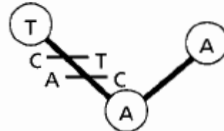
(a) Single substitution

1 change, 1 difference



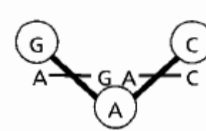
(b) Multiple substitution

2 changes, 1 difference



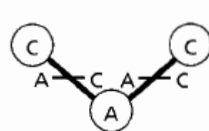
(c) Coincidental substitution

2 changes, 1 difference



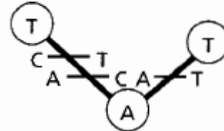
(d) Parallel substitution

2 changes, no difference



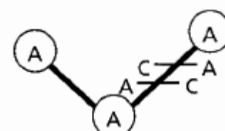
(e) Convergent substitution

3 changes, no difference



(f) Back substitution

2 changes, no difference



Page, R. and E. Holmes. 1998. Molecular Evolution: A phylogenetic Approach. Blackwell.

## Modelling nucleotide substitutions

- These dynamics can be modelled over a tree and they are incorporated into distance methods, maximum likelihood, and Bayesian inference
- Models incorporate information about the **rates at which each nucleotide is replaced** by each alternative nucleotide
  - For DNA this can be expressed as a 4 x 4 rate matrix (known as the Q matrix)
- Other model parameters may include:
  - **Site by site rate variation (aka among-site rate variation - ASRV)** - often modelled as a statistical distribution - for example a gamma distribution

## Corrections for multiple substitutions: First DNA substitution model

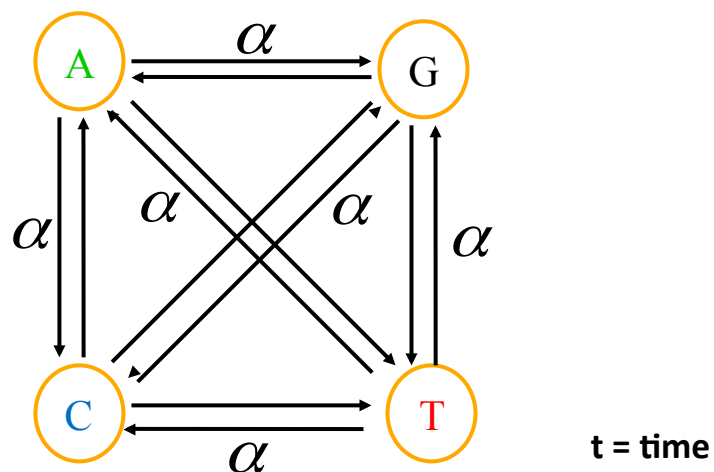
### Jukes & Cantor (1969) assumptions:

1.  $A = T = G = C$  No nucleotide bias
2. Every base changes to every other base with equal probability (no TS/TV bias)
3. All sites change with the same probability (no ASRV - among-site rate variation)

Also: probability of substitution & base composition remains constant over time/across lineages

Slide by Chris Simon 2005

### Jukes-Cantor model



- $\alpha$  = the rate of substitution ( $\alpha$  changes from A to G every t)
- The rate of substitution for each nucleotide is  $3\alpha$
- In t steps there will be  $3\alpha t$  changes

## The Q matrix

|      |   |            |            |            |            |
|------|---|------------|------------|------------|------------|
|      |   | To         |            |            |            |
|      |   | A          | C          | G          | T          |
| From | A | $-3\alpha$ | $\alpha$   | $\alpha$   | $\alpha$   |
|      | C | $\alpha$   | $-3\alpha$ | $\alpha$   | $\alpha$   |
|      | G | $\alpha$   | $\alpha$   | $-3\alpha$ | $\alpha$   |
|      | T | $\alpha$   | $\alpha$   | $\alpha$   | $-3\alpha$ |

## The Jukes-Cantor model: the simplest model

|   |            |            |            |            |
|---|------------|------------|------------|------------|
|   | A          | C          | G          | T          |
| A | $-3\alpha$ | $\alpha$   | $\alpha$   | $\alpha$   |
| C | $\alpha$   | $-3\alpha$ | $\alpha$   | $\alpha$   |
| G | $\alpha$   | $\alpha$   | $-3\alpha$ | $\alpha$   |
| T | $\alpha$   | $\alpha$   | $\alpha$   | $-3\alpha$ |

JC model: **one parameter model**

- 1) It assumes that all bases are equally frequent ( $p=0.25$ )
- 2) It assumes that all sites can change and they do so at the same rate –  $\alpha$



## The Jukes-Cantor model: the simplest model

|   | A        | C        | G        | T        |  |
|---|----------|----------|----------|----------|--|
| A | —        | $\alpha$ | $\alpha$ | $\alpha$ | JC model: <b>one parameter model</b><br>1) It assumes that all bases are equally frequent ( $p=0.25$ )<br>2) It assumes that all sites can change and they do so at the same rate – $\alpha$ |
| C | $\alpha$ | —        | $\alpha$ | $\alpha$ |  |
| G | $\alpha$ | $\alpha$ | —        | $\alpha$ |  |
| T | $\alpha$ | $\alpha$ | $\alpha$ | —        |  |

## Improvements on Jukes-Cantor

- Allow **base frequencies** to be unequal
- Allow **transitions** to be more common than **transversions**, in fact, allow separate estimates of the probability of change of **all six possible nucleotide substitutions**
- Allow the **probability of substitution to change along the molecule - ASRV**

