

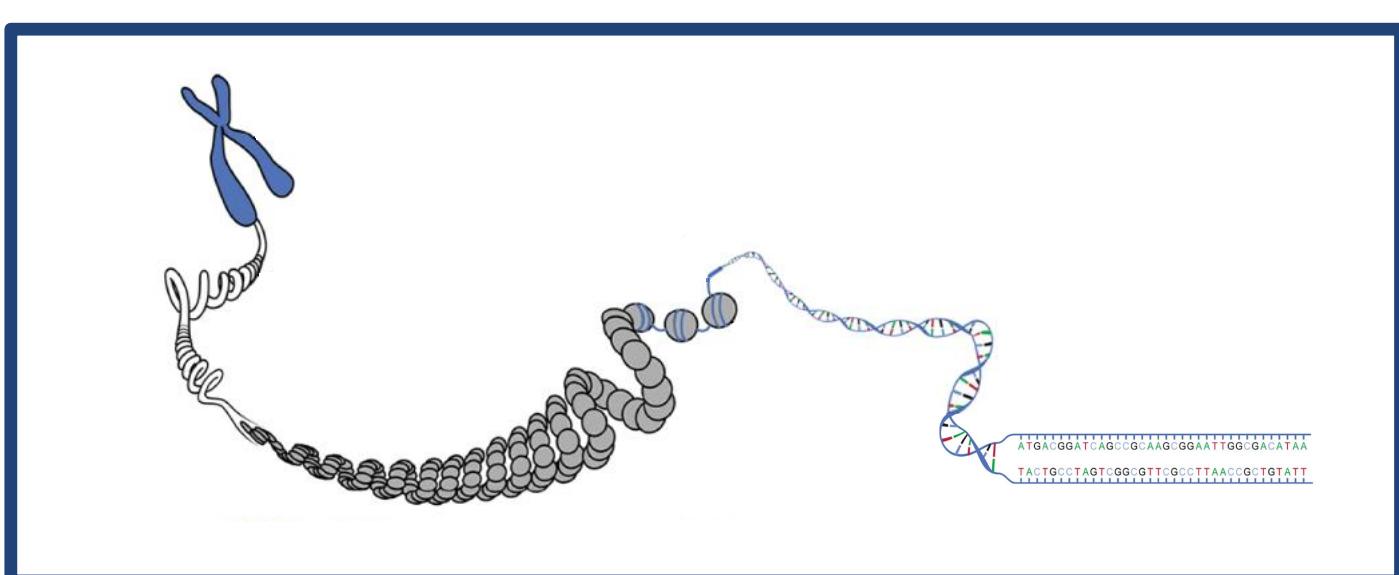
Gene-regulation project: high-throughput approaches for robust NGS data analyses

Claire Rioualen, Jacques van Helden

Theory and Approaches of Genomic Complexity (TAGC), INSERM U1090, Aix-Marseille Université, Marseille F-13288, France

Objectives & Approach

Next generation sequencing (NGS) has gotten faster and cheaper in the past few years. What we are lacking now are means to analyse the incredibly huge amount of data generated.



We have decided to handle this question. Our main focuses include:

- **reproducibility** of the biological analyses;
- **traceability** of the tools and parameters used;
- **flexibility** of the framework;
- **portability** in different operating systems.

A simple framework

```
Data description
; Sample description file
; Organism: Saccharomyces cerevisiae
; Experiment type: ChIP-seq
; ID: ENA experiment Scan Name
GSM521934 input SRR021359 SRR051938
GSM521935 chip SRR021358 SRR051929

; Design for the peak-calling of ChIP-seq
treatment control description
GSM521935 GSM521934 TBFL_occupancy

Config of the analysis
metadata:
samples: "metadata/samples.tab"
design: "metadata/design.tab"
control: "metadata/config.yml"
conf_type: "yaml"

dir:
fastq: "fastq"
genotype: "genotype"
results: "results"

tools:
trimming: "sickle"
mapping: "bowtie2"
peakcalling: "macs2 homer macs14 spp swembl"
bowtie2:
threads: "+"
max_mismatches: "1"

; Scripts for the peak-calling of ChIP-seq
; treatment - control, description
; GEM521935 - GEM521934 TBFL_occupancy
; Design for the peak-calling of ChIP-seq
; treatment control description
; GEM521935 GEM521934 TBFL_occupancy

Workflow design
include: config["gene regulation"]

# Samples
SAMPLES = read.table(config$metadata)[["sample"]][,"ID"]
TREATMENT = read.table(config$metadata)[["design"]][,"treatment"]
CONTROL = read.table(config$metadata)[["design"]][,"control"]

# Rules
include: "scripts/snakefiles/rules/bowtie2.rules"
include: "scripts/snakefiles/rules/fastqc.rules"
include: "scripts/snakefiles/rules/macs2.rules"
include: "scripts/snakefiles/rules/homer.rules"

# Flowcharts
FLOWCHARTS = expand_fastq(samples)/samples/.fastq", samples=SAMPLE_IDS)
IMPORT = expand_fastq(samples)/samples/.fastqc(samples).fastqc.html", samples=SAMPLE_IDS)
MAPPING = expand_fastq(samples)/samples/.bam", samples=SAMPLE_IDS)
PEAK_CALL = expand_fastq(samples)/samples/.bed", samples=SAMPLE_IDS)

# Workflow execution
rule all:
    all_the_required_analyses -->
    input: FLOWCHARTS, IMPORT, QUALITY_CONTROL, MAPPING, PEAK_CALL
    params: qsub = config[qsub]
```

