# Lustre: building a cluster file system for 1,000 node clusters

Phil Schwan

phil@clusterfs.com
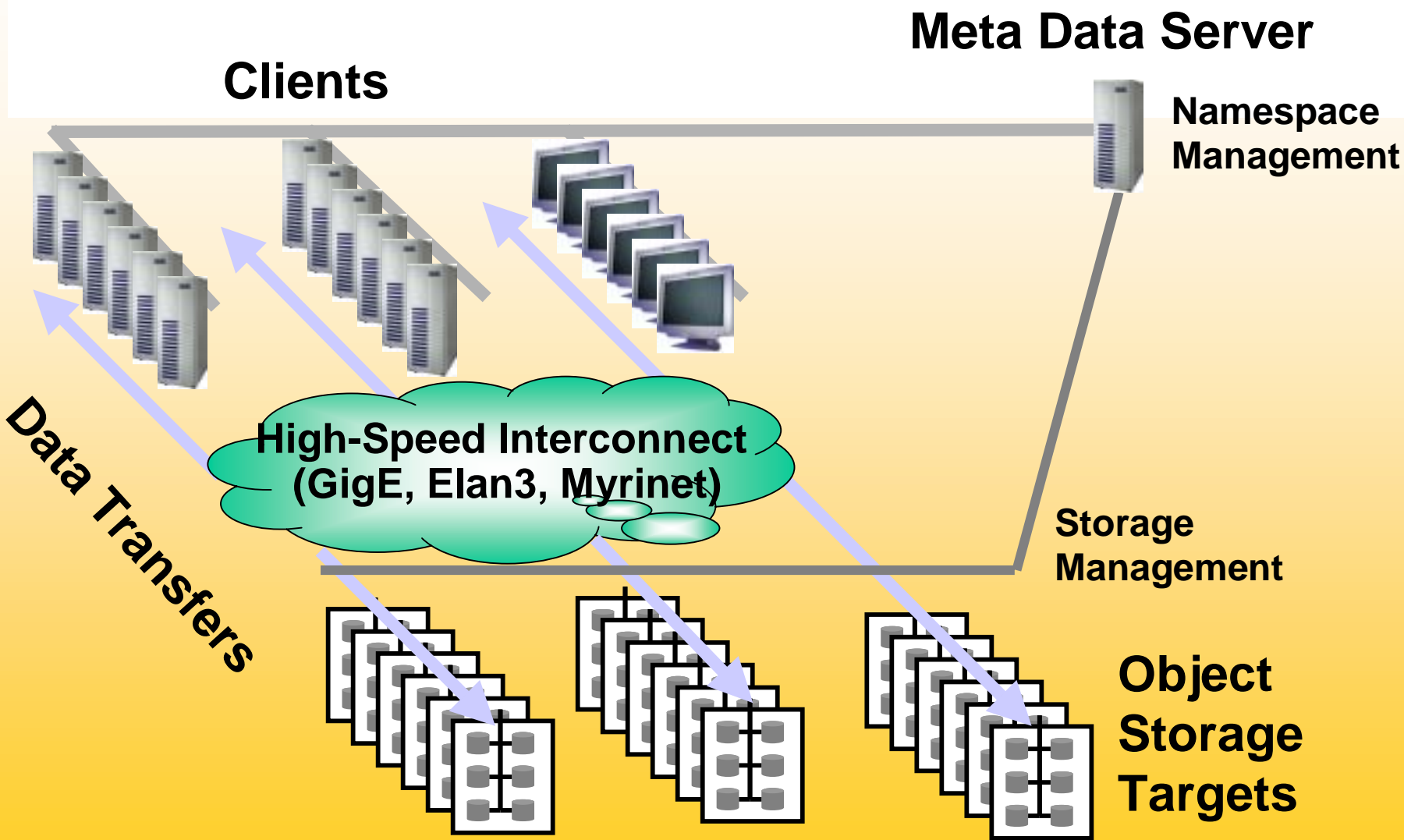http://www.clusterfs.com

**Cluster File Systems, Inc**

# Topics

- 60-second Lustre introduction
- A slice of what went well
- One critical mistake
- Questions
- Time permitting: on deck for 2003-2004

