# Modeling stopping power with time-dependent density functional theory



786,432 CPUs, 10 PFLOPS supercomputer
Argonne Leadership Computing Facility

**Andre Schleife, UIUC**

Hydrogen in Gold
(v=2.0 at. u.)

16,000 CPU-hours per simulation

Electron density (at. u.)

-0.02000
-0.01000
-0.005000
0.005000
0.01000
0.02000

Max: 0.000
Min: 0.000

Time step: 0.01 at. u.; Total time: 0.3396 fs

Real-time simulation (Ehrenfest MD), 256 gold atoms (4352 valence el.), plane-wave cutoff: 130 Ry

A. Schleife, E. Draeger, V. Anisimov, A. Correa, Y. Kanai (2013)

# Jupyter notebooks enable rapid iteration/results

```
In [35]:  @python_app
          def get_stopping_power(lattice_vector, traj_computer):
              return traj_computer.compute_stopping_power([0,0.8,0.85], lattice_vector, 1.0, abserr=0.001,
                                                          hit_threshold=2.5, full_output=1)
```

```
In [37]:  stopping_power_results = []
          for d in tqdm(dirs, desc='Submitting'):
              stopping_power_results.append(get_stopping_power(d, traj_computer))
```

Submitting  ████████████████████  100% 24/24 [00:00<00:00, 166.06it/s]

```
In [38]:  stopping_power_results = [s.result() for s in tqdm(stopping_power_results, desc='Waiting')]
```

Waiting  ████████████████████  100% 24/24 [18:47:19<00:00, 2818.33s/it]

```
In [62]:  ax = plt.subplot(111, projection='polar')
          fig = plt.gcf()

          ax.plot(angles + angles[:1], stopping_power + stopping_power[:1], marker='o')

          # Plot the 'channel value'
          ax.plot(np.linspace(0, 2*np.pi, 100), [ml_stopping_new,]*100)
          ax.set_rmax(0.25)
          ax.set_rmin(0.2)#min(stopping_power) * 0.99)

          fig.set_size_inches(4, 4)
```
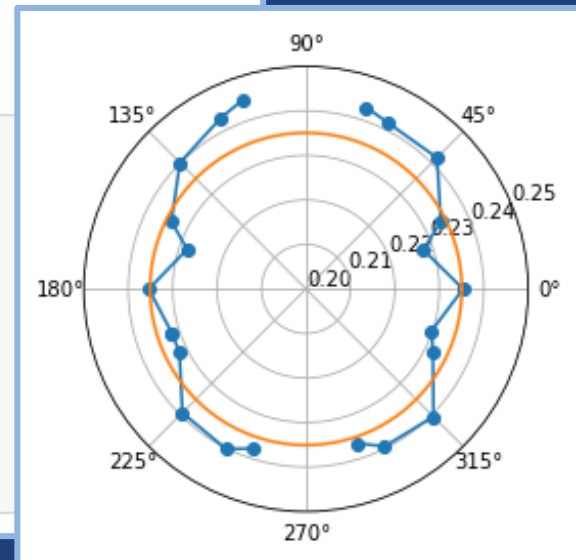
Logan Ward

