

Archaea EvoMining Results

On decade of 1970 to 1980 Archaea was recognized as new life domain, a kingdom different from Bacteria and Eucarya in an exciting first application of 16S phylogeny [woese_phylogenetic_1977]. Main differences between this kingdoms are that Archaeal DNA is not arranged in a nucleus as in Eucarya and Archaeal celular walls are not composed from peptidoglycans as in Bacteria. Archaea proteins may be highly valuable to biotechnology industry for their great stability due to Archaea extreme habitat conditions on temperature PH and salt content. Despite no Archaeal Natural products biosynthetic gene clusters (BGC's) has been reported on MiBIG Archaea do have BGC's, some seems to be acquired by horizontal gene transfer (HGT) like methano nrps {search reference}. Other Archeal natural products known are archaeosins, Diketopiperazines, Acyl Homoserine Lactones, Exopolysaccharides, Carotenoids, Biosurfactants, Phenazines and Organic Solutes but this knowledge is not comparable to Bacterial knowledge[charlesworth_untapped_2015].

Natural products biosynthetic gene clusters search is actually performed using either *high-confidence/low-novelty* or *low-confidence/high-novelty* bioinformatic approaches [medema_computational_2015]. High confidence methods compares query sequences with previously known BGC's such as nrps or PKS, examples of this algorithms are antiSMASH and clusterfinder [antismash_????]. EvoMining searches on expansions from central metabolic pathways enzyme families, it has been classified as low confidence/high novelty method. EvoMining has proved useful on Actinobacteria phylum where its use lead to Arseno-compounds discovery [cruz-morales_phylogenomic_2016]. Also on Actinobacteria antiSMASH analysis on 1245 genomes found 774 different classes of natural products, the same analysis on 876 Archaeal genomes, a full kingdom, identifies only 35 BGC's classes. So either Archaea does not have natural products BGC's or this are not yet known. Next paragraph deals with a possible approach about how natural products BGC's can be find.

Archaea resembled Bacteria in that Archaea uses horizontal gene transfer as a genic interchange mecanism, Archaeal genomes contains operons [howland_surprising_2000] and in general there is introns absence{Reference to Computational Methods for Understanding Bacterial and Archaeal Genomes}. Archaeas do have introns, but they are mainly located on genes that encodes ribosomal and transfer RNA [howland_surprising_2000]. General lack of introns allows automatic genome annotation, operons gene organization permits functional inference to a certain degree and HGT contribute to expansions on Archaeal genomes. Some phylum on Archaea has an open pangenome, and as we will show on this chapter some Archaea has central pathway expansions. Enzyme families from central pathways expansions, open pangenome and operon organization made EvoMining succesful on Actinobacteria, this lead us to think that evoMining is suitable to analyze Archaeal genomes, even more since EvoMining is a method oriented to use evolution and its not entirely based on previous knowledge of BGC's sequences if evolutionary logic behave on Archaea as on bacteria, new BGC's classes may be be found on Archaea.

EvoMining is a trade off between conserved known central metabolic function and enough expansions divergence on sequence and on clusters to divergence

Tables

Table 1: Families on Archaeabacteria

Factors	Correlation between Parents & Child
GenomeDB	876
Phylum	12
Order	23

```
label(path = "chapter3/expansion_plotArchaeas.pdf", caption = "Expansions Boxplot",label =
```



2

Central pathway expansions

Heat plot of central pathways expansions, Needs to be phylogenetically sorted.

```
## Warning: Ignoring unknown aesthetics: width
```

