# Using Docker and Apptainer
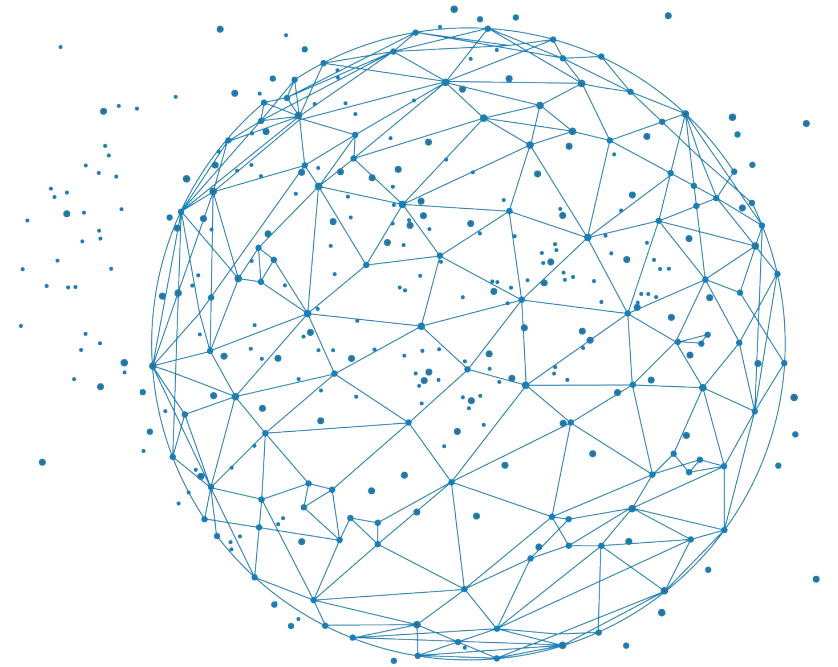
**J. Seiler**

**Based on Céline Hernandez training for FAIR Bioinfo 2023**

# Let's do an exercise first

# Exercise : putting our pipeline project on Github

- Create a Git repository for your pipeline folder
- Index and commit the following files
  - All Jupyter Notebooks
  - All R files
  - environment.yml
  - environment-linux-64.lock
- Push on Github in a new public repository

- Create a Git repository for your pipeline folder

```
$ cd pipeline
$ git init
```

- Index and commit the following files
  - All Jupyter Notebooks
  - All R files
  - environment.yml
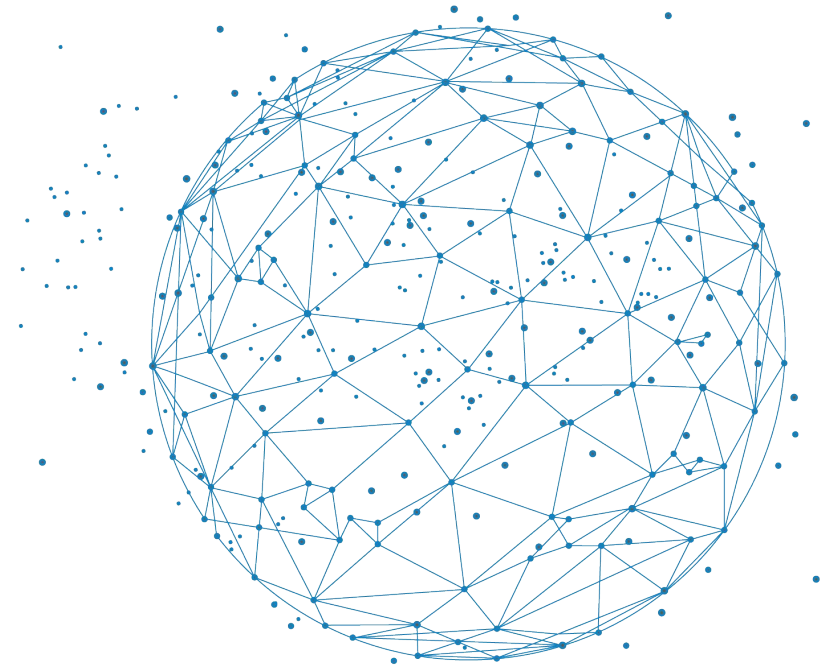  - environment-linux-64.lock

```
$ git add *.ipynb *.R environment.yml environment-linux-64.lock
$ git commit -m "Initial commit"
```

- **Push on Github in a new public repository**

```
$ git remote add origin <repo ssh url>
$ git push origin main
```

# About encapsulation

Goal : capture the system environment of applications (OS, packages, libraries,...) to control their execution.

- Hardware virtualisation (virtual machines)
- OS virtualisation (images and containers) docker
- Environment management CONDA

# Let's say we want to install RStudio...

**MacOS**

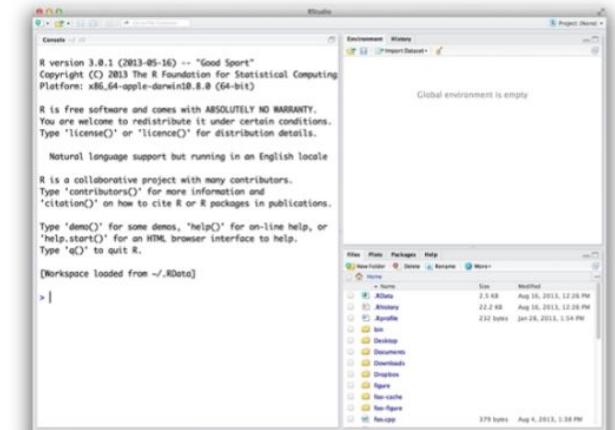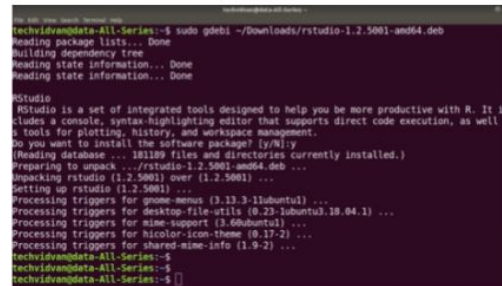**Windows**

