

Application Access to Persistent Memory – The State of the Nation(s)!

Stephen Bates, Paul Grun, Tom Talpey, Doug Voigt

Microsemi, Cray, Microsoft, HPE

The Suspects





We've Come a Long Way, Baby!





2016 Storage Developer Conference. All Rights Reserved.

Persistent Memory (PM)







Low Latency

Memory Semantics

Storage Features







- NVM Non-Volatile Memory. All types, including those that are not byte-addressable
- PM Persistent Memory. Sometimes PMEM is used but we use PM in this talk
- NVMe NVM Express. A block protocol to run over PCIe, RDMA or Fibre Channel. A SATA/SAS replacement.
- NVMP NVM Programming Model. Application-visible NVM behavior
- NVMf NVMe over Fabrics. NVMe extended over fabrics



Low Latency

| | Latency | Relative |
|--------------------|---------|-----------|
| LI Cache Read | 0.5ns | I |
| L2 Cache Read | 7ns | 14 |
| DRAM Read | I 00ns | 200 |
| The PM Opportunity | | |
| NVMe DRAM SSD Read | l Ous | 20,000 |
| NVMe NAND SSD Read | I 50us | 300,000 |
| SAS HDD | 500us | I,000,000 |

NVMe SSDs are (relatively) high latency!

PM provides persistence at memory-like speeds and semantics



Where Are We?





What is Needed?



2016 Storage Developer Conference. All Rights Reserved.

Lots of Moving Parts





Where does PM sit?

(Answer – anywhere it wants to)





2016 Storage Developer Conference. All Rights Reserved.

Rationalizing the Problem Space



Start with Consumers of NVM Services





2016 Storage Developer Conference. All Rights Reserved.

Application View



SNIA NVMP

- Describe application visible behaviors
- APIs align with OSs
- PM File System Actions
 - Map expose PM directly to applications
 - Optimized Flush make data persistent



Possible Stack for NVM Access





Optimizing Fabrics for NVM



Add persistence semantics to RDMA protocols

- Register persistent memory regions
- Completion semantics to ensure persistence, consistency
- □ Client control of persistence
- Solve the "write-hole" problem
- Lots of Initiatives underway!



Can we make this work for NVDIMMs? NVMe SSDs with CMBs?





| | NVMe |
|---|----------------------|
| ┣ | RDMA |
| | |
| | AWESOME |
| | |
| | |
| | PM |
| ┢ | RDMA |
| | |
| | AWESOME ² |

SD⁽¹⁾







...continuing down the stack









- NVMe over Fabrics Present an NVMe block device to a client over RDMA or Fibre Channel
- NVMe Controller Memory Buffers Standardize (persistent) PCIe memory on NVMe devices. NVDIMM-N on PCIe bus?
- LightNVM A low-level SSD interface that is more aligned to the underlying media (NAND)



Media

Media & Form-Factors

| | Category | Vendors | Comments | | |
|-----|-------------------------|---|---|-----------|--|
| NVM | DRAM Drop-In | Everspin Micron Toshiba/SK Hynix | DRAM like latency Super-Cap Replacement Not for bulk storage Memory Interface | PM | |
| | Storage Class Memory | Micron-Intel SanDisk Toshiba Crossbar Nantero | Faster than NAND, Cheaper/Slower than DRAM Byte Addressable Block and Memory Interfaces | | |
| | NAND | Micron Toshiba SanDisk SK Hynix Samsung | Lowest cost Slow (for NVM) Not byte addressable cheap and plentiful Block Interface | NOT PM | |



2016 Storage Developer Conference. All Rights Reserved.

PM Form Factors



NVDIMM-N

Not-NAND NVMe



NVDIMM-P









NAND NVMe

20



PM Form Factors



| Form-Factor | Media | Latency | Memory Semantics | Storage Features |
|---------------|---------------|---------|---------------------|---------------------|
| NVDIMM-N | DRAM/ MRAM | | | |
| NVDIMM-P | NAND/ PM | | | |
| Non-NAND NVMe | DRAM/ PM | 00 | | |
| NAND NVMe | NAND | | | |

Form factors impact Features (No DMA engines on a DIMM!)



PM Scenarios

Windows

2016 Storage Developer Conference. All Rights Reserved.

os Support C

- PM region as a block device (a la persistent ram disk).
- Filesystems support: direct access to the memory (e.g DAX), PM aware FS (e.g. m1fs).
- You can put your files, databases etc. on top.
- Remember we are crawling right now!
- Soon: Shared persistent memory



Libraries and Toolchains



HARD



EASY

Make it easy for applications to utilize PM, regardless of OS and ARCH!

```
section .text
global start
start:
         mov edx, len
         mov ecx, msq
         mov ebx,1
         mov eax, 4
         int. 0x80
         mov eax, 1
         int 0x80
section .data
msg db 'Hello, PMEM World!', Oxa
len equ $ - msg
```



Call to Arms





Call to Arms

- Libraries and Toolchains: NVML for non-x86, integration into glibc/gcc etc.
- Media & Form Factors: Production PM, appropriate PM form factors.
- Protocols and Interconnect: Enhancements to NVMe and RDMA, PM over Fabrics, standardization of memory channels.
- OS, new OSes?
 OS, new OSes?



Conclusions

- We are almost walking! Help out if you can
- If you want sub 10us access to persistent data then PM may be for you
- The CPU vendors have a lot of say in the interconnect but some open options exist too
- Toolchains, libraries and OSes are adapting
- New applications will complete the jigsaw and lead to revenue

