# Creating Flash-Aware Applications

Nisha Talagala

# NVM (Flash, other) is different from Disk

| 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 |

# I/O and Memory Access for Flash Aware Applications

| | |
|---|---|
| **I/O** | I/O semantics examples:<br>• Open file descriptor – open(), read(), write(), seek(), close()<br>• (New – presented today) NVM Primitives<br>• (New – presented today) NVM KV Store |
| **Memory** | Volatile memory semantics example:<br>• Allocate virtual memory, e.g. malloc()<br>• memcpy/pointer dereference writes (or reads) to memory address<br>• (Improved – presented today) Page-faulting transparently loads data from NVM into memory |

https://opennvm.github.io

# OpenNVM

Welcome to the open source project for creating new interfaces for non-volatile memory (like flash).

GNU Public License v2.0

http://www.opencompute.org/projects/storage/

# 3 Contributions to the Community

**FUSiON-iO**

## Current OpenNVM Repositories

### Flash-aware Linux swap

When working set size exceeds the capacity of DRAM, demand page from a flash-aware virtual memory subsystem.

Repository | Learn More

### Key-value interface to flash

Create NoSQL databases faster. Automate garbage collection of expired data.

Repository | Learn More

### 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

https://opennvm.github.io

# 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

# Flash Primitives: Sample Uses and Benefits

FUSiON-iO

▸ ## Databases

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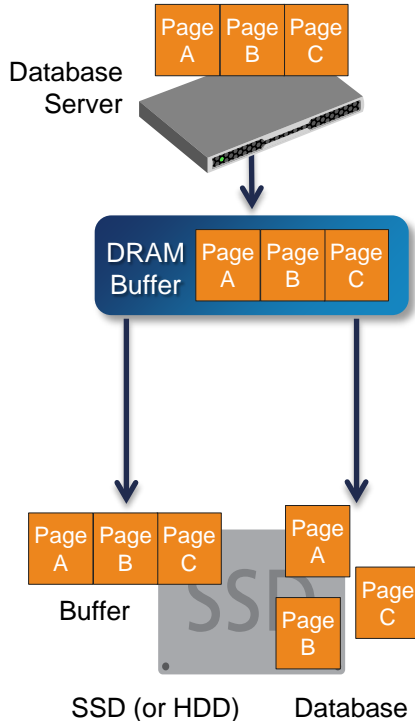provide write atomicity (MySQL
double-buffered writes, etc.)

▸ ## 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 double-write buffering.

- **50% less code in key modules**
  Atomic operations dramatically
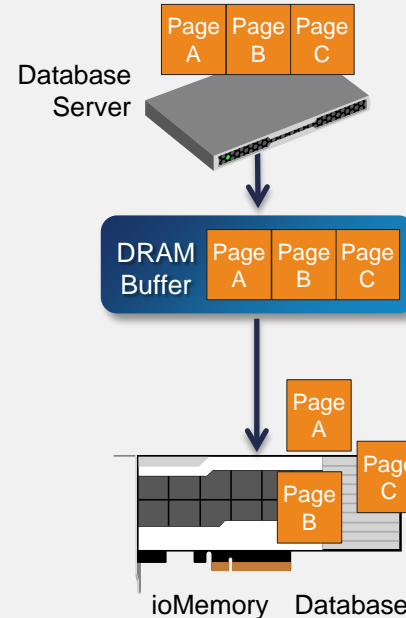  reduce application logic, such as
  journaling, built as work-arounds.

# Atomic Writes – MySQL Example

## Traditional MySQL Writes

Database Server

| Page A | Page B | Page C |

DRAM Buffer | Page A | Page B | Page C |

| Page A | Page B | Page C |

Buffer

SSD (or HDD)

Page A
Page C
Page B

Database

1 | Application initiates updates to pages A, B, and C.

2 | MySQL copies updated pages to memory buffer.

3 | MySQL writes to double-write buffer on the media.