# But the data are big, distributed…
## …and our science is collaborative

**MATERIALS DATA FACILITY**
materialsdatafacility.org

**3.2M materials data**

**1** Query

**4** Share

jupyter

**PETREL**
petrel.alcf.anl.gov

**2** Transfer

**3** Learn

**2PB, 80Gbps store**

globus connect

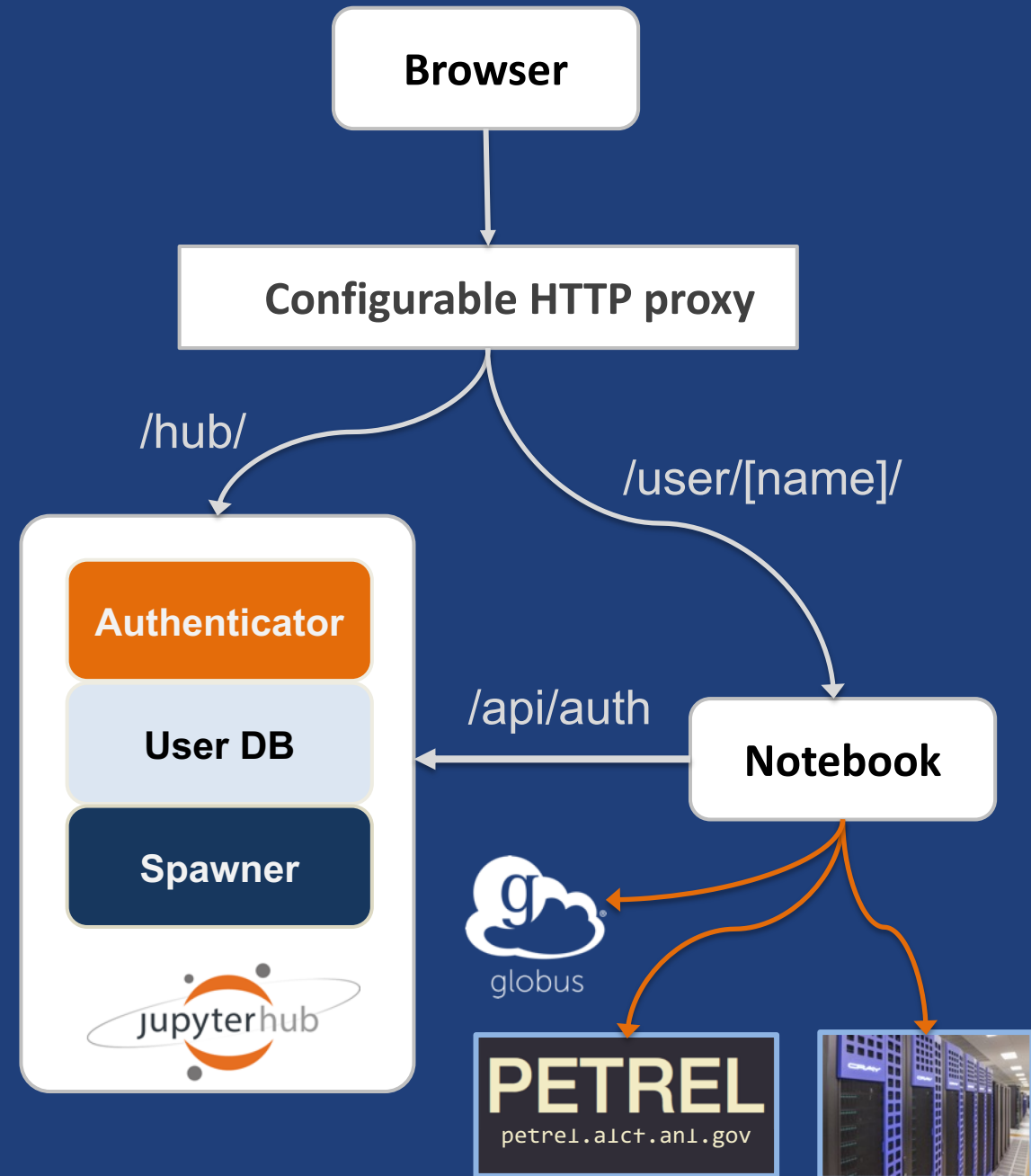globus connect

**Cooley: 290 TFLOPS**

**CRAY**

Need multi-credential, multi-service authentication and data management

# JupyterHub

- **Multi-user hub**

- **Manages multiple instances of Jupyter notebook server**
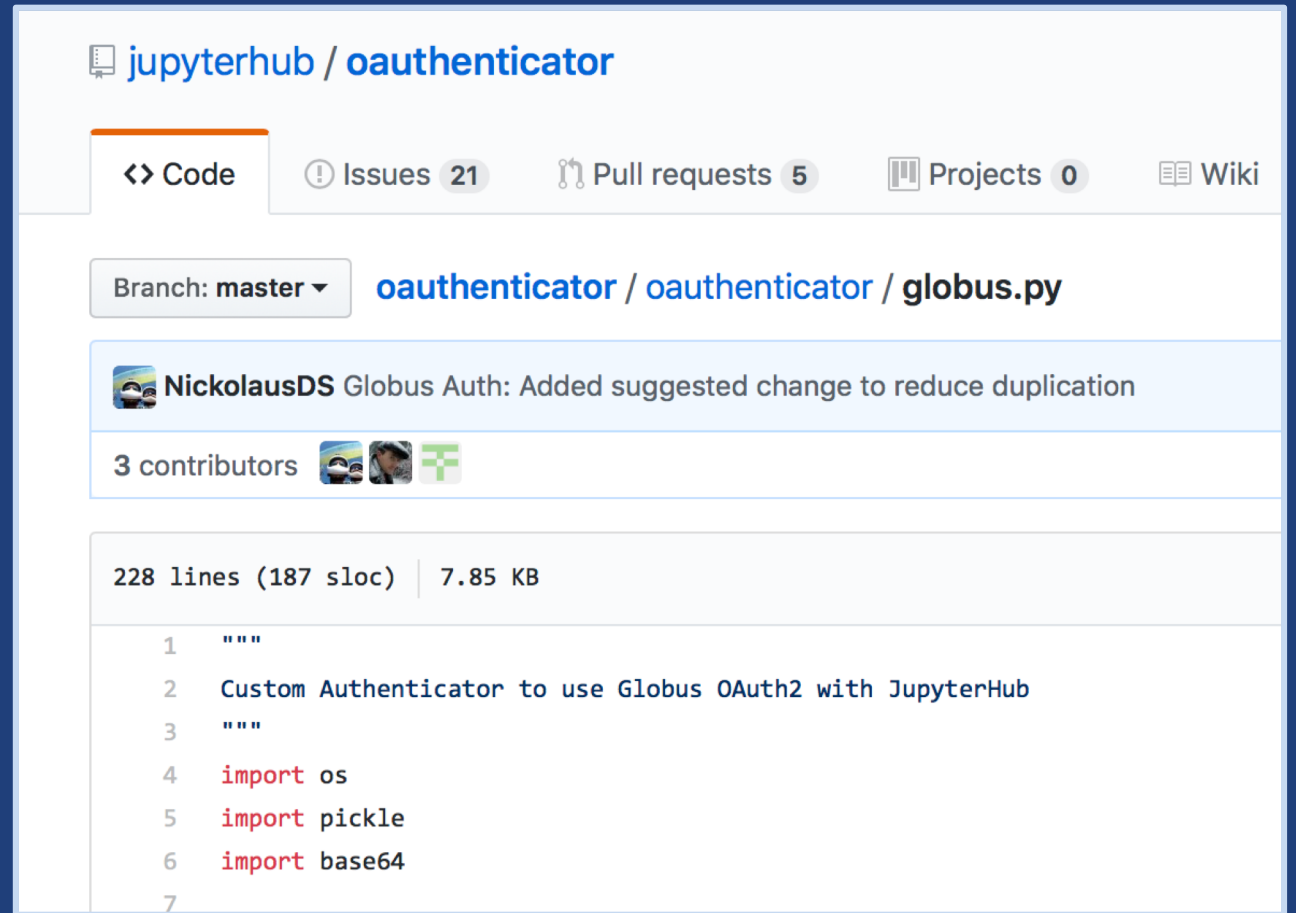
- **Configurable HTTP proxy**

