PolyStore: Exploiting Combined Capabilities of Heterogeneous Storage

Yujie Ren David Domingo Jian Zhang Paul John Rekha Pitchumani Sanidhya Kashyap Sudarsun Kannan







Applications demand storage capacity and bandwidth



https://dl.acm.org/doi/10.1109/SC.2018.00068

Evolving storage media with specialized file systems



Use them collaboratively to reduce cost and retain performance!

Hiding latency with hierarchical design philosophy

Storage devices present diverse characteristics in latency, capacity, and price



Prompt admission upon cache miss

Coarse-grained data movement

"Faster" storage absorbs write requests of new data

Pitfalls in existing heterogeneous storage systems

Pitfall I: Cannot utilize combined bandwidth of heterogeneous storage



Hierarchical device layout limits cumulative bandwidth utilization!

Pitfalls in existing heterogeneous storage systems

Pitfall 2: Cannot provide device-specific DRAM buffering policies



Lack mechanisms for applying device-specific policies for DRAM buffering!

Pitfalls in existing heterogeneous storage systems

Pitfall 3: Cannot capitalize on mature device-optimized file systems

- New tiered file systems for certain devices
 Strata [SOSP' 17], Ziggurat [FAST' 19]
- Managing in block layer with one file system
 Orthus [FAST' 21], Bcache [Linux]
- Leaving mature device-optimized FSes underutilized
 NOVA [FAST' 16], F2FS [FAST' 15], ext4, XFS [Linux]





Coupling data placement decisions within the software stack!

Naïve solution with RAID 0 for heterogeneous devices?

Combined bandwidth utilization

X Cannot capitalize on deviceoptimized file systems!

X Cannot provide device-specific DRAM buffer cache policies!



Ideal system for managing heterogeneous storage

Combined bandwidth utilization Data placement decisions above the storage software stack Reuse mature file systems **H/W** optimized tailored for storage devices file systems ✓ Flexible DRAM buffer cache "Faster" storage admission/eviction



Our solution - PolyStore

A meta-layer on top of device-optimized kernel-level file systems



Indexing data across device-optimized file systems

Indexing data across device-optimized file systems

Distribute data across heterogeneous storage devices

- Large files (e.g., streaming files) are bandwidth intensive
- Map a logical file to physical files with a scalable indexing structure
- Encode data placement and access patterns across devices



Key data structure for bandwidth utilization and flexible DRAM buffer cache!

Our solution - PolyStore

A meta-layer on top of device-optimized kernel-level file systems



Dynamic data placement mechanism for combined bandwidth

Challenge - mapping I/O from application threads to storage devices

• Storage bandwidths show diverse characteristics

 Statically identifying the optimal mapping cannot adapt to workload changes in applications



Need dynamic data placement to maximize the storage bandwidth utilization!

Epoch-based throughput monitoring and dynamic data placement

Application Threads















Epoch-based throughput monitoring and dynamic data placement



Converge to a local maximal configuration utilizing combined bandwidth!

Our solution - PolyStore

A meta-layer on top of device-optimized kernel-level file systems



Heterogeneity-aware DRAM buffering

Hide access latency for heterogeneous devices





storage

storage

Hide access latency for heterogeneous devices

• Asymmetrical performance in heterogeneous storage devices







Hide access latency for heterogeneous devices

- Asymmetrical performance in heterogeneous storage devices
- Indexing structure encodes data placement information



Hide access latency for heterogeneous devices

- Asymmetrical performance in heterogeneous storage devices
- Indexing structure encodes data placement information

 Allocate DRAM cache buffer with device-specific admission & eviction policies













 PM read has same order of speed as DRAM



Logical file

- PM read has same order of speed as DRAM
- Do not buffer data in DRAM for read-only data for PM



- PM read has same order of speed as DRAM
- Do not buffer data in DRAM for read-only data for PM
- PM write is slower than read



- PM read has same order of speed as DRAM
- Do not buffer data in DRAM for read-only data for PM
- PM write is slower than read
- Allocate DRAM buffer for PM only when data is modified





Facilitating data migration decisions across devices

• Track data hotness/coldness



- Track data hotness/coldness
- Hot data on slow device selected as victim due to memory pressure



- Track data hotness/coldness
- Hot data on slow device selected as victim due to memory pressure
- Flush hot data on slower device to faster one if space permits



- Track data hotness/coldness
- Hot data on slow device selected as victim due to memory pressure
- Flush hot data on slower device to faster one if space permits



- Track data hotness/coldness
- Hot data on slow device selected as victim due to memory pressure
- Flush hot data on slower device to faster one if space permits
- Do garbage collection for the original data block on slower device



Our solution - PolyStore

A meta-layer on top of device-optimized kernel-level file systems



More technical details in our paper!

