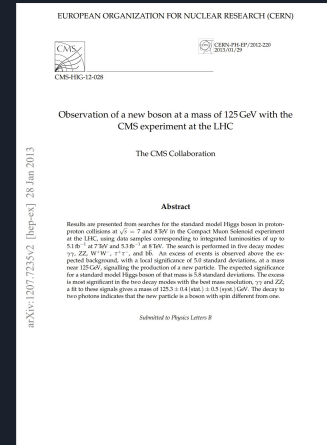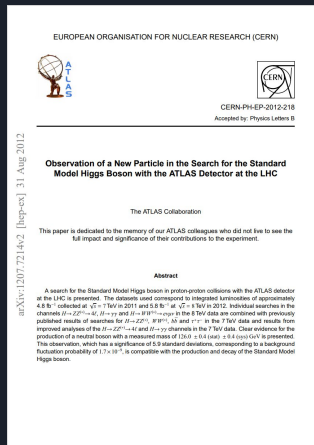# Using Python, coffea, and ServiceX to Rediscover the Higgs. Twice.

Baidyanath Kundu (Manipal Institute of Technology)

Gordon Watts (University of Washington)
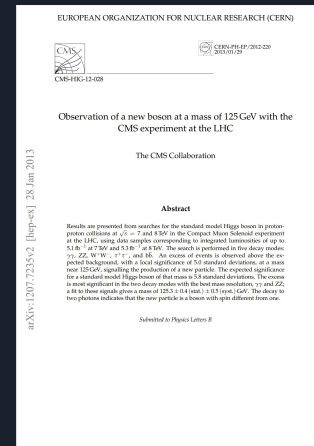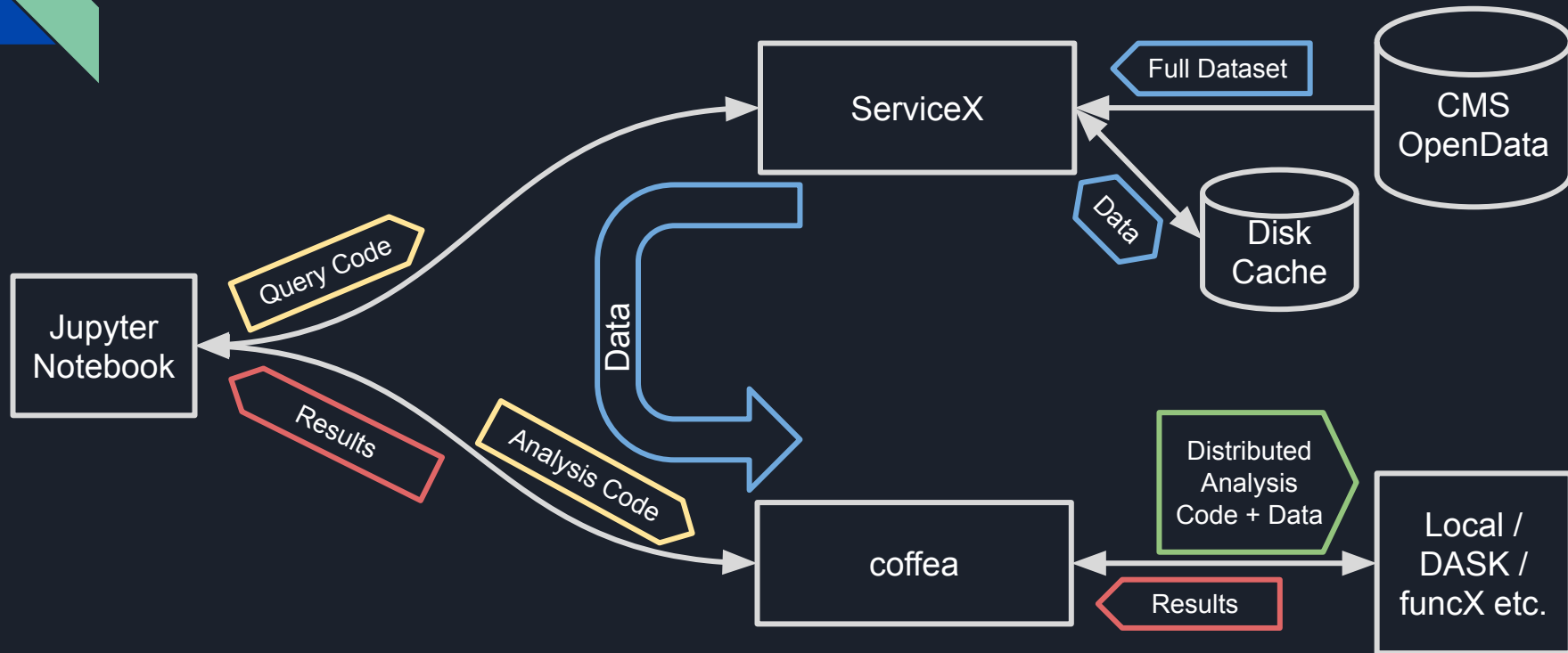
# What Higgs, Exactly?



