

What's a Workflow? From Chaos to Automation

Rousseau Baptiste









What's a Workflow?























output data























When planning an analysis, we need to consider and keep track of:

- The format of the data at all stages
- Any pre-processing steps
- Tools/software used, including which release/version
- Parameters we select when using tools





Why not a tracking table?

- Error-prone
- Non-reproducible
- No automation
- Not collaborative







Why not a script?

- Hardcoded paths
- No error handling
- Not reproducible
- Not scalable

rna_seq_workflow.sh

#!/bin/bash

Set input and output files
RAW_FASTQ="sample.fastq.gz"
TRIMMED_FASTQ="sample_trimmed.fastq.gz"
QC_DIR="qc_reports"
REFERENCE_GENOME="genome_index"
SAM_FILE="aligned.sam"

1. Run FastQC
echo "Running FastQC..."
mkdir -p \$QC_DIR
fastqc -o \$QC_DIR \$RAW_FASTQ

2. Trim reads echo "Trimming reads..." trimmomatic SE -phred33 \ \$RAW_FASTQ \$TRIMMED_FASTQ \ LEADING:3 TRAILING:3 SLIDINGWINDOW:4:15 MINLEN:36

3. Align reads
echo "Aligning reads with HISAT2..."
hisat2 -x \$REFERENCE_GENOME -U \$TRIMMED_FASTQ -S \$SAM_FILE

echo "Workflow completed!"







Why a workflow

- Automation
- Run time management
- Software management
- Portability & Interoperability
- Reproducibility





Why Workflows matter?



Why a workflow

- Automation
- Run time management
- Software management
- Portability & Interoperability
- Reproducibility

a Analysis workflow	b Traditional pipeline
Transcript expression quantification Fastq Reference Grch38 sequence Ensembl 91	Requirements Platform-specific fastQCI Salmon1 Pipeline code
Step 1: quality control fastQC v.0.11.9 Step 2: index creation Salmon v.1.3.0	Execution Re-entrance checkpoints Input data Step 1
Step 3: quantification Salmon -1 A	Step 2
QC report Transcript expression	a Software, versions, parameters

https://doi.org/10.1038/s41592-021-01254-9





Why Workflows matter?

Why a workflow **manager**

- Automation
- Run time management
- Software management
- Portability & Interoperability
- Reproducibility
- Scalability (local, cluster, cloud)
- Resuming of failed runs or steps
- Modularity
- Easy configuration



https://doi.org/10.1038/s41592-021-01254-9





Why Workflows matter?

Why a workflow manager

- Automation
- Run time management
- Software management
- Portability & Interoperability
- Reproducibility
- Scalability (local, cluster, cloud)
- Resuming of failed runs or steps
- Modularity
- Easy configuration





https://doi.org/10.1038/s41592-021-01254-9











X nextflow

- Web-based, no coding required
- Point-and-click interface
- Strong in community/shared workflows

- Python-inspired
- "Makefile for bioinformatics"
- Brief syntax sample
- Emphasize: rule-based, simple to start

- Based on Groovy (scripting language for Java)
- Good for cloud/scalability
- Popular in nf-core community





Anatomy of a Workflow





General Structure



Key Concepts:

- A workflow is made of **steps** (aka *rules* in Snakemake, *processes* in Nextflow).
- Steps are connected by **data flow**: the output of one step becomes the input of the next.
- Each step may:
 - Run a tool or script
 - Process multiple inputs
 - Produce one or more outputs











General Structure











•



Steps: the building blocks of workflows.

- A step is a self-contained unit of work.
- It declares what it needs (inputs), what it produces (outputs), and how to do it (command).

- In Snakemake \rightarrow called a **rule**.
- In Nextflow \rightarrow called a **process**.
- In Galaxy \rightarrow each **tool** is like a process, connected via a GUI.





Steps: the building blocks of workflows.

- A step is a self-contained unit of work.
- It declares what it needs (inputs), what it produces (outputs), and how to do it (command).

- In Snakemake \rightarrow called a **rule**.
- In Nextflow \rightarrow called a **process**.
- In Galaxy \rightarrow each **tool** is like a process, connected via a GUI.







