

News about Snakemake

Johannes Köster

February 5, 2013

gatk

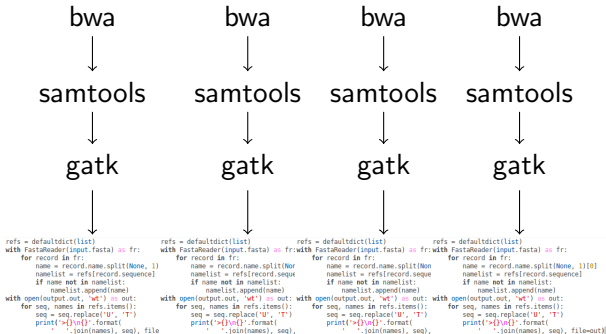
bwa

samtools

```
refs = defaultdict(list)
with FastqHeader(input.fasta) as fr:
    for record in fr:
        name = record.name.split(None, 1)[0]
        name_list = refs[record.sequence]
        if name not in name_list:
            name_list.append(name)
with open(output.out, 'w') as out:
    for seq, names in refs.items():
        seq = seq.replace('U', 'T')
        print('>{}|n{}'.format(
            __.join(names), seq), file=out]
```

bwa
↓
samtools
↓
gatk

```
refs = defaultdict(list)
with FastaReader(input.fasta) as fr:
    for record in fr:
        name = record.name.split(None, 1)[0]
        nameList = refs[record.sequence]
        if name not in nameList:
            nameList.append(name)
with open(output.out, 'w') as out:
    for seq, names in refs.items():
        seq = seq.replace('U', 'T')
        print("{}={}".format(
            __.join(names), seq), file=out])
```



bwa
↓
samtools -q1
↓
gatk
↓

```
refs = defaultdict(list)
with FastaReader(input.fasta) as fr:
    for record in fr:
        name = record.name.split(None, 1)[0]
        namelist = refs[record.sequence]
        if name not in namelist:
            namelist.append(name)
with open(output.out, 'wt') as out:
    for seq, names in refs.items():
        seq = seq.replace('U', 'T')
        print("{}|{}".format(
            __, '.join(names), seq), file=out)
```

bwa
↓
samtools -q1
↓
gatk
↓

```
refs = defaultdict(list)
with FastaReader(input.fasta) as fr:
    for record in fr:
        name = record.name.split(None, 1)[0]
        namelist = refs[record.sequence]
        if name not in namelist:
            namelist.append(name)
with open(output.out, 'wt') as out:
    for seq, names in refs.items():
        seq = seq.replace('U', 'T')
        print("{}|{}".format(
            __, '.join(names), seq), file=out)
```

bwa
↓
samtools -q1
↓
gatk
↓

```
refs = defaultdict(list)
with FastaReader(input.fasta) as fr:
    for record in fr:
        name = record.name.split(None, 1)[0]
        namelist = refs[record.sequence]
        if name not in namelist:
            namelist.append(name)
with open(output.out, 'wt') as out:
    for seq, names in refs.items():
        seq = seq.replace('U', 'T')
        print("{}|{}".format(
            __, '.join(names), seq), file=out)
```

bwa
↓
samtools -q1
↓
gatk
↓

```
refs = defaultdict(list)
with FastaReader(input.fasta) as fr:
    for record in fr:
        name = record.name.split(None, 1)[0]
        namelist = refs[record.sequence]
        if name not in namelist:
            namelist.append(name)
with open(output.out, 'wt') as out:
    for seq, names in refs.items():
        seq = seq.replace('U', 'T')
        print("{}|{}".format(
            __, '.join(names), seq), file=out)
```

GNU Make provided us with...

- a language to write rules to create each output file from input files
- wildcards for generalization
- implicit dependency resolution
- implicit parallelization
- fast and collaborative development on text files

GNU Make provided us with...

- a language to write rules to create each output file from input files
- wildcards for generalization
- implicit dependency resolution
- implicit parallelization
- fast and collaborative development on text files

but we missed...

- easy to read syntax
- simple scripting inside the workflow
- creating more than one output file with a rule
- multiple wildcards in filenames

1 Idea

2 Scheduling

3 Key Features

Snakemake Idea

Example: for samples {500, ..., 503} map reads to hg19.

