# DNA sequences I

Alignment and Sanger sequence editing, DNA saturation, basic single-gene tree construction

# DNA sequences I

- Alignment [Mafft]

- Alignment editing [BioEdit, MEGA]

- Alignment improvement [Gblocks, Trimal]

- Alignment conversion, concatenation [FASconCAT]

- Model selection [jModeltest, PartitionFinder]

- Tree reconstruction [RAxML, PAUP, MrBayes]

- Tree manipulation [FigTree, R]

# Sequence alignment

- http://mafft.cbrc.jp/alignment/server/

Strategy:
- ⦿ Auto (FFT–NS–1, FFT–NS–2, FFT–NS–i or L–INS–i; depends on data size) *Updated*

Progressive methods
- ○ FFT–NS–1 (Very fast; recommended for >2,000 sequences; progressive method)
- ○ FFT–NS–2 (Fast; progressive method)
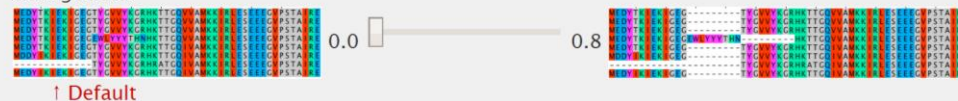- ○ G–INS–1 (Slow; progressive method with an accurate guide tree)

Iterative refinement methods
- ○ FFT–NS–i (Slow; iterative refinement method)
- ○ E–INS–i (Very slow; recommended for <200 sequences with multiple conserved domains and long gaps) Help  *Updated* (2015/Jun)
- ○ L–INS–i (Very slow; recommended for <200 sequences with one conserved domain and long gaps) Help
- ○ G–INS–i (Very slow; recommended for <200 sequences with global homology) Help
- ○ Q–INS–i (Extremely slow; secondary structure of RNA is considered; recommended for a global alignment of highly divergent ncRNAs with <200 sequences × <1,000 nucleotides; the number of iterative cycles is restricted to two, 2016/May) Help

Align unrelated segments, too? *in Alpha Testing* (2014/Mar)
If the input data is expected to be globally conserved but locally contaminated by unrelated segments, try 'Unalignlevel>0' and possibly 'Leave gappy regions'.

Unalignlevel:

0.0                     0.8

↑ Default

Gap opening penalty: 1.53   (1.0 – 5.0)
Offset value: 0.0   (0.0 – 1.0)

Score of N in nucleotide data: Example
   ↓ Long stretches of Ns tend to be gapped (excluded from the alignment).
- ⦿ (nzero) N has no effect on the alignment score.
- ○ (nwildcard) N is treated like a wildcard. *Experimental option* (2016/Apr/26)
   ↑ Try this if Ns should be aligned with usual letters.

# Sequence alignment

- *Playing with MAFFT options*



Strategy:
- FFT-NS vs Q-INS
- FFT-NS vs G-INS

Align unrelated sequences:
- 0.0 vs 0.8

Penalties:
- 1 vs 3
- 1 vs 5

Score of N in nucleotide data:
- nzero v nwildcard

Test of alignments differing by:
- locus (e.g., SSU, rbcL)
- variability
- N frequency
- …

# Sequence editing

- BioEdit - http://www.mbio.ncsu.edu/bioedit/bioedit.html



- A few tips:
  - Sequence – Manipulations – UPPERCASE
  - Sequence – Nucleic Acid – RNA->DNA
  - Alignment – Minimize alignment to mask
  - Edit – Copy sequence titles
  - Search options….

# Sequence editing

- MEGA - http://www.megasoftware.net/



- A few tips (.fas):
  - Data – Reverse Complement
  - Data – Translate/untranslate
- A few tips (.meg):
  - Distance – Compute pairwise distance
  - Phylogeny – Construct fast trees
  - Statistics (Data explorer) – Nucleotide composition
  - Highlight (Data explorer) – Mark sites

# Sequence editing

- FaBox- http://users-birc.au.dk/biopv/php/fabox/



- Edit names
- Crop and merge alignments
- Extract variable sites („Show variable sites only")
- Convert fasta to other formats (e.g., TCS input)
- Create MrBayes nexus file („Create MrBayes input file from fasta (fasta2mrbayes)"
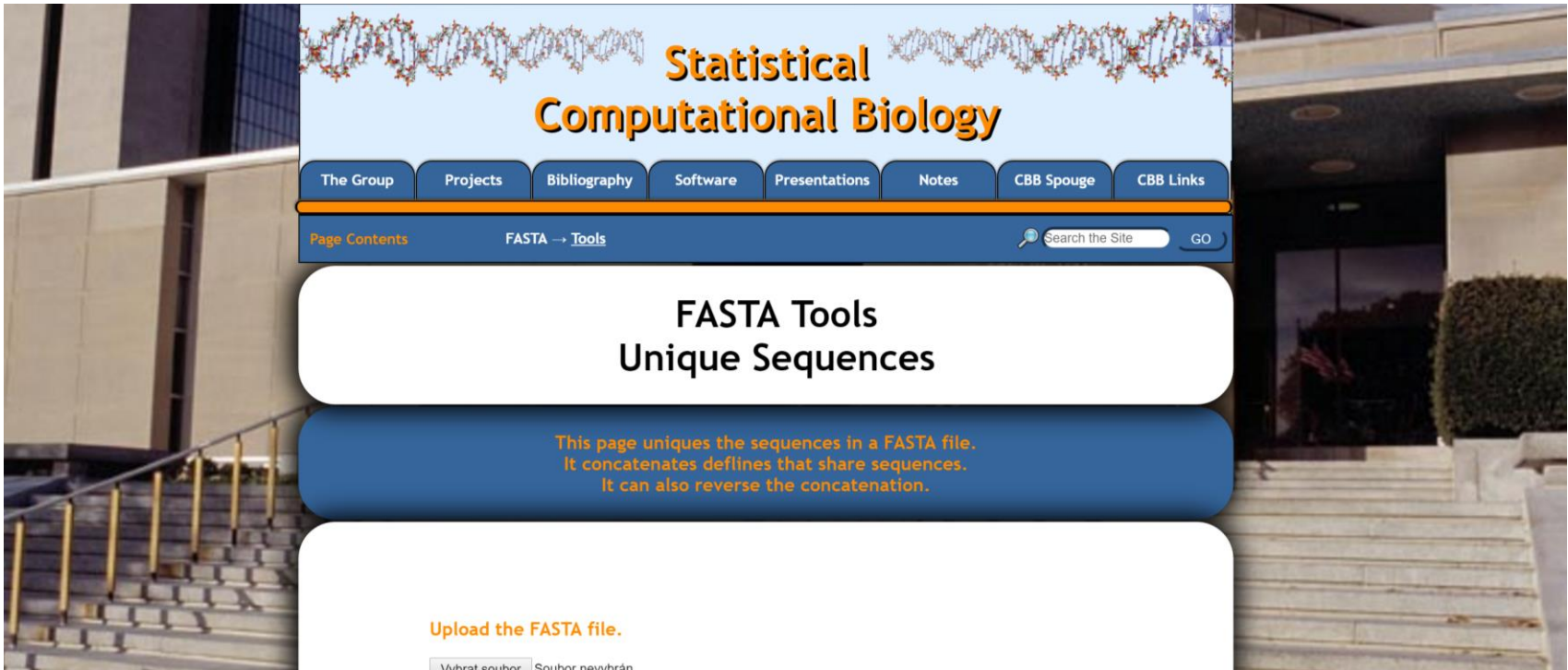
# Sequence editing

