# Lustre Networking - an overview
# LUG 2007

# Lustre Deployment Overview



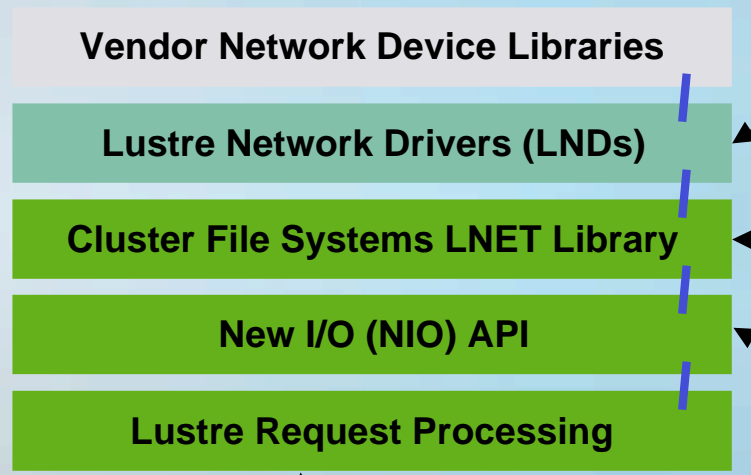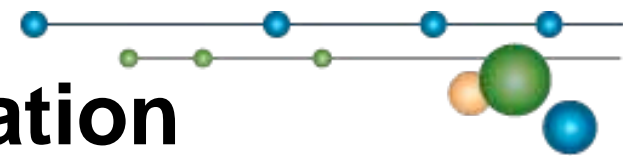Lustre Metadata Servers (MDS)

Lustre Object Storage Servers(OSS) 100's

Pool of metadata servers

MDS 1 (active)    MDS 2 (standby)

OSS 1

Elan Myrinet InfiniBand

Commodity Storage Servers

Lustre Clients 10's - 10,000's

OSS 2

OSS 3

Simultaneous support of multiple network types

Shared storage enables failover OSS

OSS 4

Router

OSS 5

GigE

OSS 6

= failover

OSS 7

Enterprise-Class Storage Arrays & SAN Fabrics

CFS

# Network features

- **Scalability - network 10,000's nodes**

- **Support for multiple networks**
  - **TCP**
  - **IB - many flavors**
  - **Elan3,4**
  - **Myricom GM, MX**
  - **Cray Seastar & RA**

- **Routing nodes between networks**

Copyright © 2006, Cluster File Systems, Inc.

