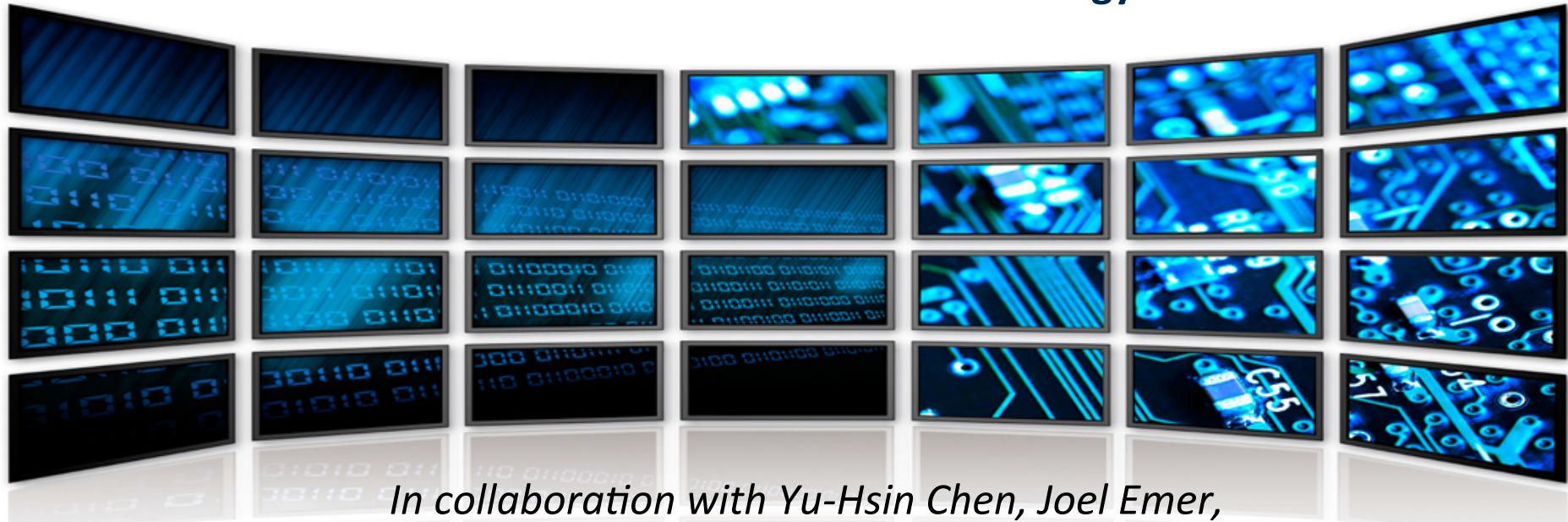


Efficient Processing of Deep Neural Networks

Vivienne Sze

Massachusetts Institute of Technology



*In collaboration with Yu-Hsin Chen, Joel Emer,
Tushar Krishna, Amr Suleiman, Tien-Ju Yang,
Zhengdong Zhang*

Contact Info

email: sze@mit.edu

website: www.rle.mit.edu/eems



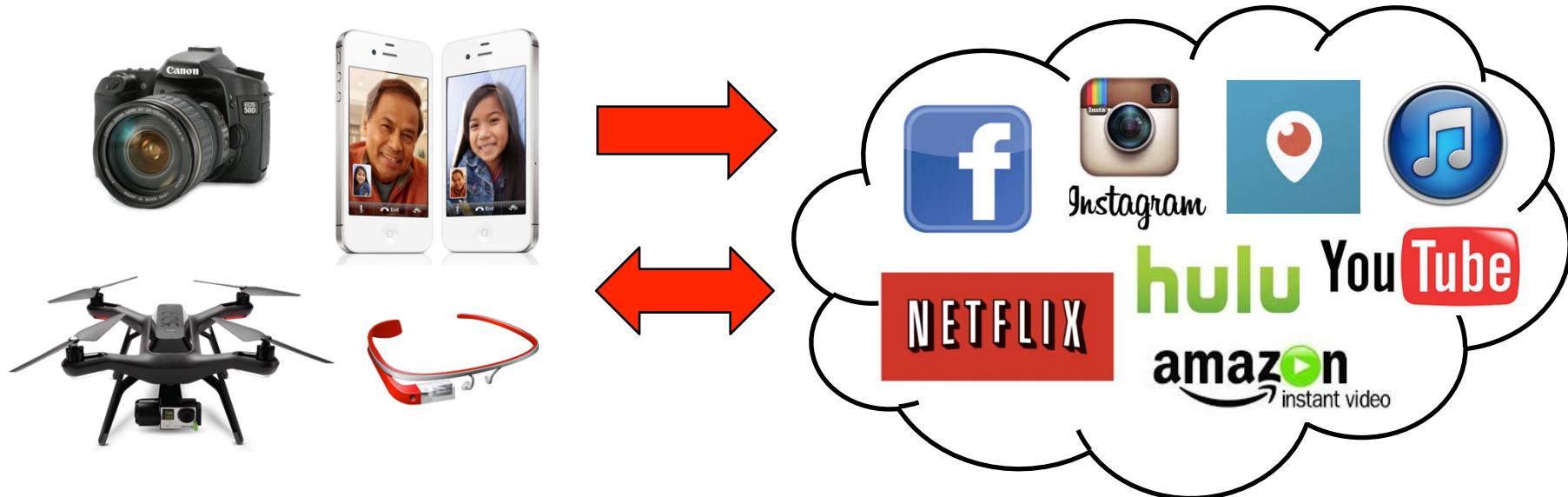
Follow @eems_mit

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**



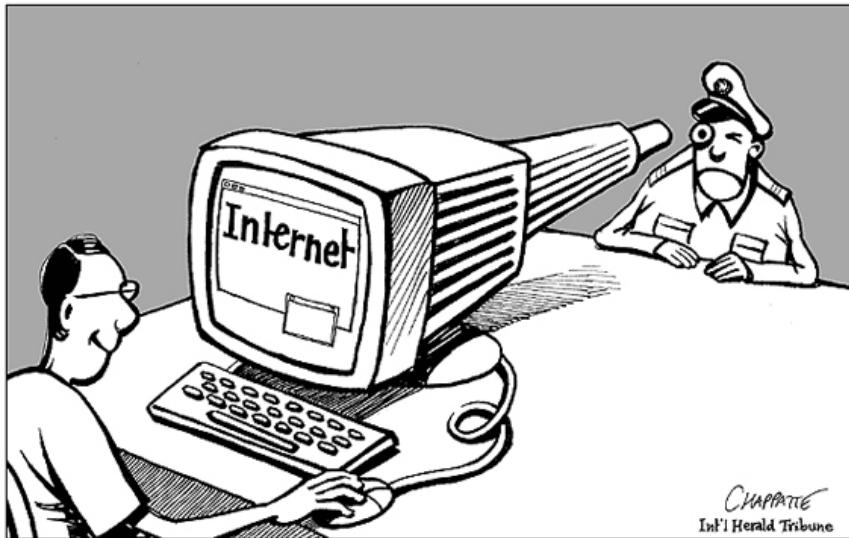
*Energy limited due
to battery capacity*

*Power limited due
to heat dissipation*

Need energy-efficient pixel processing!

