# Comparison of Microbial Comparative Genomics using Bacteriophages and Mycoplasma bacteria
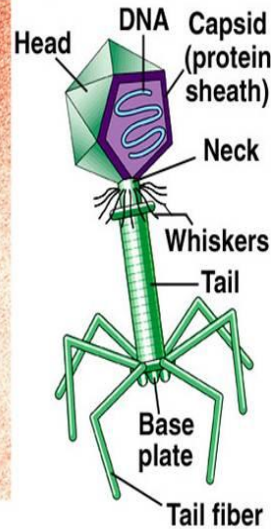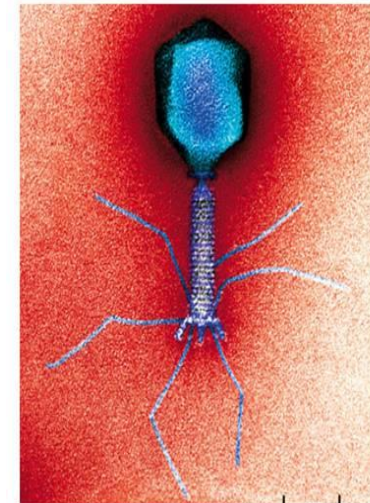
Presented by: Elizabeth Helton

# Overview

-What is a genome, gene,  and bacteriophage?

-Glimpse at Bioconductor

-What is Comparative Genomics?

-Bacteriophage Dataset

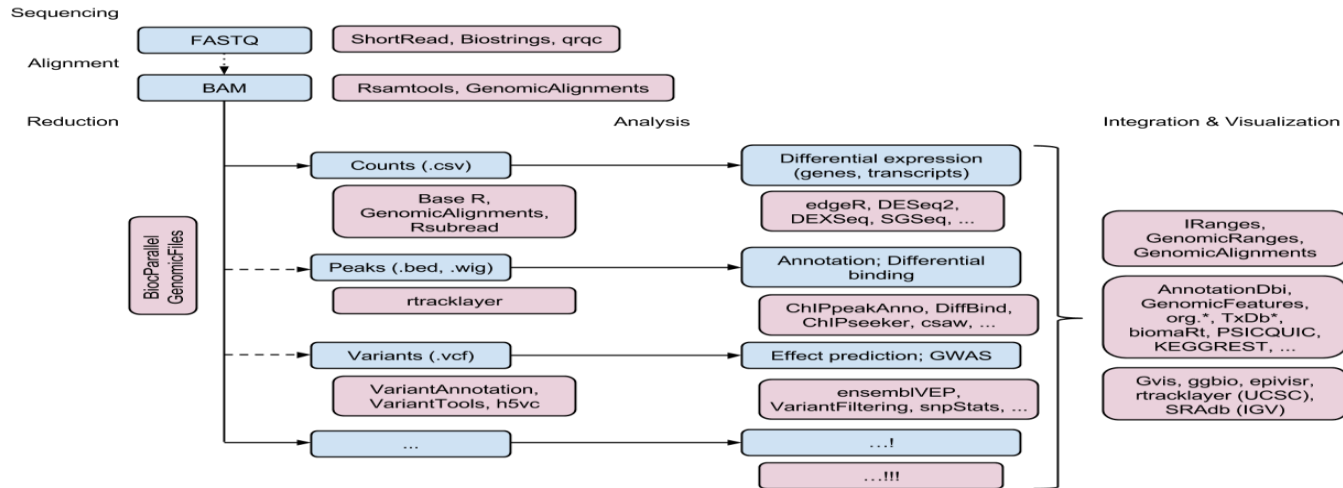- Package 'Find my Friends'

-examples

-summary

# Background Info

-Genome: Organism's complete set of DNA, which includes all of its genes and noncoding sequences

-Gene: sequence of DNA or RNA that codes for a molecule with a function (ex.proteins)

