# Lustre

## *The Inter-Galactic File System*

## Peter J. Braam

braam@clusterfs.com

http://www.clusterfilesystems.com

# Cluster File Systems, Inc

# Key requirements

- I/O throughput – 100's GB/sec
- Meta data scalability – 10,000's nodes, ops/sec, trillions of files
- Cluster recovery – simple & fast
- Storage management – snapshots, HSM
- Networking – heterogeneous networks
- Security – strong and global

**Cluster File Systems, Inc**
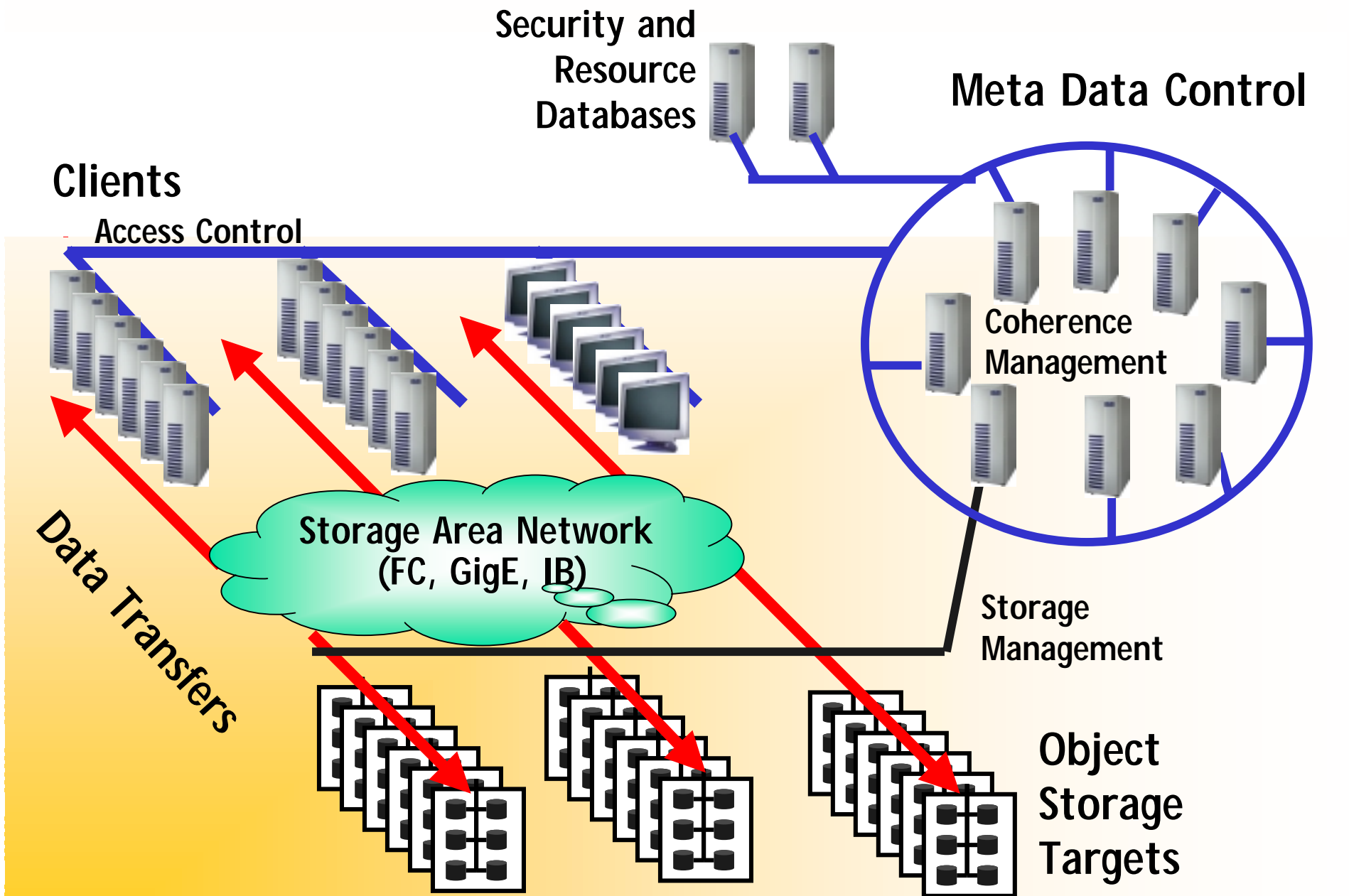
