

Ceph: A scalable, organic option for Storage-as-a-Service at CERN

Arne Wiebalck Dan van der Ster

HEPiX 2013 Spring Meeting Bologna, 19 April 2013



Department







This is a report on work in progress.



Current DSS services and demands^{CEI}

- AFS as a distributed file system
- CASTOR for tape storage / archive
- EOS for analysis
- NetApp filers for special cases
- Block storage for AI VMs
- Hadoop
- "dropbox" 式 📈

Department

Current DSS services and demands^{CE}

• AFS as a distributed file system

The new demands motivate a consolidated, generic storage system ...

Hadoop ≥ memory ≤
Hadoop ≥
Hadoop >
Hadoop >
Hadoop >
Hadoo

• "dropbox" 式 🔀



Department

Wiebalck / van der Ster -- Ceph: A scalable, organic option for Storage-as-a-Service at CERN



Never heard of Ceph?



Ceph is a distributed, open-source storage system.





Never heard of Ceph?

CERN

Ceph is a distributed, open-source storage system.

Scalability

TeraBytes to ExaBytes 10s to 10'000 machines Grow and shrink





Never heard of Ceph?

CERN

Ceph is a distributed, open-source storage system.

Scalability

TeraBytes to ExaBytes 10s to 10'000 machines Grow and shrink

Reliability

No SPOF Commodity hardware Configurable Replication





Never heard of Ceph?

CERN**IT** Department

Ceph is a distributed, open-source storage system.







Never heard of Ceph?

Ceph is a distributed, open-source storage system.





CERI

Department





Object Storage Daemons



CERN

Department

Wiebalck / van der Ster -- Ceph: A scalable, organic option for Storage-as-a-Service at CERN



Reliable, self-managing, selfhealing **object** store

RADOS Replicated Autonomic Distributed Object Store

CERN

Department

Object Storage Daemons



Wiebalck / van der Ster -- Ceph: A scalable, organic option for Storage-as-a-Service at CERN

CERN IT Department CH-1211 Geneva 23

www.cern.ch/it

Switzerland



Wiebalck / van der Ster -- Ceph: A scalable, organic option for Storage-as-a-Service at CERN

CERN

Department



Wiebalck / van der Ster -- Ceph: A scalable, organic option for Storage-as-a-Service at CERN

CERN

Department

Thinly provisioned distributed block device to be used from VMs or hosts in general; Linux kernel module or KVM/QEMU/libvirt+librbd.



CERN

Department



Wiebalck / van der Ster -- Ceph: A scalable, organic option for Storage-as-a-Service at CERN

CERN IT Department CH-1211 Geneva 23 Switzerland

www.cern.ch/it

POSIX-compliant distributed file system that ships with the Linux kernel since 2.6.34; usable via kernel module (FUSE available as well).



CERN

Department



CH-1211 Geneva 23 Switzerland www.cern.ch/it

CERN IT Department



One word about CRUSH ...



Controlled Replication Under Scalable Hashing

- algorithmic data placement
 - no central meta data server
 - clients compute where data is
 - based on CRUSH map (which is "gossip'ed")
- infrastructure-aware!
 - centers, buildings, rooms, row-of-racks, racks, hosts
 - define placement rules (e.g. all replicas in different racks)
 - address correlated failures



How does it fit with what we need?

• Provide block storage for AI VMs

- OpenStack/Cinder & Glance compatible (few options)
- GlusterFS: investigated, not there yet
- NetApp: works, but vendor lock-in/closed source
- Ceph: supported, among main use cases
- Further simplify AFS (and NFS)
 - Rely on safe backend
 - rbd (or radosgw & S3)
- Future options
 - Object backend for "dropbox"-like service?
 - CephFS as next distributed file system?



Department



Our plan



Evaluate Ceph as an option for our use cases:

- Build a (serious) prototype
 - 1 PB usable

• Test it for our use cases

- OpenStack/Cinder
- AFS backend
- OwnCloud
- CephFS (too early: not prod quality yet, SLC kernel)



What we have done so far (1)



• Set up a small scale test cluster

- O 3 MON servers, 2 clients (on 3.8.2), 1 RADOSGW (all on VMs)
- 10 OSD hosts with 22 disks each (ex-CASTOR)
- O Ceph 0.56.3 installed via RPMs on SLC6.4





Wiebalck / van der Ster -- Ceph: A scalable, organic option for Storage-as-a-Service at CERN

What we've done so far (2)

ERN **T** Department

• Setup was easy

- 250TB cluster up in ~2 days
- manual (ceph-deploy did not work for us)
- O Puppet modules exist, but not tested yet
- O RADOSGW a little more tricky

Interface tests

- RADOS: create/delete/list objects
- O RBD: on 3.8.2 client, mkfs.etx4 + mount
- O CephFS: on 3.8.2 client
- O RADOSGW: CyberDuck (list, create, delete buckets)







Wiebalck / van der Ster -- Ceph: A scalable, organic option for Storage-as-a-Service at CERN

Switzerland www.cern.ch/it

CERN IT Department CH-1211 Geneva 23

What we've done so far (3)



Functional testing

- O add/remove OSDs
- O increase/decrease replication size
- O delete object in placement group
- O corrupt object in placement group



Performance/load tests

- \odot up to 900MB/s seq read from real machine on 10GE with fio (30 threads, 4M)
- O rand read 512K: 800MB/s, 4K: 3500IO/s, writes slower: 130MB/s seq, 200IO/s rand
- O 60 fio threads on rbd device plus 60 threads from RADOS bench
- s/w upgrade under load from 0.56.3 to 0.56.4 (worked for us, but there upgrade problem reports on ML)

Community support

- O quick and helpful
- O newbie questions and kernel BUG as examples



CERN IT Department CH-1211 Geneva 23 Switzerland **www.cern.ch/it**

Wiebalck / van der Ster -- Ceph: A scalable, organic option for Storage-as-a-Service at CERN





Grizzly OpenStack/Cinder (volumes)

- create/delete volume \bigcirc
- create snapshot Ο
- create vol from snapshot Ο
- Ο attach fails (gemu/libvirt version!)
- AFS backend testing
 - rbd device as vicep partition Ο
 - volume served from there \bigcirc
- OwnCloud backend
 - connected with a patched ownCloud version (S3 libs)
 - dir creation worked, everything else failed ...







Ο

Ο

Wiebalck / van der Ster -- Ceph: A scalable, organic option for Storage-as-a-Service at CERN



CH-1211 Geneva 23 Switzerland

www.cern.ch/it

Issues identified so far

- ceph-deploy did not work for us
- CRUSH complexity
 - "2 rooms 3 replica" problem
- peering is heavy (new OSDs)
- "re-weight" apocalypse
 - too little memory per OSD
 - too many simultaneous weight changes



- flaky server caused Ceph timeouts and constant rebalancing
 - slow server can slow down the cluster
 - "budgeting problem"
 - some planning required as placement groups changes difficult



Department

Wiebalck / van der Ster -- Ceph: A scalable, organic option for Storage-as-a-Service at CERN







- The feature-rich, open-source Ceph storage system is an interesting candidate to address some of our existing and emerging storage needs.
- The initial results of our investigation look promising and justify further investments.