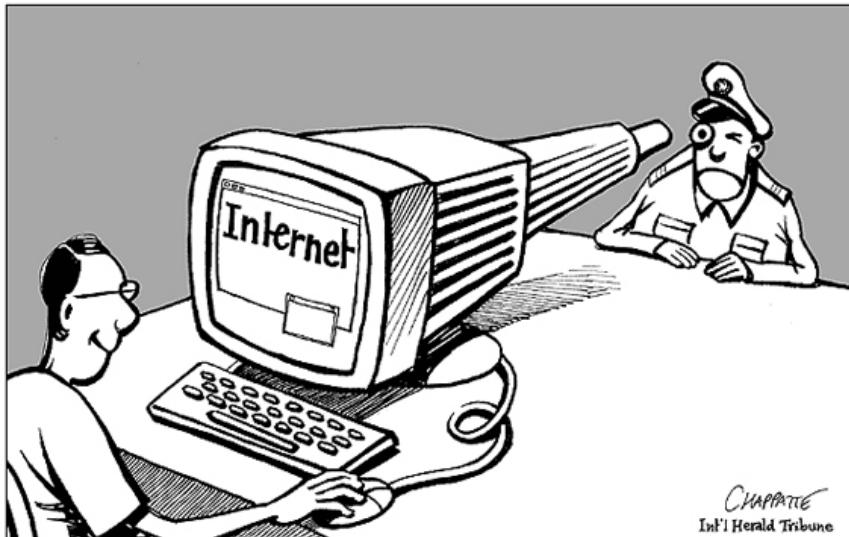
Processing at “Edge” instead of the “Cloud”

Privacy



Processing at “Edge” instead of the “Cloud”

Privacy



Latency

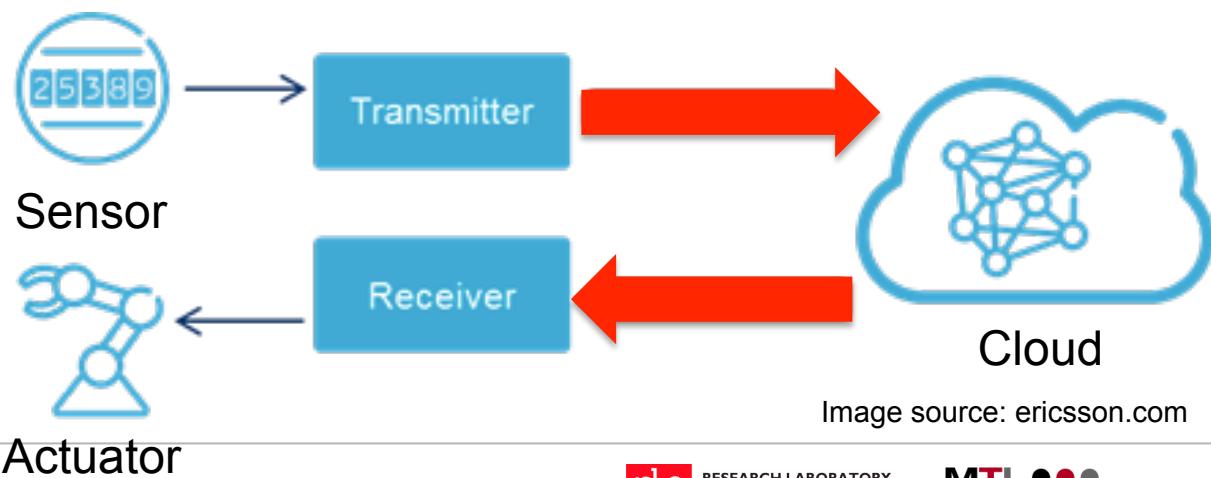
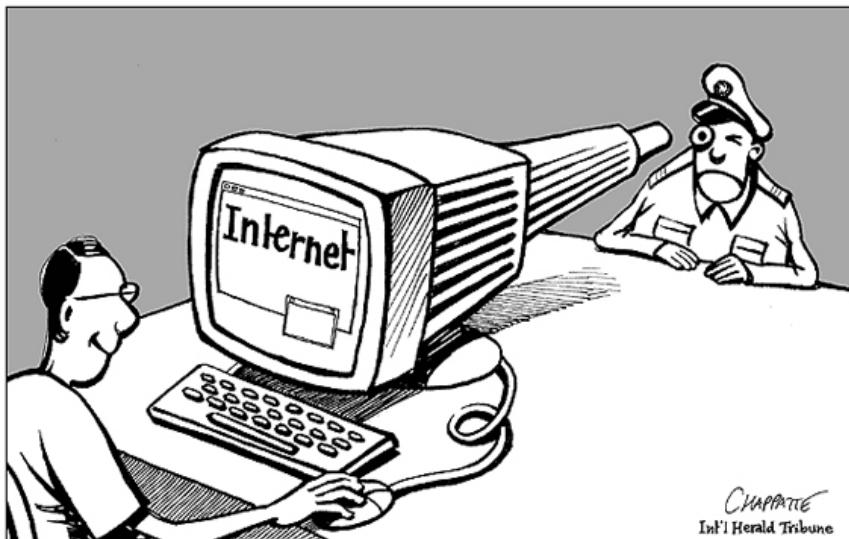


Image source: ericsson.com

Processing at “Edge” instead of the “Cloud”