# The first 3 years…

- 1999 CMU – Seagate – Stelias Computing
- 2000 Los Alamos, Sandia, Livermore:
  - need new File System
- 2001: Lustre design to meet the SGS-FS requirements?
- 2002: things moving faster
  - Lustre on MCR (1000 node Linux Cluster – bigger ones coming)
  - Lustre Hardware (BlueArc, others coming)
  - Very substantial ASCI pathforward contract (with HP & Intel)
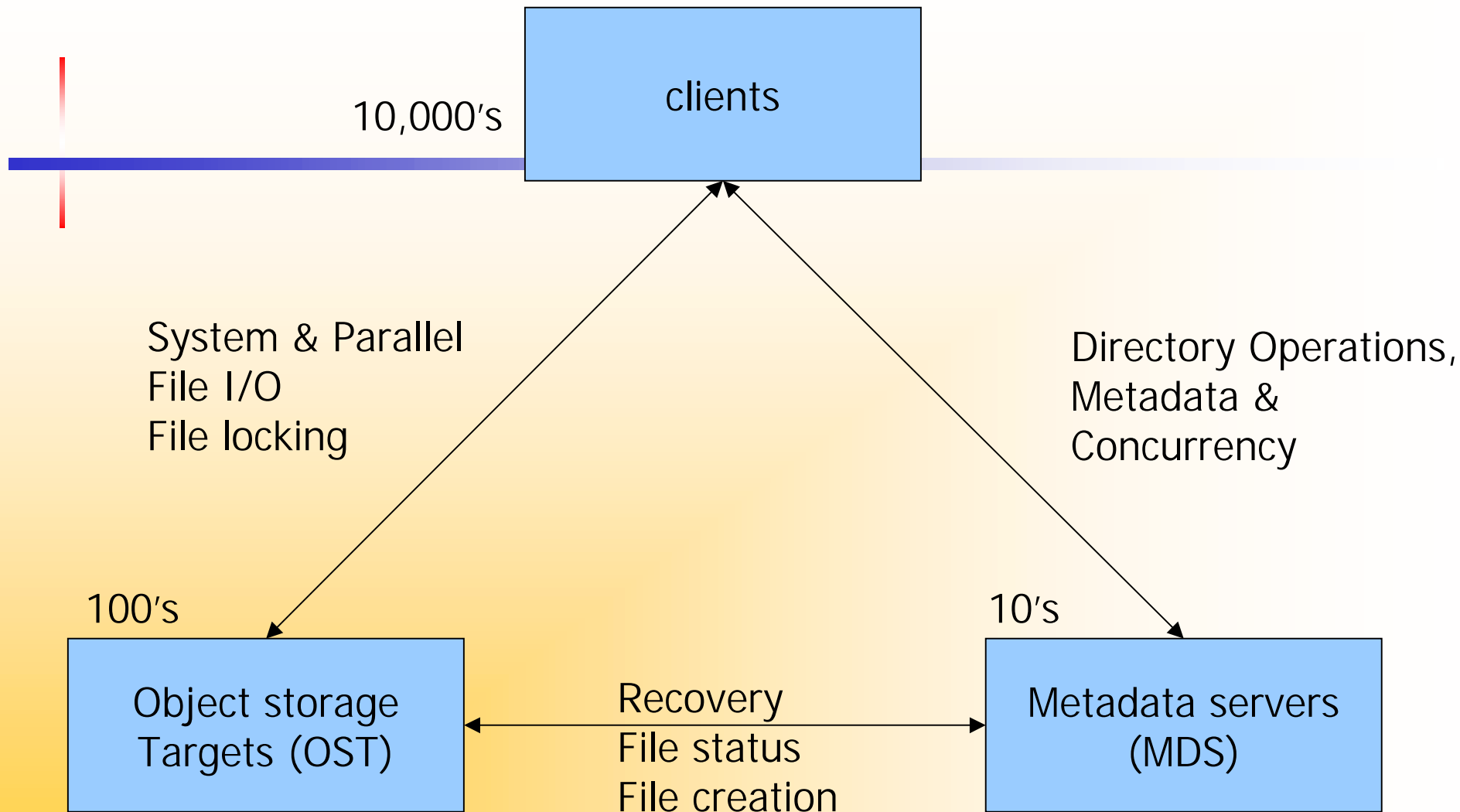
**Cluster File Systems, Inc**

