



Swiss Institute of  
Bioinformatics



# Computational comparative genomics in the era of BioGenome projects

(fast orthology and phylogeny tools)

Sina Majidian

Department of Computational Biology, University of Lausanne. SIB Swiss Institute of Bioinformatics.

November 6, 2023

# Speaker biography

2021- 2024 Postdoctoral Fellow, Switzerland.

**Comparative Genomics Lab**

(Christophe Dessimoz & Natasha Glover)

2018-2019 Guest researcher, Netherlands.

**Bioinformatics group** (Dick de Ridder)

2015-2020 PhD, Signal processing, Iran.



Iran University of  
Science and Technology



TEHRAN UNIVERSITY  
OF  
MEDICAL SCIENCES

# **Computational comparative genomics in the era of BioGenome projects**

(fast orthology and phylogeny tools)

## **Outline**

- Introduction
- Orthology inference with FastOMA
- Phylogeny inference from raw sequencing reads with read2tree

# BioGenome projects



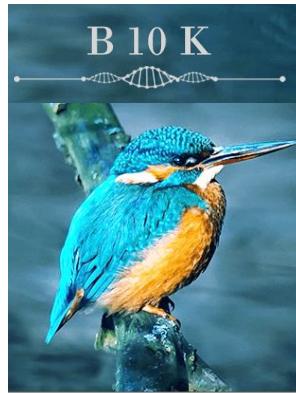
Sequence and characterize the genomes of all of Earth's eukaryotic biodiversity over ten years.



Darwin  
**TREE  
of  
LIFE**



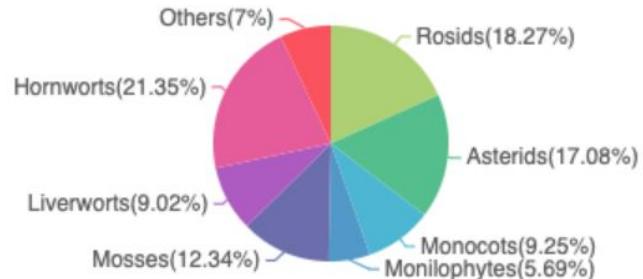
Biodiversity  
Genomics  
Europe



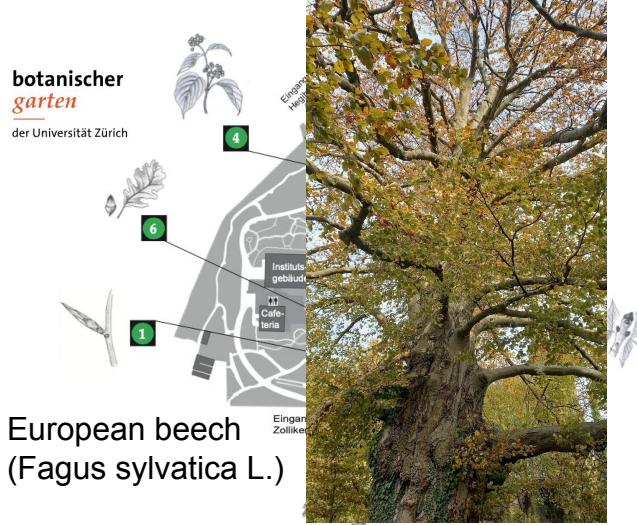
AQUATIC  
SYMBIOSIS  
GENOMICS



**10,000 Plants (10KP)** aims to sequence a member of every genus!

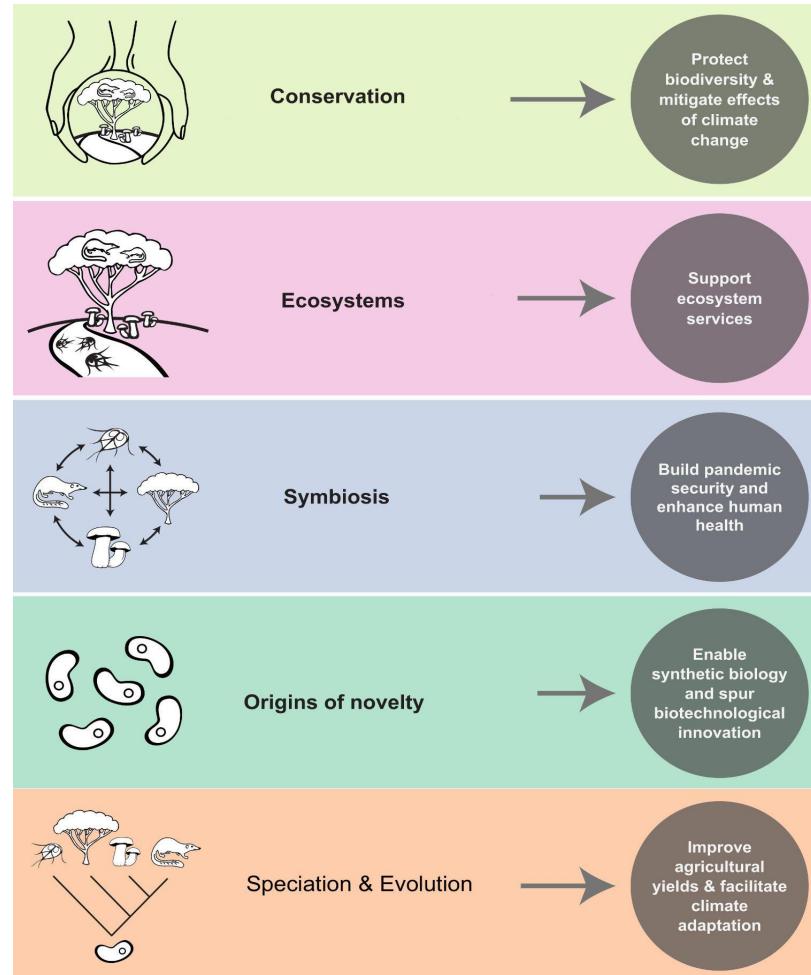


# Genome project benefits

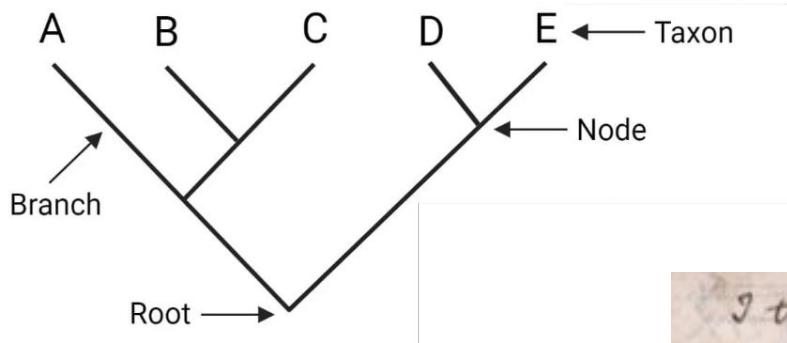


European beech  
(*Fagus sylvatica* L.)

Genomic basis for drought resistance in European beech forests threatened by climate change. Elife 2022.



# Evolutionary tree

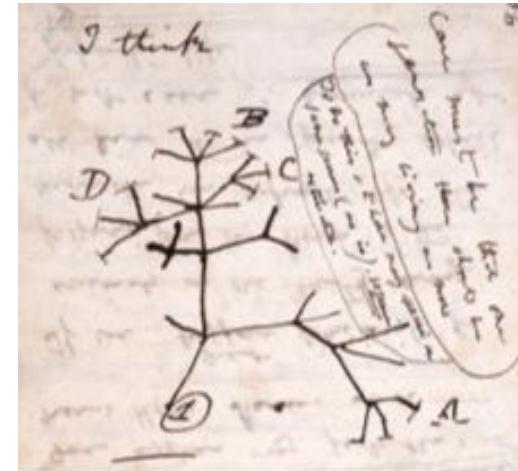


## Inferring species trees

- 1) 16s rRNA
- 2) Precomputed markers e.g. **BUSCO**

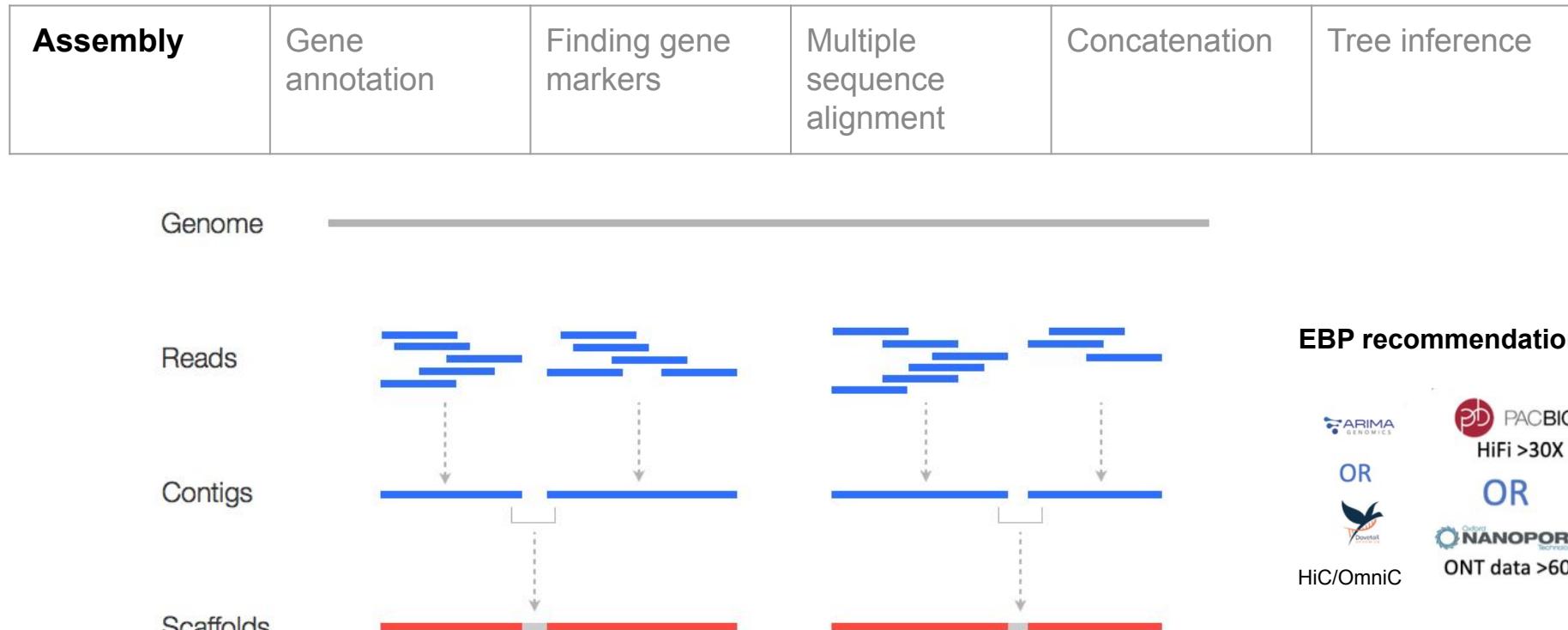
✗ not available for all clades

- 3) Gene markers



Darwin, 1837

# Inferring species trees with concatenation



# Inferring species trees with concatenation

Assembly	<b>Gene annotation</b>	Finding gene markers	Multiple sequence alignment	Concatenation	Tree inference
----------	------------------------	----------------------	-----------------------------	---------------	----------------

