# Reproducible high energy physics analyses

Diego Rodríguez Rodríguez
CERN

# High turnover of researchers



"Particle Physics author lists change with time. Here that of the first @LHCbExperiment paper in 2010. Violet: still in LHCb. Blue: left LHCb" - @PKoppenburg

# Use cases

# Three pillars

Describe → Capture → Reuse

# CERN Analysis Preservation

A platform for **preserving knowledge** and **assets** of an individual physics analysis

Capturing the elements needed to **understand** and **rerun** an analysis even several years after

Advanced **search** for high-level physics information

Applying standard **collaboration access restrictions**

# CERN Analysis Preservation

Describe    Capture



- JSONSchema

- W3C DCAT

- domain-specific-fields

- collaborative capabilities

- grabbing from Git, HTTP, XRootD

- powerful search

INVENIO

# REANA: Reusable Analyses

**Reuse**

Born as part of the CERN Analysis Preservation framework

A platform which enables the **reusability** of physics analysis

Based in **cloud technologies**

Following the **12-factor application pattern**

# Technology: REANA

# Four questions

(1) data?

(2) software?

(3) environment?

(4) workflow?

✓ EOS, CephFS

✓ Git, SVN, local machines

✓ Docker, VMs

✓ CWL, Yadage

# (1) Data

- Stored on CERN Analysis Preservation

- CephFS

- EOS

- ....



**root file**



**CSV file**

# (2) Software



```cpp
#ifndef __CINT__
#include "RooGlobalFunc.h"
#endif
#include "RooRealVar.h"
#include "RooDataSet.h"
#include "RooGaussian.h"
#include "RooChebychev.h"
#include "RooAddPdf.h"
#include "RooExtendPdf.h"
#include "TCanvas.h"
#include "TAxis.h"
#include "RooPlot.h"
using namespace RooFit ;

void fitdata(const char* input, const char* output)
{
  // Open input file with workspace (generated by rf14_wspacewrite)
  TFile *f = new TFile(input) ;

  // Retrieve workspace from file
  RooWorkspace* w = (RooWorkspace*) f->Get("w") ;

  // Retrieve x,model and data from workspace
  RooRealVar* x = w->var("x") ;
  RooAbsPdf* model = w->pdf("model") ;
  RooAbsData* data = w->data("modelData") ;

  // Fit model to data, extended ML term automatically included
  model->fitTo(*data) ;

  // Plot data and PDF overlaid
  RooPlot* xframe = x->frame(Title("Fit example")) ;
  data->plotOn(xframe) ;
  model->plotOn(xframe,Normalization(1.0,RooAbsReal::RelativeExpected)) ;

  // Overlay the background component of model with a dashed line
  model->plotOn(xframe,Components("bkg"),LineStyle(kDashed),Normalization(1.0,RooAbsReal::RelativeExpected)) ;

  // Draw the frame on the canvas
  TCanvas res("rf202_composite","rf202_composite",600,600) ;
  gPad->SetLeftMargin(0.15) ;
  xframe->GetYaxis()->SetTitleOffset(1.4) ;
  xframe->Draw();

  res.Update();
  res.SaveAs(output);
  res.Close();

}
```

# (3) Environment

**Docker** support, other technologies under investigation

Encourage the usage of **base images** i.e. *reanahub/reana-env-root6* for ROOT6 analyses

Take the most out of **image layering**

Encourage **collaboration** and **reusable images**

# reana-env-root6

```dockerfile
# Environment: ROOT6 on Ubuntu/Trusty:
FROM ubuntu:trusty
RUN apt-get update
RUN apt-get install --yes g++ cpp gcc gfortran git dpkg-dev make binutils libx11-dev libxpm-dev libxft-dev libxext-dev \
                         libssl-dev libpcre3-dev xlibmesa-glu-dev libglew1.5-dev libftgl-dev libmysqlclient-dev \
                         libfftw3-dev cfitsio-dev graphviz-dev libavahi-compat-libdnssd-dev libldap2-dev python-dev \
                         libxml2-dev libkrb5-dev libgsl0-dev libqt4-dev libx11-dev libxpm-dev

ENV ROOTSYS /usr/local
RUN git clone --quiet http://root.cern.ch/git/root.git /code/root-v6-02-12 &&\
    cd  /code/root-v6-02-12 &&\
    git checkout v6-02-12 &&\
    ./configure --all &&\
    make -j4 &&\
    make -j4 install &&\
    cd / &&\
    rm -rf /code
```

