Energy-Efficient Processing at the Edge: From Compressing to Understanding Pixels

Vivienne Sze



website: www.rle.mit.edu/eems



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

Next-Generation Video Coding (Compress Pixels)



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

Energy-Efficient Cross-Layer Design



tems technology laboratories

4

Energy-Efficient Video Compression



High Efficiency Video Coding (HEVC)

6

- HEVC achieves ~2x higher coding efficiency than H.264/AVC
- High throughput (Ultra-HD 8K @ 120fps) & low power
 Implementation-friendly features (e.g. built-in parallelism)

Size Energy		Coding Efficiency	Efficient Implementation
1.5x	Larger and Flexible Coding Block Size	Х	
	More Sophisticated Intra Prediction	Х	
2x	Larger Interpolation for Motion Comp.	Х	
2x	Larger Transform Size	Х	
	Parallel Deblocking Filter		х
	Sample Adaptive Offset	Х	
	High-Throughput CABAC	Х	х
(1994) (2003) (2013)	High Level Parallel Tools		Х

Joint algorithm and hardware design is required to address **coding efficiency, throughput and power challenges**

Efficient Hardware for HEVC

Energy-Efficient HEVC Transform

- Large transforms give coding gain: 7-9%
- Adapt to sparsity of coefficients
- Enable constant energy/pixel for all transform sizes

High-Throughput HEVC CABAC

- CABAC bottleneck in H.264/AVC
- HEVC 4x faster than H.264/AVC
- Ultra-HD 8K @ 120fps
- Also gives coding gain: 3-5%



Plii

Low Power HEVC Decoder for Wearable Devices



8



ns technology laboratories

Compression Inspired by Computer Vision



¹⁰ SIFT feature matching



SIFT features are widely used to establish correspondence between two similar images



Plii





11 SIFT is Rotate Invariant



SIFT descriptor



Canonizing the rotation

Before the matching, SIFT descriptors are normalized to the canonical pose (dominant gradient) so that patches of different orientation can be matched.



Intra Block Copy for Still Image Coding



Use one block to **predict** repetitive blocks. Only encode the **difference (residual)**.

NOT rotate invariant. Limited to screen content.





14111

¹³ *Rotate* Intra Block Copy



Repetitive structures with rotation In both screen content and camera captured images





14 Reduction of Residual Energy



HEVC

HEVC + Block Copy

HEVC + Rotate Block Copy

40% reduction of residual energy over HEVC 27% reduction of residual energy over HEVC + Block Copy

However, there is overhead in signaling the rotate angle and motion vector

First frame of *ParkScene* Sequence

chnology laboratories

|'|**|**iT

[Z. Zhang et al., ICIP 2015]



Motion Vector Prediction



Motion vectors need to be on the same rotated coordinate system

 $\frac{mv_{\theta_2}^{(2)} \text{ is encoded as } mv_{\theta_2}^{(2)} - \operatorname{round} \left(\frac{R_{\theta_2 - \theta_1}}{mv_{\theta_1}^{(1)}} \right)}{\text{Where } R_{\theta_2 - \theta_1}} = \begin{bmatrix} \cos(\theta_2 - \theta_1) & -\sin(\theta_2 - \theta_1) \\ \sin(\theta_2 - \theta_1) & \cos(\theta_2 - \theta_1) \end{bmatrix}$

Reduce average bit rate of motion vector difference by 25%

[Z. Zhang et al., ICIP 2015]





14117

HEVC + Intra Block Copy vs. HEVC + Rotate Intra Block Copy

	Sequence	Residual	BD-rate	
		reduction	reduction	
Class C	RaceHorse	23.66%	4.75	
	PartyScene	27.64%	4.65	
	BQMall	17.92%	2.70	
	BasketballDrill	22.12%	3.52	
Class D	BQSquare	30.82%	5.25	
	BasketballPass	15.44%	1.87	
	BlowingBubbles	7.59%	2.88	
	RaceHorse	28.97%	4.62	
Class E	FourPeople	18.09%	2.60	
	Johnny	12.79%	2.41	
	KristenAndSara	15.67%	2.48	
Class F	BasketballDrillText	21.15%	3.78	
screen	SlideShow	29.01%	8.03	
content	SlideEditing	19.12%	0.74	
Class C Average		22.83%	3.91	
Class D Average		20.70%	3.66	
Class E Average		15.52%	2.50	
Cla	ass F Average	23.09%	4.18	
Ov	verall Average	20.71%	3.60	

[Z. Zhang et al., ICIP 2015]

Evaluate on First Frame of JCT-VC test sequences

- Residual Energy reduction of 20.7%
- BD-rate savings of 3.6%



|'|iT

Energy-Efficient Deep Learning





Example Applications of Deep Learning

Computer Vision



Speech Recognition



Medical









Using Deep Learning for Compression

Coding Tool in Video Codec



Proposed for VVC: prediction, loop filtering, upsampling, etc.

