

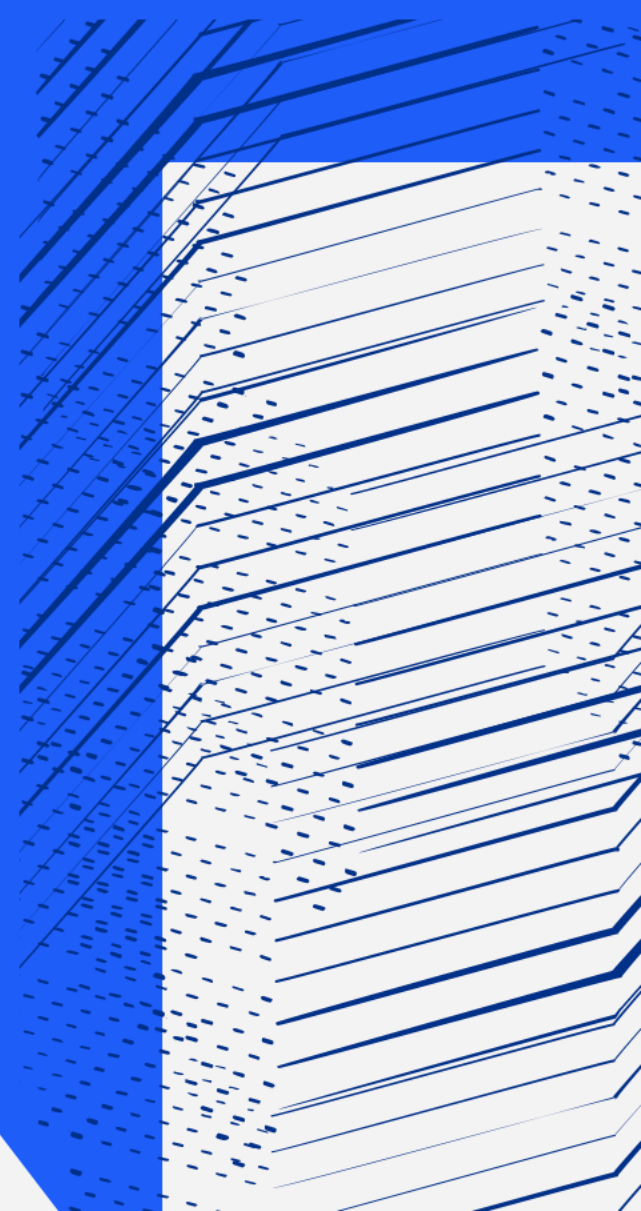


Science and
Technology
Facilities Council

Scientific Computing

Deploying and Running Ceph Clusters for Analysis Facilities

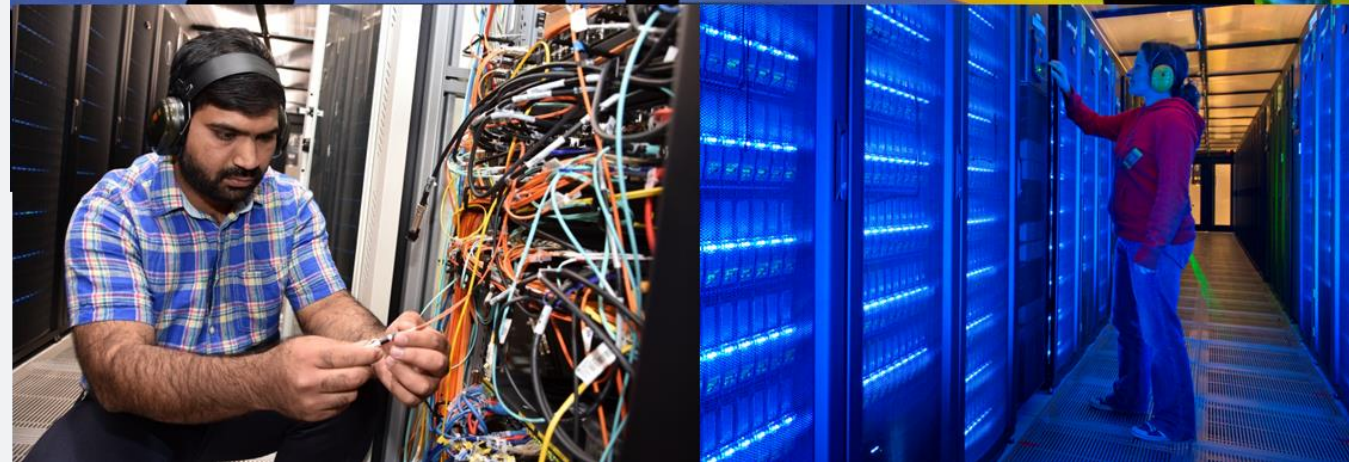
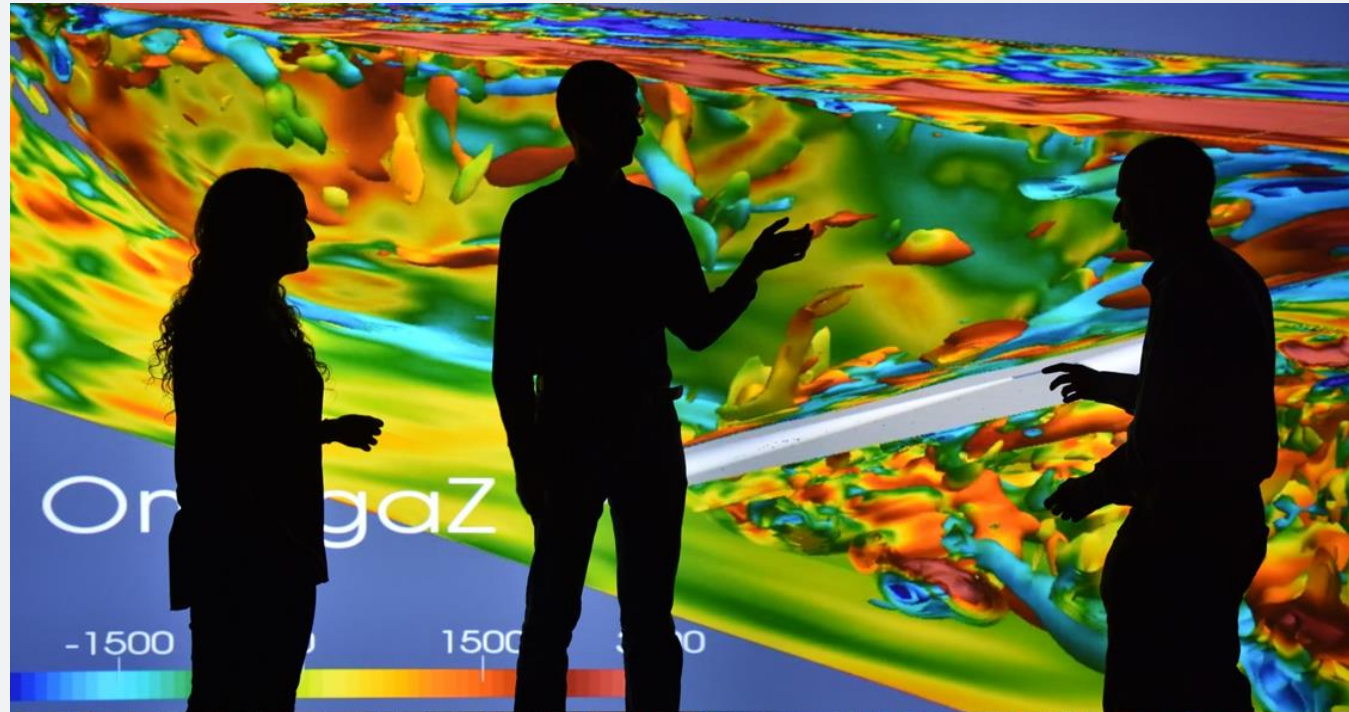
Rob Appleyard



