# nextflow

**Reproducible, Scalable, and Shareable Data Analysis Workflows**

# Writing modern workflows is complex



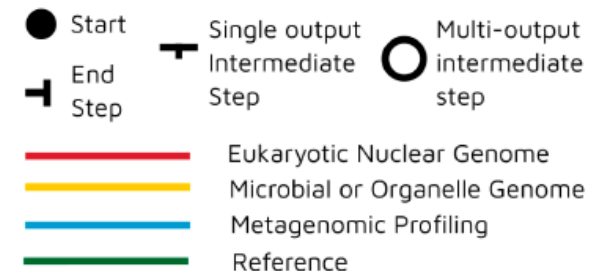nf-core/eager v2.4

Example analysis pathways

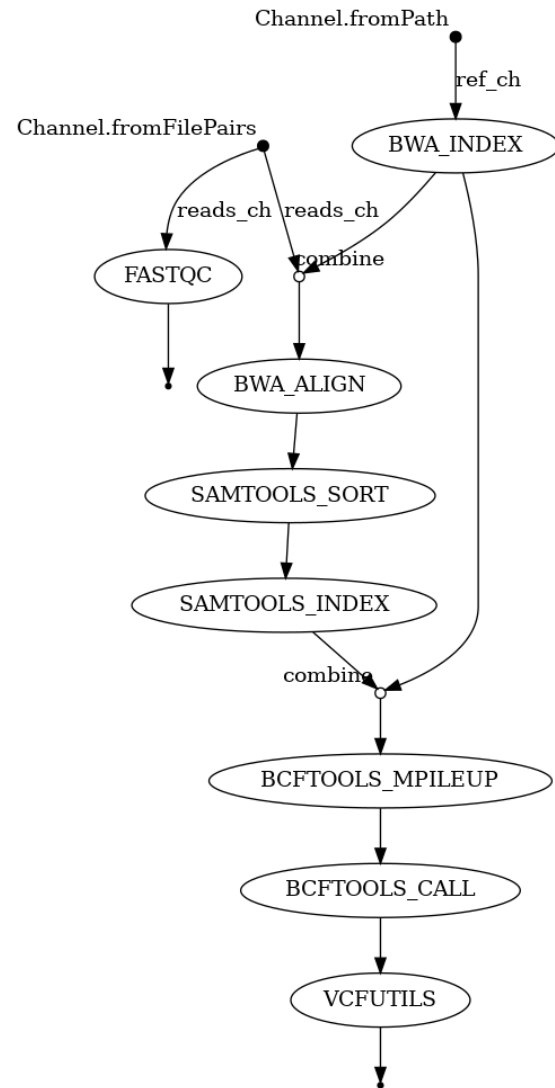state-of-the-art ancient DNA analysis pipeline

# Data analysis workflows

- **Data analysis applications perform computation to generate information from (large) datasets**

- **Embarrassingly parallel! → can spawn 100s-100k jobs over distributed cluster**

- **Mash-up of many different tools and scripts (dependencies!)**

- **Complex dependency trees and configuration → very fragile ecosystem!**

# A lot of moving parts!



- 70 processes

- 55 external scripts

- 39 software tools & libraries

Variant-Calling pipeline

# To reproduce the result of a typical computational biology experiment requires 280 hours.

## Quantifying Reproducibility in Computational Biology: The Case of the Tuberculosis Drugome

Daniel Garijo[1], Sarah Kinnings[2], Li Xie[3], Lei Xie[4], Yinliang Zhang[5], Philip E. Bourne[3*], Yolanda Gil[6*]

1 Ontology Engineering Group, Facultad de Informática, Universidad Politécnica de Madrid, Madrid, Spain, 2 Department of Chemistry and Biochemistry, University of California San Diego, La Jolla, California, United States of America, 3 Skaggs School of Pharmacy and Pharmaceutical Sciences, University of California San Diego, La Jolla, California, United States of America, 4 Department of Computer Science, Hunter College, The City University of New York, New York, New York, United States of America, 5 School of Life Sciences, University of Science and Technology of China, Hefei, Anhui, China, 6 Information Sciences Institute and Department of Computer Science, University of Southern California, LosAngeles, California, United States of America
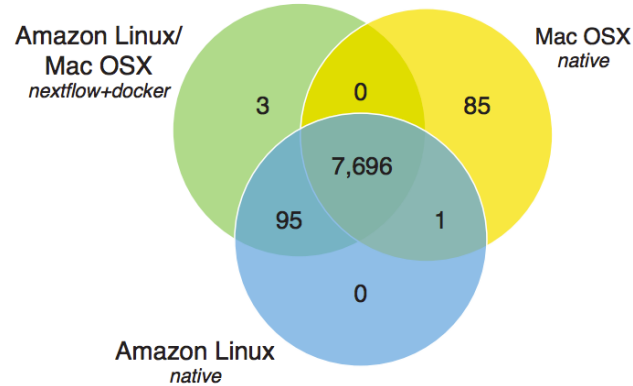
## ≈1.7 months!

The **same application**

deployed in

**different environments**

Produces

**different results**

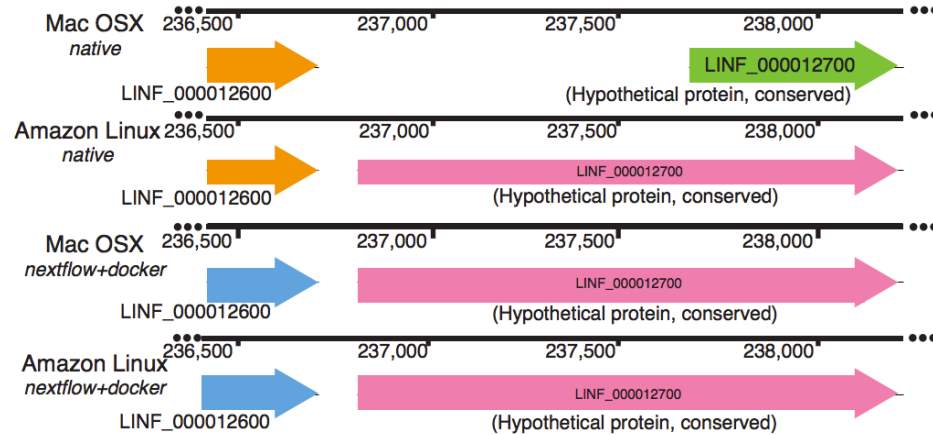# Differences in genome annotations & expressed genes



**a** Gene annotation of *Leishmania infantum* with Companion

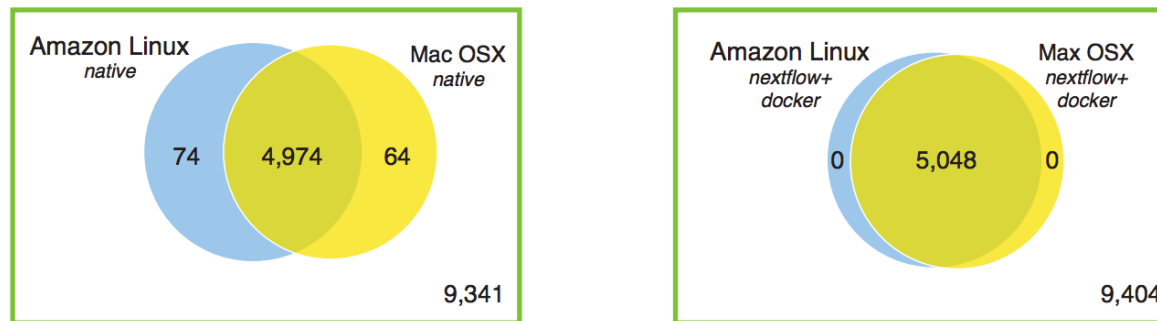**b** *Leishmania infantum* chromosome 1

**c** Transcript quantification and differential expression with Kallisto and Sleuth

VOLUME 35   NUMBER 4   APRIL 2017   NATURE BIOTECHNOLOGY

* Di Tommaso P, et al., *Nextflow enables computational reproducibility*, Nature Biotech, 2017

# Why do Reproducible Research?

An article about computational results is advertising, not scholarship.

