



Flash Memory Summit



Persistent Memory for Artificial Intelligence

Bill Gervasi

Principal Systems Architect

bilge@Nantero.com



Flash Memory Summit

Memory for AI

**Artificial
Intelligence
Variations**

**Challenges for
Data Integrity**

Agenda

**New System
Architectures**

NVRAM and AI

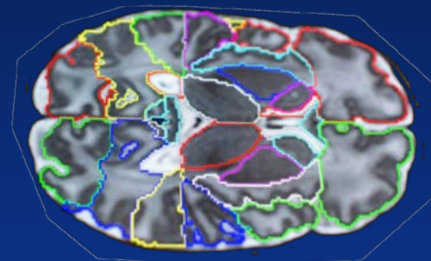


Flash Memory Summit

Every day...



...new applications for AI
emerge...



...from speech recognition to medical analysis to self driving
vehicles to purchasing tendencies to fraud detection to...





Flash Memory Summit

Stage 7 – Singularity and Transcendence

Stage 6 – Artificial SuperIntelligence

AI methods are getting smarter, too

Stage 5 – Self Aware Systems / Artificial General Intelligence

Stage 4 – Reasoning Machines

Stage 3 – Domain Specific Expertise

Stage 2 – Context Awareness and Retention

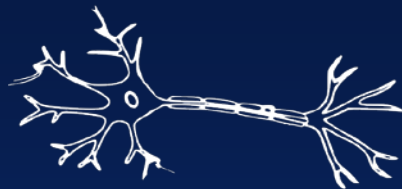
Stage 1 – Rule Based Systems

We're about mid-way into stage 4





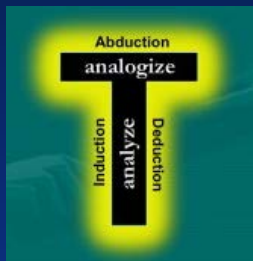
Flash Memory Summit



Neural



Fit

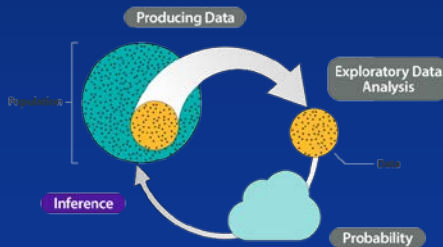


Symbolic

A wide variety of approaches



Analogy



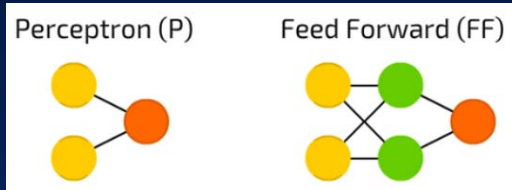
Inference

*Mueller/Massarón

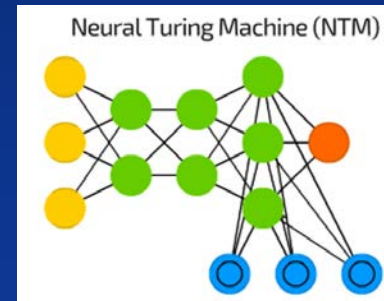
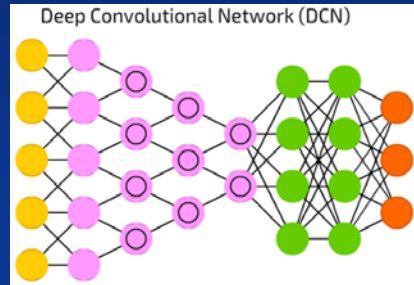
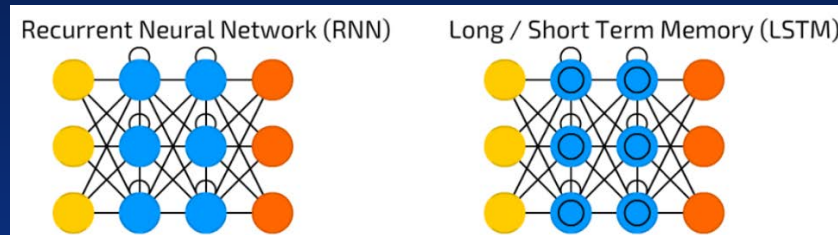


