# FAIR_bioinfo : Open Science and FAIR principles in a bioinformatics project

## How to make a bioinformatics project more reproducible

C. Hernandez[1]    T. Denecker[2]    J. Sellier[2]    G. Le Corguillé[2]
C. Toffano-Nioche[1]

[1]Institute for Integrative Biology of the Cell (I2BC)
UMR 9198, Université Paris-Sud, CNRS, CEA
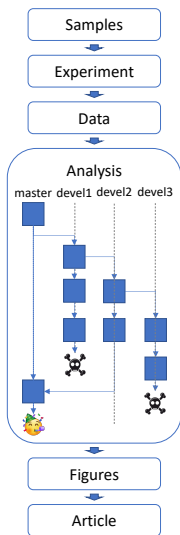91190 - Gif-sur-Yvette, France
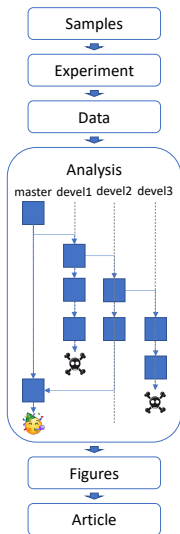
[2]IFB Core Cluster taskforce

June 2021

# Introduction

A (not-so-uncommon) nightmare

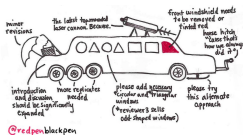# Introduction

## A (not-so-uncommon) nightmare

# Introduction

## A (not-so-uncommon) nightmare

# Introduction

A (not-so-uncommon) nightmare



What changed?

- Package
- Software
- Libraries
- Environment variables

- OS version
- Computer
- ..?

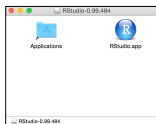# Different levels of encapsulation

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

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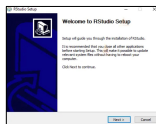- Environment management CONDA

# Encapsulation

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

MacOS

Windows

Install Rstudio ?



Use Rstudio

Unix-based

# Encapsulation

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



Host OS

Computer

# Encapsulation

RStudio

Host OS

Computer

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

# Encapsulation

RStudio

R Packages

Host OS

Computer

We started with a computer using a specific OS...
And inside this environment, we installed a new application.
Applications rely on dependencies, e.g. external libraries.

# Encapsulation

| |
|---|
| 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. . .

# Encapsulation



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...

# Encapsulation : managing environments



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

Idea : create separated environments for each application.

# Encapsulation : managing environments

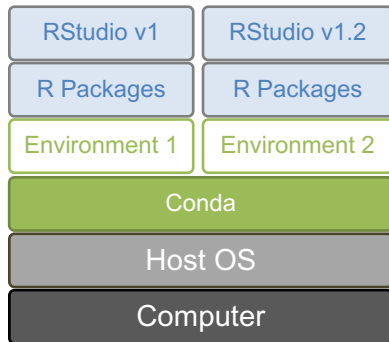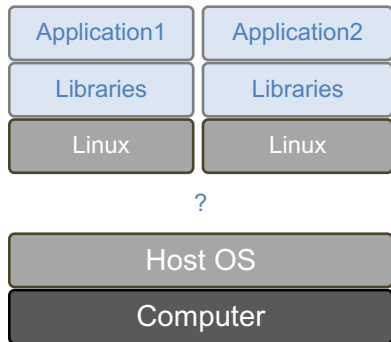| Analyse 1 | Analyse 2 |
|---|---|
| RStudio v1 | RStudio v1 |
| R Packages | R Packages |
| Environment 1 | Environment 2 |
| Conda | |
| Host OS | |
| Computer | |

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

# Encapsulation : hardware virtualisation

| Application1 | Application2 |
|---|---|
| Libraries | Libraries |
| Linux | Linux |

?

| Host OS |
|---|
| Computer |

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

# Encapsulation : hardware virtualisation

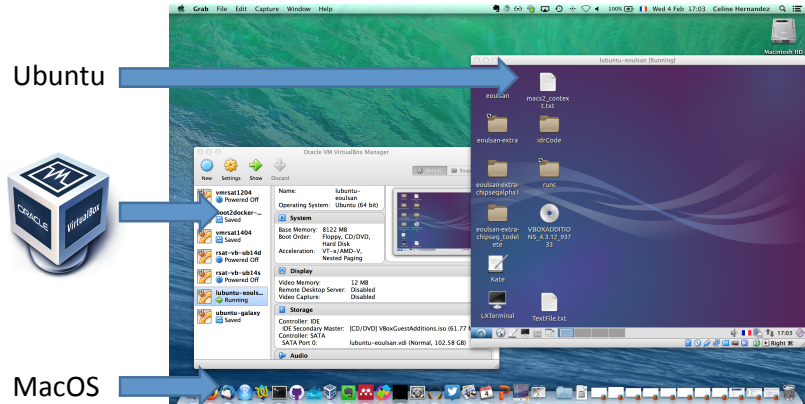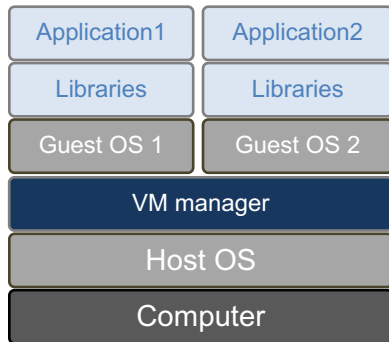| Application1 | Application2 |
|:---:|:---:|
| Libraries | Libraries |
| Guest OS 1 | Guest OS 2 |
| VM manager ||
| Host OS ||
| Computer ||

