

## Finding Orthologs and Paralogs

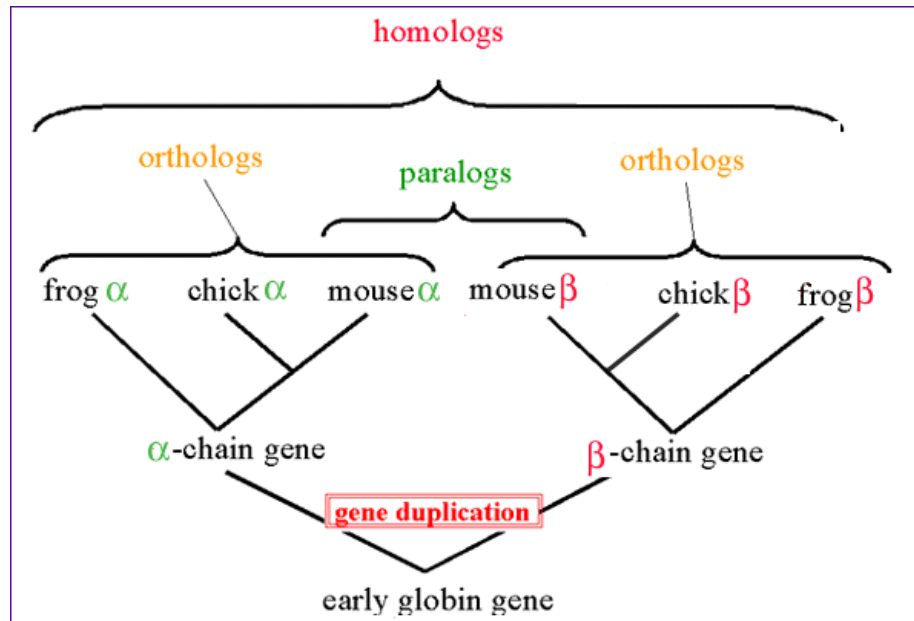


Figure 1: Orthologs

## Gene families and Orthology

Problem: How to find “same” genes across multiple species.

Genes can duplicate (Paralogs) and can be identical due to descent (Ortholog)

## Methods

- BLAST: 1 way BLAST (Gene A in Species X, what is best hit in Species Y)
- BLAST: reciprocal BLAST

## Trees can help resolve relationships

Best hits can sometimes be wrong (B) though it can be resolved with phylogenetics.

## Reciprocal Searches

- Bi-directional or Reciprocal BLAST

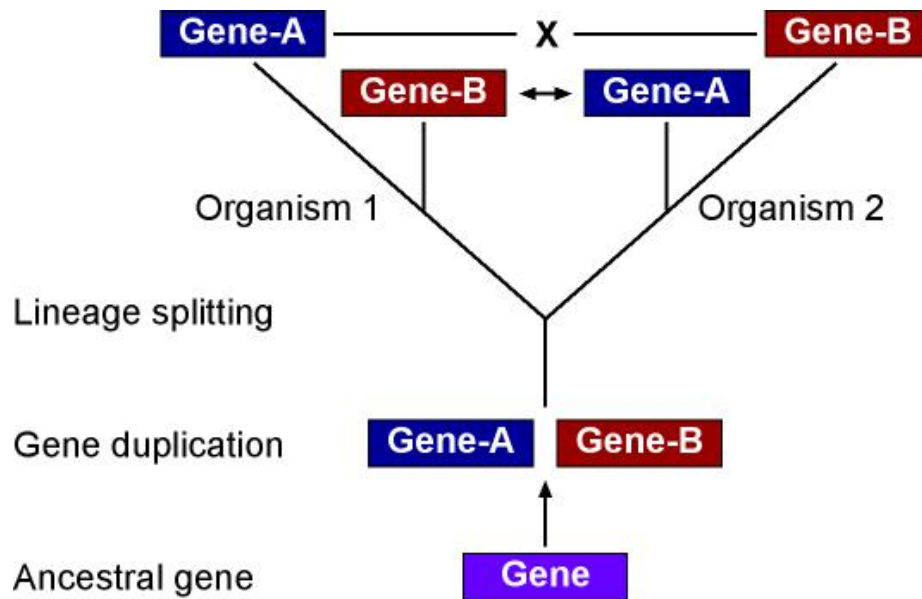


Figure 2: orthologs

## Implement Bidirectional

Method to find best top hit in one direction and the reverse.

