

# SNAKEMAKE AT ICM

---

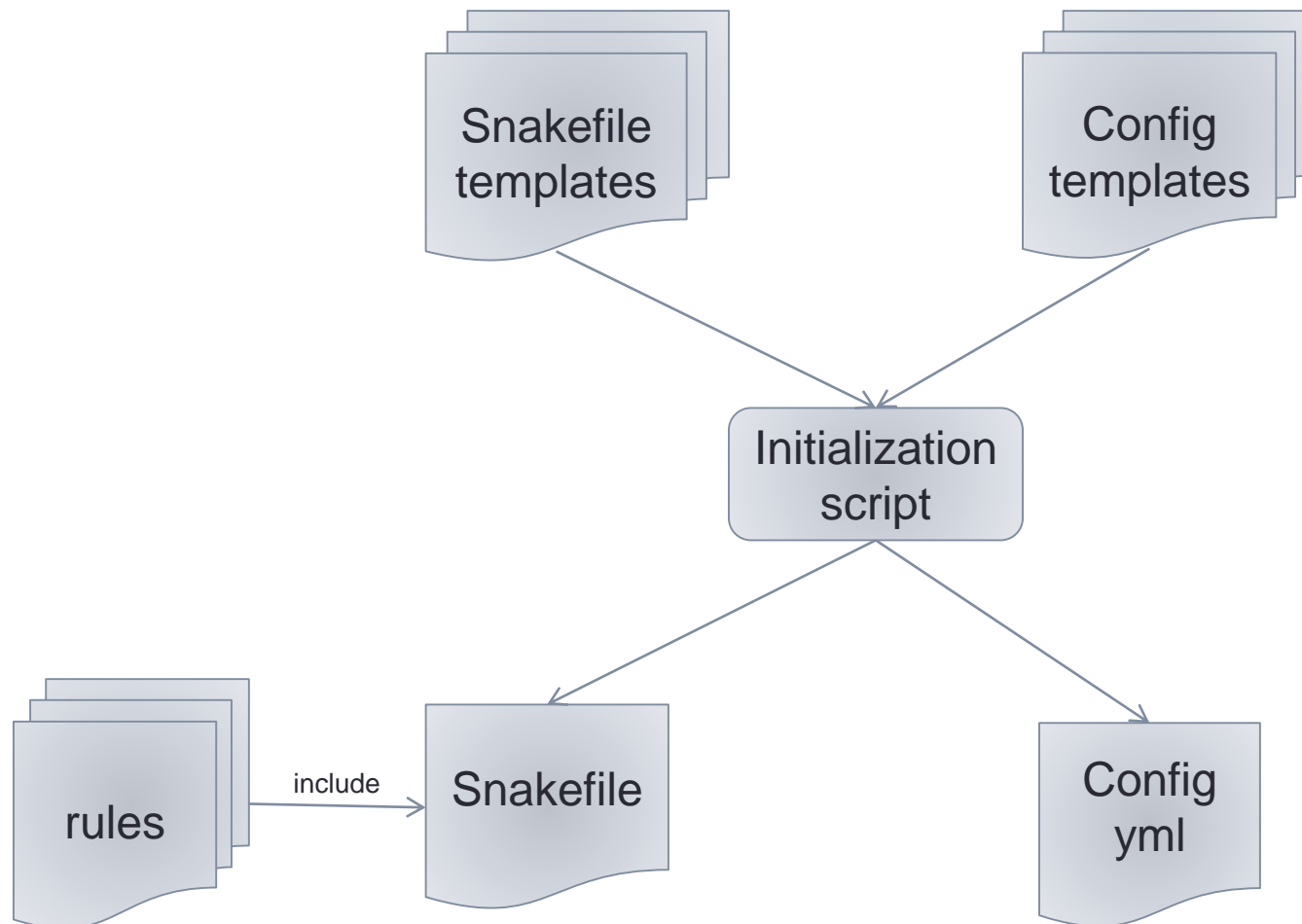
Vivien Deshaies  
Justine Guegan



# Pipelines

- WES
- Target sequencing
- RNA-seq
- WGS

# Organization



# Config file (1/3)

samples:

sample1:

run1:

fastq\_read1: fastq/run1/sample1\_1.fastq.gz

fastq\_read2: fastq/run1/sample1\_2.fastq.gz

run2:

fastq\_read1: fastq/run2/sample2\_1.fastq.gz

fastq\_read2: fastq/run2/sample2\_2.fastq.gz

include\_path: ~/git/WGS\_pipeline/include/

scripts\_path: ~/git/WGS\_pipeline/scripts/

# Config file (2/3)

gatk:

jar: /tools/GenomeAnalysisTK.jar

resources:

low:

xmx: 6g

mem: 7

high:

xmx: 19g

mem: 20

# Config file (3/3)

genome:

fasta:

/genomes/human\_g1k\_v37\_decoy.fasta

chromosomes:

[1, 2, 3, 4, 5, 6, 7, ...]

