Efficient Computing for Deep Learning, AI and Robotics

Vivienne Sze (**)**@eems_mit)

Massachusetts Institute of Technology

In collaboration with Luca Carlone, Yu-Hsin Chen, Joel Emer, Sertac Karaman, Tushar Krishna, Thomas Heldt, Trevor Henderson, Hsin-Yu Lai, Peter Li, Fangchang Ma, James Noraky, Gladynel Saavedra Peña, Charlie Sodini, Amr Suleiman, Nellie Wu, Diana Wofk, Tien-Ju Yang, Zhengdong Zhang

Slides available at <u>https://tinyurl.com/SzeMITDL2020</u>

Compute Demands for Deep Neural Networks

AlexNet to AlphaGo Zero: A 300,000x Increase in Compute



Source: Open AI (https://openai.com/blog/ai-and-compute/)

2

Compute Demands for Deep Neural Networks

Common carbon footprint benchmarks

in lbs of CO2 equivalent



Chart: MIT Technology Review

[Strubell, ACL 2019]

4 Processing at "Edge" instead of the "Cloud"



Communication

Privacy

Latency

Vivienne Sze (@eems_mit)

Computing Challenge for Self-Driving Cars

JACK STEWART TRANSPORTATION 02.06.18 08:00 AM

SELF-DRIVING CARS USE CRAZY AMOUNTS OF POWER, AND IT'S BECOMING A PROBLEM



Shelley, a self-driving Audi TT developed by Stanford University, uses the brains in the trunk to speed around a racetrack autonomously.

R NIKKI KAHN/THE WASHINGTON POST/GETTY IMAGES

(Feb 2018)

Cameras and radar generate ~6 gigabytes of data every 30 seconds.

Self-driving car prototypes use approximately 2,500 Watts of computing power.

Generates wasted heat and some prototypes need water-cooling!

Existing Processors Consume Too Much Power



< 1 Watt

> 10 Watts

Transistors are NOT Getting More Efficient

Slow down of Moore's Law and Dennard Scaling

General purpose microprocessors not getting faster or more efficient



Popularity of Specialized Hardware for DNNs

The New York Times

By CADE METZ JAN 14, 2018

8

0000



Big Bets On A.I. Open a New Frontier for Chips Start-Ups, Too. (January 14, 2018)

"Today, at least 45 start-ups are working on chips that can power tasks like speech and self-driving cars, and at least five of them have raised more than \$100 million from investors. Venture capitalists invested more than \$1.5 billion in chip start-ups last year, nearly doubling the investments made two years ago, according to the research firm CB Insights."

Power Dominated by Data Movement



Vivienne Sze (**y**@eems_mit)

[Horowitz, /SSCC 2014]

Autonomous Navigation Uses a Lot of Data

- Semantic Understanding
- High frame rate
- Large resolutions
- Data expansion



2 million pixels



10x-100x more pixels

- Geometric Understanding
- Growing map size



[Pire, RAS 2017]

11 Understanding the Environment

Depth Estimation





Semantic Segmentation



State-of-the-art approaches use Deep Neural Networks, which require up to several hundred millions of operations and weights to compute! >100x more complex than video compression

Vivienne Sze (@eems_mit)

Deep Neural Networks

Deep Neural Networks (DNNs) have become a cornerstone of AI

Computer Vision



Game Play





Medical





¹³ What Are Deep Neural Networks?



