# Oslo site report

NT1 AHM 2024-2
Darren, Maiken
Ljubljana 25.09.24-27.09.24

# Compute

# extracts from 2024-1 report - and comments

- Used up the budget for this period (as planned)
    - Seems this is not the case after all. There is 2MNOK left for hardware (compute in the first place).
- Plans: will recreate the cluster before summer going from centos 7 -> almalinux 9
    - Done, with some hickups

# Current compute state

- Finished migration from centos 7 to almalinux 9 on compute nodes, ARC-CE, ARC remote datadelivery servers, and cvmfs and atlas-db squid servers.
- Have also moved the archery update service to a new almalinux 9 NREC instance.
- Had issues with the first iteration of cluster migration, as I kept the small 8 core compute nodes as before
    - Memory seemed to run out causing all sorts of problems
    - Increasing the size to 32 core sorted out all the problems - lesson learned
- Plans to expand compute next year, hopefully keeping up with the requested cpu hours by ATLAS - we are just at the limit now - hovering around the target - and vulnerable to the site not being 100 % full at all times.
    - **probably due to jobs in the slurm queue not fitting into any partly filled node**
    - **maybe some slowness in how ARC ingests new jobs?**
    - **see next slides**

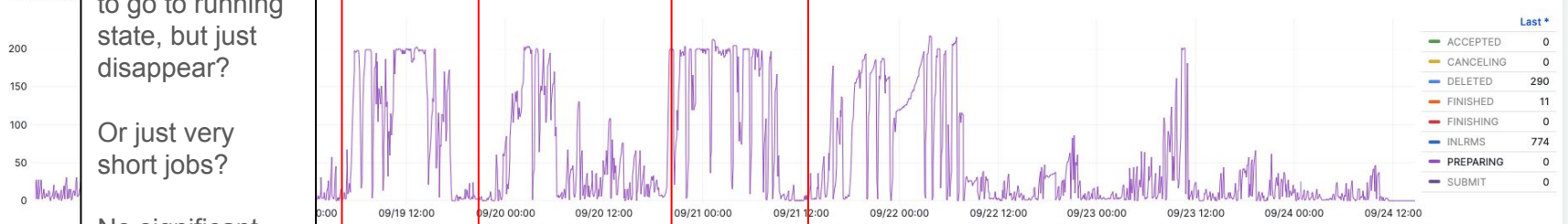# Corecount, Preparing jobs, Queueing jobs



- What is happening with the pending jobs in Slurm?

They seem not to go to running state, but just disappear?

Or just very short jobs?

