

# Upcoming Biowulf Seminars

Bldg 50, Rm 1227

- November 30, 1 - 3 pm

## **Python in HPC**

Overview of python tools used in high performance computing, and how to improve the performance of your python jobs on Biowulf

- Jan 16, 1 - 3 pm

## **Relion tips and tricks, and Parallel jobs and benchmarking**

Mechanics and best practices for submitting RELION jobs to the batch system from both the command line and via the RELION GUI, as well as methods for monitoring and evaluating the results. Scaling of parallel jobs, how to benchmark to make effective use of your allocated resources

# Making Effective Use of the Biowulf Batch System

Steven Fellini

[staff@hpc.nih.gov](mailto:staff@hpc.nih.gov)

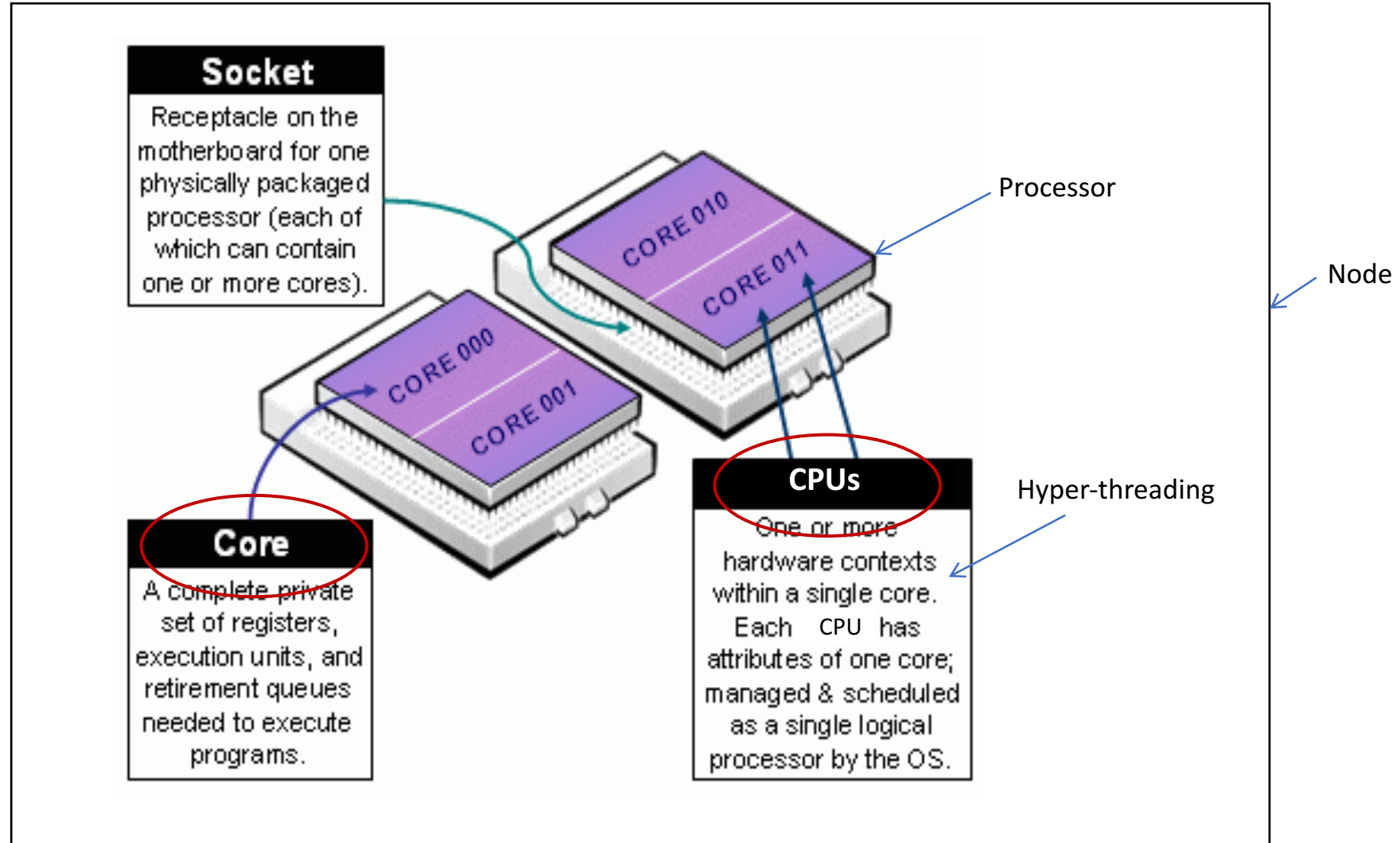
NIH HPC Staff, CIT

Oct 30, 2017

# Effective Use == Effective Resource Allocation

- Specifying resources
  - Estimating required resources
  - Allocating resources with *sbatch* and *swarm*
- Monitoring resource allocation
- Scheduling and resource allocation
- Post-mortem analysis

# Hardware Terminology Review



# Estimating Resources

- CPU
  - Check documentation (<https://hpc.nih.gov/apps/>)
  - Objective -- match CPU:Threads 1:1  
(there are exceptions, e.g., MD jobs)
- Memory
  - Run a job or swarm with a large memory allocation
  - Check actual memory usage
  - Add 10% to actual memory usage
- Time
  - Run a job or swarm with a large time allocation
  - Check actual wall time
  - Add 10% to actual wall time

# Allocating Resources with *sbatch* and *swarm*

- All jobs
  - --mem (sbatch) or -g (swarm)
  - --time (sbatch and swarm)
  - -b to bundle command lines (swarm)
- Single-threaded jobs
  - “-p 2” to load cores with 2 threads (swarm)
- Multi-threaded jobs
  - --cpus-per-task (sbatch) or -t (swarm)
  - Use \$SLURM\_CPUS\_PER\_TASK in batch script
  - OMP\_NUM\_THREADS
- Multi-node jobs
  - “Parallel Jobs and Benchmarking” Jan 16

# Monitoring Resource Allocation

- CPU
  - jobload while the job is running
  - Dashboard during or after the job
  - (No easy way to monitor GPU utilization at the moment)
- Memory
  - jobload while the job is running
  - jobhist, Dashboard or sacct during or after the job has completed
- Walltime
  - jobhist, Dashboard or sacct during or after the job has completed



# jobload

```
% jobload -u someuser
```

JOBID	TIME Elapsed / Wall	NODES	CPUS Alloc	THREADS Active	LOAD	MEMORY Used / Alloc
51863534	6-22:08:01 / 10-00:00:00	cn3095	4	4	100%	1.0 / 8.0 GB
51863535	6-22:08:01 / 10-00:00:00	cn3256	4	5	125%	0.9 / 8.0 GB
51863536	6-22:08:01 / 10-00:00:00	cn3348	4	1	25%	1.0 / 8.0 GB
51863537	6-22:08:01 / 10-00:00:00	cn3401	4	3	75%	0.9 / 8.0 GB
51881591	6-19:42:16 / 10-00:00:00	cn3097	4	1	25%	1.0 / 8.0 GB

```
% jobload -j 51874438_233
```

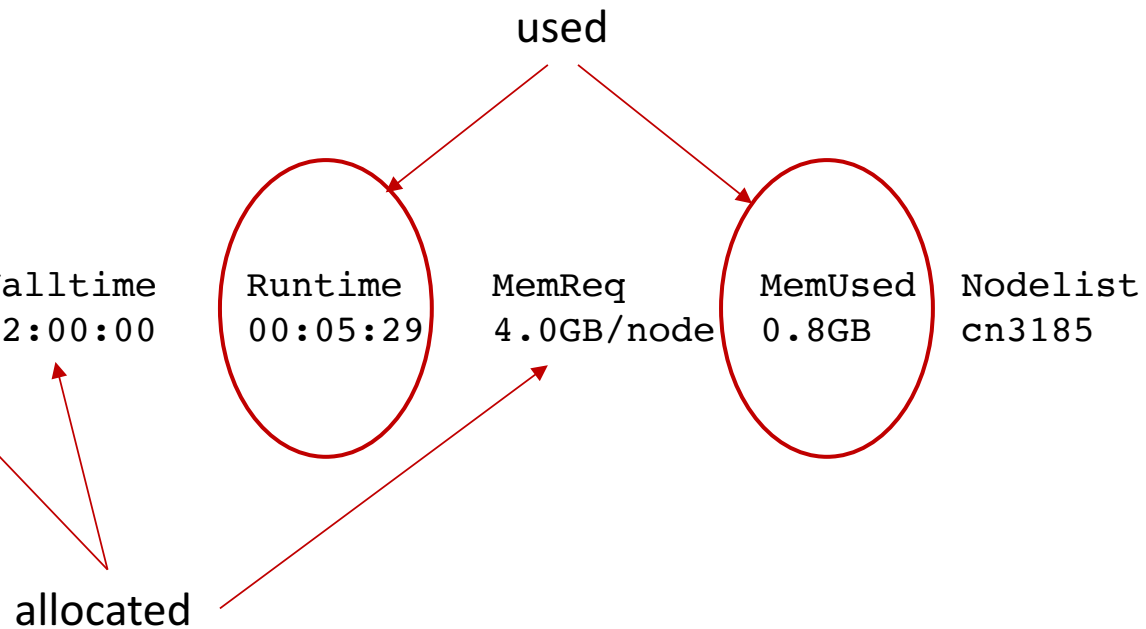
JOBID	TIME Elapsed/Wall	NODES	CPUS Alloc	THREADS Active	LOAD	MEMORY Used/Alloc
51874438_233	6-20:10:13/10-00:00:00	cn3105	2	1	50%	0.5/1.5 GB