14 Weighted Sum



Key operation is **multiply and accumulate (MAC)** Accounts for > 90% of computation

Vivienne Sze (@eems_mit)

Popular Types of Layers in DNNs

Fully Connected Layer

- Feed forward, fully connected
- Multilayer Perceptron (MLP)

Convolutional Layer

- Feed forward, sparsely-connected w/ weight sharing
- Convolutional Neural Network (CNN)
- Typically used for images

Recurrent Layer

Feedback

15

- Recurrent Neural Network (RNN)
- Typically used for sequential data (e.g., speech, language)

Attention Layer/Mechanism

- Attention (matrix multiply) + feed forward, fully connected
- Transformer [Vaswani, NeurIPS 2017]



High-Dimensional Convolution in CNN

a plane of input activations a.k.a. **input feature map (fmap)**

filter (weights)





II High-Dimensional Convolution in CNN



High-Dimensional Convolution in CNN



Sliding Window Processing

High-Dimensional Convolution in CNN



Many Input Channels (C)

AlexNet: 3 – 192 Channels (C)

Vivienne Sze (@eems_mit)

In High-Dimensional Convolution in CNN



Vivienne Sze (@eems_mit)

High-Dimensional Convolution in CNN



22 Define Shape for Each Layer



Shape varies across layers

²³ Layers with Varying Shapes

MobileNetV3-Large Convolutional Layer Configurations

Block	Filter Size (RxS)	# Filters (M)	۸) # Channels (C)					
1	3x3	16	3					
3	1x1	64	16					
3	3x3	64	1					
3	1x1	24	64					
6	1x1	120	40					
6	5x5	120	1					
6	1x1	40 120						

[Howard, *ICCV* 2019]

24 Popular DNN Models

Metrics	LeNet-5	AlexNet	VGG-16	GoogLeNet (v1)	ResNet-50	EfficientNet-B4
Top-5 error (ImageNet)	n/a	16.4	7.4	6.7	5.3	3.7*
Input Size	28x28	227x227	224x224	224x224	224x224	380x380
# of CONV Layers	2	5	16	21 (depth)	49	96
# of Weights	2.6k	2.3M	14.7M	6.0M	23.5M	14M
# of MACs	283k	666M	15.3G	1.43G	3.86G	4.4G
# of FC layers	2	3	3	1	1	65**
# of Weights	58k	58.6M	124M	1M	2M	4.9M
# of MACs	58k	58.6M	124M	1M	2M	4.9M
Total Weights	60k	61M	138M	7M	25.5M	19M
Total MACs	341k	724M	15.5G	1.43G	3.9G	4.4G
Reference	Lecun , PIEEE 1998	Krizhevsky, NeurIPS 2012	Simonyan, ICLR 2015	Szegedy, CVPR 2015	He , CVPR 2016	Tan , <i>ICML</i> 2019

DNN models getting larger and deeper

* Does not include multi-crop and ensemble

** Increase in FC layers due to squeeze-and-excitation layers (much smaller than FC layers for classification)

Efficient Hardware Acceleration for Deep Neural Networks

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 has 724M MACs
 → 2896M DRAM accesses required

26

²⁷ **Properties We Can Leverage**

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



Convolutional Reuse

(Activations, Weights) CONV layers only (sliding window)



Fmap Reuse (Activations) CONV and FC layers Input Fmaps

Filter Reuse (Weights) CONV and FC layers (batch size > 1)

Exploit Data Reuse at Low-Cost Memories





* measured from a commercial 65nm process

Farther and **larger** memories consume more power

Vivienn

Plii

²⁹ Weight Stationary (WS)



- Minimize weight read energy consumption
 - maximize convolutional and filter reuse of weights
- Broadcast activations and accumulate partial sums spatially across the PE array
- Examples: TPU [Jouppi, /SCA 2017], NVDLA

Output Stationary (OS)



- Minimize partial sum R/W energy consumption
 - maximize local accumulation
- Broadcast/Multicast filter weights and reuse activations spatially across the PE array
- Examples: [Moons, VLSI 2016], [Thinker, VLSI 2017]

Input Stationary (IS)



- Minimize activation read energy consumption
 - maximize convolutional and fmap reuse of activations
- Unicast weights and accumulate partial sums spatially across the PE array
- Example: [SCNN, ISCA 2017]

Row Stationary Dataflow



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



Row Stationary Dataflow



Dataflow Comparison: CONV Layers



[Chen, ISCA 2016]

