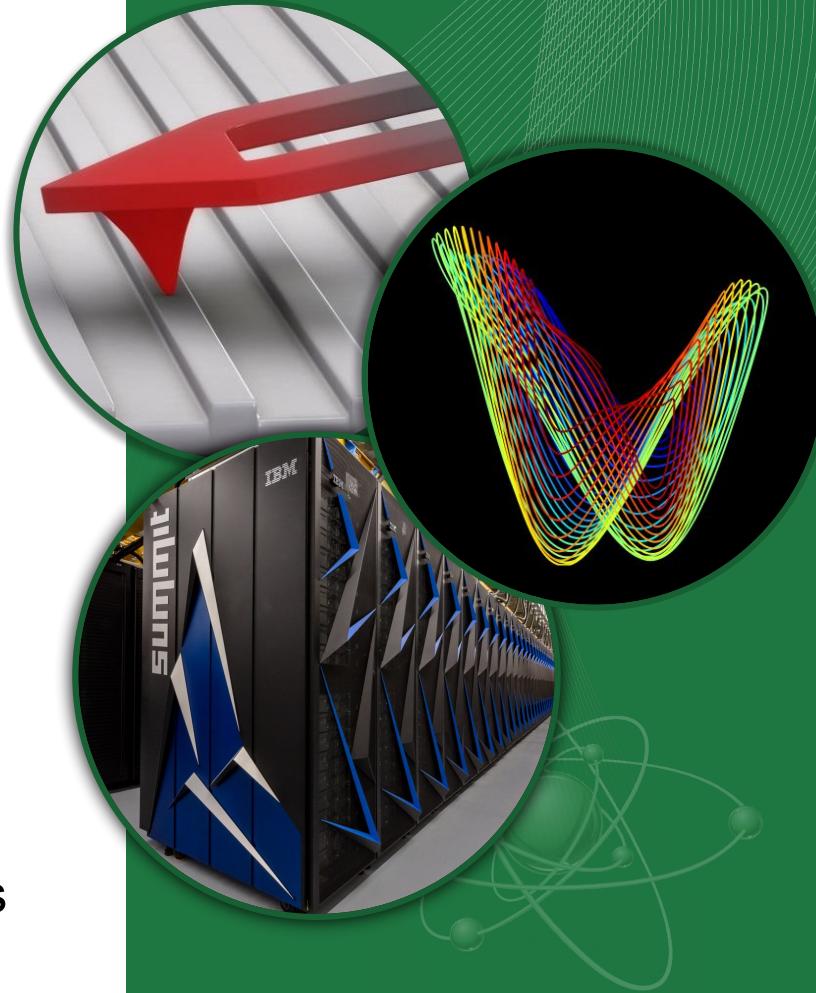


# The Jupyter Ecosystem & its Impact on Data- Analytics & Science

03/07/2019

Suhas Somnath

Advanced Data and Workflows Group  
National Center for Computational Sciences

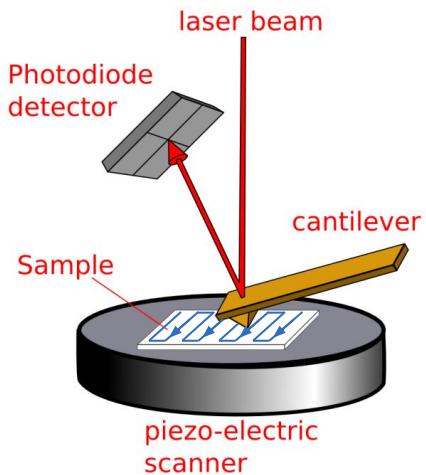


ORNL is managed by UT-Battelle  
for the US Department of Energy

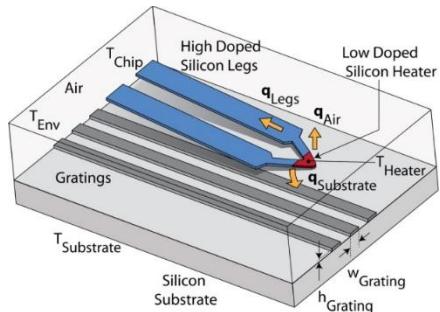
# About me...

- PhD in Mechanical Engineering (2014)
- Postdoctoral Researcher (2014 - 2017)
- Computer / Data Scientist (2017 – Present)

