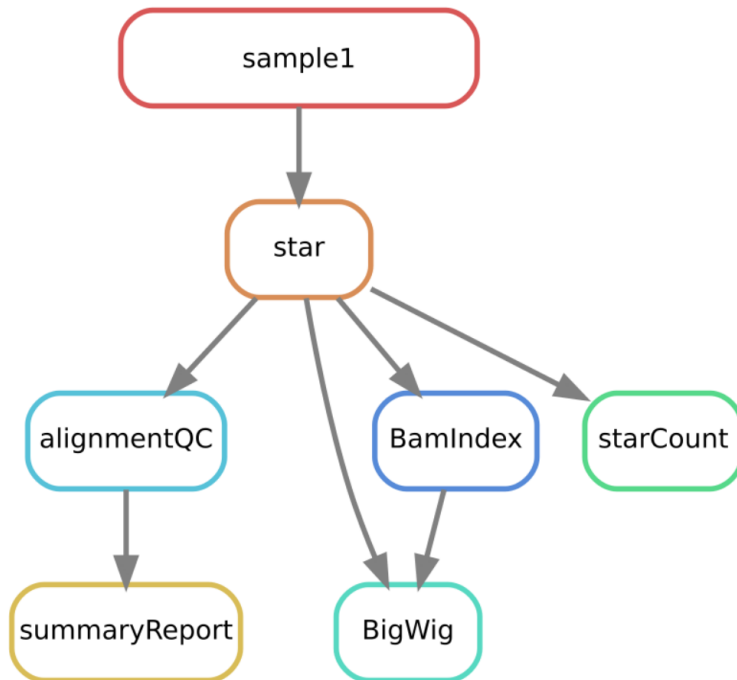


# Reproducible analysis using workflow managers

# Why using a workflow manager?

**One analysis** : multiple steps, multiple samples

Simple example, one sample



nextflow

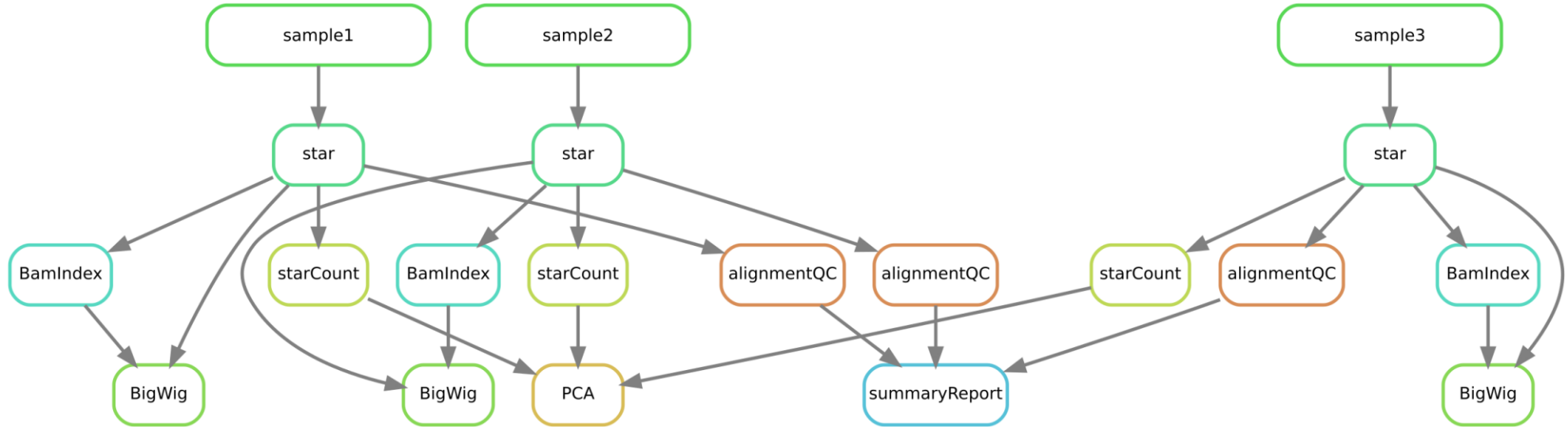


Some tasks can run **in parallel**

Some tasks have to **wait** for one or several others to complete

# Why using a workflow manager?

Simple example, more samples (3)

