



# Data Platform for collection, storage, integration, analysis and distribution

*CERN openlab Summer Student lightning talk*

Diamantis Patsidis

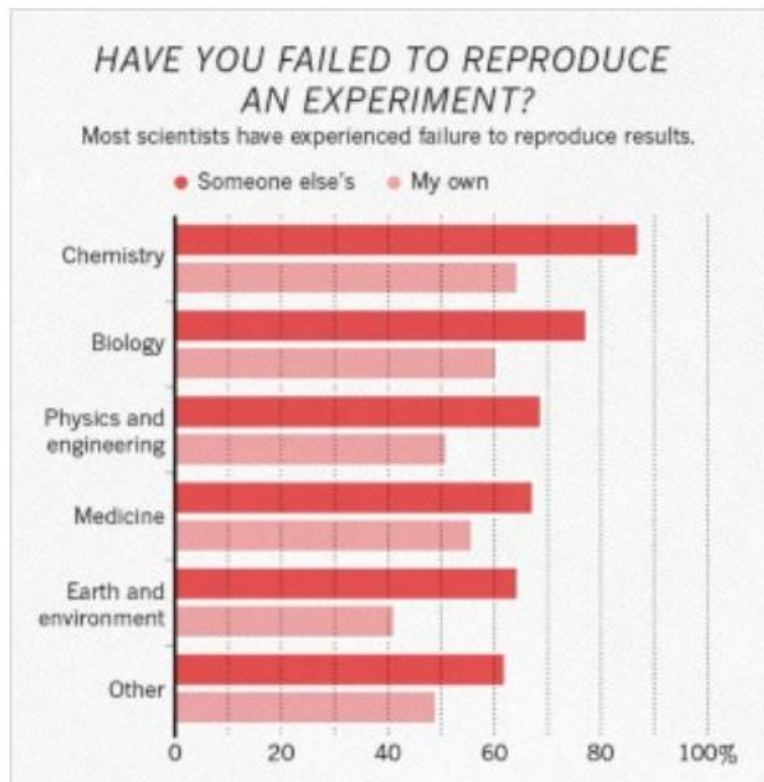
Supervisor – Anna Ferrari

06 / 09 / 2021

# Reproducibility

## What is reproducibility?

*Enable others to reproduce results without difficulty*



<https://www.nature.com/articles/533452a> (Monya Baker 2016.)

## Reproducibility rates in biomedical research

11% (Begley and Ellis, 2012, Ioannidis 2005)

22% - 49% (Freedman 2015, from 5 independent research groups)

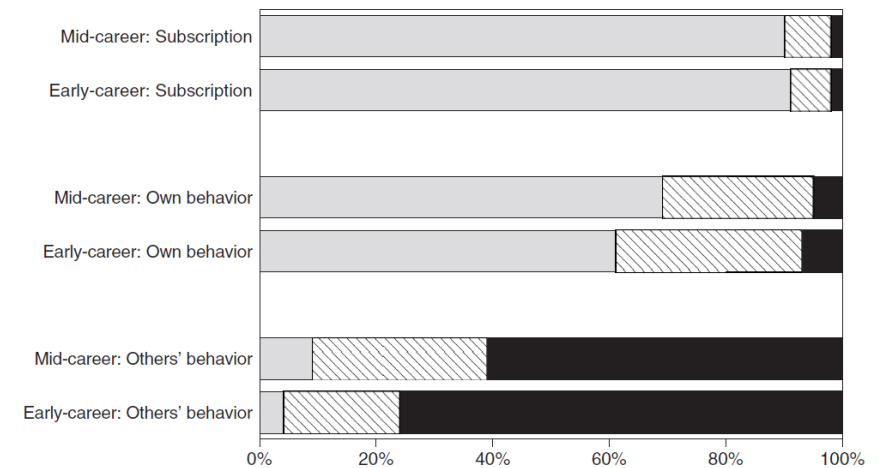


FIG. 3. Norm versus Counternorm Scores: Percent with Norm > Counternorm (dotted), Norm = Counternorm (striped), Norm < Counternorm (solid).

[Normative Dissonance in Science Results from a National Survey of US Scientists \(2007\)](#)

# Reproducibility recipe

## Four main ingredients

Data (location)

Cloud, local drive

Code (location)

GitLab, local drive

Environment (software used, ...)

