

Reproducibility and Extensibility in Scientific Research

Jessica Forde
Project Jupyter
@projectjupyter
@mybinderteam



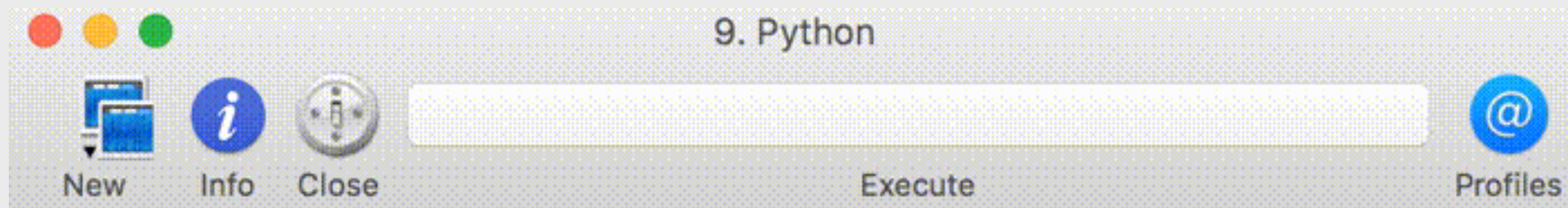
Overview

- Project Jupyter
- IPython
- Jupyter Notebook
- Architecture of JupyterHub
- The problem of reproducibility in science
- Repo2docker
- Binder
- Extending research via interactive computing



Who are we?





```
-bash-3.2$ ipython
WARNING: Attempting to work in a virtualenv. If you encounter problems, please
ninstall IPython inside the virtualenv.
Python 2.7.10 (default, Sep 23 2015, 04:34:21)
Type "copyright", "credits" or "license" for more information.

IPython 4.0.2 -- An enhanced Interactive Python.
?          -> Introduction and overview of IPython's features.
%quickref -> Quick reference.
help       -> Python's own help system.
object?    -> Details about 'object', use 'object??' for extra details.

In [1]: from postal.expand import expand_address

In [2]:
```

on-startup-files

uild_symlink
ython_analysis
*tup
0-imports.py
1-plotting.py
2-autoreload.ipynb
3-databases.py
4-pretty_printing.py
5-data_analysis.py
6-misc.py
EADME
:INSE
te.gif
ting.gif
DME.md
uirements.txt
s.gif

1

04-pretty_printing.py

06-misc.py

def fun(x):

02-autoreload.ipynb

03-databases.py

README.md

05-data_analysis.py

01-plotting.py

```
2. python  
In [8]: %pyplot  
Using matplotlib backend: TkAgg  
In [9]: █
```

~2.1 Million Notebooks on GitHub

Repositories

Developers

Trending: this month ▾

All languages

Unknown languages

C++

JavaScript

Jupyter Notebook

Matlab

Python

Other: Languages ▾

ProTip! Looking for most forked Jupyter Notebook repositories? [Try this search](#)

<https://github.com/trending/jupyter-notebook?since=monthly>

fivethirtyeight / data

★ Star

Data and code behind the articles and graphics at FiveThirtyEight

● Jupyter Notebook ★ 9,553 🍷 3,929 Built by 

★ 829 stars this month

fastai / fastai

★ Star

The fast.ai deep learning library, lessons, and tutorials

● Jupyter Notebook ★ 4,236 🍷 1,307 Built by 

★ 474 stars this month

tensorflow / probability

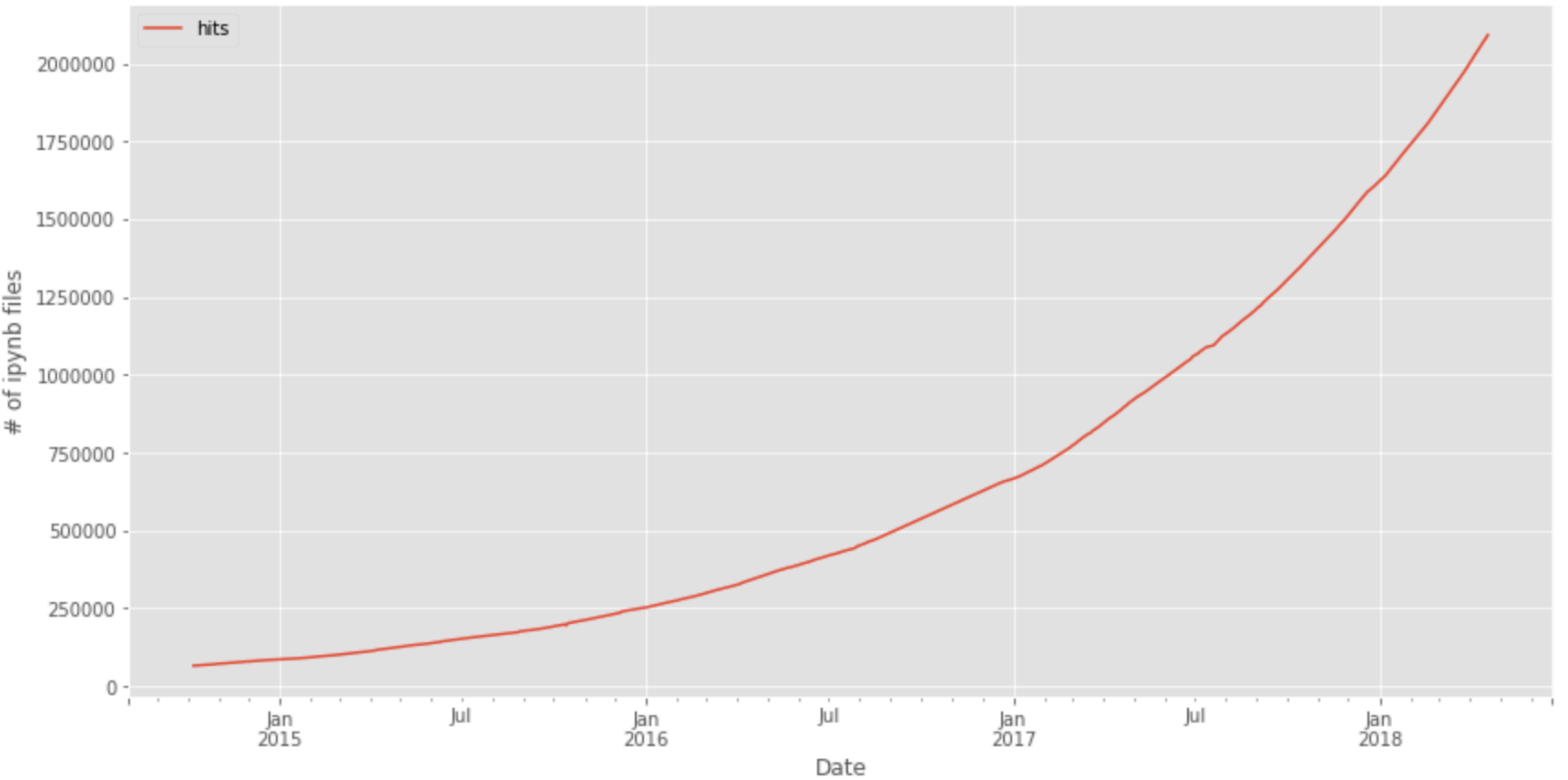
★ Star

Probabilistic reasoning and statistical analysis in TensorFlow

● Jupyter Notebook ★ 500 🍷 62 Built by 

★ 326 stars this month

GitHub search hits for 1291 days sans outliers





Project Jupyter exists to develop open-source software, open-standards, and services for interactive computing across dozens of programming languages.

Project Jupyter Mission



Project Jupyter

Project Jupyter exists to develop open-source software, open standards, and services for interactive and reproducible computing.

Jul 7, 2015 · 3 min read

GORDON AND BETTY
MOORE
FOUNDATION



THE LEONA M. AND HARRY B.
HELMSLEY
CHARITABLE TRUST

New funding for Jupyter

We are pleased to announce that the Jupyter/IPython project has received \$6M in funding from three organisations:

- The [Leona M. and Harry B. Helmsley Charitable Trust](#)
- The [Gordon and Betty Moore Foundation](#)
- The [Alfred P. Sloan Foundation](#)

The grant, which is being made both to the [University of California, Berkeley](#) and [California Polytechnic State University, San Luis Obispo](#), will support the project for three years and includes new collaborations with the [University of Southampton](#), and the [Simula Research Lab](#) in Norway.

Mission and Background

Project Jupyter's mission is to create open source tools for interactive scientific computing and data science in research, education and industry, with an emphasis on usability, collaboration and reproducibility.

This project is structured in a “3+1” format, with three main focus areas of research and development and one extra topic of ongoing work. The three focus areas are *Interactive Computing*, *Computational Narratives* and *Collaboration*. The problem of *Sustainability* will require ongoing attention but is conceptually distinct from the first three, as it doesn't focus on specific research questions or deliverables.



Newsjournal of the Society for Industrial and Applied Mathematics

[sinews.siam.org](#)

Volume 51/ Issue 2
March 2018

MARCH 2018 SIAM NEWS • 7

Jupyter: Tools for the Life Cycle of a Computational Idea

By *Min Ragan-Kelley, Carol Willing, and Jason Grout*

Computation is increasingly becoming an integral part of science and education across disciplines. The life cycle of a computational idea typically involves interactive exploration and experimentation, as well as publication and communication of results. Reproducible computation demands open research tools, good software practices, and transparent documentation of research processes and results. Project Jupyter¹ is an open community that builds open-source software tools and protocols for the life cycle of a computational idea. Two core pieces of the project are an open protocol for interactive computation and an open document format with which to record and share computational ideas. The Jupyter Notebook application builds on these to provide a powerful, interactive, computational environment.

The Jupyter Notebook

What Is a Notebook? The Jupyter

The notebook document format is free, transparent, and understandable, in keeping with its aim to facilitate open and accessible science. It is stored as a single JSON-formatted text file, making it easy to manipulate and understand using standard programming tools, without the need for Jupyter software. The notebook file format is public,² and Jupyter software is open-source under the BSD license.

Many authors communicate using Jupyter Notebook. GitHub hosts 1.4 million notebooks, and some people have written entire books as collections of notebooks, such as Jake Vanderplas's *Python Data Science Handbook*.³ Because notebook documents preserve their content structure and metadata, they are easily convertible to other formats, including plain scripts in the document's language of choice. This also makes them easy to integrate into publication pipelines via formats such as LaTeX, Markdown, and reStructuredText via Jupyter's conversion tool, `nbconvert`.⁴

Using Notebook Documents. The Jupyter Notebook server is a web-based application for interacting with notebook

cutting code cells in the same persistent kernel session, thus building and exploring a computational idea. The user can also create and interleave markdown cells to explain an idea using formatted text and mathematics (see Figure 1).

Users typically install the notebook application on their local computer, where it functions like a desktop application and works with local notebooks and data. Since it is a web application, the notebook server can also run on a remote computer and provide the notebook application via a browser, with no installation required on the user's machine.

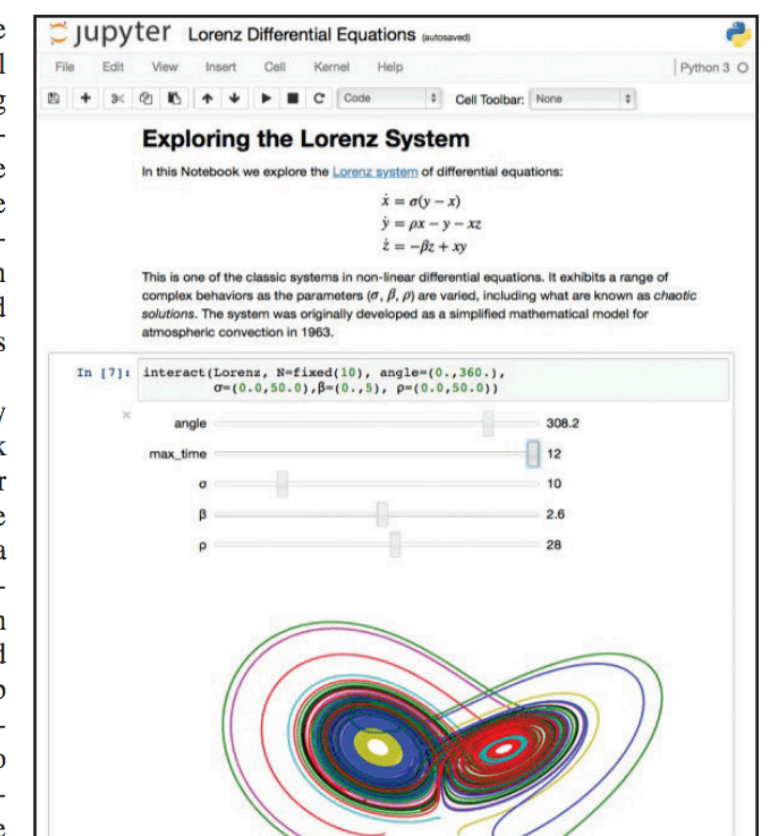


Figure 1. Jupyter notebooks exploring the Lorenz system.

output representations between front ends

We're not-for-profit

📖 README.md

Project Jupyter Governance

The purpose of this repository is to formalize the governance process that the IPython/Jupyter project has used informally since its inception in 2001. This document clarifies how decisions are made and how the various elements of our community interact, including the relationship between open source collaborative development and work that may be funded by for-profit or non-profit entities.

Table of Contents

- [Main Governance Document](#)
- [Current Steering Council and Institutional Partners](#)
- [New Subproject Incubation Process](#)
- [Process for Authoring Jupyter Related Academic Papers](#)

License of Governance Documents

To the extent possible under law, Project Jupyter has waived all copyright and related or neighboring rights to the Project Jupyter Governance documents, in accordance with the Creative Commons [CC0 license](#). This work is published from the United States. See the LICENSE.md file in this repository for details.

<https://github.com/jupyter/governance>



Jupyter Notebook



Narrative Text

Code and Visualizations

Notebook title and introduction

Sampling from the generative model

In this notebook, we will use the generative model of the HDHP (Hierarchical Dirichlet-Hawkes Process) in order to sample events. We will start with a predefined number of users, say 10, and we will attempt to model their behavior as they are posting questions in an online platform. For simplicity, our "vocabulary" will be dummy.

We start by importing all the libraries that will be required.

```
In [1]: %matplotlib inline
import datetime
import string
import hdhp
import notebook_helpers
import seaborn as sns
```

Importing external packages

Description of model parameters

Now, let us set some parameters for our model. These fall under two categories; the ones relevant to the content and then ones relevant to the time dynamics. Starting with the first set, we need to decide on:

- the vocabulary: a dummy set of 100 words, i.e. word0, word1, ... , word99.
- the minimum and maximum length of a question
- the number of words of each pattern

As far as the time dynamics is concerned, we need to set:

- α_0 : the parameters of the Gamma prior for the time kernel of each pattern
- μ_0 : the parameters of the Gamma prior for the user activity rate
- ω : the time decay parameter

Finally, in order to make the generative process more user-friendly, we can pre-set the number of patterns that our users can sample from.

```
In [2]: vocabulary = ['word' + str(i) for i in range(100)] # the "words" of our documents
doc_min_length = 5
doc_length = 10
words_per_pattern = 50

alpha_0 = (2.5, 0.75)
mu_0 = (2, 0.5)
omega = 3.5

num_patterns = 10

process = hdhp.HDHPProcess(num_patterns=num_patterns, alpha_0=alpha_0,
                           mu_0=mu_0, vocabulary=vocabulary,
                           omega=omega, words_per_pattern=words_per_pattern,
                           random_state=12)
```

Implementation of parameters

Description of need to profile data

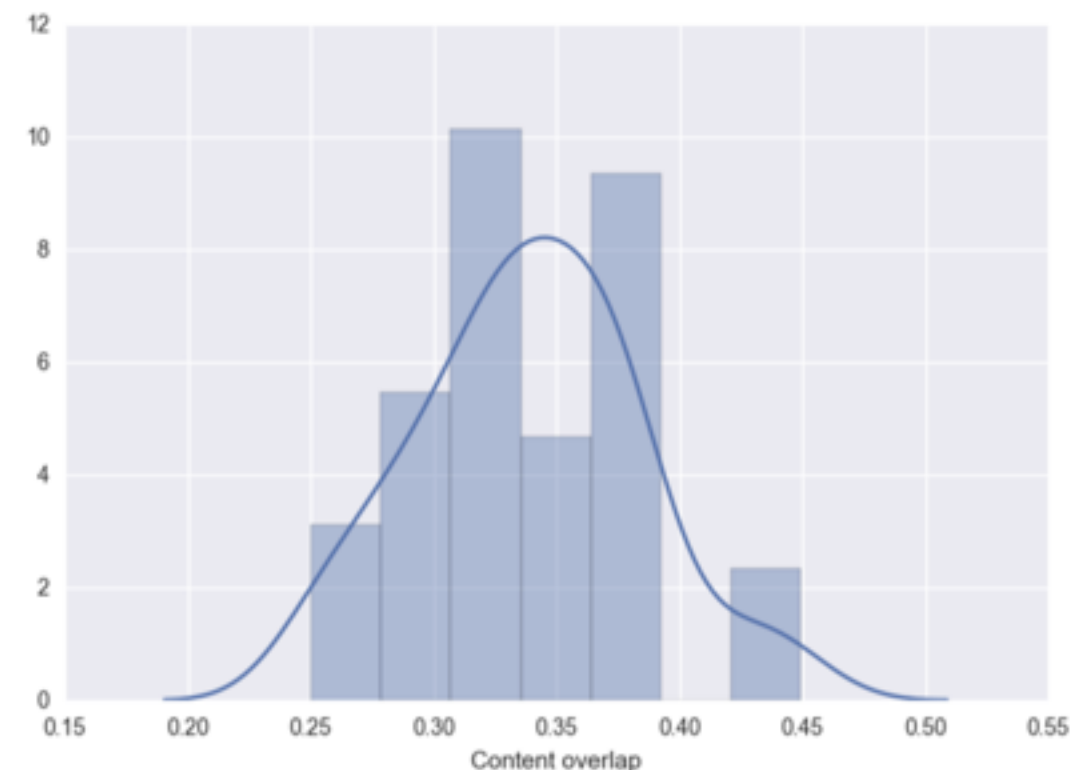
Before generating any questions, we can take a look at the patterns that we initialized our process with, and look at the content distribution of each pattern. Although each pattern has a different word distribution, we can still plot the overlap (Jaccard similarity) between the words that have non-zero probability for each pattern. Since we used a limited number of patterns, the distribution of the overlap will not be smooth.

```
In [3]: overlap = notebook_helpers.compute_pattern_overlap(process)
sns.distplot(overlap, kde=True, norm_hist=True, axlabel='Content overlap')
```

Average overlap: 0.338826769742

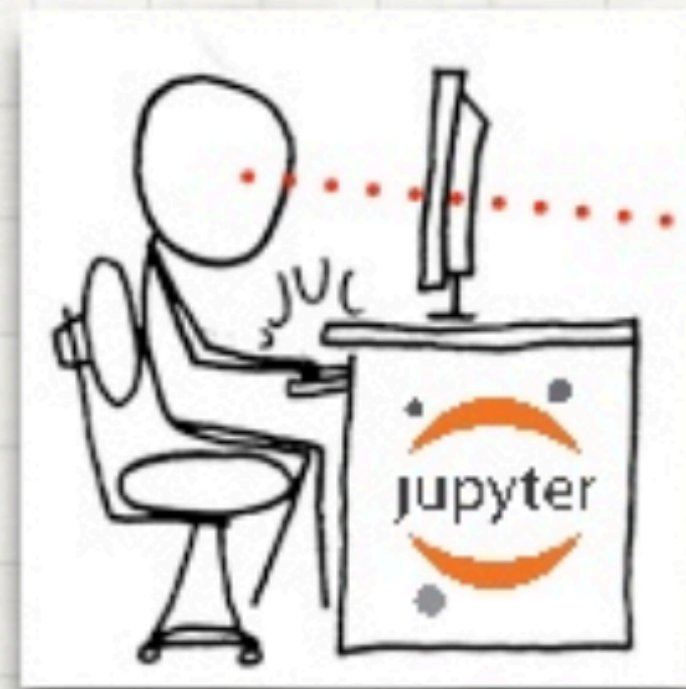
Out[3]: <matplotlib.axes._subplots.AxesSubplot at 0x10de8ca50>

