

2494-24

**Workshop on High Performance Computing (HPC) Architecture
and Applications in the ICTP**

14 – 25 October 2013

**High availability Lustre FS
implementation for genomic**

Francesco De Giorgi
eXact Lab, Trieste

The logo for 'exact' features a large, stylized 'X' on the left. The 'X' is formed by two thick, black diagonal lines that intersect. To the left of the 'X', there are several blue lines of varying lengths, each ending in a small blue dot, resembling a network or data flow diagram. The word 'exact' is written in a white, sans-serif font, with the 'x' being a blue color that matches the lines in the 'X' graphic.

exact

High availability Lustre FS
implementation for genomic

Francesco De Giorgi
francesco.degiorgi@exact-lab.it

ICTP, Trieste
October 24, 2013

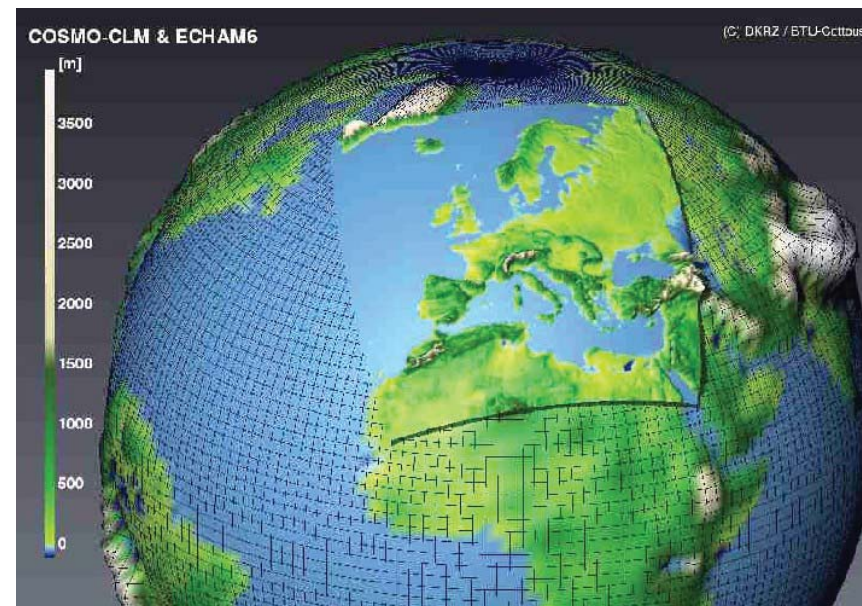
eXperience on advanced computational technologies

- spin-off of the

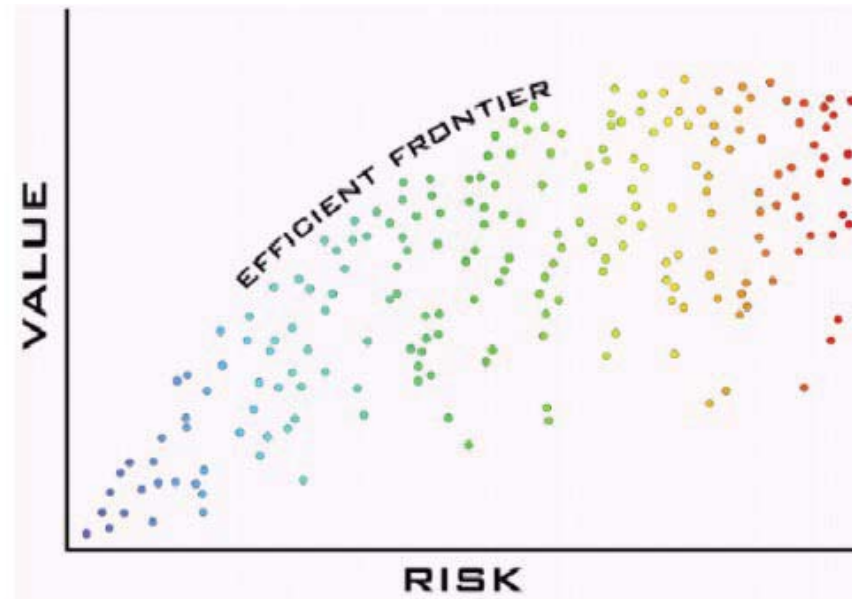


- leveraging **academic expertise**, the team provides solutions to the HPC market:
 - design and deployment of computational infrastructures
 - advanced training in HPC
 - advanced services for scientific/technical application

- HPDA: *High Performance Data Analysis*
 - tasks involving sufficient data volumes and algorithm complexity to require HPC resources
- Use cases
 - climate modeling
 - risk analysis
 - national security
 - life science



- HPDA: *High Performance Data Analysis*
 - tasks involving sufficient data volumes and algorithm complexity to require HPC resources
- Use cases
 - climate modeling
 - risk analysis
 - national security
 - life science



- HPDA: *High Performance Data Analysis*
 - tasks involving sufficient data volumes and algorithm complexity to require HPC resources
- Use cases
 - climate modeling
 - risk analysis
 - national security
 - life science

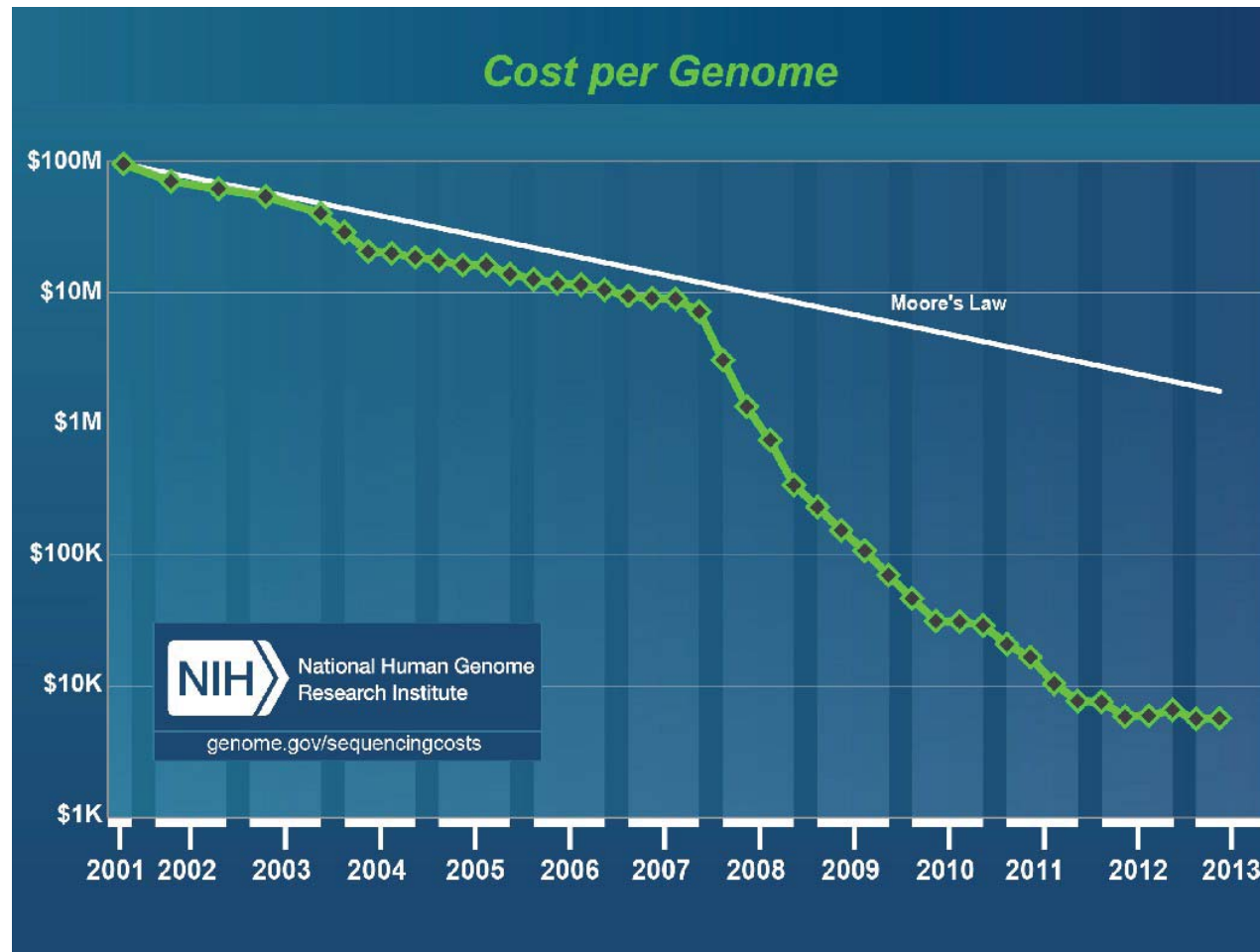


- **HPDA: *High Performance Data Analysis***
 - tasks involving sufficient data volumes and algorithm complexity to require HPC resources
- **Use cases**
 - climate modeling
 - risk analysis
 - national security
 - **life science**



Life science: DNA sequencing **eXact**

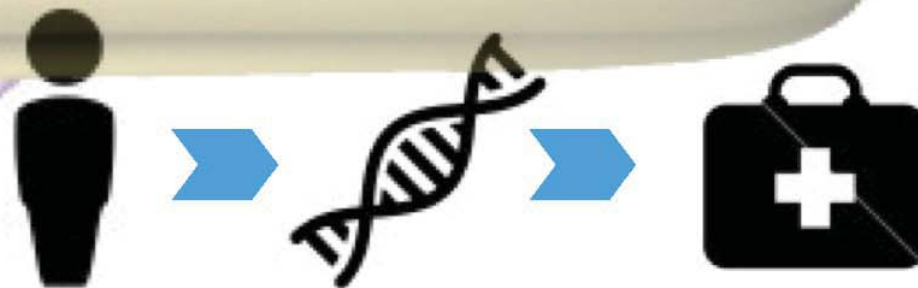
- High-throughput sequencing: a big data problem