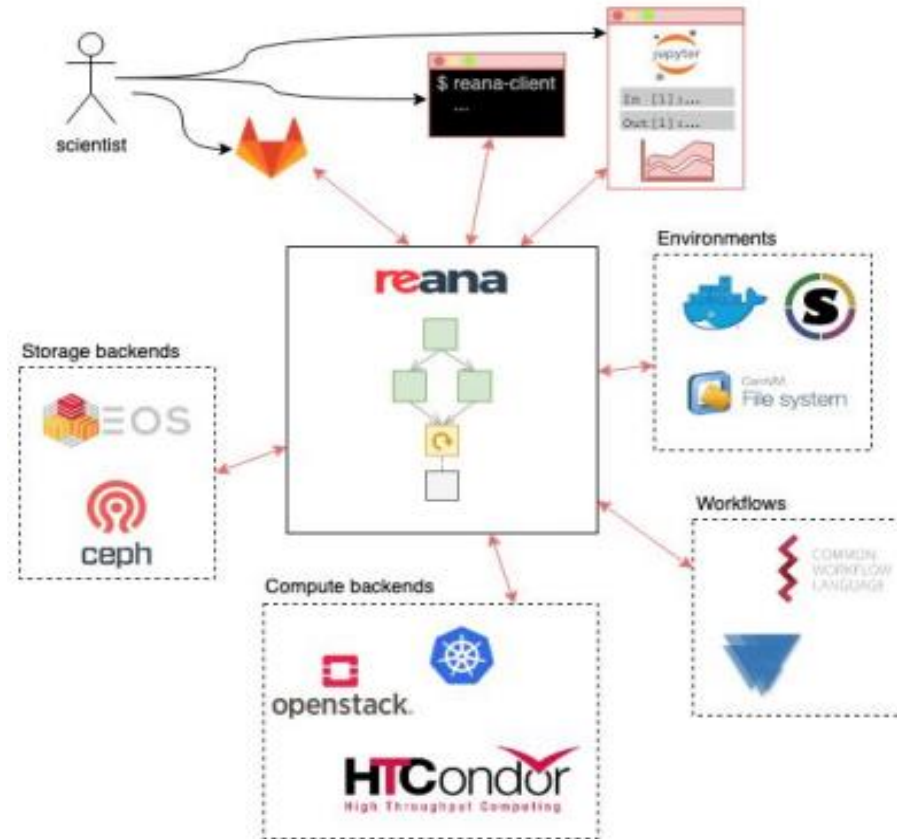
Local computer, cloud computing

Workflow (order of scripts, ...)

Bash script, README

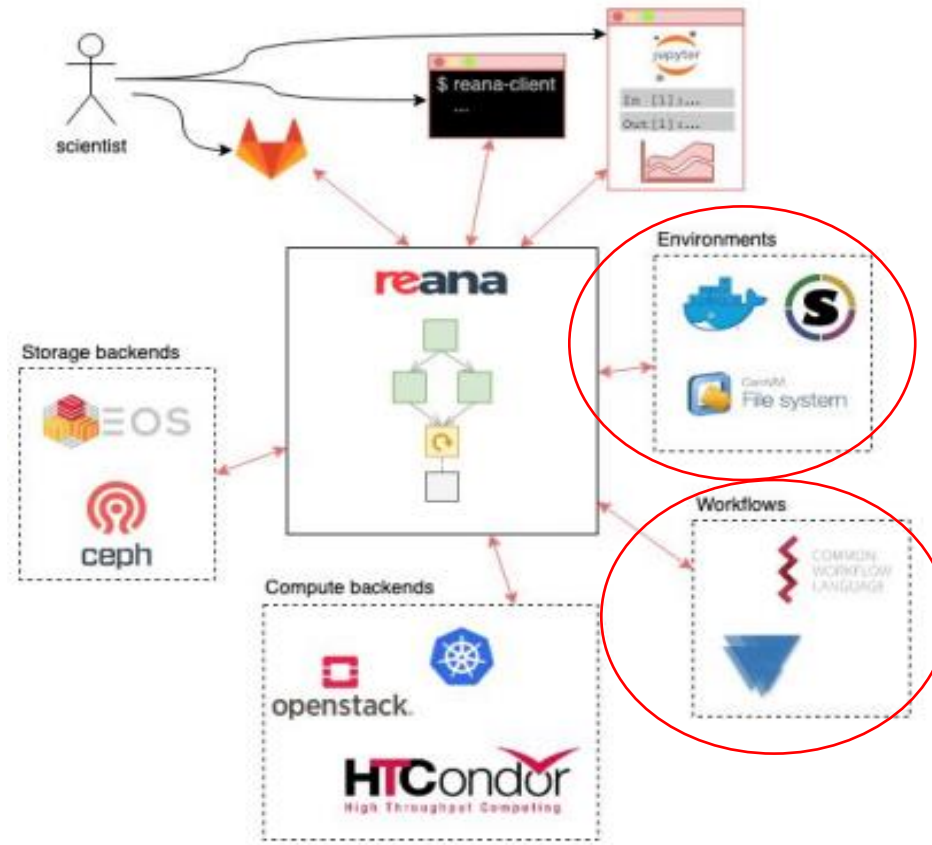
[Analysis Reproducibility with REANA on Kubernetes \(Diego Rodriguez - REANA team\)](#)

