Transparent DAX Mappings

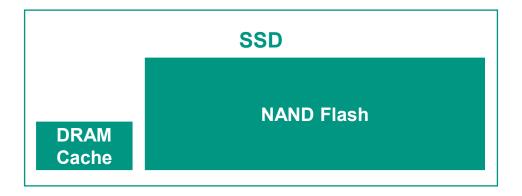
Towards Automatic Kernel Bypass with CXL-Based Hybrid SSDs

Yussuf Khalil, Daniel Habicht, Pascal Ellinger, Frank Bellosa, Javier González, Adam Manzanares, Vivek Shah

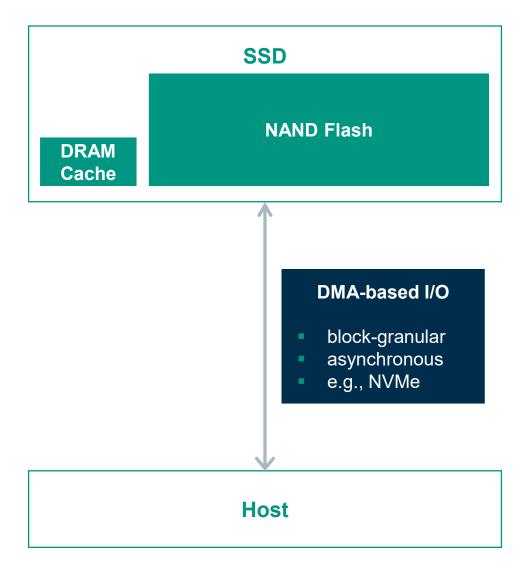




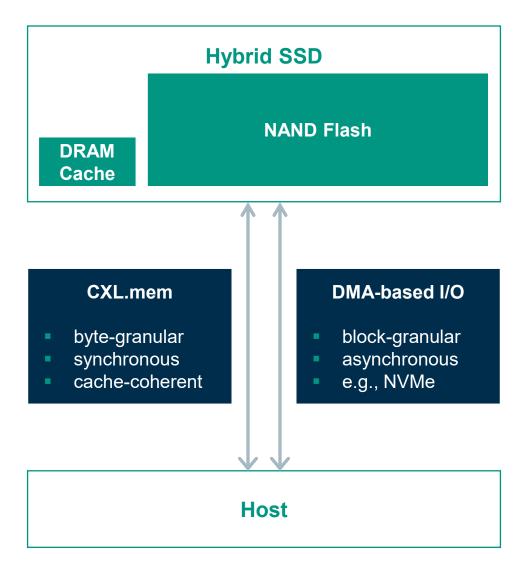
SAMSUNG



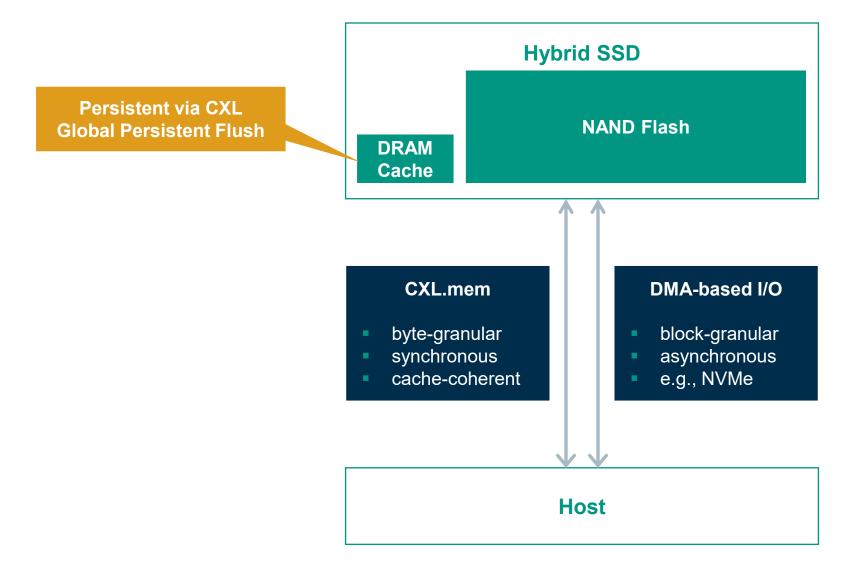




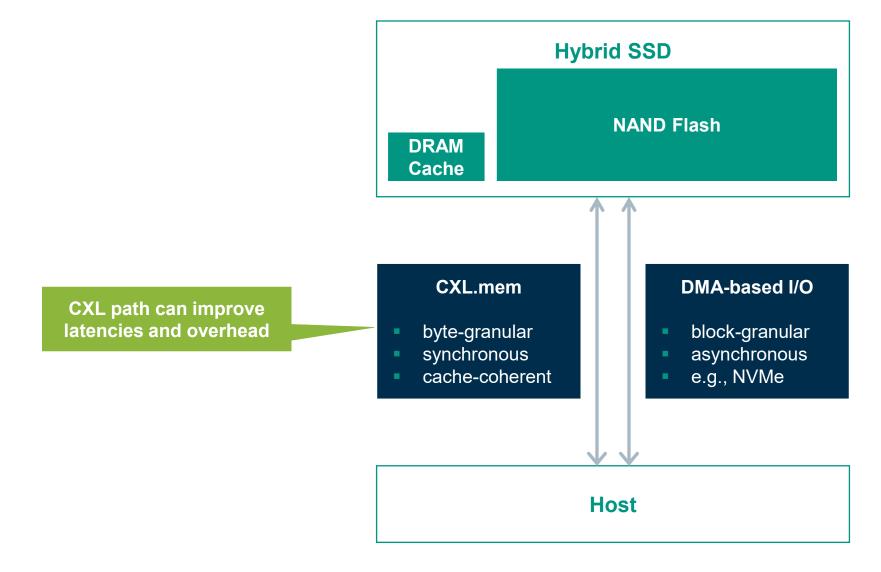




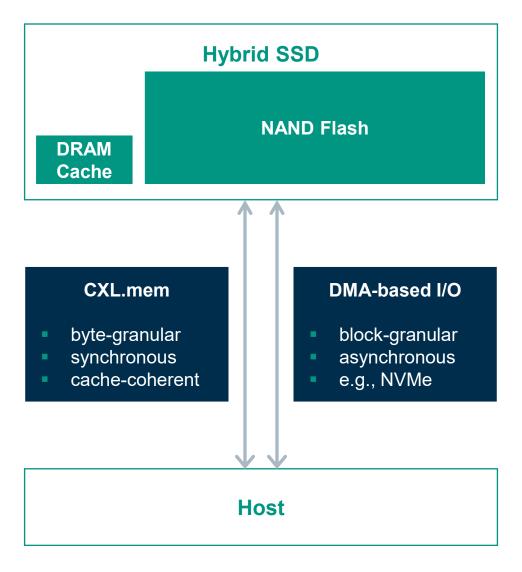




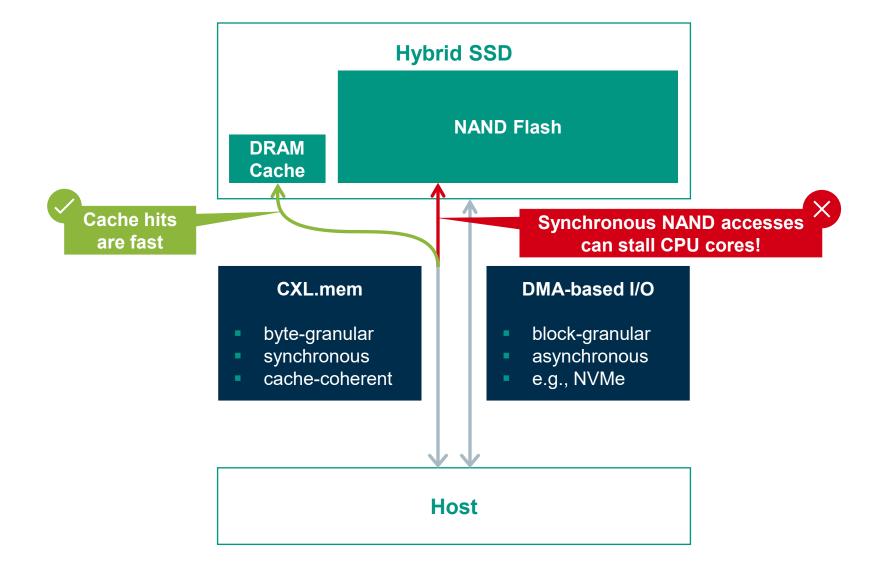




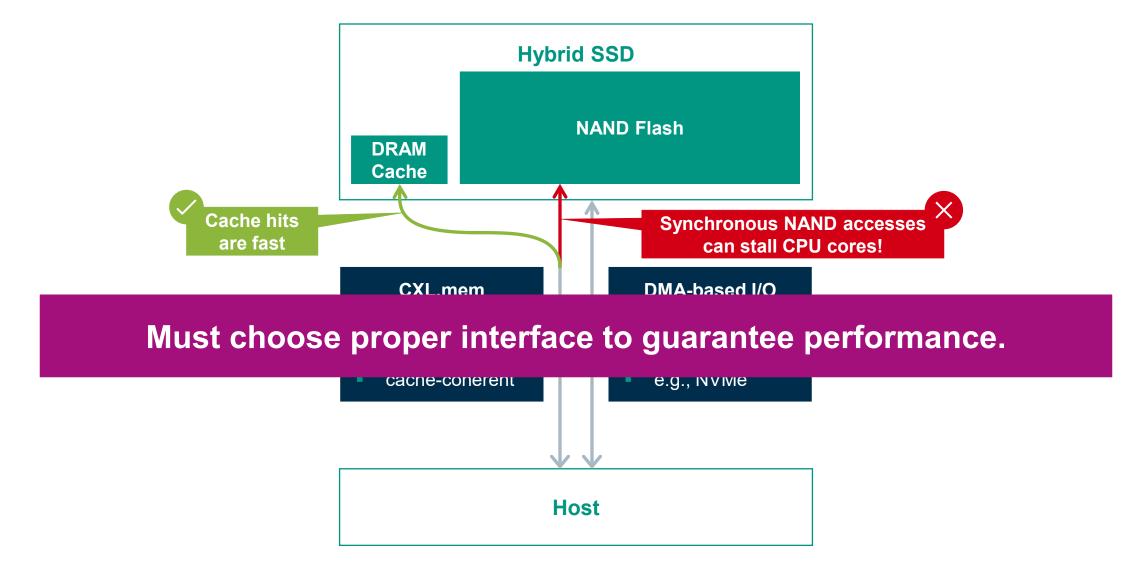




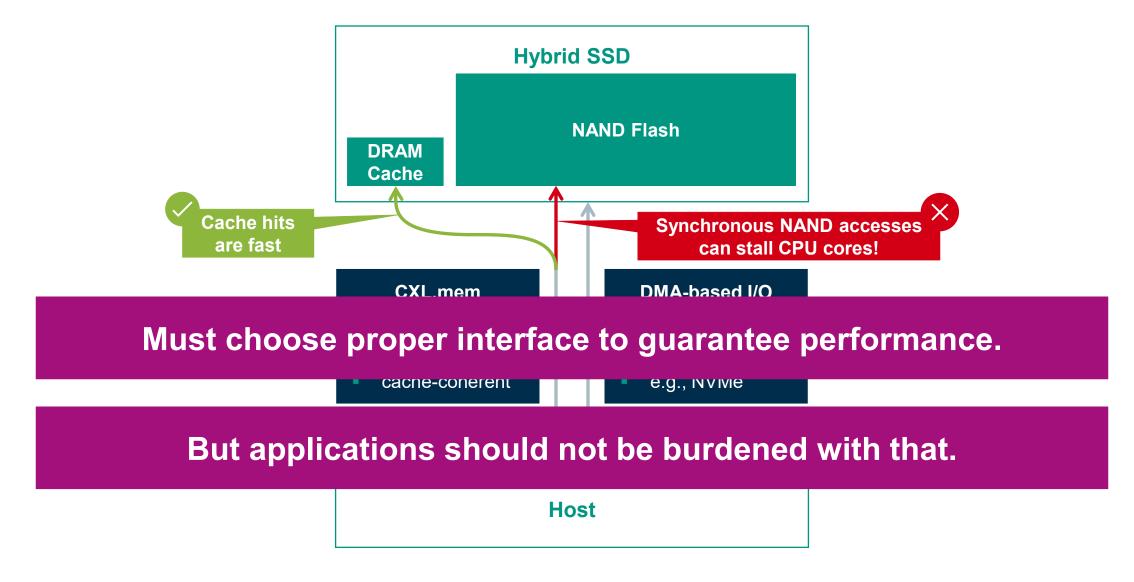


















