# ZS **Discovery**

Elucidate. Innovate. Accelerate.

# Nextflow fundamentals training

## DTU – Biosustain

Albert and Felipe - 2025 Nov, 12

# What will we cover?

» • Who am I?

» • What is Nextflow?

» • Understanding the tool

  • Basic concepts

  • Core features

» • Behind Nextflow

  • Why is a community needed?

  • The community efforts

» • Use of the tool in Industry

# Why are we here today?

**Big Data**

Experiments and datasets only get bigger and bigger

**Scalability**

With this, scalability becomes a bottleneck

**Reproducibility**

At the end, analyses must still be reproducible anywhere

# That is where Nextflow comes



## Nextflow enables reproducible computational workflows

Paolo Di Tommaso, Maria Chatzou, Evan W Floden, Pablo Prieto Barja, Emilio Palumbo & Cedric Notredame ✉

*Nature Biotechnology* **35**, 316–319 (2017) │ Cite this article

Nextflow addresses these challenges, allowing us to build efficient workflows coupled with strong supportive community behind.
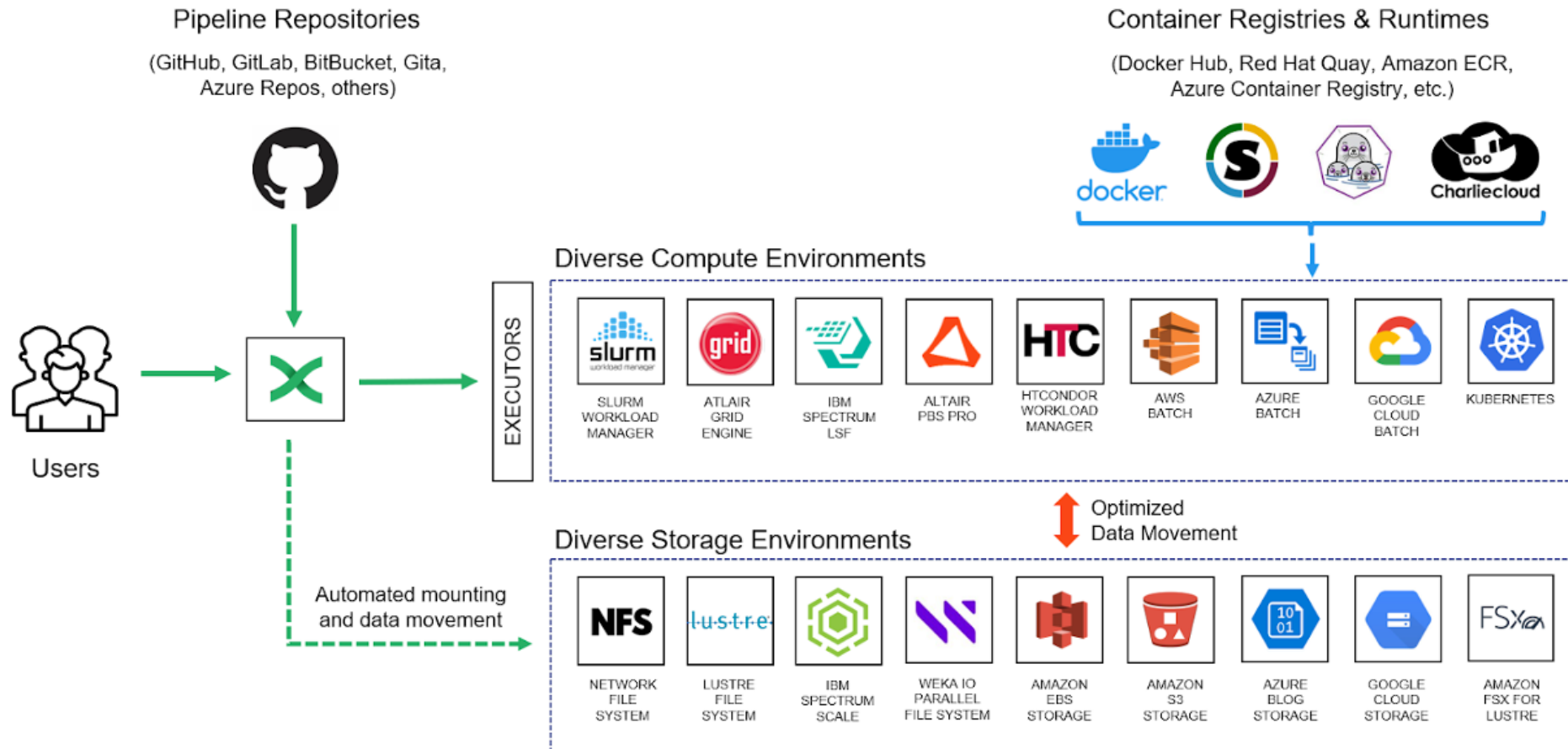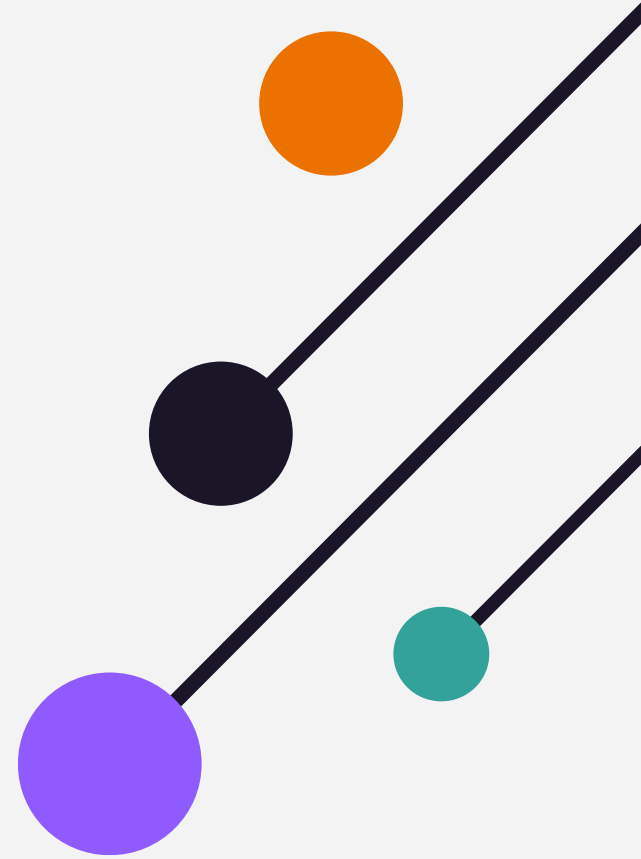
# What is Nextflow?



