

Spring 2024

BIOS 477/877

Bioinformatics and Molecular Evolution

Lecture 24

BIOS477/877 L24 - 1

1

TODAY'S TOPICS

- Phylogenetic analysis/visualization software
 - Transfer bootstrap expectation
 - Software and visualization
- Protein structure databases
 - Protein structures
 - Secondary structure assignment
 - Protein structure databases (PDB)

BIOS477/877 L24 - 2

2

Renewing Felsenstein's phylogenetic bootstrap in the era of big data

F. Lemoine^{1,2}, J.-R. Domelevo Entfellner^{3,4}, E. Wilkinson⁵, D. Corcoran¹, M. Davila Felipe¹, T. De Oliveira^{5,6} & O. Gascuel^{1,2,*}

432 | NATURE | VOL 556 | 26 APRIL 2018

Transfer distance (or R distance):
The minimum number of elements to be transferred (or removed) to transform on partition into the other

(S1, S2), (S3, S4, S5) $\xleftrightarrow{\text{Distance}=0}$ (S1, S2), (S3, S4, S5)

(S1, S2, S3), (S4, S5) $\xleftrightarrow{\text{Distance}=1}$ (S1, S2), (S3, S4, S5)

Reference (original)

tree T

S1, S2, S3, S4, S5

Distance $\Delta(b, b^*) = 1$

Bootstrap tree T*

S1, S2, S3, S4, S5

Distance $\Delta(b, b^*) = 2$

Transfer Bootstrap Expectation (TBE) for branch b:

$$TBE(b) = 1 - \frac{\phi(b, T^*)}{p-1}$$

Average from all BS trees

p: number of taxa from the smaller of the two clusters

BS: count only identical partitions (no need to be identical)

$\phi(b, T^*) = \min_{b^*} \Delta(b, b^*)$

$\Delta(b, b^*)$: Distance between branch b in tree T and branch b* in BS tree T*

- Compare the two partitions defined by branch b in tree T and all branches in BS tree T*

- Find the b* with the minimum $\Delta(b, b^*)$

BIOS477/877 L24 - 3

3

Transfer bootstrap for big phylogenies

Renewing Felsenstein's phylogenetic bootstrap in the era of big data

Robustness of Felsenstein's Versus Transfer Bootstrap Supports With Respect to Taxon Sampling

PAUL ZARAKIAN^{1*}, FREDERIC LEMOINE^{2*} AND OLIVIER GASCUEL^{1,3*}

Standard bootstrap proportion (SBP) > 70%

Standard bootstrap proportion (SBP) > 70%

Transfer bootstrap expectation (TBE) > 70%

Transfer bootstrap expectation (TBE) > 70%

- Deep branches in large phylogenies are often not supported by SBP

- TBE supports are higher without inducing falsely supported branches

BIOS477/877 L24 - 4

4

PHYLIP format for input alignment

<https://phylipweb.cit.ub.edu/phylip/>

Phylip tutorial on Canvas

Sequence number **Alignment length** (Interleaved alignment format)

```

4      125
1c2r   -----GDAAK  GEKEFN-KCK  TCHSIIAPDG  TEIVRGAKTG  PNLYGVVGRT
1ycc   TEFKAGSARK  GATLFKTRCL  QCHTVEKGG-  -----PHKVG  PNLHGIFGRH
3c2c   -----EGDAAA  GEKVS-KCL  ACHTFDQGG-  -----ANKVG  PNLFGVFENT
1etp   -----AGDAEA  GQGKVA-VCG  ACHGVDGNS-  -----PA  FN-----FPKL
AGTYPEFKYK  DSIVALLGASG  FAWTEEDIAT  YVKDPGAFKL  ERLDRKAKT
SQGAEGYSYT  DAN-----IKRN  VLWDENNMSE  YLTNPKRYIP  G-----T
AAHKDNYSYS  ESYTEMKARG  LTWTEANLAA  YVKNPKAFVL  EKSGDPKAS
AGQGERYLLK  QLQ-DIRAGS  TPGAPEGVGR  KVLLEMTGMLD  P-----
GMAFK-LAKG  GE--DVAAYL  ASVVK
KMAFGGLKKE  KDRNDLITYL  KKACE
KMTFK-LTKD  DEIENVIAYL  KTLK-
--LS-----  DQLEDIAAYF  SSQKG

```

Sequence names should be ten characters in length. Add spaces AFTER the sequence name if names are too short. Do not include space or "[O]:" in the sequence names.

