

# Energy-Efficient Hardware for Embedded Vision and Deep Convolutional Neural Networks

Vivienne Sze

Massachusetts Institute of Technology



Contact Info

email: [sze@mit.edu](mailto:sze@mit.edu)

website: [www.rle.mit.edu/eems](http://www.rle.mit.edu/eems)

# Outline

- **What is Deep Learning?**
- **How is Deep Learning being used?**
- **Why is Edge Computing important?**
- **How can we enable Deep Learning at the Edge?**

# AI and Machine Learning

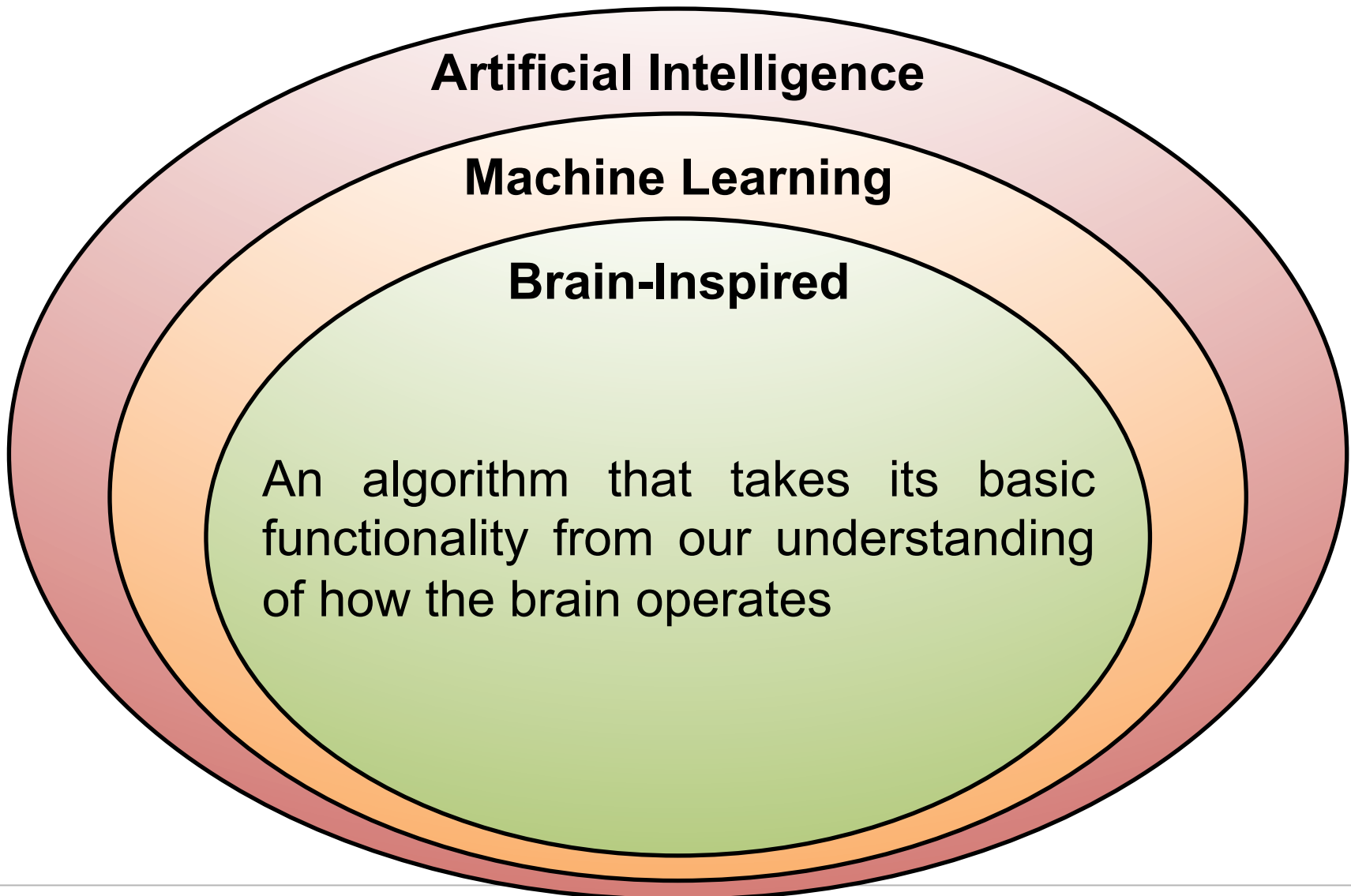
**Artificial Intelligence**

**Machine Learning**

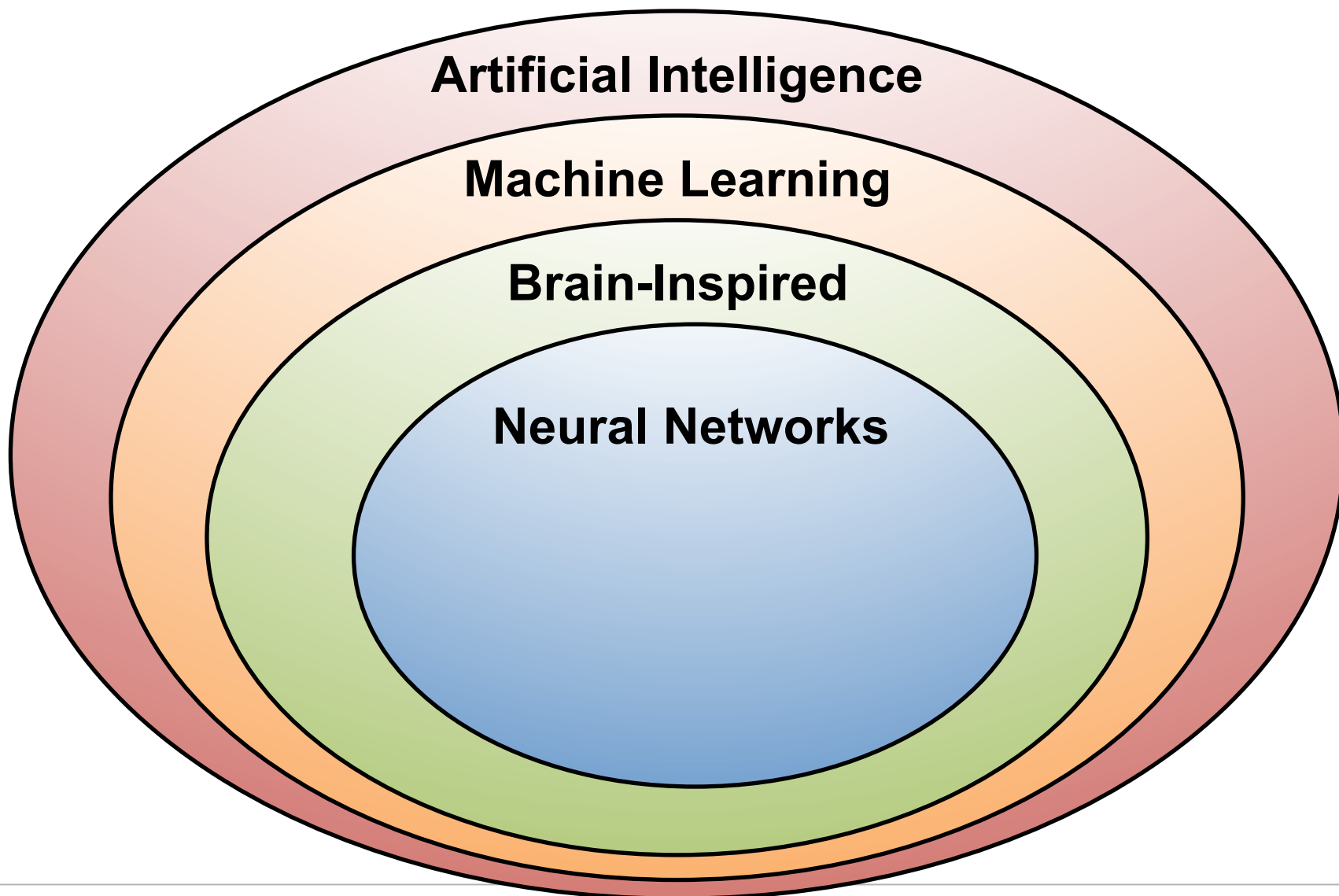
“Field of study that gives computers the ability to learn without being explicitly programmed”

– Arthur Samuel, 1959

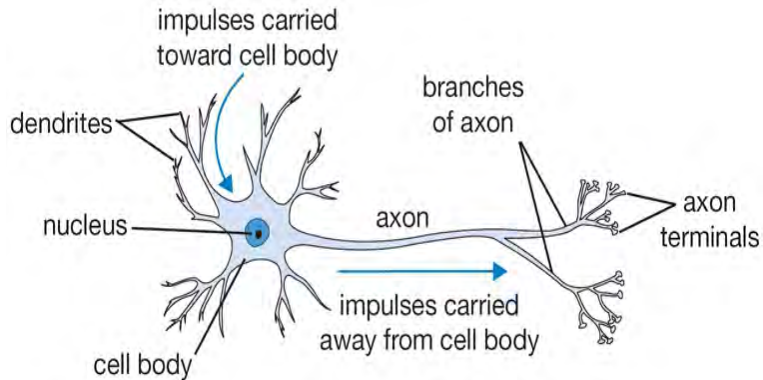
# Brain-Inspired Machine Learning



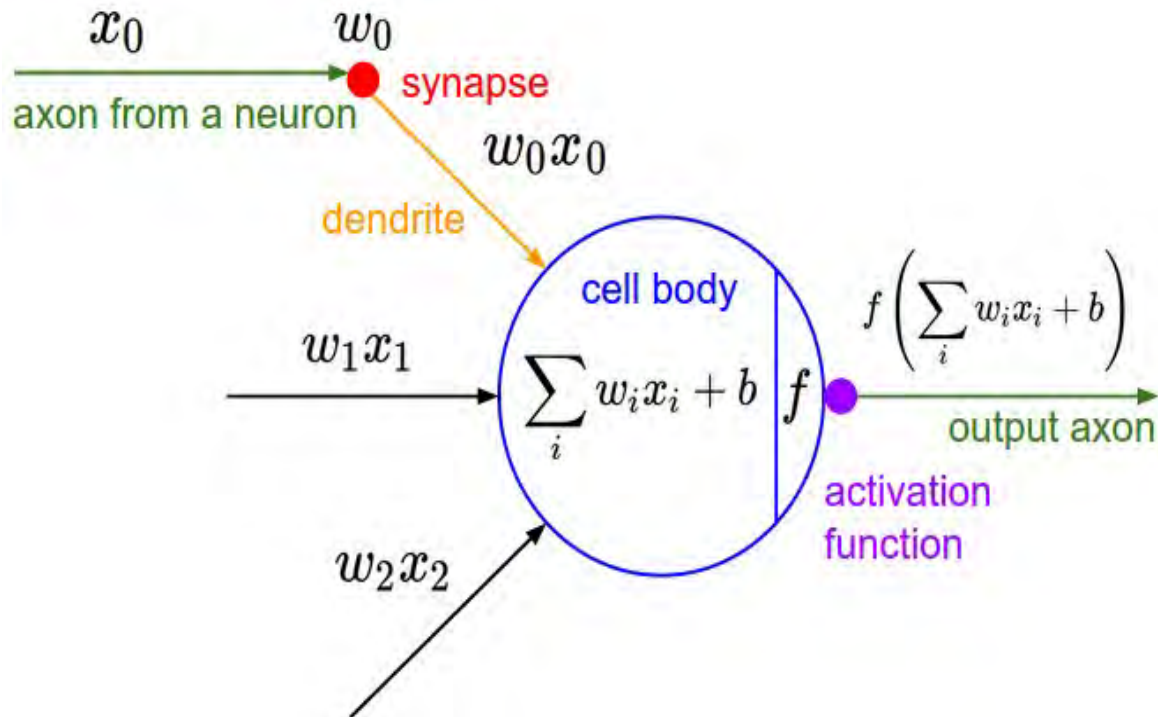
# Neural Networks



# Neural Networks: Weighted Sum

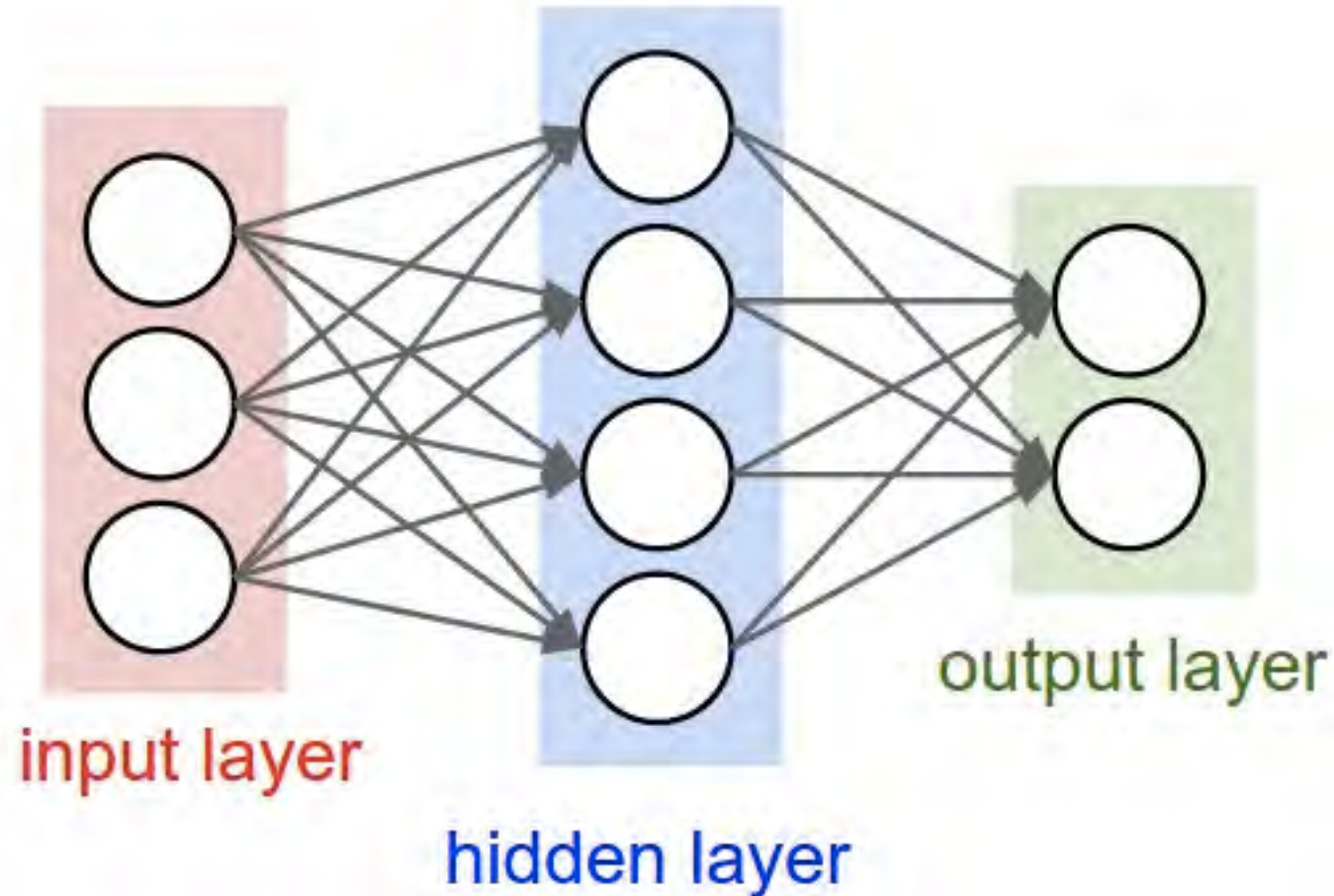


The brain contains  
 $\sim 10^{11}$  neurons connected with  
 $\sim 10^{14} - 10^{15}$  synapses

