# INTRODUCTION TO NEXTFLOW

Elixir workflows workshop - 1 Dec 2021
Paolo Di Tommaso - CTO Seqera Labs

seqeralabs

# AGENDA

- 9.30 - 9.50: Intro to Nextflow core concepts

  Paolo Di Tommaso, Seqera Labs

- 9.50 - 10.00: Deploy in the cloud with Nextflow Tower

  Evan Floden, Seqera Labs

- 10.00 - 12.00: Nextflow hands-on tutorial

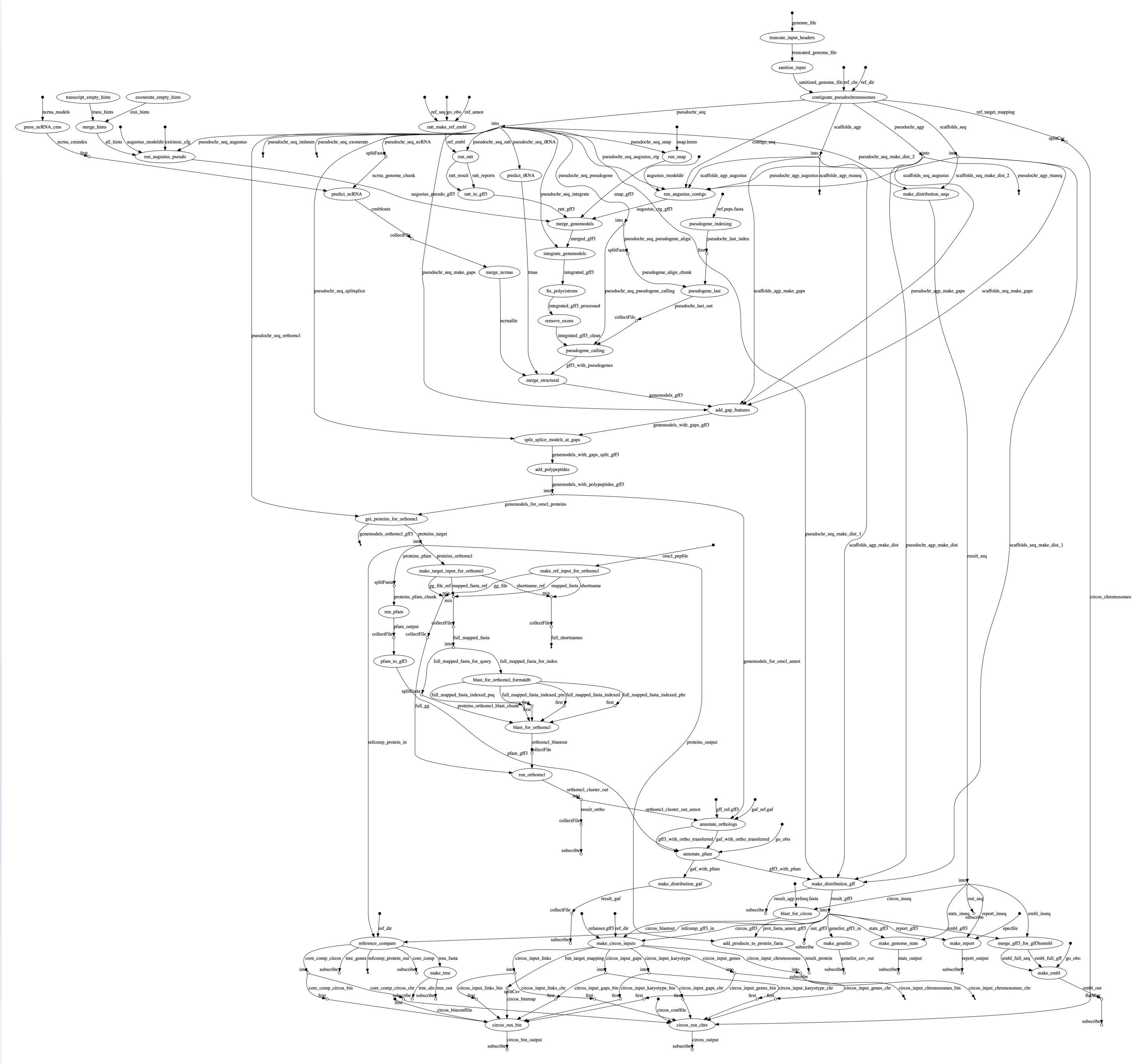  Luca Cozzuto, Centre for Genomic Regulation (CRG)