- SeqState - http://bioinfweb.info/Software/SeqState



- Coding gaps as a special state

- *File -> Load NEXUS file*
- *IndelCoder -> simple indel coding*

# Sequence editing

- FASTA Tools – Unique sequences

- https://www.ncbi.nlm.nih.gov/CBBresearch/Spouge/html_ncbi/html/fasta/uniqueseq.html
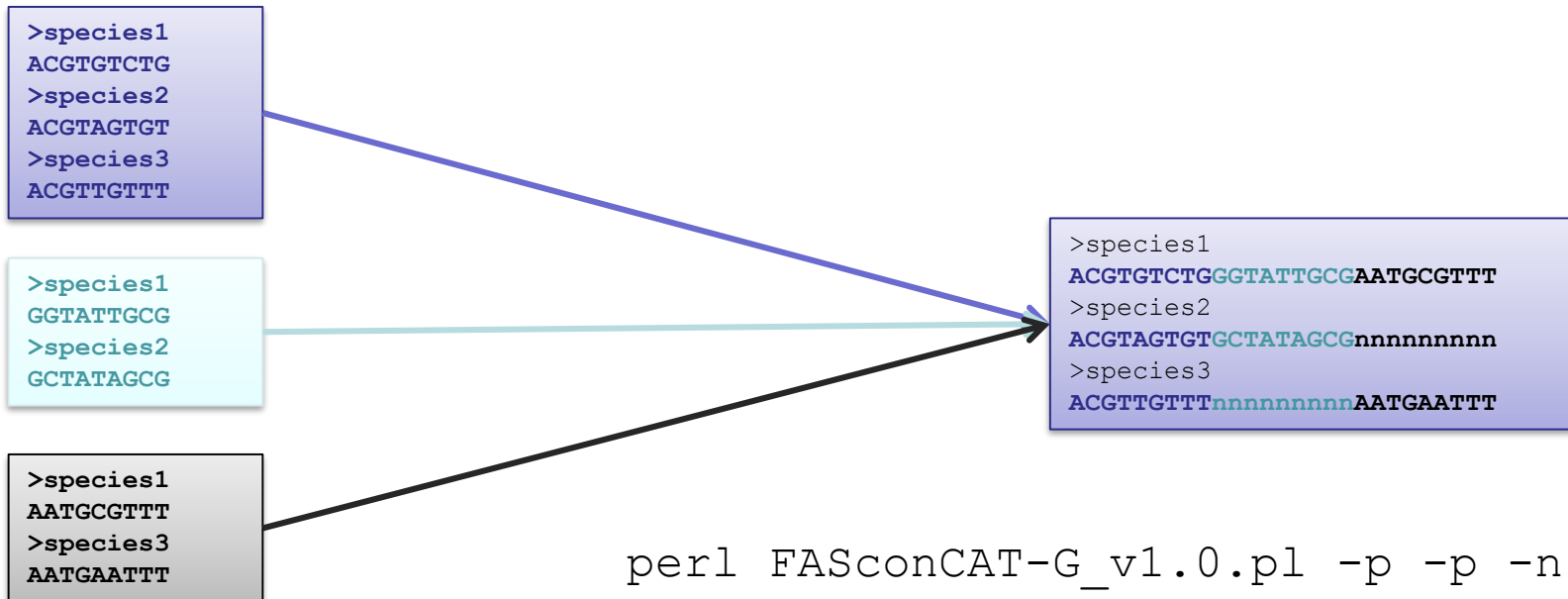
# Sequence editing

- A few tips in R:
  - Save sequence names
  - Replace sequence names
  - Minimize alignment to mask
  - Save an alignment of selected sequences
  - Concatenate alignments

# Alignment concatenations, conversion

FASconCAT – https://www.zfmk.de/en/research/research-centres-and-groups/fasconcat-g
- Perl script
- concatenating alignments (with same headers but not necessarily with all samples in all alignments)
- conversion between fasta, phylip and nexus



```
>species1
ACGTGTCTG
>species2
ACGTAGTGT
>species3
ACGTTGTTT
```

```
>species1
GGTATTGCG
>species2
GCTATAGCG
```

```
>species1
AATGCGTTT
>species3
AATGAATTT
```

```
>species1
ACGTGTCTGGGTATTGCGAATGCGTTT
>species2
ACGTAGTGTGCTATAGCGnnnnnnnnn
>species3
ACGTTGTTTnnnnnnnnnnAATGAATTT
```

perl FASconCAT-G_v1.0.pl -p -p -n -n -s

# Alignment improvement

- Gblocks:

- http://molevol.cmima.csic.es/castresana/Gblocks_server.html



*Options to play with:*

# Alignment improvement

- Gblocks



FIGURE 2. Fragment of a simulated alignment (a) and the realignment of the same sequences (after gap removal) by ClustalW (b), Mafft (c), and Probcons (d). The simulation corresponds to an asymmetric tree with divergence ×1. The blocks below each alignment represent the fragments selected by Gblocks with relaxed conditions (grey blocks) and with stringent conditions (white blocks). Positions of the alignments where more than 50% of the sequences are identical are shown with black boxes.

# Alignment improvement

FIGURE 6. Mafft alignment strategies that give rise to the statistically best topologies. When two or more strategies do not show statistical differences in Robinson-Foulds distances, all equivalent strategies are represented. The complete alignment is represented by a black block, and the relaxed and stringent Gblocks strategies by grey and white blocks, respectively.