DNA sequencing in HPC: a case study **eXact**

eXact lab's deployment of a highly available and fault-tolerant HPC infrastructure for genomic research

- For the Center for Translational Genomic and Bioinformatics
 - in San Raffaele Hospital (Milan, Italy)
- Final goal is **personalized medicine**
 - *customization of healthcare by use of genetic information*



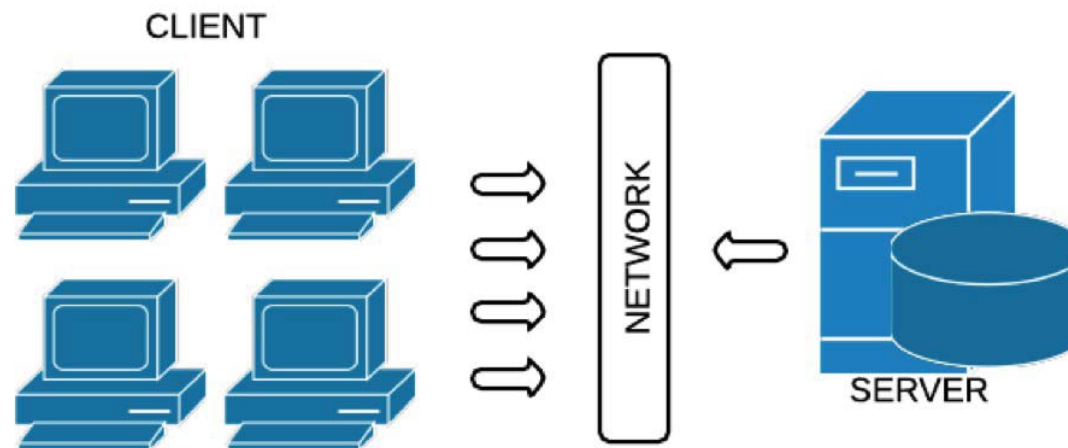
- **Huge amount of genomic data** from sequencer Illumina Hi-Seq 2000
 - To backup (~20k € per run)
 - To post-process
 - Always available
- ➔ Data from the sequencer need to be served to the computational infrastructure
- ➔ Need for a **fast, high performance, highly scalable file system**, with robust failover and recovery mechanisms

Need for a **fast, high performance, highly scalable file system**, with robust failover and recovery mechanisms

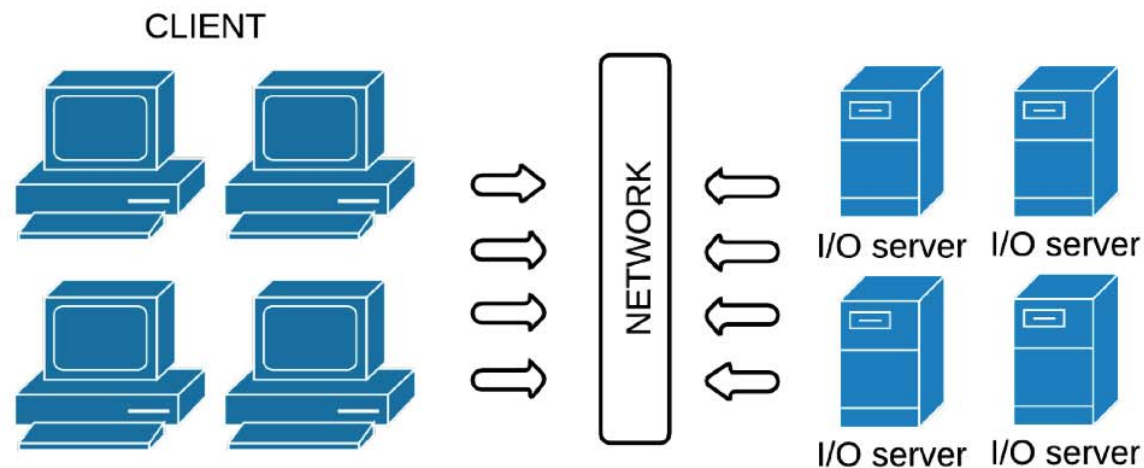
→ **Lustre File System**

- parallel and distributed
- high throughput, low latency
- scalable, redundant
- high availability features

- why a parallel and distributed file system?
- isn't NFS enough to address the requirements?
 - server is a **single point of failure** (unless data replication)
 - **not scalable**
 - **limited by network**



- a parallel FS distributes data on many servers
 - parallel I/O
 - low access latency
 - high throughput
 - scalable



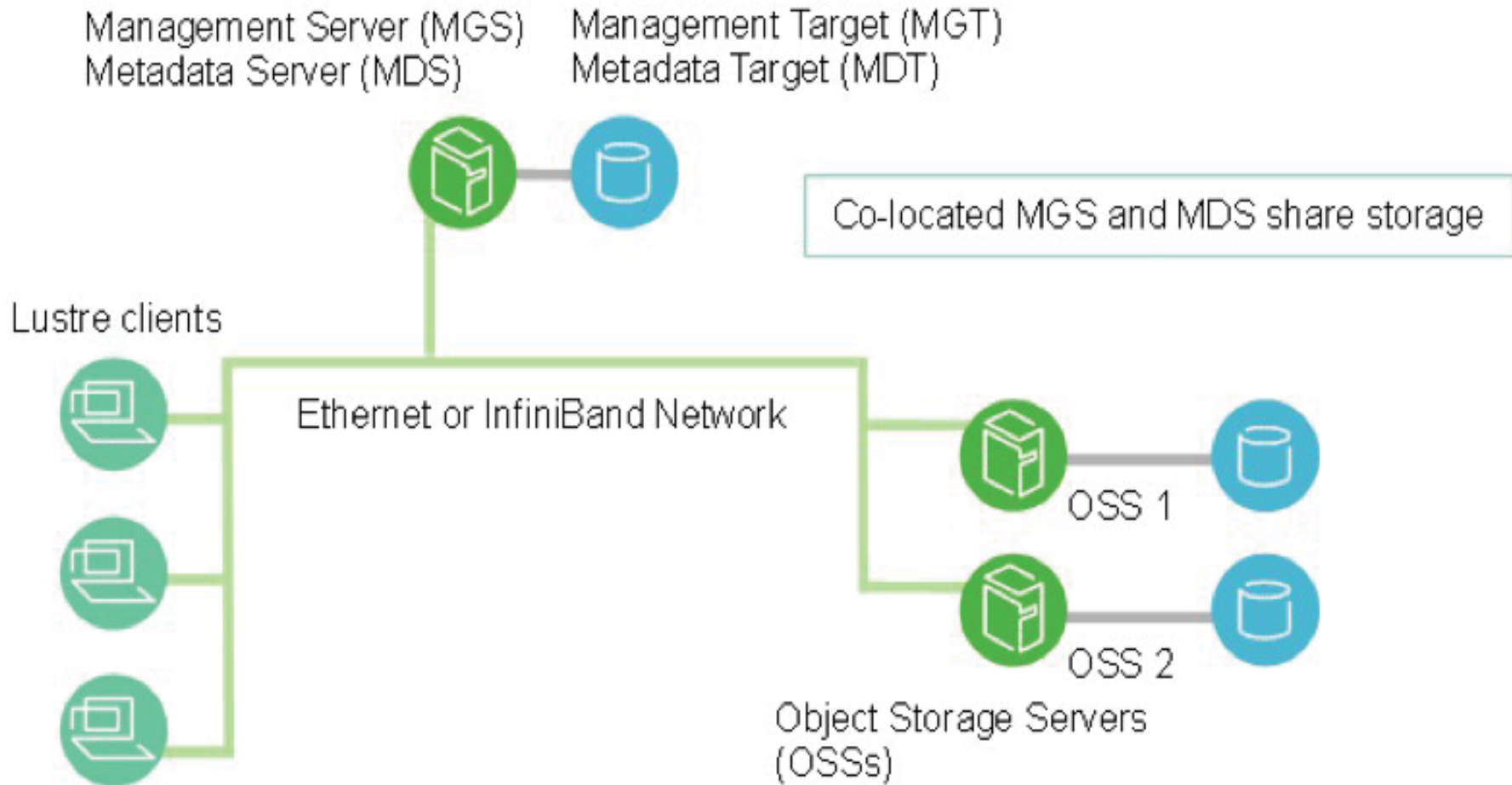
- a parallel FS distributes data on many servers
 - **pros**
 - scalability
 - performance
 - data integrity (thanks to redundancy and failover mechanism)
 - **cons**
 - needs a metadata server
 - not so easy to configure and manage as infrastructure grows

- part of Intel HPC strategy
 - road to exascale computing (by 2018)
 - 74% of CPUs in the TOP500 are Intel
- 6 on 10 in the TOP10 use Lustre
 - 60 in the TOP100
- used in heterogeneous environment
 - oil and gas, life science, meteorology, finance, ...
 - from a few Terabyte to tens of Petabyte