# QUICK BIO



- Paolo Di Tommaso

- Software engineer

- Creator & maintainer of Nextflow project

- CTO & Co-founder Seqera Labs

# GENOMICS WORKFLOWS

- Data analysis applications to extract information from large genomic datasets (TB)

- Embarrassingly parallelisation, can spawn 100s-100k jobs over distributed cluster

- Mash-up of many different tools and scripts

- Complex executions and dependencies graphs → very fragile ecosystem

- 70 tasks
- 55 custom scripts
- 39 software tools & libraries

**Complexity ⇒
instability ⇒
not reproducible result**

Steinbiss et al., *Companion parassite genome annotation pipeline*, DOI: 10.1093/nar/gkw292
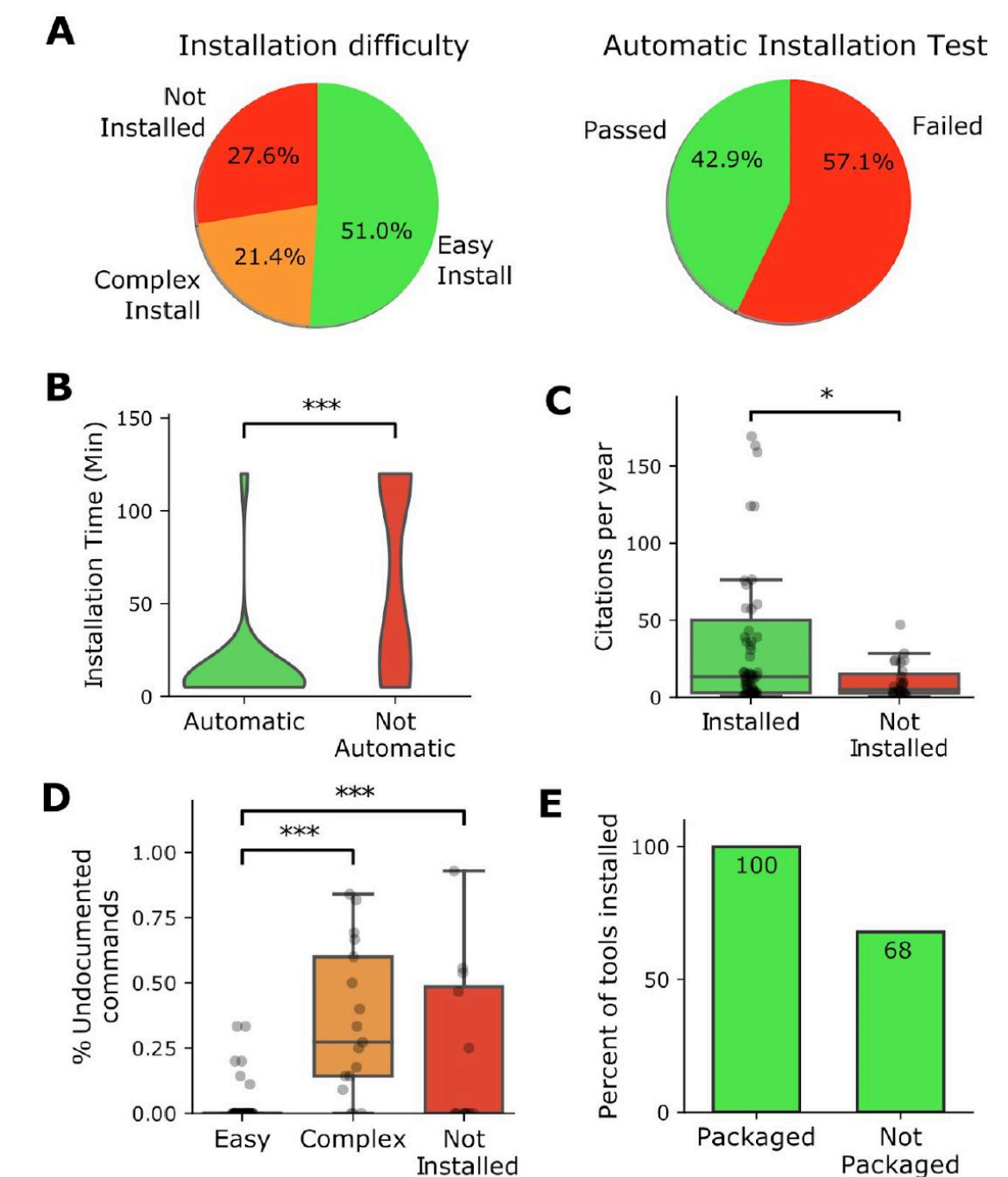
# THE PROBLEM WITH REPRODUCIBILITY