Performance (or energy) benefit with traditional POSIX read()/write()





Performance (or energy) benefit with traditional POSIX read()/write()



Performance never worse than with NVMe





Performance (or energy) benefit with traditional POSIX read()/write()



Performance never worse than with NVMe



Simple to use, no additional developer effort required





Performance (or energy) benefit with traditional POSIX read()/write()



Performance never worse than with NVMe

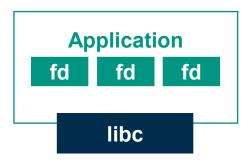


Simple to use, no additional developer effort required

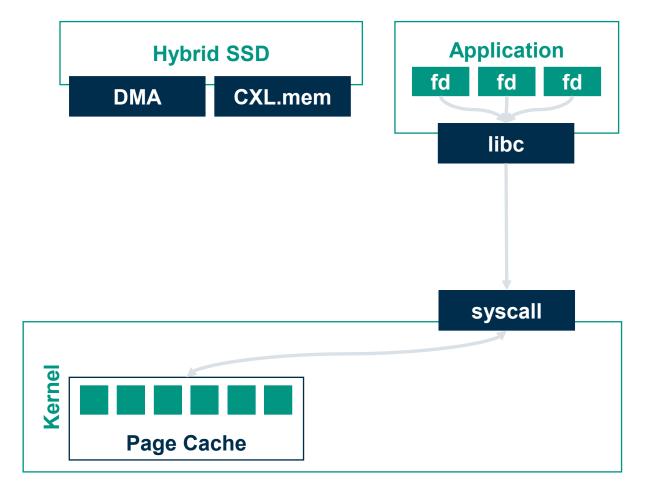
Make hybrid SSDs a drop-in replacement for NVMe SSDs.