JVET AHG: Neural Networks in Video Coding

Intra Prediction (Upsampling) [Li et al., TCSVT 2017]

End to End Auto Encoder



Challenge: Computation complexity higher than typical image processing

Deep Convolutional Neural Networks







Deep Convolutional Neural Networks





22 Deep Convolutional Neural Networks





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





²³ High-Dimensional CNN Convolution

Input Image (Feature Map)







High-Dimensional CNN Convolution

Input Image (Feature Map)



Element-wise Multiplication



²⁵ High-Dimensional CNN Convolution





²⁶ High-Dimensional CNN Convolution



Sliding Window Processing





I High-Dimensional CNN Convolution



Many Input Channels (C)





High-Dimensional CNN Convolution





High-Dimensional CNN Convolution



Image batch size: 1 – 256 (N)

l'liī



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



Properties We Can Leverage

- 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



Properties We Can Leverage

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

→ exploit **low-cost memory**



Image

Bighly-Parallel Compute Paradigms

Temporal Architecture (SIMD/SIMT)



Spatial Architecture (Dataflow Processing)





Advantages of Spatial Architecture







35 Data Movement is Expensive



Processing Engine



Data Movement Energy Cost



Maximize data reuse at lower levels of hierarchy

Bow to Map the Dataflow?



Goal: Increase reuse of input data (weights and pixels) and local partial sums accumulation

Spatial Architecture (Dataflow Processing)




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] [Google's TPU, ISCA 2017]



Output Stationary (OS)



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

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



Row Stationary Dataflow



40 Sparsity in Data

Many zeros in output fmaps after ReLU



Exploit Sparsity 41

Method 1. Skip memory access and computation



OF FLECTRONICS AT MIT

ms technology laboratories

Method 2. Compress data to reduce storage and data movement



Plii

⁴² Eyeriss: Energy-Efficient Deep Learning



278mW for AlexNet @ 30fps (batch size 4) in 65nm LP CMOS

> 10x more energy-efficient than mobile GPUs





⁴³ Features: Energy vs. Accuracy



Design of Efficient DNN Algorithms

• Popular efficient DNN algorithm approaches



... also reduced precision

- Focus on reducing number of MACs and weights
- Does it translate to energy savings?



l'lli**r**

Energy-Evaluation Methodology



45

Hardware Energy Costs of each MAC and Memory Access



Illi Energy estimation tool available at http://eyeriss.mit.edu

T MIT MIT microsystems technology laboratories massachusetts institute of technology

46 Key Observations

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







⁴⁷ Energy Consumption of Existing DNNs



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

|'||iT

[Yang et al., CVPR 2017]



Magnitude-based Weight Pruning



Reduce number of weights by removing small magnitude weights

	_
	•





Energy-Aware Pruning



l'liī

49

[Yang et al., CVPR 2017]



NetAdapt: Platform-Aware DNN Adaptation

- Automatically adapt DNN to a mobile platform to reach a target latency or energy budget
- Use **empirical measurements** to guide optimization (avoid modeling of tool chain or platform architecture)



IIIii In collaboration with Google's Mobile Vision Team



⁵¹ Latency vs. Accuracy Tradeoff with NetAdapt

 NetAdapt boosts <u>the real inference speed</u> of MobileNet by 1.7x with higher accuracy



Reference:

MobileNet: Howard et al, "Mobilenets: Efficient convolutional neural networks for mobile vision applications", arXiv 2017 **MorphNet:** Gordon et al., "Morphnet: Fast & simple resource-constrained structure learning of deep networks", CVPR 2018



I Tutorial Material on Efficient DNNs

Proceedings of EEE

Efficient Processing of Deep Neural Networks: A Tutorial and Survey

System Scaling With Nanostructured Power and RF Components Nonorthogonal Multiple Access for 5G and Beyond