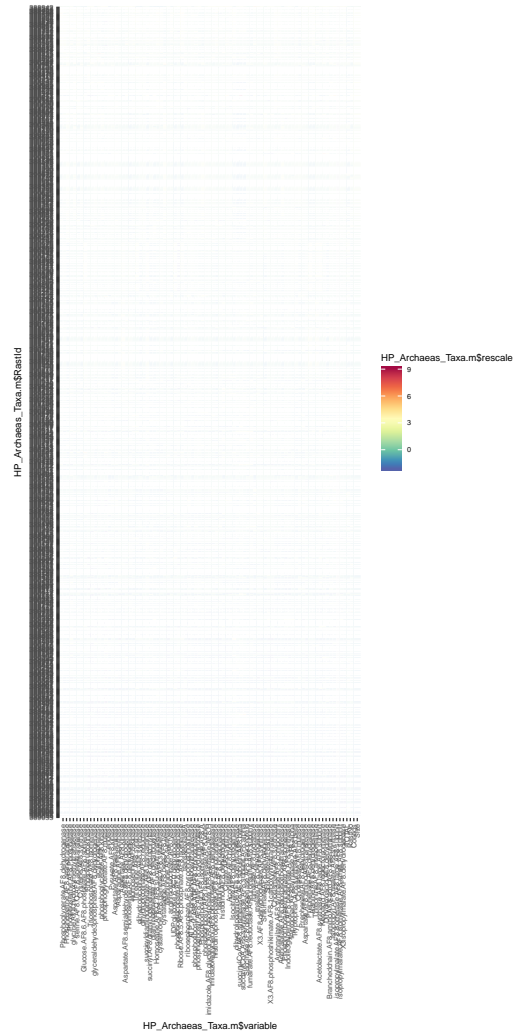


Figure 2: Archaeas Heatplot

Here is a reference to the HeatPlot: Figure 2.

Genome Size correlations

Correlation between genome size and AntiSMASH products

Genome size vs Total antismash cluster coloured by order

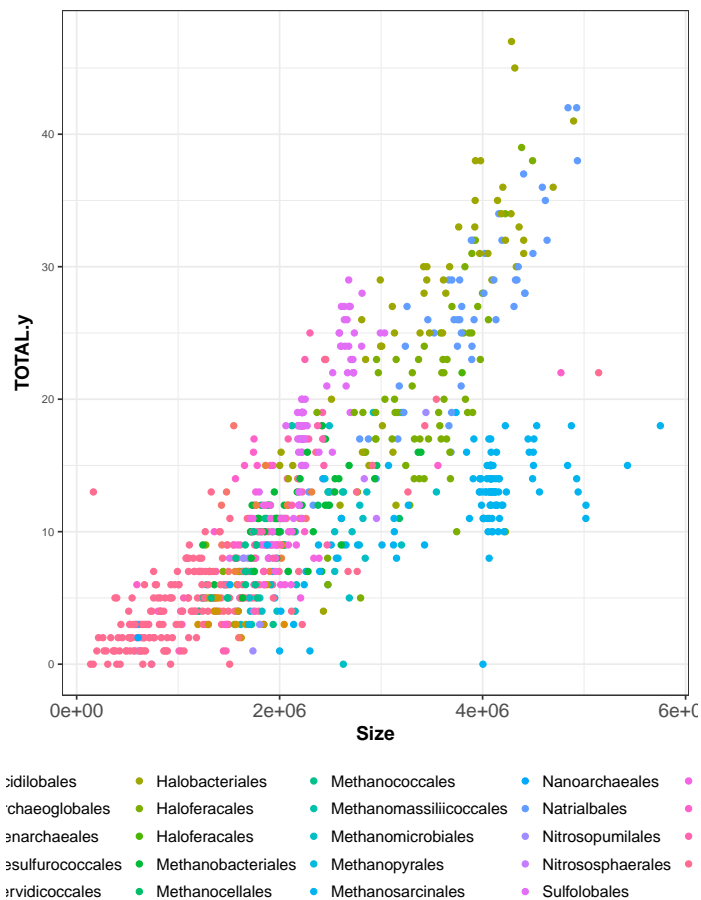


Figure 3: Correlation between Archaeas genome size and antismash Natural products detection colored by Order

Here is a reference to Genome size vs Total antismash cluster: Figure 3.

