

Workflows for reproducible, replicable, scalable, and portable science

Mark Adams, The University of Edinburgh
Tech Talk, 13 Aug 2024



mark.adams@ed.ac.uk



[@markjamesadams@genomic.social](#)



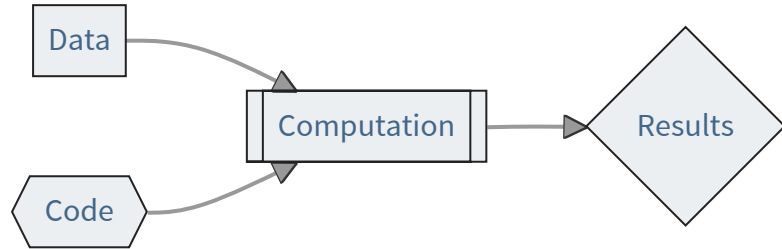
[@markjamesadams.bsky.social](#)

[X @mja](#)

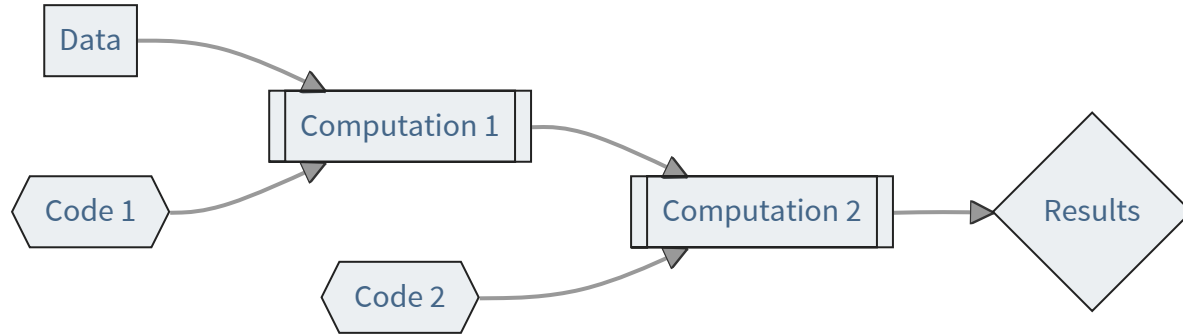
The “-ables” of workflows

- **reproducible:** same data and same code produce the same results
- **replicable:** same code runs with different data
- **scalable:** some code runs with more data and more resources
- **portable:** same code runs in different compute environments

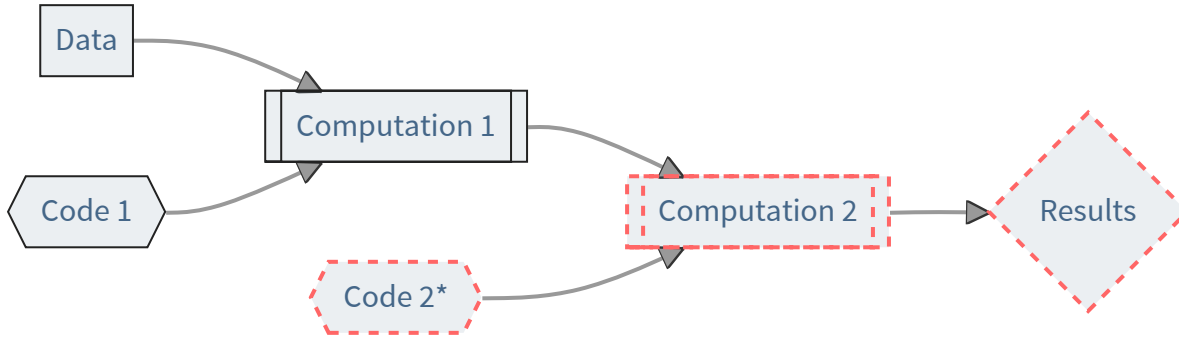
A workflow is a graph



A workflow is a graph



A workflow is a graph



When inputs change, only re-compute descendent outputs.

A workflow is a build system

Collection of notebooks.

```
thesis.qmd    symposium.qmd    workflows.qmd
```

Render them in a loop.

```
$ for QMD in *.qmd; do  
  quarto render $QMD  
done
```

Keep it DRY

Only render the notebook if the HTML doesn't exist or the notebook is newer.

```
$ for QMD in *.qmd; do
  PREFIX=$(basename $QMD .qmd)
  if [ ! -e ${PREFIX}.html ] || [ ${PREFIX}.qmd -nt ${PREFIX}.html ]; then
    quarto render $QMD
  fi
done
```

Makefiles¹

GNU Make: a tool for building programs from source code.

```
1 output : inputs
2     command
```

Rules encode relationship between inputs (“dependencies”) and outputs (“targets”)

Makefile with one rule:

Makefile

```
1 workflows.html : workflows.qmd
2     quarto render workflows.qmd
```

Run make

```
$ make workflows.html
```

```
quarto render workflows.qmd
```

```
processing file: workflows.qmd
```

```
1/3
```

```
2/3 [unnamed-chunk-1]
```

```
3/3
```

```
output file: workflows.knit.md
```

```
...
```

```
Output created: workflows.html
```

```
$ make workflows.html
```

```
make: `workflows.html' is up to date.
```

Makefile pattern rules

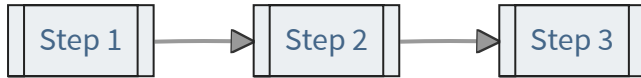
Pattern rule to render an HTML file from any Quarto notebook.

