



Nextflow Workshop

November 2021



Schedule

Day	Time	Topics
Day 1 – Nov 16, 2021	11:30 – 2:00 ET	Introduction to Nextflow Getting started with Nextflow A simple RNA-Seq pipeline Managing dependencies & containers
Day 2 – Nov 17, 2021	11:30 - 1:30 ET	Nextflow Channels Nextflow Processes Nextflow Operators Groovy basic structures and idioms DSL2 and modules
Day 3 – Nov 18, 2021	11:30 - 1:30 ET	Nextflow configuration Deployment scenarios Execution cache and resume Error handling & troubleshooting Getting started with Nextflow Tower



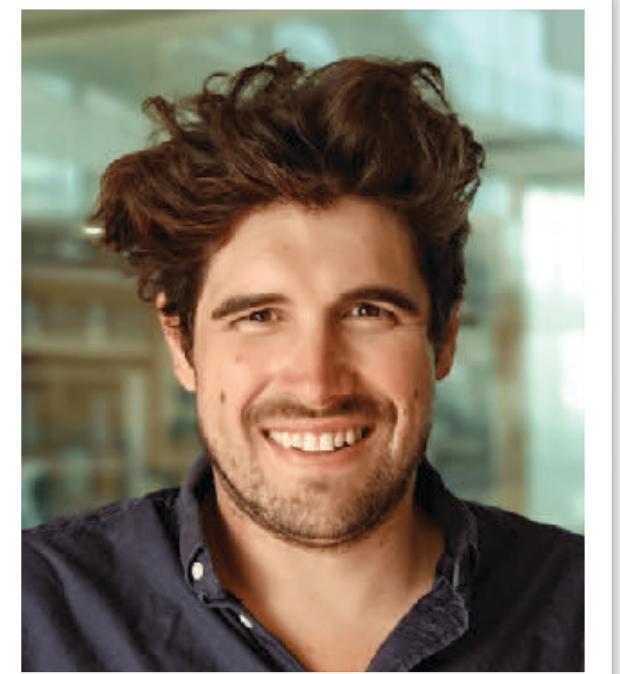
Seamless data pipelines at scale.

Big data analysis can get complicated.

Cluster & cloud configuration is difficult.

Enable your team to build, deploy and manage pipelines collaboratively on any infrastructure.

Build better pipelines, reduce costs and improve time to results.

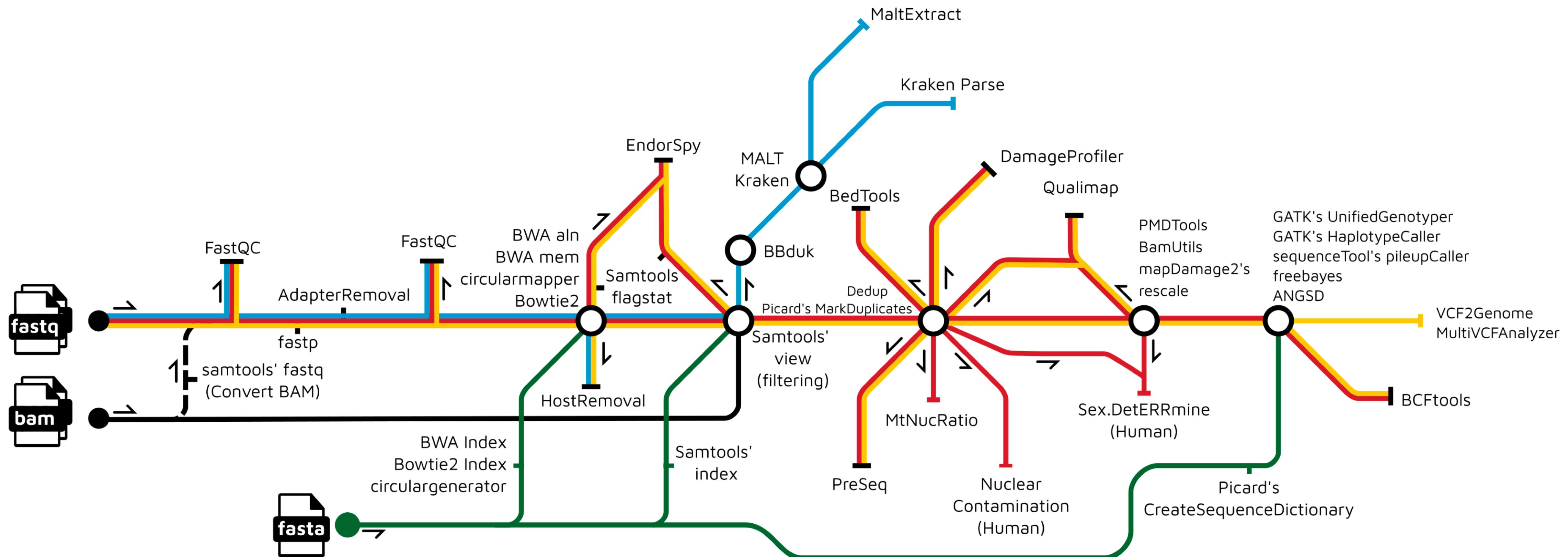


Evan Floden

CEO & Co-founder Seqera Labs

evan@seqera.io

Writing modern workflows is complex



nf-core/eager v2.4

Example analysis pathways

Legend

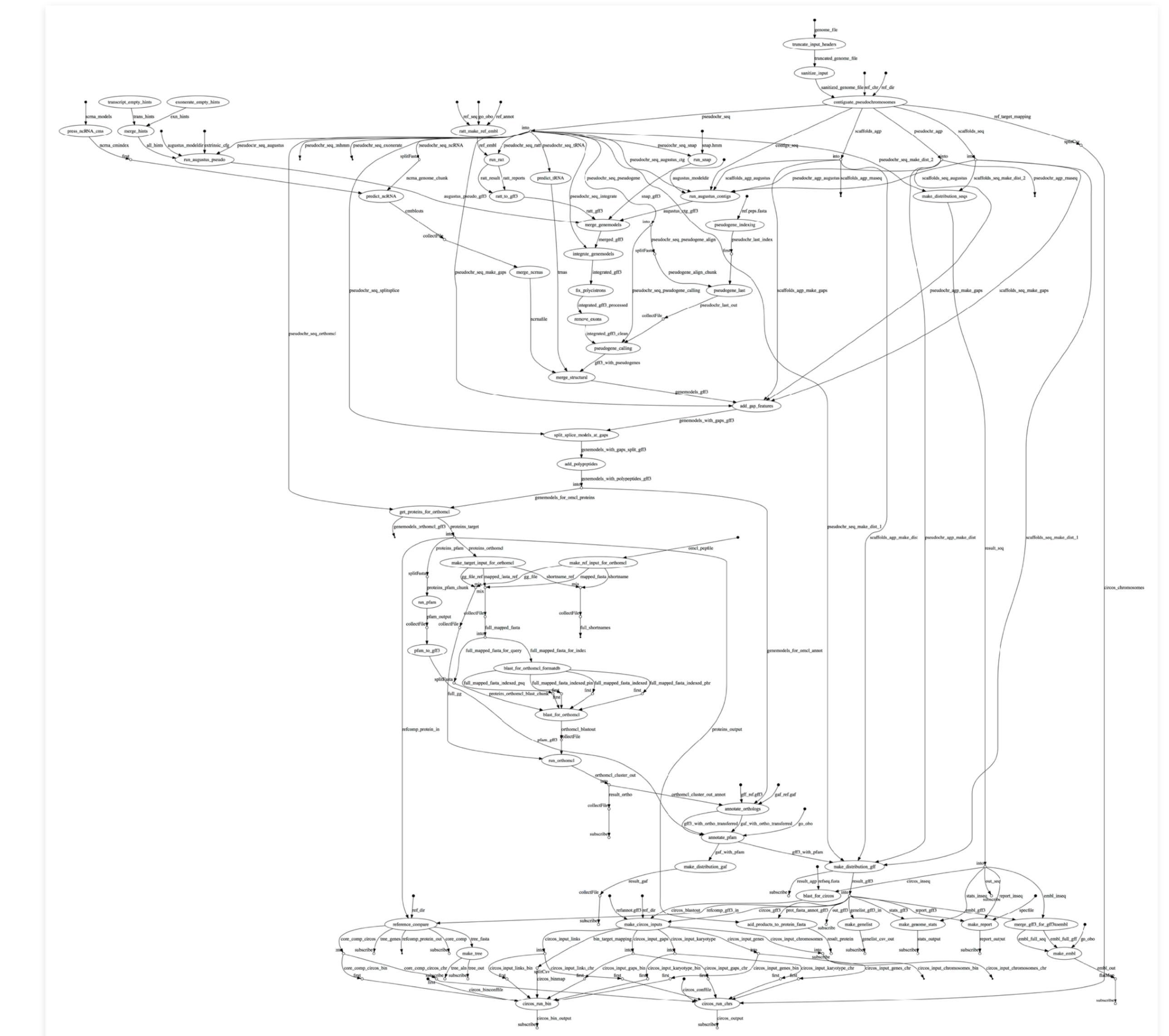
- Start
- End Step
- Single output Intermediate Step
- Multi-output intermediate step
- Eukaryotic Nuclear Genome
- Microbial or Organelle Genome
- Metagenomic Profiling
- Reference

Data analysis workflows

- Data analysis applications perform computation to generate information from (large) datasets
- Embarrassingly parallelisation → can spawn 100s-100k jobs over distributed cluster
- Mash-up of many different tools and scripts (dependancies!)
- Complex dependency trees and configuration → very fragile ecosystem

A lot of moving parts!

