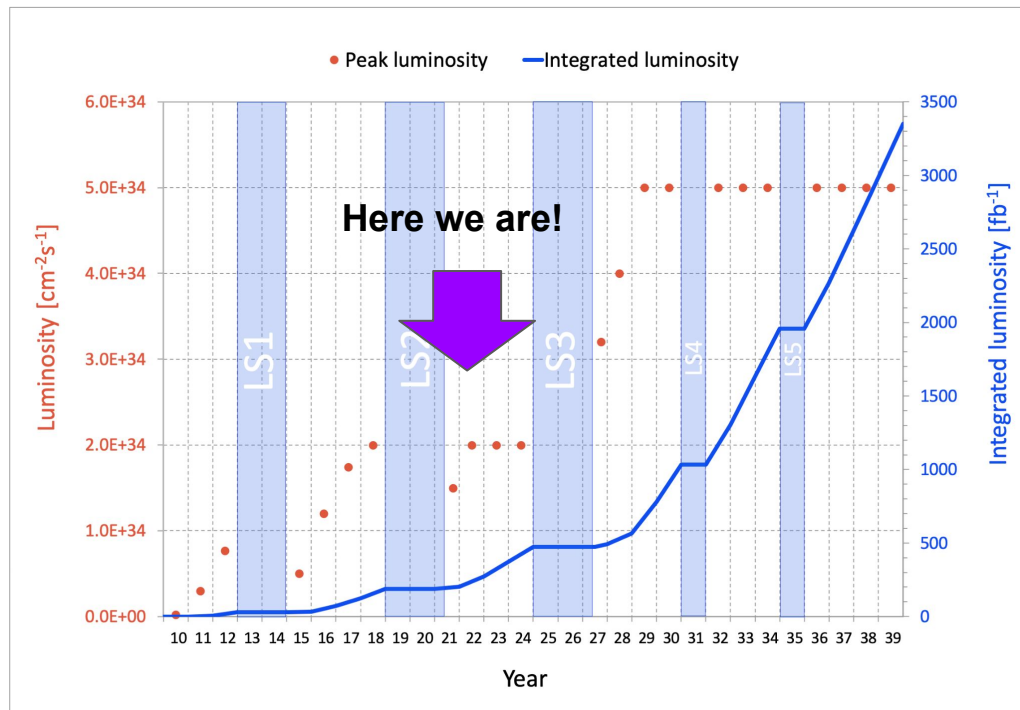# Coffea-casa: an analysis facility prototype

Mat Adamec, Ken Bloom, Oksana Shadura,
*University of Nebraska, Lincoln*

Garhan Attebury, Carl Lundstedt, John Thiltges
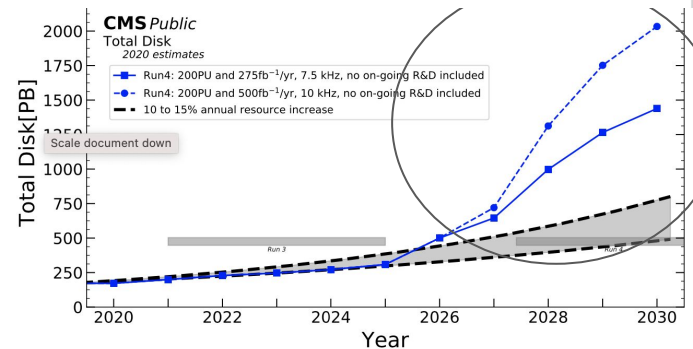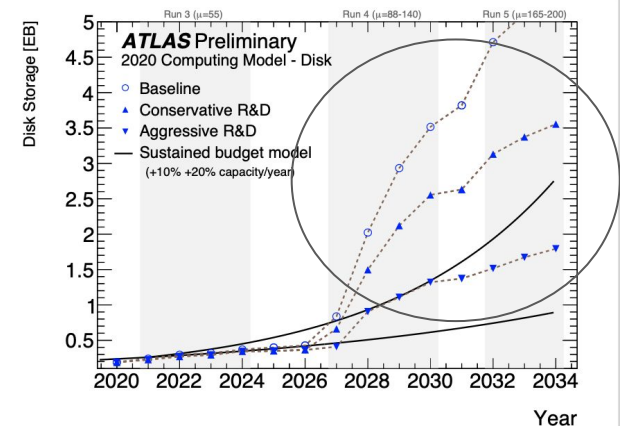*University of Nebraska Holland Computing Center*

Brian Bockelman
*Morgridge Institute*

Peak luminosity   Integrated luminosity

**Here we are!**



ATLAS Preliminary
2020 Computing Model - Disk
- Baseline
- Conservative R&D
- Aggressive R&D
- Sustained budget model
  (+10% +20% capacity/year)



