

Lecture 2: Understanding trees and DNA substitution models

Jadranka Rota and Niklas Wahlberg

Systematic Biology Group

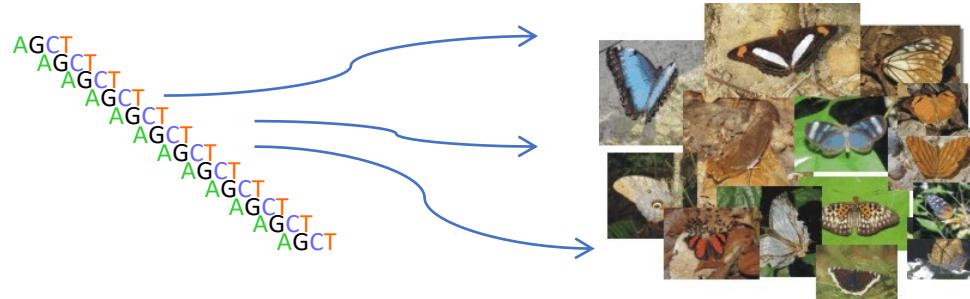
Department of Biology

Lund University



Recap: Why *molecular systematics*?

- Ease of data generation for large numbers of taxa
- Ease of generating a large number of independent data sets for given taxa
- Molecular characters behind the morphological characters we see



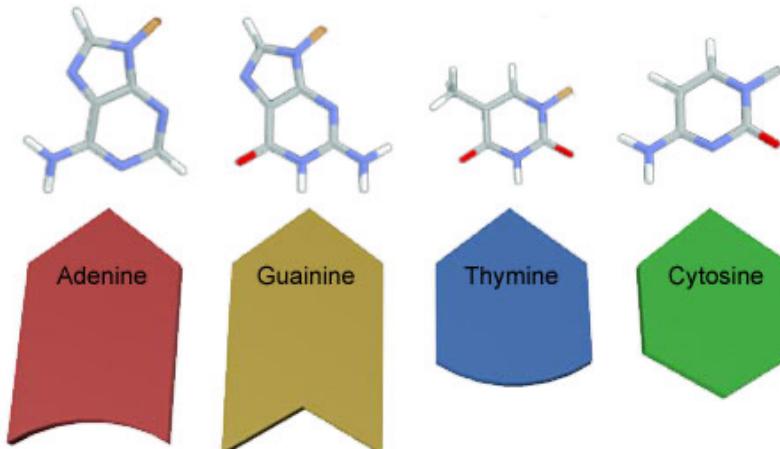
DNA as a source of information

- ▶ DNA has 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.

Homology: Definition

- Homology: similarity that is the result of inheritance from a common ancestor - identification and analysis of homologies is central to phylogenetic systematics
 - An **alignment** is a hypothesis of positional homology between bases/amino acids

Multiple Sequence Alignment

Alignment can be easy...

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 B

55 total sequences shade threshold 40 %

Mode: Edit Overwrite Selection: 0 Position: 341 Sequence Mask: None Numbering Mask: None Start ruler at: 1

MI speed slow fast

310 320 330 340 350 360 370 380 390 400

Libythea71 1 T G G A T T G C T T A T T T A A T G G A G G A T T T A G G A T T T A G G A A T T T T T A T A T
Actinote90 1 T G G A C A G T T T A C C C T C C T T C T T A A T A A G A G G A T T A G G A T T T C C T T C T T A T A T
Adelpha107 1 T G G A C A G T T T A C C C T C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A T A T
Aglais63 3 T G G A C A G T T T A C C C T C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Agraulis 24 T G G A C A G T T T A C C C T C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Amnosia101 1 T G G A C A G T T T A C C C T C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Anartia30 3 T G G A C A G T T T A C C C T C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Antanartia65 T G G A C A G T T T A C C C T C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Anthanassa12 T G G A C G T T T A C C C C C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Antirrhina109 T G G A C G T T T A C C C C C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Araschnia39 T G G A C G T T T A C C C C C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Archaeoprepa T G G A C G T T T A C C C C C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Asteroeca82 1 T G G A C G T T T A C C C C C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Caligo70 10 T G G A C A G T G A C C C C C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Calinaga64 3 T G G A C G T C A C C C C C C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Castilia76 2 T G G A C G T C A C C C C C C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Catocropte88 T G G A C G T C A C C C C C C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Catonephele6 T G G A C G T T T A C C C C C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Cercyonis8 1 T G G A C G T T T A C C C C C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Chersonesia1 T G G A C G T T T A C C C C C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Chlosyne62 1 T G G A C G T T T A C C C C C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Clossiana76 T G G A C G T C A C C C C C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Colobura68 1 T G G A C G T T T A C C C C C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Cyn thymoda T G G A C G T T T A C C C C C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Danaus108 21 T G G A C G T T T A C C C C C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Dichorragial T G G A C G T T T A C C C C C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Doleschallia T G G A C G T T T A C C C C C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Doxocopa lau T G G A C G T T T A C C C C C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Dynamine115 T G G A C G T T T A C C C C C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Eresia92 5 T G G A C G T C A C C C C C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Euclides proc T G G A C G T C A C C C C C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Euphydryas13 T G G A C G T T T A C C C C C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Euploea70 8 T G G A C G T T T A C C C C C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Euptoleta94 T G G A C G T T T A C C C C C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Gnathothric89 T G G A C G T T T A C C C C C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Greta70 9 T G G A C G T G A C C C C C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T
Hamadryas62 T G G A C G T A T C C C C C C T T C T T A A T C A A T G G A G G A T T A G G A T T T C C T T C T T A A T A T

start Poly Molecular methods to... BioEdit Sequence Align... 20:32

...or difficult

BioEdit Sequence Alignment Editor - [C:\Documents and Settings\Niklas Wahlberg\My Documents\Winclada\Geometridae\Geometrid_D2b.fst]

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

41 total sequences shade threshold 63% Position: Numbering Mask:None Start ruler at: 1

Mode: Select / Hide Selection: 0 Position: Sequence Mask: None Numbering Mask: None Start ruler at: 1

