# Choose Lustre

## Stephen Simms

Manager, High Performance File Systems

Indiana University

# Lustre is scalable – 55 PB at LLNL

# Lustre is fast – 1 TB/s at ORNL

# Lustre can support thousands of clients

# Lustre is Open

# Lustre is Open source software under GPLv2

That means it's "Free Like Beer" right?
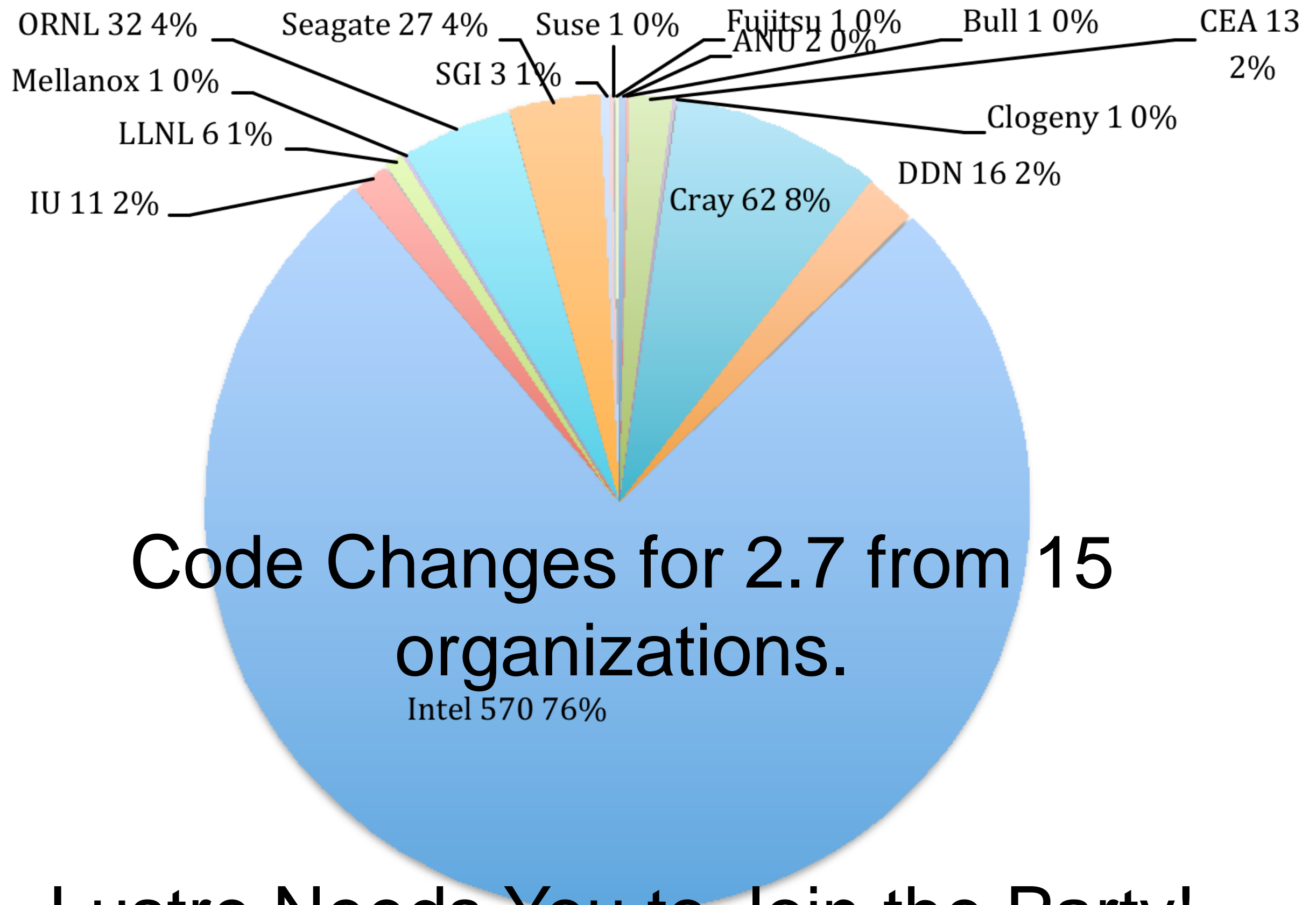
# Actually, more like a free puppy…

# It Takes Lots of Work to Maintain Lustre

- Bug Fixes
- Rigorous Testing
- Feature Development
- Maintaining Documentation
- Tree Hosting
- Code Reviews by Peers

# Many Hands Make Lighter Work

ORNL 32 4%   Seagate 27 4%   Suse 1 0%   Fujitsu 1 0%   Bull 1 0%   CEA 13 2%

Mellanox 1 0%   SGI 3 1%   ANU 2 0%   Clogeny 1 0%

LLNL 6 1%

IU 11 2%   Cray 62 8%   DDN 16 2%

Code Changes for 2.7 from 15 organizations.

Intel 570 76%

Lustre Needs You to Join the Party!

# Lustre is moving forward

- Hiccup when Lustre moved from Oracle
  - Lustre 2.0 – Fall 2010
  - Lustre 2.1 – Fall 2011

- Since then Lustre has accelerated
  - 2 major releases a year
  - Spring / Fall

