
teleost Documentation

Release latest

Mar 13, 2018

Contents:

1	Cookbook	1
2	Generalities	3
2.1	Defining variables	3
2.2	Logging functions	3
2.3	Determining sample names	3
3	Workflows	5
3.1	Core genome determination	5
3.2	Assembly and quality	6
3.3	Resistance	7
3.4	Virulence	7
3.5	Epidemiology	7
4	Indices and tables	9

CHAPTER 1

Cookbook

2.1 Defining variables

As a general rules, any `variable` referenced in this documentation must be either:

- Defined in the `yaml` config file that is passed to `snakemake` by `--configfile`
- Defined directly in the `snakemake` command by `--config variable=$value`

2.2 Logging functions

Archiving processes are defined in the file `workflows/logging.rules`. The variable `logging_folder` must be defined in the `config.yaml` or passed to `snakemake` with `--config`. Each time an effective `snakemake` run is started, a folder named with the current UTC datetime is created. A variable number of files will be copied there, so that replication of the run is possible:

- The snakefile passed to `snakemake`
- The config file
- The full command used, copied into the file `cmd.txt`
- The parameter files defining the SRA and the local samples, if they exist

The logs of every command run during the execution of the workflow will then be stored in this folder.

2.3 Determining sample names

Sample naming and matching to fastq files are handled in the file `workflows/making_sample_dataset.rules`.

2.3.1 Local samples

Local samples will be determined based on a tabulated file whose full path must be passed to the variable `local_samples` in the `config.yaml` or through `--config` on the `snakemake` command. It must contain at least two columns: *SampleName* and *ScientificName*.

Table 2.1: Local data example

SampleName	ScientificName
S10	Staphylococcus aureus
S1	Staphylococcus aureus

For each entry, there must be in the folder defined by the `link_directory` variable, two files (for paired reads) or only one (for single reads) whose filename starts by one and only one entry of the *SampleName* columns. For instance, the files `S10_001_R1_L001.fastq.gz` and `S10_001_R2_L001.fastq.gz` in the folder defined by the `link_directory` variable will be matched to the sample name `S10`. The matching is performed by using regular expressions to end the search at non alphanumeric characters or by the end of the word, thus the sample name `S1` will actually not match `S10_001_R1_L001.fastq.gz` nor `S10_001_R2_L001.fastq.gz`.

If needed, an *OldSampleName* column can be added to the file, when the read filenames and the desired new sample names can not be matched simply by testing the identity at the start of both names.

Table 2.2: Local data example with old sample names

SampleName	ScientificName	OldSampleName
S10	Staphylococcus aureus	Staur-10
S1	Staphylococcus aureus	Staur-1

In this case, the files `Staur-10_S10_L001_R1_001.fastq.gz` and `Staur-10_S10_L001_R2_001.fastq.gz` in the folder defined in `link_directory` will be matched to the sample name `S10`. Similarly, `Staur-1` will actually not match `Staur-10_S10_L001_R1_001.fastq.gz`.

2.3.2 SRA samples

SRA samples will be determined based on the tabulated file whose full path must be passed to the variable `sra_samples`. The `RunInfo` files that can be downloaded through the [SRA NCBI](#) database can be directly passed without any modification. Otherwise, four columns must be defined.

Table 2.3: SRA data example

Run	SampleName	LibraryLayout	ScientificName
ERR1140788	Mycobacterium_tuberculosis_N0145-Lineage_2	paired	Mycobacterium tuberculosis
SRR006916	Mycobacterium_tuberculosis_K21-Lineage_1	single	Mycobacterium tuberculosis

Current available workflows are implemented in the folder `workflows`. Each workflow will depend on `rules`, stored in the folder of the same name, and can also depend on other workflows. `rules` are sorted with respect to their general function in different folders.

3.1 Core genome determination

Core genomes can be calculated by three different means.

3.1.1 Ridom

cgMLST scheme from `ridom` can be extracted directly for theses species

Table 3.1: Available cgMLST schemes from ridom

Species	Taxonomy ID	Ridom ID	Reference genome assembly ID
<i>Staphylococcus aureus</i>	1280	141106	33148
<i>Mycobacterium tuberculosis</i>	1773	741110	538048
<i>Listeria monocytogenes</i>	1639	690488	264498
<i>Escherichia coli</i>	562	5064703	79781
<i>Klebsiella pneumoniae</i>	573	2187931	31388
<i>Enterococcus faecium</i>	1352	991893	526908
<i>Acinetobacter baumannii</i>	470	3956907	39528
<i>Legionella pneumophila</i>	446	1025099	30068

A bed file is constructed from the locus target file, constructing coordinates from the start and length columns of the csv file available on the [ridom website](#).

