## CEPH TUTORIAL

GridKA School 2016

Diana Gudu

August 30, 2016

Karlsruhe Institute of Technology

# INTRODUCTION ROUND

Diana Gudu

· PhD researcher in Computer Science @KIT
  · distributed multi-agent framework for trading cloud resources
· Human Brain Project
  · work on cloud storage and computing services
· MSc in Computational Science and Engineering @TU Munich
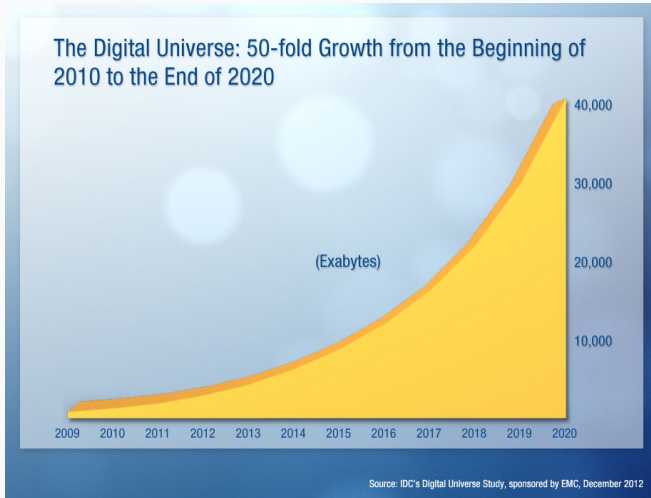· BSc in Computer Science @Polytechnic University of Bucharest
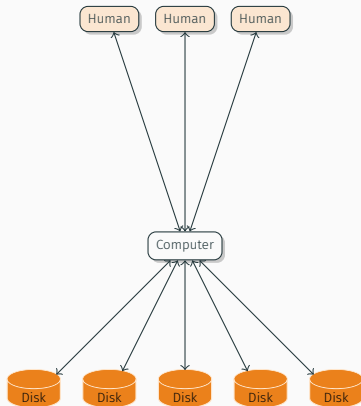
>_

# EVOLUTION OF STORAGE

The Digital Universe: 50-fold Growth from the Beginning of 2010 to the End of 2020

(Exabytes)

Source: IDC's Digital Universe Study, sponsored by EMC, December 2012

©trumpiot.co

©Oracle

# CEPH

Philosophy

- open-source

Philosophy

- open-source
- community focused

Philosophy

- open-source
- community focused
- software-defined

Philosophy

- open-source
- community focused
- software-defined
- scale-out hardware, no SPF

### Philosophy

· open-source
· community focused
· software-defined
· scale-out hardware, no SPF
· self-managing

### Philosophy

- open-source
- community focused
- software-defined
- scale-out hardware, no SPF
- self-managing
- failure is normal

## Philosophy

· open-source
· community focused
· software-defined
· scale-out hardware, no SPF
· self-managing
· failure is normal

## History

## Philosophy

· open-source
· community focused
· software-defined
· scale-out hardware, no SPF
· self-managing
· failure is normal

## History



2016
2014
2012
2010
2006
PhD thesis at UCSC    2004

## Philosophy

- open-source
- community focused
- software-defined
- scale-out hardware, no SPF
- self-managing
- failure is normal

## History



2016
2014
2012
2010
Project is open-sourced — 2006
PhD thesis at UCSC — 2004

## Philosophy

· open-source

· community focused

· software-defined

· scale-out hardware, no SPF

· self-managing

· failure is normal

## History



Included in Linux kernel — 2010

Project is open-sourced — 2006

PhD thesis at UCSC — 2004

2016

2014

2012

## Philosophy

- open-source
- community focused
- software-defined
- scale-out hardware, no SPF
- self-managing
- failure is normal

## History



Integrated into CloudStack — 2012
Included in Linux kernel — 2010
Project is open-sourced — 2006
PhD thesis at UCSC — 2004
2014
2016