# REANA



[Analysis Reproducibility with REANA on Kubernetes \(Diego Rodriguez - REANA team\)](#)

# REANA



[Analysis Reproducibility with REANA on Kubernetes \(Diego Rodriguez - REANA team\)](#)

# Defining workflows and environments

The creation of workflows and environments usually requires knowledge in specialized technologies

```
- name: prepare
dependencies: []
scheduler:
  scheduler_type: 'singlestep-stage'
  parameters:
    model: sm
    parametercard: '{workdir}/param.dat'
    inputpars: defaultparam.yml
    step: {$ref: 'preparestep.yml'}
- name: madgraph
dependencies: ['prepare','init']
scheduler:
  scheduler_type: 'singlestep-stage'
  parameters:
    outputlhe: '{workdir}/output.lhe'
    events: {stages: init, output: nevents, unwrap: true}
    paramcard: {stages: prepare, output: parcard, unwrap: true}
    step: {$ref: 'madgraph.yml'}
- name: pythia
dependencies: ['madgraph']
scheduler:
  scheduler_type: 'singlestep-stage'
  parameters:
    outputhepmc: '{workdir}/output.hepmc'
    events: {stages: init, output: nevents, unwrap: true}
    lhefile: {stages: madgraph, output: lhefile, unwrap: true}
    step: {$ref: 'pythia.yml'}
```

Three stage workflow

<https://yadage.readthedocs.io/en/latest/definingworkflows.html>

```
1 # Environment: Jupyter 1.0.0 with IPython 5.0.0 kernel on CentOS7
2 FROM centos:7
3
4 # Set system locale
5 ENV LC_ALL=en_US.UTF-8
6 ENV LANG=en_US.UTF-8
7
8 # hadolint ignore=DL3033
9 RUN yum install -y epel-release && yum clean all
10
11 # hadolint ignore=DL3033
12 RUN yum install -y \
13     gcc \
14     python3-devel \
15     python3-pip \
16     && yum clean all
17
18 # hadolint ignore=DL3013
19 RUN pip3 install --upgrade pip
20 RUN pip3 install --no-cache-dir \
21     ipython==7.16.1 \
22     jupyter==1.0.0 \
23     jupyter-client==6.1.12 \
24     jupyter-console==6.4.0 \
25     jupyter-core==4.7.1 \
26     matplotlib==3.3.4 \
27     nbconvert==6.0.7 \
28     pandas==1.1.5 \
29     papermill==2.3.3 \
30     # FIXME: Pin black to be able to create the cache folders manually and avoid errors.
31     # Related issue: https://github.com/nteract/papermill/issues/498
32     black==21.7b0 \
33     ipykernel==5.5.5
34
35 # FIXME: Create 'black' cache folder manually
36 # hadolint ignore=SC2174
37 RUN mkdir -m 770 -p /.cache/black/21.7b0/
```

