

# Molecular Phylogenetics Course

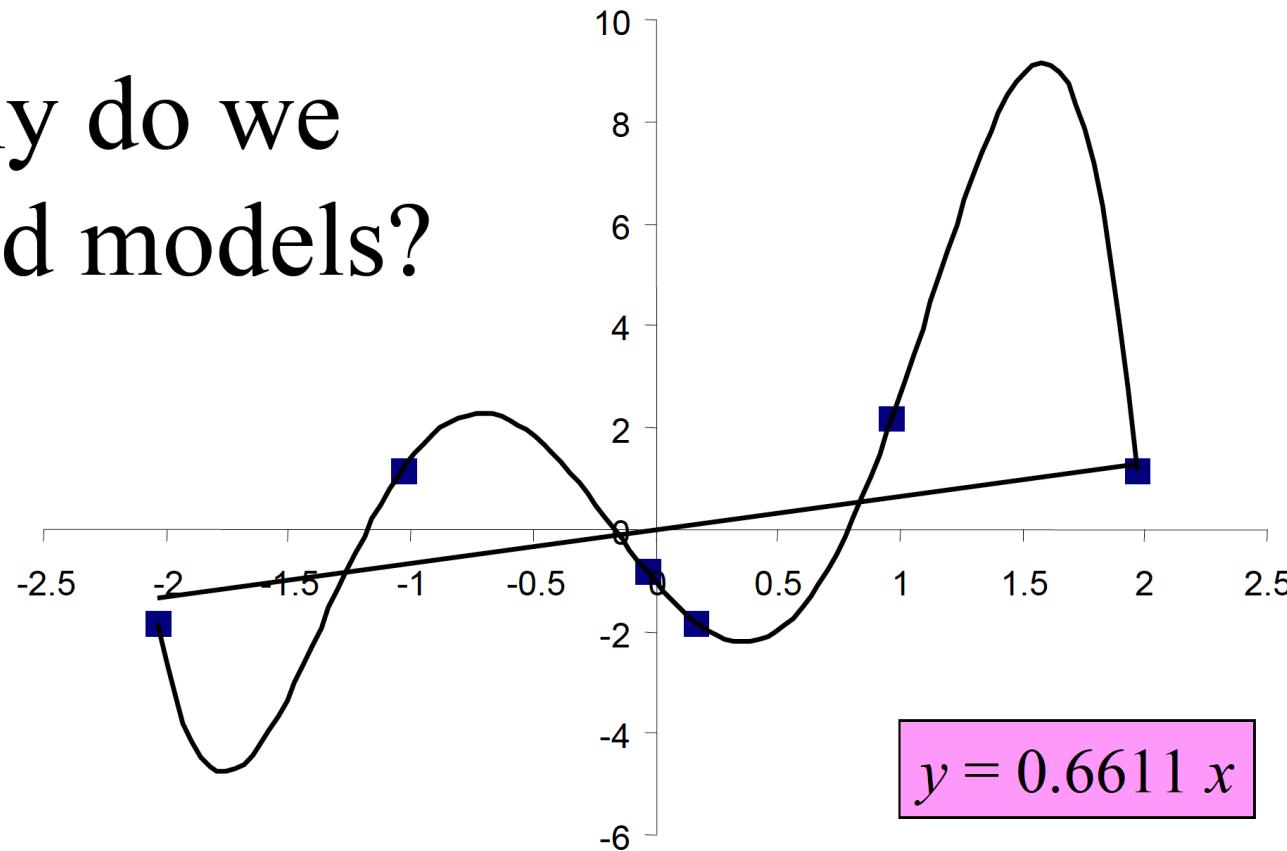
## Introduction to model-based methods

Jadranka Rota

Some slides by 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?

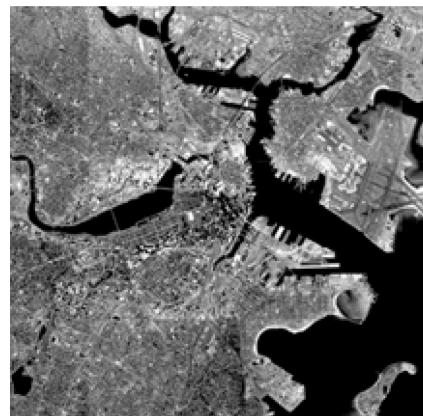


## A very *practical* MBTA subway map



17

## A very *realistic* MBTA subway map



18

▶ Which is more useful?

# 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

# Modelling evolution of DNA sequences

- 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

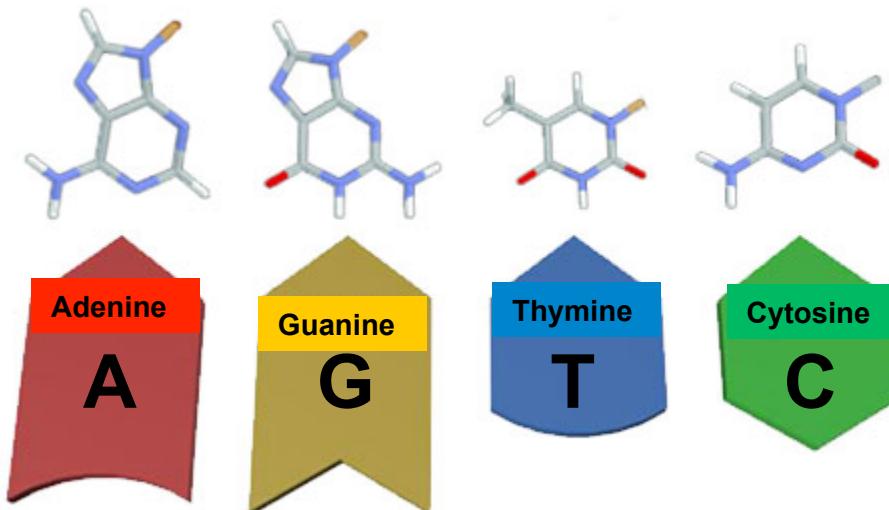
# Main Challenge

- ▶ DNA has only four characters

Purines

Pyrimidines

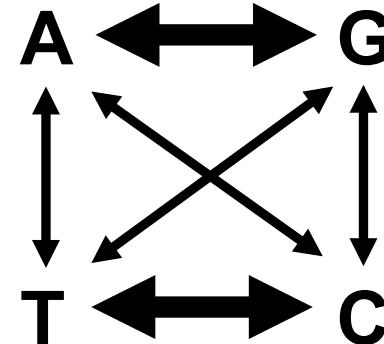
Figure B-3: The Four Nitrogenous Bases



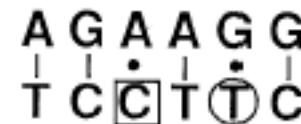
Each base has a distinct shape that can be used to distinguish it from the others.  
3D representations of the four bases are shown, with the corresponding chemical structures drawn above.

# Substitution types

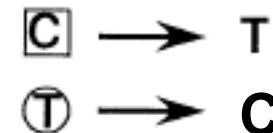
- ▶ Purines: A, G
- ▶ Pyrimidines: C, T
- ▶ Transversions
  - Pu → Pyr
  - Pyr → Pu
- ▶ Transitions – more common
  - Pu → Pu
  - Pyr → Pyr



Pur - Pyr mispairs lead to transitions



In next round of replication

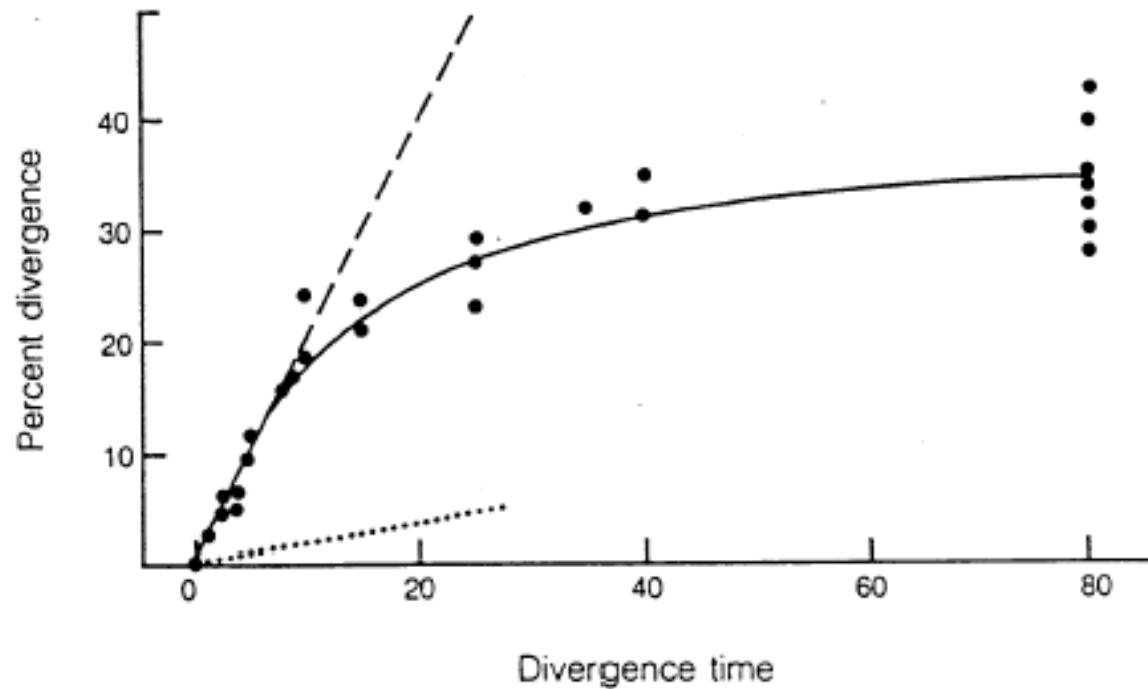


# Saturation in sequence data:

- 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

# Misleading DNA evolution

Multiple substitutions hide previous changes



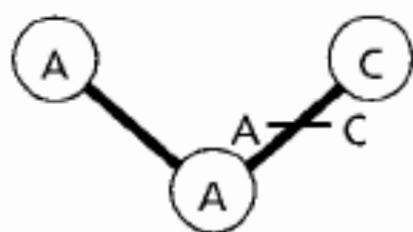
# Difference between mutation and substitution

- **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

# Types of Substitutions

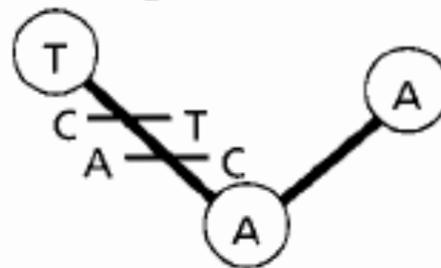
(a) Single substitution

1 change, 1 difference



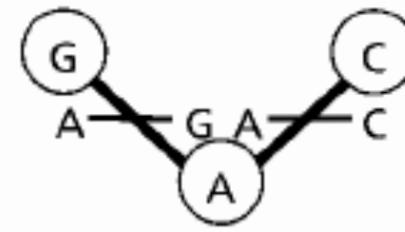
(b) Multiple substitution

2 changes, 1 difference



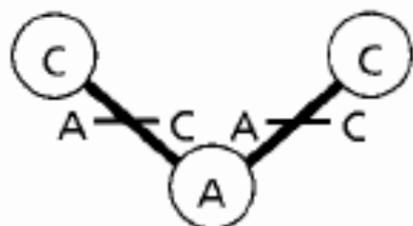
(c) Coincidental substitution

2 changes, 1 difference



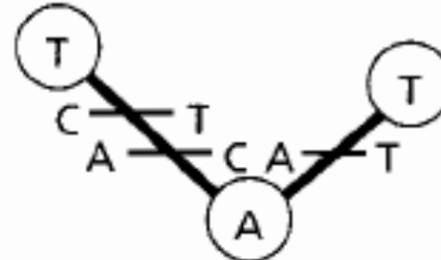
(d) Parallel substitution

2 changes, no difference



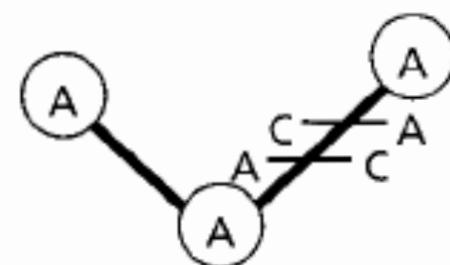
(e) Convergent substitution

3 changes, no difference



(f) Back substitution

2 changes, no difference



# Modelling evolution

- These dynamics can be modeled over a tree
- 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 - often modelled as a statistical distribution - for example a gamma distribution

