

Manchester
Bioinformatics



RNA Structure, Evolution and Phylogenetics

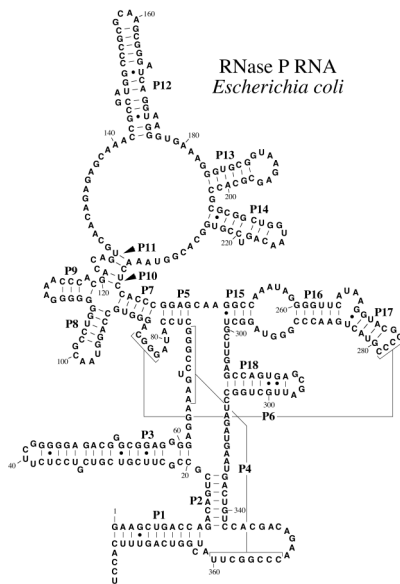
School of Biological Sciences

Paul Higgs
David Hoyle
Cendrine Hudelot
Daniel Jameson
Stephen Morgan
Nick Savill

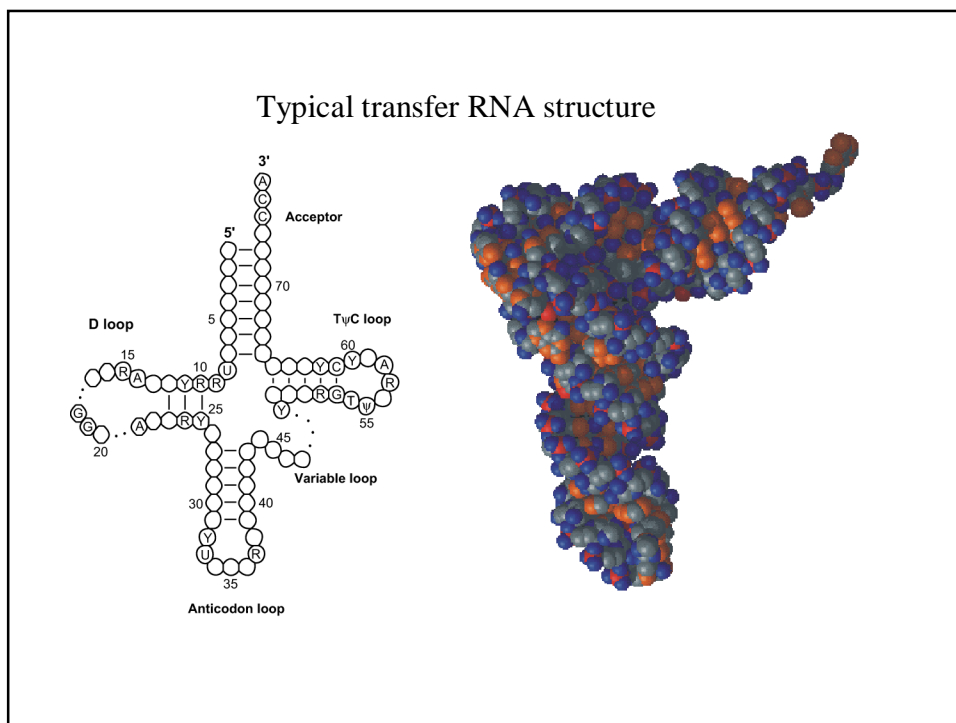
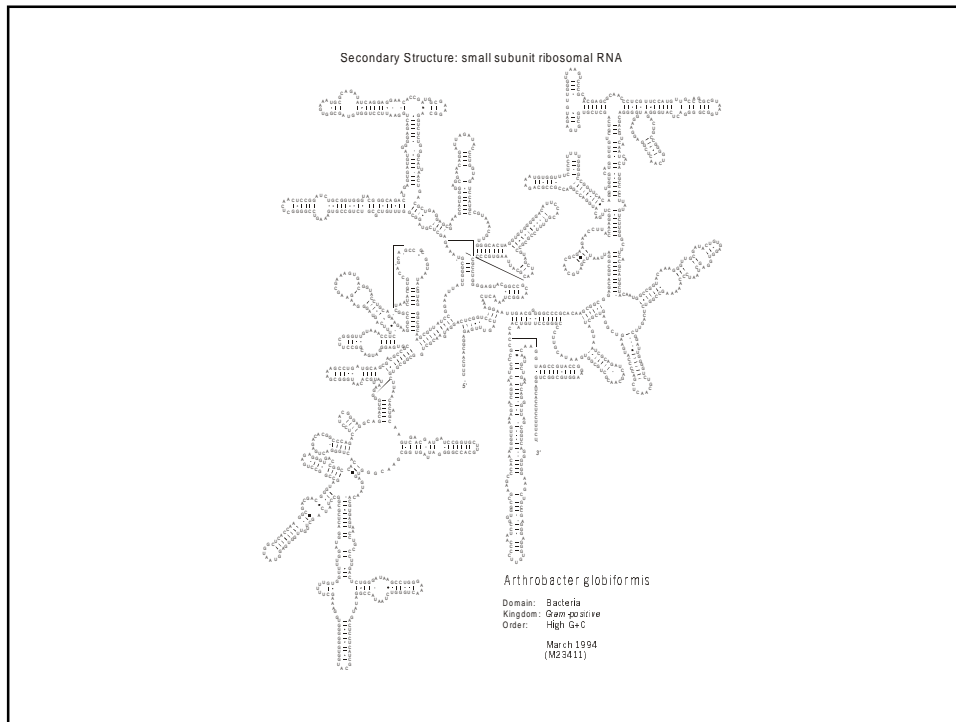
Dept of Computer Science

