

#### Persistent Buffer Cache for High-performance Storage Systems

**OpenInfra Community Days Korea 2021** 

UNIST (NECSST lab.)

Hyunsub Song





Next-generation Embedded / Computer System Software Technology

CIENCE AND ENGINEER



Next-generation Embedded / Computer System Software Technology

CIENCE AND ENGINEE

### Persistent Memory (PM)

#### **Persistent Memory Features**

- Non-volatility
  - Byte-level random access
  - Fast access time (*nanoseconds*)



Storage



NECSST

Memory

#### APRIL 2ND, 2019 by Brian Beeler

-----

#### Intel Optane DC Persistent Memory Module (PMM)



Intel has talked about Optane DC Persistent Memory Modules (PMM) publicly for over a year now, espousing the benefits of a new tier of data centric architecture that sits between DRAM and Optane DC SSDs, with sequentially slower SSD and HDD media cascading down the pyramid to tape at the archive level. The goal with persistent memory has always been to move more data closer to the CPU, offering DRAM-like latency with storage-like persistence and capacities. After a year of listening to hardware and software partners talk about the benefits of persistent memory in the lab, with the release of the second generation Intel Xeon Scalable Processors, Optane DC PMEM is now available across a wide variety of server solutions.





Restance Research Next-generation Embedded / Computer System Software Technology

### **Evolution of critical path**





#### **Evolution of critical path**



100

10,03 03 10 10 10



#### **Evolution of critical path**





#### **Studies that consider PM as storage**

PM-dedicated file system and tiered PM file system

A CO 100 100 100





Remember Next-generation Embedded / Computer System Software Technology

#### **PM Targeted File Systems**

#### Designed to reap PM performance

		PM File System (e.g., DAX, NOVA)	<b>Tiered File System</b> (e.g., Strata, Ziggurat)
SOSP 2009	"BPFS (Better I/O Through Byte-Addressable, Persistent Memory)"		
SC 2011	"SCMFS (SCMFS: A File System for Storage Class Memory)"		
EuroSys 2014	"PMFS (System Software for Persistent Memory)"		
EuroSys 2014	"Aerie (Aerie: Flexible File-System Interfaces to Storage-Class Memory)"		
EuroSys 2016	"HiNFS (A High Performance File System for Non-Volatile Main Memory)"		Disk
SOSP 2017	"NOVA (NOVA-Fortis: A Fault-Tolerant Non-Volatile Main Memory File System)"	<u>/ PM/ PM</u>	/ (SSD, HDD)
SOSP 2017	"Strata (Strata: A Cross Media File System)"	DM ES	Tiored DM ES
HotStorage 2019	"EvFS (EvFS: User-level, Event-driven File System for Non- volatile Memory)"	FIWI F3	Tiered FWIF5
FAST 2019	"Orion (Orion: A Distributed File System for Non-Volatile Main Memory and RD	MA-Capable Networks)"	
FAST 2019	"Ziggurat (Ziggurat: A Tiered File System for Non- Volatile Main Memories and D	Disks)"	
SOSP 2019	"ZoFS (Performance and Protection in the ZoFS User-space NVM File System)"		
SOSP 2019	"SplitFS (SplitFS: Reducing Software Overhead in File Systems for Persistent Mer	mory)"	
FAST 2021	"KucoFS (Scalable Persistent Memory File System with Kernel-Userspace Collabo	oration)"	
Linux kernel	"DAX (Ext4-DAX, XFS-DAX)"		

app

app



#### But...

#### PM only

- PM as end destination media
- Replace traditional storage?
  - Exception: Strata and Ziggurat

### Lengthy process to maturity

- E.g., Ext4...still in progress
- Wisdom with age







Next-generation Embedded / Computer System Software Technology

SCHOOL OF COMPUTER

CIENCE AND ENGINEER

### First Responder (FR)

PM-based cache-like layer





- FR not only acts as a cache, but also as a storage using its persistent properties



#### Goals

SCHOOL OF COMPUTER

CIENCE AND ENGINEER



- AGE BRANKER

CISSR. Next-generation Embedded / Computer System Software Technology

NECSST

#### Goals

### PM performance

- Lightweight static management
- \* Average latency for managing cache for various indexing and management policies

Our static indexing	67ns
Hashing indexing	75ns

	Activity		Radix-tree + LRU	Hash + LRU
vs.		Hashing	-	78ns
	Radix-tree	Search	35ns	162ns
	/ Hash	Insert*	19µs	19µs
Ĺ		Delete	190ns	43ns
		Touch	161ns	153ns
	LRU	Add	73ns	65ns
		Remove	88ns	82ns

