

Reproducibility for Data Pipelines and Analyses

Nextflow & nf-core



Open Science



Open Science



Open Data

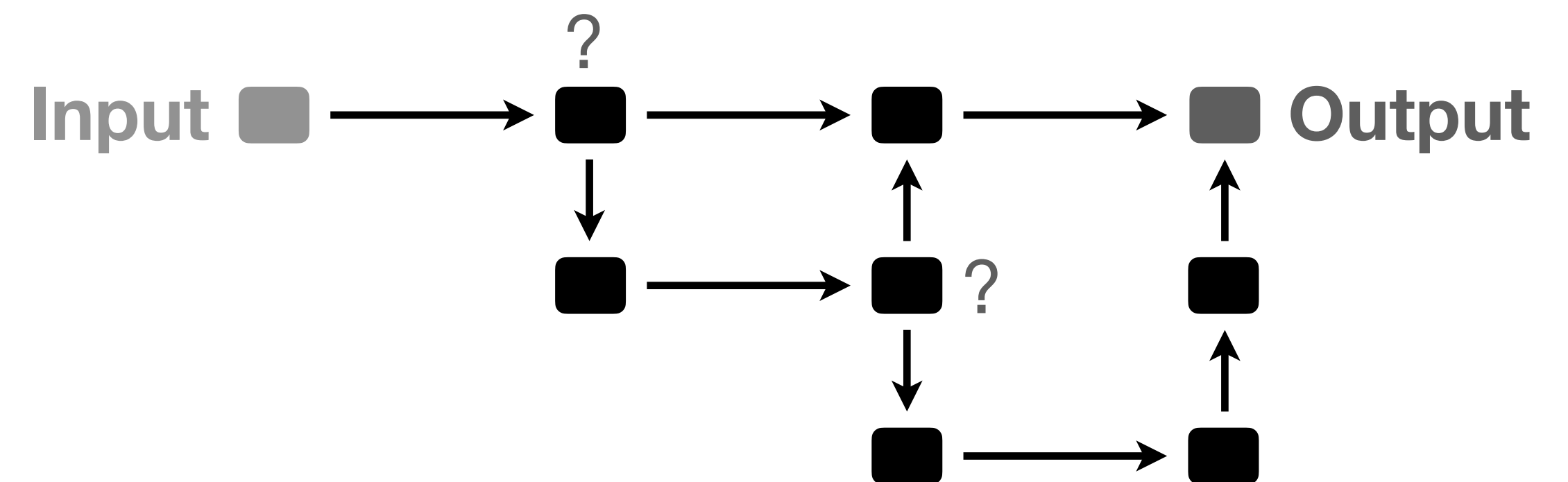
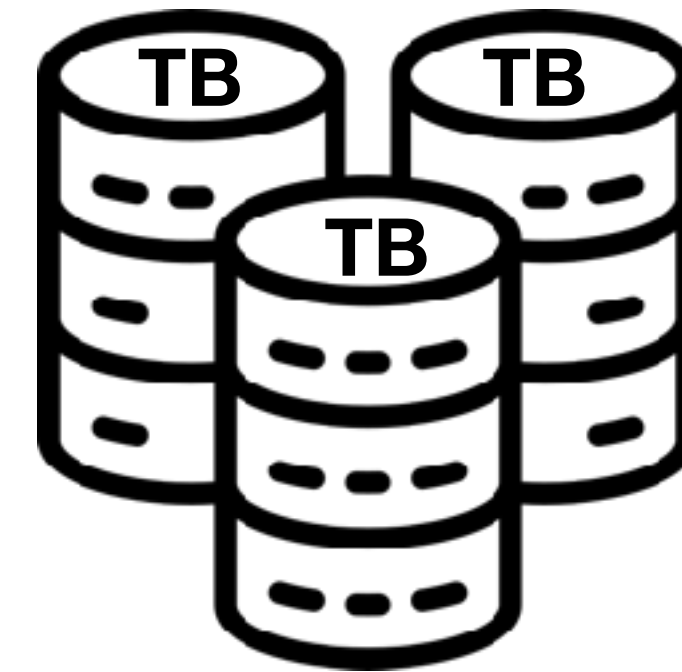


