



PacificBiosciences / pb-16S-nf

Stéphane Plaisance [VIB - Nucleomics Core, nucleomics@vib.be]

MonNov14, 2022 - version 1.0

Contents

Introduction	2
Theoretical composition of the Zymo mock community	3
Nextflow install and setup	5
Nextflow test	5
Nextflow Zymo run	6
Zymo run results	8
Examples of QIIME-View outputs	9

Introduction

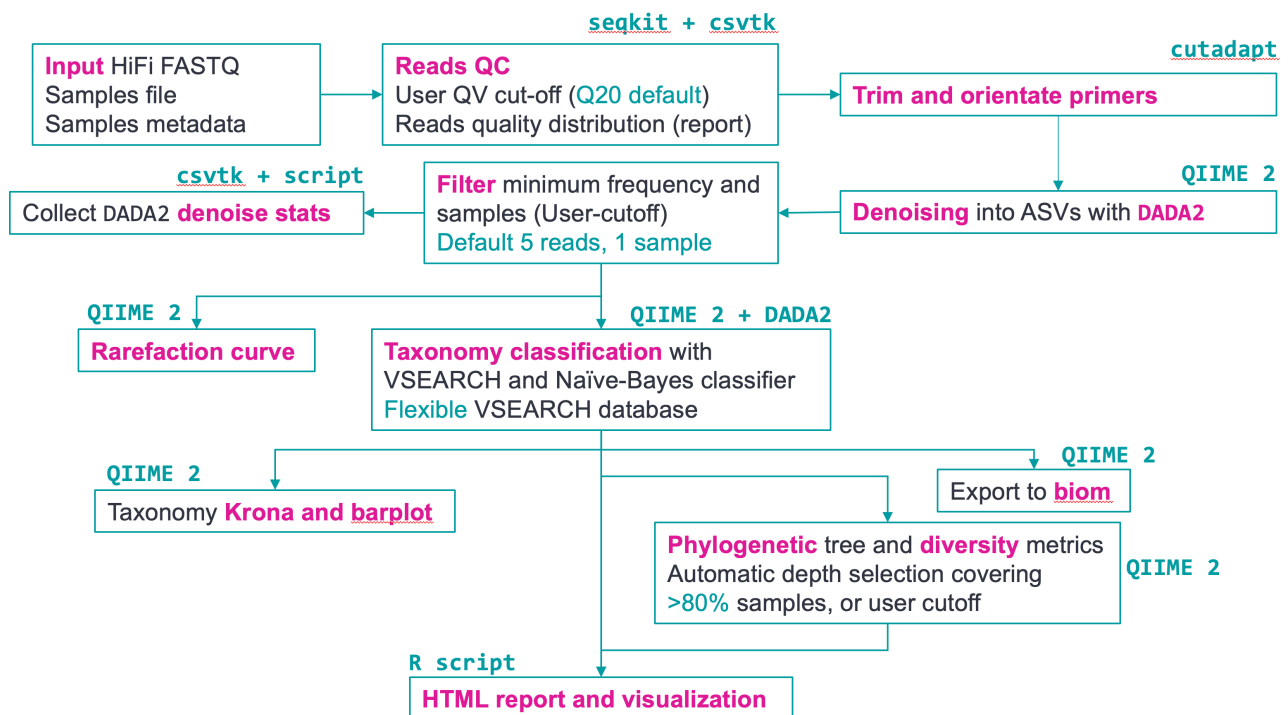
We describe here a new nextflow pipeline **pb-16S-nf** ¹ developed by **Khi Pin, Chua** (@proteinosome) as part of the Pacbio open code hosted on github and can be used to analyze data obtained with the Pacbio 16S method ². Khi Pin is actively developing this package further and was very helpful in deploying this code and correcting a few issues.

«This Nextflow pipeline is designed to process PacBio HiFi full-length 16S data into high quality amplicon sequence variants (ASVs) using **QIIME2** and **DADA2**. It provides a set of visualization through the QIIME 2 framework for interactive plotting. The pipeline generates a HTML report for the important statistics and top taxonomies» (taken from the github page).

The Nextflow pipeline depends on two text files and a matching folder of demultiplexed HiFi fastq files produced by the SMRTLink platform.

The pipeline performs a number pre-processing steps followed by DADA2 and Qiime2 commands. All of it integrated and standardized for ease of use.