Snakemake Idea

Example: for samples {500,...,503} map reads to hg19.

```
rule map_reads:
    input:  "hg19.fasta", "{sample}.fastq"
    output: "{sample}.sai"
    shell:  "bwa aln {input} > {output}"
```

Snakemake Idea

Example: for samples {500,...,503} map reads to hg19.

```
rule sai_to_bam:
    input: "hg19.fasta", "{sample}.sai", "{sample}.fastq"
    output: "{sample}.bam"
    shell:
        "bwa samse {input} | samtools view -Sbh - > {output}"
```

```
rule map_reads:
    input: "hg19.fasta", "{sample}.fastq"
    output: "{sample}.sai"
    shell: "bwa aln {input} > {output}"
```

Snakemake Idea

Example: for samples {500,...,503} map reads to hg19.

```
SAMPLES = ["500", "501", "502", "503"]
```

```
rule all:
```

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

```
rule sai_to_bam:
```

```
    input: "hg19.fasta", "{sample}.sai", "{sample}.fastq"
```

```
    output: "{sample}.bam"
```

```
    shell:
```

```
        "bwa samse {input} | samtools view -Sbh - > {output}"
```

```
rule map_reads:
```

```
    input: "hg19.fasta", "{sample}.fastq"
```

```
    output: "{sample}.sai"
```

```
    shell: "bwa aln {input} > {output}"
```

Snakemake Idea

Example: for samples {500,...,503} map reads to hg19.

```
SAMPLES = ["500", "501", "502", "503"]
```

```
rule all:
```

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

```
rule sai_to_bam:
```

```
    input: "hg19.fasta", "{sample}.sai", "{sample}.fastq"
```

```
    output: protected("{sample}.bam")
```

```
    shell:
```

```
        "bwa samse {input} | samtools view -Sbh - > {output}"
```

```
rule map_reads:
```

```
    input: "hg19.fasta", "{sample}.fastq"
```

```
    output: "{sample}.sai"
```

```
    shell: "bwa aln {input} > {output}"
```

Snakemake Idea

Example: for samples {500,...,503} map reads to hg19.

```
SAMPLES = ["500", "501", "502", "503"]
```

```
rule all:
```

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

```
rule sai_to_bam:
```

```
    input: "hg19.fasta", "{sample}.sai", "{sample}.fastq"
```

```
    output: protected("{sample}.bam")
```

```
    shell:
```

```
        "bwa samse {input} | samtools view -Sbh - > {output}"
```

```
rule map_reads:
```

```
    input: "hg19.fasta", "{sample}.fastq"
```

```
    output: temp("{sample}.sai")
```

```
    shell: "bwa aln {input} > {output}"
```

Snakemake Idea

Example: for samples {500,...,503} map reads to hg19.

```
rule all
  500.bam, 501.bam, 502.bam, 503.bam
```

```
rule sai_to_bam:
  input: "hg19.fasta", "{sample}.sai", "{sample}.fastq"
  output: protected("{sample}.bam")
  shell:
    "bwa samse {input} | samtools view -Sbh - > {output}"
```

```
rule map_reads:
  input: "hg19.fasta", "{sample}.fastq"
  output: temp("{sample}.sai")
  shell: "bwa aln {input} > {output}"
```

Snakemake Idea

Example: for samples {500,...,503} map reads to hg19.

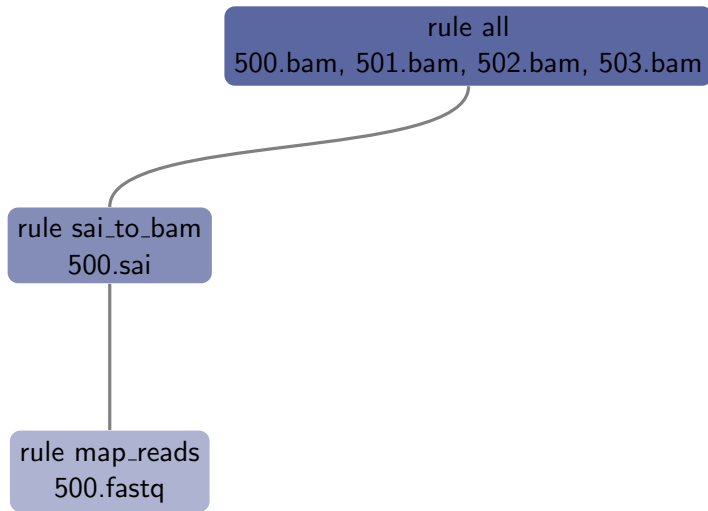
```
rule all
  500.bam, 501.bam, 502.bam, 503.bam
```

```
rule sai_to_bam
  500.sai
```