**Install Rstudio ?**
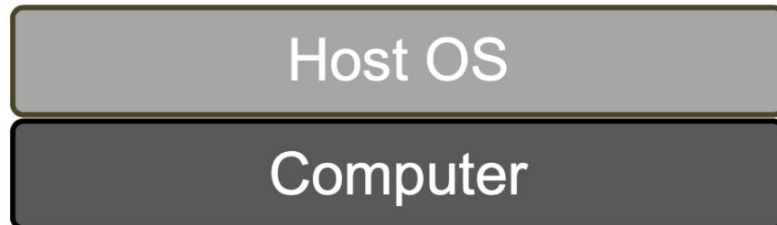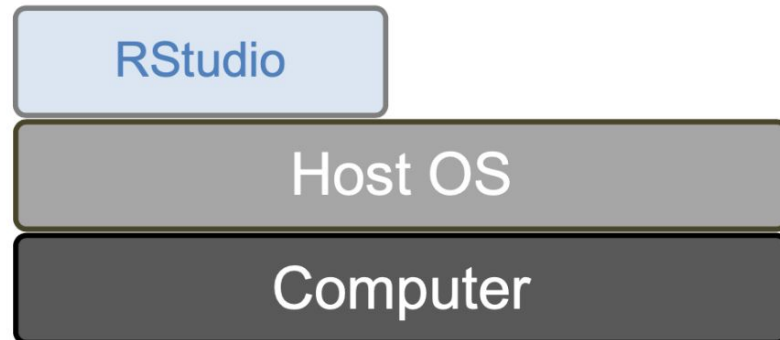
**Use Rstudio**

**Unix-based**

We started with a computer using a specific OS...

Host OS

Computer

We started with a computer using a specific OS...
And inside this environment, we installed a new application.
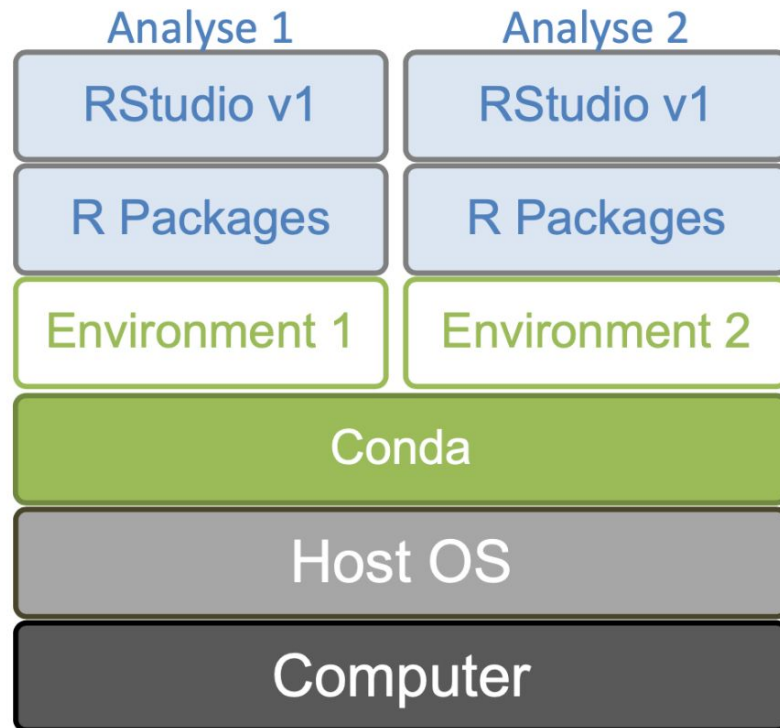
RStudio v1 | RStudio v1.2

R Packages

Host OS

Computer

Usually dependencies of different applications don't interfere.
But what if we want to test the latest version of our favourite tool?
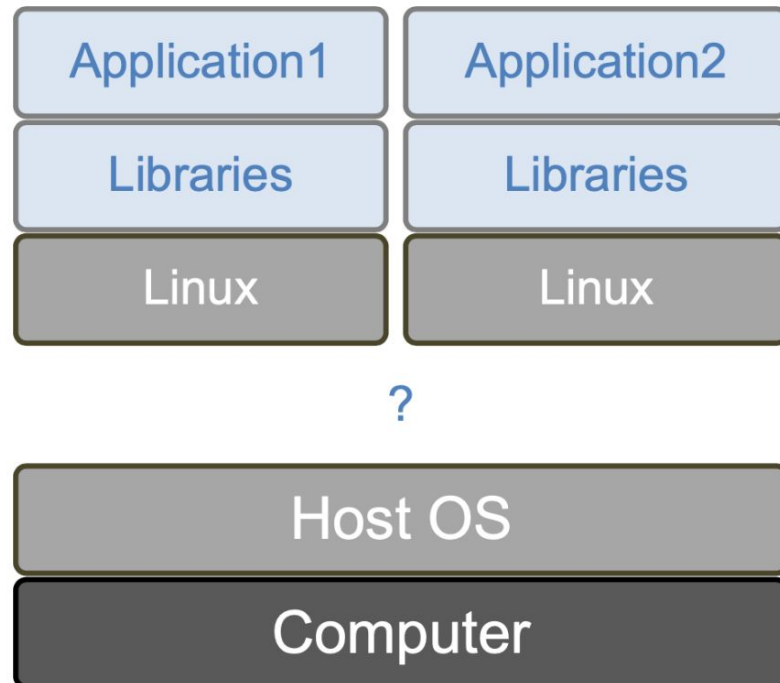There might be conflicts...

Usually dependencies of different applications don't interfere.
But what if we want to test the latest version of our favourite tool?
There might be conflicts. . .

| RStudio v1 | RStudio v1.2 |
| R Packages | R Packages |
| Environment 1 | Environment 2 |
| Conda | |
| Host OS | |
| Computer | |