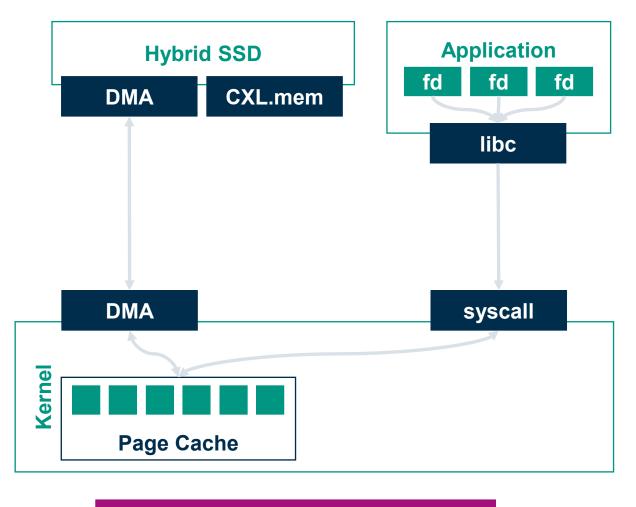






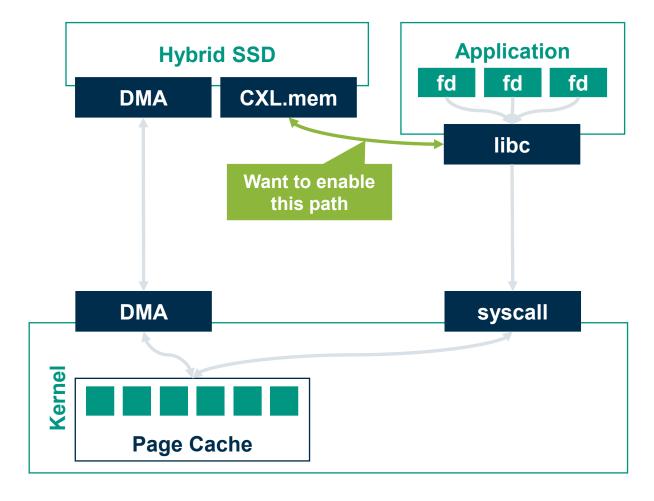




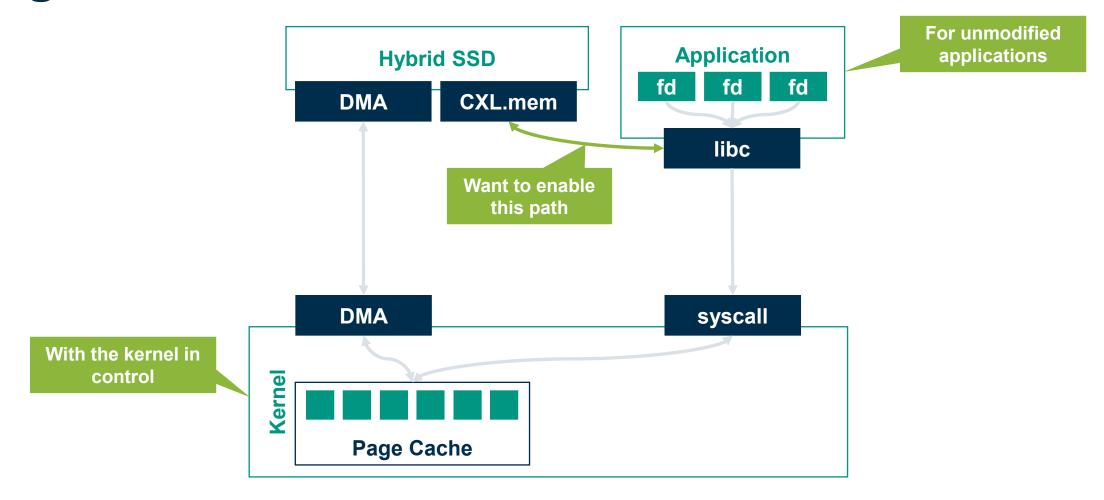


Traditional read()/write() path

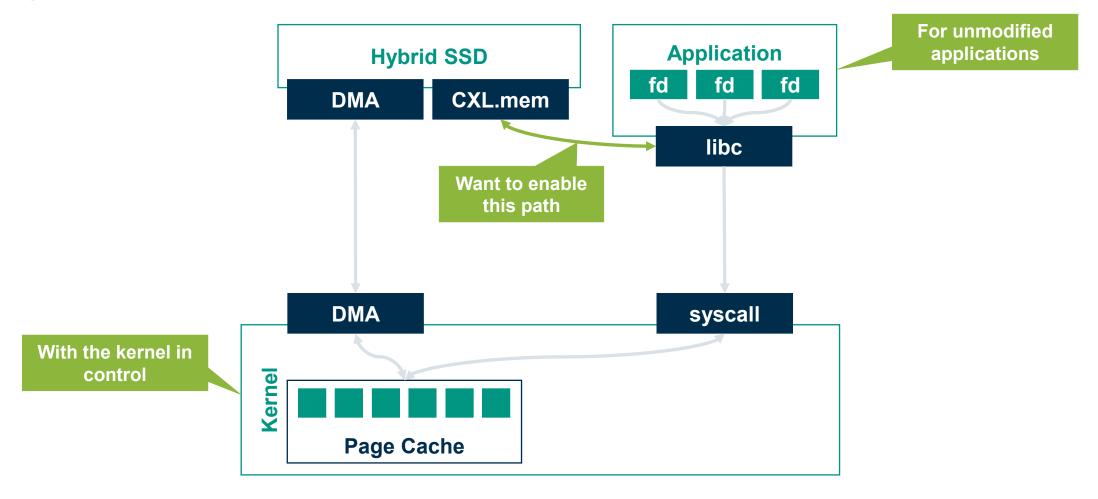








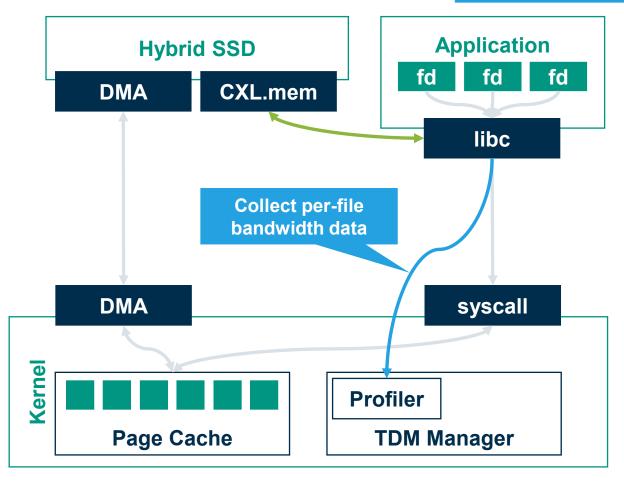




Idea: High-bandwidth files should use DAX.