Point of View: Beyond Smart Grid—A Cyber–Physical–Social System in Energy Future Scanning Our Past: Materials Science, Instrument Knowledge, and the Power Source Renaissance



Hardware Architectures for Deep Neural Networks

ISCA Tutorial

June 24, 2017

Website: http://eyeriss.mit.edu/tutorial.html



http://eyeriss.mit.edu/tutorial.html

V. Sze, Y.-H. Chen, T-J. Yang, J. Emer, "*Efficient Processing of Deep Neural Networks: A Tutorial and Survey,*" Proceedings of the IEEE, 2017



Efficient Computer Vision using Compression





Super-Resolution on Mobile Devices



Transmit low resolution for lower bandwidth

Screens are getting larger



Use **super-resolution** to improve the viewing experience of lower-resolution content (*reduce communication bandwidth*)





Complexity of Super Resolution Algorithms





8032 MACs/pixel \rightarrow ~500 GMAC/s for HD @ 30 fps

State-of-the-art super resolution algorithms use CNNs → computationally expensive, especially at high resolutions (HD or 4K)





FAST: A Framework to Accelerate SuperRes



Real-time

A framework that accelerates **any SR** algorithm by up to **15x** when running on compressed videos





Plii

Free Information in Compressed Videos







Compressed video



Block-structure

Motion-compensation

Video as a stack of pixels

Representation in compressed video

This representation can help accelerate super-resolution





58 Transfer the Super-Resolution Results



Transfer is Lightweight



Fractional Bicubic Interpolation Interpolation Skip Flag

The complexity of the transfer is comparable to bicubic interpolation. Transfer N frames, accelerate by N





•• Challenge 1: Scene Transition



Transfer will NOT work if there is a transition of scenes



Group-of-Picture (GoP) Structure

GoP structure in the compressed video provides video segmentation for free



Challenge 2: Prediction Error



Ground-truth



Non-Adaptive Artifacts when missing high frequency components of residual



FAST skips transfer on blocks with large residual

Adaptive Use lightweight metric to identify occurrence and skip transfer





⁶² Challenge 3: Blocking Artifacts



No deblocking



Overlapped Block Processing

No blocking artifacts Process each pixel multiple times

Non-Overlapped Block Processing







Process each pixel once Blocking artifacts



Challenge 3: Blocking Artifacts



FAST applies the **deblocking filter** to alleviate the blocking effect caused by **non-overlapping block division**





l'liī

64 Challenge 4: Accumulated Error



When a SR result gets transferred multiple times, the error **accumulates**



FAST estimates the accumulated error as the **accumulated Laplacian of the residual**, and stops the transfer when it exceeds a threshold





Evaluation: Accelerating SRCNN







PartyScene

RaceHorse

BasketballPass

Examples of videos in the test set (20 videos for HEVC development)





 $4 \times$ acceleration with NO PSNR LOSS. $16 \times$ acceleration with 0.2 dB loss of PSNR



Visual Evaluation









Bicubic



SRCNN with FAST

Ground-truth







⁶⁷ Visual Evaluation



SRCNN

FAST + SRCNN

Bicubic





Summary of FAST

- Transfer the SR results guided by motion vectors
- Adaptively perform the transfer by thresholding on the residue, and accumulated Laplacian
- Accelerates SR algorithms by up to 15x with minimal PSNR loss



Compressed video

Code released at www.rle.mit.edu/eems/fast

[Zhang et al., CVPRW 2017]





l'lli**r**

Enable Real-time Navigation on nano drone



Image source: Cheerson

69

Big battery

Mobile GPU

Enable energy-efficient navigation for **Search and Rescue**

[Zhang et al., RSS 2017]

http://navion.mit.edu



IIIIT In collaboration with Sertac Karaman and Luca Carlone (AeroAstro)



Localization with Visual Inertial Odometry 70

VIO determines location/orientation of drone from images and IMU (also used by headset in Augmented Reality and Virtual Reality)



Navion: Fully integrated VIO system on-chip consuming < **30mW**

http://navion.mit.edu





IIIii

71 Visual Inertial Odometry Demo





l'liī







72 Compression to Reduce Energy and Cost

Memory dominates energy and chip area.