The actual scholarship is the full software environment, code and data, that produced the result
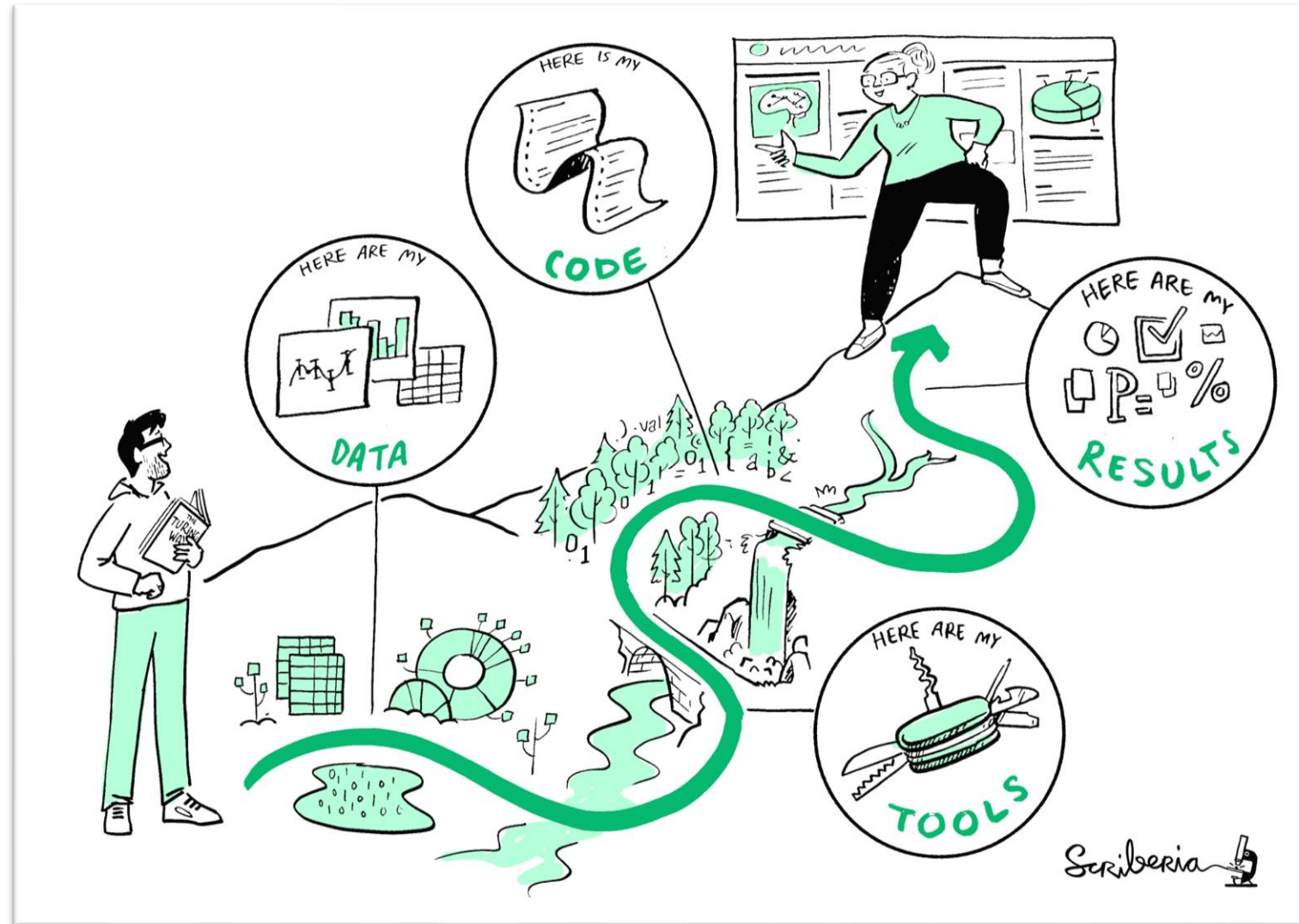- **Claerbout & Karrenbach**

Authors **need to** provide all the necessary data and the computer codes to run the analysis again, re-creating the results.



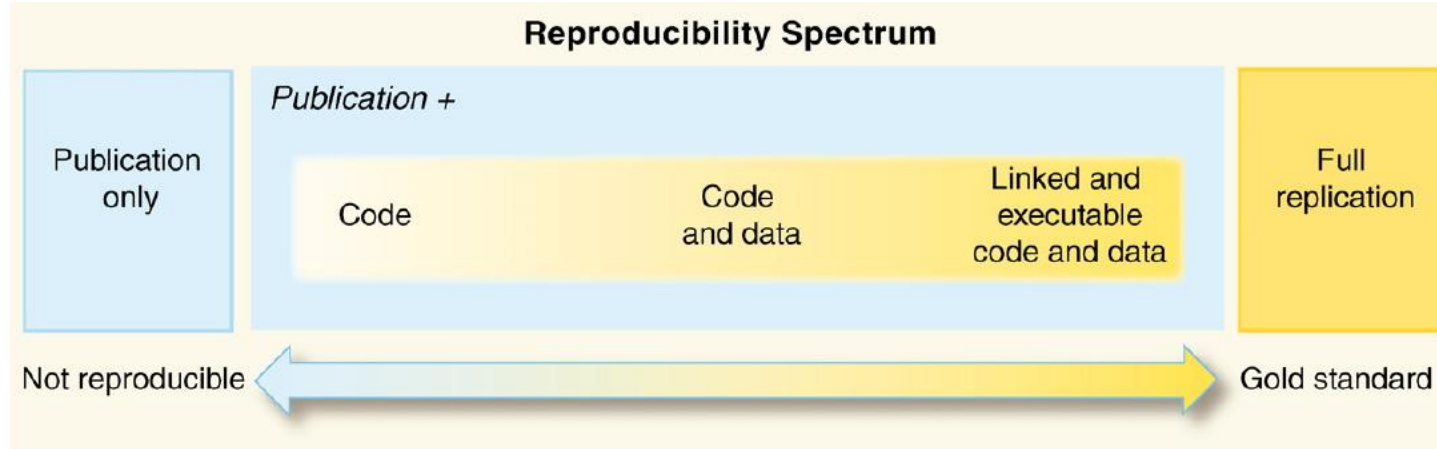*The Turing Way* project illustration by Scriberia. Used under a CC-BY 4.0 licence. DOI: 10.5281/zenodo.3332807.
https://the-turing-way.netlify.app/reproducible-research/

*covers topics related to skills, tools and best practices for research reproducibility.*

# What does this actually look like?



**Reproducibility Spectrum**

Publication only | Publication + (Code | Code and data | Linked and executable code and data) | Full replication
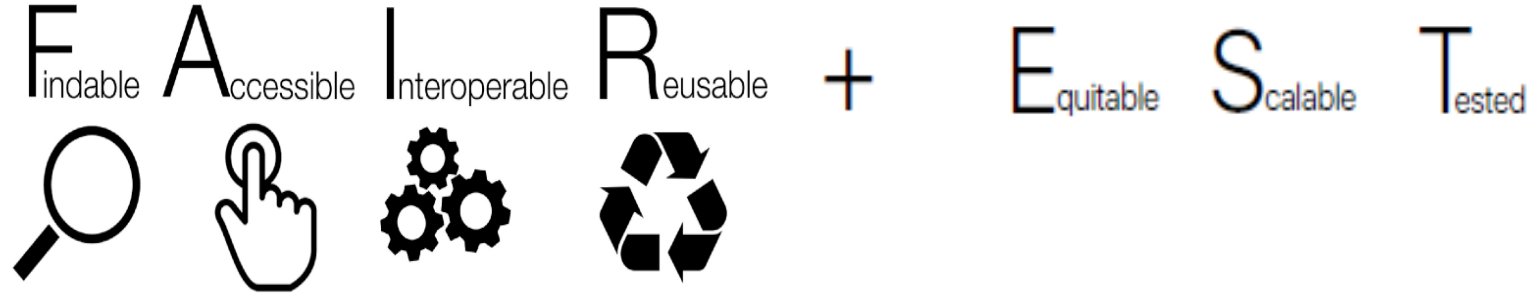
Not reproducible ←→ Gold standard