General Structure



The data flow:

- In Nextflow, channels are the way data flows between processes.
- In Snakemake, **dependencies** are defined by files a rule runs when its inputs are ready.



A simple Directed Acyclic Graph (DAG)





An example workflow

Steps:

- 1. Prepare the potatoes $\rightarrow c$
- 2. Cook the fries
- 3. Prepare the patties
- 4. Cook the patties
- 5. Assemble the burgers
- 6. serve

- \rightarrow cut_potatoes
 - \rightarrow fry_fries (after 1)
- \rightarrow shape_patties
 - \rightarrow cook_patties (after 3)
 - \rightarrow assemble_burgers (after 4)
 - \rightarrow serve (after 2 and 5)







ule name: input: <input files="" functions="" or=""/>		
output: <output files=""></output>		
params: <extra e.g.="" options="" parameters,="" tool=""></extra>	#	optional
threads: <number of="" threads="" to="" use=""></number>	#	optional
resources: <resource like="" limits="" mem_mb,="" time=""></resource>	#	optional
conda: " <path env.yaml="" to="">"</path>	#	optional
shell: <command run="" to=""/>		



r



An example workflow with Snakemake

rule name: input: <input files or functions> output: <output files> params: <extra parameters, e.g. tool options> # optional threads: <number of threads to use> # optional resources: <resource limits like mem_mb, time> # optional conda: "<path/to/env.yaml>" # optional shell:

<command to run>

.....

Snakefile SAMPLES = ["order1", "order2", "order3"] rule all: input: expand("served/{sample}.txt", sample=SAMPLES) rule cut_potatoes: output: shell: "echo 'Pommes de terre coupées pour {wildcards.sample}' > {output}" rule fry_fries: input: output: "fries/fried_{sample}.txt" shell: "echo 'Frites cuites pour {wildcards.sample}' > {output}" rule shape_patties: output: "burgers/raw_patty_{sample}.txt" shell:

"echo 'Steak haché formé pour {wildcards.sample}' > {output}"



An example workflow with **Cnextflow**

```
process name {
    tag "<short description for logs>"
    publishDir "results/", mode: 'copy'
                                                // optional
    cpus 4
                                                // optional
    memory '8 GB'
                                                // optional
    time '2h'
                                                // optional
    errorStrategy 'retry'
                                                // optional
    container 'biocontainers/tool:version'
                                                // optional
    input:
        <process input>
    output:
        <process output>
    when:
        <condition>
                                                // optional
    script:
    .....
    <command to run>
    . . . .
```





An example workflow with **X nextflow**



process name { tag " <short description="" for="" logs="">"</short>	
publishDir "results/", mode: 'copy'	// optional
cpus 4 memory '8 GB' time '2h' errorStrategy 'retry'	// optional // optional // optional // optional
container 'biocontainers/tool:version'	// optional
input: <process input=""></process>	
output: <process output=""></process>	
when: <condition></condition>	// optional
script:	
<command run="" to=""/>	
}	

1	// main.nf (Nextflow DSL2 version)
2	nevt flow enable dsl = 2
4	
	params.orders = ['order1', 'order2', 'order3']
6	Channel.fromList(params.orders).set { orders }
8	// Step 1: Cut potatoes
	process CutPotatoes {
10	tag "\\$order"
11	input:
12	val order from orders
13	output:
14	<pre>path "fries/cut_\\${order}.txt" into cut_fries_ch</pre>
15	
16	script:
17	
18	mkdir -p fries
19	echo 'Pommes de terre coupées pour \\$ {order}' > fries/cut _\\$ {order}.txt
20	
21	}
22	
23	// Step 2: Fry fries
24	process FryFries {
25	tag "\\$order"
26	input:
27	<pre>pain cut_file from cut_files_ch.map { file -> def norm = file setPercharg() replaceFirst((Acut (11))</pre>
20	<pre>def name = file.getBaseName().reptaceFirst(/*cut_/, **) [namefile]</pre>
20	lindme, filej
30	s output:
32	nath "fries/fried \\${order} txt" into fried fries ch
32	path hites/hitea_tytordel)/text into hitea_hites_en
34	script:
35	
36	order name=\\$(basename \\${cut file} sed 's/^cut //' sed 's/\.txt//')
37	mkdir -p fries
38	echo 'Frites cuites pour \\${order_name}' > fries/fried_\\${order_name}.txt
39	
40	}