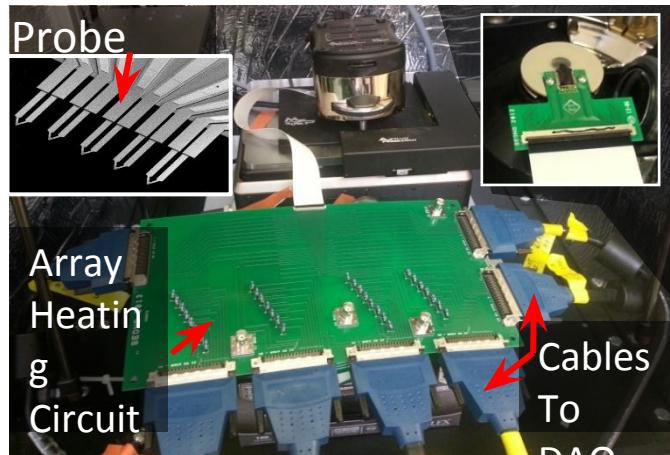
## Atomic Force Microscopy



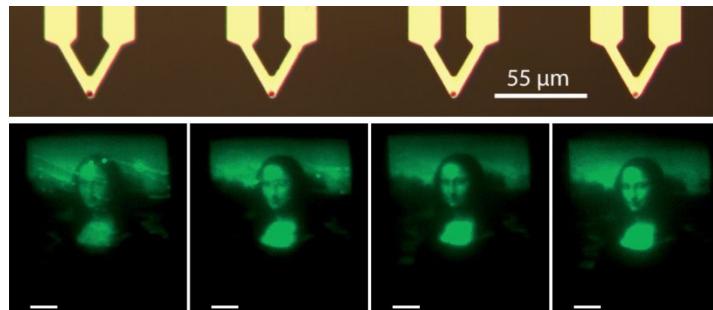
## Simulation / Modeling



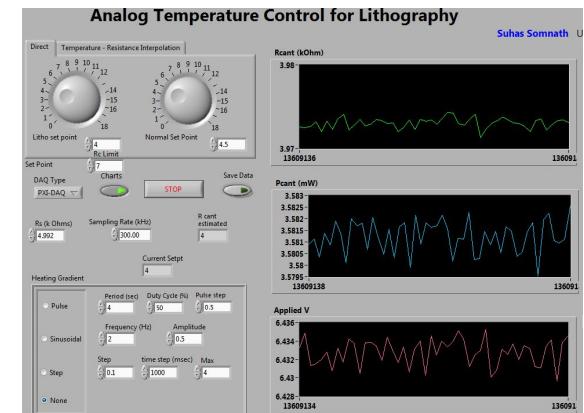
## Instrument Development



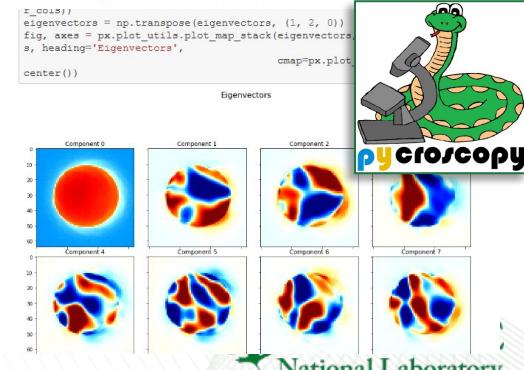
## New measurement & manufacturing techniques



## Instrumentation Software



## Data Analysis Software



# Prototyping, Disseminating, Exploring in Python

## Editor + Python Console

l/a/save.rb

```
1 module Storage
2   module API
3     class Save < Grape::API
4       version 'v1', using: :header, vendor: 'iotcc'
5       resource :save do
6         post "data" "saving data into the database"
```

Python Shell

```
File Edit Options Windows Help
Type "copyright", "credits" or "license()" for more information.

*****
Personal firewall software may warn about the connection IDLE makes to its subprocess using this computer's internal loopback interface. This connection is not visible on any external interface and no data is sent to or received from the Internet.
*****
```

```
IDLE 1.1.3
>>> import sys
>>> sys.frozen
'windows_exe'
>>> sys.version
'2.4.3 (#69, Mar 29 2006, 17:05:34) [MSC v.1310 32 bit (Intel)]'
>>>
```

## IDE

bdd\_example - [~/PycharmProjects/bdd\_example] - .../features/steps/steps.py

```
File Edit View Navigate Code Refactor Run Tools VCS Window Help
Project bdd_example features steps steps.py
  features
    steps
      steps.py
        MasterCardSupport.feature
        VisaSupport.feature
External Libraries
```

...: # TODO: Fix final stuff
...: ls = LightSource(270, 45)
...: # To use a custom hillshading mode, override the built-in shading and pass
...: # it in the rgb colors of the shaded surface calculated from "shade".
...: rgb = ls.shade(z, cmap=cm.gist\_earth, vert\_exag=0.1, blend\_mode='soft')
...: surf = ax.plot\_surface(x, y, z, rstride=1, cstride=1, facecolors=rgb,
...: linewidth=0, antialiased=False, shade=False)
...: plt.show()

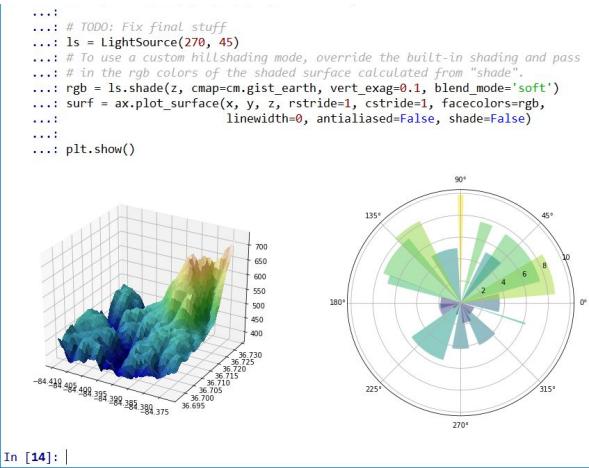
```
@given("I have card (?P<card_type>")
def step_impl(context):
    Add parameter [card_type]
    ...
    context.card_type = card_type
```

Run B features

- Done: 0 of 12 Failed: 4 (0.108 s)
- /home/filippov/bdd\_venv/bin
- Testing started at 2:47 PM
- Traceback (most recent call)
- File "/home/filippov/bdd"
match.run(runner.context)
File "/home/filippov/bdd"
self.func(context, \*arg
TypeError: step\_impl() got

Tests failed (a minute ago)