Exploit Sparsity

35

Method 1. Skip memory access and computation



<u>Method 2</u>. Compress data to reduce storage and data movement



Vivienne Sze (@eems_mit)

Everiss: Deep Neural Network Accelerator

36



[Chen, /SSCC 2016]

Exploits data reuse for **100x** reduction in memory accesses from global buffer and **1400x** reduction in memory accesses from off-chip DRAM

Overall >10x energy reduction compared to a mobile GPU (Nvidia TK1)

Eyeriss Project Website: http://eyeriss.mit.edu

Results for AlexNet

Vivienne Sze (@eems_mit) [Joint work with Joel Emer]
³⁷ Features: Energy vs. Accuracy



Vivienne Sze (@eems_mit)

[Suleiman*, Chen*, ISCAS 2017]

Illii

Energy-Efficient Processing of DNNs

A significant amount of algorithm and hardware research on energy-efficient processing of DNNs



We identified various limitations to existing approaches

Design of Efficient DNN Algorithms

• Popular efficient DNN algorithm approaches

39



... also reduced precision

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

Data Movement is Expensive





* measured from a commercial 65nm process

Energy of weight depends on **memory hierarchy** and **dataflow**

Energy-Evaluation Methodology 41



Tool available at: https://energyestimation.mit.edu/

[Yang, CVPR 2017]

E_{data}

L1 L2 L3

DNN Energy Consumption

42 Key Observations

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



Vivienne Sze (**v**@eems_mit)

[Yang, CVPR 2017]

43 Energy-Aware Pruning

Directly target energy and incorporate it into the optimization of DNNs to provide greater energy savings

- Sort layers based on energy and prune layers that consume most energy first
- EAP reduces AlexNet energy by
 3.7x and outperforms the previous work that uses magnitude-based pruning by **1.7x**



Pruned models available at <u>http://eyeriss.mit.edu/energy.html</u>

44 # of Operations vs. Latency

• # of operations (MACs) does not approximate latency well



Source: Google (https://ai.googleblog.com/2018/04/introducing-cvpr-2018-on-device-visual.html)

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)



Vivienne Sze (@eems_mit) In collaboration with Google's Mobile Vision Team

Simplified Example of One Iteration



Vivienne Sze (**y**@eems_mit)

46

[Yang, ECCV 2018]

Improved Latency vs. Accuracy Tradeoff

 NetAdapt boosts the real inference speed of MobileNet by up to 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

[Yang, ECCV 2018]

48 FastDepth: Fast Monocular Depth Estimation

Depth estimation from a single RGB image desirable, due to the relatively low cost and size of monocular cameras.

RGB

Prediction



Auto Encoder DNN Architecture (Dense Output)



Vivienne Sze (@eems_mit)

[Joint work with Sertac Karaman]

FastDepth: Fast Monocular Depth Estimation

Apply NetAdapt, compact network design, and depth wise decomposition to decoder layer to enable depth estimation at **high frame rates on an embedded platform** while still maintaining accuracy



49

[Wofk*, Ma*, *ICRA* 2019] IIIii

Many Efficient DNN Design Approaches

50



Vivienne Sze (geems_mit) [Chen*, Yang*, SysML 2018]

51 Existing DNN Architectures

- Specialized DNN hardware often rely on certain properties of DNN in order to achieve high energy-efficiency
- Example: Reduce memory access by amortizing across MAC array



Limitation of Existing DNN Architectures

- Example: Reuse and array utilization depends on # of channels, feature map/batch size
 - Not efficient across all network architectures (e.g., compact DNNs)



Limitation of Existing DNN Architectures

- Example: Reuse and array utilization depends on # of channels, feature map/batch size
 - Not efficient across all network architectures (e.g., compact DNNs)



Limitation of Existing DNN Architectures

- Example: Reuse and array utilization depends on # of channels, feature map/batch size
 - Not efficient across all network architectures (e.g., compact DNNs)
 - Less efficient as array scales up in size
 - Can be challenging to exploit sparsity