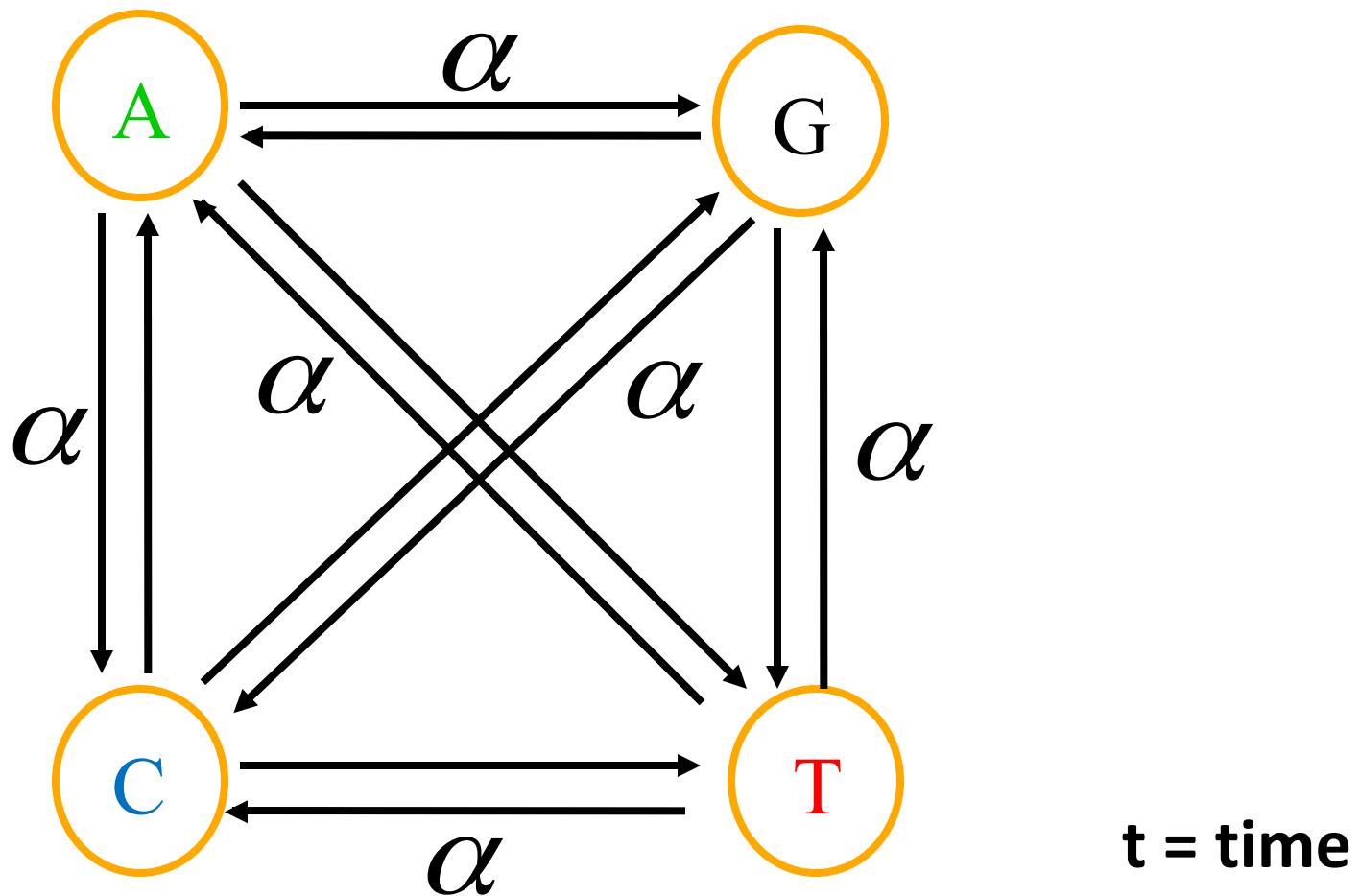
# Corrections for multiple substitutions

**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

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

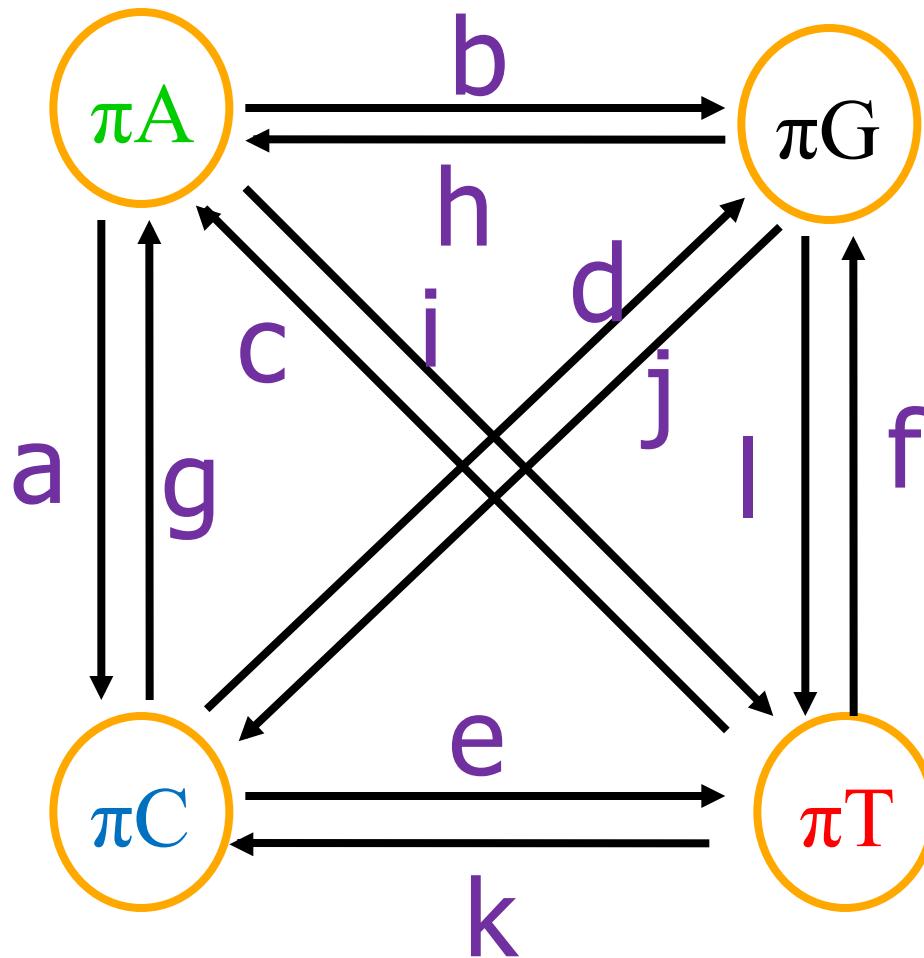
# Improvements on Jukes-Cantor

- 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**

# 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 and branch lengths**

# A general model of sequence evolution



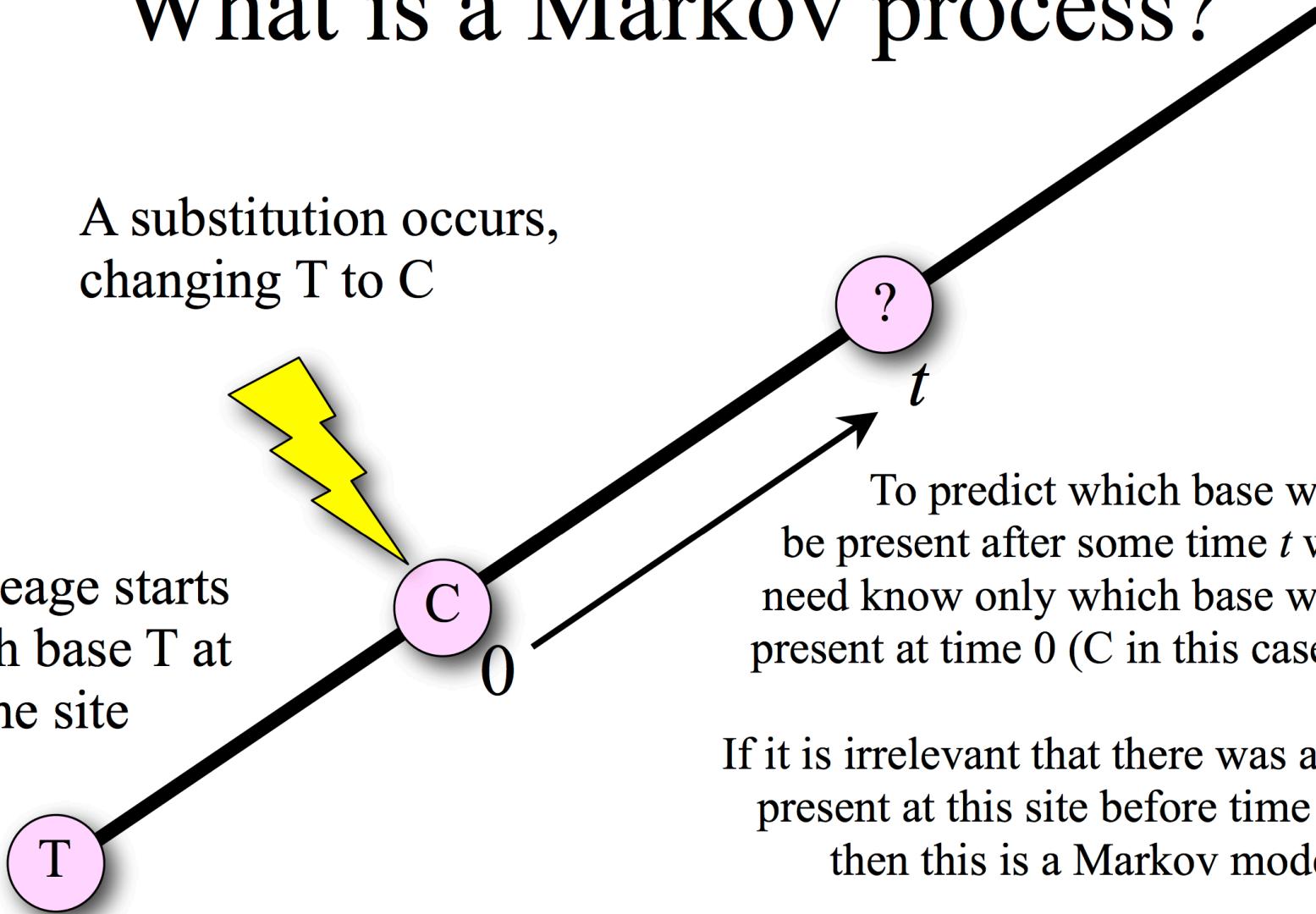
# Time-homogenous time-continuous stationary Markov models

- 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)