## Philosophy

· open-source
· community focused
· software-defined
· scale-out hardware, no SPF
· self-managing
· failure is normal

## History



2016
RedHat acquisition 2014
Integrated into CloudStack 2012
Included in Linux kernel 2010
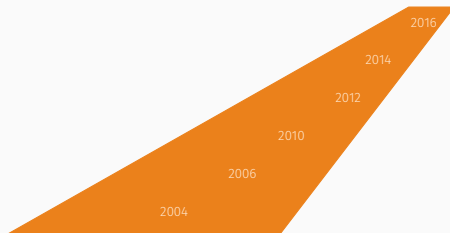Project is open-sourced 2006
PhD thesis at UCSC 2004
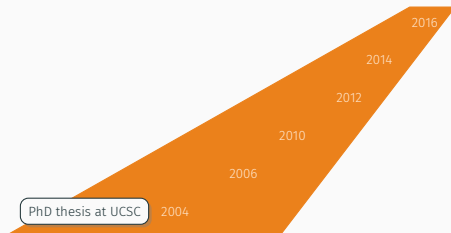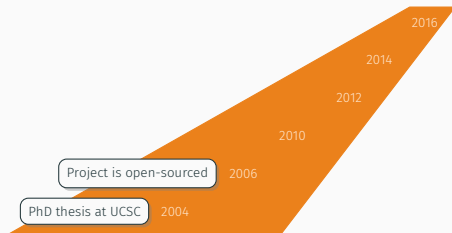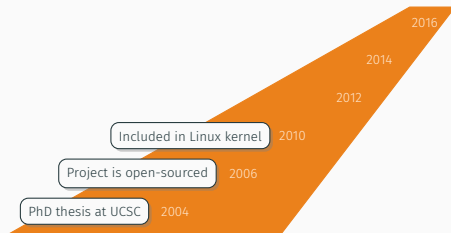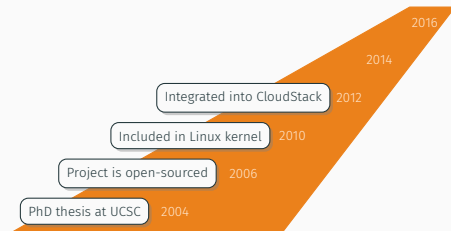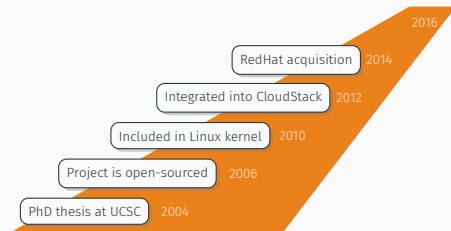
## Philosophy

- open-source
- community focused
- software-defined
- scale-out hardware, no SPF
- self-managing
- failure is normal

## History

| | |
|---|---|
| CephFS is stable | 2016 |
| RedHat acquisition | 2014 |
| Integrated into CloudStack | 2012 |
| Included in Linux kernel | 2010 |
| Project is open-sourced | 2006 |
| PhD thesis at UCSC | 2004 |

Application   REST   Host/VM   FS client

| RADOSGW | RBD | CephFS |
|---|---|---|
| bucket-based REST gateway | virtual block device | POSIX-compliant |
| S3- and Swift-compatible | Linux kernel client | Linux kernel client |
| | QEMU/KVM driver | FUSE support |

**librados**

allows apps direct access to RADOS

support for C, C++, Java, Python, Ruby, PHP

**RADOS**

A reliable, autonomous, distributed object store

consisting of self-healing, self-managing, intelligent storage nodes

### OSD

· serve objects to clients
· one per disk
· backend: btrfs, xfs, ext4
· peer-to-peer replication and recovery
· write-ahead journal

### MON

· maintain cluster state and membership
· vote for distributed decision-making
· small, odd number

## DATA PLACEMENT

How to…

…ensure *infinite* scalability?

How to...

...ensure *infinite* scalability?

· **compute** object locations instead of looking them up

How to...

