

Metagenomes in  
practice

win a prize !



# Assemblers: spades, megahit, velvet

```
spades.py --meta \  
--pe1-1 ./data/metanom_1.fastq \  
--pe1-2 ./data/metanom_2.fastq \  
-k 21,33 --only-assembler \  
-o spades_assembly --memory 4 --threads 4
```



# Metagenomic pipelines

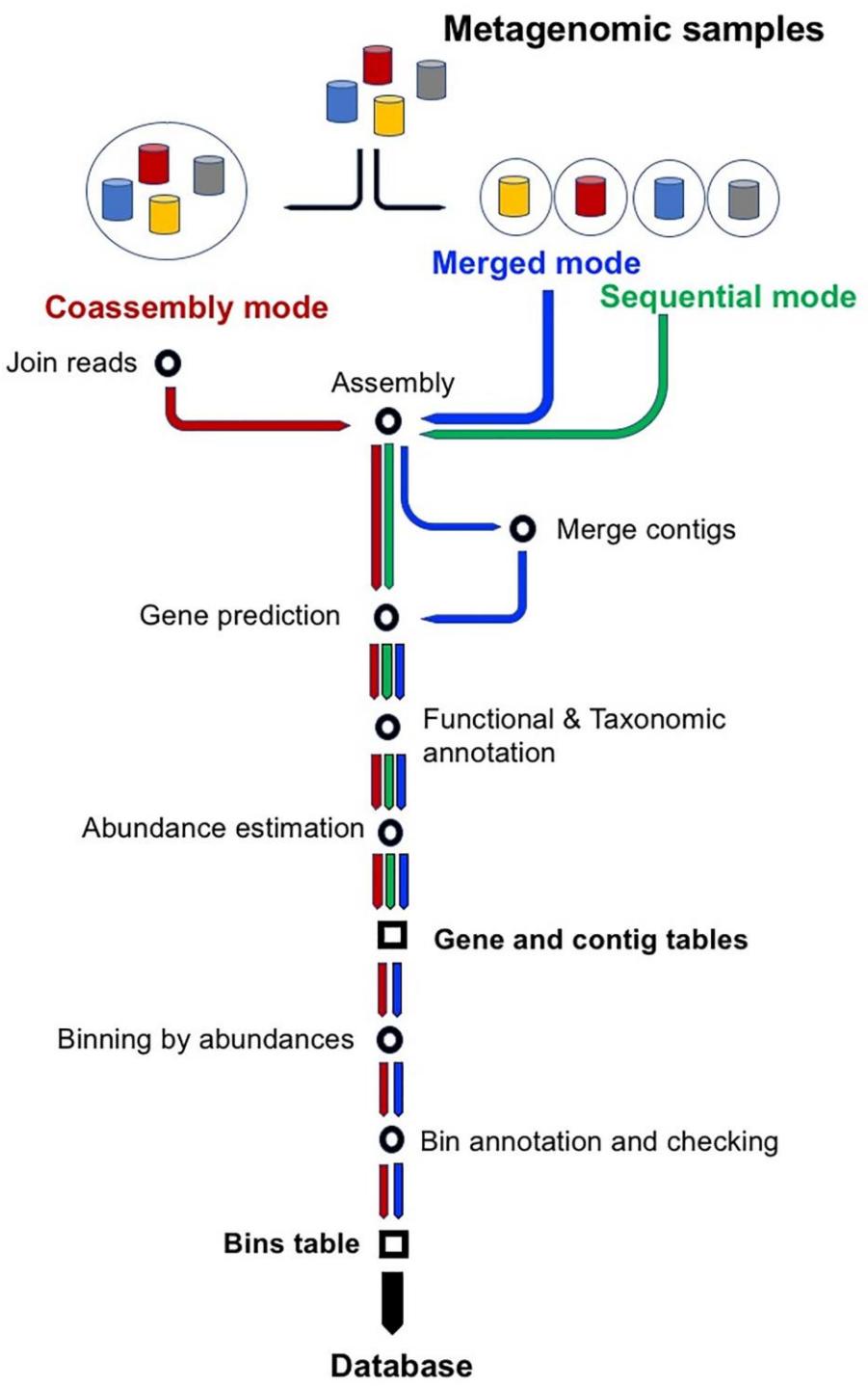
They all use the same tools, the assemblers , the same methods, usually adaptable  
Some are easier some harder but they all need cores and memory.

