



# Scalable Systems Laboratory Year 3

Rob Gardner, UChicago  
IRIS-HEP Executive Board Meeting  
February 16, 2021



Supported by National Science Foundation  
under cooperative agreement OAC-1836650



# Major SSL themes

- Facility R&D
  - *SSL to inform Tier2 evolution & analysis facilities*
  - *Multi-site, hyperconverged infrastructure*
  - *Supporting flexible, reproducible deployments*
- Supporting scalability & functional testing for IRIS-HEP and community partners
- Supporting IRIS-HEP Grand Challenges
  - *Focus on Analysis Challenge*
- Accelerated Data Delivery R&D
  - *Explore hardware acceleration at different points in the infrastructure*



quick review of some SSL  
deployments to date



# Some SSL deployments

**A diversity of  
services and  
developer  
engagements**

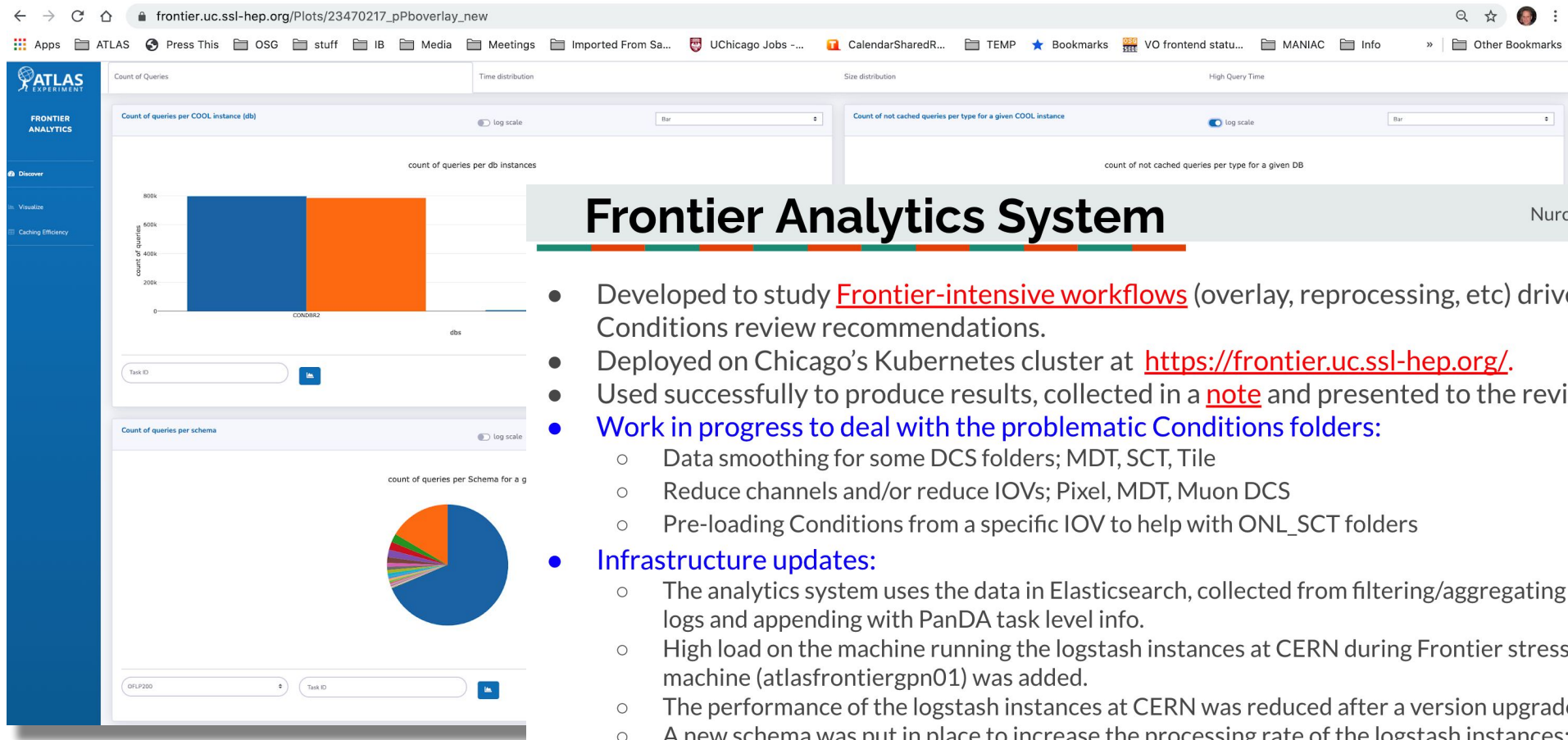
DOMA::ServiceX	Data transformation and delivery service for LHC analyses (IRIS-HEP)
DOMA::Skyhook	Programmable storage for databases, scaling Postgres with Ceph object store (IRIS-HEP)
REANA	Reusable Analysis Service (CERN development team)
CODAS Platform	JupyterLab notebooks, access to GPU resources on the Pacific Research Platform for annual summer CoDaS-HEP training event (IRIS-HEP SSC area)
Frontier Analytics	Analyze and improve data access patterns for ATLAS Conditions Data (ATLAS Distributed Computing Group)
perfSONAR Analytics	Network route visualization based on perfSONAR traces (NSF SAND project)
Parsl / FuncX	Parallel programming in Python, serverless computing with supercomputers (Computer Science)
Large-Scale Systems Group @ UChicago	Serverless computing with Kubernetes (Computer Science)
SLATE & OSG	Backfilling otherwise unused cycles on SSL with work from the Open Science Grid & ATLAS using the SLATE tools