Privacy



Communication



Image source:
www.theregister.co.uk

Latency

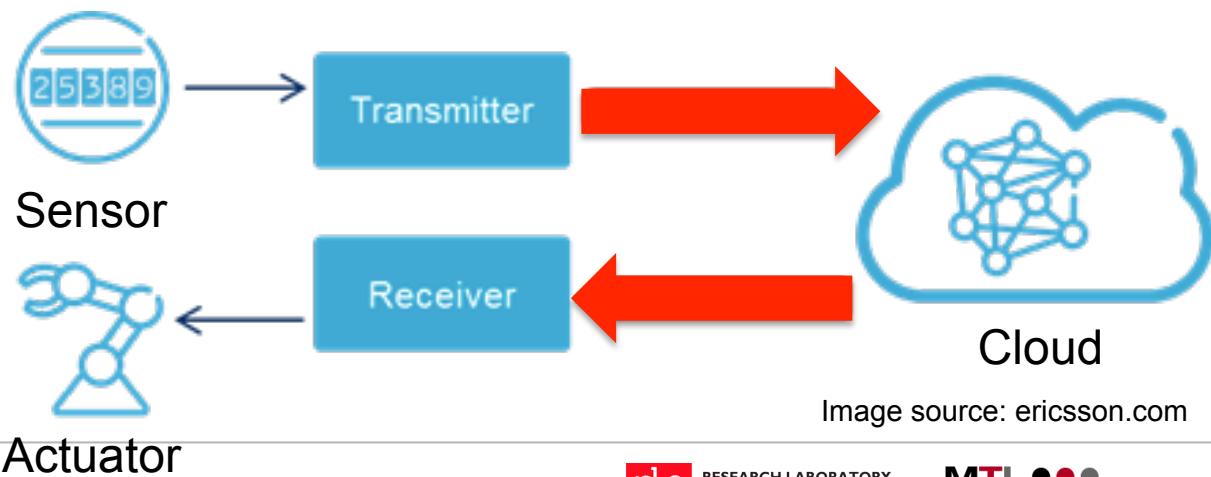
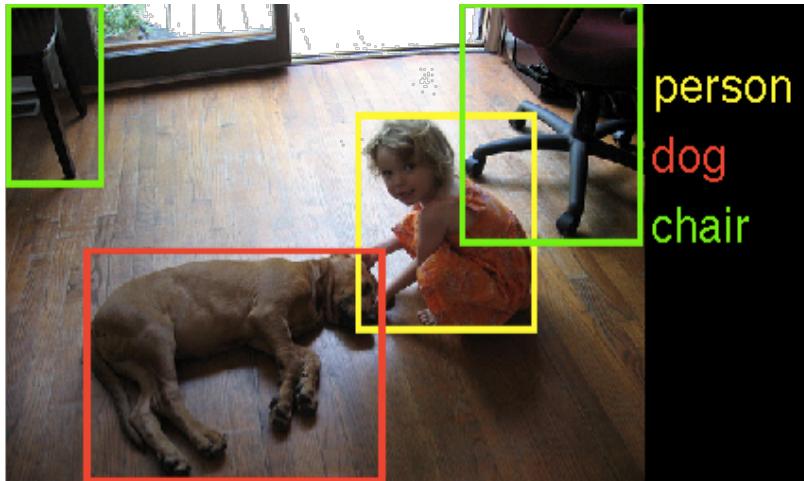