BIOS477/877 L24 - 5

5

Strict vs. relaxed PHYLIP alignment format

Strict Phylip format

```

4      20
1c2r   -----GDAAK  GEKEFN-KCK
1ycc   TEFKAGSARK  GATLFKTRCL
3c2c   -----EGDAAA  GEKVS-KCL
1etp   -----AGDAEA  GQGKVA-VCG

```

Only up to 10 characters are used as the sequence names.

No space is needed between names and sequences.

In the strict format: Names > 10 letters will be truncated!

Relaxed phylip format

```

4      20
1c2rABCDEFghi  -----GDAAK  GEKEFN-KCK
1yccXXXXXXXXTEFKAGSARK  GATLFKTRCL
3c2cYYYYYYY-----EGDAAA  GEKVS-KCL
1etpZZZZZZ-----AGDAEA  GQGKVA-VCG

```

Relaxed format allows to have longer sequence names.

At least one space is needed between names and sequences.

Used in PhyML, RAxML, etc.

Relaxed format is not supported by PHYLIP (sequence data start at 11th character)

BIOS477/877 L24 - 6

6

Maximum likelihood phylogeny: IQ-TREE

<http://www.iqtree.org>; <http://iqtree.cibiv.univie.ac.at> (other websites available)

IQ-TREE web server: fast and accurate phylogenetic trees under maximum likelihood

Server load: 6% Trifinopoulos J, Nguyen LT, von Haeseler A, Minh BQ (2016) *Mol. Acids Res.* 44 (W1): W232-W235.

Tree Inference Model Selection Analysis Results

Substitution Model Options

Substitution model: **Auto** (← Substitution model can be chosen automatically)

Rate heterogeneity: Gamma (+G) Invar. sites (+I) L₂

#Rate categories: 4 5 6

State frequency: Empirical (from data) AA model (from matrix) M₀-optimized Codon F₄L4 Codon F₃L4

Ascertainment bias correction: Yes (+ASC) No

Branch Support Analysis

Bootstrap analysis: None Ultrafast Standard

Number of bootstrap alignments: 1000 (← Default branch support: Ultrafast bootstrap + SH-aLRT branch test)

Create .ufboot file: Yes (write bootstrap trees to .ufboot file)

Minimum correlation coefficient: 0.99

Single branch tests: SH-aLRT branch test: No Yes #replicates: 1000

Approximate Bayes test: Yes No

BIOS477/877 L24 - 13

13

Maximum likelihood phylogeny: IQ-TREE

<http://www.iqtree.org>; <http://iqtree.cibiv.univie.ac.at> (other websites available)

IQ-TREE web server: fast and accurate phylogenetic trees under maximum likelihood

Server load: 6% Trifinopoulos J, Nguyen LT, von Haeseler A, Minh BQ (2016) *Mol. Acids Res.* 44 (W1): W232-W235.

Tree Inference Model Selection Analysis Results

Input Data

Alignment file: Browse... Show example >

Use example alignment: Yes

Sequence type: Auto-detect DNA Protein DNA->AA Binary Morph

Partition file: This field is optional. Browse...

#nexus
begin sets;
charset part1 = aln1.phy: 1-100; 3 201-300;
charset part2 = aln1.phy: 101-200;
charset part3 = aln2.phy: *;
charpartition mine = HKY:part1, GTR+G:part2, WAG+I+G:part3;
end;

→ Each partition can have a different substitution model (e.g., different genes or even DNA and protein alignments can be combined)