4 | Once step 3 is acknowledged, MySQL writes the updates to the actual tablespace.
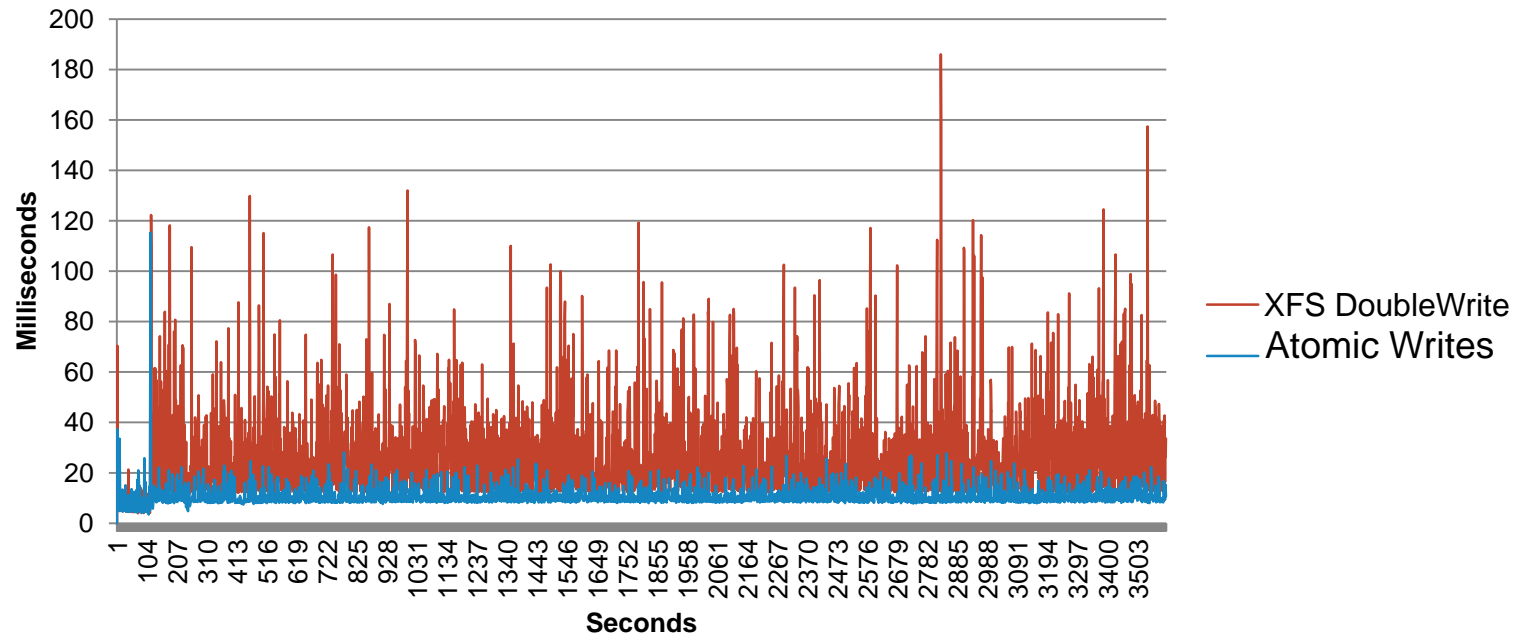
## MySQL with Atomic Writes

Database Server

| Page A | Page B | Page C |

DRAM Buffer | Page A | Page B | Page C |

Page A
Page C
Page B

ioMemory    Database

1 | Application initiates updates to pages A, B, and C.

2 | MySQL copies updated pages to memory buffer.

3 | MySQL writes to actual tablespace, bypassing the double-write buffer step due to inherent atomicity guaranteed by the (intelligent) device.