The general workflow is shown in the next figure



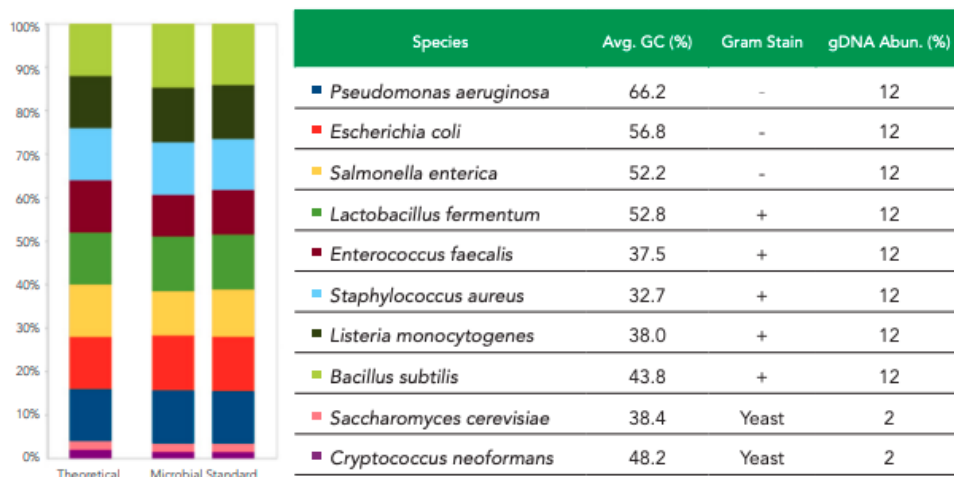
¹<https://github.com/PacificBiosciences/pb-16S-nf>

²<https://www.pacb.com/wp-content/uploads/Procedure-checklist-Amplification-of-bacterial-full-length-16S-rRNA-gene-with-barcoded-primers.pdf>

Theoretical composition of the Zymo mock community

In order to validate the workflow, we collected barcoded positive control samples from 5 Nucleomics Core 16S Sequel-IIe experiments and used them to compare pb-16S-nf results to the theoretical distribution present in the Zymo mock community ³

Defined Microbial Community



The ZymoBIOMICS® Microbial Community Standard contains three easy-to-lyse bacteria, five tough-to-lyse bacteria, and two tough-to-lyse yeasts.

Note: The copy number of rDNA genes in bacteria can vary a lot and contribute to 16S counts appear differentially affected depending on the bacterial host. This is well documented for the mock community used in this experiment in the Zymo protocol document ⁴.

Table 1: Microbial Composition

Species	Theoretical Composition (%)				
	Genomic DNA	16S Only ¹	16S & 18S ¹	Genome Copy ²	Cell Number ³
<i>Pseudomonas aeruginosa</i>	12	4.2	3.6	6.1	6.1
<i>Escherichia coli</i>	12	10.1	8.9	8.5	8.5
<i>Salmonella enterica</i>	12	10.4	9.1	8.7	8.8
<i>Lactobacillus fermentum</i>	12	18.4	16.1	21.6	21.9
<i>Enterococcus faecalis</i>	12	9.9	8.7	14.6	14.6
<i>Staphylococcus aureus</i>	12	15.5	13.6	15.2	15.3
<i>Listeria monocytogenes</i>	12	14.1	12.4	13.9	13.9
<i>Bacillus subtilis</i>	12	17.4	15.3	10.3	10.3
<i>Saccharomyces cerevisiae</i>	2	NA	9.3	0.57	0.29
<i>Cryptococcus neoformans</i>	2	NA	3.3	0.37	0.18

¹ The theoretical composition in terms of 16S (or 16S & 18S) rRNA gene abundance was calculated from theoretical genomic DNA composition with the following formula: 16S/18S copy number = total genomic DNA (g) × unit conversion constant (bp/g) / genome size (bp) × 16S/18S copy number per genome. Use this as reference when performing 16S targeted sequencing.

² The theoretical composition in terms of genome copy number was calculated from theoretical genomic DNA composition with the following formula: genome copy number = total genomic DNA (g) × unit conversion constant (bp/g) / genome size (bp). Use this as reference when inferring microbial abundance from shotgun sequencing data based on read depth.

³ The theoretical composition in terms of cell number was calculated from theoretical genomic DNA composition with the following formula: cell number = total genomic DNA (g) × unit conversion constant (bp/g) / genome size (bp)/ploidy.

