

# GridPP

UK Computing for Particle Physics

## Ceph at RAL

an update

James Adams, RAL Tier 1

Ian Collier

GridPP 35, Liverpool, 11/09/2015



- Ceph is an awesome object store
- Interest all over, Tier 1, SCD, STFC and beyond
  - ISIS using CephFS in production
  - SAGE project with CCFE Culham
  - JASMIN interested in high-performance Ceph
  - BNL running pre-production Ceph instance
- Tier1 Aiming to replace CASTOR disk with Ceph
  - CASTOR tape works well and will stay
    - Looks like it may become a layer on top of Ceph anyway
  - Trying to provide thinnest possible interface layer
    - GridFTP & Xrootd servers that talk RADOS



- Tier 1 running three clusters
  - Dev
    - Development environment for gateways (see later slides)
    - Currently being used to test upgrade to latest Ceph release
  - Cloud
    - No major changes
    - Running smoothly
    - Starting to support additional larger RBDs for VMs (eg. CVMFS data disks)
  - Grid
    - Rebuilt on all new hardware
    - Lots of development and testing work
    - Focus of rest of talk



- 2015
  - Q1 : Hardware delivered
  - Q2 : Cluster deployed, developing plugins for GridFTP & Xrootd
  - Q3 : Internal stress testing, VO functional testing
  - Q4 : Launch as pre-production service with a large VO
- 2016
  - Q1 : Launch as production service, start migrating VOs from CASTOR



- Now have “real” hardware
  - Three physical monitors
    - Dell R420 with fast disks
  - Three gateway machines
    - Dell R430 with quad 10Gbps Ethernet
    - Problems deploying these... support issues with vendor...
  - 48 storage nodes
    - 100TB-120TB with dual 10Gbps Ethernet, 64GB RAM, etc...



- All deployment using Quattor (aquilon)
- Ceph component developed with University of Ghent
  - Manages configuration files
  - Drives ceph-deploy
- Can now deploy RADOSGW with Quattor
- Still no solution for secure distribution of keys
- Adding new nodes from bare metal is now only a few steps
  - Add host to aquilon
  - Set to install
  - Reboot machine...wait...
  - Set personality to Ceph storage node



- Aiming to hit a similar storage overhead to CASTOR
  - ~83% of raw storage usable
- Erasure coding decisions
  - $k$  data chunks +  $m$  parity chunks
  - We want  $m \geq 3$
  - We want  $k$  as large as possible without making performance stinky
  - Aspiring for 14+3 or 16+3
    - For comparison... Yahoo using 8+3, Facebook using 10+4
  - Currently testing lots of combinations
    - Watch this space (presentation at HEPiX)
- Possible risk with patents and jerasure library



- One or two pools per WLCG VO
- Still thinking about approach for “small” VOs
  - One big pool or individual smaller pools?
  - ACLs are one of the issues here
  - Would really rather they used S3 (or Swift)
- Placement groups are “per pool”
- Need many placement groups
  - Overhead on cluster a factor
  - Too few and the cluster won’t balance properly
  - Too many and the cluster will never recover
  - Thousands per pool as a baseline
  - Testing this now