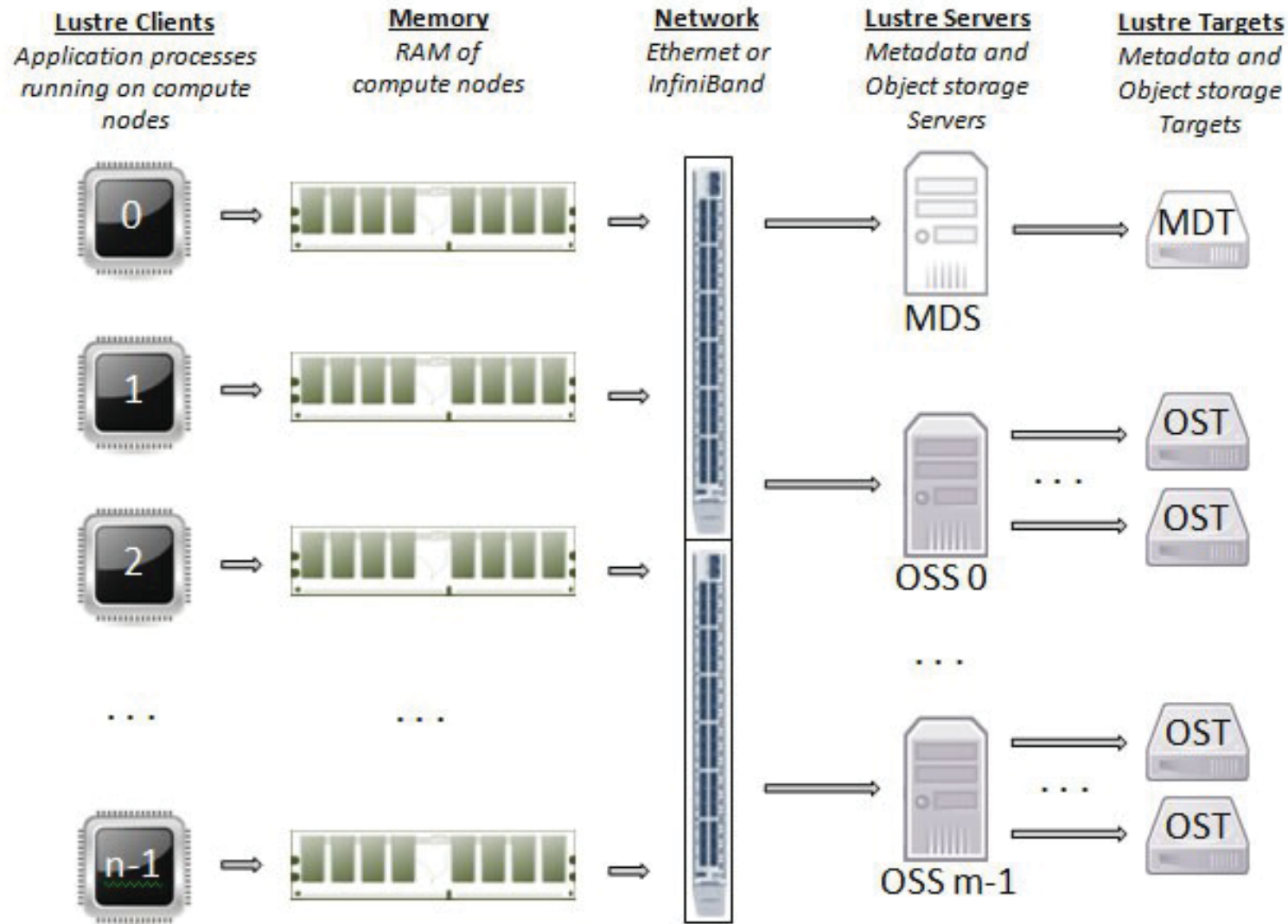
- 3 main components
 - **MDS**, *MetaData Server*
 - store metadata in the MDT, *MetaData Target*
 - **OSS**, *Object Storage Server*
 - store data files on one or more objects, each of them exists on a OST, *Object Storage Target*
 - **Lustre client**, computational nodes, desktops, an entity able to mount a Lustre file system
 - each client sees a single, coherent, synchronized namespace

- communication between servers and client is managed by Lustre networking (LNET)
- Lustre supports many high performance, low latency networks
 - InfiniBand: OpenFabric OFED (o2ib)
 - Myrinet: MX
 - any network carrying TCP traffic (GigE, 10GigE, IPoIB)
- permits RDMA, when supported by underlying networks (InfiniBand, Myrinet MX)

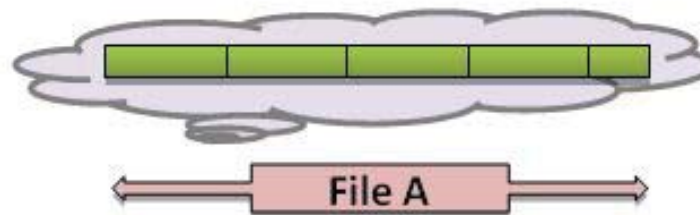
Basic architecture



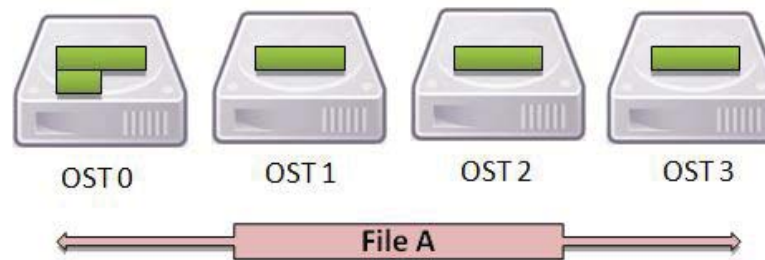
Lustre components



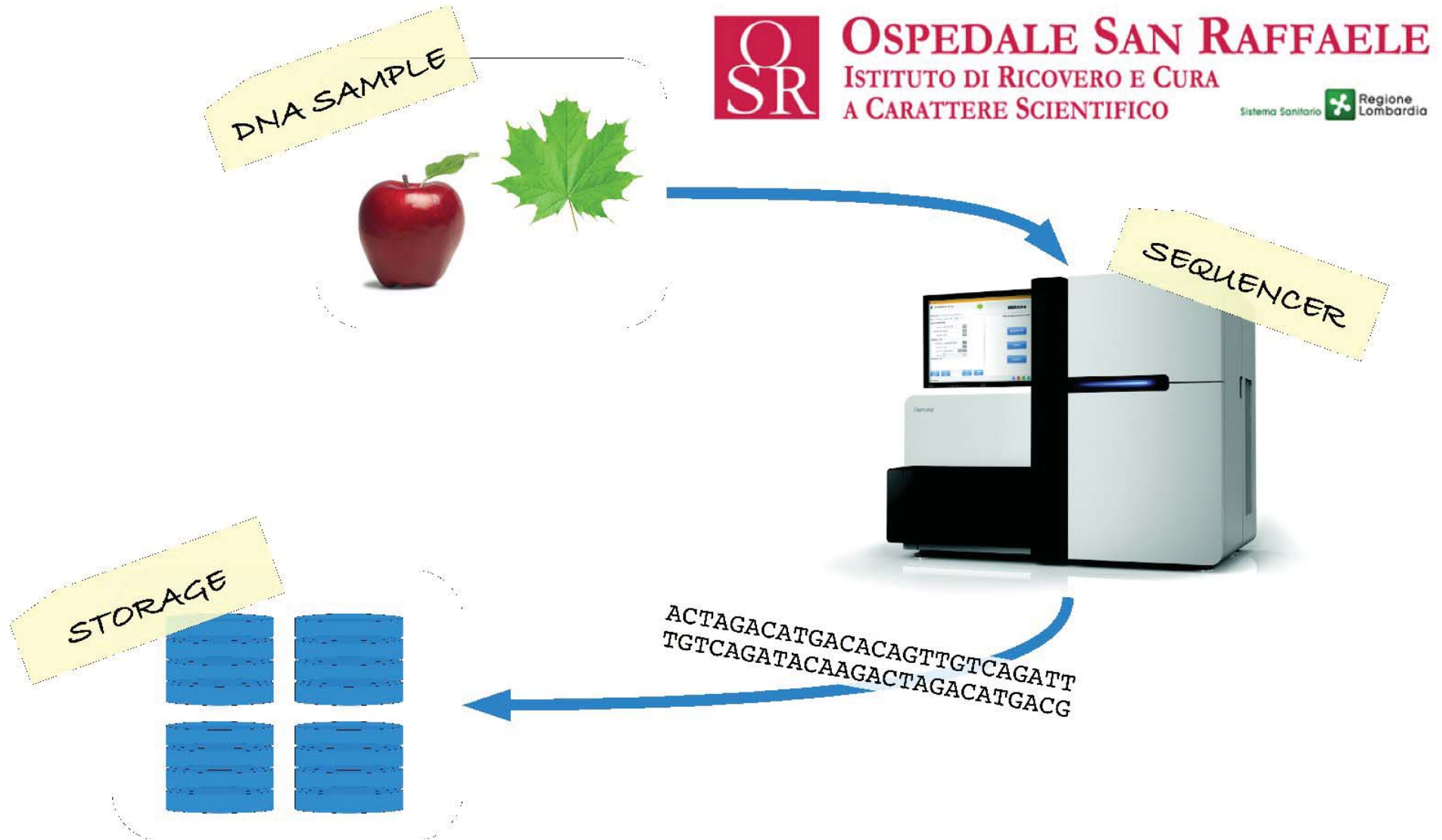
- Lustre can *stripe* a single file on multiple OSS
 - normally a file is a byte sequence...



