

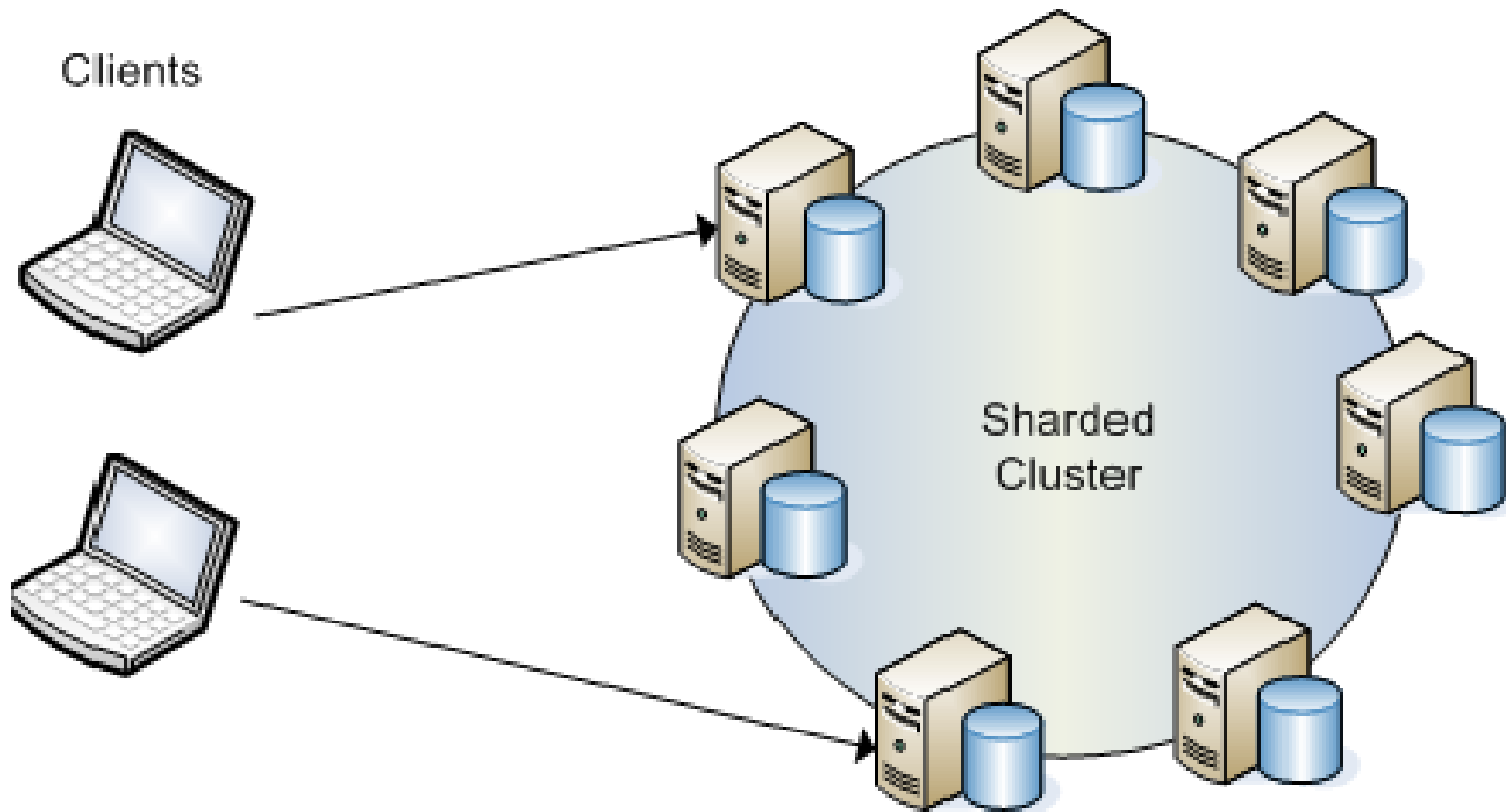
A Ceph, mint adattároló klaszter megoldás

készítette:

Szalai László,

Major Kálmán (NYME INGA)

