

Accelerating PMEM access via SPDK acceleration framework

intel®

Ziye Yang

Staff Software Engineer
Intel

2021 Dec. 15th

Agenda

1

Background & Problem

2

SPDK acceleration framework introduction

3

New PMEM bdev (pmem_accel)
development status

4

Conclusion

SPDK, PMDK, Intel® Performance Analyzers

Virtual Forum

Background & Problem

Background & Problem

- With more and more deployment of fast NICs and storage devices in data center, there is more and more pressures on CPU, e.g., many users/customers find that CPU becomes **a bottleneck** to conduct operations between DRAM and PMEM (persistent memory, e.g., Intel® Optane™ Persistent Memory).
 - Data operations “between (persistent) Memory” is different from “data operations between Memory and device”. **For the latter case, there is (R)DMA capability on device, and CPU will not be very busy.**
- How to offload the work from CPU with some possible memory offloading engines (e.g., IOAT, DSA) is a real and valid requirement.
 - For example, cloud providers would like to utilize the important CPU resources in customer visible workloads (e.g., those host CPUs should be virtualized to launch more customers' Virtual Machines)

SPDK, PMDK, Intel® Performance Analyzers

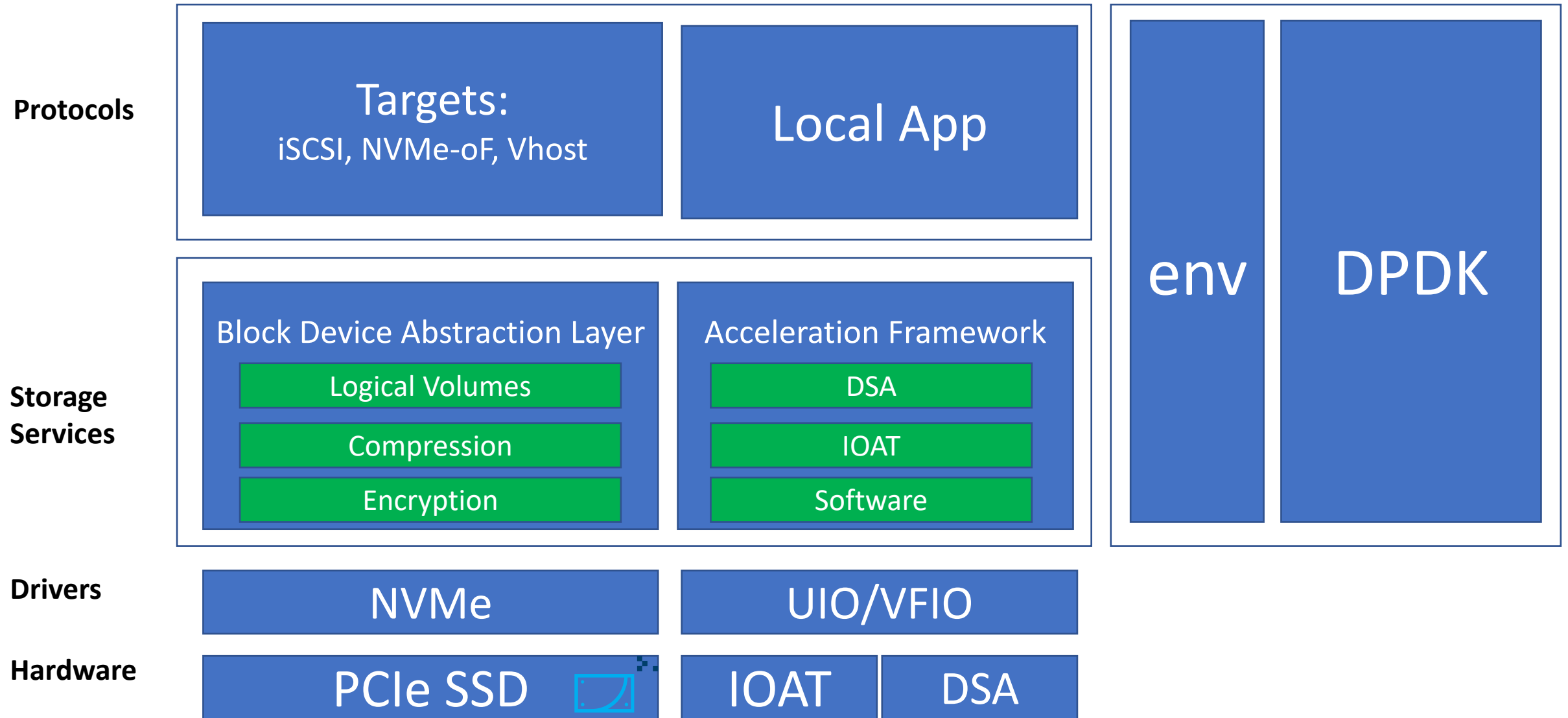
Virtual Forum

SPDK acceleration framework introduction

Some memory offloading devices

- Functionality: Support operations between (persistent) memory
- There are some memory offloading engines provided by Intel:
 - Intel IOAT: [I/O Acceleration Technology](#)
 - Intel DSA: [Intel® Data Streaming Accelerator \(Intel® DSA\)](#)
- Intel DSA is a high-performance data copy and transformation accelerator that will be integrated in future Intel® processors, targeted for optimizing streaming data movement and transformation operations common with applications for **high-performance storage, networking, persistent memory**, and various data processing applications.
 - Supported operations: Mov (e.g., memory move, crc generation...), Fill, Compare, Flush
 - Available on Intel's [Sapphire Rapids](#) Platform.
- In SPDK: How to leverage those memory offloading devices?

