# Open and Reproducible Research Services in LHC Particle Physics
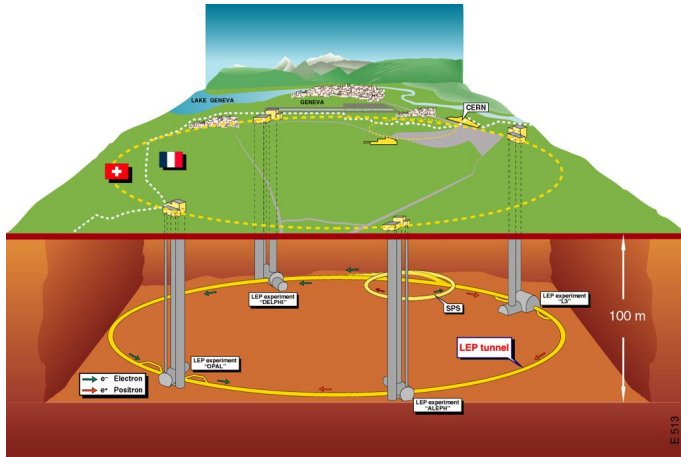
**Diego Rodríguez**
**CERN**

# CERN Large Hadron Collider
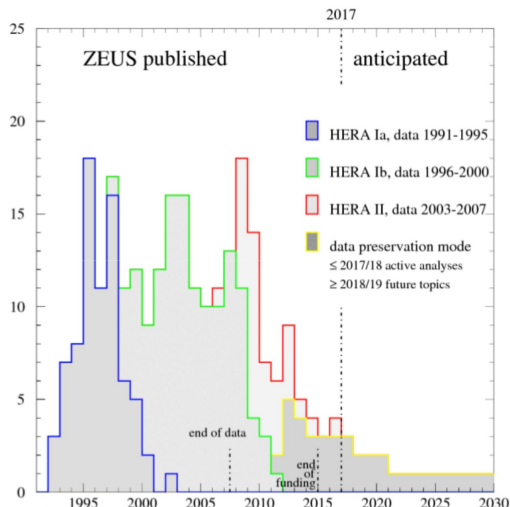


http://cds.cern.ch/record/842153



https://cds.cern.ch/record/910381

# Data and knowledge



Achim Geiser https://indico.cern.ch/event/588219

https://twitter.com/PKoppenburg/status/1301813341460066304

3

# HEP data analyses



Experimental collisions — Raw data O(PB/s) → Hardware and software trigger — Filtered data O(GB/s) → Reconstruction → Recon-structed events → Reprocessing — Analysis-formatted data O(weeks) / O(days) → Statistical analysis — O(hours) → Physics results

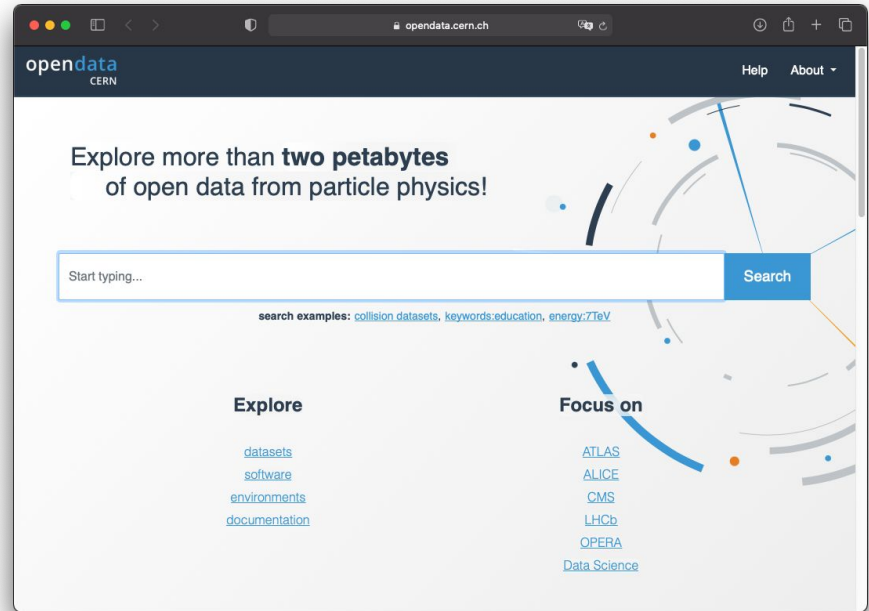Generated events — Stable particles → Interaction with simulated detector →

# CERN Open Data portal

Launched in 2014

Disseminating over 2.4 PB of data

7.500 records

900.000 files



http://opendata.cern.ch/

# Information organisation



Faceted  search

Large file download $O(GB)$

JSON Schema

Data provenance

# Education use cases



https://cds.cern.ch/record/1994217

# Research use cases

# Analysis examples

code

+

data

+

environment



http://opendata.cern.ch/record/5500

# The four questions

where is **data**?

hard drive, distributed storage

where is the **code**?