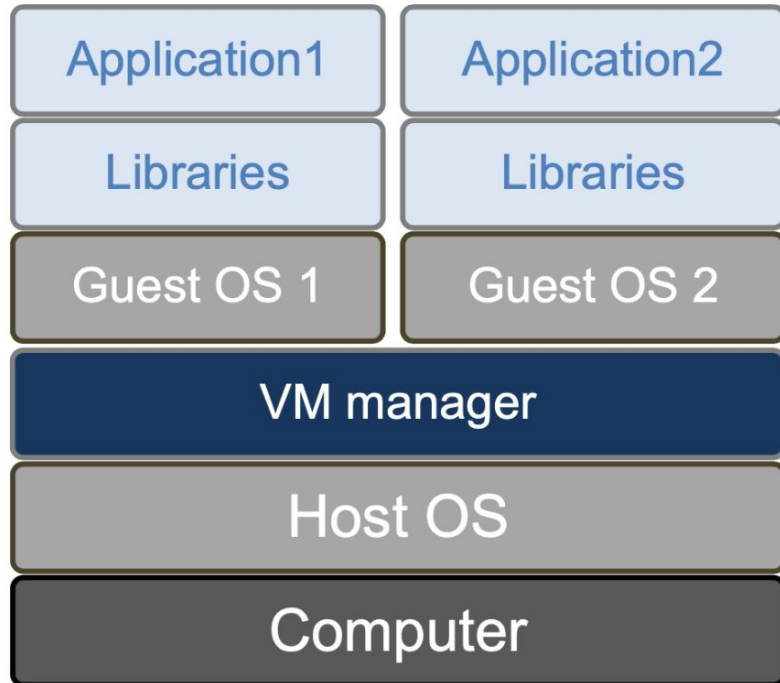Idea : create separated environments for each application.

CONDA

Idea : create separated environments for each application.
More versatile: create a new environment per analysis.

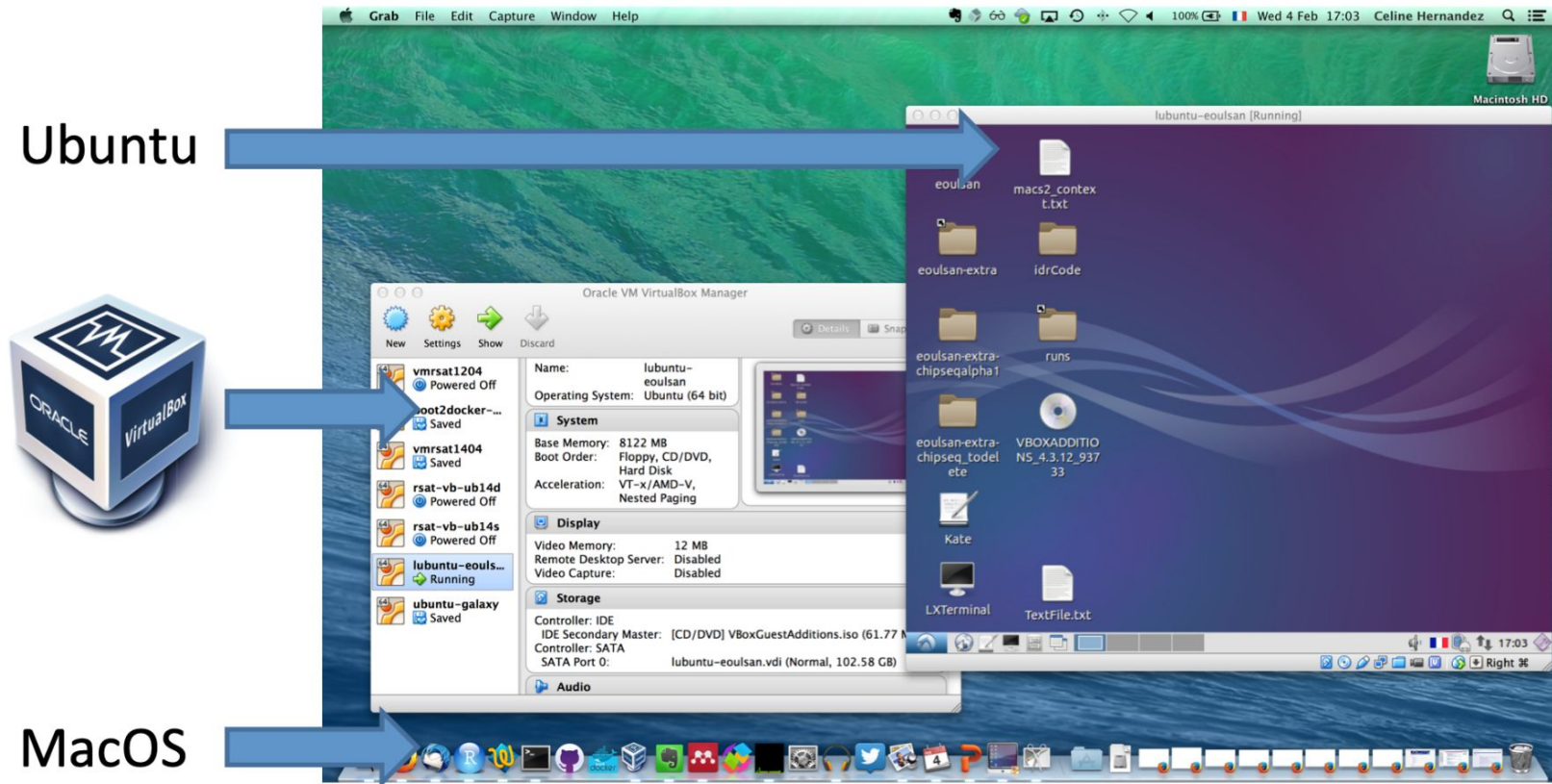But what if we want to install a software from a different OS?

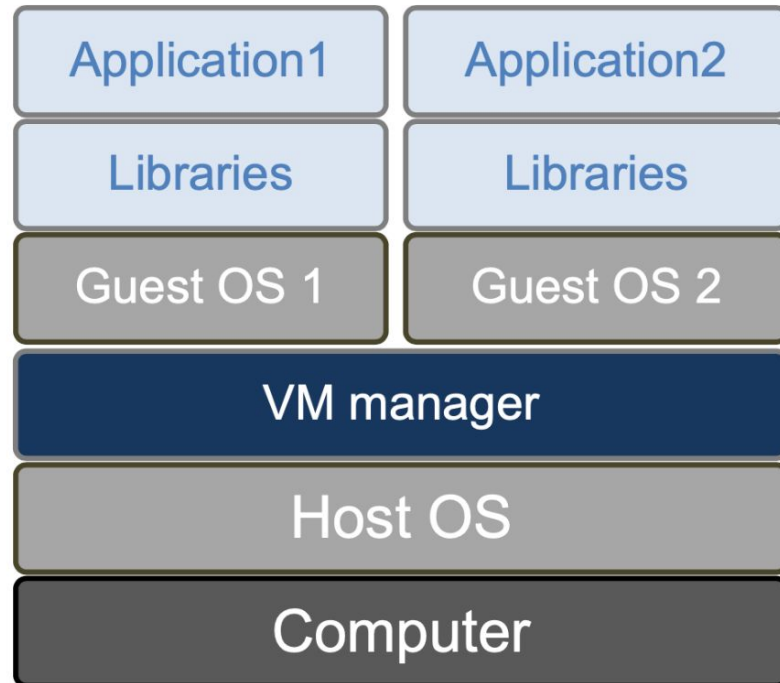| Application1 | Application2 |
|---|---|
| Libraries | Libraries |
| Guest OS 1 | Guest OS 2 |
| VM manager | |
| Host OS | |
| Computer | |