- 70 processes
- 55 external scripts
- 39 software tools & libraries



Steinbiss et al., Companion parasite genome annotation pipeline, DOI: 10.1093/nar/gkw292

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



OPEN  ACCESS Freely available online

 PLOS ONE

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

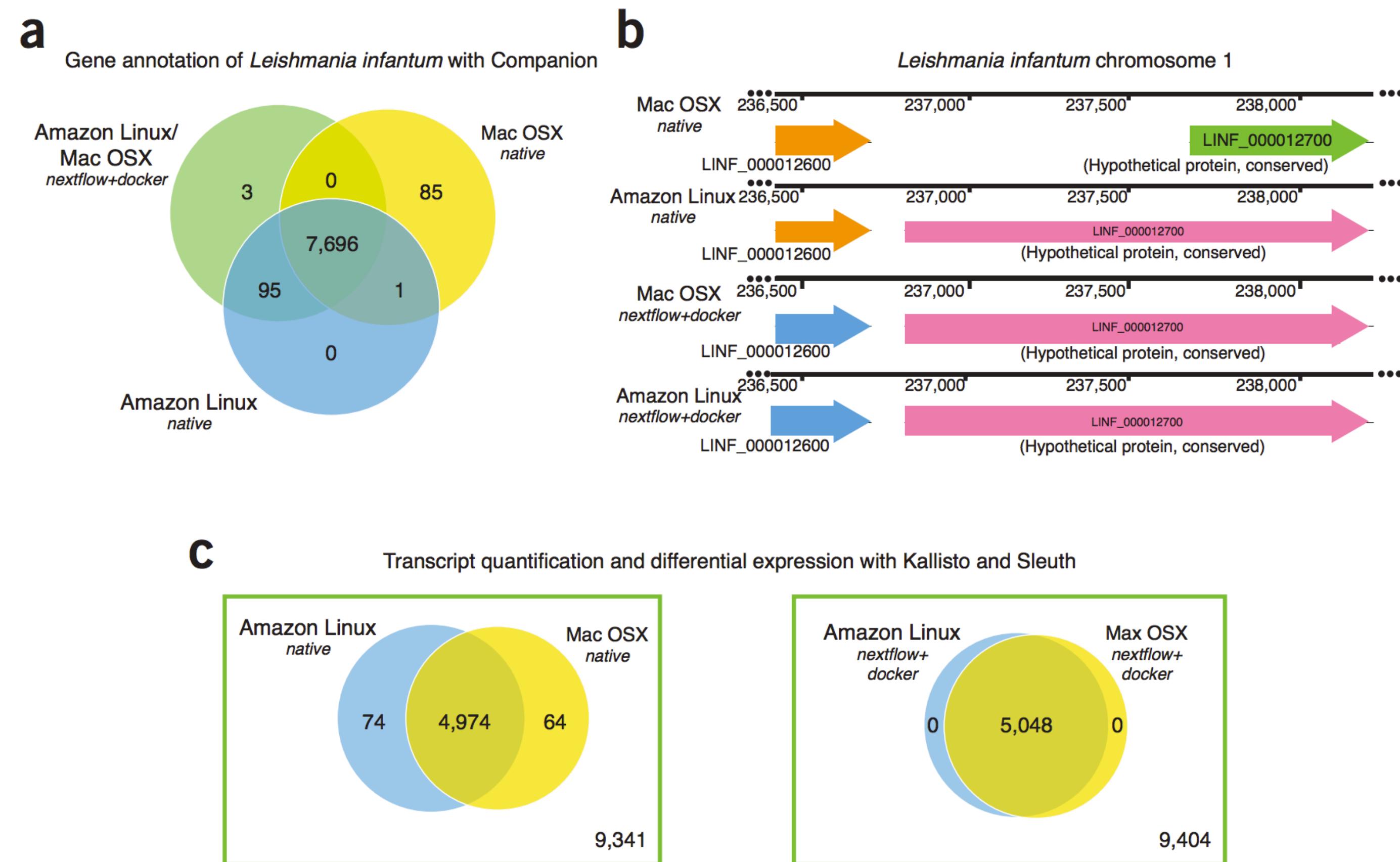
Daniel Garijo¹, Sarah Kinnings², Li Xie³, Lei Xie⁴, Yinliang Zhang⁵, Philip E. Bourne^{3*}, Yolanda Gil^{6*}

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, Los Angeles, California, United States of America

≈1.7 months!

The **same application**
deployed in
different environments
produces
different results

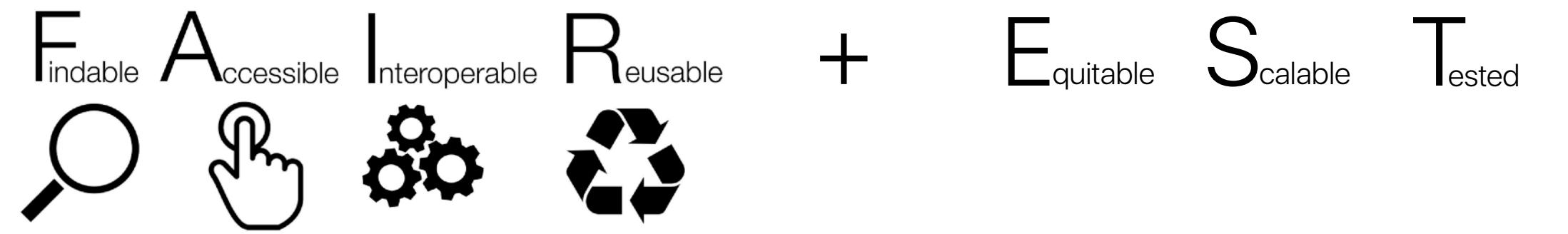
Differences in genome annotations & expressed genes



VOLUME 35 NUMBER 4 APRIL 2017 **NATURE BIOTECHNOLOGY**

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

What is our goal?



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



What is Nextflow?

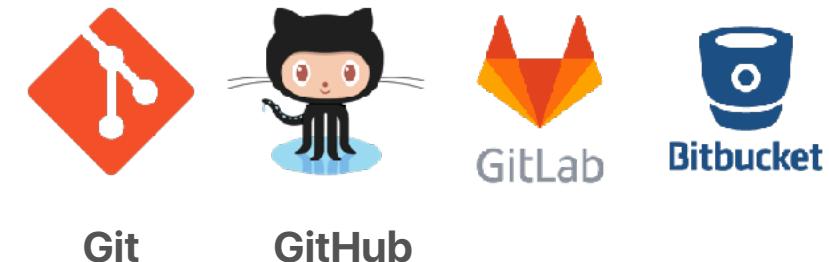
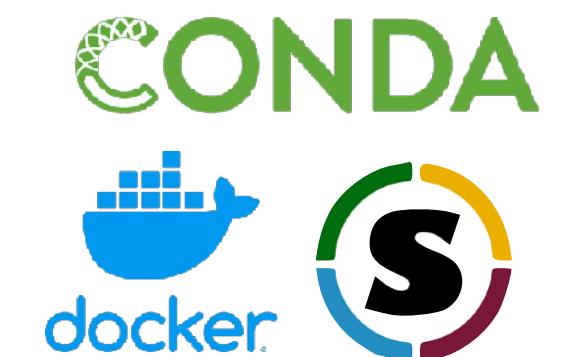
nextflow script

Write code
in any language



Define orchestration with
dataflow programming

Define software
dependencies
with containers



Version control

nextflow runtime

Orchestration of tasks to
deploy anywhere with ease

Supports all major platforms



How does it work?

- **Fast prototyping** ⇒ custom DSL that enables task composition, simplifies most use cases + general purpose programming language for corner cases
- **Easy parallelisation** ⇒ declarative reactive programming model based on dataflow paradigm, implicit portable parallelism
- **Self-contained** ⇒ functional approach, a task execution is idempotent ie. 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  
        """  
}  

```

Tasks 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  
        """  
}
```