Us and our users

STFC Scientific Computing Department (SCD)

- Part of UK Research and Innovation
- We support advanced scientific facilities, including...
 - ISIS neutron spallation source
 - Central Laser Facility
 - Rosalind Franklin Institute
 - Diamond Light Source*
 - The UK's WLCG Tier 1*



RAL Facilities – ISIS neutron source

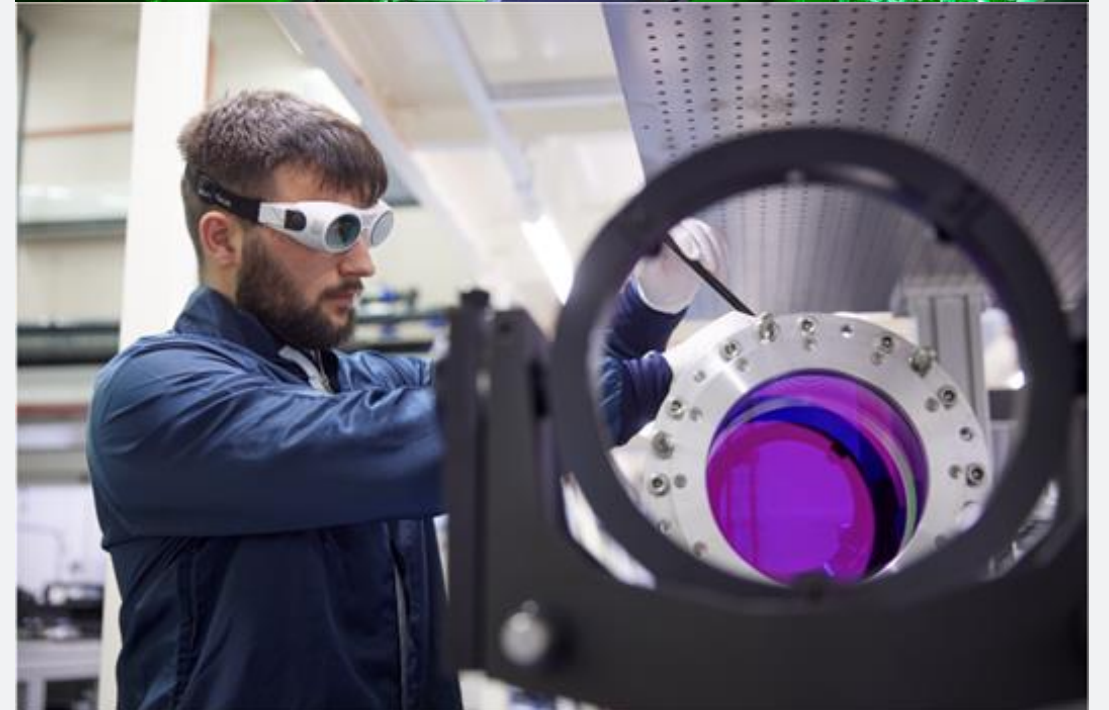
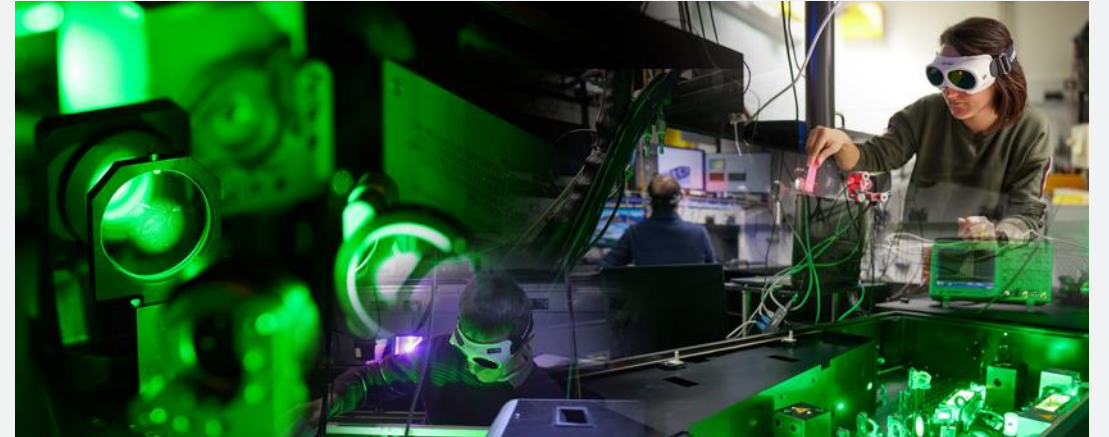
- Neutron and muon source running since 1985
- Diverse set of instruments optimised for specific analysis tasks
- Scientists bid for beamtime to conduct experiments
 - 24/7 operation during cycles
- Open for public domain research and industrial users



Image from the ISIS website: <https://www.isis.stfc.ac.uk>

RAL Facilities – Central Laser Facility (CLF)

- Specialist laser science facility with multiple instruments
 - Condensed matter
 - Life sciences
 - Plasma physics
 - Spectroscopy
- Vulcan PW laser - 10^{15}W in 10^{-12}s pulses
 - Fusion energy research
 - “Laboratory astrophysics”



RAL Facilities – Rosalind Franklin Institute (RFI)

- Technology development for health research
 - Microscopes
 - Machine learning
 - Protein libraries
 - Dynamic Imaging
- Common pattern of advances in detector technology leading to explosion in data volumes