- ...but Lustre can divide it into chunks



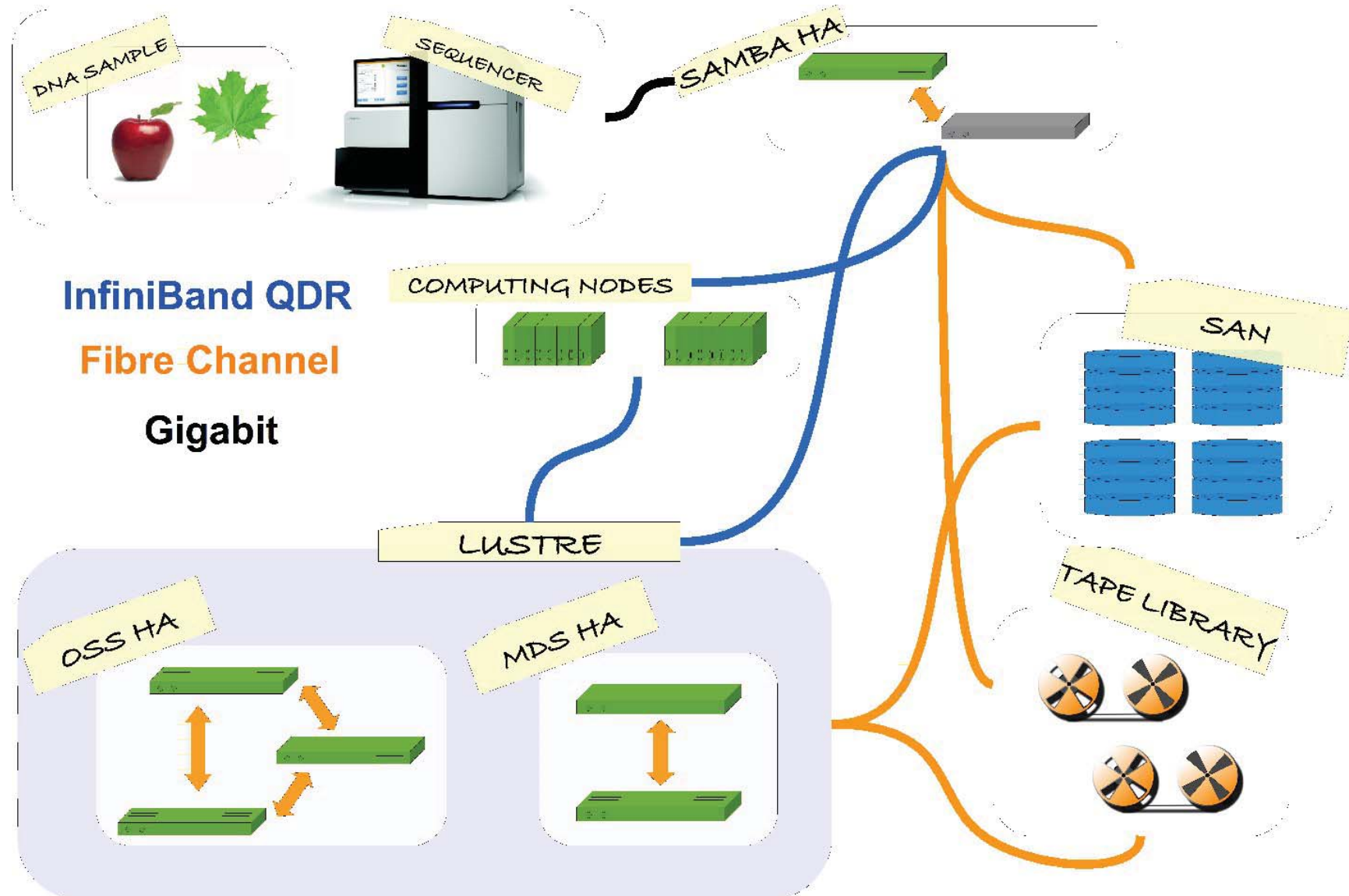
back to case study: the workflow



 **OSPEDALE SAN RAFFAELE**
ISTITUTO DI RICOVERO E CURA
A CARATTERE SCIENTIFICO

Sistema Sanitario  Regione Lombardia

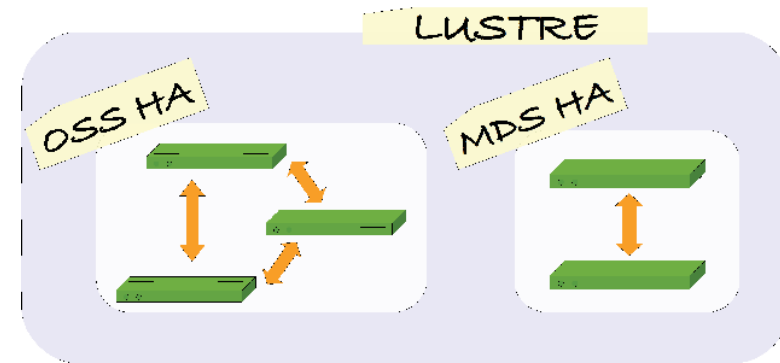
infrastructure overview



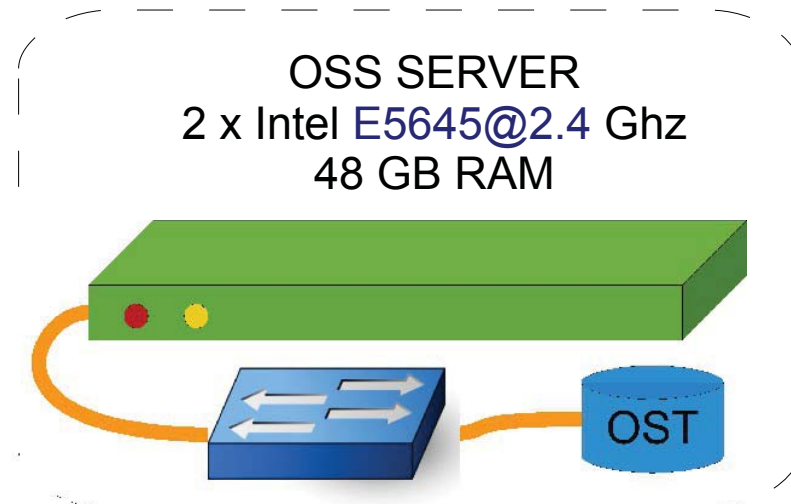
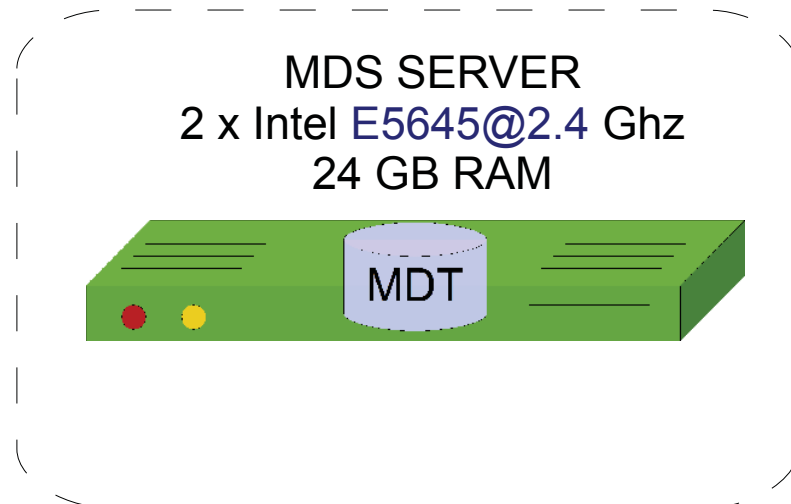
Lustre filesystem

2 Lustre filesystems

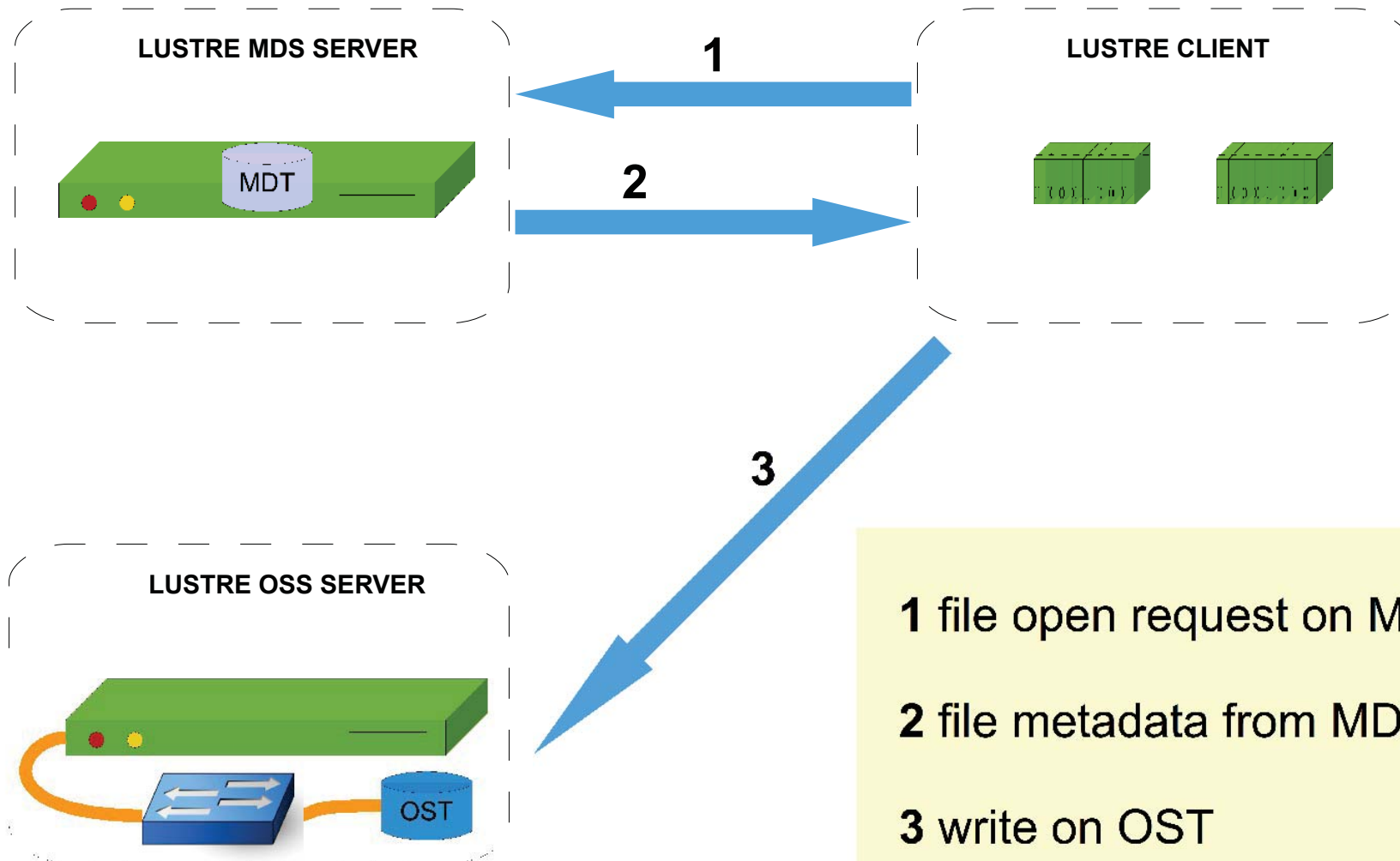
- 2 MDSs, 3 OSSs
- ~50 clients
- 60 terabytes from SAN



always available!

