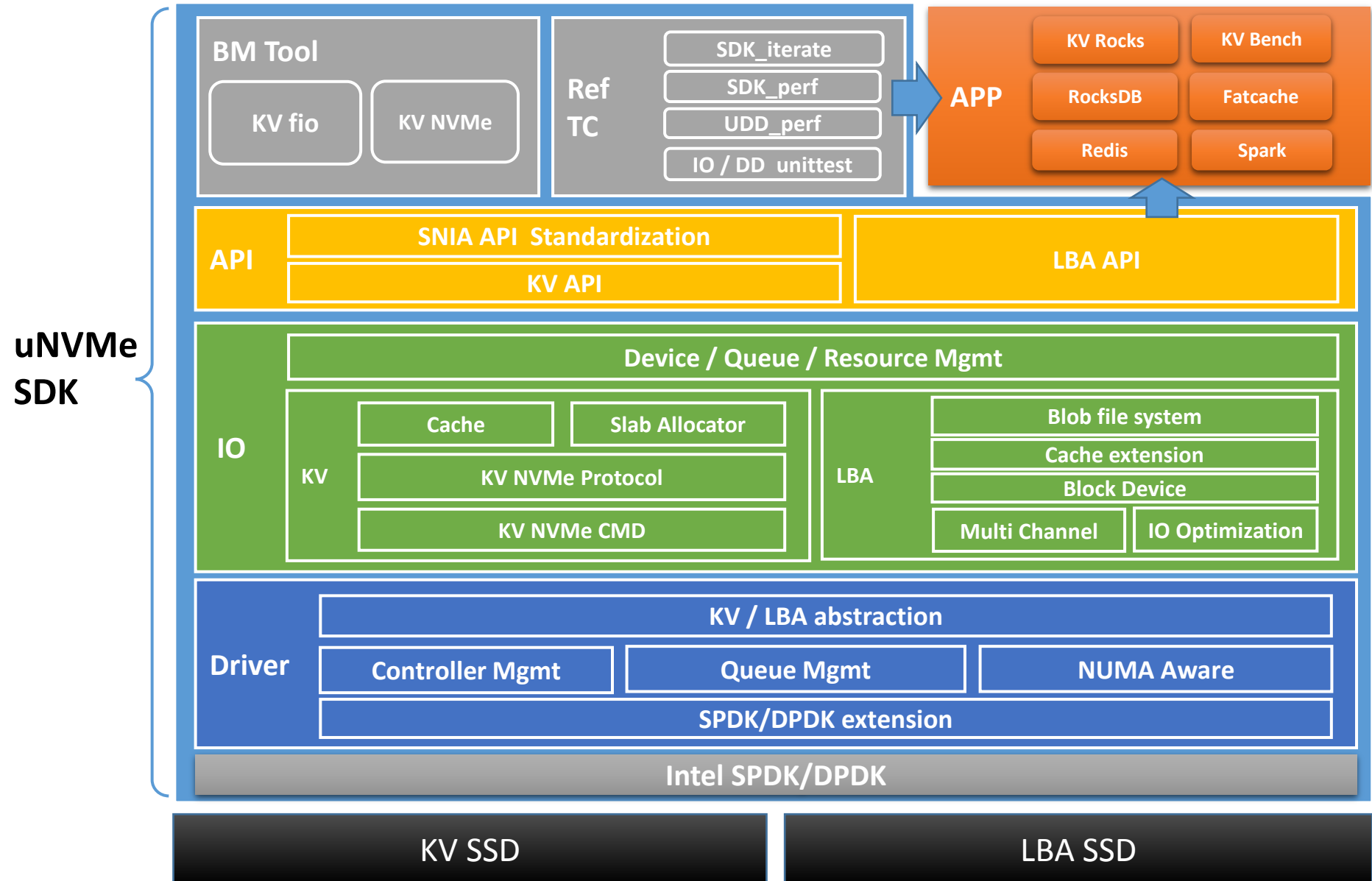


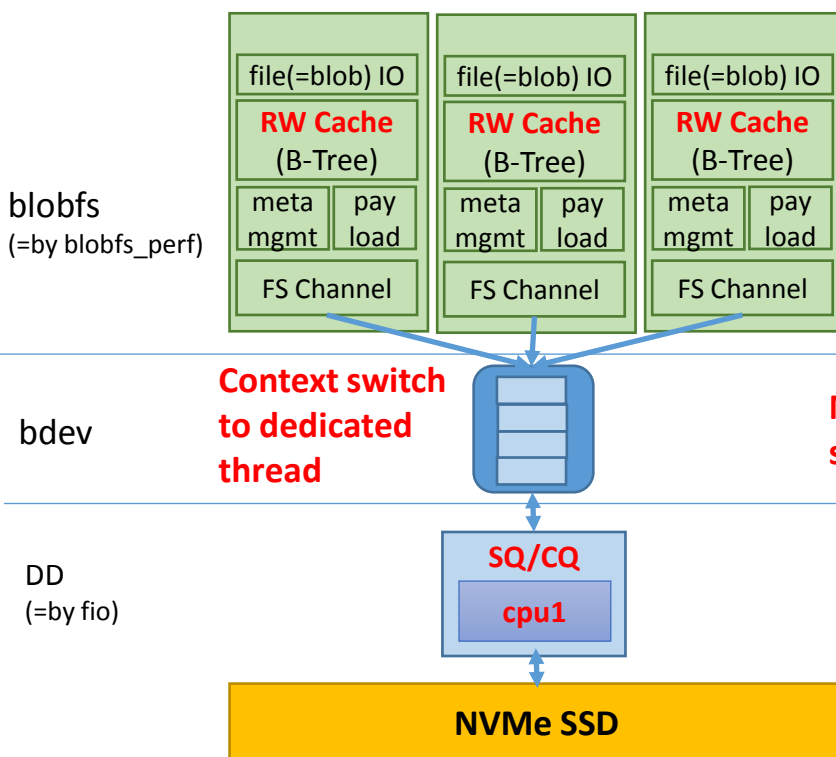
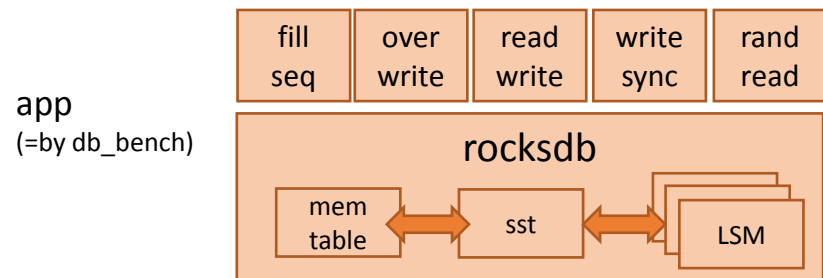
uNVMe-Blobfs