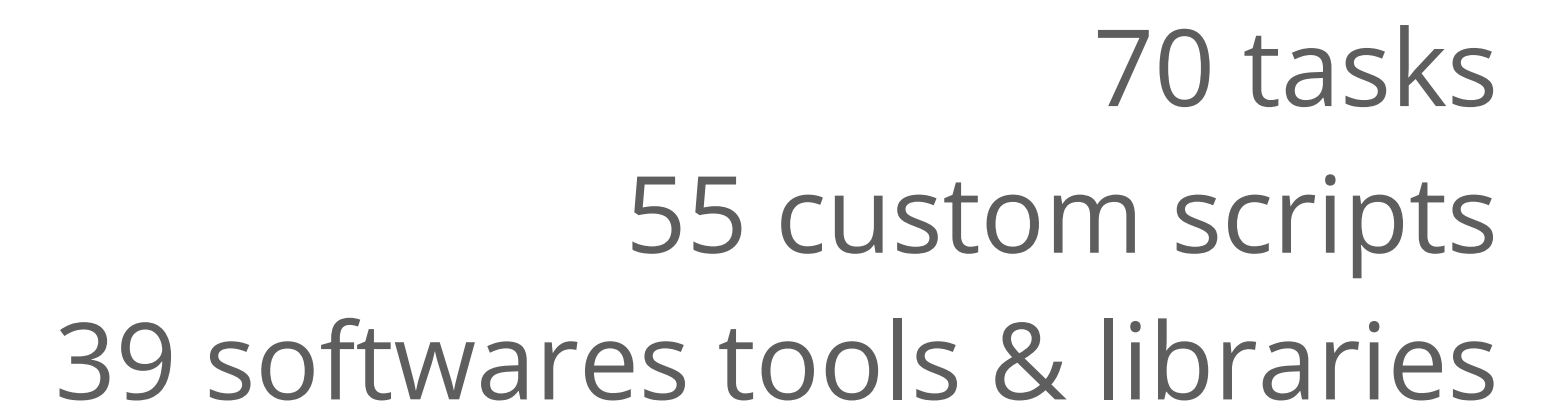
Open Data Community



Genomics


Workflow



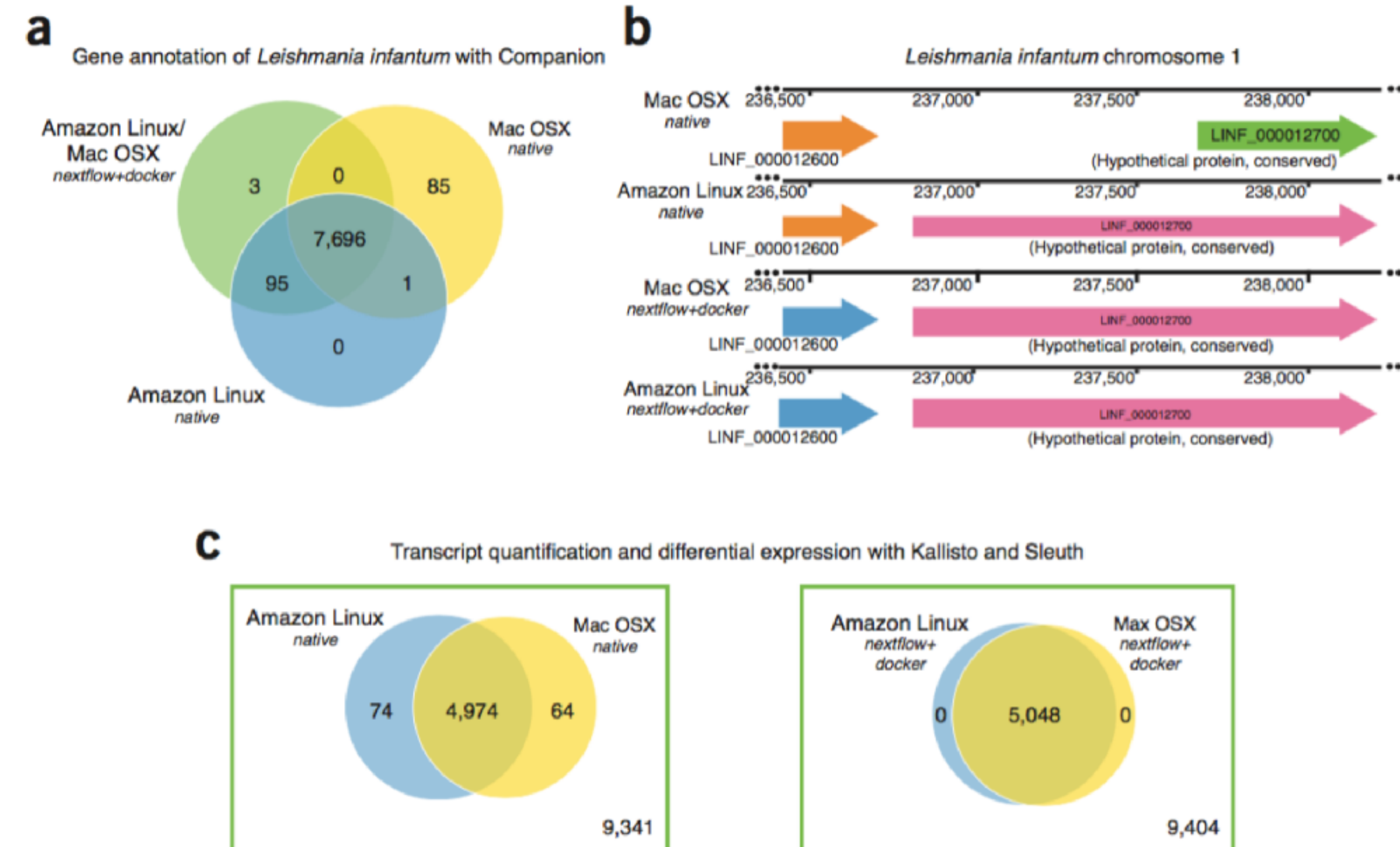


Steinbiss et al., Companion parassite
genome annotation pipeline
DOI: 10.1093/nar/gkw292