Elosztott adattárolás séma



Mik merülhetnek fel

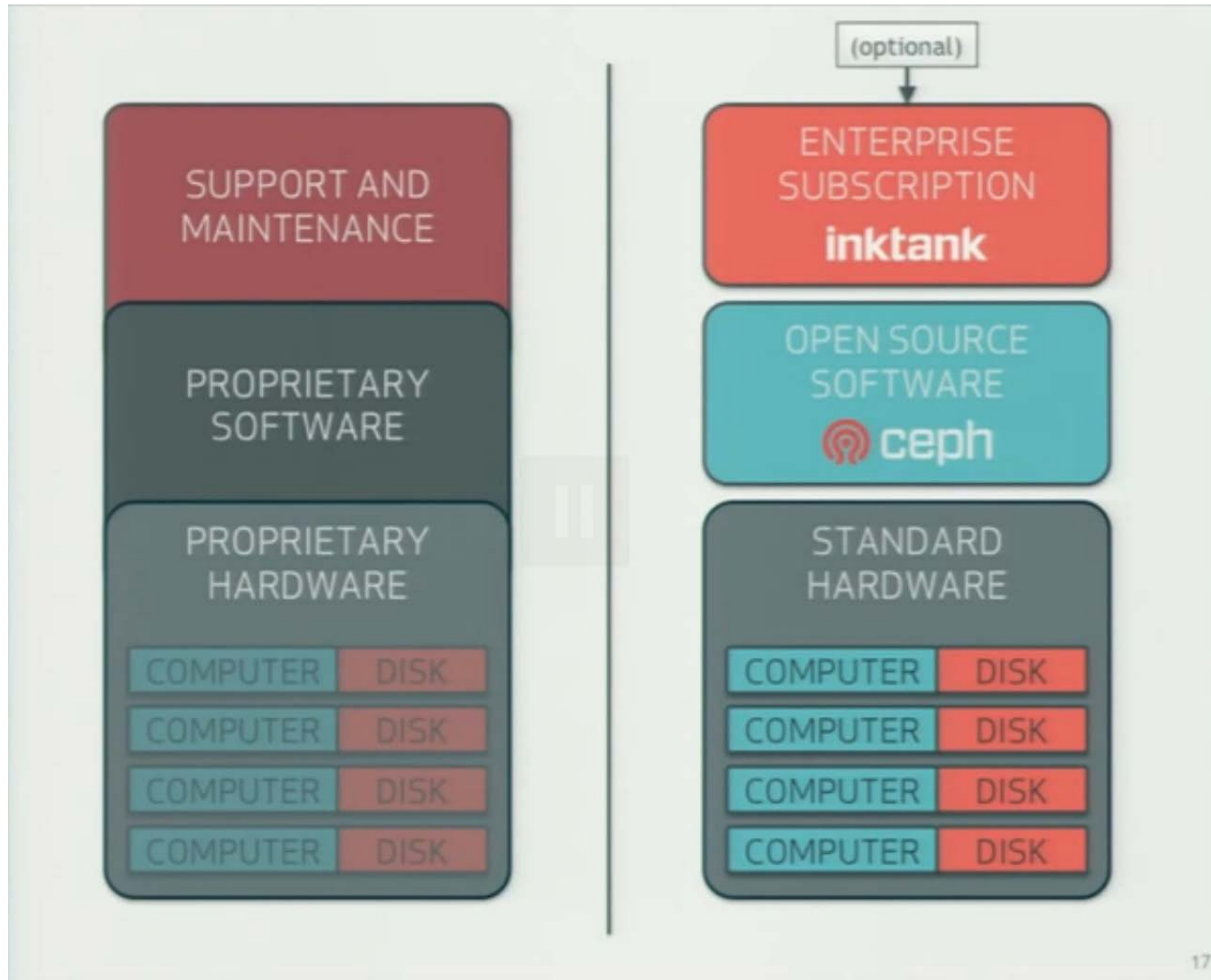
<i>If you need...</i>	Gluster	Lustre	Ceph
NAS and Scale Out NAS	✓	✓	✓
SAN			✓
Shared Filesystems	✓	✓	✓
Object Storage	✓		✓

(MooseFS, XtremFS, stb.)

