





Analog Machine Learning Accelerators

Damien Querlioz Université Paris-Saclay, CNRS, Palaiseau, France Centre de Nanosciences et de Nanotechnologies damien.querlioz@universite-paris-saclay.fr



How Far Will AI Energy Consumption Grow?



World 🗸 🛛 Business 🗸

Markets ∨ Sustainability ∨ Legal ∨ More ∨

Technology

OpenAl CEO Altman says at Davos future Al depends on energy breakthrough

Reuters

January 16, 2024 6:39 PM GMT+1 · Updated 2 months ago

М	Aa	

Sam Altman, CEO of OpenAI, attends the Asia-Pacific Economic Cooperation (APEC) CEO Summit in San Francisco, California, U.S. November 16, 2023, REUTERS/Carlos Barria/File Photo Purchase Licensing Rights [7]

DAVOS, Switzerland, Jan 16 (Reuters) - OpenAI's CEO Sam Altman on Tuesday said an energy breakthrough is necessary for future artificial intelligence, which will consume vastly more power than people have expected.

Speaking at a Bloomberg event on the sidelines of the World Economic Forum's annual meeting in Davos, Altman said the silver lining is that more climate-friendly sources of energy, particularly nuclear fusion or cheaper solar power and storage, are the way forward for AI.

Or can we make AI *dramatically* more energy efficient?



AI Energy Inefficiency Limits Its Best Applications



- Edge AI implemented in a medical implant could allow epilepsy prediction, advanced BMI...
- But currently limited to elementary AI



The AI Energy Problem Is a *Memory* Problem

- Neural networks: not a lot of arithmetics, but **Huge volume of parameters**
- In computers, GPUs, and AI accelerators developed by industry (Google TPU, Apple NPU...), memory access is extremely costly





Pedram et al , IEEE D&T 2016

The Brain Achieves High Energy Efficiency by Computing « In Memory » in Analog



10,000 times more synapses than neurons!

Brain = a gigantic memory with some computation in the middle No shared arithmetic unit



Integrating Logic & Memory Is a Considerable Challenge

- Only possible technology: Static RAM
 - Even in "5-nanometer CMOS", SRAM bit is 150x150 nanometers
 - Leaky
 - Volatile
- New Memories are coming to the rescue!

Phase Change Memory

CEA-LETI

Memristors, Oxide Resistive Memory

BL

Spin Torque Memory

WL V_{DD}

1 bit of

SRAM

_M₆

BL

Analog Machine Learning Accelerators

> The Promise of Analog In-Memory Computing

- > Analog In-Memory Computing in a System
- Non-Memristor Analog In-Memory Computing

Creative In-Memory Computing

Memristor/RRAM: an Artificial Synapse!

• TiN/HfO_x/Ti/TiN stack

High voltage: move atoms to switch memristor between low/high resistance **Low voltage:** allows reading the resistance

These New Memories Already Have a Market: Microcontroler Applications (NOR Flash replacement)

	-					
	n	f:,	20	~	•	
			IE	UI		
∕.				-	~	

All · Search

Q Newsletter Contact Wi

Products Applications Design Support Community About Infineon Careers

Press General Information Press Releases Market News Press Kits Media Pool Events Contacts

→ Home → About Infineon → Press → Market News → Infineon and TSMC to introduce RRAM technology for automotive AURIX[™] TC4x product family

Infineon and TSMC to introduce RRAM technology for automotive AURIX™ TC4x product family

Nov 25, 2022 | Market News

f in 🎔

Munich, Germany – 25 November, 2022 – Infineon Technologies AG (FSE: IFX / OTCQX: IFNNY) and TSMC today announced the companies are preparing to introduce TSMC's Resistive RAM (RRAM) Non-Volatile Memory (NVM) technology into Infineon's next generation AURIX[™] microcontrollers (MCU).

Embedded Flash microcontrollers have been the main building blocks of automotive electronic control units (ECU) since the introduction of the first engine management systems. They are essential components for clean, safe and smart cars, used in propulsion systems, vehicle dynamics control, driver assistance and body applications. They enable major innovations in the automotive space with regards to electrification, new E/E architectures and automated driving. Currently, the majority of MCU families in the market are based on embedded Flash memory technology. RRAM is a next step in embedded memories that allows to further scale to 28nm and beyond.

