# DDNの最新のLustreへの取り組みについて

**DataDirect Networks, Inc.**

Shuichi Ihara

2017/11/02

# DDN's recent Lustre activities

▶ **DDN has been contributing Lustre community**
- DDN is No.2 company of Lustre cods contribution.
- DDN developed many new Lustre features and merged them into upstream Lustre.

▶ **Presenting at Lustre conference every year. Four presentations were selected at LAD17!**
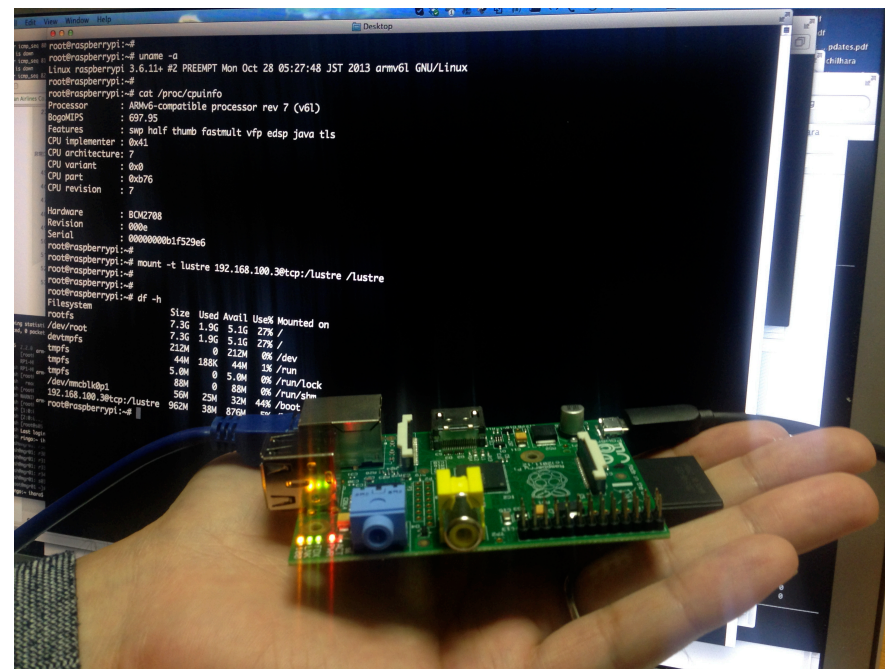- Lustre on ARM servers
- PCC(Lustre Persistent Client Cache) LU-10092
- Lustre Audit with Changelogs LU-9727
- Lustre/ldiskfs metadata performance boost LU-9796

▶ **Other activities**
- Lustre QoS (Corroboration work with University of Mainz)
  ○ Selected paper at SC17 (Tuesday, November 14th11:30am - 12pm)
- Data Archive Solution
- Lustre Integrated Policy Engine
- Lustre-ZFS
- ... others