Magnus Rattray
Howsun Jow

Supported by BBSRC

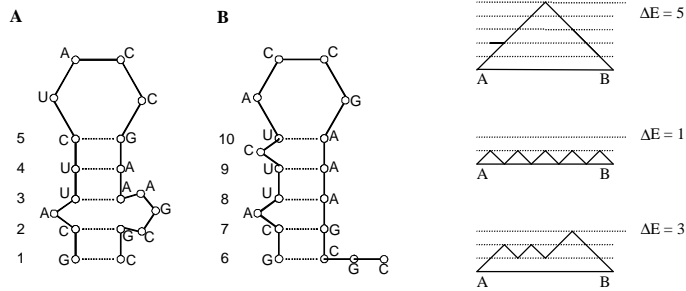


RNA Structure, Evolution and Phylogenetics

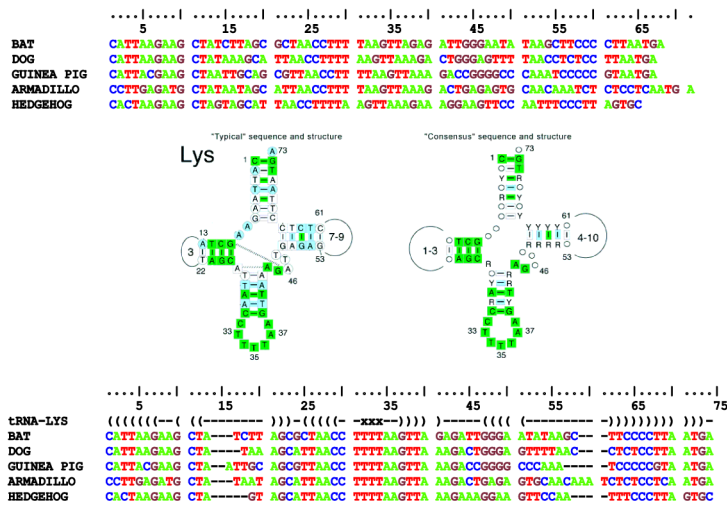


Trailer No. 2 RNA as a Disordered System

Measure Barriers Between Groundstates



The relevant barrier height controlling kinetics at low temperature is the highest point on the lowest route



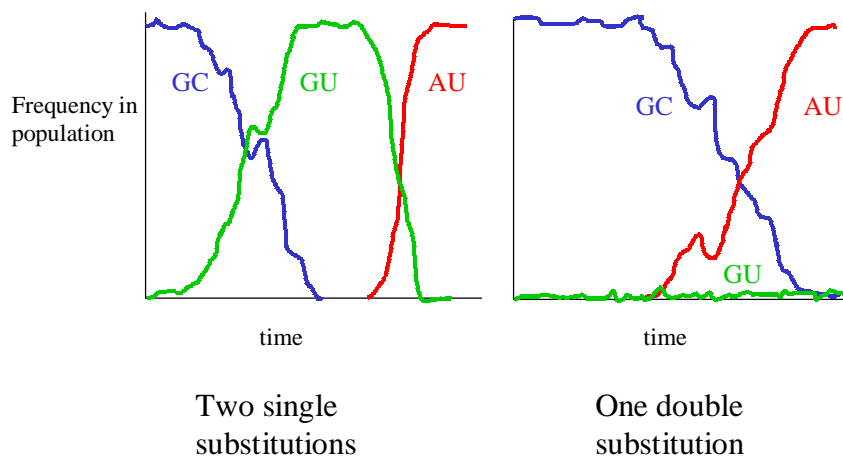
Compensatory Substitutions

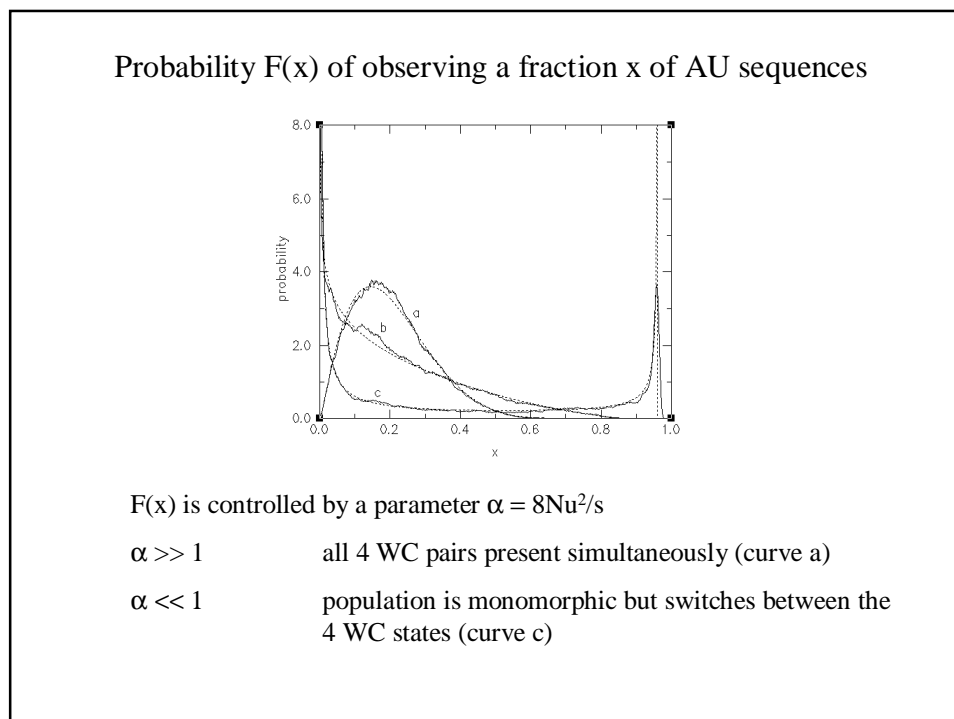
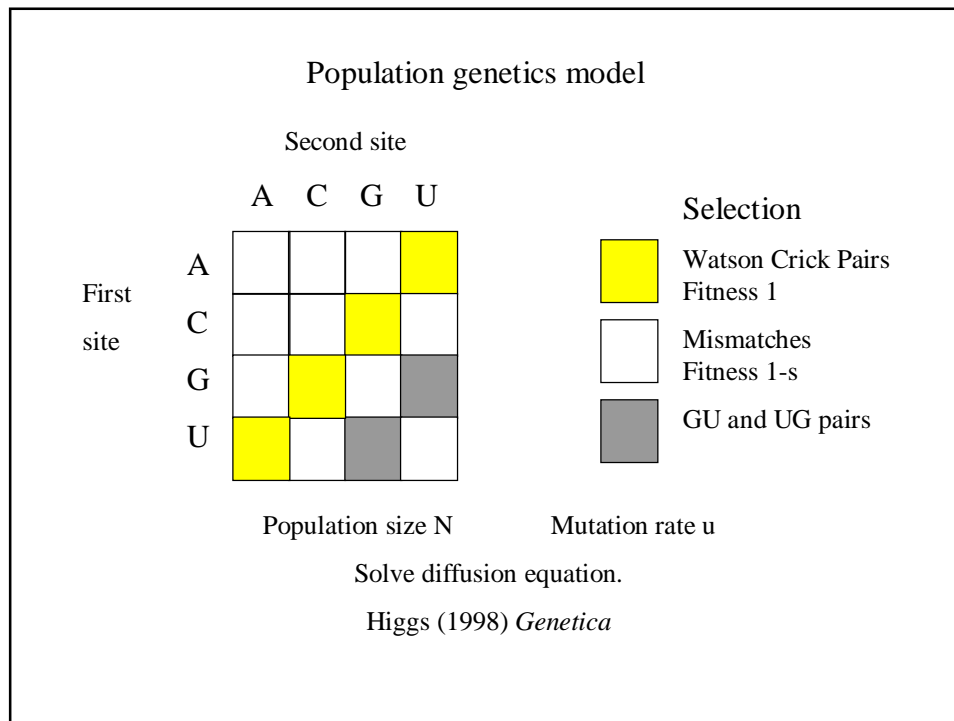
Two sides of the acceptor stem from a tRNA are shown.

