# Accelerate Ceph performance via SPDK related techniques

Ziye Yang
Intel Corporation
Oct 2015

# Legal Disclaimer

# Agenda

- Recent common requirements for Ceph

- Middle Cache tiering solution

- Building block techniques

  - I/O optimization technique
    - DPDK for storage

  - Data Processing Acceleration techniques
    - ISA-L

- Conclusion

SG
Storage Group

# Agenda

- **Recent common requirements for Ceph**

- Middle Cache tiering solution

- Building blocks techniques

  - I/O optimization technique
    - DPDK for storage

  - Data Processing Acceleration techniques
    - ISA-L

- Conclusion

SG
Storage Group

(intel)

# Recent Common requirements for Ceph

- Legacy protocol support
  - For example, iSCSI interface support for transparently migrating applications from Enterprise storage to cloud storage – Ceph
- High performance requirements for Ceph
  - Low latency for front end applications

SG
Storage Group

(intel)

# Agenda

- Recent common requirements for Ceph

- Middle Cache tiering solution

- Building block techniques
  - I/O optimization technique
    - DPDK for storage
  - Data Processing Acceleration techniques
    - ISA-L

- Conclusion

SG
Storage Group

(intel)

# Provide Mid-Tier Cache between applications and Ceph

- Protocol support:
  - Aim to target iSCSI/NFS/NVMF protocols.

- High performance
  - Provide local cache in each Mid-Tier node
  - Provide write log for data consistency.

- HA:
  - Replicate to 1+ additional nodes
  - Heartbeat for failed node detection



iSCSI/NVMF Initiator

Mid-Tier Cache
iSCSI/NVMF Target
Cache
Write Log
Ceph RBD

Mid-Tier Cache
iSCSI/NVME Target
Cache
Write Log
Ceph RBD

Ceph Cluster

# Agenda

- Recent common requirements for Ceph

- Middle Cache tiering solution

- Building blocks techniques

  - I/O optimization technique
    - DPDK for storage

  - Data Processing Acceleration techniques
    - ISA-L

- Conclusion

SG
Storage Group

(intel)

# DPDK for Storage **Overview**

- **Uses Intel® DPDK and UNS technology**
  - Optimized user space lockless polling technology in the NIC driver
  - Presents lock-light libraries and network TCP/IP services
- **Provides an enhanced software stack that optimizes iSCSI front end targets**
  - Optimizes packets in user space using lockless polling mechanisms
  - Reference software available for customer application integration to NVMe or other backend
- **Supports Linux* operating systems**
- **Enables a higher system level performance for iSCSI targets**
- **Currently available as reference software**

**User-space**

| CLOUD | NIC | Intel® DPDK NIC Driver | TCP IP (UNS) | iSCSI Target | Customer Storage App | NVMe Driver | NVMe |
|---|---|---|---|---|---|---|---|
| | | | | | | Mem Driver | DDR |
| | | | | | | CBDMA Driver | CBDMA |

**Intel® DPDK LIBRARIES**

← READ        WRITE →

Legend:
- Existing SW
- Enhanced SW
- Customer SW
- Linux* Kernel

*Other names and brands may be claimed as the property of others.

# High-Level Block Diagram



**INTEL® ETHERNET CONTROLLER**

**DPDK for Storage**

Intel® DPDK Polled Mode NIC Drivers

UNS TCP/IP (with sockets)

iSCSI

SCSI

SCSI/BDAL Translation

| BDAL API | BDAL API | BDAL API | BDAL API |
| --- | --- | --- | --- |
| malloc Block Driver | NVMe Block Driver | AHCI Block Driver | Linux* AIO Block Driver |

SCSI (sg) Generic SCSI Driver

Intel® DPDK Libraries

Copy Engine

CB-DMA Driver | nCPM Driver

CB-DMA | nCPM

NVME Controller | AHCI Controller | Linux* Kernel

intel inside™ XEON®

**Legend:**
- New Intel SW
- Existing Intel SW
- Third Party SW
- Linux* Kernel

# Intel DPDK for Storage **Benefits**

**Up to 7x Better Performance+** ······
- vs. Linux*-IO Target (LIO)
- Or 1/7 CPU overhead at same performance

**Up to 10x fewer cores utilized with the NVMe Driver** ······
- Vs. Linux* NVMe Driver

**Reduces Total Cost of Ownership** ······
- Reduce BOM costs between $80-500 by removing the need for a TOE
- Utilize free CPU cycles for other workloads