-Bacteriophage:a type of small virus that uses bacteria as a host cell, and destroys the bacteria cell
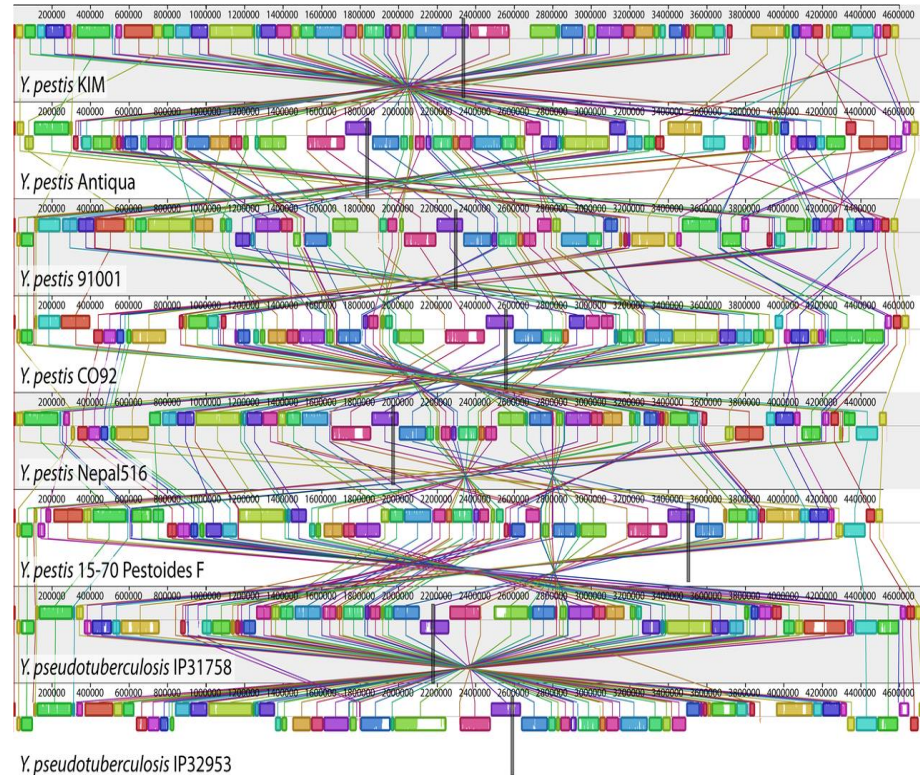
# Bioconductor

- Used for analysis, comprehension, and visual aid of genomic data. It is an open source and open developmental software program. It's primarily used in R programming. Bioconductor uses packages to solve various issues.

# Comparative genomics

-Used to compare complete genome sequences of various species

-Able to identify regions of similarity and differences between species

-Used to better understand the structure and function of human genes and come up with new ways to fight diseases

# Bacteriophage Dataset

Kalah2

- 10 Bacteriophages coming from the Mycobacterium host genus(2 of them were discovered at Webster University)
- Came from Actinobacteriophage Database
- This database shares data, pictures, protocols and analysis tools that were used in the discovery, sequencing and characterization of the phages.
- Bacteriophages Used: Bobby, Cjw1, Dori, Giles, Kalah2, Lilbit, Petra64142, ShereKhan, Spongebob, Webster2

Bobby

# Find my Friends/ comparison

-Framework for microbial comparative genomics. Defines a class system for when working with a pangenome datasets. It allows for a transparency to the underlying sequence data while being able to handle massive collections of genomes.

-Defines a set of novel algorithms that make it possible to create a high quality and speedy pangenome sequence.

| GATTCGATTAG | -> | ATT: | 2 |
|---|---|---|---|
| | | CGA: | 1 |
| | | GAT: | 2 |
| | | TAG: | 1 |
| | | TCG: | 1 |
| | | TTA: | 1 |
| | | TTC: | 1 |

$$\text{similarity} = \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\|_2 \|\mathbf{B}\|_2} = \frac{\sum\limits_{i=1}^{n} A_i B_i}{\sqrt{\sum\limits_{i=1}^{n} A_i^2} \sqrt{\sum\limits_{i=1}^{n} B_i^2}}$$

# Find My Friends Using Bacteriophages Genomes

-cdhitGrouping used to calculate pangenomes. cdhitGrouping repeatedly combines gene groups based on lower similarity thresholds. During each step the longest member in each of the gene groups becomes the model for the next step. It is best to use the lowest threshold possible to ensure that genes that are in the same group can be clustered together

```
> mypang
An object of class pgFull

The pangenome consists of 10 genes from 10 organisms
5 gene groups defined


    Core|
Accessory|==========================
==========================
Singleton|==========

Genes are translated
```

# ExpressionSet

-Views the pangenome matrix as a ExpressionSet object

> as(mypang,'ExpressionSet')

