# *Supporting a New Light Source at BNL SDCC*

*Tejas Rao (author) <raot@bnl.gov>*
*William Strecker-Kellogg (speaker) <willsk@bnl.gov>*

**HEPiX Spring 2021**

**BROOKHAVEN**
NATIONAL LABORATORY | Scientific Data and Computing Center

# NSLS-II

**Purpose**

To provide extremely bright x-rays for basic and applied research in biology and medicine, materials, chemical sciences, geosciences, environmental sciences, and nanoscience

**Sponsor**

U.S. Department of Energy (DOE), Office of Science, Office of Basic Energy Sciences
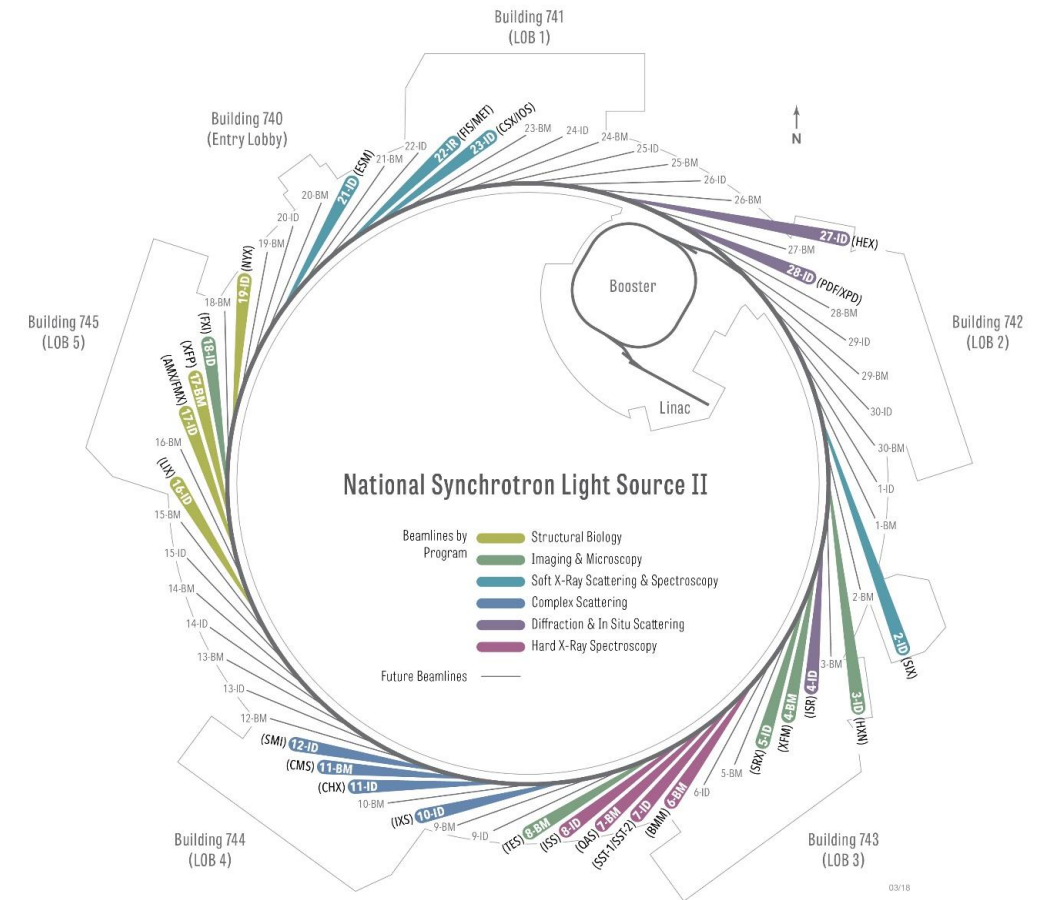
**Costs**

$912 million to design and build

**Features**

State-of-the-art, medium-energy (3 GeV) electron storage ring that produces x-rays up to 10,000 times brighter than the N

**Users**

Researchers from around the world.



