JLUG2018



Topics from ISC18, LAD & Lustre Developer Summit 2018

Shinji Sumimoto, Ph.D. Next Generation Technical Computing Unit FUJITSU LIMITED Oct. 25th, 2018





Lustre History

- SSD Based First Layer File System for the Next Generation Super-computer
 - A64FX: High Performance Arm CPU
 - SSD Based First Layer File System for the Next Generation Super-computer
 - Current Status of Lustre Based File System Development

Contribution of DL-SNAP to Lustre Community

Lustre History

- 1999: Lustre file system 研究プロジェクト開始 by Peter J. Braam
- 2001: Cluster File Systems社設立
- 2007/9: Sun社がCluster File Systems社を買収
- 2009/4: Oracle社がSun社を買収、コミュニティに不安が走る
- 2010/4: Oracle社がサポートを自社ハードに制限、コミュニティ激震
 - 欧州<u>EOFS</u>、米国国研ベース<u>OpenSFS</u>、 World Wideユーザベース<u>HPCFS</u>の3つのコミュニティに分散
- 2010/9: Whamcloud社設立
 - Oracle社からWhamcloud社に技術者が集結
- 2010/12: Oracle社がLustre開発凍結
- 2011/4: LUG2011で3つのコミュニティがOpen SFS+EOFSとして再始動
- 2012/7: Intel社がWhamcloud社を買収
- 2018/6: DDN社がIntelのLustre事業を買収、新生Whamcloud誕生
- 2019: Luster 20 years anniversary

LUG2009: Lustre 10th Anniversary





Copyright 2018 FUJITSU LIMITED

LUG 2011: Single Community





2018/6: DDN社がIntelのLustre事業を買収





ISC18: Lustre BoF







Whamcloud ベンダー中立の立場でこれからもコミュニティへの貢献 Lustre 真のクラスタファイルシステムとして正当に発展



■引き続きLustreを支持しコミュニティ発展に協力

FUITSU



A64FX: High Performance Arm CPU

- From presentation slides of Hotchips 30th and Cluster 2018
- Inheriting Fujitsu HPC CPU technologies with commodity standard ISA



A64FX Chip Overview

FUjitsu



7nm FinFET

- 8,786M transistors
- 594 package signal pins

Peak Performance (Efficiency)

- >2.7TFLOPS (>90%@dgemm)
- Memory B/W 1024GB/s (>80%@Stream Triad)

	A64FX (Post-K)	SPARC64 XIfx (PRIMEHPC FX100)
ISA (Base)	Armv8.2-A	SPARC-V9
ISA (Extension)	SVE	HPC-ACE2
Process Node	7nm	20nm
Peak Performance	>2.7TFLOPS	1.1TFLOPS
SIMD	512-bit	256-bit
# of Cores	48+4	32+2
Memory	HBM2	HMC
Memory Peak B/W	1024GB/s	240GB/s x2 (in/out)

A64FX Memory System

Extremely high bandwidth

- Out-of-order Processing in cores, caches and memory controllers
- Maximizing the capability of each layer's bandwidth



FUITSU

A64FX Core Features



- Optimizing SVE architecture for wide range of applications with Arm including AI area by FP16 INT16/INT8 Dot Product
- Developing A64FX core micro-architecture to increase application performance

	A64FX (Post-K)	SPARC64 XIfx (PRIMEHPC FX100)	SPAR64 VIIIfx (K computer)	
ISA	Armv8.2-A + SVE	SPARC-V9 + HPC-ACE2	SPARC-V9 + HPC-ACE	
SIMD Width	512-bit	256-bit	128-bit	
Four-operand FMA	our-operand FMA 🗸 Enhanced		\checkmark	
Gather/Scatter	ther/Scatter 🗸 Enhanced			
Predicated Operations ✓ Enhanced		\checkmark	\checkmark	
Math. Acceleration 🗸 Further enhanced		✓ Enhanced	\checkmark	
Compress	✓ Enhanced	\checkmark		
First Fault Load	✓ New			
FP16	✓ New			
INT16/ INT8 Dot Product	✓ New			
HW Barrier* / Sector Cache* ✓ Further enhanced		✓ Enhanced	\checkmark	