Profile plotting code



Inline plot

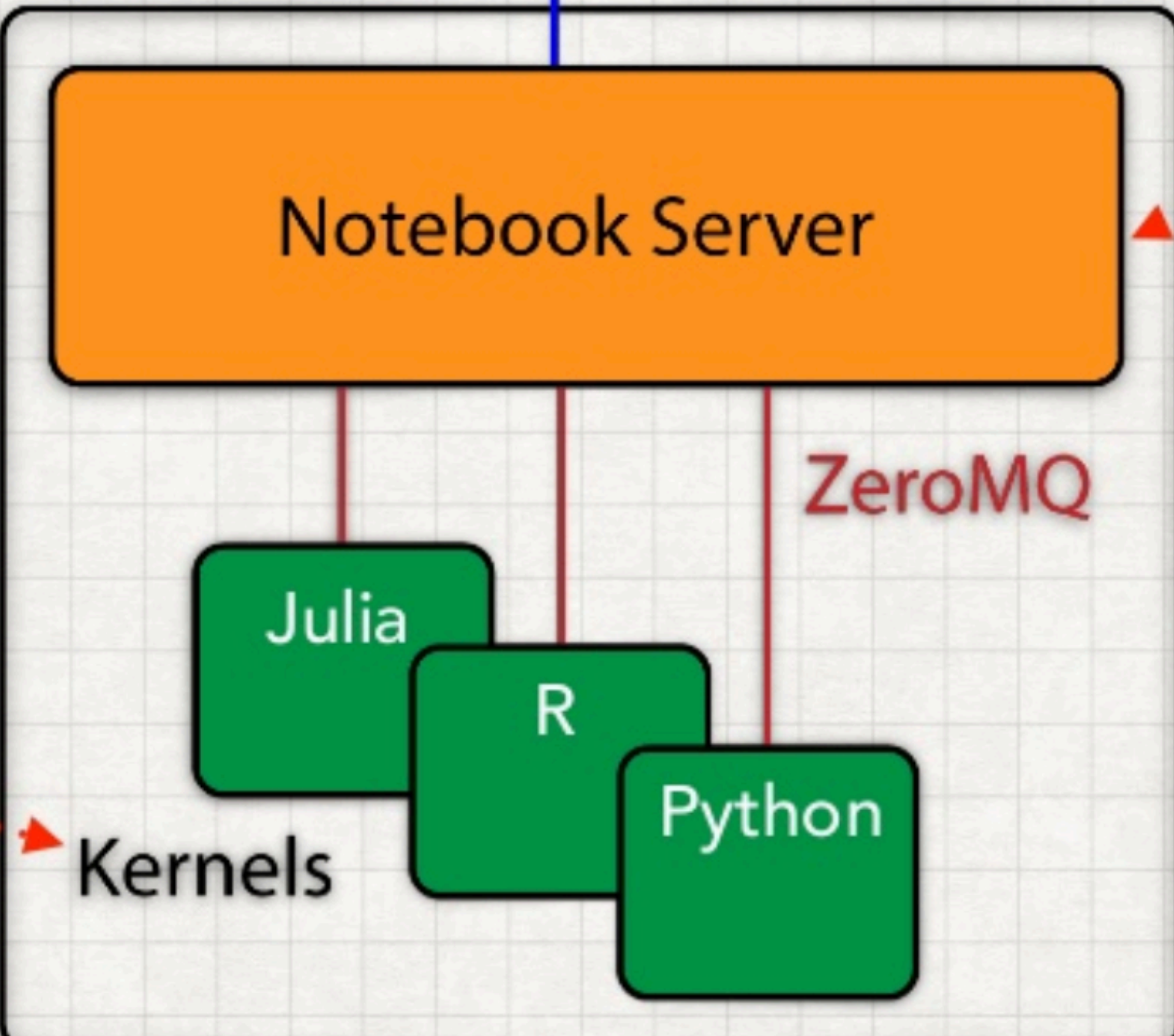
JUPYTER NOTEBOOK



Browser

View and Enter
Stuff Part

HTTP(S) | Websocket



Language Expert Thing

More than Python



Interleave Python with C++: the %%cpp magic ¶

Thanks to ROOT, it is possible to write cells in C++ within a Python notebook. This can be done using the %%cpp magic. Magics are a feature of Jupyter notebooks and when importing the ROOT module, the %%cpp magic was registered.

```
In [10]: %%cpp
cout << "This is a C++ cell" << endl;
```

This is a C++ cell

Not bad. On the other hand, ROOT offers much more than this. Thanks to its [interpreter](#) and [type system](#), entities such as functions, classes and variables, created in a C++ cell, can be accessed from within Python (and viceversa, partially).

```
In [11]: %%cpp
class A{
public:
A(){cout << "Constructor of A!" << endl;}
};
```

```
In [12]: a = ROOT.A()
```

Constructor of A!

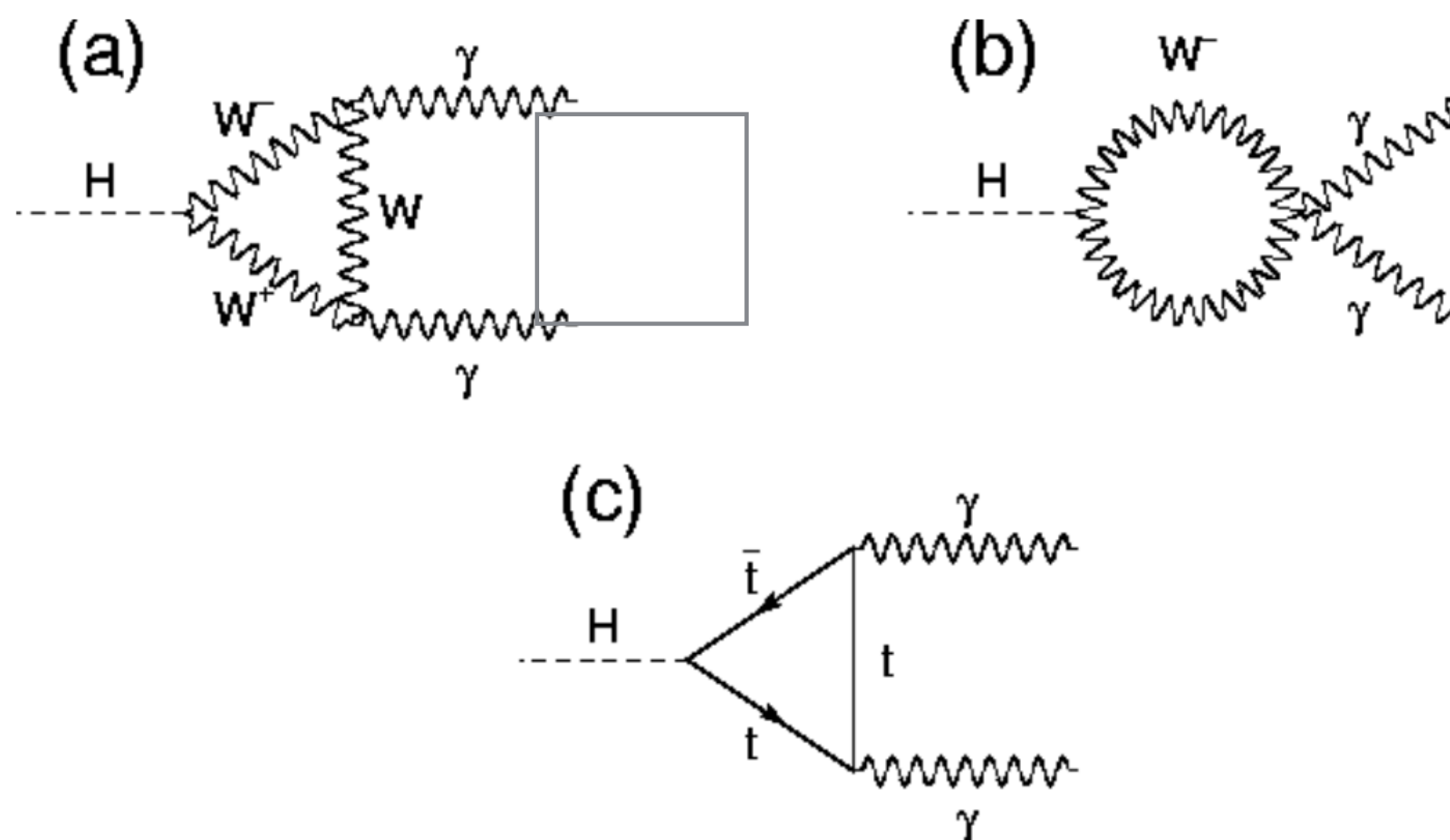
The Python and C++ worlds are so entangled that we can find back in C++ the entities created in Python. To illustrate this, from within a C++ cell, we are going to fit a function in the gauss histogram displayed above and then re-draw the canvas.

```
In [13]: %%cpp
gauss->Fit("gaus", "S");
myCanvasName->Draw();
```



Higgs decay to two photons

The Standard Model predicted the decay of the [Higgs bosons](#) into photons. The process is depicted by the diagrams below:



At the [Large Hadron Collider](#), this process has been measured. This figure shows how an Higgs boson decay looks in the CMS detector:



HiggsFit (autosaved)

Terminal

Logout

File Edit View Insert Cell Kernel Help

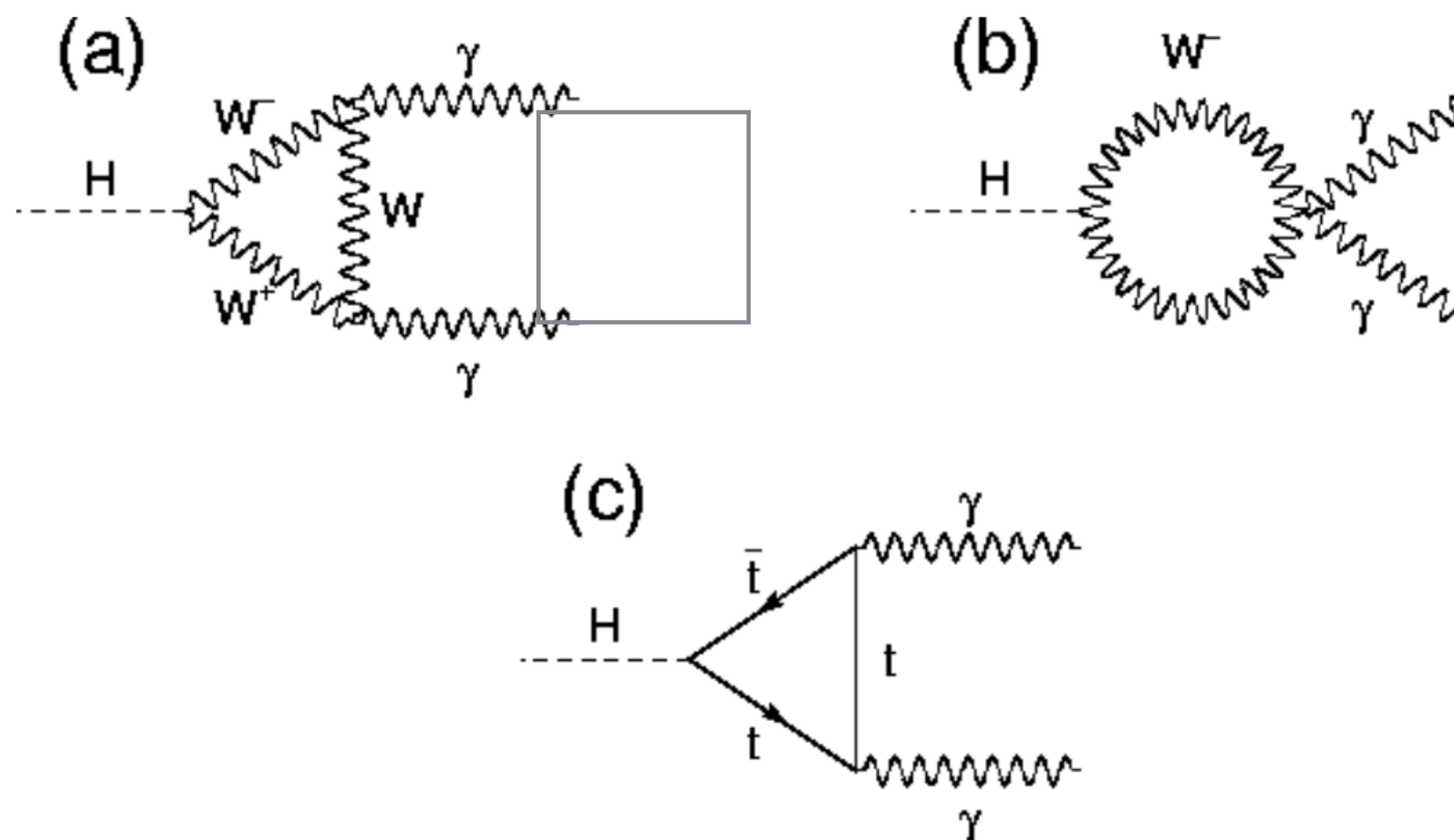
Not Trusted

ROOT C++ ●



Higgs decay to two photons

The Standard Model predicted the decay of the [Higgs bosons](#) into photons. The process is depicted by the diagrams below:



At the [Large Hadron Collider](#), this process has been measured. This figure shows how an Higgs boson decay looks in the CMS detector:



```
In [ ]: import (  
        "fmt"  
        "time"  
    )
```

```
In [ ]: // sum calculates the sum of two integers  
func sum(x, y int) int {  
    return x + y  
}
```

```
In [ ]: a, b := 3, 4
```

```
In [ ]: fmt.sp("sum(%d, %d) = %d", a, b, sum(a, b))
```

```
In [ ]: start := time.Now()  
defer func() {  
    end := time.Now()  
    fmt.Println("Interrupted an infinite loop after",  
                end.Sub(start))  
}()  
for {}
```

Jupyter kernels

Kernel Zero is [IPython](#), which you can get through [ipykernel](#), and is still a dependency of [jupyter](#). The IPython kernel can be thought of as a reference implementation, as CPython is for Python.

Here is a list of available kernels. If you are writing your own kernel, feel free to add it to the table!

Name	Jupyter/IPython Version	Language(s) Version	3rd party dependencies	Example Notebooks
Coarray-Fortran	Jupyter 4.0	Fortran 2008/2015	GFortran >= 7.1, OpenCoarrays , MPICH >= 3.2	Demo Binde demo
sparkmagic	Jupyter >=4.0	Pyspark (Python 2 & 3), Spark (Scala), SparkR (R)	Livy	Noteb Docke Image

<https://github.com/jupyter/jupyter/wiki/Jupyter-kernels>

Creating new Jupyter kernels

[Making kernels for Jupyter in the documentation.](#)

[Simple example kernel](#)

[IHaskell creator blog post](#)

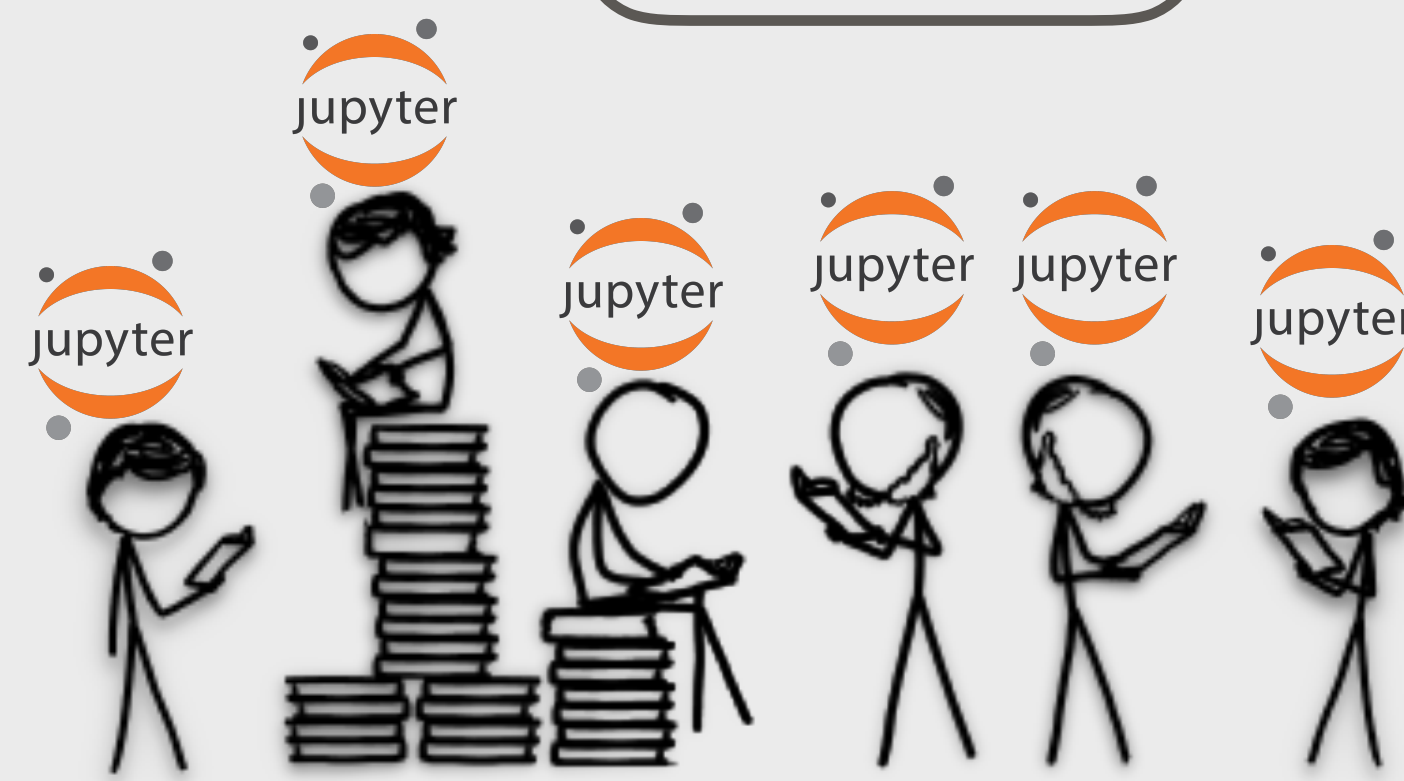
[Testing kernels against message specification \(work in progress\)](#)

[Tool to test a kernel against specification \(work in progress\)](#)

