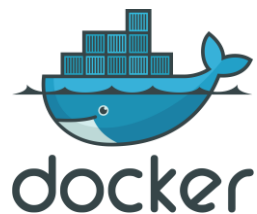# GridPP52

## *Tier1 Condor Batch Farm*

### *Tom Birkett*

STFC, Rutherford Appleton Laboratory

# Agenda

- Looking back, the past year
- Batch Farm overview
  - Fairshare logic
  - Condor defragmentation
  - Introduction of the 2023 worker nodes
- Job lifecycle
- Container image pipeline
- Condor 10 upgrade
- Docker wrapper sunsetting
- IPv6 and Docker
- ARC 7
- GPUs in Batch Farm



UK RI — Science and Technology Facilities Council

Scientific Computing

# The past year

- HTCondor upgrade
  - 9.0 -> 10.0
- Docker 27
- Decommission of site ARGUS instances
- Introduction of defragmentation of jobs on Condor pool
- IPv6 rollout across Condor components
- Optimisation of fair-share scheduling

UKRI

Science and
Technology
Facilities Council

Scientific Computing

# Batch Farm overview



- GridPP50 (Aug 2023):
  - 423 worker nodes.
  - ~ 49,000 logical CPU cores
  - HEPSPEC: ~ 630,000

- GridPP52 (Aug 2024):
  - 375 worker nodes.
  - ~ 72,000 logical CPU cores
  - HepScore: ~ 948,500
  - Job Slots: 66,296

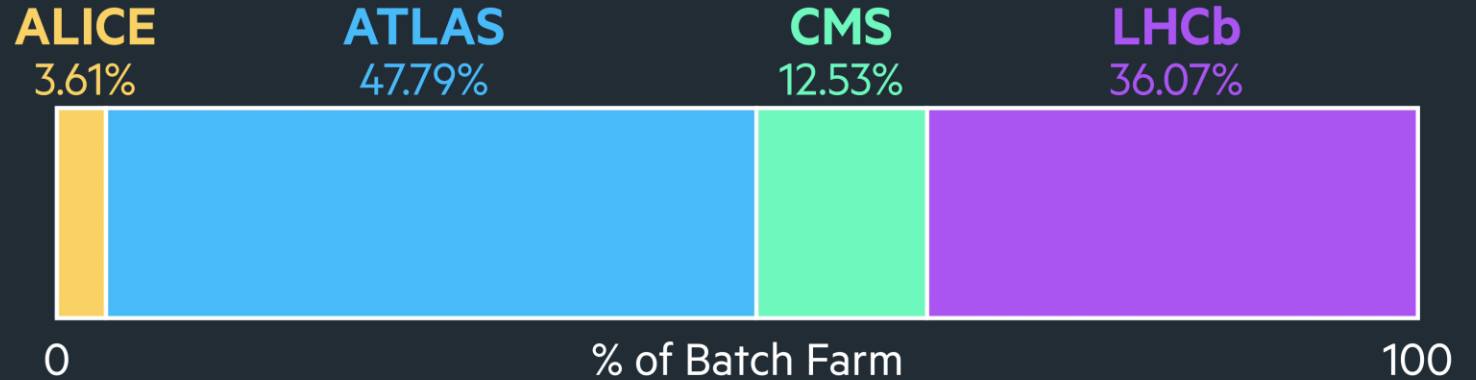UKRI Science and Technology Facilities Council

Scientific Computing

# Fairshare logic

RAL is a Tier-1 primarily for the four large LHC experiments - ATLAS, LHCb, CMS and ALICE. The graphic shows the fair-shares for the farm.