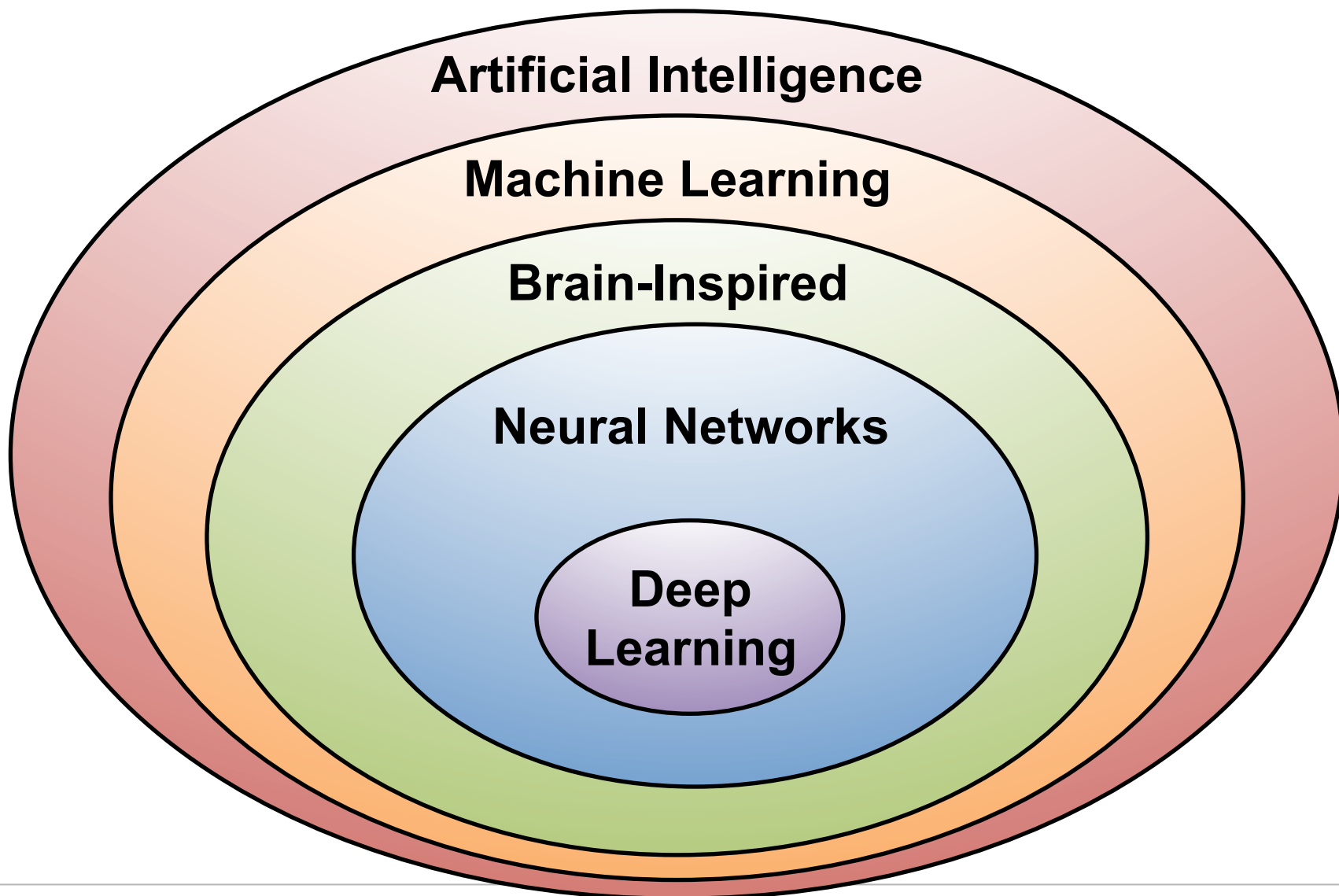




# 7 Many Weighted Sums



# Deep Learning





# 9 What is Deep Learning?

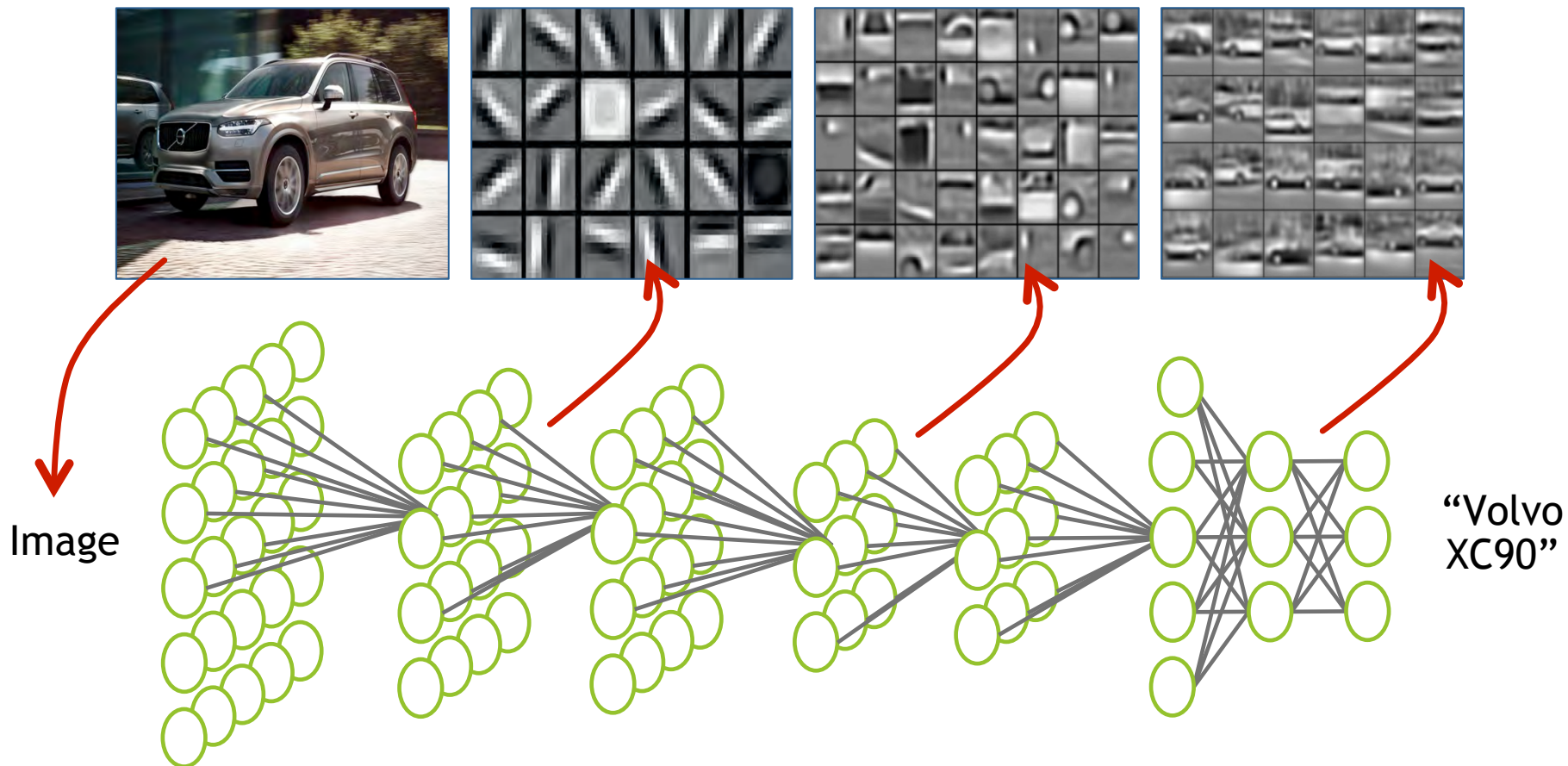


Image Source: [Lee et al., Comm. ACM 2011]

# Why is Deep Learning Hot Now?

## Big Data Availability

facebook

**350M** images uploaded per day

Walmart\*

**2.5 Petabytes** of customer data hourly

YouTube

**300 hours** of video uploaded every minute

## GPU Acceleration



## New ML Techniques



# ImageNet Challenge

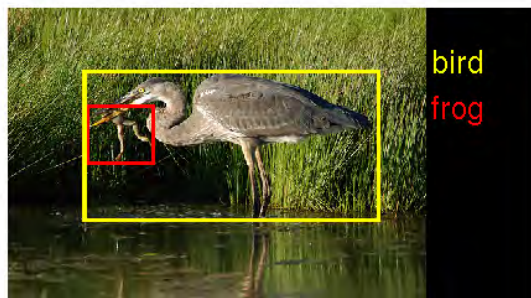
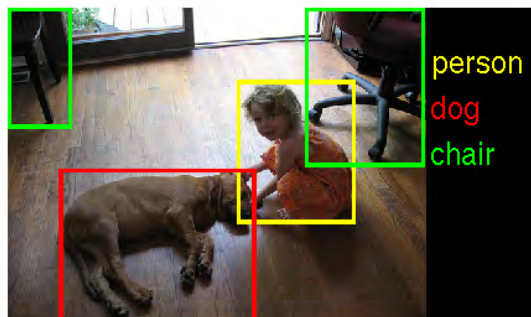
# IMAGENET

## Image Classification Task:

1.2M training images • 1000 object categories

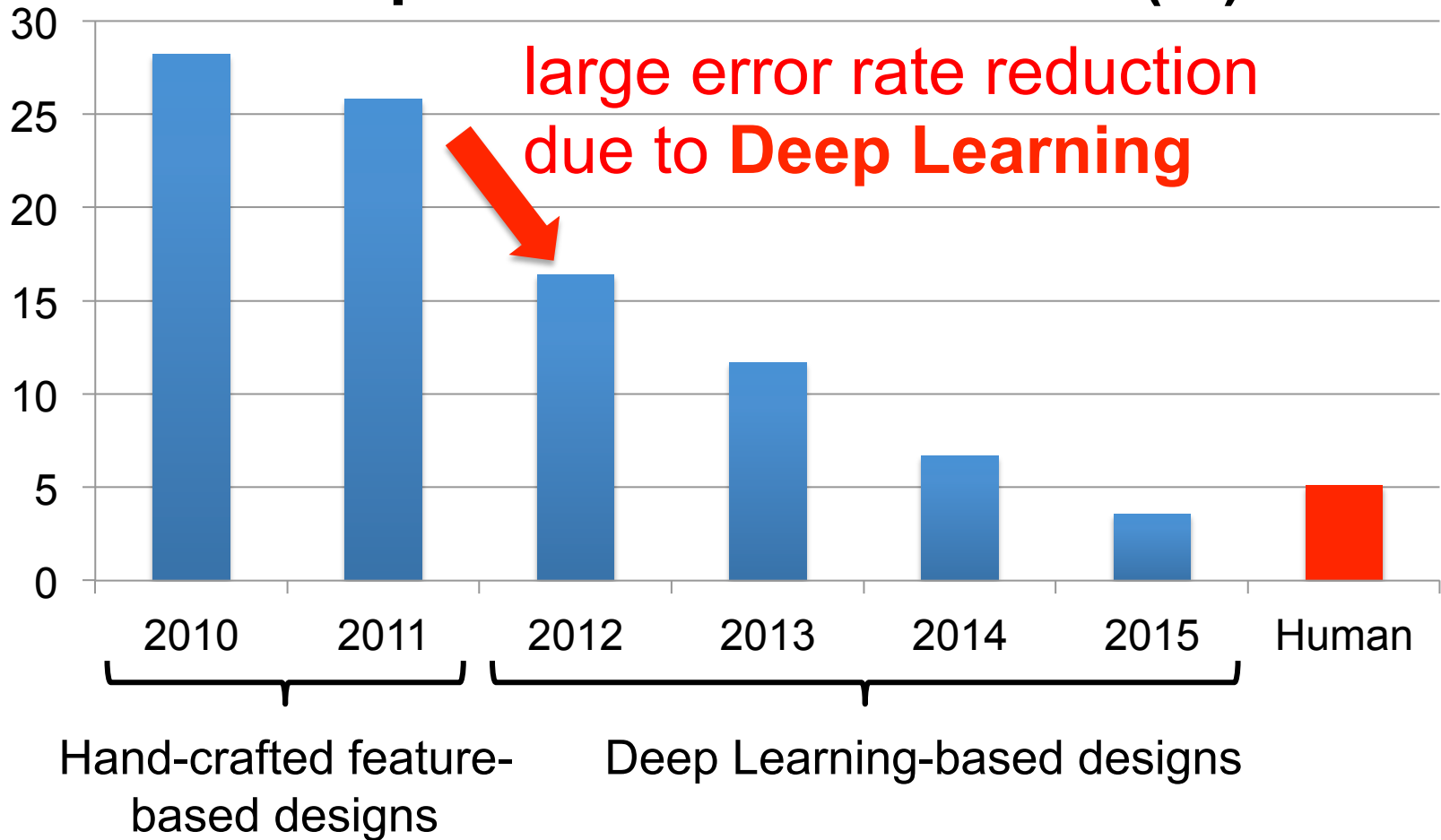
## Object Detection Task:

456k training images • 200 object categories

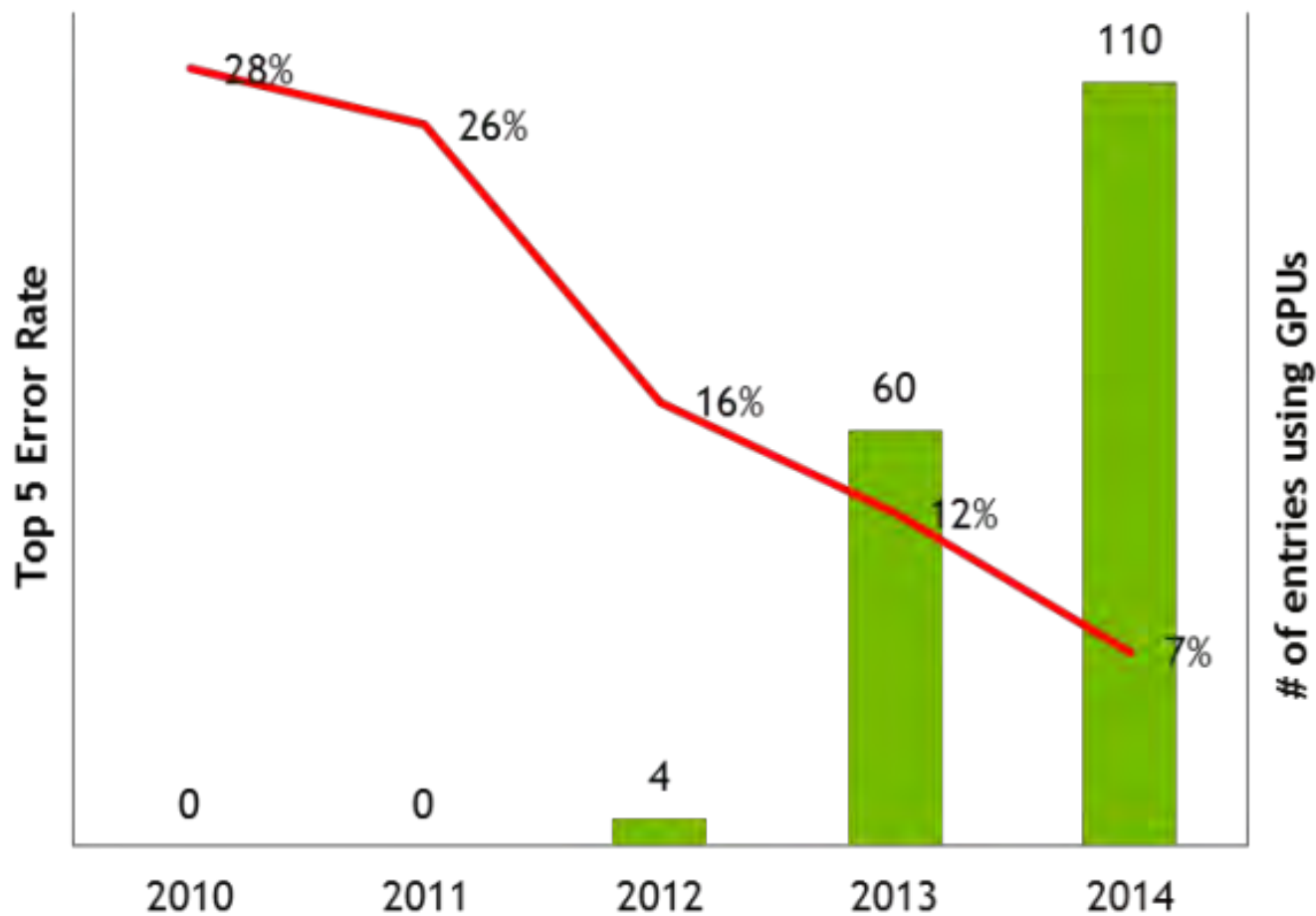


# ImageNet: Image Classification Task

## Top 5 Classification Error (%)



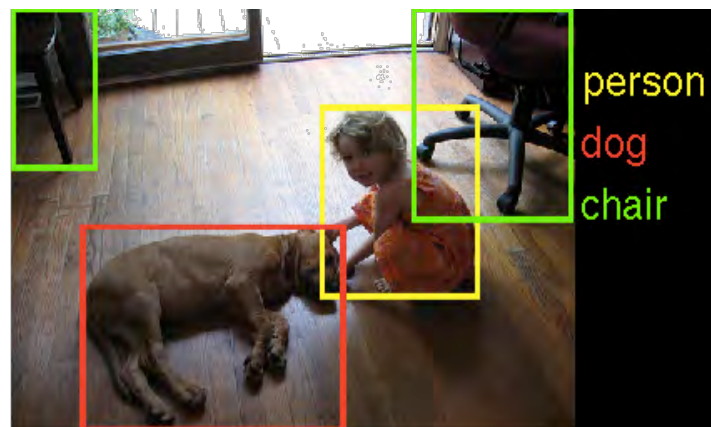
# GPU Usage for ImageNet Challenge





# Deep Learning on Images

- Image Classification
- Object Localization
- Object Detection
- Image Segmentation
- Action Recognition
- Image Generation

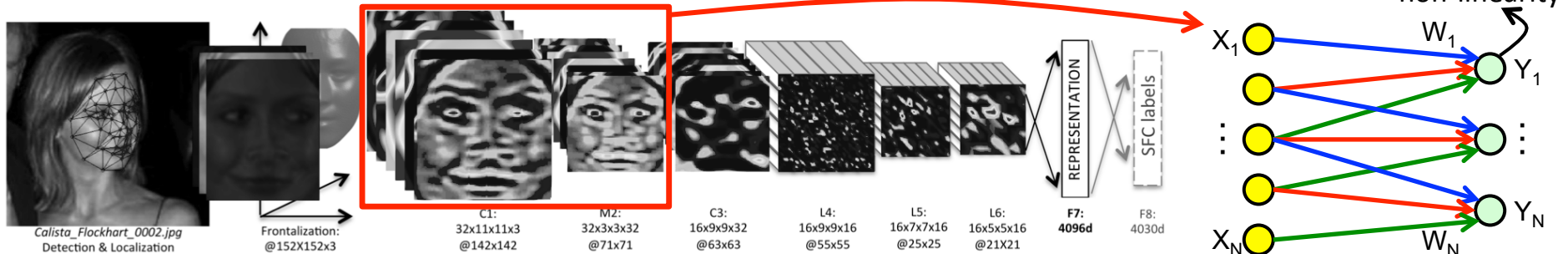




# Human or *Superhuman* Accuracy Level

- Face recognition

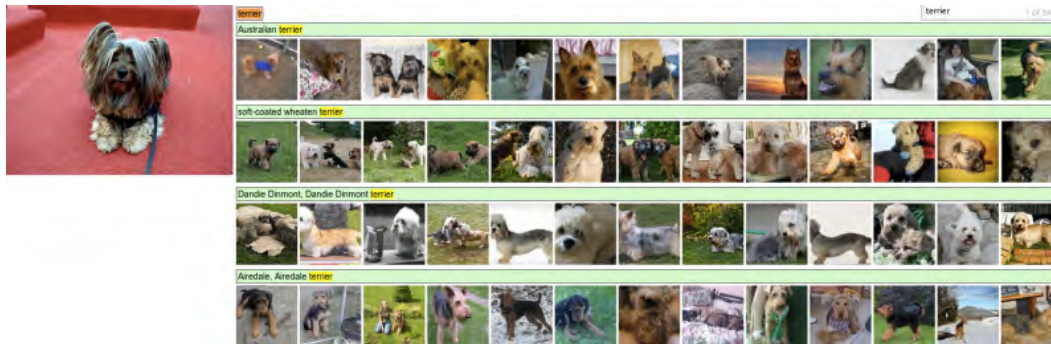
- Deep learning accuracy (97.25%) vs. Human accuracy (97.53%)



[Yaniv et al., CVPR 2014]

- Fine grained category recognition (e.g. dogs, monkeys, snakes, birds)

- Deep learning errors: 7 vs. Human errors: 28

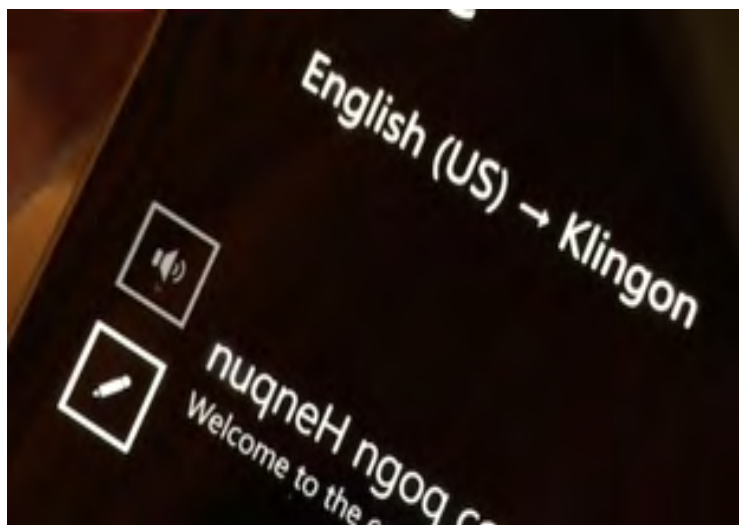


120 species of dogs

[O. Russakovsky et al., IJCV 2015]

# Deep Learning for Speech

- **Speech Recognition**
- **Natural Language Processing**
- **Speech Translation**
- **Audio Generation**