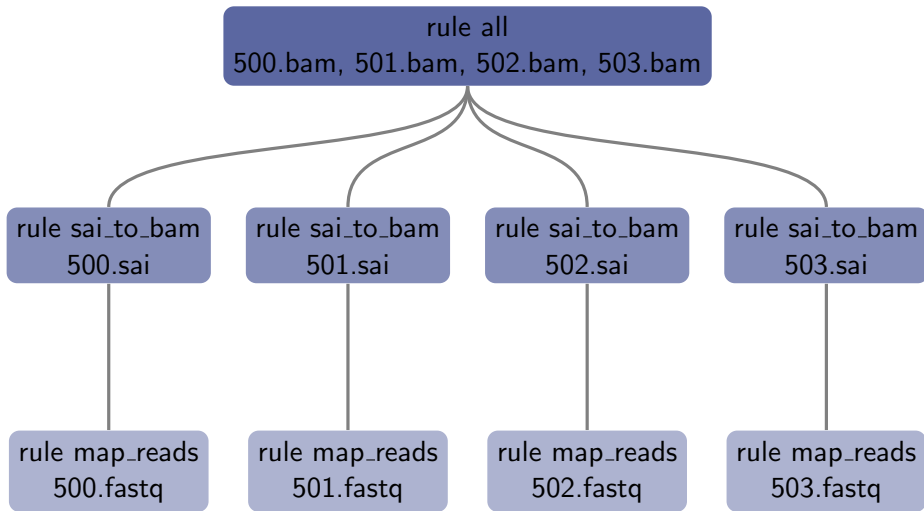
```
rule map_reads:
  input:  "hg19.fasta", "{sample}.fastq"
  output: temp("{sample}.sai")
  shell:  "bwa aln {input} > {output}"
```

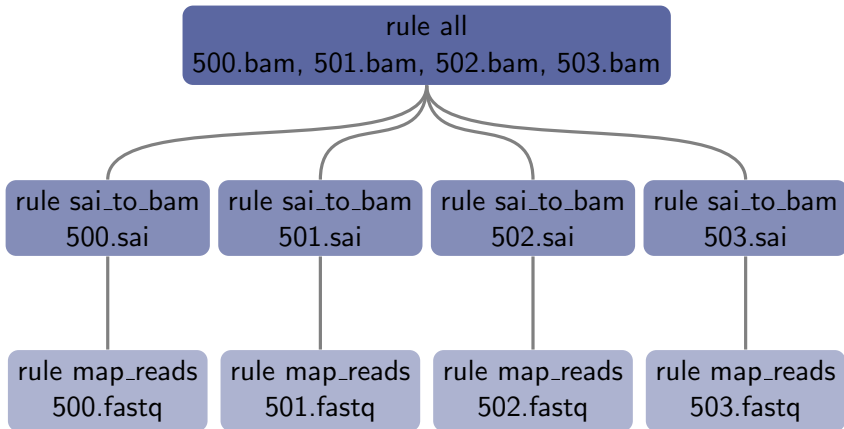

Snakemake Idea

Example: for samples {500,...,503} map reads to hg19.



Example: for samples {500,...,503} map reads to hg19.





- DAG of jobs
- each path needs to be executed serially
- two disjoint paths can be executed in parallel

an edge between two jobs A,B if input of A is matched by output of B, e.g.