Let's walk through the [code](#)

*Will write this in Python in Class*

## Clustering

- Lumping genes together based on similarity linkage
- Single-linkage means if there is a link between A-B then they are in a cluster

## Code up single-linkage

Let's look at some [code](#).

*Will write this in Python in Class*

## Issues

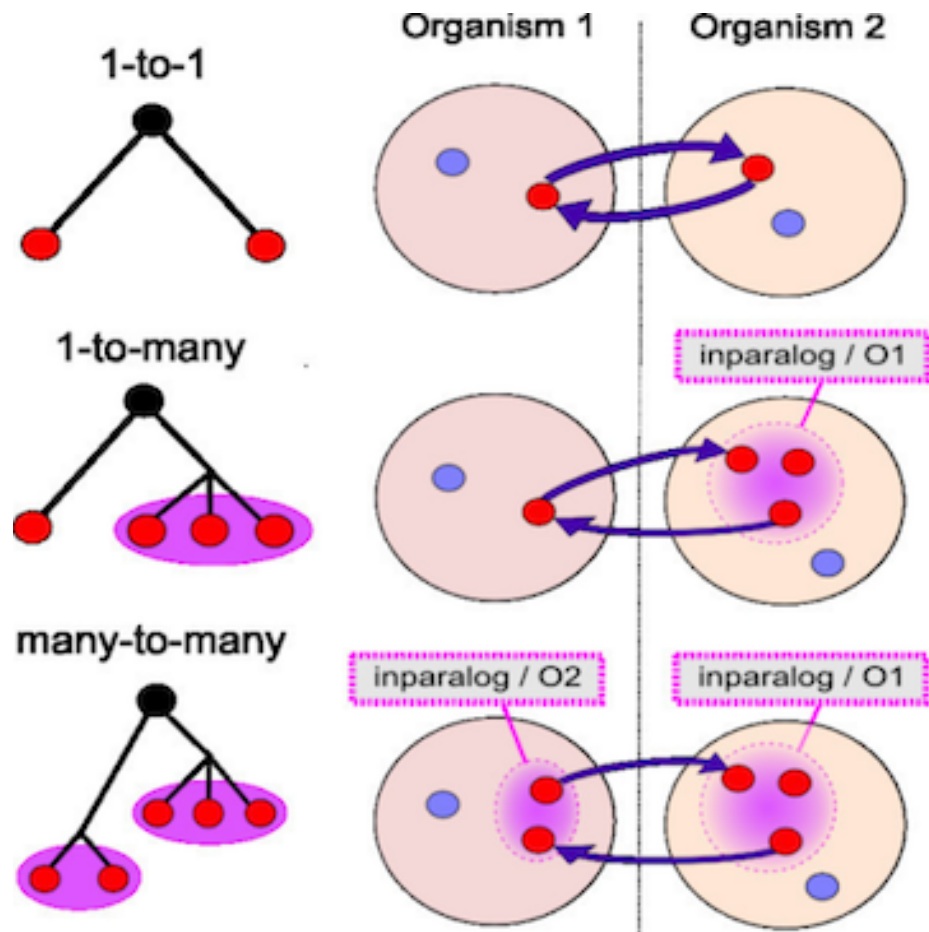


Figure 3: diagramorth

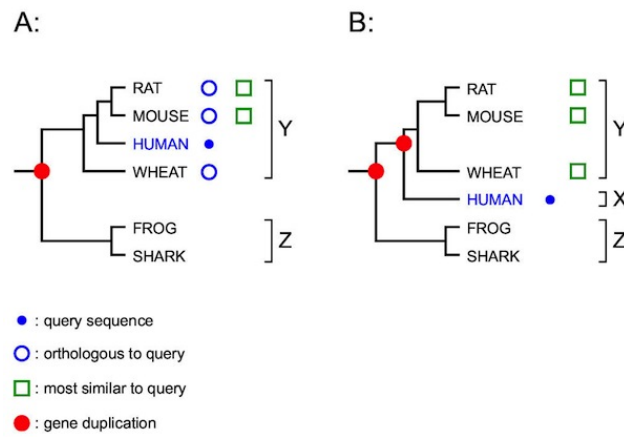


Figure 4: RIO

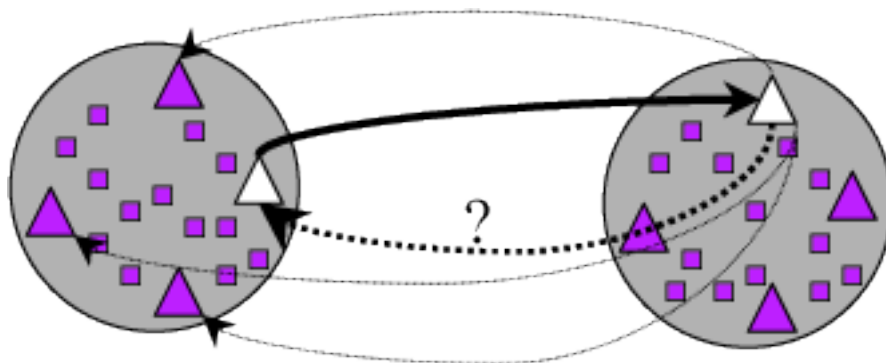


Figure 5: BRH

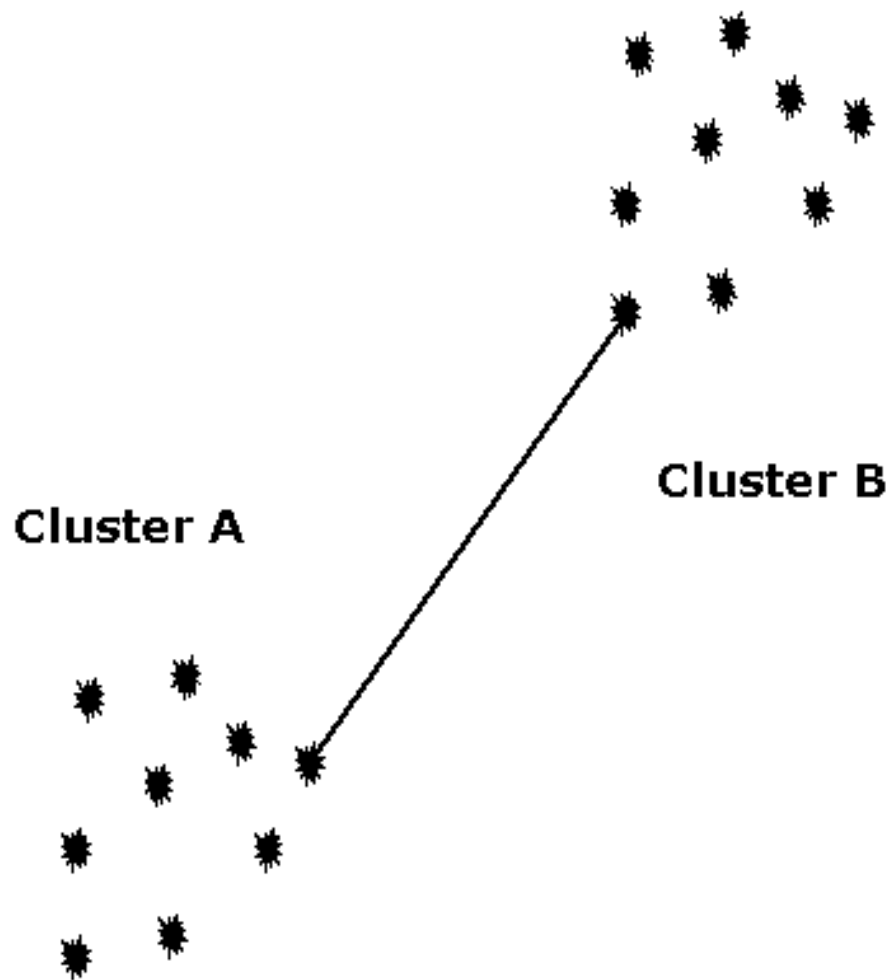


Figure 6: SingleLinkage

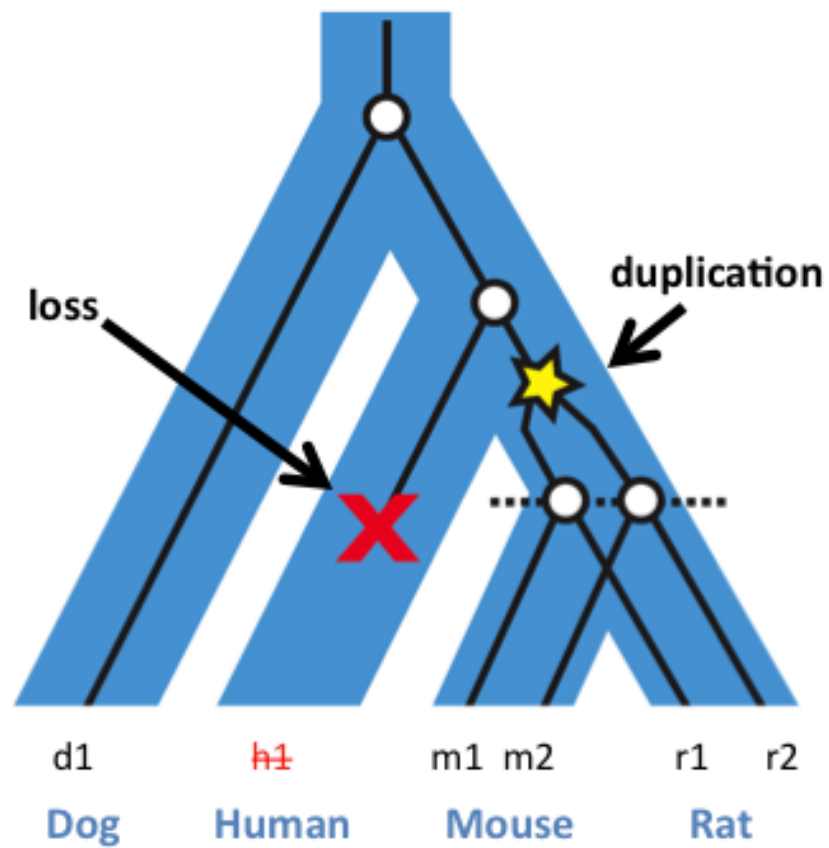


Figure 7: orthologsloss