30 40 50 60 70 80 90 100 110 120 130 140 150 160

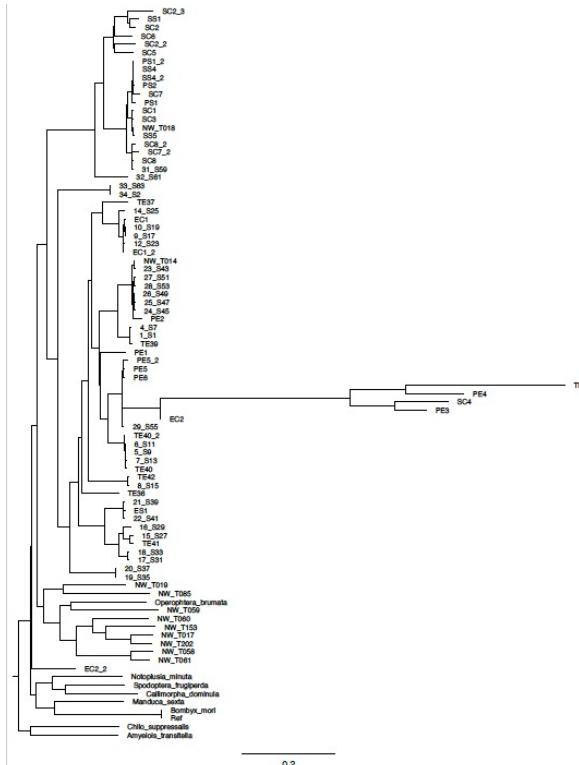
Agrotis sege G G A G G A G G G A A G G C G A G G G A C G T G G
Spodoptera l G G A G G A G G G A A G G C G A G G G A C G T G G
Drepana acuta G G A G G A G G G A A G G C G A G G G A C G T G G
Ocellia glauca G G A G G A G G G A A G G C G A G G G A C G T G G
Drepana falcata G G A G G A G G G A A G G C G A G G G A C G T G G
Drepana curvata G G A G G A G G G A A G G C G A G G G A C G T G G
Scopula immo G G A G G A G G G A A G G C G A G G G A C G T G G
Scopula ornata G G A G G A G G G A A G G C G A G G G A C G T G G
Idea basilea G G A G G A G G G A A G G C G A G G G A C G T G G
Idea strama G G A G G A G G G A A G G C G A G G G A C G T G G
Idea leucosama G G A G G A G G G A A G G C G A G G G A C G T G G
Operophtera brumata G G A G G A G G G A A G G C G A G G G A C G T G G
Operophtera brumata G G A G G A G G G A A G G C G A G G G A C G T G G
Ecliptoperla fimbriata G G A G G A G G G A A G G C G A G G G A C G T G G
Hydriomena inornata G G A G G A G G G A A G G C G A G G G A C G T G G
Hydriomena inornata G G A G G A G G G A A G G C G A G G G A C G T G G
Jodis putata G G A G G A G G G A A G G C G A G G G A C G T G G
Archaeia prochalybeata G G A G G A G G G A A G G C G A G G G A C G T G G
Colotois pennaria G G A G G A G G G A A G G C G A G G G A C G T G G
Eupithecia disjuncta G G A G G A G G G A A G G C G A G G G A C G T G G
Agriophila aurata G G A G G A G G G A A G G C G A G G G A C G T G G
Alsophilinae ae G G A G G A G G G A A G G C G A G G G A C G T G G
Lycia laponica G G A G G A G G G A A G G C G A G G G A C G T G G
Eupalus pini G G A G G A G G G A A G G C G A G G G A C G T G G
Lbd G G A G G A G G G A A G G C G A G G G A C G T G G
Eall G G A G G A G G G A A G G C G A G G G A C G T G G
Tri G G A G G A G G G A A G G C G A G G G A C G T G G
Scl G G A G G A G G G A A G G C G A G G G A C G T G G
Ppl G G A G G A G G G A A G G C G A G G G A C G T G G
Fml G G A G G A G G G A A G G C G A G G G A C G T G G
Ffl G G A G G A G G G A A G G C G A G G G A C G T G G
Pdl G G A G G A G G G A A G G C G A G G G A C G T G G
Fad1 G G A G G A G G G A A G G C G A G G G A C G T G G
Fad1 G G A G G A G G G A A G G C G A G G G A C G T G G
Ob1 G G A G G A G G G A A G G C G A G G G A C G T G G
Mrl G G A G G A G G G A A G G C G A G G G A C G T G G
Lhd G G A G G A G G G A A G G C G A G G G A C G T G G
I11 G G A G G A G G G A A G G C G A G G G A C G T G G
Ibd G G A G G A G G G A A G G C G A G G G A C G T G G
Hpl G G A G G A G G G A A G G C G A G G G A C G T G G

BioEdit Sequence Alignment Editor

start inbox for niklas... Geometridae Microsoft Power... Windows Media... BioEdit Sequence... 9:05

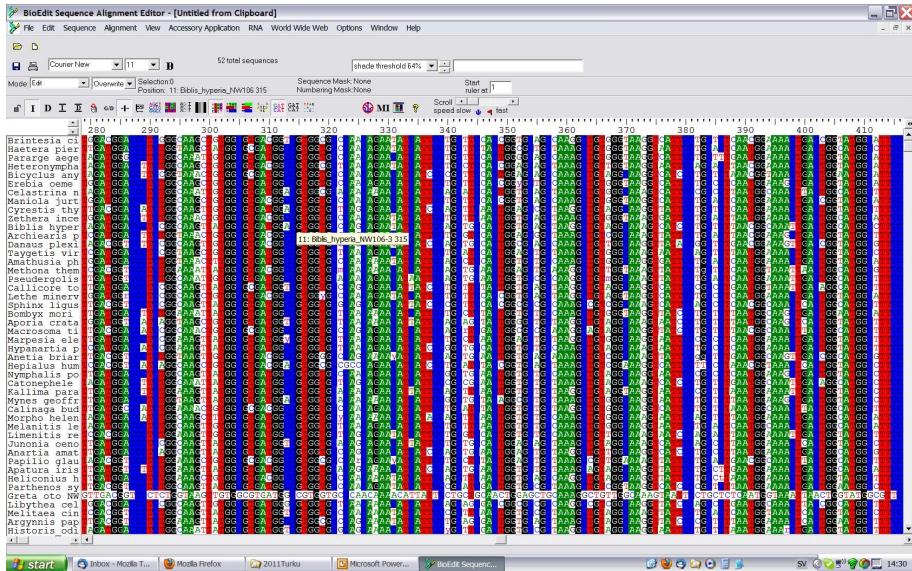
Multiple sequence alignment- goals

- To generate a concise, information-rich summary of sequence data
- Alignments can be treated as models that can be used to test hypotheses
- Does this model of events accurately reflect known biological evidence?



Multiple sequence alignment

- Manual
- Dynamic programming
- Heuristic methods
 - Progressive alignment
 - Consistency-based scoring
 - Iterative refinement methods



Manual alignment – reasons

- Might be carried out because:
 - Alignment is easy
 - There is some extraneous information (structural)
 - Automated alignment methods have encountered a local minimum problem
 - An automated alignment method can be “improved”

Protein-coding genes can often be manually aligned

BioEdit Sequence Alignment Editor - [Untitled from Clipboard]

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

Mode: Edit Overwrite Selection: 0 Position: 11: Biblis_hyperi_NW106 315 Sequence Mask:None Numbering Mask:None Start ruler at: 1

280 290 300 310 320 330 340 350 360 370 380 390 400 410

Brintesia cil GCAAGGAA GCGCAACGGC GCGCGGGT GAGGAAAGA AAGTTTCA CCGGGGAG GGGGATGGGCA GAGAAGGAA GGAAAGGAA
Haetera pier TCAAGGAA GCGCAACGGC GCGCGGGT GAGGAAAGA AAGTTTCA CCGGGGAG GGGGATGGGCA GAGAAGGAA GGAAAGGAA
Baraga aer GCGGGGAA GCGCAACGGC GCGCGGGT GAGGAAAGA AAGTTTCA CCGGGGAG GGGGATGGGCA GAGAAGGAA GGAAAGGAA
Hypopygia hyphophora AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Bicyclus any AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Erebia oeme GCA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Celestrina n AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Maniola jurt AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Cyrestis thy TCAAGGAA GCGCAACGGC GCGCGGGT GAGGAAAGA AAGTTTCA CCGGGGAG GGGGATGGGCA GAGAAGGAA GGAAAGGAA
Zerynthia p. AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Biblis hyper AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Archaeiris p. AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Danaus plexi AGAGGGAA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Taygetis vir CCA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Amathusia ph AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Methon a. AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Pseudergolina t. AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Callicore to AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Lethe minervy AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Sphinx ligus TGGGGGAA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Bombyx mori TCA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Aporia crataegi AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Machaon t. AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Marpesia ele AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Hypenartia p. AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Anetia briar AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Hepialis hur AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Nymphalis po AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Cynthia娃 AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Kallima para AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Mynes geoffr AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Calinaga bud AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Morpho helen AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Melanitis le AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Limenitis archon AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Junonia ceno AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Anartia amat AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Papilio glau AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Apatura iris AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Heliconius h AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Parthenos sylvia AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Graellsia ottocolorata AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Libythea cel AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Melitaea cin AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Argynnis pap AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA
Historis odi AGA GCGA TCGGAGAACGGG GCGGGCGGA GAGAAGGAA GGAAAGGAA

start Mozilla T... Mozilla Firefox 2011Turku Microsoft Power... BioEdit Sequenc... SV 14:30

How to align these sequences:

AGGGCTTTAA

AGGCTA

AATGGCTCTAA

GGAGCCCTAA

How to align these sequences:

A-GGGCTTAA

A--GGCT--A-

AATGGCTCTAA

GGAG-CCCTAA

How to align these sequences:

