# Thinking Fast & Slow:
# Intuition, Reasoning and Emerging Memory

## Dave Eggleston
## Intuitive Cognition Consulting

# Abstract

- Our human brain can be modeled as two distinctly different systems: a real time intuition system, and a background reasoning system.

- As we move into the AI compute era, Emerging Memory technologies play an increasingly important role in overcoming the limitations of DRAM and NAND.

- The commercialization of Emerging Memory will therefore accelerate our realization of powerful AI systems.
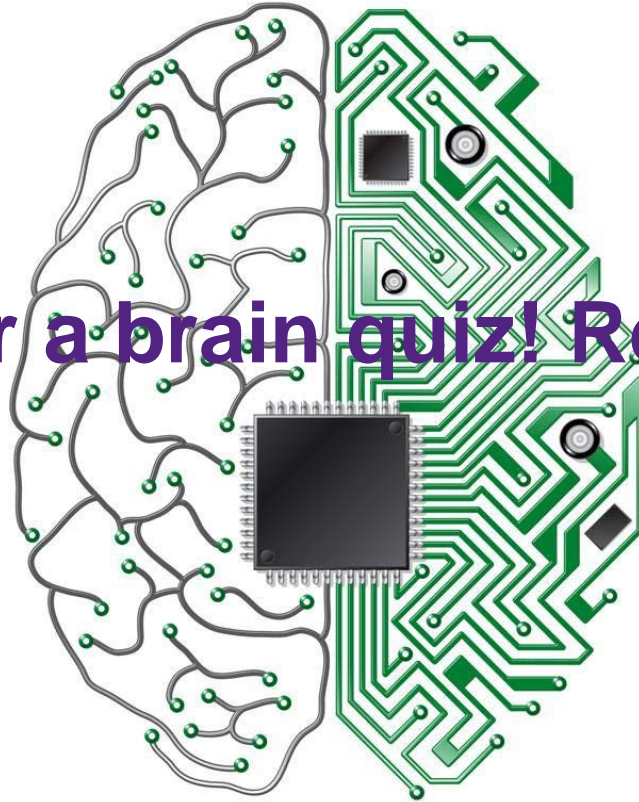
"A groundbreaking tour of the mind, and explains the two systems that drive the way we think."

"System 1 is fast, intuitive, and emotional; System 2 is slower, more deliberative, and more logical."

Daniel Kahneman is professor emeritus of psychology and public affairs at Princeton University.

# Time for a brain quiz! Ready?

SDC 18

# 17 x 24 = ?

# Intuition

# 17 x 24 = ?

# REASONING

SDC 18

# Intuition

System 1
- ❑ Lightning fast
- ❑ Automatic
- ❑ Real time
- ❑ Effortless
- ❑ Approximate

# Edge

# REASONING

SYSTEM 2
- ❑ Slow
- ❑ Interrupt driven
- ❑ Background
- ❑ Energy inefficient
- ❑ Precise

# DATACENTER

# Edge

# DATACENTER
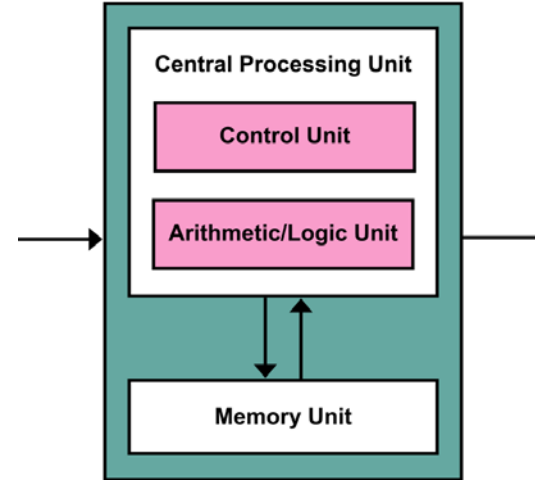


☐ Non-von Neumann architecture
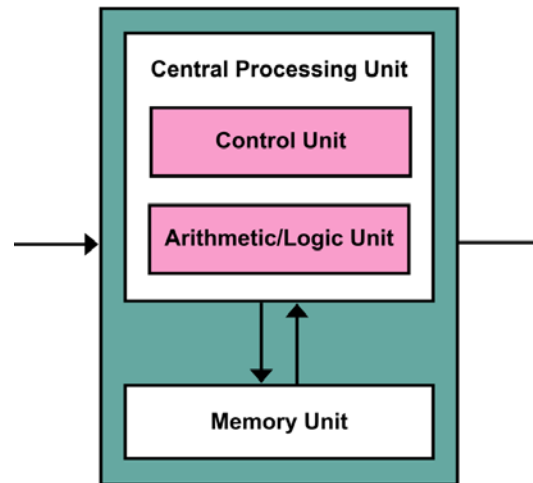
# Intuition