Idea: use virtual machines

Pros:

- Each application gets a completely different and independent environment
- Virtual machines can be transferred to another computer (using the same manager)
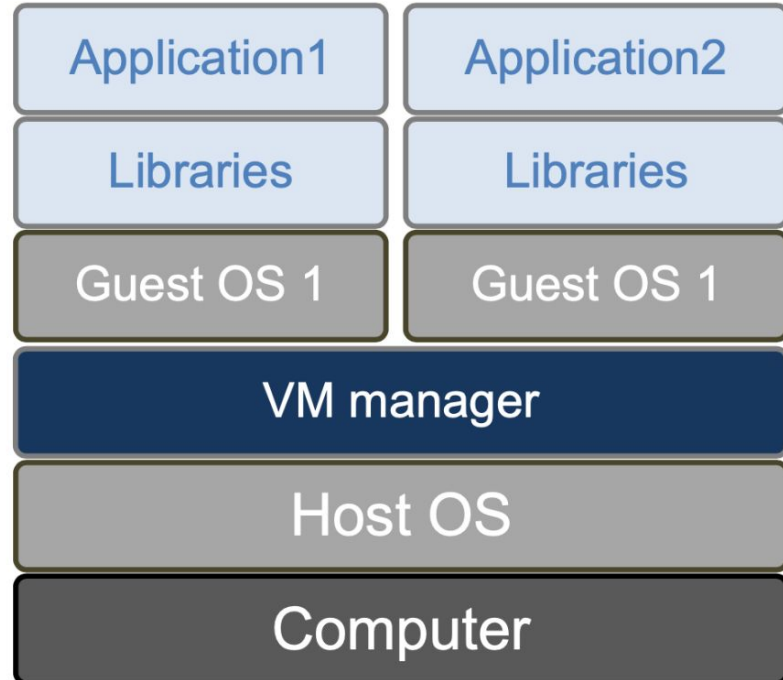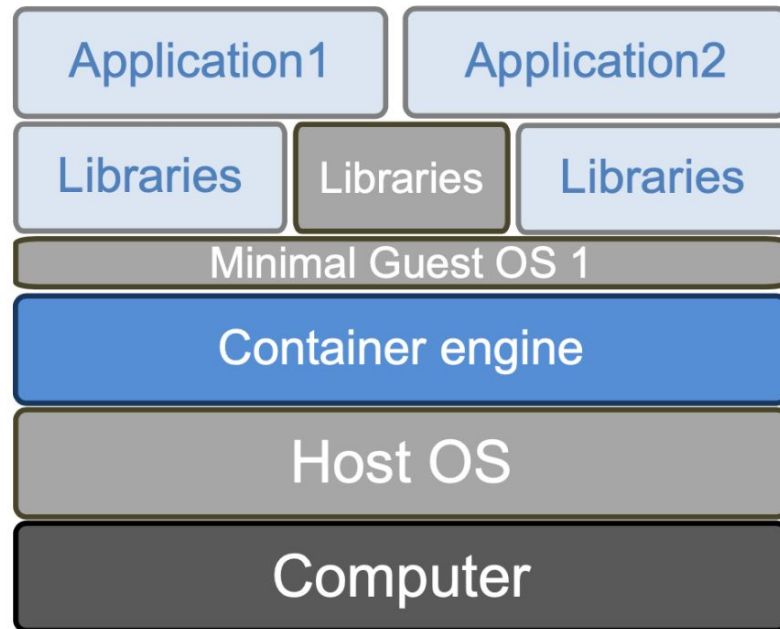
| Application1 | Application2 |
|---|---|
| Libraries | Libraries |
| Guest OS 1 | Guest OS 2 |
| VM manager | |
| Host OS | |
| Computer | |

Idea: use virtual machines
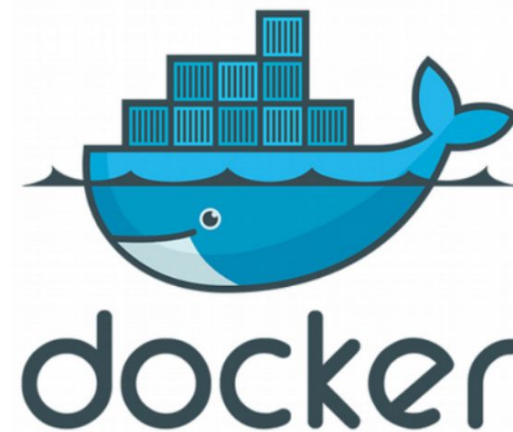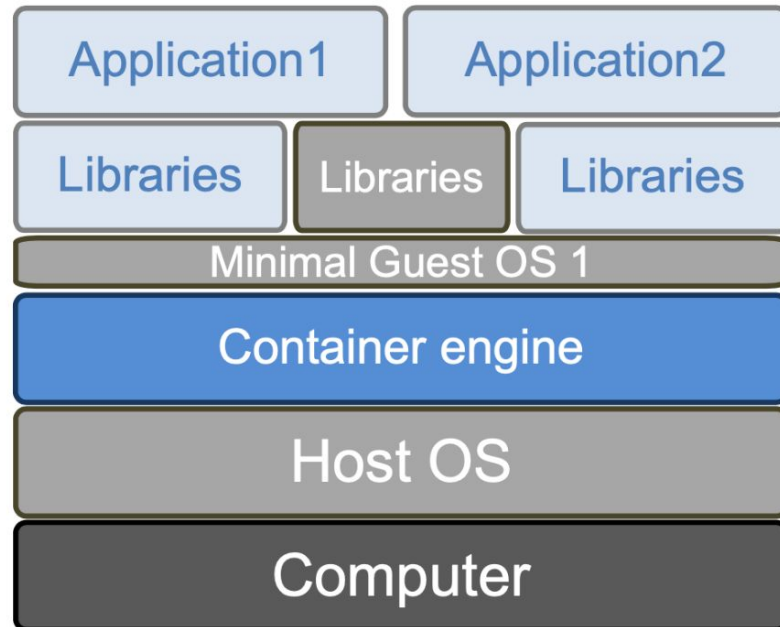Pros: transferable independent environments
Cons:

- Redundancy between VMs
- Heavy to set up
- No automation

Idea: "trick" applications into believing that they are in a different OS than the host's
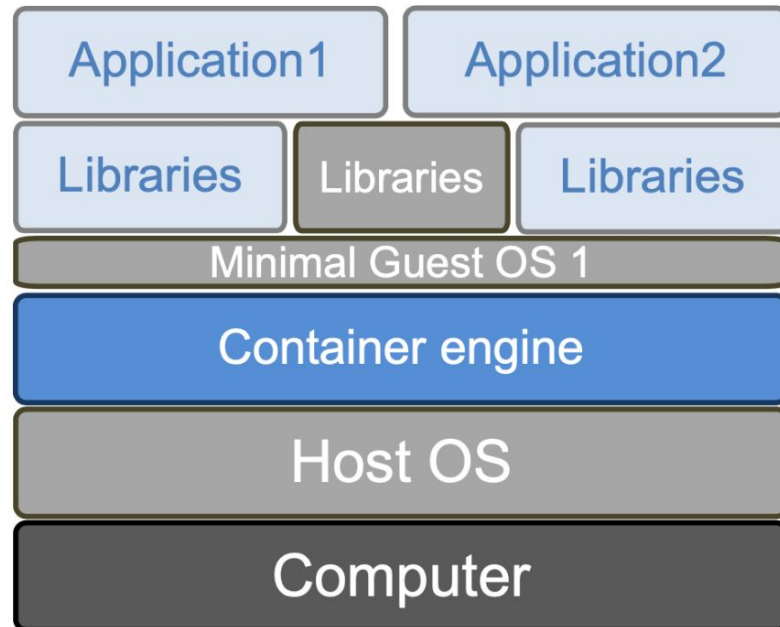Avoid redundancy.

OS virtualisation vs hardware virtualisation

Pros:

- Speed
  - ▶ Installation is faster
  - ▶ No boot time
- Lightweight
  - ▶ Minimal base OS
  - ▶ Minimal libraries and application set
- Easy sharing of applications

Cons:
- Singularity to use images on a cluster
- Changes of policies of the Docker company

## Update of the Docker Image retention policy (13/08/2020)

**What is a container image retention limit and how does it affect my account?**

Image retention is based on the activity of each individual image stored within a user account. If an image has not either been pulled or pushed in the amount of time specified in your subscription plan, the image will be tagged "inactive." Any images that are tagged as "inactive" will be scheduled for deletion. Only accounts that are on the **Free** individual or organization plans will be subject to image retention limits. A new dashboard will also be available in Docker Hub that offers the ability to view the status of all of your container images.

**What are the new container image retention limits?**

Docker is introducing a container image retention policy which will be enforced starting November 1, 2020. The container image retention policy will apply to the following plans:

- Free plans will have a 6 month image retention limit

- Pro and Team plans will have unlimited image retention
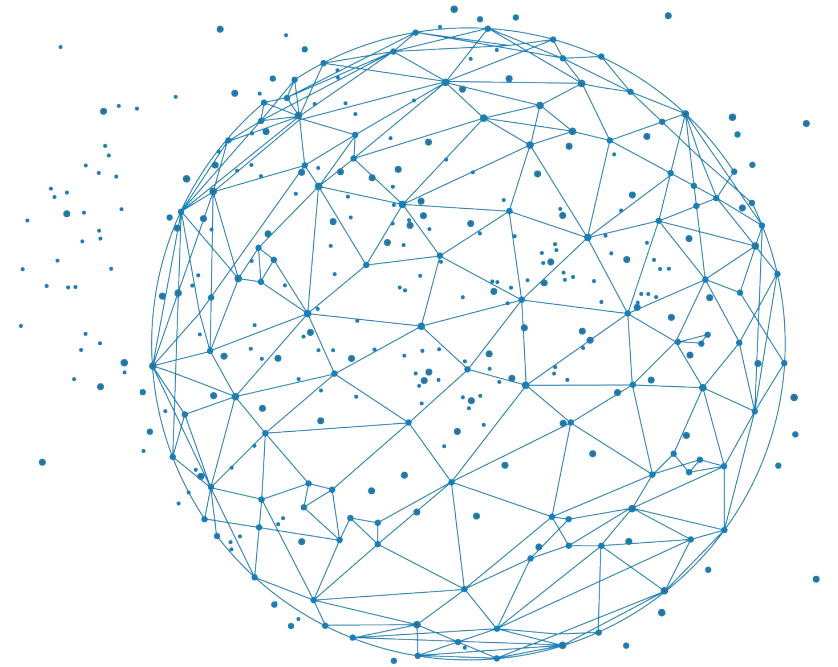
https://www.docker.com/pricing/retentionfaq

Practical Computational Reproducibility in the Life Sciences - Björn Grüning et al (2018)

# About Docker

Docker is not very "old"

- First commit January 2013
- First version March 2013
- Version 1.0 in June 2014

