# HPC3 for HPC Users

High-Level view of dates, plans, and how to get started with HPC3

# Dates to know

- July 2020 – HPC3 entered "phased production"
- Oct 2020 HPC3 entered "production"
- Nov 30, 2020 – CentOS 6 End-of-Life (EOL). Full updates stopped in 2017.
    - This is the base OS for HPC.
- Dec 31, 2020 – Projected End date for HPC
- Dec 31, 2020 – Full updates stop for CentOS7. EOL in 2024.
    - This is the OS that HPC3 runs
    - Already looking at transition to CentOS 8 in 2021/2022.
- Jan 2021 – HPC2 (HPC nodes incompatible with HPC3) Enters Service
- Dec 2024 – HPC2 Shutdown

# Accessing HPC3

ssh  hpc3.rcic.uci.edu

If you have an account on HPC, your account is *already available on* HPC3
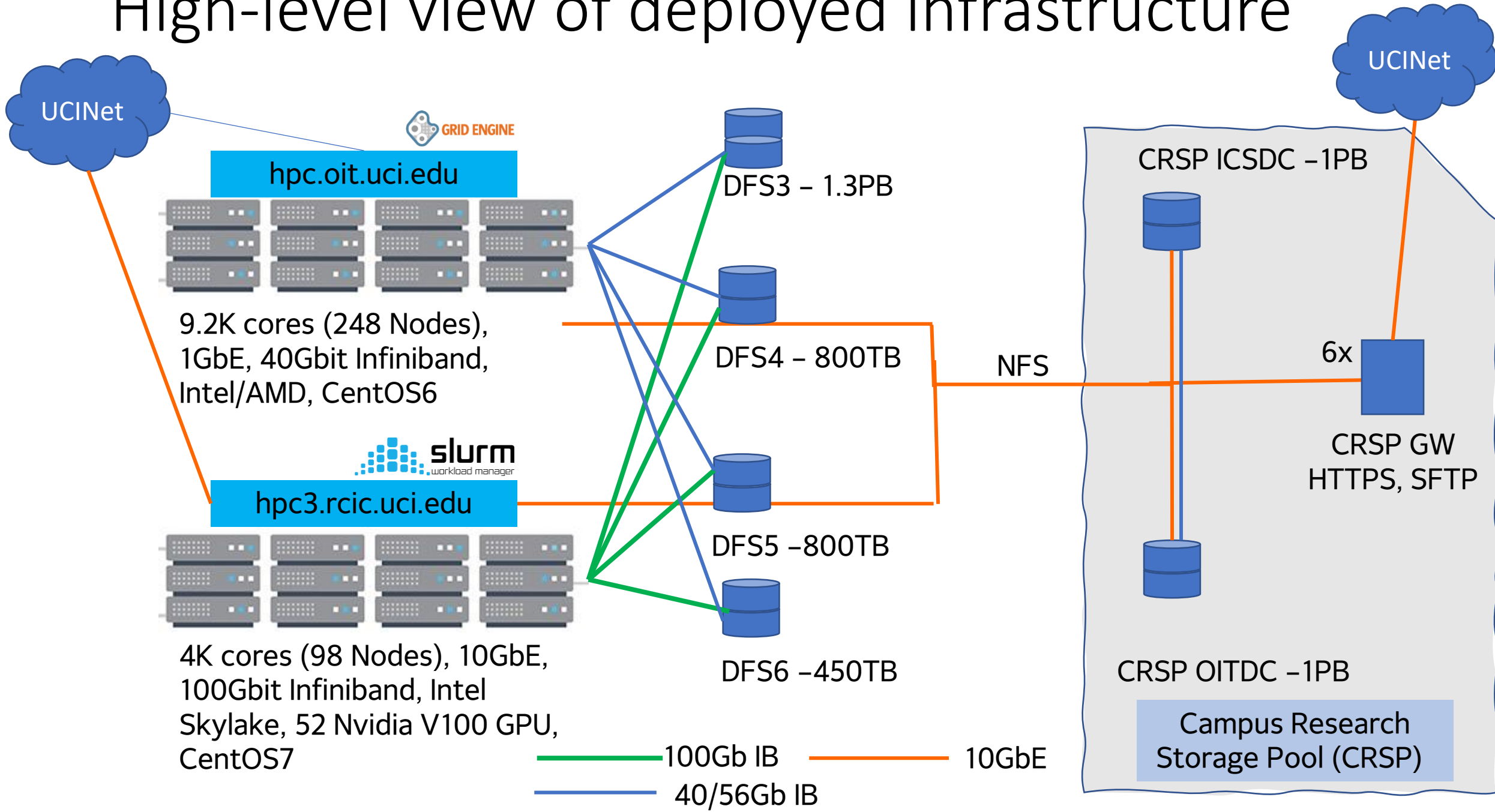
# HPC3 – Oct 2020