ExpressionSet (storageMode:lockedEnvironment)assayData: 5 features, 10 samples element names: exprs
protocolData: none
Pheno DatasampleNames: Bobby Cjw1 ... Webster2 (10 total)
varLabels: nGenes
varMetadata: labelDescription
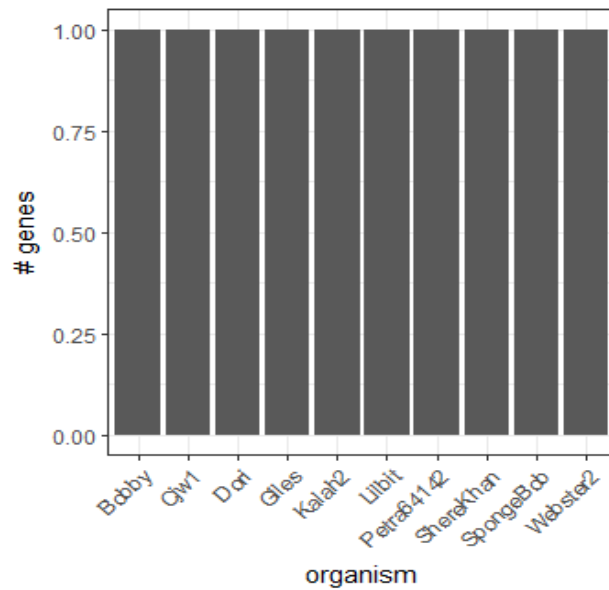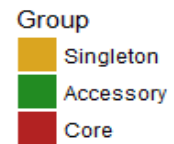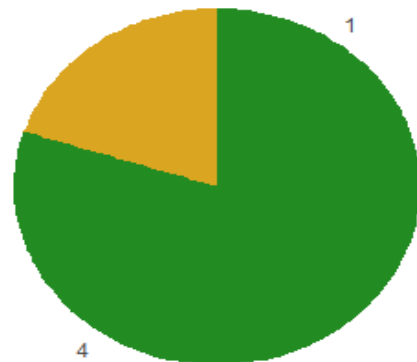featureData featureNames: OG1 OG2 ... OG5 (5total)
 fvarLabels: description group ... nGenes (7 total)
 fvarMetadata: labelDescription