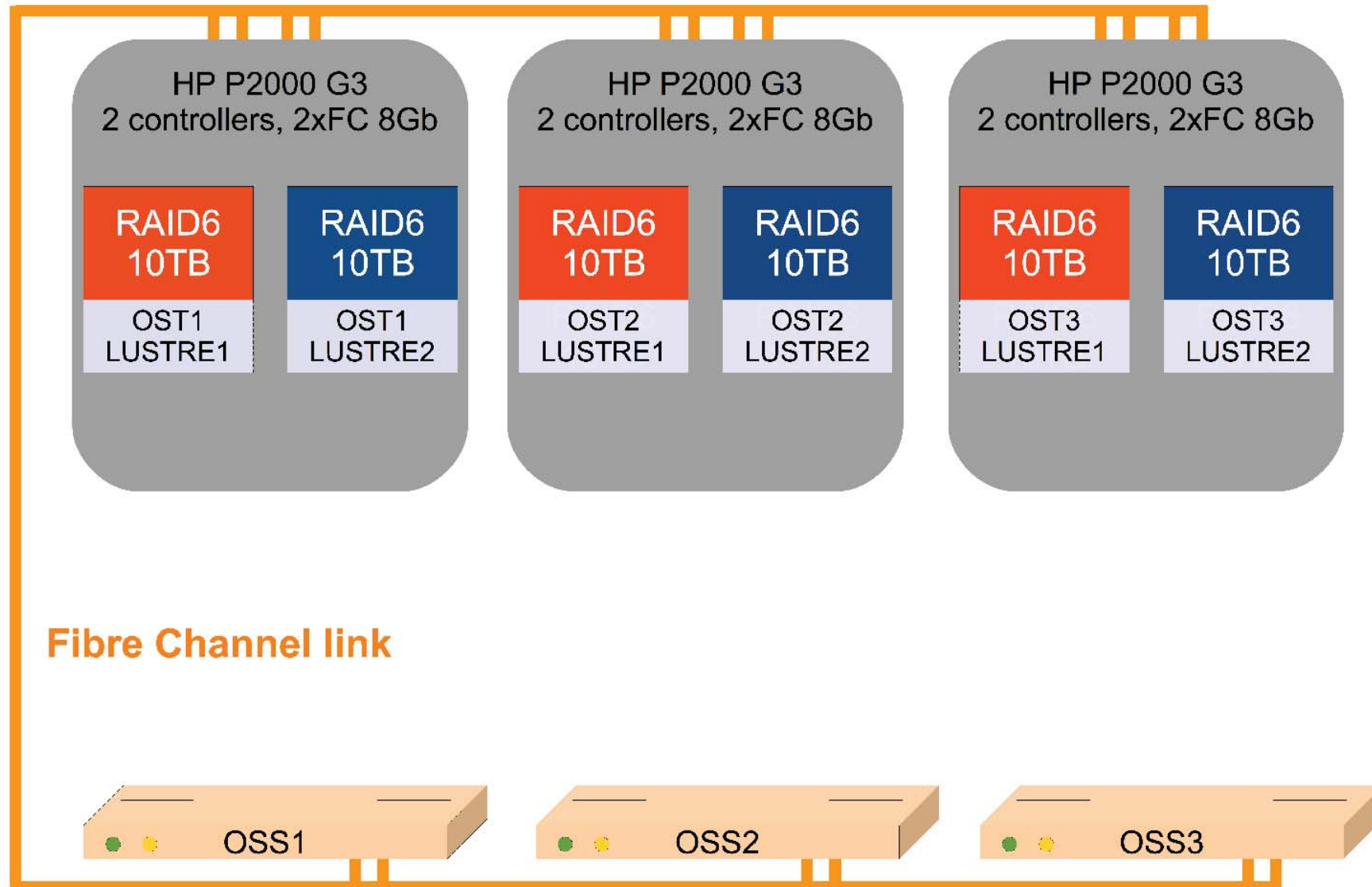


How I/O works in Lustre

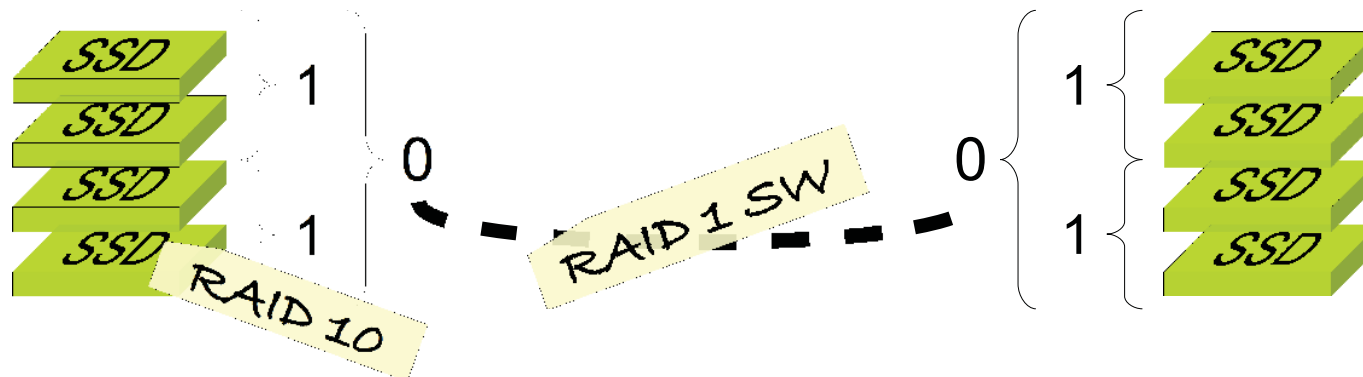
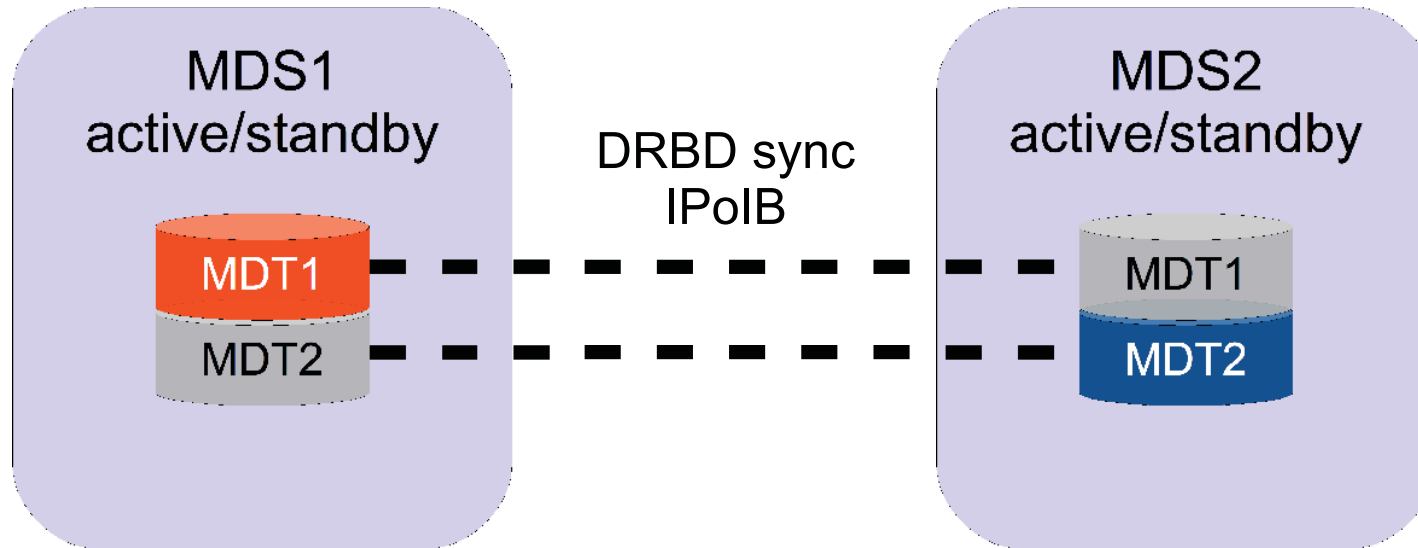



