



Flash Memory Summit

Analog In-Memory Compute using SST SuperFlash

memBraintm

Mark Reiten



Abstract

Artificial intelligence based on advanced machine learning models has gained tremendous momentum in many industries. Many machine learning optimized digital processing solutions have been introduced in the last 5 years, but none can match the power and performance advantages of analog memory-based computing devices. SST's memBrain cost/performance can be 10 times better and the power can be 100x lower than a comparably performing digital solution. This is accomplished by storing multiple levels (up to 256) per cell to represent a "weight" or "synapse" in a neural model. Multiplication is done through cell operation and addition is done by summing the output lines. Vector Matrix Multiplication is accomplished through design techniques so any existing SuperFlash process can support this new optimized compute paradigm.



Agenda

- Neural Systems and Uses
- Deep Neural Networks: *The Problem*
- Analog vs Digital compute
- What is memBrain
- memBraintm Analog Inference Designware
- Many Model Types can be supported
- Why it works
- Software Flow



Neural Systems and Uses

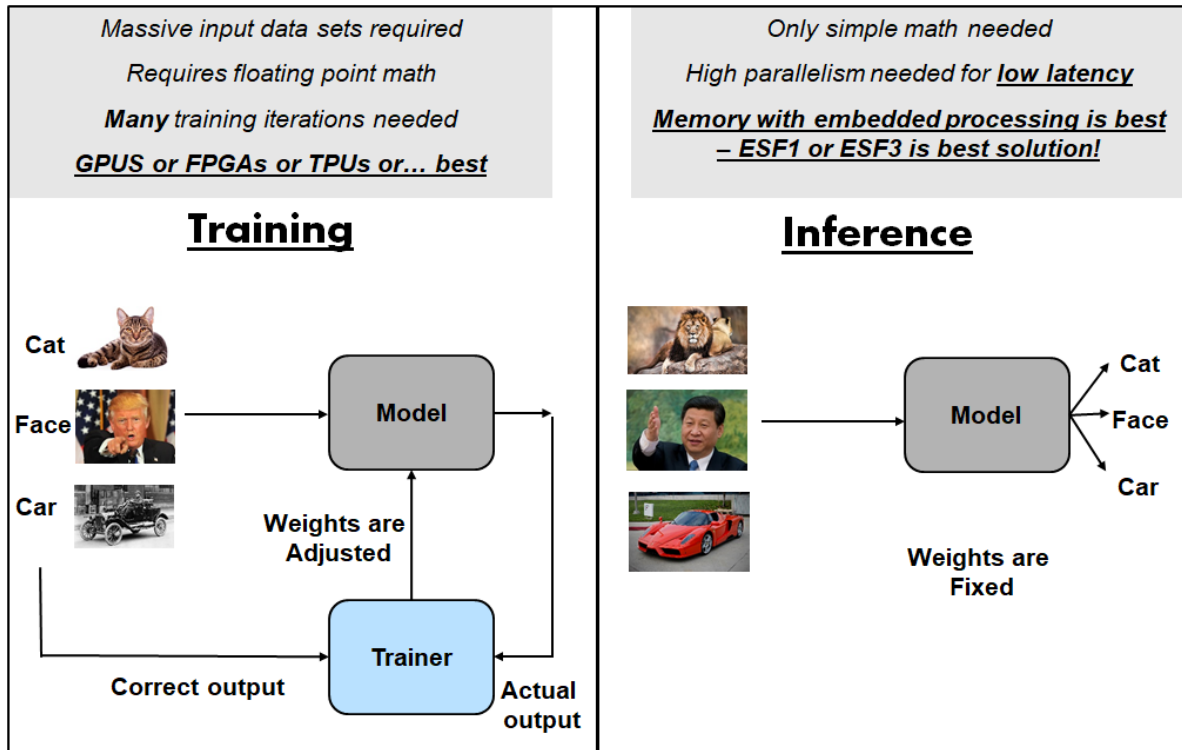


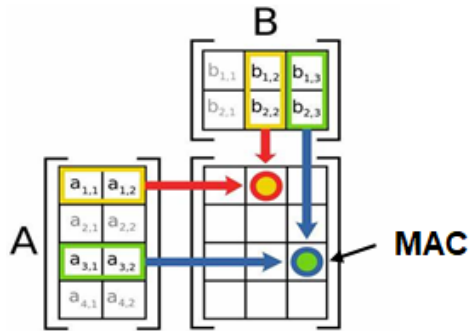
Image Processing
Object Detection and Classification
Speech recognition
Natural Language Processing
Text processing
Search
Similarity/relevance analysis
Recommendations/Ads