National Synchrotron Light Source II

# NSLS-II

- NSLS II - Newest and Brightest Synchrotron in the world, supporting a multitude of scientific research in academia, industry and national security

- 32 beamlines in Operations or in Construction, facility has a capacity of 60 beamlines

  - High variability in computing needs across beamlines

- In the next few years, data rates and capacities are expected to grow exponentially

  - Difficult to forecast exact numbers, but O(100Tb/day) not unreasonable

- In fiscal year 2019, 1755 distinct researchers used NSLS-II beamlines for their research, and publications exceeded the level of one paper per day
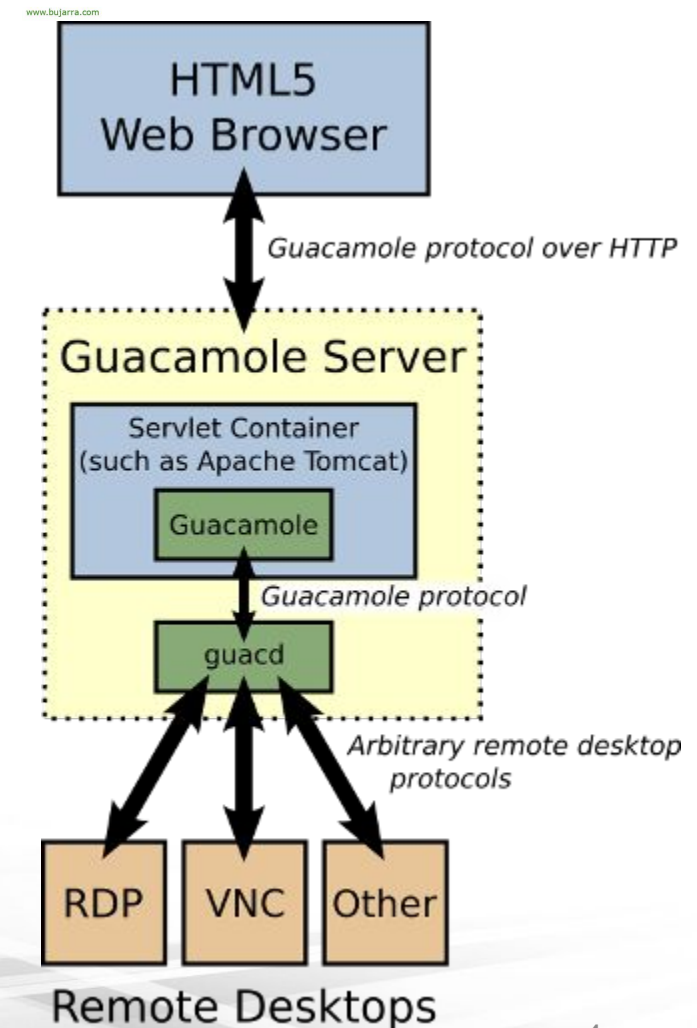
- **SDCC will support NSLS-II computing**

# Remote Desktop Gateway

Apache Guacamole™

- **Apache Guacamole** is a [clientless](clientless) remote desktop gateway. It supports standard protocols like VNC, RDP, and SSH on a single VM with 35 active users

- Apache Guacamole used by NSLS2 users to connect to beamline workstations.

    - **Authentication** is done directly by Guacamole application using SAML with BNL IDP.

    - **Authorization** is done using app-local MySQL database and defined by entries in the guacamole_connection_permission table.