³https://files.zymoresearch.com/datasheets/ds1706_zymbiomics_microbial_community_standards_data_sheet.pdf

⁴https://files.zymoresearch.com/protocols/_d6305_d6306_zymbiomics_microbial_community_dna_standard.pdf

Table 2: Strain Information

Species	NRRL Accession NO. ¹	Genome Size (Mb)	Ploidy	GC Content (%)	16/18S Copy Number	Gram Stain
<i>Pseudomonas aeruginosa</i>	B-3509	6.792	1	66.2	4	-
<i>Escherichia coli</i>	B-1109	4.875	1	46.7	7	-
<i>Salmonella enterica</i>	B-4212	4.760	1	52.2	7	-
<i>Lactobacillus fermentum</i>	B-1840	1.905	1	52.4	5	+
<i>Enterococcus faecalis</i>	B-537	2.845	1	37.5	4	+
<i>Staphylococcus aureus</i>	B-41012	2.730	1	32.9	6	+
<i>Listeria monocytogenes</i>	B-33116	2.992	1	38.0	6	+
<i>Bacillus subtilis</i>	B-354	4.045	1	43.9	10	+
<i>Saccharomyces cerevisiae</i>	Y-567	12.1	2	38.3	109 ²	Yeast
<i>Cryptococcus neoformans</i>	Y-2534	18.9	2	48.3	60 ²	Yeast

Notes:

¹ Several strains within the standard were replaced with similar strains beginning from Lot ZRC190633. This update will not affect the species composition of the standard. Refer to Appendix B to check if your product is from an older lot, and find the correct reference database to use accordingly.

² 18S rRNA gene copy numbers in a haploid genome of the two strains of *Saccharomyces cerevisiae* and *Cryptococcus neoformans* were estimated based on read depth information from mapping shotgun sequencing data.

Note that yeast have copy numbers one order of magnitude higher than bacteria, although not relevant here it will be in the case of mixed populations between yeast and bacteria and amplicons for both 16S and 18S (or ITS).

Nextflow install and setup

The nextflow pipeline is available from github and should be cloned locally on the analysis server. After download of the github repo, the first nextflow run gets 3 Docker images and downloads the classification databases.

During the first run, the docker components (N=3) will be downloaded and installed in the default Docker cache on the server

- kpinpb/pb-16s-nf-tools
- kpinpb/pb-16s-nf-qiime
- kpinpb/pb-16s-vis

```
git clone git@github.com:PacificBiosciences/pb-16S-nf.git
cd pb-16S-nf

nextflow run main.nf \
  --download-db \
  --profile docker

# a new 'databases' folder is added to the nextflow folder
# GTDB_bac120_arc53_ssu_r207_fullTaxo.fa.gz
# GTDB_ssu_all_r207.qza
# GTDB_ssu_all_r207_taxonomy.qza
# RefSeq_16S_6-11-20_RDPv16_fullTaxo.fa.gz
# silva-138-99-seqs.qza
# silva-138-99-tax.qza
# silva_nr99_v138.1_wSpecies_train_set.fa.gz
```

Nextflow test

A built-in test data can be used to validate the install as follows

```
# Create test_sample.tsv for testing
echo -e "sample-id\tabsolute-filepath\ntest_data\t$(readlink -f test_data/test_1000_reads.fastq.gz)" > test_data/test_sample.tsv

nextflow run main.nf \
  --input test_data/test_sample.tsv \
  --metadata test_data/test_metadata.tsv \
  --outdir test_results \
  --profile docker
```

The run should take only few minutes and produce a folder with intermediate data and results as discussed later in this report.

Nextflow Zymo run

Two text files need to be prepared based on the available read sets; a sample manifest and a metadata file (as standard in QIIME2):

- manifest: a sample.tsv file that relates the sample names and the full path to each fastq file
- metadata: a tsv file that relates the same sample names to sample groups or conditions used in the wet-lab experiment (info provided by the customer)

The two files used in this test run are reproduced below

- sample.tsv

sample-id	absolute-file-path
4170_bc1005-bc1096	/data/analyses/Zymo-SequelIIe-Hifi/reads/4170_bc1005-bc1096.fastq.gz
4285_bc1022-bc1107	/data/analyses/Zymo-SequelIIe-Hifi/reads/4285_bc1022-bc1107.fastq.gz
4296_bc1022-bc1060	/data/analyses/Zymo-SequelIIe-Hifi/reads/4296_bc1022-bc1060.fastq.gz
4112_bc1008-bc1075	/data/analyses/Zymo-SequelIIe-Hifi/reads/4112_bc1008-bc1075.fastq.gz
4128_bc1005-bc1107	/data/analyses/Zymo-SequelIIe-Hifi/reads/4128_bc1005-bc1107.fastq.gz

- metadata.tsv

sample_name	condition
4170_bc1005-bc1096	control
4285_bc1022-bc1107	control
4296_bc1022-bc1060	control
4112_bc1008-bc1075	control
4128_bc1005-bc1107	control

the condition column should be of type ‘categorical’ (not numeric!)

