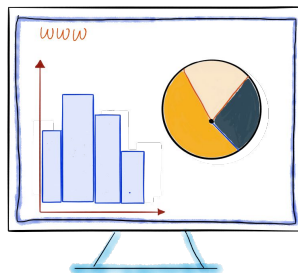# ESCAPE: a review of idea for Data + Analysis challenge using ATLAS Open Data

**Arturo Sánchez Pineda - LAPP**

26th (from 23rd) April, 2021

# Overview

**This is an attempt to describe a series of exercises to be performed during the ESCAPE Data and Analysis Challenge the next November**

- Data "multiplication" where multiple version of the same data is generated, simulating a data-augmentation process
  - Requesting data from the Datalake at higher rate than the analysis
- Writing of such "multiplied" data back to the Datalake
  - Defining different RSEs
- Exercises include the analysis of data stored in the Datalake
  - Writing back the results into the Datalake (small files of ~100's kb size each)
  - Analysis can be perform using CLI or the JupyterLab UI
- Create clear instructions for users/computers that can be part of the challenge

# Data Multiplication
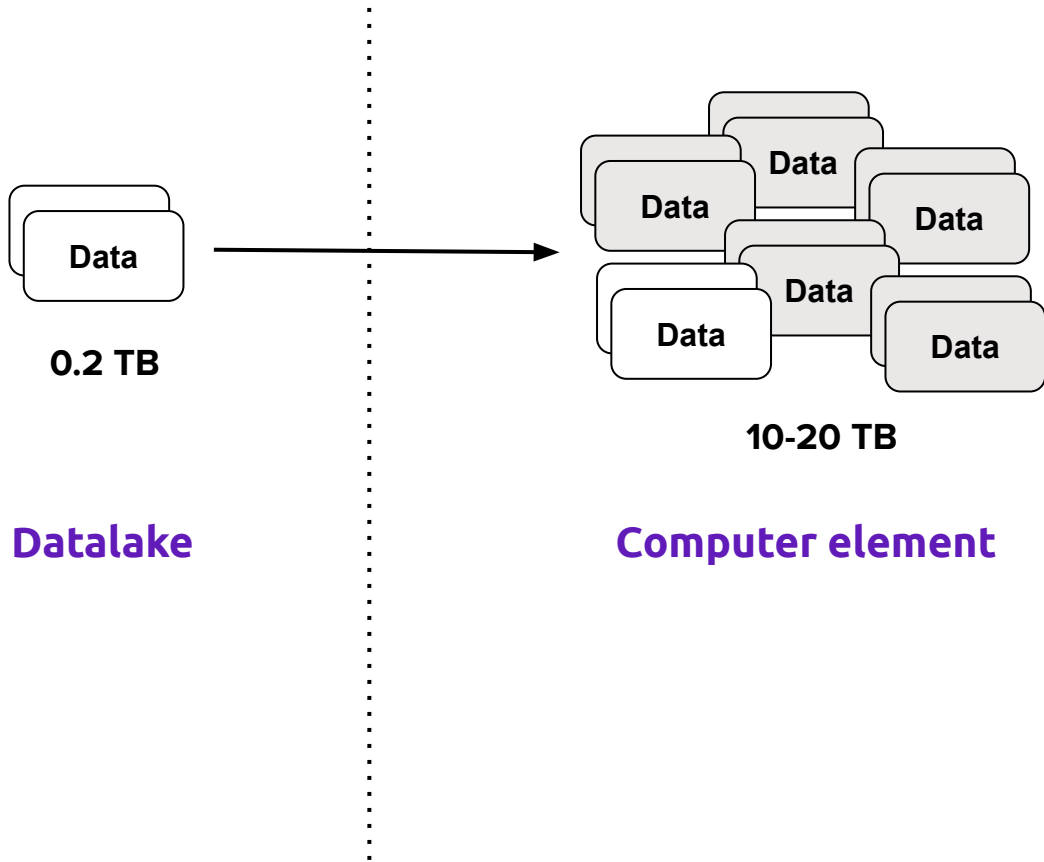
# Artificial multiplication of the data

ROOT files can be added when they share an internal structure (i.e. same trees inside). Called **"hadd"**

We can profit from that property to artificially multiply the datasets.

This process allows augmenting the data to any arbitrary value.

