

# Persistent Memory & Computational Storage

## Industry Status and Update SDC India '19



- PMEM Hardware...and the associated programing model
- What everyone already <u>should</u> know about pmem...
- What everyone forgets...
- Ways to use pmem with no app modifications
- Ways to use pmem with app modifications
- Learnings so far
- Where we're heading
- An introduction to Computational Storage

### A Fundamental Change Requires An Ecosystem

HARDWARE

STANDARD.

PLATFORMS

SOFTWARE

#### VMWare<sup>®</sup> ORACLE

Windows Server 2016

Microsoft

Windows 10 Pro for Workstations

🛆 Linux

- Linux Kernel 4.2 and later
- VMware, Oracle, SAP HANA early enablement programs



- Multiple vendors shipping NVDIMMs
- SNIA NVDIMM Special Interest Group (formed Jan'14)
- Successful demonstrations of interoperability among vendors

JEDEC JESD245B.01: Byte Addressable Energy Backed Interface (released Jul'17)
JEDEC JESD248A: NVDIMM-N Design Standard (released Mar'18)
SNIA NVM Programming Model (v1.2 released Jun'17)

 unfit ACPI NVDIMM Firmware Interface Table (v6.2 released May'I7)



- All major OEMs shipping platforms with NVDIMM support
- Requires hardware and BIOS mods

## **JEDEC-Defined NVDIMM Types**



- Host has direct access to DRAM
- NAND flash is only used for backup
- Capacity = DRAM (10's 100's GB)
- Latency = DRAM (10's of nanoseconds)
- Endurance = DRAM (effectively infinite)
- No impact to memory bus performance
- Low cost controller can be implemented
- Specifications completed and released
- Ecosystem moving into mature stage





#### **NVDIMM-P**

Host is decoupled from the media (agnostic to PM type) New protocol to "hide" non-deterministic access Capacity = PM (100's GB+) Latency = PM (>> 10's of nanoseconds) Endurance = PM (finite) Likely to impact memory bus performance Complex controller & buffer scheme likely required Specifications still under definition (2H'19 release?) No ecosystem yet, likely DDR5 timeframe

**SNIA** 

PERSISTENT MEMOR



CNTLR

© 2018 Storage Networking Industry Association. All Rights Reserved.



### Everyone should know...



#### Persistent memory…

- Allows load/store access like memory
- Is persistent like storage
- Exposed to applications using SNIA NVM TWG model

#### What isn't persistent memory:

- Something that can only speak blocks (like a disk/SSD)
- Something that is too slow for load/store access
  - > TWG's language: Would reasonably stall the CPU waiting for a load to complete















Memory Controller

PM

DRAM (cache)

DRAM (cache)

Memory Controller

PM



### Using PM as a fast SSD

- Storage APIs work as expected
- Memory-mapping files will page them into DRAM

### Using PM as DAX

- Storage APIs work as expected
- No paging (DAX stands for "Direct Access")

### Using PM as volatile capacity

- Just big main memory
- Vendor-specific feature











# **Application Modification: pmemkv**



### libpmemkv

- Experimental
- General-purpose key-value store
- Multiple pluggable engines
- Multiple language bindings
- Productization underway
- Caller uses simple API
  - But gets benefits of persistent memory

SNIA





Lots of ways to use PM without app modifications
 Try first to use existing APIs

- Example: app that can be configured for SSD tier
- Try next to use highest abstraction possible
  - Key-value store, simple block or log interfaces
- Try next to use a transaction library
  - libpmemobj

Finally, if you must program to raw mapped access

### Where we're heading



#### More transparent use cases

• Either kernel or library features, transparent to app

#### More high-level abstractions

Easier to program, less error prone

#### More support for experts as well

- More features in transaction libraries
- More language integration
- Faster remote (RPM) access

## **RPM...Some Challenges, But Usable**

#### NUMA, by definition

- Probably okay, just be aware of it
- Generally requires asynchronous operation
  - Including delayed completions
- Networks introduce unavoidable latencies
  - As long as the application can tolerate it
- Transaction model will often favor pull vs push operations
  - not necessarily native to the way application writers think

Net-net, probably can't treat remote and local PM exactly the same. Not quite transparent, but close.

SN



- Java is a very popular language on servers, especially for databases, data grids, etc., e.g. Apache projects:
  - Cassandra
     Lucene
  - Ignite
     Spark
  - HBase
     HDFS
- Want to offer benefits of persistent memory to such applications

**PM Storage Engine for Cassandra** 



- Cassandra is a popular distributed NoSQL database written in Java
- Uses a storage engine based on a Log Structured Merge Tree with DRAM and disk levels
- Could persistent memory offer Cassandra opportunities for simpler code and improved performance?









#### **Software - Persistent Memory Storage Engine**



Cassandra Pluggable Storage Engine API https://issues.apache.org/jira/browse/CASSANDRA-13474

Cassandra Persistent Memory Storage Engine https://github.com/shyla226/cassandra/tree/13981\_llpl\_engine

Low-Level Persistence Library (LLPL)

https://github.com/pmem/llpl

Java VM (JDK 8 or later)

Persistent Memory Development Kit (PMDK)

https://github.com/pmem/pmdk

Linux OS

Persistent Memory



SNIA – Persistent Memory Resource Page <u>https://www.snia.org/PM</u>

2019 Persistent Memory Summit <u>https://www.snia.org/pm-summit</u>



# **Computational Storage?**

### **Compute, Meet Data**

- Based on the premise that storage capacity is growing, but <u>storage</u> <u>architecture has remained</u> <u>mostly unchanged</u> dating back to pre-tape and floppy...
- How would you define changes to take advantage Compute at Data?



SNI

## What is Computational Storage





Distributed-Processing and Data-Driven



#### When to Computational Storage

- Large Data Transfers, PCIe is bottleneck
- Data pre/post processing & analysis
- Data can bypass the host video delivery
- Ability to move Software App to Storage

#### When to USE NVMDIMM

- Compute heavy with small data-transfer
- Small data compute in-memory compute
- Little to no parallelism



### **A New Product Category**



#### Computational Storage Device (CSx)



Computational Storage Drive (CSD)
 Computational Storage Processor (CSP)
 Computational Storage Array (CSA)





### **Possible Architectures**



### Computational Storage is a Real Market

Customers are deploying today

### Solutions exist and will continue to grow

Making the 'uniform' helps adoption

### Standardizing the host interaction is vital

We NEED more Support from Users/SW Solutions

### Working with all TWG/SIGs & initiatives is key

• Joining forces and cross-membership adds to success





# Thank You!! <u>www.SNIA.org/Computational</u> <u>www.SNIA.org/forums/sssi/NVDIMM</u>