**DDN** STORAGE

ddn.com

# Lustre on ARM Servers

ddn.com

# The Cavium ThunderX Architecture

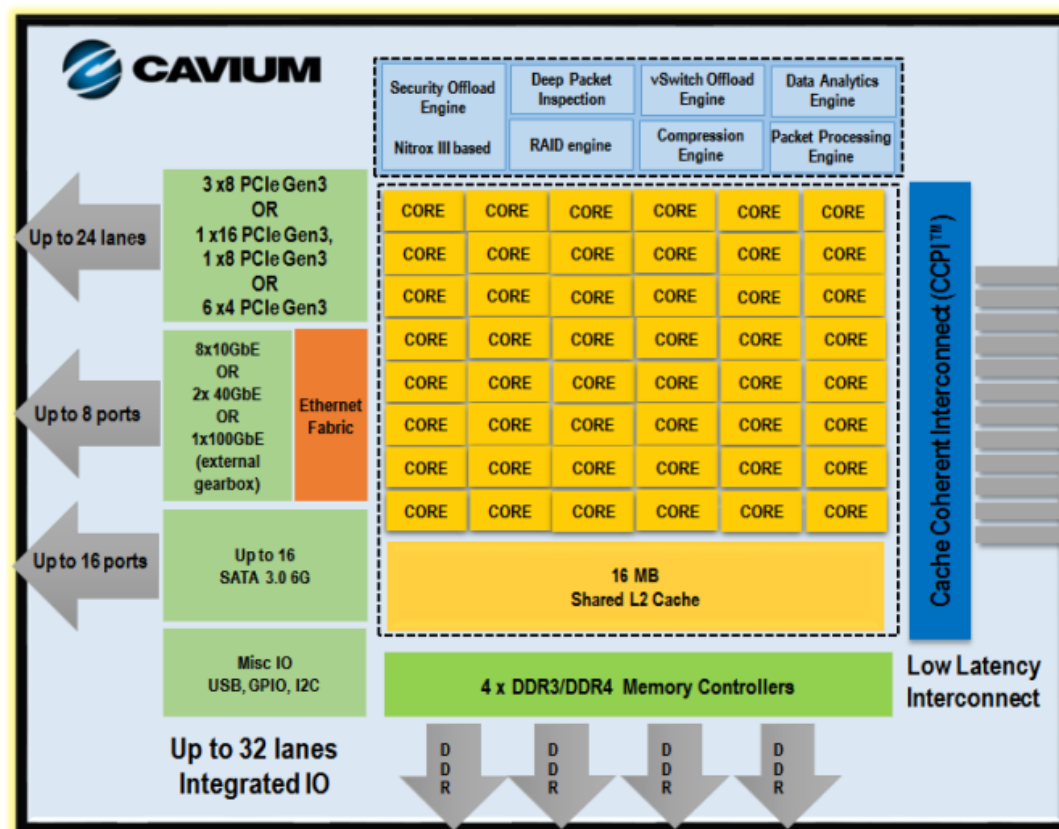▶ **SoC architecture**

- ISA: ARMV8

```
root@s167:/proc# lscpu
Architecture:          aarch64
Byte Order:            Little
                        Endian
CPU(s):                96
On-line CPU(s) list:   0-95