We can use that augmented data to run the analysis examples
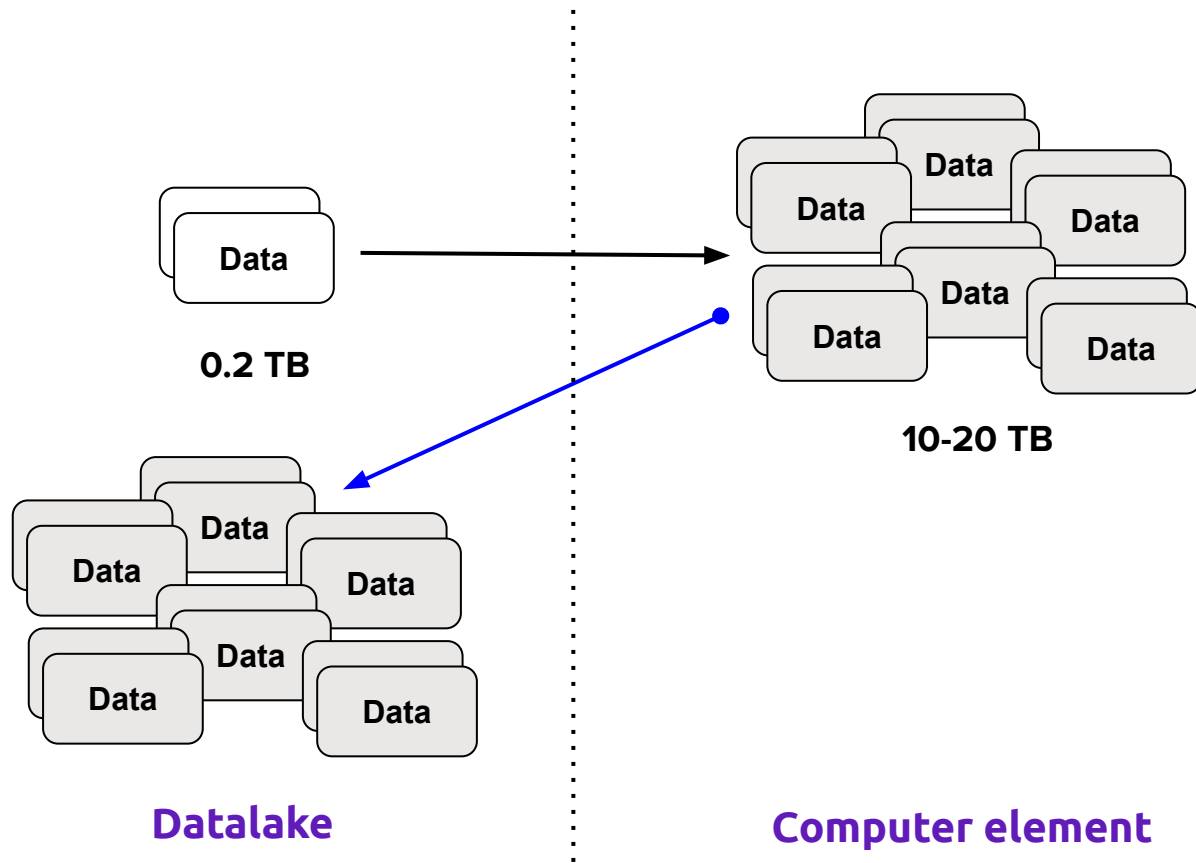● Of course, the results are meaningless.