**In summary**

- Nextflow is a **software** – Workflow orchestrator engine

- It is a Domain-specific language (DSL) built on top of Groovy

- Workflow orchestrator engine:

  - It "wraps-up" tasks as submission scripts for different computing environments and queuing systems

  - Allows use of preferred programing languages for tasks
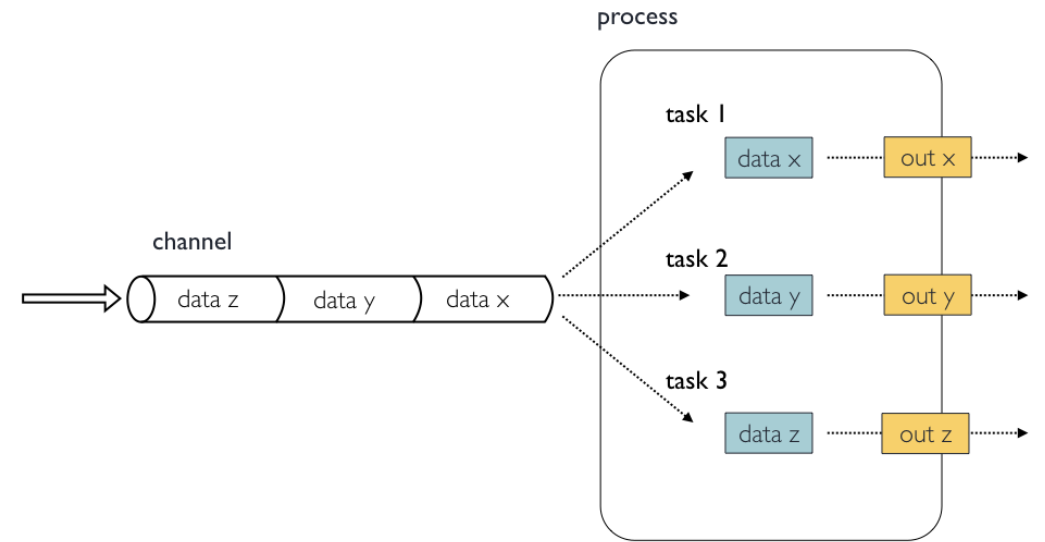
# What does it mean in practice?

# A bit about the tool before the community and its use

# Nextflow is based on few primitives

- **Process:** is every step in your pipeline, and they are executed independently isolated from each other

- **Channels:** Control the data use in processes, and allows the connection of inputs and outputs between processes to make a dependency rule

- **Workflow:** Is the final pipeline itself. If all processes and channels connected, setting the full dependency graph that sets the order of your pipeline execution

# A Nextflow script

**The .nf files** are workflow scripts

```
 main.nf
42   /*
43    * Quickly checking raw reads quality
44    */
45   process FASTQC {
46       container "quay.io/biocontainers/fastqc:0.12.1--hdfd78af_0"
47       tag "FASTQC on $sample_id"
48
49       input:
50       tuple val(sample_id), path(reads)
51
52       output:
53       path "fastqc_${sample_id}_logs"
54
55       script:
56       """
57       mkdir fastqc_${sample_id}_logs
58       fastqc -o fastqc_${sample_id}_logs -q ${reads}
59       """
60   }
61
62   workflow {
63       Channel
64           .fromFilePairs(params.reads, checkIfExists: true)
65           .set { read_pairs_ch }
66       fastqc_ch = FASTQC(read_pairs_ch)
67       fastqc_ch.view()
68   }
```

Directives

Code block

Process

Channel

Workflow
(or sub-workflow)

One can also organize separate blocks of pre-defined workflows that can be "glued" together

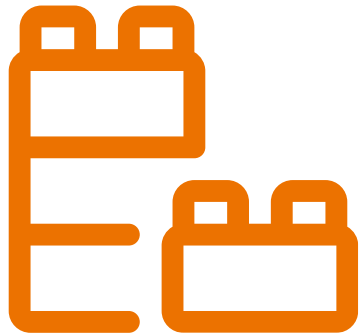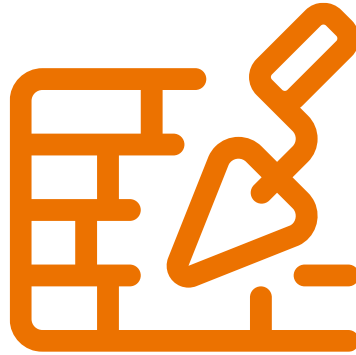--- the "sub-workflows"

# An extended example

**main.nf**

```
1   include { GREETING_WORKFLOW } from './workflows/greeting'
2   include { TRANSFORM_WORKFLOW } from './workflows/transform'
3
4   workflow {
5       names = Channel.from('Alice', 'Bob', 'Charlie')
6
7       // Run the greeting workflow
8       GREETING_WORKFLOW(names)
9
10      // Run the transform workflow
11      TRANSFORM_WORKFLOW(GREETING_WORKFLOW.out.timestamped)
12
13      // View results
14      TRANSFORM_WORKFLOW.out.upper.view { "Uppercase: $it" }
15      TRANSFORM_WORKFLOW.out.reversed.view { "Reversed: $it" }
16  }
```