Thread(s) per core:    1
Core(s) per socket:    48
Socket(s):             2
NUMA node(s):          2
L1d cache:             32K
L1i cache:             78KL2
cache:                 16384K
NUMA node0 CPU(s):     0-47
NUMA node1 CPU(s):     48-95
```
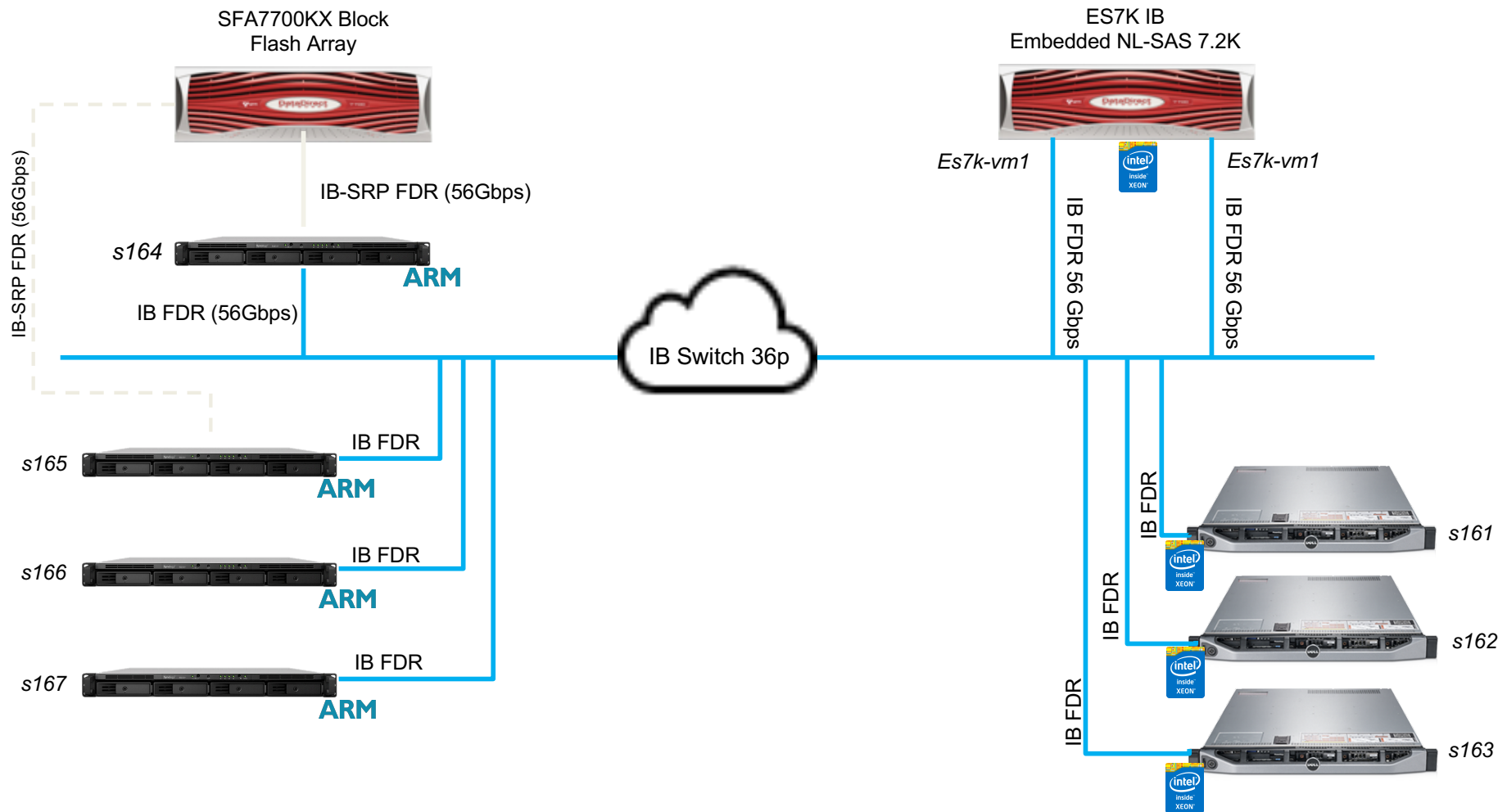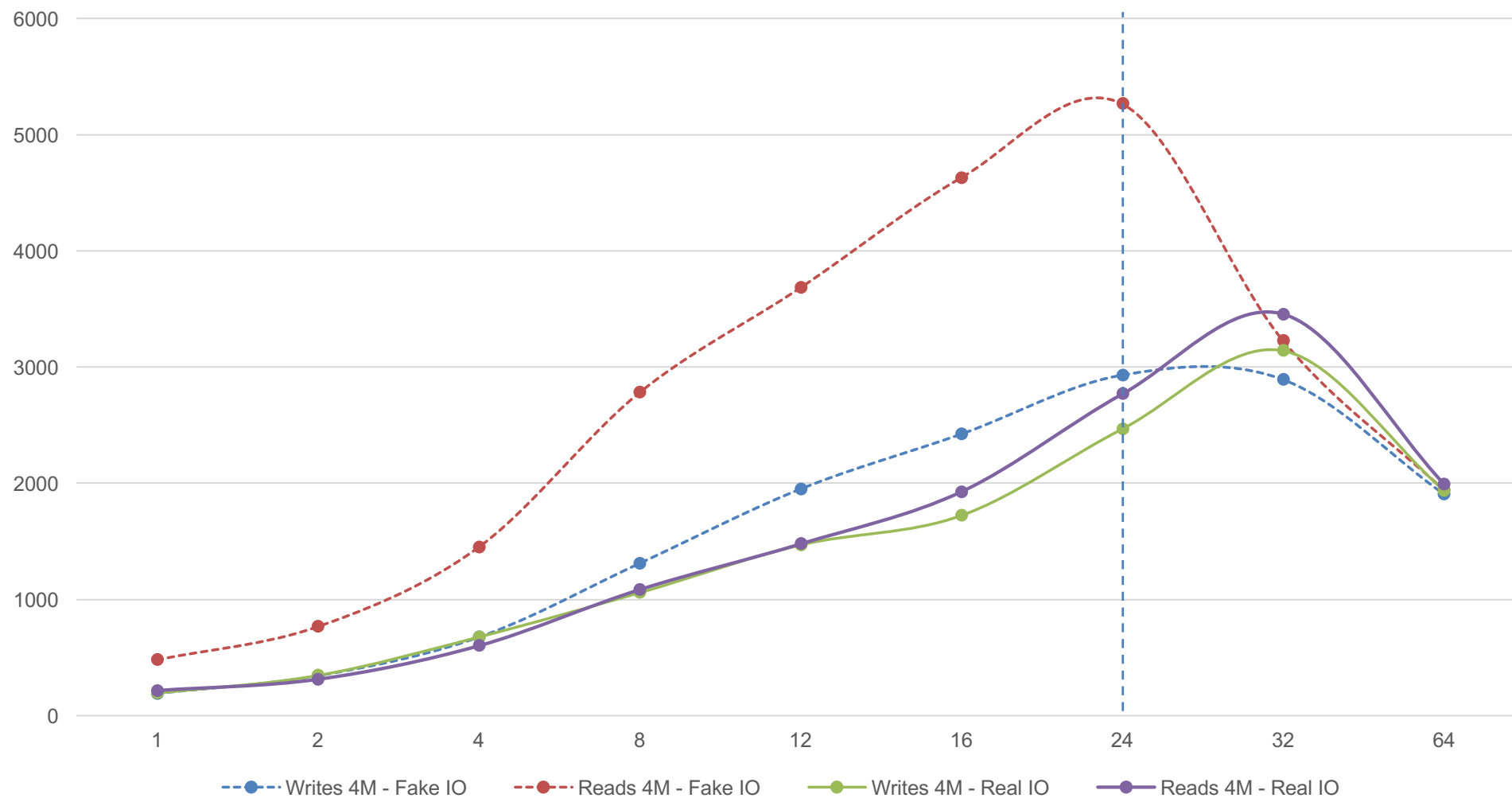
# DDN Goals evaluating ARM

