

THE LUSTRE FILE SYSTEM AND ITS IMPLEMENTATION AT LNGS

WHAT WE WILL TALK ABOUT

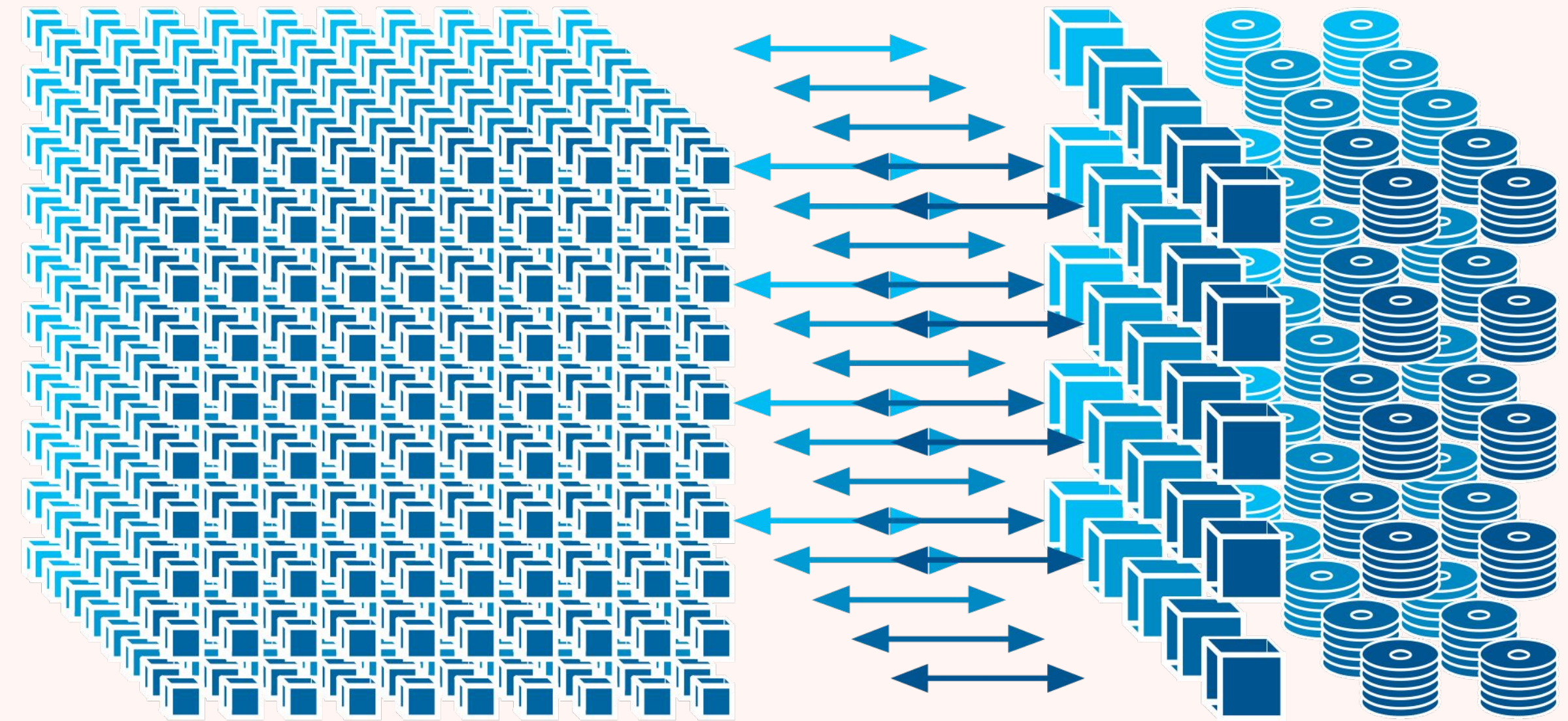
- What is Lustre?
- Lustre Components
- Lustre Architecture
- Why Lustre
- Best Practices
- Lustre implementation @ LNGS
- Future developments

·l·u·s·t·r·e[®]

WHAT IS IT

Fast, Scalable Storage for HPC

Lustre is an open-source, object-based, distributed, parallel, clustered file system



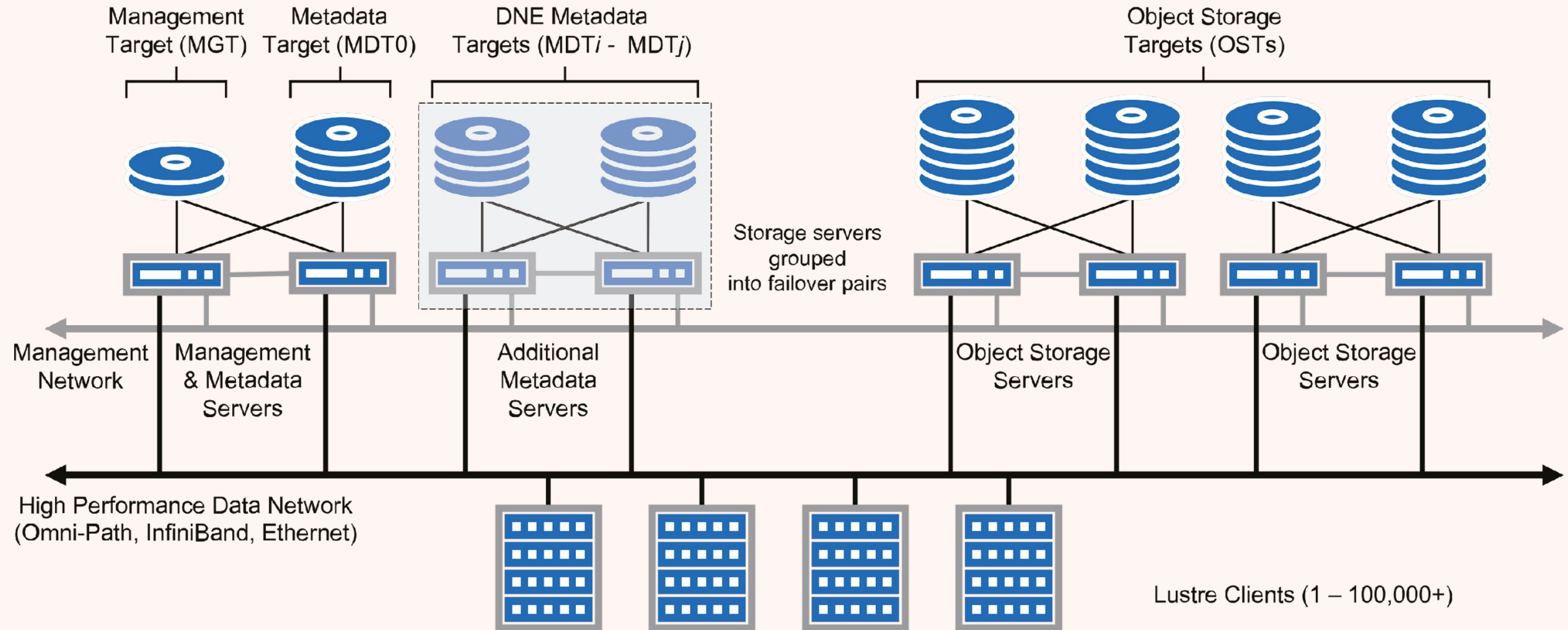
- **Lustre runs on Linux-based operating systems** and employs a client-server network architecture.
- Lustre software services are implemented entirely within the Linux kernel, as loadable modules.
- Storage is provided by a set of servers that **can scale to populations measuring up to several hundred hosts**.
- Lustre servers for a single file system instance can, in aggregate, present up to tens of petabytes of storage to thousands of compute clients simultaneously, and deliver more than a terabyte-per-second of combined throughput.

- Popular with the HPC community,
- Lustre is used to support the most demanding data-intensive applications.
- Lustre provides horizontally scalable IO for data centres of all sizes, and is deployed alongside some of the very largest supercomputers.
- The majority of the top 100 fastest computers, use Lustre for their high performance, scalable storage.
- Lustre file systems can scale from very small platforms of a few hundred terabytes up to large scale platforms with hundreds of petabytes, in a single, POSIX-compliant, name space.
- Capacity and throughput performance scale easily



