# NERSC (storage) experiences - WLCG/HSF 2020

Thanks to Glenn Lockwood, Debbie Bard, Annette Greiner, Bjoern Enders and Lisa Gerhardt for slides

Wahid Bhimji

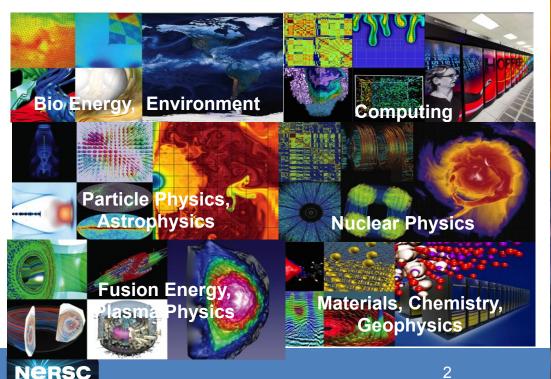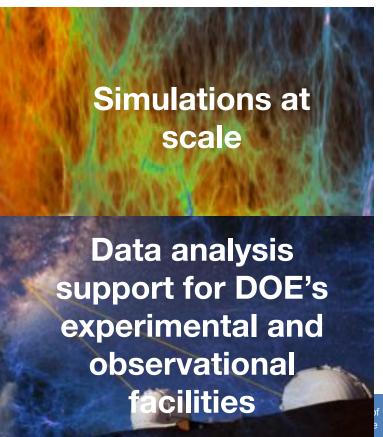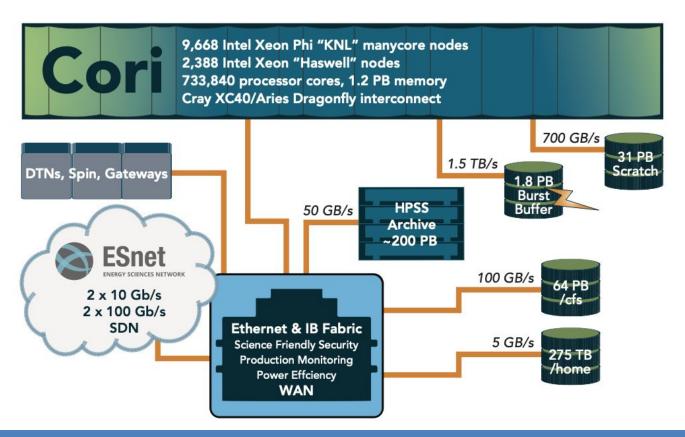Data and Analytics Group, NERSC, Berkeley Lab

November 24, 2020

# NERSC is the **mission** HPC and Data facility for the Dept Of Energy Office for Science

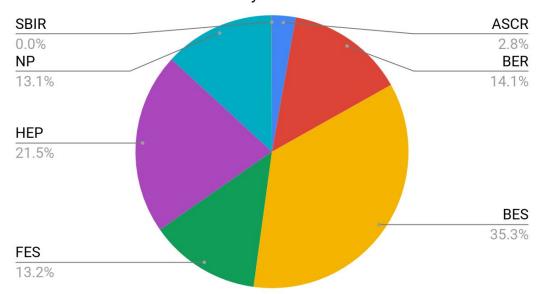**7,000+ Users, 800+ Projects** *2000+ citations /yr*



Bio Energy, Environment
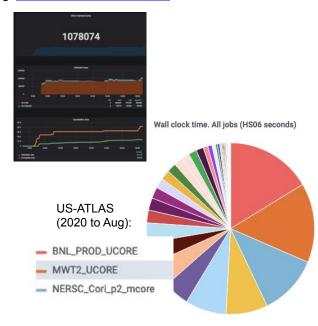
Computing

Particle Physics, Astrophysics

Nuclear Physics

Fusion Energy, Plasma Physics

Materials, Chemistry, Geophysics

Simulations at scale

Data analysis support for DOE's experimental and observational facilities

# NERSC Centre 2020

# NERSC mostly non-HEP/NP - but valuable HEP resource

E.g. NoVA '1m cores' news article



Percent of NERSC-Hours Used By Office in Allocation Year 2019



- SBIR 0.0%
- NP 13.1%
- HEP 21.5%
- FES 13.2%
- BES 35.3%
- BER 14.1%
- ASCR 2.8%

1078074

Wall clock time. All jobs (HS06 seconds)