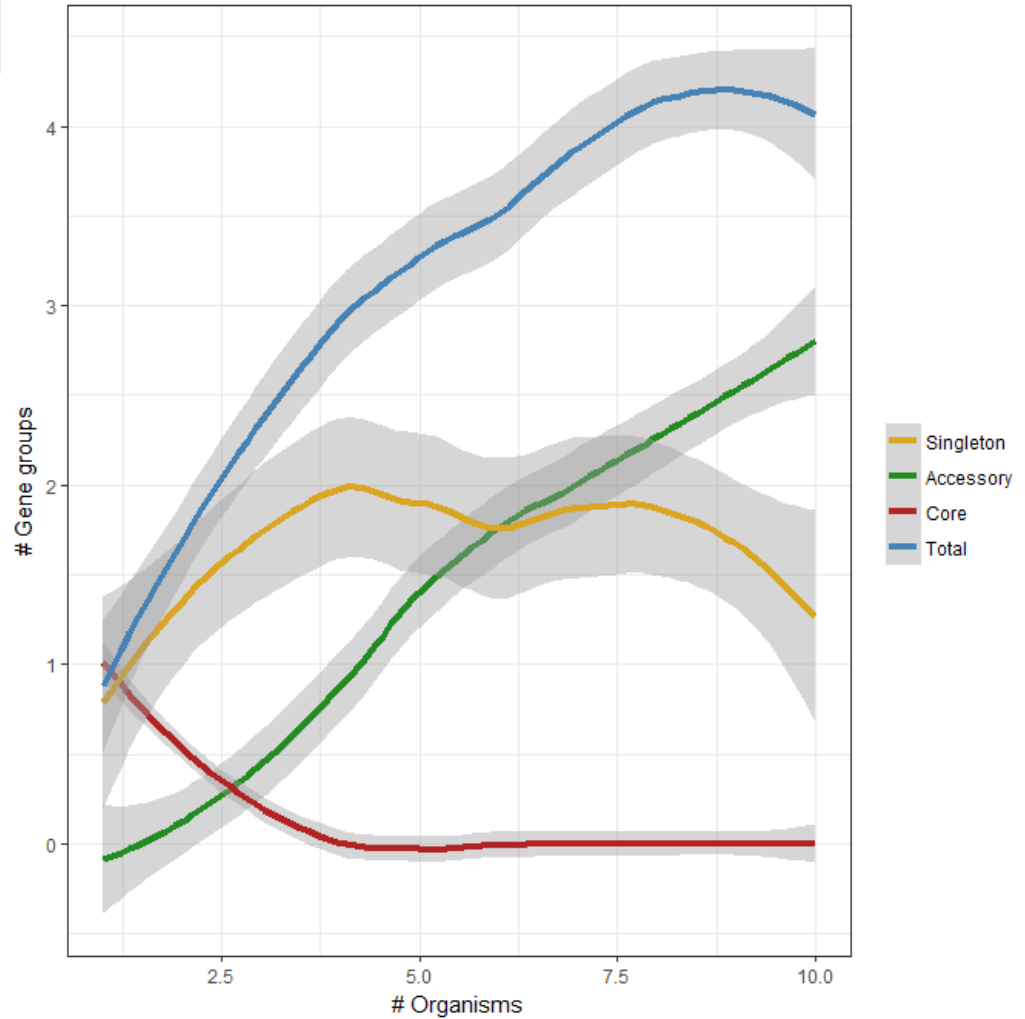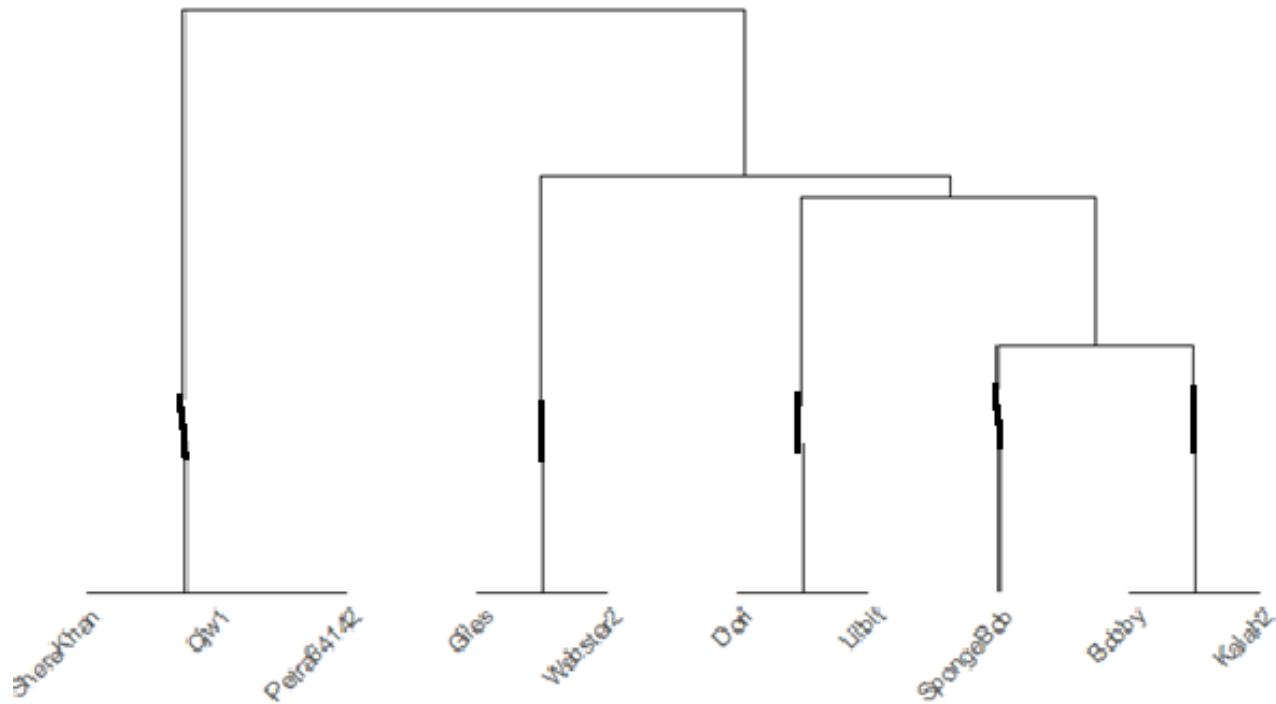experimentData: use 'experimentData(object)'

# Plot Stat

# Evolution Plot

-Views number of singleton,accessory and core genes as the amount of organisms increase