GitLab, local copy, email

what **environment** do you use?

my own laptop, remote server

what **workflow** do you use**?**

Interactive commands, bash script, README file

# Reproducible analyses



http://reana.io/



https://github.com/reanahub

# Example

```
version: 0.6.0
inputs:
  files:
    - code/gendata.C
    - code/fitdata.C
  parameters:
    events: 20000
    data: results/data.root
    plot: results/plot.png
workflow:
  type: serial
  specification:
    steps:
      - name: gendata
        environment: 'reanahub/reana-env-root6:6.18.04'
        commands:
        - mkdir -p results && root -b -q
'code/gendata.C(${events},"${data}")'
      - name: fitdata
        environment: 'reanahub/reana-env-root6:6.18.04'
        commands:
        - root -b -q 'code/fitdata.C("${data}","${plot}")'
outputs:
  files:
    - results/plot.png
```
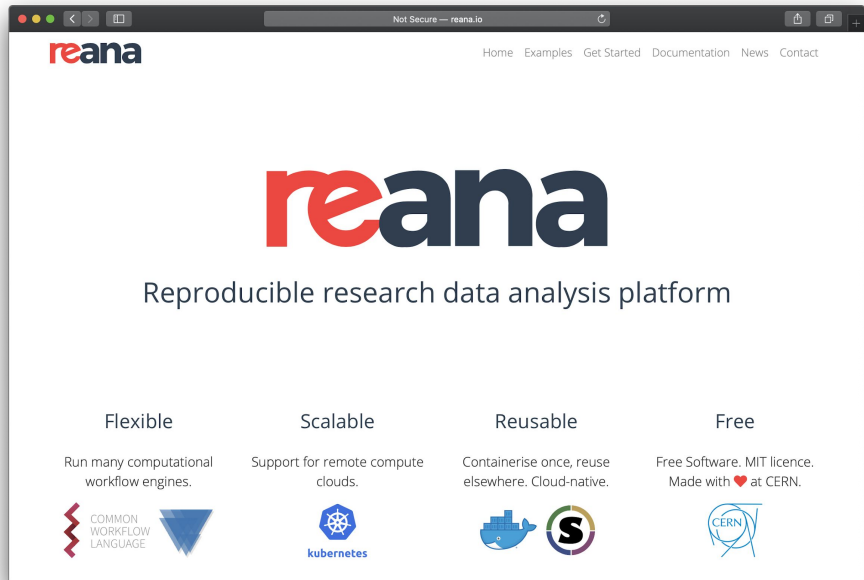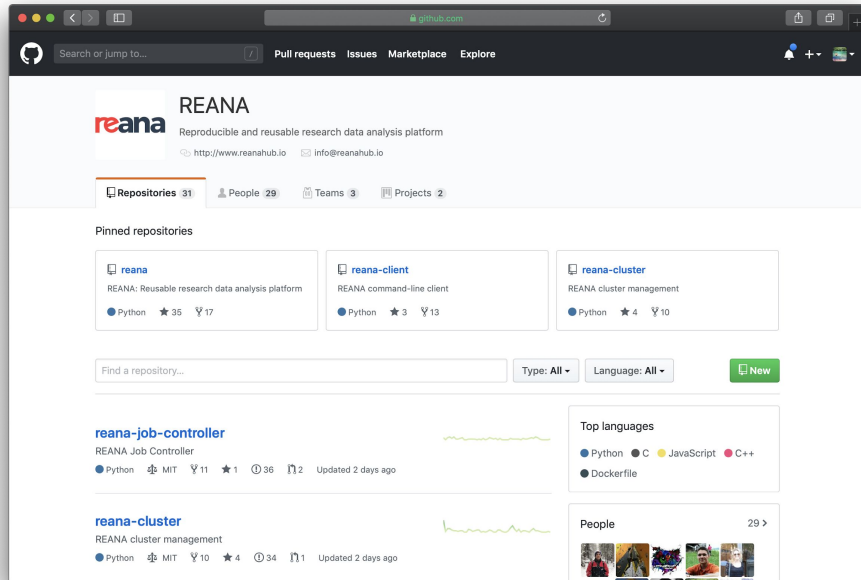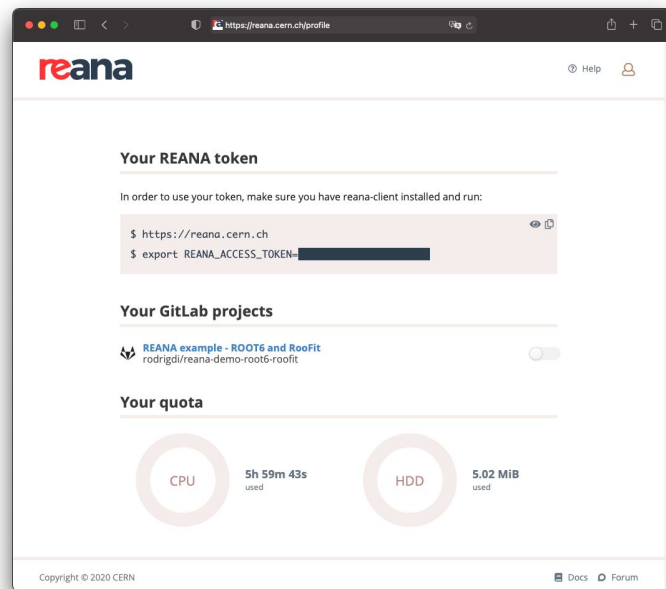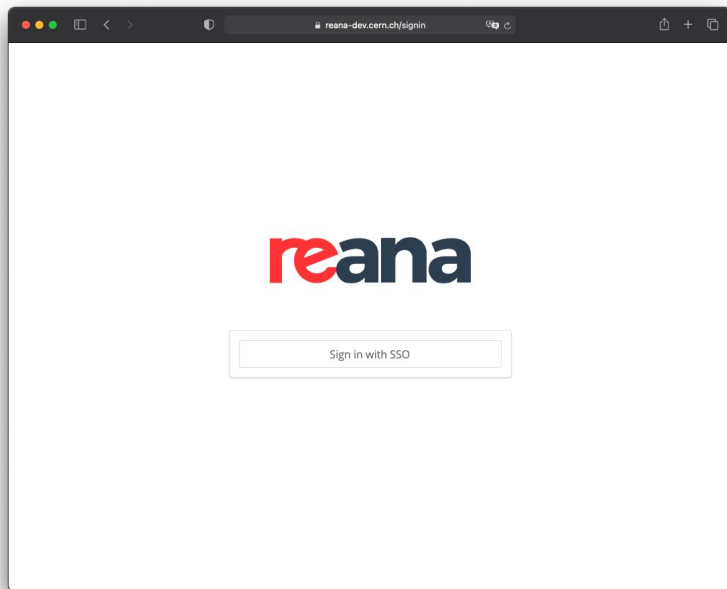
*Inputs* including: **code**, **data**, **parameters**

