# Phylogenetic trees from high-throughput sequence data analysis

**Adriana Rego**

# Meta_Microbial Workshop

Metagenomic and bioinformatic insights into microbial communities

@AdrianaRego10

adrianairego@gmail.com

ciimar

U. PORTO
FACULDADE DE CIÊNCIAS
UNIVERSIDADE DO PORTO

Sponsorship

SPECO

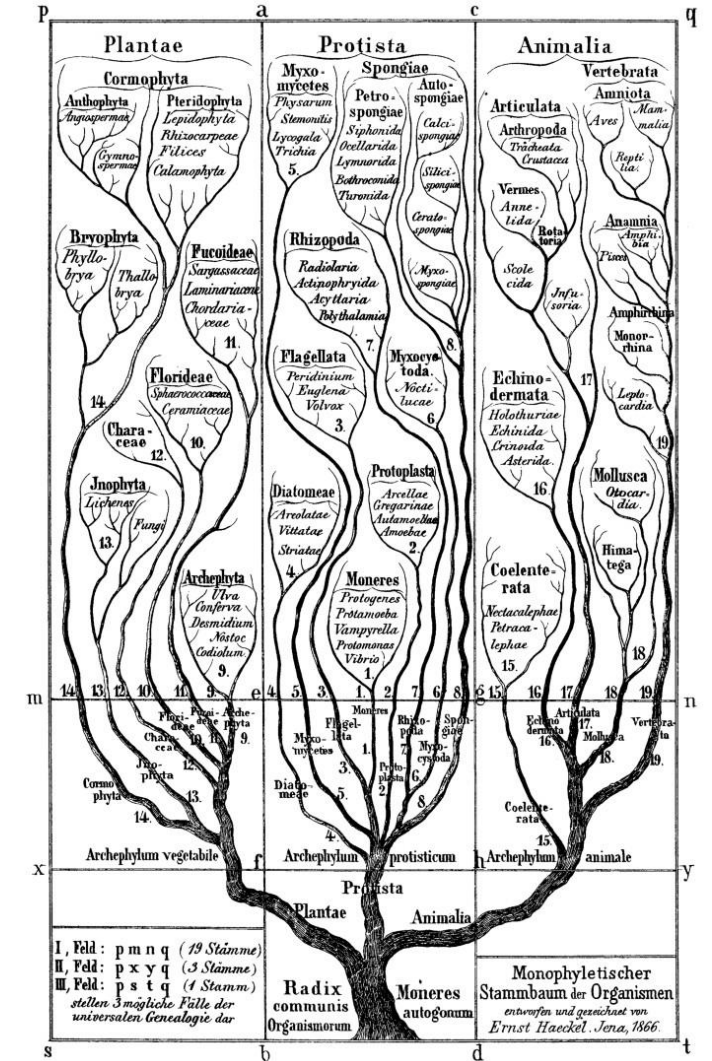BIOPORTUGAL S.A.
Químico, Farmacêutica

## Phylogenetic analysis

A **phylogenetic tree** or **evolutionary tree** is a diagrammatic representation of evolutionary relationships.

Relationships among taxa are inferred based on **homology** (inheritance from a common ancestor, commonly observed as patterns of sequence similarity).

Traditionally, phylogenies were constructed using **morphological data** only, but the advances in DNA sequencing enabled genetic information to be incorporated into phylogenetic analyses.



The earliest Tree of Life, by Ernest Haeckel (1866).

# Phylogenetic analysis

**Steps in a phylogenetic analysis**

**1** Planning

| Start with a question | Select a model (organisms or gene family) | Choose approppriate molecular markers | Collect or generate data (sequences) |

**2** Executing

| Alignment | Phylogenetic method | Phyogenetic analysis and tree reconstruction |

**3** Evaluate

| Tree evaluation | Tree interpretation and new knowledge |

# Molecular marker genes

**Examples**

16S rRNA, 23S rRNA, ITS, rpoB, gyrB, dnaK, dsrAB, amoA, amoB, mip, horA, hitA, recA, ica, frc, oxc.

A **molecular marker** is defined as any DNA sequence which shows polymorphism and can be detected using a molecular technique.

Molecular markers of bacteria should have several characteristics:

- housekeeping genes present in **all** bacterial species;
- High **polymorphism**, which make them distinguishable in different bacterial species;
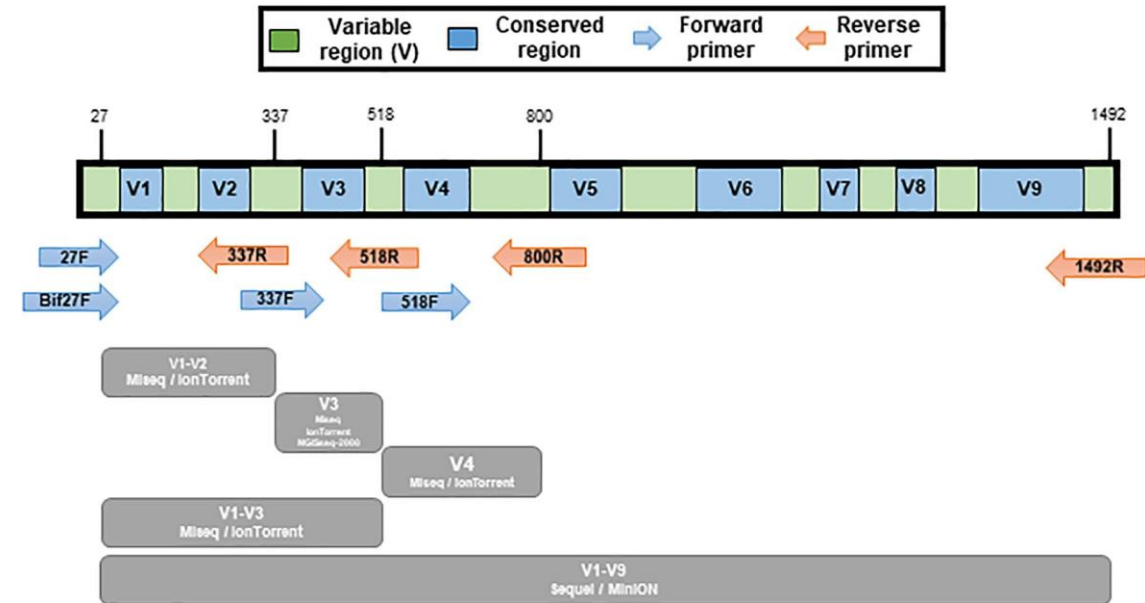- Highly **conserved** in some regions, which are easy to design appropriate primers to amplify by PCR**.**

# Marker genes in bacteria

**16S rRNA gene**

The sequence of 16S rRNA gene:
- Is **universal** in bacteria.
- Shows **evolutionary distance** and relationships between organisms and provides statistical and valid measurements for bacterial identification.
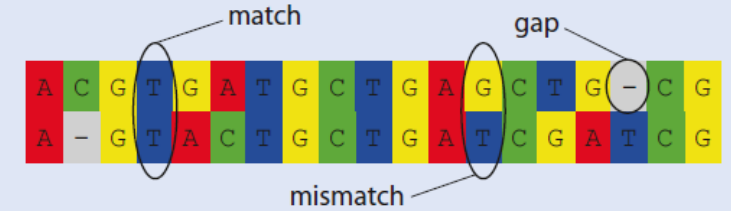- Comprises several **conversed** and **variable** regions useful to design primers.



(Adapted from Park et.al 2021)

# Alignments

For phylogenetic analyses, we should compare homologous genomic/gene positions → DNA sequences need to be aligned.
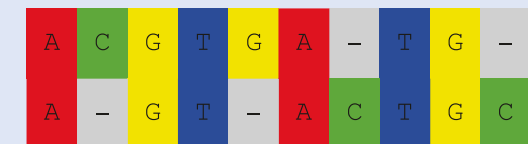
Alignments are hypotheses of **positional homologies** between nucleotides or amino acids of sequences.

The **Needleman and Wunsch algorithm** finds the optimal pairwise alignments of two sequences, which can contain matches, mismatches and gaps.



