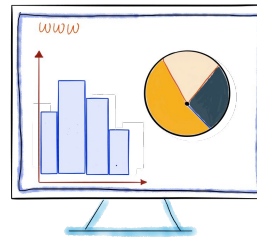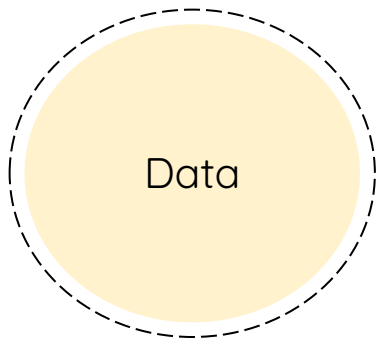# Reproducible Science in practice
## tools and ideas
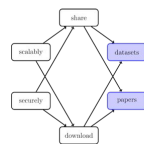
Arturo Sánchez Pineda (LAPP)
June 10, 2021 - ESCAPE (online) School

# Some current tools **per element**

**Data**

Usually, labs and experiments have dedicated data repositories for their users.

Here, I want to mention to Open Access datasets repositories as examples
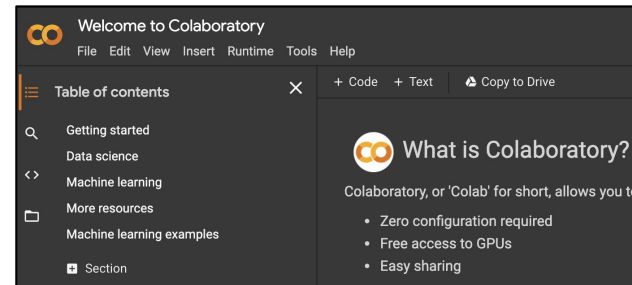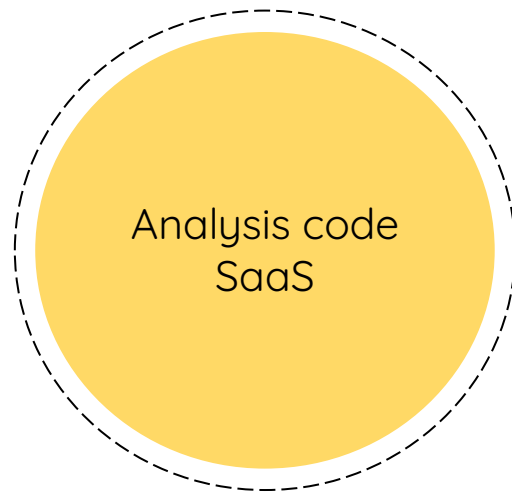
zenodo

**Academic Torrents**

Two services that can be really useful for storage and preservation of datasets (and other digital objects)

Analysis platforms/environments like Jupyter grow in popularity and started to be some of the "standard" in several domains of data analysis

Here some examples of such tools, where Jupyter is offer as a service by public or private institutions

Analysis code SaaS







A multi-user version of the notebook designed for companies, classrooms and research labs
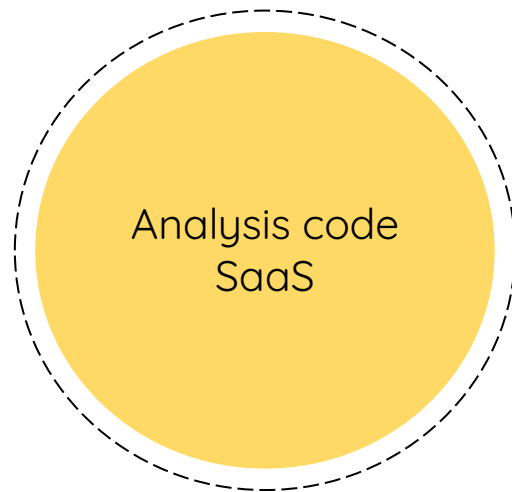
Analysis platforms/environments like Jupyter grow in popularity and started to be some of the "standard" in several domains of data analysis

Here some examples of such tools, where Jupyter is offer as a service by public or private institutions

Analysis code SaaS

IBM | IBM Developer  Topics ⌄  Products & Services ⌄  Community ⌄  Open source at IBM ⌄

Jupyter Notebook

Get Jupyter Notebook ↗
Articles
Learning Paths
Code Patterns
Podcasts
Open Project
Tutorials
Videos