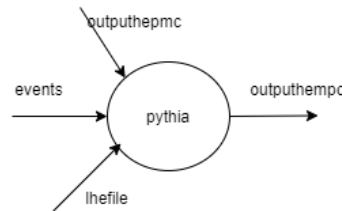
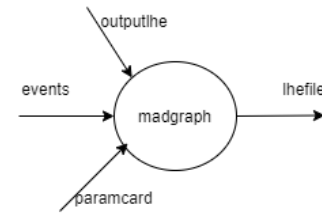
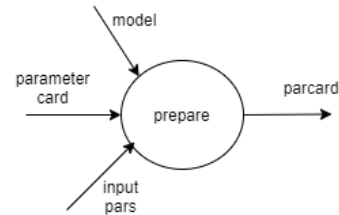
Dockerfile for Jupyter notebooks

<https://github.com/reanahub/reana-env-jupyter/blob/master/Dockerfile>

# Automating workflows

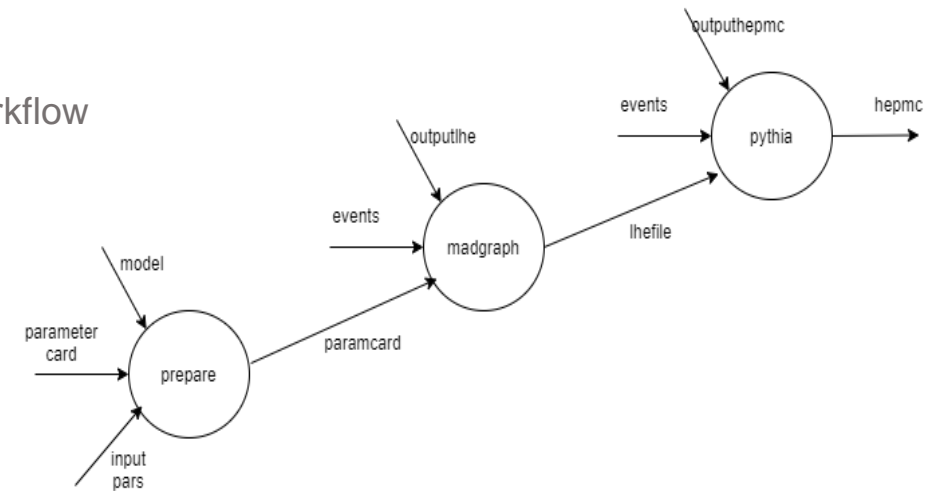
## User defines for each stage

- data/code files
- script parameters
- expected input values
- expected outputs



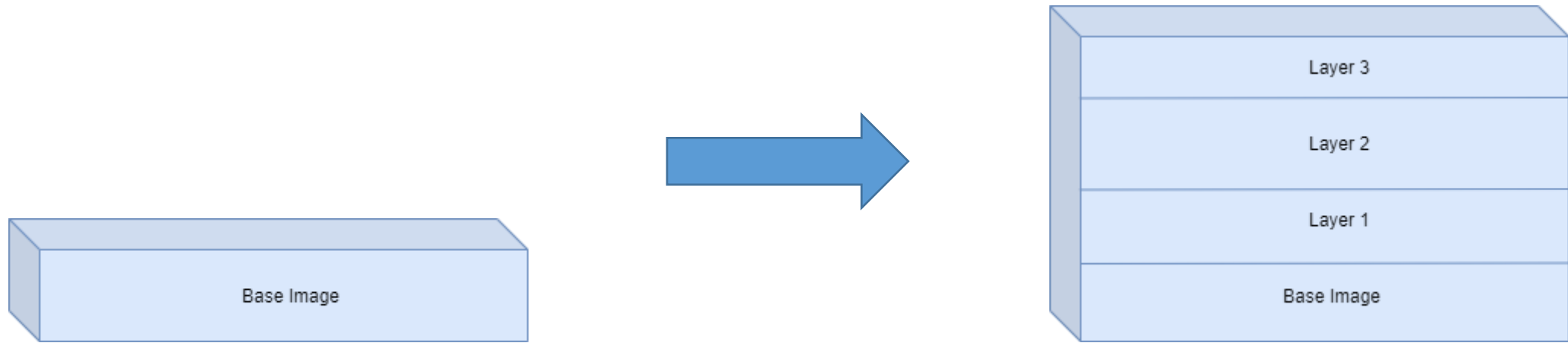
automatically define workflow

(create yaml file)



# Extent environments

Add software and packages to an environment without manually changing any files (Dockerfile)







# QUESTIONS?

*diamantis.patsidis@cern.ch*

*diampats@ee.duth.gr*