## iPython console

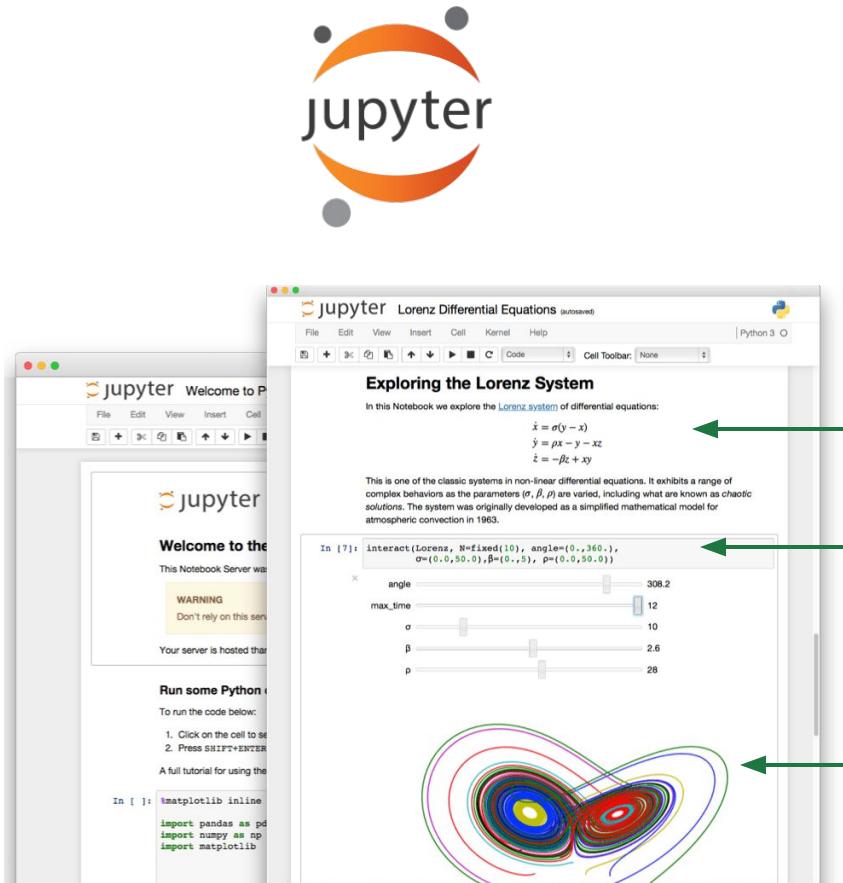


- Fast
- Cumbersome – multiple windows (code, results)
- No markdown, equations
- Hard to follow / teach
- Minimal high-level tools - auto-complete

- Very powerful
- Overkill for prototyping / teaching / learning
- No markdown, equations ...
- Cumbersome – multiple windows\*
- Hard to follow

- Light and simple
- Helpful tools like autocomplete, etc.
- In-line plots and results!
- Single document
- No markdown, equations ...
- Hard to export

# Jupyter Notebooks



- Interactive documents
- Exploratory programming
- Markdown text
  - Equations
- Code
- Images
- Interactive – slice through data, pan, move, rotate ...

# Jupyter Notebooks

<LIVE DEMO>

# Notebooks are great for

- Exploratory data analysis
  - Try different parameters / algorithms quickly
  - Keep notes
- Code prototyping
  - Develop snippets of code that can become part of a package
- Computational narrative
  - Leverage markdown capabilities to explain what, why, how, etc.
  - Cause (code + data) and effect (plots, output)

# Notebooks for Teaching

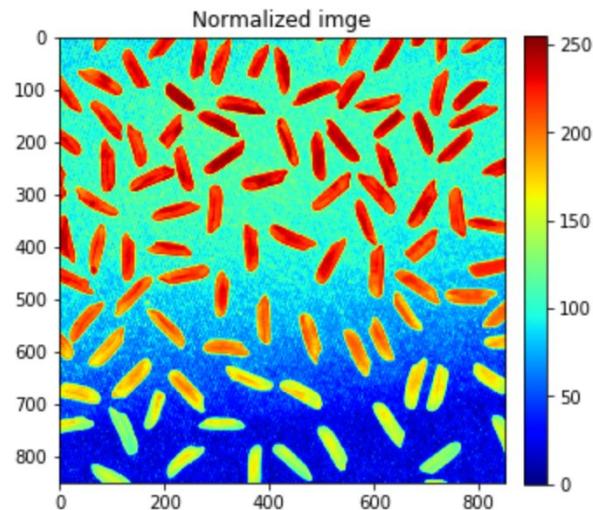


## Thresholding 1

In certain cases, we may be working on images that have two distinct phases, kinds of objects, heights, etc. **that appear different to a reasonable extent** (like the rice image) and we may want to apply separate processing steps on each part of the image.

```
In [25]: # Normalizing to [0, 255]
rice_renormalized = rice_image - rice_image.min()
rice_renormalized = 255 * rice_renormalized / rice_renormalized.max()

fig, axis = plt.subplots(figsize=(5, 5))
im_handle = axis.imshow(rice_renormalized, cmap='jet')
# Add a colorbar
cbar = plt.colorbar(im_handle, ax=axis, fraction=0.046, pad=0.04)
axis.set_title('Normalized img');
```



# Notebooks for Exploration / Prototyping

## Exploratory data analysis

```
In [15]: cropped_image = stem_image[:128, :128]
gaus_sigma = 2.0
gaussian_filtered = filters.gaussian(cropped_image, sigma=gaus_sigma)
fig, axes = plt.subplots(ncols=3, figsize=(12, 4), sharey=True)
axes[0].imshow(cropped_image, cmap='gray')
axes[0].set_title('Original image')
axes[1].imshow(gaussian_filtered, cmap='gray')
axes[1].set_title('Gaussian filtered, Sigma = {}'.format(gaus_sigma))
axes[2].imshow(cropped_image - gaussian_filtered, cmap='gray')
axes[2].set_title('Removed noise')
fig.tight_layout()
```

Original image      Gaussian filtered, Sigma = 2.0      Removed noise

Visualize the effect of filtering in the frequency domain:

The center of the image (low-frequency) remains more-or-less the same but the data outside the center (higher spatial frequencies) is drastically reduced in magnitude thereby resulting in a cleaner image. All the noise from the higher frequencies is visible in the removed noise

Thus, the Gaussian filter is essentially a low-pass filter

## Prototyping code

```
def test_index_view_with_a_future_question(self):
    """
    Questions with a pub_date in the future should not be displayed on
    the index page.
    """
    create_question(question_text="Future question.", days=30)
    response = self.client.get(reverse('polls:index'))
    self.assertContains(response, "No polls are available.",
                       status_code=200)
    self.assertQuerysetEqual(response.context['latest_question_list'], [])
```