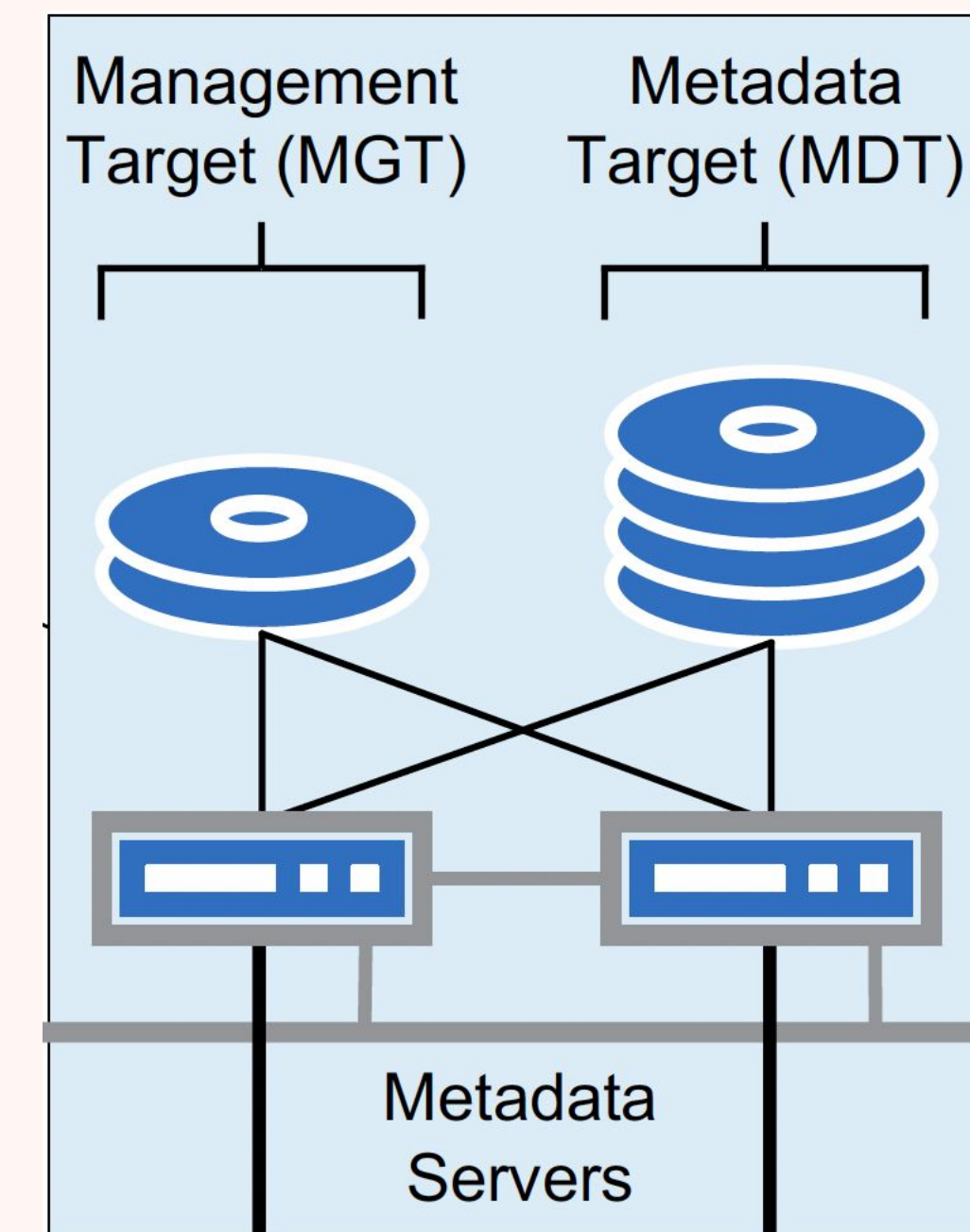
COMPONENTS AND ARCHITECTURE

Lustre File System Architecture – Server Overview



MGS + MGT: Management service, provides a registry of all active Lustre servers and clients, and stores Lustre configuration information.

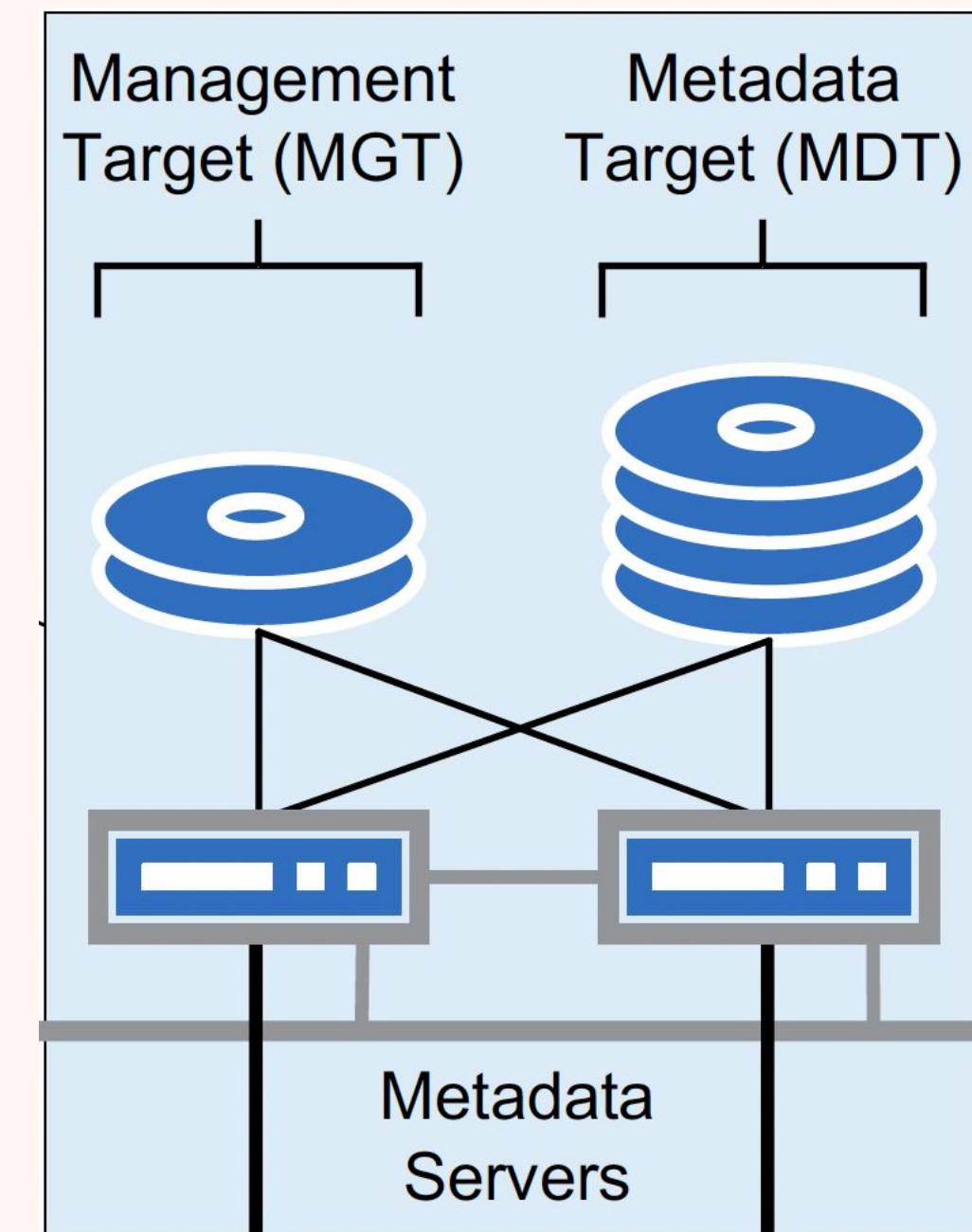
MGT is the management service storage target used to store configuration data.



MDS + MDTs: Metadata service, provides file system namespace (the file system index), storing the inodes for a file system.

MDT is the metadata storage target, the storage device used to hold metadata information persistently. Typically hosted on fast disks (SSD/NVMe)

Multiple MDS and MDTs can be added to provide metadata scaling.



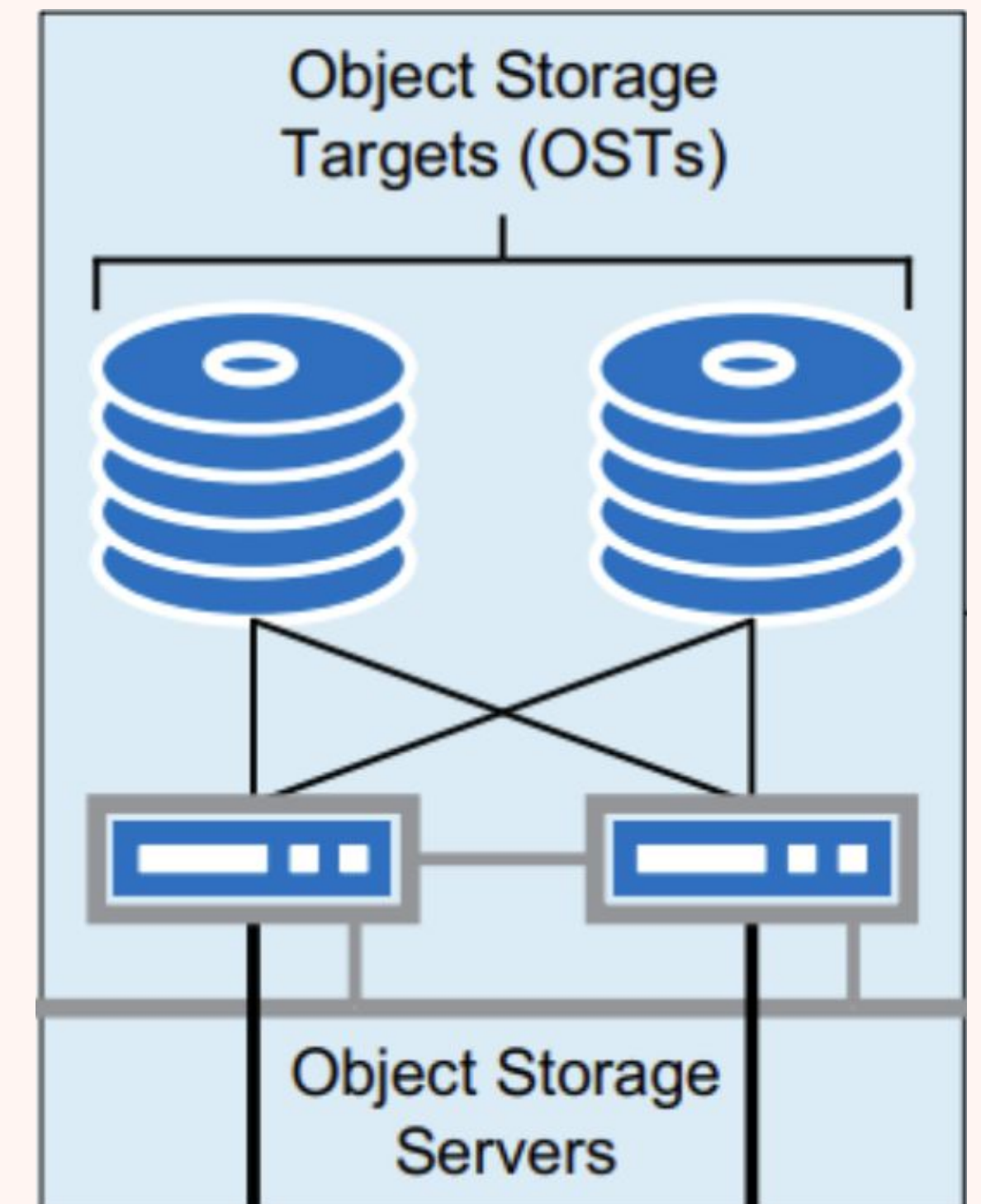
- Lustre separates metadata (inode) storage from block data storage (file content).
- All file metadata operations (creating and deleting files, allocating data objects, managing permissions) are managed by the metadata servers.
- Metadata servers provide the index to the file system.
- Metadata is stored in key-value index objects that store file system inodes: file and directory names, permissions, block data locations, extended attributes, etc.

