

FLASH IN THE DATACENTER

Nisha Talagala

Copyright [©] Fusion-io, Inc. All rights reserved.



NON VOLATILE MEMORY

Flash

- 100s GB to 10 TB per PCIe device
- Media trend increase in density, reduction of write cycles, SLC/MLC/3BPC
- 100s of thousands to millions of IOPS, GB/s of bandwidth

PCM/MRAM/STT/Other NVMs

- Still in research
- Potential of extreme performance increase



谷

HOW TO EFFECTIVELY USE FLASH?

Performance

- Closer to CPU highest bandwidth, lowest latency
- Server (compute) side flash complements storage side flash

Hierarchy of DRAM, flash, disk

Disk displacement usages

- Caches server and storage side
- Scale out and cluster file systems
 - flash in metadata server
 - storage server
- Staging, checkpoint

DRAM displacement usages

Improved paging, semi-external memory











SC11 GENERAL PARALLEL FILE SYSTEM (GPFS) DEMO



- 24 uncompressed 1080p videos (up to 6 GB/s of data)
- Five Fusion Powered GPFS-based NSD servers
- Three visualization workstations





- Five 0.5U NSD Servers, each with six core dual socket CPUs, 12 GB RAM, InfiniBand HCA, and an ioDrive2
- S Visualization workstations, with an InfiniBand HCA and an NVIDIA graphics card
- 36-port QDR InfiniBand switch



Your Server Becomes a Shared Flash Storage Appliance







- Support
 - ION Software (via Fusion-io)
 - Server (via server OEM)
 - ioDrive (via your supplier)
- PDSW 8 Supercomputing 13





NVM (FLASH, OTHER) IS DIFFERENT FROM DISK

FUSION-IO

Area	Hard Disk Drives	Flash Devices
Logical to Physical Blocks	Nearly 1:1 Mapping	Remapped at every write
Read/Write Performance	Largely symmetrical	Heavily asymmetrical. Additional operation (erase)
Sequential vs Random Performance	100x difference. Elevator scheduling for disk arm	<10x difference. No disk arm – NAND die
Background operations	Rarely impact foreground	Regular occurrence. If unmanaged - can impact foreground
Wear out	Largely unlimited writes	Limited writes
IOPS	100s to 1000s	100Ks to Millions
Latency	10s ms	10s-100s us



MULTI-QUEUE I/O IN LINUX*

- Extending Linux block I/O to support NVM performance
- Multi-queue
 - Software queues, Hardware queues
 - Per CPU issue/completion, multi-socket scaling
 - Matches inherent parallelism in NVM devices and CPUs
 - Supports upcoming queue oriented standards models
- Performance
 - 3.5x 10x increase in IOPS (from ~1M to 3.5-10M)
 - 10x 38x reduction in I/O stack latency

*Linux Block I/O: Introducing Multiqueue SSD Access on Multicore Systems Bjorling M., Axboe J., Nellans D., Bonnett P. SYSTOR 2013 University of Copenhagen and Fusion-io

DIRECT-ACCESS I/O THROUGH NATIVE INTERFACES

X



PDSW 8 – Supercomputing 13

FUSION-10



FLASH PRIMITIVES: SAMPLE USES AND BENEFITS

FUSION-ic

Databases

Transactional Atomicity: Replace various workarounds implemented in database code to provide write atomicity (MySQL double-buffered writes, etc.)

Filesystems

File Update Atomicity: Replace various workarounds implemented in filesystem code to provide file/directory update atomicity (journaling, etc.)

- 98% performance of raw writes
 Smarter media now natively understands atomic updates, with no additional metadata overhead.
- 2x longer flash media life Atomic Writes can increase the life of flash media up to 2x due to reduction in write-ahead-logging and doublewrite buffering.
- 50% less code in key modules
 Atomic operations dramatically reduce application logic, such as journaling, built as work-arounds.



MYSQL EXAMPLE: LATENCY IMPROVEMENT

K



PDSW 8 – Supercomputing 13

SION-IO

MYSQL EXAMPLE: THROUGHPUT IMPROVEMENT



谷

KEY-VALUE INTERFACE: SAMPLE USES AND BENEFITS

NoSQL Applications

Increase performance by eliminating packing and unpacking blocks, defragmentation, and duplicate metadata at app layer.

Reduce application I/O through batched operations.

