

Software Defined Storage: What Makes Ceph Unique

Federico Lucifredi Product Management Director, Ceph Storage Boston/Guadalajara, December 14th, 2015

























"STORAGE APPLIANCE"











SUPPORT AND MAINTENANCE

PROPRIETARY SOFTWARE

PROPRIETARY HARDWARE







philosophy	design
OPEN SOURCE COMMUNITY-FOCUSED	SCALABLE NO SINGLE POINT OF FAILURE
	SOFTWARE BASED
	SELF-MANAGING





CEPH STORAGE CLUSTER

A reliable, easy to manage, next-generation distributed object store that provides storage of unstructured data for applications













OSDs:

- 10s to 10000s in a cluster
- One per disk
 - (or one per SSD, RAID group...)
- Serve stored objects to clients
- Intelligently peer to perform replication and recovery tasks

Monitors:

- Maintain cluster membership and state
- Provide consensus for distributed decision-making
- Small, odd number
- These do **not** serve stored objects to clients







LIBRADOS

- Provides direct access to RADOS for applications
- C, C++, Python, PHP, Java, Erlang
- Direct access to storage nodes
- No HTTP overhead







RADOS Gateway:

- REST-based object storage proxy
- Uses RADOS to store objects
- API supports buckets, accounts
- Usage accounting for billing
- Compatible with S3 and Swift applications











RADOS Block Device:

- Storage of disk images in RADOS
- Decouples VMs from host
- Images are striped across the cluster (pool)
- Snapshots
- Copy-on-write clones
- Support in:
 - Mainline Linux Kernel (2.6.39+)
 - Qemu/KVM
 - OpenStack, CloudStack


RADOS

A reliable, autonomous, distributed object store comprised of self-healing, self-managing, intelligent storage nodes





Metadata Server

- Manages metadata for a POSIXcompliant shared filesystem
 - Directory hierarchy
 - File metadata (owner, timestamps, mode, etc.)
- Stores metadata in RADOS
- Does not serve file data to clients
- Only required for shared filesystem

What Makes Ceph Unique?

Part one: CRUSH







How Long Did It Take You To Find Your Keys This Morning?



De 25. To day is certainty a very disagree hay It' to day I am to apply for adble Dunday but for a wonder I have mission to the Bdr. It is of not had "the almost everlasting blue course a very important epoch which are the legitimate offering of in my life. It is not without " a troubled mind. It is very wident fear that I contemplate my comto my unind that a wrong thinker ing examination sti without dis or a wrong-doer can not possible thist-great distrust - of my un-· be hap so when he persistrilly kups tried ability that I shall enter a slippined neck unwilling to bow it into & upon, the arduous duyhead & receive The easy yoke of christ Ties of the Legal Profession but Man may out of the perverseness with an unbounded trust in The kindness & perfection of my Godo written revelution of thus endean great Greator of that See delights or to screen his faults from the reto take care of the creatures of proaches of conscience, but he can his Infinite Wisdom, Grand Heav. never ighone the revelation of hature enly Father, through the Inter. which is an over open book speak cession of Christ Bur Chavior & Ining in softness its is true but nottercessof that my career may be Withstanding irresistible of the Laws honorable & just; that I may es. of God, the real phovah, to whom pouse the side of Justice, & may be grory, praise, dominion & power, never be found in the ranks of Oppression & Injustice; That the Tr= top of the Horay Sozax tut 1 x2 weak, the helpless, the poor, may 81 J=0 24 4.2211 =Ut. Jufare , x2 0; 25 O X=112 Tr= Yeld XizI X ? Bofe. Thiging friend & an able & gealous and duccessful advocate; may all Jan 2nd in The Mero year of 1870. The "ends I aimest at, be my Lod's, s , and the Old year to gone at lastmy Country's & Truth's " gone with all its hope & fears - gone with all its happiness & sonow. W. Veirs Douic, fr. Montgomery County Marylande Ho be it! We, poor mortals, are rapidly B" DX2 to Toschologs Jeyfeto = 4f" y=0 7= kit drigting a down the stream of life EVZI KZ. OELSELL SEGK_ENG FOE her & Strange to say almost entirely Joshish - of the Se LA OR FLOY2 Frafact - Thay a Sundian, ever wine & LEN2, P. 21, Ner & JEY2 XIVE JERY TAXA MUZOER Ef R. OLORONZ - Toxie tokie Jack Kind Continue with as This year it is at all DEOLA Stirt Take & probet us as the did in the VEB MXIT OF XIJ TELO O TU XILLE one that has just passed for Trickoge of - 200 Toxoday appear strange - how very strange !!!

Dear Diary: Today I Put My Keys on the Kitchen Counter Barnaby, Flickr / CC BY 2.0





I Always Put My Keys on the Hook By the Door vitamindave, Flickr / CC BY 2.0

HOW DO YOU FIND YOUR KEYS WHEN YOUR HOUSE IS INFINITELY BIG AND ALWAYS CHANGING?











CRUSH

- Pseudo-random placement algorithm
 - Fast calculation, **no lookup**
 - Repeatable, deterministic
- Statistically uniform distribution
- Stable mapping
 - Limited data migration on change
- Rule-based configuration
 - Infrastructure topology aware
 - Adjustable replication
 - Weighting













What Makes Ceph Unique

Part two: thin provisioning



HOW DO YOU SPIN UP THOUSANDS OF VMs INSTANTLY AND EFFICIENTLY?