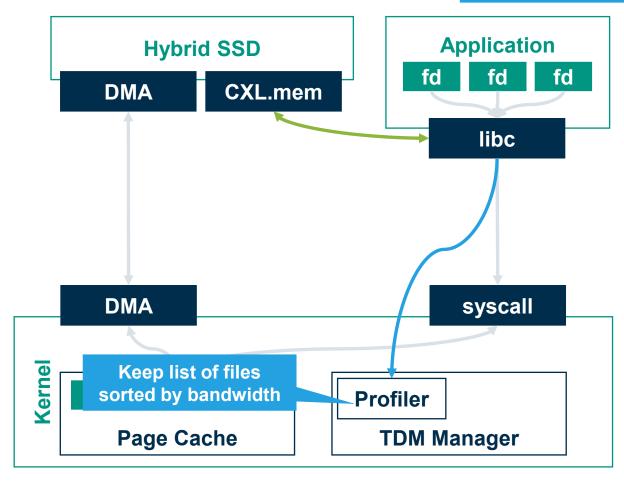




Idea: High-bandwidth files should use DAX.



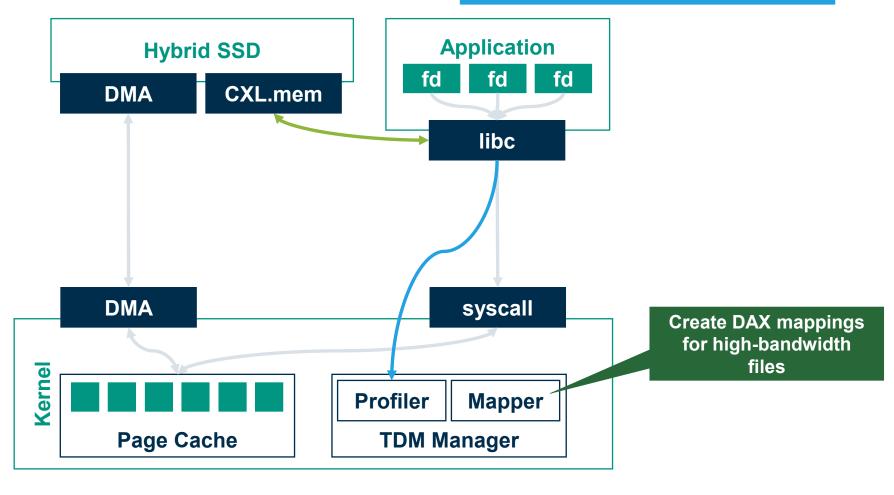




Idea: High-bandwidth files should use DAX.



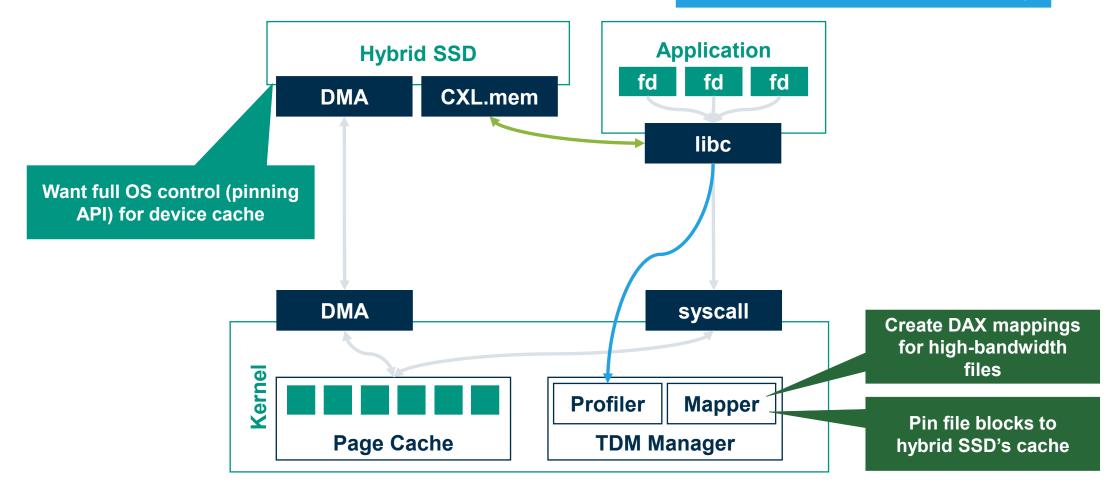




Use past bandwidth as prediction for future.