Makefile

```
1 %.html : %.qmd
2     quarto render $<
3
4 all: workflows.html symposium.html thesis.html
```

`all` rule specifies the outputs to render.

Workflows are pipelines



Unix pipes

```
bcftools view --targets-file targets.tsv dbsnp.v153.b37.vcf.gz | \
bcftools query --print-header --format '%CHROM\t%POS\t%ID\t%REF\t%ALT{0}\n' | \
gzip -c > chr_pos_rsid.tsv.gz
```

Scheduler (Sun Grid Engine) dependencies

step1.sh

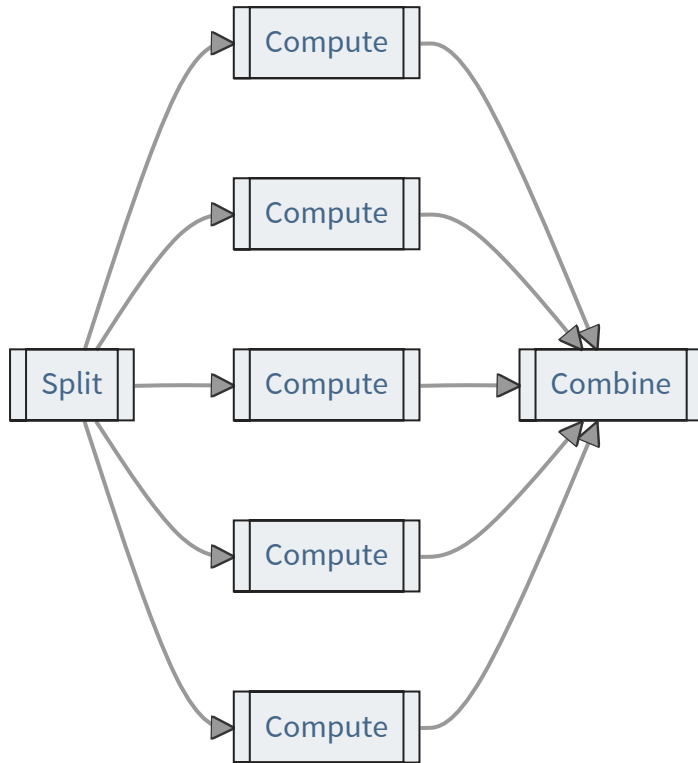
```
1  #$ -cwd
2  command1 --in $1 --out $2
```

step2.sh

```
1  #$ -cwd
2  command2 --in $1 --out $2
```

```
$ qsub -N step1 step1.sh input1 output1
$ qsub -N step2 -hold_jid step1 step2.sh output1 output2
```

Workflows are parallelisable



```
$ make all -j 3
```

```
quarto render workflows.qmd  
quarto render symposium.qmd  
quarto render thesis.qmd
```

Limitations of Makefiles

- No understanding of problem domain (need to run scripts, not just shell commands)
- Not scalable (only executes on local machine, not via HPC or cloud schedulers)
- Limited parallelisation

Limitations of schedulers

- Manually manage connecting outputs to next inputs
- Not portable between HPC environments
- Same workflow can't be developed and tested on your laptop

Workflow systems

- Provide a language to describe the pipeline
 - Declarative (CWL, WDL, Popper, Remake)
 - General purpose imperative (SciPipe, Dask, targets)
 - Domain specific imperative (bpipe, Snakemake, Nextflow)
- Manage connections and caching of inputs to outputs
- Hook into reproducible software environments (conda, Docker, Singularity)
- Scale to available resources (CPU cores, cluster nodes)

Snakemake and Nextflow

- Domain specific language (DSL)
 - Snakemake: Python
 - Nextflow: Groovy
- Workflow paradigm
 - Snakemake: **rules** determine required inputs from requested outputs
 - Nextflow: **processes** produce outputs from given inputs
- Each workflow step can be a shell command, a script (R, Python, Perl, Julia, etc), or native code [Python/Groovy]

Snakemake

Snakefile for rendering a notebook.

Snakefile

```
1 rule render:
2     input: "{notebook}.qmd"
3     output: "{notebook}.html"
4     shell: "quarto render {input}"
```

- `{notebook}` is an input/output wildcard.
- in the shell command, `{input}` is replaced with the value of `{notebook}.qmd`

Snakemake

Ask Snakemake to produce the file `workflows.html`.

Filename is matched to an output pattern, then built from the complementary input file.

```
$ snakemake -j1 workflows.html
```

```
Assuming unrestricted shared filesystem usage.
```

```
Building DAG of jobs...
```

```
Using shell: /bin/bash
```

```
Provided cores: 1 (use --cores to define parallelism)
```

```
Rules claiming more threads will be scaled down.
```

```
Job stats:
```

job	count
render	1
total	1

```
Select jobs to execute...
```

```
Execute 1 jobs...
```

```
[Wed Aug 7 07:10:02 2024]
```

localrule render:

Snakemake

Rerun the workflow

```
snakemake -j1 workflows.html
```

```
Assuming unrestricted shared filesystem usage.
```

```
Building DAG of jobs...
```

```
Nothing to be done (all requested files are present and up to date).
```

Snakemake

Ask the workflow for multiple files.

```
$ snakemake -j1 workflows.html symposium.html thesis.html
```

Assuming unrestricted shared filesystem usage.