## Reproducibility Checklist

❑ Uses Code – instead of pointing and clicking

❑ Data archived and available

❑ Open-Source – Code transparency is key to reproducibility

❑ Version tracking – using version-control software as Git and Github/Gitlab

❑ Replicate your environment – using containers or conda

❑ Analyses documented – Jupyter Notebooks or Rmarkdown

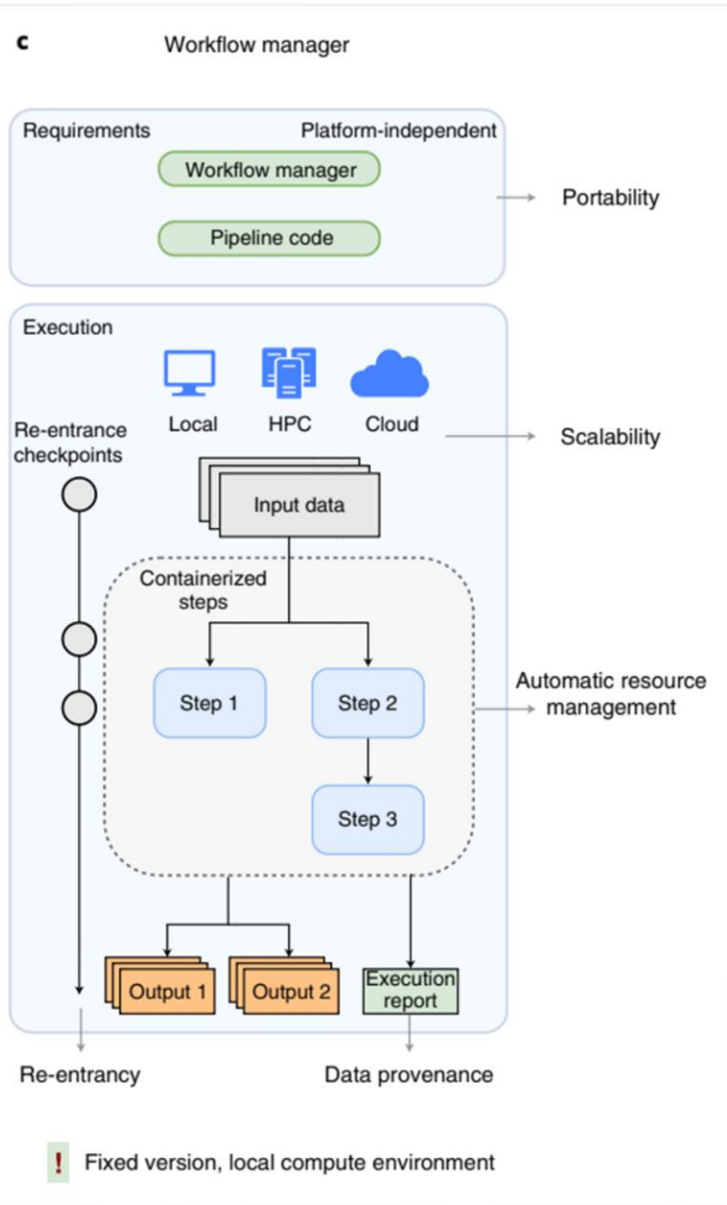❑ Automated – workflow management tool (Snakemake, CWL, Nextflow)

# What should be our goal?

F_indable A_ccessible I_nteroperable R_eusable + E_quitable S_calable T_ested

**Streamlining data analysis to make it:**

• Findable ⇒ openly available and searchable through repositories

• Accessible ⇒ can be used by everyone

• Interoperable ⇒ portable to any cluster & cloud

• Reusable ⇒ in a reproducible way

• Equitable ⇒ free of bias

• Scalable ⇒ from laptop to supercomputer

• Tested ⇒ validated using modern software practices

(WfMSs) automate computational analyses by stringing together individual data processing tasks into cohesive pipelines.

They abstract away the issues of orchestrating data movement and processing, managing task dependencies, and allocating resources within the compute infrastructure.

Wratten, L., Wilm, A. & Göke, J. Reproducible, scalable, and shareable analysis pipelines with bioinformatics workflow managers. Nat Methods 18, 1161–1168 (2021). https://doi.org/10.1038/s41592-021-01254-9

# Aspects of WfMSs relevant to bioinformatics

- **Modularity of the pipeline to enable checkpointing**

- **Scalability with respect to the number of tasks in the pipeline**

- **Robustness against failures due to data issues, resource unavailability, or aborted execution**

- **Reproducibility via logs recording data provenance and task execution**

- **Portability across compute environments**

- **Ease of development by users with a range of experience and computational knowledge.**
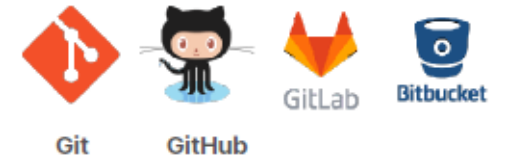
# What is Nextflow?



nextflow *script*

Write code in any language

R, Python, BASH

DATA FLOW PROGRAMMING MODEL

Define orchestration with dataflow programming

Define software dependencies with containers

CONDA, docker, S

Git, GitHub, GitLab, Bitbucket

Version control

nextflow *runtime*

Orchestration of tasks to deploy anywhere with ease

## Supports all major platforms

aws, Google Cloud, Azure

GRID ENGINE, slurm workload manager, HTCondor

Platform Computing an IBM Company, kubernetes, PBS Works

**Nextflow is a language, a runtime and a community**

https://nextflow.io/

# How does it work?

• **Fast prototyping** ⇒ custom DSL that enables task composition, simplifies most use cases + general purpose programming language for corner cases

• **Easy parallelization** ⇒ declarative reactive programming model based on dataflow paradigm, implicit portable parallelism

• **Self-contained** ⇒ functional approach, a task execution is idempotent i.e., cannot modify the state of other tasks + isolate dependencies with containers

• **Portable deployments** ⇒ executor abstraction layer + deployment configuration from implementation logic

# Task example

```
bwa mem reference.fa sample.fq \
| samtools sort -o sample.bam
```
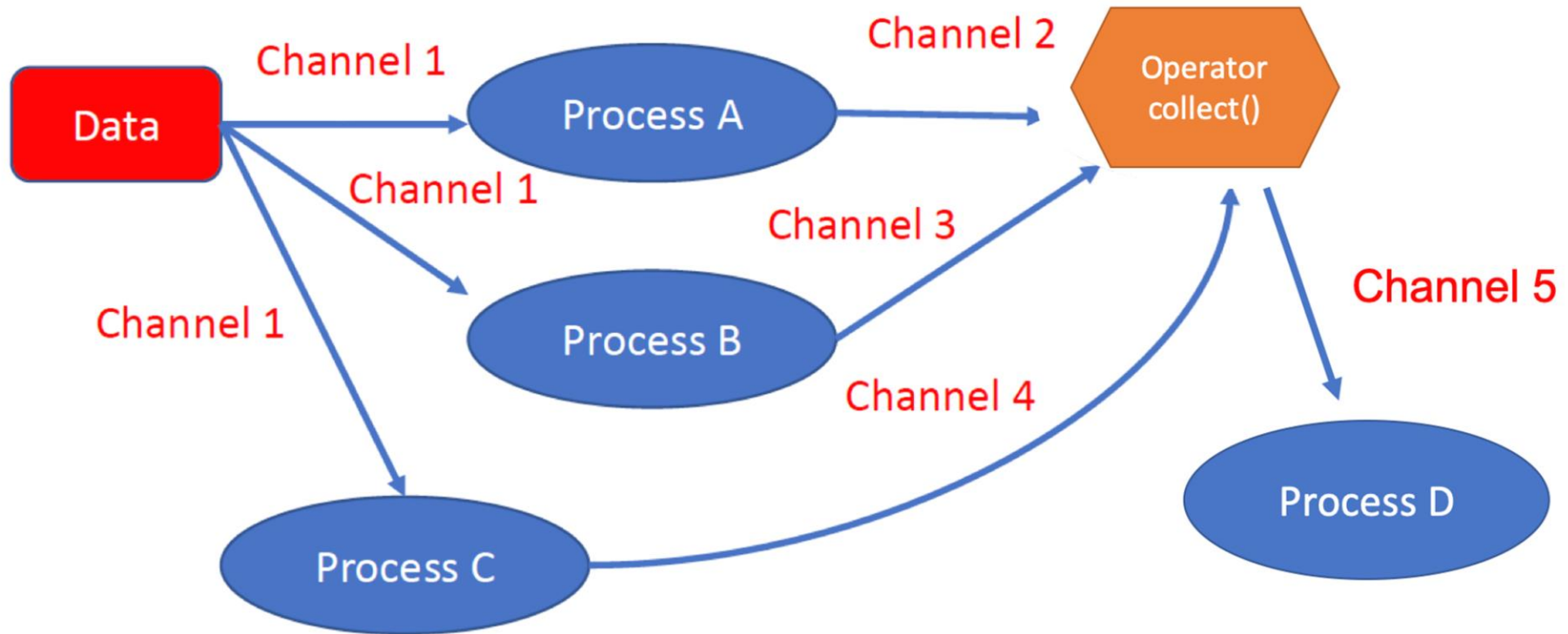
# Task example

```
process align_sample {
    input:
    file 'reference.fa' from genome_ch
    file 'sample.fq' from reads_ch
    output:
    file 'sample.bam' into bam_ch
    script:
    """

    bwa mem reference.fa sample.fq \
        | samtools sort -o sample.bam
    """
}
```

