# Status of Ceph Storage at RAL

James Adams, Bruno Canning, Shaun de Witt, Alastair Dewhurst,, George Vasilakakos

# Introduction

- For some time RAL has been working on Ceph as a replacement to disk only Castor storage.

  - We would like to retire Castor for disk only storage by 2017.

  - Ceph storage available this year is all beyond pledge, aka 'free' for any VOs that wish to test.

- Aim to provide thinnest layer possible on top of Ceph:

  - Xrootd, GridFTP and [in future] http protocols are required for LHC VOs to work.

  - S3/Swift API provided (and encouraged) for other VOs.

- Need to use Erasure Coding to keep costs inline with Castor.
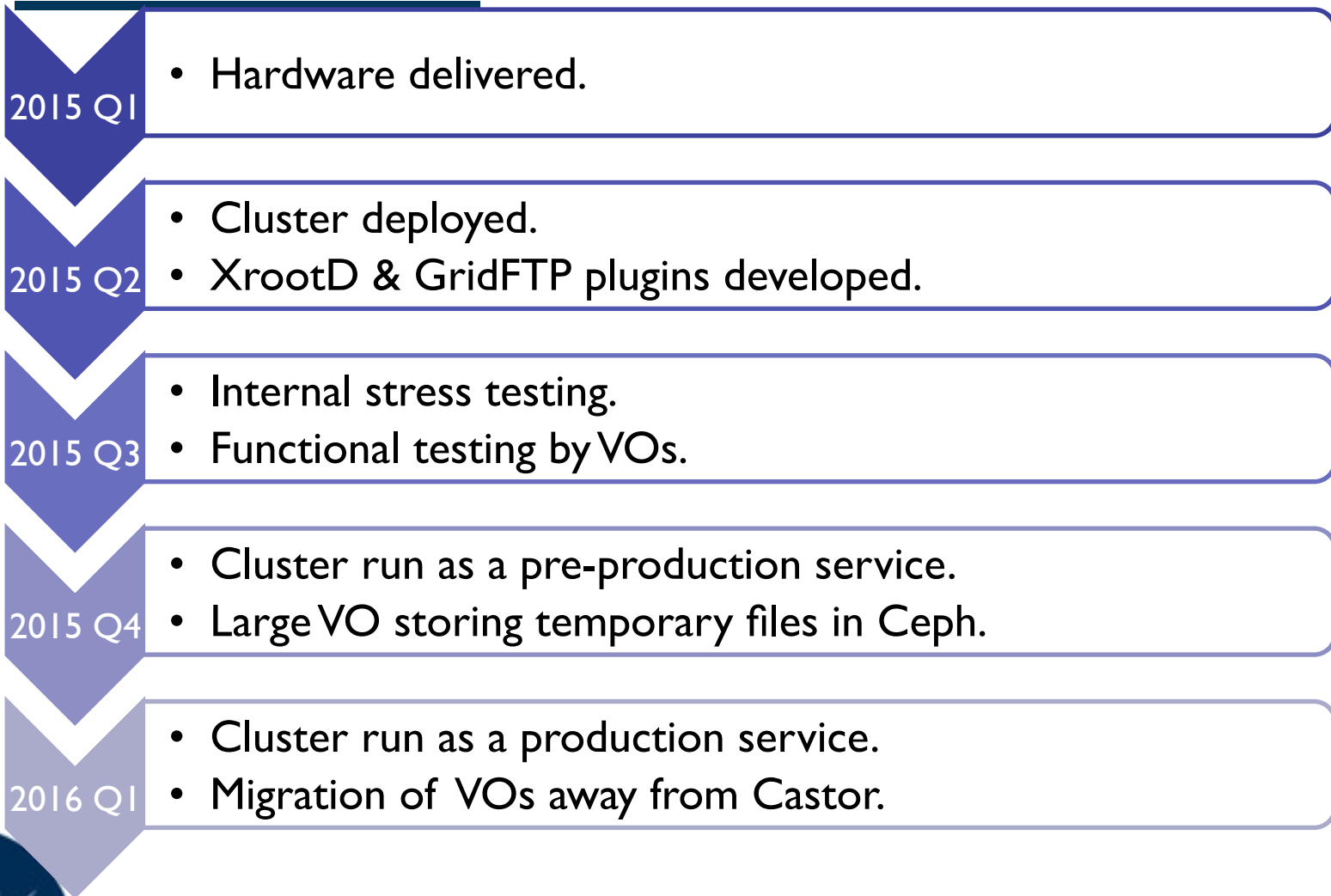
Alastair Dewhurst, 22nd July 2015

# Manpower

- Alastair Dewhurst – Project manager, responsible for making sure we deliver something the LHC VOs can use.