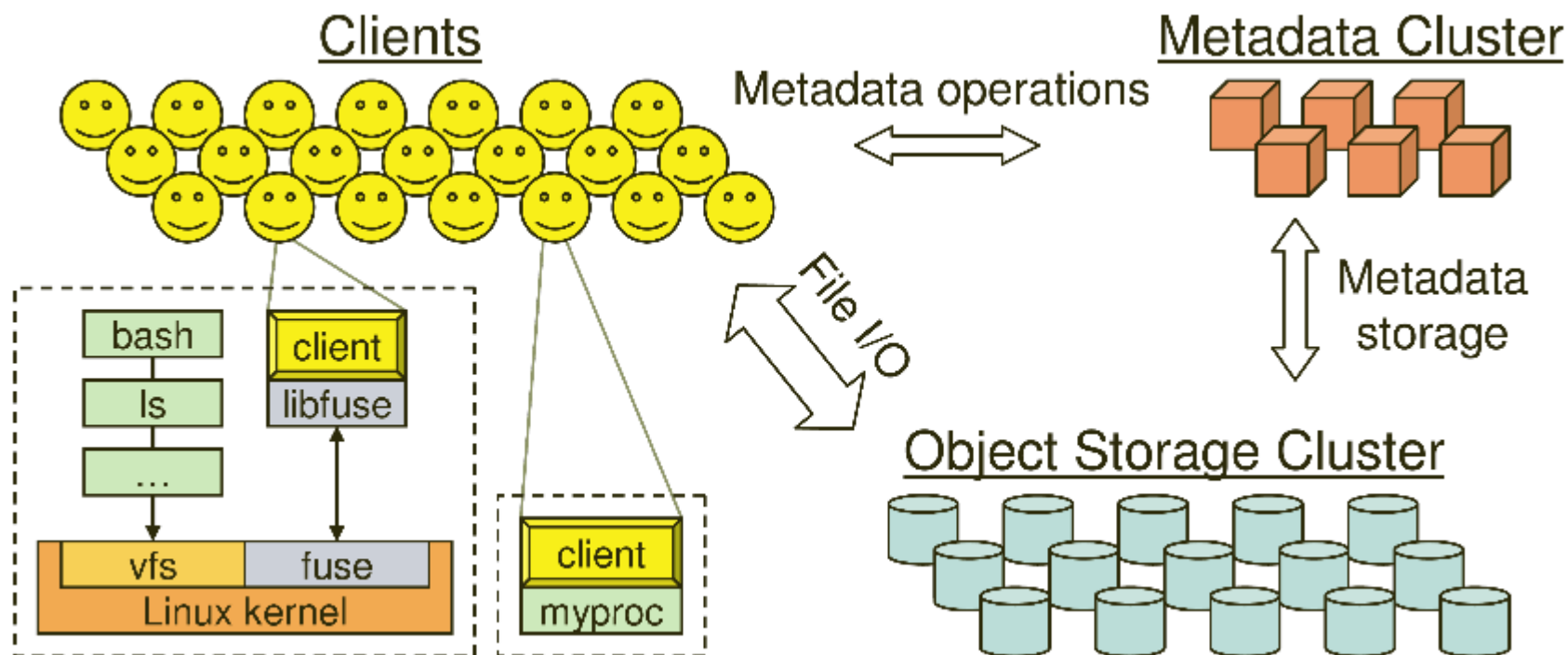
ceph, mint választott

- Antik görög kifejezés
- kezdetben doktori munka
- 2007-től három fős fejlesztőbrigád
- 2012, Inktank Storage cég
 - support

