# Introduction to Snakemake

Einar Birkeland, Researcher, Department of Pharmacy, UiO

30/10-24

# Overview

- Snakemake – what is it, and why use it?

- How it works

- Demonstration

- A couple of use cases
  - Pipeline for sequence data analysis
  - Handling of data analysis

- Additional features

- Nextflow and Snakemake

# Intro – What is Snakemake

- Workflow management system designed for reproducibility
- Simplifies defining complex pipelines in bioinformatics

- Modular and scalable approach
- Integration of existing scripts and tools

- DSL is based on Python
  - Python code seamlessly integrated

# Core concepts of Snakemake

- Snakemake workflows are <u>file-based</u> and defined by rules in a "Snakefile"
  - *Input*
  - *Output*
  - *Shell/script/run*

- Definition of software environments
  - Conda
  - Containers
  - Env-modules

- Parsing of dependencies to determine which jobs to do
  - Parse
  - Prioritize
  - Execute

# Why use Snakemake

Reproducible analysis

Why not use Snakemake (or equivalent tools)?

It takes time and effort to set up = **overhead**

However, reproducible code is
- Starting to be required by journals
- Good scientific practice

Raw data

Processing steps
- Tools (versions)
- Environments
- Settings

Data analysis
- Tools (versions)
- Environments
- Settings

Results
- Figures
- Tables
- Models, etc..

# Demo: building a simple workflow – 1 step

1. Define the file you want to create: An output file – "some_file.txt"

2. Define which file(s) is necessary to create this file (optional)

3. Define what needs to be done to create the output file – in the "shell" directive

```
workflow > ≡ Snakefile_1step.smk
1
2
3    rule create_some_file:
4        output:
5            "results/some_file.txt"
6        shell:
7            "echo 'hello' > data/some_file.txt"
```

# Demo: building a simple workflow – 2 steps

Case: We want to check out the quality of some publicly available sequencing data

1. Define the files we want to create
   1. Fastqc html files
2. Define what needs to be done to create these files
   1. Download the data
   2. Run fastqc

```
 9  rule downl
10      output
11          re
```

```
17      shell:
18          "w
19
20  rule fastq
21      output
22          ht
23          zi
24      input:
25          re
26      params:
27          outdir="data/fastqc/raw/"
28      shell:
29          "fastqc --quiet --outdir {params.outdir} {input.read}"
30
31
```

| | sample_id | url |
|---|---|---|
| 1 | sample_id→ | url |
| 2 | SRR25246600→ | ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR252/000/SRR25246600/SRR25246600 |
| 3 | SRR25246700→ | ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR252/000/SRR25246700/SRR25246700 |
| 4 | SRR25246800→ | ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR252/000/SRR25246800/SRR25246800 |
| 5 | SRR25246900→ | ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR252/000/SRR25246900/SRR25246900 |
| 6 | SRR25247000→ | ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR252/000/SRR25247000/SRR25247000 |
| 7 | SRR25247100→ | ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR252/000/SRR25247100/SRR25247100 |
| 8 | SRR25247200→ | ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR252/000/SRR25247200/SRR25247200 |
| 9 | SRR25247300→ | ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR252/000/SRR25247300/SRR25247300 |
| 10 | SRR25247400→ | ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR252/000/SRR25247400/SRR25247400 |
| 11 | SRR25247500→ | ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR252/000/SRR25247500/SRR25247500 |
| 12 | SRR25246601→ | ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR252/001/SRR25246601/SRR25246601 |
| 13 | SRR25246701→ | ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR252/001/SRR25246701/SRR25246701 |
| 14 | SRR25246801→ | ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR252/001/SRR25246801/SRR25246801 |
| 15 | SRR25246901→ | ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR252/001/SRR25246901/SRR25246901 |
| 16 | SRR25247001→ | ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR252/001/SRR25247001/SRR25247001 |
| 17 | SRR25247101→ | ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR252/001/SRR25247101/SRR25247101 |
| 18 | SRR25247201→ | ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR252/001/SRR25247201/SRR25247201 |
| 19 | SRR25247301→ | ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR252/001/SRR25247301/SRR25247301 |
| 20 | SRR25247401→ | ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR252/001/SRR25247401/SRR25247401 |
| 21 | SRR25247501→ | ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR252/001/SRR25247501/SRR25247501 |
| 22 | SRR25246602→ | ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR252/002/SRR25246602/SRR25246602 |
| 23 | SRR25246702→ | ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR252/002/SRR25246702/SRR25246702 |
| 24 | SRR25246802→ | ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR252/002/SRR25246802/SRR25246802 |
| 25 | SRR25246902→ | ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR252/002/SRR25246902/SRR25246902 |
| 26 | SRR25247002→ | ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR252/002/SRR25247002/SRR25247002 |
| 27 | SRR25247102→ | ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR252/002/SRR25247102/SRR25247102 |
| 28 | SRR25247202→ | ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR252/002/SRR25247202/SRR25247202 |
| 29 | SRR25247302→ | ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR252/002/SRR25247302/SRR25247302 |
| 30 | SRR25247402→ | ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR252/002/SRR25247402/SRR25247402 |
| 31 | SRR25247502→ | ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR252/002/SRR25247502/SRR25247502 |

```
 9   rule download_reads:
10       output:
11           reads=temp("data/raw_reads/{sample}_{R}.fastq.gz")




17       shell:
18           "wget -nc {params.url} -P {params.out_dir}"
19
20   rule fastqc:
21       output:
22           html="data/fastqc/raw/{sample}_{R}_fastqc.html",
23           zip="data/fastqc/raw/{sample}_{R}_fastqc.zip",
24       input:
25           read="data/raw_reads/{sample}_{R}.fastq.gz"
26       params:
27           outdir="data/fastqc/raw/"
28       shell:
29           "fastqc --quiet --outdir {params.outdir} {input.read}"
30
31
```