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  
        """  
}
```

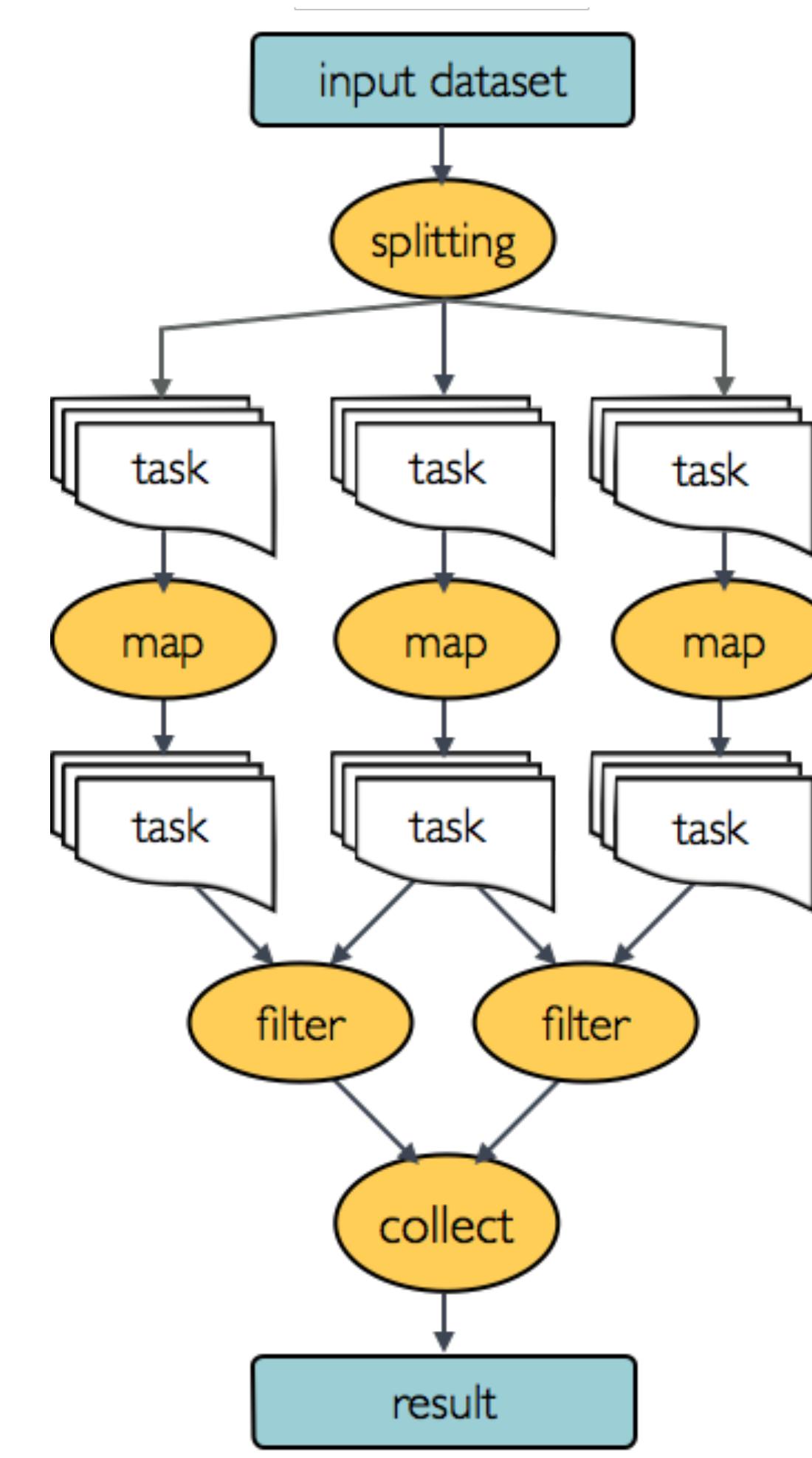
workflow

```
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  
}
```

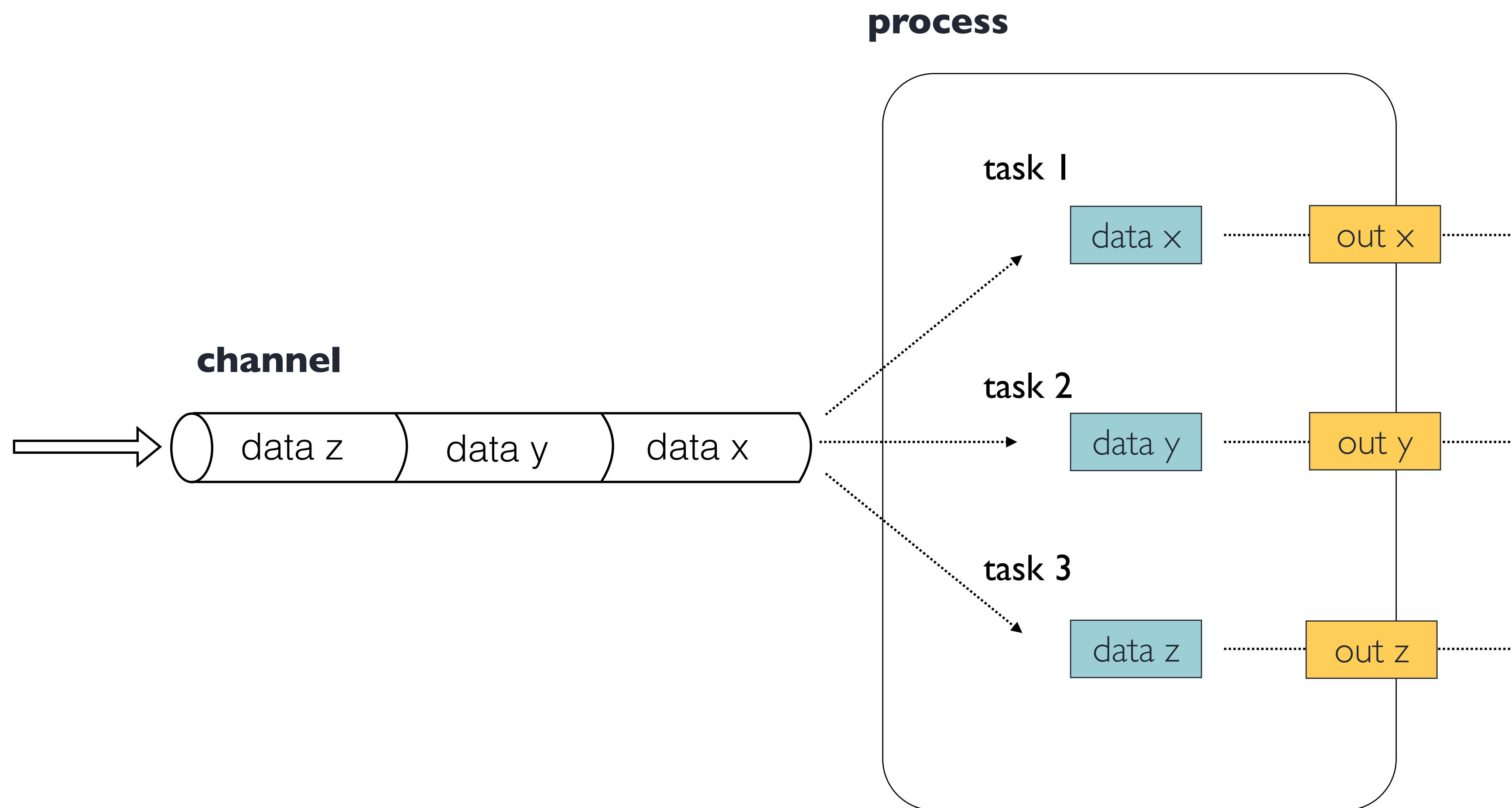
Scientist and engineers can now write complex, distributed and parallel data pipelines without requiring a degree in computer science.

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



How does parallelisation work?



How does parallelisation 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}
    """
}

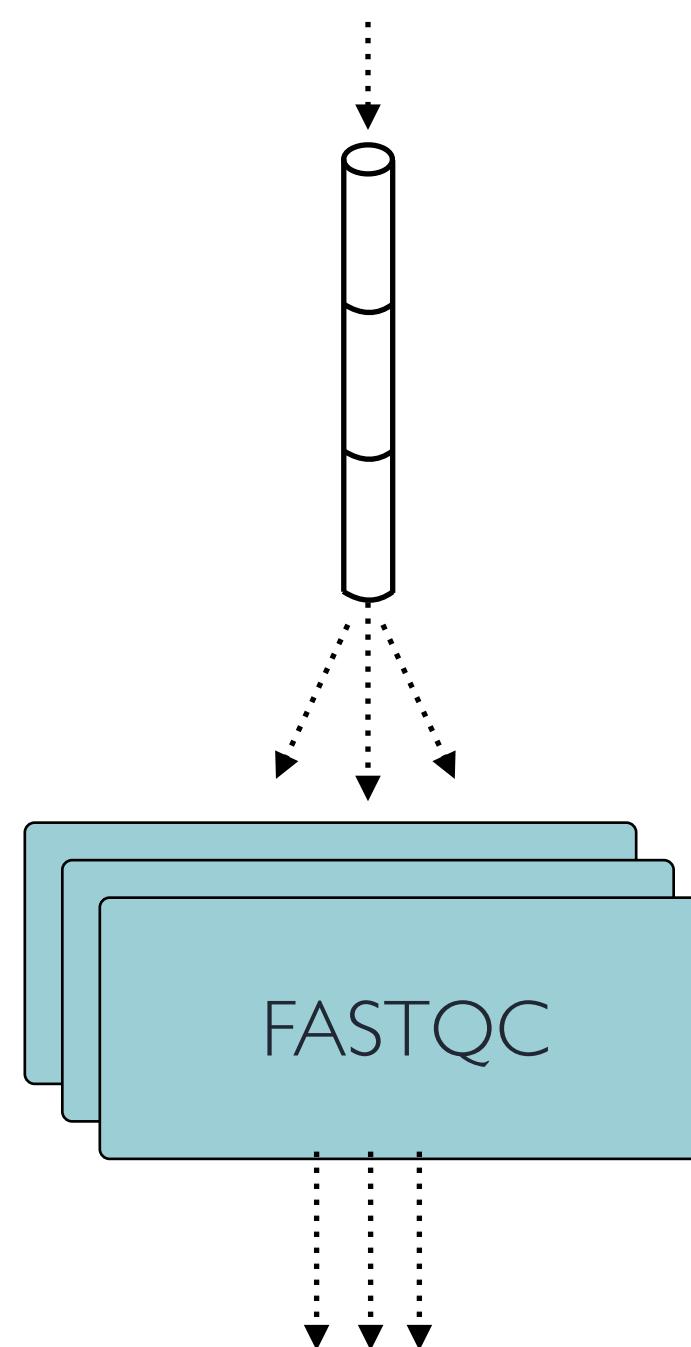
}
```