CMS Public
Total Disk
2020 estimates
- Run4: 200PU and 275fb⁻¹/yr, 7.5 kHz, no on-going R&D included
- Run4: 200PU and 500fb⁻¹/yr, 10 kHz, no on-going R&D included
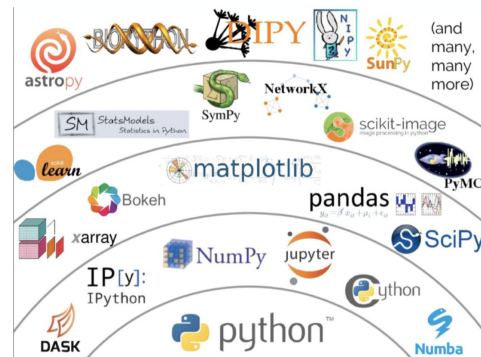- 10 to 15% annual resource increase

Our goal is driven by the desire *to bring simpler and more agile paradigms for analysis **today***, but the scale of the HL-LHC adds more complexity to the existing issue
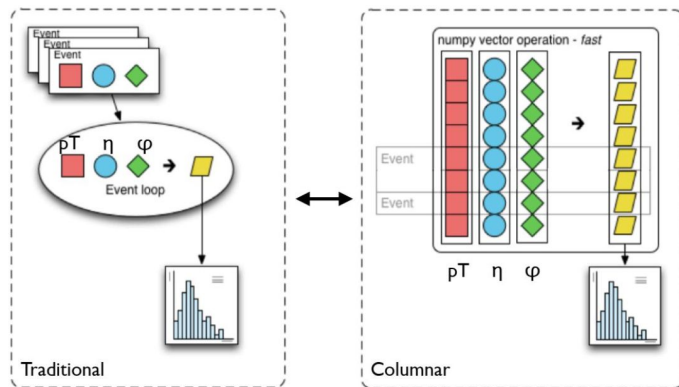
- **New pythonic ecosystem**
- Discovering the benefits of **column-oriented (columnar) data analysis**
- **Interactivity** for user data analysis
- Deliver the needed data to the processing workflow in a fine-grained approach (**data delivery services)** and **efficient storage technologies** (e.g. object stores)
- **Kubernetes (k8s) and** new concept of **"infrastructure as code"**
- **Portability** and flexibility across different environments
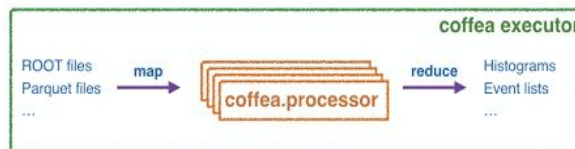- Integration with existing resources: current infrastructure is not going to be replaced in one day
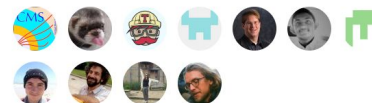
New columnar data analysis concepts!

User just needs to define a high-level wrapper around user analysis code: **the coffea processor** and coffea framework will take care of everything incl. *scaling-out*

Distributed executors!

**Coffea** developers: Lindsey Gray, Matteo Cremonesi, Bo Jayatilaka, Oliver Gutsche, Nick Smith, Allison Hall, Kevin Pedro (**FNAL**); Andrew Melo (Vanderbilt); and others