▶ **Understand if it is a viable option for mid/long term future products**

▶ **Understand what's the effort necessary to make Lustre running optimally on ARM (client and server-side)**

▶ **Understand how Lustre and general I/O behaves on ARM SoC architecture**

▶ **Contribute to the community**

**DDN STORAGE**

ddn.com

# Test Environment used for the study



SFA7700KX Block
Flash Array

ES7K IB
Embedded NL-SAS 7.2K

IB-SRP FDR (56Gbps)

IB-SRP FDR (56Gbps)

*Es7k-vm1*

*Es7k-vm1*

s164

**ARM**

IB FDR (56Gbps)

IB FDR 56 Gbps

IB FDR 56 Gbps

IB Switch 36p

s165

**ARM**

IB FDR

IB FDR

s161

s166

**ARM**

IB FDR

IB FDR

s162

s167

**ARM**

IB FDR
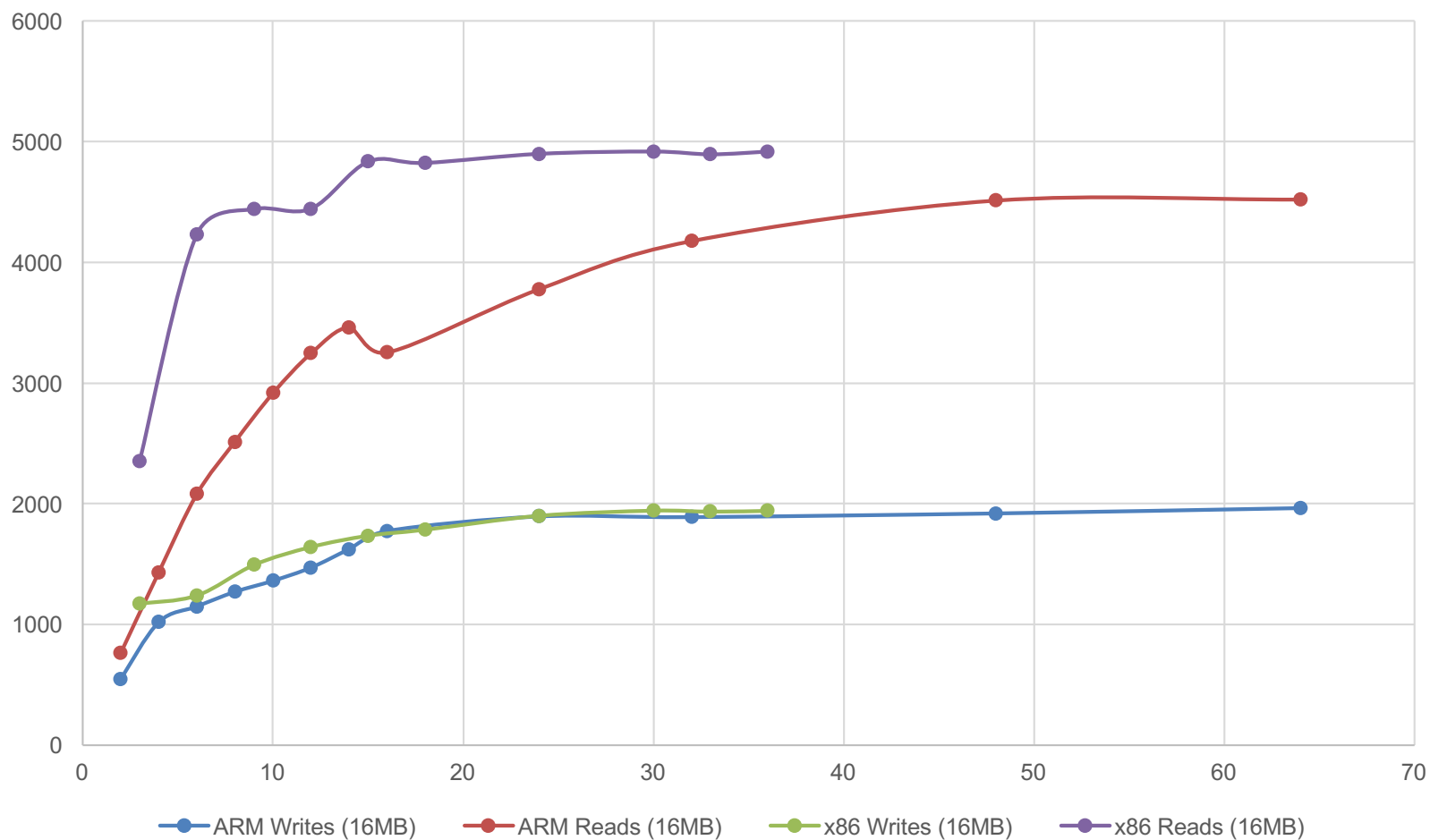
IB FDR

s163

ddn.com

DDN STORAGE

# IOR Single Client Performance – Multiple Threads



IOR Single Client Performance - 4MB RPCs - REGULAR vs FAKE IO
/mnt/arm/bin/ior.arm.mvapich -a POSIX -b 1g -r -w -F -B -t 4m -o /mnt/arm/file.out

Legend: Writes 4M - Fake IO · Reads 4M - Fake IO · Writes 4M - Real IO · Reads 4M - Real IO

DDN STORAGE

ddn.com

# ARM and x86 Clients comparison
# IOR, multiple clients - Sequential

ARM and x86 Clients - IOR Sequential Reads / Writes (ARM Server)



Legend: ARM Writes (16MB) ● ARM Reads (16MB) ● x86 Writes (16MB) ● x86 Reads (16MB) ●

DDN STORAGE

ddn.com

# PCC(Lustre Persistent Client Cache)

**DDN STORAGE**

ddn.com

# NSCC-Wuxi and the Sunway Machine Family



Sunway-I:

- CMA service, 1998

- commercial chip

- 0.384 Tflops

- 48$^{th}$ of TOP500



Sunway BlueLight:
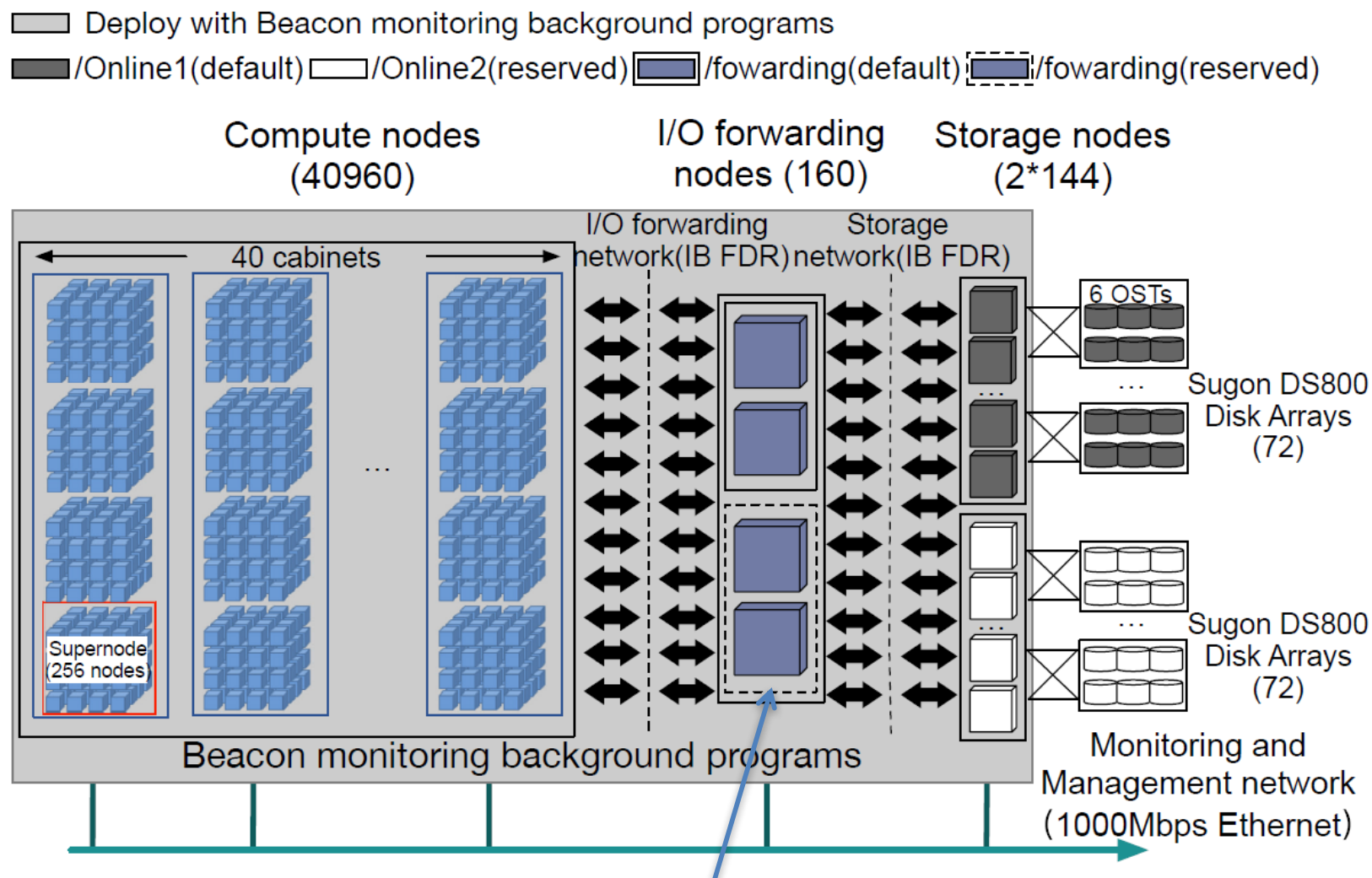
- NSCC-Jinan, 2011

- 16-core processor

- 1 Pflops

- 14$^{th}$ of TOP500



Sunway TaihuLight:

- NSCC-Wuxi, 2016

- 260-core processor

- 125 Pflops

- 1$^{st}$ of TOP500

## PCC project is collaborated by NSCC-Wuxi and DDN
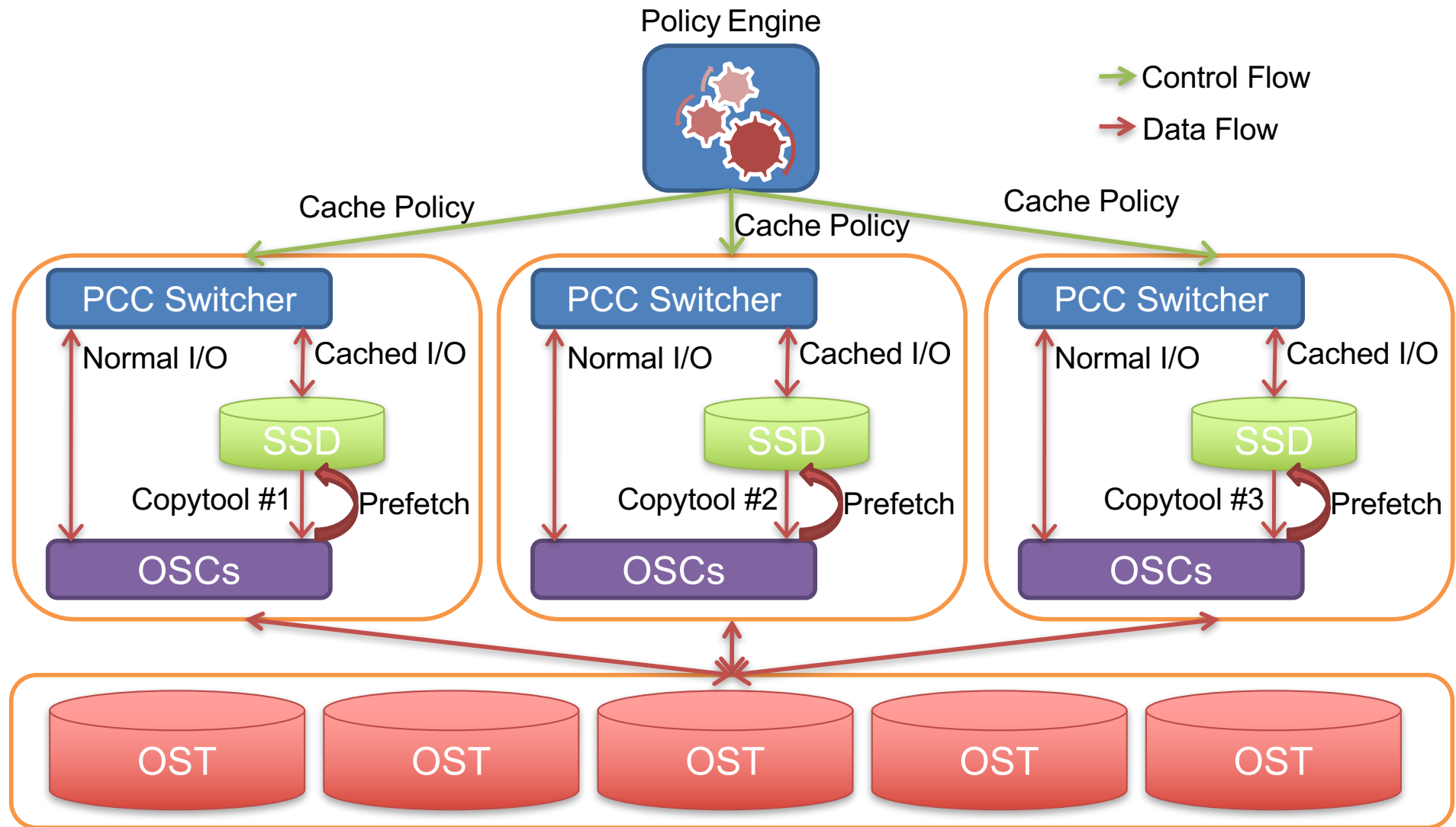
# I/O Architecture of Sunway TaihuLight



**Cache on I/O forwarding nodes (Lustre clients) should be helpful**

DDN STORAGE

ddn.com

# Why SSD cache on Lustre client?

▶ **Less overhead visible for applications**
- No network latency
- No LDLM lock and other Lustre overhead

▶ **Easier to be optimized for the best performance**
- I/O stack is much simpler
- No interference I/Os from other clients

▶ **Relatively easier than server side implementations**
- Write support for SSD cache on server side is very difficult
- Problems for write cache on server side:
  o Visibility when failover happens
  o Consistency when corruption happens

▶ **Less requirement on hardware**
  o Any kind of SSD can be used as the cache device

▶ **Reduces the pressure of OSTs**
  o Small or random I/Os are regularized to big sequential I/Os
  o Temporary files do not need to be flushed to OSTs

# Architecture of PCC

# Limitations

▶ **Not all applications are able to be accelerated by PCC**

- Locality requirements of application I/Os
  ○ Applications shall not access the cached file through multiple clients
  ○ But no inconsistency will happen even the application writes the cached file on a remote client
- Capacity of each local cache is limited
  ○ Size of a cached file is limited to the available space of the local cache
  ○ The total cached data on a single client is limited

▶ **Files can not be partly cached**

- Partial cache can be implemented if HSM supports partial archive/restore

▶ **The total PCC clients are limited to 32 Today**

- Only 32 different archive numbers are supported by Lustre
- This upper limitation can be raised in the future

ddn.com

# Lustre Audit with Changelogs

ddn.com

DDN
STORAGE

# Need for audit in Lustre

▶ **Support of rich security features:**

- authentication with Kerberos
- mandatory access control with SELinux
- isolation
- etc.

