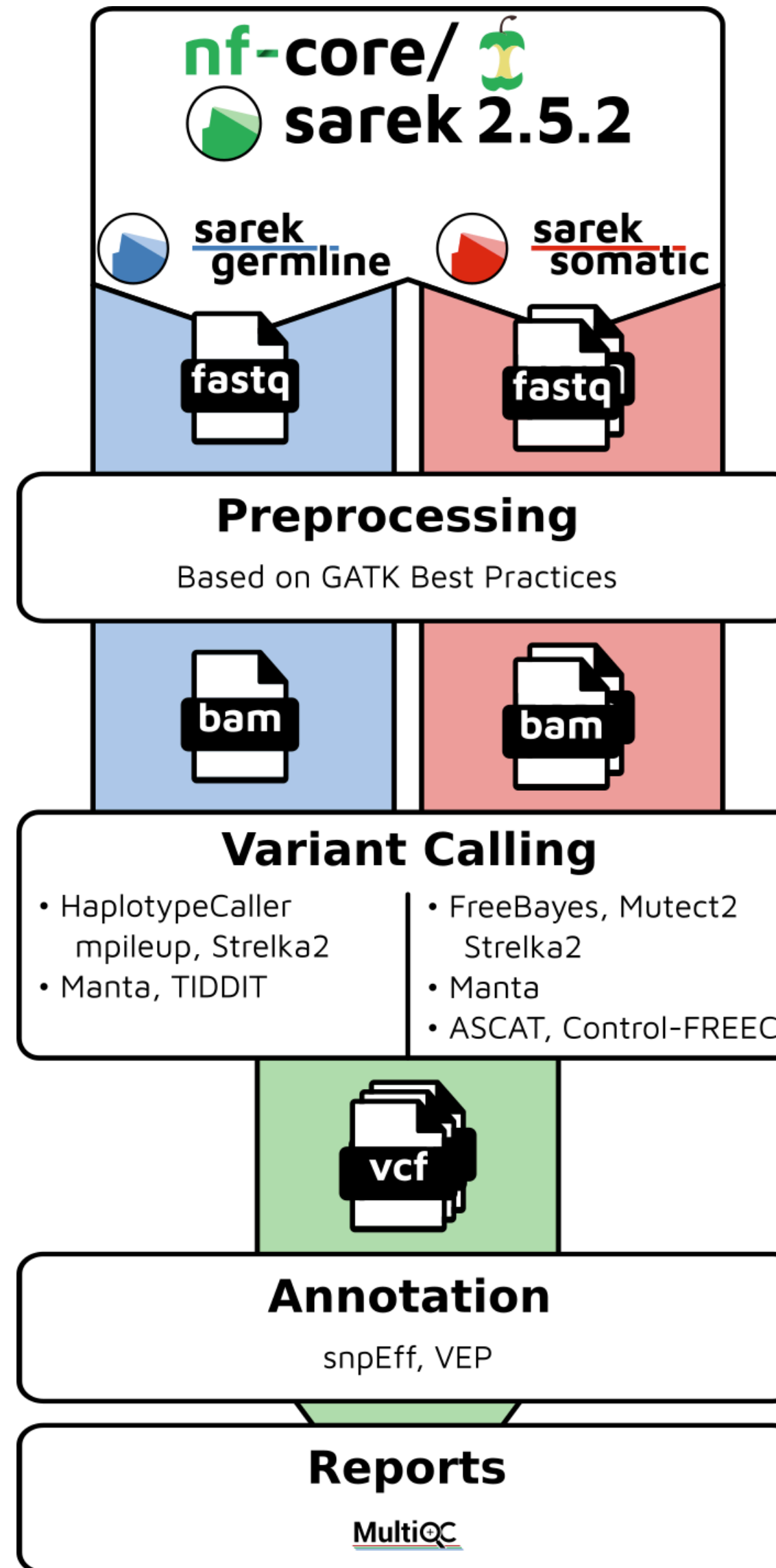# CZI EOSS

December 2020

# Project Mission

- Deploy anywhere, securely

- Enable best-practise software engineering

- Operate from GUI to API

seqeralabs

# NGS Variant Calling

https://github.com/nf-core/sarek



seqeralabs

# Parasite Genome Annotation

https://github.com/sanger-pathogens/
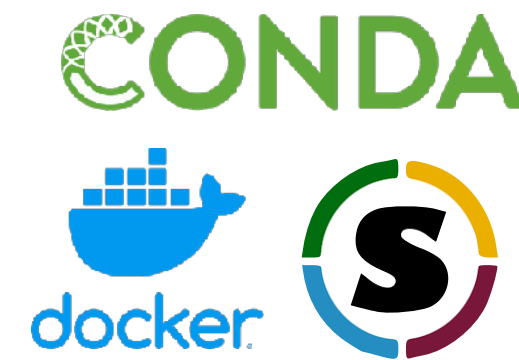companion