Contributors 32

+ 21 contributors

**Today (event size)**

| MINIAOD** ~ 35kB | NANOAOD*** ~ 1 kB |
|---|---|

**HL-LHC (event size)**

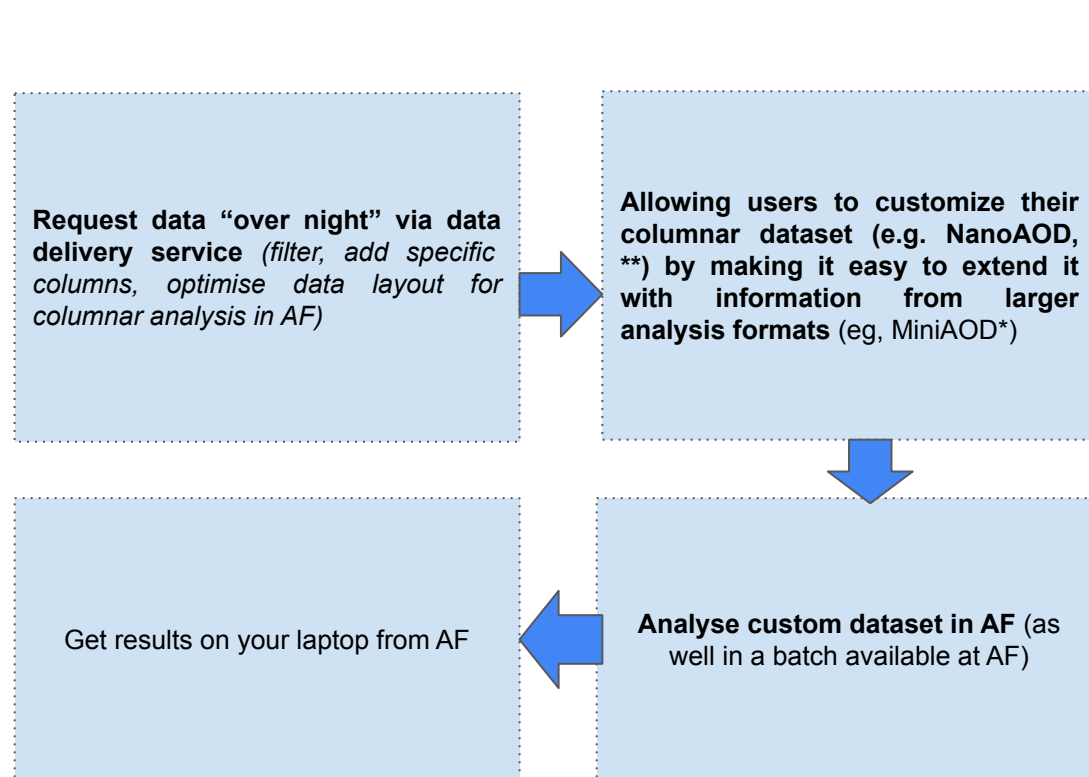| **MINIAOD ~ 250kB** | **NANOAOD ~ 2 kB** |
|---|---|

If to switch to more compact data format - **"NanoAOD"**, *data volumes should be manageable,* but this format will likely be missing data elements necessary for any given analysis.

*Idea:* If we use a sample small enough to be used for interactive analysis in AF: NANOAOD, for example it could be useable as driver for data delivery service to add objects from MINIAOD overnight!

*Here on CMS data formats as an example, shown ideas are applicable for other experiments
**MiniAOD - c++ class hierarchy data format
**NanoAOD - compact, Ntuple like data format readable by bare ROOT

**The AF expected to assist with 90% of analyses** using NanoAOD** by merging parts or derived from MiniAOD* into NanoAOD** (automatically, without the intervention of the end-user)

**Request data "over night" via data delivery service** (filter, add specific columns, optimise data layout for columnar analysis in AF)

**Allowing users to customize their columnar dataset (e.g. NanoAOD, **) by making it easy to extend it with information from larger analysis formats** (eg, MiniAOD*)

**Analyse custom dataset in AF** (as well in a batch available at AF)
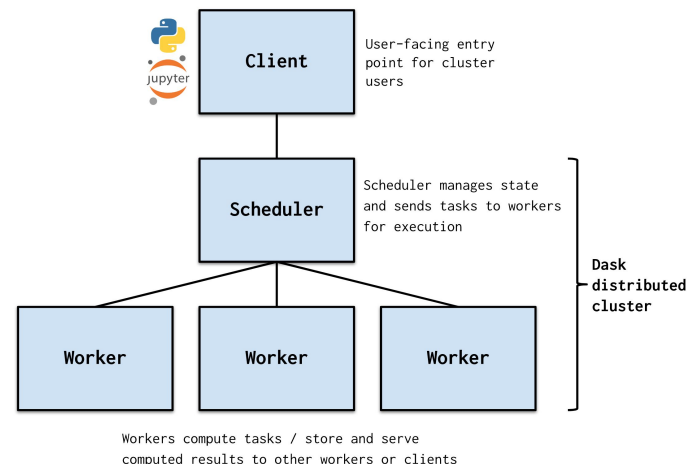
Get results on your laptop from AF