# Alignment improvement

- SOAP - http://ueg.ulb.ac.be/SOAP/

# Alignment improvement

- Trimal - http://trimal.cgenomics.org
  - automated removal of spurious sequences or poorly aligned regions

```
trimal -in example1 -out output6 -htmlout output6.html -gappyout
```

```
Selected Residue / Sequence
Deleted Residue / Sequence
                    10        20        30        40        50        60
            =========+=========+=========+=========+=========+=========+
   Sp8      ----------GLGKV---IVY-GIVLGTKS-DQFSNWVVWL-----FPWNGLQIHMMGII
   Sp17     --------FAYTAPD---LLLIGFLLKTVA-T-FG--DTWF-----QLWQGLDLNKMPVF
   Sp10     ----------DPAVL----FV--IMLGTIT-K-FS--SEWF-----FAWLGLEINMMVII
   Sp26     AAAAAAAA----ALL---TYL-GLFLGTDY-----EN---FAAAAANAWLGLEINMMAQI
   Sp33     ----------PTIL---NIA-GLHMETDI-N-FS--LAWF-----QAWGGLEINKQAIL
   Sp6      ----------ASGAI---LTL-GIYLFTLC-AVIS--VSWY-----LAWLGLEINMMAII
```
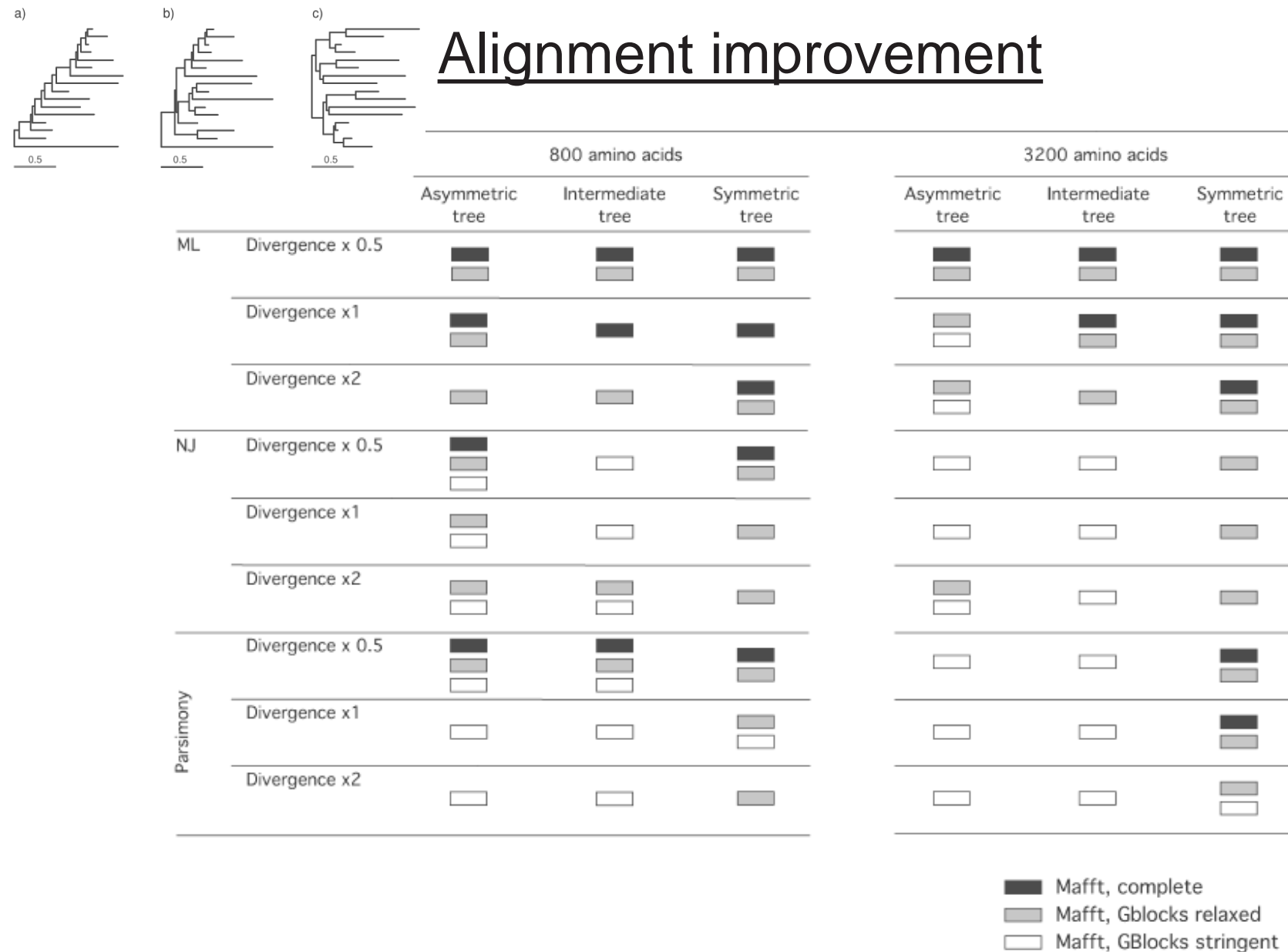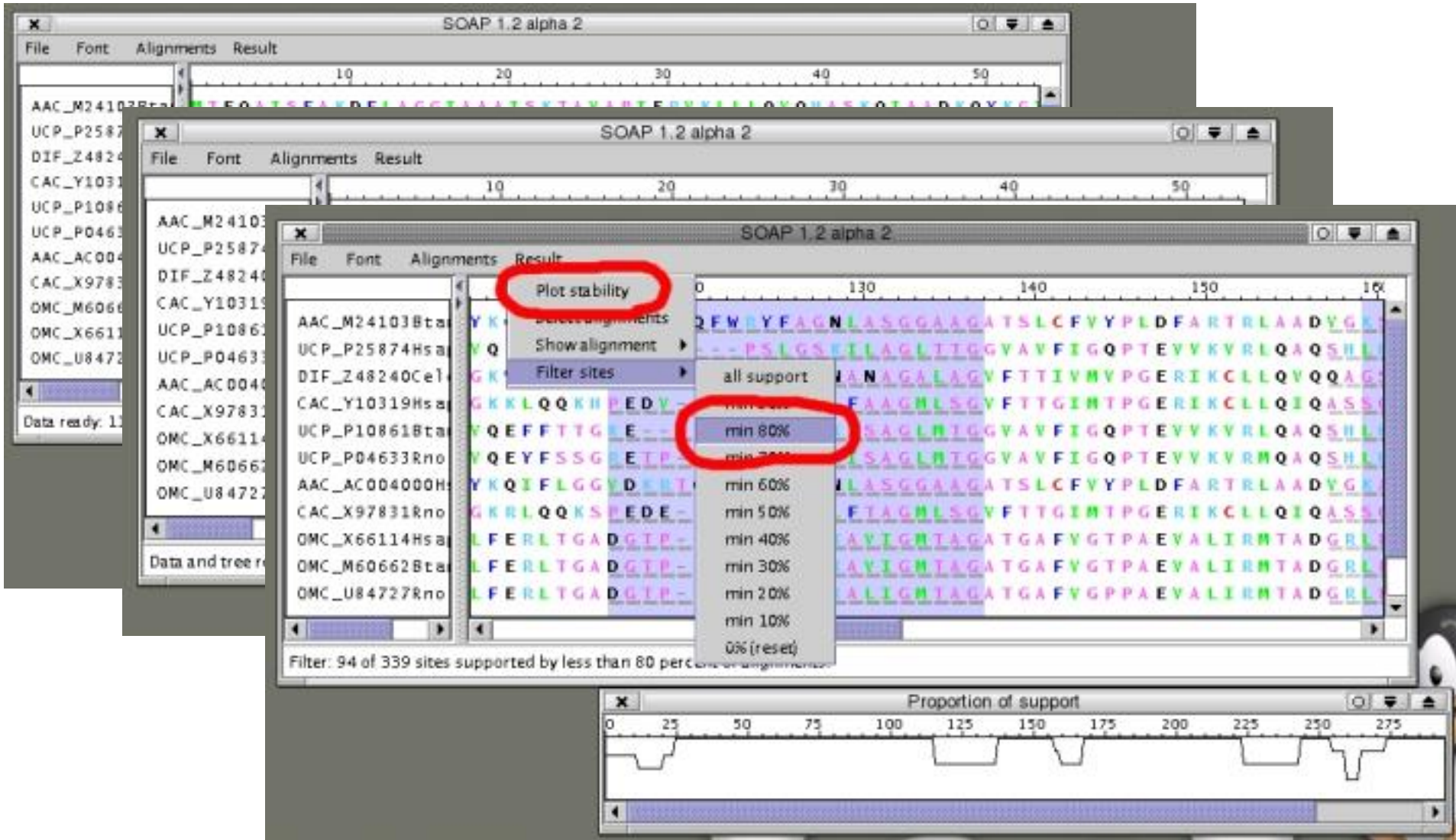
```
trimal -in example1 -out output7 -htmlout output7.html -strict
```

```
Selected Residue / Sequence
Deleted Residue / Sequence
                    10        20        30        40        50        60
            =========+=========+=========+=========+=========+=========+
   Sp8      ----------GLGKV---IVY-GIVLGTKS-DQFSNWVVWL-----FPWNGLQIHMMGII
   Sp17     --------FAYTAPD---LLLIGFLLKTVA-T-FG--DTWF-----QLWQGLDLNKMPVF
   Sp10     ----------DPAVL----FV--IMLGTIT-K-FS--SEWF-----FAWLGLEINMMVII
   Sp26     AAAAAAAA----ALL---TYL-GLFLGTDY-----EN---FAAAAANAWLGLEINMMAQI
   Sp33     ----------PTIL---NIA-GLHMETDI-N-FS--LAWF-----QAWGGLEINKQAIL
   Sp6      ----------ASGAI---LTL-GIYLFTLC-AVIS--VSWY-----LAWLGLEINMMAII
```
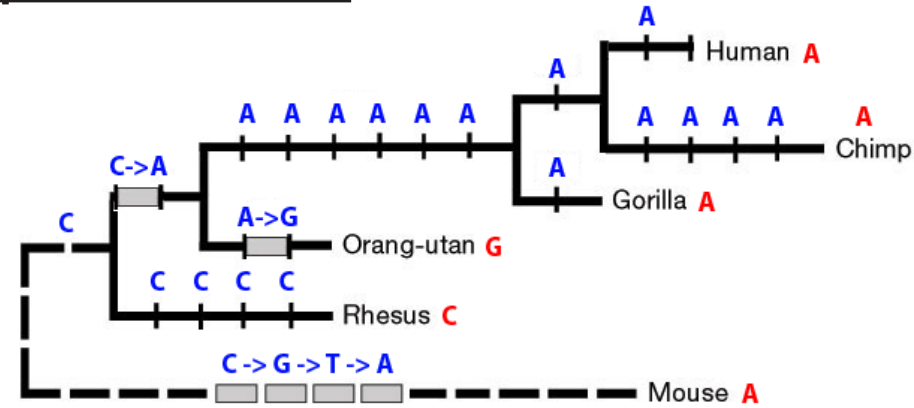
# Alignment improvement



