



Lesson 3 - Cloud infrastructure – storage and data repositories

Milan Brož

Software engineer

Storage engineering

Storage cost

OpenStack storage types

Software-defined storage concepts

Data persistence and redundancy

Virtualization

Distributed storage

Security

Q & A



Storage

- Capacity
- Availability, Reliability
- Data integrity, Redundancy
- Performance
- Scalability
- Security

=> Cost



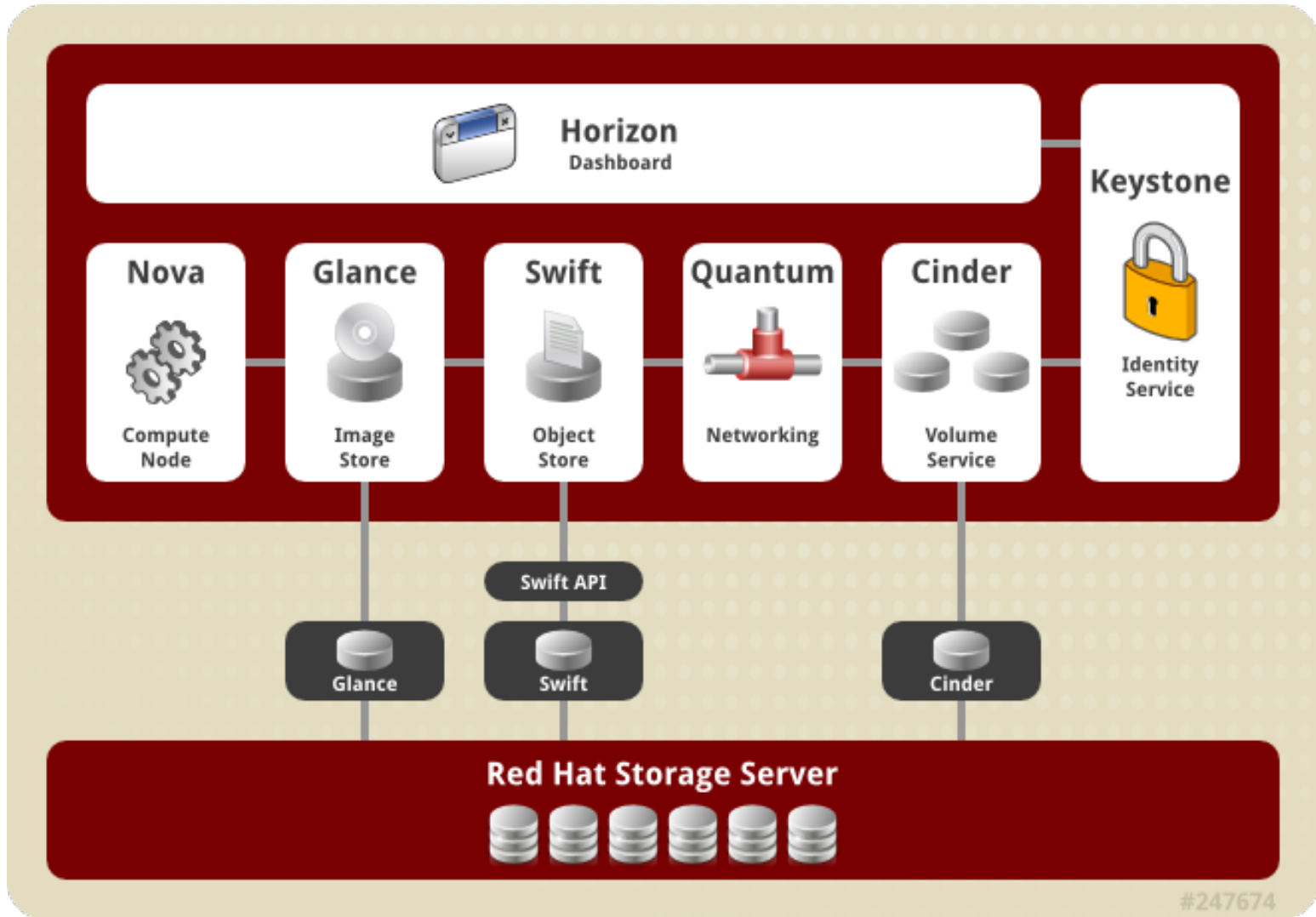
Manageability



Storage in OpenStack as an example



Persistent Storage – Example



Ephemeral storage

- Disappears when VM is terminated
- Temporary data ~ computing clusters
- Visible locally