Finding genomic coordinates of coding regions

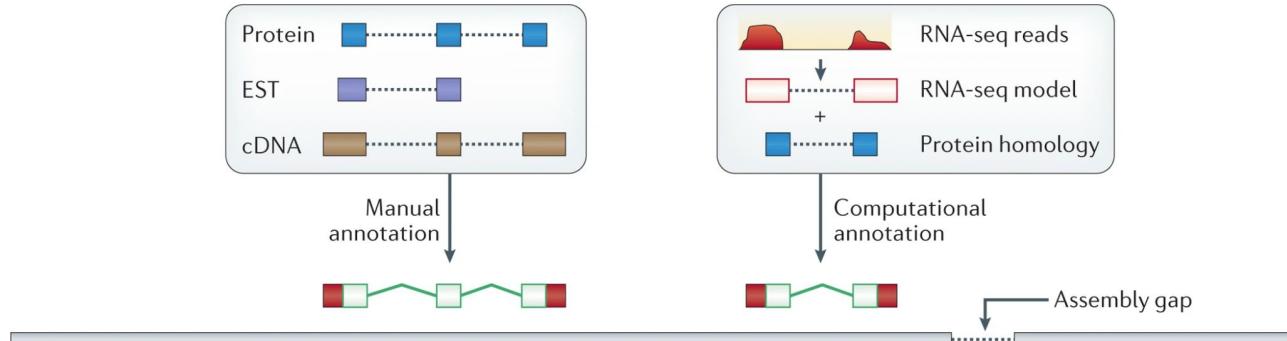


illumina<sup>®</sup>  
RNAseq  
->5 tissue types

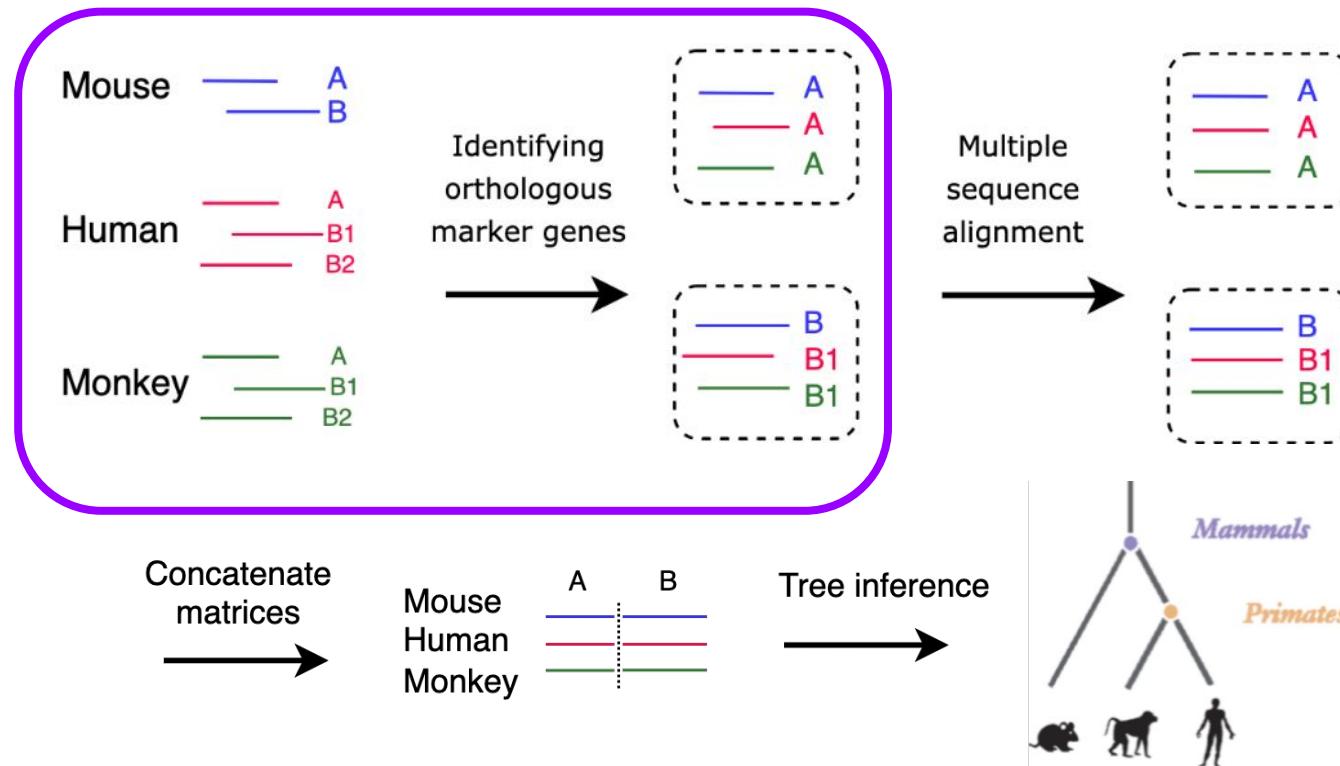
AND/OR

PACBIO<sup>®</sup>  
ISOseq

→ proteome



Assembly	Gene annotation	Finding gene markers	Multiple sequence alignment	Concatenation	Tree inference
----------	-----------------	----------------------	-----------------------------	---------------	----------------



# Homology

homologous regions or regions of common ancestry.

## · INS\_HUMAN

Protein<sup>i</sup> | Insulin

Gene<sup>i</sup> | INS



	Description	Scientific Name	Max Score	Total Score	Query Cover	E value	Per. Ident	Acc. Len	Accession
<input checked="" type="checkbox"/>	insulin isoform UB [Homo sapiens]	Homo sapiens	226	226	100%	7e-74	100.00%	153	QMS45324.1
<input checked="" type="checkbox"/>	insulin [Gorilla gorilla gorilla]	Gorilla gorilla go...	226	226	100%	8e-74	100.00%	153	XP_004050475.2
<input checked="" type="checkbox"/>	Homo sapiens insulin [synthetic construct]	synthetic construct	224	224	100%	2e-73	100.00%	111	AAP36446.1
<input checked="" type="checkbox"/>	insulin_preproprotein [Homo sapiens]	Homo sapiens	223	223	100%	2e-73	100.00%	110	NP_000198.1
<input checked="" type="checkbox"/>	insulin isoform X2 [Pongo abelii]	Pongo abelii	222	222	100%	1e-72	99.09%	110	XP_024110665.1
<input checked="" type="checkbox"/>	INS [synthetic construct]	synthetic construct	221	221	100%	1e-72	99.09%	110	AKI70564.1
<input checked="" type="checkbox"/>	insulin isoform U1 [Homo sapiens]	Homo sapiens	225	225	100%	1e-72	100.00%	204	QMS45321.1
<input checked="" type="checkbox"/>	insulin_preproprotein [Pan troglodytes]	Pan troglodytes	221	221	100%	2e-72	98.18%	110	NP_001008996.1

GORGO03436	M A L W M R - L L P L L A L L A L - W G P D P A A A F V N Q H L C G S H L V E A L Y L V C G E R G F F - Y T P - K T R R E A E D L Q -
HUMAN03911	M A L W M R - L L P L L A L L A L - W G P D P A A A F V N Q H L C G S H L V E A L Y L V C G E R G F F - Y T P - K T R R E A E D L Q -
PONAB02480	M A L W M R - L L P L L A L L A L - W G P D P A - A F V N Q H L C G S H L V E A L Y L V C G E R G F F - Y T P - K T R R E A E D L Q -
NOMLE29895	M A L W M R - L L P L L A L L A L - W G P D P A P A F V N Q H L C G S H L V E A L Y L V C G E R G F F - Y T P - K T R R E A E D P Q -
AOTNA33386	M A L W M H - L L P L L A L L A L - W G P E P A P A F V N Q H L C G P H L V E A L Y L V C G E R G F F - Y A P - K T R R E A E D L Q -
CEBIM30277	M A L W M H - L L P L L G L L A L - W G P E P A P A F V N Q H L C G P H L V E A L Y L V C G E R G F F - Y A P - K T R R E A E D L Q -
SAIBB37889	M A L W M H - L L P L L A L L A L - W G P E P A P A F V N Q H L C G P H L V E A L Y L V C G E R G F F - Y A P - K T R R E A E D P Q -
MICMU00461	M A L W T R - L L P L L A L L A L - W G P E P A P A F V N Q H L C G S H L V E A L Y L V C G E R G F F - Y T P - K S R R E V E D A Q -
OTOGA04835	M A V W M R - L L P L L A L L A L - W G P E P A P A F V N Q H L C G S H L V E A L Y L V C G E R G F F - Y T P - K A R R D T E P D Q -
TUPBE03840	M A L W T C - F L P L L T L L A L - W G P E P A P A F V N Q H L C G S H L V E A L Y L V C G E R G F F - Y T P - K T R R E V E D S Q -
CANLF05920	M A L W M R - L L P L L A L L A L - W A P A P T R A F V N Q H L C G S H L V E A L Y L V C G E R G F F - Y T P - K A R R E V E D L Q -
VULUV17698	M A L W M R - L L P L L A L L A L - W A P A P T R A F V N Q H L C G S H L V E A L Y L V C G E R G F F - Y T P - K A R R E V E D L Q -
SURSU35615	M A L W T R - L L P L L A L L A L - W A P A P A R G F V N Q H L C G S H L V E A L Y L V C G E R G F F - Y T P - K A R R E A G D L Q -
MYOLU19828	M A L W T R - L L P L L A L L A L - W A P A P A Q A F N H E H L C G E D L V D I M T I I C G D Q G F F - K N P - K A A R E L P D P Q -
CERSS29869	M A L W T R - L L P L L A L L A L - W S P A P T R A F V N Q H L C G S H L V E A L Y L V C G E R G F F - Y T P - K A R R E A E D P Q -

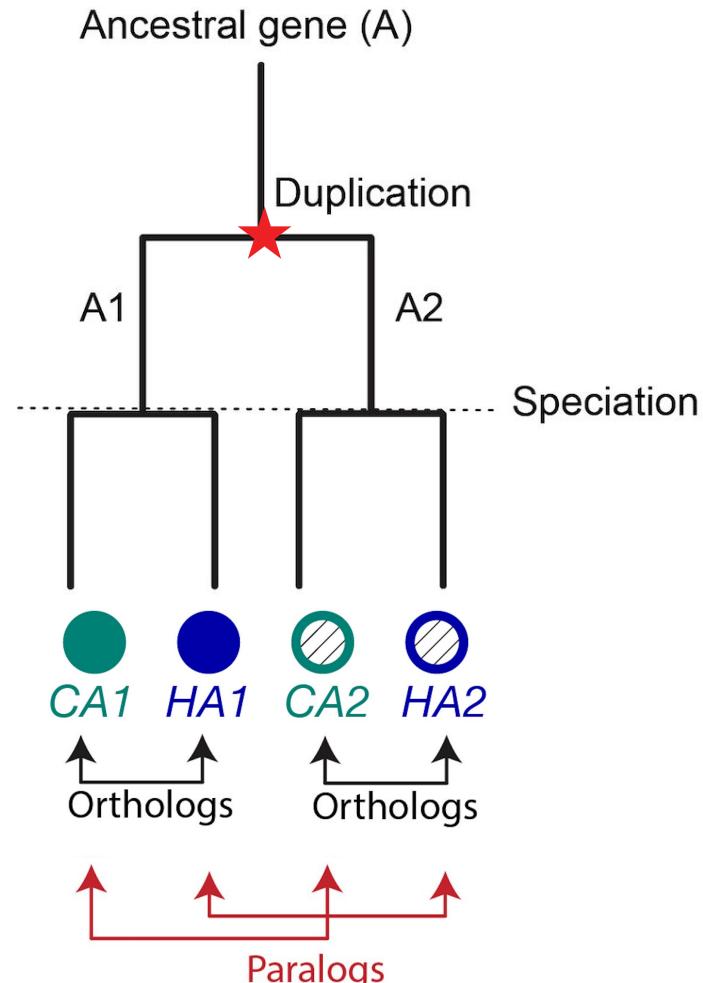
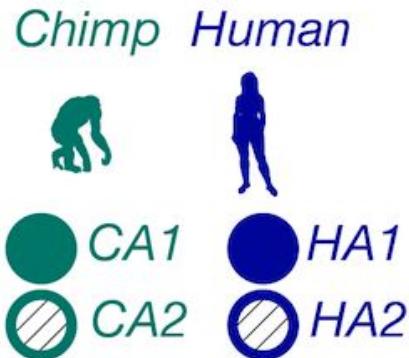
Similarity implies common ancestry.

# Orthology vs. paralogy

Two classes of homologous genes:

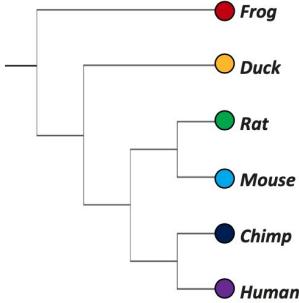
- **Orthology:** A relation between pairs of genes that started diverging via evolutionary **speciation**.
  - useful for species tree inference.
  - tend to have conserved functions.
- **Paralogy:** A relation between pairs of genes that started diverging via gene **duplication**.

# Orthology vs. paralogy

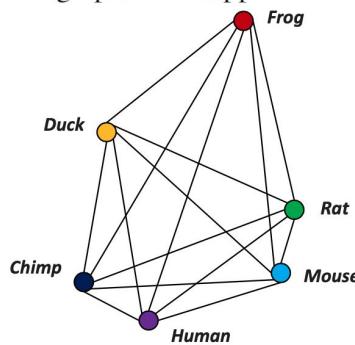


# Orthology inference methods

tree-based approaches



graph-based approaches

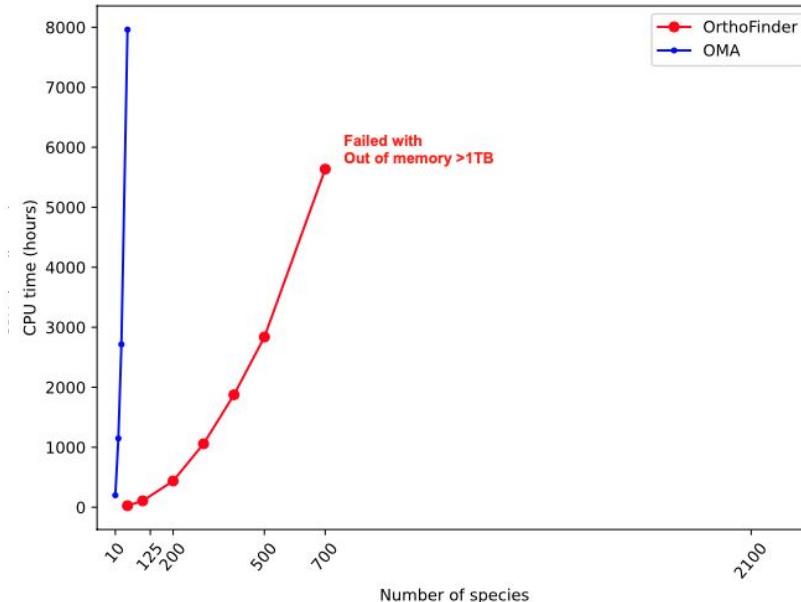


✗ All-vs-all comparison

Method/resource	# species	
<b>Oma</b> orthologous matrix browser	2,851	Graph based
<b>OrthoDB</b>	20,110	Graph based
<b>DB phylome</b>	6,000	Tree based
<b>EggNOG 6.0.0</b>	12,535	Graph + tree
<b>OrthoFinder 2</b>	-	Graph + tree

# Is orthology inference possible for >1000 species ?

- UniProt reference eukaryotic proteomes



We need new concepts  
and methods to handle  
data from BioGenome  
projects!