# In other words …
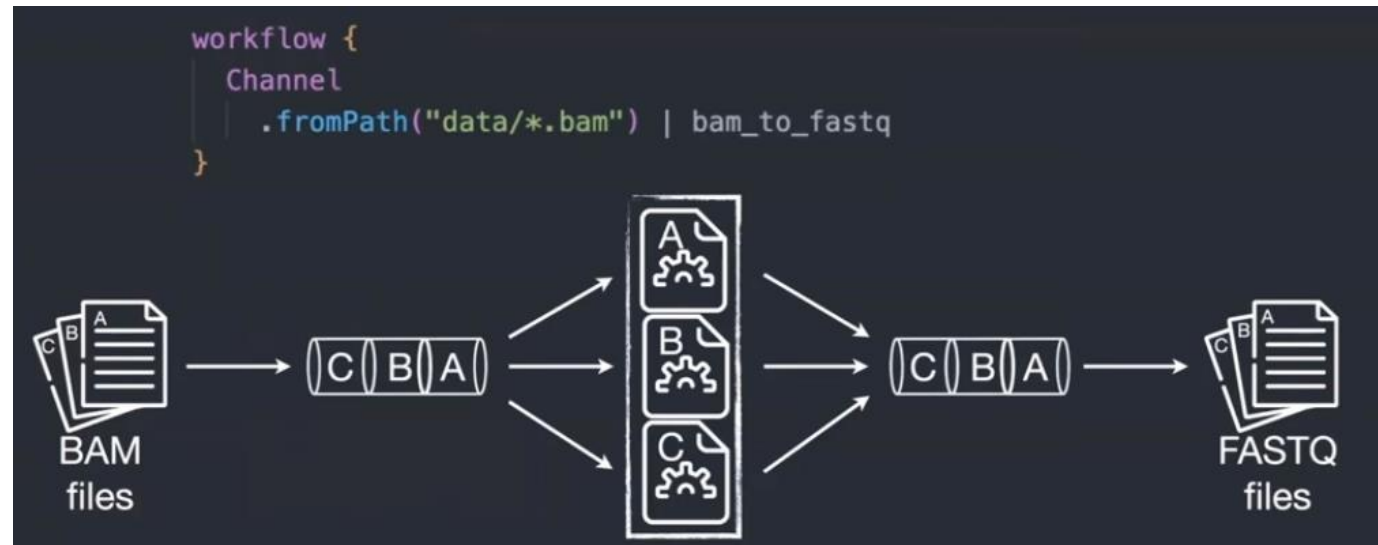


Processes / Modules

Sub-workflows

Pipeline / workflow

# Nextflow – Interesting features

- Implicit **parallelism** (tasks in a process are run by default in parallel)
  - Every module is a "startable" task as long as its input channel exists

- **Re-entrancy** (resume partial runs, do not need to rerun the entire pipeline when something went wrong, it starts from wherever it stopped)

- **Reusability** (reuse different modules, subworkflows, written and containerized by the Nextflow community)
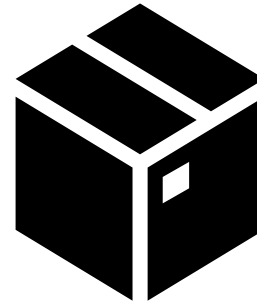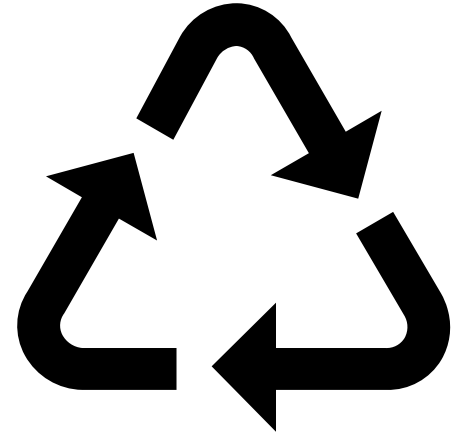
```
workflow {
    Channel
        .fromPath("data/*.bam") | bam_to_fastq
}
```
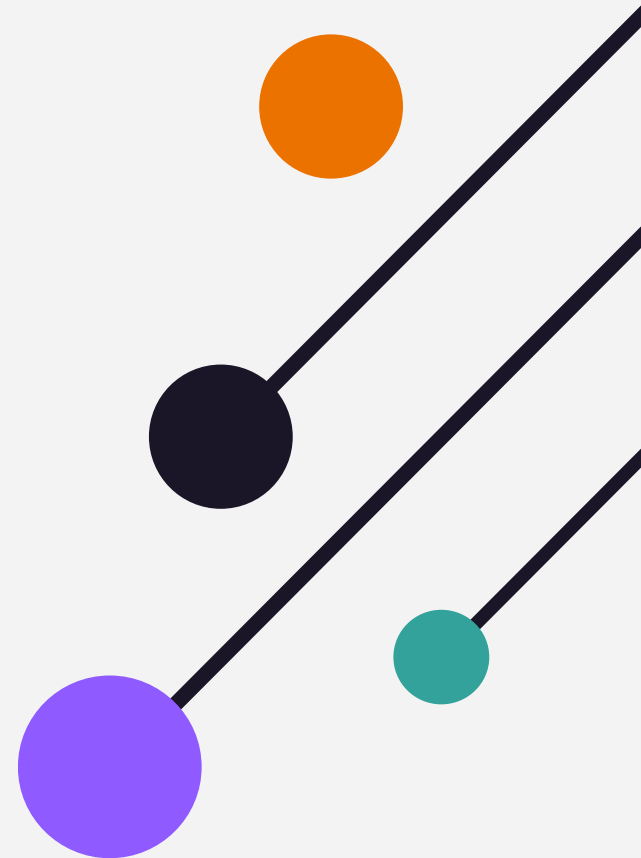


**main.nf**

```
1   include { GREETING_WORKFLOW } from './workflows/greeting'
2   include { TRANSFORM_WORKFLOW } from './workflows/transform'
3
4   workflow {
5       names = Channel.from('Alice', 'Bob', 'Charlie')
6
7       // Run the greeting workflow
8       GREETING_WORKFLOW(names)
9
10      // Run the transform workflow
11      TRANSFORM_WORKFLOW(GREETING_WORKFLOW.out.timestamped)
12
13      // View results
14      TRANSFORM_WORKFLOW.out.upper.view { "Uppercase: $it" }
15      TRANSFORM_WORKFLOW.out.reversed.view { "Reversed: $it" }
16  }
```

