Jin-Soo Kim (jinsoo.kim@snu.ac.kr) Systems Software & Architecture Lab. Seoul National University

Fall 2020

Solid-State Drives (SSDs)



Anatomy of an SSD

Samsung 850 Evo DRAM SSD Controller NAND Flash

(Messy) Storage Interfaces



Moving Closer to the Processor



Serial ATA (SATA)

- Primary internal storage interconnect for desktop and mobile PCs
 - Evolved from (Parallel) ATA
 - More than I.I billion SATA drives shipped during 2001-2008
 - Market share (as of 2008): Desktop (99%), Mobile PC (97.7%), Enterprise (27.6%)
- Serial, point-to-point, half duplex
- Why SATA?
 - Lower pin count (cost, space), Lower voltage support (5V \rightarrow 0.7V)
 - Higher performance: SATA 3 600MB/s @ 6Gbps
 - Simple drive configuration (no slave)
 - Greater reliability (CRC/packet)
 - Migration to servers (hot plug, NCQ, ...)

SATA NCQ

- Enqueue up to 32 commands in the drive
- Process them in an out-of-order fashion

Native Command Queuing

Requested Read: A, B, C, D NCQ Reordered Read: B, D, A, C

Legacy Command Non-Queued

Requested Read: A, B, C, D Non-reordered Read: A, B, C, D



D

Complete

(1.25 revolutions)

В

Α

С



NVMe (NVM Express)

- The industry standard interface for high-performance NVM storage
 - NVMe I.0 in 2011 by NVM Express Workgroup
 - NVMe 1.2 in 2014
- PCIe-based
- Lower latency
 - Direct connection to CPU
 - No HBA (Host Bus Adapter) required: reduced power and cost
- Scalable bandwidth
 - IGB/s per lane (PCle Gen3)
 - Up to 32 lanes



NVMe Overview

- Deep queue: 64K commands per queue, up to 64K queues
- Streamlined command set: only 13 required commands
- One register write to issue a command ("doorbell")
- Support for MSI-X and interrupt aggregation



NVMe SSD Form Factors



SSD Internals



The Unwritten Contract

Several assumptions are no longer valid

Assumptions	Disks	SSDs
Sequential accesses much faster than random	\bigcirc	\bigotimes
No write amplification	\bigcirc	\bigotimes
Little background activity	\bigcirc	\bigotimes
Media does not wear down	\odot	\bigotimes
Distant LBNs lead to longer access time	\bigcirc	\bigotimes

The Multi-streamed Solid-State Drive

(J.-U. Kang et al., HotStorage, 2014)

Some of slides are borrowed from the authors' presentation.

Effects of Write Patterns

Previous write patterns (= current state) matter



Stream



The Multi-streamed SSD

Mapping data with different lifetime to different streams



Working Example

• High GC efficiency \rightarrow Performance improvement



For effective multi-streaming, proper mapping of data to streams is essential!

Architecture



Case Study: Cassandra



Cassandra's Write Patterns

Write operations when Cassandra runs



Mapping #1: Conventional

Just one stream ID (= conventional SSD)



Mapping #2: Multi-App

Separate application writes (ID I) from system traffic (ID 0)



Mapping #3: Multi-Log

Use three streams; further separate Commit Log



Mapping #4: Multi-Data

Give distinct streams to different tiers of SSTables



Results: Conventional

- Cassandra's normalized update throughput
 - Conventional "TRIM off"



Results: Conventional with TRIM

- Cassandra's normalized update throughput
 - Conventional "TRIM on"



Results: Multi-App

- Cassandra's normalized update throughput
 - "Multi-App" (System data vs. Cassandra data)



Results: Multi-Log

- Cassandra's normalized update throughput
 - "Multi-Log" (System data vs. Commit-Log vs. Flushed data)



Results: Multi-Data

- Cassandra's normalized update throughput
 - "Multi-Data" (System data vs. Commit-Log vs. Flushed data vs. Compaction Data)