# What is Nextflow?

## nextflow *script*

Write code
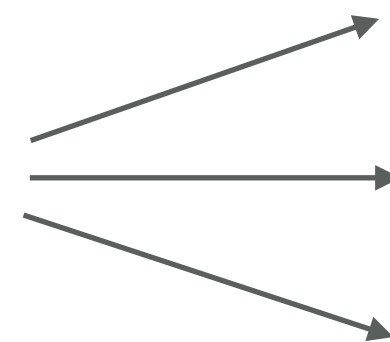in any language.

Define orchestration with
dataflow programming.

Define software
dependencies
with containers.

Version
control.

## nextflow *runtime*

Orchestration of tasks to
deploy anywhere with ease.

seqeralabs

# Task example

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

# Task example

```
process align_sample {

    input:
    path 'reference.fa' from genome_ch
    path 'sample.fq' from reads_ch

    output:
    path 'sample.bam' into bam_ch


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

}
```

seqeralabs

# 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
    """
}
```

# How does it work?

- **Fast prototyping** ⇒ custom DSL that enables tasks 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

seqeralabs

# DSL2

A major revision of the Nextflow DSL

- Pipeline modularisation

- Component reuse

- Fluent definition of recurrent implementation patterns

seqeralabs

# Nextflow syntax - DSL 2

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.

seqeralabs

# Centralised cluster orchestration

computer cluster



NFS/Lustre

submit jobs

cluster node

cluster node

cluster node

cluster node

cluster node

- Nextflow orchestrates workflow execution submitting jobs to a compute schedular

- Can run in the head node or a compute node

- Requires a shared storage to exchange data between tasks

- Ideal for corse-grained parallelism

# Cloud orchestration



AWS / Google Cloud

submit jobs

EC2 + Spot

VMs

VMs

VMs

VMs

VMs

AWS Batch

S3 bucket
object store

nextflow

- Nextflow orchestrates workflow execution via AWS Batch

- Launched 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.

seqeralabs

kubernetes

Azure

aws

# PORTABILITY

# PORTABILITY



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

# PORTABILITY



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

# Enterprise adoption

# Tower for Nextflow

AWS Reference Deployment

Internet

AWS Cloud

AWS Batch → Nextflow → Compute

Storage

EKS

Back-end Container

Security/OAuth

Controllers

Services

JPA / Hibernate

Redisson

Front-end Container

NGINX web server

Angular

Aurora Database

ElastiCache

*Support available for all cloud or on-premise infrastructure.*

# Sharing and collaboration



Share workflow executions with colleagues who can follow along live on the progress of the pipeline or catch up on the results at any time in the future.

# API Release

## Nextflow Tower API 1.0.0

Nextflow Tower service API

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

### ⬆ actions

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

### collaborators

| | | |
|---|---|---|
| **GET** | /collaborators/workflow/{workflowId} | List the collaborators of the workflow for the authenticated user |
| **POST** | /collaborators | Add a collaborator |
| **DELETE** | /collaborators/{collaboratorId}/workflow/{workflowId} | Delete a collaborator |

seqeralabs

# Pipelines Feature



Preconfigured workflows with input fields based on JSON schema.

# Extended documentation

# CZI EOSS Goals

**WP1**

Goal 1.1          Reach 8,000 monthly Nextflow users

Goal 1.2          Reach 150 active nf-core Slack members

Goal 1.3          Sustain community support and advocacy activities

Goal 1.4          Appoint project positions

**WP2**

Goal 2.1          Improve Nextflow scalability and support for public clouds

Goal 2.2          Expand the support for web-based usage of nf-core pipelines

Goal 2.3          Nextflow kernel for Jupyter notebooks

Goal 2.4          GA4GH API compliant TES & WES executors

Goal 2.5          Migrate existing nf-core pipelines to Nextflow DSL2

Goal 2.6          Introduce module-level testing for nf-core

Goal 2.7          nf-core template for DSL2 pipelines

**WP3**

Goal 3.1          Nextflow website refresh: learning, community and support

Goal 3.2          Four new community Nextflow training events

Goal 3.3          Ten expanded bursaries for 2020 Nextflow Community Conference

Goal 3.4          Establish 2 nf-core hackathons and 2 user workshops

Goal 3.5          Tutorial videos on nf-core website

seqeralabs

nextflow + nf-core

https://nf-co.re

Community efforts to collect production ready analysis pipelines built with Nextflow

https://nf-co.re

**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

Chan Zuckerberg Initiative