



Swiss Institute of  
Bioinformatics



# Orthology inference at scale with FastOMA



Sina Majidian

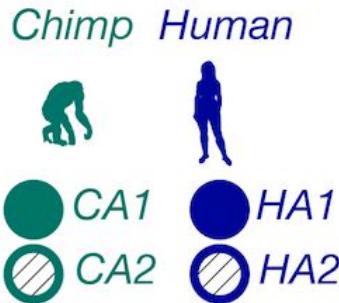
27 July 2023



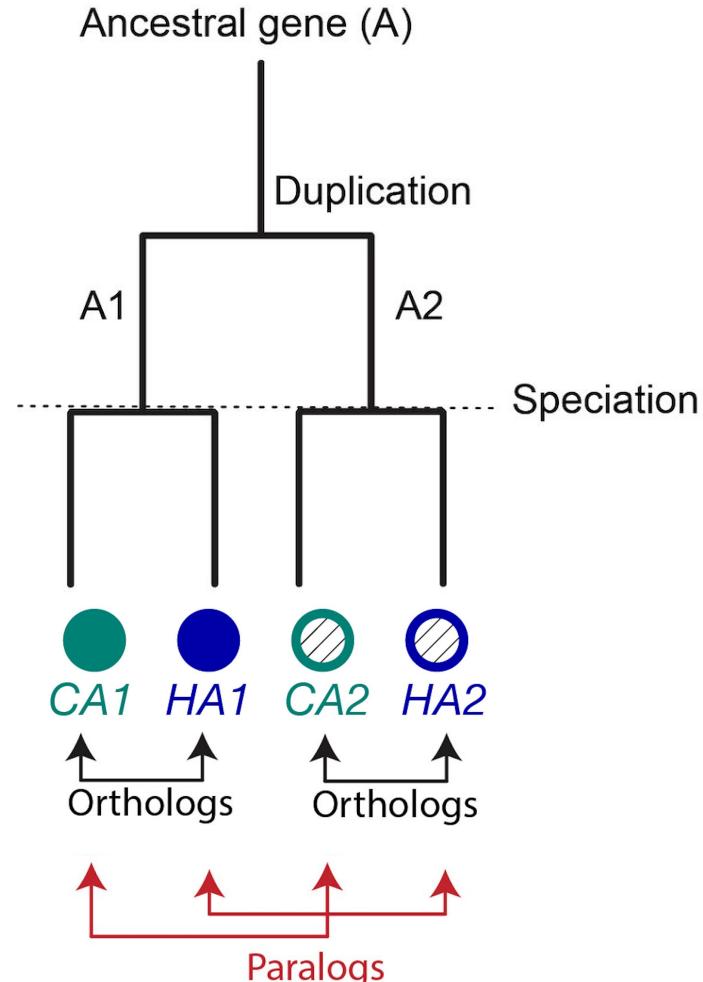
@DessimozLab  
@SinaMajidian

# Introduction

- Orthology vs paralogy



- Informative for
  - studying gene evolution
  - phylogeny inference, ...
- Found using all-vs-all alignment
  - Not scalable**

