

Model Selection

GENE 606/ENTO 606/ WFSC 646

Updated 02/20/12

Suggested readings (Older)

- Ripplinger, J., and J. Sullivan. 2008. Does choice in model selection affect maximum likelihood analysis? *Systematic Biology* 57:76-85.
- Abdo, Z. 2005. Accounting for uncertainty in the tree topology has little effect on the decision-theoretic approach to model selection in phylogeny estimation. *Molecular Biology and Evolution* 22: 691-703
- Sullivan, J., Z. Abdo, P. Joyce and D.L. Swofford (2005). Evaluating the performance of a successive approximations approach to parameter optimization in maximum-likelihood phylogeny estimation. *Molecular Biology and Evolution* 22(6): 1386-1392.
- Sullivan, J. and P. Joyce (2005). Model Selection in Phylogenetics. *Annual Review of Ecology and Systematics* 36: 445-466.
- Minin V. 2003. Performance-based selection of likelihood models for phylogeny estimation *Systematic Biology* 52: 674-683
- Posada, D. and T.R. Buckley (2004) Model selection and Model Averaging in Phylogenetics: Advantages of Akaike Information Criterion and Bayesian Approaches Over Likelihood Ratio Tests. *Systematic Biology* 53(5): 793-808.
- Posada D (2008) jModelTest: Phylogenetic model averaging. *Molecular Biology and Evolution* **25**, 1253-1256.

Suggested readings (Newer)

- Kedzierska AM, Drton M, Guigó R, Casanellas M (2011) SPIn: Model Selection for Phylogenetic Mixtures via Linear Invariants. Mol Biol Evol. doi:10.1093/molbev/msr259
- Thi Nguyen MA, Gesell T, von Haeseler A (2012) ImOSM: Intermittent Evolution and Robustness of Phylogenetic Methods. Mol Biol Evol 29:663-673
- Thi Nguyen MA, Klaere S, von Haeseler A (2011) MISFITS: Evaluating the Goodness of Fit between a Phylogenetic Model and an Alignment. Mol Biol Evol 28:143-152

Review of models

- **Reversible** (GTR family): rate matrix, among-site rate variation, base frequencies
- **Non-reversible:**
 - accommodate base frequency changes in different parts of tree
 - accommodate changes in rates of sites in different parts of tree
- **Non-independence of sites**
 - Codon models (To cover in Amino-Acid Analyses Lecture)
 - rRNA models
- **Partitioned models** (To cover in Data Partitions Lecture)
 - Assume different models for each partition (e.g. different genes, codons, stem vs loop)

Why and when is the model important

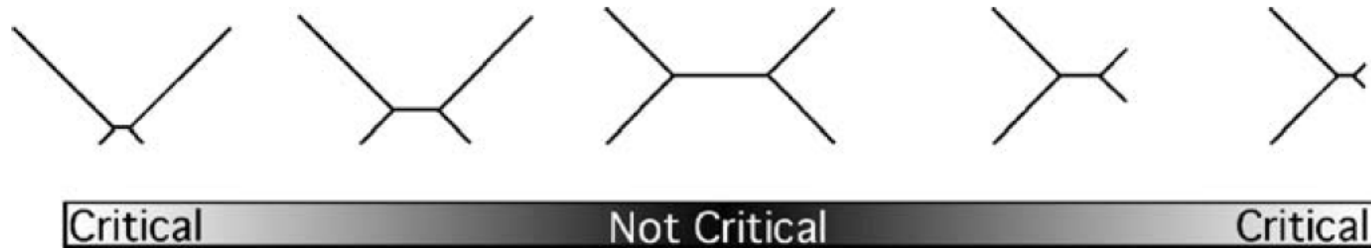
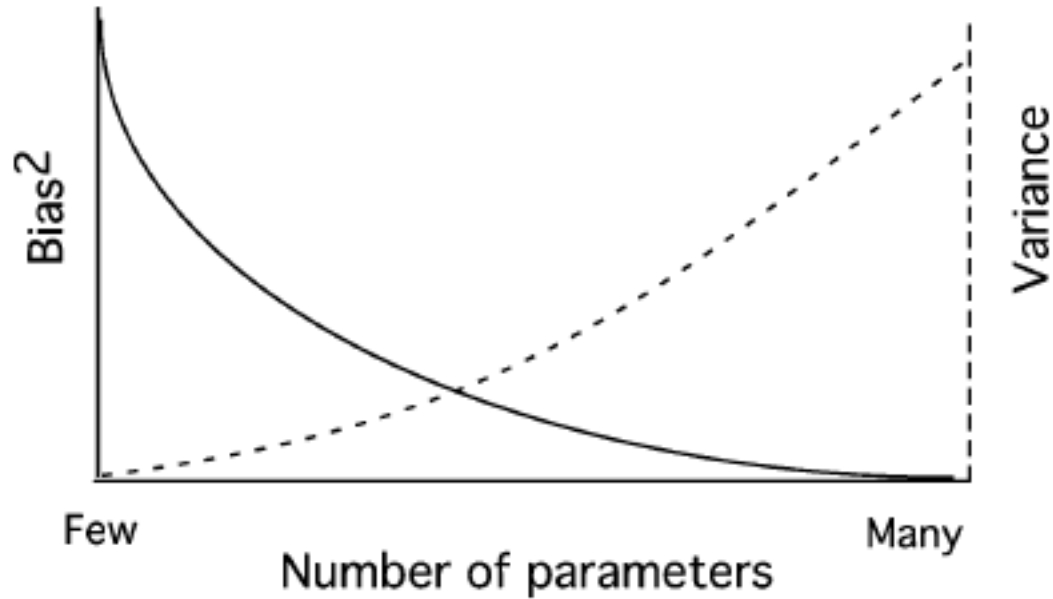