DNA substitution model AA substitution model

BIOS477/877 L24 - 14

14

Phylogenetic method programs/websites

- **PhyML** <http://www.atgc-montpellier.fr/phyml/> (3.0)
 - Smart Model Selection is included (no need to choose a substitution model).
 - Standalone or web interface (aLRT SH-like, standard bootstrap, and transfer bootstrap)
- **RAxML** <http://sco.h-its.org/exelixis/web/software/raxml/index.html>
 - Randomized Axelerated Maximum Likelihood (includes rapid bootstrap)
- **RAxML-NG** <https://github.com/ambrojo/raxml-ng/releases>
 - Faster and more accurate than RAxML, includes transfer bootstrap
- **IQ-TREE** <http://www.iqtree.org/> (IQ-TREE 2 also available)
 - A fast and effective stochastic algorithm for estimating ML phylogenies
 - Includes model selection and ultrafast bootstrap; IQ-TREE 2 includes transfer bootstrap
- **FastTree 2** <http://microbesonline.org/fastree/>
 - Approximately-maximum likelihood phylogenetic trees (standalone only)
- **Booster** <https://booster.pasteur.fr/> (transfer bootstrap with PhyML and FastTree)
- **Los Alamos databases and tools** <https://www.hiv.lanl.gov/content/sequence/HIV/HIVTools.html> (PhyML, IQ-TREE)
- **MrBayes** Bayesian Inference of Phylogeny <https://mbiweb.cden.itsi.edu/MrBayes/index.html>
- **CIPRES** (Cyberinfrastructure for Phylogenetic Research) <http://www.phylip.org/>
 - High performance parallel codes: RAxML-NG, MrBayes, GARLI, Ball-Phy, PAUP, IQ-Tree, etc.
- **NGPhylogeny.fr** <https://ngphylogeny.fr/>
 - includes PhyML, MrBayes, FastME, TNT, etc., also BMGE, Gblocks, etc. for filtering
- **Phylogeny.fr** <http://www.phylogeny.fr>

BIOS477/877 L24 - 15

15

More phylogeny programs/websites

- **MEGA X** <http://www.megasoftware.net/>
 - UPGMA, NJ, ME, MP, ML (include visualization)
- **Phylo3.698** <https://phylo3p.github.io/phylo/>
 - Includes: FM, UPGMA, NJ, ME, MP, ML
 - Rphyloip, an R interface for Phyloip: <http://www.phytools.org/Rphyloip/>
 - Includes Phyloip, PhyML, MrBayes; also ModelTree
- **Phyloip** <http://phyloip.bioinformatics.fsu.edu>
 - includes Phyloip, PhyML, MrBayes; also ModelTree
- **SeaView** <http://douda.prabi.fr/software/seaview>
 - MSA (Muscle, Clustal Ω) and phylogeny (NJ, BioNJ, MP, PhyML, including transfer bootstrap)
- **MAFFT** <https://mafft.cbrc.jp/alignment/server/>
 - Includes NJ (distance: JTT, WAG, etc.)
- **DAMBE** <http://dambe.bio.uottawa.ca/DAMBE/dambe.aspx>
 - UPGMA, NJ, FastME, FM, MP, ML (many substitution models; visualization)
- **CRAN** (R projects for phylogenetics) <https://cran.r-project.org/web/views/Phylogenetics.html>
- **Ape** (R package for Analysis of Phylogenetics and Evolution) <https://cran.r-project.org/web/packages/ape/index.html>

BIOS477/877 L24 - 16

16

Newick tree format with node labels

branch lengths

((1yce:0.87136,3c2c:0.29363):0.10012,1etp:1.57312,1c2r:0.62829); **No node label**

((1yce:0.87136,3c2c:0.29363)X:0.10012,1etp:1.57312,1c2r:0.62829); **With node label**

((1yce,3c2c)X,1etp,1c2r); **Only with node label, no branch length**

((1yce:0.87136,3c2c:0.29363)50:0.10012,1etp:1.57312,1c2r:0.62829); **Using node labels to show bootstrap values**

BIOS477/877 L24 - 17

17

MEGA Molecular Evolutionary Genetics Analysis

(<https://www.megasoftware.net>)

Original ML tree with bootstrap values Consensus tree with bootstrap values

(((HBB_HUMAN:0.04045989,HBB_HORSE:0.1532071)0)1.0000:0.6439084,(HBA_HUMAN:0.02935227,HBA_HORSE:0.10205070)0.7800:0.19434686)0.9200:0.388280 04.GLBS_PETMA:0.65024328)0.8000:0.31847481,MYG_PHYCA1:0.4418589,LGB2_LUPLU:1.59298882);

MEGA includes bootstrap values as node labels. No branch lengths are included in the consensus tree.