- Quality check trimming
- Assembly
- Normalization
- Annotation
- Binning
- Other functionalities

Pipelines
SqueezeMeta
Atlas
nf-core
bioBakery4
JAMS
WGSA2

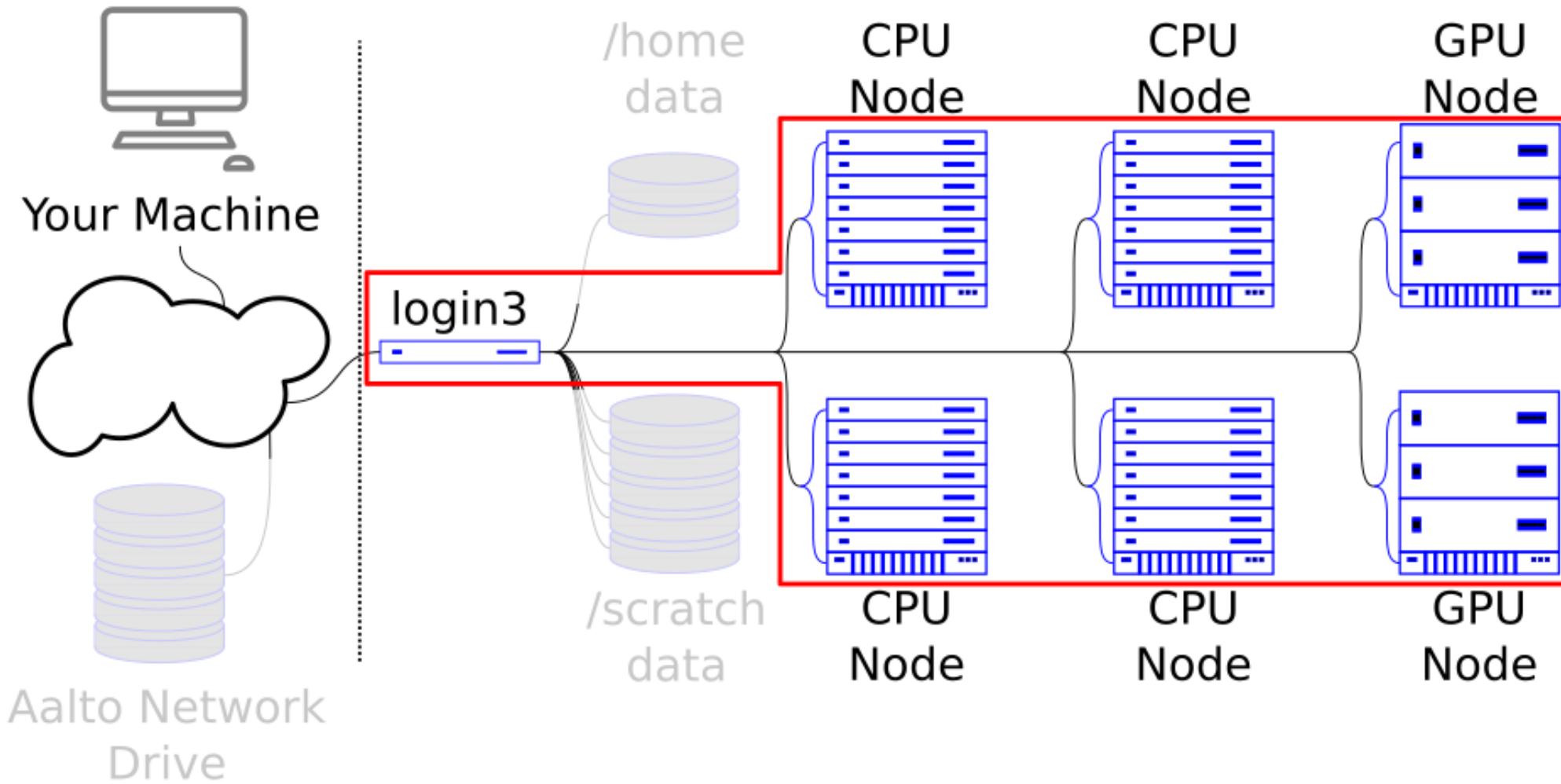
# SqueezeMeta

	<b>MG-Rast (Meyer et al., 2008)</b>	<b>Anvio (Eren et al., 2015)</b>	<b>Smash community (Arumugam et al., 2010)</b>	<b>Humann (Abubucker et al., 2012)</b>	<b>fmap (Kim et al., 2016)</b>	<b>MetaWrap (Uritskiy et al., 2018)</b>	<b>Samsa2 (Westreich et al., 2018)</b>	<b>IMP (Narayanasamy et al., 2016)</b>	<b>SqueezeMeta</b>
Assembly	No	No	Yes	No	No	Yes	No	Yes	Yes
Data source	Reads or contigs	Contigs	Contigs	Reads	Reads or contigs	Contigs	Reads (RNA)	Reads	Reads
Gene prediction	Yes	Yes	Yes	No	No	No	No	Yes	Yes
Function assignment	Yes	Yes	Yes	Yes	Yes	No	Yes	Yes	Yes
RNA assignment	Yes	Yes	No	No	No	No	Yes	Yes	Yes
Taxonomic assignment	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Gene abundances	Yes	Yes	No	Yes	Yes	No	Yes	Yes	Yes
Metagomic comparison	Yes	Yes	Yes	Yes	Yes	No	Yes	Yes	Yes
Co-assembly	No	No	No	No	No	Yes	No	Yes	Yes
Binning	No	Support	No	No	No	Yes	No	Yes	Yes
Bin validation	No	Yes	No	No	No	No	No	No	Yes
Local Installation	No	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes



1. Assembly
2. RNA prediction and classification
3. ORF (CDS) prediction
4. Homology searching against taxonomic and functional databases
5. Hmmer searching against Pfam database
6. Taxonomic assignment of genes
7. Functional assignment of genes
8. Blastx on parts of the contigs with no gene prediction or no hits
9. Taxonomic assignment of contigs, and check for taxonomic disparities
10. Coverage and abundance estimation for genes and contigs
11. Estimation of taxa abundances
12. Estimation of function abundances
13. Merging of previous results to obtain the ORF table
14. Binning with different methods
15. Binning integration with DAS tool
16. Taxonomic assignment of bins, and check for taxonomic disparities
17. Checking of bins with CheckM
18. Merging of previous results to obtain the bin table
19. Merging of previous results to obtain the contig table
20. Prediction of kegg and metacyc pathways for each bin
21. Final statistics for the run

# Running Squeezemeta in SLURM



# Running Sqeezemeta in SLURM

To run the pipeline in basic mode you need to download:

**397G** ./db

of databases

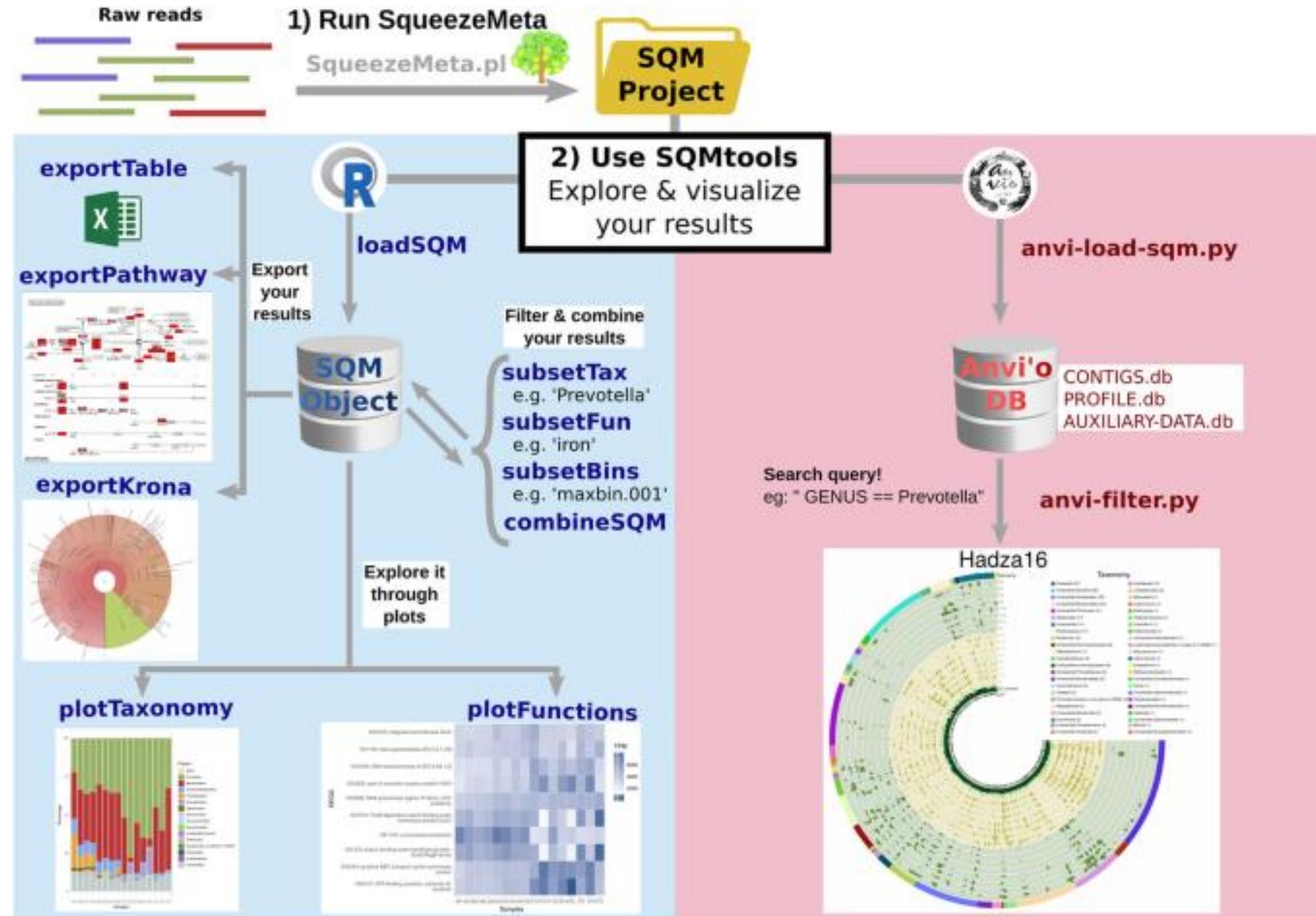
- `#!/bin/bash`
- `#SBATCH -J squeeze.sh`
- `#SBATCH --cpus-per-task=256`
- `#SBATCH --mem=512gb`
- `#SBATCH --time=05:00:00`
- `#SBATCH --constraint=cal`
- `#SBATCH --error=pipe.%J.err`
- `#SBATCH --output=pipe.%J.out`
- module load squeezemeta/1.6.3
- `SqueezeMeta.pl -m coassembly \`
- `-s samples.txt -f data \`
- `-t 128 -p metagenome_assembly`