* Utilizing AArch64 implementation-defined system registers

A64FX Chip Level Application Performance



- Boosting application performance up by micro-architectural enhancements, 512-bit wide SIMD, HBM2 and semi-conductor process technologies
 - > 2.5x faster in HPC/AI benchmarks than that of SPARC64 XIfx tuned by Fujitsu compiler for A64FX micro-architecture and SVE



A64FX Kernel Benchmark Performance (Preliminary results)

A64FX TofuD Overview

Halved Off-chip Channels

Power and Cost Reduction

Increased Communication Resources

- TNIs from 2 to 4
- Tofu Barrier Resources

Reduced Communication Latency

Simplified Multi-Lane PCS

Increased Communication Reliability

Dynamic Packet Slicing: Split and Duplicate

	Tofu K.comp	Tofu2 FX100	TofuD
Data rate (Gbps)	6.25	25.78	28.05
# of signal lanes per link	8	4	2
Link bandwidth (GB/s)	5.0	12.5	6.8
# of TNIs per node	4	4	6
Injection bandwidth per node (GB/s)	20	50	40.8





TofuD: 6D Mesh/Torus Network

Six coordinates: (X, Y, Z) × (A, B, C)

- X, Y and Z: sizes are depends on the system size
- A, B and C: sizes are fixed to 2, 3, and 2 respectively
- Tofu stands for "torus fusion"



FUJITSU

TofuD: Packaging – CPU Memory Unit

Two CPUs connected with C-axis

- $\blacksquare X \times Y \times Z \times A \times B \times C = 1 \times 1 \times 1 \times 1 \times 1 \times 2$
- Two or three active optical cable cages on board

Each cable is shared by two CPUs



FUITSU

TofuD: Packaging – Rack Structure

FUJITSU

Rack

- 8 shelves
- 192 CMUs or 384 CPUs

Shelf

- 24 CMUs or 48 CPUs
- $\blacksquare X \times Y \times Z \times A \times B \times C = 1 \times 1 \times 4 \times 2 \times 3 \times 2$

Top or bottom half of rack

4 shelves

 $\blacksquare X \times Y \times Z \times A \times B \times C = 2 \times 2 \times 4 \times 2 \times 3 \times 2$





TofuD: Put Latencies & Throughput& Injection Rate Fujirsu

TofuD: Evaluated by hardware emulators using the production RTL codes
Simulation model: System-level included multiple nodes

	Communication settings	Latency
Tofu	Descriptor on main memory	1.15 µs
	Direct Descriptor	0.91 µs
Tofu2	Cache injection OFF	0.87 µs
	Cache injection ON	0.71 µs
TofuD	To/From far CMGs	0.54 µs
	To/From near CMGs	0.49 µs

	Put throughput	Injection rate		
Tofu	4.76 GB/s (95%)	15.0 GB/s (77%)		
Tofu2	11.46 GB/s (92%)	45.8 GB/s (92%)		
TofuD	6.35 GB/s (93%)	38.1 GB/s (93%)		



Next Generation File System Design

Next Generation File System Structure and Design
Next-Gen 1st Layer File System Overview

Pros:

Stable Application Performance for Jobs

Cons:

- Requiring three times amount of storage which a job needs
- Pre-defining file name of stage-in/out processing lacks of usability
- Data-intensive application affects system usage to down because of waiting prestaging-in/out processing



Next-Gen File System Requirement and Issues Fujirsu

Requirements

- ■10 times higher access performance
- ■100 times larger file system capacity
- Lower power and footprint

Issues

How to realize 10 times faster and 100 times larger file access at a time?