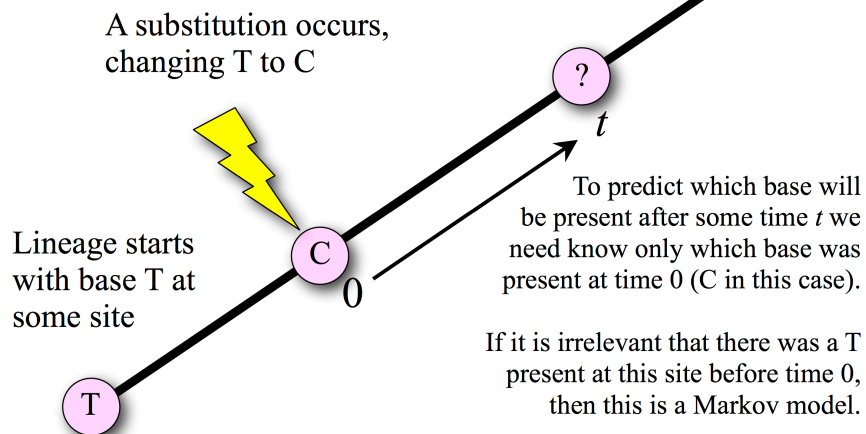
## Parameters we are interested in

- The mean instantaneous **substitution rate**  
=the general mutation rate + rate of fixation in population
- The relative **rates of substitution between each nucleotide**
- The average **frequencies of each base** in the dataset
- Topology (part of the model) and **branch lengths**

## Time-homogenous time-continuous stationary Markov models: Assumptions

- Rate of change from base  $i$  to base  $j$  is independent of the base that occupied a site prior to  $i$  (Markov property)

# What is a Markov process?

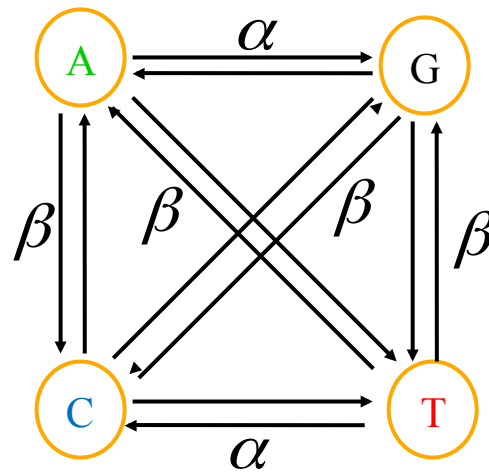


Paul O. Lewis (2014 Woods Hole Workshop in Molecular Evolution)

## Time-homogenous time-continuous stationary Markov models

- Rate of change from base  $i$  to base  $j$  is independent of the base that occupied a site prior to  $i$  (Markov property)
- Substitution rate does not change over time (homogeneity)
- Relative frequencies of A, G, C, and T are at equilibrium (stationarity)
- Rate of change from base  $i$  to base  $j$  is identical to the rate of change from base  $j$  to base  $i$  (time reversibility)

### Kimura (1980) model: K2P



$\alpha$  = transitions     $\beta$  = transversions

### The Kimura model has 2 parameters

|   | A        | C        | G        | T        |
|---|----------|----------|----------|----------|
| A | —        | $\beta$  | $\alpha$ | $\beta$  |
| C | $\beta$  | —        | $\beta$  | $\alpha$ |
| G | $\alpha$ | $\beta$  | —        | $\beta$  |
| T | $\beta$  | $\alpha$ | $\beta$  | —        |

K2P model is more realistic, but still

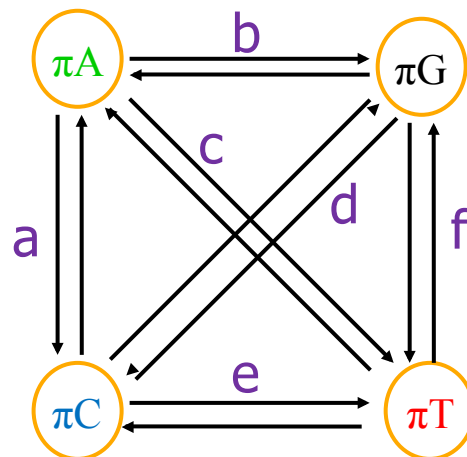
- 1) It assumes that all bases are equally frequent ( $p=0.25$ )
- 2) There are two substitution types (transitions –  $\alpha$  and transversions –  $\beta$ )

## The Hasegawa-Kishino-Yano model

|   | A              | C              | G              | T              |
|---|----------------|----------------|----------------|----------------|
| A | —              | $\pi_C \beta$  | $\pi_G \alpha$ | $\pi_T \beta$  |
| C | $\pi_A \beta$  | —              | $\pi_G \beta$  | $\pi_T \alpha$ |
| G | $\pi_A \alpha$ | $\pi_C \beta$  | —              | $\pi_T \beta$  |
| T | $\pi_A \beta$  | $\pi_C \alpha$ | $\pi_G \beta$  | —              |

HKY model:  
 1) Base frequencies are allowed to vary:  $\pi_A, \pi_C, \pi_G, \pi_T$   
 2) There are two substitution types (transitions -  $\alpha$  and transversions -  $\beta$ )

