# Using the IARC nextflow bioinformatics pipelines
# Day 2

Matthieu Foll
IARC course: May 23-24 2018

International Agency for Research on Cancer
Lyon, France

# Agenda

**Wednesday 23 May**

| | |
|---|---|
| 09:00-10:00 | Introduction to bioinformatics pipelines, nextflow, docker, Github and the IARC organization |
| 10:00-10:30 | Practical application: running your first pipeline |
| 10:30-11:00 | Break |
| 11:00-11-30 | The hidden structure of nextflow: work folder and configuration |
| 11:30-12:30 | Practical application: configuring, crashing, resuming and debugging pipelines |

**Thursday 24 May**

| | |
|---|---|
| 09:00-09:30 | Introduction to HPC clusters and running pipelines on a cluster. |
| 09:30-10:30 | Practical application: trace and visualise pipeline execution with log files. |
| 10:30-11:00 | Break |
| 11:00-11h30 | Introduction to the nextflow language: understanding what the pipelines are doing |
| 11:30-12:30 | Practical application: advanced usages toward reproducibility (choosing a container, Github releases and branches, modifying a pipeline etc.) |

# IARC bioinformatics platform

International Agency for Research on Cancer bioinformatics platform

⊙ Lyon, France   ⌘ http://www.iarc.fr   ✉ follm@iarc.fr

## Pinned repositories

Customize pinned repositories

### ≡ IARC-nf

List of IARC bioinformatics nextflow pipelines

★ 4   ⑂ 2

### ≡ needlestack

Multi-sample somatic variant caller

🔵 R   ★ 16   ⑂ 4

### ≡ mutspec

Mutation Spectra Analysis

🔵 Perl   ★ 5

### ≡ alignment-nf

Whole Exome/Whole Genome Sequencing alignment pipeline

🟠 Groovy   ★ 4

### ≡ RNAseq-nf

RNAseq analysis pipeline

🔴 HTML   ★ 2   ⑂ 1

### ≡ mutect-nf

Mutect pipeline with Nextflow

🔴 HTML   ★ 3   ⑂ 1

Search repositories...    Type: All ▾    Language: All ▾    📖 New

## needlestack

Multi-sample somatic variant caller

docker   bam-files   pipeline   nextflow   ngs   regression

🔵 R   ★ 16   ⑂ 4   Updated 5 days ago

## RNAseq_analysis_scripts

Scripts for RNA seq analysis

🔵 R   ★ 1   Updated 14 days ago

### Top languages

🔵 R   🔴 HTML   🔵 Perl   🟢 Shell
🟠 Groovy

### Most used topics    Manage

nextflow   ngs   pipeline

bam-files   alignment

# Our philosophy

- "Do It Once, Do It Right, And Use It Everywhere"

- "Keep it simple, stupid" (KISS principle):

  - most systems work best if they are kept simple

  - simplicity should be a key goal in design

  - code easier to maintain and to understand

# Our design

- Too much automation is not for us:

  - Hard to read, to maintain and to keep modular

  - eg: we prefer to have one alignment pipeline; one variant calling pipeline; one annotation pipeline; one QC pipeline.

- One pipeline = one GitHub repository

- Docker containers hosted on DockerHub, compatible with Singularity

- CircleCI for tests and deployment

- Standardised readme, params, help etc. (one shared template)

- Use GitHub issues and releases

- Master branch ← beta branch ← dev branch

# In practice

- Entry point: GitHub group

  - https://github.com/IARCbioinfo

- One central repo references all nextflow pipelines:

  - https://github.com/IARCbioinfo/IARC-nf

  - List pipelines with a short description

  - One pipeline = one repo, ends with "-nf"

  - Common instructions to use the pipelines (install nextflow, configuration, basic usage, docker…)

- A "template-nf" nextflow "hello-world" repo

# High Performance Computing (HPC) cluster

Users

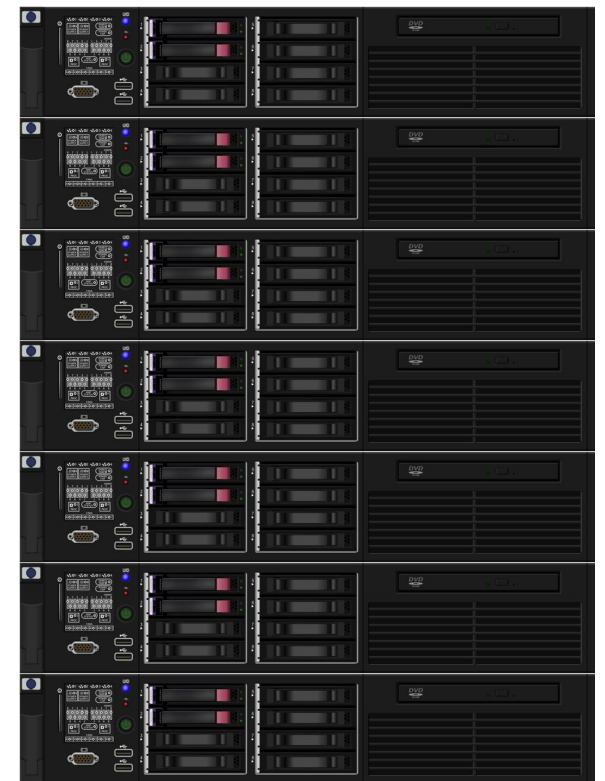Computing nodes

Network switch

Management node

Storage nodes

Backup

# Job scheduler

- Computer program for controlling unattended background program execution of jobs

- Choose host based on:

  - Compute resource availability

  - Execution time allocated to user

  - Number of simultaneous jobs allowed for a user

  - Job priority

  - …

- SGE, SLURM, PBS, Torque, **LSF** etc.