US-ATLAS (2020 to Aug):

- BNL_PROD_UCORE
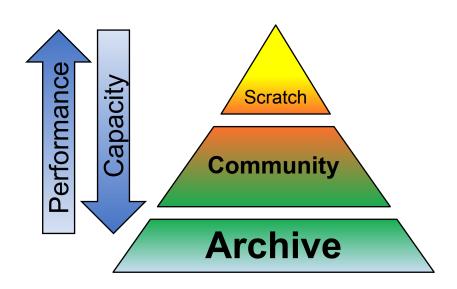- MWT2_UCORE
- NERSC_Cori_p2_mcore

# NERSC perspective on storage hierarchy



- **Scratch (weeks – months)**
  - Mounted on only one HPC system
  - User data purged after 4-12 weeks
- **Community (months – years)**
  - Mounted center-wide (HPCs, web, k8s)
  - Quotas
  - User data archived at project end
- **Archive (years – decades)**
  - Not "mounted" anywhere (object-like)
  - No effective quota

# NERSC's storage hierarchy - today



Performance ↑

Capacity ↓

- Burst Buffer
- Scratch
- Community
- Archive

**1.8 PB**
**1.5 TB/s**

**30 PB**
**700 GB/s**

**64 PB**
**150 GB/s**

**230 PB**
**~50 GB/s**

NERSC

U.S. DEPARTMENT OF ENERGY | Office of Science

# NERSC Systems Roadmap



**NERSC-7:**
Edison
Multicore
CPU

**2013**

**NERSC-8: Cori**
Manycore CPU
NESAP Launched:
transition applications to
advanced architectures

**2016**

**NERSC-9: Perlmutter**
CPU and GPU nodes
Continued transition of
applications and support for
complex workflows

**2020**

**NERSC-10:**
Exa system

**2024**

**NERSC-11:**
Beyond
Moore

**2028**
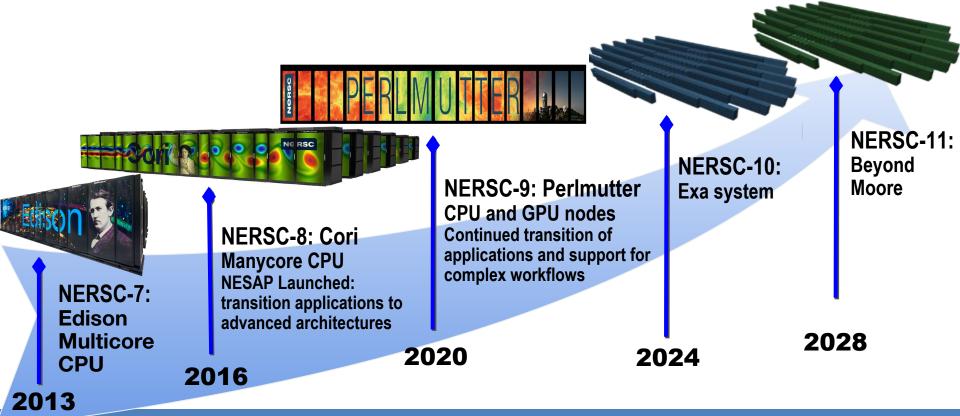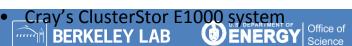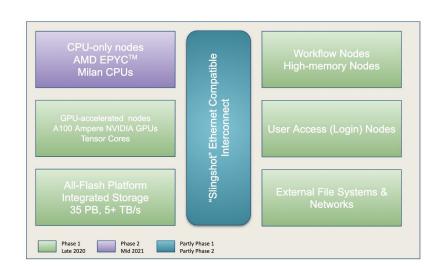
# Perlmutter