The Infineon AURIX TC4x microcontroller products combine performance extension with latest trends in virtualization, security, and networking features to enable the next generation of software-defined vehicles and new E/E architectures. TSMC and Infineon successfully created the basis for introduction of RRAM in the automotive domain. This will put AURIX microcontrollers on a broader Image: Stress stres

ST to add FD-SOI/ePCM 32bit MCU for embedded applications

In H2, ST will sample an STM32 MCU for industrial applications made on an 18nm Samsung FD-SOI process with embedded phase change (ePCM) memory. Volume production is planned for H2 2025.

We repurpose them as *magic memory* for AI!

Memristors Are Very Different From SRAM/DRAM

- Read is just as good as SRAM/DRAM
- Write is slower and write endurance is limited
 - You need to move atoms!
- BER before ECC is typically ~ 10^{-6}

Analog In-Memory Computing Performs Neural Network Inference Very Naturally

A matrix of analog memristors naturally implements a layer of neural network with **Ohm's and Kirchhoff's laws**!

Memristor conductance G = synaptic weight w

In-Principle Energy Efficiency Is Astonishing

• Multiplication $E = U \times i \times t = U^2 \times t / R$ 0.1 V 1 ns 100 k Ω 0.1 fJ

• Accumulation E=0

0.05 fJ / operation

How Do We Measure Efficiency?

Currently, everybody is using

TOPS/W

- = Tera Operations / Second / Watt
- = Tera Operations / Joule

In principle, analog IMC could be 20,000 TOPS/W

How Do We Measure Efficiency?

AI Accelerator Survey and Trends

Albert Reuther, Peter Michaelas, Michael Jones, Vijay Gadepally, Siddharth Samsi, and Jeremy Kepner *MIT Lincoln Laboratory Supercomputing Center* Lexington, MA, USA {reuther.pmichaleas,michael.jones,vijayg,sid,kepner}@ll.mit.edu

(2023 Data)

Analog In-Memory Computing with Memristors

The Challenge of Analog In-Memory Computing: Memristors Act as Random Variables

Esmanhotto et al, 2200145, Advanced Intelligent Systems, 2022

Fighting the Random Character of Memristors: *Program and Verify*

Most Analog IMC Uses Differential Schemes

- Weight is difference between two memristor conductance
- ~3b / Memristor ->4b weight
- Sufficient for almost all neural networks

Memristors Are For <u>Weight Stationary</u> Hardware

- Analog programming operations are long
- Write endurance is limited (e.g., one million)

Scaling Up the Concept

2Mb compute-in memory array (TSMC technology)

1 m	and the same	Link to	der her	 and the second
H				
1				
E				
1				5
H				
9				

UNIVERSITE PARIS-SACLAY

Check for updat

A CMOS-integrated compute-in-memory macro based on resistive random-access memory for AI edge devices

Cheng-Xin Xue^{1,4}, Yen-Cheng Chiu^{1,4}, Ta-Wei Liu¹, Tsung-Yuan Huang¹, Je-Syu Liu¹, Ting-Wei Chang¹, Hui-Yao Kao¹, Jing-Hong Wang¹, Shih-Ying Wei¹, Chun-Ying Lee¹, Sheng-Po Huang¹, Je-Min Hung¹, Shih-Hsih Teng¹, Wei-Chen Wei¹, Yi-Ren Chen¹, Tzu-Hsiang Hsu¹, Yen-Kai Chen¹, Yun-Chen Lo¹, Tai-Hsing Wen¹, Chung-Chuan Lo¹, Ren-Shuo Liu¹, Chih-Cheng Hsieh¹, Kea-Tiong Tang¹, Mon-Shu Ho², Chin-Yi Su³, Chung-Cheng Chou³, Yu-Der Chih³ and Meng-Fan Chang [©]¹⊠

2M Resistive Rams, 1T1R Used as binary memory Very complex periphery These periphery circuits require precise calibration

146 TOPS/W (22nm CMOS)

Xue et al, Nature Electronics, 4, 81 (2021)

