Joint Design of Algorithms and Hardware for Energy-Efficient DNNs

Vivienne Sze



NIPS 2016 EMDNN Workshop

Phir



s technology laboratories

Video is the Biggest Big Data

Over 70% of today's Internet traffic is video Over 300 hours of video uploaded to YouTube <u>every minute</u> Over 500 million hours of video surveillance collected <u>every day</u>



Need energy-efficient pixel processing!



Energy-Efficient Pixel Processing



Goal: Increase coding efficiency, speed and energy-efficiency

Energy-Efficient Computer Vision & Deep Learning (Understand Pixels)



Goal: Make computer vision as ubiquitous as video coding

I Typical Constraints on Video Coding

- Area cost
 - Memory Size 100-500kB
- Power budget
 - < 1W for smartphones</p>
- Throughput
 - Real-time 30 fps
- Energy
 - ~1nJ/pixel









MIT Object Detection Chip [VLSI 2016]

Eyeriss: Energy-Efficient Hardware for DCNNs

Yu-Hsin Chen, Tushar Krishna, Joel Emer, Vivienne Sze, ISSCC 2016 / ISCA 2016











Deep Convolutional Neural Networks

Modern *deep* CNN: can be over **100** CONV layers





Deep Convolutional Neural Networks





Deep Convolutional Neural Networks





Convolutions account for more than 90% of overall computation, dominating **runtime** and **energy consumption**





Input Image (Feature Map)













Sliding Window Processing







Many Input Channels (C)











Image batch size: 1 – 256 (N)



ms technology laboratories

Large Sizes with Varying Shapes

AlexNet¹ Convolutional Layer Configurations

Layer	Filter Size (R)	# Filters (M)	# Channels (C)	Stride
1	11x11	96	3	4
2	5x5	256	48	1
3	3x3	384	256	1
4	3x3	384	192	1
5	3x3	256	192	1

Layer 1



34k Params 105M MACs Layer 2





307k Params 224M MACs



885k Params 150M MACs





- Operations exhibit high parallelism
 - → high throughput possible



- Operations exhibit high parallelism
 → high throughput possible
- Memory Access is the Bottleneck



* multiply-and-accumulate



- Operations exhibit high parallelism
 → high throughput possible
- Memory Access is the Bottleneck



Worst Case: all memory R/W are **DRAM** accesses

Example: AlexNet [NIPS 2012] has 724M MACs
 → 2896M DRAM accesses required



- Operations exhibit high parallelism
 → high throughput possible
- Input data reuse opportunities (up to 500x)

→ exploit **low-cost memory**



Fmap

In Highly-Parallel Compute Paradigms

Temporal Architecture (SIMD/SIMT)



Spatial Architecture (Dataflow Processing)





Advantages of Spatial Architecture







22 How to Map the Dataflow?



partial sums accumulation



14117

23

Energy-Efficient Dataflow

Yu-Hsin Chen, Joel Emer, Vivienne Sze, ISCA 2016

Maximize data reuse and accumulation at RF





24 Data Movement is Expensive



Processing Engine



Data Movement Energy Cost



Maximize data reuse at lower levels of hierarchy

25 Weight Stationary (WS)



- Minimize weight read energy consumption
 - maximize convolutional and filter reuse of weights
- Examples:

[Chakradhar, ISCA 2010] [nn-X (NeuFlow), CVPRW 2014] [Park, ISSCC 2015] [Origami, GLSVLSI 2015]



Output Stationary (OS)



- Minimize partial sum R/W energy consumption
 - maximize local accumulation
- Examples:

[Gupta, *ICML* 2015] [ShiDianNao, *ISCA* 2015] [Peemen, *ICCD* 2013]





27 No Local Reuse (NLR)



- Use a large global buffer as shared storage
 - Reduce **DRAM** access energy consumption
- Examples:

[DianNao, ASPLOS 2014] [DaDianNao, MICRO 2014] [Zhang, FPGA 2015]



Row Stationary: Energy-efficient Dataflow





































- Maximize row convolutional reuse in RF
 - Keep a filter row and fmap sliding window in RF
- Maximize row psum accumulation in RF





2D Convolution in PE Array









35 2D Convolution in PE Array







36 2D Convolution in PE Array





2D Convolution in PE Array







Convolutional Reuse Maximized



Filter rows are reused across PEs horizontally



Convolutional Reuse Maximized



Feature map rows are reused across PEs diagonally





Maximize 2D Accumulation in PE Array



Partial sums accumulate across PEs vertically





41 CNN Convolution – The Full Picture



Map rows from **multiple fmaps, filters** and **channels** to same PE to exploit other forms of reuse and local accumulation

Evaluate Reuse in Different Dataflows

Weight Stationary

- Minimize movement of filter weights

Output Stationary

- Minimize movement of partial sums

No Local Reuse

- Don't use any local PE storage. Maximize global buffer size.