## The General Time-Reversible model



## The General Time-Reversible model (GTR)

|   | A         | C         | G         | T         |
|---|-----------|-----------|-----------|-----------|
| A | —         | $\pi_C a$ | $\pi_G b$ | $\pi_T c$ |
| C | $\pi_A a$ | —         | $\pi_G d$ | $\pi_T e$ |
| G | $\pi_A b$ | $\pi_C d$ | —         | $\pi_T f$ |
| T | $\pi_A c$ | $\pi_C e$ | $\pi_G f$ | —         |

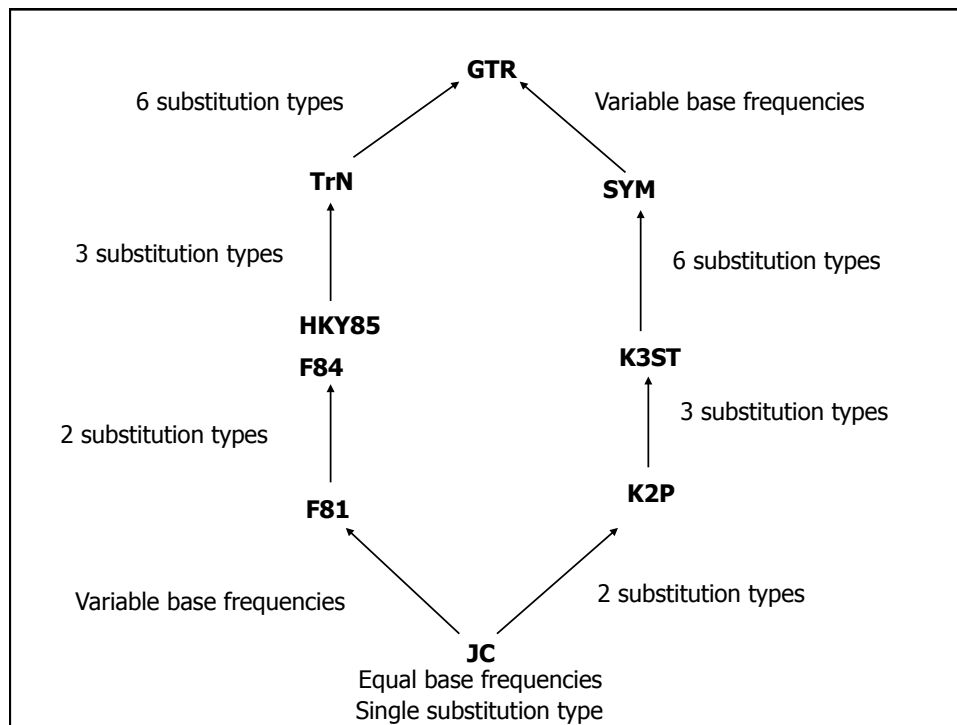
GTRmodel:

- 1) Base frequencies are allowed to vary:  $\pi_A$ ,  $\pi_C$ ,  $\pi_G$ ,  $\pi_T$
- 2) There are six substitution types:  $a$ ,  $b$ ,  $c$ ,  $d$ ,  $e$ ,  $f$

## The most commonly used models

- Almost all models used are special cases of one model:
  - The general time reversible model - GTR

ACAGGTGAGGCTCAGCCAAATTTGAGCTTTGTCGATAGGT



## Models

- **Model parameters can be:**
  - **estimated** from the data (using a likelihood function)
  - can be **pre-set** based upon assumptions about the data (for example that for all sequences all sites change at the same rate and all substitutions are equally likely - e.g. the Jukes-Cantor model)
  - *wherever possible avoid assumptions which are violated by the data because they can lead to incorrect trees*

## Modelling among-site rate variation (ASRV)

- All of the models so far assume that the **rate of change is the same for every position** in the alignment
- Biggest difference in substitution rate between variable and “invariable” sites
- Two classes of “invariable sites”
  - Highly restricted “not free to vary”
  - not observed to vary but in fact variable
    - due to convergence or reversal
    - % invariable sites can’t be calculated by simple sequence comparison.

Slide by Chris Simon 2005

ASRV, Yang 1996, TREE 11(9):367-372

## Why is modelling ASRV important?

- Protein-coding genes – 1<sup>st</sup>, 2<sup>nd</sup>, 3<sup>rd</sup> codon positions evolve differently from each other
- RNA molecules – stems and loops
- Introns vs. exons