# Orthologous groups

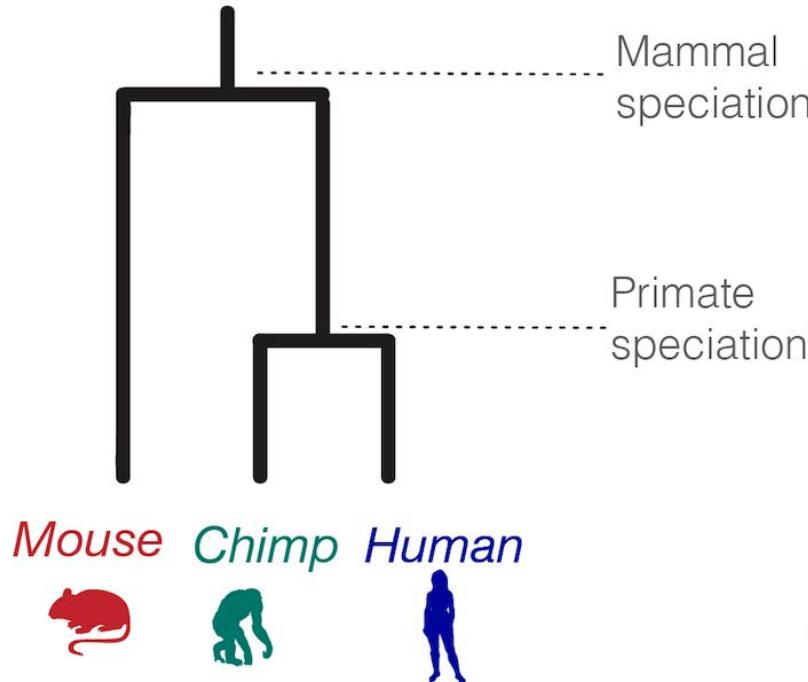
- **Orthologous groups:** Groups of genes which are all orthologous to each other.
- **Hierarchical Orthologous Group (HOG):** Groups of genes descended from a common ancestral gene at a specific taxonomic level

# HOG

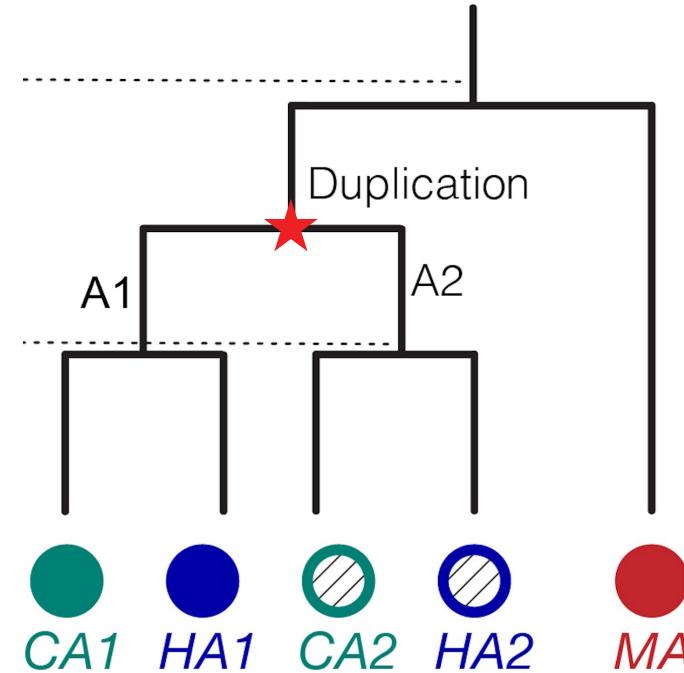
Hierarchical Orthologous Group

- Genes descended from a common ancestral gene at a specific taxonomic level

Species tree



Gene tree

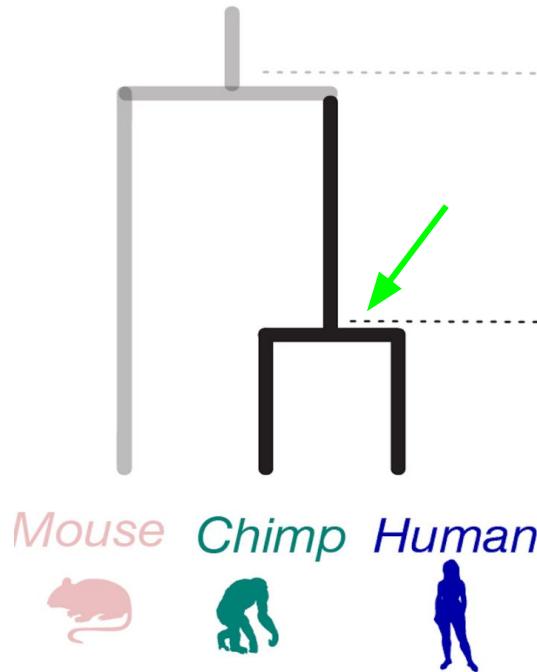


# HOG

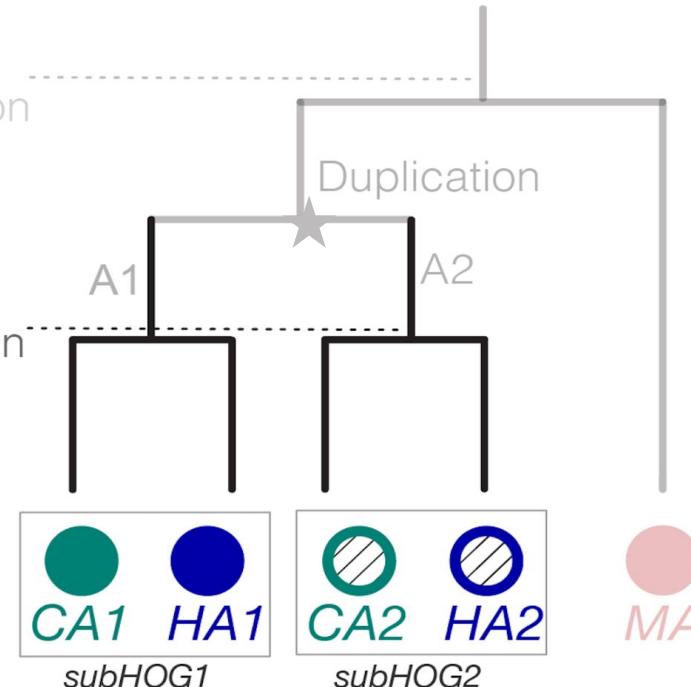
Hierarchical Orthologous Group

- Genes descended from a common ancestral gene at a specific taxonomic level

Species tree



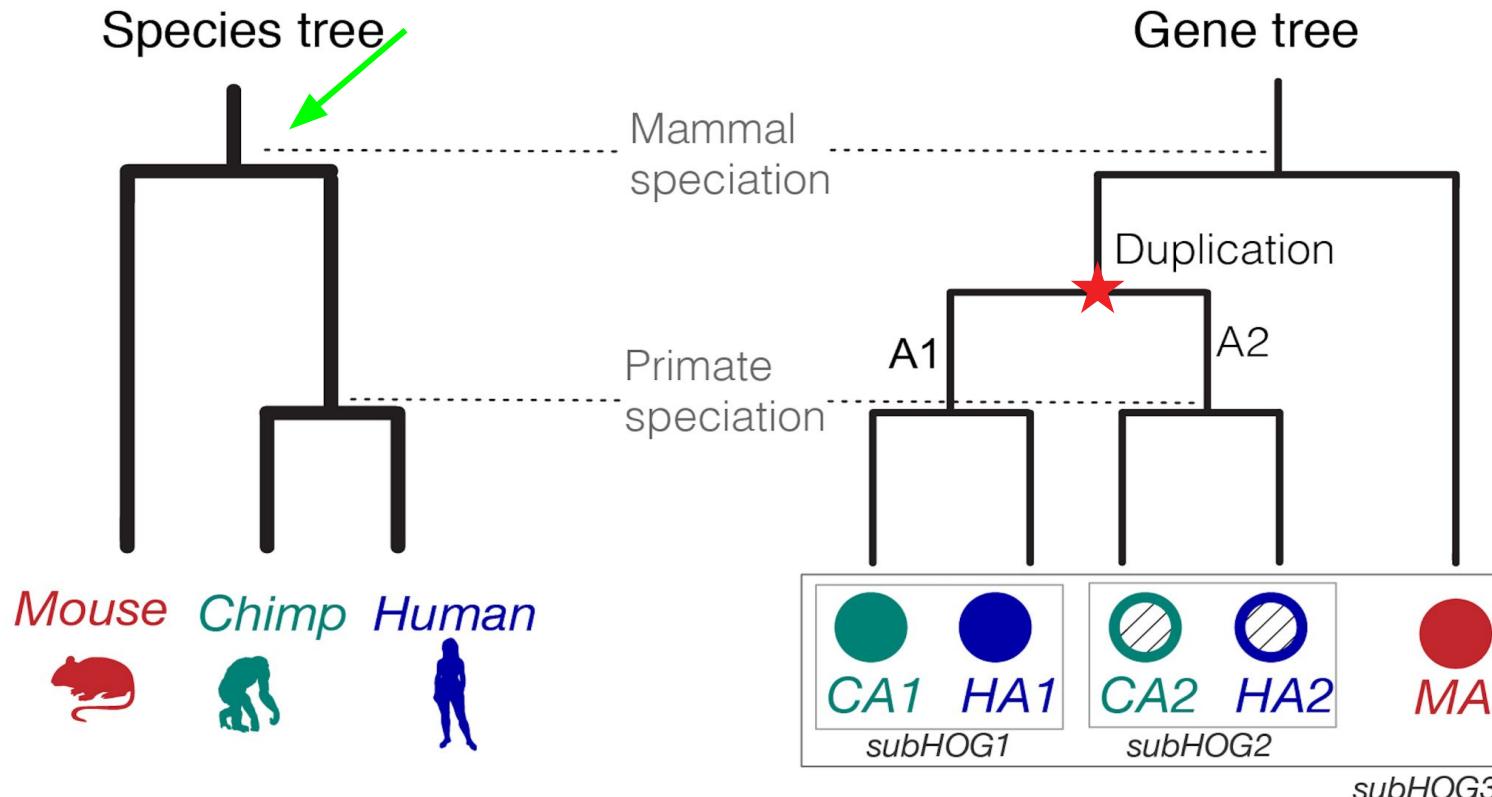
Gene tree



# HOG

Hierarchical Orthologous Group

- Genes descended from a common ancestral gene at a specific taxonomic level

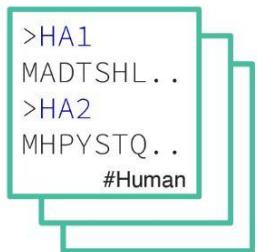


# Outline

- Introduction
- Orthology inference with FastOMA
- Phylogeny inference from raw sequencing reads with read2tree

# FastOMA, our new tool

Input Proteomes

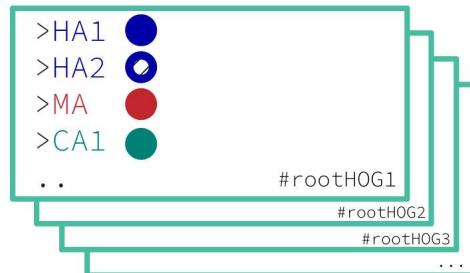


OMAmer

Mapping sequences  
on OMA gene families  
based on k-mers

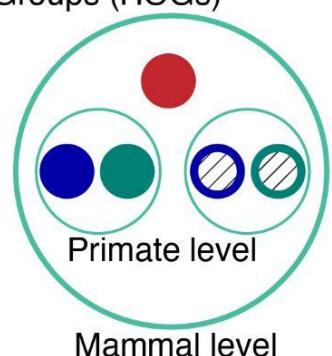


Root HOGs (Gene families)