Nextflow enables reproducible computational workflows

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

NATURE BIOTECHNOLOGY VOLUME 35 NUMBER 4 APRIL 2017

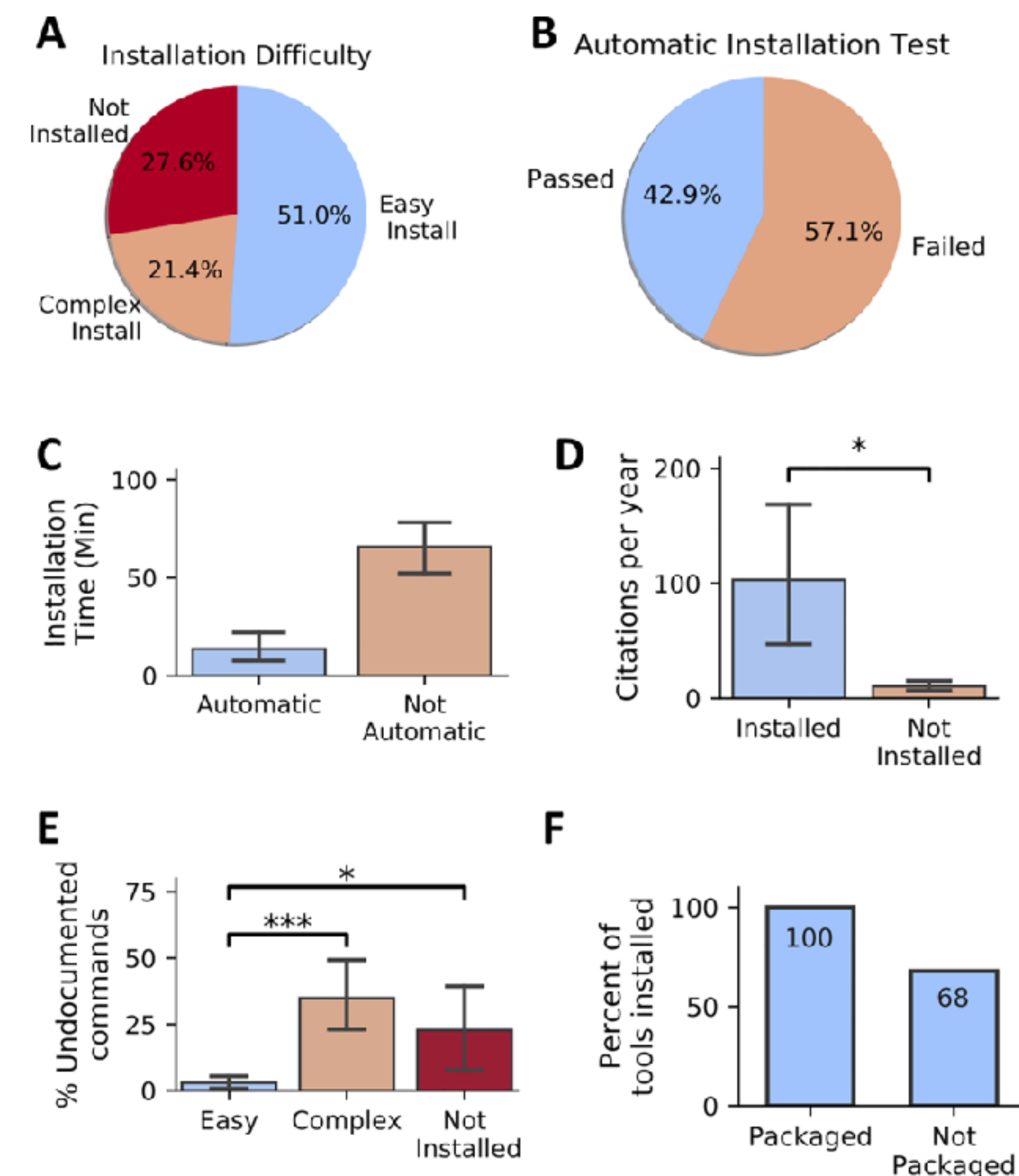


Challenges and recommendations to improve the installability and archival stability of omics computational tools (2019)

Serghei Mangul, et al. Plos Biology <https://doi.org/10.1371/journal.pbio.3000333>

We found that 28% of all omics software resources are currently not accessible through URLs published in the paper.

Among the tools selected 49% were difficult to install or could not be installed at all.

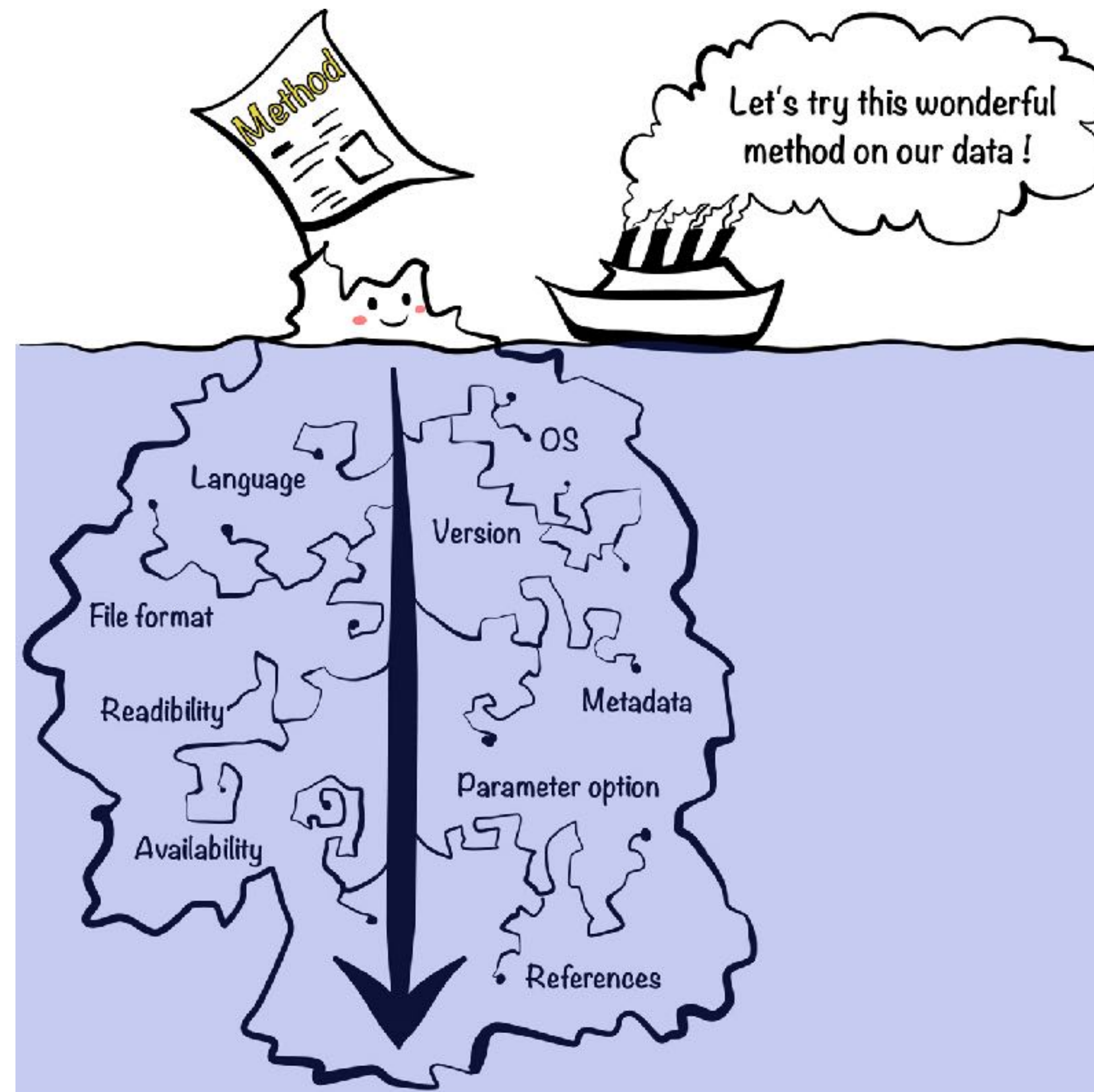


Comparison of the Companion pipeline annotation of *Leishmania infantum* genome executed across different platforms *

Platform	Amazon Linux	Debian Linux	Mac OSX
<i>Number of chromosomes</i>	36	36	36
<i>Overall length (bp)</i>	32.032.223	32.032.223	32.032.223
<i>Number of genes</i>	<u>7.781</u>	<u>7.783</u>	<u>7.771</u>
<i>Gene density</i>	236,64	<u>236,64</u>	<u>236,32</u>
<i>Number of coding genes</i>	7.580	<u>7.580</u>	<u>7570</u>
<i>Average coding length (bp)</i>	1.764	<u>1.764</u>	<u>1.762</u>
<i>Number of genes with multiple CDS</i>	113	<u>113</u>	<u>111</u>
<i>Number of genes with known function</i>	4.147	<u>4.147</u>	<u>4.142</u>
<i>Number of t-RNAs</i>	<u>88</u>	<u>90</u>	88

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





Kim et al. Experimenting with reproducibility: a case study of robustness in bioinformatics, *GigaScience*, Volume 7, Issue 7, July 2018. <https://doi.org/10.1093/gigascience/giy077>

