

Snakemake workflows deployed on virtual environments: a promising way of integrating high-throughput data in RegulonDB

Claire Rioualen¹, Lucie Khamvongsa¹, Alberto Santos-Zavaleta², Mishael Sánchez-Pérez², Jocelyn Brayet³, Julio Collado-Vides², Jacques van Helden¹

1. **Technologies avancées pour le génome et la clinique** (TAGC), INSERM U1090, Aix-Marseille Université, Marseille F-13288, France
2. **Computational Genomics Research program**, Center for Genomics Sciences - UNAM, México
3. **Institut Curie**, PSL Research University, Mines Paris Tech, Inserm, U900, F-75005, Paris, France.

Next-generation sequencing (NGS) has become a mainstream technology in genomics, and it has gotten increasingly cheaper and faster to obtain genomic data. However, the development of reliable tools for the analysis of the huge amount of data generated is still lagging behind¹.

Snakemake, a flexible environment for workflow development

- **python** library for building workflows
- inherits concepts from the **GNU make** software:
 - **target** files or operations to be performed
 - **rules** describing how to produce these targets
- operations can be done in **python**, **R**, **shell** languages
- dependencies between rules are defined by their **inputs** and **outputs**
- **wildcards** can be used for automatization

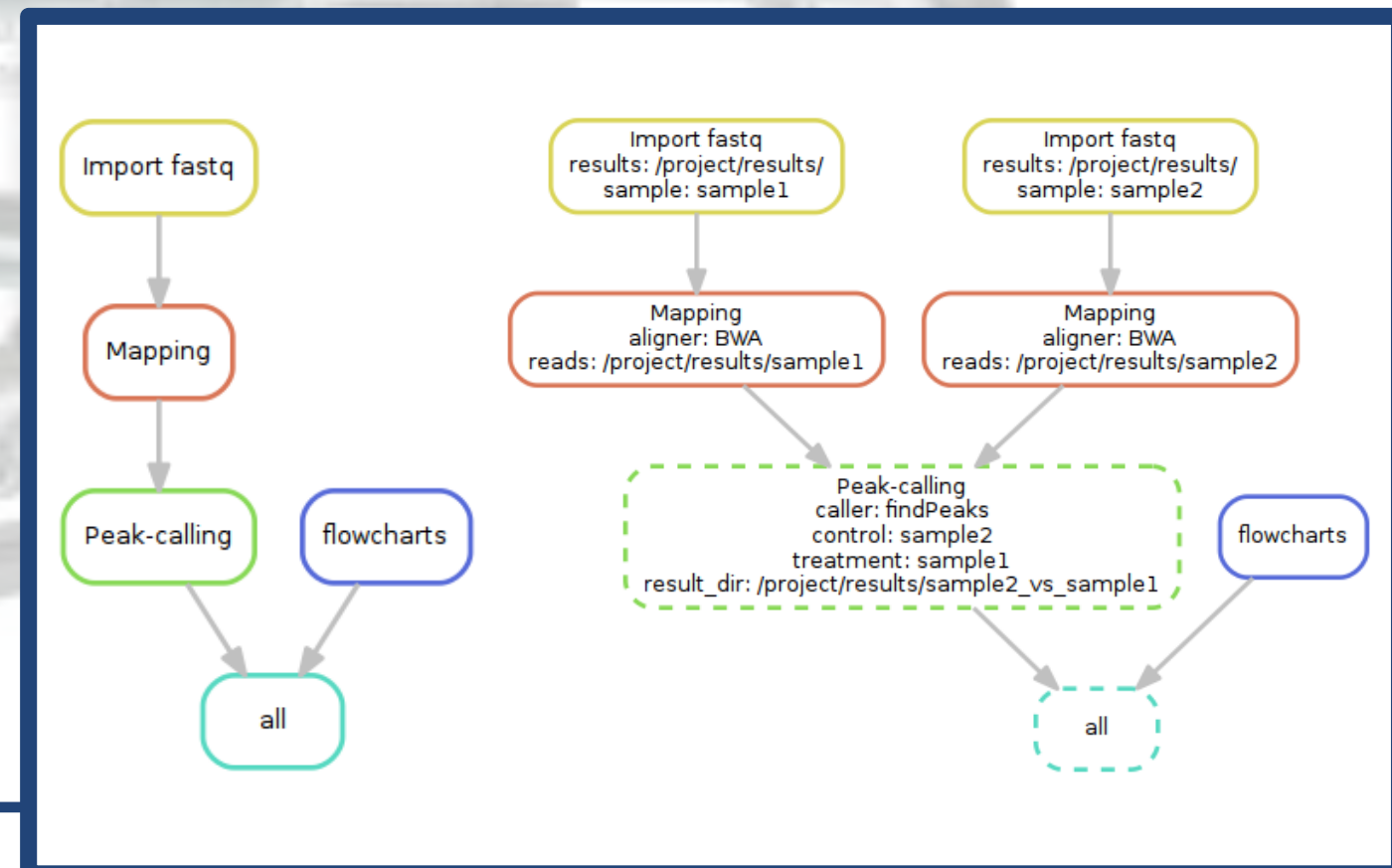
Workflow example

```
## workflow.py
SAMPLES = ["sample1", "sample2"]
CHIP = "sample2"
CONTROL = "sample1"

rule all:
    input:
        expand("{sample}.bam", sample = SAMPLES)
        expand("{chip}_vs_{control}",
            chip = CHIP, control = CONTROL)
    output:
        chip = CHIP, control = CONTROL

rule mapping:
    input:
        "{file}.fastq"
    output:
        "{file}.bam"
    shell:
        "bowtie {input} > {output}"

rule peak_calling:
    input:
        chip = "{chip}.bam", control = "{chip}.bam"
    output:
        "{file}.bam"
    shell:
        "macs2 {input.chip} {input.control} > {output}"
```