- 1 file open request on MDT
- 2 file metadata from MDT
- 3 write on OST

Lustre components: OSTs and OSSs



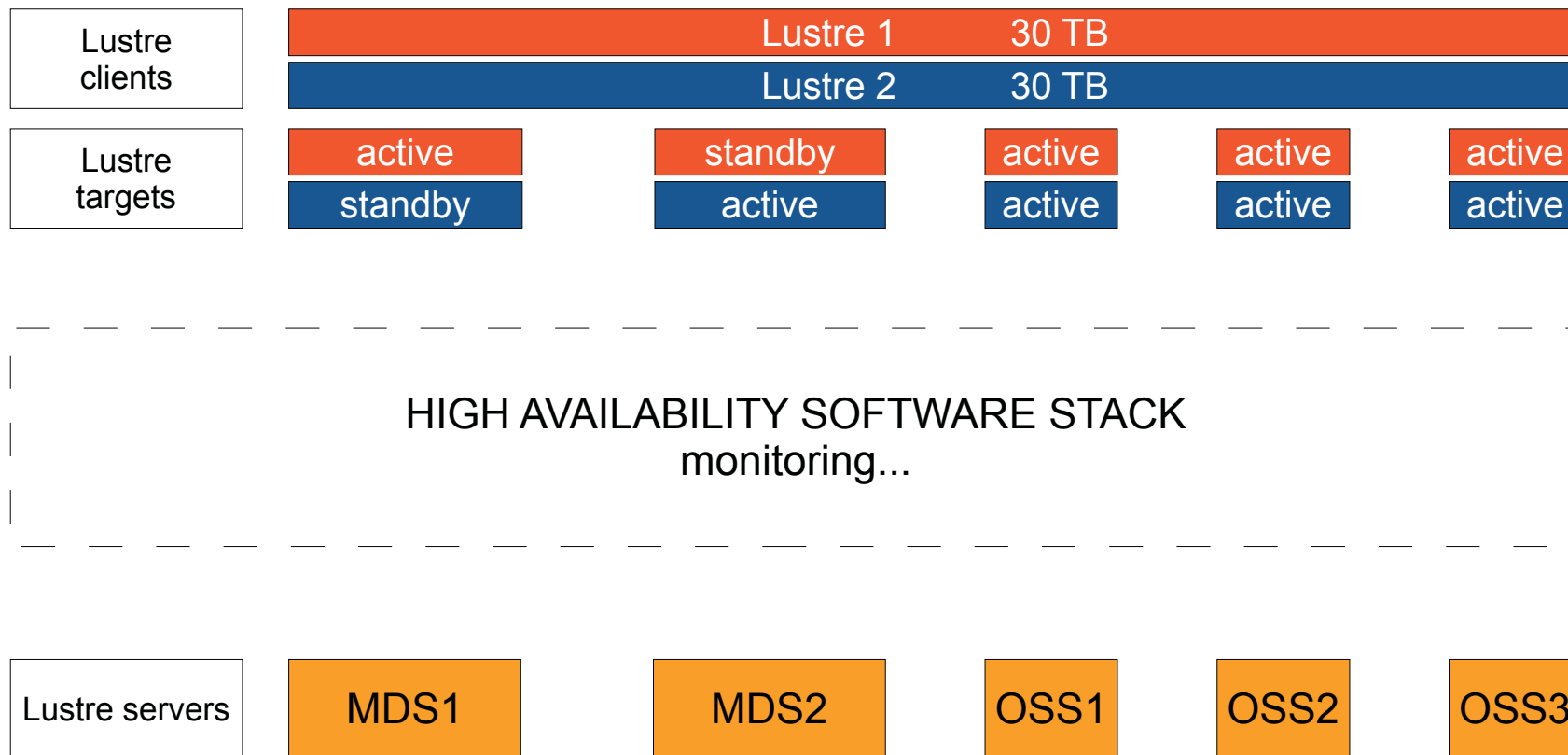
Lustre components: MDSs and MDTs



 { HP 100GB 3G SATA MLC LFF (3.5-inch)
SC Enterprise Mainstream Solid State Drive – PCI-e attached