98 nodes
3920 cores
52 Nvidia V100 GPUs

- HPE Apollo 2600 CPU Nodes  (4 nodes/2U)
    - 2 x 20 Core Intel Skylake
    - 192GB RAM (384GB + 768GB variants)
    - 10GbE
    - 100Gb EDR Infiniband
    - 2 X SSDs/node
- HPE DL380 GPU Nodes  (1 node/2U)
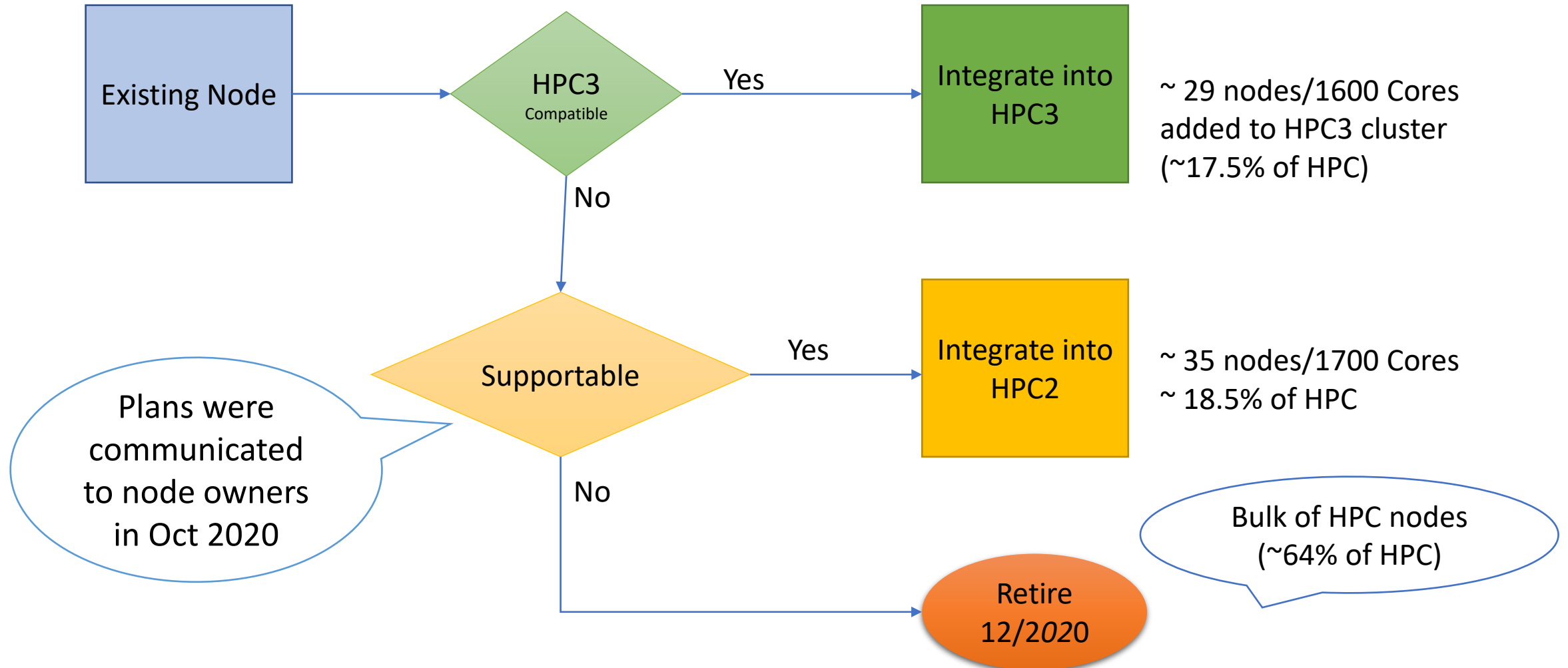    - As above +
    - 4 Nvidia V100 GPUs