We developed a public library of re-usable rules⁵ which can be combined into different workflows: trimming, mapping, peak-calling, FastQC reports, motif search, IGV sessions... Ongoing is the development of rules to handle RNA-seq data and pipelines combining both types of data.

ChIP-seq pipeline: E. coli study case compared with literature and RegulonDB data



RegulonDB (<http://regulondb.ccg.unam.mx/>) is the primary database on transcriptional regulation in *Escherichia coli* K-12 containing knowledge manually curated from original scientific publications, experimental evidence, complemented with computational predictions.

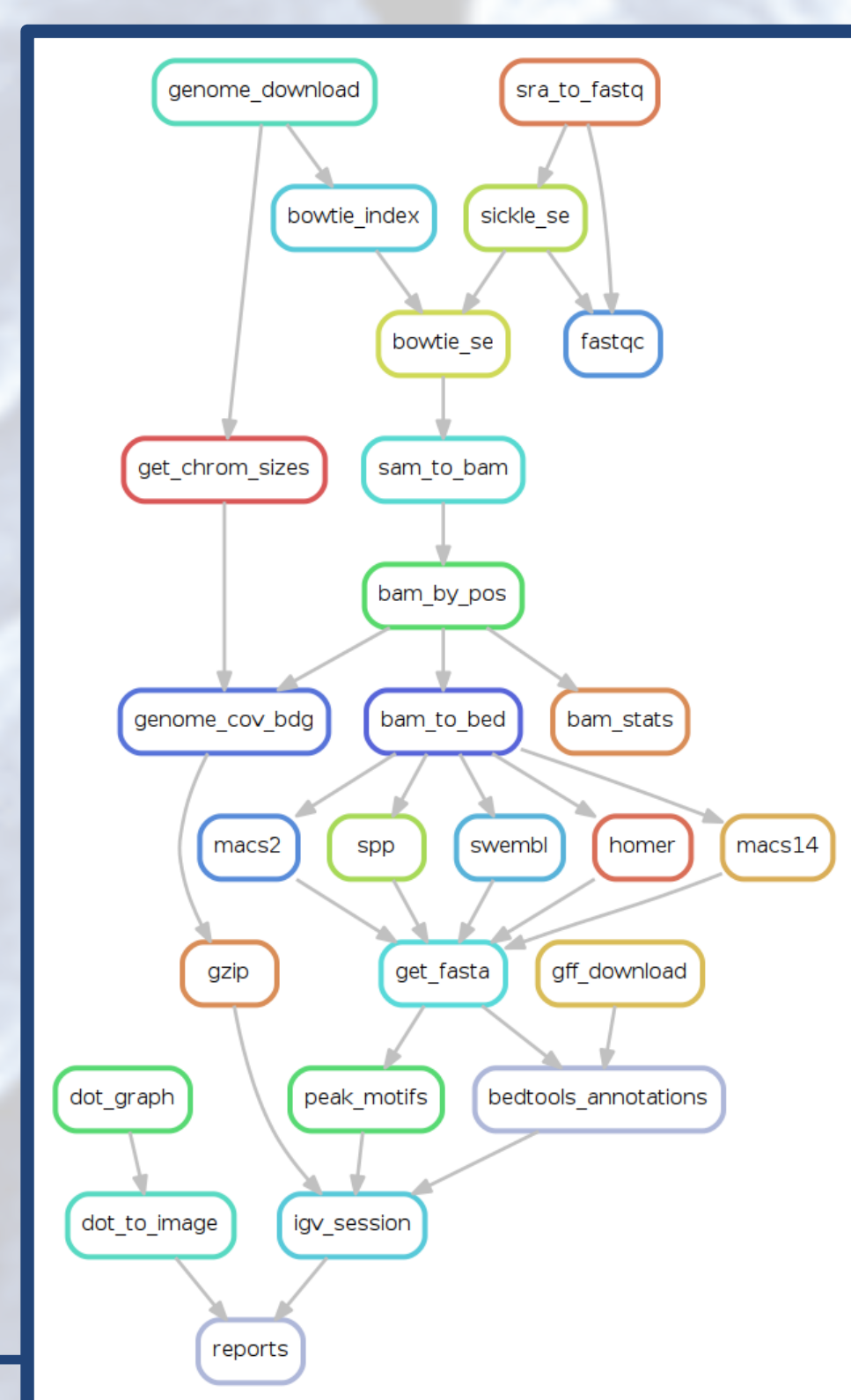
This database could be further developed by integrating high-throughput data from experiments such as ChIP-seq or RNA-seq. Thus we decided to confront our pipeline to the data curated in RegulonDB, and with the results published by Myers et al.³ which combined ChIP-seq and RNA-seq to detect direct target genes of **FNR** in aerobic and **anaerobic conditions**.

We ran the following workflow on samples GSM1010219 (FNR IP ChIP-seq Anaerobic A) and GSM1010224 (Anaerobic input DNA), from GEO subseries GSE41187.

Reads were trimmed using Sickle, and Bowtie was chosen for the alignment after Bowtie2, BWA and subread showed similar mapping rates with twice as much computing time.

Applying all the operations above to samples GSM1010219 and GSM1010224 took about **24mn** to complete on a virtual machine with 32 Go RAM.

