## Persistent Memory in 2020: Introduction



## Dave Eggleston Intuitive Cognition Consulting



November 2020

Source: Intel





## DRAM! (But...)







1. More cores, more channels, DDR bandwidth/core has... slowed down?!

Santa Clara, CA November 2020

Source: Micron







## DRAM! (But...)

HPC System 1		HPC System 2	
Cores	56	Cores	112
Memory	112GB	Memory	1.8TB
Memory Power	50W	Memory Power	~700W!

2. More cores, more memory = Order of magnitude more Power!

Santa Clara, CA November 2020

Source: The Next Platform, Intuitive Cognition Consulting



## H

## DRAM! (But...)

#### Possible DDR5 Server Subsystem



Different types of DIMMs will be common in DDR5

3. More cores, more channels = Substantially more DDR board complexity!

Santa Clara, CA November 2020

Source: JEDEC

## Can PM deliver where DRAM is falling short?

Bandwidth
Power
Complexity

. Cost/GB

And.

Plus what application(s) does PM accelerate?

## PM Media Options



PM media choices have been broad, but seeing consolidation

Source: GLOBALFOUNDRIES



## PM Media: Who is Doing What



- Only PCRAM has shown the characteristics necessary to deliver the PM capacity, cost, and performance.
- Memory giants are now consolidating around PCRAM.
- MRAM is a good embedded NVM for foundries – just not suitable for high capacity PM.

<u>Source</u>: Industry announcements & scuttlebutt



## PM Media: Optane Series 100 to 200

	100 Series	200 Series	% Improvement		
DIMM Capacities	128/256/512 GB	+			
DDR Frequencies	2666, 2400, 2133,	2666 MT/sec			
	1866 MT/sec				
Performance shown for 256GB DIMM @ 15W (Best performing DIMM)					
Endurance 100% Writes	363 PBW	497 PBW	37%		
256B					
Endurance 100% Writes	91 PBW	125 PBW	37%		
64B					
BW 100% Read	6.8 GB/sec	8.1 GB/sec	19%		
256B					
BW 100% Read	1.75 GB/sec	2.03 GB/sec	16%		
64B					
BW 100% Write	2.3 GB/sec	3.15 GB/sec	37%		
256B					
BW 100% Write	0.58 GB/sec	0.79 GB/sec	36%		
64B					

200 series Optane focused on endurance and BW improvements

Santa Clara, CA November 2020

Source: Intel, Intuitive Cognition Consulting



## PM Value Propositions: TCO, Throughput, Speed

## Intel Optane PMem Delivering Real World Benefits



Santa Clara, CA November 2020

Source: Intel



## DB App: PM + DRAM Outperforms DRAM alone

### **DRAM-like Performance**



#### Sysbench QPS on MySQL



## AI/ML App: PM acceleration

#### Case Study: AI / Machine Learning – Facebook's DLRM

#### Background

Customer Type: AL / ML Customer

#### **Business Challenge:**

Dynamic and Scalable Production Inferencing

#### Platform:

 Innovative Big Memory Computing platform for leveraging persistent memory for realtime, AI/ML and Advanced Analytics and extensible to all memory – Centric workloads.

#### Software:

 Software Defined Architecture extracting performance benefits of cutting edge hardware supporting workload portability to truly compute anywhere with the memory speeds.





#### Result:

- Customer has state of the art AI/ML Big Memory Platform that is can scale and deliver performance when Data is Greater than Memory
- Achieved flexible software defined platform Big Memory Computing capabilities and poised for future dynamic model and data growth

## Intervention 10x inferencing acceleration, NO App rewrite!

Santa Clara, CA November 2020

Source: Penguin Computing



DDR (now) and CXL (future) for PM attachment

Santa Clara, CA November 2020

Source: SMART Modular



## PM Form Factor: DDR DIMM to E1.S







E1.S form factor as the PM successor to the DDR DIMM

Santa Clara, CA November 2020

Source: SMART



## PM: \$/GB versus DRAM

PMEM DRAM 1 x 512GB \$13.86/GB 1 x 256GB \$7.02/GB \$18.94/GB 1 x 128GB \$4.00/GB \$13.67/GB 1 x 64GB \$7.65/GB 1 x 32GB \$8.43/GB 1 x 16GB \$9.37/GB

August 2020 prices from online resellers. Prices vary widely.

PM priced ~30% to 50% of DRAM (\$/GB basis)

Santa Clara, CA November 2020

Source: The Next Platform, MemVerge



## **PM Introduction Summary**

- DRAM challenges ahead
- PM can address bandwidth, power, complexity, cost
- PM media consolidation around PCRAM
- 2<sup>nd</sup> generation PM is here now with improvements
- DB and AI/ML applications accelerated using PM + DRAM
- PM priced attractively versus DRAM
- PM form factor from DDR DIMM to E1.S



# PERSISTENT MEMORY

## Stay tuned for details at <u>www.snia.org/pm-summit</u>

- 9<sup>th</sup> annual Summit April 21-22, 2021
- Deep dive into the latest memory and storage developments
- New for 2021
  - Expanded to two days 30+ sessions
  - Virtual stream live or watch on demand

## Thank You!



## Everything You Need To Know For Success