Snakemake¹, a flexible environment for workflow development

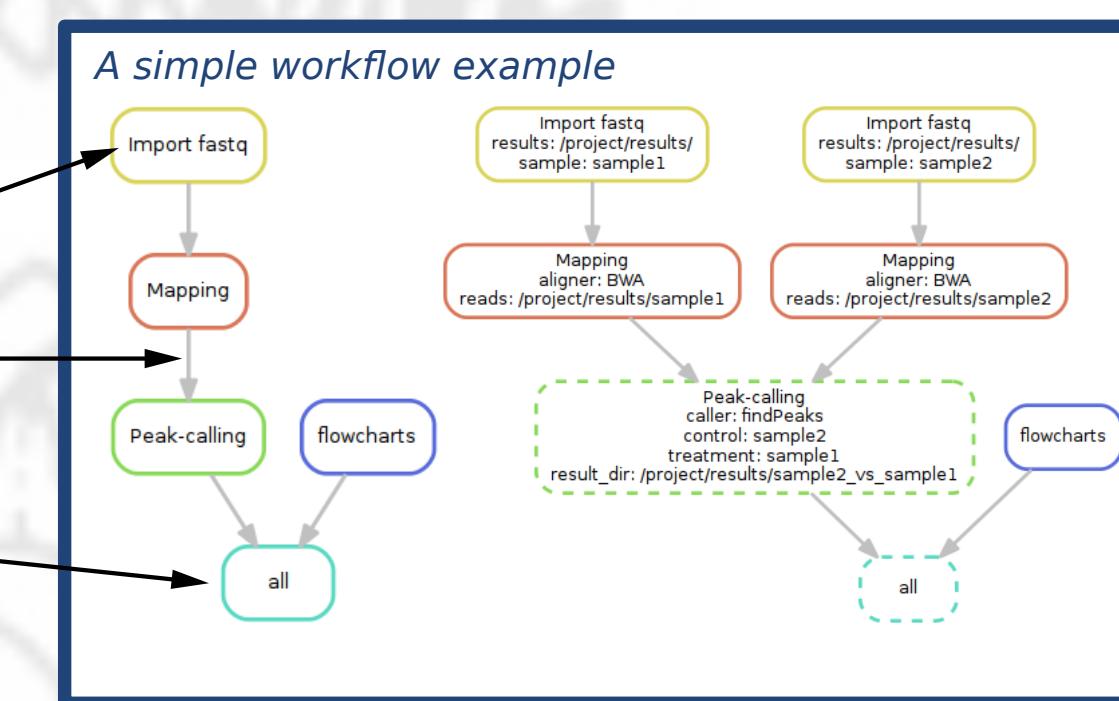
Python library, GNU make concepts

Rules are recipes to perform an operation (python, R or shell)