Example target file:

```
snakemake --snakefile $pipeline_folder/workflows/core_genome/make_ridom.rules
core_genomes/Staphylococcus_aureus/ridom/33148.bed
```

 will create the BED file defining the core genomic regions in the genome of the assembly ID 33148 (*Staphylococcus aureus* COL)

3.1.2 Enterobase

cgMLST scheme from [enterobase](#) is extracted for *Salmonella enterica*:

Table 3.2: Available cgMLST schemes from enterobase

Species	Taxonomy ID	Enterobase ID	Reference genome assembly ID	Scheme
<i>Salmonella enterica</i>	28901	SALwgMLST	359488	cgMLSTv1

A bed file for the reference genome [359488](#), based on the locus tag present in this genome is constructed. For instance, over the 3002 locus of the *Salmonella* cgMLSTv1, 69 come from a different genome than the reference 359488.

```
snakemake --snakefile $pipeline_folder/workflows/core_genome/make_enterobase.
rules core_genomes/Salmonella_enterica/ridom/359488.bed
```

 will create the BED file defining the core genomic regions in the genome of the assembly ID 359488 (*Salmonella enterica* subsp. *enterica* serovar Typhimurium str. D23580)

3.1.3 ParSNP

For species unavailable on either resource, core genome can be calculated using *parsnp* and the complete genomes of the species available on RefSeq. As ParSNP is not available on bioconda, the binary must be downloaded from the [ParSNP website](#) and placed in your \$PATH.

3.2 Assembly and quality

Aggregates rules for assembling genomes and performing various quality control checks. Required parameters:

- `cov_cutoff`: contigs whose coverage is below this cutoff will be excluded from the final assembly
- `adapter_file_name`: look for the adaptor for this library preparation kit (possible [values](#))
- `adapter_removal_param1`, `adapter_removal_param2`, `adapter_removal_param3`: parameters for adapter trimming ([reference](#))
- `minimum_quality_base`: leading and trailing bases below this quality will be removed
- `minimum_read_length`: reads shorter than this threshold after trimming will be discarded (be careful when using reads from SRA!)

Deliverables:

- `quality/multiqc/self_genome/multiqc_report.html`: quality control report based on the results of **fastqc**, **trimmomatic**, **qualimap**, **quast** and **prokka** for every sample
- `samples/{sample_name}/annotation/`: folder containing all annotation files from the *prokka* software

3.3 Resistance

Depends on the *Assembly and quality* workflow.

Required parameters:

- `resistance_prediction_softwares`: list of software for genetic resistance assessment. Possible values: `mykrobe` and `rgi`.

Deliverables:

- `samples/{sample_name}/annotation/resistance/rgi.tsv`: results files for RGI
- `samples/{sample_name}/annotation/resistance/mykrobe.tsv`: results file for mykrobe

3.4 Virulence

Depends on the *Assembly and quality* workflow.

Required parameters:

- `virulence_factors`: file with list of uniprot accession of virulence factors. An example is available in the folder `data/staph/db/`

Deliverables:

- `virulence_summary.xlsx`: summary of virulence proteins found in every samples.

3.5 Epidemiology

Depends on the *Assembly and quality* workflow (for determining the Sequence Types).

Required parameters:

- `minimum_coverage_for_calling`: minimum of coverage for considering a genomic position when counting differences between samples. Any position (SNP or non-SNP when compared to the reference) having a lower coverage will be masked
- `minimum_alternate_fraction_for_calling`: minimum ratio of observations favouring a SNP over observations not favouring a SNP. Any SNPs not meeting this criteria will also be masked

Deliverables:

- `typing/{snp_caller}/core_{ridom or enterobase}/{reference_genome}/bwa/distance_snp_mst_no_st.svg`: Minimum spanning tree of the distance in snps between every sample over the core genome as defined by `ridom` or `enterobase`. Available species and values for reference genomes are listed in the files in `data/core_genome_dbs/`. If the species under consideration has a multiple locus sequence type available, `typing/{snp_caller}/core_{ridom or enterobase}/{reference_genome}/bwa/distance_snp_mst_with_st.svg` can be generated with the ST of each sample.
- `phylogeny/{snp_caller}/core_{ridom or enterobase}/{reference_genome}/bwa/phylogeny_no_st.svg`: A phylogeny based on the alignments of the core SNPs, using RAxML. Available species and values for reference genomes are listed in the files in `data/core_genome_dbs/`. If the species under consideration has a multiple locus sequence type available, `phylogeny/{snp_caller}/core_{ridom or enterobase}/{reference_genome}/bwa/phylogeny_with_st.svg` can be generated with the ST of each sample.

- `quality/multiqc/mapping_to_{reference_genome}/multiqc_report.html`: **multiqc** report of **qualimap**, **fastqc** and **trimmomatic** of every samples when mapping against the reference. Check for quality control.

CHAPTER 4

Indices and tables

- `genindex`
- `modindex`
- `search`