*MiniAOD* - c++ class hierarchy data format
**NanoAOD* - compact, Ntuple like data format readable by bare ROOT
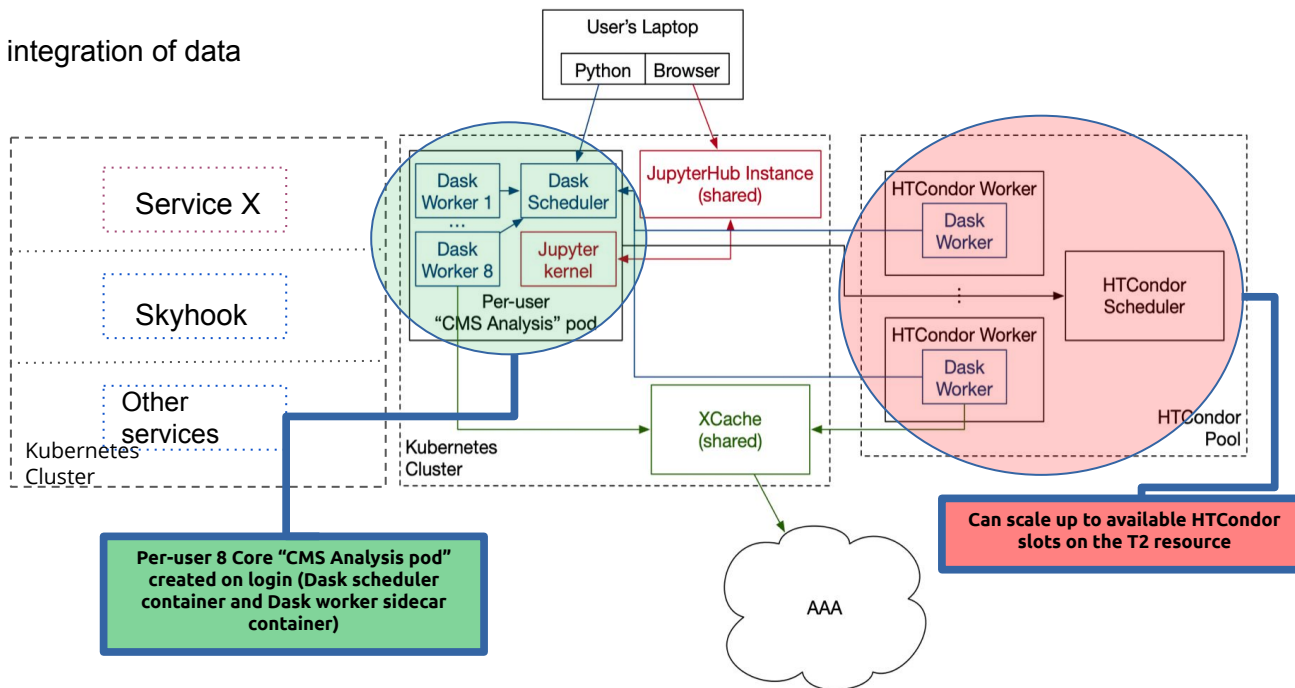
# Why we use Dask?

- Dask provides a task-management computational framework in Python based on the manager-worker paradigm
- Dask exposes lower-level APIs letting to build custom systems for in-house applications (!)
- Integrates with HPC clusters, running a variety of schedulers including SLURM, LSF, SGE and *HTCondor via "dask-jobqueue"*
- ***This allows us to create a user-level interactive system via queueing up in the batch system***

**Dask can be used inside Jupyter or you can simply launch it through Jupyter and connect directly from your laptop**

Client — User-facing entry point for cluster users

Scheduler — Scheduler manages state and sends tasks to workers for execution

Worker   Worker   Worker

Dask distributed cluster

Workers compute tasks / store and serve computed results to other workers or clients