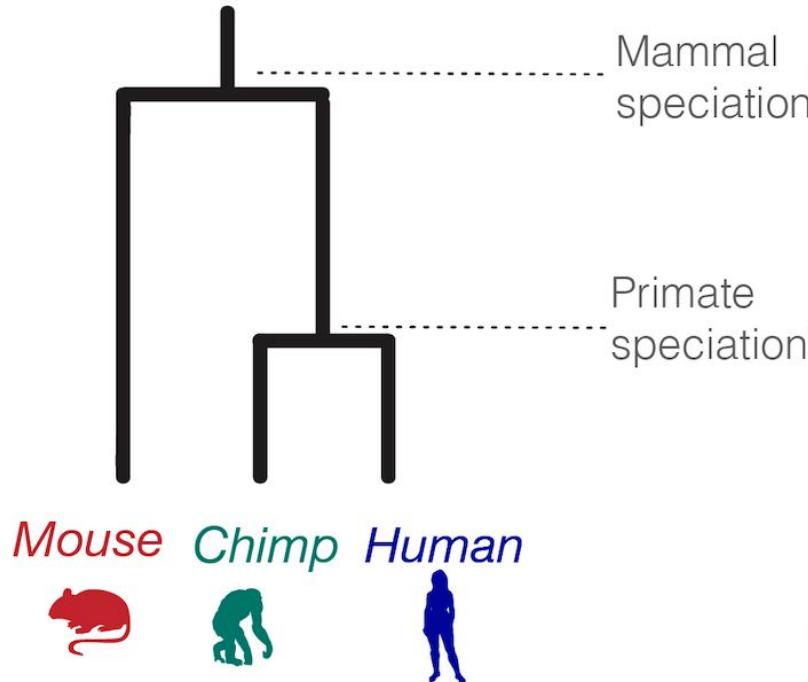


# HOG

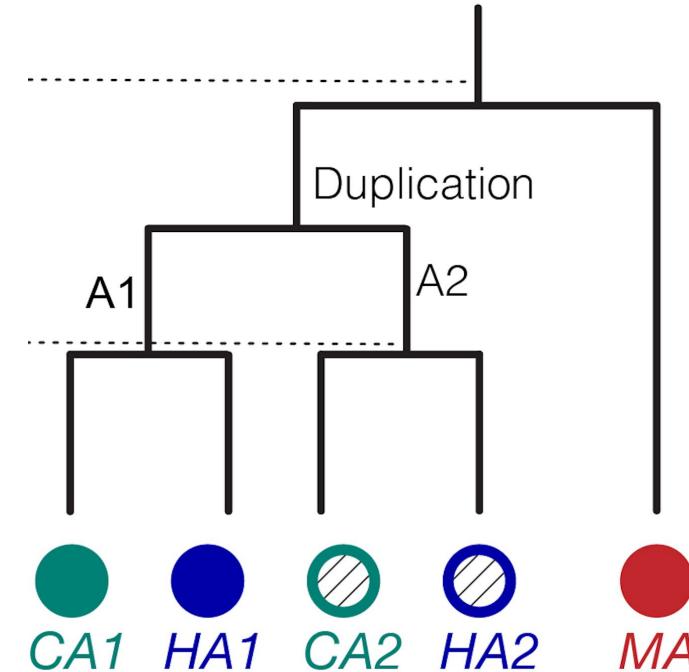
Hierarchical Orthologous Group

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

Species tree



Gene tree

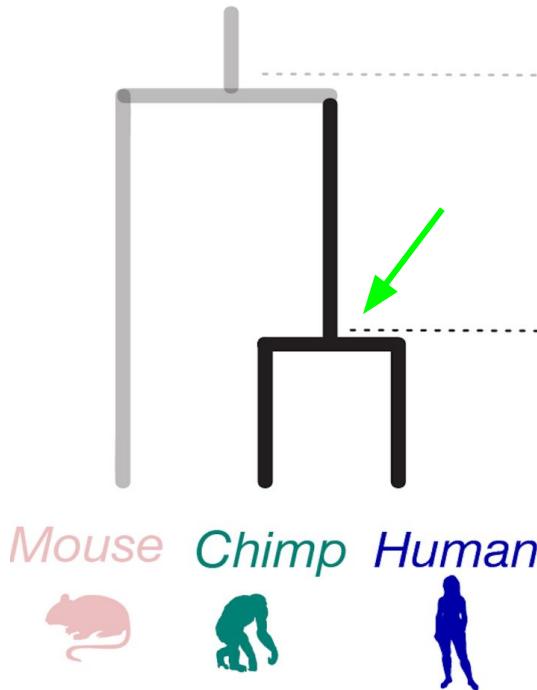


# HOG

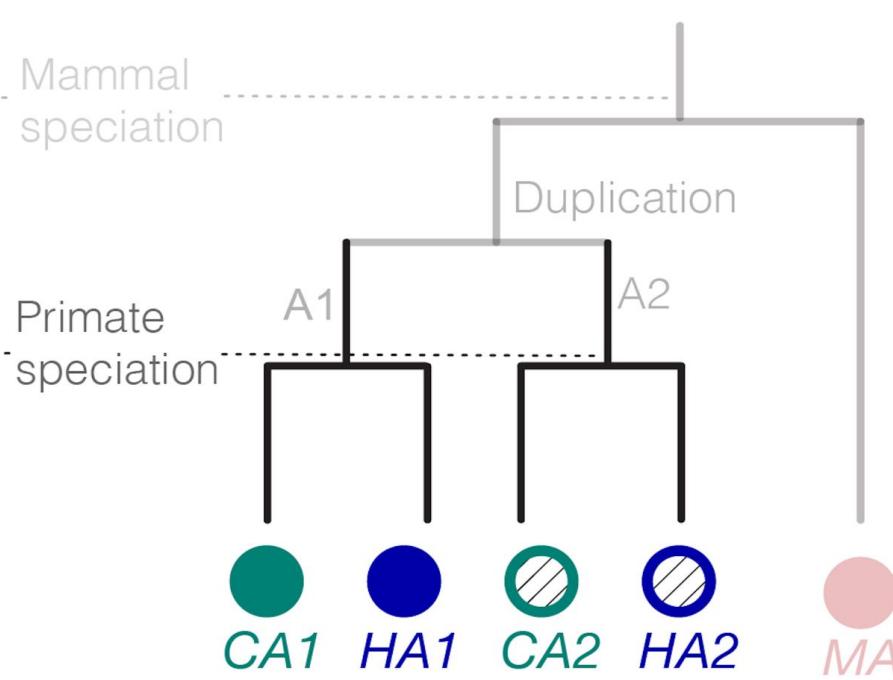
Hierarchical Orthologous Group

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

Species tree



Gene tree

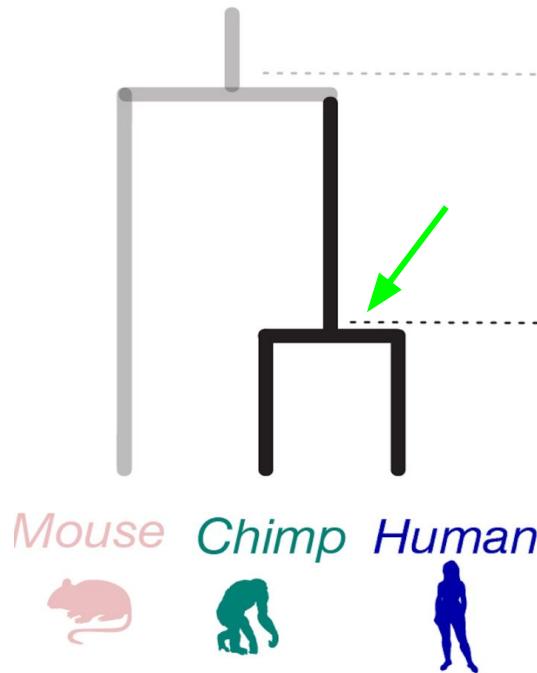


# HOG

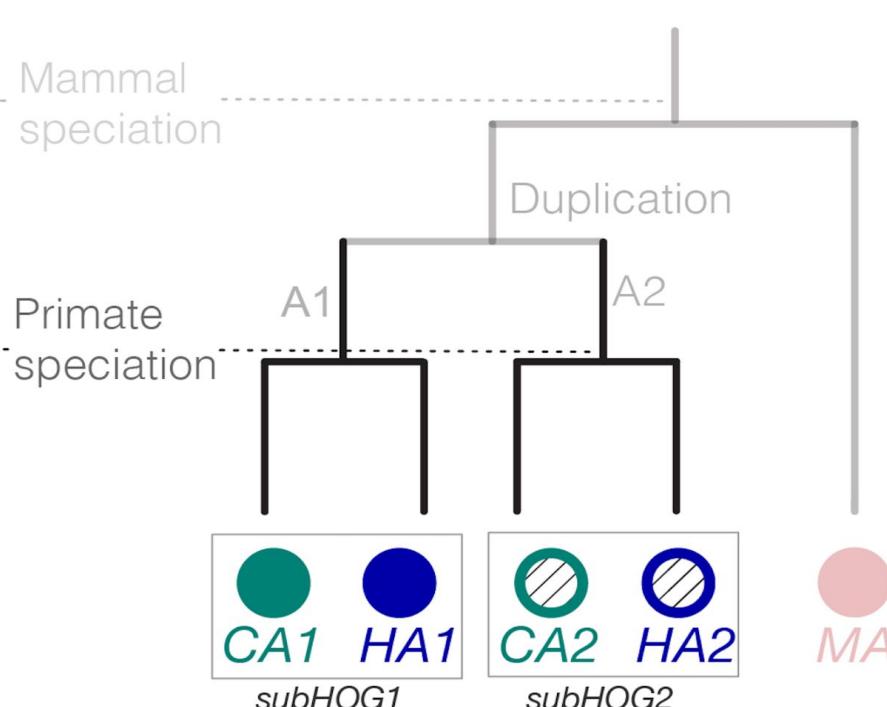
Hierarchical Orthologous Group

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

Species tree



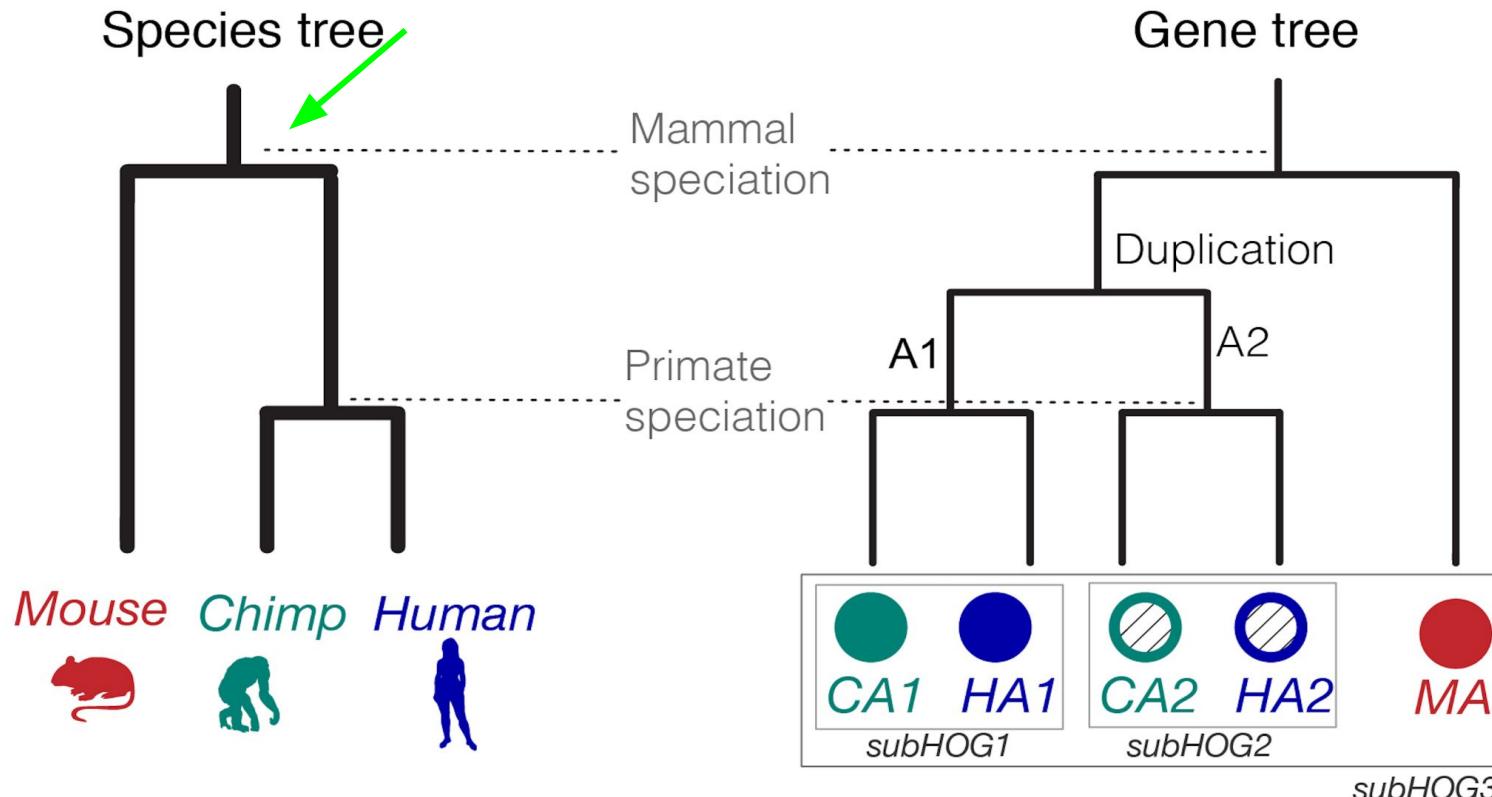
Gene tree