**A comprehensive analysis of the usability and archival stability of omics computational tools and resources**

We found that 26% of all omics software resources are currently *not accessible* through URLs published in the paper.
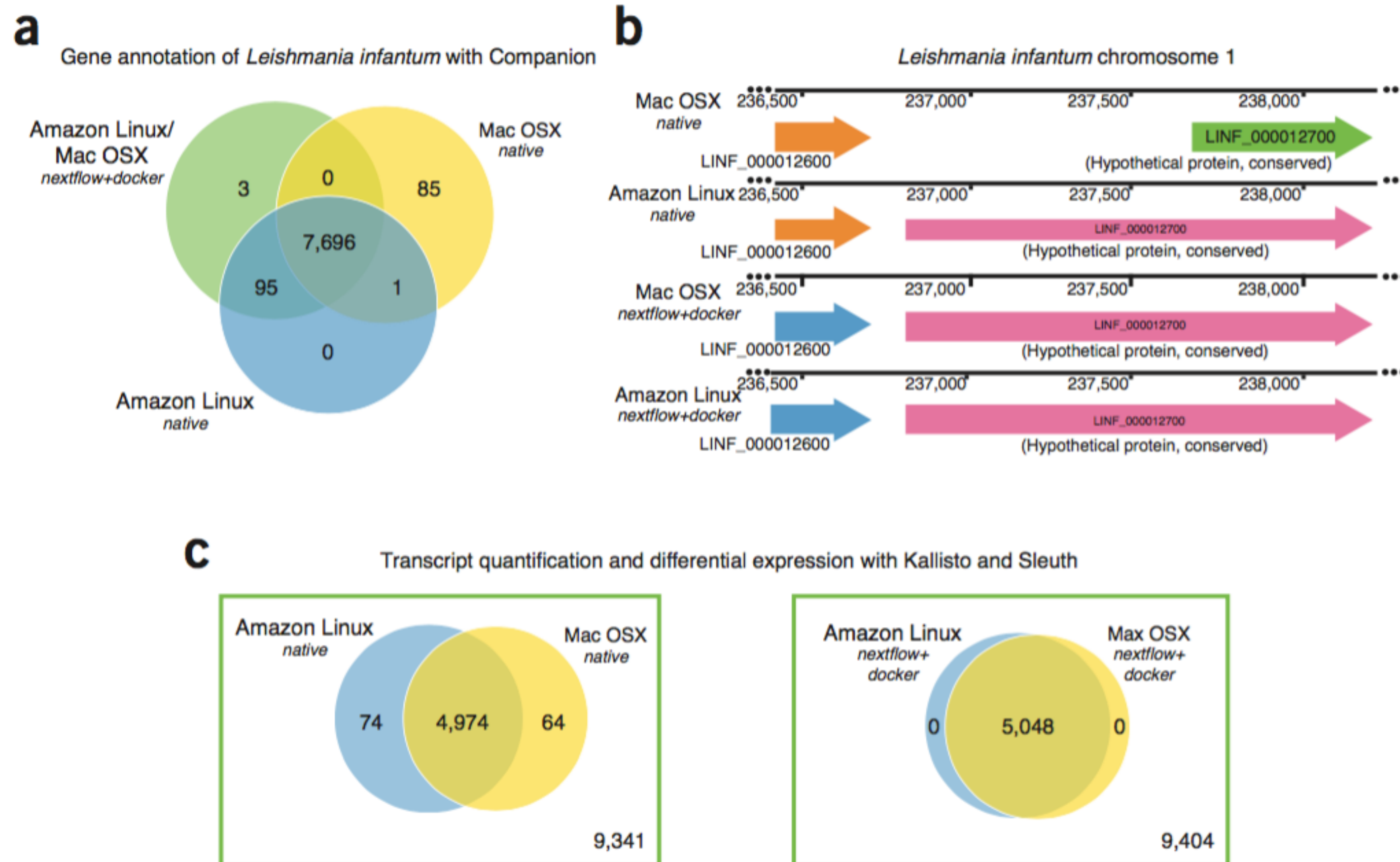
Among the tools selected 49% were deemed "difficult to install," and *28% of the tools failed to be installed* due to problems in the implementation.