Need Flexible Dataflow

 Use flexible dataflow (Row Stationary) to exploit reuse in any dimension of DNN to increase energy efficiency and array utilization



Example: Depth-wise layer

Need Flexible NoC for Varying Reuse

- When reuse available, need **multicast** to exploit spatial data reuse for energy efficiency and high array utilization
- When reuse not available, need **unicast** for high BW for weights for FC and weights & activations for high PE utilization
- An all-to-all satisfies above but too expensive and not scalable



⁵⁷ Hierarchical Mesh





Vivienne Sze (**y**@eems_mit)

[Chen, JETCAS 2019]

Eyeriss v2: Balancing Flexibility and Efficiency

Efficiently supports

58

- Wide range of filter shapes
 - Large and Compact
- Different Layers
 - CONV, FC, depth wise, etc.
- Wide range of sparsity
 - Dense and Sparse
- Scalable architecture

🛚 v1.5 & MobileNet 🔎 v2 & MobileNet 📮 v2 & sparse MobileNet



Speed up over Eyeriss v1 scales with number of PEs

# of PEs	256	1024	16384
AlexNet	17.9x	71.5x	1086.7x
GoogLeNet	10.4x	37.8x	448.8x
MobileNet	15.7x	57.9x	873.0x

Over an order of magnitude faster and more energy efficient than Eyeriss v1

[Chen, JETCAS 2019]

Looking Beyond the DNN Accelerator for Acceleration

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*)

FAST: A Framework to Accelerate SuperRes



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

Vivienne Sze (@eems_mit)

⁶² Free Information in Compressed Videos







Compressed video

Pixels

Block-structure

Motion-compensation

Video as a stack of pixels

Representation in compressed video

This representation can help accelerate super-resolution

```
Vivienne Sze ( @eems_mit)
```



⁶³ Transfer is Lightweight

Bicubic

Interpolation Interpolation



Skip Flag

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

Vivienne Sze (@eems_mit)

Fractional

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

Vivienne Sze (@eems_mit)

Visual Evaluation



SRCNNFAST +BicubicSRCNNSRCNN

Look **beyond** the DNN accelerator for opportunities to accelerate DNN processing (e.g., structure of data and temporal correlation)

Code released at <u>www.rle.mit.edu/eems/fast</u>

Vivienne Sze (@eems_mit)

Beyond Deep Neural Networks

Visual-Inertial Localization

Determines location/orientation of robot from images and IMU



Localization at Under 25 mW

First chip that performs *complete* Visual-Inertial Odometry

68

Front-End for camera (Feature detection, tracking, and outlier elimination)

Front-End for IMU (pre-integration of accelerometer

and gyroscope data)

Back-End Optimization of Pose Graph

Consumes **684× and 1582×** less energy than mobile and desktop CPUs, respectively



Navion Project Website: <u>http://navion.mit.edu</u> [Zhang et al., RSS 2017], [Suleiman et al., VLSI 2018]

Vivienne Sze (@eems_mit)

[Joint work with Sertac Karaman]

Key Methods to Reduce Data Size

Navion: Fully integrated system – no off-chip processing or storage



Use **compression** and **exploit sparsity** to reduce memory down to 854kB

Where to Go Next: Planning and Mapping 70

Robot Exploration: Decide where to go by computing Shannon Mutual Information



Vivienne Sze (@eems mit)

⁷¹ Challenge is Data Delivery to All Cores

Process multiple beams in parallel



Data delivery from memory is limited



Specialized Memory Architecture

Break up map into **separate memory banks** and novel storage pattern to minimize read conflicts when processing different beams in parallel.

Diagonal Banking Pattern

Memory Access Pattern



Compute the mutual information for an **entire map** of 20m x 20m at 0.1m resolution **in under a second** \rightarrow a 100x speed up versus CPU for 1/10th of the power.