Overview

An open-source web application that supports interactive data science and scientific computing across all programming languages

Jupyter Notebooks are open-source web applications that let you create and share documents that contain live code, equations, visualizations and narrative text.

jupyterhub

A multi-user version of the notebook designed for companies, classrooms and research labs
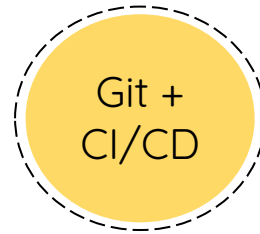
Microsoft Azure Notebooks

**Learn more about all the notebooks experiences from Microsoft and GitHub**

The Azure Notebooks preview has ended. You can enjoy powerful, integrated Jupyter notebooks with the following products and services from Microsoft and GitHub.

Microsoft

+

jupyter

ESCAPE
European Science Cluster of Astronomy & Particle physics ESFRI research infrastructures

**ATLASSIN**
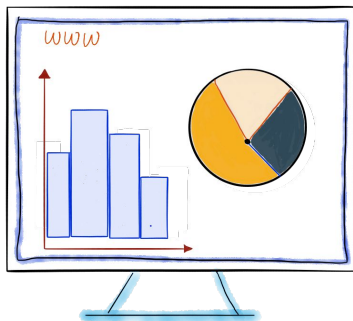**Bitbucket**

**GitLab**

**GitHub**

Git + CI/CD

There are several companies that allow the creation and hosting of Git repositories (you are using one of those right now)

But you can also self-hosted one of those instances.
They also come with a lot of functionalities like CI/CD

ESCAPE
European Science Cluster of Astronomy &
Particle physics ESFRI research infrastructures

The computing infrastructure is, for example, your laptop/desktop machine. There the needed OS, software and tools are installed to perform the analysis (like what we are doing during the school)

But you can also get the needed environment using Virtual Machines or containers
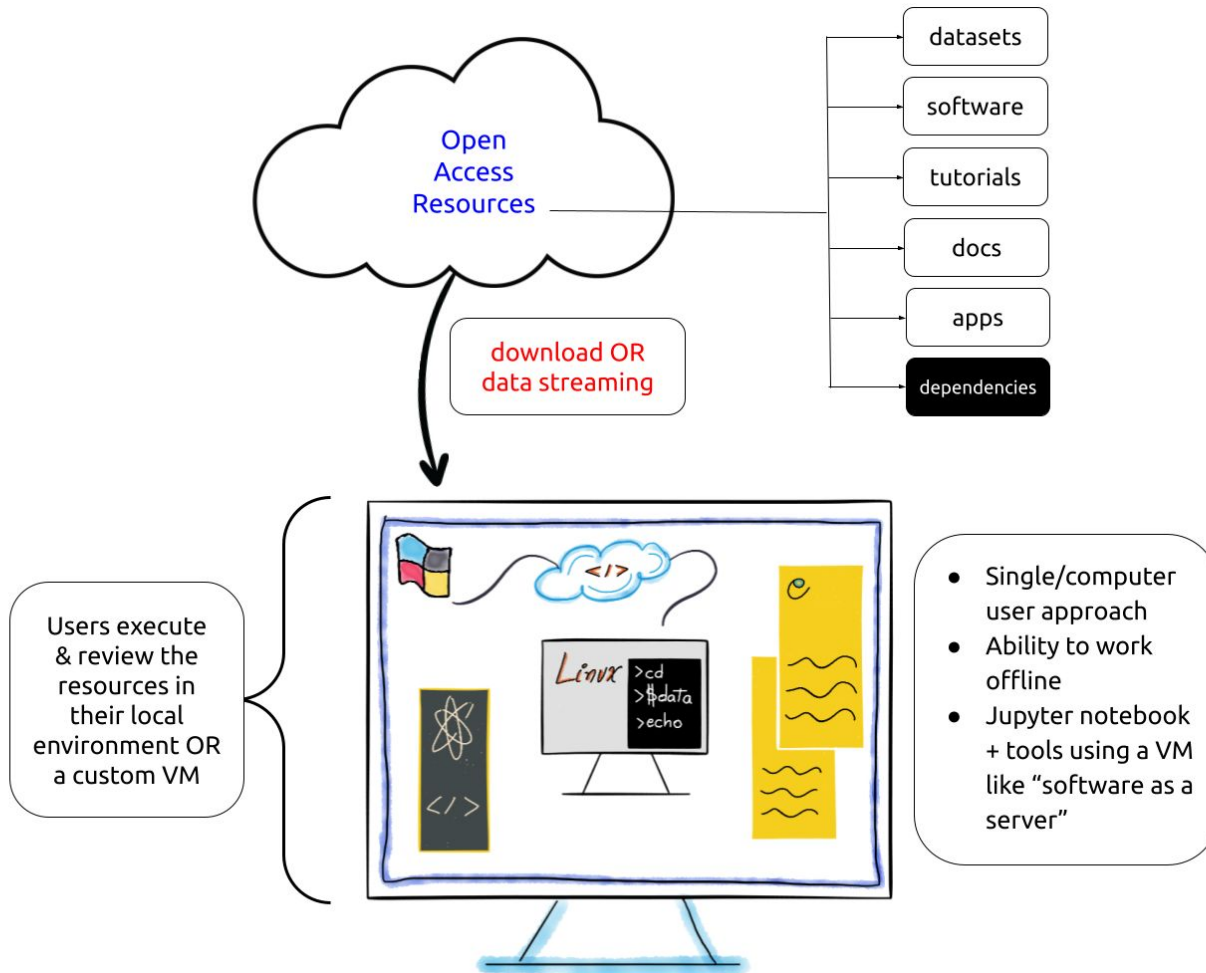