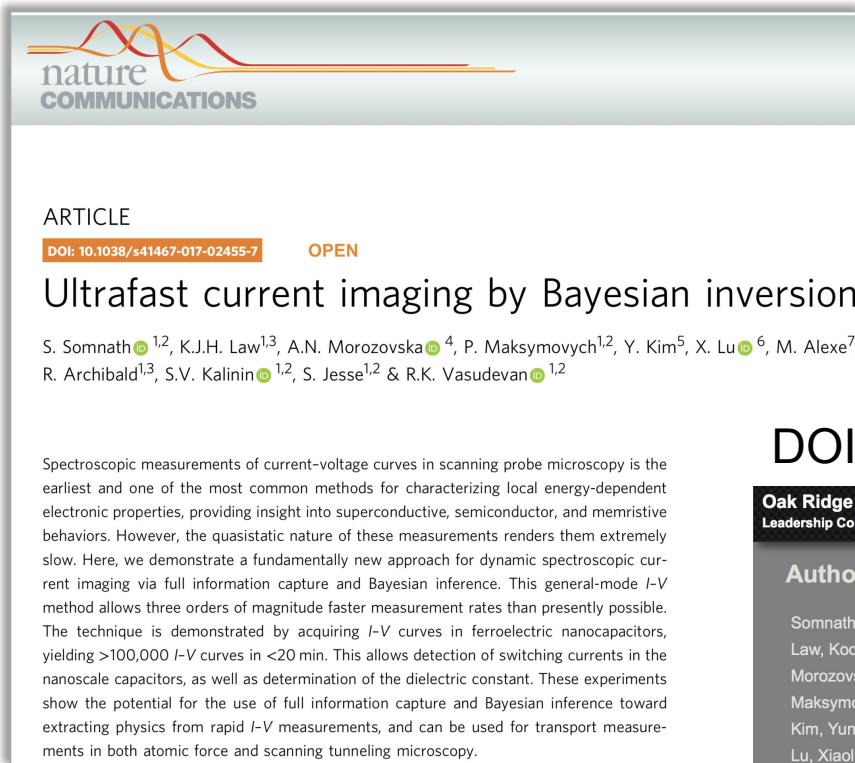
- Try parameters / algorithms quickly
- Keep notes

Develop snippets of code that can become part of a package

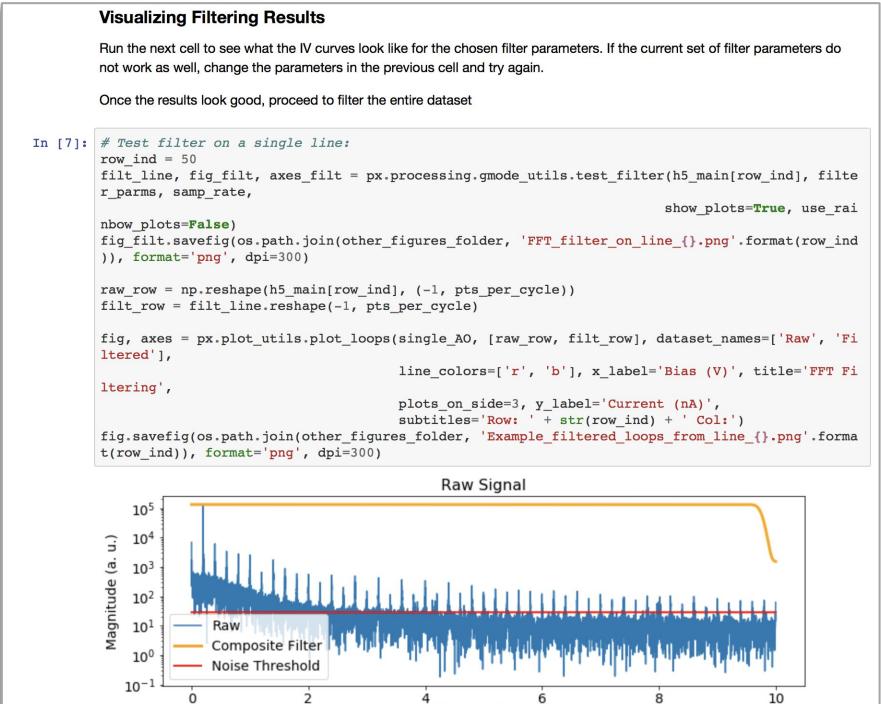
# Towards Open & Reproducible Analytics / Science

Academic papers accompanied with:

- Jupyter notebooks that show all analysis (raw data → figures)
- Data with DOI number
- Software environment → containers



Jupyter notebook associated with paper



DOI associated with data (raw → paper figures)

Oak Ridge National Laboratory 10.13139/OLCF/1410993

Download

## Authors

Somnath, Suhas	somnaths@ornl.gov
Law, Kody	lawkj@ornl.gov
Morozovska, Anna	anna.n.morozovska@gmail.com
Maksymovych, Petro	maksymovychp@ornl.gov
Kim, Yunseok	ykim943@gmail.com
Lu, Xiaoli	xliu@xidian.edu.cn
Alexe, Marin	M.Alexe@warwick.ac.uk