(Adapted from Bleidorn 2017)

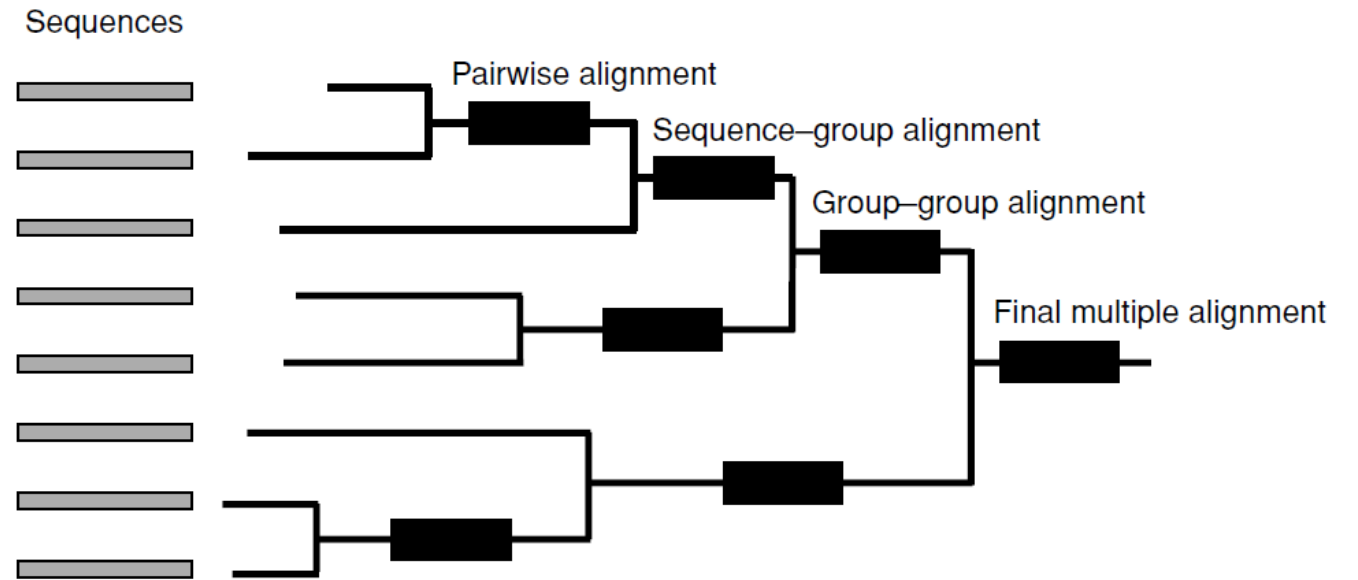# Alignments

**Multiple sequence alignments (MSA)**

Progressive MSA



(Adapted from Des Higgins and Philippe Lemey in *The Phylogenetic Handbook*)

# Alignments

For HTS data, are needed **fast** and accurate **multiple sequence alignments** (MSA).

Some popular **MSA** programs for HTS:

MAFFT  - https://mafft.cbrc.jp/alignment/software/
MUSCLE  - https://drive5.com/muscle5/
T-COFFEE  - https://tcoffee.crg.eu/apps/tcoffee/index.html

# Alignments

**MAFFT** is a MSA program for unix-like operating systems.  It offers a range of multiple alignment methods, L-INS-i (accurate; for alignment of <~**200** sequences), FFT-NS-2 (fast; for alignment of <~**30,000** sequences).

-command line and online version (limited to 100 sequences)
- L-INS-i is one of the most accurate MSA methods currently available.
- suitable for SSU rRNA alignments

MAFFT version 7
Multiple alignment program for amino acid or nucleotide sequences

Download version
    Mac OS X
    Windows
    Linux
    Source
**Online version**
  **Alignment**
    mafft --add
    Merge
    Phylogeny
    Rough tree
Merits / limitations
Algorithms
Tips
Benchmarks
Feedback
  Follow

This service was unstable due to maintenance, 18:00 – 21:00, May 23, JST.

To avoid overload, try a light-weight option, for MSA of full-length SARS-CoV-2 genomes (2020/Apr).

For a large number of short sequences, try an experimental service.

Experimental service for aligning raw reads (2019/Aug)

If you need an MSA of only a specific region, then try extracting the region first (2022/Oct). *New!*

Multiple sequence alignment and NJ / UPGMA phylogeny

**Input:**
Paste protein or DNA sequences in fasta format.  Example

Usage

```
% mafft [arguments] input > output
```

# Alignments

**MUSCLE v5** performs accurate and faster MSAs, capable of escalating to larger datasets.

- on large datasets, Muscle v5 is 20-30% more accurate than MAFFT
- No graphical user interface

## Muscle5

MUSCLE has been cited by
### 51,053 papers
Google scholar
Last updated 22 Jun 2023
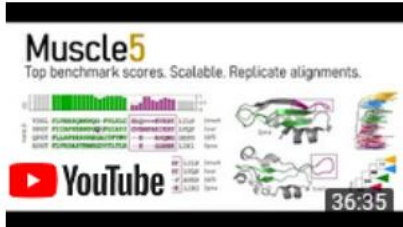
**Download**

**Documentation**

**Support and feedback**

**MUSCLE v3**

### Next-generation MUSCLE
Muscle v5 is a major re-write of MUSCLE based on new algorithms.

### Highest accuracy, scalable to thousands of sequences
Compared to previous versions, Muscle v5 is much more accurate, is often faster, and scales to much larger datasets. At the time of writing (late 2021), Muscle v5 has the highest scores on multiple alignment benchmarks including Balibase, Bralibase, Prefab and Balifam. It can align tens of thousands of sequences with high accuracy on a low-cost commodity computer (say, an 8-core Intel CPU with 32 Gb RAM). On large datasets, Muscle v5 is 20-30% more accurate than MAFFT and Clustal-Omega.

### Alignment ensembles
Muscle v5 can generate ensembles of high-accuracy alternative alignments. All replicates have equal average accuracy on benchmark test, including the MSA made with default parameters. By comparing results of downstream analysis (trees, structure prediction...) on different replicates, you can assess the effects of alignment errors on your study.

Muscle5
Top benchmark scores. Scalable. Replicate alignments.

▶ YouTube                                                    36:35

## Examples

```
muscle -align seqs.fa -output aln.afa
```

## Alignments

**T-Coffee** can be used to align sequences or to combine the output of your favourite alignment methods (Clustal, Mafft, Probcons, Muscle...) into one unique alignment (M-Coffee).

-web interface and command line option

**T COFFEE**

Home    History    Tutorial    References    Contacts    Projects

T-Coffee
*Aligns DNA, RNA or Proteins using the default T-Coffee*

Sequences input
*Paste or upload your set of sequences in FASTA format*

Sequences to align
Click here to use the sample file

- OR - Click here to upload a file

| mode | **Expresso** | |
|------|----------|---|
| Seq. | All | `t_coffee foo.seq -mode expresso` |
| Accuracy | High | |

# Alignments

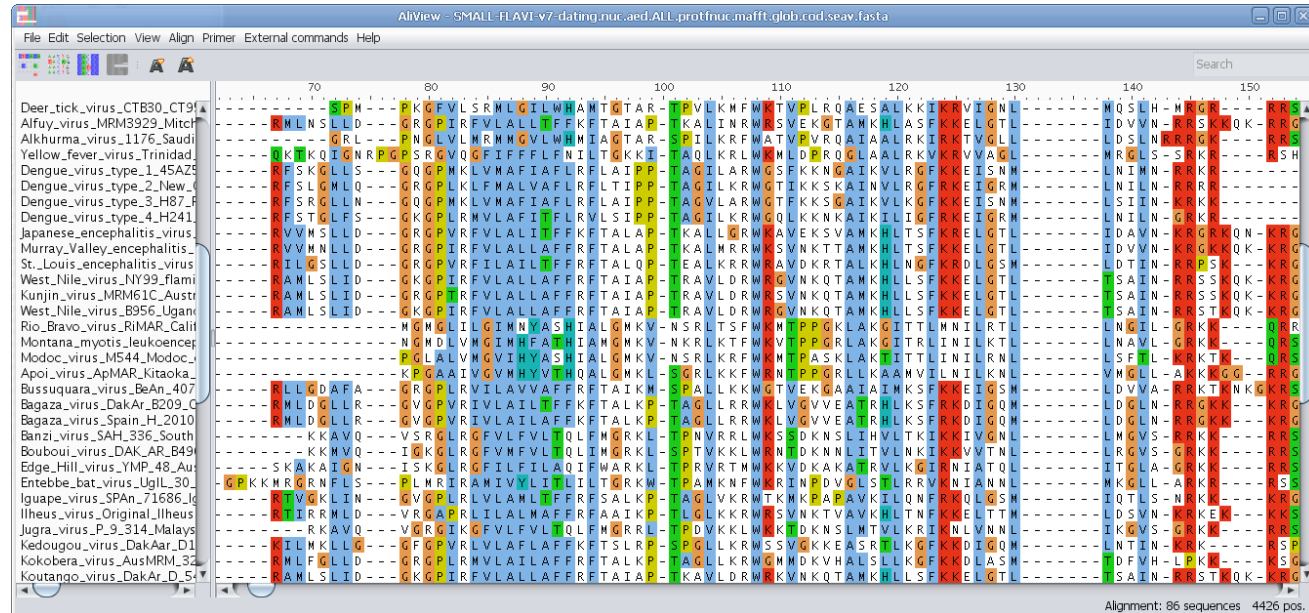**Visualization of MSA** from HTS data

**AliView** - https://ormbunkar.se/aliview/
- one of the fastest alignment viewer and editor
- created to work with large alignments.

