

NVMFS Benchmark

Introduction

Makes MySQL flash-aware and improves latency performance and consistency while increasing storage efficiency of Fusion ioMemory PCIe application accelerators

Replaces traditional double-write operation with single atomic write for lower, more consistent latency and reduced flash wear

New NVM Compression algorithm delivers benefits of compression without performance penalty

Tools: TPC-C, Warehouse numbers: 1000, storage: FusionIO SX300 1.6TB, OS: CentOS
6.5 kernel: 2.6.32-431.el6.x86_64
InnoDB buffer pool: 64G
sync_binlog = 0
innodb_flush_log_at_trx_commit=2
10.0.15-MariaDB-log MariaDB Server(supports atomic writes)
File system: NVMFS ,ext4
warmup time: 120s
runtime: 3600s
Threads: 32 ~ 512
Machine: Dell R720 24core Intel(R) Xeon(R) CPU E5-2630 v2 @ 2.60GHz

How to install NVMFS

Step	Action
1	Login as root.
2	Disable <code>auto_attach</code> by issuing the following command: <pre># modprobe iomemory-vsl4 auto_attach=0</pre>
3	Detach the device if it is already attached: <pre># fio-detach /dev/fctx</pre> (where <code>fctx</code> is <code>fct1</code> or <code>fct2</code> , etc.)
4	Format the device with Fusion ioMemory SDK options enabled use -A for atomic writes -P for Persistent TRIM -e for sparse address <pre>fio-format -APye -b 512 /dev/fctx</pre>
5	Attach the device for I/O operations. <pre># fio-attach /dev/fctx</pre> This should result in device (<code>fiox</code>) appearing under <code>/dev</code> .

Step	Action
1	Load the kernel module nvme nvme_core using the modprobe command as follows: <ul style="list-style-type: none">• modprobe nvme nvme_core
2	Verify that the kernel module has loaded properly using the lsmod command as follows: <ul style="list-style-type: none">• lsmod grep nvme

#mount

/dev/sda2 on / type ext4 (rw)

proc on /proc type proc (rw)

sysfs on /sys type sysfs (rw)

devpts on /dev/pts type devpts (rw,gid=5,mode=620)

tmpfs on /dev/shm type tmpfs (rw)

/dev/sda1 on /boot type ext4 (rw)

/dev/sda4 on /home type ext4 (rw)

none on /proc/sys/fs/binfmt_misc type binfmt_misc (rw)

/dev/fioa on /storage/sas type NVMe (rw,noatime)

How to configure MariaDB

Partial write operations

When InnoDB writes to the filesystem, there is generally no guarantee that a given write operation will be complete (not partial) in cases of a poweroff event, or if the operating system crashes at the exact moment a write is being done.

Without detection or prevention of partial writes, the integrity of the database can be compromised after recovery.

innodb_doublewrite - an imperfect solution

Since its inception, InnoDB has had a mechanism to detect and ignore partial writes via the [InnoDB Doublewrite Buffer](#) (also `innodb_checksum` can be used to detect a partial write).

Doublewrites, controlled by the `innodb_doublewrite` system variable, comes with its own set of problems. Especially on SSD, writing each page twice can have detrimental effects (write leveling).

Atomic write - a faster alternative to innodb_doublewrite

A better solution is to directly ask the filesystem to provide an atomic (all or nothing) write guarantee. Currently this is only available on the NVMFS (previously called directFS) filesystem on FusionIO devices that provide atomic write functionality. This functionality is supported by MariaDB's XtraDB and InnoDB storage engines with [MDEV-4338](#).

Enabling Atomic Writes

To use atomic writes instead of the doublewrite buffer, add:

```
innodb_use_atomic_writes = 1
```

to the my.cnf config file.