OSS + OSTs:

Object Storage service, provides bulk storage of data. Typically hosted on capacitive disks (HDD)

Files can be written in stripes across multiple object storage targets (OSTs). Striping delivers scalable performance and capacity for files.

OSS are the primary scalable service unit that determines overall aggregate throughput and capacity of the file system.



- Lustre object storage servers write the data content out to persistent storage.
- Object servers can be written to concurrently, and individual files can be distributed across multiple objects on multiple servers.
- This allows very large files to be created and accessed in parallel by processes distributed across a network of computer infrastructure.
- Block data is stored in binary byte array data objects on a set of storage targets.
- A single Lustre "file" can be written to multiple objects across many storage targets.

Lustre filesystem backend

Lustre servers have a choice of backend storage target formats, either **LDISKFS** (derived from EXT4), or **ZFS**

Lustre servers using **LDISKFS** storage require patches to the Linux kernel. These patches are to improve performance, or to enable instrumentation useful during the automated test processes in Lustre's software development lifecycle. The list of patches continues to reduce as kernel development advances, and there are initiatives underway to completely remove customized patching of the Linux kernel for Lustre servers.

Lustre servers using **ZFS OSD** storage and Lustre clients do not require patched kernels.

The Lustre client software is being merged into mainstream Linux kernel, and is available in kernel-staging.

Clients:

- Lustre clients mount each Lustre file system instance using the Lustre Network protocol (LNet). They present a POSIX-compliant file system to the OS.
- Applications use standard POSIX system calls for Lustre IO, and do not need to be written specifically for Lustre.
- Clients do not access storage directly: all I/O is sent over a network.
- Clients aggregate the metadata name space and object data to present a coherent POSIX file system to applications. The Lustre client software splits IO operations into metadata and block data, communicating with the appropriate services to service IO transactions. This is the **key concept of Lustre's design** – **separate small, random, IOPS-intensive metadata traffic from the large, throughput-intensive, streaming block IO.**

Lustre Networking

- Lustre is a network-based file system, all IO transactions are sent using network RPCs using the **LNet** protocol.
- Clients have no local persistent storage and are sometimes diskless.
- Supports many different network technologies, including OPA, IB, Ethernet.

Lustre Networking

LNet can aggregate IO across independent interfaces, enabling network multipathing. It has native support for TCP/IP networks as well as RDMA networks such as Intel Omni-Path Architecture (OPA) and InfiniBand.

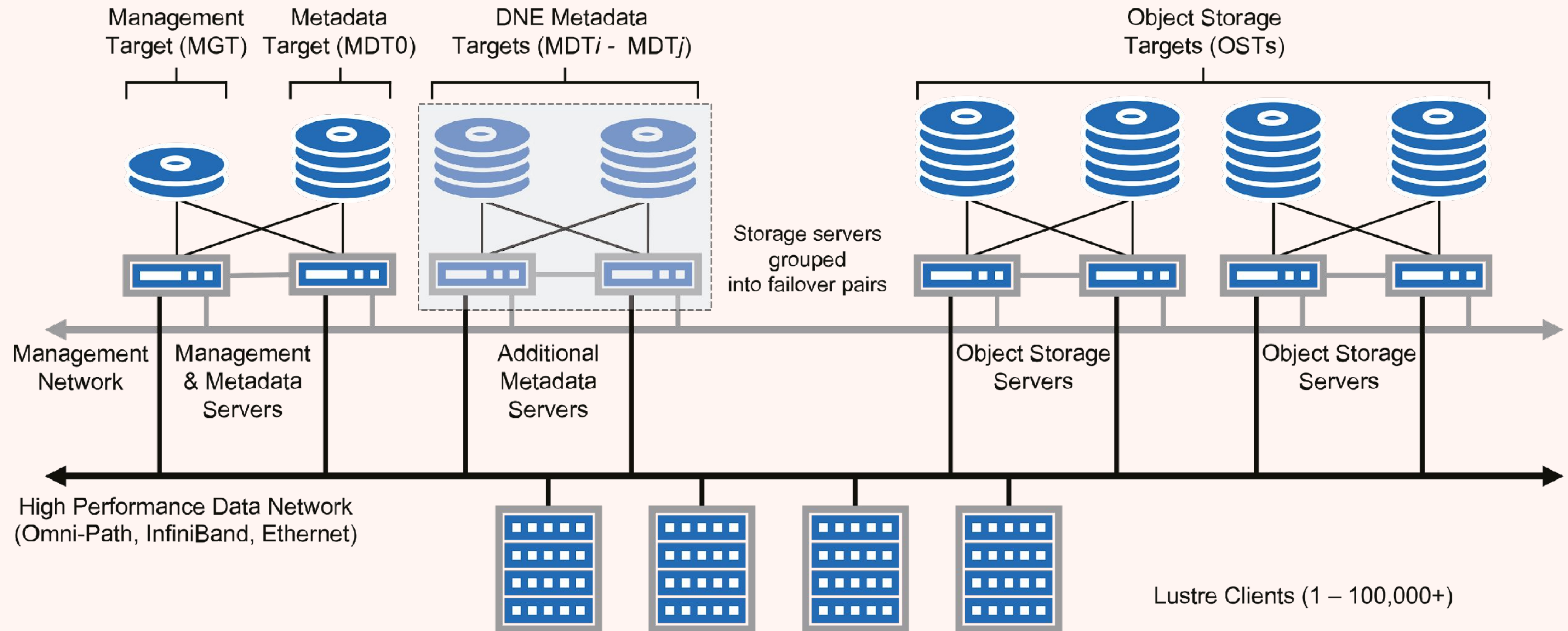
Servers and clients can be multi-homed, and traffic can be routed using dedicated machines called LNet routers.

Lustre network (LNet) routers provide a gateway between different LNet networks. Multiple routers can be grouped into pools to provide performance scalability and to provide multiple routes for availability.

Lustre Networking

- Most Lustre protocol actions are initiated by clients, the most common activity in the Lustre protocol is for a client to initiate an RPC to a specific target.
- A server may also initiate an RPC to the target on another server, e.g. an MDS RPC to the MGS for configuration data; or an RPC from MDS to an OST to update the MDS's state with available space data.
- Object storage servers never communicate with other object storage servers, all coordination is managed via the MDS or MGS.
- OSS do not initiate connections to clients or to an MDS.
- An OSS is relatively passive: it waits for incoming requests from either an MDS or Lustre clients.

Lustre File System Architecture – Server Overview



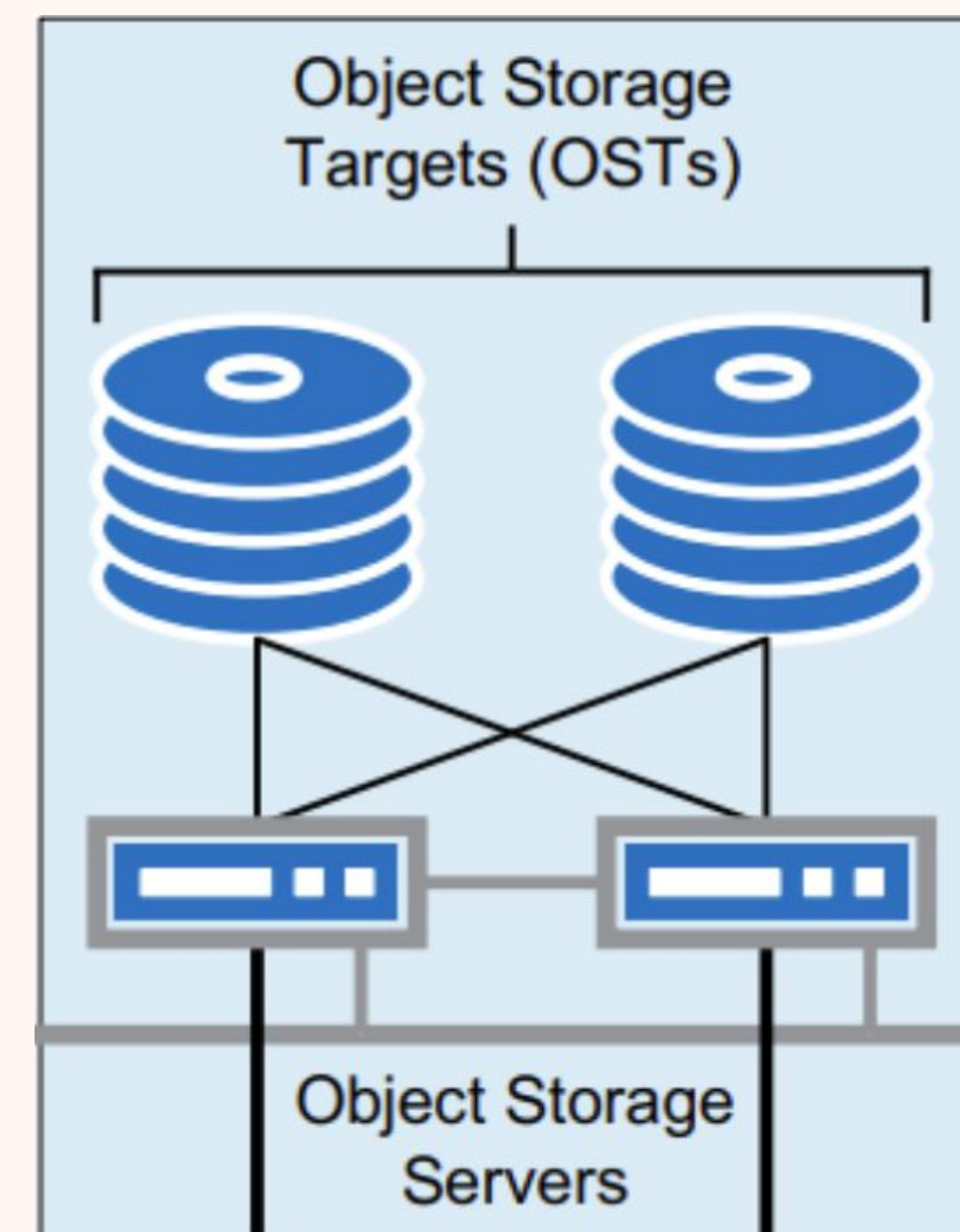
Lustre and High Availability

Service availability / continuity is sustained using a High Availability failover resource management model, where multiple servers are connected to shared storage subsystems and services are distributed across the server nodes.