Substitution saturation

- Some positions in alignment were changed multiple times
- Due to only 4 nucleotide types, the noise is stochastically incresing by time
- Saturated positions can represent the majority of variability in data
- A big problem for MP analyses!

1) Saturation curves:
- Comparison of sequence distances accounted under the simple and more complex substitution models

# Alignment improvement

2) Site stripping
- Deletion of saturated alignment positions

# Model selection

- *p-distance*: the simplest distance measure: the proportion of sites that differ between two sequences (saturation!).
- *substitution models* – models of how frequently different mutations occur in the DNA to precisely estimate the evolutionary distance between organisms
- Different combinations of base frequencies and nucleotide substitutions



Nature Reviews | Genetics

# Model selection

- Jukes-Cantor (JC69): equal base frequencies, all substitutions equally likely
  (nst=1, rate classification: aaaaaa)

|   | A | T | C | G |
|---|---|---|---|---|
| A | - | $\alpha$ | $\alpha$ | $\alpha$ |
| T | $\alpha$ | - | $\alpha$ | $\alpha$ |
| C | $\alpha$ | $\alpha$ | - | $\alpha$ |
| G | $\alpha$ | $\alpha$ | $\alpha$ | - |

- Felsenstein (F81): variable base frequencies, all substitutions equally likely
  (nst=1, rate classification: aaaaaa)

|   | A | T | C | G |
|---|---|---|---|---|
| A | - | $\alpha \prod_T$ | $\alpha \prod_C$ | $\alpha \prod_G$ |
| T | $\alpha \prod_A$ | - | $\alpha \prod_C$ | $\alpha \prod_G$ |
| C | $\alpha \prod_A$ | $\alpha \prod_T$ | - | $\alpha \prod_G$ |
| G | $\alpha \prod_A$ | $\alpha \prod_T$ | $\alpha \prod_C$ | - |

- Kimura (K80): equal base frequencies, one transition rate and one transversion rate
  (nst=2, rate classification: abaaba)

|   | A | T | C | G |
|---|---|---|---|---|
| A | - | $\beta$ | $\beta$ | $\alpha$ |
| T | $\beta$ | - | $\alpha$ | $\beta$ |
| C | $\beta$ | $\alpha$ | - | $\beta$ |
| G | $\alpha$ | $\beta$ | $\beta$ | - |