A Realization « Truer » to the Original Concept

Article

A compute-in-memory chip based on resistive random-access memory

https://doi.org/10.1038/s41586-022-04992-8 Received: 27 July 2021 Accepted: 17 June 2022 Weier Wan^{1,2}, Rajkumar Kubendran^{2,3}, Clemens Schaefer⁴, Sukru Burc Eryilmaz¹, Wenqiang Zhang⁵, Dabin Wu⁵, Stephen Deiss², Priyanka Raina¹, He Qian⁵, Bin Gao⁵, Siddharth Joshi^{2,4}, Huaqiang Wu⁵, H.-S. Philip Wong¹, & Gert Cauwenberghs²

Nature (2022)

Current mode:

 $I_{out} = \Sigma_i V_i G_i$

MVM output dynamic range varies with models

0 200

Output current (µA)

CNN

LSTM

-200

-400

Voltage mode:

Normalize dynamic range

0

Output voltage (V)

CNN

LSTM

400 -0.2 -0.1

 $\Sigma_i G_i$

0.1

0.2

10-45 TOPS/W in 130nm CMOS

The Samsung Approach

Article A crossbar array of magnetoresistive memory devices for in-memory computing Nature 2022

- Implements Binarized Neural Networks
- MRAM
- Uses Resistance (and not Conductance) as a synaptic weight

The Samsung Approach Requires Complex Periphery Circuitry

 Need for calibration, compensation, and a quite complex analog to digital scheme using time (TDC)

405 TOPS/W in 28 nm CMOS

Analog Machine Learning Accelerators

The Promise of Analog In-Memory Computing

- > Analog In-Memory Computing in a System
- Non-Memristor Analog In-Memory Computing

Creative In-Memory Computing

How to Go From an Array to a Full Neural Network

Fig. 8. Dataflows for DNNs.

In digital ML accelerators, weight stationary is not the most researched option.

Sze et al, Proc. IEEE (2017)

Where to Set the Boundary Between Analog and Digital?

The Path Forward

- Facing diminishing performance improvements in traditional building blocks
- Further gains must come from applicationspecific architectures
- Opportunity: Blur boundaries between analog and digital processing further
 - Don't think analog or digital, think information processing
 - Minimize data conversions, data movement, embrace analog compute
- Not a new idea, but more relevant than ever

Boris Murmann, BIOCAS, 2019

• **My word of caution**: many computer arch. overestimate the cost of ADC. Use ADCs specifically designed for IMC (e.g., CCO)

A Full System with Mostly Analog Routing

SwitchBoard

h

Article

An analog-AI chip for energy-efficient speech recognition and transcription

https://doi.org/10.1038/s41586-023-06337-5 Received: 13 December 2022

S. Ambrogio¹²⁵, P. Narayanan¹, A. Okazaki², A. Fasoli¹, C. Mackin¹, K. Hosokawa², A. Nomura², T. Yasuda², A. Chen¹, A. Friz¹, M. Ishii², J. Luquin¹, Y. Kohda², N. Saulnier³, K. Brew³, S. Choi³, I. Ok³, T. Philip³, V. Chan³, C. Silvestre³, I. Ahsan³, V. Narayanan⁴, H. Tsai¹ & G. W. Burr¹

Nature 2023

Full system Mostly analog routing

C

• 35M phase change memories

27 CORS UNIVERSITE

A Full System with Mostly Digital Routing

nature electronics

Article

https://doi.org/10.1038/s41928-023-01010-1

A 64-core mixed-signal in-memory compute chip based on phase-change memory for deep neural network inference

Received: 27 May 2023	Manuel Le Gallo 🕲 ¹⁶ 🖂, Riduan Khaddam-Aljameh ^{1,6} , Milos Stanisavljevic ^{1,6} ,	
Accepted: 10 July 2023	Athanasios Vasilopoulos © ¹ , Benedikt Kersting ¹ , Martino Dazzi ¹ , Geethan Karunaratne © ¹ , Matthias Brändli ¹ , Abhairai Singh ¹ , Silvia M, Müller ² ,	
Published online: 10 August 2023	Julian Büchel ¹ , Xavier Timoneda ¹ , Vinay Joshi ¹ , Malte J. Rasch ³ , Urs Egger ¹ ,	
Check for updates	Angelo Garofalo ¹ , Anastasios Petropoulos ⁴ , Theodore Antonakopoulos ⁴ , Kevin Brew ⁶ , Samuel Choi ⁵ , Injo Ok ⁵ , Timothy Philip ⁵ , Victor Chan ⁵ ,	
	Claire Silvestre ⁵ , Ishtiaq Ahsan ⁵ , Nicole Saulnier ⁵ , Vijay Narayanan ³ ,	

