

Bioinformatics Workflows with Snakemake

Johannes Köster

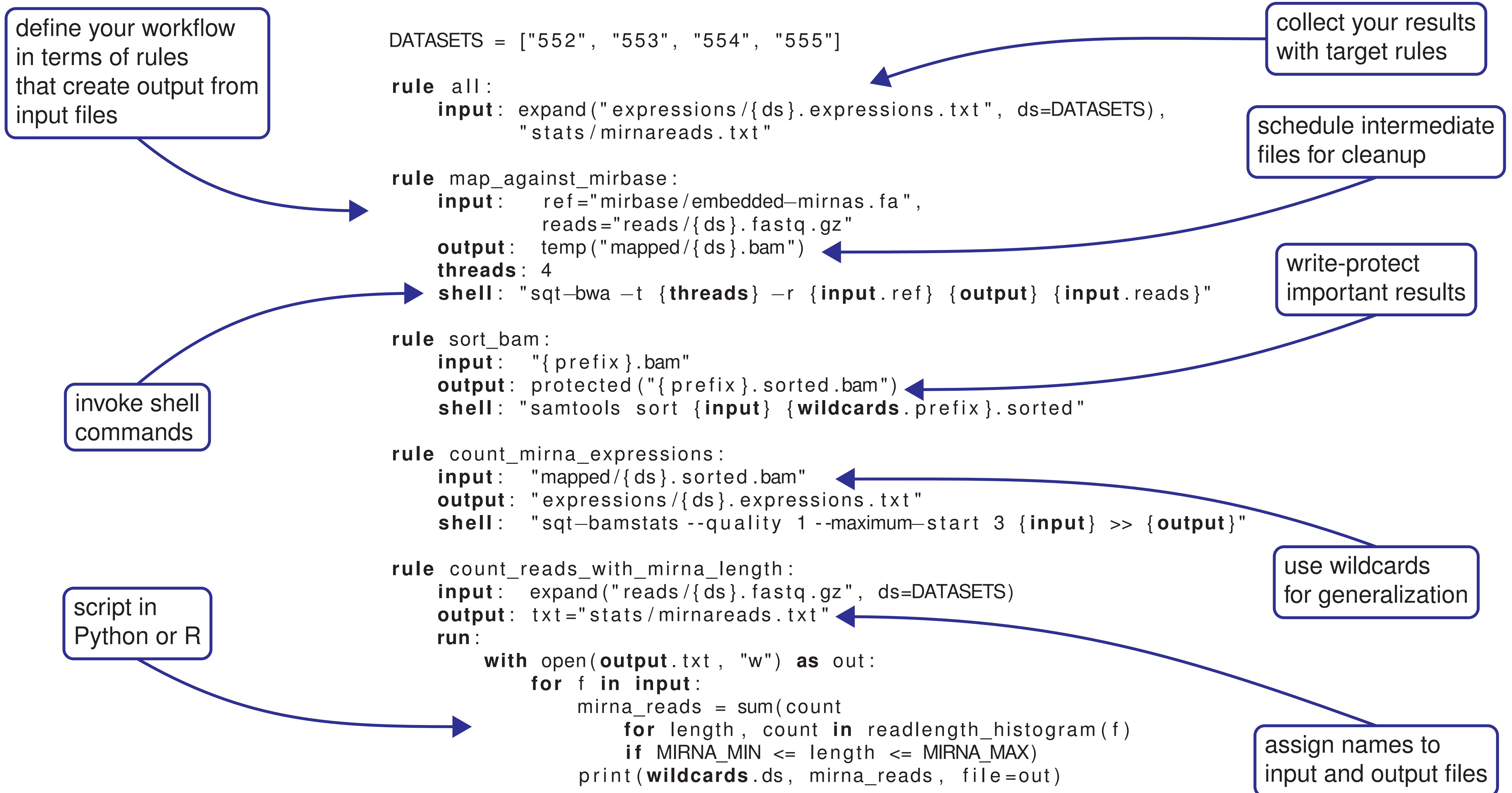
Genome Informatics, Institute of Human Genetics, University of Duisburg-Essen
Paediatric Oncology, University Hospital Essen

CONTRIBUTION

In most bioinformatics analyses, many different files are created by the combination of various tools. A sophisticated workflow management eases the inclusion of new datasets and parameter changes and ensures reproducibility. Snakemake is a workflow engine that combines an easy to read text-based definition language with a versatile execution environment that is scalable from single core machines over multi-core servers to clusters.