# LSF most useful commands

- **bsub command** to launch a job. Useful options:

  - **-I**: Interactive job

  - **-oo:** output file name (overwrite)

  - **-eo**: error file name (overwrite)

  - **-J**: job name

  - **-m**: choose the hosts (or list of hosts, e.g. bsub -m "cn08 cn09 cn10")

  - **-n**: number of CPUs

  - **-R "rusage[mem=XXX]"** **-M XXX**: ask for XXX MB of RAM

- **bjobs** to monitor jobs. Useful options:

  - **-w** for the full name

  - **-l** option for more details

  - **-u username** for a given user (or **-u all** for all users)

**Nextflow is doing this for you!!!**

International Agency for Research on Cancer

World Health Organization

# Example

```
samtools view -H NA06984.bam
```

with docker becomes:

```
docker run -it --rm -v $PWD:$PWD -w $PWD --entrypoint /bin/bash samtools_img -c
                                            "samtools view -H NA06984.bam"
```

with LSF becomes:

```
bsub -oo out.txt -eo err.txt -n 1 -R "rusage[mem=1000]" -M 1000 -J samtools
                                        "samtools view -H NA06984.bam"
```

with docker and LSF becomes:

```
bsub -oo out.txt -eo err.txt -n 1 -R "rusage[mem=1000]" -M 1000 -J samtools
 "docker run -it --rm -v $PWD:$PWD -w $PWD --entrypoint /bin/bash samtools_img -c
                                        \"samtools view -H NA06984.bam\" "
```

Hum… but we use singularity on our cluster… The syntax is actually different!
(`singularity exec …`)

**Nextflow is doing this for you!!!**

# Running on a cluster



```
[x140083:nf_coverage_demo follm$ ssh follm@jupiter.iarc.fr
[follm@jupiter.iarc.fr's password:
 Last login: Fri Mar  3 11:25:06 2017 from 10.10.13.47
[[follm@hn ~]$ cd /data/follm/nf_coverage_demo/
[[follm@hn nf_coverage_demo]$ nextflow run plot_coverage.nf --bam_folder BAM/ --bed TP53.bed
 N E X T F L O W  ~  version 0.23.4
 Launching `plot_coverage.nf` [clever_ramanujan] - revision: 66d7be595f
[warm up] executor > local
 [1e/ce2c42] Submitted process > coverage (3)
 [f3/a818f6] Submitted process > coverage (1)
 [2a/31fbc9] Submitted process > coverage (2)
 [dc/ff61ac] Submitted process > coverage (6)
```

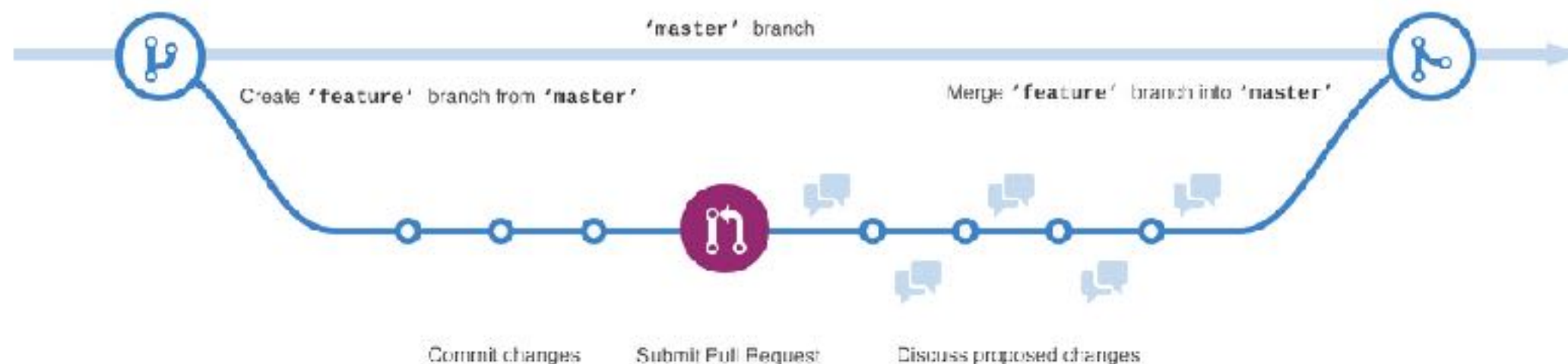**Running on the head-node: don't do that!!!**

```
sftp://10.10.156.1//mnt/beegfs/follm/nf_coverage_demo/nextflow.config
1    process.executor='lsf'
2
```

```
[[follm@hn nf_coverage_demo]$ echo "process.executor='lsf'" > nextflow.config
[[follm@hn nf_coverage_demo]$ nextflow run plot_coverage.nf --bam_folder BAM/ --bed TP53.bed
 N E X T F L O W  ~  version 0.23.4
 Launching `plot_coverage.nf` [gigantic_hawking] - revision: 66d7be595f
[warm up] executor > lsf
 [37/af473d] Submitted process > coverage (4)
 [ee/1ce1f9] Submitted process > coverage (1)
 [fe/84c2d2] Submitted process > coverage (6)
```

International Agency for Research on Cancer

World Health Organization

Control queue size with `-qs` option

# Speaking GitHub language

- **Branches**: master branch considered definitive, use other branches to experiment.

- Changes are called **commit**: one can **push** or **pull** commits to repository or a branch. Commits have unique IDs (hash).

- **Pull requests**: ask someone to **merge** the commits you did into their branch.

- **Releases** are user friendly commit IDs with **tag** names

- **Issues** are used to discuss bugs, features etc.



