



Flash Memory Summit

NVM Usage in the AI Era

Dave Eggleston

Intuitive Cognition Consulting

THINKING, FAST AND SLOW



DANIEL
KAHNEMAN

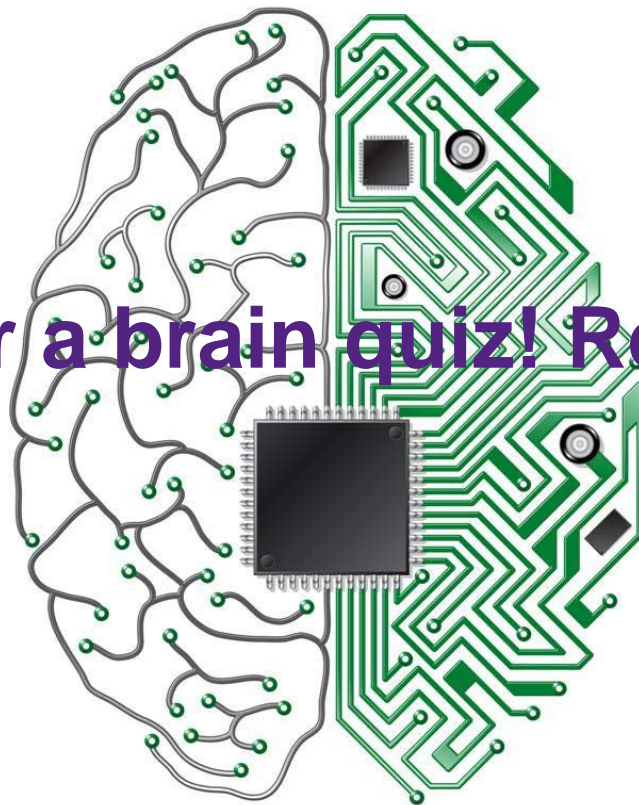
WINNER OF THE NOBEL PRIZE IN ECONOMICS

“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?






$$17 \times 24 = ?$$







Intuition

$$17 \times 24 = ?$$

REASONING



Intuition

System 1

- Lightning fast
- Automatic
- Real time
- Effortless
- Approximate

Edge

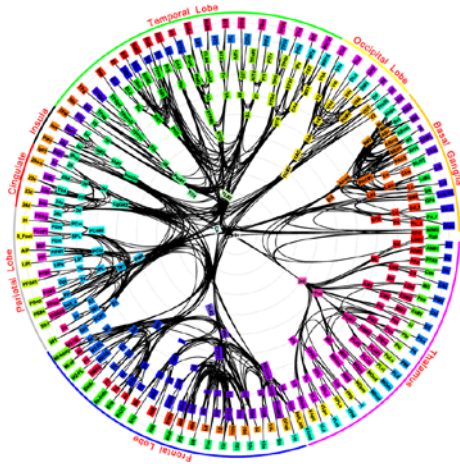
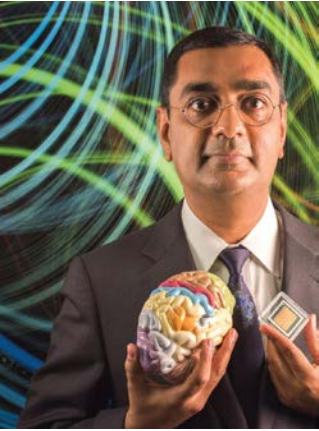
REASONING

SYSTEM 2

- Slow
- Interrupt driven
- Background
- Energy inefficient
- Precise

DATACENTER

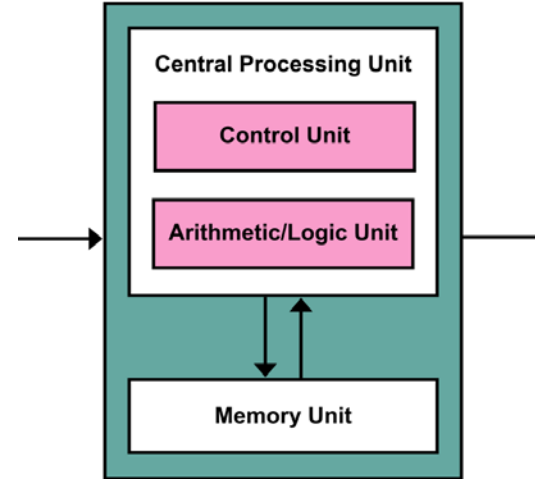
Edge



❑ Non-von Neumann architecture

Inference

DATACENTER



❑ VON NEUMANN ARCHITECTURE

Inference

TRAINING



AI Bottlenecks, Wants, and Solutions

Edge Inference



MCU Bottleneck

Ultra low power & Speed

NVM?

DATACENTER Inference



Memory Bottleneck

Parallelism

NVM?

DATACENTER TRAINING



I/O Bottleneck

Fast & Coherent Checkpointing

NVM?



