

Accelerating Time to Insight

Performance optimized emerging
system architectures

Balint Fleischer

Senior Director

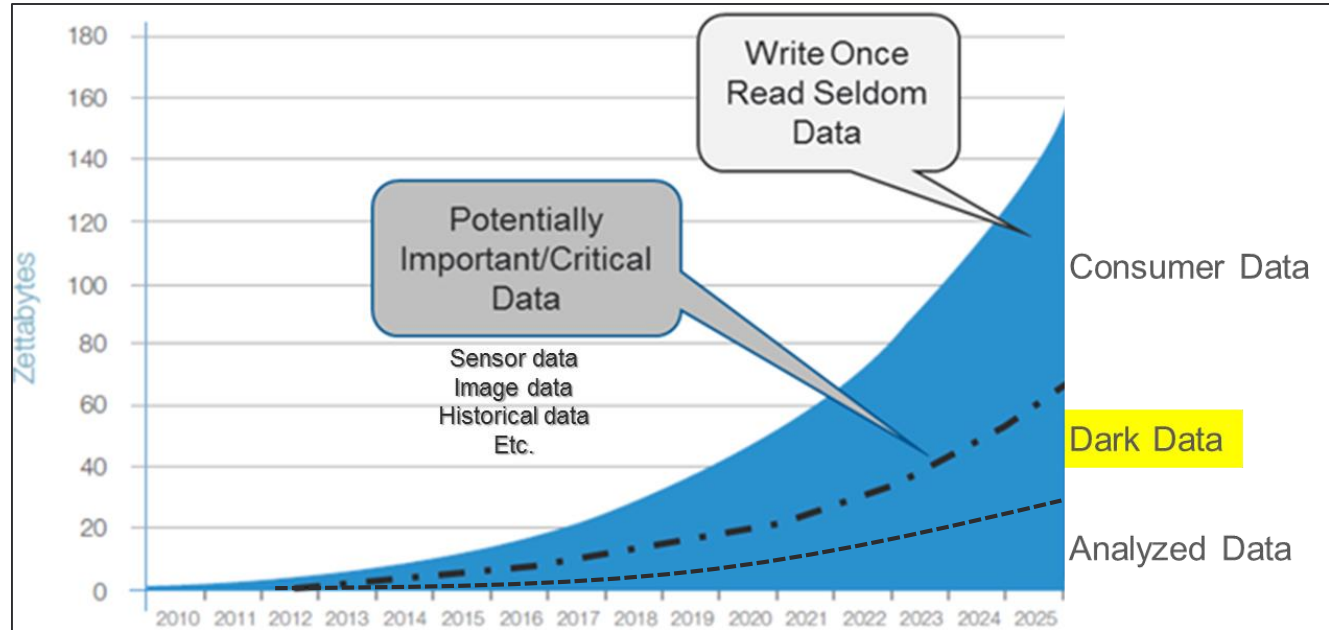
Advanced Computing Solutions

September 25th, 2019

©2018 Micron Technology, Inc. All rights reserved. Information, products, and/or specifications are subject to change without notice. All information is provided on an "AS IS" basis without warranties of any kind. Statements regarding products, including regarding their features, availability, functionality, or compatibility, are provided for informational purposes only and do not modify the warranty, if any, applicable to any product. Drawings may not be to scale. Micron, the Micron logo, and all other Micron trademarks are the property of Micron Technology, Inc. All other trademarks are the property of their respective owners.



Rapidly growing *Dark Data*



Rapidly growing Data sets from many sources

+

Increasing complex algorithms

+

Need for faster Time to Insight

+

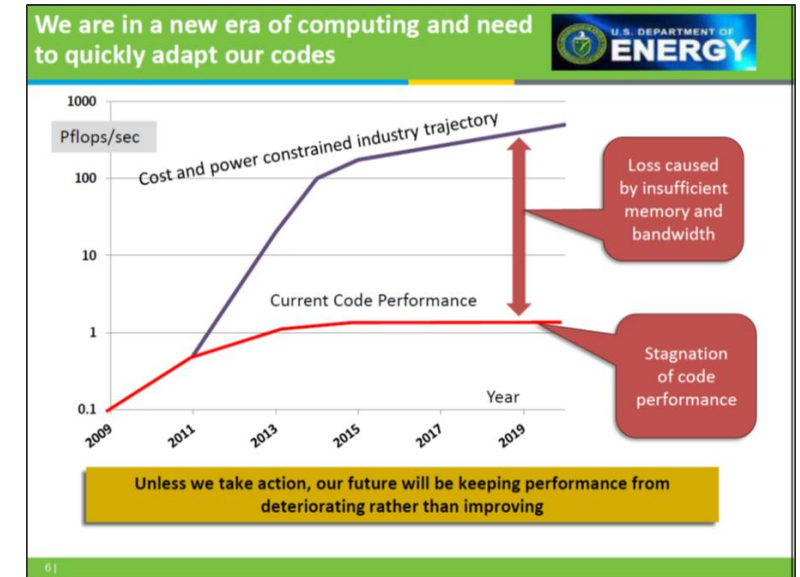
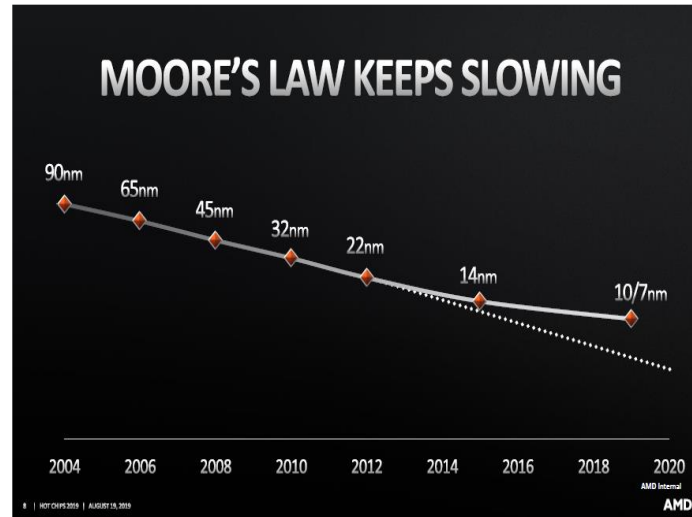
Affordability challenges



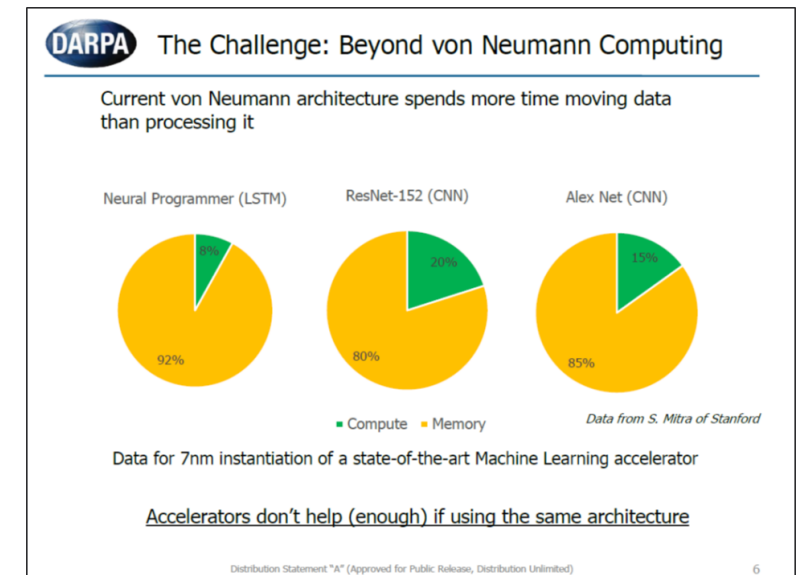
Dark Data

Estimated data
1 Zettabyte = 10^{21} (one billion, trillion) bytes

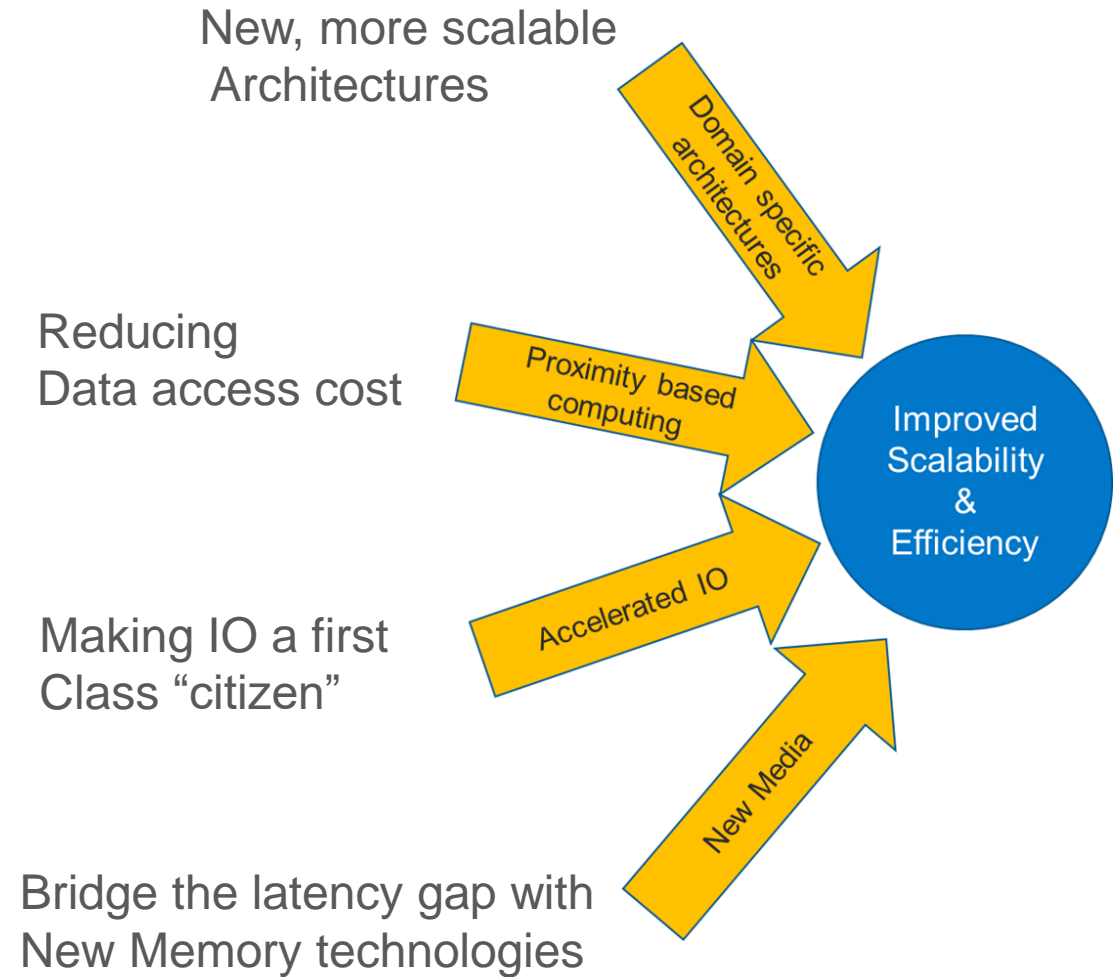
Perfect storm:
Moore's law is slowing,
Dennard scaling is ending and von Neumann architecture became a bottleneck