- James Adams has joined the project and is technical coordinator (i.e. decides how to do things).

- Shaun de Witt – Line managing George and Bruno and offering advice on xrootd/GridFTP.

- George Vasilakakos – 100% of time on Ceph.

- Bruno Canning – 30 – 60% of time on Ceph (depending on Castor commitments).

- Ian Johnson – 3 month project on GridFTP plugin (now finished).

Alastair Dewhurst, 22nd July 2015

# Timeline

**2015 Q1**
- Hardware delivered.

**2015 Q2**
- Cluster deployed.
- XrootD & GridFTP plugins developed.

**2015 Q3**
- Internal stress testing.
- Functional testing by VOs.

**2015 Q4**
- Cluster run as a pre-production service.
- Large VO storing temporary files in Ceph.

**2016 Q1**
- Cluster run as a production service.
- Migration of VOs away from Castor.

Alastair Dewhurst, 22nd July 2015

# Cloud Cluster

- Cloud cluster is providing RBD storage for cloud machines.

    - The Cloud cluster is a pre-production service.

- No significant changes to setup since CHEP/HEPiX talks.

- It is seeing increasing used throughout RAL as anyone can create machines easily.

    - Being used to run WN for testing new configurations.

- Stable running…

    - No news is good news! ☺

# Grid Cluster Setup

- 3 Physical Monitors (Dell R430)

  - Each located in a different rack with storage nodes.

- 3 Gateways (Dell R430 + 2 x 10GB/s network links)

- This years hardware has been delivered and passed acceptance testing:

  - 21 x 120TB and 26 x 100TB storage nodes

    - 64GB memory, 2 x 6 x 2.4GHz CPU, 2 x 10GB/s network cards.

    - RAID Card – Allows nodes to be used in Castor.

    - Single SSD – Purchased before we understood journaling.
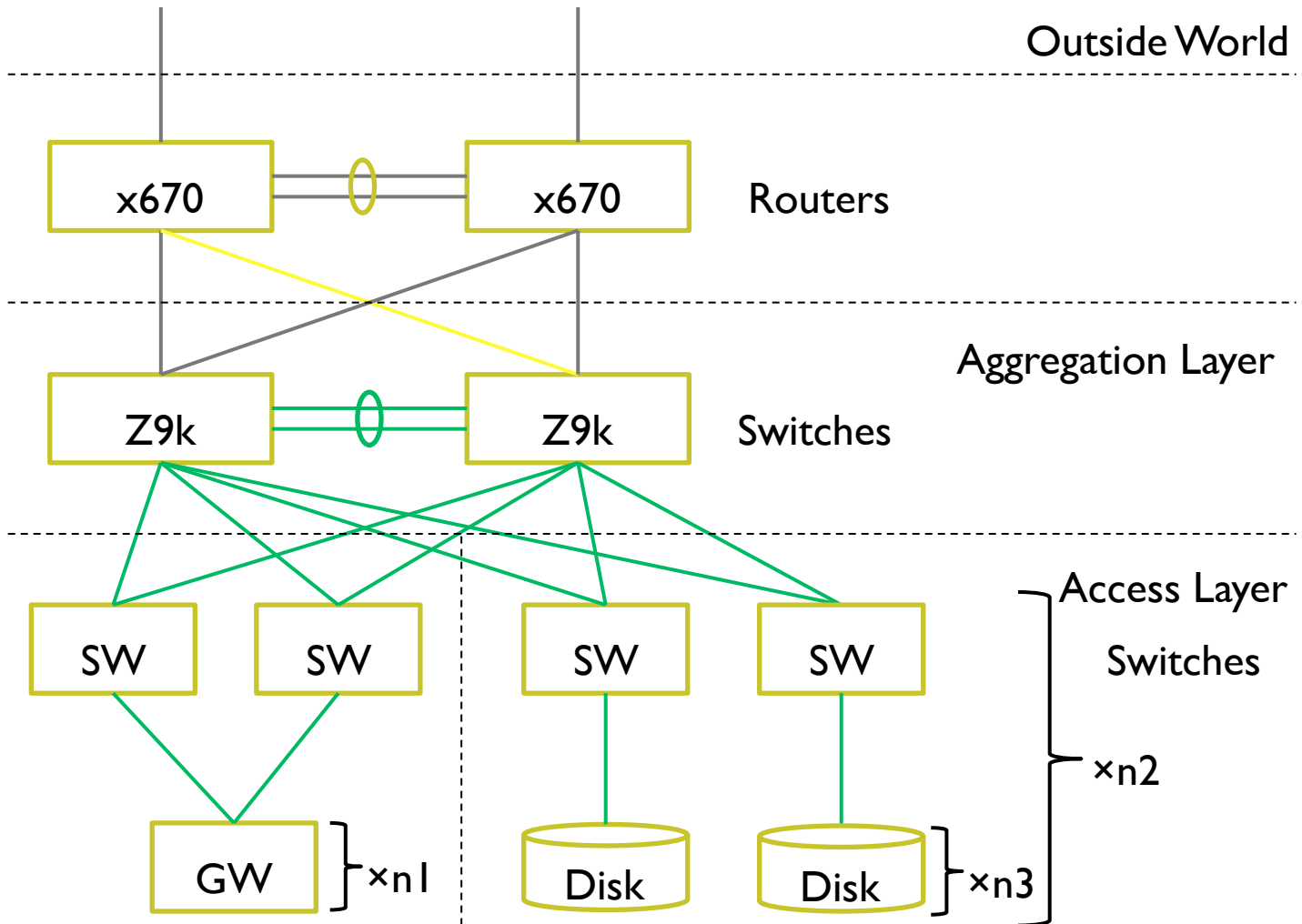
# Storage nodes