# Modular Network Implementation

**Vendor Network Device Libraries**

**Lustre Network Drivers (LNDs)**

**Cluster File Systems LNET Library**

**New I/O (NIO) API**

**Lustre Request Processing**

**Support for multiple network types Including routing API**

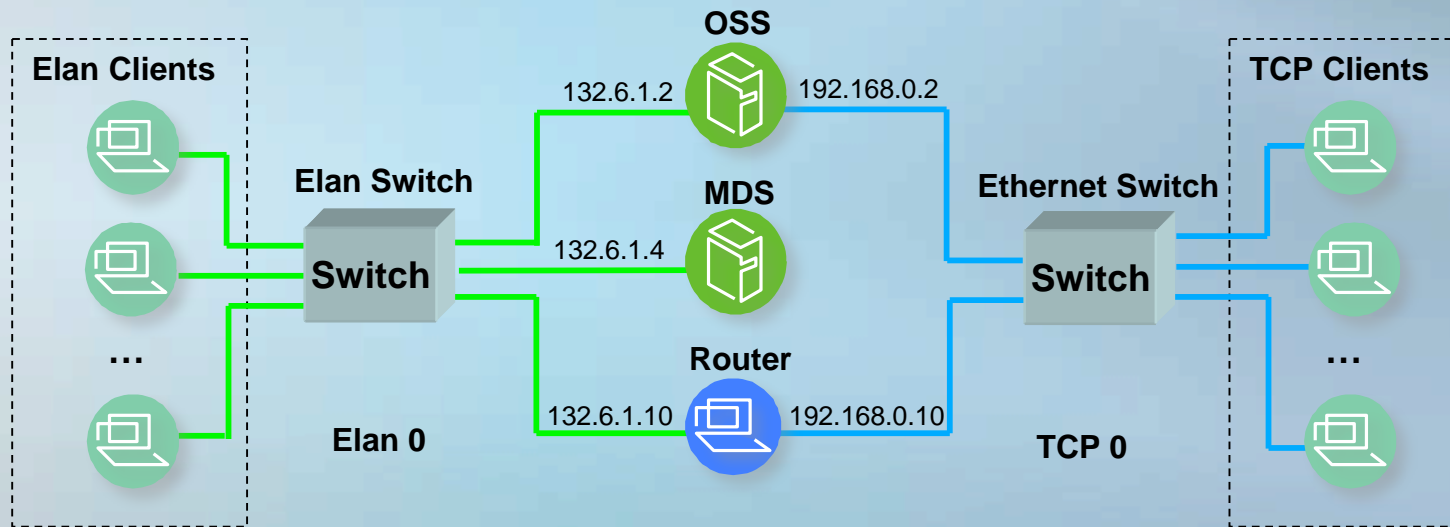**Similar to Sandia Portals with some new and different features**

**Move small and large buffers Use RDMA Generate events**

Zero-copy marshalling libraries
Service framework and request dispatch
Connection and address naming
Generic recovery infrastructure

Key:

| | Protocol |
|---|---|
| **Portable Lustre component** | |
| **Not portable** | |
| **Not supplied by CFS** | |

