

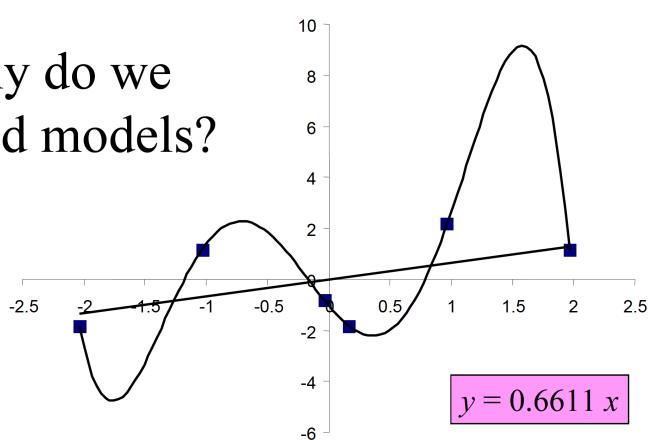
Molecular Phylogenetics Course

# Introduction to model-based methods

Jadranka Rota

Some slides by Paul Lewis and Chris Simon (University of Connecticut, USA)

Why do we  
need models?

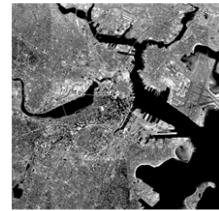


A very *practical* MBTA subway map



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

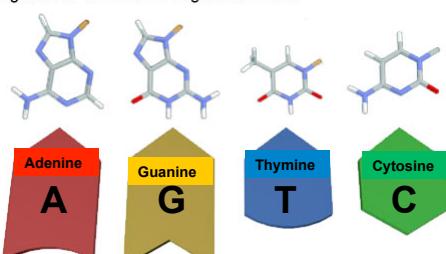
## Main Challenge

- ▶ DNA has only four characters

Purines

Pyrimidines

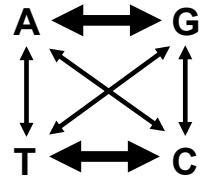
Figure B-3: The Four Nitrogenous Bases



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
- ▶ Transitions – more common
  - Pu → Pu
  - Pyr → Pyr



Pur - Pyr mispairs lead to transitions

A G A A G G  
T C C T T C

In next round of replication

C → T  
T → C

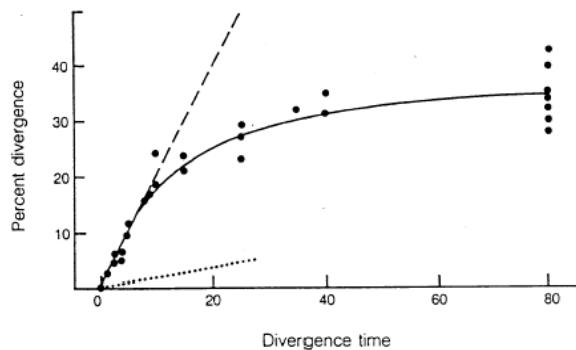
Slide by Chris Simon 2005

## 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 distances

## 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

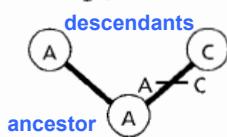
- **Substitutions** = mutational changes observed in populations
- **Mutations** = not all observed in populations, randomly distributed
  - 1) removed by proof reading enzymes
  - 2) cause death of cell, gamete, embryo

Slide by Chris Simon 2005

## Types of Substitutions

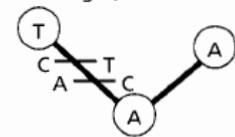
(a) Single substitution

1 change, 1 difference



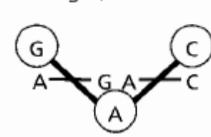
(b) Multiple substitution

2 changes, 1 difference



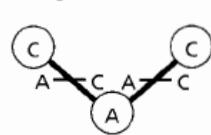
(c) Coincidental substitution

2 changes, 1 difference



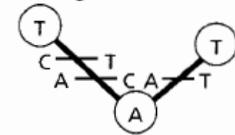
(d) Parallel substitution

2 changes, no difference



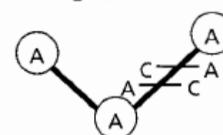
(e) Convergent substitution

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

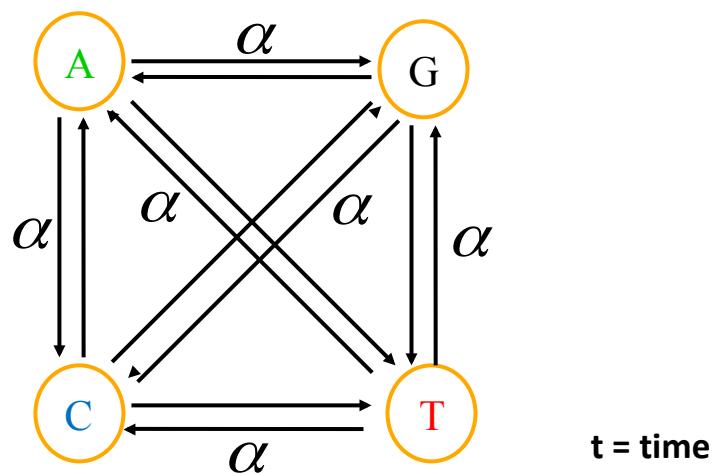
Jukes & Cantor (1969) assumptions:

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	-3α	α	α	α
	C	α	-3α	α	α	
	G	α	α	-3α	α	
	T	α	α	α	-3α	

## The Jukes-Cantor model: the simplest model

	A	C	G	T
A	-3α	α	α	α
C	α	-3α	α	α
G	α	α	-3α	α
T	α	α	α	-3α

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	-	$\alpha$	$\alpha$	$\alpha$
C	$\alpha$	-	$\alpha$	$\alpha$
G	$\alpha$	$\alpha$	-	$\alpha$
T	$\alpha$	$\alpha$	$\alpha$	-