⇒ **Audit as a proof of security in place**

▶ **Lustre outside of traditional HPC field**

▶ **e.g. Life science**

- data privacy is crucial

⇒ **Audit as a regulation compliance**

# Audit with SELinux

| Pros | Cons |
|---|---|
| • integrated logging and auditing facility<br><br>• proven | • on client side<br>• need to consolidate |

ddn.com

DDN
STORAGE

# Audit with Changelogs

| Pros | Cons |
|------|------|
| • integrated in Lustre<br>• centralized<br>• transactional | • lacks some info |

ddn.com

DDN
STORAGE

# Audit with Changelogs

▶ **Lustre activity as seen by MDS**

- file system namespace
- file metadata

▶ **Store in Changelog records**

- internal Lustre files

▶ **Read from audit nodes**

- dedicated clients

```
5 01CREAT 15:44:32.385864793 2017.07.18 0x0 t=[0x200000402:0x3:0x0]
        ef=0x1 p=[0x200000402:0x2:0x0] fileA
```

DDN
STORAGE

# Lustre needs for proper audit

▶ **Identify subject of action**
  - uid/gid
  - NID

▶ **Record all actions**
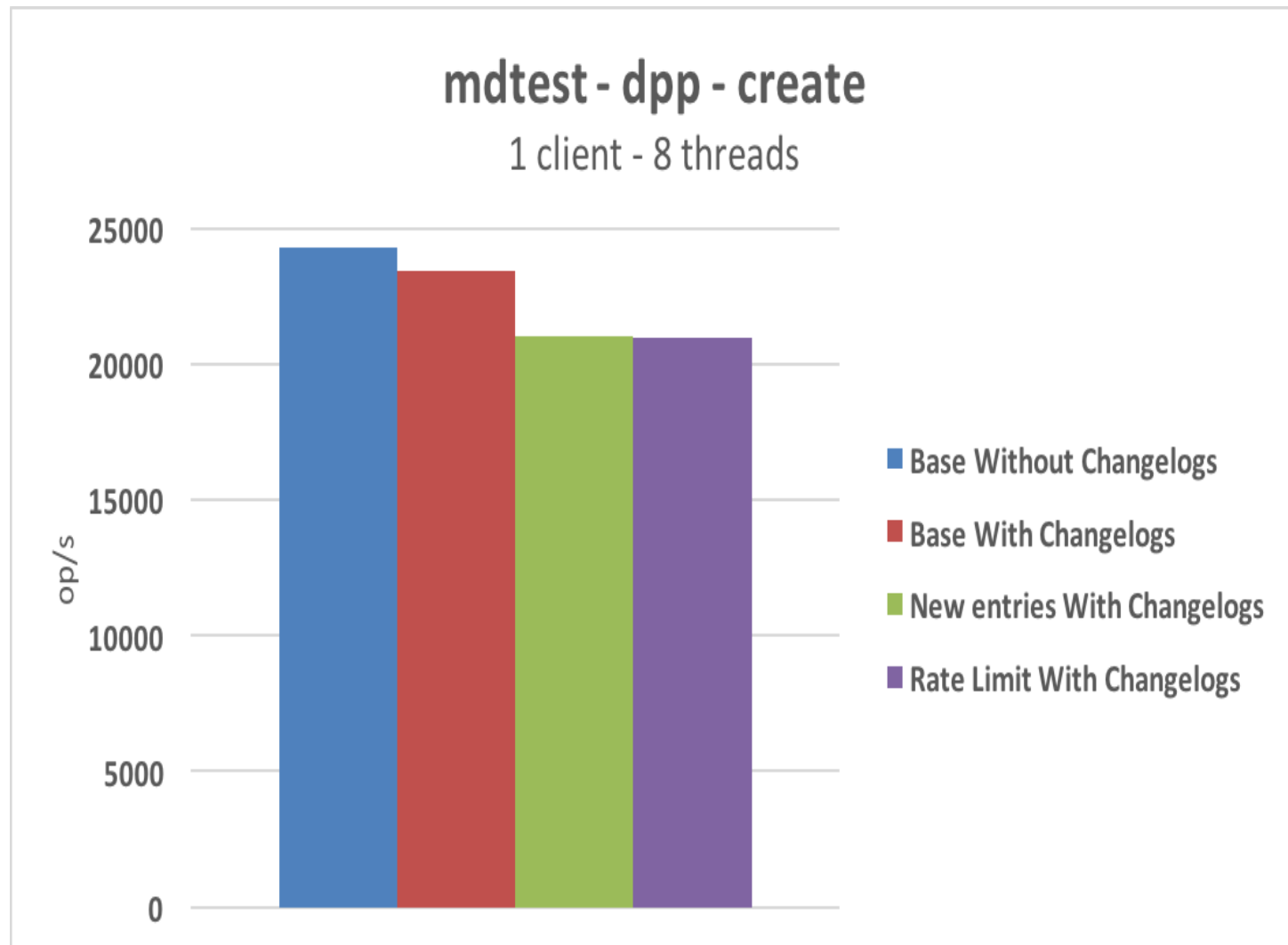  - open
  - close
  - xattr
  - denied accesses

ddn.com

# Audit with Changelogs: impact study

▶ **Changelogs space consumption evaluation**

|  | # changelog entries | changelog size |
|---|---|---|
| After 10 000 files created | 30000 | 3755824 |
| After 10 000 files read | 50000 | 6096448 |
| After 10 000 files removed | 60000 | 7461440 |

▶ **MDT**

**DDN STORAGE**

# Audit with Changelogs: impact study



mdtest - dpp - create
1 client - 8 threads

- Base Without Changelogs
- Base With Changelogs
- New entries With Changelogs
- Rate Limit With Changelogs

ddn.com

DDN STORAGE

# Lustre Metadata Performance improvement

# Why is metadata performance important?

▶ **Lustre is general purpose filesystem for Big data**

- 1 Million files per job are quite common with life science application

- AI/Machine learning type of workload requires small file access with low latency. Metadata performance is one of key factors of it.