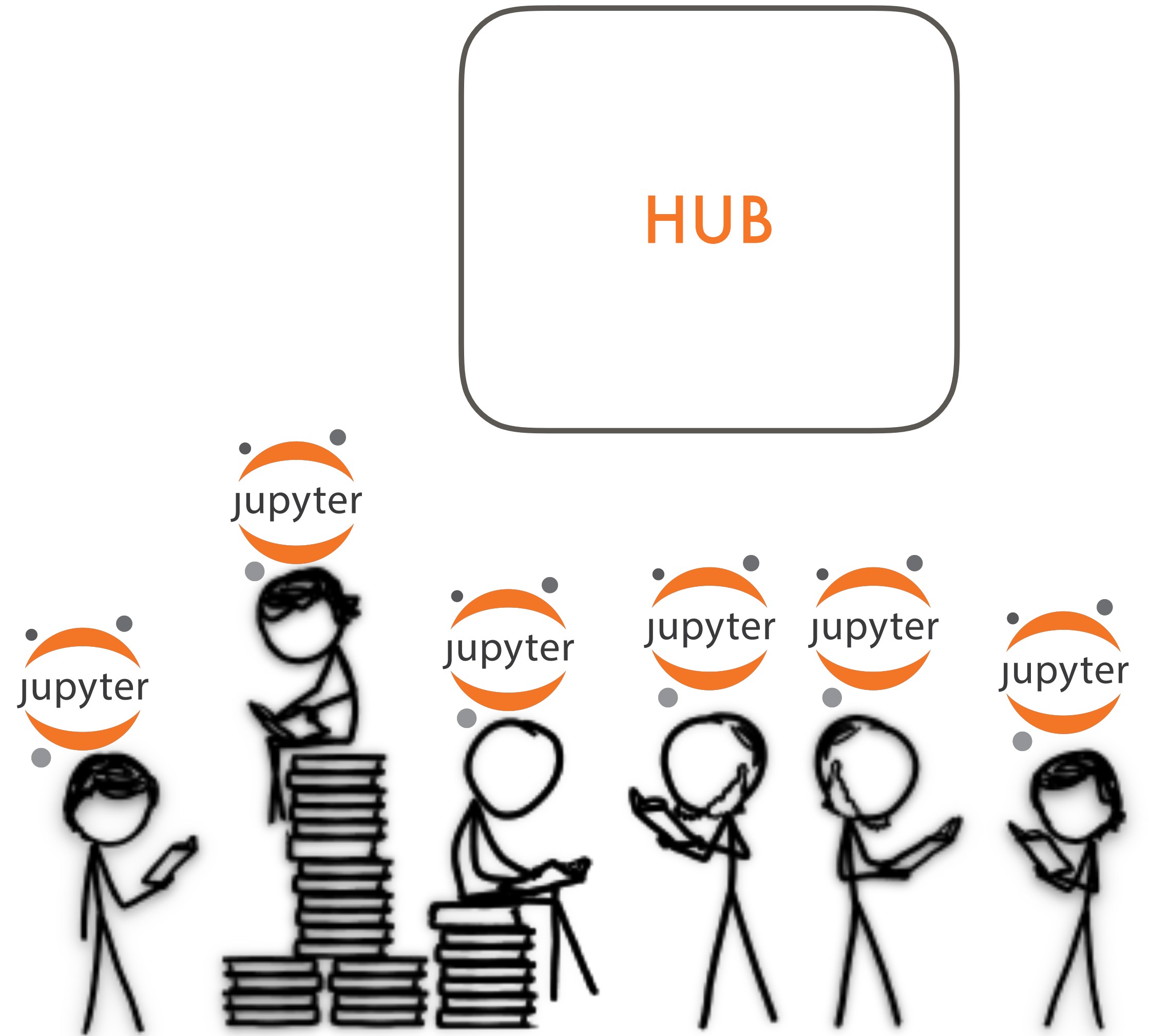
More than Notebooks



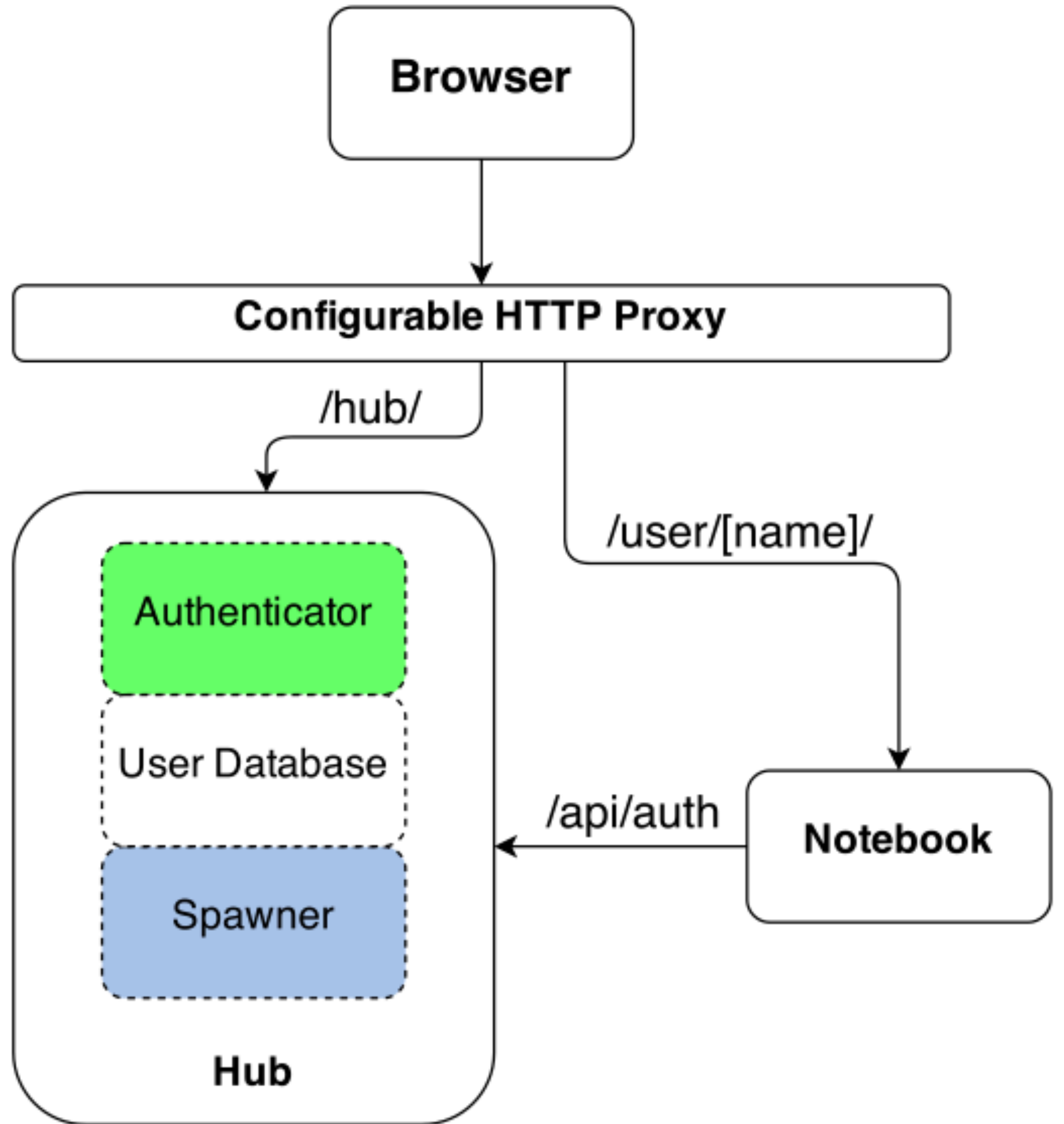
jupyterhub



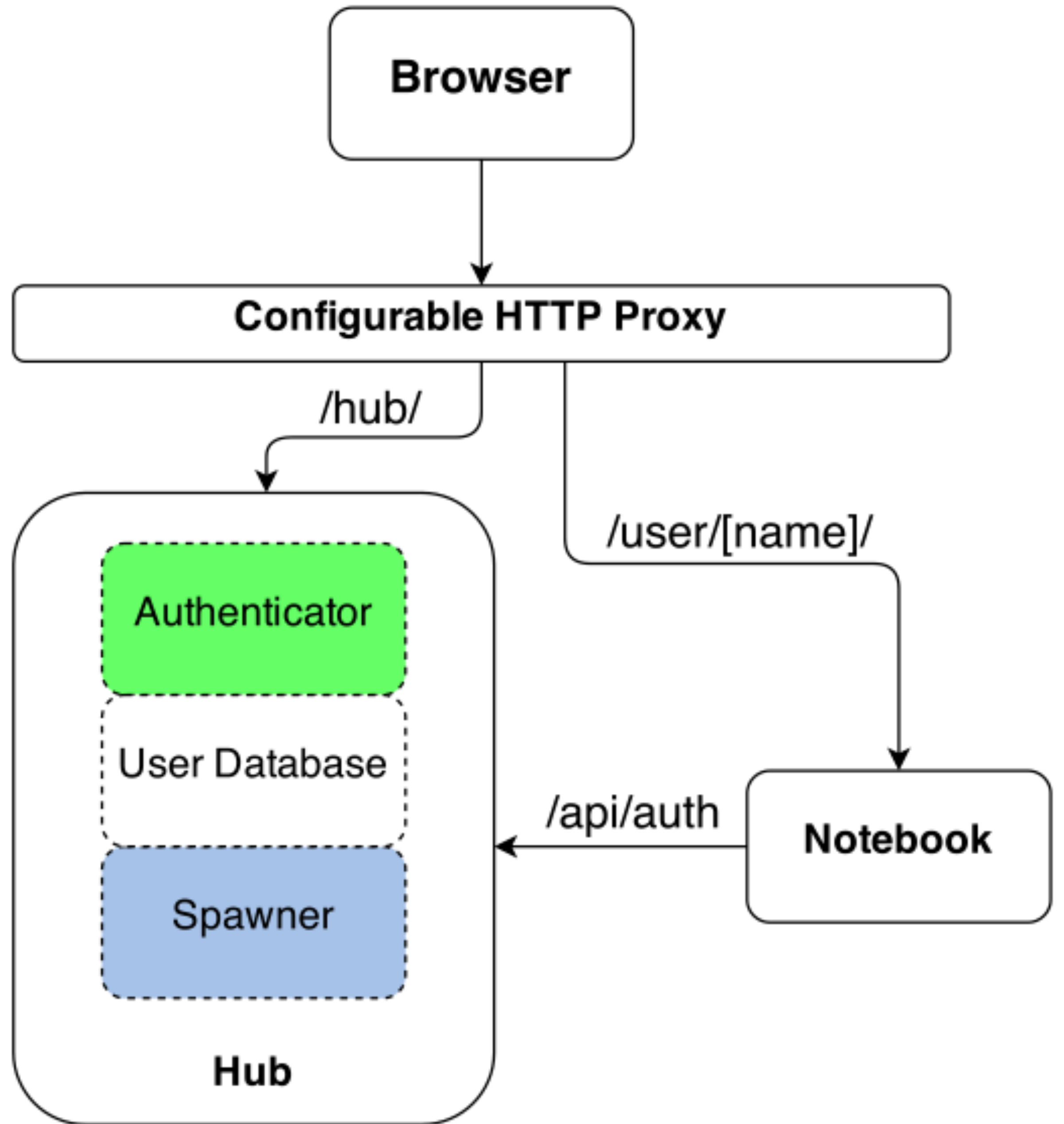
- JupyterHub provides a single user Jupyter Notebook server for each person in a group
- Similar to SWAN



- **Hub**: manages user accounts, authentication, and coordinates Single User Notebook Servers using a Spawner
- **Proxy**: routes HTTP requests to the Hub and Single User Notebook Servers.
- **Spawner**: starts a **single-user notebook servers** when a user logs in



- Hub launches a proxy
- Proxy forwards all requests to Hub by default
- Hub handles login, and spawns single-user servers on demand
- Hub configures proxy to forward url prefixes to the single-user notebook servers



Kubernetes as Spawner

[Next »](#)



Zero to JupyterHub with Kubernetes

A tutorial to help install and manage JupyterHub with Kubernetes

Quick search

Zero to JupyterHub

JupyterHub is a tool that allows you to quickly utilize cloud computing infrastructure to manage a hub that enables users to interact remotely with a computing environment that you specify.

JupyterHub offers a useful way to standardize the computing environment of a group of people (e.g., for a class of students or an analytics team), as well as allowing people to access the hub remotely.

This growing collection of information will help you set up your own JupyterHub instance. It is in an early stage, so the information and tools may change quickly. If you see anything that is incorrect or have any questions, feel free to reach out at the issues page.

Creating your JupyterHub

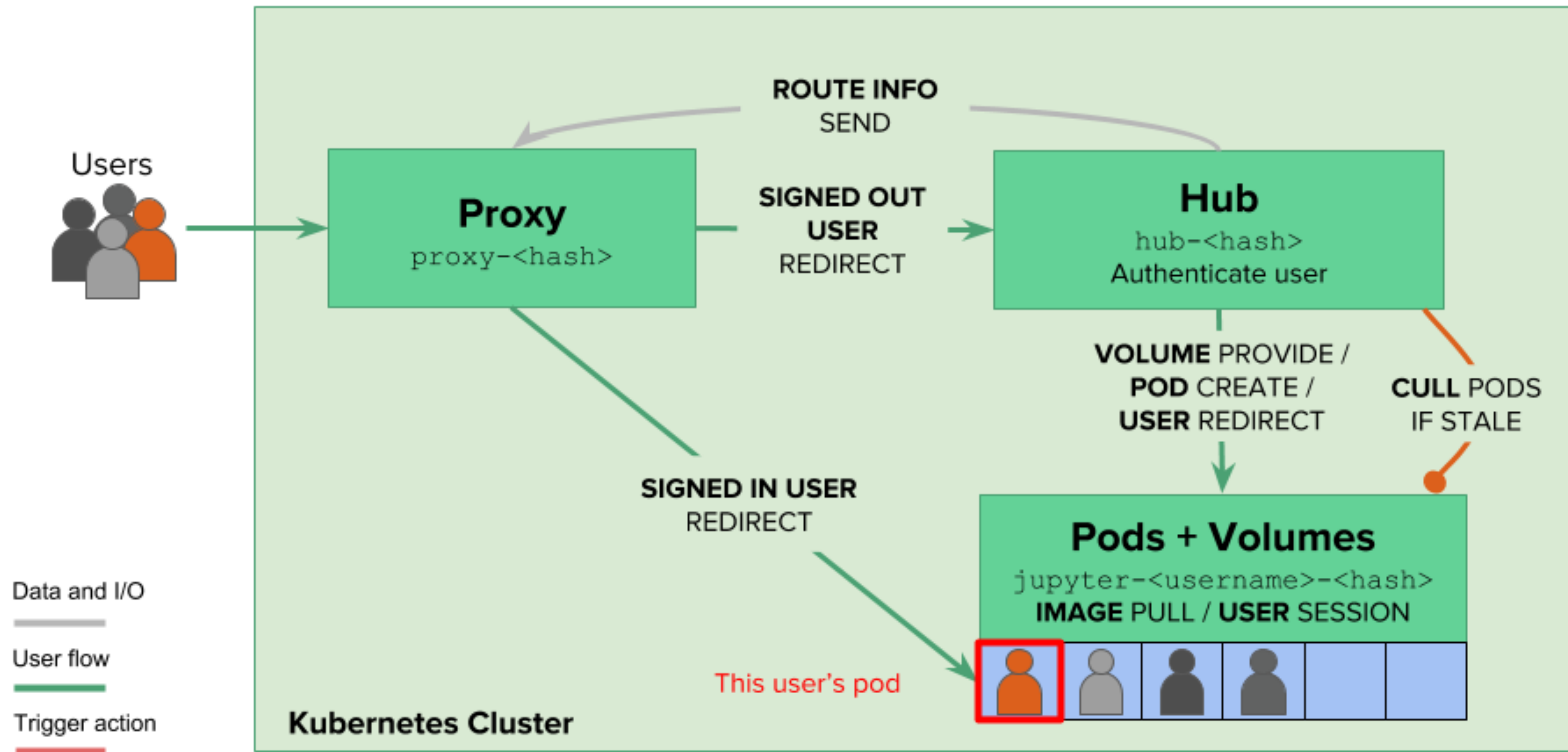
This tutorial starts from “step zero” and walks through how to install and configure a complete JupyterHub deployment in the cloud. Using Kubernetes and the JupyterHub Helm chart provides sensible defaults for an initial deployment.

Kubernetes as Spawner

JupyterHub Architecture
(high-level details)

Cloud Volumes
Provides persistent storage

Image Registry
Provides environment images



Scaling JupyterHub



KubeCon + CloudNativeCon North America 2017 has ended

Thursday, December 7 • 3:50pm - 4:25pm

Large Scale Teaching Infrastructure with Kubernetes - Yuvi Panda, Berkeley University

[Sign up](#) or [log in](#) to save this to your schedule and see who's attending!

<http://sched.co/CU7L>



Tweet



Share

Data Science & Programming literacy is an important aspect of literacy in the 21st century, but teaching these skills at scale is quite difficult. At UC Berkeley, we are trying - our 'Foundations of Data Science' course has no pre-requisites, and routinely attracts more than a 1000 students from across majors.



Scaling JupyterHub

PAWS

Phabricator project: [#paws](#)

PAWS: A Web Shell (PAWS) is a [Jupyter notebooks](#) deployment that has been customized to make interacting with Wikimedia wikis easier. It allows users to create and share documents that contain live code, visualizations such as graphs, rich text, etc. The user created notebooks are a powerful tool that enables data analysis and scientific research, and also transforms the way in which programmers write code - by enabling an exploratory environment with a quick feedback loop, and a low barrier for entry through it's easy to use graphical interface.

Sign in with your wiki account and tada!

Contents [\[hide\]](#)

- [Usage](#)
- [Documentation](#)
- [Other notes](#)
- [See also](#)



Scaling JupyterHub

Pangeo: JupyterHub, Dask, and XArray on the Cloud

This work is supported by [Anaconda Inc](#), the NSF EarthCube program, and UC Berkeley BIDS

A few weeks ago a few of us stood up pangeo.pydata.org, an experimental deployment of JupyterHub, Dask, and XArray on Google Container Engine (GKE) to support atmospheric and oceanographic data analysis on large datasets. This follows on [recent work](#) to deploy Dask and XArray for the same workloads on super computers. This system is a proof of concept that has taught us a great deal about how to move forward. This blogpost briefly describes the problem, the system, then describes the collaboration, and finally discusses a number of challenges that we'll be working on in coming months.

The Problem

Atmospheric and oceanographic sciences collect (with satellites) and generate (with simulations) large datasets that they would like to analyze with distributed systems. Libraries like Dask and XArray already solve this problem computationally if scientists have their own clusters, but we seek to expand access by deploying on cloud-based systems. We build a system to which people can log in, get Jupyter Notebooks, and launch Dask clusters without much hassle. We hope that this increases access, and connects more scientists with more cloud-based datasets.



Kubeflow

The Kubeflow project is dedicated to making **deployments** of machine learning (ML) workflows on [Kubernetes](#) simple, portable and scalable. Our goal is **not** to recreate other services, but to provide a straightforward way to deploy best-of-breed open-source systems **for ML** to diverse infrastructures. Anywhere you are running Kubernetes, you should be able to run Kubeflow.

This repository contains the manifests for creating:

- A [JupyterHub](#) to create and manage interactive [Jupyter notebooks](#). [Project Jupyter](#) is a non-profit, open-source project to support interactive data science and scientific computing across all programming languages.
- A [TensorFlow Training Controller](#) that can be configured to use either CPUs or GPUs and dynamically adjusted to the size of a cluster with a single setting
- A [TensorFlow Serving](#) container to export trained TensorFlow models to Kubernetes

This document details the steps needed to run the Kubeflow project in any environment in which Kubernetes runs.