# jobhist

```
# jobhist 52102264_67
Jobid      Partition
52102264_67 norm
```

State	Nodes	CPUs	Walltime	Runtime	MemReq	MemUsed	Nodelist
COMPLETED	1	2	02:00:00	00:05:29	4.0GB/node	0.8GB	cn3185

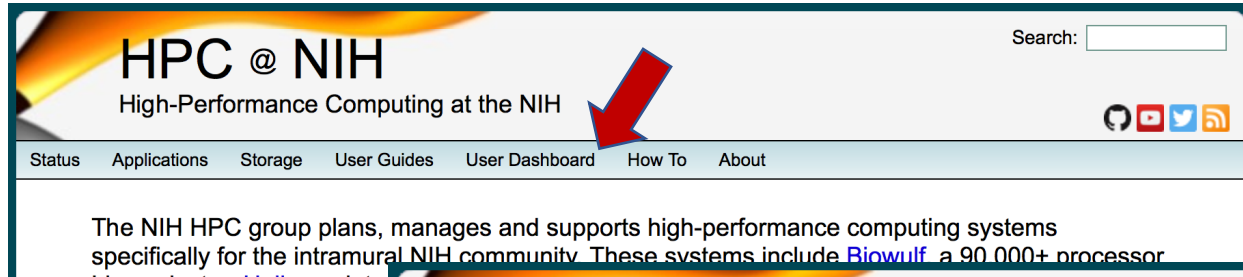


# sacct

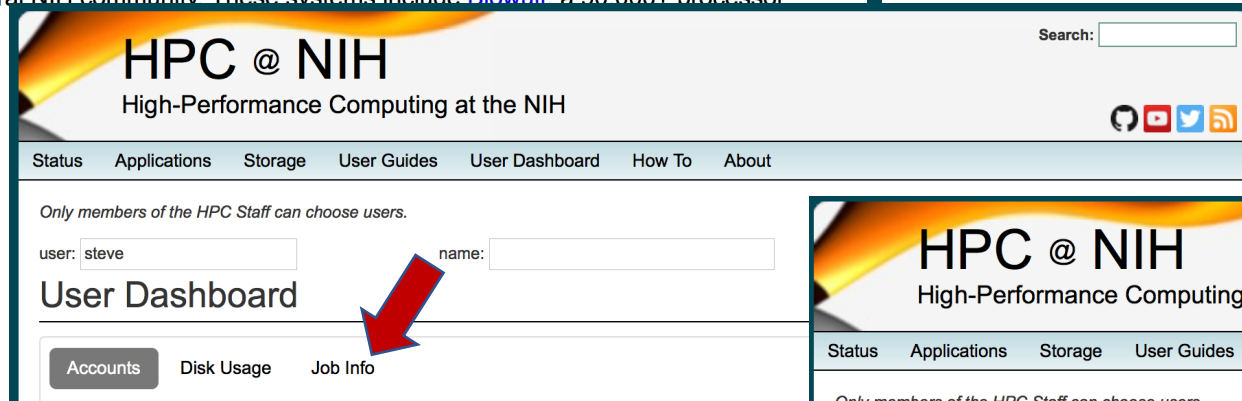
```
% sacct --format=Jobname,AllocCPUS,AllocNodes,ReqMem,MaxRSS,Elapsed -j 52102332
```

JobName	AllocCPUS	AllocNodes	ReqMem	MaxRSS	Elapsed
tbss_2_reg	2	1	4Gn	815152K	00:05:29
batch	2	1	4Gn	815152K	00:05:29

# Using Your Dashboard to Monitor Jobs



<https://hpc.nih.gov>



The screenshot shows the HPC @ NIH User Dashboard with the 'Job Info' tab selected. It displays a table of job information with columns for jobid, jobname, submit, start, end, and state. The table is filtered to show 10 entries. A search bar is also present.

jobid	jobname	submit	start	end	state
51819363	sinteractive	Oct 19, 8:35:19	Oct 19, 8:35:26	Oct 19, 16:35:46	TIMEOUT
51819128	sinteractive	Oct 19, 8:22:15	Oct 19, 8:22:19	Oct 19, 8:31:57	COMPLE...
51762779	bash	Oct 18, 12:38:08	Oct 18, 12:38:55	Oct 18, 12:39:37	CANCEL...
52557544	sbatch	Oct 26, 8:29:11	Oct 26, 8:29:15	n/a	RUNNING
52557242	sbatch	Oct 26, 8:27:02	Oct 26, 8:27:02	Oct 26, 8:27:04	FAILED

<https://hpc.nih.gov/dashboard/>

# Scheduling and Resource Allocation

- Scheduling is determined by job priority
- Priority is determined by Fairshare value of user
- Fairshare is determined by recent cpu and memory allocations of running jobs
- Unnecessarily long time allocation will prevent jobs from being backfilled

# Why are my jobs pending?

- 'freen' shows free CPUs but not free memory or disk
- Other jobs have higher priority (sprio)
- Nodes are reserved for higher-priority jobs

# Consequences of...

	Specifying more resources than needed	Specifying fewer resources than needed
CPU	Wasted CPU resources, possibly unnecessary scheduling delays	Job runs a little/a lot slower
Memory	Wasted memory resources, possibly unnecessary scheduling delays	Job is “Killed” by the kernel
Time	Possibly unnecessary scheduling delays	Job is killed by the batch system

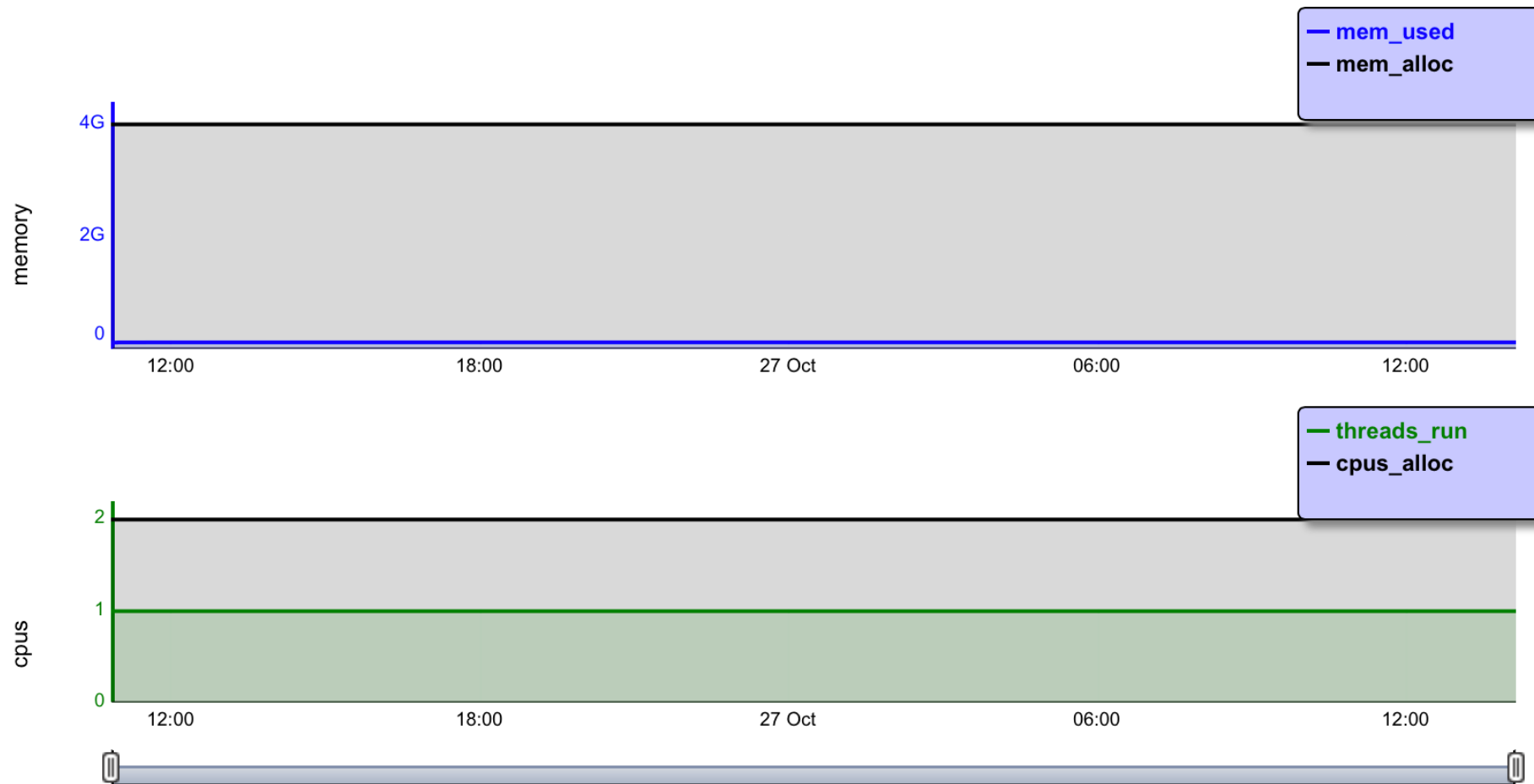
# Post-mortem of jobs using user Dashboard