**Cluster File Systems, Inc**

# (very) Basics

- GPLed cluster file system for Linux
  - Stable on 2.4.x, making rapid progress for 2.6
- Aims for POSIX compliance
- Layering of object protocols
- Distributed lock manager
- Usually separate metadata and file data servers

**Cluster File Systems, Inc**

Clients

Meta Data Server

Namespace Management

Data Transfers

High-Speed Interconnect (GigE, Elan3, Myrinet)

Storage Management

Object Storage Targets

Cluster File Systems, Inc

# What went well

Distributed Lock Manager

**Cluster File Systems, Inc**

# Simplicity

- Based on VAX DLM concepts

- Built from scratch

- Why not use the IBM DLM?
  - 1/10th of the size (4,000 lines of code)
  - We don't need most complicated features
  - We do need different extensions

**Cluster File Systems, Inc**

# DLM extensions

- Extent locks
  - A new lock type, with an extra field
  - Policy for automatic extent growth
  - Common case: one user, one lock
- Lots of file systems don't manage this well
  - One lock per file – no concurrency
  - One lock per block/page – billions of locks
  - Locks?  Where?

**Cluster File Systems, Inc**

# DLM extensions part 2

- Intent locking
    - Allows the DLM to make policy decisions
    - Grant a lock in a low-concurrency situation
    - Perform the operation in a high-concurrency situation
- More on intent locks coming up

**Cluster File Systems, Inc**

# What went well

Scaling metadata to 1,000 nodes

**Cluster File Systems, Inc**

# Scaling metadata to 1,000 nodes

- Consider: 2,000 processes on 1,000 nodes
- All create one or more files simultaneously
- In the same directory

- This is not contrived—some LLNL science runs do this every hour

**Cluster File Systems, Inc**

# Metadata option #1: lock the directory

- Take a write lock on the parent directory
- Check to see if the file exists
- Add the new directory entry


- Very efficient for the single-user case
- Easy to implement: mimics the VFS code
- A complete disaster for our 1,000-node example

**Cluster File Systems, Inc**

# Metadata option #2: raw calls

- Execute operations entirely on the server
- Don't return locks to clients, only a status code


- Avoids lock ping-pong
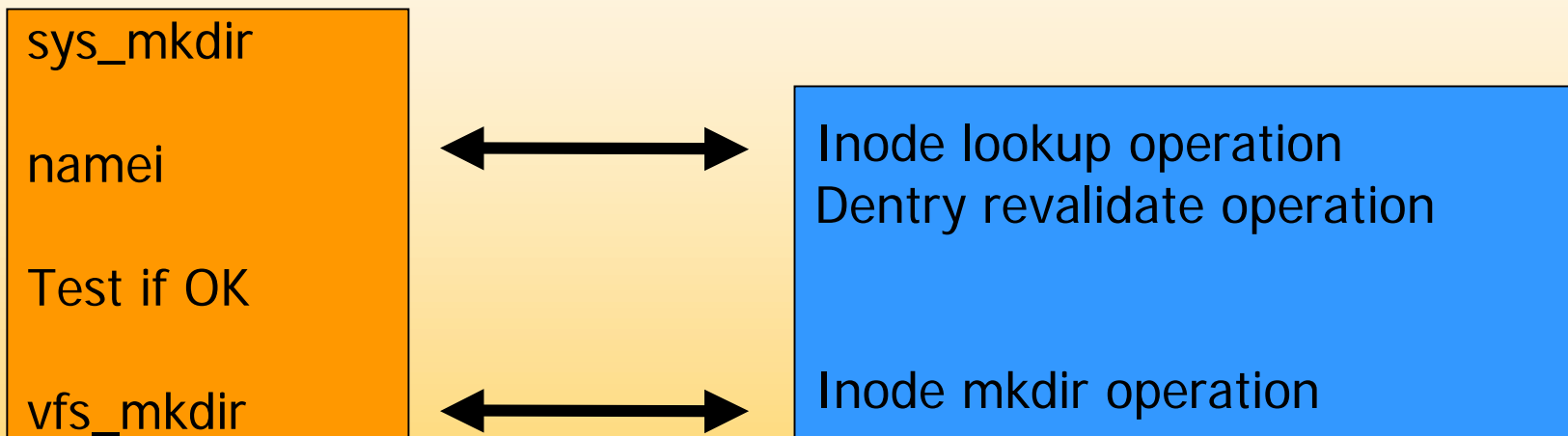- File creations can take one RPC
- Not very good for the single-user case