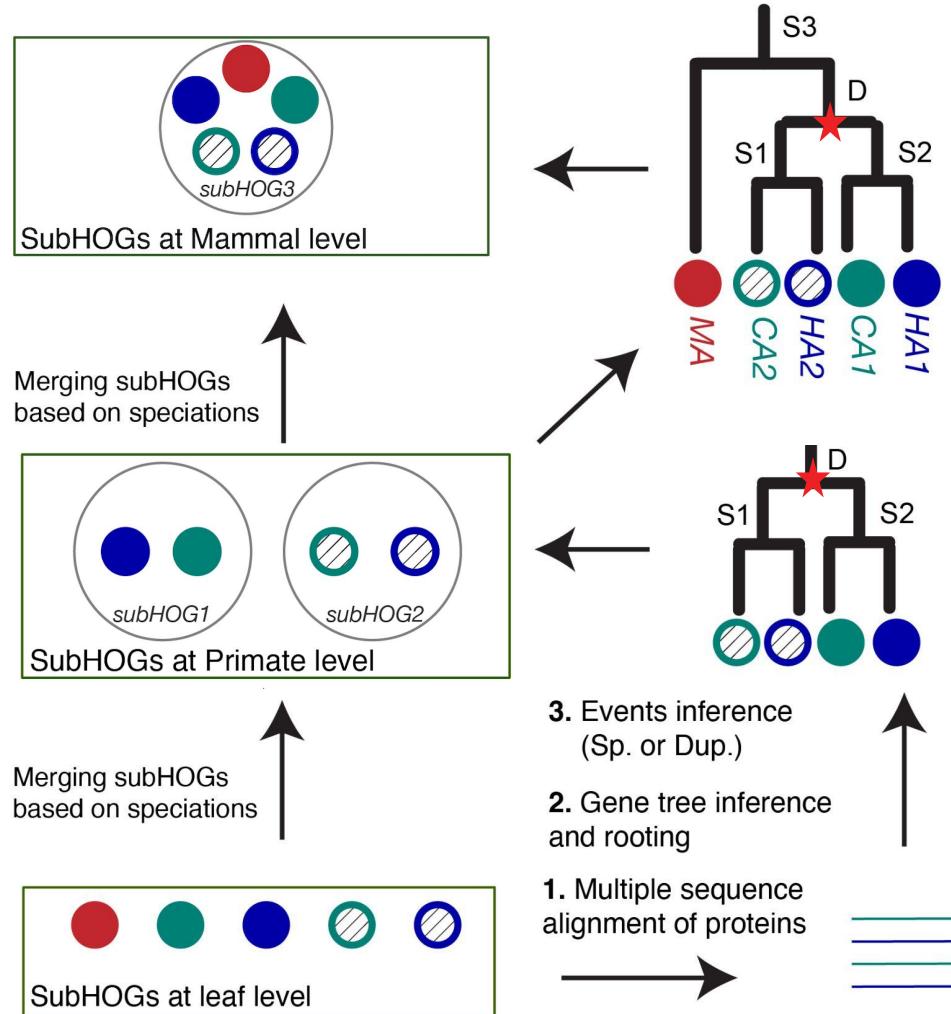
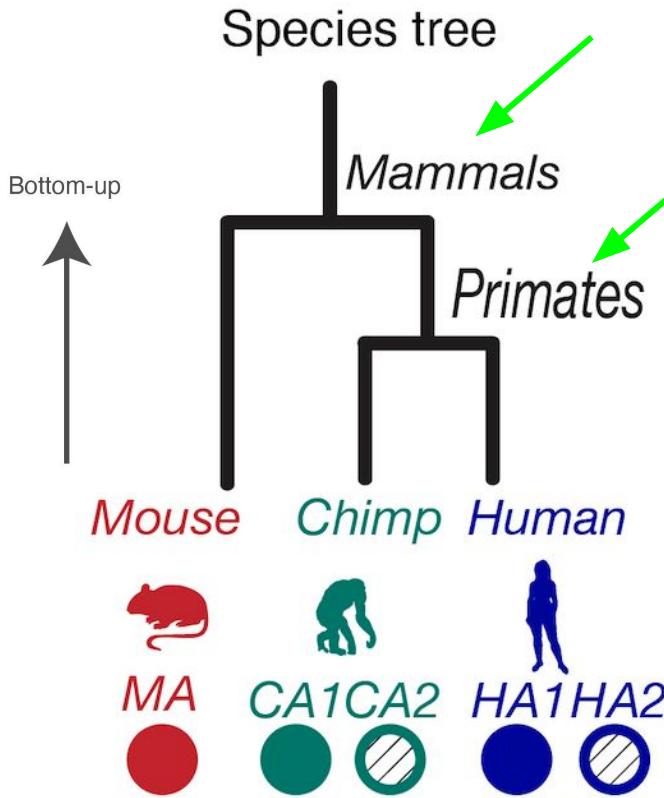


SubHOG/event  
inference  
(in parallel)

Hierarchical Orthologous Groups (HOGs)



# HOG inference

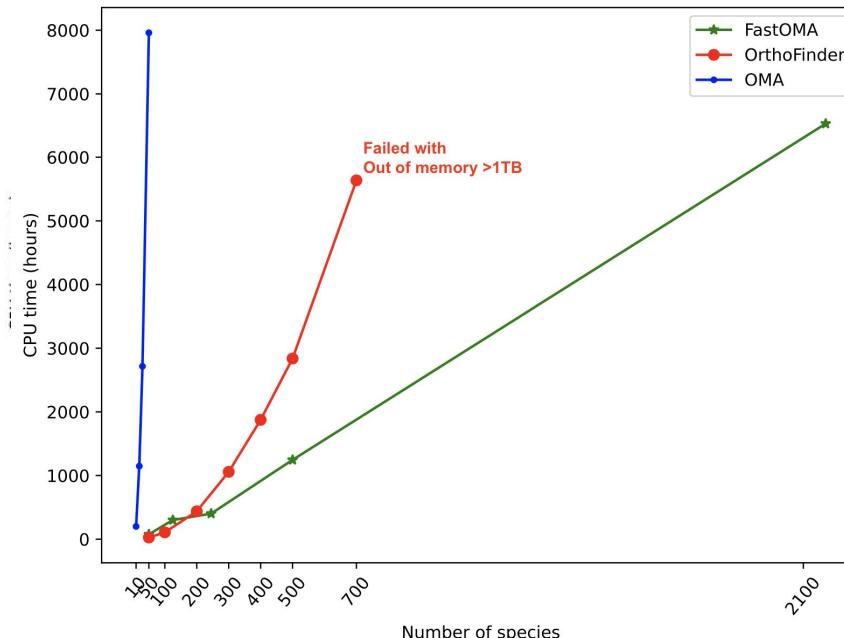


# Orthology inference with FastOMA

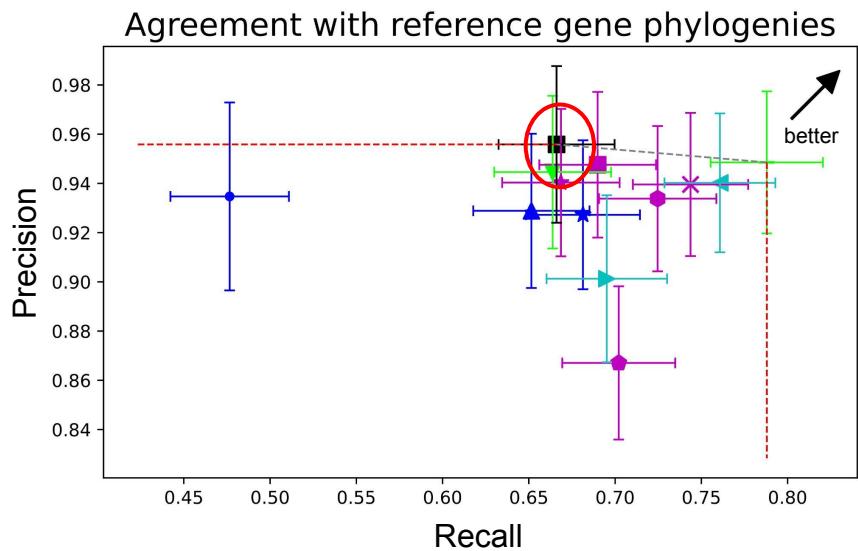
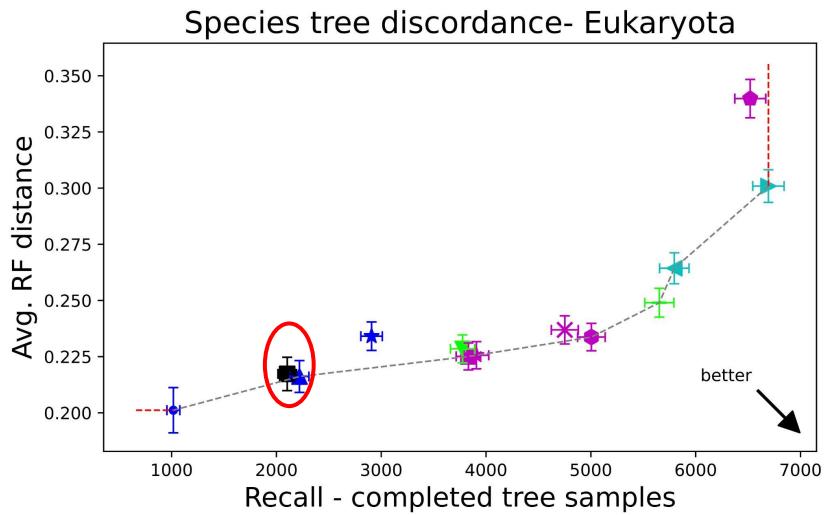


[github.com/DessimozLab/  
FastOMA](https://github.com/DessimozLab/FastOMA)

- UniProt reference proteomes
- 2180 eukaryotic species
- in a single day using 300 CPUs



# Quest for orthologs benchmarking



- ▲ OMA\_Pairs
- OMA\_Groups
- ★ OMA\_GETHOGs
- ✗ Domainoid+

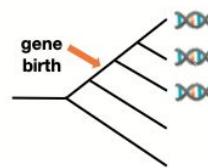
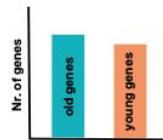
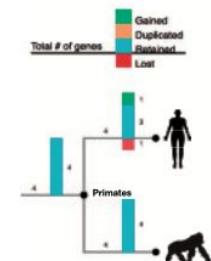
- InParanoid\_Xenfix
- ◆ OrthoMCL
- Ortholnspector 3

- ★ sonicparanoid
- ◀ PANTHER
- ◆ Ensembl\_Compara

- ▼ Hieranoid\_2
- Orthofinder
- FastOMA

State-of-the-art methods

## PyHAM



Finding  
taxonomically  
restricted  
genes

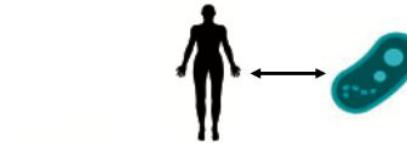
Phylo-  
stratigraphy

Elucidating  
gene gains  
and losses



Prediction  
of gene  
function

HOGprop

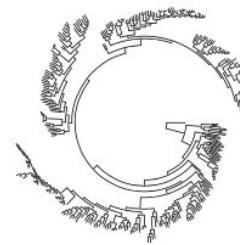


Verification  
of function  
conservation

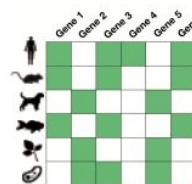
ORTHOLOGY  
APPLICATIONS

Phylo-  
genomics

Finding the  
best model  
systems



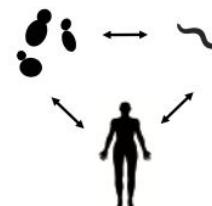
Phylo-  
genetic  
profiling



HOGprof



A. Nicheperovich, ... S. Majidian. OMAMO: orthology-based alternative model organism selection. *Bioinformatics*, 2022.



## Take-home messages (1/2)

- New genome assemblies are becoming available.
- Orthology is at the heart of studying evolution.
- Scalable methods are needed to handle huge amount of data.
- FastOMA is a scalable and reliable method for orthology inference.

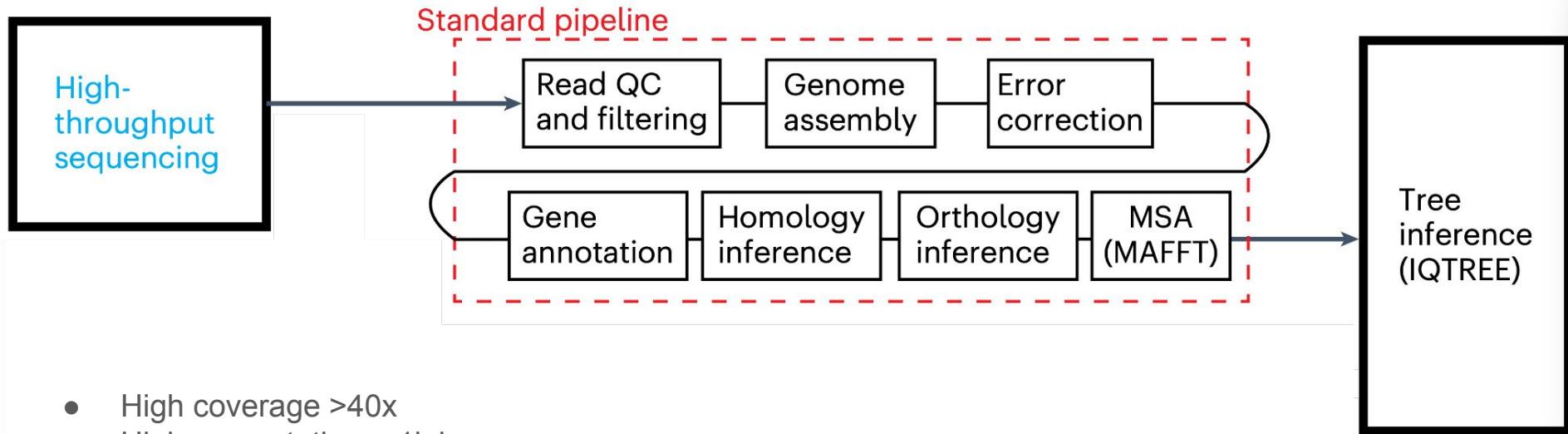
The basis of this approach is proteomes.