# High-level view of deployed infrastructure



UCINet

UCINet

GRID ENGINE

hpc.oit.uci.edu

9.2K cores (248 Nodes),
1GbE, 40Gbit Infiniband,
Intel/AMD, CentOS6

slurm
workload manager

hpc3.rcic.uci.edu

4K cores (98 Nodes), 10GbE,
100Gbit Infiniband, Intel
Skylake, 52 Nvidia V100 GPU,
CentOS7

DFS3 – 1.3PB

DFS4 – 800TB

DFS5 –800TB

DFS6 –450TB

NFS

CRSP ICSDC –1PB

6x

CRSP GW
HTTPS, SFTP

CRSP OITDC –1PB

Campus Research
Storage Pool (CRSP)

100Gb IB          10GbE

40/56Gb IB

# What Happens to Existing HPC Hardware? Maintain Two Clusters: HPC2 and HPC3

Existing Node → HPC3 Compatible

**Yes** → Integrate into HPC3 — ~ 29 nodes/1600 Cores added to HPC3 cluster (~17.5% of HPC)

**No** → Supportable

**Yes** → Integrate into HPC2 — ~ 35 nodes/1700 Cores ~ 18.5% of HPC

**No** → Retire 12/2020

Plans were communicated to node owners in Oct 2020

Bulk of HPC nodes (~64% of HPC)

# Differences Between HPC2 and HPC3

- Generally run the same OS and application stack
- HPC2 is pure condo (owner queues only. No free/pub access)
  - Hardware that is not compatible with HPC3
  - Will run to ~ end of 2024. (earlier if hardware won't update to CentOS8)
  - CANNOT be expanded.
- HPC3
  - No "owner" queues
  - Instead, Free queues and "Accounted" (charged to an account)
  - Expandable
  - Granted Cycles, Purchased Hardware, Purchased Cycles

# Expected Hardware State of HPC3 in 2020

- By Jan 31,2021
  - 5900 Cores, 135 nodes  (includes purchases and compatible HPC nodes)
  - No core more than 2 years old
  - FDR/EDR Infiniband
- By Jun 30, 2021
  - Expect to purchase another O(2000) cores from UCI funds
  - Unknown number of Faculty purchases.

- No firm GPU expansion plans
  - We are looking at vendors who might supply less expensive, "gamer" GPUS in servers.

# Notable Changes from HPC to HPC3

- CentOS6
- SGE Scheduler
- Checkpointing
- Owner Queues, Free, Pub
- No job hour accounting
- Free queues migrated/checkpointed jobs

HPC

→

- CentOS7
- Slurm Scheduler
- ~~Checkpointing~~
- Standard, High-Memory, GPU, and Free Queues (partitions)
- **Accounted** jobs charge a "bank of hours". *Cannot* be killed
- **Free** jobs cost nothing, but *can be* killed by accounted jobs. No restart.

HPC3

# Transitions can be "difficult"

- We'll do our level best to make it smooth, but there is a number of things that will cause some bumps.
  - **Home areas are different**. DFS3/4/5/6 systems are available on both HPC and HPC3
  - The HPC Software Stack was 9+ years in the making.
    - We have already transitioned (updated, rebuilt, harmonized) a good fraction of it. See the *Software Environment* user guide https://rcic.uci.edu/hpc3/software-tutorial.html
  - ALL User-compiled software needs to be rebuilt. The OS changed.
  - New Slurm scheduler.
    - Users must transition their batch scripts from SGE to SLURM. This is straight-forward in most cases. See the *Slurm Batch Jobs* user guide https://rcic.uci.edu/hpc3/slurm.html
  - New sharing policy
    - Accounted vs. Free Jobs. See the *Reference* user guide https://rcic.uci.edu/hpc3/hpc3-reference.html
    - No checkpoint/restart (technology on HPC is dead, BLCR was last developed in 2013)
  - New GPUs are available
    - Still a limited resource, but there is more availability (by request)
- We, RCIC, are going to have very full plates as we move through this, we ask for patience.