**_We are easily bridging K8s resources with UNL Tier2 resources, while providing interactive environment!_**

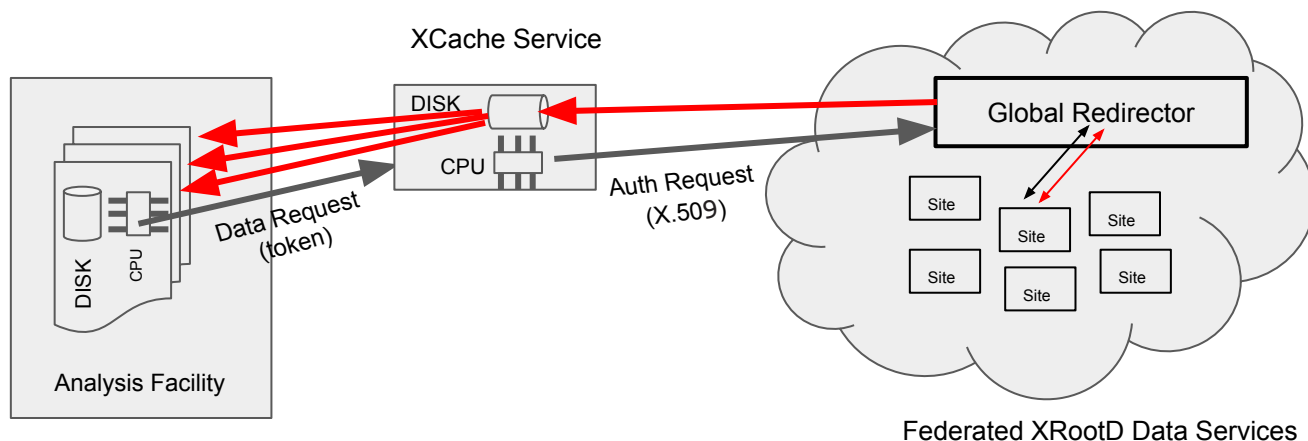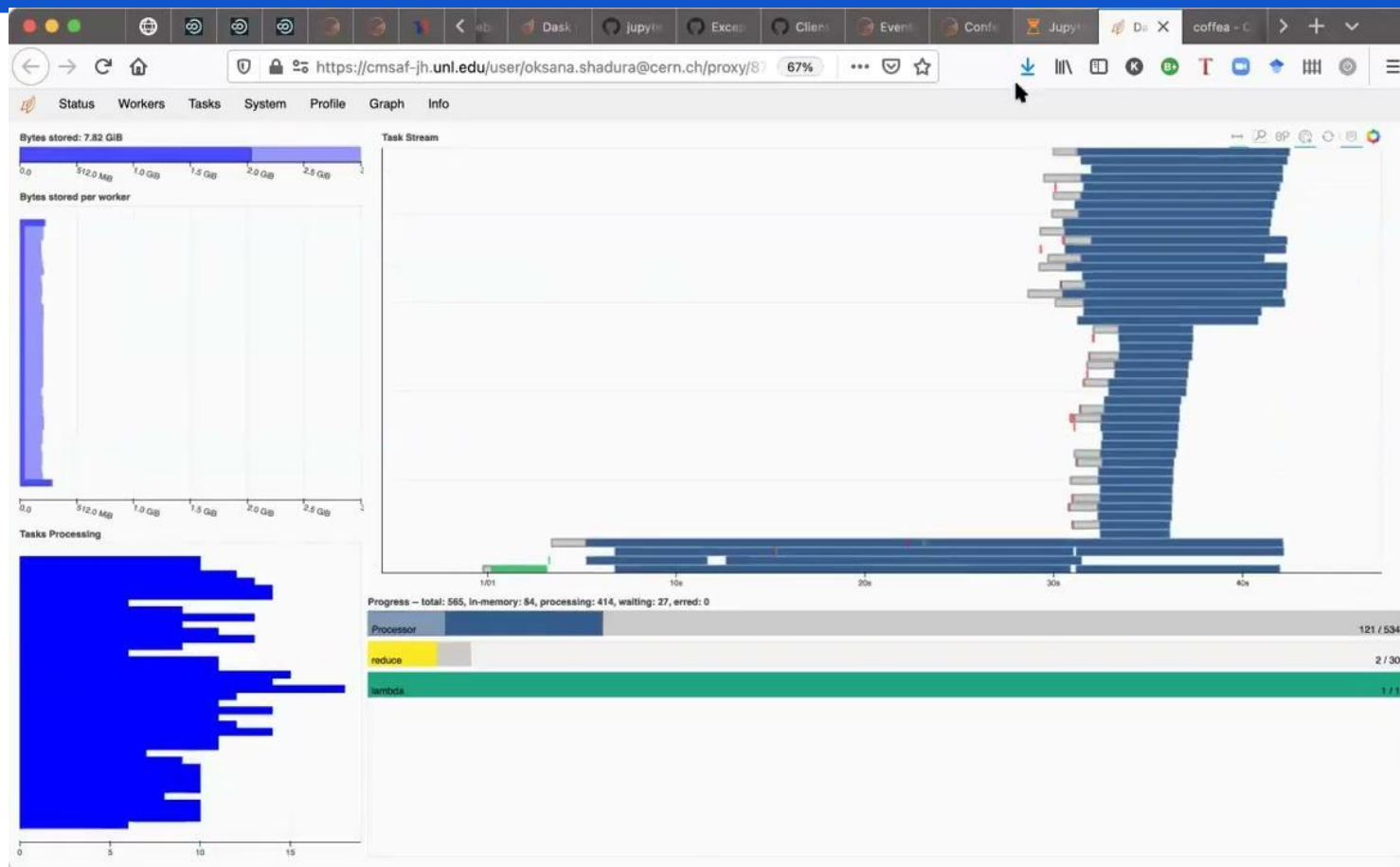On-going work on integration of data delivery services



Per-user 8 Core "CMS Analysis pod" created on login (Dask scheduler container and Dask worker sidecar container)