(OpenMPDK-uNVMe)

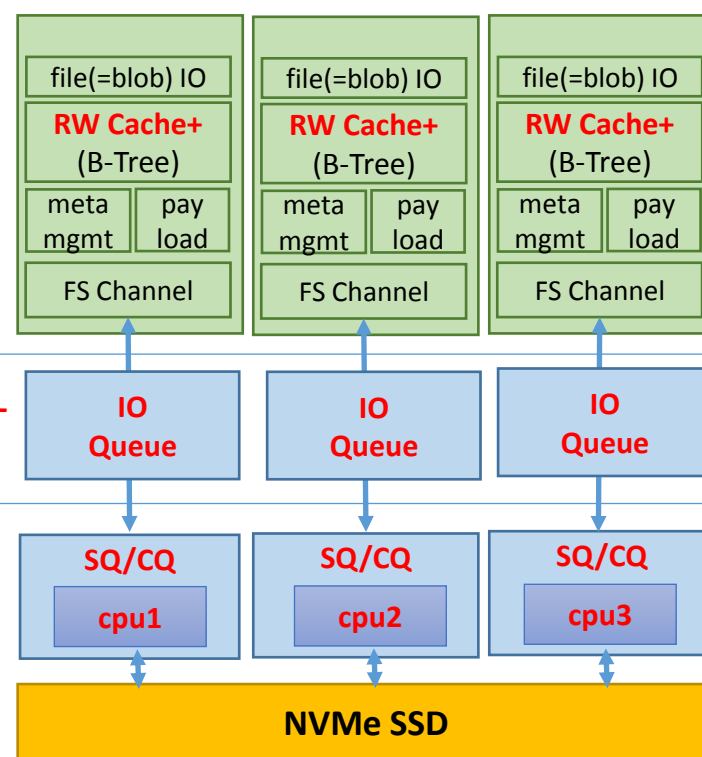
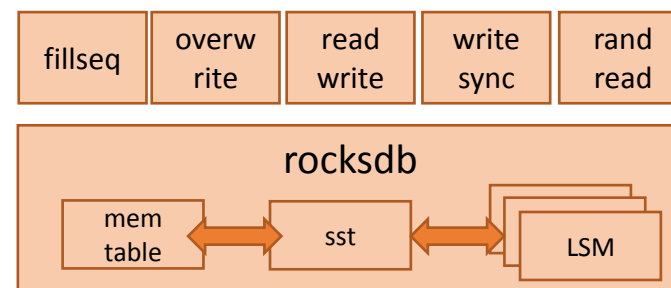
Kyungsan Kim / SW Dev Team / Memory Business Unit