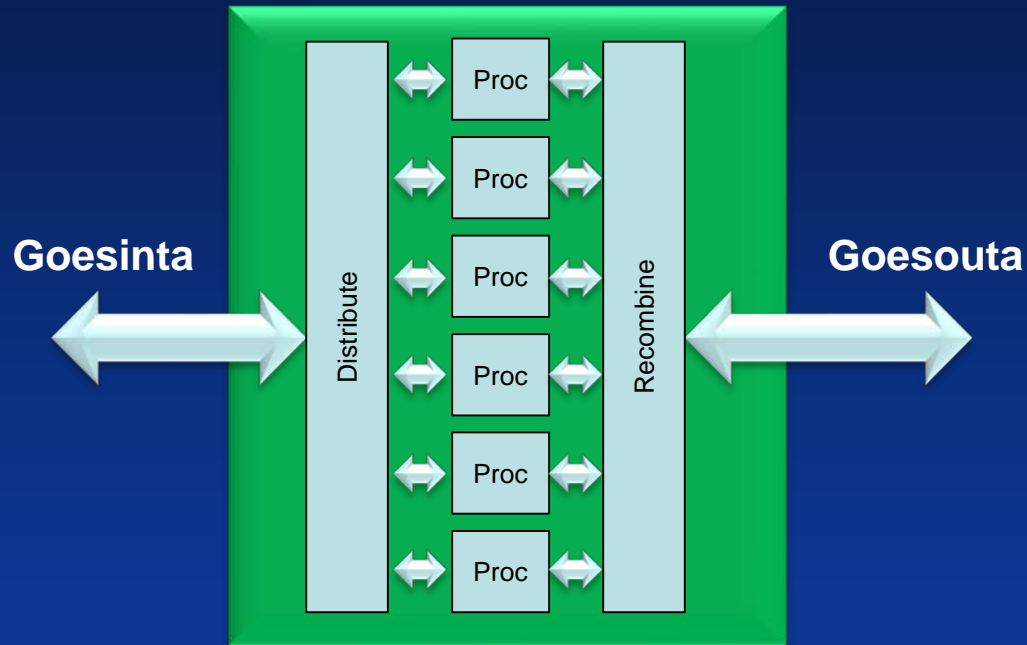


Flash Memory Summit



Each generation adds a new complexity in self modification





AI Accelerator Characteristics

SIMD (single instruction multiple data) rules

Wide array of simple processing elements

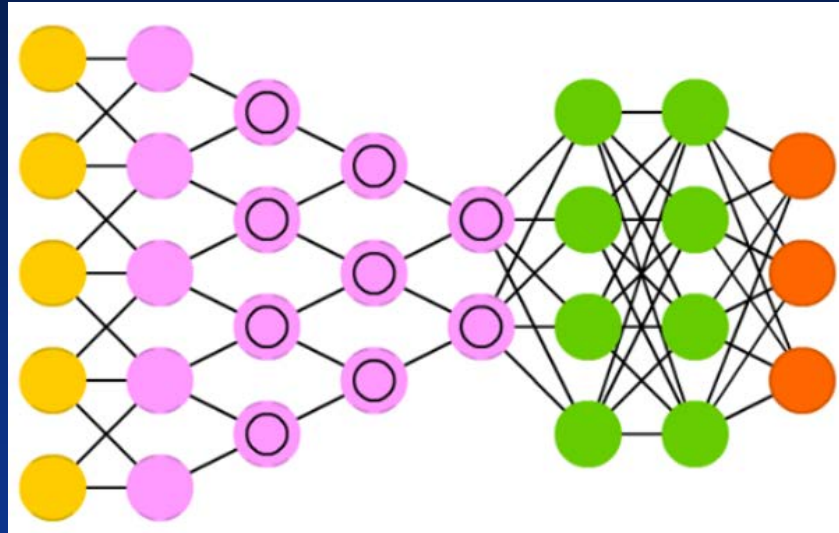
Reduced floating point precision

Tuned for matrix operations



Flash Memory Summit

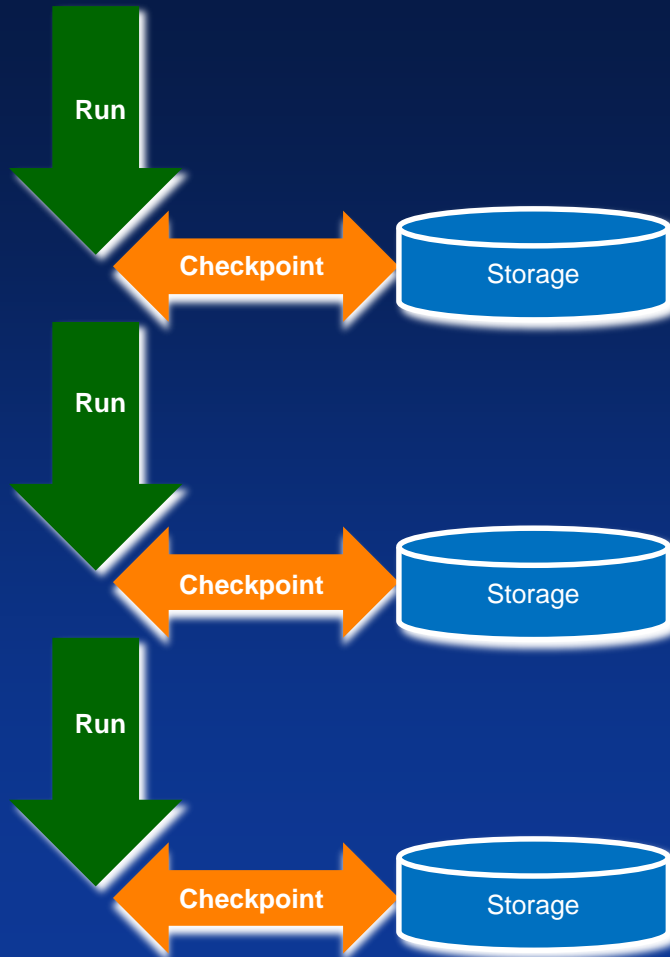
All that trapped data



To avoid losing learned data, systems must periodically checkpoint



Flash Memory Summit



Ah, checkpointing!

Great way to burn lots of power and waste performance for no good reason



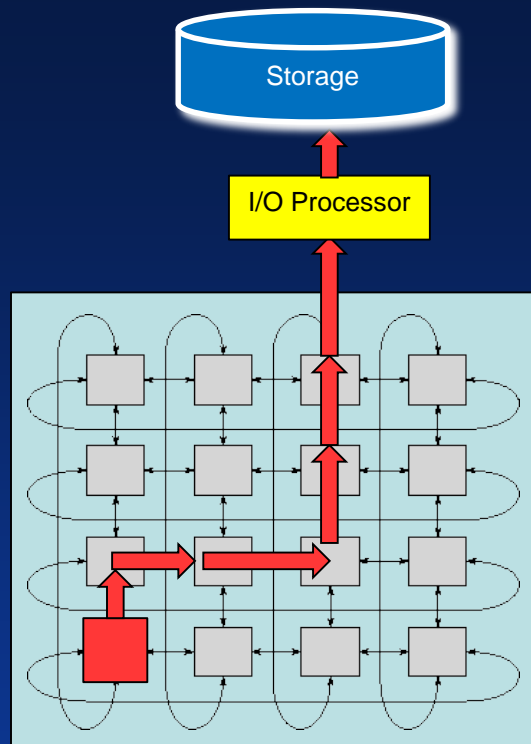
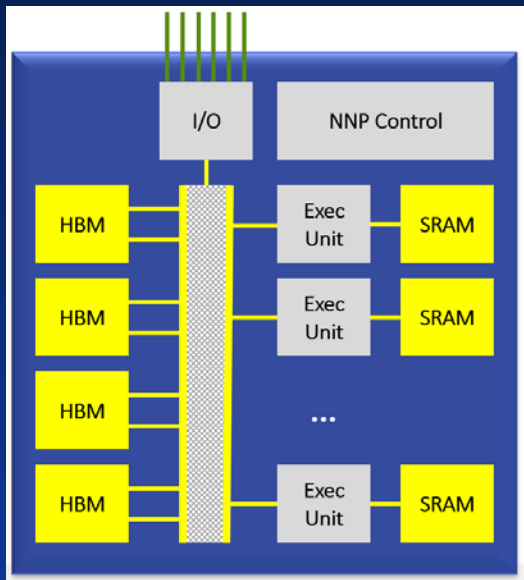
So now let's take our cute little AI engine and expand it into a mesh

Oh, look, some data I've learned I don't want to lose!

Let's save it!

See how much fun it is to checkpoint?

Wasn't that great?