Computing IaaS

Examples of this VM usage

**Host Machine**

**Use as a Server**

A classical example is to host a VM that we can use as a private "server" isolating and preserving the working environment

Open Access Resources

datasets

software

tutorials

docs

apps

dependencies

download OR data streaming

Users execute & review the resources in their local environment OR a custom VM

Linux
>cd
>#data
>echo

- Single/computer user approach
- Ability to work offline
- Jupyter notebook + tools using a VM like "software as a server"

ESCAPE
European Science Cluster of Astronomy & Particle physics ESFRI research infrastructures

# The JupyterLab **UI**

# The JupyterLab UI

A well-known tool for all of us (data analysis and visualisation) is the Jupyter notebook.

JupyterLab is a suite of tools and features that allow interacting with multiple elements in a single view. And do the computation, of course.

https://jupyterlab.readthedocs.io/en/stable/

Repositories and
some stored datasets

Computer
facility

> < run analysis >

> ...

>command 1 ...

>command 2 ...

Analysis
development

# An example of JupyterLab

(a 90 sec video)

Repositories and
some stored datasets

Data analysis notebooks

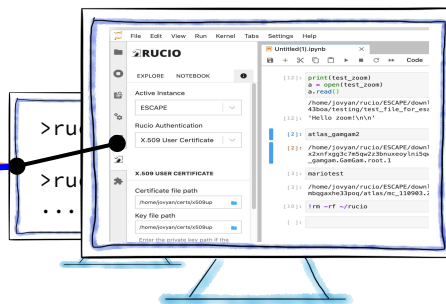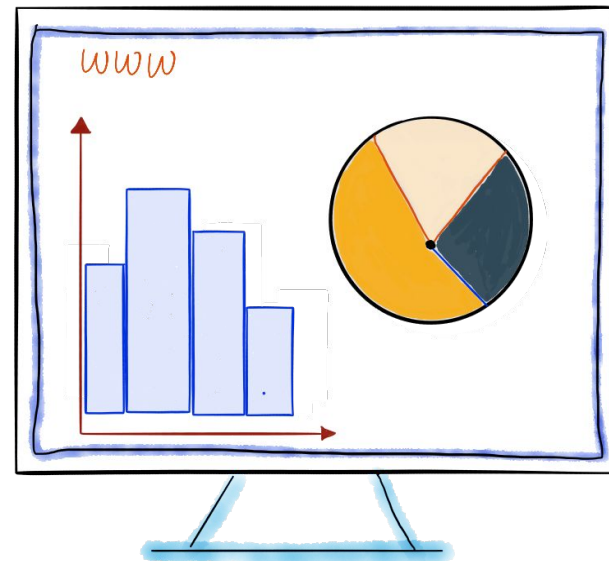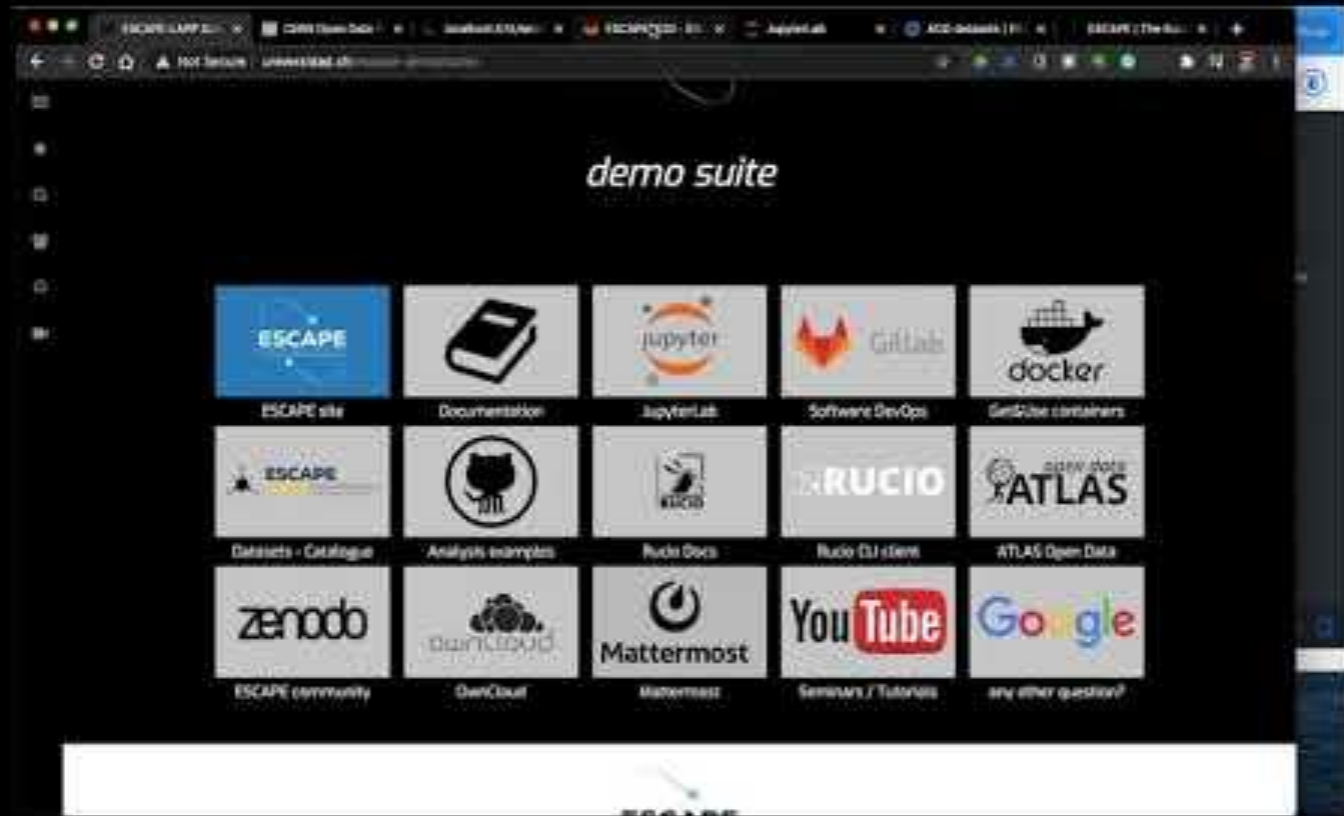Analysis results and
visualisation

# An example of JupyterLab

(a 150 sec video)

# Containers using **Docker**

Containers allow the preservation and reproducibility of software environments and applications

As an example, we can have a Docker container and execute it in our machine or in a remote cloud... inside a VM :)



Computing IaaS

Containerized Applications

| App A | App B | App C | App D | App E | App F |

Docker

Host Operating System

Infrastructure

| Virtual Machine | Virtual Machine | Virtual Machine |
| App A | App B | App C |
| Guest Operating System | Guest Operating System | Guest Operating System |