**Data**

**0.2 TB**

**Datalake**

**Data**
**Data**
**Data**
**Data**
**Data**
**Data**

**10-20 TB**

**Computer element**

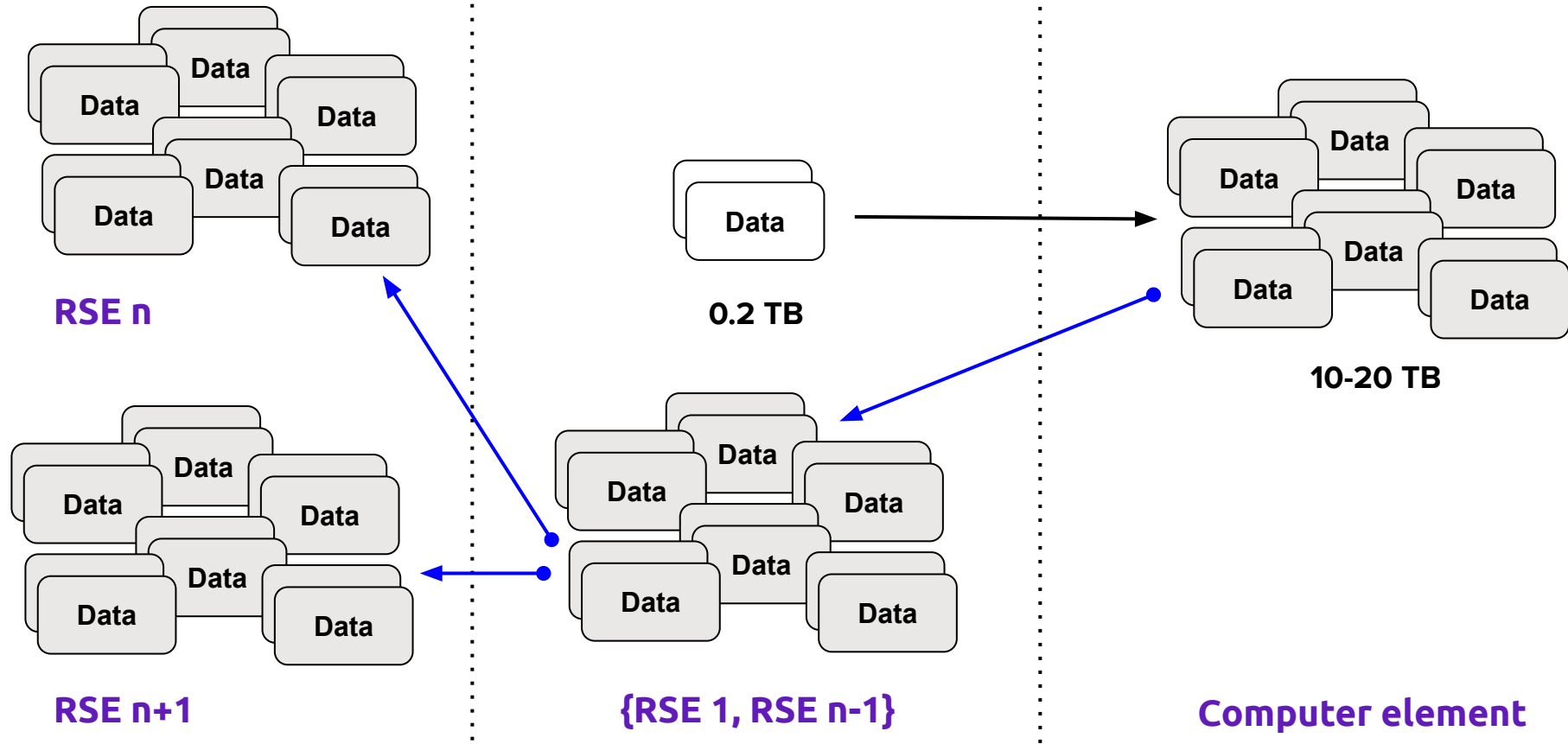# Artificial multiplication of the data

The "augmented" data
- More and/or larger files
- The process to add multiples files also use computing

After (or during the process), the code automatically can write such new data back to the Datalake
- We can also replicate in multiples RSEs as part of the challenge's tests

**Data**

**0.2 TB**

**Data** **Data** **Data** **Data** **Data** **Data**

**10-20 TB**

**Data** **Data** **Data** **Data** **Data** **Data**

**Datalake**

**Computer element**

# Artificial multiplication of the data



**RSE n**

**0.2 TB**

**10-20 TB**

**RSE n+1**
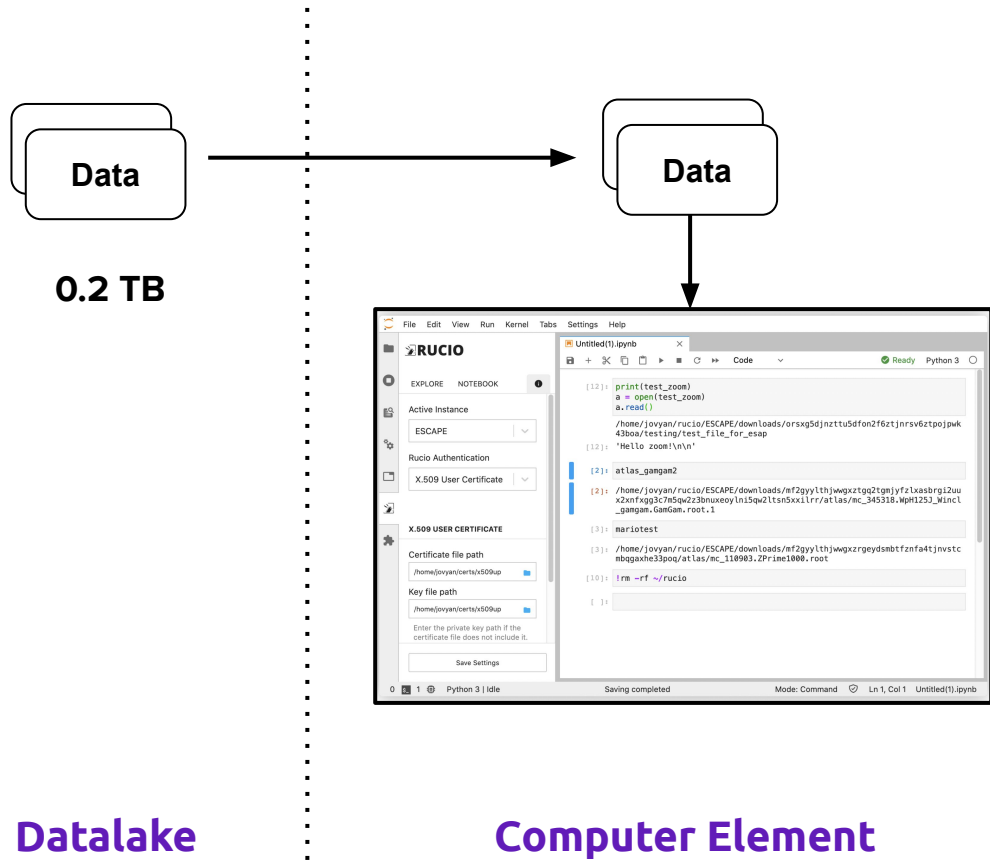
**{RSE 1, RSE n-1}**

**Computer element**

# Data Analysis

# Analysis examples

We can use the current ATLAS Open Data analysis examples to retrieve and use datasets from the Datalake

