

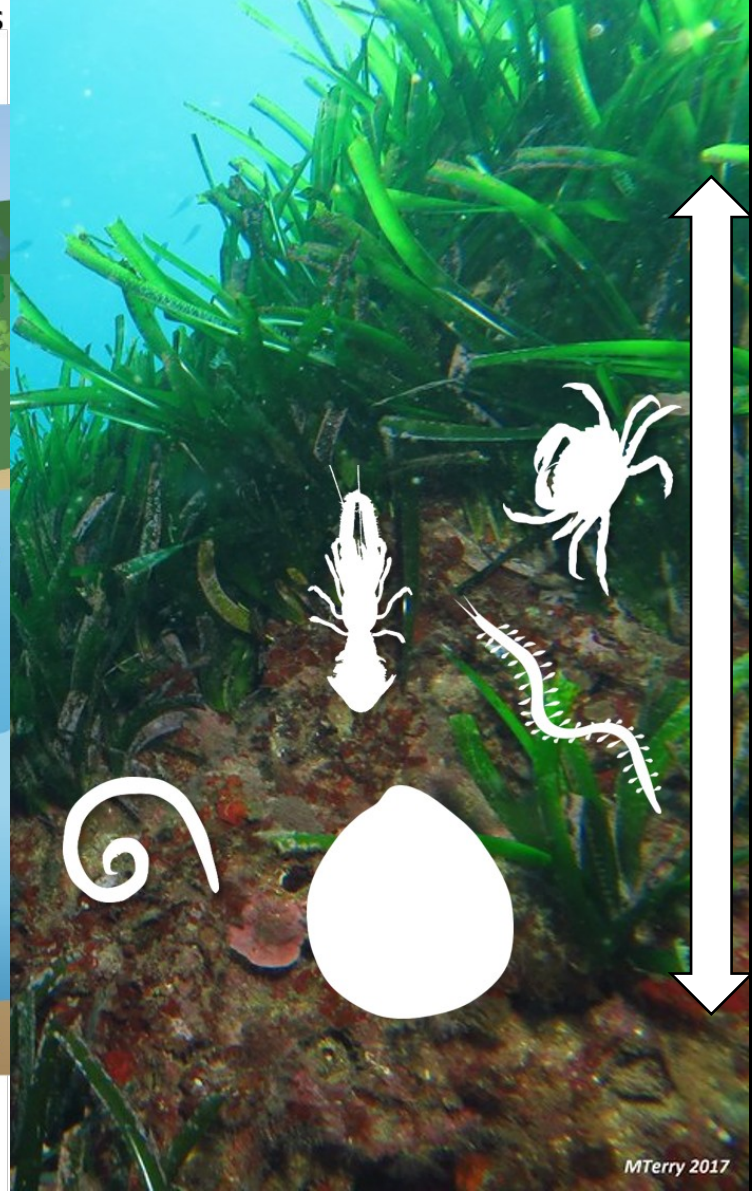
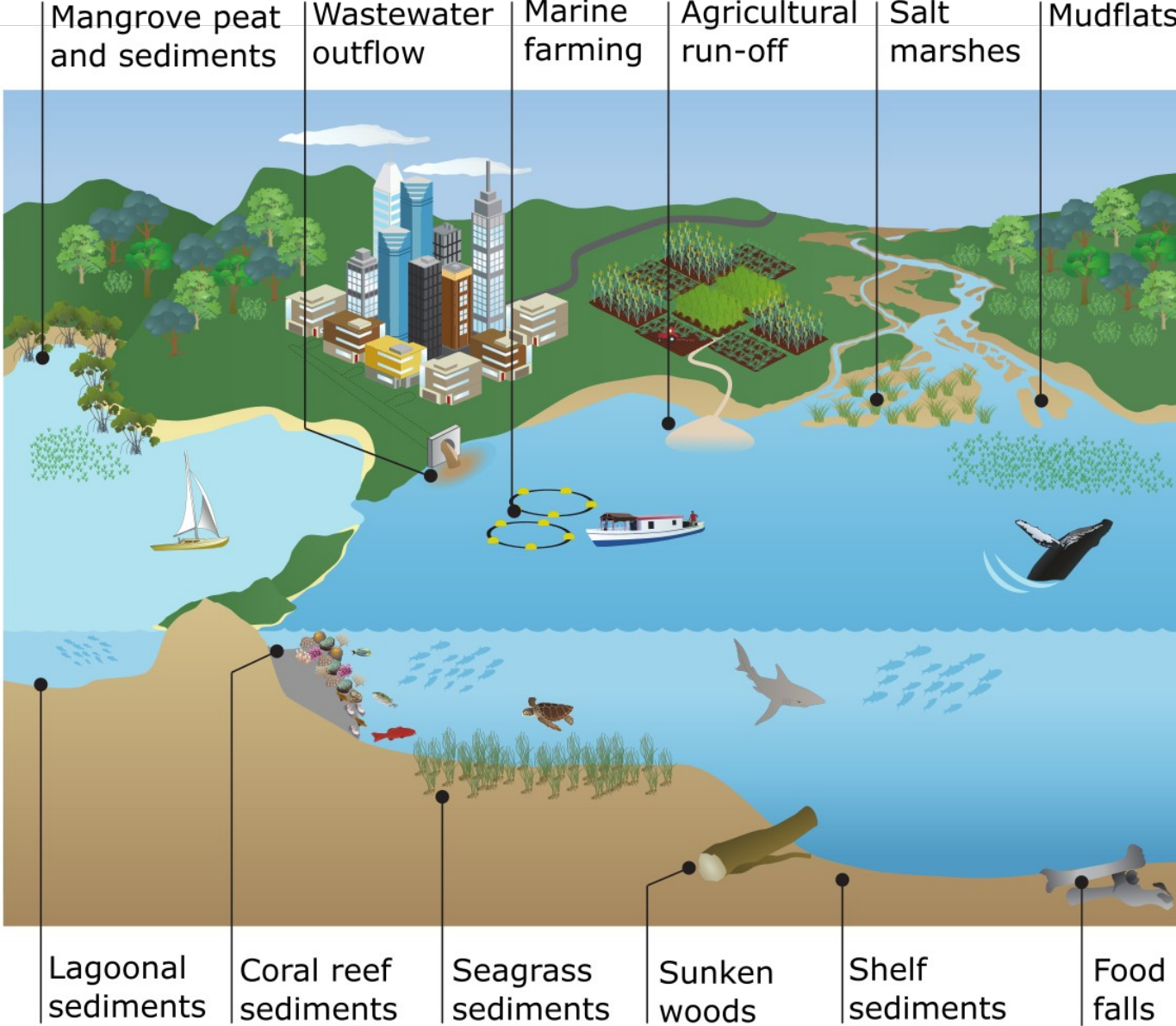
Building a network of invertebrate-microbe associations from seagrass beds