Reproducibility





Genomic analysis of elongated skulls reveals extensive female-biased immigration to Early Medieval Bavaria

Krishna R. Veeramah^a, Andreas Rott^{b,1}, Melanie Groß^{c,1}, Lucy van Dorp^d, Saioa López^e, Karola Kiršanow^c, Christian Sell^c, Jens Blöcher^c, Daniel Wegmann^{f,g}, Vivian Link^{f,g}, Zuzana Hofmanová^{f,g}, Joris Peters^{b,h}, Bernd Trautmann^b, Anja Gairhosⁱ, Jochen Haberstroh^j, Bernd Pääffgen^k, Garrett Hellenthal^d, Brigitte Haas-Gebhardⁱ, Michaela Harbeck^{b,2,3}, and Joachim Burger^{c,2,3}

^aDepartment of Ecology and Evolution, Stony Brook University, Stony Brook, NY 11794-5245; ^bState Collection for Anthropology and Palaeoanatomy, Bavarian Natural History Collections, 80333 Munich, Germany; ^cPalaeogenetics Group, Institute of Organismic and Molecular Evolution, Johannes Gutenberg University Mainz, 55099 Mainz, Germany; ^dUCL Genetics Institute, Department of Genetics, Evolution and Environment, University College London, WC1E 6BT London, United Kingdom; ^eCancer Institute, University College London, WC1E 6DD London, United Kingdom; ^fDepartment of Biology, University of Fribourg, 1700 Fribourg, Switzerland; ^gSwiss Institute of Bioinformatics, 1700 Fribourg, Switzerland; ^hArchaeoBioCenter and Institute for Palaeoanatomy, Limnology Research and the History of Veterinary Medicine, Ludwig Maximilian University, 80539 Munich, Germany; ⁱBavarian State Archaeological Collection, 80539 Munich, Germany; ^jBavarian State Department of Monuments, 80339 Munich, Germany; ^kInstitute of Prehistoric and Protohistoric Archaeology, Ludwig Maximilian University, 80539 Munich, Germany

Edited by Eske Willerslev, University of Copenhagen, Copenhagen, Denmark, and approved January 30, 2018 (received for review November 21, 2017)

Modern European genetic structure demonstrates strong correlations with geography, while genetic analysis of prehistoric humans has indicated at least two major waves of immigration from outside the continent during periods of cultural change. However, population-level genome data that could shed light on the demographic processes occurring during the intervening periods have been absent. Therefore, we generated genomic data

to form in the 5th century AD, and that it emanated from a combination of the romanized local population of the border province of the former Roman Empire and immigrants from north of the Danube (2). While the Baiuvarii are less well known than some other contemporary groups, an interesting archaeological feature in Bavaria from this period is the presence of skeletons with artificially deformed or elongated skulls (Fig. 1A).

POPULATION
BIOLOGY

PNAS / Richard Goerg / Getty / The Atlantic

The Scientific Paper Is Obsolete

Here's what's next.



Atlantic Monthly, April 5 2018

The more sophisticated science becomes, the harder it is to communicate results. Papers today are longer than ever and full of jargon and symbols. They depend on chains of computer programs that generate data, and clean up data, and plot data, and run statistical models on data. These programs tend to be both so sloppily written and so central to the results that it's contributed to a replication crisis, or put another way, a failure of the paper to perform its most basic task: to report what you've actually discovered, clearly enough that someone else can discover it for themselves.

- James Somers



Reproducibility:
Producing similar results
with the same data



Reproducibility:

A software problem

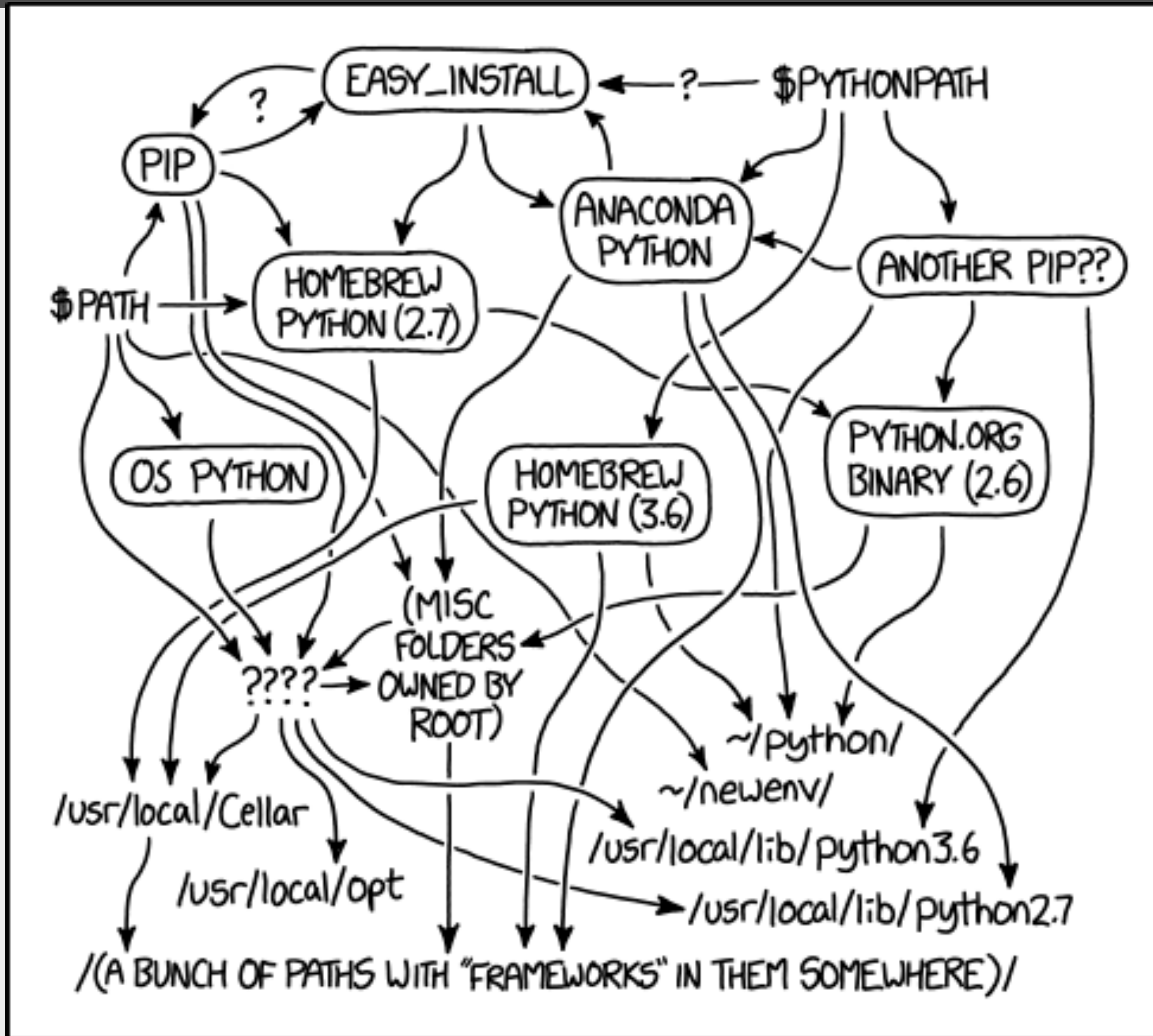


Technical Solutions

- GitHub
- Open Science Framework
- CodaLab
- RunMyCode
- Research Journals



This week in xkcd



MY PYTHON ENVIRONMENT HAS BECOME SO DEGRADED THAT MY LAPTOP HAS BEEN DECLARED A SUPERFUND SITE.

Reproducible scientific software pipelines

- Repository would contain:
 - Data
 - Dependencies
 - Hyperparameters
 - Scripts to run the jobs on similar hardware
 - Analysis code



Reproducible scientific software pipelines

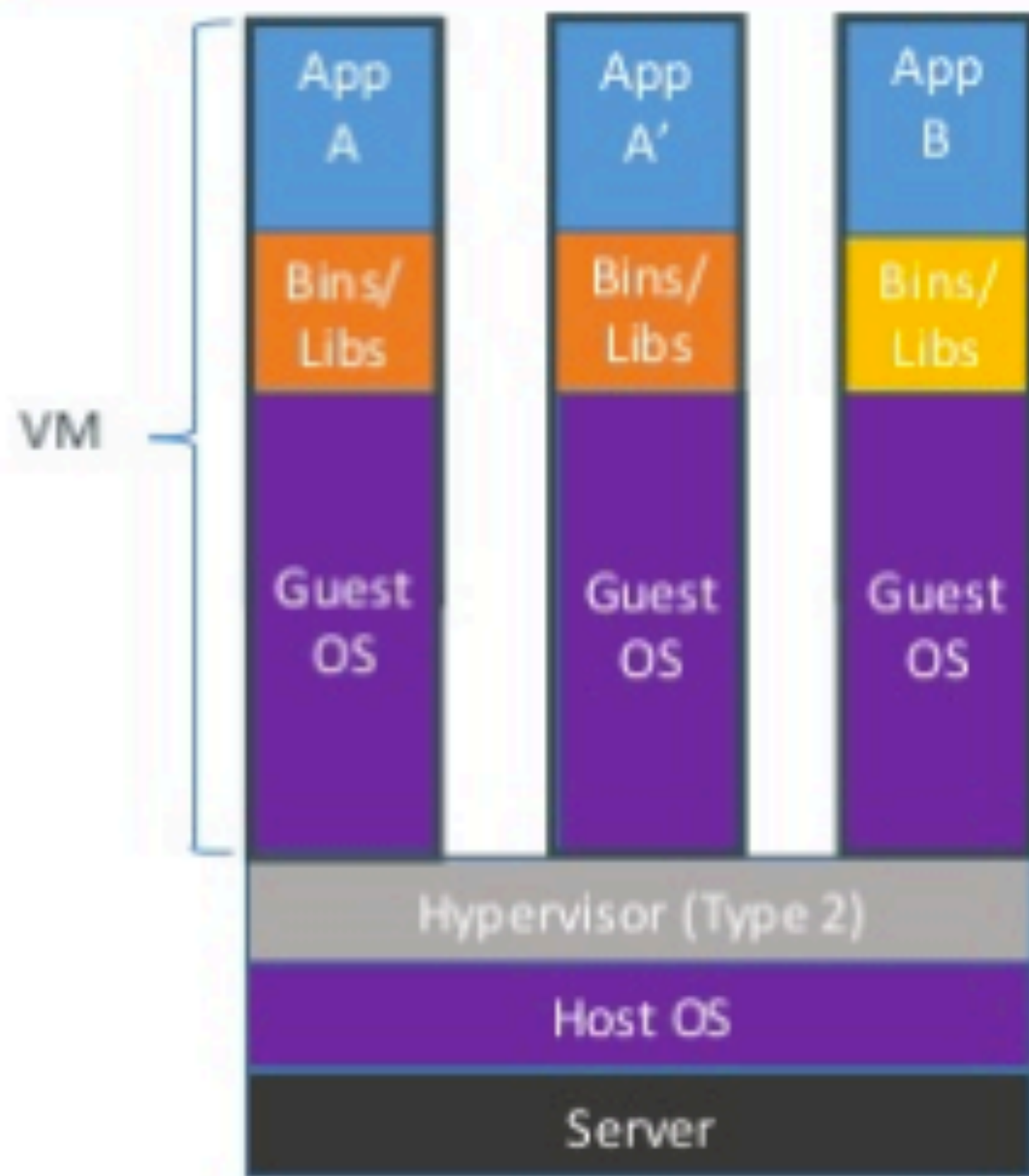
- Repository would contain:
 - Data
 - Dependencies
 - Hyperparameters
 - Scripts to run the jobs on similar hardware
 - Analysis code

**No mention of OS or
lower-level software**



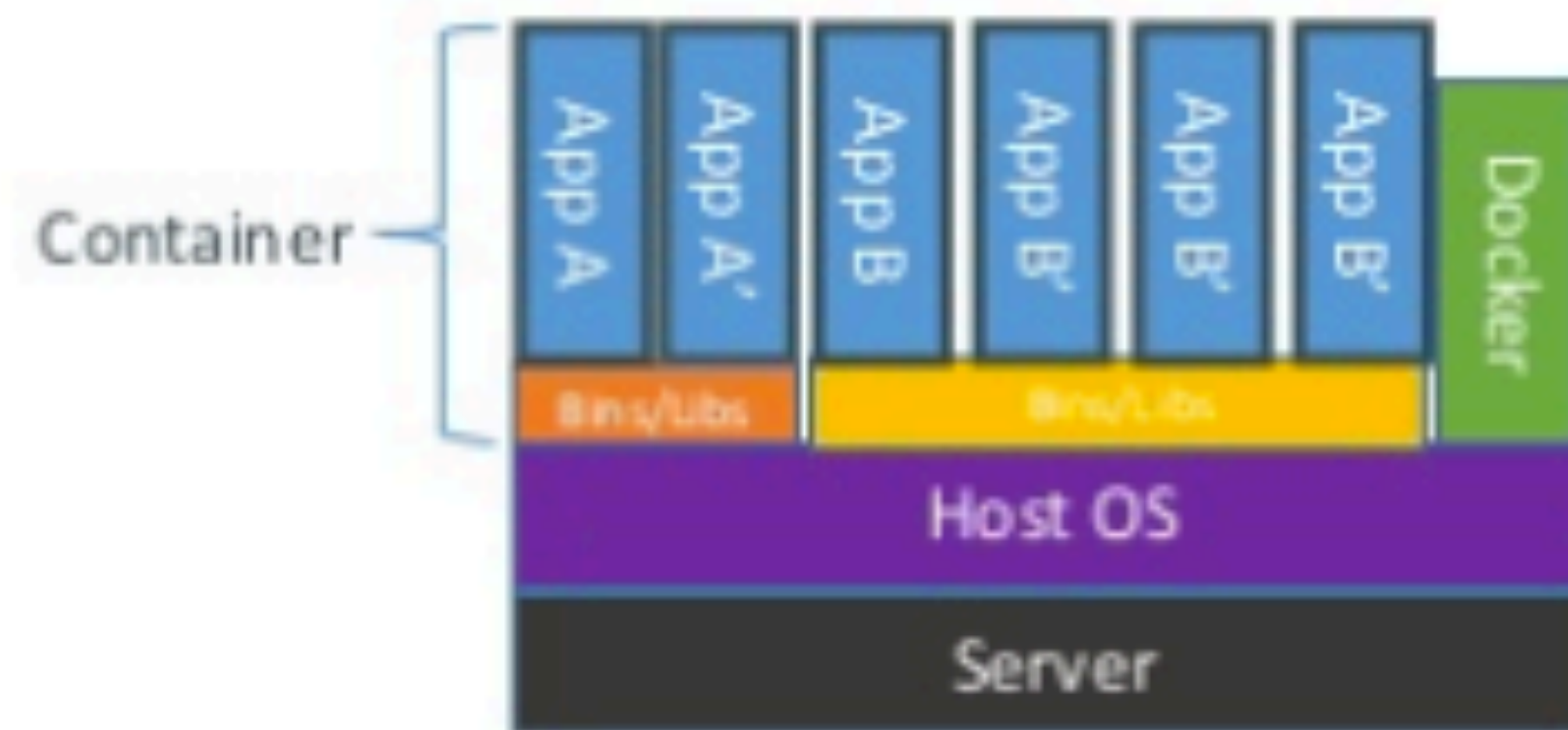
Docker for reproducible software

Containers vs. VMs



Containers are isolated, but share OS and, where appropriate, bins/libraries

...result is significantly faster deployment, much less overhead, easier migration, faster restart



Dockerfiles from GitHub Repos

jupyter-repo2docker

build passing docs passing

jupyter-repo2docker takes as input a repository source, such as a GitHub repo. It then builds, runs, and/or pushes Docker images built from that source.

See the [repo2docker documentation](#) for more information.

Pre-requisites

1. Docker to build & run the repositories. The [community edition](#) is recommended.
2. Python 3.4+.

Supported on Linux and macOS. [See documentation note about Windows support.](#)



Using repo2docker

Note that Docker needs to be running on your machine for this to work.

Example:

```
jupyter-repo2docker https://github.com/norvig/pytudes
```

After building (it might take a while!), it should output in your terminal something like:

```
Copy/paste this URL into your browser when you connect for the first time,  
to login with a token:
```

```
http://0.0.0.0:36511/?token=f94f8fabb92e22f5bfab116c382b4707fc2cade56ad1ace0
```



repo2docker as reproducibility unit-test

- repo2docker looks for configuration files to determine how to build the docker image
- Typically describe dependencies and other instructions to create the environment
- configs in either root or folder called `binder`




Supported config files

- **Dockerfile**: full environment setup
- **environment.yml**: conda
- **requirements.txt**: pip
- **REQUIRE**: Julia
- **apt.txt**: Debian packages
- **postBuild**: custom install script
- **runtime.txt**: Python runtime



Example repo

Branch: master ▾ **GAN_tutorial / environment.yml**


 **mickypaganini** Update environment.yml

1 contributor

9 lines (8 sloc) | 139 Bytes

```
1 name: binder
2 dependencies:
3   - pytorch
4   - torchvision
5   - matplotlib==2.0.2
6   - numpy==1.14.1
7   - ipython==6.2.1
8   - scikit-learn==0.19.1
```

Branch: master ▾ **GAN_tutorial / runtime.txt**

 **mickypaganini** Create runtime.txt

1 contributor

2 lines (1 sloc) | 11 Bytes

```
1 python-2.7
```



Installing LaTeX

Branch: master ▾

latex / apt.txt

Find file

Copy path



choldgraf updating to latex symbol

f152e6d on Dec 12, 2017

1 contributor

8 lines (7 sloc) | 126 Bytes

Raw

Blame

History



```
1 texlive-latex-base
2 texlive-latex-recommended
3 texlive-science
4 texlive-latex-extra
5 texlive-fonts-recommended
6 dvipng
7 ghostscript
```



Installing LaTeX

Branch: master ▾

latex / apt.txt

Find file

Copy path



choldgraf updating to latex symbol

f152e6d on Dec 12, 2017

1 contributor

8 lines (7 sloc) | 126 Bytes

Raw

Blame

History

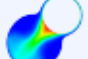


```
1 texlive-latex-base
2 texlive-latex-recommended
3 texlive-science
4 texlive-latex-extra
5 texlive-fonts-recommended
6 dvipng
7 ghostscript
```






Using binder folder and postBuild

Branch: master ▾ [jupyterlab-demo](#) / [binder](#) / [postBuild](#) Find file Copy path

 **ian-r-rose** Remove workspace. b0cd871 on Feb 16

1 contributor

Executable File | 23 lines (19 sloc) | 733 Bytes Raw Blame History   


```
1 set -ex
2
3 # Set up r kernel
4 export R_USER_LIB=$HOME/.R/library
5 mkdir -p $R_USER_LIB
6 echo "R_LIBS_USER='$R_USER_LIB'" > $HOME/.Renviro
7 cat $HOME/.Renviro
8 R -e "install.packages(c('repr', 'IRdisplay', 'evaluate', 'crayon', 'pbdZMQ', 'devtools', 'uuid', 'digest'), lib='$R_USER_LIB',
9 R -e "devtools::install_github('IRkernel/IRkernel', lib='$R_USER_LIB')"
10 R -e "IRkernel::installspec()"
11
12 jupyter nbconvert --to notebook --execute --ExecutePreprocessor.timeout=60 --stdout notebooks/R.ipynb > /dev/null;
13
14 invoke build --env-name=root --no-kernel
15 invoke demofiles
16 invoke talk -t demo
17 rm -rf demofiles
18 rm -rf notebooks
19 rm -rf narrative
20 rm -rf slides
21 rm demo/notebooks/Julia.ipynb
22 jupyter lab clean
```

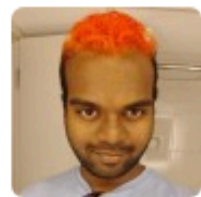
<https://github.com/jupyterlab/jupyterlab-demo>



Future: repo2docker with cloud services

Easy way to run a single repo on your own cloud instance #37

 **Open** yuvipanda opened this issue on Jul 4, 2017 · 1 comment



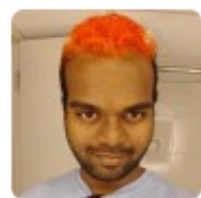
yuvipanda commented on Jul 4, 2017 • edited ▾

Member



We want to provide really easy ways for users to run a github repo on their own AWS or gcloud or OpenStack whatever instance. It should ideally allow them to do persistent long term work on it, including pushing stuff back to GitHub. We can use cloud-init + packer to make this happen across clouds without too much work on our part.