Dependencies between rules inferred by inputs and outputs

Final target defines operations to be performed

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



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

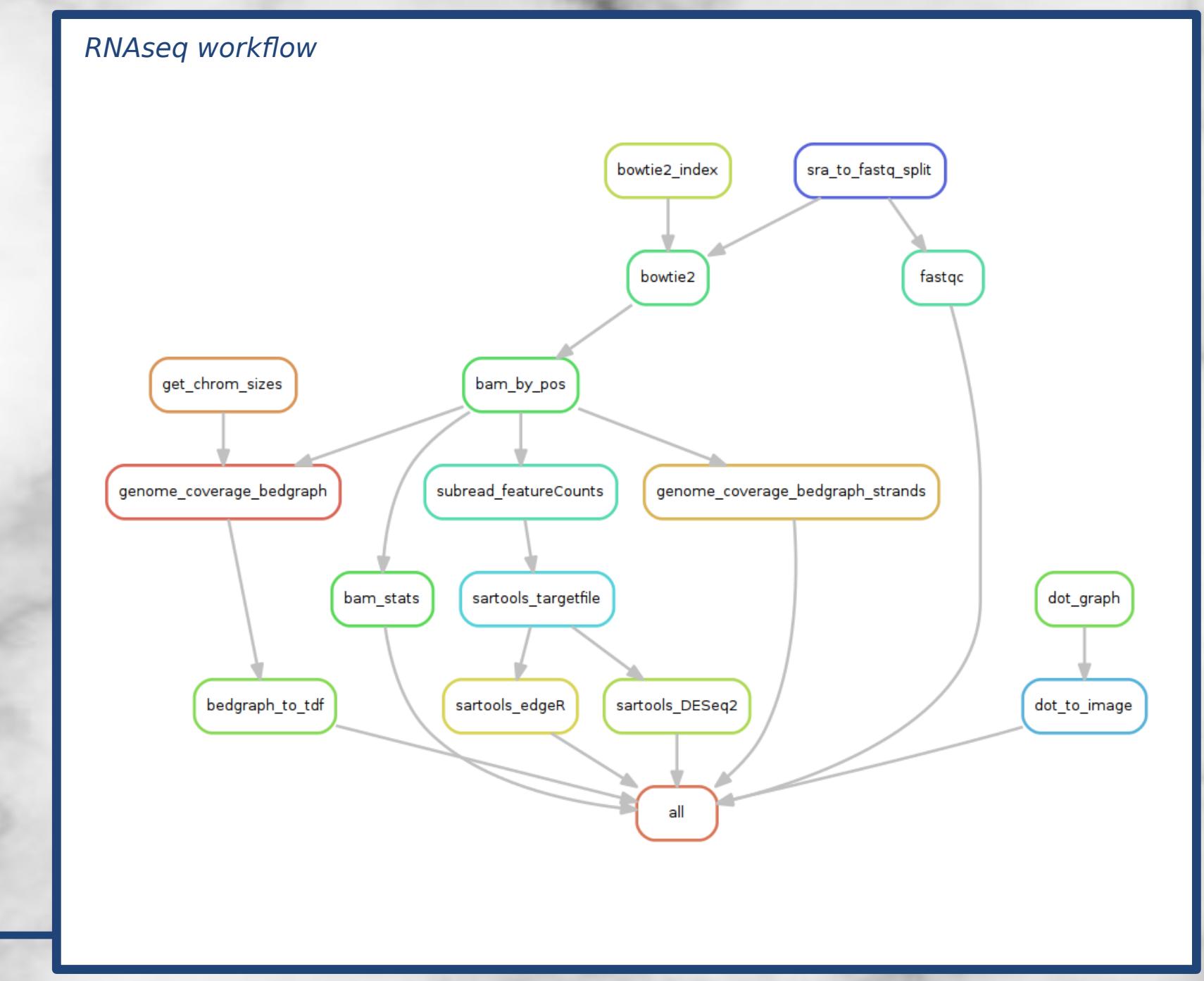
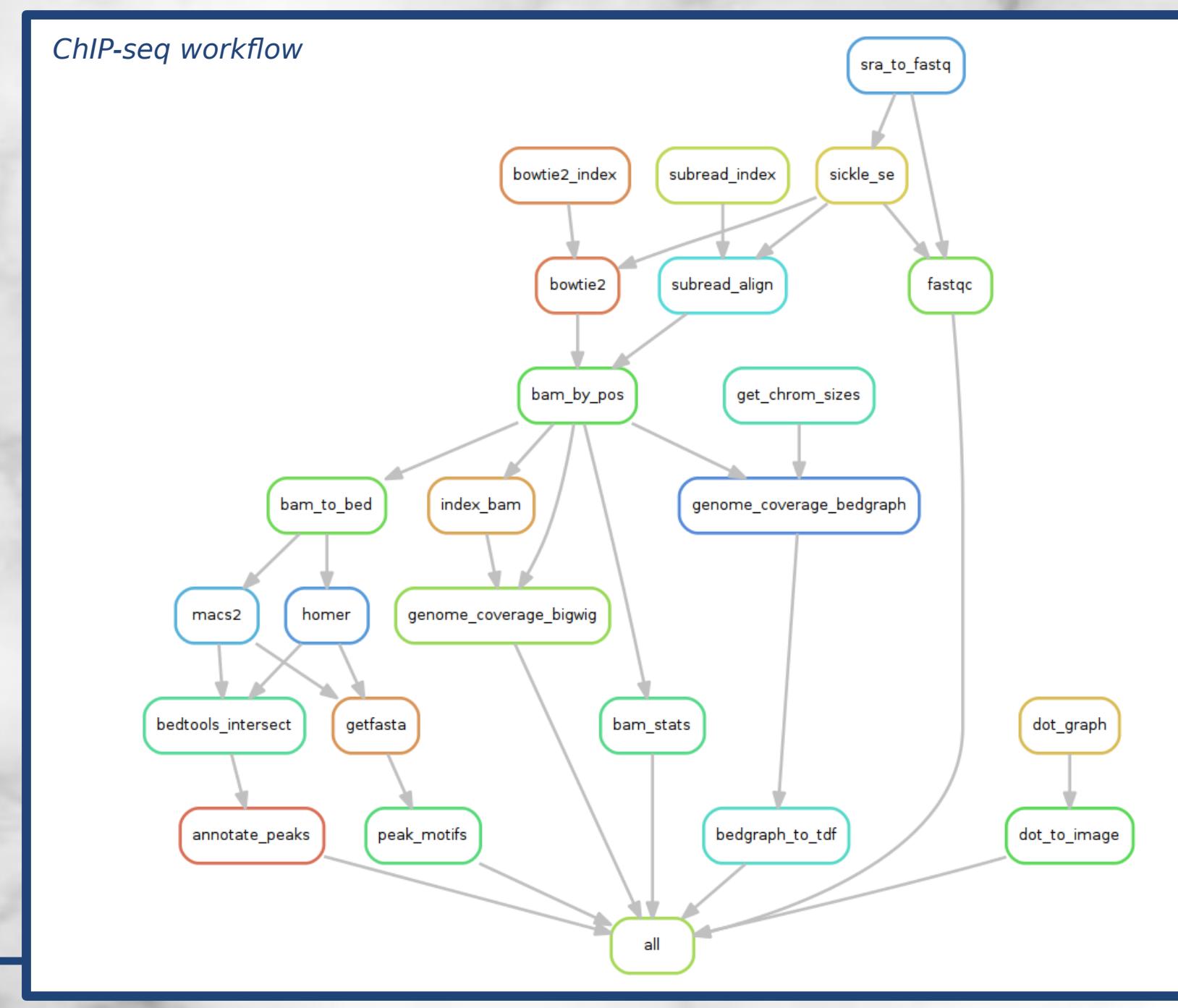
rule all:
    input:
        expand("{sample}.bam", sample = SAMPLES)
        expand("{chip}.vs_{control}", chip = CHIP, control = CONTROL)
    output:
        "{file}.bam"
    shell:
        "bowtie2 {input} > {output}"