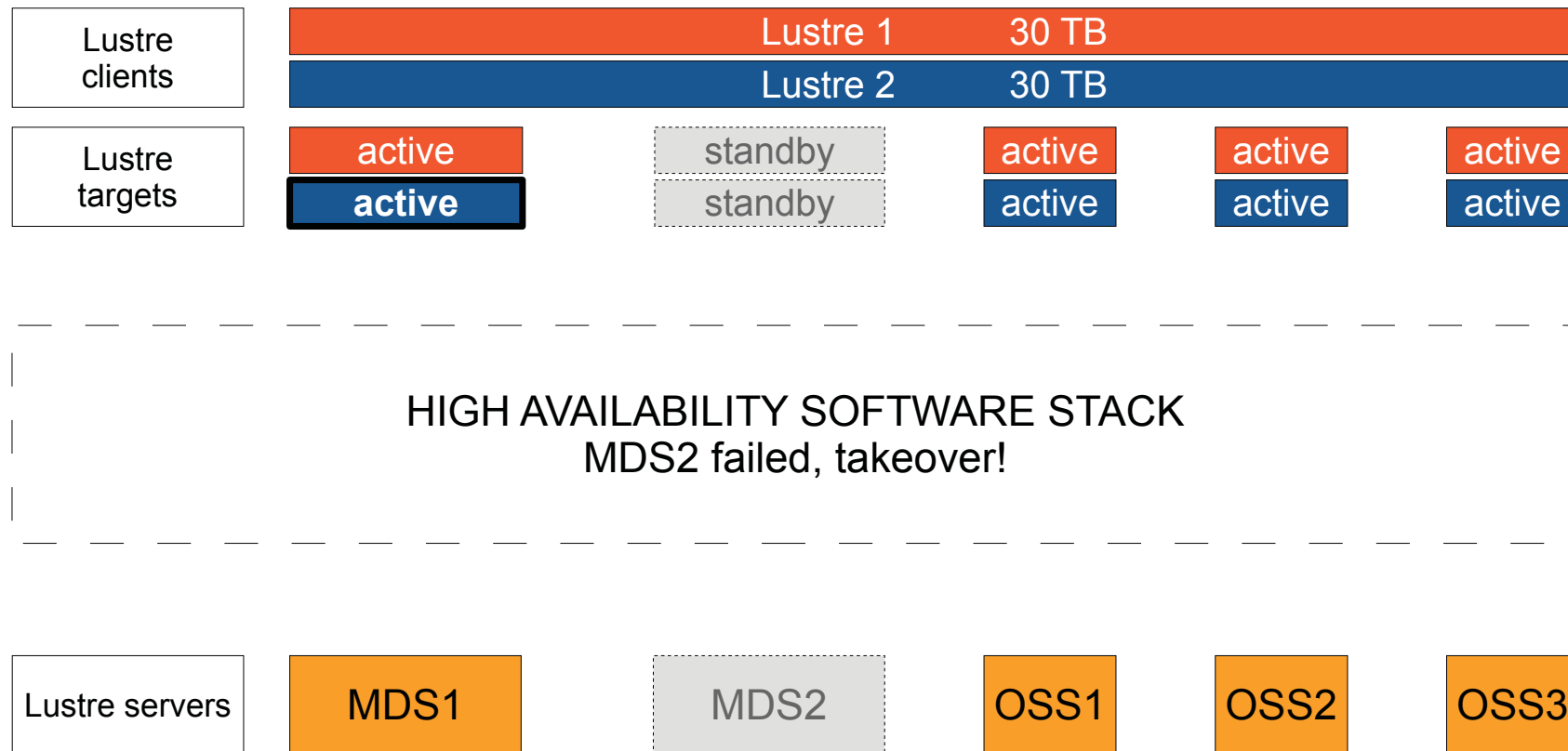
Lustre high availability

- In production, no failures

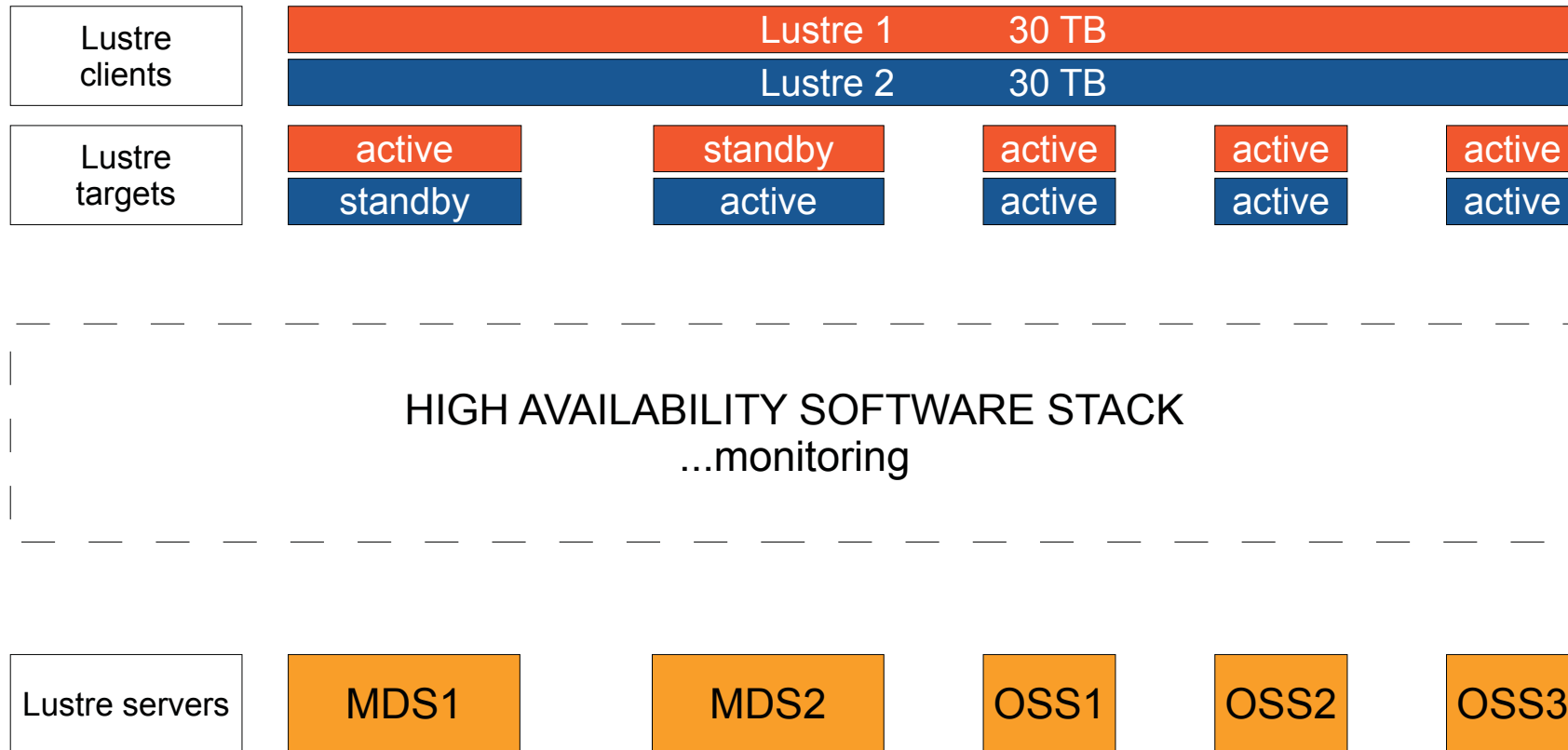


MDS2 failure

- MDS1 will take over the service of MDS2

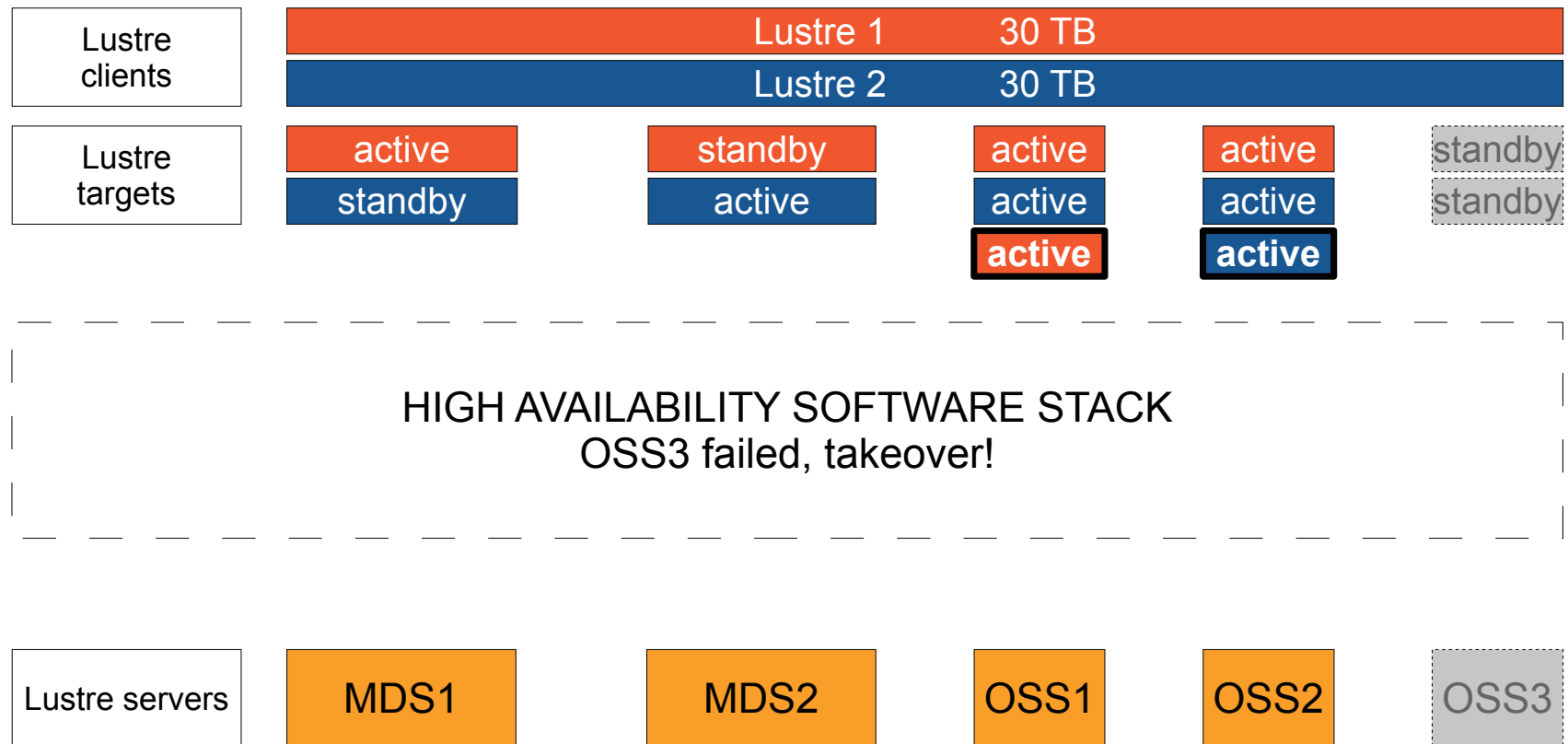


- MDS2 recovers its services



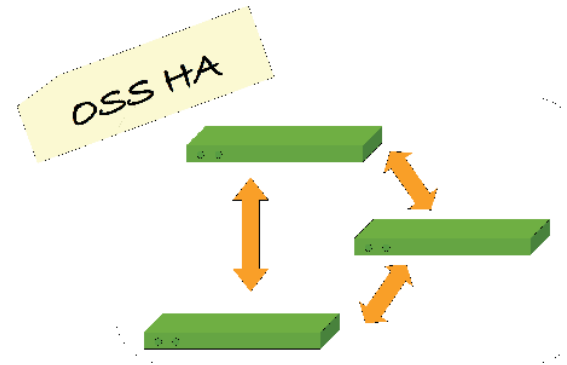
OSS3 failure

- OSS1 and OSS2 will take over its service

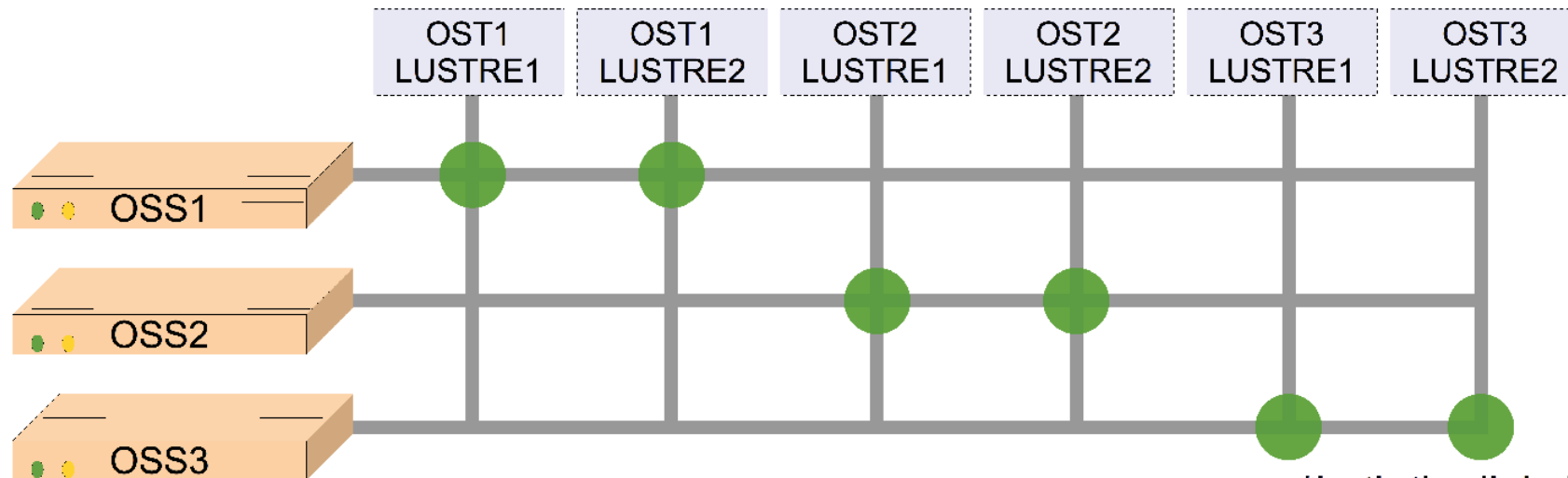


- Failures

- Power*
- Fibre channel*
- InfiniBand*



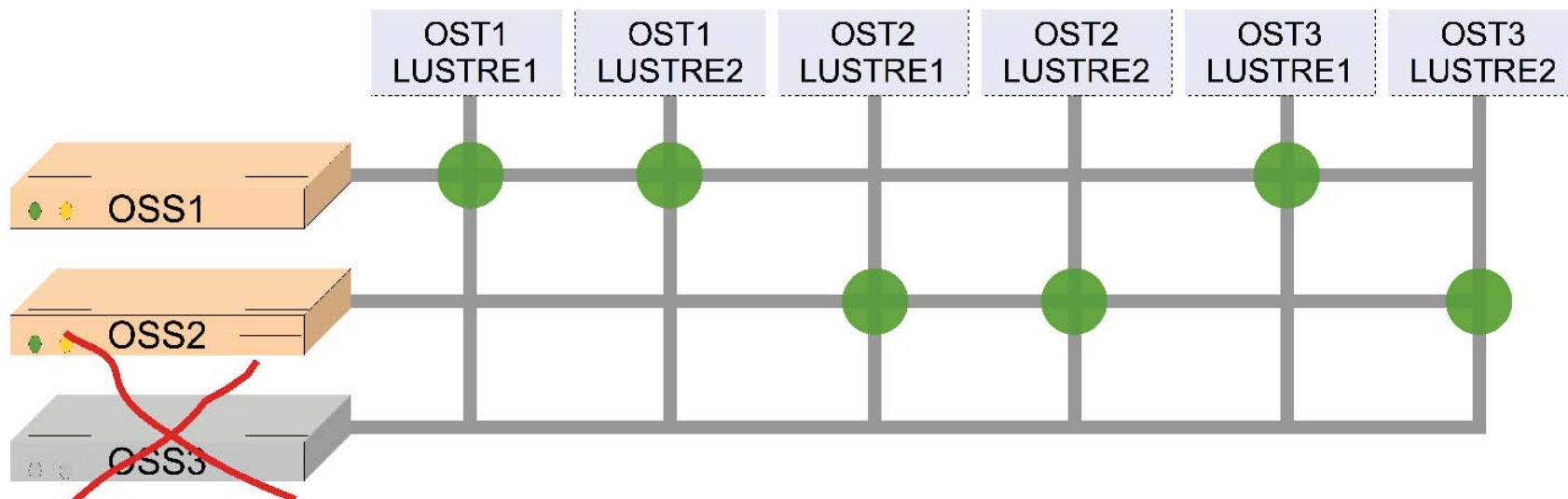
- In production, each OSS server mounts 2 OSTs



*both the links!