Note that in a real experiment, the conditions will describe more sample groups than just ‘control’

The following code was run to start the analysis:

```
# the full list of full path fastq can be obtained with:
# find fastq_folder -name "*.fastq.gz" -exec readlink -f {} \;

# use >= 32 cpu for good performance
cpu=32
infolder=<path-to-indata>
sample_file=${infolder}/sample.tsv
metadata_file=${infolder}/metadata.tsv
outfolder=<path-to-outdata>

nextflow run main.nf \
  --input ${sample_file} \
  --metadata ${metadata_file} \
  --outdir ${outfolder} \
  --dada2_cpu ${cpu} \
  --vsearch_cpu ${cpu} \
  --cutadapt_cpu ${cpu} \
  --profile docker
```

The nextflow pipeline produces live output and stores all log files for inspection as well as trouble-shooting

The pipeline executes the following tasks using here default parameters

Launching `main.nf` [big_torvalds] DSL2 - revision: 6990708c9f

Parameters set for pb-16S-nf pipeline for PacBio HiFi 16S

=====

Number of samples in samples TSV: 176

Filter input reads above Q: 20

Trim primers with cutadapt: Yes

Forward primer: AGRGTTYGATYMTGGCTCAG

Reverse primer: AAGTCGTAACAAGGTARCY

Minimum amplicon length filtered in DADA2: 1000

Maximum amplicon length filtered in DADA2: 1600

maxEE parameter for DADA2 filterAndTrim: 2

minQ parameter for DADA2 filterAndTrim: 0

Pooling method for DADA2 denoise process: pseudo

Minimum number of samples required to keep any ASV: 1

Minimum number of reads required to keep any ASV: 5

Taxonomy sequence database for VSEARCH: /opt/biotools/pb-16S-nf/databases/GTDB_ssu_all_r207.qza

Taxonomy annotation database for VSEARCH: /opt/biotools/pb-16S-nf/databases/GTDB_ssu_all_r207.taxonomy.qza

Skip Naive Bayes classification: false

SILVA database for Naive Bayes classifier: /opt/biotools/pb-16S-nf/databases/silva_nr99_v138.1_wSpecies_train_set.fa.gz

GTDB database for Naive Bayes classifier: /opt/biotools/pb-16S-nf/databases/GTDB_bac120_arc53_ssu_r207_fullTaxo.fa.gz

RefSeq + RDP database for Naive Bayes classifier: /opt/biotools/pb-16S-nf/databases/RefSeq_16S_6-11-20_RDPv16_fullTaxo.fa.gz

VSEARCH maxreject: 100

VSEARCH maxaccept: 100

VSEARCH perc-identity: 0.97

QIIME 2 rarefaction curve sampling depth: null

Number of threads specified for cutadapt: 80

Number of threads specified for DADA2: 80

Number of threads specified for VSEARCH: 80

Script location for HTML report generation: /opt/biotools/pb-16S-nf/scripts/visualize_biom.Rmd

Container enabled via docker/singularity: true

Version of Nextflow pipeline: 0.4

executor > Local (534)

```

[-      ] process > pb16S:write_log                -
[-      ] process > pb16S:QC_fastq (176)           -
[-      ] process > pb16S:cutadapt (176)           -
[-      ] process > pb16S:QC_fastq_post_trim (176) -
[-      ] process > pb16S:collect_QC               -
[-      ] process > pb16S:prepare_qiime2_manifest  -
[-      ] process > pb16S:import_qiime2            -
[-      ] process > pb16S:demux_summarize           -
[-      ] process > pb16S:dada2_denoise            -
[-      ] process > pb16S:filter_dada2             -
[-      ] process > pb16S:dada2_qc                 -
[-      ] process > pb16S:qiime2_phylogeny_diversity -
[-      ] process > pb16S:dada2_rarefaction         -
[-      ] process > pb16S:class_tax                 -
[-      ] process > pb16S:dada2_assignTax           -
[-      ] process > pb16S:export_biom              -
[-      ] process > pb16S:barplot_nb               -
[-      ] process > pb16S:barplot                 -
[-      ] process > pb16S:html_rep                 -
[-      ] process > pb16S:krona_plot               -

```

Zymo run results

After running the 5 Zymo samples, the standard output can be inspected and part of it shared with the customer

The main output folder has the following standard structure:

```
cutadapt_summary
dada2
filtered_input_FASTQ
import_qiime
nb_tax
parameters.txt
results
summary_demux
trimmed_primers_FASTQ
```

The **results** folder contains symbolic links to all final key files and can be forwarded to the customer as-is

```
alpha-rarefaction-curves.qzv
best_tax_merged_freq_tax.tsv
best_taxonomy.tsv
best_taxonomy_withDB.tsv
best_tax.qza
dada2_qc.tsv
dada2_stats.qzv
dada2_table.qzv
feature-table-tax.biom
feature-table-tax_vsearch.biom
krona.qzv
merged_freq_tax.qzv
phylogeny_diversity
rarefaction_depth_suggested.txt
reads_QC
samplefile.txt
stats.tsv
tax_export
taxonomy_barplot_nb.qzv
taxonomy_barplot_vsearch.qzv
taxonomy_vsearch.qza
visualize_biom.html
vsearch_merged_freq_tax.tsv
```

The results of this run are shared next to this report in the **Zymo-SequellIe-Hifi_results_local** folder to allow more exploration of this typical data.