Or

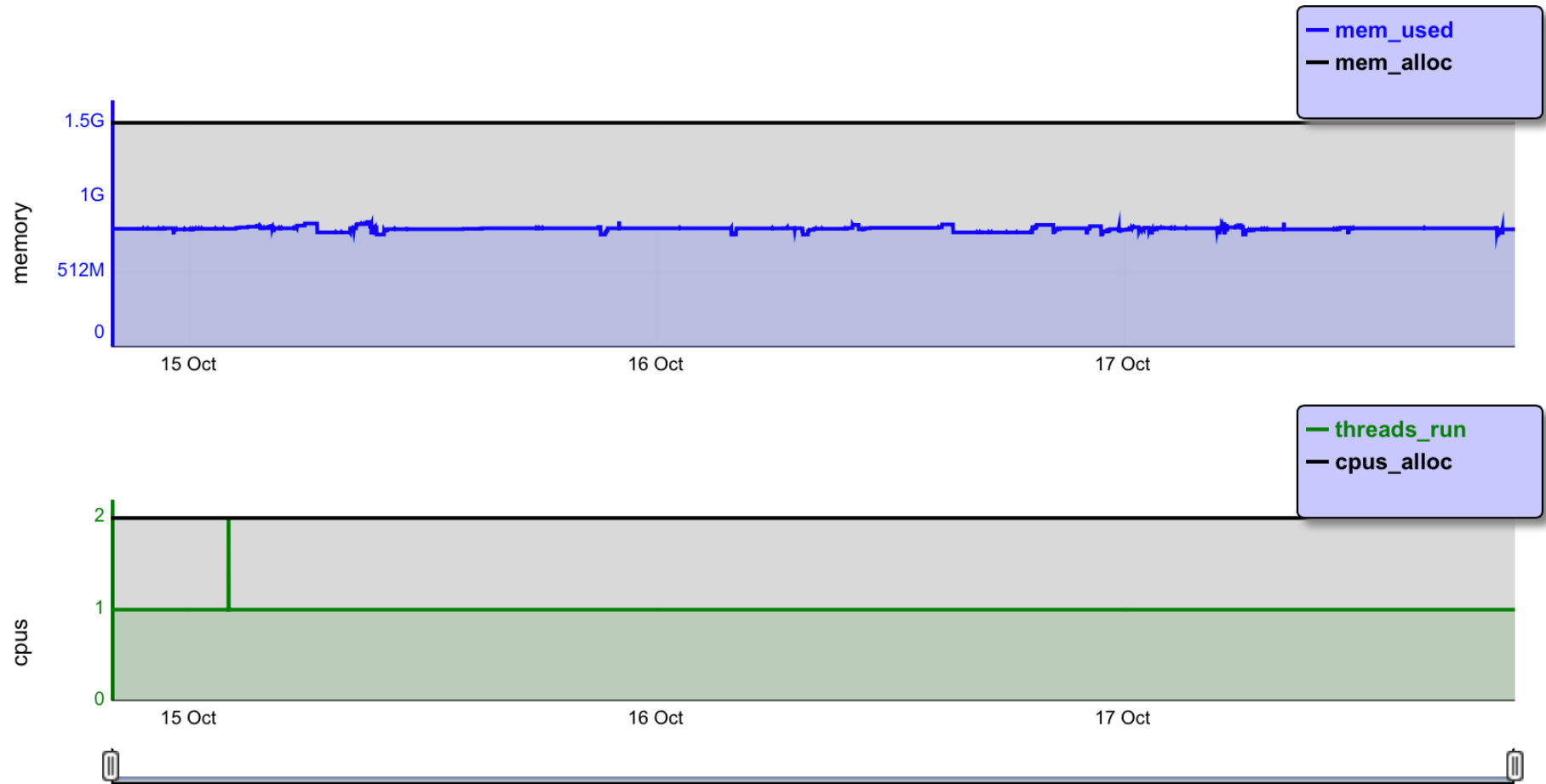
The Good,  
the Bad,  
and the Ugly...