Image source: ericsson.com

Example Applications of Machine Learning

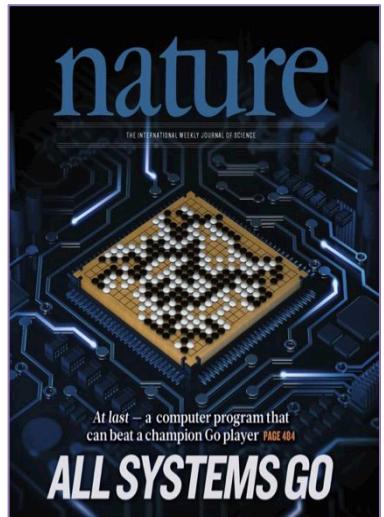
Computer Vision



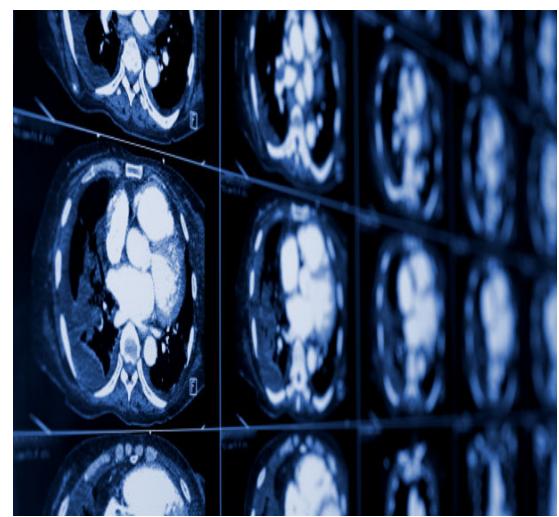
Speech Recognition



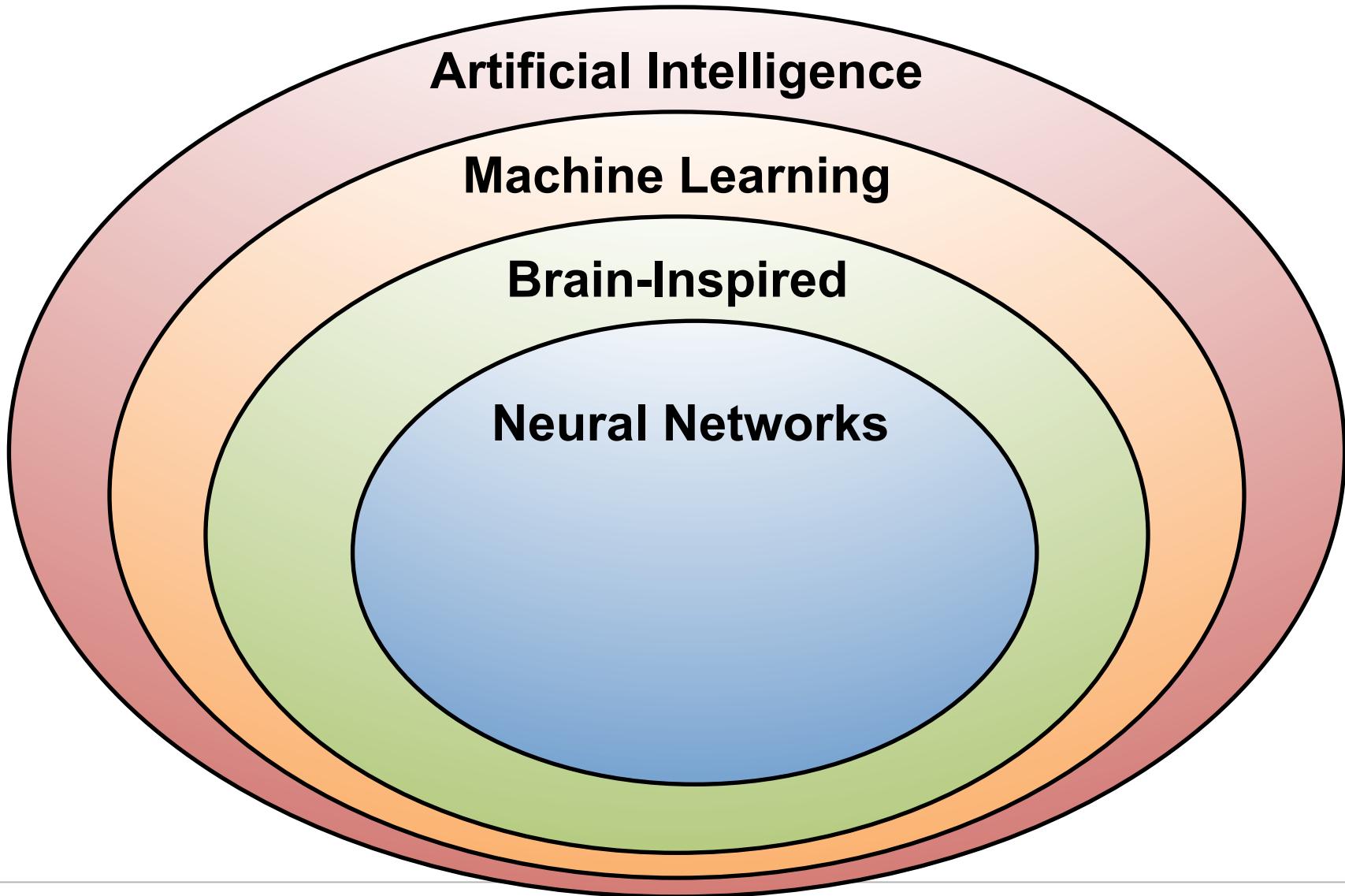
Game Play



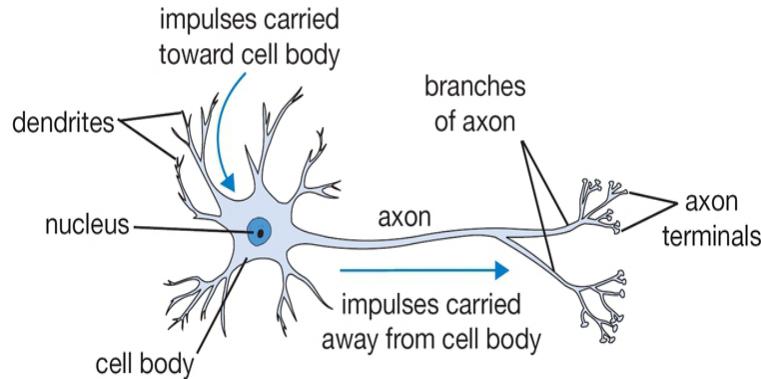
Medical



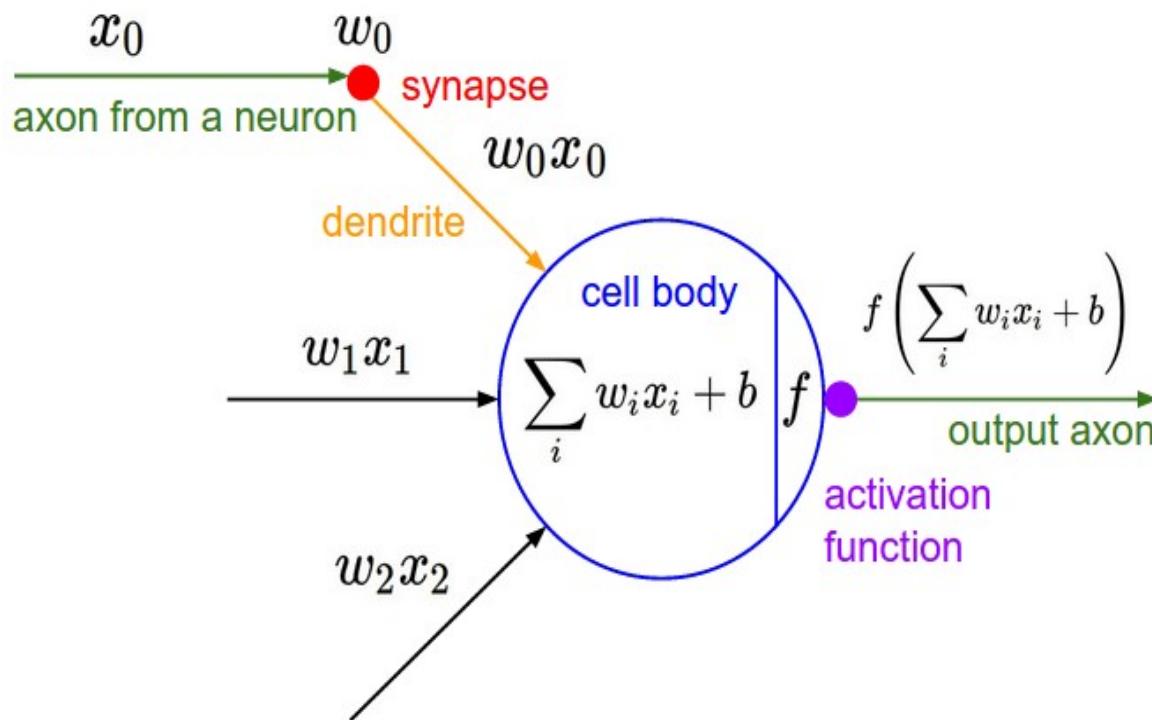