Vanilla SPDK



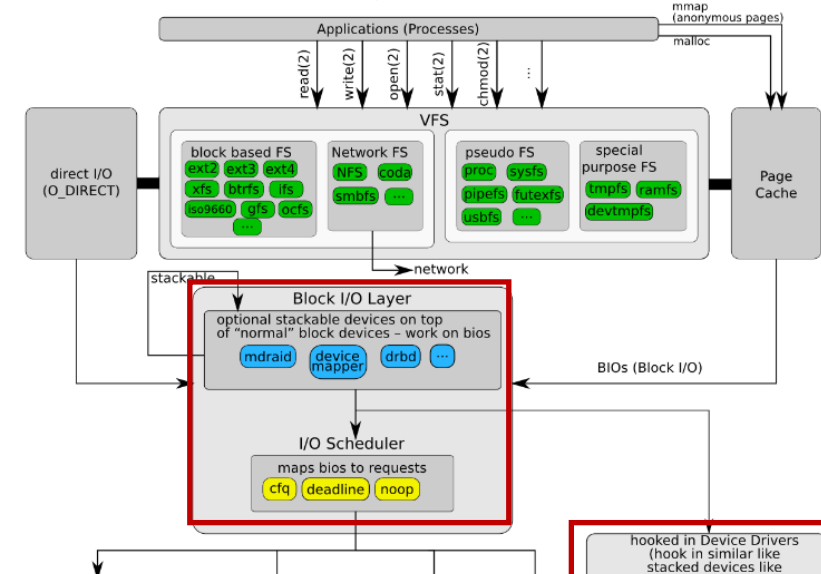
uNVMe-MQ



- Around 2013 Oct, Linux kernel introduced a big change in storage IO stack, **block-mq**, that reflects prevalent of multi-core system with high speed storage

