Exceptional service in the national interest





Ultra-Efficient Neural Algorithm Accelerator Using Processing With Memory

Matthew J. Marinella, Sapan Agarwal, Robin Jacobs-Gedrim, Alex Hsia, David Hughart, Elliot Fuller, Steve J. Plimpton, A. Alec Talin, and Conrad D. James

Sandia National Laboratories

*matthew.marinella@sandia.gov



Sandia National Laboratories is a multi-program laboratory managed and operated by Sandia Corporation, a wholly owned subsidiary of Lockheed Martin Corporation, for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-AC04-94AL85000.

Why do we need more efficient computers?

- Google Deep Learning Study
 - 16000 core, 1000 machine GPU cluster
 - Trained on 10 million 200x200 pixel images
 - Training required 3 days
 - Training dataset size: no larger than what can be trained in 1 week
- What would they like to do?
 - ~2 billion photos uploaded to internet per day (2014)
 - Can we train a deep net on <u>one day of image data</u>?
 - Assume 1000x1000 nominal image size, linear scaling (both assumptions are unrealistically optimistic)
 - Requires 5 ZettaIPS to train in 3 days (ZettaIPS=10²¹ IPS; ~5 billion modern GPU cores)
 - World doesn't produce enough power for this!
 - Data is increasing exponentially with time
- Need >10¹⁶-10¹⁸ instruction-per-second on 1 IC
 - Less than 10 fJ per instruction energy budget

Feature 1

Feature 2

Feature 3

Feature 4

Feature 5

Feature 6



Q. Le, IEEE ICASSP 2013

Evolution of Computing Machinery In Sandia Laborator



Inspired by Hasler and Marr, Frontiers in Neuroscience, 2013

Metal Oxide Resistive RAM (ReRAM)

- Example: Sandia TiN/Ta/TaOx/TiN device
- Starts as insulating MIM structure
- Forming: remove O²⁻ → soft breakdown
- Bipolar resistance modulation
- Excellent memory attributes: Switching in less than 1ns, less than 1 pJ demonstrated, scaling to 5nm, >10¹² write cycles





Crossbar Theoretical Limits





Exascale-computations per sec on one chip!

10¹² active devices/chip x 10⁶ cycle per

- In order to not melt the chip, entire area must be limited to ~100W
- Allowed energy per operation = P x t/op $= 100W / 10^{18} = 10^{-16} = 100 aJ/operation$
- 10nm line capacitance = 10 aF
- Can charge line to 1V with 10 aJ

interest per second (1 M-op):

Drawback: "only" ~100B transistors/chip



How does a crossbar perform a useful computation per device?



Electronic Vector Matrix Multiply



Mapping Backprop to a Crossbar



Vector Matrix Multiply, Rank 1 Update: Key kernel used in many algorithms

Integrate current to get an analog output value

Sandia

National

Accelerator Architecture







Neural Inspired Computational Elements



 I_0

 $I_0 - \Delta I$

8



Experimental devices have several nonidealities: Write Variability, Write **Nonlinearity, Asymmetry, Read Noise**

Circuits also have A/D, D/A noise, parasitics





- Use as a neuromorphic weight requires precise analog tuning
- Dataset requires 1000 repeated SET and RESET pulses
- Nominal pulse values
 - SET: +1V 10ns RT/PW/FT
 - RESET: -1V 10ns RT/PW/FT
 - READ: 100 mV 1 ms RT/PW/FT $^{\vee}$



Repeated Pulsed Cycling





TaOx ReRAM in Backprop Training





ТаОх	Large Images	Small Images	File Types
10 ns	84.45%	71.40%	77.67%
100 ns	78.48%	89.48%	67.78%
1 us	71.48%	71.84%	56.33%

How can training accuracy be improved?

Li-Ion Synaptic Transistor for Analog Computation (LISTA)





500 nm anode/gate current-collector electrolyte/insulator

G-V for LISTA Transistor



E. Fuller et al, Adv Mater, accepted 2017

Analog State Characterization



Neural Inspired Computational Elements

Sandia National

aboratories

LISTA-device Performance for Backprop Algorithm



See Poster for Detail!

E. Fuller et al, Adv Mater, accepted 2017

Electrochemical Neuromorphic Organic Device (eNode)





van de Burgt ... Saleo, Nature Mater., 2017 in press

Electrochemical Neuromorphic Organic Device (eNode)





van de Burgt ... Saleo, Nature Mater., 2017 in press

Circuit-Level Improvement



- Allows much closer to ideal with high variability TaOx device
- LISTA achieves essentially perfect accuracy
- Requires tradeoff of energy/latency for accuracy – exact tradeoff depends on algorithm reqs.