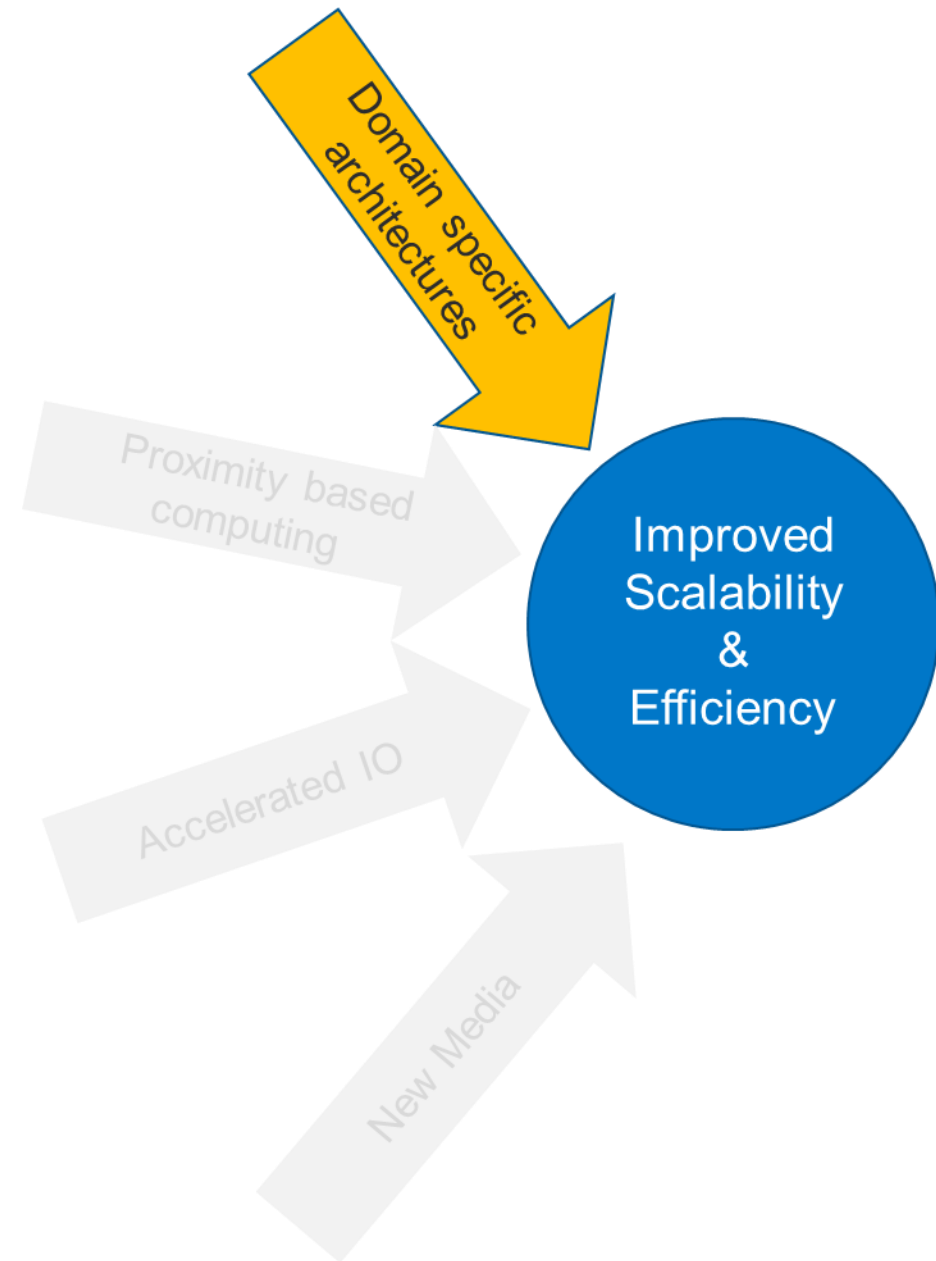
B. Meisner, The Bump in the road to ExaFlops and Rethinking LINPACK, HPC User Forum, June 2014



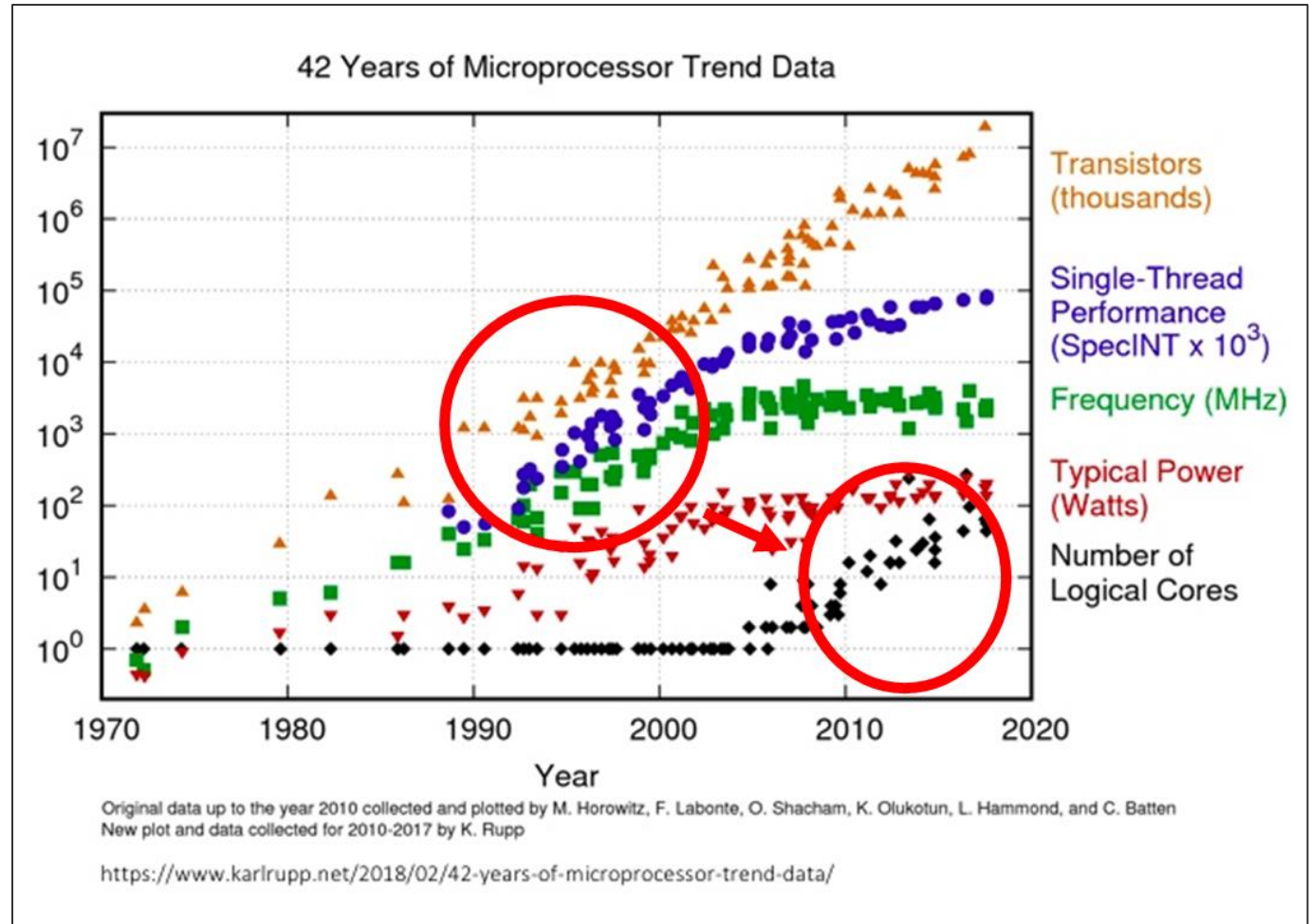
**Large,
industry wide
approach but
without a
master plan
to insure
continued
scaling**



**New, Non Von
Neumann
Architectures
Scale well (for
a while) with
process
technology
improvements**



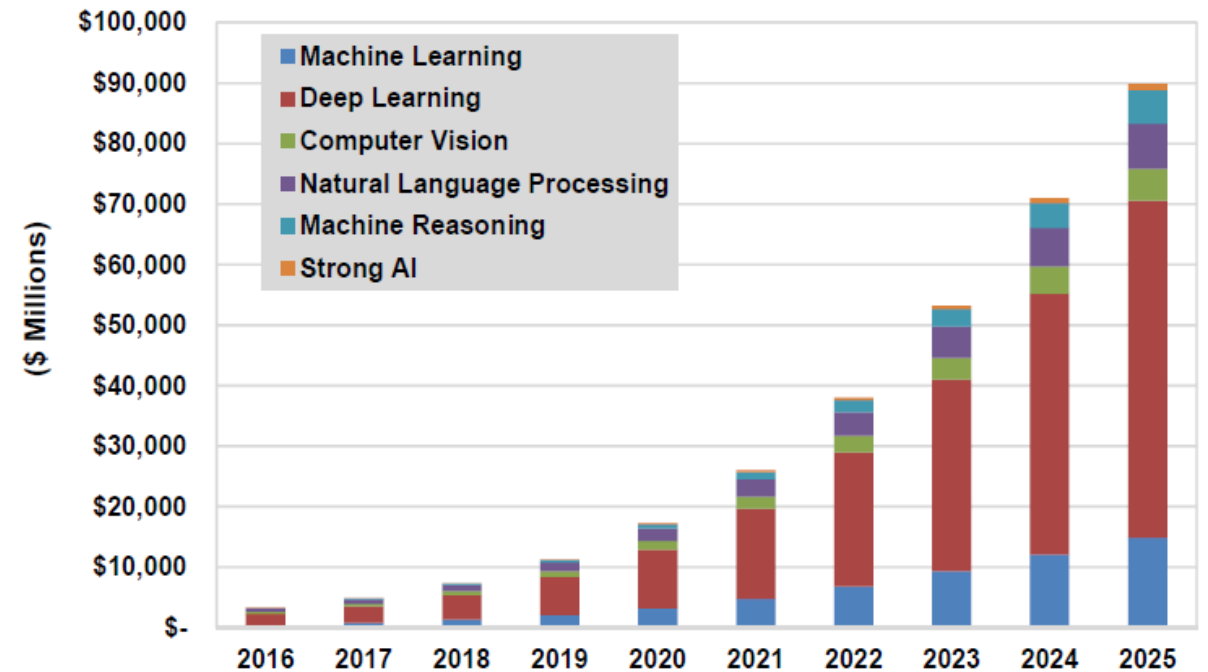
General
purpose CPU
performance
CAGR is
declining



Machine Learning + Deep Learning emerging as the key application for data analytics

A rich target for
Domain specific
computing

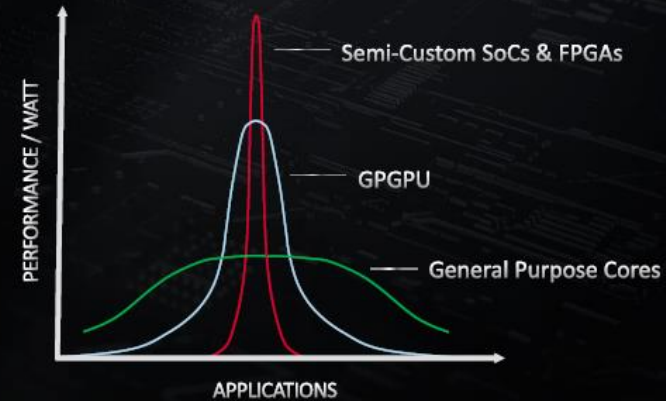
Chart 3.5 Annual Artificial Intelligence Revenue by Technology, World Markets: 2016-2025