Figure 2 The effect of topology on robustness. At the center of the continuum, phylogenetics signal is strong and model choice is not critical (i.e., maximum likelihood is robust to violations of model assumptions). In the Felsenstein zone (*left*), model selection is critical, as is also the case for the inverse Felsenstein zone (*right*).

Model Selection

- All models are wrong, but some are useful (Box 1976)
- Model selection is a way of approximating, not identifying, full reality
- Statistical model selection is based on the parsimony principle; hypotheses should be kept as simple as possible
- Increasing the number of parameters will increase the fit between the model and the data (increase the likelihood), but at a cost
- Trade-off between bias and variance:
 - Bias: distance between the average estimate and truth
 - Variance: spread of the estimates around the truth

Bias vs Variance



Model Selection Strategies

- Likelihood Ratio Tests ([jModeltest](#))
- Akaike Information Criterion ([jModeltest](#))
- Bayesian Information Criterion ([jModeltest](#))
- Performance-Based model selection ([DT-ModSel](#) and [jModeltest](#))
- Bayes factors (“manually”)

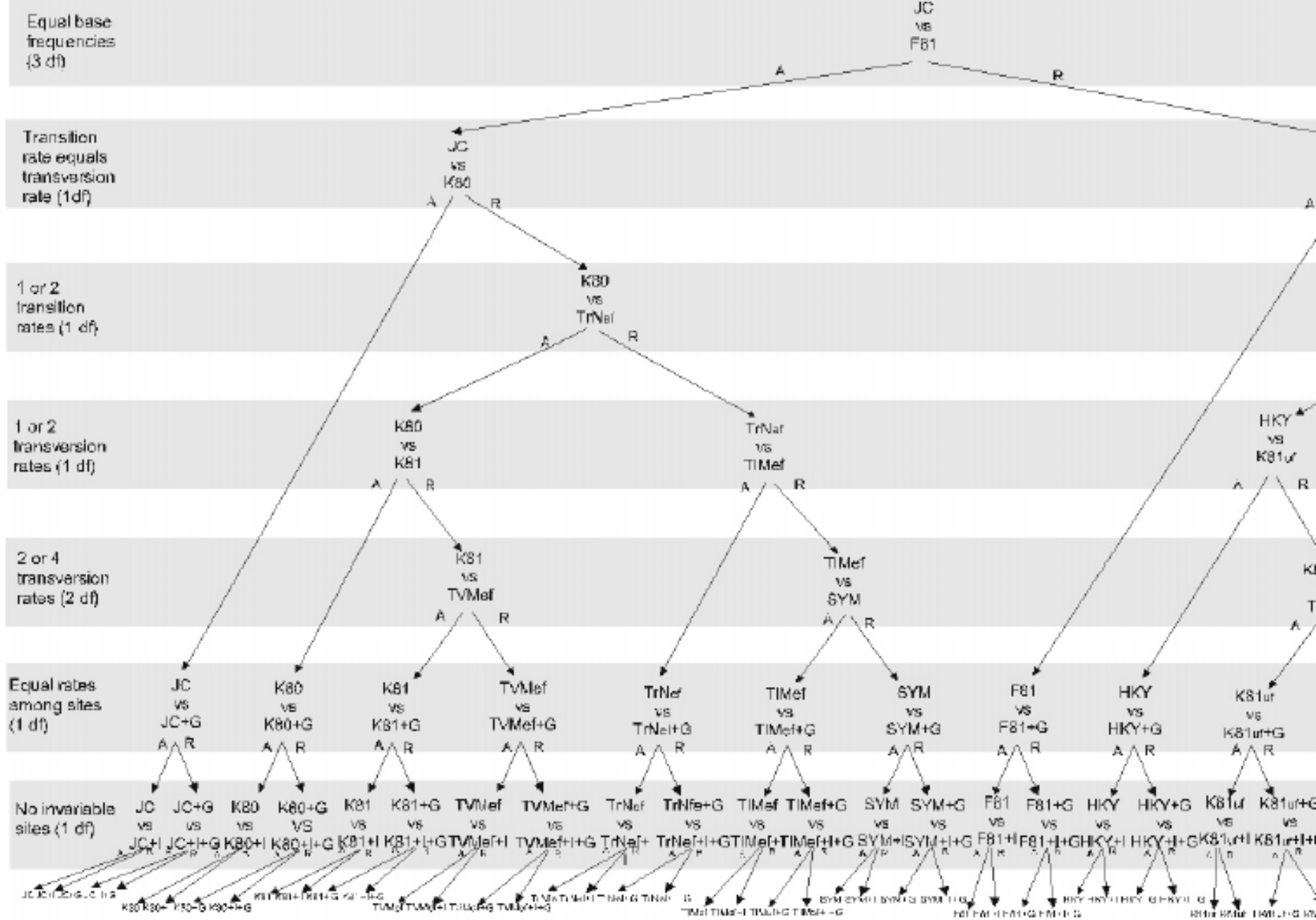
Likelihood Ratio Test (LRT)