IDEA

In Snakemake, a workflow is defined in terms of rules, that create output files from input files by applying a command line tool or executing Python or R:



FEATURES

Run the whole workflow (no arguments) or create specific targets (rule or filenames):

```
$ snakemake mapped/552.bam
```

Execute jobs in parallel up to a given limit of usable cores:

```
$ snakemake --cores 40
```

To run Snakemake on a cluster, only a shared filesystem and a qsub-like submission command is needed. Here, parameters of the rules can be forwarded to the clustering engine:

```
$ snakemake --cluster "qsub -pe threaded {threads}"
```

Summarize all involved files, listing pending updates based on modification dates, tool versions and implementation changes:

```
$ snakemake --summary | less
```

Visualize the workflow via graphviz dot:

```
$ snakemake --dag | dot | display
```

Force certain rules to re-run:

```
$ snakemake --forcerun count_mirna_expressions
```

Prioritize a target in the scheduler (for example an urgent dataset):

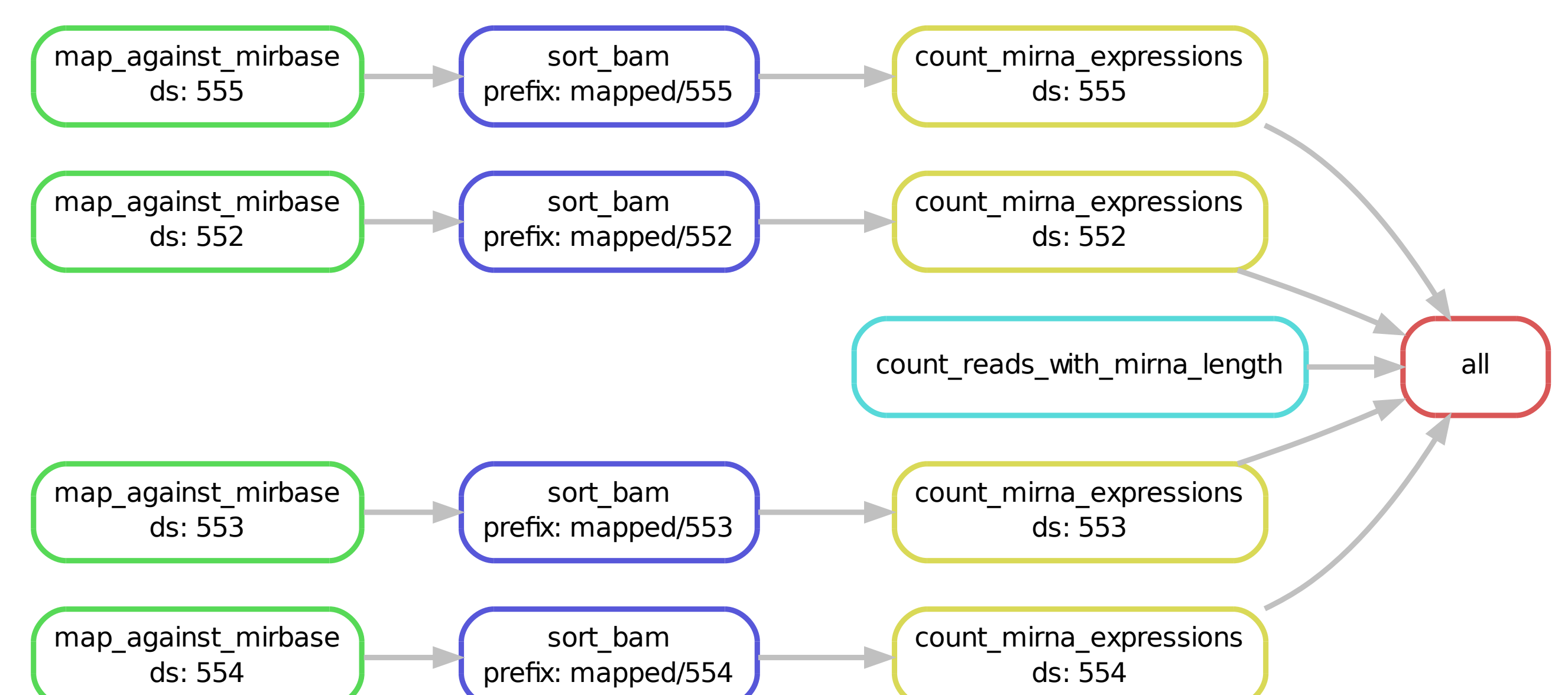
```
$ snakemake --prioritize expressions/552.expressions.txt
```

Finally, Snakemake provides facilities for

- dynamic updates depending on a rule's output,
- portable HTML5 reports that provide a semantically connected collection of generated figures and tables for collaborators,
- creation of rule libraries that can be included into workflows.

FROM RULES TO JOBS

Dependencies between the rules are determined automatically by matching input files against output files. This yields a directed acyclic graph of jobs (DAG) that represents the execution plan of Snakemake. Jobs are only executed if their output files are not present, at least one input file is newer than an output file, or if execution is forced via the command line.



Jobs that lie on the same path are dependent on each other. Independent jobs can be executed in parallel. Snakemake provides a scheduling mechanism that maximizes the usage of given cores with respect to threads, priorities and input file sizes by solving a knapsack problem.

Johannes Köster, Sven Rahmann. "Snakemake - A scalable bioinformatics workflow engine". Bioinformatics 2012.

visit Snakemake at
<https://bitbucket.org/johanneskoester/snakemake>