(Source: Tractica)

On a given technology node, Domain Specific Architectures deliver greater performance on targeted applications

OPTIMIZING SYSTEM PERFORMANCE WITH ACCELERATED COMPUTING



23 | HOT CHIPS 2019 | AUGUST 19, 2019

CHART FOR ILLUSTRATIVE PURPOSES

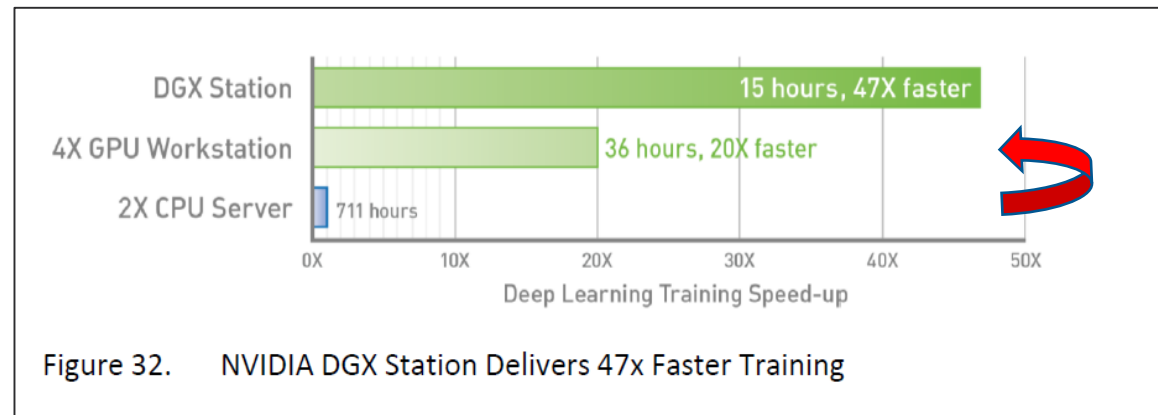
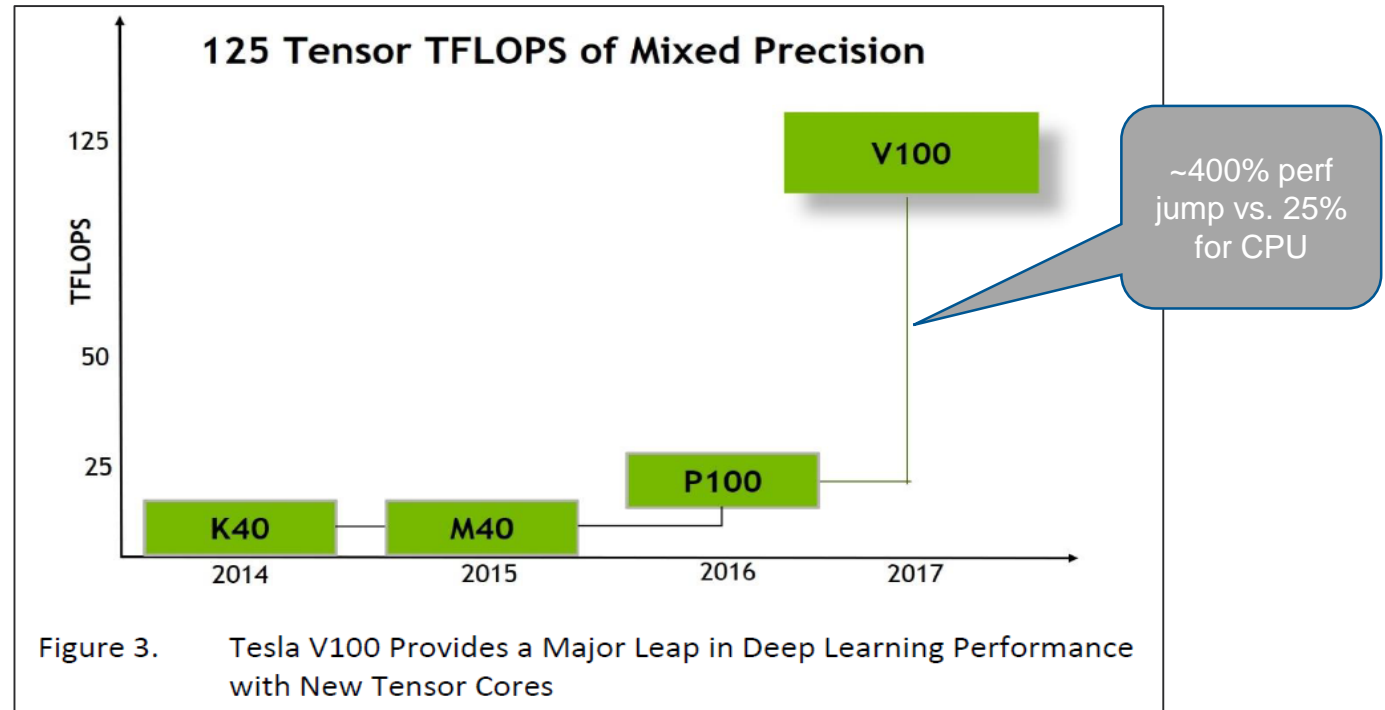
AMD

Pattern for domain specific architectures

- Simpler, lower performance, more efficient processing elements
- High degree of parallelism (1000s of processing elements)
- Highly optimized on die data movement
- New, application specific ISA
- Custom compiler

An example of
Domain Specific
Architecture
delivering greater
performance

Optimized AI
processors deliver
even better gain



There are over 300
AI processor
designs worldwide
innovating from
Low End to High
End



The Hailo-8 achieves 26 TOPS (Image: Hailo)

EETimes 08.29.10 Details of Hailo AI Edge Accelerator Emerge

Claimed performance
2.8 TOPS @ 1.6W
~2x faster vs. CPU
1/15th power vs. CPU

Cerebras Wafer Scale Engine

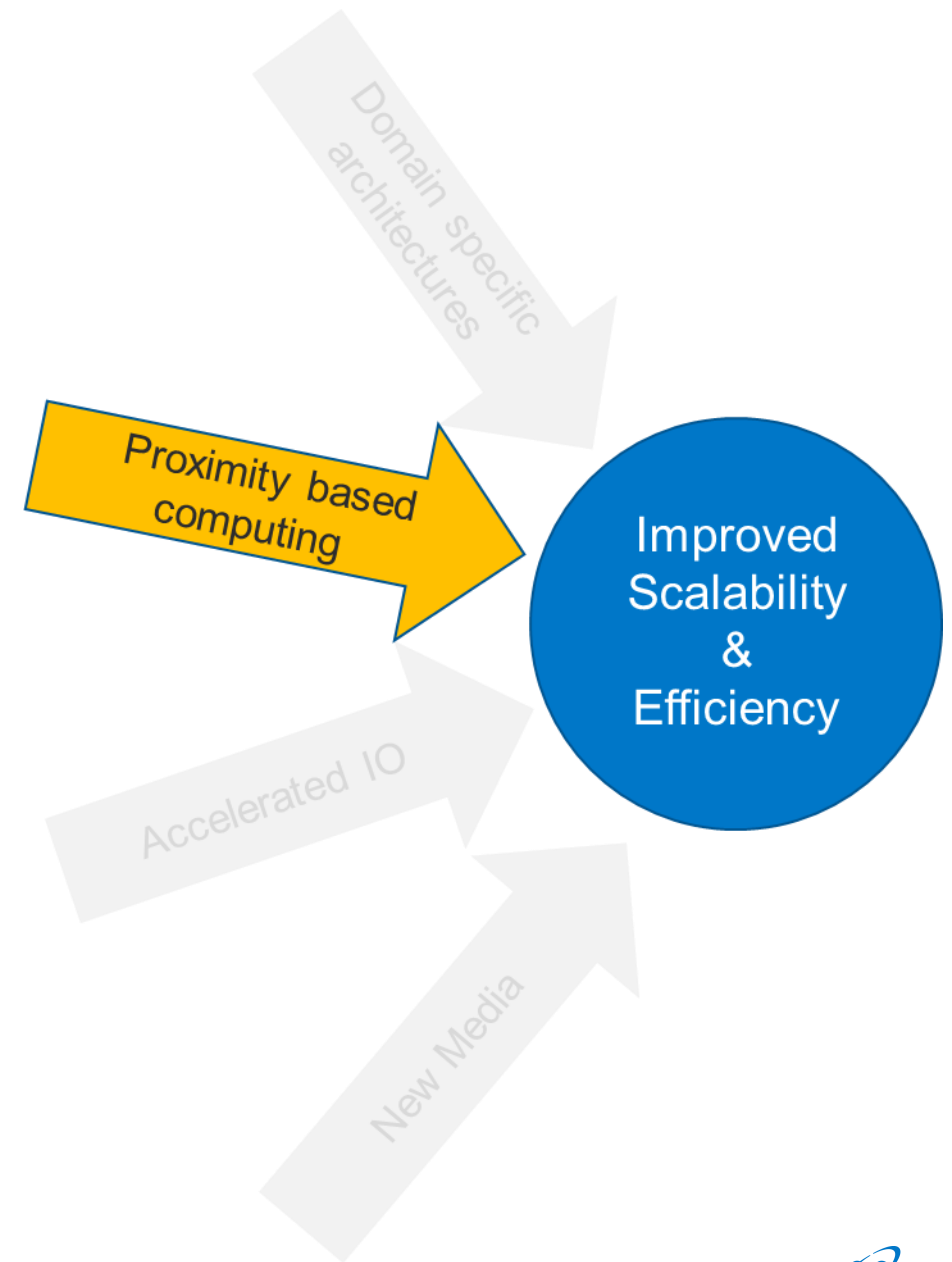
Cerebras WSE
1.2 Trillion Transistors
46,225 mm² Silicon