Vivienne Sze (@eems_mit)

72

[Joint work with Sertac Karaman]
73 Monitoring Neurodegenerative Disorders



Dementia affects 50 million people worldwide today (75 million in 10 years) [World Alzheimer's Report]

Mini-Mental State Examination (MMSE)

Q1. What is the year? Season? Date?

Q2. Where are you now? State? Floor?

Q3. Could you count backward from 100 by sevens? (93, 86, ...)



Agrell et al. *Age and Ageing,* 1998.

- Neuropsychological assessments are time consuming and require a trained specialist
- Repeat medical assessments are sparse, mostly qualitative, and suffer from high retest variability

⁷⁴ Use Eye Movements for *Quantitative* Evaluation

Eye movements can be used to quantitatively evaluate severity, progression or regression of neurodegenerative diseases

High-speed camera



Phantom v25-11

Substantial head support





SR EYELINK 1000 PLUS

Reulen et al., Med. & Biol. Eng. & Comp, 1988.

Clinical measurements of saccade latency are done in constrained environments that rely on specialized, costly equipment.

Measure Eye Movements Using Phone



Vivienne Sze (@eems_mit)

75

[Saavedra Peña, EMBC 2018] [Lai, ICIP 2018]

Illii

⁷⁶ Looking For Volunteers for Eye Reaction Time



If you are near or on MIT Campus and interested in volunteering your eye movements for this study, please contact us at

volunteer-eye-movement@mit.edu

Icom Power 3D Time of Flight Imaging

- Pulsed Time of Flight: Measure distance using round trip time of laser light for each image pixel
 - Illumination + Imager Power: 2.5 20 W for range from 1 8 m
- Use computer vision techniques and passive images to estimate changes in depth without turning on laser
 - CMOS Imaging Sensor Power: < 350 mW</p>



Vivienne Sze (**v**@eems_mit)

[Noraky, *ICIP* 2017]

78 Results of Low Power Depth ToF Imaging



RGB Image

Depth Map Ground Truth Depth Map Estimated

Mean Relative Error: 0.7% Duty Cycle (on-time of laser): 11%

[**Noraky**, *ICIP* 2017]



- Efficient computing extends the reach of AI beyond the cloud by reducing communication requirements, enabling privacy, and providing low latency so that AI can be used in wide range of applications ranging from robotics to health care.
- Cross-layer design with specialized hardware enables energy-efficient AI, and will be critical to the progress of AI over the next decade.

Today's slides available at <u>https://tinyurl.com/SzeMITDL2020</u>

Additional Resources

Overview Paper

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

Book Coming Spring 2020!



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



More info about **Tutorial on DNN Architectures** http://eyeriss.mit.edu/tutorial.html



For updates EEMS Mailing List

Follow @eems_mit

Additional Resources



MIT Professional Education Course on "Designing Efficient Deep Learning Systems" <u>http://shortprograms.mit.edu/dls</u>

Next Offering: July 20-21, 2020 on MIT Campus

Additional Resources

Talks and Tutorial Available Online

https://www.rle.mit.edu/eems/publications/tutorials/





YouTube Channel EEMS Group – PI: Vivienne Sze



Acknowledgements





Joel Emer



Thomas Heldt



Sertac Karaman

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



Vivienne Sze (memory @eems_mit) Mailing List: http://mailman.mit.edu/mailman/listinfo/eems-news

References

Energy-Efficient Hardware for Deep Neural Networks

- Project website: <u>http://eyeriss.mit.edu</u>
- Y.-H. Chen, T. Krishna, J. Emer, V. Sze, "Eyeriss: An Energy-Efficient Reconfigurable Accelerator for Deep Convolutional Neural Networks," IEEE Journal of Solid State Circuits (JSSC), ISSCC Special Issue, Vol. 52, No. 1, pp. 127-138, January 2017.
- 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.
- Y.-H. Chen, T.-J. Yang, J. Emer, V. Sze, "Eyeriss v2: A Flexible Accelerator for Emerging Deep Neural Networks on Mobile Devices," IEEE Journal on Emerging and Selected Topics in Circuits and Systems (JETCAS), June 2019.
- Eyexam: <u>https://arxiv.org/abs/1807.07928</u>