# Demo: building a simple workflow – 2 steps

Case: We want to check out the quality of some publicly available sequencing data

1. Define the files we want to create
    1. Fastqc html files
2. Define what needs to be done to create these files
    1. Download the data
    2. Run fastqc

Additional considerations

Python integration

1. Use pandas to read a file with information about samples
2. Used a lambda function to retrieve information from the pandas dataframe
3. Imported fastqc as a module

# Demo: building a simple workflow – 4 steps

Case: We want to check out the quality of some publicly available sequencing data, and run an R script on the multiqc output

1.  Define the files we want to create
    1.  **Multiqc** html file
    2.  A plot based on multiQC output
2.  Define what needs to be done to create these files
    1.  Download the data
    2.  Run fastqc
    3.  Run multiqc on all fastqc output
    4.  Call an R script that produces the output we want

**Define the samples we want to use**

**Define a function to get url**

**"rule all" specifies the target file**

```python
1   import pandas as pd
2
3   samples_file = "resources/samples.tsv"
4
5   sampleTable = pd.read_csv(samples_file, sep="\t")
6
7   samples = sampleTable["sample_id"].tolist()
8   samples = samples[:2]
9
10  def get_url(wildcards):
11      tmp_url = sampleTable[ sampleTable["sample_id"] == wildcards.sample]["url"].iloc[0]
12      return  tmp_url + "_" + wildcards.R + ".fastq.gz"
13
14  rule all:
15      input:
16          "results/downstream_analysis/plot_for_manuscript_1.png"
17
18  rule download_reads:
19      output:
20          reads=temp("data/raw_reads/{sample}_{R}.fastq.gz")
21      params:
22          url=lambda w: get_url(w),
23          out_dir="data/raw_reads"
24      threads:
25          1
26      shell:
27          "wget -nc {params.url} -P {params.out_dir}"
28
29  rule fastqc:
30      output:
31          html="results/fastqc/raw/{sample}_{R}_fastqc.html",
```

```
32                 zip="results/fastqc/raw/{sample}_{R}_fastqc.zip",
33             input:
34                 read="data/raw_reads/{sample}_{R}.fastq.gz"
35             params:
36                 outdir="results/fastqc/raw/"
37             threads:
38                 1
39             envmodules:
40                 "FastQC/0.11.9-Java-11"
41             shell:
42                 "fastqc --quiet --outdir {params.outdir}  {input.read}"
43
44     rule multiqc:
45         output:
46                 html="results/multiqc/raw/multiqc.html",
47                 general_stats="results/multiqc/raw/multiqc_data/multiqc_general_stats.txt"
48         input:
49                 fastqc=expand("results/fastqc/raw/{sample}_{R}_fastqc.zip", sample = samples, R = [1,2])
50         params:
51                 indir="results/fastqc/raw/",
52                 outdir="results/multiqc/raw/",
53                 outname="multiqc"
54         envmodules:
55                 "MultiQC/1.12-foss-2021b"
56         shell:
57                 "multiqc --force -o {params.outdir} -n {params.outname} {params.indir}"
58
59     rule downstream_analysis:
60         output:
61                 plot="results/downstream_analysis/plot_for_manuscript_1.png"
62         input:
63                 general_stats="results/multiqc/raw/multiqc_data/multiqc_general_stats.txt"
64         envmodules:
65                 "R-bundle-Bioconductor/3.15-foss-2022a-R-4.2.1"
66         script:
67                 "scripts/analyze_data.R"
68
```