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

Virtual environments

Virtualization allows running analyses in a **controlled environment**, ensuring **reproducibility** and **portability**.

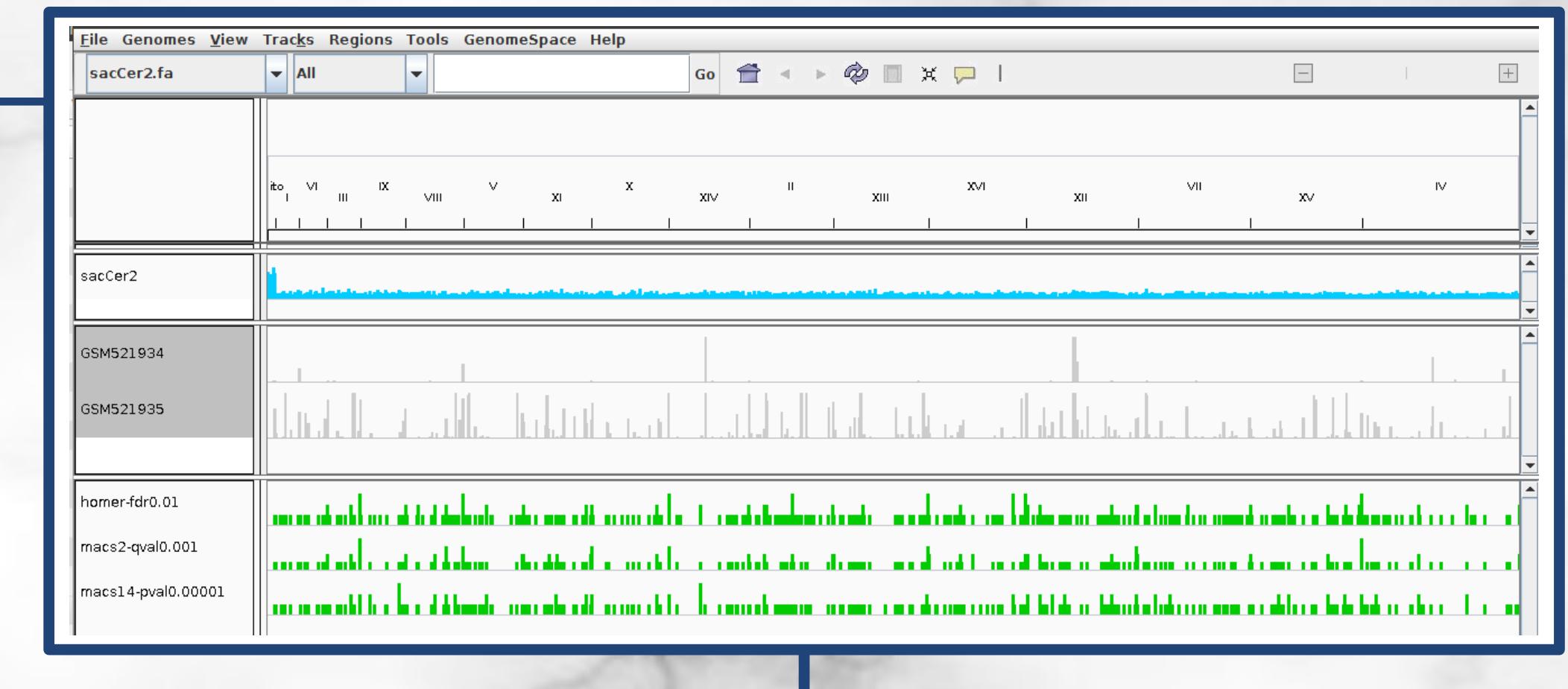