- Hasegawa-Kishino-Yano (HKY): variable base frequencies, one transition rate and one transversion rate
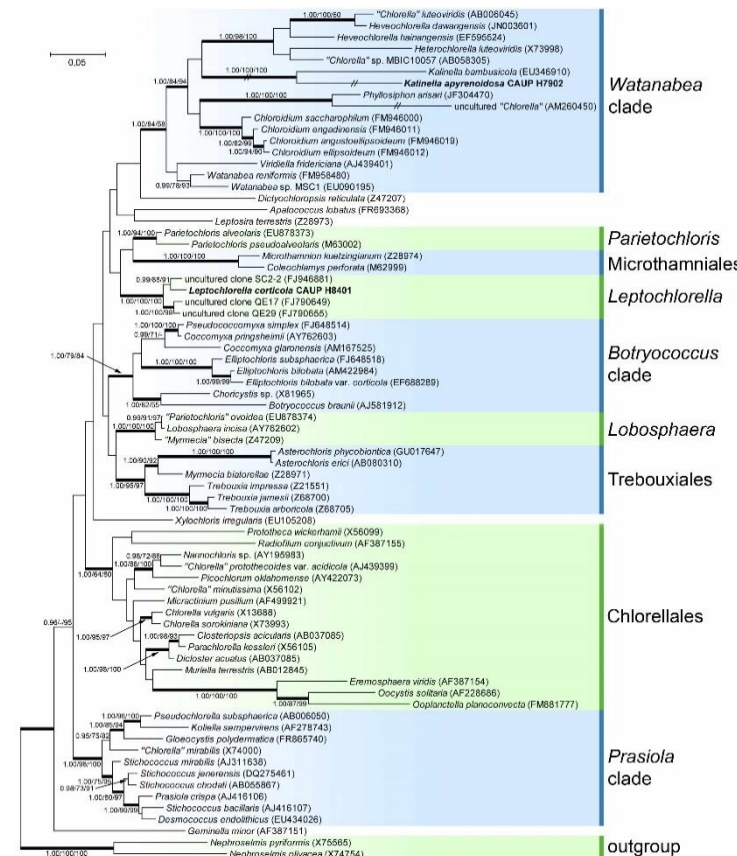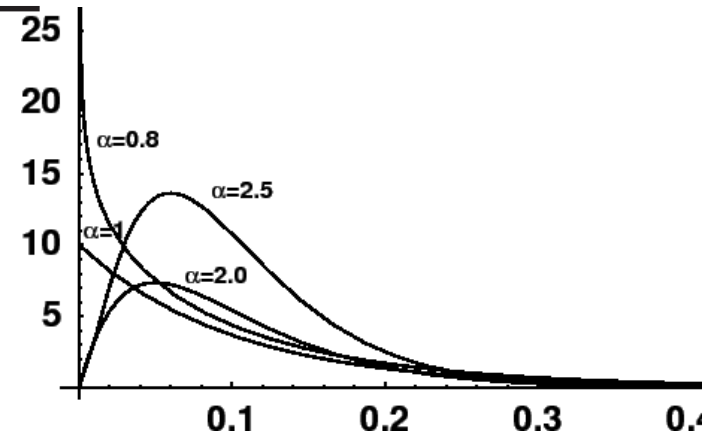  (nst=2, rate classification: abbbba)

|   | A | T | C | G |
|---|---|---|---|---|
| A | - | $\beta \prod_T$ | $\beta \prod_C$ | $\alpha \prod_G$ |
| T | $\beta \prod_A$ | - | $\beta \prod_C$ | $\beta \prod_G$ |
| C | $\beta \prod_A$ | $\beta \prod_T$ | - | $\beta \prod_G$ |
| G | $\beta \prod_A$ | $\beta \prod_T$ | $\beta \prod_C$ | - |

- General time reverdible (GTRvariable base frequencies, symmetrical substitution matrix
  (nst=6, rate classification: abcdef)

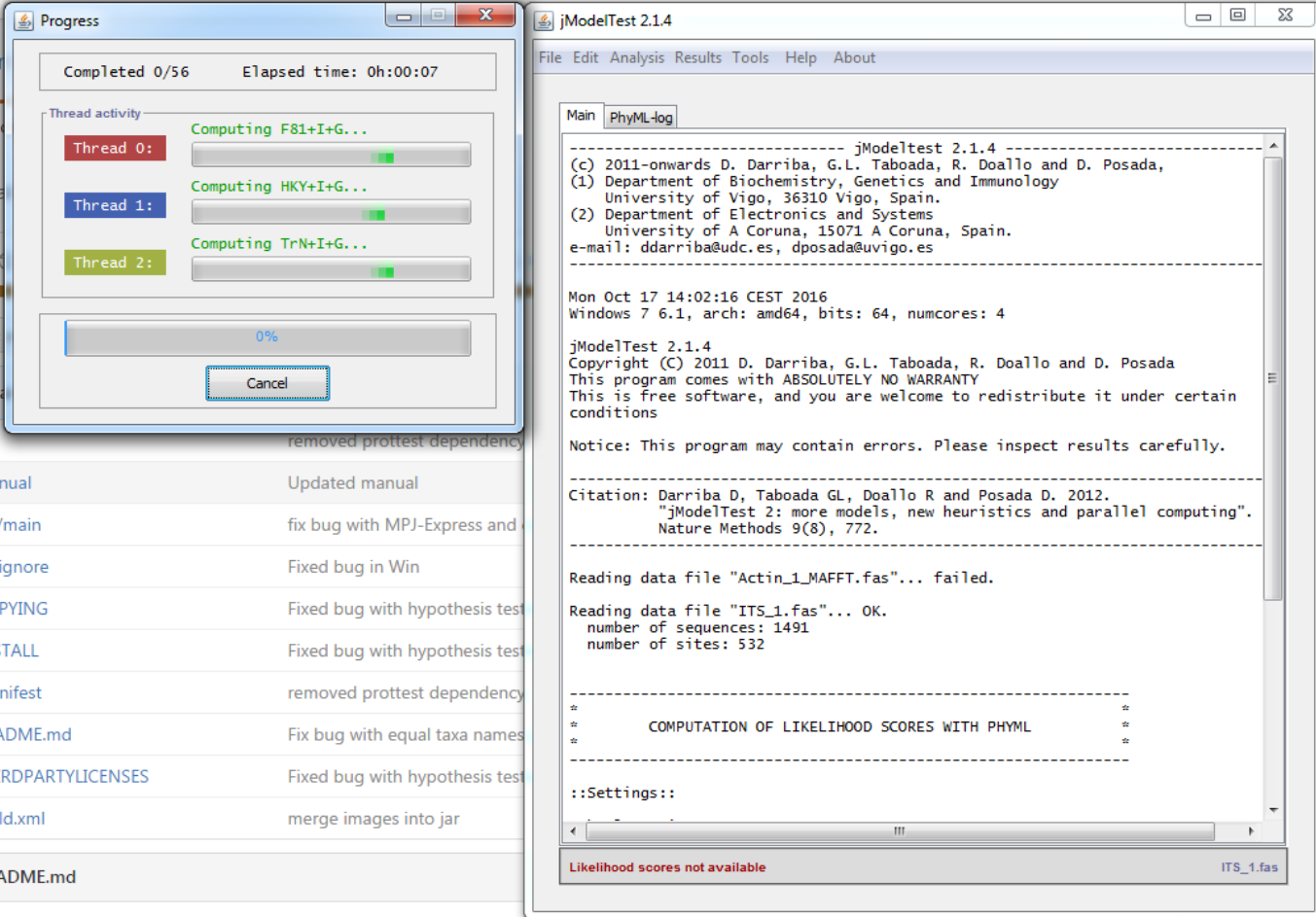|   | A | T | C | G |
|---|---|---|---|---|
| A | - | $\alpha \prod_T$ | $\beta \prod_C$ | $\gamma \prod_G$ |
| T | $\alpha \prod_A$ | - | $\delta \prod_C$ | $\varepsilon \prod_G$ |
| C | $\beta \prod_A$ | $\delta \prod_T$ | - | $\zeta \prod_G$ |
| G | $\alpha \prod_A$ | $\varepsilon \prod_T$ | $\zeta \prod_C$ | - |

# Model selection

- Gamma distribution (Γ): models a variability in substitution rates on different alignment positions. Usually, a model is simplified to 4 α categories