## Tools to go after Orthologous and Paralogous sequences

- OrthoFinder

## Steps to build orthologs on cluster

We will take 3 datasets of annotated Cyanobacteria, download and run analysis to generate Ortholog table.

```
#!/usr/bin/bash
#SBATCH --ntasks 16 --mem 8G -p short
module load ncbi-blast
module load orthofinder
module load miniconda2
CPU=8

mkdir -p cyanobacteria
cd cyanobacteria
curl -L -O ftp://ftp.ensemblgenomes.org/pub/bacteria/release-45/fasta/bacteria_10_collection
curl -L -O ftp://ftp.ensemblgenomes.org/pub/bacteria/release-45/fasta/bacteria_0_collection
curl -L -O ftp://ftp.ensemblgenomes.org/pub/bacteria/release-45/fasta/bacteria_4_collection

# uncompress files and name them all *.fasta
for file in *.fa.gz
do
    m=$(basename $file .pep.all.fa.gz)
    pigz -dc $file > $m.fasta
done

cd ..

orthofinder.py -a $CPU -f cyanobacteria
```

## Ortholog results

The output file by default will be the date of the analysis. Opening the file `cyanobacteria/Results_XXX/Orthogroups.txt` but I made a [folder](#) in the examples you look over. Here's one [table](#)

## Format

GroupName\tSp1\_Gene1, Sp1\_Gene2\tSp2\_Gene1, Sp2\_Gene2\tSp3\_Gene1, Sp3\_Gene2

Cyanobacterium\_aponinum\_pcc\_10605.ASM31767v1 Nostoc\_punctiforme\_pcc\_73102.ASM2002v1  
OG00000000 EKQ66605, EKQ66611, EKQ66662, EKQ66782, EKQ66954, EKQ66984, EKQ67084, EK