Data Mining
Business Pattern extraction
Bioinformatics
Hybrid Applications
Visual scene understanding
Personal Assistants (speech)
Game playing
Robotics (image + speech)



Deep Neural Networks: *The Problem*

DNNs require vast numbers of Multiply-Accumulate operations (MACs)



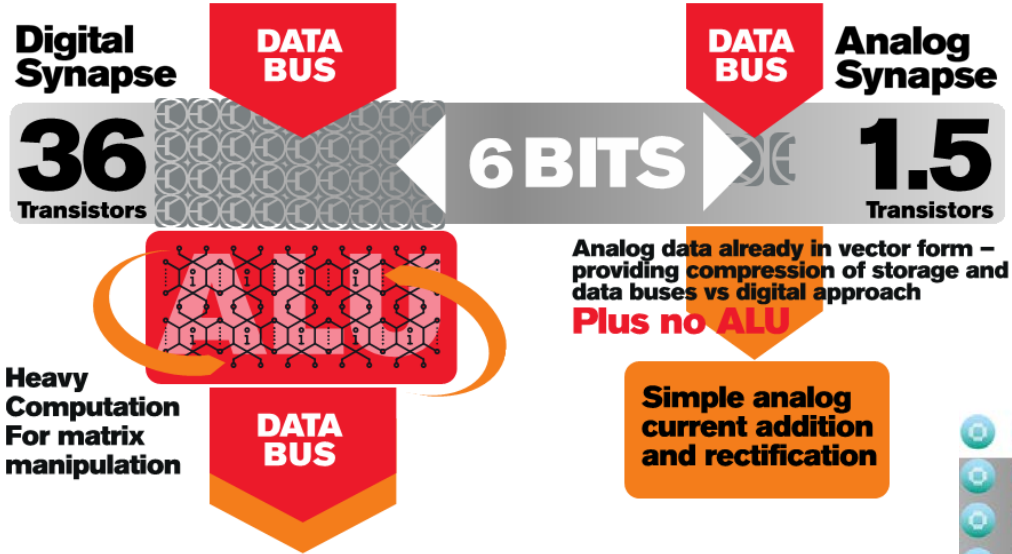
Network	Weights	MACs
Alexnet	61 M	725 M
ResNet 50	23 M	3.5 B
VGG-19	46 M	22 B

Vast numbers of MAC operations favors keeping weights in local storage

Cannot fit these into a stand alone digital edge processor



Analog vs Digital Compute



Analog advantages:

Data and compute compression

No Data thrashing!

In neural networks, involving simple operations on large data sets, time needed to read and write computation registers will dwarf actual computation time.

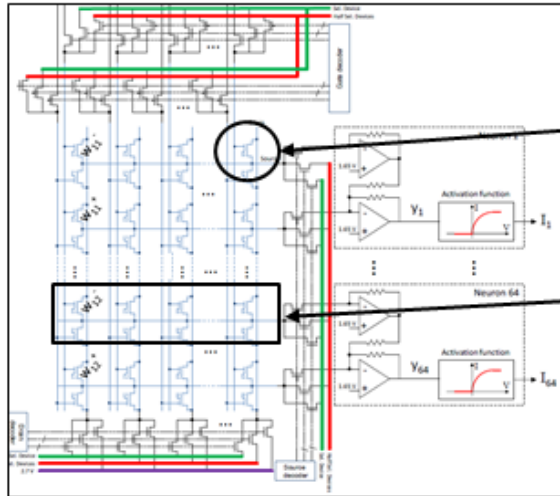




What is memBrain?

memBrain

solves the compute problem by storing the weights in eFlash and using analog cell operation to perform the MAC operations inside the storage array



Each cell stores up to 8 bits in 1.5 transistors
Multiply happens through cell operation