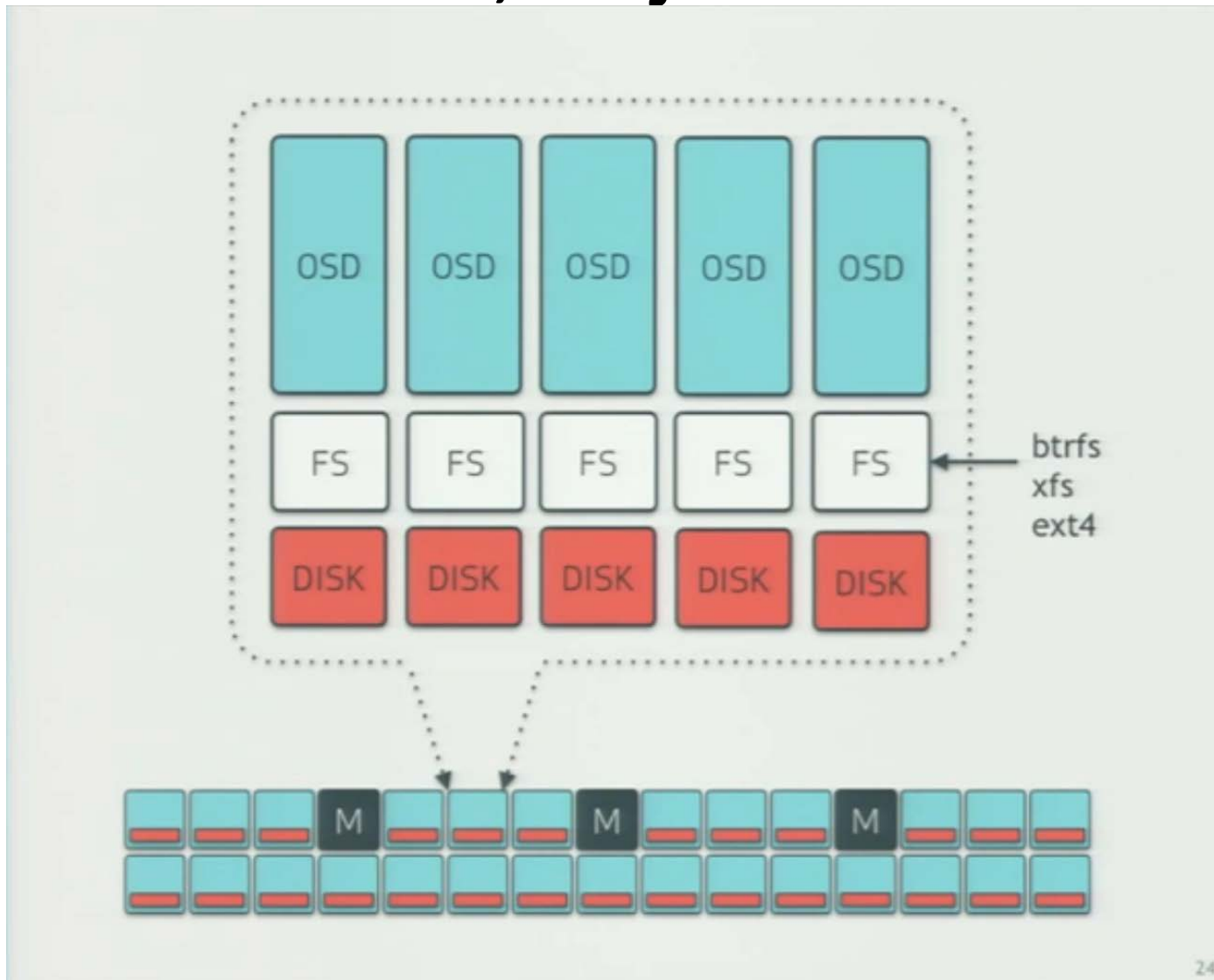
Filozófia



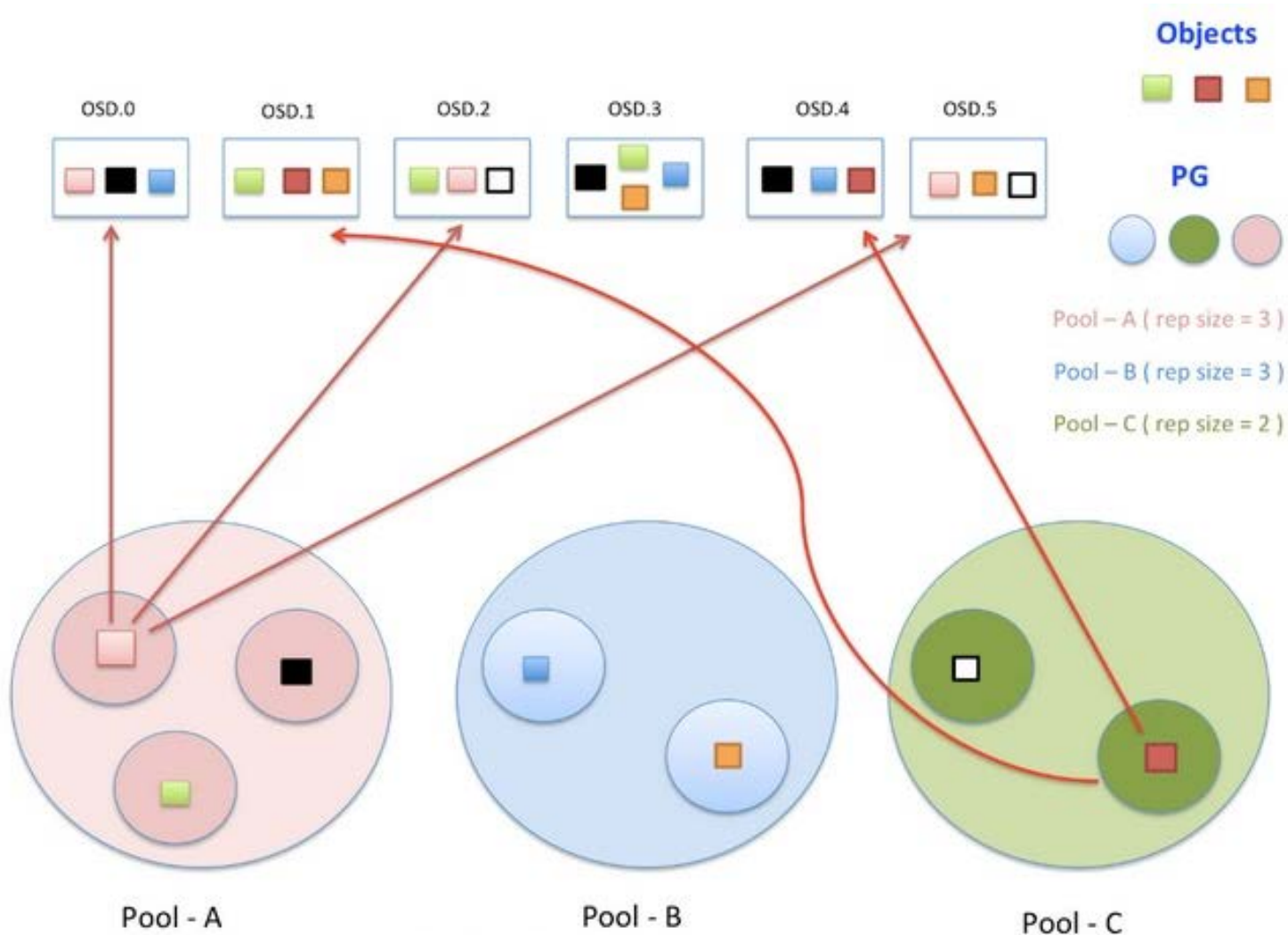
Hogyan, miként?