# Deep Learning on Games

## Google DeepMind AlphaGo

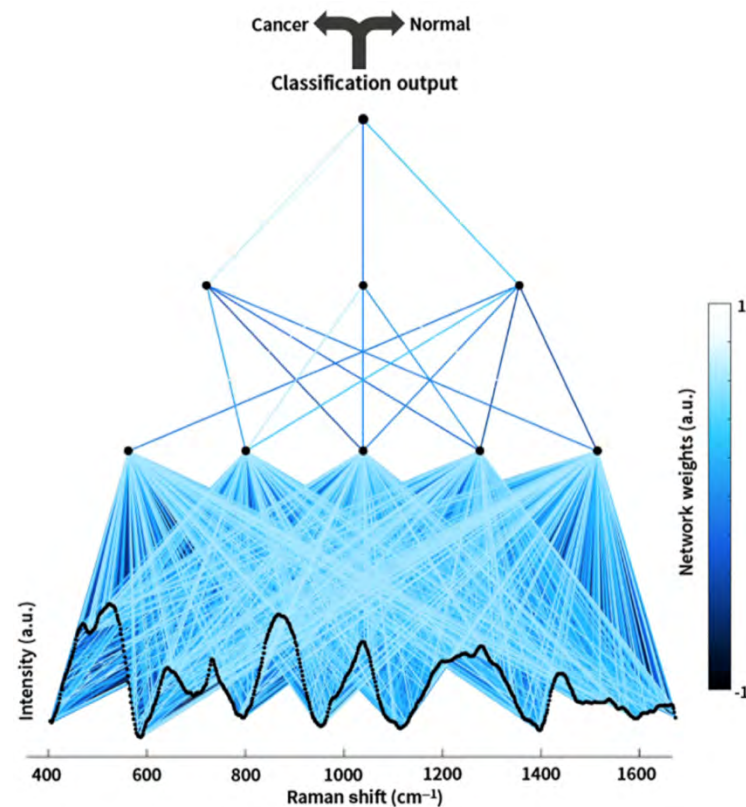
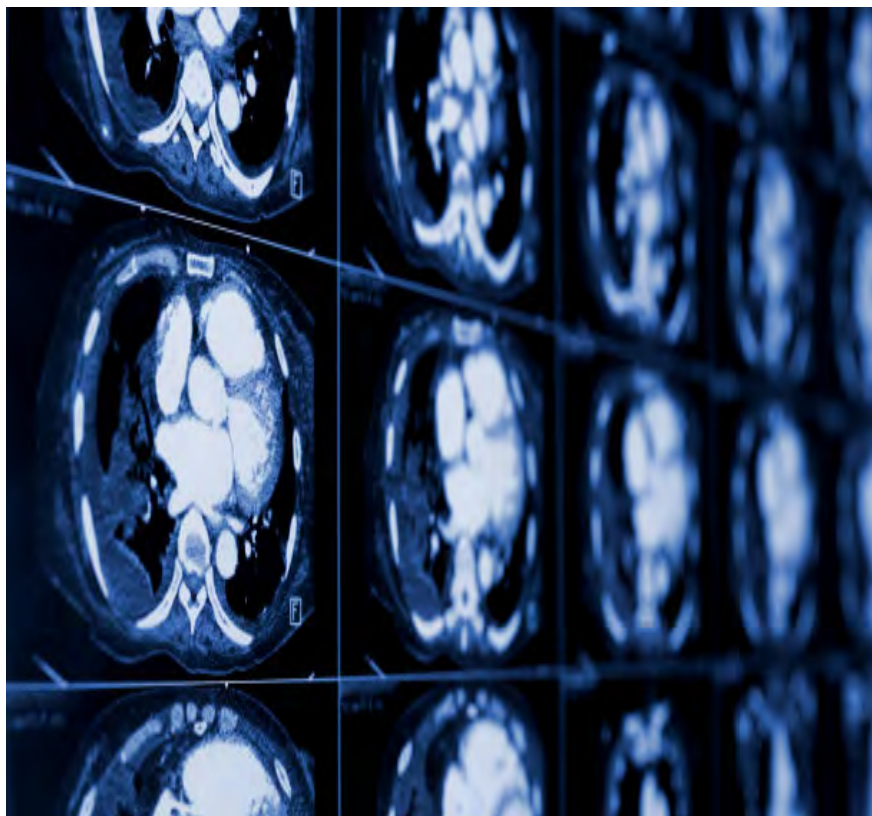
*Go is exponentially more complex than chess ( $10^{170}$  legal positions)*



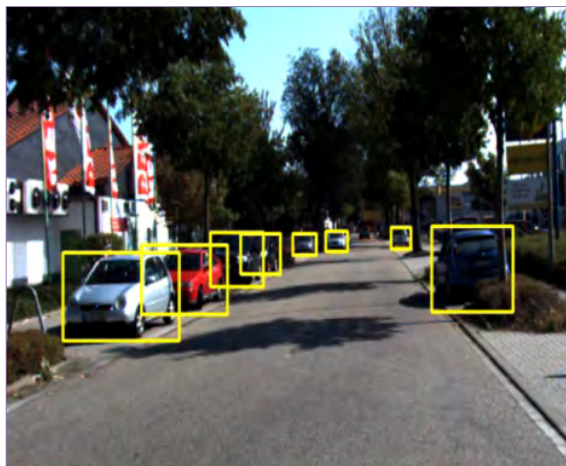


# Medical Applications of Deep Learning

- **Brain Cancer Detection**



# Deep Learning for Self-driving Cars



# Other Emerging Applications

- **Medical** (Cancer Detection, Pre-Natal)
- **Finance** (Trading, Energy Forecasting, Risk)
- **Infrastructure** (Structure Safety and Traffic)
- Weather Forecasting and Event Detection

**This talk will focus on image classification**

<http://www.nextplatform.com/2016/09/14/next-wave-deep-learning-applications/>

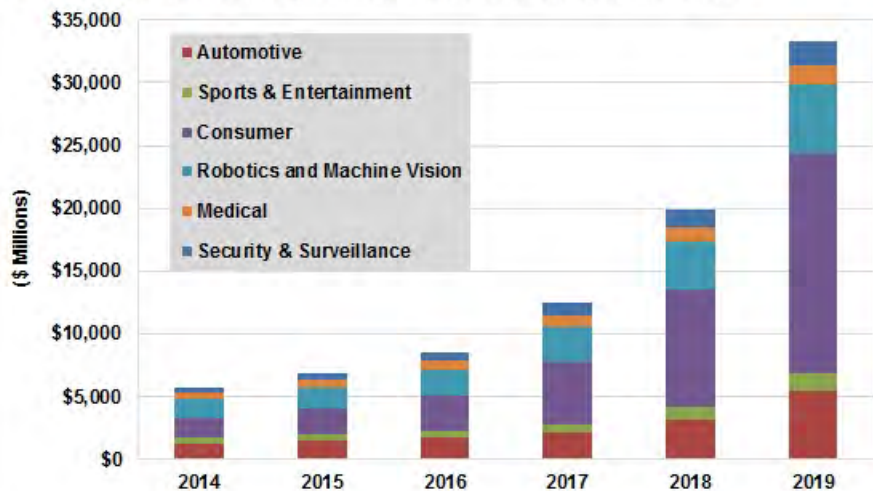


# Opportunities

## \$500B Market over 10 Years!



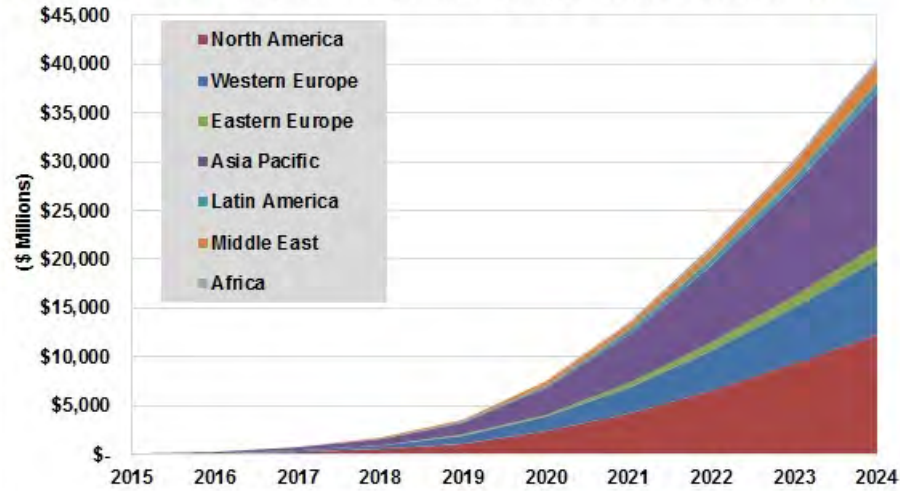
Computer Vision Revenue by Vertical Market, World Markets: 2014-2019



Source: Tractica



Cumulative Deep Learning Software Revenue by Region, World Markets: 2015-2024



Source: Tractica

# Opportunities

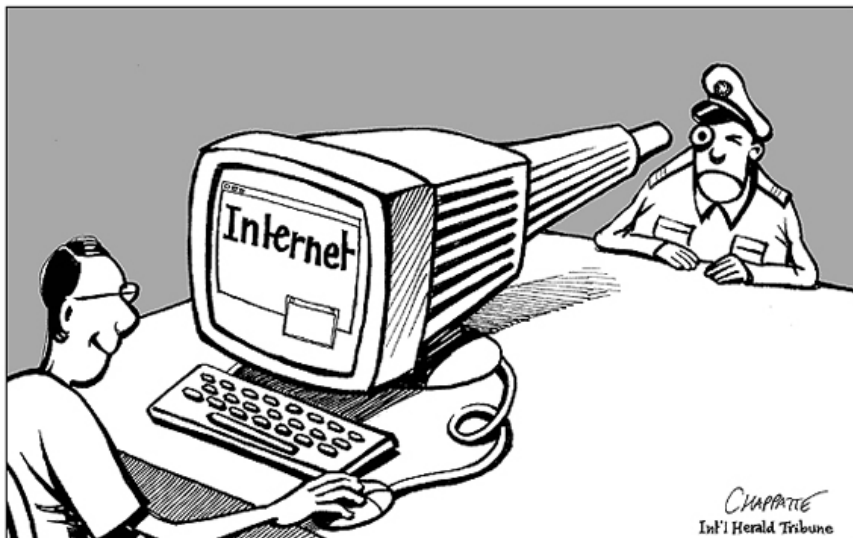
From EE Times – September 27, 2016

”Today the job of training machine learning models is limited by compute, if we had faster processors we’d run bigger models...in practice we train on a reasonable subset of data that can finish in a matter of months. We could use improvements of several orders of magnitude – 100x or greater.”

– Greg Diamos, Senior Researcher, SVAIL,  
Baidu

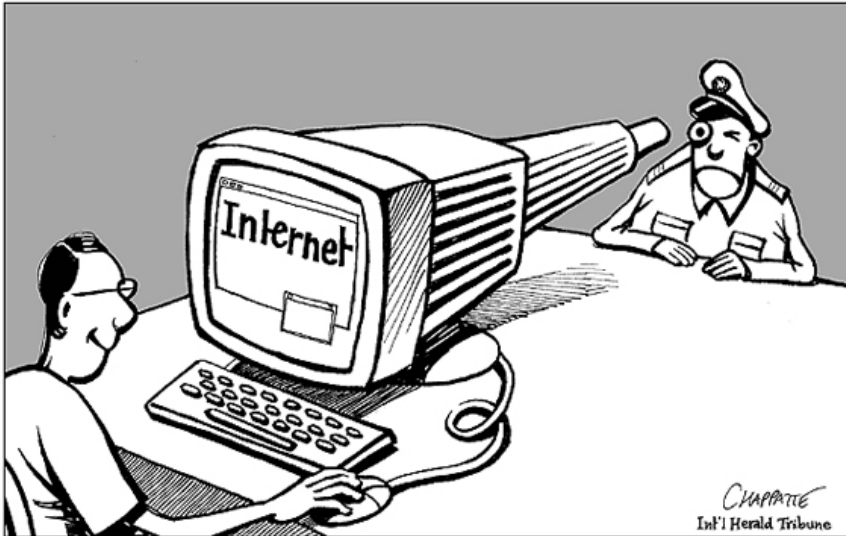
# Processing at “Edge” instead of the “Cloud”

## Privacy



# Processing at “Edge” instead of the “Cloud”

## Privacy



## Latency



Sensor



Actuator

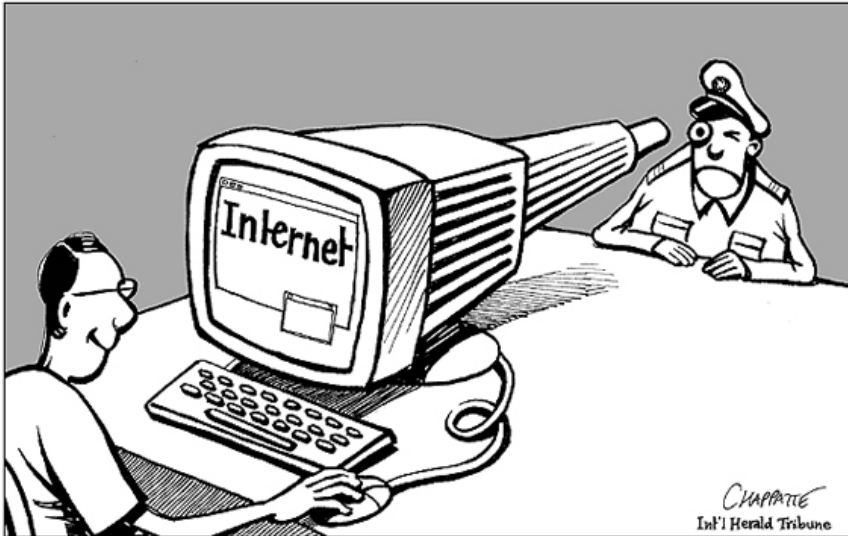


Cloud

Image source: ericsson.com

# Processing at “Edge” instead of the “Cloud”

## Privacy



## Communication



Image source: [www.theregister.co.uk](http://www.theregister.co.uk)

## Latency



Sensor



Actuator



Cloud

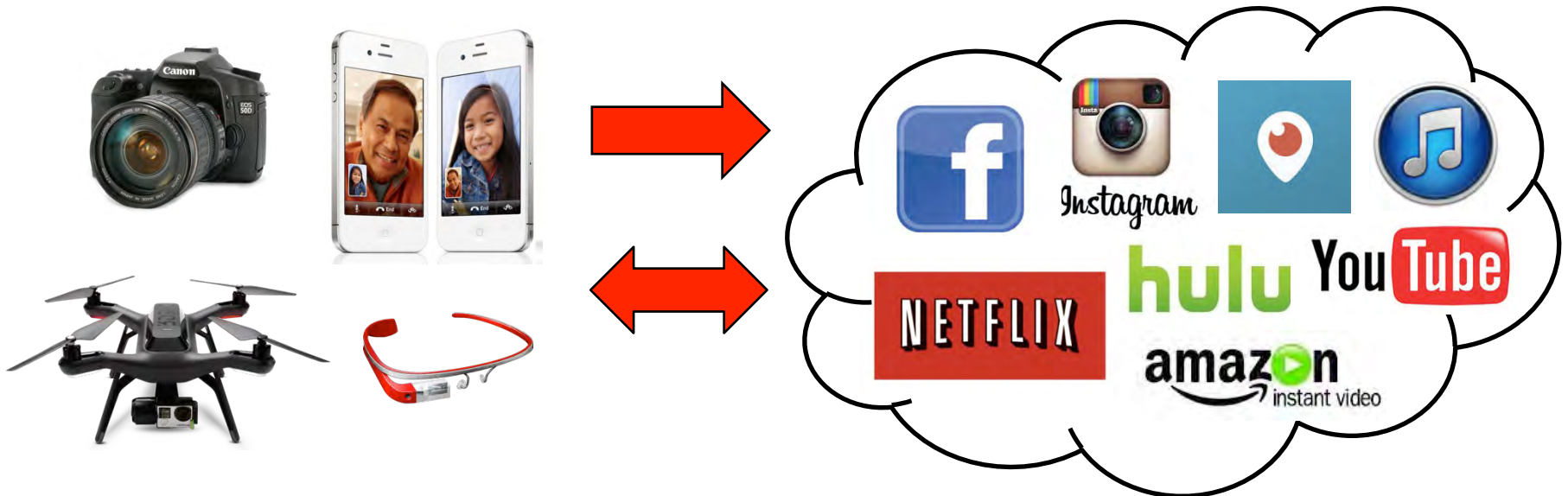
Image source: [ericsson.com](http://ericsson.com)

# Video is the Biggest Big Data

Over 70% of today's Internet traffic is video

Over 300 hours of video uploaded to YouTube **every minute**

Over 500 million hours of video surveillance collected **every day**



*Energy limited due  
to battery capacity*

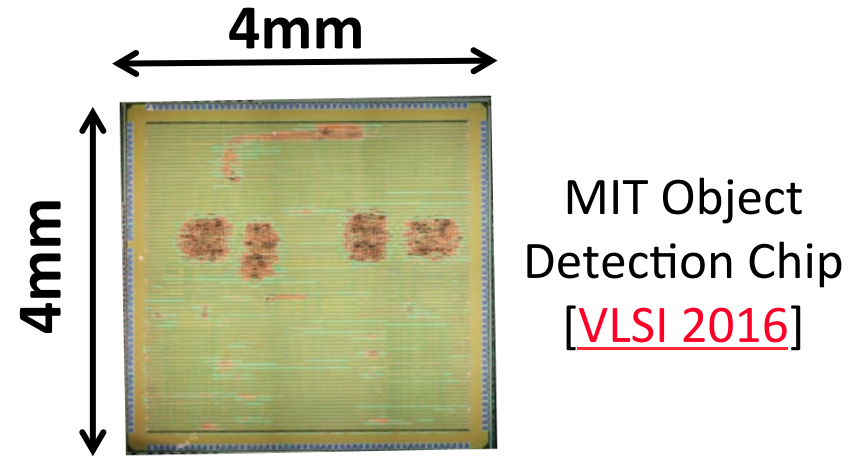
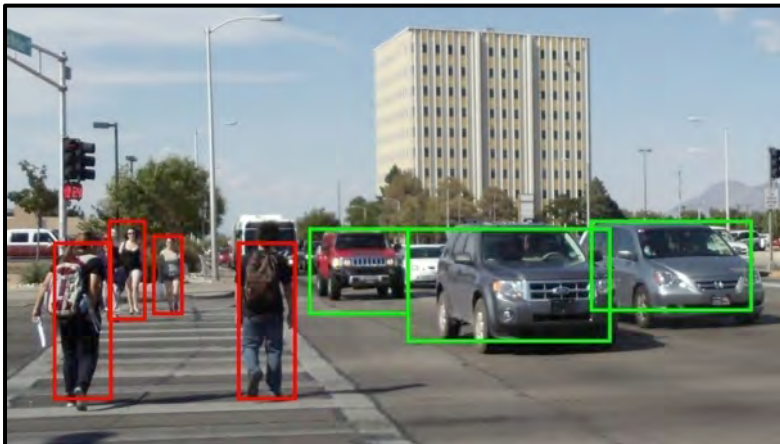
*Power limited due  
to heat dissipation*

Need energy-efficient pixel processing!