Genome size vs Total antimash cluster detected splitted by order

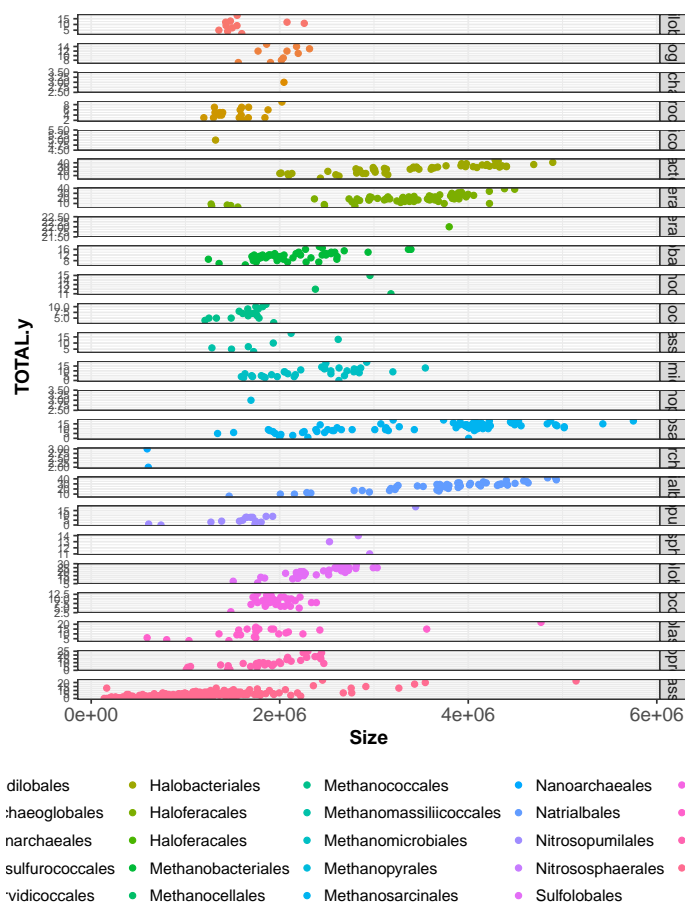


Figure 4: Correlation between Archaeas genome size and antimash Natural products detection grided by Order

Here is a reference to Correlation between genome size and antimash Natural products detection grided by Order plot: Figure 4.