• Limitations of Existing Efficient DNN Approaches

- Y.-H. Chen*, T.-J. Yang*, J. Emer, V. Sze, "Understanding the Limitations of Existing Energy-Efficient Design Approaches for Deep Neural Networks," SysML Conference, February 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.
- Hardware Architecture for Deep Neural Networks: <u>http://eyeriss.mit.edu/tutorial.html</u>

References

• Co-Design of Algorithms and Hardware for Deep Neural Networks

- 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), 2017.
- Energy estimation tool: <u>http://eyeriss.mit.edu/energy.html</u>
- T.-J. Yang, A. Howard, B. Chen, X. Zhang, A. Go, V. Sze, H. Adam, "NetAdapt: Platform-Aware Neural Network Adaptation for Mobile Applications," European Conference on Computer Vision (ECCV), 2018.
- D. Wofk*, F. Ma*, T.-J. Yang, S. Karaman, V. Sze, "FastDepth: Fast Monocular Depth Estimation on Embedded Systems," IEEE International Conference on Robotics and Automation (ICRA), May 2019. <u>http://fastdepth.mit.edu/</u>

• Energy-Efficient Visual Inertial Localization

- Project website: <u>http://navion.mit.edu</u>
- 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.
- A. Suleiman, Z. Zhang, L. Carlone, S. Karaman, V. Sze, "Navion: A 2mW Fully Integrated Real-Time Visual-Inertial Odometry Accelerator for Autonomous Navigation of Nano Drones," IEEE Journal of Solid State Circuits (JSSC), VLSI Symposia Special Issue, Vol. 54, No. 4, pp. 1106-1119, April 2019.

References

86

• Fast Shannon Mutual Information for Robot Exploration

- Z. Zhang, T. Henderson, V. Sze, S. Karaman, "FSMI: Fast computation of Shannon Mutual Information for information-theoretic mapping," IEEE International Conference on Robotics and Automation (ICRA), May 2019.
- P. Li*, Z. Zhang*, S. Karaman, V. Sze, "High-throughput Computation of Shannon Mutual Information on Chip," Robotics: Science and Systems (RSS), June 2019.

• Low Power Time of Flight Imaging

- J. Noraky, V. Sze, "Low Power Depth Estimation of Rigid Objects for Time-of-Flight Imaging," IEEE Transactions on Circuits and Systems for Video Technology (TCSVT), 2019.
- J. Noraky, V. Sze, "Depth Estimation of Non-Rigid Objects For Time-Of-Flight Imaging," IEEE International Conference on Image Processing (ICIP), October 2018.
- J. Noraky, V. Sze, "Low Power Depth Estimation for Time-of-Flight Imaging," IEEE International Conference on Image Processing (ICIP), September 2017.

Monitoring Neurodegenerative Disorders Using a Phone

- H.-Y. Lai, G. Saavedra Peña, C. Sodini, T. Heldt, V. Sze, "Enabling Saccade Latency Measurements with Consumer-Grade Cameras," IEEE International Conference on Image Processing (ICIP), October 2018.
- G. Saavedra Peña, H.-Y. Lai, V. Sze, T. Heldt, "Determination of saccade latency distributions using video recordings from consumer-grade devices," IEEE International Engineering in Medicine and Biology Conference (EMBC), 2018.
- H.-Y. Lai, G. Saavedra Peña, C. Sodini, V. Sze, T. Heldt, "Measuring Saccade Latency Using Smartphone Cameras," IEEE Journal of Biomedical and Health Informatics (JBHI), March 2020.

Vivienne Sze (**v**@eems_mit)