-AGGGCTTAA

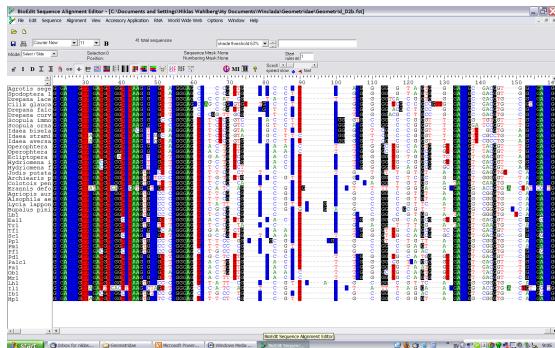
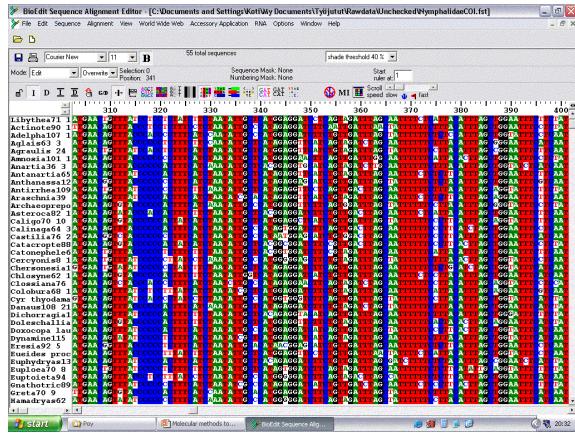
-A-GGC--TA-

AATGGCTCTAA

-GGAGCCCTAA

Multiple sequence alignment

- Is not easy! How to be objective?
- Dynamic programming
- Heuristic methods
 - Progressive alignment
 - Consistency-based scoring
 - Iterative refinement methods



Dynamic programming

- For two sequences, the best alignment can be found by scoring all possible pairs of aligned nucleotides and penalizing gaps
- An optimality criterion
- Time and computer memory needed grows exponentially with number of sequences
- Becomes impossible to align more than 4 sequences of modest length
- Fails to fully exploit phylogeny and does not incorporate an evolutionary model

Heuristics: Progressive alignment

- Devised by Feng and Doolittle in 1987
- A heuristic method and as such is not guaranteed to find the ‘optimal’ alignment
- Requires $n-1+n-2+n-3\dots n-n+1$ pairwise alignments as a starting point
- Most successful implementation is Clustal and MAFFT

Bottom line

- Alignments are extremely important in phylogenetics
- A bad alignment means many wrong statements of homology, which means pure rubbish as output
- A good alignment can be hard to attain

The Tree

Finding the optimal trees

Numbers of possible trees for N taxa

1	1
2	1
3	1
4	3
5	15
6	105
7	945
8	10395
9	135135
10	2027025
11	34459425
12	654729075
13	13749310575
14	316234143225
15	7905853580625
16	... dwarfed by rare giant elliptical galaxies, which can be 20 times more massive. By measuring the number and luminosity of observable galaxies, astronomers put current estimates of the total stellar population at roughly 70 billion trillion (7×10^{22}). ... 0002000870762850625
20	221643095476699771875 (2×10^{20})
50	3×10^{74}

How can
we find
the most
optimal
tree?

<https://skyandtelescope.org/astronomy-resources/how-many-stars-are-there/>

Finding optimal trees - exact solutions

- Exact solutions can only be used for small numbers of taxa
- **Exhaustive search** examines all possible trees
- **Branch and bound** does not examine all trees, but will find optimal tree(s)
- Typically used for problems with 10–20 taxa

Finding optimal trees - heuristics

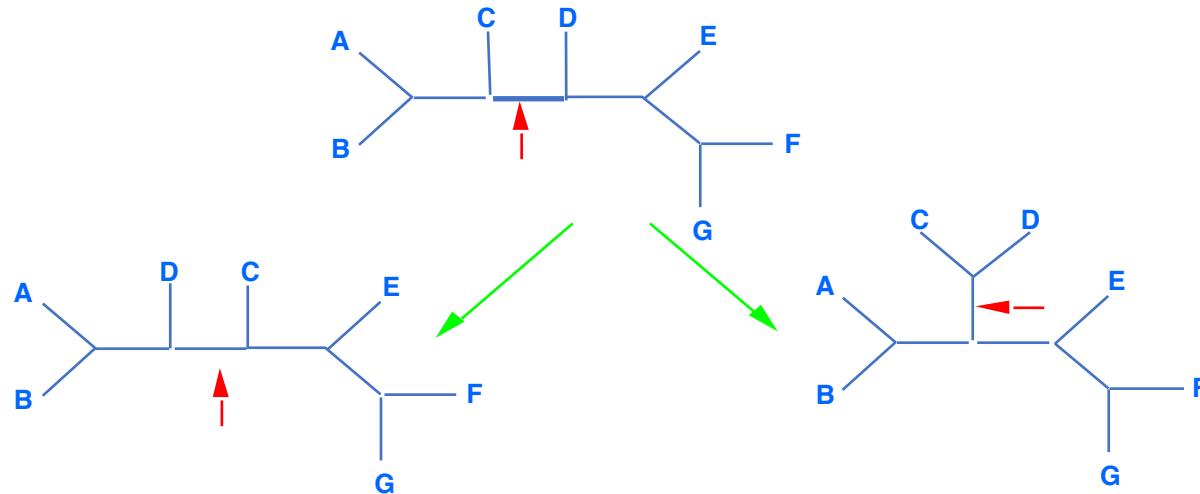
- The number of possible trees increases faster than exponentially with the number of taxa
- Exhaustive searches impractical for many data sets (an NP-complete problem)
- Heuristic methods are used to search tree space for optimal trees by building or selecting an initial tree and swapping branches to search for better ones
- The trees found are not guaranteed to be optimal – they are best guesses

Finding optimal trees – branch swapping

- Nearest neighbor interchange (NNI)
- Subtree pruning and regrafting (SPR)
- Tree bisection and reconnection (TBR)

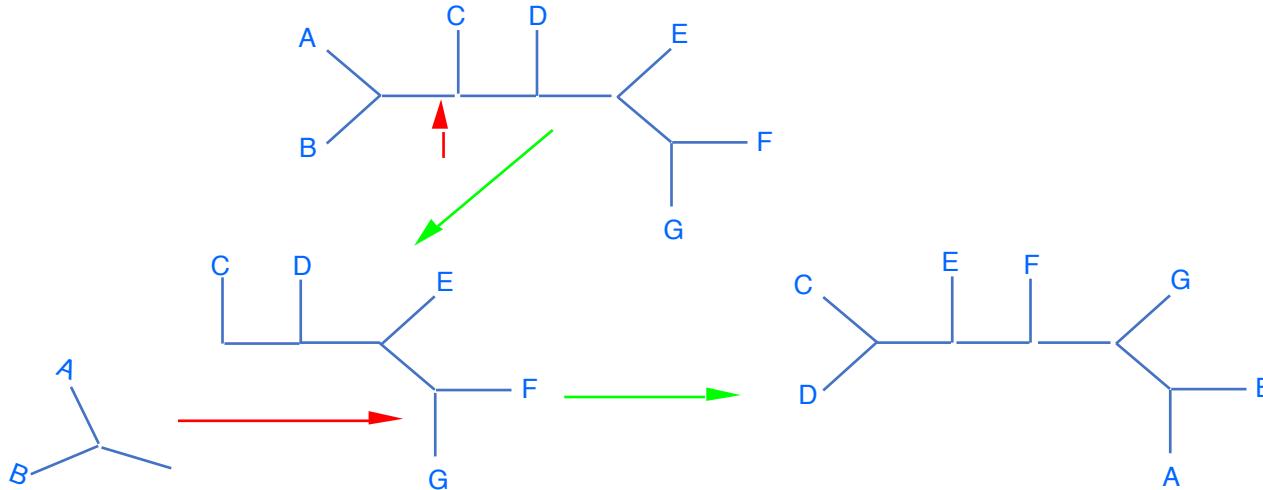
Moving through treespace

Nearest neighbor interchange (NNI)