**Cluster File Systems, Inc**

# What went wrong

Metadata Intent Locking

**Cluster File Systems, Inc**

# Current Linux VFS

VFS

FS

| | |
|---|---|
| sys_mkdir | |
| namei | Inode lookup operation |
| | Dentry revalidate operation |
| Test if OK | |
| vfs_mkdir | Inode mkdir operation |

**Cluster File Systems, Inc**

# We added "intents" to lookups

**VFS**

**FS**

sys_mkdir
namei
  intent mkdir

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

Test if OK
  no:
d_intent_release

⟷ Release lock

vfs_mkdir

⟷ Inode mkdir operation

d_intent_release

⟷ Release lock

**Cluster File Systems, Inc**

# What's the point?

- VFS code prefers to lock directories
    - The intent code reorganizes around that
- Not all metadata loads are alike
    - Locking directories is terrible for concurrent updates
    - Server execution is terrible for the common single-client case

- Intents give the option to the file system

**Cluster File Systems, Inc**

# Intents gone wrong

- Juggling too many things
  - Needed too many locks to safely use the VFS code
- Too many corner cases
  - Server view of execution must *exactly* match the client's view


- There was a simpler solution right around the corner…

**Cluster File Systems, Inc**

# Intents become "raw" operations

- Still gives the lock manager an opportunity to choose

- If contention is low…
  - Server returns a *write-back lock*
- If contention is high…
  - Server executes for us, sends a return code
  - Returns no locks at all
  - We skip **all** client VFS code

**Cluster File Systems, Inc**

# What went well

Object Protocols

**Cluster File Systems, Inc**

# Lustre file I/O in brief

- Object protocols were a no-brainer
- No shared-block file system will scale to 1,000 nodes
  - Shared disk much too expensive
  - Locking for block allocation
- Lustre storage targets manage object and block allocation

**Cluster File Systems, Inc**

# Lustre file I/O in brief

- Very simple object-based protocol.  For example:

    - lock(object id, start, end) → returns lock handle
    - write(object id, offset, data, length)
    - unlock(lock handle)

**Cluster File Systems, Inc**

# Yes, but does it *work*?

**Cluster File Systems, Inc**

# What went so-so

Debugging

**Cluster File Systems, Inc**

# Debugging

- An extensive logging system
  - The log output is frequently more than the size of the I/O
  - Full debug is gigabytes for a simple test
- Tools to filter and contextualize the logs
- Using the logs requires immense understanding

*"also, if you think the stuff in the Matrix about people learning to "read" is lame, then you haven't watched Phil read Lustre debug logs." – Jacob*

**Cluster File Systems, Inc**

# Tools

- Linus hates them, but we thrive on good debug tools
- The first time we write a piece of code, we test it in UML under GDB
- We make extensive use of *mcore*, *netdump*, and *crash* on the real hardware
  - Saves more time than I care to count
- Working on kgdb-over-UDP extension
  - Turns out someone already started

**Cluster File Systems, Inc**

# Debug issues

- Sometimes we let our tools slip
  - Improving our tools almost immediately improves our work
  - The tools on ia64 were *terrible* for a long time
- It's not a trivial system
- It's still a 1,000-node state machine

- Overall, we get a B for debugging

**Cluster File Systems, Inc**

# The real world

- 3 of the top 8 supercomputers in the world run Linux.  Lustre runs on all 3.
  - LLNL MCR: 1,100-node ia32 cluster (#3)
  - LLNL ALC: 950-node ia32 cluster (#6)
  - PNNL EMSL: 950-node ia64 cluster (#8)
- Installing in 2003-2004:
  - NCSA: 1,000 nodes
  - SNL/ASCI Red Storm: 8,000 nodes
  - LANL Pink: 1,000 nodes
- Chosen for ASCI PathForward SGS file system

**Cluster File Systems, Inc**