CFS

# Routing - an example
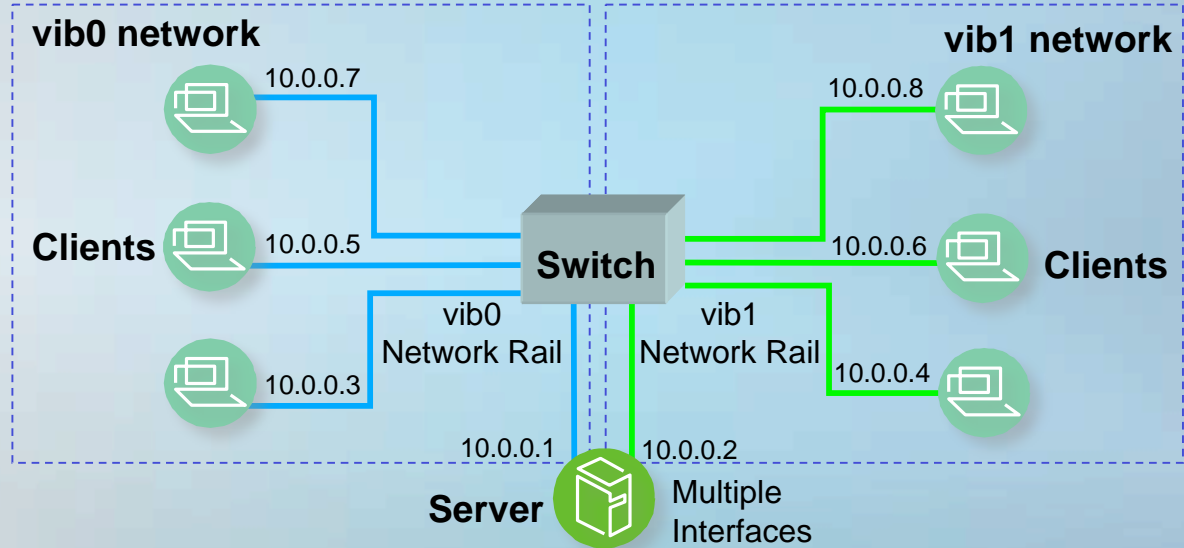


## Configuration:

options lnet 'ip2nets="tcp0 192.168.0.*; elan0 132.6.1.*"'
  'routes="tcp0 [2,10]@elan0; elan0 192.168.0.[2,10]@tcp0"'
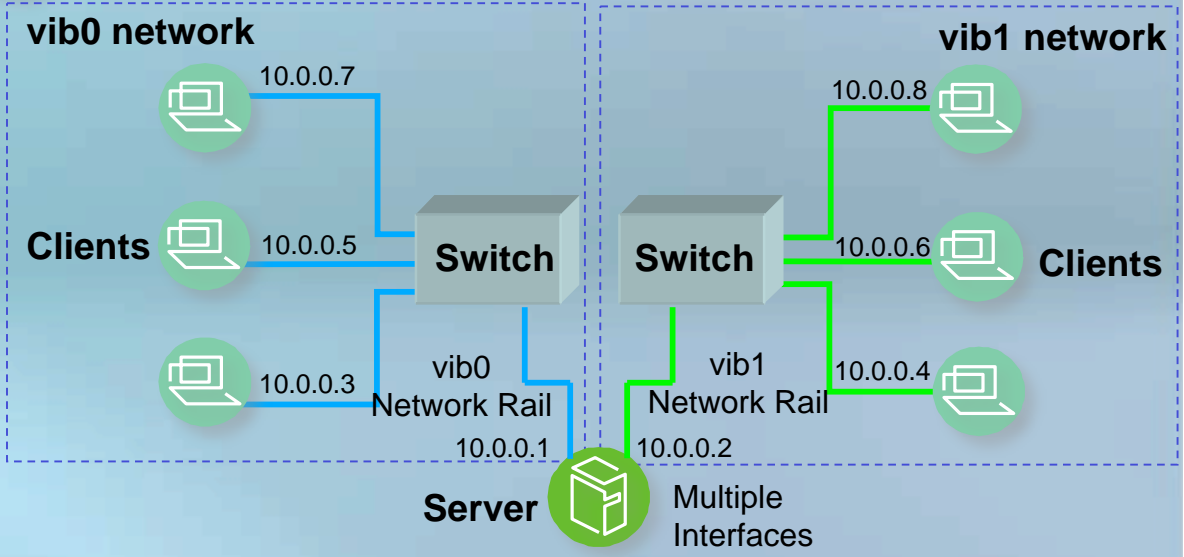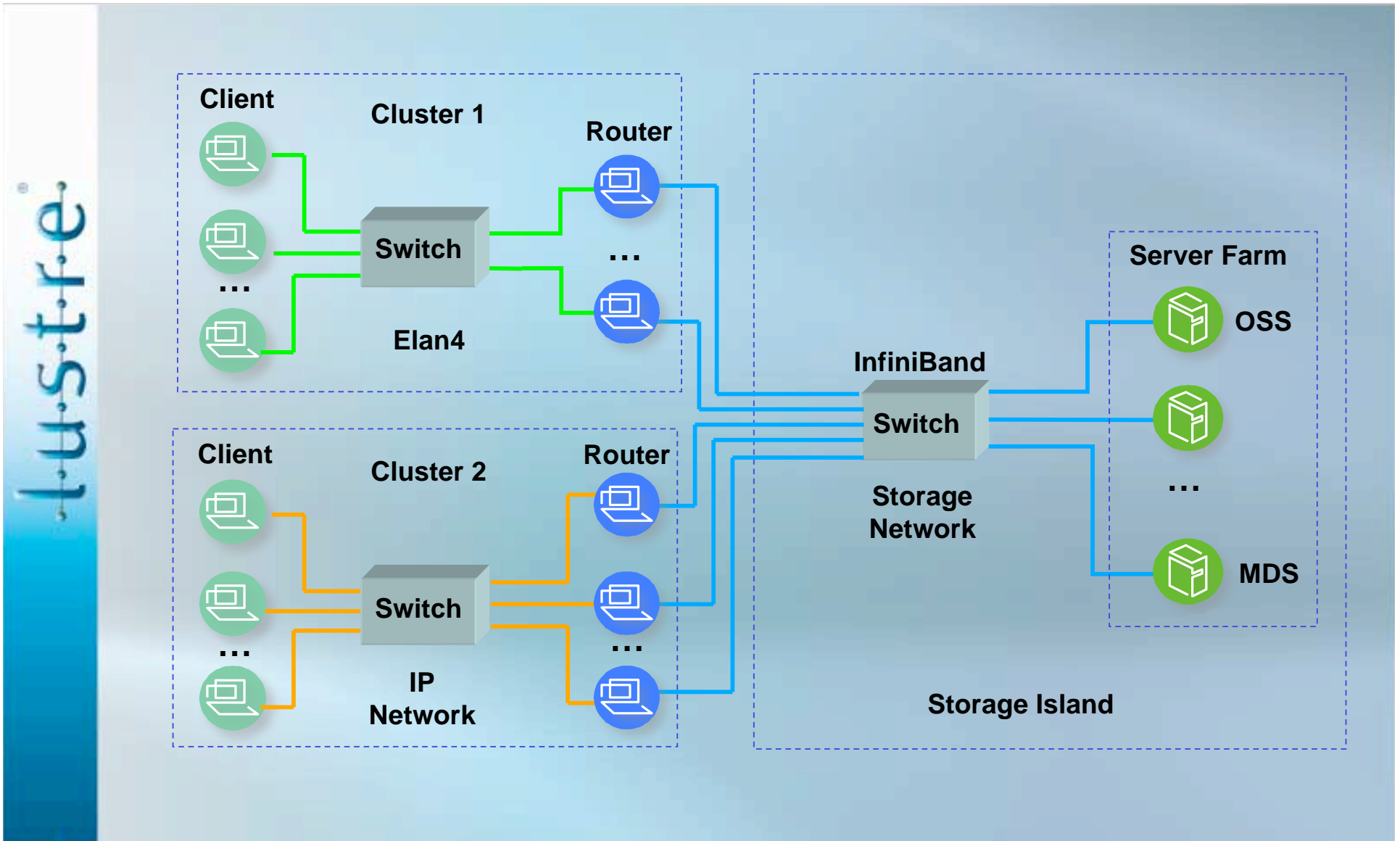
# Multiple interfaces and LNET