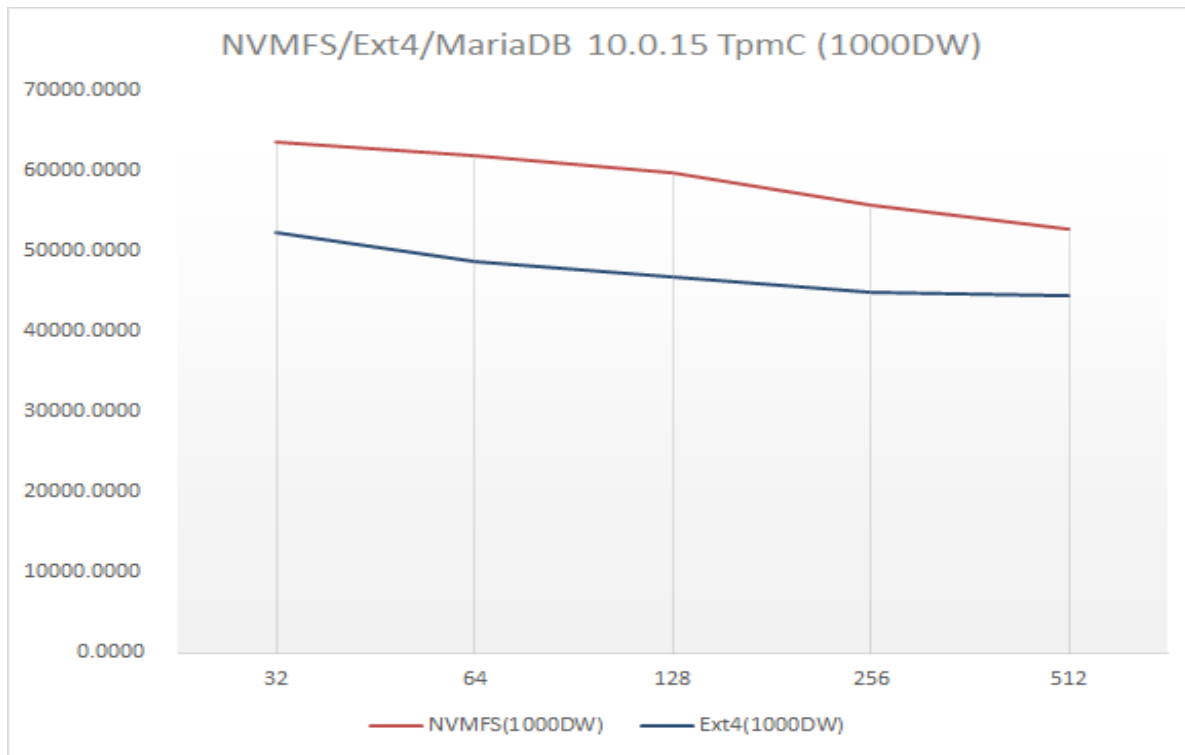
About innodb_use_atomic_writes

The following happens when `innodb_use_atomic_writes` is switched ON

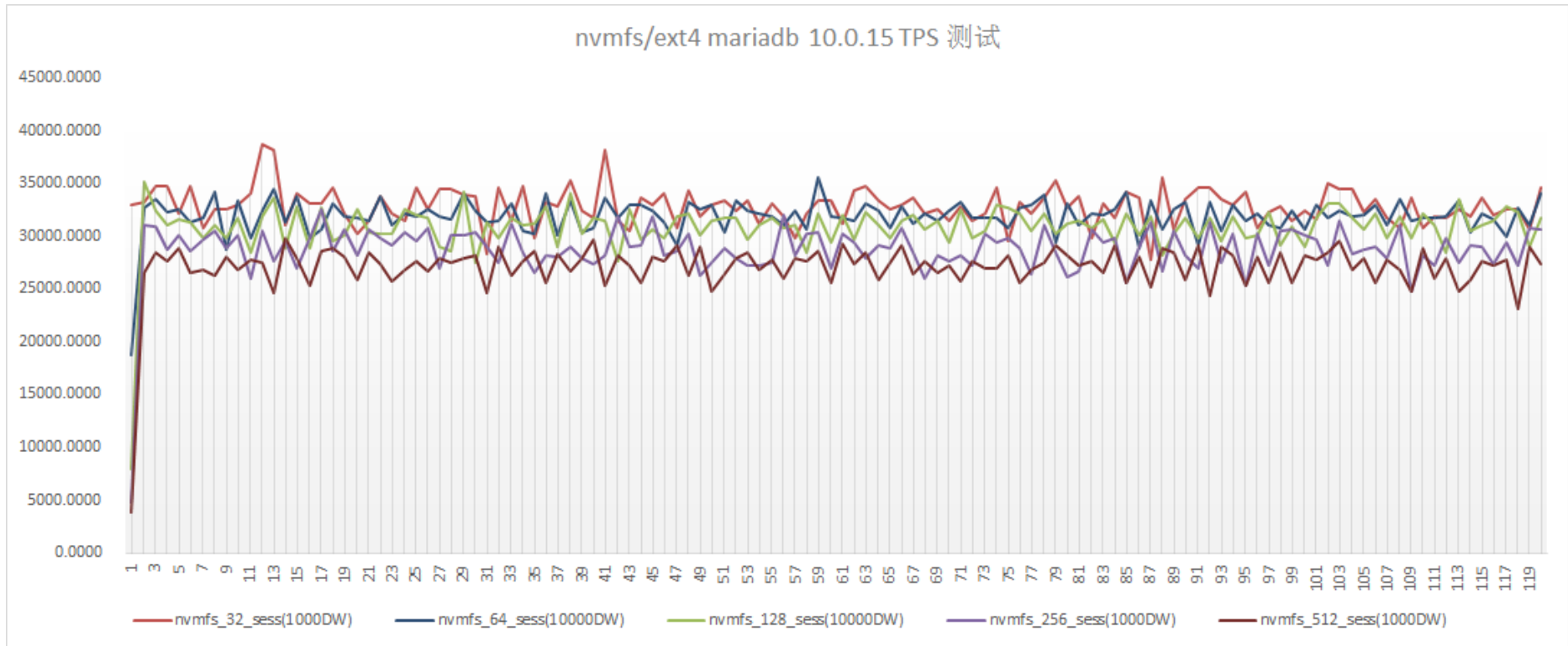
- if [innodb_flush_method](#) is neither `O_DIRECT`, `ALL_O_DIRECT`, or `O_DIRECT_NO_FSYNC`, it is switched to `O_DIRECT`
- [innodb_use_fallocate](#) is switched `ON` (files are extended using `posix_fallocate` rather than writing zeros behind the end of file)
- Whenever an InnoDB datafile is opened, a special `ioctl()` is issued to switch on atomic writes. If the call fails, an error is logged and returned to the caller. This means that if the system tablespace is not located on an atomic write capable device or filesystem, InnoDB/XtraDB will refuse to start.
- if [innodb_doublewrite](#) is set to `ON`, `innodb_doublewrite` will be switched `OFF` and a message written to the error log.