## Goal: Liberate the notebook!

- **Tokens for remote services**
- **APIs for remote actions, e.g. data management via Globus service**

# Securing JupyterHub with Globus Auth plugin

- **Existing OAuth framework**

- **Can restrict IdP**

- **Custom scopes**

- **Tokens passed into notebook environment**



**github.com/jupyterhub/oauthenticator**

# Securing JupyterHub with Globus Auth

Visit https://developers.globus.org/ to set up your app. Ensure *Native App* is unchecked and make sure the callback URL looks like:

```
https://[your-host]/hub/oauth_callback
```

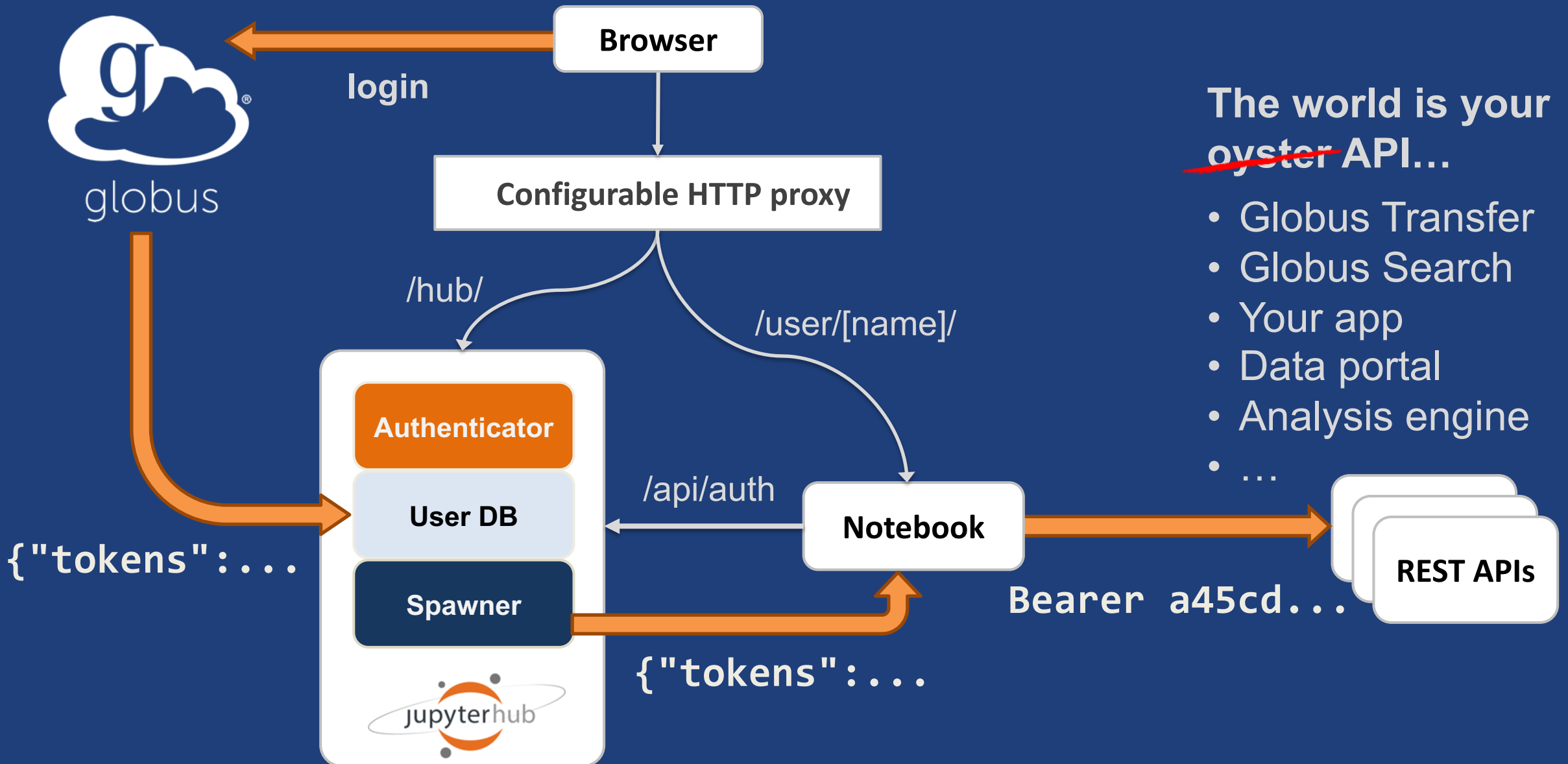Set scopes for authorization and transfer. The defaults include:

```
openid profile urn:globus:auth:scope:transfer.api.globus.org:all
```

Set the above settings in your `jupyterhub_config`:

```python
# Tell JupyterHub to create system accounts
from oauthenticator.globus import LocalGlobusOAuthenticator
c.JupyterHub.authenticator_class = LocalGlobusOAuthenticator
c.LocalGlobusOAuthenticator.enable_auth_state = True
c.LocalGlobusOAuthenticator.oauth_callback_url = 'https://[your-host]/hub/oauth_callback'
c.LocalGlobusOAuthenticator.client_id = '[your app client id]'
c.LocalGlobusOAuthenticator.client_secret = '[your app client secret]'
```