$$\delta = 2(\ln L_1 - \ln L_0)$$

- Can **only** be used to evaluate **nested** hypotheses
- L_1 more complex model
- Test statistic evaluated under assumption of asymptotic convergence to X^2 (d.f. = diff in # parameters)

Hierarchical (h)LRT in phylogenetics

1. Infer a phylogenetic tree with another method (parsimony or distance)
2. Estimate the likelihood of that tree under different models of the GTR family (56 models in Modeltest program; 88 models in jModeltest)
3. Conduct LRT in a hierarchical fashion



Potential weaknesses of hLRT

- Dependence on initial estimate of topology
 - Use of initial trees has little effect on model chosen, but
 - Very poor trees can yield very poor model estimates
- Arbitrary order of comparison can have effect on which model is selected
- Can only compare nested hypotheses

Akaike Information Criteria (AIC)

$$AIC_i = -2 \ln L_i + 2k_i$$

- For a particular model i :
 - L_i = max log likelihood
 - k = number of parameters
- prefer the model with the smallest AIC
- provides a measure of fit between model and data and includes a penalty for overparameterization
- Small sample sizes ($n / k < 40$), use AIC_c

Bayesian Information Criteria (BIC)

$$\text{BIC}_i = -2 \ln L_i + k_i \ln n$$

- For a particular model i :
 - L_i = max log likelihood
 - k = number of parameters
 - n = sample size (# of characters ???)
- provides a measure of fit between model and data and penalizes for overparameterization (more heavily than AIC, especially with large n)

Performance Based-DT

- ranks models on the basis of the weighted expected error in branch-length estimates
- weights are derived from the BIC
- focuses on the fact that both the tree topology and the branch lengths (the rate of evolution \times the time between each node or speciation event in the tree) are critical.
- Minin et al. 2003: incorporates some measure of phylogenetic performance. “Asses models through a penalty or loss function, related to how dissimilar the branch length estimates are across models, and pick the model with the minimum posterior loss”.

Comparison of model selection approaches

- hLRT, AIC, and BIC all use an initial topology
- hLRT can only compare nested hypotheses, while AIC and BIC can compare multiple nested and non-nested hypotheses
- AIC and BIC outcome does not depend on the order of comparisons, while hLRT does
- AIC and BIC allow assessment model selection uncertainty, and estimation of phylogenies and model parameters using all available models (model-averaged inference or multi-model inference).

Comparison of model selection approaches

Table 2
Model Selection Strategies Implemented in jModelTest

	Hierarchical Likelihood Ratio Tests	Dynamical Likelihood Ratio Tests	Akaike Information Criterion	Bayesian Information Criterion	Performance- Based Selection
Abbreviation	hLRTs	dLRTs	AIC	BIC	DT
Base tree	Fixed	Fixed	Fixed, optimized	Fixed, optimized	Fixed, optimized
Nesting requirement	Yes	Yes	No	No	No
Simultaneous comparison	No	No	Yes	Yes	Yes
Selection uncertainty	No	No	Yes	Yes	Yes ^a
Parameter importance	No	No	Yes	Yes	Yes ^a
Model averaging	No	No	Yes	Yes	Yes ^a

^a DT weights are simply the rescaled reciprocal DT scores. This is a gross implementation very likely to change.