**vib0 network**

Clients

10.0.0.7

10.0.0.5

10.0.0.3

**Switch**

vib0 Network Rail

**vib1 network**

10.0.0.8

10.0.0.6  Clients

10.0.0.4

vib1 Network Rail

10.0.0.1

10.0.0.2

**Server**  Multiple Interfaces

**Support through:**
- **multiple Lustre networks**
- **on one or two physical networks**
- **requires clients to load balance**

**vib0 network**

Clients

10.0.0.7

10.0.0.5

10.0.0.3

**Switch**

vib0 Network Rail

**Switch**

vib1 Network Rail

**vib1 network**

10.0.0.8

10.0.0.6  Clients

10.0.0.4

10.0.0.1

10.0.0.2

**Server**  Multiple Interfaces

# Site wide file systems

Copyright © 2006, Cluster File Systems, Inc.

# Router features

- **Redundant routers**
- **Sophisticated buffer level load balancing**
- **Failed routers are avoided**
- **Failed routers are pinged & recoverable**

- **Future router features may be:**
  - Control plane - adjust policies
  - Dynamic router addition
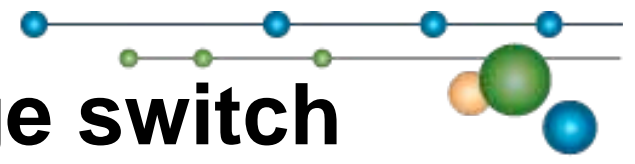  - All of these require LNET access to the management node