**also: in discussions with Coffea team to provide infrastructure when appropriate**



# Frontier Analytics on SSL (ATLAS)

in production



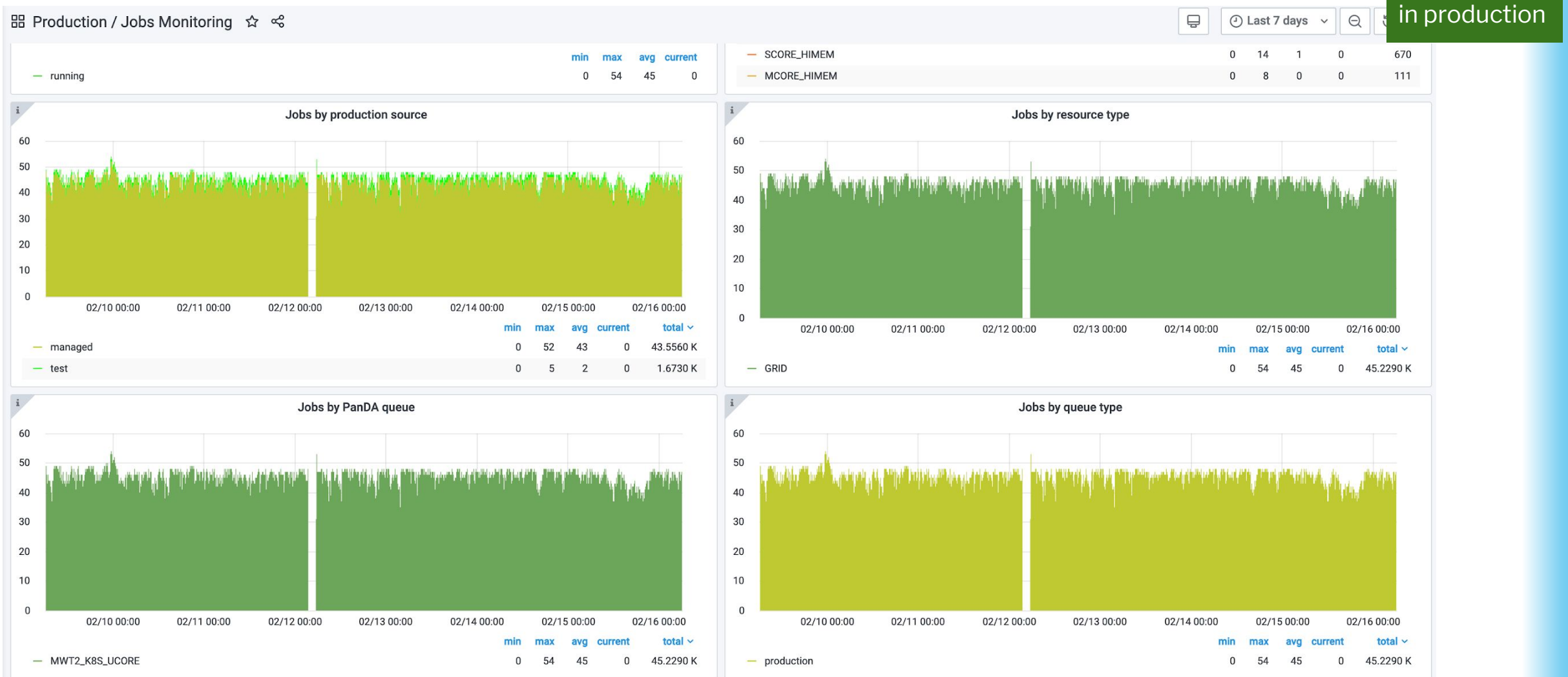
## Frontier Analytics System

Nurcan, Julio, Andrea, Elizabeth

- Developed to study Frontier-intensive workflows (overlay, reprocessing, etc) driven by the Run-3 Conditions review recommendations.
- Deployed on Chicago's Kubernetes cluster at <https://frontier.uc.ssl-hep.org/>.
- Used successfully to produce results, collected in a note and presented to the review committee.
- **Work in progress to deal with the problematic Conditions folders:**
  - Data smoothing for some DCS folders; MDT, SCT, Tile
  - Reduce channels and/or reduce IOVs; Pixel, MDT, Muon DCS
  - Pre-loading Conditions from a specific IOV to help with ONL\_SCT folders
- **Infrastructure updates:**
  - The analytics system uses the data in Elasticsearch, collected from filtering/aggregating the Frontier server logs and appending with PanDA task level info.
  - High load on the machine running the logstash instances at CERN during Frontier stress test in March. A new machine (atlasfrontiergpn01) was added.
  - The performance of the logstash instances at CERN was reduced after a version upgrade. It was downgraded.
  - A new schema was put in place to increase the processing rate of the logstash instances; group lines of each query and create multiple output files that serve as input to final logstash filtering instances.
- **Plan is to run the next round of Frontier stress tests in the upcoming months after the adjustment of the problematic Conditions folders is in place gradually and use the analytics system to check on the improvements.**