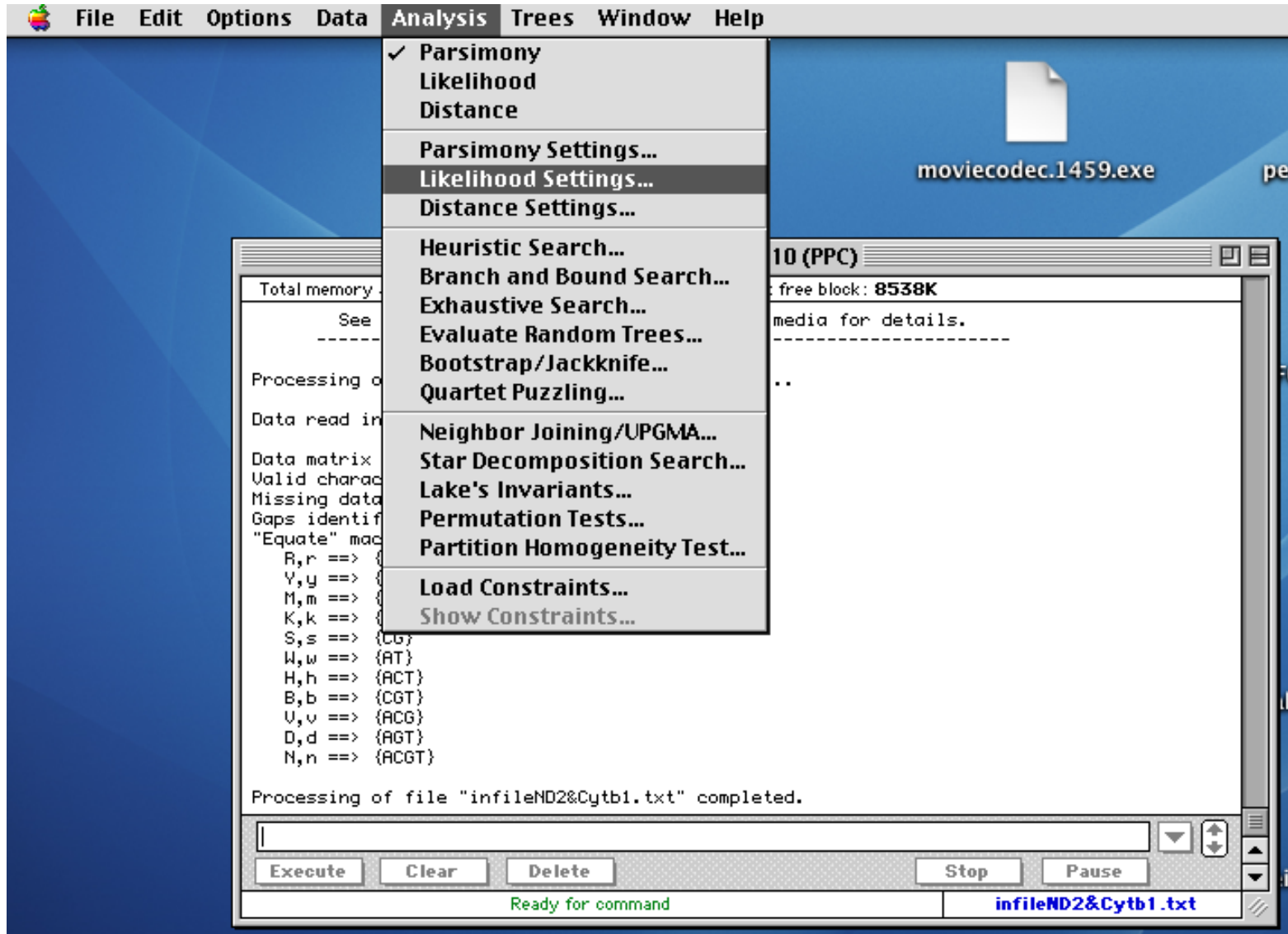
From Posada 2008

Implementation of model selection approaches

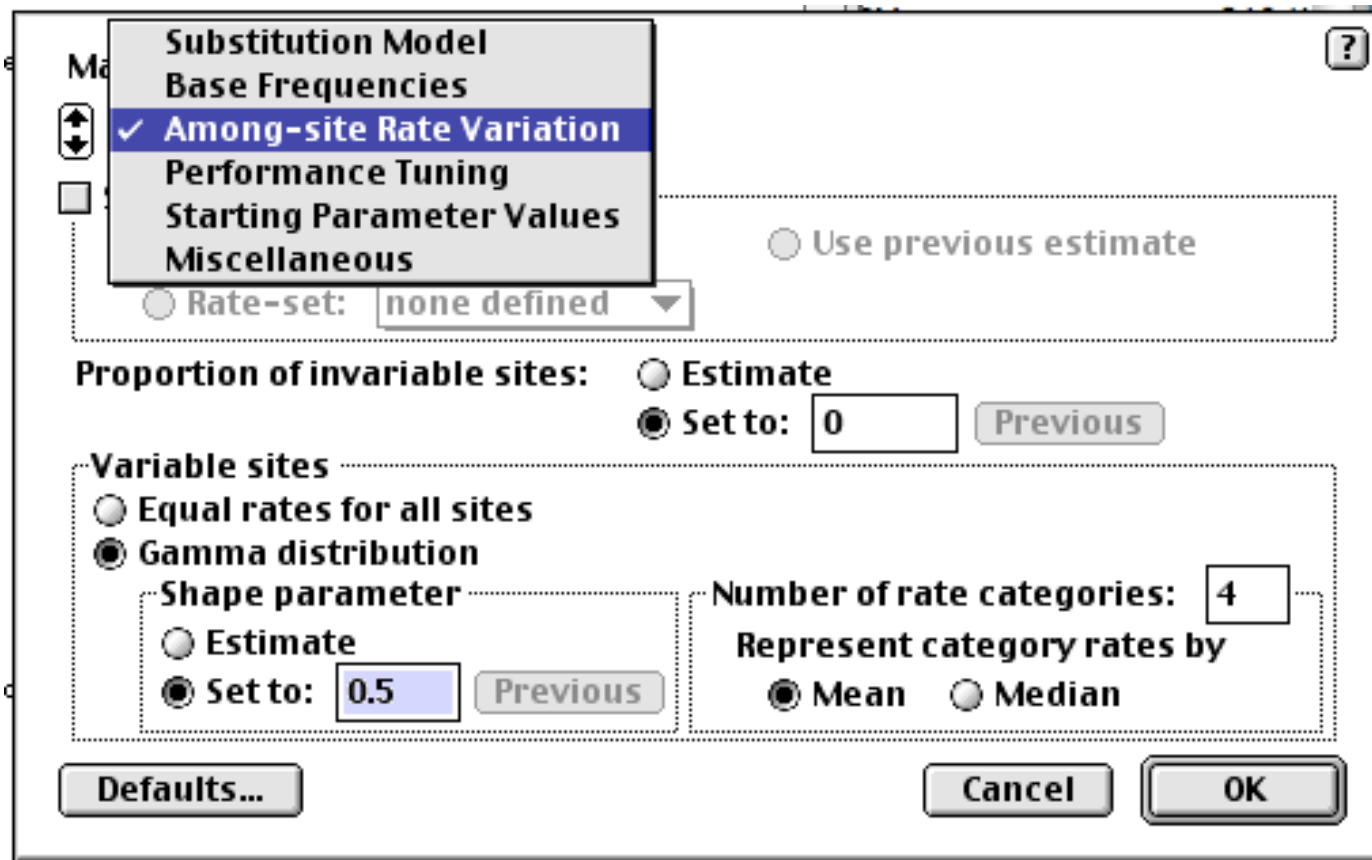
- jModeltest <http://darwin.uvigo.es/software/modeltest.html>
- A command file for PAUP is already available within the Modeltest package (not used any more)
 - Infers a NJ tree
 - Estimates likelihoods and parameter estimates for each of 56 models (output: model.scores)
 - Use model.scores file as input for Modeltest