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 –  $\alpha$

## 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

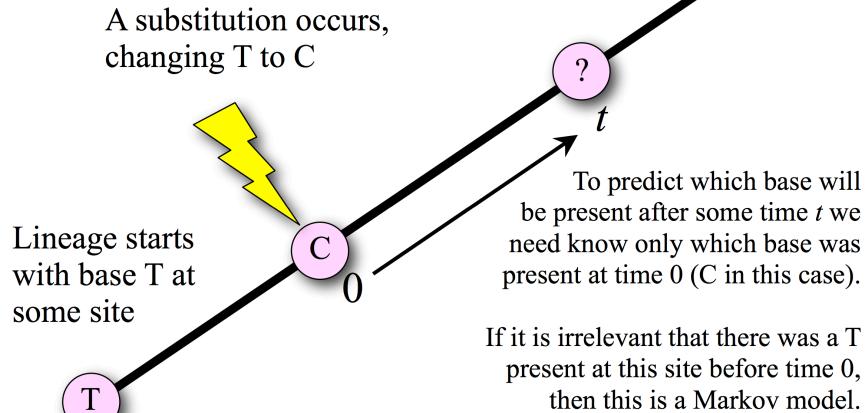
## Parameters we are interested in

- 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?

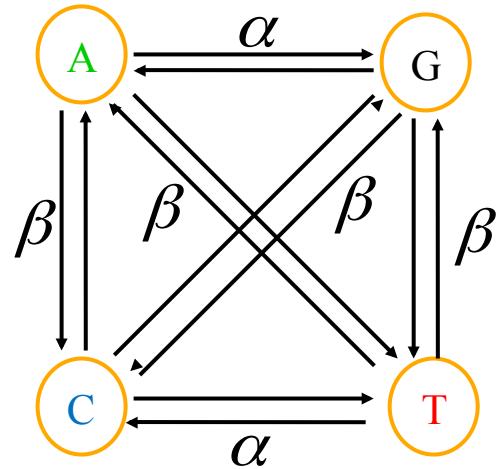


Paul O. Lewis (2014 Woods Hole Workshop in Molecular Evolution)

### 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	-	$\beta$	$\alpha$
C	$\beta$	-	$\alpha$
G	$\alpha$	$\beta$	-
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 -  $\beta$ )

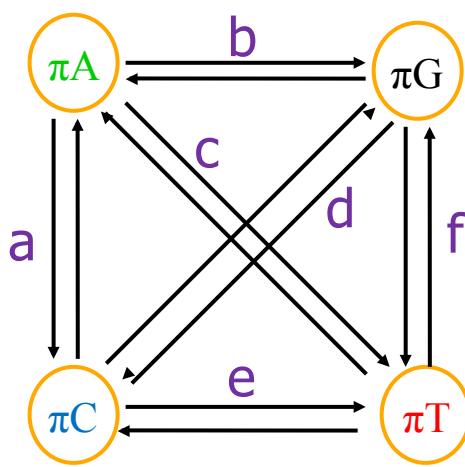
## The Hasegawa-Kishino-Yano model

	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	$\pi_A c$	$\pi_C e$	$\pi_G f$	—

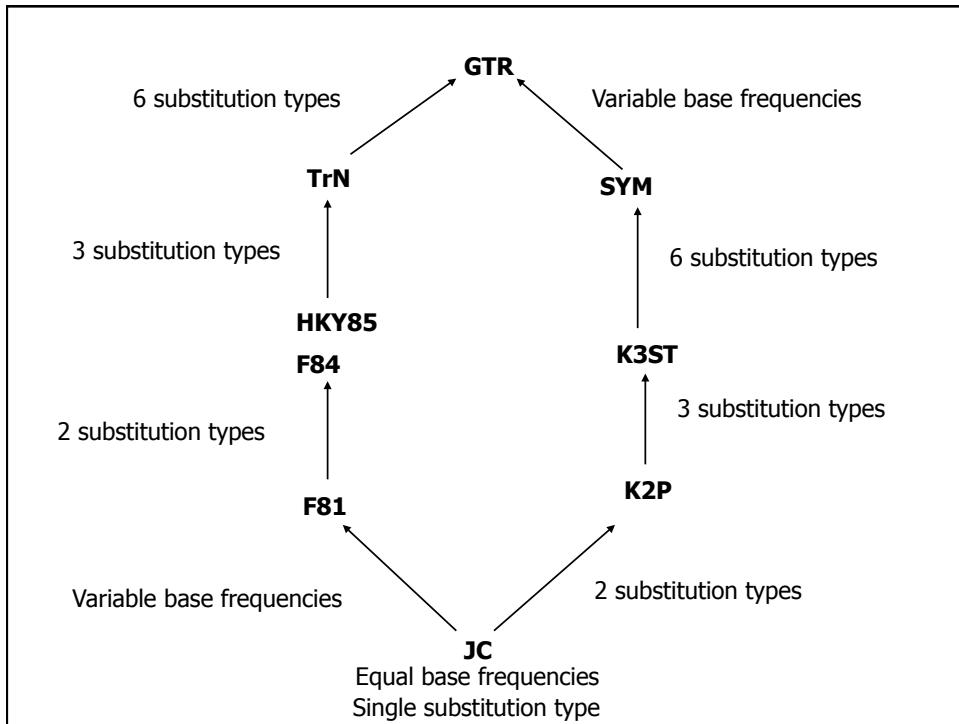
GTRmodel:

- 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

ACAGGTGAGGCTCAGCCAATTTGAGCTTTGTCGATAGGT



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

## Typical pattern of variation among codon positions

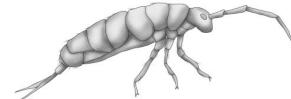
E.g. in mtDNA in Collembola

**56.7% of all variable sites are located in third positions**

**1st 27.9%    2nd 15.4%    3rd 56.7%**

**96.9% of all third positions are variable**

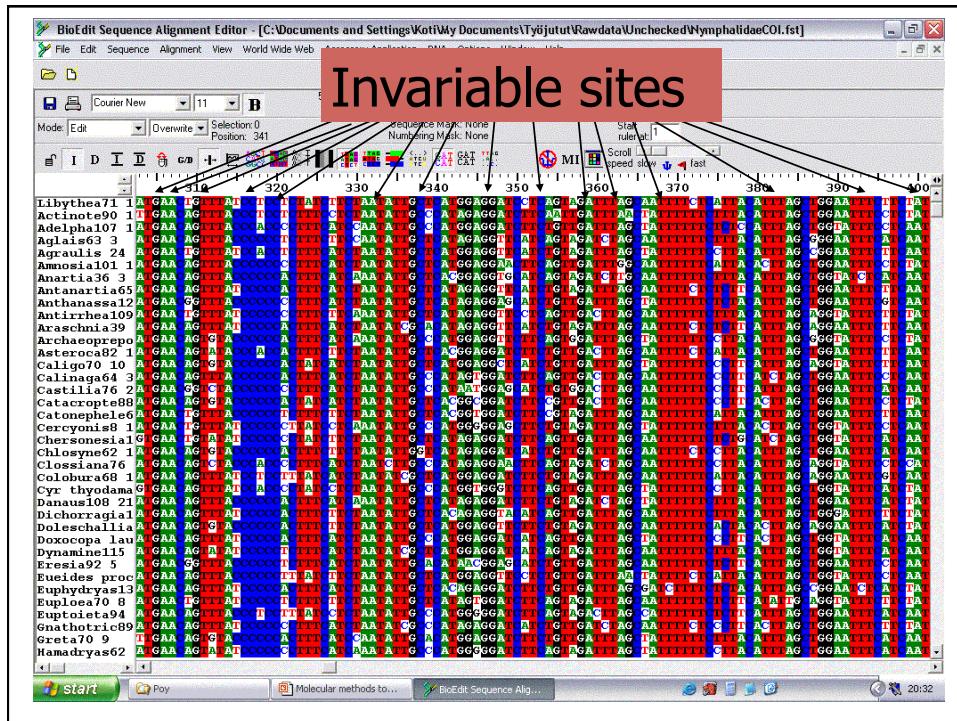
**1st 47.8%    2nd 26.3%    3rd 96.9%**



**Frati et al. 1997. J. Mol. Evol.**

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+ $I+G$ )

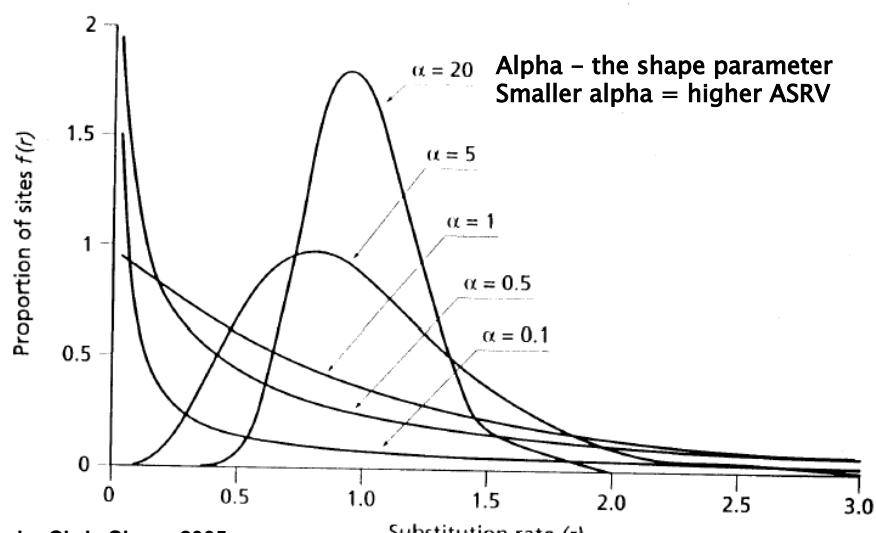
## Modelling ASRV in variable sites

- 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

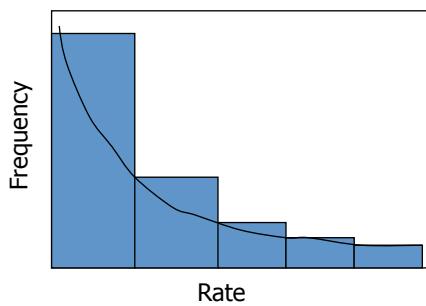
### Gamma distribution: Relative substitution rates for different $\alpha$ values



Slide by Chris Simon 2005

## Gamma distribution computationally costly

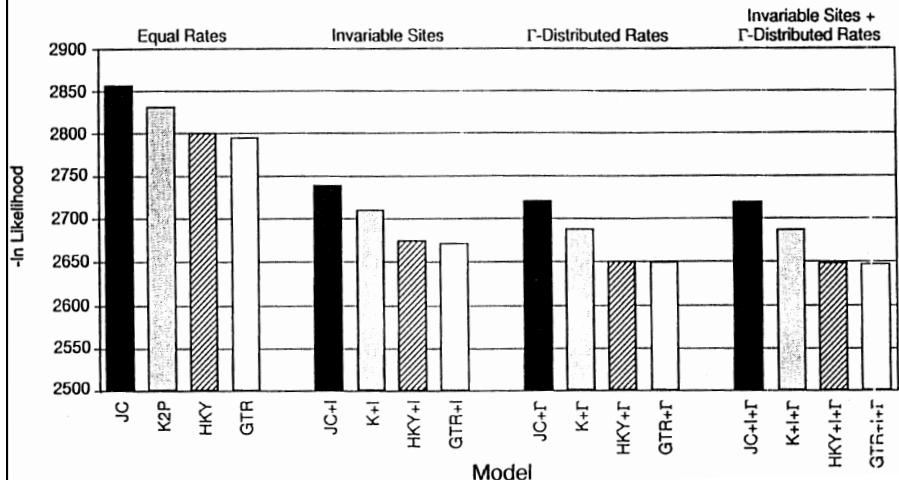
- Computational difficulties in using continuous distribution
- Most programs use discrete categories



