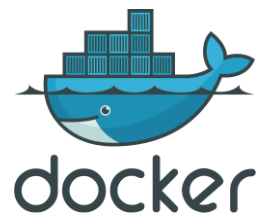


GridPP49

Tier1 Condor Batch Farm

Tom Birkett

STFC, Rutherford Appleton Laboratory



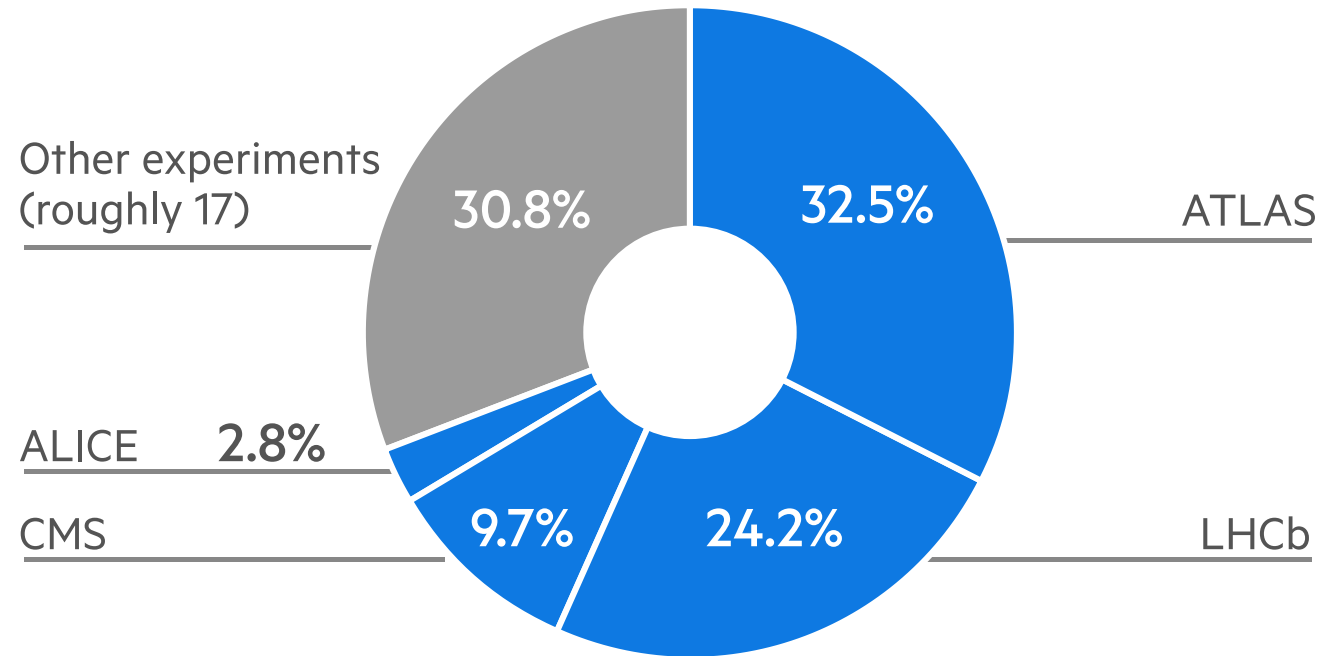
Agenda

- Introduction
 - What is the Batch Farm?
 - Who are the stakeholders?
- Job lifecycle
- High level activities from last year
- Upcoming projects



Overview

- RAL is a Tier-1 primarily for the four large LHC experiments. The graphic shows the fair-shares for the farm.
- RAL also supports roughly 17 other experiments