- The following slides are images of the GUI windows of Paup* for selecting different models

Likelihood Settings



ML settings submenus



JC

Maximum likelihood options: ?

↕ Base Frequencies ▼

☐ Use empirical frequencies

☒ Assume equal frequencies

☐ Estimate

☐ Set to: A = 0.25
Previous C = 0.25
G = 0.25
T = 0.25

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

Two-parameter model variant for unequal base frequencies

☐ Hasegawa-Kishino-Yano (1985)

☐ Felsenstein (1984)

Defaults... Cancel OK

F81

Maximum likelihood options: ?

↕ Base Frequencies ▼

☒ Use empirical frequencies
☐ Assume equal frequencies
☐ Estimate
☐ Set to:

Previous

A = 0.25429
C = 0.32574
G = 0.13325
T = 0.28672

Two-parameter model variant for unequal base frequencies

☐ Hasegawa-Kishino-Yano (1985)
☐ Felsenstein (1984)

Defaults... Cancel OK

	A	C	G	T
A		$\mu\pi_C$	$\mu\pi_G$	$\mu\pi_T$
C	$\mu\pi_A$		$\mu\pi_G$	$\mu\pi_T$
G	$\mu\pi_A$	$\mu\pi_C$		$\mu\pi_T$
T	$\mu\pi_A$	$\mu\pi_C$	$\mu\pi_G$	

K2P

Maximum likelihood options:

Substitution Model

☐ All rates equal ("1 ST")

☒ Ti rate ≠ tv rate ("2 ST")

Ti/tv ratio: ☐ Estimate ☒ Set to:

	A	C	G	T
A		β	α	β
C	β		β	α
G	α	β		β
T	β	α	β	

☐ General time-reversible ("6 ST")

Rate matrix: ☐ Estimate

☐ Set to:

	C	G	T
A			
C			
G			

HKY85

Maximum likelihood options:

Base Frequencies ▼

☐ Use empirical frequencies
☐ Assume equal frequencies
☒ Estimate
☐ Set to: A =
 C =
 G =
 T =

Two-parameter model variant for unequal base frequencies

☒ Hasegawa-Kishino-Yano (1985)
☐ Felsenstein (1984)

	A	C	G	T
A		$\mu\pi_C$	$\mu\pi_{TC}$	$\mu\pi_T$
C	$\mu\pi_A$		$\mu\pi_G$	$\mu\pi_{GT}$
G	$\mu\pi_{AC}$	$\mu\pi_C$		$\mu\pi_T$
T	$\mu\pi_A$	$\mu\pi_{TC}$	$\mu\pi_G$	

F84

Maximum likelihood options:

☒ Base Frequencies

☐ Use empirical frequencies
☐ Assume equal frequencies
☒ Estimate
☐ Set to:

A =
 C =
 G =
 T =

Two-parameter model variant for unequal base frequencies
☐ Hasegawa-Kishino-Yano (1985)
☒ Felsenstein (1984)

	A	C	G	T
A		$\mu\pi_C$	$\mu\pi_G(1 + \frac{\kappa}{\pi_R})$	$\mu\pi_T$
C	$\mu\pi_A$		$\mu\pi_G$	$\mu\pi_T(1 + \frac{\kappa}{\pi_Y})$
G	$\mu\pi_A(1 + \frac{\kappa}{\pi_R})$	$\mu\pi_C$		$\mu\pi_T$
T	$\mu\pi_A$	$\mu\pi_C(1 + \frac{\kappa}{\pi_Y})$	$\mu\pi_G$	