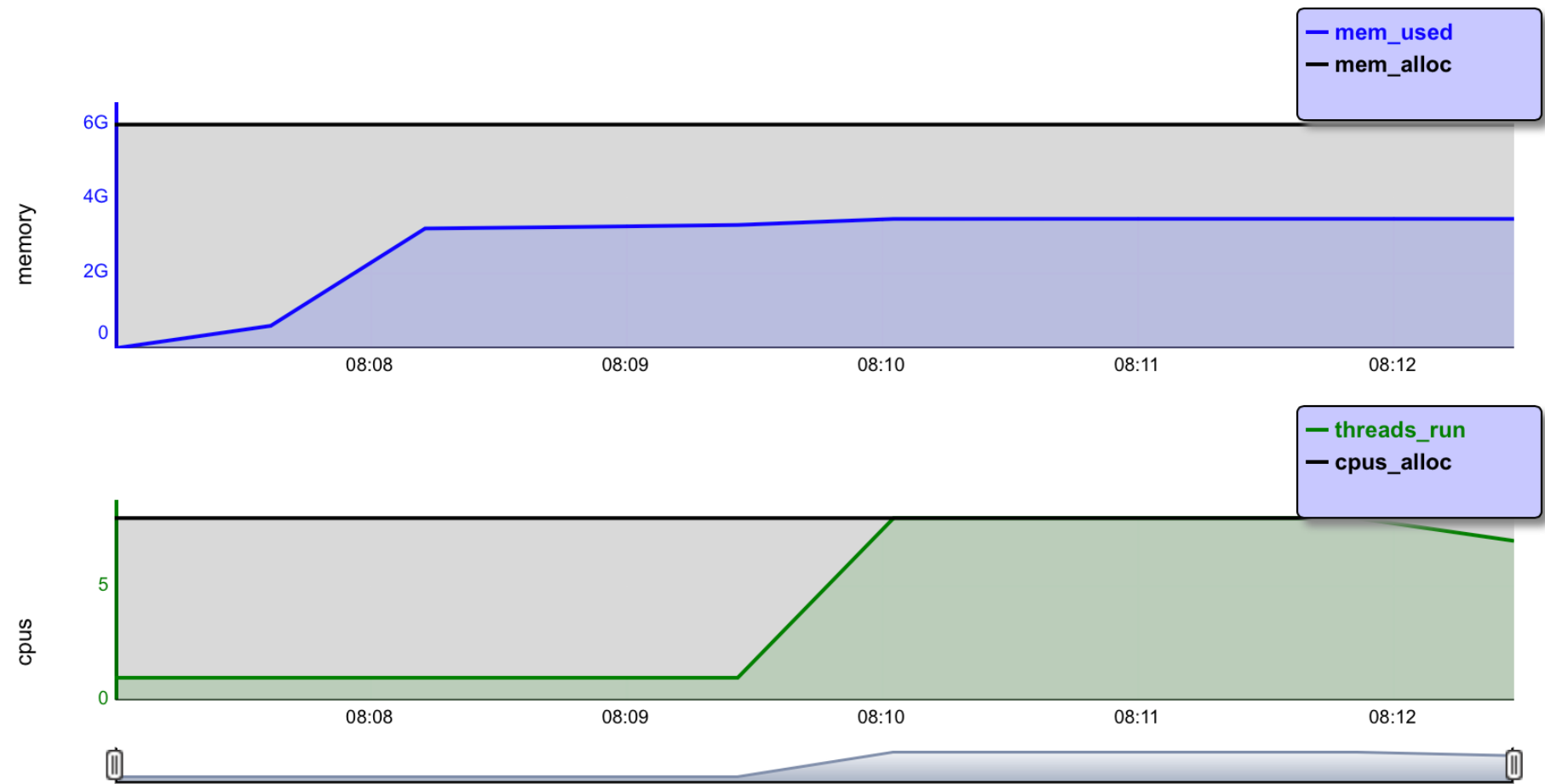




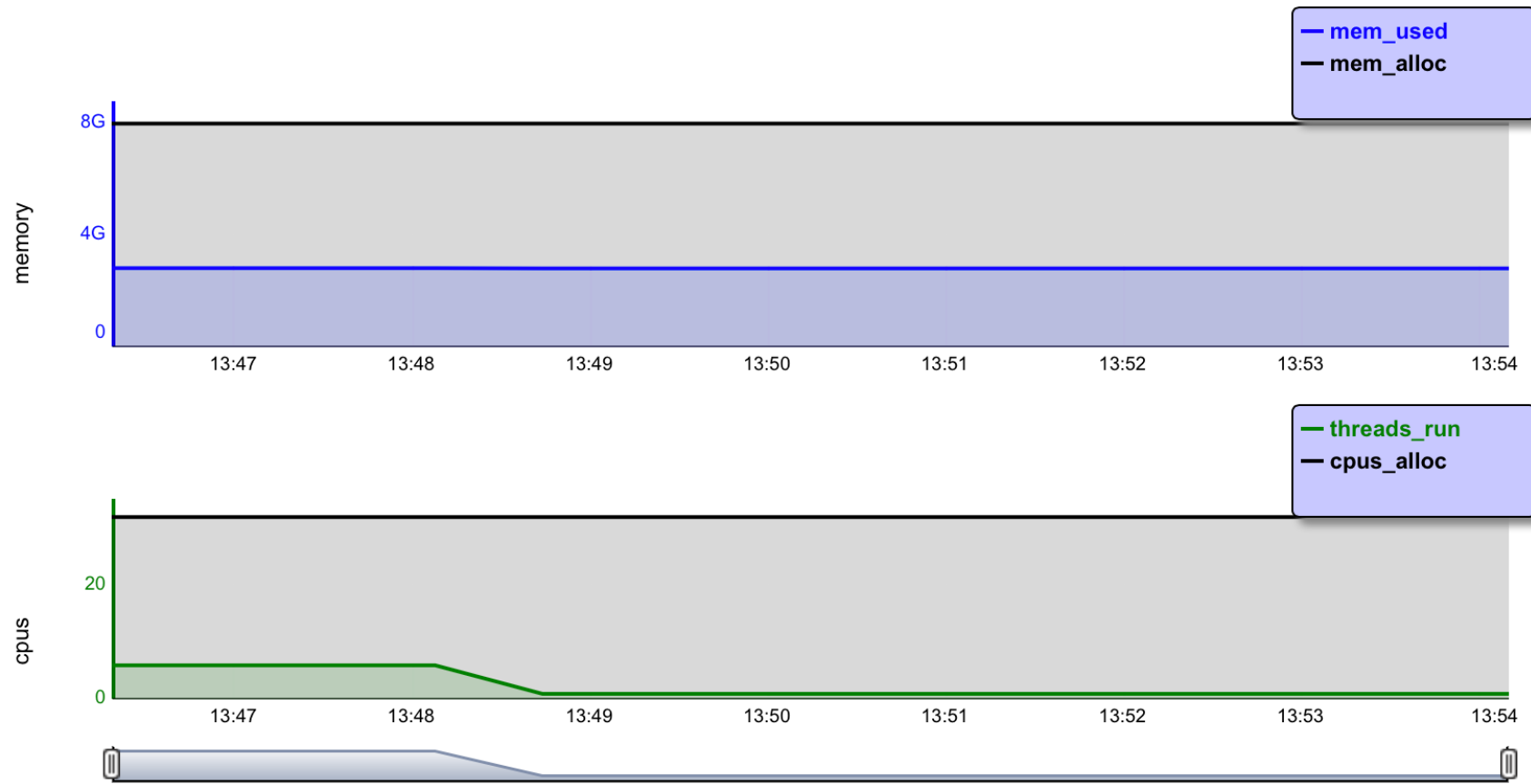
Comment: job is running with default allocations for CPU and memory  
Recommendation: if a subjob of a large swarm, try “-p 2”



