

TODO

- Caching?

Lustre

Paul Bienkowski

2bienkow@informatik.uni-hamburg.de

Proseminar “Ein-/Ausgabe - Stand der Wissenschaft”

2013-06-10

Outline

- 1 Introduction
- 2 The Project
 - Goals and Priorities
 - History
 - Who is involved?
- 3 Lustre Architecture
 - Network Architecture
 - Data Storage and Access
 - Software Architecture
- 4 Performance
 - Throughput Examples
 - Scalability
 - ???
- 5 Conclusion
- 6 References

What is Lustre

parallel, scaling, for clusters, based within linux kernel...

The Project

1 Introduction

2 The Project

- Goals and Priorities
- History
- Who is involved?

3 Lustre Architecture

- Network Architecture
- Data Storage and Access
- Software Architecture

4 Performance

- Throughput Examples
- Scalability
- ???

5 Conclusion

6 References

Goals

2007: performance > features > stability

“it’s a science project”

2010: stability > performance > features

used in high-performance production environments

History

- started as a research project in 1999 by Peter Braam
- Braam founds **Cluster File Systems**
- 1.0 released in 2003
- **Sun Microsystems** acquires Cluster File Systems in 2007
- **Oracle Corporation** acquires Sun Microsystems in 2010
- Oracle ceases Lustre development, many new Organizations continue development, including **Xyratec**, **Whamcloud**, and more
- in 2012, **Intel** acquires Whamcloud
- in 2013, Xyratec purchases the original Lustre trademark from Oracle

Who is involved?

Oracle *no development*, only pre-1.8 support

Intel funding, preparing for *exascale computing*

Cray funding, development (Titan Supercomputer)

Xyratex hardware bundling

OpenSFS (Open Scalable File Systems) “keeping Lustre open”

EOFS (EUROPEAN Open File Systems) (community collaboration)

FOSS Community many joined one of the above to help development
(e.g. Braam works for Xyratex now)

DDN, Dell, NetApp, Terascale, Xyratex

storage hardware bundled with Lustre

Who is involved?

Supercomputers

Titan & Co. use it!

Lustre Architecture

1 Introduction

2 The Project

- Goals and Priorities
- History
- Who is involved?

3 Lustre Architecture

- Network Architecture
- Data Storage and Access
- Software Architecture

4 Performance

- Throughput Examples
- Scalability
- ???

5 Conclusion

6 References

Network Structure

What data is stored where? (graph)

Metadata Server

Where are they? How to access?

Object Storage Server

Contain OSTs (Object storage targets)

Network Capabilities

How is data transferred?

Protocol stack (TCP, ...)

Different network types (ethernet, infiniband, ...)

Failover

Failover mechanism and typical setups (graphs)

Why are is a failover mechanism cool? Live-Upgrades!

Excursion: INodes

... because MDS do something similar (metadata records)
(graph)
compare this on next slide

Metadata

how metadata is stored in the MDS

what metadata is stored?

how metadata is fetched from the MDS

Striping

[repeat] what is striping (RAID 0)

why do they use it in lustre – speed advantage

Software Architecture

what software is running where?

Interversion Compatibility

Sun “guarantees” [citation needed] compatibbility between minor versions
→ on-line upgrade-ability using failover systems

Idiskfs - Customized ext3

why we need a customized filesystem to work ON TOP of

Kernel patching (serverside)

just tell them the kernel needs to be patched (2.6.*) and what that means

Patchfree Client

How can clients access the data? (lustre-fs, liblustre, NFS)
(kernel-independent)
even NFS, that works everywhere!

Limitations

Server

- very platform dependent
- needs compatible kernel

Client

- all linux kernel >2.6 supported
- NFS for Windows, MacOS
- even FUSE support on the way

Performance

1 Introduction

2 The Project

- Goals and Priorities
- History
- Who is involved?

3 Lustre Architecture

- Network Architecture
- Data Storage and Access
- Software Architecture

4 Performance

- Throughput Examples
- Scalability
- ???

5 Conclusion

6 References

Throughput Examples

Yes, the speeds add up!

There are systems with 5000 OSS.

Up to 160 OSS / file.

16 OST/OSS is quite normal.

1 TiB/OST

Scalability

Just multiply. Works as long as your network supports it (e.g. InfiniBand > WiFi ...)

???

???

Maybe I can find some more data samples – where is the interesting stuff?

Conclusion

...

References

...