Automatic **trimming of MSA**
**trimAl** - http://trimal.cgenomics.org/trimal

- automated removal of spurious sequences or poorly aligned regions from a MSA
- Command line and webserver available

**AliView**





trimAl
A tool for automated alignment trimming

# Phylogenetic analysis

Most commonly used **phylogenetic methods** for HTS data:

**Maximum-parsimony** – Character-based method - Which tree has the fewest mutations?
- fast, appropriate for very similar sequences and a small number of sequences

**Maximum-Likelihood** – Probabilistic method - Which tree has the highest probability given the observed alignment?
- Suitable for very dissimilar sequences

For **HTS** data, popular softwares include:
RaxML - https://github.com/amkozlov/raxml-ng
FastTree - http://www.microbesonline.org/fasttree/
PhyML - http://www.atgc-montpellier.fr/phyml/
IQ-TREE -  http://www.iqtree.org/

## Phylogenetic analysis

**Software for phylogenetic inference**

**RaxML** - https://github.com/amkozlov/raxml-ng

RAxML-NG is a phylogenetic tree inference tool which uses maximum-likelihood (ML) optimality criterion.

-developed for handling large datasets with low memory consumption and advanced search algorithms

-command line version

**Usage example**

```
./raxml-ng --search1 --msa testDNA.fa --model GTR+G
```

## Phylogenetic analysis

**Software for phylogenetic inference**

**IQ-TREE** - http://www.iqtree.org/

A fast and effective algorithm to infer phylogenetic trees by ML.

- IQ-TREE compares favorably to RAxML and PhyML in terms of likelihoods with similar computing time.

- supports datasets with thousands of sequences or millions of alignment sites
- -command line and webserver available (http://iqtree.cibiv.univie.ac.at/)



**IQ-TREE**

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

| Server load: 69% | Trifinopoulos J, Nguyen LT, von Haeseler A, Minh BQ (2016) *Nucl. Acids Res.* 44 (W1): W |

**Tree Inference** | **Model Selection** | **Analysis Results**

**For a quick start, take a look at the** tutorial **for the IQ-TREE web server.**
**Please visit the** IQ-TREE homepage **for more information or if you want to download the main software.**
Data Privacy Statement: All your personal data are strictly confidential and will not be shared with any third parties. Your dat will be automatically deleted after 180 days.

Input Data

**Alignment file :** [                    ] [Browse...] [Show example >]

**Use example alignment:** ☐ Yes                                                      [?]

**Sequence type:** ⦿ Auto-detect ○ DNA ○ Protein ○ Codon   [?]
○ DNA->AA ○ Binary ○ Morphology

**Partition file:** [This field is optional.] [Browse...] [Show example >]
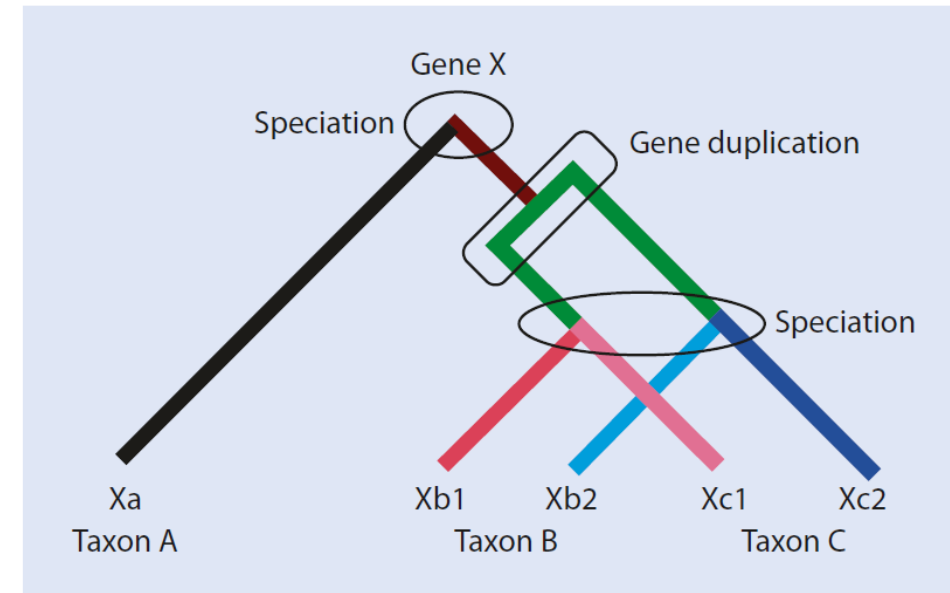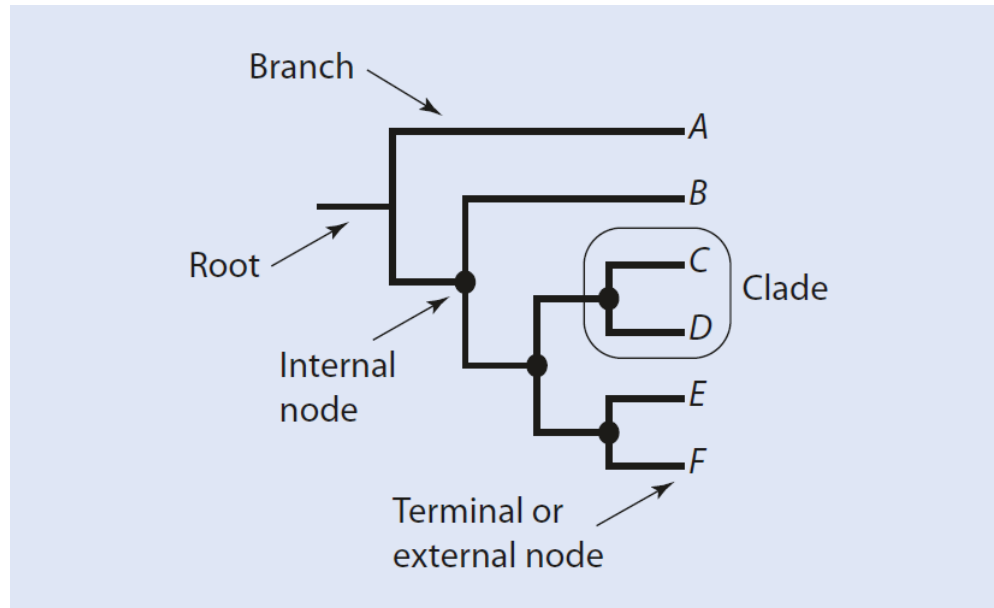
Partition type: ⦿ Edge-linked                                            [?]
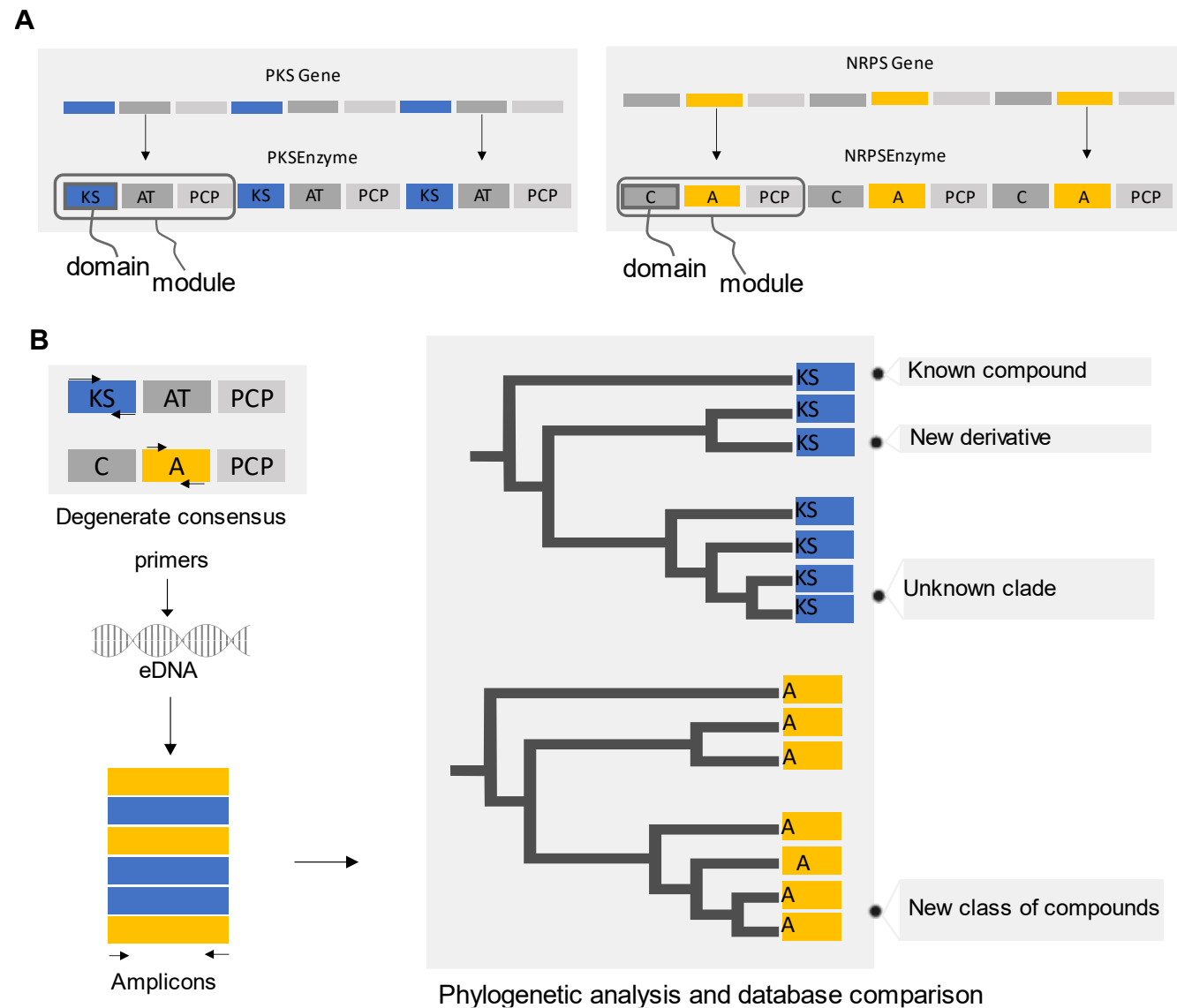○ Edge-unlinked