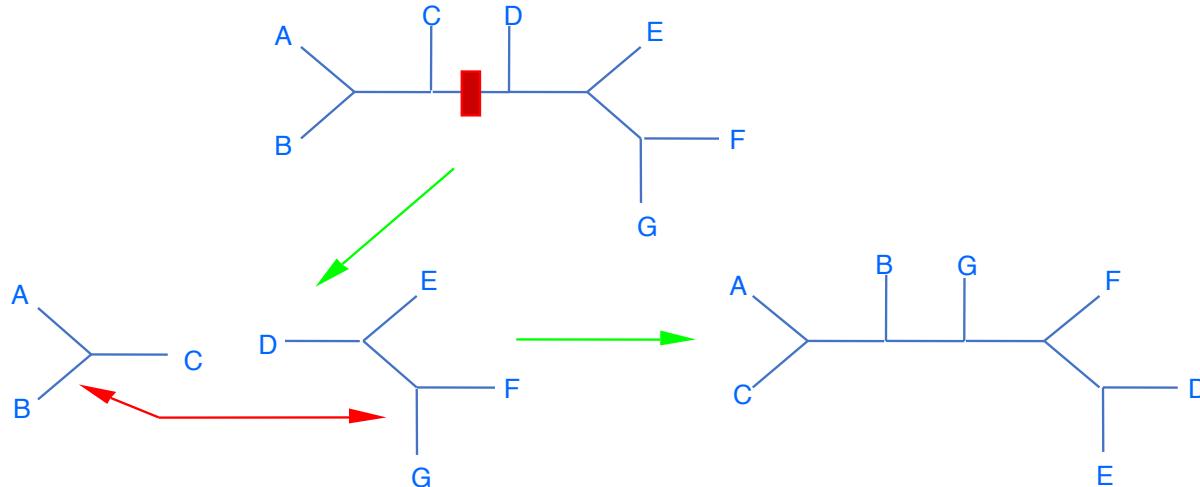
Moving through treespace

Subtree pruning and regrafting (SPR)



Moving through treespace

Tree bisection and reconnection (TBR)



Consensus methods

Multiple optimal trees

- Many methods can yield multiple equally optimal trees
- We can further select among these trees with additional criteria, but
- Typically, relationships common to all the optimal trees are summarised with *consensus trees*

Consensus methods

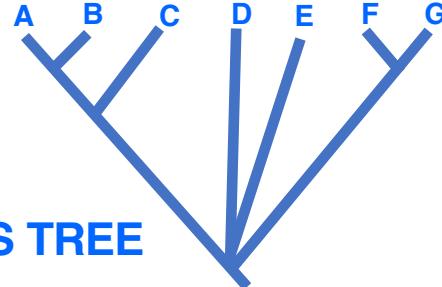
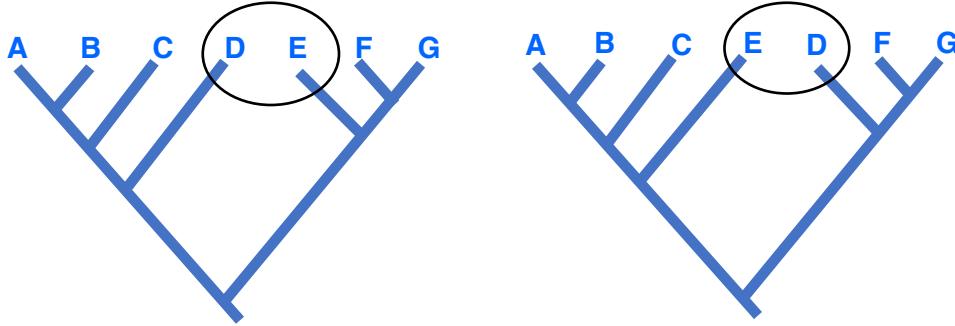
- A consensus tree is a summary of the agreement among a set of fundamental trees
- There are many consensus methods that differ in:
 1. the kind of agreement
 2. the level of agreement
- Consensus methods can be used with multiple trees from a single analysis or from multiple analyses

Strict consensus methods

- Strict consensus methods require **agreement across all the fundamental trees**
- They show only those relationships that are supported by the parsimonious interpretation of the data
- This method produces a consensus tree that includes all and only those full splits found in all the fundamental trees
- Other relationships (those in which the fundamental trees disagree) are shown as unresolved polytomies

Strict consensus methods

TWO FUNDAMENTAL TREES



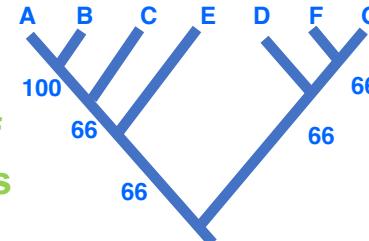
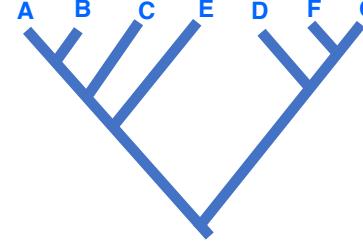
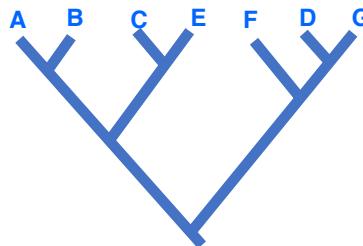
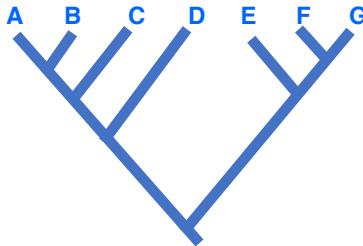
STRICT CONSENSUS TREE

Majority-rule consensus methods

- Majority-rule consensus methods require agreement across a majority of the fundamental trees
- May include relationships that are not supported by the most parsimonious interpretation of the data
- This method produces a consensus tree that includes all and only those full splits found in a majority (>50%) of the fundamental trees
- Other relationships are shown as unresolved polytomies
- Of particular use in bootstrapping

Majority rule consensus

THREE FUNDAMENTAL TREES

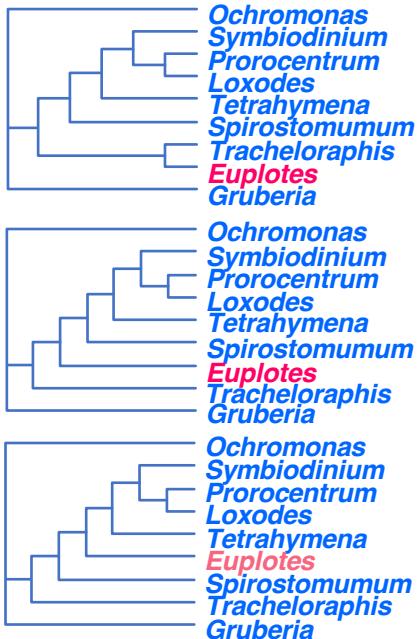


Numbers indicate frequency of
clades in the fundamental trees

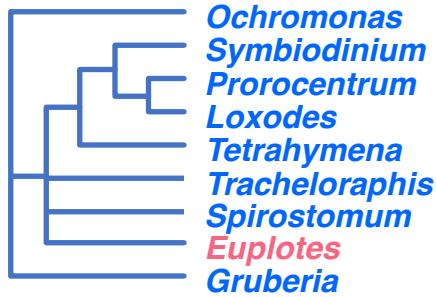
MAJORITY-RULE CONSENSUS TREE

Consensus methods

Three fundamental trees

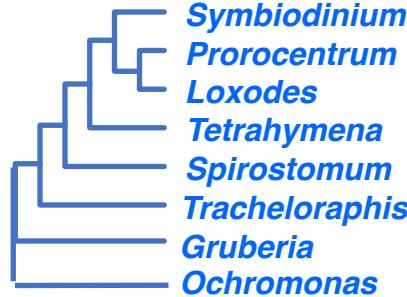


Strict (component)

