**Purpose**

This policy outlines the guidelines for utilizing the Wright GPU computing cluster at Carnegie Mellon University, a resource designed to support computationally intensive tasks within the Statistics & Data Science, Philosophy, and CMIST departments.

**Head Node**

- **Function:** Users submit jobs and SLURM distributes the job to the corresponding node.
- **Access:** Access to the head node is restricted to authorized users.
- **Usage:** The head node should primarily be used for administrative tasks, such as job submission, monitoring, and configuration. Users submit jobs and SLURM distributes the job to the assigned partition.
- **Restriction**: Users are not permitted to run large jobs on the head node. Automatic measures are taken to ensure users do not overuse the head node.

**Compute Nodes**
- **Function:** Nodes dedicated to performing computational tasks utilizing SLURM.
- **Access:** Direct access to the computing nodes is restricted. Users must utilize SLURM to harness computational resources.
- **Usage:** Users submit batch jobs or can open interactive sessions with nodes utilizing SLURM.
- **Available Nodes and Partitions:**
    - Node 1 (n01) - Statistics & Data Science
        - Partition - statds
        - Partition - statds_condo
    - Node 2 (n02)  - Philosophy
        - Partition - phil
        - Partition - phil_condo
    - Node 3 (n03)  - CMIST
        - Partition - cmist
        - Partition - cmist_condo
- **Hardware:**
    - 8 NVIDIA L40 GPUs
    - 48 GB GDDR6 with ECC
    - 64 CPUs with 4 threads (256 available cores)
    - 1TB of RAM
- **Job Submission:** Jobs should be submitted using the SLURM workload manager.
- **Resource Allocation:** SLURM will allocate resources based on job requirements and priorities.
- **Restrictions:** Users cannot ssh into compute nodes directly unless they have actively running jobs. Once their job ends, any active ssh connections and processes spawned off of the session will terminate.

## User Accounts and Access

- **Account Creation:** Users must obtain authorization from their department to create accounts on the cluster. User accounts are then assigned to a SLURM account (statds, phil, cmist) that will give them access to their assigned partition.
  - [Wright Access Request Form](#)
- **Password Management:** Users are responsible for maintaining the security of their passwords.
- **Access Control:** Access to the cluster will be controlled through username and password based authentication.

## Job Submission and Execution

- **General Usage:** Users are only allowed to run jobs on their assigned node which reflects their department. If a user submits a job to an unassigned node, the job will not run.
- **Shared (Condo) Usage:** Users can submit jobs on a non-affiliated node by using a condo partition. Priority is given to users based on department; condo partitions have lower priorities, so jobs submitted to them can be pre-empted. For example, if a user submits a job to a non-affiliated node via the condo partition, and then a user of that node submits a job that needs those resources, the condo job will be preempted.
- **SLURM Batch Jobs:** Users are only permitted to run scripts via batch jobs in SLURM. Job scripts should be well-structured and documented to facilitate reproducibility. Users must specify required resources and the duration of their job. By default, batch jobs have a maximum time limit of 48 hours. Node owners can request for longer MAX job times if required. If a user submits a job longer than the defined MAX job time, their job will not run.
- **SLURM Interactive Sessions:** Users are only permitted to interact directly with compute nodes by using interactive sessions. Users must specify required resources and the length (in time) of their job. Interactive sessions are limited to 2 hours.
- **Job Monitoring:** Users can monitor the status of their jobs using SLURM commands and "[util.stat.cmu.edu/slurmstats](http://util.stat.cmu.edu/slurmstats)".
- **Job Submit Plugin:** The submission of jobs can be programmatically controlled and monitored via the job submit plugin. The plugin is controlled with a Lua script. The script allows for custom logs, actions, and alerting. An example of its usage is disallowing interactive jobs that request time over a predetermined limit. Users are alerted of the error and can modify their jobs accordingly.
- **GPU Requirement:** Users must request GPUs by explicitly requesting the number they need. If a user requests no GPUs, they aren't given access to any.

**Example Condo Model Usage:** A statds user submits a job to the phil_condo partition and uses all remaining available resources on that node. A philosophy user submits a job to the phil partition with no resources available. Since the phil user is using the phil partition, which takes priority over the phil_condo partition, the statds user's job will be preempted (the consumed computing/GPU resources are cleared) and the philosophy user's job will run. The statds user's job will run again once the other job completes.

**Cluster Usage Policies**

- **Fair Usage:** User jobs can also be prioritized based on different job factors. These are defined by the multifactor plugin, which is defined in the SLURM configuration file. The time a job has been waiting in the queue, the user's history/frequency of job submission, and the size of a job can all determine if a job can be prioritized over others.

*Current Fair ShareMultifactor Job Priority Plugin Configuration*

PriorityWeightAge=1000
PriorityWeightFairshare=10000
PriorityWeightJobSize=1000
PriorityWeightPartition=1000
PriorityWeightQOS=1500
PriorityUsageResetPeriod=MONTHLY

**Example of Fair Share:** Node 1 has 8 GPUs total, but only 4 available. User A submits a job to Node 1 requesting 8 GPUs and a runtime of 1 hour. Since there aren't enough GPUs available, the job will be placed into a pending state until 8 GPUs become available. User B submits multiple jobs that use 1 GPU and a runtime of 48 hours each. Since the GPU resources are available, User B's job(s) will take priority over User A's. However, over time, user B's job(s) will become less prioritized over User A's single job. This is because as User A's job sits and waits, the prioritization weight will increase due to its age and requested runtime only being 1 hour. User B's prioritization will decrease as the number of job submissions grows. This will eventually allow User A's job to run. User usage statistics are reset weekly (can be configured to monthly).

- **Resource Limits:** Users can request all available resources the node offers but a job will only run if those resources are available. This includes both regular and condo partitions.
- **Prohibited Activities:** Running jobs directly on a compute node, interfering with actively running jobs, and manipulating configuration settings are prohibited.

**Maintenance and Support**

- **Scheduled Downtime:** The cluster may experience scheduled downtime for maintenance or upgrades. All maintenance periods will be announced to users via "wright-cluster-info@andrew.cmu.edu".
- **Support Services:** Support services will be available to assist users with cluster usage and troubleshooting by contacting "wright-cluster-info@andrew.cmu.edu".

**Policy Enforcement**

- **Monitoring:** Cluster usage will be monitored to ensure compliance with this policy.
- **Consequences:** Violations of the policy may result in restrictions or loss of access to the cluster.

By following this policy, users can leverage the Wright Cluster to efficiently execute their research and computational tasks.