# MySQL Example: Latency Improvement



**2-4x Latency Improvement on Percona Server**
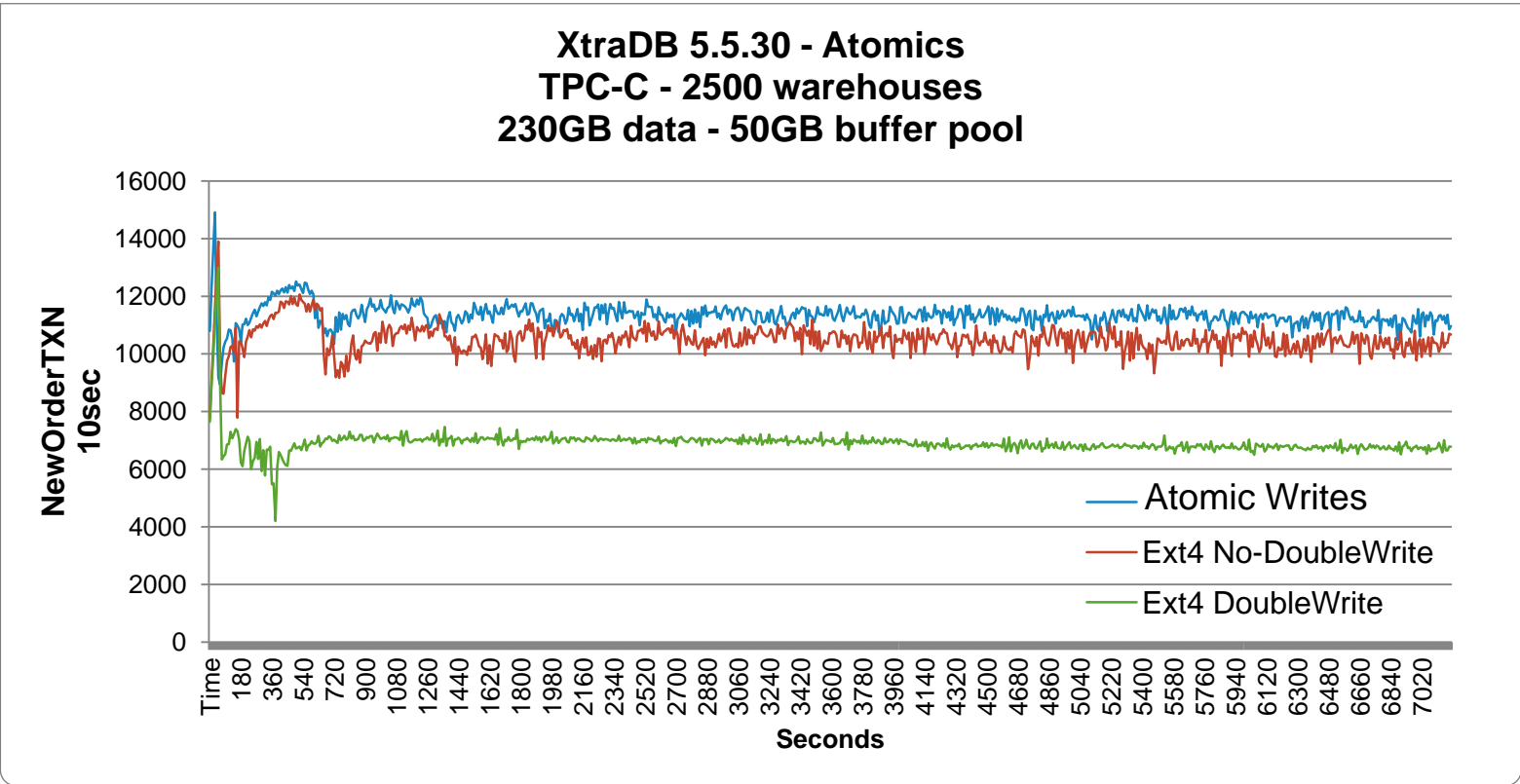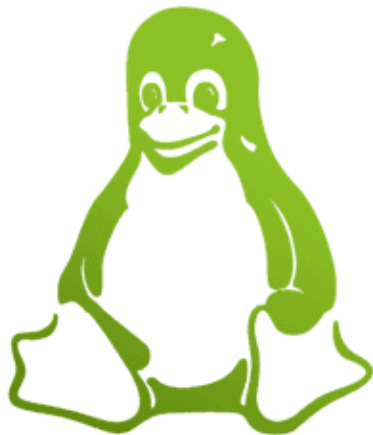
Sysbench 99% Latency
OLTP workload

Legend: XFS DoubleWrite, Atomic Writes

Y-axis: Milliseconds
X-axis: Seconds

# MySQL Example: Throughput Improvement

**FUSION-IO**

## 70% Transactions/sec Improvement on MariaDB Server

**XtraDB 5.5.30 - Atomics**
**TPC-C - 2500 warehouses**
**230GB data - 50GB buffer pool**



Legend:
- Atomic Writes
- Ext4 No-DoubleWrite
- Ext4 DoubleWrite

Y-axis: NewOrderTXN 10sec (0 – 16000)

X-axis: Seconds (Time, 180, 360, 540, 720, 900, 1080, 1260, 1440, 1620, 1800, 1980, 2160, 2340, 2520, 2700, 2880, 3060, 3240, 3420, 3600, 3780, 3960, 4140, 4320, 4500, 4680, 4860, 5040, 5220, 5400, 5580, 5760, 5940, 6120, 6300, 6480, 6660, 6840, 7020)

