# Ceph

scalable, unified storage
files, blocks & objects

# Storage system

# Open Source

LPGL2
no copyright assignment

# Incubated by DreamHost

started by Sage Weil at UC Santa Cruz, research group partially funded by tri-labs
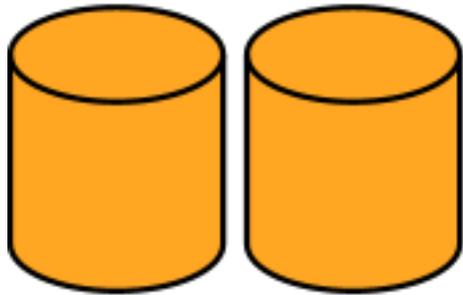
# 50+ contributors

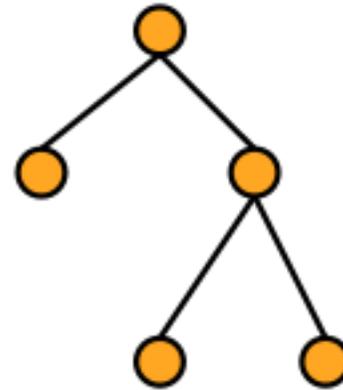around the world

# Commodity hardware

# No SPoF

# No bottlenecks

# Smart storage

peers detect, gossip, heal
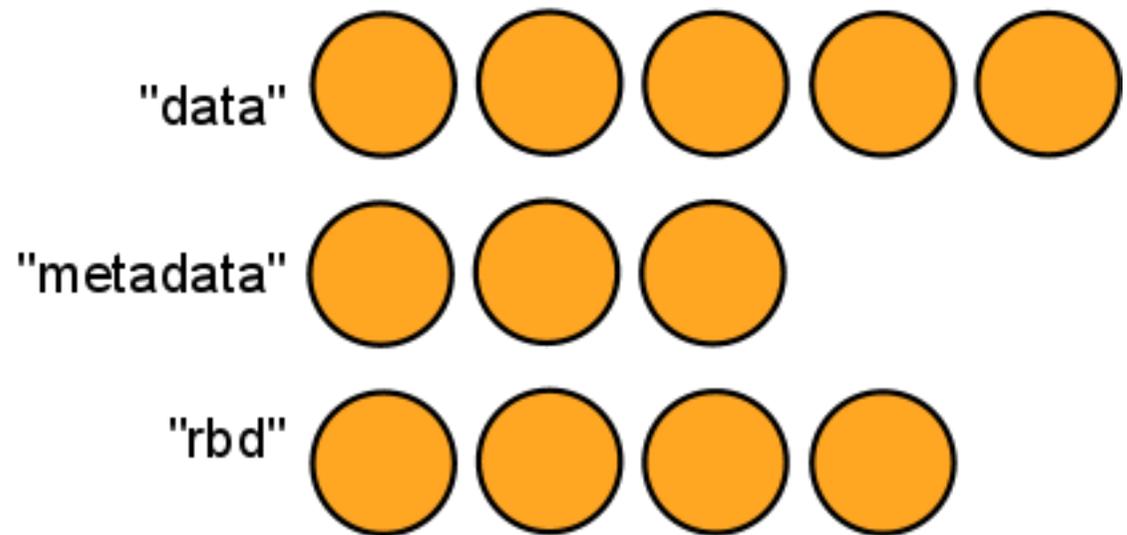
# Block devices

# Distributed file system

# Objects

"data"

"metadata"

"rbd"
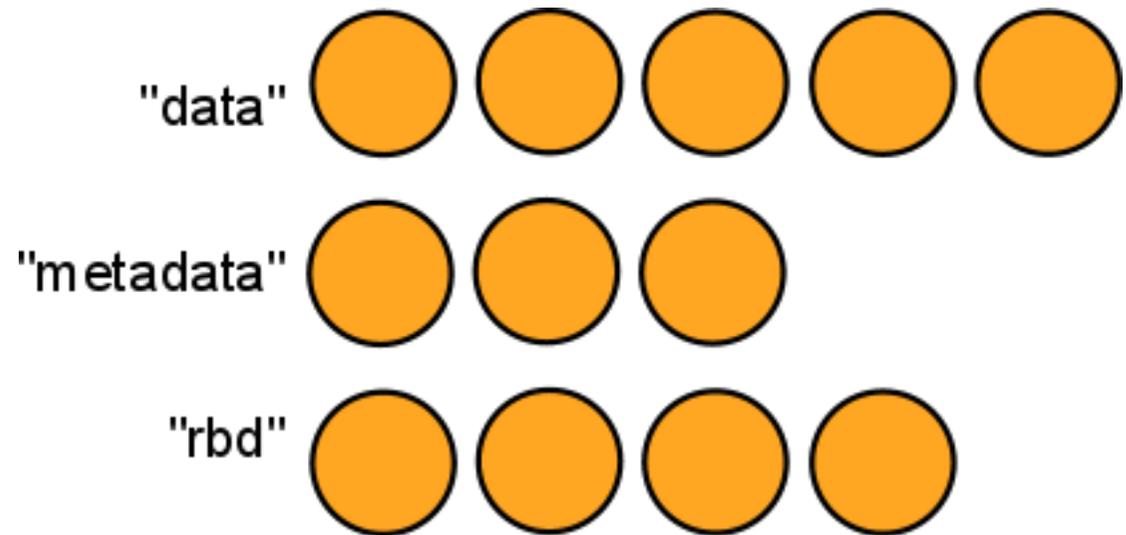
# Monitors

# pool, name

→ data (bytes),

metadata: key=value, k2=v2, ...