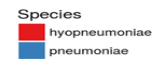-Can be biased toward order of organisms

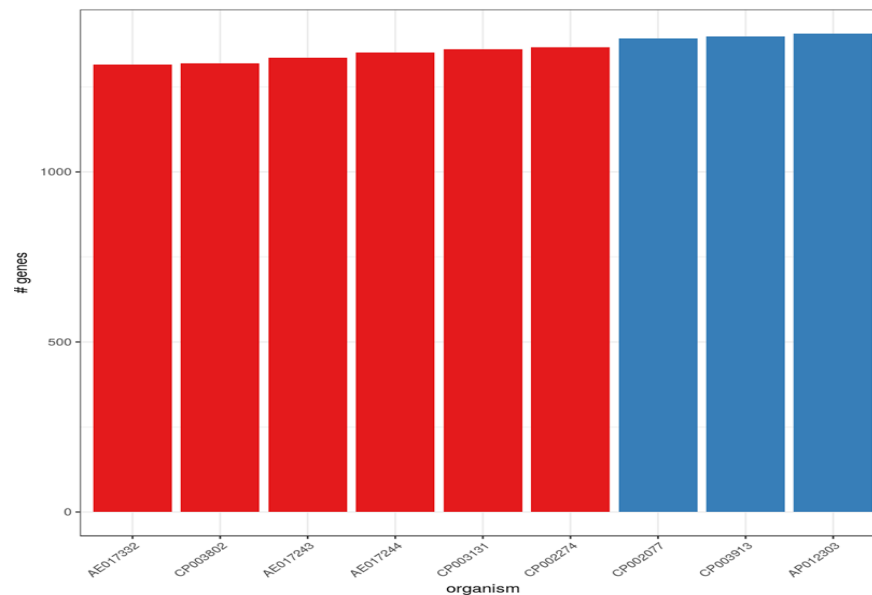# Kmer heatplot

-Comparison of Kmer
values to each organism

# Dendrogram

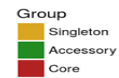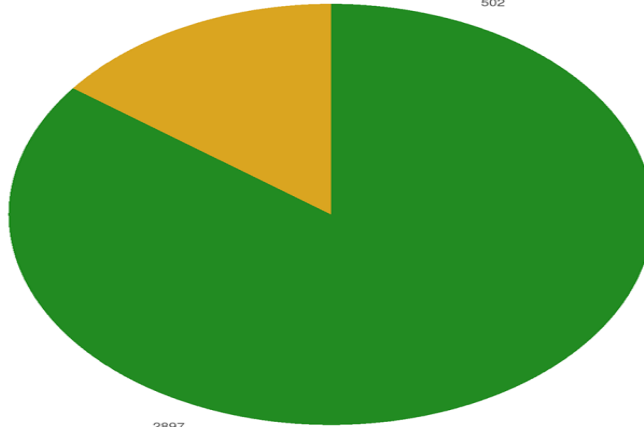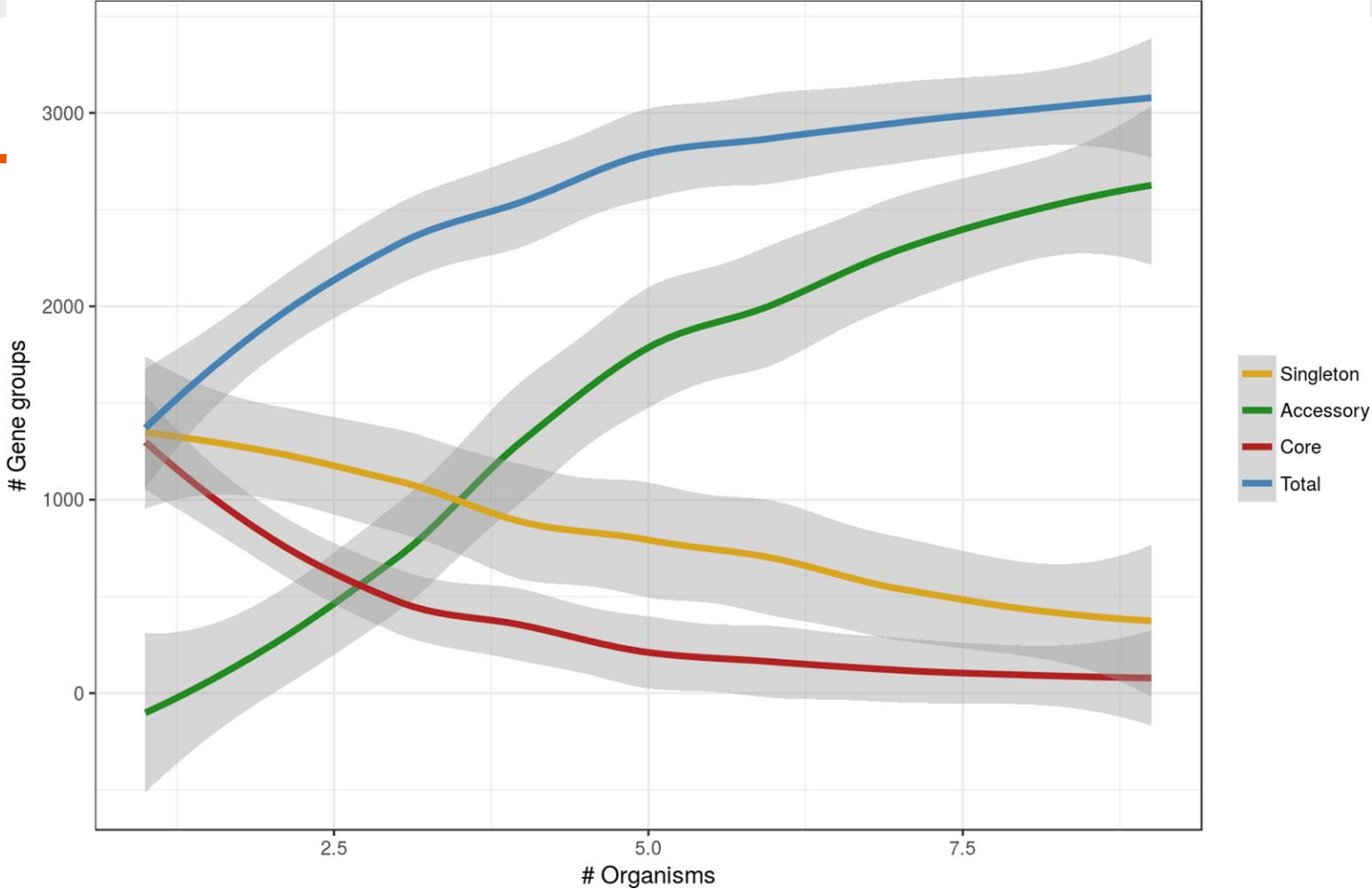# FindMyFriends Using Mycoplasma