Nature Electronics 2023

Pier Andrea Francese¹, Evangelos Eleftheriou¹ & Abu Sebastian ¹

Full system Mostly digital routing (massive ADC/DAC)

10 TOPS/W

« But the TOPS/W Are Not That High! »

 These are early prototypes, optimized for functionality not for perf

TOPS/W Can Be a Very Misleading Unit

- « Operation » is ill-defined
- These are for peak conditions, which are rarely reached in practice
- The difference between real-life and peak conditions varies extensively depending on the different types of hardware

EETimes

HOME NEWS V PERSPECTIVES DESIGNLINES V PODCAST EDUCATION V STORE

DESIGNLINES | AI & BIG DATA DESIGNLINE

TOPS: The Truth Behind a Deep Learning Lie

By Ludovic Larzul, Mipsology 06.25.2021 🔲 0

Jan Werth Apr 26, 2021 · 11 min read · 💿 Listen Y () in Ø 🖓

When "TOPS" are Misleading

Neural accelerators are often characterized with the performance feature "TOPS" — Trillion operations per second. But that alone is not enough. It is important to know how these accelerators work and what else should be considered when making a comparison.

AUTO, SECURITY & PERVASIVE COMPUTING

Lies, Damn Lies, And TOPS/Watt

249 f 20 ¥ 7 in 20 ≺ Shares Questions you need to ask to make sure you understand AI hardware performance.

JANUARY 7TH, 2019 - BY: GEOFF TATE

<u>.</u>...

 ${f T}$ here are almost a dozen vendors promoting inferencing IP. but none of them

Aron Kirschen Aug 12, 2020 · 10 min read · O Listen

Y G 🖬 🖉 🖓

Why TOPS/W is a bad unit to benchmark nextgen AI chips

Real Energy Consumption Can Be Dramatically Higher than Peak

NVIDIA A100 GPU ResNet50

Courtesy of David Novo (Univ Montpellier/CNRS)

How I View Architecture Research

How I View Architecture Research

Analog Machine Learning Accelerators

> The Promise of Analog In-Memory Computing

Analog In-Memory Computing in a System

Non-Memristor Analog In-Memory Computing \triangleright

Creative In-Memory Computing

From Unconventional to More Conventional?

- For cultural reason, many analog IMC ideas are developed for emerging technology
- There's a movement to adapt these ideas to more established technology

IMC with Flash Memory

NAND: Mainstream data storage (SSD)
 Page (16KB*N) read/write (slower, 10-100us)
 Big block (> 4MB) erase (GC, write amplification)
 Rely on controller managements. Intensive ECC
 3D charge-trapping memories (>300 layers), with MLC/TLC/QLC

- Byte addressable; Small-unit P/E
- No controller needed. Almost no ECC
- 2D FG stops at 45nm node → Going to 3D charge-trapping NOR
- NAND flash constrains are challenging for IMC
- NOR Flash memory can be used for IMC, but does not scale to advanced CMOS

Prospects of Computing In or Near Flash Memories

Hang-Ting Lue, Chun-Hsiung Hung, Keh-Chung Wang and Chih-Yuan Lu

Macronix International Co., Ltd 16 Li-Hsin Road, Hsinchu Science Park, Hsinchu, Taiwan. E-mail: <u>htlue@mxic.com.tw</u>

IEDM 2024

Toward 3D IMC?

Prospects of Computing In or Near Flash Memories

Hang-Ting Lue, Chun-Hsiung Hung, Keh-Chung Wang and Chih-Yuan Lu

Macronix International Co., Ltd 16 Li-Hsin Road, Hsinchu Science Park, Hsinchu, Taiwan. E-mail: <u>htlue@mxic.com.tw</u>

IEDM 2024

NAND: Mainstream data storage (SSD)

WL

WL

WL30

WL3

GSL

- Page (16KB*N) read/write (slower, 10-100us)
- Big block (> 4MB) erase (GC, write amplification)
- Rely on controller managements. Intensive ECC
- 3D charge-trapping memories (>300 layers), with MLC/TLC/QLC

- NOR: classic legacy for code Flash
- Random access read (faster, ~100ns)
- Byte addressable; Small-unit P/E
- No controller needed. Almost no ECC
- 2D FG stops at 45nm node → Going to 3D charge-trapping NOR

Figure 4 Current 3D NOR adopts CMOS under Array (CuA) process integration. In the future, we will go for Cu hybrid bonding to connect separated CMOS chip with 3D array (like 3D NAND).