- Users following Lustre releases
  - Majority using 2.5 in production

# Alright, alright, stop the marketing

Lustre = Linux +Cluster

Lustre is a *parallel* distributed file system

- High performance filesystem used by >60 of the top 100 supercomputers in the world
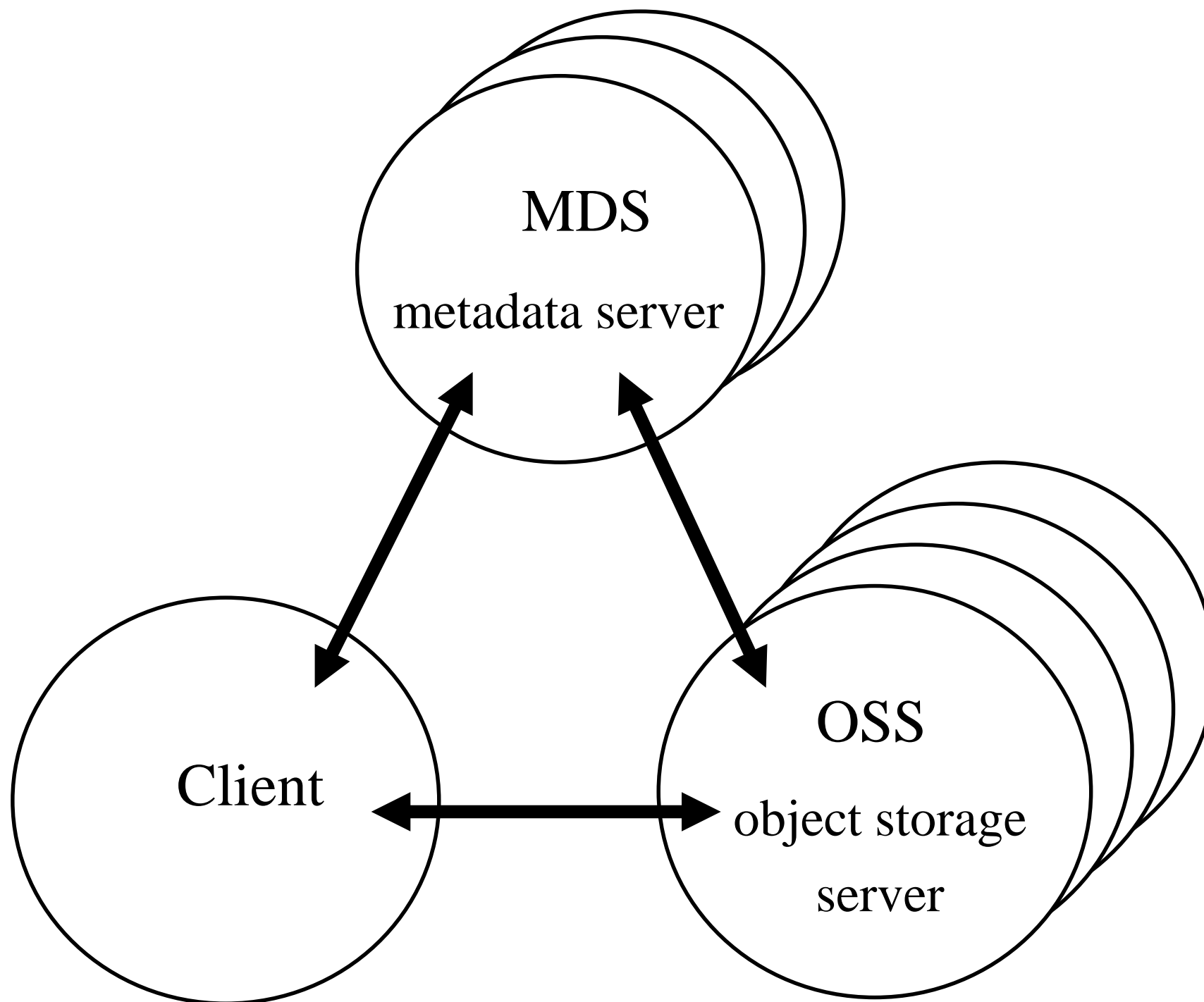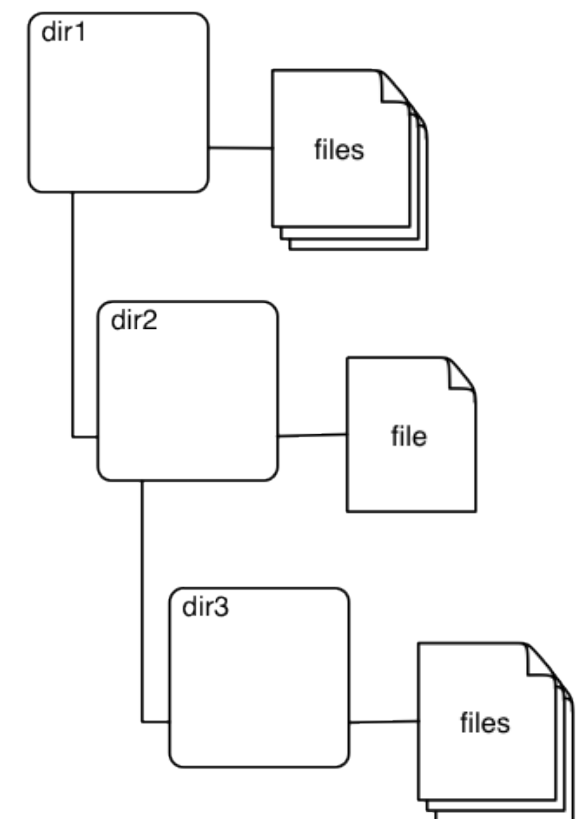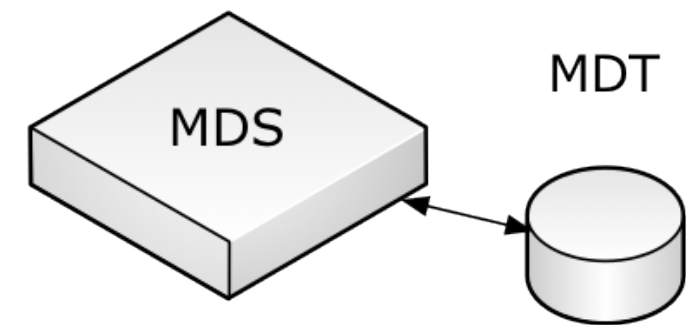- POSIX compliant – behaves like other file systems
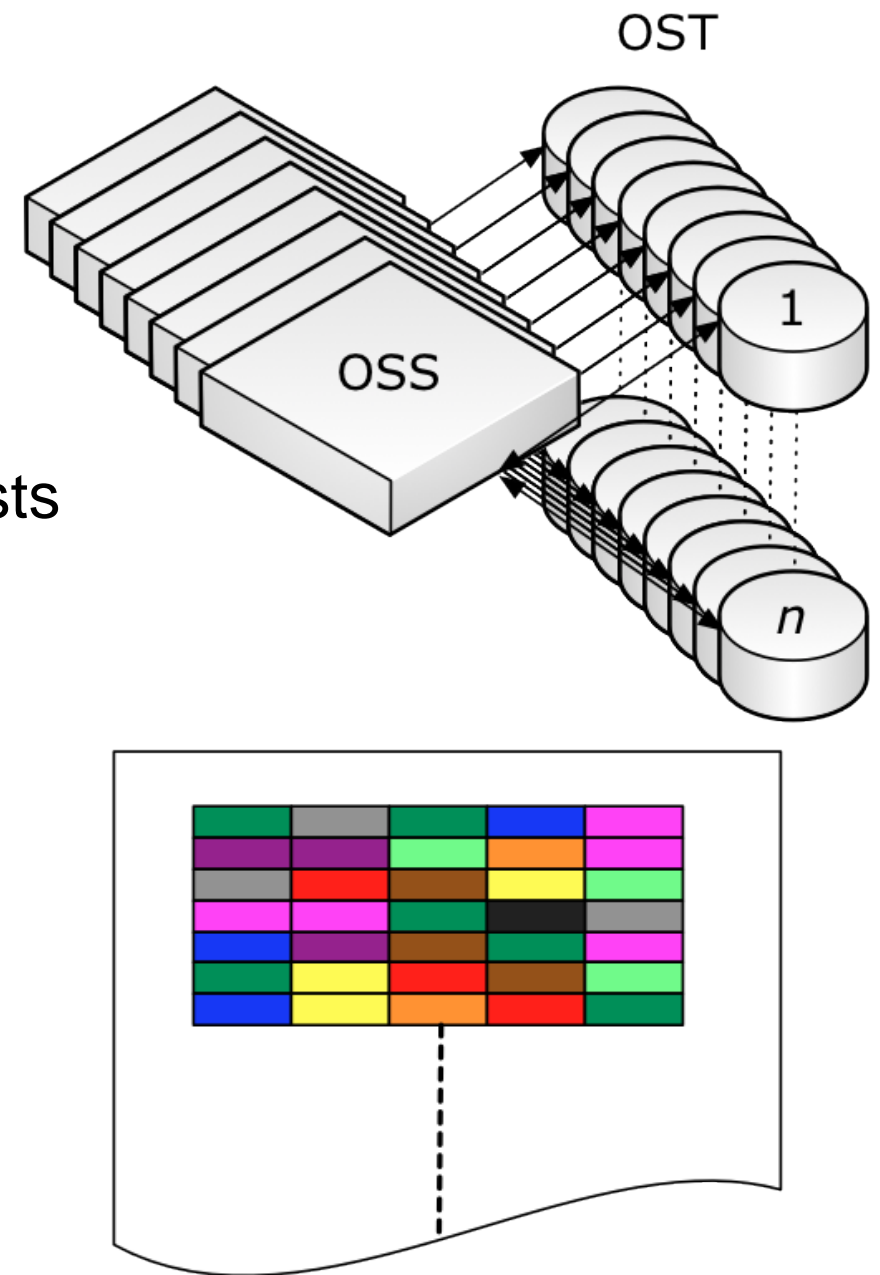
# Lustre: The Players