# 2nd Contribution: Linux Fast-Swap

**Flash-aware Linux swap**

When working set size exceeds the capacity of DRAM, demand page from a flash-aware virtual memory subsystem.
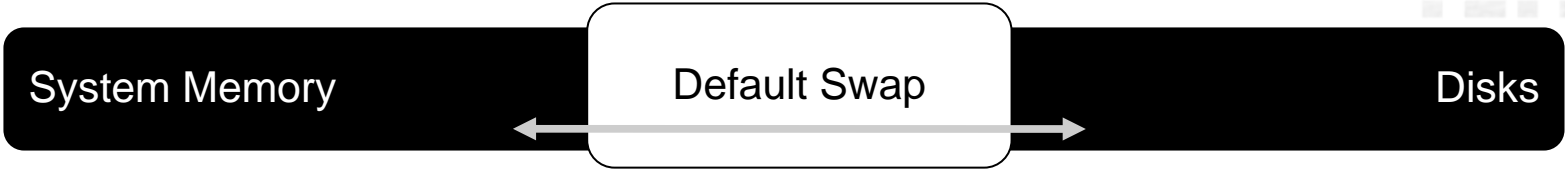
Repository | Learn More

On GitHub:

- Documentation

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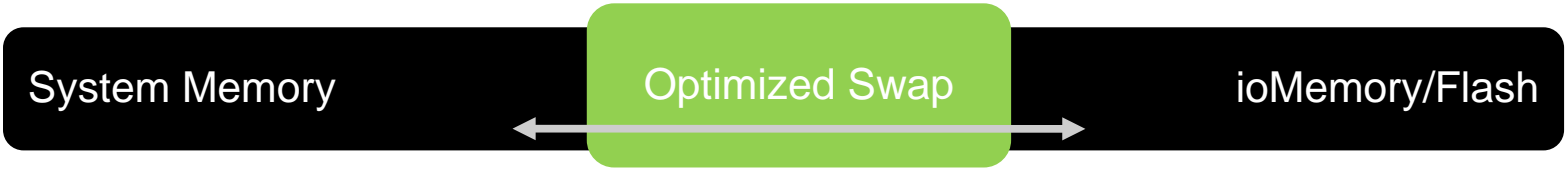
- Benchmarking utility

https://opennvm.github.io

# Improving Linux Swap (Demand-paging)

**FUSiON-iO**

| System Memory | Default Swap | Disks |
|---|---|---|

Originally designed as a last resort to prevent OOM (out-of-memory) failures
- Never tuned for high-performance demand-paging
- Never tuned for multi-threaded apps
- Poor performance

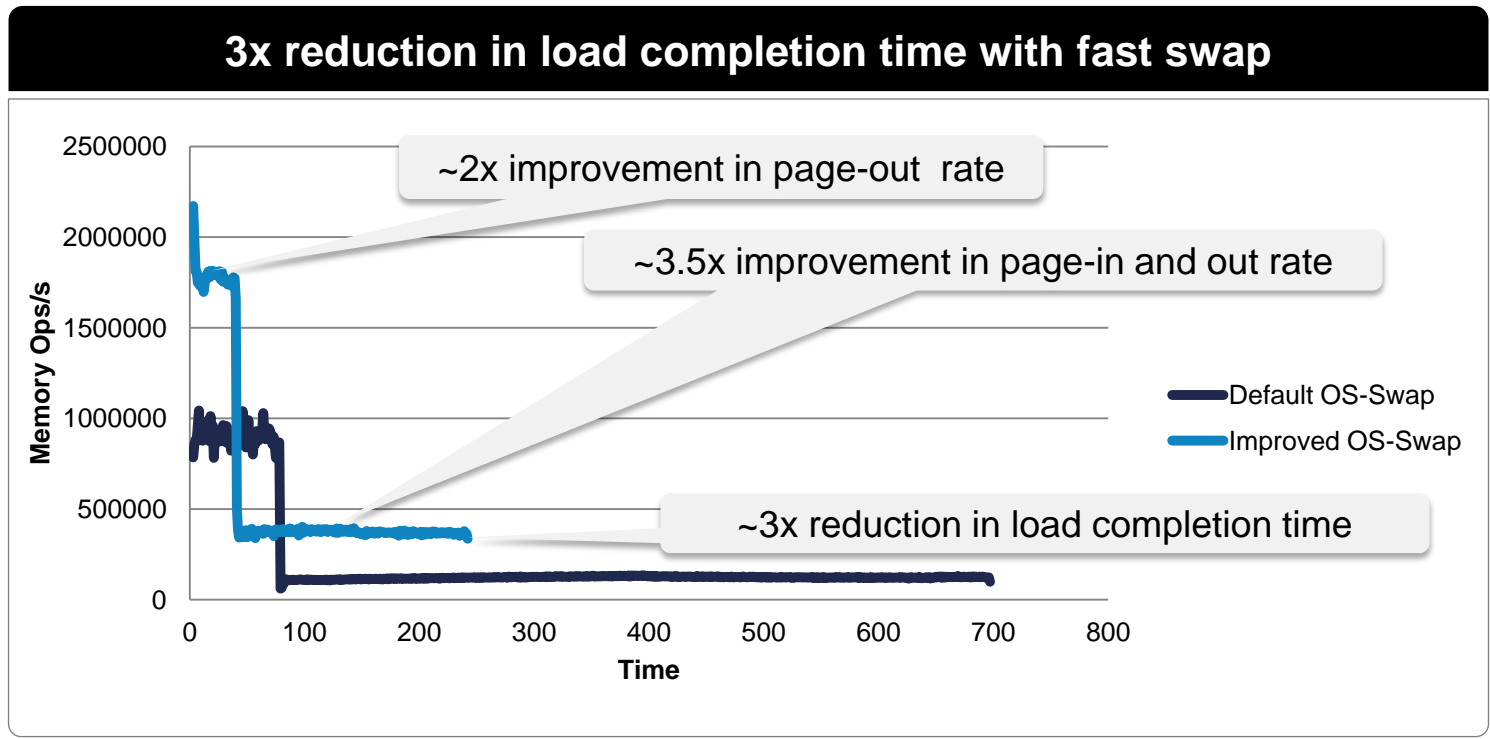| System Memory | Optimized Swap | ioMemory/Flash |
|---|---|---|