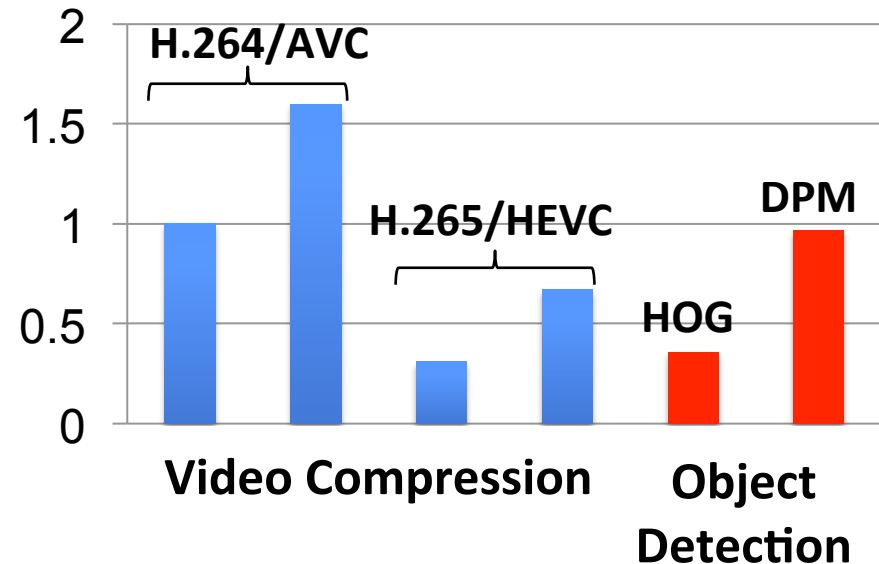


# Typical Constraints on Video Coding

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



## Energy

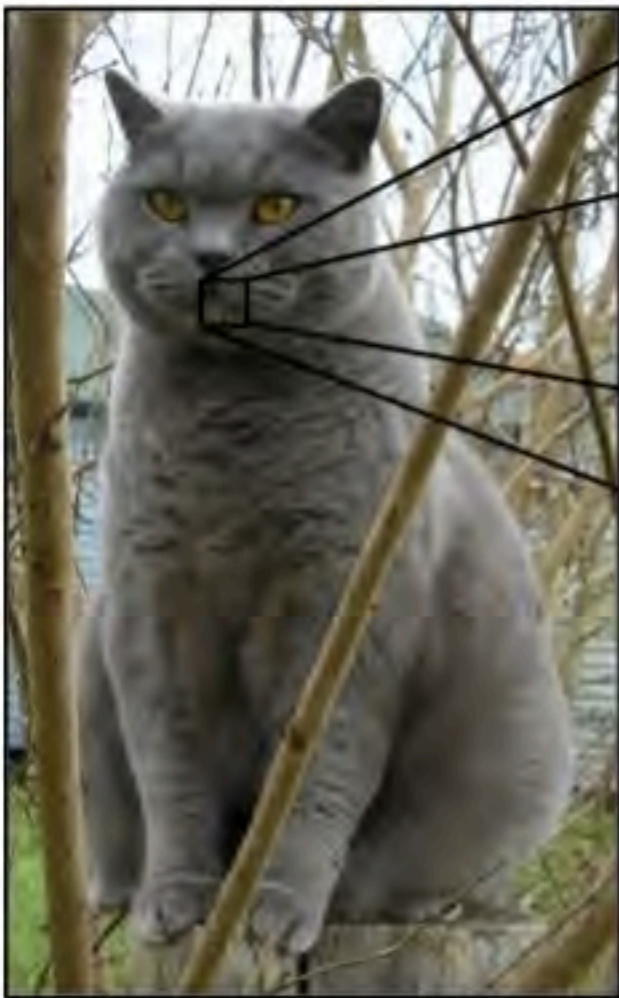


# Why is Vision Difficult?



Cat

# Why is Vision Difficult?



08	02	22	97	38	15	00	40	00	75	04	05	07	78	52	22	50	77	81	50
49	49	99	40	17	81	18	57	60	87	17	40	98	43	69	44	04	56	62	00
81	49	81	73	55	79	14	29	93	71	40	67	52	03	30	03	49	13	36	65
52	70	95	23	04	60	31	42	89	74	68	36	01	32	36	71	37	02	36	91
22	31	16	71	51	67	83	59	41	92	36	54	22	40	40	28	66	33	13	80
26	47	30	00	59	03	45	02	44	75	33	33	78	36	84	20	35	17	12	50
32	98	81	28	64	23	67	10	26	38	40	47	59	34	70	66	18	38	64	70
67	26	20	68	02	62	12	20	95	63	94	39	63	08	40	91	66	49	94	21
24	55	38	05	66	73	99	26	97	17	78	78	96	83	14	88	34	89	63	72
21	36	23	09	75	00	76	44	20	45	35	14	00	61	33	97	34	31	33	95
72	17	53	26	22	75	31	67	15	94	03	80	04	62	16	14	09	53	56	92
16	39	05	42	96	35	31	47	55	58	88	24	00	17	54	24	36	29	85	57
06	56	00	48	35	71	89	07	05	44	44	37	44	60	21	58	51	54	17	58
19	80	81	68	05	94	47	69	28	73	92	13	86	52	17	77	04	89	55	40
04	52	08	83	97	38	99	16	07	97	57	32	16	26	26	79	33	27	98	86
77	24	40	67	57	62	20	72	03	46	33	67	46	55	12	32	63	93	53	69
04	42	16	73	50	85	39	11	24	84	72	18	08	46	29	32	40	62	76	36
20	69	06	41	72	30	23	88	04	00	89	69	82	67	39	85	74	04	36	16
20	73	35	29	78	31	90	01	74	31	49	71	49	30	11	16	23	57	05	54
01	70	84	71	83	51	34	49	16	92	33	48	61	43	52	01	89	19	67	48

What the computer sees

**Computer vision requires more processing than video compression**

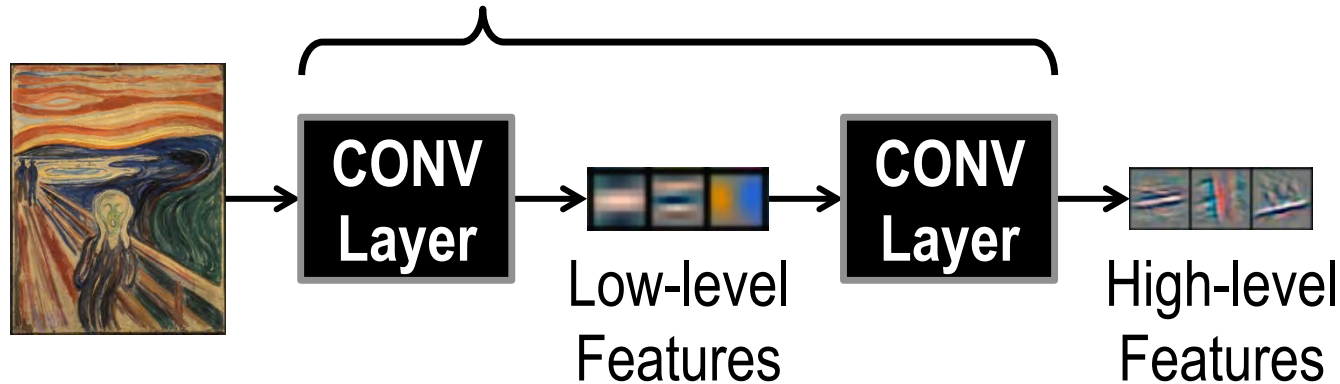
# Eyeriss: Energy-Efficient Hardware for DCNNs

Yu-Hsin Chen, Tushar Krishna, Joel Emer, Vivienne Sze, ISSCC 2016 [[paper](#)] / ISCA 2016 [[paper](#)]



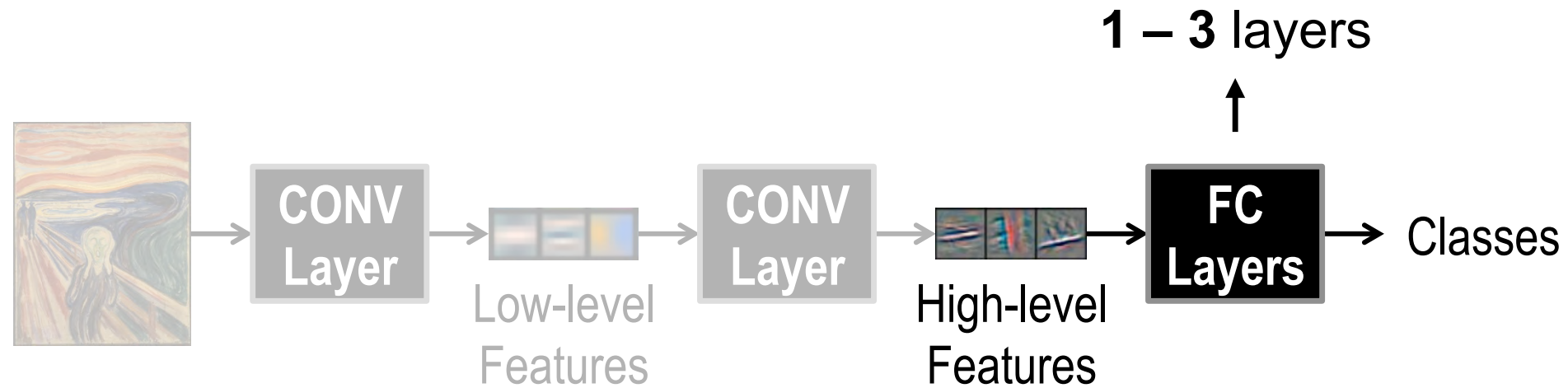
# Deep Convolutional Neural Networks

Modern *deep* CNN: up to **1000** CONV layers

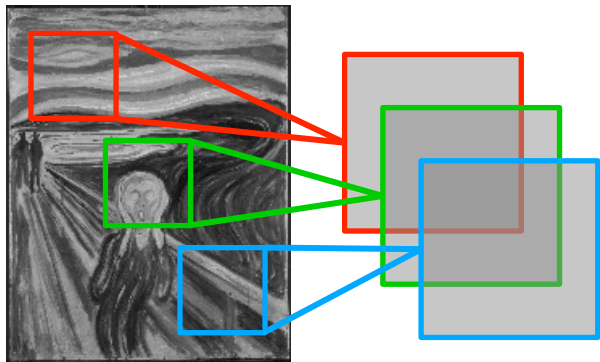
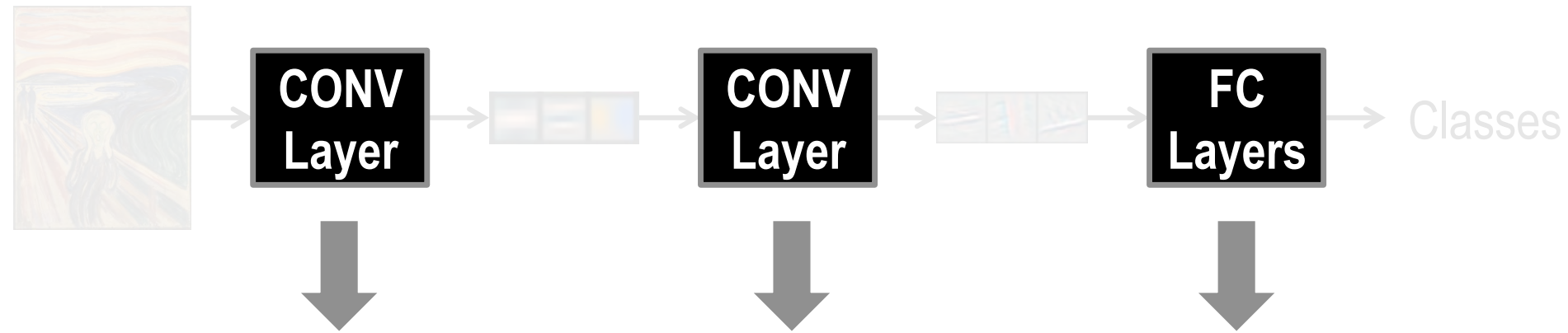




# Deep Convolutional Neural Networks



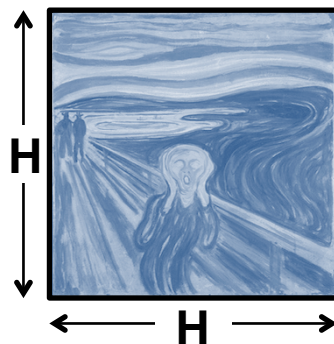
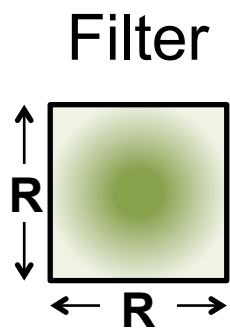
# Deep Convolutional Neural Networks



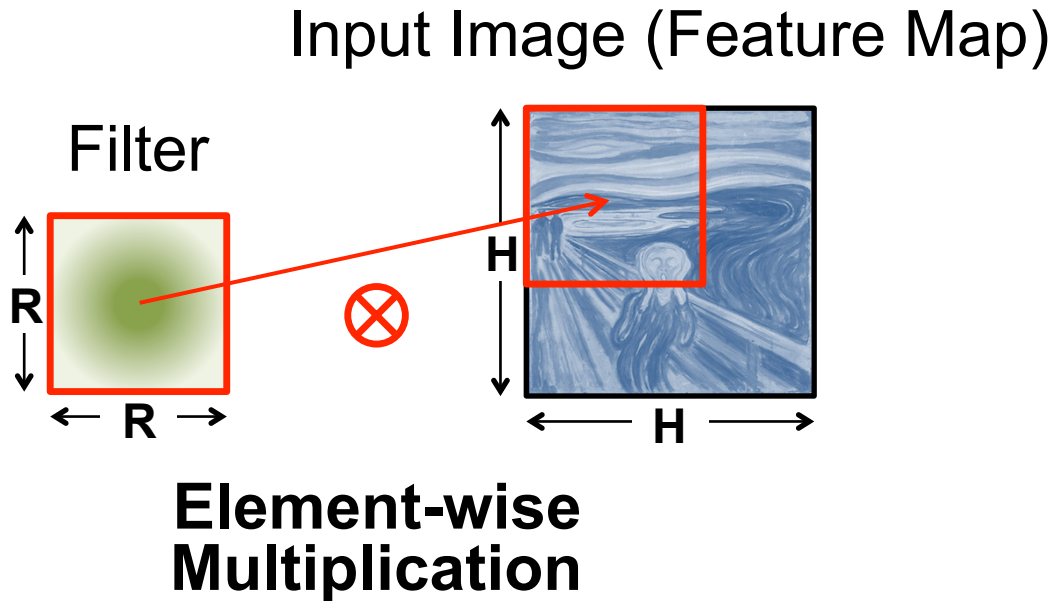
**Convolutions** account for more than 90% of overall computation, dominating **runtime** and **energy consumption**

# High-Dimensional CNN Convolution

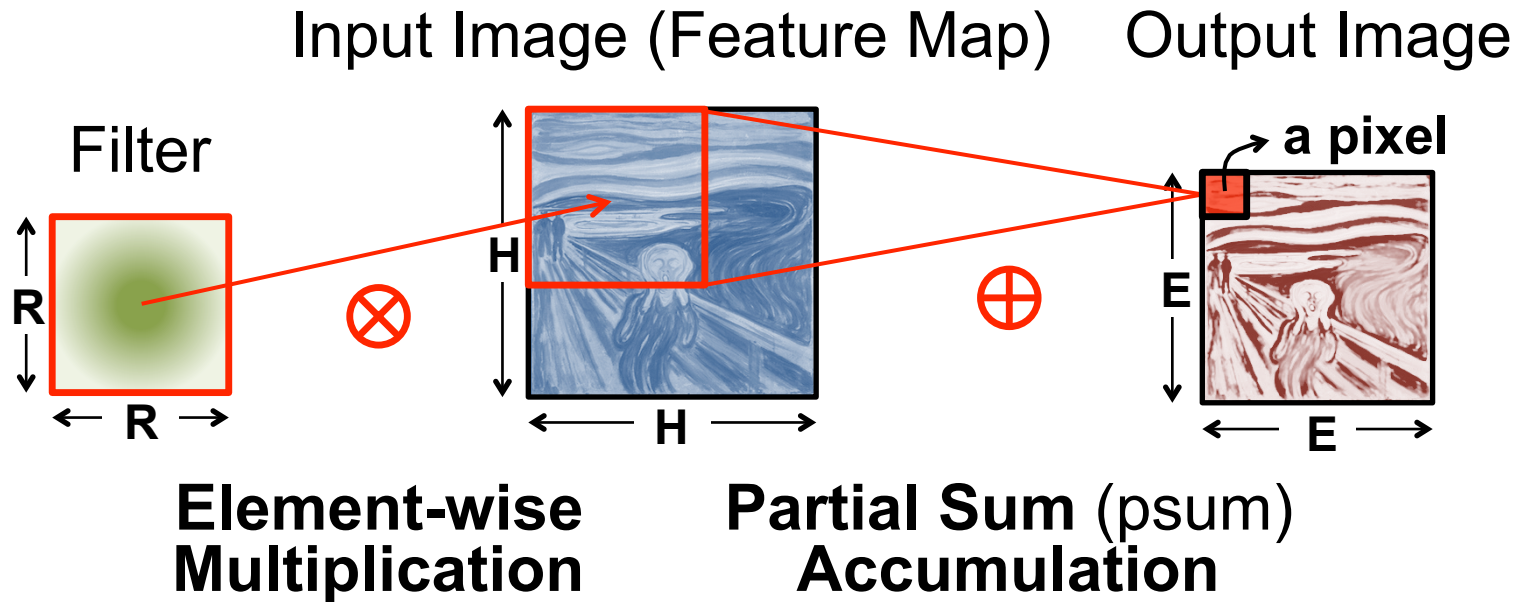
Input Image (Feature Map)