# Lustre Architecture - MDS

- **Metadata Server (MDS)**
  - Node(s) that manage namespace, file creation and layout, and locking.
    - Directory operations
    - File open/close
    - File status
    - File creation
    - Map of file object location
  - Relatively expensive serial atomic transactions to maintain consistency
- **Metadata Target (MDT)**
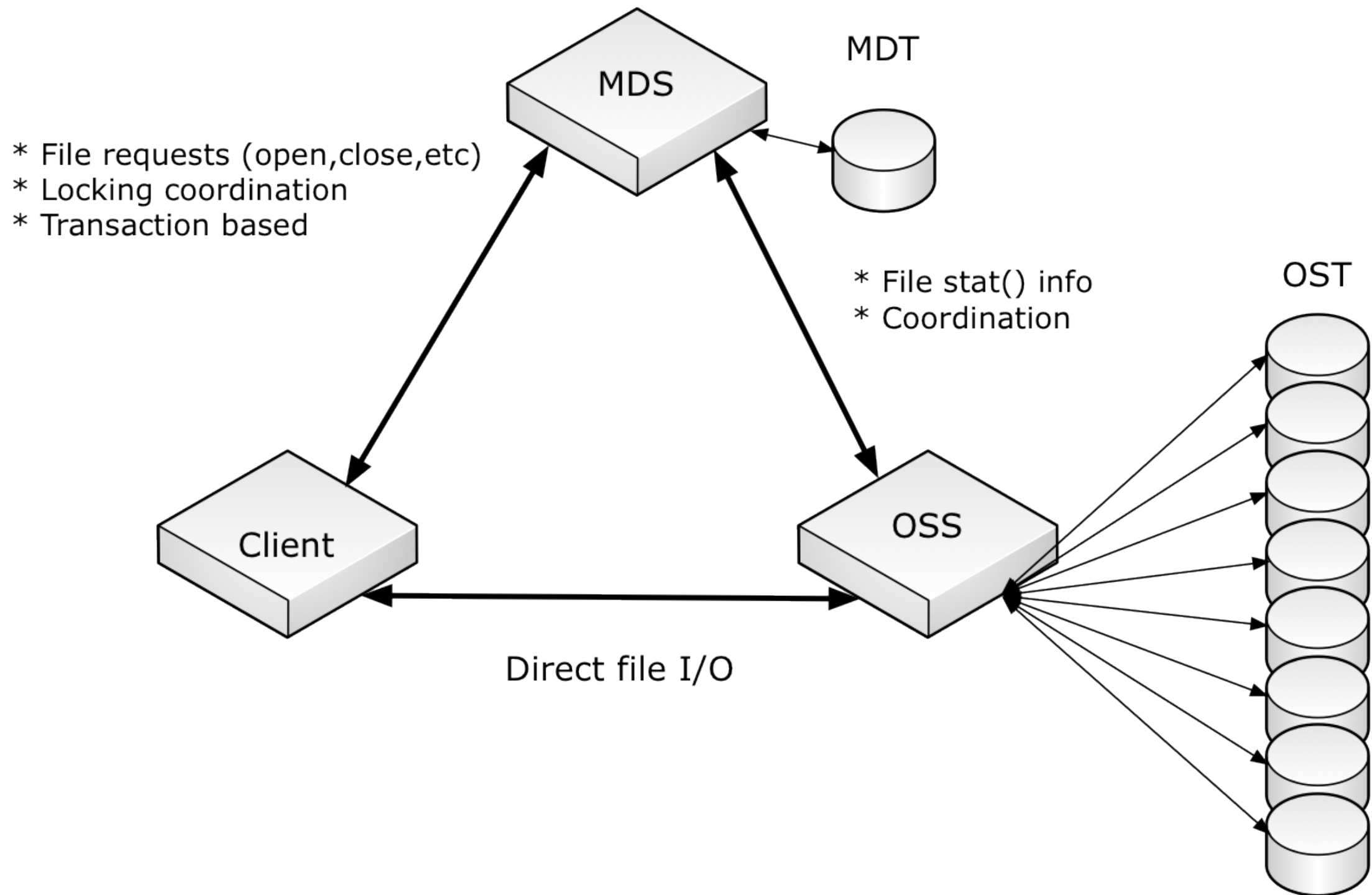  - Block device that stores metadata