☐ VON NEUMANN ARCHITECTURE

# REASONING

# Edge

Store Dr. Moda for now; we'll come back to discuss Edge & Intuition
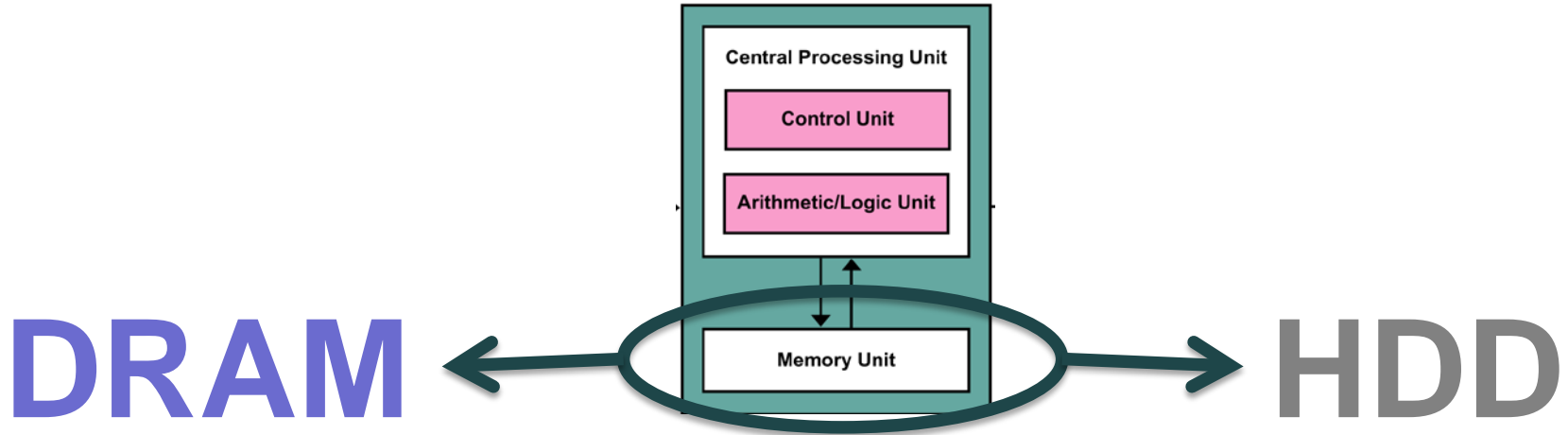
☐ Non-von Neumann architecture

## Intuition

# DATACENTER



Central Processing Unit
- Control Unit
- Arithmetic/Logic Unit
- Memory Unit

☐ VON NEUMANN ARCHITECTURE

# REASONING

# Reasoning

**DRAM** ← → **HDD**

Central Processing Unit

Control Unit

Arithmetic/Logic Unit

Memory Unit

- Once upon a time, long ago…

# Reasoning



**DRAM** ← **NAND**

- And for a while we were happy! ☺ ☺ ☺

# Moore's Law is slowing – but still need low cost bits



- Cost gap between DRAM and NAND continues to increase
- Need cost-effective emerging memory to fill this gap
- Trajectory for DRAM prices for the next 5 years uncertain

Source: IDC

# Moore's Law is slowing – but still need low cost bits



**DRAM costs too much!**

Legend:
- DRAM ASP
- NAND ASP

Chart y-axis values: 1000000, 100000, 10000, ... 1, 0.1, 0.01

Annotations on chart: CAGR -25%, CAGR -28%, **Emerging memory**

X-axis: 1990, 2000, 2010, **2018**, 2020

- Cost gap between DRAM and NAND continues to increase
- Need cost-effective emerging memory to fill this gap
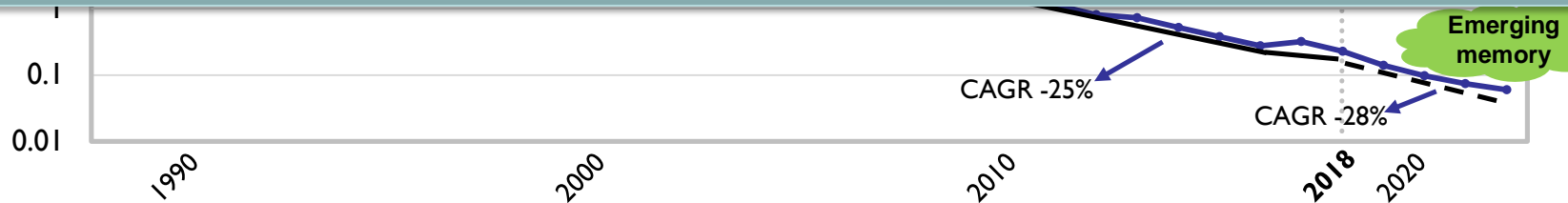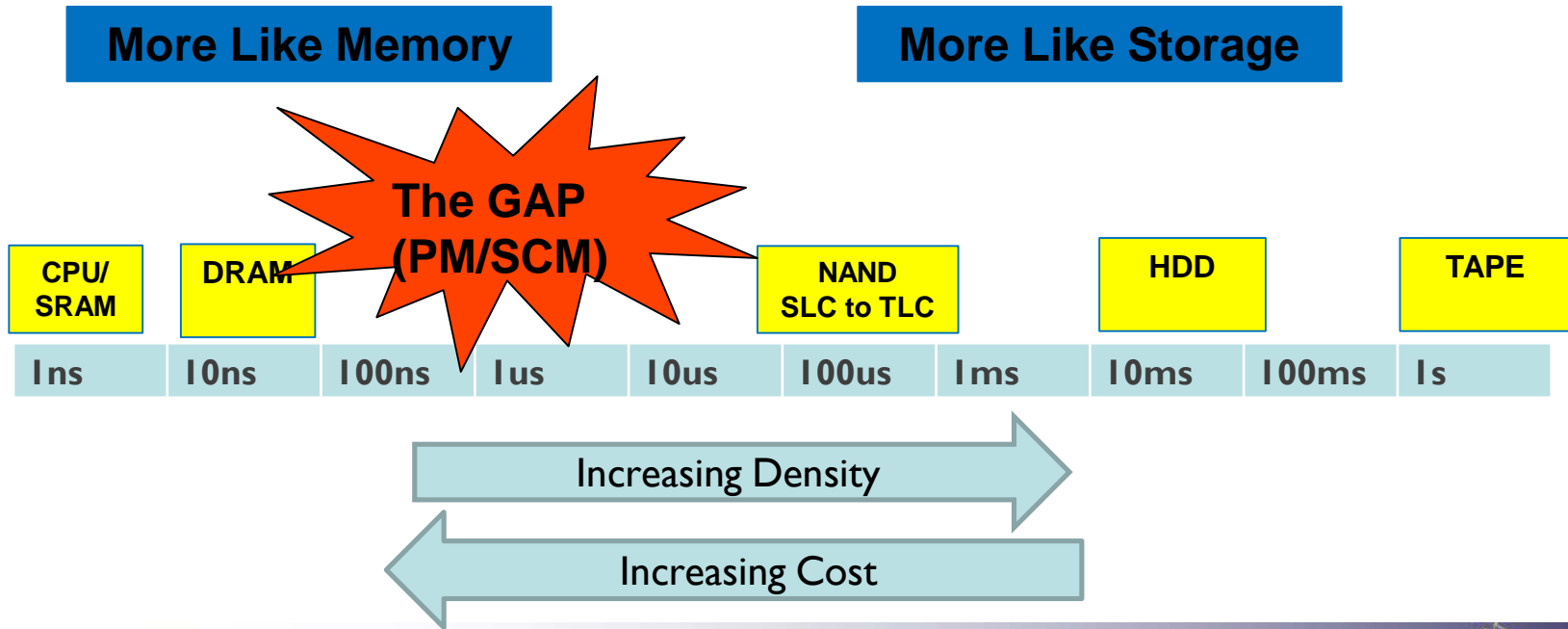- Trajectory for DRAM prices for the next 5 years uncertain