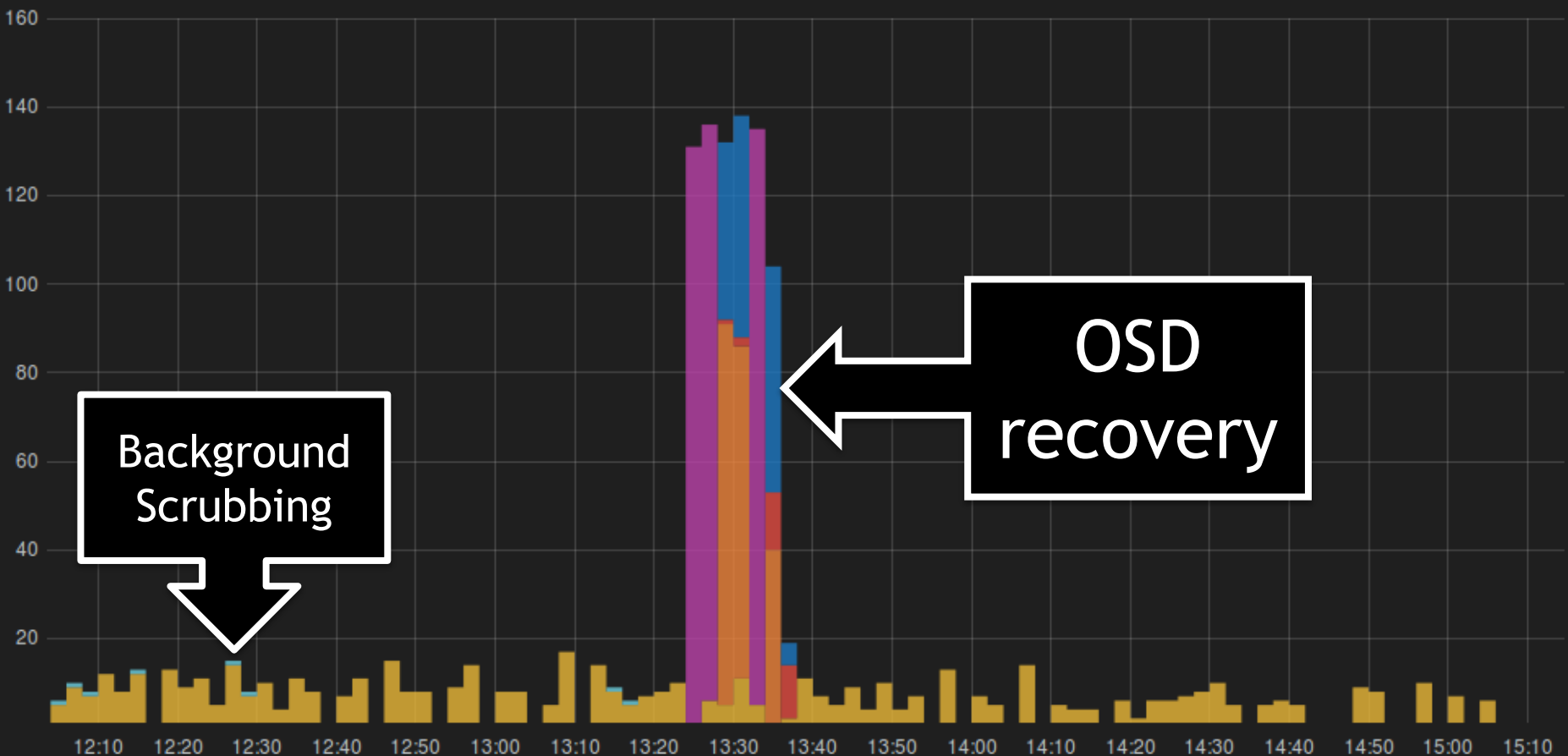


- RADOS gateway
  - Mostly interested in S3 support
    - Sub-URLs rather than subdomains (we don't control our DNS)
  - IPV6 works well (yay!)
- GridFTP gateway
  - GSI proxies working
  - Working on VOMS proxies now
  - Being packaged for release soon™
- Xrootd gateway
  - Part of Xrootd 4.2 release



- Dashboards
  - Build by Nuffield placement student Ignacy Debicki
  - InfluxDB + Grafana
  - Metric collection scripts on GitHub
    - <https://github.com/stfc/ceph-InfluxDB-metricsCollector>
  - Can now watch cluster doing its thing