# How to interpret a phylogenetic tree?



(Adapted from Bleidorn 2017)

## Functional genes



**A**

PKS Gene

PKSEnzyme

| KS | AT | PCP | KS | AT | PCP | KS | AT | PCP |

domain
module

NRPS Gene

NRPSEnzyme

| C | A | PCP | C | A | PCP | C | A | PCP |

domain
module

**B**

| KS | AT | PCP |

| C | A | PCP |

Degenerate consensus

primers

eDNA

Amplicons

KS — Known compound
KS — New derivative
KS
KS
KS — Unknown clade
KS

A
A
A

A
A
A — New class of compounds
A

Phylogenetic analysis and database comparison

**Biosynthetic domains (KS and A)** are highly conserved and have proven to be very informative in a phylogenetic context.
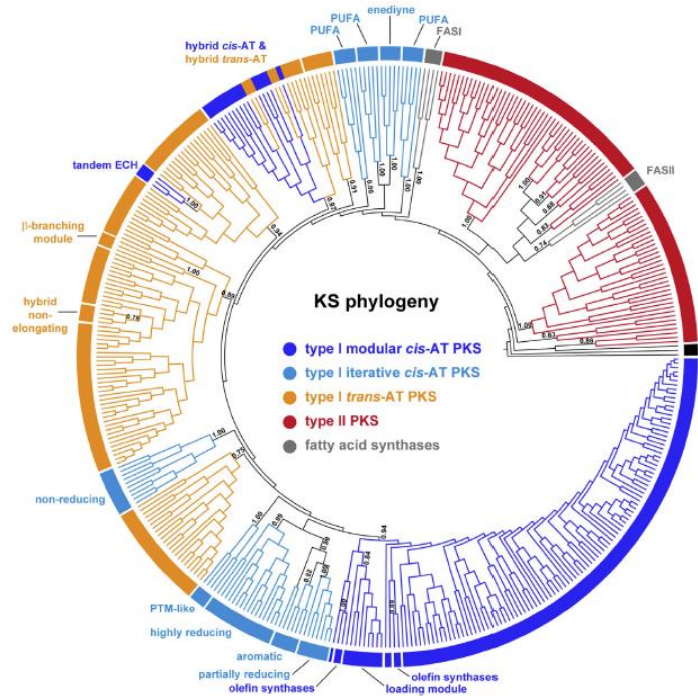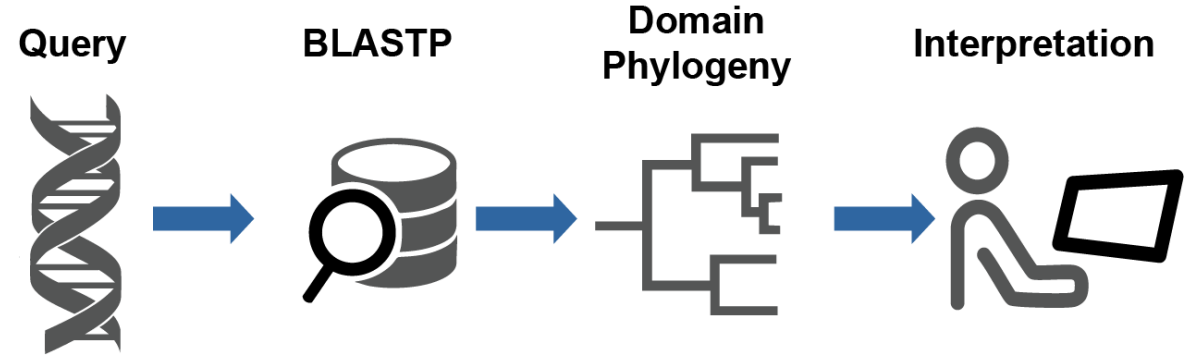
Biosynthetic gene clusters (BGCS) encoding for structurally related metabolites are predicted to share **common ancestry** and exhibit high sequence **identity** among these conserved domains.

(Rego 2023)

## Functional genes

**NAPDOS2.0**

NaPDoS2, the second-generation Natural Product Domain Seeker, rapidly detects and classifies **ketosynthase (KS)** and **condensation (C)** domains from genomic, metagenomic, or PCR amplicon sequence data.

