

Snakemake : What? Why?

- Snakemake is lightweight
- Snakemake is flexible
- Snakemake use both classical *Bash* and *Python* syntax
- Your analysis process don't have to be distorted
- Snakemake has a vast community
- Looks like *Make* but it's not : make the most of *Make* and *Python*
- Snakemake is actually usefull beyond reproducibility
- Usable on cluster with minimum effort



Snakemake : Philosophy

- Looks like a Python script
- The workflow is a set of entangled rules (inputs/outputs/shell)



The dependency between rules is based on I/O

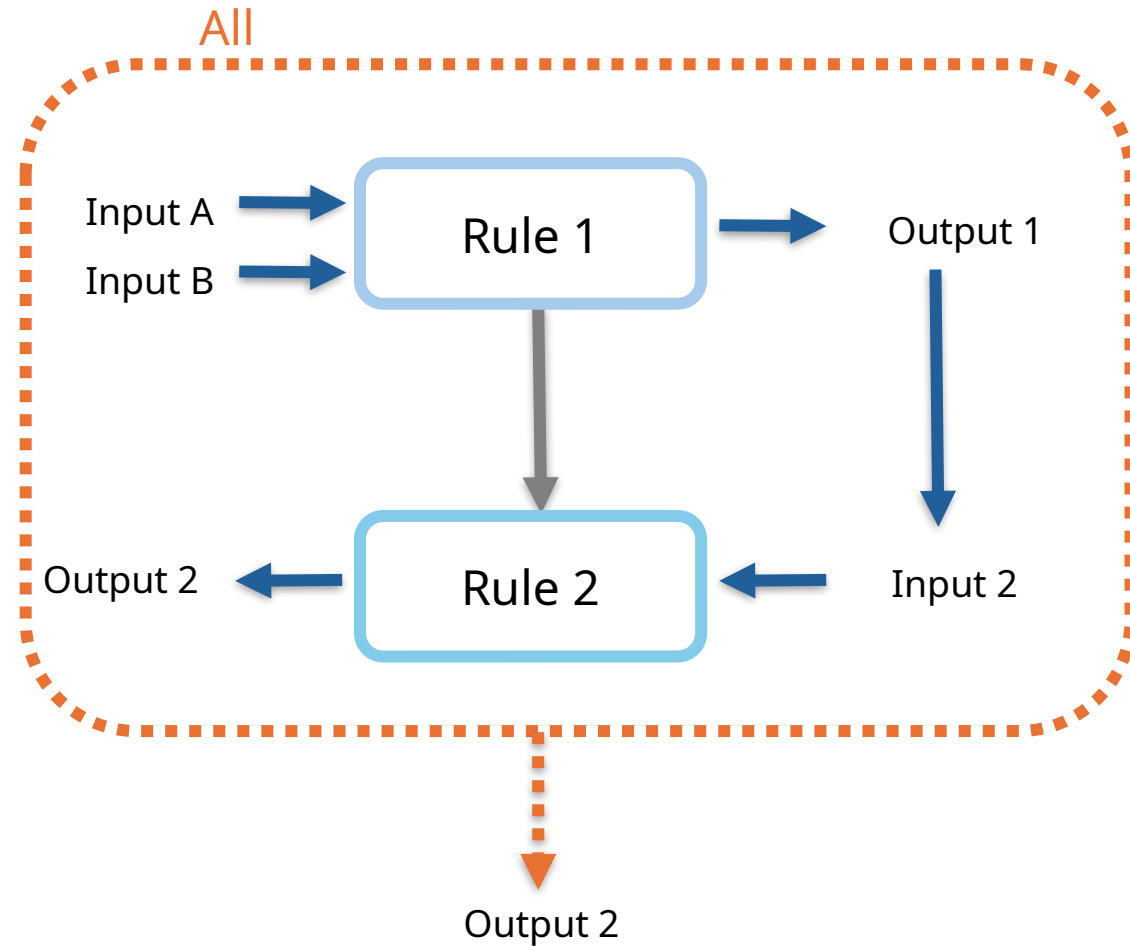
- Computed by Snakemake based on
- Snakefile

```
1 rule rule_1:
2     input:
3         input_1 = 'file_1.txt'
4     output:
5         output_1 = 'file_1.out'
6     shell:
7         """
8         So some shell command here
9         I.e. echo "Hello World" > file_1.out
10    """
```

- Snakemake rules have to be defined in *Snakefile*
- Launch pipeline : `snakemake -c 1`
- Launch pipeline in dry mode : `snakemake -n`
- Plot workflow steps : `snakemake --dag | dot -T pdf > dag.pdf`
- Re-run all jobs the output of which is recognized as incomplete : `snakemake -c 1 --&erun-incomplete`

Snakemake : My first workflow

```
1 rule rule_1:
2     input:
3         input_A = 'file_a.txt'
4         input_B = 'file_b.txt'
5
6     output:
7         output_1 = 'file_1.out'
8     shell:
9         """
10        So some shell command here
11        I.e. cat {input.input_A} {input.input_B} > file_1.out
12        """
13 rule rule_2:
14     input:
15         input_2 = 'file_1.out'
16
17     output:
18         output_2 = 'words_total.out'
19     shell:
20         """
21        So some shell command here
22        I.e. wc -w {input.input_2} > words_total.out
23        """
24 rule all:
25     input:
26         'words_total.out'
```