Largest GPU
21.1 Billion Transistors
815 mm² Silicon

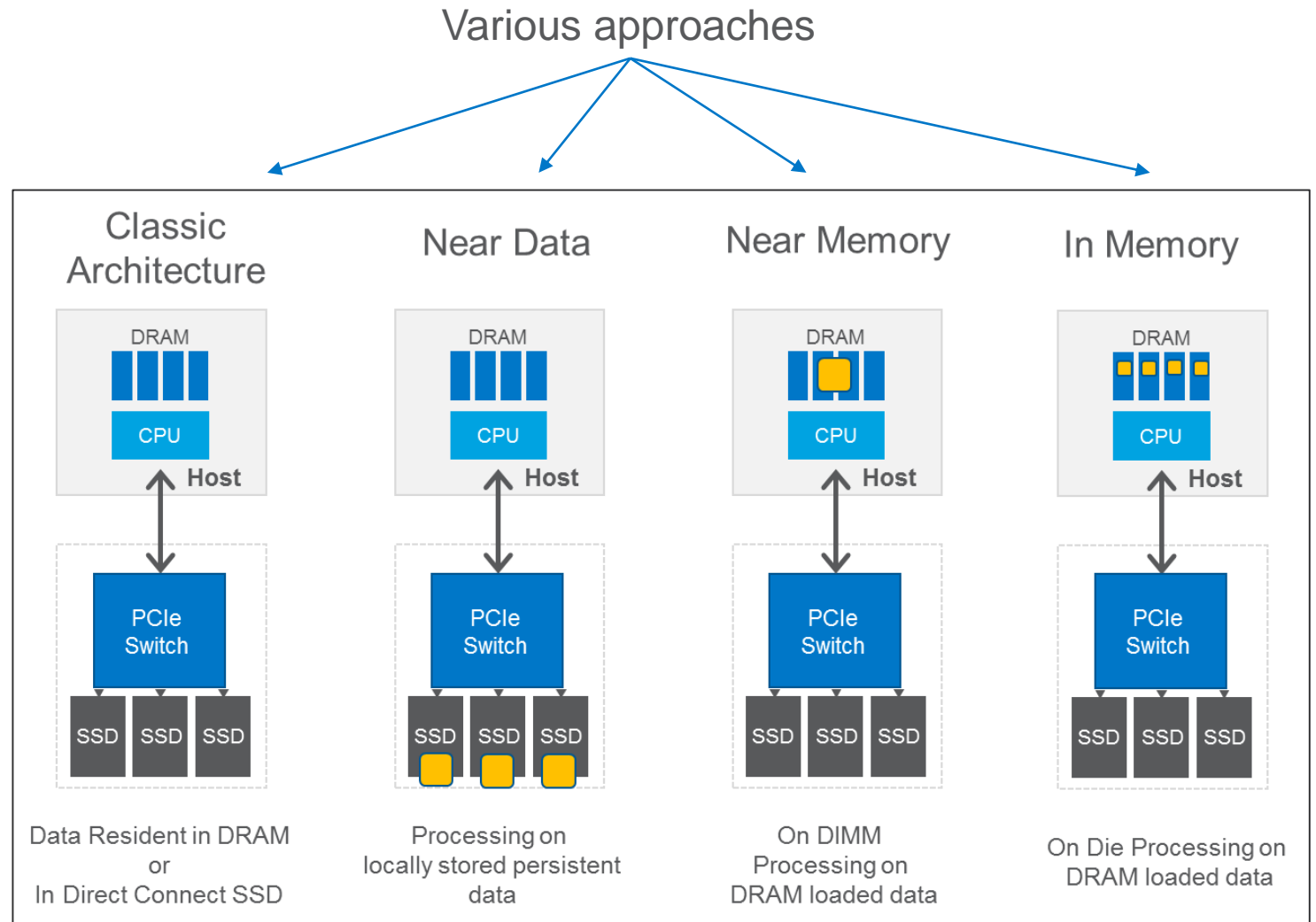
Hot Chips 2019, Cerebras Presentation

“Cost” of IO is critical to performance scaling and energy consumption

The goal is improving Bandwidth, reducing Latency & reducing Data movement

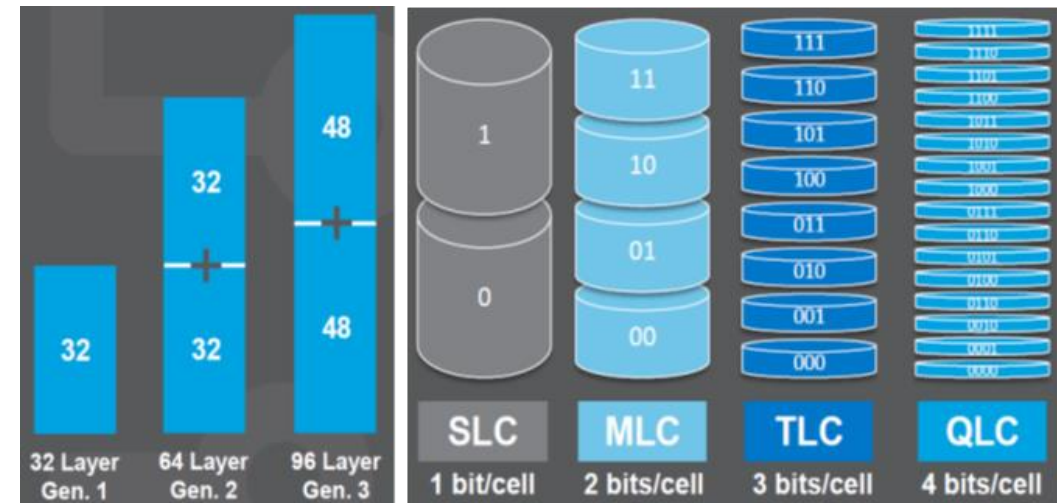
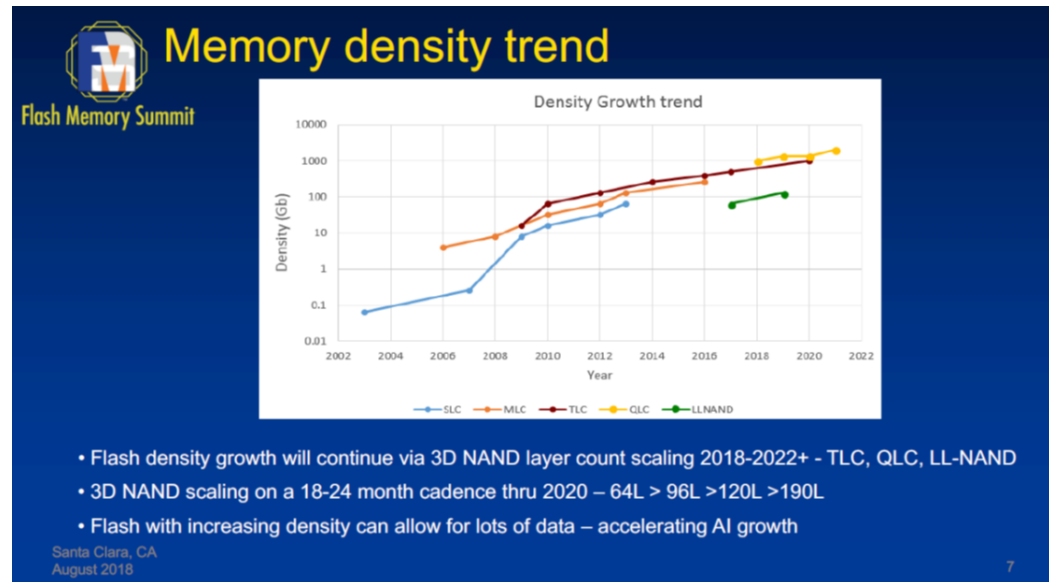


Improving Data Proximity is a vehicle to address IO cost



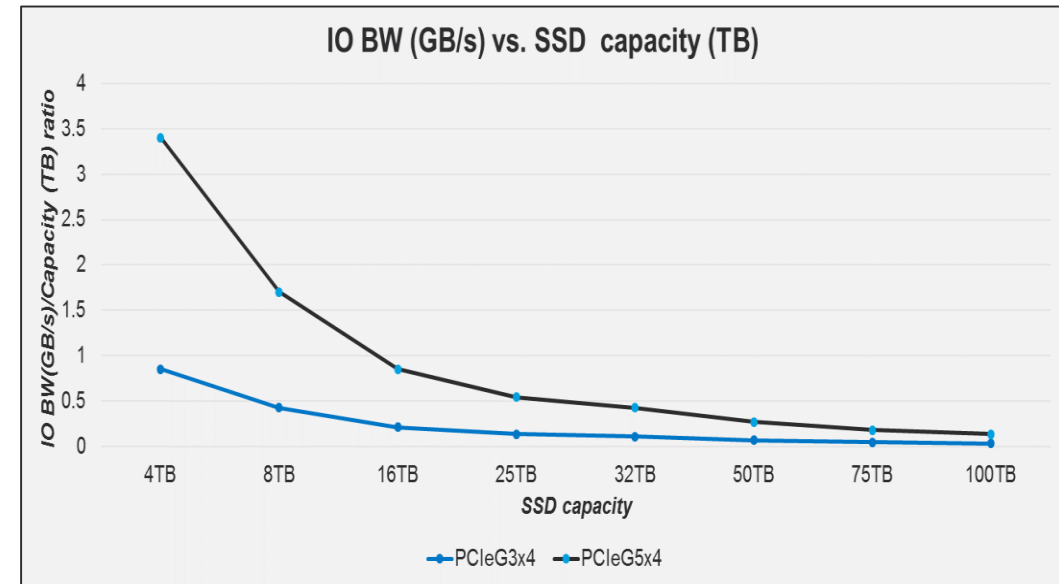
Processing Element

NAND die capacity continues to grow enabling more dense SSDs and storage solutions



Difficult to Beat the NAND Cost Structure

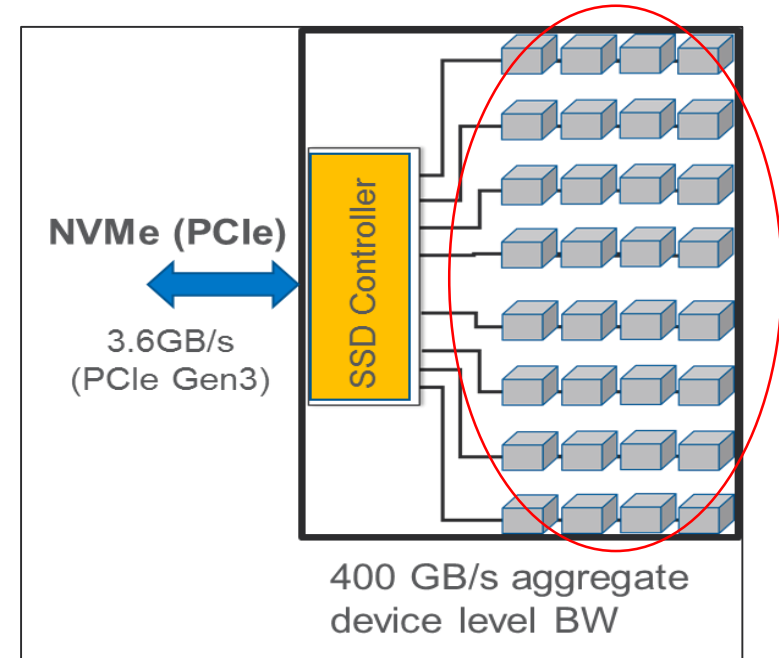
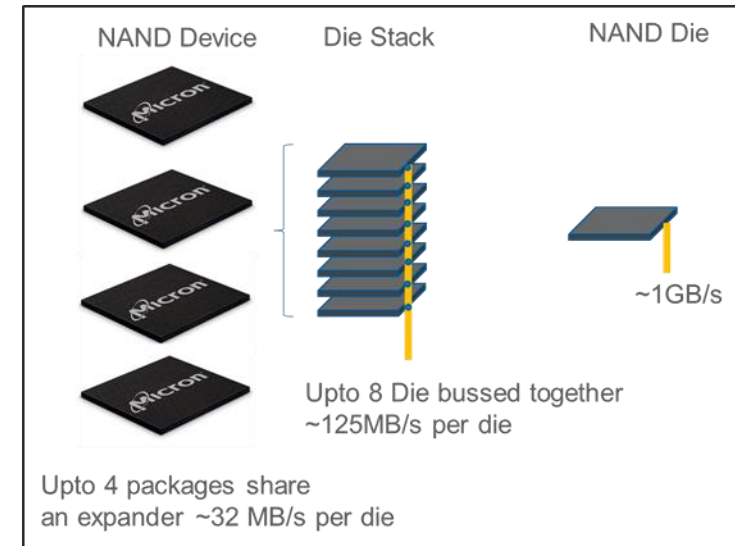
SSD IO Bandwidth is lagging SSD Capacity growth