Both experiments have released their Run 1 Higgs Discovery Data on CERN Open Data
Including Source Code To Reproduce Analyses!

# What Higgs, Exactly?

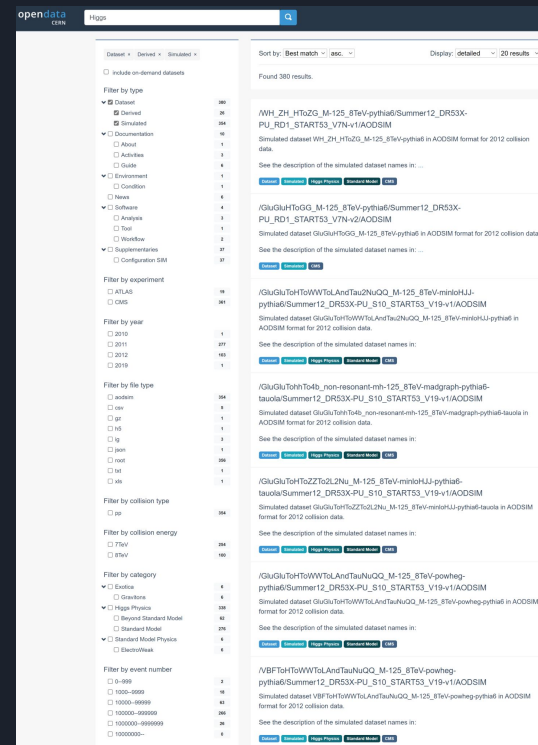Can we reproduce their results using the ServiceX, func_adl, coffea tool chain?
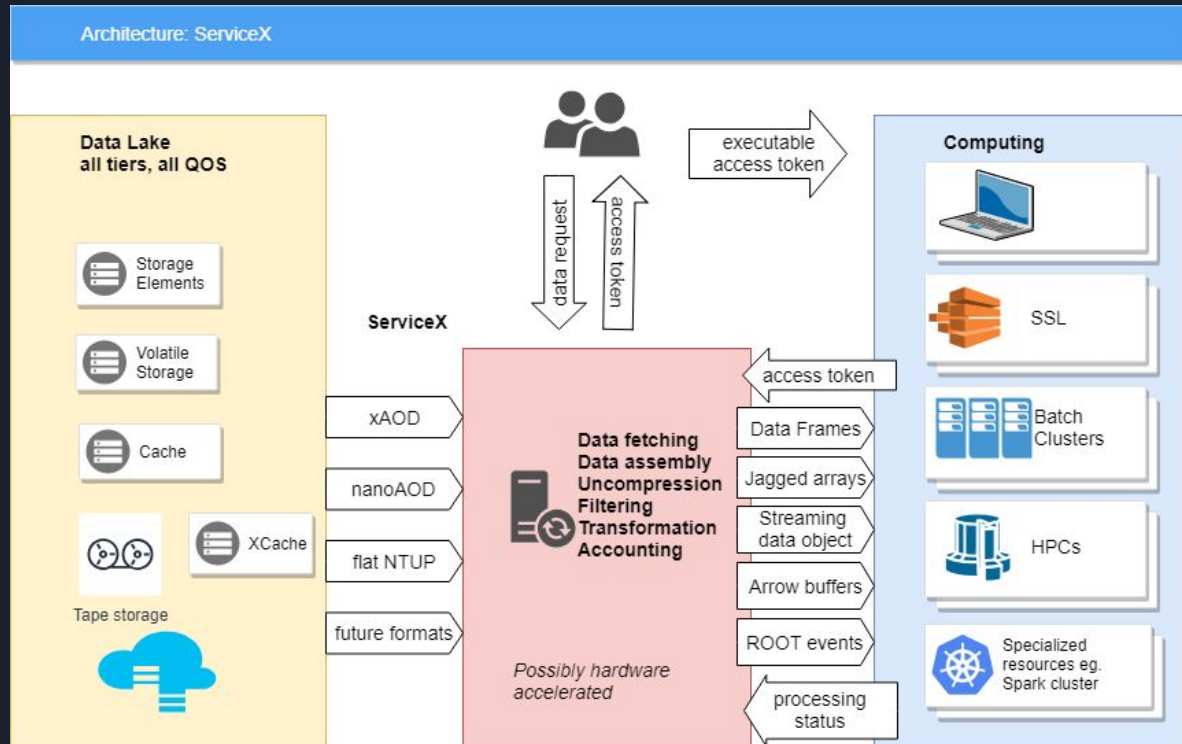
# Analysis In Python