# SSL as Harvester k8s target (ATLAS)

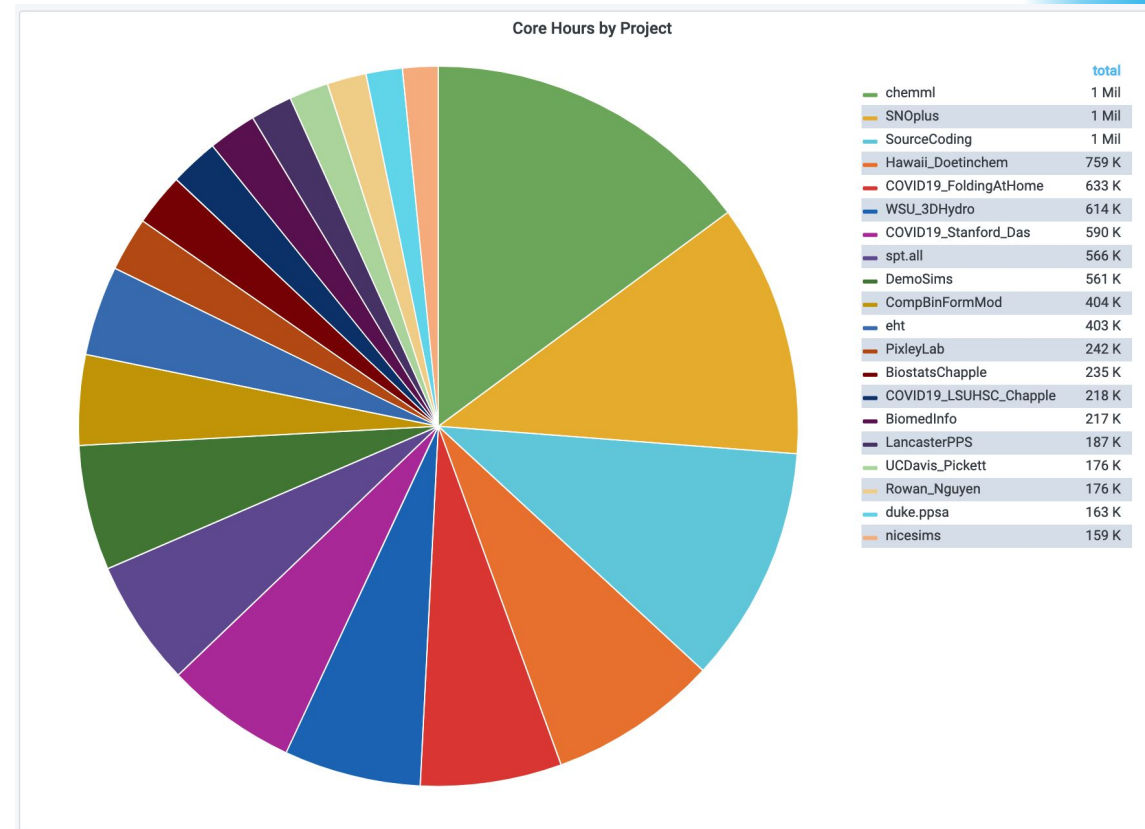
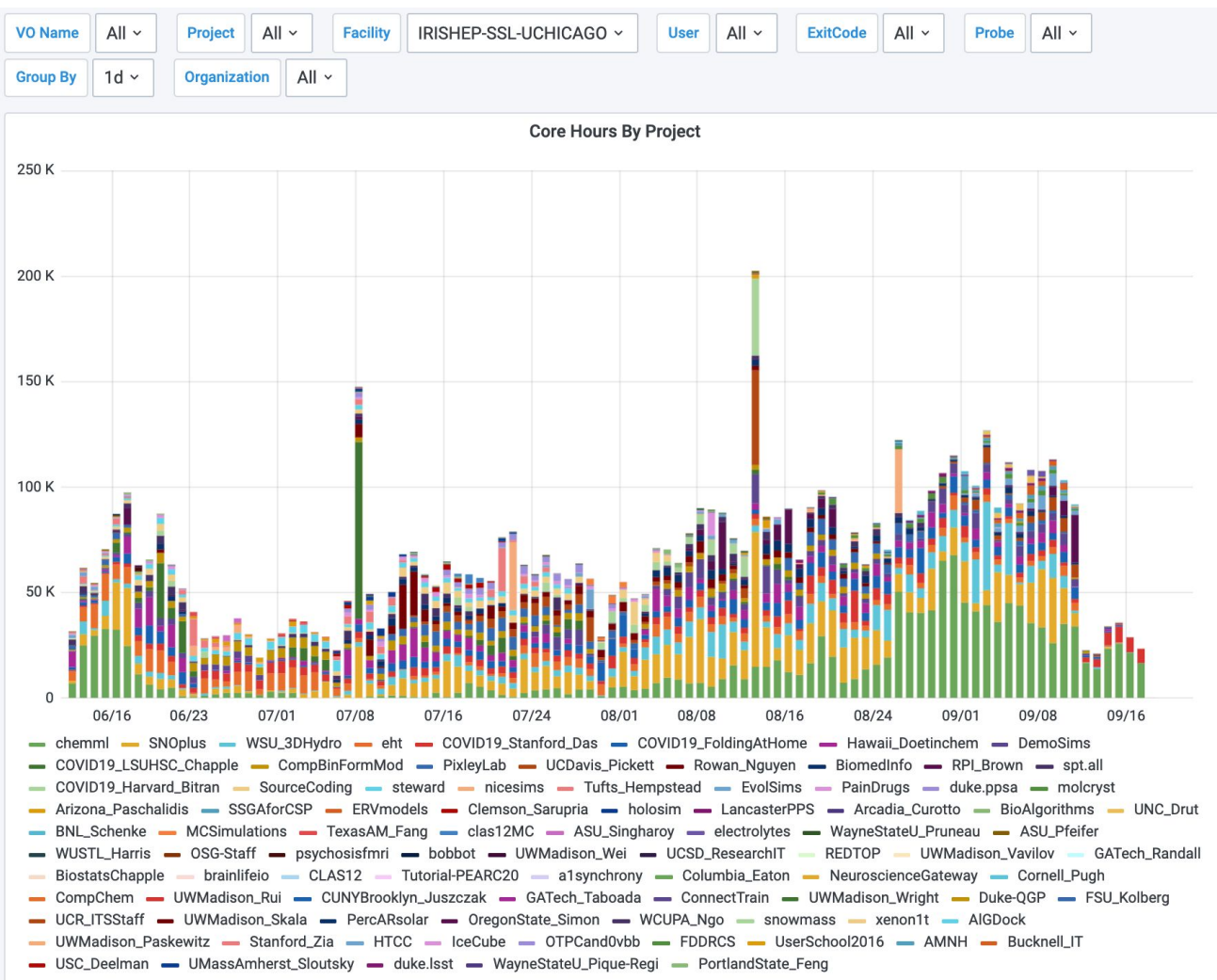