Benchmark Result

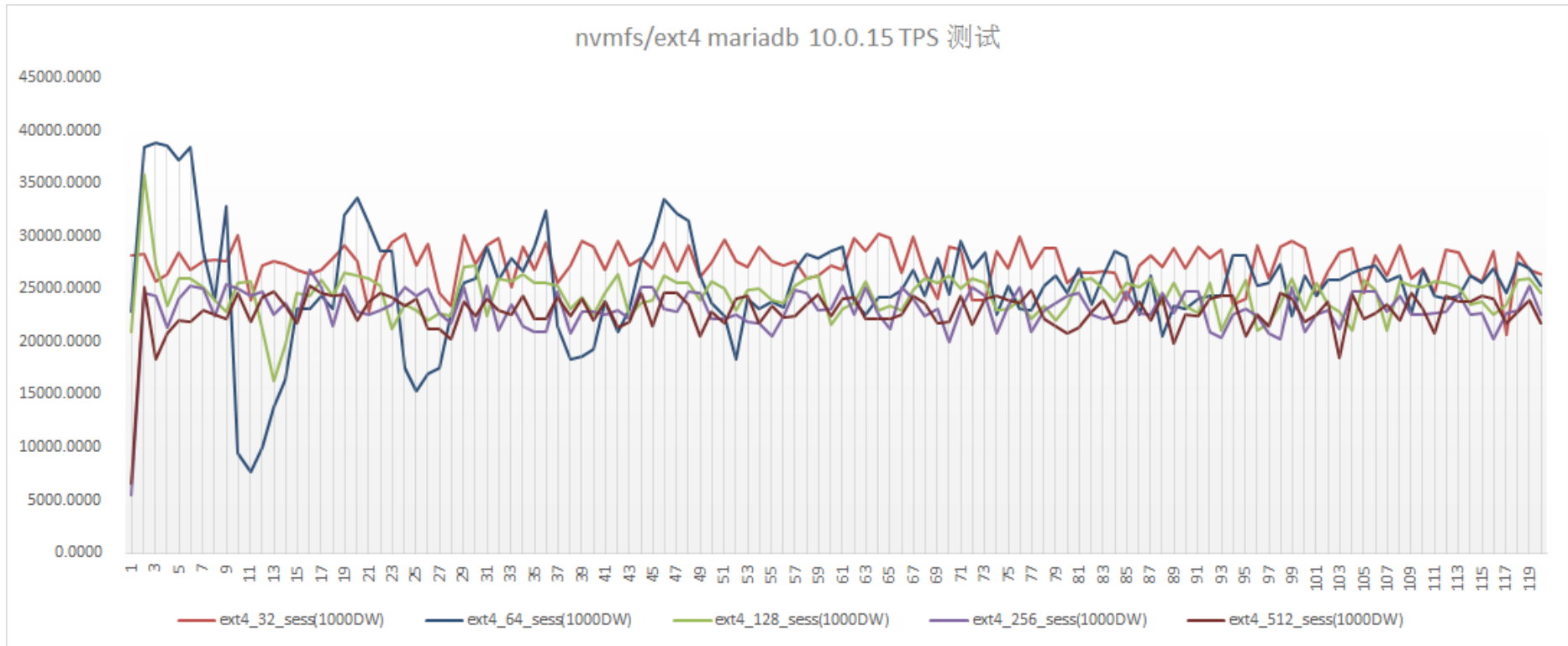
TpmC of all threads



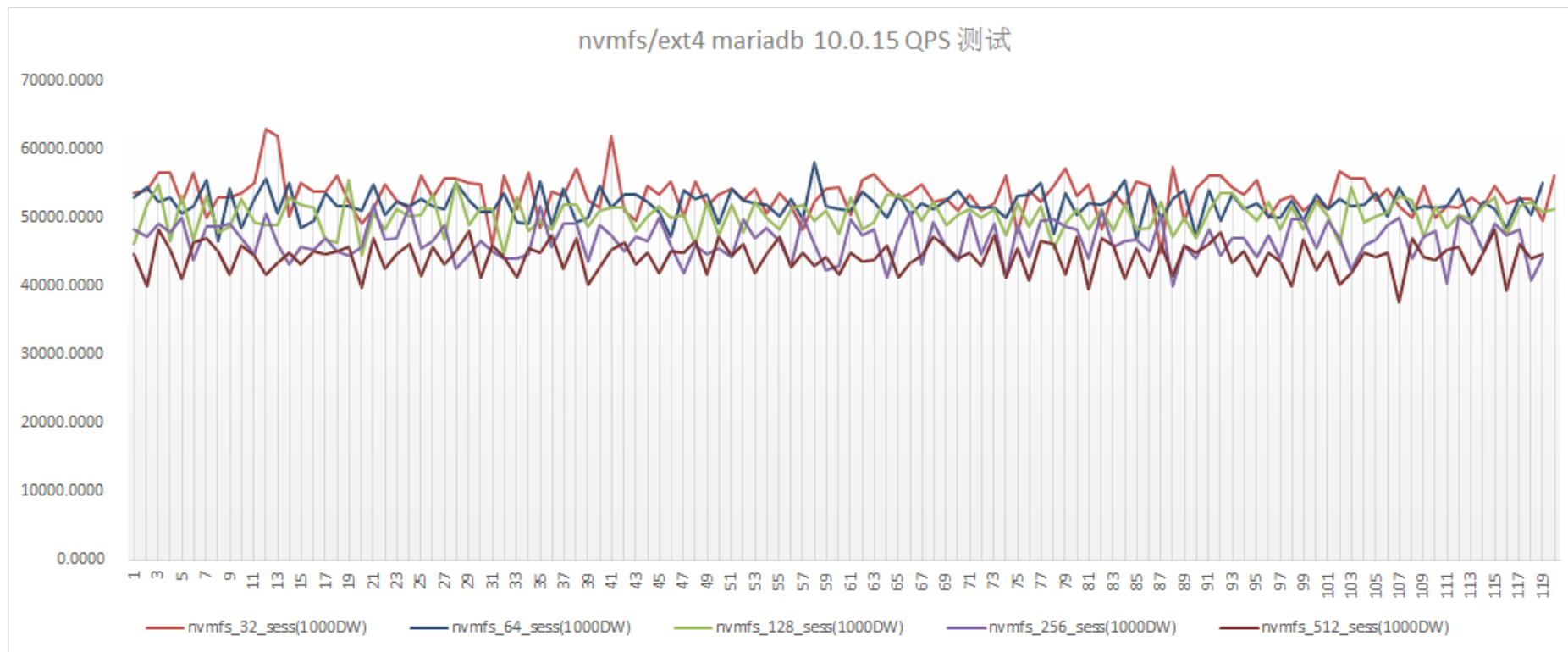
TPS of NVMFS