All files ending with *.qzv* are QIIME2 visualization files that can be fed to the online QIIME2-Viewer (<https://view.qiime2.org/>) to create and customize plots or tables.

Files with extension *.qza* are QIIME2 objects that can be reloaded in QIIME2 to proceed in the analysis while files with extension *.tsv* are data files that can be used for further analysis (eg. in *R*).

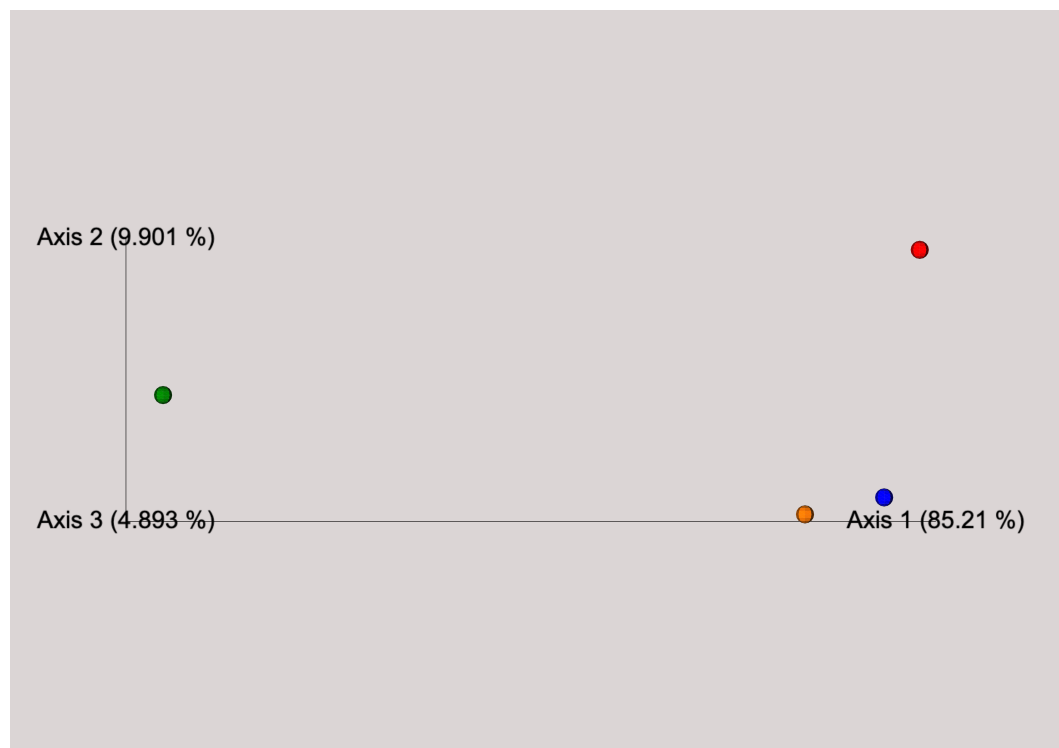
The main file present in the results folder is the RMarkdown converted document **visualize_biom.html** reporting all QC metrics and key findings through the user browser. Tables in that file are live and can be filtered.