RNA codon table

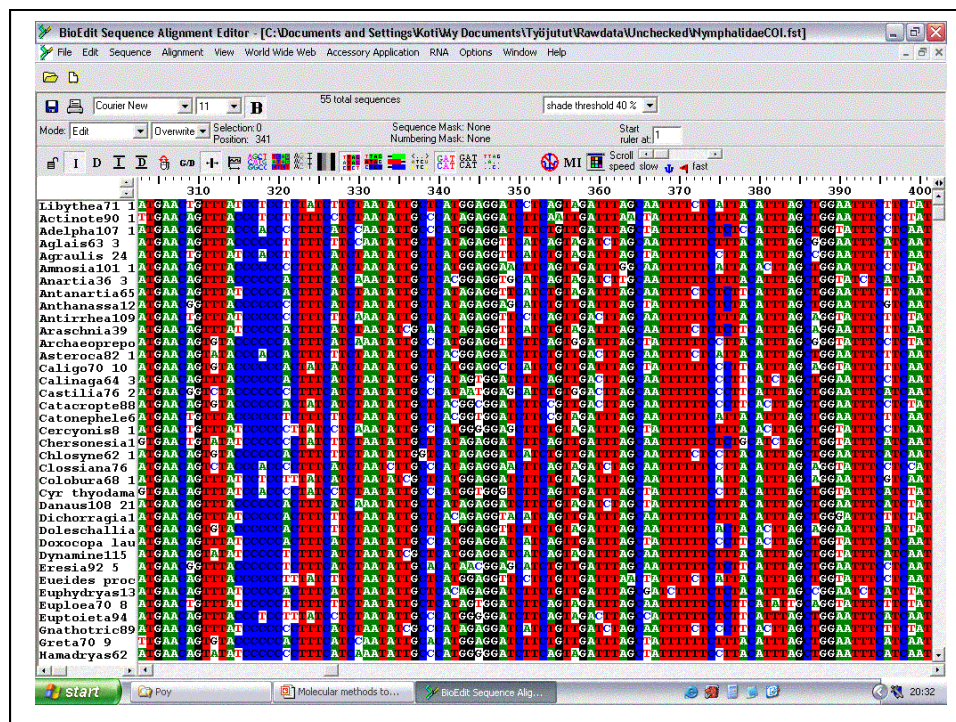
| 1st position | 2nd position             |                          |                            |                           | 3rd position     |
|--------------|--------------------------|--------------------------|----------------------------|---------------------------|------------------|
|              | U                        | C                        | A                          | G                         |                  |
| U            | Phe<br>Phe<br>Leu<br>Leu | Ser<br>Ser<br>Ser<br>Ser | Tyr<br>Tyr<br>stop<br>stop | Cys<br>Cys<br>stop<br>Trp | U<br>C<br>A<br>G |
| C            | Leu<br>Leu<br>Leu<br>Leu | Pro<br>Pro<br>Pro<br>Pro | His<br>His<br>Gln<br>Gln   | Arg<br>Arg<br>Arg<br>Arg  | U<br>C<br>A<br>G |
| A            | Ile<br>Ile<br>Ile<br>Met | Thr<br>Thr<br>Thr<br>Thr | Asn<br>Asn<br>Lys<br>Lys   | Ser<br>Ser<br>Arg<br>Arg  | U<br>C<br>A<br>G |
| G            | Val<br>Val<br>Val<br>Val | Ala<br>Ala<br>Ala<br>Ala | Asp<br>Asp<br>Glu<br>Glu   | Gly<br>Gly<br>Gly<br>Gly  | U<br>C<br>A<br>G |

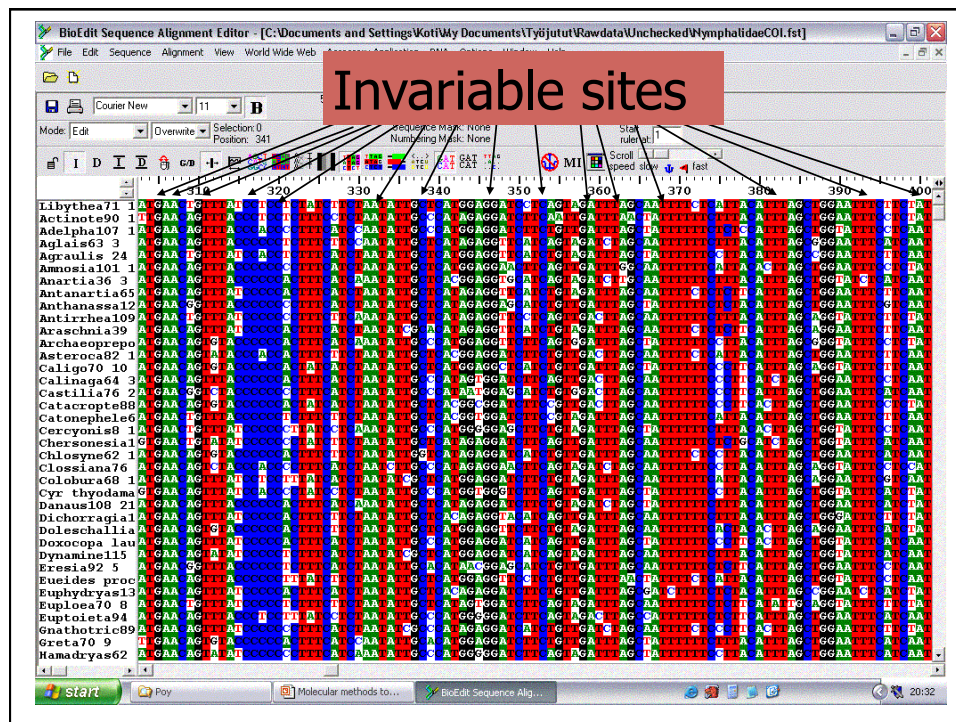
Amino Acids

Ala: Alanine  
Arg: Arginine  
Asn: Asparagine  
Asp: Aspartic acid  
Cys: Cysteine  
Gln: Glutamine  
Glu: Glutamic acid  
Gly: Glycine  
His: Histidine  
Ile: Isoleucine  
Leu: Leucine  
Lys: Lysine  
Met: Methionine  
Phe: Phenylalanine  
Pro: Proline  
Ser: Serine  
Thr: Threonine  
Trp: Tryptophane  
Tyr: Tyrosine  
Val: Valine



Slide by Chris Simon 2005





## Modelling among-site rate variation (ASRV)

- The most common additional parameters are:
  - A correction for the **proportion of sites** which are **invariable** (parameter  $I$ )
  - A correction for **variable site rates** at those sites which can change (parameter gamma,  $G$ )
- All models can be supplemented with these parameters (e.g. GTR+ $I$ + $G$ , HKY+ $I$ + $G$ )