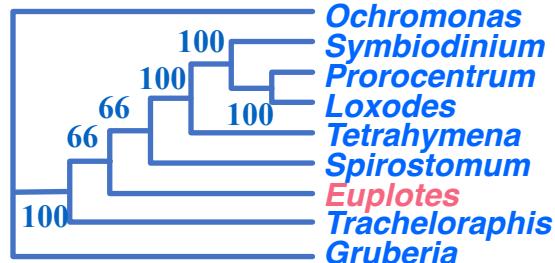


Strict reduced cladistic

Euplates excluded



Majority-rule



Consensus methods – use

- Currently majority-rule methods mainly used
 - Bootstrapping
 - Bayesian methods
- Reduced methods can be useful to identify problem taxa (rogue taxa)
 - E.g. RogueNaRok
- Strict methods mainly used in parsimony analyses
 - Rarely used with molecular data

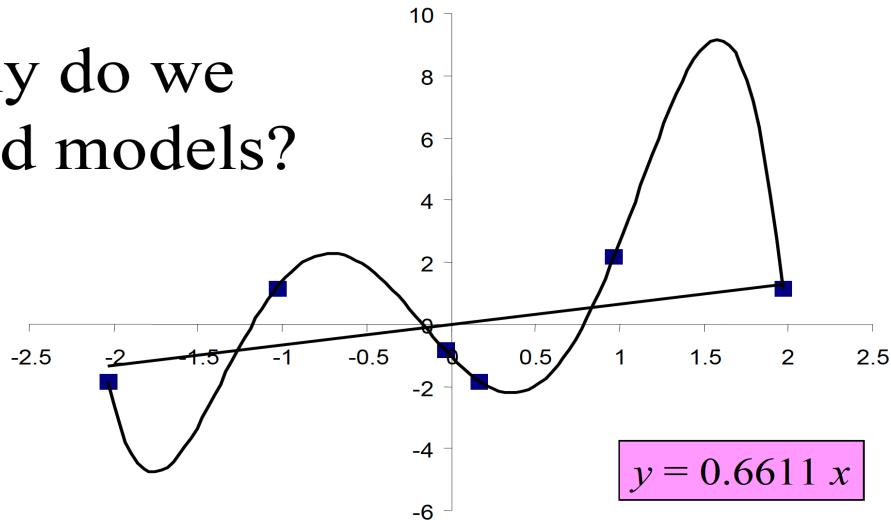
Take home messages

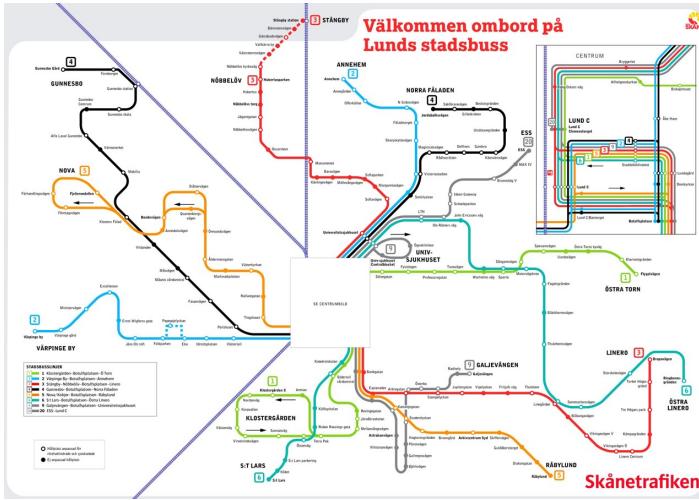
- Statements of homology are the basis of phylogenetics
- Alignments of molecular sequences are very strong statements of positional homology
- Finding an optimal tree is not a trivial task

Modelling DNA Sequence Evolution

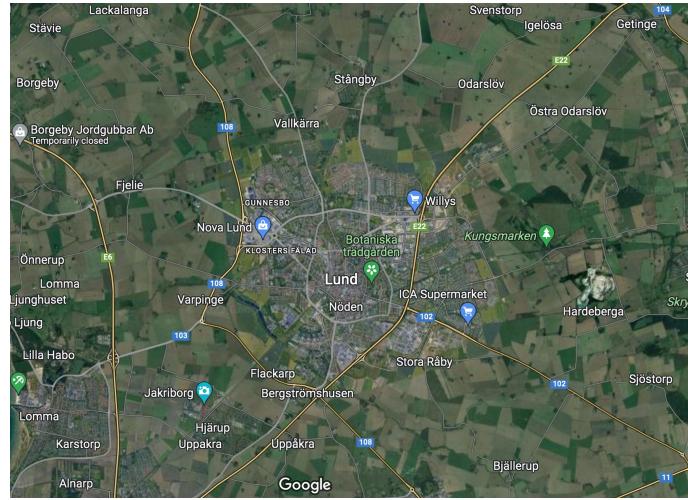
$$y = -1.5972 x^5 + 23.167 x^4 - 126.18 x^3 + 319.17 x^2 - 369.22 x + 155.67$$

Why do we
need models?





A simplified map of bus routes in Lund



A realistic map of Lund

- Which one would you use to get around Lund by bus?

Models: an overview

- In general, models help us **predict** the **future based on our observations**
- With **more parameters**, models have a **better fit** to the data (observations)
- Underparamaterized models: poor fit to the observed data
- Overparameterized models: poor prediction of future observations
- Choosing best models based on different criteria
 - Likelihood ratio tests, AIC, BIC, Bayes factors

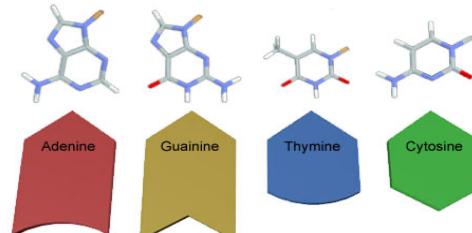
What do we model in DNA sequence evolution?

- Nucleotide substitutions
 - The rate at which each nucleotide is replaced by each alternative nucleotide

What is the challenge?

- DNA has only four characters

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.

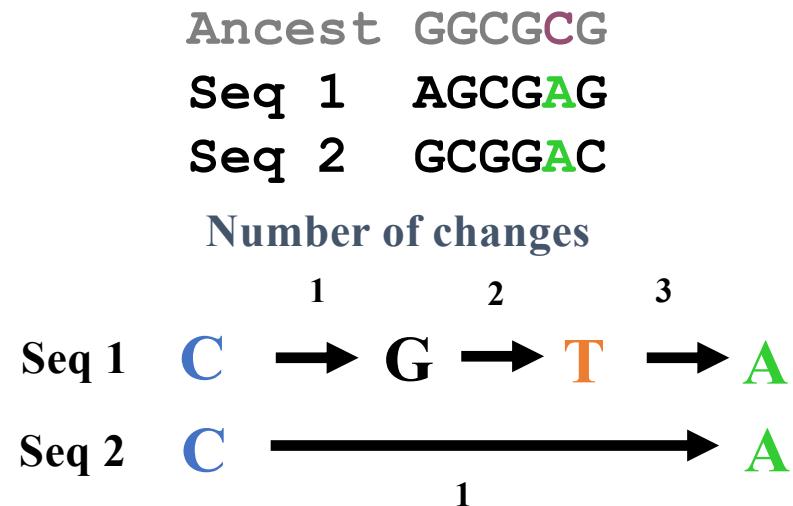
Saturation in sequence data

- Saturation is due to **multiple changes 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 proteins often codon position 3)
- In severe cases the data become essentially random and all information about relationships can be lost

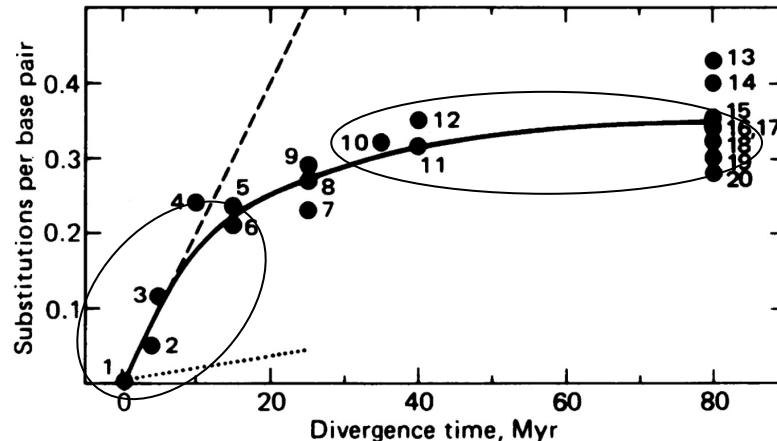
Multiple changes at a single site
- hidden changes