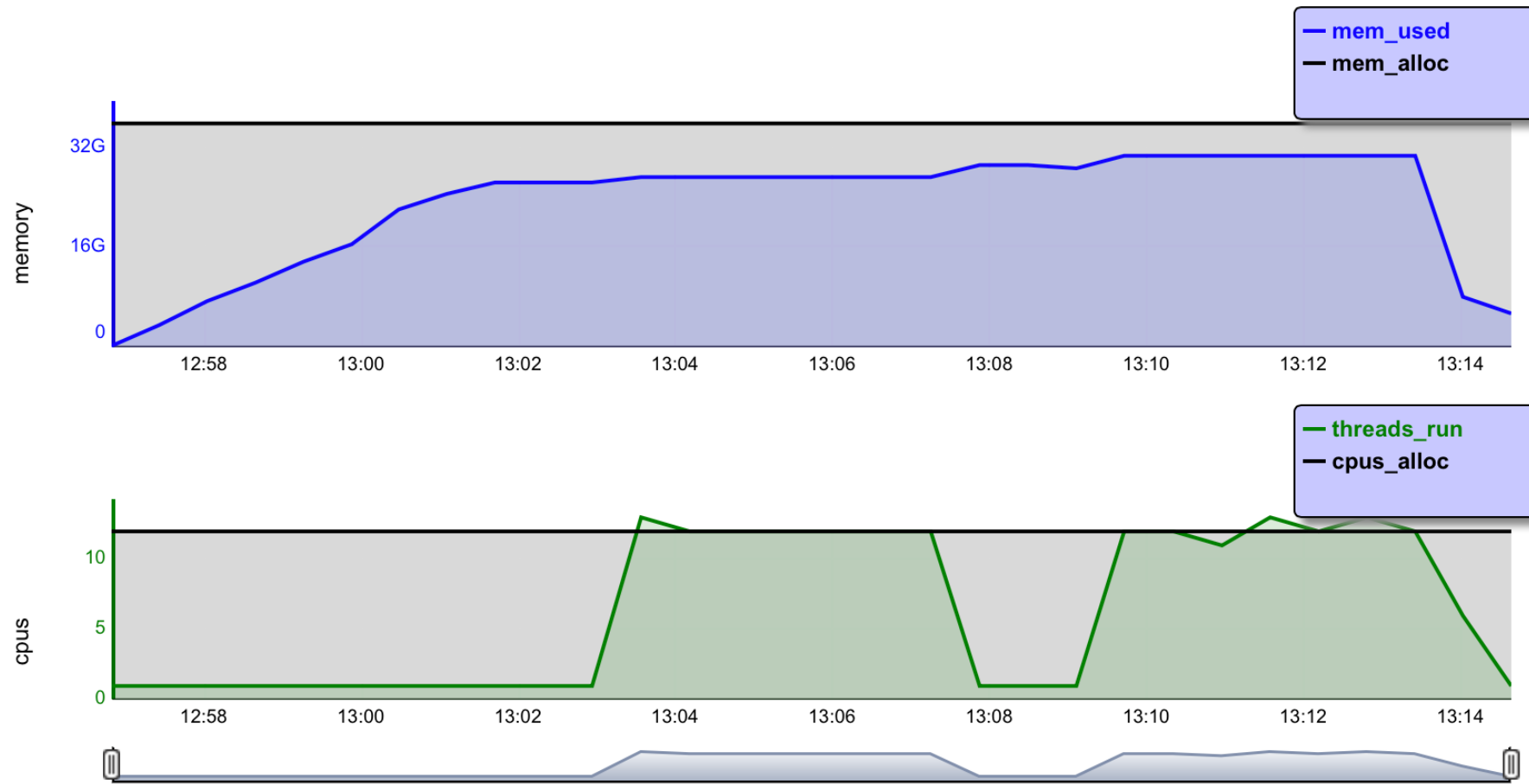
A+



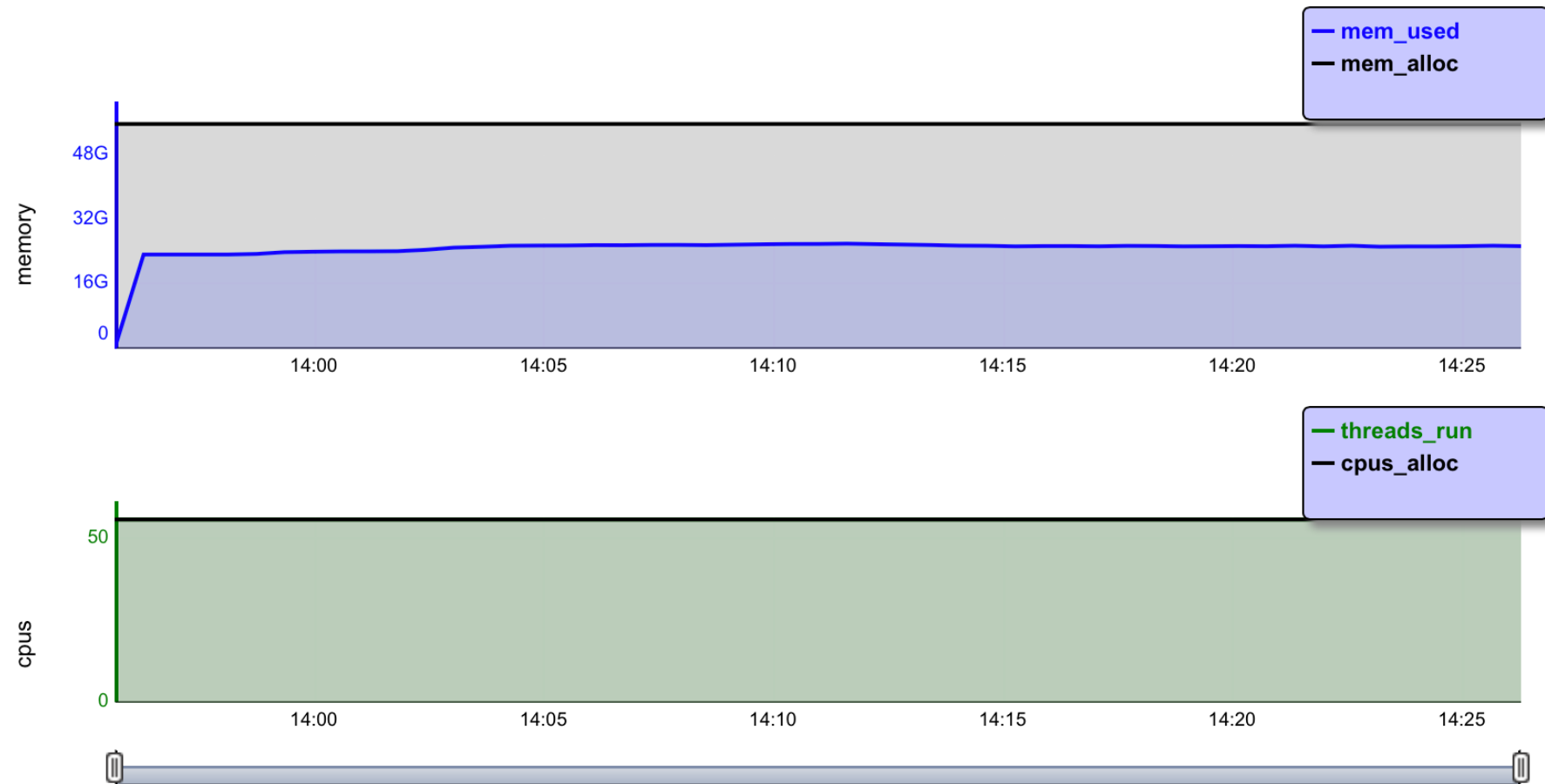
Another A+



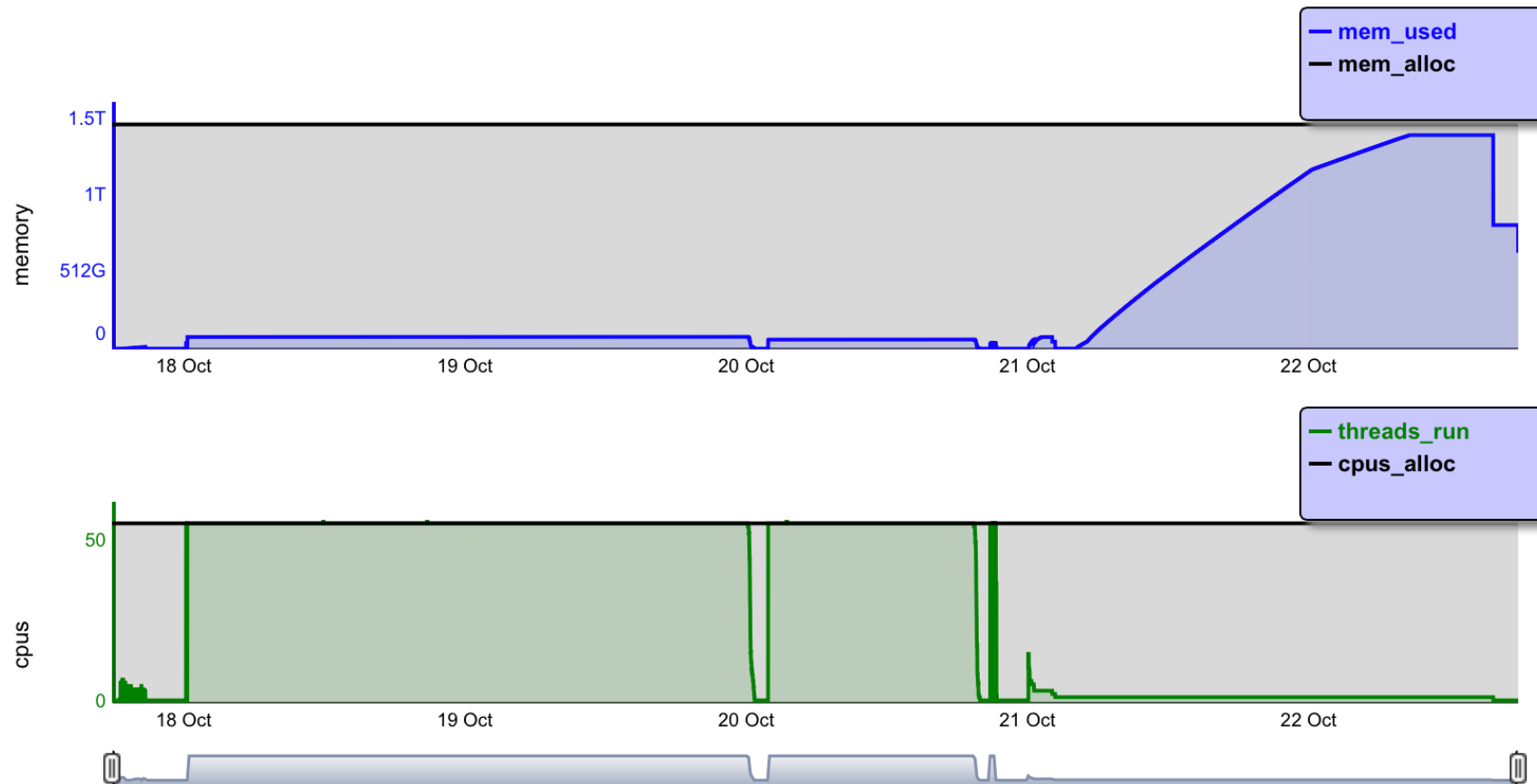
Recommendation: reduce CPU allocation



A+

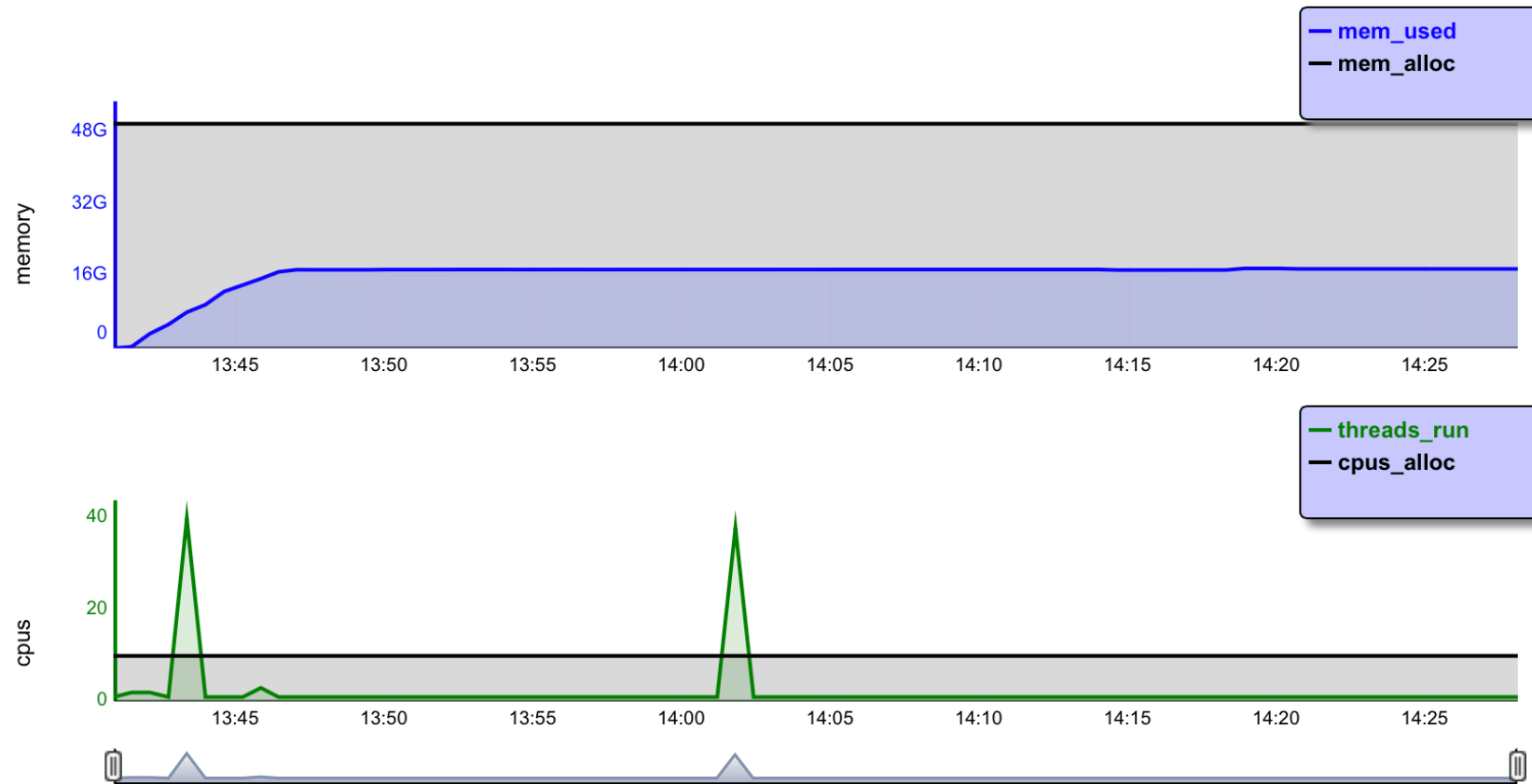


Comment: perfect CPU utilization; underutilized memory but entire node is allocated due to cpu allocation