GTR

Maximum likelihood options:

Substitution Model

☐ All rates equal ("1 ST")
☐ Ti rate ≠ tv rate ("2 ST")
 Ti/tv ratio: ☐ Estimate ☐ Set to:
 Previous

☒ General time-reversible ("6 ST")
 Rate matrix: ☐ Estimate
 Previous ☒ Set to:

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

	C	G	T
A	1	1	1
C		1	1
G			1

Defaults... Cancel OK

GTR

Maximum likelihood options:

Substitution Model

☐ All rates equal ("1 ST")
☐ Ti rate ≠ tv rate ("2 ST")
 Ti/tv ratio: ☐ Estimate ☐ Set to:

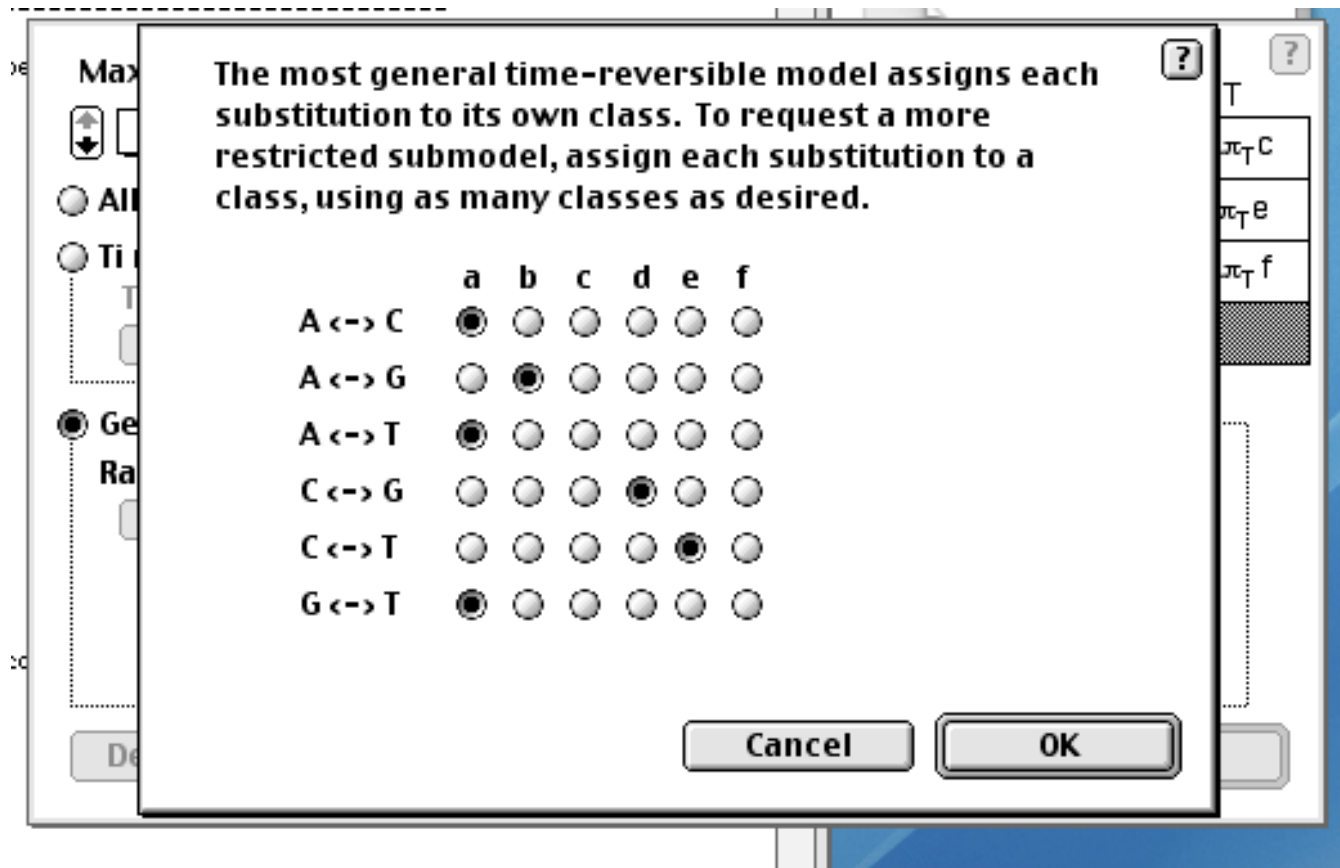
☒ General time-reversible ("6 ST")
 Rate matrix: ☒ Estimate
 ☐ Set to:

	A	C	G	T
A		$\mu\pi_{C A}$	$\mu\pi_{G A}$	$\mu\pi_{T A}$
C	$\mu\pi_{A C}$		$\mu\pi_{G C}$	$\mu\pi_{T C}$
G	$\mu\pi_{A G}$	$\mu\pi_{C G}$		$\mu\pi_{T G}$
T	$\mu\pi_{A T}$	$\mu\pi_{C T}$	$\mu\pi_{G T}$	