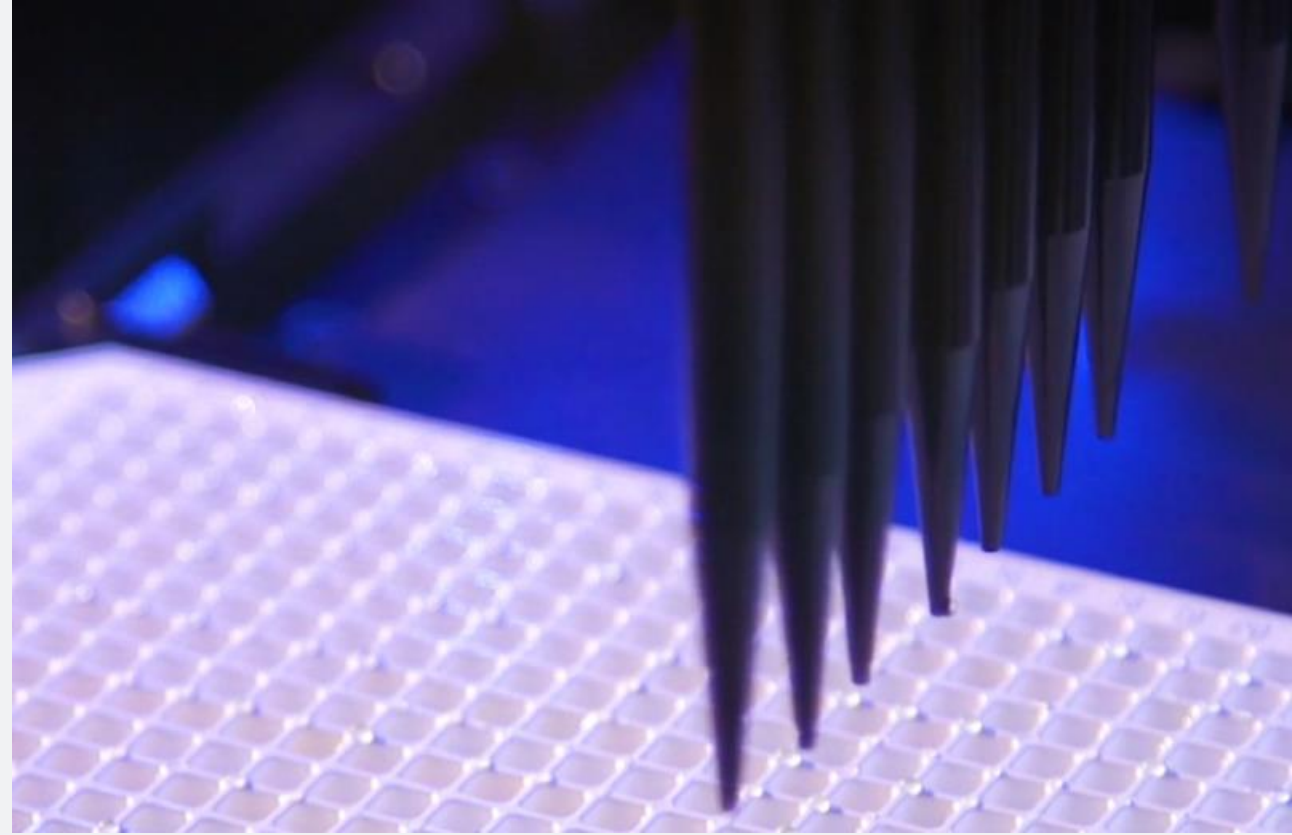


Image from the RFI website: <https://www.rfi.ac.uk/>

Experimental Use cases

- User-facing experimental facilities at RAL – ISIS, Diamond, CLF
 - Users turn up for a tightly-constrained period of experimental time and wish to guide their work with live data analysis.
 - Analysis tasks are user-defined
- Users who want to make large datasets securely available to external institutes.
- Organisational private cloud hosts a large, diverse collection of VMs
 - Needs high-performance data storage for VM images and working data.
- WLCG Tier 1 requires very high capacity, high-throughput and low cost data access to local batch, large amounts of inter-site I/O, and an interface to the tape system.

Our Ceph Services

Sirius

250TiB triple-replicated SSD RBD block storage

Arided

400TiB triple-replicated SSD CephFS

Deneb

5.5PiB erasure-coded HDD CephFS

Echo

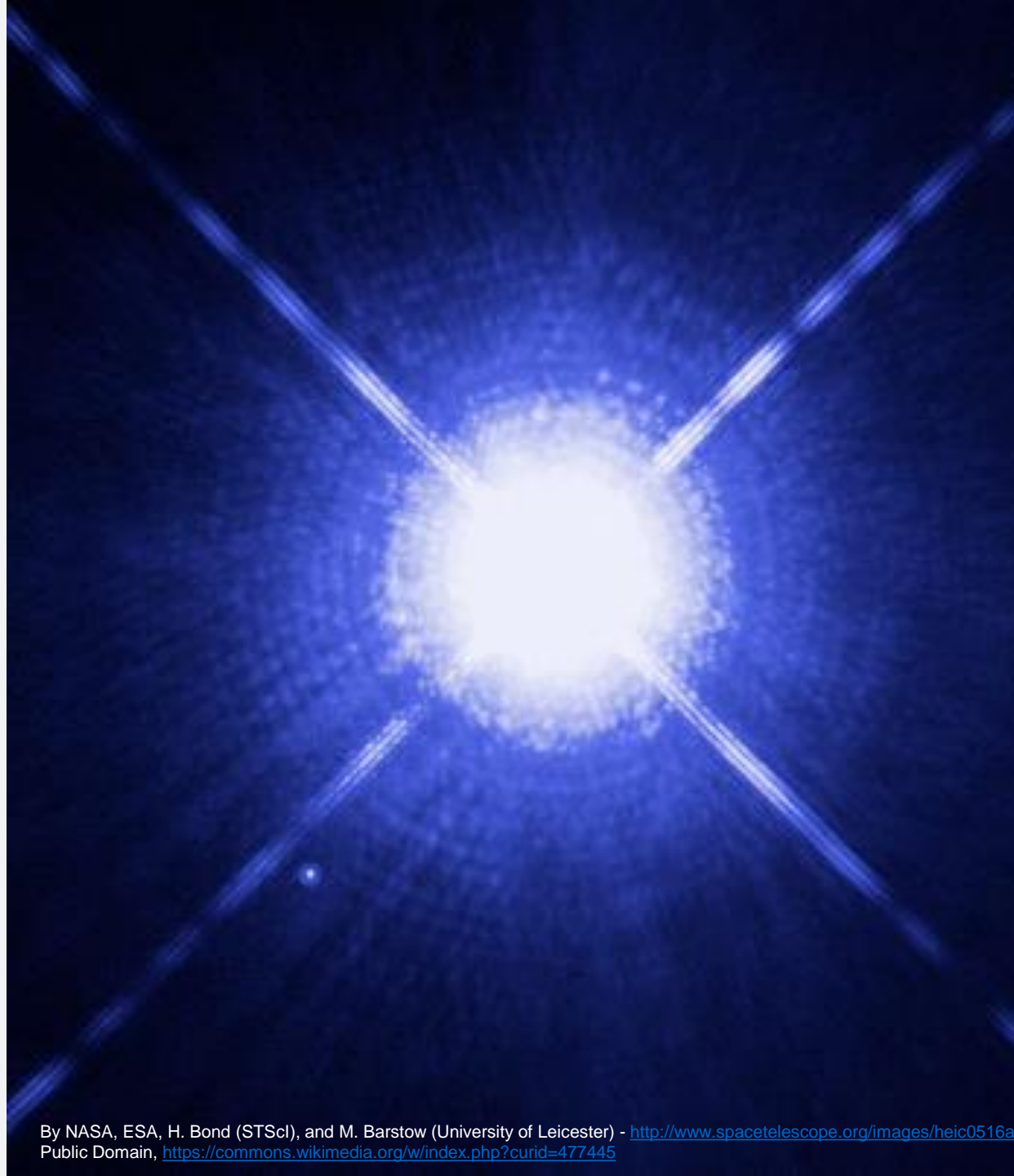
57PiB erasure-coded HDD object storage presenting XrootD, GridFTP, S3 and SWIFT interfaces.