Results: GC Overheads

Cassandra's GC overheads



Results: Latency

- Cassandra's cumulated latency distribution
 - Multi-streaming improves write latency
 - At 99.9%, Multi-Data lowers the latency by 53% compared to Normal





- Mapping application and system data with different lifetimes to SSD streams
 - Higher GC efficiency, lower latency
- Multi-streaming can be supported on a state-of-the-art SSD and coexist with the traditional block interface
- Standardized in TIO SCSI (SAS SSDs) in 2015
- Standardized in NVMe 1.3 in 2017

Open-Channel SSD (OCSSD)

Why OCSSD?

I/O Isolation

Enable I/O isolation between tenants by allocating your SSD into separate parallel units.



Predictable Latency

No more guessing when an IO completes. You know which parallel unit is accessed on disk.



Data Placement & I/O Scheduling

Manage your non-volatile memory as a block device, through a filesystem or inside your application.

OCSSD Architecture



Zoned Namespace (ZNS) SSDs

Source: Matias Bjørling, How Zoned Namespaces Improve SSD Lifetime, Throughput, and Latency, FMS, 2020.

Zoned Storage Model

- Zones are laid out sequentially in an NVMe namespace
- The zone size is fixed and applies to all zones in the namespace
 - e.g., 512 MiB
- The command set inherits the NVMe Command Set
 - Built upon the conventional block interface (Read, Write, Flush and other commands)
 - Adds rules to collaborate on host and device data placement





Writing to a Zone

- "Sequential Write Required"
 - Must be written sequentially
 - Must be reset if written to again
- Each zone has a set of associated attributes:
 - Write pointer
 - Zone Starting LBA
 - Zone Capacity
 - Zone state
- Very similar to writing zones with host-managed SMR HDDs

Empty Zone





Partially Written Zone

Reading from a Zone

- Writes are required to be sequential within a zone
- Reads may be issued to any LBA within a zone and in any order





SMR HDDs and ZNS SSDs

- Host-managed SMR HDDs
 - Implements the SMR (ZAC/ZBC) specifications
 - ZAC: Zoned Device ATA Command Set in T13/SATA
 - ZBC: Zoned Block Commands in T10/SAS



NVMe ZNS SSDs

- Implements the Zoned Namespace Command Set specification
- Aligns with ZAC/ZBC to allow interoperability
- A single unified software stack support both storage types
 - Utilizes the already mature Linux storage stack built for SMR HDDs

Zoned Storage Software Stack

Linux Kernel's Zoned Subsystem (5.9)

User Space User Space Applications Applications without Zone Support without Zone Applications with Zone Support Applications with Zone Support Support Linux Kernel Regular File-systems /w Regular File-Systems /w File-Systems File-Systems Zone Support Zoned Storage (ext4,xfs) (f2fs, btrfs) (blobfs) (zonefs) Device Device Mapper Mapper (dm-zoned) (ftl) **Block Layer** bdev SCSI/ATA **NVMe** SCSI/ATA **NVMe** Ultrastar Ultrastar' DC ZN540 DC HC650 ZBC/ZAC SMR HDDs & NVMe ZNS SSDs 20₁₁

SPDK (20.10)

Biscuit: A Framework for Near-Data Processing of Big Data Workloads

(B. Gu et al., ISCA, 2016)

Some of slides are borrowed from the authors' presentation.