**What if there is no high quality assembly?**

# Outline

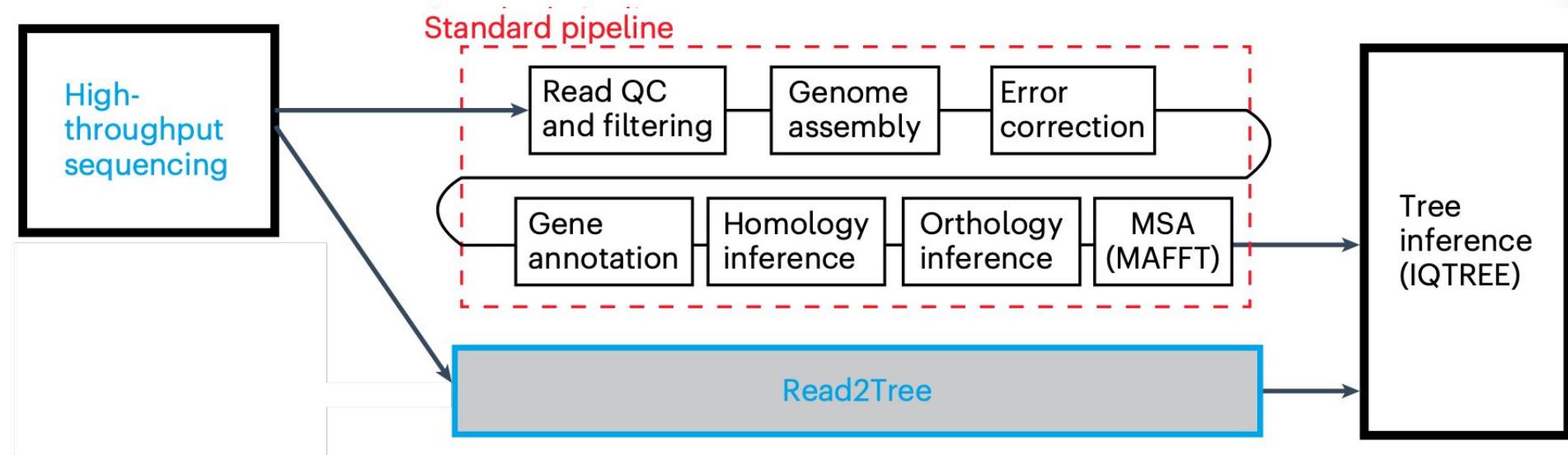
- Introduction
- Orthology inference with FastOMA
- Phylogeny inference from raw sequencing reads with read2tree

# Phylogenetic tree: current status



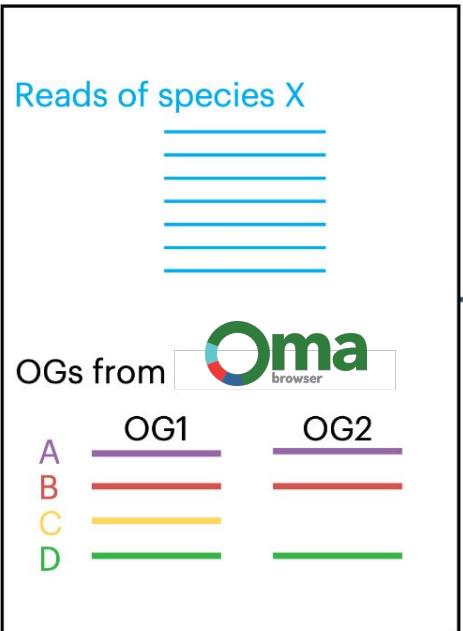
- High coverage >40x
- High computation > 1k hours
- High cost
- Orthogonal techs
- Expertise to close gaps, spots errors/contamination, annotation

# Read2Tree: A faster approach



# Read2Tree: How it works?

Input



Read2Tree

1. Align OGs (MAFFT)



2. Map reads



3. Build consensus



4. Select best

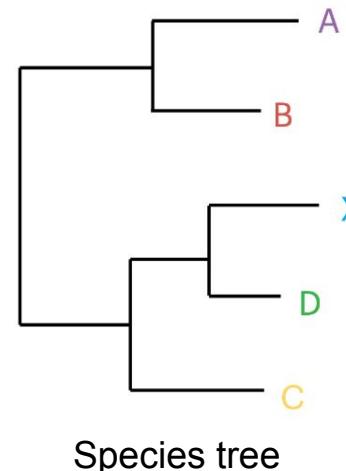


5. Add to align and concatinate



Output

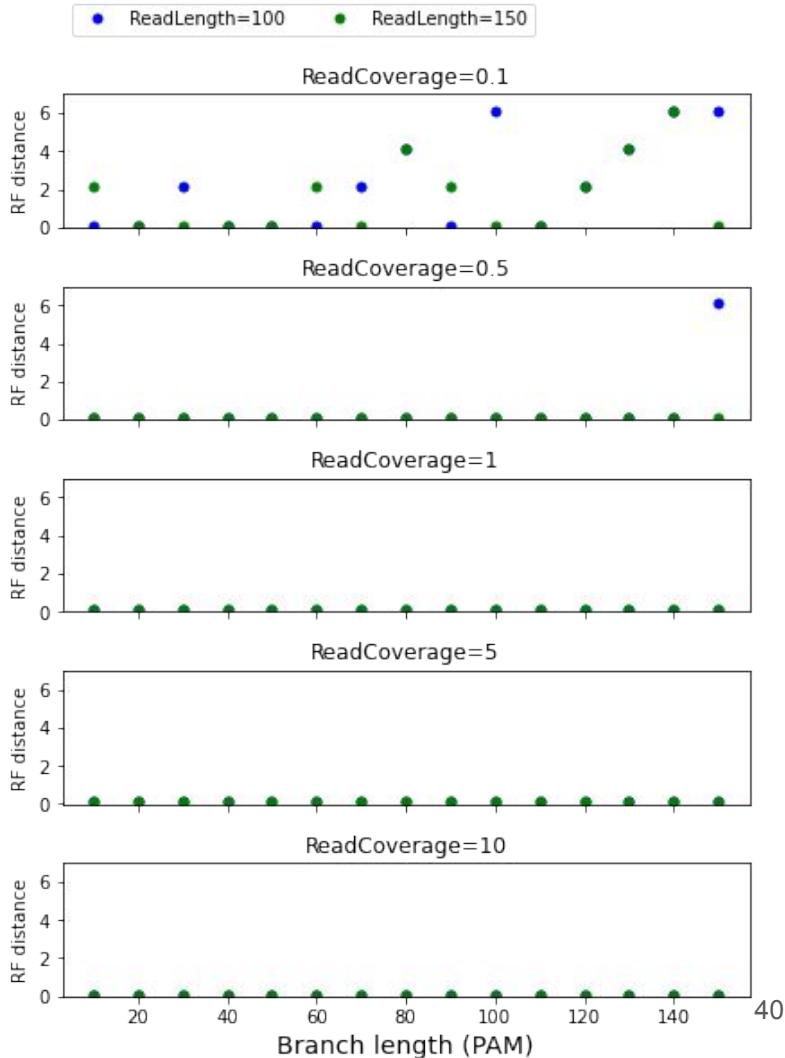
Infer tree (IQTREE)



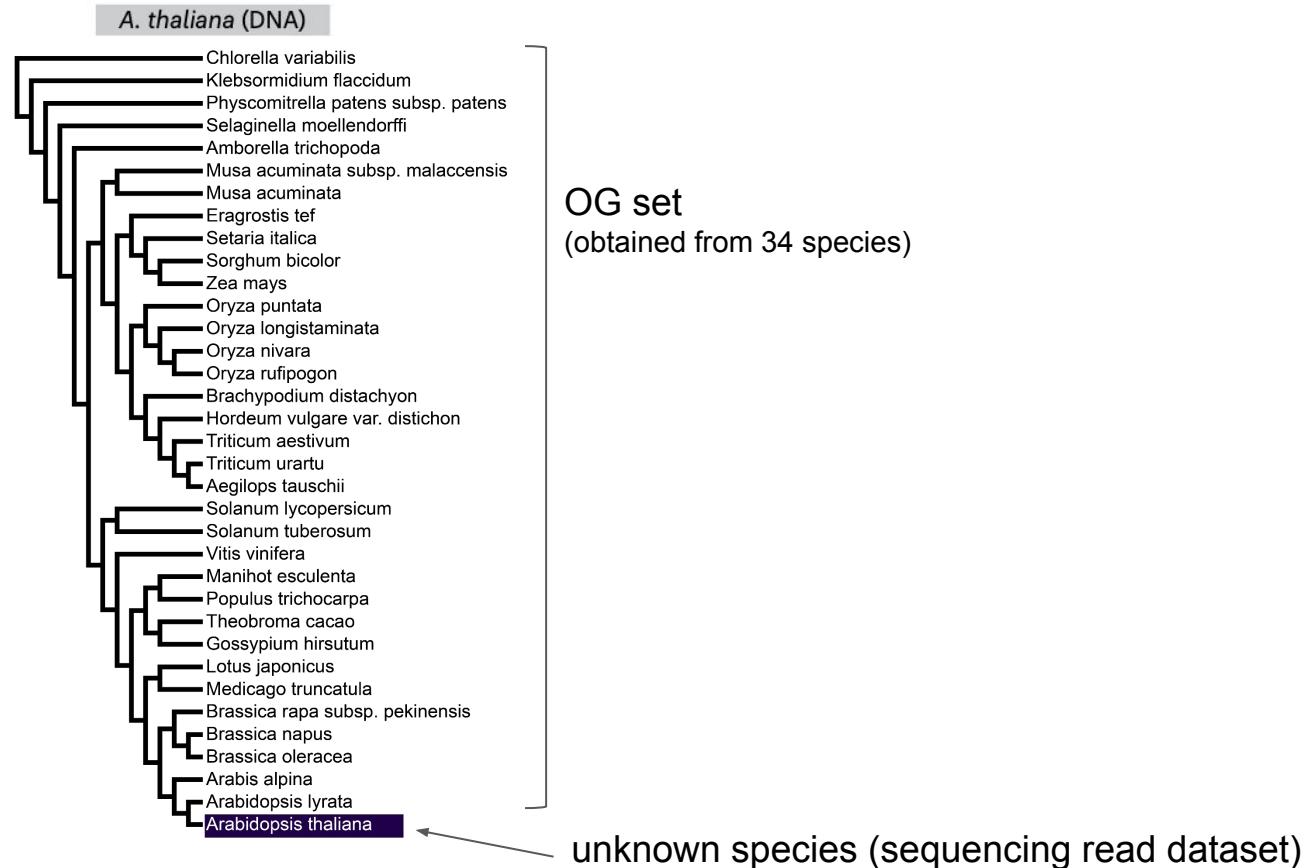
OG: Orthologous Groups (marker genes)

# Read2Tree: inferring simulated species tree

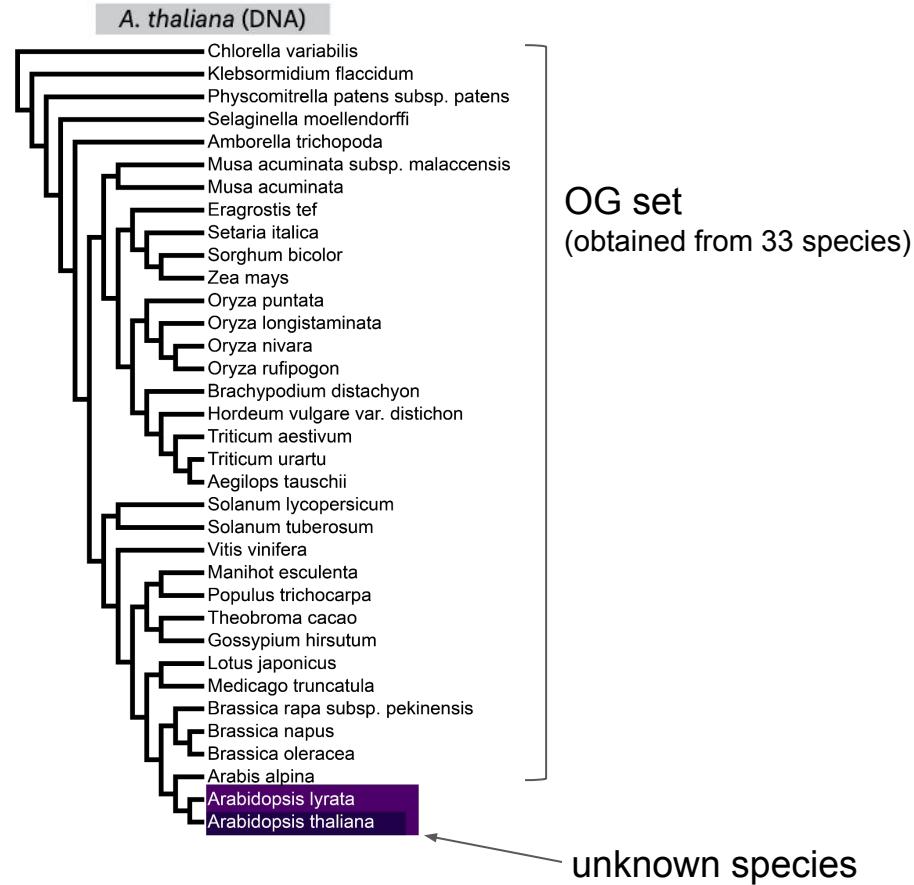
- 15 genomes containing 100 genes
- Simulated illumina DNA reads
  - Coverage values 0.1-10
- OG set based on 14 species
- Inferring the tree for one species
- Comparing inferred tree with true one
  - in terms of Robinson–Foulds (RF) distance
- 15 different branch length values