Seq 1	AGCGAG
Seq 2	GCGGAC

Multiple changes at a single site
- hidden changes



“Multiple hits” or saturation



<https://www.pnas.org> › doi › pnas.76.4.1967

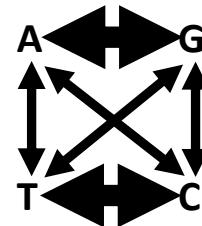
Rapid evolution of animal mitochondrial DNA. - PNAS

by WM Brown · 1979 · Cited by 4306 — Rapid evolution of animal mitochondrial DNA. W M Brown, M George, Jr, and A C WilsonAuthors Info & Affiliations. April 1, 1979. 76 (4) 1967-1971.

Brown et al. 1979. PNAS 76:1967

Substitution types

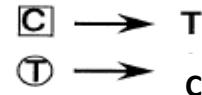
- Purines: A, G
- Pyrimidines: C, T
- Transversions
 - Pu --> Pyr
 - Pyr --> Pu
- Transitions – more common
 - Pu --> Pu
 - Pyr --> Pyr



Pur - Pyr mispairs lead to transitions

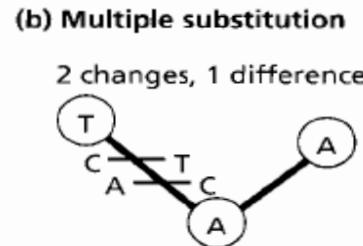


In next round of replication



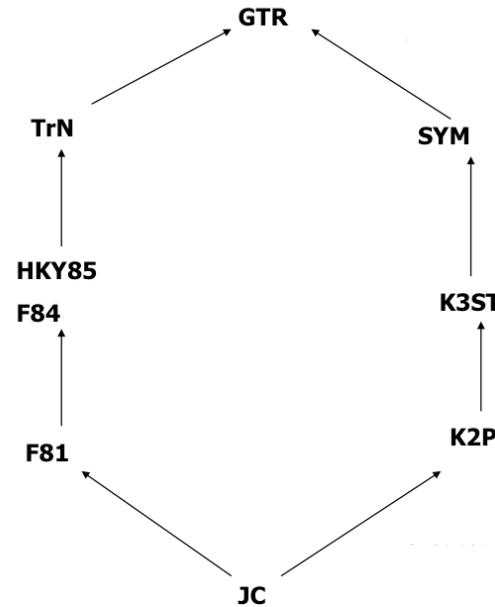
Saturation in sequence data:

- 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



Saturation in sequence data (cont.)

- 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



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 4×4 rate matrix (known as the Q matrix)
- Other model parameters may include:
 - Site by site rate variation (aka among-site rate variation – ASRV)

Corrections for multiple substitutions:

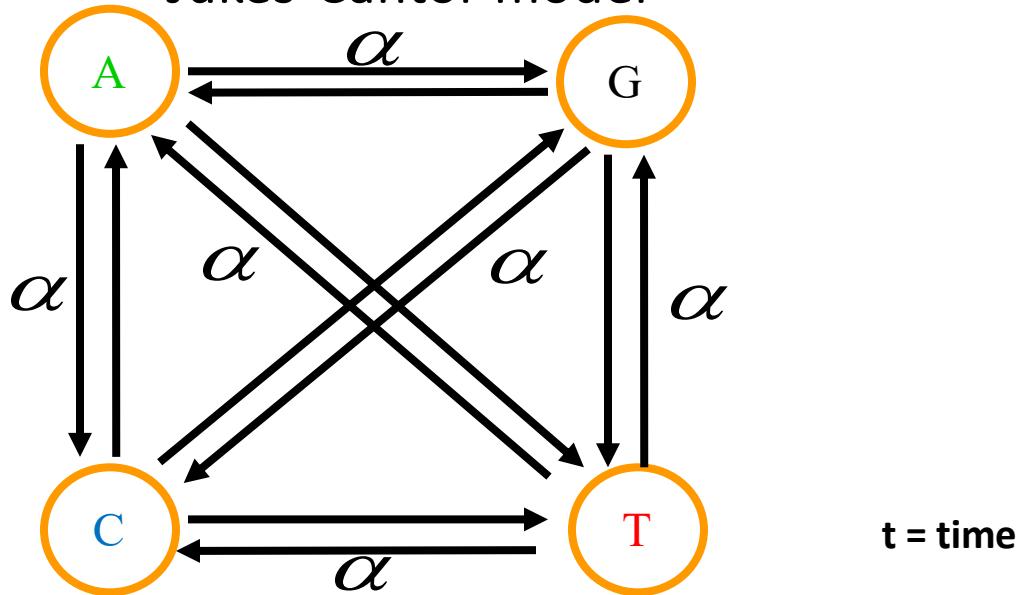
First DNA substitution model

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

Jukes-Cantor model



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

The Q matrix

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

From

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 of α

The Jukes-Cantor model: the simplest model

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

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 of α

Improvements on Jukes-Cantor

- Allow **base frequencies** to be unequal to accommodate e.g. sequences such as these

AAACCTGGATTACCGAGATTAAAGCGATATATTGCAATGC

34% A

17% C

29% T

20% G

- 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**

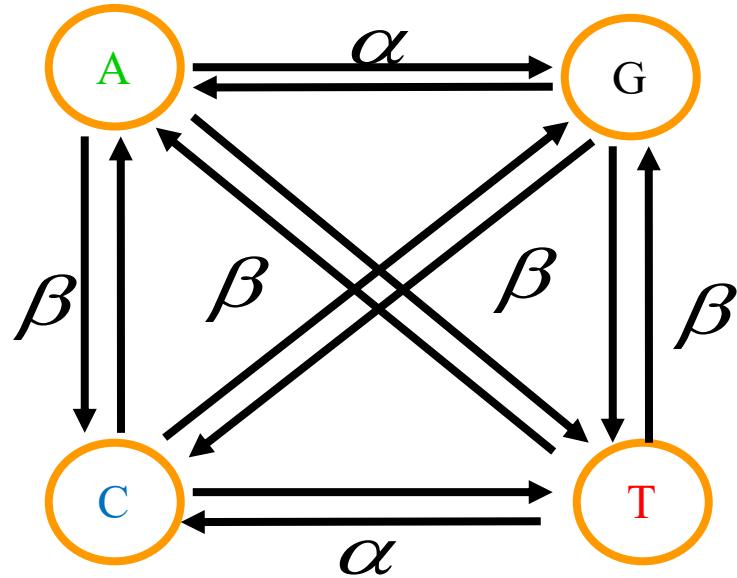
		2nd position				3rd position
1st position	U	C	A	G		
U	Phe	Ser	Tyr	Cys	U	
	Phe	Ser	Tyr	Cys	C	
	Ile	Ser	stop	stop	A	
	Leu	Ser	stop	Trp	G	
C	Leu	Pro	His	Arg	U	
	Leu	Pro	His	Arg	C	
	Leu	Pro	Gln	Arg	A	
	Leu	Pro	Gln	Arg	G	
A	Ile	Thr	Asn	Ser	U	
	Ile	Thr	Asn	Ser	C	
	Ile	Thr	Lys	Arg	A	
	Met	Thr	Lys	Arg	G	
G	Val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

Amino Acids

Ala: Alanine
Arg: Arginine
Asp: Aspartic acid
Asn: Asparagine
Cys: Cysteine
Ile: Isoleucine

Gln: Glutamine
Glu: Glutamic acid
Lys: Lysine
Leu: Leucine
Phe: Phenylalanine
Pro: Proline
Ser: Serine
Thr: Threonine
Val: Valine
Trp: Tryptophan