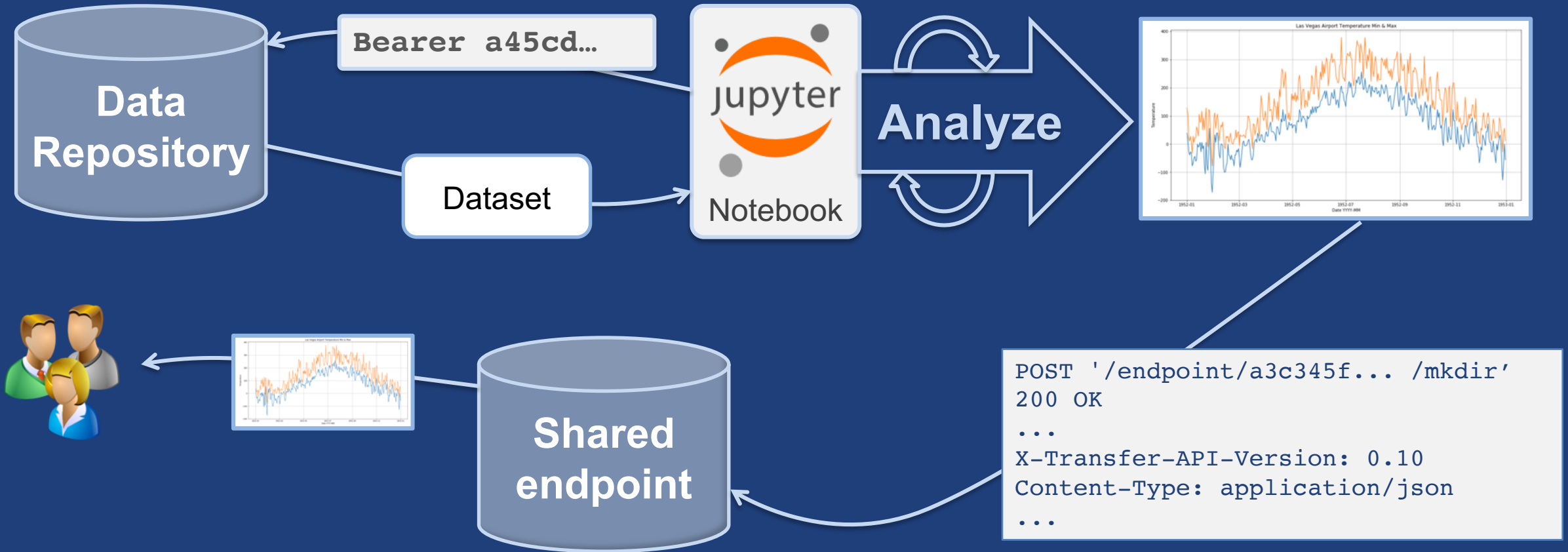**github.com/jupyterhub/oauthenticator#globus-setup**

# Tokens in Jupyter notebooks

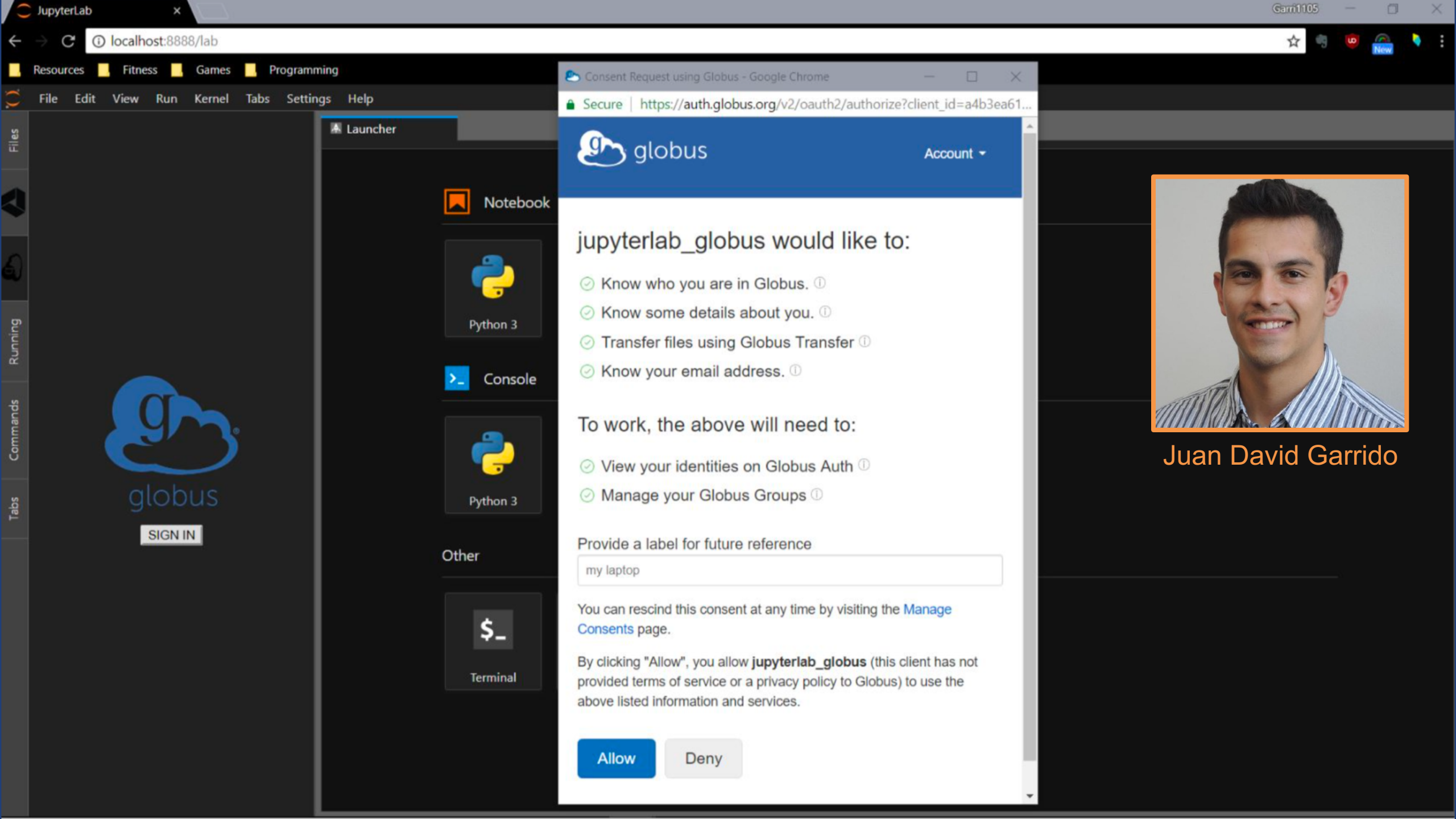# Automated data analysis/results distribution