Correlation between genome size and Central pathway expansions

Genome size vs Total central pathway expansion coloured by order

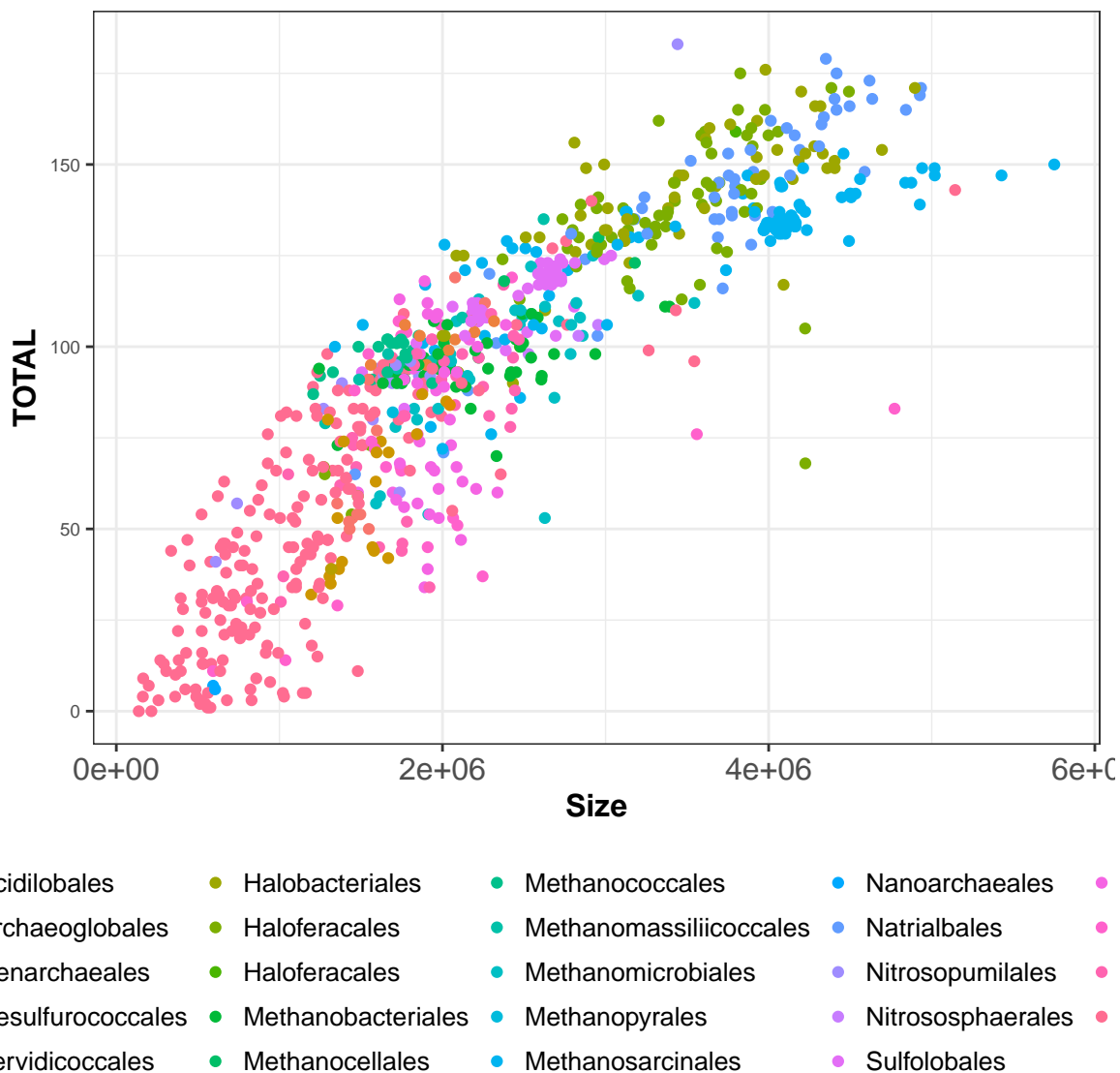


Figure 5: Correlation between Archaeas genome size and central pathway expansions

Here is a reference to the size vs Total central pathway expansion plot: Figure 5.