# Approach

- Initially Linux focused

- Was given blank sheet

- Learn from successes

  - GPFS on ASCI White

  - TUX web server, DAFS protocol

  - Sandia Portals Networking

  - Use existing disk file systems: ext3, XFS, JFS

- New protocols

  - InterMezzo, Coda

**Cluster File Systems, Inc**

# Big Lustre picture

**Cluster File Systems, Inc**

Security and Resource Databases

Meta Data Control

Clients

Access Control

Coherence Management

Data Transfers

Storage Area Network (FC, GigE, IB)

Storage Management

Object Storage Targets

**Cluster File Systems, Inc**

Lustre System

clients

10,000's

System & Parallel
File I/O
File locking

Directory Operations,
Metadata &
Concurrency

100's

Object storage
Targets (OST)

Recovery
File status
File creation

10's

Metadata servers
(MDS)

**Cluster File Systems, Inc**

# Ingredient 1: object storage

**Cluster File Systems, Inc**
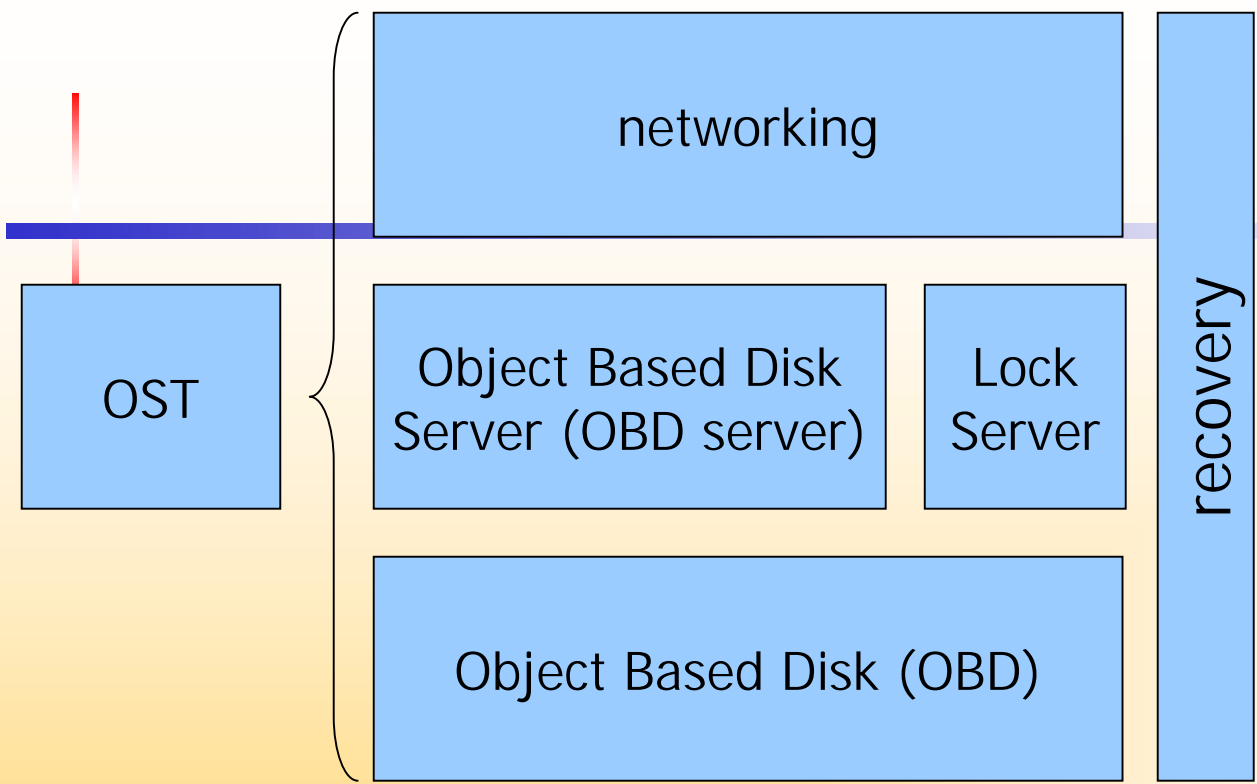
# What is Object Based Storage?

- Object Based Storage Device
  - More intelligent than block device
- Speak storage at "inode level"
  - create, unlink, read, write, getattr, setattr
  - iterators, security, almost arbitrary processing
- So...
  - Protocol allocates physical blocks, no names for files
- Requires
  - Management & security infrastructure