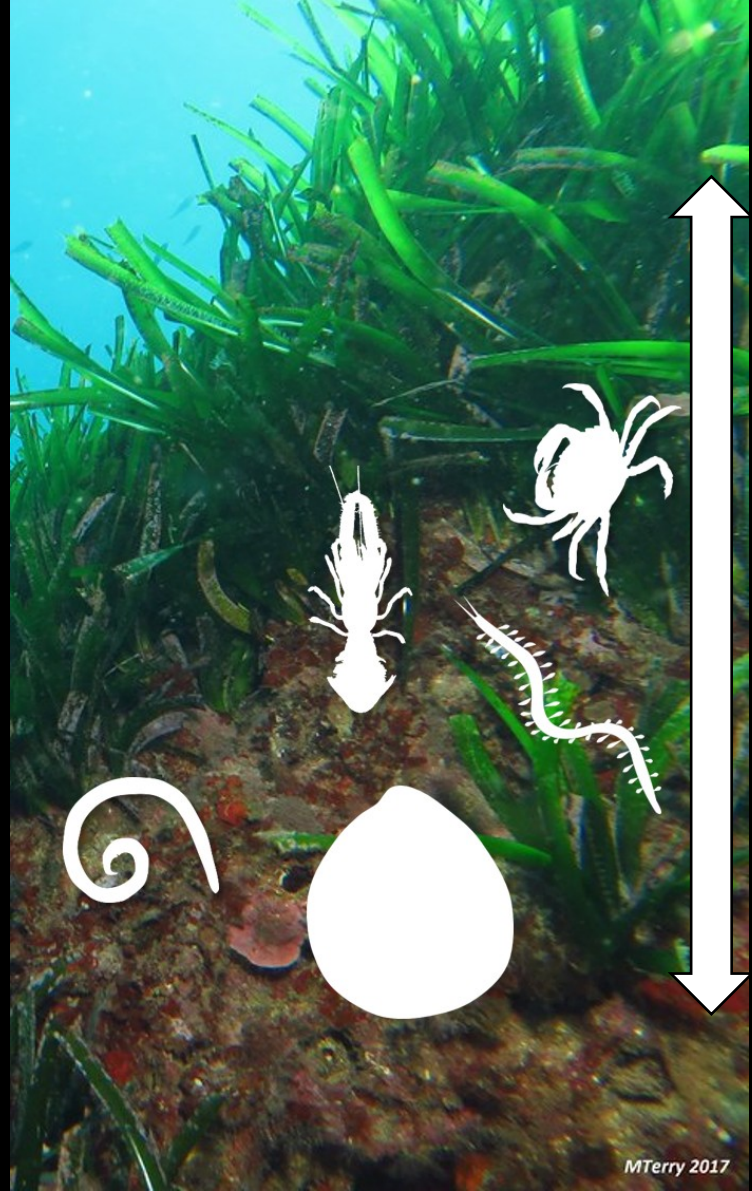
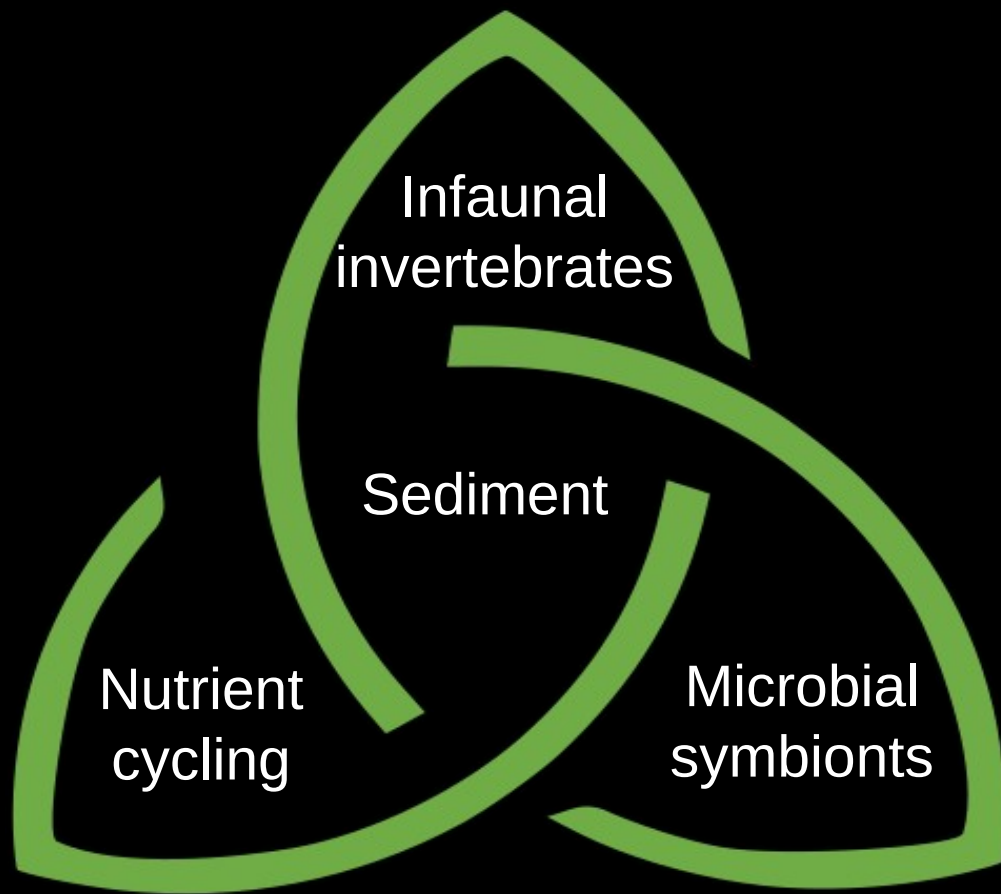
Luis Montilla