ceph

Sirius

- Underpinning infrastructure for private cloud
- The SCD Cloud is a general-purpose site resource, plus some external users
 - Utility VMs
 - CPU/GPU compute
- 250 TiB usable with 3*replication (750TiB raw)
 - Pure NVMe storage
 - Originally specced on HDDs
 - ...but this ran out of IOPS.
- RADOS Block Device (RBD) access only
 - High performance scratch/cache space for cloud.
- Typical storage node spec:
 - 8*4TB read-intensive NVMe SSD
 - 32-core AMD EPYC 7502P
 - 128GB RAM
 - 25Gb Networking



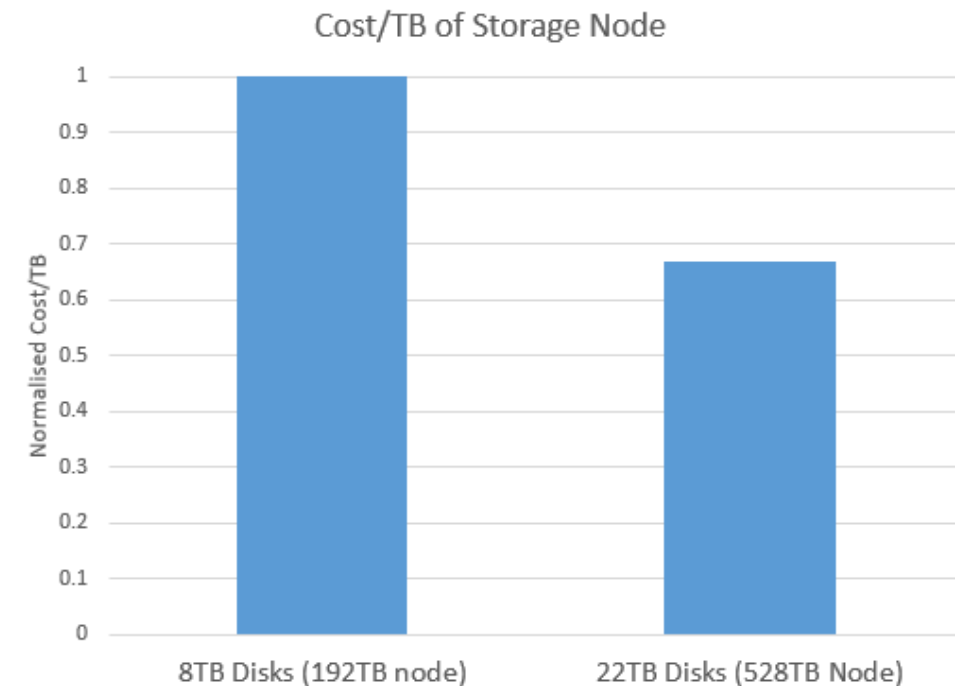
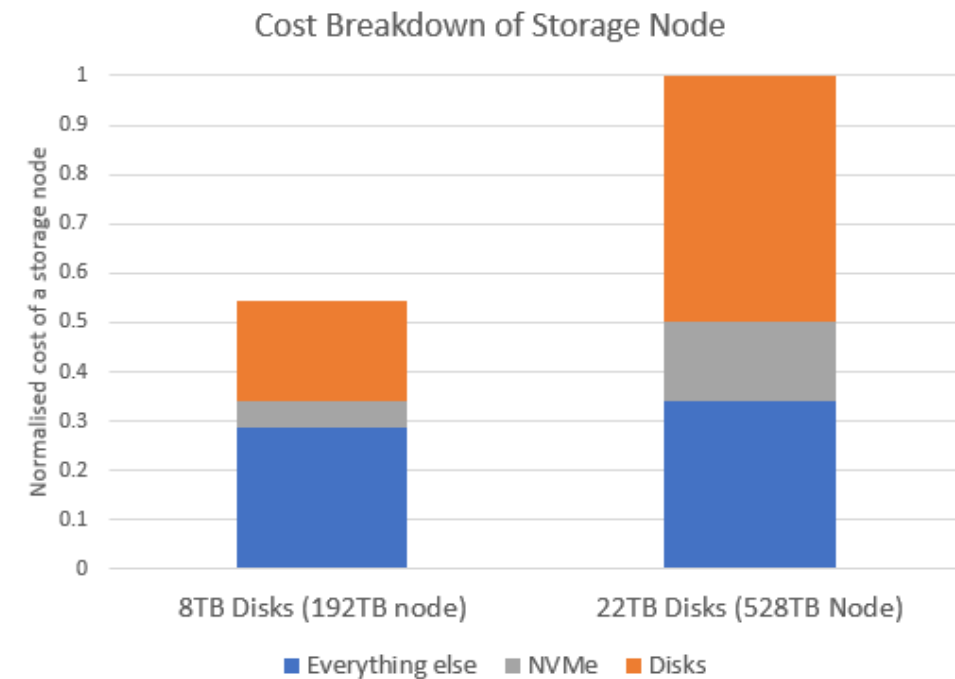
Deneb

- Genesis of project:
 - Multiple user requests for large, sharable FS for experimental data cache
 - STFC user with a self-managed and EOL CephFS cluster that needed a replacement
- CephFS advantages
 - Scalable, highly-available, mountable POSIX file system
 - Scientists can collaborate on a shared file system
- 5.5PB usable with 8+3 erasure coding (7.5PiB raw).
- Throughput is very low (~200MB/s on a 60-node cluster), but latency is noticed by users
- Shared between multiple large user communities
 - ISIS
 - Rosalind Franklin Institute
 - RAL Central Laser Facility



Deneb hardware

- Typical storage node spec:
 - 20*8TB SATA HDDs
 - 3*1.6TB NVMe for RocksDBs (3%)
 - 2*16C/32T Intel Xeon Silver 4216 @ 2.1GHz
 - 128GB RAM
 - 25Gb Networking
- Open question – this is probably not the optimum hardware
 - 8TB HDD is not an obvious storage medium to buy in bulk
 - Bigger disks are *much* more cost effective
 - SSDs are *much* faster
- Plan for this year's procurement...
 - Buy mostly our standard hardware...
 - ...plus one test node with 22TB disks + 4% NVMe
 - This host offers >30% better cost/TB vs 8TB disks



Arided

- 2019:
 - Sirius is low on space - users are placing bulk data in VM images...
 - Sirius is very expensive, so let's make something better for that use case
- 400TiB usable with 3*replication, on SATA SSDs (1.2PiB raw)
- Manila - Shared File Systems as a Service
 - Users can self-provision native CephFS shares and mount them on their VMs
- Why not use Deneb?
 - Cloud supports extremely diverse use cases
 - 3 * replication is most flexible approach
 - Small files imply overhead with EC
 - Most consistent performance
- Typical storage node spec:
 - 24*4TB SATA SSDs
 - 2*24-core AMD EPYC 7413
 - 256GB RAM
 - 25Gb Networking