Machine Learning and Neural Networks



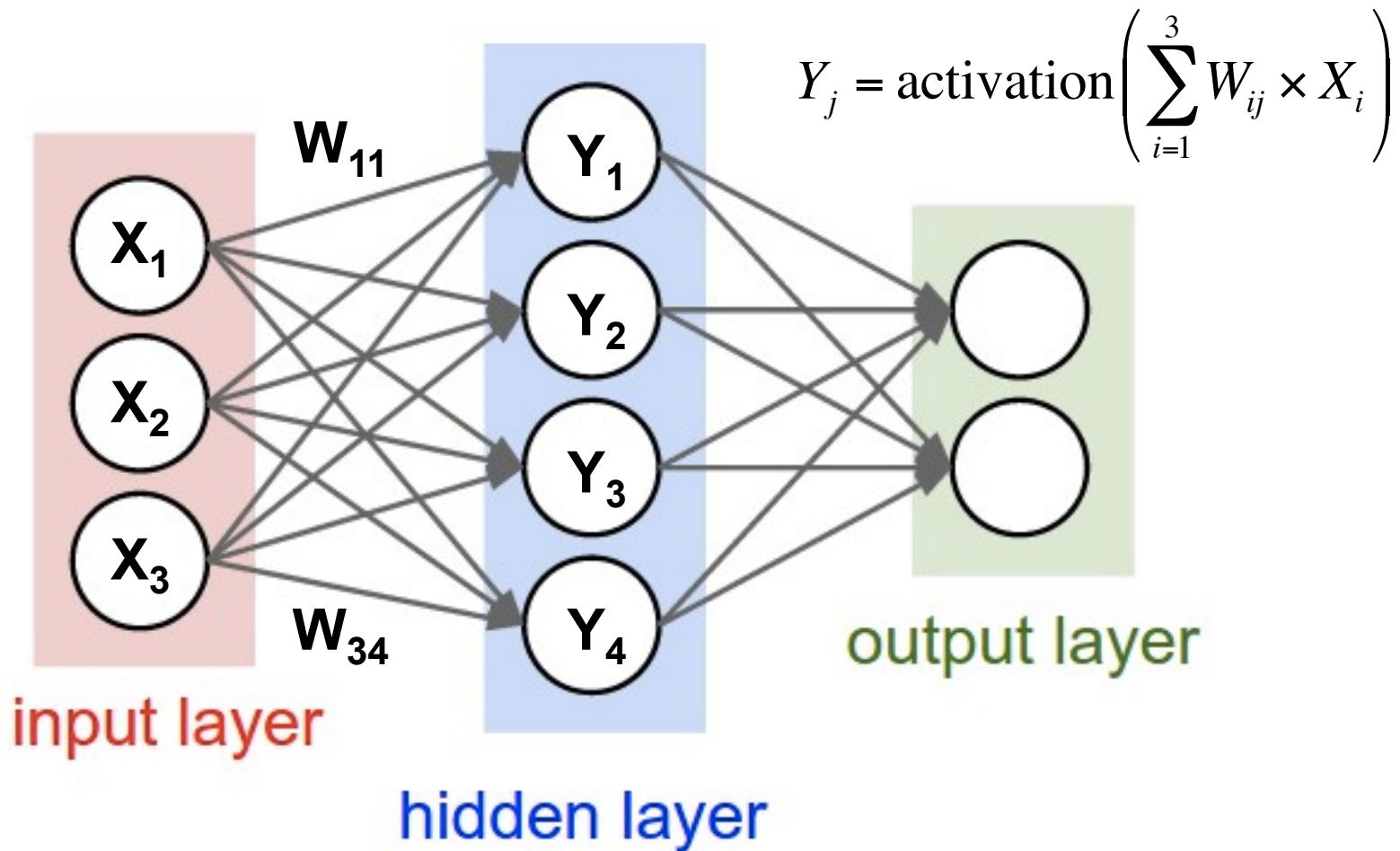
Neural Networks: Weighted Sum



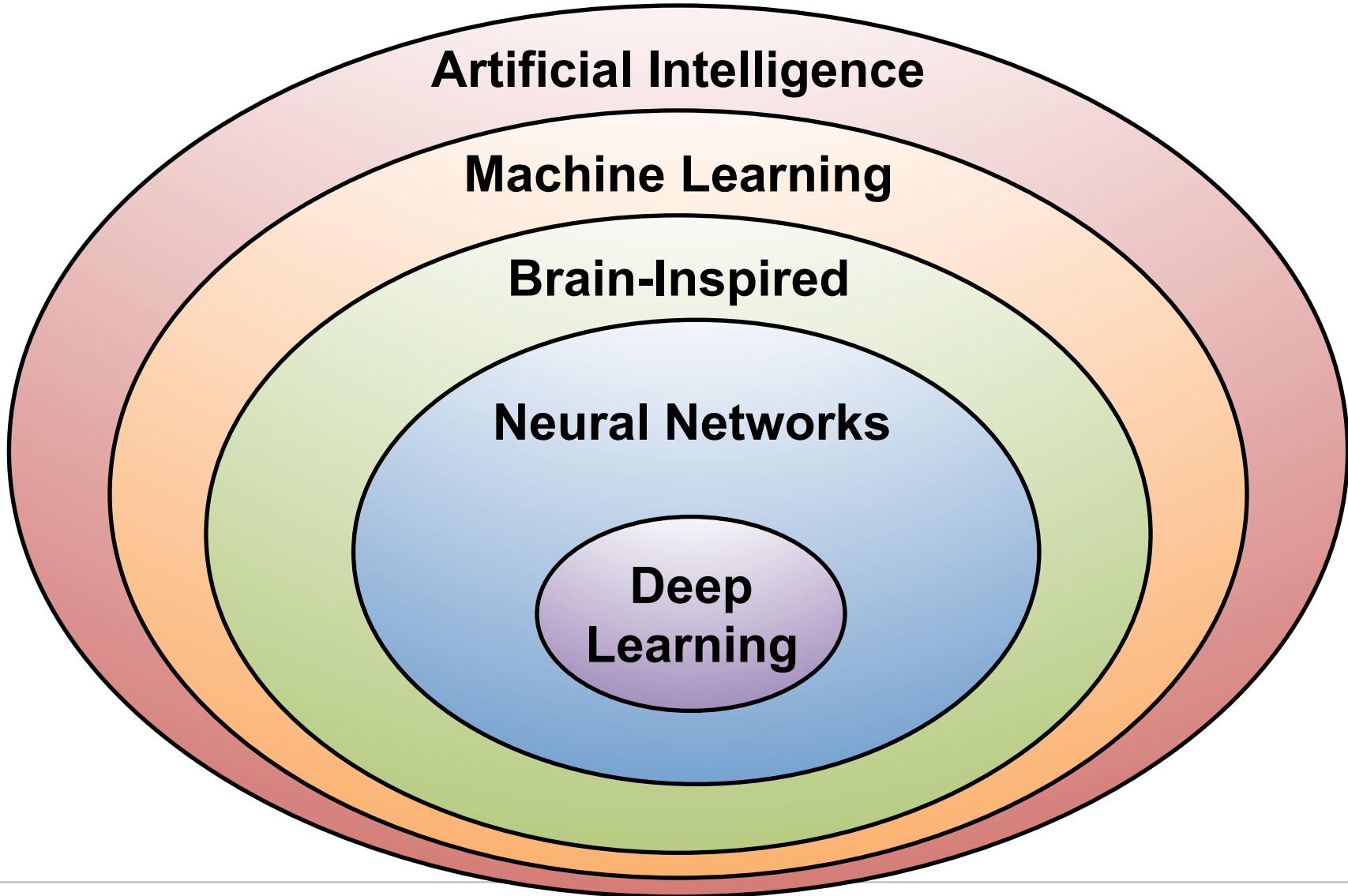
The brain contains
~ 10^{11} neurons connected with
~ $10^{14} – 10^{15}$ synapses



Many Weighted Sums



Deep Learning



What is Deep Learning?