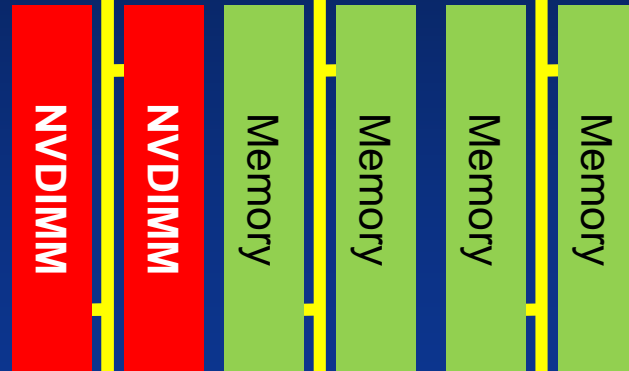
Flash Memory Summit

I/O



CPU \$

Memory Control



Trend is to move non-volatility closer to the CPU



Flash Memory Summit

Fortunately, a new wave of technology is coming

**Memory
Class
Storage**



Flash Memory Summit

Speed of a DRAM

Power-off data persistence

Unlimited write endurance

**Memory
Class
Storage**

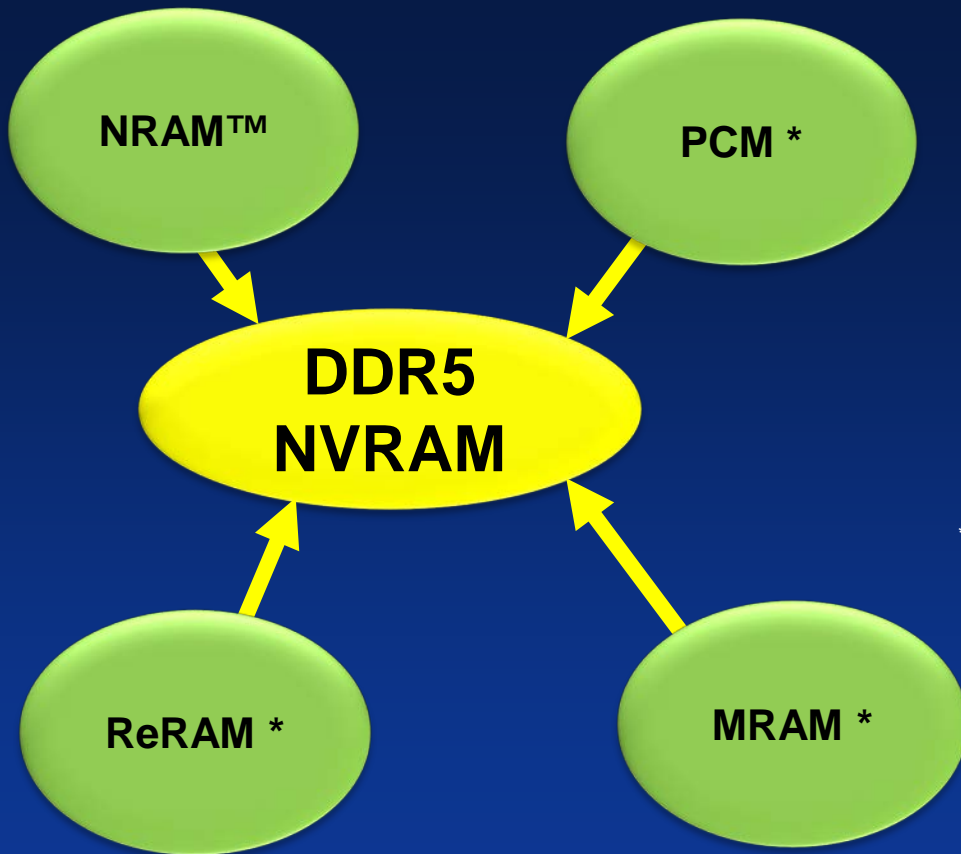


**PERSISTENT DRAM
DROP-IN
REPLACEMENT**



Flash Memory Summit

JEDEC STANDARD
Addendum No. 1 to JESD79-5, DDR5 Non-Volatile RAM (NVRAM)
JESD79-5-1
JC-42 Item #1856.10
September 2018
JEDEC SOLID STATE TECHNOLOGY ASSOCIATION
JEDEC