# High-Dimensional CNN Convolution

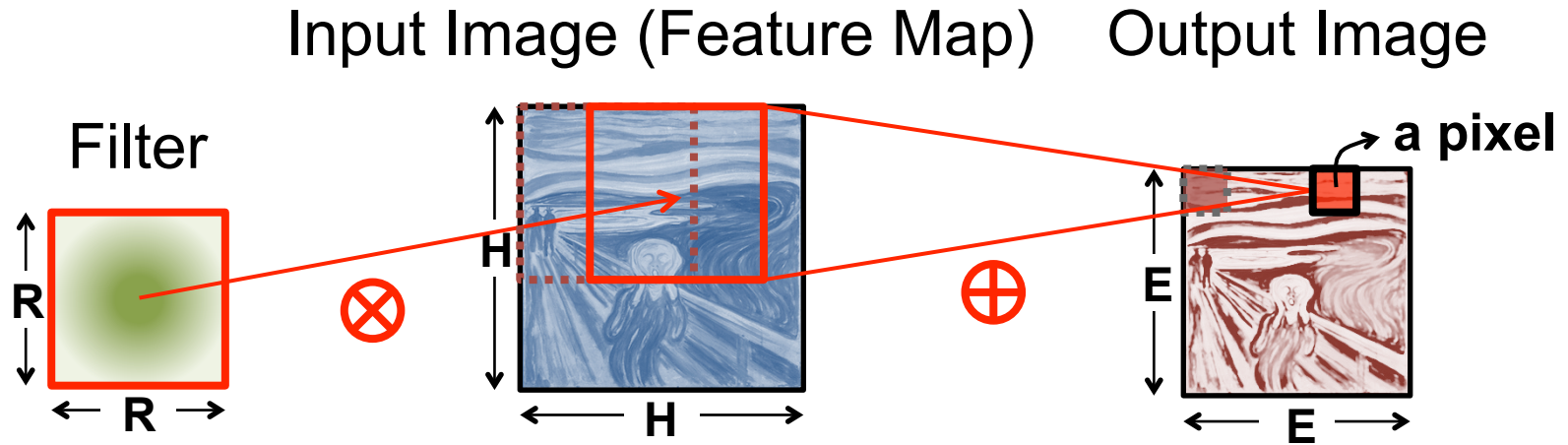


# High-Dimensional CNN Convolution



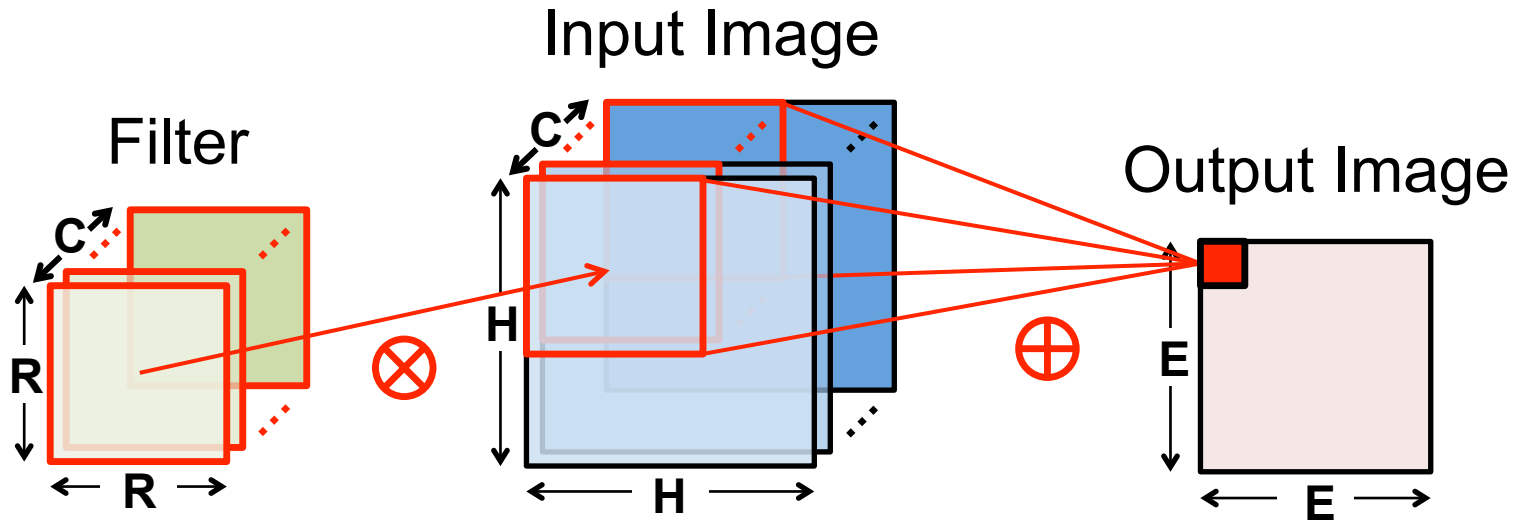


# High-Dimensional CNN Convolution



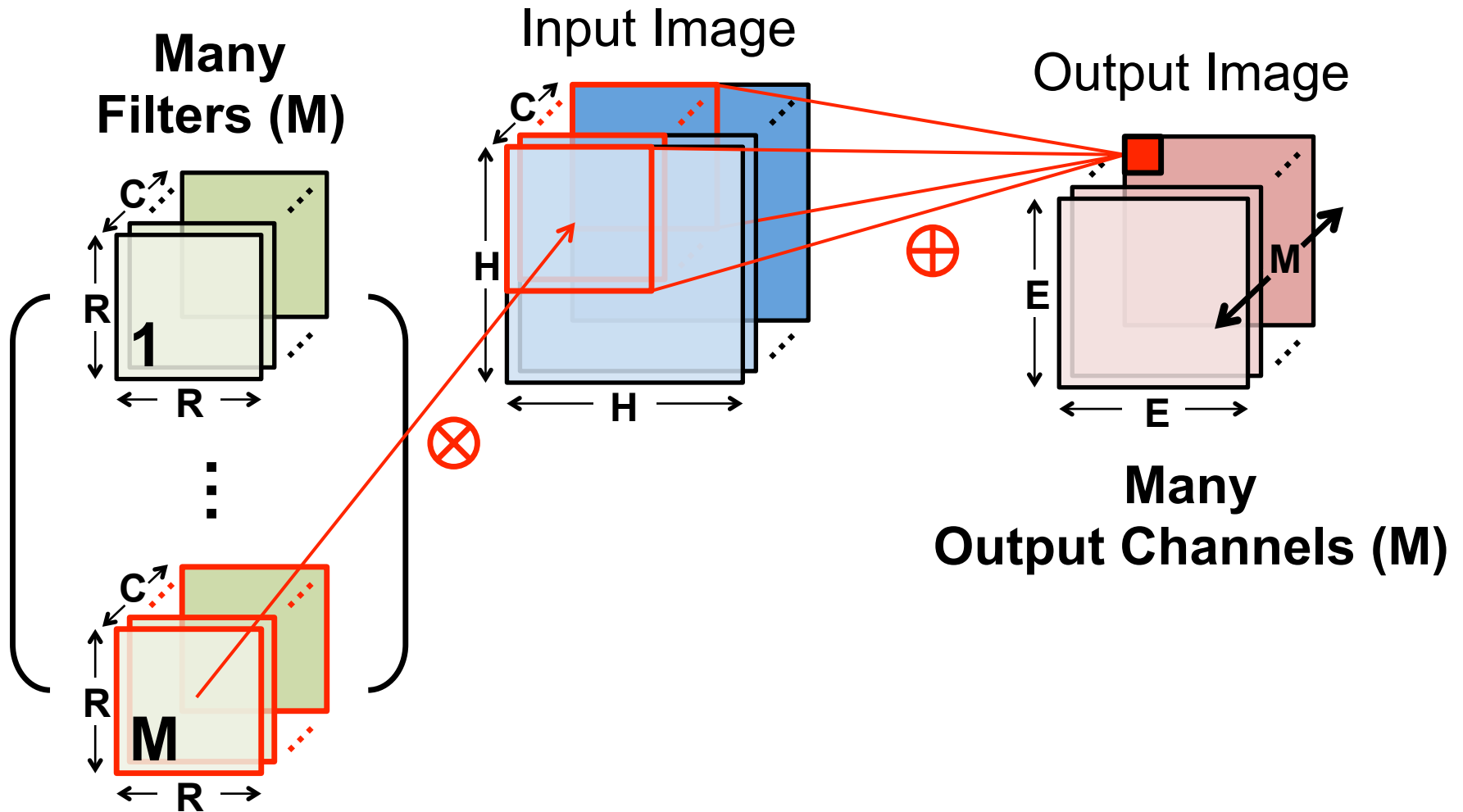
**Sliding Window Processing**

# High-Dimensional CNN Convolution

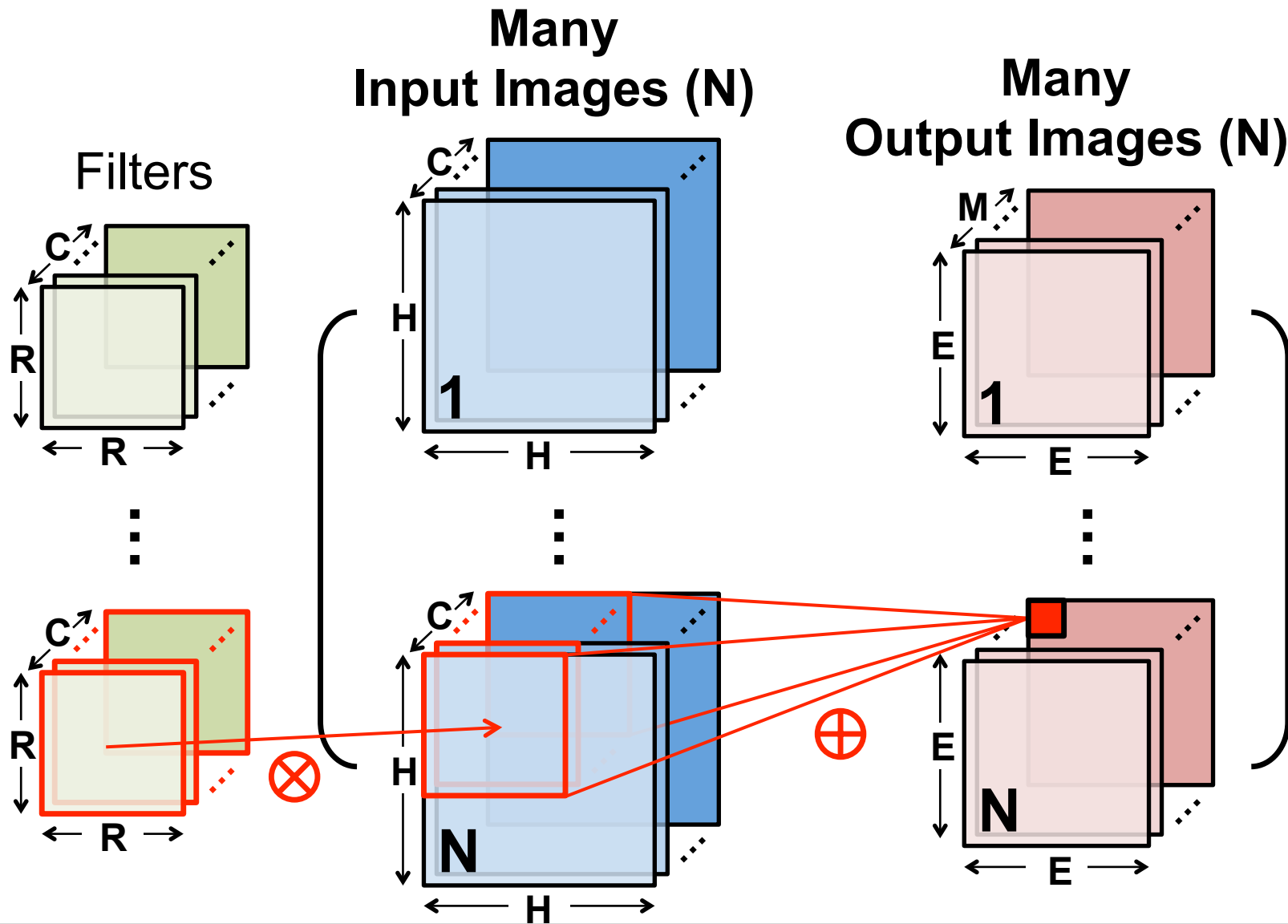


**Many Input Channels (C)**

# High-Dimensional CNN Convolution



# High-Dimensional CNN Convolution

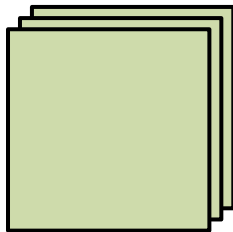


# Large Sizes with Varying Shapes

## AlexNet<sup>1</sup> 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



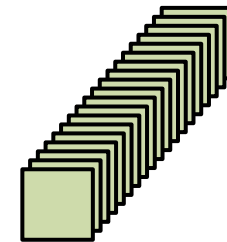
34k Params  
105M MACs

Layer 2



307k Params  
224M MACs

Layer 3



885k Params  
150M MACs

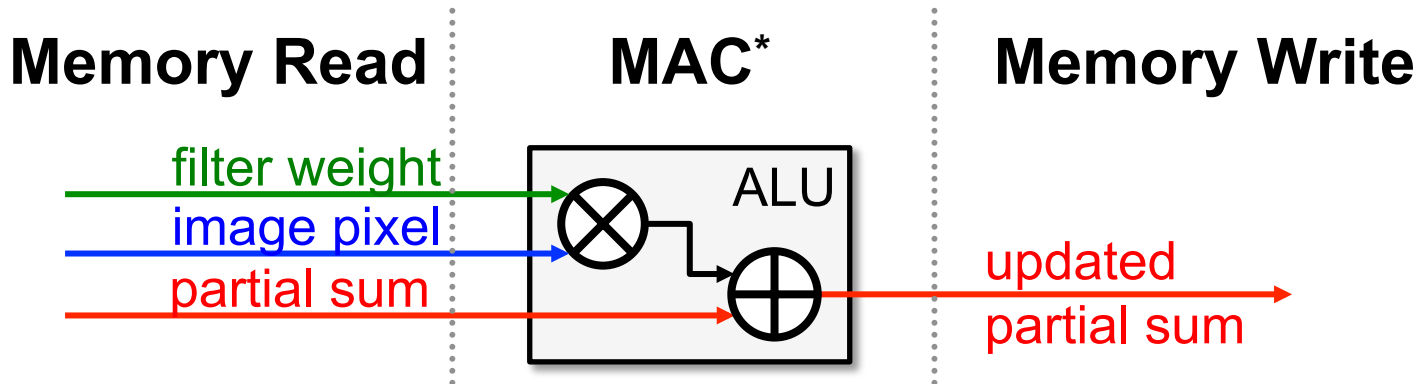


# Properties We Can Leverage

- Operations exhibit **high parallelism**  
→ **high throughput** possible

# Properties We Can Leverage

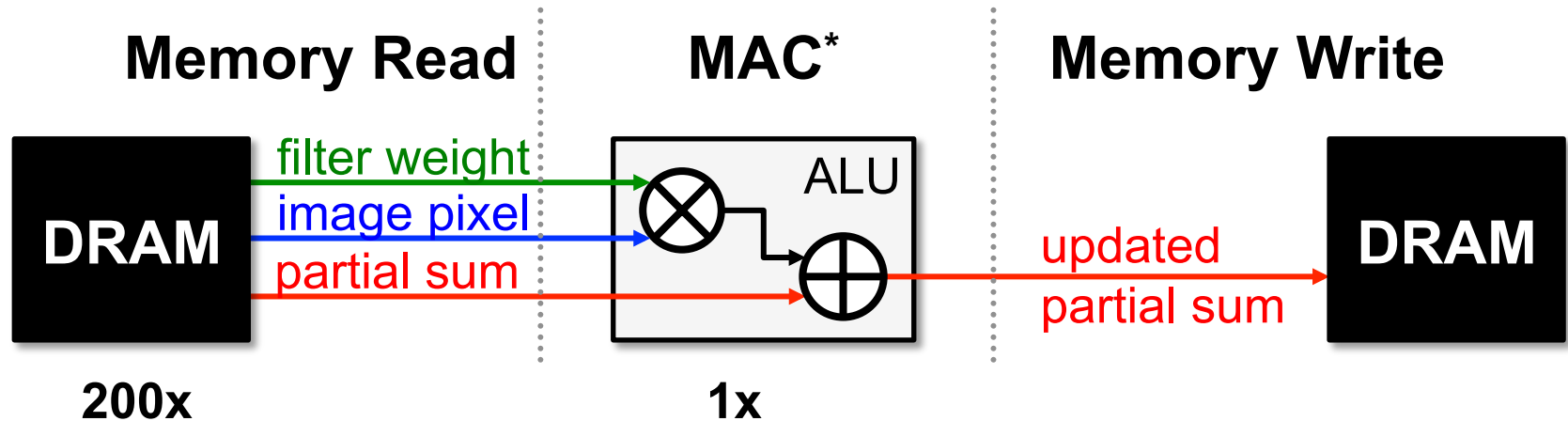
- Operations exhibit **high parallelism**  
→ **high throughput** possible
- Memory Access is the Bottleneck



\* multiply-and-accumulate

# Properties We Can Leverage

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

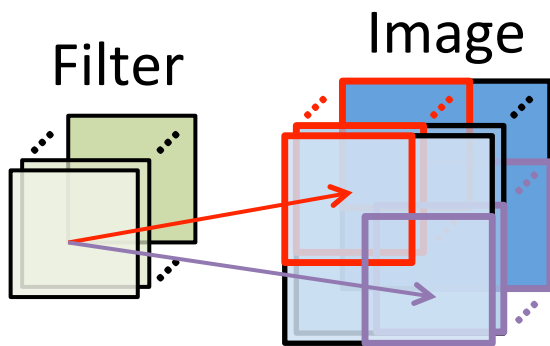


Worst Case: all memory R/W are **DRAM** accesses

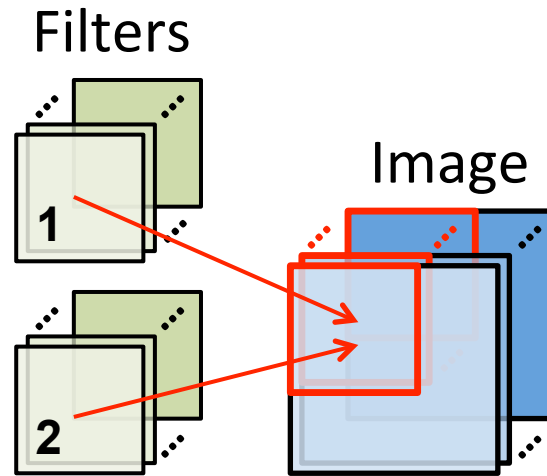
- Example: AlexNet [NIPS 2012] has **724M** MACs  
→ **2896M** DRAM accesses required

# Properties We Can Leverage

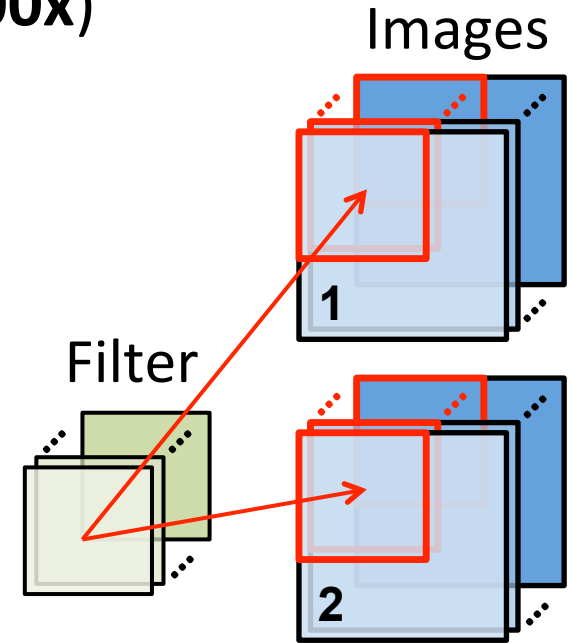
- Operations exhibit **high parallelism**  
→ **high throughput** possible
- **Input data reuse** opportunities (**up to 500x**)  
→ exploit **low-cost memory**



**Convolutional  
Reuse**  
(pixels, weights)



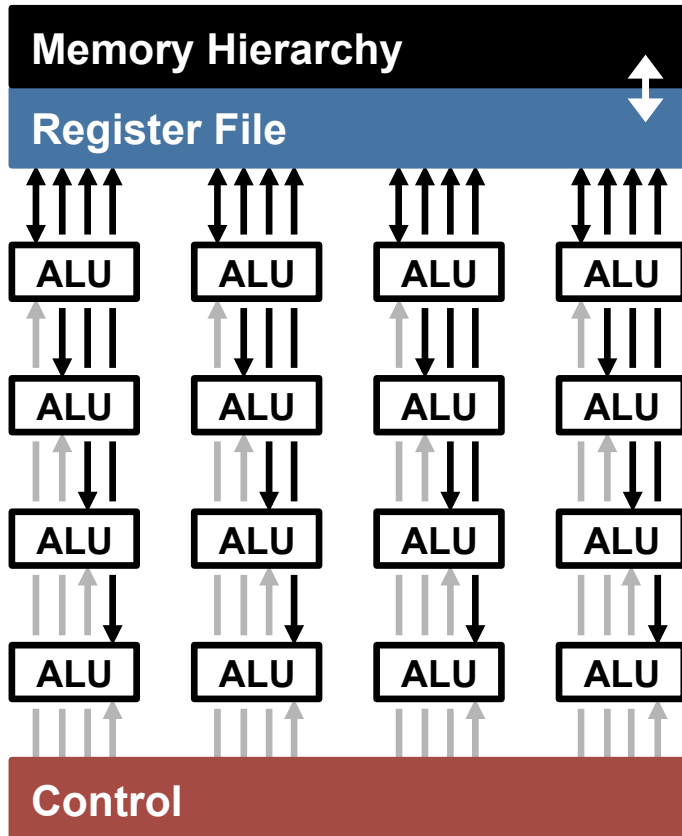
**Image  
Reuse**  
(pixels)



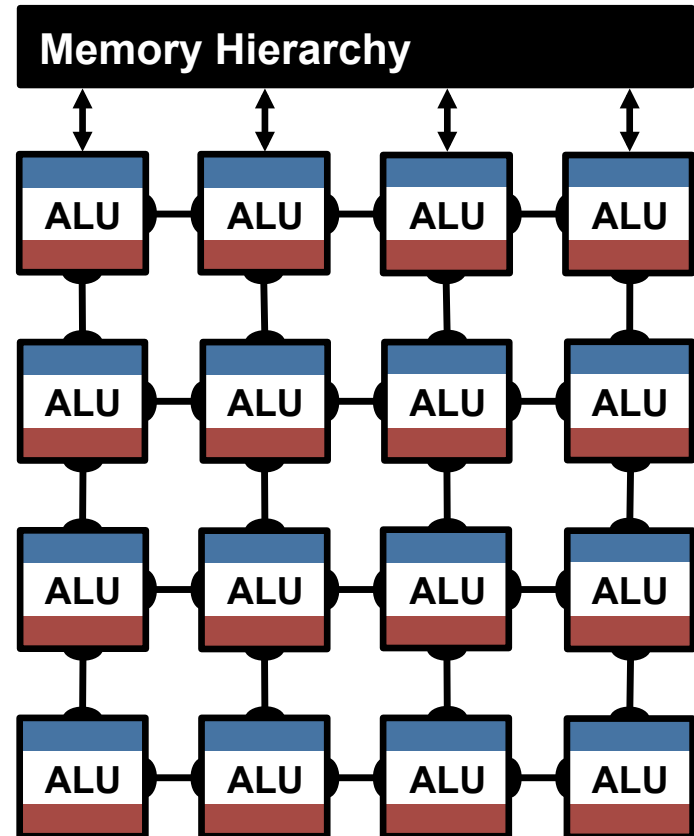
**Filter  
Reuse**  
(weights)

# Highly-Parallel Compute Paradigms

## Temporal Architecture (SIMD/SIMT)



## Spatial Architecture (Dataflow Processing)



# Advantages of Spatial Architecture

Temporal Architecture  
(SIMD/SIMT)

## Efficient Data Reuse

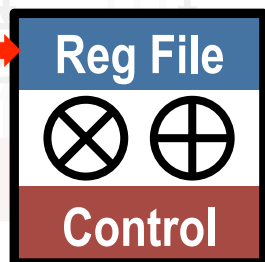
Distributed local storage (RF)