# Lustre Architecture - OSS

- **Object Storage Server (OSS)**
  - Multiple nodes that manage network requests for file objects on disk.

- **Object Storage Target (OST)**
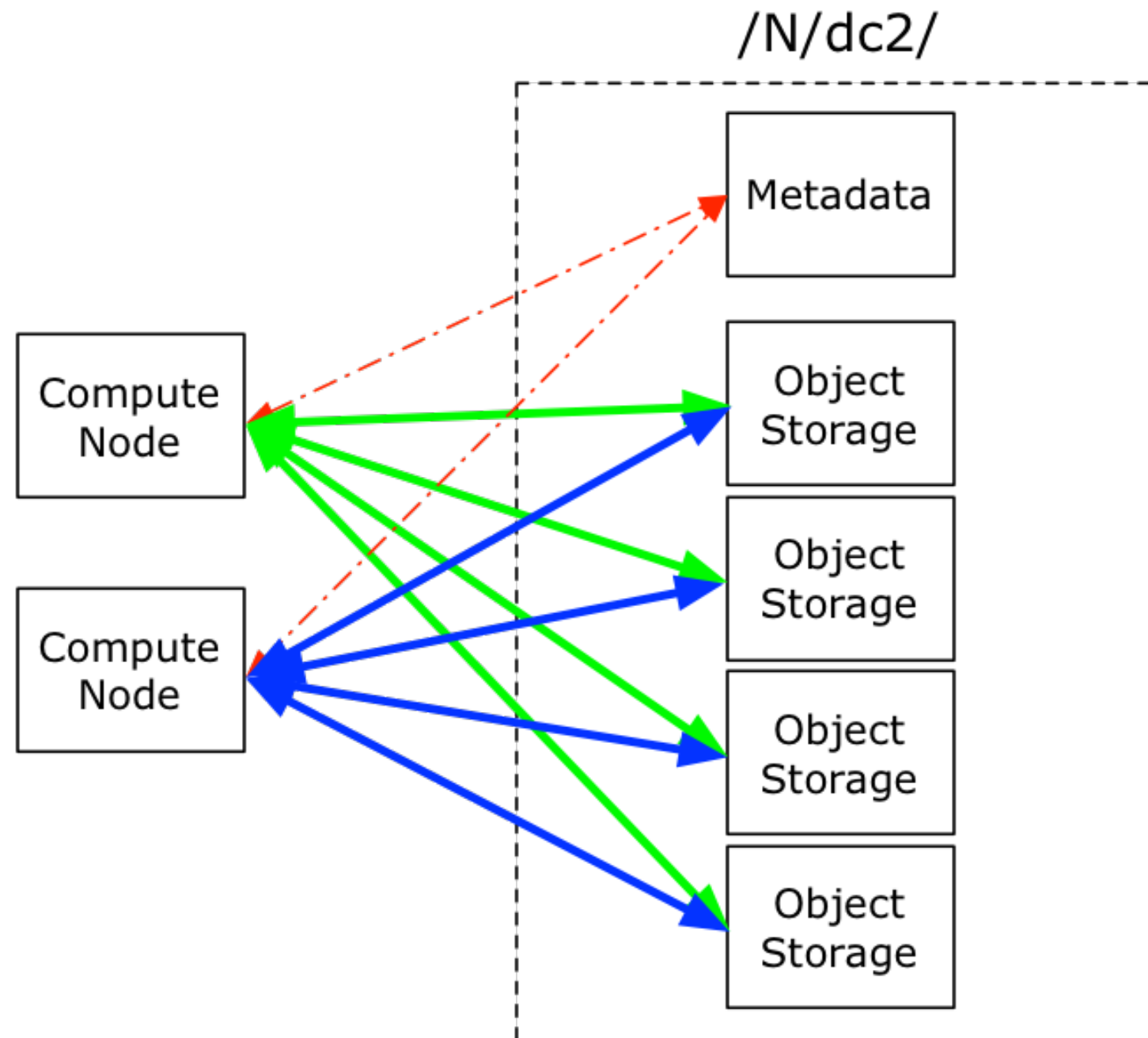  - Block device that stores file objects