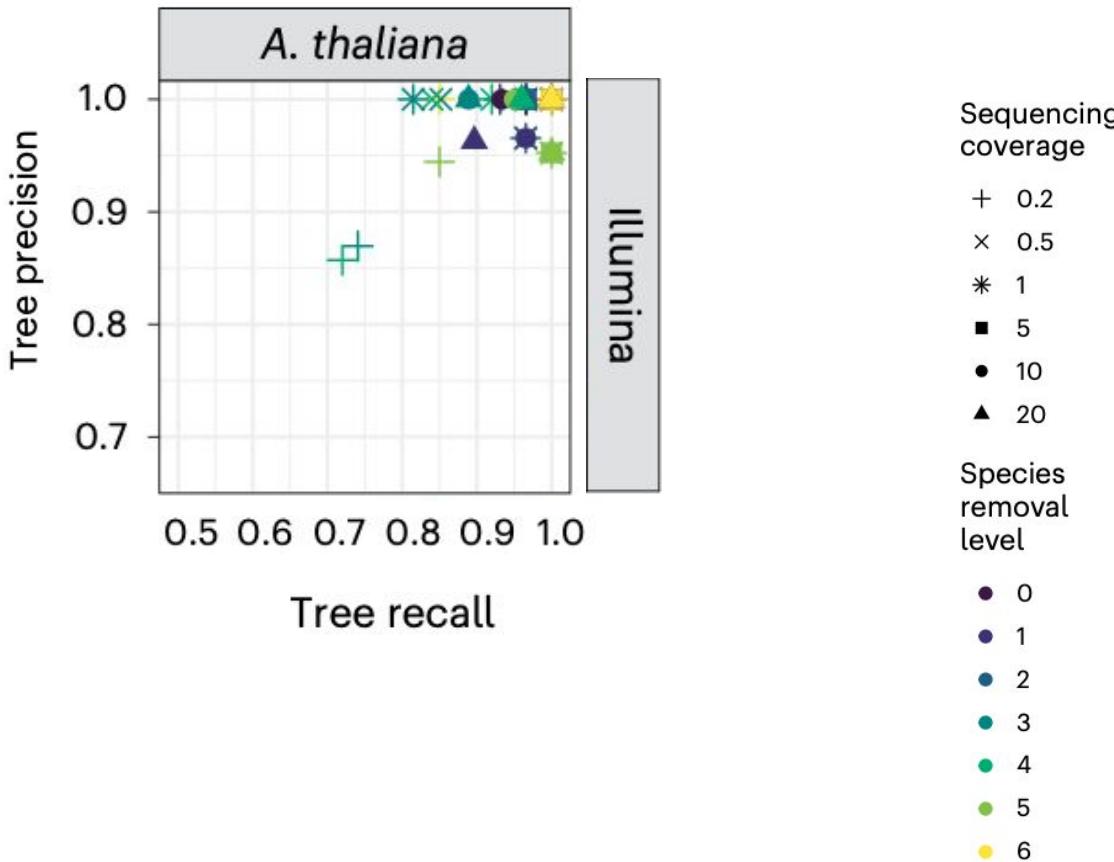
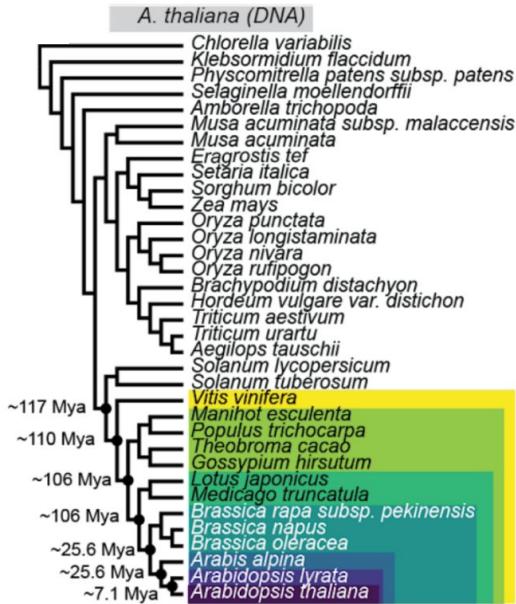
# Read2Tree: Benchmarking



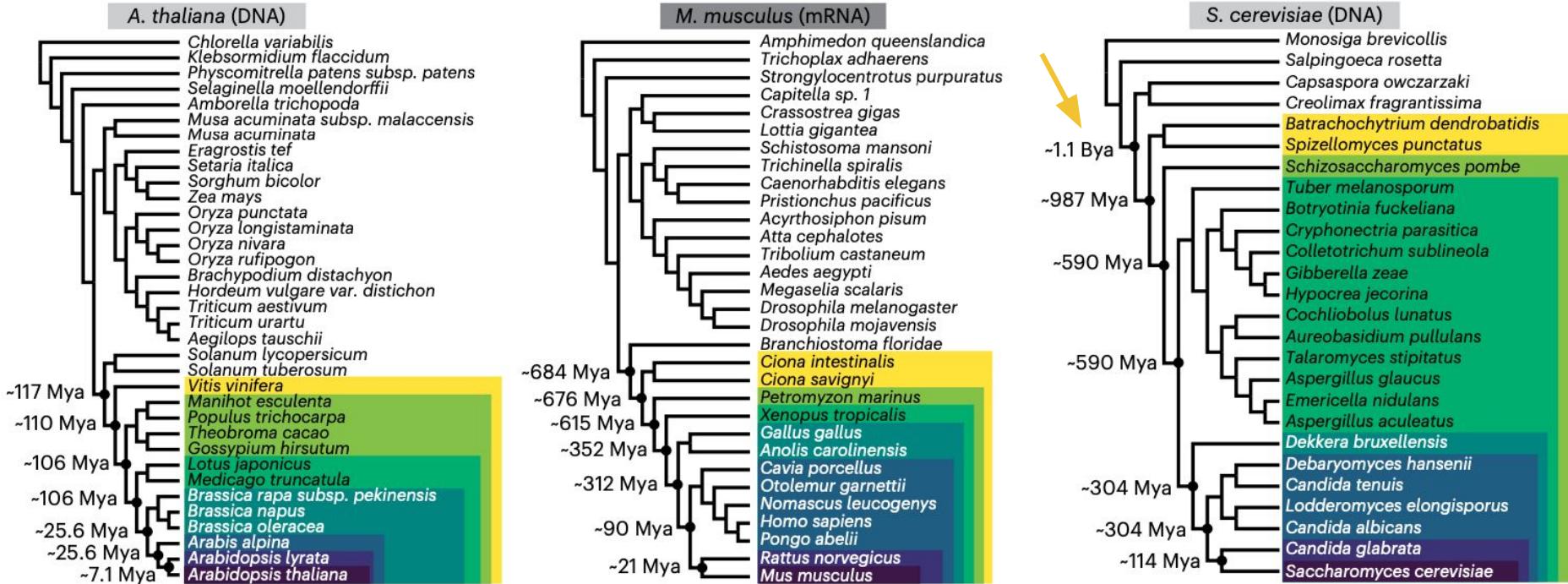
# Read2Tree: Benchmarking



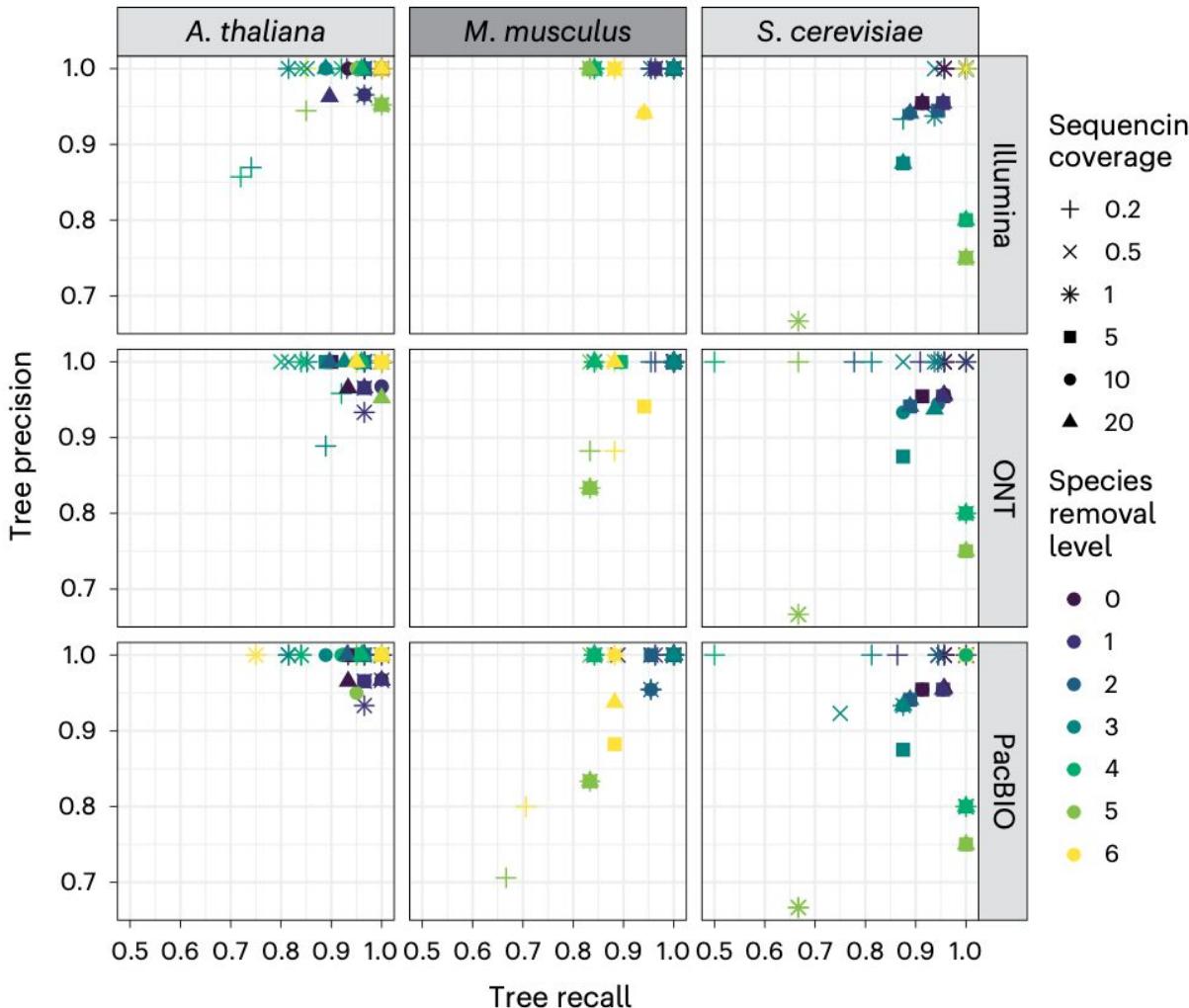
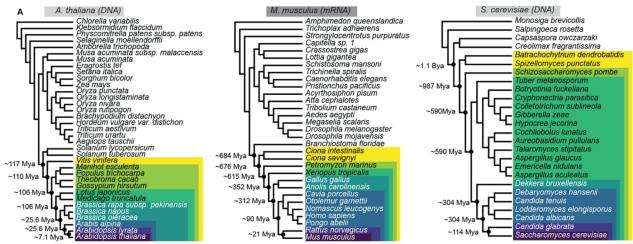
# Reconstructed tree accuracy