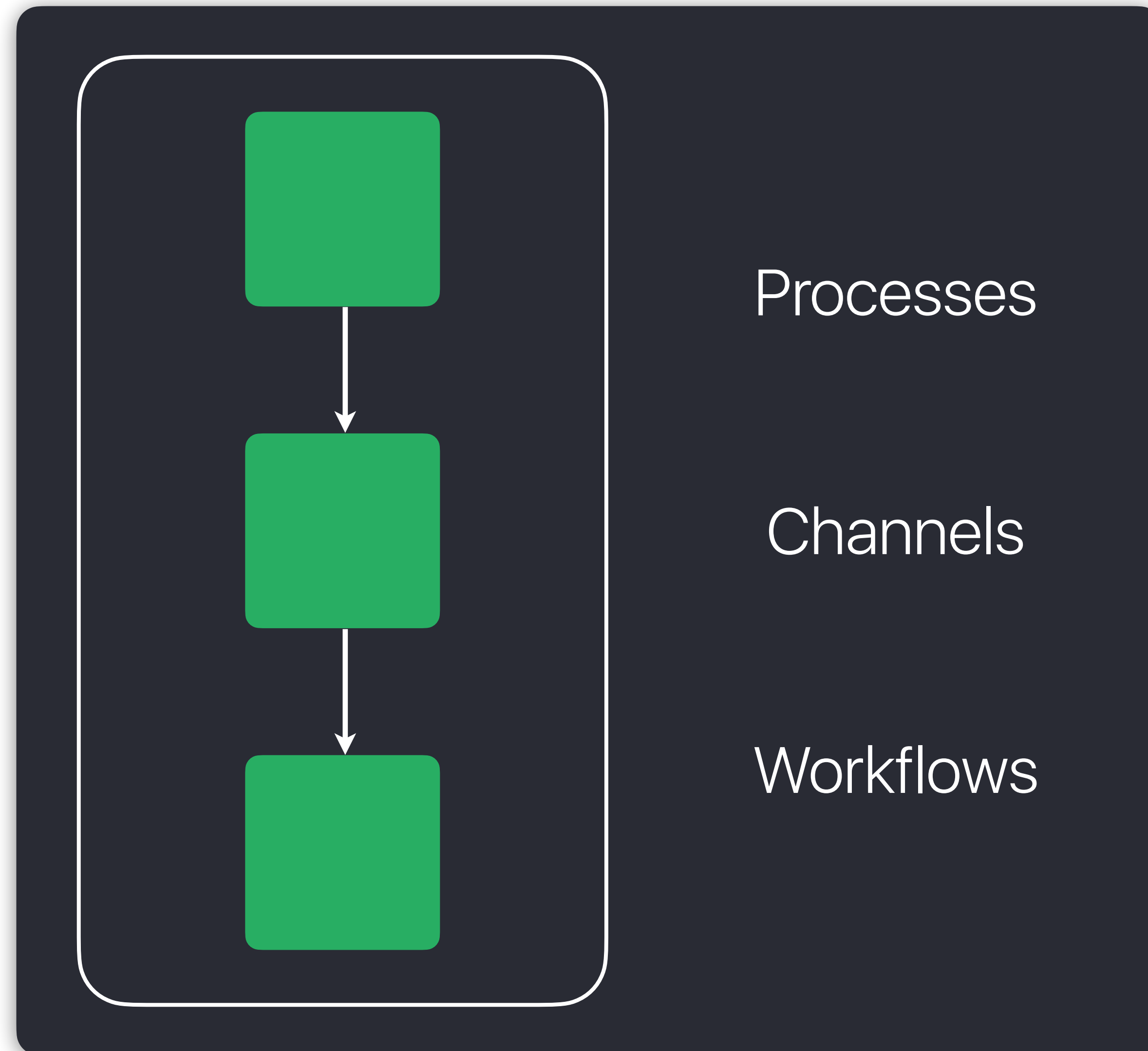
nextflow



nextflow

Language

nextflow



nextflow

```
#!/usr/bin/env nextflow
process fastqc {
    input:
    path input

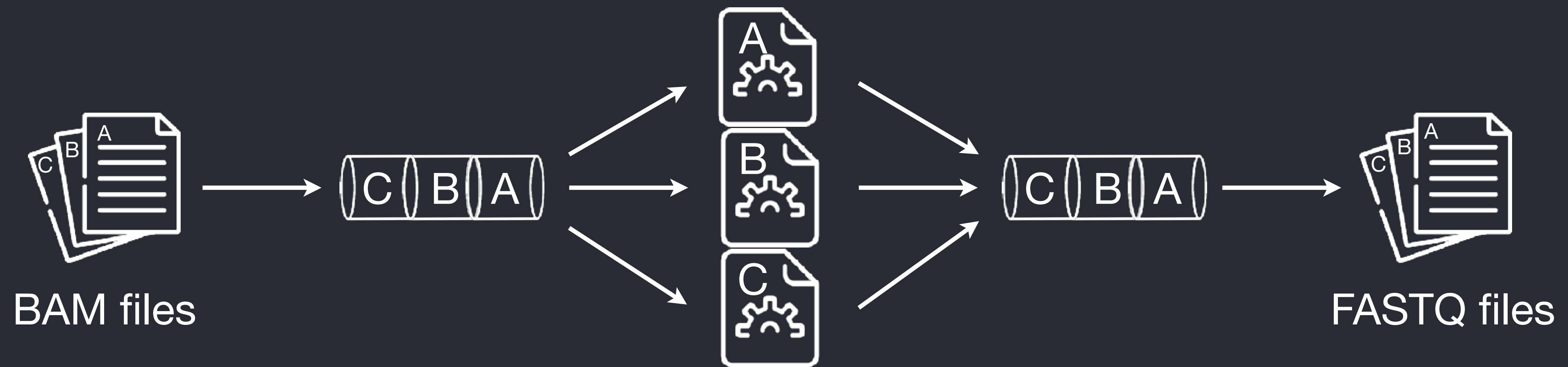
    output:
    path "*_fastqc.{zip,html}"

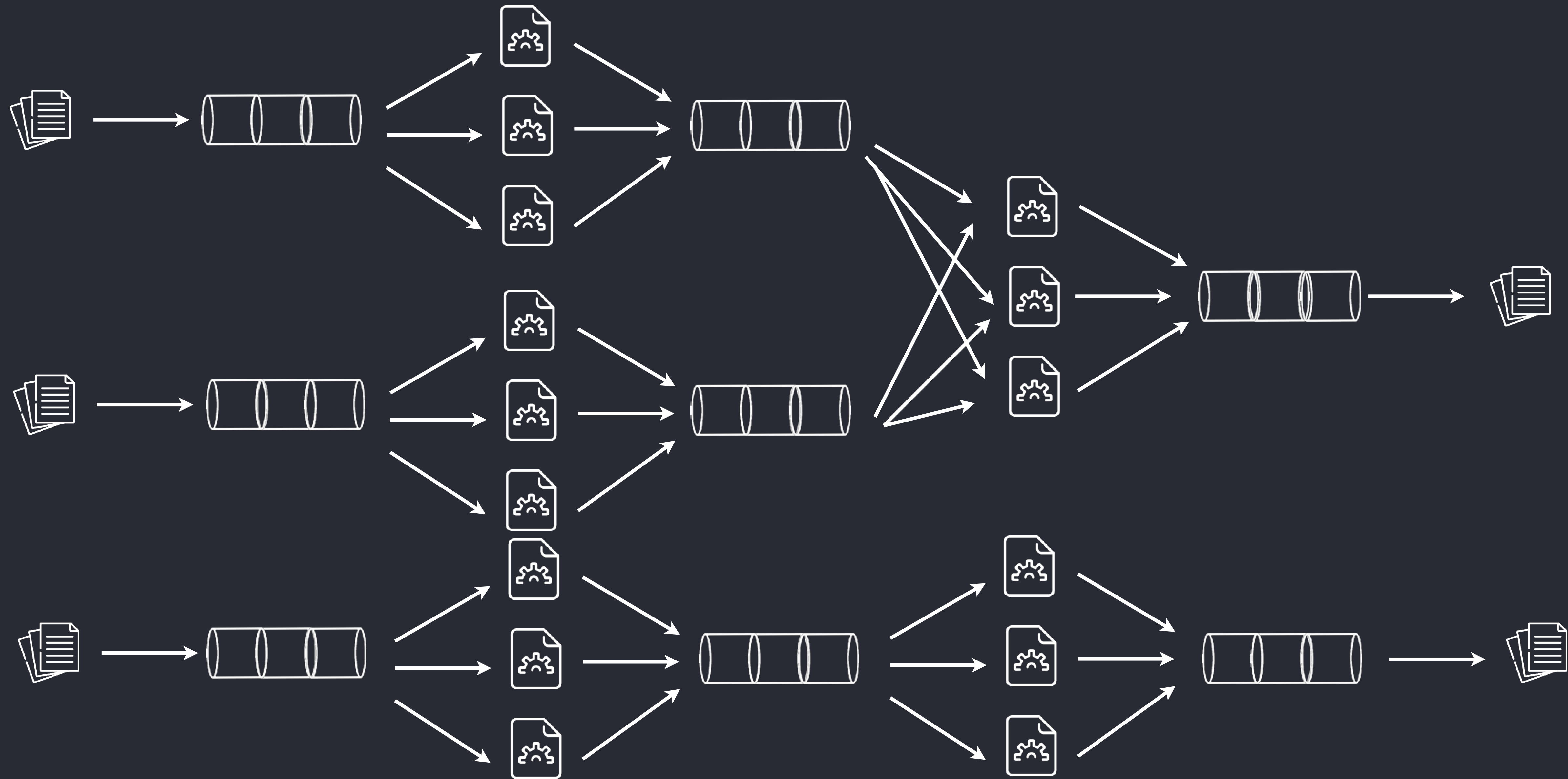