# k8s deployment of OSG workers

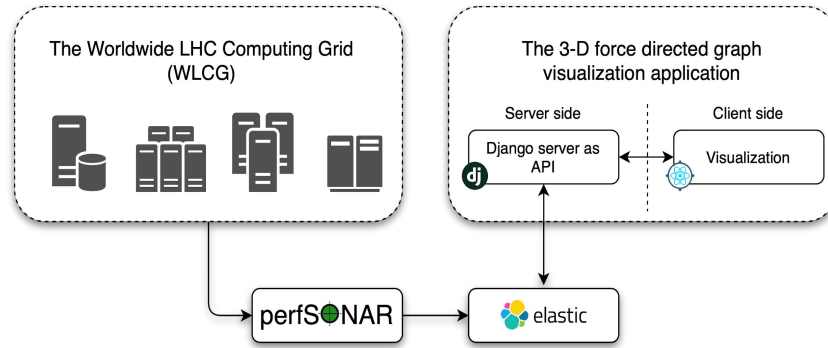
in production





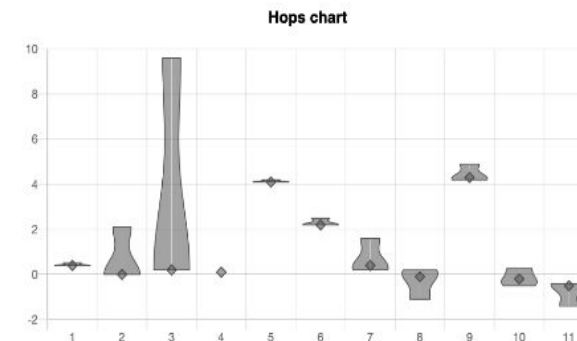
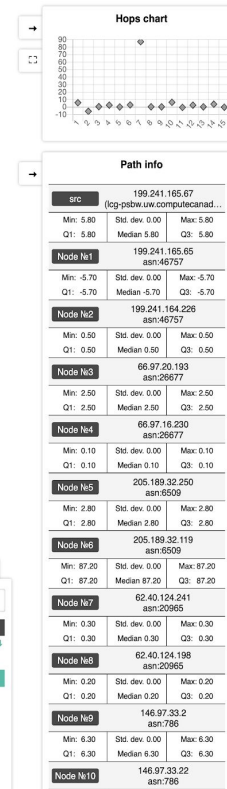
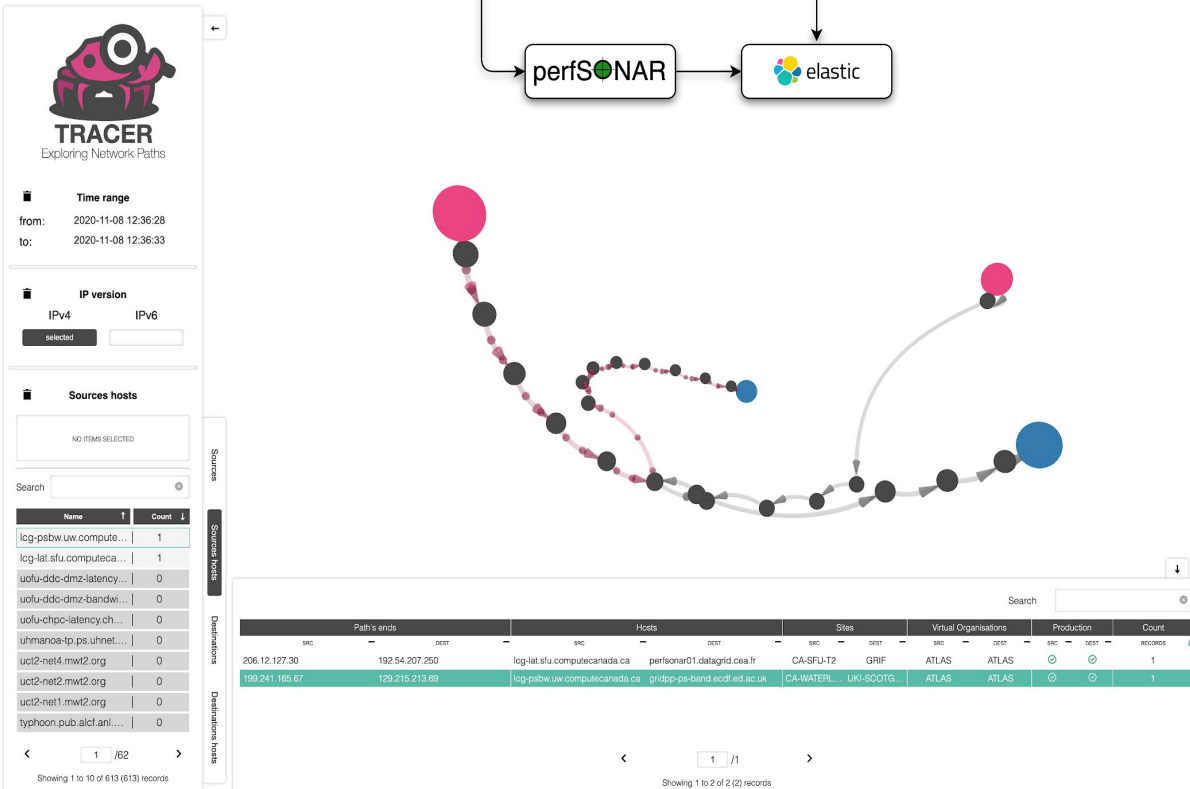
# PerfSONAR Network traces TRACe-ExploRer (NSF SAND)

in production



A visualization platform for researchers and network engineers to:

- better understand the topology of our RENs
- identify non-reliable or/and non-optimal network paths ( eg. routing loops, rapidly changing routes)







# Hosting OSG services (on River)

in production

```
$ slate instance list --group osg-ops
```

Name	Cluster	ID
osg-hosted-ce-amnh-ares	uchicago-river-v2	instance_AOUk0iliGjg
osg-hosted-ce-amnh-hel	uchicago-river-v2	instance_-5tbF_slj3k
osg-hosted-ce-amnh-mendel	uchicago-river-v2	instance_KlZKGY-i5hM
osg-hosted-ce-asu-dell-m240	uchicago-river-v2	instance_-T9qcccY3e0
osg-hosted-ce-clarkson-acres	uchicago-prod	instance_omQQTlH2-XU
osg-hosted-ce-computecanada-cedar	uchicago-river-v2	instance_EjL5pbnc594
osg-hosted-ce-fsu-hnpgrid	uchicago-river-v2	instance_j1D_dZ3_jrI
osg-hosted-ce-gsu-acore	uchicago-river-v2	instance_qU078rDyrSw
osg-hosted-ce-my-cluster	osgcc	instance_6YYWSOQ2Ahk
osg-hosted-ce-nd-caml-gpu	uchicago-river-v2	instance_5OP48zv5rek
osg-hosted-ce-psc-bridges	uchicago-river-v2	instance_SW6qKz9cFVA
osg-hosted-ce-sdsc-triton-stratus	uchicago-river-v2	instance_0waaDcR5wiU
osg-hosted-ce-sut-ozstar	chtc-tiger	instance_5qgv0f6m9oQ
osg-hosted-ce-tcnj-elsa	uchicago-river-v2	instance_hLptxaFKTjI
osg-hosted-ce-tufts-cluster	chtc-tiger	instance_DNA800VQAIs
osg-hosted-ce-uci-gpatlas	uchicago-river-v2	instance_twHw8lU6_Zg
osg-hosted-ce-uconn-xanadu	uchicago-river-v2	instance_fI2zEGUbdWg
osg-hosted-ce-ucsd-comet	uchicago-river-v2	instance_o6038H1CDIg
osg-hosted-ce-usf-sc	uchicago-river-v2	instance_4riG7c9yTFA
osg-hosted-ce-uwm-nemo	uchicago-river-v2	instance_Zrk8YgF3yK8
osg-hosted-ce-wsu-grid	uchicago-river-v2	instance_oVjOC0nHnN8

osg-hosted-ce-tacc-frontera

osg-hosted-ce-tacc-stampede2

osg-hepcloud-ops

uchicago-river-v2

instance\_ZKBN7nCW1hI

osg-hepcloud-ops

uchicago-river-v2

instance\_MEsPx3\_7fqQ

A number of hosted compute elements have been deployed and are in production operation on the SSL in support of OSG (temporarily until PATH infrastructure is deployed at UC)

Demonstrating stable operation on flexible K8s substrate



infrastructure



# Infrastructure Planning

- Split SSL into stable production and development clusters (**done**)
- Supporting imminent release of ServiceX release candidate 1 (**done/on-going**)
- Supporting development of Skyhook integration from DOMA area on development Rook (Ceph) object storage system (**Rook deployed**)
- Invite k8s cluster partners



# Creating an SSL cluster pattern

- Want to establish a baseline for an SSL-like cluster
  - *More work for operators, but simplifies user expectations*
- Cherry pick best features from friendly clusters:
  - *GitOps (such as CERN, UW-Madison)*
  - *Sealed Secrets (such as UW-Madison)*
  - *Rook (such as UChicago, CERN, PRP)*
  - *User Portal (such as PRP)*
  - *Prometheus/Grafana (such as UChicago, CERN, PRP)*
  - *Federation / Admiralty ? (such as PRP)*
  - *SLATE (such as UW-Madison, UChicago)*