# What is a Markov process?

A substitution occurs,  
changing T to C

Lineage starts  
with base T at  
some site



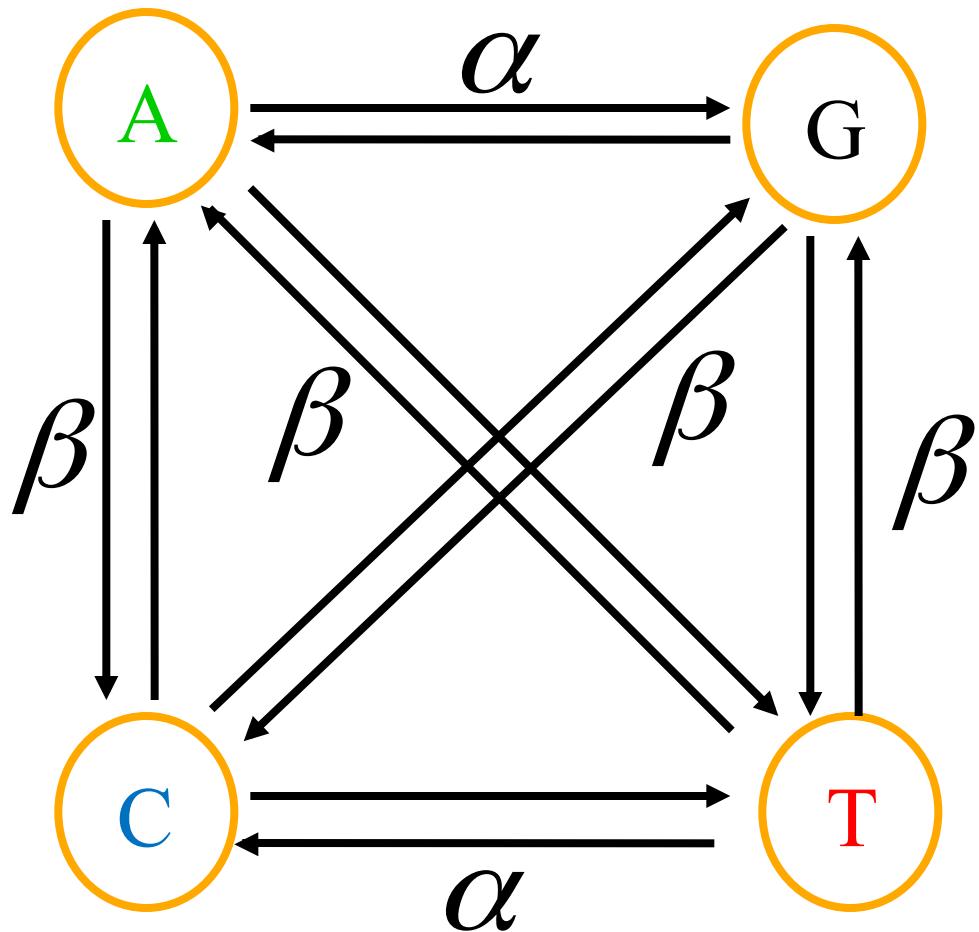
To predict which base will be present after some time  $t$  we need know only which base was present at time 0 (C in this case).

If it is irrelevant that there was a T present at this site before time 0, then this is a Markov model.

# 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)

# 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 -  $\beta$ )

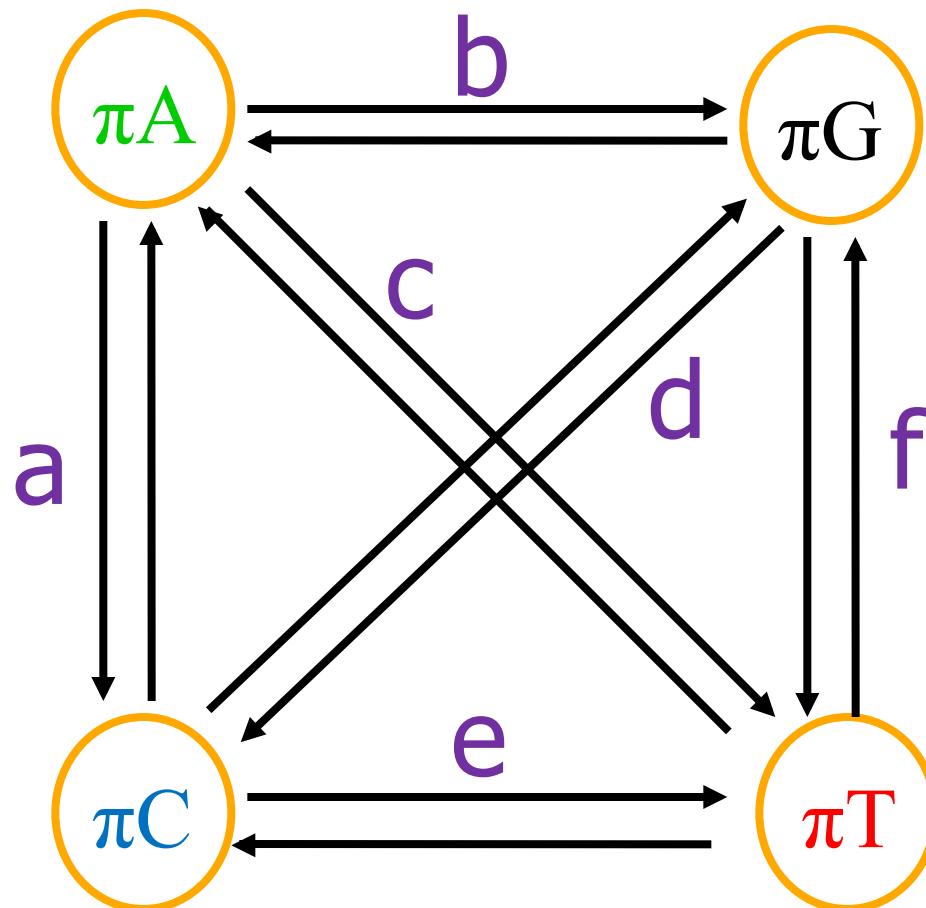
# The Hasegawa-Kishino-Yano model

	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	$\pi_A c$	$\pi_C e$	$\pi_G f$	—

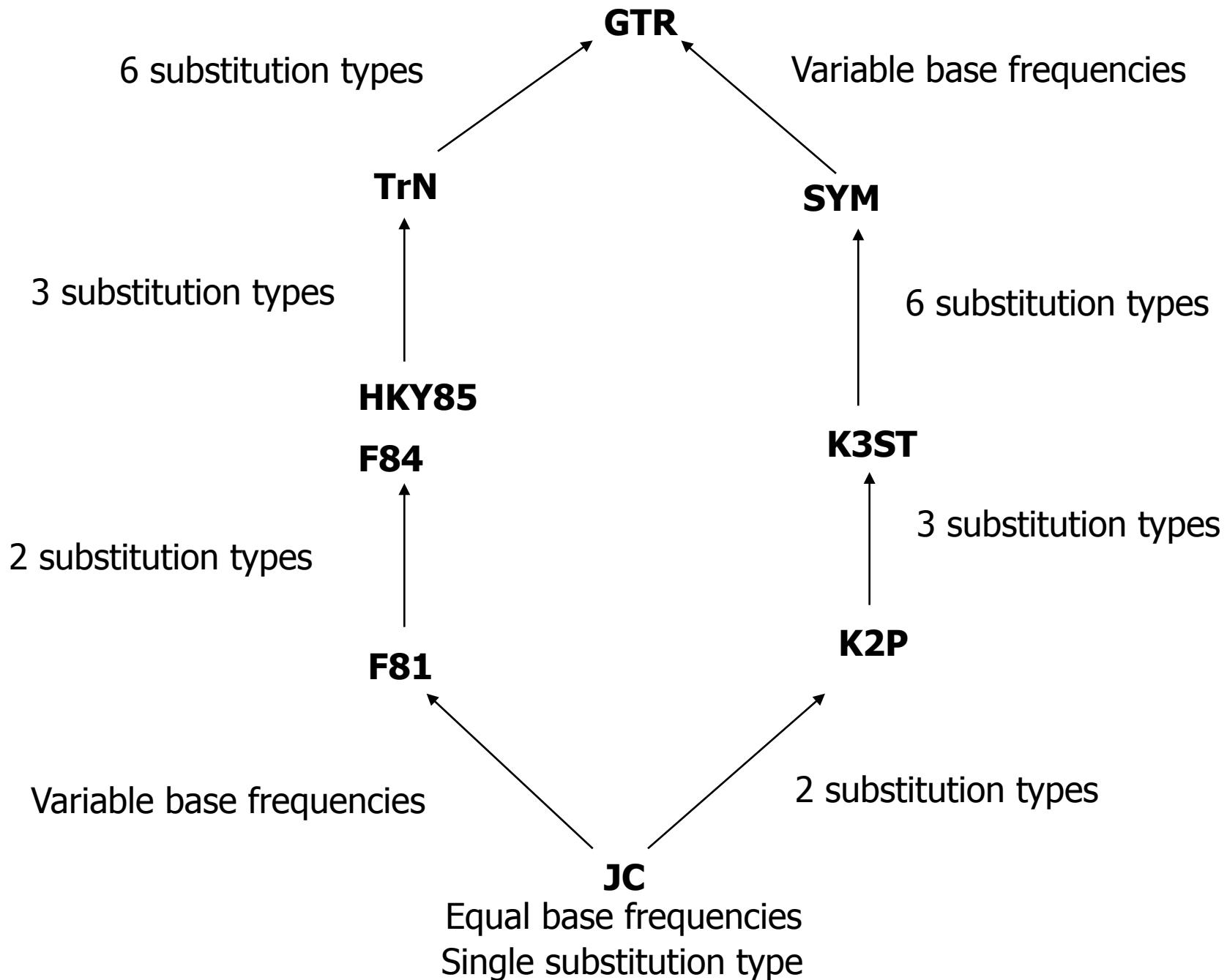
GTRmodel:

- 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

ACAGGTGAGGGCTCAGCCAATTGAGCTTTGTCGATAGGT



# 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)

- 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.

# Third codon position not completely degenerate: synonymous and non-synonymous changes

	AGA								UUA	
	AGG								UUG	
GCA	CGA						GGA		AUA	CUA
GCC	CGC						GGC		AUC	CUC
GCG	CGG	GAC	AAC	UGC	GAA	CAA	GGG	CAC	CUG	AAA
GCU	CGU	GAU	AAU	UGU	GAG	CAG	GGU	CAU	CUU	AAG
Ala	Arg	Asp	Asn	Cys	Glu	Gln	Gly	His	Ile	Leu
										Lys

		AGC								
		AGU								
		CCA	UCA	ACA			GUA			
		CCC	UCC	ACC			GUC	UAA		
	UUC	CCG	UCG	ACG			UAC	GUG	UAG	
AUG	UUU	CCU	UCU	ACU	UGG	UAU	GUU	UGA		
Met	Phe	Pro	Ser	Thr	Trp	Tyr	Val	Stop		

Different probability of change in 2 vs. 3 vs. 4 vs. 6-codon families; “Universal” Genetic Code (nuclear)

# Typical pattern of variation among codon positions

E.g. in Collembola

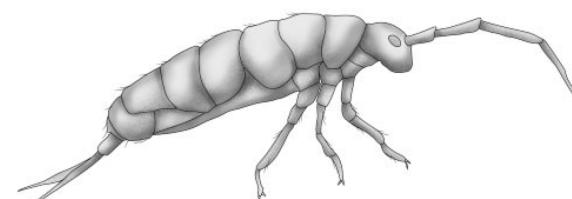
**56.7% of all variable sites are located in third positions**

**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. JME**



BioEdit Sequence Alignment Editor - [C:\Documents and Settings\Koti\My Documents\Työjutut\Rawdata\Unchecked\NymphalidaeCOI.fst]

File Edit Sequence Alignment View World Wide Web Accessory Application RNA Options Window Help



