# subMIT Overview

Josh Bendavid
Basic Computing Services (subMIT) Review
June 21, 2024

# Introduction

- subMIT system provides an interactive login pool + scale-out to batch resources
  - Home directories
  - Convenient software environment (Alma Linux 9 native, docker/singularity images, conda)
  - SSH or Jupyterhub access
  - Local batch system with O(1000) cores, >50 GPU's
  - Additional storage for software installation/development, large datasets
  - Convenient access to larger external resources (OSG, CMS Tier-2 and Tier-3, LQCD Cluster, EAPS)
- User support is a key feature of the system
  - Beyond basic troubleshooting
    - Help users make optimal use of the available resources
    - Expert advice on designing/improving workflows
    - Customize and evolve system configuration to accommodate user needs as appropriate

# Introduction

- Storage and networking
  - 500TB of spinning disks
  - Local storage (1TB/user), 10's of TB for larger group datasets
  - 40TB of ultra-fast NVME storage with room for future expansion
  - Fast networking: 100 Gbps ethernet
    - RoCE (RDMA over Converged Ethernet) has been partially tested/commissioned, should be possible for MPI applications
- System is located in B24 basement, with 100gbps uplink
- Additional resources recently integrated
  - More disk storage (100TB contributed from ABRACADABRA)
  - Integration of existing computing resources from research groups
  - Purchase of several large core count/high memory machines by research groups for additional computing resources and to support specialized workflows and/or R&D where large single node scaling is useful
    - Current "high density" template, Dual AMD EPYC 192 core/384 thread with 1.5TB of memory

# Introduction: subMIT Website



- ● Website (with User's Guide/Instructions):
  https://submit.mit.edu/
  - ○ Overview and general information
  - ○ Direct JupyterHub access
  - ○ User's Guide:
    https://submit.mit.edu/submit-users-guide/

# Introduction: Project Organization

- Formally the project is organized as ***Basic Computing Services*** in the Physics Department
  - **Project Team:** Implementation/Operations/Maintenance of the system
  - **Users Group:** Contact point between the user community and the project team, forum for user feedback, requests, information flow to and from users
  - **Steering Committee:** Faculty oversight, funding, etc
  - See https://submit.mit.edu/?page_id=6

# Users Group In Practice

- Regular meetings (every few months)
  - Advertised and open to the broader community
  - Topical presentations from project team, Users Group representatives, or other users or community members
  - Forum for feedback and information flow between the user community and the project team
  - Regular timeslot: Tuesday 10:00-11:00 EST
  - Last meeting agenda: https://indico.mit.edu/event/1050/
- Users Group representatives
  - Identified representatives from research groups across the department
  - Attend the monthly meetings
  - Provide feedback from your groups/community
  - Distribute information/news from the project team

# Users Group Representatives

- Users group has been formed (JB as coordinator)
- Current Users Group representative (associated faculty/group)
  - Yin Lin (Phiala Shanahan)
  - Siddharth Mishra-Sharma (Jesse Thaler)
  - Prajwal Mohan Murthy (Bob Redwine)
  - Kaliroë Pappas (LNS Neutrino/Dark Matter)
  - Amer Al-Hiyasat(Julien Tailleur)
  - Yitian Sun (Tracy Slatyer)
  - Molly Taylor (LNS Heavy Ion Group)

# Users Group Meetings

- E.g. presentation from April Users Group meeting on Visualization of Geant Simulations on subMIT
- https://indico.mit.edu/event/1050/

**Crossover in finite sized systems**





**Users Group Meeting**

📅 Tuesday Apr 30, 2024, 10:00 AM → 11:00 AM America/New_York
📍 Kolker Room (26-414) (MIT)

Description   https://mit.zoom.us/j/96743699673?pwd=b3h2Q3c3cVQwYW12blhMUG5SWXZCZz09

**10:00 AM** → 10:15 AM   **subMIT status and updates**   ⏱15m
Speaker: Joshua Bendavid (Massachusetts Institute of Technology)
📄 subMIT status User'...

**10:15 AM** → 10:30 AM   **Studying Diffusion with Conservation of Center of Mass Using subMIT**   ⏱15m
Speaker: Sunghan Ro (MIT)
📄 24_04_SubMIT_v2.p...

⏱15m
...pas (MIT laboratory for nuclear science), Molly Taylor (Massachusetts Institute of Technology), Prajwal Mohan Murthy (MIT LNS), ...ma (MIT), Sunghan Ro (MIT), Yin Lin (Massachusetts Institute of Technology), Yitian Sun (Massachusetts Institute of Technology)

⏱5m

... thanks subMIT!