Due to structure conservation alignment is possible in widely different species.

	1234567	7654321
	((((())))))
Bacillus subtilis	G G CUC G G	C CGAG C C
Escherichia coli	G C CC G GA	U C CG G GC
Saccharomyces cerevisiae	G C G GAU U	A AU U CG C
Drosophila melanogaster	G C CGAAA	UUUC G GC
Homo sapiens	G C CGAAA	UUUC G GC

Two possible mechanisms of compensatory substitutions





Analysis of RNA sequence databases

	tRNA mitoch.	tRNA general	tRNA archaea	Rnase P	SSU rRNA	
G+C average	0.339	0.532	0.636	0.594	0.545	
G+C helical regions	0.448	0.681	0.829	0.730	0.674	
Frequencies	GC	0.266	0.372	0.473	0.385	0.352
	CG	0.121	0.260	0.320	0.296	0.298
	AU	0.257	0.128	0.057	0.117	0.122
	UA	0.233	0.142	0.077	0.104	0.173
	GU	0.046	0.043	0.031	0.050	0.020
	UG	0.030	0.025	0.020	0.022	0.021
	MM	0.046	0.030	0.022	0.026	0.014
Number of sequences	884	754	64	84	455	
Number of pairs	21	21	21	80	296	

Selection for thermodynamically stable structures

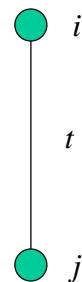
Higgs (2000) *Quart. Rev. Biophysics*

A Substitution Rate Model for RNA Stems
for use with Molecular Phylogenetics

$P_{ij}(t)$ = probability of being in state j at time t
given that ancestor was in state i at time 0.

$$\frac{dP_j}{dt} = \sum_k P_k r_{kj}$$

States label pairs (AU, GC etc.) not single bases.



Model 7A is a General Reversible 7-state Model

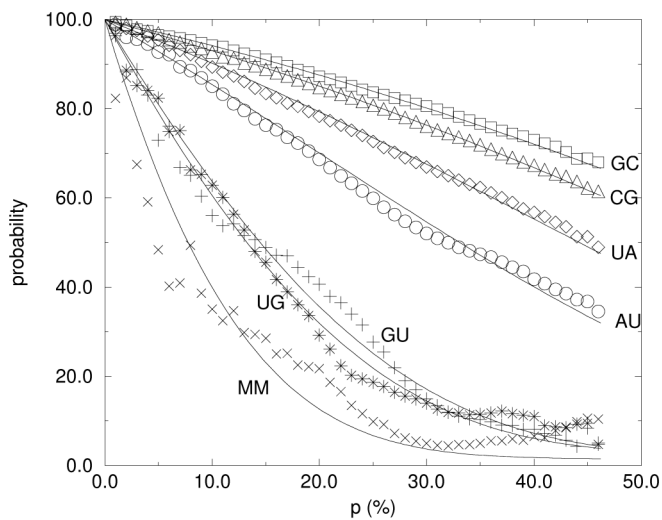
7 frequencies π_i + 21 rate parameters α_{ij}