# Experiment with the demo notebook

- **Login into our JupyterHub*: jupyter.demo.globus.org**

- **Launch (spawn) a notebook server; get tokens**

- **Access Globus APIs; download some data**

- **"Analyze" data (generate plot)**

- **PUT results (graph) on an HTTPS endpoint**

*zero-to-jupyterhub.readthedocs.io

# Futures…

Browse data on local storage

localhost:8888/lab

Resources    Fitness    Games    Programming

File  Edit  View  Run  Kernel  Tabs  Settings  Help

Launcher

File Manager

Globus Tutorial Endpoint 1

/~/

select none    Sort

Transfer

New folder

Rename

Delete

Notebook

Python 3

Console

Python 3

Transfer data to local storage

Other

Terminal    Text Editor

localhost:8888/lab

Resources  Fitness  Games  Programming

File  Edit  View  Run  Kernel  Tabs  Settings  Help

Launcher ✕ | sensitivity_intro ✕

Markdown ∨

Python 3 ○

Search

Select Index ▾

Select Index
Ramses
MDF
Kasthuri

Search remote databases

# A practical introduction to sensitivity analysis

**Leif Rune Hellevik**, Department of Structural Engineering, NTNU

Date: **Jul 13, 2018**

```
# ipython magic
%matplotlib
%auto
```

In [2]:
```
# plot conf
import matpl
import matpl
plt.style.us
# import sea
matplotlib.rcParams['lines.linewidth'] = 3
fig_width, fig_height = (7.0,5.0)
matplotlib.rcParams['figure.figsize'] = (fig_width, fig_height)
```

In [3]:
```
# import modules
import numpy as np
from numpy import linalg as LA
import chaospy as cp
from sensitivity_examples_nonlinear import generate_distributions
from monte_carlo import generate_sample_matrices_mc
from monte_carlo import calculate_sensitivity_indices_mc
import pandas as pd
from operator import index
```

## Introduction

This practical introduction to sensitivity analysis is based on the presentation and examples found in [saltelli_global_2008]. To give the reader an even better hands on experience of the topic, we have integrated the computations in a python notebook format.

Many sensitivity analyses reported in the literature are based on derivatives at set point or point of interest. Indeed such approaches are based on the fact that the derivative of $\partial Y_i/\partial X_j$ of quantity of interest $Y_i$ as a function of an input variable $X_j$ can be thought of as the mathematical definition of the sensitivity of $Y_i$ versus $X_j$.