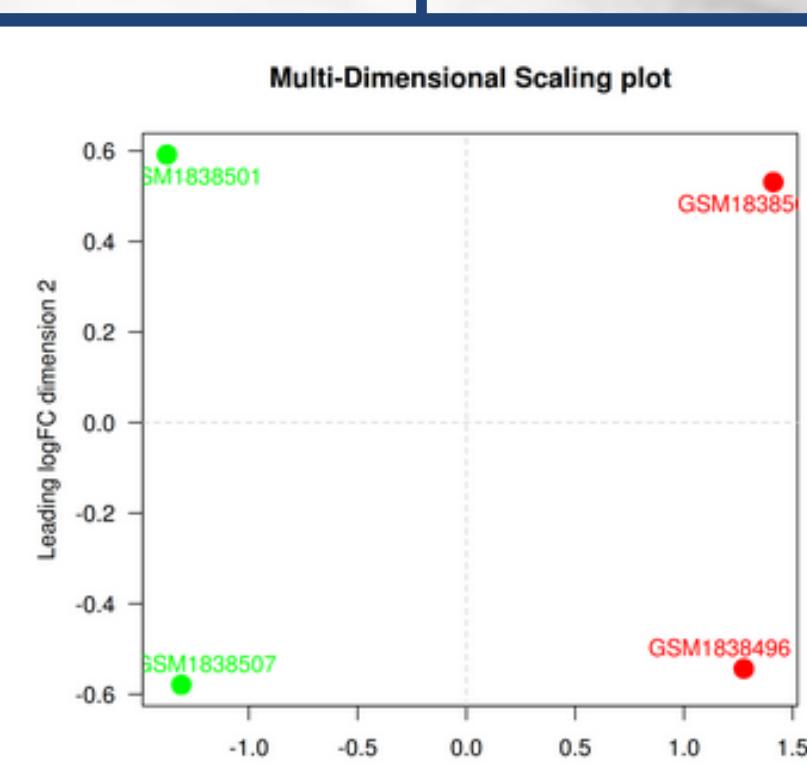
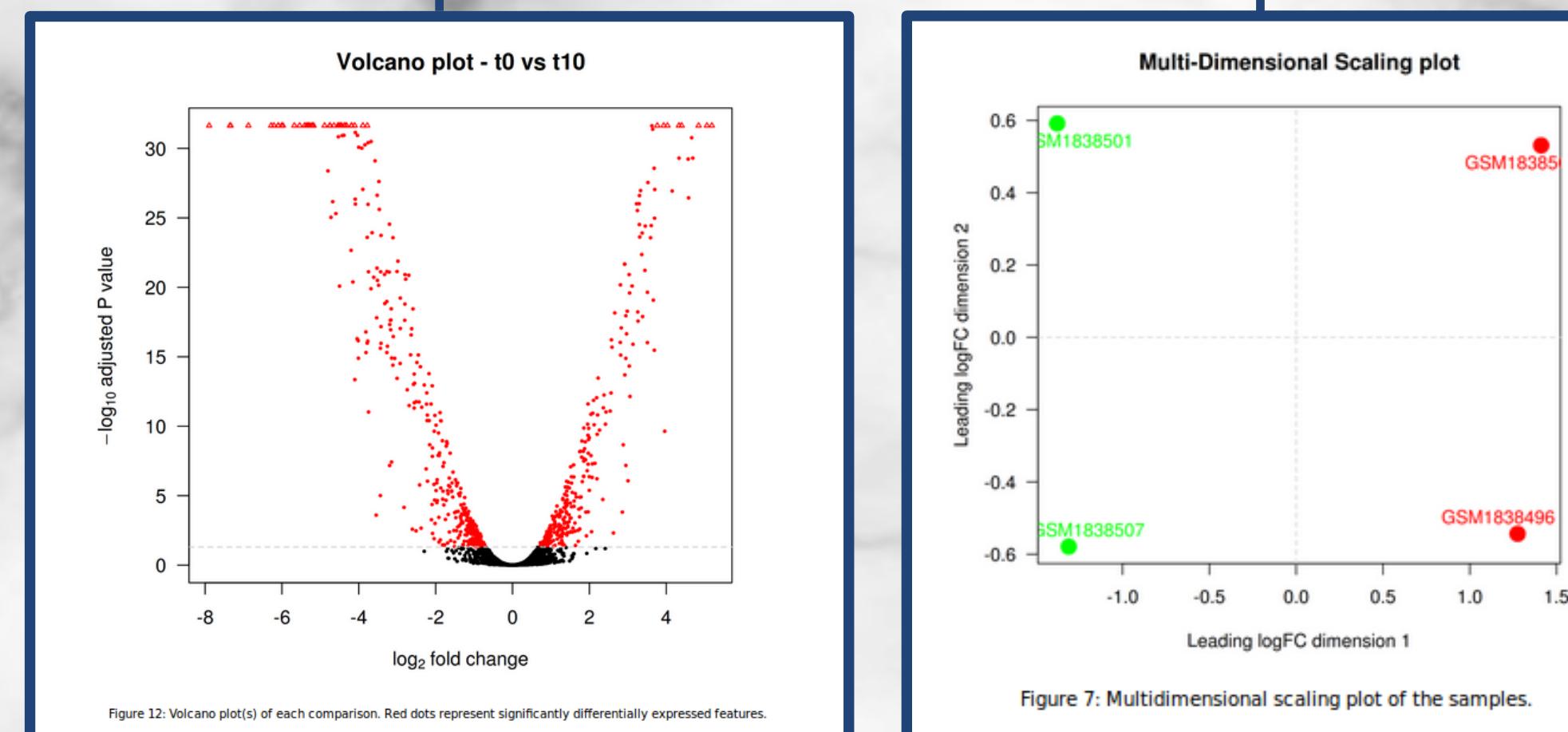
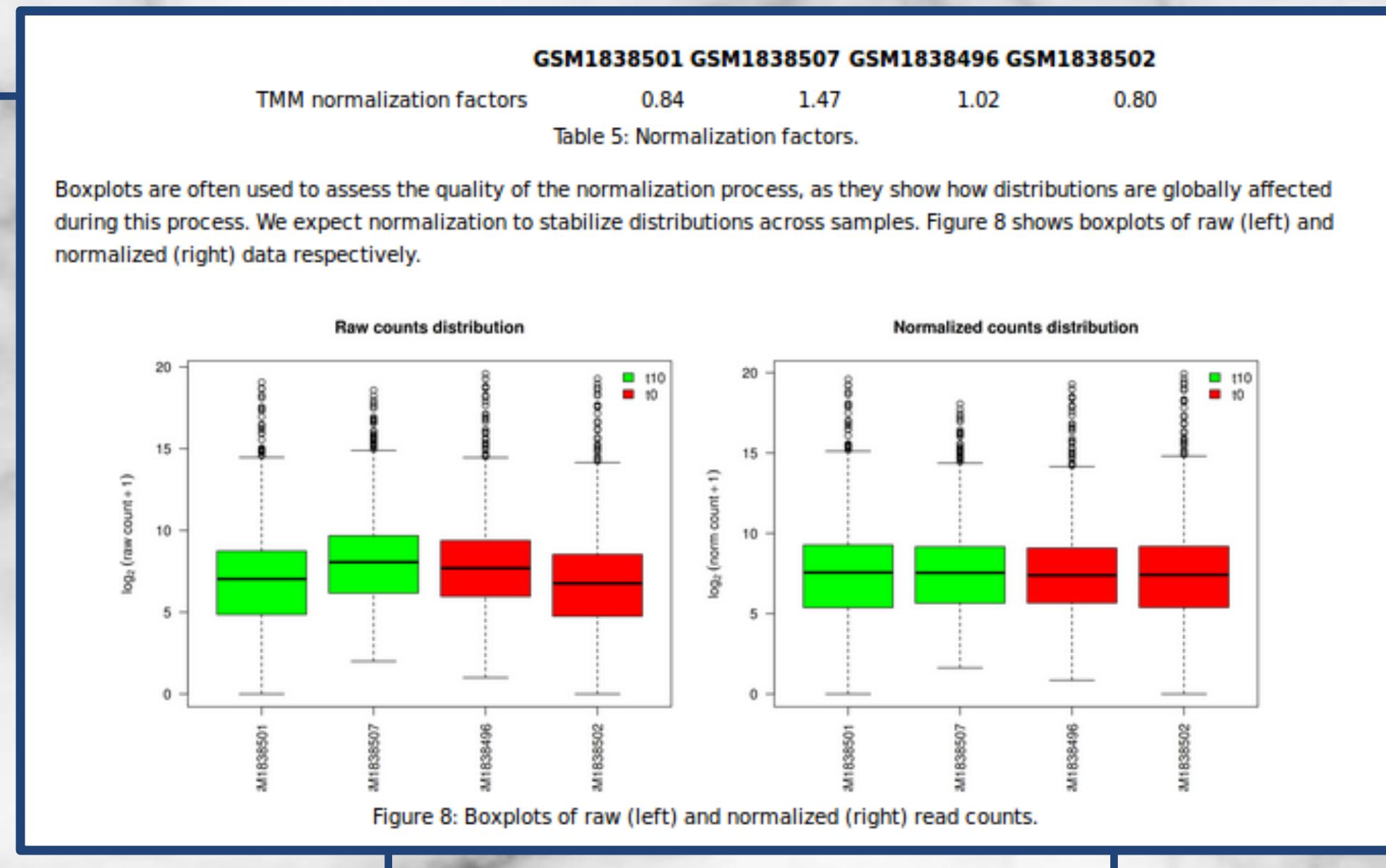