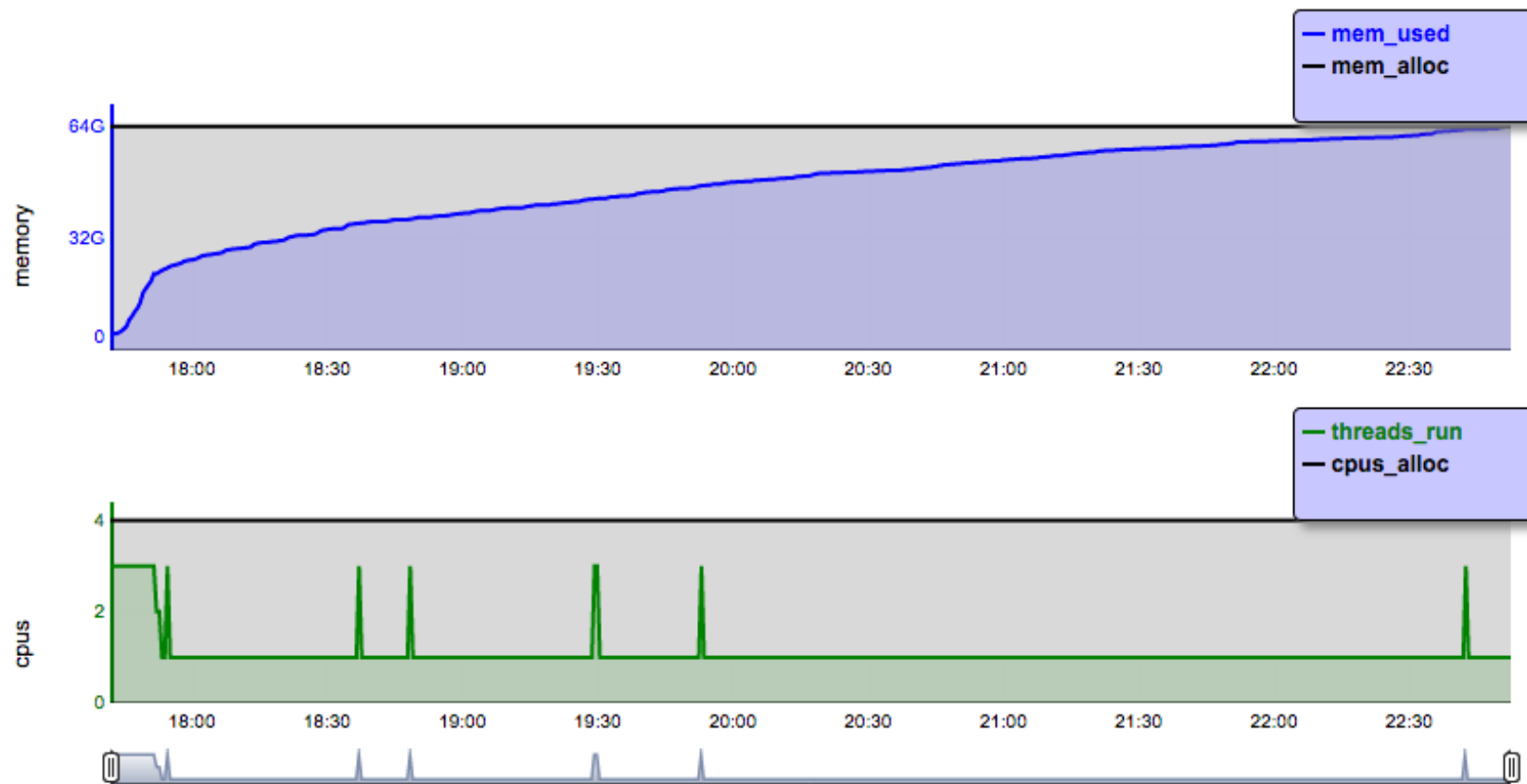


Comment: good overall utilization; possibly split into two jobs with a dependency, and with differing resource allocations

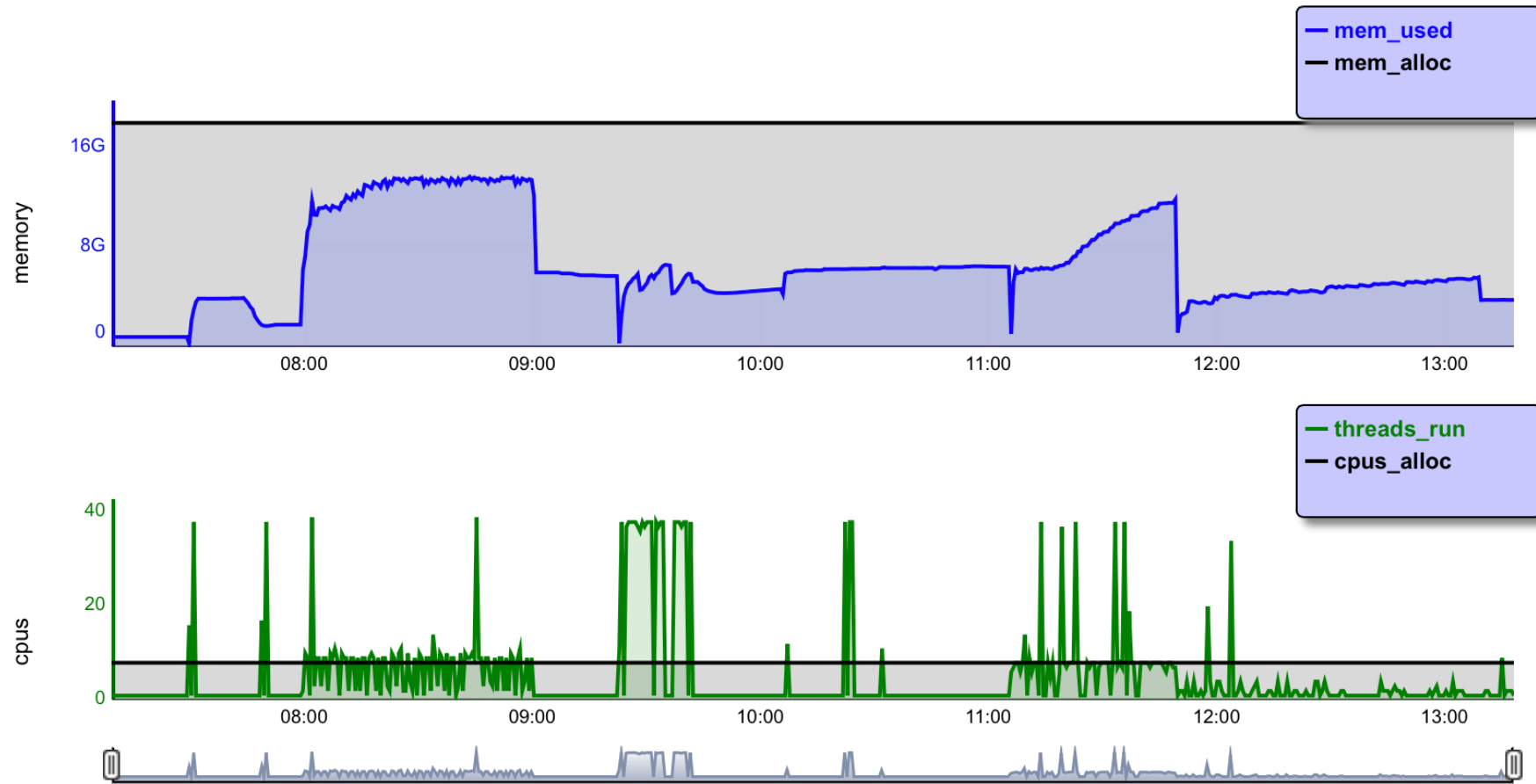




Comment: 8 CPUs too little/too much, 2 would do  
Recommendation: could run in half the memory