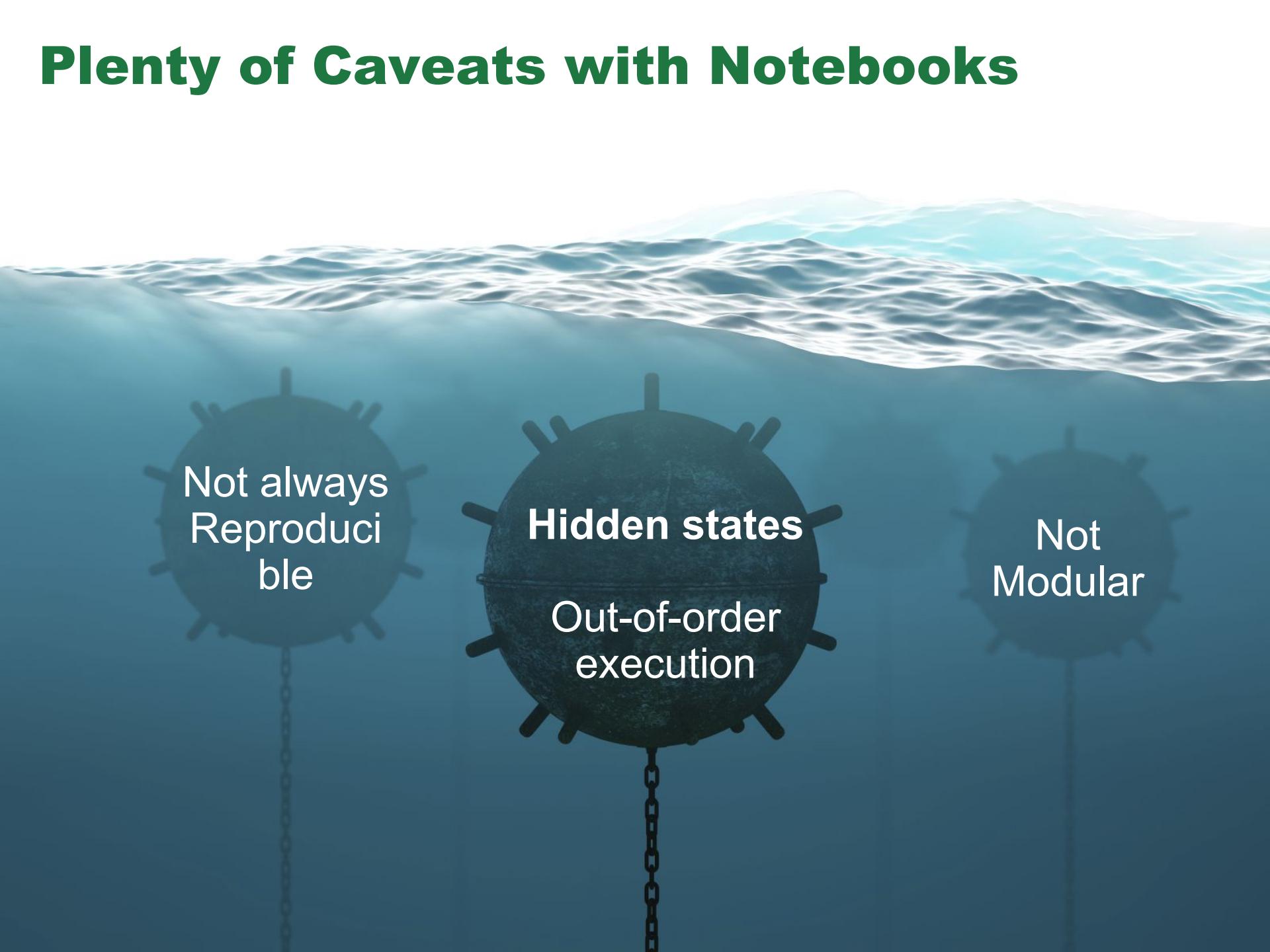
# Notebooks for Narratives

- Example notebooks that supplement academic articles
  - [https://pycroscopy.github.io/pycroscopy/papers\\_conferences.html#journal-papers-using-pycroscopy](https://pycroscopy.github.io/pycroscopy/papers_conferences.html#journal-papers-using-pycroscopy)

# Is it all smooth sailing with Notebooks?



# Plenty of Caveats with Notebooks



Not always  
Reproduc-  
ible

**Hidden states**  
**Out-of-order  
execution**

Not  
Modular

Notebooks encourage



**BAD  
HABITS**

A large, glowing red neon sign with the words "BAD HABITS" in a bold, sans-serif font. The letters are outlined in bright red neon light, which is reflected on the surface below. The sign is set against a dark, textured background of a brick wall.

# Hidden State

jupyter madness

The screenshot shows a Jupyter Notebook interface with the following code execution history:

- In [1]: `def f(x): return x + 2`
- In [3]: `y = f(2)`
- In [4]: `y == 4`
- Out[4]: `False`
- In [5]: `print(y)`

The output for In [5] is the number 5.

if you look at  
the numbers,  
it's clear that  
*something else*  
was run  
between the  
first two lines



@joelgrus #jupytercon

# Questionable Reproducibility

```
In [1]: import torch
import torch.nn as nn
from torch.autograd import Variable

import pandas as pd
import random
import string
import numpy as np

import sys, os

import torch.utils.data as data

os.environ["CUDA_VISIBLE_DEVICES"] = '0'

all_characters = string.printable
number_of_characters = len(all_characters)

artists = [
'ABBA',
'Ace Of Base',
'Aerosmith',
'Avril Lavigne',
'Backstreet Boys',
'Bob Marley',
'Bon Jovi',
'Britney Spears',
'Bruno Mars',
'Coldplay',
'Def Leppard',
'Depeche Mode',
'Ed Sheeran',
```

pytorch 0.3? 0.4? 0.4.1?

I suppose I could guess  
based on the commit date

Will it work on my machine?  
• Hardware  
• Operating system  
• Python version?

@joelgrus #jupytercon

# Questionable Reusability

```
In [1]: import torch
import torch.nn as nn
from torch.autograd import Variable

import pandas as pd
import random
import string
import numpy as np

import sys, os

import torch.utils.data as data

os.environ["CUDA_VISIBLE_DEVICES"] = '0'

all_characters = string.printable
number_of_characters = len(all_characters)

artists = [
'ABBA',
'Ace Of Base',
'Aerosmith',
'Avril Lavigne',
'Backstreet Boys',
'Bob Marley',
'Bon Jovi',
'Britney Spears',
'Bruno Mars',
'Coldplay',
'Def Leppard',
'Depeche Mode',
'Ed Sheeran',
```

Why are you doing things  
this way?  
Instructions?

Can I use my data instead?