  ...ensure *infinite* scalability?
   · **compute** object locations instead of looking them up
  ...ensure deterministic placement?

How to...

   ...ensure *infinite* scalability?

     · **compute** object locations instead of looking them up

   ...ensure deterministic placement?

     · **pseudo** random algorithm

How to…

…ensure *infinite* scalability?
· **compute** object locations instead of looking them up

…ensure deterministic placement?
· **pseudo** random algorithm

…minimize data movement?

How to...

...ensure *infinite* scalability?

· **compute** object locations instead of looking them up

...ensure deterministic placement?
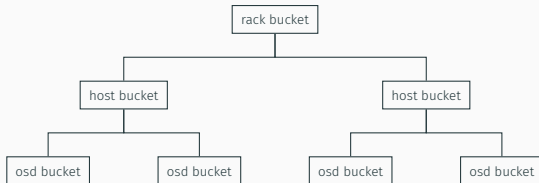
· **pseudo** random algorithm

...minimize data movement?

· stable mapping, e.g. consistent hashing
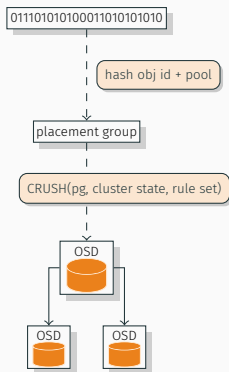
Controlled Replication Under Scalable Hashing
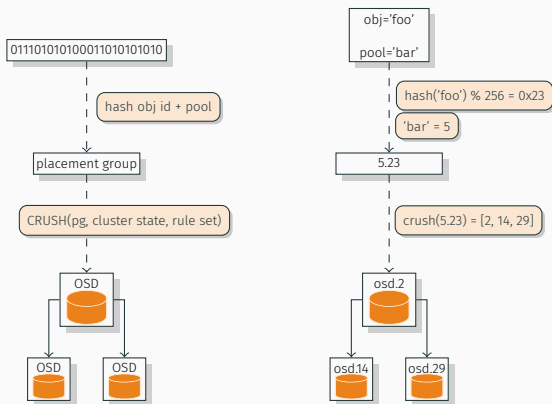
· pseudo-random placement algorithm
· repeatable, deterministic
· statistically uniform distribution
· stable mapping: minimal data migration
· rule-based configuration, topology aware

Controlled Replication Under Scalable Hashing

· pseudo-random placement algorithm
· repeatable, deterministic
· statistically uniform distribution
· stable mapping: minimal data migration
· rule-based configuration, topology aware

### Pools

- logical groups for storing objects
- manage
  - # PGs
  - # replicas
  - ruleset
  - permissions
  - snapshots

## Pools

· logical groups for storing objects
· manage
  · # PGs
  · # replicas
  · ruleset
  · permissions
  · snapshots

## Placement groups

· fragments of logical groups
  · aggregate objects within a pool for scalability
· 1 PG over several OSDs ( # replicas)
· several PGs per OSD ( $\approx 100$ )
· more PGs
  $\implies$ data durability
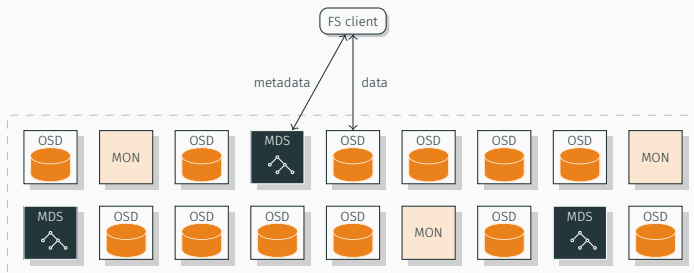  $\implies$ even distribution
  $\implies$ resource overhead

# CEPH CLIENTS

- direct access to RADOS for applications
- C, C++, Python, Java, Erlang, PHP
- native socket access, no HTTP overhead

- RESTful API
- unified object namespace
- S3 and Swift compatible
- user database and access control
- usage accounting, billing