**Cluster File Systems, Inc**

networking

OST

Object Based Disk
Server (OBD server)

Lock
Server

recovery

Object Based Disk (OBD)

alternatives

Ext2 OBD
(raw inodes)

OBD Filter

File system
Ext3, Reiser, XFS, JFS,…

**Object
Storage
Target**

**Cluster File Systems, Inc**

# How does object storage help?

**Cluster File Systems, Inc**

# File – I/O

- Open file on metadata system
- Get information
  - What objects
  - What storage controllers
  - What part of the file
  - Striping pattern
- Use connection to storage controllers you need
  - Do logical object writes to OST
  - From time to time OST updates MDS with new file sizes

**Cluster File Systems, Inc**

# I/O bandwidth requirements

- Required: 100's GB/sec

- Consequences:
  - Saturate 100's – 1000's of storage controllers
  - Block allocation must be spread over cluster
  - Lock management must be spread over cluster

- This almost forces object storage controller approach

**Cluster File Systems, Inc**

# Ingredient 3: metadata handling
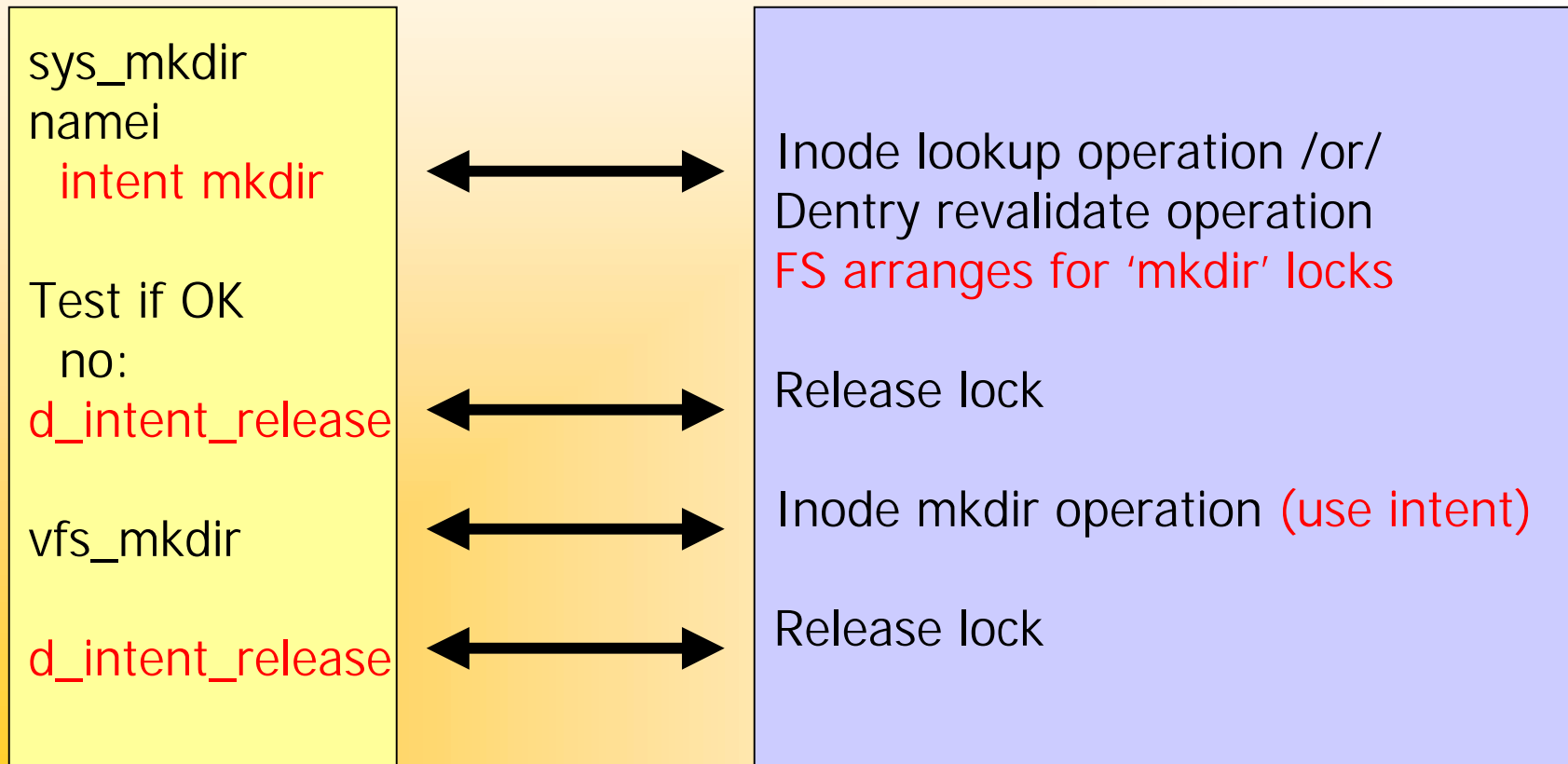
**Cluster File Systems, Inc**

# Intent based locks & Write Back caching

- Protocol adaptation between clients and MDS

- Low concurrency - write back caching

  - On client in memory updates with delayed replay on MDS

- High concurrency

  - Want single network request per transaction, no lock revocations

  - Intent based locks – lock includes all info to complete transaction

**Cluster File Systems, Inc**

# Linux VFS changes: intent lookups

**VFS**

**FS**

sys_mkdir
namei
  intent mkdir

Test if OK
 no:
d_intent_release

vfs_mkdir

d_intent_release

Inode lookup operation /or/
Dentry revalidate operation
FS arranges for 'mkdir' locks

Release lock

Inode mkdir operation (use intent)