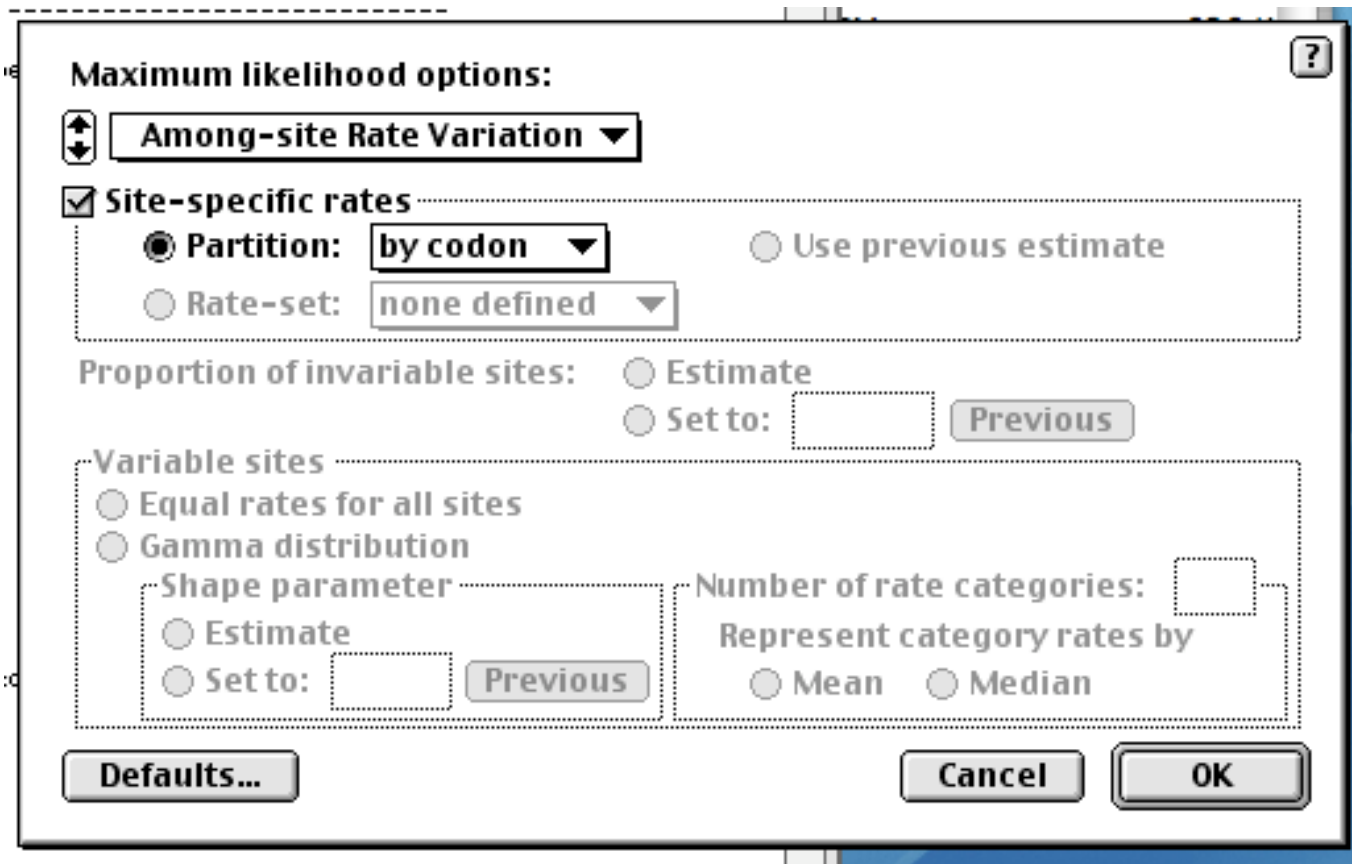
C G T
 A
 C
 G

1

GTR submodels



Among Site Rate Variation (by codon)



A screenshot of a software dialog box titled "Maximum likelihood options:". The dialog box has a question mark icon in the top right corner. It contains several sections for configuring site rate variation models. The "Among-site Rate Variation" section is expanded, showing options for "Site-specific rates" (checked), "Partition" (set to "by codon"), and "Rate-set" (set to "none defined"). There are also options for "Proportion of invariable sites" (Estimate or Set to) and "Variable sites" (Equal rates for all sites or Gamma distribution). The Gamma distribution section includes a "Shape parameter" (Estimate or Set to) and a "Number of rate categories" (set to 4). The "Represent category rates by" section has options for "Mean" and "Median". At the bottom, there are buttons for "Defaults...", "Cancel", and "OK".

Maximum likelihood options: ?