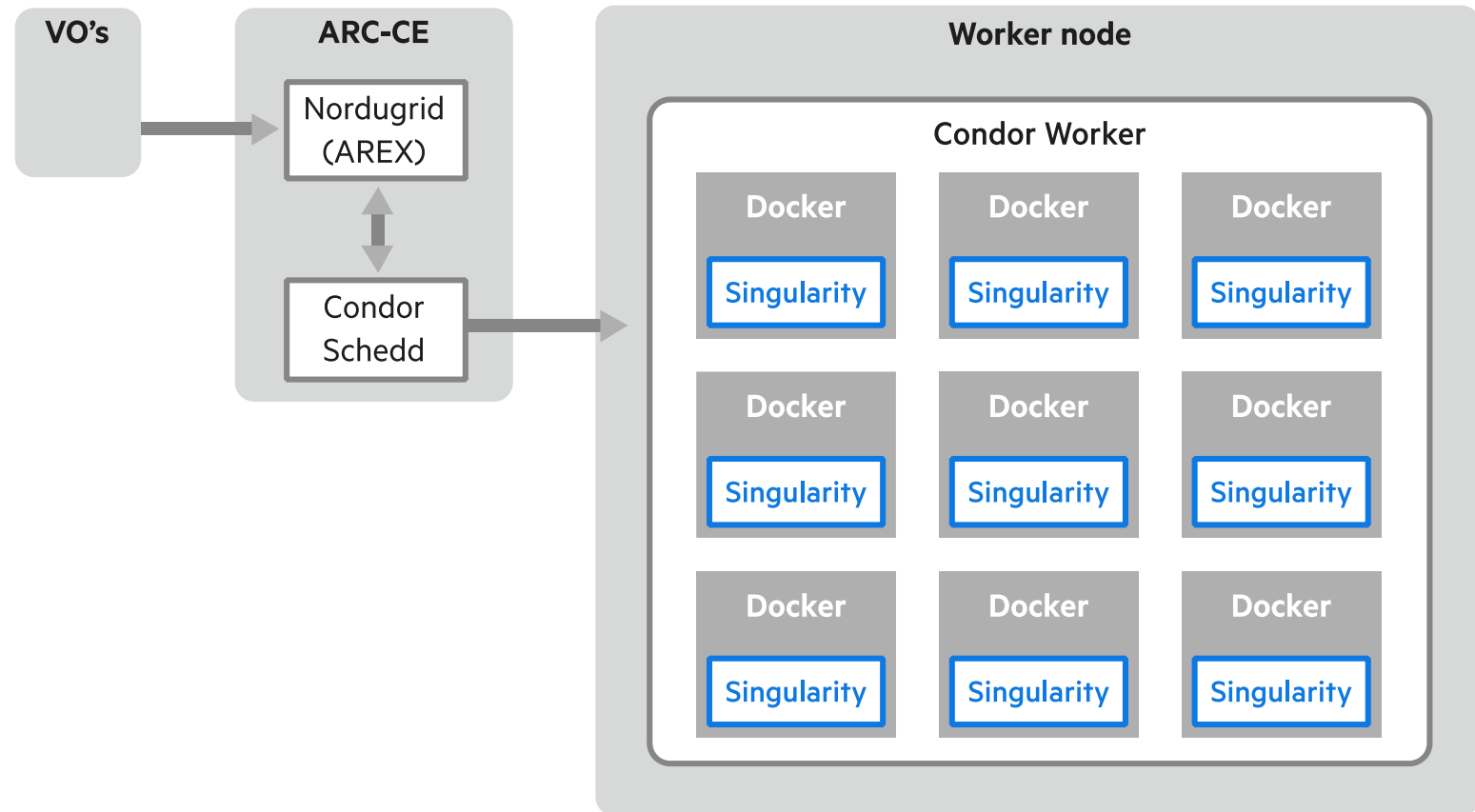


Batch Farm background

- Batch Farm at RAL:
 - Currently 423 worker nodes in production.
 - ~ 49,000 logical CPU cores
 - HEPSPEC: ~ 630,000
- RAL runs a mix of multicore and single core jobs

Job lifecycle

- Experiment submits job
- Job description is read by ARC-CE and given to Condor Schedd
- Worker (Startd) picks up job from schedd, starts Docker container
- Job inside container runs Singularity / Apptainer once payload has been downloaded.