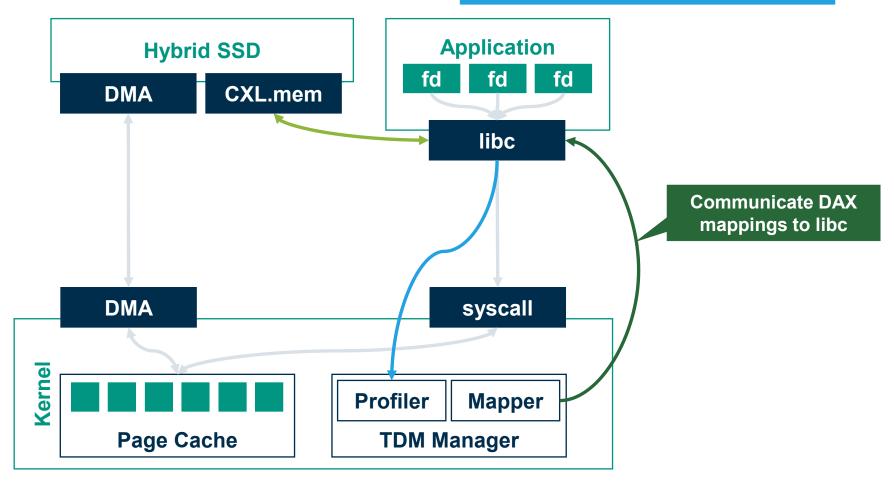
TDM = Transparent DAX Mapping



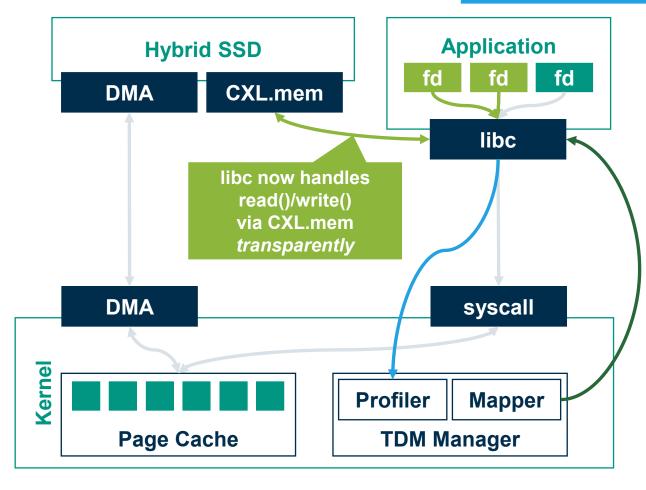
Use past bandwidth as prediction for future.



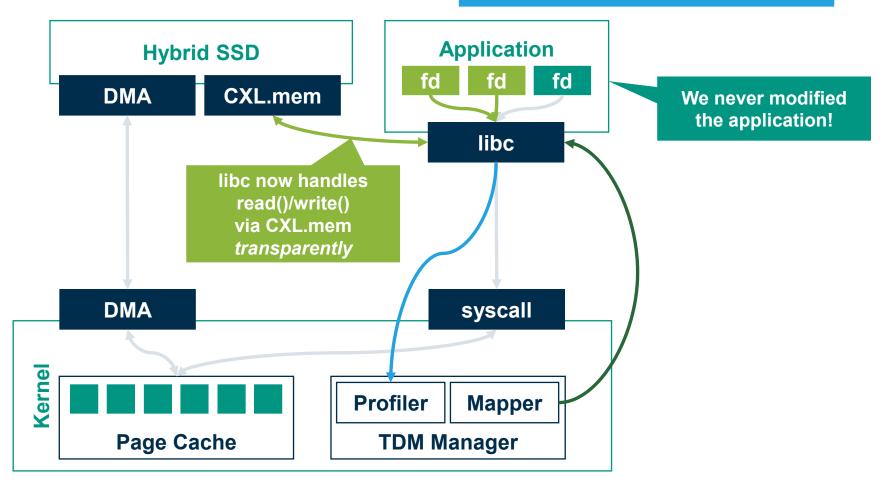




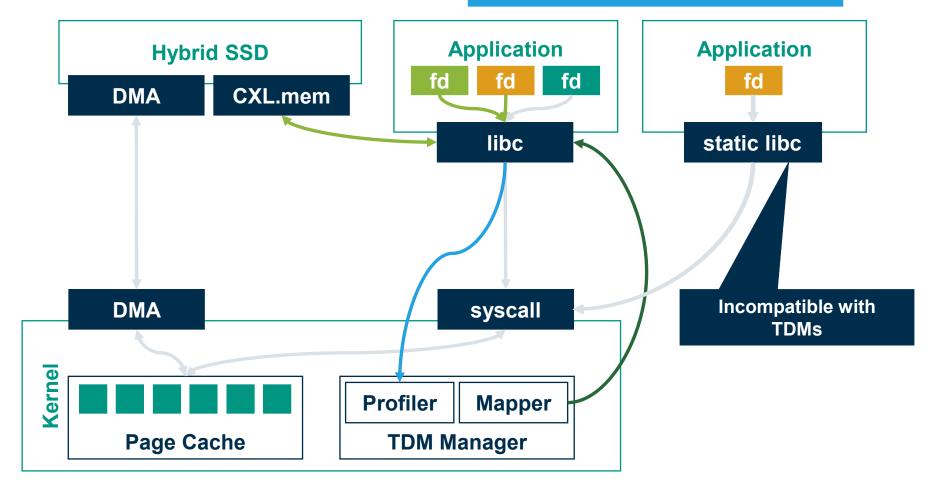






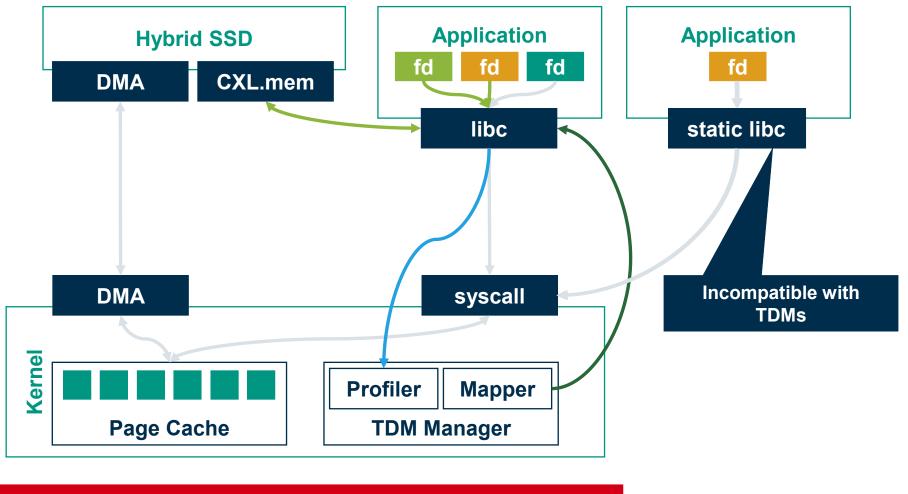






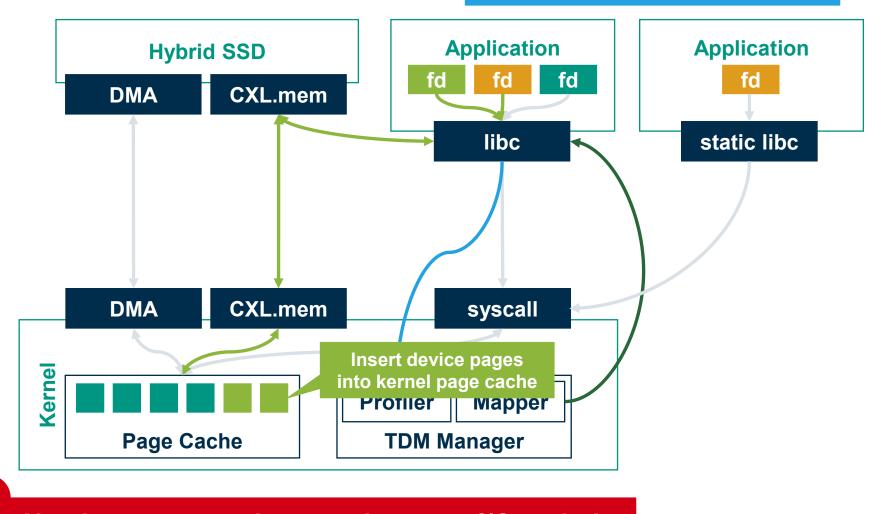


TDM = Transparent DAX Mapping



Need to ensure coherence between I/O paths!