Source: IDC

# The Latency Spectrum and Gaps ~ 2015

**More Like Memory**

**More Like Storage**

**The GAP (PM/SCM)**

| CPU/ SRAM | DRAM | | NAND SLC to TLC | HDD | TAPE |

| 1ns | 10ns | 100ns | 1us | 10us | 100us | 1ms | 10ms | 100ms | 1s |

Increasing Density →

← Increasing Cost

# The Latency Spectrum and Gaps ~ 2015

| More Like Memory | | More Like Storage |

## NAND is too slow!

| SRAM | | | | | SLC to TLC | | | |
|------|------|------|------|------|------|------|------|------|------|
| 1ns | 10ns | 100ns | 1us | 10us | 100us | 1ms | 10ms | 100ms | 1s |

Increasing Density →

← Increasing Cost

# Limitations

**DRAM**
Cost

**NAND**
Latency

# Emerging Memory Targets



Cost
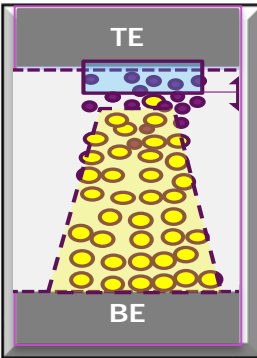1/3rd DRAM

Latency
Read <1μs
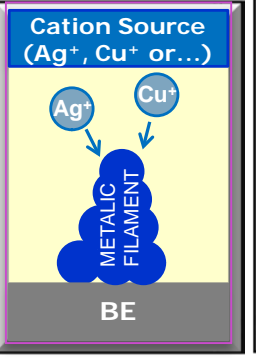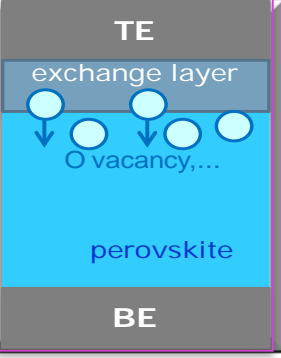
# Emerging Memory Targets



Does such an Emerging Memory even exist?

## Cost
## 1/3ʳᵈ DRAM

## Latency
## Read <1µs

# Switching Mechanisms

| Filamentary ReRAM | | | Interfacial ReRAM | Bulk Transition | | |
|---|---|---|---|---|---|---|
| Oxygen vacancy migration | Thermo-Chemical Fuse/ antifuse | Electro-Chemical ECM | Schottky or Tunnel Barrier | Phase Change PCM | Tunnel Magneto resistance | Electronic MIT (Mott) |
| BIPOLAR | UNIPOLAR | BIPOLAR | BIPOLAR | UNIPOLAR | BIPOLAR | UNIPOLAR |
| TMO: HfO2, TaO2 | TMO: NiO2 | CBRAM: Cu, Ag based | Memristor VMCO, PCMO, TiO2 | Chalcogenide alloys: GST | STT RAM (CoFeB, MgO) | VO2, NbO2 |