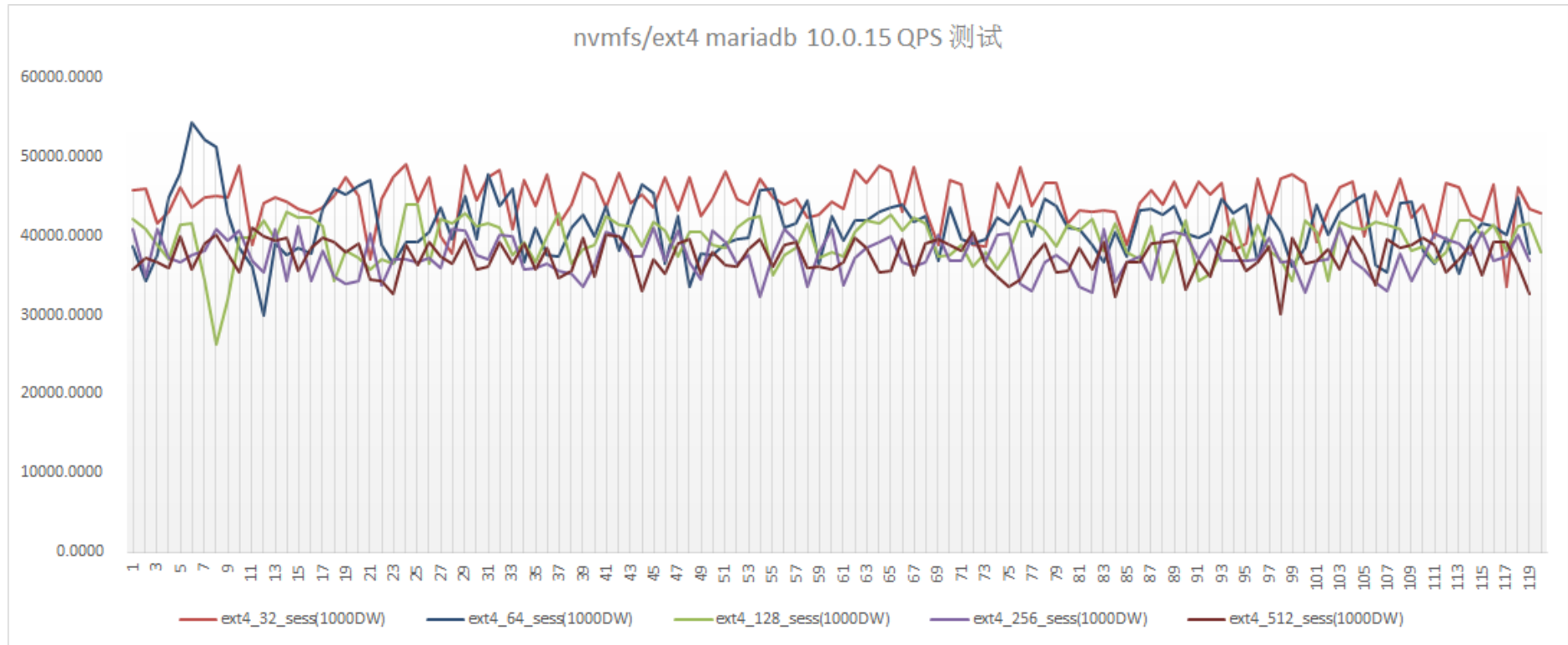
TPS of Ext4



QPS of NVMFS



QPS of Ext4



Benchmark by SanDisk

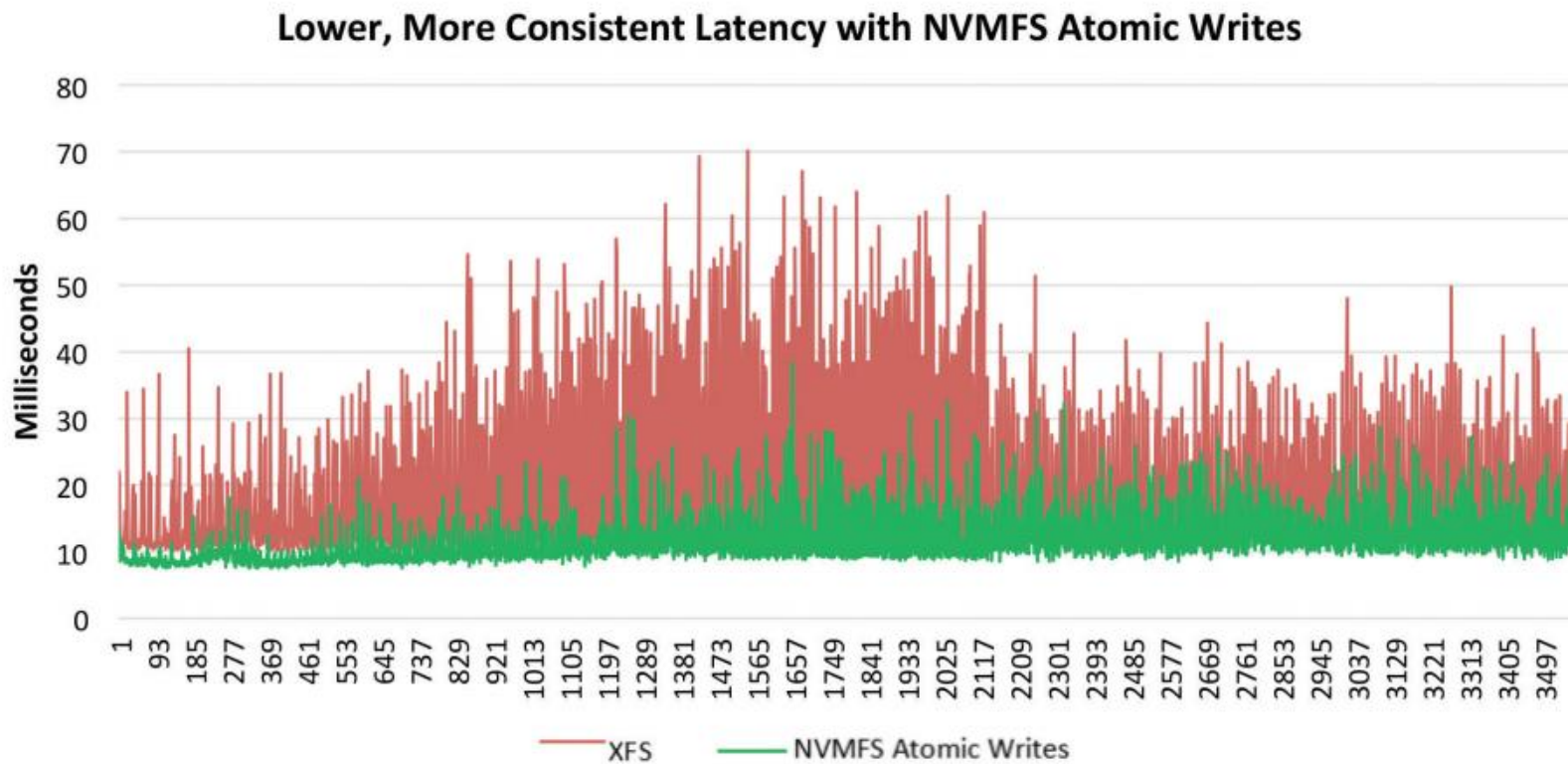


Figure 1: XFS vs. NVMFS Atomic Writes (Sysbench - MariaDB 10.0.15, 4000 OLTP TXN injection/second, 99% latency, 220 GB data - 10 GB buffer pool)*