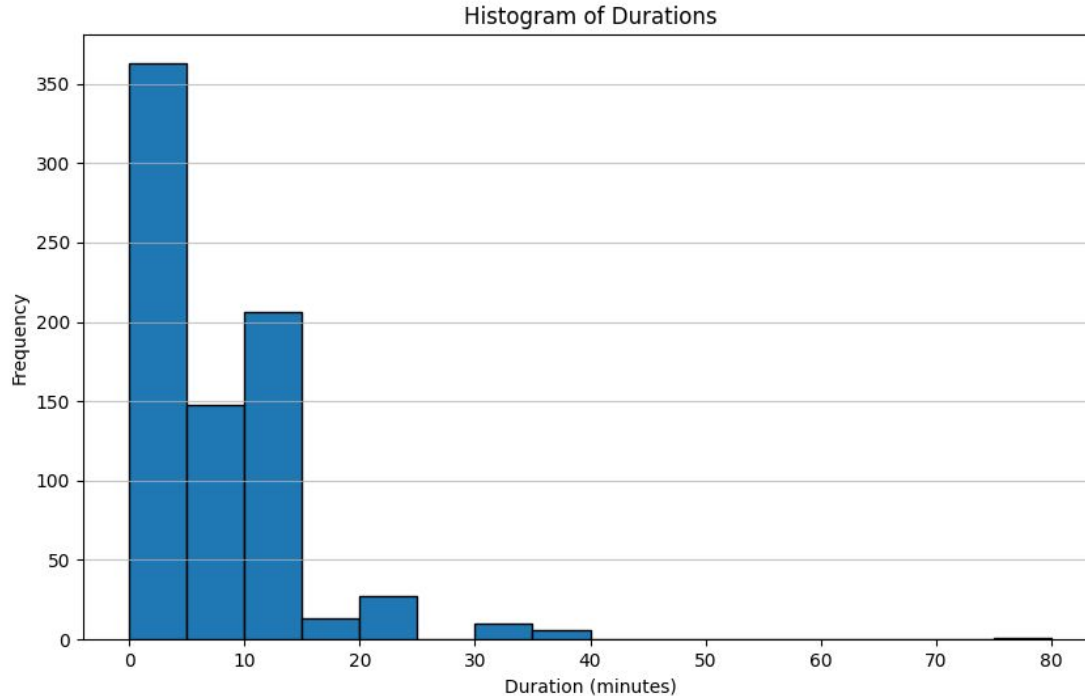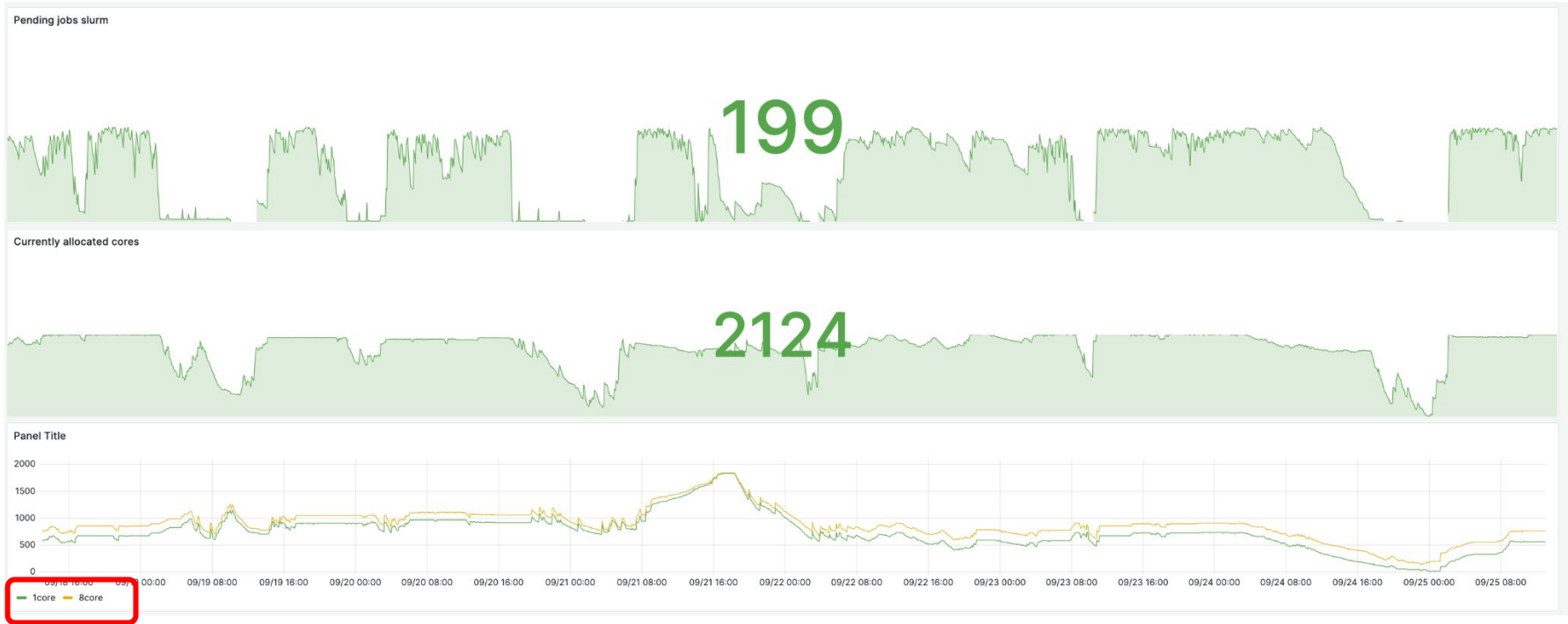No significant number of cancelled jobs according to Slurm.

# Datastaging is typically fast - not the reason for delays filling up cluster



Histogram of Durations

# Does or does this not indicate spikes of short jobs in the core-gaps?

# Seem to have mostly 1 core jobs - so even having 200 pending in Slurm queue does not help much to fill up empty slots

# Preliminary conclusion

- Seems I should just try to allow having many more prepared jobs ready - so maybe 500-1000 jobs, instead of the typically around 200 jobs queued (pending) in the slurm.
- Tweak datadelivery configuration
    - currently I have maxprepared = 768, maxdelivery = 256
    - this is number of files though!
    - data delivery servers 16 x 16 cores = 256 = maxdelivery
    - Turn up to maxprepared=2560 and maxdelivery=512 to see if that helps

# Disk

**State of disk**

Old disk is still in production and carrying most of the UiO weight.
New disk is working for everything but dCache :(
Solutions being investigated
- Short term - possibly partition the disk pools into virtual machines to limit the number of CPU cores to spin lock against
- Remove immediate sync on write from disk pools
- Fix large buffering in Darren's new code.

Actively moving UiO disk operations to UiO storage department as soon as we have a solution for performance issues on dCache threading.

Disk is delivering 2.1PB "reliably"
We will deliver 4+PB "reliably" soon

# Tape

State of tape

We are delivering 10+PB today with no plans of expansion during this cycle.

(UiO internal "ugh" stuff)
We don't have full control over system updates and need to resolve this.

By moving to UiO storage department for operations, we should have better control over outage windows on TSM.

# Network
# no change!