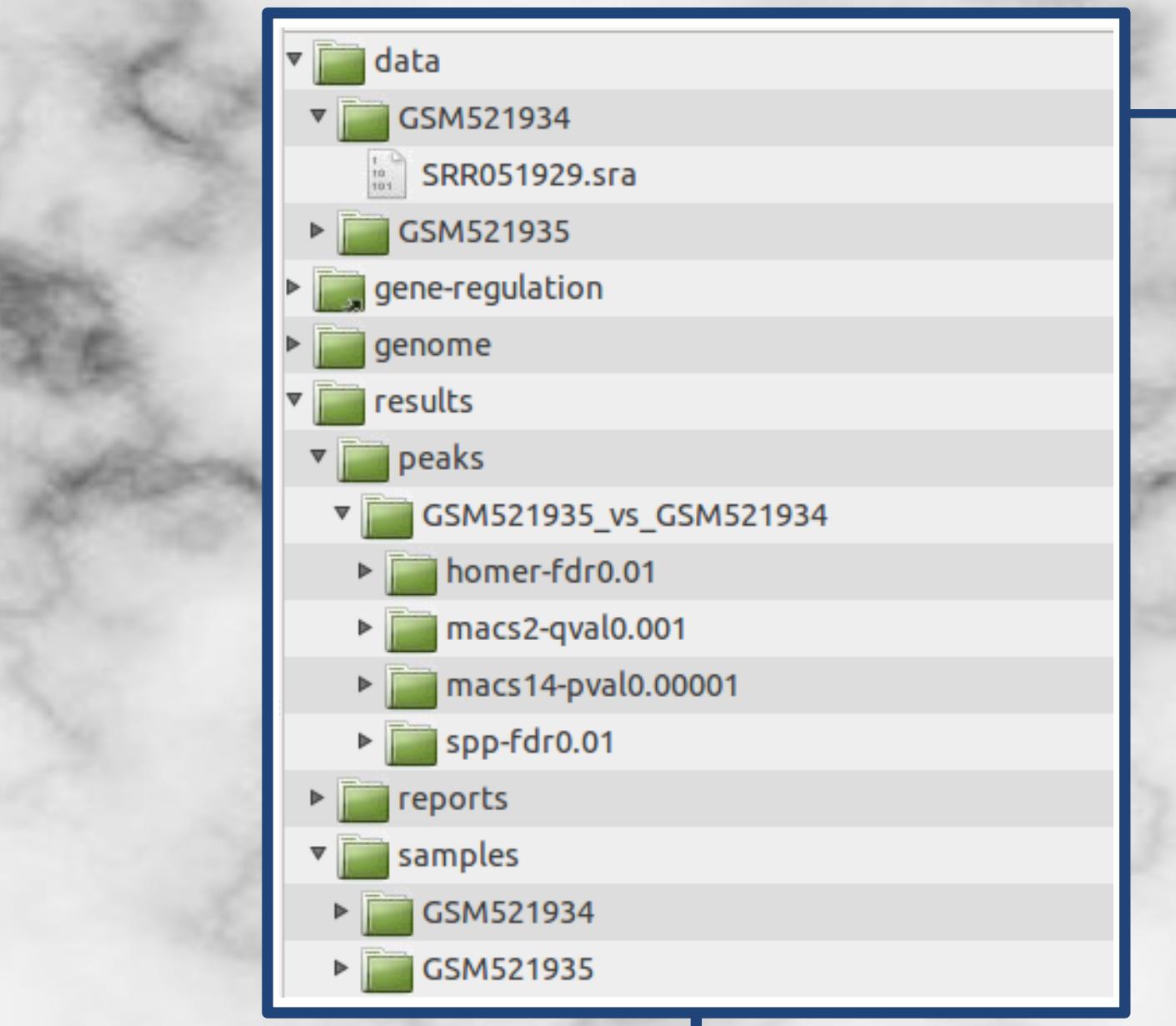
Several environments were developed, and **tutorials** written accordingly.