Genome size vs Total central pathway expansion grided by order

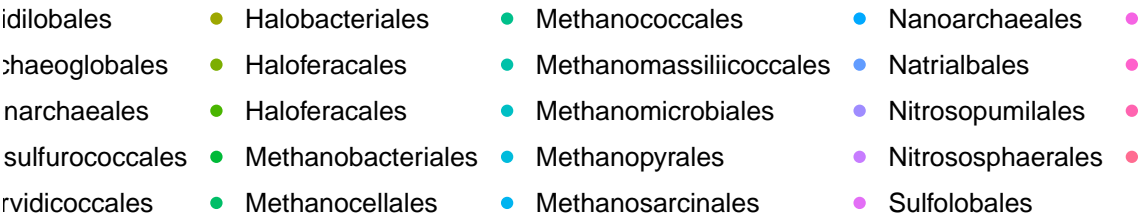
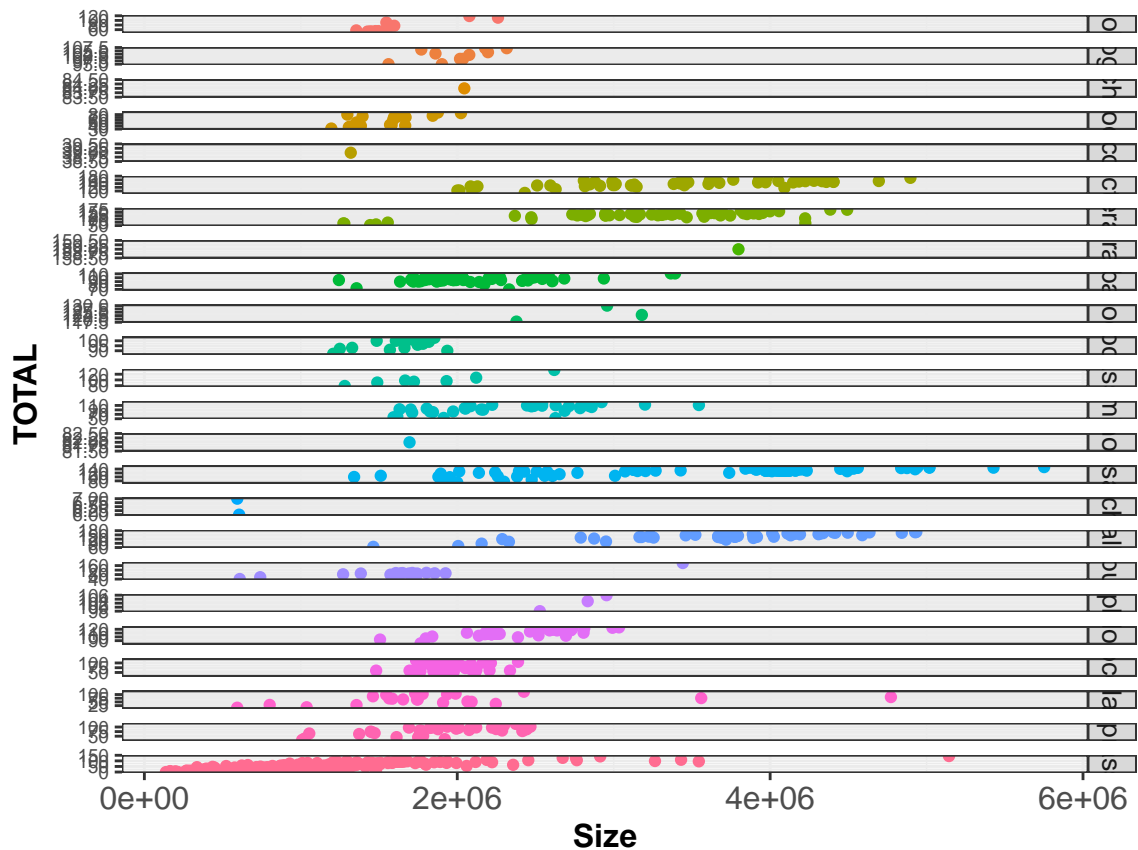


Figure 6: Correlation between Archaeas genome size and central pathway expansions grided by order

Here is a reference to the Genome size vs Total central pathway expansion grided by order plot: Figure 6.

Correlation between genome size and each of the central pathway families. Data are coloured by metabolic family instead of coloured by taxonomical order. This treatment allows to answer how different metabolic families grow when genome size grows.

Also I want to add form given by taxonomical order.

```
## Warning: The shape palette can deal with a maximum of 6 discrete values  
## because more than 6 becomes difficult to discriminate; you have  
## 24. Consider specifying shapes manually if you must have them.  
## Warning: Removed 64823 rows containing missing values (geom_point).
```

Genome size vs Total central pathway expansion coloured by metabolic Family

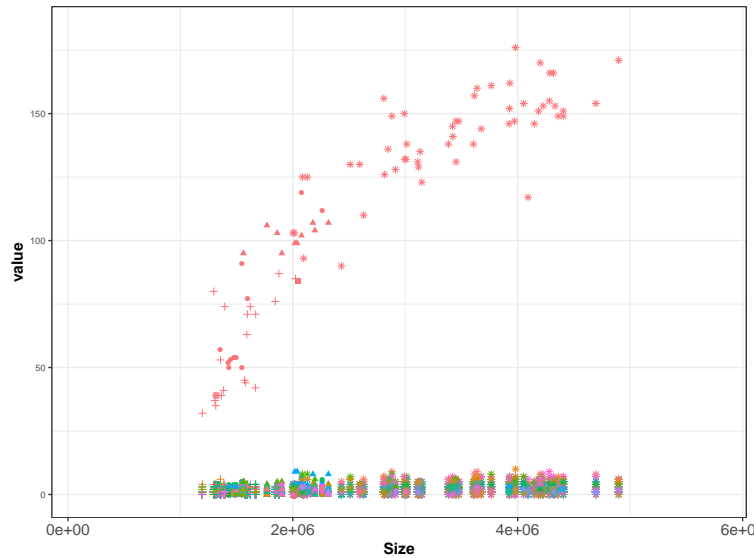


Figure 7: Correlation between Archaeas Genome size vs Total central pathway expansion coloured by metabolic Family