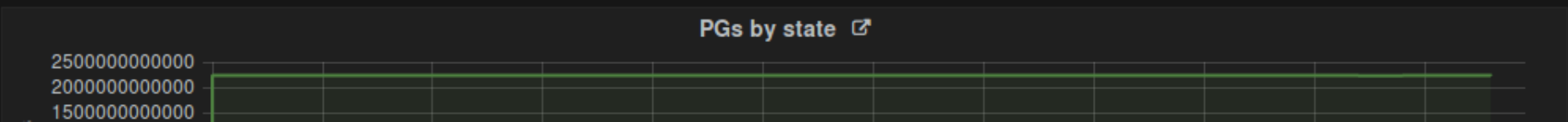
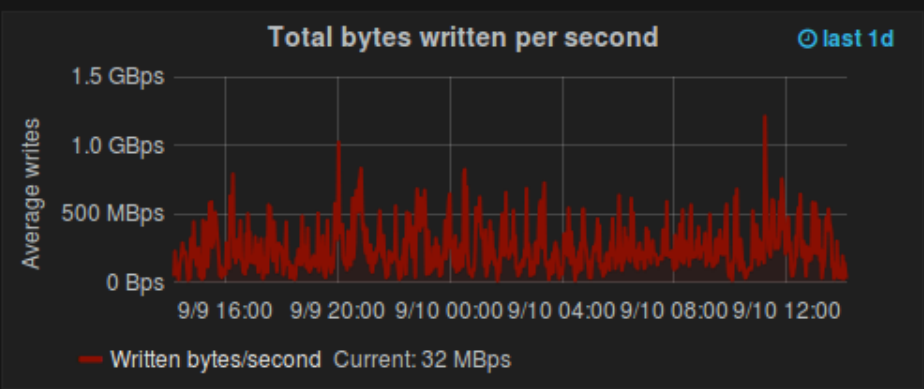
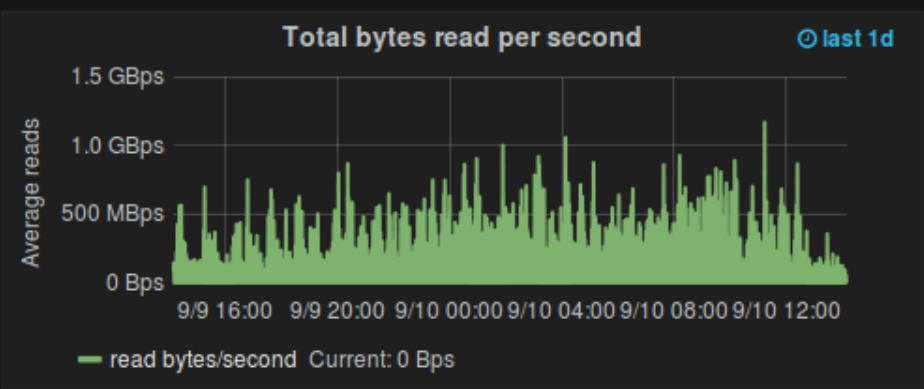
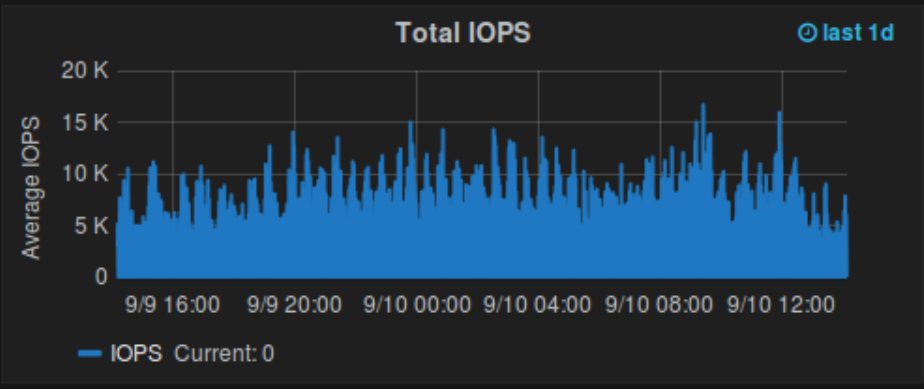
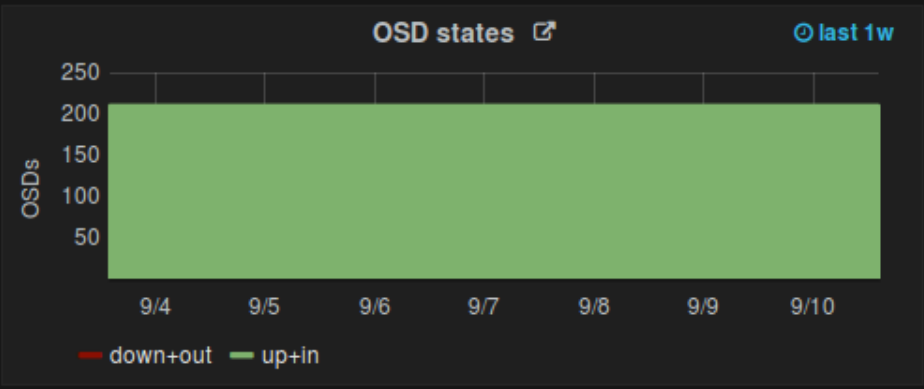