# Nextflow enables reproducible computational workflows

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

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

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

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

# MAIN CHALLENGES

Scalability & parallelisation

Enable portability

Guarantee reproducibility

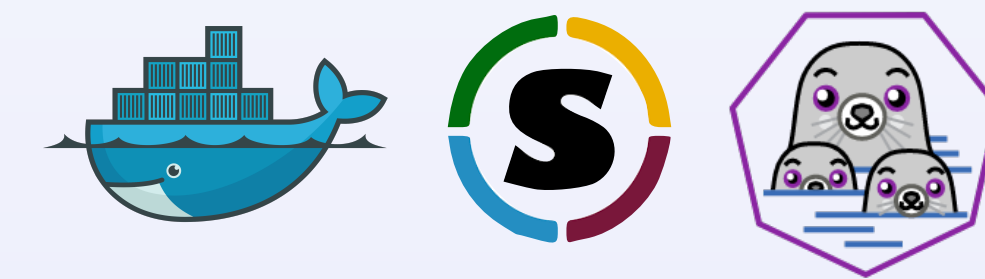# THE nextflow FUNDAMENTALS FOR SCALEABLE GENOMIC WORKFLOWS

## Write code in any language
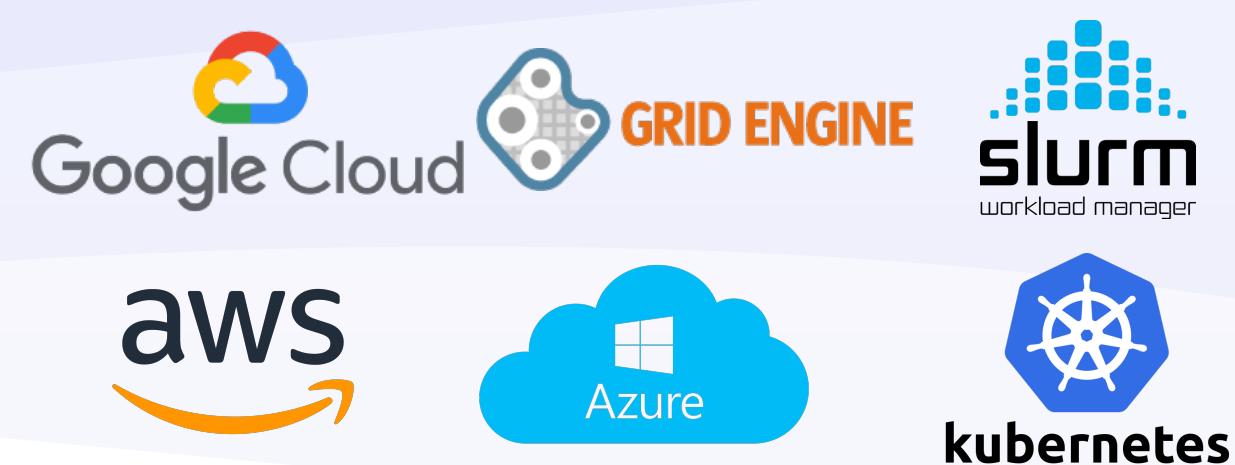


## Declarative parallelisation



## Isolate dependencies with containers



## Deploy everywhere



## Version control



## Open source

# SIMPLE EXAMPLE

```
#
# one-liner to convert a BAM file to a FASTQ file
#


 samtools view file.bam \
    | awk 'BEGIN {FS="\t"} {print "@" $1 "\n" $10 "\n+\n" $11}' > file.fq
```

# NEXTFLOW PROCESS

```
process bam_to_fastq {
  input:
    path "file.bam"

  output:
    path "file.fq"

  script:
   """
    samtools view file.bam \
    | awk 'BEGIN {FS="\t"} {print "@" $1 "\n" $10 "\n+\n" $11}' > file.fq
   """
}
```

# NEXTFLOW WORKFLOW

```
process bam_to_fastq {
    input:
      path "file.bam"

    output:
      path "file.fq"

    script:
      """
      samtools view file.bam \
        | awk 'BEGIN {FS="\t"} {print "@" $1 "\n" $10 "\n+\n" $11}' > file.fq
      """
}


workflow {

    channel.fromPath("/data/sample.bam") | bam_to_fastq

}
```