Building DAG of jobs...

Using shell: /bin/bash

Provided cores: 1 (use `--cores` to define parallelism)

Rules claiming more threads will be scaled down.

Job stats:

job	count
render	2
total	2

Select jobs to execute...

Execute 2 jobs...

Nextflow

Pipeline is defined as channels that flow into and out of processes.

notebooks.nf

```
1  params.notebook = "*.qmd"
2
3  workflow {
4      QMD_CH = Channel.fromPath(params.notebook)
5      HTML_CH = RENDER(QMD_CH)
6  }
7
8  process RENDER {
9
10     publishDir "."
11
12     input:
13     path qmd
14
15     output:
16     path("${qmd.baseName}.html")
17
```

Nextflow

Run workflow process by specifying the input parameters.

```
$ nextflow run notebooks.nf --notebook workflows.qmd -resume

N E X T F L O W ~ version 22.10.6
Launching `notebooks.nf` [lonely_jang] DSL2 - revision: 02a9f1b9b3
executor > local (1)
[c8/e651a4] process > RENDER (1) [100%] 1 of 1 ✓
```

Rerun the workflow

```
nextflow run notebooks.nf --notebook workflows.qmd -resume

N E X T F L O W ~ version 22.10.6
Launching `notebooks.nf` [reverent_venter] DSL2 - revision: 02a9f1b9b3
[c8/e651a4] process > RENDER (1) [100%] 1 of 1, cached: 1 ✓
```

Nextflow

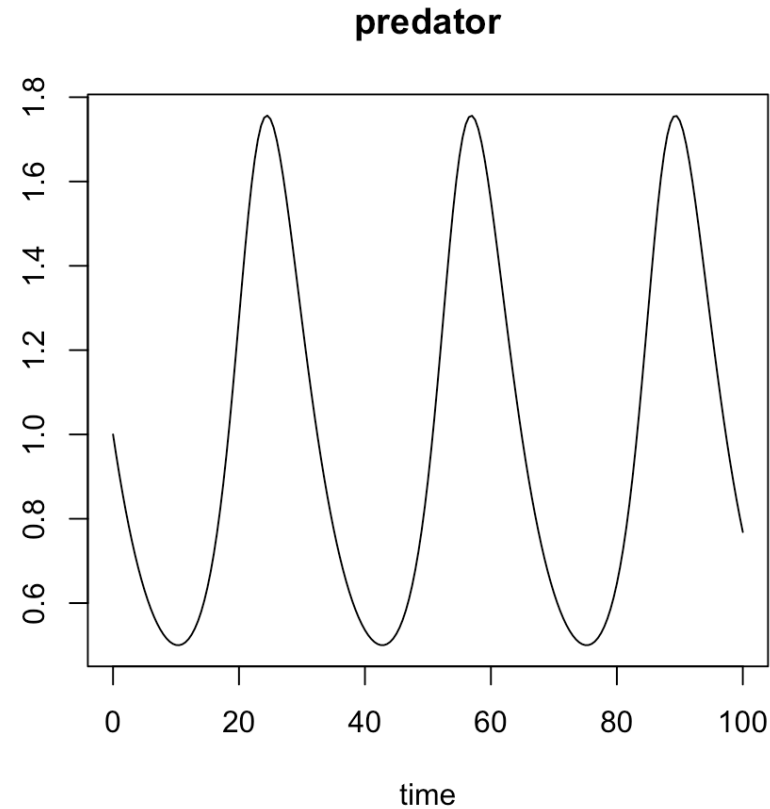
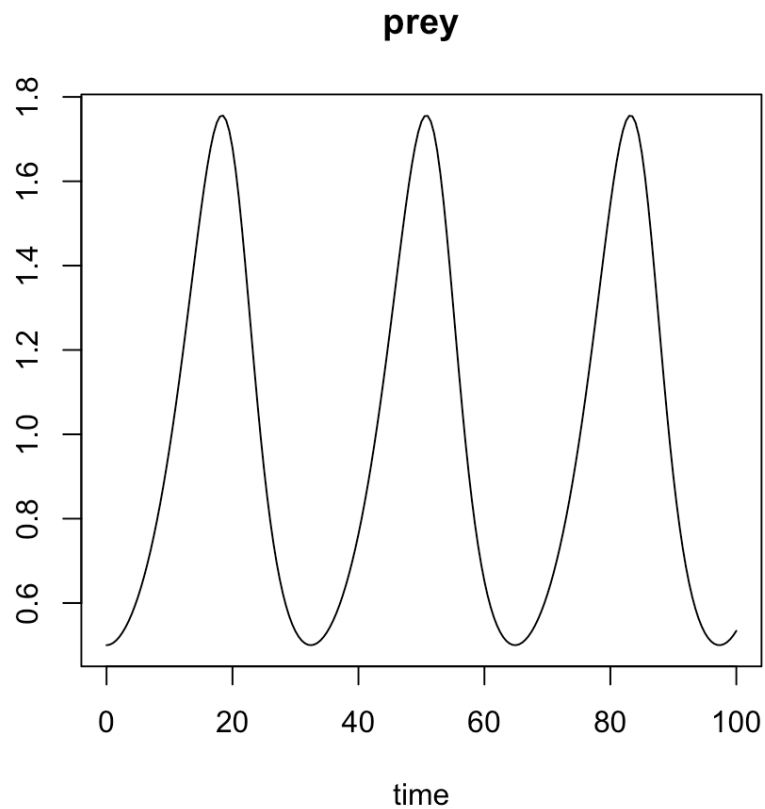
Run workflow by specifying multiple files for the input parameter.

```
$ nextflow run -resume notebooks.nf --notebook "{workflows,symposium,thesis}.qm

N E X T F L O W ~ version 22.10.6
Launching `notebooks.nf` [pedantic_heisenberg] DSL2 - revision: 02a9f1b9b3
executor > local (2)
[9e/fcbb56] process > RENDER (1) [100%] 3 of 3, cached: 1 ✓
```