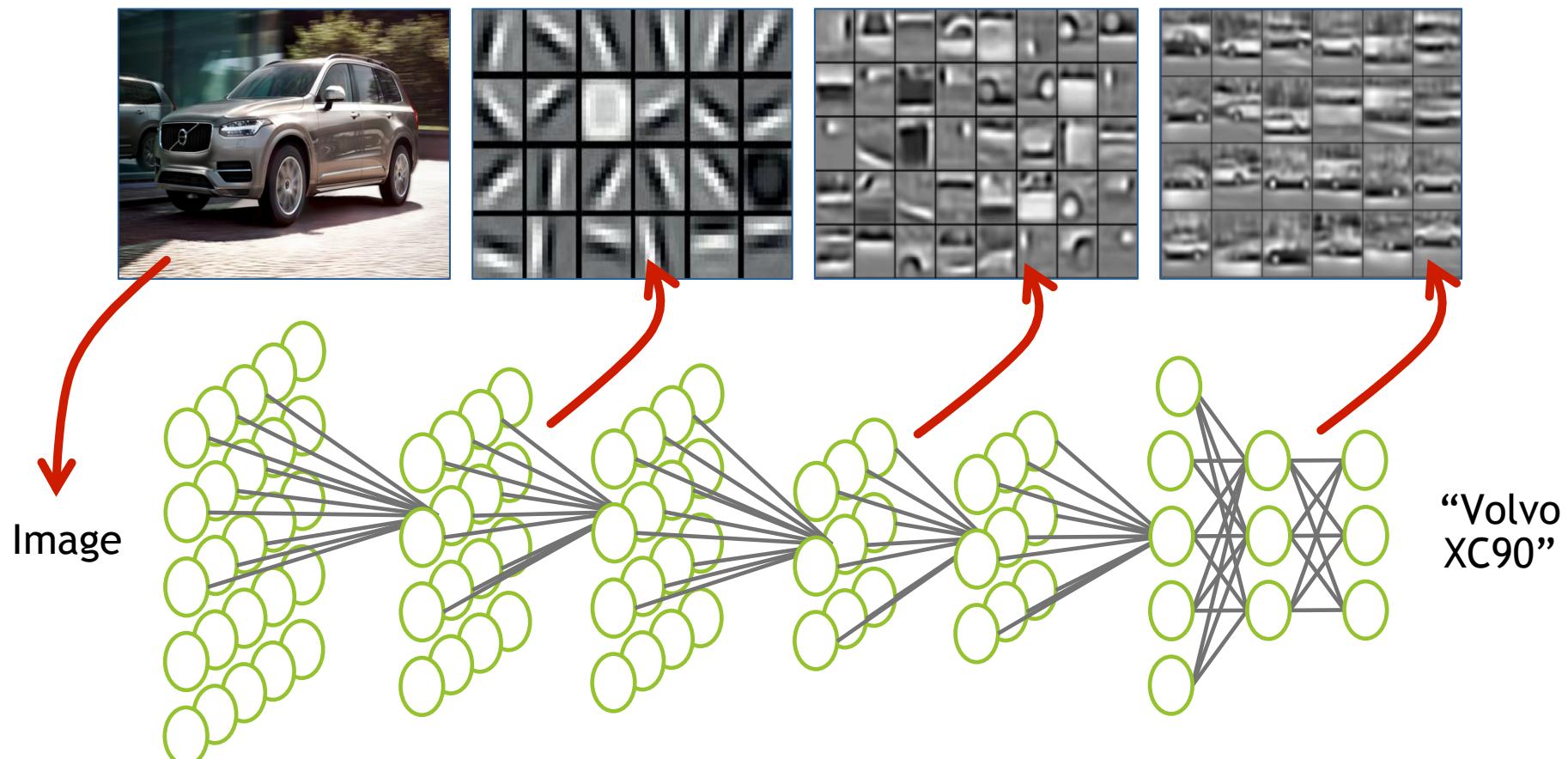


Image Source: [Lee et al., Comm. ACM 2011]

Why is Deep Learning Hot Now?