# Simplest Lustre System



MDT

MDS

* File requests (open,close,etc)
* Locking coordination
* Transaction based
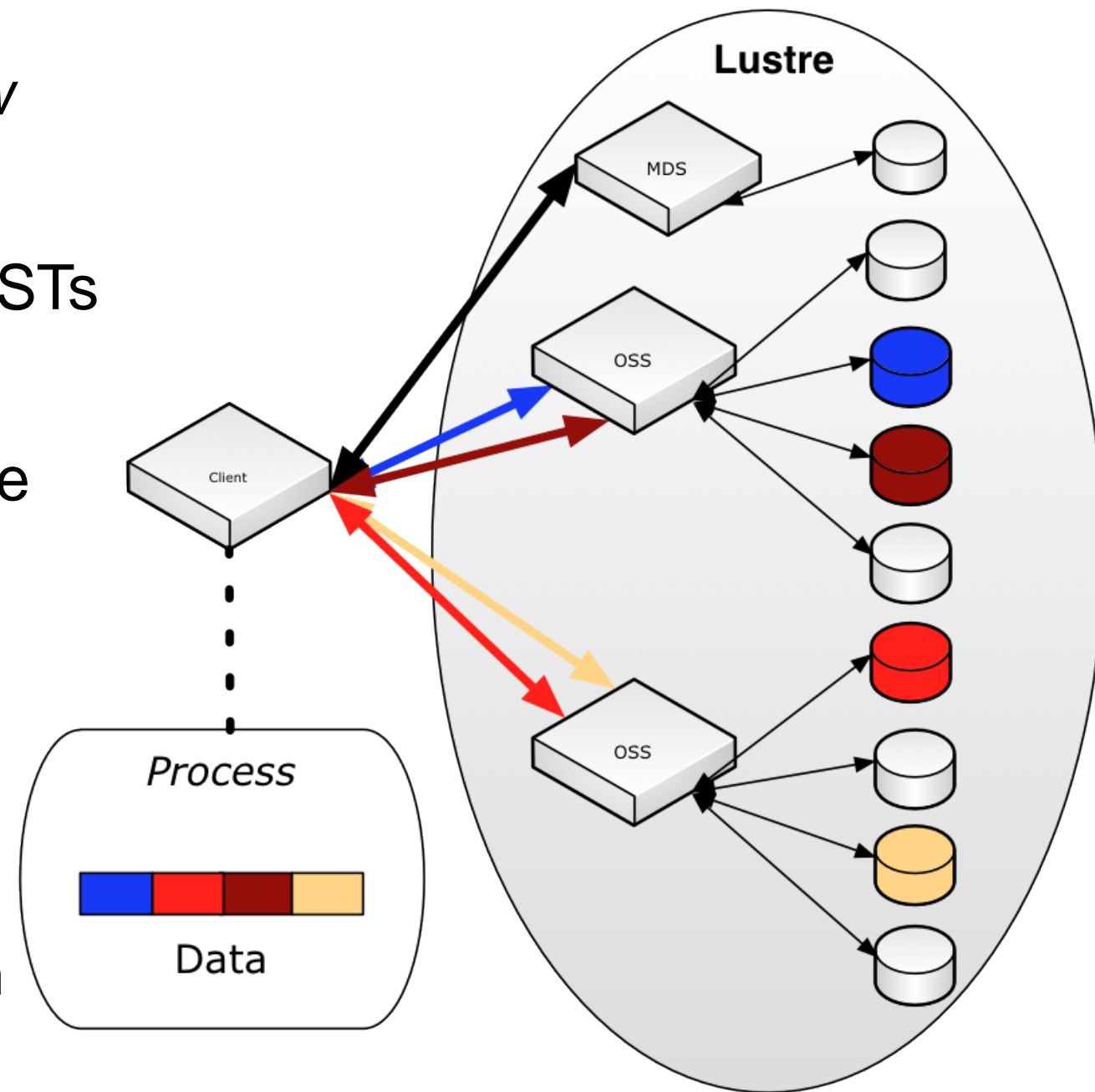
* File stat() info
* Coordination