Tuned for flash (leverages native characteristics)
- O(1) algorithm for swap_out – reduce algorithm time and leverage fast random I/O
- Per CPU reclaim – greater throughput for multi-threaded environments
- Intelligent read-ahead on swap-in – cut legacy, disk-era cruft for rotational latency

# Fast Swap - Performance

**3x reduction in load completion time with fast swap**

Chart: Memory Ops/s vs Time

- ~2x improvement in page-out rate
- ~3.5x improvement in page-in and out rate
- ~3x reduction in load completion time

Legend:
- Default OS-Swap
- Improved OS-Swap

# 3rd Contribution: Key-Value Interface



**Key-value interface to flash**

Create NoSQL databases faster. Automate garbage collection of expired data.

Repository | Learn More

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 (30 Aug)

https://opennvm.github.io

# 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.
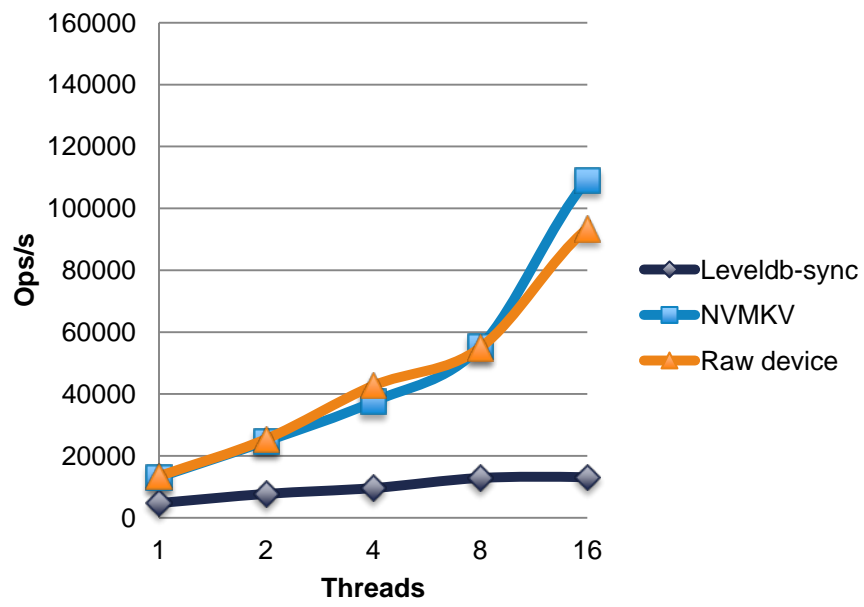
- **Near performance of raw device**
  Smarter media now natively understands a key-value I/O interface with lock-free updates, crash recovery, and no additional metadata overhead.