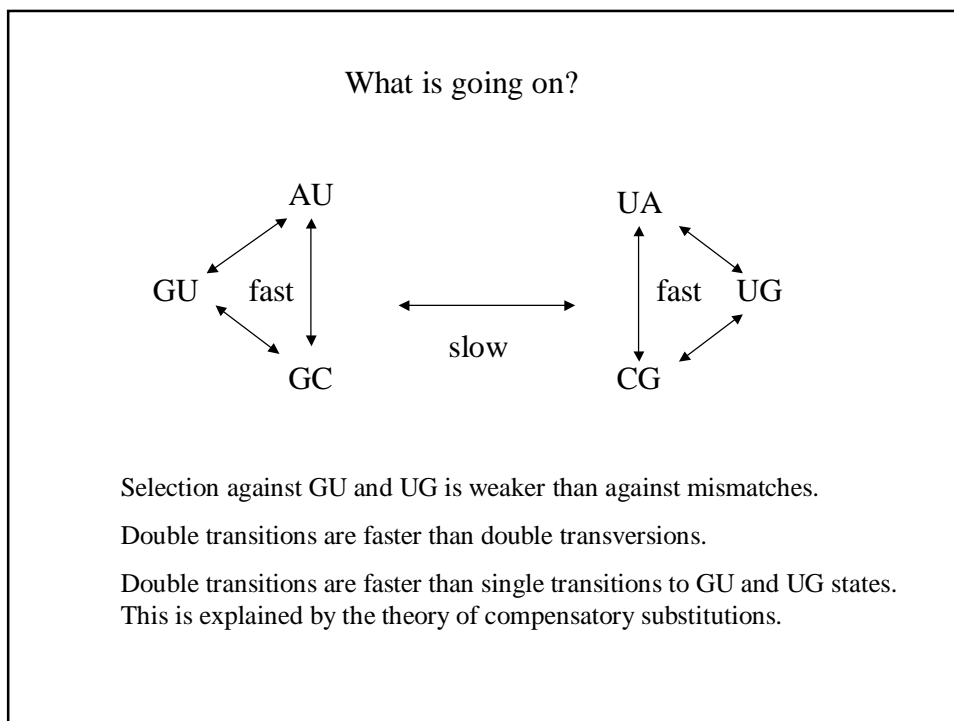
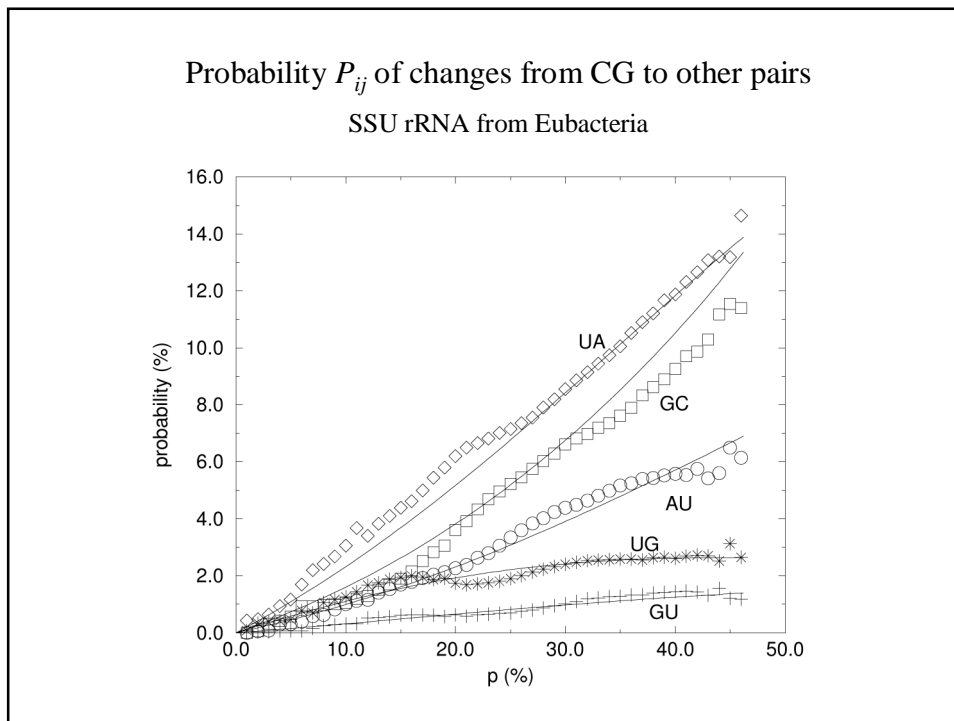
- 2 constraints = 26 free parameters

		1	2	3	4	5	6	7
		AU	GU	GC	UA	UG	CG	MM
1	AU	*	$\pi_2\alpha_{12}$	$\pi_3\alpha_{13}$	$\pi_4\alpha_{14}$	$\pi_5\alpha_{15}$	$\pi_6\alpha_{16}$	$\pi_7\alpha_{17}$
2	GU	$\pi_1\alpha_{12}$	*	$\pi_3\alpha_{23}$	$\pi_4\alpha_{24}$	$\pi_5\alpha_{25}$	$\pi_6\alpha_{26}$	$\pi_7\alpha_{27}$
3	GC	$\pi_1\alpha_{13}$	$\pi_2\alpha_{23}$	*	$\pi_4\alpha_{34}$	$\pi_5\alpha_{35}$	$\pi_6\alpha_{36}$	$\pi_7\alpha_{37}$
4	UA	$\pi_1\alpha_{14}$	$\pi_2\alpha_{24}$	$\pi_3\alpha_{34}$	*	$\pi_5\alpha_{45}$	$\pi_6\alpha_{46}$	$\pi_7\alpha_{47}$
5	UG	$\pi_1\alpha_{15}$	$\pi_2\alpha_{25}$	$\pi_3\alpha_{35}$	$\pi_4\alpha_{45}$	*	$\pi_6\alpha_{56}$	$\pi_7\alpha_{57}$
6	CG	$\pi_1\alpha_{16}$	$\pi_2\alpha_{26}$	$\pi_3\alpha_{36}$	$\pi_4\alpha_{46}$	$\pi_5\alpha_{56}$	*	$\pi_7\alpha_{67}$
7	MM	$\pi_1\alpha_{17}$	$\pi_2\alpha_{27}$	$\pi_3\alpha_{37}$	$\pi_4\alpha_{47}$	$\pi_5\alpha_{57}$	$\pi_6\alpha_{67}$	*