Manila from a user's perspective

STFC Cloud Terms SLA FAQs stfc admin vwa13372

Project / Share / Shares

Shares

Filter + CREATE SHARE DELETE SHARES

Displaying 57 items

<input type="checkbox"/>	Name	Description	Metadata	Size	Status	Protocol	Visibility	Share Network	Share Group	Actions
<input type="checkbox"/>	cephshare309	-		1GiB	Available	CEPHFS	private	-	-	EDIT SHARE
<input type="checkbox"/>	cephshare309	-		1GiB	Available	CEPHFS	private	-	-	EDIT SHARE
<input type="checkbox"/>	cephshare308	-		1GiB	Available	CEPHFS	private	-	-	EDIT SHARE
<input type="checkbox"/>	cephshare307	-		1GiB	Available	CEPHFS	private	-	-	EDIT SHARE
<input type="checkbox"/>	cephshare306	-		1GiB	Available	CEPHFS	private	-	-	EDIT SHARE
<input type="checkbox"/>	cephshare305	-		1GiB	Available	CEPHFS	private	-	-	EDIT SHARE
<input type="checkbox"/>	cephshare304	-		1GiB	Available	CEPHFS	private	-	-	EDIT SHARE
<input type="checkbox"/>	cephshare303	-		1GiB	Error	CEPHFS	private	-	-	MANAGE RULES

Create Share

Share Name: tom-test-share

Description: [Empty]

Share Protocol: NFS

Size (GiB): 10

Share Type: cephfsnfstype

Availability Zone: [Empty]

Share Group: [Empty]

Metadata: [Empty]

Make visible for all

Description: Select parameters of share you want to create.

Metadata: One line - one action. Empty strings will be ignored. To add metadata use: key=value

Share Limits

Total Gibibytes: 964 of inf GiB Used

Number of Shares: 57 of inf Used

CANCEL CREATE



Access management

Share Overview

Name	cephshare309
ID	06825d6e-062b-4aaf-a773-b19b0f63e0db
Status	Available
Export locations	Path: 130.246.223.45:6789,130.246.223.46:6789,130.246.223.47:6789:/volumes/_r
	Preferred: False
	Is admin only: False
	Share Replica ID: 9e593cba-fc47-4e6a-aec4-4b74530f1e87
Visibility	private
Availability zone	nova
Size	1 GiB
Protocol	CEPHFS
Share type	Name: cephfstype
	ID: 4013e155-e731-468d-aaa9-141b065961e8

Share Rules: cephshare309

+ ADD RULE

DELETE RULES

Displaying 2 items

<input type="checkbox"/>	Access Type	Access to	Access Level	Status	Access Key	Actions
<input type="checkbox"/>	cephx	cephx-auth-id-test	rw	active	AQCfAXRcD0fX0hAAII9DLZdAoEjWnR0z037SDw==	DELETE RULE
<input type="checkbox"/>	cephx	cephx-auth-id-test-ro	ro	active	AQCfAnRc/xfBMhAAINyZHP8d3IZvX7Wy1mZyg==	DELETE RULE

Mounting Shares

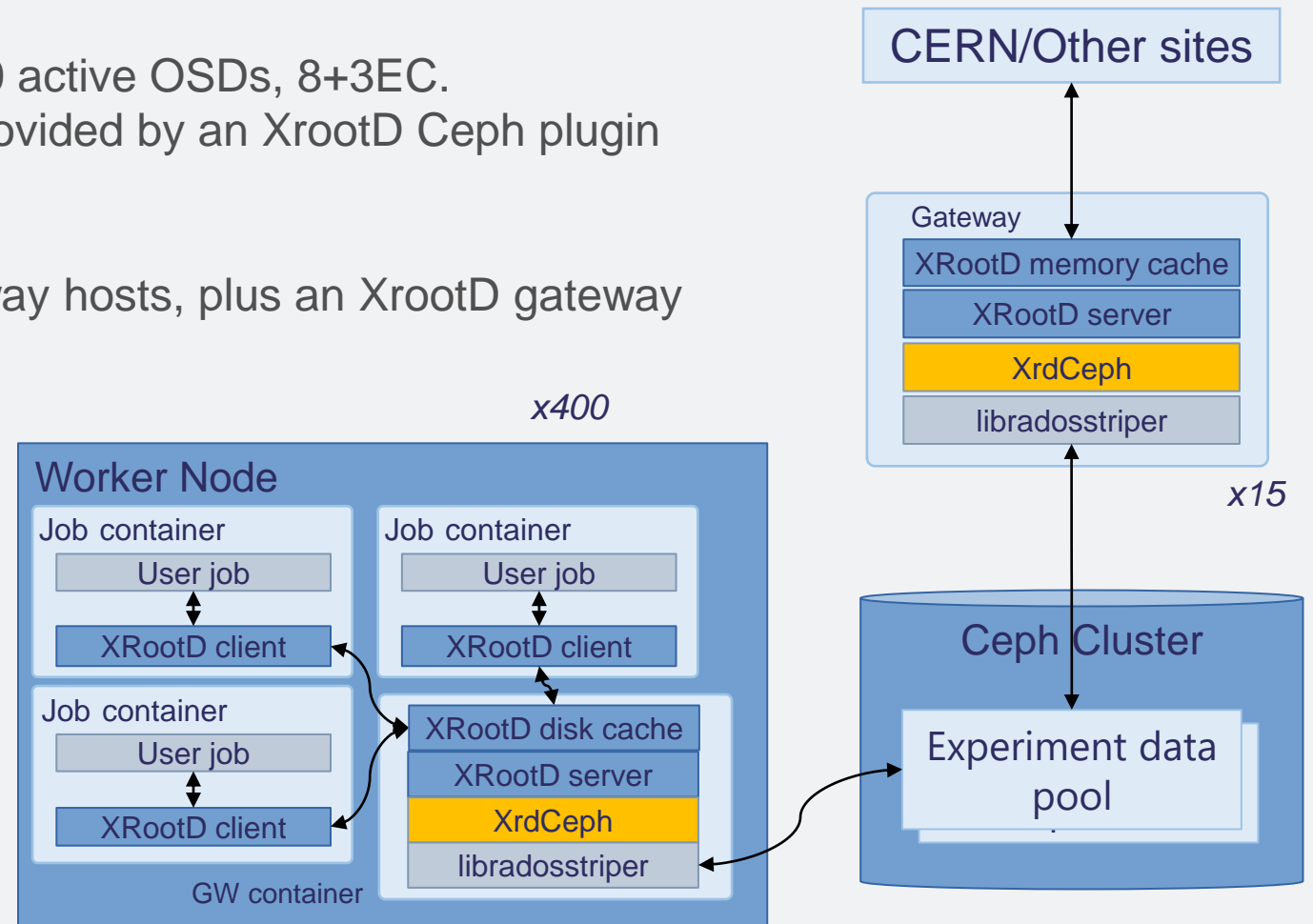
CephFS:

```
[root@host ~]# mount -t ceph [monitor IPs]:/volumes/_nogroup/c712ad9e-45a5-4cd9-8e9b-e28882ea4ba1 /root/mountdir  
> -o name=tomtestuser1,secret=*****==  
[root@host ~]# df -h  
Filesystem  
...  
[monitor IPs]:/volumes/_nogroup/c712ad9e-45a5-4cd9-8e9b-e28882ea4ba1 54T 37G 54T 1% /root/mountdir  
[root@host ~]#
```