# Task Composition

```
process align_sample {
    input:
    file 'reference.fa' from genome_ch
    file 'sample.fq' from reads_ch
    output:
    file 'sample.bam' into bam_ch
    script:
    """

    bwa mem reference.fa sample.fq \
        | samtools sort -o sample.bam
    """

}
```

```
process index_sample {
    input:
    file 'sample.bam' from bam_ch
    output:
    file 'sample.bai' into bai_ch
    script:
    """

    samtools index sample.bam
    """
}
```

# Workflow

# Nextflow Syntax

task

```
process QUANT {
    input:
    path index
    tuple val(pair_id), path(reads)

    output:
    path pair_id

    script:
    """
    salmon quant -i $index \
            -1 ${reads[0]} \
            -2 ${reads[1]} \
            -o $pair_id
    """
}
```
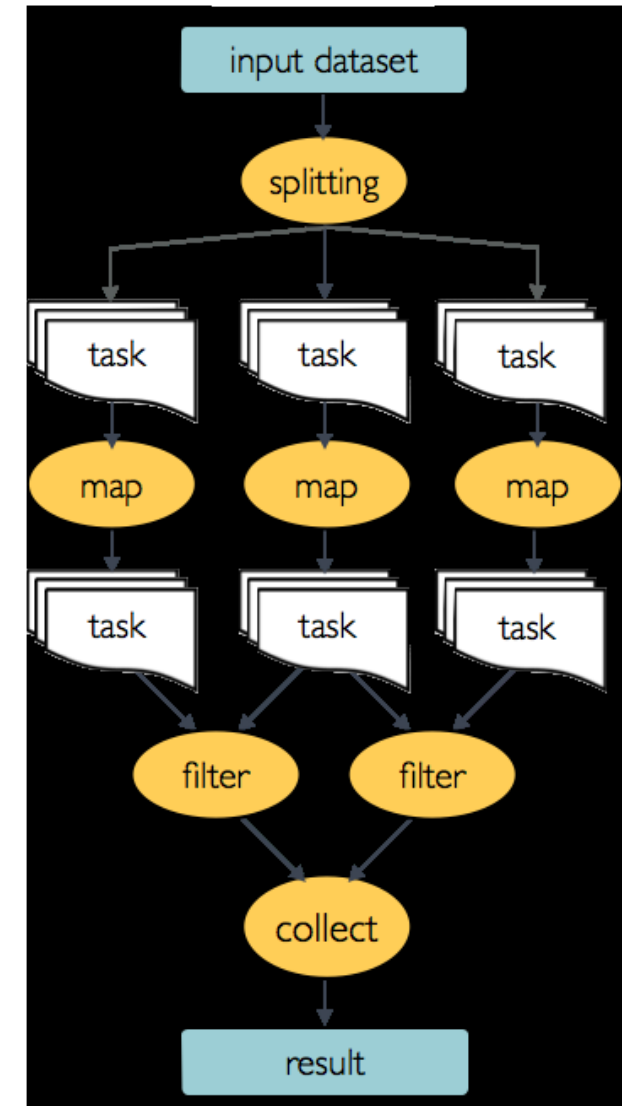
workflow

```
params.outdir = 'results'

include { INDEX } from './index'
include { QUANT } from './quant'
include { FASTQC } from './fastqc'

workflow RNASEQ {
    take:
        transcriptome
        read_pairs_ch

    main:
        INDEX(transcriptome)
        FASTQC(read_pairs_ch)
        QUANT(INDEX.out, read_pairs_ch)

    emit:
        QUANT.out | concat(FASTQC.out) | collect
}
```
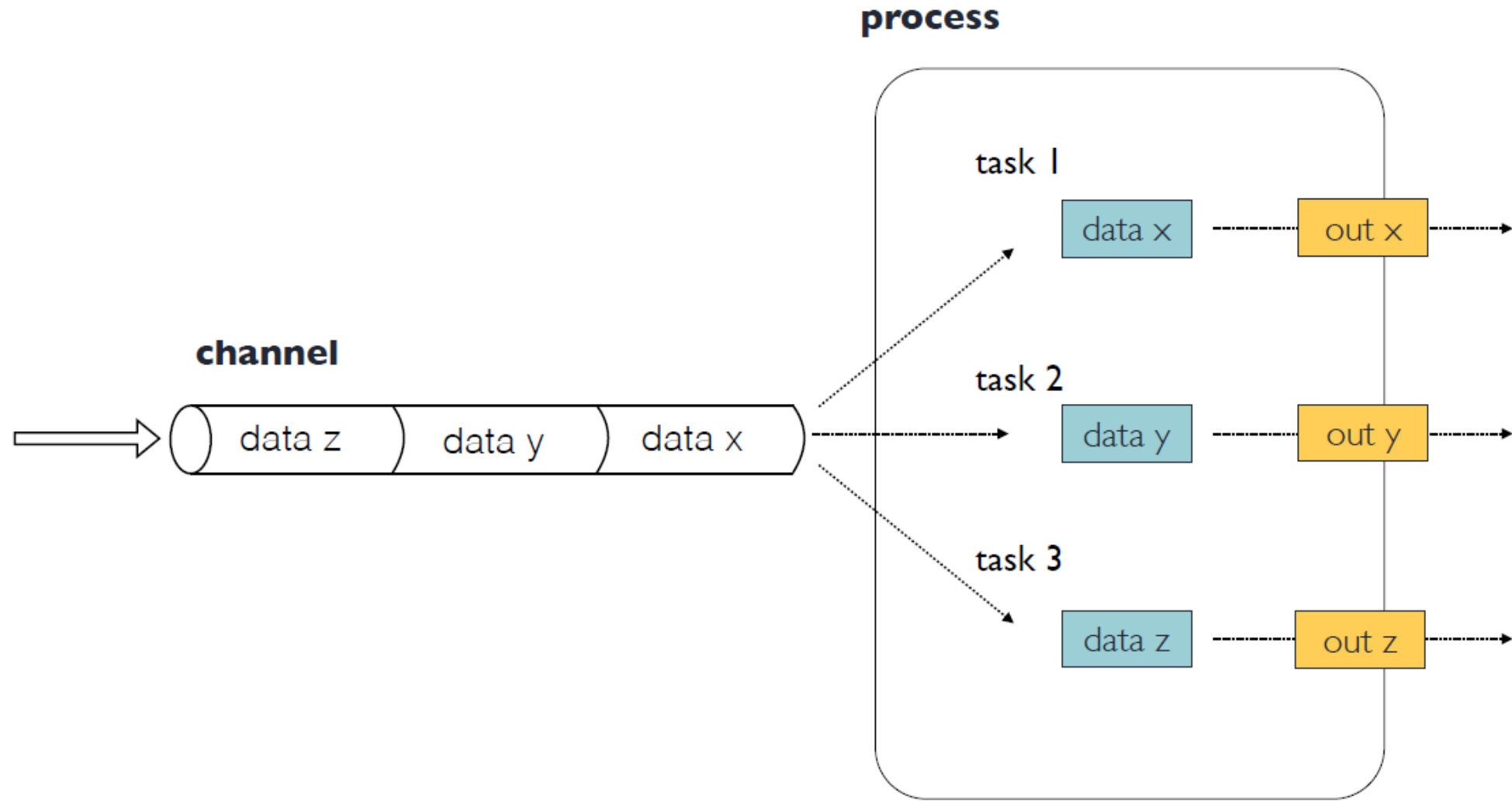
**Scientists and engineers can now write complex, distributed and parallel data pipelines without requiring a degree in computer science.**

# Dataflow concepts

- **Declarative computational model for parallel process executions**

- **Processes wait for data, when an input set is ready the process is executed**

- **They communicate by using dataflow variables i.e., async FIFO queues called channels**

- **Parallelization and tasks dependencies are implicitly defined by process in/out declarations**