*Workflow* including: **steps/commands**, **container images**

*Outputs* including: **plots**, **data etc...**

# Example

# Example

# Example

# Technology: repository



https://inveniosoftware.org

The technology behind:
**https://www.hepdata.net/**
**https://inspirehep.net/**
**https://zenodo.org**

🐍 Python

𝓕𝓵𝓪𝓼𝓴 Flask

⚛️ React

{∕} JSONSchema

📦 Redis

🔴 Elasticsearch

🐘 PostgreSQL

# Technology: REANA

- ○ Cloud-native application
- ○ Extensible
  - ■ Storage backends
  - ■ Compute backends
  - ■ Container technologies
  - ■ Workflow engines

# Try them out!

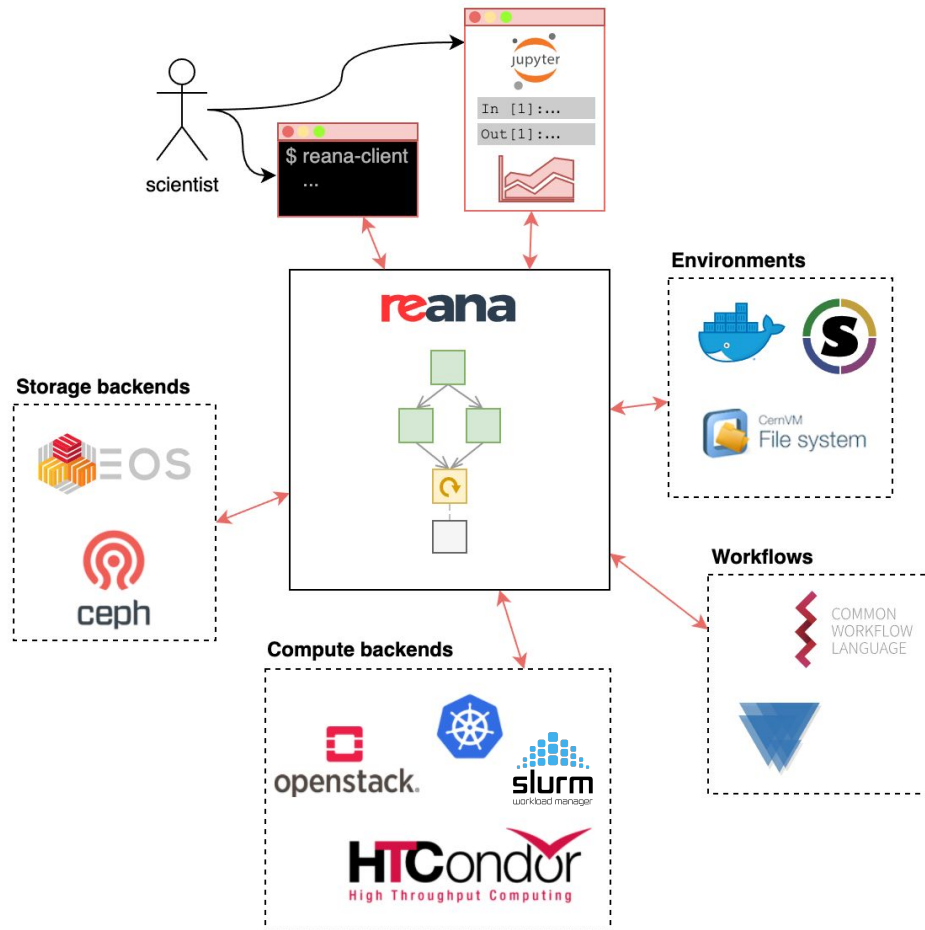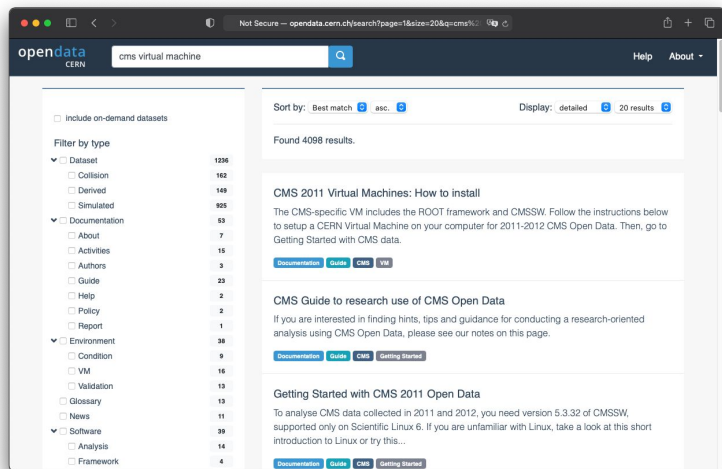## Use [opendata.cern.ch](opendata.cern.ch)



## Install REANA on premises/locally
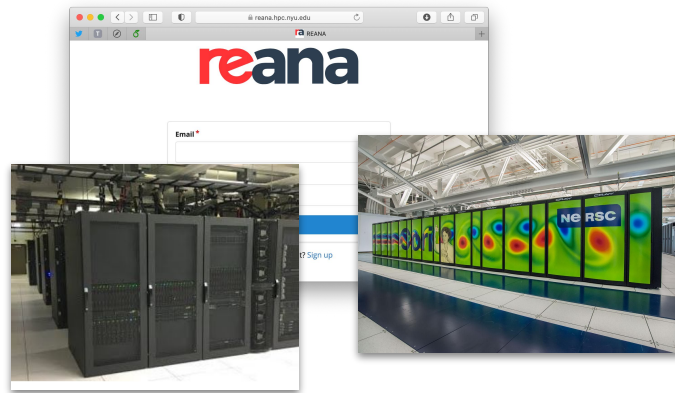


```
$ helm repo add reanahub \
    https://reanahub.github.io/reana

$ helm repo update

$ helm install reana reanahub/reana
```

Helm repository at [http://reanahub.github.io/reana](http://reanahub.github.io/reana), more documentation at [http://docs.reana.io/development/deploying-at-scale/](http://docs.reana.io/development/deploying-at-scale/)

# What's next

# Get in touch

## CERN Open Data

https://opendata.cern.ch

https://github.com/cernopendata

https://forum.opendata.ch/

https://gitter.im/cernopendata/opendata.cern.ch

https://twitter.com/cernopendata

## REANA

https://www.reana.io

https://github.com/reanahub/reana

https://forum.reana.io/

https://gitter.im/reanahub/reana

https://twitter.com/reanahub

https://docs.reana.io