



Flash Memory Summit

Persistent Memory Is the Answer to Today's Data Center Challenges

Jung Yoon, Ranjana Godse – IBM Supply Chain Engineering
Andrew Walls – IBM Flash Systems

Santa Clara, CA
August 2018



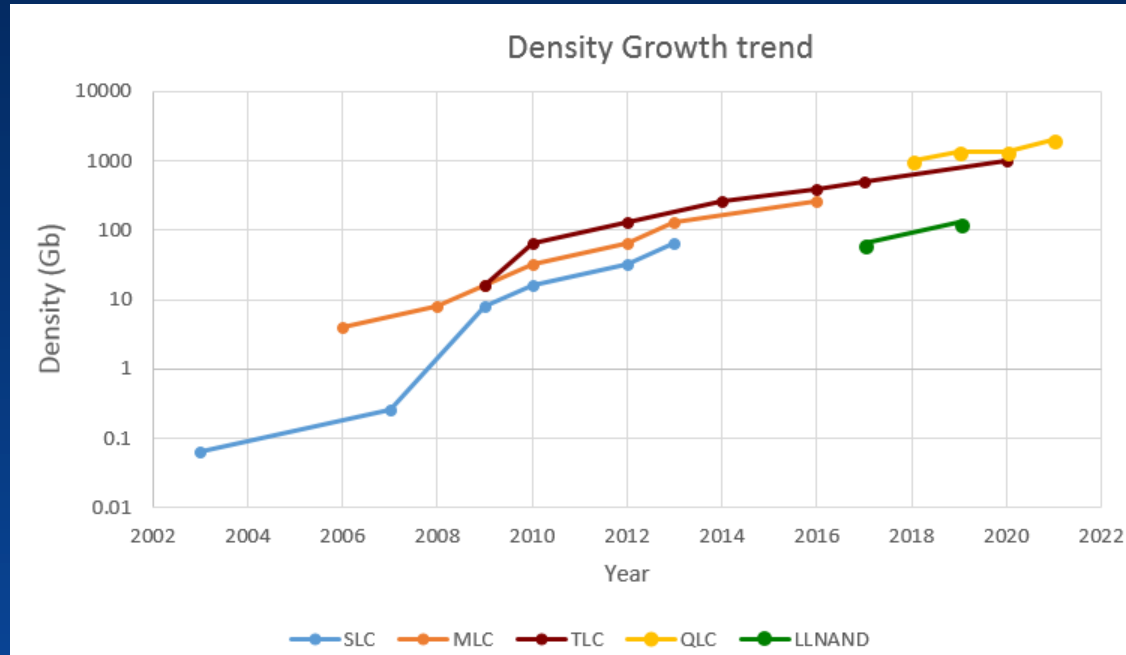
Agenda

Flash Memory Summit

- 3D NAND Scaling and cost reduction has fueled a disruption in the data center
- Persistent Memory and its arduous journey into the Data Center
- Use Case #1 Persistence
- Use Case #2 External Storage Consistency Enhancer
- Use Case #3 Hot Tier
- Use Case #4 Super Fast Tier 0
- What will Change
- Where are things headed