Adattárolás, objektumszintű



Hol az adat?



Copyright © Karan Singh , All rights reserved

Networkshop 2014

Felépítő kockák



Monitors:

- Maintain cluster map
- Provide consensus for distributed decision-making
- Must have an odd number
- These do **not** serve stored objects to clients



OSDs:

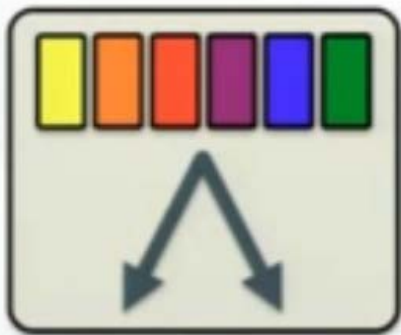
- One per disk (recommended)
- At least three in a cluster
- Serve stored objects to clients
- Intelligently peer to perform replication tasks
- Supports object classes



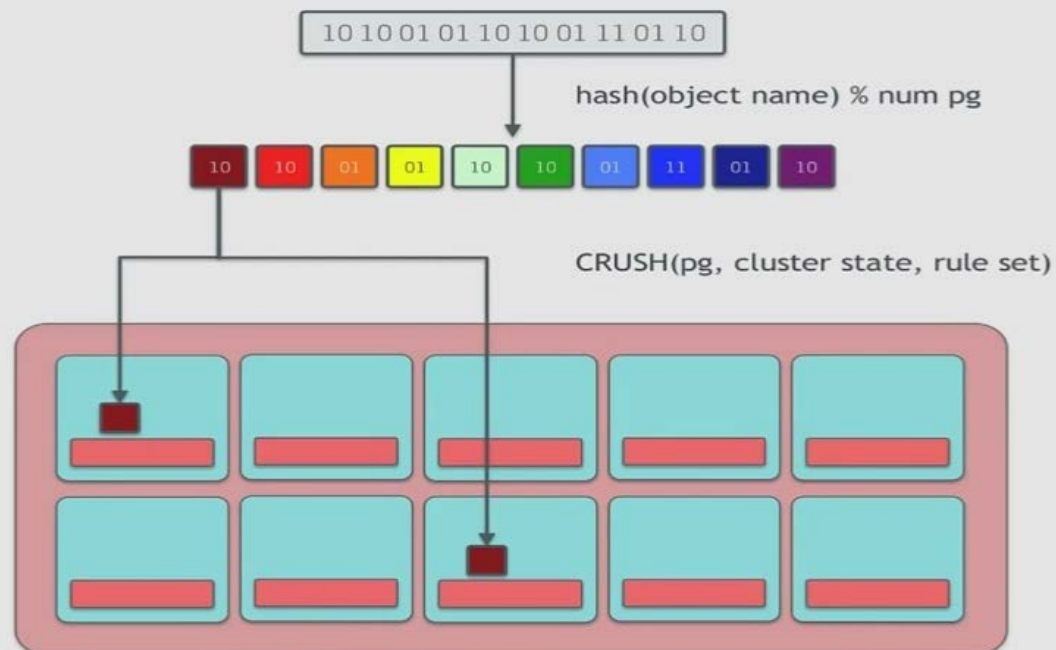
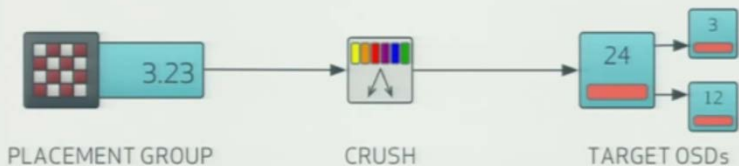
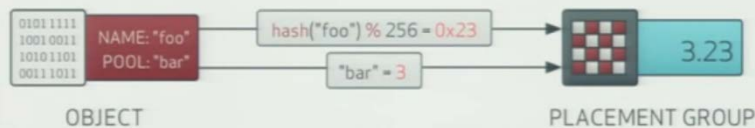
Metadata Server