Query of very large data sets
Will take an increasingly long time

NAND stacks and SSDs are designed for capacity scaling

Device level aggregate
BW is ~100x vs. SSD
External IO Bandwidth
Die level Bandwidth is
~1000x

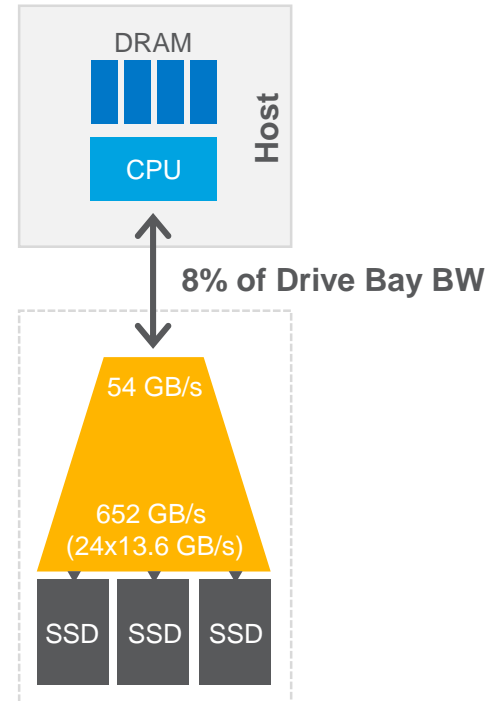


Moving processing into SSD can benefit from the high Embedded Bandwidth

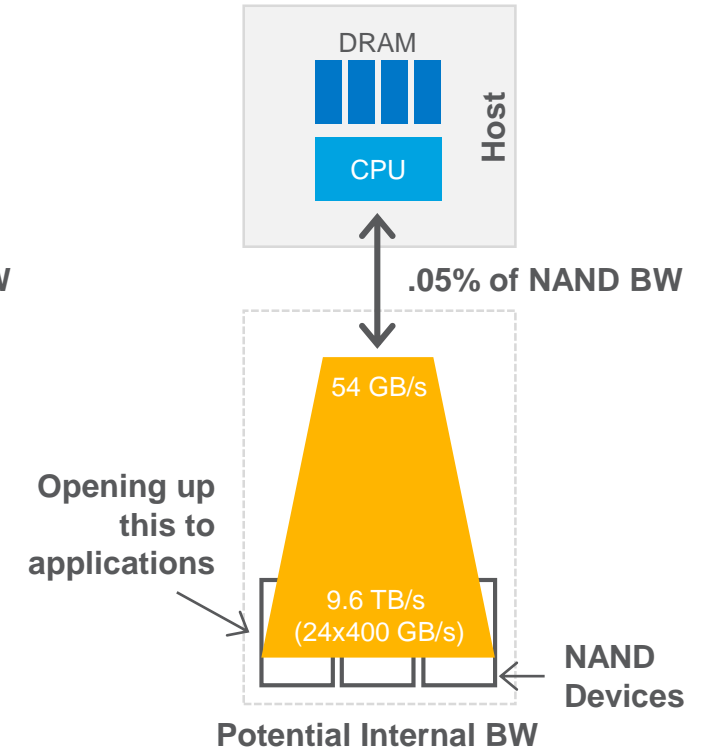
Time to Insight on very large shardable data will also benefit from Massive parallelism

* For simplicity single socket CPU Hosts are shown

652 GB/s aggregate IO BW at SSD connector



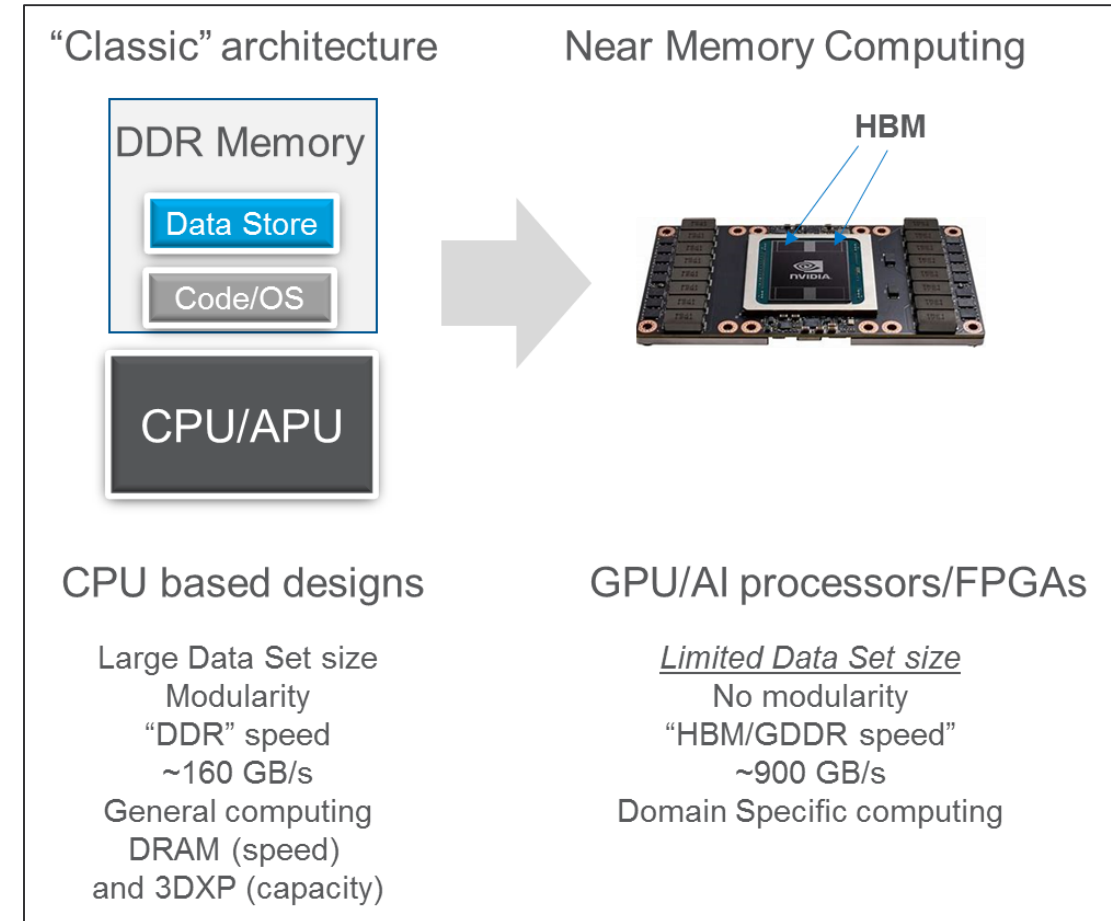
9.6 TB/s aggregate untapped RAW internal BW
Usable BW can be as much as ~25%



Parallel Processing coupled with Near Data Computing enables faster time to insight

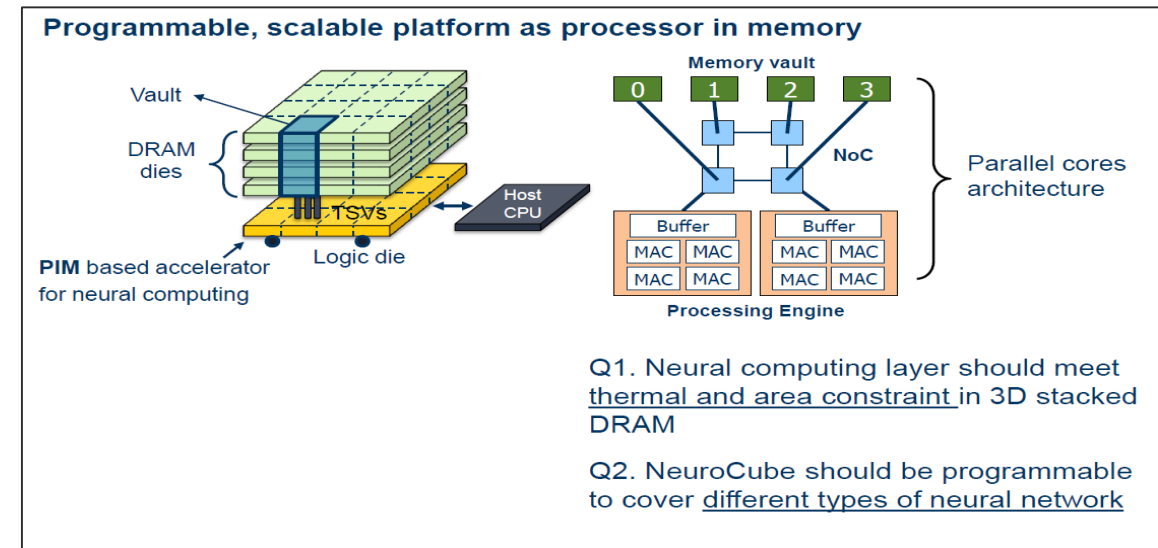
Note: In both cases there are 24 drives per JBOF, identical PCIe G5 NVMe interface.

Near Memory Computing can Deliver 5-10x better Bandwidth to applications in a 2.5D packaging technology



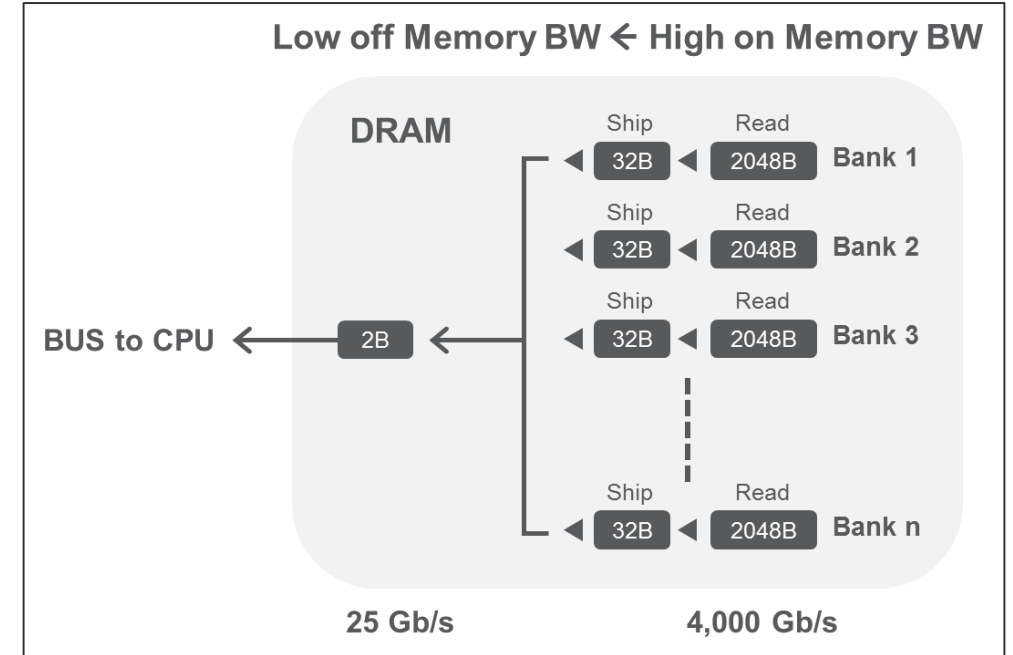
Tighter integration
of processing and
memory can lead to
even greater
Performance at
lower Power