RAL also supports 20 other experiments
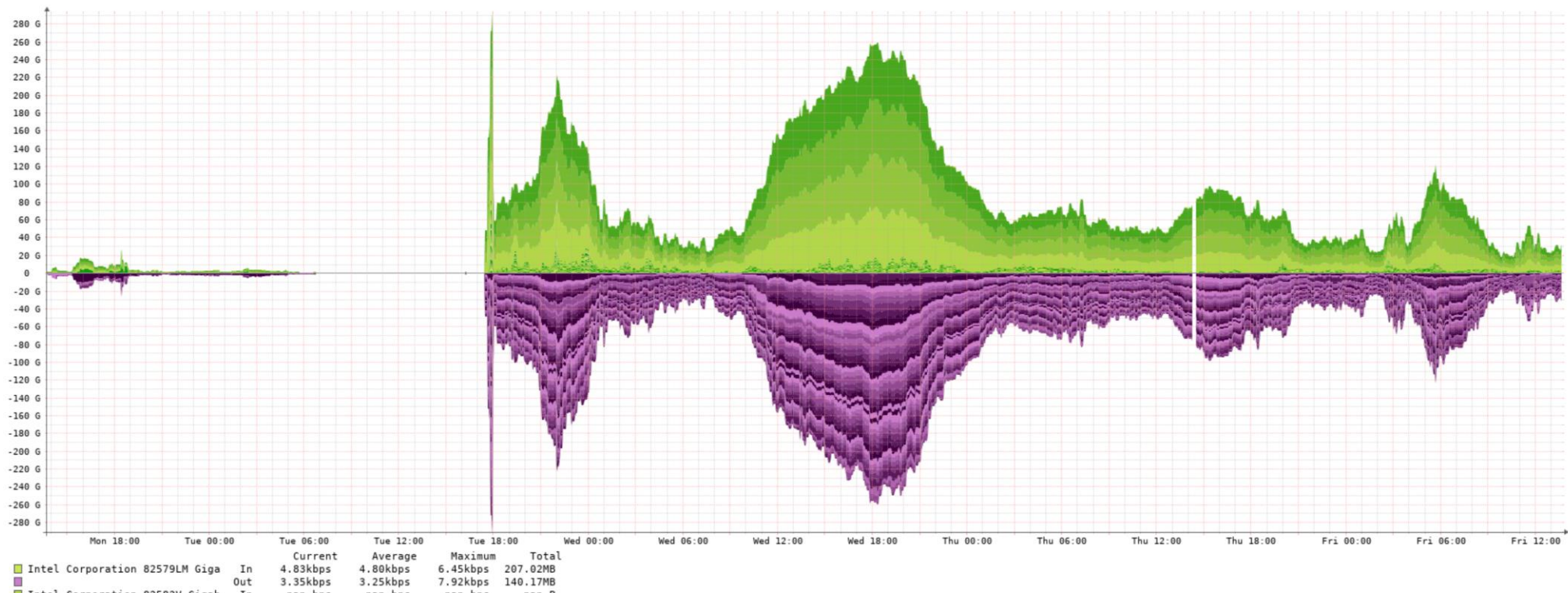
## Scenario 1
### Four LHC experiments

| ALICE | ATLAS | CMS | LHCb |
|-------|-------|-----|------|
| 3.61% | 47.79% | 12.53% | 36.07% |

0     % of Batch Farm     100

# Condor Defragmentation

DAEMON_LIST=$(DAEMON_LIST), DEFRAG

DEFRAG_CANCEL_REQUIREMENTS=(Cpus >= 16)

DEFRAG_DRAINING_MACHINES_PER_HOUR=30

DEFRAG_INTERVAL=600

DEFRAG_MAX_CONCURRENT_DRAINING=40

DEFRAG_MAX_WHOLE_MACHINES=-1

DEFRAG_REQUIREMENTS=(RalCluster == "wn-2022-lenovo" || RalCluster == "wn-2023-lenovo") && (StartJobs =?= true) && (PartitionableSlot =?= true) && (Cpus == 0)