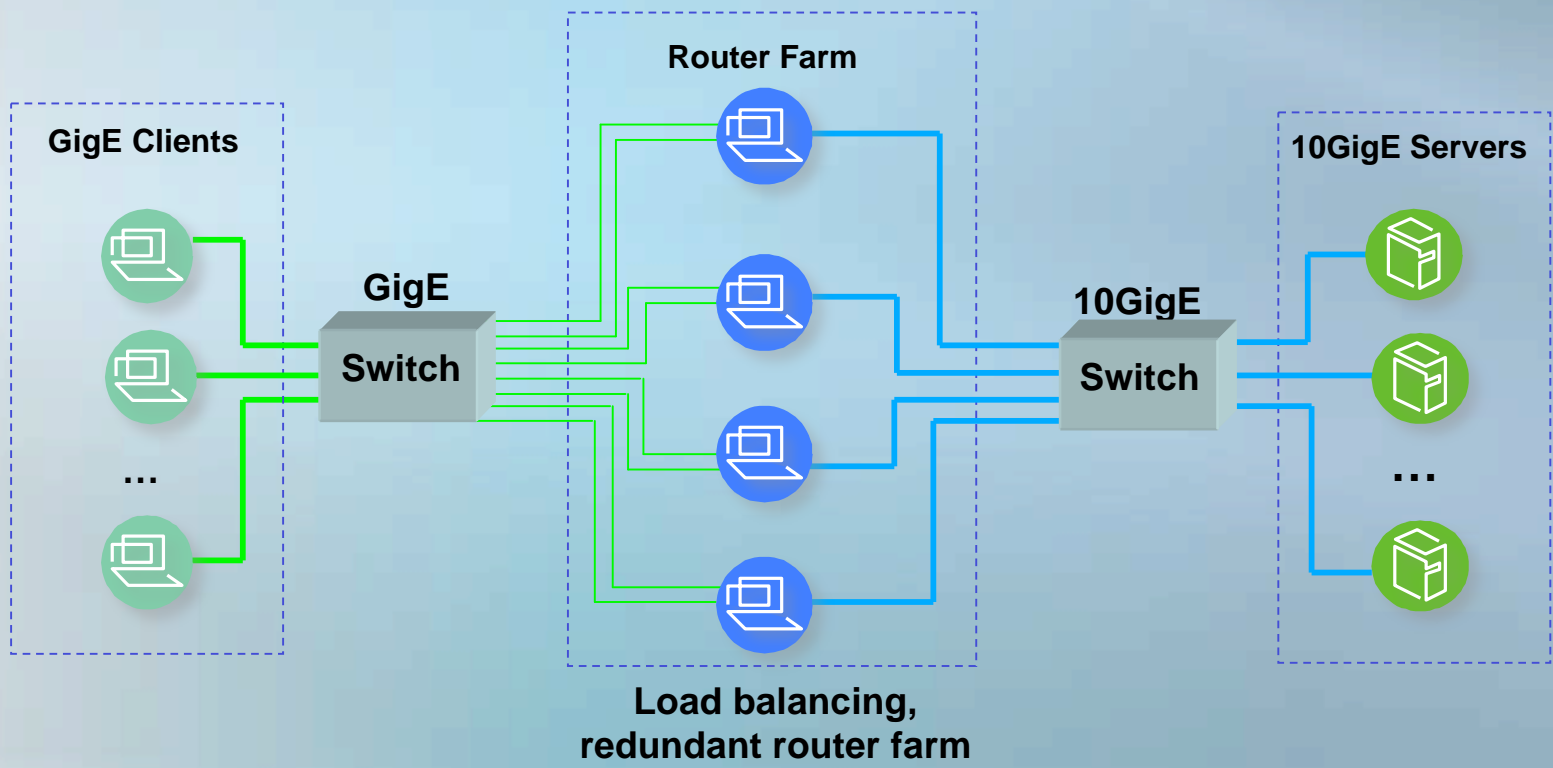Copyright © 2006, Cluster File Systems, Inc.