- Cray Shasta System providing 3-4x capability of Cori

- GPU-accelerated and CPU-only nodes
  - Large CPU-only partition providing capability similar to Cori
  - Services for complex workflows
  - Optimized data software enabling analytics and ML at scale

- GPU nodes: 4 NVIDIA A100 "Ampere" GPUs each w/Tensor Cores, NVLink-3 and High-BW memory + 1 AMD "Milan" CPU
  - Over 6000 GPUs
  - Unified Virtual Memory support improves programmability

- Cray "Slingshot" - High-performance, scalable, low-latency Ethernet- compatible network
  - Capable of Terabit connections to/from the system

- Single-tier All-Flash Lustre based HPC file system
  - 6x Cori's bandwidth
  - Cray's ClusterStor E1000 system



Phased delivery
1st phase: Early 2021
2nd phase: Summer 2021



BERKELEY LAB   U.S. DEPARTMENT OF ENERGY | Office of Science   NeRSC

# NERSC's storage hierarchy - soon



**35 PB**
**5.0 TB/s**

**128-192 PB**
**300-450 GB/s**

**300-500 PB**
**~50 GB/s**

Performance

Capacity

Scratch

Community

Archive

# Data Management at NERSC



Performance →

Capacity →

Scratch

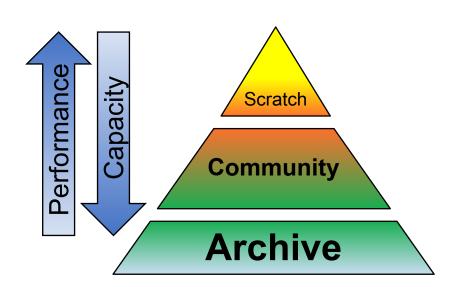Community

Archive

**30% of Globus data movement is purely internal**

**Community File System:** Intended for sharing scientific data between groups of NERSC users. 27 PB and 3 billion inodes used

Very useful, but has some issues
- Group quotas that fill up
- Permissions drift
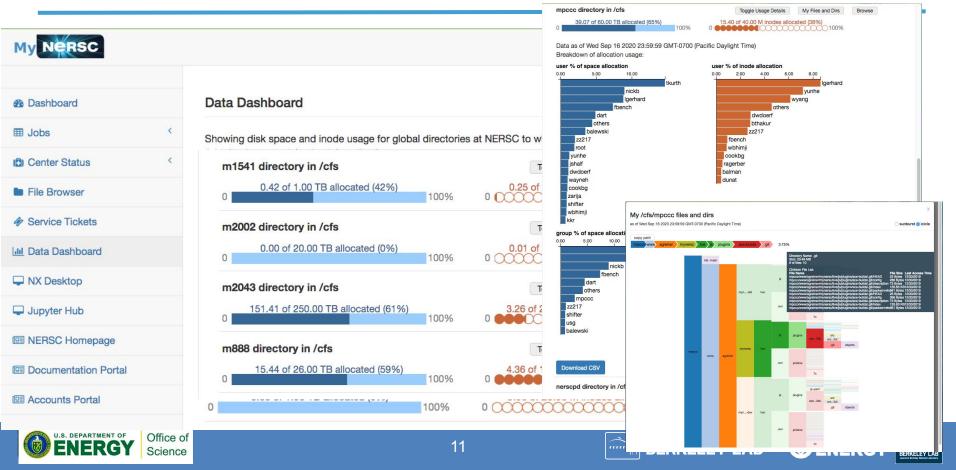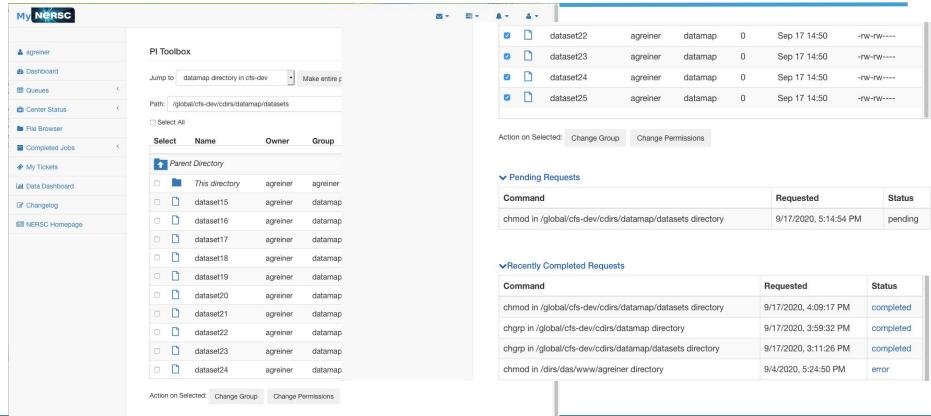- Data migration between the tiers