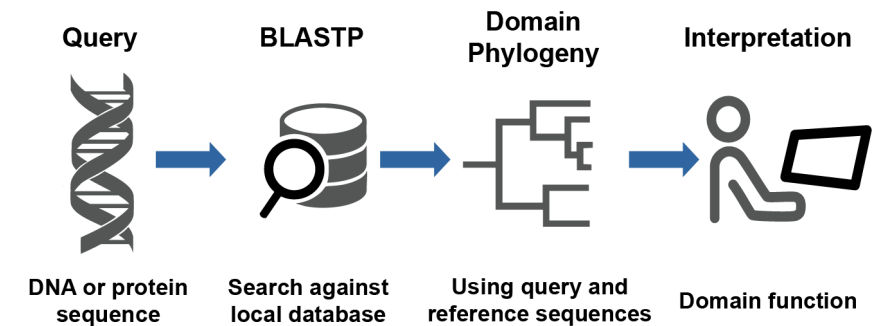https://npdomainseeker.sdsc.edu/napdos2/

## Functional genes

NAPDOS2.0

**Main applications of NAPDOS approach**

Poorly assembled (fragmented) genomes and metagenomes
Amplicons
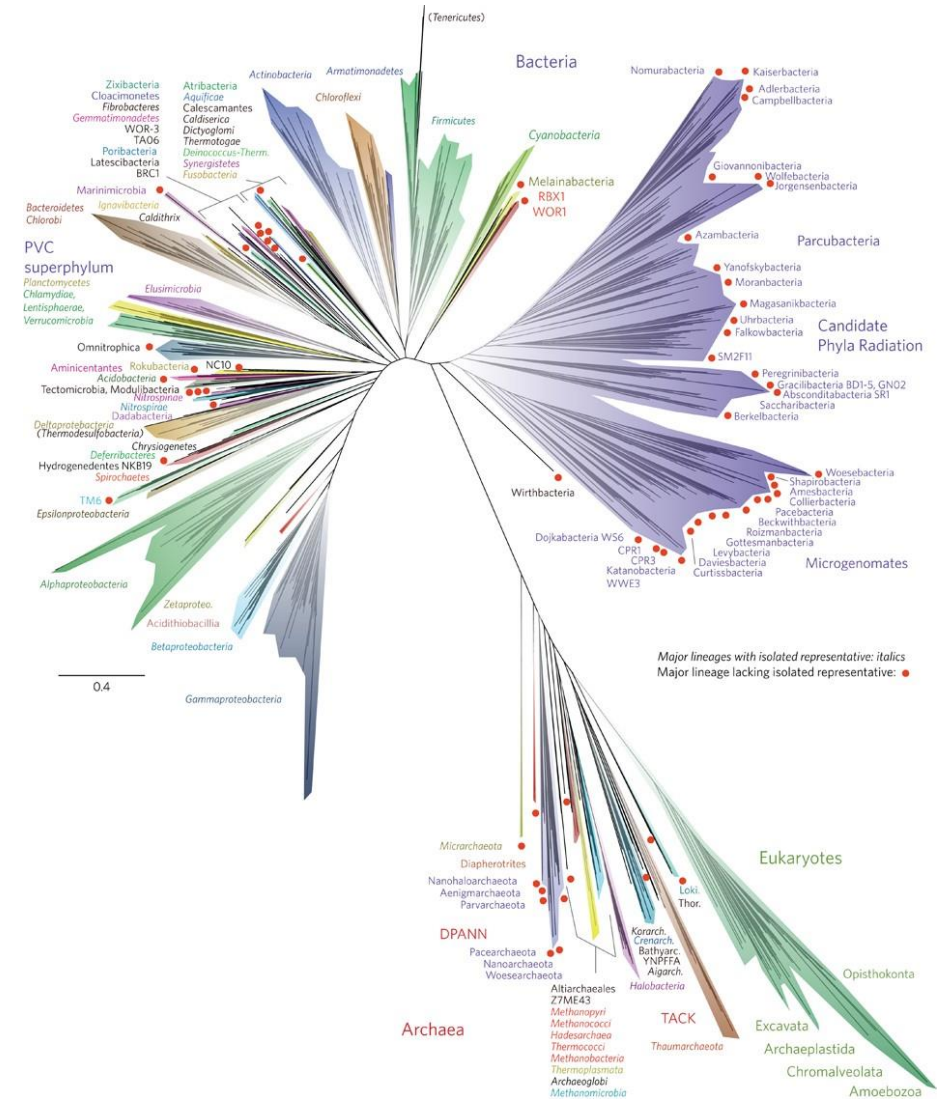Fast estimate of biosynthetic potential
New clades = new functionality

Limitations: comparison to database, new diversity might not be detected

| Query | BLASTP | Domain Phylogeny | Interpretation |
|---|---|---|---|
| DNA or protein sequence | Search against local database | Using query and reference sequences | Domain function |

# Phylogenomics

**Phylogenomics** involves the reconstruction of evolutionary relationships by comparing sequences of **whole genomes** or **portions** of genomes.
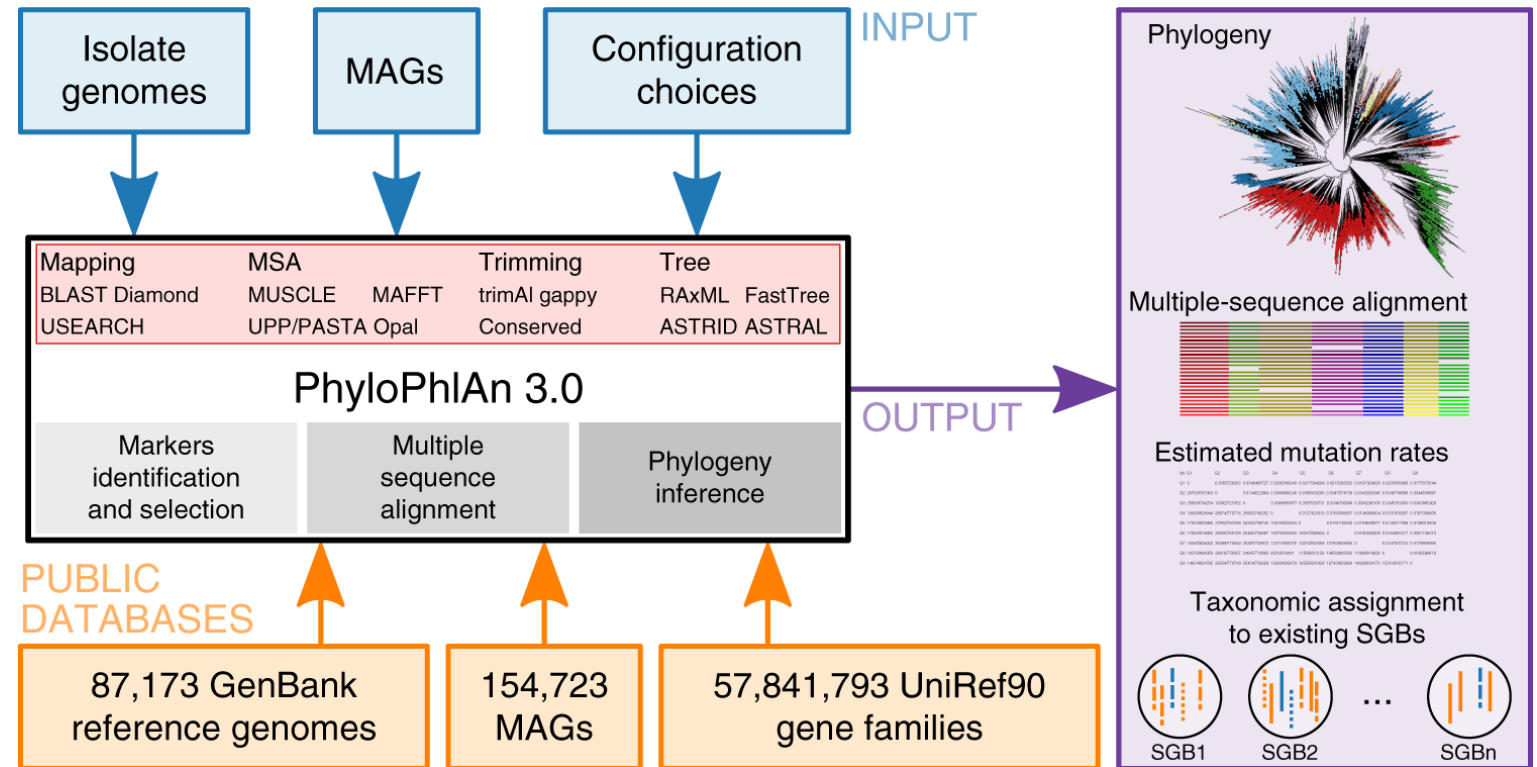
Genome-scale phylogenetic datasets yield an increase in the **statistical confidence** of inferred relationships, often yielding maximally **supported** species trees.



A current view of the tree of life, encompassing the total diversity represented by **sequenced genomes** (Hug et al. 2016).