# (4) Workflow

Structured computational workflows over free-text READMEs

Testable approach

Support for Yadage workflows

Support for CWL workflows

# Yadage workflow example

```yaml
stages:
  - name: gendata
    dependencies: ['init']
    scheduler:
      scheduler_type: singlestep-stage
      parameters:
        events: {stages: init, output: events, unwrap: true}
        outfilename: '{workdir}/data.root'
      step:
        process:
          process_type: 'interpolated-script-cmd'
          script: root -b -q 'gendata.C({events},"{outfilename}")'
        publisher:
          publisher_type: 'frompar-pub'
          outputmap:
            data: outfilename
        environment:
          environment_type: 'docker-encapsulated'
          image: johndoe/reana-demo-root6-roofit
  - name: fitdata
    dependencies: ['gendata']
    scheduler:
      scheduler_type: singlestep-stage
      parameters:
        data: {stages: gendata, output: data, unwrap: true}
        outfile: '{workdir}/plot.png'
      step:
        process:
          process_type: 'interpolated-script-cmd'
          script: root -b -q 'fitdata.C("{data}","{outfile}")'
        publisher:
          publisher_type: 'frompar-pub'
          outputmap:
            plot: outfile
        environment:
          environment_type: 'docker-encapsulated'
          image: johndoe/reana-demo-root6-roofit
```



**Lukas Heinrich** <u>**diana-hep/yadage**</u>

# CWL workflow example

ATLAS full chain analysis example



sig



mc



data

# How does it work?

# Deploy REANA locally

```
(RCLUSTER-2.7) ➜  reana-cluster git:(master) ✗ minikube start --kubernetes-version="v1.6.4"

Starting local Kubernetes v1.6.4 cluster...
Starting VM...
Moving files into cluster...
Setting up certs...
Starting cluster components...
Connecting to cluster...
Setting up kubeconfig...
Kubectl is now configured to use the cluster.
```

# Deploy REANA locally

```
(RCLUSTER-2.7) → reana-cluster git:(master) x reana-cluster init
[INFO] Validating REANA cluster specification file: /Users/rodrigdi/reana/reana-cluster/reana-cluster.yaml
[INFO] /Users/rodrigdi/reana/reana-cluster/reana-cluster.yaml is a valid REANA cluster specification.
[INFO] Cluster type specified in cluster specifications file is 'kubernetes'
[INFO] Creating a ReanaBackend object for Kubernetes interaction.
[INFO] Connecting to Kubernetes at https://192.168.99.100:8443
[INFO] Writing deployable REANA cluster configuration to ./cluster_config/
Init complete
```

# Deploy REANA locally

```
(RCLUSTER-2.7) → reana-cluster git:(master) x reana-cluster get server
[INFO] Validating REANA cluster specification file: /Users/rodrigdi/reana/reana-cluster/reana-cluster.yaml
[INFO] /Users/rodrigdi/reana/reana-cluster/reana-cluster.yaml is a valid REANA cluster specification.
[INFO] Cluster type specified in cluster specifications file is 'kubernetes'
[INFO] Creating a ReanaBackend object for Kubernetes interaction.
external_name: None
internal_ip: None
external_ip_s: 192.168.99.100
ports: ['31201']
```

# Example analysis



```
(reana-client) ➜ reana-demo-helloworld git:(master) tree
.
├── README.rst
├── code
│   └── helloworld.py
├── environment
│   └── Dockerfile
├── inputs
│   └── names.txt
├── outputs
├── reana.yaml
└── workflow
    └── yadage
        └── workflow.yaml

6 directories, 6 files
```

