

Snakemake

Johannes Köster
Genome Informatics, Institute of Human Genetics,
Faculty of Medicine, University Duisburg-Essen

April 10, 2014

Structure

① Motivation

② Basic Idea

③ Advanced Features

Outline

① Motivation

② Basic Idea

③ Advanced Features

Motivation

What we liked about GNU Make:

- text based
- rule paradigm
- lightweight

And what not:

- cryptic syntax
- limited scripting
- multiple output files
- scalability

Snakemake

- hook into python interpreter
- pythonic syntax for rule definition
- full python scripting
- scalability
- workflow specific functionality beyond Make basics
- stable community:

Outline

① Motivation

② Basic Idea

③ Advanced Features

Syntax

```
SAMPLES = "500 501 502 503".split()
```

```
# require a bam for each sample
```

```
rule all:  
  input:  
    expand("{sample}.bam", sample=SAMPLES)
```

```
# map reads
```

```
rule map:  
  input:  
    "reference.bwt",  
    "{sample}.fastq"  
  output:  
    "{sample}.bam"  
  threads: 8  
  shell:  
    "bwa mem -t {threads} {input} | " # refer to threads and input files  
    "samtools view -Sbh - > {output}" # refer to output files
```

```
# create an index
```

```
rule index:  
  input:  
    "reference.fasta"  
  output:  
    "reference.bwt"  
  shell:  
    "bwa index {input}"
```

Basic Usage

perform a dry-run

```
$ snakemake -n
```

execute the workflow using 8 cores

```
$ snakemake -j 8
```

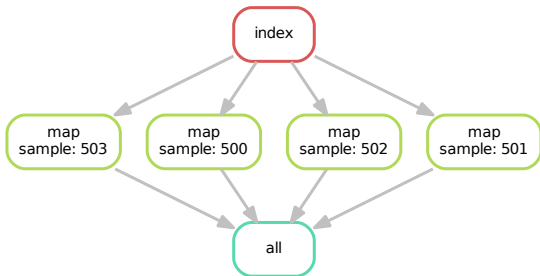
execute the workflow on a cluster (with up to 20 jobs)

```
$ snakemake -j 20 --cluster "qsub -pe threaded {threads}"
```


Visualization

visualize the DAG of jobs

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



Outline

① Motivation

② Basic Idea

③ Advanced Features

Advanced Syntax

```
SAMPLES = "500 501 502 503".split()
```

```
rule all:  
  input:  
    expand("{sample}.bam", sample=SAMPLES)
```

map reads with peanut

```
rule map:  
  input:  
    "reference.hdf5",  
    "{sample}.fastq"  
  output:  
    "{sample}.bam"  
  threads: 8  
  resources: gpu=1 # define an additional resource  
  version: shell("peanut --version")  
  shell:  
    "peanut map -t {threads} {input} | "  
    "samtools view -Sbh - > {output}"
```

create an index with peanut

```
rule index:  
  input:  
    "reference.fasta"  
  output:  
    "reference.hdf5"  
  shell:  
    "peanut index {input} {output}"
```

Scheduling

Maximize the number of running jobs with respect to

- priority
- number of descendants
- input size

while not exceeding

- provided cores
- provided resources

A multi-dimensional knapsack problem.

Sub-Workflows

```
SAMPLES = "500 501 502 503".split()
```

```
# define subworkflow
```

```
subworkflow:  
  workdir: "../mapping"
```

```
rule all:  
  input:  
    expand("{sample}/results.xprs", sample=SAMPLES)
```

```
# estimate transcript expressions
```

```
rule express:  
  input:  
    REF,  
    mapping("{sample}.bam") # refer to output of subworkflow  
  output:  
    "{sample}/results.xprs"  
  shell:  
    "express {input} -o {wildcards.sample}"
```

HTML5 Reports

```
from snakemake.utils import report
```

```
rule report:
```

```
input:
```

```
T1="results.csv",  
F1="plot.pdf"
```

```
output:
```

```
html="report.html"
```

```
run:
```

```
report(""
```

```
=====  
Some Title  
=====
```

```
See table T1_, display  
some math
```

```
.. math::  
|cq_0 - cq_1| > {MDIFF}
```

```
"", output.html, **input)
```

snakemake report

Combining expressions of lncRNAs measured by qPCR with HuEx exon arrays

HuEx and qPCR datasets and their combination

The datasets were given: An lncRNA assay using qPCR [T1] and 274 primary tumors analyzed with Affymetrix HuEx exon arrays within the MIC.

ENSEMBL Gene IDs for lncRNAs were extracted from given lncRNA qPCR assay description [T2]. For these, accessions were derived from the ENSEMBL lncRNA annotation track. Affymetrix HuEx ID probes that had either been not been identified, and combined to make junctions [T3]. Here, each row depicts an lncRNA given as ENSEMBL ID together with all the HuEx probes that should measure the expression of said lncRNA.