## Inter-PE Communication

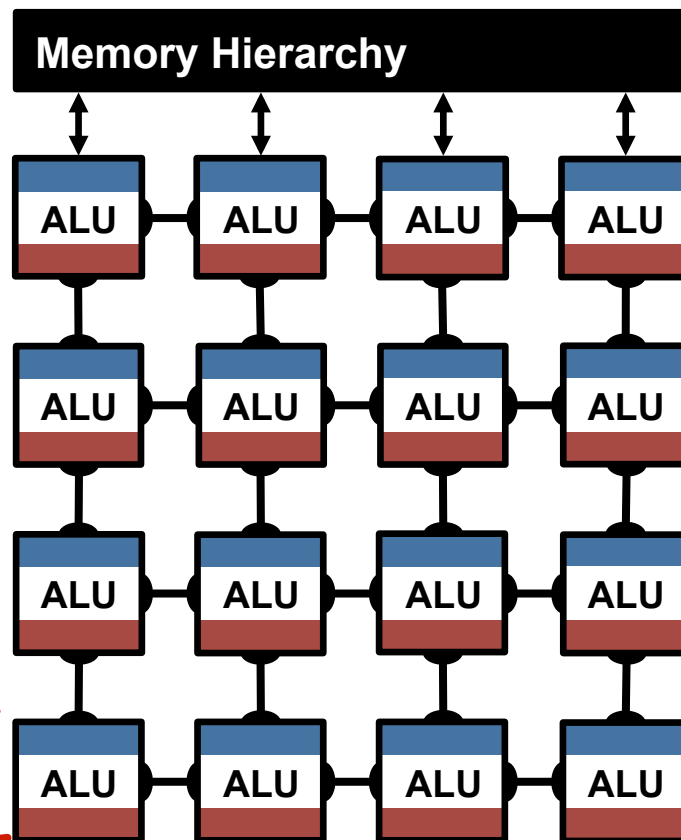
Sharing among regions of PEs

Processing  
Element (PE)

0.5 – 1.0 kB



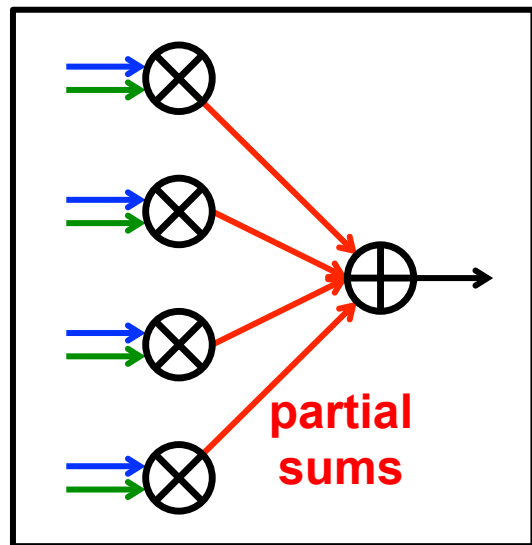
Spatial Architecture  
(Dataflow Processing)



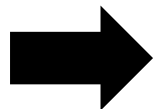


# How to Map the Dataflow?

## CNN Convolution

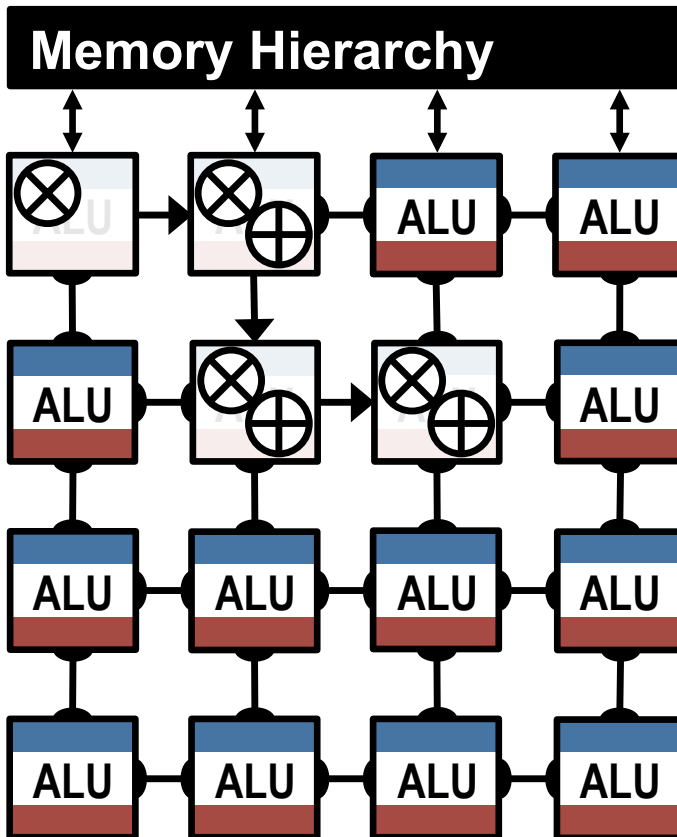


?



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

## Spatial Architecture (Dataflow Processing)

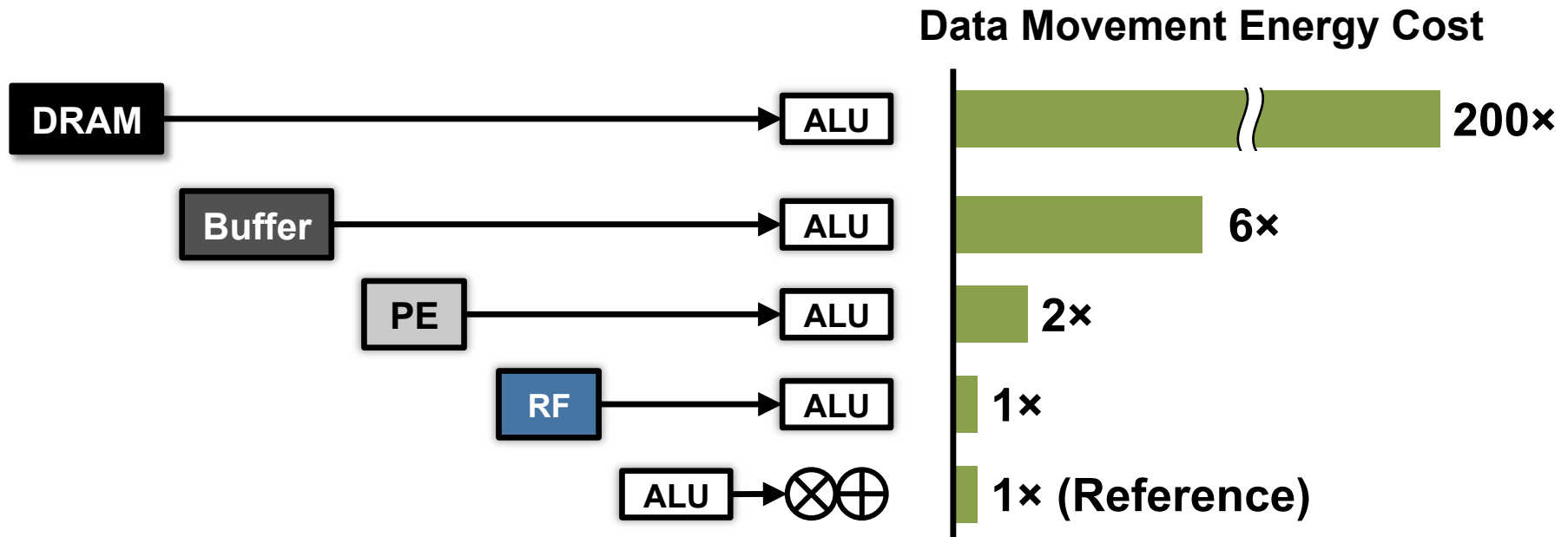
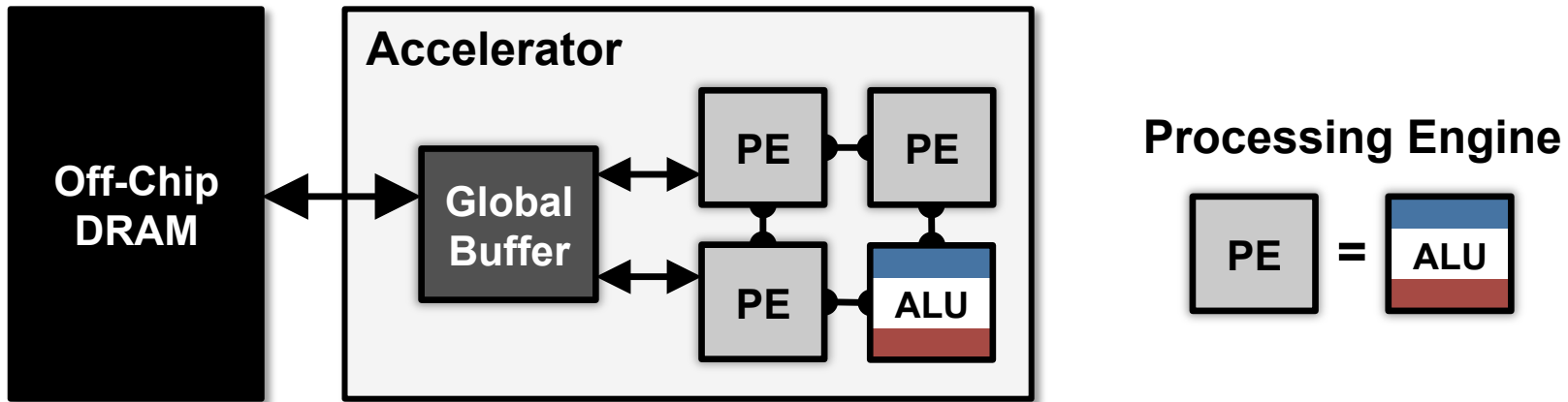


# Energy-Efficient Dataflow

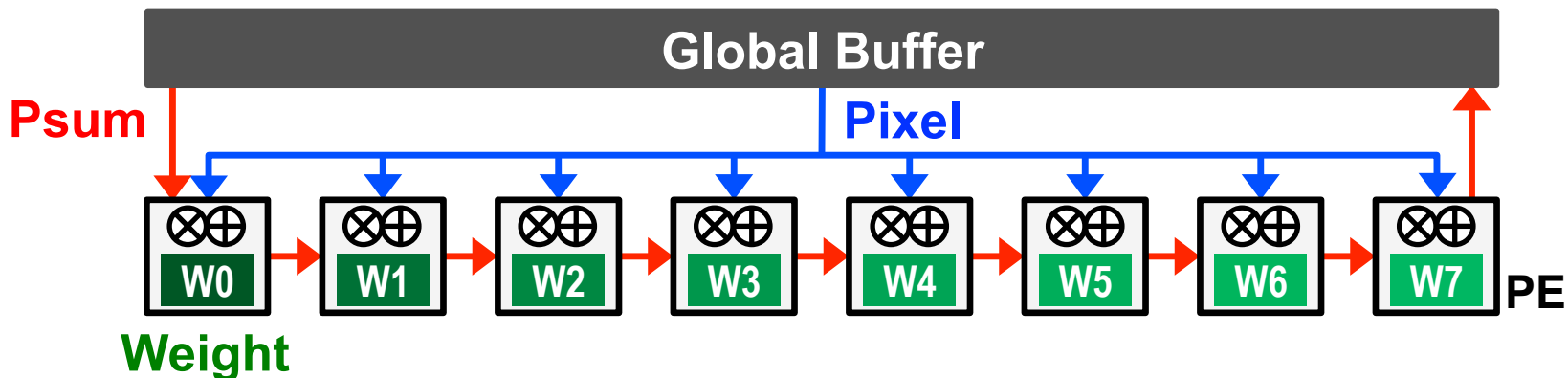
Yu-Hsin Chen, Joel Emer, Vivienne Sze, [ISCA 2016](#)

**Maximize data reuse and accumulation at RF**

# Data Movement is Expensive



Maximize data reuse at lower levels of hierarchy

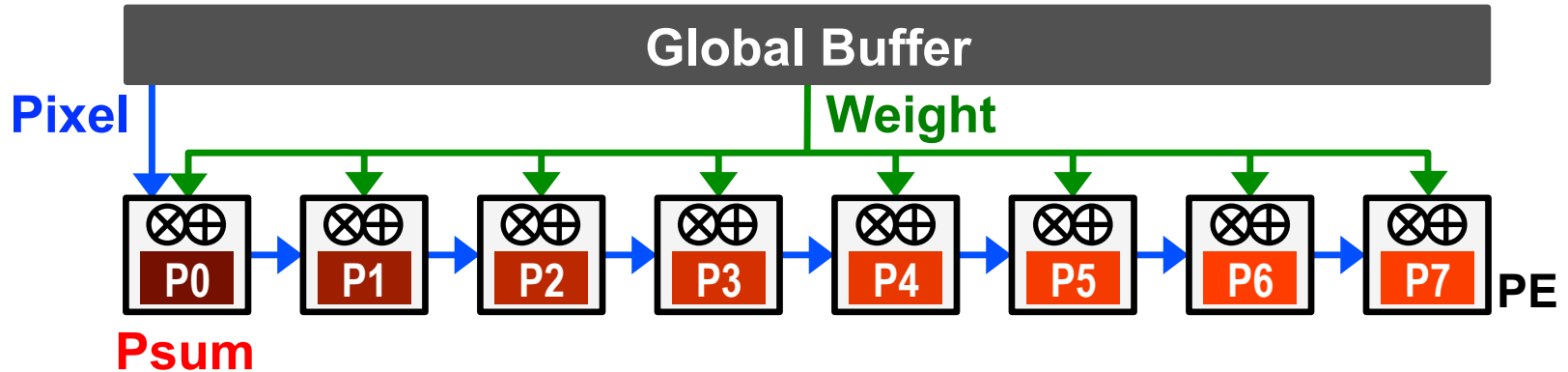


- Minimize **weight** read energy consumption
  - maximize convolutional and filter reuse of weights

• **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

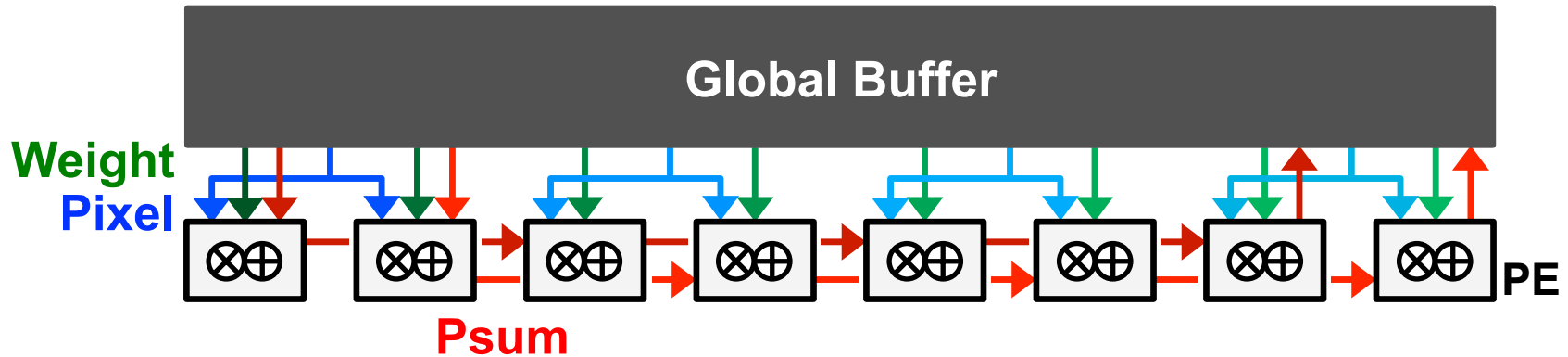
- **Examples:**

[Gupta, *ICML* 2015]

[ShiDianNao, *ISCA* 2015]

[Peemen, *ICCD* 2013]

# No Local Reuse (NLR)

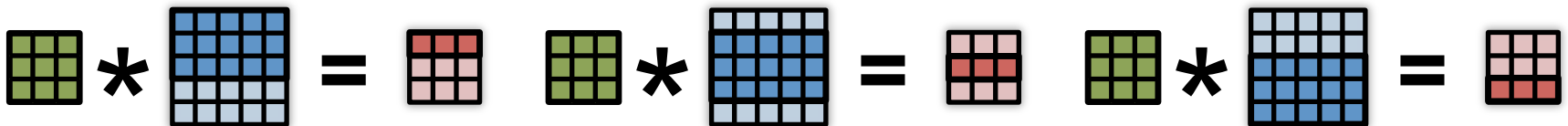
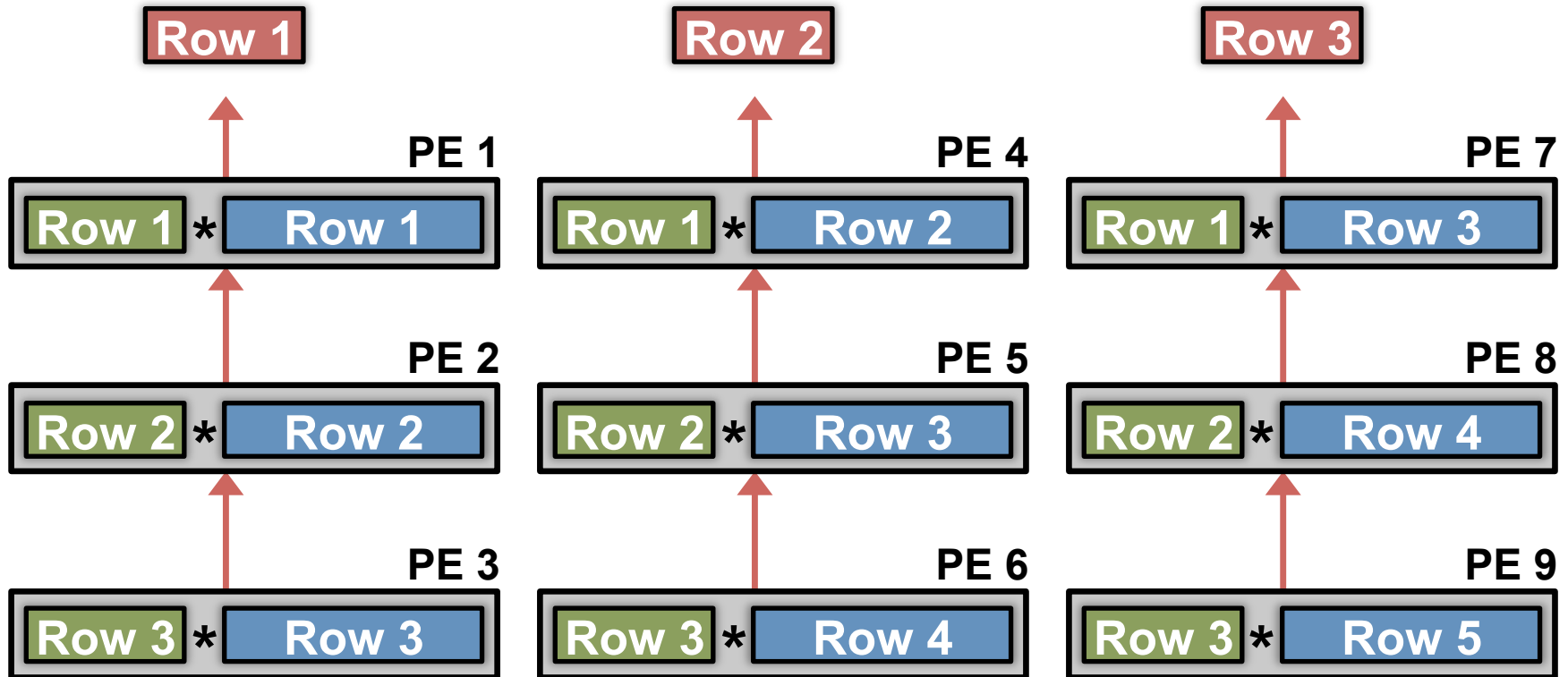


- Use a **large global buffer** as shared storage
  - Reduce **DRAM** access energy consumption
- **Examples:**