Hypervisor

Infrastructure

ESCAPE
European Science Cluster of Astronomy & Particle physics ESFRI research infrastructures
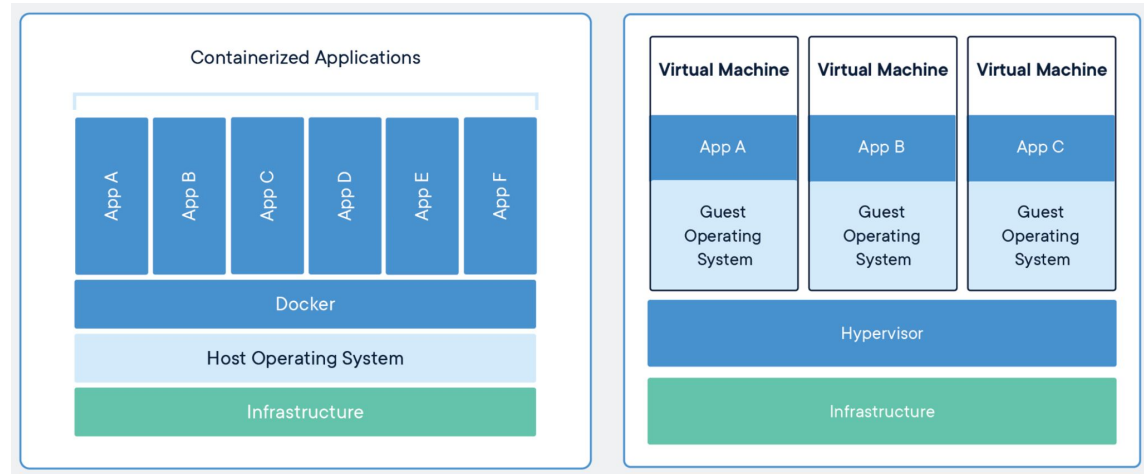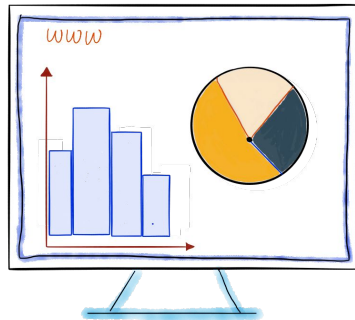
Containers allow the preservation and reproducibility of software environments and applications