Probability of remaining in same state P_{ii}

SSU rRNA sequences from Eubacteria





Analysis of RNA Substitution Rates

		tRNA mitoch.	tRNA general	tRNA archaea	Rnase P	SSU rRNA
Mutabilities	GC	0.67	0.49	0.45	0.65	0.55
	CG	0.84	0.83	0.89	0.60	0.66
	AU	0.86	1.46	4.01	1.46	1.40
	UA	0.77	1.24	1.78	1.09	0.93
	GU	2.44	1.96	1.85	1.72	3.92
	UG	3.32	5.01	3.00	2.84	4.36
	MM	2.32	0.99	0.86	5.24	7.84
Double transitions / Double transversions		4.7	1.7	2.3	3.1	2.1
Double transitions / Transitions to GU or UG		1.6	2.0	8.9	3.6	2.8

Thermodynamic properties influence Evolutionary properties

Model Selection

Savill, Hoyle & Higgs (2001) *Genetics*

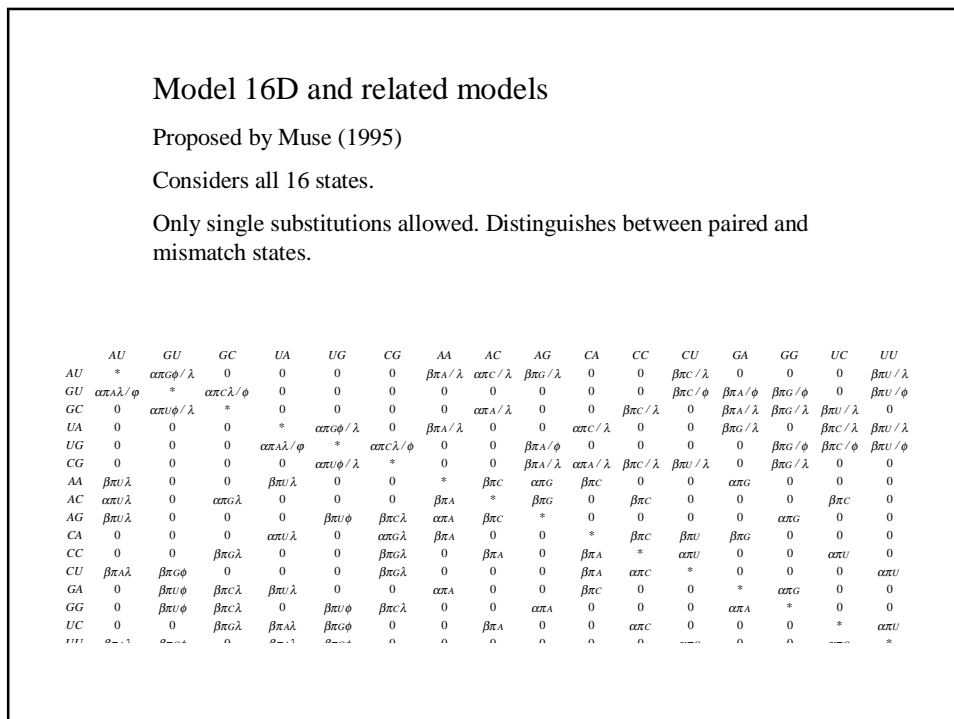
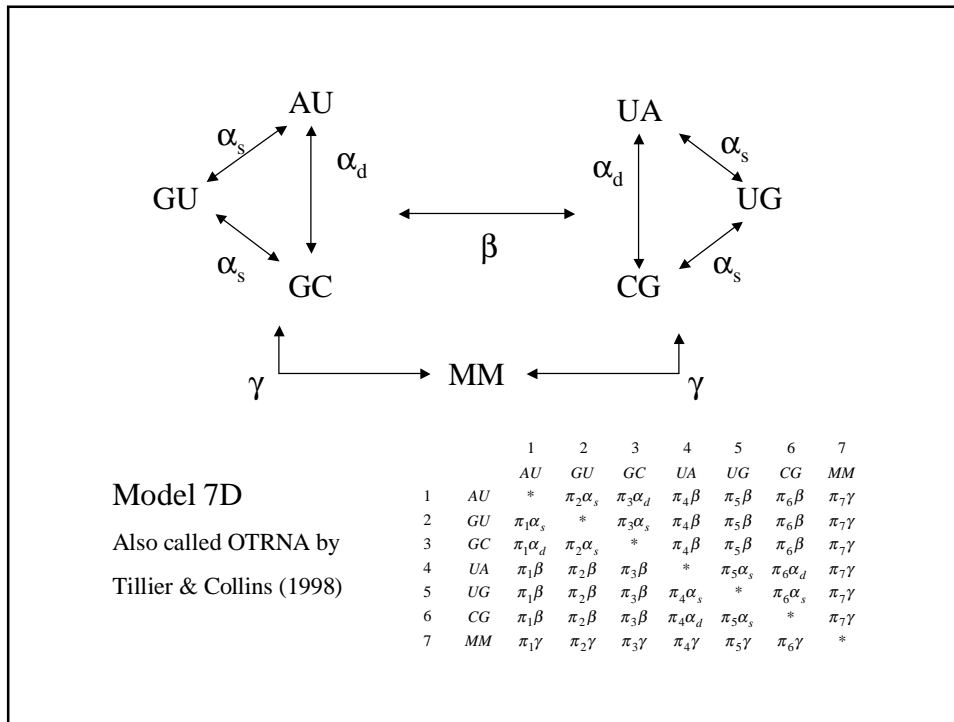
How complex must a model be to explain the data?