**Free Source Code** ······
- Customizable source code available as reference
- Evaluation source available upon request

**User Space Implementation** ······
- Portable/Upgradable and permissive licensing
- Requires Software License Agreement for full product use

# Intel DPDK for Storage
## Full **Packaging** and **Contents**

## Library Package Includes:

- Intel DPDK | UNS | Optimized Storage Stacks as reference software
- User space support code (written in C):
  - POSIX compliant
  - Demo/Usage, Unit test (functional correctness), Basic performance
- API manuals – may include links or copy key papers
- Release.txt (release notes, version, and library serial IDs)
- Linux* Support

## Source Agreement

- Source available under Restricted Use License Agreement Confidential (RULAC)
- Source code available under source license agreements (SLA)

# Intel® DPDK for Storage:
# Case Study1:  Performance Comparison with LIO

(intel)

# Intel® Xeon® Processor E5-2620v2-iSCSI Read/Write: 4 KB Data (performance per/core)

## NVM Express Backend

- 4 KB-Random-100% Write
- 4 KB-Random-100% Read
- 4 KB-Random-70% Read 30% Write

### PERFORMANCE

IO/s (in thousands)

| | DPDK for Storage | DPDK for Storage 2 Core | LIO 6 Core |
|---|---|---|---|

### PERFORMANCE/CORE

IO/s (in thousands)

| | DPDK for Storage | LIO |
|---|---|---|

**Up to 650% increase in max performance per core[+]**

SG
Storage Group

# Intel® DPDK for Storage
## Case Study2:  User space NVME driver(SPDK) Benefit

# 4 KB Random Read Performance: **4 x NVMe Drives**
## Single-Core Intel® Xeon® Processor



DPDK for Storage NVMe driver delivers up to **6x** performance improvement vs. Kernel NVMe driver with a single-core Intel® Xeon® processor

SG
Storage Group

# 4 KB Random Read Performance: **1-4 NVMe Drives**
# Single-Core Intel® Xeon® Processor



**DPDK for Storage NVMe driver scales linearly in performance from 1 to 4 NVMe drives with a single-core Intel® Xeon® processor**

SG
Storage Group

# Agenda

- Recent common requirements for Ceph

- Middle Cache tiering solution

- Building blocks techniques
  - I/O optimization technique
    - DPDK for storage
  - Data Processing Acceleration techniques
    - ISA-L

- Conclusion

(intel)

# Benefits of using **Intel ISA-L**

**Intel ISA-L**
enables Storage OEMs to obtain more performance from Intel CPUs and reduce investment in developing their own optimizations

Up to **7X** BANDWIDTH for Hash functions compared to OpenSSL algorithms

Allows **MAXIMUM** UTILIZATION of additional cores

Up to **4X** BANDWIDTH improvement on compression compared to zlib

**FASTER** TTM/ **LESS** RESOURCES than developing optimizations from scratch

Allows Intel to DEVELOP **OPTIMIZATIONS** that use new architectural enhancements that are TTM