# Switching Mechanisms

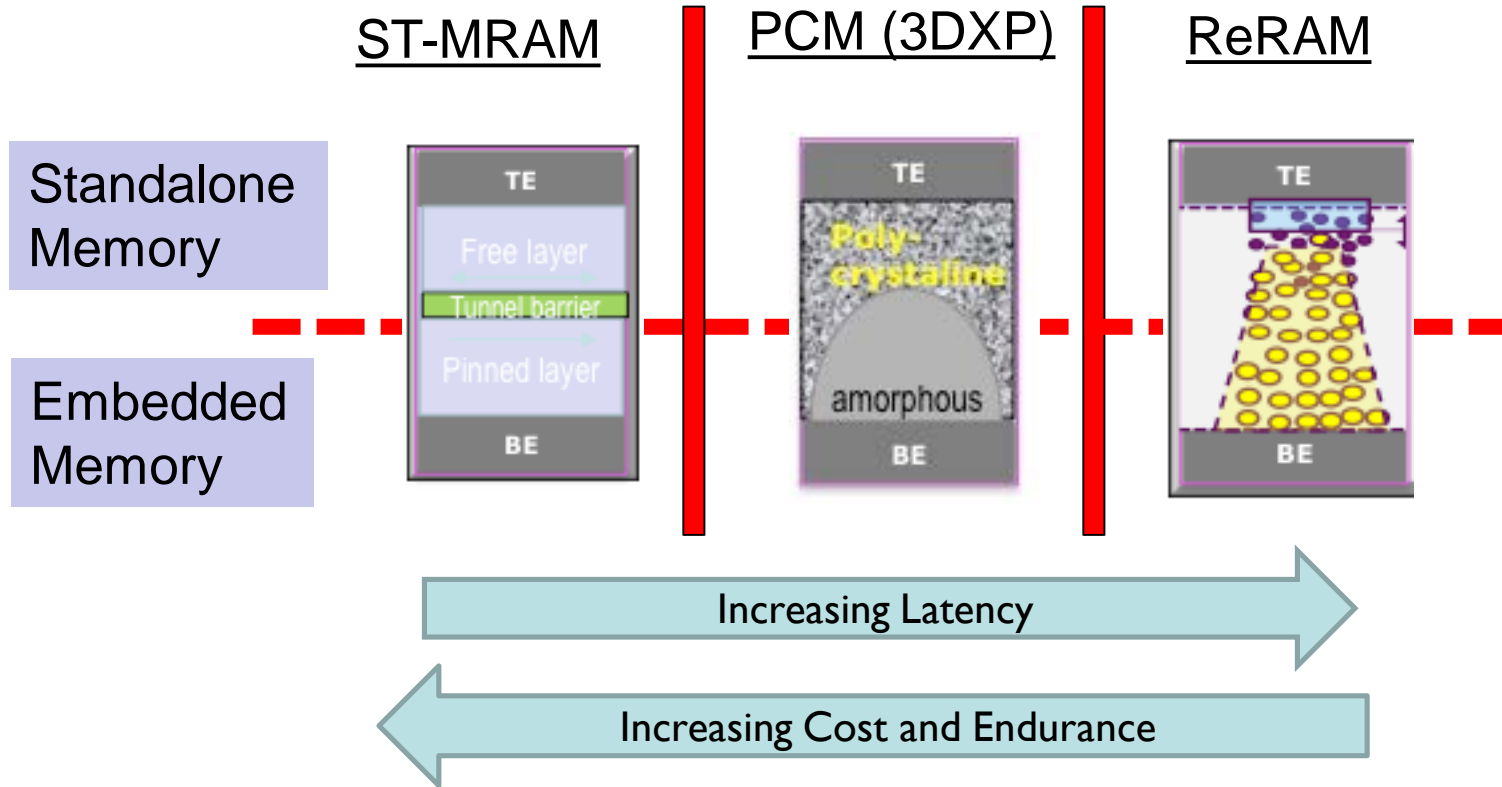| Filamentary ReRAM | | | Interfacial ReRAM | Bulk Transition | | |
|---|---|---|---|---|---|---|
| Oxygen vacancy | Thermo-Chemical | Electro-Chemical | Schottky or Tunnel | Phase Change | Tunnel Magneto | Electronic MIT |
| BE | BE | BE | BE | BE | BE | BE |
| BIPOLAR | UNIPOLAR | BIPOLAR | BIPOLAR | UNIPOLAR | BIPOLAR | UNIPOLAR |
| TMO: HfO2, TaO2 | TMO: NiO2 | CBRAM: Cu, Ag based | Memristor VMCO, PCMO, TiO2 | Chalcogenide alloys: GST | STT RAM (CoFeB, MgO) | VO2, NbO2 |

(Too) many switching mechanisms!

# Classifying the What

ST-MRAM          PCM (3DXP)          ReRAM

Standalone
Memory

Embedded
Memory



Increasing Latency

Increasing Cost and Endurance

# WHO is doing what?

| ST-MRAM | PCM (3DXP) | ReRAM |
|---|---|---|

**Standalone**

EVERSPIN® TECHNOLOGIES

(intel)
Micron®

SONY
adesto™ TECHNOLOGIES

**Embedded**

GLOBALFOUNDRIES®
SAMSUNG  tsmc  UMC

ST

CROSSBAR
Panasonic

SDC 18

# Let's focus on the two shipping technologies!



ST-MRAM

PCM (3DXP)

ReRAM

Standalone

Embedded

# WHO is doing what?



ST-M... ReRAM

EVE... TECH...