Can scale up to available HTCondor slots on the T2 resource

- **Authentication inside the system is independent of grid credentials**
  - Coffea-casa facility uses **OpenID Connect (OIDC)**
  - **Enabled token authentication** for HTCondor:
    - Generated a token for authentication with HTCondor, required for Dask scale-out to the larger resources
  - **Generated a data access token for authentication with a local XRootD server**
  - Generated X.509 credentials (including a CA, host certificate, and user certificate) for use in Dask for TLS as well for user communication to Dask scheduler endpoint
- Security: TLS enabled communication between workers and scheduler by default
- Kubernetes pod customization 'hook' to create secrets for services
- Highly customized **"CMS Analysis" Docker container(s)**
- All features are **incorporated into a Helm chart** (Kubernetes packaging format)

- ***CoffeaCasaCluster:* extending HTCondorCluster integration for Dask**
  - To handle the customizations needed for the Coffea-casa environment, we developed the *CoffeaCasaCluster* object, an extension of ***Dask-jobqueue***'s *HTCondorCluster* object.
  - *CoffeaCasaCluster* ensures the Dask worker starts with the appropriate Docker container in the HTCondor batch system with appropriate configurations and with the firewall ports configured correctly.

- For speeding up data access Nebraska Tier-2 hosts an **XCache service with 90TB of cache space**
- Access data hosted by an HEP experiment:
  - *no GSI credential within the facility,* **the auto-generated data access token can be used to authenticate with an proxy service based on XRootD/XCache**

# Coffea-casa @ UNL - demo



Link to video

# Inviting first users



tHq analysis (116M Events, 34 GB, 78 Files)
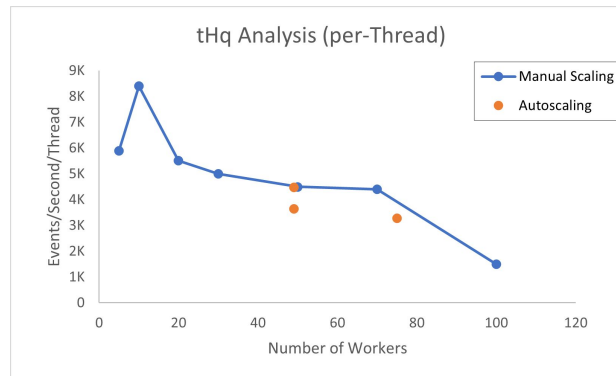
Mat Adamec (UNL undergrad)

Zora Che (BU CS undergrad)

FTAnalysis (740M Events, 81 GB, 1377 Files)

Top quark analyses using the Coffea framework
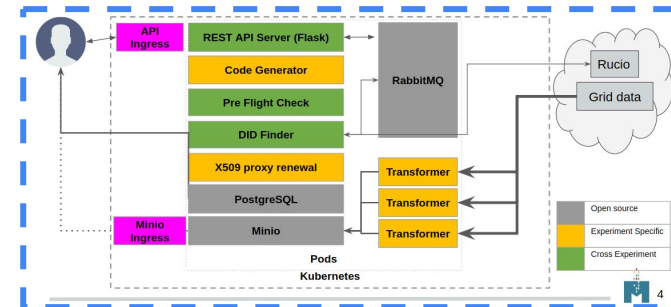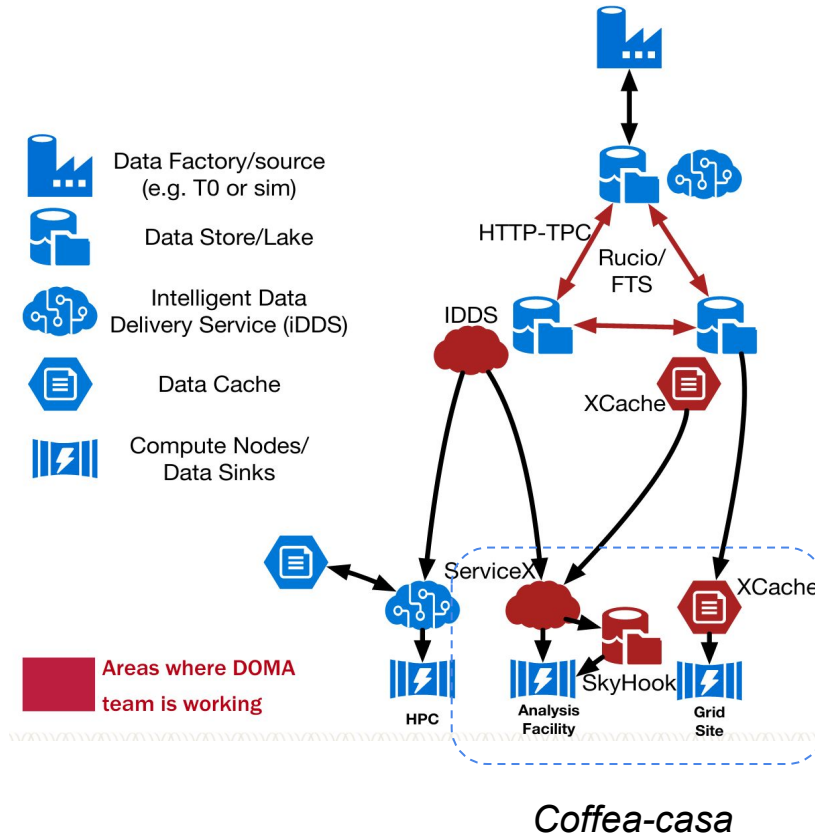https://github.com/TopEFT/topcoffea