Persistent storage

- Data always available (no dependency on instance)
- Can be shared among resources / instances



Object store

- binary objects of various length
- REST API

Block (volume) storage

- Block (sector-level) devices
- Can be backed by a file image

Shared file-system storage

- mounted to a directory



Object store = SWIFT

stateless swift-proxy

Block (volume) storage = CINDER

Backend Cinder drivers (LVM, GPFS, EMC, ...)

Shared File-system = MANILA

Backend Manila drivers (Ceph, GlusterFS, NFS, ...)

Image service = GLANCE

deduplication, clones, ...



Generic Storage Concepts



Software Defined Storage (SDS)

- "Commodity hardware with abstracted storage logic"
- Policy-based management of storage
- Virtualization
- Resource management
- Similar concept as Software Defined Network (SDN)
Note: distributed storage is mostly about networking!
- Thin provisioning, deduplication, replication, snapshots,
...

SDS definition differs among vendors!



Hardware and low-level storage protocols

- **Physical storage**
 - Rotational drives / hard disk drives (HDD)
 - Flash / SSD drives
 - Persistent Memory (byte-addressable!)
 - Tapes, magneto-optical drives, ...
- **Block-oriented storage access protocols**
 - "Small Computer System Interface" (SCSI), Serial Attached SCSI (SAS)
 - Serial ATA (SATA)
 - Fibre channel (FC) (not only fiber-optic)
 - InfiniBand (IB)



Storage connectivity through network

- **Direct-Attached Storage (DAS)**

- local, host-attached

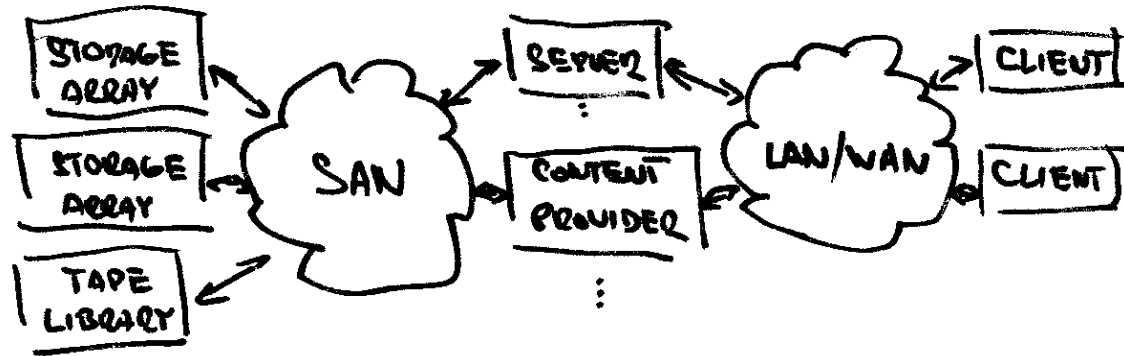
- **Network-Attached Storage (NAS)**

- remote storage device
- communication protocol
 - usually over IP-based network
 - high-level: NFS, CIFS, HTTP, ...
 - low-level: iSCSI (SCSI over IP), FC (point-to-point), Network Block Device (NBD)



Storage connectivity through network

- Storage Area Network (SAN)
 - private network
 - switched fabric
 - communication protocol
 - Fibre Channel
 - InfiniBand
 - FC over Ethernet (FCoE)



- **High availability (HA)** – assuring access to resources (data)
 - Service-level agreement (SLA)
 - common 9s levels