We also have a responsibility to do this as securely as possible, rather than leaving random botnet invitations across the internet.



yuvipanda commented on Jul 4, 2017

Member



Note that this will be only for single-users - we should also have this for binderhub itself (which is multiuser).



willingc added **enhancement** **needs: discussion** labels on Dec 13, 2017



repo2docker as reproducibility unit-test

- making an image from a repository requires:
 1. version of the base docker image
 2. version of repo2docker itself
 3. versions of the libraries installed by the repository



repo2docker as reproducibility unit-test

- making an image from a repository requires:
 1. version of the base docker image
 2. version of repo2docker itself
 3. versions of the libraries installed by the repository
- repo2docker deterministically controls 1 & 2
- user controls 3



repo2docker as reproducibility unit-test

- Because repo2docker is deterministic, we can determine whether researchers have fully described the software environment used to load, analyze, and visualize their data
- Version numbers especially important
- If the configuration files don't create an environment that can run the code provided, the repository is **not reproducible** as currently configured



Reproducibility:
Producing similar results
with the same data



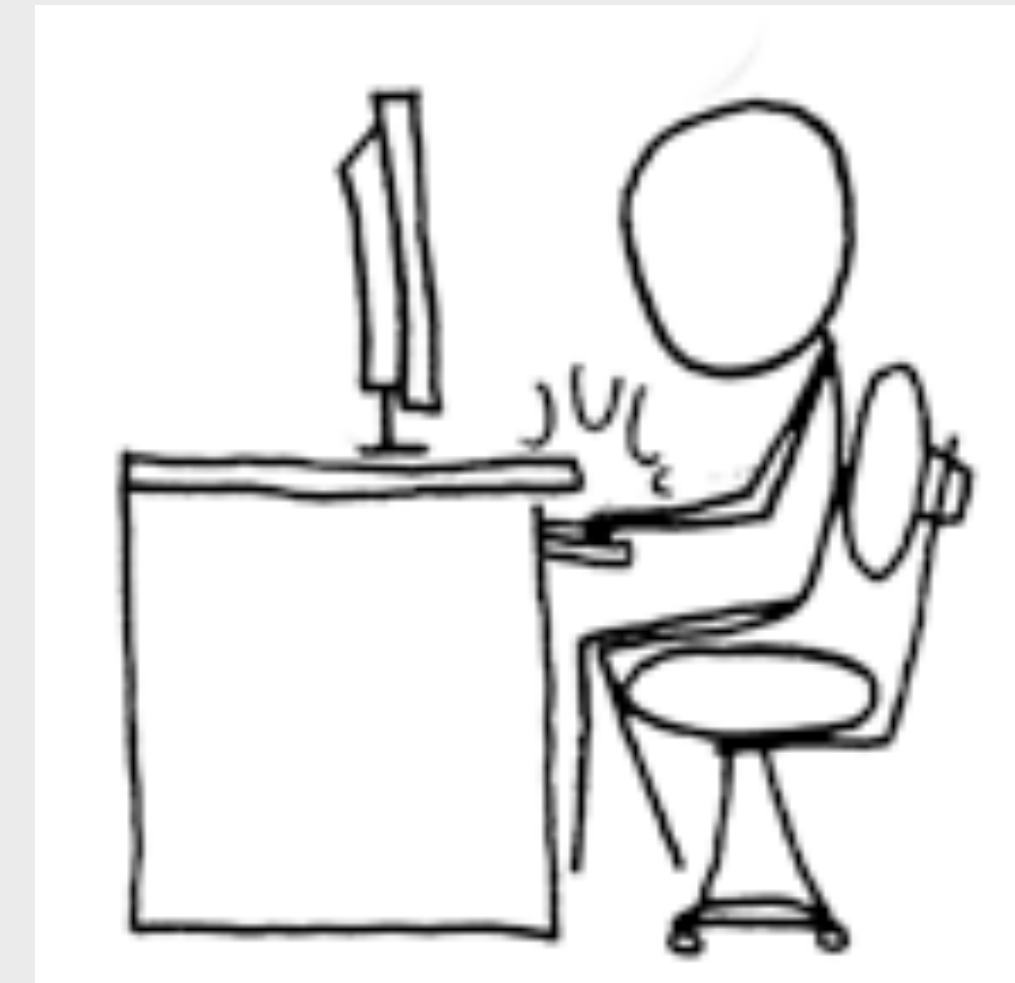
Author



Author



Other scientist



Reproducibility:

Producing similar results
with the same data
independently



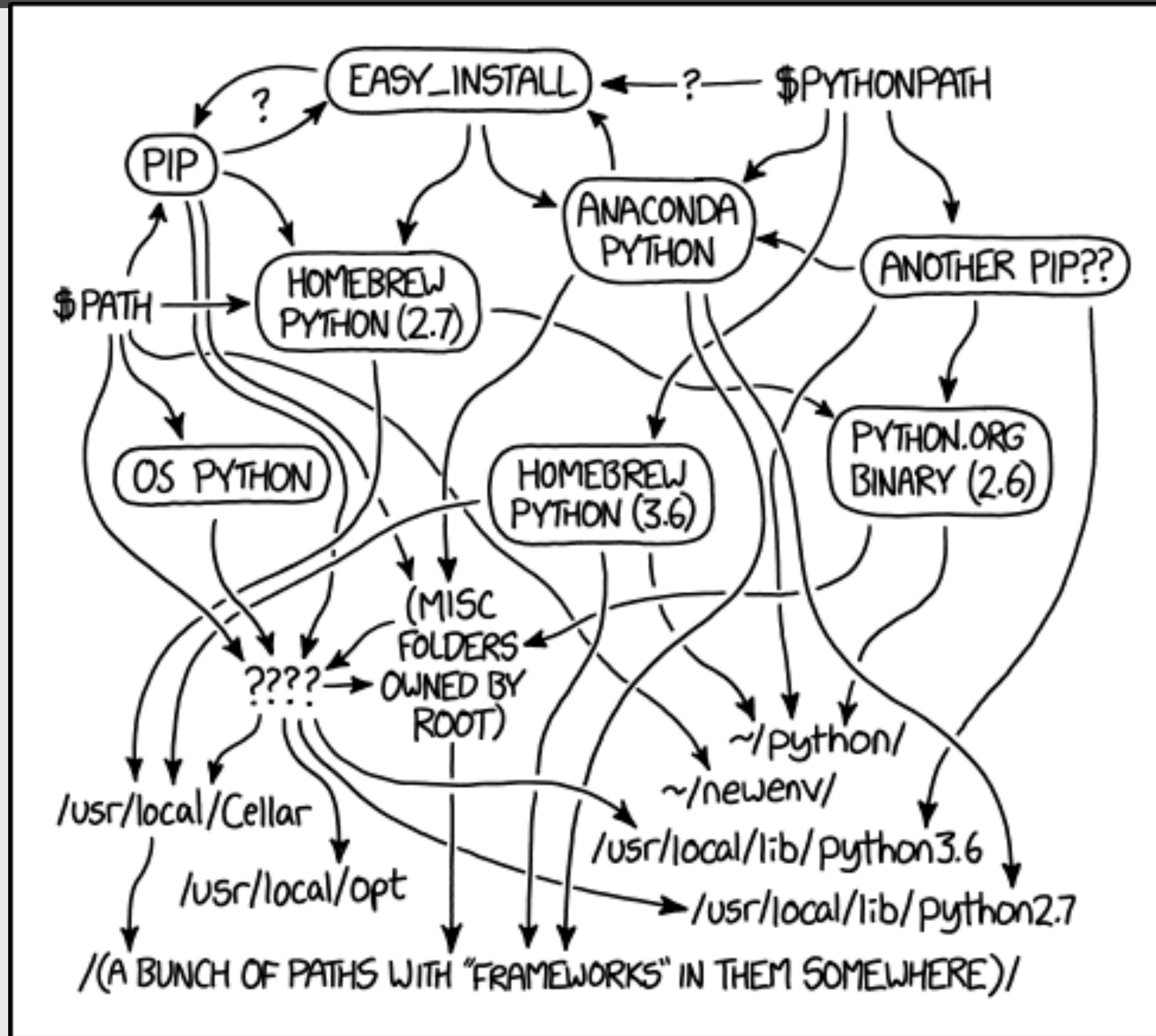
repo2docker:
Produces the same
environment as the scientist



docker opens the
doors to the
laboratory itself



At a minimum, we should specify our environment



MY PYTHON ENVIRONMENT HAS BECOME SO DEGRADED THAT MY LAPTOP HAS BEEN DECLARED A SUPERFUND SITE.



```
conda env export > environment.yml
```



```
pip freeze > requirements.txt
```

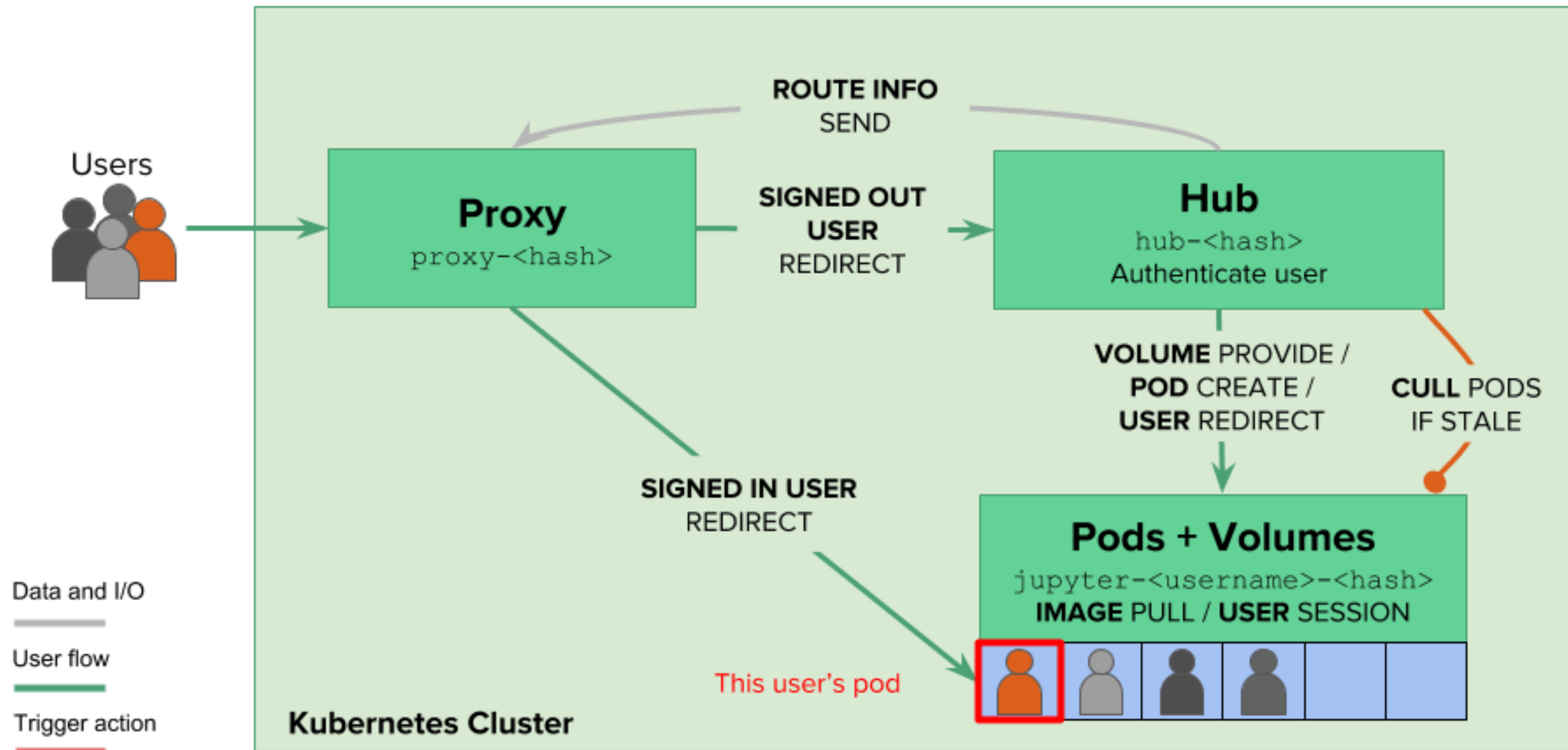


Docker for pods

JupyterHub Architecture
(high-level details)

Cloud Volumes
Provides persistent storage

Image Registry
Provides environment images



JupyterHub, repo2docker as a service



Binder:

Public repositories using
JupyterHub and repo2docker





Toward publishing reproducible computation with Binder

Binder makes it easy to include an interactive version of your analysis, with the supporting data and code, alongside a published paper.



[LABS](#) May 13, 2016



Turn a GitHub repo into a collection of interactive notebooks powered by Jupyter and Kubernetes.

Have a repo full of Jupyter notebooks? With Binder, you can add a badge that opens those notebooks in an executable environment, making your code immediately reproducible by anyone, anywhere.

100% free and open source. Browse examples. Read the FAQ. (Currently in testing, let us know if you run into trouble!)

1

Tell us your GitHub repo

user/project OR github url

This should contain Jupyter notebooks. If one of them is called index.ipynb it will be where your Binder starts. Any extra folders or files (e.g. data) will be included. See an example repo that uses Binder.



Introducing Binder 2.0 – share your interactive research environment

The Project Jupyter team shares its reboot of Binder, the tool that allows researchers to make their GitHub repositories executable by others.



LABS Nov 30, 2017



Turn a GitHub repo into a collection of interactive notebooks

Have a repository full of Jupyter notebooks? With Binder, open those notebooks in an executable environment, making your code immediately reproducible by anyone, anywhere.

Build and launch a repository

GitHub repo or URL

Git branch, tag, or commit

Path to a notebook file (optional)

File ▾

launch

Copy the URL below and share your Binder with others:

Fill in the fields to see a URL for sharing your Binder.



Copy the text below, then paste into your README to show a binder badge: `launch binder`

Build and launch a repository

GitHub repo or URL

Git branch, tag, or commit

Path to a notebook file (optional)

File ▼

launch

Copy the URL below and share your Binder with others:



Copy the text below, then paste into your README to show a binder badge: `launch binder`

Already built!

Launching

[https://mybinder.org/gh/
username/repo/branch](https://mybinder.org/gh/username/repo/branch)



GitHub (GitLab and others
also available)



[https://mybinder.org/gh/
username/repo/branch](https://mybinder.org/gh/username/repo/branch)



[https://mybinder.org/gh/
username/repo/branch?
urlpath=lab](https://mybinder.org/gh/username/repo/branch?urlpath=lab)



Files

code

Name	Last Modified
Jets.ipynb	seconds ago
Toy.ipynb	29 minutes ago
jets.py	29 minutes ago

Running

Commands

Cell Tools

Tabs

Jets.ipynb

Code Python 3

```
In [1]: %matplotlib inline
import matplotlib.pyplot as plt
plt.rcParams["figure.figsize"] = (6,6)

import numpy as np
np.random.seed = 777
```

Prepare data

```
In [2]: import h5py

# Get data from http://www.igb.uci.edu/~pfbaldi/p

f = h5py.File("hepjets/test_no_pile_5000000.h5",
X_no_pile = f["features"].value
y_no_pile = f["targets"].value.ravel()

f = h5py.File("hepjets/test_pile_5000000.h5", "r"
X_pile = f["features"].value
y_pile = f["targets"].value.ravel()
```

```
In [3]: from sklearn.cross_validation import train_test_

X = np.vstack((X_no_pile, X_pile))
y = np.concatenate((y_no_pile, y_pile)).ravel()
z = np.zeros(len(X))
z[len(X_no_pile):] = 1

strates = np.zeros(len(X))
strates[(y==0)&(z==0)] = 0
strates[(y==0)&(z==1)] = 1
strates[(y==1)&(z==0)] = 2
strates[(y==1)&(z==1)] = 3
```

paper.pdf

Search icons

Learning to Pivot with Adversarial Networks

Gilles Louppe
New York University
g.louppe@nyu.edu

Michael Kagan
SLAC National Accelerator Laboratory
makagan@slac.stanford.edu

Kyle Cranmer
New York University
kyle.cranmer@nyu.edu

Abstract

Several techniques for domain adaptation have been proposed to account for differences in the distribution of the data used for training and testing. The majority of this work focuses on a binary domain label. Similar problems occur in a scientific context where there may be a continuous family of plausible data generation processes associated to the presence of systematic uncertainties. Robust inference is possible if it is based on a pivot – a quantity whose distribution does not depend on the unknown values of the nuisance parameters that parametrize this family of data generation processes. In this work, we introduce and derive theoretical results for a training procedure based on adversarial networks for enforcing the pivotal property (or, equivalently, fairness with respect to continuous attributes) on a predictive model. The method includes a hyperparameter to control the trade-off between accuracy and robustness. We demonstrate the effectiveness of this approach with a toy example and examples from particle physics.

1 Introduction

Machine learning techniques have been used to enhance a number of scientific disciplines, and they have the potential to transform even more of the scientific process. One of the challenges of applying machine learning to scientific problems is the need to incorporate systematic uncertainties, which affect both the robustness of inference and the metrics used to evaluate a particular analysis strategy.

In this work, we focus on supervised learning techniques where systematic uncertainties can be associated to a data generation process that is not uniquely specified. In other words, the lack of systematic uncertainties corresponds to the (rare) case that the process that generates training data is unique, fully specified, and an accurate representative of the real world data. By contrast, a common situation when systematic uncertainty is present is when the training data are not representative of the real data. Several techniques for domain adaptation have been developed to create models that are more robust to this binary type of uncertainty. A more generic situation is that there are several plausible data generation processes, specified as a family parametrized by continuous nuisance parameters, as is typically found in scientific domains. In this broader context, statisticians have for long been working on robust inference techniques based on the concept of a pivot – a quantity whose distribution is invariant with the nuisance parameters (see e.g., (Degroot and Schervish, 1975)).

Assuming a probability model $p(X, Y, Z)$, where X are the data, Y are the target labels, and Z are the nuisance parameters, we consider the problem of learning a predictive model $f(X)$ for Y conditional on the observed values of X that is robust to uncertainty in the unknown value of Z . We introduce a flexible learning procedure based on adversarial networks (Goodfellow et al., 2014) for enforcing that $f(X)$ is a pivot with respect to Z . We derive theoretical results proving that the procedure converges towards a model that is both optimal and statistically independent of the nuisance parameters (if that model exists) or for which one can tune a trade-off between accuracy and robustness (e.g., as driven by a higher level objective). In particular, and to the best of our knowledge, our contribution is the first solution for imposing pivotal constraints on a predictive model, working regardless of the

pure-python fitting/limit-setting/interval estimation

HistFactory-style

DOI [10.5281/zenodo.1169739](https://doi.org/10.5281/zenodo.1169739) build **passing** coverage **96%** health **97%** docs **master** pypi package **0.0.8** launch **binder**

The HistFactory p.d.f. template [[CERN-OPEN-2012-016](#)] is per-se independent of its implementation in ROOT and sometimes, it's useful to be able to run statistical analysis outside of ROOT, RooFit, RooStats framework.

This repo is a pure-python implementation of that statistical model for multi-bin histogram-based analysis and its interval estimation is based on the asymptotic formulas of "Asymptotic formulae for likelihood-based tests of new physics" [[arxiv:1007.1727](#)]. The aim is also to support modern computational graph libraries such as PyTorch and Tensorflow in order to make use of features such as autodifferentiation and GPU acceleration.