# Updated Web Site

While documentation is the last thing to be built, we're made some really good progress.

- [https://rcic.uci.edu/hpc3/hpc3-reference.html](https://rcic.uci.edu/hpc3/hpc3-reference.html)

- Reference to many of the details of HPC3
- Slurm guide, definition of basic queues
- The software environment

# How are the queues changing?

- Instead of "owner" queues there are **Slurm Accounts**
- The CPU hours allocated to a job are charged to an account
  - 1 core/hr == 1 unit
- GPU hours need a special GPU account
- Are there "no cost" free queues available (like free, free64 on HPC?)
  YES:
    - *free partition*, cost 1 core/hr = 0 units
    - *free-gpu partition, cost* 1 GPU/hr = 0 units
- I like free, what's the catch?
  - **A free job can be killed at any time** to make room for a non-free (allocated) job.
  - These are neither checkpointed nor automatically requeued, it's a user responsibility to resubmit a job.

# Two Job Classes on HPC3

- Accounted
  - A bank of core hours is debited after a job is completed
  - Job CANNOT be pre-empted by another job
- Free
  - No core hours are debited (you need an account with at least 1 hour of credit to run, this is a SLURM requirement)
  - Job CAN be pre-empted (killed) by an accounted job. But not by another free job.

**HPC3 Queues**

| Table 1. HPC3 Available Queues | | | | | |
|---|---|---|---|---|---|
| **Partition** | **Default memory/core** | **Max memory/core** | **Default / Max runtime** | **Cost** | **Jobs preemption** |
| **CPU Partitions** | | | | | |
| **standard** | 4.5GB | 4.5GB | 2day / 14day | 1 / core-hr | No |
| **free** | 4.5GB | 18.0GB | 1day / 3day | 0 | Yes |
| **debug** | 4.5GB | 18.0GB | 15min / 30min | 1/core-hr | No |
| **highmem** | 9.0GB | 9.0GB | 2day / 14day | 1/core-hr | No |
| **hugemem** | 18.0GB | 18.0GB | 2day / 14day | 1/core-hr | No |
| **GPU Partitions** | | | | | |
| **gpu** | 4.5GB | 4.5GB | 2day / 14day | 1/core-hr, 32/GPU-hr | No |
| **free-gpu** | 4.5GB | 9.0GB | 1day / 3day | 0 | Yes |
| **gpu-debug** | 4.5GB | 9.0GB | 15min / 30min | 1/core-hr, 32/GPU-hr | No |

# Let's talk accounted jobs

- Every User on HPC3 is given a one-time 1000 CPU hour account
  - If you do not specify an account for your batch job, **THIS is the one that gets charged.**
  - **If you do not pay attention, you will run out of hours quickly.**

- Most users on HPC3 are affiliated with a research group/class.  Research groups have separate accounts.  **You need to explicitly charge a specific account**

- **Free jobs require a Slurm account with positive balance.**  That account is not charged, but a job will not start if the balance is negative.
  - As long as your personal Slurm account has positive balance, your free jobs will start.

# Some tips on effectively utilizing your allocation

- We HIGHLY recommend:
  - Use free for exploratory work
  - Use accounted for Production work
- I'm an "owner", what's in it for me?
  - We compute what your node could *theoretically* deliver if used exclusively by you and then credit 95% (annually) of that to your lab Slurm account.
    - Example: 40-core system is credited with 332800 core hours
- I bought just one node, can I use more than 1 node?
  - Most definitely!
  - If the job is accounted, it won't be killed.
  - Allows you to expand to something larger for shorter periods of time.

