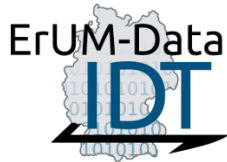


Opportunistic extension of a local compute cluster with NEMO resources for HEP workflows

Stefan Kroboth, Michael Böhler, Anton J. Gamel, Benjamin Rottler, Markus Schumacher

Physikalisches Institut, Albert-Ludwigs-Universität Freiburg

14th Annual Meeting of the Helmholtz Alliance
"Physics at the Terascale"
24.11.2021



Bundesministerium
für Bildung
und Forschung



ATLAS-BFG (Tier 2/3)

- ▶ Compute center for the ATLAS experiment
- ▶ ≈ 3200 cores (25 kHS06) / ≈ 2.5 PB storage (dCache)
- ▶ ATLAS production/analysis jobs
- ▶ Jobs by users of local HEP research groups
- ▶ Scheduler: SLURM



bwForCluster NEMO

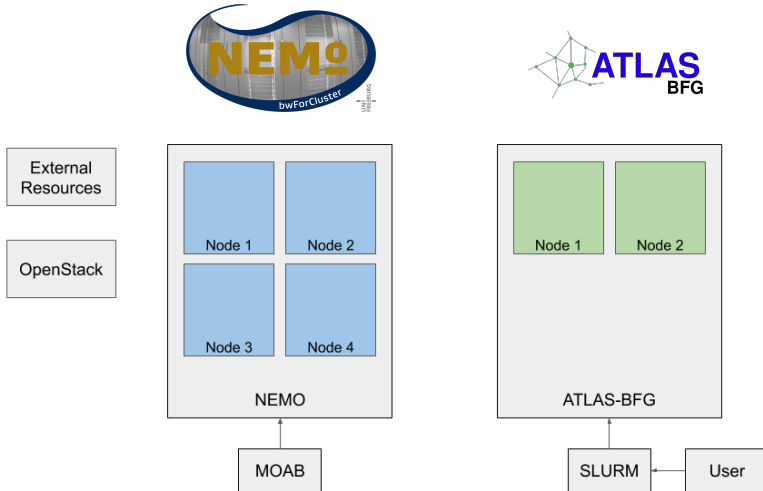
- ▶ HPC cluster at Freiburg University
- ▶ ≈ 18000 cores (340 kHS06) / ≈ 800 TB storage (BeeGFS)
- ▶ Different software setup than ATLAS-BFG
- ▶ Scheduler: MOAB



COBaID/TARDIS

- ▶ Opportunistically integrates resources from NEMO into ATLAS-BFG
- ▶ Based on demand and availability

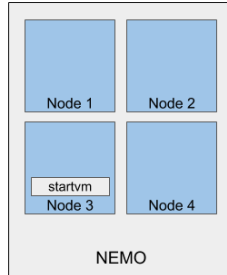




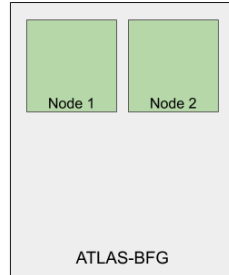


External
Resources

OpenStack

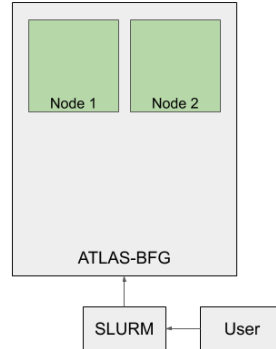
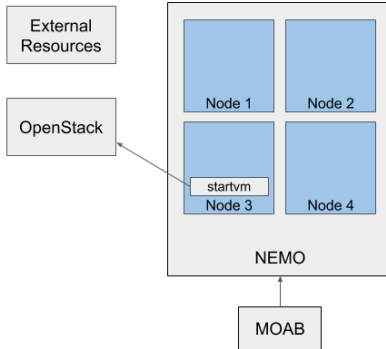