How does parallelisation work?

```
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}  
        """  
}  
  
}
```

Implicit parallelism

```
Channel.fromPath("data/*.fastq")
```



CWL vs. Nextflow

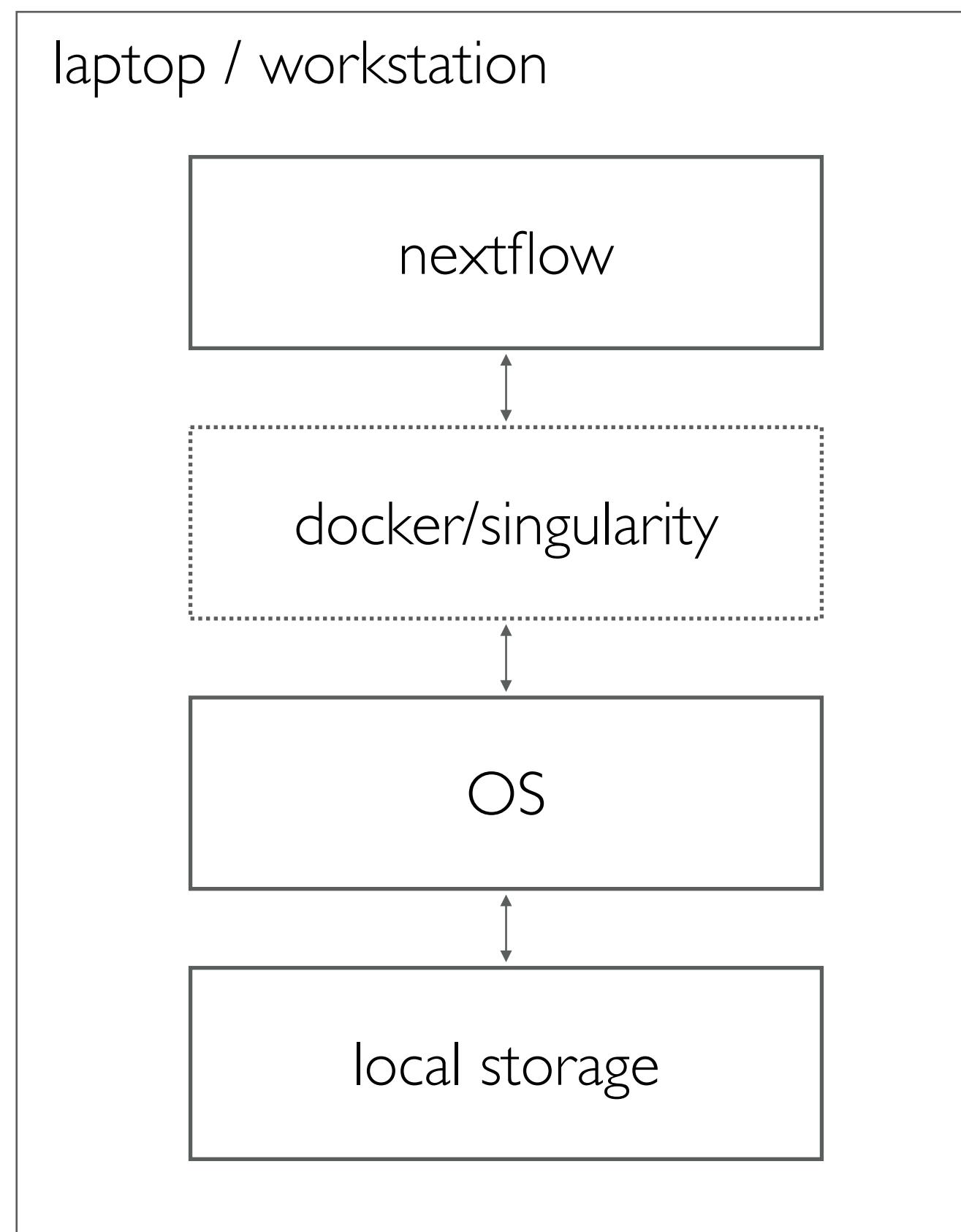
- Language specification
- Declarative meta-language (YAML/JSON)
- Verbose
- Committee driven
- Many vendors/implementations (and specification version)
- Language + app. runtime
- DSL on top of a general purpose programming lang.
- Concise, fluent
- Community driven
- Single implementation, quick iterations