## 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

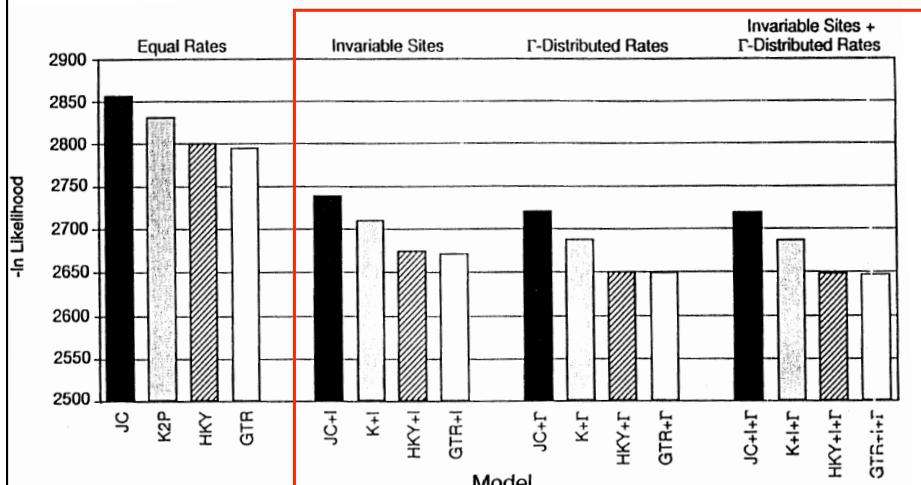
## ASRV >> fit improvement than by other parameters



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

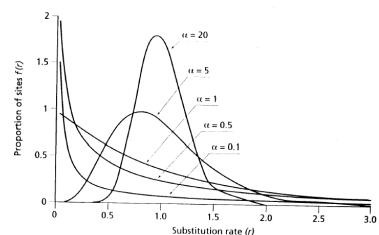
Slide by Chris Simon 2005

## 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

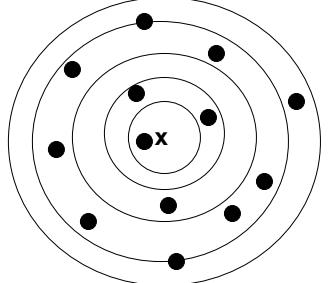
## Parameters in models of DNA evolution

- **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

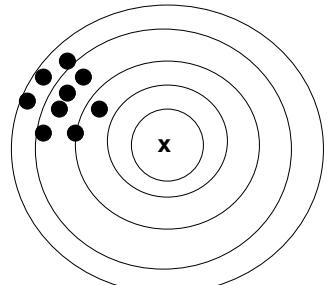
Slide by Chris Simon 2005

## 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



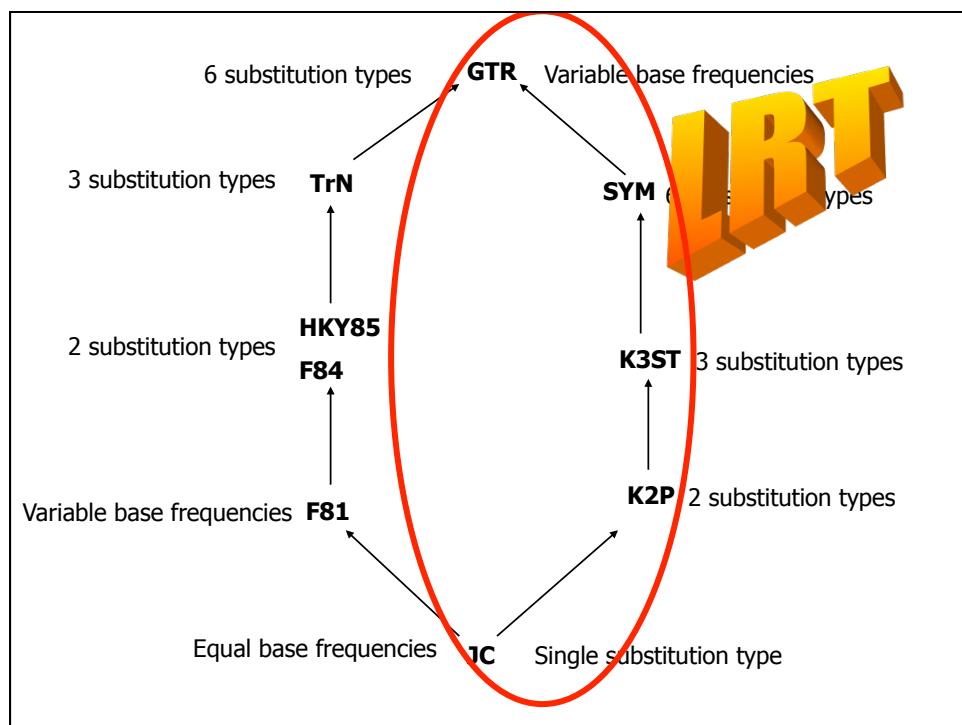
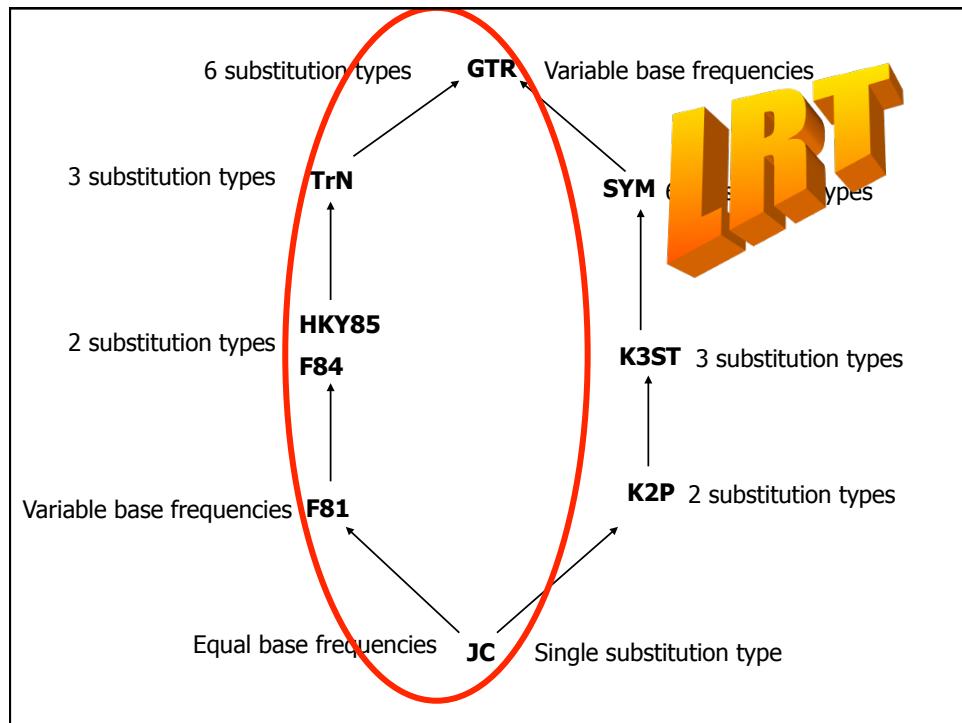
Dave Swofford's Target Analogy