- **3x throughput on same SSD**
  Early benchmarks comparing against synchronous levelDB show over 3x improvement.

- **Up to 3x capacity increase**
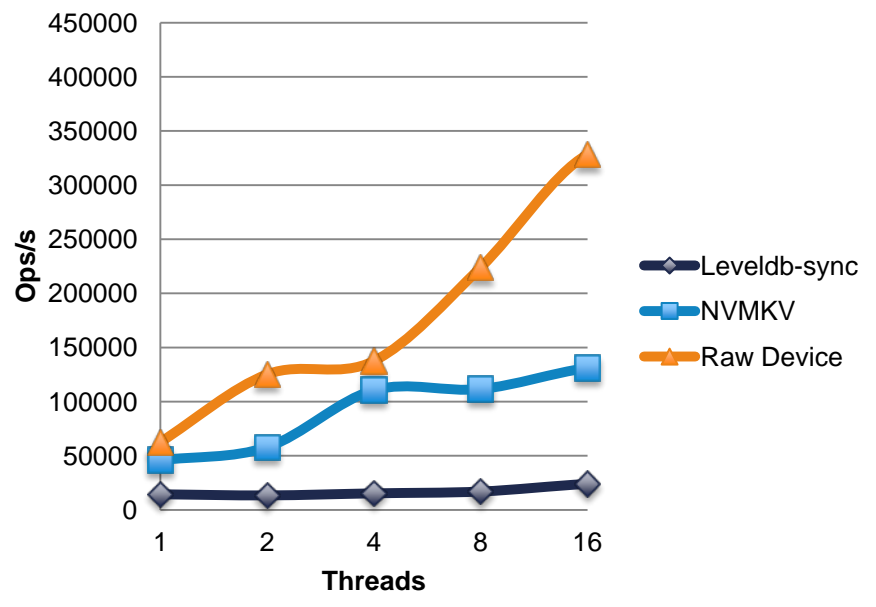  Dramatically reduces over-provisioning through coordinated garbage collection and automated key expiry.

# Key-Value Interface - Performance

FUSiON-iO

## Key-Value get/put vs. Raw read/write vs. levelDB read/write



GET v READ — Ops/s vs Threads (1, 2, 4, 8, 16). Legend: Leveldb-sync, NVMKV, Raw device.

PUT v WRITE — Ops/s vs Threads (1, 2, 4, 8, 16). Legend: Leveldb-sync, NVMKV, Raw Device.

# **OpenNVM, Standards, and Consortiums**

- opennvm.github.io

  - ▸ Primitives API specifications, sample code

  - ▸ Linux swap kernel patch and benchmarking tools

  - ▸ key-value interface API library, sample code, benchmark tools

- INCITS SCSI (T10) active standards proposals:

  - ▸ SBC-4 SPC-5 Atomic-Write
    http://www.t10.org/cgi-bin/ac.pl?t=d&f=11-229r6.pdf

  - ▸ SBC-4 SPC-5 Scattered writes, optionally atomic
    http://www.t10.org/cgi-bin/ac.pl?t=d&f=12-086r3.pdf

  - ▸ SBC-4 SPC-5 Gathered reads, optionally atomic
    http://www.t10.org/cgi-bin/ac.pl?t=d&f=12-087r3.pdf

- SNIA NVM-Programming TWG draft guide:
  http://snia.org/forums/sssi/nvmp

# Apps Using OpenNVM technology

MariaDB

Learn More »
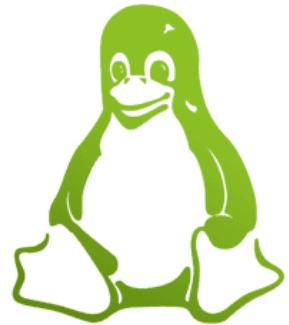
PERCONA
SERVER

Learn More »

https://opennvm.github.io

# Join us at opennvm.github.io



## Current OpenNVM Repositories

**Flash-aware Linux swap**

When working set size exceeds the capacity of DRAM, demand page from a flash-aware virtual memory subsystem.

[ Repository ] [ Learn More ]

**Key-value interface to flash**

Create NoSQL databases faster. Automate garbage collection of expired data.

[ Repository ] [ Learn More ]

**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 ]

# THANK YOU

fusionio.com  |  REDEFINE WHAT'S POSSIBLE