From a transparent data organization...

A variety of visualization tools

Intuitive file organization provides accessibility for biologists and bioinformaticians



A view of the ChIP-seq peaks detected by the workflow, under the Integrative Genome Browser (IGV)

Visualization in genomic context

Quality control

FastQC reports, SAMTools Flagstat

Discussion & Perspectives

We have developed a set of **reusable** workflows and tools for NGS analysis. For this purpose we chose to use the **Snakemake** workflow engine, which was able to fulfill our expectations:

Modularity: one can pick up any set of rules in order to build a custom pipeline.

Recycling: the rules can be used in several workflows, by several people.

Benchmarking: a number of alternative software are already available for benchmarking the peak-calling and the mapping steps.

Flexibility: config file allows tuning parameters to cater workflows to different species (bacteria, plants, mammals...) or conditions.

Collaboration: several developers can share their rules & workflows.

Portability: pipelines can be run on servers, PCs and virtual environments.

Accessibility: all the material is versioned & available in a public GitHub repository.

The number of Snakemake users is growing in France, which allowed us to create a **community** called "**CoBRAS**" (Communauté de Bioinformaticiens Rassemblés Autour de Snakemake), and to organize meetings in order to exchange methods and good practices.

Our work is available under several platforms, including GitHub and the Docker Hub. It was also used in a **protocol published** recently.

Upcoming developments include a **collaboration** with the **University of Mexico**, funded by the **NIH**, which purpose is to integrate public NGS data to enrich the RegulonDB database on *E. coli* regulation.

Protocol
Published Protocols
Volume 14(2) of the series Methods in Molecular Biology pp 297-322
Date 25 April 2016

RSAT:Plants: Motif Discovery in ChIP-Seq Peaks of Plant Genomes

Jamie A. Castro-Morales, Claire Rioualen, Bruno Contreras-Moreira, Jacques van Helden, et al.



rioualen / gene-regulation

PUBLIC | AUTOMATED BUILD
rioualen/gene-regulation

Acknowledgments: This work is funded by France Génomique, NIH grant GM0110597 & FOINS-CONACYT - Fronteras de la Ciencia 2015 - ID 15. I would like to thank Lucie Khamvongsa, Jocelyn Brayet, Bruno Contreras-Moreira, Julio Collado-Vides & all his team for their great contributions to this work.