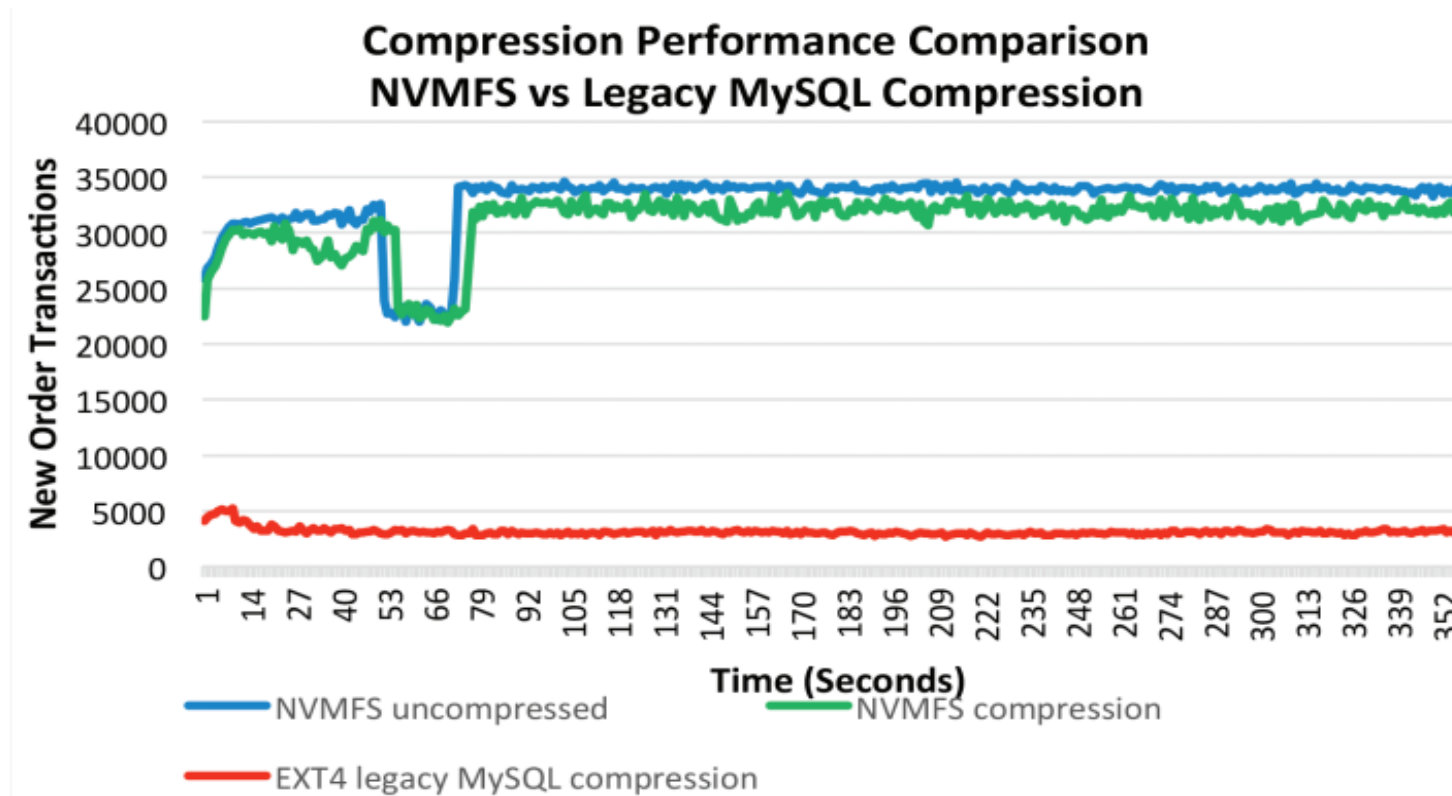


Figure 2: NVMFS compression on MariaDB workload (TPC-C-like benchmark: 1000 warehouses - 75 GB MySQL Buffer pool, MariaDB 10.0.15)*

In My test, MySQL has a better performance on ext4 filesystem which reaches about 5000 TpmC. I wonder if SanDisk really test ext4 filesystem on Flash storage

Details of this Benchmark

Totally NVMFS has a great performance improvement – TPS +20% QPS +25% (and SanDisk's Test shows NVMFS also has a better performance than XFS)

http://www.vmcd.org/docs/NVMFS_INTRO.pdf

http://www.vmcd.org/docs/NVMFS_User_Guide.pdf

<http://www.vmcd.org/docs/ext4.log>

<http://www.vmcd.org/docs/nvmfs.log>