BIOS477/877 L24 - 18

18

Visualization of phylogeny

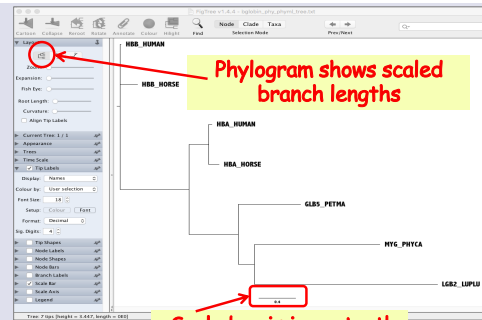
- **FigTree** (Macintosh, Windows, Linux/Unix) <http://tree.bio.ed.ac.uk/software/figtree/>
- **iTOL** (Web service) <http://itol.embl.de/>
- **TreeViewer** (Macintosh, Windows, Linux/Unix) <https://treeviewer.org/>
- **IcyTree**(Web service) <https://icytree.org/>
- **Phylo.io** (Web service) <http://phylo.io/index.html>
- **PRESTO** (Web service) <http://www.atgc-montpellier.fr/presto/>
- **EvoView v2** (Web service) <http://www.evolgenius.info/evolview/>
- **Phylogenetic tree (newick) viewer** (Web service) <http://cctoolkit.org/treeview/>

BIOS477/877 L24 - 19

19

FigTree (<http://tree.bio.ed.ac.uk/software/figtree/>)

((HBB_HUMAN:0.01371329,HBB_HORSE:0.18948102,((HBA_HUMAN:0.03403350,HBA_HORSE:0.11224574)0.847000:0.21449710,(GLBS_PETMA:0.67390367,(MTG_PHYCA:1.36289368,LGB2_LUPLU:1.98798207)0.000000:0.17306486)0.963000:0.53642455)1.000000:0.75001414);

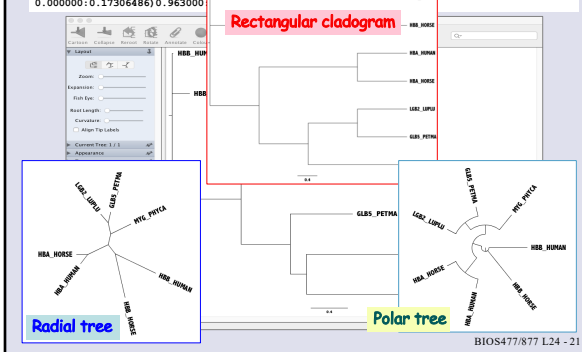


BIOS477/877 L24 - 20

20

FigTree

((HBB_HUMAN:0.01371329,HBB_HORSE:0.18948102,((HBA_HUMAN:0.03403350,HBA_HORSE:0.11224574)0.847000:0.21449710,(GLBS_PETMA:0.67390367,(MTG_PHYCA:1.36289368,LGB2_LUPLU:1.98798207)0.000000:0.17306486)0.963000:0.53642455)1.000000:0.75001414);

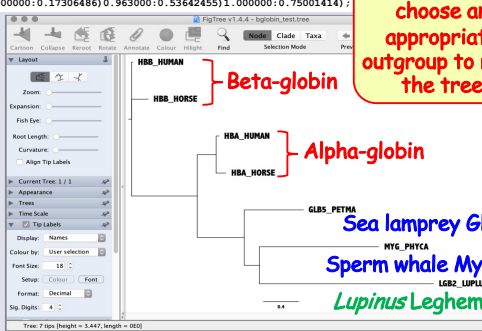


BIOS477/877 L24 - 21

21

FigTree

((HBB_HUMAN:0.01371329,HBB_HORSE:0.18948102,((HBA_HUMAN:0.03403350,HBA_HORSE:0.11224574)0.847000:0.21449710,(GLBS_PETMA:0.67390367,(MTG_PHYCA:1.36289368,LGB2_LUPLU:1.98798207)0.000000:0.17306486)0.963000:0.53642455)1.000000:0.75001414);