## PhyloPhlan

**PhyloPhlAn** can reconstruct strain-level phylogenies using clade-specific informative **phylogenetic markers**, and can also scale to very-large phylogenies comprising **>17,000** microbial species.



https://huttenhower.sph.harvard.edu/phylophlan/

## Visualization of phylogenetic/phylogenomic trees

**Web Portals**

Interactive Tree of Life viewer (iTOL) – https://itol.embl.de/
Evolview - https://www.evolgenius.info/help/

**Software/libraries**

FigTree - http://tree.bio.ed.ac.uk/software/figtree/
GraphLan - https://huttenhower.sph.harvard.edu/graphlan/
ggtree: an R package for visualization of phylogenetic trees with their annotation data-
https://github.com/YuLab-SMU/ggtree

https://en.wikipedia.org/wiki/List_of_phylogenetic_tree_visualization_software

**Visualization of phylogenetic/phylogenomic trees**

**Web Portals**

Interactive Tree of Life viewer (**iTOL**) – https://itol.embl.de/

- Web-based
- User friendly (interactive interface)
- Highly costumizable



Welcome to iTOL v6

**Interactive Tree Of Life** is an online tool for the display, annotation and management of phylogenetic and other trees.

Manage and visualize your trees directly in the browser, and annotate them with various datasets.

## Visualization of phylogenetic/phylogenomic trees

**Web Portals**

Interactive Tree of Life viewer (**iTOL**) – https://itol.embl.de/

## Visualization of phylogenetic/phylogenomic trees

**Web Portals**

Evolview -
https://www.evolgenius.info/help/

- Web-based
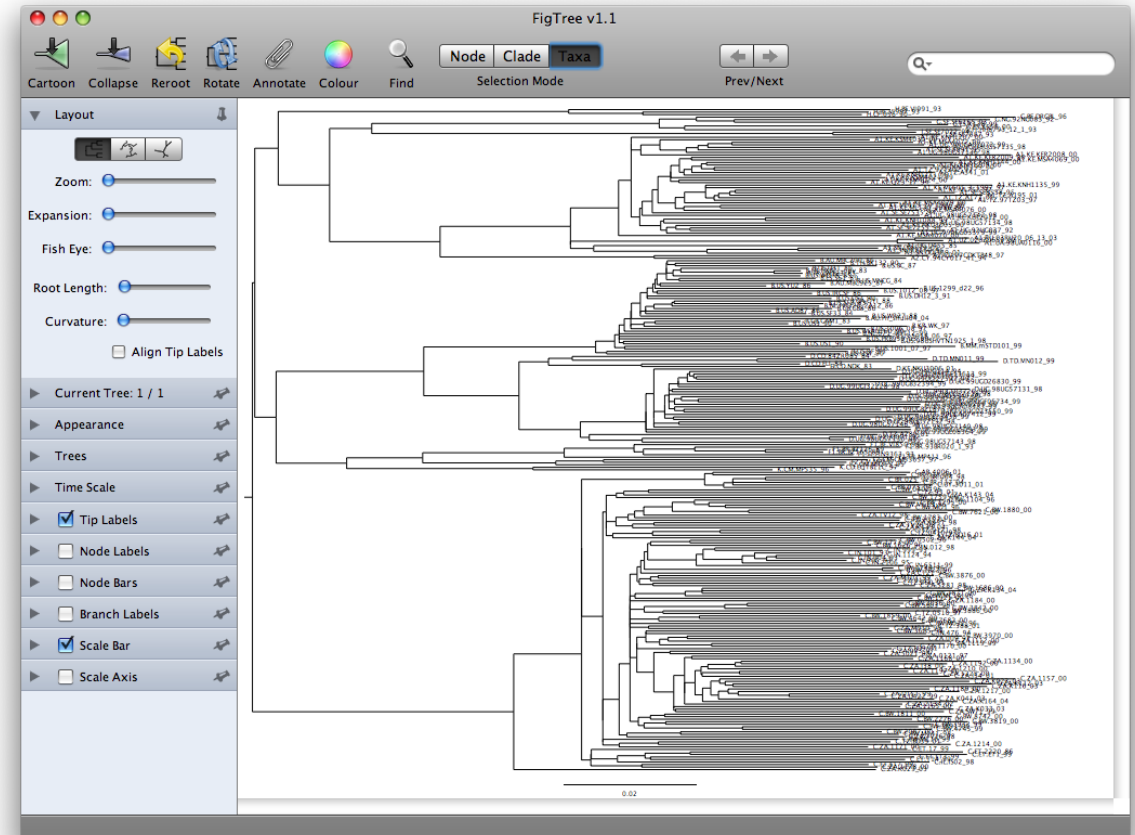- User friendly (interactive interface)

## Visualization of phylogenetic/phylogenomic trees

**Software**

FigTree
- graphical viewer of phylogenetic trees and a program for producing publication-ready figures.
-available for Mac, Windows and Linux
-interactive interface
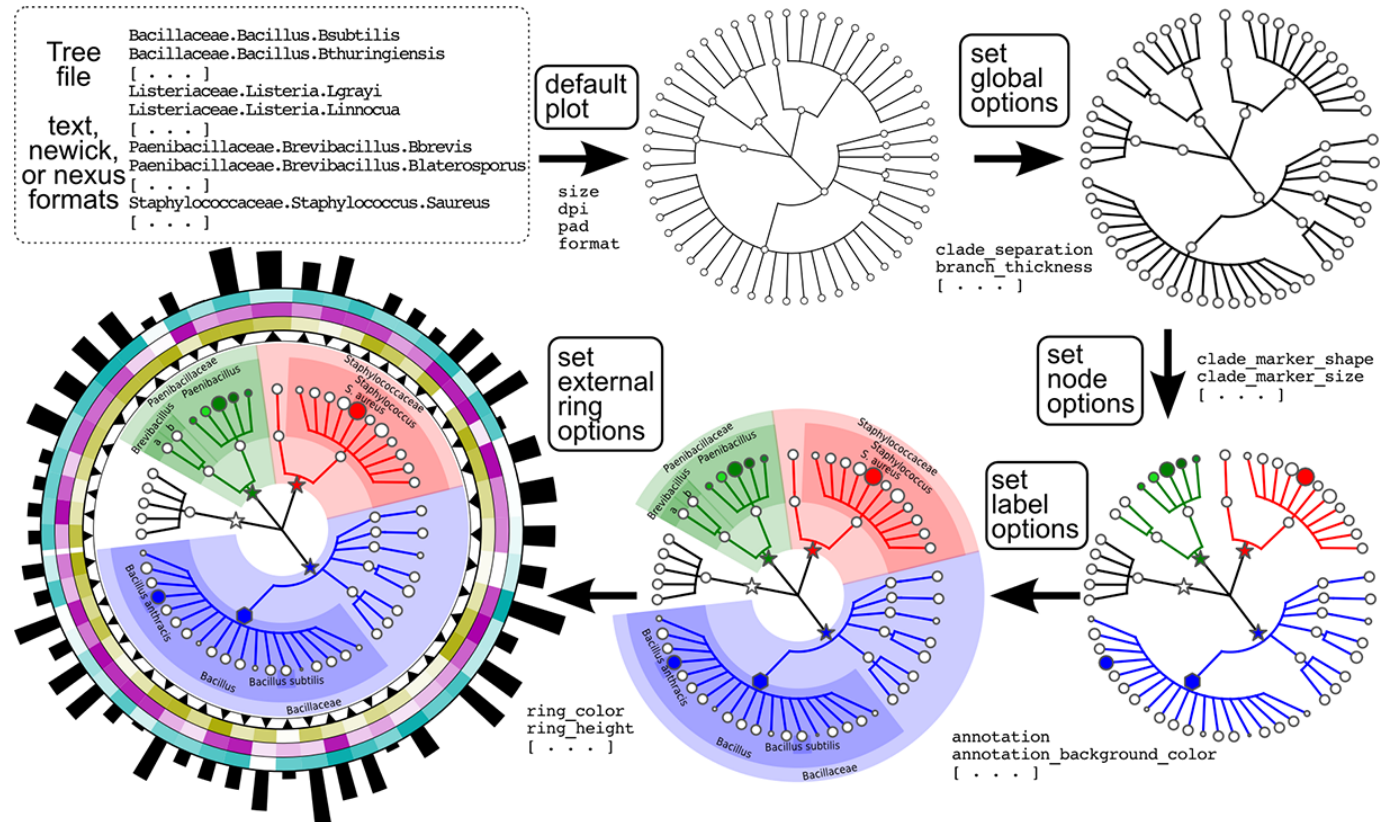
## Visualization of phylogenetic/phylogenomic trees