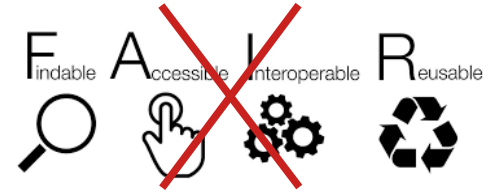


Some tasks are done **for each sample**

Some **combine** the information from several samples

→ Usually more samples, without workflow manager you may get **lost...**

# Why using a workflow manager?



And/or end up with something like that...

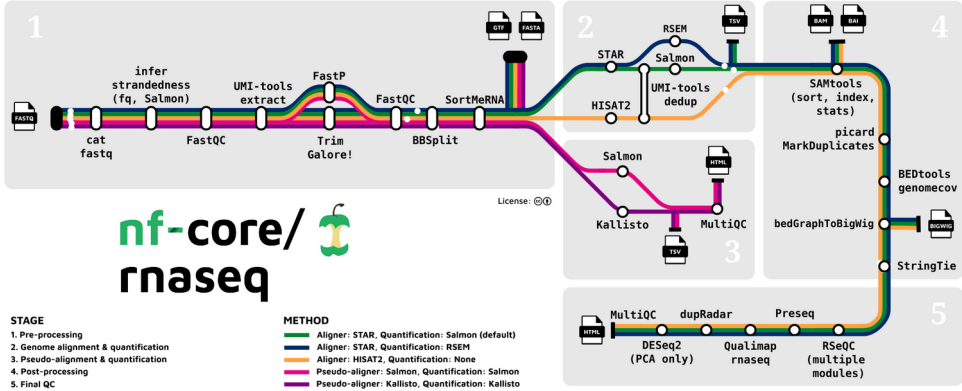
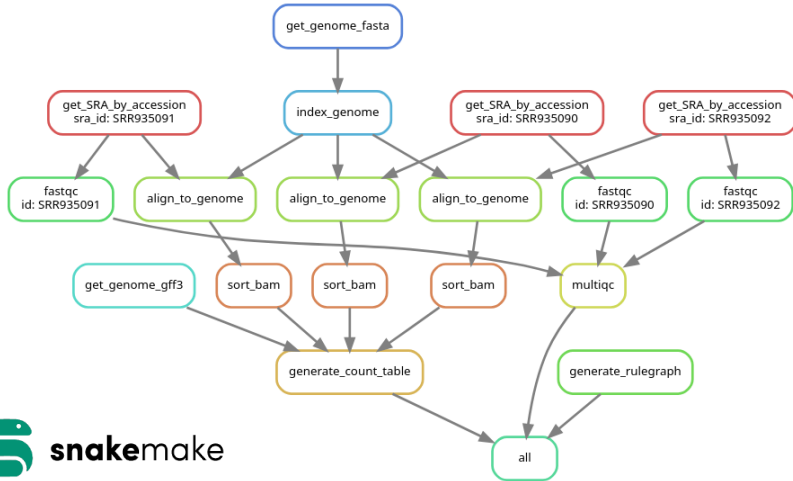
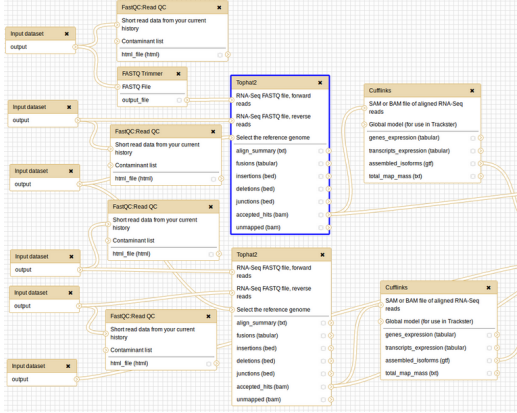
|                  | one               |    |                | safe           |    |                  | s/v2 |                |    |  |
|------------------|-------------------|----|----------------|----------------|----|------------------|------|----------------|----|--|
|                  | B1                | B2 | B3             | B4             | B5 | B6               | C1   | C2             | C3 |  |
| QC               | ✓                 | ✓  | <del>✓</del> ✓ | <del>✓</del> ✓ | ✓  | ✓                | ✓    | ✓              | ✓  |  |
| Mapping<br>23/01 | ✓<br>23/01        | ✓  | ✓              | <del>✓</del> ✓ | ✓  | ✓ <sup>2-6</sup> | ✓    | ✓              | ✓  |  |
| bw (Feb)         | ✓                 | ✓  | ✓              | ✓              | ✓  | <del>✓</del> ✓   | ✓    | <del>✓</del> ✓ | ✓  |  |
| Count            | ✓<br><del>✓</del> | ✓  | ✓              | ✓              | ✓  | ✓                | ✓    | ✓<br>OK        | ✓  |  |
| Edger            |                   |    |                | OK             |    |                  |      |                |    |  |
|                  |                   |    |                | summary        |    |                  | OK   |                |    |  |