Energy and Latency Analysis



	SRAM	Digital ReRAM	Analog ReRAM
Equivalent Area ~450 1k × 1k matrices	400 mm ²	32 mm ²	11 mm ²
Total Time per 1-layer cycle 3 Ops: 2 reads, 1 write	~ 100 μs Transpose read dominated	~ 60 μs Update dominated: 10 ns write	~ 5 μs Temporal coding dominated: 256 levels
Total Energy per 1-layer cycle	~ 1000 nJ Multiply dominated	~ 700 nJ Multiply dominated	~ 15 nJ
Matrix Storage Area	95%	50%	17%
Periphery Area	5%	50%	100% crossbar on top of periphery
Matrices per 400 mm ² Chip	~450	~5,500	~15,000
Energy / Operation	330 fJ	230 fJ	5 fJ
Operations/Second	14 TeraOps/S	270 TeraOps/S	9,000 TeraOps/S

Energy and Latency Analysis



		SRAM	Digital ReRAM	Analog ReRAM Crossbar
Matrix Storage 1024×1024	Area	800,000 μm²	35,000 μm²	Array: 4,300 μm ² Periphery: 8,460 μm ²
[Values are per- array]	Read	30 nJ / 15 μs	15 nJ / 4 μs	~ 3 nJ / ~ 1.5 μs
	Read Transpose	300 nJ / 65 μs	15 nJ / 4 μs	~ 3 nJ / ~ 1.5 μs
	Write	30 nJ / 15 μs	50 nJ / 45 μs	3 nJ / ~ 1.5 μs
Multiply Accumulators [256 in parallel]	Area	19,000 μm²		FREE
	Run: 1M ops	200 nJ / 4 μs		
Output LUT	Area	1,400 μm²		Uses Digital Methods
	Read	1 nJ / 1 μs		
Input/Output Buffers	Area	13,000 μm²		Uses Digital Methods
[8 bits]	Per Run	~ 0.1 nJ		
Vector Cache*	Area	11 ,250 μm²	500 μm²	
16 entries 1024x8-bit [Values are per 1024x8 vector]	Read	~ 0.1 nJ / ~ 0.2	μs ~ 1 nJ / 4 ns	Uses Digital Methods
	Write	~ 0.1 nJ / ~ 0.2	μs ~ 1 nJ / 50 ns	

Conclusion



- Dennard (constant power density) scaling has ceased and Moore's law is slowing
- New paradigms like neuromorphic computing will be required for sub-fJ computing
- We now require a device through system design mentality
 - Motivation behind CrossSim
 - See poster for more detail on CrossSim
- Oxide-based resistive memory offers intriguing device options for both eras
- Novel LISTA and eNode devices, offer significant potential in the development of a low energy neural accelerator
 - See LISTA and eNode posters for more detail on these

Thank you!





Acknowledgements



- This work is funded by Sandia's Laboratory Directed Research and Development as part of the Hardware Acceleration of Adaptive Neural Algorithms Grand Challenge Project
- Many shared ideas among collaborators:
 - Alberto Saleo, Yoeri van de Burgt, Stanford
 - Engin Ipek, U Rochester
 - Dave Mountain, Mark McLean, US Government
 - Stan Williams, John Paul Strachan, HPL
 - Dhireesha Kudithipudi, RIT
 - Jianhua Yang, U Mass
 - Hugh Barnaby, Mike Kozicki, Sheming Yu, ASU
 - Tarek Taha, U Dayton
 - Paul Franzon, NC State University
 - Sayeef Salahuddin, UC Berkeley
 - Dozens of others...
- We are especially interested in collaborations on cross-sim!





Electronic Numerical Integrator And Computer Developed by US Army/U Penn – 1946 150 kW, 357 FLOPs 400 J/FLOP (10 bit)



Where Are we Today?



- Single Unit: Nvidea Tesla P100 GPU
 - Most advanced GPU processor specs, released late 2016
 - Target's deep learning and neural applications
 - 20 TFLOPs 16 bit peak performance w/ peak power dissipation of 300W
 - 70 GFLOPs/watt or about 15 pJ/FLOP (16 bit)
- Supercomputer: Sunway TaihuLight (China)
 - Top supercomputer in the world
 - ShenWei processor
 - 90 PFLOPs peak, 15 MW power
 - 6 GFLOPs/W or about 170 pJ/FLOP
- Need >1000x improvement to tackle internet-scale problems



Basics of Neural Networks



Simple Network: Backpropagation



Another Analogy



Mathematical







TaOx ReRAM in Backprop Training





Data set	# Training Examples	# Test Examples	Network Size
UCI Small Digits[1]	3,823	1,797	64×36×10
File Types[2]	4,501	900	256×512×9
MNIST Large Digits[3]	60,000	10,000	784×300×10

1 µs pulses

ReRAM Measurements

- DC Current-voltage "loops" sweeps are not time-controlled
 - Excessive heating and early wearout
 - Do not provide info on dynamics
- Physical switching < 10ns</p>
- Need pseudo RF setup to measure
 - Ground/signal, conductor backed
 - Agilent B1530 module
 - 10 ns RT/FT, 10 ns PW
 - 1 V nominal, ~140 mV overshoot





Effect of Pulse Width and Edge Time



- Shorter pulses may be employed to lower conductance switching range
- Linearity qualitatively similar across Pulse Width (PW) and Edge Time (ET)
 - Best for SET at 100 ns
 - Best for RESET at 1 us
- Relative conductance change increased with shorter Pulse Width / Edge Time

Nominal Pulse Voltage Values: SET: +1 V RESET: -1 V

Analog Core: Forward Propagation







Analog Core: Back Propagation