# How does parallelization work?

# How does parallelization work?

```
samples_ch = Channel.fromPath('data/sample.fastq')

process FASTQC {

  input:
  file reads from samples_ch
  output:
  file 'fastqc_logs' into fastqc_ch


  script:
  """

  mkdir fastqc_logs
  fastqc -o fastqc_logs -f fastq -q ${reads}
  """

}
```
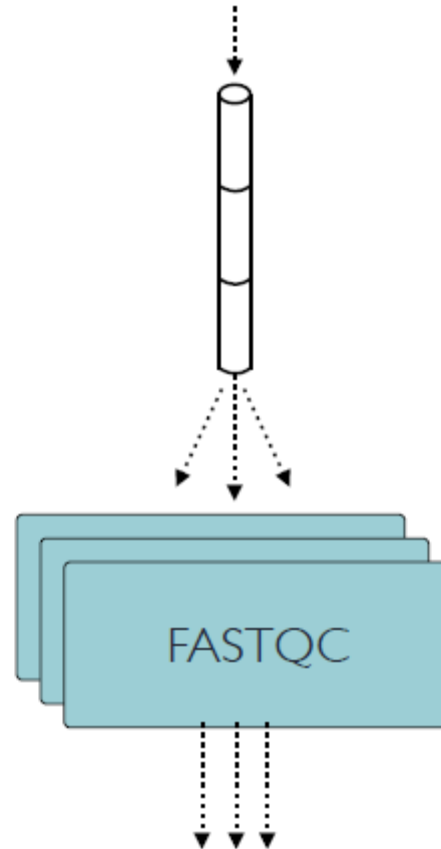
```
samples_ch = Channel.fromPath('data/*.fastq')

process FASTQC {

  input:
  file reads from samples_ch
  output:
  file 'fastqc_logs' into fastqc_ch


  script:
  """

  mkdir fastqc_logs
  fastqc -o fastqc_logs -f fastq -q ${reads}
  """

}
```

```
reads_ch = Channel.fromFilePairs( "data/*_{1,2}.fastq.gz", checkIfExists: true )




[SRR2584863, [/data/SRR2584863_1.fastq.gz, /data/SRR2584863_2.fastq.gz]]
[SRR2584866, [/data/SRR2584866_1.fastq.gz, /data/SRR2584866_2.fastq.gz]]
[SRR2589044, [/data/SRR2589044_1.fastq.gz, /data/SRR2589044_2.fastq.gz]]
```
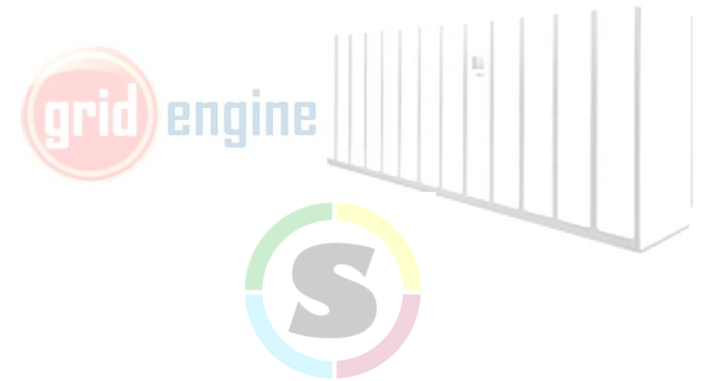
# Implicit parallelism

# CWL vs Nextflow

- Language specification

- Declarative meta-language (YAML/JSON)

- Verbose

- Committee driven

- Many vendors/implementations (and specification version)

- Language + app. runtime

- DSL on top of a general purpose programming lang.

- Concise, fluent

- Community driven

- Single implementation, quick iterations

# Snakemake vs Nextflow

- Command line oriented tool

- Pull model

- Rules defined using file name patterns

- Compute DAG ahead

- Built-in support for Singularity/Docker

- Custom scripts for cluster deployments

- No support for source code management system

- Python based

- Command line oriented tool

- Push model

- Can manage any data structure

- Compute DAG at runtime

- All major container runtimes

- Built-in support for clusters and cloud

- Built-in support for Git/GitHub, etc., manage pipeline revisions
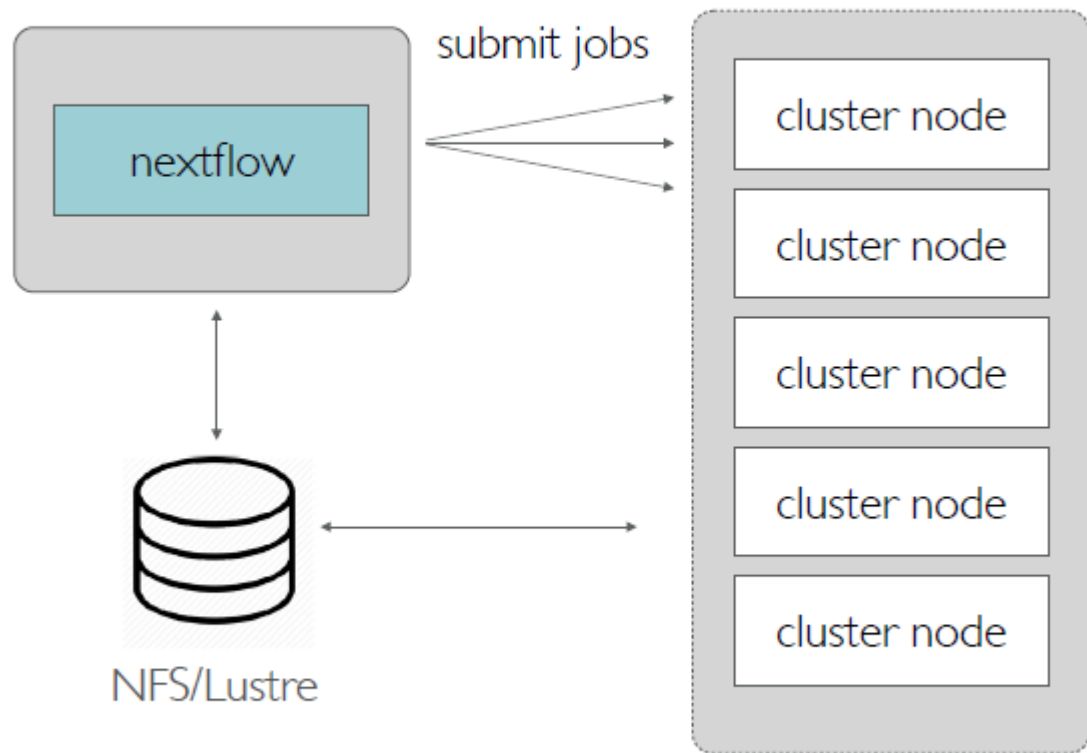
- Groovy/JVM based

# Deployment Scenarios

# Local execution

laptop / workstation

nextflow

docker/singularity

OS

local storage

- **Common development scenario**

- **Dependencies can be managed using a container runtime**

- **Parallelization is managed by spawning posix processes**

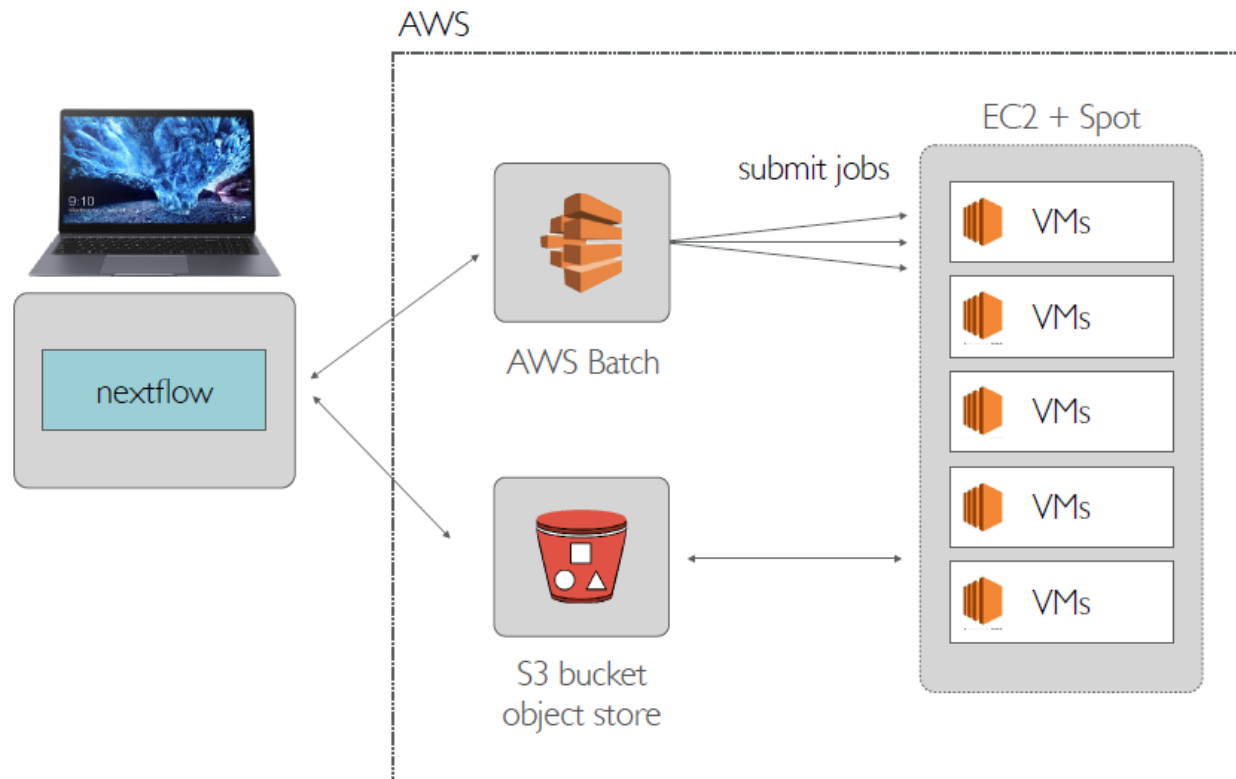- **Can scale vertically using fat server / shared mem. machine**