# IMPLICIT PARALLELISM

```
process bam_to_fastq {
  input:
    path "file.bam"

  output:
    path "file.fq"

  script:
    """
    samtools view file.bam \
    | awk 'BEGIN {FS="\t"} {print "@" $1 "\n" $10 "\n+\n" $11}' > file.fq
    """
}


workflow {

  channel.fromPath("/data/*.bam")  | bam_to_fastq

}
```
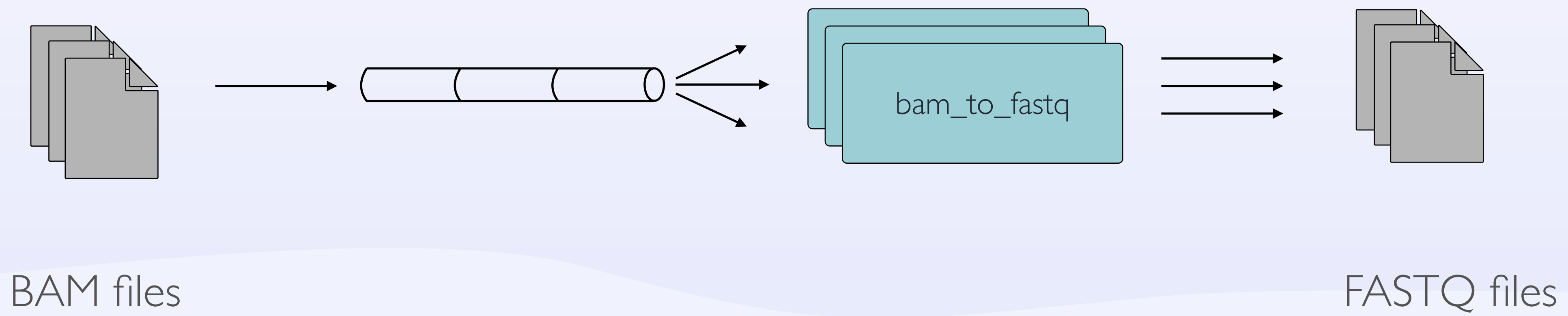
# IMPLICIT PARALLELISM

```
channel.fromPath("/data/*.bam")
```



BAM files

bam_to_fastq

FASTQ files

# CONTAINERISATION

Config file

Nextflow

job    job    job
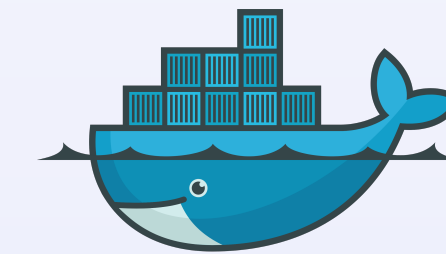
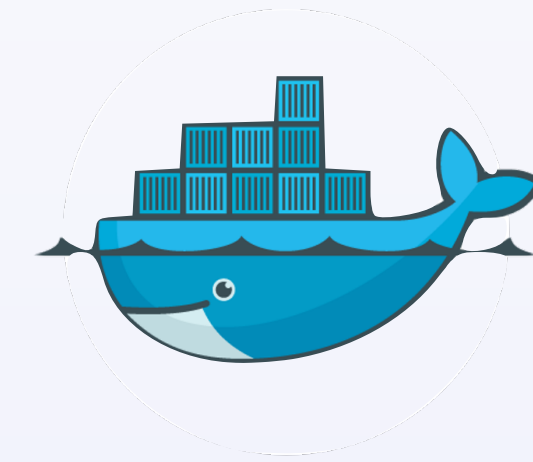- Nextflow envisioned the use of software containers to fix computational reproducibility
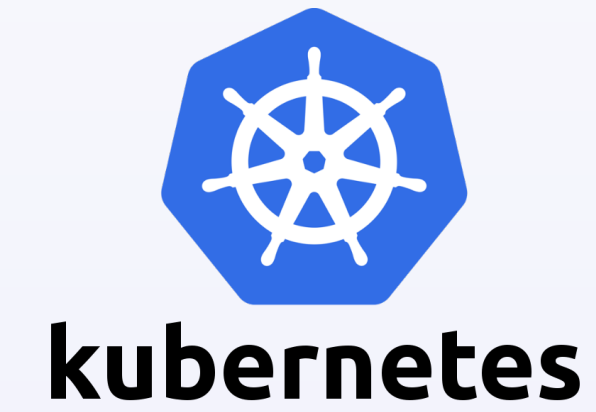
- Mar 2014, support for Docker
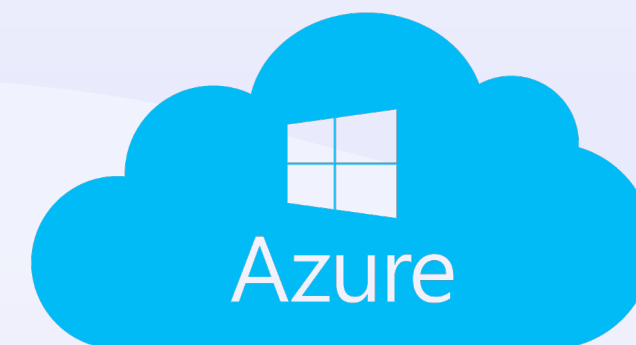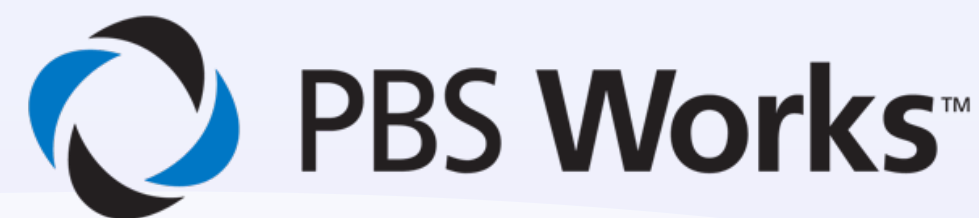
- Dec 2016, support for Singularity