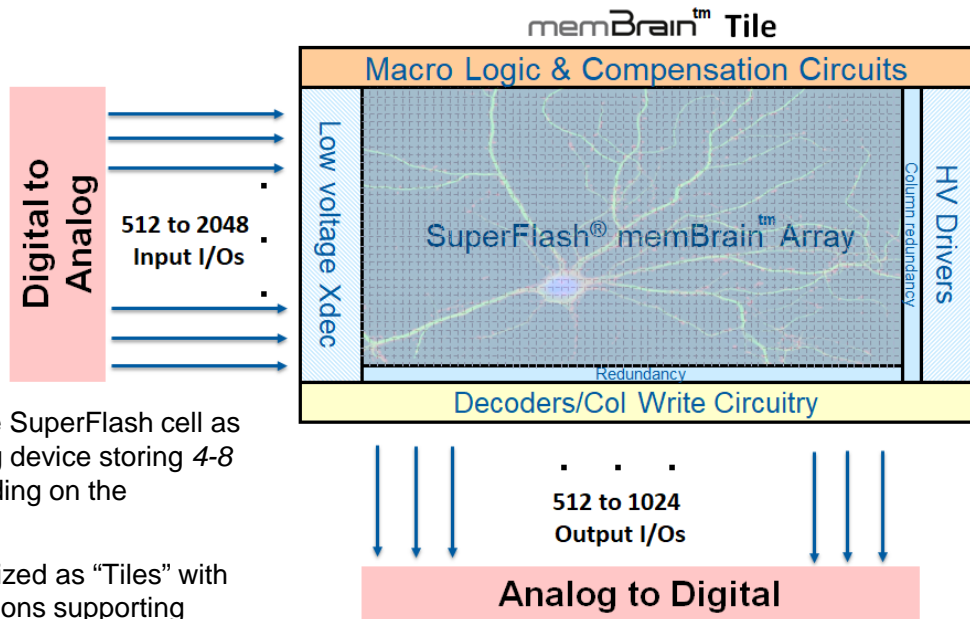
Output line functions as a Neuron
Summation happens along output

Compared to 48 SRAM transistors for 8 bits!

Power is ~.3pj per MAC !!



memBrain™ Analog Inference Designware



memBrain uses the SuperFlash cell as a multi-level analog device storing 4-8 bits per cell depending on the application

memBrain is organized as “Tiles” with wide I/O configurations supporting **massively parallel** multiply/accumulate operations

512x512 tile full frame cycle time 10-30us

Depends on D-A and A-D power

Energy is **0.3pJ per MAC** with D/A+A/D @ 30us frame cycle time

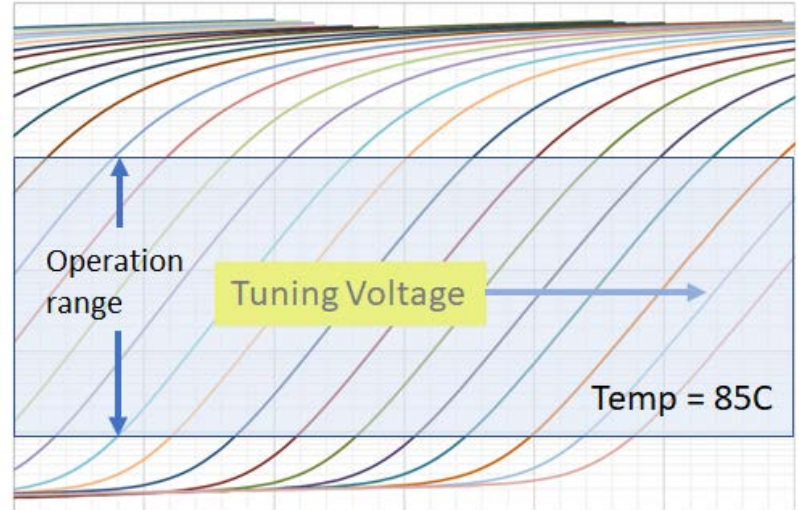
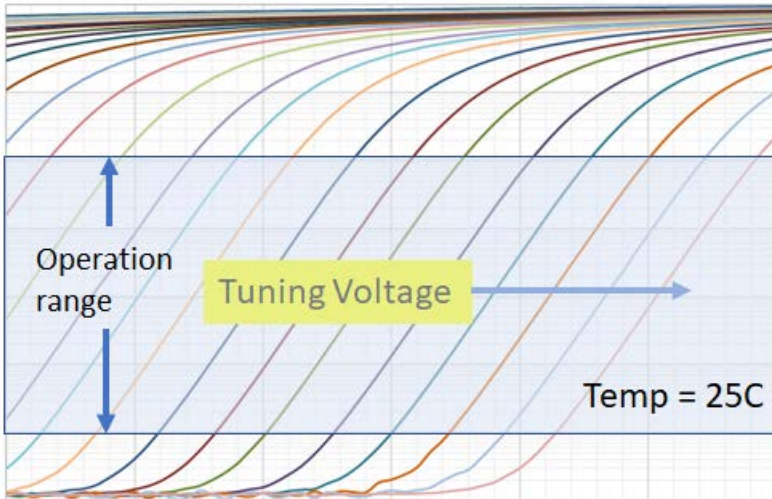
Area with D to A input and A to D output blocks = **.48 mm² on 40nm for 512x512 Tile**