OST

Client

OSS

Direct file I/O
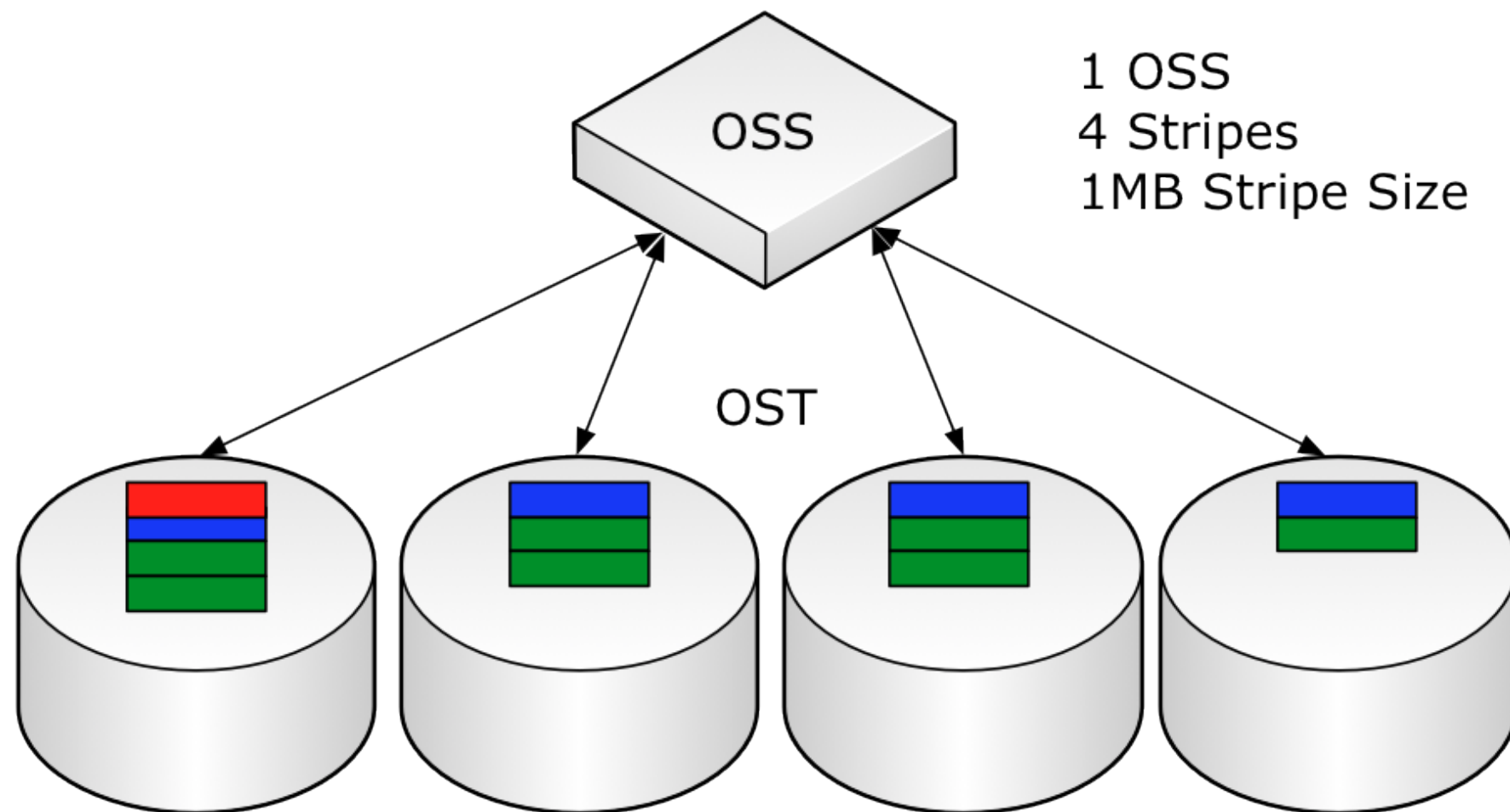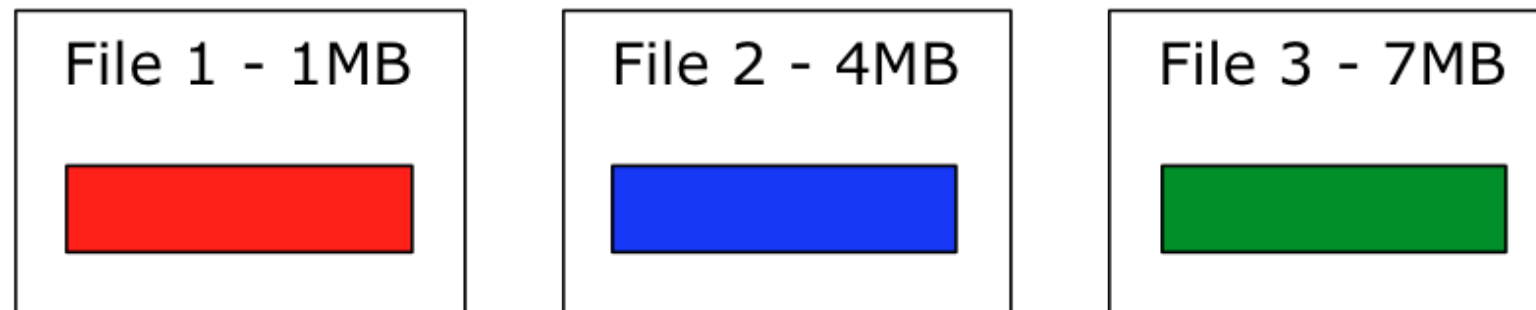
OpenSFS.