- Manages metadata for a POSIX-compliant 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

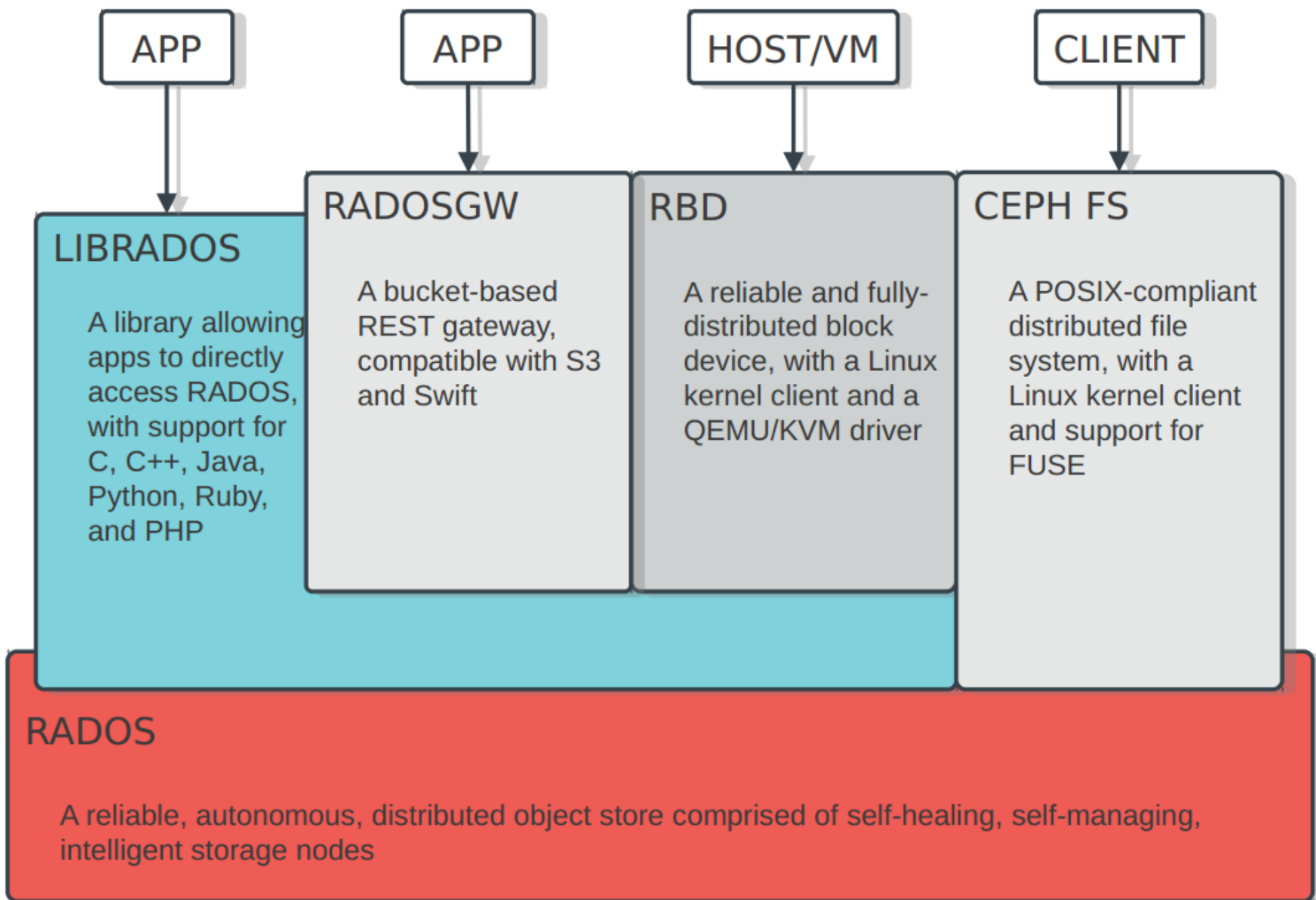
CRUSH algorithmus



- CRUSH
 - Pseudo-random placement algorithm
 - Ensures even distribution
 - Repeatable, deterministic
 - Rule-based configuration
 - Replica count
 - Infrastructure topology
 - Weighting