Standalone

Embedded

GLOBAL...

SAMSUNG ts...

SONY

...desto™
...HNOLOGIES

CROSSBAR

...anasonic

CNT/NRAM
NEW MEMORY B...
FE HFOX/HRAM
OTHER "SCM"

2018 Storage Developer Conference. © Intuitive Cognition Consulting All Rights Reserved.

# Does such an Emerging Memory even exist?

| | **Cost**<br>1/3rd of DRAM | **Latency**<br>Read <1us |
|---|---|---|
| PCM<br>(3DXP) | 🖐️ | 👎 |
| ST-MRAM | 👎 | 👍 |

Uh oh…

That's not good.

# A Word to the Wise: ~~Replacements~~

- There are <u>no</u> 1:1 memory replacements.

- Stop. looking. for. Them.

- Moving from DRAM+NAND world into one with <u>Combinations</u> of memory.

Just one word.
Combinations.

# Let's use memory combinations!

| | **Cost** 1/3rd of DRAM | **Latency** Read <1us | Winning Combination |
|---|---|---|---|
| PCM (3DXP) | 👍 | 👍 | +DRAM DIMMs for reduced latency! |
| ST-MRAM | 👍 | 👍 | +NAND for reduced SSD cost! |

# 2018 products featuring Emerging Memory

# FMS18: Intel Optane DIMM (3DXP on DRAM bus)

Big and Affordable Memory

High Performance Storage

Direct Load/Store Access

Native Persistence

128, 256, 512GB

DDR4 Pin Compatible

Hardware Encryption

High Reliability

**Now shipping samples
broad developer engagement**

# FMS18: Intel Optane DIMM (3DXP on DRAM bus)

Technology and Architecture and Latency

Technology and Architecture shifts

"Closer to the processor" trend

- Drive latency
- Controller latency
- Software latency

http://ieeexplore.ieee.org/document/8003284/

Really Close to Zero!*

time on calendar scale

*illustrative purposes only – not measured or part of referenced pap

See page 2 "Notices and Disclaimer

- Arch, software, hardware total effort
- Reduces Optane read latency to a few μs
- DRAM operates as "near memory", Optane operates as "far memory"
- Intel controlled

# FMS18: IBM Flash Core Module (MRAM+NAND)
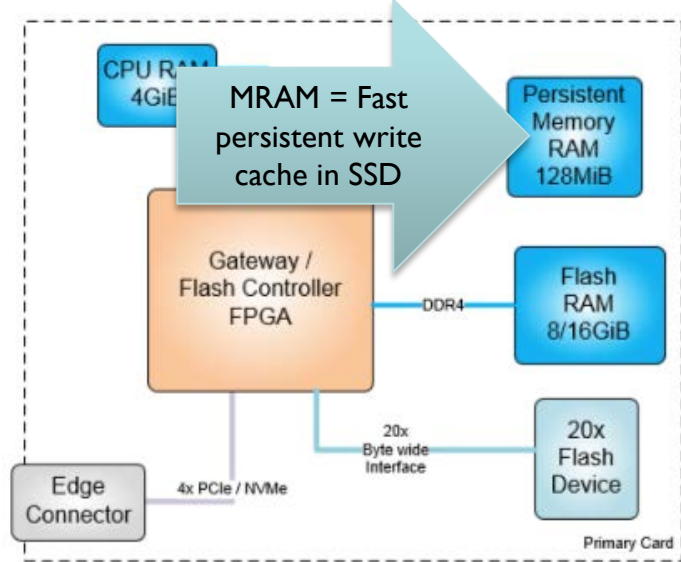
## Introducing The IBM FlashCore Module

IBM **FlashCore™** technology delivers key differentiators

- Built in, performance neutral **hardware compression** and **encryption**
- Using **64 layer 3DTLC NAND**
- Enterprise data **reliability**
- Cognitive Algorithms for **Wear Levelling, Health binning, Heat segregation** and media management
- Intelligent media management that **keeps settings ideal** to keep performance consistent.
- **Endurance** without latency penalty
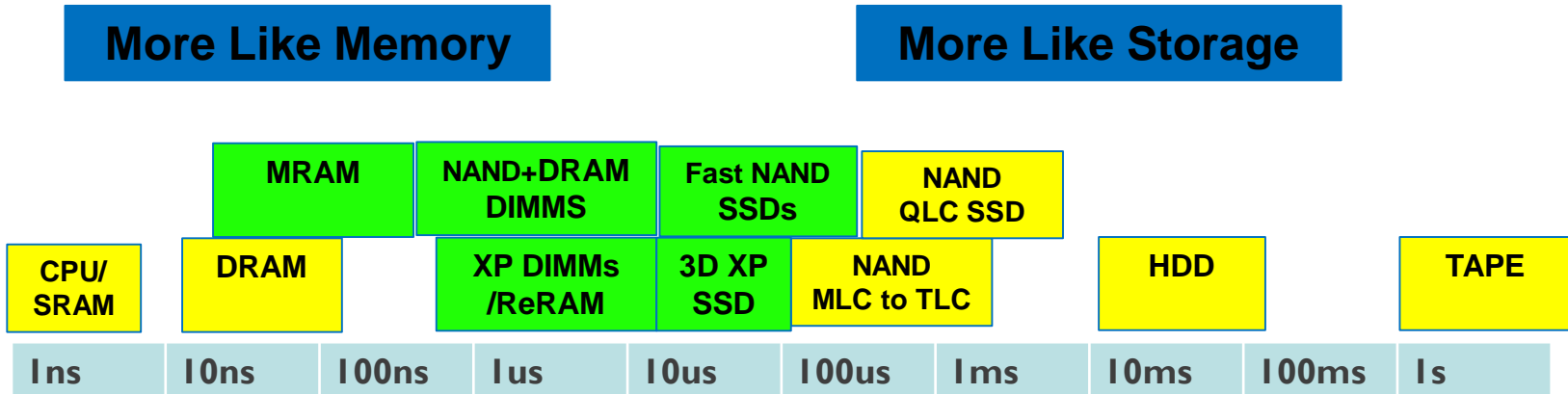- **FIPS 140** certification