- Analysis can be notebooks or analysis frameworks
- They can take from a few minutes (e.g. 5-30 min)
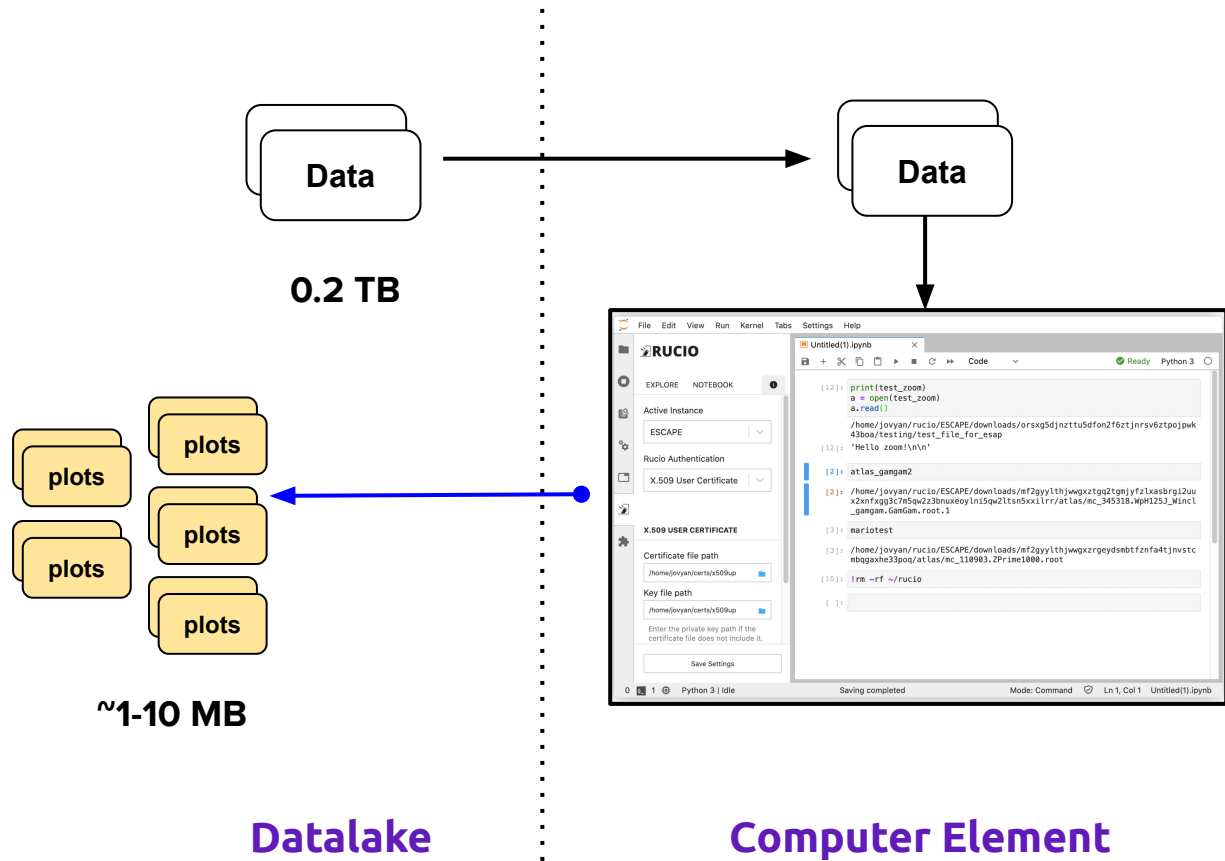- To several (e.g. 4) hours