AI Bottlenecks, Wants, and Solutions

Edge Inference



MCU Bottleneck

Ultra low power & Speed

NVM?

DATACENTER Inference



Memory Bottleneck

Parallelism

NVM?

DATACENTER TRAINING



I/O Bottleneck

Fast & Coherent Checkpointing

NVM?



AI Bottlenecks, Wants, and Solutions

Edge Inference



MCU Bottleneck

Ultra low power & Speed

NVM?

DATACENTER Inference



Memory Bottleneck

Parallelism

NVM?

DATACENTER TRAINING



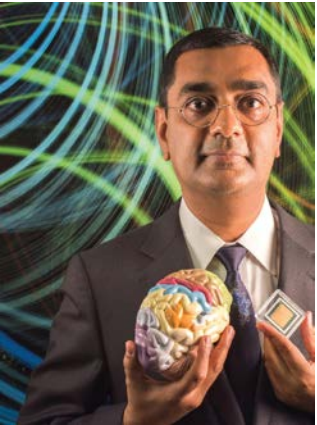
I/O Bottleneck

Fast & Coherent Checkpointing

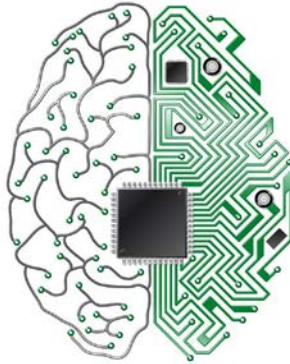
NVM?



Flash Memory Summit



Edge Inference



- Brain operates on <20 Watts
- Von Neumann inference >20 GigaWatts!
- Want non-von Neumann architecture (low power)
- Want real-time inference (speed)
- MCU lacks matrix math capability
- Want a MAC accelerator to do matrix math
- Use analog NVM for MAC acceleration!

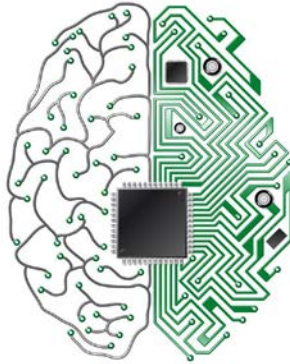
MCU Bottleneck

Ultra low power &
Speed

Analog NVM-based
MAC acceleration



Edge Inference



- Brain operates on <20 Watts
- Von Neumann inference >20GigaWatts!
- Must have non-von Neumann architecture
- Want real-time inference (speed)
- MCU lacks matrix math capability
- Want a MAC accelerator to do matrix math
- Use analog NVM for MAC acceleration!

MCU Bottleneck

Ultra low power & Speed

Analog NVM-based MAC acceleration

MYTHIC

Analog Computation in Flash Memory for Datacenter-scale AI Inference in a Small Chip

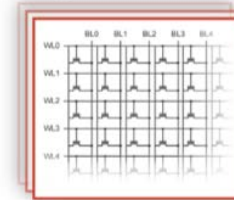
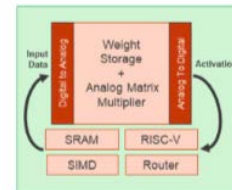
Dave Fick, CTO/Founder
Mike Henry, CEO/Founder

MYTHIC

© 2018 Mythic AI/Infer/accelerated

Mythic Mixed-Signal Computing

Single Tile



Made possible with Mixed-Signal Computing on embedded flash

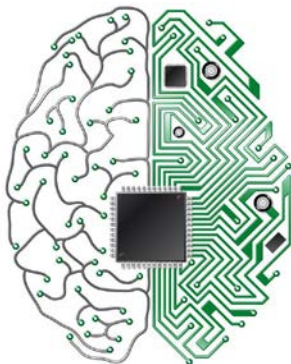
© 2018 Mythic AI/Infer/accelerated



Flash Memory Summit



Edge Inference



- Brain operates on <20 Watts
- Von Neumann inference >20GigaWatts!
- Must have non-von Neumann architecture
- Want real-time inference (speed)
- MCU lacks matrix math capability
- Want a MAC accelerator to do matrix math
- Use analog NVM for MAC acceleration!