- Proportion of invariable sites (I): existence of a majority of invariable sites may negatively affect the estimated genetic distances. I model is particularly important in joint occurrence of very short and long branches

- Covarion (cov): models a variability in nucleotide substitution rates across the phylogenetic tree

# Model selection

- jModelTest: https://github.com/ddarriba/jmodeltest2

# Best partitioning

- Partition Finder (Lanfear et al.) – selecting best-fit partitioning schemes and models of evolution
  - for all partitions simultaneously
  - merge partitions with same model into one
  - requires Python 2.7
  - alignment in phylip format
  - configuration file

```
# ALIGNMENT FILE #
alignment = test.phy;

# BRANCHLENGTHS #
branchlengths = linked;

# MODELS OF EVOLUTION #
models = all;
model_selection = aicc;

# DATA BLOCKS #
[data_blocks]
Gene1_pos1 = 1-789\3;
Gene1_pos2 = 2-789\3;
Gene1_pos3 = 3-789\3;

# SCHEMES #
[schemes]
search = greedy;
```
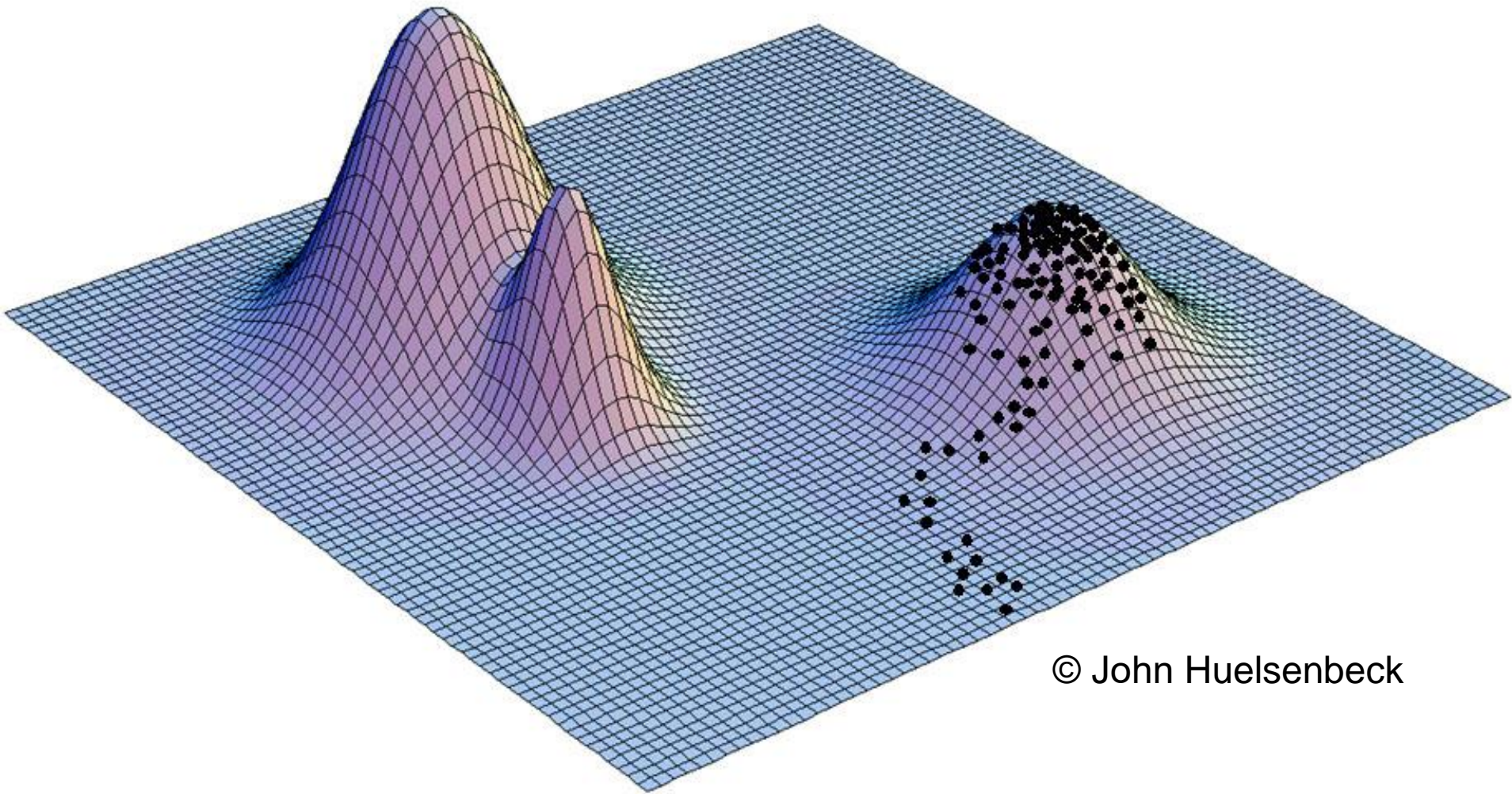
# Bayesian inference - MrBayes



© John Huelsenbeck

# Bayesian inference - MrBayes

```
begin mrbayes;


charset rbcL1      = 3-1179 \ 3;
charset rbcL2      = 1-1180 \ 3;
charset rbcL3      = 2-1181 \ 3;
charset ITS12      = 1182-1525 1683-1902;
charset RNA        = 1526-1682;

partition marker = 5:rbcL1,rbcL2,rbcL3,ITS12,RNA;

set partition = marker;

lset applyto=(1)   nst=6  rates=gamma;

lset applyto=(2,5) nst=1  rates=equal;

lset applyto=(3,4) nst=2  rates=gamma;

prset applyto=(1,3,4)  statefreqpr=dirichlet(1,1,1,1);

prset applyto=(2,5)  statefreqpr=fixed(equal);

prset applyto=(all) ratepr=variable;

unlink statefreq=(all) revmat=(all) tratio=(all) shape=(all)
pinvar=(all);

mcmc ngen=5000000 samplefreq=100 nchains=4;
end;
```

```
begin mrbayes;

charset ITS1          = 1-152;
charset ITS2          = 319-519;
charset RNA           = 153-318;
charset intron1       = 520-727;
charset exon          = 728-850;
charset intron2       = 851-1142;

partition vse = 6:ITS1,ITS2,RNA,intron1,exon,intron2;

set partition = vse;

lset applyto=(all)   nst=mixed  rates=gamma;

prset applyto=(all)  statefreqpr=dirichlet(1,1,1,1);

prset applyto=(all) ratepr=variable;

unlink statefreq=(all) revmat=(all) tratio=(all) shape=(all) pinvar=(all);

mcmc ngen=5000000 samplefreq=100 nchains=4;
end;
```