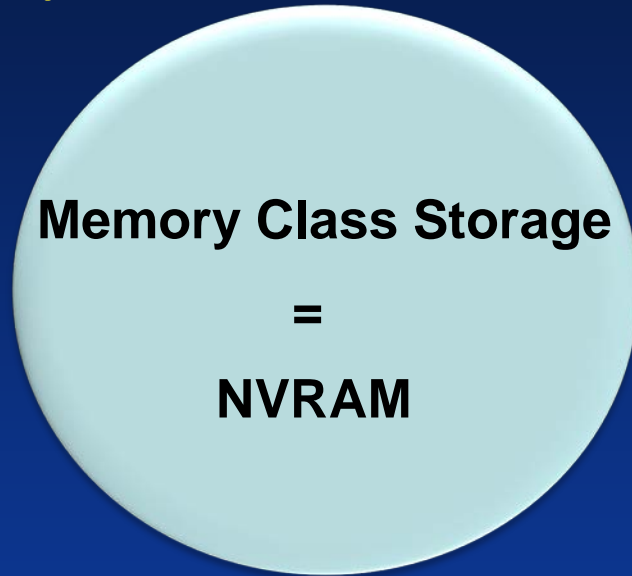
* Future generation devices

Details in my talk
on Thursday 8:30

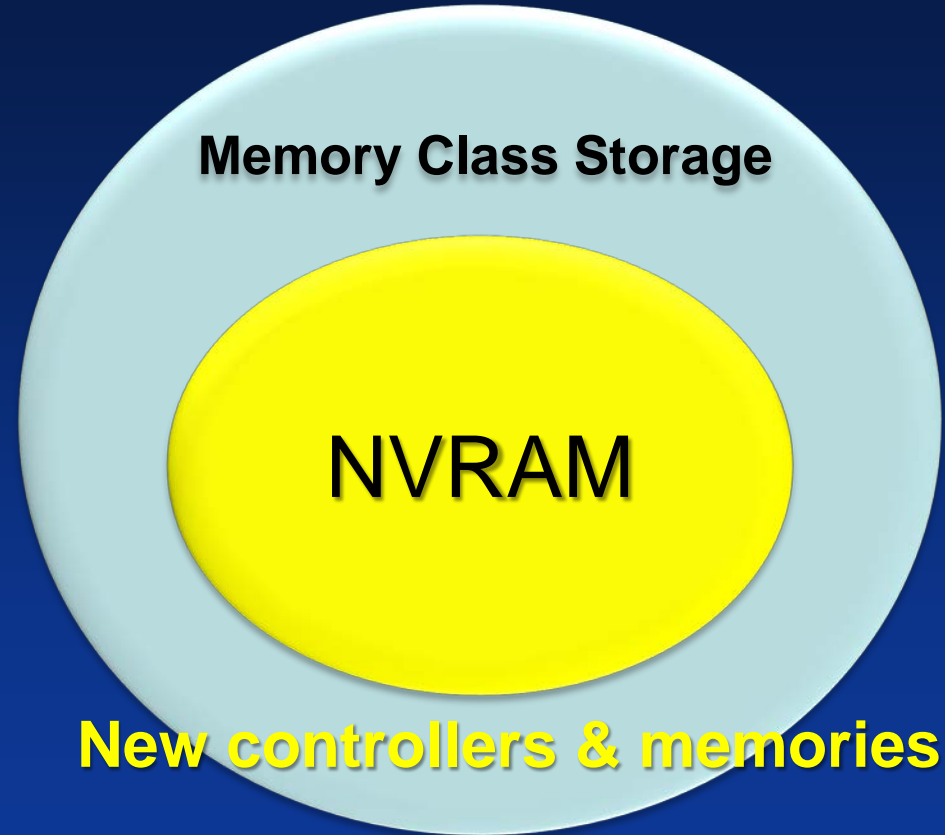


Flash Memory Summit

In the future



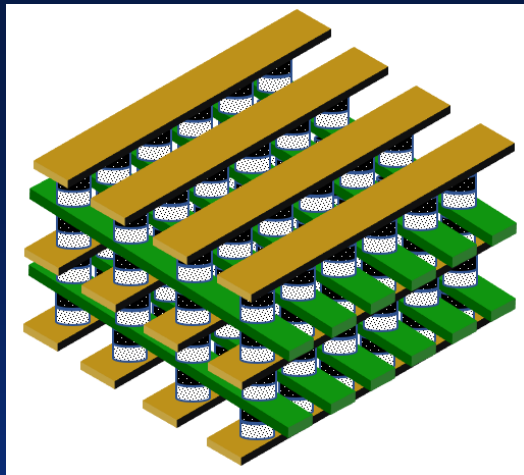
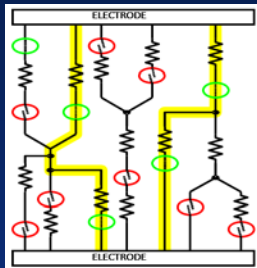
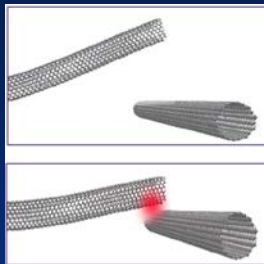
For now...



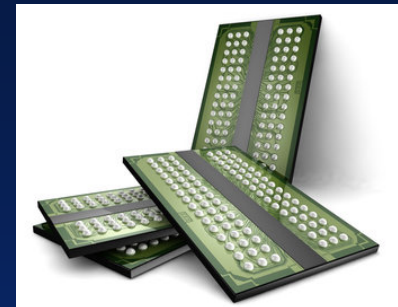
New controllers & memories



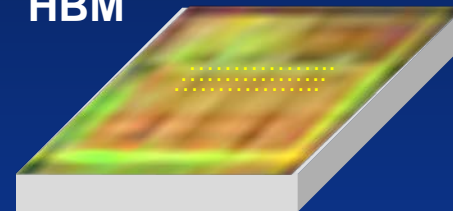
Flash Memory Summit



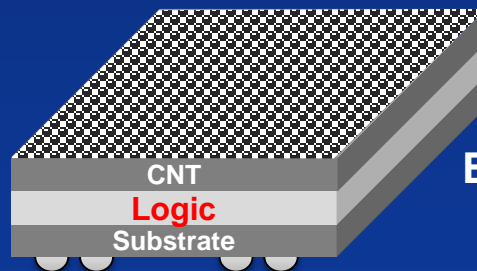
DDR4
DDR5



HBM



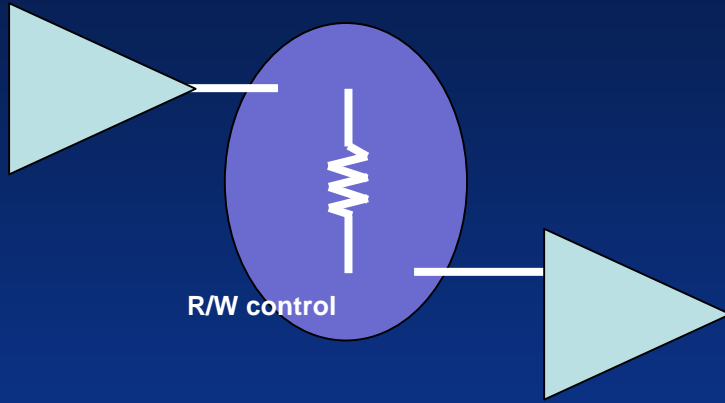
Nantero NRAM is a persistent memory using carbon nanotubes to build resistive arrays which can be arranged in a DRAM compatible device or deposited directly on circuits