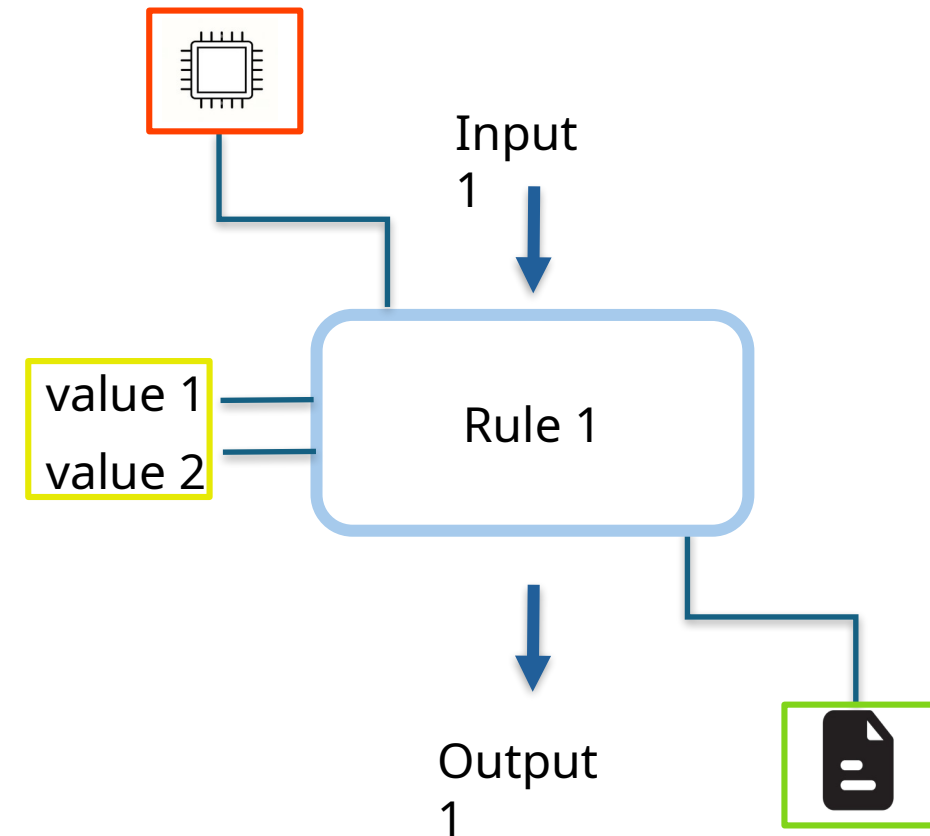


Launch pipeline : `snakemake --cores 1`

Snakemake : Pimp my rule

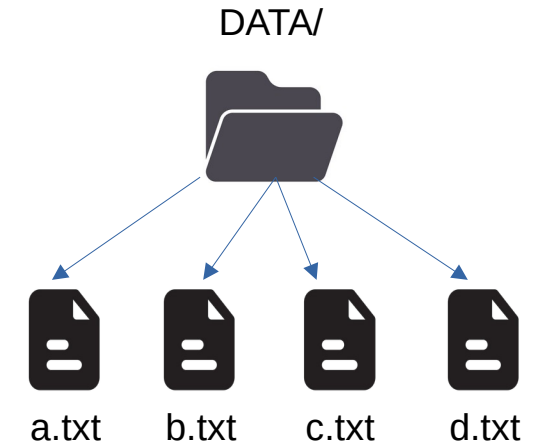
```
rule generic_rule:
    input:
        input_file = 'file_a.txt'
    output:
        output_file = 'file_a.out'
    threads: 2
    resources:
        mem_mb=100
    params:
        param1 = 'value1'
        param2 = 'value2'
    conda : 'envs/myenv.yaml'
    message: 'This is a generic rule using {params.param1} and {params.param2}'
    log: 'logs/{input.inputfile}.log'

    shell:
        """
        So some shell command here
        I.e. command --arg1 {params.param1} --args2 {params.param2} --threads {threads} \
        {input.input_file} > {output.output_file}
        """
```



Snakemake : Make my workflow generic

```
1 from os import join
2 DATA_FOLDER = "DATA/"
3 RESULTS_FOLDER = "RESULTS/"
4 SAMPLES = ["A", "B", "C", "D"]
5
6 rule rule_1:
7     input:
8         my_input = expand(join(DATA_FOLDER, "{sample}.txt"), sample = SAMPLES)
9     output:
10        my_output = expand(join(REULTS_FOLDER, "{sample}.cut.txt"), sample = SAMPLES)
11    shell:
12        """
13        cut -f 1 {input.my_input} > {output.my_output}
14        """
15
16 rule rule_2:
17     input:
18         my_input_bis = expand(join(REULTS_FOLDER, "{sample}.cut.txt"), sample = SAMPLES)
19     output:
20         my_last_output = expand(join(REULTS_FOLDER, "{sample}.counts"), sample = SAMPLES)
21     shell:
22         """
23         wc -l {input.my_input_bis} > {output.my_output_bis}
24         """
25
26 rule all:
27     input:
28         expand(join(DATA_FOLDER, "{sample}.counts"), sample = SAMPLES)
```



a.cut.txt
b.cut.txt
c.cut.txt
d.cut.txt

Rule 1

a.counts
b.counts
c.counts
d.counts

Rule 2

Use wildcards to make rules adaptable to various file names