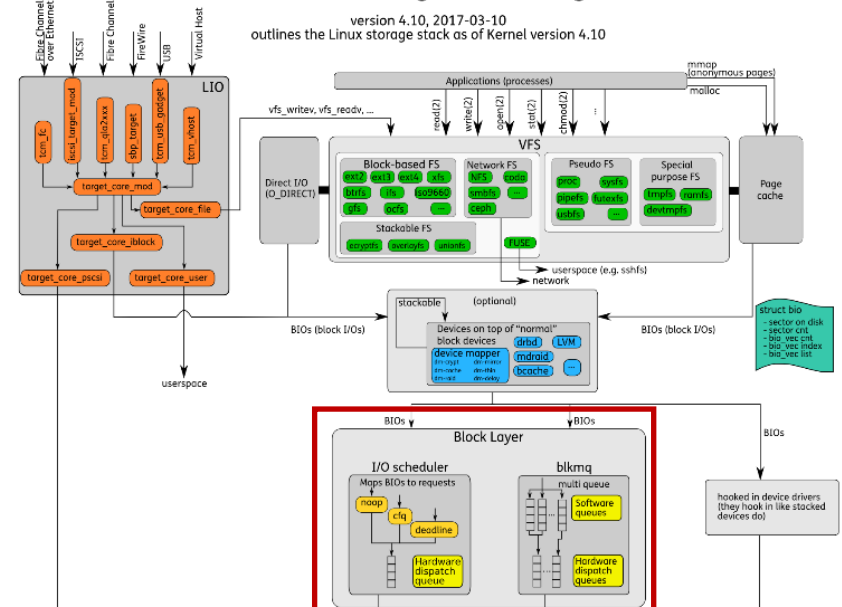
The Linux I/O Stack Diagram

version 1.0, 2012-06-20
outlines the Linux I/O stack as of Kernel version 3.3



The Linux Storage Stack Diagram

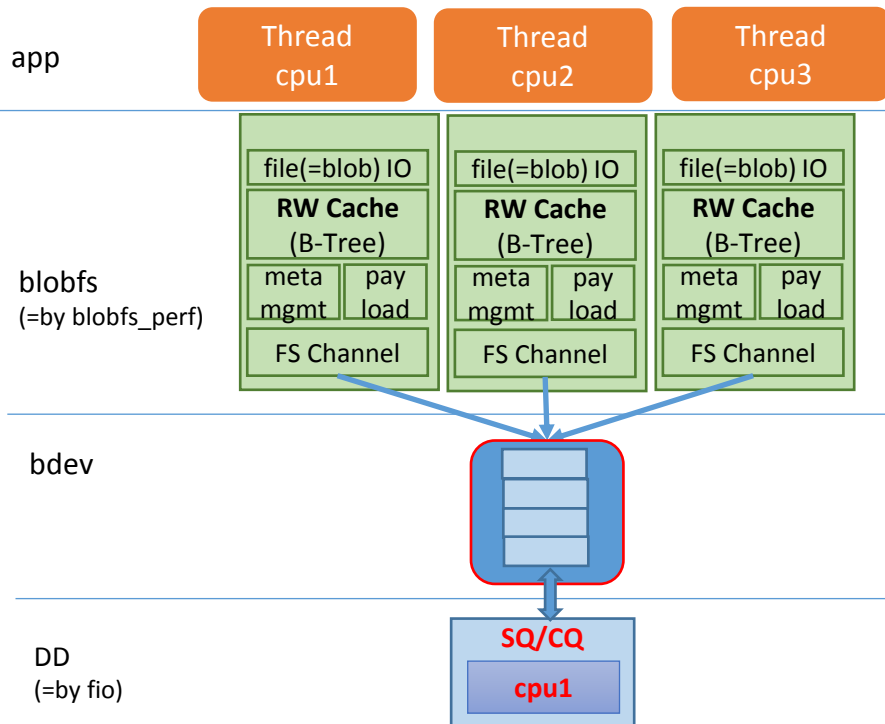
version 4.10, 2017-03-10
outlines the Linux storage stack as of Kernel version 4.10



Driver	Device Name	Supported Devices	blk-mq Since Kernel Version
null_blk	/dev/nullb* ^[8]	none (test drivers)	3.13 (git commit [8])
virtio-blk	/dev/vd*	Virtual guest drivers (e.g. under KVM ^{[7][8]})	3.13 (git commit [8])
mtip32xx	/dev/rssd*	Micron RealSSD PCIe	3.16 (git commit [9])
scsi (scsi_mq)	/dev/sd*	e.g. SAS and SATA SSDs/HDDs	3.17 (git commit [9])
NVMe	/dev/nvme*	e.g. Intel SSD DC P3600 DC P3700 Series ^[9]	3.19 (git commit [9])
rbd	/dev/rdb*	RADOS Block Device (Ceph)	4.0 (git commit [10])
ubi/block	/dev/ubiblock*		4.0 (git commit [10])
loop	/dev/loop*	Loopback-Device	4.0 (git commit [10])
dm / dm-mpath		request-based device mapper targets (derzeit ist dies ausschließlich dm-multipath)	planned for 4.1 ^[10]