Echo for the WLCG

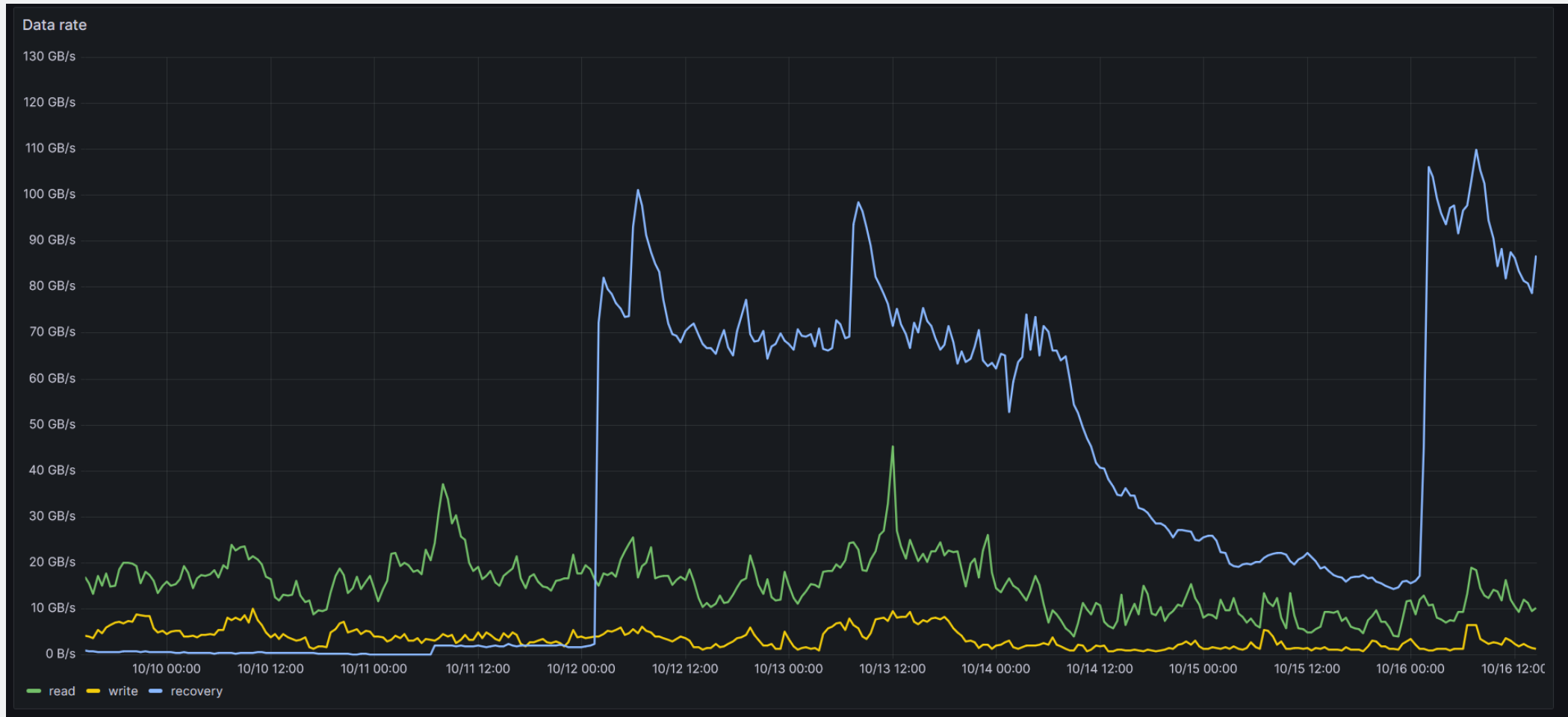


- First production Ceph cluster at RAL
- 57PiB usable (78PiB raw), 300 hosts, 6600 active OSDs, 8+3EC.
- Low-level (RADOS object store) access provided by an XrootD Ceph plugin developed in-house
 - All files are stored as 64MB stripes
- Dedicated cluster of external XrootD gateway hosts, plus an XrootD gateway process on every worker node.
- Typical storage node spec (last year):
 - 24*18TiB HDD
 - 2*16C/32T Xeon Silver 4216@2.1GHz
 - 192GB RAM
 - 25Gb Networking
- Typical gateway host spec:
 - 2*16C/32T Xeon Silver 4214R@3.5GHz
 - 192GiB RAM
 - 25Gb Networking



An architecture diagram of how Echo is accessed, from [Tom Byrne's presentation at Cephalocon 2023](#)