**Software**

**GraPhlAn**

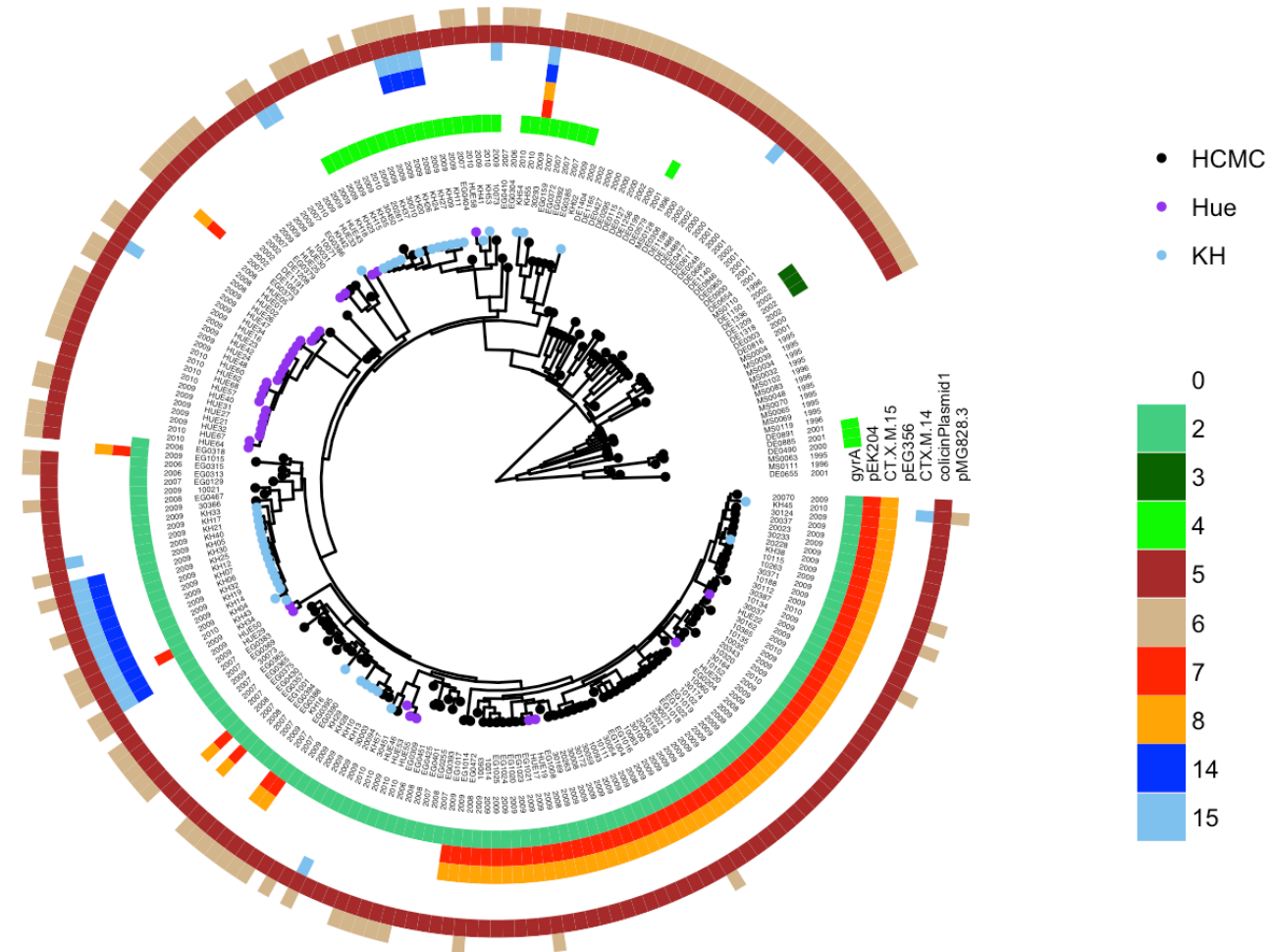https://github.com/biobakery/graphlan

- is a software tool for producing high-quality circular representations of taxonomic and phylogenetic trees.

  - publication-ready representations
  - command line

## Visualization of phylogenetic/phylogenomic trees

**Software**

R package **ggtree** - https://github.com/YuLab-SMU/ggtree
- an extension of the 'ggplot2' plotting system
- visualization and annotation of phylogenetic trees and other tree-like structures with their annotation data.
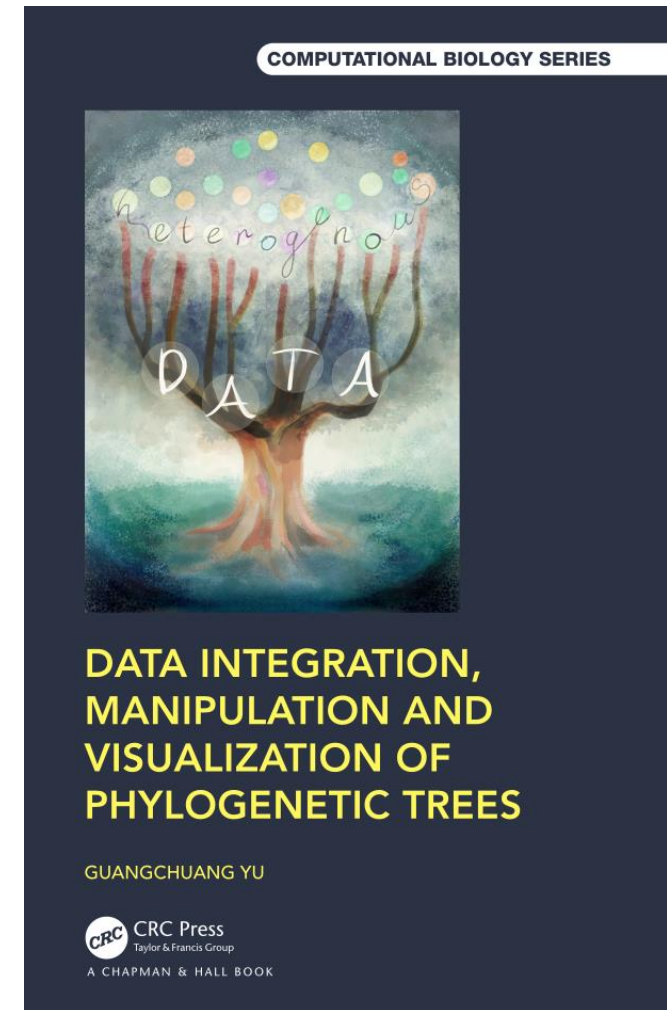
## Visualization of phylogenetic/phylogenomic trees

**Resources for visualization and annotation of phylogenetic trees using R**

This book (https://yulab-smu.top/treedata-book/index.html) is meant as a guide for data integration, manipulation and visualization of phylogenetic trees using a suite of R packages, tidytree, treeio, ggtree and ggtreeExtra.