# Router uses

- **Site wide file systems**
  - Build a storage island
  - Add routers with multiple interface types
  - Connect clusters with different interconnects

- **Accessing fast servers**
  - Build a small fast server network
  - Attach routers for large slower client network
  - Utilize server bandwidth without switches

# Routers act as 10GigE - 1Gige switch



Router Farm

GigE Clients

GigE Switch

10GigE Servers

10GigE Switch

Load balancing,
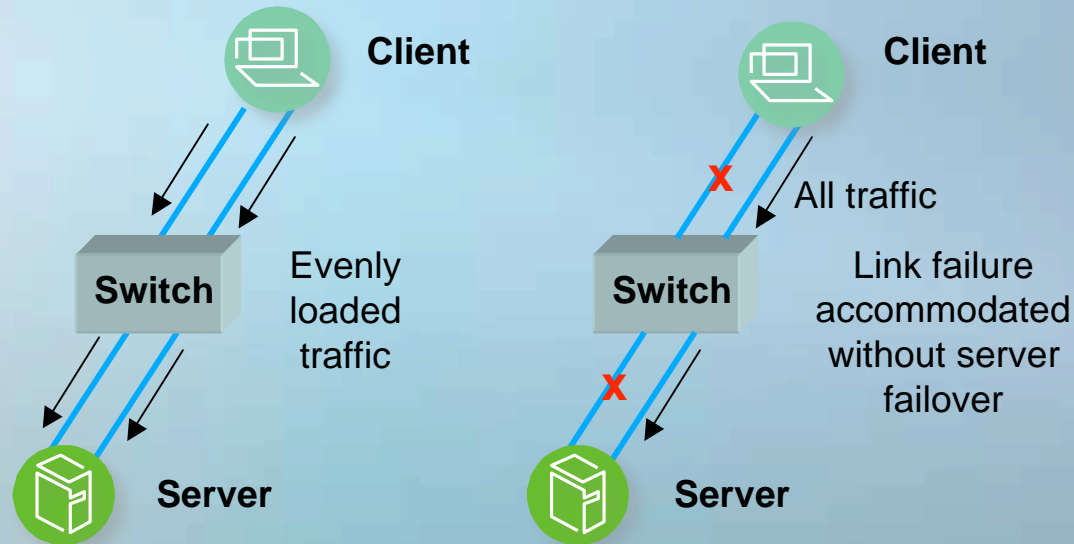redundant router farm

CFS

# Multiple interface handling

- Future work
- Desirable features
  - Link level load balancing
  - Link failover
  - N-to-K link handling