- **HA resources access**

- on-demand
- active/passive
- active/active

UPTIME (%)	DOWNTIME (%)	DOWNTIME PER YEAR	DOWNTIME PER WEEK
98	2	7.3 days	3 hr 22 minutes
99	1	3.65 days	1 hr 41 minutes
99.8	0.2	17 hr 31 minutes	20 minutes 10 sec
99.9	0.1	8 hr 45 minutes	10 minutes 5 sec
99.99	0.01	52.5 minutes	1 minute
99.999	0.001	5.25 minutes	6 sec
99.9999	0.0001	31.5 sec	0.6 sec



Generic Storage Concepts Data Protection and Redundancy



- **Data integrity protection**
 - random error detection (parity) / correction
- **Erasur codes** – Forward Error Correction (FEC)
 - Redundancy
 - RAID (Redundant Array of Independent Disks)
 - Erasure coding in distributed storage
- **Backup and disaster recovery**
 - **"RAID is not a backup!"**
 - File corruption, bugs (disk, controller, OS, application, ...)
 - Admin error, malware
 - Catastrophic failure (datacentre fire)
 - Offline and off-site backup replica



Common non-RAID and RAID disk configurations

- **JBOD** – "Just a Bunch of Disks" (collection of disks, no redundancy)
- **RAID-0** – striping (for performance, no redundancy, no parity)
- **RAID-1** – mirroring (no parity)
- **RAID-5** – block-level striping + distributed parity (XOR)
- **RAID-6** – block-level striping + double distributed parity
- **RAID-10** – nested RAID example (1+0: striping over mirrored drives)
- **RAIDZ** (in ZFS) – similar to RAID-5, dynamic stripes, self-healing
- **MAID** (Massive Array of Idle Disks) – "Write once, read occasionally"
- ...
- **Degraded mode**
 - RAID-5 (RAID-6 soon): large drives reconstruction time, fail during rebuild
- **Hardware RAID vs software RAID vs "fake RAID"** (processing in fw/driver)