# The reproducible journey

How to make my analysis reproducible in a few years / by someone else ?



# Workflow managers



# Pros

**Reproducible** and **portable** analysis, controlled environment

Task parallelisation -> efficient and **scalable** (HPCs)

Resuming of failed runs or steps

Modular → library of reusable blocks (any language), very efficient to **benchmark** methods and parameters

**Easy configuration** for bioinformaticians and **biologists**

# Cons

**Learning effort ...**

# Controlled environments

Most of the tools in your workflow need some requirements.

Some need lot of requirements, ie Python or R libraries in specific version.

Often, those dependencies are **not compatible** with what you already have on our system or between them.

→ need for **independant environments**

- **Conda** environments
- **Containers** : Docker, Apptainer (ex Singularity)
- **Modules**



# Conda

Automatically get packages from repositories called « channels »

**Warning** : Anaconda curated channels are **not free** !

« *Utilizing Miniconda to pull package updates from the Anaconda Public Repository without a commercial licence [...] is considered a violation of the Terms of Service* »

**Teaching** is the only exception, doing public research is not !

- Don't use *defaults* channel which includes *main/default* and *r*.
- Don't use Anaconda or Miniconda softwares but **Mamba** or **Miniforge**

Use open source and free channels such as **bioconda** or **conda-forge**

# Conda / Mamba

## Create a new environment

```
mamba create --no-default-packages -n myPythonEnv python=3.12.7
mamba create -n jupyterEnv -c conda-forge jupyterlab
```

yaml file to define your environment

-> define channels (order matters!) and packages

```
mamba env create -f nanopore_env.yaml
```

## Work in an environment

```
$ conda activate nanopore
$ samtools --help
```

nanopore\_env.yaml

```
name: nanopore
channels:
  - epi2melabs
  - bioconda
  - conda-forge
  - nodefaults
dependencies:
  - samtools
  - htlib=1.14
  - modbam2bed
  - epi2melabs
```

# Conda limits

## ⊖ Reproducible ?

- Environments built few months/years apart from the same recipe are **different**
- Some can't be rebuilt (packages not on the repos anymore, changes in librairies)

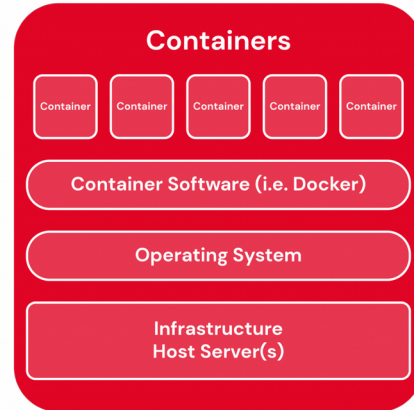
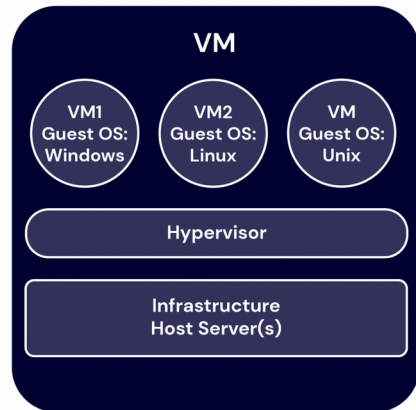
→ Use of containers solves this issue (but need to store a big file)