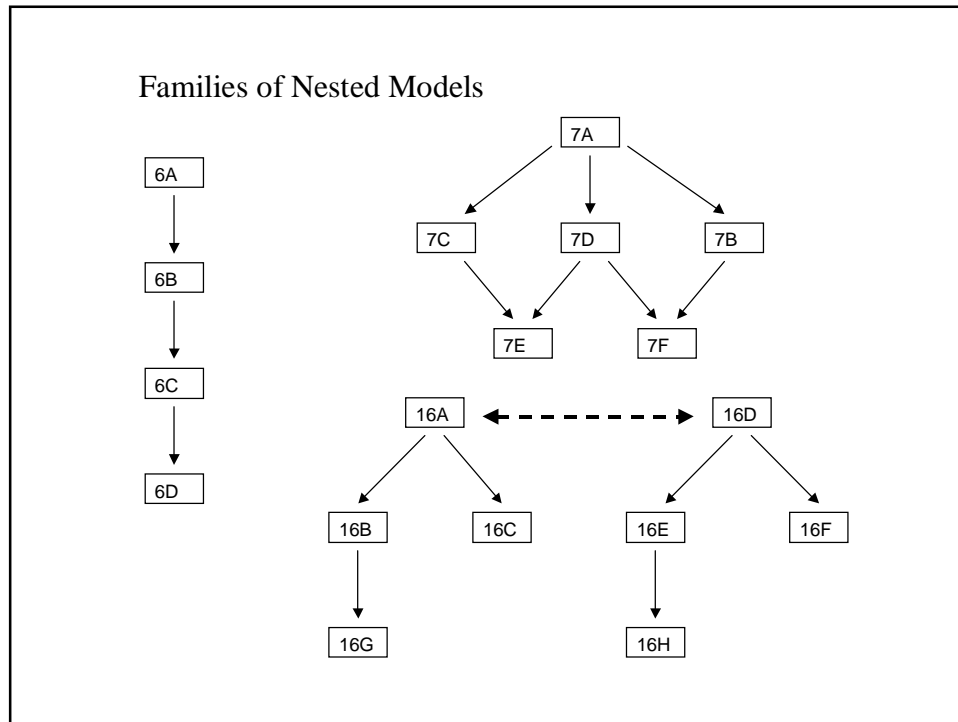
How many parameters are justified?

Many other rate models have been proposed:

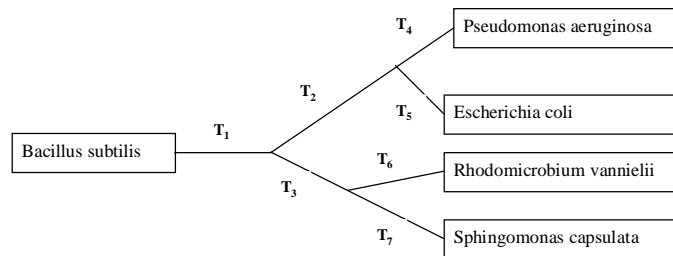
- Schöniger & von Haeseler
- Muse
- Rzhetsky
- Tillier & Collins

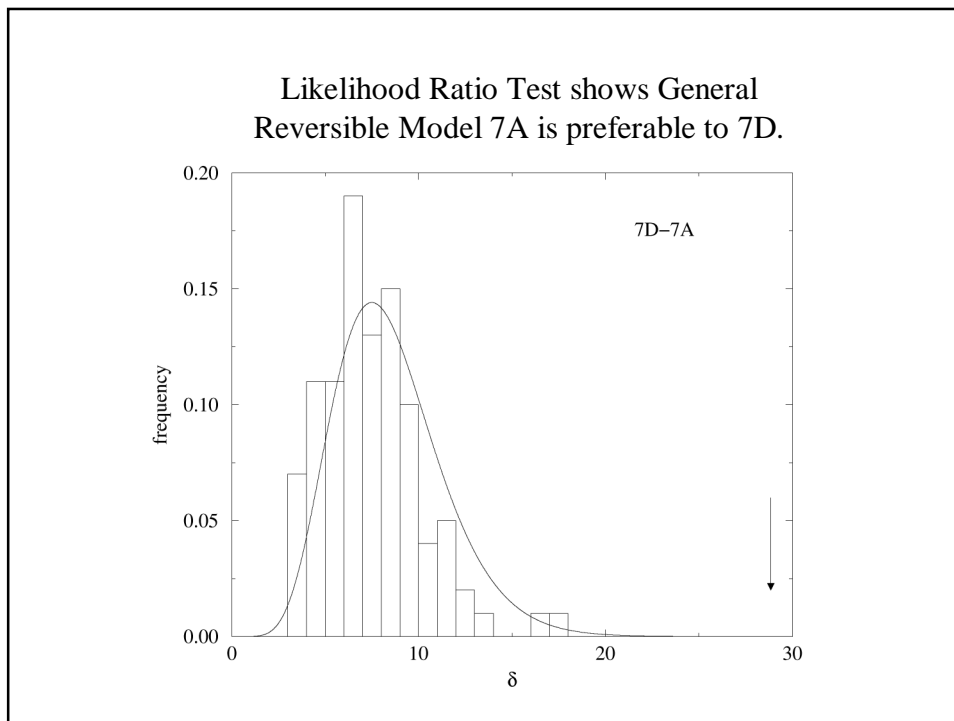
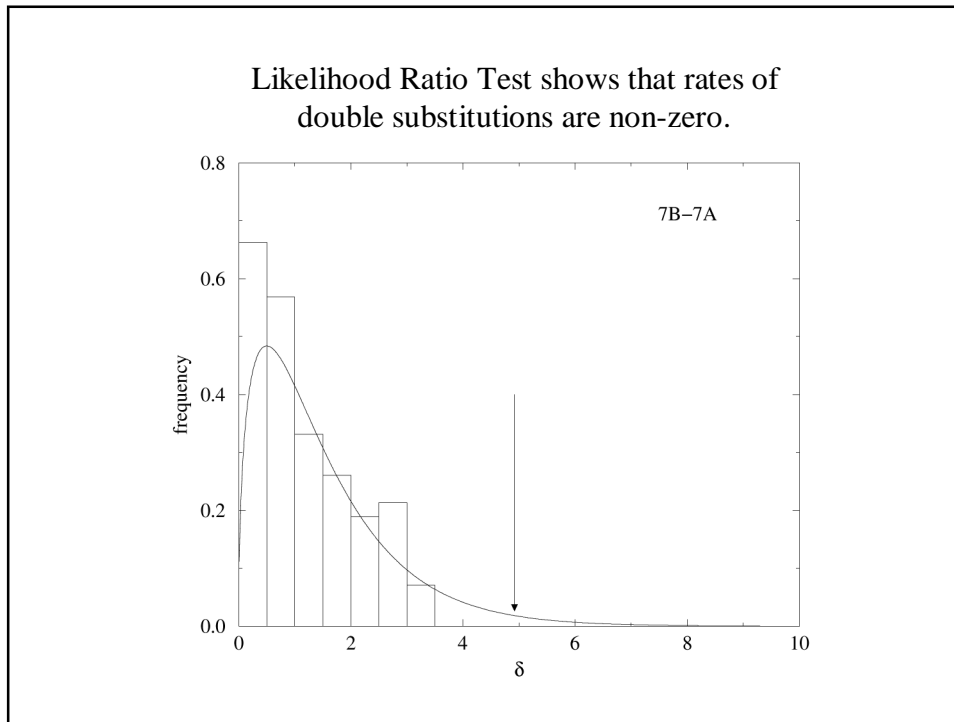
Use likelihood-based statistical tests to distinguish between models.