Performance per silicon area and power are orders of magnitude better than optimized digital solutions



Why it Works

ESF3 Cell analog operation varies the floating gate potential with repeatable consistent separation between steps across array and temp

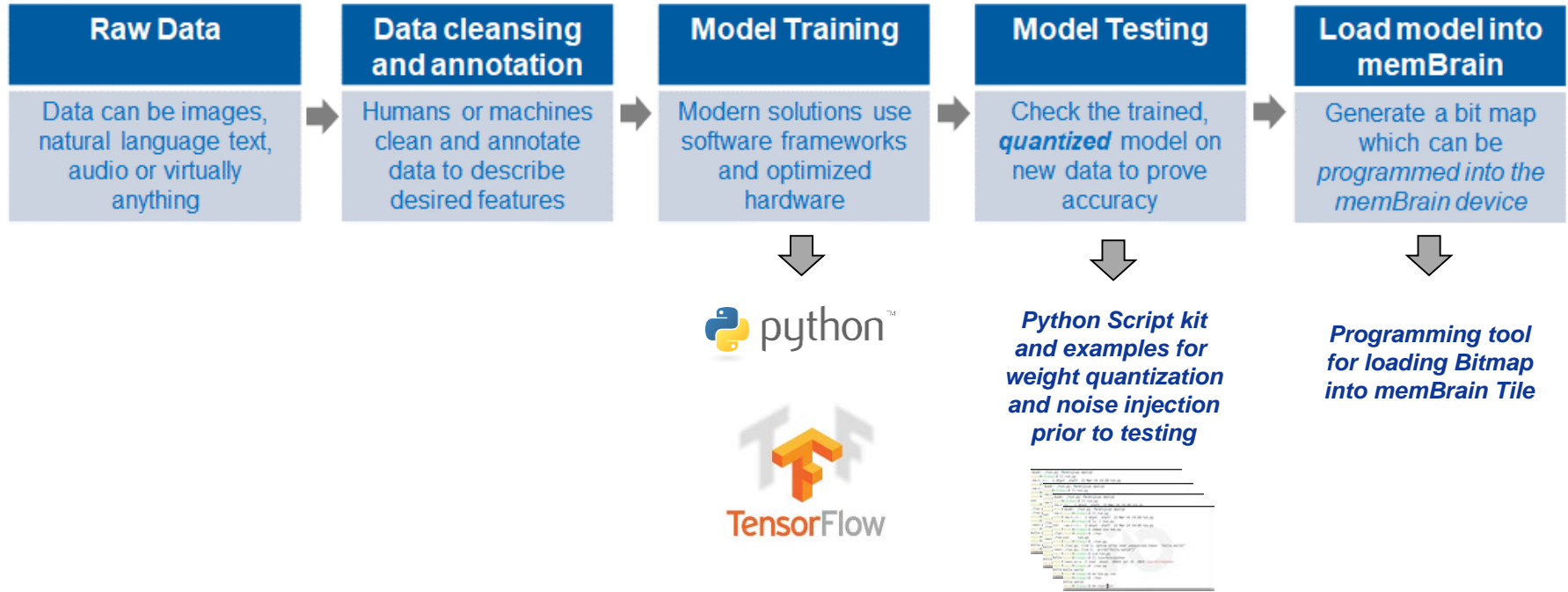


ESF3 Cell data retention is excellent
Noise is low

Noise and temperature can be compensated for



Software Flow





Flash Memory Summit

Thank You