Also, write back the outputs



**0.2 TB**

**Datalake**

**Computer Element**

# Analysis examples

The outputs of the analysis can be upload to the Datalake

- The outputs are small; they are plots that can also be store in ROOT files
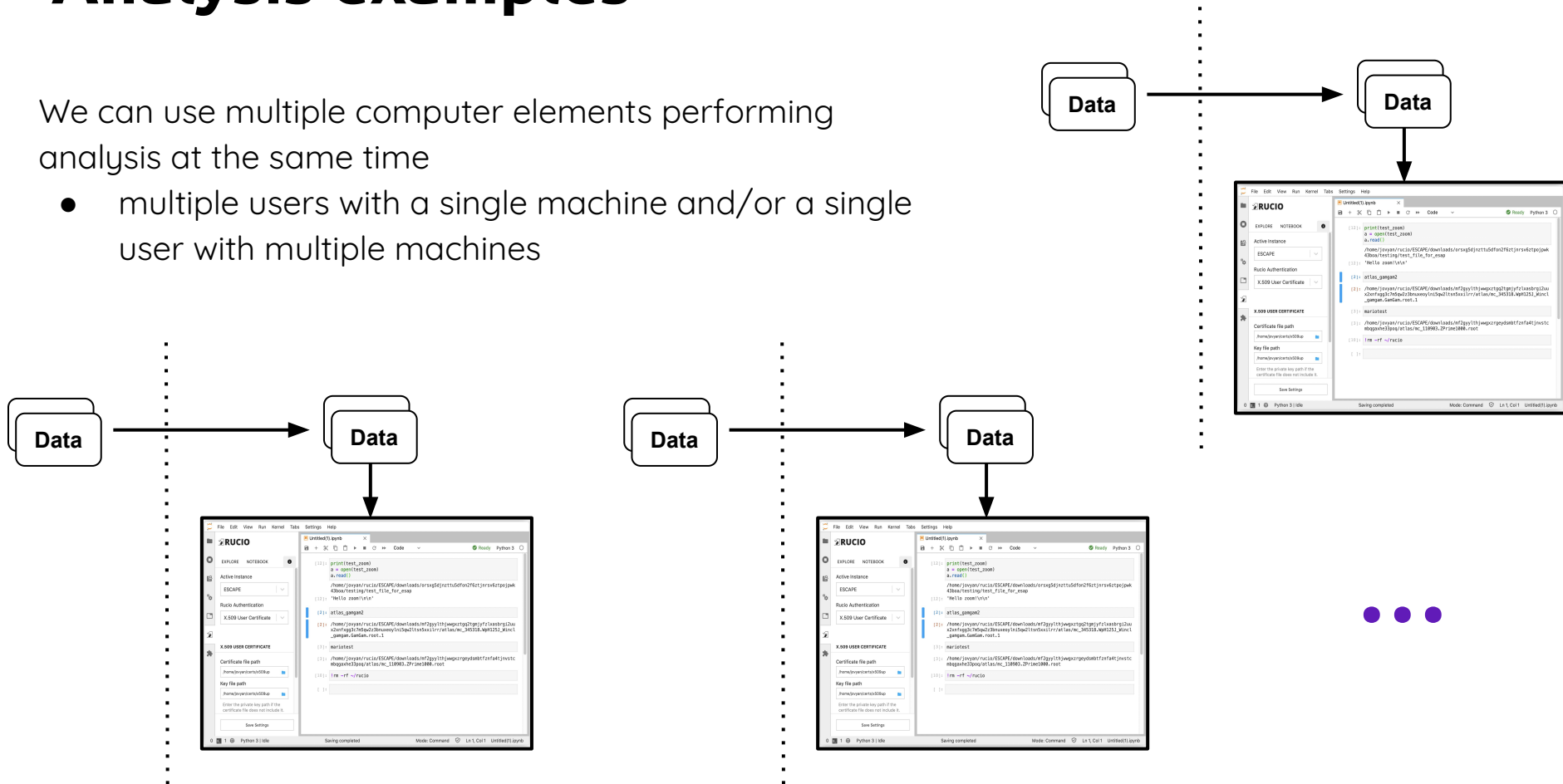- No intention to upload single PNG files to The Datalake

**Data**

**0.2 TB**

**plots**
**plots**
**plots**
**plots**
**plots**
**plots**

**~1-10 MB**

**Data**



**Datalake**                    **Computer Element**

# Analysis examples

We can also run the analysis examples over the "multiplied" data

- This can help to simulate longer analysis that can last several hours (e.g. ~8-12 hours)
- In case this kind of "stress" is useful in this challenge



**0.2 TB**

**~TB**

**~10's MB**

**Datalake**

**Computer Element**

# Analysis examples

We can use multiple computer elements performing analysis at the same time

- multiple users with a single machine and/or a single user with multiple machines

# Summary

This set of slides intends to look for feedback on the possible activities to be performed during the next ESCAPE data and analysis challenge

Many of the processes will require the creation of scripts to "automatise" some of those processes

As one of the previous step, I would like to replicate those ~220GB of open data to as many sites as possible (already existing in Fr&It RSEs)

Another essential part of the job will be to prepare instructions so others can replicate the same exercises and simulate a more realistic scenario: multiple users.

# Backup

# ATLAS Open Data datasets in the Datalake

- ROOT yet need to be importable from a notebook

  - It is deployed for testing in [DockerHub](DockerHub)

- Add more datasets to the Datalake

  - All the 13 TeV and 8 TeV ATLAS Open Data samples

  - 16 datasets → 940 samples (ROOT files)

  - < 200 GB

  - Scope used: **ATLAS_OD_EDU** (for **ATLAS O**pen **D**ata for **EDU**cation)

  - Source of the datasets:

    [http://opendata.atlas.cern/samples-13tev/](http://opendata.atlas.cern/samples-13tev/) & [http://opendata.atlas.cern/samples-8tev/](http://opendata.atlas.cern/samples-8tev/)

- Another [set of 10 ROOT files](set of 10 ROOT files) to come (dedicated Jet MC samples) → 1 dataset, ~21 GB.

# ATLAS Open Data → C++ examples framework

**To run C++ analyses**

More computational-complex particle physics analysis examples using the existing publicly available data

More in Opendata.atlas.cern - documentation 13 TeV - physics

Also use PROOF, adding a parallel component to the examples.

# A view to the current container

current **`.gitlab-ci.yml`** to build & deploy

**write / update code**    **Add, commit, and push**

**GitLab deploys and publish the container**

commit  push

commit  push

✔  ✔

**Existing container**

container registry at
**`in2p3.fr`** and **`docker.io`**

updates goes in the **`.gitlab-ci.yml`**

## Container CI / CD

- The series of resources is package in a single container
- The CI setup automatically handles the publication of the container