- Plans to scale via load-balancing sessions through a proxy to a pool of backends
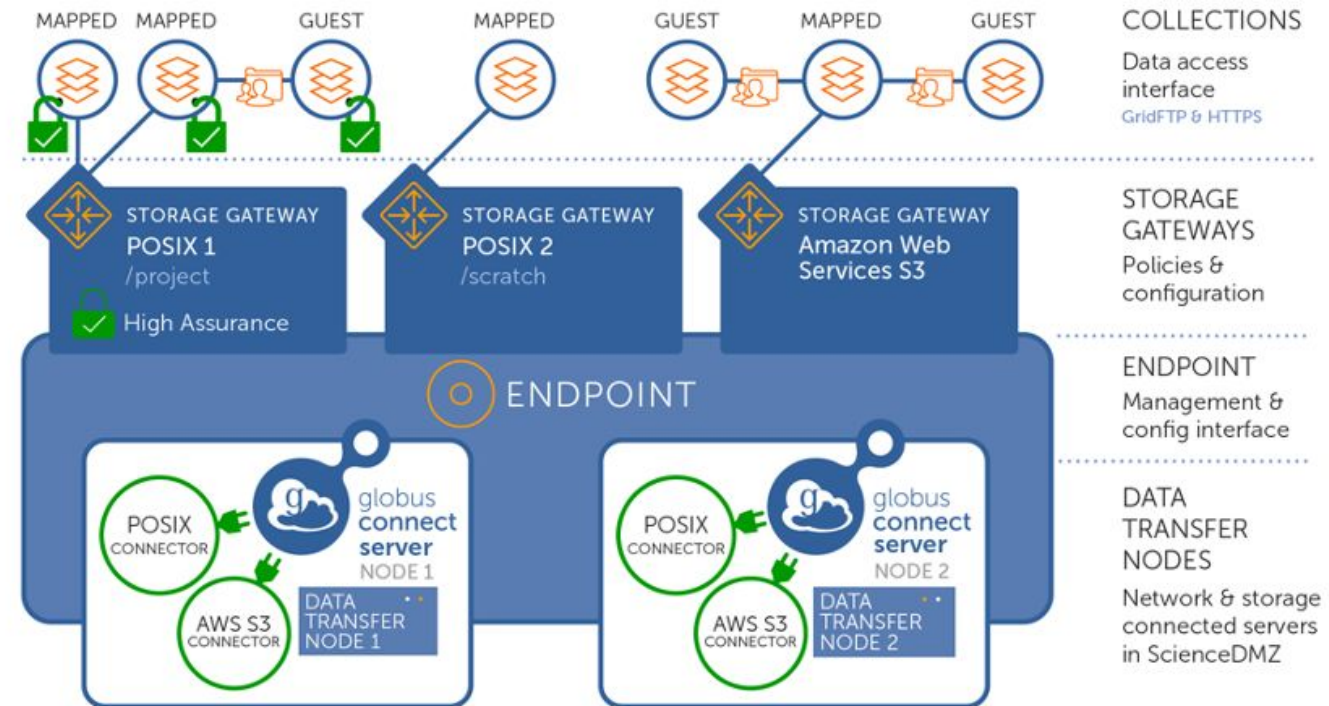
# Lustre storage

- Currently, Lustre filesystem with 4PB usable capacity serving all functional beamlines

- Lustre is back-ended by 2 X NetApp E5760 and 8 Dell R740 servers

  - HW Raid 6 (8+2), 420 drives total, FC connected

- Aggregate sequential performance is 42GB/sec (reads or writes)

- Using pacemaker/corosync for metadata(MDS) and data (OSS) failover

- Lustre also exported via **Globus**/Samba/SFTP/S3

# Globus Data transfers

➢ Lustre data filesystem and NFS home directories available worldwide via Globus endpoint "NSLS2"

➢ Oauth based access management architecture.

➢ Muti – DTN endpoints for load balancing and redundancy.

➢ Configurable Identity mapping system.

    ○ Map Identities from different domains to local accounts.

➢ Sharing of files is through mapped collections or guest shared collections.