Courier New

11

55 total sequences

shade threshold 40 %

Mode: Edit

Overwrite ▾ Selection: 0  
Position: 34

Sequence Mask: None  
Numbering Mask: None

Start  
ruler at:



start



Page



## Molecular methods to

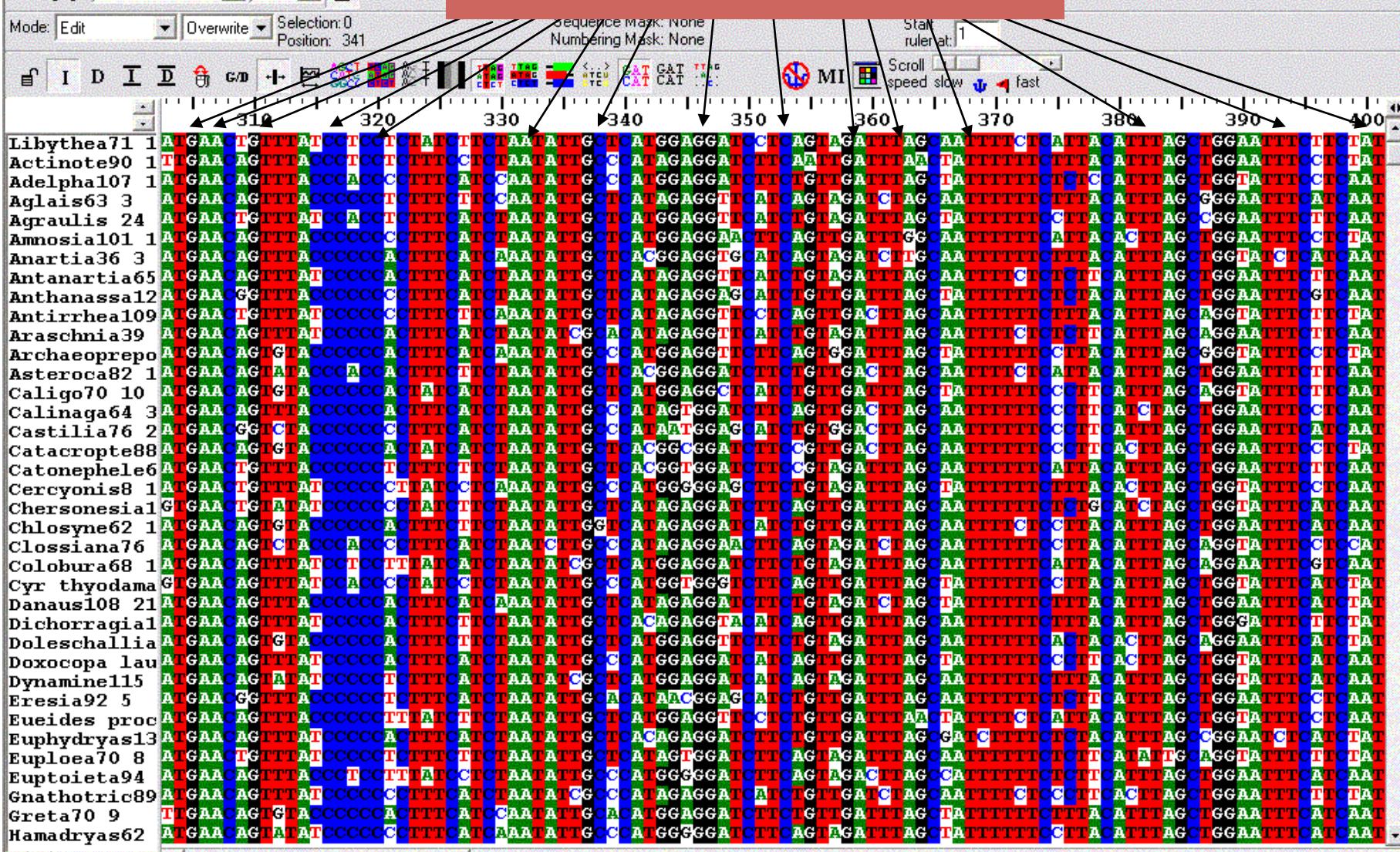


 BioEdit Sequence Align.



20:32

# Invariable sites



# Modelling among-site rate variation (ASRV)

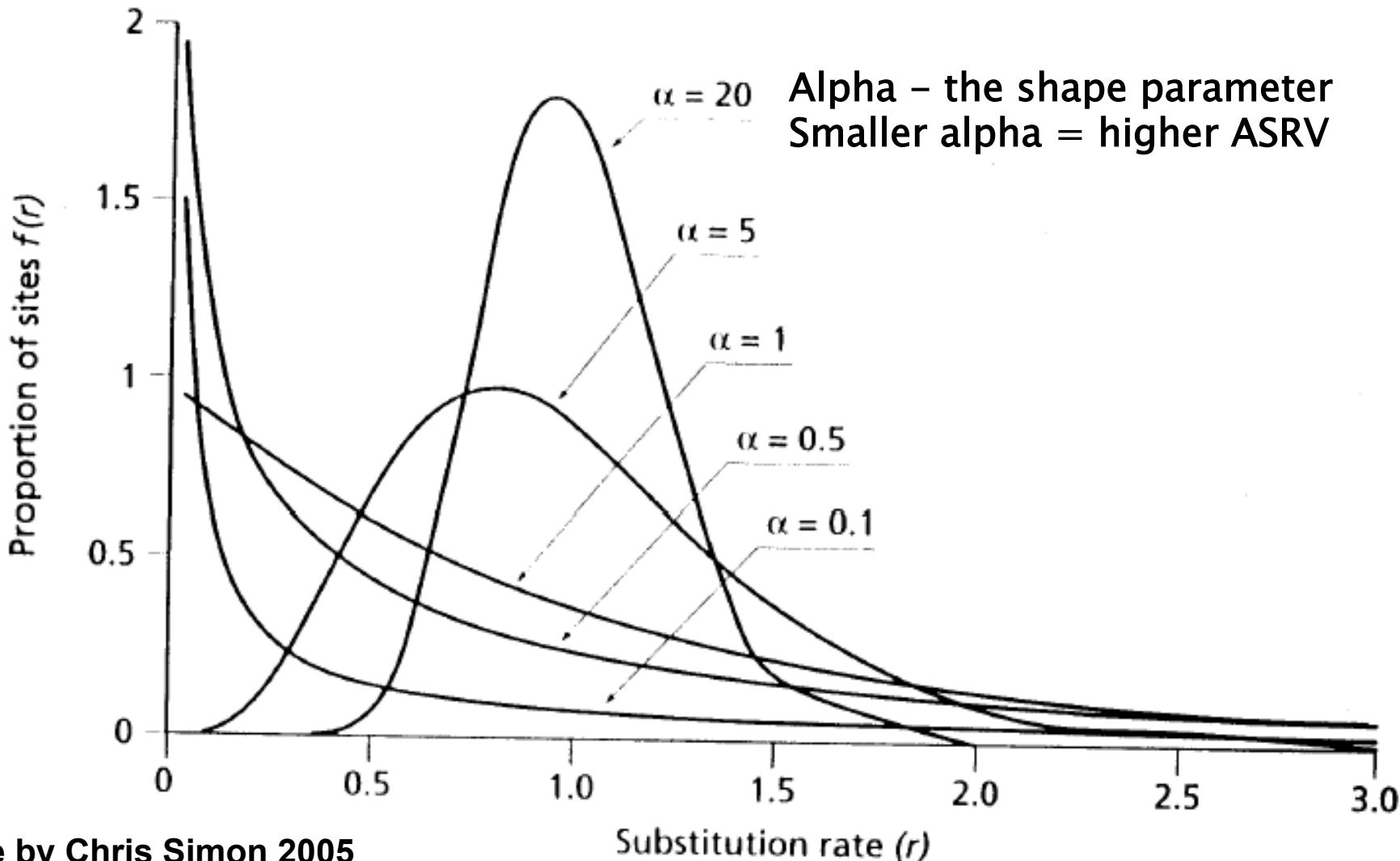
- **The most common additional parameters are:**
  - A correction for the proportion of sites which are **invariable** (parameter  $I$ )
  - A correction for variable site rates at those sites which can change (parameter gamma,  $G$ )
- All models can be supplemented with these parameters (e.g. GTR+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

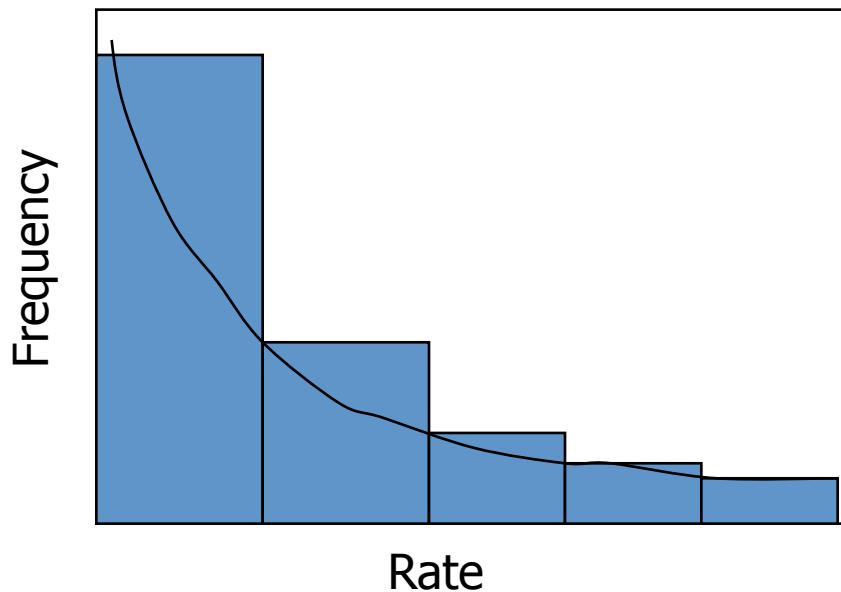
# Gamma distribution:

## Relative substitution rates for different $\alpha$ values



# Gamma distribution computationally costly

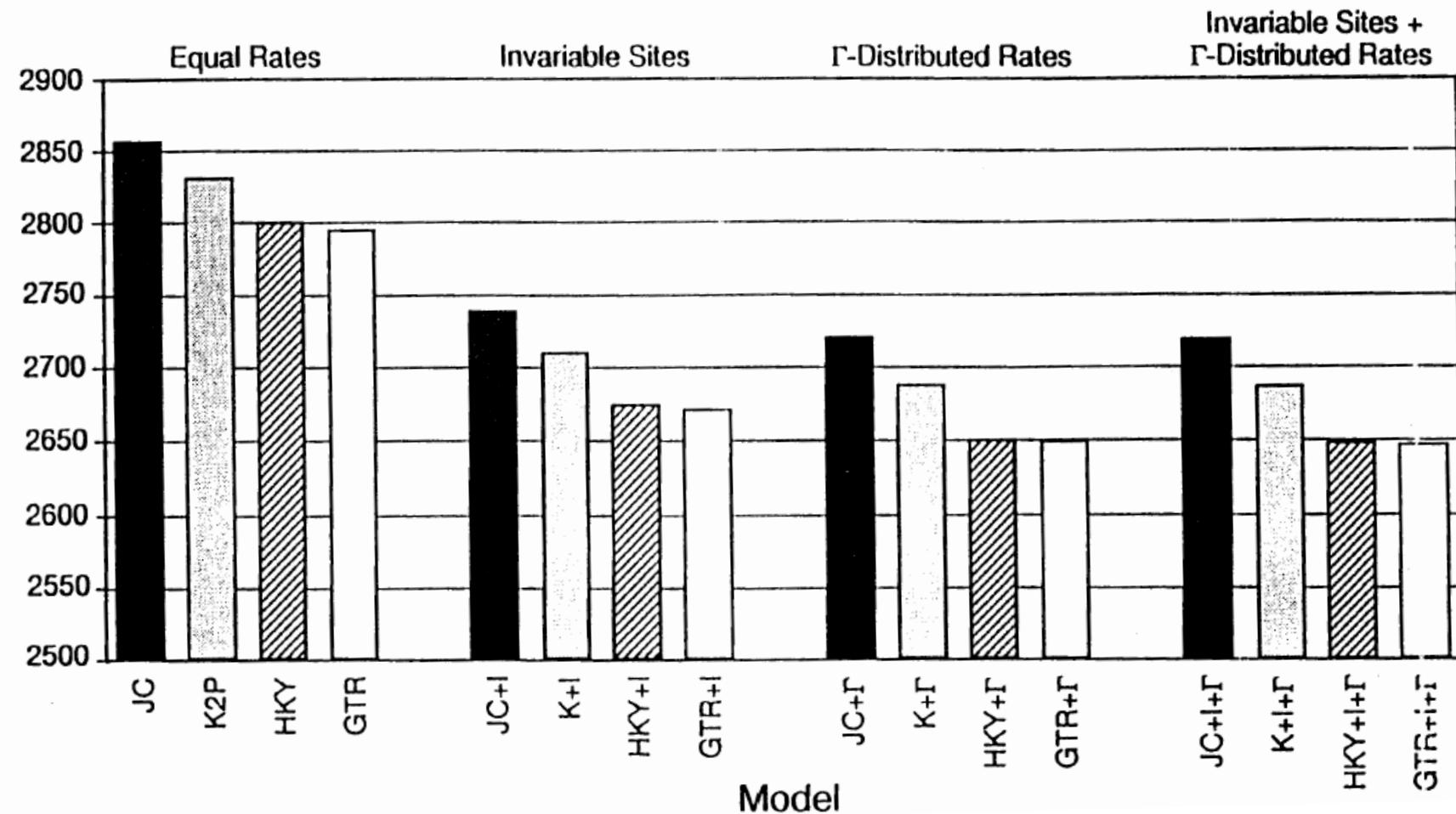
- Computational difficulties in using continuous distribution
- Most programs use discrete categories