Frédéric & Berkay's job

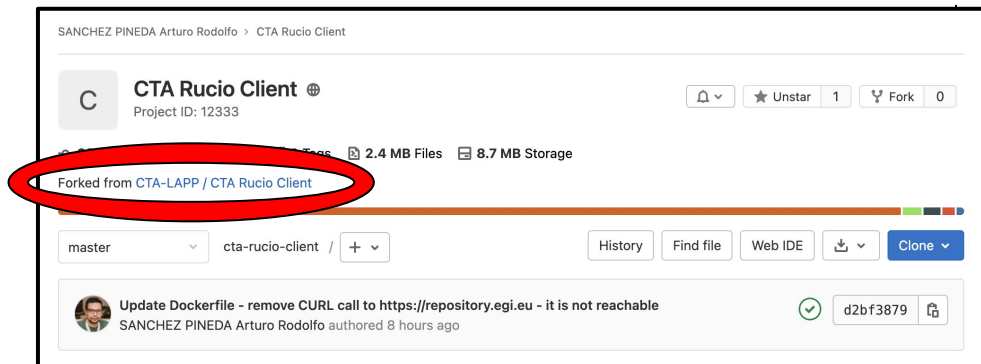**Several developments and deployment already in place**

- The compendium of resources includes the current rucio client + JupyterLab + RUCIO extension, proxy & authentication, …
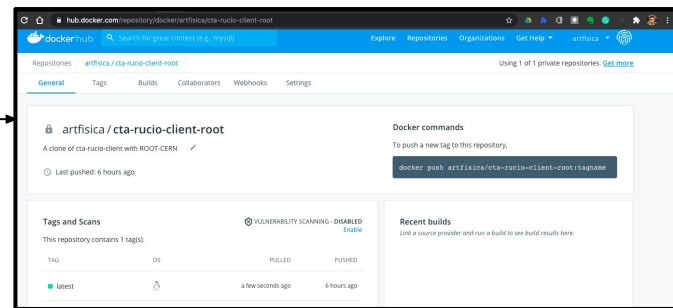
# A view to the current container

current **`.gitlab-ci.yml`** to build & deploy

write / update code    Add, commit, and push

GitLab deploys and publish the container



updates goes in the **`.gitlab-ci.yml`**
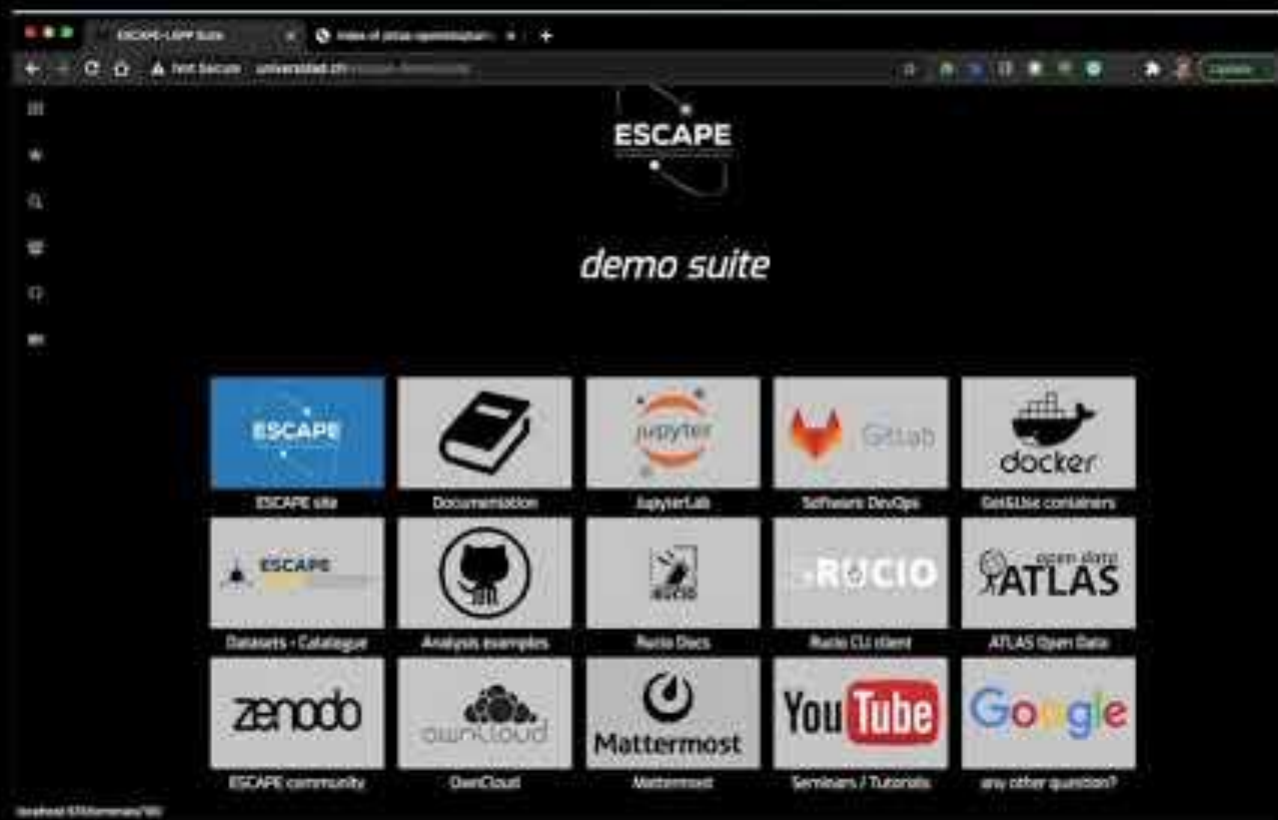
## Container CI / CD

- The series of resources is package in a single container
- The CI setup automatically handles the publication of the container