The past year

- HTCondor upgrade
 - 8.8 -> 9.0
- Singularity upgrade to Apptainer
- "Zombie" containers
- SL7 to Rocky 8 for Condor Startds and Managers
- Docker 23

Upgrading to Apptainer

- What drove the Singularity/Apptainer upgrade?
 - Apptainer moved into Linux Foundation
 - Current running version: Singularity 3, straightforward upgrade path
 - Support Singularity 3 is being sunsetted
 - Continued security fixes and improvements
- https://apptainer.org/docs/admin/main/singularity_migration.html

Rocky 8

Reasons for upgrading to Rocky 8:

- Condor 9.0 modules written in Python 3 had better support on EL8
- The new 2021 generation of Batch Farm Worker-Nodes use the AMD EPYC 7763 64-Core Processor. Rocky 8 gave us optimal compatibility with this chipset. (Using ML branch kernel)
- Security support of the OS until 2029
- Risk of complications to jobs is reduced due to use of containerised job environments. Jobs run in containerisation and will remain running in a EL7 environment while the host OS ran Rocky 8
- Mature OS, Rocky 9 had only just been released and was untested.

Upgrading Docker / "Zombie" containers

- End of December 2022, RAL experienced an influx of containers that caused the workernode to hang.
- After investigation, two issues were identified:
 - Incorrect configuration of VOMS endpoint leading to lockup of CVMFS client
 - Bug in Docker causing containers to remain if process in container has a PID of 0

Condor 9 upgrade

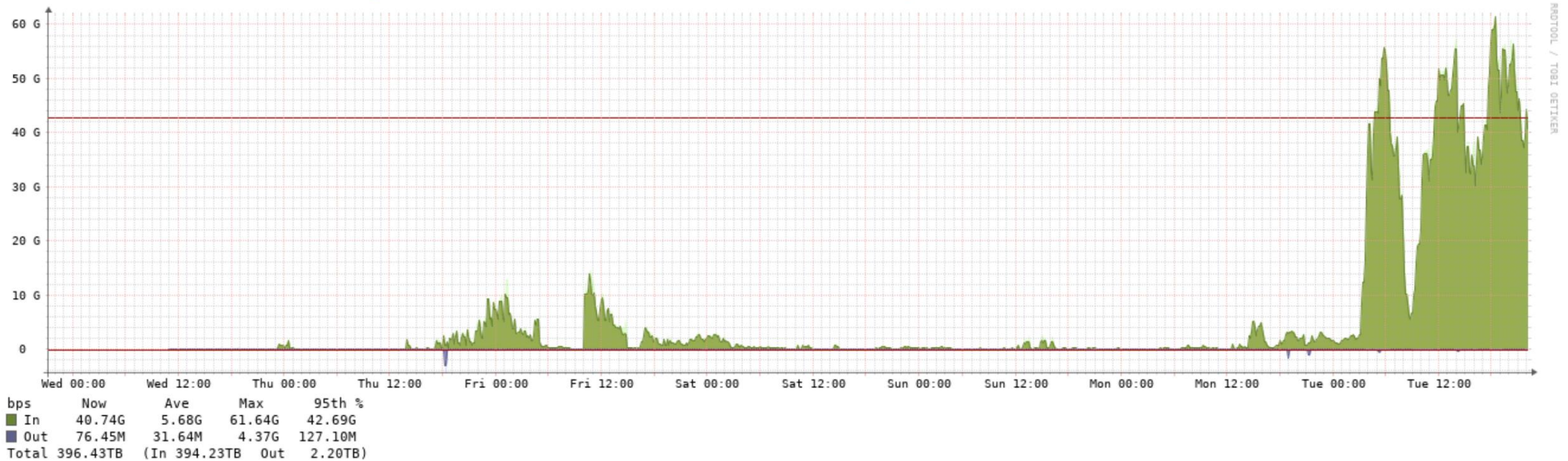
- The Workernodes and Central Managers upgraded to Condor 9
- Improved security
 - Daemon communication encrypted by default
- LTS support at time of upgrade
- Token support for daemon authentication
- New `condor_watch_q` tool that efficiently provides live job status updates
- Submitter ceilings allow administrators to set limits on the number of running jobs per user across the pool

LHCONE

Advertised the 2020 and 2021 generation of workernodes to the LHCONE. Go-live this time last week

From To

[Hide Legend](#) | [Show Previous](#) | [Show RRD Command](#) | [Zoom to Port Speed](#) | [To show trend, set to future date](#)