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)

# Encapsulation : hardware virtualisation

Ubuntu

MacOS

# Encapsulation : hardware virtualisation

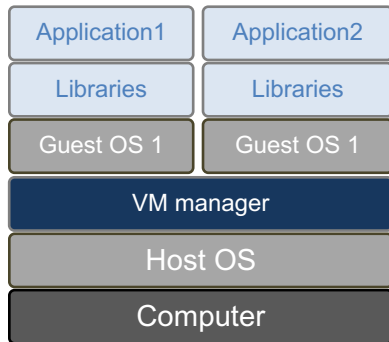| | |
|---|---|
| Application1 | Application2 |
| Libraries | Libraries |
| Guest OS 1 | Guest OS 2 |
| VM manager | |
| Host OS | |
| Computer | |

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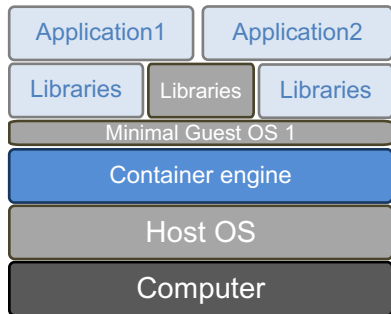environments
Cons:

- Redundancy between VMs
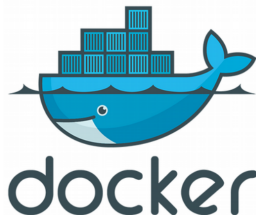- Heavy to set up
- No automation
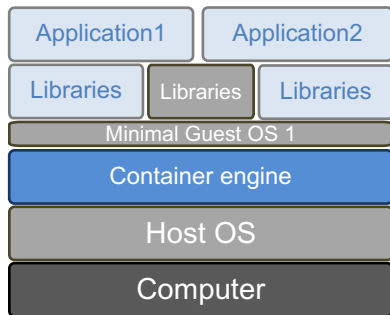
# Encapsulation : OS virtualisation

# Encapsulation : OS virtualisation



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

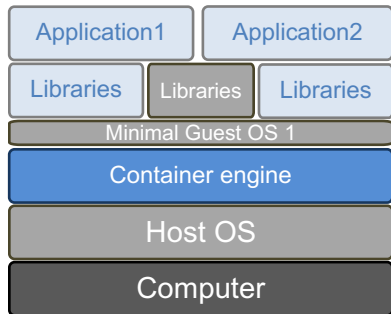# Encapsulation : OS virtualisation



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

# Encapsulation : OS virtualisation

| Application1 | Application2 |
|---|---|
| Libraries | Libraries | Libraries |

Minimal Guest OS 1

Container engine

Host OS

Computer

Cons:

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

# Docker policy

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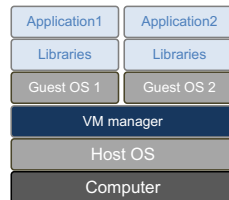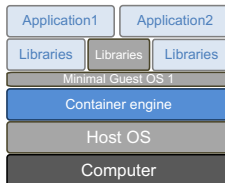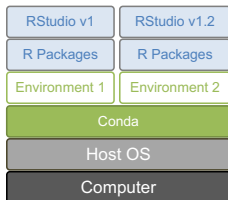
**What are the new container image retention limits?**

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:
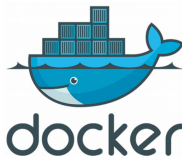
- Free plans will have a 6 month image retention limit

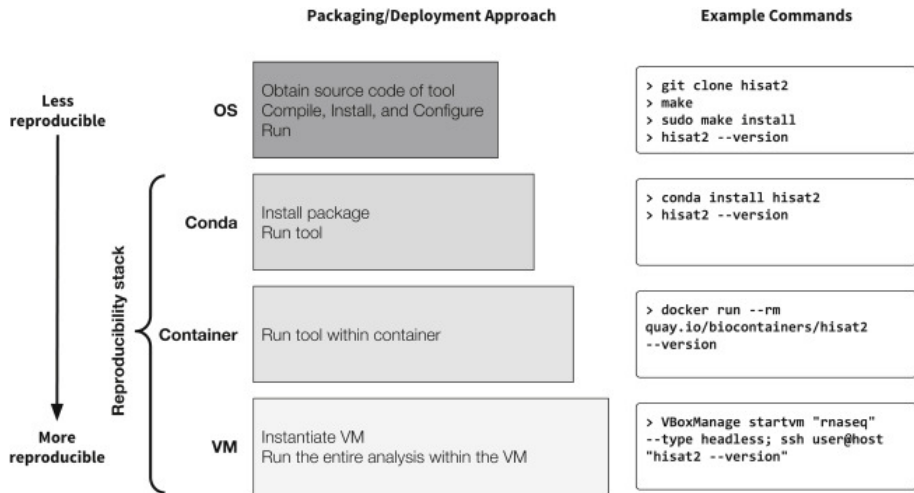- Pro and Team plans will have unlimited image retention

https://www.docker.com/pricing/retentionfaq

# Encapsulation

# Encapsulation and reproducibility stack



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