| Fixed Globus Metrics | Globus Internal Trans... | Globus Internal Data Tra... | Globus Collaboration En... |
|---|---|---|---|
| **30,377,622.13** Total Data Transferred [GB] | **13,965,183.829** Total Data Transferred Internally [GB] | **11,926,864.919** Total Data Transferred Internally to/from Cori Scratch [GB] | **816,667.051** Total Data Transferred Using Collab Endpoint [GB] |
| **785** Total Users | | | |
| **31.95** Average Unique Users Per Day | **145** Total Users | **99** Total Users | **10** Total Users |

U.S. DEPARTMENT OF ENERGY | Office of Science

BERKELEY LAB

U.S. DEPARTMENT OF ENERGY

BERKELEY LAB

# Data Dashboard: Helping Users Share Their Space

# PI Toolbox: Permission Wrangling
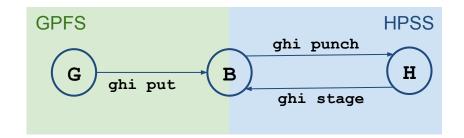
# **G**PFS-**H**PSS **I**nterface

Use HPSS using the familiar file system interface

| | |
|---|---|
| `ghi ls` | show what file system the files are currently on, files are marked 'G' for GPFS, 'H' for HPSS and 'B' for both |
| `ghi put` | copy the files to HPSS, makes files dual residents on both files systems |
| `ghi stage` | move the files back from HPSS to GPFS |
| `ghi punch` | move the files to HPSS (leaving a stub behind on GPFS) |
| `ghi pin` | keep the files from being removed from GPFS |
| `ghi lock` | keep the files from being removed using rm or otherwise modified |

- Users interact with a file system directory that's tied to HPSS behind the scenes

- GHI puts the files optimally into HPSS on your behalf, no need to htar small files together or break things into 500GB sized chunks

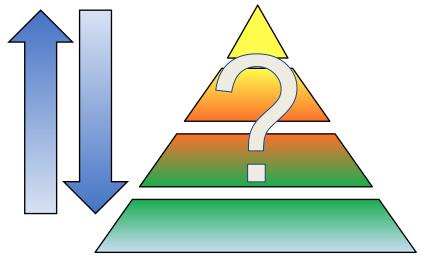- Shared namespace, so deleting a file from GPFS will remove it from HPSS
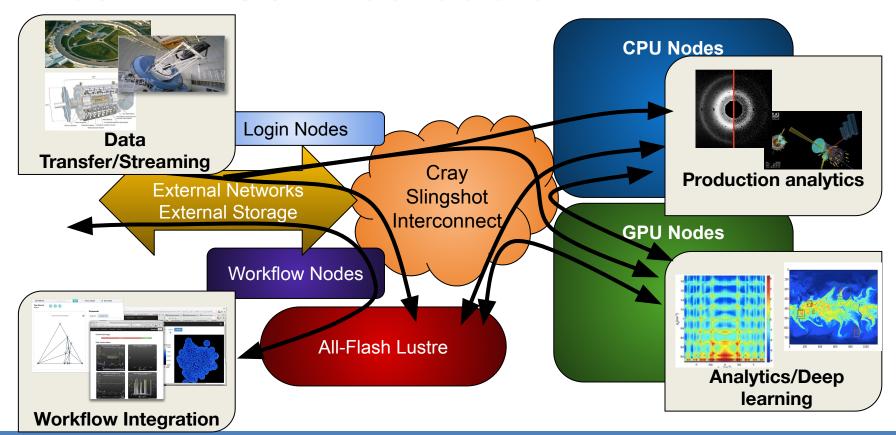
# NERSC's storage future