# Lustre Parallel I/O
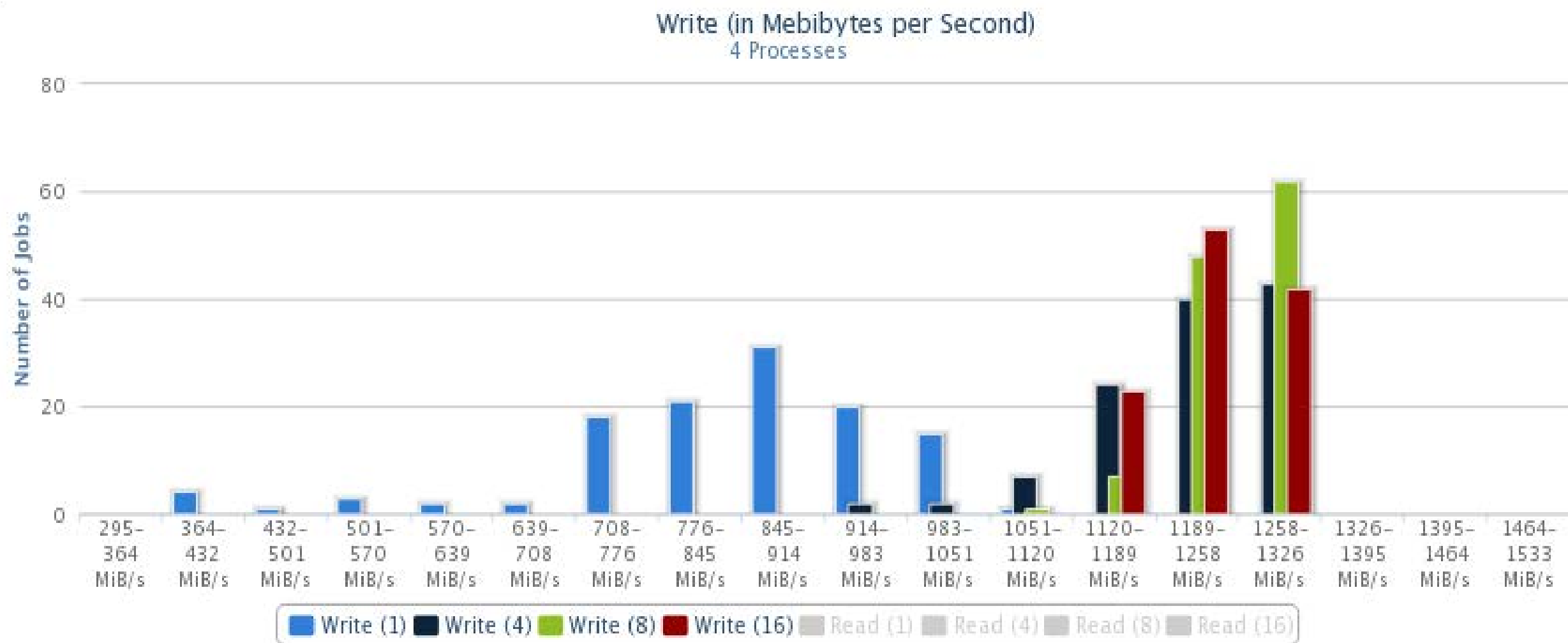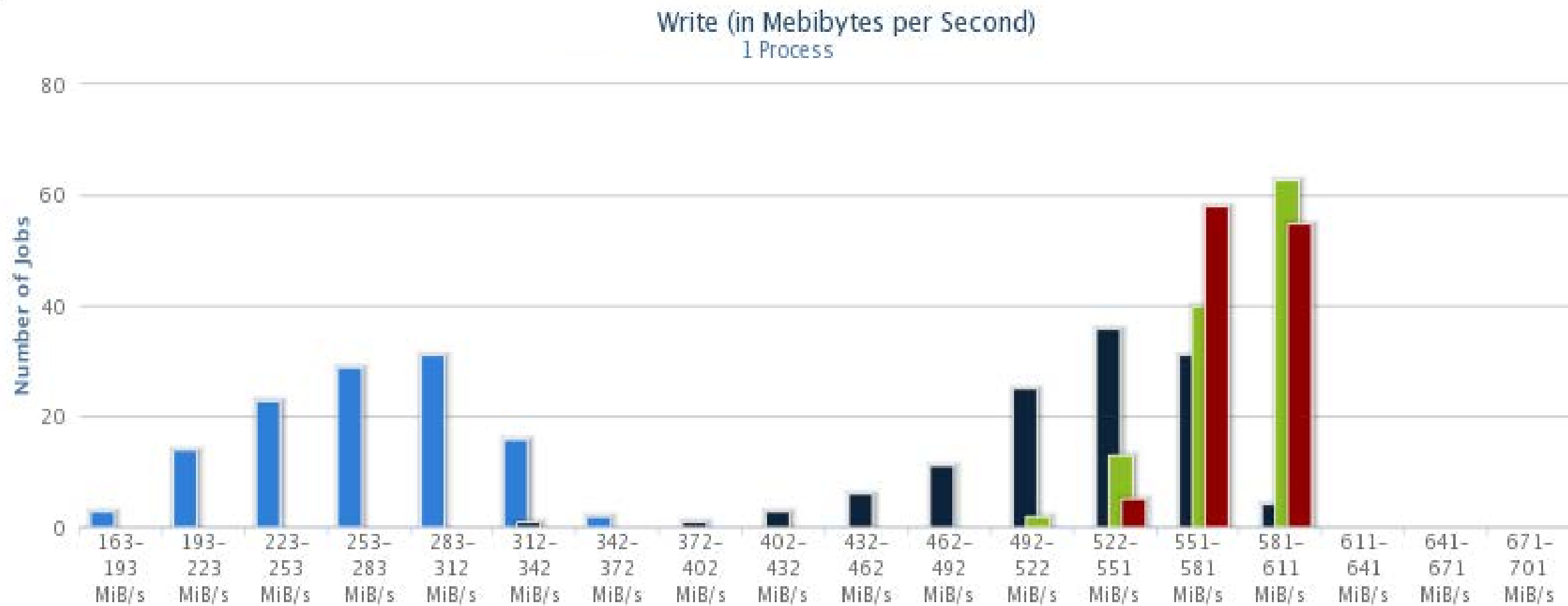
# Striping Data

- Lustre allows you to control *how* data is written, if you want
  - *Stripe* data across multiple OSTs
    - can stripe files OR directories
  - Can increase I/O performance with reading and writing
  - With DNE2 metadata can be *Striped* across multiple MDTs
- Striping analogous to RAID 0
- Default striping set by sysadmin

# Striping Example

Write (in Mebibytes per Second)
1 Process

Write (in Mebibytes per Second)
4 Processes

Write (1)　Write (4)　Write (8)　Write (16)　Read (1)　Read (4)　Read (8)　Read (16)