Roberta Piredda

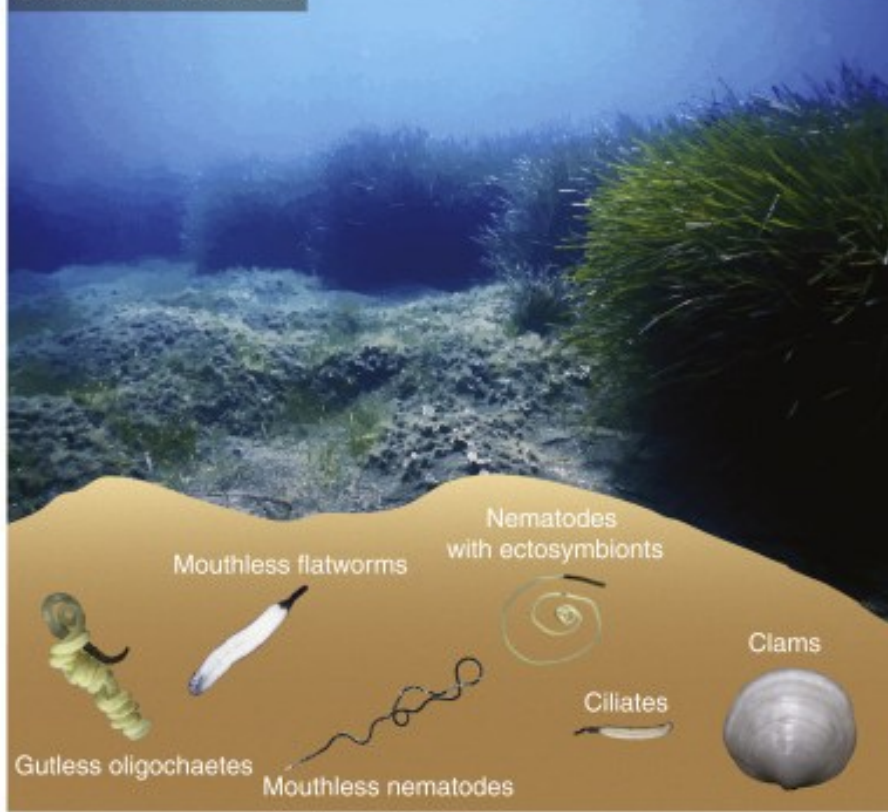
Ulisse Cardini







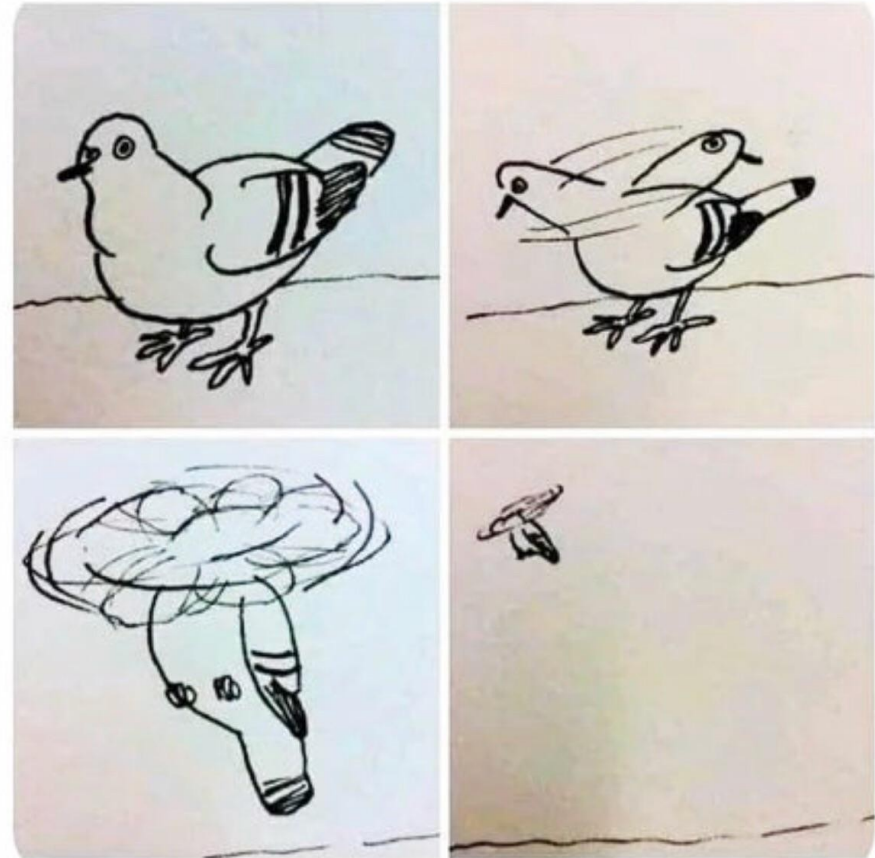
Shallow-water sediments



Current Biology

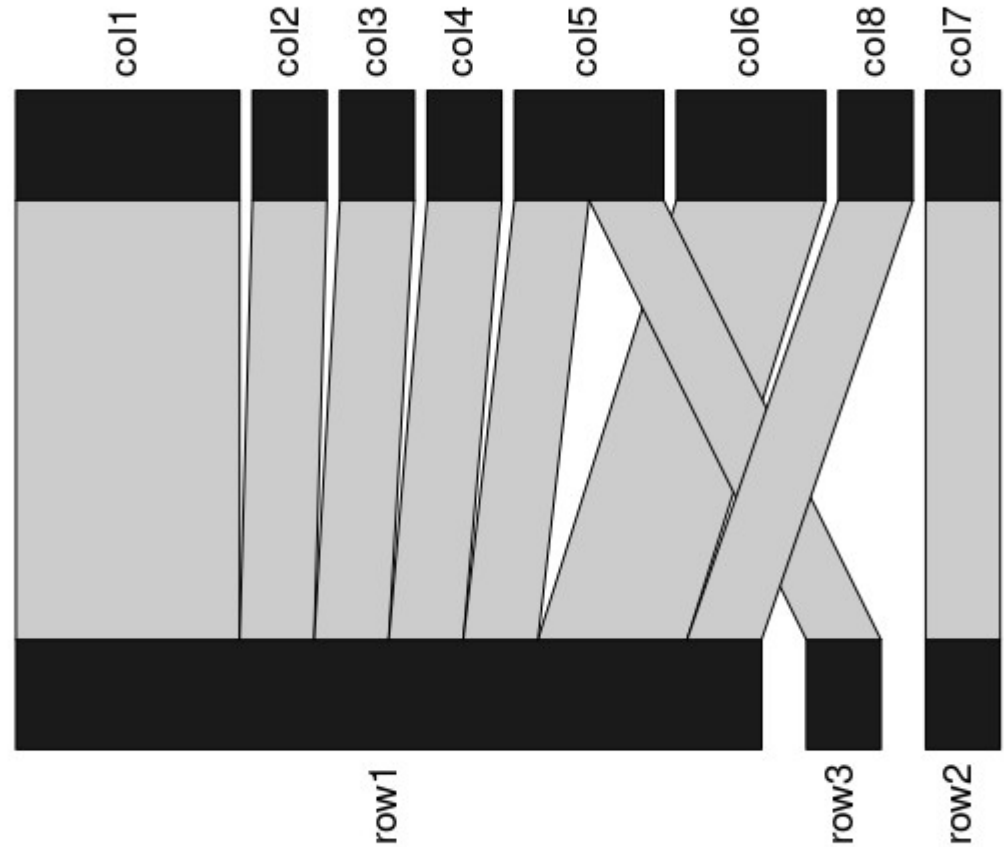
What patterns would we see in a more extensive data set?