# Parsimony - PAUP

```
[maximum parsimony block]
begin paup;
log start=yes file=MP.log replace=yes;
set autoclose=yes warnreset=no increase=auto;
set criterion=parsimony;
hsearch;
savetrees brlens=yes file=treeMP.tre replace=yes;
contree /strict=yes majrule=yes treefile=contree.tre replace=yes;
log stop;
end;
```

```
[bootstrap MP block]
log start=yes file=MPboot.log replace=yes;
bootstrap search=heuristic nreps=100 conlevel=50;
savetrees from=1 to=1 file=MPboot.tre savebootp=nodelabels maxdecimals=1 replace=yes;
log stop;
end;
```

# ML - RAxML

## Basic command line parameters
-m          substitution model
-p          random seed
-t          starting tree (if not specified, parsimony tree is
                  generated using randomized stepwise addition)
-s          input file (phylip or fasta)
-#          number of replicates
-n          suffix for resulting files

## Bootstrapping
a) find best ML tree (best-scoring tree)
**raxmlHPC -m GTRGAMMA -p 12345 -# 20 -s dna.phy -n bestML**
- generate 20 trees, best saved to RAxML_bestTree.bestML
b) compute bootstrap replicates
**raxmlHPC -m GTRGAMMA -p 12345 -b 12345 -# 100 -s dna.phy -n boot**
- generate 100 bootstrap matrices
- trees generated to RAxML_bootstrap.boot
c) bootstrapu values mapped onto best ML tree
**raxmlHPC -m GTRGAMMA -p 12345 -f b -t RAxML_bestTree.bestML -z**
**RAxML_bootstrap.boot -n finalboot**
- two files generated: RAxML_bipartitions.finalboot (bootstrap values as nodes) and
RAxML_bipartitionsBranchLabels.finalboot (bootstrap values above branches)

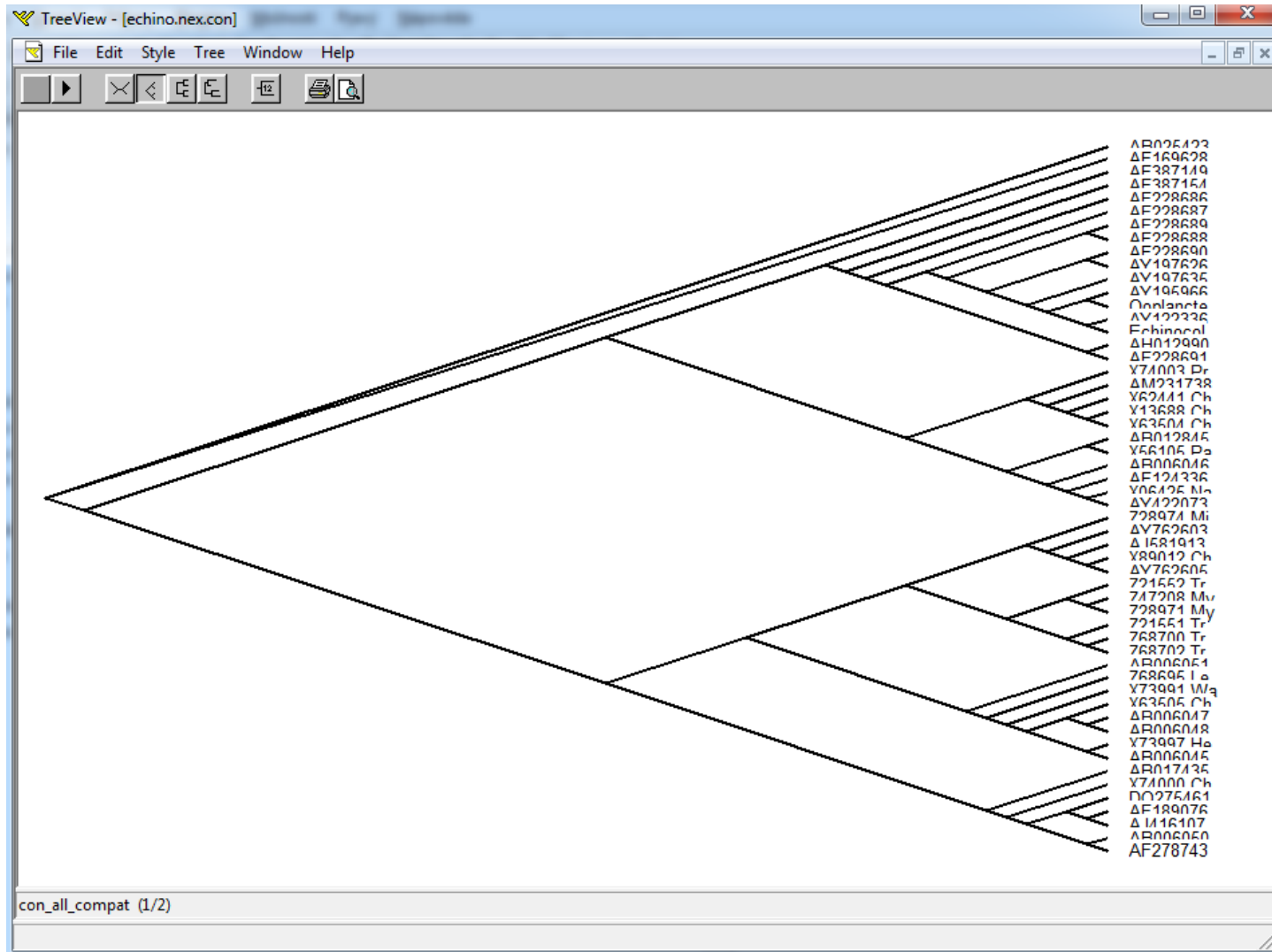## Rapid bootstrap
- much faster than standard bootstrap
- complete analysis (ML search+ bootstrapping) in one step
**raxmlHPC -f a -m GTRGAMMA -p 12345 -x 12345 -# 100 -s dna.phy -n rbs**
- RAxML_bipartitions.rbs – best ML tree with mapped bootstrap replicates
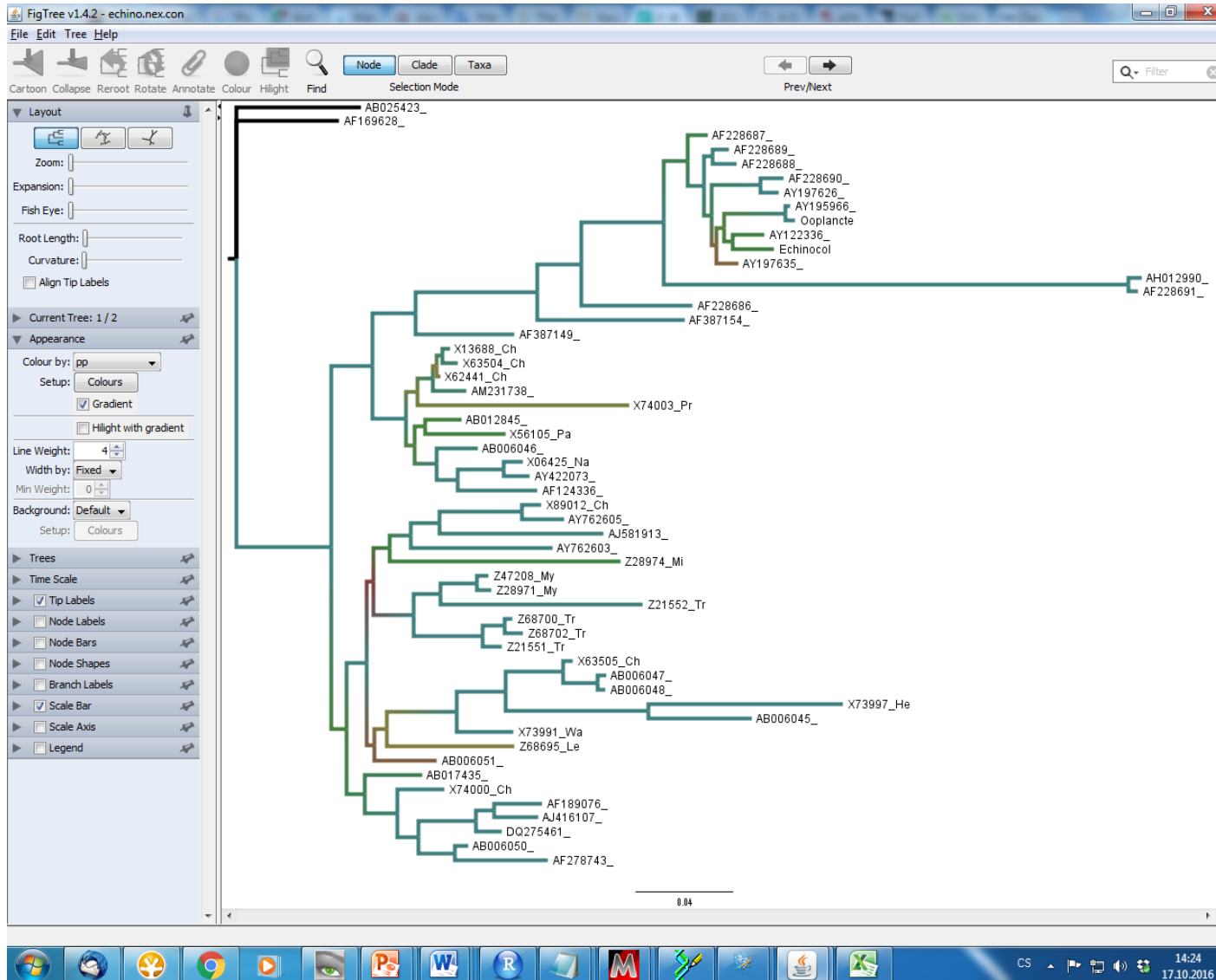
# Tree Manipulation