- Jan 2020, support for Podman

# MULTI-PLATFORM

# EXECUTION REPORT

## Nextflow Report    Summary  Resources  Tasks                    [trusting_cuvier]

# Nextflow workflow report

`[trusting_cuvier]` *(resumed run)*

Workflow execution completed successfully!

**Run times**
Fri Apr 27 23:19:53 CEST 2018 - Sat Apr 28 03:18:15 CEST 2018 (completed a day ago, duration: **3h 58m 21s**)

| 5329 succeeded | 2849 cached | |
|---|---|---|

**Nextflow command**

```
nextflow run main.nf -profile crg --std_align=true --default_align=true --align_method=CLUSTALO,MAFFT --
tree_method=CLUSTALO,MAFFT_PARTTREE --seqs=/users/cn/egarriga/datasets/homfamClustalo/seqs/*.fa --
refs=/users/cn/egarriga/datasets/homfamClustalo/refs/*.ref -with-report -with-trace -resume -bg
```

| | |
|---|---|
| **CPU-Hours** | 156.6 (31.5% cached, 4.6% failed) |
| **Launch directory** | /nfs/users2/cn/egarriga/projects/dpa_cp |
| **Work directory** | /nfs/users2/cn/egarriga/projects/dpa_cp/work |
| **Project directory** | /nfs/users2/cn/egarriga/projects/dpa_cp |
| **Script name** | main.nf |
| **Script ID** | 6ff267a42e50448d41927a6e5a9787fc |
| **Workflow session** | 087c9bc8-e488-4311-88aa-961138c42fd6 |
| **Workflow profile** | crg |
| **Workflow container** | cbcrg/regressive-msa:v0.2.4 |
| **Container engine** | singularity |
| **Nextflow version** | version 0.28.2, build 4782 (06-04-2018 12:25 UTC) |

# EXECUTION REPORT

# TIMELINE CHART

## Processes execution timeline

Launch time: 15 Jun 2016 15:03
Elapsed time: 49m 9s

| Process | Duration / Memory |
|---|---|
| downloadReference (1) | 5.7s / 549.2 GB |
| downloadSRA (1) | 59.9s / 1.2 GB |
| decompressReference (1) | 1.5s / - |
| indexReference (1) | 5.7s / 80.9 MB |
| extractSRA (1) | 1m 11s / 150.5 MB |
| trim (1) | 36.1s / 8.1 GB |
| mergeTrimEnds (1) | 40.2s / 31.6 MB |
| filterKMC (1) | 2m 14s / 8.6 GB |
| filterKHMER (1) | 38m 31s / 8.9 GB |
| alignReads_kmc (1) | 10m 17s / 2 GB |
| sortAlignment_kmc (1) | 1m 52s / 538.9 MB |
| indexAlignment_kmc (1) | 8.8s / 36.2 MB |
| callVariants_kmc (1) | 20m 16s / 89.2 MB |
| alignReads_khmer (1) | 1m 23s / 1.8 GB |
| sortAlignment_khmer (1) | 22s / 524.4 MB |
| indexAlignment_khmer (1) | 2.2s / 36.2 MB |
| callVariants_khmer (1) | 5m 21s / 89.1 MB |

