

Stochastic Dataflow Computing Wojciech Romaszkan, Tianmu Li, Jiyue Yang, Rahul Garg, Albert Li, Di Wu, Kang Wang, Sudhakar Pamarti, Puneet Gupta, University of California, Los Angeles Foundation Required for Novel Compute (FRANC) Artificial Intelligence

Background

Al in Edge Devices

- Edge AI is highly desirable due to lower latency and improved privacy and security
- Edge AI should support a variety of applications under tight, often time-varying, energy constraints
- Memory Bottleneck: Currently, large AI models incur high memory access energy and latency costs



Stochastic Computing (SC)

- Numbers represented as average of random binary streams
 - Compact multiply-and-accumulate (MAC) units enable massive parallelization
 - Longer streams improve compute accuracy
- SC is perfect for edge AI
 - Faster, smaller, and more energy efficient
 - Enables runtime adaptation of energy, latency, and accuracy



Challenges

- How to deliver high AI inference accuracy without sacrificing the parallelization benefits of SC?
- How to maintain scalability/programmability without sacrificing energy efficiency ?
- How to scale memory bandwidth to match SC's compute parallelism? Conventional memory is not dense enough.

Approach

Improving the Accuracy of SC Inference

- SC-aware AI training accounts for SC induced errors completely closes the SC-fixed point accuracy gap
- Complemented by controlled stream randomness, hybrid SCfixed point addition and other techniques



Scalable, Programable SC Architecture

- Fully digital, programmable architecture (with an instruction set)
- Massive compute parallelization:
- MAC/mm²: ~96k (SC) vs ~0.5k (fixed-point) Computation skipping based stochastic pooling to save 4X-9X
- energy/latency with no accuracy loss
- Run-time programmable accuracy vs energy/latency



Authors' Ow

Coupling SC with Magneto-Electric RAM

- Voltage controlled magnetic tunneling junction (VC-MTJ) based memory as cheap, high density, non-volatile, on-chip storage 5x smaller and 3x less energy than other on-chip memories
- End of project goal: full integration of SC and VC-MTJ based onchip storage



CMOS 14nm prototypes currently in the foundry