Big Data Availability

GPU Acceleration

New ML Techniques



350M images uploaded per day



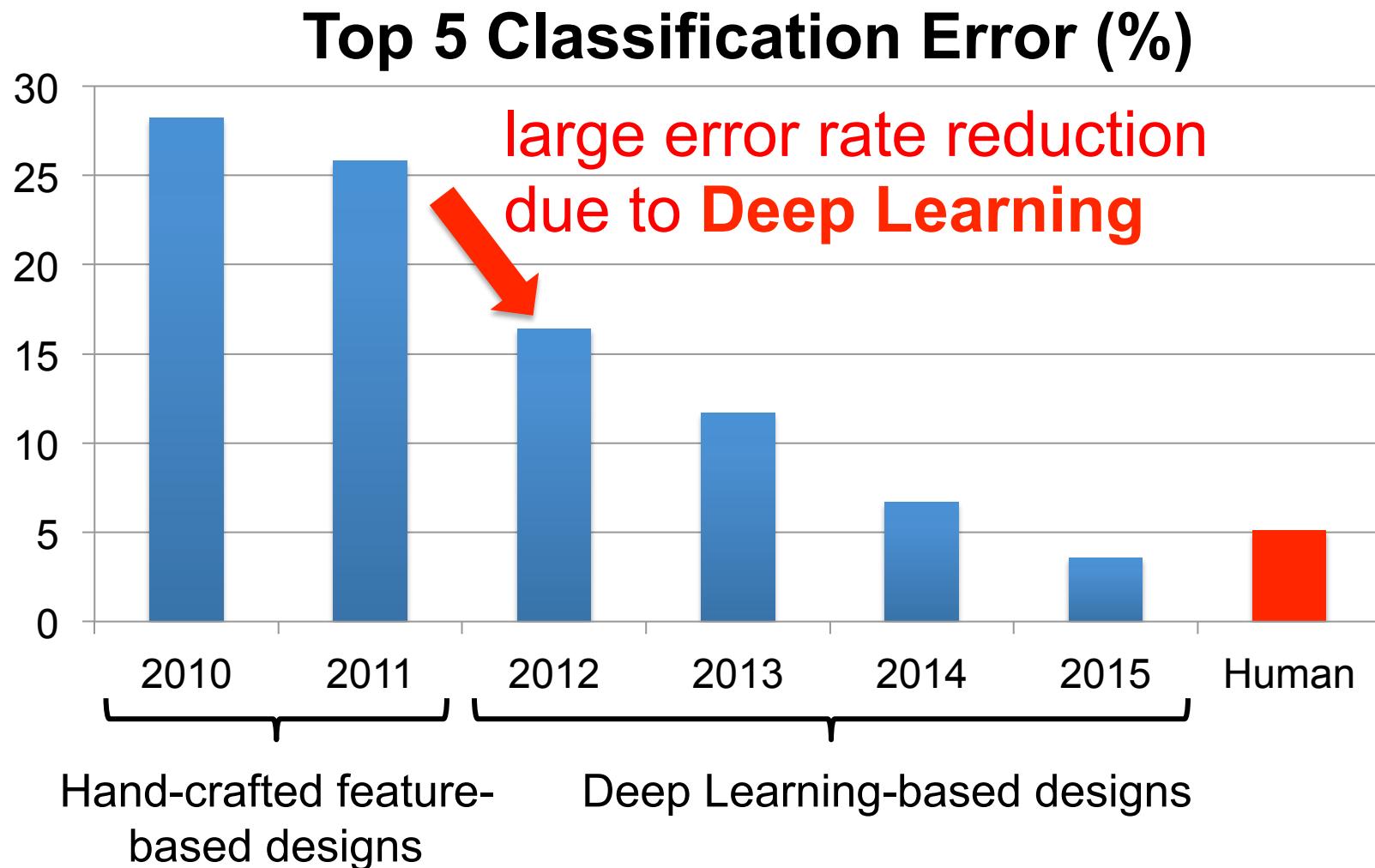
2.5 Petabytes of customer data hourly



300 hours of video uploaded every minute

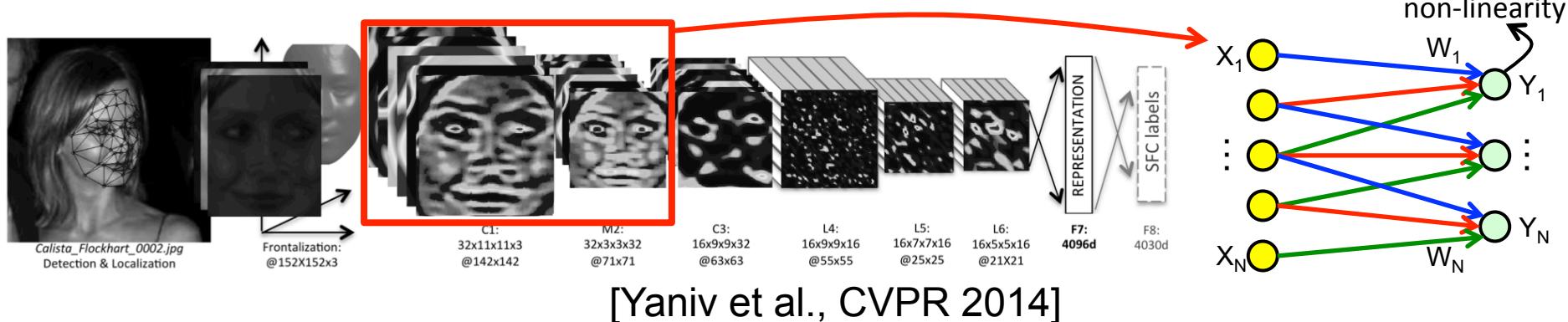


ImageNet: Image Classification Task

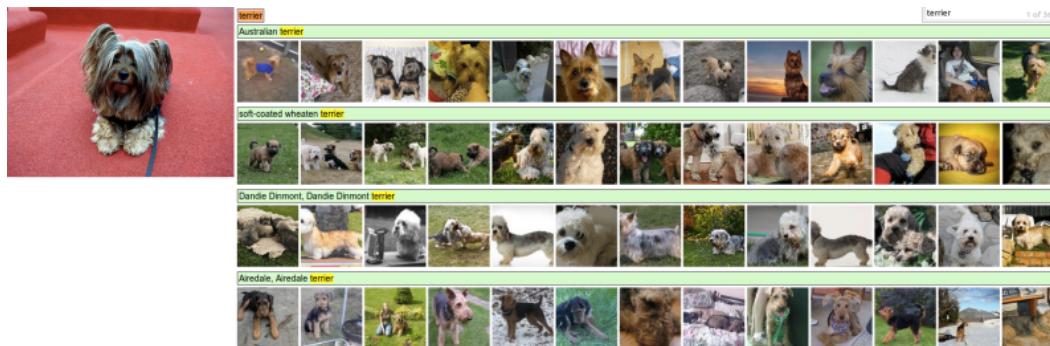


Human or *Superhuman* Accuracy Level

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



- Fine grained category recognition (e.g. dogs, monkeys, snakes, birds)
 - Deep learning errors: 7 vs. Human errors: 28