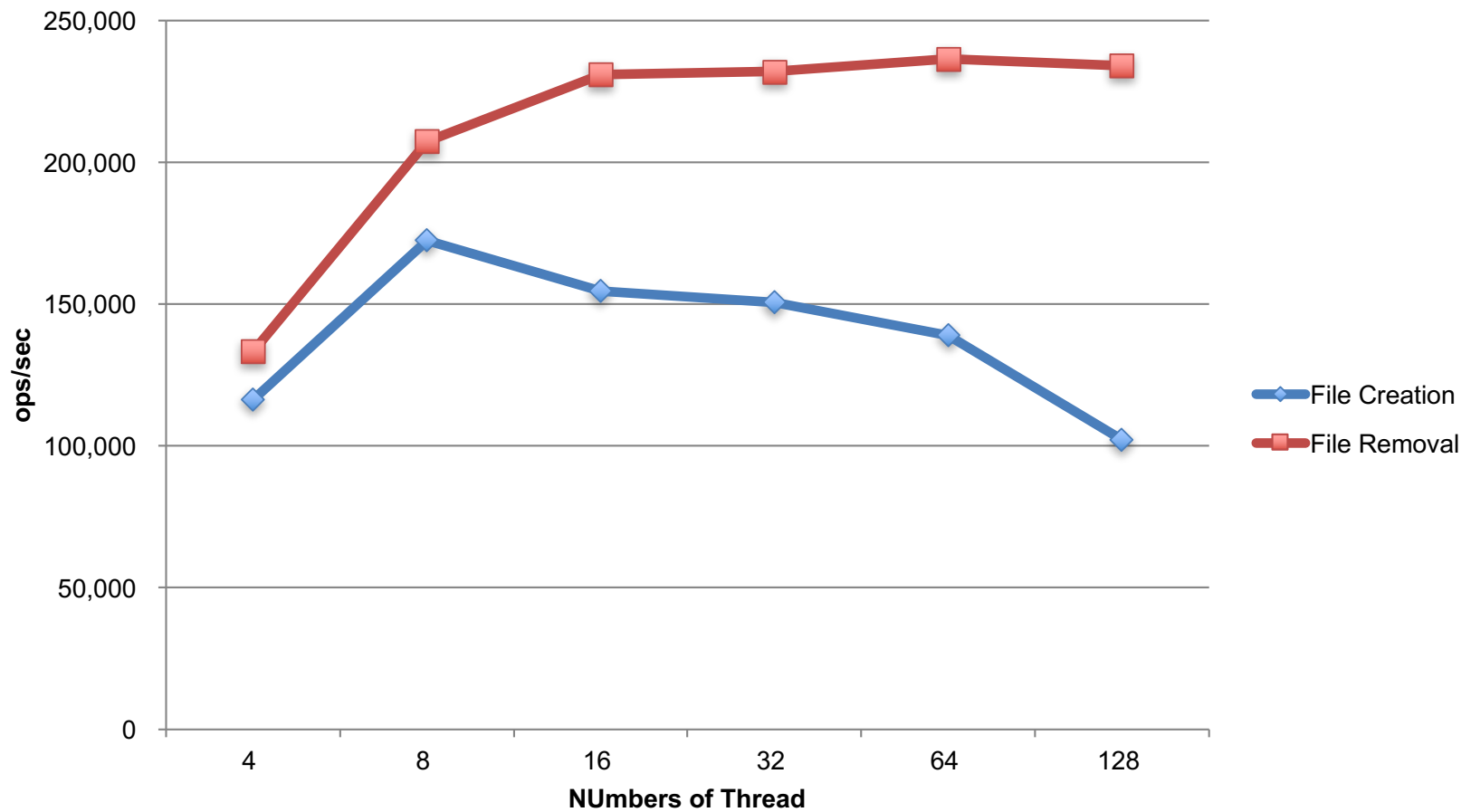- Lustre metadata performance has been performing well.

▶ **Vertical and Horizontal scale**

- 28 (and 32) CPU cores/socket is available Today.

- DNE helps Horizontal scale out Metadata, but needs to understand your single MDS metadata performance first.

DDN STORAGE

ddn.com

# MDS-Survey
# RHEL7.3/Lustre-2.10/ldiskfs

**MDS-Survey(File Creation and Unlink)**
**RHEL7.3/Lustre-2.10.1RC/ldiskfs (Quota Enabled)**

# A problem on File creation under concurrency

▶ **Profiled with perf-tools during mdtest to ldiskfs/ext4**

- Collected CPU costs for all functions in ext4 and jbd2
- Found heavy lock contentions on group spinlock

| FUNC | TOTAL_TIME(us) | COUNT | AVG(us) |
|---|---|---|---|
| ext4_create | 1707443399 | 1440000 | 1185.72 |
| **_raw_spin_lock** | **1317641501** | **180899929** | **7.28** |
| jbd2__journal_start | 287821030 | 1453950 | 197.96 |
| jbd2_journal_get_write_access | 33441470 | 73077185 | 0.46 |
| ext4_add_nondir | 29435963 | 1440000 | 20.44 |
| ext4_add_entry | 26015166 | 1440049 | 18.07 |
| ext4_dx_add_entry | 25729337 | 1432814 | 17.96 |
| ext4_mark_inode_dirty | 12302433 | 5774407 | 2.13 |