Read (in Mebibytes per Second)
1 Process

Read (in Mebibytes per Second)
4 Processes

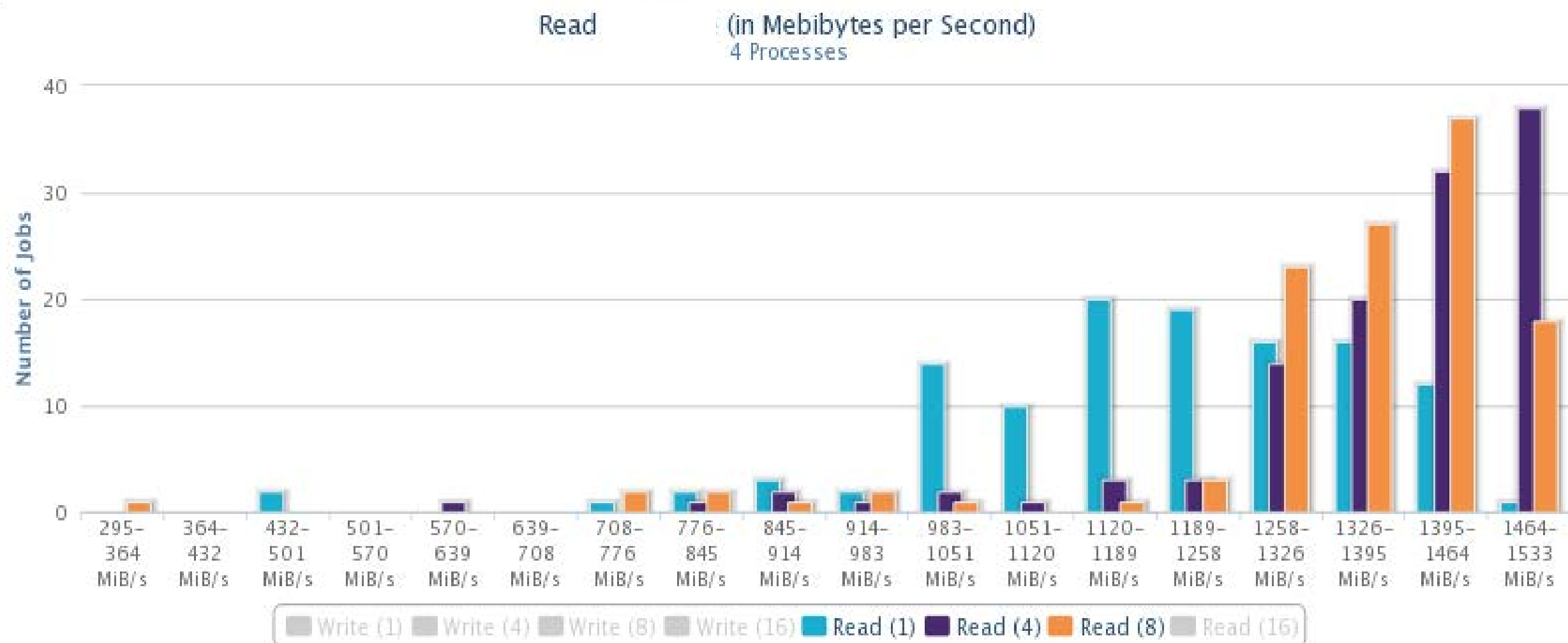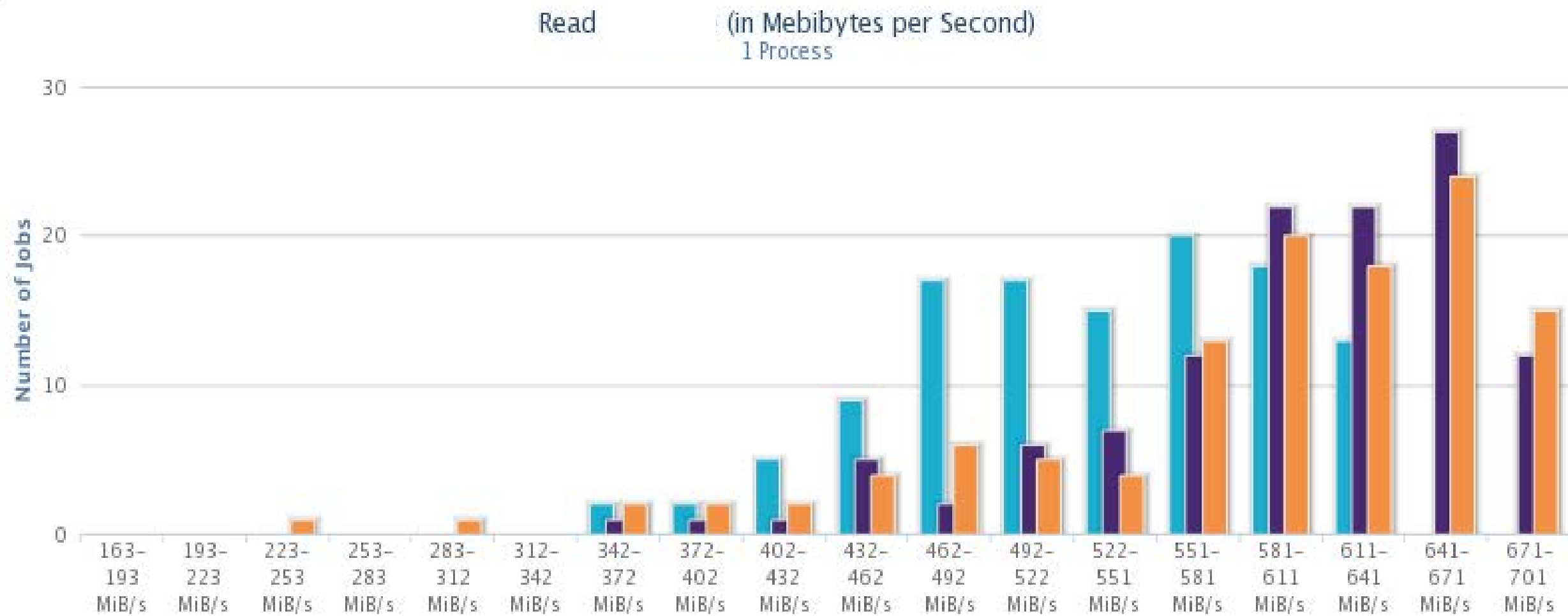Write (1)  Write (4)  Write (8)  Write (16)  Read (1)  Read (4)  Read (8)  Read (16)

# Lustre Striping

| Advantages |
| --- |
| **Bandwidth** – file objects are striped across OSTs, so bandwidth is the aggregate I/O rate |
| **File Size** – file objects striped across OST can have a total size larger than available space on any single OST |

| Disadvantages |
| --- |
| **User Overhead** – Time and thought required to understand your I/O patterns and create stripe layout for directories and files |
| **System Overhead** – Additional stripes means more OST lookups to determine the size of the file (more time) |

- Striping will not benefit *ALL* applications

OpenSFS.

# Take Home Message

## Choose Lustre!

It scales – size, speed, clients

It's open, growing, and needs your help

It gives users powerful options

Tools available to help with installation

Filets, Chops, Removes household odors

Act now and no salesman will visit your home

OpenSFS

# Thank You for Your Time and Attention

**Open Scalable File Systems, Inc.**
3855 SW 153rd Drive
Beaverton, OR 97006
Ph: 503-619-0561
Fax: 503-644-6708
admin@opensfs.org

www.opensfs.org