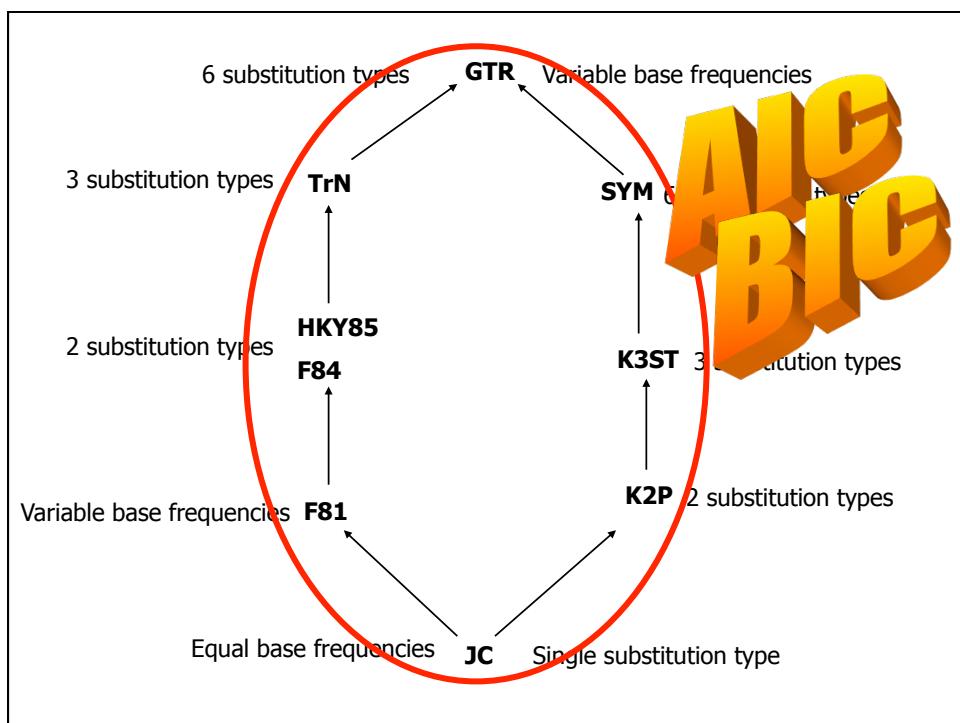
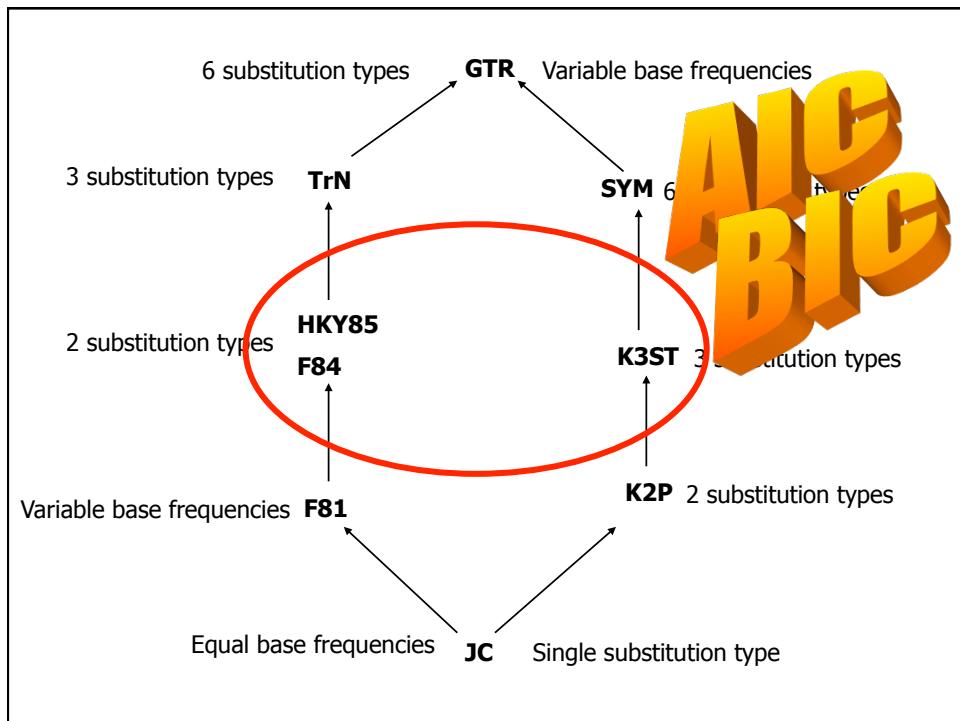
- Many parameters, higher error variance but clustered around the true value (higher accuracy, lower precision)
- Few parameters, lower error variance but may not be centered around the mean (lower accuracy, higher precision)