Row Stationary



Evaluate Reuse in Different Dataflows

Weight Stationary

- Minimize movement of filter weights

Output Stationary

- Minimize movement of partial sums

No Local Reuse

- Don't use any local PE storage. Maximize global buffer size.

Row Stationary

Evaluation Setup

- Same Total Area
- AlexNet
- 256 PEs
- Batch size = 16



Dataflow Comparison: CONV Layers



RS uses 1.4× – 2.5× lower energy than other dataflows



Dataflow Comparison: CONV Layers





Dataflow Comparison: FC Layers



RESEARCH LABORATORY OF ELECTRONICS AT MIT





Energy-Efficient Accelerator

Yu-Hsin Chen, Tushar Krishna, Joel Emer, Vivienne Sze, ISSCC 2016

Exploit data statistics

Exercise Deep CNN Accelerator

Data Compression Saves DRAM BW

Apply Non-Linearity (ReLU) on Filtered Image Data

Zero Data Processing Gating

- Skip PE local memory access
- Skip MAC computation
- Save PE processing power by 45%

⁵⁵ Chip Spec & Measurement Results¹

Technology TSMC 65nm LP 1P9M		4000 um		
On-Chip Buffer	108 KB			4000 μm
# of PEs	168		3'3'3'3'3'3'3'3'3'3 3'3	
Scratch Pad / PE	0.5 KB	Glo	bal	Spatial Array
Core Frequency	100 – 250 MHz	Bu	ifer	(168 PEs)
Peak Performance	33.6 – 84.0 GOPS			
Word Bit-width	16-bit Fixed-Point			
Natively Supported CNN ShapesFilter Width: 1 – 32 Filter Height: 1 – 12 Num. Filters: 1 – 1024 Num. Channels: 1 – 1024 Horz. Stride: 1 – 12				

AlexNet: For 2.66 GMACs [8 billion 16-bit inputs (**16GB**) and 2.7 billion outputs (**5.4GB**)], only requires **208.5MB** (buffer) and **15.4MB** (DRAM)

4000 µm

56 Comparison with GPU

	This Work	NVIDIA TK1 (Jetson Kit)	
Technology	65nm	28nm	
Clock Rate	200MHz	852MHz	
# Multipliers	168	192	
On-Chip Storage	Buffer: 108KB Spad: 75.3KB	Shared Mem: 64KB Reg File: 256KB	
Word Bit-Width	16b Fixed	32b Float	
Throughput ¹	34.7 fps	68 fps	
Measured Power	278 mW	Idle/Active ² : 3.7W/10.2W	
DRAM Bandwidth	127 MB/s	1120 MB/s ³	

- 1. AlexNet Convolutional Layers Only
- 2. Board Power
- 3. Modeled from [Tan, SC11]

Demo of Image Classification on Eyeriss

https://vimeo.com/154012013

Integrated with BVLC Caffe DL Framework

Summary of Eyeriss Deep CNN

- Eyeriss: a reconfigurable accelerator for state-of-the-art deep CNNs at below 300mW
- Energy-efficient dataflow to reduce data movement
- Exploit data statistics for high energy efficiency
- Integrated with the Caffe DL framework and demonstrated an image classification system

Learn more about Eyeriss at

http://eyeriss.mit.edu

58

Features: Energy vs. Accuracy

59

60

Designing Energy-Efficient CNNs using Energy-Aware Pruning

Tien-Ju Yang, Yu-Hsin Chen, Vivienne Sze, arXiv 2016

Key Metrics for Embedded DNN

- Accuracy → Measured on Dataset
- Speed \rightarrow Number of MACs
- Storage Footprint \rightarrow Number of Weights
- Energy \rightarrow ?

Energy-Evaluation Methodology

62

Hardware Energy Costs of each MAC and Memory Access

Energy estimation tool to be released on http://eyeriss.mit.edu

63 Key Observations

- Number of weights *alone* is not a good metric for energy
- All data types should be considered

Energy Consumption of Existing DNNs

Deeper CNNs with fewer weights do not necessarily consume less energy than shallower CNNs with more weights

64

[Yang et al., arXiv 2016]

Magnitude-based Weight Pruning

Reduce number of weights by **removing small magnitude weights**

65

Energy-Aware Pruning

3.7x reduction in AlexNet / 1.6x reduction in GoogLeNet

66

[Yang et al., arXiv 2016]

- Energy-Efficient Approaches
 - Minimize data movement
 - Balance flexibility and energy-efficiency
 - Exploit sparsity with joint algorithm and hardware design
- Joint algorithm and hardware design can deliver additional energy savings (directly target energy)
- Linear increase in accuracy requires exponential increase in energy

Acknowledgements: This work is funded by the DARPA YFA grant, TSMC University Shuttle Program, MIT Center for Integrated Circuits & Systems, and gifts from Intel, Nvidia and Texas Instruments.