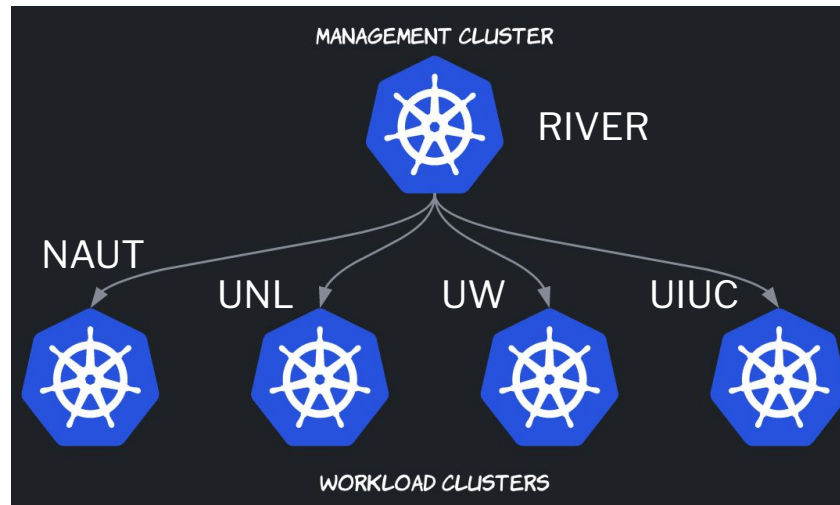
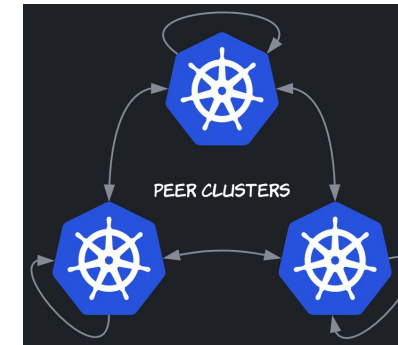
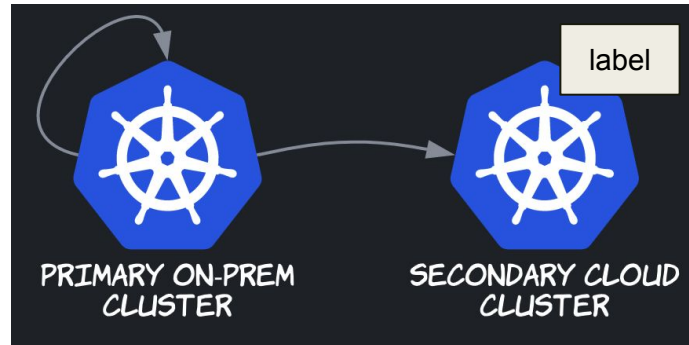
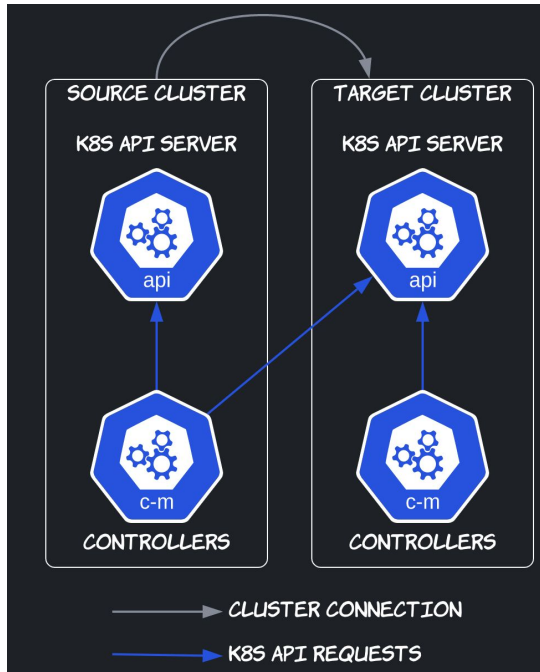
# Infrastructure R&D for potential Distributed Analysis

- Federation topology
- Authentication & authorization
- Scheduling applications

**There are obvious questions about how much to delegate to the infrastructure, how to manage access, which "applications" would benefit.**



# Potential topologies



Three preliminary meetings with interested **SSL cluster partners** to discuss the concept and challenges (cf <https://indico.cern.ch/category/11518/>)





# scalability testing



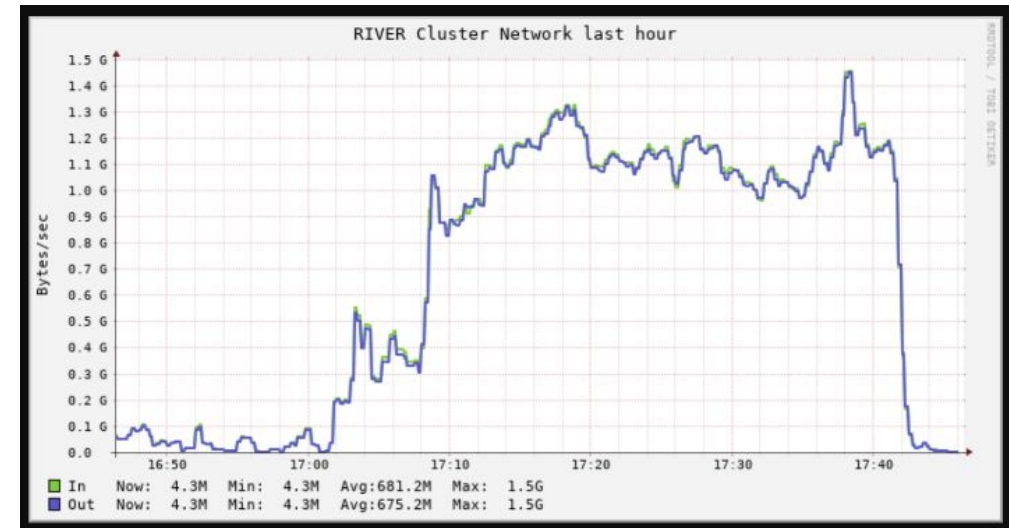
# Previous ServiceX Scalability Testing on SSL

- 10 TB xAOD input sample
  - Request 100 columns from 7 collections (~30% of file)
- Scaling ServiceX transformers (stable up to 1000)
  - Output rates of ~ 300 MB/s, total transform times < 30 minutes
- Finding bottlenecks
  - Have found/fixed issues with dCache I/O, Rucio access reporting, file transfer instabilities, race conditions in service

Queue 9c6a46d0-dedf-4490-b06d-a32c1c048b85

Overview

Queued messages last hour ?





# ServiceX Scalability Plan in 2021

## Phase space (432 tests)

- Instances (ATLAS, CMS, xAOD, uproot)
- Input dataset sizes (1, 10, 100 TB)
- Input origins (UC, US, grid wide)
- Concurrency (1, 5, 10 requests)
- Output sizes (0.1, 1, 10% of input size)
- Data analysis (concurrent, delayed)

Hard to do without a scripted testing and fully instrumented ServiceX and clients. Luckily instrumenting it is relatively easy to do.