What does it support


Implemented variations:

- HistoSys
- OverallSys
- ShapeSys
- NormFactor
- Multiple Channels
- Import from XML + ROOT via [uproot](#)
- ShapeFactor
- StatError

<https://github.com/diana-hep/pyhf>

pytorch_tests_onoff x

Secure | https://hub.mybinder.org/user/diana-hep-pyhfl-1djehjf/notebooks/examples/notebooks/pytorch_tests_onoff.ipynb ☆ ⓘ

jupyter pytorch_tests_onoff Last Checkpoint: 8 minutes ago (unsaved changes)  Logout

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3

📁 + ✂️ 📄 📄 ⬆️ ⬇️ ▶️ Run 🛑 ↺ ▶️ Code 🗣️

```
In [5]: %pylab inline

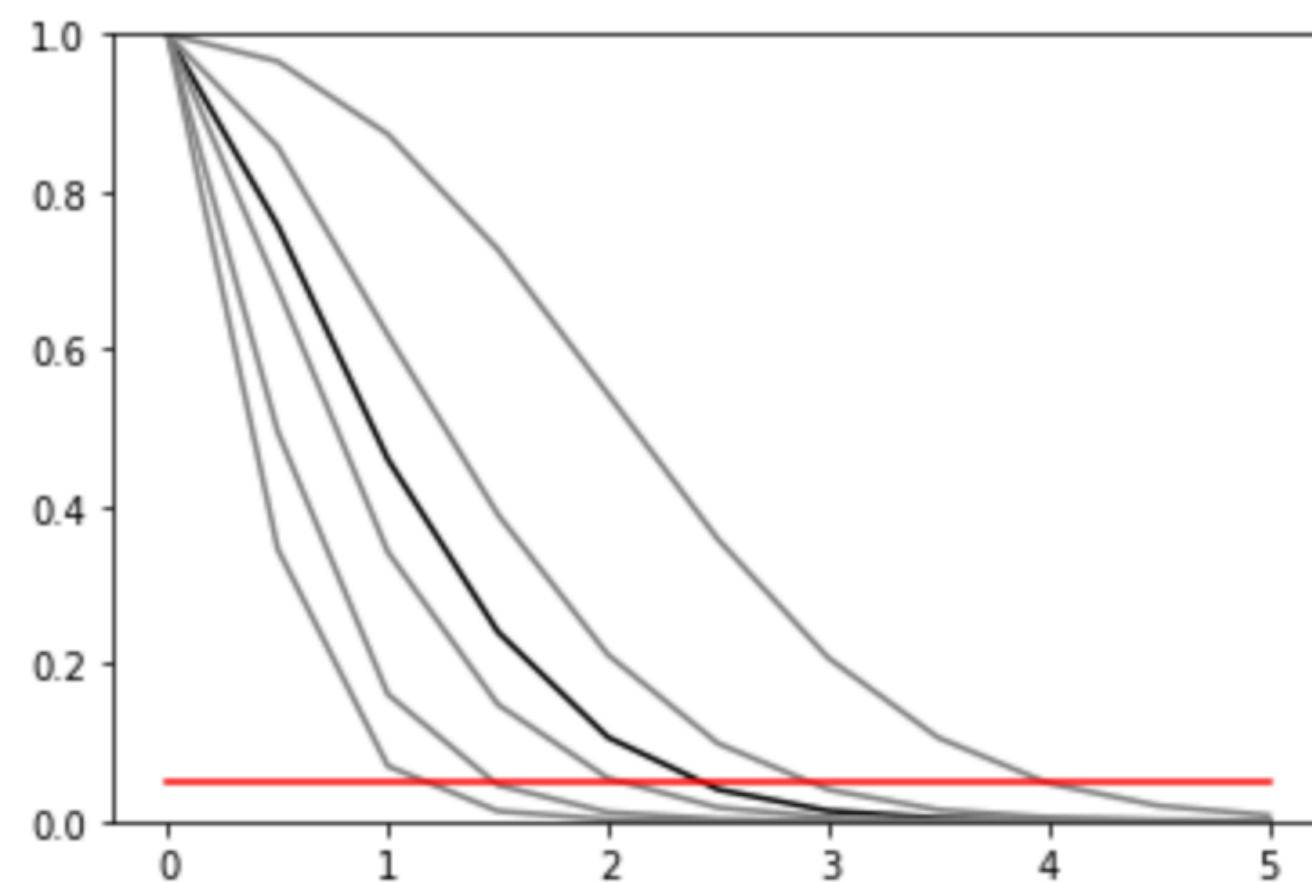
testmus = np.linspace(0,5,11)
results = []
for mu in testmus:
    results.append(runOnePoint(mu, data, pdf, init_pars, par_bounds))

obs = [1./r[-2:][0] for r in results ]
exp = [[1./r[-2:][1][i] for r in results] for i in range(5) ]

def plot_results(testmus,cls_obs, cls_exp, test_size = 0.05):
    plt.plot(testmus,cls_obs, c = 'k')
    for i,c in zip(range(5),['grey','grey','grey','grey','grey']):
        plt.plot(testmus,cls_exp[i], c = c)
    plt.plot(testmus,[test_size]*len(testmus), c = 'r')
    plt.ylim(0,1)

plot_results(testmus,obs,exp)
```

Populating the interactive namespace from numpy and matplotlib



https://mybinder.org/v2/gh/diana-hep/pyhfl/master?urlpath=tree/examples/notebooks/pytorch_tests_onoff.ipynb

pyhf

`pyhf` is a work-in-progress standalone implementation of the HistFactory p.d.f. template and an implementation of the test statistics and asymptotic formulae described in the paper by Cowan, Cranmer, Gross, Vitells: *Asymptotic formulae for likelihood-based tests of new physics* [[arxiv:1007.1727](https://arxiv.org/abs/1007.1727)].

Models can be defined using JSON specification, but existing models based on the XML + ROOT file scheme are readable as well.

The Demo

The input data for the statistical analysis was build generated using the containerized workflow engine [yadage](https://yadage.org/) (see demo from KubeCon 2018 [[youtube](https://www.youtube.com/watch?v=9v8v8v8v8v8)]). Similarly to Binder this utilizes modern container technology for reproducible science. Below you see the execution graph leading up to the model input data at the bottom.

```
In [3]: import base64
from IPython.core.display import display, HTML
anim = base64.b64encode(open('workflow.gif', 'rb').read()).decode('ascii')
HTML(''.format(anim))
```

Out[3]:

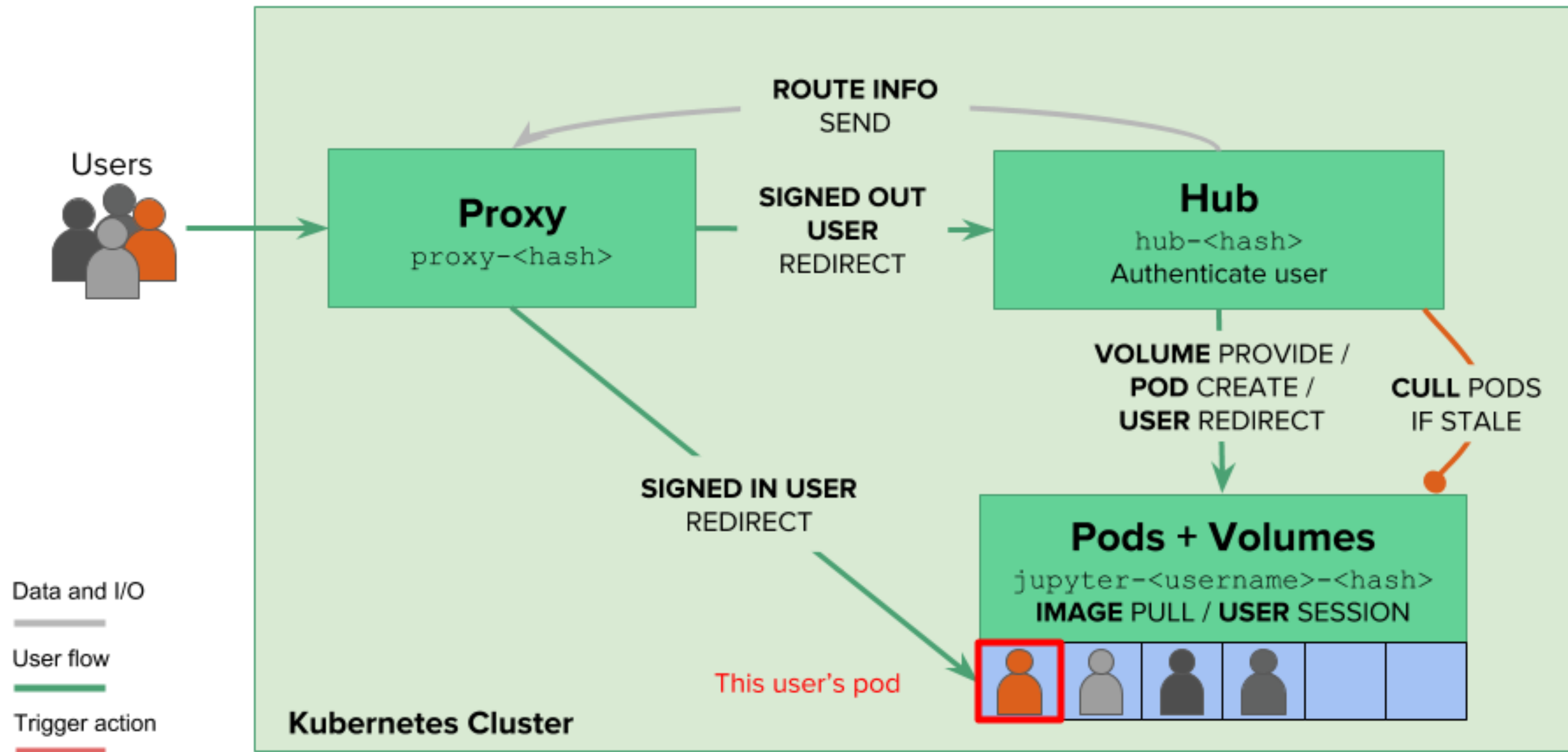


JupyterHub

JupyterHub Architecture
(high-level details)

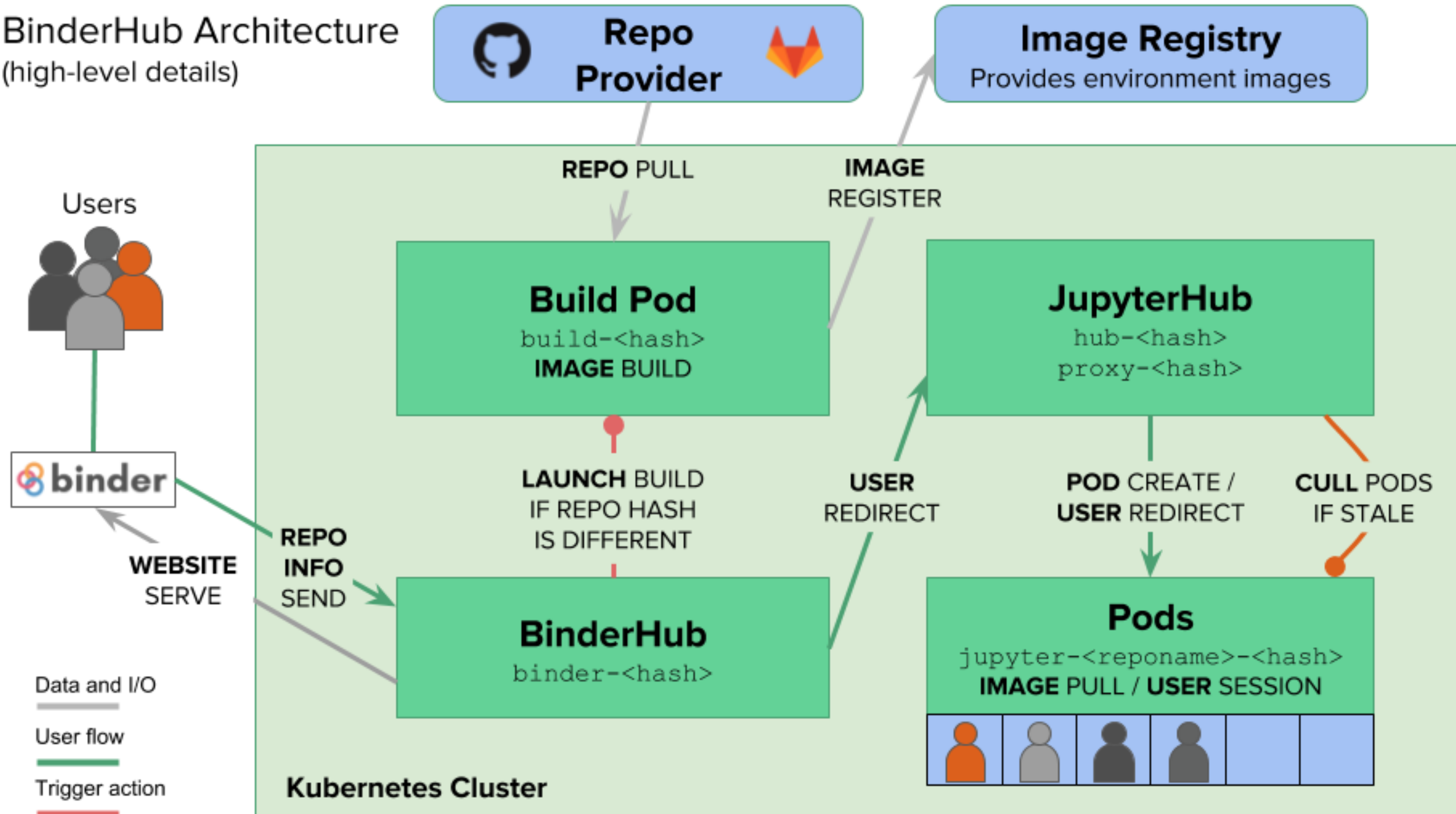
Cloud Volumes
Provides persistent storage

Image Registry
Provides environment images



JupyterHub + repo2docker = binder

BinderHub Architecture
(high-level details)



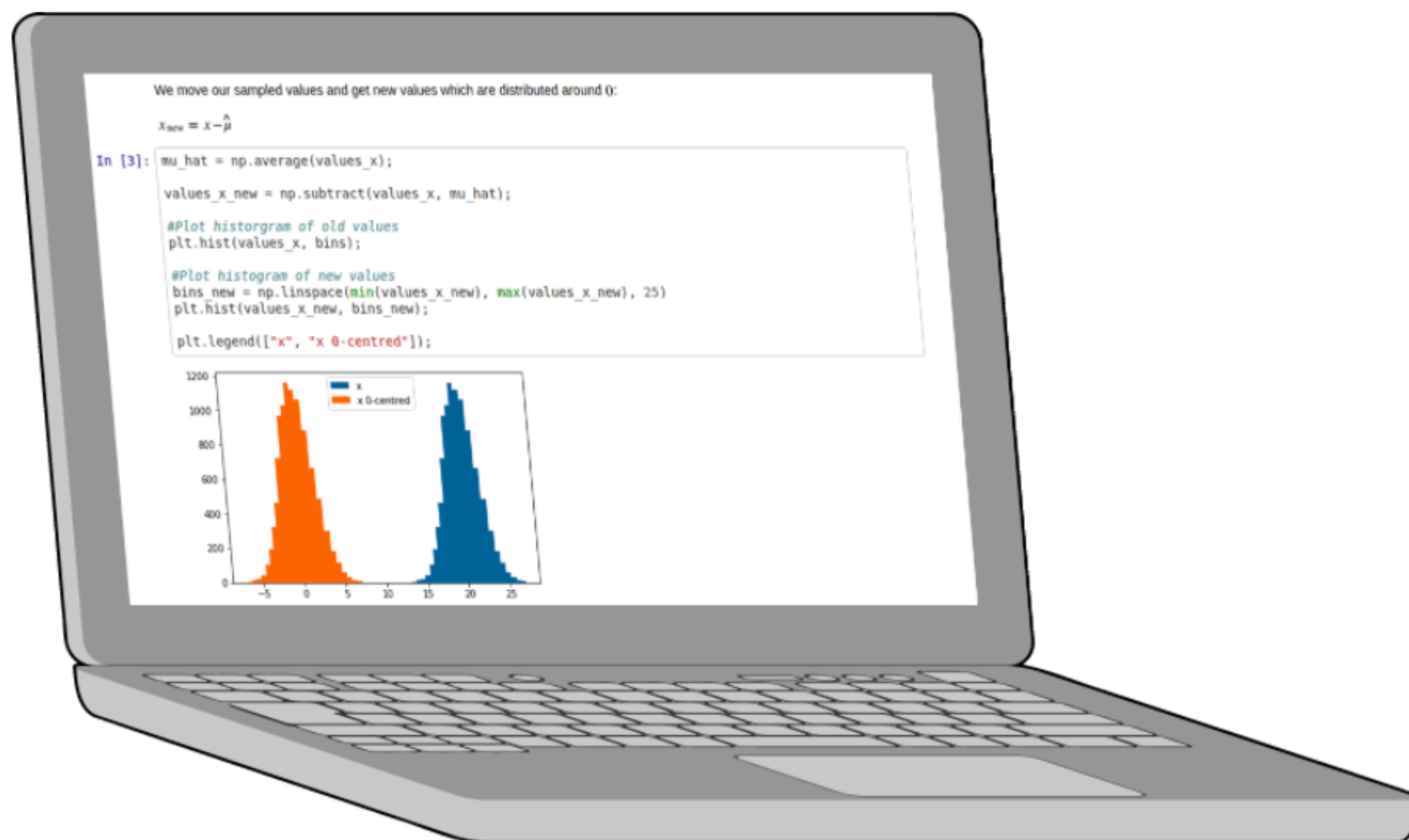


GESIS Notebooks (beta)

Home

Your Server

Binder



Open Research Computing for the Social Sciences

Your Server

Create your own JupyterHub

[Next »](#)



Zero to JupyterHub with Kubernetes

A tutorial to help install and manage JupyterHub with Kubernetes

Quick search

Zero to JupyterHub

JupyterHub is a tool that allows you to quickly utilize cloud computing infrastructure to manage a hub that enables users to interact remotely with a computing environment that you specify.

JupyterHub offers a useful way to standardize the computing environment of a group of people (e.g., for a class of students or an analytics team), as well as allowing people to access the hub remotely.

This growing collection of information will help you set up your own JupyterHub instance. It is in an early stage, so the information and tools may change quickly. If you see anything that is incorrect or have any questions, feel free to reach out at the issues page.

Creating your JupyterHub

This tutorial starts from “step zero” and walks through how to install and configure a complete JupyterHub deployment in the cloud. Using Kubernetes and the JupyterHub Helm chart provides sensible defaults for an initial deployment.

Create your own BinderHub

Zero to BinderHub

A guide to help you create your own BinderHub from scratch.

- 1. Create your cloud resources
 - 1.1. Setting up Kubernetes on Google Cloud
 - 1.2. Install Helm
 - 1.3. Set up the container registry
- 2. Set up BinderHub
 - 2.1. Preparing to install
 - 2.2. Create **secret.yaml** file
 - 2.3. Create **config.yaml**
 - 2.4. Install BinderHub
 - 2.5. Connect BinderHub and JupyterHub
 - 2.6. Try out your BinderHub Deployment
 - 2.7. Increase your GitHub API limit
- 3. Tear down your Binder deployment
 - 3.1. Contracting the size of your cluster
 - 3.2. Deleting the cluster <https://binderhub.readthedocs.io/>





BinderHub

A project to build and serve Binders

 Watch

41

Table of Contents

[1. Create your cloud resources](#)

[2. Set up BinderHub](#)

[3. Tear down your Binder deployment](#)

[1. The BinderHub Architecture](#)

5. BinderHub Deployments

BinderHub is open-source technology that can be deployed anywhere that Kubernetes is deployed. The Binder community hopes that it will be used for many applications in research and education. As new organizations adopt BinderHub, we'll update this page in order to provide inspiration to others who wish to do so.

If you or your organization has set up a BinderHub that isn't listed here, please [open an issue](#) on our GitHub repository to discuss adding it!

5.1. GESIS Institute for Social Sciences

Deployed on bare-metal using `kubeadm`.