# ASRV: Yang discrete model

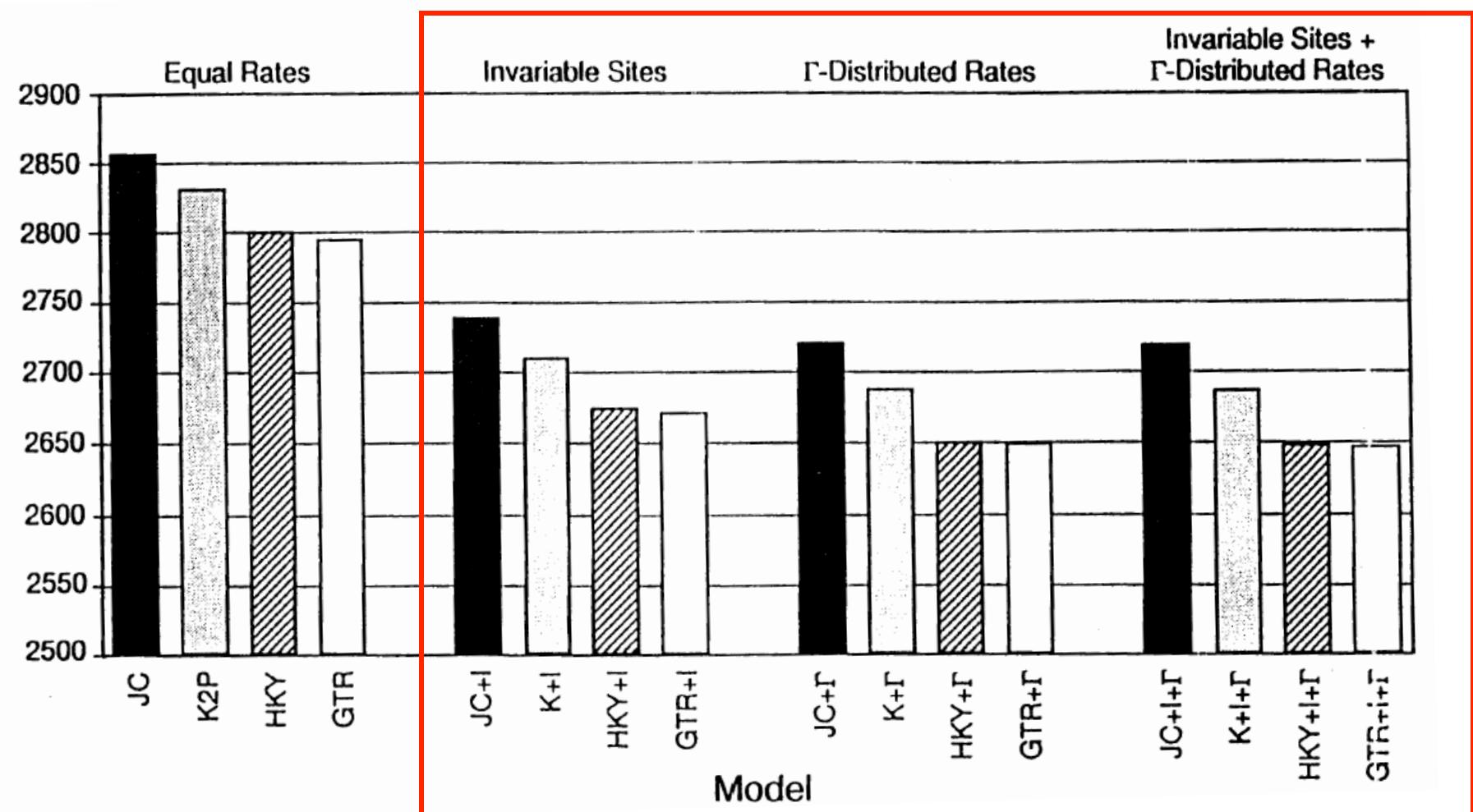
- Continuous data divided into “n” discrete rate classes (generally 4)
- If  $\alpha < 0.2$  Yang recommends more rate classes
- Less computer intensive than obtaining likelihoods by integrating over the continuous gamma distribution

# ASRV >> fit improvement than by other parameters



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

# ASRV >> fit improvement than by other parameters



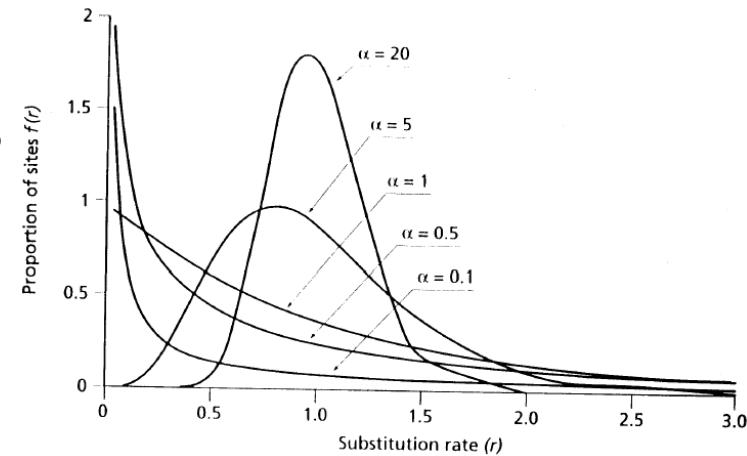
Frati, Simon, Sullivan, Swofford. 1997. JME 44:145-158

# 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
  - Implemented in IQ-TREE

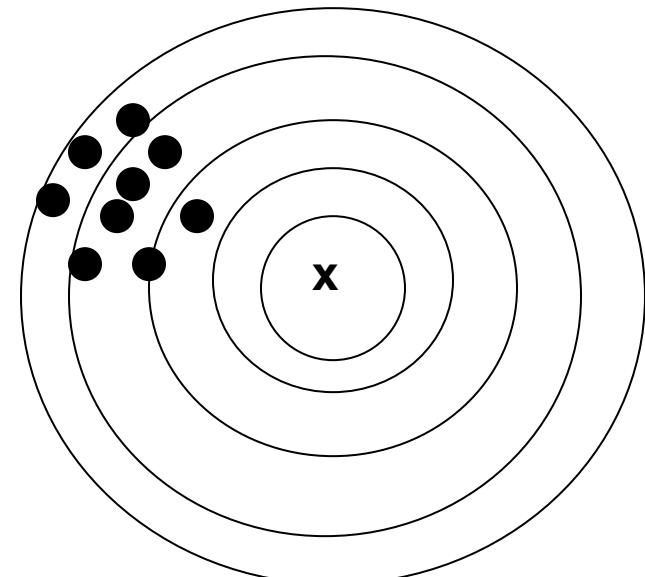
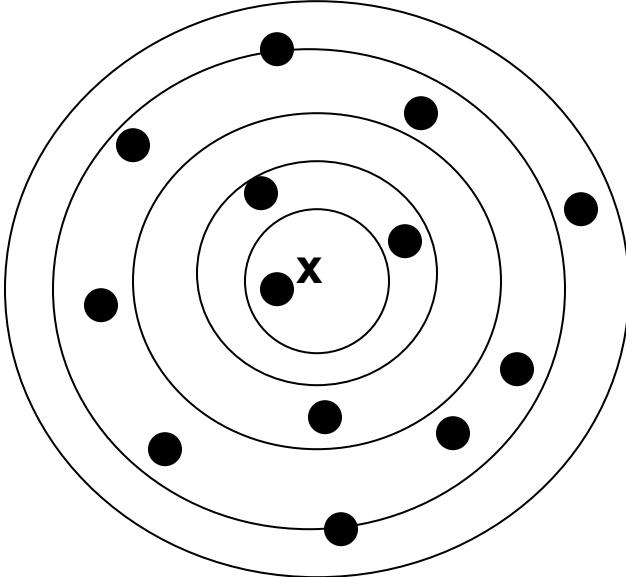


# Parameters in models of DNA evolution

- **Numbers of parameters estimated:**
  - Substitutions (up to 5; 1 fixed, 5 estim.)
  - Base composition (1 fixed, 3 estim.)
  - 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

