# An Overview of Jupyter Notebooks at SDCC (with examples)
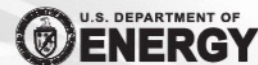
Ofer Rind

SDCC Technology Meeting

September 27, 2018

**70** YEARS OF **DISCOVERY**

A CENTURY OF SERVICE

U.S. DEPARTMENT OF **ENERGY**   **BROOKHAVEN** NATIONAL LABORATORY
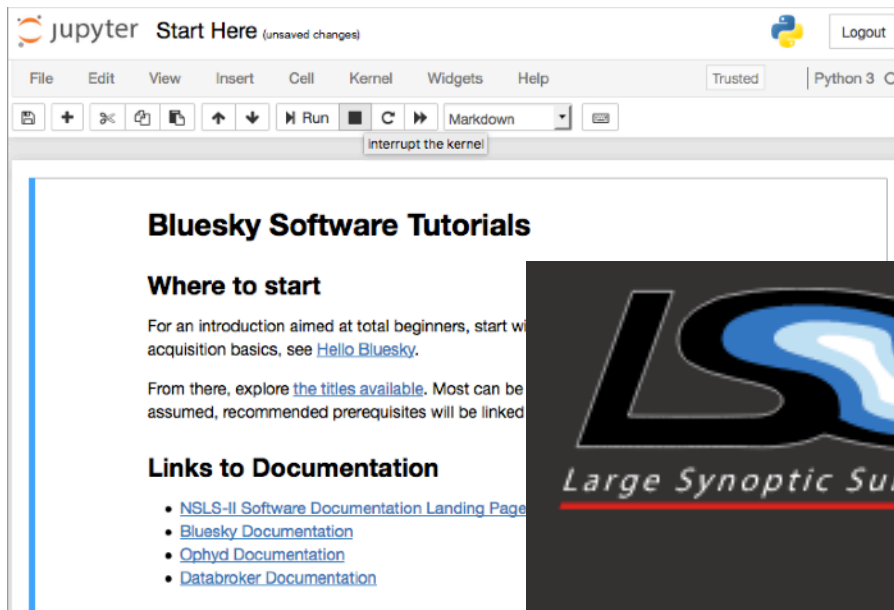
# Data Analysis As A Service

- Jupyter Notebooks (iPython)

  - Provide a flexible, standardized, platform independent interface through a web browser

  - Well-suited for *interactive* analysis

  - No local software to install

  - Many language extensions (kernels and tools available

  - Easy to share, reproduce, document results and create tutorials

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

- From the facility point of view:  Can we layer this service atop existing resources, without building a new dedicated infrastructure, such as a specialized cluster? (cf. CERN Swan)

# Growing in Popularity Across Scientific Community

# Some terminology

- **Jupyter notebook**: web-based application suitable for capturing the whole computation process: developing, documenting, and executing code, as well as communicating the results

- **Jupyterlab**: next-generation web-based user interface

# Some terminology

- **Jupyterhub**: multi-user hub, spawns, manages, and proxies multiple instances of the single-user Jupyter notebook server

# Current Setup at BNL

- Three Jupyterhub servers deployed on RHEV
  - <u>Jupyter01</u>: Primary HTC development environment with access to HTCondor queue
  - <u>Jupyter02</u>: Primary HPC development environment with access to IC via Slurm
  - <u>Jupyter03</u>: Currently testing OAuth/Keycloak implementation
- Access via ssh tunnel through firewall to Jupyterhub https proxy
- Kerberos auth to Jupyterhub server - use your SDCC account
  - Transparent setup leverages PAM stack
- Implemented a couple of experiment-specific custom kernels

# Examples

- Slurm - Higgs Test (Tensorflow) - IC python environment

- Belle-II (Bilas Pal) - Custom python3 kernel with BASF2 environment

- ATLAS GridScan (Viviana Cavaliere) - Custom python2 kernel with ATLAS environment

- ATLAS ML (JianCong Zeng) - Python3 kernel, with scikit-hep uproot and ML packages

- HTCondor API - Kristy Li/Will Strecker-Kellogg

# Questions For Discussion

- Are people currently using Jupyter notebooks? How are they being used?

- Are there use cases at the SDCC, including at scale?

- What services would users like SDCC to provide?