[Suleiman et al., VLSI 2018]

Apply various **compression techniques** to reduce on-chip storage cost by 4.1x. Entire VIO system is fully integrated on chip (20mm²).

http://navion.mit.edu






- Video is perhaps the biggest of the 'big data' being collected and transmitted.
- Moving from compressing to understanding pixels at the edge increasingly desirable due to privacy/security and latency constraints. However, energy significantly limited at edge.
- Co-design of algorithms and hardware can enable energyefficient video coding, computer vision and deep learning such that they can efficiently operate on edge devices such as smartphones and autonomous vehicles/drones.
- Bridging the gap between video coding, computer vision and deep learning plays an important role in overcoming many of the challenges faced by next generation of edge devices.



74 Acknowledgements



Research conducted in the **MIT Energy-Efficient Multimedia Systems Group** would not be possible without the support of the following organizations:



75 References

Energy-Efficient Video Coding

- M. Tikekar, C.-T. Huang, V. Sze, A. Chandrakasan, "Energy and Area-Efficient Hardware Implementation of HEVC Inverse Transform and Dequantization," IEEE International Conference on Image Processing (ICIP), pp. 2100-2104, October 2014.
- Y.-H. Chen, V. Sze, "A Deeply Pipelined CABAC Decoder for HEVC Supporting Level 6.2 High-Tier Applications," IEEE Transactions on Circuits and Systems for Video Technology (TCSVT), Vol. 25, No. 5, pp. 856-868, May 2015.
- M. Tikekar, V. Sze, A. Chandrakasan, "A Fully-Integrated Energy-Efficient H.265/HEVC Decoder with eDRAM for Wearable Devices," IEEE Symposium on VLSI Circuits (VLSI-Circuits), June 2017.

Compression Inspired by Computer Vision

 Z. Zhang, V. Sze, "Rotate Intra Block Copy for Still Image Coding," IEEE International Conference on Image Processing (ICIP), September 2015.



76 References

Energy-Efficient Deep Learning

- Y.-H. Chen, T. Krishna, J. Emer, V. Sze, "Eyeriss: An Energy-Efficient Reconfigurable Accelerator for Deep Convolutional Neural Networks," IEEE International Conference on Solid-State Circuits (ISSCC), pp. 262-264, February 2016.
- Y.-H. Chen, J. Emer, V. Sze, "Eyeriss: A Spatial Architecture for Energy-Efficient Dataflow for Convolutional Neural Networks," International Symposium on Computer Architecture (ISCA), pp. 367-379, June 2016.
- A. Suleiman, Z. Zhang, V. Sze, "A 58.6mW Real-time Programmable Object Detection with Multi-Scale Multi-Object Support Using Deformable Parts Models on 1920×1080 Video at 30fps," IEEE Symposium on VLSI Circuits (VLSI-Circuits), pp. 184-185, June 2016.
- A. Suleiman*, Y.-H. Chen*, J. Emer, V. Sze, "Towards Closing the Energy Gap Between HOG and CNN Features for Embedded Vision," IEEE International Symposium of Circuits and Systems (ISCAS), Invited Paper, May 2017.
- T.-J. Yang, Y.-H. Chen, V. Sze, "Designing Energy-Efficient Convolutional Neural Networks using Energy-Aware Pruning," IEEE Conference on Computer Vision and Pattern Recognition (CVPR), July 2017.
- T.-J. Yang, A. Howard, B. Chen, X. Zhang, A. Go, V. Sze, H. Adam, "NetAdapt: Platform-Aware Neural Network Adaptation for Mobile Applications," arXiv, April 2018.
- V. Sze, Y.-H. Chen, T.-J. Yang, J. Emer, "Efficient Processing of Deep Neural Networks: A Tutorial and Survey," Proceedings of the IEEE, vol. 105, no. 12, pp. 2295-2329, December 2017.





77 References

• Efficient Computer Vision using Compression

- Z. Zhang, V. Sze, "FAST: A Framework to Accelerate Super-Resolution Processing on Compressed Videos," CVPR Workshop on New Trends in Image Restoration and Enhancement, July 2017.
- A. Suleiman, Z. Zhang, L. Carlone, S. Karaman, V. Sze, "Navion: A Fully Integrated Energy-Efficient Visual-Inertial Odometry Accelerator for Autonomous Navigation of Nano Drones," IEEE Symposium on VLSI Circuits (VLSI-Circuits), June 2018.
- Z. Zhang*, A. Suleiman*, L. Carlone, V. Sze, S. Karaman, "Visual-Inertial Odometry on Chip: An Algorithm-and-Hardware Co-design Approach," Robotics: Science and Systems (RSS), July 2017.