File   Edit   View   Run   Kernel   Tabs   Settings   Help

Launcher ×    sensitivity_intro ×

Markdown ∨                                        Python 3 ○

Files

Search

MDF ▼

*

∨ Filters

alloy_db_test_5748_v1
*Material:* Al2Au1
*Elements:* Al,Au
*Files:* 3

ohmic_si_c_contacts_v1
*Material:* Ti6Au2C4
*Elements:* C,Ti,Au
*Files:* 4

ohmic_si_c_contacts_v1
*Material:* Ti6Ir2C4
*Elements:* Ir,C,Ti
*Files:* 2

ohmic_si_c_contacts_v1
*Material:* Ti8Au2C6
*Elements:* C,Ti,Au
*Files:* 2

ohmic_si_c_contacts_v1
*Material:* Ti9Au6C6
*Elements:* C,Ti,Au
*Files:* 2

ohmic_si_c_contacts_v1
*Material:* Ti8Ir2C6
*Elements:* Ir,C,Ti
*Files:* 2

ohmic_si_c_contacts_v1
*Material:* Ti8Au2C6
*Elements:* C,Ti,Au
*Files:* 2

ab_initio_solute_database_v3
*Material:* Mg96

# A practical introduction to sensitivity analysis

**Leif Rune Hellevik**, Department of Structural Engineering, NTNU

**Date:** **Jul 13, 2018**

```
In [1]:   # ipython magic
          %matplotlib notebook
          %load_ext autoreload
          %autoreload 2
```
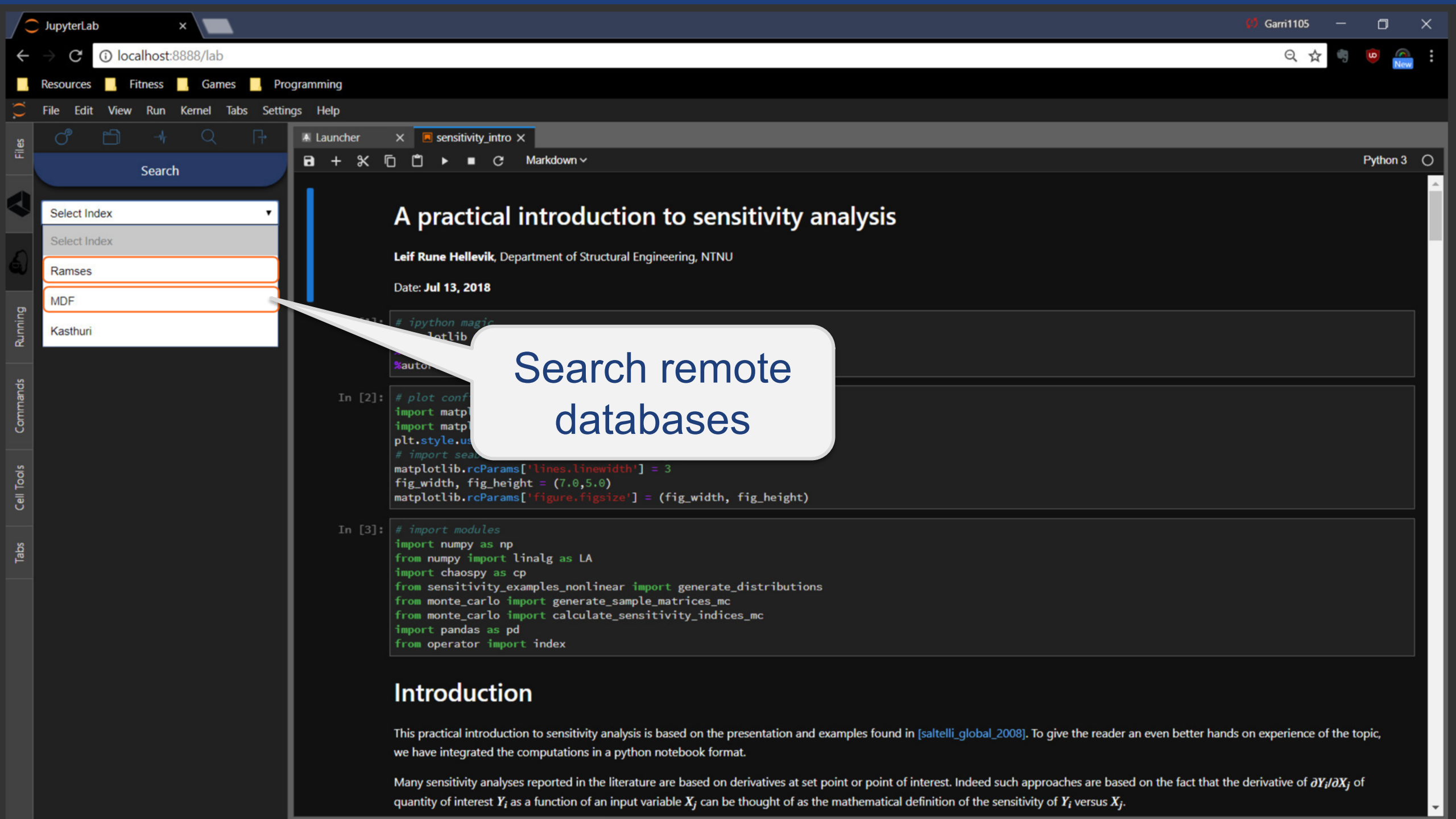
```
In [2]:   # plot configuration
          import matplotlib
          import matplotlib.pyplot as plt
          plt.style.use("ggplot")
          # import seaborn as sns # sets another style
          matplotlib.rcParams['lines.linewidth'] = 3
          fig_width, fig_height = (7.0,5.0)
          matplotlib...
```

```
In [3]:   #
          import
          from numpy
          import ch
          from sens
          from mont
          from monte
          import pandas as pd
          from operator import index
```

Select materials data

## Introduction

This practical introduction to sensitivity analysis is based on the presentation and examples found in [saltelli_global_2008]. To give the reader an even better hands on experience of the topic, we have integrated the computations in a python notebook format.

Many sensitivity analyses reported in the literature are based on derivatives at set point or point of interest. Indeed such approaches are based on the fact that the derivative of $\partial Y_i/\partial X_j$ of quantity of interest $Y_i$ as a function of an input variable $X_j$ can be thought of as the mathematical definition of the sensitivity of $Y_i$ versus $X_j$.