- 21 Cluster Vision Storage nodes:

  - 36 x 4TB drives

  - 2x Intel(R) Xeon(R) CPU E5-2630 v3 @ 2.40GHz5 (16 cores, 32 threads)

  - 64GB memory

- 26 OCF Storage nodes:

  - 24 x 5TB drives

  - 2 x Intel(R) Xeon(R) CPU E5-2620 v2 @ 2.10GHz (12 cores, 24 threads)

  - 64GB memory

- Both generations have 64GB memory and 2x 10GB/s network cards.

- The Cluster vision machines have a single SSD drive in them.

# Network (theory)

Outside World

x670 — x670  Routers

Aggregation Layer

Z9k — Z9k  Switches

Access Layer

SW  SW  SW  SW

Switches

×n2

GW  ×n1

Disk  Disk  ×n3

Alastair Dewhurst, 22nd July 2015

# Network (currently)

- The network switches for the gateway machines are being re-used from the 2010 procurement.

  - Oldest generation which works with 10GB/s interfaces.

  - Need to retire some machines before we have access to second switch.

  - Still connected to Tier 1 network via central switch (not directly to mesh).

- We intend to have a rebalancing network.

  - Will make use of the second network card in all storage nodes.

  - Money in this years budget for switches.

  - Martin Bly is currently evaluating what to buy.

Alastair Dewhurst, 22nd July 2015

# Gateway

- We intend our gateway machines to be identically configured and to run all plugins.

- Our Grid Cluster has a RADOS Gateway supporting S3/SWIFT.

  - Using Civetweb as easiest to configure and looks to have best future support within Ceph community.

  - Encouraging users to try it as extremely easy to support.

- We have tried to join the ATLAS Event Service, testing.

- Will take part in FTS3 S3 testing.

  - We are keeping track of WLCG http task force.

Alastair Dewhurst, 22nd July 2015

# libRadosStriper

- Sebastien Ponce (CERN), has contributed the LibRadosStriper to Ceph mainline.

- Available in Giant:

  - Bug meant files over 2GB couldn't be written.

- It was officially added to Hammer:

  - Fixed 2GB bug.

  - It doesn't currently work in Hammer due to a bug introduced when LTTng was introduced.

  - Sebastien has fixed code and this has been merged with mainline – still waiting for it to become part of official patch.

# XrootD

- Sebastien Ponce is also developing XrootD plugin for Ceph.

  - Needs to be "bullet proof" as he intends for it to be used in a future version of Castor (with Ceph underneath).

- No new testing recently as waiting for bug fixes.

  - Sebastien has been improving throughput by aligning chunk sizes.

- As well as putting XrootD plugin on Gateways we intend to make each WN an XrootD 'Gateway'.

  - Vast majority of jobs from WN use XrootD

  - Will mean WN talk directly to cluster (not via Gateway)

  - Working on making an RPM and adding to WN configuration.

- Shaun and George V have been working on authentication.

  - Currently Gridmap file but looking at adding voms awareness.

Alastair Dewhurst, 22nd July 2015

# GridFTP

- In January at XrootD workshop, a GridFTP plugin was created by Sebastien Ponce also using libRadosStriper.

