

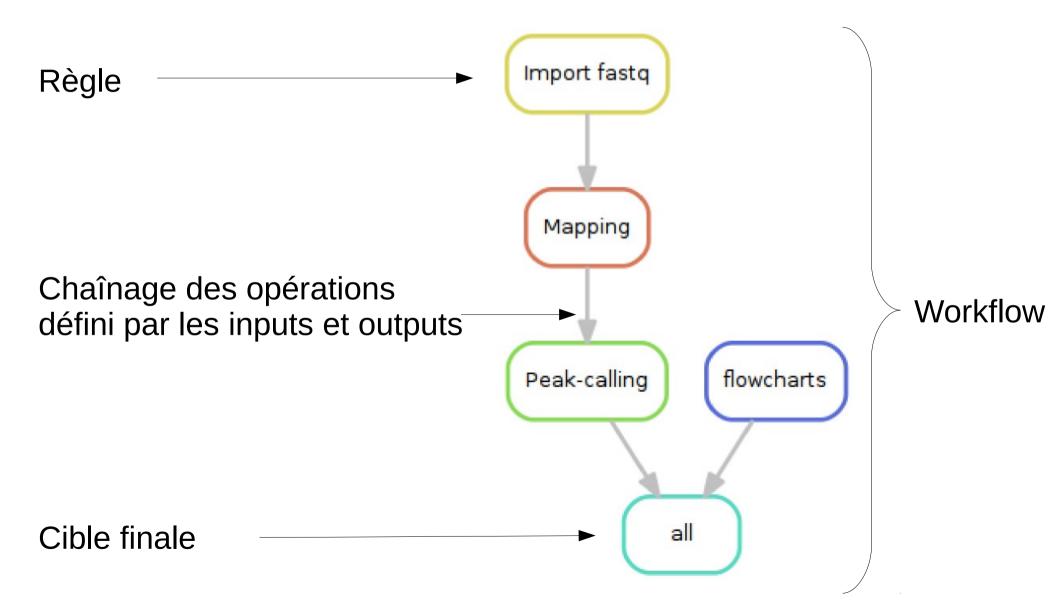




# Développement de workflows pour le NGS Exécution dans un environnement virtuel Claire Rioualen

# Principes de Snakemake

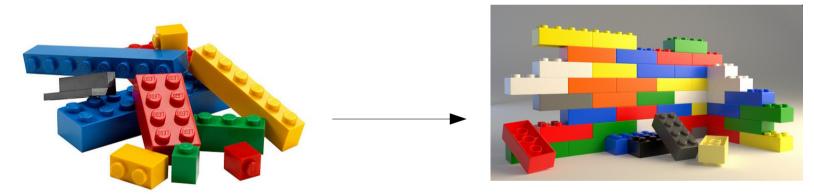




# Gene-regulation : développement de pipelines flexibles et reproductibles pour le NGS

- Pool de règles recyclables et interchangeables
  - Mapping (bowtie, subread, tophat...)
  - Peak-calling (macs2, homer, swembl...)

- ...



- Workflows recyclables et customisables
  - ChIP-seq
  - RNA-seq

- Organisation des métadonnées par projet ou dataset :
  - 1 fichier "workflow.py"

```
1## Includes
 2 RULES = os.path.join("gene-regulation/scripts/snakefiles/rules")
 4 include: os.path.join(RULES, "bowtie2_index.rules")
 5 include: os.path.join(RULES, "bowtie2.rules")
 6 include: os.path.join(RULES, "fastqc.rules")
 7 include: os.path.join(RULES, "homer.rules")
 8 include: os.path.join(RULES, "macs2.rules")
 9 include: os.path.join(RULES, "spp.rules")
11 ## Samples
12 SAMPLE IDS = read table(config["metadata"]["samples"])['ID']
14 ## Experimental design
15 DESIGN = read table(config["metadata"]["design"])
16 TREATMENT = DESIGN['treatment']
17 CONTROL = DESIGN['control']
24 IMPORT = expand(FASTQ_DIR + "/{samples}/{samples}, fastq", samples=SAMPLE_IDS)
26 ## Quality report
27 QC = expand(FASTQ DIR + "/{samples}/{samples} fastqc/{samples} fastqc.html", samples=SAMPLE IDS)
29 ## Mapping
30 ALIGNER ---- = config["tools"]["mapping"].split()
31 INDEX ----- = expand(GENOME DIR + "/{aligner}/" + config["genome"]["fasta file"], aligner=ALIGNER)
32 MAPPING .... = expand(SAMPLE DIR + "/{samples}/{samples} {aligner}.bam", samples=SAMPLE IDS, aligner=ALIGNER)
34 ## Peak-calling
35 PEAKCALLER = config["tools"]["peakcalling"].split()
36 PEAKS .... = expand(expand(PEAKS DIR + "/{treat} vs {control}/{{peakcaller}}.bed", zip, treat=TREATMENT, control=CONTROL), peakcaller=PEAKCALLER)
38 rule all:
39 ....input: \
   ..... IMPORT, \
42 · · · · · · INDEX, · \
              MAPPING. \
```