Future work / plans

- HTCondor upgrade
 - 9.0 -> 10.0
- Docker 23
- ARC 7 (Tokens)
- Removal of Docker wrapper.
- Rebuild fair-share config.
 - Dedicated partition for ATLAS
- IPv6

Docker 23 Rollout

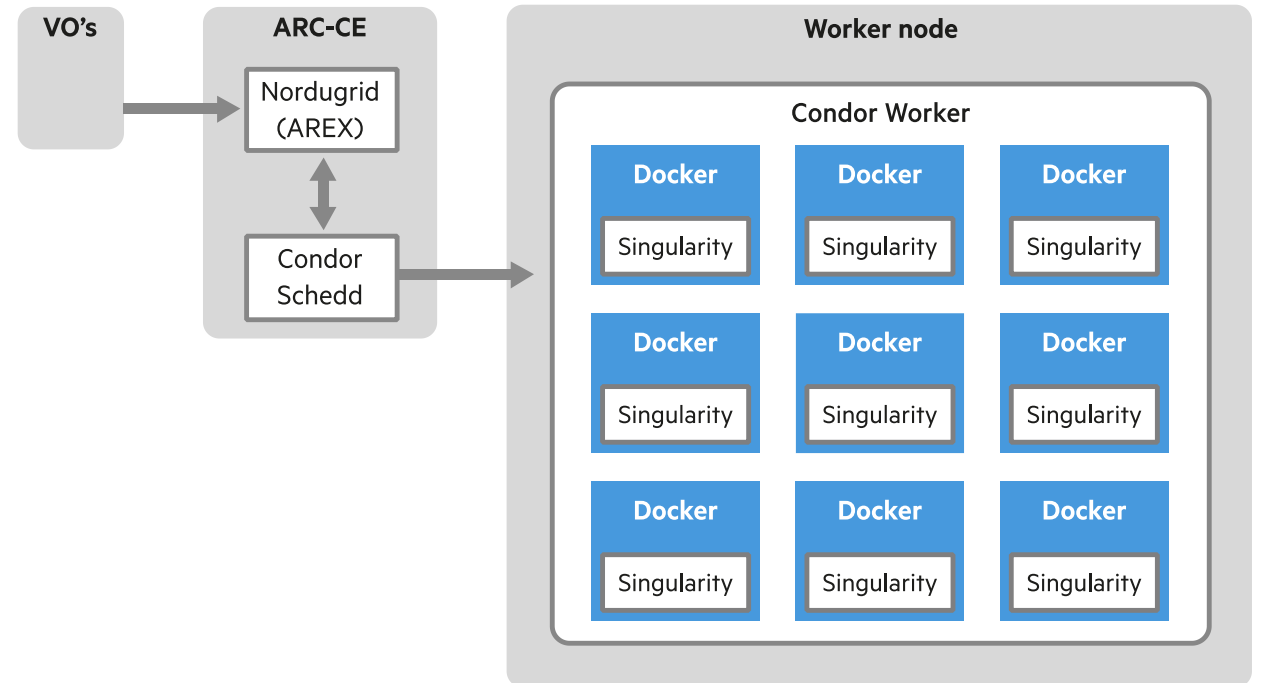
Worker-nodes need to be drained / empty of jobs.

Half of farm is drained in preparation of upgrade.

Once empty, commence upgrade of Docker 23

Bring nodes back slowly and monitor.

Repeat for remaining half of farm



Docker Wrapper sunseting

- The `docker.py` wrapper used at RAL, builds the docker run command that is executed by Condor.
- These include:
 - PANDA environment variables
 - Apptainer environment variables
 - Security opts
 - Memory reservations
- To remove a layer of complexity, the function of the `docker.py` wrapper can now be achieved using native ClassAd language in Condor

IPv6 on Workernodes and Docker

- Docker natively supports IPv6. Implantation is a topic of debate.
 - Docker does not support SLAAC assignment out the box, this can be limiting
 - Requirement to statically assign IPv6 subnet in Docker config (/etc/docker/daemon.json):

```
{  
  "ipv6": true,  
  "fixed-cidr-v6": "2001:db8:1::/64"  
}
```