## Modelling ASRV in variable sites

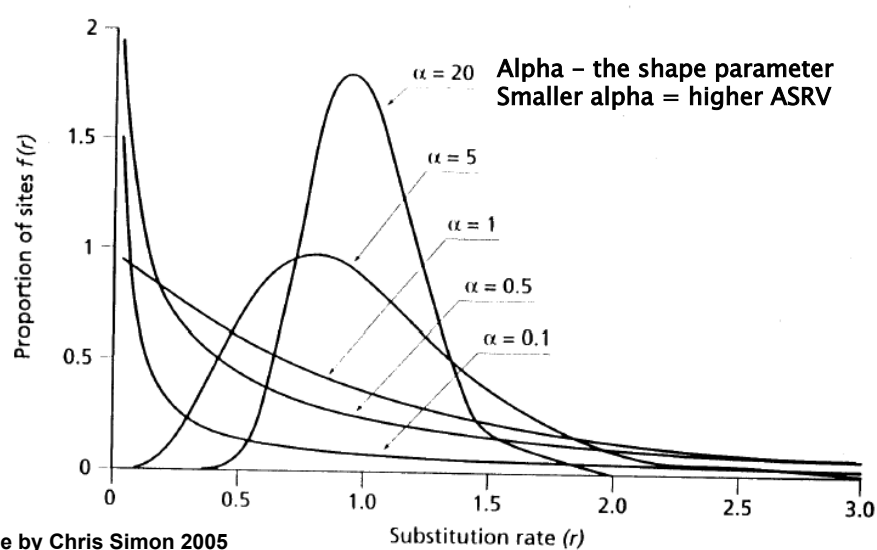
- ASRV in variable sites commonly modelled with a gamma distribution
- **Alpha** – the shape parameter of this distribution

Slide by Chris Simon 2005

ASRV, Yang 1996, TREE 11(9):367-372

## Gamma distribution:

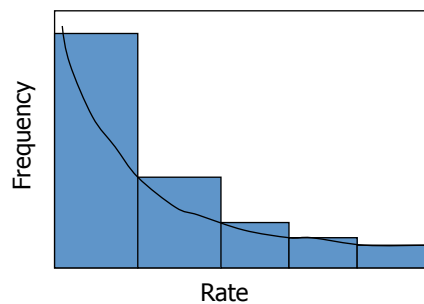
Relative substitution rates for different  $\alpha$  values



Slide by Chris Simon 2005

## Gamma distribution computationally costly

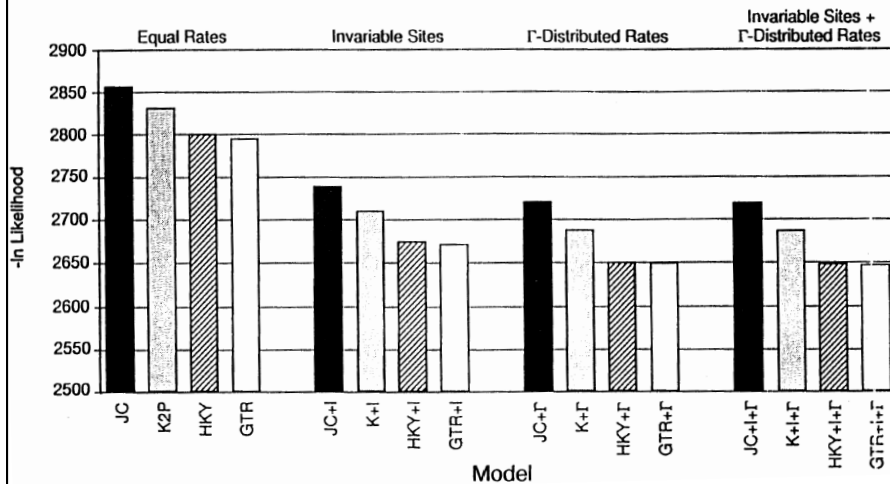
- Computational difficulties in using continuous distribution
- Most programs use discrete categories



## ASRV: Yang discrete model

- Continuous data divided into “n” discrete rate classes (generally 4)
- If  $\alpha < 0.2$  Yang recommends more rate classes
- Less computer intensive than obtaining likelihoods by integrating over the continuous gamma distribution

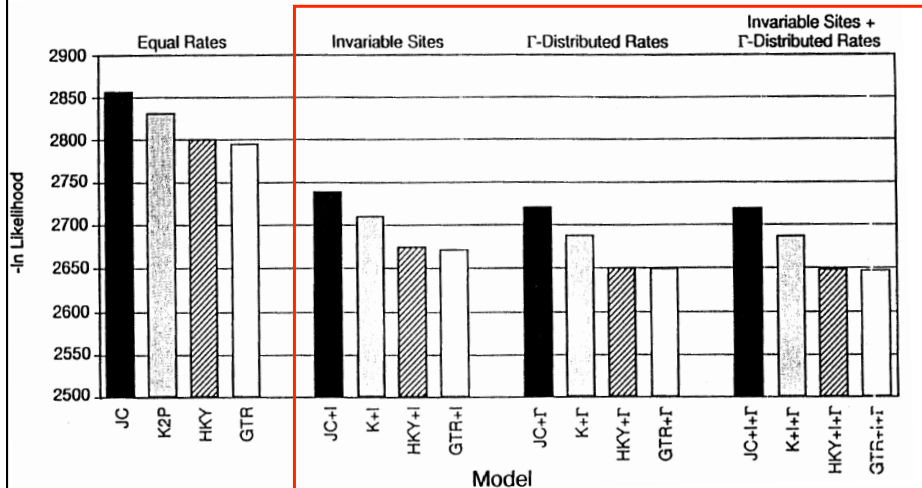
## ASRV >> fit improvement than by other parameters