- Organisation des métadonnées par projet ou dataset :
  - 1 fichier "workflow.py"
  - 1 fichier "config.yml"

```
2 - author: "Claire Rioualen"
 3 - qsub: - " - V - m - a - - d - . "
    genome:
   ····organism: "Saccharomyces cerevisiae"
   --- fasta file: "sacCer2.fa"
7 ... gff3 Tile: "sacCer2.gff3"
    ... gtf file: "sacCer2.gtf"
10 metadata:
11 · · · · samples: "metadata/samples.tab"
12 ... design: "metadata/design.tab"
13 .... configfile: "metadata/config.yml"
14 · · · seq type: "se"
16 dir:
17 · · · · fastq: "fastq"
18 ... genome: "genome"
19 · · · results: "results"
21 tools:
22 .... trimming: "sickle"
23 ... mapping: "bowtie2"
24 --- peakcalling: "macs2 homer macs14 spp swembl"
26 sickle:
27 .... threshold: "20"
28
29 bowtie2:
30 .... threads: "5"
31 ... max_mismatches: "1"
32
33 macs14:
34 ... mfold: "5,30"
35 ... keep dup: "auto"
   bandwidth: "300"
```

- Organisation des métadonnées par projet ou dataset :
  - 1 fichier "workflow.py"
  - 1 fichier "config.yml"
  - 1 fichier "samples.tab"

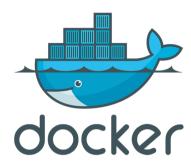
```
1; Sample description file
2; Organism: Saccharomyces cerevisiae
3; Experiment type: ChIP-seq
4;
5 ID→Condition→ ENA_experiment→Scan Name
6 GSM521934→ input→ SRX021359→ SRR051930
7 GSM521935→ chip → SRX021358→ SRR051929
8
9
```

- Organisation des métadonnées par projet ou dataset :
  - 1 fichier "workflow.py"
  - 1 fichier "config.yml"
  - 1 fichier "samples.tab"
  - 1 fichier "design.tab"

```
1; Design for the peak-calling of ChIP-seq analysis
2;
3 treatment → control
4 GSM521935 → GSM521934
```

#### Développement d'environnements virtuels

- Contrôler l'environnement d'exécution d'un workflow
- Fournir une solution "clés en main" avec tous les programmes et dépendances inclus
  - Système d'exploitation (Ubuntu 14.04)
  - Librairies et packages divers
  - Outils d'analyse NGS
  - Dépôt git gene-regulation
- Tutoriels
  - Gene-regulation tutorials













## Exécution dans un environnement virtuel



- Création d'une appliance Gene-regulation sur le cloud IFB
- Création d'une instance de Gene-regulation
- Connexion d'un disque virtuel :

```
root@vm0103:~/mydisk/GSE20870-analysis# ll -R
total 16
drwxr-xr-x 4 root root 4096 Dec 1 16:01 ./
drwxr-xr-x 3 root root 4096 Jun 2 13:32 ../
drwxr-xr-x 4 root root 4096 Jun 2 12:52 data/
lrwxrwxrwx 1 root root 21 Dec 1 15:59 gene-regulation -> /root/gene-regulation/
drwxr-xr-x 2 root root 4096 Dec 1 16:02 genome/
./data:
total 16
drwxr-xr-x 4 root root 4096 Jun 2 12:52 ./
drwxr-xr-x 2 root root 4096 Jun 2 12:52 GSM521935/
./data/GSM521934:
total 81192
drwxr-xr-x 2 root root
                        4096 Jun 2 12:52 ./
drwxr-xr-x 4 root root
                        4096 Jun 2 12:52 .../
-rw-r--r-- 1 root root 83043426 Jun 2 12:52 SRR051929.sra
./data/GSM521935:
total 94124
drwxr-xr-x 2 root root
                        4096 Jun 2 12:52 ./
                        4096 Jun 2 12:52 ../
drwxr-xr-x 4 root root
-rw-r--r-- 1 root root 96275406 Jun 2 12:52 SRR051930.sra
total 29652
drwxr-xr-x 2 root root
                        4096 Dec 1 16:02 ./
                        4096 Dec 1 16:01 ../
-rw-r--r-- 1 root root 12360636 Jun 2 13:33 Saccharomyces cerevisiae.R64-1-1.30.dna.genome.fa
rw-r--r-- 1 root root 6673011 Jun 2 13:38 Saccharomyces cerevisiae.R64-1-1.30.qff3
-rw-r--r-- 1 root root 11268860 Jun 2 13:38 Saccharomyces cerevisiae.R64-1-1.30.gtf
root@vm0103:~/mydisk/GSE20870-analysis#
```