➢ Shared site-license with central BNL IT

# Rucio Testing

- Rucio was tested as a data management solution for NSLS-II

  - Test instance of Rucio stood up at SDCC

  - Data registration integrated into bluesky suitcase

  - Documents created at beamline serialized and compressed then registered in Rucio

  - Replication rules in Rucio initiate and verify transfer to SDCC storage using Globus as the

    transfer agent

# NFS storage

- Dedicated NetApp A300 all *NVME* flash storage for home directories and software filesystems

  - In our experience NetApp all flash ONTAP storage are very reliable and performs well

- Protocols used are SMB/NFSv3/iSCSI

- AUTH_SYS authentication enabled which are required for supplemental groups with 16+ members

  - Server-side verification of user/group maps via AD

  - Multiple Centrify NIS servers used for directory searches (users/groups)

- Consistent performance with latency < 1ms and 180K random IOPS

# Filesystem Auditing

 Lustre and NFS filesystem user **AUTH_SYS** and rely on the client to state UID/GID which have potential risks.

- Kerberos or Shared secret key authentication not desired as it add complexity and it not convenient to use

- FS is mounted on user-controlled workstations so security is a concern

 Lustre changelog records real-time events that change the file system namespace or file metadata

- Contains all information necessary for auditing purposes

 Certain nodes can have audit logs disabled to avoid changelog flooding for nodes such as backup nodes and HSM

  agents

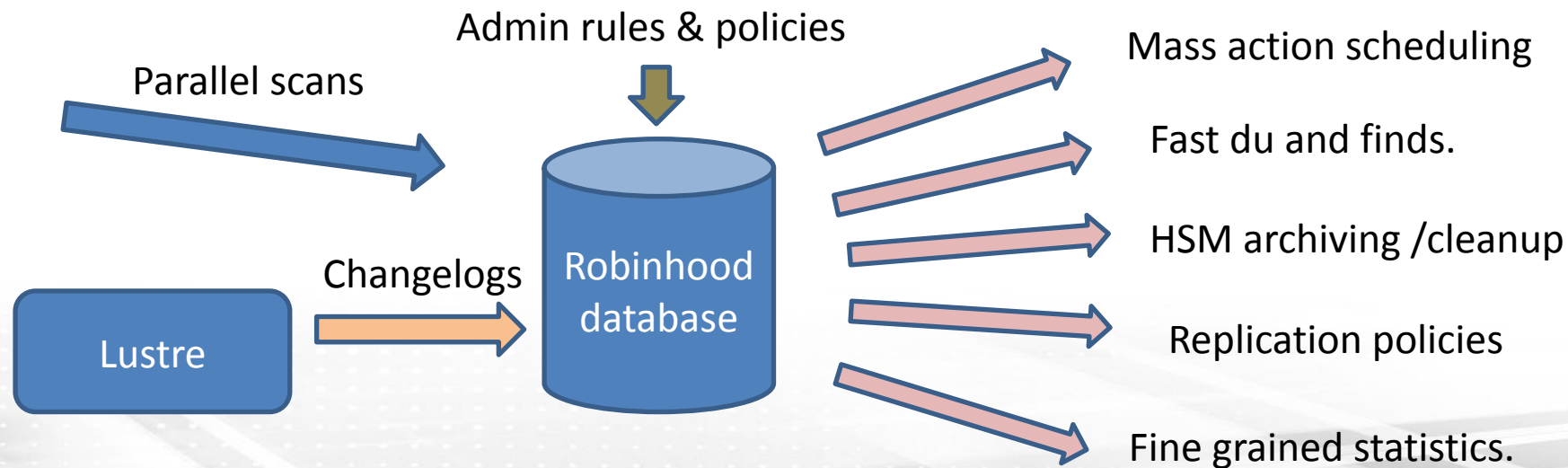Event   Timestamp                                        File identifier              Client