## ⊖ Huge **number of files** in Conda envs (typically 10-20K...)

# Containers

Container image ~ virtual computer disk with all of the necessary software, libraries and configuration to run one or more applications

!= virtual machine



A container made in one OS works only on that OS (e.g. Windows or Linux)

<https://www.backblaze.com/blog/vm-vs-containers/>

# Containers

You can use **existing containers**



Docker Hub

<https://hub.docker.com/>



**APPTAINER**

Singularity Hub (2016-2021)



<https://oras.land>



Sylabs

<https://sylabs.io/>

From Docker

```
$ apptainer pull tensorflow.sif docker://tensorflow/tensorflow:latest
```

From Shub

```
$ apptainer pull apptainer-images.sif shub://vsoch/apptainer-images
```

From supporting OCI (Open Containers Initiative) registry

```
$ apptainer pull image.sif oras://url/to/gitlab/container/registry:latest
```

# Build a container

Need to be **super user** !

## Apptainer definition file

```
BootStrap: docker
From: ubuntu

%post
# update apt
apt update
# git install
apt install -y git
# clone git repo and install
git clone --recursive
https://github.com/MBoemo/DNAscent.git
cd DNAscent
git checkout 4.0.3
make

%runscript
# Define the run command
exec echo "$@"
```

```
$ sudo apptainer build dnascent_v4.0.3.sif Singularity_4.0.3
```

## From a **base image**

1. Pull an existing image as a base
2. Install packages
3. Make the image

Singularity\_4.0.3

# Build a container

Need to be **super user** !

```
Bootstrap: docker
From: mambaorg/micromamba:1.5.9
```

```
%files
# Copy files for building
env.yaml /setupfile/env.yaml
```

```
%post
# Create conda env
cd /setupfile
micromamba env create -f env.yaml -n env_name
micromamba clean --yes --tarballs
```

```
# Export FAIR Files
micromamba list -n env_name > env_installed_packages.txt
micromamba env export --no-build -n env_name > env_export.yaml
```

```
%environment
# set Conda env bin in the PATH
export PATH=/opt/conda/envs/env_name/bin:$PATH
```

```
%runscript
# Define the run command
echo "Container was created $NOW"
exec echo "$@"
```

Apptainer\_def

## Containing a **Conda** environment

1. Pull an existing image as a base
2. Create Conda env from a yaml file defining the environment
3. Export path to conda bin
4. Define what does the « run » command

```
$ sudo apptainer build image.sif Apptainer_def
```

# Build a container

- From a **Docker** image

```
$ sudo aptainer build alpine.sif docker://alpine
```

- Online on **Sylabs** cloud (<https://cloud.sylabs.io/>)



## Singularity Container Services

Singularity Container Services makes containerization easy. Use it to build, share, and secure your performance intensive applications

```
$ singularity build -r library://demo/demo/alpine docker://alpine
INFO: Starting build...
INFO: Creating SIF file...
2.7MiB / 2.7MiB [=====] 100 % 33.1 MiB/s 0s
INFO: Build complete: library://demo/demo/alpine
```



Get Started Today for Free

 Sign In

 Sign Up



# Use a container

**pull** : get the image

**shell** : open a shell inside the container

**run** : run the runsript contained in the container

**exec** : use tools contained in the container

```
$ aptainer pull docker://ghcr.io/apptainer/lolcow
```

```
$ aptainer shell lolcow_latest.sif
```

```
Apptainer> ls /bin
```

```
$ aptainer run lolcow_latest.sif
```

```
< Fri Oct 11 11:40:34 CEST 2024 >
```

```
-----
 \   ^__^
  \  (oo)\_______
     (__)\       )\/\
        ||----w |
        ||     ||
```

```
$ aptainer exec lolcow_latest.sif cowsay moo
```

```
< moo >
```

```
-----
 \   ^__^
  \  (oo)\_______
     (__)\       )\/\
        ||----w |
        ||     ||
```