# Multiple interface features



Client

Switch

Evenly loaded traffic

Server

Client

X

All traffic

Switch

X

Link failure accommodated without server failover
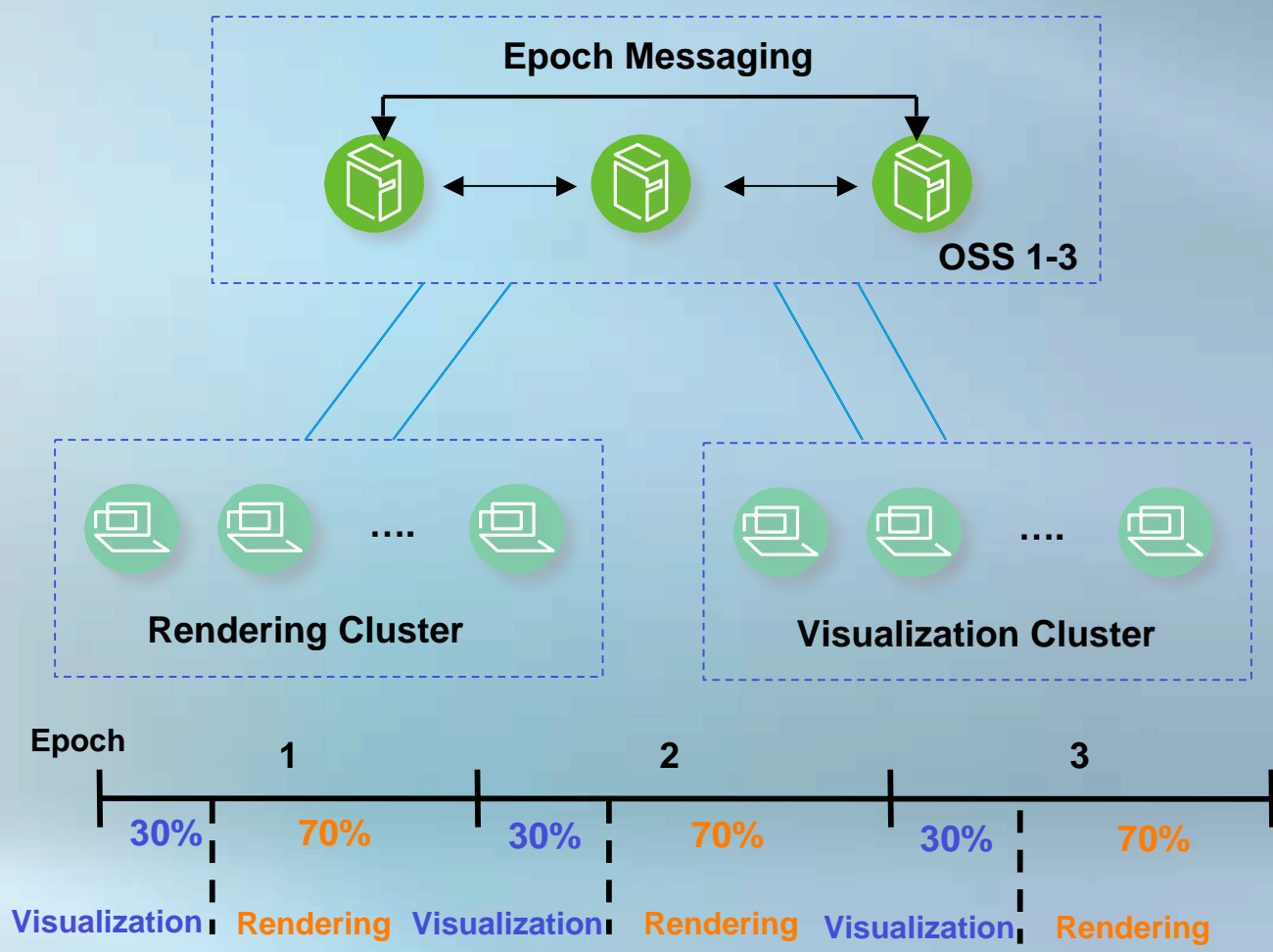
Server

# Server level load balancing

- Multiple clusters will compete for storage

- Coordinate access to the servers
  - Policy

- Lustre's task is:
  - Not to decide policy
  - Enable policies

Copyright © 2006, Cluster File Systems, Inc.

# Server level Load balancing

**Epoch Messaging**

**OSS 1-3**

**Rendering Cluster** .... **Visualization Cluster**

| Epoch | 1 | | 2 | | 3 | |
|---|---|---|---|---|---|---|
| | 30% | 70% | 30% | 70% | 30% | 70% |
| | Visualization | Rendering | Visualization | Rendering | Visualization | Rendering |