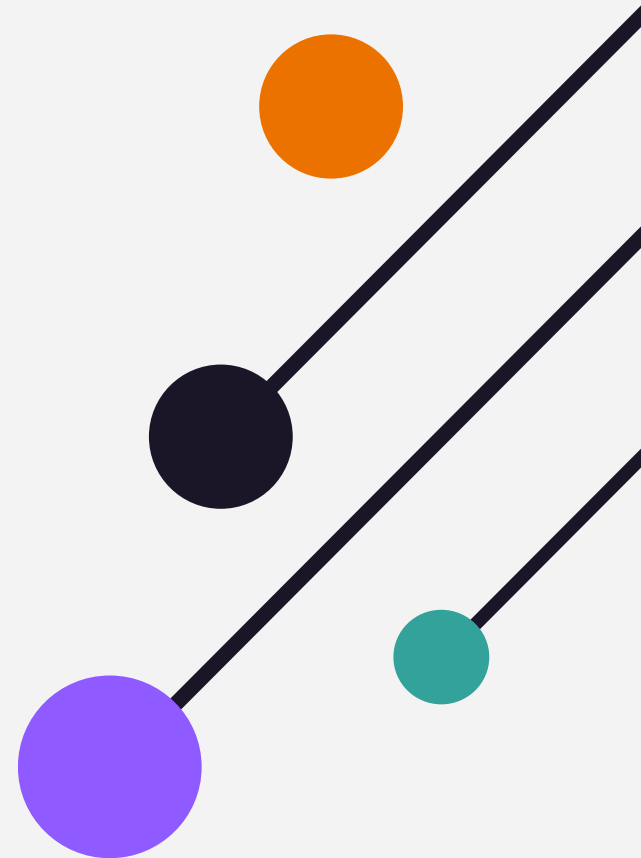# How Nextflow address the challenges of today

- **Reproducible** between runs
  - integration with code management tools
  - all packages downloaded, organized in containers, and control over computing environment
- **Portable** between systems
  - you can write the code in your laptop and can run everywhere (HPC, cloud)
  - works with most of computing environments
- **Scalable**
  - it can be run for 10 on your laptop or thousands of samples in an HPC or the cloud
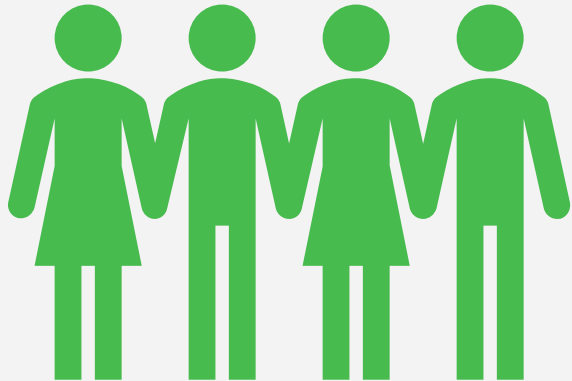- **Integration** of existing tools, systems, and industry standards
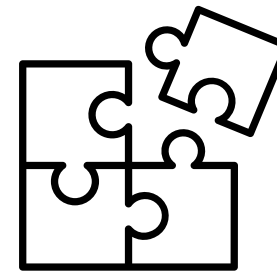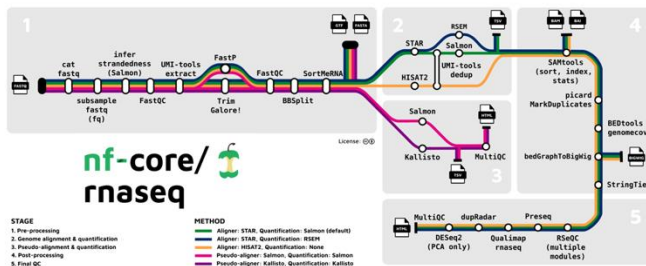
# Break

# Behind Nextflow

# The community itself is what makes Nextflow great



- Driven by #OpenScience and #OpenSource

- Driven by the need to know software, params, options, method used, being sure of what is done

- Nextflow advocates for an open community with multiple channels to interact and connect people, forums, Slack channels, conferences, hackathons regularly organize

- **Nf-core**

**nf-core**

A global community effort to collect a curated set of open-source analysis pipelines built using Nextflow.
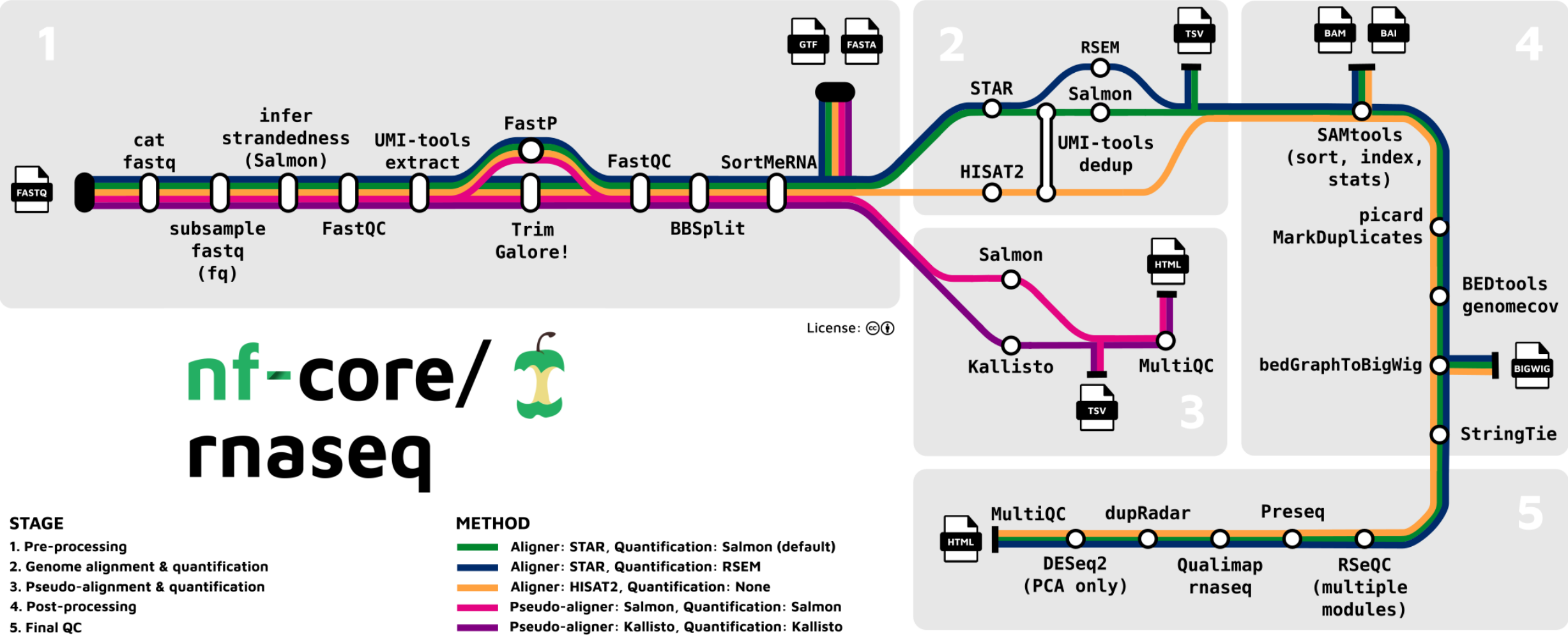https://nf-co.re/


