

Preliminary/Draft

Project Summary

Funding of \$150,000 has been procured by Harold Scheraga for updates and additions to the Matrix cluster over the next four years. It is projected that \$50,000 annually will be spent on this project in three of the four years. The current cluster is based on an older operating system and technology, and will need structural upgrades in the first year support these updates and additions. Once these structural upgrades are in place, it should be easier to add additional computational nodes in future years of this project. Year 1 of this project will address the headnode and other infrastructure, along with some computational node additions. Future years investments will be primarily focused on the addition of computational nodes, but infrastructure updates could occur as deemed necessary. Another possibility would be to spend more in Year 1 on new infrastructure and computational nodes, and less in future years.

Current Matrix cluster summary:

- Headnode – Purchased late 2010
 - Storage
 - OS drives: 2x160GB Raid 1 = 160GB (includes scratch)
 - Data drives: 4x3TB Raid 10 = 6TB (Includes user & application files)
 - Hot spare data drive: 1x3TB
 - OS - Fedora 13 (released May 5, 2010)
 - Basic UPS power protection
 - Currently backed up via Cornell CIT's EZBackup service for file snapshots and disaster recovery
- Computational CPU nodes – 87 of varying age and specifications (3-6 years old)
- Computational GPU nodes – 4 approximately 3 years old – not integrated with existing cluster
- 3 racks, approximately 5 network switches, approximately 15-20 power strips, network cabling

There are a number of reliability and supportability changes that need to be considered for the new Matrix upgrade. Updates needed include:

- Headnode
- Operating System and cluster software
- Computational nodes
- Storage
- Backup method / system
- Rack space
- Networking
- Power

Some guidance on preferences in these areas are needed as outlined below.

Operating System and cluster software

In order to add new computing nodes to Matrix, an up-to-date Operating System will be needed, because building a new kernel in Fedora 13 is not practical.

- The new CPU's supported in the current OS / kernel.
- ChemIT deployment experience & support is now with CentOS, with Werewolf provisioning, and Torque & Maui scheduler software. ChemIT has built 4 clusters using CentOS in the past year.

- Fedora support in Chemistry is repair-only at this point.
- ChemIT recommends using CentOS 6.5 for the operating system.

Headnode

A new headnode will be installed, built up, and tested while leaving the current headnode and cluster in operation until the new cluster is ready to go. Use of newer hardware is advised, to improve performance and reliability of the cluster. The specs of the headnode are dependent on storage decisions (below).

- Once the new cluster is confirmed by researchers as operational, there will be a cut-over date, after which existing nodes will be converted to the new cluster.
- The old headnode will be retained for a period of time (approx.. 30-60 days) as a reference only, to assist with application configuration or debugging if needed. After this time, the old headnode can be converted to an additional compute node.
- The storage drives from the old head node may be retained longer or archived if desired – in this case, substitute drives will need to be purchased.
- ChemIT recommends using 1U single computer headnode of the similar specs and base configuration as the computing nodes if possible, with additional memory and disk space. But the final determination must be informed by storage decisions (below).
 - ChemIT recommends 32 GB RAM (or more)

Storage

- ChemIT strongly recommends separating “system storage” from “user storage” of files / long term storage. This will simplify the headnode design (reducing its cost and risks), improve performance, and provide more storage and backup options.
 - Most computational clusters require users to only store files for current computations on the cluster / headnode, and then move their files to separate long-term storage when complete. This separate storage system can be integrated to make things easy / seamless for users.
- There are several cost-effective options for external storage.
 - Headnode storage –The headnode storage should be of the highest performance (solid-state drive (SSD)), and provide space for
 1. System - operating system, applications, compute node provisioning
 2. Scratch space
 3. Some temporary data
 - User file storage – as this storage is more static, performance is less critical, which helps with cost, and configurations. Some options include
 1. Local NAS (Network Attached Storage device)
 - Can provide large, expandable storage, without impacting the operation or configuration of the headnode. (iSCSI or NFS mount)
 - Performance similar to local or Cornell network storage.
 - Scale & pricing depends on needs – large systems (10-100TB are very expensive, but small systems (1-24TB) are cost effective and expandable as needed.
 - Device and drives will need to be replaced periodically, so should be budgeted accordingly.