# 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



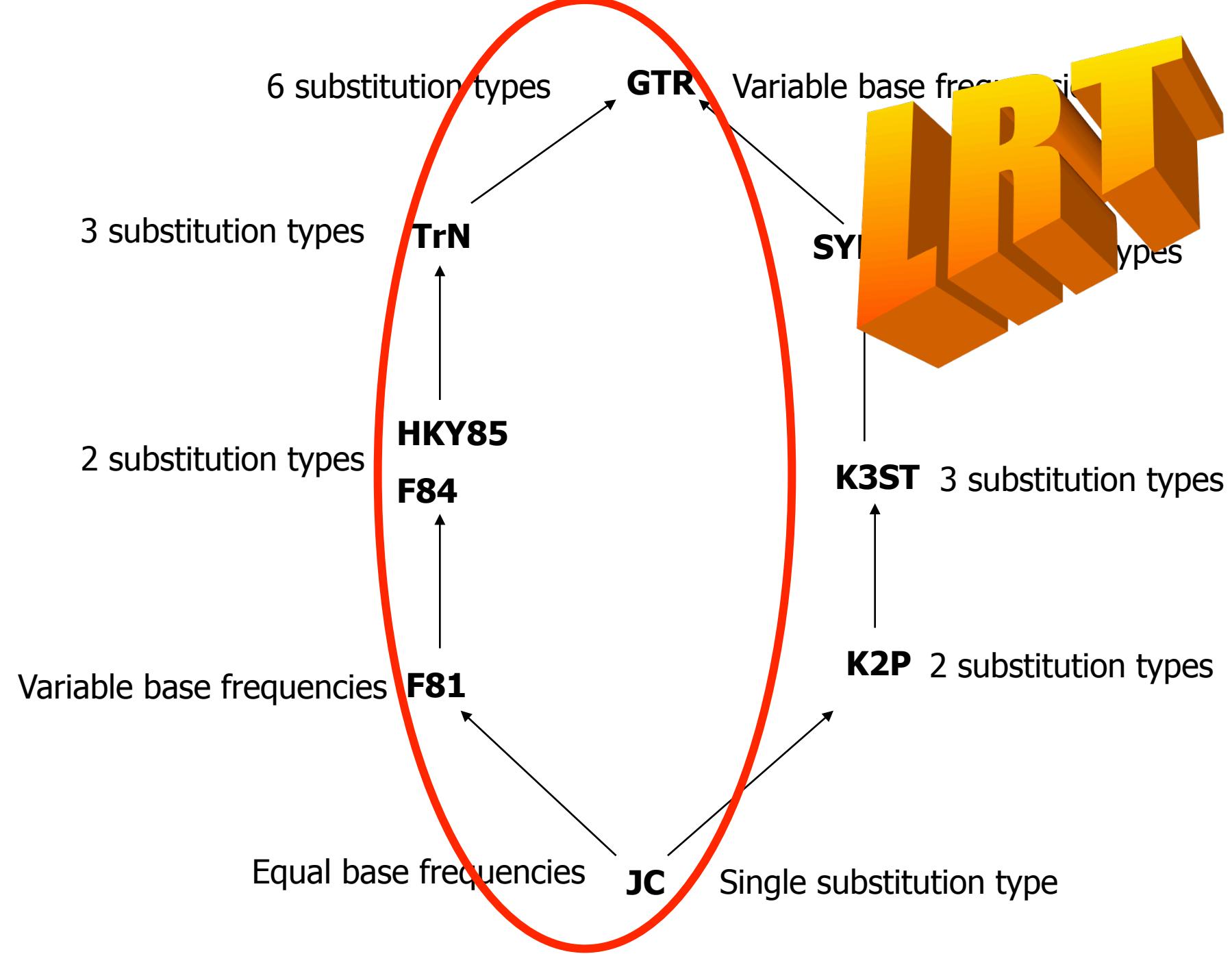
**Swofford's Target Analogy**

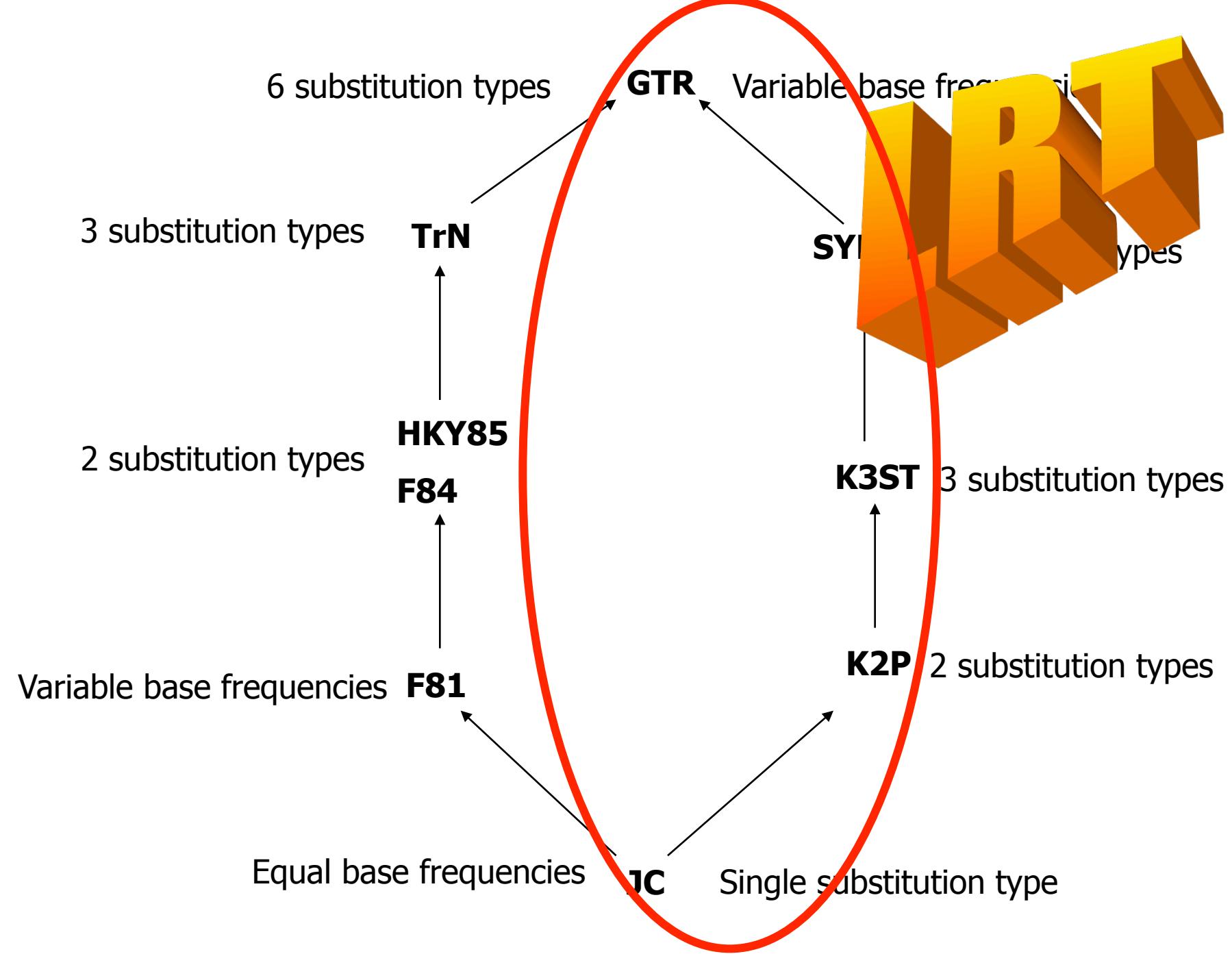
## **Trade-off between highly parameterized models & model error variance**

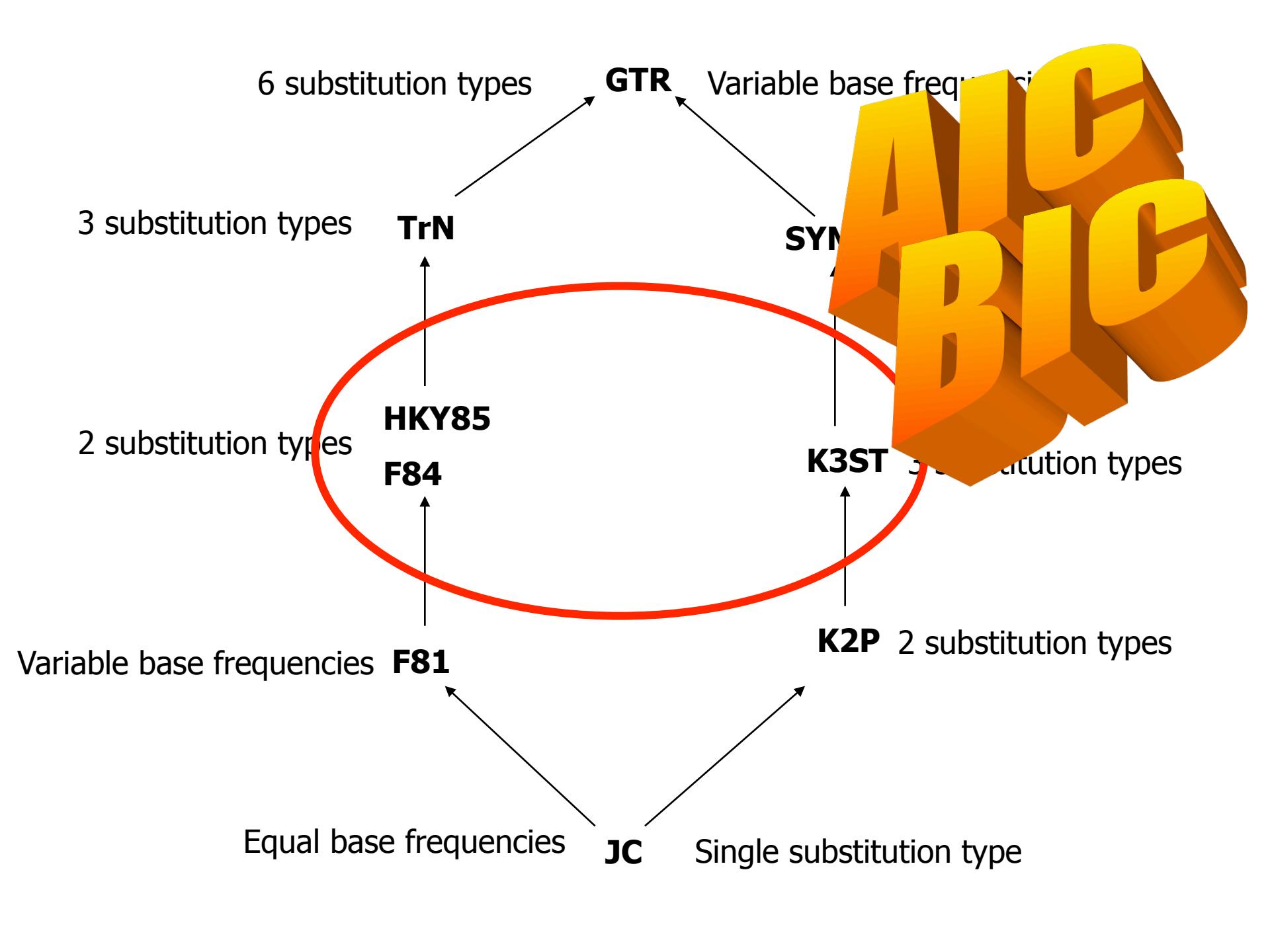
- **Many parameters, higher error variance but clustered around the true value (higher accuracy)**
- **Few parameters, lower error variance but may not be centered around the mean (lower accuracy)**

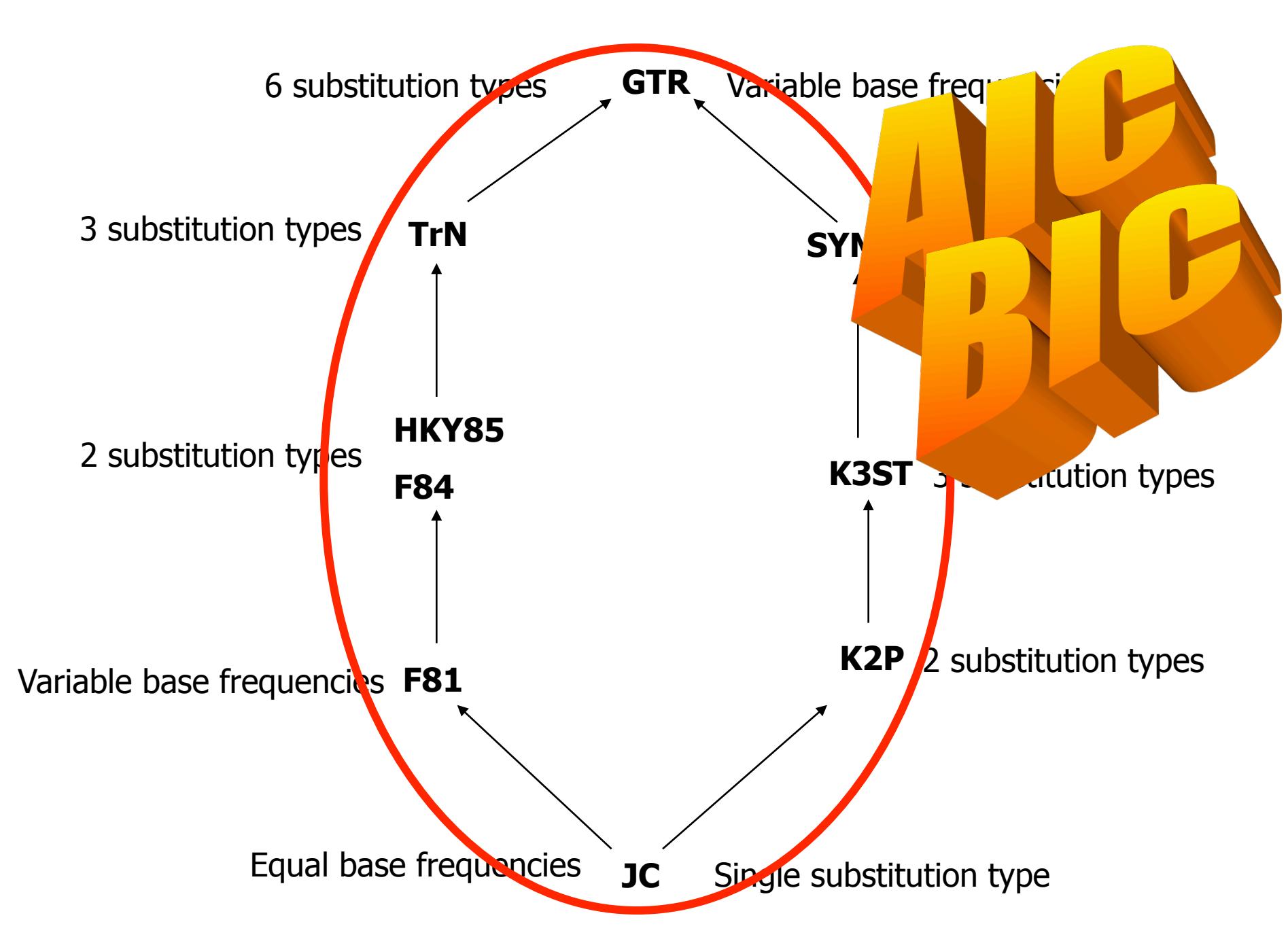
# Choosing between models

- 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 both tests, one needs to **compute the likelihood of the model**
- This will be covered in the 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
- $AIC(M) = -2 \cdot \text{Log(Likelihood}(M)) + 2 \cdot K(M)$ 
  - $K(M)$  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
- $AIC_c$  – correction for finite sample size – usually used