Among-site Rate Variation ▼

☒ Site-specific rates

☒ Partition: **by codon** ▼ ☐ Use previous estimate

☐ Rate-set: **none defined** ▼

Proportion of invariable sites: ☐ Estimate ☐ Set to: Previous

Variable sites

☐ Equal rates for all sites

☐ Gamma distribution

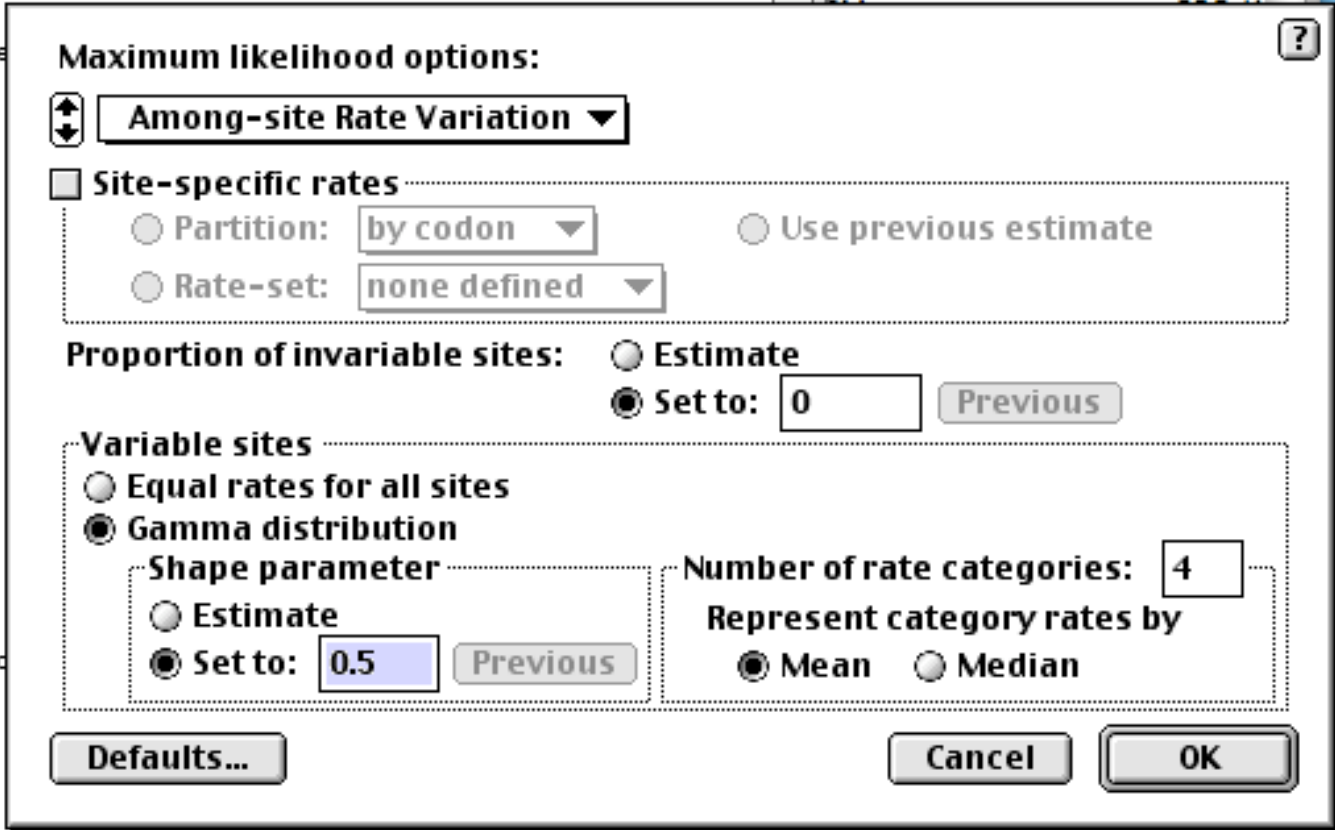
Shape parameter: ☐ Estimate ☐ Set to: Previous

Number of rate categories:

Represent category rates by: ☐ Mean ☐ Median

Defaults... Cancel OK

Among Site Rate Variation (Gamma and PINVAR)



A screenshot of a software dialog box titled "Maximum likelihood options:". The dialog box contains several sections for configuring evolutionary models. At the top, there is a dropdown menu for "Among-site Rate Variation". Below it is a checkbox for "Site-specific rates" with associated options for partitioning and rate sets. The "Proportion of invariable sites" section includes radio buttons for "Estimate" and "Set to" with a text input field. The "Variable sites" section has radio buttons for "Equal rates for all sites" and "Gamma distribution", with the latter having a "Shape parameter" sub-section with "Estimate" and "Set to" options. To the right of the "Gamma distribution" section is a "Number of rate categories" input field and a "Represent category rates by" section with "Mean" and "Median" radio buttons. At the bottom are buttons for "Defaults...", "Cancel", and "OK".

Maximum likelihood options: ?

Among-site Rate Variation ▼

☐ Site-specific rates

☐ Partition: by codon ▼ ☐ Use previous estimate

☐ Rate-set: none defined ▼

Proportion of invariable sites: ☐ Estimate ☒ Set to: 0 Previous

Variable sites

☐ Equal rates for all sites

☒ Gamma distribution

Shape parameter

☐ Estimate ☒ Set to: 0.5 Previous

Number of rate categories: 4

Represent category rates by

☒ Mean ☐ Median

Defaults... Cancel OK