**4.8TBu, 9.6TBu, 19.2TBu** capacity options with up to 3:1 compression

CPU RAM 4GiE

MRAM = Fast persistent write cache in SSD

Persistent Memory RAM 128MiB

Gateway / Flash Controller FPGA

DDR4

Flash RAM 8/16GiB

20x Byte wide Interface

20x Flash Device

Edge Connector

4x PCIe / NVMe

Primary Card

# The Latency Spectrum and Gaps ~ Now

| More Like Memory | | | | | More Like Storage | | | |
|---|---|---|---|---|---|---|---|---|

| | MRAM | NAND+DRAM DIMMS | Fast NAND SSDs | NAND QLC SSD | | | | |
| CPU/ SRAM | DRAM | XP DIMMs /ReRAM | 3D XP SSD | NAND MLC to TLC | | HDD | | TAPE |
| 1ns | 10ns | 100ns | 1us | 10us | 100us | 1ms | 10ms | 100ms | 1s |

# The Latency Spectrum and Gaps ~ Now

**More Like Memory**                    **More Like Storage**

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | MRAM | NAND+DRAM DIMMS | Fast NAND SSDs | NAND QLC SSD | | | | | |
| CPU/ SRAM | DRAM | XP DIMMs /ReRAM | 3D XP SSD | NAND MLC to TLC | | HDD | | TAPE | |
| 1ns | 10ns | 100ns | 1us | 10us | 100us | 1ms | 10ms | 100ms | 1s |

## Several memory combinations reduce latency!

# Other interesting 2018 "stuff"

# FMS18: Toshiba XL-Flash (Low Latency SSD)

- 10x reduced latency vs. TLC
- Still not 1us ☹
- For low latency SSDs, attached to compute nodes
- Samsung Z-NAND, Intel Optane SSDs are competitors

# FMS18: JEDEC NVDIMM-P

- Emerging memory and DRAM on the same DDR bus
- Open standard
- Non-deterministic behavior allowed
- Will compete with Intel Optane DIMMs
- Backed by all major memory companies
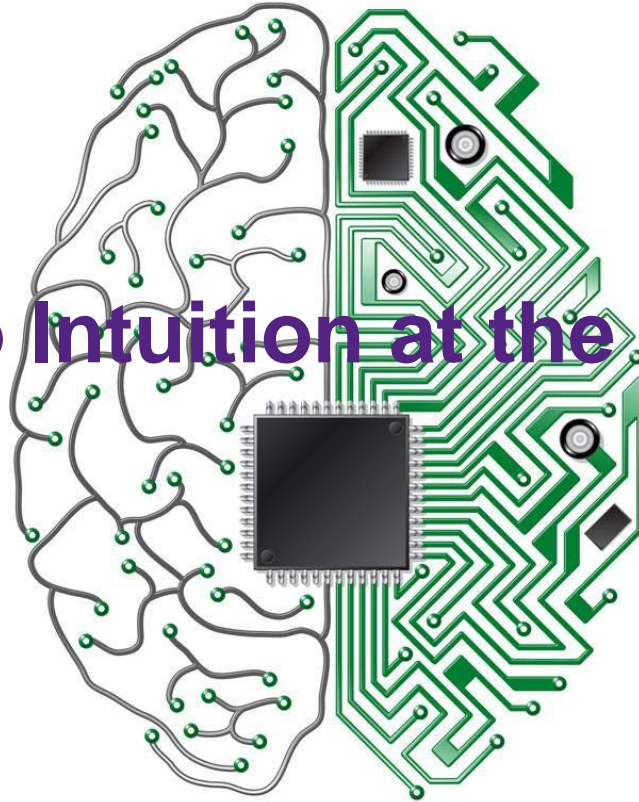
# Summary: Your 2018 SSD/DIMM Watch List

- Everspin MRAM in low latency SSDs
- Intel 3DXP in lower cost server Optane DIMMs
- Toshiba XL NAND in lower latency SSDs
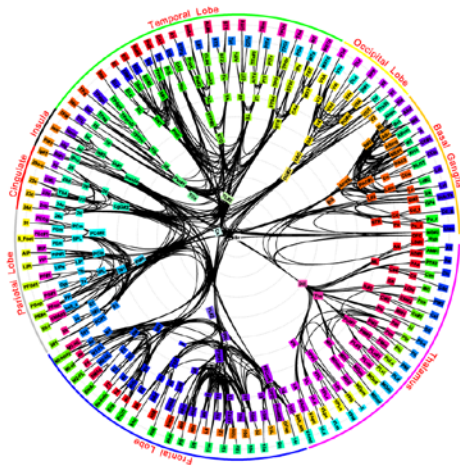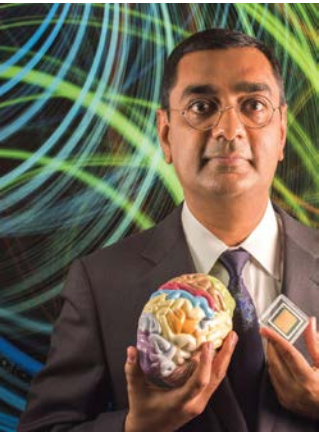- JEDEC NVDIMM-P in open standard DIMMs