'master' branch

Create 'feature' branch from 'master'      Merge 'feature' branch into 'master'

Commit changes      Submit Pull Request      Discuss proposed changes

# Nextflow and GitHub

- Any branch, tag or commit ID can be used to specify the revision you want to execute using the
`-r` option:
```
nextflow run nextflow-io/hello -r mybranch
nextflow run nextflow-io/hello -r v1.1
```

- Be careful with docker version used:

  - Docker also has `tags` and hash IDs

  - If well done, a particular pipeline version will point to the right container automatically using the config file

  - If not you can manually specify the docker tag using:
  ```
  nextflow run -with-docker user/repo:tag
  ```

# Social coding

## The histogram would be much nicer in pink #1

[Edit] [New issue]

⊘ Open  **mfoll** opened this issue a minute ago · 1 comment

**mfoll** commented a minute ago  [Member]  +☺ ✏

*No description provided.*

🏷 👤 **mfoll** added the `enhancement` label a minute ago

👤 👤 **mfoll** self-assigned this a minute ago

**mfoll** commented 33 seconds ago  [Member]  +☺ ✏ ✕

Thanks for your feedback, this will be implemented in v1.1

**Assignees** ⚙
👤 mfoll

**Labels** ⚙
`enhancement`

**Projects**
None yet

**Milestone** ⚙
No milestone

🔖 👤 **mfoll** added a commit that referenced this issue a minute ago

○ 👤 changed hist color to pink ⋯                                    c75e7e0

🚫 👤 **mfoll** closed this just now

---

## changed hist color to pink                                    **Browse files**
Asked in #1

⑂ **master**

---

👤 **mfoll** committed on **GitHub** 13 seconds ago        1 parent e89aba0   commit c75e7e0a678905530fc71a9f02bdf4446da42d82

---

📄 Showing **1 changed file** with **1 addition** and **1 deletion**.        Unified  Split

---

2 ⬛⬛⬜⬜⬜ plot_coverage.nf                                                    View

⚕ @@ -57,7 +57,7 @@ process plot {

| 57 |   | ''' | 57 |   | ''' |
| 58 |   | #!/usr/bin/env Rscript | 58 |   | #!/usr/bin/env Rscript |
| 59 |   | pdf("coverage.pdf") | 59 |   | pdf("coverage.pdf") |
| 60 | - | hist(read.table("all_average.txt")[,1]) | 60 | + | hist(read.table("all_average.txt")[,1],col="pink") |
| 61 |   | dev.off() | 61 |   | dev.off() |
| 62 |   | ''' | 62 |   | ''' |
| 63 |   | } | 63 |   | } |

**International Agency for Research on Cancer**

🌐 **World Health Organization**

🏷 v1.1
⌖ c75e7e0

# v1.1

👤 **mfoll** released this 2 minutes ago

Changed histogram color to pink

## Downloads

📄 **Source code** (zip)

📄 **Source code** (tar.gz)

🏷 v1.0
⌖ e89aba0

# v1.0

👤 **mfoll** released this 7 minutes ago · **1 commit** to master since this release

Update nextflow.config

## Downloads
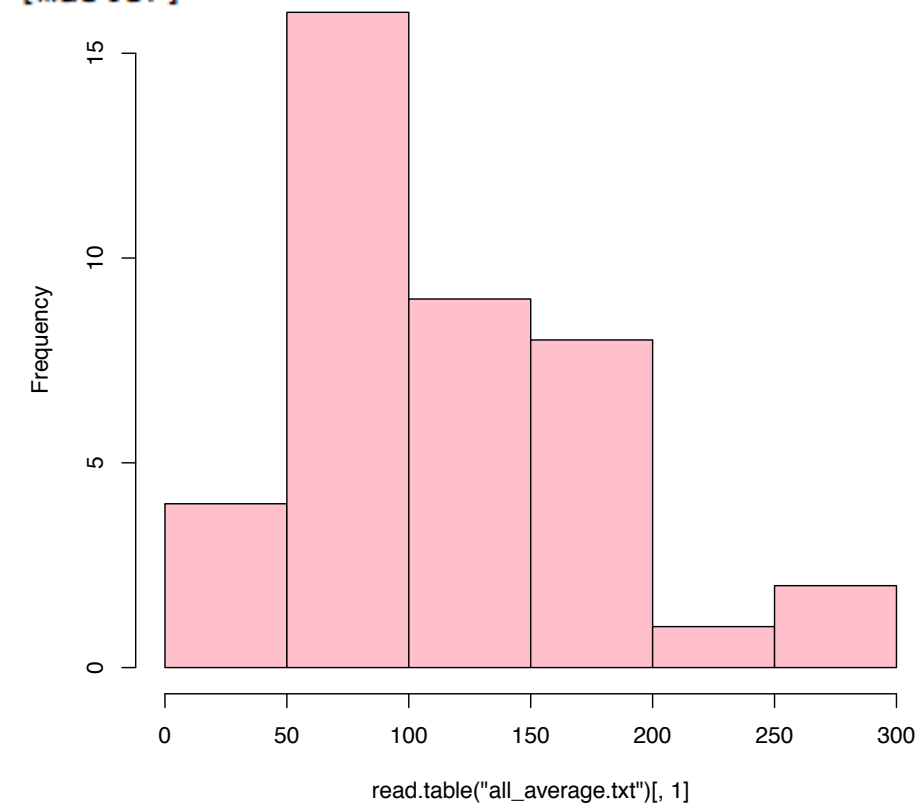
📄 **Source code** (zip)

📄 **Source code** (tar.gz)

# Running specific versions

[x140083:nf_coverage_demo follm$ nextflow run iarcbioinfo/nf_coverage_demo -latest --bam_folder BAM/ --bed TP53.bed
N E X T F L O W ~ version 0.23.4
Pulling iarcbioinfo/nf_coverage_demo ...
 Fast-forward
Launching `iarcbioinfo/nf_coverage_demo` [tiny_panini] - revision: c75e7e0a67 [master]
[warm up] executor > local
[88/be7ca4] Submitted process > coverage (2)
[68/25f37e] Submitted process > coverage (7)

**Histogram of read.table("all_average.txt")[, 1]**



Frequency

read.table("all_average.txt")[, 1]