No automatic profiling, files selected manually



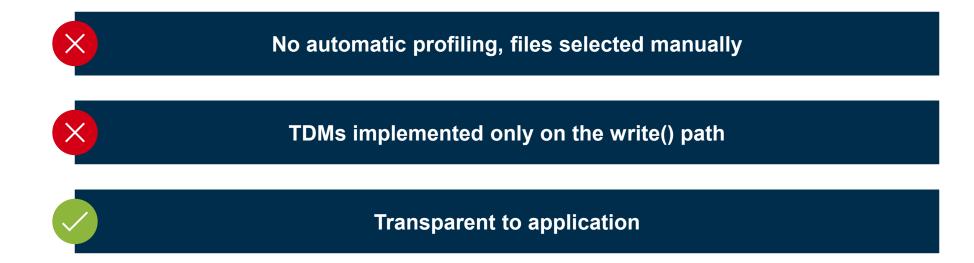


No automatic profiling, files selected manually

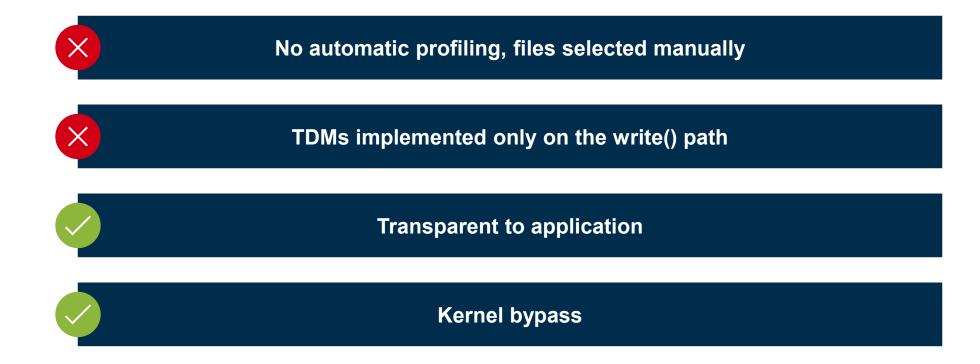


TDMs implemented only on the write() path



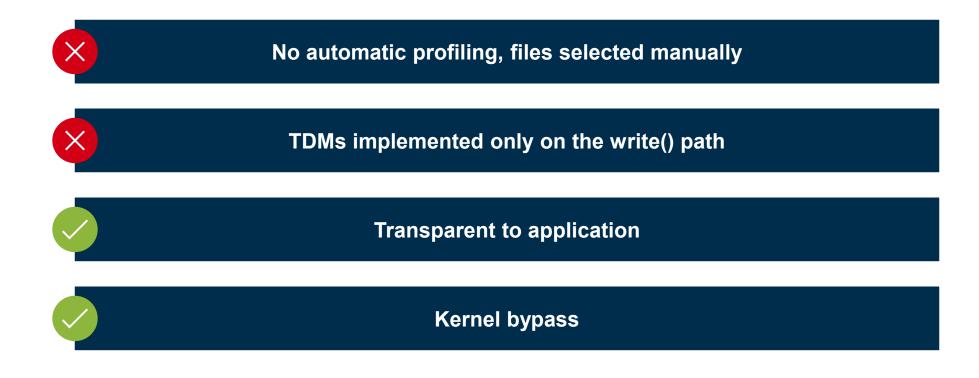








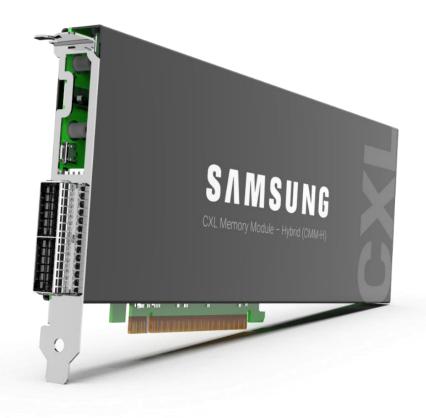
Preliminary Implementation



This is early-stage work.

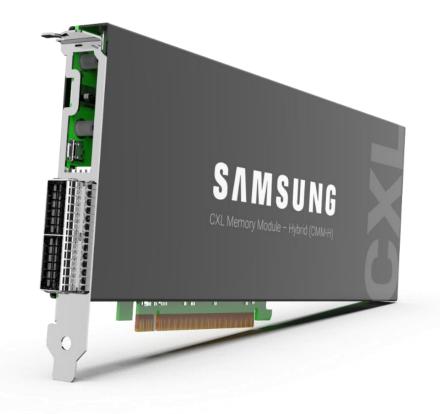


- FPGA-based CXL 2.0 x8 frontend
 - Entire storage capacity accessible via CXL.mem
- Internal Samsung PM9A3 960 GB NVMe SSD
 - Via internal PCle 4.0 x4 bus
- 48 GiB DRAM cache



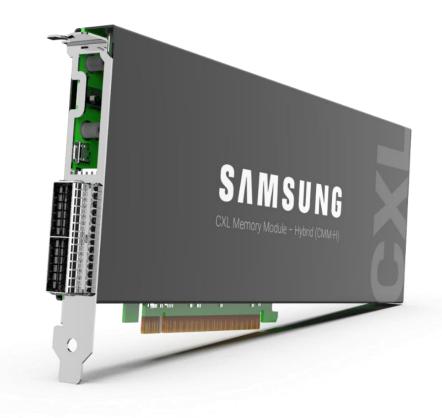


- FPGA-based CXL 2.0 x8 frontend
 - Entire storage capacity accessible via CXL.mem
- Internal Samsung PM9A3 960 GB NVMe SSD
 - Via internal PCle 4.0 x4 bus
- 48 GiB DRAM cache
- Internal LRU cache manager
 - Cache management via prefetch()/evict() API
 - No explicit cache pinning currently possible



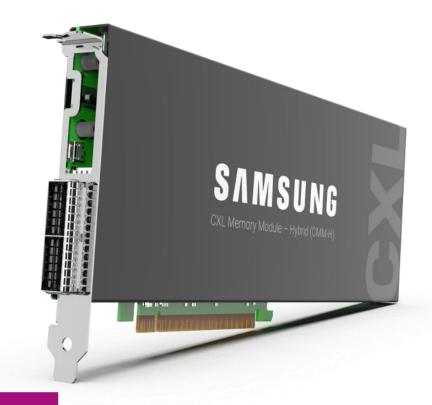


- FPGA-based CXL 2.0 x8 frontend
 - Entire storage capacity accessible via CXL.mem
- Internal Samsung PM9A3 960 GB NVMe SSD
 - Via internal PCle 4.0 x4 bus
- 48 GiB DRAM cache
- Internal LRU cache manager
 - Cache management via prefetch()/evict() API
 - No explicit cache pinning currently possible
- DMA interface not yet available
 - Kernel uses CXL as well



