

ZUFS

How to Develop an Efficient PM-based File System in User Space

Dr. Golander Amit, NetApp

Storage Media Generations





PM marries the best of both worlds:



PM Hardware



SCM (Storage Class Memory) **PM** (Persistent Memory) Byte-addressable Media Byte-addressable Device @ Near-memory speed @ Near-memory speed, on memory bus - (page) 0.1us DRAM Cache **Byte Addressable Block Addressable** CPU 0.1us <1us Capacity PM 1000 FLASH Persistent 100 Volatile SDXP STT-MRAM 10 NVDIMM-I PCIe Controller NIC **NVMe** 100us 1 SAS HBA 0.1us 1us 100us 1000us Latency 5000us HDD SSD 1000us

Rounded latency numbers & under typical load

NetApp



© 2018 SNIA SDC. All Rights Reserved.

Background



Plexistor (acquired by NetApp)

- PM-based FS pioneer since 2013
- Contributing some of our IP

Our PM-based FS approach:

- Support legacy applications
 & Enable NPM (e.g. SPDK)
- Feature rich
- Integrate with NetApp product portfolio



Kernel Vs. User Space FS Implementation

Kernel	User space
Fast (shortest path)	Portable
	Resilient (contained)
	Simpler to add functionality & Debug
	Fewer licensing restrictions

The gap: Near-memory speed Kernel-to-User bridge





Why not extend FUSE to PM?

SDCE SNIAEMEA FEBRUARY 2018 TEL AVIV, ISRAEL STORAGE DEVELOPER CONFERENCE

FUSE architecture is great for HDDs and ok(ish) for SSDs, but not suitable for PM

\$/G	HDD	Flash	PM Memory Latency
ψ, O	TCP FUSE		RDMA ?
Design Assumptions		FUSE	ZUFS
	Typical medias	Built for HDDs & extended to Flash	Built for PM/NVDIMMs and DRAM
	SW Perf. goals	 Secondary (High latency media) Async I/O Throughput 	SW is the bottleneckLatency is everything
	SW caching	Slow media -> Rely on OS Page Cache	Near-memory speed media -> Bypass OS Page Cache
	Access method	I/O only	I/O and mmap (DAX)
	Cost of redundant copy / context switch	Negligible	The bottleneck -> Avoid copies, queues & remain on core
	Latency penalty under load	100s of µs	3-4 µs

ZUFS Features & Architecture

Low latency & Efficient

- Core & L1 cache affinity
- Zero data copy

Manages devices

- Optimal pmem access
- NUMA aware
- Data mover to lower tier devices
- Page table mapping supports
 I/O & DAX semantics
- Misc
 - Async hook available
 - System service







FUSE Vs ZUFS Penalty (PM, DirectIO)



- Measured on
 - Dual socket, XEON 2650v4 (48HT)
 - DRAM-backed PMEM type
- Random 4KB DirectIO write access





ZUFS is a Kernel-to-User bridge designed for PM

Enables NetApp solutions

Open Source, being contributed upstream

- Hope to accelerate PM adoption and innovation
- https://github.com/NetApp/zufs-zus & zufs-zuf
- You're welcome to use, review and contribute code



Thank you

© 2018 SNIA SDC. All Rights Reserved.