SPDK high level architecture



SPDK acceleration framework

SPDK
Application
or
Library

→
spdk_idxd_submit_fill()
spdk_idxd_submit_copy()
spdk_idxd_submit_crc32c()

→
spdk_ioat_submit_fill()
spdk_ioat_submit_copy()

→
spdk_accel_submit_fill()
spdk_accel_submit_copy()
spdk_accel_submit_crc32c()

DSA Public Low-Level Library

IOAT Public Low-Level Library

Acceleration Framework

DSA Module

IOAT Module

Software Module

SDPK DSA Feature support (publicly available from 20.10)

- Functionality: Support operations between (persistent) memory
 - ***Copy, Fill, CRC32C, Dual-cast, Compare***
- Optional batching support, application configurable batch sizes
- (May be deprecated in the future to make batch supporting inside the library)
- Dedicated work queues
- 2 pre-defined group/engine configurations selectable via API, easily changed in code as new/different configs are identified
- All DSA units are ganged together and collectively used, load balanced amongst SPDK threads
- Dynamic load balancing and flow control as threads come and go
- No interrupts or locking
- No dependency on other DSA software

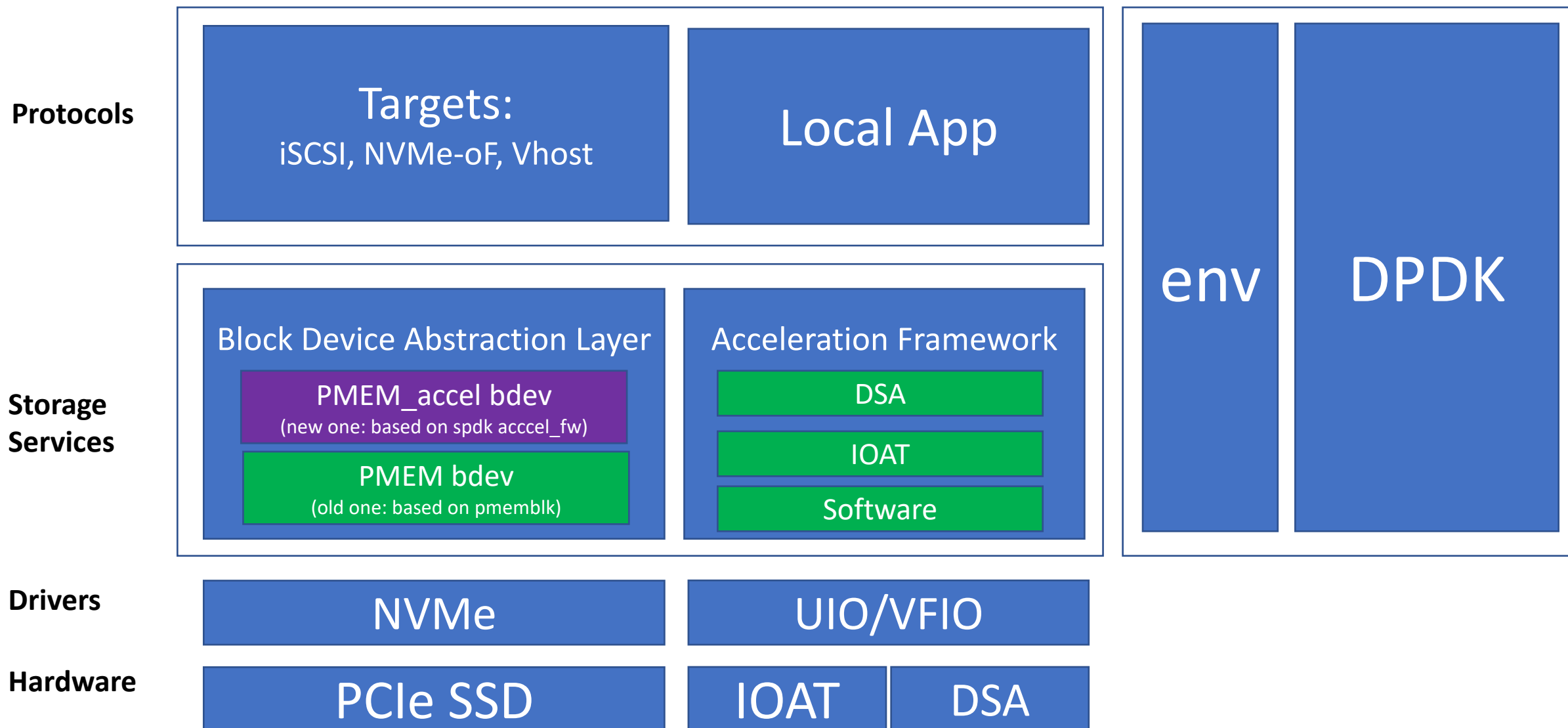
New PMEM bdev (pmem_accel) development status