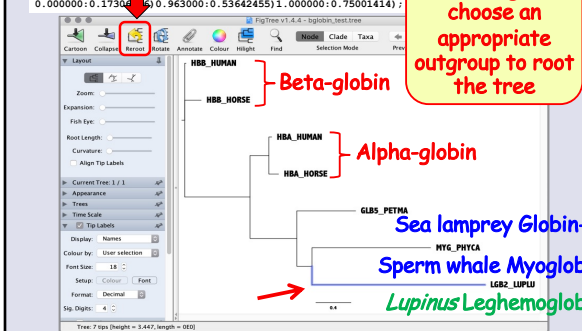


BIOS477/877 L24 - 22

22

FigTree

((HBB_HUMAN:0.01371329,HBB_HORSE:0.18948102,((HBA_HUMAN:0.03403350,HBA_HORSE:0.11224574)0.847000:0.21449710,(GLBS_PETMA:0.67390367,(MTG_PHYCA:1.36289368,LGB2_LUPLU:1.98798207)0.000000:0.17306486)0.963000:0.53642455)1.000000:0.75001414);

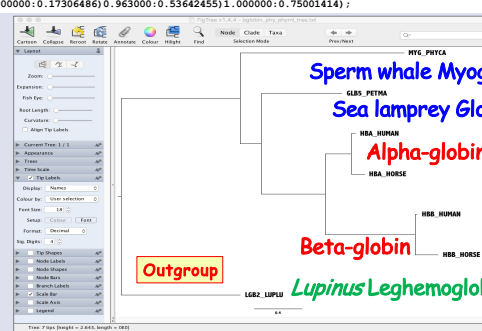


BIOS477/877 L24 - 23

23

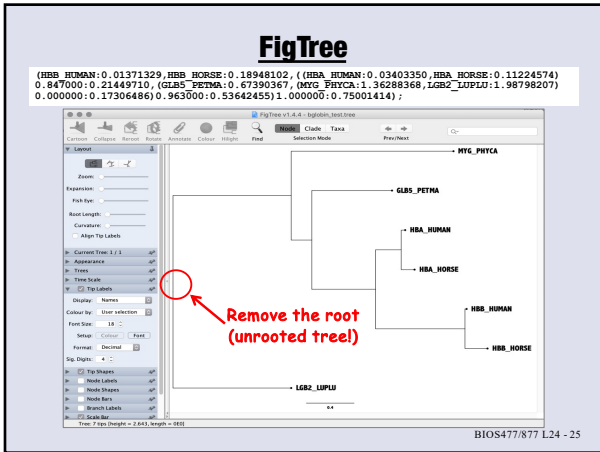
FigTree

((HBB_HUMAN:0.01371329,HBB_HORSE:0.18948102,((HBA_HUMAN:0.03403350,HBA_HORSE:0.11224574)0.847000:0.21449710,(GLBS_PETMA:0.67390367,(MTG_PHYCA:1.36289368,LGB2_LUPLU:1.98798207)0.000000:0.17306486)0.963000:0.53642455)1.000000:0.75001414);