"500.bam" matches "{sample}.bam"
 \Leftrightarrow
"500.bam" $\in L(".+\.bam")$

In case of ambiguity:

- Constrain wildcards: "{sample, [0-9]+}.bam"
- Order rules: ruleorder: sai_to_bam > sort_bam

execute the set of jobs E^* among all $E \subseteq J$ that maximizes under lexicographical order

$$\sum_{j \in E} (p_j, i_j)$$

such that

$$\sum_{j \in E} t_j \leq l$$

- J set of jobs ready to execute
- T provided cores
- l idle cores
- t_j threads of job j
- p_j priority of job j
- i_j input size of job j

```
rule plot_coverage_histogram:
    input: "{sample}.bam"
    output: hist="{sample}.coverage.pdf"
    run:
        plt.hist(np.fromiter(
            shell("samtools mpileup {input} | cut -f4",
                iterable=True),
            dtype=int))
        plt.savefig(output.hist)
```

R Rules

```
rule plot_coverage_histogram:  
  input: ...  
  output: ...  
  run:  
    R("""  
    # some R code  
    """)
```

Dynamic Output

In case of unknown number of output files...
dynamically update the DAG to process these

Dynamic Output

In case of unknown number of output files...
dynamically update the DAG to process these

```
rule all:  
  input: dynamic("{cluster}.pdf")
```

```
rule plot:  
  input: "{cluster}.csv"  
  output: "{cluster}.pdf"  
  shell: "gnuplot ..."
```

```
rule cluster:  
  input: ...  
  output: dynamic("{cluster}.csv")  
  shell: "cluster ..."
```

Dynamic Output

In case of unknown number of output files...
dynamically update the DAG to process these

```
rule all
dynamic("{cluster}.pdf")
```

```
rule plot:
  input: "{cluster}.csv"
  output: "{cluster}.pdf"
  shell: "gnuplot ..."
```

```
rule cluster:
  input: ...
  output: dynamic("{cluster}.csv")
  shell: "cluster ..."
```

Dynamic Output

In case of unknown number of output files...
dynamically update the DAG to process these

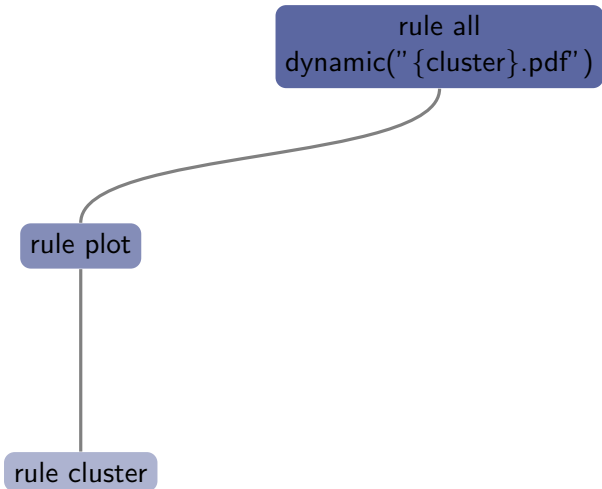
```
rule all  
dynamic("{cluster}.pdf")
```

```
rule plot
```

```
rule cluster:  
  input: ...  
  output: dynamic("{cluster}.csv")  
  shell: "cluster ..."
```

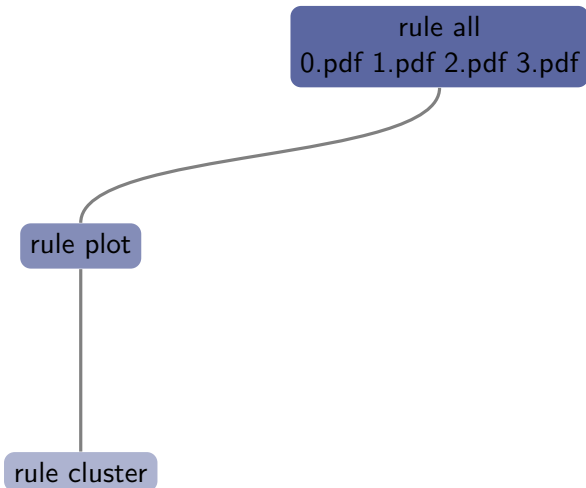
Dynamic Output

In case of unknown number of output files...
dynamically update the DAG to process these