Working with scripts

Example script that simulates a Lotka-Volterra model



Snakemake

Snakefile

```
1 rule lv:
2     output: "lv/prey{prey}-predator{pred}.png"
3     script: "scripts/lv.R"
```

Snakemake passes a special `snakemake` object to the script.

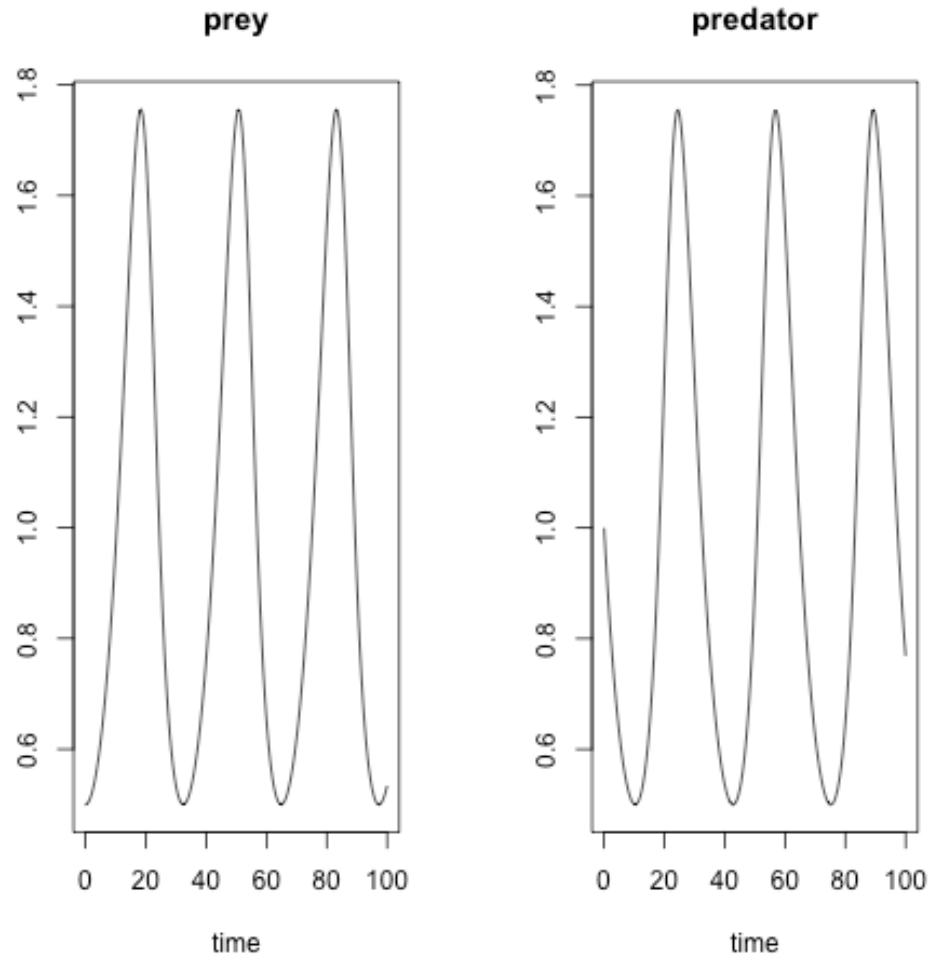
scripts/lv.R

```
1 library(simecol)
2
3 # parameterise simulation
4 # get parameter values from file wildcards
5 # prey and predator initial values
6 data(lv)
7 init(lv) <- c(preym = as.numeric(snakemake@wildcards$prey),
8               predm = as.numeric(snakemake@wildcards$pred))
9
10 # run simulation
11 simObj <- sim(lv)
12
13 # make and save plot to specified output
14 png(snakemake@output[[1]])
```

```
15 plot(simObj)
16 dev.off()
```

Snakemake

```
$ snakemake -j1 lv/prey0.5-predator1.0.png
```



Nextflow

lv.nf

```
1 workflow {
2     PREY_CH = Channel.of(params.prey)
3     PRED_CH = Channel.of(params.pred)
4     LV(PREY_CH, PRED_CH)
5 }
6
7 process LV {
8
9     publishDir 'lv', mode: 'copy'
10
11     input:
12     val prey
13     val pred
14
15     output:
16     path("prey${prey}-pred${pred}.png")
17
18     script:
19     """
```