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

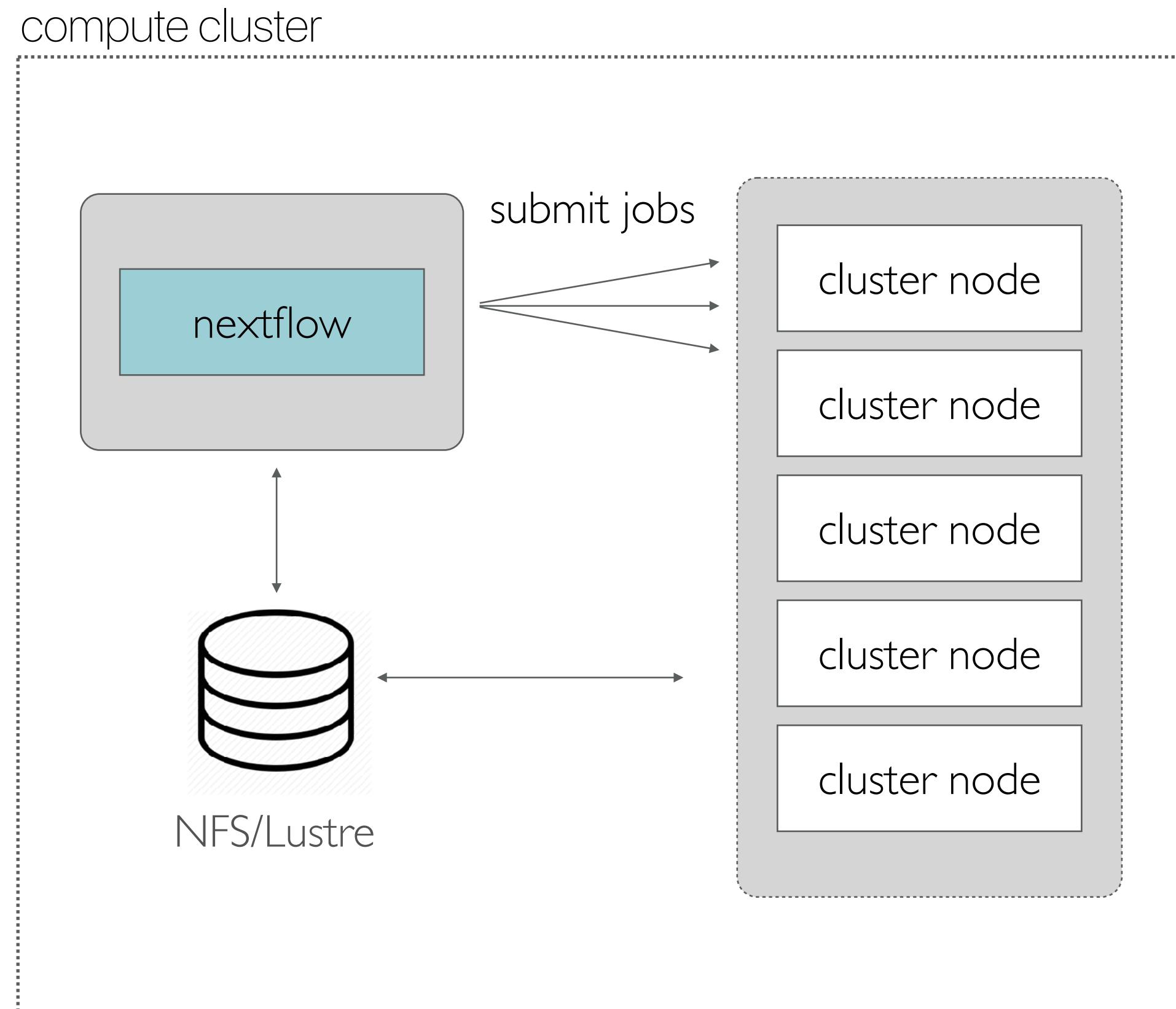
Deployment Scenarios

Local execution



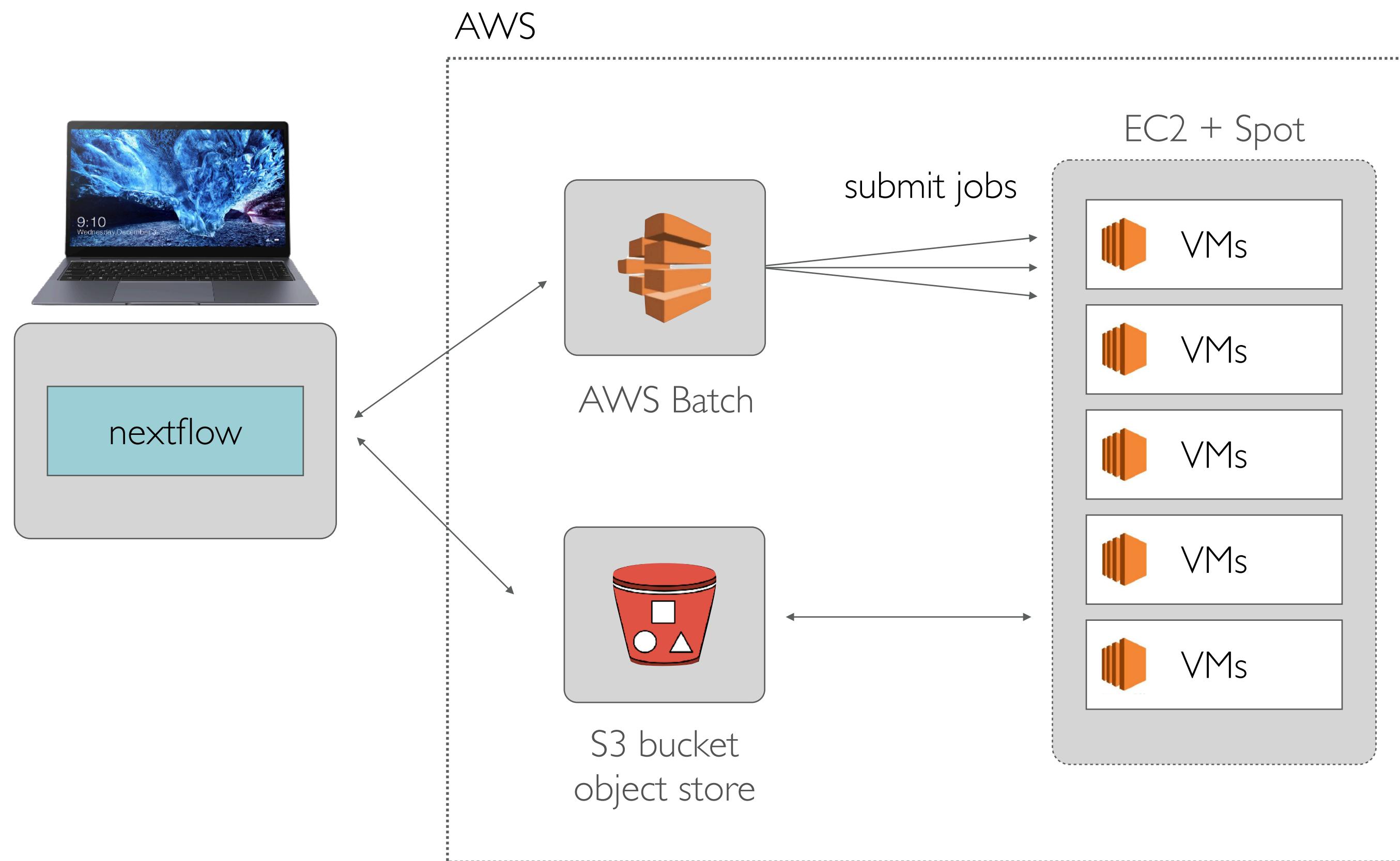
- Common development scenario
- Dependencies can be managed using a container runtime
- Parallelisations is managed spawning posix processes
- Can scale vertically using fat server / shared mem. machine

Centralised cluster orchestration



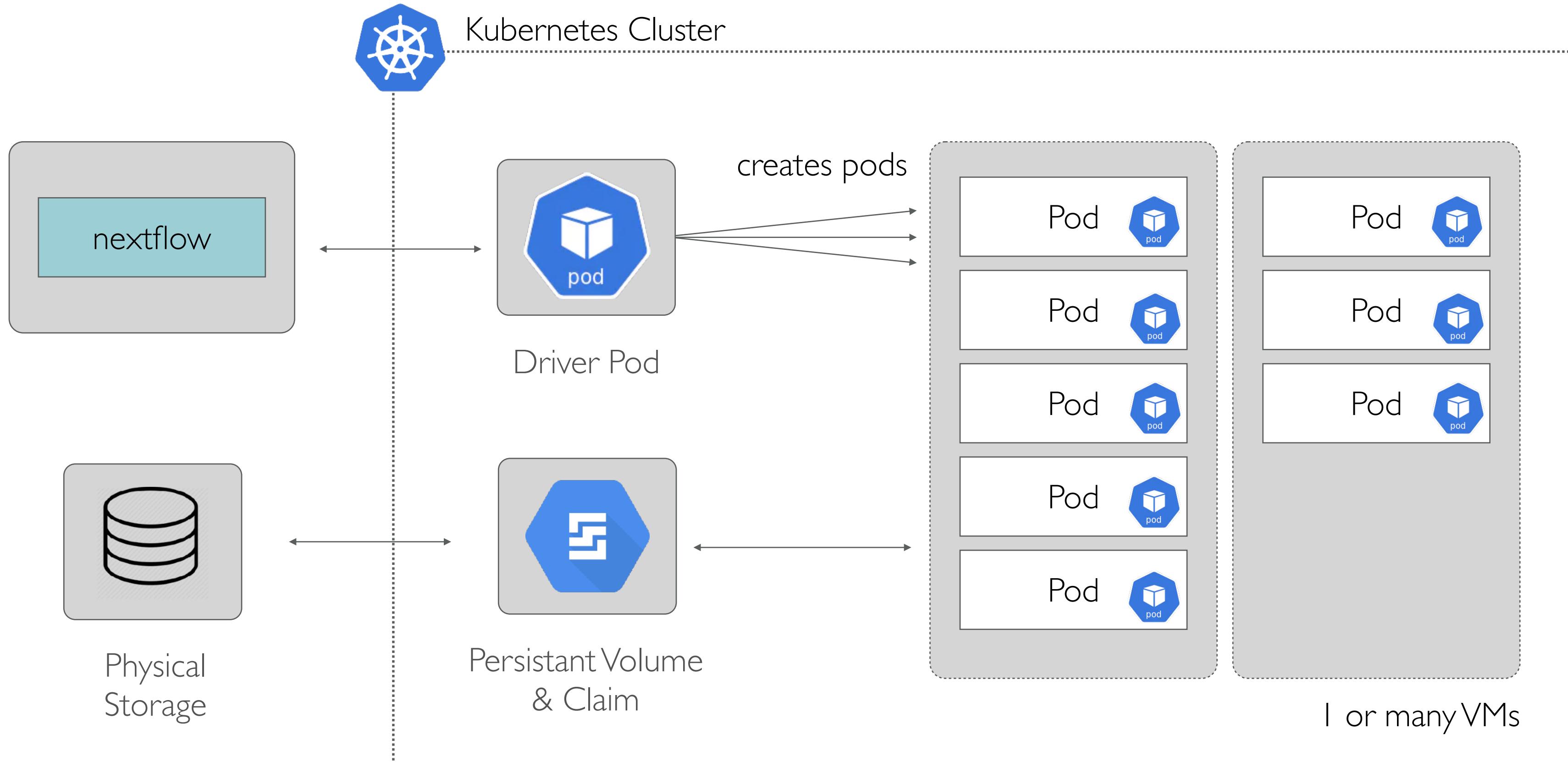
- Nextflow orchestrates workflow execution submitting jobs to a compute scheduler.
- Can run in the head node or a compute node.
- Requires a shared storage to exchange data between tasks.
- Ideal for coarse-grained parallelism.

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.

Cloud native orchestration



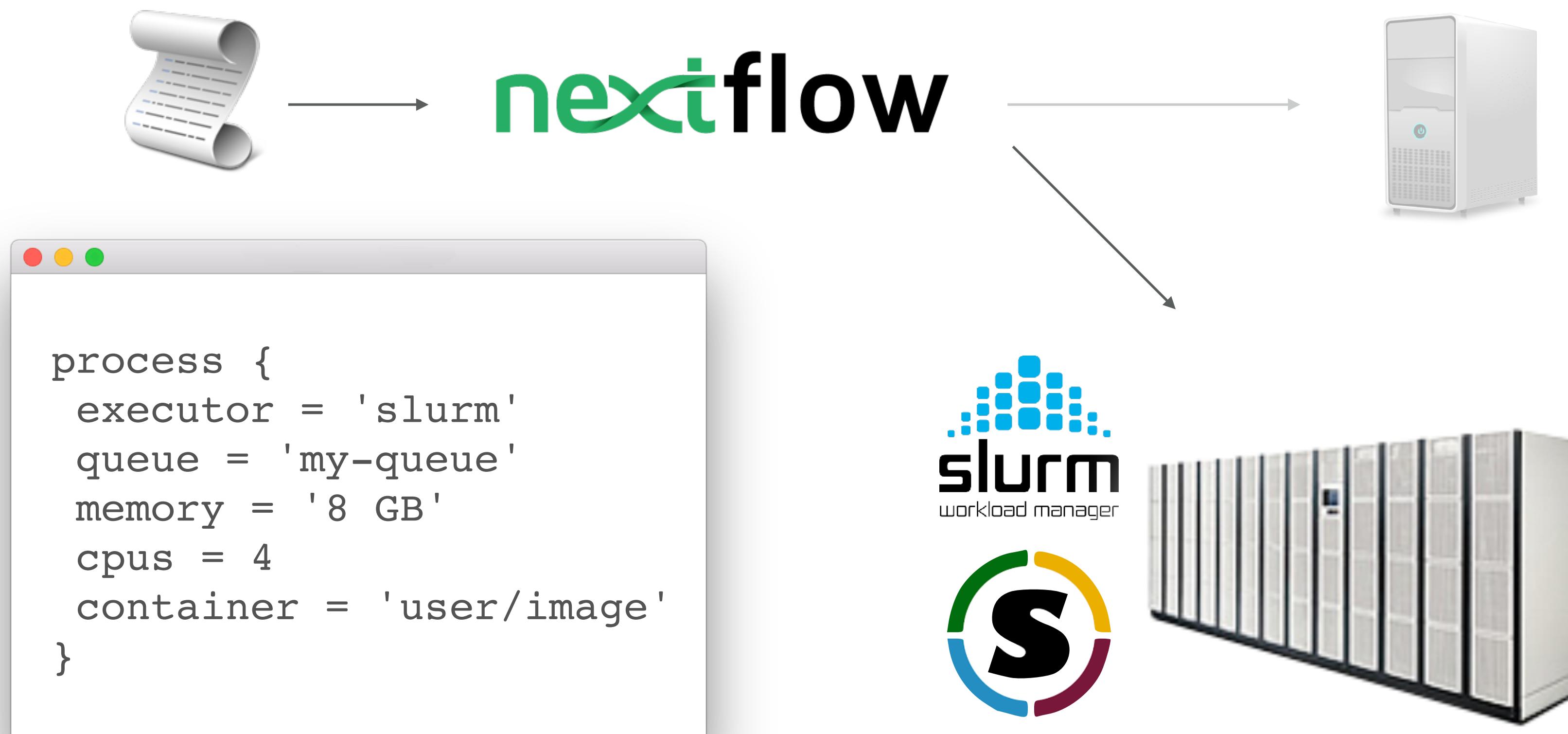
- Nextflow orchestrates workflow execution
- Each task is executed as a kubernetes pod.
- Transfer of data between tasks via a persistent volume claim.