## Monitors:

### Service side:

- ServiceX steps timings; transformers (CPUtime, walltime, memory); Bandwidth to components (server, MinIO, clients)
- MinIO, RMQ, PostgreSQL

### Client side:

- Time to complete
  - *Transformation only*
  - *Full analysis*
- CPU/Wall, memory



# Expected ServiceX Scale testing timeline & needed SSL readiness

## Instrumentation tasks: (complete by June 2021)

- . Prepare data collection in Elasticsearch
- . Instrument ServiceX components
- . Instrument clients
- . Create visualizations to understand results

## Scripting tasks: (complete by July 2021)

- . Dataset identification
- . Replicate to SSL staging locations
- . Provide user access (ATLAS/CMS)
- . Provide ability to execute transformation requests
- . Provide ability for users to test their analysis

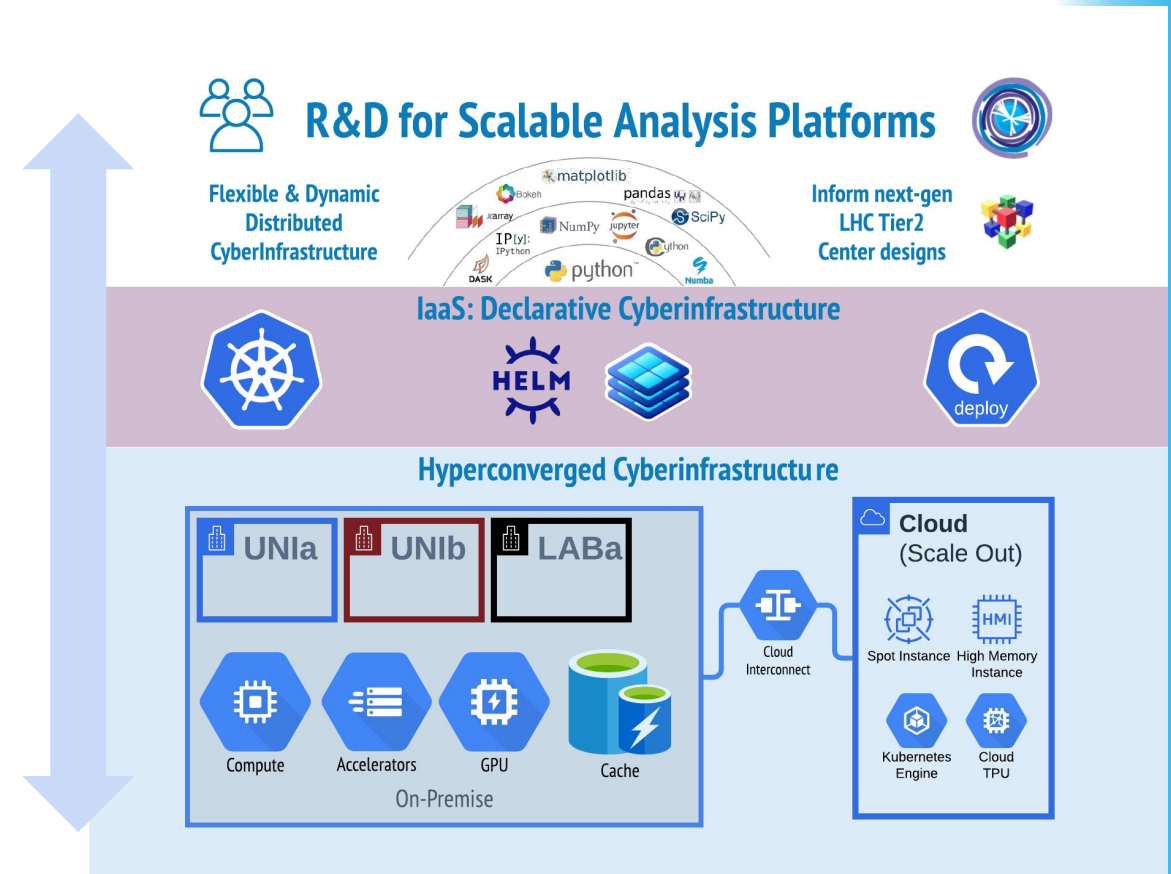


analysis grand challenge  
accelerated data delivery



# Supporting the Grand Analysis Challenge

- Fits with declarative CI
- **Milestone:**
  - *Prepare K8s clusters and other prerequisites including monitoring & analytics*
  - *Scale: 200 TB, 1500 core*
  - *With AS & DOMA, execute challenge*
- Initial focus will be on deploying identified application set from IRIS-HEP Analysis Systems blueprint







- ✓ 1. ServiceX / SkyHook
- ✓ 2. FuncX for FaaS
- ✓ 3. REANA
- ✓ 4. and a generic JupyterHub entry point that can access customized containers with the necessary software