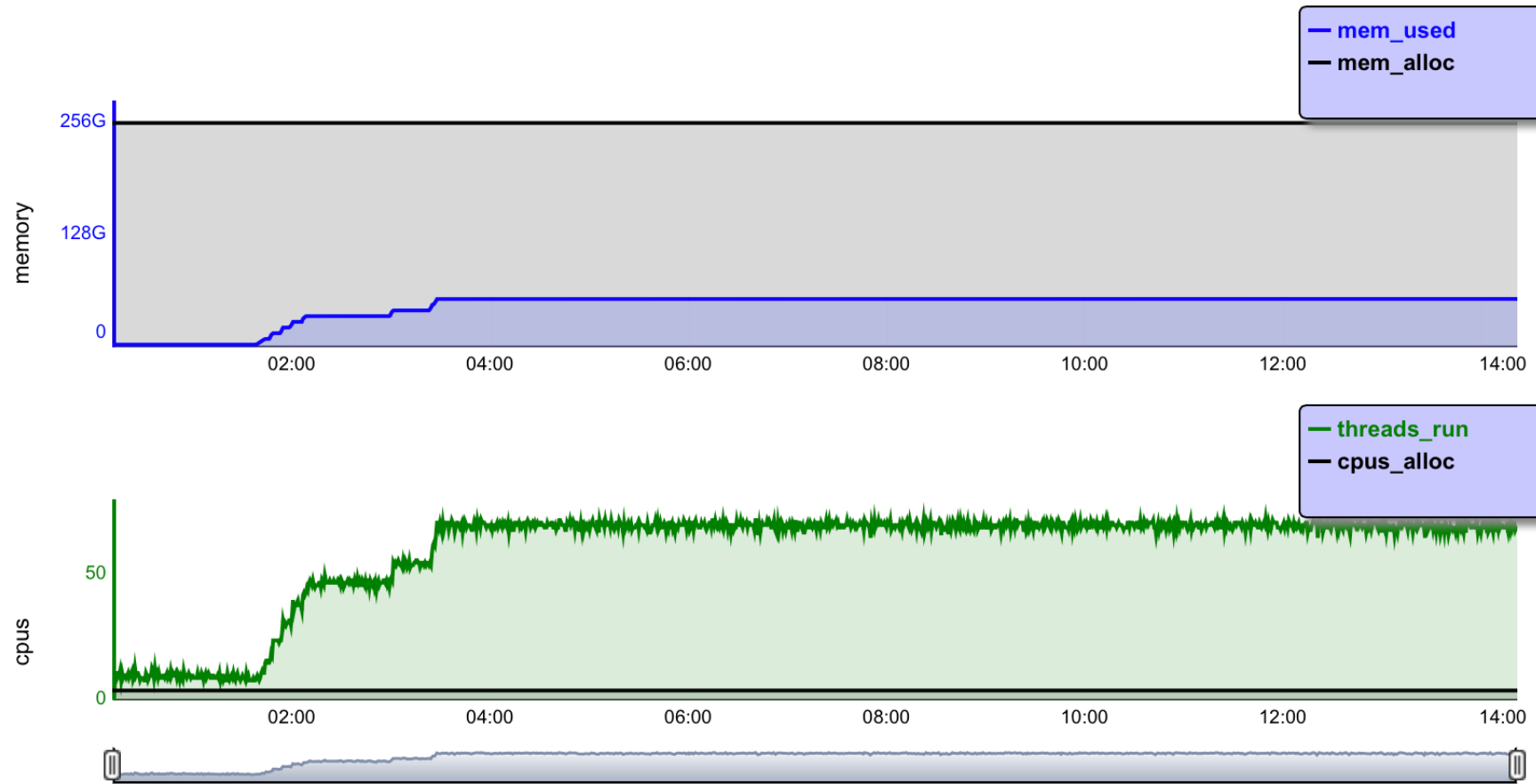


The memory use likely exceeded memory allocation



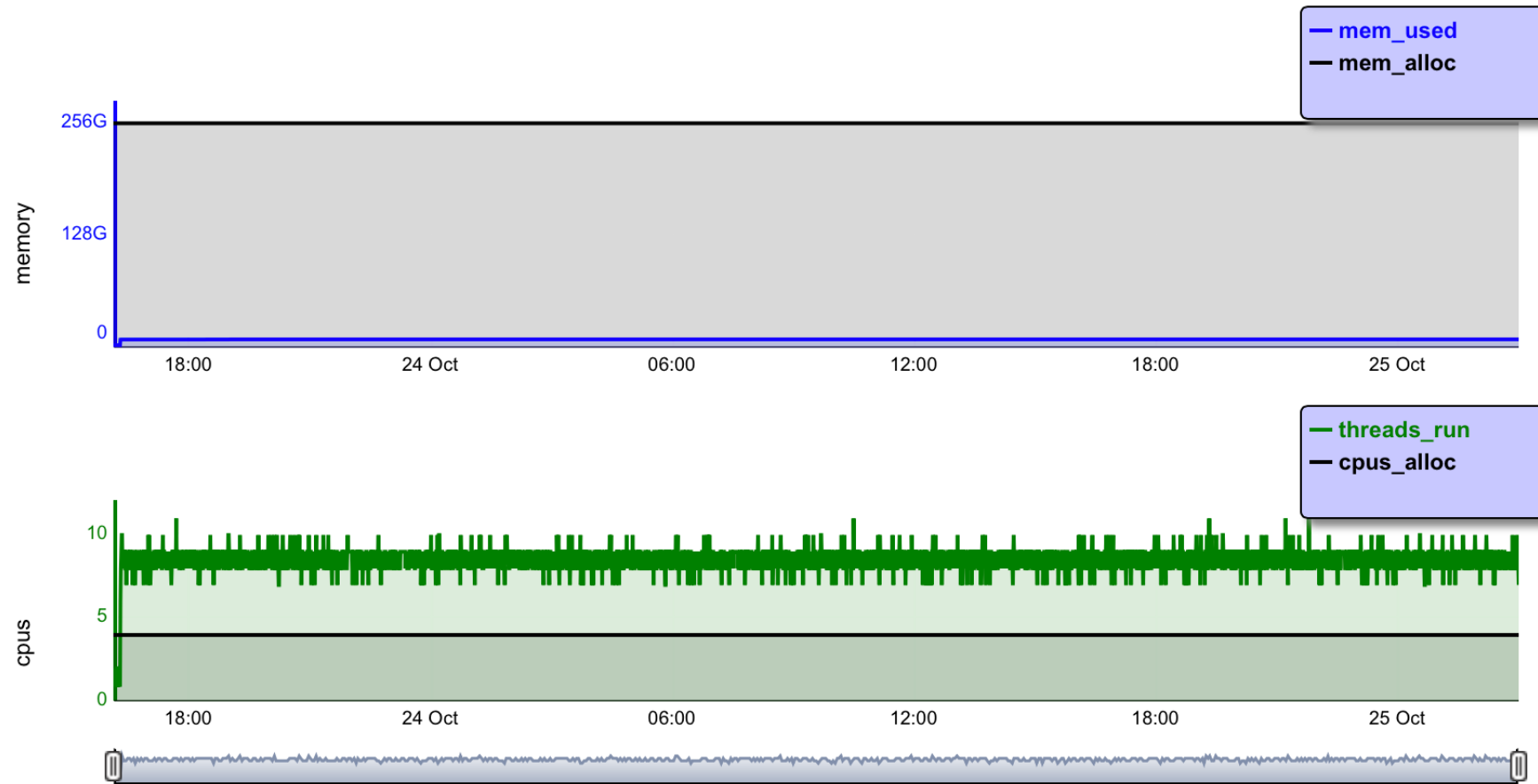
Comment: good memory utilization

Recommendation: increasing CPU allocation probably won't help



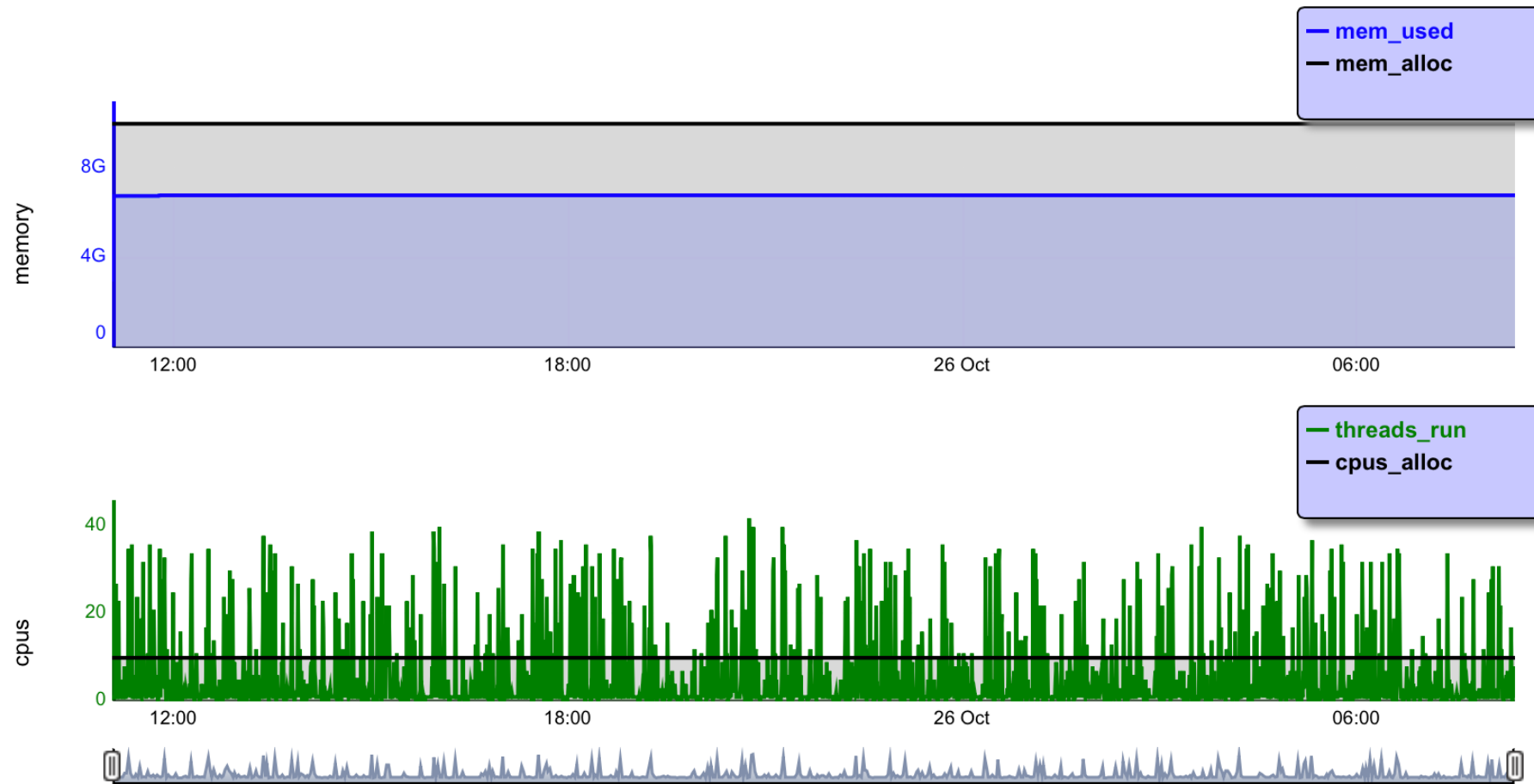
Comment: CPUs badly overloaded

Recommendation: could run in less than half the memory



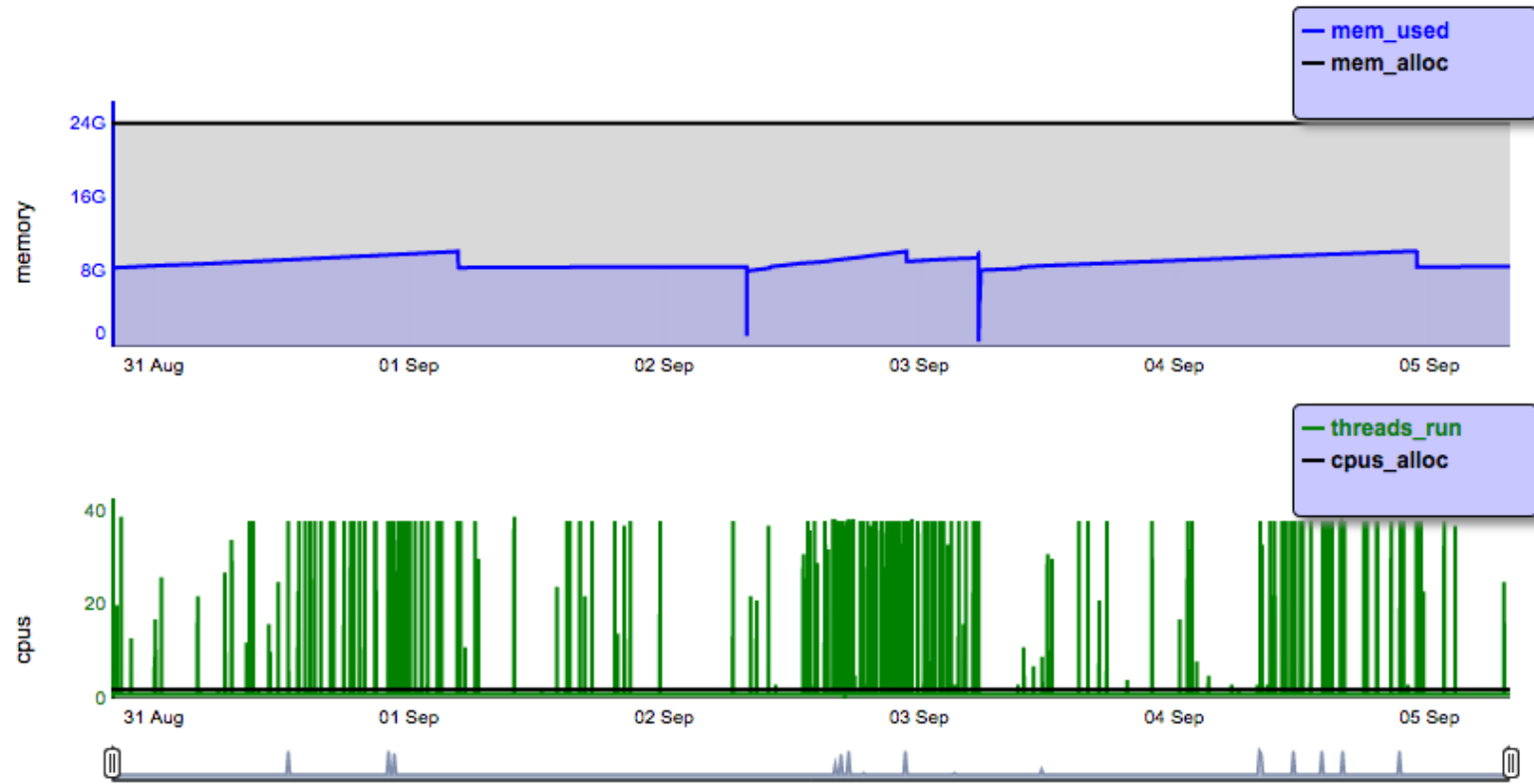
Comment: CPUs overloaded 200%

Recommendation: 256 GB memory allocated, MBs used



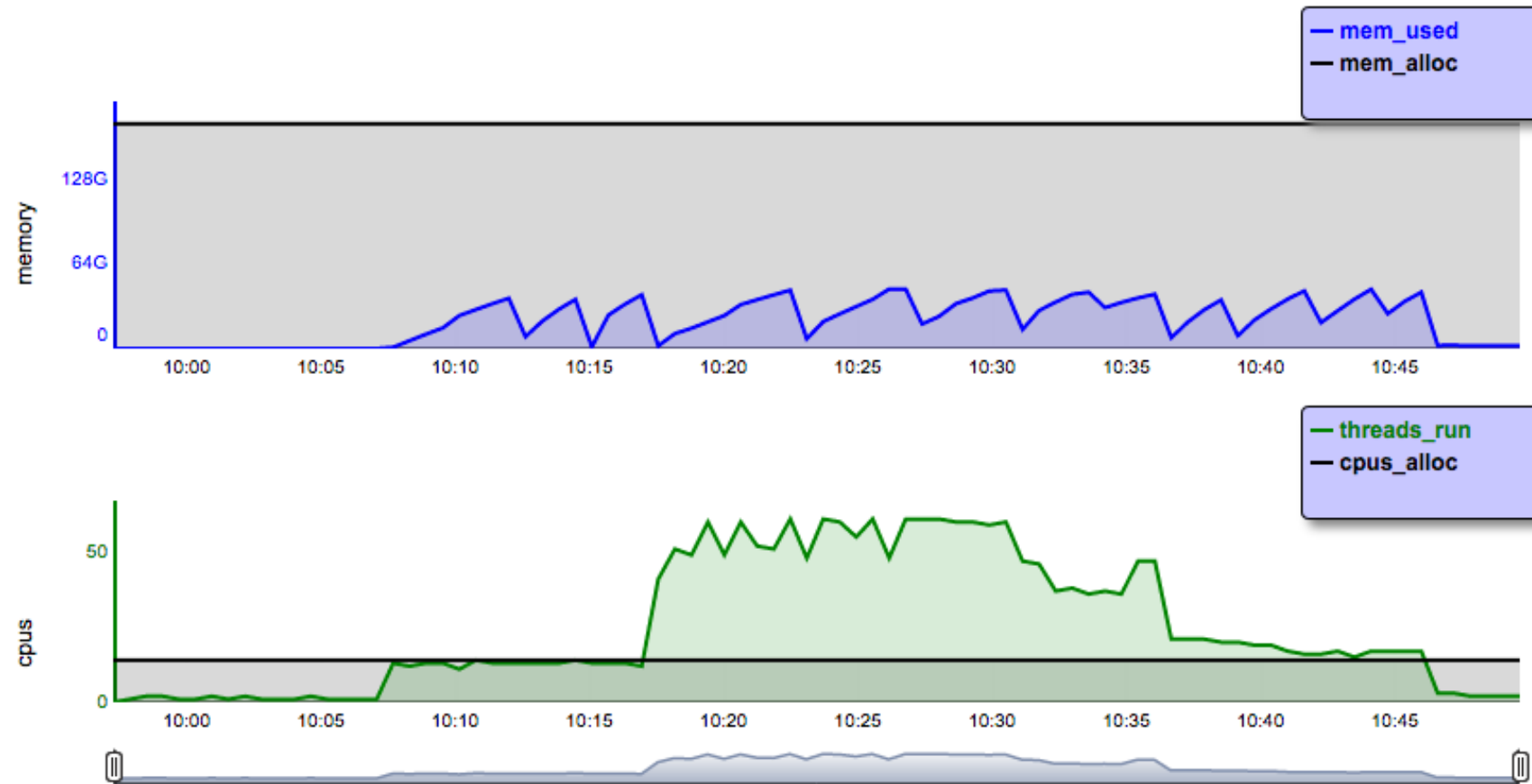
Comment: good memory utilization

Recommendation: might run faster with 32 CPUs?



Recommendation: increasing CPU count might improve?





Recommendation: allocating 56 CPUs would likely help