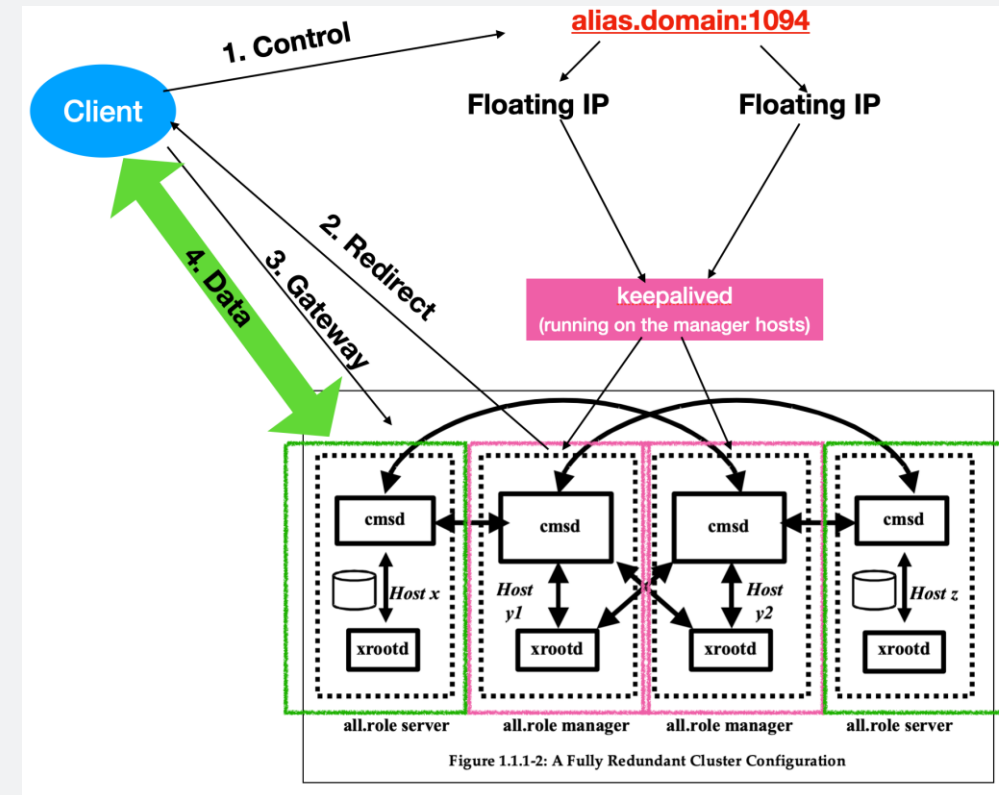
Echo internal throughput



Echo-specific Developments



- 2022 – XrootD external gateways were a problem
 - Not enough of them
 - Primitive load balancing (round-robin DNS) meant overloaded/broken hosts became black holes
 - Internally developed Ceph XrootD plugin needed to catch up
- Now – much less so
 - Just deployed 10 new gateway hosts
 - Implemented a cmsd redirector for the cluster
 - Lots of work on optimising XrootD for Ceph
- Planned in 2024:
 - Move cluster's failure domain to rack level
 - Resilience
 - Ease of maintenance
 - SSD-only pools

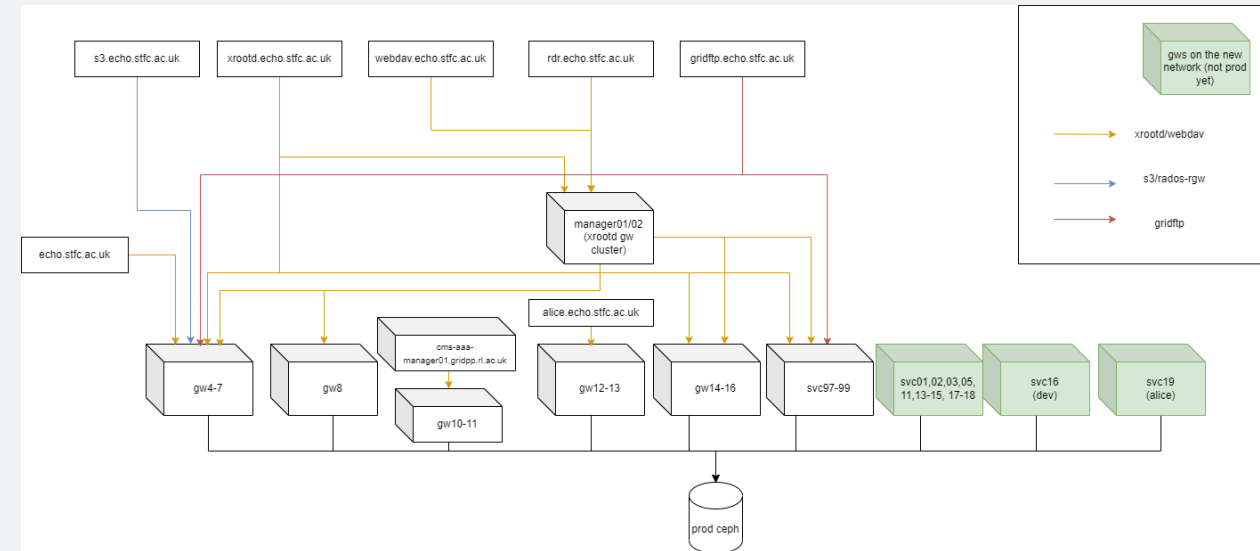


A diagram of our cmsd implementation (from Jyothish Thomas's GridPP talk* and the XrootD documentation)

Echo – More Gateways



- Implemented XrootD load balancing – great!
 - Much higher average utilisation, more bandwidth, better load balance
- New problem!
 - WN gateways are read-only
 - Job results are written using ‘external’ gateways
 - Inter-site transfers can crowd this out
 - So have some dedicated gateways for this...
- Complexity creep...
 - Need a ground-up rethink of what we want these systems to do and how we should manage them

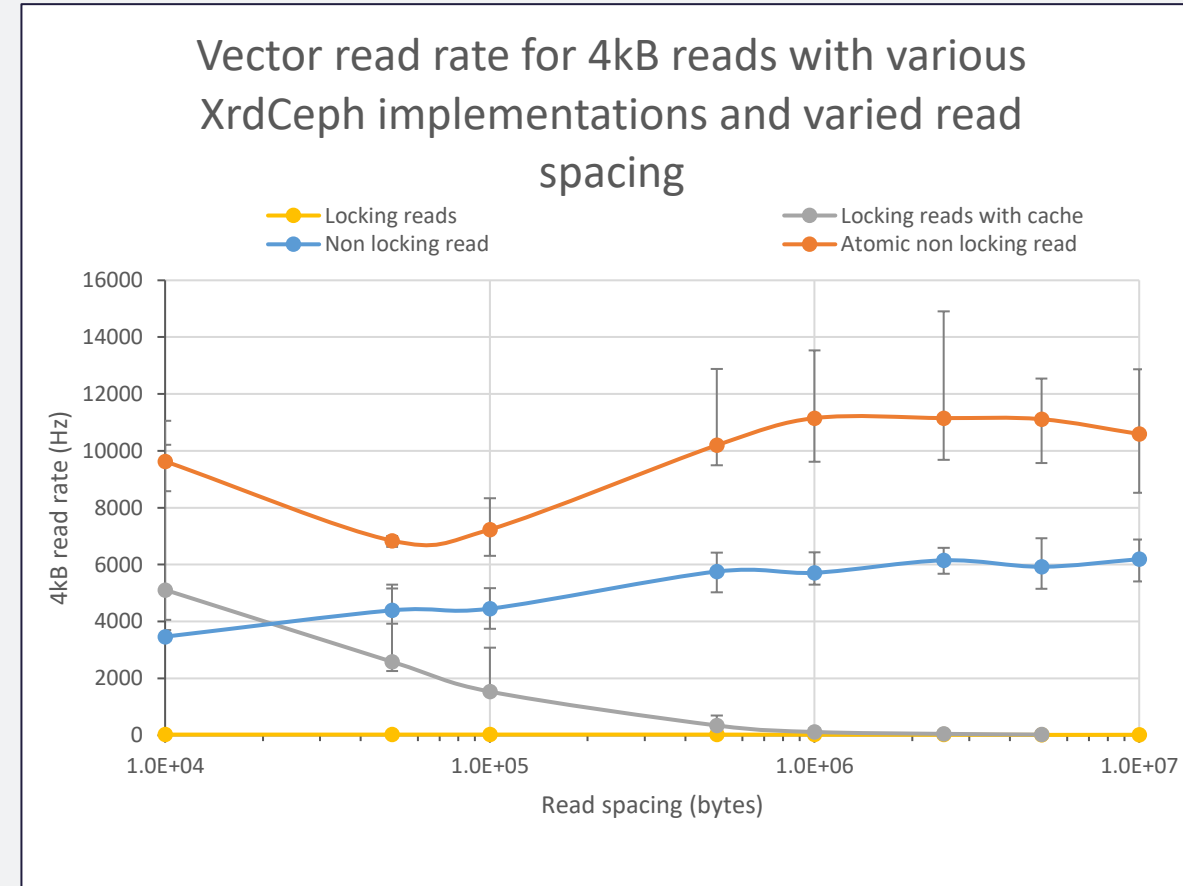


An internal diagram of what gateway host does what.