- [Deployment repository](#)
- [BinderHub / JupyterHub links](#)

[« 4. BinderHub API Documentation](#)

[6. Configuration and Source Code Ref...](#)

The scope of the problem

An empirical analysis of journal policy effectiveness for computational reproducibility

Victoria Stodden^{a,1}, Jennifer Seiler^b, and Zhaokun Ma^b

^aSchool of Information Sciences, University of Illinois at Urbana–Champaign, Champaign, IL 61820; and ^bDepartment of Statistics, Columbia University, New York, NY 10027

Edited by David B. Allison, Indiana University Bloomington, Bloomington, IN, and accepted by Editorial Board Member Susan T. Fiske January 9, 2018 (received for review July 11, 2017)

A key component of scientific communication is sufficient information for other researchers in the field to reproduce published findings. For computational and data-enabled research, this has often been interpreted to mean making available the raw data from which results were generated, the computer code that generated the findings, and any additional information needed such as workflows and input parameters. Many journals are revising author guidelines to include data and code availability. This work evaluates the effectiveness of journal policy that requires the data and code necessary for reproducibility be made available postpublication by the authors upon request. We assess the effectiveness of such a policy by (i) requesting data and code from authors and (ii) attempting replication of the published findings. We chose a random sample of 204 scientific papers published in the journal *Science* after the implementation of their policy in February 2011. We found that we were able to obtain artifacts from 44% of our sample and were able to reproduce the findings for 26%. We find this policy—author remission of data and code postpublication upon request—an improvement over no policy, but currently insufficient for reproducibility.

reproducible research | data access | code access | reproducibility policy | open science

computational reproducibility of published results. We use a survey instrument to test the availability of data and code for articles published in *Science* in 2011–2012. We then use the scientific communication standards from the 2012 Institute for Computational and Experimental Research in Mathematics (ICERM) workshop report to evaluate the reproducibility of articles for which artifacts were made available (11). We then assess the impact of the policy change directly, by examining articles published in *Science* in 2009–2010 and comparing artifact ability to our postpolicy sample from 2011–2012. Finally, we discuss possible improvements to journal policies for enabling reproducible computational research in light of our results.

Results

We emailed corresponding authors in our sample to request the data and code associated with their articles and attempted to replicate the findings from a randomly chosen subset of the articles for which we received artifacts. We estimate the artifact recovery rate to be 44% with a 95% bootstrap confidence interval of the proportion [0.36, 0.50], and we estimate the replication rate to be 26% with a 95% bootstrap confidence interval [0.20, 0.32].



The scope of the problem

An empirical analysis of journal policy effectiveness for computational reproducibility

Victoria Stodden^{a,1}, Jennifer Seiler^b, and Zhaokun Ma^b

^aSchool of Information Sciences, University of Illinois at Urbana–Champaign, Champaign, IL 61820; and ^bDepartment of Statistics, Columbia University, New York, NY 10027

Edited by David B. Allison, Indiana University Bloomington, Bloomington, IN, and accepted by Editorial Board Member Susan T. Fiske January 9, 2018 (received for review July 11, 2017)

A key component of scientific communication is sufficient information for other researchers in the field to reproduce published findings. For computational and data-enabled research, this has often been interpreted to mean making available the raw data from which results were generated, the computer code that generated the findings, and any additional information needed such as workflows and input parameters. Many journals are revising author guidelines to include data and code availability. This work evaluates the effectiveness of journal policy that requires the data and code necessary for reproducibility be made available postpublication by the authors upon request. We assess the effectiveness of such a policy by (i) requesting data and code from authors and (ii) attempting replication of the published findings. We chose a random sample of 204 scientific papers published in the journal *Science* after the implementation of their policy in February 2011. We found that we were able to obtain artifacts from 44% of our sample and were able to reproduce the findings for 26%. We find this policy—author remission of data and code postpublication upon request—an improvement over no policy, but currently insufficient for reproducibility.

reproducible research | data access | code access | reproducibility policy | open science

computational reproducibility of published results. We use a survey instrument to test the availability of data and code for articles published in *Science* in 2011–2012. We then use the scientific communication standards from the 2012 Institute for Computational and Experimental Research in Mathematics (ICERM) workshop report to evaluate the reproducibility of articles for which artifacts were made available (11). We then assess the impact of the policy change directly, by examining articles published in *Science* in 2009–2010 and comparing artifact ability to our postpolicy sample from 2011–2012. Finally, we discuss possible improvements to journal policies for enabling reproducible computational research in light of our results.

Results

We emailed corresponding authors in our sample to request the data and code associated with their articles and attempted to replicate the findings from a randomly chosen subset of the articles for which we received artifacts. We estimate the artifact recovery rate to be 44% with a 95% bootstrap confidence interval of the proportion [0.36, 0.50], and we estimate the replication rate to be 26% with a 95% bootstrap confidence interval [0.20, 0.32].



We emailed corresponding authors in our sample to request the data and code associated with their articles and attempted to replicate the findings from a randomly chosen subset of the articles for which we received artifacts. We estimate the artifact recovery rate to be 44% with a 95% bootstrap confidence interval of the proportion [0.36, 0.50], and we estimate the replication rate to be 26% with a 95% bootstrap confidence interval [0.20, 0.32].

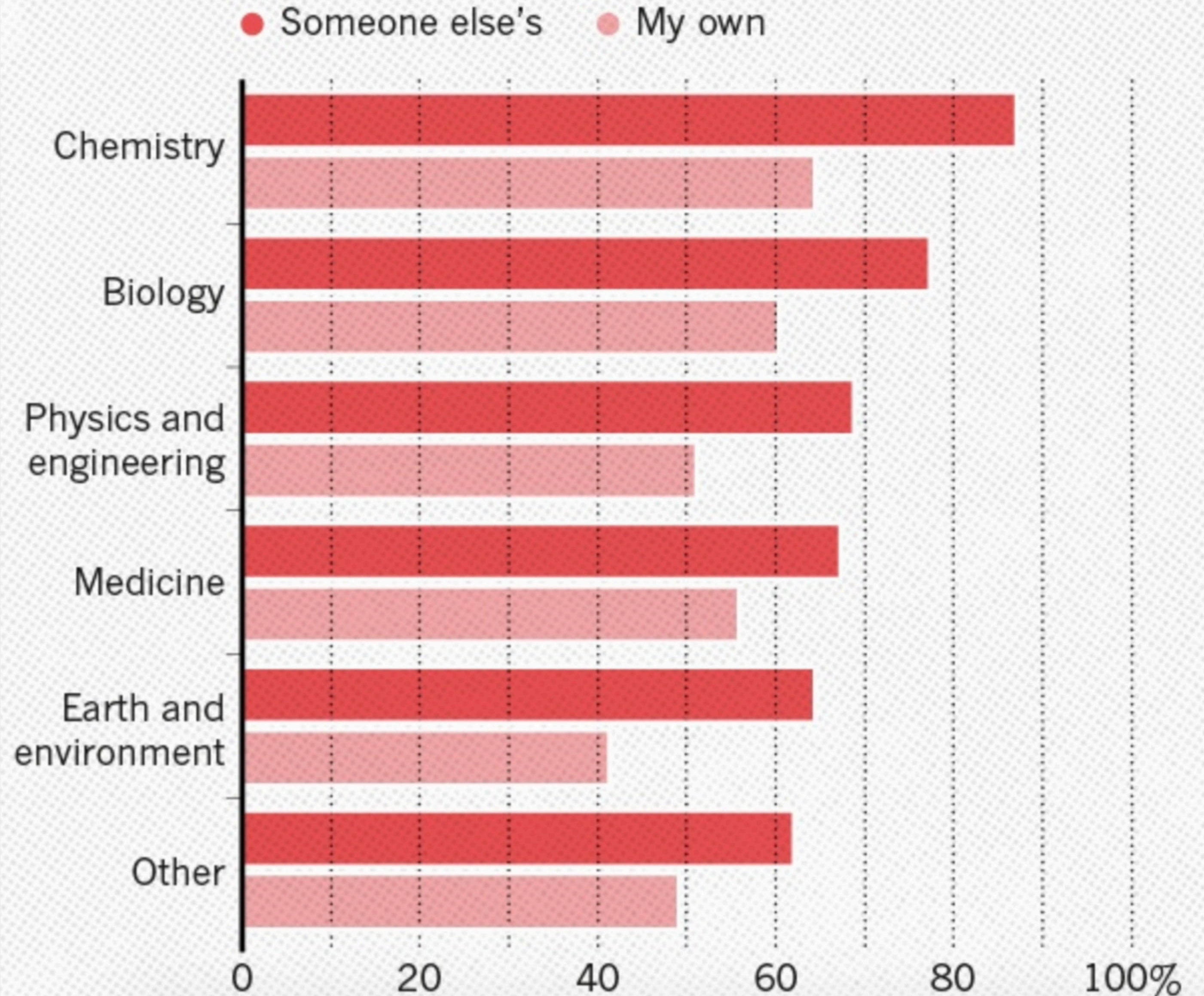
- Stodden et al.



Moyna Baker,
2016 Nature
(n=1,500)

HAVE YOU FAILED TO REPRODUCE AN EXPERIMENT?

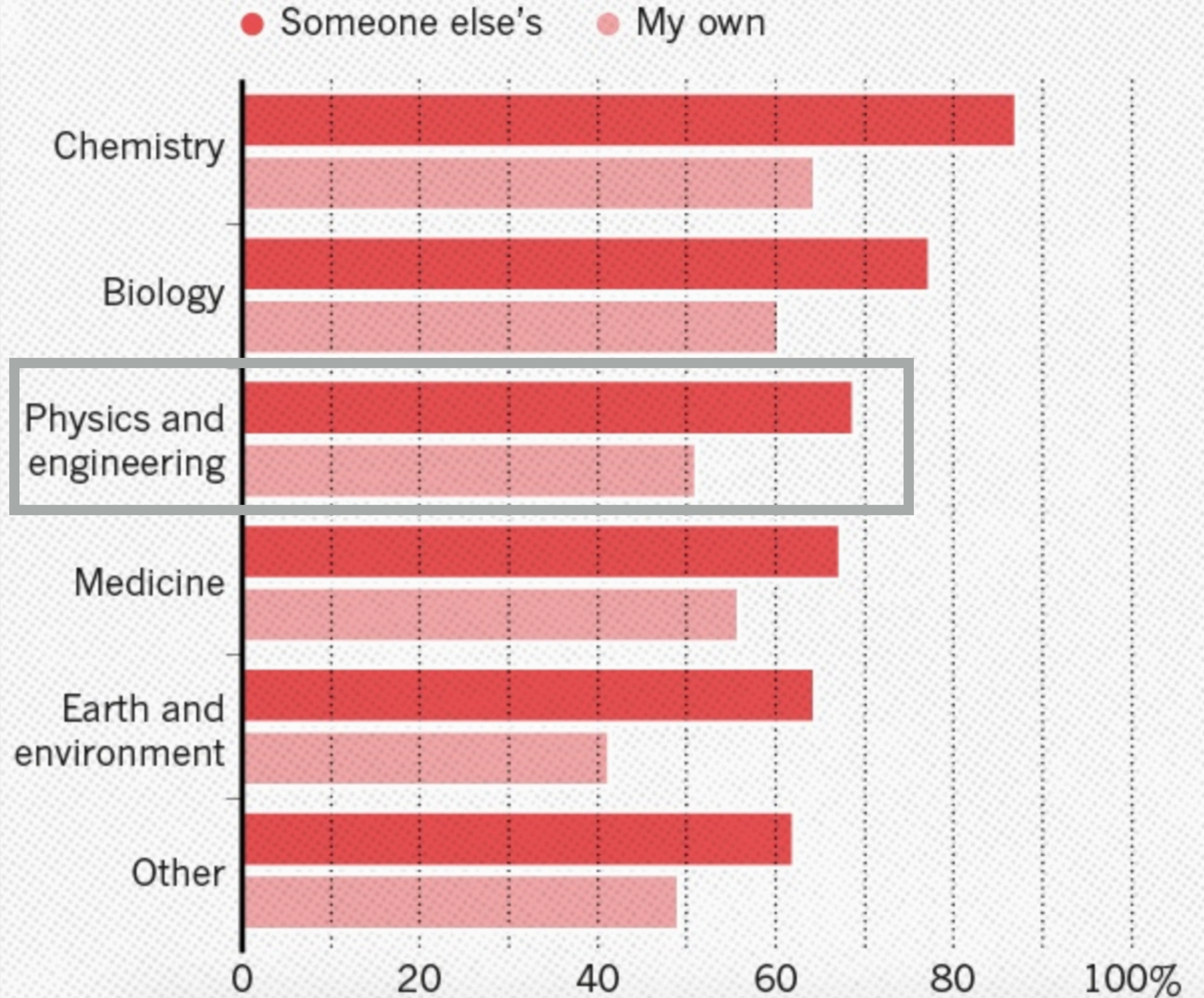
Most scientists have experienced failure to reproduce results.



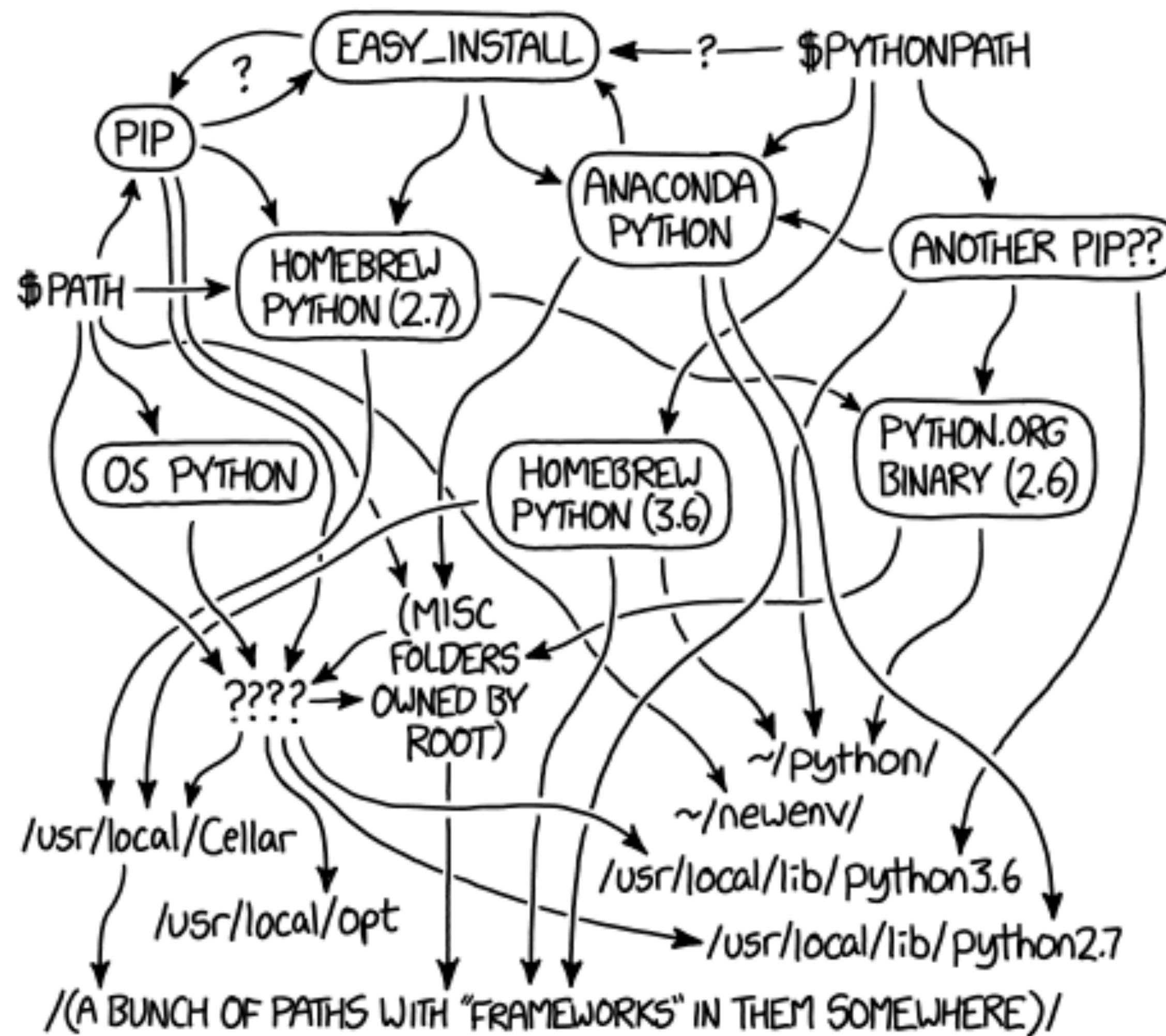
Moyna Baker,
2016 Nature

HAVE YOU FAILED TO REPRODUCE AN EXPERIMENT?

Most scientists have experienced failure to reproduce results.



GitHub repo's aren't enough



MY PYTHON ENVIRONMENT HAS BECOME SO DEGRADED THAT MY LAPTOP HAS BEEN DECLARED A SUPERFUND SITE.



Report what you've actually discovered, clearly enough that someone else can discover it for themselves.

- James Somers



The maintainer's mind

Install pyhf from master to get notebooks to run on binder

Edit

#151

 Open jzf2101 wants to merge 3 commits into `diana-hep:master` from `jzf2101:master`

 Conversation 5

 Commits 3

 Files changed 2

+1 -1 



The maintainer's mind



matthewfeickert commented 18 hours ago • edited ▾

Collaborator



Hi @jzf2101, wow I can't tell you how happy this makes me to see someone from the Jupyter team fixing up the Binder env I wrote! 😊

I had also noticed this behavior (had to open up a terminal in Binder and run `python setup.py install` first) but I hadn't gotten around to fixing it yet. Thanks so much for doing this for me!

LGTM, so unless @lukasheinrich has any comments I think we can squash and merge this in.



lukasheinrich commented 18 hours ago

Collaborator



Thanks @jzf2101 for the PR. I agree, the binder from the repo, should take the code from the repo, not from PyPI. I think perhaps we should install in "editable" mode e.g. `pip install -e .`

That way it should be possible to edit the pyhf code in-place, reload the python modules and use the edited version. Other than that this LGTM.



The scope of the problem

An empirical analysis of journal policy effectiveness for computational reproducibility

Victoria Stodden^{a,1}, Jennifer Seiler^b, and Zhaokun Ma^b

^aSchool of Information Sciences, University of Illinois at Urbana–Champaign, Champaign, IL 61820; and ^bDepartment of Statistics, Columbia University, New York, NY 10027

Edited by David B. Allison, Indiana University Bloomington, Bloomington, IN, and accepted by Editorial Board Member Susan T. Fiske January 9, 2018 (received for review July 11, 2017)

A key component of scientific communication is sufficient information for other researchers in the field to reproduce published findings. For computational and data-enabled research, this has often been interpreted to mean making available the raw data from which results were generated, the computer code that generated the findings, and any additional information needed such as workflows and input parameters. Many journals are revising author guidelines to include data and code availability. This work evaluates the effectiveness of journal policy that requires the data and code necessary for reproducibility be made available postpublication by the authors upon request. We assess the effectiveness of such a policy by (i) requesting data and code from authors and (ii) attempting replication of the published findings. We chose a random sample of 204 scientific papers published in the journal *Science* after the implementation of their policy in February 2011. We found that we were able to obtain artifacts from 44% of our sample and were able to reproduce the findings for 26%. We find this policy—author remission of data and code postpublication upon request—an improvement over no policy, but currently insufficient for reproducibility.

reproducible research | data access | code access | reproducibility policy | open science

computational reproducibility of published results. We use a survey instrument to test the availability of data and code for articles published in *Science* in 2011–2012. We then use the scientific communication standards from the 2012 Institute for Computational and Experimental Research in Mathematics (ICERM) workshop report to evaluate the reproducibility of articles for which artifacts were made available (11). We then assess the impact of the policy change directly, by examining articles published in *Science* in 2009–2010 and comparing artifact ability to our postpolicy sample from 2011–2012. Finally, we discuss possible improvements to journal policies for enabling reproducible computational research in light of our results.

Results

We emailed corresponding authors in our sample to request the data and code associated with their articles and attempted to replicate the findings from a randomly chosen subset of the articles for which we received artifacts. We estimate the artifact recovery rate to be 44% with a 95% bootstrap confidence interval of the proportion [0.36, 0.50], and we estimate the replication rate to be 26% with a 95% bootstrap confidence interval [0.20, 0.32].



Scientific communication has become software engineering



Always code and comment in such a way that if someone a few notches junior picks up the code, they will take pleasure in reading and learning from it.