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

## Need to know the likelihood of a model

- For these tests, one needs to [compute the likelihood of the model](#)
- Covered in next lecture
- For now, assume we know the likelihood of the models we want to compare
- Comparison tools:
  - Likelihood ratio test (LRT)
  - Akaike information criterion (AIC) and corrected AIC ( $AIC_c$ )
  - Bayesian information criterion (BIC)

## Likelihood ratio test (LRT)

$$LR = 2 * (\ln L_1 - \ln L_0)$$

Alternative hypothesis  
More parameter-rich

Null hypothesis  
Less parameter-rich

- LRT statistic approximately follows a chi-square distribution
- Degrees of freedom equal to the number of extra parameters in the more complex model

## Akaike Information Criterion

- A measure of the **relative quality of statistical models for a given dataset** ([Wikipedia definition](#))
  - It deals with the trade-off between the goodness of fit and the complexity of the model
- **AIC( $M$ ) =  $-2 * \text{Log}(\text{Likelihood}(M)) + 2 * K(M)$** 
  - $K(M)$  is the number of estimable parameters of model  $M$
- Given a dataset, models can be ranked according to their AIC
- The model with the lowest AIC is selected
- $\text{AIC}_c$  – correction for finite sample size – usually used

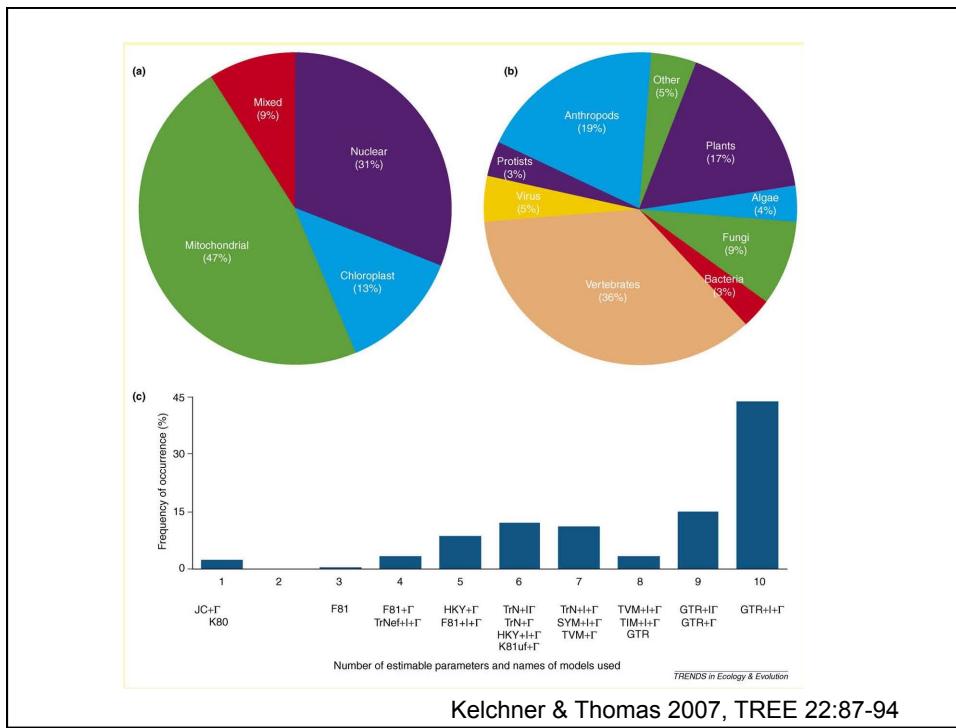
## Bayesian Information Criterion

- BIC also takes into account sample size  $n$
- $BIC(M) = -2 \cdot \text{Log}(\text{Likelihood}(M)) + K(M) \cdot \text{Log}(n)$ 
  - $K(M)$  is the number of estimable parameters of model  $M$  and  $n$  is the number of characters
- The model with the lowest BIC is selected

## 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>

Output from jModelTest						
* CORRECTED AKAIKE INFORMATION CRITERION (AICc)						
* -lnL = 124050.4448						
-----						
Sample size: 7705.0						
Model selected:						
Model = GTR+I+G						
partition = 012345						
K = 299						
freqA = 0.3118						
freqC = 0.1859						
freqG = 0.1721						
freqT = 0.3302						
R(a) [AC] = 1.8584						
R(b) [AG] = 8.1223						
R(c) [AT] = 3.8422						
R(d) [CG] = 1.9446						
R(e) [CT] = 13.2353						
R(f) [GT] = 1.0000						
p-inv = 0.5420						
gamma shape = 1.0570						
-----						
<b>Best model</b>						
Model	-lnL	K	AICc	delta	weight	cumWeight
GTR+I+G	124050.44479	299	248723.116454	0.000000	1.000000	1.000000
HKY+I+G	124507.54811	295	249628.667552	905.551098	2.30e-197	1.000000
SYM+I+G	124637.30469	296	249890.343721	1167.227268	0.00e+000	1.000000
GTR+G	125146.28705	298	250912.636212	2189.519758	0.00e+000	1.000000
SYM+G	125604.64648	295	251822.864292	3099.747838	0.00e+000	1.000000
HKY+G	125827.80977	294	252267.028447	3543.911993	0.00e+000	1.000000
K80+I+G	126897.82784	292	254402.741487	5679.625033	0.00e+000	1.000000
K80+G	127883.70157	291	256372.328272	7649.211818	0.00e+000	1.000000
GTR+I	128411.14327	298	257442.348652	8719.232198	0.00e+000	1.000000
SYM+I	128662.58435	295	257938.740032	9215.623578	0.00e+000	1.000000
HKY+I	129622.81043	294	259857.029767	11133.913313	0.00e+000	1.000000
F81+I+G	129925.12106	294	260461.651027	11738.534573	0.00e+000	1.000000
F81+G	130912.74863	293	262434.744325	13711.627871	0.00e+000	1.000000
K80+I	131130.71024	291	262866.345612	14143.229158	0.00e+000	1.000000
JC+I+G	131880.87114	291	264366.667412	15643.550958	0.00e+000	1.000000
JC+G	132773.35353	296	266149.472099	17426.355645	0.00e+000	1.000000
F81+I	143289.73580	293	287188.718665	38465.602211	0.00e+000	1.000000
JC+I	144715.21061	296	290033.186259	41310.069805	0.00e+000	1.000000
GTR	146171.15092	297	292960.199774	44237.083321	0.00e+000	1.000000
SYM	146261.55346	294	293134.515827	44411.399373	0.00e+000	1.000000
HKY	148503.72377	293	297616.694605	48893.578151	0.00e+000	1.000000
K80	149762.98834	290	300128.741719	51405.625265	0.00e+000	1.000000
F81	155229.75311	292	311066.592027	62343.475573	0.00e+000	1.000000
JC	156259.70450	289	313120.014529	64396.898076	0.00e+000	1.000000