An example workflow with **= Galaxy**

💶 Galaxy	Analyze Data	Workflow V	/isualize -	Shared Data 🔻	Admin	User -	Help 🗸				Using 3	30.2 GB	
Tools	Ne	w workflow						ć	> 6	ľ	۵		
search tools	8							Name					
								New wor	kflow				
Inputs								Version					
Get Data								Version	0, 0 steps	(active	e)	\$	
Expression Tools								Annotatio	n				
Differential Expression								New ann	otation				
IGV screenshots												4	
Collection Operations								These note workflow is	es will be s viewed.	visible	when th	IS	
Lift-Over								Tags					
Text Manipulation								4					
Convert Formats								Apply tags and find ite	to make ems with t	t easy he sam	to searc ne tag.	h for	
Filter and Sort													
Join, Subtract and Group													
Fetch Alignments/Sequence	es												
<		100% +	F								A s	simple text dataset	







.



Parallelisation:







Parallelisation:

• When two steps can be applied simultaneously.





Parallelisation:

• When two steps can be applied simultaneously.







Parallelisation:

- When two steps can be applied simultaneously.
- When the same step can be applied to multiple samples





Parallelisation:

- When two steps can be applied simultaneously.
- When the same step can be applied to multiple samples

Inputs





Parallelisation:

- When two steps can be applied simultaneously.
- When the same step can be applied to multiple samples

Inputs

INSTITUT FRANCAIS DE BIOINFORMATIQUE



FRANCE



Parallelisation:

- When two steps can be applied simultaneously.
- When the same step can be applied to multiple samples

Inputs







Parallelisation:

- When two steps can be applied simultaneously.
- When the same step can be applied to multiple samples

Inputs







Parallelisation:

- When two steps can be applied simultaneously.
- When the same step can be applied to multiple samples







Parallelisation:

- When two steps can be applied simultaneously.
- When the same step can be applied to multiple samples

- Workflow engines know what can be done at the same time, they optimize execution.
- You don't need to write loops, the workflow engine handles it.





Operators (Nextflow):

- They allow dynamic and reactive workflows
- Operators let you transform channels
 - Filter (filter, first, unique)
 - Combine (map, groupTuple, collect, flatten)
 - Process text (splitSV, splitFasta)
 - Fork (multiMap, branch)
 - Math (count, min, max)







Takeaways



Key takeaways

What is a Workflow?

- A workflow is a formal way to chain multiple processing steps.
- It ensures automation, reproducibility, and scalability.

Why Use a Workflow System?

- No more handwritten tracking
- Each step runs only when needed
- Automatic error detection and resuming
- Tracks software versions, parameters, and inputs/outputs
- Parallel execution = faster processing of many samples

When Should You Consider One?

- When your analysis involves different steps
- When you need to **reuse**, **share**, or **scale** your work
- When you want to run the same analysis on multiple datasets
- When you're collaborating with others or aiming for FAIR principles





Key takeaways

Which Workflow manager?

- Galaxy
 - non-coders
 - web interface
 - community workflows
- Snakemake
 - simple, readable workflows
 - great for Python users
- Nextflow
 - scalable and cloud-ready
 - strong for pipelines & HPC













Trainings:

- IFB WF4bioinfo training <u>https://moodle.france-bioinformatique.fr/course/view.php?id=29</u>
- Galaxy https://training.galaxyproject.org/
- Snakemake <u>https://snakemake.readthedocs.io/en/stable/tutorial/tutorial.html</u>
- Nextflow <u>https://training.nextflow.io/2.1.7/fr/</u>

Communities:

- WorkflowHub <u>https://workflowhub.eu/</u>
- Galaxy <u>https://galaxyproject.org/community/</u>
- nf-core (Nextflow) <u>https://nf-co.re/</u>

Tools:

- Seqera AI (Snakemake) <u>https://seqera.io/ask-ai/chat</u>
- GENIAC (Nextflow) https://geniac.readthedocs.io/en/latest/intro.html







Thank you for your attention !