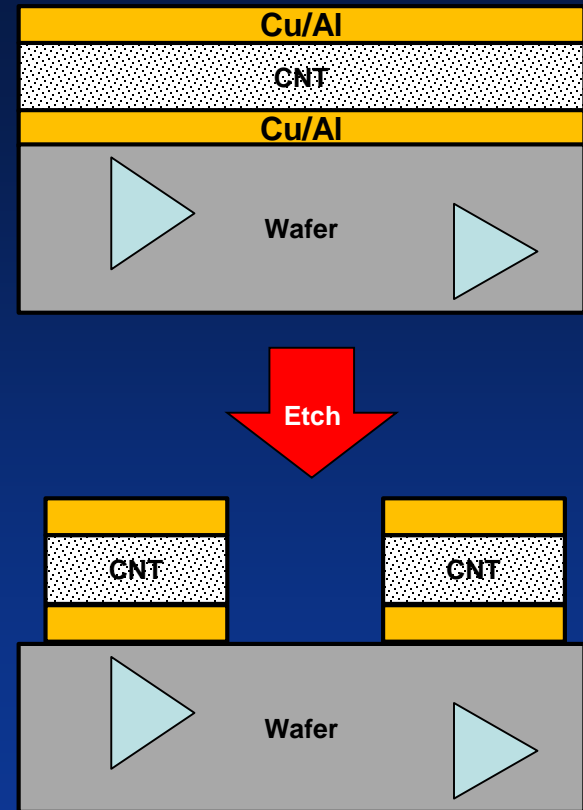
Embedded



Flash Memory Summit



NRAM may be applied on top of driver circuits, integrated into the logic





Flash Memory Summit

I/O

Storage-less systems enabled



CPU \$

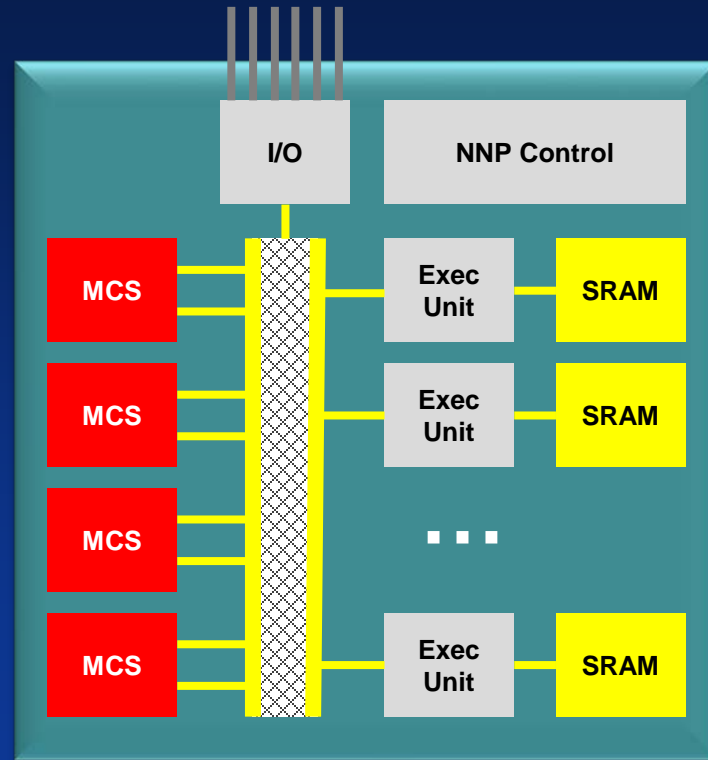
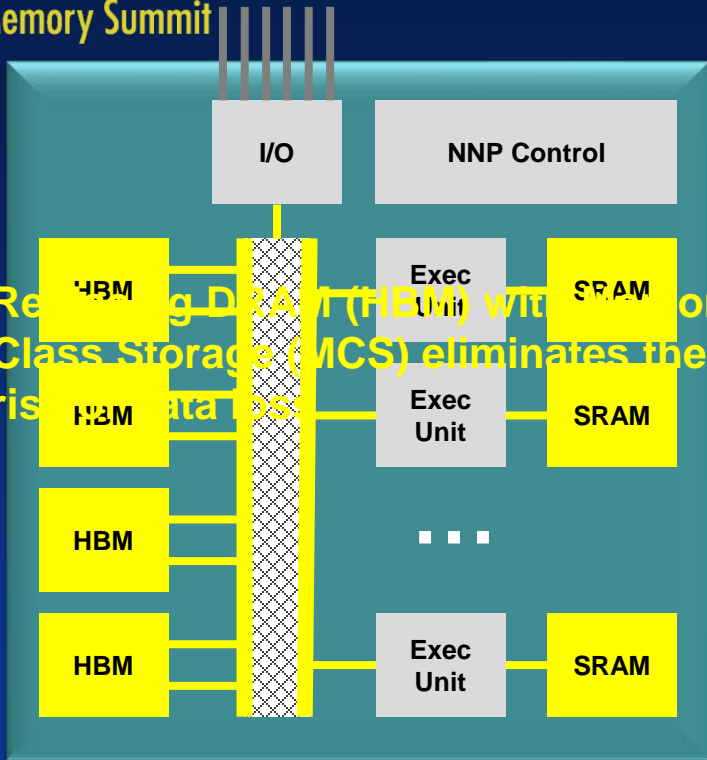
Memory Control

- Memory Class Storage
- Memory Class Storage
- Memory Class Storage
- Memory Class Storage
- Memory Class Storage
- Memory Class Storage



Flash Memory Summit

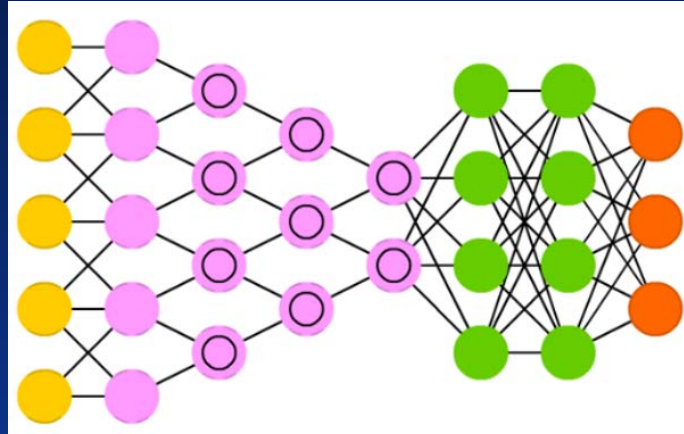
Replacing DRAM (HBM) with Memory Class Storage (MCS) eliminates the risk of data loss





Flash Memory Summit

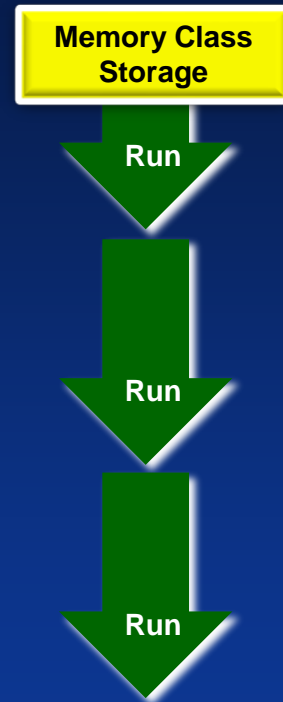
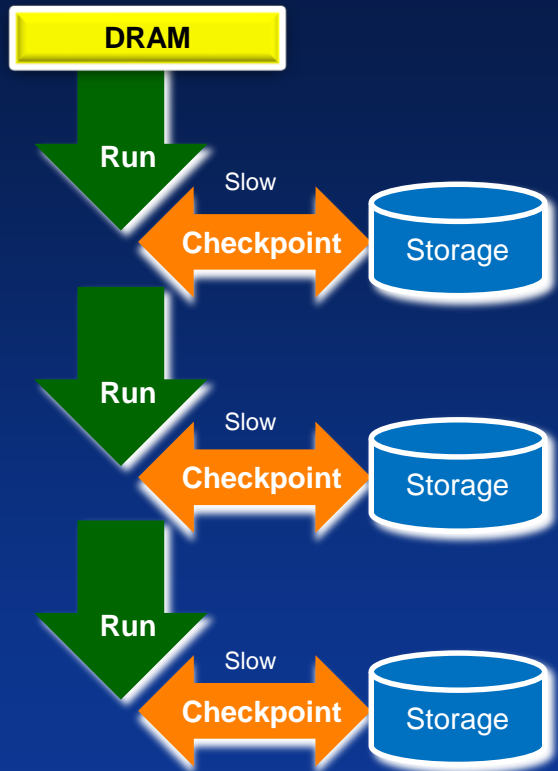
All that trapped data...