# Exécution dans un environnement virtuel

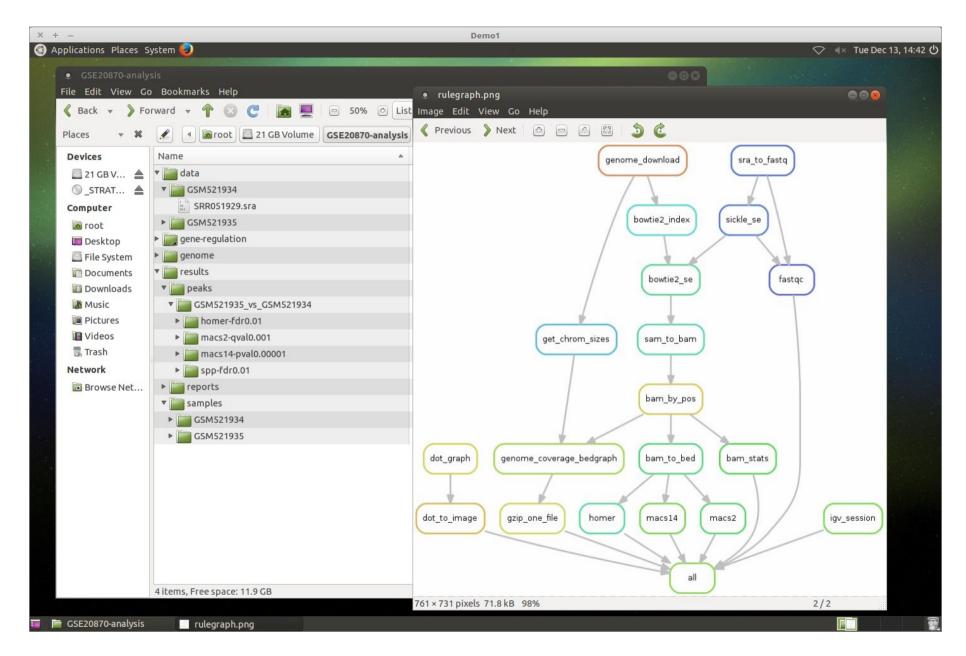


- Création d'une appliance sur le cloud IFB
- Création d'une instance de Gene-regulation
- Connexion d'un disque virtuel
- Exécution du workflow:

# Aperçu des résultats (bureau distant)



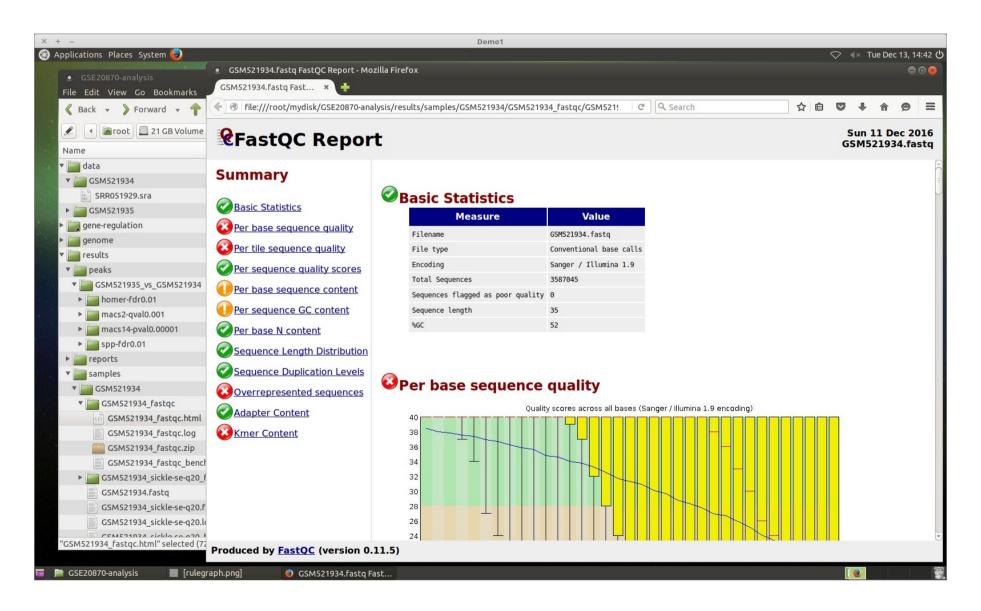








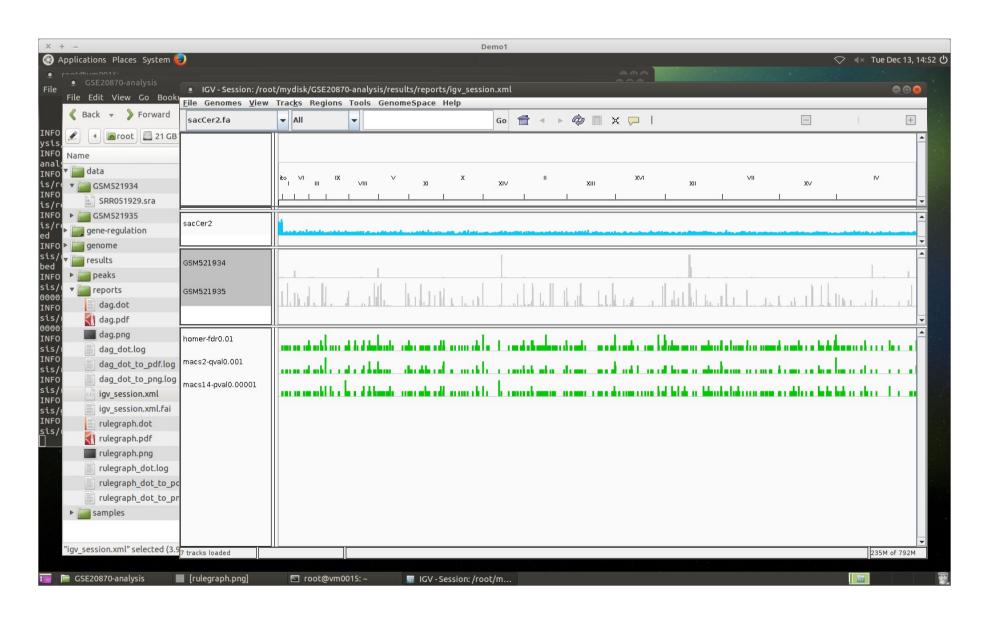




# Aperçu des résultats (bureau distant)



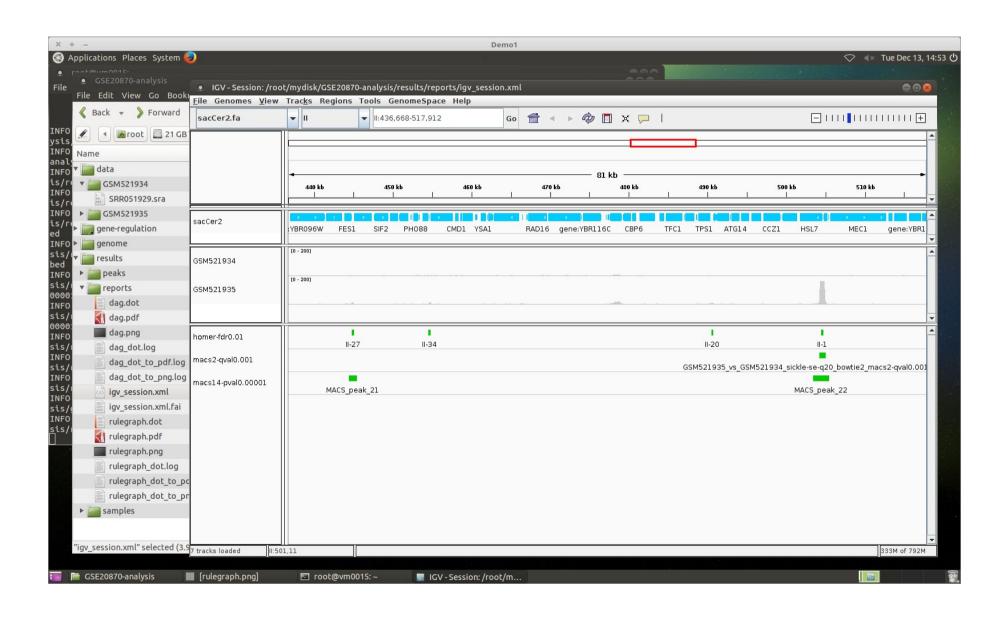




# Aperçu des résultats (bureau distant)











- Jacques van Helden
- Denis Puthier
- Lucie Khamvongsa-Charbonnier
- Jaime Castro-Mondragon



- Julio Collado-Vides
- Alberto Santos-Zavaleta
- Mishael Sánchez-Pérez







Jocelyn Brayet