- Various demands on storage e.g.
  - Write-heavy high-bw HPC
  - Read-heavy high-iops Deep learning and analytics
  - Metadata management
  - Workflow chaining
- Various storage innovations e.g.
  - Object-store access
  - Filesystems on demand
  - Very high IOPS NVRam
- Challenges include
  - Stability and reliability
  - Multi-user HPC integration

# Wider NERSC infrastructure

# Wider NERSC data services

Scientific Instruments

Scientist Interaction

Modified from O'Reilly Blog

**Data transfer**
Globus
GridFTP
Xrootd

**Portals**
Newt
Django
Jupyter

**Analytics**
Python
ROOT
R
Matlab

**File Format**
HDF5
NetCDF
ROOT

**Storage**
Lustre
GPFS
DataWarp
HPSS

**Workflow**
Fireworks
Taskfarmer
Dask

**Batch Processing**
SLURM
Spark

**Machine Learning**
TensorFlow
PyTorch
Scikit-Learn

**Databases**
MongoDB
Posgres
MySQL

**Visualization**
Visit
ParaView
Matplotlib

Y LAB

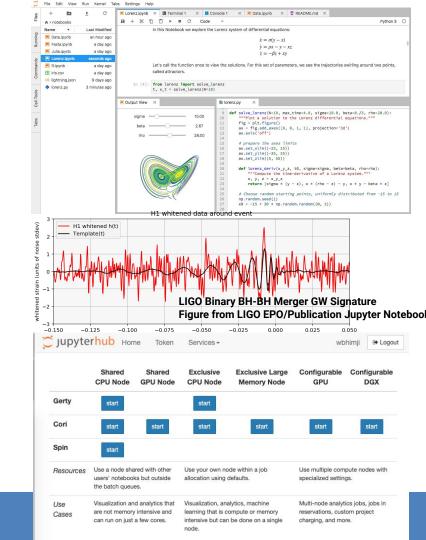U.S. DEPARTMENT OF ENERGY | Office of Science

# Jupyter Notebooks for HPC



- **Jupyter growing in popularity at NERSC and broader community**
  - **> 700 unique users @NERSC**
- **NERSC's goal - Enable exploratory data analytics, deep learning, workflows, through Jupyter on HPC**
- **New features in 2019-2020:**
  - **Access to Cori compute nodes**
  - **Access to Cori GPUs**

**Go to [https://jupyter.nersc.gov](https://jupyter.nersc.gov)**

**New User Training talk**

LIGO Binary BH-BH Merger GW Signature
Figure from LIGO EPO/Publication Jupyter Notebook

# Conclusions

- NERSC supports an increasing number of data-rich projects across difference science domains

- Upcoming upgrades to center systems and storage
  - All-flash scratch tier with Perlmutter
  - Large capacity "community" filesystem