Near Memory Computing example

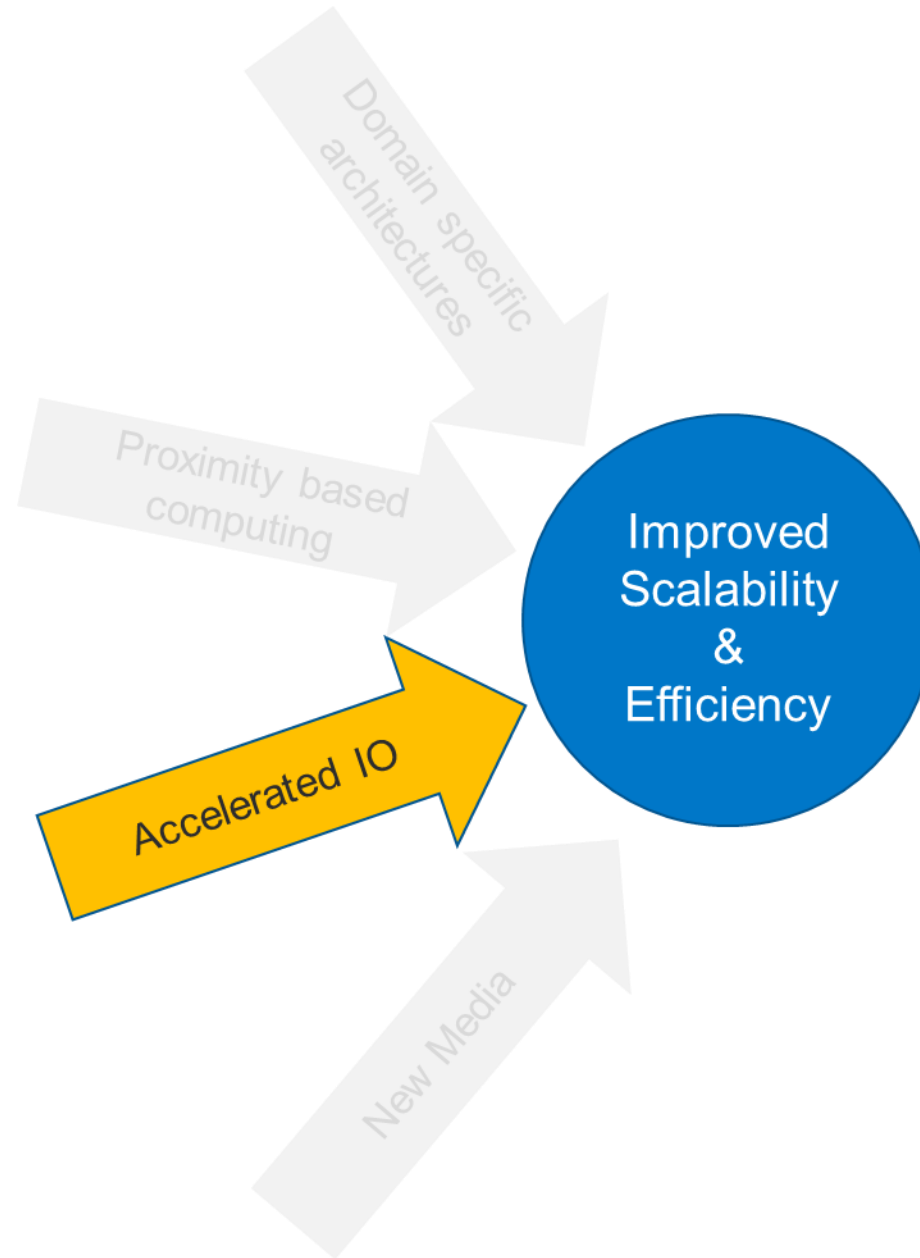


[Kim et al., NeuroCube, ISCA 2016]

**On Die integration of
compute and memory
can take advantage of
~160x on Die bandwidth
on smaller Data sets**



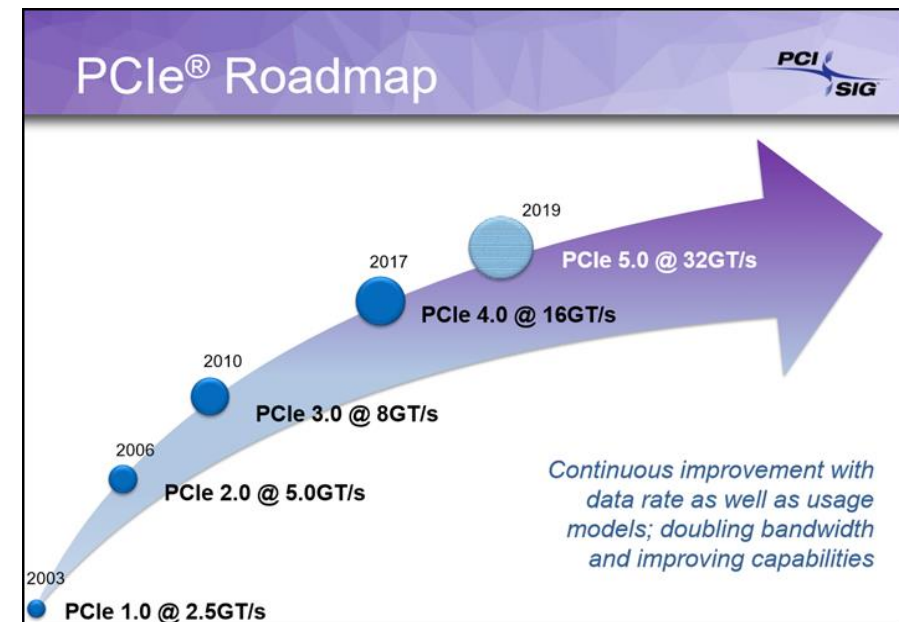
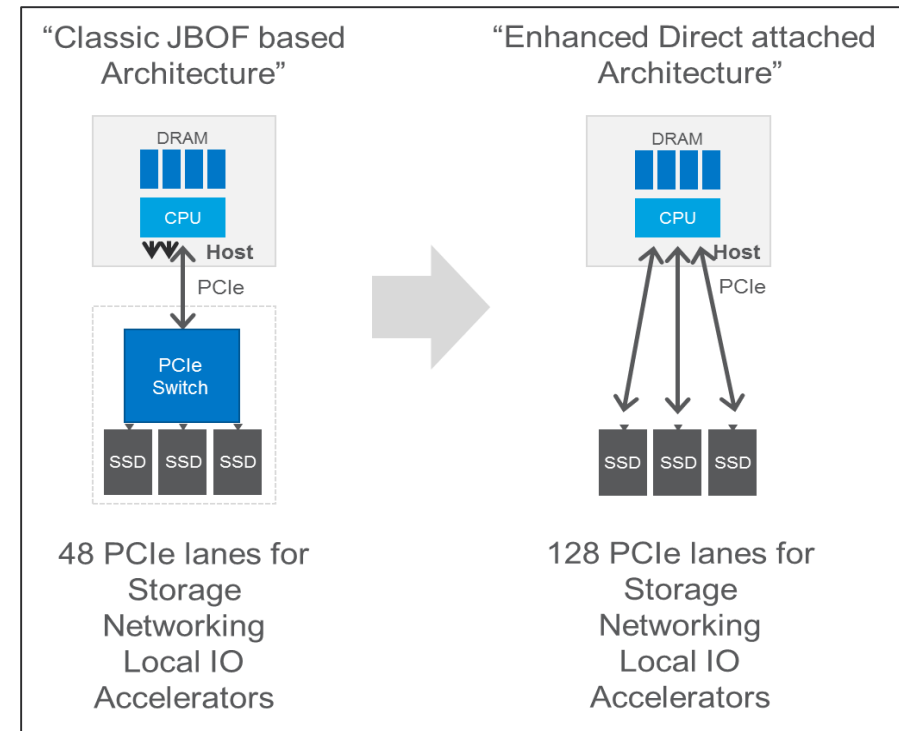
Large investment into improving Platform IO



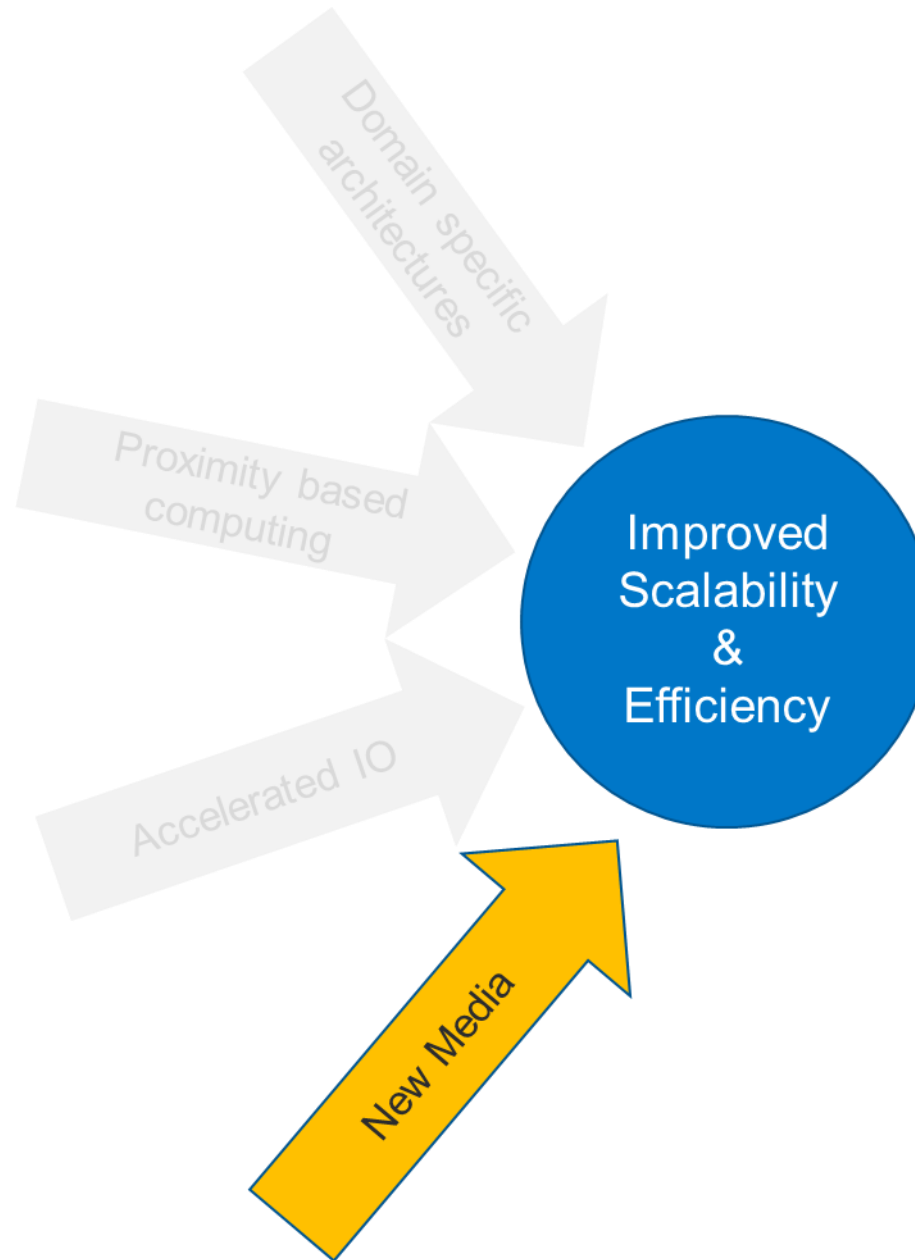
Increasing integrated IO lane count to grow connectivity

