

The Edward S. Rogers Sr. Department of Electrical & Computer Engineering UNIVERSITY OF TORONTO

Abstract

Laconic:

Deep Neural Network inference accelerator

- •Targets primarily CNNs, can do any layer type •Exploits effectual bit content for activations and weights
- •Term-serial accelerator

20.5x faster than data-parallel accelerators **2.6x** more energy efficient **26%** area cost

Motivation

- ineffectual accelerators exploit • Previous computations including zero skipping, precision variability, and zero bit skipping to improve the performance and energy efficiency of CNNs.
- •We show that policy "At+Wt", eliminating the ineffectual computations at the bit level for both the activations and the weights, has the highest potential speedup.



Conventional bit-parallel accelerator Execution time does not vary with data bit content

activations

weights

Laconic: Term-Serial Accelerator

activations

weights



Laconic Processing Element

- Inputs:
- 16 Activations, 1 term/cycle
- 16 Weights, 1 term/cycle
- Multiplication is done term-serially
- Reduces the products to a single output
- Supports cascading for smaller layers



Laconic Deep Learning Computing

Sayeh Sharify, Mostafa Mahmoud, Alberto Delmas Lascorz, Milos Nikolic, Andreas Moshovos

Baseline



- More area and energy efficient adder tree
- Divides the output of histogram stage to 6 groups, G0, G1, ..., G5.
- Outputs within the same group have no overlapping bits that are "1".
- Concatenates the outputs within a group to calculate their sum.



Laconic Architecture







Enhanced Adder Tree

 $G_0 = \{N^{30}, N^{24}, N^{18}, N^{12}, N^6, N^0\}$ $G_1 = \{N^{31}, N^{25}, N^{19}, N^{13}, N^7, N^1\}$ $G_2 = \{N^{26}, N^{20}, N^{14}, N^8, N^2\}$ $G_3 = \{N^{27}, N^{21}, N^{15}, N^9, N^3\}$ $G_4 = \{N^{28}, N^{22}, N^{16}, N^{10}, N^4\}$ $G_5 = \{N^{29}, N^{23}, N^{17}, N^{11}, N^5\}$

Datapath Layout



- 65nm TSMC
- 128W Tile:

- **752 GOPS**

Results Relative to Baseline

Throughput:



Energy efficiency:



S: Sparse





