

# Snakemake

[https://github.com/nichollskc/IMPC\\_analysis](https://github.com/nichollskc/IMPC_analysis)

<https://snakemake.readthedocs.io/en/stable/index.html>

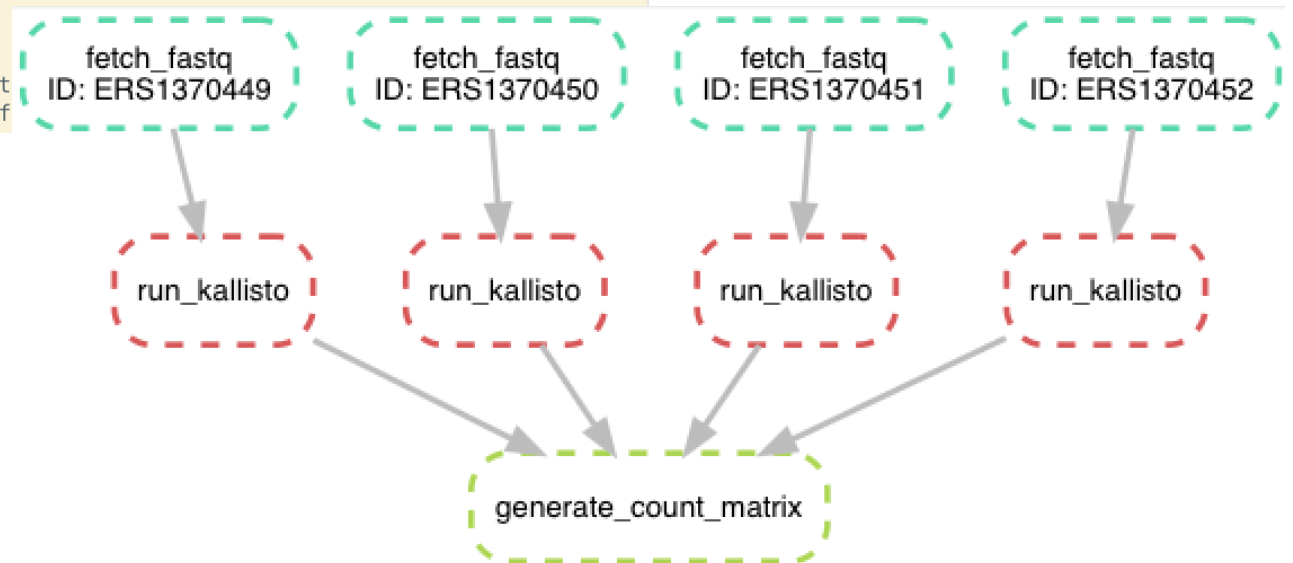
# Idea of snakemake – defines dependencies

- Python version of Make
- Rules define inputs, outputs and commands to go from inputs to outputs
- Re-runs a rule if the inputs are older than the outputs

```

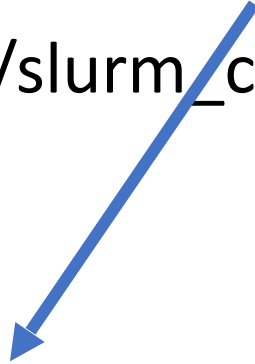
1 configfile: "config.yml"
2
3 rule fetch_fastq:
4     params:
5         url=lambda wildcards: config['FASTQ_URLS'][wildcards.ID]
6     output:
7         temp("data/fastq/{ID}.fastq.gz")
8     shell:
9         "wget -O {output} --no-verbose {params.url} 2>&1"
10
11 rule run_kallisto:
12     input:
13         "data/fastq/{ID}.fastq.gz"
14     output:
15         "data/kallisto/{ID}/abundance.tsv"
16     shell:
17         "kallisto quant --index={config[INDEX_FILE]} --output-dir=$(dir
18         name {output}) --single -l 50 -s 2 {input} 2>&1"
19
20 rule generate_count_matrix:
21     input:
22         expand("data/kallisto/{ID}/abundance.tsv",
23             ID=list(config['FASTQ_URLS'].keys()))
24     output:
25         df="data/tpm.tsv"
26
27 run:
28     import worker
29     df = worker.generate_count_matrix()
30     df.to_csv(output.df)

```



snakemake data/tpm.tsv

```
snakemake --cluster "<COMMAND>" --jobs 100  
--cluster-config cluster/slurm_config.json
```



# Nice features

- Temporary/protected outputs
- Special status for log files
- Job groups – run jobs in one SLURM job
- R/python scripts can access snakemake variables e.g. list of inputs
- Different cluster settings for different jobs
  
- Excellent documentation

<https://snakemake.readthedocs.io/en/stable/index.html>