Next-Gen. File System Design

FUJITSU

- K computer File System Design
 - How should we realize High Speed and Redundancy together?
 - Introduced Integrated Two Layered File System.
- Next-Gen. File System/Storage Design
 - Another trade off targets: Power, Capacity, Footprint
 - Difficult to realize single Exabyte and 10TB/s class file system in limited power consumption and footprint.
 - Additional Third layer Storage for Capacity is needed:



Next Gen. File System Design

FUjitsu

- Introducing three level hierarchical storage.
 - 1st level storage: Accelerating application file I/O performance (Local File System)
 - 2nd level storage: Sharing data using Lustre based file system (Global File System)
 - 3rd level storage: Archive Storage (Archive System)
- Accessing 1st level storage as file cache of global file system and local storage
 - File cache on computing node is also used as well as 1st level storage



Application's Access Pattern and SSD Cache Effects Fujirsu

Comparison of Effective Pattern for SSD based storage

	Distributed	Single Shared	Single Shared	I/O	
	Files	Files (1)	Files (2)	Master	
File Reading	Processes				
File , Writing	Processes				
File Read:	Rereading Case∶⊚	Rereading Case :	Rereading Case∶⊚	Rereading Case∶⊚	
Effects	Non Rereading ∶ ×	Non Rereading : ×	Non Rereading : ×	Non Rereading : ×	
File Write:	Rewriting Case:	Rewriting Case:		Rewriting Case :	
Effects	Non Rewriting : O	Non Rewriting : O		Non Rewriting : O	

Data Sharing in a Job on SSD

- Write-Read in a process and among processes are effective to use SSD
- For Persistent Files: File cache of global file system should be shared among processes
- For Temporary Files: Two types of temporary file systems are effective to use SSD
 - Temporary Local System (in a process)
 - Temporary Shared File System (among processes)



Next Gen. Layered File System Requirements Fujirsu

Application views:

- Local File System: Application Oriented File Accesses(Higher Meta&Data I/O)
- Global File System: Transparent File Access
- Archive System: In-direct Access or Transparent File Access(HSM)
- Transparent File Access to the Global File System
 - Local File System Capacity is not enough as much as locating whole data of Global File System
 - File Cache on node memory and Local File System enables to accelerate application performance

	Meta Perf.	Data BWs	Capacity	Scalability	Data Sharing in a Job	Data Sharing among Jobs
Local File System	Ø	Ø	×	Ø	Ø	×
Global File System	\bigcirc	\bigcirc	0	0	×	Ô
Archive System	×	×	Ø	×	×	×

Next-Gen 1st Layer File System Overview

FUĴĨTSU

Goal: Maximizing application file I/O performance

Features:

- Easy access to User Data: File Cache of Global File System
- Higher Data Access Performance: Temporary Local FS (in a process)
- Higher Data Sharing Performance: Temporary Shared FS (among processes)
- Now developing LLIO(Lightweight Layered IO-Accelerator) Prototype



LLIO Prototype Implementation

- Two types of Computing Nodes
 - Burst Buffer Computing Node(BBCN)
 - Burst Buffer System Function with SSD Device
 - Computing Node(CN)
 - Burst Buffer Clients: File Access Request to BBCN as burst buffer server



File Access Sequences using LLIO (Cache Mode)



FUJITSU

LLIO Prototype I/O Performance





Evaluated on IA servers using Intel P3608

Higher I/O performance than those of NFS, Lustre
Utilizing maximum physical I/O device performance by LLIO



Contribution of DL-SNAP to Lustre Community

Our Strategy for contributing DL-SNAP:

- We are ready to contribute our current DL-SNAP code for Lustre 2.6
- We have made the LU-ticket for the DL-SNAP (LU-11512, Oct. 2018)

What is DL-SNAP?



- DL-SNAP is designed for user and directory level file backups.
- Users can create a snapshot of a directory using lfs command with snapshot option and create option like a directory copy.
- The user creates multiple snapshot of the directory and manage the snapshots including merge of the snapshots.
- DL-SNAP also supports quota to limit storage usage of users.

DL-SNAP Use-case 1



Avoiding file deletion or corruption by file operation





DL-SNAP Use-case 2



Maintaining large database with partially different data
Updating database by an application using DL-SNAP



FUJTSU

shaping tomorrow with you