mycoPan
## An object of class pgFullLoc
##
## The pangenome consists of 12247 genes from 9 organisms
## 3141 gene groups defined
##     Core|
##Accessory|==================================================:
## Singleton|======
## Genes are translated
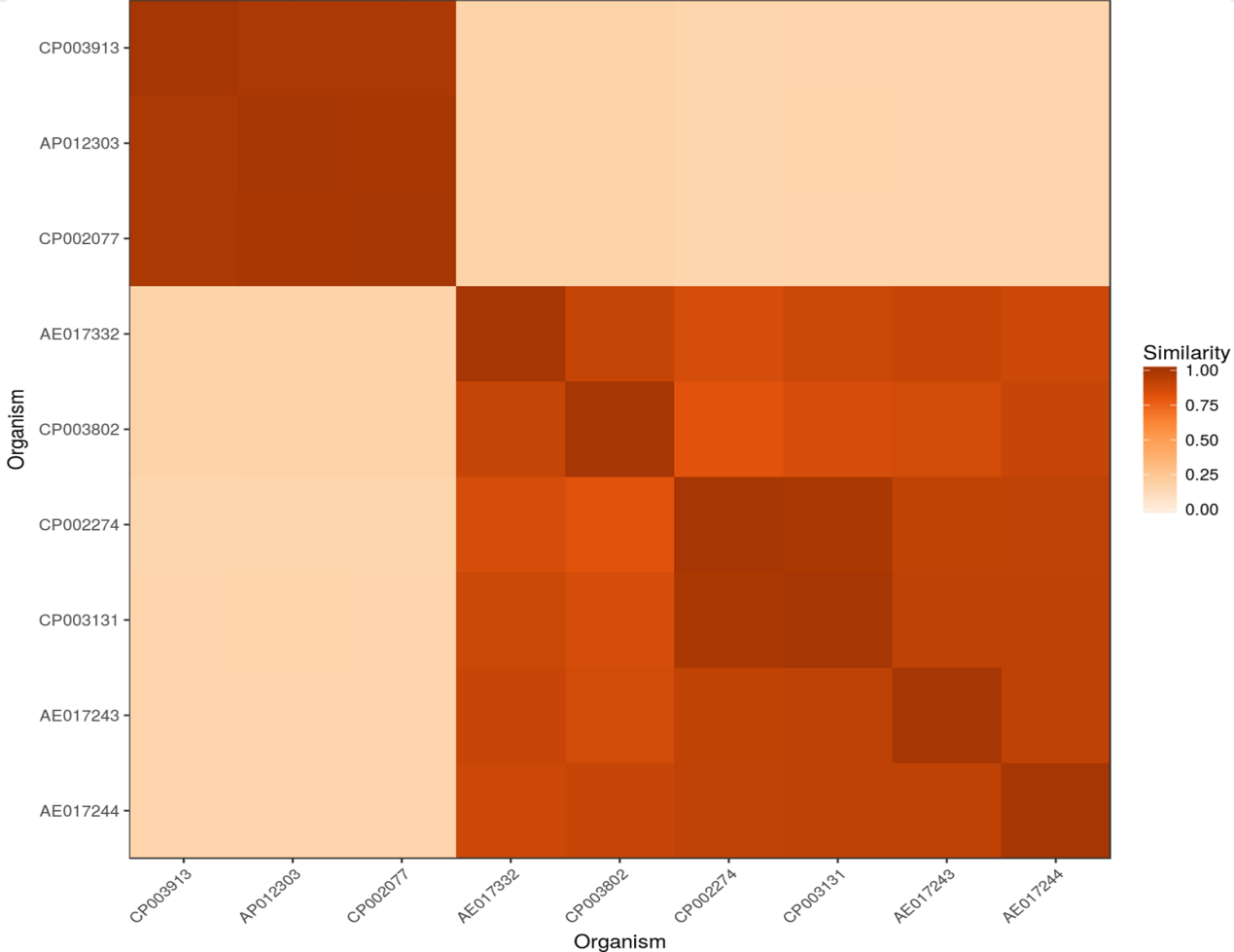
# Pangenome as ExpressionSet

```
## ExpressionSet (storageMode: lockedEnvironment)
## assayData: 3399 features, 9 samples
##   element names: exprs
## protocolData: none
## phenoData
##   sampleNames: AE017243 AE017244 ... CP003913 (9 total)
##   