The Future of

Hardware Technologies for Computing

N3XT 3D MOSAIC, Illusion Scaleup, Co-design

Subhasish Mitra



Department of EE and Department of CS

Stanford University

Thanks: Students, Sponsors, Collaborators



Subhasish Mitra

Relays, Vacuum Tubes, Discrete Transistors, ICS

TURNING POTENTIAL INTO REALITIES: THE INVENTION OF THE INTEGRATED CIRCUIT

[Kilby Nobel Lecture 2000]

"yields ... too low to be profitable"

"best [devices] ... not made with semiconductors"

"elegant devices messed up with all the other stuff"

Abundant-Data Computing

Many walls simultaneously



Nano<u>Systems</u>



Data Explosion & Memory Wall

Both Von Neumann <u>&</u> non-Von Neumann architectures

"Ideally ... desire an indefinitely large memory capacity such that any particular ... word would be immediately available. ... It does not seem possible physically to achieve such a capacity. We are therefore forced ... hierarchy of memories"

[Burks, Goldstine, Von Neumann, 1946]

Compute-in-X

Computation immersed in memory (& sensors)

Computing Today





N3XT 3D MOSAIC

MOnolithic / Stacked / Assembled IC



N3XT 3D MOSAIC

MOnolithic / Stacked / Assembled IC



N3XT 3D \supset 3D Folding



N3XT 3D: Many Technologies



Subhasish Mitra

Carbon Nanotube FETs (CNFETs): Many Innovations



Resistive RAM (RRAM)



Subhasish Mitra

Ultra-dense Monolithic 3D



Lab NanoSystems Hardware



Lab to Fab

Analog Devices

SkyWater Tech. Foundry



3D NanoSystem



[Shulaker Nature 17]

HD Computing: Brain-Inspired ⊃ Neural Nets



Non-volatile Computing System: RRAM + Silicon CMOS



Non-volatile Computing System: RRAM + Silicon CMOS



[Wu ISSCC 19]

Multiple bits-per-cell RRAM System



$1TnR \times Multiple bits-per-cell RRAM$



[Hsieh IEDM 19, IEEE EDL 21]

CHIMERA: RRAM Edge AI Inference & Training



CHIMERA: RRAM All Memory On-chip Dataflow



New dataflow: large benefits



N3XT 3D MOSAIC

MOnolithic / Stacked / Assembled IC



Dream: All Memory + Compute On-chip



Off-Chip Memory Accesses Costly

Large EDP overheads vs. Dream Chip



Illusion System



Illusion Ideal for AI

Illusion \approx Dream 1.1 \times Dream EDP

Illusion Energy ≤ 1.05 ×

Dream Energy

Illusion Exec. Time ≤ 1.05 ×

Dream Exec. Time

(measured for AI inference)

[Giordano Symp. VLSI Circuits 21, Radway Nature Electronics 21]

Subhasish Mitra

Hardware-proven backed by theory



6-CHIMERA chip Illusion system 8-chip Illusion system

Many Illusion Mappings: Energy/Exec. Time Tradeoffs

Illusion system using 6 CHIMERA chips: 12 MByte ResNet-18



Both mappings ≤ 1.1× Dream EDP

[Giordano Symp. VLSI Circuits 21]

CHIMERA Illusion Benefits vs. Off-chip Flash

12 MByte ResNet-18 Edge AI: 1 inference per second

(e.g., security, activity monitoring, healthcare, wildlife)

6-CHIMERA chip Illusion System



6 CHIMERA chips

12 MBytes RRAM overall

Minimal inter-chip messages





Based on [Giordano, Symp. VLSI Circuits 2021], off-chip Flash: 0.6 GBytes/s, 2 nJ/Byte

Illusion ⊃ Traditional Parallel Computing



Traditional parallel: 10 MByte chip-to-chip messages

(12 MByte RestNet-18)

Traditional parallel: e.g., [Zimmer Symp. VLSI Circuits 19, Shao MICRO 19]



CHIMERA Illusion Benefits vs. Traditional Parallel

Illusion: [Giordano, Symp. VLSI Circuits 2021], traditional parallel: [Shao MICRO 2019]

Illusion Ideal for AI

Wide Variety of AI tasks

CNNs, DLRMs, LSTMs, Transformers, single/multiple inferences, training, ...



1,024 × Workload Growth





Quadratically reduce Illusion total message cost

Illusion Scaleup

Maintain 1.1× Dream EDP despite growing Dream Chips



Illusion Scaleup

Illusion Scaleup is Fungible



Many NanoSystems Opportunities

Co-design: device + circuit + arch. + algorithm

□ Dense compute + thermal

□ New software optimizations

Co-design Examples

multiple abstraction layers cooperate for large benefits

- □ **Illusion scaleup**: tech + circuit + arch + algorithm
- **RRAM edge AI training**: tech + arch + algorithm
- **Efficient hyperdimensional computing**: tech + circuit + arch + algorithm
- **RRAM all memory on-chip dataflow:** tech + arch
- □ **1TnR + multiple bits-per-cell RRAM**: tech + circuit
- **Imperfection-immune carbon nanotube VLSI:** tech + circuit



RRAM Edge AI Incremental Training Challenging

New Low-Rank Training (LRT)

| 16-bit Chip To Chip (C2C) Image: Read of the second | | LRT hardware results: iso-accurate vs. Stochastic Gradient Descent (SGD) |
|---|--------------------------------|---|
| | RRAM weight update steps | 101× fewer vs. SGD |
| | Energy Delay Product | 340× better vs. SGD |
| | Endurance (20 samples/min.) | 10 years (LRT + ENDURER*) vs. 2 weeks (SGD) |
| | On-chip SRAM capacity | 37× smaller vs. SGD |

[Giordano Symp. VLSI Circuits 21] *ENDURER: [Aly Proc. IEEE 19, Wu ISSCC 19]

Conclusion

NanoSystems today

Industrial fabs: Carbon nanotube FETs + RRAM + monolithic 3D

□ *N3XT 3D MOSAIC* + Illusion scaleup key

- Computation immersed in memory
- Large benefits over growing problem sizes, ideal for AI

Co-design of the "right" kind

Big opportunities for NanoSystems