\* Insert includes mechanism to find empty blocks

#### Ensure durability/consistency

• Static protocol naturally fulfills this







#### Internal components of FR



- Chunk: Actual data is stored
- Tag: Some file's information and the status of chunk are stored
  - Key used for indexing: *key* mod Floor(N/2)
  - Bits in Status flag: V (Valid), N (New)



#### SIL NECSST

m Software Research Next-generation Embedded / Computer System Software Technology

#### Design

#### Layout of FR



- Static placement/replacement scheme in FR
  - Every files have destined location within FR with no PM allocator
    - $\rightarrow$  Can result in higher miss rate and collision



#### Design



- Stride: To eliminate invasion in chunk as much as possible
  - Files are positioned apart from each other by stride
- Periodic Flush: To reduce penalty (for clean chunk, there is no penalty for collision)
  - Data is written to chunk, is flushed in the background

#### NECSS1

Material Booksaire Research Next-generation Embedded / Computer System Software Technology

### Indexing

• *Key* mod N  $\rightarrow$  *key* mod Floor(N/2)

#### Cases

- S[C(0, 0), C(0, 0)]: the slot is empty
- S[C(1, 1), C(0, 0)]: only one chunk has valid data
- S[C(1, 1), C(1, 0)]: both chunks have valid data





#### Case 1





Case 2







20 / 47

Case 3





#### Failure recovery in FR

#### Key point

• If a fault occurs at any step, the recovered state will either be one of the initial states of the figures or have already completed the intended write



#### Failure recovery in FR



Case 3-2



#### Failure recovery in FR



Case 3-2



#### **Performance evaluation**

SCHOOL OF COMPUTER

CIENCE AND ENGINEE

#### System configuration

	Description
CPU	Intel(R) Xeon(R) Gold 6242 CPU @ 2.80GHz (16 cores)
DRAM	Samsung 32GB 2666MHz DDR4 RDIMM×4 (128GB)
PM	Intel Optane DC Persistent Memory (128GB)
Storage	Samsung V-NAND SSD 860 EVO (1TB)
OS	Linux Ubuntu 18.04.3 LTS (64bit) kernel v4.18

#### Description of experimental comparison

Notation	Description	Configuration		
		РМ	DRAM	Backing storage
FR-X	FR applied to Ext4 (X is period value, e.g., 10ms)	128GB		1TB (SSD)
Ext4	Traditional block-based file system		128GB	1TB (SSD)
DM-WC	DM-Writecache applied to Ext4	128GB	128GB	1TB (SSD)
DAX	PM-aware file system developed based on Ext4			128GB (PM)
NOVA	PM-aware file system			128GB (PM)





#### Benchmarks

Filebench	R:W	Mean file size	# of files	Key-value store	Data set size	R:W
Fileserver	1:2	128KB	200K	YCSB-A,-F	4GB	1:1
Varmail	1:1	32KB	800K	YCSB-B,-D,-E	4GB	19:1
OLTP	1:1	1.5GB	20	YCSB-C	4GB	1:0

- Filebench
  - Fileserver: write-intensive workload *without* fsync() calls
  - Varmail and OLTP: have considerable number of fsync() calls
- YCSB (record selection for -D is Latest, while all others are Zipfian)
  - Application: RocksDB
  - -A, -F: write-intensive workloads
  - -C: read-only workload
  - B, -D, -E: read-intensive workloads



#### Overall performance



#### Observations

- [Filebench] For Varmail and OLTP, FR is roughly 94x and 8x better than Ext4 and roughly 4.5x and 1.5x better than DM-WC

FR performance is slightly lower than NOVA, while DAX suffers for Varmail

- [YCSB] FR, NOVA, and DAX show similar performance

Ext4 and DM-WC perform worst for YCSB-A, -B, and -C



#### Overall performance



#### Observations

- [Filebench] For Varmail and OLTP, FR is roughly 94x and 8x better than Ext4 and roughly 4.5x and 1.5x better than DM-WC

FR performance is slightly lower than NOVA, while DAX suffers for Varmail

- [YCSB] FR, NOVA, and DAX show similar performance

Ext4 and DM-WC perform worst for YCSB-A, -B, and -C



#### Overall performance



#### Observations

- [Filebench] For Varmail and OLTP, FR is roughly 94x and 8x better than Ext4 and roughly 4.5x and 1.5x better than DM-WC

FR performance is slightly lower than NOVA, while DAX suffers for Varmail

- [YCSB] FR, NOVA, and DAX show similar performance

Ext4 and DM-WC perform worst for YCSB-A, -B, and -C



#### **Effect of PM size**

SCHOOL OF COMPUTER

CIENCE AND ENGINEER

- Performance results for PM size of 2<sup>x</sup>GB, where x is value of points in x-axis
  - Normalized to the performance of FR when *x* = 7 (128GB)