Individual storage targets are managed as active-passive failover resources, and multiple resources can run in the same HA configuration for optimal utilisation. If a server develops a fault, then any Lustre storage target managed by the failed server can be transferred to a surviving server that is connected to the same storage array.

Failover is completely application-transparent: system calls are guaranteed to complete across failover events.

All Lustre server types (MGS, MDS and OSS) support failover.



Lustre and High Availability

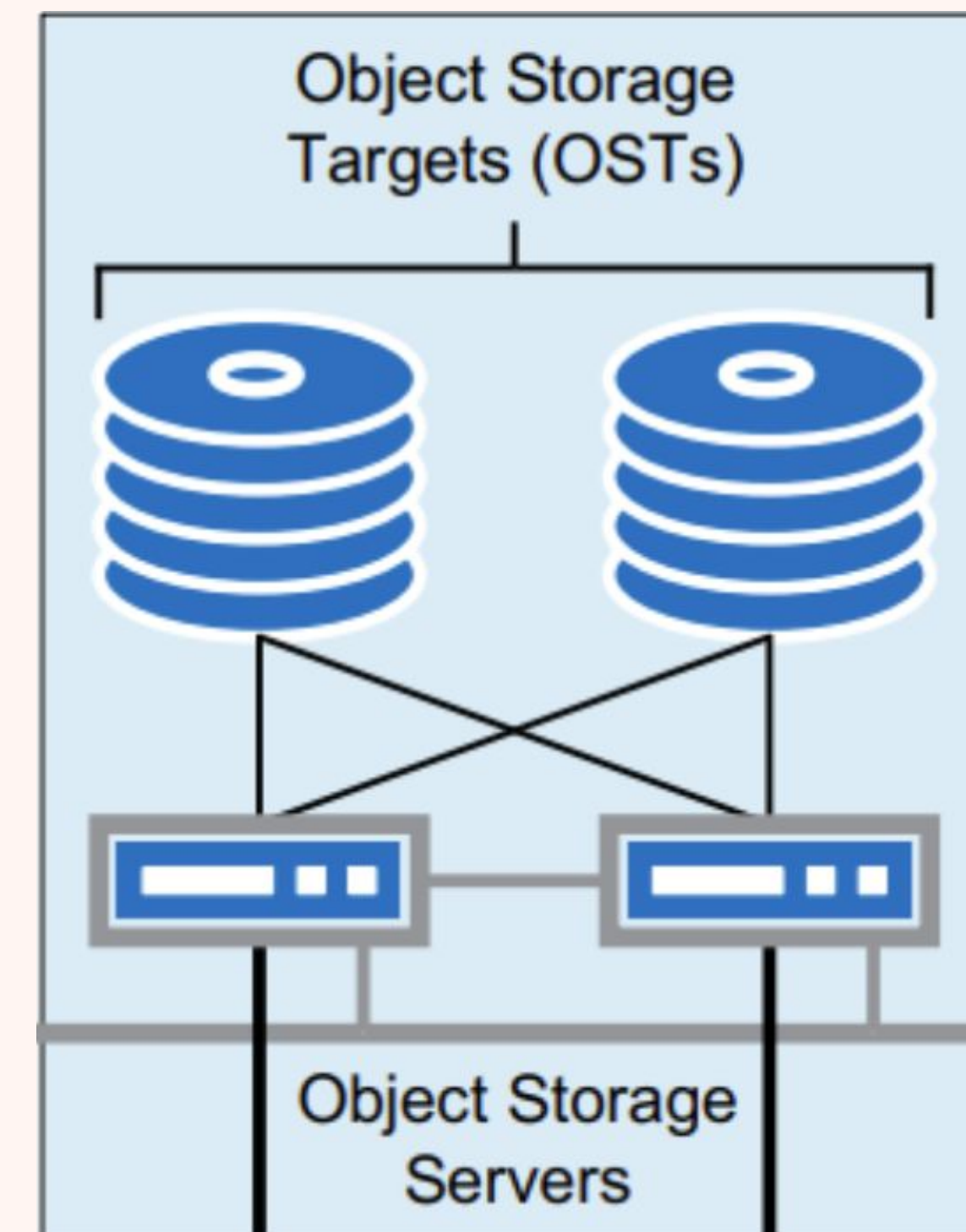
The most common blueprint employs two servers joined to shared storage in an HA clustered pair topology.

While HA clusters can vary in the number of servers, a two-node configuration provides the greatest overall flexibility as it represents the smallest storage building block that also provides high availability.

A single Lustre file system can scale linearly based on the number of building blocks.

The minimum HA configuration for Lustre is a metadata and management building block that provides the MDS and MGS services, plus a single object storage building block for the OSS services.

Using these basic units, one can create file systems with hundreds of OSSs as well as several MDSs, using HA building blocks to provide a reliable, high-performance platform.



Lustre Storage Scalability

	Value using LDISKFS backend	Value using ZFS backend	Notes
Maximum stripe count	2000	2000	Limit is 160 for ldiskfs if "ea_inode" feature is not enabled on MDT
Maximum stripe size	< 4GB	< 4GB	
Minimum stripe size	64KB	64KB	
Maximum object size	16TB	256TB	
Maximum file size	31.25PB	512PB*	
Maximum file system size	512PB	8EB*	
Maximum number of files or subdirectories per directory	10M for 48-byte filenames. 5M for 128-byte filenames.	2 ⁴⁸	
Maximum number of files in the file system	4 billion per MDT	256 trillion per MDT	
Maximum filename length	255 bytes	255 bytes	
Maximum pathname length	4096 bytes	4096 bytes	Limited by Linux VFS

Lustre Storage Scalability

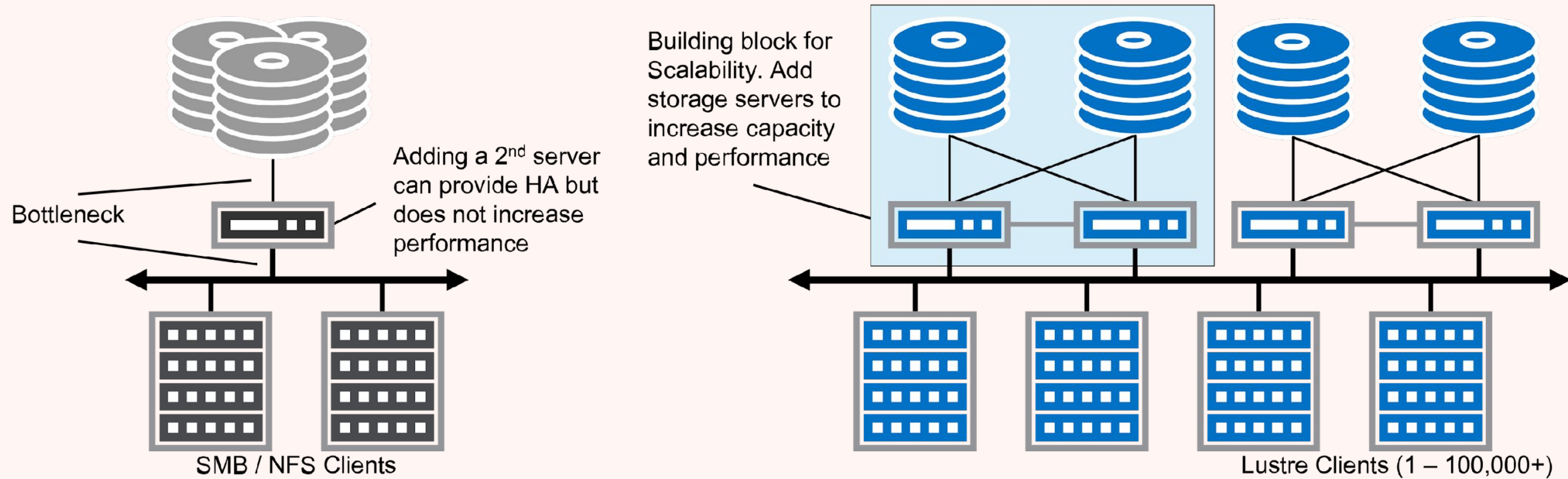
The values here are intended to reflect the potential scale of Lustre installations.

The largest Lustre installations used in production today are measured in 10's of petabytes, with projections for 100's of petabytes within 2 years.

Typical Lustre system deployments range in capacity from 1 – 60PB, in configurations that can vary widely, depending on workload requirements.

The values depicted for ZFS file system and file system size are theoretical, and in fact represent a conservative projection of the scalability of ZFS with Lustre.

Traditional Network File Systems vs Lustre



Distributed client, central server, e.g. NFS, SMB
 Performance and scalability is limited by single server bandwidth
 All clients use a single network path to one server
 Adding a new storage server creates a separate file system
 Storage becomes fragmented into silos

Distributed client, distributed servers (Lustre)
 Scalable Metadata and Object Storage services
 Scale performance and capacity linearly as client population grows
 Client network scalability is not limited by the file system
 All clients can access all storage servers simultaneously and in parallel
 Single coherent name space across all servers. Can have hundreds of network storage servers in a single file system

Where Lustre differentiates itself from other network file systems, such as NFS or SMB, is in its ability to seamlessly scale both capacity and performance linearly to meet the demands of data-intensive applications with minimal additional administrative overhead.