# EDITORS

main.nf — ~/projects/callings-nf

main.nf

```
48  log.info """\
49  C A L L I N G S  -  N F    v 1.0
50  ================================
51  genome   : $params.genome
52  reads    : $params.reads
53  variants : $params.variants
54  blacklist: $params.blacklist
55  results  : $params.results
56  gatk     : $params.gatk
57  """
58
59  /*
60   *  Parse the input parameters
61   */
62
63  GATK          = params.gatk_launch
64  genome_file   = file(params.genome)
65  variants_file = file(params.variants)
66  blacklist_file = file(params.blacklist)
67  reads_ch      = Channel.fromFilePairs(params.reads)
68
69
70  /**********
71   * PART 1: Data preparation
72   *
73   * Process 1A: Create a FASTA genome index (.fai) with samtools for GATK
74   */
75  process '1A_prepare_genome_samtools' {
76    tag "$genome.baseName"
77
78    input:
79        file genome from genome_file
80
81    output:
82        file "${genome}.fai" into genome_index_ch
83
84    script:
85    """
86    samtools faidx ${genome}
87    """
88  }
89
90
91    * Process 1B: Create a FASTA genome sequence dictionary with Picard for GATK
92    */
93  process '1B_prepare_genome_picard' {
94    tag "$genome.baseName"
95
96    input:
97        file genome from genome_file
98    output:
99        file "${genome.baseName}.dict" into genome_dict_ch
100
101   script:
102   """
103   PICARD=`which picard.jar`
104   java -jar \$PICARD CreateSequenceDictionary R= $genome O= ${genome.baseName}.dict
105   """
106  }
```

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

---

main.nf

pcpu.sh    dd.sh    data    main.nf

```
47  log.info """\
48  C A L L I N G S  -  N F    v 1.0
49  ================================
50  genome    : $params.genome
51  reads     : $params.reads
52  variants  : $params.variants
53  blacklist: $params.blacklist
54  results   : $params.results
55  gatk      : $params.gatk
56  """
57
58
59  /*
60   *  Parse the input parameters
61   */
62
63  GATK          = params.gatk_launch
64  genome_file   = file(params.genome)
65  variants_file  = file(params.variants)
66  blacklist_file = file(params.blacklist)
67  reads_ch      = Channel.fromFilePairs(params.reads)
68
69
70  /**********
71   * PART 1: Data preparation
72   *
73   * Process 1A: Create a FASTA genome index (.fai) with samtools for GATK
74   */
75
76  process '1A_prepare_genome_samtools' {
77    tag "$genome.baseName"
78
79    input:
80        file genome from genome_file
81
82    output:
83        file "${genome}.fai" into genome_index_ch
84
85    script:
86    """
87    samtools faidx ${genome}
88    """
89  }
90
91
92  /*
93   * Process 1B: Create a FASTA genome sequence dictionary with Picard for GATK
94   */
95
96  process '1B_prepare_genome_picard' {
97    tag "$genome.baseName"
98
99    input:
100       file genome from genome_file
101   output:
102       file "${genome.baseName}.dict" into genome_dict_ch
103
104   script:
105   """
106   PICARD=`which picard.jar`
107   java -jar \$PICARD CreateSequenceDictionary R= $genome O= ${genome.baseName}.dict
108   """
109   }
110
```

⑂ master*   ↻ 0 ↓ 9↑      ⊗ 0 ⚠ 0                    Ln 1, Col 1   Spaces: 4   UTF-8   LF   Nextflow

# SUMMARY

- Data analysis reproducibility is hard and it's often underestimated.

- Nextflow does not provide a magic solution but enables best-practices and provides support for community and industry standards.

- It strictly separates the application logic from the configuration and deployment logic, enabling self-contained workflows.

- Applications can be easily deployed across different environment in a reproducible manner with a single command.

- The functional/reactive model allows applications to scale with ease.