Each of these reads probes summarizes the expression of one lncRNA. We calculated and normalized the expressions for the given 274 primary tumors. This was done with the Affymetrix Power Tools implementation of RMA with default parameters. Figure F1 shows the histogram of log₂ expressions. It remains to be investigated if the RMA normalization has successfully removed batch effects since the tumor data comes from different labs.

Estimation of regulated lncRNAs in the qPCR dataset

We estimate the consistency between the two controls by calculating the log fold change and throwing away all lncRNAs that exceed a threshold of 3 in this field. For the remaining lncRNAs the log change between treatment and the mean of the two controls is calculated.

Table T2 shows unregulated lncRNAs sorted by strength of fold change. Table T3 shows the same for downregulated lncRNAs. Figure F1 shows the histogram of log fold changes.

Counting tumors expressing the regulated lncRNAs in the HuEx dataset

We only consider those lncRNAs that can be measured by exonarray probes (see Table T3). For these we calculate from the HuEx data described above the number of tumors with a minimum predicted expression of 6. Further, we test whether the mean of the log₂ exonarray expression is significantly higher than 6 using a one-sided test: the p-values computed after Bonferroni correction can be found in the column `prob` expressed in the tables T4 and T5. Finally, the Pearson correlation of ENSG00000120012 with each of the candidates is computed, since a negative correlation with ENSG00000120012 should point towards expression in the tumors (corrected p-values in the column `prob` computed). Tables T6 and T7 show the results for up- and downregulated lncRNAs. Figure F2 and F3 show histograms of the obtained counts. In theory, positive fold changes in the qPCR assay should correspond to underexpression with ENSG00000120012 in the microarray, and negative fold changes in the qPCR assay should correspond to positive correlation in the microarray. Figure F4 and F5 show heatmaps of tumor data in general.

Mathematical background

The provided qPCR analysis yielded cq_0 values that are on a logarithmic scale compared to the real molecule counts. This is because each qPCR cycle in theory doubles the amount of molecules. Since this rate is not reached in practice we assume a factor of 2.8 here. Consequently a fold change on these log-scaled values has to be computed as subtraction instead of a quotient. Further, a non-logarithmic fold change of cq_0 corresponds to $\ln(cq_0)$ in logarithmic scale. This is however only applicable to not normalized cq_0 values which are not available. In consequence the log fold change provided here should only be compared relatively to each other.

We denote the different cq_0 values as $cq_0^1, cq_0^2, \dots, cq_0^k$ for a given lncRNA i . The consistency between the two controls (see above) is now computed as

$$|cq_{0,1} - cq_{0,2}| > 8.5$$

The log fold change between treatment and the mean of controls is calculated as

$$cq_{0,i} - \frac{1}{2} (\ln(cq_{0,1} + cq_{0,2}))$$

Since the latter cq_0 values are all logarithmic, the mean here corresponds to the geometric mean of the real molecule counts. This is intended since a log-normal distribution of the mean by the higher cq_0 value.

A one-sided one sample t-test was used to determine whether the mean of exonarray expression is significantly higher than 6. The test is applicable here independently of the actual distribution of provided expressions because the number of probesets is sufficiently large such that the mean follows a gaussian distribution.

[T1] [lncdiffchange.html#tq1](#)
[T2] [lncdiffchange.html#tq2](#)
[T3] [lncdiffchange.html#tq3](#)
[T4] [lncdiffchange.html#tq4](#)
[T5] [lncdiffchange.html#tq5](#)
[T6] [lncdiffchange.html#tq6](#)
[T7] [lncdiffchange.html#tq7](#)
[T8] [lncdiffchange.html#tq8](#)
[T9] [lncdiffchange.html#tq9](#)
[T10] [lncdiffchange.html#tq10](#)
[T11] [lncdiffchange.html#tq11](#)
[T12] [lncdiffchange.html#tq12](#)
[T13] [lncdiffchange.html#tq13](#)
[T14] [lncdiffchange.html#tq14](#)
[T15] [lncdiffchange.html#tq15](#)

Data Provenance

Summarize output file status

```
$ snakemake --summary
```

file	date				rule	version	status	plan	
500.bam	Thu	Apr	10	10:55:17	2014	map	1.0	ok	no update
501.bam	Thu	Apr	10	10:55:17	2014	map	1.0	ok	no update
502.bam	Thu	Apr	10	10:55:17	2014	map	1.0	updated input files	update pending
503.bam	Thu	Apr	10	10:55:17	2014	map	0.9	version changed to 1.0	no update

Trigger updates:

```
# update files with changed versions
```

```
$ snakemake -R `snakemake --list-version-changes`
```

```
# update files with changed code
```

```
$ snakemake -R `snakemake --list-code-changes`
```

Conclusion

Snakemake is a Make-like workflow system providing

- a readable syntax
- sophisticated scripting with python
- scalability from single-core to cluster
- support for hybrid computing
- data provenance
- modularization capabilities

Roadmap:

- DRMAA support
- a workflow or rule library

<http://bitbucket.org/johanneskoester/snakemake>

Köster, J., Rahmann, S., Snakemake – a scalable bioinformatics workflow engine.
Bioinformatics 2012.