What if I only want to use  
some portions?  
Can't import (easily)

# Typical excuse...



[REDACTED]  
@ [REDACTED]

Follow



Data science code doesn't need to  
follow the rules of good software  
engineering, because data science is  
not about creating software but about  
experimenting with building prototypes  
of models. ➤ Great tip from  
@jeremyphoward

4:18 AM - 2 May 2018

4 Retweets 26 Likes



@joelgrus  
#jupytercon

# Other challenges

- Importing notebooks in other notebooks quickly becomes ugly
- Same with notebooks sharing state with others
- Git diffs are filled with JSON changes / cruft

# Notebooks are not good for

- Teaching / following code
  - Out-of-order execution of code
  - Deleted / updated cells
- Software development
  - Not modular / reusable
  - Challenging / tedious to unit test
  - Importing notebooks quickly becomes ugly
  - Git often catches JSON changes rather than code / results
  - Collaboration is tedious
- Reproducibility
  - Dependencies not captured

# Not the end of the world for Jupyter

- Follow best practices
- Use notebooks as they were intended

## Ten Simple Rules for Reproducible Research in Jupyter Notebooks

Adam Rule, Amanda Birmingham, Cristal Zuniga, Ilkay Altintas,  
Shih-Cheng Huang, Rob Knight, Niema Moshiri, Mai H. Nguyen, Sara  
Brin Rosenthal, Fernando Pérez, Peter W. Rose

arXiv preprint arXiv:1810.08055. 2018 Oct 13.

<https://github.com/jupyter-guide/ten-rules-jupyter>

# 1. Tell a Story for an Audience

- Leverage markdown to tell a compelling story
- Why are you doing things this way?
- What do the results mean and why?



## Saltelli's algorithm step by step

### Step 1: sample matrix creation

For Saltellis Algorithm we need to create two different sample matrices  $A, B$  each of the size  $M \times P$ :

$$\mathbf{A} = \begin{bmatrix} z_1^{(A,1)} & \dots & z_i^{(A,1)} & \dots & z_P^{(A,1)} \\ \vdots & & & & \vdots \\ z_i^{(A,M)} & \dots & z_i^{(A,M)} & \dots & z_P^{(A,M)} \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} z_1^{(B,1)} & \dots & z_i^{(B,1)} & \dots & z_P^{(B,1)} \\ \vdots & & & & \vdots \\ z_i^{(B,M)} & \dots & z_i^{(B,M)} & \dots & z_P^{(B,M)} \end{bmatrix}.$$

In addition we create  $P$  additional matrices  $C_i$  of the size  $M \times P$  compound of matrix  $A$  and matrix  $B$ . In a matrix  $C_i$  all columns will be have the same values as the  $B$  matrix, except the  $i$ -th column, which will have the values of  $A$ :

$$\mathbf{C}_i = \begin{bmatrix} z_1^{(B,1)} & \dots & z_i^{(A,1)} & \dots & z_P^{(B,1)} \\ \vdots & & & & \vdots \\ z_i^{(B,M)} & \dots & z_i^{(A,M)} & \dots & z_P^{(B,M)} \end{bmatrix}$$

This was implemented in the method: `A, B, C = generate_sample_matrices_mc(number_of_samples, number_of_parameters, joint_distribution, sample_method)`.

```
In [8]: # sample matrices
def generate_sample_matrices_mc(Ns, number_of_parameters, jpdf, sample_method='R'):

    xtot = jpdf.sample(2*Ns, sample_method).transpose()
```

## 2. Document the Process, not just Results

- Document before you forget!
- Name, contact, date(s) ...
- Both successful and failed experiments
- Where you got a block of code from
- Avoid manipulating figures etc. using external software

### 3. Add Divisions to Make Steps Clear

- Avoid one-line and very long cells
- Each cell -> one meaningful step
- Label cells using markdown above
- Split long notebooks into multiple notebooks

# 4. Modularize Code

- Avoid multiple copies of code with alterations
  - Wrap into function
  - Turn into module / package

# 5. Record Dependencies

- Explicitly print out dependencies using extensions
  - Watermark (<https://github.com/rasbt/watermark>)
- Write out requirements.txt
  - Easier to deploy containers. E.g. – Binder (more later)

# 6. Use Version Control

- Track modifications, source of bugs, etc.
- Notebook-specific diffing tool, e.g. – nbdime
  - Understands notebook structure
  - Presents meaningful diffs – minimal JSON metadata
  - <https://github.com/jupyter/nbdime>

# 7. Build a Pipeline (Experimentation -> Production)

- Generalize notebooks for reusability
  - Accept different input data
  - Different versions of packages
- Cleanly organized. From the top:
  - All imports
  - Key variable declarations
  - Preparatory steps – e.g. – data cleaning
- Re-run entire notebook in one shot to catch bugs
- Parametrize notebooks (e.g. – Papermill)
  - Run via command line / interactively
- Link with Makefile – noninteractive execution
- Write tests

# 8. Share and Explain your Data

- Notebook without data can be useless
- Prepare data
  - Anonymize sensitive data
  - Abridge very large datasets
  - Open format
- Share data (or portion of) publicly
  - Figshare (<https://figshare.com/>)
  - Zenodo (<https://zenodo.org/>)
  - Get Digital Object Identifier (DOI) to uniquely identify data

## 9. Enable Notebooks to be Read, Run, Explored

- Share notebook on GitHub or similar
- Add license info (reuse)
- Add supplements:
  - ReadMe
  - Static HTML / PDF exports of completed notebooks
    - What it will look like if everything worked perfectly
  - Alternatively view through NbViewer
  - Upload / link to related academic articles
- Zero-install environments (containers)
  - Binder
  - Docker

## 10. Contribute to Reproducible & Open Research

- Make reproducibility part of the standard process
  - Not an afterthought
- Crowdsource testing / readability
- Consider:
  - Open-sourcing code
  - Contributing to other projects

# Binder

1

## Enter your repository information

Provide in the above form a URL or a GitHub repository that contains Jupyter notebooks, as well as a branch, tag, or commit hash. Launch will build your Binder repository. If you specify a path to a notebook file, the notebook will be opened in your browser after building.