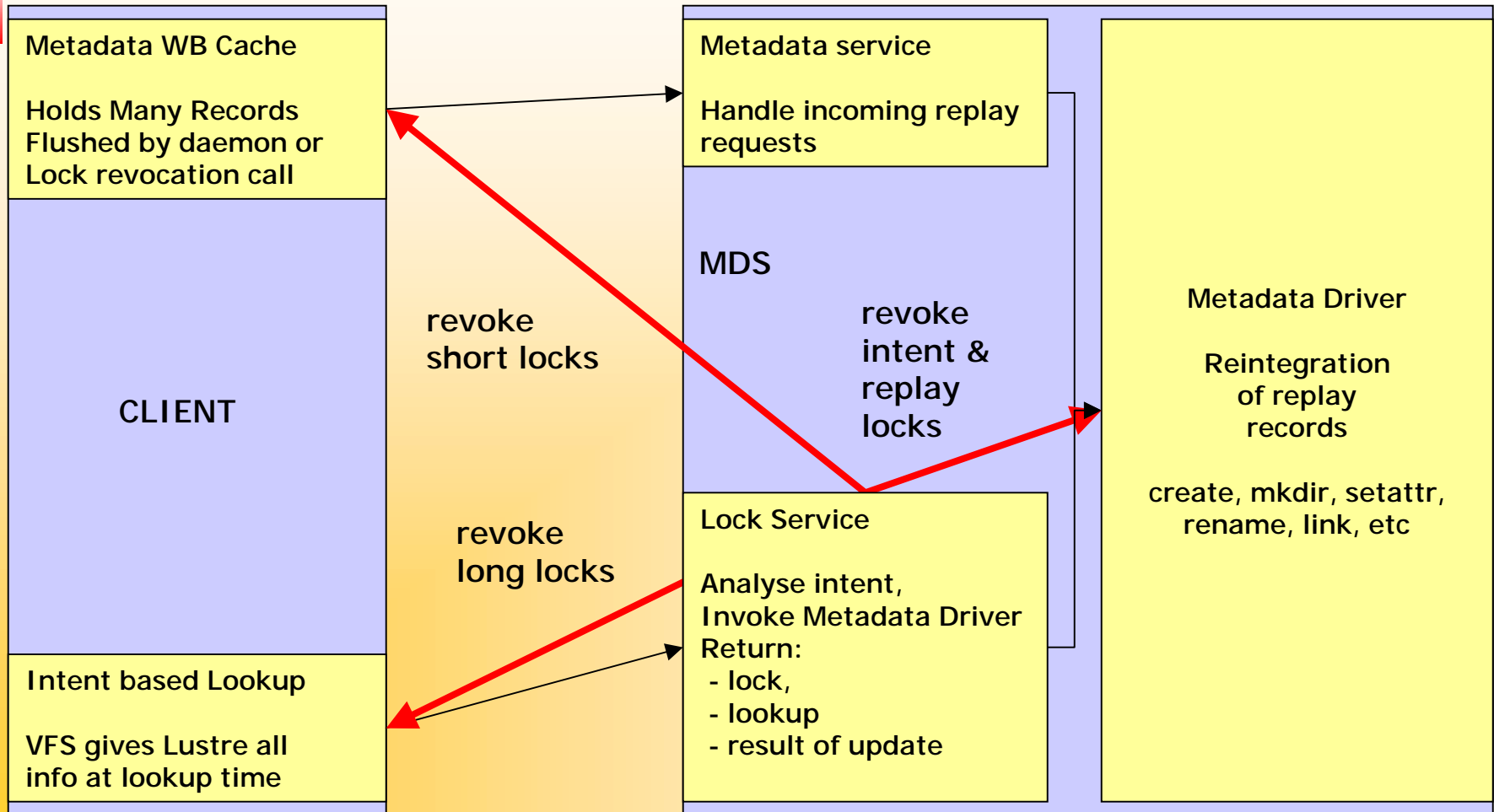Release lock

**Cluster File Systems, Inc**

# Two types of metadata locks:

- **Long locks –**
  - Lock whole pathname, help with concurrency
  - e.g. locking the root directory is BAD
    - so lock /home/peter & /home/phil separately
- **Short Locks**
  - Lock a directory subtree -help for delegation
  - e.g. a single lock on /home/phil is GOOD

**Cluster File Systems, Inc**

# Metadata updates

**Metadata WB Cache**

**Holds Many Records**
**Flushed by daemon or**
**Lock revocation call**

**CLIENT**

**Intent based Lookup**

**VFS gives Lustre all**
**info at lookup time**

revoke
short locks

revoke
long locks

**Metadata service**

**Handle incoming replay**
**requests**

**MDS**

revoke
intent &
replay
locks