    script:
    """
    fastqc -q $input
    """
}

workflow {
    Channel.fromPath("*.fastq.gz") | fastqc
}
```

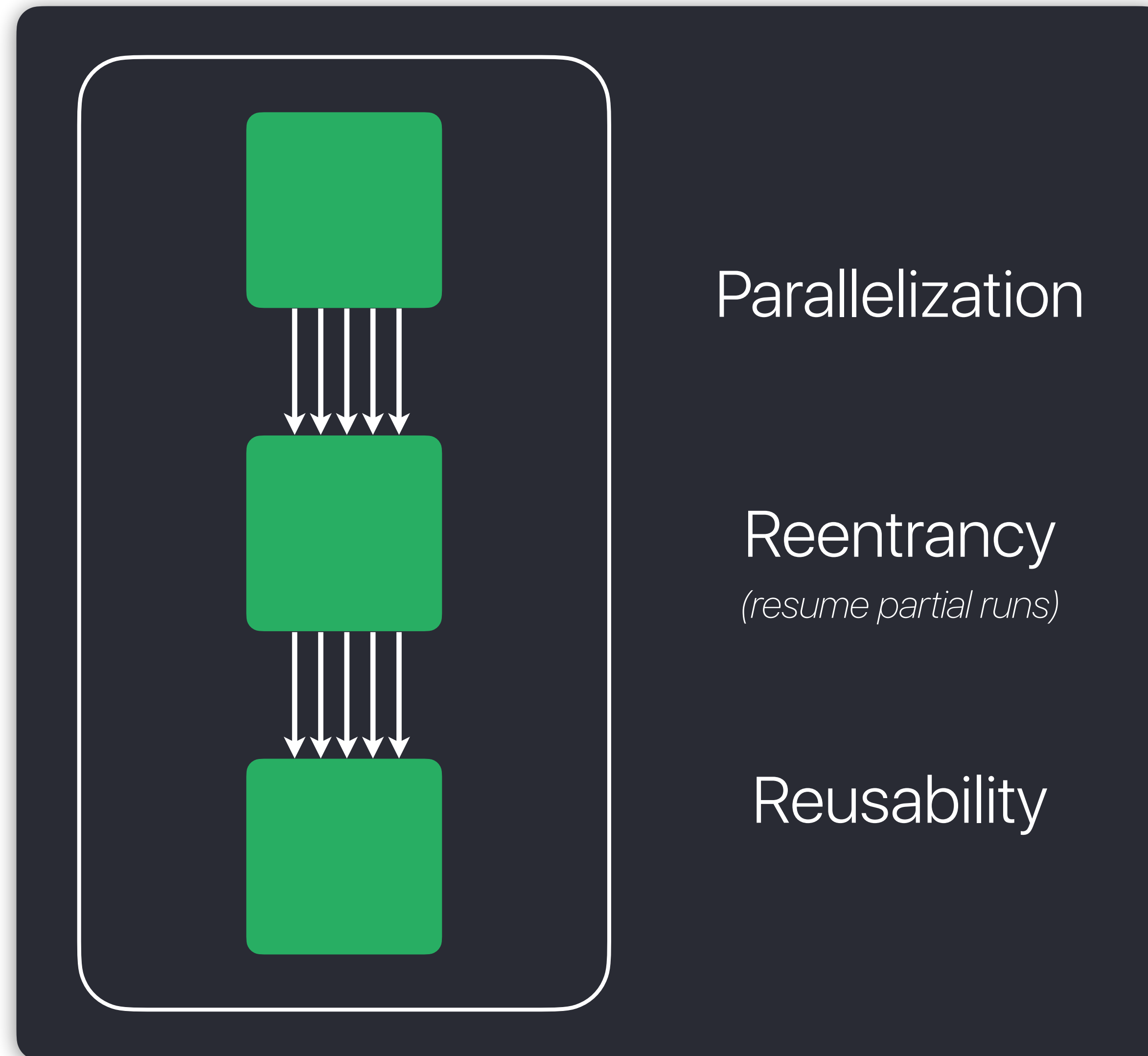
Implicit parallelism

```
workflow {  
  Channel  
  | .fromPath("data/*.bam") | bam_to_fastq  
}
```