JupyterLab ✕ | Garri1105

localhost:8888/lab

Resources | Fitness | Games | Programming

File Edit View Run Kernel Tabs Settings Help

Files

Search

MDF ▼

◁ Overview

ohmic_si_c_contacts_v1
▼ Object
  ▼ crystal_structure: Object
    number_of_atoms: 16
    space_group_number: 194
    volume: 185.2659999999991
  ▼ dft: Object
    converged: true
    cutoff_energy: 400
    exchange_correlation_functional: "PA
  ▼ files: Array[2]
    ▼ 0: Object
      data_type: "ASCII text"
      filename: "CONTCAR_cutoff-400eV"
      globus: "globus://e38ee745-6d04-11
      length: 2074
      mime_type: "text/plain"
      sha512: "376136fe3e7d772e848fb4665
      url: "https://e38ee745-6d04-11e5-b
    ▶ 1: Object
  ▼ material: Object
    composition: "Ti8Ir2C6"
    ▶ elements: Array[3]
  ▼ mdf: Object
    ingest_date: "2018-03-28T15:22:18.99
    mdf_id: "5abbb32b34a2263dfa3e3352"
    parent_id: "5abbb32a34a2263dfa3e321c"
    resource_type: "record"
    scroll_id: 310
    source_name: "ohmic_si_c_contacts_v1"