- Same contentions exist in the upstream kernel

DDN STORAGE

ddn.com

# Fix lock contentions in upstream kernel

▶ **Fixed and merged upstream kernel (4.14)**
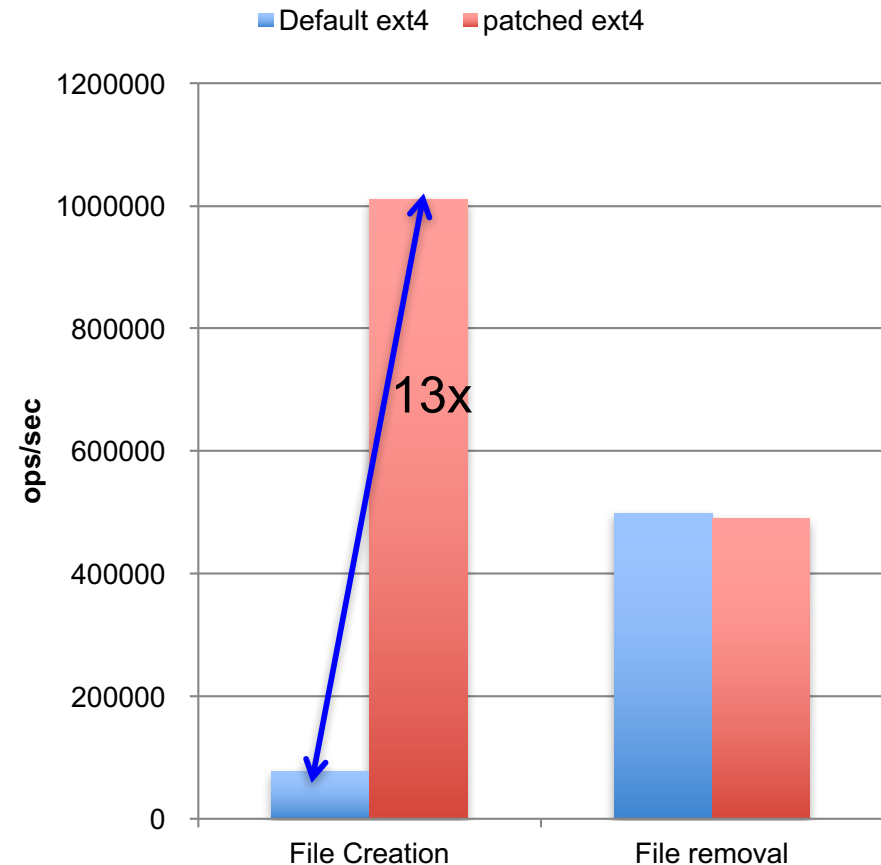
Wang Shilong (2):
ext4: cleanup goto next group
ext4: reduce lock contention in __ext4_new_inode

▶ **13x performance improvement on file creation**

• Run mdtest to ext4 directly
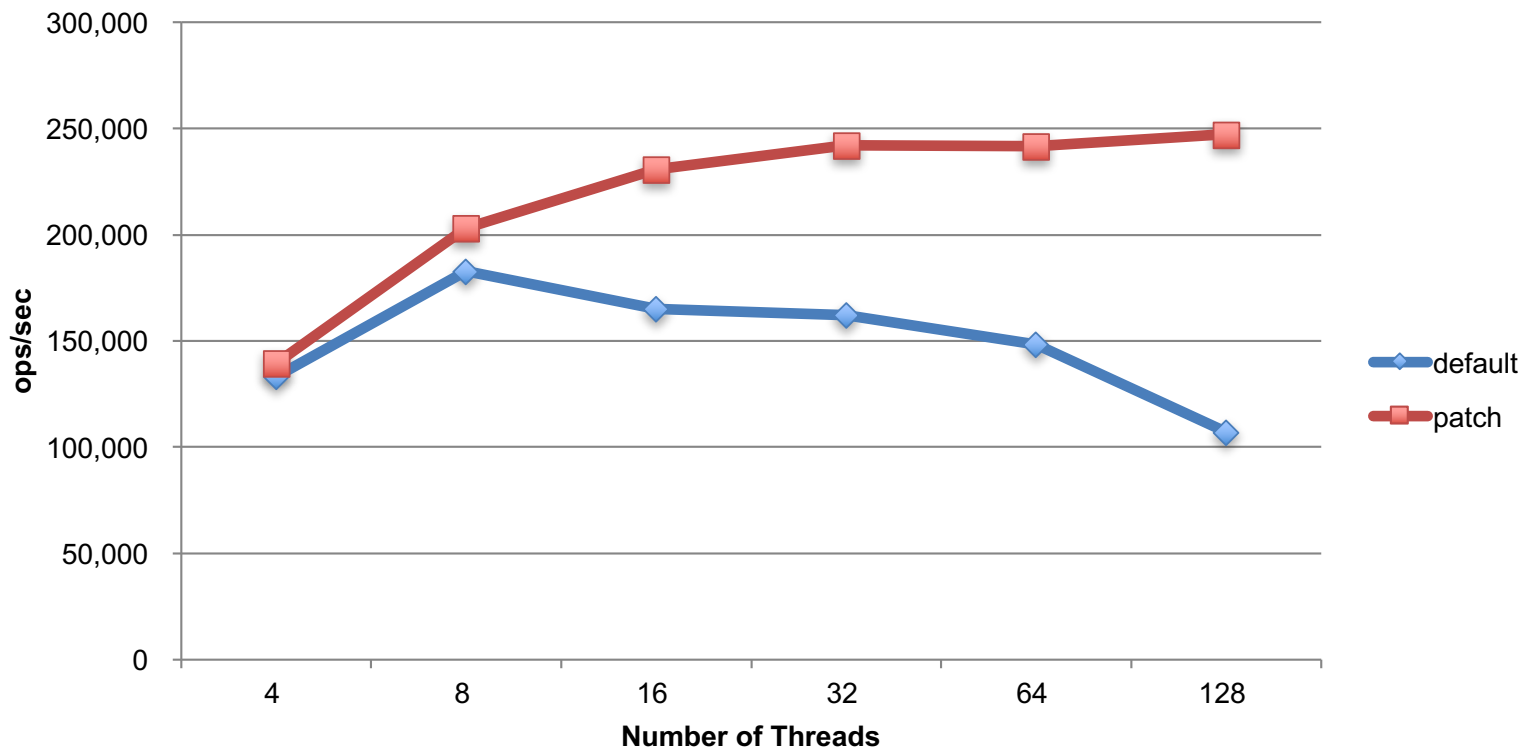• Unique directory operations
• Quota disabled

### mdtest to ext4 (linux-4.13-rc5)



■ Default ext4   ■ patched ext4

13x

ddn.com

# mds-survey on patched ldiskfs

▶ **LU-9796: speedup file creation under heavy concurrency**

▶ **Ported patches to ldiskfs for RHEL7 kernel**

**File Creation :mds-survey on ldiskfs
1 x MDS and 1 x MDT(2 x RAID1 SSD)**

ddn.com

# Conclusion

▶ **DDN keeps investment to Lustre and contributions to Lustre community**
- DDN Lustre R&D in Japan and China
- Our most of developed new features comes from valuable customer feedbacks!

▶ **Deliver adaption and optimizations for new hardware and new technology in advance**
- Performance lab is located in Tokyo
- Various early testing, performance optimization are ongoing

▶ **Welcome Co-research and Collaboration**
- Not only co-research, but also Alpha/Beta testing and feedback are much appreciate!

ddn.com