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

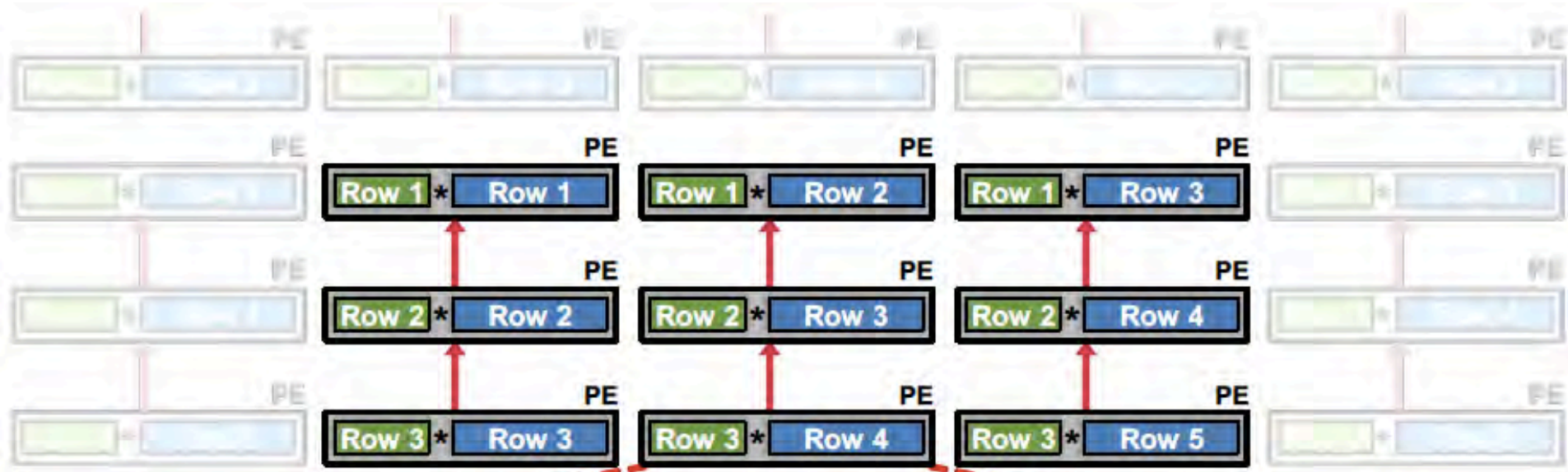


# Row Stationary Dataflow



Optimize for **overall energy efficiency** instead  
for only a certain data type

# CNN Convolution – The Full Picture



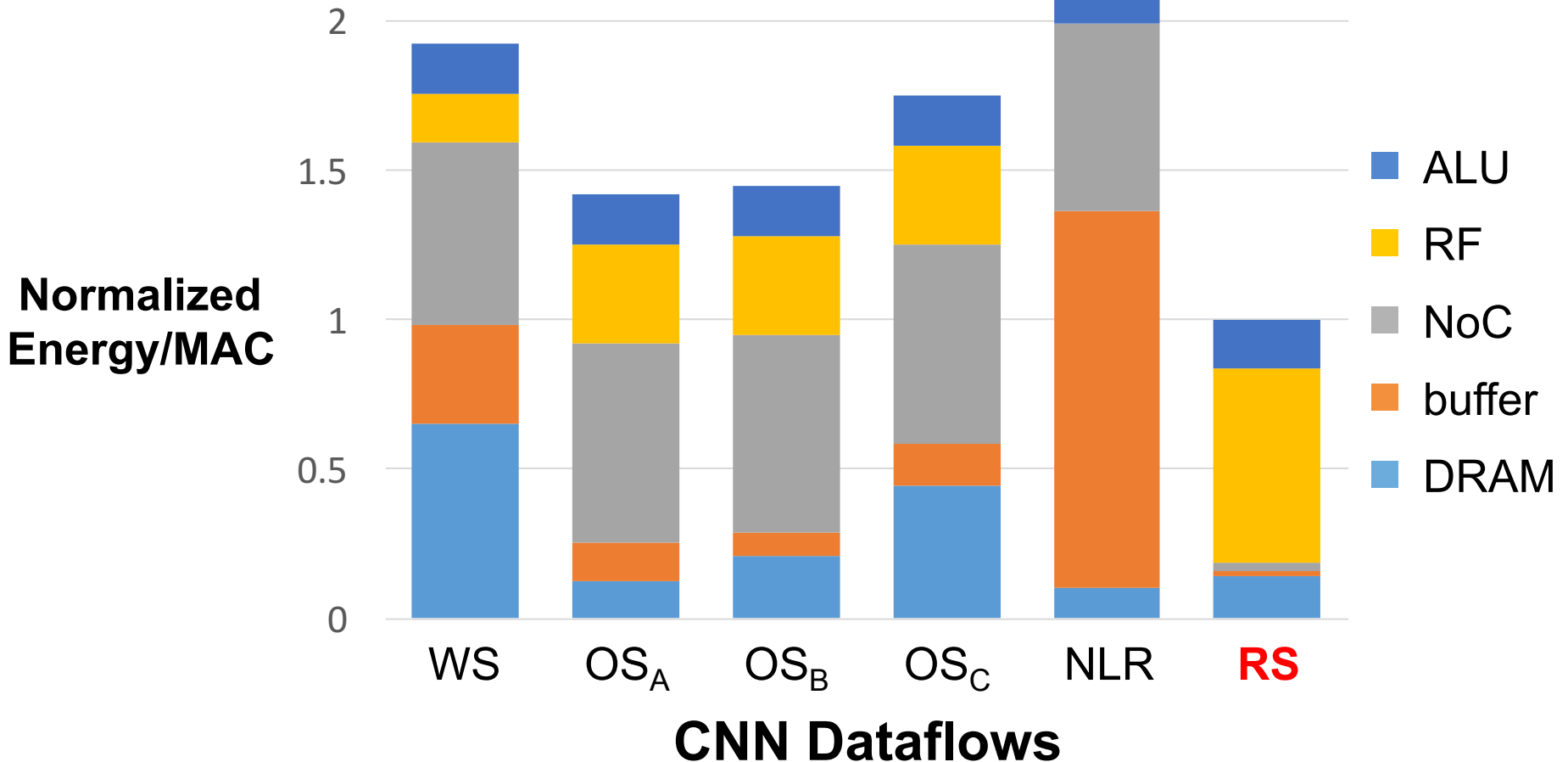
Multiple **images**:  $\text{Filter 1} * \text{Image 1 \& 2} = \text{Psum 1 \& 2}$

Multiple **filters**:  $\text{Filter 1 \& 2} * \text{Image 1} = \text{Psum 1 \& 2}$

Multiple **channels**:  $\text{Filter 1} * \text{Image 1} = \text{Psum}$

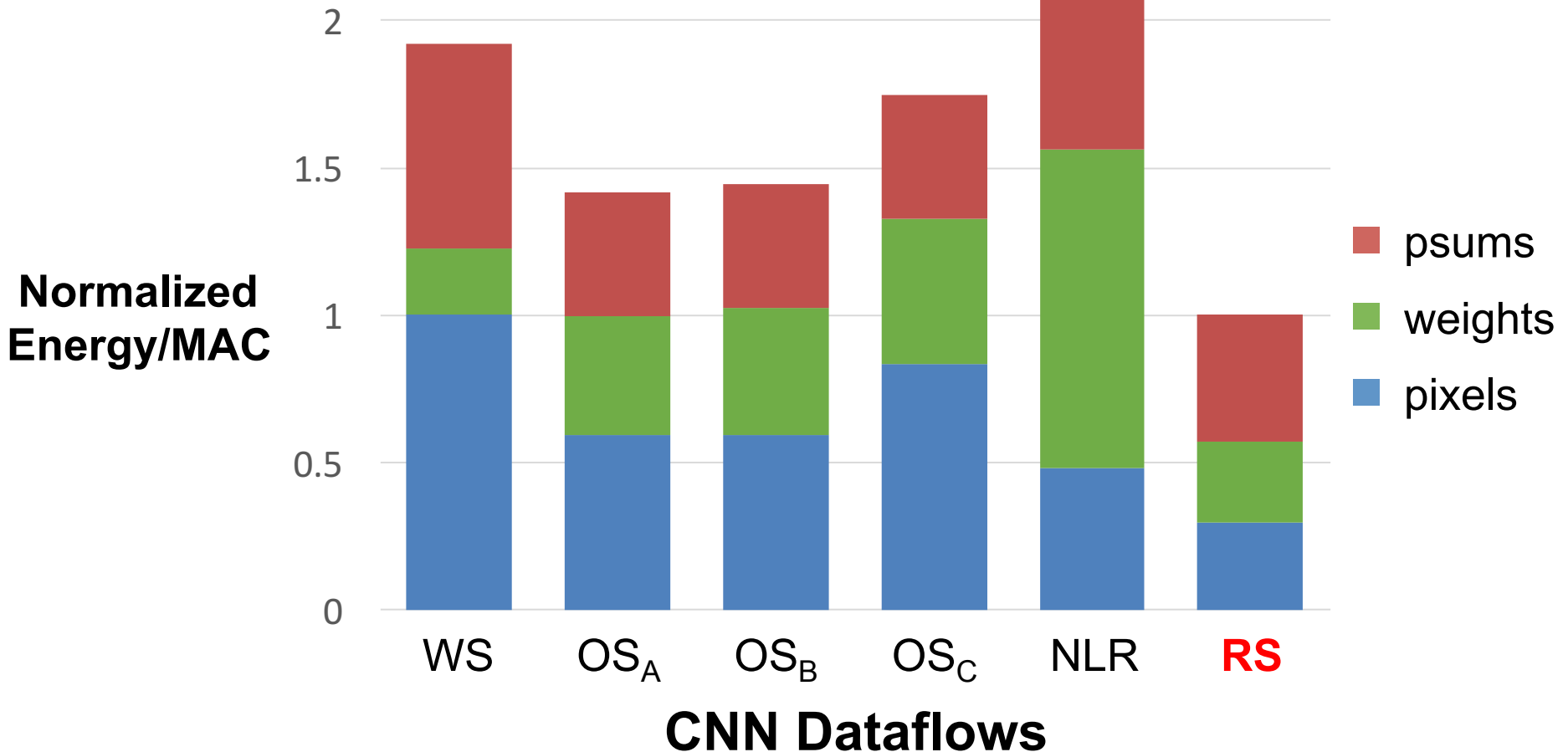
Map rows from **multiple images**, **filters** and **channels** to same PE to exploit other forms of reuse and local accumulation

# Dataflow Comparison: CONV Layers



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

# Dataflow Comparison: CONV Layers



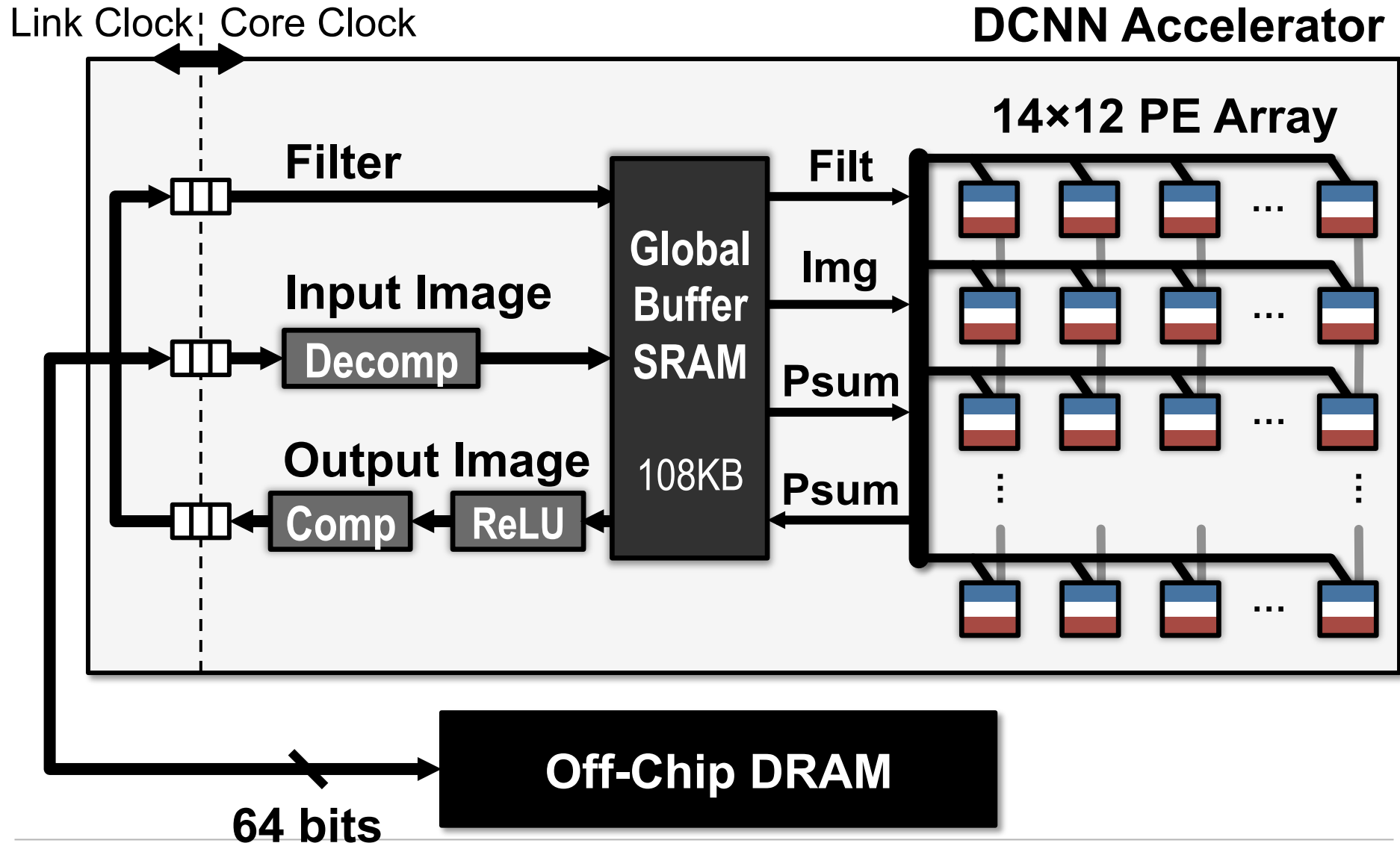
RS optimizes for the best **overall** energy efficiency

# Energy-Efficient Accelerator

Yu-Hsin Chen, Tushar Krishna, Joel Emer, Vivienne Sze, [ISSCC 2016](#)

## Exploit data statistics

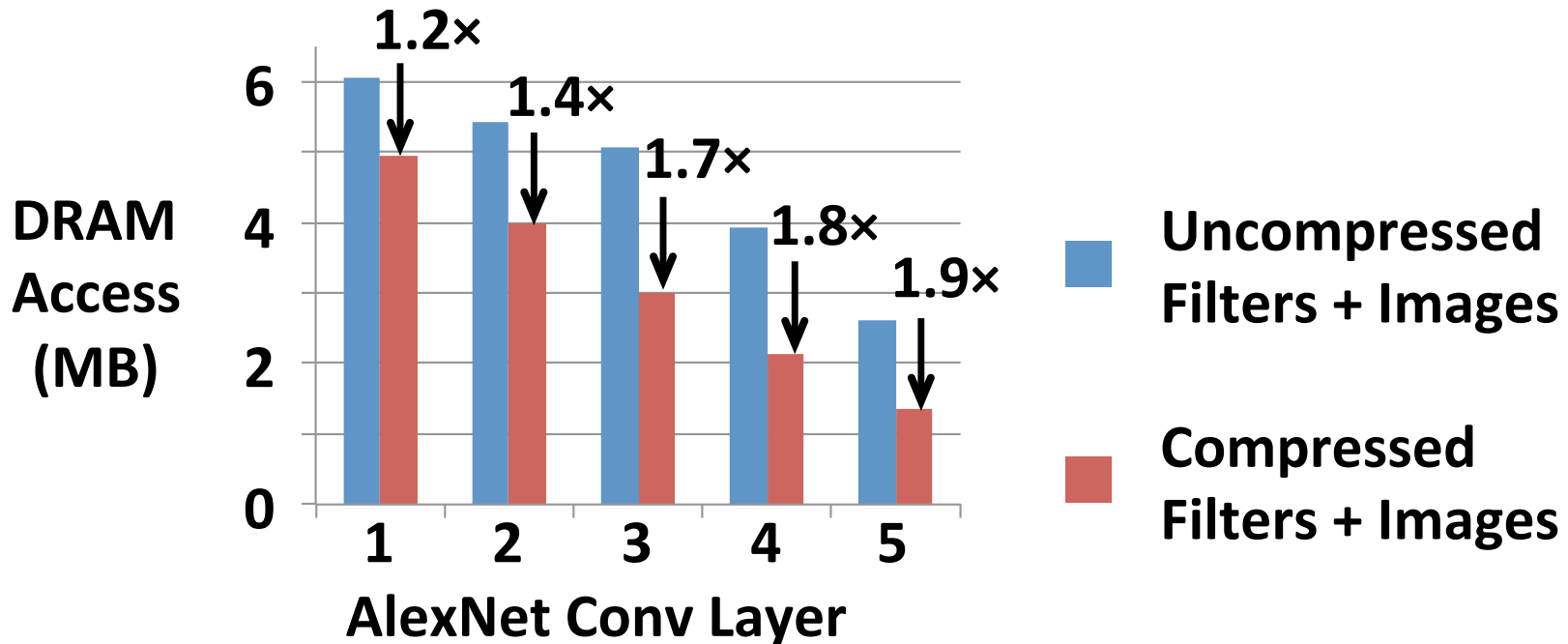
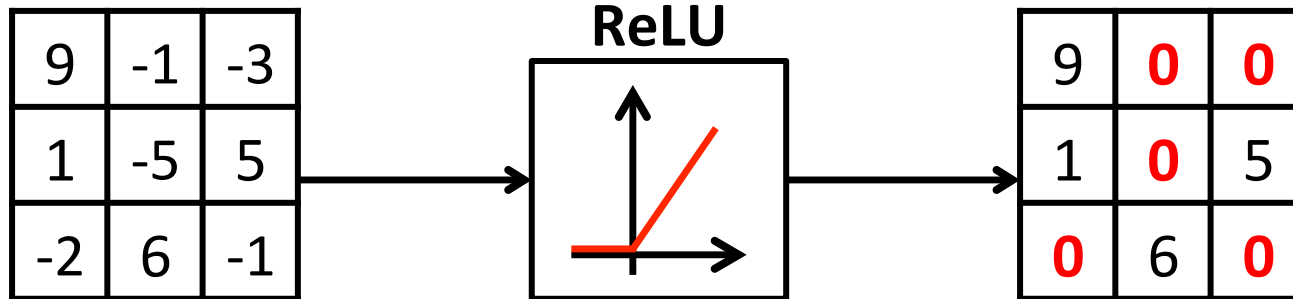
# Eyeriss Deep CNN Accelerator