Faster PCIe will enable higher speed devices

Silicon photonics to extend reach



Emerging Media technologies to improve capacity, latency and access methods

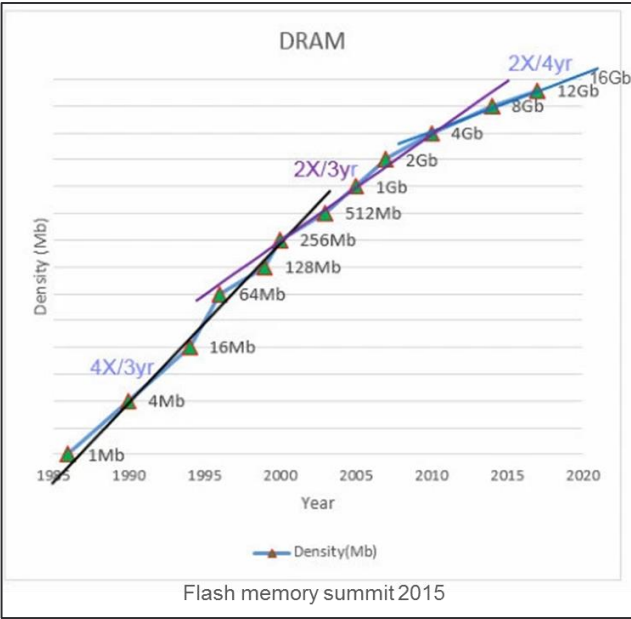


Emerging Memory research is intensifying

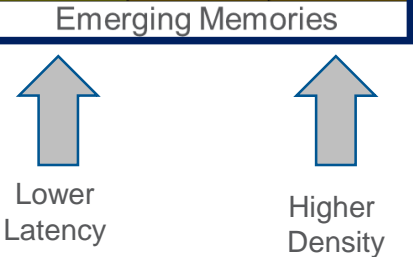
Difficult to beat
DRAM Performance
& Energy

But

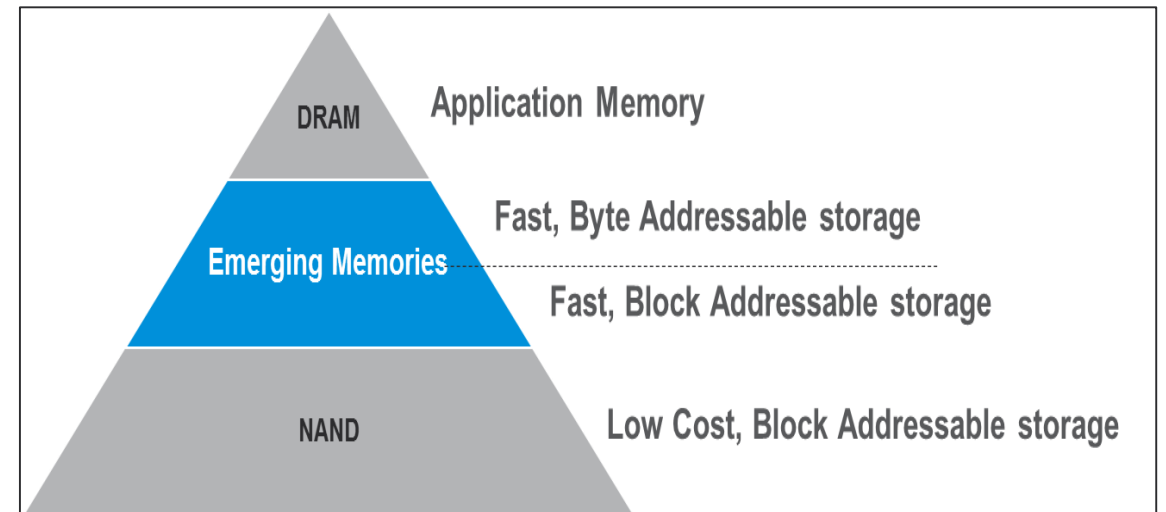
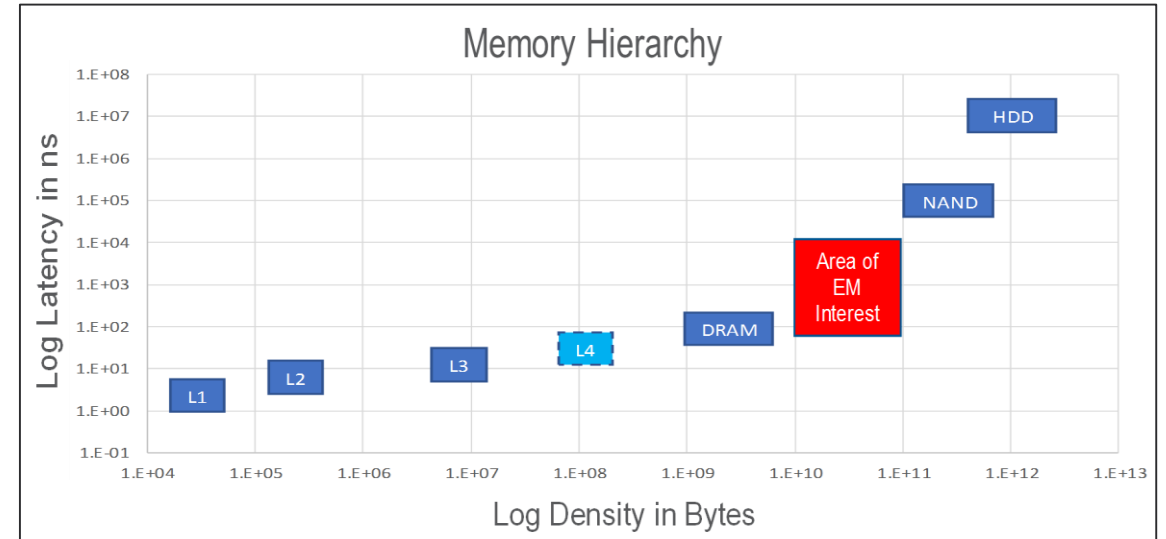
DRAM density growth
is slowing



	DRAM	STTRAM	PCM/ 1T1R RRAM	Cross point RRAM	NAND
Read Latency	20ns	~ 50ns	~100ns-200ns	~100ns-200ns	~10us
Write Latency	20ns	~ 50ns	~1us	~1us	~10us
Read Endurance	>1e15	>10 ¹¹	>10 ⁷	>10 ⁷	>10 ⁷
Write Endurance	>1e15	>10 ¹¹	>10 ⁶	>10 ⁶	2K-100K
Write/Read Energy/bit	<10pJ/bit	~25pJ/bit	~100-200 pJ/bit	~100-200 pJ/bit	> 100pJ/bit



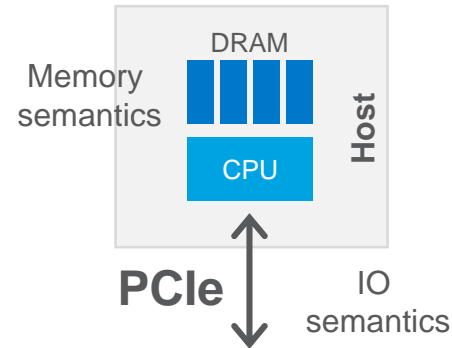
**Emerging
Memories can fill
data access
latency gaps and
enable new
storage models**



Attaching EM
requires server
architecture changes

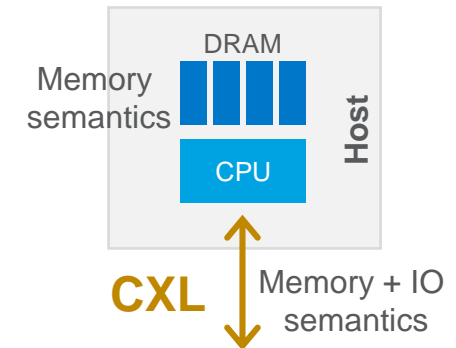
CXL is the emerging
Standard for EM
attach