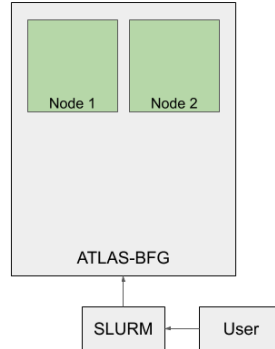
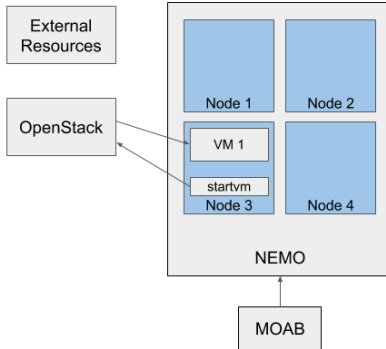
MOAB

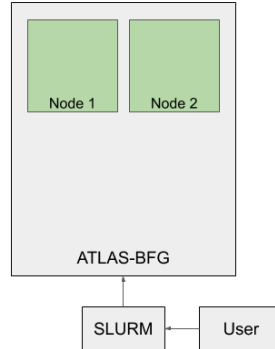
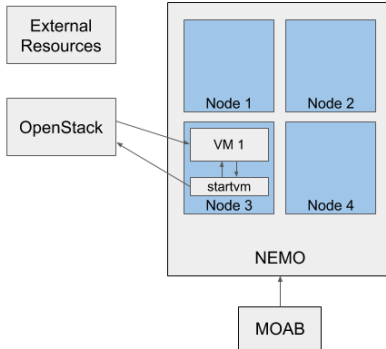


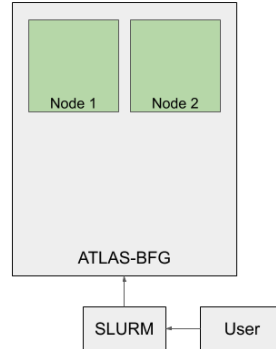
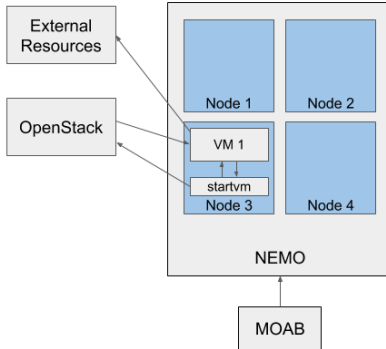
SLURM