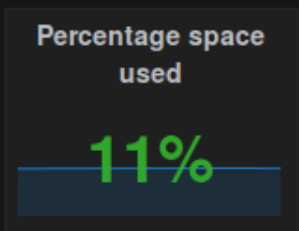
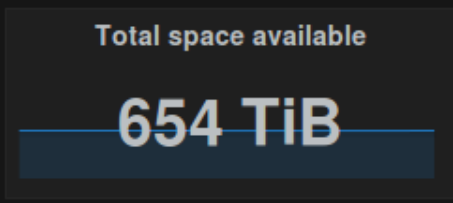
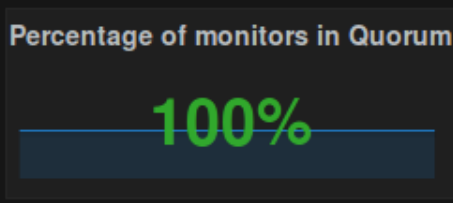
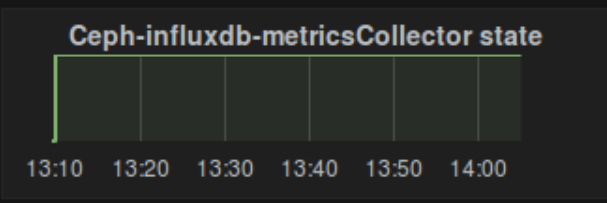
Background Scrubbing