Examples of QIIME-View outputs

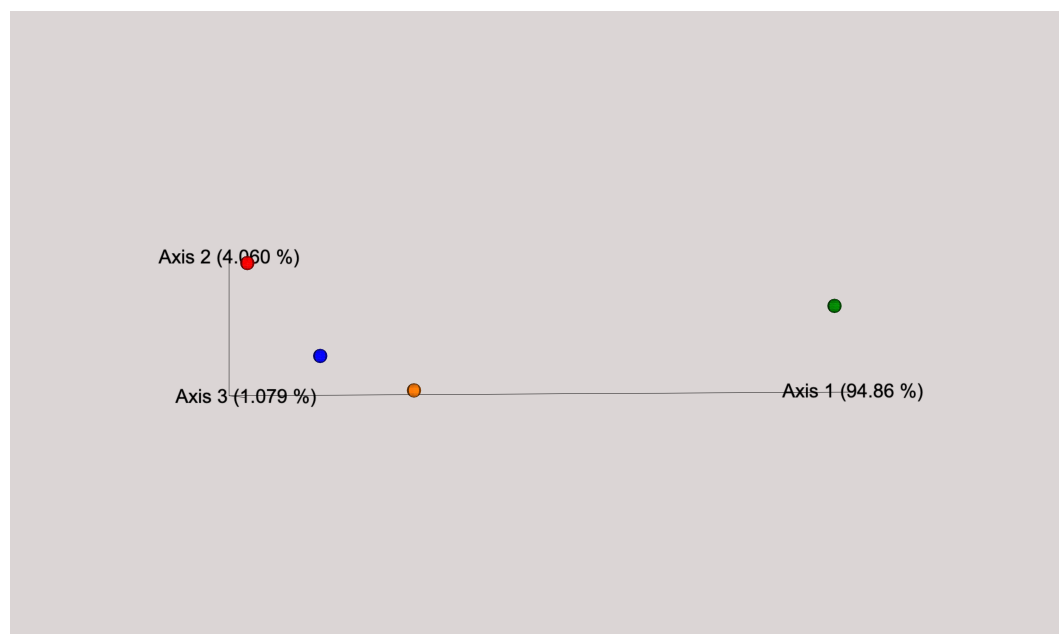
The *.qzv* files present in the results folder allow plotting using the **Qiime viewer site** ⁵

As illustration we show below two of the plots produced by the pipeline and showing multi dimensional principal component analyses results

- Bray Curtis



- Weighted unifrac

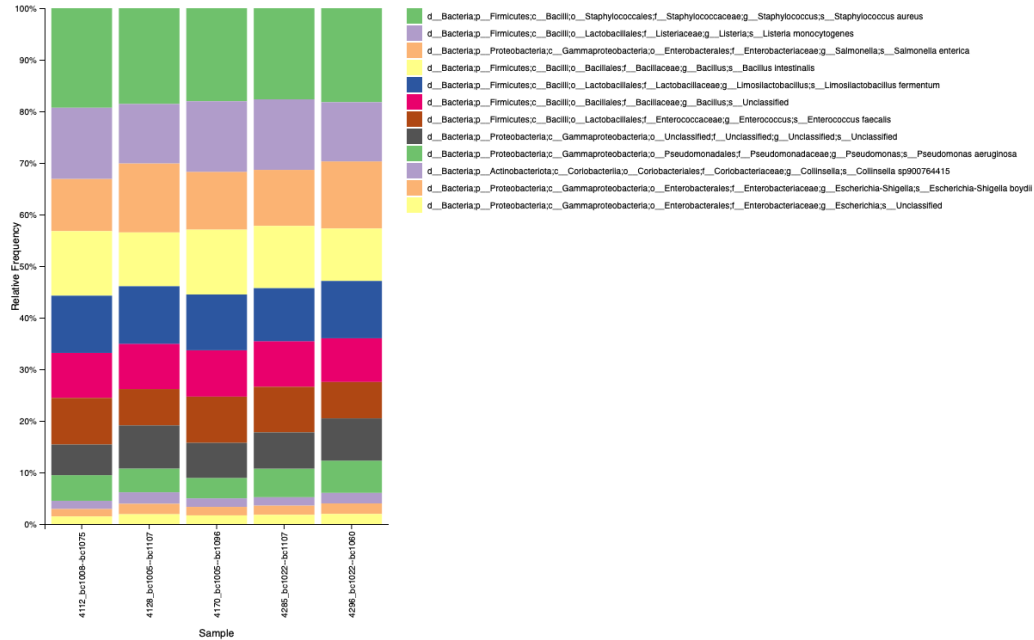


The viewer can also produce pictures for other results present in the results folder among which the two classification