2

## We build a Docker image of your repository

Binder will search for a dependency file, such as requirements.txt or environment.yml, in the repository's root directory ([more details on more complex dependencies in documentation](#)). The dependency files will be used to build a Docker image. If an image has already been built for the given repository, it will not be rebuilt. If a new commit has been made, the image will automatically be rebuilt.

3

## Interact with your notebooks in a live environment!

A [JupyterHub](#) server will host your repository's contents. We offer you a reusable link and badge to your live repository that you can easily share with others.

# Teaching with Notebooks before Binder

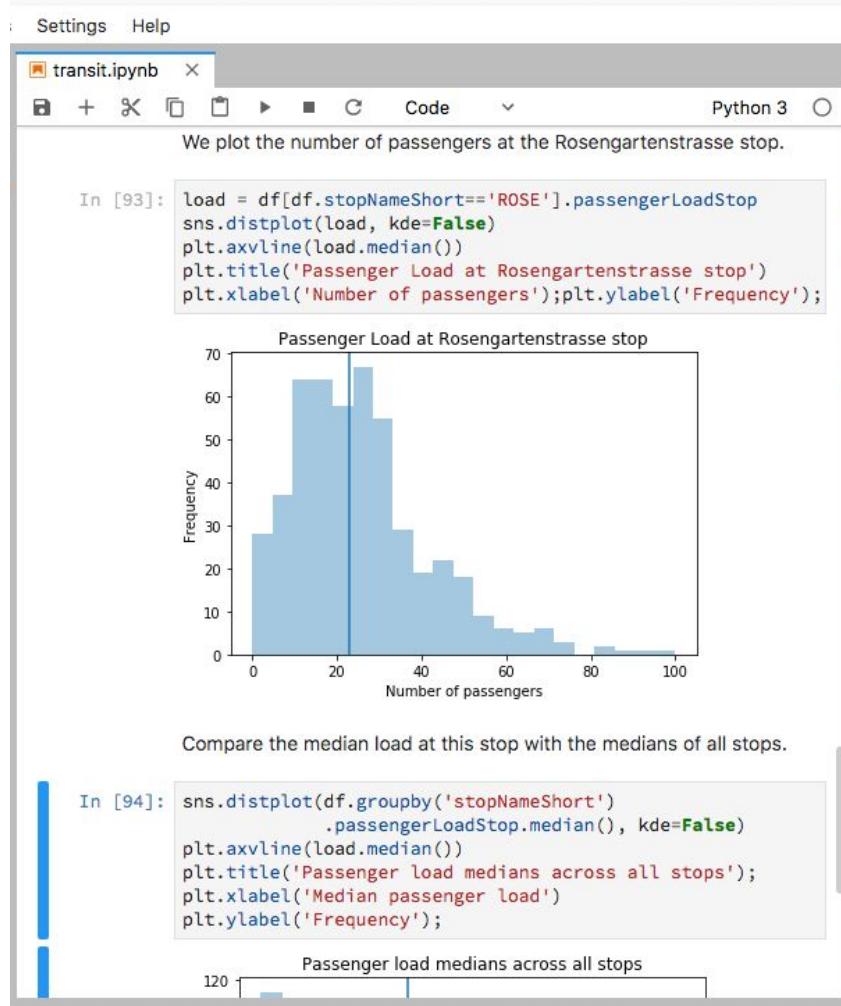
- Before the workshop:
  - Write clean notebooks
  - Try notebooks on multiple configurations manually!
  - Write good install instructions
  - Hope attendees will install requirements
- At the workshop:
  - 50 – 75 % attendees have not installed software
  - Slow conference Wi-Fi!
  - Someone is running some untested configuration
    - Struggle with installation. Repeat.
  - Schedule off by an hour!

# Teaching with Notebooks using Binder

- Before the workshop:
  - Write clean notebooks
  - ~~Try notebooks on multiple configurations manually!~~
  - Configure Binder with requirements.txt etc.
  - ~~Hope attendees will install requirements~~
- At the workshop:
  - Point everyone towards the Binder link
  - Everyone is up and running in < 5 mins
- Other options – Azure Notebooks, Google Collaboratory, ....

# Jupyter Lab - Notebook on steroids

Live-updating visualization panel



Customizable  
tab areas

Terminals, scripts, notebooks,...