But its adoption was fast

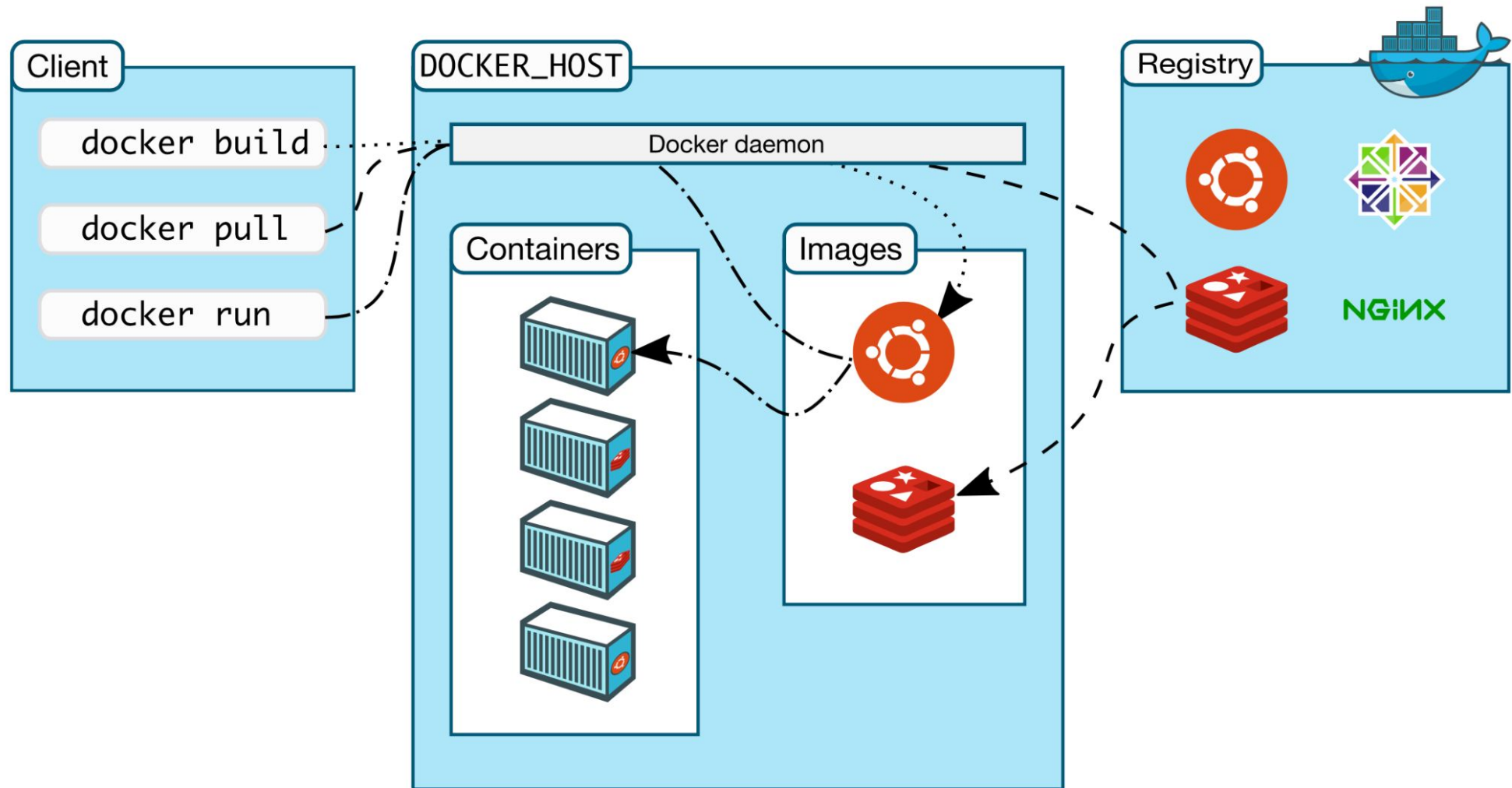- Officially packaged in Ubuntu since 2014 (v14.04)

## Image



- Set of libraries and functions
- Fixed. Cannot be modified
- Can be stored/shared online
- Can be automatically built

## Container



- "Active image"
- Can be modified (interactive)
- Can be turned into an image
- One image, many containers

(https://docs.docker.com/get-started/overview/)

# What is Docker?



(https://hub.docker.com/)

# Encapsulating our pipeline tools into Docker

**What we have now ?**

- A Git repository with our Notebook
- A environment lock file to recreate quickly a conda environment and get all dependencies

**But**

- Notebooks require a JupyterLab server…

What if we could shared a runnable JupyterLab environment including our Notebook and all its dependencies ?

Introducing…



`https://mybinder.org`

Binder is an online service that lets you share your notebook(s) in a interactive environment.

**How does it work ?**

BinderHub (the backend of Binder) do the following operation:

- Fetch your repo from GitHub
- Analyse the contents
- Build a Docker image based on your repo
- Launch that Docker image in the cloud
- Connect you to it via your browser

# Get the HTTPS url of your repository on Github

# Go to https://mybinder.org

# Paste your Github repository URL and click Launch

# The build process is running…

Binder is picking up our environment.yml file automatically !

# What if I don't have a Notebook. Can I still share my work as a Docker image ?

Introducing…    **Github Container Registry**

# Let's create a shell version of our Notebook

```
$ jupyter nbconvert my_first_notebook.ipynb --to script
[NbConvertApp] Converting notebook my_first_notebook.ipynb to script
[NbConvertApp] Writing 990 bytes to my_first_notebook.sh
```

We need to **customize our script** to make it **more portable**

Add the **bash shebang** at the beginning of the my_first_notebook.sh script :

```
#!/usr/bin/env bash
```

Set a variable to retrieve the directory containing the script (we need this to call the R script)

```
SCRIPT_DIR=$( cd -- "$( dirname -- "${BASH_SOURCE[0]}" )" &> /dev/null && pwd )
```

Use the variable in the R script call

```
R < $SCRIPT_DIR/Deseq2.r --no-save
```

## Add the new shell script to your repository

```
$ git add my_first_notebook.sh
$ git commit -m "add bash version of the pipeline"
```

Github can help us create automatically a Docker image based on a Dockerfile located in our repository.

The step to do this are :

- Create a Dockerfile and push it Github
- Enable the Publish Docker Container Github action

## Create a new docker folder in your pipeline folder

```
$ mkdir docker
```

# Create a new file named `Dockerfile` in the `docker` folder

```
FROM mambaorg/micromamba:1.5.6
ADD ../environment.yml .
RUN micromamba install -y -n base -f environment.yml && \
    micromamba clean --all --yes

USER root
RUN mkdir /opt/pipeline
ENV PATH="/opt/pipeline:$PATH"
ADD ../my_first_notebook.sh /opt/pipeline/
ADD ../Deseq2.r /opt/pipeline
RUN chmod +x /opt/pipeline/my_first_notebook.sh

USER $MAMBA_USER
```

Create an image based on micromamba distribution
Add the environment.yml file
Install packages listed in the environment.yml in the base environment
and cleanup everything to make the image as light as possible