# Centralized cluster orchestration



- **Nextflow orchestrates workflow execution submitting jobs to a compute scheduler.**
- **Can run in the head node or a compute node.**
- **Requires a shared storage to exchange data between tasks.**
- **Ideal for coarse-grained parallelism.**

# Cloud batch orchestration



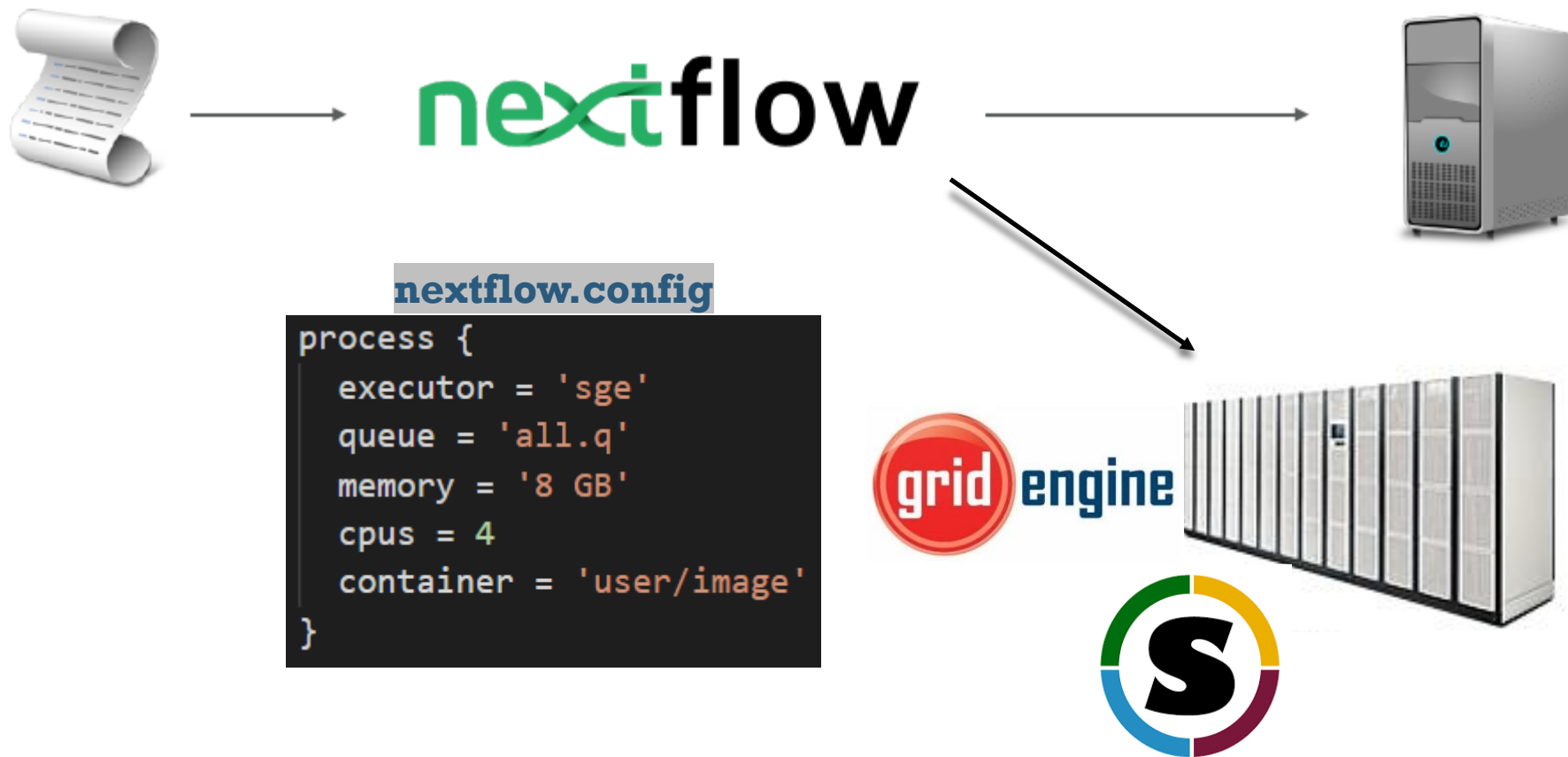- **Nextflow orchestrates workflow execution via AWS Batch**

- **Launch workflow from anywhere into the cloud**

- **Transfer of data between local environment and cloud storage**

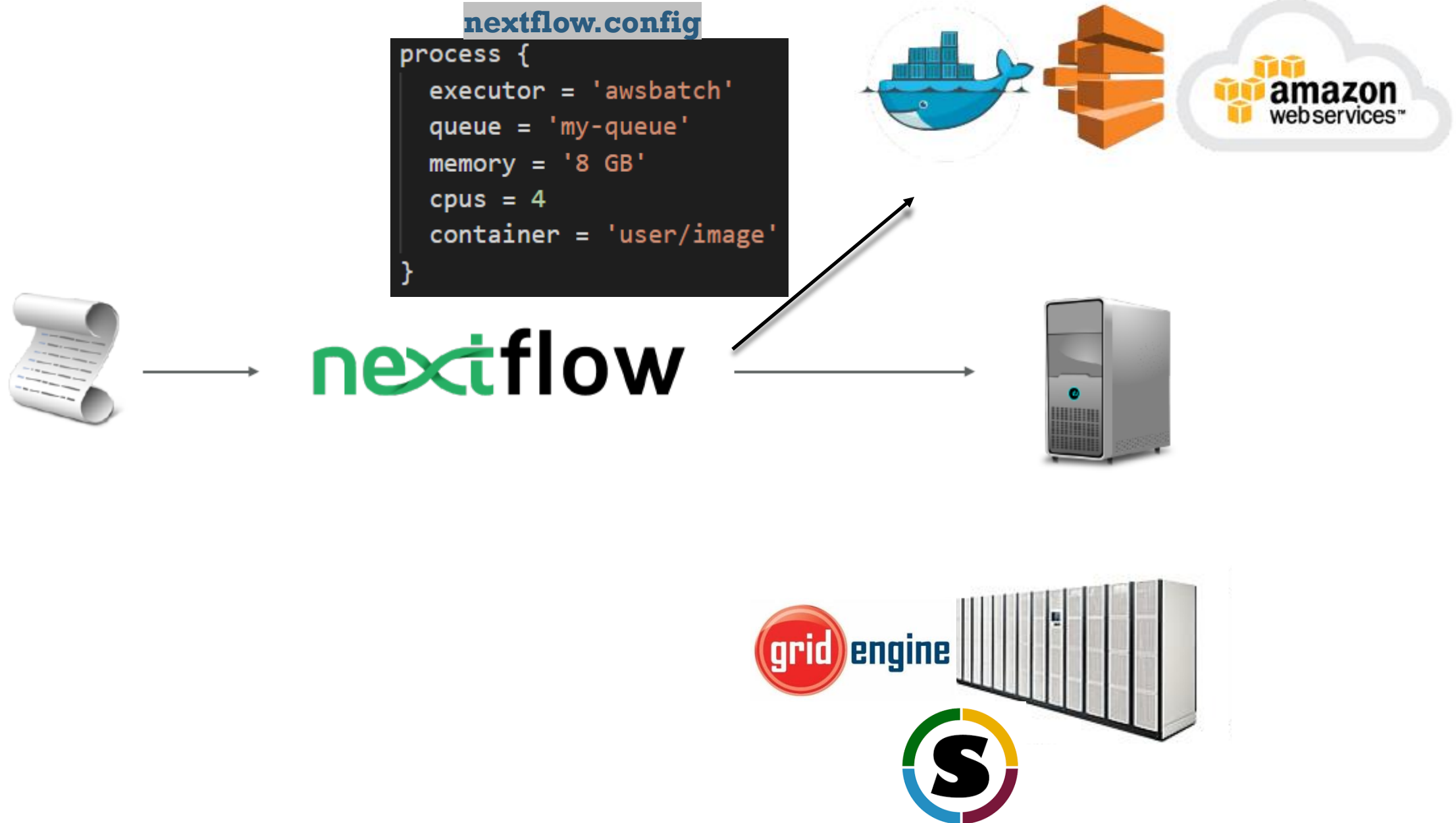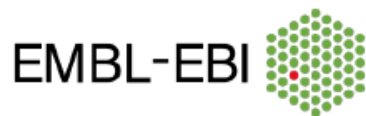- **Requires a shared object storage to exchange data between VMs.**

local

# Portability



**nextflow.config**

```
process {
  executor = 'sge'
  queue = 'all.q'
  memory = '8 GB'
  cpus = 4
  container = 'user/image'
}
```