Tudáslista

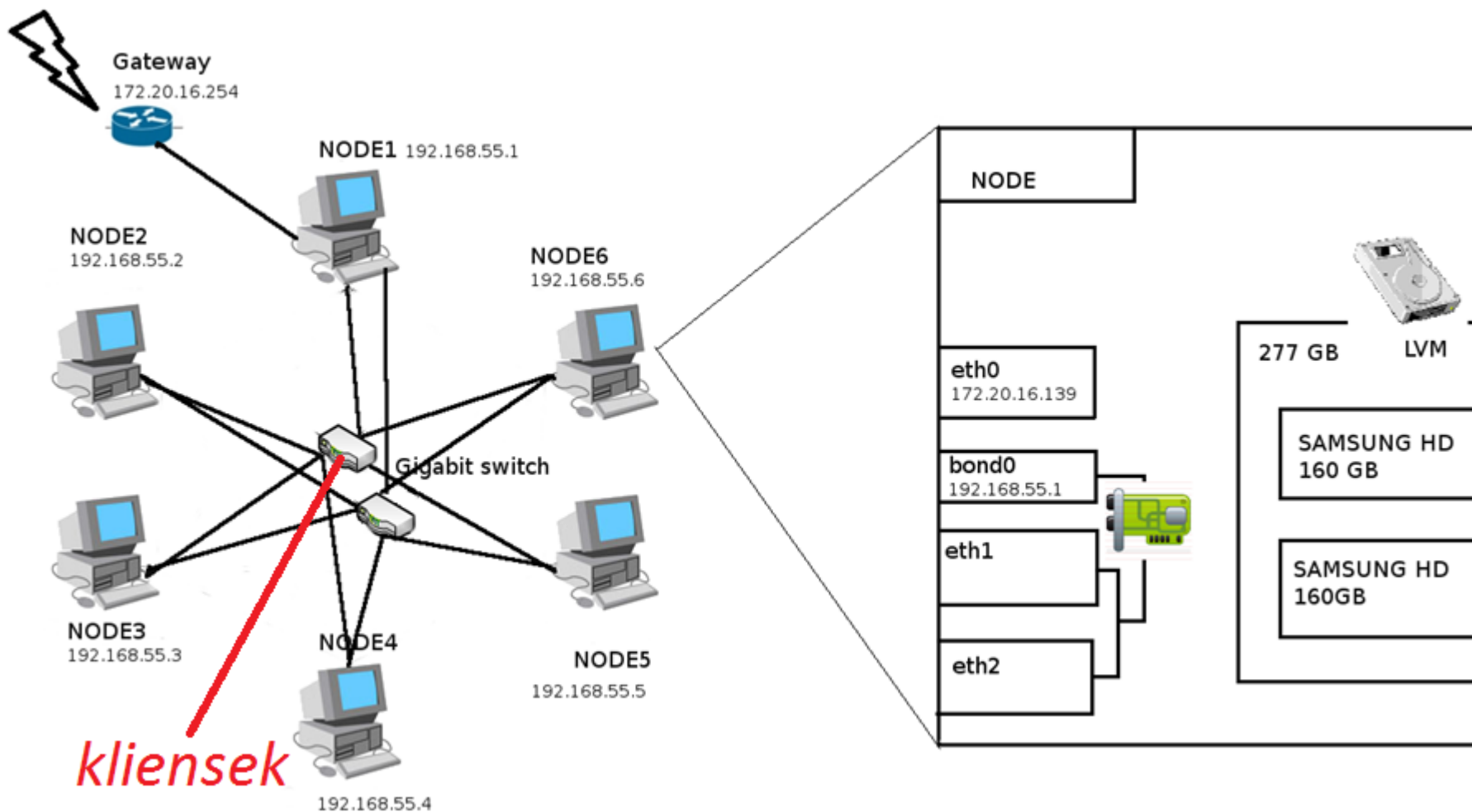


Hogy is kezdtük #1

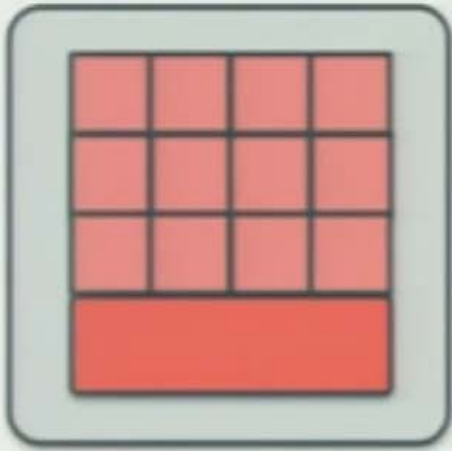
- Régebbi Core2Duo asztali gépek
- 2 GB RAM
- 2 x 160 GB SATA2 HDD
- 2 Gigabit NIC

- 2 db Gigabit Office Switch, CAT5 UTP kábelek

Hogy is kezdtük #2, topológia



Blokkeszköz!



RADOS Block Device:

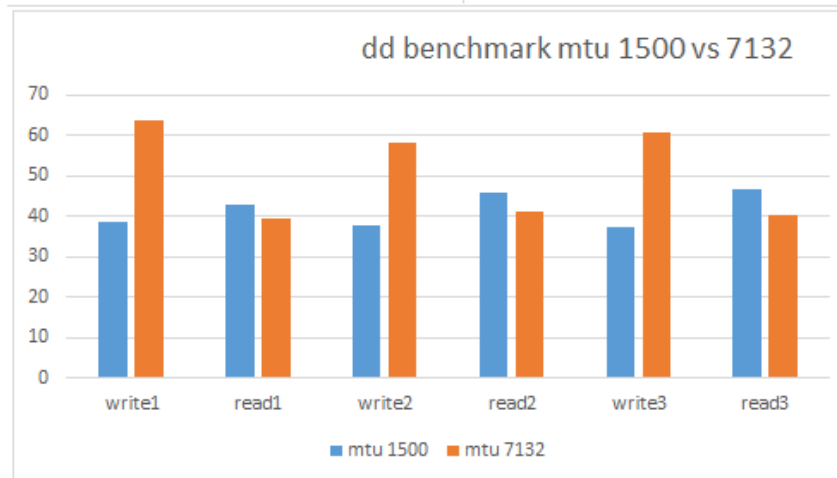
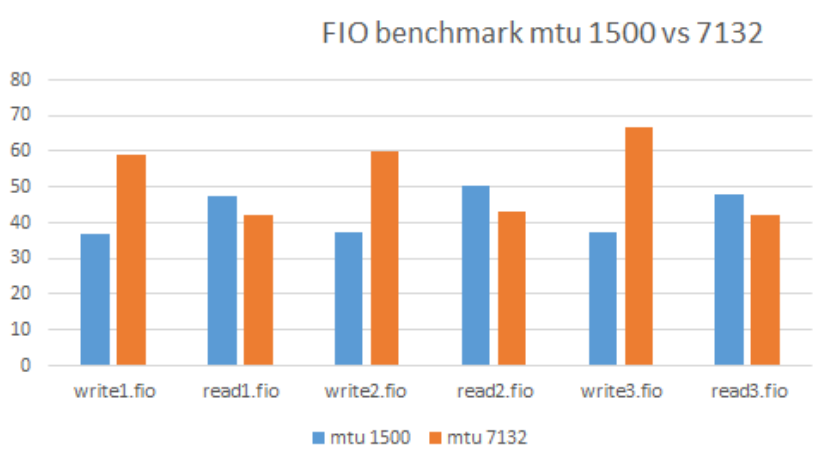
- Storage of disk images in RADOS
- Decouples VMs from host
- Images are striped across the cluster (pool)
- Snapshots
- Copy-on-write clones
- Support in:
 - Mainline Linux Kernel (2.6.39+)
 - Qemu/KVM, native Xen coming soon
 - OpenStack, CloudStack, Nebula, Proxmox

És filerendszer!






- Komplettn metadata tárolás
- POSIX „kompatibilitás”
- NFS szerű csatolás és működés
- Jól működő flock() és fcntl() kezelés
- Nem ajánlott egyelőre produktív üzembn



Sebességtesztek



Egyéb tesztek

- Failover
 - RBD 
 - CephFS 
- Cluster startup, recovery 
- Snapshot (brtfs) 
- Live migration 
- Performance tuning 