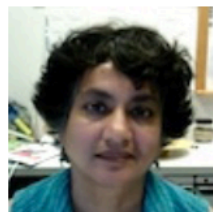
staff@hpc.nih.gov



Steve Bailey



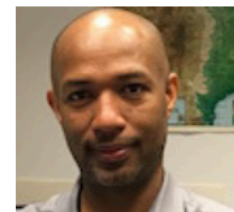
Steven Fellini, Ph.D.



Susan Chacko,  
Ph.D.

Picture  
unavailable

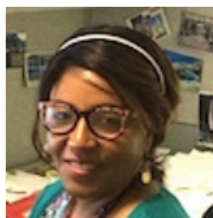
Afif Elghraoui



Ainsley Gibson



David Hoover, Ph.D.



Patsy Jones

Picture  
unavailable

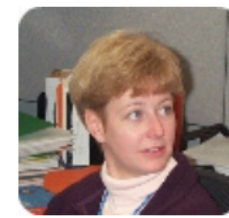
Charles Lehr



Jean Mao, Ph.D.



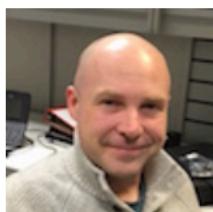
Tim Miller



Charlene Osborn



Mark Patkus



Dan Reisman



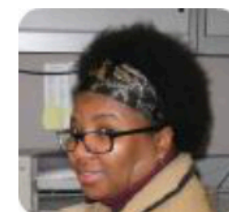
Wolfgang Resch,  
Ph.D.



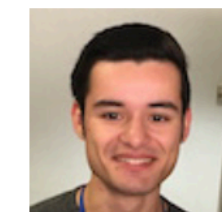
Jerez Te, Ph.D.



Rick Troxel



Sylvia Wilkerson



Michael Harris  
(intern)