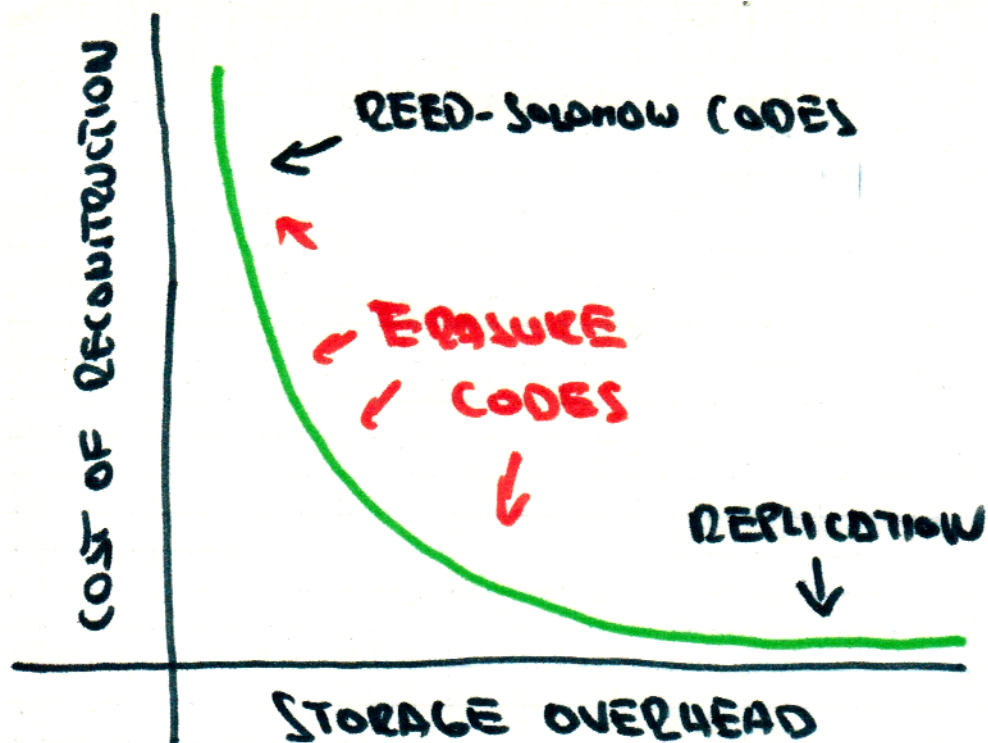


- **Data protection is trade-off**
 - Storage overhead
 - Reconstruction cost
 - Reliability
- Still active research ...
 - From simple XOR (RAID) to Galois Field arithmetic – $GF(2^x)$
 - Reed-Solomon codes, Pyramid codes
 - Bit-Matrix codes
 - ...



Erasure coding – Data protection

Erasure codes trade-off and efficient solution



Generic Storage Concepts Virtualization



Storage pool

set of disks, blocks, ... allocatable area for data

- **Pre-allocated**

- partition table, logical volume in Logical Volume Manager

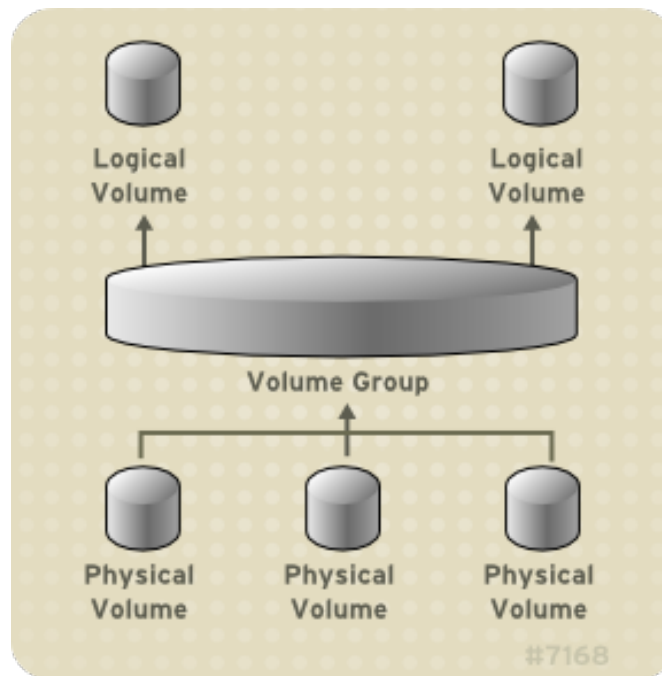
- **On-demand allocated**

- **Thin provisioning** (only blocks in use are allocated)
- Flexible allocation
- Used in snapshots
- Possible over-allocation (sharing "unallocated" space)