## 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 (covered in tutorials)

## 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 (covered in tutorials)

### Output from PartitionFinder

**STRATEGY 1**

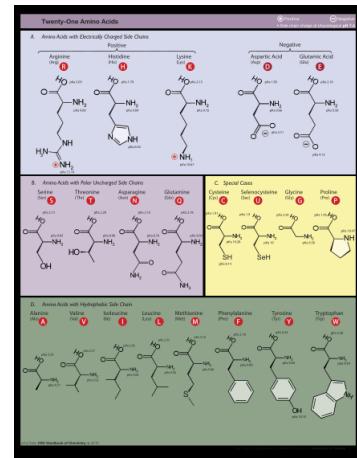
Best partitioning scheme			
Scheme Name	:	step_10	
Scheme lnL	:	-118756.435840	
Scheme BIC	:	239483.118009	BIC score
Number of params	:	352	
Number of sites	:	7705	
Number of subsets	:	8	
Subset	Best Model	# sites	subset id
1	GTR+G	714	d117c135876d3e868828f25d09953a9f
2	F81+I	4290	4a675e7e540e0cb009621eb52d7552b78
3	GTR+G	440	5d196738fa5467fb2f903ee5b3d1bb3
4	GTR+G	637	065fff28d8fb017945d1d664ee88c2e
5	SYM	252	428b6c6e493caaddb383eeaddffda1f08
6	SYM+G	413	c4bdbe6a705db90b2cd3ea803712b39
7	HKY+G	228	88f6044a648ae5cebd0d58a60b0ec5c5
8	GTR+G	731	69ca434e659726db47f5700ecf6cdee5

**STRATEGY 2**

Best partitioning scheme			
Scheme Name	:	separate	
Scheme lnL	:	-119493.545898	
Scheme BIC	:	242206.807081	BIC score comparison
Number of params	:	359	Strategy 1 BIC = 239 403 ✓
Number of sites	:	7705	Strategy 2 BIC = 242 200
Number of subsets	:	7	
Subset	Best Model	# sites	subset id
1	GTR+G	408	88d9e408a7429cf1f7c04b3c6ad405b
2	GTR+I+G	4040	6147b066b0346034e59c738ca8abab6
3	GTR+G	407	2078ac99fe345782a3e9b948ae6783
4	GTR+G	342	a9de0e1cf843a93e4545a0e53c47027b
5	GTR+G	733	0e945f1386af47731957496eda7faeef
6	GTR+G	368	067651d33be225f5683a0ee6ce036017
7	GTR+G	430	6ecc87f27a392fcae101b8b0c99c64e0

# 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](https://doi.org/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

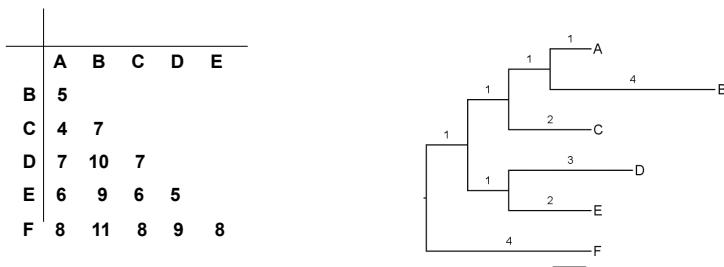
# Inferring phylogenies

- **Distance methods**
  - A clustering method using pairwise distances between sequences (neighbour joining)
- **Discrete characters**
  - Using an optimality criterion to choose the best tree
    - Maximum parsimony (Ockham's razor)
      - Best explanation is the simplest one (the one that minimizes the number of substitutions)
      - Doesn't perform as well as model-based methods on molecular data
      - Still used for morphological characters
    - Maximum likelihood
    - Bayesian inference

## Distance methods

## Distance methods involve two stages

- Stage 1: calculate the evolutionary distance between pairs of sequences (using a DNA substitution model)
- Stage 2: use the distances to construct a tree that describes those evolutionary distances



## Distance Methods

**Distance Estimates:** estimation of the divergence between two sequences deriving from a common ancestor.

- it is a measure of (dis)similarity between sequences
- branch lengths are proportional to the distance
- if we assume a molecular clock the distance is directly proportional to time

**Distance can be expressed as a proportion of sites that differ between two sequences:**

98 base pairs (bp), 15 bp differ, or D = 15/98 = 0.153 or 15.3%

```

Antirrhoea109 A GAA TGT TTA TCCCCC CTTT TTT TAA TAT TTG TCA AGAGG TCTT AGTT GACT TAG CAA TTTT TCTT TTA ATTATAG TGG TA TTT TTT TAA
Araschnia39 A GAA AG TGT TCCCCC AC TTTT AT TAA TAT TTG TCA AGAGG TCTT AGTT GACT TAG CAA TTTT TCTT TTA ATTATAG TGG TA TTT TTT TAA

```

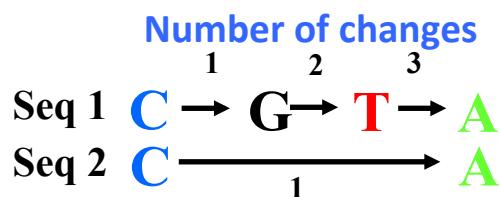
## Distance matrix:

	1	2	..	n
1				
2	0,33			
:				
n	0,23	0,63		

- > Direct measure of distance underestimates the true distance
- Remember multiple hits!