# CERN opendata



Over 2 PB of data, 1000's of datasets

# ServiceX: Overview

# coffea

A light-weight processing framework for processing event columnar data.

- Processing files
- Distributed processing in a farm (DASK, slurm, Spark, etc.)
- Flexible data model that makes column-wise data look like row-wise.



In the end, many parts of coffea should disappear into dedicated packages... until then, missing gaps in functionality have a home here.

documentation

# Notebooks

We have 3 notebooks to work through!

⭘ Introduction to func_adl, ServiceX, and coffea

⭘ The ATLAS Higgs Discovery Dataset

⭘ The CMS Higgs Discovery Dataset



[Tutorial GitHub Repo](#)

# Team Effort

- Data and MC: ATLAS & CMS Collaborations for all the data, and making it public!

- ServiceX: Ben Galewsky and Andrew Eckart and Suchandra Thapa and everyone else on the team

- CMS Run 1 AOD Transformer: Baidyanath Kundu

- Uproot/Flat ROOT backend: Mason Proffit

- IRIS-HEP for supporting a large fraction of these people in one way or another.

- CMS Awkward Code: Brian Cruz

# Running This Yourself...

[Tutorial GitHub Repo](#)

➡️ Instructions are on the repo

➡️ We are happy to give out the ServiceX end-points (unofficial atm)

➡️ Some integration with coffea is in beta

➡️ Not yet obvious that CERN Open Data portal can handle the data load!

➡️ Binder is not supported

## Using

You can find the final notebooks used in the talk in the talks directory. The notebooks directory contains practice notebooks used to develop concepts for the talk. They are not necessarily well documented.

### ServiceX for the demo

☰ README.md

```
api_endpoints:
  - endpoint: http://xxx.org
    type: open_uproot
  - endpoint: http://yyy.org
    type: cms_run1_aod

backend_types:
  - type: open_uproot
    return_data: parquet
  - type: cms_run1_aod
    return_data: root
```

Please get in touch with us to get the address of the open instances running `ServiceX`.

### Setting up the environment

Setup your environment:

1. This has been run under python 3.9.6. It should work with anything that is 3.7 or greater.
2. Check out this repository locally, and check out the coffea patched repository locally.
3. For the `coffea` repository, check out the branch pr_servicex_flat_root_files. For this package use the head.
4. `python -m venv .venv`, and activate the new environment.
5. `pip install -r requirements.txt`
6. In the root directory of the checked out `coffea` package, run `pip install -e[servicex]`.

From there you can start `jupyter-lab`.

### Running on binder

It is not currently possible to run on `binder` as `ServiceX` uses a non-standard port to download data.