librados (C)
libradospp (C++)
Python
PHP
*your favorite language here*

# Smart client

talk to the cluster, not to a gateway
compound operations
choose your consistency (ack/commit)

# Pools

replica count,
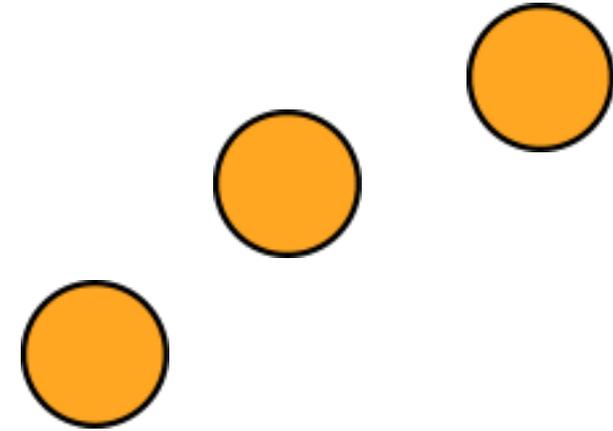access control,
placement rules,

...

zone
row
rack
host
disk

# CRUSH

deterministic placement algorithm
no lookup tables for placement
DC topology and health as input
balances at scale

# Autonomous

others say: expect failure
we say: expect balancing
failure, expansion, replica count, ...

# btrfs / ext4 / xfs / *

really, anything with xattrs
btrfs is an optimization
can migrate one disk at a time
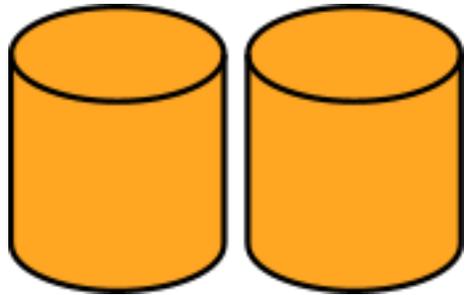
# process per X

X = disk, RAID set, directory
tradeoff: RAM & CPU vs fault isolation

# RADOS gateway

adds users, per-object access control
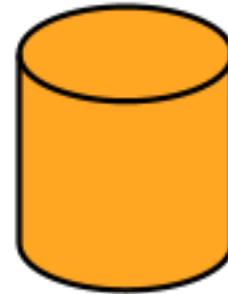HTTP, REST, looks like S3 and Swift

# i <3 boto

use any s3 client
just a different hostname
we'll publish patches & guides

# RBD

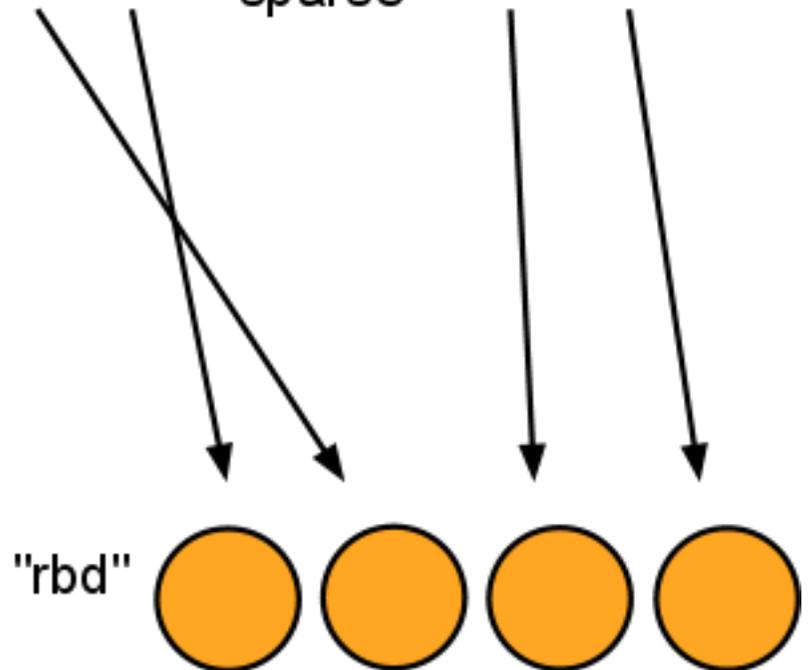RADOS Block Device

# Block device

# Chunks
(4MB, configurable)

sparse

# Objects

"rbd"

# Live migration

one-line patch to libvirt
don't assume everything is a filename

# Snapshots

cheap, fast
rbd create *mypool/myimage@mysnap*

# Copy on Write

layering aka base image
soon

# /dev/rbd0
# /dev/rbd/*

rbd map *imagename*

# QEmu/KVM driver

no root needed
shorter codepath

# Ceph Distributed Filesystem

# mount -t ceph

or FUSE

# High Performance Computing

# libcephfs

no need to mount, no FUSE
no root access needed
also from Java etc
Samba, NFS etc gateways

# Hadoop shim

replaces HDFS,
avoids NameNode and DataNode

devops
devops
devops

# Chef cookbooks

Open Source
on Github
soon

# Barclamp

Open Source
on Github
soon

# devving to help ops

new store node
hard drive replacement
docs, polish, QA

# Questions?

ceph.newdream.net

github.com/NewDreamNetwork

tommi.virtanen@dreamhost.com

P.S. we're hiring!