Native file-editors with live-preview

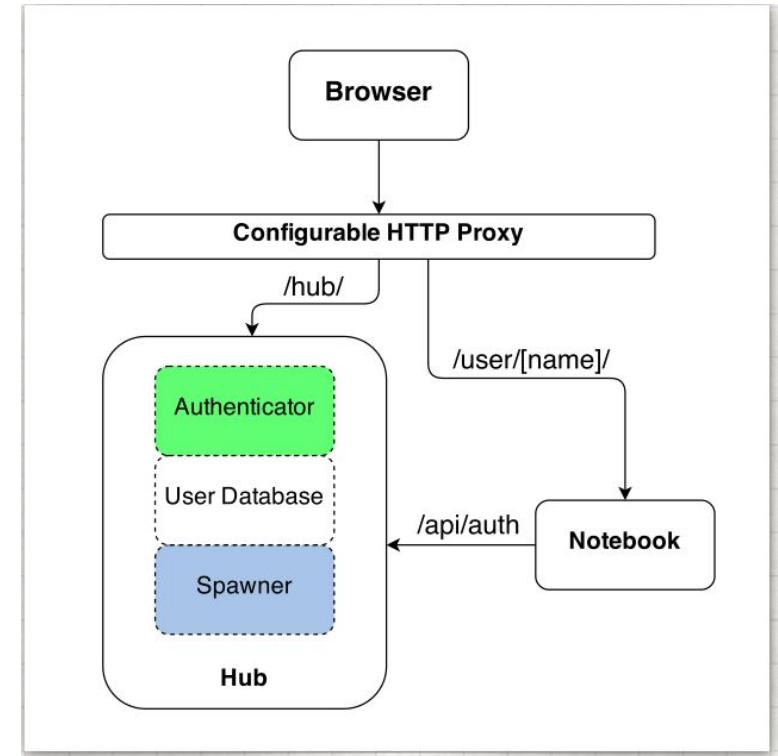


# Jupyter Lab - Notebook on steroids

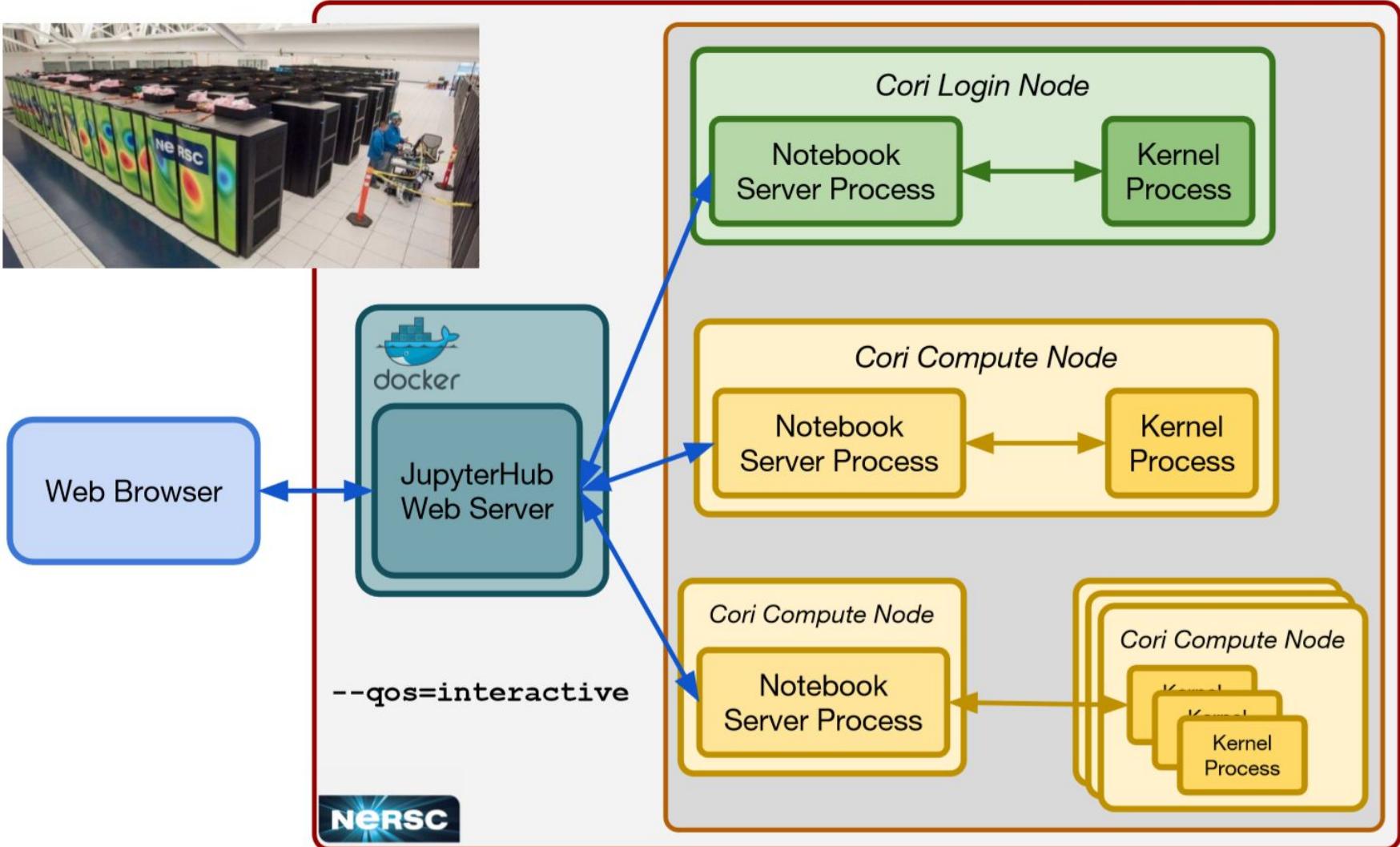
- Drag-and-drop to reorder notebook cells and copy them between notebooks.
- Run code blocks interactively from text files (.py, .R, .md, .tex, etc.).
- Link a code console to a notebook kernel to explore code interactively without cluttering up the notebook with temporary scratch work.
- [https://jupyterlab.readthedocs.io/en/stable/getting\\_started/overview.html](https://jupyterlab.readthedocs.io/en/stable/getting_started/overview.html)

# Jupyter Hub -> Multi-user notebook server

- Standardized + managed environment
  - Users don't have to install anything
  - Access server via simple webpage
- Standardized authentication (e.g. – Google)
- Optionally – make compute + data available to all
- Great for classes, labs



# Jupyter + Cluster / HPC



<https://www.exascaleproject.org/event/jupyter/>

# Summary

- Jupyter notebooks - interactive web documents with code, markdown, results, plots, interactive widgets
- Great for:
  - Computational narrative tool
  - Exploratory programming
- Bad for:
  - Out-of-order / deleted cells - hidden states
  - Does not address reproducibility
  - Poor modularity and reusability
- 10 best-practices
- Beyond Notebooks:
  - Jupyter Lab - customizable panels, code editors, etc.
  - Jupyter Hub - Multi-user notebooks
  - Links to high-performance computation

# Thank you!

[somnaths@ornl.gov](mailto:somnaths@ornl.gov)

[https://docs.google.com/presentation/d/1hZjKeu2X6v  
orHZxyQE0PDMI8bts3iNXxf1e7DG8X3h4/edit?usp=s  
haring](https://docs.google.com/presentation/d/1hZjKeu2X6v<br/>orHZxyQE0PDMI8bts3iNXxf1e7DG8X3h4/edit?usp=s<br/>haring)