DEFRAG_SCHEDULE=graceful
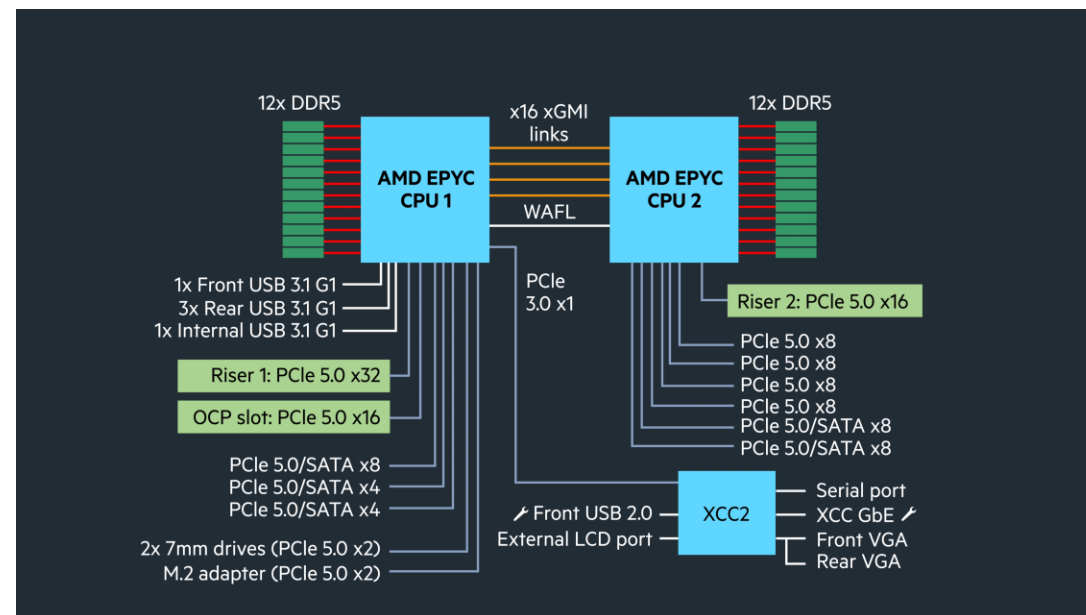
DEFRAG_UPDATE_INTERVAL=300

# Introduction of 2023 worker nodes

- Partial activation of 2023 worker nodes on Tuesday, the 20th.

- Full activation of 2023 worker nodes on Wednesday, the 21st.

- Peak:
  - 257.8 Gbps
  - 257.7 Gbps
  - Total aggregate: 515.58 Gbps



| | Current | Average | Maximum | Total |
|---|---|---|---|---|
| Intel Corporation 82579LM Giga In | 4.83kbps | 4.80kbps | 6.45kbps | 207.02MB |
| Out | 3.35kbps | 3.25kbps | 7.92kbps | 140.17MB |