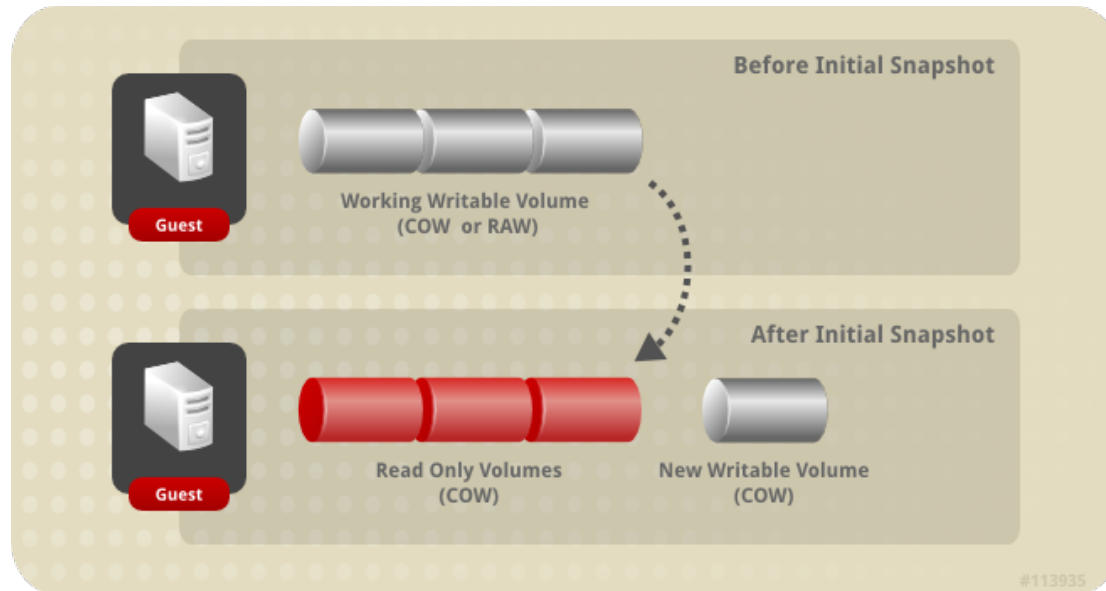


Storage Pool

- Volume Group: Linux Logical Volume Manager



- **Snapshot of storage** in specific time
 - Allows quick revert to older state (recovery)
- **Copy on Write (COW)** principle
 - delayed copy to snapshot (before origin write)
 - write to origin => need to copy the changed block first



Template

- Application of deduplication + snapshots (+ thin provisioning)
- **Virtual machine template**
 - base operating system
 - common configuration (networking, firewall, ...)
 - common applications (webservers, user packages, ...)
- One base image, only changes are stored

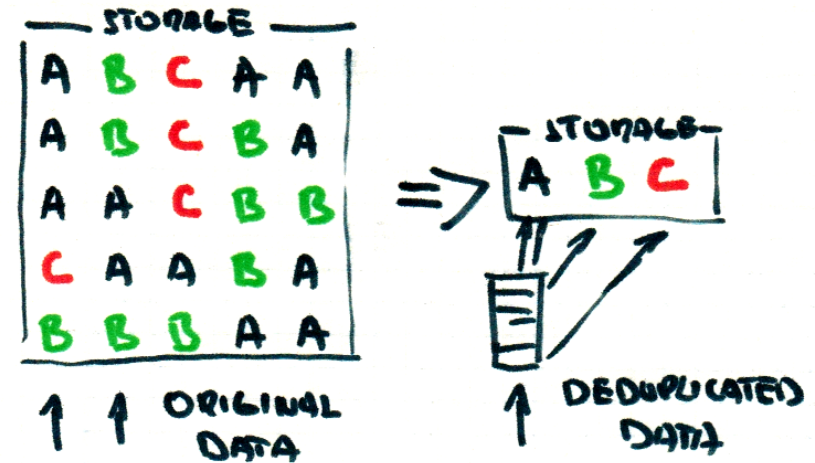
- Application containers + template
 - used in Docker



Deduplication / Compression

• Deduplication

- avoid to store repeated data
- file or block level
- space-efficient, stateless mode
- deduplication performance
- data corruption amplification



• Compression

- more generic algorithms
- special case: zeroed blocks



- **Tiered storage**

- Several layers of storage in one chain
- Different performance, availability, recovery requirements
- Cache (REST API)

- **Virtualization of drivers**

- virtio, pass-through device



Generic Storage Concepts

Distributed Storage



Clustered => cooperating nodes

Distributed => storage + network

Distributed storage transparency

- Access (same as local)
- Location (any node)
- Failure (self-healing)



