

# Switch to Perlmutter

The good, the bad and the ugly

Sergiu Weisz  
sergiu.weisz@cern.ch

# Quick refresh

- In 2019 we started Cori's integration
- Being run by hand since 2020
- Running custom version of the software

# Issues with Cori

- Needed to be run by hand, no way to run as a service
- CVMFS java version did not work
- Could not run newer version of MonaLisa
- General lack of availability
- Network bottlenecks (debugged later)

# New toy: Perlmutter

- More CPU nodes (3072 > 2388), 2 x 64 cores, 512GB RAM
- 1536 GPU nodes with 4 x A100, 1 x 64 cores, 256GB RAM
- More straightforward network
  - Slingshot works over Ethernet, no transceivers needed
- No more service nodes
  - They can't guarantee service nodes availability

# CE move to Spin

- Spin is a Rancher based Kubernetes cluster manager
- CE will run in container
- Can't submit jobs through SLURM, need to use SuperFacility API

# Conditions for running CE

- Connection from Central Services to CE
  - Can be done through ingress
- Connection from nodes to CE for MonaLisa (udb/8884) and from CS to CE (tcp/1094)
  - WNs have access to Spin machines on all ports
- Shared file system access would be a big plus
- Connection from WNs to outside

# What we need to do

- Build a container that has what's needed to run CE
  - Need to have CVMFS
- Write the interface to use SuperFacility API for jobs submission

# What we need inside the CE container

- CVMFS (done)
- Host certs (done)
  - Can't do mkdir in `/root` in container
  - Should move to persistent volume claim `.config`
- Access to CFS to use renewed hostcerts
  - They require you to run root-less, but for this, we need additional setup
  - Costin: do we need a UNIX username associated with the ID?
  - Can't do useradd in init script

# What we need for SuperFacility API

- Write authentication flow into the code (Done)
- Wrap the existing SLURM BQ code in authenticated HTTP requests
  - Will be using the command execution interface, which is discouraged by NERSC
- Needed API token for write/execute two weeks ago
  - Requires more blood sacrifices

# Conclusion

- Two working threads
  - CE requires some more tuning work to get up and running
  - The SuperFacility API looks like it will do the job well
  
- The Good:
  - CE should work seamless after first setup, and won't need outside intervention
- The Bad:
  - Don't have a way to easily debug for outside users
- The Ugly
  - Bureaucracy, bureaucracy everywhere