Working with resources

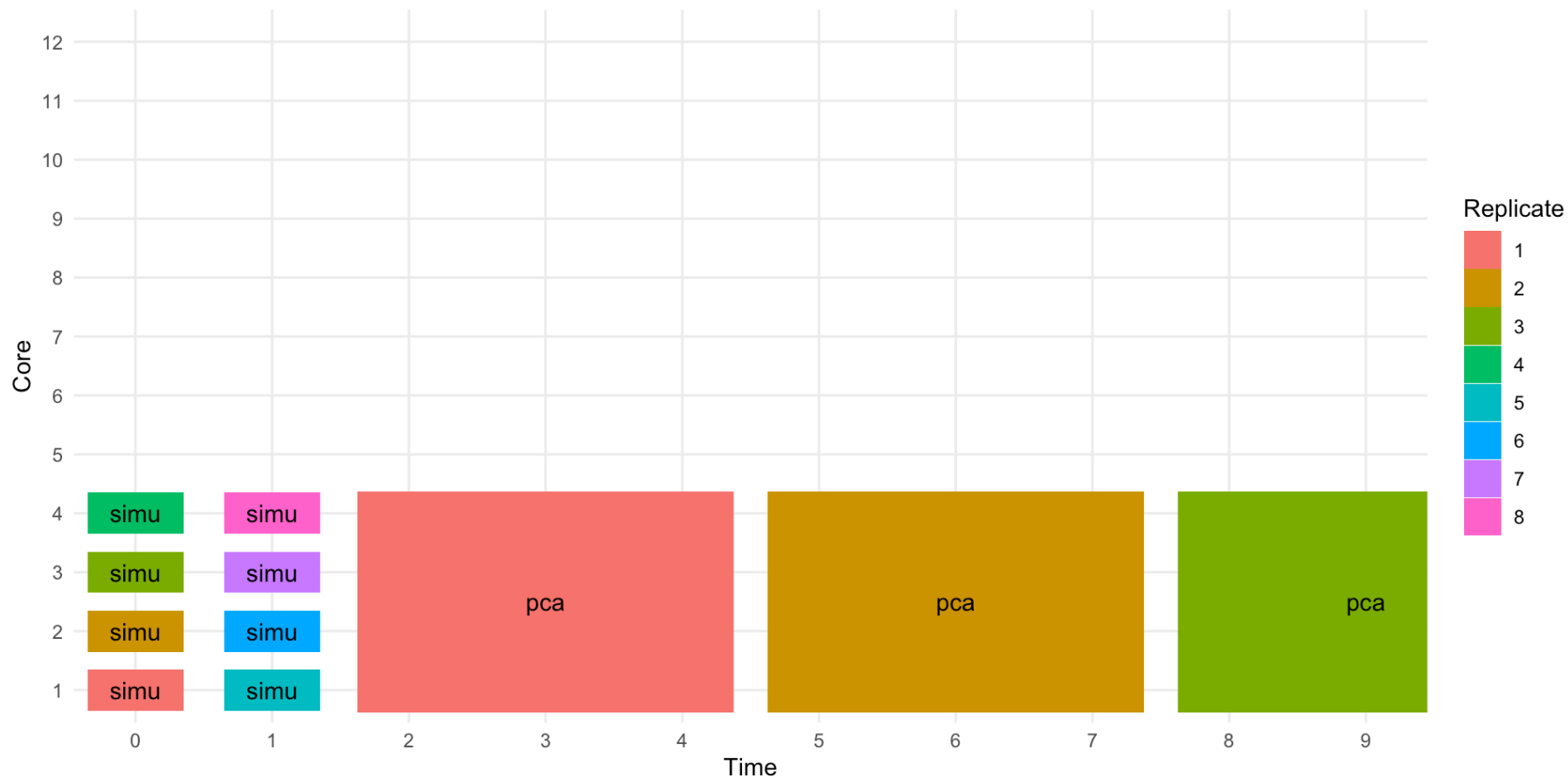
Snakefile

```
1 # Simulate 100k genotypes for 2k samples
2 rule simu:
3     output: multiext("sim/{bfile}", ".pgen", ".psam", ".pvar")
4     threads: 1
5     resources:
6         mem_mb=1000
7     shell: ""
8     plink2 --dummy 2000 100000 --out sim/{wildcards.bfile} \
9     --threads {threads} --memory {resources.mem_mb}
10    ""
11
12 # run principal components analysis on genetic data
13 rule pca:
14     input: multiext("sim/{bfile}", ".pgen", ".psam", ".pvar")
15     output: multiext("pca/{bfile}", ".eigenval", ".eigenvec")
16     threads: 4
17     resources:
18         mem_mb=8000
19     shell: ""
```

Generate and analyze 8 simulated datasets, utilising 4 cores

```
$ snakemake -c4 pca/sim-{1,2,3,4,5,6,7,8}.{eigenvec,eigenval}
```