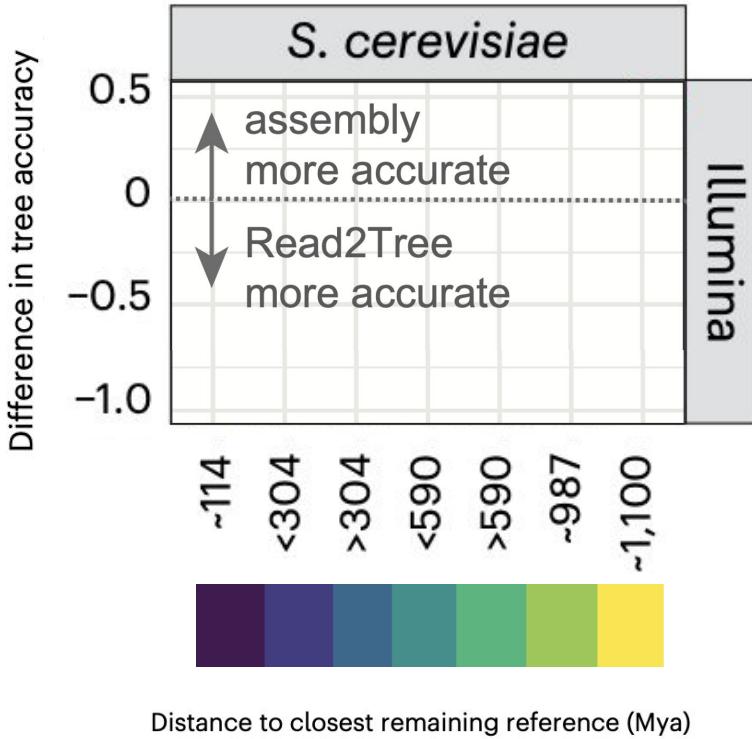
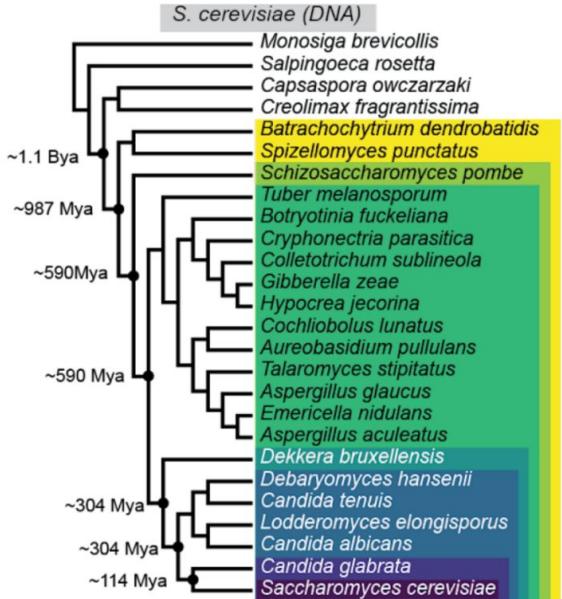
# Read2Tree: Benchmarking



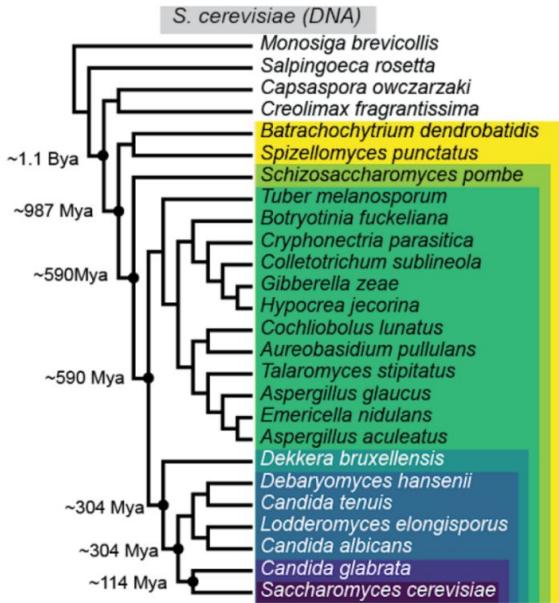
# Reconstructed tree accuracy



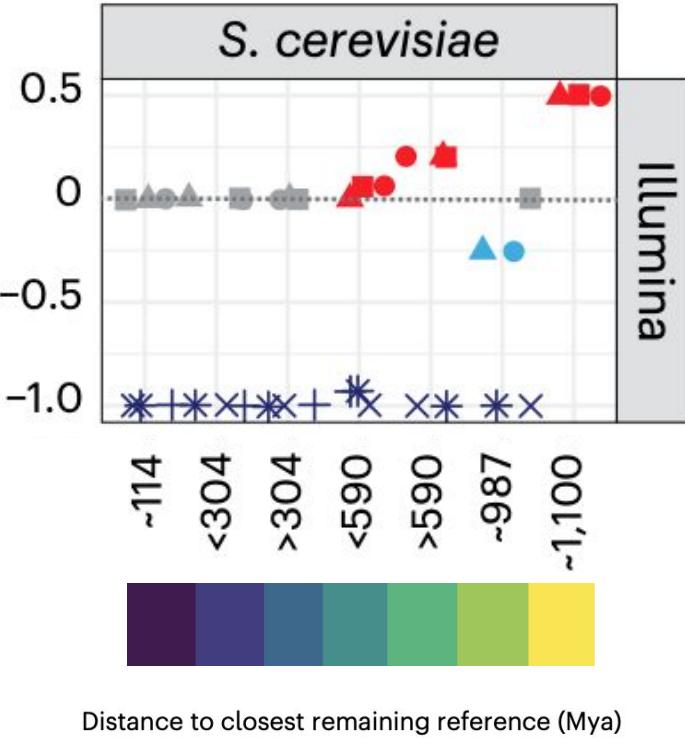
# Read2Tree more accurate than assembly for tree inference



# Read2Tree more accurate than assembly for tree inference



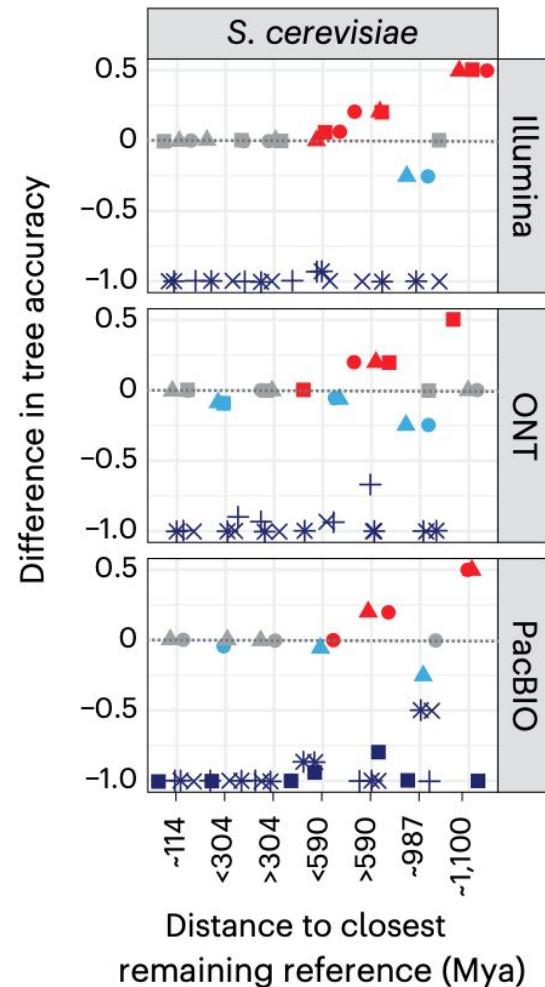
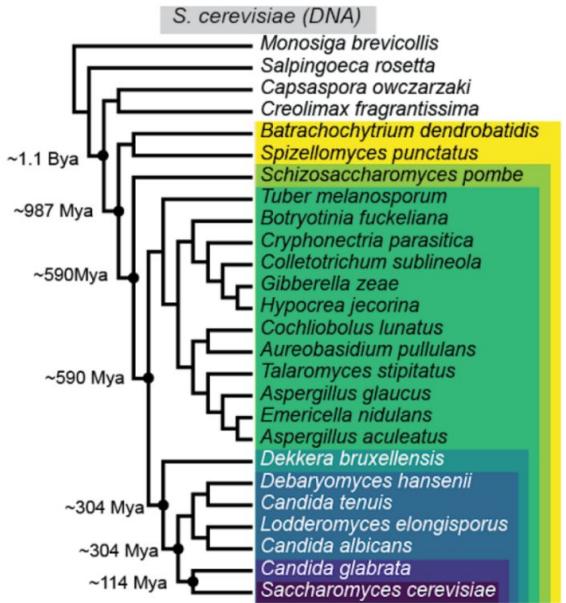
Difference in tree accuracy



- Assembly better
- Read2Tree better
- Equal
- Only Read2Tree applicable  
(due to low coverage)

Sequencing +0.2 \*1 •10 ▽60  
coverage ×0.5 ■5 ▲20

# Read2Tree more accurate than assembly for tree inference



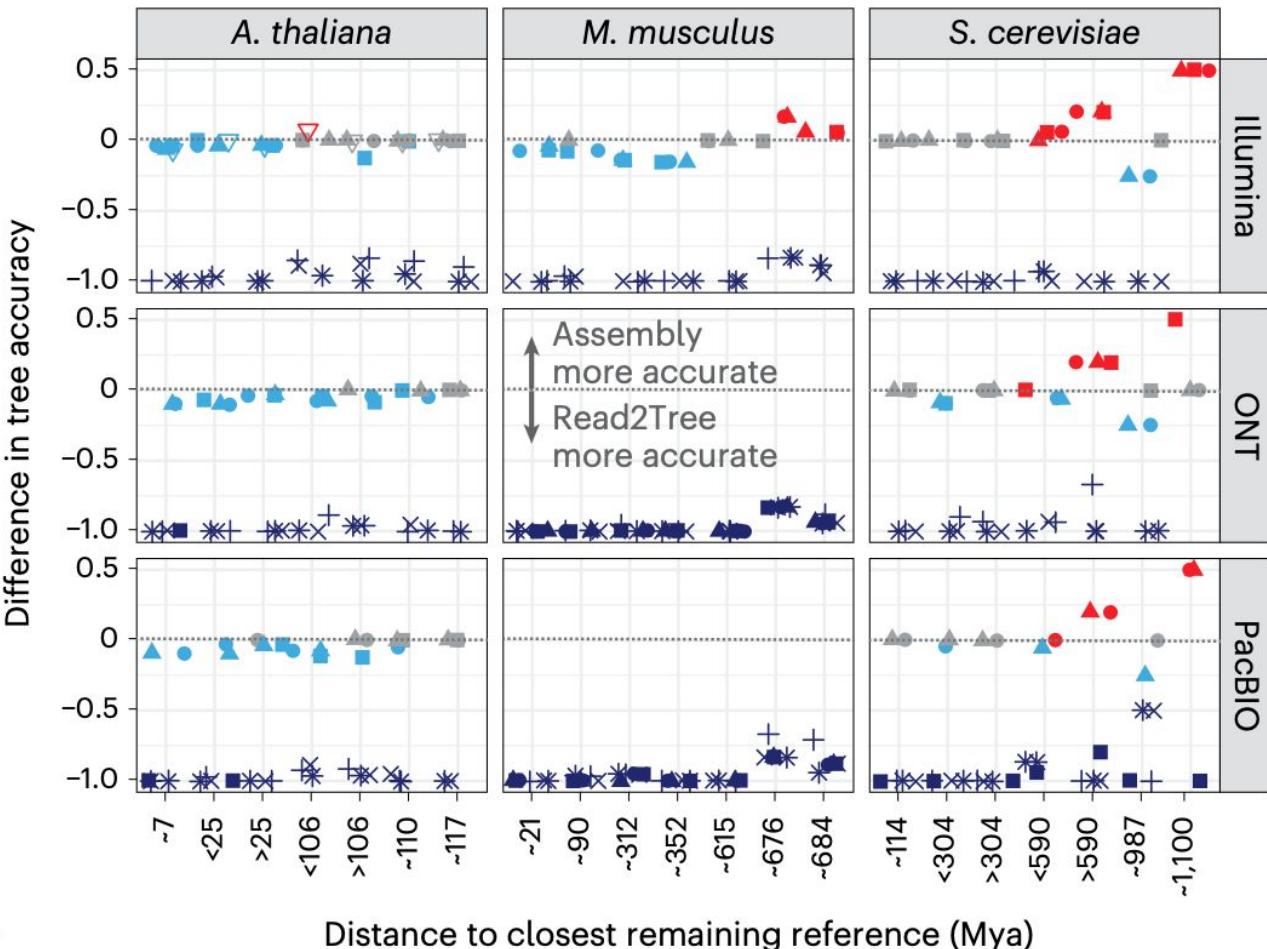
- Assembly better
- Read2Tree better
- Equal
- Only Read2Tree applicable  
(due to low coverage)