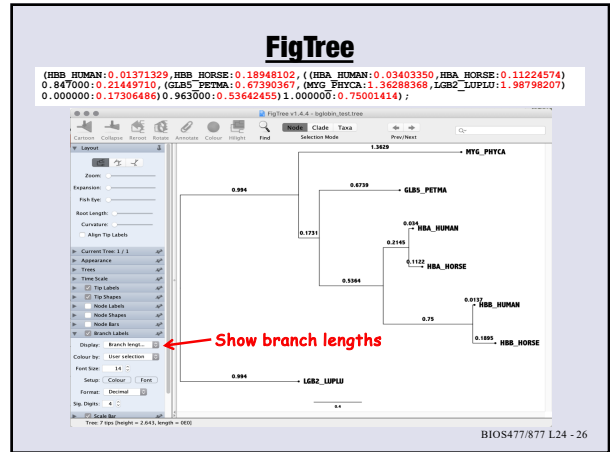


BIOS477/877 L24 - 24

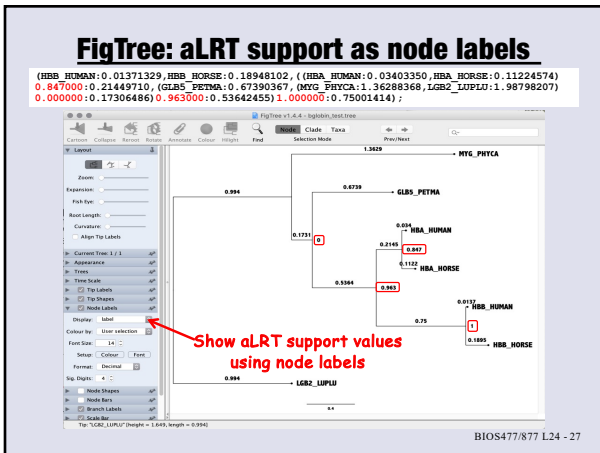
24



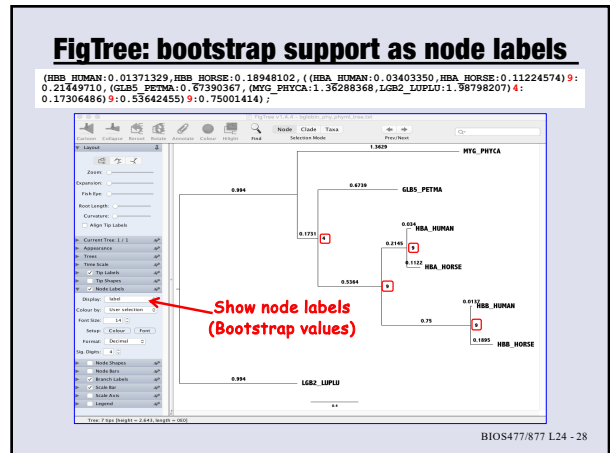
25



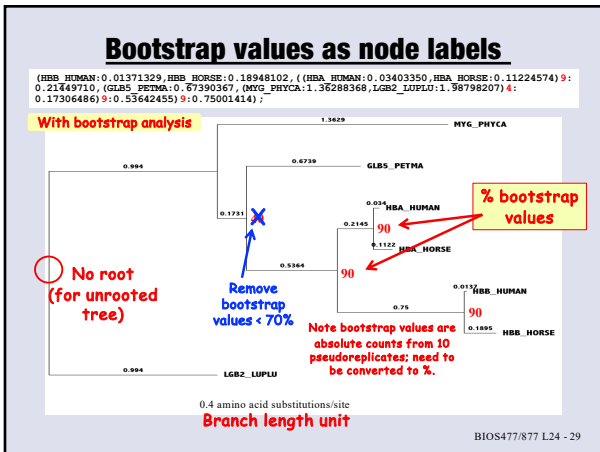
26



27



28



29

IQ-Tree and iTOL (<http://itol.embl.de/>)

IQ-Tree Newick format tree file:

(IRF4 Homo sapiens:0.0261131273,IRF4 Pteropus alecto:0.0178822599,(((IRF4 Canis lupus dingo:0.0220230929,IRF4 Bos taurus:0.0459977745)88.5/100:0.0079686246,(IRF4 Vombatus ursinus:0.0699729028,(IRF4 Ornithorhynchus anatinus:0.0336223152,(((IRF4 Gallus gallus:0.0482000904,IRF4 Alligator mississippiensis:0.0658074429)92.4/99:0.0181052721,((IRF4 Mauremys mutica:0.0000026997,IRF4 Dermochelys coriacea:0.0103795045)42.4/95:0.0058138194,IRF4 Pelodiscus sinensis:0.0125162652)94.9/100:0.0190368037)92.99:0.0244166567,(((IRF4 Varanus komodoensis:0.0302789819,((IRF4 Zootoca vivipara:0.0024867736, ...

IQ-Tree can provide two branch support values (e.g., ultrafast bootstrap/SH-aLRT)

iTOL can use both information

BIOS477/877 L24 - 30

30

IQ-Tree and iTOL (<http://itol.embl.de/>)

ITOL control panel

Define multiple threshold ranges

Operator: AND

bootstrap1 80 to 100

bootstrap2 80 to 100

Update thresholds

Disable thresholds

Branch metadata display

Multiple ranges can be defined

Size of the symbol can be changed based on the range (e.g., minimum: 80% maximum: 100%)

BIOS477/877 L24 - 31

31

IQ-Tree and iTOL (<http://itol.embl.de/>)

More examples available on the tree gallery

BIOS477/877 L24 - 32

32

Phylogenetic methods: pros and cons

- Criticisms to **distance methods**
 - Depend on distance estimation method
 - Summarizing a set of sequences by a pairwise distance matrix loses information
- Criticisms to **maximum parsimony methods**
 - “**Long branch attraction**” problem
 - If the internal branch is short relative to the terminal branches, by chance 1 and 3 may acquire the same nucleotide independently
 - ((1,3),(2,4)) may become the MP tree
 - No multiple hit correction
- Criticisms to **maximum likelihood methods**
 - Require an explicit model of evolution
 - Which model to use?