# I need more memory! (?)

- Access to highmem and hugemem queues is limited to
  - Groups who purchased these nodes
  - Others who demonstrated that the standard queue doesn't meet their needs
- There are fewer highmem/hugemem nodes, and we try to keep them available.
- Free queue jobs can request up to 18GB/core.
- `seff jobID`
  - After a job has finished, you can see what it utilized

- Sample Job has good CPU utilization! (17.2 core-days in 10.5 hours, ~64%)
- Could have easily fit on a standard memory node

```
$ seff 1356600
Job ID: 1356600
Cluster: hpc3
User/Group: user1/user1
State: COMPLETED (exit code 0)
Nodes: 1
Cores per node: 40
CPU Utilized: 11-06:08:24
CPU Efficiency: 63.84% of 17-15:10:40 core-walltime
Job Wall-clock time: 10:34:46
Memory Utilized: 13.99 GB
Memory Efficiency: 3.98% of 351.56 GB
```

# Sharing of Slurm Accounts?

- Every user is given a personal Slurm account, these cannot be shared with others
- "Lab"/"Group" SLURM bank accounts are designed TO BE SHARED among users within the same lab/group

  - NEVER share your Unix login account/password!.

- <u>Account coordinators</u> – people with special privileges on a Slurm Account.
  - Can add/remove/limit other users
  - RCIC must add the Unix login (UCNetID) and create that user's personal Slurm account
  - https://rcic.uci.edu/hpc3/slurm.html#_account_access_control

# Use **sbank** to find out which accounts you can access

```
login-x$ sbank balance statement
User          Usage |        Account        Usage | Account Limit Available (CPU
hrs)

---------- --------- + --------------- --------- + ------------- ---------
panteater *        0 |        PANTEATER        0 |         1,000     1,000
panteater *        0 |        UCI-PI_LAB   24,990 |       100,000    75,010
user17             0 |        UCI-PI_LAB   24,990 |       100,000    75,010
uci-pi             0 |        UCI-PI_LAB   24,990 |       100,000    75,010
user760            0 |        UCI-PI_LAB   24,990 |       100,000    75,010
user6004           0 |        UCI-PI_LAB   24,990 |       100,000    75,010
user349       21,416 |        UCI-PI_LAB   24,990 |       100,000    75,010
user5          3,574 |        UCI-PI_LAB   24,990 |       100,000    75,010
```

user `panteater` can charge to two accounts: 1) the personal account and 2) the `UCI-PI_LAB` account.

# What is allocated to users

- Every user: 1000 core hours. This is a **one-time** allocation.

- Faculty and researchers who function as PIs: ~ 200K core hours/year.
  - Must specifically request this allocation
  - Fulfilled from UCI-purchased nodes.
  - First allocation is for 6 months.
  - Will adjust future allocations depending on requests, available hardware, utilization

- 50 GB of home area space

- 1 TB of temporary space in /pub/<ucnetid>

# Evolving policy

- Goal
  - If a user wants to run a single core, allocated job in the standard queue, there is a very high probability that a core is available.
- How?
  - Cluster is not over-allocated
  - Average is 80% should run allocated jobs, 20% free.
    - There will be times when we cannot achieve the above goal (congestion/abuse)
  - Overzealous users can easily run their accounts dry.
  - We can put per-user "clamps" on utilization for problematic users
- General policies help
  - No more than 400 active jobs/user
  - Fair-share.  When lots of queued jobs across the cluster user who is utilizing more cores has their queue priority lowered.

# Various Limits

- Limits on usage are required to *balance* <u>fair access</u> and <u>large resource usage</u>, for example
  - # active cores, # queued jobs, # active jobs, #consumables (e.g. GPUs)
  - <u>Most</u> users will likely never run into limits
- Principles
  - No single group should have active more than about ~25% of any specific resource
    - largest # of active cores should be O(1000)
  - Number of Queued jobs should be reasonable
    - about 3000 queued jobs/user
  - Exceptions need to be made for
    - Specific requests to run larger jobs for relatively short periods
    - More severe limits on jobs that adversely affect the system (usually I/O related)
    - Small number of groups that have purchased MORE than 25% of a specific resource (e.g. huge memory nodes)

# Talking to RCIC and to Each Other

- **How do I ask for help/talk to RCIC?**
  - Same as HPC: send email to hpc-support@uci.edu
    This automatically creates a help ticket
  - Read that fine website: https://rcic.uci.edu


- **What about talking to RCIC and other Users@UCI?**
  - Join  the new! Google group
    https://groups.google.com/a/uci.edu/g/rcic-users
  - Chat with us on Slack: https://rcicos.slack.com/