Switch to root user
Create a /opt/pipeline folder
Add the /opt/pipeline folder to the image $PATH
Add the pipeline script to the /opt/pipeline folder
Also add the Deseq2.r file
Make sure the script is executable

Switch back to mamba user

# Add the docker folder to the repository

```
$ git add docker
$ git commit -m "add support for docker"
```

## Push the last commits to Github

```
$ git push origin main
```

Github let you define actions that will be triggered automatically when your repository changes.
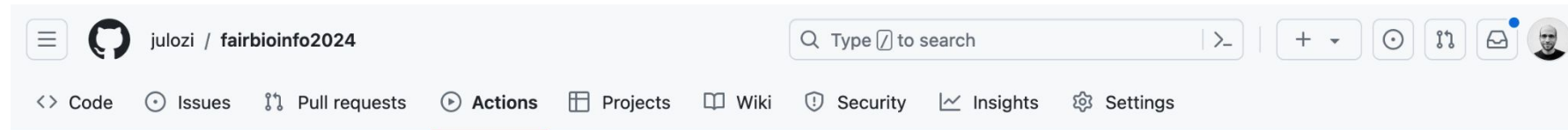
Actions are configure through **YAML files** called workflow.

Github proposes a set of pre-defined actions ready to use.

The **Publish Docker Container** action is a workflow that will build a Docker image based on a Dockerfile found at the root of your repository. It will then publish it in the Docker registry provided by Github.

Go to your repository on Github and move to the **Actions** tab
You should see the **Publish Docker Container** action suggested
Click its **Configure** button

Github gives us a publish-docker.yml workflow ready to be pushed on our repository.

We need to make **two changes** in this file :

First, the image signing tool proposed by Github is outdated (this is a bug).
We need to upgrade it :

```
# Install the cosign tool except on PR
# https://github.com/sigstore/cosign-installer
- name: Install cosign
    if: github.event_name != 'pull_request'
    # uses: sigstore/cosign-installer@6e04d228eb30da1757ee4e1dd75a0ec73a653e06 #v3.1.1
    uses: sigstore/cosign-installer@e1523de7571e31dbe865fd2e80c5c7c23ae71eb4 #v3.4.0
    with:
    # cosign-release: 'v2.1.1'
    cosign-release: 'v2.2.3'
```

By default, the workflow is looking for a Dockerfile at the root of the repository. However we have created our Dockerfile in a docker folder.

We fix the job to use the docker/Dockerfile :

```
- name: Build and push Docker image
  id: build-and-push
  uses: docker/build-push-action@0565240e2d4ab88bba5387d719585280857ece09 # v5.0.0
  with:
  context: .
  file: ./docker/Dockerfile
  push: ${{ github.event_name != 'pull_request' }}
  tags: ${{ steps.meta.outputs.tags }}
  labels: ${{ steps.meta.outputs.labels }}
  cache-from: type=gha
  cache-to: type=gha,mode=max
```

# Click the green **Commit changes…** button and validate the commit in the modal window.
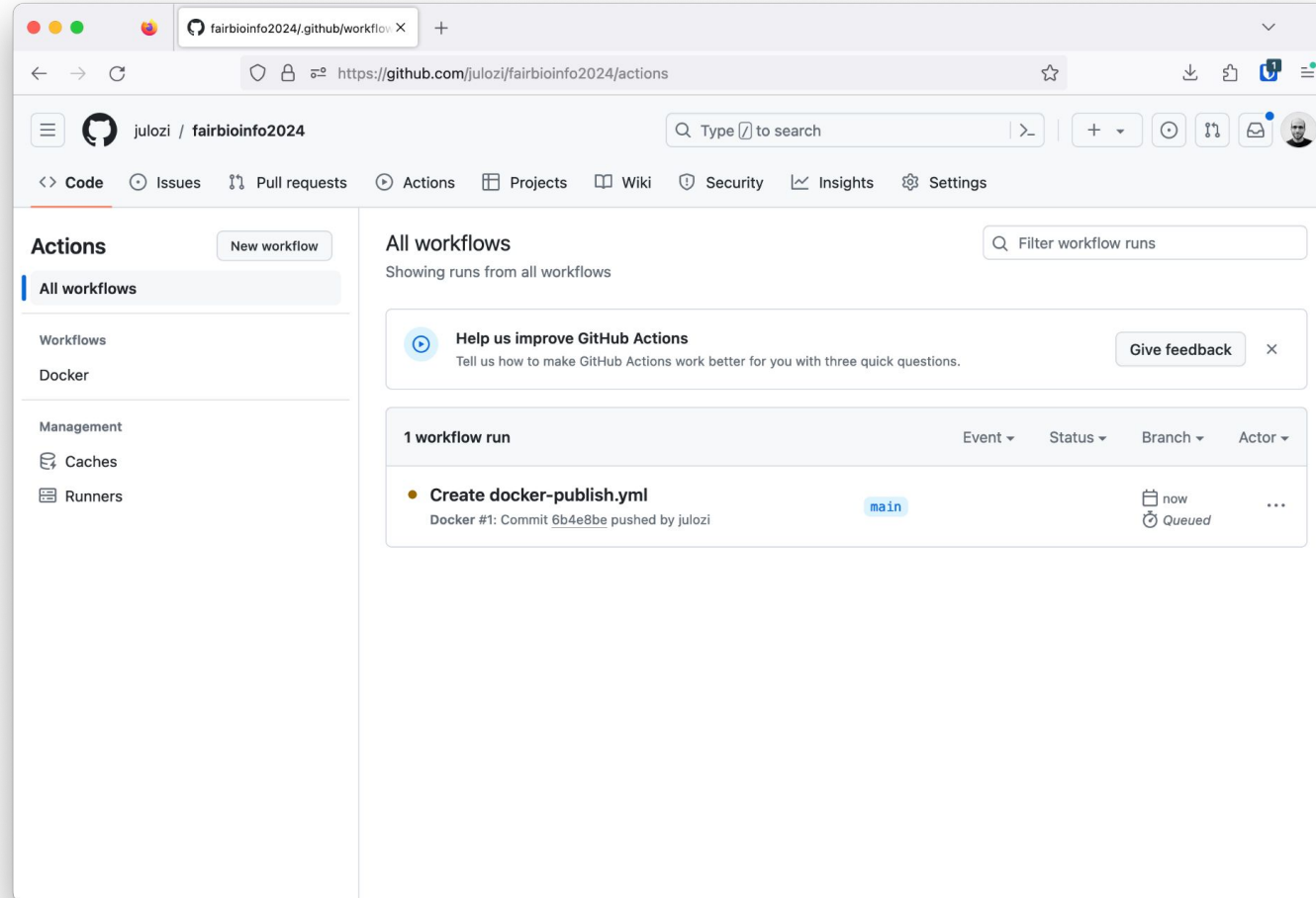
Github will automatically trigger the new action workflow.
You can find it back in the **Actions** tab.