Purpose

- Design and implement a new bdev (named as pmem_accel) based on ***persistent memory*** and ***spdk accel_fw*** to demonstrate that hardware-based memory offloading engine (e.g., IOAT/DSA) can help improve the performance.
- For this new bdev, we will not consider some complicated I/O usage scenarios (e.g., I/O atomicity for large size, sequence order, which should be addressed and guaranteed by upper layer logic).
- Our purpose is to leverage IOAT/DSA in spdk accel_fw to reduce the CPU participation overhead.
- Ideally, such offloading should be provided by PMDK, but currently PMDK library has some limitations, i.e.,
 - ***It only provides the synchronized API interface.*** (This will be an ongoing work inside PMDK)
 - ***It does not have the ability to plugin memory offloading engines inside PMDK.***

PMEM_accel bdev

(based on PMEM device and spdk accel_fw)



Two possible PMEM integration in SPDK to implement pmem_accel bdev

- We use PMEM in **App Direct** (AD) mode for persistency purpose. And the PMEM usage in AD mode can be divided into
 - Fsdax mode
 - Format the pmem device with “fsdax” choice. Users get block devices, e.g., /dev/pmem0, and users can create file system (especially kernel supported file system) upon that.
 - Dax Mode
 - Format the pmem device with “dax” choice . Then we get char devices, e.g., /dev/dax0.0, then you directly operate on such device.

Supported combinations

- Using FSDAX mode,
 - For CPU execution, you can choose to use libpmemobj/libpmem or not use.
- Using DAX mode,
 - For CPU execution, you can use libpmem or not use.
 - Limitation of using libpmem, pmem_map_file API's limitation on the file size.
- Data persistency while using DSA/IOAT in spdk accel_fw
 - For using DSA, we need to enable the ***persistency support by revising SPDK's idxd and accel library*** since our purpose is not having data lost while there is sudden power lost)
 - For using IOAT, ***the write durability cannot be fully guaranteed while there is power lost.***

Accessibility \ PMEM mode	FSDAX (mount choice: -o DAX)	FSDAX	DAX
CPU	Y	Y	Y
IOAT	N	Y	Y
DSA	N	Y	Y

PMEM_accel bdev summary

- Current decision (According to analysis in previous pages):
 - We choose to use **DAX mode** directly. But we want to continue discussing with customers whether DAX mode can fully satisfy customers' requirements. And we provide the ongoing [patch](#) in public to notify our progress.
 - In this usage, file system can not be viewed directly on the char device. The pmem devices can be partitioned into different character devices. A SPDK process will directly use one or several /dev/dax* devices directly.
- Ideally, we think that "**PMEM + FSDAX + -O DAX**" should be the perfect solution.
 - But it requires the kernel support especially in file system in Linux kernel , and it may not be available for a long time.
 - We will not use PMEM + FSDAX, because it can not provide the data durability while using offloading devices because offloading devices operate on page caches .
 - **Data path: Cache + DRAM -> BLOCK device -> PMEM**

SPDK, PMDK, Intel® Performance Analyzers

Virtual Forum

Conclusion

PMEM_accel bdev summary

- We realize the overhead/bottleneck of CPU to conduct data operations between (persistent) memory while performing storage workloads.
 - Data operations “between (persistent) Memory” is different from “data operations between Memory and device”. **For the latter case, there is (R)DMA capability on device, and CPU is not very busy.**
 - So, additional memory operation offloading engines (e.g., IOAT/DSA) are needed to use in order to save CPU resources.
- We introduce SPDK acceleration framework (spdk **accel_fw**) to leverage different memory offloading engines to address such problems and give the usage on PMEM device.
- Next step: We will continue improving spdk **accel_fw** and leverage it to solve different problems in storage aspect.
- For performance report will be shown in [spdk website](#) in the future.

SPDK, PMDK, Intel® Performance Analyzers

Virtual Forum

Q & A

Thanks!

The Intel logo is centered in the upper half of the image. It features the word "intel" in a white, lowercase, sans-serif font. A small blue square is positioned above the letter "i". To the right of the word "intel" is a registered trademark symbol (®). The background is a blue-tinted photograph of server racks in a data center.

intel®

SPDK, PMDK, Intel® Performance Analyzers

Virtual Forum