# HOG

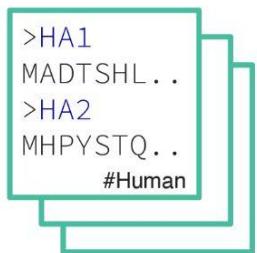
Hierarchical Orthologous Group

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



# FastOMA, our new tool

Input Proteomes

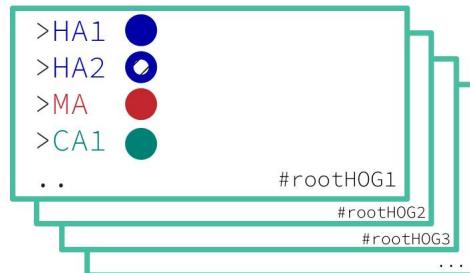


OMAmer

Mapping sequences  
on OMA gene families  
based on k-mers

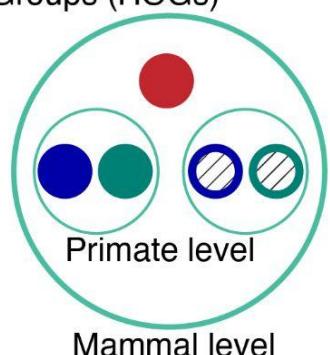


Root HOGs (Gene families)

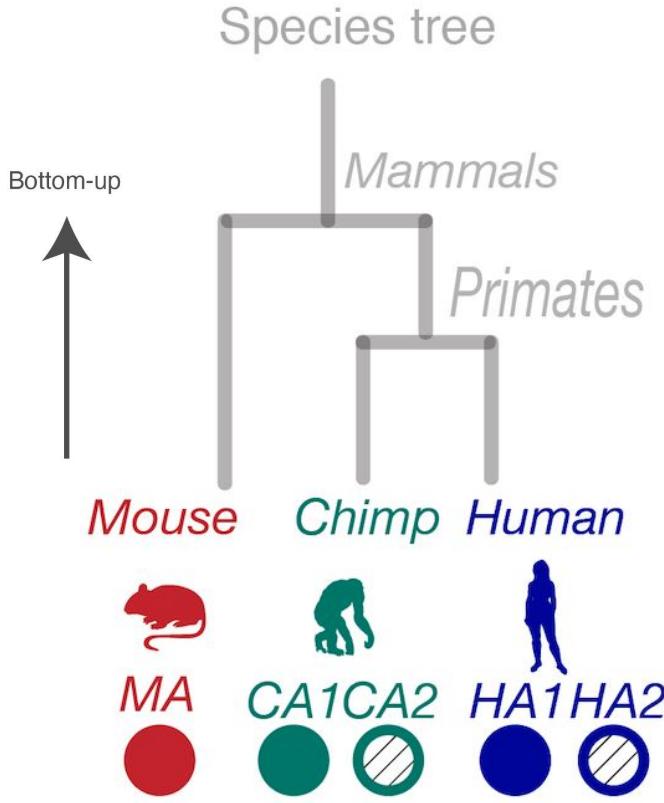


SubHOG/event  
inference  
(in parallel)

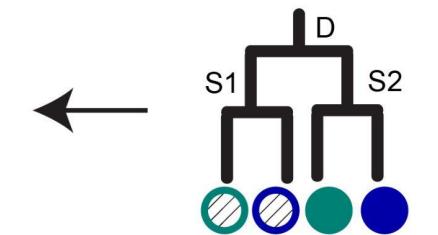
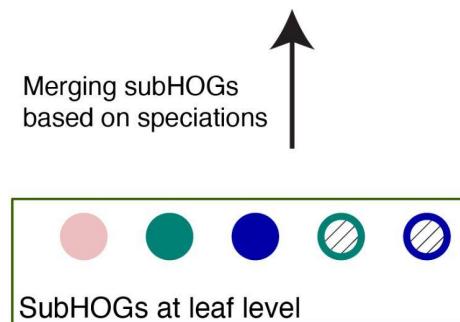
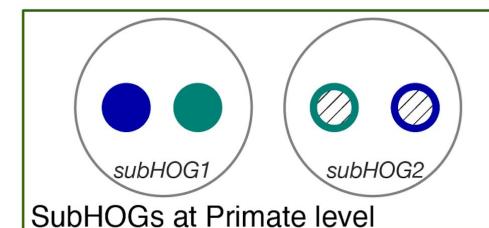
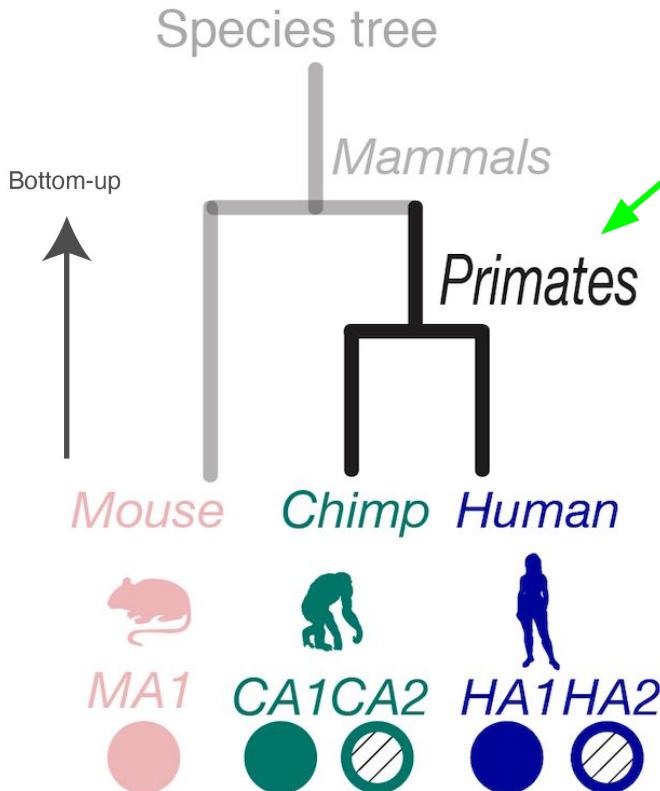
Hierarchical Orthologous Groups (HOGs)