Evaluation

- Can PolyStore utilize the combined bandwidth of storage devices?
- Can PolyStore improve performance for real-world applications?



Multi-thread benchmark access private 2GB files (O_DIRECT flag)



- Scalable indexes distribute data across PM and NVMe
- Dynamic data placement effectively utilize the combined bandwidth

32 benchmark threads access private 2GB files (O_DIRECT flag)



Sequential append

32 benchmark threads access private 2GB files (O_DIRECT flag)



32 benchmark threads access private 2GB files (O_DIRECT flag)



Static coarse-grained data placement

NVMe

PM

32 benchmark threads access private 2GB files (O_DIRECT flag)



Static coarse-grained data placement

NVMe

PM

32 benchmark threads access private 2GB files (O_DIRECT flag)





32 benchmark threads access private 2GB files (O_DIRECT flag)



\bigcirc	Benchmark	Τł	nr	ea	ds	5	
Initial	\$ \$ \$ \$ \$		Ş	Ş	Ş	Ş	Ş
	Finish		Ş	Ş	Ş	Ş	Ş
			Ş	Ş	Ş	Ş	Ş
			Ş	Ş	Ş	Ş	Ş
V				ne sous-stati	anve 1		(intel)
	PM			N	١V	M	е

32 benchmark threads access private 2GB files (O_DIRECT flag)



C	Benchmark	Thr	ea	ds	;		
Initial	\$ \$ \$ \$ \$	Ş	Ş	Ş	Ş	Ş	
	Finish	ş	Ş	Ş	Ş	Ş	
PM bandwidth unutilized!	PM	Ş	Ş	Ş	Ş	Ş	
	bandwidth unutilized!	Ş	Ş	Ş	Ş	Ş	
				28554 1		(intel)	
	PM	1	Ν	١V	M	е	

32 benchmark threads access private 2GB files (O_DIRECT flag)



Benchmark Threads							
Initial	\$ \$ \$ \$ \$	Ş	Ş	Ş	Ş	Ş	
	Finish	ş	Ş	Ş	Ş	Ş	
	PM	ş	Ş	Ş	Ş	Ş	
	bandwidth unutilized!	Ş	Ş	Ş	Ş	Ş	
				201112 		(intel)	
	PM	T	Ν	١V	M	е	

32 benchmark threads access private 2GB files (O_DIRECT flag)



Sequential append

32 benchmark threads access private 2GB files (O_DIRECT flag)



Sequential append

PolyStore data indexes + dynamic data placement

32 benchmark threads access private 2GB files (O_DIRECT flag)



PolyStore data indexes + dynamic data placement

PM

32 benchmark threads access private 2GB files (O_DIRECT flag)



PolyStore data indexes + dynamic data placement

NVMe

PM

32 benchmark threads access private 2GB files (O_DIRECT flag)



PM

PolyStore data indexes + dynamic data placement

NVMe

32 benchmark threads access private 2GB files (O DIRECT flag)



PolyStore data indexes + dynamic data placement

٦.		
Ľ-	SOLD-STOTE BRINE	linte
1	N۷	′Me

PM

\$ \$ \$ \$ \$ \$

<u>}</u> } <u>}</u> } <u>} </u>

333

32 benchmark threads access private 2GB files (O_DIRECT flag)



C	Benchmark	Threads
Initial	\$ \$ \$ \$ \$	\$ \$ \$ \$ \$
	Finish	\$ \$ \$ \$ \$
	\$ \$	\$ \$ \$
	\$ \$ \$	\$ \$
ent	PM	NVMe

PolyStore data indexes + dynamic data placement

32 benchmark threads access private 2GB files (O_DIRECT flag)



Enabling DRAM buffer cache for all approaches



Enabling DRAM buffer cache for all approaches



• The OS page cache cannot handle device characteristics

Enabling DRAM buffer cache for all approaches



• The OS page cache cannot handle device characteristics

Enabling DRAM buffer cache for all approaches



• The OS page cache cannot handle device characteristics

Using OS page cache

Enabling DRAM buffer cache for all approaches



• The OS page cache cannot handle device characteristics

PolyStore DRAM buffer cache

Enabling DRAM buffer cache for all approaches



- The OS page cache cannot handle device characteristics
- PolyStore DRAM buffer cache
- PolyStore provide flexible data admission/eviction control

Enabling DRAM buffer cache for all approaches



- The OS page cache cannot handle device characteristics
- PolyStore DRAM buffer cache
- PolyStore provide flexible data admission/eviction control





Hierarchical design philosophy no longer suits bandwidth-intensive applications

PolyStore - a meta layer on top of mature device-optimized file systems

- Scaling device bandwidth horizontally with dynamic data placement
- DRAM buffering mechanism adapting to device characteristics

https://github.com/RutgersCSSystems/PolyStore