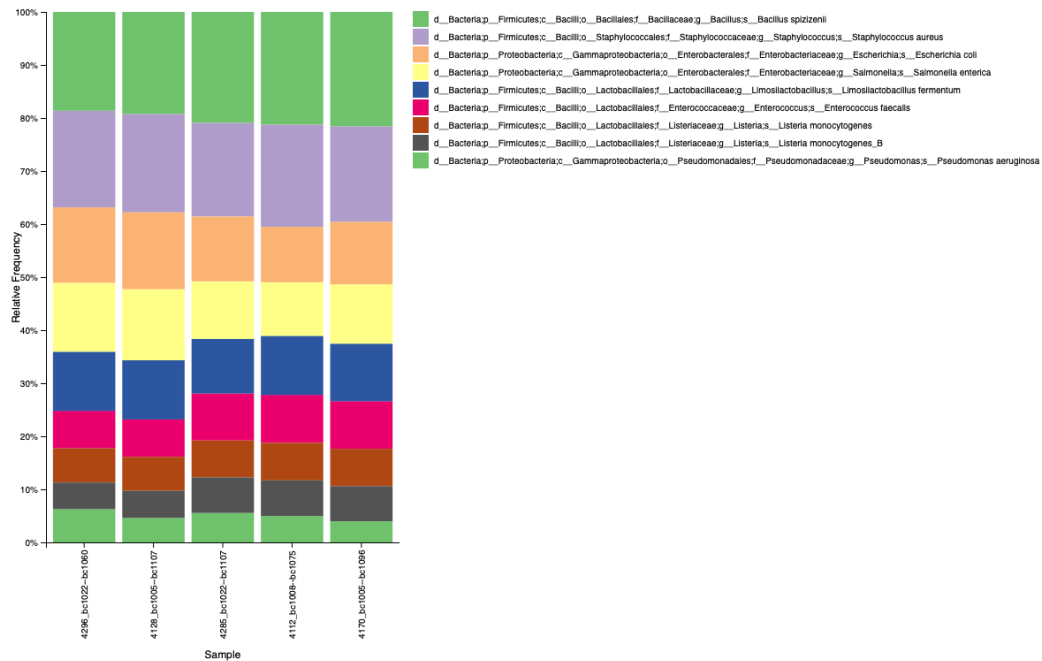
⁵<https://view.qiime2.org/>

outcomes produced by the workflow for our 5 Zymo samples. The *Vsearch* classification is sometimes more accurate and is based only on Vsearch best hits while the *Naive Bayes Classifier* is based on multiple search results and may be more complete but may sometimes include absent species (see doc on github)).

- taxonomy_barplot_nb.qzv

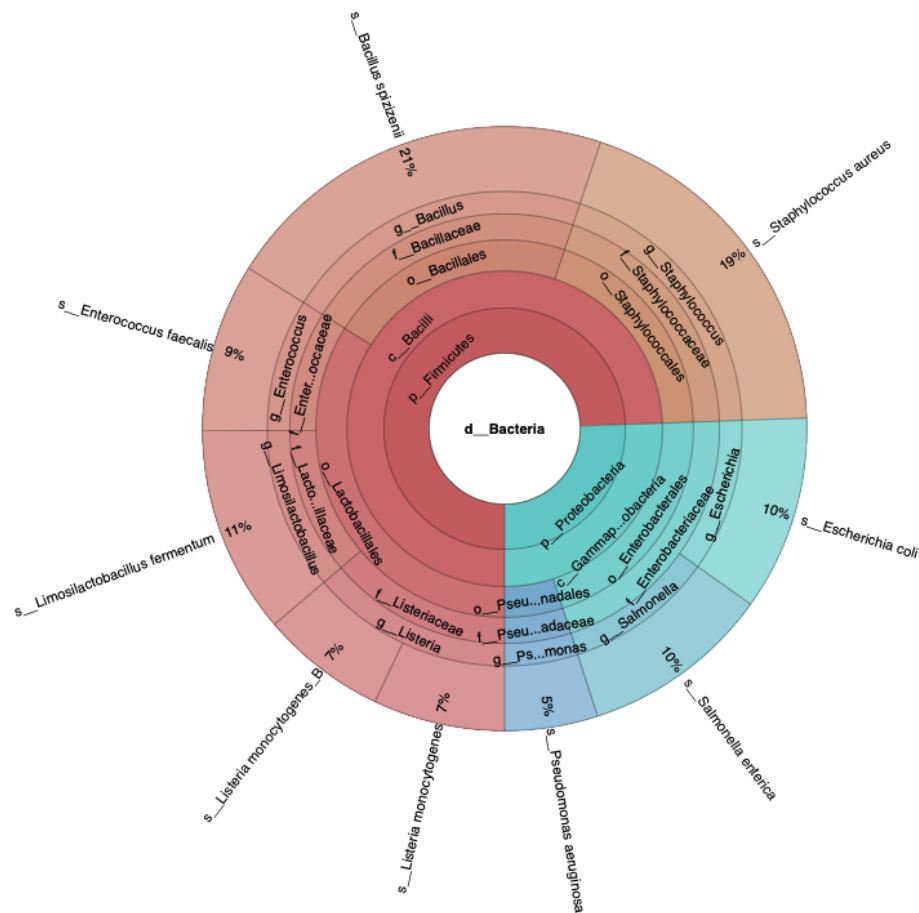


- taxonomy_barplot_vsearch.qzv



When compared to the theoretical distribution shown in the first part of this report, experimental results are nicely concordant.