# HOG inference

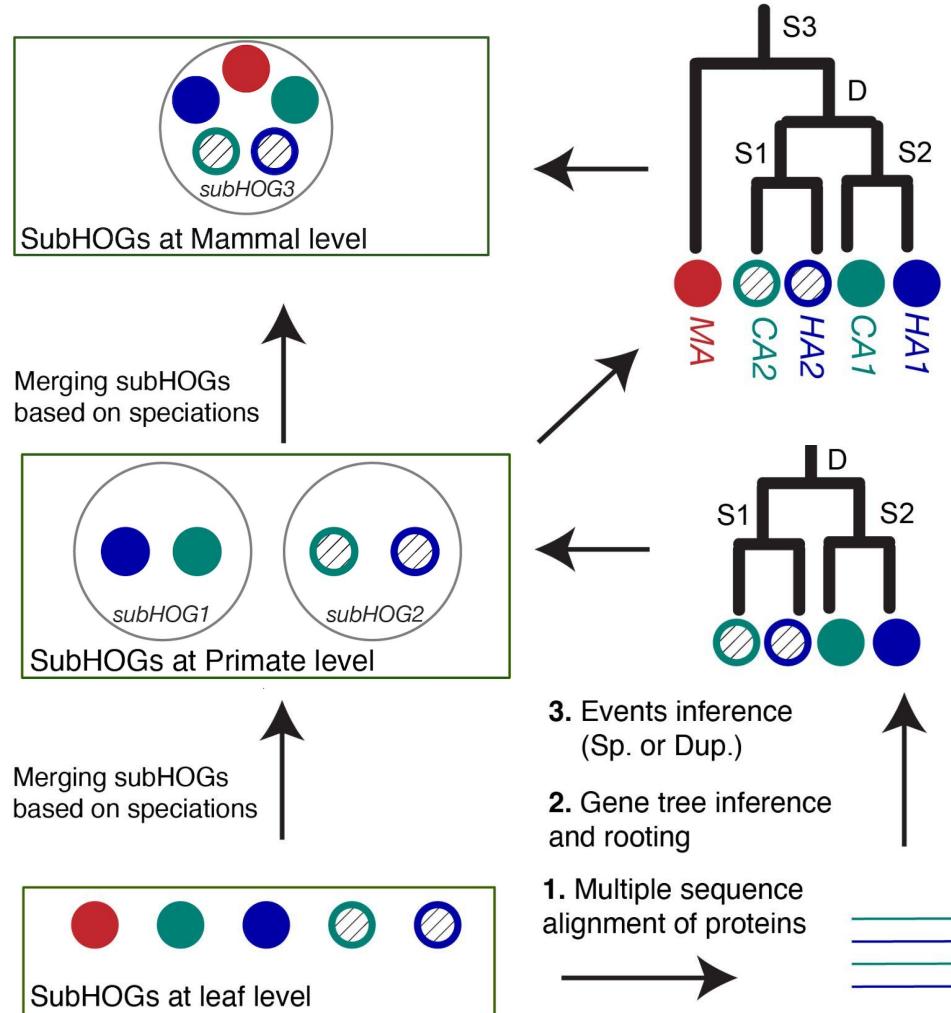
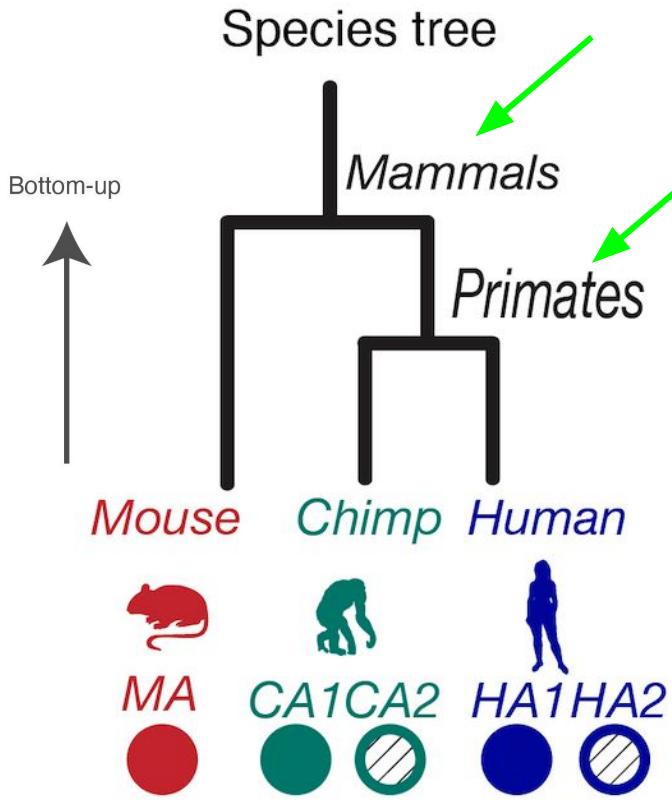