Frati, Simon, Sullivan, Swofford. 1997. JME 44:145-158

Slide by Chris Simon 2005

## ASRV >> fit improvement than by other parameters



Frati, Simon, Sullivan, Swofford. 1997. JME 44:145-158

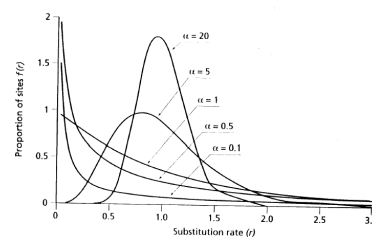
Slide by Chris Simon 2005

## Difficulties in estimating ASRV

- The parameters  $I$  and  $G$  covary!
- $(I + G)$  can be estimated, but the values of  $I$  and  $G$  are not easily teased apart
- Parameter  $G$  takes  $I$  into account,  $I$  not needed (in many/most? datasets)

## Another method for modelling ASRV

- Gamma distribution is always unimodal
  - Not necessarily the case in our dataset!
- Flexible rate heterogeneity across sites model
  - Probability distribution free model so that you can find the distribution that fits your data (FreeRate Model)
  - Implemented in IQ-TREE



Kalyaanamoorthy et al. 2017 (Nature Methods) doi:10.1038/nmeth.4285

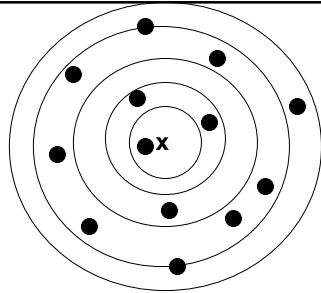
## **Parameters in models of DNA evolution**

- **Numbers of parameters estimated:**
  - Substitutions (up to 5; 1 fixed, 5 estimated)
  - Base composition (1 fixed, 3 estimated)
  - Among-site-rate variation
    - Gamma shape parameter = 1 parameter
    - Invariant sites = 1 parameter
    - Gamma + I = 2 parameters
  - Partitioned models – add up parameters of each partition

Slide by Chris Simon 2005

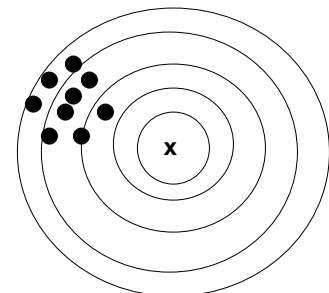
## **Models can be made more parameter rich to increase their realism**

- **But the more parameters estimated, the more time needed, and the more sampling error accumulates**
  - One might have a realistic model but large sampling errors
  - Realism comes at a cost in time and precision!
  - Fewer parameters may give an inaccurate estimate, but more parameters decrease the precision of the estimate
  - In general use the simplest model which fits the data



### Trade-off between highly parameterized models & model error variance

- Many parameters, higher error variance but clustered around the true value (higher accuracy, lower precision)
- Few parameters, lower error variance but may not be centered around the mean (lower accuracy, higher precision)