30 / 47

#### Dynamic workload



- Standard workloads do not capture the dynamics of real-world workloads
  - In terms of pattern
    - Standard workloads: Working set does not change with time
    - Real-world workloads: Working set grows and shrinks as time evolves
  - In terms of operation generation
    - Standard workloads: No change in the access intensities of working set over time
    - Real-world workloads: Access intensities are also vary with time
- $\rightarrow$  Need for **dynamic workload** that is more representative of real-world workloads



Working set of standard workload (Fileserver)



### Dynamic workload

We devise synthetic workloads using I/O testing tool FIO

	Description
File size	1~14GB [whole numbers only: 1 to 2 files of each]
Distribution	Pareto 0.1/0.5 [2/2], Zipf 0.2/1.2 [5/1], N [6], R [5]
IO unit	4KB [14], 8KB [7]
Read:Write	1:0 [4], 19:1 [9], 1:1 [8]
Fsync	$0 \sim 5\%$ [whole numbers only: 2 to 5 files of each]
Intensity	Request interval: $1\mu$ s $\sim$ 1ms (randomly distributed)

(a) Characteristics of 21 files used to generate synthetic workload (FIO)



(b) Working set (FIO-12)

#### **Configuration: Total IO size is 575GB**

- FIO-6: 6 files, 50GB working set
- FIO-12: 12 files, 100GB working set



System Software Research Next-generation Embedded / Computer System Software Technology

### **Dynamic workload**

Performance results



#### **Observations**

- FR provides more than 9x higher aggregate throughput and ended over 3x faster than Ext4
- FR is providing immediately durable in-order semantics
- For NOVA and DAX, cannot run as dataset is larger than PM size





Next-generation Embedded / Computer System Software Technology

SCHOOL OF COMPUTER

CIENCE AND ENGINEER

### FR on other storage platforms





#### **Step 1: FR on Ceph client**

FR is applied to the client server running 3 OSDs as block device

	Description	Ļ	Λ_
CPU	Intel(R) Xeon(R) Gold 6242 CPU @ 2.80GHz (16 cores)	 n	4
DRAM	Samsung 32GB 2666MHz DDR4 RDIMM×4 (128GB)	d	
PM	Intel Optane DC Persistent Memory (512GB)	<u>d</u>	
Storage	Samsung V-NAND SSD 860 EVO (1TB)	ň	3 -
OS	Linux Ubuntu 18.04.3 LTS (64bit) kernel v4.18	0	-





NECSST



#### **Observations**

- Media performance: NVMe SSD > SATA SSD > network-connected OSDs > HDD
- FR and Ext4 have performance changes depending on the media performance
- DM-WC still shows similar performance



RIGS

SCHOOL OF COMPUTER

#### FR module can be applied everywhere!!!



### What if you apply FR to Ceph?

**NECSST** 

- High performance
- High reliability









Next-generation Embedded / Computer System Software Technology

ceph

SCHOOL OF COMPUTER SCIENCE AND ENGINEERIN

### Application point

 Object Storage Daemon (OSD)





\* "File Systems Unfit as Distributed Storage Backends: Lessons from 10 Years of Ceph Evolution," In Proc. SOSP'19

Architecture of Ceph RADOS Block Device (RBD)



- Application point: "BlueStore"
  - OSD -> Storage backend (BlueStore)





System Software Research Next-generation Embedded / Computer System Software Technology

#### Cache in BlueStore

- BlueStore's own write-through cache in user space
- Uses the 2Q and LRU algorithm



Next-generation Embedded / Computer System Software Technology



- FR in BlueStore
  - Critical path can be performed with PM performance





System Software Research Next-generation Embedded / Computer System Software Technology

#### **Step 2: In progress**

# FR on Ceph

- Analysis of target
  - BlueStore, OnodeCache code
- Implementation of FR
  - Modularizing of FR in user-level (C  $\rightarrow$  C++)
  - Allocation user-level PM pool
  - Applying FR to BlueStore
- Debugging
- Validation of FR effect
- Devising ways to overcome shortcoming





Next-generation Embedded / Computer System Software Technology

SCHOOL OF COMPUTER

CIENCE AND ENGINEE

#### Summary

#### First Responder (FR)

- PM-based cache-like layer
- Keep legacy storage system and storage media "as-is"
- PM performance
  - Respond quickly with in-order semantics
  - Hide traditional I/O stack overhead
- Ensure durability/consistency
  - Protocol implemented with static management



## Thank you!!!