**Available at https://github.com/reanahub/reana-demo-helloworld**
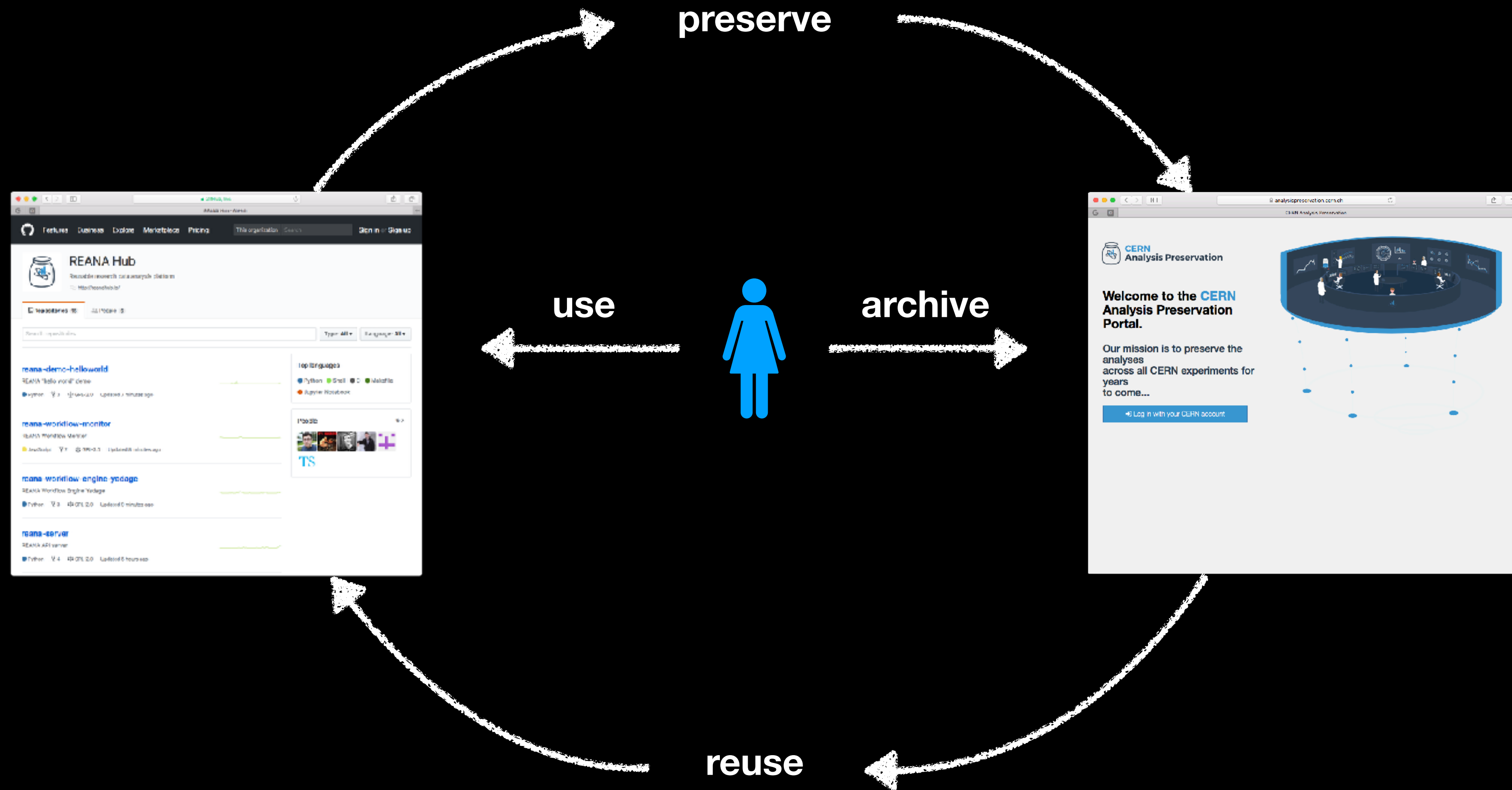
# Run on REANA

```
(reana-client) ➜  reana-demo-helloworld git:(master) export REANA_SERVER_URL=http://192.168.99.100:31201
(reana-client) ➜  reana-demo-helloworld git:(master) reana-client ping
[INFO] REANA Server URL ($REANA_SERVER_URL) is: http://192.168.99.100:31201
[INFO] Connecting to http://192.168.99.100:31201
[INFO] Server is running.
```
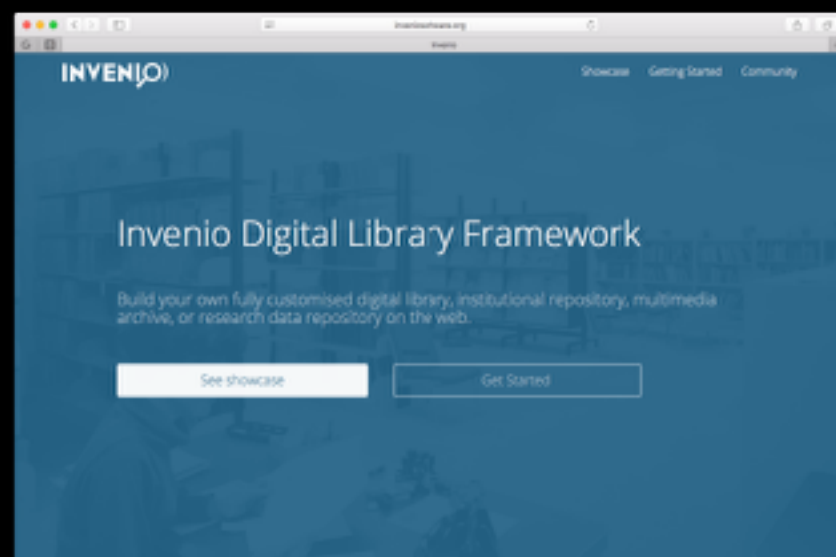
# Run on REANA

```
$ reana-client workflow create
$ reana-client code upload helloworld.py
$ reana-client inputs upload names.txt
$ reana-client workflow start
$ reana-client workflow status
# wait until the workflow finishes
$ reana-client outputs list
$ reana-client outputs download helloworld/greetings.txt
```