**Several tools and updates added**

- Mainly ROOT + some dependencies and extra tools...
- Jupyter conf file to handle the usage of the rucio extension (Muhammad feedback, see later)
- From JupyterLab-3 the widgets are installed using ipywidgets instead of labextension
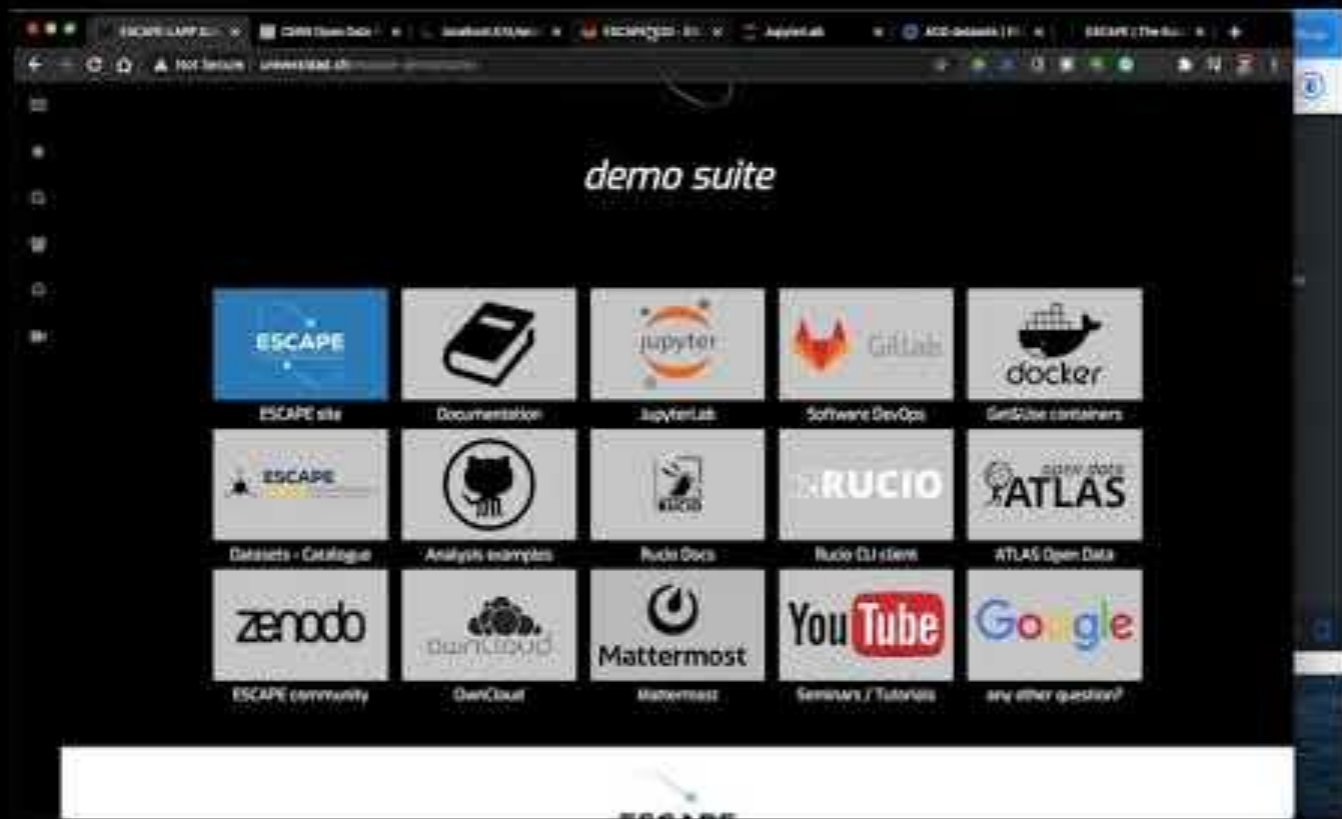
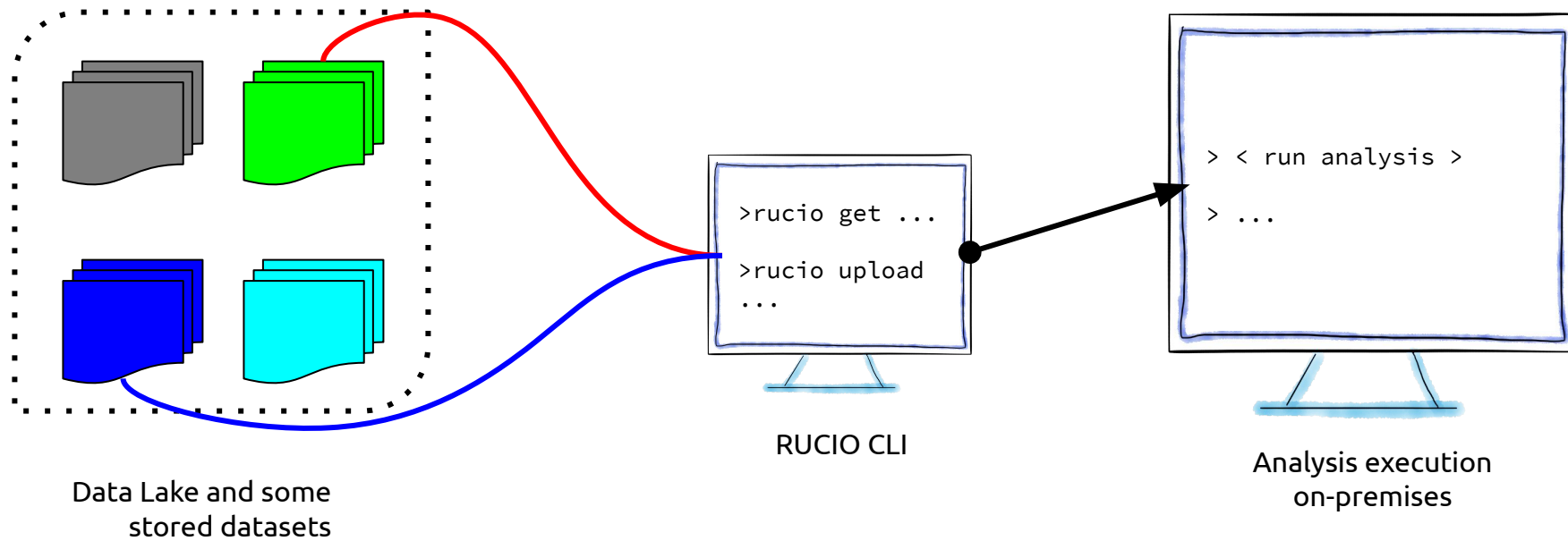# The RUCIO CLI client