Portability



Portability



Portability



→ **nextflow**

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



**Configuration
decoupling**
is the key to
**portable
deployments**

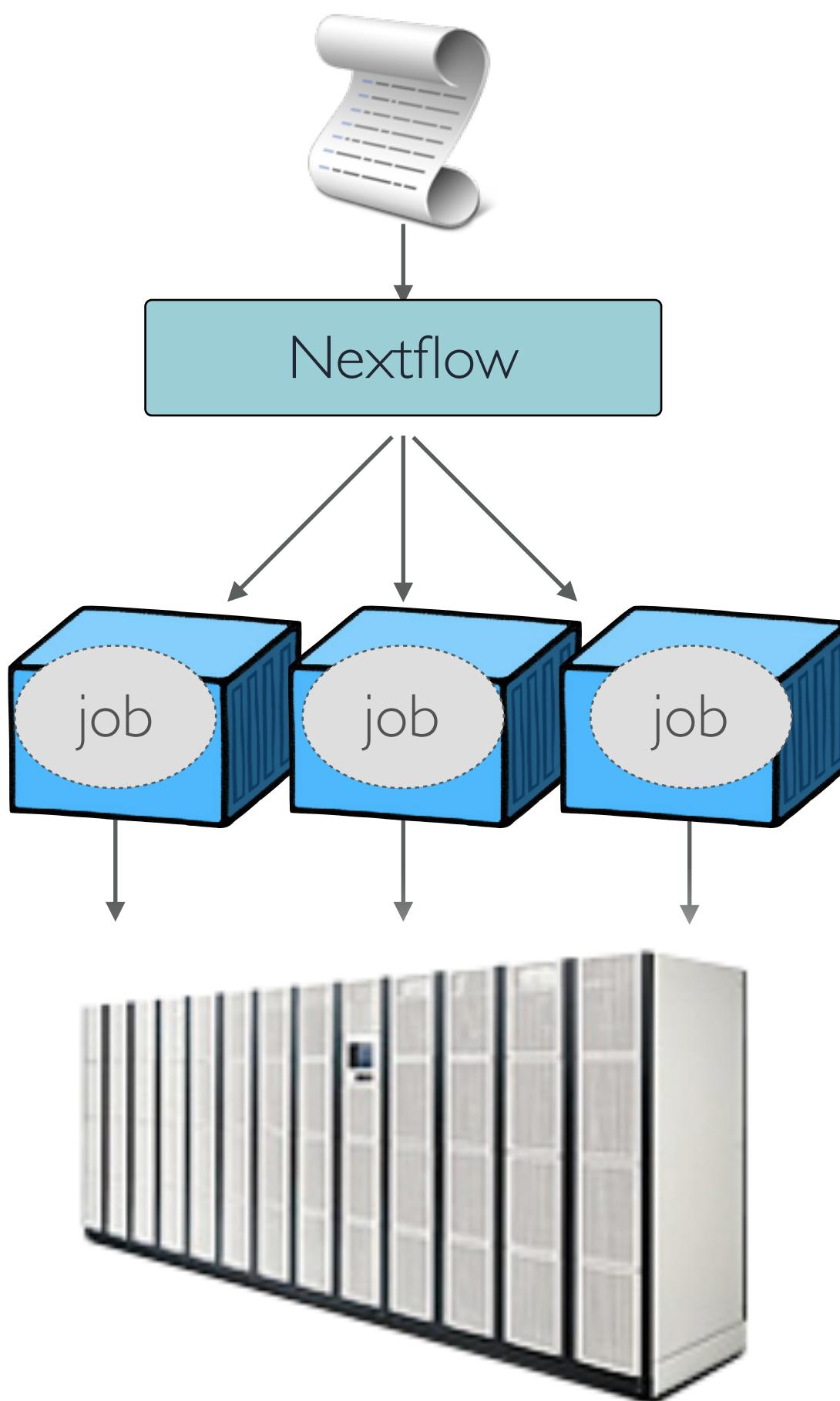
An aerial photograph of a bustling port. In the foreground, a large white cargo ship is docked at a quay, its deck covered with a variety of shipping containers in red, blue, green, and white. A massive yellow gantry crane stands prominently over the ship, its long arm extended towards the container deck. The background is filled with a dense grid of shipping containers stacked high in several yards, stretching into the distance. Other smaller ships and port infrastructure are visible in the far background.

Containerisation

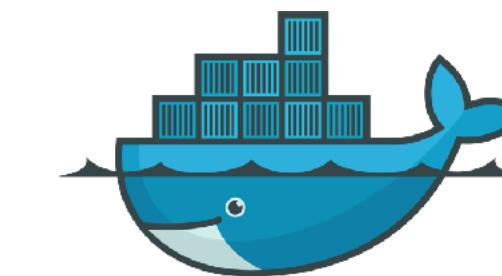
Container vs. VM

- Lighter: MB vs GB
- Faster startup: ms/sec vs minutes
- Virtualise a process/application instead of a 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.

Containerisation



- Nextflow envisioned the use of software containers to fix computational reproducibility
- Mar 2014 (ver 0.7), support for Docker
- Dec 2016 (ver 0.23), support for Singularity
- 2020 support for Podman



When use containers?

Always!