Additional plots can be produced like the Krona classification shown next (sample bc1008-bc1075)



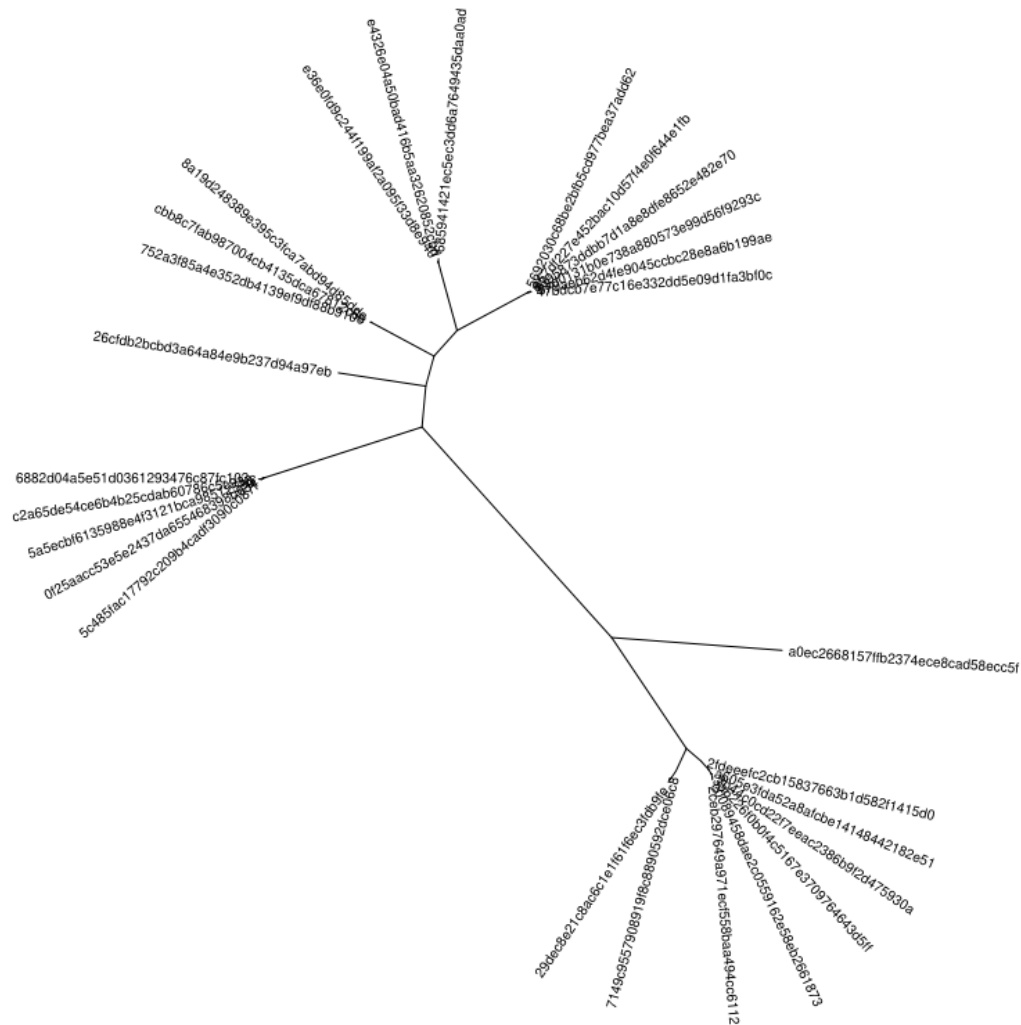
Various tables (.qzv) can be converted to pretty tables in the viewer as well.

As example, the content of **dada2_stats.qzv** (this table is also present in the html report)

sample-id #q2types	input numeric	filtered numeric	percentage of input passed filter numeric	denoised numeric	non-chimeric numeric	percentage of input non-chimeric numeric
4112_bc1008--bc1075	33316	27023	81.11	26898	26890	80.71
4128_bc1005--bc1107	21887	16156	73.82	15004	14163	64.71
4170_bc1005--bc1096	34400	29109	84.62	28728	28654	83.3
4285_bc1022--bc1107	21090	18261	86.59	17923	17891	84.83
4296_bc1022--bc1060	54492	42356	77.73	39248	36322	66.66

A phylogeny tree can be converted to a picture using a tree viewer ⁶ and the file *phylotree_mafft_rooted.nwk* (note that, in the current version of the pipeline, the tree shows the ASV labels rather than the stain names)

⁶<https://github.com/arklumpus/TreeViewer>



A little *R* magic can replace the Feature.ID used for tip labels with genus and species extracted from the taxonomy results.