nextflow



nextflow

Language

Software

Compute

nextflow



AWS CodeCommit



Azure Repos

Software

Compute

nextflow



git



GitHub



Bitbucket



GitLab



Gitea



AWS CodeCommit



Azure Repos



docker.



podman



Singularity



CONDA

Compute

nextflow



git



GitHub



Bitbucket



GitLab



Gitea



AWS CodeCommit



Azure Repos



docker



podman



Singularity



CONDA

SGE



Microsoft Azure



slurm
workload manager



openstack



LSF

PBS



Google Cloud



kubernetes



Reproducibility

```
process INDEX {  
  
    input:  
    | path transcriptome  
    output:  
    | path 'salmon_index'  
  
    script:  
    """"  
    salmon index --threads $task.cpus -t $transcriptome -i salmon_index  
    """"  
}
```



Reproducibility

```
process INDEX {  
  conda "bioconda::salmon=1.9.0"  
  input:  
    path transcriptome  
  output:  
    path 'salmon_index'  
  
  script:  
    """"  
    salmon index --threads $task.cpus -t $transcriptome -i salmon_index  
    """"  
}
```



Reproducibility

```
process INDEX {  
  container "nextflow/rnaseq-nf"  
  input:  
  | path transcriptome  
  output:  
  | path 'salmon_index'  
  
  script:  
  """"  
  salmon index --threads $task.cpus -t $transcriptome -i salmon_index  
  """"  
}
```



Reproducibility

```
process INDEX {  
  executor "slurm"  
  input:  
    path transcriptome  
  output:  
    path 'salmon_index'  
  
  script:  
    """"  
    salmon index --threads $task.cpus -t $transcriptome -i salmon_index  
    """"  
}
```