- Deploying tools to ease data management within the center and beyond

**Perlmutter**

**>16x MDS + >270 OSS**
1x AMD Rome
2x Slingshot NICs
24x 15 TB NVMe

**>1,500 GPU nodes**
1x AMD Milan
4x NVIDIA A100
4x Slingshot NICs

**Slingshot**
200 Gb/s
2-level dragonfly

**24x Gateway nodes**
2x AMD Rome
2x Slingshot NICs
2x Mellanox CX6 VPI NICs

**Community
(GPFS)**

SAN

SAN

**>3,000 CPU nodes**
2x AMD Milan
1x Slingshot NIC

**2x Ethernet routers**
400 Gb/s/port
> 10 Tb/s routing

WAN

**Experimental &
Observational
Facilities**
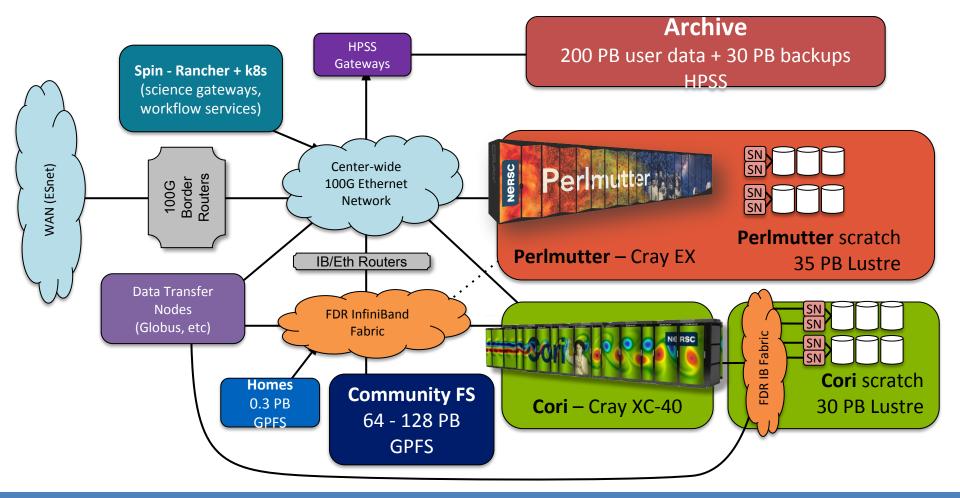
# Federated Identity (FedID) allows a person to use a single digital identity across multiple organizations

- Simplifies cross-facility workflows
- Users have fewer, more familiar, passwords and login pages
- NERSC has fewer support tickets (eg, password resets)
- Home institution manages account lifecycles
- NERSC still manages local authorization
- Core technology is well-established and mature
- *Policy/trust decisions were the bulk of our analysis*

# Spin: Container Services for Science

**Many projects need more than HPC.**

***Spin is a platform for services.***

Users deploy their **science gateways, workflow managers, databases, and other network services** with Docker containers.

- *Access HPC file systems and networks*
- *Use public or custom software images*
- *Orchestrate complex workflows*
- *Secure, scalable, and managed*

**Some projects using Spin:**

| | |
|---|---|
| Track and compare analyses of nightly sky surveys | science gateway |
| Classify and store reusable earth sciences data | data repository |
| Manage production genomic workflows and data at scale | science gateway |
| Process real-time events for dark matter detection | workflow manager |
| Explore materials properties or build simulated materials | science gateway |

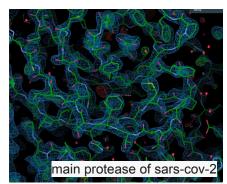# LCLS-II is using NERSC for real-time data analysis

- Several experiments at the LCLS-II (x-ray free electron laser at SLAC) are now using NERSC for real-time data analysis for materials science and Covid-19 research



Say hello to Tethrene!

- Can analyze a 5 minute experiment in ~3 minutes for feedback to beamline staff, transferring 15TB/day to NERSC

  o **Real-time** data analysis using real-time queue and advanced reservations

  o Used services running on **Spin** to orchestrated jobs/parameters/results in real time between several concurrent remote users



main protease of sars-cov-2



ESnet data rate copying data from LCLS to NERSC -- spikes are runs being transferred in real time

# Collaborative Distributed Data Analysis with **Spin**



Incoming data

Monitor runs

Monitor analysis

**Science!**

Say hello to Tethrene!