Echo XrootD Optimisations



- Specific user jobs were very problematic
 - Small vector reads (a job requesting many, very small, pieces of a large file)
- Problems:
 - Excessive caching led to massive read amplification
 - Lots of back-and-forth to single-threaded OSD processes before a read started (locks)
- Major software upgrades (credit to Jyothish Thomas and our XrootD development team, see talks linked below)
 - We can exploit the fact that WLCG Echo files are immutable - locks are redundant - remove them!
 - Clients can choose preferred behaviour – copy-to-scratch or vector read
 - Switch to vector reads was initially problematic – overloaded OSD transaction rate capacity
 - Needed to re-add a buffering layer
 - Now working well in production



The impact of our Vector read change – from [Tom Byrne's Cephalocon talk](#)

Echo for non-WLCG users

- Echo provides an S3 endpoint
 - Low cost storage on a widely-used interface
 - Popular with internal users - RFI
 - Easy to manage
- But is Echo the right way to supply this?
 - Echo is focussed on the WLCG
 - Vast majority of volume/throughput
 - Procurement, management style, access patterns
 - Huge cluster + bespoke XrootD access software is already complex
- How can we ensure that we provide an appropriate level of service with all the features that the users want?



Amazon S3 Logo from [Wikipedia](#)

Tier-1 Storage Plans – Long Term

HDD

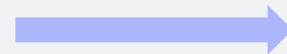
As HDDs become more archival, why not just use tape?



We expect to have production SSD endpoints by the end of GridPP7 (~2027).

Focusing Swift-HEP work on QoS specifically towards SSD.

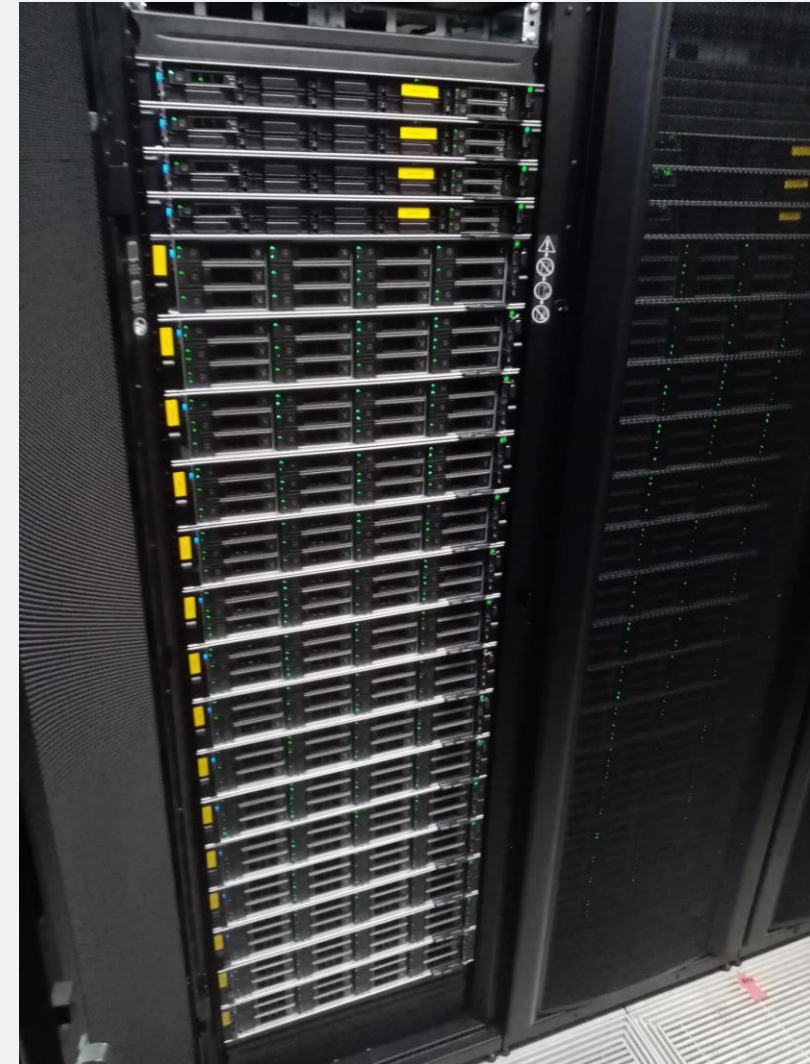
Tape



SSD

Future Plans – management

- Standardise on a modern OS and a modern version of Ceph
 - First Rocky 8, then Ceph Pacific
- Ceph clusters are proliferating, but our management approach is dated
- RAL site has a full config management system with automated package management, but our interactions with Ceph are largely manual
 - Except Arided, which uses cephadm – proof of concept
- Project to switch all clusters to proper orchestration tooling
 - Probably cephadm
 - Echo is the tricky one due to scale
- Deprecate Echo's dedicated network for internal rebalancing.
 - Extra network interface adds a lot of complication, doesn't solve bottleneck



Conclusions

- RAL talks a lot about Echo
 - But have standardised on Ceph for fulfilling many other computing requirements
 - General effort to standardise solutions provided to site users.
- Future
 - Expand use cases, support new users
 - Clusters should (mostly) be organised by use case rather than user community

The UKRI logo consists of the letters 'UK' stacked above 'RI' in a white, bold, sans-serif font, set against a dark blue square background.

Science and
Technology
Facilities Council

Scientific Computing

Questions?

RAL Facilities – Diamond Light Source (DLS)

- X-ray synchrotron running since 2007
- 32 specialised beamlines placed at tangents to the beamline
 - Diffraction, spectroscopy
- Tape archival also in SCD
- 10% of beamtime is available for commercial users, the rest is for public domain science



Image from the Diamond website: <https://www.diamond.ac.uk/Home/About.html>

Abstract

- Deploying and Running Ceph Clusters for Analysis Facilities
- The RAL Scientific Computing Department provides support for several large experimental facilities. These include, among others, the ISIS neutron spallation source, the Diamond X-Ray Synchrotron, the Rosalind Franklin Institute, and the RAL Central Laser Facility. We use several Ceph storage clusters to support the diverse requirements of these users.
- These include Deneb, a petabyte-scale CephFS cluster, Sirius, a pure-NVMe cluster used to provide the underlying storage for STFC's private cloud, our WLCG-focussed Echo cluster which also provides S3 and SWIFT access, and Arided, a new SSD cluster providing mountable CephFS storage to our private cloud. While all of these services use Ceph to provision the storage, each has a different architecture and usage profile. In particular, Arided has been deployed with the cephadm cluster management system, a first at RAL.
- This paper will provide an outline of these services, their development and deployment, how they are used, their hardware requirements and loadings, our experiences of supporting them as production services. We will discuss the expected development roadmaps for these services for the remainder of 2023 and going into 2024, and also provide an update on recent changes to the Echo service and its XrootD interface.