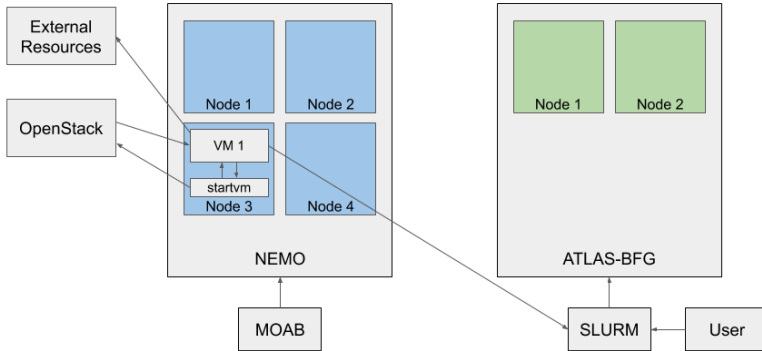
User

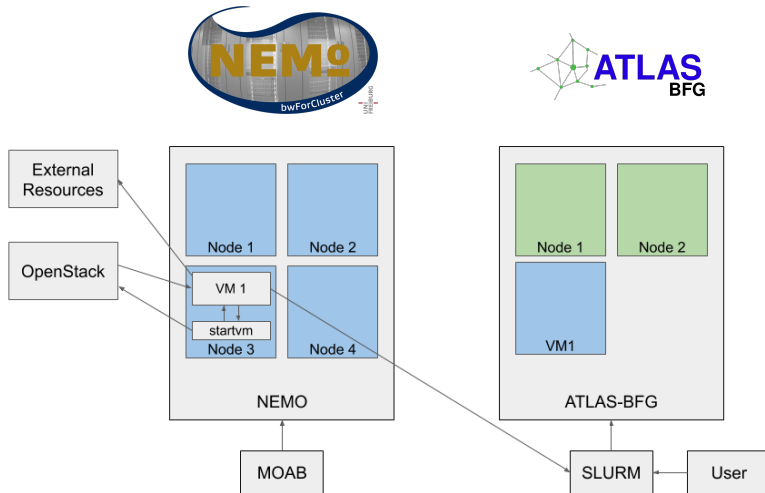


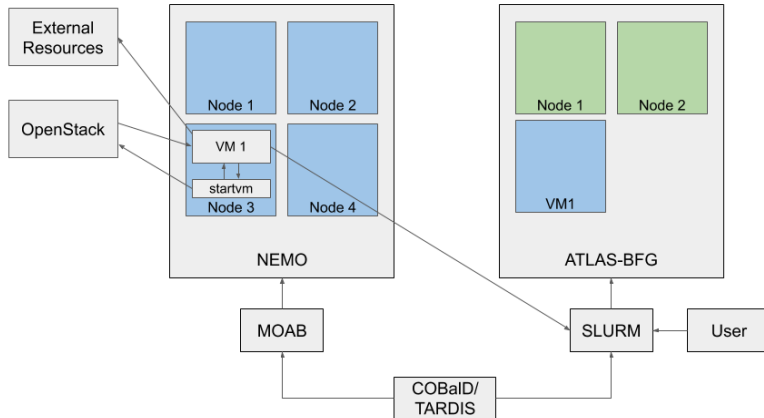












COBaID

- ▶ COBaID is an Oppportunistic Balancing Daemon
- ▶ Decides whether to request or release resources



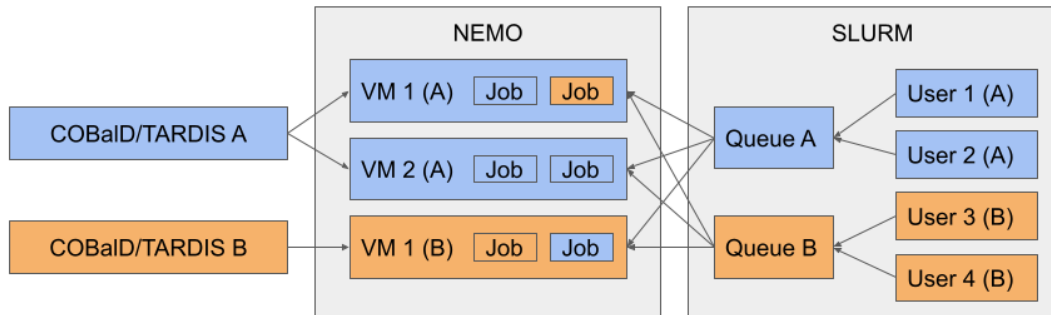
TARDIS

- ▶ Transparent Adaptive Resource Dynamic Integration System
- ▶ Interaction with the batch systems
- ▶ Resource integration



Main developers from group of Prof. Günter Quast (KIT)
<https://github.com/MatterMiners>

- ▶ Four local HEP research groups (A to D) with a share in NEMO
- ▶ Each served with its own COBaID/TARDIS instance
- ▶ Each has its own SLURM queue/partition
- ▶ Efficient use of resources due to sharing VMs across HEP groups

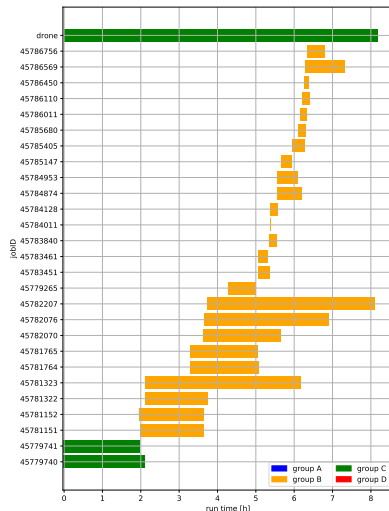


- ▶ Monitoring (Grafana, Elasticsearch, Prometheus, Vector)

- ▶ Four local HEP research groups
- ▶ Each served with its own COBaID
- ▶ Each has its own SLURM queue
- ▶ Efficient use of resources due to