Reproducible

Between runs

Portable

Between systems

Scalable

Everywhere

nextflow



nextflow

nf-core 





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

<https://nf-co.re>

nf-core



73
PIPELINES

<https://nf-co.re>



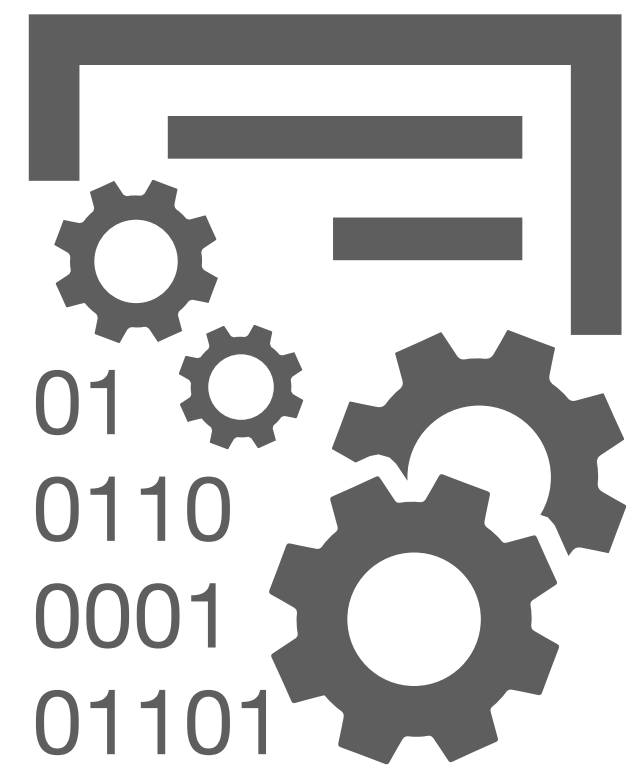
TOOLS

Running pipelines

Writing pipelines

Testing / automation

<https://nf-co.re>



708

MODULES

24

SUB-WORKFLOWS

<https://nf-co.re>

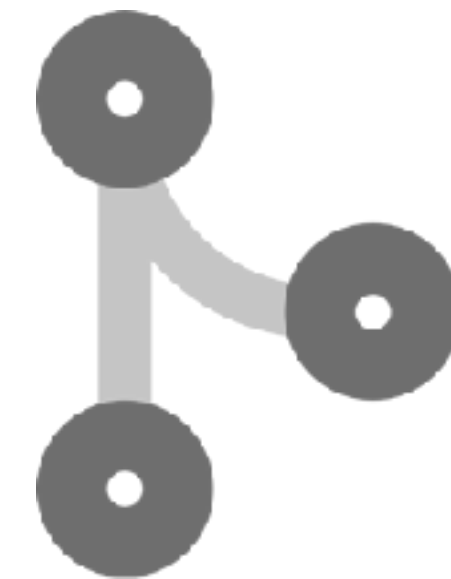
nf-core



Develop with
the community



Use a common
template



Collaborate,
don't duplicate

<https://nf-co.re>

4135

Slack users

527

GitHub organisation
members

1583

GitHub contributors

3294

Twitter followers

90

Repositories

11.61K

Pull Requests

32.37K

Commits

5.21K

Issues


<https://nf-co.re>



<https://nf-co.re>

Correspondence | Published: 13 February 2020

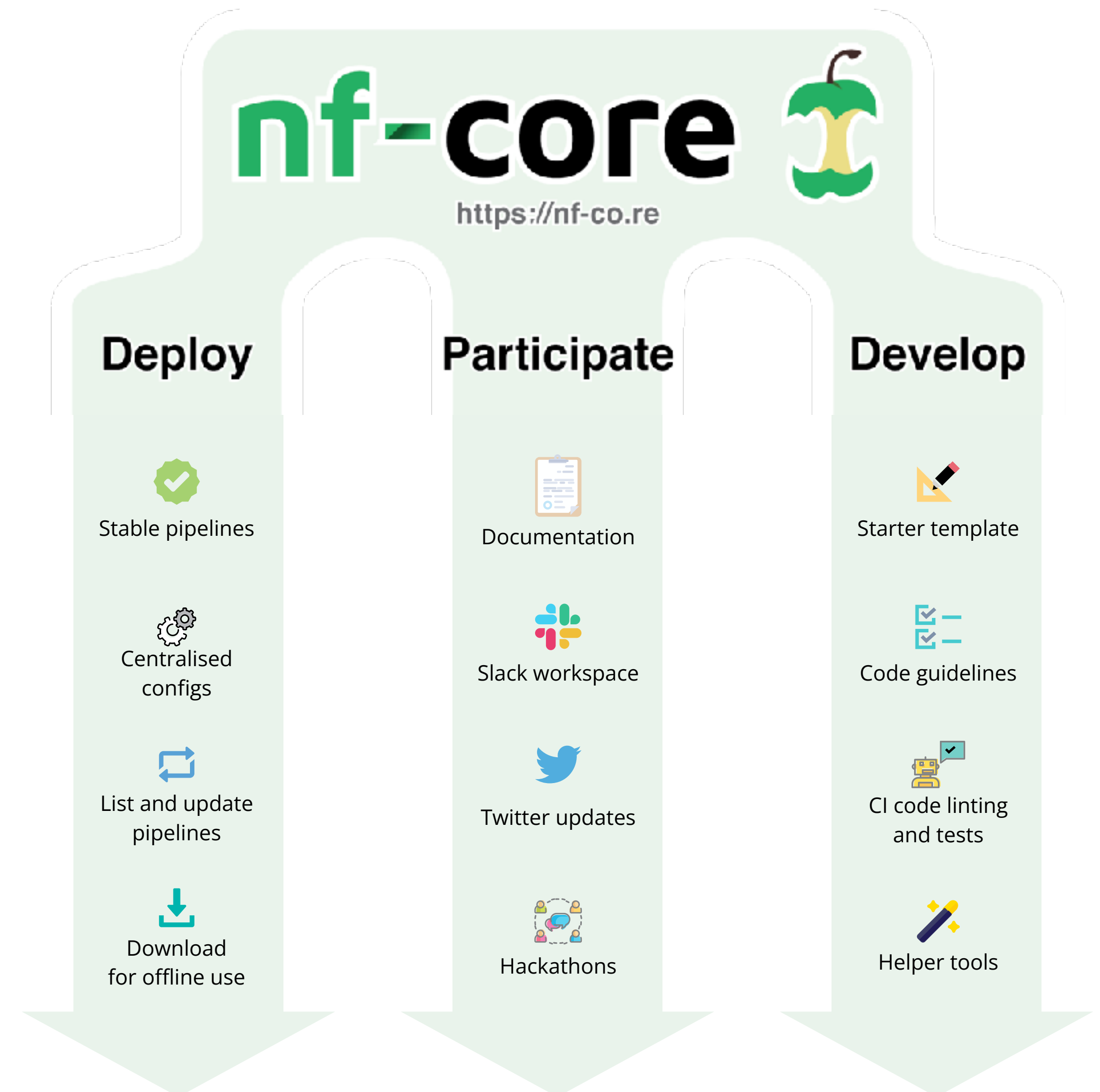
The nf-core framework for community-curated bioinformatics pipelines

Philip A. Ewels, Alexander Peltzer, Sven Fillinger, Harshil Patel, Johannes Alneberg, Andreas Wilm, Maxime Ulysse Garcia, Paolo Di Tommaso & Sven Nahnsen 

Nature Biotechnology **38**, 276–278(2020) | [Cite this article](#)

3253 Accesses | **3** Citations | **172** Altmetric | [Metrics](#)

To the Editor — The standardization, portability and reproducibility of analysis pipelines are key issues within the bioinformatics community. Most bioinformatics pipelines are designed for use on-premises; as a result, the associated software dependencies and execution logic are likely to be tightly coupled with proprietary computing environments. This can make it difficult or even impossible for others to reproduce the ensuing results, which is a fundamental requirement for the validation of scientific findings. Here, we introduce the nf-core framework as a means for the development of collaborative, peer-reviewed, best-practice analysis pipelines (Fig. 1). All nf-core pipelines are written in Nextflow and so inherit the ability to be executed on most computational infrastructures, as well as having native support for container technologies such as Docker and Singularity. The nf-core community (Supplementary Fig. 1) has developed a suite of tools that automate pipeline creation, testing, deployment and synchronization. Our goal is to provide a framework for high-quality bioinformatics pipelines that can be used across all institutions and research facilities.