MCU Bottleneck

Ultra low power & Speed

Analog NVM-based MAC acceleration

SYNTIANT

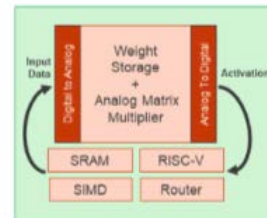


Flash Memory Summit 2019
Santa Clara, CA

© Intuitive Cognition Consulting
Dave Eggleston

Mythic Mixed-Signal Computing

Single Tile



Made possible with
Mixed-Signal Computing
on embedded flash

MYTHIC

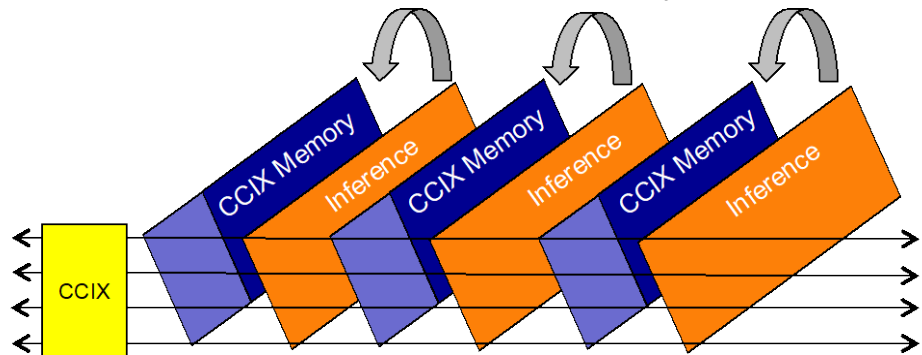
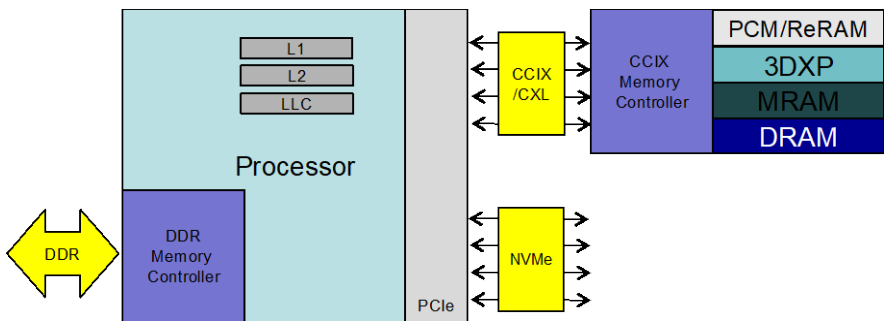


DATACENTER Inference

Memory Bottleneck
Parallelism
CCIX Memory with Inference engines



- Hold trained model in memory
- Inference bottlenecked by CPU-DDR
- Want parallelism for speed
- Multiple sets of CCIX Memory and Inference engines
- CCIX enables load/store peer-peer sharing
- Avoid CPU-DDR bottleneck
- Use DRAM or NVM based CCIX Memory!





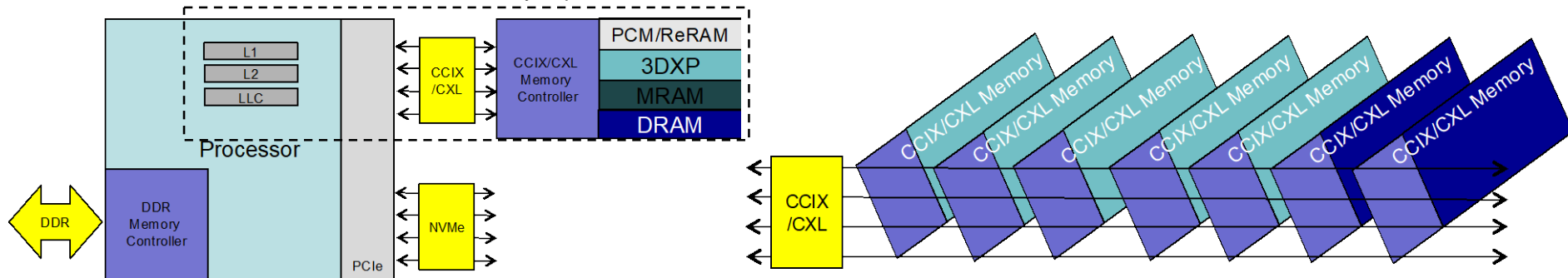
DATACENTER TRAINING



I/O Bottleneck
Fast & Coherent Checkpointing
CCIX/CXL coherent NVM expansion

- Want more memory
- Want fast & coherent NVM checkpointing
- Want fast rebuild time
- Avoid 5us+ latency on I/O
- Avoid using precious DDR slots
- Coherent memory expansion on CCIX/CXL!
- Mix NVM and DRAM on CCIX/CXL!

Cache coherent memory expansion





AI Bottlenecks, Wants, and Solutions

Edge Inference



MCU Bottleneck

Ultra low power & Speed

Analog NVM based MAC acceleration

DATACENTER Inference



Memory Bottleneck

Parallelism

CCIX Memory with Inference engines

DATACENTER TRAINING



I/O Bottleneck

Fast & Coherent Checkpointing

CCIX/CXL coherent NVM expansion



Flash Memory Summit

Now on to the AI NVM Experts!



SYNTIANT



Qualcomm

Dave Eggleston
Intuitive Cognition Consulting
Technology & Business Strategy

Email: dave@in-cog.com

Twitter: [@NVM_DaveE](https://twitter.com/NVM_DaveE)

LinkedIn:
[linkedin.com/in/deggleston/](https://www.linkedin.com/in/deggleston/)