nf-core/rnaseq pipeline diagram



**>100 pipelines** available to process and analyse many different data types

**Modules**: >1,000 reusable components that can be integrated into pipelines

**Subworkflows**: >70 pre-assembled combinations of modules aimed at streamlining commonly used workflows
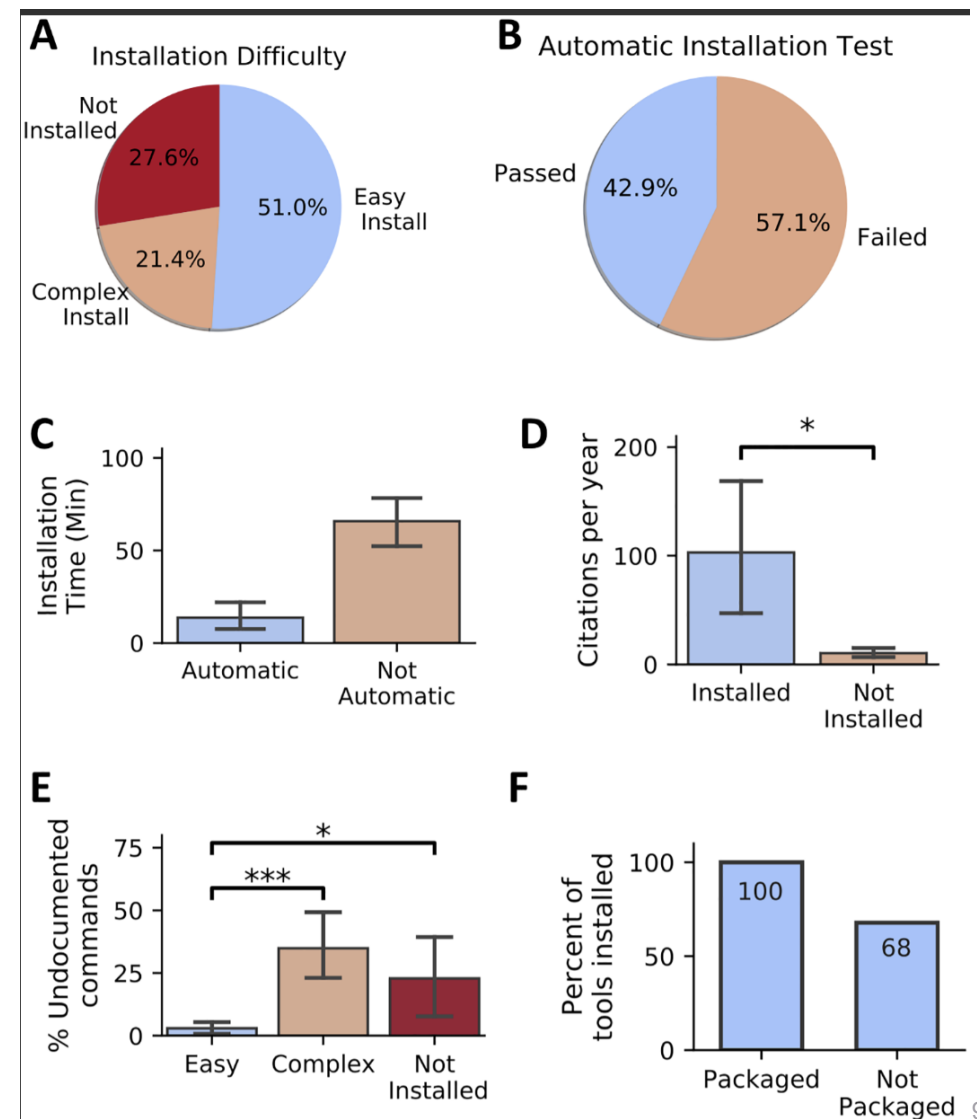
# A nextflow pipeline in metro map