Kimura (1980) model: K2P



α = transitions

β = transversions

The Kimura model has 2 parameters

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

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 – α and transversions - β)

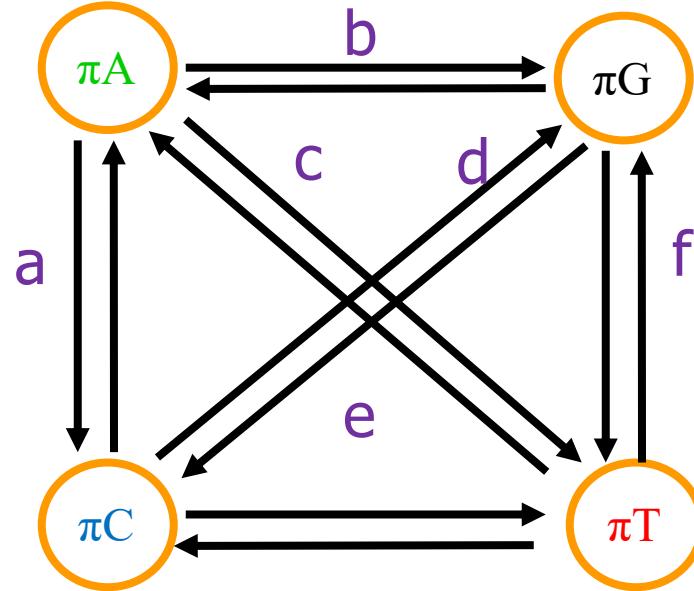
The Hasegawa-Kishino-Yano model

	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: π_A , π_C , π_G , π_T
- 2) There are two substitution types (transitions – α and transversions – β)

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	$\pi_A c$	$\pi_C e$	$\pi_G f$	—

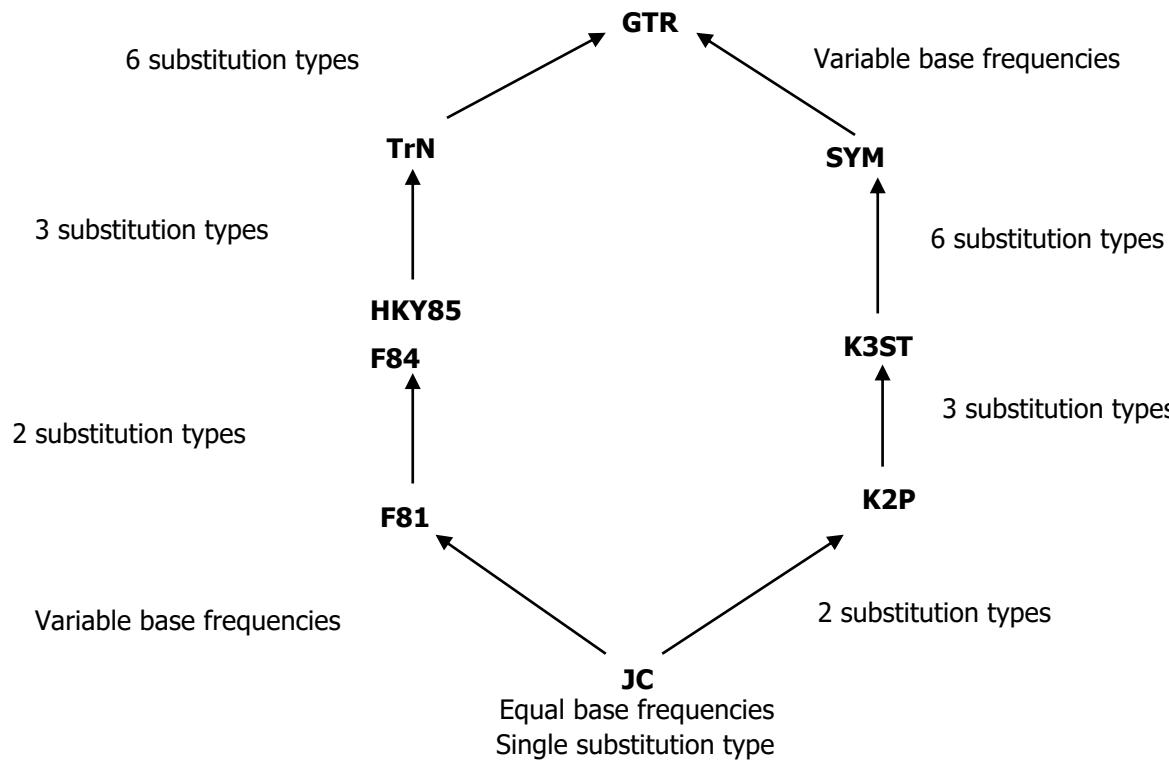
GTR model:

- 1) Base frequencies are allowed to vary: π_A , π_C , π_G , π_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

ACAGGTGAGGCTCAGCCAATTTGAGCTTTGTCGATAAGGT



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
- Variable vs. 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

REVIEWS

Among-site rate variation and its impact on phylogenetic analyses

Ziheng Yang

<https://www.sciencedirect.com/science/article/pii/0890623896900111>

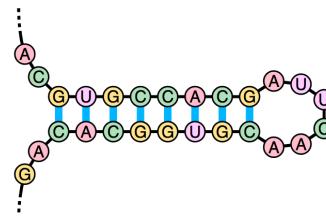
Among-site rate variation and its impact on phylogenetic ...

by Z Yang · 1996 · Cited by 1342 — Recent analyses show that failure to account for rate variation can have drastic effects, leading to biased dating of speciation events, biased...

Yang (1996) TREE 11(9): 367–372

Why is modelling ASRV important?

- Protein-coding genes – 1st, 2nd, 3rd 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	Ser	Tyr	Cys	U
	Phe	Ser	Tyr	Cys	C
	Leu	Ser	stop	stop	A
	Leu	Ser	stop	Trp	G
C	Leu	Pro	His	Arg	U
	Leu	Pro	His	Arg	C
	Leu	Pro	Gln	Arg	A
	Leu	Pro	Gln	Arg	G
A	Ile	Thr	Asn	Ser	U
	Ile	Thr	Asn	Ser	C
	Ile	Thr	Lys	Arg	A
	Met	Thr	Lys	Arg	G
G	Val	Ala	Asp	Gly	U
	Val	Ala	Asp	Gly	C
	Val	Ala	Glu	Gly	A
	Val	Ala	Glu	Gly	G

Ala: Alanine
Arg: Arginine
Asn: Asparagine
Asp: Aspartic acid
Cys: Cysteine

Gln: Glutamine
Glu: Glutamic acid
Gly: Glycine
His: Histidine
Ile: Isoleucine

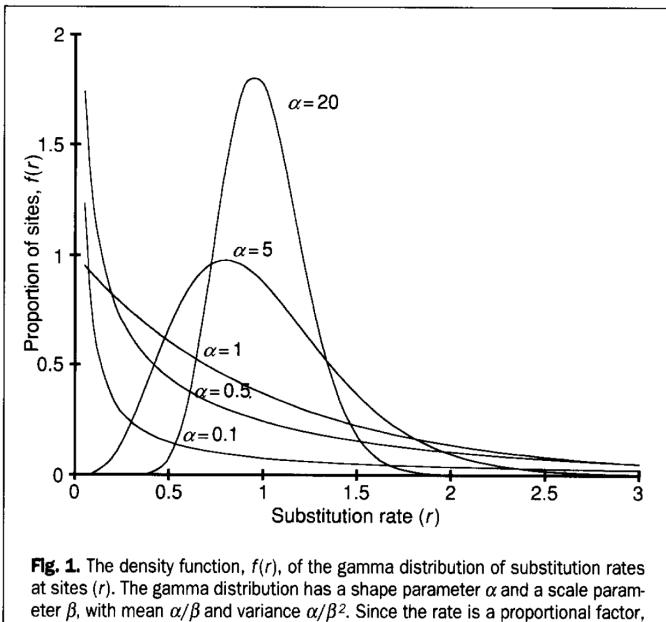
Leu: Leucine
Lys: Lysine
Met: Methionine
Phe: Phenylalanine
Pro: Proline