# HOG inference



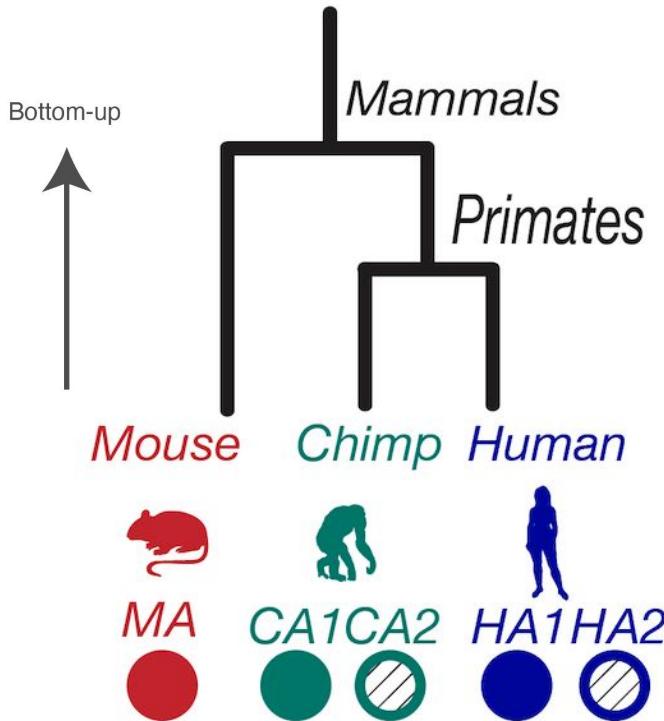
Merging subHOGs based on speciations

# HOG inference



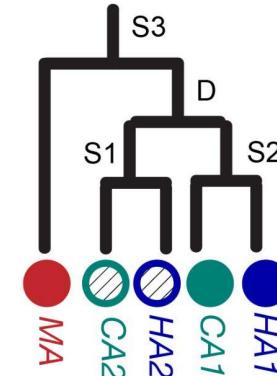
# HOG inference

Species tree

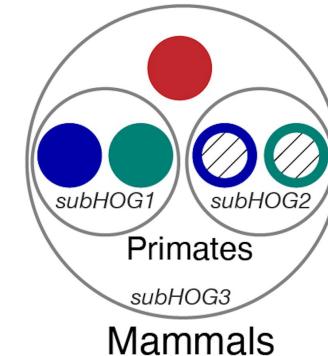


## HOG in orthoXML

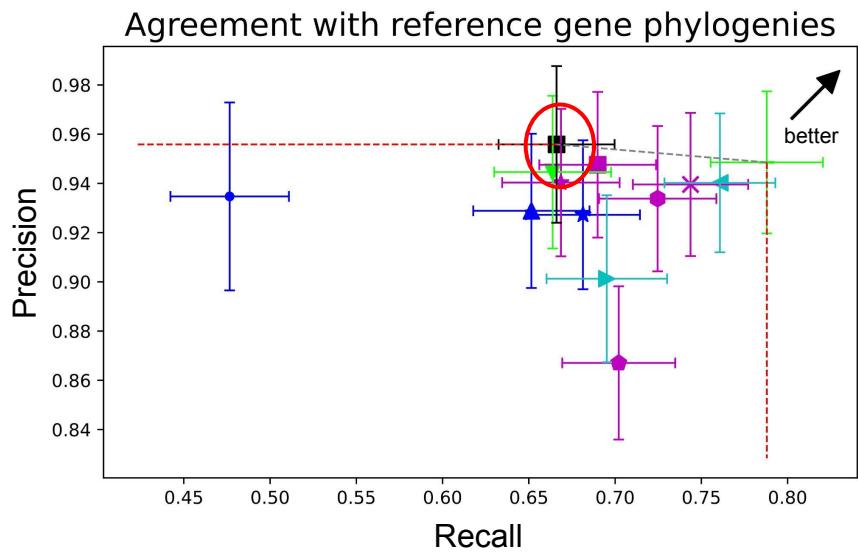
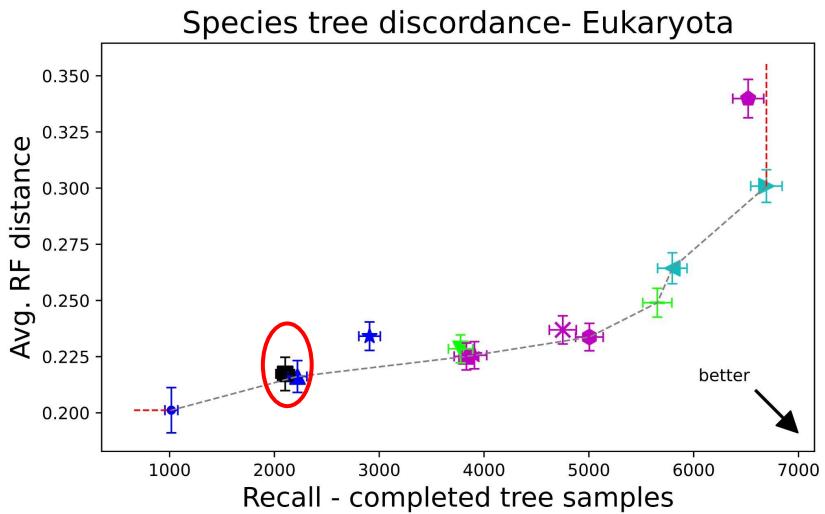
```
<orthoGroup3 Mammals>
  MA
  <paraGroup>
    <orthoGroup1 Primate>
      HA1
      CA1
    </orthoGroup1>
    <orthoGroup2 Primate>
      HA2
      CA2
    </orthoGroup2>
  </paraGroup>
</orthoGroup3>
```



## Nested structure of HOG



# Quest for Orthologs Benchmarking set



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

- InParanoid\_Xenfix
- ◆ OrthoMCL
- Ortholnspector 3

- ★ sonicparanoid
- ◀ PANTHER
- ◆ Ensembl\_Compara

- ▼ Hieranoid\_2
- Orthofinder
- FastOMA

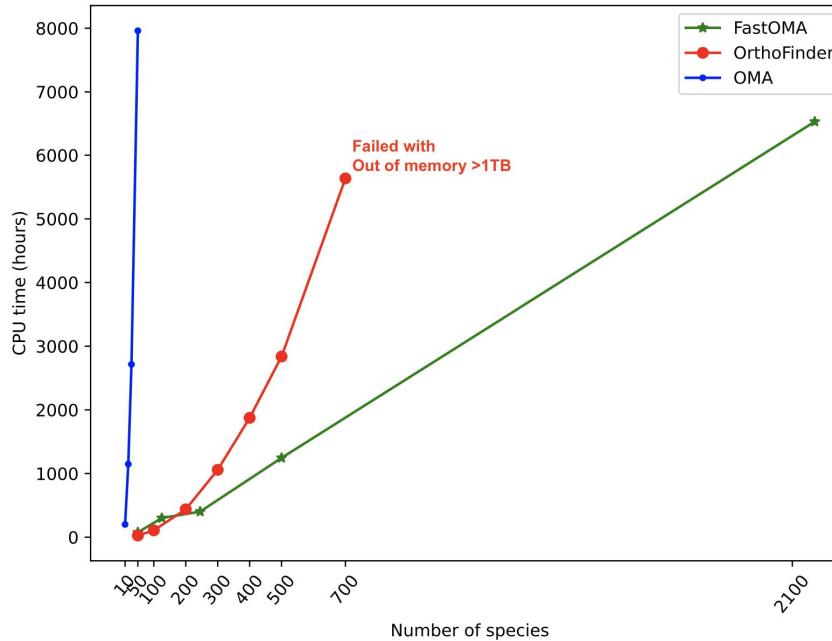
State-of-the-art methods

# Orthology inference for Eukaryote dataset

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



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



# Thank you !

Our lab is hiring postdocs!  
[lab.dessimoz.org](http://lab.dessimoz.org)