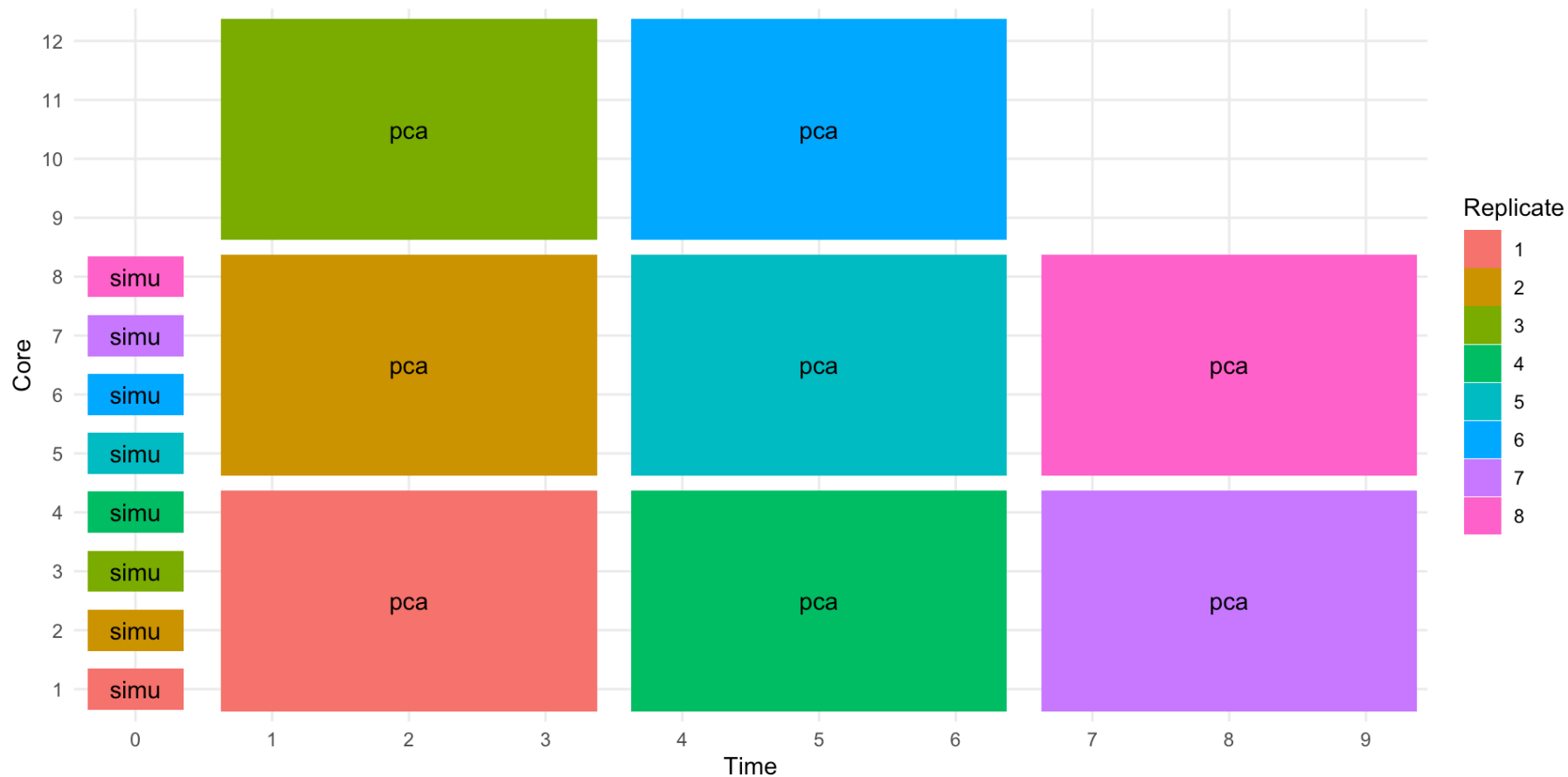
Core occupancy over time



Generate and analyze 8 simulated datasets, utilising **12** cores

```
$ snakemake -c12 pca/sim-{1,2,3,4,5,6,7,8}.{eigenvec,eigenval}
```

Core occupancy over time



Schedulers

Amazon Web Services, Azure, Google Cloud, HTCondor, Kubernetes, LSF, Torque, SGE, SLURM, etc

- Snakemake:
 - Profiles: github.com/Snakemake-Profiles
 - Plugins: snakemake.github.io/snakemake-plugin-catalog/
- Nextflow
 - Executors: nextflow.io/docs/latest/executor.html
 - Institution cluster configs: github.com/nf-core/configs

Working with schedulers

```
$ snakemake -jl --profile sge pca/sim-{1,2,3,4,5,6,7,8}.{eigenvec,eigenval}
```

```
Building DAG of jobs...
```

```
Using shell: /usr/bin/bash
```

```
Provided cluster nodes: 100
```

```
Job stats:
```

