# A Cloud-based Grid Computing Site

Jonathan Woithe [1]    Martin Sevior [2]    Paul Jackson [1]    David Dossett [2]
Marcus Ebert [3]

[1]University of Adelaide, Australia
[2]University of Melbourne, Australia
[3]University of Victoria, Canada

# Outline

# Grid Computing Overview

▶ Grid Computing: provide Compute Elements (CE) and Storage Elements (SE) for High Energy Physics experiments

▶ Resources distributed across the globe. For the ATLAS experiment:

- Storage: Disk: 400 PB disk, 600 PB tape

- CPU hours: 300 million per month

- Average data transfer throughput: 50 GB/s

▶ AU-Melbourne Grid site: CE and SE for two projects, ATLAS and Belle II

# Motivation

▶ Grid sites traditionally use bare metal and specialised filesystems

▶ Institutional research computing infrastructure is increasingly cloud-based
  → Need to use what is readily available

▶ Want to use industry standard interfaces
  → Avoid esoteric filesystems which require domain-specific knowledge

▶ Exploit economies of scale from cloud resource providers for Grid computing

▶ Easily increase compute and storage as funding allows and demand grows

# Intrastructure Description

- ▶ VMs provided by Melbourne Research Cloud (MRC)

- ▶ Orchestration by OpenStack

- ▶ Server configuration managed by Ansible, tracked in git

- ▶ All VMs run AlmaLinux 9

▶ 750 TB of S3 compatible object store

- Not a traditional filesystem

- Each "file" is an object in a database

- The object's "key" is interpreted as its filesystem path

- No explicit objects for filesystem directories

- ▶ Use a single bucket for flexibility

- ▶ Belle II and ATLAS have separate key namespaces

  - • Gives illusion of separate top-level directories

- ▶ XRootD serves data from S3 though davs, https, root
  (https://xrootd.org/)
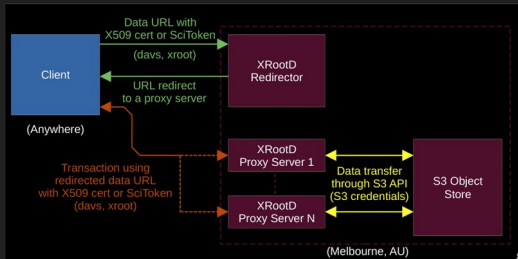
- ▶ Enabled by the XRootD-s3 work at SLAC
  (https://cds.cern.ch/record/2857626/files/ATL-SOFT-SLIDE-2023-125.pdf)

▶ XRootD redirector VM (2-core, 8 GB RAM)

  • Authenticates incoming request

  • Redirects request to a proxy server

▶ XRootD proxy server VM (8-core, 32 GB RAM)

  • Authenticates access request, serves requested resource

▶ Currently have 2 proxy servers. Add more to increase throughput as needed.

# Infrastructure Description

Storage: XRootD

▶ Storage Resource Reporting (SRR) JSON file

- Defines Belle II and ATLAS storage shares, space usage and capacities

- Generated hourly by bash script using s3cmd on an XRootD proxy server

▶ Adler32 checksums (managed on an XRootD proxy server)

- Evaluated on first request by python script using boto3 library for S3 access

- Checksum stored as metadata attribute on S3 object for later reference

# Intrastructure Description

▶ Third Party Copy

- Data copy between two SEs initiated by a third party

- Processed by XRrootD proxy servers

- `davs://` transfers: handled by internally by XRootD

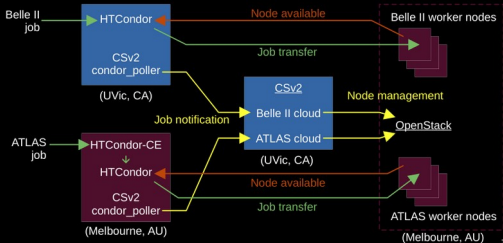- `root://` transfers: use pipelined xrdcp and s3cmd processes

▶ Slightly different architectures used for Belle II and ATLAS

▶ Cloud resources managed by Cloud Scheduler v2 (CSv2) instance at UVic
(`https://csv2.heprc.uvic.ca`)

# Intrastructure Description

Compute

▶ HTCondor controller host VM at UVic

▶ Jobs submitted to HTCondor by local DIRAC site-director
(No need for authentication by HTCondor-CE)

▶ CSv2 monitors HTCondor, starts/stops worker VMs to match demand

▶ Worker VM setup by CSv2 via cloud-init, notifies HTCondor when ready

▶ HTCcondor runs job on selected worker VM (8-core, 32 GB RAM)

▶ A MRC VM runs HTCondor and HTCondor-CE (8-core, 32 GB)

  • Token authentication is used

▶ Jobs submitted by PanDA to HTCondor-CE on the HTCondor host

▶ After authorisation, jobs are sent to HTCondor by HTCondor-CE

▶ CSv2 processes proceed as for Belle II

▶ Belle II storage is operational (400 TB, 19 TB used as of 8 Oct 2025)

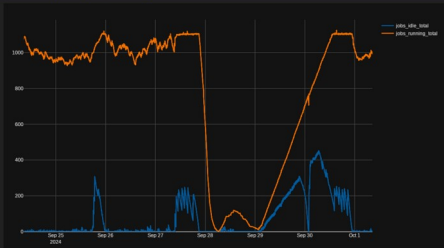▶ Belle II compute is operational (900 vCPUs)

# Current status

## Belle II - TPC matrix



(from https://people.na.infn.it/~spardi/tpc-davs-latest.html)

# Current status

▶ ATLAS storage is operational (350 TB, 271 TB used as of 8 Oct 2025)

▶ ATLAS compute is operational (432 vCPUs)

- Looking to add additional vCPUs to ATLAS pool

# Benchmarks

|                | Within cloud | In Australia |
|----------------|:------------:|:------------:|
| davs:// read   | 108 MB/s     | 40 MB/s      |
| davs:// write  | 123 MB/s     | 74 MB/s      |
| Checksum calc  | 3.2 s        | 3.4 s        |
| Checksum fetch | 0.72 s       | 0.98 s       |
| s3 read        | 213 MB/s     | n/a          |
| s3 write       | 165 MB/s     | n/a          |
| root:// read   | 6.6 MB/s     | 5.9 MB/s     |
| root:// write  | 132 MB/s     | 70 MB/s      |

Read/write tests used gfal-copy, checksum tests used gfal-sum. s3 tests were run on an XRootD proxy server.
Results are the average of 5 tests, each using a 1 GB test file.

# Challenges

- ▶ Invisible application firewalls

- ▶ Slow `root://` read

- ▶ Read/write speed variability, particularly outside Australia

- ▶ Shift to AlmaLinux 9 environment

# Future plans

▶ Increase storage and compute resources as funding allows.

- Tentatively planning for an additional 1 PB, mostly directed towards ATLAS

- Add 1000 vCPUs to ATLAS pool

▶ Monitor transfers for Belle-II and ATLAS, add extra proxy servers as needed

# Conclusions

▶ We have built a grid site using cloud storage and compute in the MRC

▶ The "AU-Melbourne" Grid site is in production

- CE and SE resources are provided for ATLAS and Belle-II

▶ It is possibly the first production Grid site with cloud-based CE and SE