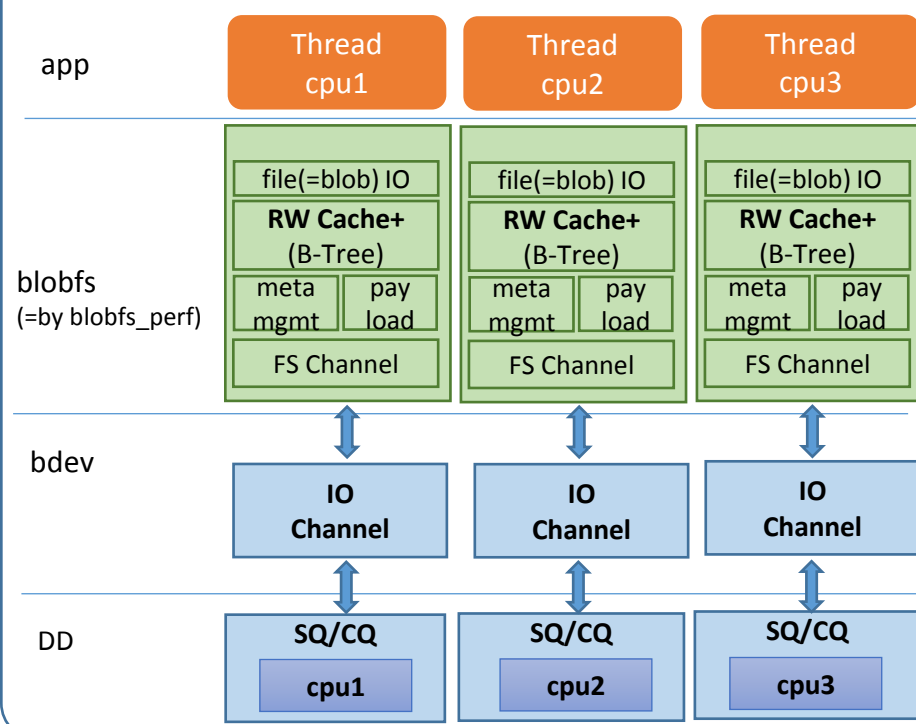
• SPDK Limitation

- blobfs and bdev layer has 3 limitations in that it is still designed to utilize **single block layer queue on a dedicated CPU**. Therefore,
 - ✓ 1) Request queue contention : on every single IO request for queue manipulation.
 - ✓ 2) Scalability : doesn't make use of multi core and multi-queue SSD capability, unlikely current kernel IO stack.
 - ✓ 3) Remote memory access : single cpu on bdev layer forces remote memory access across CPU sockets. Hence, Blobfs is not NUMA-aware.



• uNVMe MQ

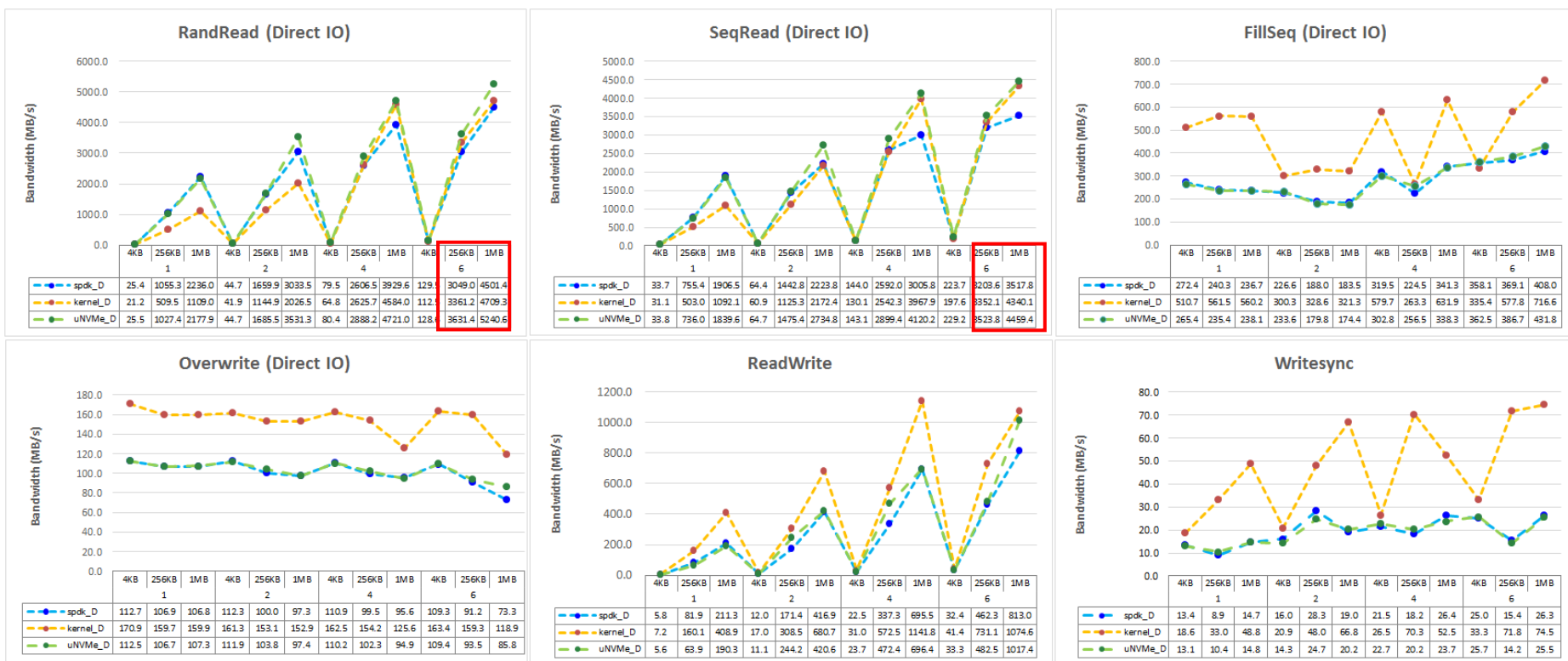
- ✓ Support **per CPU queue** to avoid the architectural inefficiencies
- ✓ On this release, we have applied MQ only on **READ** requests