[x140083:nf_coverage_demo follm$ nextflow run iarcbioinfo/nf_coverage_demo -r v1.0 --bam_folder BAM/ --bed TP53.bed
N E X T F L O W ~ version 0.23.4
Launching `iarcbioinfo/nf_coverage_demo` [amazing_almeida] - revision: e89aba0637 [v1.0]
[warm up] executor > local
[8a/1ffa6b] Submitted process > coverage (8)
[37/221eee] Submitted process > coverage (3)
[d8/7199ee] Submitted process > coverage (5)

# Nextflow outputs

- stdout (what you see on the screen)

- work folder

- log file: `.nextflow.log`

- history: `.nextflow/history`

- Optional:

  - trace

  - timeline

  - report

  - email notification

# History and log

```
1   2017-03-03 10:39:57 2s   backstabbing_shockley   OK   dd542716cd16e9973b69b8b5d76f4d70   7ccf3989-b998-4c34-a2d4-ba5235c4ac83   nextflow run plot_coverage_1.nf --bam_folder BAM/ --bed TP53.bed
2   2017-03-03 10:40:43 2.7s compassionate_bohr  OK   dd542716cd16e9973b69b8b5d76f4d70   c921a9c9-71b5-4279-831d-035da53f0f5d   nextflow run plot_coverage_1.nf --bam_folder BAM/ --bed TP53.bed
3   2017-03-03 10:55:59 3.5s reverent_jones OK   66d7be595f210c609f1629809093c5ad   316692d7-bebb-42e9-911c-a13d90624d95   nextflow run plot_coverage_2.nf --bam_folder BAM/ --bed TP53.bed
4   2017-03-03 10:56:11 3.2s modest_shirley OK   6f32c834925fa85243ce5a62cdb04703   9705a819-95f8-46d8-be5c-2e5414e1f9a9   nextflow run plot_coverage_2.nf --bam_folder BAM/ --bed TP53.bed
5   2017-03-03 11:07:01 3.7s nauseous_bose  OK   66d7be595f210c609f1629809093c5ad   10519e2d-0cf2-481a-91f1-70fa27b1198a   nextflow run plot_coverage.nf --bam_folder BAM/ --bed TP53.bed
6   2017-03-03 11:13:01 2.9s high_blackwell ERR  66d7be595f210c609f1629809093c5ad   37d8599c-f6b1-4743-ac52-7d183d6df682   nextflow run plot_coverage.nf --bam_folder BAM/ --bed TP53.bed
```

```
1    Mar-03 11:16:51.945 [main] DEBUG nextflow.cli.Launcher - S> /usr/local/bin/nextflow run plot_coverage.nf -with-trace -with-timeline --bam_folder BAM/ --bed TP53.bed
2    Mar-03 11:16:52.039 [main] INFO  nextflow.cli.CmdRun - N E X T F L O W  ~  version 0.23.4
3    Mar-03 11:16:52.048 [main] INFO  nextflow.cli.CmdRun - Launching `plot_coverage.nf` [cranky_leavitt] - revision: 66d7be595f
4    Mar-03 11:16:52.318 [main] DEBUG nextflow.Session - Session uuid: a88b516d-6f44-4041-b9f2-82e8acR4b17e
5    Mar-03 11:16:52.318 [main] DEBUG nextflow.Session - Run name: cranky_leavitt
6    Mar-03 11:16:52.320 [main] DEBUG nextflow.Session - Executor pool size: 8
7    Mar-03 11:16:52.338 [main] DEBUG nextflow.cli.CmdRun -
8      Version: 0.23.4 build 4170
9      Modified: 24-02-2017 09:38 UTC (10:38 CEST)
10     System: Mac OS X 10.12.1
11     Runtime: Groovy 2.4.7 on Java HotSpot(TM) 64-Bit Server VM 1.8.0_51-b16
12     Encoding: UTF-8 (UTF-8)
13     Process: 30522@x140883.local [10.10.13.47]
14     CPUs: 8 - Mem: 16 GB (47.1 MB) - Swap: 1 GB (938.2 MB)
15   Mar-03 11:16:52.345 [main] DEBUG nextflow.Session - Work-dir: /Users/follm/nf_coverage_demo/work [Mac OS X]
16   Mar-03 11:16:52.345 [main] DEBUG nextflow.Session - Script base path does not exist or is not a directory: /Users/follm/nf_coverage_demo/bin
17   Mar-03 11:16:52.428 [main] DEBUG nextflow.Session - Session start invoked
18   Mar-03 11:16:52.434 [main] DEBUG nextflow.processor.TaskDispatcher - Dispatcher > start
19   Mar-03 11:16:52.435 [main] DEBUG nextflow.trace.TraceFileObserver - Flow starting -- trace file: /Users/follm/nf_coverage_demo/trace.txt
20   Mar-03 11:16:52.438 [main] DEBUG nextflow.script.ScriptRunner - > Script parsing
21   Mar-03 11:16:52.536 [main] DEBUG nextflow.script.ScriptRunner - > Launching execution
22   Mar-03 11:16:52.555 [main] DEBUG nextflow.Channel - files for syntax: glob; folder: BAM/; pattern: *.bam; options: null
23   Mar-03 11:16:52.629 [main] DEBUG nextflow.processor.ProcessFactory - << taskConfig executor: null
24   Mar-03 11:16:52.629 [main] DEBUG nextflow.processor.ProcessFactory - >> processorType: 'local'
25   Mar-03 11:16:52.634 [main] DEBUG nextflow.executor.Executor - Initializing executor: local
26   Mar-03 11:16:52.636 [main] INFO  nextflow.executor.Executor - [warn up] executor > local
27   Mar-03 11:16:52.641 [main] DEBUG n.processor.LocalPollingMonitor - Creating local task monitor for executor 'local' > cpus=8; memory=16 GB; capacity=8; pollInterval=100ms; dumpInterval=5m
28   Mar-03 11:16:52.644 [main] DEBUG nextflow.processor.TaskDispatcher - Starting monitor: LocalPollingMonitor
29   Mar-03 11:16:52.644 [main] DEBUG n.processor.TaskPollingMonitor - >>> barrier register (monitor: local)
30   Mar-03 11:16:52.646 [main] DEBUG nextflow.executor.Executor - Invoke register for executor: local
31   Mar-03 11:16:52.694 [main] DEBUG nextflow.Session - >>> barrier register (process: coverage)
32   Mar-03 11:16:52.703 [main] DEBUG nextflow.processor.TaskProcessor - Creating operator > coverage -- maxForks: 8
33   Mar-03 11:16:52.722 [main] DEBUG nextflow.processor.ProcessFactory - << taskConfig executor: null
34   Mar-03 11:16:52.723 [main] DEBUG nextflow.processor.ProcessFactory - >> processorType: 'local'
35   Mar-03 11:16:52.723 [main] DEBUG nextflow.executor.Executor - Initializing executor: local
36   Mar-03 11:16:52.723 [main] DEBUG nextflow.Session - >>> barrier register (process: mean)
37   Mar-03 11:16:52.729 [main] DEBUG nextflow.processor.TaskProcessor - Creating operator > mean -- maxForks: 8
38   Mar-03 11:16:52.850 [main] DEBUG nextflow.processor.ProcessFactory - << taskConfig executor: null
39   Mar-03 11:16:52.850 [main] DEBUG nextflow.processor.ProcessFactory - >> processorType: 'local'
40   Mar-03 11:16:52.850 [main] DEBUG nextflow.executor.Executor - Initializing executor: local
41   Mar-03 11:16:52.851 [main] DEBUG nextflow.Session - >>> barrier register (process: plot)
42   Mar-03 11:16:52.864 [main] DEBUG nextflow.processor.TaskProcessor - Creating operator > plot -- maxForks: 8
43   Mar-03 11:16:52.865 [main] DEBUG nextflow.script.ScriptRunner - > Await termination
44   Mar-03 11:16:52.865 [main] DEBUG nextflow.Session - Session await
45   Mar-03 11:16:52.990 [Actor Thread 4] INFO  nextflow.Session - [cd/081600] Submitted process > coverage (2)
46   Mar-03 11:16:52.990 [Actor Thread 5] INFO  nextflow.Session - [5e/f753d5] Submitted process > coverage (3)
47   Mar-03 11:16:52.990 [Actor Thread 9] INFO  nextflow.Session - [9c/0cda8b] Submitted process > coverage (7)
48   Mar-03 11:16:52.990 [Actor Thread 8] INFO  nextflow.Session - [4d/8bd5f0] Submitted process > coverage (6)
49   Mar-03 11:16:52.991 [Actor Thread 3] INFO  nextflow.Session - [31/bf2185] Submitted process > coverage (1)
```