Here is a reference to the Genome size vs Total central pathway expansion coloured by metabolic Family plot: Figure 7.

Future Work: Genome size vs Total central pathway expansion grided by metabolic Family For clarity I need to also grid and group by Metabolic Pathway

Here is a reference to Genome size vs Total central pathway expansion grided by metabolic Family plot: ??.

Natural products

Natural products recruitments from EvoMining heatplot

We can see natural products recruitment after central pathways expansions colored by their kingdom. Natural products recruited by metabolic family, colored by phylogenetic origin.

Recruitments after central pathways expansions coloured by Kingdom

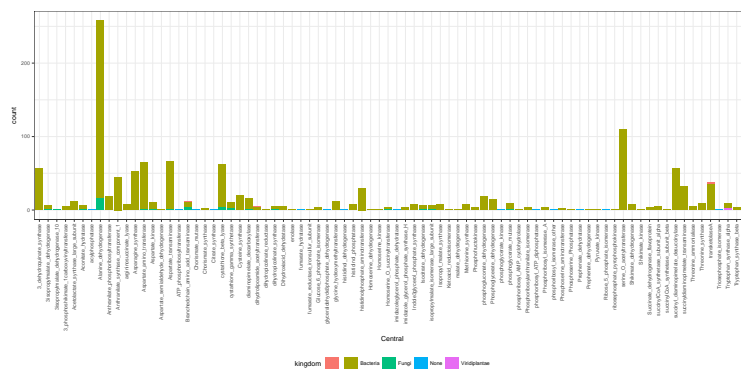


Figure 8: Archaeas Recruitmens on central families coloured by kingdom

Here is a reference to Recruitments after central pathways expansions colourd by Kingdom plot: Figure 8.

$$\text{II} \quad \vdash \quad f \quad \vdash \quad \text{D} \quad \vdash \quad \vdash \quad f \quad \vdash \quad \vdash \quad 1 \quad \vdash \quad 1 \quad \vdash \quad \vdash \quad \vdash \quad 1 \quad \vdash \quad 1 \quad \vdash \quad 1 \quad \vdash \quad 1 \quad \vdash \quad \text{E} \quad \vdash \quad 0$$

Archaeas AntiSMASH

Taxonomical diversity on Archaeasbacteria Data

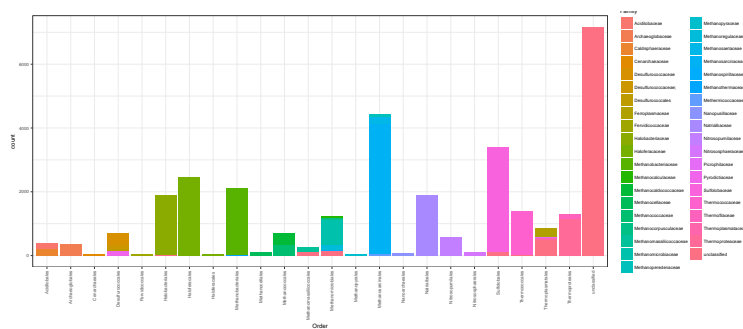


Figure 10: Archaeas Diversity

Here is a reference to Recruitments after central pathways expansions coloured by taxa plot: Figure 10.