Open Source Community

55K+

monthly
downloads

150K+

lines of
code

10,000+

active developers
/month

1.6K+

stars on
github

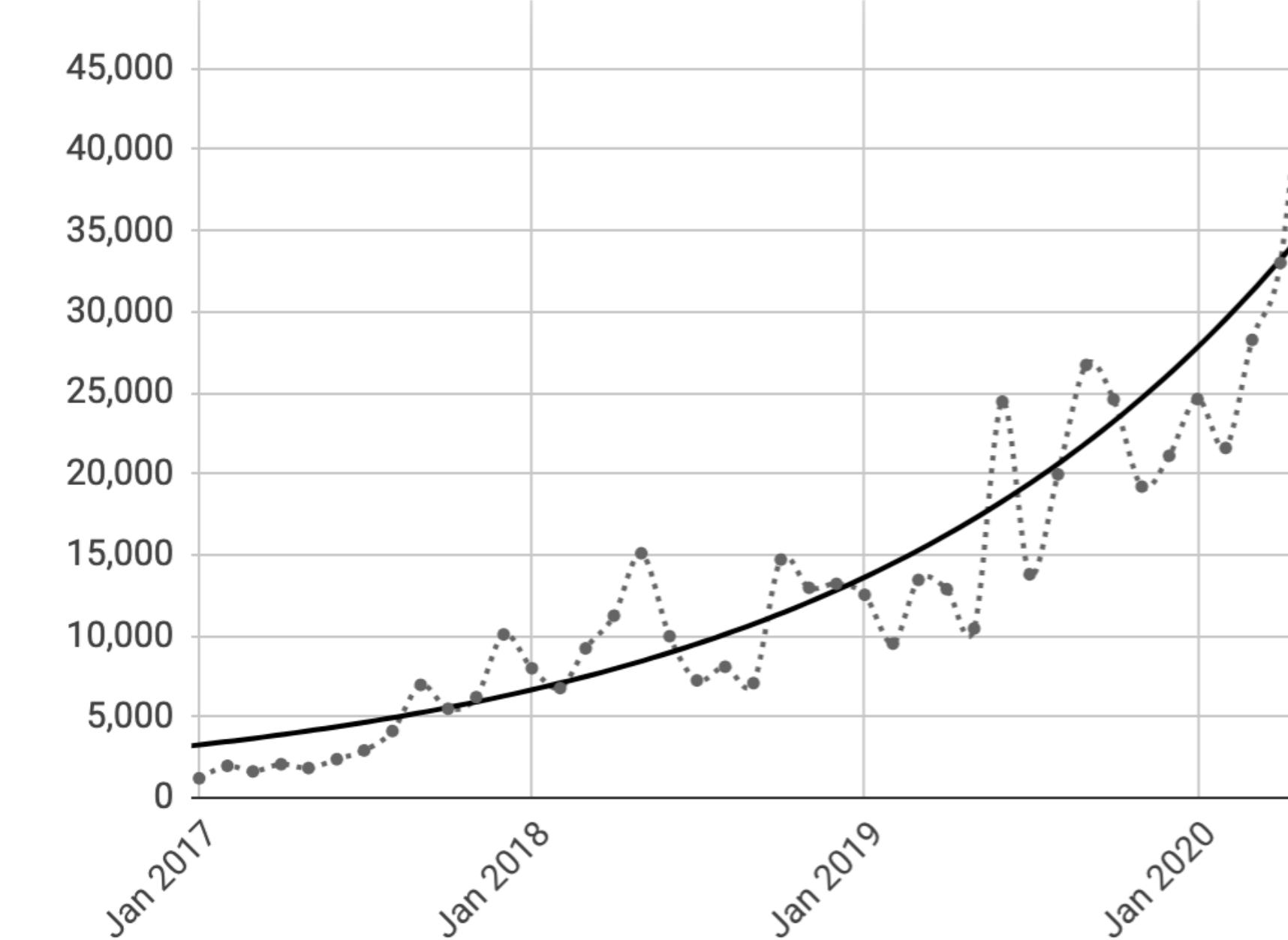
35

international
Workshops

100+

contributors

monthly downloads



Over 2 million downloads

Editors



main.nf — ~/projects/callings-nf

```
main.nf
```

```
48 log.info """\n49 CALLINGS - NF v 1.0\n50 ======\n51 genome : $params.genome\n52 reads : $params.reads\n53 variants : $params.variants\n54 blacklist: $params.blacklist\n55 results : $params.results\n56 gatk : $params.gatk\n57 """\n58\n59 /*\n60 * Parse the input parameters\n61 */\n62\n63 GATK = params.gatk_launch\n64 genome_file = file(params.genome)\n65 variants_file = file(params.variants)\n66 blacklist_file = file(params.blacklist)\n67 reads_ch = Channel.fromFilePairs(params.reads)\n68\n69\n70 //*****\n71 * PART 1: Data preparation\n72 *\n73 * Process 1A: Create a FASTA genome index (.fai) with samtools for GATK\n74 */\n75 process '1A_prepare_genome_samtools' {\n76 tag "$genome.baseName"\n77\n78 input:\n79     file genome from genome_file\n80\n81 output:\n82     file "${genome}.fai" into genome_index_ch\n83\n84 script:\n85 """\n86     samtools faidx ${genome}\n87 """\n88 }\n89\n90 /*\n91 * Process 1B: Create a FASTA genome sequence dictionary with Picard for GATK\n92 */\n93 process '1B_prepare_genome_picard' {\n94 tag "$genome.baseName"\n95\n96 input:\n97     file genome from genome_file\n98 output:\n99     file "${genome.baseName}.dict" into genome_dict_ch\n100\n101 script:\n102 """\n103     PICARD=`which picard.jar`\n104     java -jar \$PICARD CreateSequenceDictionary R= $genome O= ${genome.baseName}.dict\n105 """\n106 }
```

main.nf* 93:1 LF UTF-8 Nextflow master 4 files

main.nf

```
pcpu.sh dd.sh data main.nf
```

```
48 log.info """\n49 CALLINGS - NF v 1.0\n50 ======\n51 genome : $params.genome\n52 reads : $params.reads\n53 variants : $params.variants\n54 blacklist: $params.blacklist\n55 results : $params.results\n56 gatk : $params.gatk\n57 """\n58\n59 /*\n60 * Parse the input parameters\n61 */\n62\n63 GATK = params.gatk_launch\n64 genome_file = file(params.genome)\n65 variants_file = file(params.variants)\n66 blacklist_file = file(params.blacklist)\n67 reads_ch = Channel.fromFilePairs(params.reads)\n68\n69\n70 //*****\n71 * PART 1: Data preparation\n72 *\n73 * Process 1A: Create a FASTA genome index (.fai) with samtools for GATK\n74 */\n75 process '1A_prepare_genome_samtools' {\n76 tag "$genome.baseName"\n77\n78 input:\n79     file genome from genome_file\n80\n81 output:\n82     file "${genome}.fai" into genome_index_ch\n83\n84 script:\n85 """\n86     samtools faidx ${genome}\n87 """\n88 }\n89\n90\n91 /*\n92 * Process 1B: Create a FASTA genome sequence dictionary with Picard for GATK\n93 */\n94\n95 process '1B_prepare_genome_picard' {\n96 tag "$genome.baseName"\n97\n98 input:\n99     file genome from genome_file\n100 output:\n101     file "${genome.baseName}.dict" into genome_dict_ch\n102\n103 script:\n104 """\n105     PICARD=`which picard.jar`\n106     java -jar \$PICARD CreateSequenceDictionary R= $genome O= ${genome.baseName}.dict\n107 """\n108 }\n109\n110
```

Ln 1, Col 1 Spaces: 4 UTF-8 LF Nextflow

Widespread enterprise adoption



Growing community content

Production ready analysis
pipelines built with Nextflow.

Available Pipelines

Can you think of another pipeline that would fit in well? [Let us know!](#)

Search keywords Filter: Released 28 Under development 17 Archived 4 Sort: Last Release Alphabetical Stars Display:

Pipeline Name	Version	Last Published	Stars	Description
nf-core/diaproteomics ✓	1.2.2	2 days ago	5	Automated quantitative analysis of DIA proteomics mass spectrometry measurements.
nf-core/dualrnaseq ✓	1.0.0	2 weeks ago	3	Analysis of Dual RNA-seq data - an experimental method for interrogating host-pathogen interactions through simultaneous RNA-seq.
nf-core/mag ✓	1.2.0	2 weeks ago	42	Assembly and binning of metagenomes
nf-core/ampliseq ✓	1.2.0	3 weeks ago	54	16S rRNA amplicon sequencing analysis workflow using QIIME2
nf-core/sarek ✓	2.7	1 month ago	91	Analysis pipeline to detect germline or somatic variants (pre-processing, variant calling and annotation) from WGS / targeted sequencing
nf-core/eager ✓	2.3.1	Published 1 month ago	45	A fully reproducible and state-of-the-art ancient DNA analysis pipeline

Chan
Zuckerberg
Initiative



A community effort to collect a curated set of analysis pipelines built using Nextflow.

VIEW PIPELINES

Deploy

- Stable pipelines
- Centralized configs
- List and update pipelines
- Download for offline use

Participate

- Documentation
- Slack workspace
- Twitter updates
- Hackathons

Develop

- Starter template
- Code guidelines
- CI code linting and tests
- Helper tools

nextflow

data pipelines at scale

massively scalable pipelines across cluster & cloud.

The world's leading workflow software
for genomics, biopharma and life sciences.



features.



Powerful HPC execution engines

Deploy across cloud & clusters effortlessly

Out-of-the-box support for AWS and GCP plus schedulers including SLURM, LSF & Grid Engine.



Portable & reproducible

Containers without the hassle

Docker and Singularity integrations encapsulate all pipeline dependencies across environments.



Language agnostic

Write pipelines in your language

The flexibility to develop, port and reuse code allows you to do things your way.

nextflow tower

delivering discovery.

Manage, optimize and launch data analysis pipelines from a secure command-post.



Real time monitoring, anywhere

[Keep up-to-date with progress.](#)

Track workloads and retrieve a full history of executions across mobile and desktop browsers.



Launch workflows

[Deploy on local, cluster, hybrid, cloud & serverless.](#)

Execution of workflows via the GUI launchpad or API and submit to user-defined environments.



Team & organization management

[Collaborate and share securely.](#)

User management via enterprise providers such as LDAP and 3rd party authentication methods.



Advanced cloud-native storage

[Reduce data transfer times.](#)

Significantly speed up workloads through the use of vendor-specific systems such as Lustre FSx.