# Report (-with-report)

Nextflow workflow report

[ecstatic_booth]

Workflow execution completed successfully!

**Run times**

Tue May 22 18:31:22 CEST 2018 - Tue May 22 18:33:27 CEST 2018 (completed 3 hours ago, duration: **2m 6s**)

187 succeeded

**Nextflow command**

```
nextflow run iarcbioinfo/nf_coverage_demo -latest -r v1.2 -with-docker --bam_folder data_test/BAM/BAM_multiple/ --bed
data_test/BED/TP53_exon2_11.bed -with-report
```

| | |
|---|---|
| **CPU-Hours** | 0.1 |
| **Launch directory** | /Users/follm |
| **Work directory** | /Users/follm/work |
| **Project directory** | /Users/follm/.nextflow/assets/iarcbioinfo/nf_coverage_demo |
| **Script name** | plot_coverage.nf |
| **Script ID** | 1bfa2de3de08d52bac30fdde640551c8 |
| **Workflow session** | 1239Gc92-954e-4f90-B7eG-35547d3f77ec |
| **Workflow repository** | https://github.com/IARCbioinfo/nf_coverage_demo.git , revision v1.2 (commit hash 19c019913af89b4b8f6d93d115b136ef23cd0d74 ) |
| **Workflow profile** | standard |
| **Workflow container** | iarcbioinfo/nf_coverage_demo |
| **Container engine** | docker |
| **Nextflow version** | version 0.29.1, build 4804 (10-05-2018 07:47 UTC) |