Monitor experiment

Submit jobs

cctbx.xfel

# Machine-readable supercomputers: the Superfacility API

**Vision: all NERSC interactions are callable;
backend tools assist large or complex operations.**

**Endpoints currently prototyped:**

| | |
|---|---|
| `/accounting` | retrieve allocation info for a user or project |
| `/auth` | obtain OAuth2 authentication tokens (JWTs) |
| `/callbacks` | register callbacks for asynchronous/chained operations |
| `/file` | browse, upload, and download files |
| `/health` | retrieve system health status |
| `/jobs` | submit jobs and check job status |
| `/transfer` | move data with Globus or between NERSC storage tiers |
| `/reservations` | submit and manage future compute reservations |

**Superfacility API** 1.0
[ Base URL: /api/v1 ]
/api/v1/swagger.json

API access to NERSC

**SFapi**

**auth** JWT token creation, verification and revocation

POST `/auth/login`

POST `/auth/revoke`

POST `/auth/verify`

**file** basic file browsing, upload and download of small files to and from NERSC

PUT `/file/{machine}/{path}`

GET `/file/{machine}/{path}`

**accounting** Get accounting information about the user's projects

GET `/accounting/projects`

GET `/accounting/projects/{repo_name}/jobs`

GET `/accounting/roles`

**callbacks/callbacks** Manage workflow reservations at NERSC

GET `/callbacks/callbacks/` This api requires authentication

POST `/callbacks/callbacks/` This api requires authentication

**https://api.nersc.gov/**

ESnet ENERGY SCIENCES NETWORK

CRD COMPUTATIONAL RESEARCH DIVISION

NERSC

# The Superfacility API: sustainable, scalable automation

- Less user/staff DIY: simpler, standardized tooling (Python, etc)
  - Stable refactor target for established projects
  - Easier on-ramp for new projects
- Fit (not fight) standard software design patterns
  - Shared libraries and API calls
  - Authentication and security models built on OAuth2 Standard and JSON Web Tokens (JWTs)

*using the API from a Jupyter notebook to check Cori status*

Before we start any computing, let's check whether Cori is up.

```
[3]: health_cori = api("health/resource_statuses/cori", data={"notes":"false", "outages":"true"}, as_form=True)[0]
     print("Cori is %s" % health_cori['status'])

     Cori is active
```
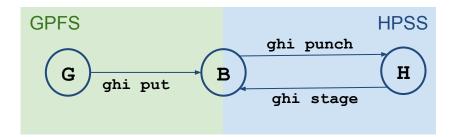
We can also take a look into the future to better plan our work around planned outages.

```
[4]: planned_outages = [o for o in health_cori['outages'] if o['status'].lower()=='planned']
     print(planned_outages)   #make this nicer

     [{'startdate': '2020-05-20T05:00:00', 'enddate': '2020-05-20T19:00:00', 'description': 'Scheduled Maintenance', 'notes': 'ExVivo and CGPU resources will be unavailab
     le during this maintenance.', 'status': 'Planned', 'swo': 'true', 'identifier': 'QXg7SbWP3KAeG0mwkyQS', 'updatedate': None}]
```

Jupyter  SFapi

# **G**PFS-**H**PSS **I**nterface

exemplary use cases



- **Archiving Complex Directory Structures**
  - Experiments [at the ALS at LBNL] often have complex sets of [microscopy image] data, i.e. large volumes of image data in tens of MBs along with a few kB-sized text file.
  - The files are organized in a complex directory structure that must be maintained.
  - Use a single command (`ghi put`) to archive the entire directory into HPSS.
- **Large volumes of Infrequently Accessed Data**
  - 100TBs of data only accessed a few times a year for reanalysis (DESI)
  - Use `ghi put` to put the data into HPSS
  - In between analyses, frees up disk space by using `ghi punch` to move most data off of GPFS but still leaves a browsable directory structure behind
  - Selectively retrieves some (or all) files with `ghi stage`