# Bayesian Information Criterion

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

# Modeltesting 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.

# Output from jModelTest

```
*          *
*      CORRECTED AKAIKE INFORMATION CRITERION (AICc) *
*
```

Sample size: 7705.0

Model selected:

Model = GTR+I+G

partition = 012345

-lnL = 124050.4448

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

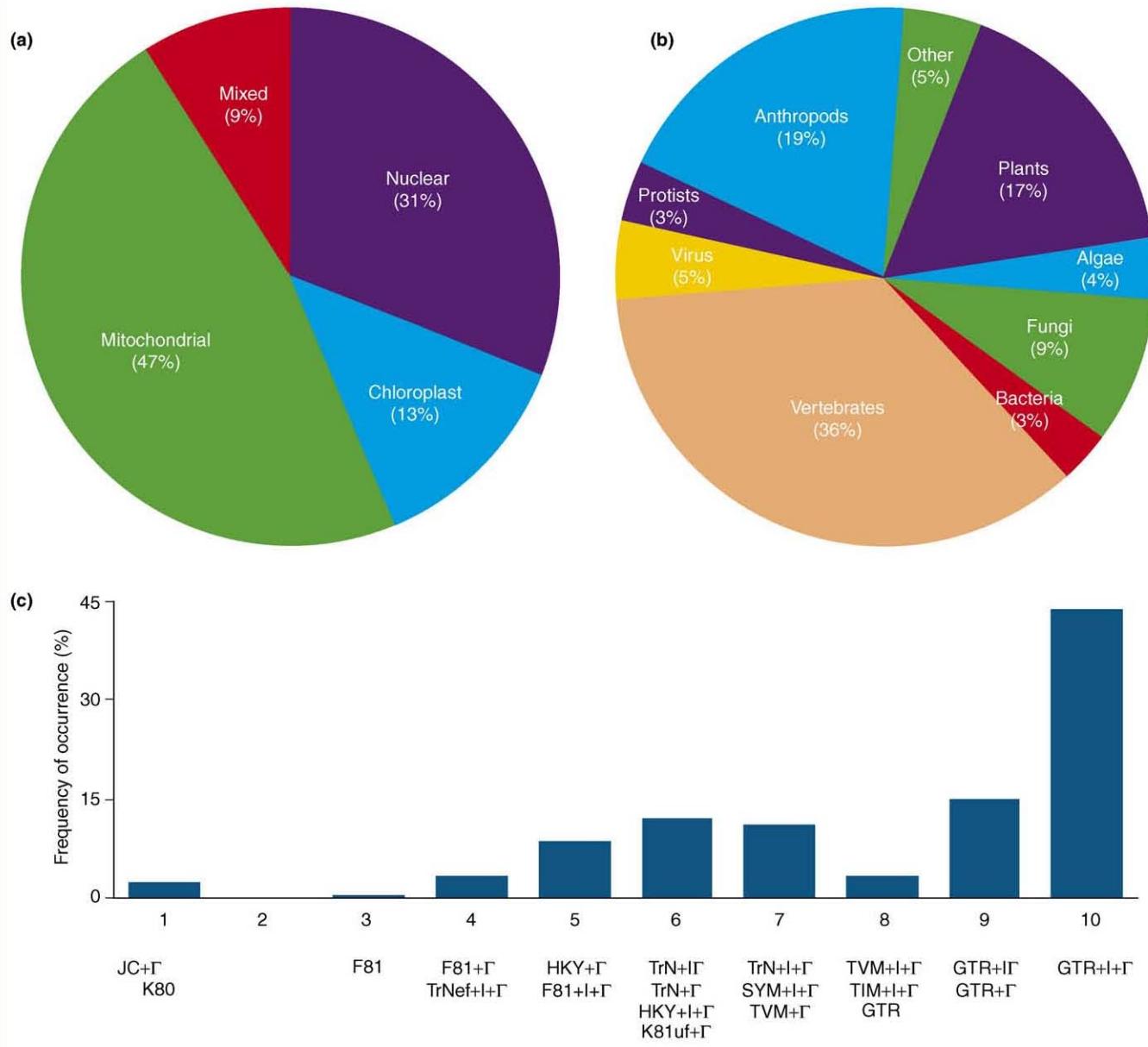
R(f) [GT] = 1.0000

p-inv = 0.5420

gamma shape = 1.0570

Best model

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	290	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	290	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 - ModelFind implemented

# Partitioned models

- **A priori separation of characters into different partitions**
  - E.g.: different codon position, RNA, introns/exons, etc.
- **Each partition analyzed with a different model**
- **Allows for heterogeneity across data subsets in overall rate and in substitution model parameters, plus in some programs also possible to unlink topology and branch lengths**
- **“Different data subsets can thus have independent branch lengths or even different topologies.”**  
**(Ronquist and Huelsenbeck, 2003:1573)**

# Output from PartitionFinder

Best partitioning scheme

```
Scheme Name      : step_10
Scheme lnL       : 119126.425049
Scheme BIC       : 239403.118009
Number of params : 352
Number of sites  : 7705
Number of subsets: 8
```

BIC score

Subset	Best Model	# sites	subset id	Partition names
1	GTR+G	714	d117c135876d3e868828f25d0953a9bd	Partition_11, Partition_14, Partition_8, Partition_12, Partition_10
2	F81+I	4290	4a675e7e540ecb009621eb52d7552b70	Partition_1
3	GTR+G	440	5d196738faf5467fb2f903ee5b3d1bb3	Partition_13, Partition_16, Partition_17, Partition_15
4	GTR+G	637	065fff28d8d1b017945d1d664ee88c2e	Partition_6, Partition_9, Partition_7
5	SYM	252	428b6c6e493caaddb383eeadddfdaf08	Partition_2
6	SYM+G	413	c4bdbe6a705db90b2dc3dea803712b39	Partition_4, Partition_5
7	HKY+G	228	88f604a648ae5dceb0d58a600bec5ca5	Partition_3
8	GTR+G	731	69ca434e659726db47f5700ecf6cdee5	Partition_18

Best model for each partition

STRATEGY 1

Best partitioning scheme

```
Scheme Name      : separate
Scheme lnL       : -119493.545898
Scheme BIC       : 242200.007081
Number of params : 359
Number of sites  : 7705
Number of subsets: 7
```

STRATEGY 2

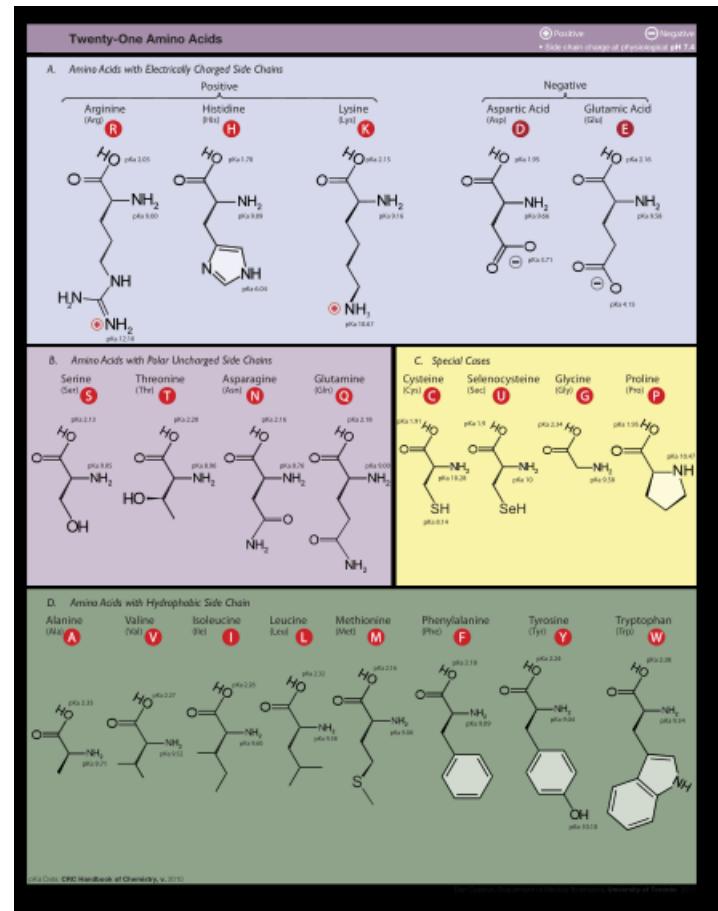
BIC score comparison

Strategy 1 BIC = 239403 ✓  
Strategy 2 BIC = 242200

Subset	Best Model	# sites	subset id	Partition names
1	GTR+G	485	88a99148e17429cfbf7c84b3c6ad405b	Partition_3
2	GTR+I+G	4940	6147b066b9346034e59c2738ca8adab6	Partition_1
3	GTR+G	407	207a8ca99fe345782ac3e9b948ae6783	Partition_4
4	GTR+G	342	a9de06c1f843a93e4545a0c53c47027b	Partition_5
5	GTR+G	733	0c045f1386af47731957496eda7faeed	Partition_2
6	GTR+G	368	067651d33be225f5683a0ee6ce036017	Partition_6
7	GTR+G	430	6ecc87f27a392fcac101b8b0c99c64e0	Partition_7

# Protein models

- 21 amino acids – 21x21 matrix too big for estimation
- Models are based largely on empirical aa replacement with matrices from different taxa
- Examples: JTT, WAG, MtREV, Blosum62



# Recommended reading

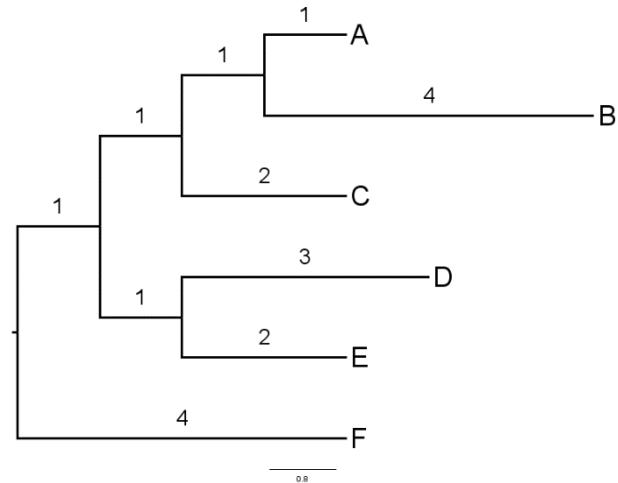
- Hoff et al. 2016. [Does the choice of nucleotide substitution models matter topologically?](#) BMC Bioinformatics 17: 143. doi.org/10.1186/s12859-016-0985-x
- Kainer & Lanfear. 2015. [The Effects of Partitioning on Phylogenetic Inference.](#) Molecular Biology and Evolution, 32(6), 1611–1627. doi.org/10.1093/molbev/msv026

# Distance methods

# Distance methods involve two stages

- Stage 1: calculate the evolutionary distance between pairs of sequences
- Stage 2: use the distances to construct a tree that describes those evolutionary distances

	A	B	C	D	E
B	5				
C	4	7			
D	7	10	7		
E	6	9	6	5	
F	8	11	8	9	8



# 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 bp of which 15 differ, or  $D = 15/98 = 0.153$

```

Antirrhina109 ATGAACTGTTTATCCCCCCCCTTTCTTCAGAATTAATGCTCATAGAGGGTTCCCTAGITGACTAGCAATTTTTTCTTTAACATTTAGCAGGTATTTCTTCAT
Araschnia39 ATGAACTGAGTTTATCCCCCCCACCTTTCTATCAATAATCGCACATAGAGGGTTCTAGTAGATTTAGCAATTTTCTCTCTTCAATTTAGCAGGAATTTCAT