- storage of disk images in RADOS
- images are striped across the cluster
- decoupling of VMs from host
- thin provisioning
  - physical storage only used once you begin writing
- snaphots, copy-on-write clones
- support in Qemu, KVM

- files striped over RADOS objects
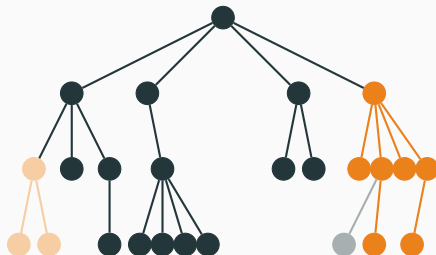- strong consistency for POSIX semantics
  - use O_LAZY to relax consistency

- files striped over RADOS objects
- strong consistency for POSIX semantics
  - use O_LAZY to relax consistency

### Metadata Server

- manages metadata for POSIX-compliant filesystem
  - directory hierarchy
  - file metadata: owner, timestamps, mode etc
- stores metadata in RADOS
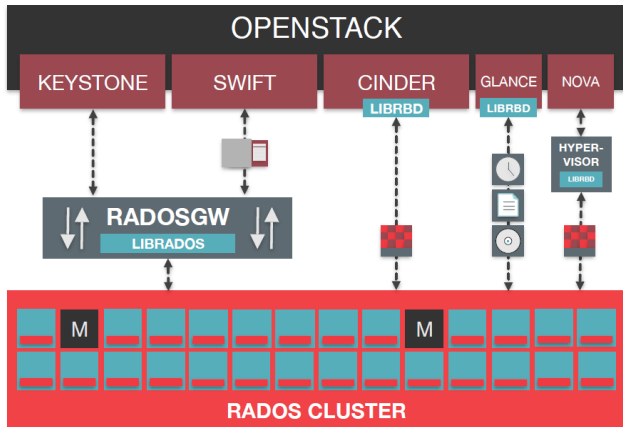- multiple MDS for HA and load balancing

# FINAL REMARKS

- CRUSH magic
- clustered, dynamic metadata management
- thin provisioning of block storage
- unified storage (object, block, file)

| | HDFS | iRODS | Ceph | GlusterFS | Lustre |
|---|---|---|---|---|---|
| Architecture | centralized | centralized | distributed | decentralized | centralized |
| Naming | index | database | CRUSH | EHA | index |
| API | CLI, FUSE, REST, API | CLI, FUSE, API | FUSE, mount, REST | FUSE, mount | FUSE |
| Fault detection | fully connect. | P2P | fully connect. | detected | manually |
| System availability | no failover | no failover | high | high | failover |
| Data availability | replication | replication | replication | RAID-like | no |
| Placement strategy | auto | manual | auto | manual | no |
| Replication | async. | sync. | sync. | sync. | RAID-like |
| Cache consistency | WORM, lease | lock | lock | no | lock |
| Load balancing | auto | manual | manual | manual | no |

- separate cluster and public networks
- SSD journals: accelerate bursts and random writes
- memory: 1-2 GB per OSD, CPU: 1.5 GHz per OSD
- redundancy
  - replication: increased read performance, capacity impact
  - erasure coding: better space efficiency, high CPU overhead, ~~RBD~~
  - cache tiering: write-back overlay pool; combine replication and erasure coding

©RedHat

- docs: ceph.com/docs
- wiki: wiki.ceph.com
- mailing lists
  - ceph-users: ceph-users@ceph.com
  - ceph-devel: ceph-devel@vger.kernel.org
- IRC channels (server irc.oftc.net)
  - #ceph
  - #ceph-devel
- github: github.com/ceph/ceph
- get involved: ceph.com/community/contribute

## TUTORIAL

- deploy a Ceph cluster
- basic operations with the storage cluster
- data placement: CRUSH
- Ceph Filesystem
- block storage: RBD
- advanced topics: erasure coding
- troubleshooting challenge