COBaID/TARDIS A

COBaID/TARDIS B



SLURM

User 1 (A)

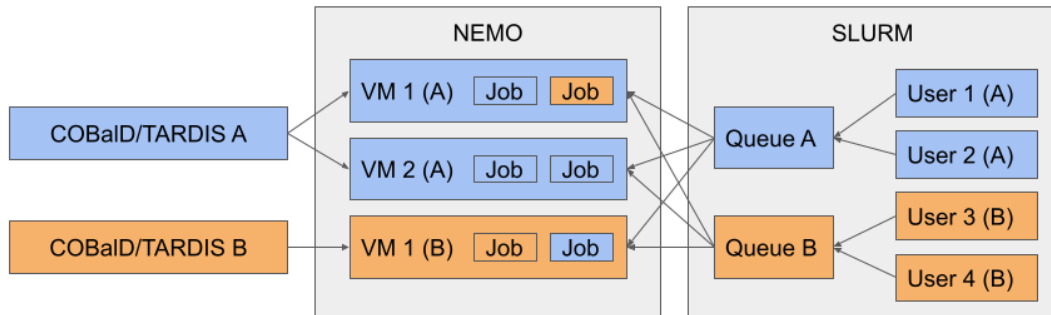
User 2 (A)

User 3 (B)

User 4 (B)

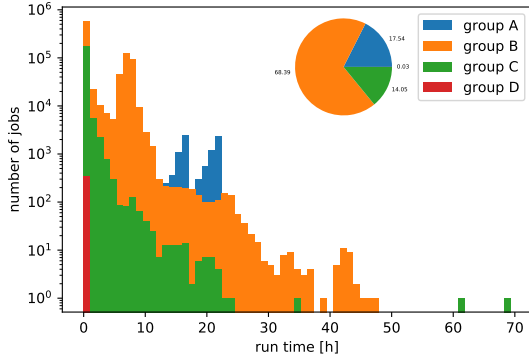
- ▶ Monitoring (Grafana, Elasticsearch, Prometheus, vector)

- ▶ Four local HEP research groups (A to D) with a share in NEMO
- ▶ Each served with its own COBaID/TARDIS instance
- ▶ Each has its own SLURM queue/partition
- ▶ Efficient use of resources due to sharing VMs across HEP groups

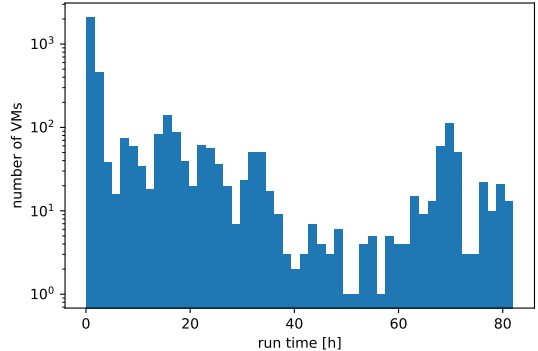


- ▶ Monitoring (Grafana, Elasticsearch, Prometheus, Vector)