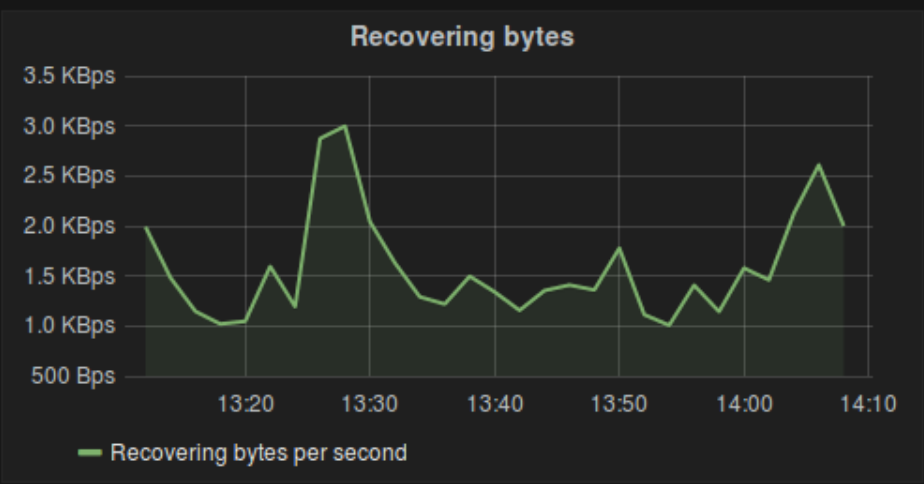
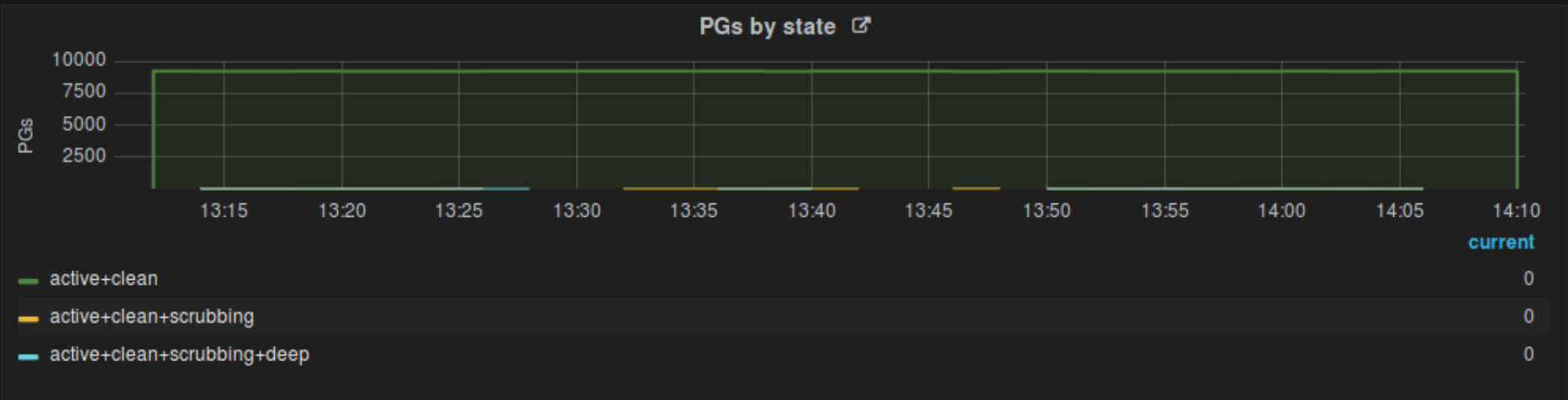
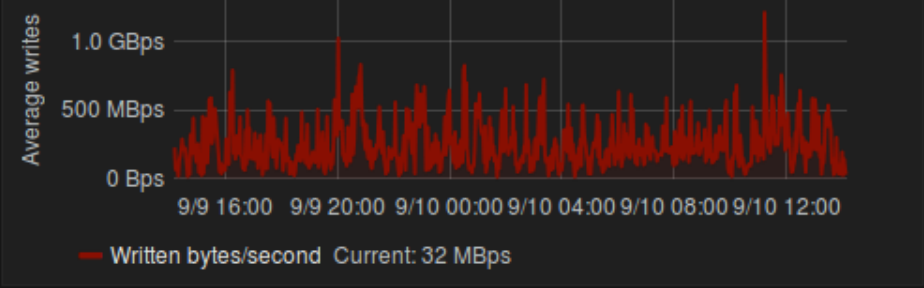
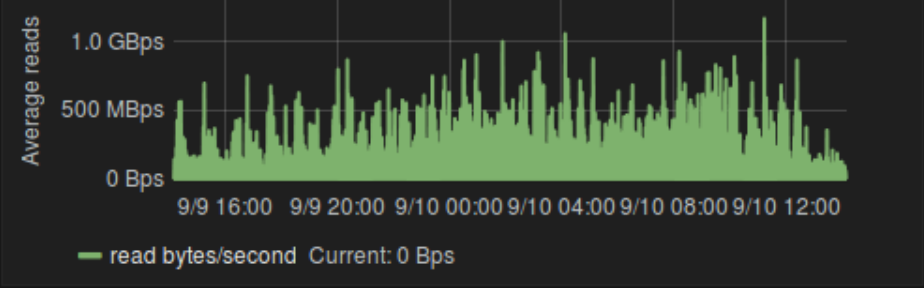
OSD recovery

State	Count
active+clean	8181
active+clean+scrubbing	11
active+clean+scrubbing+deep	0
active+degraded	0
active+recovering+degraded	0
active+recovery_wait+degraded	0
active+undersized+degraded	0



cluster: ceph-cloud





- Grid cluster is main focus of work now
  - Lots of testing to determine optimal configuration
- Start testing with VOs in October
  - Functional testing
- Production service early next year
  - Migration VO by VO
- Cloud storage cluster is operational now