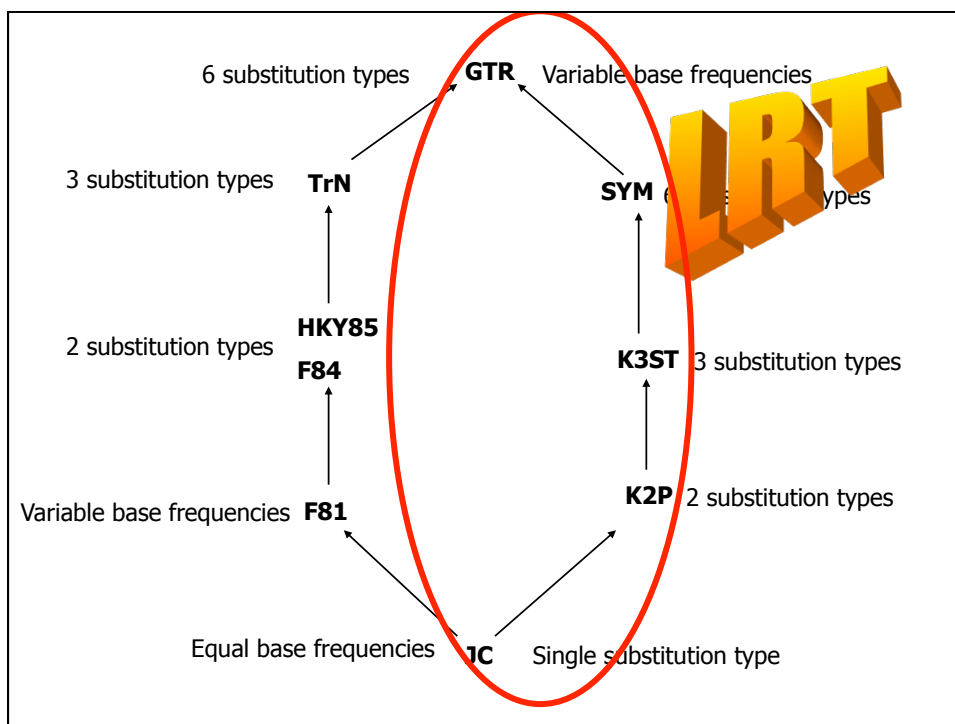
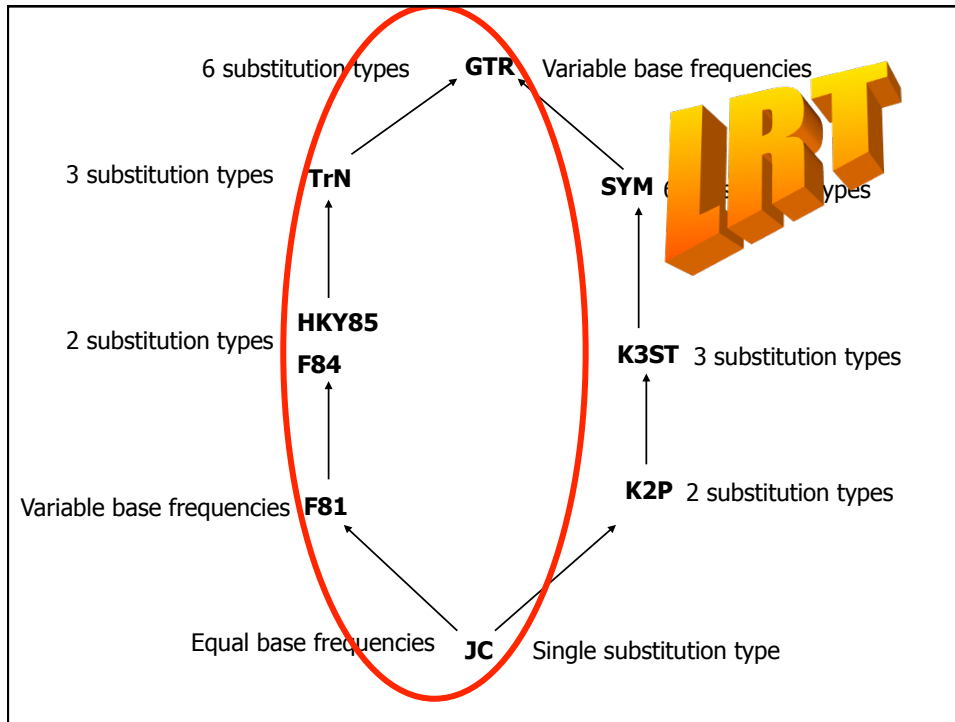
Dave Swofford's Target Analogy

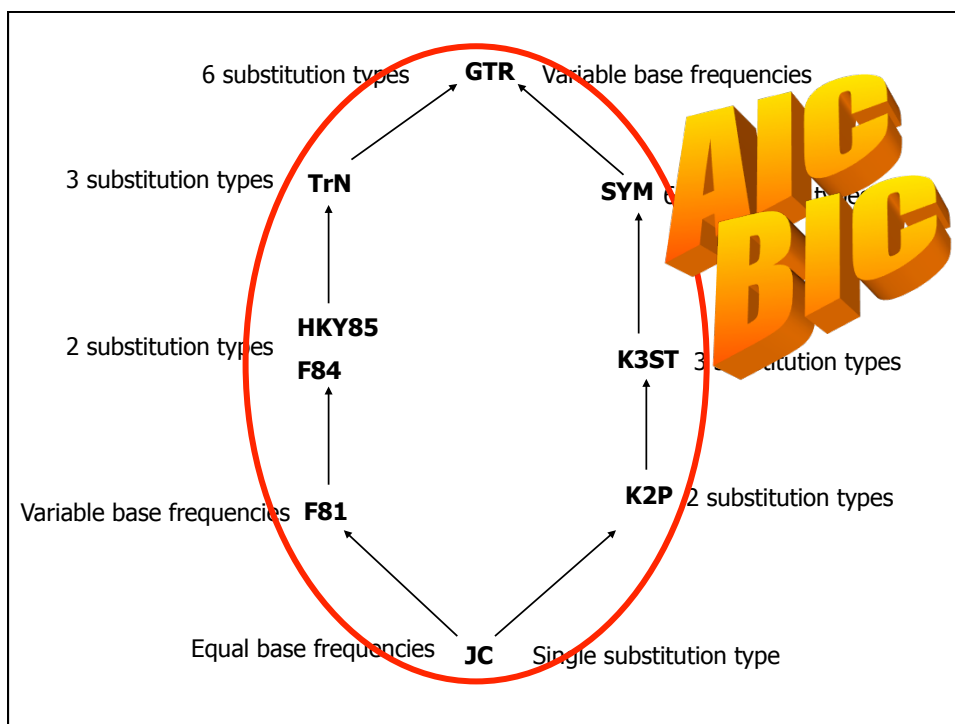
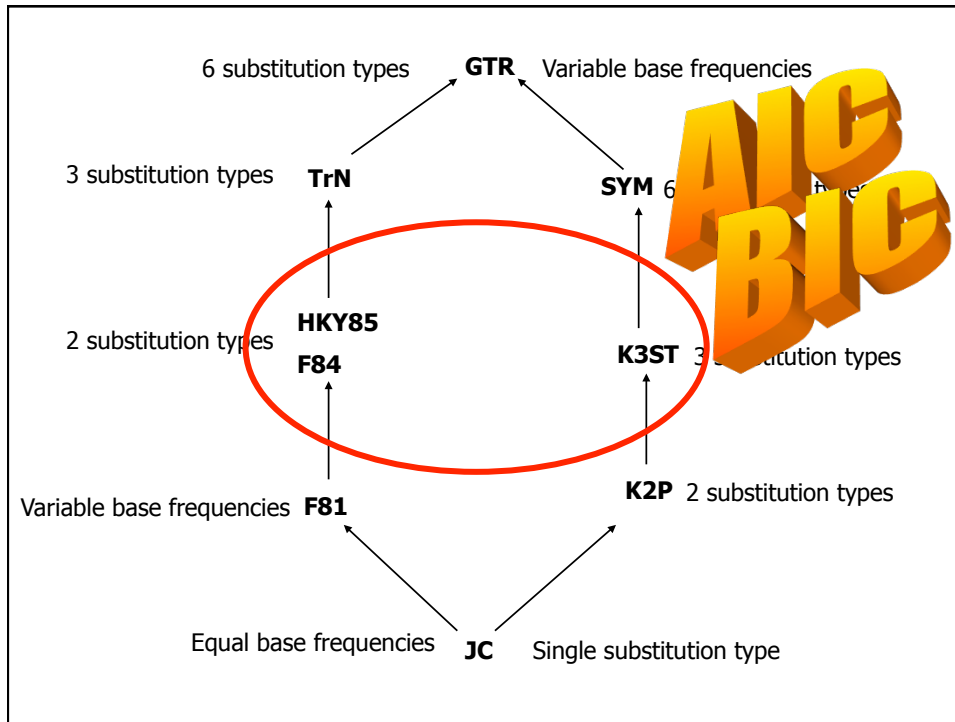
Slide by Chris Simon 2005

