

NARESH SHANBHAG

UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN



MRAM-BASED DEEP IN-MEMORY ARCHITECTURES

This research was developed with funding from the Defense Advanced Research Projects Agency (DARPA). The views, opinions and/or findings expressed are those of the author and should not be interpreted as representing the official views or policies of the Department of Defense or the U.S. Government.

ΤΗΕ ΤΕΑΜ

Princeton

University of Illinois at Urbana-Champaign

Raytheon MS

GLOBALFOUNDRIES

Mike Burkland Ajey Jacob Pavan Hanumolu Naveen Verma Naresh Shanbhag (UIUC) (UIUC) (Princeton) [DoD systems] [devices, foundry] [systems & ckts] [mixed-signal ICs] [systems & ckts] S. Soss, D. L.-Y. Chen J. Broussard A. Patil, H. Hua, S. Gonugondla, K.-H. Kim P. Deaville Brown L. Suantak N. Gaul, B. Paul, B. Zhang J. Goulding B. Lanza S. Cole T. Gangwer R. Kelley C. Contreras Distribution A. Approved for public release; distribution is unlimited

REALIZING ARTIFICIAL INTELLIGENCE @ THE EDGE

Al in the Cloud



[source: pexels.com]

- client-server model
- efficiency challenged
- security + privacy issues



THE DATA MOVEMENT PROBLEM



fundamental question

how do we design **intelligent autonomous machines** that **operate at the limits** of **energy-delay-accuracy**?



Systems on Nanoscale Information fabriCs

www.sonic-center.org

[2013-'17]



(STARnet Program by Semiconductor Research Corporation & DARPA)

Shannon-Inspired Statistical Computing for the Nanoscale Era

By NARESH R. SHANBHAG^(D), *Fellow IEEE*, NAVEEN VERMA, *Member IEEE*, YONGJUNE KIM^(D), *Member IEEE*, AMEYA D. PATIL, *Student Member IEEE*, AND LAV R. VARSHNEY^(D), *Senior Member IEEE*

Proceedings of IEEE, Special Issue on *non von Neumann Computing*, January 2019. Distribution A. Approved for public release; distribution is unlimited

SHANNON-INSPIRED MODEL OF COMPUTING



1 use information-based metrics e.g., mutual information $I(Y_o; \hat{Y})$

- 2 design low SNR fabrics, e.g., deep in-memory architecture (DIMA)
- 3 develop statistical error-compensation (SEC) techniques

FRANC OBJECTIVE



THE DEEP IN-MEMORY ARCHITECTURE (DIMA)



SRAM DIMA PROTOTYPES

100X EDP reduction over von Neumann equivalent* @ iso-accuracy

* 8b fixed-function digital architecture with identical SRAM size



Multi-functional inference processor (65nm CMOS) Random forest processor (65nm CMOS)





On-chip training processor (65nm CMOS) Fully (128) row-parallel compute (130nm CMOS)

[Feb. JSSC'18]

[ESSCIRC'17, JSSC July'18]

[ISSCC'18, JSSC Nov.'18]

[VLSI'16, JSSC'17]

MRAM OPPORTUNITIES & CHALLENGES





MTJ resistance distribution



Opportunities:

 $\begin{array}{l} \mbox{high density} \rightarrow \mbox{high on-chip compute density} \\ \mbox{non-volatility} \rightarrow \mbox{ultra low-power duty-cycled operation} \\ \mbox{in production} \rightarrow \mbox{enables system prototyping \&} \\ \mbox{reduces time to DoD availability} \end{array}$

Challenges:

high conductance \rightarrow large array currents small TMR \rightarrow high sensitivity readout huge $R_{MTJ} - R_{MOS}$ limits cell compute models high cell density \rightarrow severe pitch-matching const. high MTJ process variation restricts bit-line SNR

MRAM-BASED DEEP IN-MEMORY ARCHITECTURES

four DIMAs proposed – two selected for prototyping



- $M \times N$ single-shot matrix vector multiply
- Functional read: binary dot products on BLs
- Compute cell: 2T-2MTJ (1b XNOR/AND)
- BL voltage sensing via 4-bit SAR ADCs
- SNR vs. energy trade-off due to sensing circuits Distribution A. Approved for put

- $M \times N$ single-shot matrix vector multiply
- Functional read: 5-bit × 4-bit dot products on SLs
- Compute cell: 1T-1MTJ (4 × 1 multiplication)
- SL current sensing via time-based 5-bit ADCs
- e to sensing circuits SNR vs. energy trade-off due to sensing circuits Distribution A. Approved for public release; distribution is unlimited

EDP GAINS OVER VON NEUMANN EQUIVALENT



SHANNON-INSPIRED COMPUTE MODELS



- relaxes MRAM-based DIMA's compute SNR requirements → enables substantial energy savings
- two methods: stochastic data-driven hardware resilience & coded DIMA

STOCHASTIC DATA-DRIVEN HARDWARE RESILIENCE



CODED DIMA - ENHANCING DIMA'S COMPUTE SNR

MTJ variation aware coding enables HIGH system SNR in spite of LOW circuit SNR



Distribution A. Approved for public release; distribution is unlimited

CURRENT STATUS

- verified & quantified RMS's mission capability enabled by MRAM-based DIMA
- MRAM-based DIMA demonstrates > 2 orders-of-magnitude EDP reduction

DIMA prototype designs taped out



MRAM-DIMA	
circuit_char_blk (1x1)	

- (2019) [single-bank] validate EDP gains & calibrate models
- (2020) [multi-bank] enable scaling and system integration

Technology transition via SRAM-based DIMA



[H. Jia, arXiv:1811.04047, Princeton]

- DIMA technology transition to RMS for algorithm prototyping & evaluation
- facilitates future MRAM-based DIMA integration by RMS

NEXT STEPS & SUMMARY

[Srivastava, et al., ISCA'18] decision A/D & RDL Cross BL processor (CBLP) Improvement 6T SRAM bit cel BLP BLP BLP I-DAC VBIAS,O X[1] Factor | WL Self-biasing Bit cell bit-cell Replica WL RESET upsized to drive WL capacitace Support Template K-Nearest Col. mux Matching Neighbour filtering NeighbourRegression Vector Neural 12 12 11 Network mux & buff Energy 4-Bank Throughput Energy *≰ K*-bus

Multi-functioned DIMAs for Diverse Applications

Parameterized DIMAs for Platform-design Tools

Summary

Mission: To realize > 200X in EDP gains in DoD workloads by integrating **MRAM device** within **Deep In-memory Architectures** using **Shannon-inspired compute models**

FRANC Program: "to provide the foundation for new materials technology and new integration approaches to be exploited in pursuit of novel compute architectures"

ERI ELECTRONICS RESURGENCE INITIATIVE SUMMIT

2019 | Detroit, MI | July 15 - 17