Sequencing +0.2 \*1 •10 ▽60  
coverage ×0.5 ■5 ▲20

# Read2Tree more accurate than assembly for tree inference

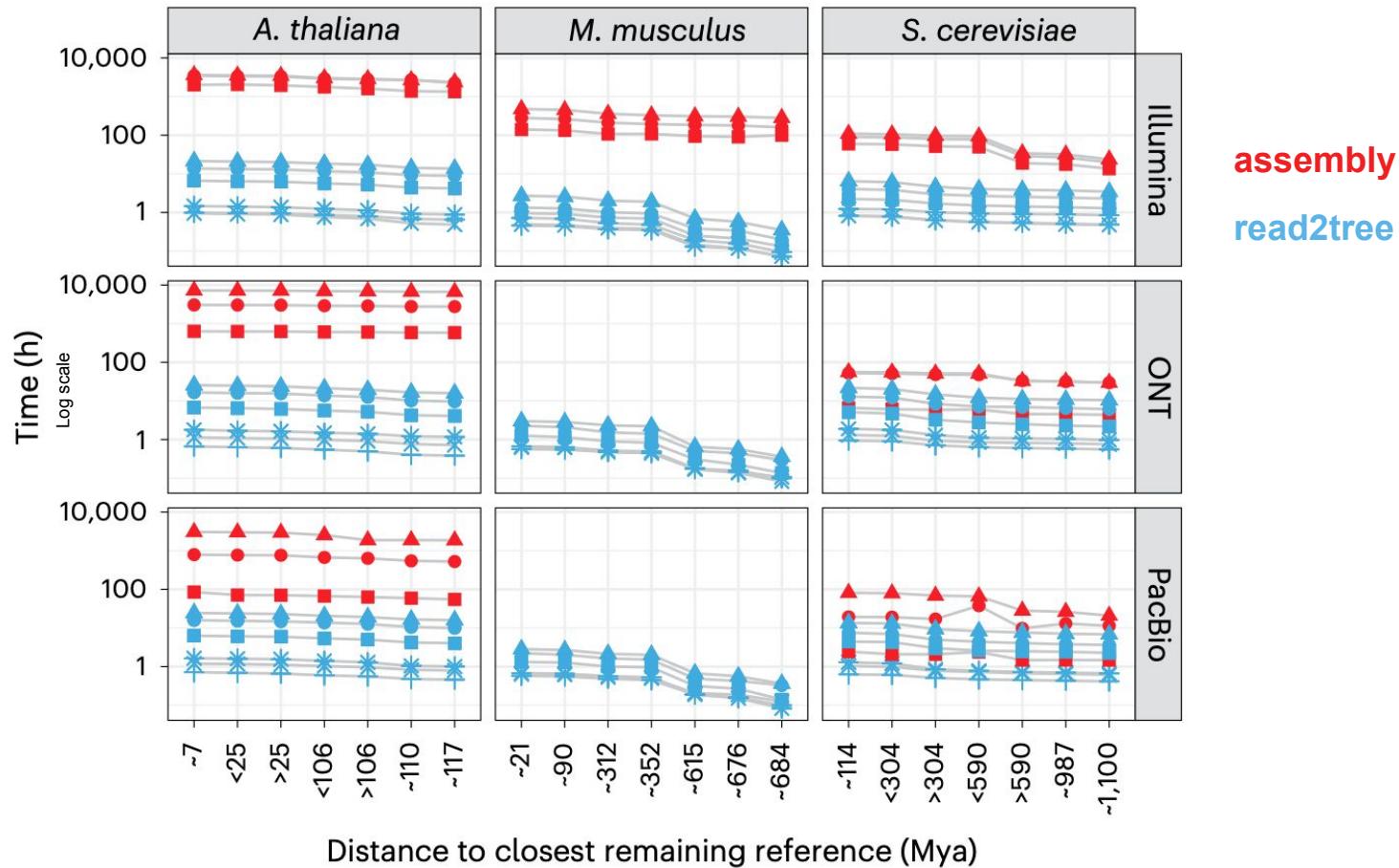
Sequencing coverage +0.2 \*1 •10 ▽60  
x0.5 ■5 ▲20

- Assembly better
- Read2Tree better
- Equal
- Only Read2Tree applicable

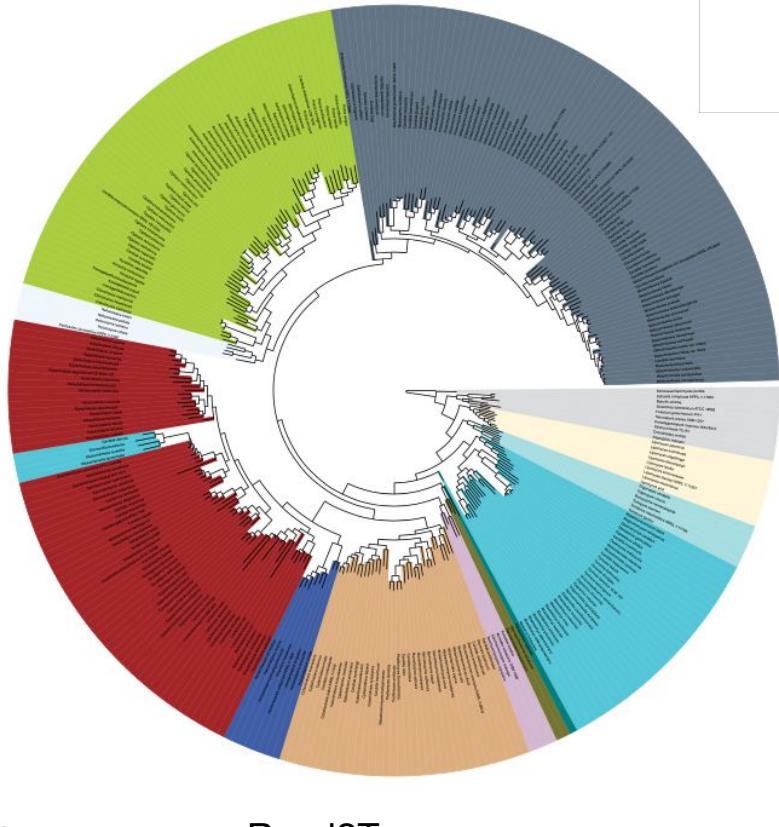


MASH results are not shown.

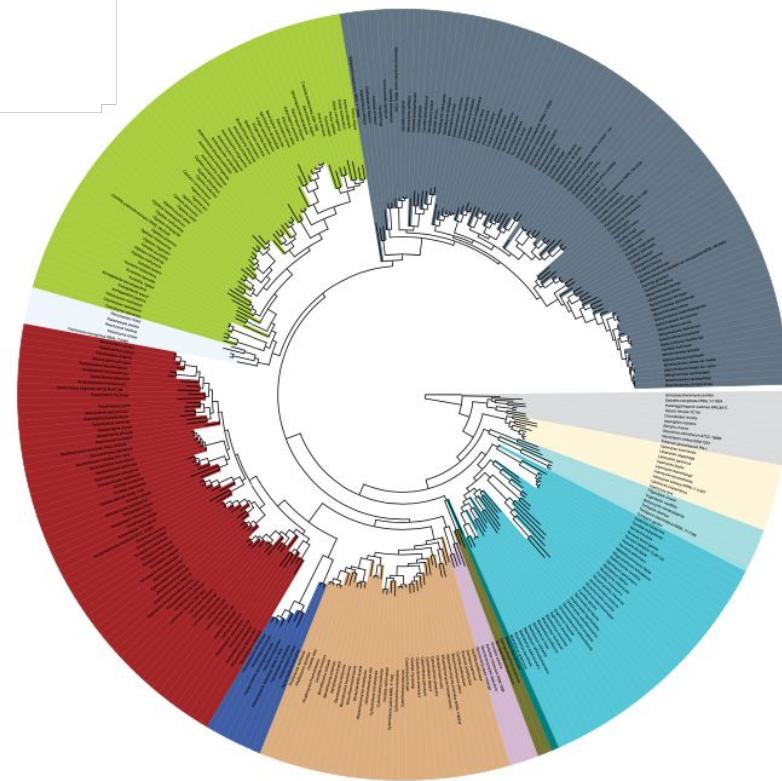
# Read2Tree is orders of magnitude faster



# Read2Tree reproduces state-of-the-art yeast phylogeny



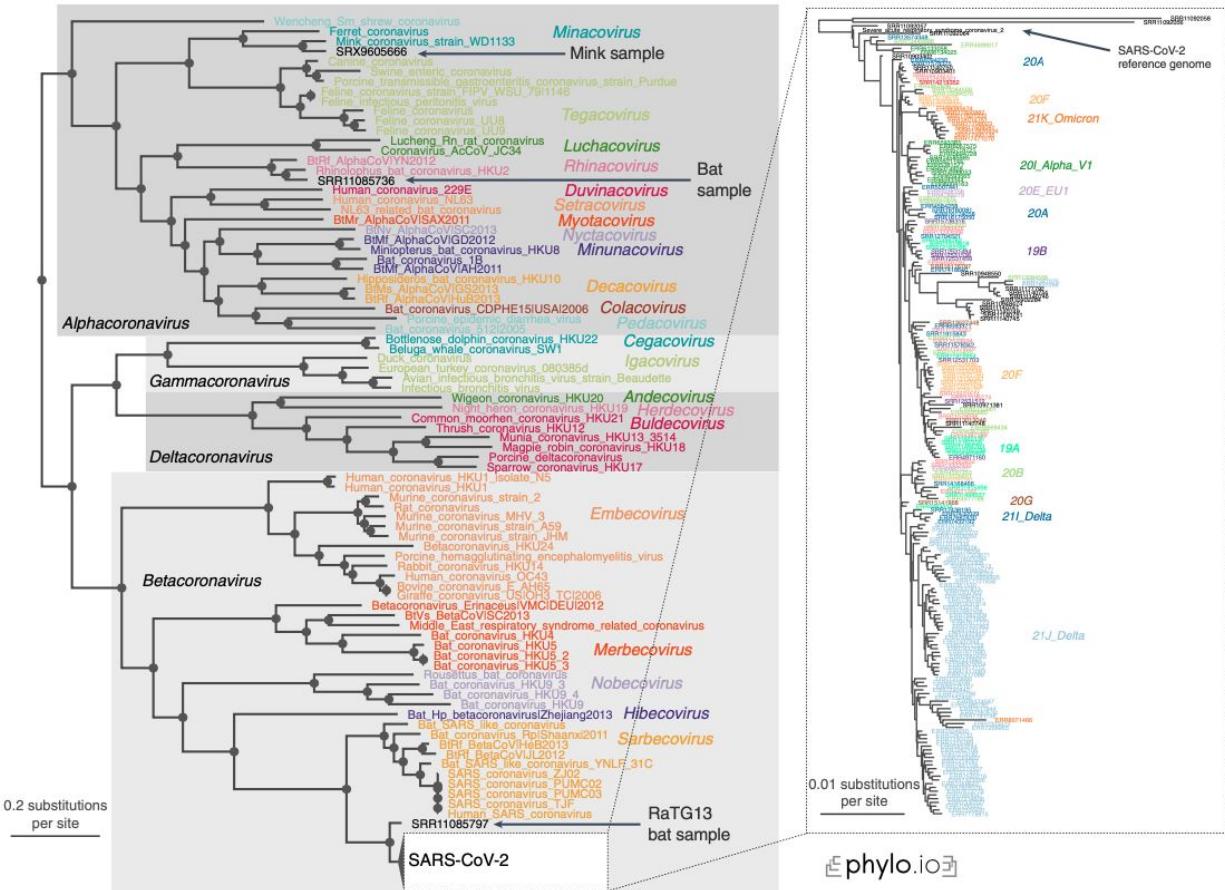
Read2Tree



Shen et al. Cell 2018

# Read2Tree: Coronavirus

- Rapid identification and placement of viruses
- Inferring tree using raw sequencing data
  - Data download takes most time!
- Main genera
  - Alpha-, Beta-, Gamma- and Delta- subgenera



# Read2Tree: Summary

- Rapid phylogenetic tree reconstruction
- Scales from low coverage to large numbers of samples
- Eases comparative genomics from small to large labs
  - Potentially removes biases along the way
  - Low coverage



[github.com/DessimozLab/  
read2tree](https://github.com/DessimozLab/read2tree)

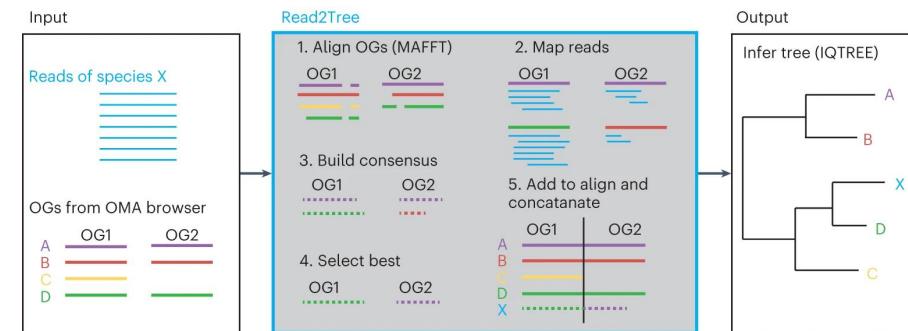
- Future directions:
  - Metagenomics samples (multiple samples)
  - Single cell genomics applications

## Inference of phylogenetic trees directly from raw sequencing reads using Read2Tree

Received: 18 April 2022

Accepted: 16 March 2023

David Dylus <sup>1,2\*</sup>, Adrian Altenhoff <sup>2,3</sup>, Sina Majidian <sup>1,2</sup>,  
Fritz J. Sedlazeck <sup>4,5</sup> & Christophe Dessimoz <sup>1,2,6,7</sup>



# Thank you !



UNIL | Université de Lausanne



Swiss Institute of  
Bioinformatics

Zürcher Hochschule  
für Angewandte Wissenschaften