- WGS analysis → pipeline parallelization by chromosome

rule merge\_clean\_bam:

input:

```
bams = expand(
    "tmp_split/{{sample}}_{chrom}_sorted_markup_recal.bam",
    chrom=config["genome"]["chromosomes"]
)
```

# Rules

## rule example:

input: "merged\_reads/{sample,[A-Za-z0-9\-\-]+}\_merged.bam"

output: "output"

log:

out = "logs/example.log",

err = "logs/example.err"

benchmark: "benchmarks/example.txt"

params:

gatk\_jar = config["gatk"]["jar"],

xmx = config["gatk"]["resources"]["high"]["xmx"],

genome = config["genome"]["fasta"]

threads: 1

resources:

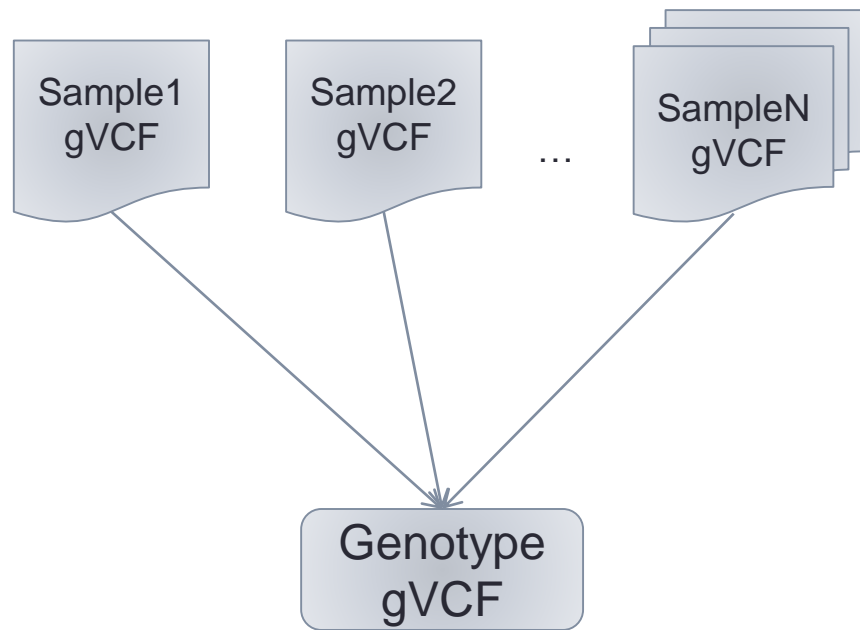
mem = config["gatk"]["resources"]["high"]["mem"]

shell:

"commande"

" > {log.out} 2> {log.err}"