Smash diversity

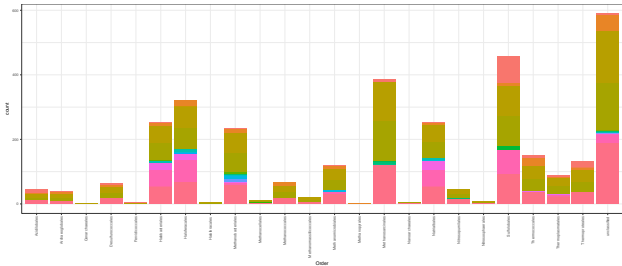


Figure 11: Archaeas Smash Taxonomical Diversity

Here is a reference to Recruitments after central pathways expansions coloured by taxa plot: Figure 11.

Total central pathway expansions by genome vs Total antimash cluster detected splitted by order

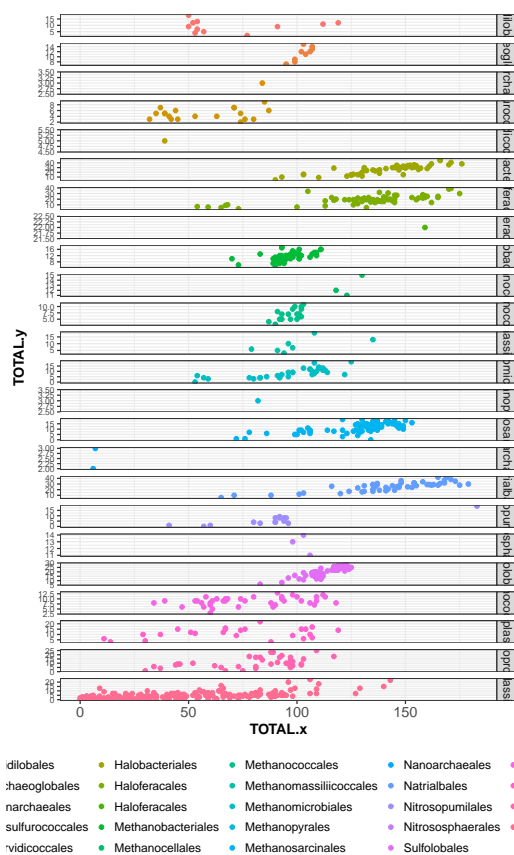


Figure 13: Correlation between Archaeas central pathway expansions and antimash Natural products detection

Here is a reference to the expansions vs antimash NP's clusters splitted by order plot Figure 13.

AntisMash vs Expansions by taxonomic Family

Natural products colored by family

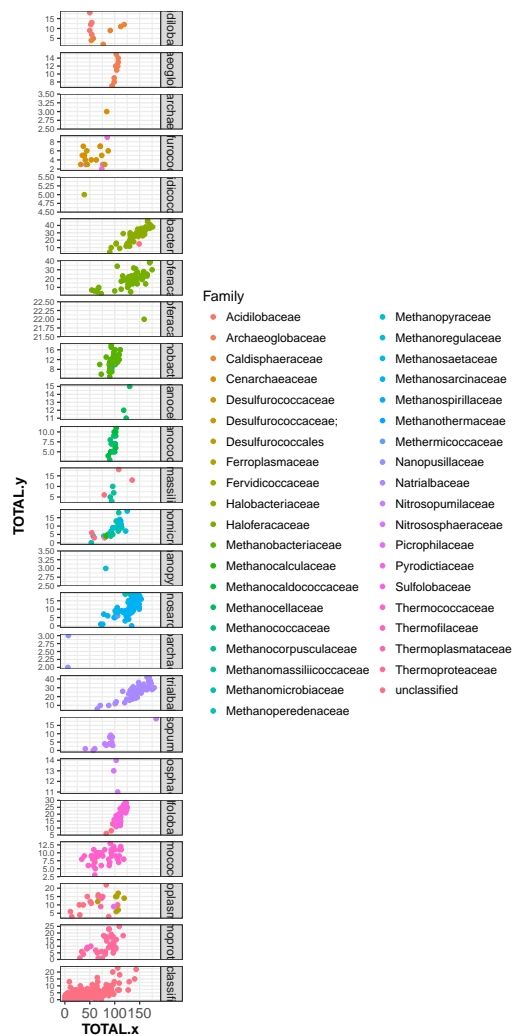


Figure 14: Archaeas Natural products by family

Here is a reference to the Natural products colored by family plot Figure 14.