When your program is a complete mess, but it does its job



Network of associations
between seagrass
macrofauna and microbial
symbionts

The first ingredient: a list of
291 dominant invertebrates
from Mediterranean
seagrasses



Search query: Genus + one group of symbionts (e.g. Bacteria)

Items: 1 to 20 of 177

<< First < Prev Page 1 of 9 Next > Last >>

☐ [Loripes lacteus gill symbiont clone 1B 16S ribosomal RNA gene, partial sequence](#)

1. 1,527 bp linear DNA

Accession: GQ853556.1 GI: 261343262

[Taxonomy](#)

[GenBank](#) [FASTA](#) [Graphics](#)

☐ [Loripes lacteus gill symbiont clone 1A 16S ribosomal RNA gene, partial sequence](#)

2. 1,569 bp linear DNA

Accession: GQ853555.1 GI: 261343261

[Taxonomy](#)

[GenBank](#) [FASTA](#) [Graphics](#)

☐ [Loripes lacteus gill symbiont clone 2C 16S ribosomal RNA gene, partial sequence](#)

3. 1,503 bp linear DNA

Accession: FJ752447.1 GI: 261208344

[Taxonomy](#)

[GenBank](#) [FASTA](#) [Graphics](#)

☐ [Loripes lacteus gill symbiont clone 1C 16S ribosomal RNA gene, partial sequence](#)

4. 1,538 bp linear DNA

Accession: FJ752446.1 GI: 261208343

[Taxonomy](#)

[GenBank](#) [FASTA](#) [Graphics](#)

Find items

Search details

("Loripes"[Organism] OR loripes[All Fields]) AND ("Bacteria"[Organism] OR "Bacteria Latreille et al. 1825"[Organism] OR bacteria[All Fields])

Search

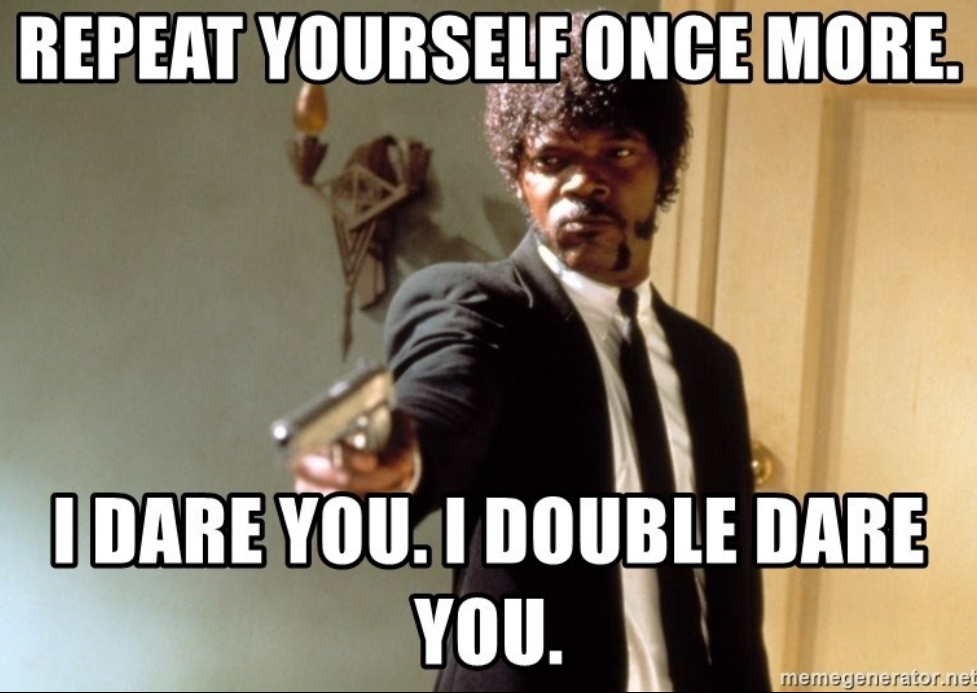
See more...

Recent activity

[Turn Off](#) [Clear](#)

Your browsing activity is empty.

Manually download the records
for 291 genera of invertebrates
(?)



rOpenSci: The *rentrez* package

`rentrez` provides functions that work with the [NCBI Eutils](#) API to search, download data from, and otherwise interact with NCBI databases.



The first setback:

We needed a specific
qualifier

GenBank ▾

Loripes lacteus gill symbiont clone 1B 16S ribosomal RNA gene,

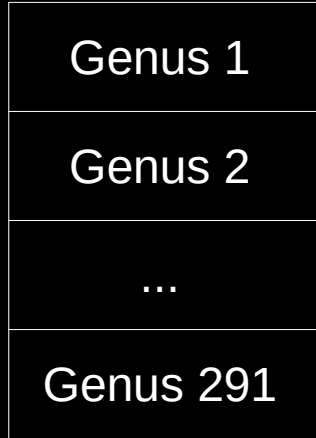
GenBank: GQ853556.1

[FASTA](#) [Graphics](#)

Go to: ☐

LOCUS	GQ853556	1527 bp	DNA	linear	BCT 31-AUG-2010
DEFINITION	Loripes lacteus gill symbiont clone 1B 16S ribosomal RNA gene, partial sequence.				
ACCESSION	GQ853556				
VERSION	GQ853556.1				
KEYWORDS	.				
SOURCE	Loripes lacteus gill symbiont				
ORGANISM	Loripes lacteus gill symbiont Bacteria; Proteobacteria; Gammaproteobacteria; sulfur-oxidizing symbionts.				
REFERENCE	1 (bases 1 to 1527)				
AUTHORS	Mausz,M., Schmitz-Esser,S. and Steiner,G.				
TITLE	Identification and comparative analysis of the endosymbionts of Loripes lacteus and Anodontia fragilis (Bivalvia: Lucinidae)				
JOURNAL	Unpublished				
REFERENCE	2 (bases 1 to 1527)				
AUTHORS	Mausz,M., Schmitz-Esser,S. and Steiner,G.				
TITLE	Direct Submission				
JOURNAL	Submitted (26-AUG-2009) Department for Evolutionary Biology, University of Vienna, Institute for Organismal Systems Biology, Althanstrasse 14, Vienna A-1040, Austria				
FEATURES	Location/Qualifiers				
source	1..1527 /organism="Loripes lacteus gill symbiont" /mol_type="genomic DNA" /host="Loripes lacteus isolate 52" /clone="1B" /product="16S ribosomal RNA"				
rRNA	<1..>1527				
ORIGIN	1 agagtttgat catggctcag attgaacgct ggcggaggcc taacacatgc aagtcgagcg 61 gtaacagggg gagcttgctc tctgctgacg agcggcggac ggggtgcgtaa cacgtaggaa				