Transfer

■ sensitivity_intro ✕

💾 ➕ ✂ 📋 📋 ▶ ■ C    Markdown ⌄                                    Python 3 ○

# A practical introduction to sensitivity analysis

**Leif Rune Hellevik**, Department of Structural Engineering, NTNU

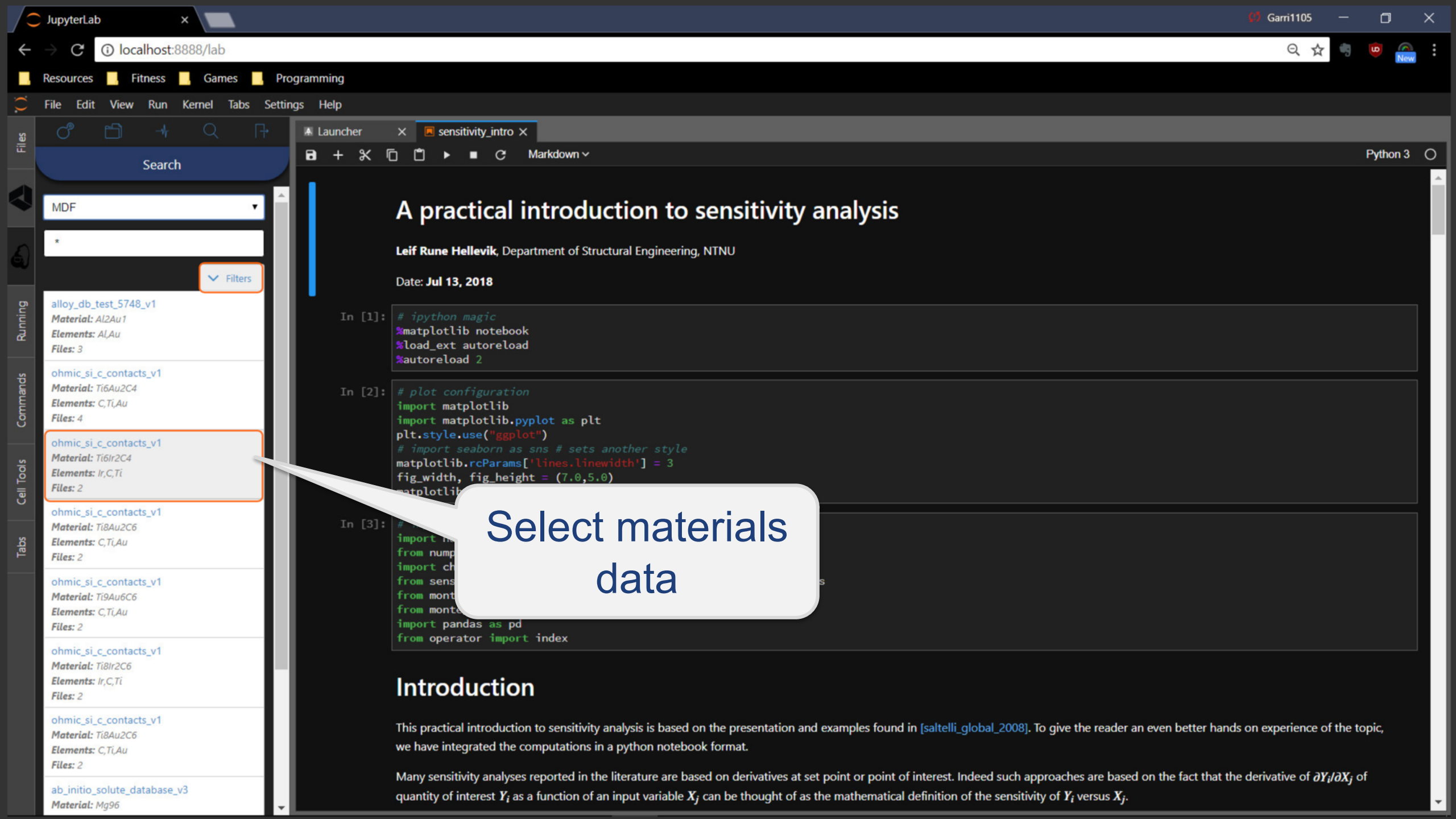Date: **Jul 13, 2018**

In [1]:
```
# ipython magic
%matplotlib notebook
%load_ext autoreload
%autoreload 2
```
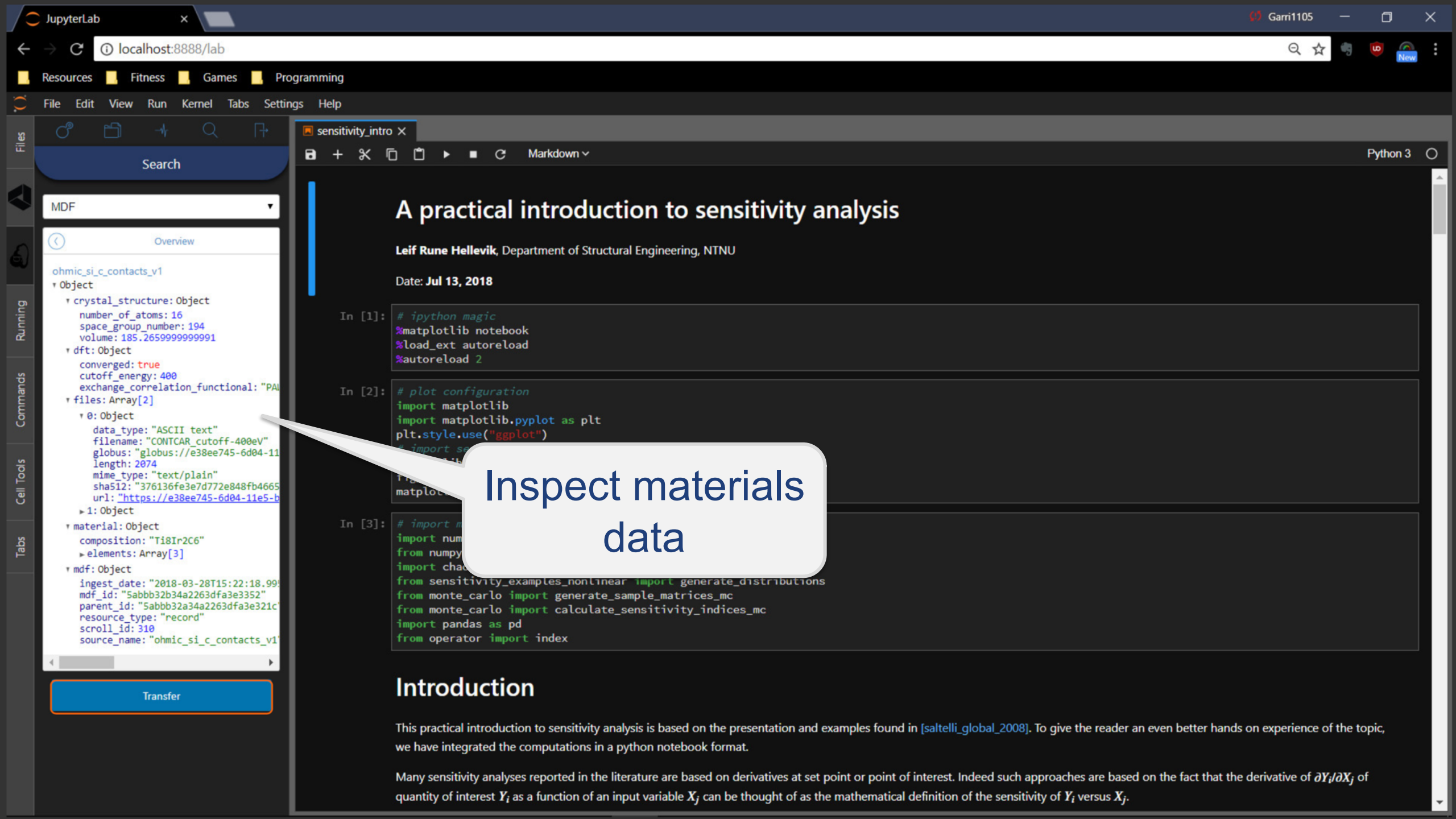
In [2]:
```
# plot configuration
import matplotlib
import matplotlib.pyplot as plt
plt.style.use("ggplot")
# import se
matplot
```

In [3]:
```
# import n
import num
from numpy
import cha
from sensitivity_examples_nonlinear import generate_distributions
from monte_carlo import generate_sample_matrices_mc
from monte_carlo import calculate_sensitivity_indices_mc
import pandas as pd
from operator import index
```

Inspect materials data

## Introduction

This practical introduction to sensitivity analysis is based on the presentation and examples found in [saltelli_global_2008]. To give the reader an even better hands on experience of the topic, we have integrated the computations in a python notebook format.

Many sensitivity analyses reported in the literature are based on derivatives at set point or point of interest. Indeed such approaches are based on the fact that the derivative $\partial Y_i/\partial X_j$ of a quantity of interest $Y_i$ as a function of an input variable $X_j$ can be thought of as the mathematical definition of the sensitivity of $Y_i$ versus $X_j$.

# Incorporate seamless parallel computing

- **(Data) science apps require…**
  - Interactivity
  - Scalability (more than a desktop)
  - Reproducibility (publish code, docs)

- **Solution: JupyterHub + Parsl**
  - Interactive computing environment
  - Notebooks for publication
  - Can run on dedicated hardware

**Parsl** **Python parallel library**

- **Tasks exposed as functions (Python, bash)**
- **Python code to glue functions together**
- **Globus for auth and data movement**

**parsl-project.org**

```python
@python_app
def compute_features(chunk):
    for f in featurizers:
        chunk = f.featurize_dataframe(chunk, 'atoms')
    return chunk

chunks = [compute_features(chunk)
          for chunk in np.array_split(data, chunks)]
```

# Containerized data science ecosystem

Users select container to execute custom Jupyter environment

Uniform IAM platform scalable for distributed objects

Registry for container discovery and referencing

**Auth**
globus

**Search**
globus

container metadata

**jupyterhub**

**Container Registry**

**Notebook Server**

Containers used for both Jupyter notebook server and compute nodes

Container definitions tracked in version control systems

container recipes

Containers staged to local file systems

**Transfer**
globus

Compute

Compute

Compute

Compute

containers

**ALCF Petrel**

**Supercomputer**

Containers used other tasks; analysis, ML, etc.