...can just stay there!



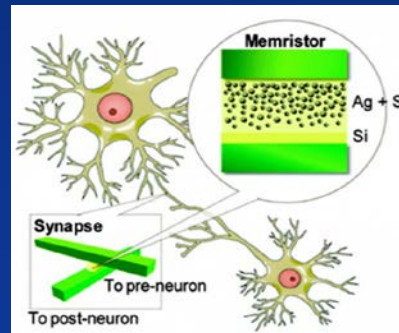
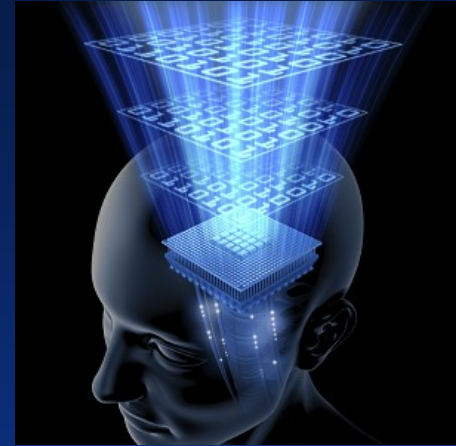
Flash Memory Summit



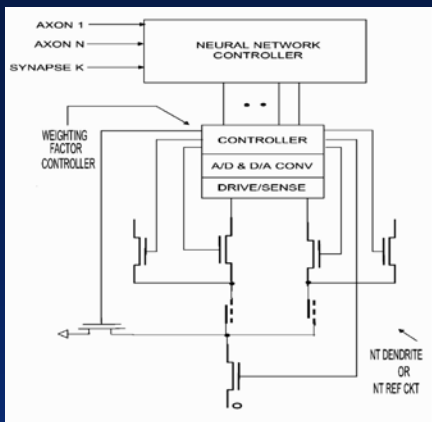


Flash Memory Summit

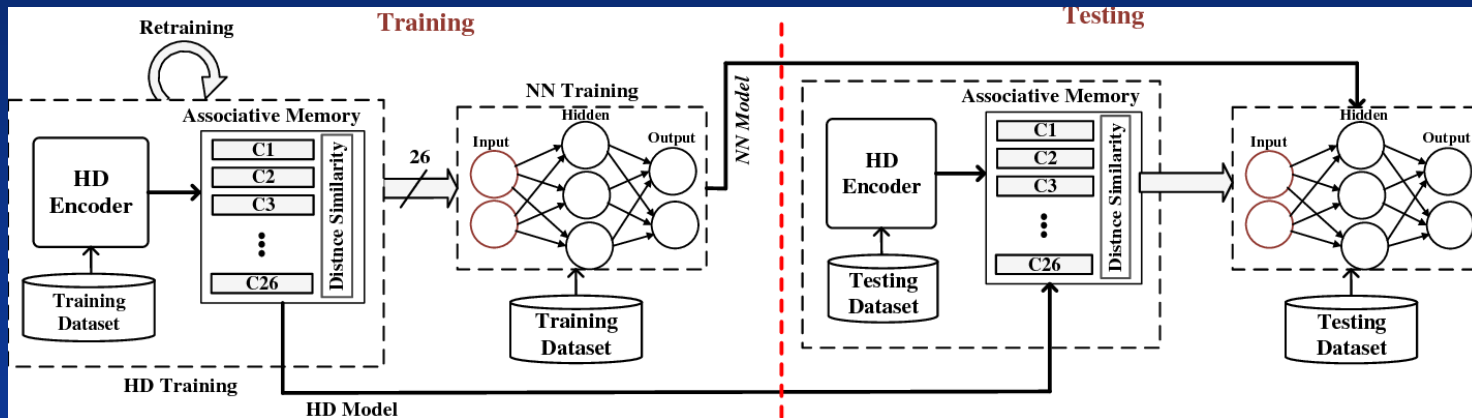
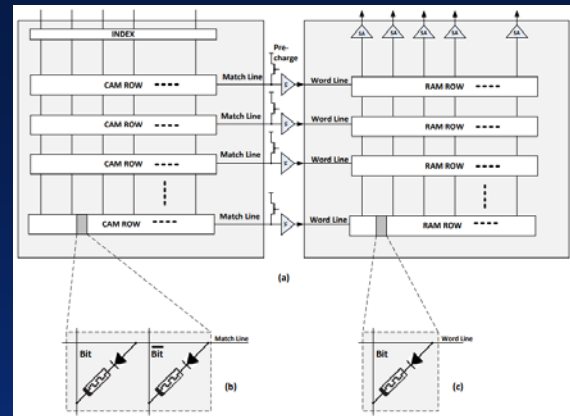
There is a HUGE gap between
AI research work and what's
being built