# Portability

nextflow.config

```
process {
  executor = 'awsbatch'
  queue = 'my-queue'
  memory = '8 GB'
  cpus = 4
  container = 'user/image'
}
```

# Configuration Decoupling
## is the key to
## portable deployments

# Widespread enterprise adoption

# Container vs. VM



- Lighter: MB vs GB

- Faster startup: ms/secs vs minutes

- Virtualize a process/application instead of OS/Hardware

- Immutable: don't change over time, thus guarantee replicability over executions.

- Composable: the output of one container is directly consumable as input by another container.

- Transparent: they are created with a well-defined automated procedure.

**Nextflow Supports:**

Conda environments

Docker containers

Shifter Containers

Singularity containers

Podman containers

Charliecloud containers

When use containers?
**Always!**

# Caching and Checkpointing

A key feature of Nextflow is **re-entrancy** which is the ability to restart a pipeline after an error from the last successful process.

```
nextflow run word_count.nf --input 'data/untrimmed_fastq/*.fastq.gz' -resume
```

```
N E X T F L O W  ~  version 21.04.3
Launching `word_count.nf` [chaotic_knuth] - revision: 65cc76d410
[dd/c19734] process > NUM_LINES (5) [100%] 6 of 6, cached: 6 ✓
SRR2589044_2.fastq.gz4428360
SRR2584863_1.fastq.gz6213036
SRR2589044_1.fastq.gz4428360
SRR2584866_1.fastq.gz11073592
SRR2584866_2.fastq.gz11073592
SRR2584863_2.fastq.gz6213036
```

# Reporting

# Reporting

## Tasks

This table shows information about each task in the workflow. Use the search box on the right to filter rows for specific values. Clicking headers will sort the table by that value and scrolling side to side will reveal more columns.

Values shown as: [ Human readable ▾ ]
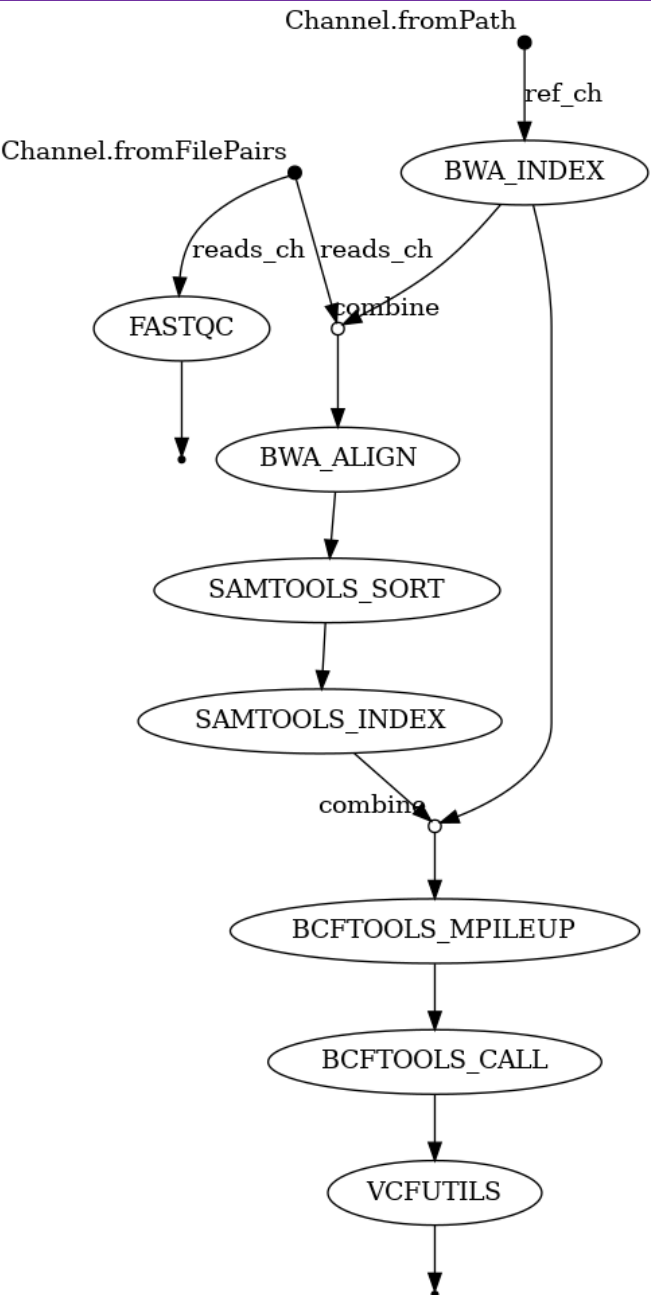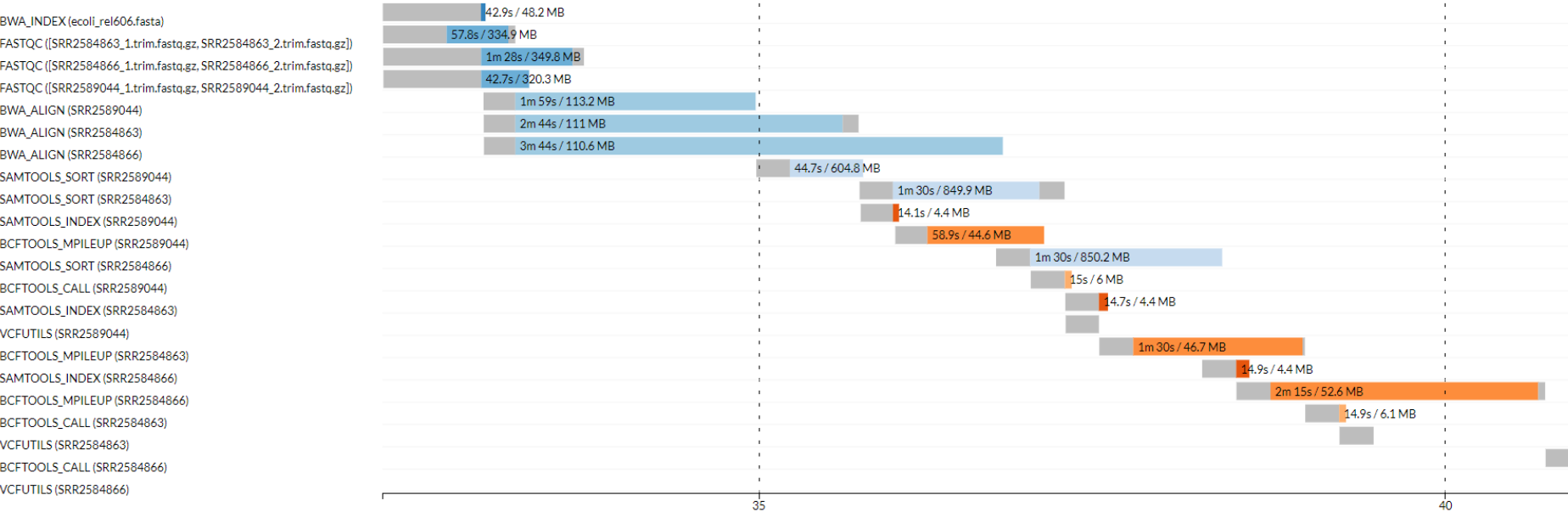
Show [ 25 ▾ ] entries

Filter: [ Metrics ] [ Metadata ] [ All ]  Search: [          ]

| task_id ↑↓ | process ↑↓ | tag ↑↓ | status ↑↓ | allocated cpus ↑↓ | %cpu ↑↓ | allocated memory ↑↓ | %mem ↑↓ | vmem ↑↓ | rss ↑↓ | pea |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | BWA_INDEX | ecoli_rel606.fasta | COMPLETED | 1 | 96.6 | - | 0.0 | 31.230 MB | 6.836 MB | 69.0 |
| 2 | FASTQC | [SRR2584863_1.trim.fastq.gz, SRR2584863_2.trim.fastq.gz] | COMPLETED | 1 | 120.6 | - | 0.1 | 3.313 GB | 327.680 MB | 3.33 |