International

World
Organ

# Trace (-with-trace)

| | task_id | hash | native_id | name | status | exit | submit | duration | realtime | %cpu | rss | vmem | rchar | wchar |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | | | | | | | | | | | | | | |
| 2 | 2 | cd/081600 | 30542 | coverage (2) | COMPLETED | 0 | 16:53.0 | 419ms | 225ms | 11.40% | 7.3 MB | 4.7 GB | 0 | 0 |
| 3 | 4 | 9e/0cda8b | 30543 | coverage (7) | COMPLETED | 0 | 16:53.0 | 446ms | 397ms | - | - | - | - | - |
| 4 | 3 | 5e/f753d5 | 30546 | coverage (3) | COMPLETED | 0 | 16:53.0 | 459ms | 411ms | - | - | - | - | - |
| 5 | 6 | dc/087bc6 | 30541 | coverage (5) | COMPLETED | 0 | 16:53.0 | 467ms | 417ms | - | - | - | - | - |
| 6 | 5 | c6/83374d | 30539 | coverage (4) | COMPLETED | 0 | 16:53.0 | 550ms | 283ms | 2.30% | 6.6 MB | 4.7 GB | 0 | 0 |
| 7 | 1 | 31/bf2185 | 30544 | coverage (1) | COMPLETED | 0 | 16:53.0 | 641ms | 591ms | - | - | - | - | - |
| 8 | 8 | d7/18814e | 30545 | coverage (8) | COMPLETED | 0 | 16:53.0 | 695ms | 646ms | - | - | - | - | - |
| 9 | 11 | f0/67de05 | 30949 | coverage (11 | COMPLETED | 0 | 16:53.4 | 240ms | 227ms | - | - | - | - | - |
| 10 | 14 | 2e/dce17f | 30928 | coverage (14 | COMPLETED | 0 | 16:53.4 | 352ms | 80ms | 0.70% | 1 MB | 2.3 GB | 0 | 0 |
| 11 | 9 | 32/05246f | 31055 | coverage (9) | COMPLETED | 0 | 16:53.5 | 342ms | 84ms | 0.70% | 1 MB | 2.3 GB | 0 | 0 |
| 12 | 16 | 2c/01e5ea | 31195 | coverage (16 | COMPLETED | 0 | 16:53.7 | 276ms | 127ms | 0.70% | 8.6 MB | 4.7 GB | 0 | 0 |
| 13 | 17 | a5/8e5a2e | 31288 | coverage (17 | COMPLETED | 0 | 16:53.8 | 231ms | 127ms | - | - | - | - | - |
| 14 | 12 | ee/c77b04 | 31191 | coverage (12 | COMPLETED | 0 | 16:53.7 | 340ms | 95ms | 0.70% | 1 MB | 2.3 GB | 0 | 0 |
| 15 | 18 | 5d/0f162d | 31449 | mean (1) | COMPLETED | 0 | 16:53.9 | 158ms | 74ms | - | - | - | - | - |
| 16 | 15 | 72/ddf049 | 31146 | coverage (15 | COMPLETED | 0 | 16:53.6 | 420ms | 219ms | 10.80% | 9.6 MB | 4.7 GB | 0 | 0 |
| 17 | 10 | 95/26c5e0 | 30954 | coverage (10 | COMPLETED | 0 | 16:53.5 | 634ms | 542ms | - | - | - | - | - |
| 18 | 19 | 65/d41f91 | 31517 | mean (2) | COMPLETED | 0 | 16:54.0 | 176ms | 142ms | - | - | - | - | - |
| 19 | 21 | 3b/800cd4 | 31566 | mean (3) | COMPLETED | 0 | 16:54.0 | 181ms | 161ms | - | - | - | - | - |
| 20 | 23 | 0e/e9abc8 | 31599 | mean (4) | COMPLETED | 0 | 16:54.1 | 173ms | 127ms | - | - | - | - | - |
| 21 | 20 | 4d/b7bb4e | 31534 | coverage (18 | COMPLETED | 0 | 16:54.0 | 269ms | 237ms | - | - | - | - | - |
| 22 | 22 | 36/e28f53 | 31583 | coverage (19 | COMPLETED | 0 | 16:54.0 | 231ms | 217ms | - | - | - | - | - |
| 23 | 7 | 4d/8bd5f0 | 30540 | coverage (6) | COMPLETED | 0 | 16:53.0 | 1.4s | 358ms | 2.50% | 5.2 MB | 4.7 GB | 0 | 0 |
| 24 | 28 | ac/c45043 | 31852 | mean (6) | COMPLETED | 0 | 16:54.3 | 169ms | 157ms | - | - | - | - | - |
| 25 | 25 | de/2ed75c | 31715 | coverage (21 | COMPLETED | 0 | 16:54.1 | 346ms | 188ms | 10.60% | 10.5 MB | 4.7 GB | 0 | 0 |
| 26 | 24 | e0/450f34 | 31649 | coverage (20 | COMPLETED | 0 | 16:54.1 | 438ms | 386ms | - | - | - | - | - |
| 27 | 27 | ff/c08e08 | 31807 | coverage (22 | COMPLETED | 0 | 16:54.2 | 307ms | 153ms | 0.00% | 10.7 MB | 4.7 GB | 0 | 0 |
| 28 | 31 | a0/337b22 | 32049 | mean (8) | COMPLETED | 0 | 16:54.4 | 295ms | 242ms | - | - | - | - | - |
| 29 | 34 | d5/04b9a6 | 32117 | mean (9) | COMPLETED | 0 | 16:54.5 | 204ms | 104ms | - | - | - | - | - |
| 30 | 33 | a1/b9754f | 32116 | coverage (25 | COMPLETED | 0 | 16:54.5 | 246ms | 144ms | - | - | - | - | - |
| 31 | 13 | 1a/2bbf4f | 30935 | coverage (13 | COMPLETED | 0 | 16:53.4 | 1.4s | 1.2s | 100.70% | 2.7 MB | 4.7 GB | 0 | 0 |

International Agency for Research on Cancer

World Health Organization

# Practical

- Run pipelines on the cluster using Singularity `-with-singularity`

- Explore all nextflow outputs/logs, including in the `work` directory

- Try to run without docker/singularity

- Try to run v1.0, still with docker/singularity

- Try different combinations of pipeline version and container version

- Look at the help of the `nextflow` commands to find useful options (in particular `run -latest, -qs`, `-bg`)

- Look at the `clone`, `drop`, `list` and `pull` commands to manage pipelines

- Try to modify a pipeline (feel free to send me a PR!)

- Check if a new version of nextflow itself is available using `nextflow self-update`

# Dataflow programming

- Traditionally, a program is a series of operations happening in a specific order ("sequential programming")

- Dataflow programming emphasizes the movement of data and models programs as a series of connections:

  - Explicitly defined inputs and outputs connect operations, which function like black boxes.

  - An operation runs as soon as all of its inputs become valid. Thus, dataflow languages are inherently parallel.

https://en.wikipedia.org/wiki/Dataflow_programming

# Our first process

```
1
2    // Defines pipeline parameters
3    params.bam_folder = null
4    params.bed = null
5
6    // The bed file
7    bed = file(params.bed)
8
9    // Creates the `bam` channel
10   bam = Channel.fromPath( params.bam_folder+'/*.bam' )
11
12   // Step 1. launch bedtools software to calculate coverage at each position of the bed
13 ▼ process coverage {
14
15       input:
16       file bam
17       file bed
18
19       output:
20       file 'coverage.txt' into coverage
21
22       shell:
23       '''
24       bedtools coverage -d -a !{bed} -b !{bam} > coverage.txt
25       '''
26 ⌐ }
```