# Back to Intuition at the Edge!

# Edge



☐ Non-von Neumann architecture

# Intuition

# DATACENTER

## Let's talk about Intuition at the Edge now!

☐ VON NEUMANN ARCHITECTURE

# REASONING

# Edge



## Intuition

☐ Non-von Neumann architecture

- Intuition system in the human brains use <20 watts of power
- Using von-Neumann architecture for Intuition >20 Gigawatts!
- Highly networked, local compute nodes
- Trained neural nets (NN) perform lightning fast intuition
- Embedded Emerging Memory used to hold weights inside NN
- Sum the weights using analog combine
- Most efficient implementation is analog memory (6+ levels per cell)
- Doesn't require 7nm/5nm process!

# Hot Chips 2018: Edge Intuition SoC

**Mythic Mixed-Signal Computing**



□ Mythic currently uses embedded NOR Flash in analog mode to hold weights
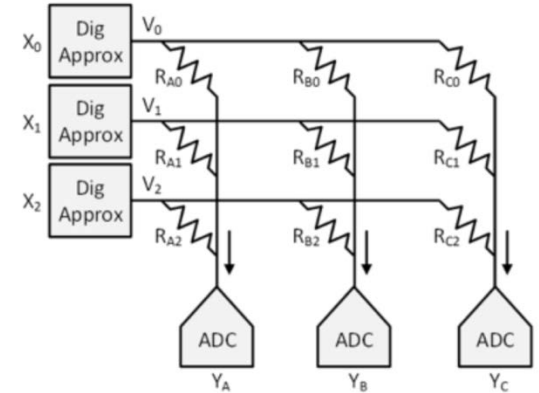
# Hot Chips 2018: Edge Intuition SoC

**Mythic Mixed-Signal Computing**

**Single Tile**

Tiles Connected in a Grid

Example DNN Mapping (Post-Silicon)

**Scene**

$X_0$  Dig Approx  $V_0$  $R_{A0}$  $R_{B0}$  $R_{C0}$

$X_1$  Dig Approx  $V_1$  $R_{A1}$  $R_{B1}$  $R_{C1}$

## Lightning fast, low power, trained NN, analog memory

SIMD    Router

**Expandable Grid of Tiles**

© 2018 Mythic. All rights reser

MYTHIC

☐ Mythic currently uses embedded NOR Flash in analog mode to hold weights

# FMS18: Future Edge Intuition SoCs use ReRAM?

**Image Capture**

**Facial Detection (HOG)**

**Preprocessing**

**Neural Network (FaceNet Inception Model)**

**Representation (128 dimensions)**

**Classification of 100,000s Identities stored in ReRAM in one iteration**
Non-Volatile, Instant-on

Pattern Recognition | Pattern Recognition | ..... | Pattern Recognition | Pattern Recognition

128

CROSSBAR ReRAM

Identity 1 | Identity 2 | ..... | Identity M | Spare

## Simultaneous Processing with Deterministic Performance

- Parallel comparison against all identities
- If no match, new identity created (learning)
- Classification performed in one cycle independent of number of identities

CROSSBAR

# FMS18: Future Edge Intuition SoCs use ReRAM?

| Image Capture | Facial Detection (HOG) | Preprocessing | Neural Network (FaceNet Inception Model) | entation ensions) |
|---|---|---|---|---|

**Classification of 100,000s Identities stored in ReRAM in one iteration**
Non-Volatile, Instant-on

Identity 1 | Identity 2 | Identity M
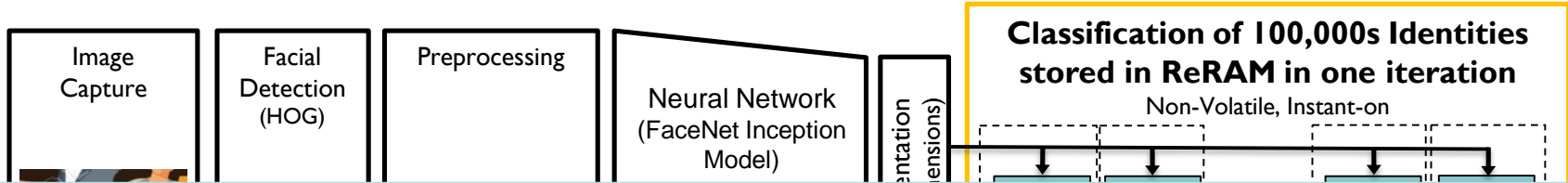
## Embedded ReRAM as analog memory in NN