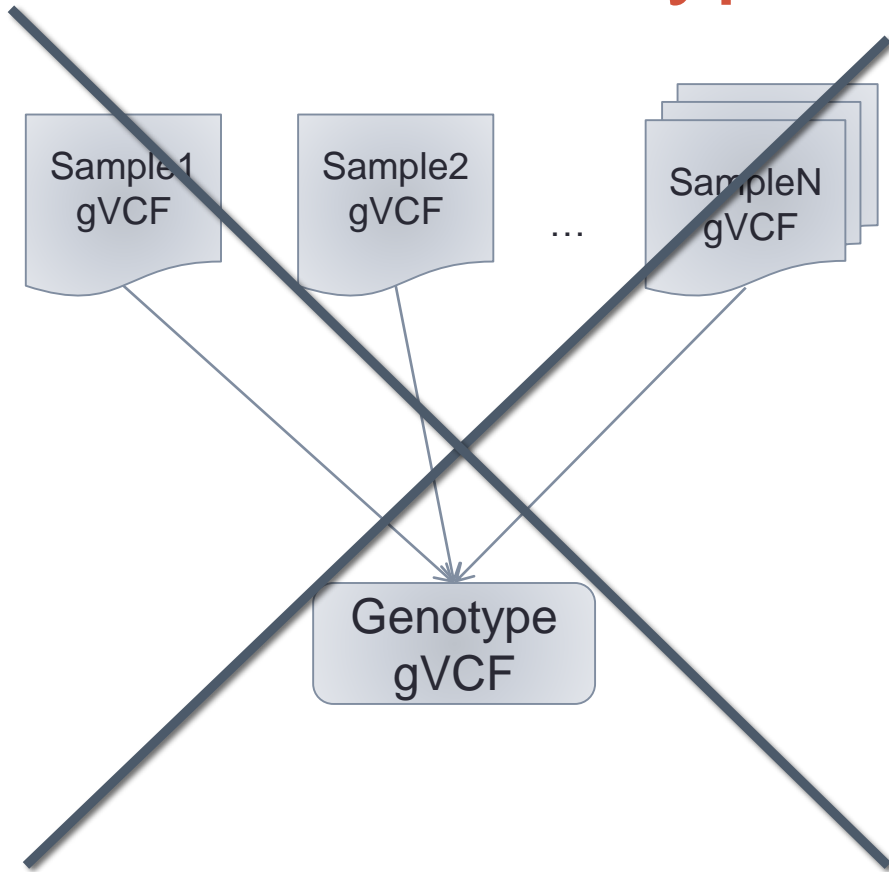
# GATK GenotypeGVCF



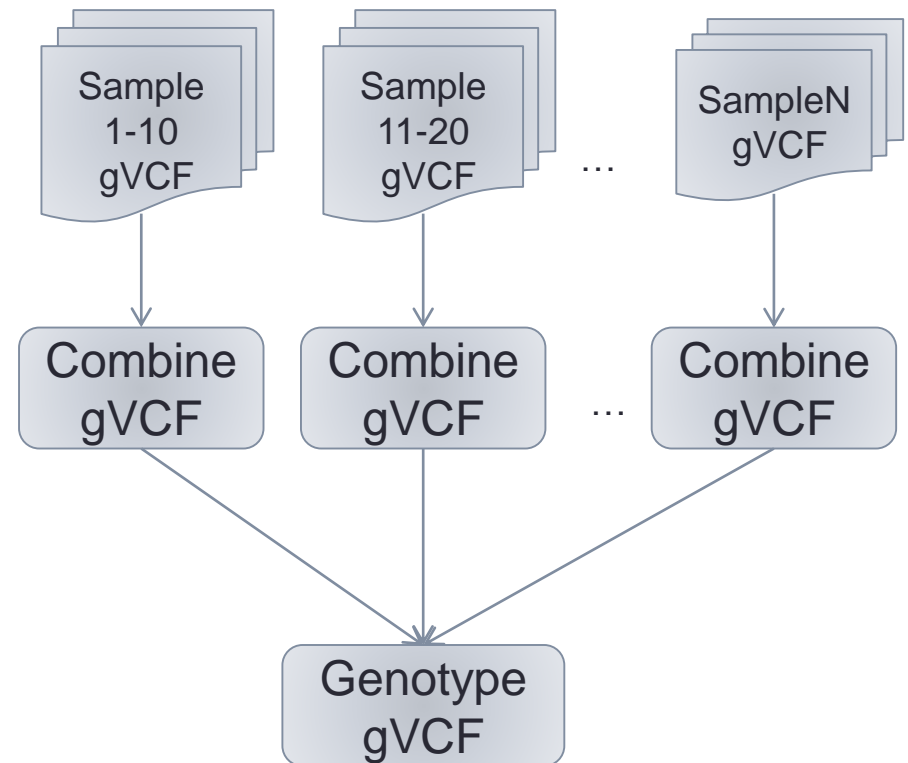
- With 70 samples ~ 100 GB



# GATK GenotypeGVCF



- With 70 samples ~ 100 GB



- With 70 samples ~ 20 GB

# Rule : merge\_gvcf

```
sampleGroupList = get_sample_group_list(  
    list(config["samples"].keys()),  
    config["number_of_sample_by_gvcf"]  
)
```

rule merge\_gvcf:

```
    input:  
        lambda wildcards: list(  
            ["variants/{sample}.g.vcf.gz".format(sample = sample) for sample in  
sampleGroupList[int(wildcards.n)]]  
        )  
    output: "variants/sample_group.{n}.g.vcf.gz"
```

rule genotype\_call:

```
    input:  
        expand("variants/sample_group.{n}.g.vcf.gz", n=range(0, len(sampleGroupList)))
```

# Snakemake launching



- Local :

```
snakemake -p --cores=20 --resources mem=60
```

- Slurm :

```
snakemake -j200 -p  
--cluster "sbatch --mem={resources.mem}000 -c {threads}"
```

# Rulegraph



# Perspectives of evolution



- Scripts conservation
  - Bash scripts created by snakemake are currently automatically deleted
- Log slurm / snakemake
  - Cluster errors are in : slurm-`{jobid}`.out
  - Command errors are in snakemake output

# Thanks to

- Justine Guegan
- Ivan Moszer
- Romain Daveau
- Vincent Perlberg
- Ludovic Prevost