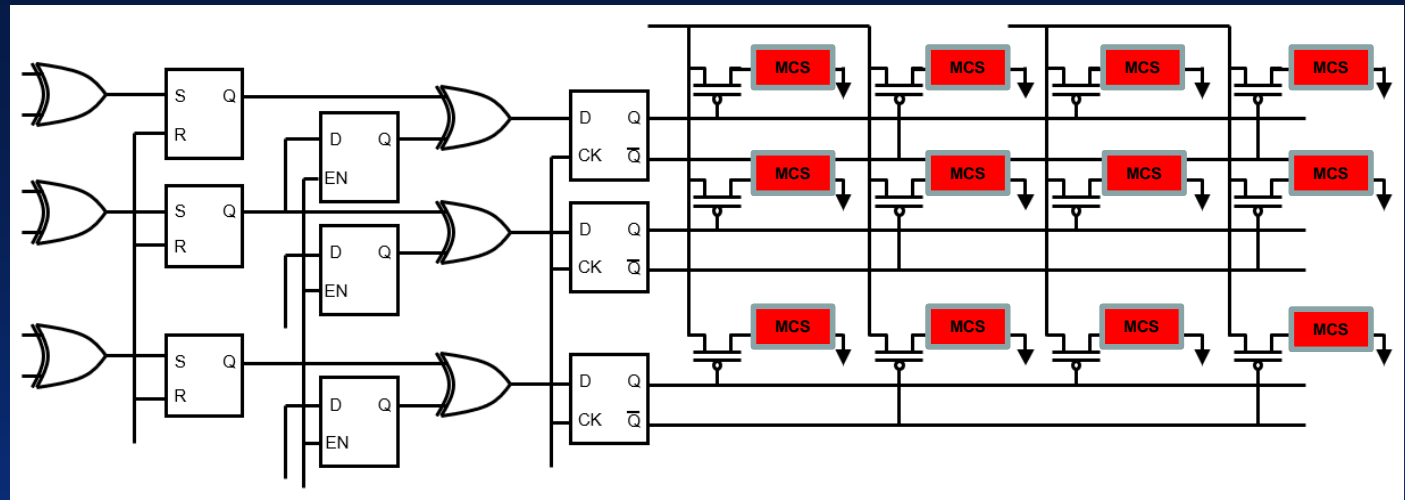


Modeling the brain
leads to embedded
memory



AI embedded memories are often not von Neumann





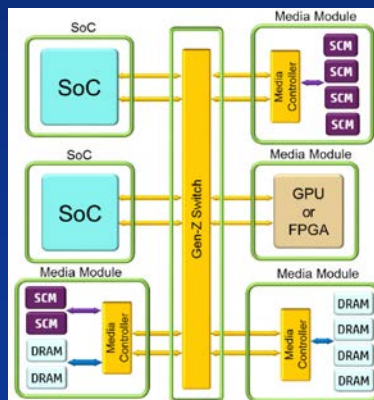
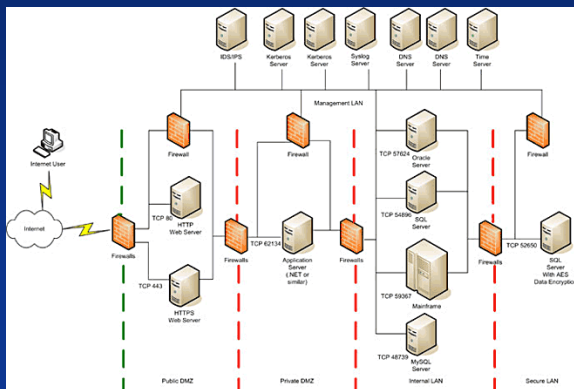
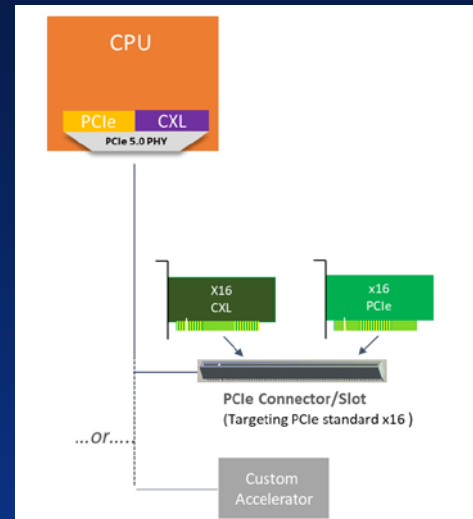
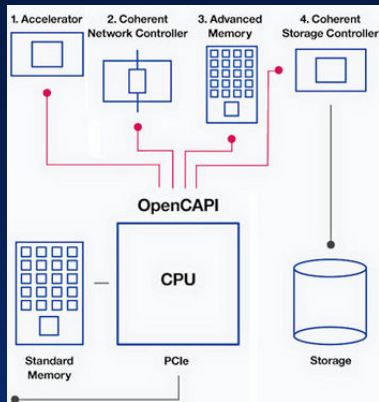
Embedded memories suffer from read/write time issues when used as random access (including checkpointing)