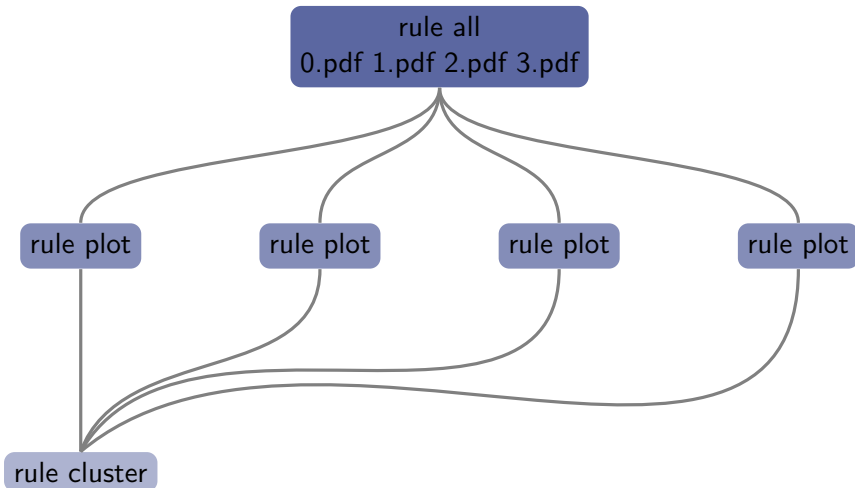
Dynamic Output

In case of unknown number of output files...
dynamically update the DAG to process these



Dynamic Output

In case of unknown number of output files...
dynamically update the DAG to process these



HTML Reports

rule report:

```
input: "table.csv", "plot.pdf"
```

```
output: "report.html"
```

```
run:
```

```
report("""
```

```
=====
```

```
Report of some project
```

```
=====
```

```
Some text containing a formula
```

```
:math:  $\sum_j \ln E t_j \leq I$ 
```

```
and embedding a table F1_ and a figure F2_.
```

```
""", output[0], F1=input[0], F2=input[1])
```

snakemake report

Report of some project

Some text containing a formula $\sum_{j \in E} t_j \leq I$ and embedding a table [F1](#) and a figure [F2](#).

2013-01-31

Combining expressions of lncRNAs measured by qPCR with HuEx exon arrays

HuEx and qPCR datasets and their combination

Two datasets were given: An lncRNA assay using qPCR [1] and 274 primary tumors analysed with Affymetrix HuEx exon arrays within the NHC.

ENSEMBL Gene IDs for lncRNAs were extracted from given lncRNA qPCR assay description [1]. For these, exonic loci were derived from the ENSEMBL hg19 v09 annotation track. Affymetrix HuEx 1.0 probes that lie within these loci were identified, and combined to meta-probesets [2]. Here, each row depicts an lncRNA given as ENSEMBL ID together with all the HuEx probes that should measure the expression of one of its exons.

Each of these meta-probesets 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. 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 fold-change and throwing away all lncRNAs that exceed a threshold in this test. For the remaining lncRNAs the fold-change between treatment and the mean of the two controls is calculated.

Table 14 shows upregulated lncRNAs sorted by strength of fold-change, Table 15 shows the same for downregulated lncRNAs. Figure 1 shows the histogram of fold-changes.

Counting tumors expressing the regulated lncRNAs in the HuEx dataset

We assume lncRNAs with an absolute logarithmic fold-change greater than 0.69 to be regulated.

We only consider those lncRNAs that can be measured by exonarray probes (see Table 16). For these we calculate from the HuEx setup described above the number of tumors with a minimum probeset expression of 6. Table 17 shows the results, Figure 2 shows a histogram of the observed counts.

Mathematical background

The provided qPCR analysis yielded -cq values that are on a logarithmic scale compared to the real molecule counts. This is because each PCR cycle in theory doubles the amount of molecules. Since this rate is not reached in practice we assume a factor of 1.8 here.

Consequently a fold-change on these logspaced values has to be computed as subtraction instead of a quotient. Further, a non-logarithmic foldchange of f corresponds to $\log_{1.8} f$ in logarithmic scale. It has to be noted that at the moment it is not yet clear whether the normalization applied before has an effect on the assumed factor.