When the build is finished successfully, the Docker image can be found in Packages section of the repository.

# Click on the main branch link to get the image URL

# Click on the main branch link to get the image URL



You can now run the pipeline on any Docker enabled computer :

```
docker run -t -i \
ghcr.io/julozi/fairbioinfo2024:main \
my_first_notebook.sh
```

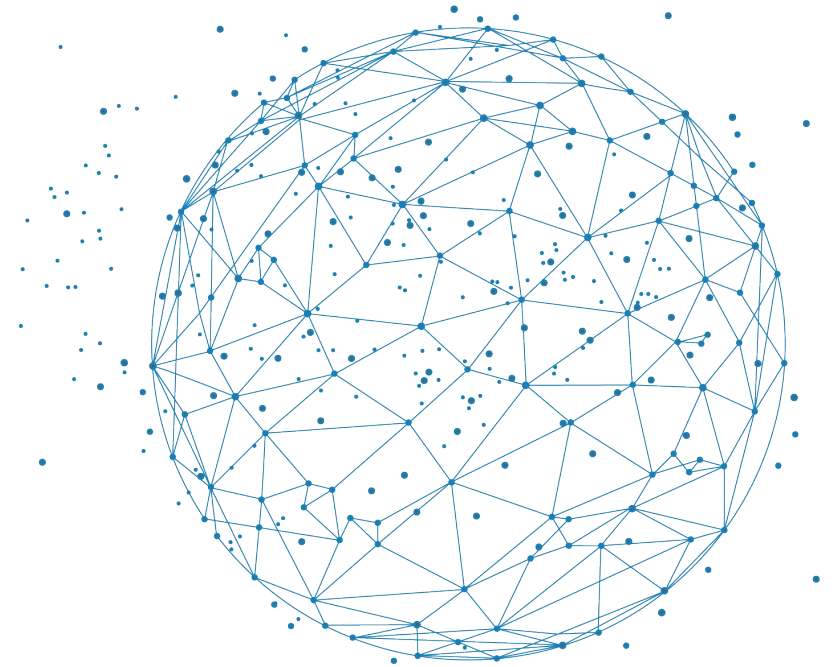*Use the* `--platform linux/amd64` *option if you are working on Apple Silicon*
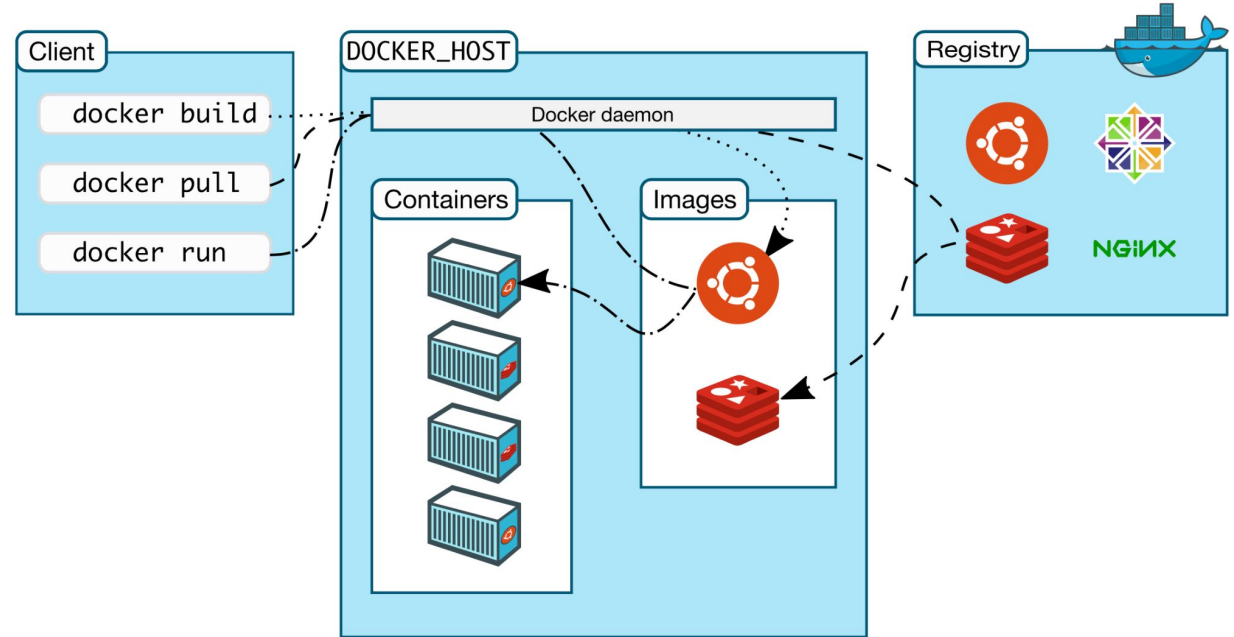
# Apptainer :
# container on HPC cluster

Docker requires a **Docker Host** to run containers.

The Docker Host is a system daemon that run as root and can access to a reserved part of the hardware resources.

**This is not compatible with an HPC cluster where hardware resources are already managed by a job scheduler (SLURM)**



(https://docs.docker.com/get-started/overview/)

Apptainer is an open source container platform designed to run complex applications on high-performance computing (HPC) clusters in a simple, portable, and reproducible way.

An Apptainer container image is a file

An Apptainer running container is a user process

Apptainer has its own image definition format which is different than Dockerfile.

However, it is possible to build an Apptainer image directly from a Docker image URL.

Let's create an Apptainer image file from our pipeline Docker image :

```
$ mkdir apptainer
$ cd apptainer
$ apptainer build pipeline.sif docker://<your docker image URL>
```

This will create a pipeline.sif file ready to be used.

You can now run the pipeline script contained in the image :

```
$ apptainer run pipeline.sif my_first_notebook.sh
```