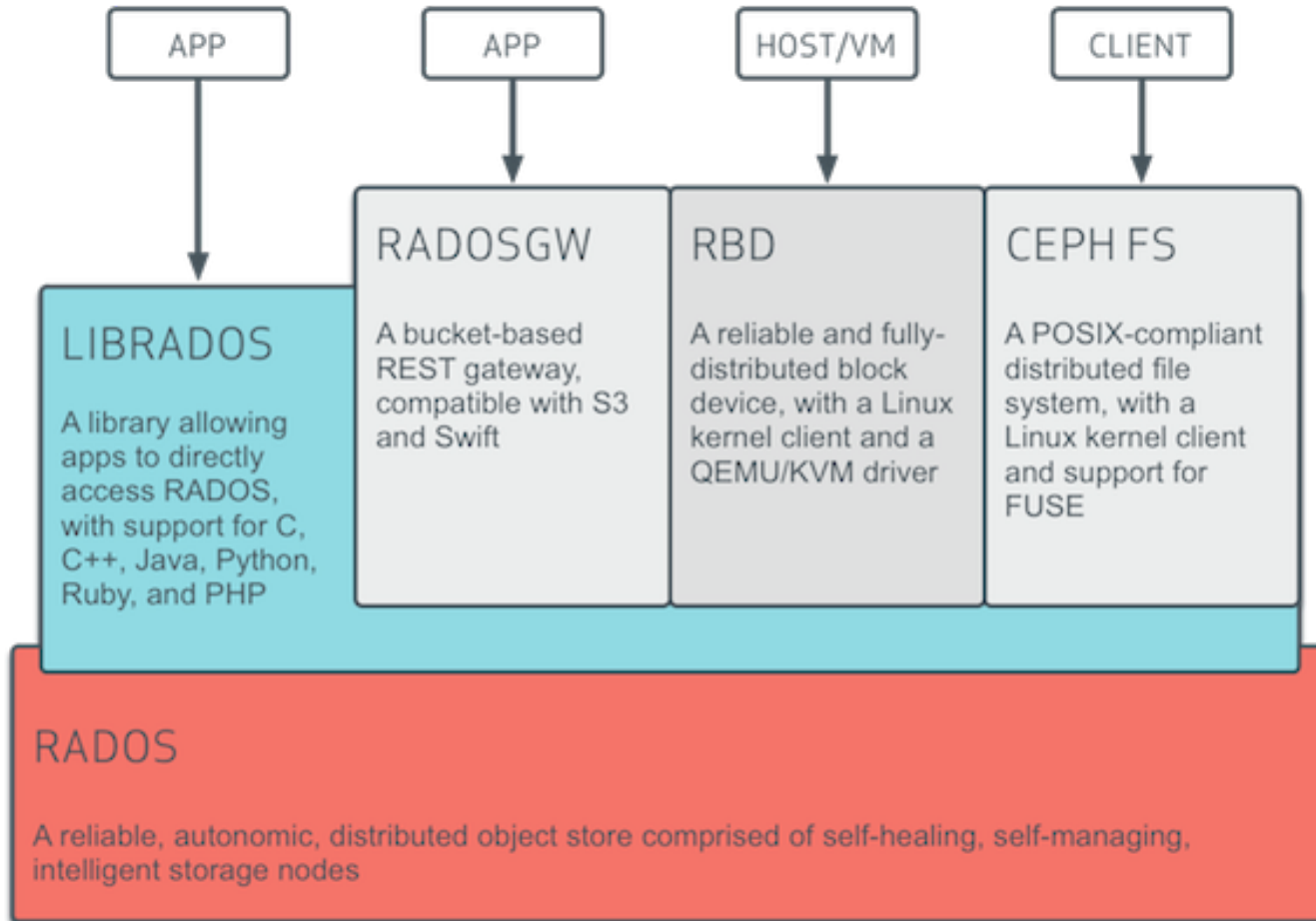
Distributed storage / file-systems examples

- Ceph, GlusterFS (Red Hat)
- General Parallel File System – GPFS (IBM)
- Hadoop File-System HDFS (Apache)
- Windows Distributed File-System (Microsoft)
- GoogleFS / GFS (Google)
- Isilon (EMC²)

...