Same envmodule as before for FastQC

Envmodule specific for MultiQC

Envmodule specific for R + bioconductor

```
 32                       zip="results/fastqc/raw/{sample}_{R}_fastqc.zip",
 33              input:
 34                  read="data/raw_reads/{sample}_{R}.fastq.gz"
 35              params:
 36                  outdir="results/fastqc/raw/"
 37              threads:
 38                  1
 39              envmodules:
 40                  "FastQC/0.11.9-Java-11"
 41              shell:
 42                  "fastqc --quiet --outdir {params.outdir}  - {input.read}"
 43
 44      rule multiqc:
 45          output:
 46              html="results/multiqc/raw/multiqc.html",
 47              general_stats="results/multiqc/raw/multiqc_data/multiqc_general_stats.txt"
 48          input:
 49              fastqc=expand("results/fastqc/raw/{sample}_{R}_fastqc.zip", sample = samples, R = [1,2])
 50          params:
 51              indir="results/fastqc/raw/",
 52              outdir="results/multiqc/raw/",
 53              outname="multiqc"
 54          envmodules:
 55              "MultiQC/1.12-foss-2021b"
 56          shell:
 57              "multiqc --force -o {params.outdir} -n {params.outname} {params.indir}"
 58
 59      rule downstream_analysis:
 60          output:
 61              plot="results/downstream_analysis/plot_for_manuscript_1.png"
 62          input:
 63              general_stats="results/multiqc/raw/multiqc_data/multiqc_general_stats.txt"
 64          envmodules:
 65              "R-bundle-Bioconductor/3.15-foss-2022a-R-4.2.1"
 66          script:
 67              "scripts/analyze_data.R"
 68
```

"expand" combines variables into a list

"script" directive allows integration btw snakemake and R/python/Julia... etc

# Additional features of Snakemake

Environments:

- Envmodules
- Conda
  - YAML specification
- Containers

Scalability:

- Parallelization
- Cluster compatability
  - Slurm, others
- Cloud computing
  - AWS, google cloud, etc

- Modularity
- Logging
- Benchmarking
- DAG generation: Visualization of dependencies

```
workflow > envs >  ! bowtie2.yaml
 1    name: bowtie2
 2    channels:
 3      - bioconda
 4      - conda-forge
 5      - defaults
 6    dependencies:
 7      - bbmap=39.01
 8      - bedtools=2.31.0
 9      - bowtie2=2.5.1
10      - samtools=1.17
11
```

Collaboration/community

- Snakemake wrappers
  - Best practices for common tasks

- Snakemake workflows
  - Collection of all snakemake workflows on github

# Find out more

[Snakemake homepage](#)

- [Paper](#)
- Tutorials
  - [Slides](#)
  - Online [tutorial](#)
- [Wrappers](#)
- [Workflows](#)

# Nextflow or Snakemake?

| Pipeline complexity | Snakemake Pros/Cons | Nextflow Pros/Cons |
|---|---|---|
| Simple | Pros:<br>Easy syntax, quick setup,<br>File-based management | Pros:<br>Containerization,<br>Parallelization |
| | Pros:<br>Limited flexibility,<br>Container handling | Cons:<br>Groovy syntax complexity, setup overhead |
| Intermediate | Pros:<br>Modularity, DAG visualization, cluster compatability | Pros:<br>Dynamic workflows, resource management,<br>Modularity |
| | Cons:<br>Limited branching,<br>Manual resource setting | Cons:<br>Steeper learning curve,<br>Debugging complexity |
| Complex | Pros:<br>Python integration<br>DAG visualization, cloud support | Pros:<br>Adaptive workflows, strong container/cloud support |
| | Cons:<br>Rigid execution, limited dynamic handling | Cons:<br>Resource-intensive, Groovy complexity |

# Comparison to Nextflow

- **Syntax/Language**
  - Snakemake: Python-like syntax, more accessible for beginners
  - Nextflow: Groovy-based DSL may have a steeper learning curve
- **Best Use Cases**
  - Snakemake: Local/small-scale workflows, quick customizations
  - Nextflow: Large-scale, distributed workflows
- **Reproducibility with Containers**
  - Both tools support Docker and Singularity for reproducibility
- **Scalability**
  - Top-down rule dependency resolving in Snakemake vs. bottom-up in Nextflow
    - Snakemake may be slower for large and complex workflows
  - Snakemake less easily portable to large-scale computational resources