# Why do we need a Community effort?



## Challenges and recommendations to improve the installability and archival stability of omics computational tools
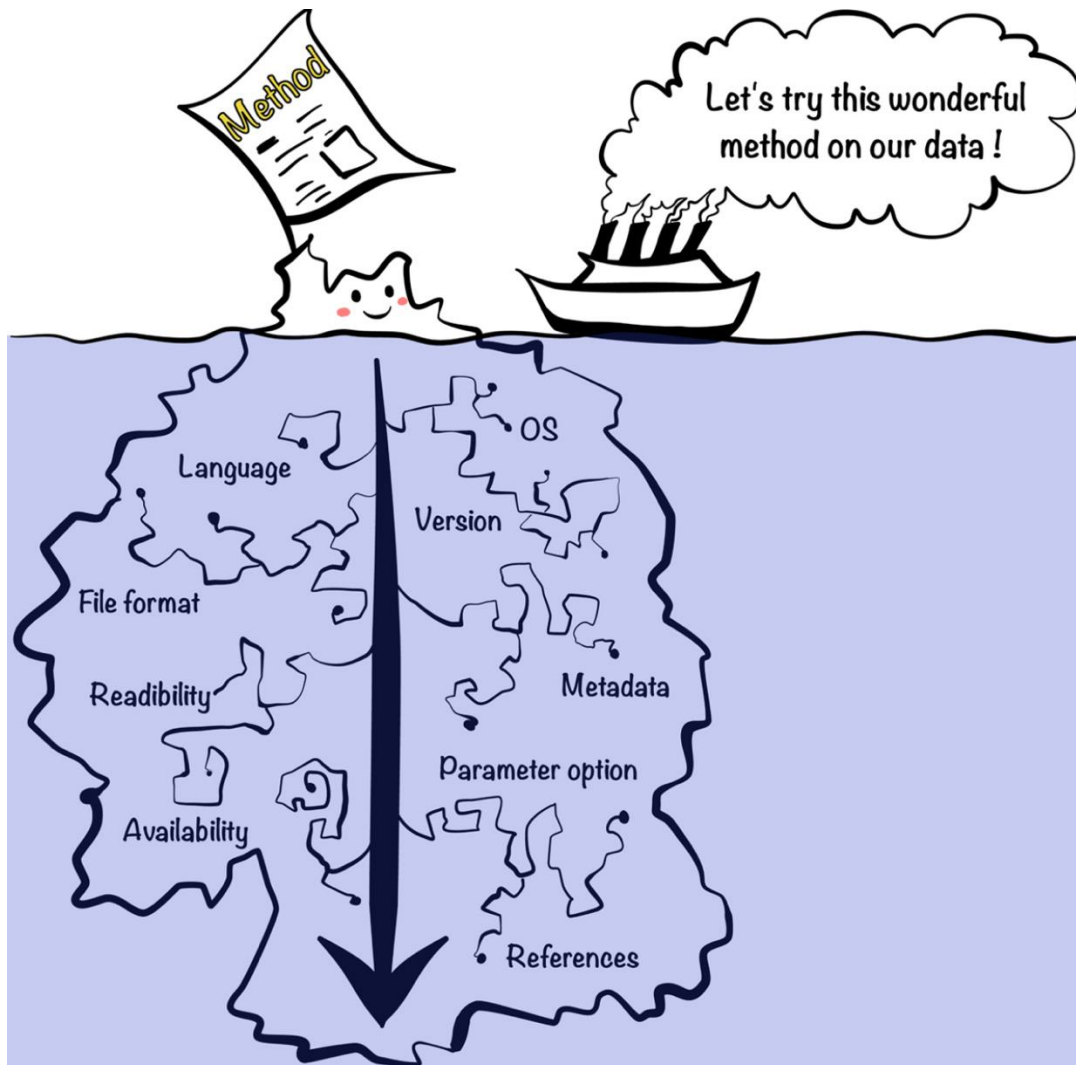
Serghei Mangul co ✉, Thiago Mosqueiro co, Richard J. Abdill, Dat Duong, Keith Mitchell, Varuni Sarwal, Brian Hill, Jaqueline Brito, Russell Jared Littman, Benjamin Statz ※, Angela Ka-Mei Lam, Gargi Dayama, Laura Grieneisen, [ ... ], Ran Blekhman [ view all ]

- 28% of all omics software resources were not accessible through URLs published

- Among the tools found, 49% were difficult to install or could not be installed at all!

9

# Why do we need a Community effort?



**Experimenting with reproducibility: a case study of robustness in bioinformatics** 🔓

Yang-Min Kim ✉, Jean-Baptiste Poline, Guillaume Dumas

*GigaScience*, Volume 7, Issue 7, July 2018, giy077,
https://doi.org/10.1093/gigascience/giy077

- "First we tried to rerun the analysis with the code and the data provided by the authors. Second we reimplemented the whole method in a python package…"

# Nf-core helps addressing some of the challenges

nf-core

Modules are open source and maintained by the community --- **shared ownership**

All modules are containerized with public containers --- **portable and repeatable**

All steps have automated CI/CD testing to ensure results are reproduced even after changes --- **reproducibility**
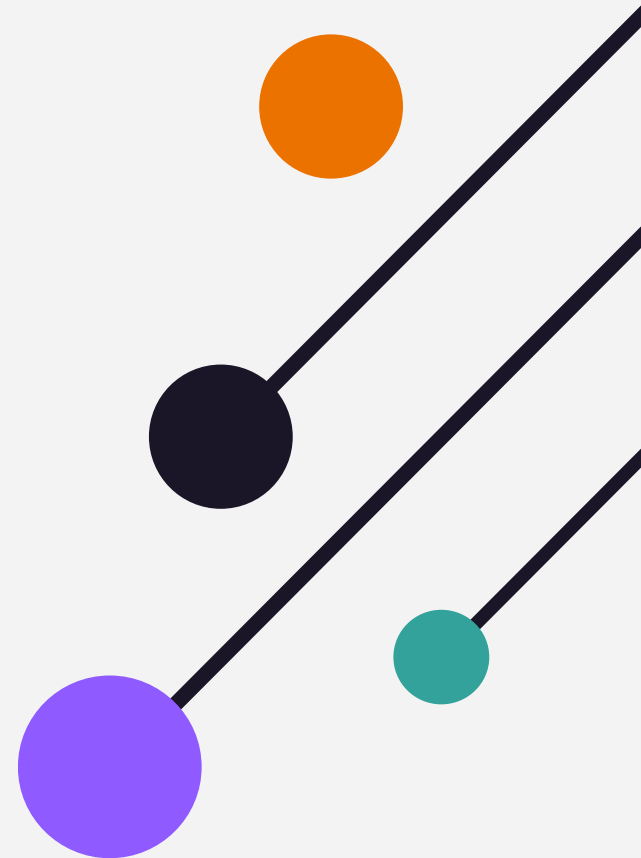
Not only modules, but the same applies to shared sub-workflows and even pipelines --- **reproducibility**
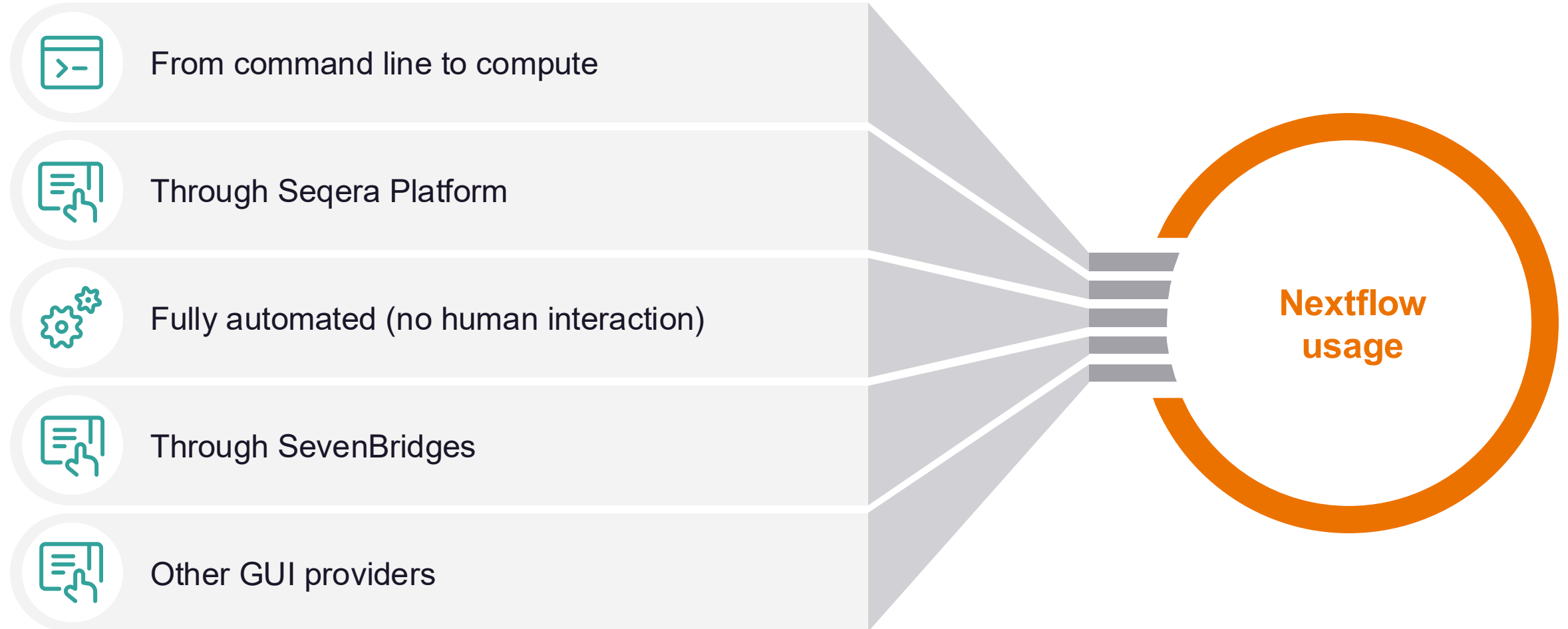
**Even when very careful, reproducibility is challenging.**