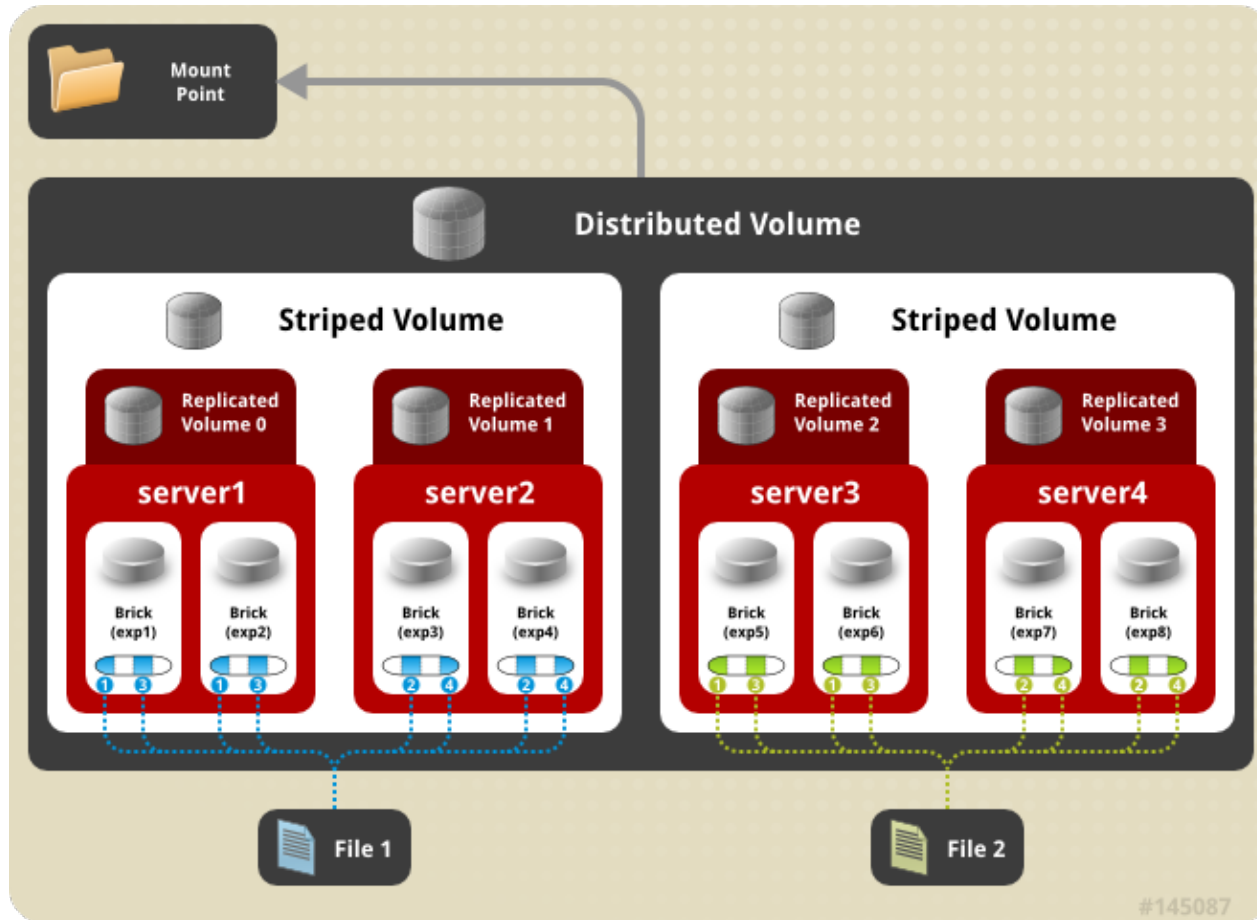


CEPH – Distributed storage



GlusterFS – Distributed storage

Example of access of **GlusterFS** resources



Generic Storage Concepts Security



- **Security policies**
- **Confidentiality**
 - Storage encryption (at-rest)
 - Data connection encryption (in-transit)
 - Key management
- **Authentication**
- **Integrity** (in cryptography sense – authenticated encryption)
- **Access control, permissions**
- **Secure data disposal / destruction**
- **Audit**

...



Encryption on client side

- "End-to-End" encryption
- Lost Efficiency for deduplication/compression

Encryption on server side

- Partially lost confidentiality for clients
(server has access to decrypted data)

Data at-rest – combination of ...

- Full disk encryption
- Filesystem encryption
- Object store encryption