# Running the pipeline

# Adding a second process

```
19        output:
20        file 'coverage.txt' into coverage
21
22        shell:
23        '''
24        bedtools coverage -d -a !{bed} -b !{bam} > coverage.txt
25        '''
26 ⌐  }
27
28    // Step 2. launch custom awk script to calculate the mean coverage
29 ▼  process mean {
30
31        input:
32        file coverage
33
34        output:
35        stdout average
36
37        shell:
38        '''
39        awk '{ sum += $6 } END { if (NR > 0) print sum / NR }' !{coverage}
40        '''
41 ⌐  }
```

# Running the pipeline



```
[x140083:nf_coverage_demo follm$ nextflow run plot_coverage_2.nf --bam_folder BAM/ --bed TP53.bed
N E X T F L O W  ~  version 0.23.4
Launching `plot_coverage_2.nf` [modest_shirley] - revision: 6f32c83402
[warm up] executor > local
[f6/cc4a86] Submitted process > coverage (7)
[e7/1c6c05] Submitted process > coverage (1)
[1f/08005d] Submitted process > coverage (4)
[37/0db07c] Submitted process > coverage (8)
[ba/7092ed] Submitted process > coverage (3)
[1b/e48edf] Submitted process > coverage (2)
[e3/c18377] Submitted process > coverage (12)
[89/fde6b3] Submitted process > coverage (5)
[d3/456b82] Submitted process > coverage (6)
[b6/59320b] Submitted process > coverage (11)
[c9/319f64] Submitted process > coverage (10)
[e5/be57ba] Submitted process > coverage (9)
[9a/ed4a42] Submitted process > coverage (13)
[28/243d0a] Submitted process > coverage (14)
[77/6d9dee] Submitted process > coverage (15)
[c4/2c85af] Submitted process > coverage (16)
[17/f2b669] Submitted process > coverage (17)
[a8/e2f0a9] Submitted process > coverage (18)
[87/7e3a9b] Submitted process > coverage (19)
[c5/7489dd] Submitted process > mean (1)
[25/8f2f62] Submitted process > mean (2)
[56/d40ff8] Submitted process > mean (3)
[be/02c2d1] Submitted process > mean (4)
[ce/0d34de] Submitted process > coverage (21)
```

~/nf_coverage_demo/work/c5/7489ddc158fe4aecc3eeeb4f6495e3/.command.sh

```
1  #!/bin/bash -ue
2  awk '{ sum += $6 } END { if (NR > 0) print sum / NR }' coverage.txt
3
```

▼ ☐ c5
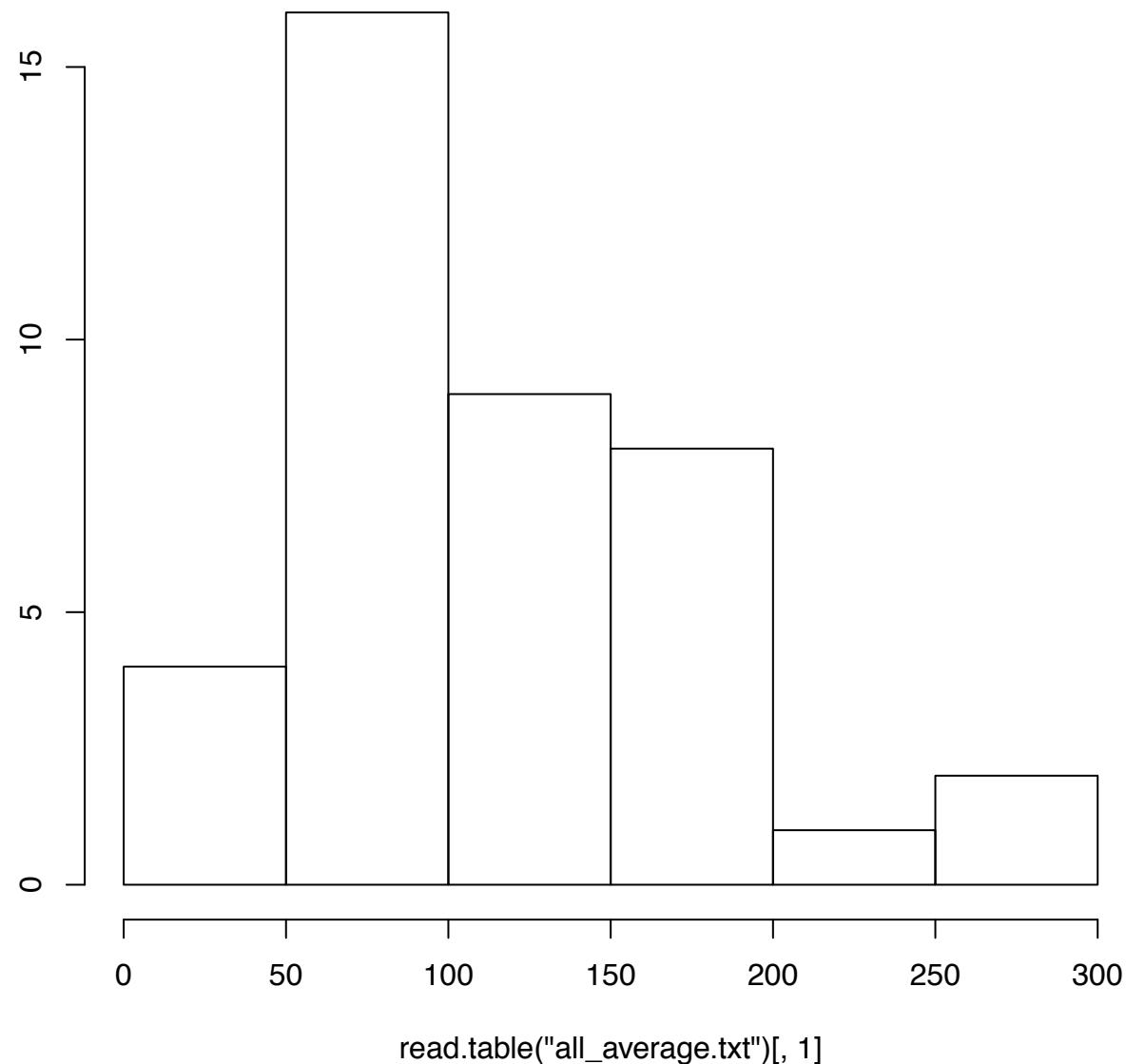　　▼ 📁 7489ddc158fe4aecc3eeeb4f6495e3
　　　　📄 coverage.txt
　　　　📄 .exitcode
　　　　📄 .command.sh
　　　　📄 .command.run
　　　　📄 .command.out
　　　　📄 .command.log
　　　　📄 .command.err
　　　　📄 .command.begin