**Lock Service**

**Analyse intent,**
**Invoke Metadata Driver**
**Return:**
**- lock,**
**- lookup**
**- result of update**

**Metadata Driver**

**Reintegration**
**of replay**
**records**

**create, mkdir, setattr,**
**rename, link, etc**

**Cluster File Systems, Inc**

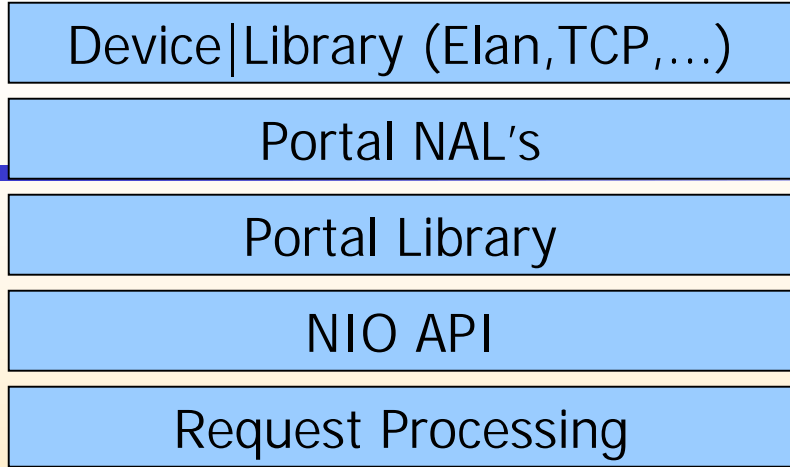# Ingredient 3: Storage Networking

**Cluster File Systems, Inc**

# Lustre networking

- Currently runs over
  - TCP,
  - Quadrics
  - Myrinet (almost)
- Other networks we are looking at:
  - SAN's
  - I/B
  - NUMA interconnects (@ GB/sec)
  - SCTP

**Cluster File Systems, Inc**

# Portals

- Sandia Portals message passing
    - simple message passing API
    - support for remote DMA
    - support for plugging in device support
    - Network Abstraction Layers

**Cluster File Systems, Inc**

Device|Library (Elan,TCP,…)

Portal NAL's

Portal Library

NIO API

Request Processing

now: Elan & IP
soon: Sandia, GM

Sandia's API
CFS improved impl.