- Code for the Maintainer



As scientists, we are now all becoming open source maintainers



As scientists, we are now all becoming open source maintainers
Welcome to the club :D





Genomic analysis of elongated skulls reveals extensive female-biased immigration to Early Medieval Bavaria

Krishna R. Veeramah^a, Andreas Rott^{b,1}, Melanie Groß^{c,1}, Lucy van Dorp^d, Saioa López^e, Karola Kiršanow^c, Christian Sell^c, Jens Blöcher^c, Daniel Wegmann^{f,g}, Vivian Link^{f,g}, Zuzana Hofmanová^{f,g}, Joris Peters^{b,h}, Bernd Trautmann^b, Anja Gairhosⁱ, Jochen Haberstroh^j, Bernd Pääffgen^k, Garrett Hellenthal^d, Brigitte Haas-Gebhardⁱ, Michaela Harbeck^{b,2,3}, and Joachim Burger^{c,2,3}

^aDepartment of Ecology and Evolution, Stony Brook University, Stony Brook, NY 11794-5245; ^bState Collection for Anthropology and Palaeoanatomy, Bavarian Natural History Collections, 80333 Munich, Germany; ^cPalaeogenetics Group, Institute of Organismic and Molecular Evolution, Johannes Gutenberg University Mainz, 55099 Mainz, Germany; ^dUCL Genetics Institute, Department of Genetics, Evolution and Environment, University College London, WC1E 6BT London, United Kingdom; ^eCancer Institute, University College London, WC1E 6DD London, United Kingdom; ^fDepartment of Biology, University of Fribourg, 1700 Fribourg, Switzerland; ^gSwiss Institute of Bioinformatics, 1700 Fribourg, Switzerland; ^hArchaeoBioCenter and Institute for Palaeoanatomy, Limnology Research and the History of Veterinary Medicine, Ludwig Maximilian University, 80539 Munich, Germany; ⁱBavarian State Archaeological Collection, 80539 Munich, Germany; ^jBavarian State Department of Monuments, 80339 Munich, Germany; ^kInstitute of Prehistoric and Protohistoric Archaeology, Ludwig Maximilian University, 80539 Munich, Germany

Edited by Eske Willerslev, University of Copenhagen, Copenhagen, Denmark, and approved January 30, 2018 (received for review November 21, 2017)

Modern European genetic structure demonstrates strong correlations with geography, while genetic analysis of prehistoric humans has indicated at least two major waves of immigration from outside the continent during periods of cultural change. However, population-level genome data that could shed light on the demographic processes occurring during the intervening periods have been absent. Therefore, we generated genomic data

to form in the 5th century AD, and that it emanated from a combination of the romanized local population of the border province of the former Roman Empire and immigrants from north of the Danube (2). While the Baiuvarii are less well known than some other contemporary groups, an interesting archaeological feature in Bavaria from this period is the presence of skeletons with artificially deformed or elongated skulls (Fig. 1A).

POPULATION BIOLOGY

PNAS / Richard Goerg / Getty / The Atlantic

The Scientific Paper Is Obsolete

Here's what's next.



Atlantic Monthly, April 5 2018

Here's what's next





BinderHub

A project to build and serve Binders

 Watch

41

Table of Contents

[1. Create your cloud resources](#)

[2. Set up BinderHub](#)

[3. Tear down your Binder deployment](#)

[1. The BinderHub Architecture](#)

5. BinderHub Deployments

BinderHub is open-source technology that can be deployed anywhere that Kubernetes is deployed. The Binder community hopes that it will be used for many applications in research and education. As new organizations adopt BinderHub, we'll update this page in order to provide inspiration to others who wish to do so.

If you or your organization has set up a BinderHub that isn't listed here, please [open an issue](#) on our GitHub repository to discuss adding it!

5.1. GESIS Institute for Social Sciences

Deployed on bare-metal using `kubeadm`.

- [Deployment repository](#)
- [BinderHub / JupyterHub links](#)

[« 4. BinderHub API Documentation](#)

[6. Configuration and Source Code Ref...](#)

Create your own BinderHub

Zero to BinderHub

A guide to help you create your own BinderHub from scratch.

- 1. Create your cloud resources
 - 1.1. Setting up Kubernetes on Google Cloud
 - 1.2. Install Helm
 - 1.3. Set up the container registry
- 2. Set up BinderHub
 - 2.1. Preparing to install
 - 2.2. Create **secret.yaml** file
 - 2.3. Create **config.yaml**
 - 2.4. Install BinderHub
 - 2.5. Connect BinderHub and JupyterHub
 - 2.6. Try out your BinderHub Deployment
 - 2.7. Increase your GitHub API limit
- 3. Tear down your Binder deployment
 - 3.1. Contracting the size of your cluster
 - 3.2. Deleting the cluster <https://binderhub.readthedocs.io/>




 [jupyterhub](#) / [binderhub](#)


 Unwatch ▾ 41


 Star 432

 Fork 65


 Code

 Issues 112

 Pull requests 12

 Projects 1

 Wiki

 Insights

 Settings

Federation support #63

Edit

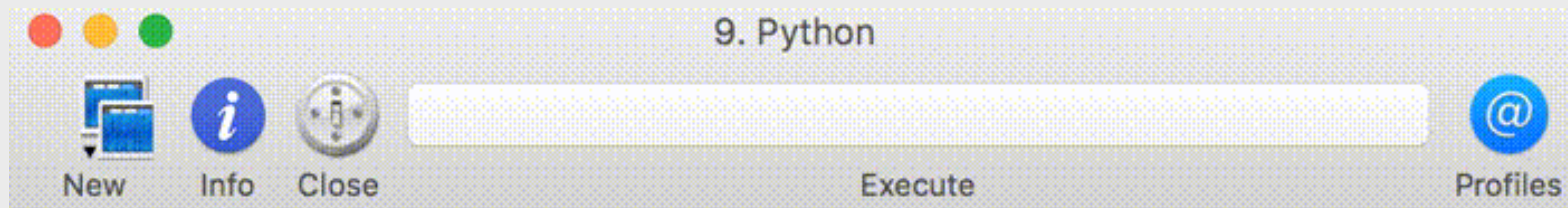
New issue

 Open

yuvipanda opened this issue on Jun 24, 2017 · 10 comments

Show, don't tell





```
-bash-3.2$ ipython
WARNING: Attempting to work in a virtualenv. If you encounter problems, please
ninstall IPython inside the virtualenv.
Python 2.7.10 (default, Sep 23 2015, 04:34:21)
Type "copyright", "credits" or "license" for more information.

IPython 4.0.2 -- An enhanced Interactive Python.
?          -> Introduction and overview of IPython's features.
%quickref -> Quick reference.
help       -> Python's own help system.
object?    -> Details about 'object', use 'object??' for extra details.

In [1]: from postal.expand import expand_address

In [2]:
```


dtak / rrr

Unwatch 5

Star 17

Fork 4

Code

Issues 2

Pull requests 0

Projects 0

Wiki

Insights

Code/figures in Right for the Right Reasons <https://www.ijcai.org/proceedings/201...>

34 commits

2 branches

0 releases

2 contributors

MIT

Toy Colors.ipyn x

Code Python 3

color_dataset_generator for more details.

We will train a multilayer perceptron to classify these images, explore which rule(s) it implicitly learns, and constrain it to use only one rule (or neither).

Let's first load our dataset:

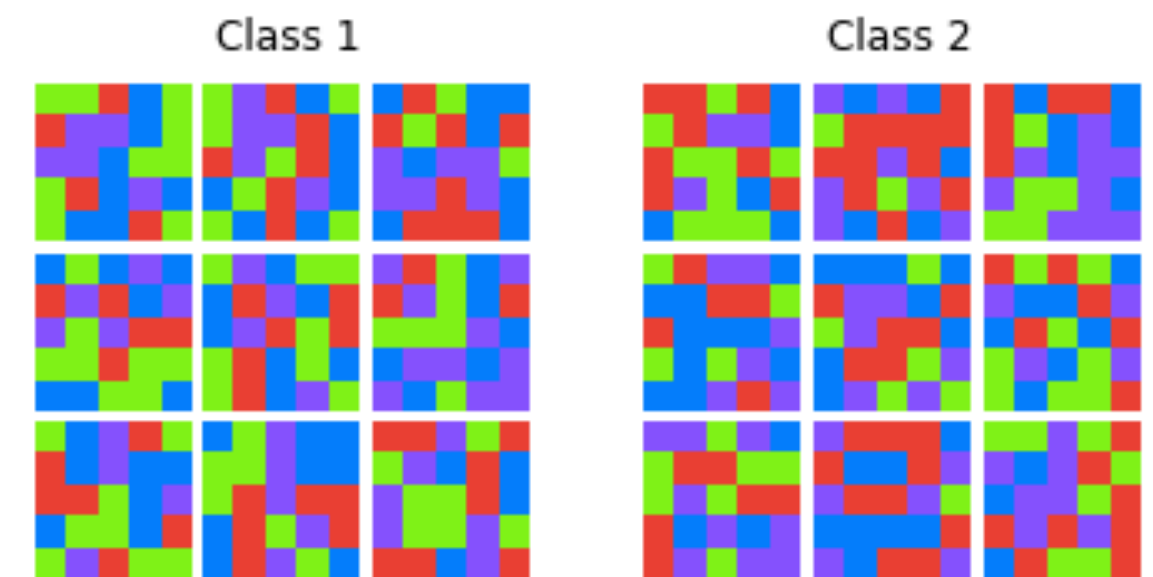
```
In [2]: X, Xt, y, yt = generate_dataset(cachefile='../data/toy-col...
E1 = np.array([ignore_rule2 for _ in range(len(y))])
E2 = np.array([ignore_rule1 for _ in range(len(y))])
```

```
In [3]: print(X.shape, Xt.shape, y.shape, yt.shape, E1.shape, E2.s...
(20000, 75) (20000, 75) (20000,) (20000,) (20000, 75) (200...
00, 75)
```

Understanding the Dataset

Let's just examine images from each class quickly and verify that in class 1, the corners are all the same color and the top-middle three pixels are all different (none of which should hold true in class 2):

```
In [4]: plt.subplot(121)
plt.title('Class 1')
image_grid(X[np.argwhere(y == 0)[:9]], (5,5,3), 3)
plt.subplot(122)
plt.title('Class 2')
image_grid(X[np.argwhere(y == 1)[:9]], (5,5,3), 3)
plt.show()
```

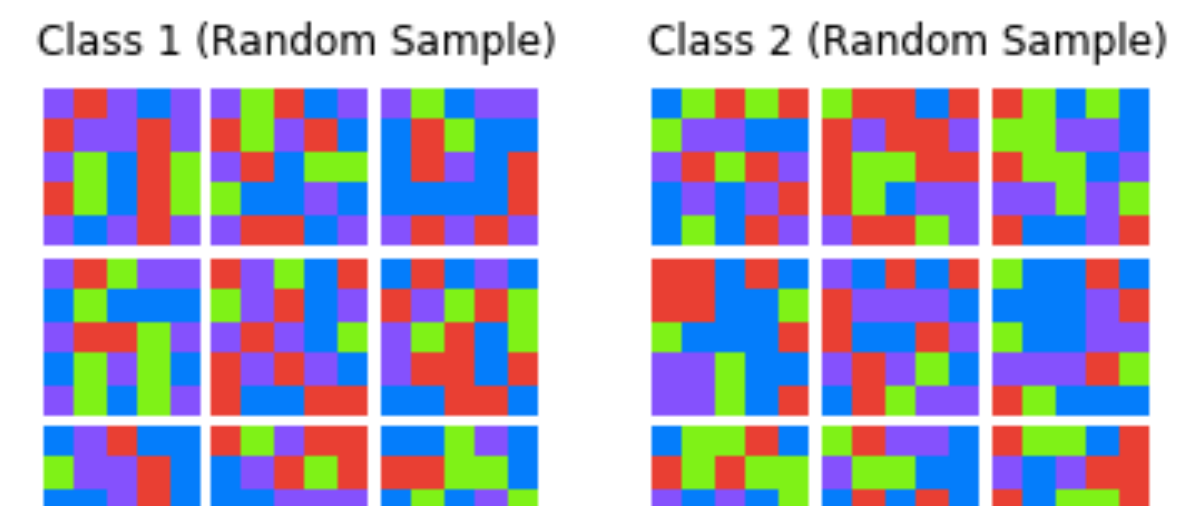


toy_colors.py x

```
1 import numpy as np
2 import os
3
4 colors = {
5     'r': np.array([255, 0, 0], dtype=np.uint8),
6     'o': np.array([255, 128, 0], dtype=np.uint8),
7     'y': np.array([255, 255, 0], dtype=np.uint8),
8     'g': np.array([0, 255, 0], dtype=np.uint8),
9     'b': np.array([0, 128, 255], dtype=np.uint8),
10    'i': np.array([0, 0, 255], dtype=np.uint8),
11    'v': np.array([128, 0, 255], dtype=np.uint8)
12 }
13
14 imglen = 5
15 imgshape = (imglen, imglen, 3)
16 topleft = [0,[0]]
17 topright = [0,[imglen-1]]
18 botleft = [imglen-1,[0]]
19 botright = [imglen-1,[imglen-1]]
20
21 ignore_rule1 = np.zeros(imgshape)
```

Cell: Toy Colors x

```
In [19]: import numpy as np
plt.subplot(121)
plt.title('Class 1 (Random Sample)')
image_grid(X[np.argwhere(y == 0)\
    [np.random.choice(X[np.argwhere(y == 0)].shape[0], 9,
        replace=True)]], (5,5,3), 3)
plt.subplot(122)
plt.title('Class 2 (Random Sample)')
image_grid(X[np.argwhere(y == 1)\
    [np.random.choice(X[np.argwhere(y == 1)].shape[0], 9,
        replace=True)]], (5,5,3), 3)
```



What if?

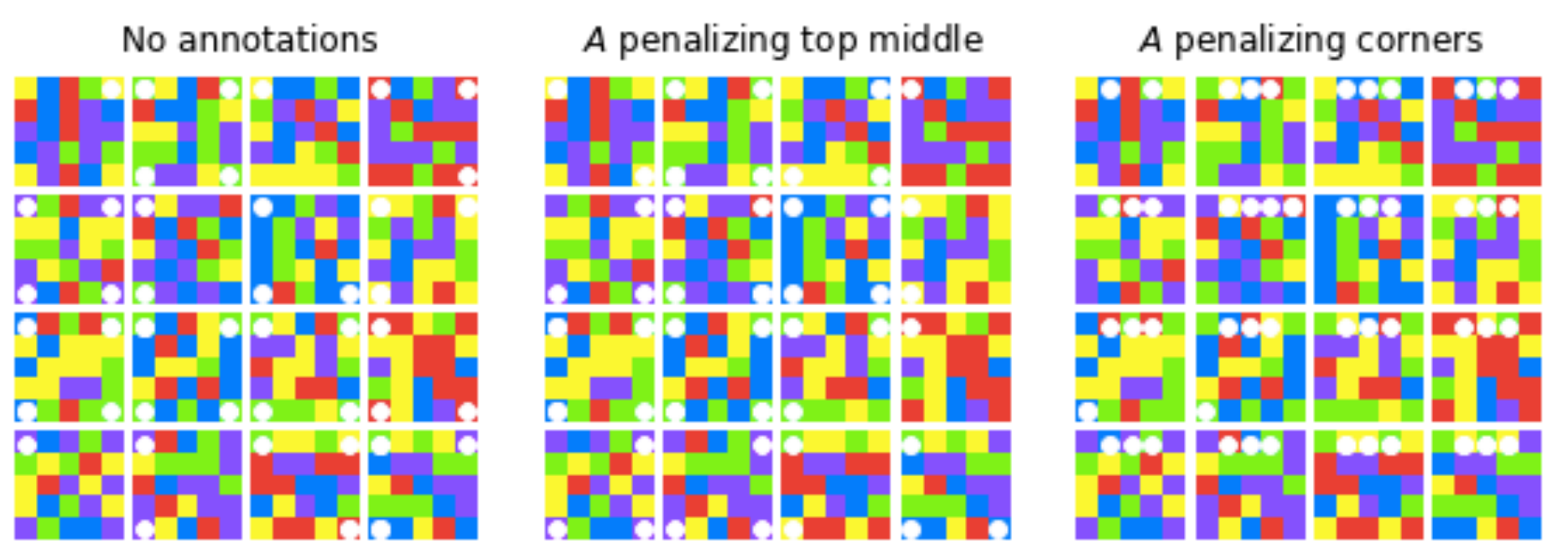


Out[7]: 0.9987000000000000

```
In [8]: # Train a model constrained to use the second rule
mlp_rule2 = MultilayerPerceptron(l2_grads=1000)
mlp_rule2.fit(X, y, E2)
mlp_rule2.score(Xt, yt)
```

Out[8]: 0.9880499999999999

```
In [19]: # Visualize largest weights
with figure_grid(1,3, rowwidth=8) as g:
    g.next()
    explain(mlp_plain, 'No annotations')
    g.next()
    explain(mlp_rule1, '$A$ penalizing top middle')
    g.next()
    explain(mlp_rule2, '$A$ penalizing corners')
```



Notice that when we explicitly penalize corners or the top middle, the model appears to learn the *other* rule perfectly. We haven't identified the pixels it does treat as significant in any way, but they are significant, so the fact that they show up in the explanations means that the explanation is probably an accurate reflection of the model's implicit logic.

When we don't have any annotations, the model does identify the top-middle pixels occasionally, suggesting it defaults to learning a heavily but not completely corner-weighted combination of the rules.

What happens when we forbid it from using either rule?

```
7 'y': np.array([255, 255, 0], dtype=np.uint8),
8 'g': np.array([0, 255, 0], dtype=np.uint8),
9 'b': np.array([0, 128, 255], dtype=np.uint8),
10 'i': np.array([0, 0, 255], dtype=np.uint8),
11 'v': np.array([128, 0, 255], dtype=np.uint8)
12 }
13
14 imglen = 5
15 imgshape = (imglen, imglen, 3)
16 topleft = [0,[0]]
17 topright = [0,[imglen-1]]
18 botleft = [imglen-1,[0]]
19 botright = [imglen-1,[imglen-1]]
20
21 ignore_rule1 = np.zeros(imgshape)
22 for corner in [topleft, topright, botleft, botright]:
23     ignore_rule1[corner] = 1
24 ignore_rule1 = ignore_rule1.ravel().astype(bool)
25
26 ignore_rule2 = np.zeros(imgshape)
27 ignore_rule2[[0,[1,2,3]]] = 1
28 ignore_rule2 = ignore_rule2.ravel().astype(bool)
29
30 def random_color():
31     return colors[np.random.choice(['r','g','b','v','y'])]
32
33 def Bern(p):
34     return np.random.rand() < p
35
36 def any_repeats(row):
37     n_unique = len(set(tuple(c) for c in row))
38     return n_unique < len(row)
39
40 def ensure_class_0_rules_apply(img):
41     # Rule 1
42     img[topleft] = img[botright]
43     img[topright] = img[botright]
44     img[botleft] = img[botright]
45
46     # Rule 2
47     toprow = img[0]
48     while any_repeats(toprow[1:-1]):
49         toprow[1+np.random.choice(imglen-2)] = random_color()
50
51 def ensure_class_1_rules_apply(img):
```

Contribute, get involved

- jupyter.org
- Twitter:
 - [@mybinderteam](https://twitter.com/mybinderteam)
 - [@projectjupyter](https://twitter.com/projectjupyter)



Contribute, get involved

- GitHub:
 - github.com/jupyterhub
 - github.com/jzf2101
- We read issues!
- We mark issues for new contributors as 'good first issue' or 'help wanted'
- github.com/jupyter/governance for code of conduct



Contribute, get involved

- Gitter
 - gitter.im/jupyterhub/jupyterhub
 - gitter.im/jupyterhub/binder
 - gitter.im/jupyter/jupyter



Thanks to all the maintainers



jupyter.org

And the contributors too!



jupyter.org



Acknowledgements

GORDON AND BETTY
MOORE
FOUNDATION



Alfred P. Sloan
FOUNDATION

THE LEONA M. AND HARRY B.
HELMSLEY
CHARITABLE TRUST