~/nf_coverage_demo/work/c5/7489ddc158fe4aecc3eeeb4f6495e3/.command.out

```
1    61.4442
2
```

# The final process

```
34        output:
35        stdout average
36
37        shell:
38        '''
39        awk '{ sum += $6 } END { if (NR > 0) print sum / NR }' !{coverage}
40        '''
41   ⌐  }
42      // collect output of all means to a single file
43      all_average = average.collectFile(name: 'all_average.txt')
44
45      // Step 3: plot histogram of mean coverage using a custom R script
46   ▾  process plot {
47
48        input:
49        file all_average
50
51        output:
52        file 'coverage.pdf'
53
54        publishDir '.', mode: 'move'
55
56        shell:
57        '''
58        #!/usr/bin/env Rscript
59        pdf("coverage.pdf")
60        hist(read.table("all_average.txt")[,1])
61        dev.off()
62        '''
63   ⌐  }
```

This is our final result

# Running the pipeline

# Result

**Histogram of read.table("all_average.txt")[, 1]**



read.table("all_average.txt")[, 1]

```
--/nf_coverage_demo/plot_coverage.nf ▾
1
2      // Defines pipeline parameters
3      params.bam_folder = null
4      params.bed = null
5
6      // The bed file
7      bed = file(params.bed)
8
9      // Creates the 'bam' channel
10     bam = Channel.fromPath( params.bam_folder+'/*.bam' )
11
12     // Step 1. launch bedtools software to calculate coverage at each position of the bed
13  ▾  process coverage {
14
15         input:
16         file bam
17         file bed
18
19         output:
20         file 'coverage.txt' into coverage
21
22         shell:
23         '''
24         bedtools coverage -d -a !{bed} -b !{bam} > coverage.txt
25         '''
26  ∟  }
27
28     // Step 2. launch custom awk script to calculate the mean coverage
29  ▾  process mean {
30
31         input:
32         file coverage
33
34         output:
35         stdout average
36
37         shell:
38         '''
39         awk '{ sum += $6 } END { if (NR > 0) print sum / NR }' !{coverage}
40         '''
41  ∟  }
42     // collect output of all means to a single file
43     all_average = average.collectFile(name: 'all_average.txt')
44
45     // Step 3: plot histogram of mean coverage using a custom R script
46  ▾  process plot {
47
48         input:
49         file all_average
50
51         output:
52         file 'coverage.pdf'
53
54         publishDir '.', mode: 'move'
55
56         shell:
57         '''
58         #!/usr/bin/env Rscript
59         pdf("coverage.pdf")
60         hist(read.table("all_average.txt")[,1])
61         dev.off()
62         '''
63  ∟  }
```
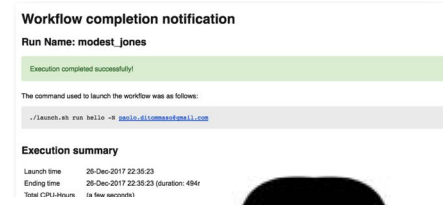
# Full life cycle



**Nextflow** runs each *process* of the pipeline inside the container on the available *executors* (workstation, HPC, cloud)

**4**

**3 Nextflow** downloads the docker container with all the *software* from DockerHub

**5** If circleci checks have passed, a new **docker** container is build and hosted on DockerHub
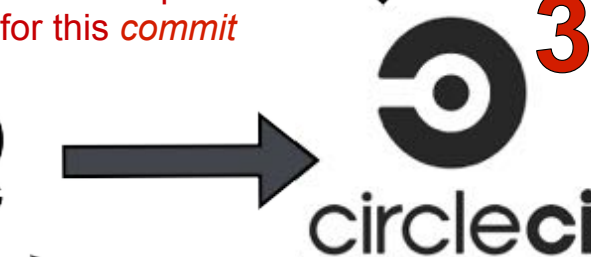
**5 Nextflow** sends an execution *report*

**1** A **user** runs the pipeline using *nextflow*. If something goes wrong he opens an *issue*.

**2 Nextflow** downloads the pipeline from github

**2 Github** tracks the *changes* and creates a unique identifier for this *commit*

**3 Circleci** starts running *tests* in the cloud to check if the changes haven't broken something and are still producing expected results on some test datasets

**1** A **developer** modifies the pipeline and pushes the changes to github (*commit*).

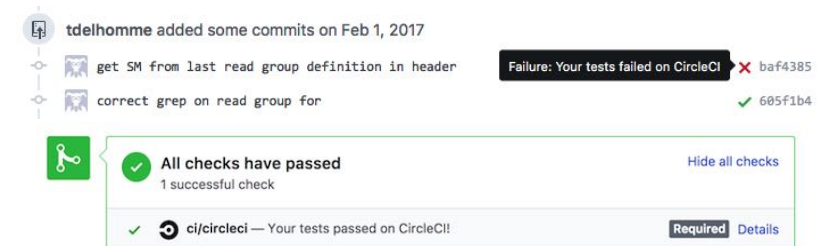This is the only manual action he has to perform as all the following steps are triggered automatically

**4 Circleci** sends feedback on the pipeline github webpage about the results of the test

# Final word

Don't try to twist nextflow to fit your current habits, give it a chance to actually change your habits