To increase capacity and throughput, add more servers with the required storage. Lustre will automatically incorporate new servers and storage, and the clients will leverage the new capacity automatically.

New capacity is automatically incorporated into the pool of available storage.

WHY

·l·u·s·t·r·e[®]·

High-performance

Deploy and enjoy from few MB/s to a few TB/s parallel file I/O or hundreds of thousands per second parallel metadata I/O.

Has been the filesystem of choice for the majority of the top 10 supercomputers for more than a decade

Proven in Mission-Critical Environments

Running in the most demanding environments - autonomous driving, genomics, oil & gas, etc.

Initially developed under the United States Department of Energy's Accelerated Strategic Computing Initiative Path Forward project.

Scalability

Capable of scaling as your requirements grow - hundreds of petabytes, billions of files, and tens of thousands of clients.

Commercial Lustre offerings available for all the major cloud providers, providing options for scaling beyond the traditional on premises model.

Flexibility

Supports a wide range of architectures, clients, networks.

From InfiniBand on the network to ARM on the CPU, Lustre has a history of embracing and incorporating new technologies.

Due to its modular design, it has the ability to leverage ZFS for JBOD installations, as well as hardware RAID scenarios.

Community-driven

Available under GPL v2 open source license and developed by a worldwide community under an open development model.

Worldwide organizations like OpenSFS and EOFS help to interface between Lustre users and vendors.



BEST PRACTICES

- Avoid Using `ls -l`
 - Use `ls` by itself if you just want to see if a file exists
 - Use `ls -l filename` if you want the long listing of a specific file
- Avoid Having a Large Number of Files in a Single Directory
- Avoid Repetitive “`stat`” Operations
- Avoid Repetitive Open/Close Operations



IMPLEMENTATION AT LNGS

- At LNGS we are building a new HPC data center with fundings coming from different sources.
- We chose Lustre as our main POSIX storage backend, while using CEPH for object and block storage and maybe as file storage for some use cases.
- We presently have a very basic Lustre setup at LNGS with two servers acting as OSS, MGS and MDS and two SAN boxes.
- In the following slides we will describe how they are configured

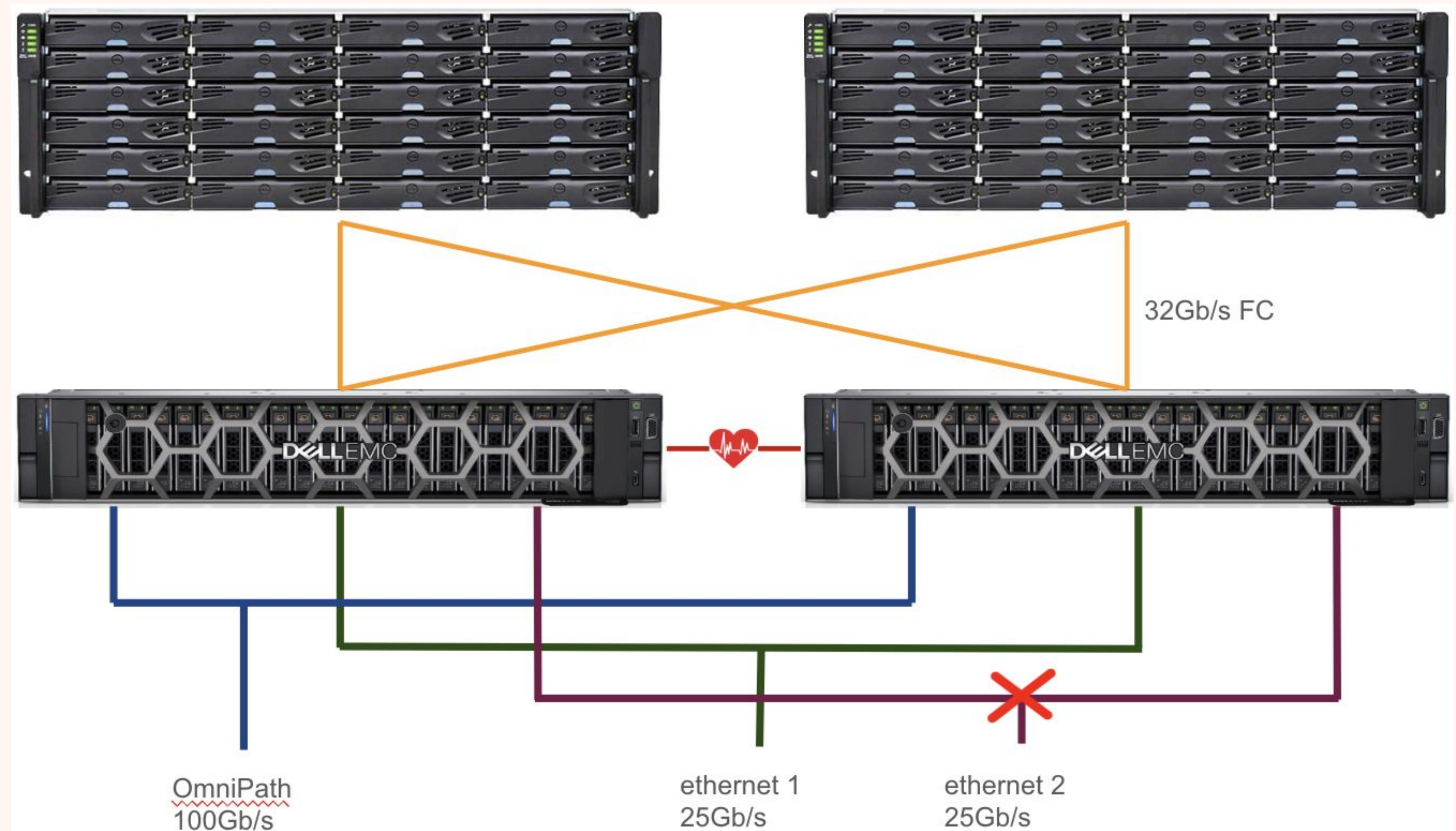


Hardware

This figure represents the physical layout the Lustre storage @LNGS as it is now.

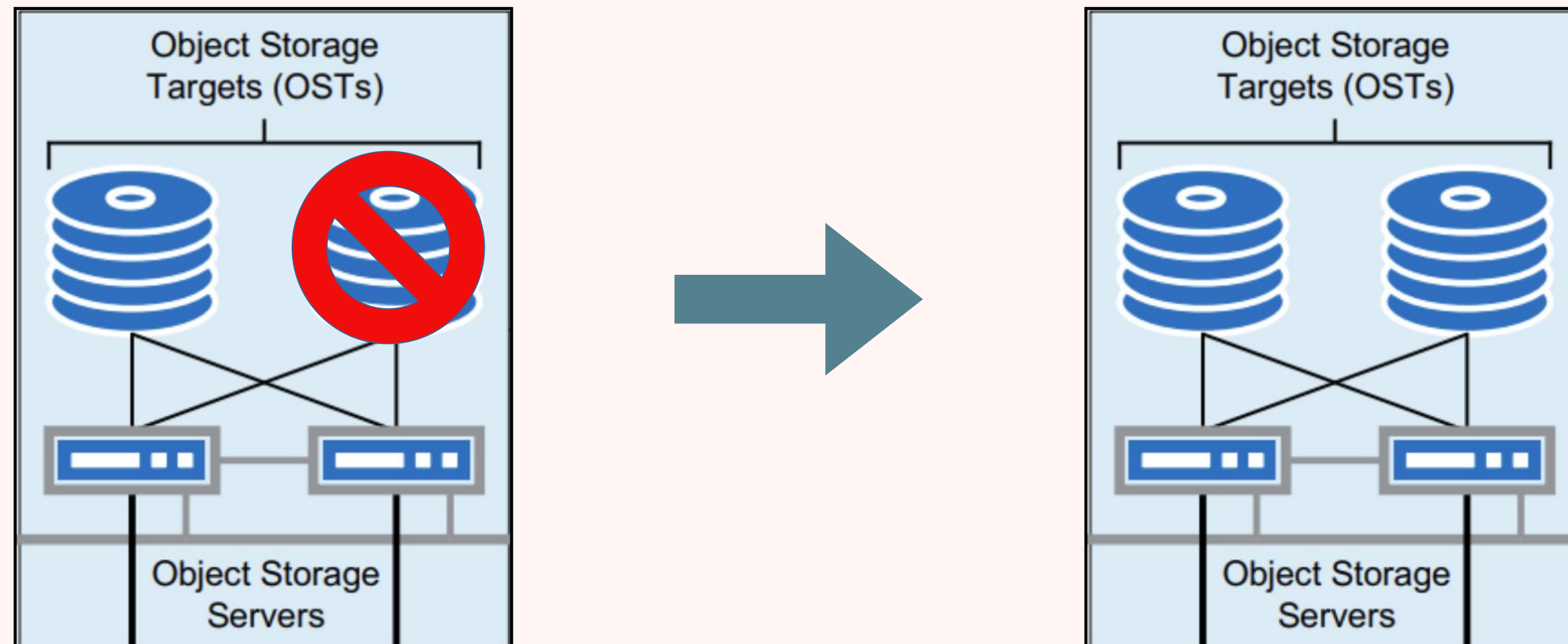
It is coherent with the reference architecture described before.

Such layout assumes that the SANs will not fail while the servers may fail.