Replacing embedded memory with Memory Class Storage removes this barrier



Flash Memory Summit

Computing networks are evolving to fabrics

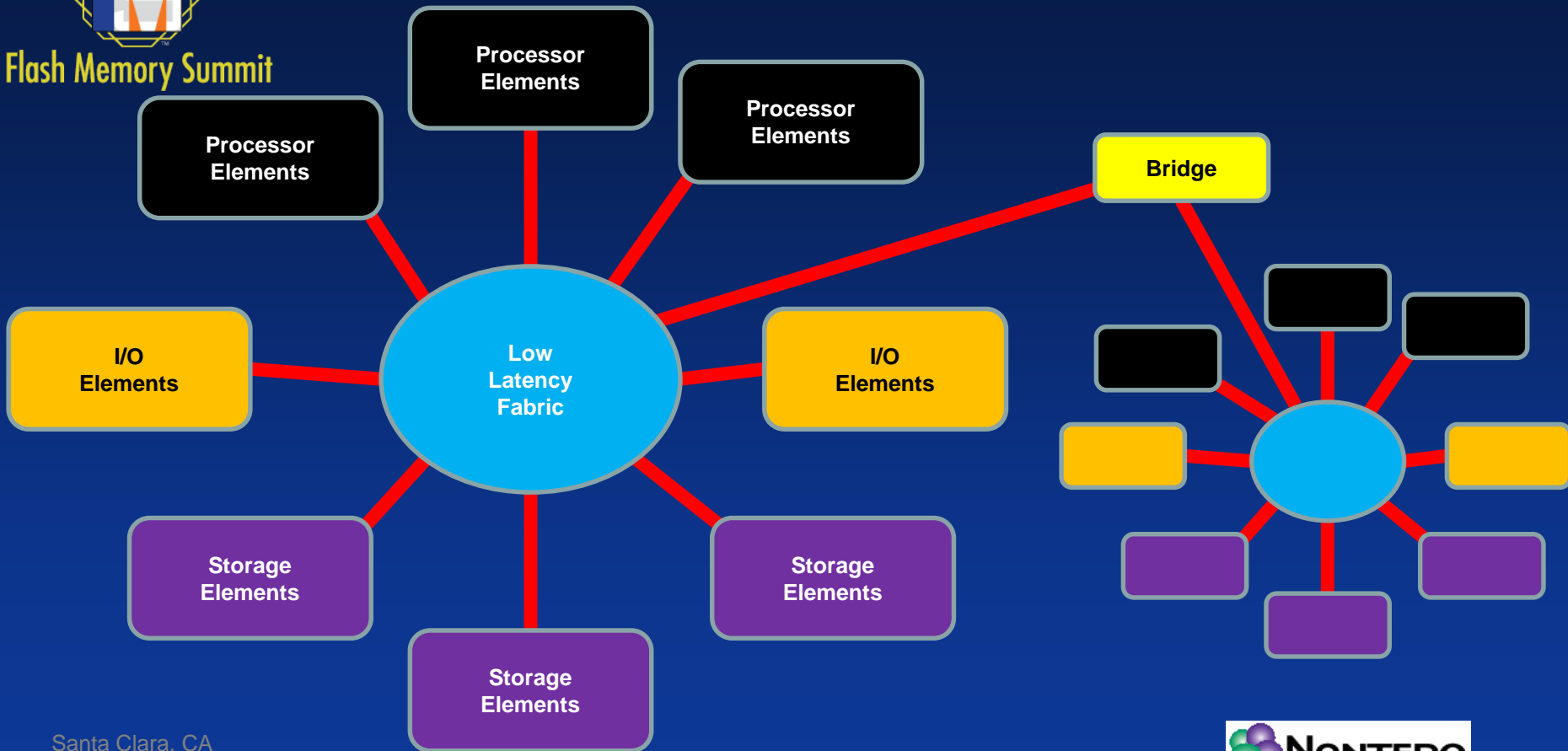


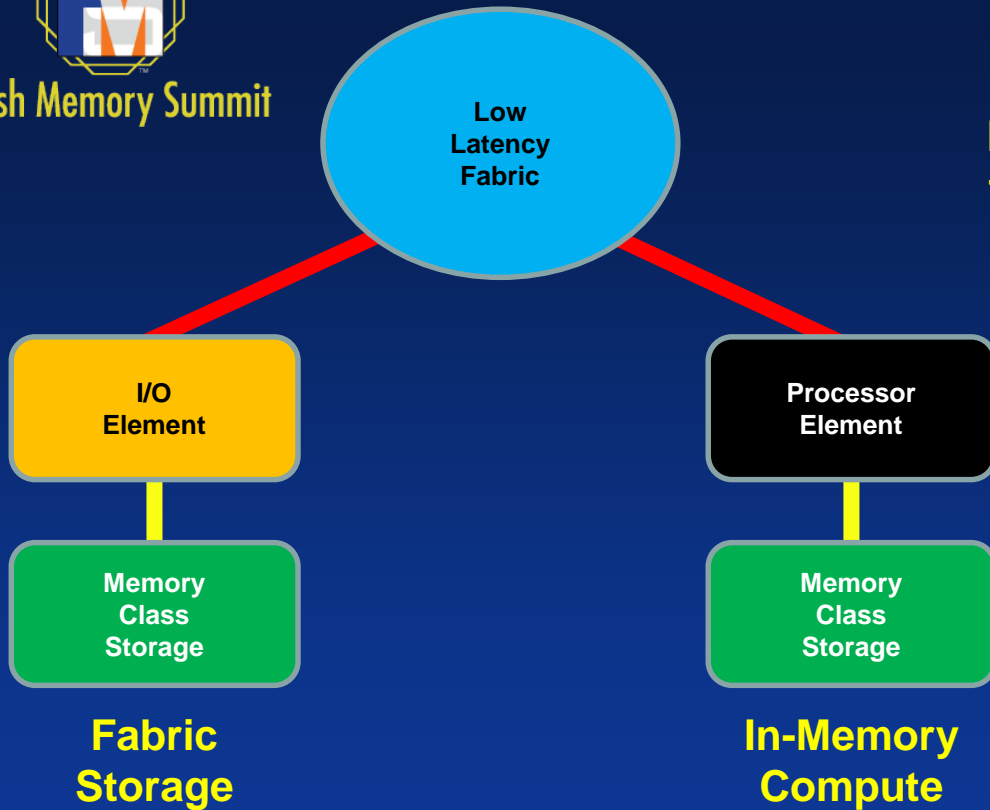
Heterogeneous processing in a unified space





Flash Memory Summit





DDR DRAM protocols do nothing for these architectures but get in the way

Likely place for new Memory Class Storage interface protocols to emerge



Flash Memory Summit

Number of Devices Needed to Achieve Gen-Z Throughput

Gen-Z

DDR4-3200 x8 DDR5-6400 x8 LPDDR5-6400 x64 GDDR6-18000 x32 HBM3-4800 x1024

GT/s	Lanes	GB/s	3.2 GB/s	6.4 GB/s	51.2 GB/s	72 GB/s	615 GB/s
25	64	320	100x	50x	7x	5x	1x
25	128	640	200x	100x	13x	9x	2x
32	64	400	125x	63x	8x	6x	1x
32	128	800	250x	125x	16x	12x	2x
56	64	700	219x	110x	14x	10x	2x
56	128	1400	219x	110x	28x	20x	3x
112	64	1400	438x	220x	28x	20x	3x
112	128	2800	875x	438x	55x	39x	5x





Flash Memory Summit

**Number of
Devices
Needed**

**High System
Power**

Problems with Standalone Memories

**Little Control
Over # GB**

**Lots of
Wasted Data
Access**



Flash Memory Summit

Gen-Z

GT/s	Lanes	GB/s
25	64	320
25	128	640
32	64	400
32	128	800
56	64	700
56	128	1400
112	64	1400
112	128	2800

NRAM-6400 x4096

3200 GB/s*

1x

1x

1x

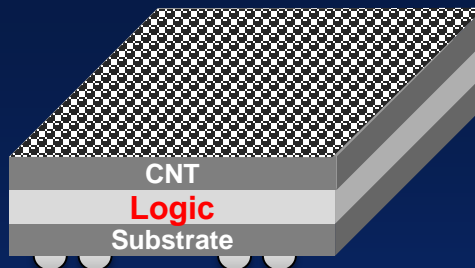
1x

1x

1x

1x

1x



Single chip fabric storage possible

In-Memory computing enabled

Significantly lower power

Access granularity up to the controller:

Avoid unused data accesses



Flash Memory Summit

**Embedded
Persistent Memory
changes how you
view AI architecture**



**Embedded
Persistent Memory
enables new
implementations**

Santa Clara, CA
August 2019





Flash Memory Summit

**Common
Memory Pools
Used**

**AI Uses SIMD
Data Flow**

**Data
Checkpointing is
Horrid for AI**

Agenda

**New Applications
for Fabric Based
Systems Coming**

**NVRAM Solves
Many AI
Memory Issues**



Flash Memory Summit

Questions?

Bill Gervasi

Principal Systems Architect

bilge@Nantero.com