(a 90 sec video, mainly for new users)
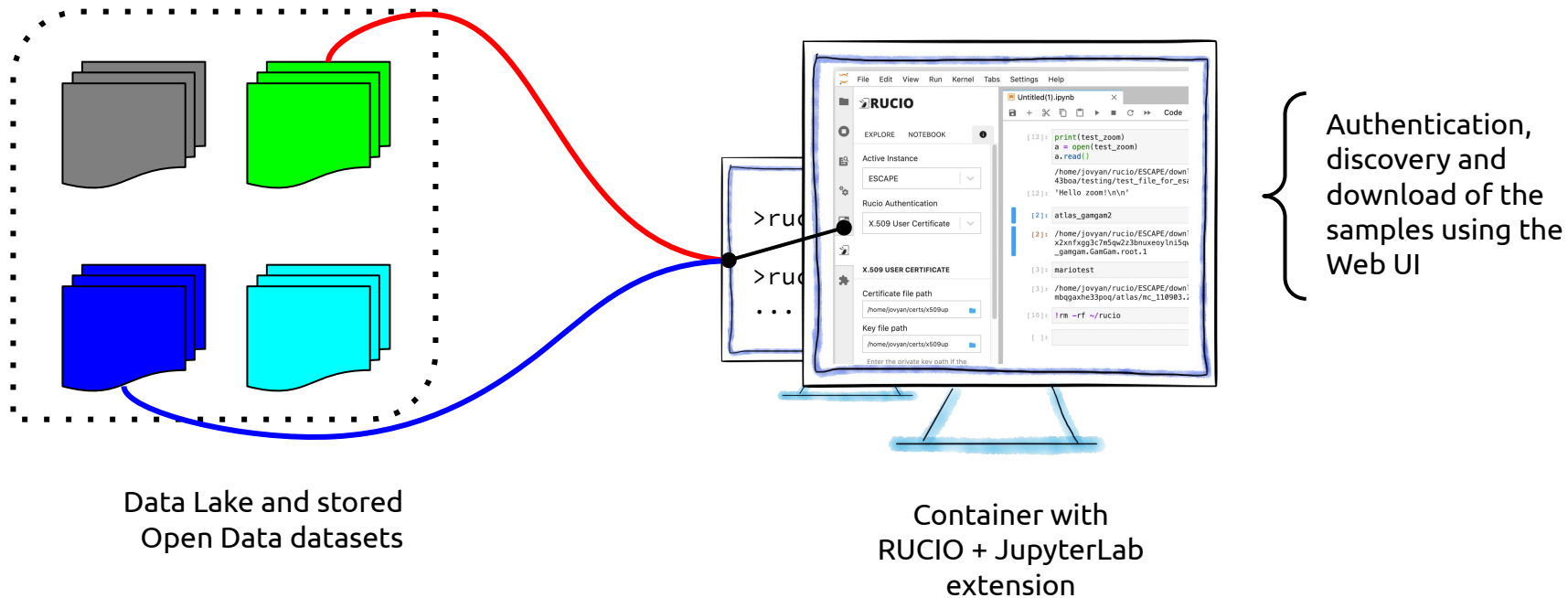
Ongoing developments with **JupyterLab & RUCIO extension**

(a 150 sec video)

More tools to finish to integrate in the container, like more kernels, PROOF, CVMFS

Data Lake and some
stored datasets

RUCIO CLI

```
>rucio get ...

>rucio upload
...
```

```
> < run analysis >

> ...
```

Analysis execution
on-premises

**CLI interaction with samples**

Authentication, discovery and download of the samples using the Web UI

Data Lake and stored Open Data datasets

Container with RUCIO + JupyterLab extension

**RUCIO+JupyterLab (container) interaction for users**