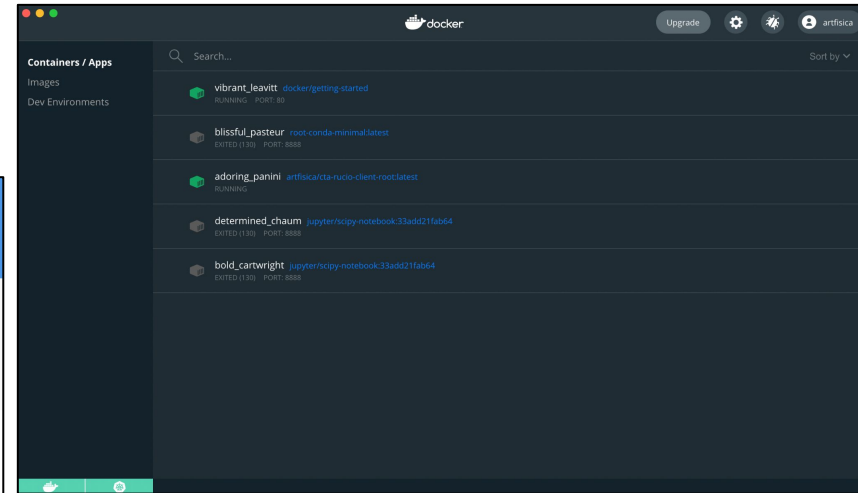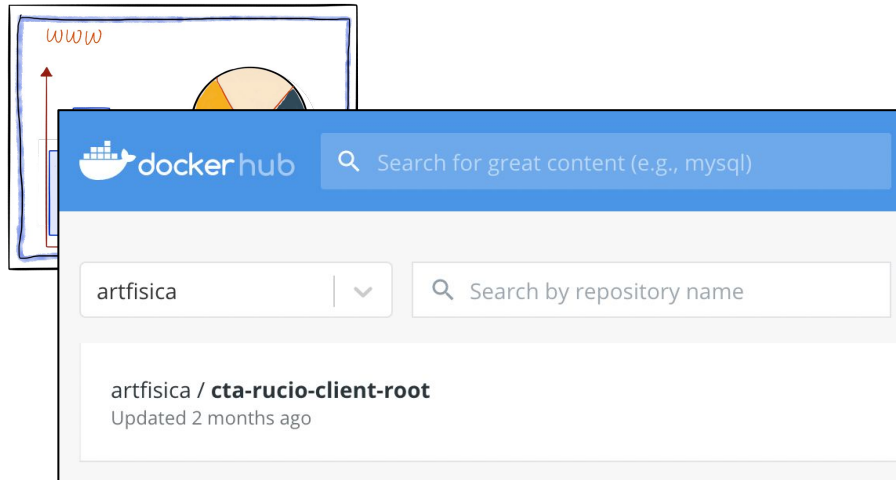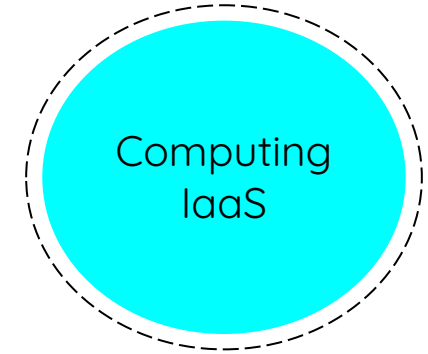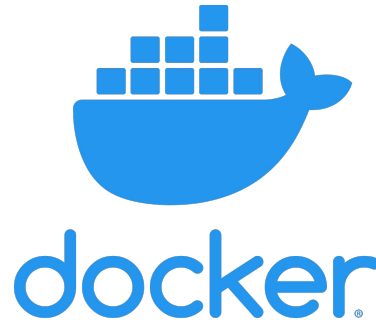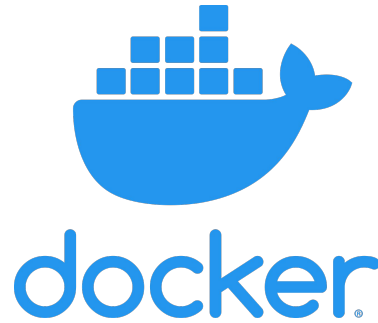As an example, we can have a Docker container and execute it in our machine or in a remote cloud… inside a VM :)
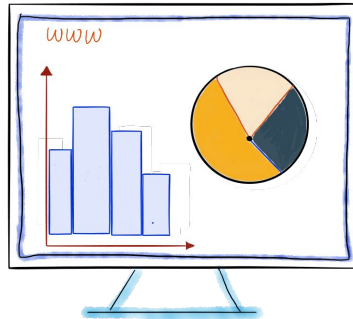


Computing IaaS

For example, you can use a Docker container to deploy a JupyterLab instance in your computer

Same can be done with JupyterHub

Computing IaaS



Selecting an Image →

## Jupyter Docker Stacks

Jupyter Docker Stacks

⭐ Star  6,013

Jupyter Docker Stacks are a set of ready-to-run Docker images containing Jupyter applications and interactive computing tools. You can use a stack image to do any of the following (and more):

### Navigation

Notebook server in a local Docker container
for a team using JupyterHub
ockerfile

```
arturosanchezpineda@lappm-p936 ~ % docker run -p 8888:8888 jupyter/scipy-notebook:33add21fab64
Unable to find image 'jupyter/scipy-notebook:33add21fab64' locally
33add21fab64: Pulling from jupyter/scipy-notebook
345e3491a907: Already exists
57671312ef6f: Already exists
5e9250ddb7d0: Already exists
e7754708c251: Pull complete
87ea6bc379a7: Pull complete
4f4fb700ef54: Pull complete
fd071e899aca: Pull complete
877d1df1399e: Pull complete
4c423a439cbd: Pull complete
c938f6886e4d: Downloading [===================>       ]  77.41MB/91.14MB
497d6a83aa43: Download complete
7cd70cc13eeb: Download complete
527c1354eec7: Download complete
8283241e9dbb: Download complete
6a2bfbad551f: Downloading [====================>      ]  116.6MB/287.1MB
941ce72a2d0b: Download complete
1e0c07ad01b8: Downloading [=====================>     ]  113.4MB/168.1MB
dc290e06bbe4: Waiting
da39b0750d3c: Waiting
4567a9b613fa: Waiting
```