Retrieve genbank files



Store data

Write a loop - > Repeat this step for each 291 genera x 5 groups (bacteria, archaea, virus, fungi, eukaryota)

```
> gb_file
Object of class 'efetch'
LOCUS      GQ853556                      1527 bp    DNA      linear    BCT 31-AUG-2010
DEFINITION Loripes lacteus gill symbiont clone 1B 16S ribosomal RNA gene,
            partial sequence.
ACCESSION  GQ853556
VERSION    GQ853556.1
KEYWORDS   .
SOURCE     Loripes lacteus gill symbiont
  ORGANISM Loripes lacteus gill symbiont
            Bacteria; Proteobacteria; Gammaproteobacteria; sulfur-oxidizing
            symbionts.
REFERENCE  1 (bases 1 to 1527)
  AUTHORS  Mausz,M., Schmitz-Esser,S. and Steiner,G.
...
EFetch query using the 'nuccore' database.
Query url: 'https://eutils.ncbi.nlm.nih.gov/entrez/eutils/efetch.fcgi?efetch...'
Retrieval type: 'gbwithparts', retrieval mode: 'text'

> |
```

Retrieve genbank files

Identify objects with 0 entries, and objects
with >900 entries

You were not in the database

I am, with 0 entries

210 genera didn't
return any results,
leaving us with only
28% of the original
list.



Retrieve genbank files

Write files to disk

Also, manually downloading files with >900 entries

The specific element we wanted isn't easily retrievable

Also, the gb objects include genetic sequences

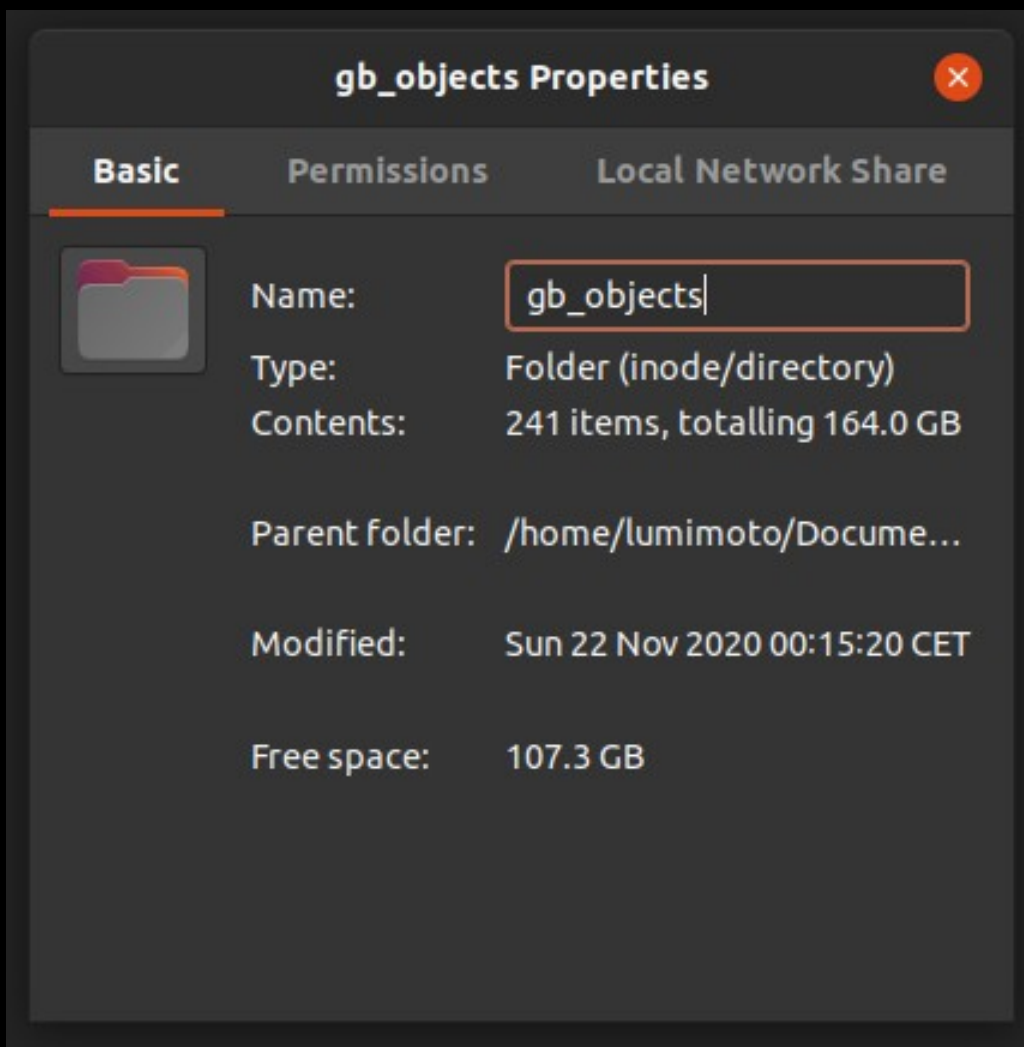
(maybe using regex)

```

gb_file$content
[1] "LOCUS      GQ853556      1527 bp      DNA      linear      BCT 31-AUG-2010\ndefinition  Loripes lacteus gill symbiont clone 1B 16S ribosomal
RNA gene,\n      partial sequence.\naccession      GQ853556.1\nversion      GQ853556.1\nkeywords      \nsource      Loripes lacteus gill symbiont\nORIGIN      Bacteria; Proteobacteria; Gammaproteobacteria; sulfur-oxidizing\n      symbionts.\nREFERENCE
1 (bases 1 to 1527)\nauthors      Mausz,M., Schmitz-Esser,S. and Steiner,G.\n      title      Identification and comparative analysis of the endosymbionts
of\n      Loripes lacteus and Anodontia fragilis (Bivalvia: Lucinidae)\n      journal      Unpublished\nREFERENCE      2 (bases 1 to 1527)\nauthors      Mausz,M., Schmitz-Esser,S. and Steiner,G.\n      title      Direct Submission\n      journal      Submitted (26-AUG-2009) Department for Evolutionary Biology,\n      University of Vienna, Institute for Organismal Systems Biology,\n      Althanstrasse 14, Vienna A-3492, Austria\nfeatures      Loc
ation/Qualifiers\n      source      1..1527\n      /organism=\"Loripes lacteus gill symbiont\"\n      /mol_type=\"geno
mic DNA\"\n      /host=\"Loripes lacteus isolate 52\"\n      /db_xref=\"taxon:682670\"\n      /clone=\"1B
\"\n      rRNA      <1..>\n      /rRNA=\"16S ribosomal RNA\"\nORIGIN      \n      1 agagtttgat catggctcag attgaacgct ggcgg
aggcc taacacatgc aagtcgagcg\n      61 gtaacagggg gagcttgctc tctgctgacg agcgccggac gggcgctgaa cagctaggaa\n      121 tctgcccggt agtgggggat agccccgaga aa
tccggatt aattaccgat acgcccaga\n      181 ggggaaagca ggggattcgt tctttcga gattggacct tgcgtattg gatgacgctg\n      241 cgctcgatta gcttctgttg ggggtaaagg
ctccaagaag caacgctgcg tagctgctct\n      301 gaggagatga tcagccacac tgggactgag acagcgccca gactctcag gaggcgagca\n      361 gtgggggaata tgcacaattg ggggaaa
ccc tgatgcagcc atgcccgcgtg tgtgaagaag\n      421 gctctagggt tgtaaaagcac ttacgacgag gaggaaaagg ttgtgattaa tactcaacag\n      481 ctgtgacgtt actcgcagaa gaag
caccgg ctaactcgt gccagcagcc gcggaataac\n      541 ggagggtgcg agcgtaatac ggaattactg ggcgtaaac gtcgtaggc ggctgcctaa\n      601 gttggatgtg aaagccccgg g
cttaacctg ggaactgcat ccaaaactgg gcggctagag\n      661 tgcggaagag gagtgtgtaa tttcctgtgt agcggtgaaa tgcgtagata taggaaggaa\n      721 caccagtggc gaaggcgac
a cctctggtctg acactgacgc tgaggtacga aagcgtggg\n      781 agcaaacagg attagatacc ctggtagtc acgccgtaaa cgatgctact tagtgtctcg\n      841 gagtcttgta tctctg
gtaa cgagactaac gcgataagta gacgcctgg ggagtagcgc\n      901 cgcaaggtta aactcaaat gaattgacgg gggccgcac aagcggtgga cgatgtggtt\n      961 taattcgaag caa
cgcaag aaccttacct ggccttgaca tcttgcaaat cctttagaga\n      1021 tagaggagtg ccttcgggaa cgcagagaca ggtgctgcat ggctgtcgtc agctcgtgtc\n      1081 gtgagatgtt
gggttaagtc cgcgaacgag gcgaaccctt gtctcagtt accagcacgt\n      1141 tatggtgggc actctgggga gactgccggt gacaaacgg aggaagggtg ggaacgacgtc\n      1201 aagtcat
cat ggccttacg gcaggggcta cacacgtgct acaatggtgc atacagacgg\n      1261 ttgccaagcc cgagggtgga gctaactgta gaaagtgcat cgtagtcgg attggagtct\n      1321 gcaa
ctcgac tcatatgag cgaatcgtc agtaactgtg aatcagaagt tcacggtgaa\n      1381 tagtctccc ggctttgata acacgccccc tcaaccattg ggagtggtgt gctccagaag\n      1441 t
gttagctct aacctctctt ttcgaagggg gaaggcgaat caccacgga gtattcatga\n      1501 ctggggtgaa gtcataacaa gtagcccn\n      /n\n"
```