Ser: Serine
Thr: Threonine
Trp: Tryptophane
Tyr: Tyrosine
Val: Valine

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 among-site rate variation with Gamma distribution



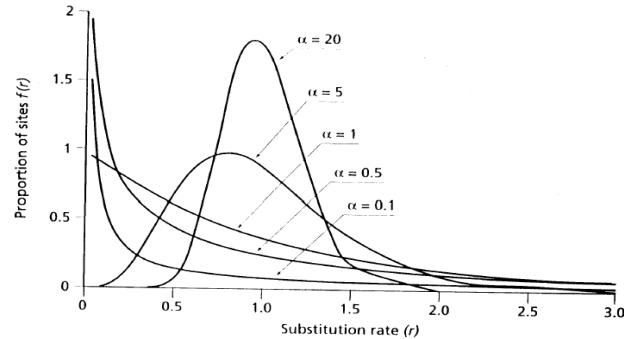
Gamma distribution:
Relative substitution rates for
different α values

Fig. 1 from Yang 1996:
Alpha – the shape parameter of the gamma distribution

Smaller alpha = higher ASRV

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

Modelling ASRV leads to greater improvement in fit than other parameters

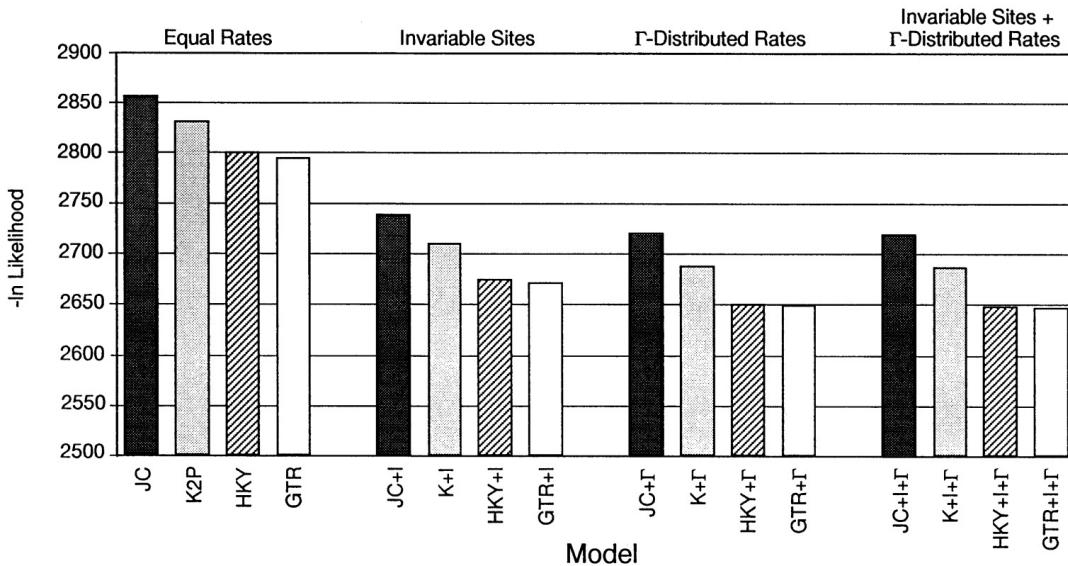


Fig. 4 from Frati et al. 1997. J. Mol. Evol. 44:145-158

Modelling ASRV leads to greater improvement in fit than other parameters

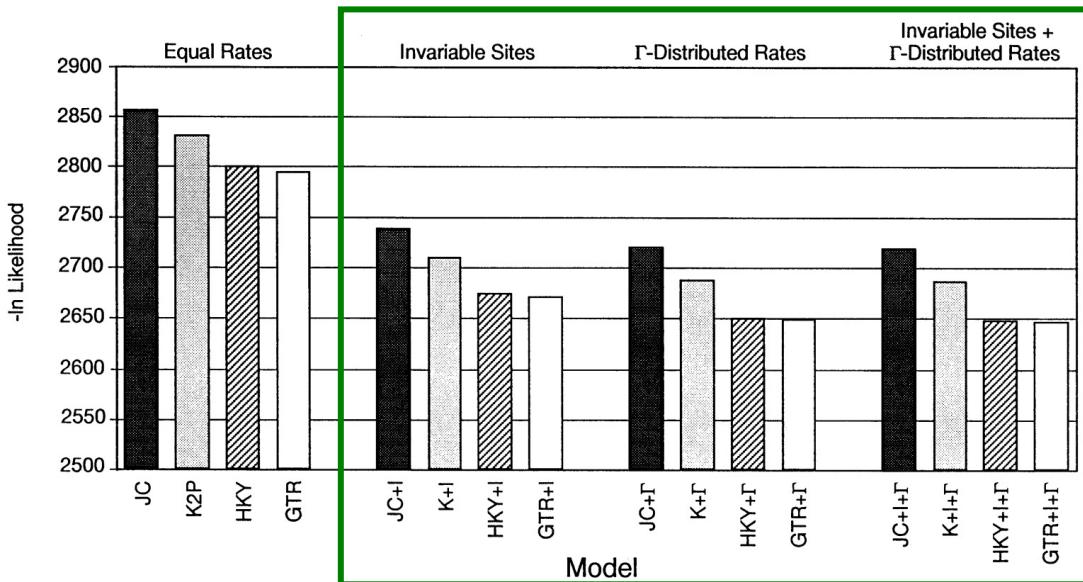
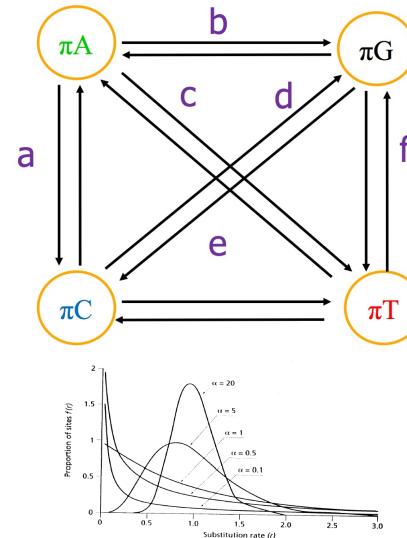


Fig. 4 from Frati et al. 1997. J. Mol. Evol. 44:145-158

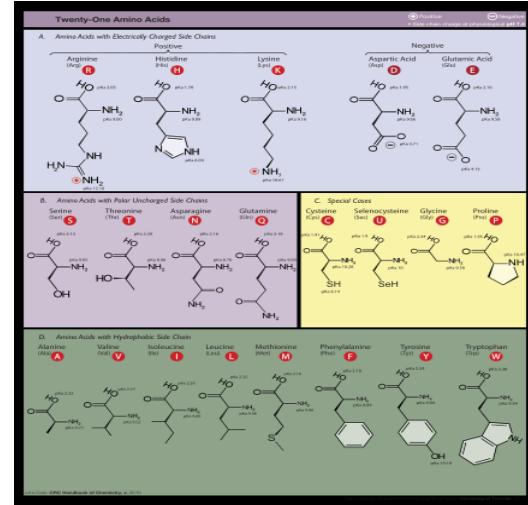
Parameters in models of DNA evolution

- Numbers of parameters estimated:
 - Base composition
 - 1 fixed, 3 estimated
 - Substitutions
 - up to 5; 1 fixed, 5 estimated
 - Among-site-rate variation
 - Gamma shape parameter = 1 parameter
 - Invariant sites = 1 parameter
 - Gamma + I = 2 parameters



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



Inferring phylogenies: methodological overview

- **Distance methods**
 - A clustering method using pairwise distances between sequences (e.g. neighbour joining)
- **Discrete characters**
 - Using an optimality criterion to choose the best tree
 - Maximum parsimony (Occam'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 – disadvantages

- Prone to systematic errors
- Problems with missing data
- Generally outperformed by Maximum Likelihood and Bayesian 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))