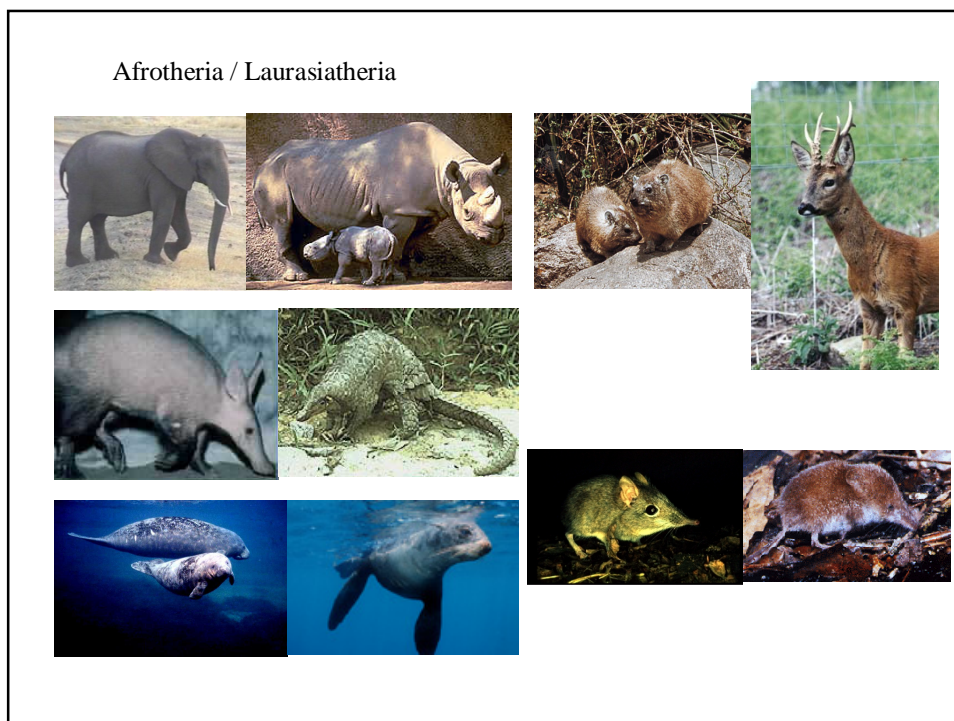
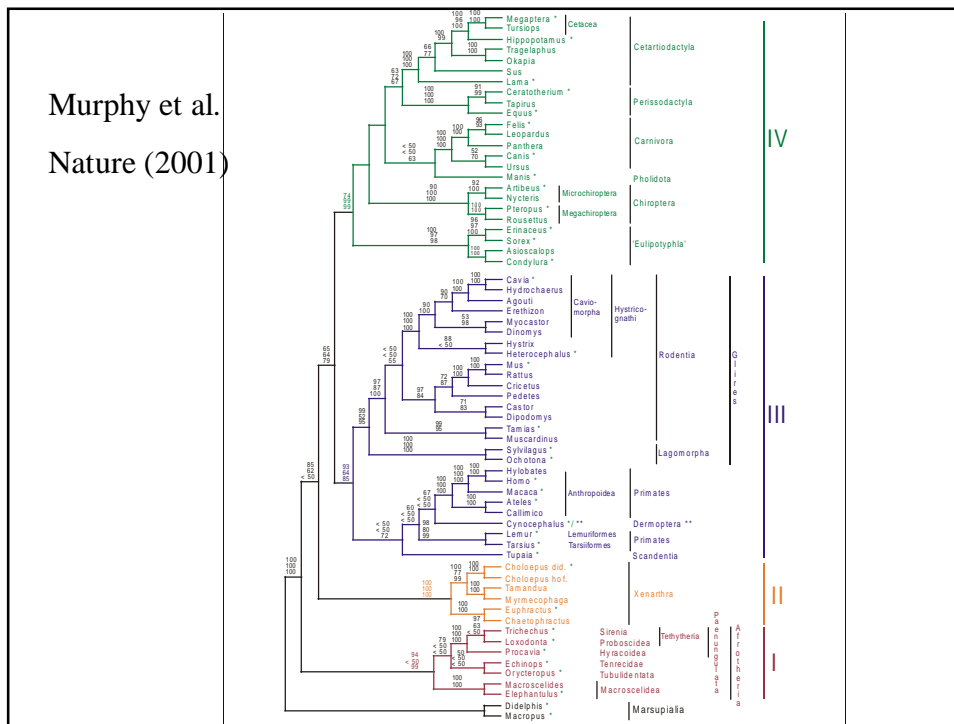