Retrieve genbank files

Write files to disk



Retrieve genbank files

Write files to disk

Merge & transform

```
> gb_file
Object of class 'efetch'
LOCUS      GQ853556                1527 bp    DNA        linear    BCT 31-AUG-2010
DEFINITION Loripes lacteus gill symbiont clone 1B 16S ribosomal RNA gene,
            partial sequence.
ACCESSION  GQ853556
VERSION    GQ853556.1
KEYWORDS   .
SOURCE     Loripes lacteus gill symbiont
            ORGANISM  Loripes lacteus gill symbiont
                        Bacteria; Proteobacteria; Gammaproteobacteria; sulfur-oxidizing
                        symbionts.
REFERENCE  1 (bases 1 to 1527)
AUTHORS    Mausz,M., Schmitz-Esser,S. and Steiner,G.
...
EFetch query using the 'nuccore' database.
Query url: 'https://eutils.ncbi.nlm.nih.gov/entrez/eutils/efetch.fcgi?efetch...'
Retrieval type: 'gbwithparts', retrieval mode: 'text'

> |
```

	Taxon 1	Taxon 2	Taxon 3
Sample 1			
Sample 2			

Retrieve genbank files

The easiest choice was to use a third-party software

Write files to disk

Gbk2fas (Göker et al. 2010)

Merge & transform

<http://www.goeker.org/mg/clustering/>

Auxiliary programs

- [gbk2fas](#): Like the old version, but many more options. 27 distribution files. Lists with product names can be created, as well as [m4](#) files.
- [gbk2fas, old version](#): use this program to convert sequence data. It can be used to (1) adapt the FASTA headers to each user's needs; (2) write the headers in alignments, tree files, etc. later on by full name. For example, reference partitions for clustering optimization with [Göker et al. 2010](#).

Retrieve genbank files

Write files to disk

Merge & transform

Auxiliary programs

- [gbk2fas](#): Like the old version, but many more options. 27 distribution files. Lists with product names can be created, as well as [m4](#) macros.
- [gbk2fas, old version](#): use this program to convert sequence data from one to (1) adapt the FASTA headers to each user's needs; (2) write the headers in alignments, tree files, etc. later on by full name and instance, reference partitions for clustering optimization with [Göker et al. 2010](#).



Most recent version
(Linux)



Old version (Linux)



New version
(Windows)

Retrieve genbank files

Write files to disk

Merge & transform

Reimport to R

The screenshot displays the RStudio environment with the following components:

- Source Panel:** Contains R code for reading and specifying columns of two datasets. The first dataset, `genera_master`, is read from `output/host_sym_raw.txt` using `read_delim()` with tab delimiters. The second dataset, `host_attr`, is read from `data/genera_attr.csv` using `read_csv()`. Both sections include a `cols()` function call to specify column types: `col_character()` for text and `col_double()` for numerical data.
- Console Panel:** Shows the execution progress of the code, including a completion message: `|=====| 100% 2061 MB`.
- Environment Panel:** Lists the objects in the global environment:
 - `genera_ma...`: 1750892 observations of 24 variables.
 - `host_attr`: 290 observations of 10 variables.

Retrieve genbank files

Write files to disk

Merge & transform

Reimport to R

Data cleaning (hosts)

- 1) Non-marine organisms
- 2) Corals, sponges, algae

worms

repo status **Active** CRAN **OK** R-check **passing** codecov **84%** downloads **2854/month** CRAN **0.4.2**

`worms` is a R client for the World Register of Marine Species

- World Register of Marine Species (WoRMS) <http://www.marinespecies.org/>
- WoRMS REST API docs: <http://www.marinespecies.org/rest/>

See the taxize book (<https://taxize.dev>) for taxonomically focused work in this and similar packages.