Users already running their analysis on Coffea-casa:

- Some of the examples are **done by undergraduate students**
- **Approachable** (even with all complexity of system behind) and **interactive**

*tHq results require some investigation*

Data Factory/source (e.g. T0 or sim)

Data Store/Lake

Intelligent Data Delivery Service (iDDS)

Data Cache

Compute Nodes/ Data Sinks

**Areas where DOMA team is working**

HTTP-TPC

Rucio/ FTS

IDDS

XCache

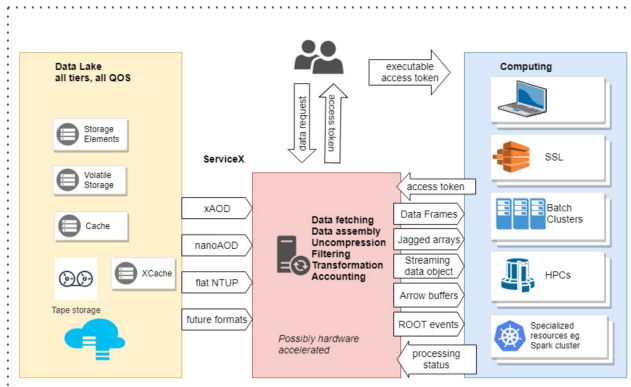ServiceX

XCache

HPC

Analysis Facility

Grid Site

SkyHook

*Coffea-casa*

*Servicex*

*Skyhook*

# Future work: data delivery services @ Coffea-casa

## ServiceX



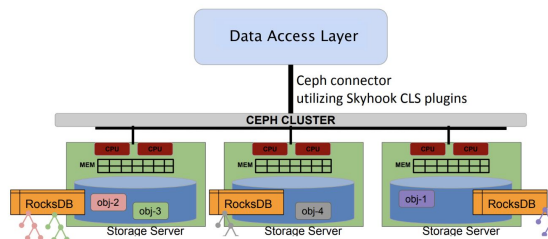**ServiceX provides user level ntuple production**
- Converts experiment-specific datasets to columns (e.g. NanoAOD, DAOD)
- Enable simple cuts or simple derived columns and fields (*heavy-weight analysis will still happen via some separate processing toolchain (like CMS CRAB)*)
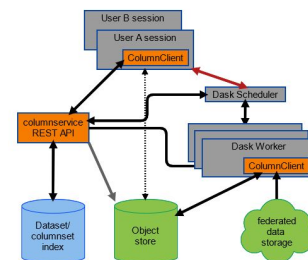
***[Check the talk of KyungEon Choi@ vCHEP2021]***

