Bit Prudent In-Cache Acceleration of Deep Convolutional Neural Networks

Xiaowei Wang, Jiecao Yu, Charles Augustine, Ravi Iyer, Reetuparna Das

> M-Bits Research Group UNIVERSITY OF MICHIGAN





Background











A Convolutional Layer

3D Filters (M)

each filter: C channels each channel: RxS weights Input Activations (C channels)

Output Activations (M channels)





Mapping CNN to Neural Cache



Mapping CNN to Neural Cache





Mapping CNN to Neural Cache





11 Way 2 Way 3 Way 20

Put it together



① Filter Loading



2 Input Loading











Motivation for This Work



Sparse Models – Hard to gain speedup



Unable to skip computation in SIMD model

Outline

Background

- Sparsity-Aware Architecture
 - Coalescing channels
 - Overlapping filters
- Low-Precision Architecture
- Evaluation Methodology
- Results

Architecture for Sparsity: Coalescing



Dynamic Input Coalescing



Coalescing Unit Architecture – 8 channel example



Coalescing Unit Architecture – 8 channel example



Input Loading for Coalescing



Input-loading-aware Structured Pruning



Input-loading-aware Structured Pruning

$$\sum_{r=1}^{R} \sum_{s=1}^{S} \sum_{m=1}^{M} (w_{r,s,m,c=1})^2$$

M=4 C=256

Channel masks are the same across all arrays



Outline

- Background
- Sparsity-Aware Architecture
 - Coalescing channels
 - Overlapping filters
- Low-Precision Architecture
- Evaluation Methodology
- Results

Architecture for Sparsity: Overlapping



Architecture for Sparsity: Overlapping



Outline

- Background
- Sparsity-Aware Architecture
- Low-Precision Architecture
- Evaluation Methodology
- Results



Step 1: Multiply Zero Weights



Step 2: Calculate 1's complement



Step 3: Partial sum accumulation (2's complement conversion)



MAC cycle goes down with Low-Precision Arch.

Bit-width Configuration	MAC cycles	
8-bit weight/	94 12	
Ternary weight/		
4-bit activation		
Binary weight/		
4-bit activation	0	

Evaluation Methodology

CNN Models

- AlexNet and Inception v3 convolutional layers
- Performance-aware pruning
- 8-bit quantization for weights and inputs

Training Methods

 modified TensorFlow tool-chain

Configurations

- CPU-base
- GPU-base
- Neural \$
- SPR-coal/SPR-olap
- LPR-2b/LPR-1b

		CPU (2 sockets)	GPU (1 card)	Neural Cache
	Processor	Intel Xeon E5- 2597 v3, 2.6GHz, 28 cores, 56 threads	Nvidia Titan Xp, 1.6GHz, 3840 CUDA cores	2.5GHz Compute SRAM, 1032192 Bit-serial ALUs
	On-chip memory	78.96 MB	9.14 MB	35 MB
	Off-chip memory	64 GB DRAM	12 GB DRAM	64 GB DRAM
,	Profiler / Simulator (Performance)	TensorFlow tfprof	TensorFlow tfprof	Cycle accurate simulator + C Microbench
	Profiler / Simulator (Energy)	Intel RAPL Interface	NVIDIA System Management Interface	SPICE simulation + Intel RAPL Interface

Performance/Accuracy Trade-off (AlexNet)



Performance/Accuracy Trade-off (AlexNet)



Performance/Accuracy Trade-off (Inception V3)



Energy Efficiency for All Conv. Layers

AlexNet

Inception V3



Summary

Improving on a highly parallel in-cache accelerator by exploiting ...

- 1) **sparsity** in filter weights:
 - a) offline coalescing of weights + online coalescing of inputs
 - b) overlapping filters
- 2) low-precision weights:

efficient MAC algorithm

Overall Efficiency



622x over CPU 52x over GPU 2.6x over Neural Cache Accuracy



0.4% accuracy loss

Area Overhead



2% of CPU chip

Bit Prudent In-Cache Acceleration of Deep Convolutional Neural Networks

Xiaowei Wang, Jiecao Yu, Charles Augustine, Ravi Iyer, Reetuparna Das

Thank you! Questions & Answers



M-Bits Research Group UNIVERSITY OF MICHIGAN