Today



PCIe Gen4
IO devices
SSD
Accelerators

~2022

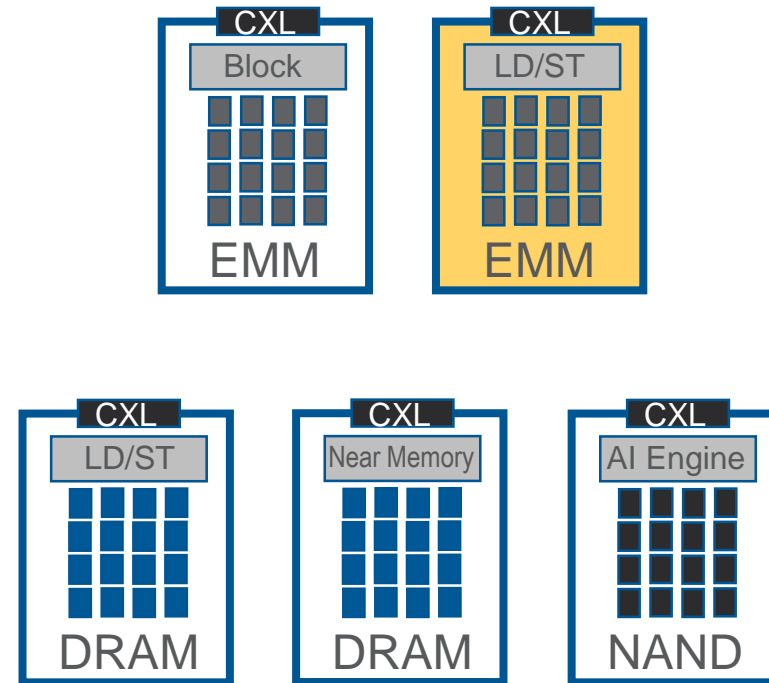


PCIe Gen5
2x IO Bandwidth
IO devices
SSD
Emerging Memory
Coherent Accelerators

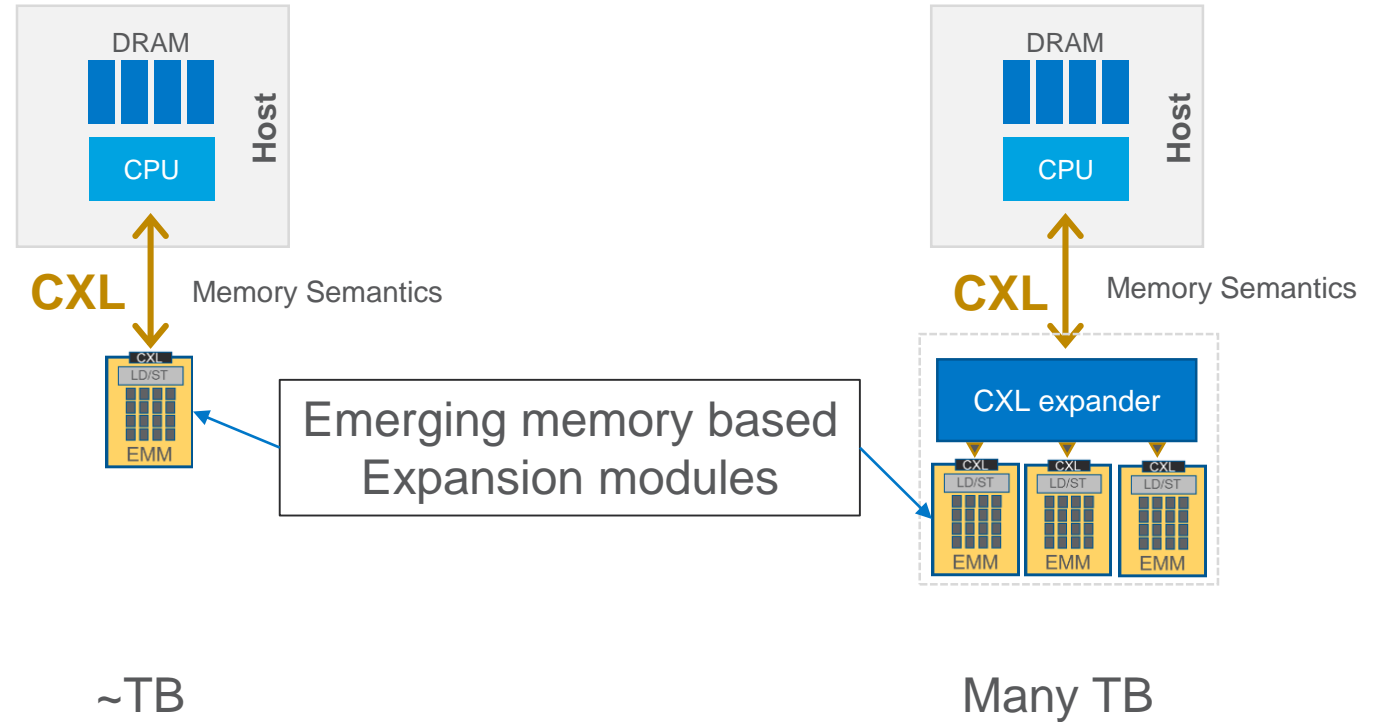
* For simplicity sake single CPU Hosts are shown

**CXL enables
innovations ranging
from different
memory types to
heterogeneous
computing**

Some possible memory examples

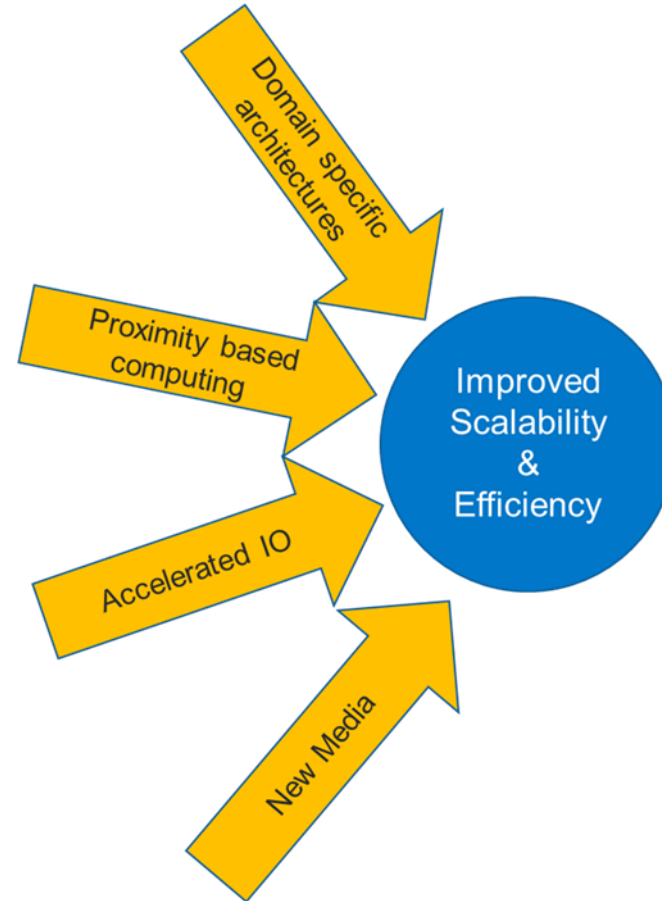


Scaling Emerging Memory capacity is critical to address use case requirements

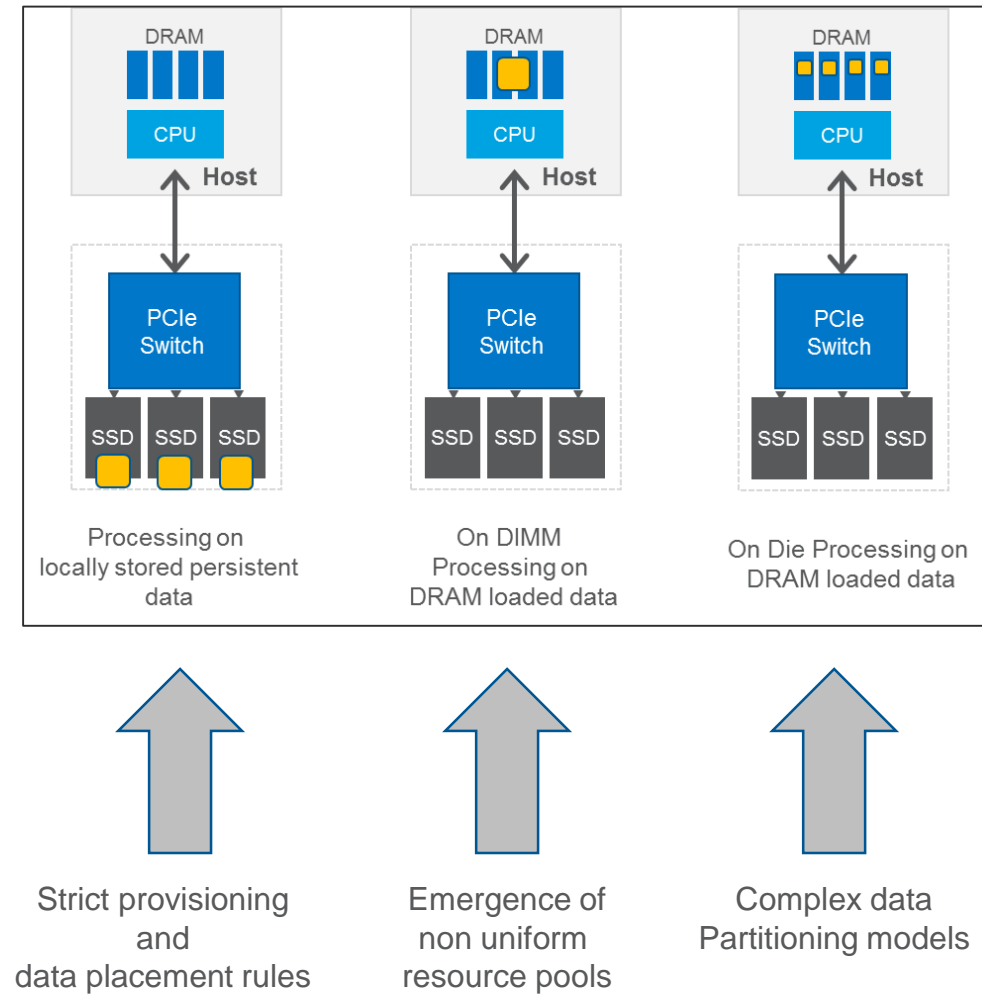


* For simplicity, single CPU Hosts are shown

The collection of these new technologies and architectures will impact how we build future Data Centers and Systems

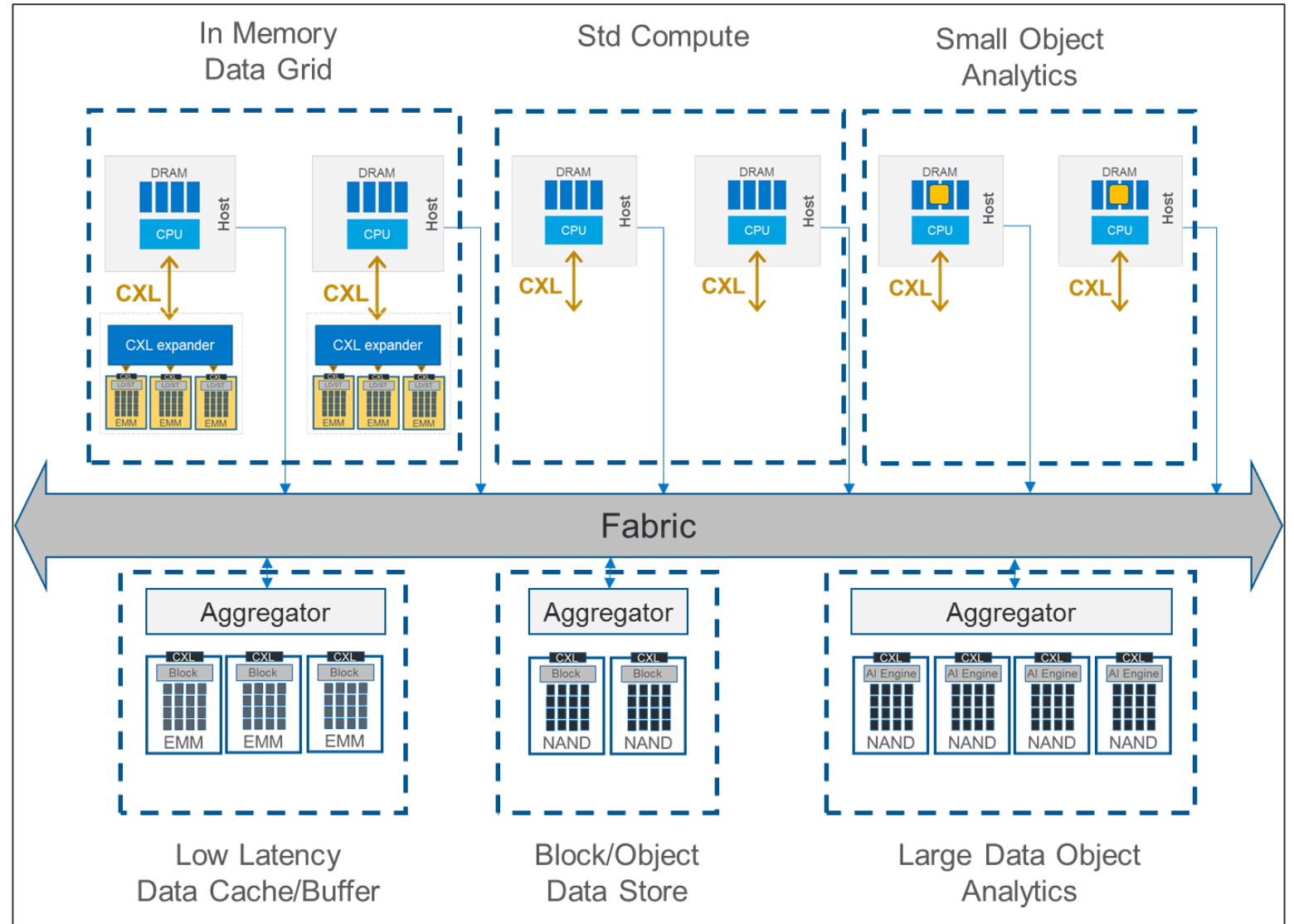


Co locating Compute and Data requires change in data placement, provisioning and load balancing strategies



* For simplicity single socket CPU Hosts are shown

Heterogeneous Resource Pools Example



* For simplicity single socket CPU Hosts are shown

Summary

Exciting times.

The Data Growth and the need for Faster Insight drives transitioning from decades old architectures to a new, emerging model utilizing breakthrough technologies

Buckle your seatbelt!

Thank you!



bfleischer@micron.com