SWISS NATIONAL SCIENCE FOUNDATION



Swiss Institute of  
Bioinformatics



UNIL | Université de Lausanne







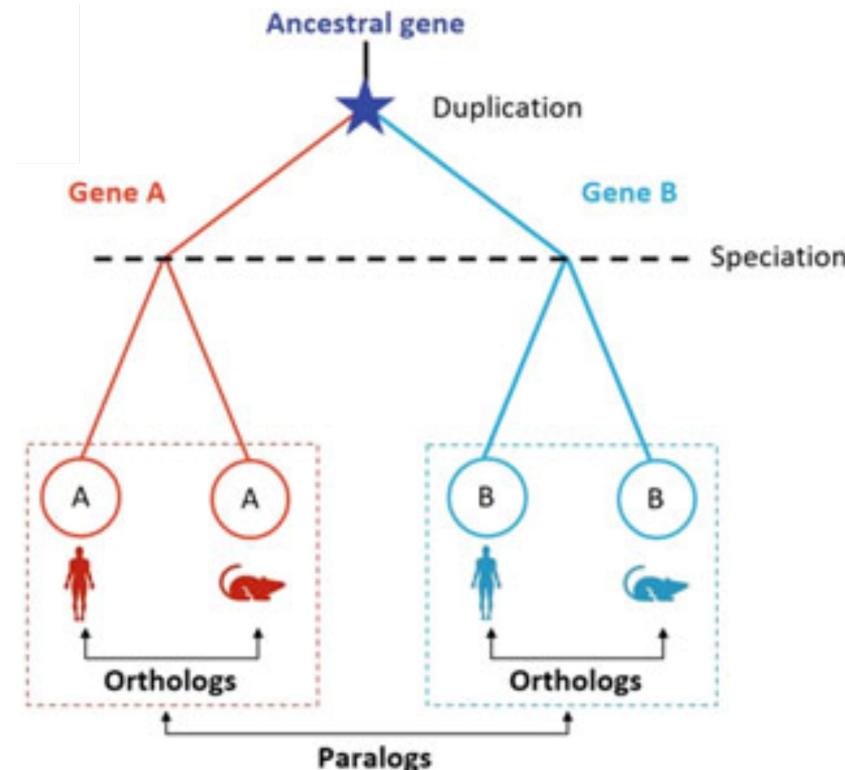
# Outline

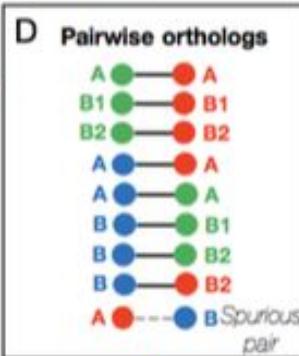
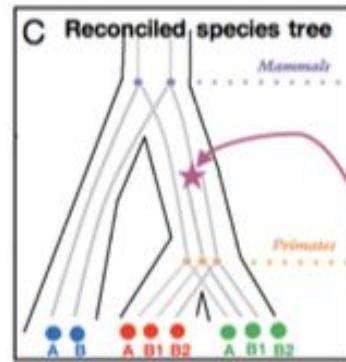
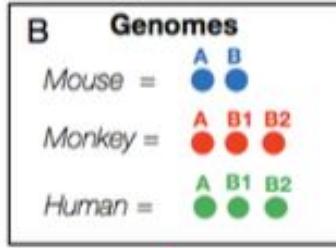
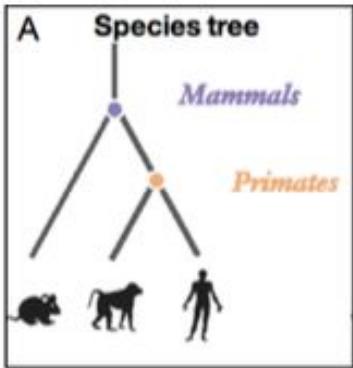
- Introduction to orthology
- FastOMA
- OMA use cases
  - Phylogenetic profiling
  - Gene copy number
- Evolution of phenotypes
- Open questions



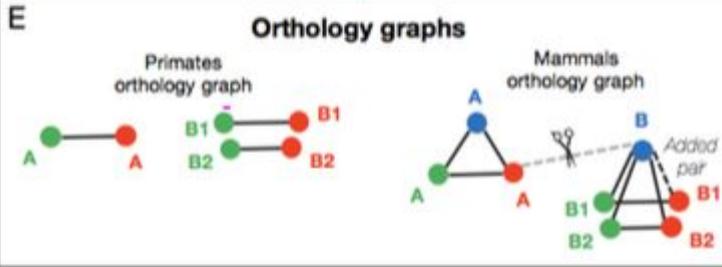
# Introduction

- Homologs
  - two proteins that descend from the same ancestor.
- Orthologs
  - two proteins in two species evolved from one gene in the same common ancestor
  - diverged from a speciation event
  - generally with same function
- Found using all-vs-all alignment
- Informative for studying gene family evolution, phylogeny inference, ...



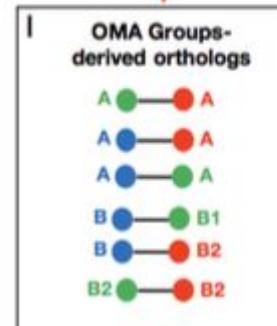
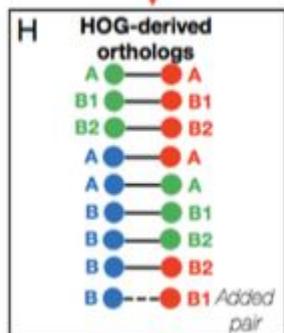
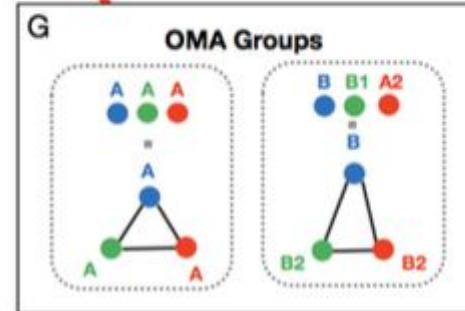
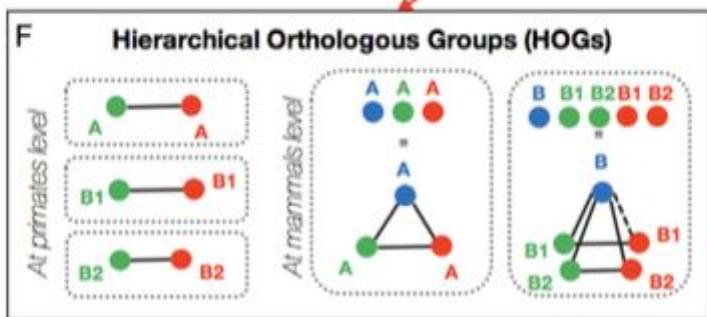


● Gene  
 — Orthologous relation  
 - - - Missing orthologous relation  
 - - - - Spurious orthologous relation



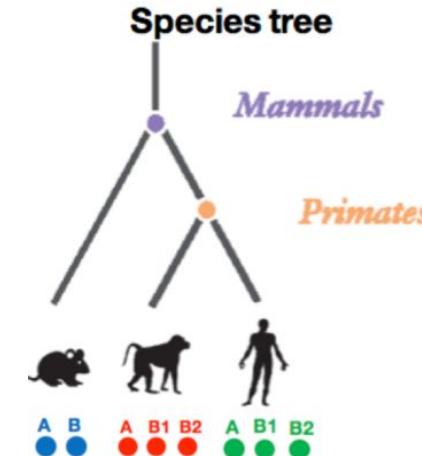
CC

clique

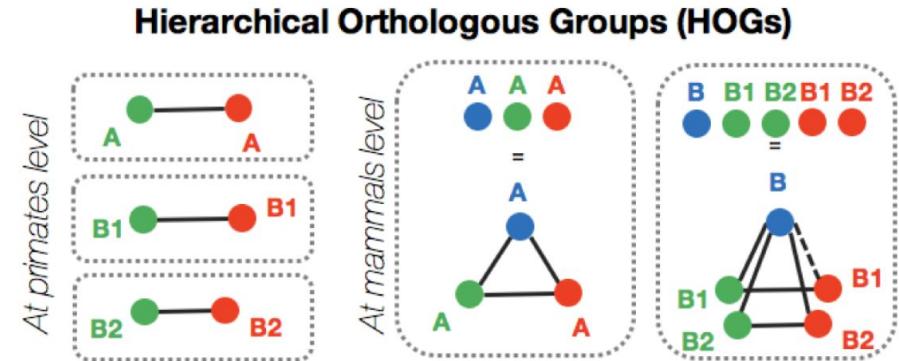


# HOG Hierarchical Orthologous Group

- sets of genes descended from a common ancestral gene at a specific taxonomic level



- Informative for
  - map gene loss and duplication onto species trees
  - Reconstruct ancestral genomes
  - Predict of gene function in new species
  - Identify genetic changes enabling organisms to adapt to different environments



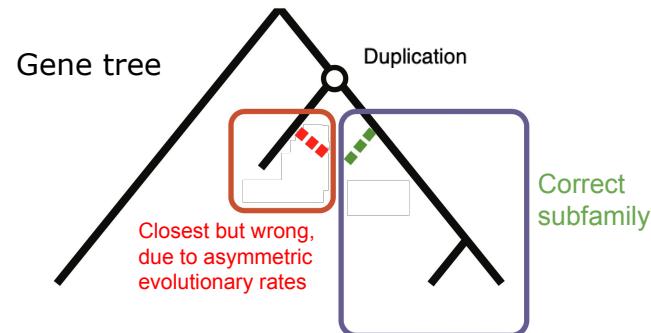
# Basis of the proposed pipeline (FastOMA)

- Map proteins to the database of subfamily of genes
- Traditionally achieved by finding the closest sequence (by BLAST or DIAMOND).