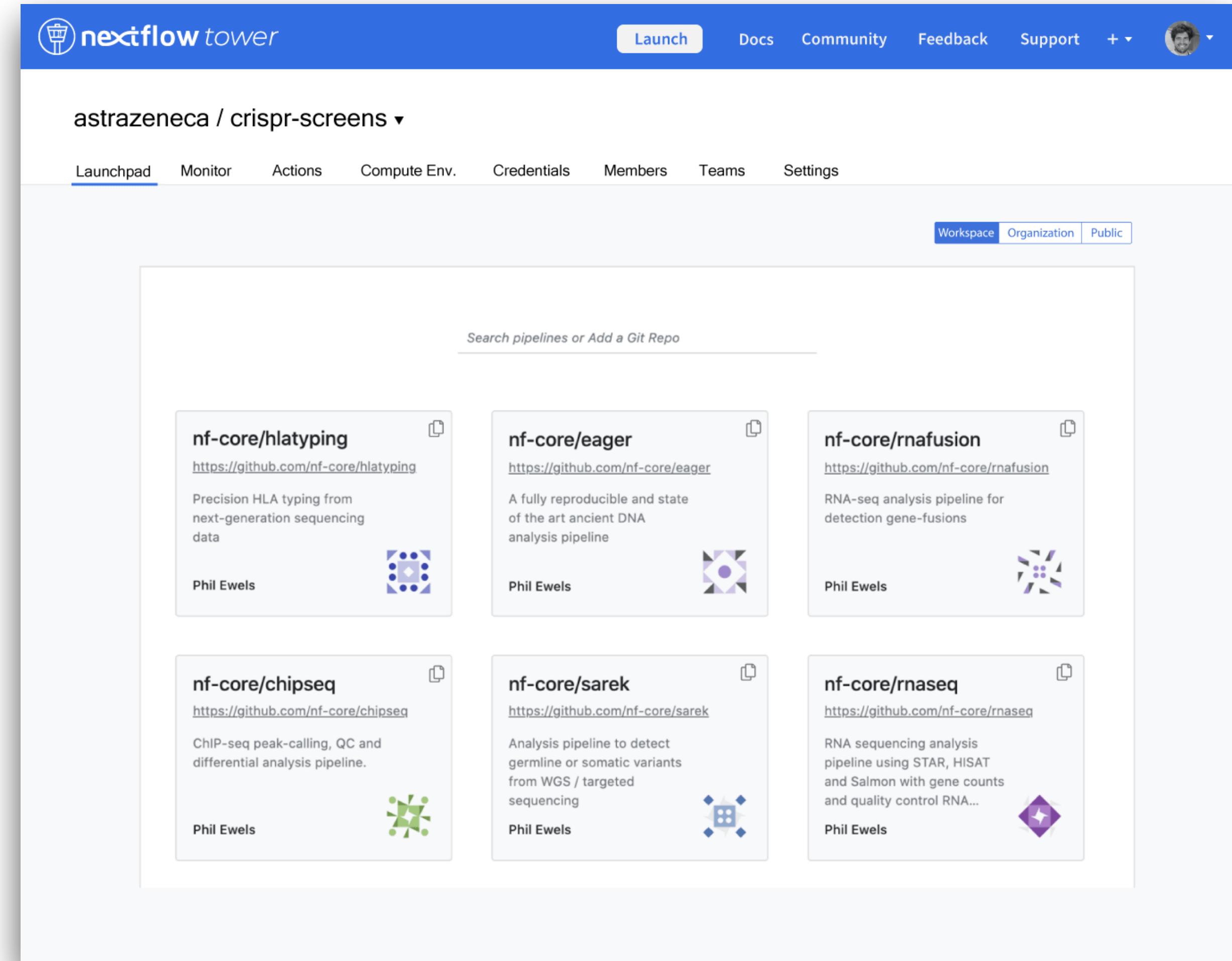


Accounting & budgeting

[Understand and control project costs.](#)

Estimates of compute costs and granular permissioning of resources across teams.

Build pipelines that run out-of-the-box. Anywhere.



Preconfigured pipelines with input data and forms based on JSON schema.

Manage cloud and on-prem environments in one place.

New Compute Environment

Name

A name of your choice to identify this compute environment

Platform

Select the platform

-  Amazon Batch
-  Google Life Sciences
-  IBM LSF
-  Slurm Workload Manager
-  Amazon EKS
-  Google GKE
-  Kubernetes

Deploy in any infrastructure with native support
for all major cloud, Kubernetes & HPC platforms.

Provide teams with the resources they need.

The screenshot shows the Nextflow Tower web application interface. At the top, there is a blue header bar with the Nextflow Tower logo, navigation links for Launch, Docs, Community, Feedback, Support, and a user profile icon. Below the header, the page title is "AstraZeneca plc" with a location pin icon and the URL "https://www.astrazeneca.com". The main content area has a navigation bar with tabs: Workspaces, **Pipelines**, Compute Env., Credentials, Members, Teams, and Settings. The Pipelines tab is selected. A search bar labeled "Search pipeline..." is present. A prominent "Create pipeline" button is located in the top right corner of the main content area. Two pipeline entries are listed:

- Hello pipe**
Repository: <https://github.com/pditommaso/hello>
Description: Some pipeline description
Last event (never) Delete
- Hello pipe**
Repository: <https://github.com/pditommaso/hello>
Description: Some pipeline description
Workspace: crispr-screens
Last event (never) Delete

Role-based access for resource management across teams and organizations,
plus security, compliance and cost control best practices.

Powerful API for integrations & automations

Nextflow Tower API 1.0.0

Nextflow Tower service API

Email: info@seqera.io
URL: https://seqera.io

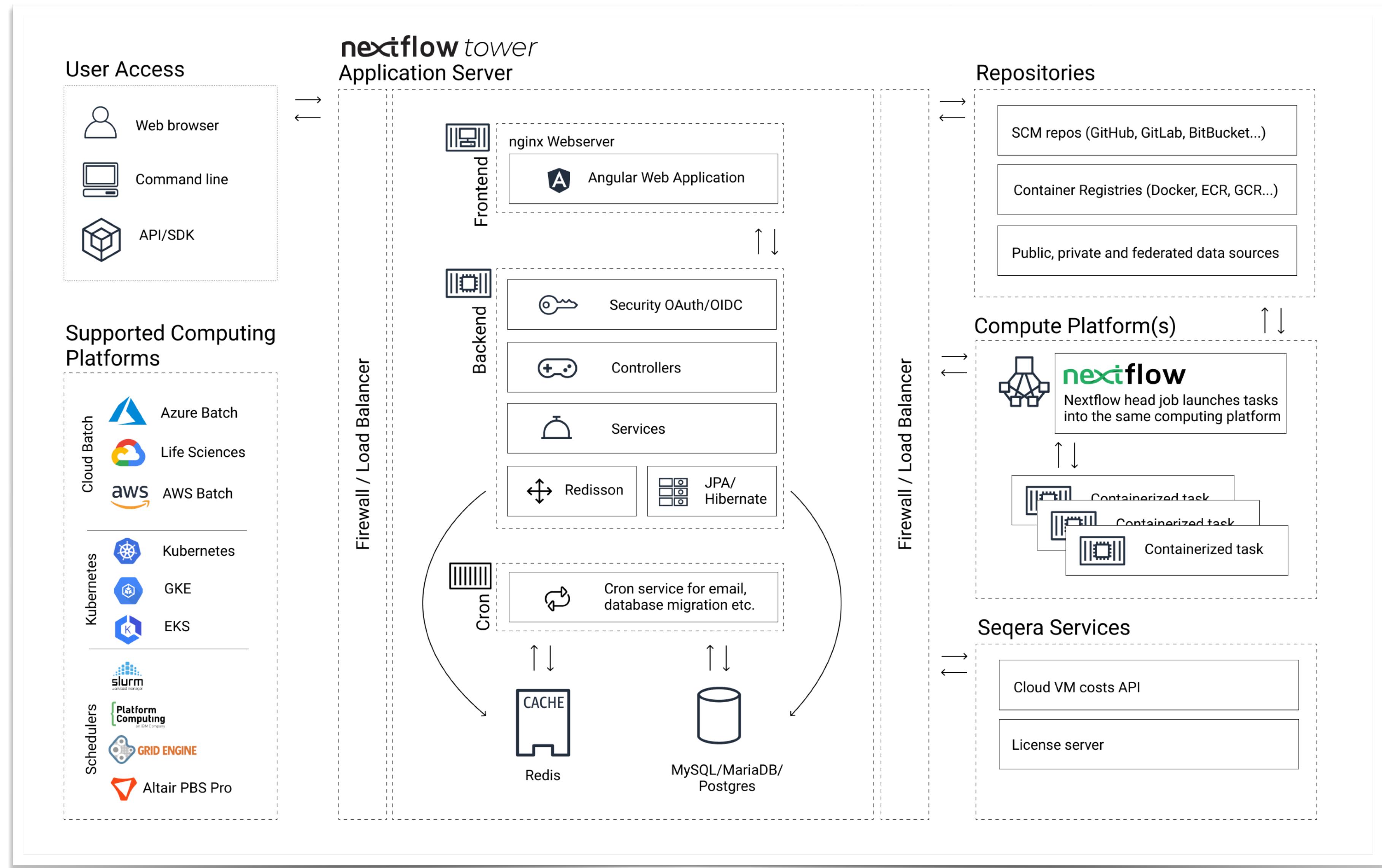
↑ actions

<code>GET /actions</code>	List the available Pipeline actions for the authenticated user
<code>GET /actions/types</code>	List the supported event types that can trigger a pipeline action
<code>GET /actions/{actionId}</code>	Describe an existing pipeline action
<code>PUT /actions/{actionId}</code>	Update a pipeline action
<code>POST /actions</code>	Create a new pipeline action
<code>POST /actions/{actionId}/pause</code>	Toggle the pause status of an existing pipeline action
<code>POST /actions/{actionId}/launch</code>	Trigger the execution of a Tower Launch action
<code>DELETE /actions/{actionId}</code>	Delete a pipeline action

collaborators

<code>GET /collaborators/workflow/{workflowId}</code>	List the collaborators of the workflow for the authenticated user
<code>POST /collaborators</code>	Add a collaborator
<code>DELETE /collaborators/{collaboratorId}/workflow/{workflowId}</code>	Delete a collaborator

Full stack solution for data pipelines at scale



Let's get started with the workshop!

<https://www.seqera.io/training>