# Bonus round

# Want iSCSI?

export an RBD
potential SPoF & bottleneck
not a good match for core Ceph
*your product here*

# s3-tests

unofficial S3 compliance test suite
run against AWS, codify responses

- `squeeze` (Debian 6.0)
- `lenny` (Debian 5.0)
- `oneiric` (Ubuntu 11.11)
- `natty` (Ubuntu 11.04)
- `maverick` (Ubuntu 10.10)

**Todo:** http://ceph.newdream.net/debian/dists/ also has `lucid` (Ubuntu 10.04), should that be removed?

Whenever we say *DISTRO* below, replace that with the codename of your operating system.

Run these commands on all nodes:

```
wget -q -O- https://raw.github.com/NewDreamNetwork/ceph/master/keys/release.asc \
| sudo apt-key add -

sudo tee /etc/apt/sources.list.d/ceph.list <<EOF
deb http://ceph.newdream.net/debian/ DISTRO main
deb-src http://ceph.newdream.net/debian/ DISTRO main
EOF

sudo apt-get update
sudo apt-get install ceph
```

**Todo:** For older distributions, you may need to make sure your apt-get may read .bz2 compressed files. This works for Debian Lenny 5.0.3: `apt-get install bzip2`

**Todo:** Ponder packages; ceph.deb currently pulls in gceph (ceph.deb Recommends: ceph-client-tools ceph-fuse libcephfs1 librados2 librbd1 btrfs-tools gceph) (other interesting: ceph-client-tools ceph-fuse libcephfs-dev librados-dev librbd-dev obsync python-ceph radosgw)

## Red Hat / CentOS / Fedora

### Status as of 2011-09

# Teuthology

study of cephalopods
multi-machine dynamic tests
Python, gevent, Paramiko

cluster.only('osd').run(args=['uptime'])

roles:
- [mon.0, mds.0, osd.0]
- [mon.1, osd.1]
- [mon.2, osd.2]
- [client.0]

```yaml
tasks:
- ceph:
- trashosds:
    op_delay: 1
    chance_down: 10
- kclient:
- workunit:
    all:
      - suites/bonnie.sh
```
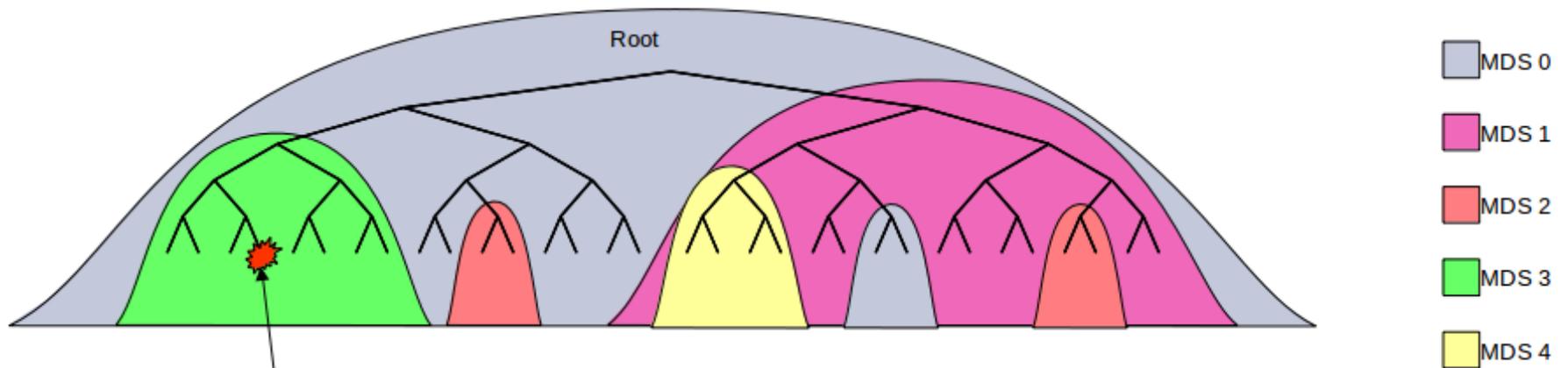
# ceph-osd plugins

SHA-1 without going over the network
update JSON object contents

Root

MDS 0
MDS 1
MDS 2
MDS 3
MDS 4

Busy directory fragmented across many MDS's

# Questions?

ceph.newdream.net

github.com/NewDreamNetwork

tommi.virtanen@dreamhost.com

P.S. we're hiring!