“Host unknown” = plants and vertebrates

Impostor



Retrieve genbank files

Write files to disk

Merge & transform

Reimport to R

Data cleaning (hosts)

Data cleaning (symbionts)



Lowest common taxonomic resolution

Raw file: 1,750,892 rows

Cleaned data set: 39,036 rows (2,2%)



igraph – The network analysis package

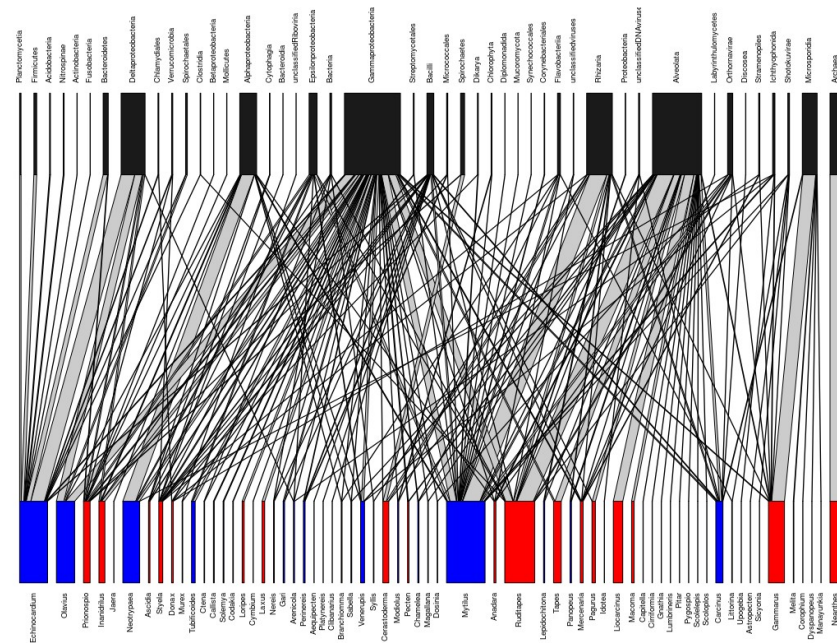
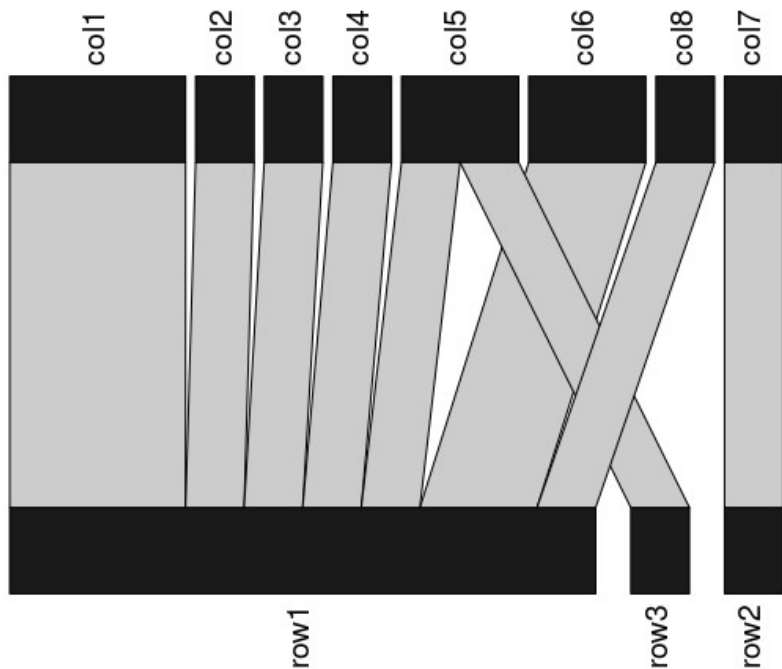
igraph is a collection of network analysis tools with the emphasis on efficiency, portability and ease of use. igraph is open source and free. igraph can be programmed in R, Python, Mathematica and C/C++.

igraph R package

python-igraph

IGraph/M

igraph C library



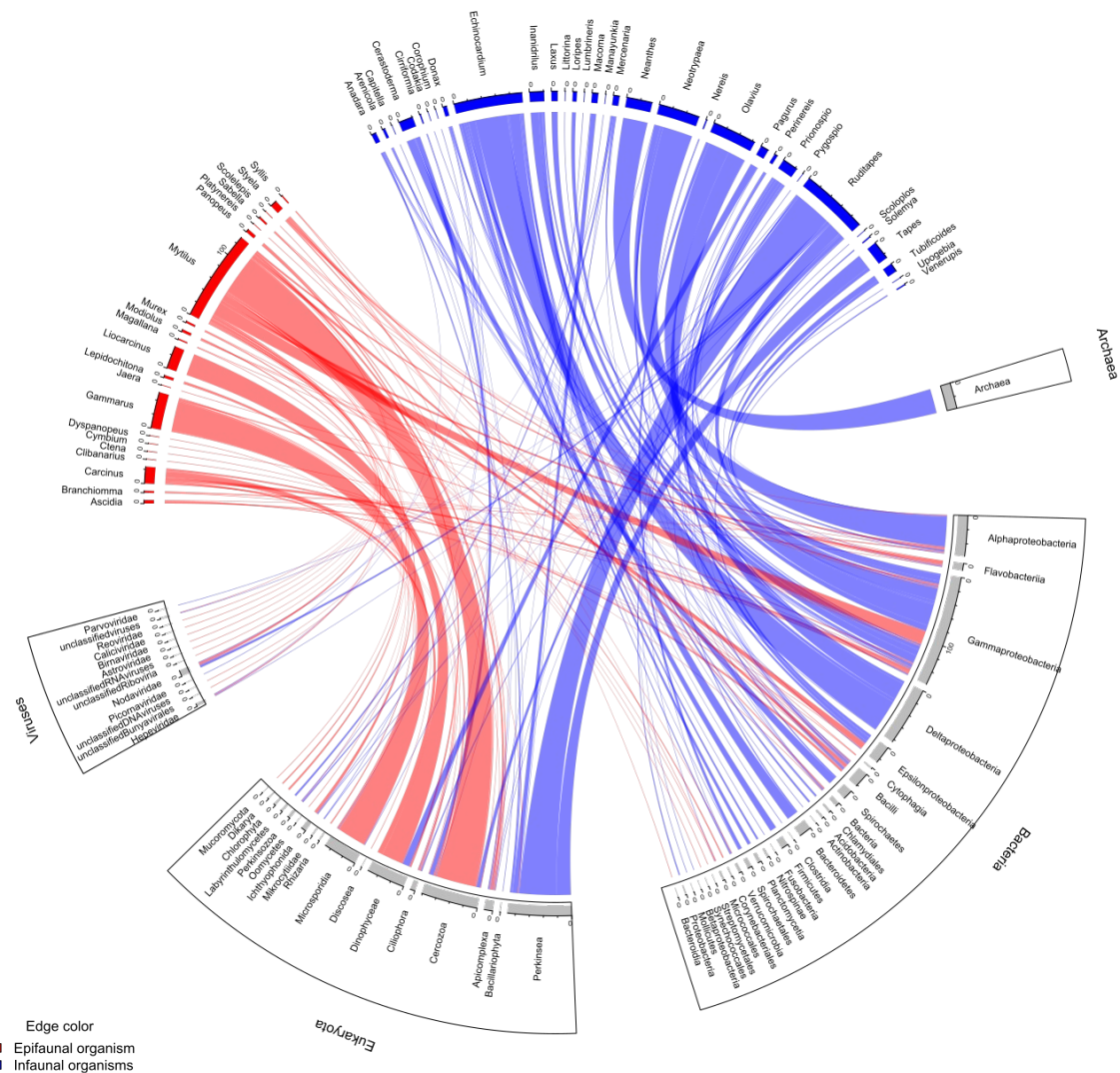
112 nodes	53 invertebrate genera	21 epifaunal
		32 infaunal
	59 microbial groups	28 Bacteria
		17 Eukaryota
		13 viruses
		1 Archea