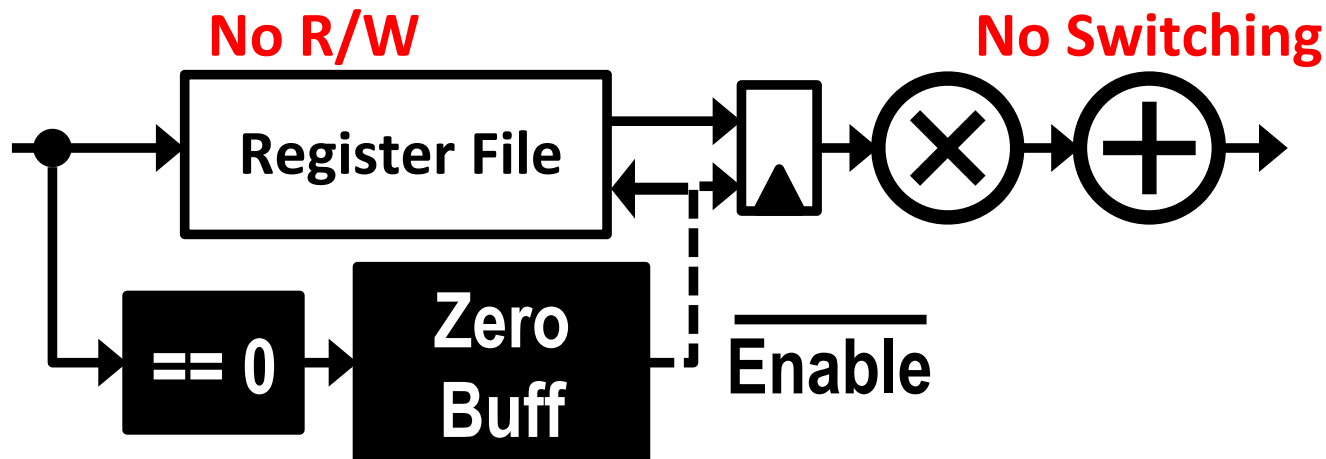
# Data Compression Saves DRAM BW

Apply Non-Linearity (**ReLU**) on Filtered Image Data



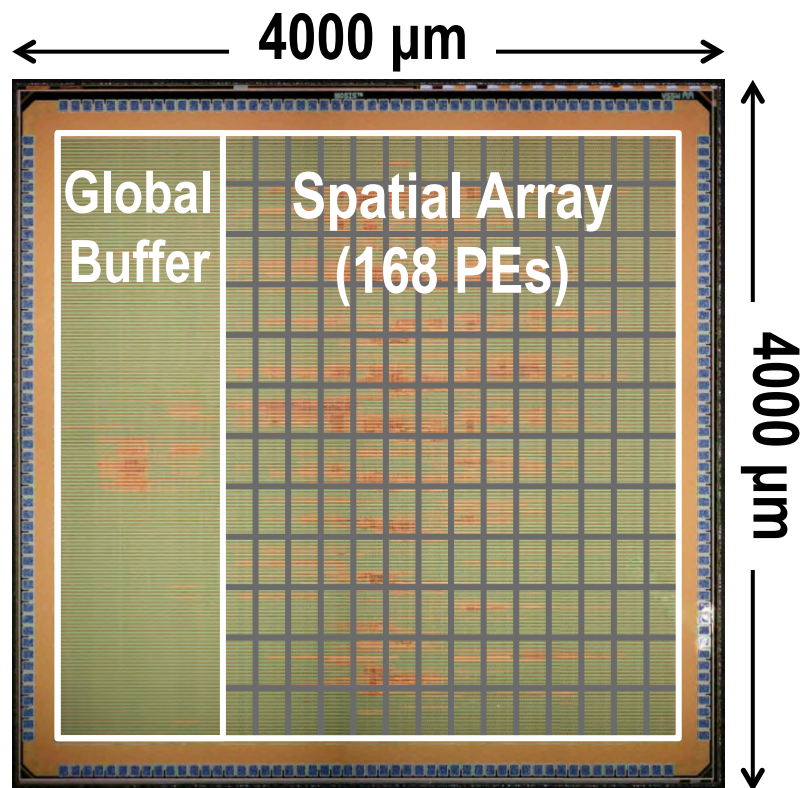
# Zero Data Processing Gating

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



# Eyeriss Chip Spec & Measurement Results

<b>Technology</b>	TSMC 65nm LP 1P9M
<b>On-Chip Buffer</b>	108 KB
<b># of PEs</b>	168
<b>Scratch Pad / PE</b>	0.5 KB
<b>Core Frequency</b>	100 – 250 MHz
<b>Peak Performance</b>	33.6 – 84.0 GOPS
<b>Word Bit-width</b>	16-bit Fixed-Point
<b>Natively Supported CNN Shapes</b>	Filter Width: 1 – 32 Filter Height: 1 – 12 Num. Filters: 1 – 1024 Num. Channels: 1 – 1024 Horz. Stride: 1–12 Vert. Stride: 1, 2, 4



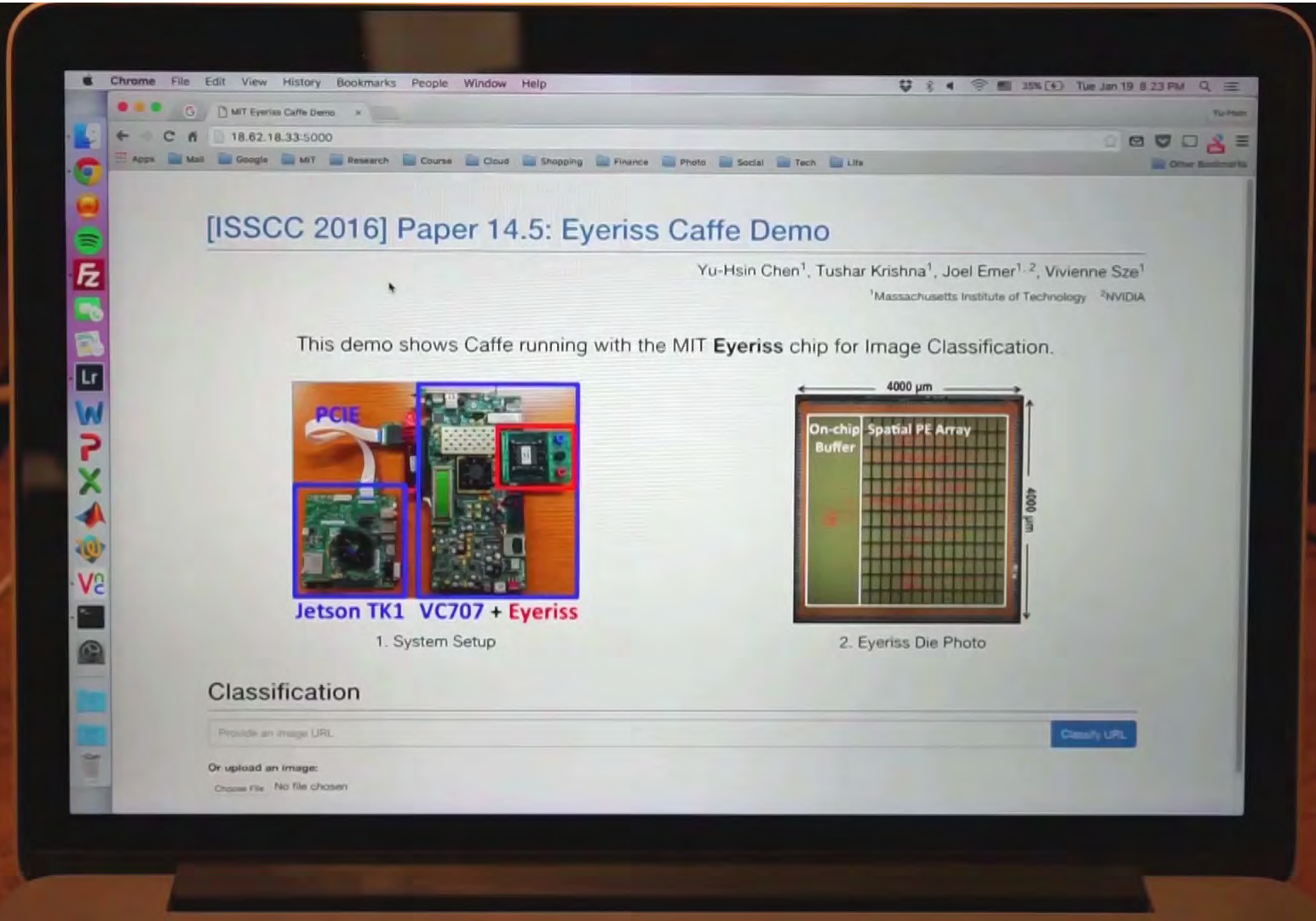
AlexNet: For 2.66 GMACs [8 billion 16-bit inputs (**16GB**) and 2.7 billion outputs (**5.4GB**)], only requires **208.5MB** (buffer) and **15.4MB** (DRAM)

# Comparison with GPU

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

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

More info about **Eyeriss** and  
**Tutorial on DNN Architectures** at

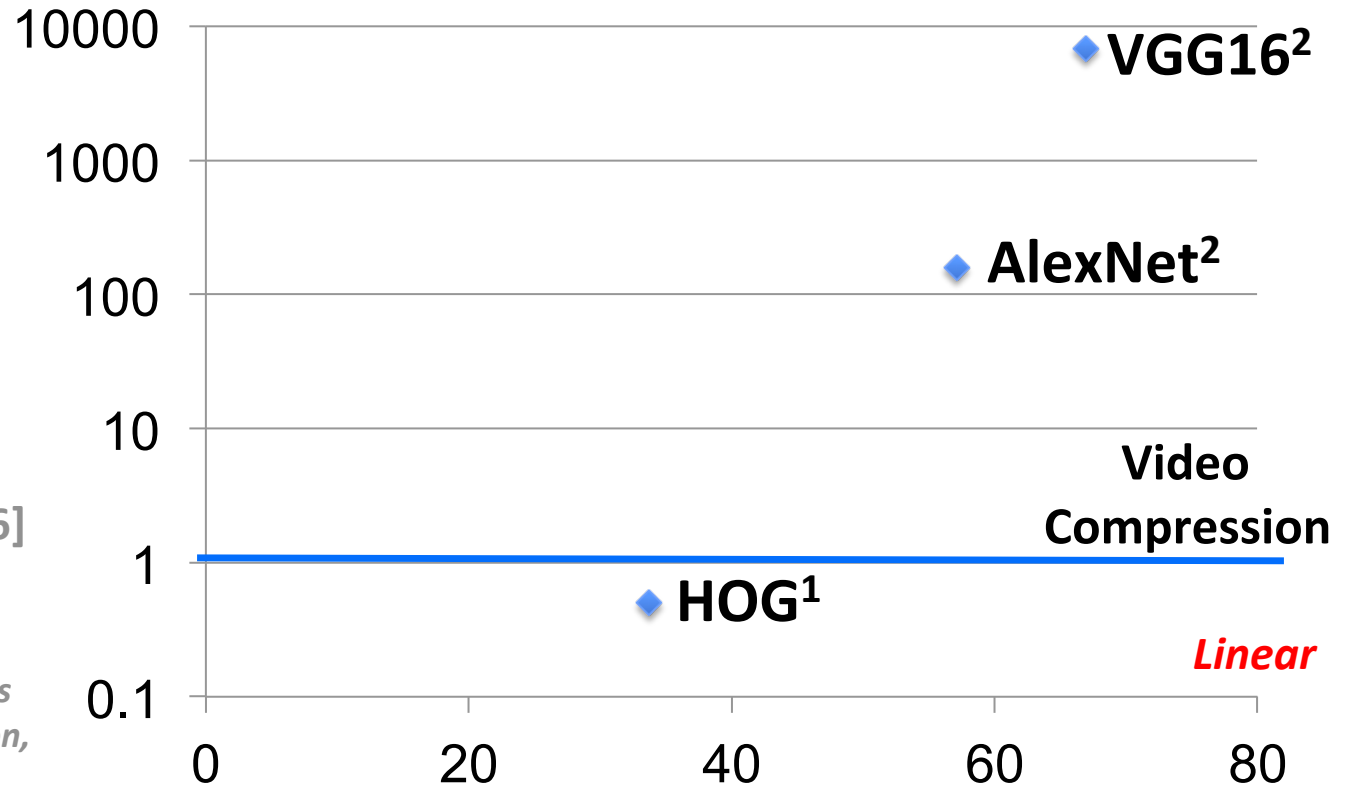
<http://eyeriss.mit.edu>



# Features: Energy vs. Accuracy

*Exponential*

Energy/  
Pixel (nJ)



*Measured in 65nm\**

- [Suleiman, VLSI 2016]
- [Chen, ISSCC 2016]

*\* Only feature extraction. Does not include data, augmentation, ensemble and classification energy, etc.*

**Accuracy (Average Precision)**

*Measured in on VOC 2007 Dataset*

- DPM v5 [Girshick, 2012]
- Fast R-CNN [Girshick, CVPR 2015]

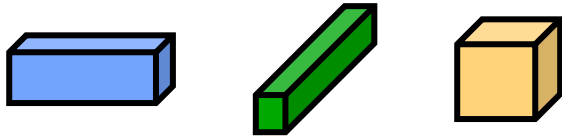


# Designing Energy-Efficient CNNs using Energy-Aware Pruning

Tien-Ju Yang, Yu-Hsin Chen, Vivienne Sze, [CVPR 2017](#)

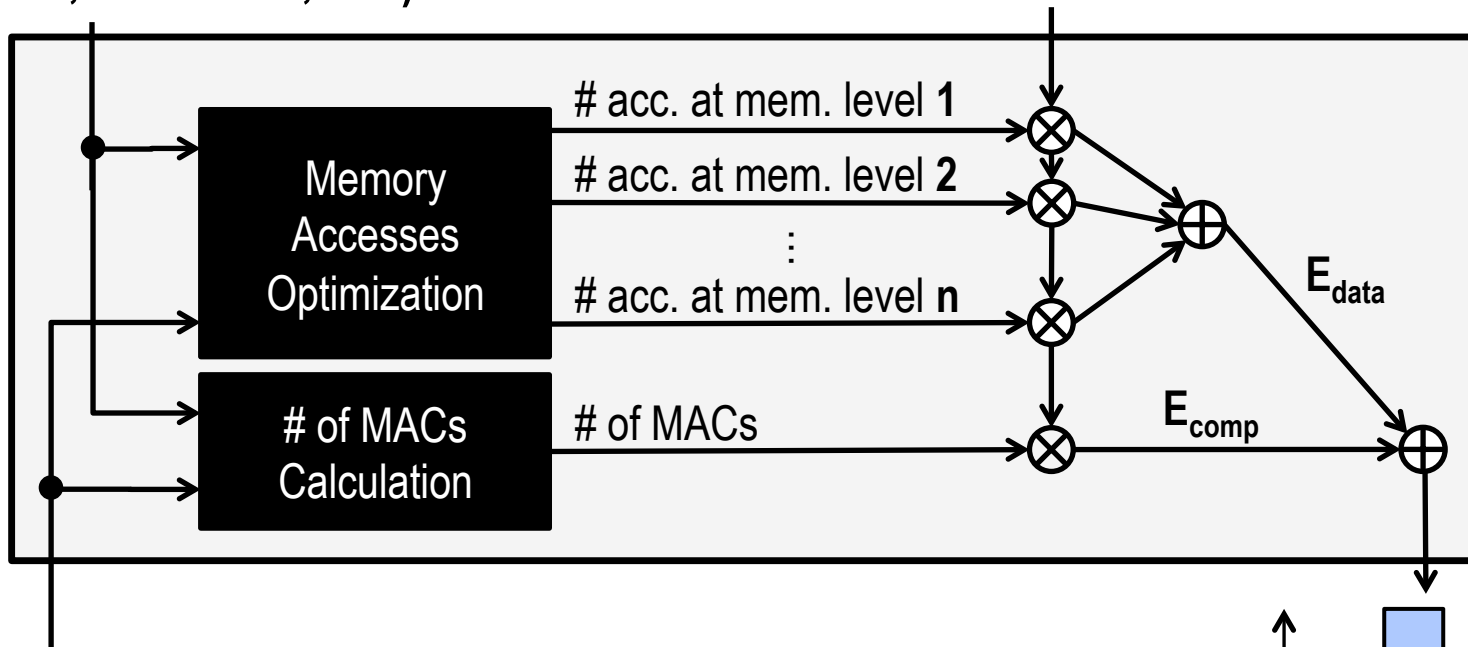


# Energy-Evaluation Methodology



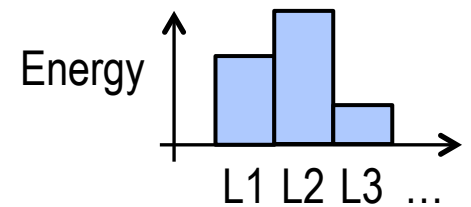
**CNN Shape Configuration**  
(# of channels, # of filters, etc.)

**Hardware Energy Costs of each  
MAC and Memory Access**



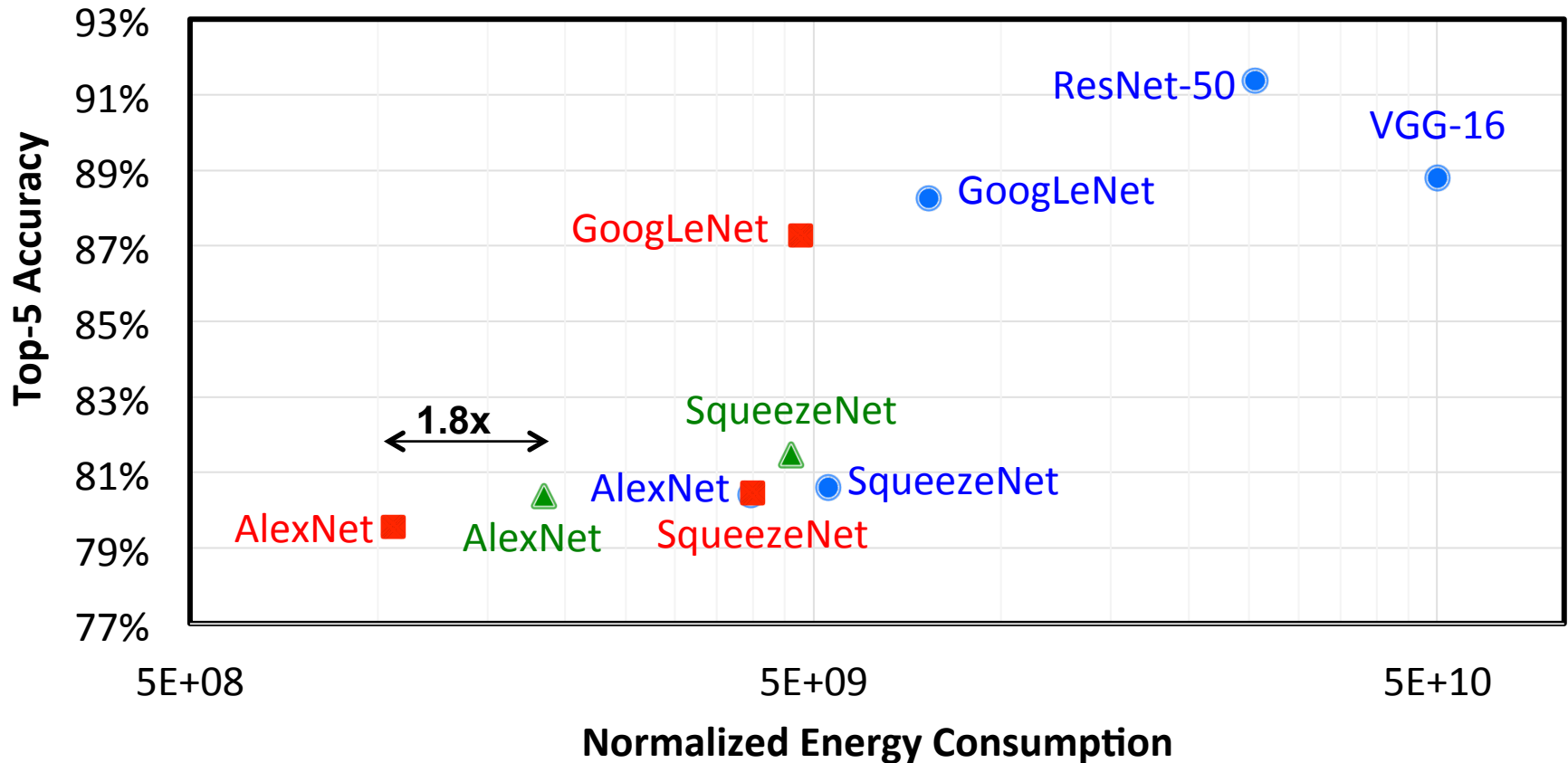
**CNN Weights and Input Data**

[0.3, 0, -0.4, 0.7, 0, 0, 0.1, ...]



**CNN Energy Consumption**

# Energy-Aware Pruning



● Original DNN    ▲ Magnitude-based Pruning    ■ Energy-aware Pruning (This Work)

Remove weights from layers in order of highest to lowest energy  
**3.7x reduction in AlexNet / 1.6x reduction in GoogLeNet**

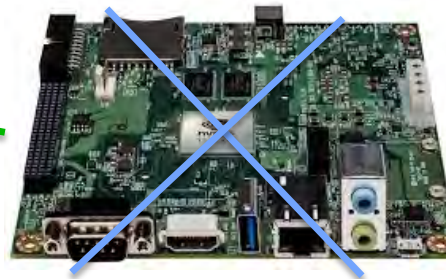
# Enable real-time navigation on nanoDrone



Image source: Cheerson



Big battery



Mobile GPU

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



# Acknowledgements



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





# References

More info about **Eyeriss** and  
**Tutorial on DNN Architectures** at  
<http://eyeriss.mit.edu>

More info about research in the **Energy-Efficient  
Multimedia Systems Group @ MIT**  
<http://www.rle.mit.edu/eems>

For updates



Follow @eems\_mit

<http://mailman.mit.edu/mailman/listinfo/eems-news>