Exploring Energy Efficient Analog Neural Network Accelerator Designs

Joel Emer & Vivienne Sze

October 2022

Typical Deep Neural Network Accelerator Organization



Objective: Achieve high energy efficiency and processing speed

Design Space - Digital + Dense

Dense Dense Dense Image Scratch Pad GRS (N) (12x16b REG) Instructions 2-stage pipelined Zero Accumulate Input Psum Buffer NFU-2 NFU-. multiplier Ţ→ŚĦŢ≠⊗ŀ Scratch Pad (225x16b SRA ≦ Inst. Partial Sum RISC-V Scratch Pad PE (24x16b REG) GRS (W) GRS (S) GPIO Eyeriss [JSSC2017] Simba [MICRO2019] 🌟 DianNao [ASPLOS2014]

Digital-Compute Accelerator Designs



Classification of Important Design Aspects

Dataflow

Processing schedule for uncompressed data movement and compute in time and space

Architecture

Detailed characterization of hardware, e.g., storage sizes and access and compute widths

Technology

Assessment of costs at device level, e.g., area, energy, latency

Mapping Choices

Energy-efficiency of peak-perf mappings of a single architecture/problem



480,000 mappings shown

Spread: 19x in energy efficiency

Only 1 is optimal, 9 others within 1%

A model needs a mapper to evaluate a DNN workload on an architecture

6,582 mappings have min. DRAM accesses but vary 11x in energy efficiency

A mapper needs a good cost model to find an optimal mapping

Source: Parashar, Timeloop

Timeloop: Analytical Modeling for Tensor Accelerators



Processing In Memory (PIM)

Activation is input voltage (V_i) Weight is resistor conductance (G_i)



- Reduce data movement by moving compute into memory
- Compute with memory storage elements

Analog Compute

- Activations, weights and/or partial sums are encoded with analog voltage, current, or resistance
- Increased sensitivity to circuit non-idealities
- A/D and D/A circuits to interface with digital domain
- Leverage emerging memory device technology

eNVM:[Yu, PIEEE 2018], SRAM:[Verma, SSCS 2019]

Conventional vs Processing in Memory (PIM)



Design Space + Analog/Dense



Digital-Compute Accelerator Designs

Analog-Compute Accelerator Designs



Pieces to Evaluate a PIM Design

- PIM Hardware Components
 - Components determine the basic characteristics of the system
 - Choice of ADC, DAC, memristors
- System Architecture
 - Architecture design sets data movement and reuse opportunities
 - Choices of crossbar dimensions, buffer hierarchy, data movement
- Workload
 - Workload characteristics set interaction with the architecture
- Mapping
 - Workloads can be mapped to architectures in many ways, each with energy, throughput, and utilization tradeoffs.



Timeloop: Analytical Modeling for Tensor Accelerators

Workload



Timeloop + Memristive Crossbars



Exploiting Sparsity



Design Space + Digital Sparse



Digital-Compute Accelerator Designs

Analog-Compute Accelerator Designs



Classification of Important Design Aspects

Dataflow

Processing schedule for uncompressed data movement and compute in time and space

Sparse Accelerations

Impact of sparse acceleration features on workloads with sparse data

Architecture

Detailed characterization of hardware, e.g., storage sizes and access and compute widths

Technology

Assessment of costs at device level, e.g., area, energy, latency

Sparseloop: Analytical Modeling for Sparse Tensor Accelerators



format

gating

skipping

Available at <u>sparseloop.mit.edu</u> (MICRO 2022 Distinguished Artifact Award: Nellie Wu)

Design Space + Analog/Sparse



17

Timeloop + Memristive Crossbars



Conclusion

- Timeloop+Accelergy is a model:
 - that generates performance and energy projections...
 - allowing analysis of design-space tradeoffs...
 - for a wide range of architectures, including...
 - for designs with digital and analog components...
 - as well as for designs that exploit sparsity...
 - that is fast enough to search for optimal mappings...
 - and should facilitate design-space searches...

Questions