Memory density trend

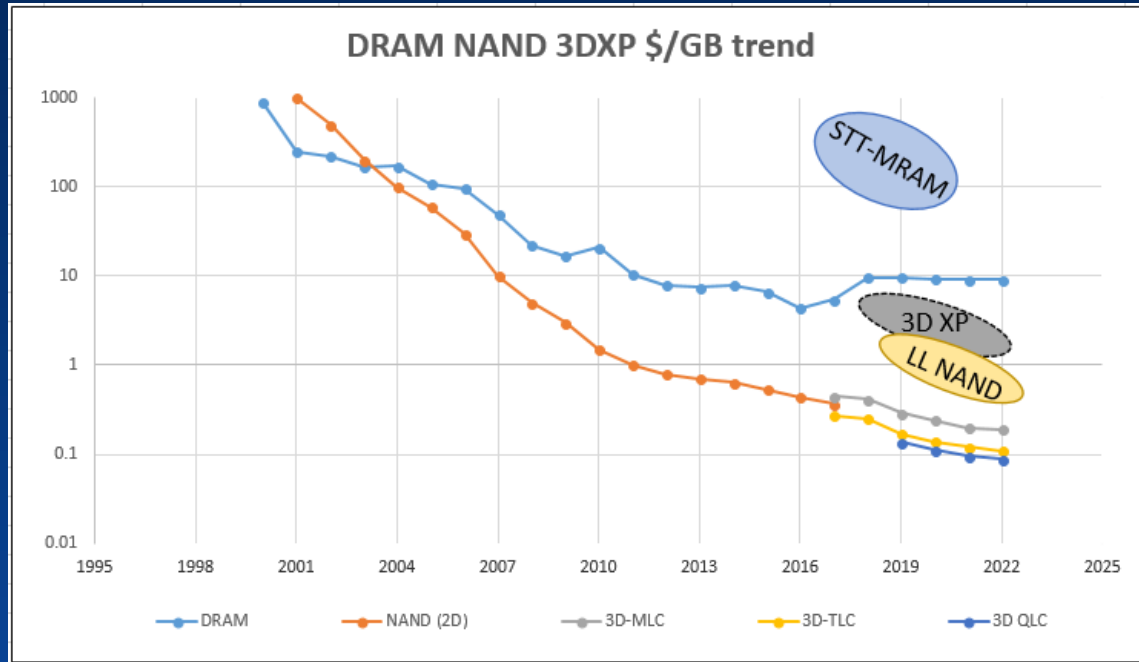


- Flash density growth will continue via 3D NAND layer count scaling 2018-2022+
- Significant fab investments ongoing for new flash Fabs in 2018-20
- 3D NAND layer count scaling provides clear path for flash density growth – TLC, QLC, LL-NAND
- 3D NAND scaling on a 18-24 month cadence thru 2020 – 64L > 96L > 120L > 190L



Flash Memory Summit

\$/GB trend for DRAM, NAND and SCM



- Flash Market Demand growth – driven by high density 3D NAND, Enterprise SSD, Mobile applications growth
- 3D NAND bit cost reduction achieved by 3D Cell layer increases
- Strong penetration of 3D TLC in Cloud Datacenter & Enterprise Storage in 2018-19.

Santa Clara, CA
August 2018



Flash Memory Summit

Persistent memory and talk of Persistent memory is Everywhere!

- MRAM in embedded, other places. Even in the new IBM NVMe Flash Core module
- 3DXP Optane – Fast Tier
- 3DXP on DRAM bus
- LL NAND
- Spin Torque, RRAM, Phase Change
- Exotics



Flash Memory Summit

Persistence

- DRAM is one of the most important inventions of all time
-- BUT --
- Super caps and batteries are a pain!
 - (No offense to any super cap providers or battery providers in the audience)
- MRAM is providing a small persistent embedded memory for many applications
- Used in New FS9100 FCM
- 3DXP could eventually get rid of batteries in fast write caches
- However, not the highest priority.





Flash Memory Summit

Storage Consistency Enhancer

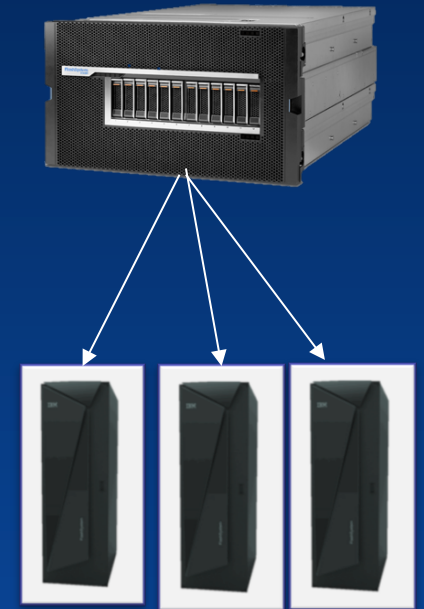
- Flash has significantly reduced average response time
- Snapshots, Data Reduction metadata (LSA), and other structures being out of cache can cause latency spikes
- Storing metadata, bitmaps and other structures in persistent fast storage can make response time consistency



Flash Memory Summit

Hot Tier

- Flash entered the data center via this mechanism
- Still prevalent in thousands of data centers around the world
- Topologies
 - In the storage array
 - As a tier 0 in a virtualized storage system
 - In the server as an HBA or NVMe card





Flash Memory Summit

Super Fast Tier 0

- Today's tier 0 with flash are $\sim 100\mu\text{S}$
- NVMe over Fabrics can reduce that to $\sim 70\mu\text{S}$
- Optane/LL Nand - $\sim 25\mu\text{S}$
- 3DXP as DRAM – Less





Flash Memory Summit

How will Fast Tier 0 be used with Persistent Memory

- FS 900 Tier 0 has been very successful
- Relational database accelerators
- Metadata for Scaleout filesystems
- Fast Write Buffers for HPC
- Smart Tiering in Analytics
- Still can not execute out of it unless on DRAM bus



Flash Memory Summit

Applications that can take advantage

- Throughput and tiers in Machine learning and Deep Learning and analytics for data at rest
- Metadata
- Real Time Analytics



Flash Memory Summit

What will Change

- 3DXP on the DRAM bus
- Metadata
- Real Time Analytics