## Simultaneous Processing with Deterministic Performance

- Parallel comparison against all identities
- If no match, new identity created (learning)
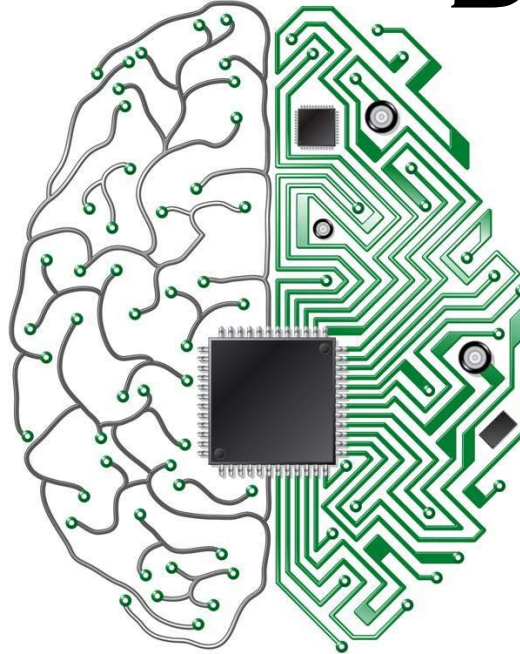- Classification performed in one cycle independent of number of identities

**CROSSBAR**

# Why do I care? What is Emerging Memory enabling for AI?

## Edge

- Analog memory in NN
- Lightning fast intuition at low power

## DATACENTER

- TBs of memory
- Better reasoning, based on more data

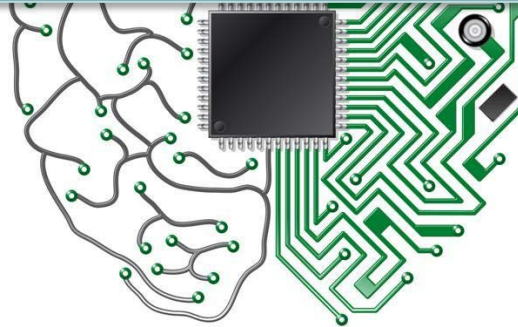# Why do I care? What is Emerging Memory enabling for AI?

## Edge

## DATACENTER

Emerging Memory accelerates AI!

- Lightning fast intuition at low power

- Faster reasoning, based on more data

# Take-away points:

- Your brain has two distinct systems
- Reasoning system needs reduced latency and cost
- There are NO 1:1 memory replacements
- Emerging memory combos with NAND and DRAM
- Intuition system needs new architecture and low power
- Emerging memory utilized as analog weight in neural nets
- Emerging memory accelerates AI systems

# Register NOW for SNIA 2019 PM Summit!

PERSISTENT MEMORY

SNIA

PM SUMMIT

JANUARY 24, 2019 | SANTA CLARA, CA

https://www.snia.org/events/persistent-memory-summit/persistent-memory-summit-2019-registration
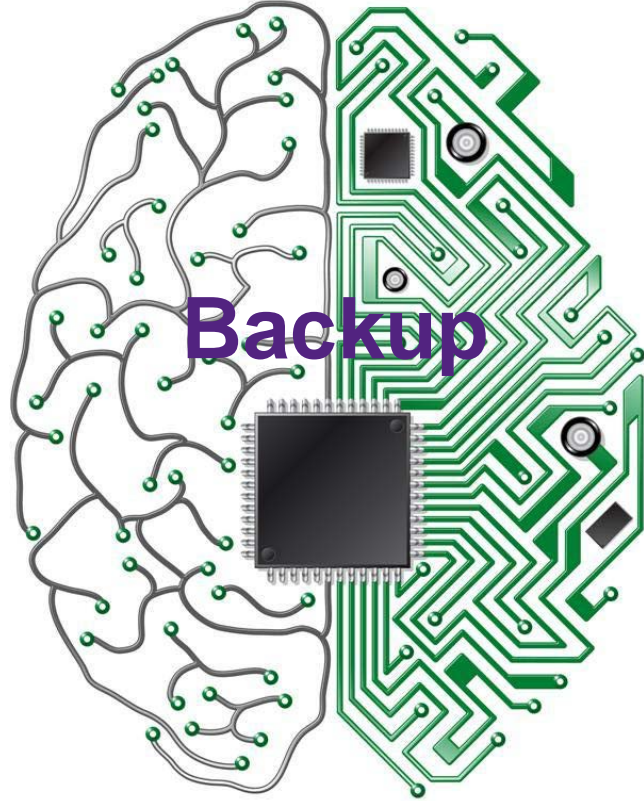
# SDC 18

September 24-27, 2018
Santa Clara, CA

Dave Eggleston
Intuitive Cognition Consulting
Technology & Business Strategy

Email: dave@in-cog.com
Twitter: @NVM_DaveE
LinkedIn:
linkedin.com/in/deggleston/

**Backup**

SDC 18

# Talk Outline

1. Demonstrate the human brain has two distinct systems
2. Discuss the best fit compute architecture to model each AI system
3. Articulate how DRAM and NAND are applied to the compute architectures
4. Identify the key limitations of DRAM and NAND
5. Present and classify some Emerging Memory alternatives
6. Discuss the Emerging Memory system enhancements
7. Identify who is doing what (by when) in the Emerging Memory landscape
8. Articulate the unique challenges in realizing an AI intuition system
9. Propose how Emerging Memory may solve some intuition problems
10. Point to the future of AI systems based on Emerging Memory