Selected trees from EvoMining

Phosphoribosyl_isomerase_3 family
Figure from EvoMining



Figure 15: Phosphoribosyl isomerase A EvoMiningtree

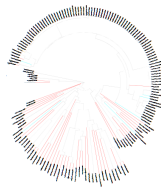


Figure 16: Phosphoribosyl isomerase other EvoMiningtree

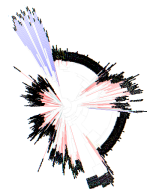


Figure 17: Phosphoribosyl anthranilate isomerase EvoMiningtree

Footnotes and Endnotes

You might want to footnote something.¹ The footnote will be in a smaller font and placed appropriately. Endnotes work in much the same way. More information can be found about both on the CUS site or feel free to reach out to data@reed.edu.

Bibliographies

Of course you will need to cite things, and you will probably accumulate an armful of sources. There are a variety of tools available for creating a bibliography database (stored with the .bib extension). In addition to BibTeX suggested below, you may want to consider using the free and easy-to-use tool called Zotero. The Reed librarians have created Zotero documentation at <http://libguides.reed.edu/citation/zotero>. In addition, a tutorial is available from Middlebury College at <http://sites.middlebury.edu/zoteromiddlebury/>.

R Markdown uses *pandoc* (<http://pandoc.org/>) to build its bibliographies. One nice caveat of this is that you won't have to do a second compile to load in references as standard L^AT_EX requires. To cite references in your thesis (after creating your bibliography database), place the reference name inside square brackets and precede it by the "at" symbol. For example, here's a reference to a book about worrying: [Molina1994]. This Molina1994 entry appears in a file called `thesis.bib` in the `bib` folder. This bibliography database file was created by a program called BibTeX. You can call this file something else if you like (look at the YAML header in the main .Rmd file) and, by default, is placed in the `bib` folder.

For more information about BibTeX and bibliographies, see our CUS site (<http://web.reed.edu/cis/help/latex/index.html>)². There are three pages on this topic: *bibtex* (which talks about using BibTeX, at <http://web.reed.edu/cis/help/latex/bibtex.html>), *bibtexstyles* (about how to find and use the bibliography style that best suits your needs, at <http://web.reed.edu/cis/help/latex/bibtexstyles.html>) and *bibman* (which covers how to make and maintain a bibliography by hand, without BibTeX, at <http://web.reed.edu/cis/help/latex/bibman.html>). The last page will not be useful unless you have only a few sources.

If you look at the YAML header at the top of the main .Rmd file you can see that we can specify the style of the bibliography by referencing the appropriate csl file. You can download a variety of different style files at <https://www.zotero.org/styles>. Make sure to download the file into the `csl` folder.

Tips for Bibliographies

- Like with thesis formatting, the sooner you start compiling your bibliography for something as large as thesis, the better. Typing in source after source is mind-numbing enough; do you really want to do it for hours on end in late April? Think of it as procrastination.
- The cite key (a citation's label) needs to be unique from the other entries.
- When you have more than one author or editor, you need to separate each author's name by the word "and" e.g. `Author = {Noble, Sam and Youngberg, Jessica},.`
- Bibliographies made using BibTeX (whether manually or using a manager) accept L^AT_EX markup, so you can italicize and add symbols as necessary.
- To force capitalization in an article title or where all lowercase is generally used, bracket the capital letter in curly braces.
- You can add a Reed Thesis citation³ option. The best way to do this is to use the `phdthesis` type of citation, and use the optional "type" field to enter "Reed thesis" or "Undergraduate thesis."

¹footnote text

²@reedweb2007

³@noble2002

Anything else?

If you'd like to see examples of other things in this template, please contact the Data @ Reed team (email data@reed.edu) with your suggestions. We love to see people using *R Markdown* for their theses, and are happy to help.