Flowchart generated by snakemake for the ChIP-seq workflow

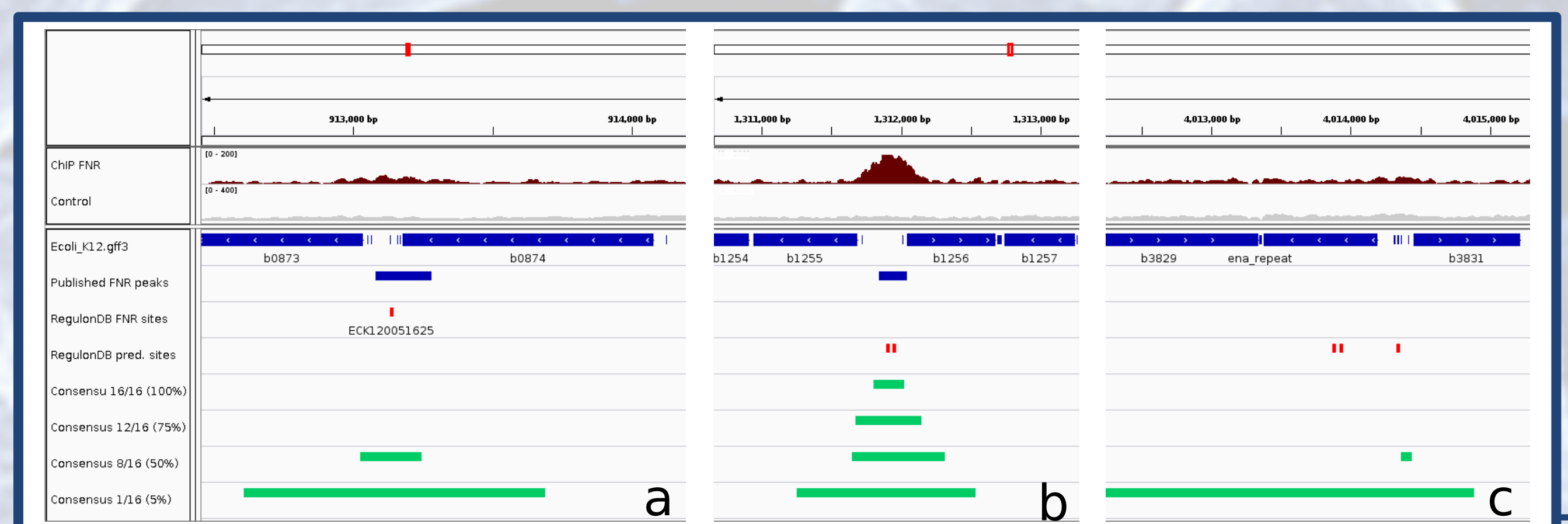


Thanks to the Snakemake environment, we were able to run 5 different peak-callers, using a variety of parameters.

Myers et al.³ ran 3 peak-callers: CisGenome, Mosaics and NCIS.

IGV session generated by the workflow

We kept 16 bed files whose peak count was between 100 and 500, and generated "consensus peak sets" using several thresholds. These peak sets were compared with published peaks³ and with the motifs found in RegulonDB⁴.



- 49 peaks correlated with known sites from RegulonDB, out of 79 sites manually checked (a)
- 29 peaks came as **new evidence** for RegulonDB predicted sites in both published data³ and our pipeline (b)
- 16 predicted sites were supported by our **consensus peaks only** (c), meaning that many peak-callers couldn't back them up.
- More than 200 peaks were not associated with known or predicted sites in RegulonDB. These peaks will be further investigated by running a motif search algorithm available in our pipeline⁹ and combining **RNA-seq** data to check for **expression changes**

References

1. Bailey T, et al. (2013) Practical Guidelines for the Comprehensive Analysis of ChIP-seq Data. PLoS Comput Biol 9(11): e1003326. doi:10.1371/journal.pcbi.1003326
2. Johannes Köster and Sven Rahmann. Snakemake: a scalable bioinformatics workflow engine. Bioinformatics (2012) 28 (19): 2520-2522. doi: 10.1093/bioinformatics/bts480
3. Myers KS, et al. Genome-scale Analysis of *Escherichia coli* FNR Reveals Complex Features of Transcription Factor Binding. PLoS Genet (2013) 9(6): e1003565. doi:10.1371/journal.pgen.1003565 - Data: <http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE41187>
4. Gama-Castro S et al. RegulonDB version 9.0: high-level integration of gene regulation, coexpression, motif clustering and beyond. J.Nucleic Acids Res. 2016 Jan 4;44(D1):D133-43. doi: 10.1093/nar/gkv1156. Epub 2015 Nov 2.
5. GitHub Gene-regulation: <https://github.com/rioualen/gene-regulation>
6. Docker image: <https://hub.docker.com/r/rioualen/gene-regulation>
7. Institut Français de Bioinformatique: <http://france-bioinformatique.fr>
8. VirtualBox software: <https://www.virtualbox.org/>
9. Medina-Rivera A et al. (2015) RSAT 2015: Regulatory Sequence Analysis Tools. Nucleic Acids Res. 2015 (Web Server issue) in press.

Acknowledgments: This work is funded by France Génomique. Collaboration with the RegulonDB team was possible thanks to NIH grant GM0110597 & FOINS-CONACYT - Fronteras de la Ciencia 2015 - ID 15