job	count
-----	-----
pca	8
simu	8
total	16

```
Select jobs to execute...
```

Working with software environments

Package managers, containers, and modules

Snakefile

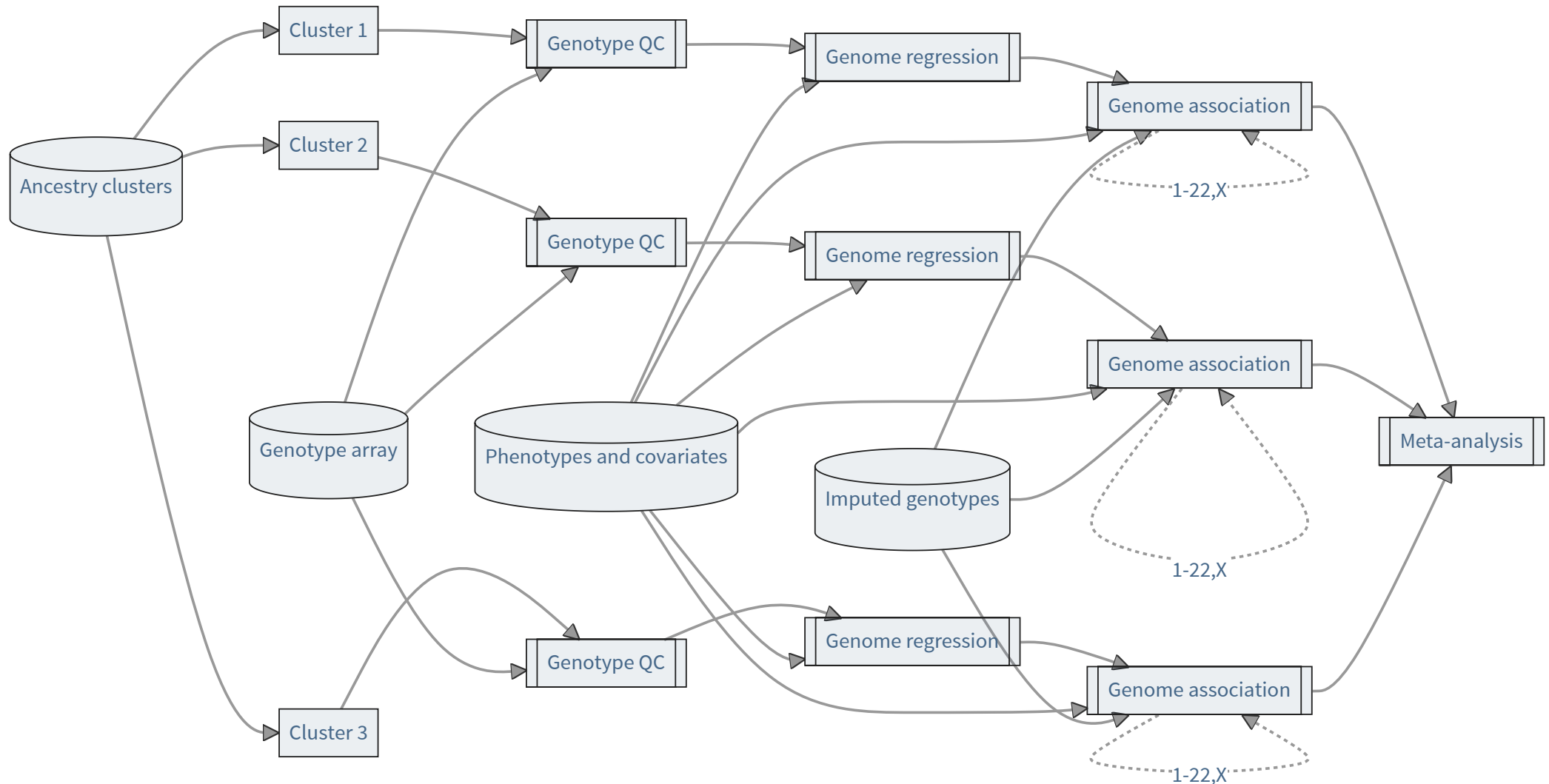
```
1 rule pca:
2     conda:
3         "envs/plink.yaml"
4     container:
5         "quay.io/biocontainers/plink2:2.00a2.3--hf22980b_0"
6     envmodule:
7         "igmm/apps/plink/2.00a3LM"
```

envs/plink.yaml

```
1 channels:
2     - conda-forge
3     - bioconda
4 dependencies:
5     - plink2 =2.00a2.3
```

Examples

Multi-ancestry genetic association study (UK Biobank, All of Us, etc)



Inputs

ukb-regenie-hrc.nf

```
1  params.bt = null // binary phenotypes file
2  params.qt = null // quantitative phenotypes file
3  params.keep = "rf_hgdplkg_clusters.keep"
4  params.remove = "PGC.remove"
5  params.bfile = "autosome.{bed,bim,fam}"
6  params.pfile = "ukb_imp_v3.qc.{pgen,psam,pvar}"
7  params.clusters = "ukb_randomforest_clusters.tsv"
8  params.covar = "ukb_randomforest_clusters.covar"
9  params.covar_list = "PC1,PC2,PC3,PC4,PC5,PC6"
10 params.covar_cat_list = "sex,genotyping"
11 params.min_cases = 80
```

Parse inputs from CSV/TSV file

```
ukb-regenie-hrc.nf
```

```
1 // genetic similarity clusters
2 CLUSTERS_CH = Channel
3   .fromPath(params.clusters, checkIfExists: true)
4
5 // parse cluster file to get names of each cluster
6 CLUSTER_NAMES_CH = CLUSTERS_CH
7   .splitCsv(sep: "\t", skip: 1, header: ['fid', 'iid', 'cluster'])
8   .map { it -> it.cluster }
9   .unique()
```

- `splitCsv()`: parse delimited file. Each row becomes an item.
- `map()`: get value of `cluster` column for each item.
- `unique()` output unique items.

Perform genotype QC for each ancestry cluster

ukb-regenie-hrc.nf

```
1 // Genotype QC
2 BFILE_CLUSTERS_CH = BFILE_CH
3   .combine(KEEP_CH)
4   .combine(REMOVE_CH)
5   .combine(CLUSTERS_CH)
6
7 QC_CH = QC(BFILE_CLUSTERS_CH, CLUSTER_NAMES_CH)
```

ukb-regenie-hrc.nf

```
1 process QC {
2     tag "QCing genotypes ${cluster}"
3
4     cpus = 1
5     memory = 16.GB
6     time = '1h'
7
8     input:
9     tuple val(bfile), path(bedbimfam), path(keep), path(remove), path(cluster)
10    each cluster
11
12    output:
13    tuple val(cluster), path("${cluster}.snplist"), path("${cluster}.id")
```