2. Cornell central file service (SFS). This is high-availability storage, with local and off-site backup included in the price, as well as a snapshot for versioning if desired (uses 20% of storage).
 - The service is subscription based, and can expand to meet needs.
 - Connection is via GB Ethernet, and performance is excellent.
 - Both NFS and NTFS configurations are available.
 - Backup and disaster recovery is included with this service.
 - CIT maintains all hardware and upgrades.
3. Separate file server box
 - We can build a custom dedicated file server, separate from the head node.
 - Basic functionality is similar to a NAS appliance
 - The NAS may be more cost effective, flexible, and robust.
4. On-board storage (on headnode) [CURRENTLY BEING DONE]
 - Limited physical space for disks
 - Limited file versioning
 - Lowers system performance and reliability, as it runs on the same disk controller as the system is using.
 - Requires additional backup system
5. User's local computers – Many clusters require users to move their data back to their local computers, and provide their own data storage and backup. While cost effective, this can be harder to use, and less reliable for group data. This is not a likely option for Matrix.
 - ChemIT recommends using a NAS or Cornell SFS service. Additional work will be needed to compare price and performance options.

Backup

“Backup” covers several areas. As we learned in the Fall of 2013, on-system copies are best for quick shadow / snapshot versions of files, but not for disaster recovery.

Types of backups

1. File Versioning – a local disk (fast) or external (disk or tape) incremental copy of files
 2. Operating System Backup – back up snapshot of the operating system, applications, and configurations, but not user data. Used to restore system operation to a recent state. Can be done to a second hard drive, or a backup system. May require system downtime.
 3. File backup – an incremental copy of user's current files, providing redundancy and disaster recovery abilities – to be able to restore the current state. It is desired to have a local copy for ease of access, as well as an off-site copy for disaster recovery. The off-site copy should at least be in another building, preferably another site.
 4. Archive – a long term copy of important data which may need to be referenced in the future.
- Backup of user data is most important for the cluster. Backup options for user include:
 1. EZBackup – Cornell automated backup service. (Fee based) [CURRENTLY BEING DONE]
 - Creates daily copies of files (3 versions standard).
 - Stored on campus at CIT, and an off-campus copy at Weill Medical in NY City, providing true disaster recovery capability.
 - Long-term archive snapshots can be created. (Separate fee)
 2. Use of Cornell SFS for storage. (Fee based)

- Includes EZBackup.
- Versioning option on by default (uses 20% of storage).
- 3. External Backup appliance – may be a 2nd NAS box, or a backup appliance (disk or tape based)
 - This solution could be fairly good, but would not provide off-site backups.
 - Hardware would need to be maintained and replaced periodically, so should be budgeted accordingly
- 4. On-system storage only – this is not recommended.

Compute nodes

- Specs for CPU & RAM have been provided
 - 64 GB RAM in each node is recommended
- Disk size & type (Magnetic or SSD)
- Form factor (twin or quad)
- Number of nodes is based on overall configuration and funding.

Networking

Network for this expansion is planned to stay with current 1GB Ethernet. Network needs will include a switch and cabling. The existing cluster hardware will also be re-cabled, simplifying the layout

Power

- UPS – An Uninterruptible power supply is needed on the headnode to keep it running during short power outages, and to allow the system to do a controlled shutdown in case of a longer outage. This helps minimize disk damage and corruption, which frequently occur due to power outages.
- Distribution Strips – power strips are needed, depending on node configuration and power requirements.

Rack

- An additional computer rack will be needed in year 1, and should provide enough space for all 3 phases of the expansion.

Additional information & decisions needed:

- Project roles & responsibilities
- Project Scheduling – includes critical dates, budget availability, work assignments, testing needs
- Application installation and maintenance roles & responsibilities
- Storage management
- Funding / allocations for any repairs or replacements.
- Maintenance scheduling / agreement