# The output of Sqezemeta

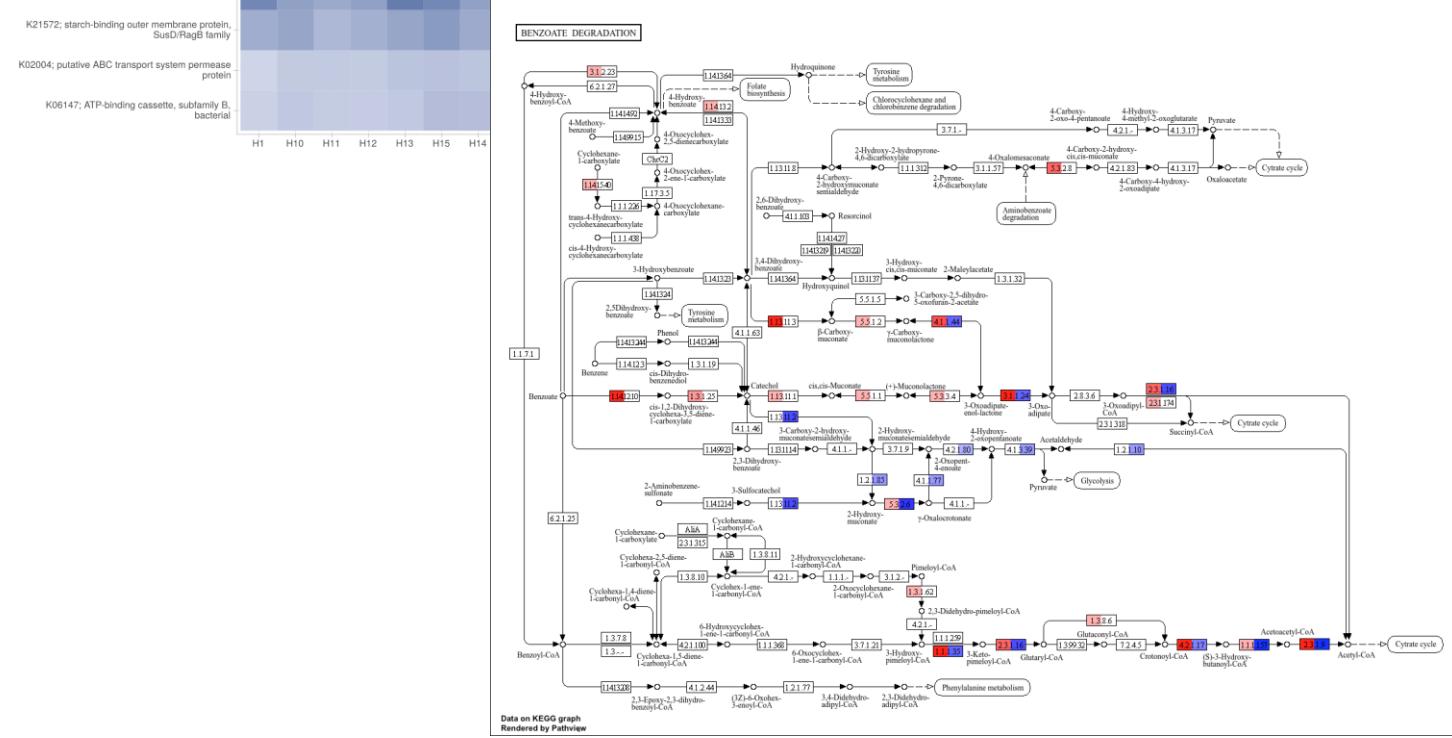
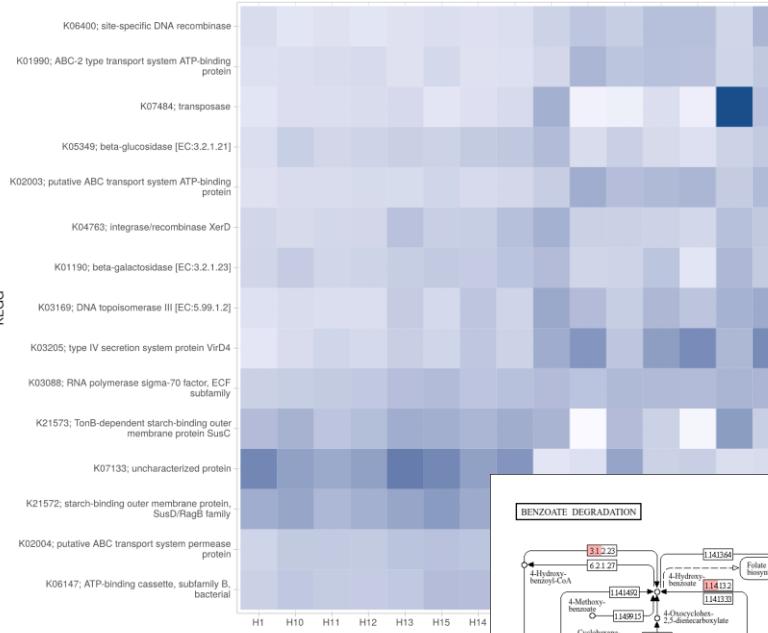
- Logs
- Assembly
- Bins
- ...Tables

```
drwxr-xr-x  9 adoss  staff  288 Mar  5 17:25 do not open
-rwxr-xr-x@ 1 adoss  staff  168 Mar  5 17:47 spades.sh
[adoss@Adams-Laptop meta_exercise % cd ..
[adoss@Adams-Laptop meta_exc % ll
total 2220432
drwxr-xr-x 14 adoss  staff          448 Mar  5 17:24 data not used
-rw-r--r--@ 1 adoss  staff          253 Mar  5 15:00 meta_exercise.Rproj
drwxr-xr-x  7 adoss  staff          224 Mar  5 17:47 meta_exercise
-rw-r--r--  1 adoss  staff  1136851089 Mar  5 17:29 meta_exercise.zip
drwxr-xr-x@ 14 adoss  staff          448 Mar  5 17:24 quast_megahit
drwxr-xr-x@ 14 adoss  staff          448 Mar  5 17:24 spades_megahit
drwxr-xr-x  8 adoss  staff          256 Mar  6 07:26 squeezemeta_results
-rw-r--r--  1 adoss  staff          195 Mar  5 17:37 thingstoinstall.txt
[adoss@Adams-Laptop meta_exc % cd squeezemeta_results
[adoss@Adams-Laptop squeezemeta_results % ll
total 64
drwxr-xr-x 15 adoss  staff      480 Mar  6 07:25 metagenome_assembly
-rw-r--r--@ 1 adoss  staff      458 Mar  5 12:47 pipe.10654274.err
-rw-r--r--@ 1 adoss  staff  17376 Mar  5 12:58 pipe.10654274.out
-rw-r--r--@ 1 adoss  staff      126 Mar  5 11:19 samples.txt
-rw-r--r--@ 1 adoss  staff      308 Mar  5 11:22 squeeze.sh
[adoss@Adams-Laptop squeezemeta_results % less samples.txt
```