8

# Storage breakdown

- Several different storage areas are available covering different use cases
  - /home/submit/<username>
    - Home directories (nfs server), redundant disk array with backups
    - 5GB quota
    - Use for software development and (small) critical data
  - /work/submit/<username>
    - Work directory (nfs server), no backups (but redundant disk array)
    - 50GB quota
    - Use for software installation (conda or docker/singularity images)
  - /data/submit/<username>
    - Large distributed disk system, no backups, but redundancy against disk failure ("erasure coding")
    - 1TB user quota, larger quotas available in dedicated group directories
    - Store large datasets here
  - /scratch/submit/<username>
    - Fast NVMe SSD array
    - Commissioned by several groups for high performance data analysis
  - /cvmfs/
    - Read-only distributed storage for distributing software, singularity images, etc
    - Several CERN-related repositories are available
    - Local repository /cvmfs/cvmfs.cmsaf.mit.edu where additional software or data can be added if needed
- Flexible tiered storage system, can accommodate a wide range of user needs
- Larger datasets encouraged to use shared group space, but quotas can be increased when needed

# Interactive Use: Terminal or JupyterHub



- Interactive Jupyter session available directly from website with touchstone authentication (subMIT account still required)
- SLURM is used to efficiently share resources between interactive and batch use
- Primary usage is research, but also being used for classroom exercises

# Communication Channels

- User support mailing list: submit-help@mit.edu
- Experimental large language model application under development for interactive user support and to augment support ticket handling
  - Joint project with College of Computing, with dedicated funding
  - More discussion later + dedicated talk at LLM workshop on Friday https://indico.mit.edu/event/759/
- Slack workspace: https://mit-submit.slack.com
  - "help-desk" channel
- Monthly Users Group Meetings
  - Open for discussion
  - Open for user contributions: full set of Users Group representatives can be contacted at submit-usersgroup@mit.edu
- Annual subMIT workshop
  - February 2024 workshop: https://indico.mit.edu/event/956/
- In addition to direct interaction with the subMIT project team, users are encouraged to discuss with Users Group representative from their own group or "nearby" group

# Linux Distribution Upgrade

- Previous CentOS 7 distribution reaches EOL for maintenance updates in June 2024
- Decision by Red Hat to reorganize CentOS project and releases disrupted the logical upgrade path from CentOS 7->8
- Decision was taken to upgrade from CentOS 7 to Alma Linux 9, considering:
  - Ease of transition
  - Support lifetime
  - Functionality
  - Direction being taken at other universities and labs (CERN, Fermilab, etc)
- Discussion included Users Group and broader community
- System has been fully upgraded, with user facing services (ssh, Jupyter, batch queues) switched to Alma 9 by default in April
- Dedicated documentation to ease user transition
  - https://submit.mit.edu/submit-users-guide/future/alma.html
- Well-supported and documented use of containers to keep older software environments available where needed

# Mass Storage Upgrade

- Current mass storage system (/data/submit) is 500TB of spinning disks in a Gluster distributed filesystem
- **Users experience throughput bottlenecks for high performance analysis**
  - unable to effectively leverage the throughput of a large number of disks in parallel
- **Some other performance and operational issues related to user access patterns (system responsiveness/reliability problems for users)**
  - Large number of small files, large number of files in a single directory
  - Technical protections and user education to mitigate
- A number of limitations in maintenance and flexibility, plus suboptimal failure modes (files appear to have vanished when they are only temporarily offline)
- These issues also drive extra demand for scarce/expensive NvME storage (/scratch area)
- Migration in progress to higher performance CephFS

# Mass Storage Upgrade

- Migration in progress to higher performance CephFS
- Test system was deployed using spare disks for initial testing/planning
- Storage servers hosting existing gluster system have been upgraded to Alma 9 as well
- Production Ceph system deployed in parallel on existing servers, 216TB of spinning disks from spares + disks borrowed from Tier 2 project
- Additional disks have been ordered (20x20TB) from funds allocated for storage upgrade hardware
- Migration plan:
  - Commissioning and performance tests being performed on 200TB system now
  - Additional disks to be incorporated into cluster
  - Migration of user data from gluster to CephFS
  - Decommissioning of gluster and incorporation of disks into Ceph
- Estimated complete migration by the end of July
- Net result will be a faster and more robust mass storage system, with almost 2x the existing capacity

# Mass Storage Upgrade



- Monitoring and maintenance capabilities of new system already far superior

15

# Today's Review

- Indico page with timetable and slides:
  - https://indico.mit.edu/event/1073/