- TreeView - http://en.bio-soft.net/tree/TreeView.html/

# Tree Manipulation

- FigTree - http://tree.bio.ed.ac.uk/software/figtree/
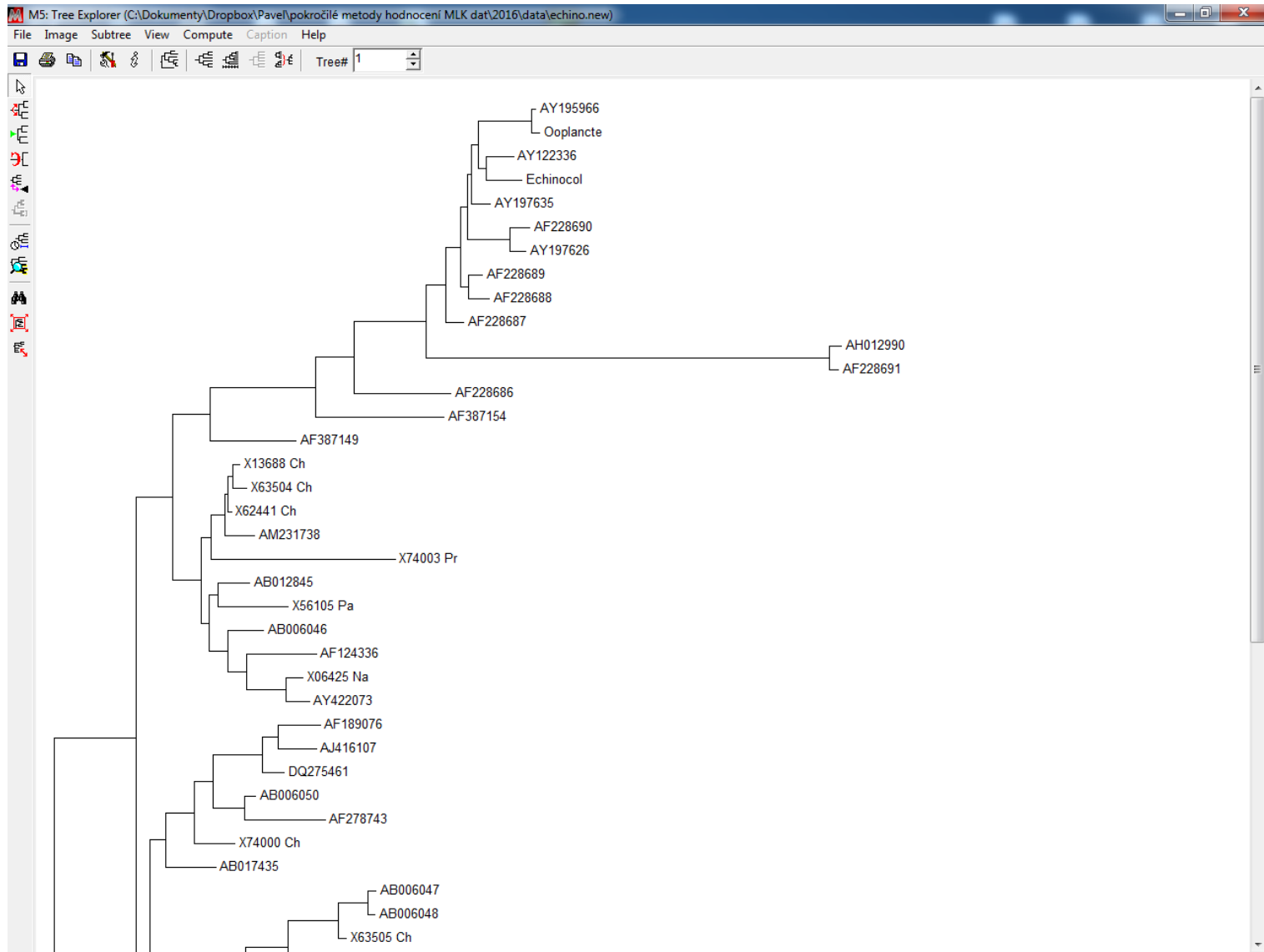
# Tree Manipulation

# Excercises:

Compare phylogenetic trees using:

- alignments with different amount of missing data

- alignments obtained by different MAFFT options

- original and minimized alignments after cleaning poorly aligned regions

- indel coding

- simple and complex evolutionary models

- partitioned and un-partitioned concatenated datasets

- different inferences (BI, ML, parsimony)