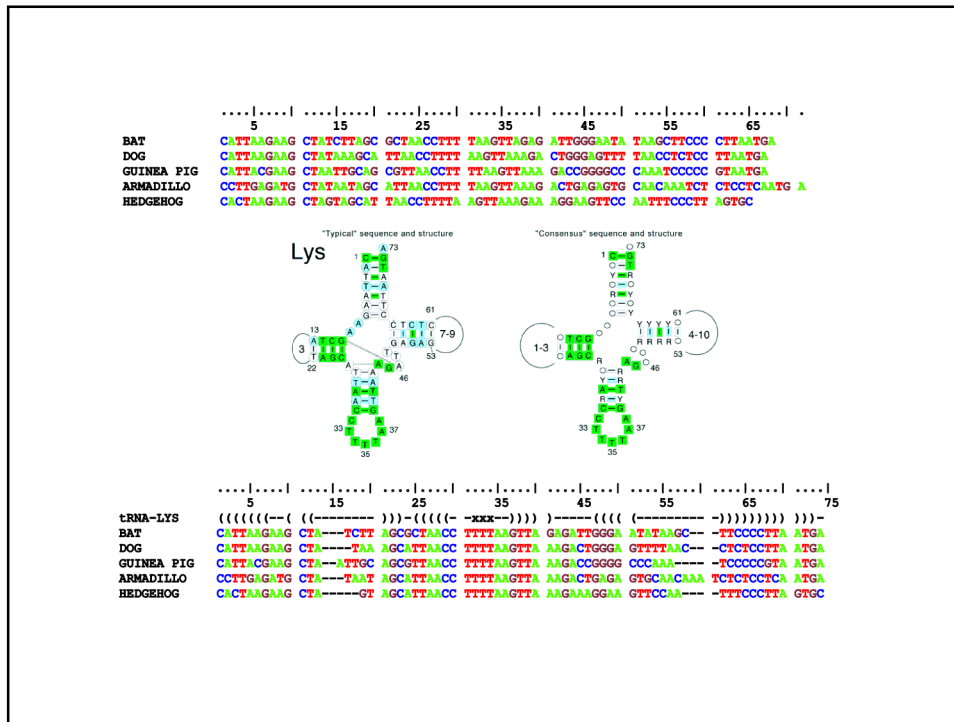


Test models using “known” phylogeny of five bacteria.
Small sub-unit ribosomal RNA sequences.









RNA Phylogeny Package

Soon to be released ...

Uses a range of evolutionary models designed specifically for RNA paired regions (as well as standard single site models).

Secondary structure is specified in the alignment as bracket notation.

Exhaustive Maximum Likelihood search for small numbers of species or small numbers of clusters.

Statistical tests to distinguish between proposed trees: Kishino-Hasegawa (important to consider correlation between sites)

Markov Chain Monte Carlo methods for tree searching.

Sub-trees in the Laurasiatheria Group

Carnivora

((Harbour Seal, Gray Seal), Dog), Cat)

Perissodactyla

((Donkey, Horse), Indian Rhino, White Rhino)

Cetartiodactyla

((Blue Whale, Sperm Whale), Hippo), (Cow, Sheep), Pig)

Chiroptera

(Long-tailed Bat, Fruit-eating Bat), (Flying Fox, Flying Fox)

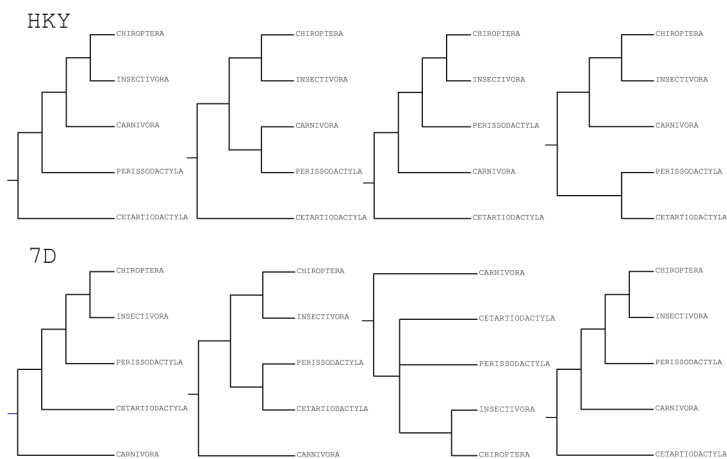
Insectivora

(Hedgehog, Mole)

Outgroup: Rabbit

Maximum Likelihood Cluster Arrangements

tRNA Data Set



Preliminary Results with MCMC

43 mammals with completely sequenced mitochondrial genomes

SSU + LSU + 14 tRNAs on the same strand

Study paired regions using the 7D model

Well-defined groups: 100% support

e.g. Primates, Marsupials, Rodents, Carnivores.

Four large-scale clades well-defined

Afrotheria (100%), Laurasiatheria (97%),
Primates+Glires (86%), Xenarthra (100%)

Insectivores+Bats (87%)

A few wandering species: e.g. Pig

Manchester
Bioinformatics



Thoughts for the day...

Thermodynamic parameters affect evolutionary parameters.

Rate models for phylogenies need to be tailored to the sequences studied.