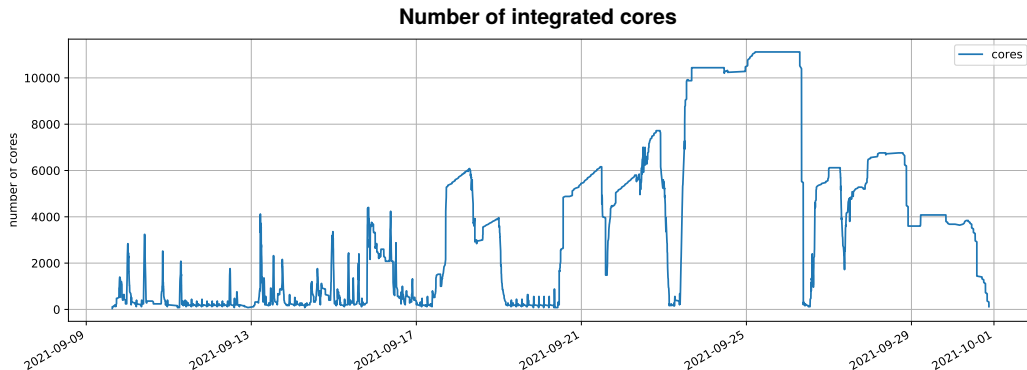
Job runtime per group



VM runtime

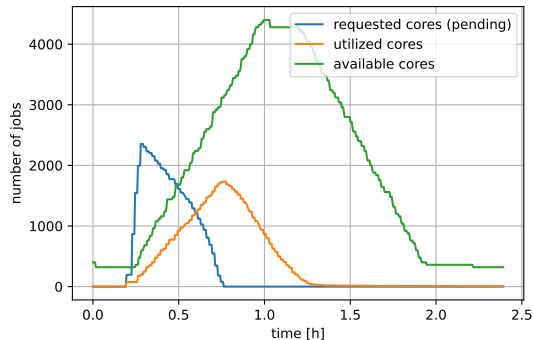


- ▶ Cluster usage varies heavily across HEP groups
- ▶ Long running and fully utilized VMs are preferable for efficiency

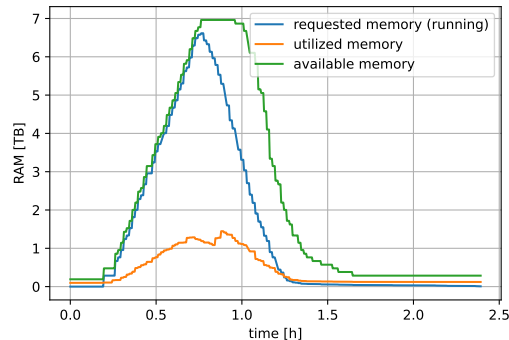


- ▶ Resources are only integrated when needed and released otherwise
- ▶ Efficient use of resources across cluster boundaries

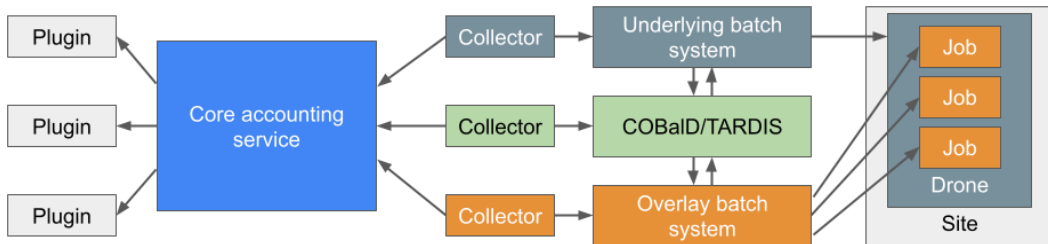
CPU cores



RAM



- ▶ Some CPU cores running idle because of RAM requirements of jobs
- ▶ Steepness of increase in available resources could be a measure for the responsiveness of the setup
 - ▶ Influenced by: COBaID/TARDIS settings, NEMO utilization, ...



- ▶ Opportunistic resources are hidden behind other resources
- ▶ Track usage of all resources
- ▶ Can be used for
 - ▶ Accounting with an experiment such as ATLAS/CMS/...
 - ▶ Inter-site billing between sites sharing resources
 - ▶ Matching fairshare/priority to provided resources
- ▶ Currently in development
 - ▶ Core services and prototype of framework for collector and plugins ready
 - ▶ Prototype of collector plugin for COBaID/TARDIS ready
 - ▶ Next step: test fairshare/priority matching in production

Conclusion

- ▶ Integrated NEMO resources into ATLAS-BFG in an opportunistic manner transparent to the users
- ▶ Successfully in operation for more than two years
- ▶ Providing an additional $2 \cdot 10^6$ CPU hours since beginning of 2020

Outlook

- ▶ Integrate other resources such as Cloud providers (bwCloud)
- ▶ Accounting
 - ▶ Enable prototype in production environment: adjust group priority on ATLAS-BFG (fairshare) based on provided resources per group on NEMO
 - ▶ Develop further items necessary for a complete accounting framework