# Use a container

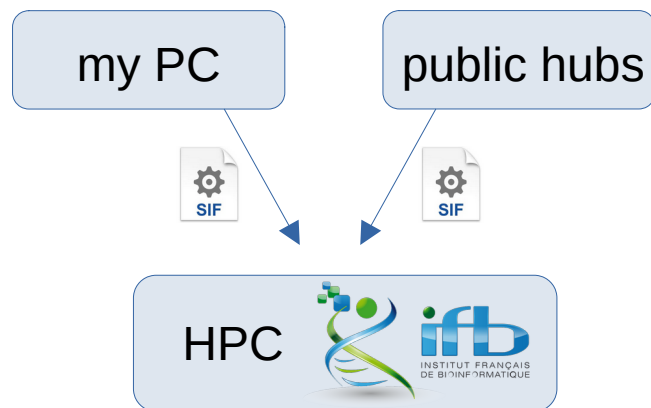
- Direct command line

```
$ aptainer shell image.sif
```

```
$ aptainer run image.sif
```

```
$ aptainer exec image.sif hostname
```

- In sbatch script
- In workflows



nextflow



# In a workflow

One container per tool

A single container for all tools

Go to <https://scrumblr.ethibox.fr/workflow> and add pros and cons of both approaches

# On the cluster : Modules



Tools ready to  
use for everyone !



↑  
`module load fastqc/0.11.9`



**BIOCONDA**<sup>®</sup>

by default



if a licence must be accepted  
if not in Bioconda and hard to integrate or urgent  
if a Docker container exists

# I need a new module...



Forum IFB

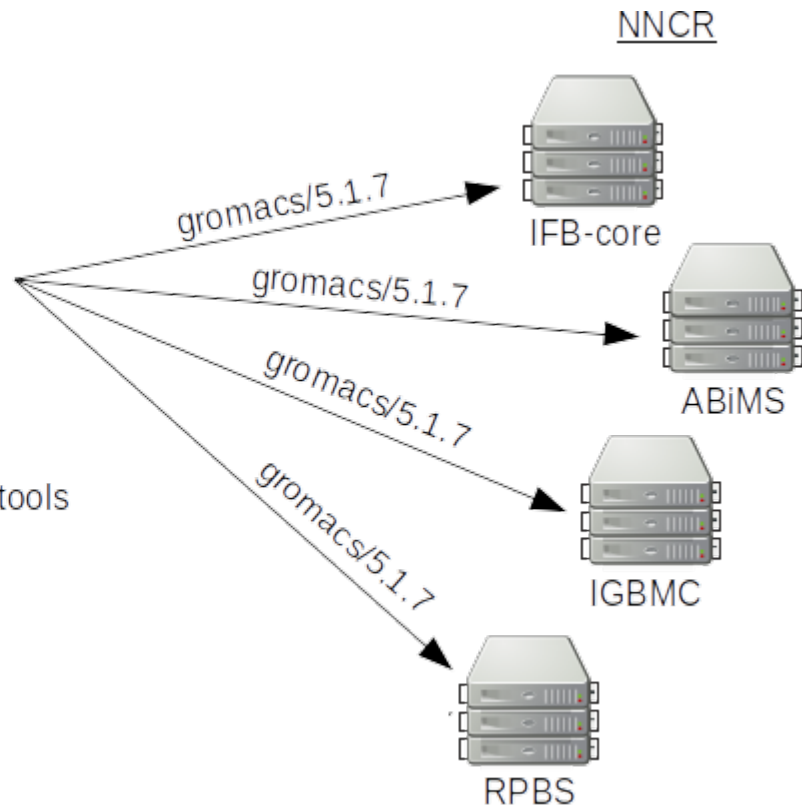
[community.france-bioinformatique.fr](http://community.france-bioinformatique.fr)

Gromacs 5.1.7 ?



Gitlab IFB

[gitlab.com/ifb-elixirfr/cluster/tools](http://gitlab.com/ifb-elixirfr/cluster/tools)



# Acknowledgement

Thanks to Laurent Jourden who initiated the container presentation !

Let's connect to IFB core cluster !