# Intel® ISA-L Packaging and Contents

**Source Code Library**

**Single Core Low-Level Functions**
(OS independent functions)

**Supports 64-bit, Intel® Xeon® and Atom Processor**
E5-2600/2400 and Atom C2000 product family forward

**Gen Over Gen Function Updates**
to take advantage of new processor features

# Intel® ISA-L **Functions**

XOR (RAID 5), **P+Q** (RAID 6), Reed Solomon Erasure Code

**00..0** | **Data**

**Divisor** $n+1$ bits

Remainder

**CRC** $n$ bits

**Sender**

**CRC-T10**, **CRC-IEEE** (802.3), **CRC32-iSCSI**

**CRC** | **Data**

**CRC** | **Data**

**Divisor**

**Remainder**

Zero, accept
Non-zero, reject

**Receiver**

## DATA PROTECTION

## DATA INTEGRITY

# PERFORMANCE OPTIMIZING

## CRYPTOGRAPHIC HASHING

Dog

06d80e7
b0C50bs
49a509t
b49f249
24e8c8o
05x84q4

**Multi-Buffer:** SHA-1, SHA-256, SHA-512, MD5

## COMPRESSION "DEFLATE"

**IGZIP:** Fast Compression

## ENCRYPTION

plaintext

$e_B$ **Public** encryption key
$d_B$ **Private** encryption key

plaintext

Encryption Algorithm → Decryption Algorithm

Ciphertext

**Sender**

**Receiver**

**XTS-AES 128**,
**XTS-AES 256**

# Hashing Usage: Data Deduplication Optimizations (Fix Size)

**DEDUPLICATION ENGINE**

**DATA PROCESSING**

**010110010**
**0101010101**
**01110101**
**101010101**

**INCOMING DATA STREAM**

**Data Chunking**

| 0010 | 1010 | 0101 | 1100 | 1100 | 1101 | 1010 | 0010 |

Intel ISA-L
Multi-buffer Hashing
Algorithms

**SHA-1, SHA-256, SHA-512, MD5**

Intel ISA-L
Hashing Function
Stitching Algorithm

**Multi-hash-sha1+murmur3_128**

**Indexing**

| A | B | C | A | D | B | E | D |

| 0010 | 1010 | 0101 | 0010 | 1100 | 1010 | 1101 | 1100 |

**Store Data**

| A | B | C | D | E |

Up To **7X** Performance Over OpenSSL Algorithms

**Key:**
Intel ISA-L
3rdParty

# Hashing Usage: Data Deduplication Optimizations (Dynamic Size)

**INCOMING DATA STREAM**

0010
10101
101
0101

**DATA PROCESSING**

Intel ISA-L Rolling hash fingerprinting

**DEDUPLICATION ENGINE**

**Data Chunking**

| 001 | 01010 | 001 | 01010 | 1101 | 001 | 1100 | 1100 |

**Indexing**

| A | B | A | A | C | B | D | C |
|---|---|---|---|---|---|---|---|
| 001 | 01010 | 001 | 001 | 1100 | 01010 | 1101 | 1100 |

**Store Data**

| A | B | C | D |

**DATA PROCESSING**

Intel ISA-L
Multi-buffer Hashing Algorithms

**SHA-1, SHA-256, SHA-512, MD5**

Intel ISA-L
Hashing Function Stitching Algorithm

**Multi-hash-sha1+murmur3_128**

Up To **7X** Performance Over OpenSSL Algorithms

**Key:**
Intel ISA-L
3rdParty

# Intel ISA-L provides a solution to deploy Erasure Code (EC) with better performance, so that data replication can be done faster with half the space of other methods.

- Support any Matrixes: Vandermonde Reed-Solomon EC, Cauchy Reed-Solomon EC
- Support the different EC strategies: Local Reloadable Code EC, Regeneration Code EC, Hitchhiker Code EC



Erasure Code — Storage Capacity Needed

3X Replication — Storage Capacity Needed

D1 D2 D3 D4 D5 D6 D7 D8 D9

Reconstruct D1 and D4

P1 D2 D3 P2 D5 D6 D7 D8 D9

**DATA PROCESSING**

Intel ISA-L EC(9+3) Encode

**~10X** Performance Over Traditional Lookup Table Code

P1 P2 P3

Intel ISA-L EC(9+3) Decode

**DATA PROCESSING**

**~10X** Performance Over Traditional Lookup Table Code

D1 D4

**Key:** Intel ISA-L

**Source:** "Erasure Code and Intel® Intelligent Storage Acceleration Library"
http://www.intel.com/content/www/us/en/storage/erasure-code-isa-l-solution-video.html

SG Storage Group

(intel)

# Solving Real-World Problems: Qihoo 360

DEPLOYED
**INTEL-ISAL-BASED**
HDFS Raid for **INTEL-XEON-BASED** Cold Storage

EC Encode SPEEDS
**45X FASTER**
THAN JAVA VRS

EC Decode SPEEDS
**36X FASTER**
THAN JAVA VRS

REDUCED C O S T S
BY **25%~30%**

**Source:** Case Study "Intel and Qihoo 360 Internet Portal Datacenter - Big Data Storage Optimization Case Study"
https://software.intel.com/en-us/articles/intel-and-qihoo-360-internet-portal-datacenter-big-data-storage-optimization-case-study

# Solving Real-World Problems: Alibaba

**DEPLOYED**
**INTEL-ISAL-BASED** Sheepdog Erasure Code for **INTEL-ATOM-C2000-BASED** Cold Storage

Data Recovery SPEEDS
## 4X FASTER
THAN Sheepdog ZFEC

## 5X
CPU utilization reduction

(intel)

# Conclusion

- In this presentation, we introduce the storage optimization techniques provided by Intel for accelerating the Ceph performance:

- I/O optimization technique: DPDK for storage (SPDK)

- Data Processing Acceleration: ISA-L

- These kinds of building block techniques can help customers to accelerate the Ceph performance on IA platform

Q & A ?