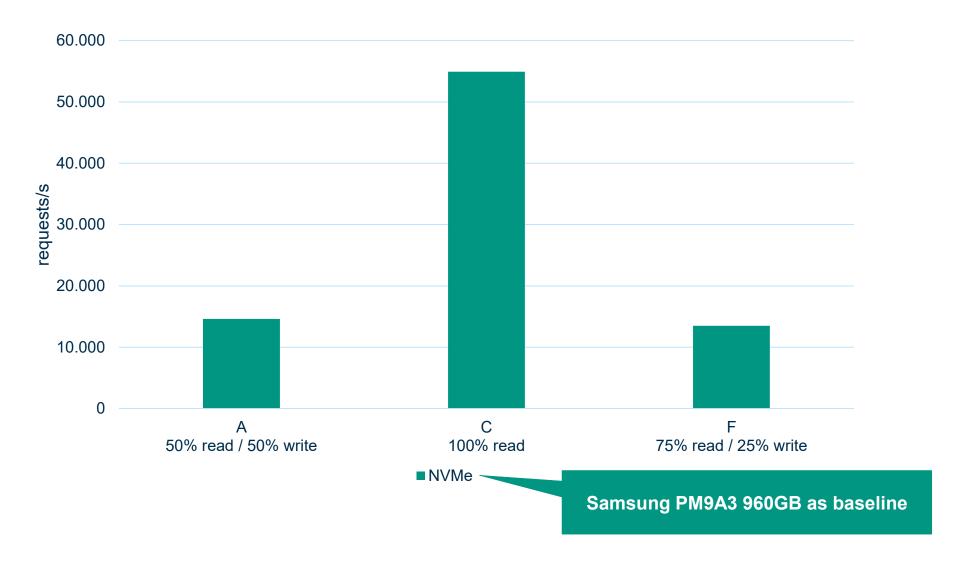


- FPGA-based CXL 2.0 x8 frontend
 - Entire storage capacity accessible via CXL.mem
- Internal Samsung PM9A3 960 GB NVMe SSD
 - Via internal PCle 4.0 x4 bus
- 48 GiB DRAM cache
- Internal LRU cache manager
 - Cache management via prefetch()/evict() API
 - No explicit cache pinning currently possible
- DMA interface not yet available
 - Kernel uses CXL as well

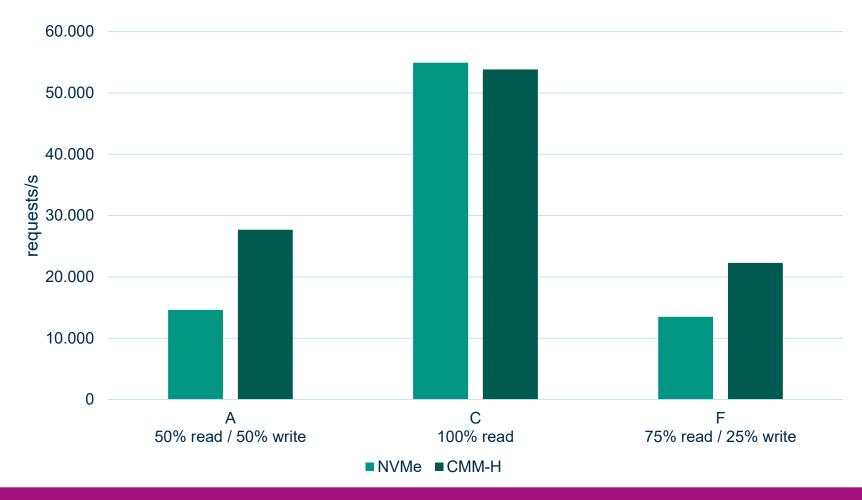


Current implementation focuses on kernel bypass functionality.



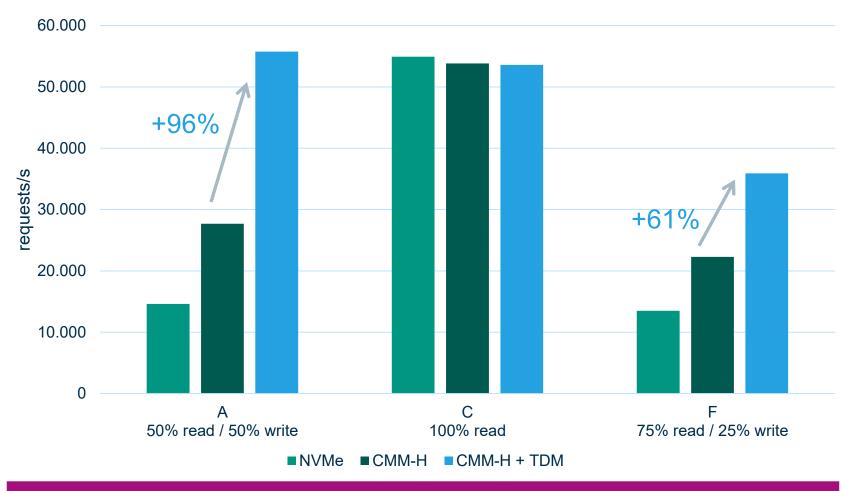






Without TDMs, CMM-H is already faster thanks to reduced latencies.

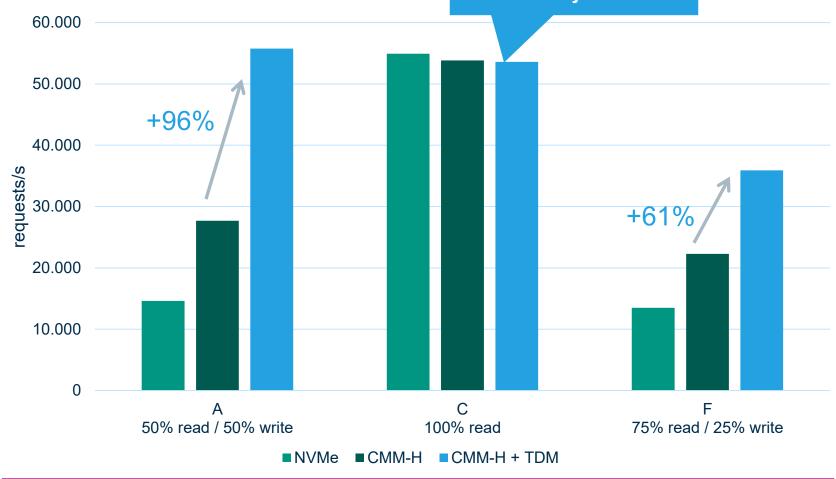




TDMs further increase throughput by up to 96%.