varLabels: nGenes Id ... GenBankDivision (14 total)
##   varMetadata: labelDescription
## featureData
##   featureNames: OG1 OG2 ... OG3399 (3399 total)
##   fvarLabels: description group ... nGenes (7 total)
##   fvarMetadata: labelDescription
## experimentData: use 'experimentData(object)'
```
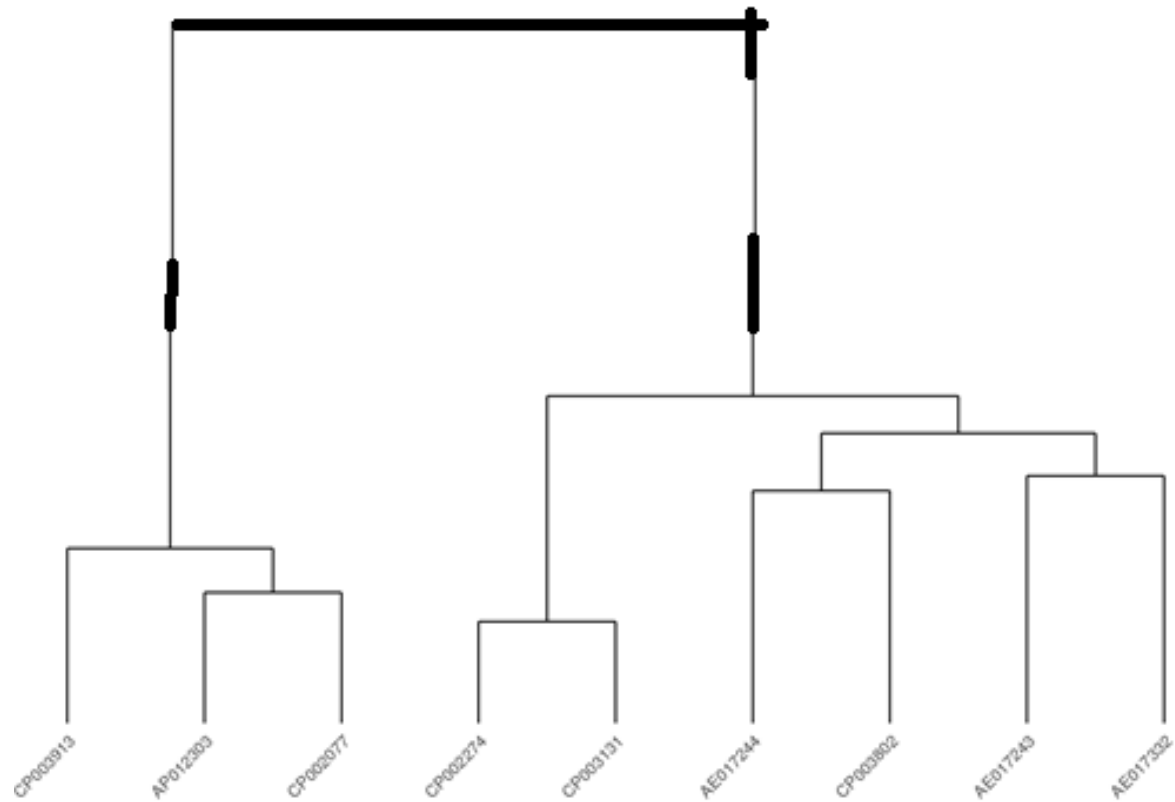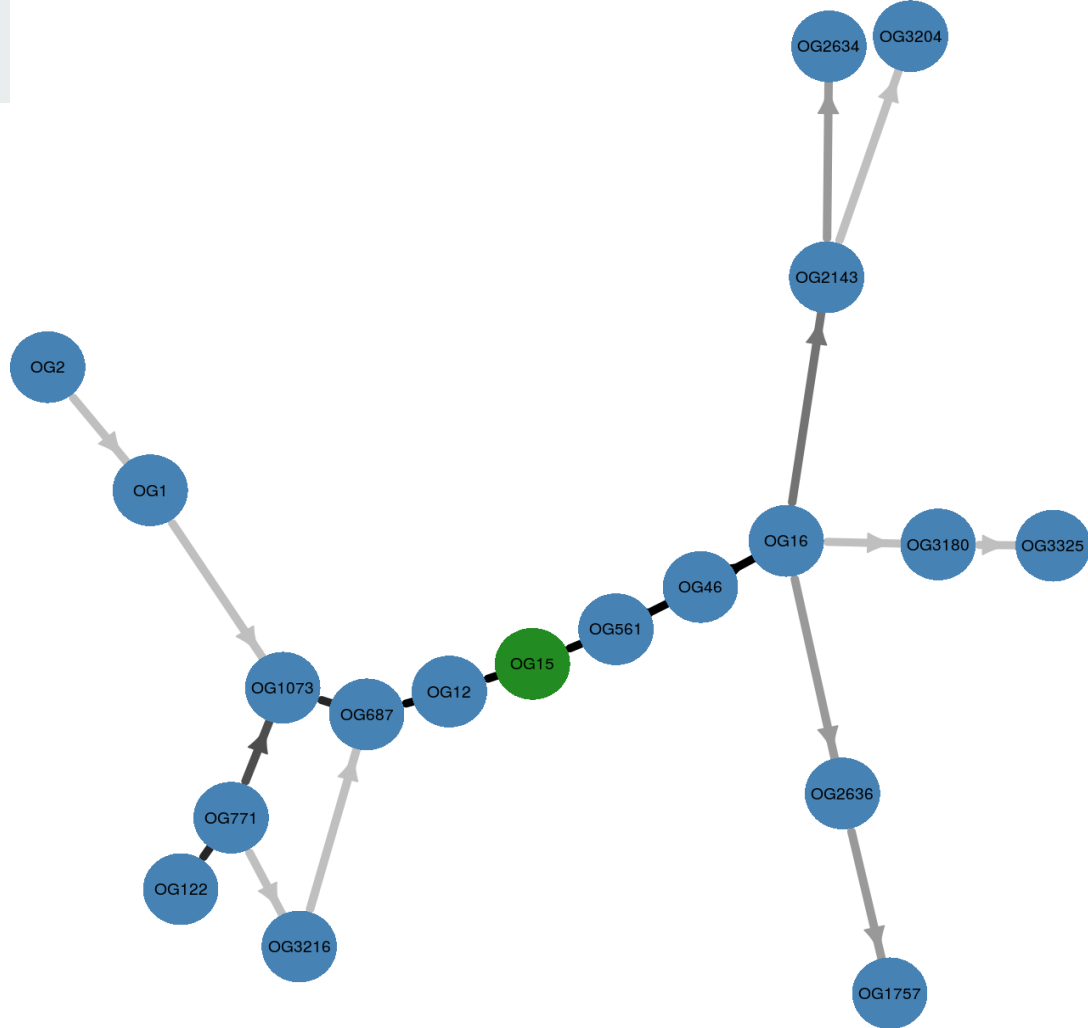
# Evolution Plot

# Kmer Similarity Graph

# Dendogram of Pangenome

# Neighborhood

# References

Pictures on genomes: Google Images

"Actinobacteriophages." *The Actinobacteriophage Database* , 28 Nov. 2017, phagesdb.org/.

Pedersen, Thomas Lin. "FindMyFriends." *Bioconductor*, 2003, bioconductor.org/packages/release/bioc/html/FindMyFriends.html.

Pedersen, Thomas Lin. "Creating Pangenomes Using FindMyFriends." *Bioconductor*, 30 Oct. 2017, www.bioconductor.org/packages/devel/bioc/vignettes/FindMyFriends/inst/doc/FindMyFriends_intro.html.

NIH. "Comparative Genomics Fact Sheet." *National Human Genome Research Institute (NHGRI)*, 3 Nov. 2015, www.genome.gov/11509542/comparative-genomics-fact-sheet/.