Reduce overprovisioning due to lack of coordination between two-layers of garbage collection (application-layer and flash-layer). Some top NoSQL applications recommend over-provisioning by 3x due to this.

- Near performance of raw device Smarter media now natively understands a key-value I/O interface with lock-free updates, crash recovery, and no additional metadata overhead.
- **3x throughput on same SSD** Early benchmarks comparing against synchronous levelDB show over 3x improvement.
- Up to 3x capacity increase
 Dramatically reduces over-provisioning through coordinated garbage collection and automated key expiry.



谷

MEMORY-ACCESS THROUGH NATIVE INTERFACES



PDSW 8 – Supercomputing 13

FUSION-10



GRAPH500* AND DI-MMAP**

- Traversing massive graphs
 - "Using 2.56TB of Fusion-io NAND flash to access data using memory semantics, LLNL's new Graph500 algorithm can process graphs 8x larger than before with only a 50% performance degradation compared to an all DRAM system."
 - Results: 55.6 MTEPS (Million Traversed Edges Per Second)
 4 x 640GB Fusion-io MLC
- DI-MMAP: Accelerated mmap for highly concurrent apps
 - 3-5x improvement in mmap performance

* Graph500: Traversing massive graphs with NAND flash; Pearce, Gokhale, & Amato (LLNL) **DI-MMAP: A High Performance Memory Map Runtime for Data Intensive Applications; Van Essen, Hsieh, Ames, Gokhale (LLNL)







谷

COMPARING I/O AND MEMORY ACCESS SEMANTICS

I/O	 I/O semantics examples: Open file descriptor – open(), read(), write(), seek(), close() (New) Write multiple data blocks atomically, nvm_vectored_write() (New) Open key-value store – nvm_kv_open(), kv_put(), kv_get(), kv_batch_*()
Memory Access (Volatile)	 Volatile memory semantics example: Allocate virtual memory, e.g. malloc() memcpy/pointer dereference writes (or reads) to memory address (Improved) Page-faulting transparently loads data from NVM into memory
Memory Access (Non- Volatile)	Non-volatile memory semantics example: (New) Allocate and manage persistent memory

PDS₩ 8^{3,20}Supercomputing 13

FUSION-10





1ST CONTRIBUTION: FLASH PRIMITIVES



Flash programming primitives

Use built-in characteristics of the Flash Translation Layer to perfrom journal-less updates (more performance and less flash wear = lower TCO)

Repository Learn More

On GitHub:

- API specifications, such as:
 - nvm_atomic_write()
 - nvm_batch_atomic_operations()
 - nvm_atomic_trim()
- Sample program code

https://opennvm.github.io

PDSW 8 – Supercomputing 13

usion-io



2ND CONTRIBUTION: LINUX FAST-SWAP



https://opennvm.github.io

On GitHub:

- Documentation
- Experimental Linux kernel with virtual memory swap patch (3.6 kernel)
- Benchmarking utility

PDSW 8 – Supercomputing 13

SiON-iO

3RD CONTRIBUTION: KEY-VALUE INTERFACE



On GitHub:

- API specifications, such as: nvm_kv_put()
 - nvm_kv_get()
 - nvm_kev_batch_put()
 - nvm_kv_set_global_expiry()
- Sample program code
- Benchmarking utility
- Source code for flash optimized key value store

SiON-iO



PDSW 8.20 Supercomputing 13

谷

OPENNVM, STANDARDS, AND CONSORTIUMS

- opennvm.github.io
 - Primitives API specifications, sample code
 - Linux swap kernel patch and benchmarking tools
 - key-value interface API library and code, sample usage code, benchmark tools
- INCITS SCSI (T10) active standards proposals:
 - SBC-4 SPC-5 Atomic-Write <u>http://www.t10.org/cgi-bin/ac.pl?t=d&f=11-229r6.pdf</u>
 - SBC-4 SPC-5 Scattered writes, optionally atomic <u>http://www.t10.org/cgi-bin/ac.pl?t=d&f=12-086r3.pdf</u>
 - SBC-4 SPC-5 Gathered reads, optionally atomic <u>http://www.t10.org/cgi-bin/ac.pl?t=d&f=12-087r3.pdf</u>
- SNIA NVM-Programming TWG draft guide: <u>http://snia.org/forums/sssi/nvmp</u>