**UKRI**
**Science and Technology Facilities Council**

Scientific Computing
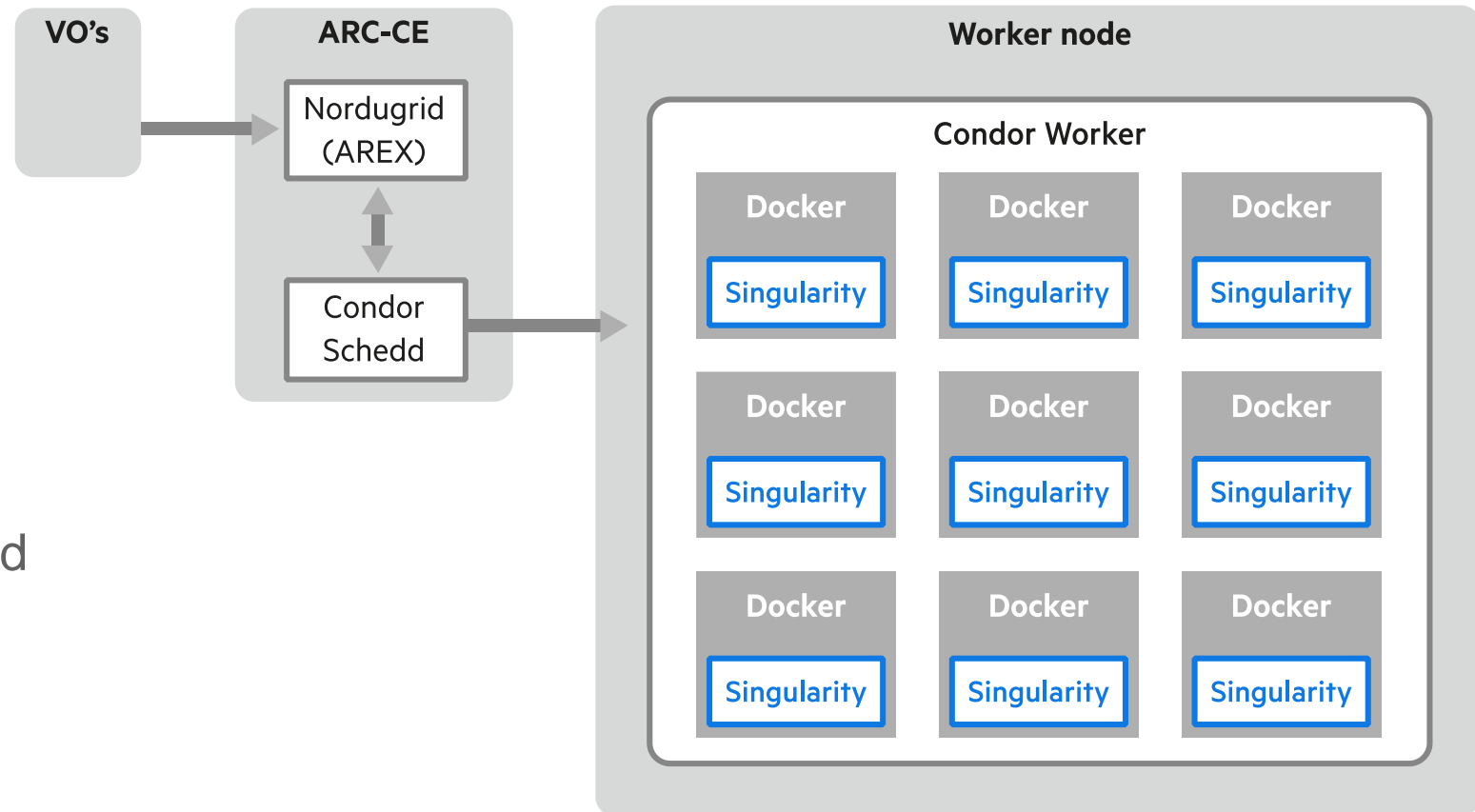
# Lenovo 2023's (SR645 V3)

- Look into better ways to use our current worker nodes to help us reach Net Zero goals

- Lenovo SR645 V3 specs:
  - CPU: 2 x 96 core (192 threads)
  - RAM: 1536GB
  - Disk: 2 x 5961GiB NVMe (11.6TiB total)

UKRI
Science and Technology Facilities Council

Scientific Computing

# 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

- RAL has two job queues:
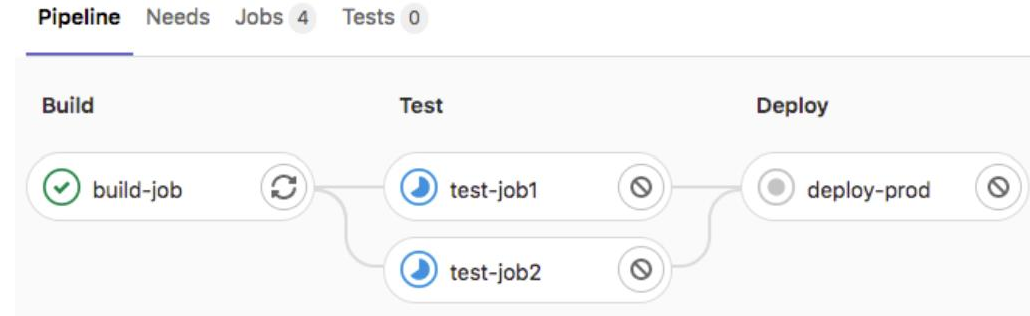  - `EL8`: Rocky 8
  - `EL9`: Rocky 9

# Container image pipeline (Automation)

Tier1 Grid Services graduate, Earl Talavera is working on automating container image builds used on the Batch Farm.