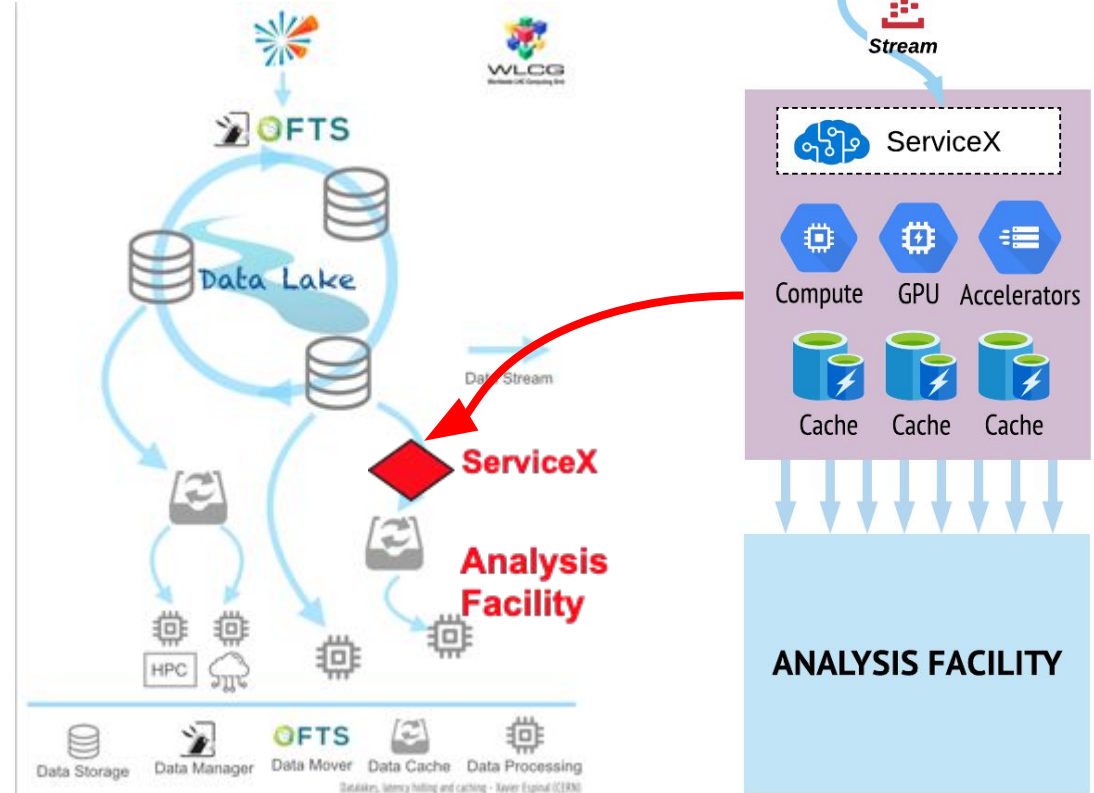
# Accelerated Data Delivery

High level goal is to explore (hardware) accelerating data delivery from ROOT format to columnar formats

Two milestones

- **Opportunity assessment** which profiles baseline transformer performance, evaluates technology options (Dec 2020)
- If cost/benefit relative to ServiceX on [cluster baseline](#) indicates, build **prototype system and benchmark its performance** (July 2021)

Unfortunately our graduate student decided to explore other CS research areas. So this activity is delayed until we find a new CS student.





# SSL Year 3 Summary

- SSL is a flexible resource for a variety of testing opportunities within IRIS-HEP **and partners**
- Now have production and development SSL clusters and **new cluster partners identified**
- Move towards an SSL cluster pattern using collective experience from the community
- Support functional and scalability tests, analysis grand challenge planning
  - *Prepare resources, provide access, instrument with analytics*

Questions?



Extra



# Cluster Partners

## SSL-UChicago-RIVER-dev: development RIVER cluster

- k8s version: v1.17.11
- cvmfs: yes
- Local filesystem: Rook
- Resources:
  - CPU: 432 cores (9 machines, 48 cores per machine)
  - Storage: 3.2 TB (rook/ceph)
- Pod types available: CPU
- Gaining access: same as SSL-UChicago-RIVER

## UNL

- K8s version 1.18.6
- CEPH storage (small, old hardware, ~ 4 TB only)
- CMS OAuth enabled access (to the JH resource)
- Jupyterhub powered with job scaling to the htcondor pool of our T2
- Hosts:

			cores	RAM
red-kube-vm00[1,2,3]	masters	VMs living on c07[14,16,18].shor	2	8GB
red-kube-c07[24,26,28,30]	workers	R710s with disks for Rook.io	24	96GB
red-kube-c10[35,36,37]	workers	Sun X2200	8	32GB
red-kube-c69[21-26]workers		Sun X2200	8	24GB
red-kube-c69[27-30]workers		4-in-2 Supermicro	16	64GB
red-kube-c6931	workers	1U Supermicro	8	32GB

## PRP Nautilus

- k8s version: 1.18.6
- Cvmfs: yes
- Local filesystem: Rook, SeaweedFS, BeeGFS
- Resources:
  - CPUs: 7000 cores
  - GPUs: 500+
  - Storage: 2.5+PB (rook/ceph), 2PB BeeGFS
- Pod types available: Any (CPU, GPU, FPGA)
- Gaining access: <https://ucsd-prp.gitlab.io/userdocs/start/toc-start/>
- OpenNSA controller L2

## UW-Tiger

- k8s version: v1.19.1
- Resources: constantly expanding, when I can get to it, lots of old, with more old and new on the way
  - As of Oct 13, 2020
    - Limited non shared persistent local storage
    - 368 cores (8x 40 old 1x48 new)
  - Soon:
    - 1200+ more old cores
    - 30 TB rook/ceph
    - 8x 48 new cpu
  - Slightly later:
    - More newer cores
    - 1.5 PB rook/ceph
- Gaining access: email BrianB and JeffP



# Cluster Partners

## UIUC-Boneyard

- k8s version: 1.18.8
- Compute hardware (~400 cores):
  - 22 Dell PowerEdge R410 (mixture of 16/24 cores, 2.54/2.8 GHz, 23.5/49.5 GB)
  - 4 Dell PowerEdge R710 (mixture of 16/24 cores, 2.5/2.7 GHz, 24.7 GB)
- Storage hardware (~80 TB)
  - 5 Dell MD1200 (~80 TB)
  - 2 Dell R510 w/12 3.5" drives
- Current configuration (subset of total to be brought online, as above)
  - 1 node (boneyard.ncsa.illinois.edu) for login & kubectl operation
  - 1 hardware/OS management node
  - 3 k8s control panel nodes
  - 3 k8 compute nodes
  - 1 storage nodes (w/ two connected enclosures for a total of 30 TB operational)
- Successfully deployed funcX service via k8, more testing underway
- Pod types available: CPU

## SSL-UChicago-RIVER: production RIVER cluster (production SSL cluster) - **for production phase**

- k8s version: v1.16.7
- cvmfs: yes
- nearby xcache: xcache.mwt2.org
- Local filesystem: Rook
- Resources:
  - CPU: 2784 cores (58 machines, 48 cores per machine)
  - GPU: none (note the fiona has been moved back to the UChicago cluster)
  - Storage: 23 TB (rook/ceph)
- Pod types available: CPU
- Gaining access: contact UC admins