■ uNVMe-MQ

- developing multi queue block layer over SPDK blobfs
- Around 30% **gain on Seq/Rand Read workload** at Blob filesystem layer (*RocksDB db_bench*)
- Applying uNVMe-MQ on rocksdb workload on many testbed

<Rocksdb Performance comparison - SPDK, kernel and uNVMe MQ Direct IO>

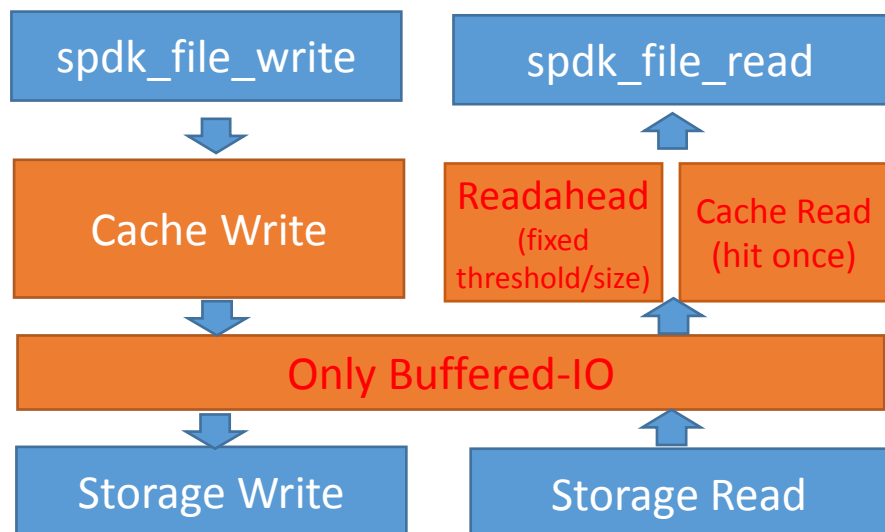


System: i7-7700k 4core 8threads @4.50GHz, 16GB DRAM, Ubuntu 16.04.3(kernel 4.9.82), PM1725a 1.6TB
 RocksDB: v5.6.1, 16B key, 3K/255KB/1023KB value, disable block cache, repeat 5 time on readseq, fillseq from N threads
 SPDK: v18.04.1

- Please refer Evaluation Guide for Test configuration (thread = 1/2/4/6 and so on)

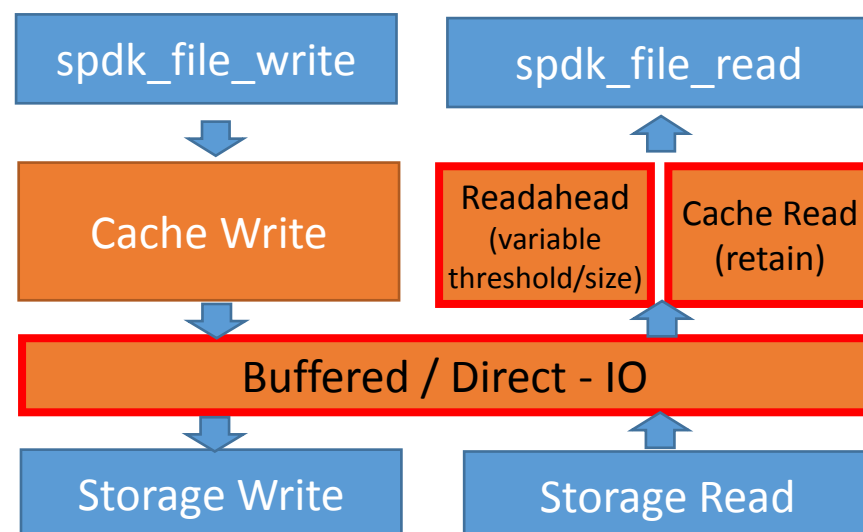
- **SPDK Limitation**

- ✓ **Readahead**: only when the amount of read requests beyond **consecutive 128KB**, it is judged as sequential read and **prefetch 512KB** in size constantly, otherwise prefetch is never triggered
- ✓ **Reclaim Policy**: **when a cached data hits** on a read request, **it is reclaimed instantly**, thus when the read data is revisited, it necessarily leads to disk read
- ✓ **Buffered IO only**: blobfs **always occupies GB size cache** memory