## Models correct for unobserved changes

- All models include a correction for multiple substitutions at the same site
  - All (except Logdet distances) can be modified to include a gamma correction for site rate heterogeneity (among site rate variation)



**Distance can be expressed as a proportion of sites that differ between two sequences:**

Antirrhinum109	A	T	G	A	T	C	T	T	T	A	T	G	A	T	T	T	A	A	T	A	T	G	G	A	T	T	T	A
Araschnia39	A	G	A	G	T	T	C	C	C	A	T	T	A	A	T	G	A	G	A	G	A	T	G	A	T	T	T	A
Astochiopus29	A	T	G	A	T	T	A	T	T	A	T	G	A	T	T	T	A	T	A	T	G	A	T	T	T	T	A	
Astrocopus82	A	G	A	G	T	T	C	C	C	A	T	T	A	A	T	G	A	G	A	G	A	T	G	A	T	T	T	A
Caligo70 10	T	G	A	G	T	T	C	C	C	A	T	T	A	A	T	G	A	G	A	G	A	T	G	A	T	T	T	A
Calinaga64 3	T	G	A	G	T	T	C	C	C	A	T	T	A	A	T	G	A	G	A	G	A	T	G	A	T	T	T	A
Castilla76 2	T	G	A	G	T	T	C	C	C	A	T	T	A	A	T	G	A	G	A	G	A	T	G	A	T	T	T	A
Catacropte88	A	G	A	G	T	T	C	C	C	A	T	T	A	A	T	G	A	G	A	G	A	T	G	A	T	T	T	A
Catonephele6	A	G	A	G	T	T	C	C	C	A	T	T	A	A	T	G	A	G	A	G	A	T	G	A	T	T	T	A
Cercyonis8 1	A	G	A	G	T	T	C	C	C	A	T	T	A	A	T	G	A	G	A	G	A	T	G	A	T	T	T	A
Chesoniad1	G	G	A	G	T	T	C	C	C	A	T	T	A	A	T	G	A	G	A	G	A	T	G	A	T	T	T	A
Chlosyne62 1	T	G	A	G	T	T	C	C	C	A	T	T	A	A	T	G	A	G	A	G	A	T	G	A	T	T	T	A
Clossiana76	T	G	A	G	T	T	C	C	C	A	T	T	A	A	T	G	A	G	A	G	A	T	G	A	T	T	T	A
Colobura68 1	T	G	A	G	T	T	C	C	C	A	T	T	A	A	T	G	A	G	A	G	A	T	G	A	T	T	T	A

Dissimilarities matrix:

	1	2	..	n
1				
2	0.33			
:				
n	0.23	0.63		

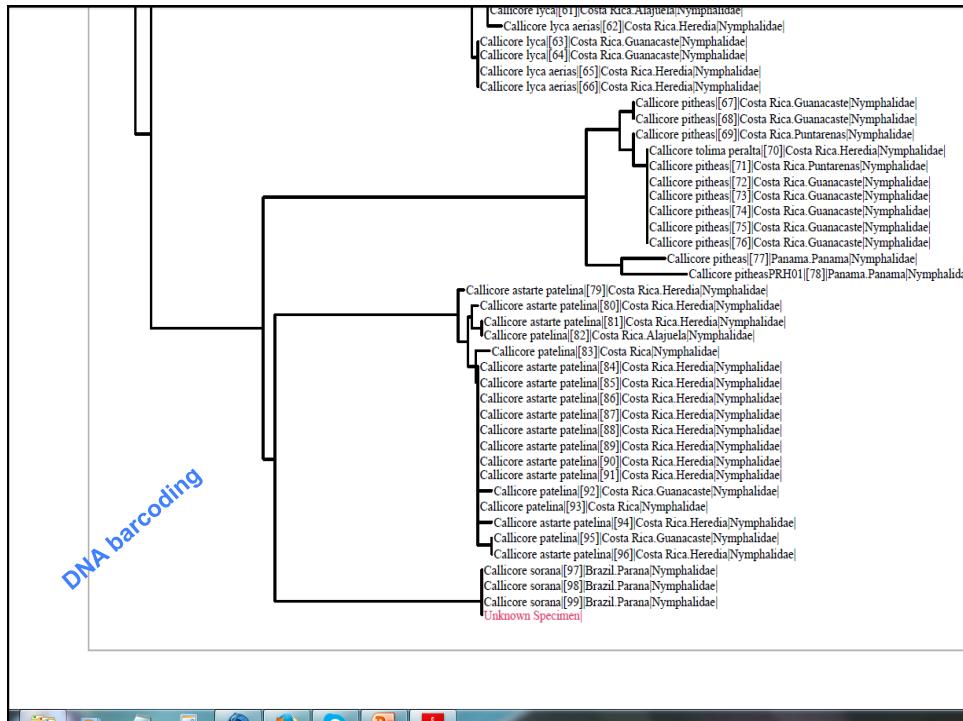
Evolutionary distance matrix:

	1	2	..	n
1				
2	0.35			
:				
n	0.24	0.66		

*Correction for multiple substitutions* →

## Distances - advantages

- Computationally fast
- A large number of models are available with many parameters - improves estimation of distances
- Great for getting a quick tree in data checking/exploration phase



## Distances – disadvantages (1/2)

- Prone to systematic errors
- Problems with missing data
- Rate variations in different parts of a tree are intractable for distance measures
- Information on variation in characters is lost once sequence differences are converted to distances

## **Distances – disadvantages (2/2)**

- **Generally outperformed by Maximum Likelihood methods in choosing the correct tree in computer simulations**
  - See e.g. Ogden & Rosenberg (2006) *Multiple Sequence Alignment Accuracy and Phylogenetic Inference*. *Syst. Biol.* 55(2): 314–328 (DOI: [10.1080/10635150500541730](https://doi.org/10.1080/10635150500541730))