The Gitlab CI/CD Project involves a fork of the GitHub repository of grid-worker nodes and utilises GitLab's CI/CD pipeline tools to automate the build of the grid-worker nodes docker images and pushes it into Harbor authenticating via a Robot Account.

It is set on a monthly schedule executing on the first day of the month, with the builds occurring inside an OpenStack VM where the building of images is set as jobs on multiple GitLab runners concurrently executing.

Pipeline   Needs   Jobs 4   Tests 0

Build          Test              Deploy

build-job      test-job1         deploy-prod
               test-job2

GitLab   openstack.   HARBOR

Science and Technology Facilities Council

Scientific Computing

https://github.com/stfc/grid-workernode

# Condor 10 upgrade

- RAL upgraded to Condor 10 - October 2023

- Highlight features:
  - GSI Authentication method has been removed (X.509 proxies are still handled by HTCondor) (HTCONDOR-697)
    - From Condor 10, you must specify `USE_VOMS_ATTRIBUTES=True` to continue using x509 ClassAds on both Startd and Schedd
  - Add the ability to select a particular model of GPU when the execution points have heterogeneous GPU cards installed or cards that support nVidia MIG (HTCONDOR-953)
  - For IDTOKENS, signing key is not required on every execution point (HTCONDOR-638)

UK RI
Science and Technology Facilities Council

Scientific Computing

# Docker Wrapper sunsetting

- The `docker.py` wrapper used at RAL, builds the docker run command that is executed by Condor. https://github.com/stfc/ral-htcondor-tools/blob/master/docker.py

- 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. Implementation methods vary, RAL configures IPv6 in the following way:

- `/etc/docker/daemon.json`:

```json
{

 "bip": "10.20.0.1/20",

 "experimental": true, (Not required since Docker 27)

 "ip6tables": true, (Default since Docker 27)

 "live-restore": true,

 "log-driver": "local",

 "storage-driver": "overlay2"

}
```

Science and
Technology
Facilities Council

Scientific Computing

# IPv6 Docker networking continued

▪ Setup Docker network with IPv6 config:

docker network create \
--driver=bridge \
--subnet=172.28.0.0/16 \
--ip-range=172.28.4.0/23 \
--gateway=172.28.5.254 \
--ipv6 \
--subnet=**2001:db8:54:3:82f6:dd83:1000::/112** \
--ip-range=**2001:db8:54:3:82f6:dd83:1000:8000/113** \
--gateway=**2001:db8:54:3:82f6:dd83:1000:1** \
ralworker

▪ Host IPv6 network config:

IPv4:
inet xxx.xxx.221.131/26

IPv6:
inet6 **2001:db8:54:3:82f6:dd83::/64**
inet6 fe80::63f:72ff:fed4:4688/64

Science and
Technology
Facilities Council

Scientific Computing

# 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

- Containerised 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

- Release candidate 1 now available!
  http://www.nordugrid.org/arc/arc7/common/repos/testing-repository.html

UK RI
Science and
Technology
Facilities Council

Scientific Computing

# GPUs in Batch Farm

- Day 1:
  - Configure Condor to be GPU aware [https://htcondor-wiki.cs.wisc.edu/index.cgi/wiki?p=HowToManageGpus](https://htcondor-wiki.cs.wisc.edu/index.cgi/wiki?p=HowToManageGpus) (HTCondor 10+ supports GPU awareness)
  - Create new ARC queue for GPU jobs
- Day 2:
  - Integrate GPU's in STFC cloud with Coyote: [https://github.com/stfc/coyote_scripts](https://github.com/stfc/coyote_scripts)
  - Work with community to understand accounting for GPU work

UKRI
Science and Technology Facilities Council

Scientific Computing

# Future work / plans

- HTCondor upgrade
  - 10.0 -> 23.0
- ARC 7 (Tokens)
- Removal of Docker wrapper.
- GPU's in Batch Farm
- Production environment for running jobs in STFC Cloud

# Questions?

[Thomas.Birkett@stfc.ac.uk](mailto:Thomas.Birkett@stfc.ac.uk)