BIOS477/877 L24 - 33

33

What causes phylogenetic incongruence

Steensyk *et al.* (2023); also see Kapli *et al.* (2021)

Driver of incongruence	Factor
Incomplete lineage sorting	Biological
Horizontal gene transfer	Biological
Hybridization or introgression and recombination	Biological
Natural selection	Biological
Sampling (taxon and locus)	Analytical, stochastic error
Insufficient number of genes or divergent sites	Analytical, stochastic error
Erroneous orthologue detection	Analytical, systematic error
Model misspecification	Analytical, systematic error
Multiple sequence alignment errors	Analytical, treatment error
Excessive trimming	Analytical, treatment error
Inappropriate character recoding	Analytical, treatment error

Incongruent gene vs. species trees

Insufficient taxon sampling
Fast evolving lineage
Outgroup selection

Long-branch attraction
Inadequate model complexity

BIOS477/877 L24 - 34

34

Guidelines for phylogenetic analysis

Anisimova *et al.* (2013); Kapli *et al.* (2020)

- **Data**
 - DNA or protein? Coding or noncoding?
 - Quality of data (e.g., sequencing error, gene annotation)
 - Identification of **homologues** (orthologues only? paralogues?)
 - Carefully choose the **outgroup** sequence(s)
- Reconstruct **multiple sequence alignment**
 - Which method?
 - Data trimming/filtering (filter out unreliable alignment regions) **But with caution!**
- Selection of **substitution model**
 - Simplified vs. over-fitting (model testing can be done)
- Phylogenetic inference
 - Analyze data **combined** (with partitions) or **individually**
 - Methods (**ML, Bayesian, or NJ** for speed)
 - Branch support statistics (**bootstrap, etc.**)

BIOS477/877 L24 - 35

35

Incorporating alignment uncertainty into Felsenstein's phylogenetic bootstrap to improve its reliability

Jia-Ming Chang^{1,*}, Evan W. Floden², Javier Herrero^{1,3}, Olivier Gascuel⁴, Paolo Di Tommaso² and Cedric Notredame^{2,4,*}

Super-MSA: Concatenation of 7 alignments (ClustalW, DCA, Dialign2, Mafft, Muscle, ProbCons, and T-Coffee)

SBOOT: Bootstrap analysis is done using Super-MSA

SBOOT outperformed individual alignments for phylogenetic accuracy.

BIOS477/877 L24 - 36

36

Torsion angles

- The conformation of the backbone can be described by the **torsion angles** (or **dihedral angles** or **rotation angles**): ϕ and ψ
- The entire path of the backbone is known if ϕ and ψ are specified
- Some values of ϕ and ψ are more likely than others
- Due to steric interference between non-bonded atoms

BIOS477/877 L24 - 43

43

Secondary structure

➤ **Alpha helix:** $\phi = -60$ and $\psi = -45$

Hydrogen bonds stabilize the alpha helix

BIOS477/877 L24 - 44

44

Secondary structure

➤ **Beta sheets**

BIOS477/877 L24 - 45

45

Secondary structure

➤ **Beta turns**

- Proline and glycine occur frequently in beta turns

BIOS477/877 L24 - 46

46

Ramachandran plot

- The sterically allowed values for ϕ and ψ can be calculated (colored regions)
- α : right-handed **alpha helix**
- α_L : left-handed **alpha helix**
- $\uparrow\uparrow$: **parallel beta sheet**
- $\uparrow\downarrow$: **antiparallel beta sheet**
- C**: collagen helix

BIOS477/877 L24 - 47

47

Secondary structure conformation

	ϕ	ψ	H-bond pattern
Right handed α -helix	-57	-47	$i + 4$
π -helix	-57	-70	$i + 5$
3_{10} helix	-49	-26	$i + 3$
Parallel β -sheet	-119	113	
Antiparallel β -sheet	-139	135	

BIOS477/877 L24 - 48

48