# Reusability    Preservation



preserve

use    archive

reuse

# Challenges

**Social**
- adopting structured computational workflow specifications
- publish or perish culture
- scientific benefit vs cost of preservation

**Data**
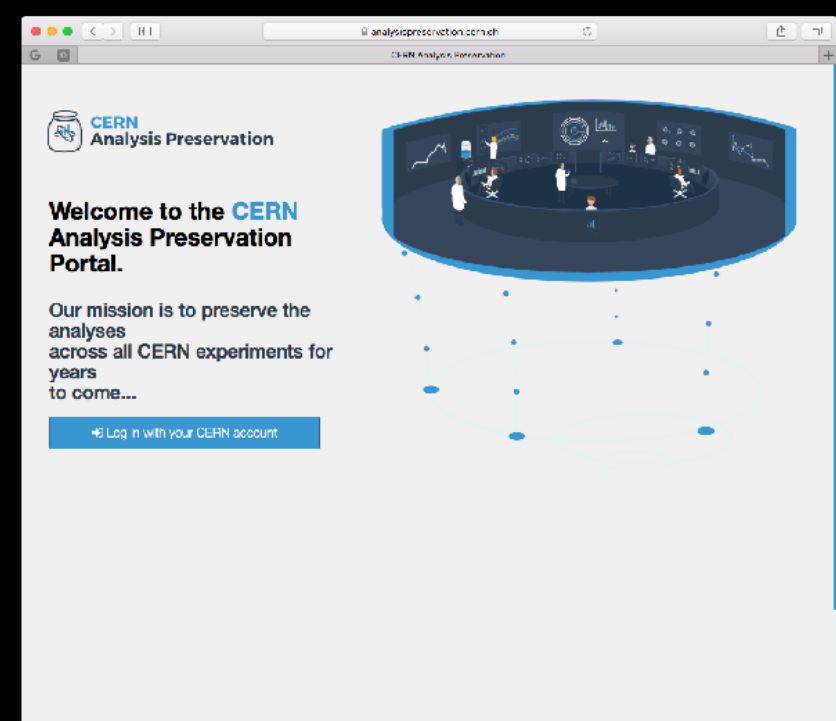- Ever-increasing data size?

**Software**
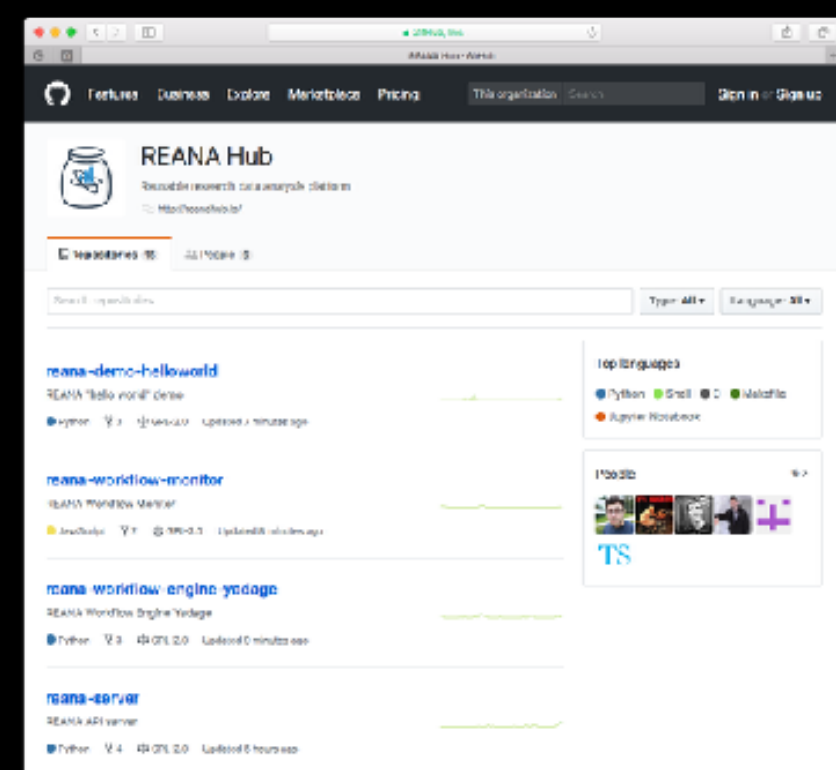- Ever-changing computing technology?

🔗 http://inveniosoftware.org

🐙 http://github.com/inveniosoftware

🐦 @inveniosoftware

✉️ info@inveniosoftware.org



🔗 http://analysispreservation.cern.ch

🐙 http://github.com/cernanalysispreservation

✉️ analysis-preservation-support@cern.ch



🔗 http://reanahub.io

🐙 http://github.com/reanahub

🐦 @reanahub

✉️ info@reanahub.io

# Questions?