Hardware

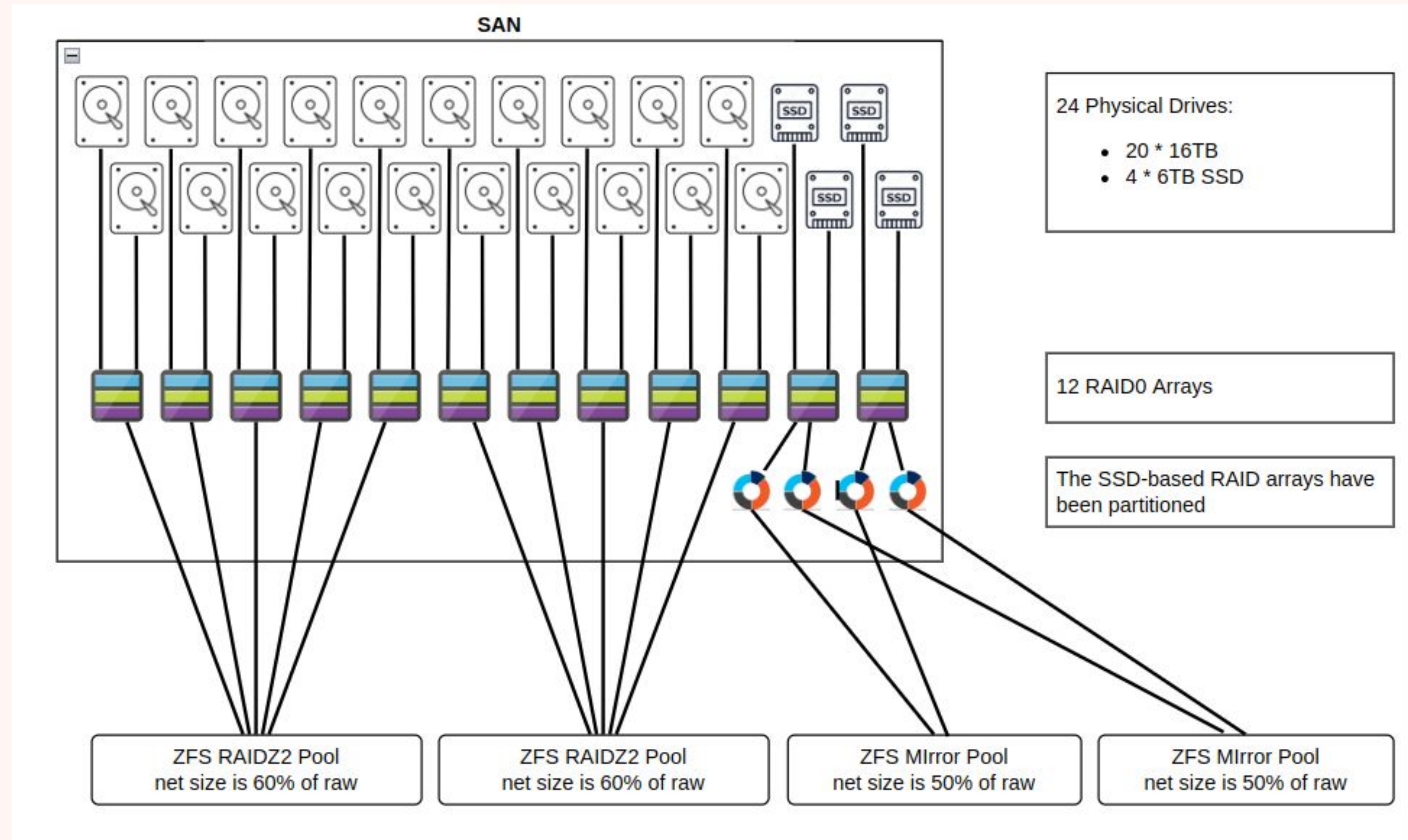
We actually started with only one SAN and later added a second one. This helped us understand better how lustre can scale horizontally in time



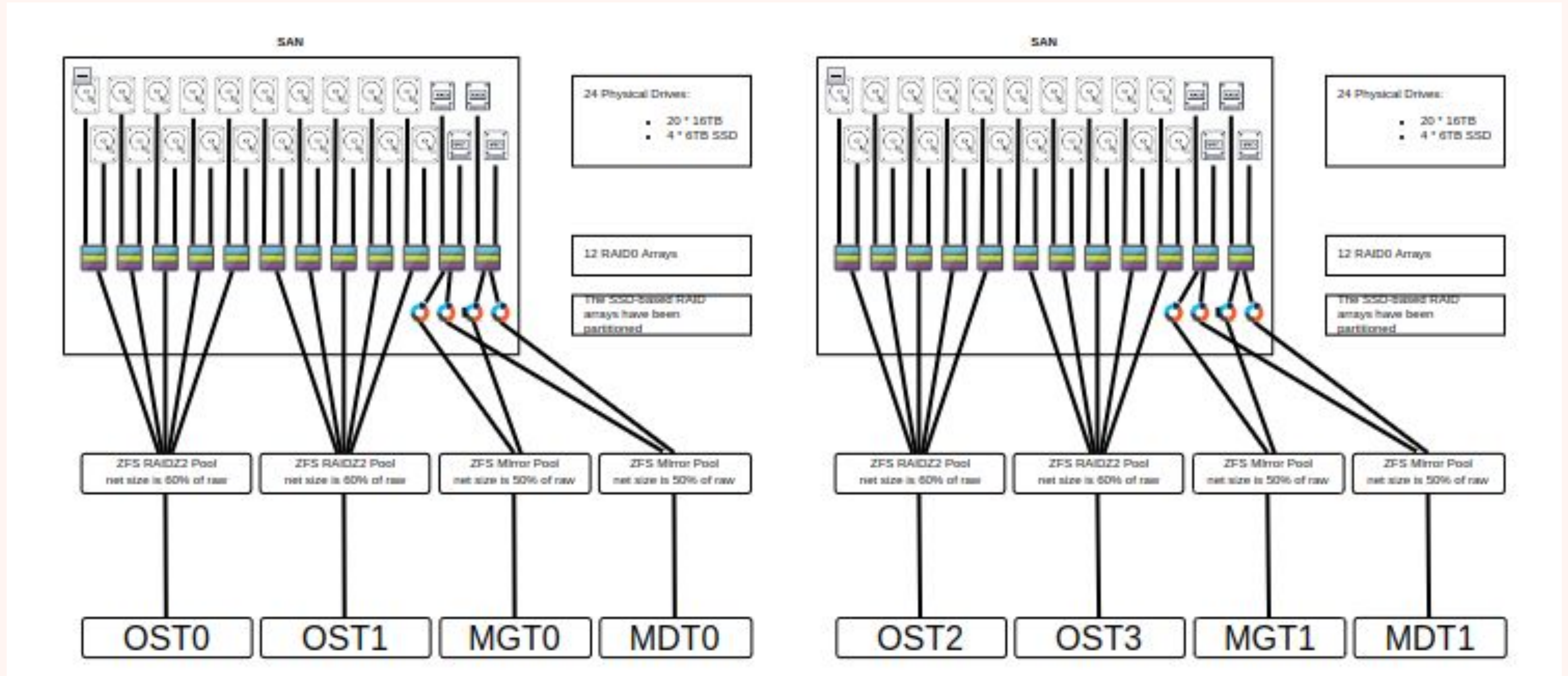
Disk Layout

This is how we configured our SANs

- 4* RAID0 on SSD as MGTs and MDTs
- 20*RAID0 on HDD as OSTs
- OSTs are grouped in RAIDZ2 ZFS pools (similar to RAID 6)
- MDTs and MGTs are grouped in ZFS Mirror pools
- LUSTRE “sees” 4 83TB OSTs



How physical disks map to Lustre storage services



Network Layout

The LNGS HPC cluster is based on an OmniPath high performance network. Some newer hosts presently only connect to the cluster via ethernet

- **OmniPath NIC**

This is used by most of our HPC cluster nodes that are Omni-Path capable

- **Ethernet NICs**

The first NIC is used by some HPC cluster nodes that are not Omni-Path capable.

We were not able to make the second ethernet connection work despite an apparently correct configuration. The second ethernet connection was meant to make lustre available to OpenStack VMs

Software

We use Lustre v 2.4 with kernel on rocky linux 8.6

The kernel needed to be rebuilt, not an easy job

We used an existing puppet module to setup the Lustre cluster, but we needed to adapt it to our needs

What does the lustre puppet module do

- it installs all needed packages on the lustre servers and clients
- on servers, it installs and configures the ZFS module
- on clients and server, it installs and configures the lustre module

```
package {[  
  'lustre',  
  'kmod-lustre-osd-zfs',  
  'lustre-osd-zfs-mount',  
  'kmod-lustre',  
  'lustre-resource-agents',  
]:}
```

```
yum install lustre kmod-lustre-osd-zfs  
lustre-osd-zfs-mount kmod-lustre  
lustre-resource-agents
```

What does the lustre puppet module do

- it configures the ZFS pools for MGTS, MDTs, OSTs

```
/usr/sbin/zpool create -f \
-o recordsize=1024k \
-o dnodesize=auto \
-o canmount=off \
-o multihost=on \
-o cachefile=none \
-o ashift=12 \
-o compression=lz4 \
hpc4dr01-ost0 \
raidz2 \
/dev/disk/by-id/scsi-3600d0231000c3f325473143f3c8a05d4 \
/dev/disk/by-id/scsi-3600d0231000c3f32176d676f3a195104 \
/dev/disk/by-id/scsi-3600d0231000c3f323a1607b96f639e66 \
/dev/disk/by-id/scsi-3600d0231000c3f326fdf34302cbc7781 \
/dev/disk/by-id/scsi-3600d0231000c3f322dcbbef21cb06b3b
```

```
[root@stor-2 ~]# zpool status hpc4dr01-ost0
pool: hpc4dr01-ost0
state: ONLINE
scan: scrub repaired 0B in 01:56:24 with 0 errors on Wed Jan 3 07:55:14 2024
config:

NAME                                STATE      READ WRITE CKSUM
hpc4dr01-ost0
  raidz2-0
    scsi-3600d0231000c3f325473143f3c8a05d4  ONLINE    0     0     0
    scsi-3600d0231000c3f32176d676f3a195104  ONLINE    0     0     0
    scsi-3600d0231000c3f323a1607b96f639e66  ONLINE    0     0     0
    scsi-3600d0231000c3f326fdf34302cbc7781  ONLINE    0     0     0
    scsi-3600d0231000c3f322dcbbef21cb06b3b  ONLINE    0     0     0

errors: No known data errors
```

What does the lustre puppet module do

- it formats the lustre filesystems
- it creates the lustre mount points
- it creates a cron job for a periodic scrubbing of the volumes