Together, via open science and shared ownership, the community ensures all pieces are tested and reproducible all together, while being extremely portable.

# Nextflow in Industry

# Widely adopted, but in many different ways



From command line to compute

Through Seqera Platform

Fully automated (no human interaction)

Through SevenBridges

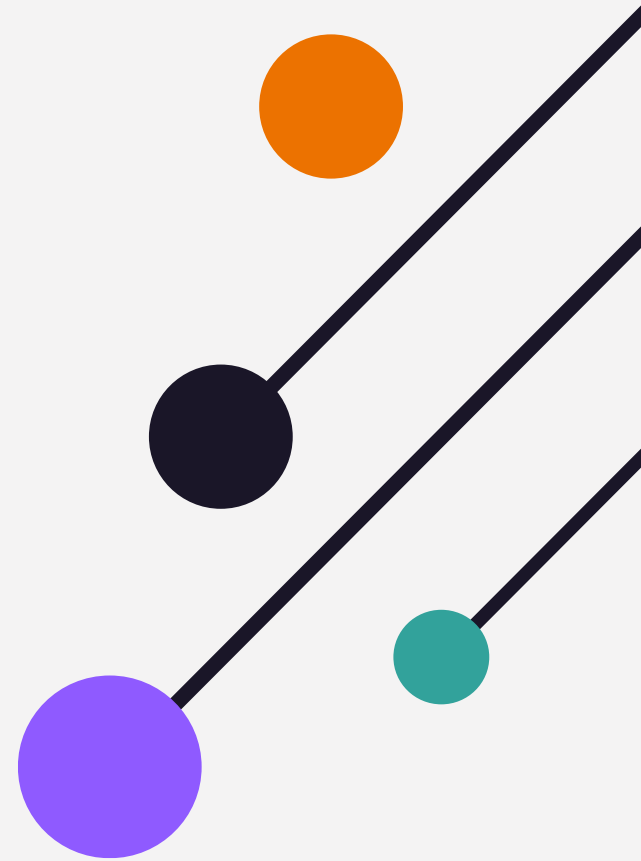Other GUI providers

**Nextflow usage**

# Anyways ...

**nf-core**

- Some use as is, directly from nf-core git and contribute there