# SQMtools



lv1	lv2	lv3	type	rows/names	columns	data
\$orfs	Stable		dataframe	orfs	misc. data	misc. data
	Sabund		numeric matrix	orfs	samples	abundances (reads)
	Sbases		numeric matrix	orfs	samples	abundances (bases)
	Scov		numeric matrix	orfs	samples	coverages
	Scpm		numeric matrix	orfs	samples	coverages per million of reads
	Stpm		numeric matrix	orfs	samples	tpm
	Sseqs		character vector	orfs	(n/a)	sequences
	Stax		character matrix	orfs	tax. ranks	taxonomy
\$contigs	Stable		dataframe	contigs	misc. data	misc. data
	Sabund		numeric matrix	contigs	samples	abundances (reads)
	Sbases		numeric matrix	contigs	samples	abundances (bases)
	Scov		numeric matrix	contigs	samples	coverages
	Scpm		numeric matrix	contigs	samples	coverages per million of reads
	Stpm		numeric matrix	contigs	samples	tpm
	Sseqs		character vector	contigs	(n/a)	sequences
	Stax		character matrix	contigs	tax. ranks	taxonomies
\$bins	Stable		dataframe	bins	misc. data	misc. data
	Slength		numeric vector	bins	(n/a)	length
	Sabund		numeric matrix	bins	samples	abundances (reads)
	Spercent		numeric matrix	bins	samples	percentages (reads)
	Sbases		numeric matrix	bins	samples	abundances (bases)
	Scov		numeric matrix	bins	samples	coverages
	Scpm		numeric matrix	bins	samples	coverages per million of reads
	Stax		character matrix	bins	tax. ranks	taxonomy
Staxa	Ssuperkingdom	Sabund	numeric matrix	superkingdoms	samples	abundances (reads)
	Spercent	Sabund	numeric matrix	superkingdoms	samples	percentages (reads)
	Sphylum	Sabund	numeric matrix	phyla	samples	abundances (reads)
	Spercent	Sabund	numeric matrix	phyla	samples	percentages (reads)
	Sclass	Sabund	numeric matrix	classes	samples	abundances (reads)
	Spercent	Sabund	numeric matrix	classes	samples	percentages (reads)
	Sorder	Sabund	numeric matrix	orders	samples	abundances (reads)
	Spercent	Sabund	numeric matrix	orders	samples	percentages (reads)
	Sfamily	Sabund	numeric matrix	families	samples	abundances (reads)
	Spercent	Sabund	numeric matrix	families	samples	percentages (reads)
	Sgenus	Sabund	numeric matrix	genera	samples	abundances (reads)
	Spercent	Sabund	numeric matrix	genera	samples	percentages (reads)
	SSpecies	Sabund	numeric matrix	species	samples	abundances (reads)
	Spercent	Sabund	numeric matrix	species	samples	percentages (reads)
Sfunctions	SKEGG	Sabund	numeric matrix	KEGG ids	samples	abundances (reads)
	Sbases	numeric matrix	KEGG ids		samples	abundances (bases)
	Scov	numeric matrix	KEGG ids		samples	coverages
	Scpm	numeric matrix	KEGG ids		samples	coverages per million of reads
	Stpm	numeric matrix	KEGG ids		samples	tpm
	Scopy_number	numeric matrix	KEGG ids		samples	avg. copies
	SCOG	Sabund	numeric matrix	COG ids	samples	abundances (reads)
	Sbases	numeric matrix	COG ids		samples	abundances (bases)
	Scov	numeric matrix	COG ids		samples	coverages
	Scpm	numeric matrix	COG ids		samples	coverages per million of reads
	Stpm	numeric matrix	COG ids		samples	tpm
	Scopy_number	numeric matrix	COG ids		samples	avg. copies
	SPFAM	Sabund	numeric matrix	PFAM ids	samples	abundances (reads)
	Sbases	numeric matrix	PFAM ids		samples	abundances (bases)
Stotal_reads	Scov	numeric matrix	PFAM ids		samples	coverages
	Scpm	numeric matrix	PFAM ids		samples	coverages per million of reads
	Stpm	numeric matrix	PFAM ids		samples	tpm
	Scopy_number	numeric matrix	PFAM ids		samples	avg. copies
	Sproject_name	character vector	(empty)	(n/a)		total reads
	Ssamples	character vector	(empty)	(n/a)		project name
	Stax_names_long	Ssuperkingdom	character vector	short taxa names	(n/a)	long taxa names
	Sphylum	character vector	short taxa names	(n/a)		long taxa names
Stax_names_short	Sclass	character vector	short taxa names	(n/a)		long taxa names
	Sorder	character vector	short taxa names	(n/a)		long taxa names
	Sfamily	character vector	short taxa names	(n/a)		long taxa names
	Sgenus	character vector	short taxa names	(n/a)		long taxa names
	SSpecies	character vector	short taxa names	(n/a)		long taxa names
	Stax_names_short	character vector	long taxa names	(n/a)		short taxa names
	SKEGG_names	character vector	KEGG ids	(n/a)		KEGG names
	SKEGG_paths	character vector	KEGG ids	(n/a)		KEGG hierarchy
SCOG_names	SCOG_paths	character vector	COG ids	(n/a)		COG names
	SCOG_paths	character vector	COG ids	(n/a)		COG hierarchy
Sext_annot_sources		character vector	(empty)	(n/a)		External databases



<https://github.com/jtamames/SqueezeMeta/wiki/Using-R-to-analyze-your-SQM-results>

## Preparation:

- Download the data
- Install SQMtools library in R
- Install “pathview” library in R (NOT “pathviewr”)
- BiocManager::install("pathview")
- Install assembler - spades and quast in one conda environment
- `conda create -n genome_assembly -c bioconda \ spades quast`
- `conda activate genome_assembly`

## Tasks:

- Assemble the data using spades
- Check the quality of the assembly using quast
- Backup: Check which assembly is better using quast
- Read the results of SqueezeMeta pipeline in R and :
  - Check the bins and decide which are good and which not
  - Choose your 3 favorite function, plot taxonomy related to them and pathways