## Processes execution timeline

Launch time: 16 Jan 2022 19:32
Elapsed time: 9m 2s
Legend: job wall time / memory usage (RAM)

BWA_INDEX (ecoli_rel606.fasta) — 42.9s / 48.2 MB
FASTQC ([SRR2584863_1.trim.fastq.gz, SRR2584863_2.trim.fastq.gz]) — 57.8s / 334.9 MB
FASTQC ([SRR2584866_1.trim.fastq.gz, SRR2584866_2.trim.fastq.gz]) — 1m 28s / 349.8 MB
FASTQC ([SRR2589044_1.trim.fastq.gz, SRR2589044_2.trim.fastq.gz]) — 42.7s / 320.3 MB
BWA_ALIGN (SRR2589044) — 1m 59s / 113.2 MB
BWA_ALIGN (SRR2584863) — 2m 44s / 111 MB
BWA_ALIGN (SRR2584866) — 3m 44s / 110.6 MB
SAMTOOLS_SORT (SRR2589044) — 44.7s / 604.8 MB
SAMTOOLS_SORT (SRR2584863) — 1m 30s / 849.9 MB
SAMTOOLS_INDEX (SRR2589044) — 14.1s / 4.4 MB
BCFTOOLS_MPILEUP (SRR2589044) — 58.9s / 44.6 MB
SAMTOOLS_SORT (SRR2584866) — 1m 30s / 850.2 MB
BCFTOOLS_CALL (SRR2589044) — 15s / 6 MB
SAMTOOLS_INDEX (SRR2584863) — 14.7s / 4.4 MB
VCFUTILS (SRR2589044)
BCFTOOLS_MPILEUP (SRR2584863) — 1m 30s / 46.7 MB
SAMTOOLS_INDEX (SRR2584866) — 14.9s / 4.4 MB
BCFTOOLS_MPILEUP (SRR2584866) — 2m 15s / 52.6 MB
BCFTOOLS_CALL (SRR2584863) — 14.9s / 6.1 MB
VCFUTILS (SRR2584863)
BCFTOOLS_CALL (SRR2584866)
VCFUTILS (SRR2584866)

35    40

Channel.fromPath
ref_ch
Channel.fromFilePairs
BWA_INDEX
reads_ch  reads_ch
combine
FASTQC
BWA_ALIGN
SAMTOOLS_SORT
SAMTOOLS_INDEX
combine
BCFTOOLS_MPILEUP
BCFTOOLS_CALL
VCFUTILS

# Open-Source Community



55K+ monthly downloads

10,000+ active developers /month

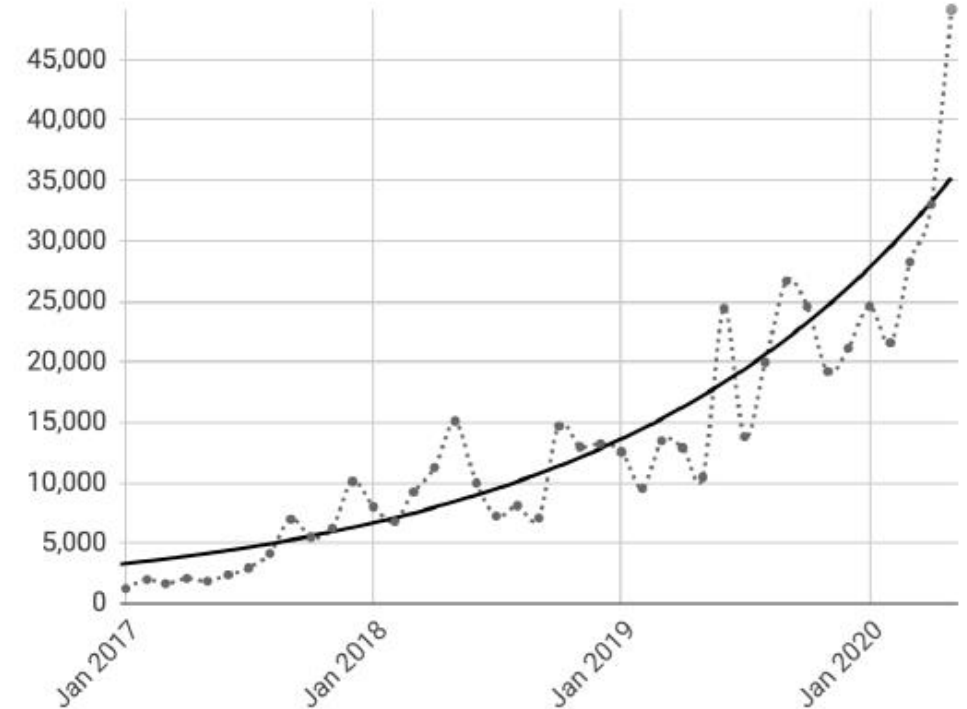35 international Workshops

150K+ lines of code

1.6K+ stars on github

100+ contributors

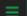monthly downloads

**Over 2 million downloads**

nf-core

# Growing community content
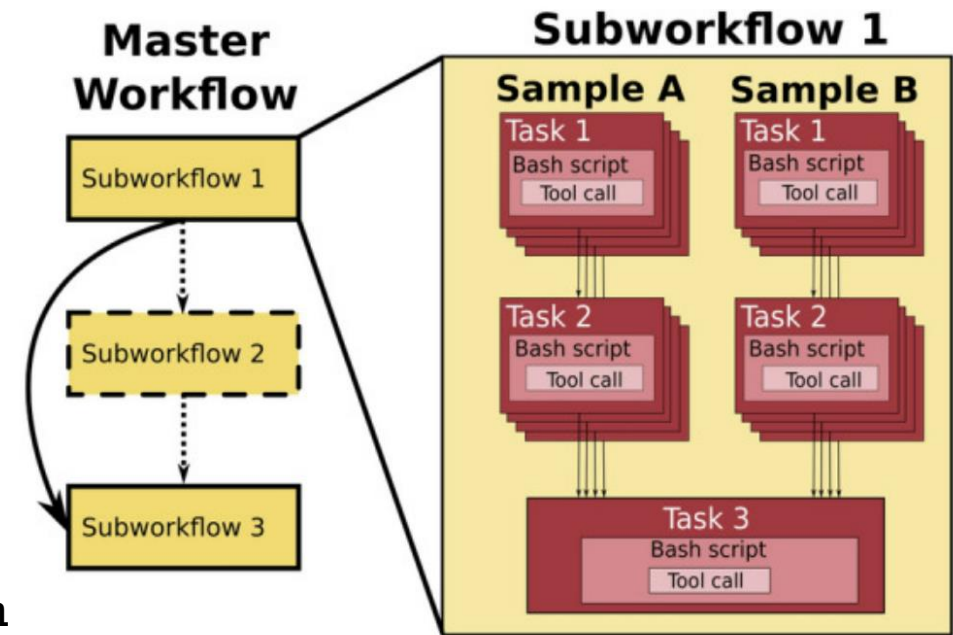
Production ready analysis pipelines built with Nextflow.

# DSL2 Modules & Sub-Workflows = Re-Usability!

- Modularity is a very important design principle for production bioinformatics workflows.
- The core idea is to build a library of reusable modules (tasks or sub-workflows) and assemble them into various master workflows.

This enables -

(1) performing different analyses without having to refactor the entire workflow

(2) check-pointing and restart of a workflow run from a task in the middle of analysis if needed

(3) customizing runtime environments and compute resources which may vary between analysis stages.



```
=========================================================================
    Include Modules
=========================================================================
*/

include { FASTQC }                              from "./modules/fastqc" addParams(OUTPUT: "${params.outdir}/fastqc")
include { BWA_INDEX }                           from "./modules/bwa_index" addParams(OUTPUT: "${params.outdir}/bwa_index")
include { BWA_ALIGN }                           from "./modules/bwa_align" addParams(OUTPUT: "${params.outdir}/bwa_align")
include { SAMTOOLS_SORT; SAMTOOLS_INDEX }       from "./modules/samtools" addParams(OUTPUT: "${params.outdir}/sorted_bam")
include { BCFTOOLS_MPILEUP; BCFTOOLS_CALL; VCFUTILS } from "./modules/bcftools" addParams(OUTPUT: "${params.outdir}/vcf")
```