We denote the different -cq values as α_{tr}^i , α_{tr}^i , and α_{tr}^i for a given lncRNA i . The consistency between the two controls (see above) is now computed as

$$|\alpha_{tr}^i - \alpha_{tr}^i| > 0.5$$

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

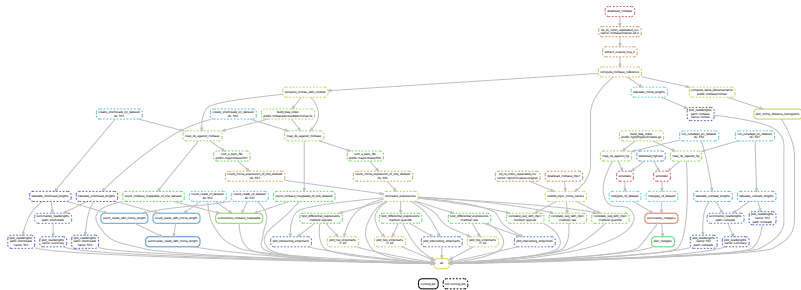
$$\alpha_{tr}^i - \frac{1}{2} (\alpha_{tr}^i + \alpha_{tr}^i)$$

Since the latter -cq values are still logarithmic, the mean here corresponds to the geometric mean of the real molecule counts. This is intended since it avoids domination of the mean by the higher -cq value.

2013-02-01

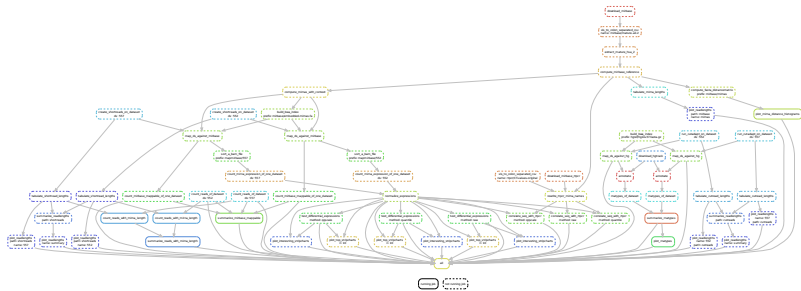
Using the dot language of graphviz:

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

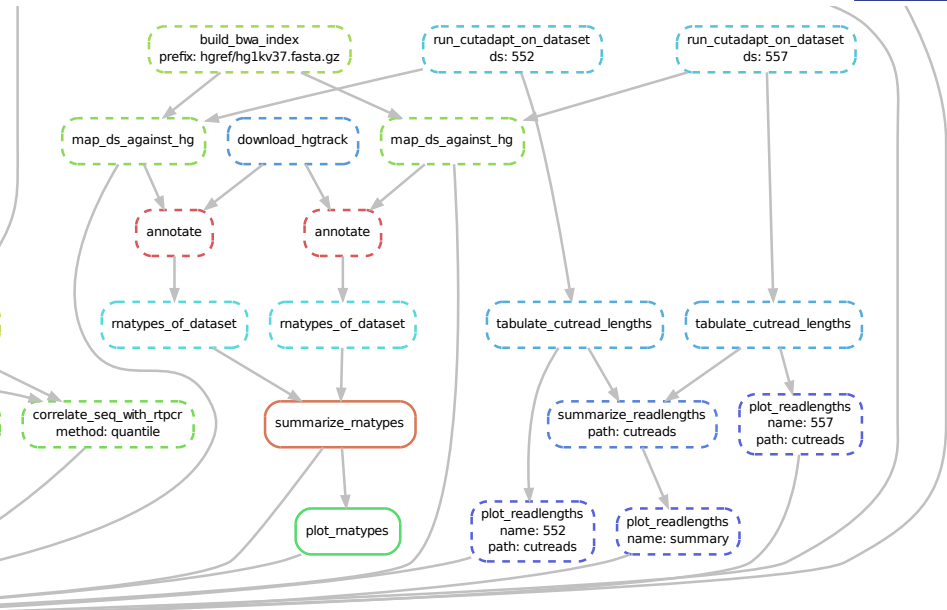


Using the dot language of graphviz:

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



Visualization



Conclusion

Snakemake is a new workflow system that provides:

- an easy pythonic textual representation
- multiple wildcards in filenames
- dynamic update of job DAG
- implicit parallelization and dependency resolution
- job scheduling considering threads, priorities and input size
- cluster and batch support

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

depends on Python ≥ 3.2

Early adopters:

Shirley Liu's lab, Broad Institute
Genome of the Netherlands Project