- **uNVMe-cache**

- ✓ **Readahead** : altered blobfs and provide an API to **configure best-fit prefetch threshold and size** depending on IO size out of an application
- ✓ **Reclaim Policy** : altered blobfs and provide an API to **determine reclaim policy**
- ✓ **Support Direct IO** : Altered blobfs and provide API to **support DIRECT IO** as well as Buffered IO as a selective way for cache supporting application

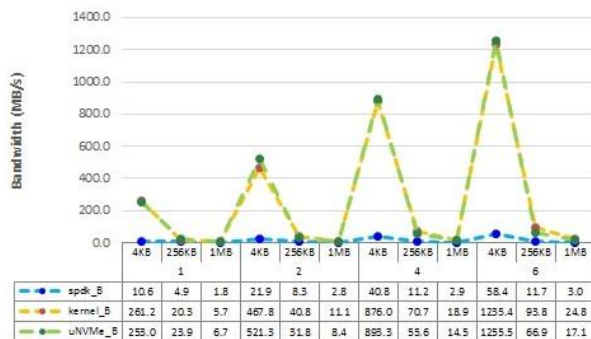


uNVM read cache

- uNVM-read cache shows performance gain against SPDK on db_bench workloads

<Rocksdb Performance comparison - SPDK, kernel and uNVM MQ Buffered IO>

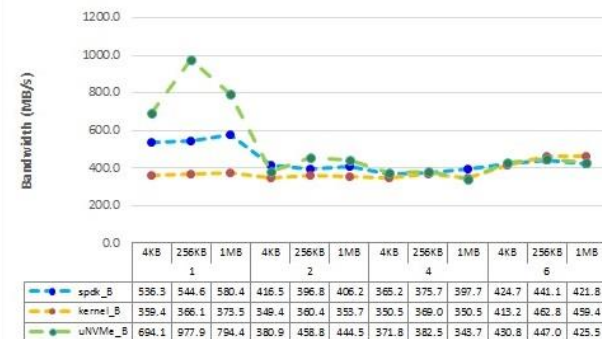
RandRead (Buffered IO)



SeqRead (Buffered IO)



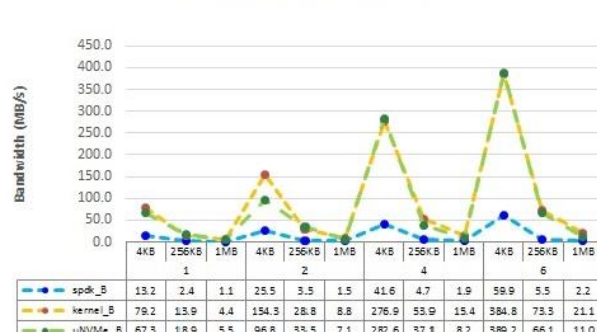
FillSeq (Buffered IO)



Overwrite (Buffered IO)



ReadWrite (Buffered IO)



Writesync (Buffered IO)



System : xeon E5-2667 16cores 32threads @3.60GHz , 256GB DRAM, Debian 9.1(kernel 4.9), PM983 4TB

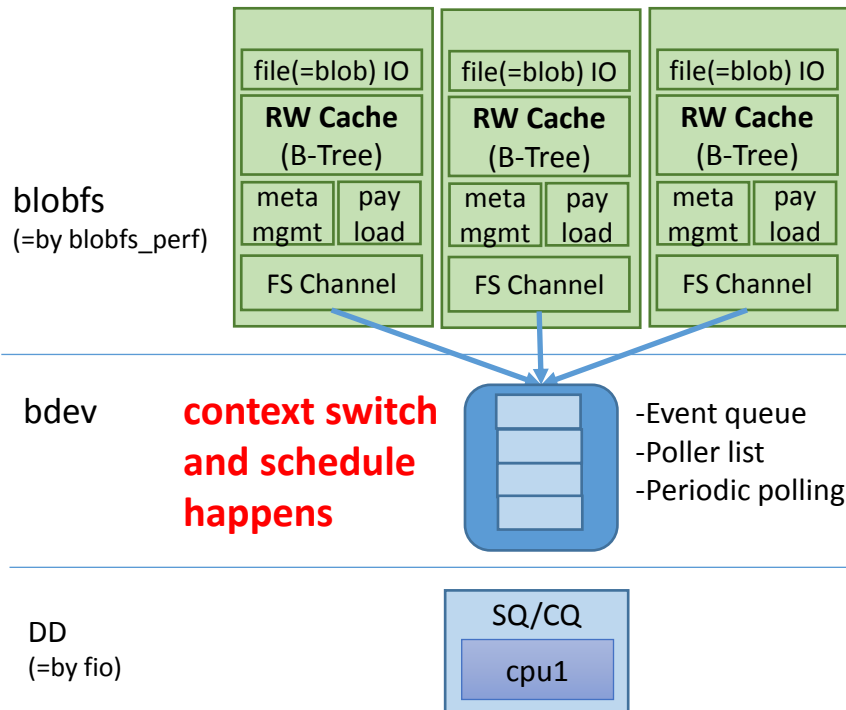
RocksDB: v5.6.1, 16B key, 1000B value, 1M keys for each thread, disable block cache, repeat 5 time on readseq, fillseq from N threads