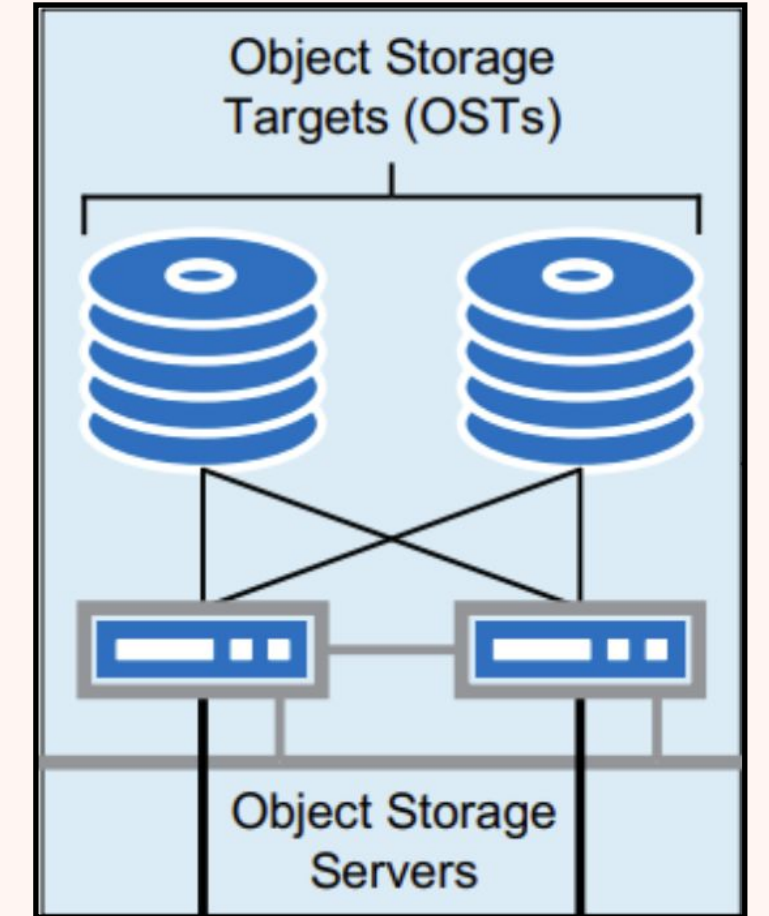
```
/usr/sbin/mkfs.lustre \  
--backfstype=zfs \  
--fsname=hpc4dr01 \  
--ost \  
--index=0 \  
--mgsgnode 10.24.52.14@tcp0:10.27.52.14@o2ib \  
--mgsgnode 10.24.52.15@tcp0:10.27.52.15@o2ib \  
--servicenode 10.24.52.14@tcp0:10.27.52.14@o2ib \  
--servicenode 10.24.52.15@tcp0:10.27.52.15@o2ib \  
hpc4dr01-ost0/ost0
```

What does the lustre puppet module do

- it optionally installs the nagios plugins for monitoring (although we do not use them)
- it configures an HA cluster with corosync and pacemaker

Lustre cluster HA

- both servers “see” all the disks
- normally each server only imports some of the zfs pools
- each server mounts the filesystems present on the imported zfs pools
- when a server fails the remaining one becomes in charge of all resources (there might be a performance loss but the system stays active), after making sure that the malfunctioning server will not come back alive
- the clients must know all the servers so that they only contact the ones that are active



```
[statio@hpc4dr-login ~]$ grep lustre /etc/fstab
10.24.52.14@tcp0:10.24.52.15@tcp0:/hpc4dr01 /mnt lustre defaults
/mnt/home /home lustre bind 0 0
/mnt/data /data lustre bind 0 0
/mnt/software /software lustre bind 0 0
/mnt/projects /projects lustre bind 0 0
[statio@hpc4dr-login ~]$
```

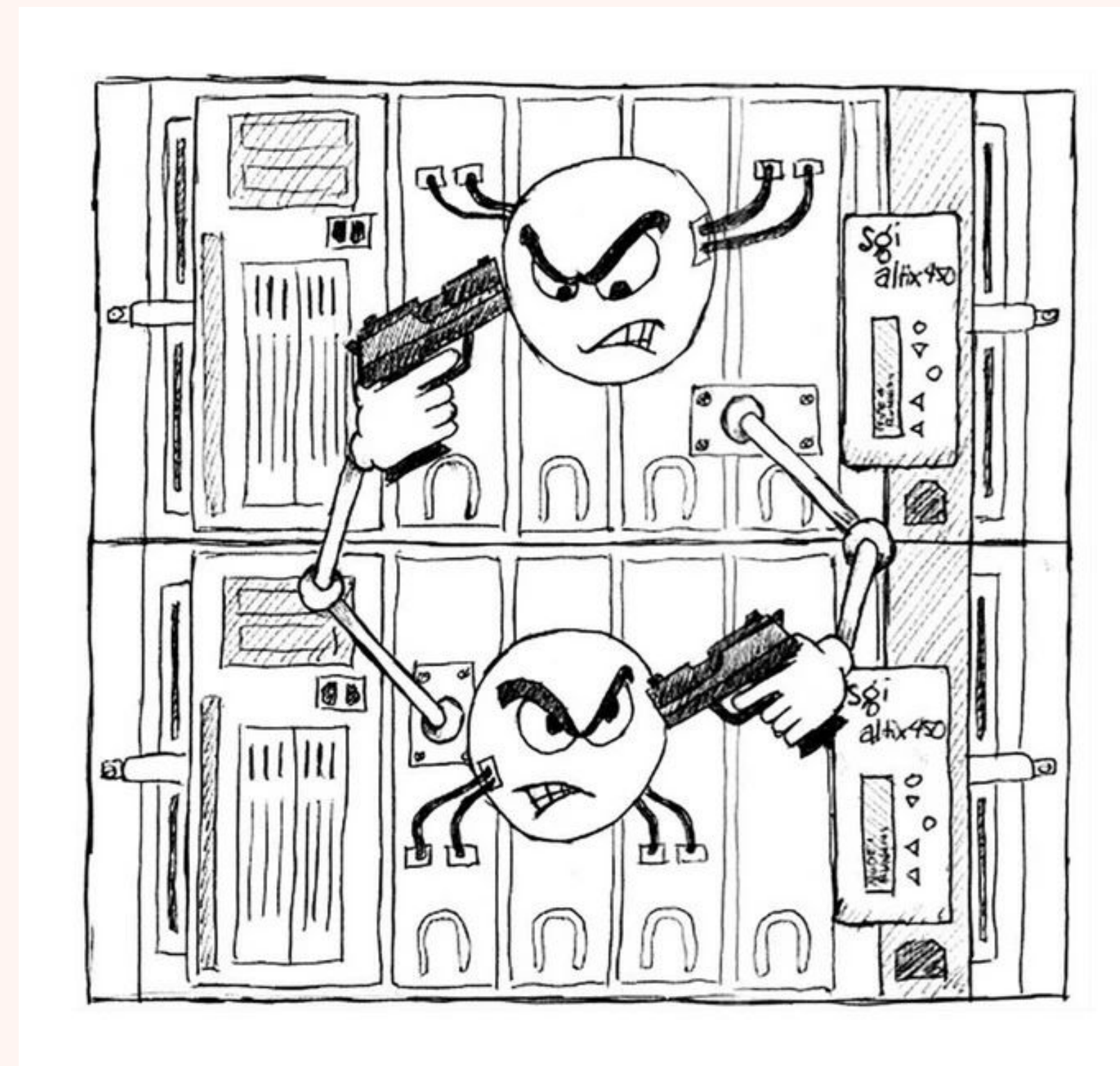
STONITH = Shoot The Other Node In The Head

STONITH ("Shoot The Other Node In The Head" or "Shoot The Offending Node In The Head"), is a technique for [fencing](#) in [computer clusters](#)

Fencing is the isolation of a failed [node](#) so that it does not cause disruption to a computer cluster. As its name suggests, STONITH fences failed nodes by resetting or powering down the failed node.

Multi-node error-prone contention in a cluster can have catastrophic results, such as if both nodes try writing to a shared storage resource. STONITH provides effective, if rather drastic, protection against these problems. A STONITH decision can be based on various decisions which can be customer specific plugins.

Google's inclusive language developer documentation discourages usage of this term.



Lustre cluster HA

- This is the pacemaker status on a healthy cluster
- Resources are evenly distributed among the storage servers
- No failures are reported

```
[root@stor-1 ~]# pcs status
Cluster name: lustre
Status of pacemakerd: 'Pacemaker is running' (last updated 2024-11-24 06:03:40 +01:00)
Cluster Summary:
 * Stack: corosync
 * Current DC: stor-1 (version 2.1.4-5.el8_7.2-dc6eb4362e) - partition with quorum
 * Last updated: Sun Nov 24 06:03:41 2024
 * Last change: Thu Nov 21 10:45:13 2024 by hacluster via crmd on stor-2
 * 2 nodes configured
 * 16 resource instances configured

Node List:
 * Online: [ stor-1 stor-2 ]

Full List of Resources:
 * ipmi-poweroff-stor-1      (stonith:fence_ipmilan):      Started stor-2
 * ipmi-poweroff-stor-2      (stonith:fence_ipmilan):      Started stor-1
 * Resource Group: hpc4dr01_OST0:
 * ZFS_hpc4dr01_OST0 (ocf::heartbeat:ZFS):      Started stor-2
 * lustre_hpc4dr01_OST0 (ocf::lustre:Lustre):      Started stor-2
 * Resource Group: hpc4dr01_OST1:
 * ZFS_hpc4dr01_OST1 (ocf::heartbeat:ZFS):      Started stor-1
 * lustre_hpc4dr01_OST1 (ocf::lustre:Lustre):      Started stor-1
 * Resource Group: hpc4dr01_OST2:
 * ZFS_hpc4dr01_OST2 (ocf::heartbeat:ZFS):      Started stor-2
 * lustre_hpc4dr01_OST2 (ocf::lustre:Lustre):      Started stor-2
 * Resource Group: hpc4dr01_OST3:
 * ZFS_hpc4dr01_OST3 (ocf::heartbeat:ZFS):      Started stor-1
 * lustre_hpc4dr01_OST3 (ocf::lustre:Lustre):      Started stor-1
 * Resource Group: hpc4dr01_MDT0:
 * hpc4dr01_ZFS_MDT0 (ocf::heartbeat:ZFS):      Started stor-2
 * hpc4dr01_lustre_MDT0 (ocf::lustre:Lustre):      Started stor-2
 * Resource Group: hpc4dr01_MDT1:
 * hpc4dr01_ZFS_MDT1 (ocf::heartbeat:ZFS):      Started stor-1
 * hpc4dr01_lustre_MDT1 (ocf::lustre:Lustre):      Started stor-1
 * Resource Group: MGT:
 * ZFS_MGT (ocf::heartbeat:ZFS):      Started stor-1
 * lustre_MGT (ocf::lustre:Lustre):      Started stor-1

Daemon Status:
corosync: active/enabled
pacemaker: active/enabled
pcsd: active/enabled
```