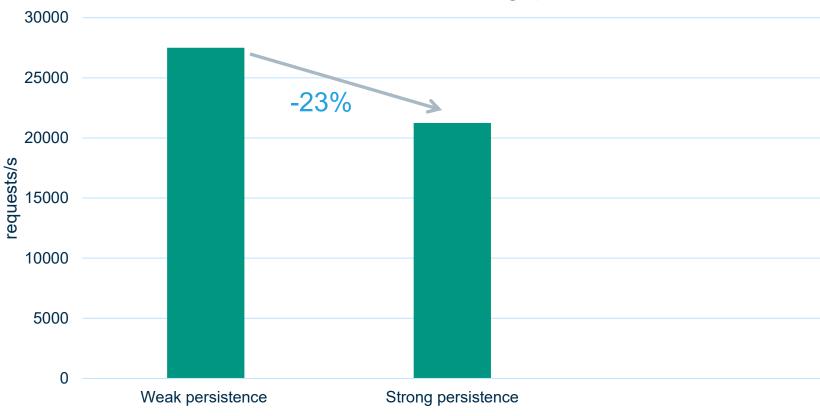
No difference expected for read-only workloads



TDMs further increase throughput by up to 96%.



Geometric Mean Throughput

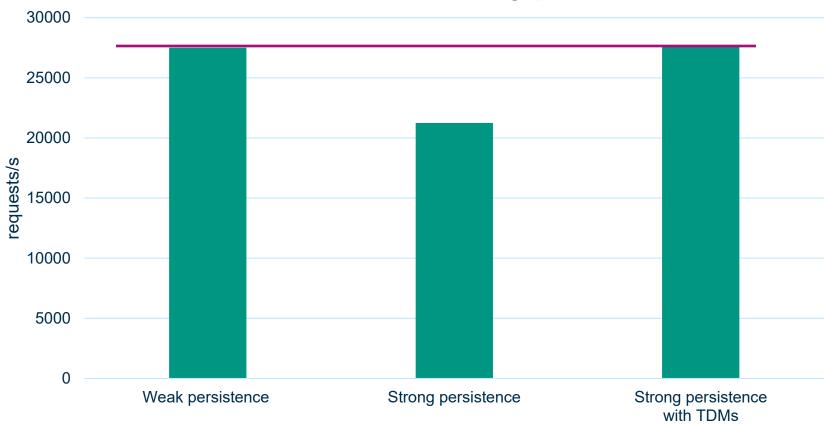




Strong persistence = fsync() after every write



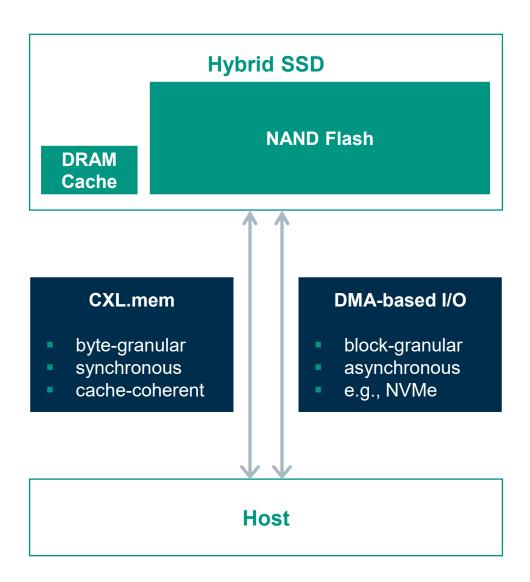
Geometric Mean Throughput



TDMs make persistence cheap.

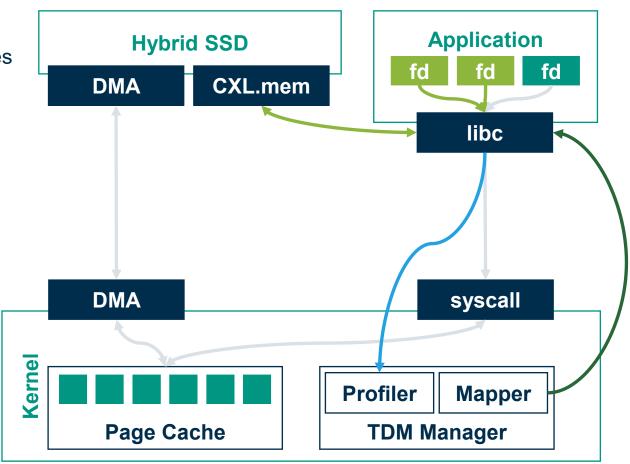


Hybrid SSDs warrant novel resource management strategies



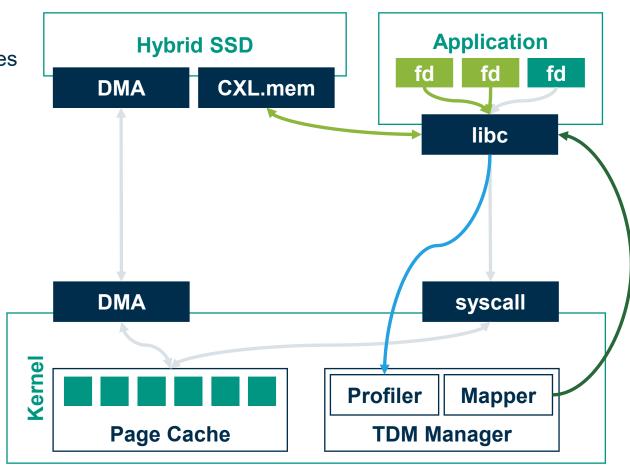


- Hybrid SSDs warrant novel resource management strategies
- Transparent DAX Mappings (TDMs)
 - ...put the OS in control of the device's cache
 - ...avoid synchronous NAND accesses
 - ...enable direct access for high-bandwidth files
 - ...employ lightweight and automatic profiling
 - ...do not require application modifications





- Hybrid SSDs warrant novel resource management strategies
- Transparent DAX Mappings (TDMs)
 - ...put the OS in control of the device's cache
 - …avoid synchronous NAND accesses
 - ...enable direct access for high-bandwidth files
 - ...employ lightweight and automatic profiling
 - ...do not require application modifications
- Preliminary implementation shows potential of approach
 - Up to 96% increased throughput in Valkey
 - Strong persistence becomes cheap





- Hybrid SSDs warrant novel resource management strategies
- Transparent DAX Mappings (TDMs)
 - ...put the OS in control of the device's cache
 - …avoid synchronous NAND accesses
 - ...enable direct access for high-bandwidth files
 - ...employ lightweight and automatic profiling
 - ...do not require application modifications
- Preliminary implementation shows potential of approach
 - Up to 96% increased throughput in Valkey
 - Strong persistence becomes cheap
- Our work makes the case for a pinning API in hybrid SSDs