Join the community



<https://nf-co.re/join>



nextflow *tower*



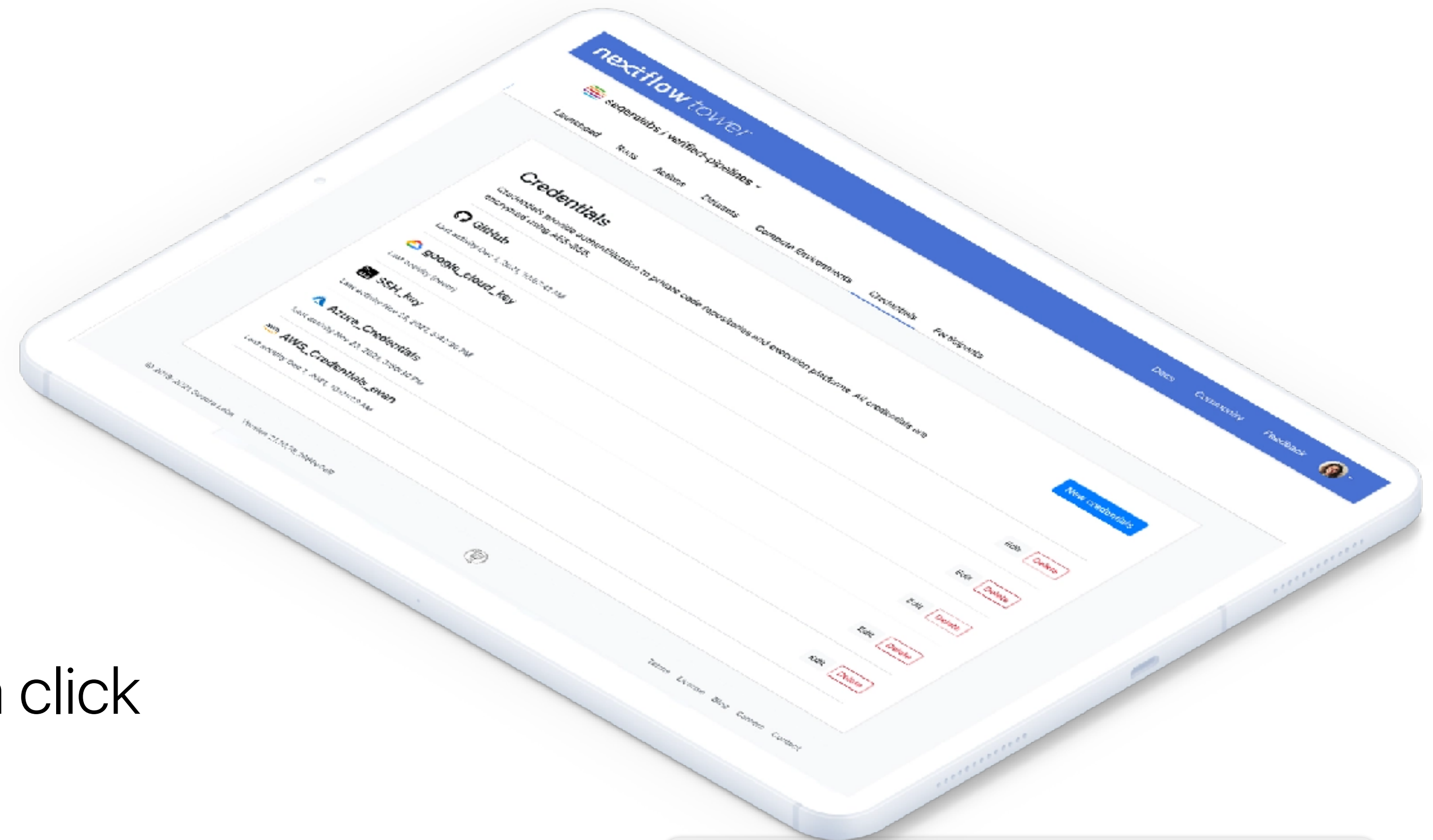


Intuitive launchpad interface

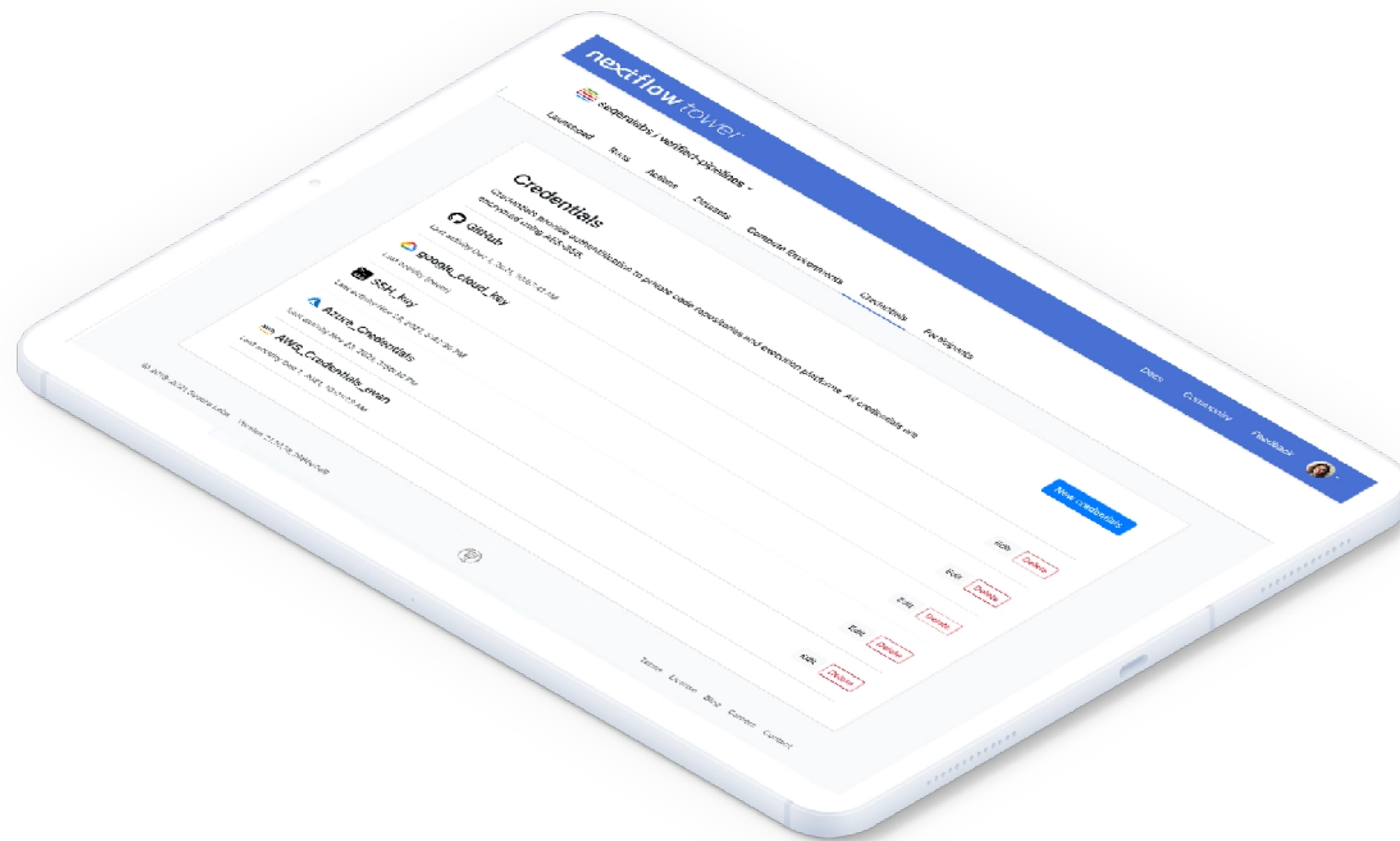
Launch, manage, and monitor

Share runs and work in teams

Create cloud infrastructure with a click



<https://tower.nf>



<https://tower.nf>

Marcel Ribeiro-Dantas, PhD



seqeralabs

<https://seqera.io>

<http://mribeirodantas.xyz>



mribeirodantas

mribeirodantas@seqera.io



mribeirodantas

**Chan Zuckerberg
Initiative**



nextflow SUMMIT 2022

<https://summit.nextflow.io>

Nextflow / nf-core training

13-16 March 2023

nf-core hackathon

27-29 March 2023

<https://nf-co.re>

<https://nextflow.io>