ESCAPE
European Science Cluster of Astronomy &
Particle physics ESFRI research infrastructures

# Installing Docker and **reproduce notebooks**

arturos@cern.ch

# A concrete example
**reproducible analysis**

# reana

## Reproducible research data analysis platform

| Flexible | Scalable | Reusable | Free |
|---|---|---|---|
| Run many computational workflow engines. | Support for remote compute clouds. | Containerise once, reuse elsewhere. Cloud-native. | Free Software. MIT licence. Made with ❤ at CERN. |

https://reanahub.io/

# Reproducible analyses

Lesson on reproducible analyses and reusable containerised scientific workflows

https://awesome-workshop.github.io/reproducible-analyses/

# Last **comments**

# Common Services and Infrastructure

## Computing

- Academic / Dedicated allocated Computer infrastructure
  - Includes local resources
  - Includes HPC and super computers
- Public and Commercial Cloud Computing (IaaS)
- Volunteering Computing over Ethernet or the Internet
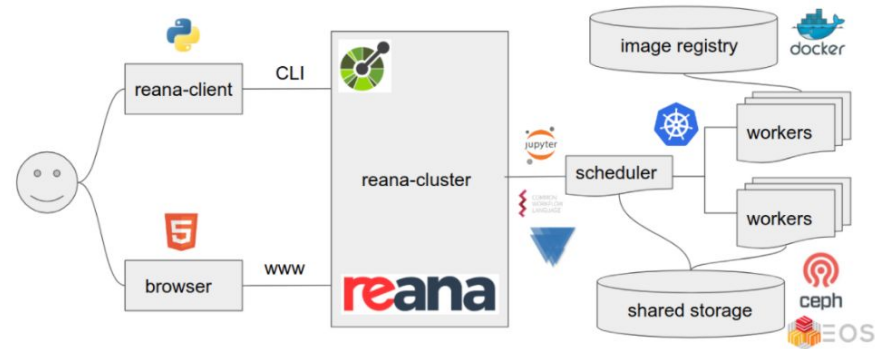- SysAdmins become part of the experiment.

## Monitoring

- Open Source tech and tools to keep track of process and experiments.
- Also to monitor in an automatic way vast datasets with the help of Machine Learning (ML) and Artificial Intelligence (AI)
- Services are deployed like "Monitoring as a Service" MaaS).

## Storage

- Multiple and interconnected storage facilities that can be costless for small and medium experiments
  - Includes volunteering and academic resources
- Use software coming for large experiments for data structure and file systems.

## Bookkeeping

- Different than storage, booking relies in informatic tools and protocols to track the production, usage and results of data.
- Also relevant for production chain when delivering components (hardware and software) to others.
- Reproduction of results.

## Software

- Software design,production, pipelines, CI/CD is vital for any scientific and academic endeavour.
- Tools for efficient code development and also Open Source and industrial quality frameworks.
- Creation of solutions that last as the experiments evolve
- SaaS will be crucial for institutions in the region.

ESCAPE
European Science Cluster of Astronomy & Particle physics ESFRI research infrastructures

In my view, **reproducibility** refers to a series of principles, techniques, tools and practical considerations that allow the documentation, recording and preservation of data analysis pipelines — enhancing the possibilities of collaborations across borders and increasing the probabilities of replicating results by others (and yourself) in the future.

Reproducibility involves using standard and well-established protocols to ensure that your code will survive outside your computer, the passing of time and that others will be able to use it as a starting point for new analysis.

# Thanks!

# Arturo Sánchez Pineda



**arturos@cern.ch**

Linkedin

**/arturo-sanchez-pineda/**

 **@artfisica**

https://twitter.com/Arturo_RSP

**I am post-doctoral fellow at LAPP-CNRS, France. Member of the ESCAPE and ATLAS groups.**
I studied Fundamental Physics and System Engineering in the Universidad de Los Andes, Venezuela, with a PhD in Fundamental and Applied Physics from Università di Napoli "Federico II", Italy.
I was previously a postdoctoral fellow at Physics Department at Università di Udine and an Associate at INFN, Italy. Also, an ATLAS TDAQ System Administrator at CERN, Switzerland, and Research Associate at the High Energy, Cosmology and Astroparticle Section at ICTP, Italy.
← **And I do a lot of outreach :)**