SPDK: v18.04.1, blobfs cache 4GB uNVM: blobfs cache 4GB , RETAIN_CACHE=1, PREFETCH_CTL=1, 4KB prefetch_threshold for 4KB block_size, 128KB prefetch_threshold for 256KB/1MB block_size

- Please refer Evaluation Guide for Test configuration (thread = 1/2/4/6 and so on)

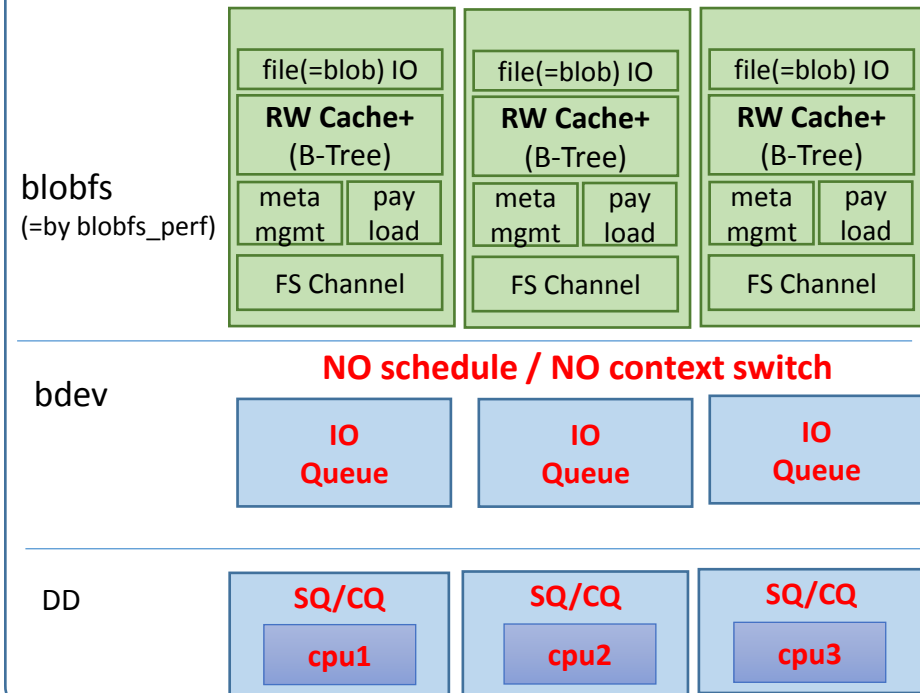
• SPDK Limitation

- ✓ Context-switch happens at blobfs to block device layer on every single blob IO request.
- ✓ Then, Block IO thread performs sort of scheduling : event queue, polling, interval polling . And big-loop of the IO thread consumes a certain amount of CPU cycle.



• uNVMe path

- ✓ Plan to provide IO path from top application top to bottom driver without IO schedule and context-switch



■ Reuse of uNVMe-Rocksdb Phase1

- db_bench
- Workload: fillseq / randread / seqread / overwrite / writesync / readwrite
- ✧ *For testing direct / buffered IO, please refer test configuration of Evaluation Guide*

■ Newly introduce blobfs layer BM Tool: Blobfs_perf

- Performs file IO on blobfs perf
- Seqwrite / seqread / randread/
- Buffered / Direct-IO
- Control read-cache reclaim
- Control prefetch size and prefetch threshold
- Multi-threads
- IO size and block size
- Working on vanilla SPDK / uNVMe

THE NEXT CREATION STARTS HERE

Placing **memory** at the forefront of future innovation and creative IT life