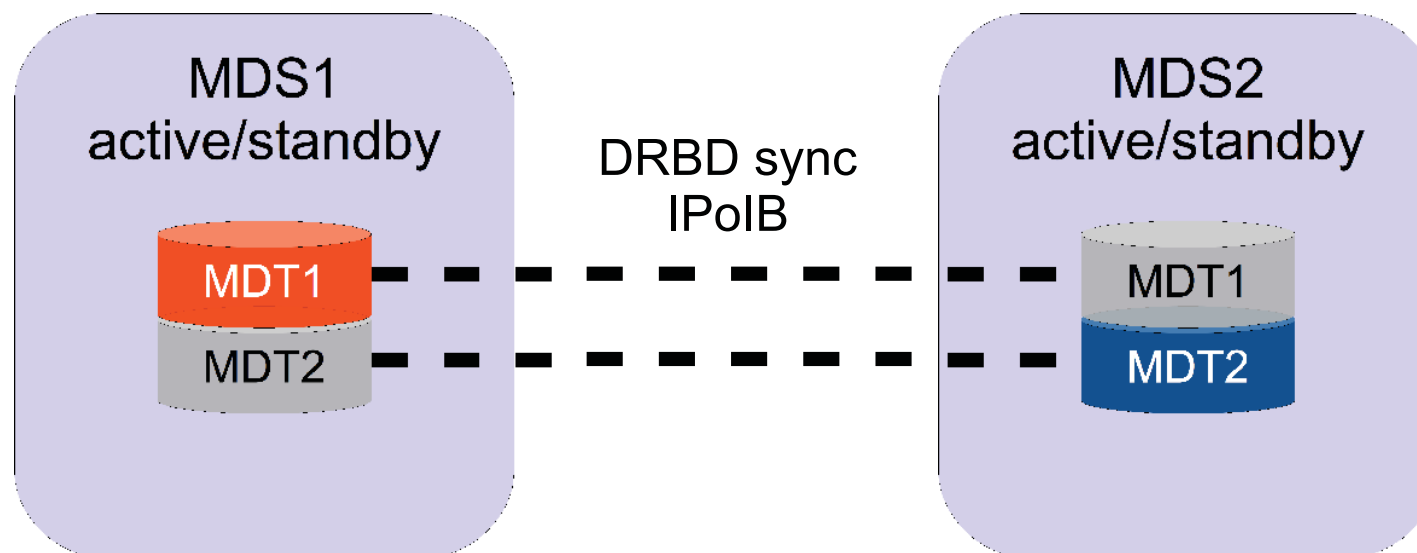
High availability on OSSs

- If OSS3 fails
 - the HA software will acknowledge the failure
- OSS2, OSS1 receive a new OST each



High availability on MDSs

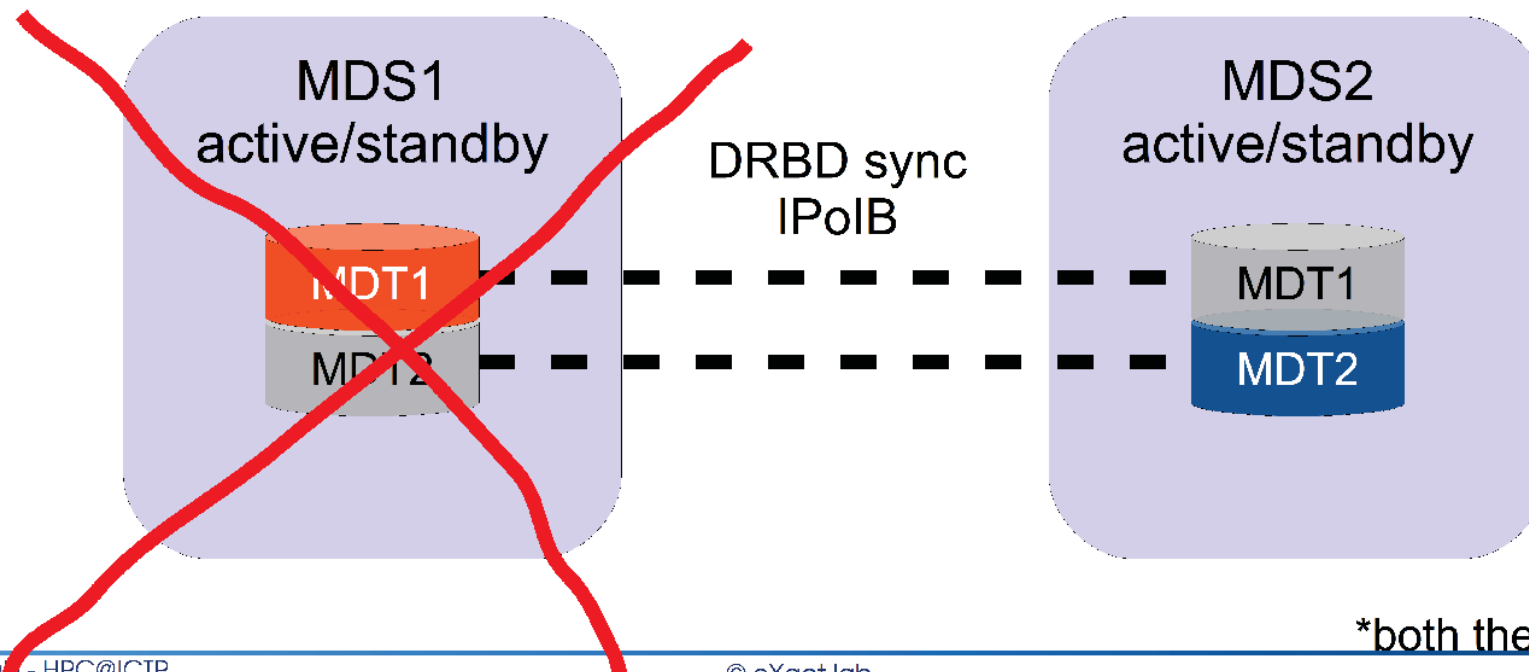
- In normal condition
 - MDTs are replicated between MDSs
 - only one replica is active for Lustre client



*both the links!

MDS1 failure

- data integrity **MUST** be ensured
- MDS1 *irresponsive* \neq *isn't accessing my data*
- clients must not be able to reach MDS1!

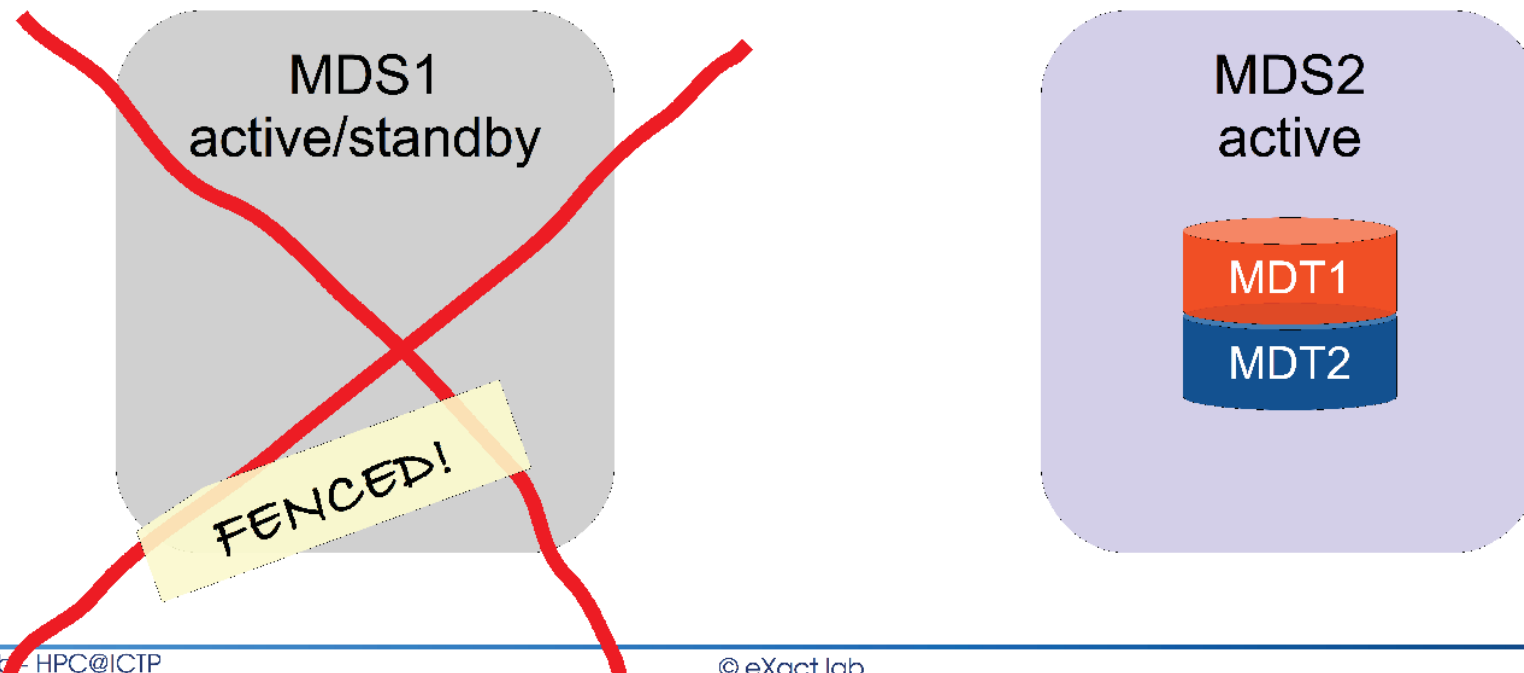


*both the links!

how to ensure data integrity **eXact**



- STONITH aka *Shoot The Other Node In The Head!*
 - MDS1 fails
 - MDS2 takes over its services
 - **MDS2 powers off MDS1 !**



High availability tests

- Unplug → *failover*
 - Power*
 - InfiniBand*
 - Fibre Channel (OSS)*
 - InfiniBand + Fibre Channel
- Replug → *failback*



DOWNTIME = ~80s

➔ Completely transparent for clients!

High Availability Lustre
FS implementation
for Genomic



info@exact-lab.it

www.exact-lab.it