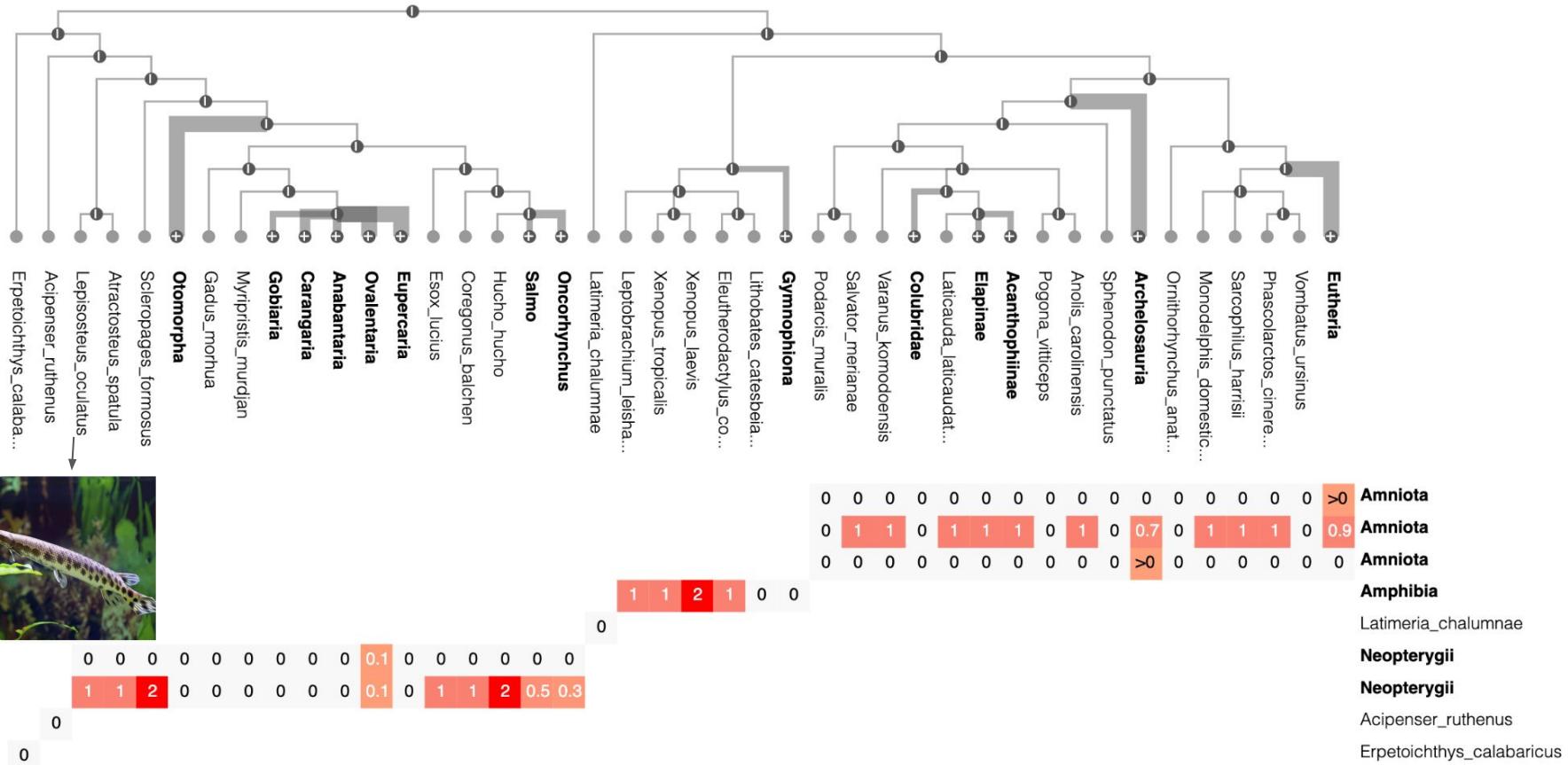
✗ the closest sequence might belong to a different subfamily (not ortholog)

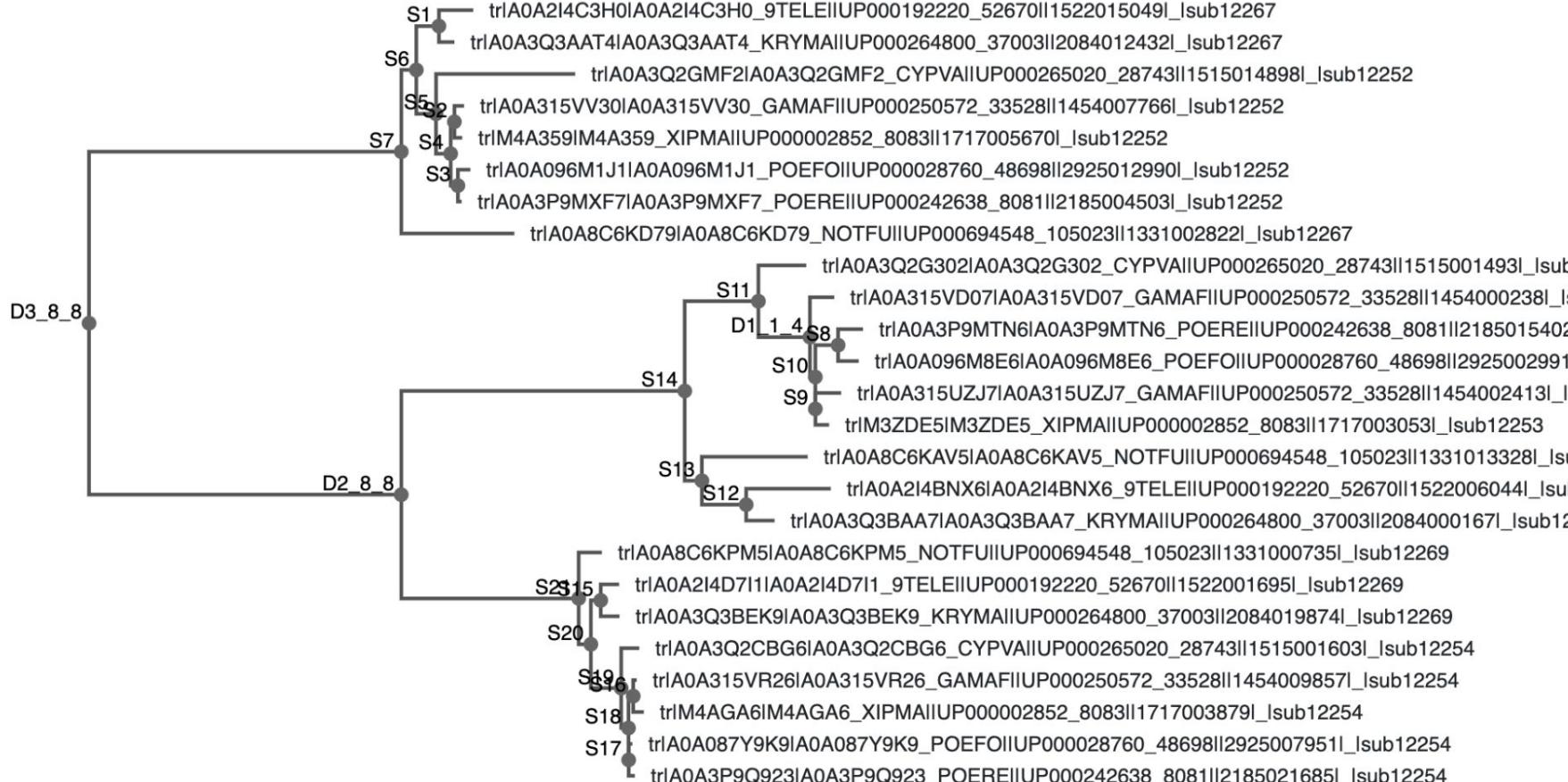
Our solution is OM Amer:  
A subfamily-level classifier  
using subfamily-informed k-mers.