= 144





What Makes Ceph Unique?

Part three: clustered metadata

lrwxrwxrwx 1 root root 26 Apr 26 17:54 libgssapi_krb5.so -> mit-krb5/libgssapi krb5.s lrwxrwxrwx 1 root root 21 Apr 26 17:55 libgssapi_krb5.so.2 -> libgssapi_krb5.so.2.2 -rw-r--r-- 50 root root 216824 Jul 31 2012 libgssapi krb5.so.2.2 13 Apr 26 17:54 libgs.so.8 -> libgs.so.8.71 lrwxrwxrwx 1 root root -rw-r--r-- 17 root root 9478048 Jan 25 2011 libgs.so.8.71 lrwxrwxrwx 1 root root 21 Apr 26 17:55 libgssrpc.so -> mit-krb5/libgssrpc.so lrwxrwxrwx 1 root root 16 Apr 26 17:55 libgssrpc.so.4 -> libgssrpc.so.4.1 -rw-r--r-- 50 root root 115352 Jul 31 2012 libgssrpc.so.4.1 21832 Sep 8 2010 libgthread-2.0.a -rw-r--r-- 50 root root -rw-r--r-- 50 root root 972 Sep 8 2010 libgthread-2.0.la lrwxrwxrwx 1 root root 26 Apr 26 17:55 libgthread-2.0.so -> libgthread-2.0.so.0.2400 26 Apr 26 17:55 libgthread-2.0.so.0 -> libgthread-2.0.so.0.240 lrwxrwxrwx 1 root root 17704 Sep 8 2010 libgthread-2.0.so.0.2400.2 -rw-r--r-- 50 root root 4096 Apr 26 18:00 libgtk2.0-0 drwxr-xr-x 2 root root 9275282 Oct 14 2010 libgtk-x11-2.0.a -rw-r--r-- 49 root root 981 Oct 14 2010 libgtk-x11-2.0.la -rw-r--r-- 49 root root lrwxrwxrwx 1 root root 26 Apr 26 17:55 libgtk-x11-2.0.so -> libgtk-x11-2.0.so.0.2000 26 Apr 26 17:55 libgtk-x11-2.0.so.0 -> libgtk-x11-2.0.so.0.200 lrwxrwxrwx 1 root root -rw-r--r-- 49 root root 4319784 Oct 14 2010 libgtk-x11-2.0.so.0.2000.1 15 Apr 26 17:55 libgvc.so -> libgvc.so.5.0.0 lrwxrwxrwx 1 root root lrwxrwxrwx 1 root root 15 Apr 26 17:55 libgvc.so.5 -> libgvc.so.5.0.0 -rw-r--r-- 49 root root 504424 Jul 5 2010 libgvc.so.5.0.0 16 Apr 26 17:55 libgvpr.so -> libgvpr.so.1.0.0 lrwxrwxrwx 1 root root 16 Apr 26 17:55 libgvpr.so.1 -> libgvpr.so.1.0.0 lrwxrwxrwx 1 root root -rw-r--r-- 50 root root 482856 Jul 5 2010 libgvpr.so.1.0.0 -rw-r--r-- 50 root root 267948 Apr 13 2009 libHalf.a lrwxrwxrwx 1 root root 16 Apr 26 17:55 libHalf.so -> libHalf.so.6.0.0 lrwxrwxrwx 1 root root 16 Apr 26 17:55 libHalf.so.6 -> libHalf.so.6.0.0 -rw-r--r-- 50 root root 269992 Apr 13 2009 libHalf.so.6.0.0 -rw-r--r-- 50 root root 52850 Nov 1 2009 libhistory.a

POSIX Filesystem Metadata Barnaby, Flickr / CC BY 2.0







three metadata servers







??





















DYNAMIC SUBTREE PARTITIONING
Getting Started With Ceph

Have a working cluster up quickly.

Read about the latest version of Ceph.

• The latest stuff is always at http://ceph.com/get

Deploy a test cluster using ceph-deploy.

• Read the quick-start guide at http://ceph.com/qsg

Deploy a test cluster on the AWS free-tier using Juju.

• Read the guide at http://ceph.com/juju

Read the rest of the docs!

• Find docs for the latest release at http://ceph.com/docs

Getting Involved With Ceph

Help build the best storage system around!

Most project discussion happens on the mailing list.
Join or view archives at <u>http://ceph.com/list</u>

IRC is a great place to get help (or help others!)Find details and historical logs at <u>http://ceph.com/irc</u>

The tracker manages our bugs and feature requests.

Register and start looking around at <u>http://ceph.com/tracker</u>

Doc updates and suggestions are always welcome.

• Learn how to contribute docs at http://ceph.com/docwriting

Ceph Hammer (v0.94.x) Best Ceph ever.

- 1. Rados Performance enhancements: All Flash environments
- 2. Simplified RGW deployment
- 3. RGW Object Versioning and Bucket Sharding
- 4. RBD Mandatory Locking, Object Maps, Copy on Read
- 5. CephFS Snapshot improvements

and many more. See https://ceph.com/releases/v0-94-hammer-released/





Questions?

Federico Lucifredi PM Director, Ceph

federico@redhat.com @0xF2

redhat.com | ceph.com