circlize: circular visualization in R

build passing CRAN 0.4.12 downloads 780K

Circular layout is an efficient way for the visualization of huge amounts of information. Here the circlize package provides an implementation of circular layout generation in R as well as an enhancement of available software. The flexibility of this package is based on the usage of low-level graphics functions such that self-defined high-level graphics can be easily implemented by users for specific purposes. Together with the seamless connection between the powerful computational and visual environment in R, circlize gives users more convenience and freedom to design figures for better understanding complex patterns behind multi-dimensional data.

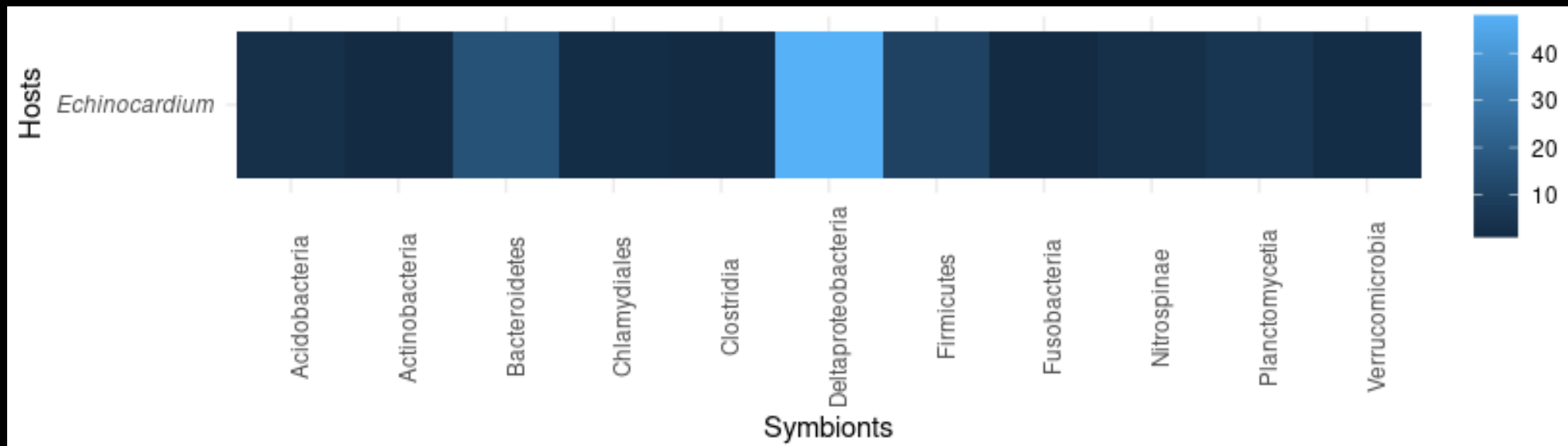
112 nodes	53 invertebrate genera	21 epifaunal
		32 infaunal
	59 microbial groups	28 Bacteria
		17 Eukaryota
		13 viruses
		1 Archaea

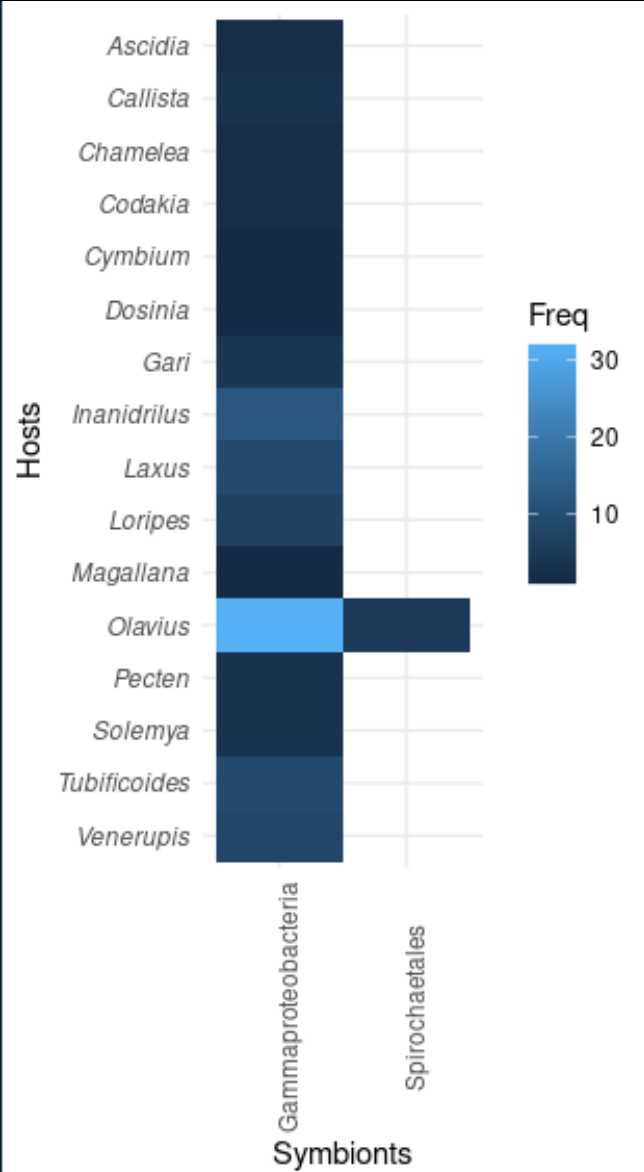


Examples of obtained communities



© Hans Hillewaert





Next steps

Expand list to include macrofauna
from other marine ecoregions

Add additional attributes if possible

Optimize code



build error build failing downloads 103/month CRAN not published

biofiles - an interface to GenBank/GenPept files in R

biofiles provides interfacing to GenBank/GenPept or Embl flat file records. It includes utilities for reading and writing GenBank files, and methods for interacting with annotation and sequence data.

Thanks for your attention
Questions?

Luismmontilla.com
luismiguel.montilla@szn.it