**LRS policy allocates 30% of each epoch time slice to visualization and 70% to rendering.**
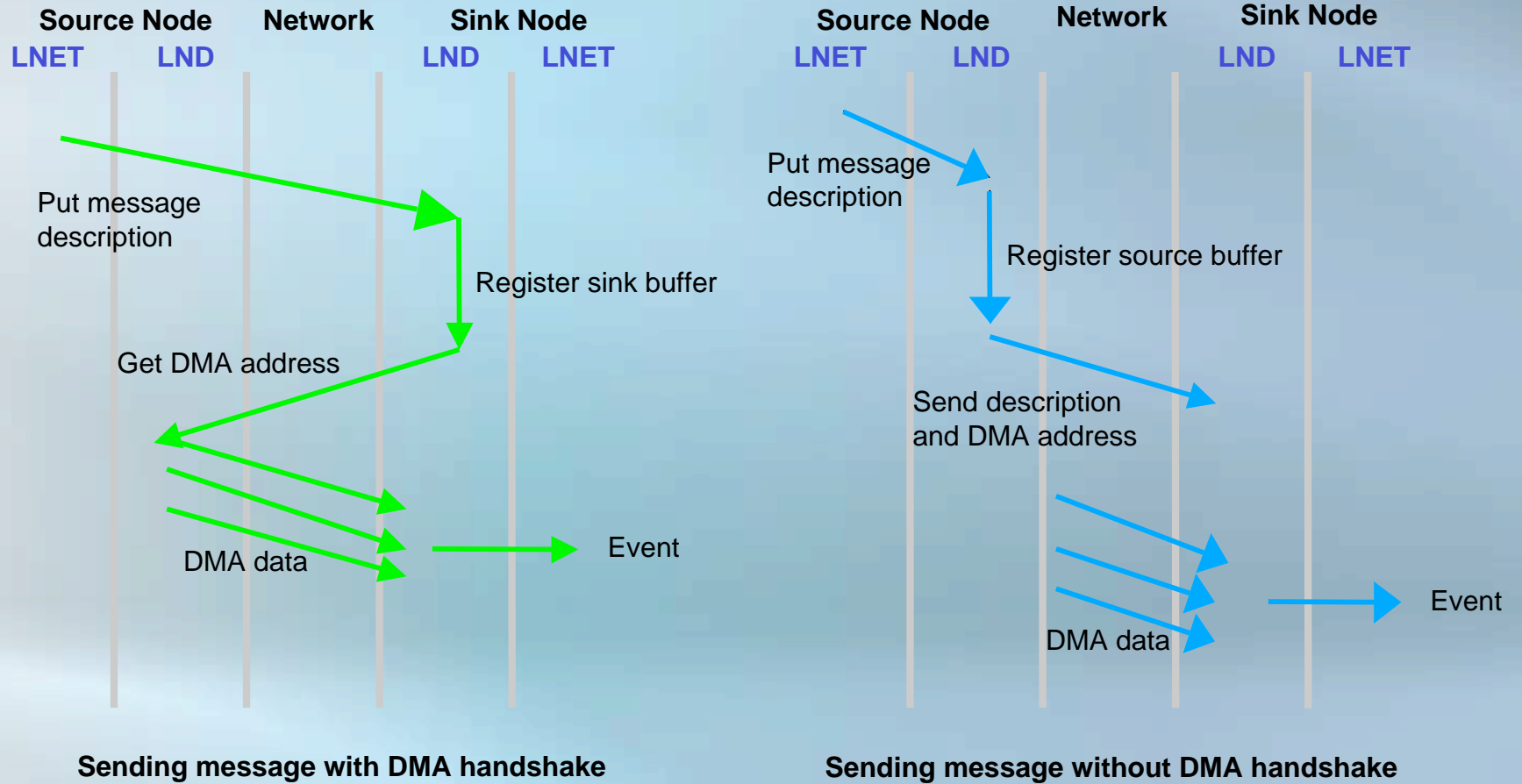
CFS

# Interrupt free asynchronous IO

- If clients register RDMA buffers before IO
  - Data can be sent or drained while they compute
  - No interrupts on the clients

- Requires LNET changes

Copyright © 2006, Cluster File Systems, Inc.

**Source Node**  **Network**  **Sink Node**
LNET    LND        LND    LNET

Put message
description

Register sink buffer

Get DMA address

DMA data                Event

**Sending message with DMA handshake**

**Source Node**  **Network**  **Sink Node**
LNET    LND        LND    LNET

Put message
description

Register source buffer

Send description
and DMA address

DMA data            Event

**Sending message without DMA handshake**

Copyright © 2006, Cluster File Systems, Inc.

# Thank you

CFS