```

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

	Antirrhinum109	Araschnia39	Archaeoprepus	Asterocaea82	Caligo70	Calinaga64	Castilia76	Catacropte88	Catonephele6	Cercyonis8	Chersonesia1	Chlosyne62	Clossiana76	Colobura68	
Antirrhinum109	ATGAACTGTTTATCCCCCCTTTTCATAGAGGGTTCTCAGTTGACCTTAGAAATTTCCTTACATTAGAGGGTTCATCTGTAGATTAGCAATTTCCTCTCTCATTTAGCAGGAATTTCCTCAAT	ATGAAACAGTTTATCCCCCCTTTCACTAAATATCAGACATAGAGGGTTCTCAGTTGACCTTAGAAATTTCCTCTCTCATTTAGCAGGAATTTCCTCAAT	ATGAAACAGTGTACCCCCCCTTTCACTAAATATGCCCCATGGAGGTTCATCTGAGTTAGCTATTTCCTGTTGACCTTAGAAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	ATGAAACAGTATAACCCCTCTCTCTAAATATGCTCACGGAGGATCTTCTGTTGACCTTAGAAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	ATGAAACAGTGTAACCCCCCCTTTCACTAAATATGCTCATGGAGGCTCATCTGTTGAGTTAGCTATTTCCTGTTGACCTTAGAAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	ATGAAACAGTTAACCCCCCCTTTCACTAAATATGCCCCATAGTGGATCTTCAGTTGACCTTAGAAATTTCCTCTCATCTAGCTGGAAATTTCCTCAAT	ATGAAACGGTCTAACCCCCCCTTTCACTAAATATGCCCCATAGTGGAGCAGCTGTGGACCTTAGAAATTTCCTCTCATTTAGCTGGAAATTTCCTCAAT	ATGAAACAGTGTACCCCCCCTTTCACTAAATATGCTCACGGCGGGATCTTCCGTTGACCTTAGAAATTTCCTCTCATTTAGCTGGAAATTTCCTCAAT	ATGAAACAGTGTACCCCCCCTTTCACTAAATATGCTCACGGTGGATCTTCCTGTTGACCTTAGAAATTTCCTCTCATTTAGCTGGAAATTTCCTCAAT	ATGAACTGTTTATCCCCCCTTTCACTAAATATGCCCCATGGGGGGAGCTTCCTGTTGAGTTAGCTATTTCCTTACACTTAGCTGGAAATTTCCTCAAT	ATGAACTGTTTATCCCCCCTTTCACTAAATATGCTCATAGAGGATCTTCAGTTGATTTAGCAATTTCCTCTGCTAGCTGGAAATTTCCTCAAT	ATGAAACAGTGTACCCCCCCTTTCACTAAATATGCTCATAGAGGATCTCTGTTGAGTTAGCTATTTCCTTACACTTAGCTGGAAATTTCCTCAAT	ATGAAACAGTGTACCCCCCCTTTCACTAAATATGCTCATAGAGGAGCTTCAGTTGATTTAGCTATTTCCTTACACTTAGCTGGAAATTTCCTCAAT	ATGAAACAGTTAACCCCCCCTTTCACTAAATATGCTCATGGAGGATCTTCAGTTGATTTAGCTATTTCCTTACACTTAGCTGGAAATTTCCTCAAT	ATGAAACAGTTAACCCCCCCTTTCACTAAATATGCTCATGGAGGATCTTCAGTTGAGTTAGCTATTTCCTTACACTTAGCTGGAAATTTCCTCAAT

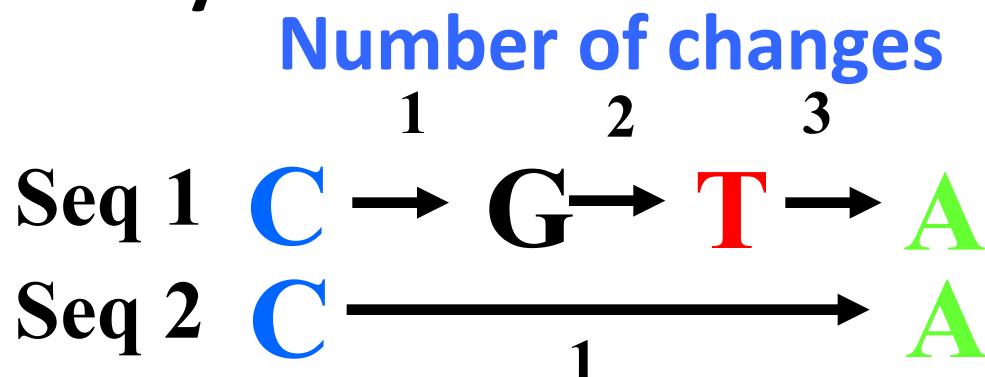
Distance matrix:

	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



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

Antirrhinum109	ATGAACTGTTTAT	TCCCCCCTTTTC	TCTAAATGCTCATAGAGGTTCCTCA	AGTTGACCTTAGAATTTTCTTTAACATTAGAGGCTTCATCTGTAGATTAGC	AAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	TAGGGTATTTCCTTCAT	TAT
Araschnia39	ATGAAACAGTTTAT	TCCCCCCTTTTC	TCTAAATGCTCATAGAGGTTCCTCA	AGTTGACCTTAGAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	TATTTTCTCTCATTTAGCAGGAATTTCCTCAAT	TAT	
Archaeorepo	ATGAAACAGTGTACCCCCC	ACTTTCACTCTAAATATTGCCATGGAGGTTCCTCA	AGTTGACCTTAGAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	TATTTTCTCTCATTTAGCAGGAATTTCCTCAAT	TATTTTCTCTCATTTAGCAGGAATTTCCTCAAT	TAT	
Asterocaea82	1 ATGAAACAGTATAACCC	ACCACTTTCTTCTCTAAATATTGCCATGGAGGTTCCTCA	AGTTGACCTTAGAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	TATTTTCTCTCATTTAGCAGGAATTTCCTCAAT	TATTTTCTCTCATTTAGCAGGAATTTCCTCAAT	TAT	
Caligo70	10 ATGAAACAGTGTACCCCCC	ACTATCACTCTAAATATTGCCATGGAGGTTCCTCA	AGTTGACCTTAGAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	TATTTTCTCTCATTTAGCAGGAATTTCCTCAAT	TATTTTCTCTCATTTAGCAGGAATTTCCTCAAT	TAT	
Calinaga64	3 ATGAAACAGTTTACCCCCC	ACTTTCACTCTAAATATTGCCATAGTGGATCTTCAGTTGACCTTAGA	AAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	AAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	AAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	AAT	
Castilia76	2 ATGAAACGGTCTACCCCCC	CCTTTCACTCTAAATATTGCCATAGTGGATCTTCAGTTGACCTTAGA	AAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	AAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	AAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	AAT	
Catacropte88	ATGAAACAGTGTACCCCCC	ACTATCACTCTAAATATTGCCATAGTGGATCTTCAGTTGACCTTAGA	AAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	AAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	AAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	TAT	
Catonephele6	ATGAAACAGTTTACCCCCC	CTTCTTCTCTAAATATTGCCATAGTGGATCTTCAGTTGACCTTAGA	AAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	AAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	AAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	AAT	
Cercyonis8	1 ATGAAACAGTTTATCCCCC	CTTATCTCTAAATATTGCCATGGAGGTTCCTCTGTAGATTAGC	TATTTTCTCTCATTTAGCAGGAATTTCCTCAAT	TATTTTCTCTCATTTAGCAGGAATTTCCTCAAT	TATTTTCTCTCATTTAGCAGGAATTTCCTCAAT	AAT	
Chersonesia1	GTGAACTGTTATCCCCC	CTTATCTCTAAATATTGCCATAGAGGATCTTCAGTTGATTTAGC	AAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	AAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	AAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	AAT	
Chlosyne62	1 ATGAAACAGTGTACCCCCC	ACTTTCTTCTCTAAATATTGCCATAGAGGATCTTCAGTTGATTTAGC	AAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	AAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	AAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	AAT	
Clossiana76	ATGAAACAGTCTACCCC	ACCCCTTTCTCTAAATATTGCCATAGAGGATCTTCAGTTGATTTAGC	AAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	AAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	AAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	CAT	
Colobura68	1 ATGAAACAGTTTATCTCTT	CTATCTCTAAATATTGCCATAGGGAGGTCTTCAGTTGATTTAGC	AAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	AAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	AAATTTCCTCTCATTTAGCAGGAATTTCCTCAAT	AAT	

Dissimilarities matrix:

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

*Correction for  
multiple  
substitutions*

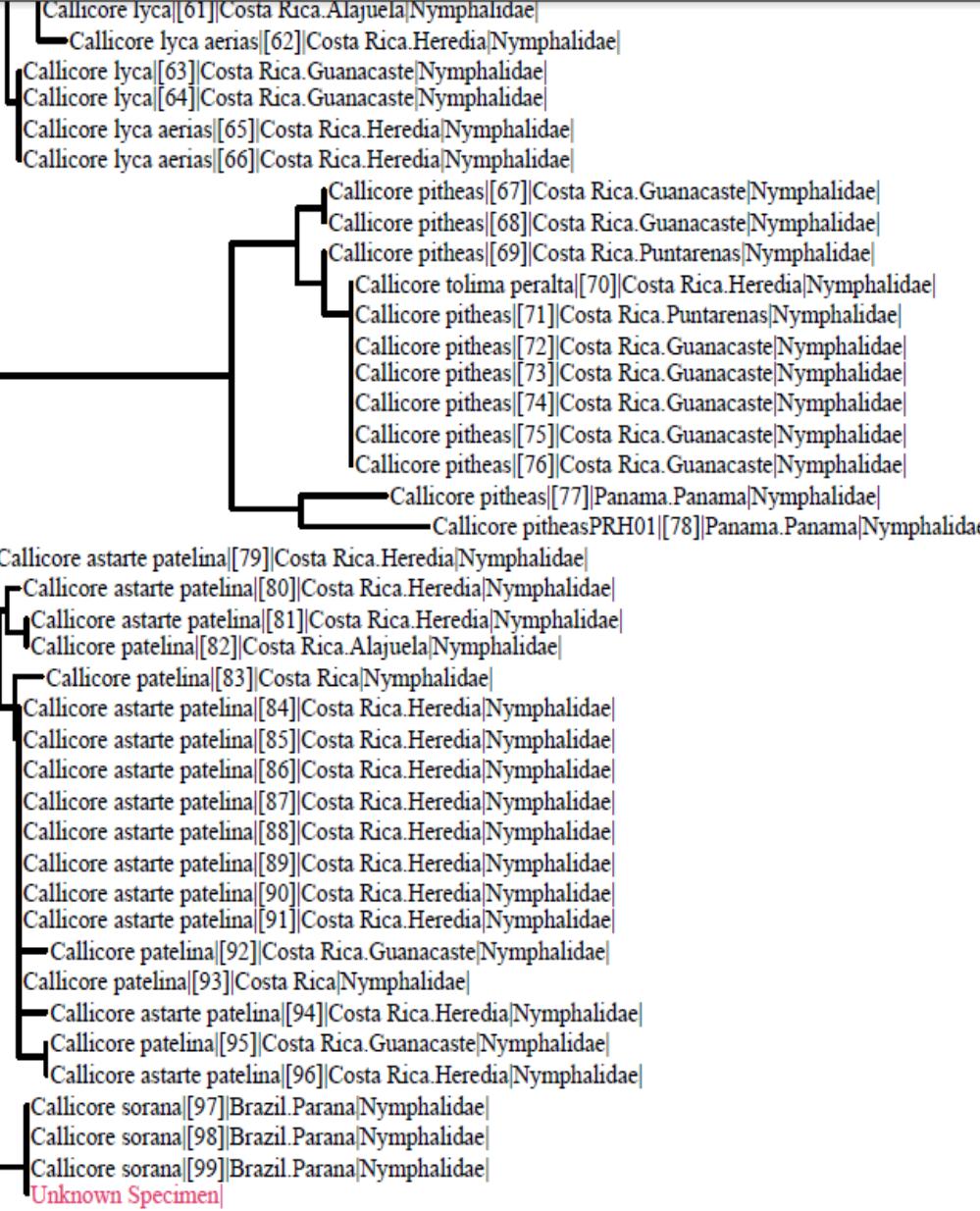


	1	2	..	n
1				
2	0.35			
:				
n	0.24	0.66		

# **Distances - advantages**

- Fast when using clustering algorithms
- A large number of models are available with many parameters - improves estimation of distances

DNA barcoding



# **Distances - disadvantages**

- **Distance estimates are only correct if model used is correct**
  - But that's also the case for ML and BI
  - 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**

- **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))