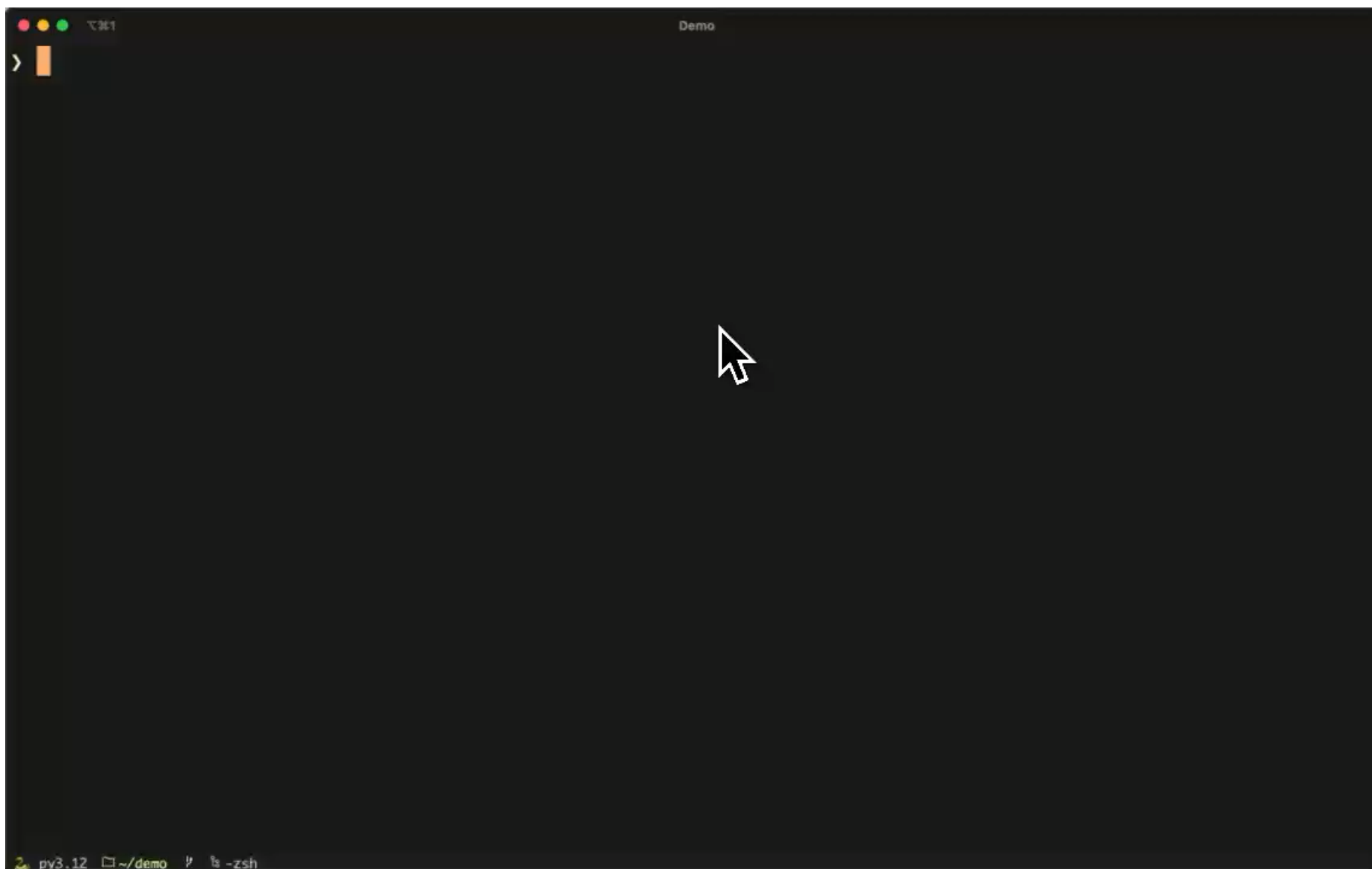- Adding modules, pipelines, sub-workflows, etc. – Being part of it

- Some have internal private services so they "wrap-up extra modules around

- Using the shared modules to benefit from the thorough community testing and validation
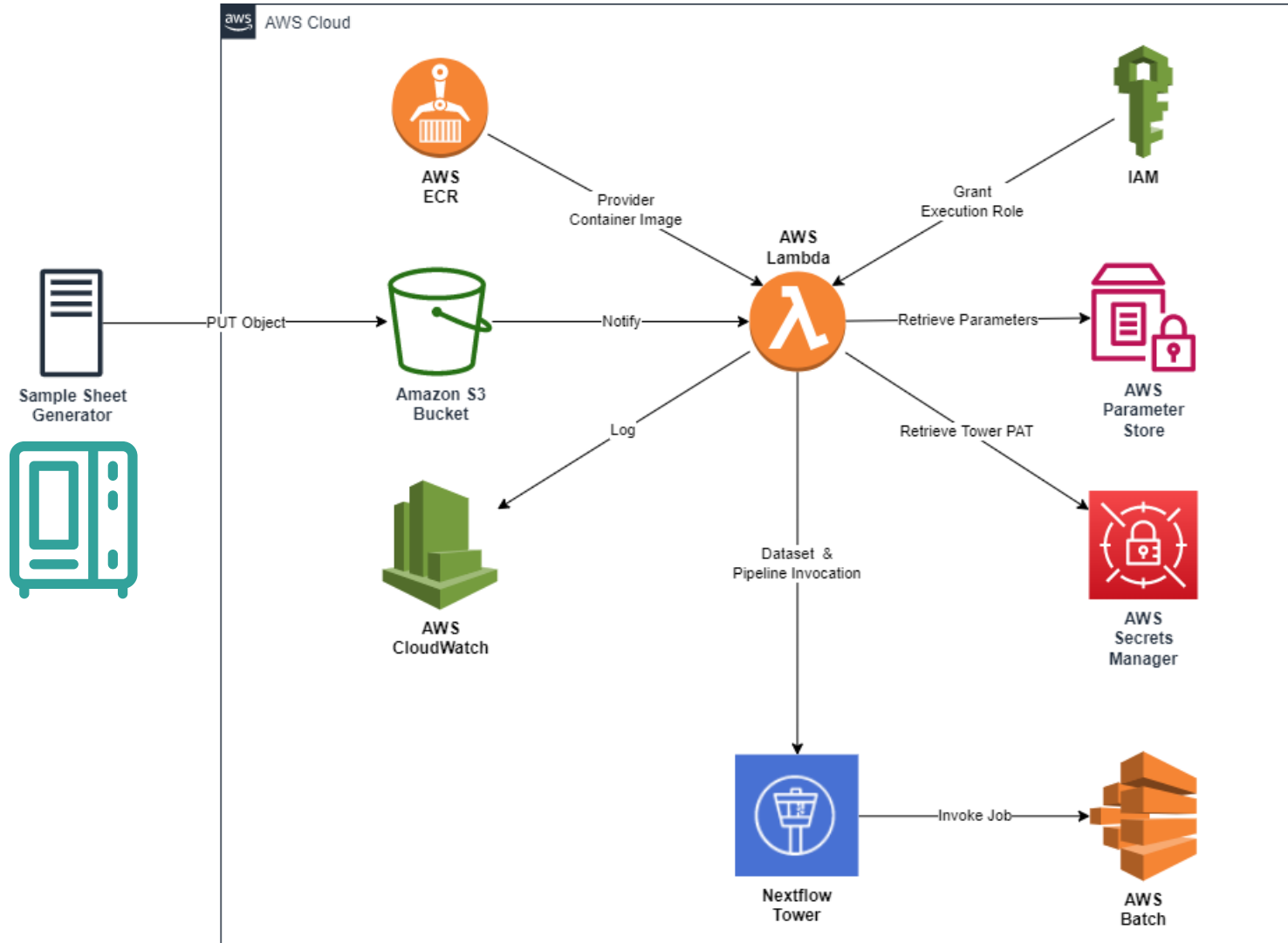
# Some examples of deployment and industry case study

# Via Seqera Platform

# Via automated settings



**Machines upload data that triggers the engine**

# ZS **Discovery**

Elucidate. Innovate. Accelerate.

# Thank you!

Impact where it matters.