ARC 7

- Support for Tokens!
- Remove hack to get tokens functional.
 - ARC 6 hard codes a requirement to expect x509 attributes. This can be worked around by crafting a custom RTE
 - The workaround requires informing VO's, not sustainable
- REST-based webmonitor
- Some replacement for ganglia
- Containerized version of the ARC-CE REST only server
- Intelligent Garbage collector for controldir leftovers
- Allow possibility to keep the *.errors file of jobs for sysadmin to investigate problems

ARC 7 – Obsolete components

- Server-side gridftp interface (both for job submission and as a storage element)
- ACIX
- Support for LRMS-ssh
- Python LRMS scripts
- Support for EMIES as a job interface
- Support for python2
- Use C++11
- LDAP to be marked as DEPRECATED
- Support for EGIS (in the client)
- Support for every interface in the ARC7 client except the REST

Smaller features, options to be removed from ARC7:

- Glue1 schema
- Fake site-bdii
- Argus support
- Nordugridmap-like gridmapfile generation

Fairshare refactor

Refactor fair-share config to better support ATLAS multicore / singlecore workloads

User Name	Effective Priority	Real Priority	Priority Factor
group_DTEAM_OPS.ops.ops011@gridpp.rl.ac.uk	0.50	0.50	1.00
group_DTEAM_OPS.ops.ops002@gridpp.rl.ac.uk	0.50	0.50	1.00
group_DTEAM_OPS.ops.ops044@gridpp.rl.ac.uk	0.50	0.50	1.00
group_HIGHPRIO.sum_tests.cmssgm@gridpp.rl.ac.uk	0.98	0.98	1.00
group_ALICE.alice.alice043@gridpp.rl.ac.uk	50000.00	0.50	100000.00
group_NONLHC.bio.bio017@gridpp.rl.ac.uk	50000.00	0.50	100000.00
group_NONLHC.lsst.tlsst008@gridpp.rl.ac.uk	50065.20	0.50	100000.00
group_LHCB.lhcb.lhcb052@gridpp.rl.ac.uk	50160.00	0.50	100000.00
group_ATLAS.atlas.atlassgm@gridpp.rl.ac.uk	50202.20	0.50	100000.00
group_NONLHC.dune.dune014@gridpp.rl.ac.uk	50978.30	0.51	100000.00
group_NONLHC.fermilab.fermi006@gridpp.rl.ac.uk	724843.50	724.84	1000.00
group_NONLHC.bio.bio045@gridpp.rl.ac.uk	2737132.25	27.37	100000.00
group_NONLHC.snoplus.snopluspilot010@gridpp.rl.ac.uk	2955285.50	29.55	100000.00
group_ATLAS.atlas_pilot_multicore.tatls015@gridpp.rl.ac.uk	6267528.50	62.68	100000.00
group_NONLHC.dune.dune004@gridpp.rl.ac.uk	10367403.00	103.67	100000.00
group_NONLHC.enmr.enmr022@gridpp.rl.ac.uk	14225772.00	142.26	100000.00
group_NONLHC.ligo.ligo020@gridpp.rl.ac.uk	31176220.00	311.76	100000.00
group_ALICE.alice.alicesgm@gridpp.rl.ac.uk	52123320.00	521.23	100000.00
group_ATLAS.prodatls.patls002@gridpp.rl.ac.uk	67385560.00	673.86	100000.00
group_ATLAS.atlas_pilot.tatls015@gridpp.rl.ac.uk	613508736.00	6135.09	100000.00
group_ATLAS.prodatls_multicore.patls002@gridpp.rl.ac.uk	712699520.00	7127.00	100000.00
group_CMS.cms_pilot_multicore.ttcms032@gridpp.rl.ac.uk	1.03849e+09	10384.85	100000.00
group_LHCB.lhcb_pilot.tlhcb006@gridpp.rl.ac.uk	1.58274e+09	15827.43	100000.00

Number of users: 23



Science and
Technology
Facilities Council

Scientific Computing

Questions?

Thomas.Birkett@stfc.ac.uk



Science and
Technology
Facilities Council

Scientific Computing