Lustre cluster HA

- This is the pacemaker status on a degraded cluster
- Resources are unevenly distributed
- Some failures are reported
- Although with reduced performance, the cluster is functional

```

pcsd: active/enabled
[root@stor-2 ~]# pcs status
Cluster name: lustre
Status of pacemakerd: 'Pacemaker is running' (last updated 2024-11-18 05:35:28 +01:00)
Cluster Summary:
* Stack: corosync
* Current DC: stor-1 (version 2.1.4-5.el8_7.2-dc6eb4362e) - partition with quorum
* Last updated: Mon Nov 18 05:35:29 2024
* Last change: Wed Oct 9 10:39:20 2024 by root via cibadmin on stor-1
* 2 nodes configured
* 16 resource instances configured

Node List:
* Online: [ stor-1 stor-2 ]

Full List of Resources:
* ipmi-poweroff-stor-1 (stonith:fence_ipmilan): Started stor-2
* ipmi-poweroff-stor-2 (stonith:fence_ipmilan): Started stor-1
* Resource Group: hpc4dr01_OST0:
* ZFS_hpc4dr01_OST0 (ocf::heartbeat:ZFS): Started stor-2
* lustre_hpc4dr01_OST0 (ocf::lustre:Lustre): Started stor-2
* Resource Group: hpc4dr01_OST1:
* ZFS_hpc4dr01_OST1 (ocf::heartbeat:ZFS): Started stor-2
* lustre_hpc4dr01_OST1 (ocf::lustre:Lustre): Started stor-2
* Resource Group: hpc4dr01_OST2:
* ZFS_hpc4dr01_OST2 (ocf::heartbeat:ZFS): Started stor-2
* lustre_hpc4dr01_OST2 (ocf::lustre:Lustre): Started stor-2
* Resource Group: hpc4dr01_OST3:
* ZFS_hpc4dr01_OST3 (ocf::heartbeat:ZFS): Started stor-2
* lustre_hpc4dr01_OST3 (ocf::lustre:Lustre): Started stor-2
* Resource Group: hpc4dr01_MDT0:
* hpc4dr01_ZFS_MDT0 (ocf::heartbeat:ZFS): Started stor-2
* hpc4dr01_lustre_MDT0 (ocf::lustre:Lustre): Started stor-2
* Resource Group: hpc4dr01_MDT1:
* hpc4dr01_ZFS_MDT1 (ocf::heartbeat:ZFS): Started stor-2
* hpc4dr01_lustre_MDT1 (ocf::lustre:Lustre): Started stor-2
* Resource Group: MGT:
* ZFS_MGT (ocf::heartbeat:ZFS): Started stor-1
* lustre_MGT (ocf::lustre:Lustre): Started stor-1

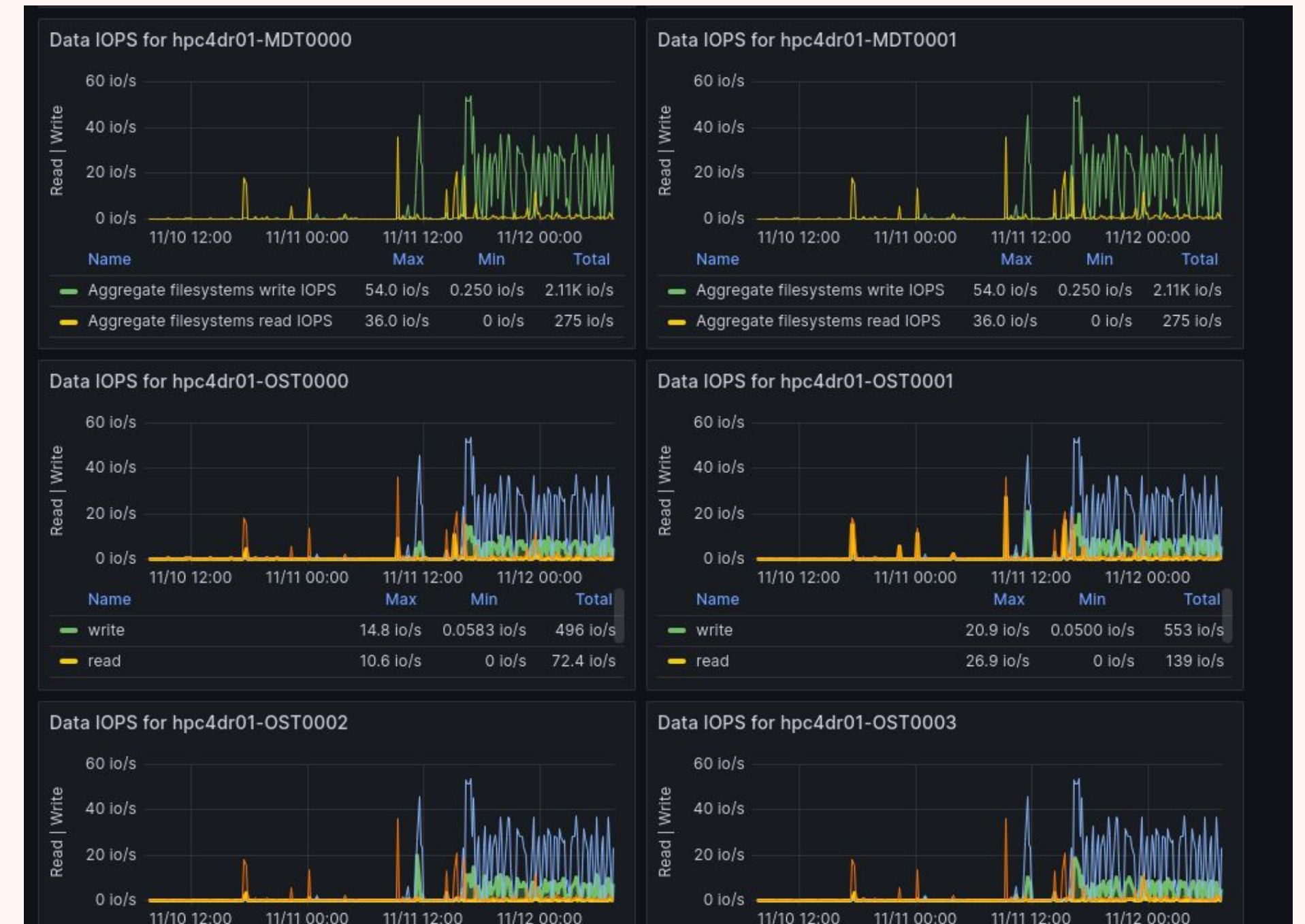
Failed Resource Actions:
* lustre_hpc4dr01_OST0_start_0 on stor-1 'error' (1): call=76, status='complete', last-rc-change='Sat Sep
* lustre_hpc4dr01_OST2_start_0 on stor-1 'error' (1): call=77, status='complete', last-rc-change='Sat Sep
* hpc4dr01_lustre_MDT0_start_0 on stor-1 'error' (1): call=78, status='complete', last-rc-change='Sat Sep

Daemon Status:
corosync: active/enabled
pacemaker: active/enabled
pcsd: active/enabled
[root@stor-2 ~]#

```

Monitoring

- we use check-mk for monitoring the infrastructure (others may use zabbix or nagios)
 - server health
 - SAN health
 - server load
 - disk space
 - network connectivity
 - services health
 -
- we use a prometheus exporter + prometheus + grafana for monitoring the lustre cluster
- we wrote a custom exporter for monitoring disk usage per user, still needs to be fully integrated in grafana





FUTURE DEVELOPMENTS

- A Lustre infrastructure is difficult to deploy and maintain, particularly if a whole data center relies on it as its main storage backend
- It requires a deep knowledge of the software that one can only reach by dedicating to it a large amount of time, but we have little time and many items to work on
- We were not able to understand some aspects of the system and reach a fully satisfactory setup
- At LNGS we decided to move to a Lustre system supported by DDN

- <https://wiki.lustre.org/>
- <https://github.com/guilbaults/puppet-lustre>
- <https://hpcf.umbc.edu/general-productivity/lustre-best-practices/>
- <https://en.wikipedia.org/wiki/STONITH>

POSIX

Released in the late 1980s, POSIX (Portable Operating System Interface) is a family of standards created to make sure that applications developed on one UNIX flavor can run on other UNIXes. The POSIX standard describes how system calls must behave. One particular section of the standard defines the semantics (behavior) of a POSIX compatible file system. Today, many of the UNIX systems that existed back in the 80s and 90s have become irrelevant. On the other hand, Linux became ubiquitous; one of the reasons being the POSIX compatibility. Being POSIX compatible made it easier to move from a UNIX to LINUX.

For a file system to be compatible with POSIX, it must:

Implement strong consistency. For example, if a write happened before a read, the read must return the data written. Have atomic writes, where a read either returns all data written by a concurrent write or none of the data but is not an incomplete write.

Implement certain operations, like random reads, writes, truncate, or fsync.

Control access to files using permissions (see here) and implement calls like chmod, chown, and so on to modify them.