## Skyhook DM



The **Skyhook DM is converting event data from ROOT files to the internal object-store format**
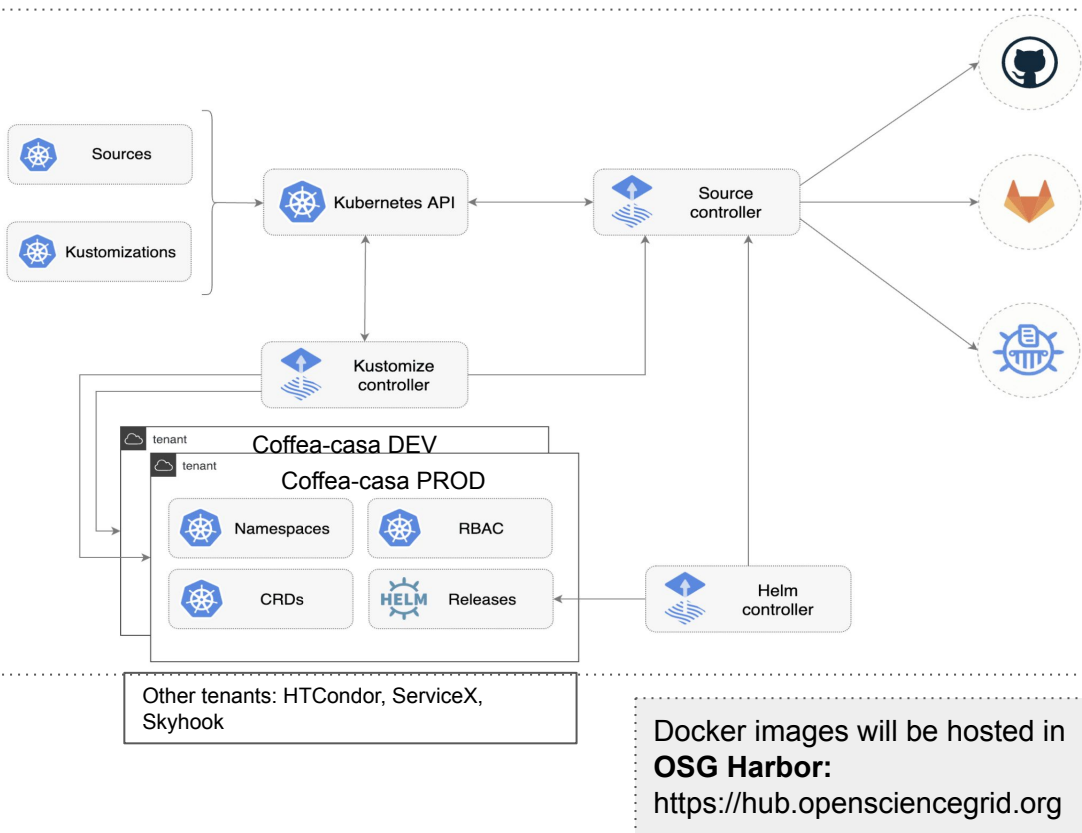
- Provides *extension of Ceph*-side C++ plugins transition from on-disk format to desired memory format
- Uses Dask workers to distribute data to clients
- Data delivered as Arrow tables or via Arrow Dataset API

## Columnservice



**Columnservice (FNAL)** is a multi-tenant service for caching columnar data that removes the need to curate skims and re-run expensive algorithms.

**GitOps** combines ***Git with Kubernetes programmatic service deployment*** properties and serves as an operating model for developing and delivering Kubernetes-based infrastructure and applications.

**We expect t**o package the core infrastructure (e.g., removing the site-specific passwords and secret keys) as a Helm chart, which will support different configuration such as **opendata coffea-casa,** **ATLAS coffea-casa or maybe generic af-casa!**

Other tenants: HTCondor, ServiceX, Skyhook

Docker images will be hosted in **OSG Harbor:** https://hub.opensciencegrid.org

- **Two AF facilities** with the possible outcome of adding more sites as soon as we gain experience

**CMSAF @T2 Nebraska**
**"Coffea-casa"**
**https://cmsaf-jh.unl.edu**

**Elastic AF @ Fermilab**

**Developed by:** Burt Holzman, Maria Acosta (FNAL)

# Conclusions

- The prototype analysis facility at Nebraska, **Coffea-casa, serves as an effective prototype and demonstration of several technologies under development for use in HL-LHC analysis;**
- Coffea-casa demonstrates features such as efficient data access services, notebook interfaces, token authentication, and automatic, external cluster scaling;
- We believe an critical future feature will be access to a **"column service":** the facility can be used to serve that "column" from a remote site;
- **Initial users have been testing the facility** to provide feedback and the team plans to distribute Coffea-casa products artifacts for use at other Kubernetes-based sites.

# Thank you!

Coffea-casa webpage
GH discussions
coffea-casa-dev@cern.ch

# Deploying XCache as a direct dependency Coffea-casa



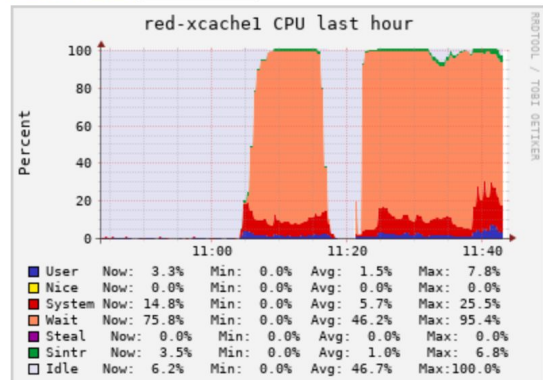- *XRootd file-based caching proxy* (known as Proxy File Cache (PFC) within XRootd code and documentation)
- **Already give visible performance improvement**
- Should be beneficial for datasets caching **in case of skimming**
- We are working on XCache Helm chart for easier deployment