# Extra Slides

# Example: Creating a Local Kernel Environment

```
[-bash-4.2$ cd ~/.local/share/jupyter/kernels/ATLAS/
[-bash-4.2$ ls
kernel.json   logo-64x64.png   setup.sh
[-bash-4.2$ cat kernel.json
{
 "argv": [
   "/usatlas/u/rind/.local/share/jupyter/kernels/ATLAS/setup.sh",
   "-f",
   "{connection_file}"
 ],
 "display_name": "ATLAS test",
 "language": "python"
}
[-bash-4.2$ cat setup.sh
#! /usr/bin/env bash
#
# Set path
#export PATH=/u0b/software/anaconda2/bin:$PATH
#export PATH=/usr/bin:/etc:/usr/sbin:/usr/ucb:$HOME/bin:/usr/bin/X11:/sbin:.

export ATLAS_LOCAL_ROOT_BASE=/cvmfs/atlas.cern.ch/repo/ATLASLocalRootBase
export ALRB_localConfigDir=$HOME/localConfig
source ${ATLAS_LOCAL_ROOT_BASE}/user/atlasLocalSetup.sh --quiet
source ${ATLAS_LOCAL_ROOT_BASE}/utilities/oldAliasSetup.sh root --rootVersion=6.08.06-HiggsComb-x86_64-slc6-gcc49-opt

# python will be in the anaconda2 directory
exec /u0b/software/anaconda2/bin/python -m ipykernel_launcher $@
```

# HTCondor API

- Provide access to distributed computing through familiar APIs (python's threading, multiprocessing, asyncio, etc…)
  - "I'd like to submit and manage a job or cluster of jobs"

```
[willsk@willsktop condor (master)]$ python
Python 2.7.15 (default, May  9 2018, 11:18:37)
[GCC 7.3.1 20180303 (Red Hat 7.3.1-5)] on linux2
Type "help", "copyright", "credits" or "license" for more informati
on.
>>> import job
>>> j = job.Job({"Executable": "/bin/sleep", "Arguments": '2000'})
>>> print j.submit()
78
>>> j.status
1
>>> j.hold()
[ TotalChangedAds = 1; TotalSuccess = 1; TotalBadStatus = 0; TotalP
ermissionDenied = 0; TotalError = 0; TotalNotFound = 0; TotalJobAds
 = 1; TotalAlreadyDone = 0 ]
>>> j.status
5
>>> j.remove()
[ TotalChangedAds = 1; TotalSuccess = 1; TotalBadStatus = 0; TotalP
ermissionDenied = 0; TotalError = 0; TotalNotFound = 0; TotalJobAds
 = 1; TotalAlreadyDone = 0 ]
>>> j._jobdata
>>>
```

```
[willsk@willsktop condor (master)]$ cat tests/jdfs/sleep.job
Executable      = /bin/sleep
Arguments       = 10
Log             = sleep.$(cluster).log

Queue 4
[willsk@willsktop condor (master)]$ python
Python 2.7.15 (default, May  9 2018, 11:18:37)
[GCC 7.3.1 20180303 (Red Hat 7.3.1-5)] on linux2
Type "help", "copyright", "credits" or "license" for more information.
>>> import job
>>> c = job.JobCluster.from_jdf('tests/jdfs/sleep.job')
>>> c.status
{0L: 1L, 1L: 1L, 2L: 1L, 3L: 1L}
>>> c.status
{0L: 2L, 1L: 2L, 2L: 2L, 3L: 2L}
>>> c.wait()
```

# HTCondor API

- At a higher level, abstract away the batch job layer
  - I'd like to run over a dataset

```
>>> data = (basename(x) for x in glob.glob('tests/*.dat'))
>>> c = CondorMapper('script.sh', data)
>>> c.set_on_complete(lib.merge)
>>> c.wait()
>>> []
```

- Serialize function, ship off to jobs, serialize output, gather…

- Early stage of development - see Will Strecker-Kellogg for details

- Also looking at dask-jobqueue

```python
import numpy, collections

def logistic(r):
    x = 0.5
    d = collections.deque(maxlen=10)
    for _ in xrange(50000):
        x = x * r * (1.0 - x)
        d.append(x)
    return list(d)

map(logistic, numpy.arange(1.01,3.99,0.01))

condormap(logistic, numpy.arange(1.01,3.99,0.002))
```