18662 01CREAT 12:57:01.783242245 2021.02.28 0x0 t=[0x200000bfd:0x23f2:0x0] ef=0xf u=0:0 nid=10.42.0.3@tcp p=[0x200000bfd:0x1:0x0] 483361
18663 06UNLNK 12:57:04.832648165 2021.02.28 0x1 t=[0x200000bfd:0x23dd:0x0] ef=0xf u=0:0 nid=10.42.0.3@tcp p=[0x200000bfd:0x1:0x0] 483361
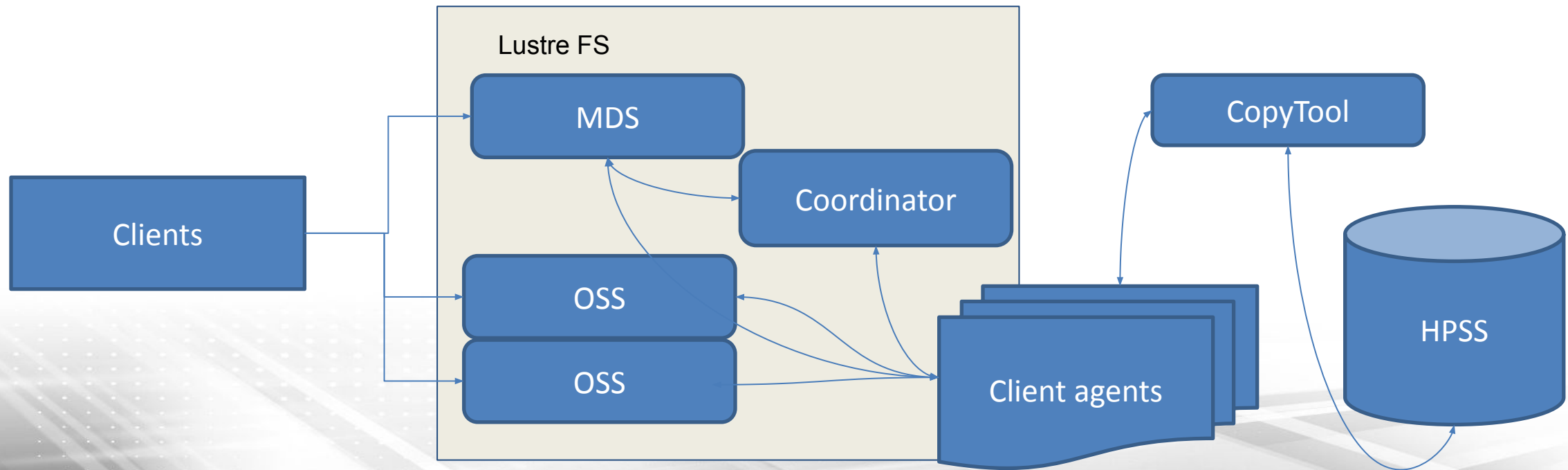
# Robinhood Policy Engine

- MySQL database ingesting all Lustre changelogs from MDS

- Aggregate statistics and customized reports can be made available in few seconds

- Complete Life cycle management of data including HSM data.

- Applying mass actions to filesystem entries quickly using various criteria and actions.

- Lustre changelogs used to keep Robinhood database updated in near real time.

Admin rules & policies

Parallel scans

Lustre

Changelogs

Robinhood database

Mass action scheduling

Fast du and finds.

HSM archiving /cleanup

Replication policies

Fine grained statistics.