**ggtree book** – https://guangchuangyu.github.io/ggtree-book/short-introduction-to-r.html
Practical examples of annotation and visualization of phylogenetic trees using ggtree.
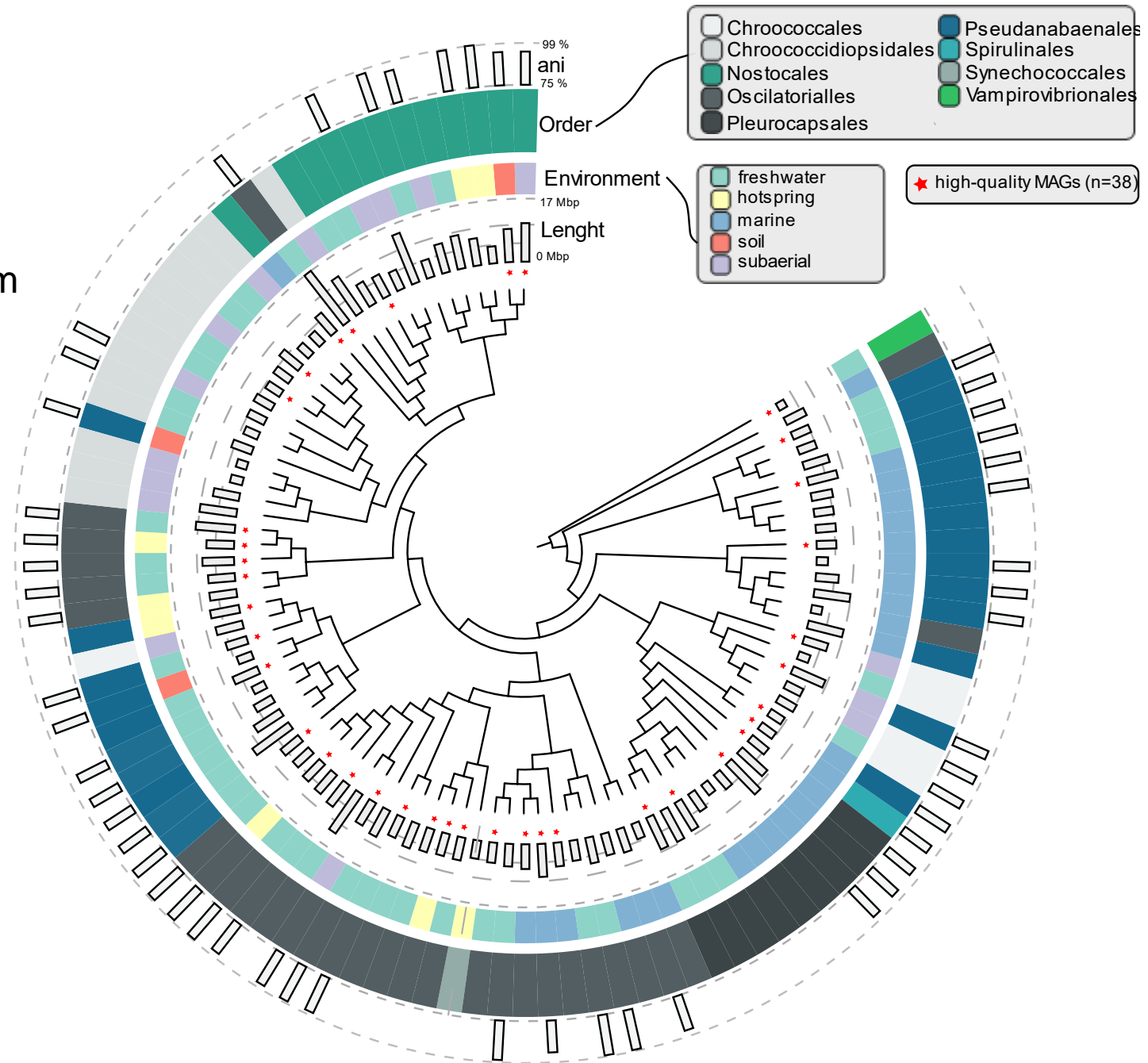


**Chapter 4 Visualization and annotation of phylogenetic trees: *ggtree***

**Pratical example**

**Cyanobacteria MAGs** - Samples collected from the environment
- Which Cyanonacteria orders are present?
- How novel these MAGs are?
- How are they distributed by environment type?
- What about genome lenght and completeness and contamination?



Legend — Order:
- Chroococcales
- Chroococcidiopsidales
- Nostocales
- Oscilatorialles
- Pleurocapsales
- Pseudanabaenales
- Spirulinales
- Synechococcales
- Vampirovibrionales

Environment:
- freshwater
- hotspring
- marine
- soil
- subaerial

★ high-quality MAGs (n=38)

ani 99 % / 75 %
Order
Environment
17 Mbp
Lenght 0 Mbp

**Pratical example**

**Workflow followed:**

Tree reconstruction - **Phylophlan**

Tree Edition – **iTOL** - include previously obtained
metadata (MAG/genome size,
completeness/contamination of MAG, ANI to reference
genomes) as datasets for iTOL

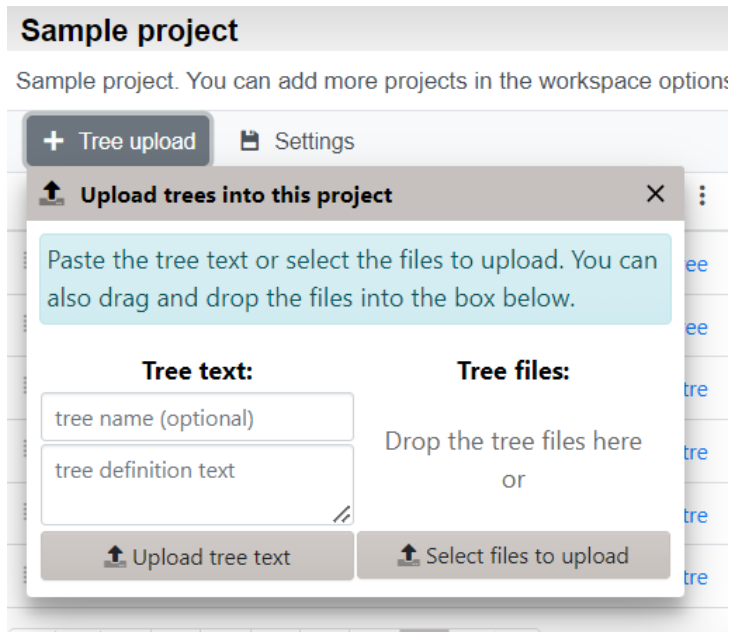**Import datasets to iTOL**
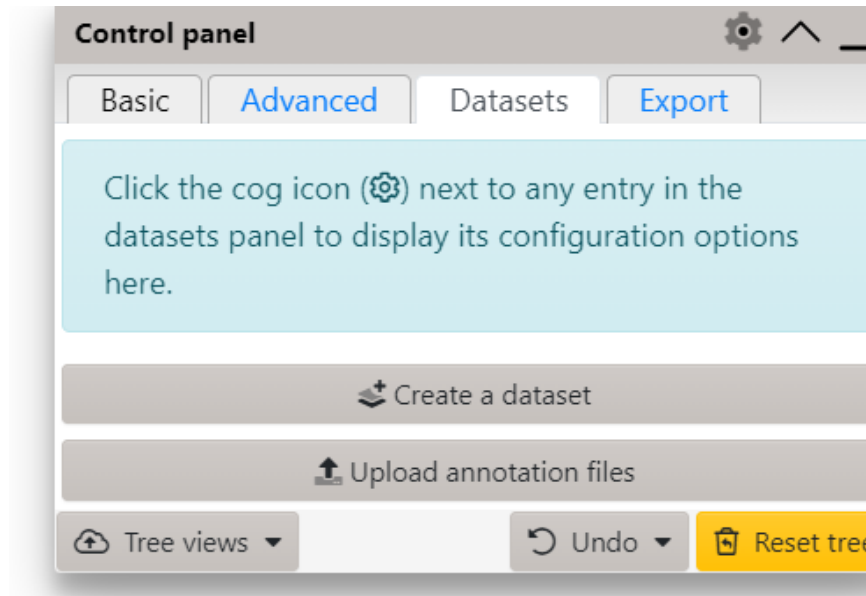
**Annotation template files – help page**

https://itol.embl.de/help.cgi

**itol.toolkit R package**

https://tongzhou2017.github.io/itol.toolkit/index.html

**Pratical example**

**Import tree into iTOL**

**Import datasets into iTOL**

**Control panel for datasets**

**Bibliography**

Christoph Bleidorn

# Phylo-genomics

An Introduction

Springer

The
**Phylogenetic
Handbook**

A Practical Approach to
Phylogenetic Analysis and
Hypothesis Testing

**Second Edition**

Edited by Philippe Lemey,
Marco Salemi and
Anne-Mieke Vandamme

## Glossary

**Clade** Monophyletic group in a phylogenetic tree, thereby representing at least two terminals (which share a common ancestor).

**Maximum likelihood (ML)** Likelihood-based optimality criterion to find the best tree in a phylogenetic analysis through the computation of probabilities of character evolution given an explicit evolutionary model.

**Maximum parsimony (MP)** Optimality criterion which selects the phylogenetic tree(s) minimizing the total number of character state changes.

**Monophyletic group** Group containing a (hypothetical) ancestor and all of its descendants.

**Orthology** Pairs of homologous genes which have emerged through a speciation event are called orthologs.

**Root** Point of a topology where it is hypothetically connected to the remaining tree of life. Rooted trees are used to polarize character evolution.

**Questions?**

# Thank you!

# Any questions?

✉ adrianairego@gmail.com