- Ian Johnson worked on improving this from April to July.

- Update was presented at end of June:

  - https://indico.cern.ch/event/402898/

  - Since then Ian has aligned chunk sizes and seen a 10 fold increase in performance.

- Remaining issues:

  - While globus-url-copy has good performance, FTS transfers remain slow.

  - GridFTP plugin adds a "/" at the beginning of object names.

Alastair Dewhurst, 22nd July 2015

# Monitoring

- We have started to integrate our Ceph instances with our existing Nagios and Ganglia monitoring.

- Will be using RAL elastic search cluster.

  - So far have it parsing GridFTP logs.

  - Need to add Ceph logs but working on getting correct amount of information.

- What about Calamari?

  - Non trivial to setup.  Doesn't do everything we need.

| Host ▲▼ | Service ▲▼ | | Status ▲▼ | Last Check ▲▼ | Duration ▲▼ | Attempt ▲▼ | |
|---|---|---|---|---|---|---|---|
| gceph-mon-1 | Check CEPH Health | ✖ | OK | 15:50:04 | 3d 5h 53m 53s | 1/3 | HEALTH OK |
| | Check CEPH MON | ✖ | OK | 15:50:04 | 3d 5h 53m 53s | 1/3 | MON OK |
| gceph-mon-2 | Check CEPH Health | ✖ | OK | 16:04:04 | 3d 5h 6m 57s | 1/3 | HEALTH OK |
| | Check CEPH MON | ✖ | OK | 15:46:04 | 3d 5h 21m 55s | 1/3 | MON OK |
| gceph-mon-3 | Check CEPH Health | ✖ | OK | 16:04:04 | 3d 5h 6m 57s | 1/3 | HEALTH OK |
| | Check CEPH MON | ✖ | OK | 16:04:04 | 3d 5h 6m 57s | 1/3 | MON OK |
| gdss489 | Check CEPH OSD | ✖ | OK | 15:46:04 | 3d 7h 2m 23s | 1/3 | OSD OK |
| gdss490 | Check CEPH OSD | ✖ | OK | 15:44:05 | 3d 7h 8m 11s | 1/3 | OSD OK |

Alastair Dewhurst, 22nd July 2015

# Dashboard

- We decided that Ganglia did not provide sufficient methods to visualize metrics.

  - Looking at Grafana + InfluxDB.

- Nuffield student (A-level) Ignacy Debicki will be working on making a dashboard for the next month.



Alastair Dewhurst, 22nd July 2015

# Erasure Coding

- How are we planning on storing data?

  - 3 replicas is too expensive (we need <30% overhead)

  - Have to use Erasure Coding (EC)

- EC breaks data into 'k' chunks and creates 'm' parity chunks.

  - Can lose any 'm' OSDs without losing data.

- Use case is such that EC should work well.

  - LHC VOs write objects once and read them a few times.

  - EC does not support partial writes.

Alastair Dewhurst, 22nd July 2015

# EC benchmarking

- CERN have run tests on EC benchmarking:

  - https://cds.cern.ch/record/2015206?ln=en

  - Aim to repeat tests on our cluster and also try 16 + x EC.

- Early testing has shown Grid Cluster can get line speed into whatever node 'ceph bench' is run on.

  - Waiting for Gateways to be setup so we can scale up testing.

- Early testing has also shown that increasing object size significantly reduces performance of 'ceph bench'.

  - 4MB is default, by 16MB performance drops by 30%.

  - 'ceph bench' fails on object sizes over 100MB

Alastair Dewhurst, 22nd July 2015

# Data loss concern?

- It is inevitable that at some point we will lose data.

  - Probably combination of hardware + human error.

- What happens if we lose a placement group?

  - We lose [cluster size] / [# placement groups] amount of data.

- Objects are striped across placement groups:

  - With object size of 4MB a 1GB file would be split across 250 placement groups.

  - Losing 4MB of a 1GB file means it is completely lost.

  - Can lose 250 times as much data as Ceph thinks it has lost.

  - Larger files (normally useful data) are disproportionately affected.

# SL7

- Ceph is designed for up to date operating systems.

    - A significant chunk of the community developed stuff does not work "out of the box" with SL6.

    - We intend to move to SL7 soon™.

- 'Working' SL7 only became available recently via our configuration management system.

Alastair Dewhurst, 22nd July 2015

# Summary

- Grid Cluster is up on new hardware.

- Plugins are functionally working.

- Stress testing is starting.