# Robinhood and Hierarchical Storage Management

- ◻ Will be using Lustre CopyTool to archive files into IBM HPSS in the future
  - ○ CopyTool: A user-space program managing data between Lustre and the HSM
    - ■ Uses HPSS-API directly
- ◻ External PolicyEngine, like Robinhood, decides what to archive, what needs to be released from Lustre
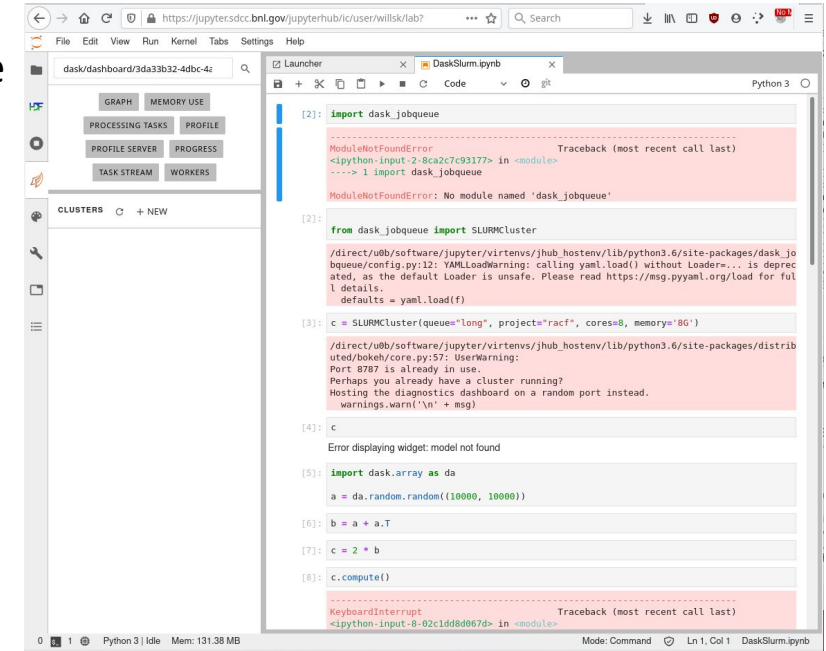- ◻ Trigger file migration depending on policy rules (for Lustre HSM)

# Active Directory & Centrify

- ➢ NSLS2 preferred to use BNL's Active Directory for account management, over SDCC's existing LDAP/IPA setup
  - ○ Wanted to be able to utilize the same system for authorization/authentication at their beamlines
    - ■ Including access to Windows workstations
  - ○ NSLS2 also required fine-grained control over access to particular systems
- ➢ Centrify DirectControl chosen as the solution for integrating Linux hosts with AD
  - ○ Full agent/installation only running on hosts where direct authentication needed
    - ■ *PAM auth like SSH-logins*
    - ■ CentrifyDC Provides PAM and NSS modules
  - ○ Other hosts setup to bind to adnisd NIS server for nss passwd/group maps - authorization only
    - ■ For example on batch nodes and IDP-web-authenticated systems like Jupyter
  - ○ Fully puppetized CentrifyDC and Centrify NIS client deployment

# Jupyterhub

☐ Migrating existing small condor-backed jupyter platform to a new Slurm-backed one

with ActiveDirectory users and login based on IDPs with MultiFactor Auth

    ○ *MFA is a requirement for web-based "interactive" access @ BNL*

☐ **Delegated Administration**

    ○ Allow NSLS2-designated software librarians to maintain the whole stack

        ■ Jupyterhub runs unprivileged

           ● Use of sudo-based "*batchspawner*" allows hub to run with

            minimal privilege -- only needs to be able to:

            `sudo -u $USER sbatch $PAYLOAD`

        ■ Separate JupyterLab environment, and the kernel environments

        managed on NSLS-2 filesystems

☐   All behind our proxies providing a common URL namespace and SSL termination

# Slurm Cluster

➢ 30 nodes cluster

    ○ 12 nodes with 2x V100 GPUs

    ○ 18 nodes CPU only (Cascade Lake)

    ○ EDR IB connected (single switch)

    ○ <u>768GB</u> RAM / <u>48</u> Physical Cores

➢ Hardware / Batch Configuration

    ○ 2 master nodes (slurmctld, slurmdb)

    ○ 2 interactive (ssh) submit nodes

    ○ 2 Partitions

        ■ **debug** partition 1 GPU node (10 minutes)

        ■ **long** partition 25 nodes (14 CPU, 11 CPU+GPU) (4 days)

➢ **User Administration**

    ○ Delegated to an NSLS2 staff administrator who has the slurm admin role

        ■ Able to modify the slurm cluster as they see fit (modify/add/remove qos, accounts, users)