Move small & large buffers
Generate events

0-copy marshalling libraries
service framework
client request dispatch
connection & address naming
generic recovery infrastructure

# Lustre Network Stack

**Cluster File Systems, Inc**

# Ingredient 4: Storage Management

**Cluster File Systems, Inc**

# Components of OB Storage

- Storage Object Device Drivers
  - **Class driver** – attach driver to interface
  - **Targets, clients** – remote access
  - **Direct drivers** – to manage physical storage
  - **Logical drivers** – for intelligence & storage management
  - **Object storage "applications"** – eg. the file system
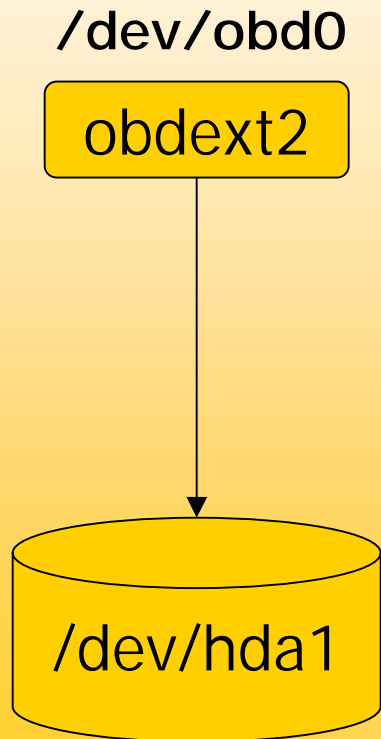
**Cluster File Systems, Inc**

# Examples of logical modules

- Storage management:
    - System software, trusted
    - Often inside the standard data path,
    - Often involves iterators
    - Eg: security, snapshots, versioning data migration, raid
- Lustre offers active disks
    - almost arbitrary intelligence can be loaded into OST driver stack
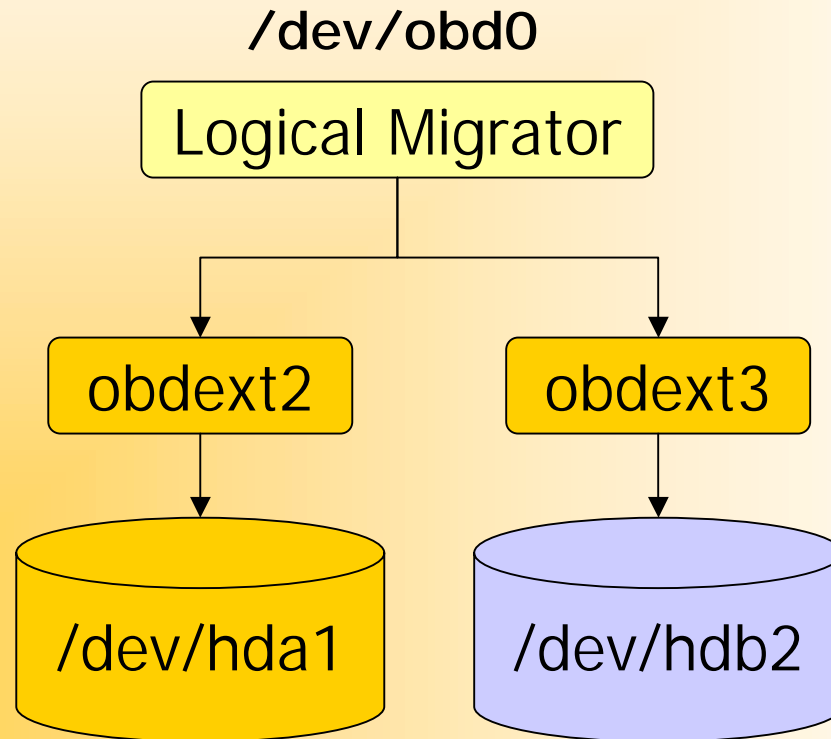
**Cluster File Systems, Inc**

# Example of management: hot data migration:

**Key principle**: dynamically switch object device types

Before…                    During…                    After…

/dev/obd0                  /dev/obd0                  /dev/obd0

| obdext2 |        | Logical Migrator |        | obdext3 |

                    | obdext2 |    | obdext3 |

/dev/hda1          /dev/hda1    /dev/hdb2          /dev/hdb2

**Cluster File Systems, Inc**

# Conclusions

- We think Lustre can run well on BlueGene

**Cluster File Systems, Inc**

# Cluster File Systems

- Small scale service company
  - Open Source development
  - contract work for Government labs
  - some consulting and collaboration with industry
- Extremely specialized and extreme expertise
  - we only do file systems and storage
- Investments etc
  - Please visit "Save the Children"
  - no thank you – it's perfectly possible to go forward without

**Cluster File Systems, Inc**