SRAM-Based IMC

- SRAM can be used for both digital and analog IMC
- Most exciting analog designs rely on switched cap principles

Fully Row/Column-Parallel In-memory Computing SRAM Macro employing Capacitor-based Mixed-signal Computation with 5-b Inputs

Jinseok Lee, Hossein Valavi, Yinqi Tang and Naveen Verma Princeton University, Princeton, NJ, USA (jinseokl@princeton.edu)

VLSI 2021

• Positioning: edge or HPC?

Memory-Centric Computing: Recent Advances in Processing-in-DRAM (Invited)

Onur Mutlu, Ataberk Olgun, Geraldo F. Oliveira, Ismail E. Yuksel ETH Zürich

IEDM 2024

DRAM-Based IMC

Fig. 1: An example of performing the MAJority-of-three operation (i.e., MAJ3 (A, B, C)) (a) and the NOT operation (i.e., dst=NOT (src)) in Ambit [111]. In (a), we focus on DRAM cell and sense amplifier operations (**①**). Initially, cells A, B, C, and bitline have voltage levels of GND, VDD, VDD, and VDD/2, respectively (**①**). We first perform a triple-row activation (TRA) to simultaneously activate cells A, B, and C (**①**). When the wordlines of all three cells are raised simultaneously charge sharing

 DRAM based-IMC is possible but destructive read, and small signals -> most adapted for digital

Analog Machine Learning Accelerators

> The Promise of Analog In-Memory Computing

- > Analog In-Memory Computing in a System
- Non-Memristor Analog In-Memory Computing

Creative In-Memory Computing

Memristor Imperfections Can Naturally Produce a *Bayesian* Neural Network!

In Bayesian models, everything is a random variable that follows specific probability distributions Memristors actually act as a random variable that follow specific probability distributions!

Our concept: Bayesian models can be a "better" way to exploit memristors

Bonnet et al, Nature Communications 14, 7530 (2023)

Bayesian Neural Networks

Assuming NN was trained to recognize « cats » and « dogs »

Memristor-Based Bayesian Neural Networks

We program 50 memristor-neural networks (each with two layers). We apply same input to them

Bonnet et al, Nature Communications 14, 7530 (2023)

We get 50 outputs:

their dispersion tells about the *certainty* of the neural network

Neural network trained with Variational Inference incorporating a specific "technological loss"

(CNrs)

43

Memristor-Based Bayesian Neural Networks

We program 50 memristor-neural networks. We apply same input to them

Bonnet et al, Nature Communications 14, 7530 (2023)

Memristor-Based Bayesian Neural Networks

We program 50 memristor arrays, and we apply the same input to them

Bonnet et al, Nature Communications 14, 7530 (2023)

Fully Experimental Arrythmia Recogniton

Traditional Neural Network

80% classification accuracy

Simulation Conventional NN (float32)

Bayesian Neural Network

79% classification accuracy

Easily recognizes unknown types of arrythmia

Bonnet et al, Nature Communications 14, 7530 (2023)

46 CNTS UNIVERSITE

The Grand Challenge of On-Chip Learning

- Current in-memory computing AI accelerators are focused on inference, which makes sense
- On-chip learning/adaptation would also have tremendous prospects

Backpropagation, the Canonical Method for Training Networks, Is Not Adapted for In-Memory Computing

 $C = -\ln a$

$$\frac{\partial C}{\partial w_i} = (a-1)h_i$$
$$\frac{\partial C}{\partial v_{ij}} = (a-1)w_ih_i(1-h_i)x_j$$
$$\frac{\partial C}{\partial u_{kl}} = (a-1)\sum_i w_ih_i(1-h_i)v_{kl}g_k(1-g_k)x_l$$

To update weight u_{kl}, you need calculation that involves information about the *whole* network!