Prepend key to phenotype channels

ukb-regenie-hrc.nf

```
1 // binary
2 if(params.bt != null) {
3     BT_CH = Channel
4         .fromPath(params.bt, checkIfExists: true)
5         .map { it -> ["bt", it] }
6 } else {
7     BT_CH = Channel.empty()
8 }
9
10 // quantitative
11 if(params.qt != null) {
12     QT_CH = Channel
13         .fromPath(params.qt, checkIfExists: true)
14         .map { it -> ["qt", it] }
15 } else {
16     QT_CH = Channel.empty()
17 }
```

Specify different option flags for each phenotype

ukb-regenie-hrc.nf

```
1 STEP1_FLAGS_CH = Channel
2   .of(["bt", "--bt --minCaseCount ${params.min_cases}"],
3       ["qt", ""])
4   )
5
6 STEP2_FLAGS_CH = Channel
7   .of(["bt", "--bt --af-cc --firth --approx --pThresh 0.01 --minCaseCount $
8       ["qt", ""])
9   )
10
11 FLAGS_CH = STEP1_FLAGS_CH
12   .join(STEP2_FLAGS_CH)
13   .map { it -> [ it[0], ["step1": it[1], "step2": it[2]] ] }
```



```
[[ "bt" ] [ "step1": "--bt --minCaseCount 80", "step2": "--bt --af-cc --firth -
--approx --pThresh 0.01 --minCaseCount 80" ] ]
[[ "qt" ] [ "step1": "", "step2": "" ] ]
```

Combine phenotypes, program flags, and other input files

```
ukb-regenie-hrc.nf
```

```
1      PHENO_FLAGS_BFILE_CH = PHENO_COUNT_CH
2          .combine(FLAGS_CH, by: 0)
3          .combine(BFILE_CH)
4          .combine(QC_CH)
5          .combine(COVAR_CH)
6          .combine(COVAR_LIST_CH)
7          .combine(COVAR_CAT_LIST_CH)
```

Dynamic resource allocation

Genome regression step

ukb-regenie-hrc.nf

```
1 process STEP1 {
2     tag "${cluster}-${traits}-${pheno}"
3
4     // increase time on each attempt
5     // allocate 16GB of memory for each phenotype
6     cpus = 8
7     memory = { 8.GB + n_pheno * 16.GB }
8     time = { 6.hour * task.attempt }
9
10    input:
11    tuple val(cluster), val(traits), val(pheno), val(n_pheno), path(phenos)
```

Parallelise across chromosomes

Genome association step, then merge

ukb-regenie-hrc.nf

```
1 CHR_CH = Channel.of(1..23)
2 STEP2_CH = STEP2(STEP1_FLAGS_PFILE_CH, CHR_CH)
3 GWAS_CH = MERGE(STEP2_CH)
```

ukb-regenie-hrc.nf

```
1 process STEP2 {
2     tag "${cluster}-${traits}-${pheno}"
3
4     cpus = 8
5     memory = 16.GB
6     time = '24h'
7
8     input:
9     tuple val(cluster), val(traits), val(pheno), path(phenos), val(flags),
10    each chr
11
12    output:
13    tuple val(cluster), val(traits), val(pheno), path("step2_${cluster}-${t
```

Multiple levels of parallelisation

- Each input file can contain multiple phenotypes
- Workflow accepts multiple phenotype files
- Separate analysis for each ancestry cluster
- Separate analysis for each chromosome
- Regression/association steps run across multiple cores

Process UK Biobank release

Turns UKB release object into a DuckDB database

github.com/mja/ukb-release-nf

```
$ nextflow run duckdb.nf -c custom.config -resume --enc ukb12345.enc --key k123
```

- Decrypts release file
- Downloads data dictionary
- Determines which fields are in the release file
- Converts fields for each category
- Processes fields through R
- Combines tables into SQL database

Major depressive disorder GWAS meta-analysis

From the Psychiatric Genomics Consortium

[github.com/psychiatric-genomics-consortium/mdd-wave3-
meta](https://github.com/psychiatric-genomics-consortium/mdd-wave3-meta)

Determine inputs/outputs from configuration file

config.yaml

```
1 sumstats:
2   daner_mdd_GenScot.eur.hg19.SCID_0721a: data/sumstats/daner_mdd_genscot.SC
3   text_mdd_FinnGen.eur.hg38.R5_18032020: data/sumstats/Depressio_FinnGen_R5
4   text_mdd_UKBB.eur.hg19.MD_glm_202107: data/sumstats/ukb_mdd.202107.md.eu
```

meta.smk

```
1 # Copy summary statistics listed in config.yaml under sumstats
2 # with key FORMAT_COHORT.POP.hgNN.RELEASE
3 rule stage_sumstats:
4   input: lambda wildcards: config["sumstats"][wildcards.cohort]
5   output: "resources/sumstats/{cohort}.gz"
6   log: "logs/sumstats/stage/{cohort}.log"
7   shell: "ln -sv {input} {output} > {log}"
8
9 # Harmonize names of all summary statistics listed under sumstats in config
10 rule sumstats:
11   input: expand("resources/sumstats/{sumstats}.gz", sumstats=config["sums
```


Conclusions

- Workflow as a record of analysis steps: **reproducibility**
- Workflow as a re-runnable pipeline: **replicability**
- Workflow as a resource manager: **scalability**
- Workflow as a compute environment: **portable**

Code for this talk: github.com/mja/tech-talk-workflows

List of pipeline frameworks:

github.com/pditommaso/awesome-pipeline