## Choosing between models

- Tools to determine whether the model can estimate parameters from the data
- When models are nested
  - Likelihood ratio test (LRT)
- When models are not nested
  - Akaike Information Criterion (AIC)
  - Bayesian Information Criterion (BIC)







## Estimation of substitution model parameters

- Yang (1995) has shown that parameter estimates are reasonably stable across tree topologies provided trees are not “**too wrong**”
- Thus one can obtain a tree using a quick method and then estimate parameters on that tree
- These parameters can then be used to calculate the likelihood of a model for model comparison

## Model-testing programs

- **Modeltest**
  - Posada & Crandall. 1998. MODELTEST: testing the model of DNA substitution. *Bioinformatics* 14(9): 817-818.
- **jModeltest**
  - Darriba et al. 2012. jModelTest 2: more models, new heuristics and parallel computing. *Nature Methods* 9(8), 772.
- **PartitionFinder**
  - Lanfear et al. 2016. PartitionFinder 2: new methods for selecting partitioned models of evolution for molecular and morphological phylogenetic analyses. *MBE* 34(3), 772 – 773.
- **ModelFinder built into IQ-Tree**
  - S. Kalyaanamoorthy, B.Q. Minh, T.K.F. Wong, A. von Haeseler, and L.S. Jermiin (2017) ModelFinder: Fast Model Selection for Accurate Phylogenetic Estimates, *Nature Methods*, 14:587–589. <https://doi.org/10.1038/nmeth.4285>

## Model testing easier nowadays

- **Bayesian statistical framework**
  - MrBayes has a model jumping feature
  - It samples over all possible models based on their probabilities
  - No longer necessary to test for which model is optimal
- **Maximum Likelihood framework**
  - IQ-Tree - ModelFinder implemented

## Partitioned models (1/2)

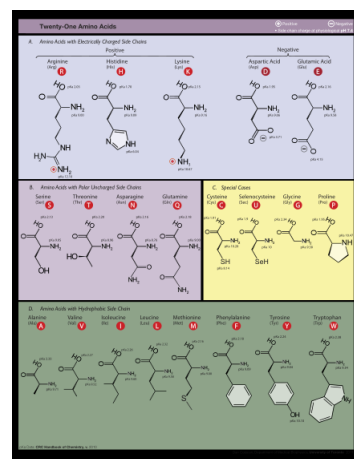
- Today's datasets tend to be large, including hundreds or thousands of genes
- Unrealistic to have the same model for the whole dataset (**underparameterization**)
- Modelling DNA substitution for separate sections of the data (partitions)
  - E.g. different genes, codon positions, introns/exons, etc.
- To avoid **overparameterization**, partitions with similar properties can be merged

## Partitioned models (2/2)

- This approach allows us to accommodate heterogeneity across data subsets in overall rate and in substitution model parameters
- In some programs also possible to unlink topology and branch lengths so that each data subset evolves differently from each other
- Built into IQ-Tree

## Models of amino acid substitution

- Empirical and mechanistic models
- **Empirical models:** based on empirical AA replacement with matrices from different taxa
  - 20 amino acids – 20x20 matrix too big for estimation
  - Examples: JTT, WAG, LG, MtREV (for mitochondria), Blosum62
- **Mechanistic models:**
  - e.g. codon models (61x61 matrix)
  - Tend to outperform empirical models BUT
  - Computationally very intensive



## Recommended reading

- Christoph Bleidorn (2017) [Phylogenomics: An Introduction](#) (DOI: 10.1007/978-3-319-54064-1)
- Hoff et al. 2016. [Does the choice of nucleotide substitution models matter topologically?](#) BMC Bioinformatics 17: 143. doi.org/10.1186/s12859-016-0985-x
- Kainer & Lanfear. 2015. [The Effects of Partitioning on Phylogenetic Inference](#). Molecular Biology and Evolution, 32(6), 1611–1627. doi.org/10.1093/molbev/msv026