Learning in the Brain Is Mysterious

- Synapses are physical objects that have only access to local information
- An old idea (Hebb): « cells that fire together wire together »

 Multiple theoretical works show that Hebbian learning rules (e.g., STDP) are not as powerful as backprogation

EqProp: a Type of Neural Networks Grounded in Physics and with « Brain-like » Learning

Learns with Hebbian-like learning but is equivalent to backpropagation

Scellier & Bengio, fnins 2017

Energy $E = \frac{1}{2} \sum_{i} s_{i}^{2} - \frac{1}{2} \sum_{i,j} w_{ij} \rho(s_{i}) \rho(s_{j})$

 s_i : neuron states, ρ : sigmoid function

The neural network is a dynamical system that goes naturally toward its energy minimum

Local: depends only on s_i nearest-neighbors Similar to leaky integrate-and-fire neuron

50

When the network has converged, we get the output a

If the Output a Is Not the Right Value, EqProp then "Nudges" the Network Toward It

Scellier & Bengio, fnins 2017

Error $C = \frac{1}{2}(a-1)^2$

Nudged Energy $F = E + \beta C$

$$\frac{da}{dt} + a = \left[\sum_{j} w_{ij}\rho(s_j)\right]\rho'(a) + \beta(a-1)$$
« Nudging »

The perturbation of the output *a* propagates throughout the network that reaches a new equilibrium

S_{nudged}

Equilibrium Propagation

• Change the weight w_{ij} by

•
$$\rho(s_{nudged,j})\rho(s_{nudged,i}) - \rho(s_{free,j})\rho(s_{free,i})$$

➢ If neurons i and j had MORE Hebbian correlation during the nudged phase, then INCREASE their connection w_{ii}

 \succ Otherwise, DECREASE w_{ij}

The Local Learning Rule of Equilibrium Propagation Leads to High Recognition Rates

Mathematical equivalence of EP gradients with Backpropagation Through Time

M. Ernoult, J. Grollier, D. Querlioz, Y. Bengio, B. Scellier, NeurIPS (2019)

Collaboration

EP scales to CIFAR-10

A. Laborieux, M. Ernoult, B. Scellier, Y. Bengio, J. Grollier, D. Querlioz, fnins 15, 129 (2021)

EP trains Binary Neural Networks

- Binary synapses (CIFAR 10)
- Binary synapses and neurons (MNIST)
- Ternary gradients

J. Laydevant, M. Ernoult, D. Querlioz, J. Grollier, CVPR (2021)

EqProp Can Train a Quantum Ising Machine

• EqProp adapts very naturally to physical systems that compute

Laydevant, Markovic and Grollier, Nature Communications 15, 3671 (2024)

Conclusion

- Analog in-memory computing: huge potential to reinvent electronics for cognitivetype tasks and AI
- Memristors and related devices are rich devices that can be used in a variety of ways
- They are particulary adapted for Bayesian approaches, paving the way toward trustworthy AI
- Next ARCHI grand challenges: scaling and solving learning
- Very important & creative time for micro/nano-electronics research.
 Considerable benefits from algorithm/electronics/technology research

Acknowledgments

- Kamel-Eddine Harabi
- Clément Turck
- Tifenn Hirtzlin
- Atreya Majumdar
- Marie Drouhin
- Jacques-Olivier Klein
- Maxence Ernoult
- Axel Laborieux
- Adrien Renaudineau
- Théo Ballet

- Elisa Vianello
- Tifenn Hirtzlin
- Eduardo Esmanhotto
- Djohan Bonnet
- Niccolo Castellani
- François Andrieu

- David Novo
- Pascal Benoit
- Paul Delestrac
- Aymen Romdhane
- Bruno Lovison-Franco

- Julie Grollier
- Jérémie Laydevant
- Maxence Ernoult
- Marie Drouhin

- Yoshua Bengio
- Benjamin Scellier

• Jean-Michel Portal • Jean-Pierre Walder

• Marc Bocquet

٠

Fadi Jebali

- Eloi Muhr
- Mathieu-Coumba Faye

HOME ABOUT COMMITTEE

57

Cross-disciplinary Conference on Memory-Centric Computing (CCMCC)

OCTOBER 8-10, 2025

DRESDEN, GERMANY

Thank you for your attention!

Postdoc positions available

damien.querlioz@universite-paris-saclay.fr https://sites.google.com/site/damienquerlioz/