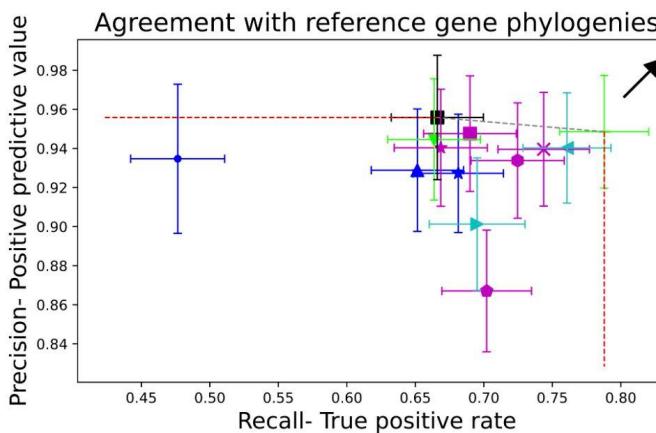
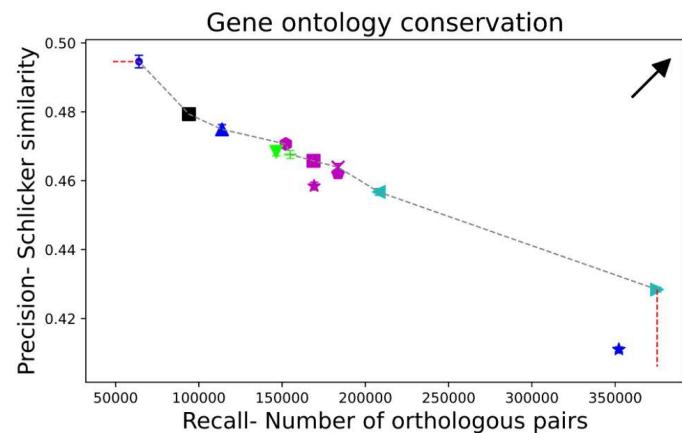
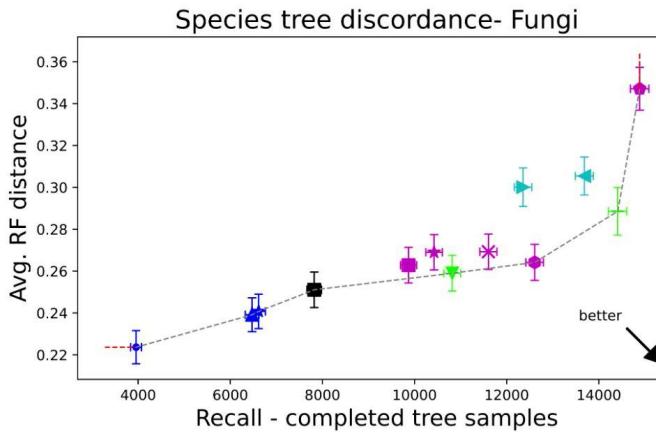
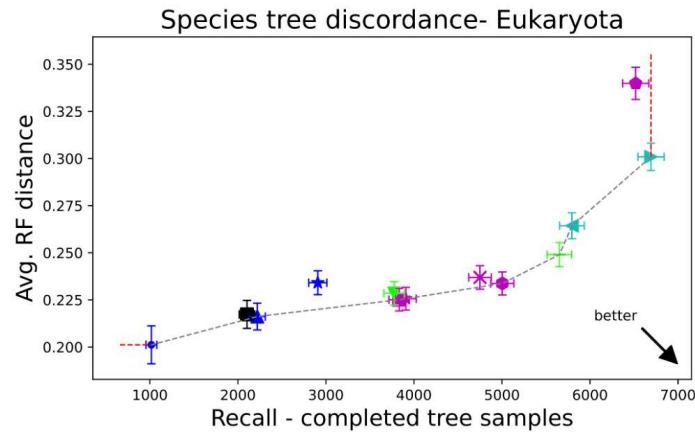


# Versican Core Protein

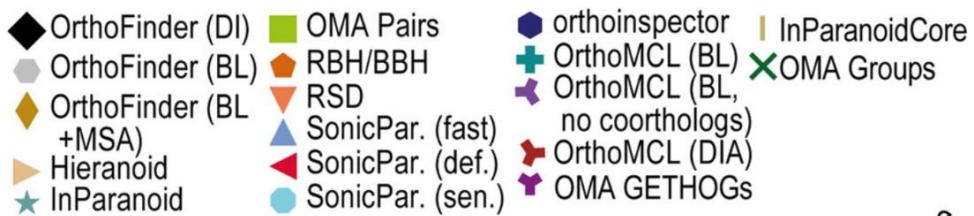
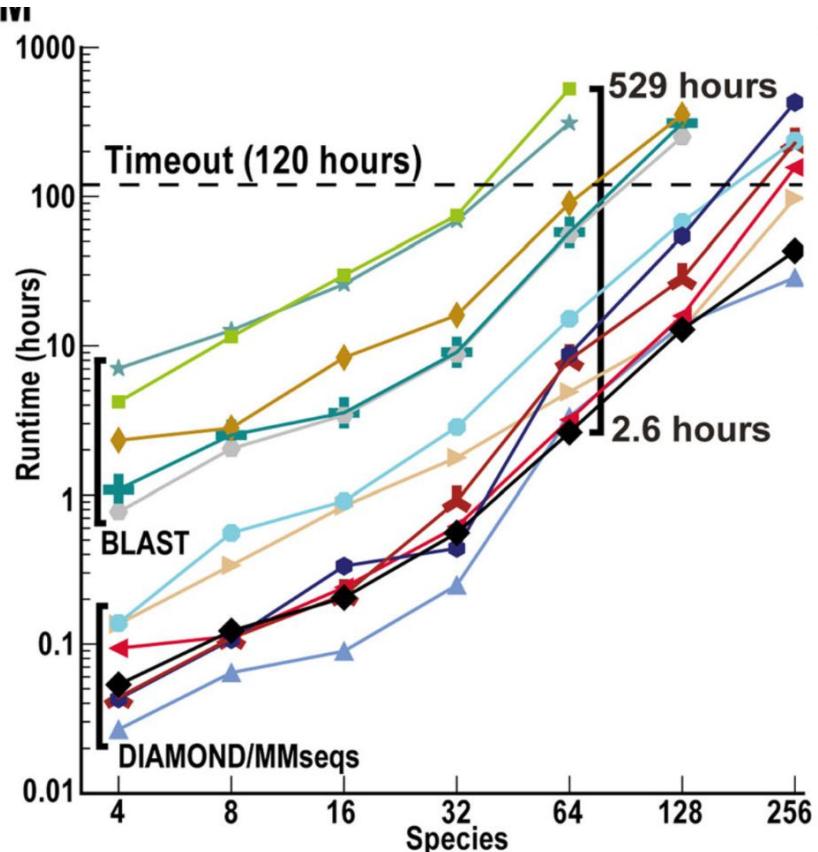
403 proteins  
from 367 species







- |               |                     |                   |               |
|---------------|---------------------|-------------------|---------------|
| ▲ OMA_Pairs   | ■ InParanoid_Xenfix | ★ sonicparanoid   | ▼ Hieranoid_2 |
| ● OMA_Groups  | ▲ OrthoMCL          | ▲ PANTHER         | — Orthofinder |
| ★ OMA_GETHOGS | ● OrthoInspector 3  | ▲ Ensembl_Compara | ■ FastOMA     |
| ✗ Domainoid+  |                     |                   |               |



N

	Orthologs	Multi-species orthogroups	Rooted gene trees	Gene duplications	Rooted species tree
Hieranoid *	✓	✓	✓		
InParanoid	✓				
OMA	✓	✓			✓
OrthoFinder	✓	✓	✓	✓	✓
OrthoMCL	✓	✓	✓	✓	✓
Orthoinspector	✓	✓			
RBH/BBH	✓				
RSD	✓				
SonicParanoid	✓	✓			

\* requires rooted species tree