Moving Data

- In memory hierarchy
- Move data toward ALU to remedy long latency while accessing high-locality data



- In computation hierarchy
 - Move computation toward memory to remedy long latency while accessing lowlocality data



Hardware

SSD with a user-programmable NDP (Near-Data Processing) framework



[Inside of PM1725]

Item	Description	
Host interface	PCIe Gen.3 ×4 (3.2 GB/s)	
Protocol	NVMe 1.1	
Device density	1 TB	
SSD architecture	Multiple channels/ways/cores	
Storage medium	Multi-bit NAND flash memory	
Compute resources	Compute resources Two ARM Cortex R7 cores	
for Biscuit	@750MHz with MPU	
On-chip SRAM	< 1 MiB	
DRAM	\geq 1 GiB	

Hardware pattern matcher on each flash memory channel

Runtime

- Cooperative multithreading
 - A limited form of multithreading (fiber as a scheduling unit)
 - Less context switch overhead
 - Safe resource sharing without locking
- Shared nothing architecture
 - All data transmission among threads through I/O ports
 - Enforced by the programming model and APIs
- Dynamic loader for user programs
 - User program as position-independent code (PIC)
 - Symbol relocation to locate each program in a separate address space

Biscuit System Architecture



Development Process



TPC-H Results

- Data analytics on MariaDB
 - TPC-H dataset at a scale factor of 100 (160GiB)
 - 6.1x on average
 - 3.6x speedup for running all queries (two days vs. 13 hours)



Key-Value SSD (KVSSD)

Block: Parking Lot Structure

A driver (host) is responsible for parking (data management)



Object: Valet Parking

A parking facility (storage) is responsible for parking (data management)



KV Stores Common in Systems at Scale



KVSSD Prototype

Samsung KV-PM983



NGSFF KV SSD





Key Value Store is everywhere!



KVSSD Design

- Key size: up to 255B
- Value size: up to 2MB
- <u>https://github.com/OpenMPDK/KVSSD</u>



KVSSD Software Stack

RocksDB vs. KVSSD

- RocksDB
 - Originated by Facebook and actively used in their infrastructure
 - Most popular embedded NoSQL database
 - Persistent Key-Value store
 - Optimized for fast storage (e.g., SSD)
 - Uses Log-Structured Merge (LSM) Tree architecture
- KV Stacks on KVSSD
 - Benchmark tool directly operates on KVSSD through KV Stacks

RocksDB vs. KVSSD: Evaluation

Block SSD Clien	nt: kvbe	KV SSD		• Bet _ _ • IO	tter Performance Lean software stacks Overhead moved to device Ffficiency
RocksDB		Key Value API		-	Reduction of host traffic to devices
Filesystem	VS.	Key Value ADI	KV Sta	cks	
Block Driver		KV Driver		Hardware	Intel(R) Xeon(R) CPU E5-2630 v3 @ 2.40GHz 96 GB RAM PM983(Block) & KV-PM983 SSD
				Software	Ubuntu 16.04 RocksDB v5.0.2 on XFS 50M records, 16B Key, 4KB value
PM983		KV-PM983			

RocksDB vs. KVSSD: Results

- Random PUT performance
 - 8x more QPS (Query Per Second) with KV Stacks than RocksDB on block SSD
 - 90+% less traffic goes from host to device with KV SSD than RocksDB on block device

* Workload: 100% random put, 16-byte keys of random uniform distribution, 4KB-fixed values on single PM983 and KV-PM983 in a clean state

Computational Storage

SNIA: Computational Storage Instances

SNIA: Computational Storage Devices (CSx)

NGD Newport SSD

Flash Memory

Source: J. Do et al., "Cost-effective, Energy-efficient, and Scalable Storage Computing for Large-scale AI Applications, ACM ToS, 2020.

4190.568 Advanced Operating Systems | Fall 2020 | Jin-Soo Kim (jinsoo.kim@snu.ac.kr)

NGD Newport SSD: Software Stack

Source: J. Do et al., "Cost-effective, Energy-efficient, and Scalable Storage Computing for Large-scale AI Applications, ACM ToS, 2020.

4190.568 Advanced Operating Systems | Fall 2020 | Jin-Soo Kim (jinsoo.kim@snu.ac.kr)

Eideticom NoLoad CSP

- NVMe computational accelerators
 - Compression
 - Encryption
 - Erasure coding
 - Deduplication
 - Data analytics
 - AI and ML

Samsung SmartSSD

- PM983F
- Xilinx FPGA