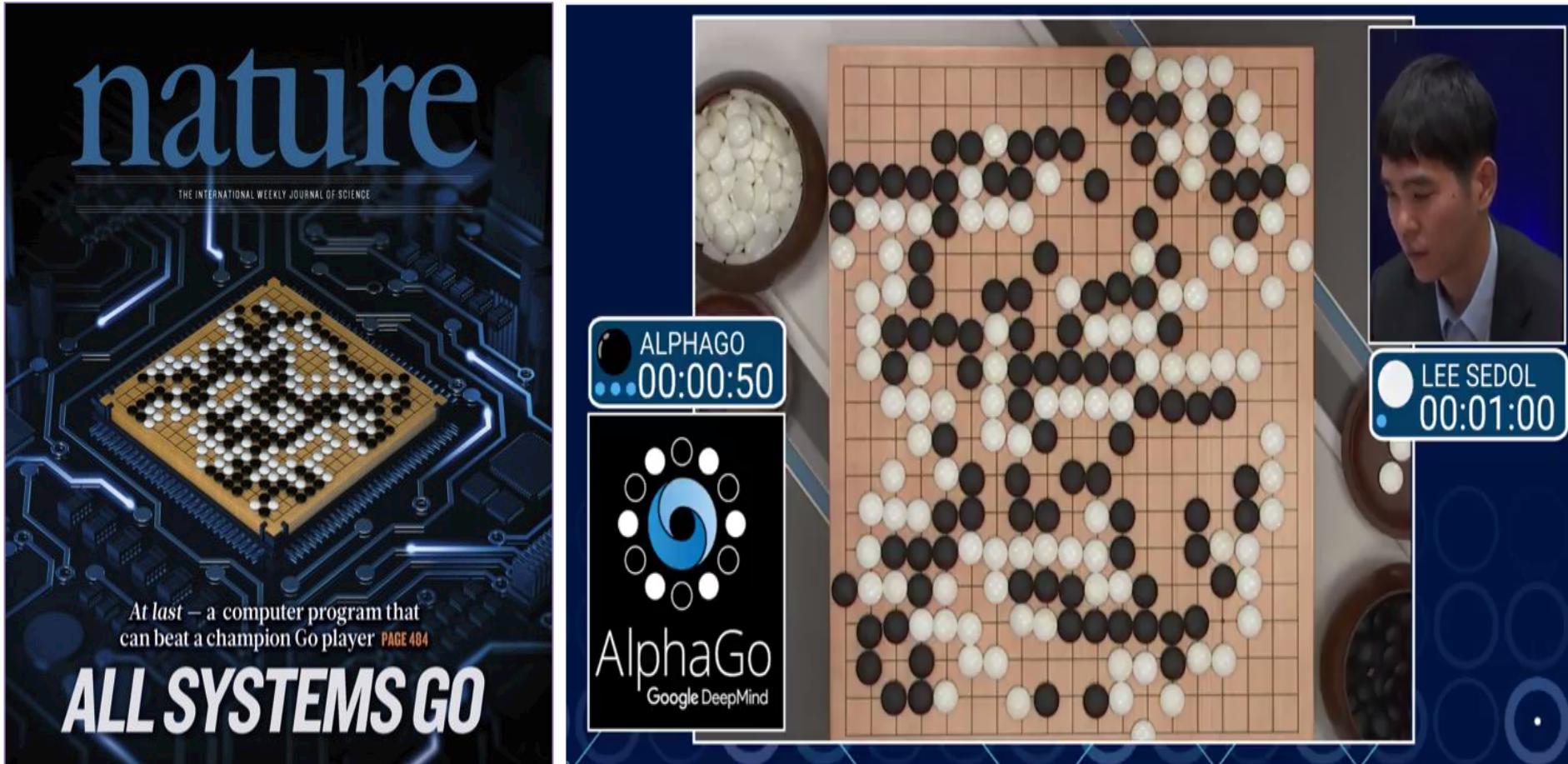
120 species of dogs

[O. Russakovsky et al., IJCV 2015]

Deep Learning on Games

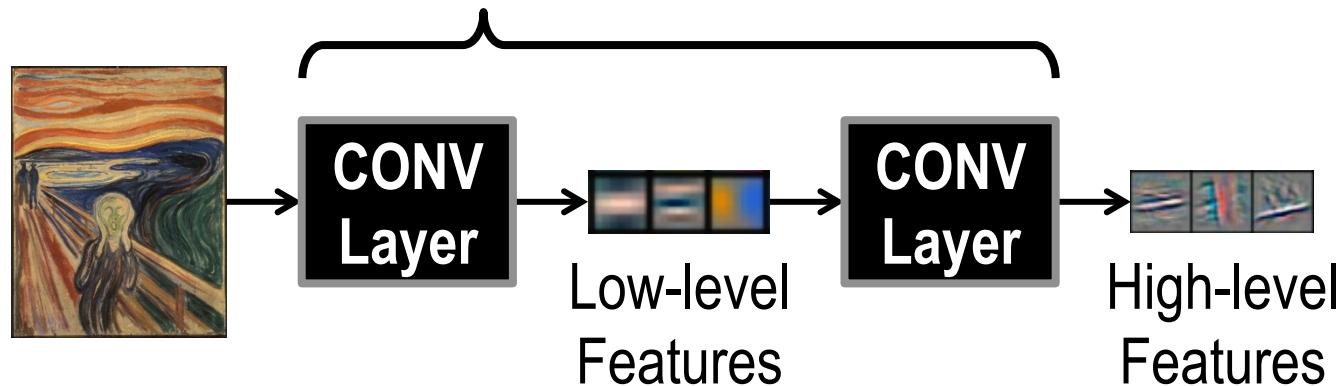
Google DeepMind AlphaGo

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

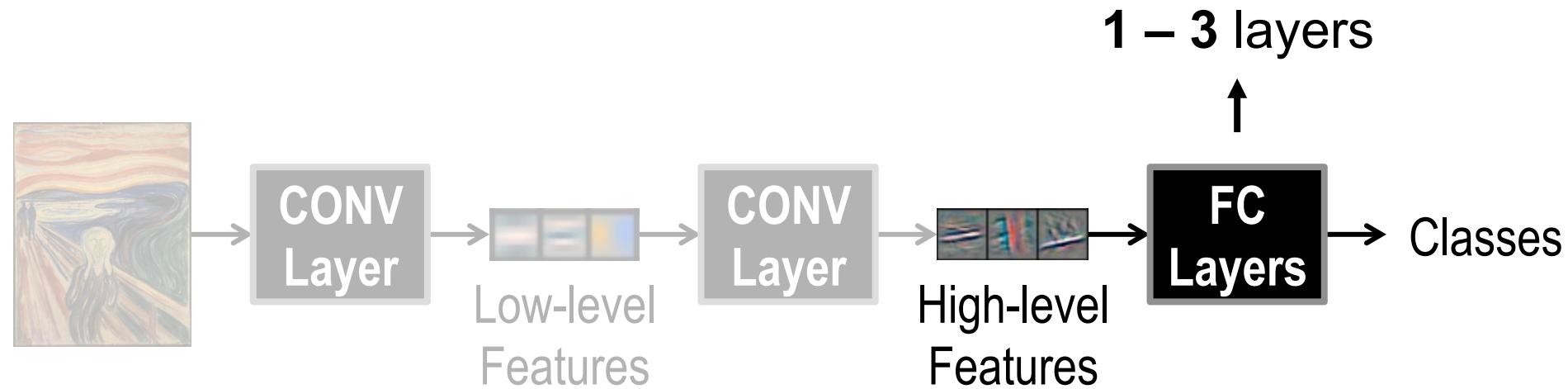


Deep Convolutional Neural Networks

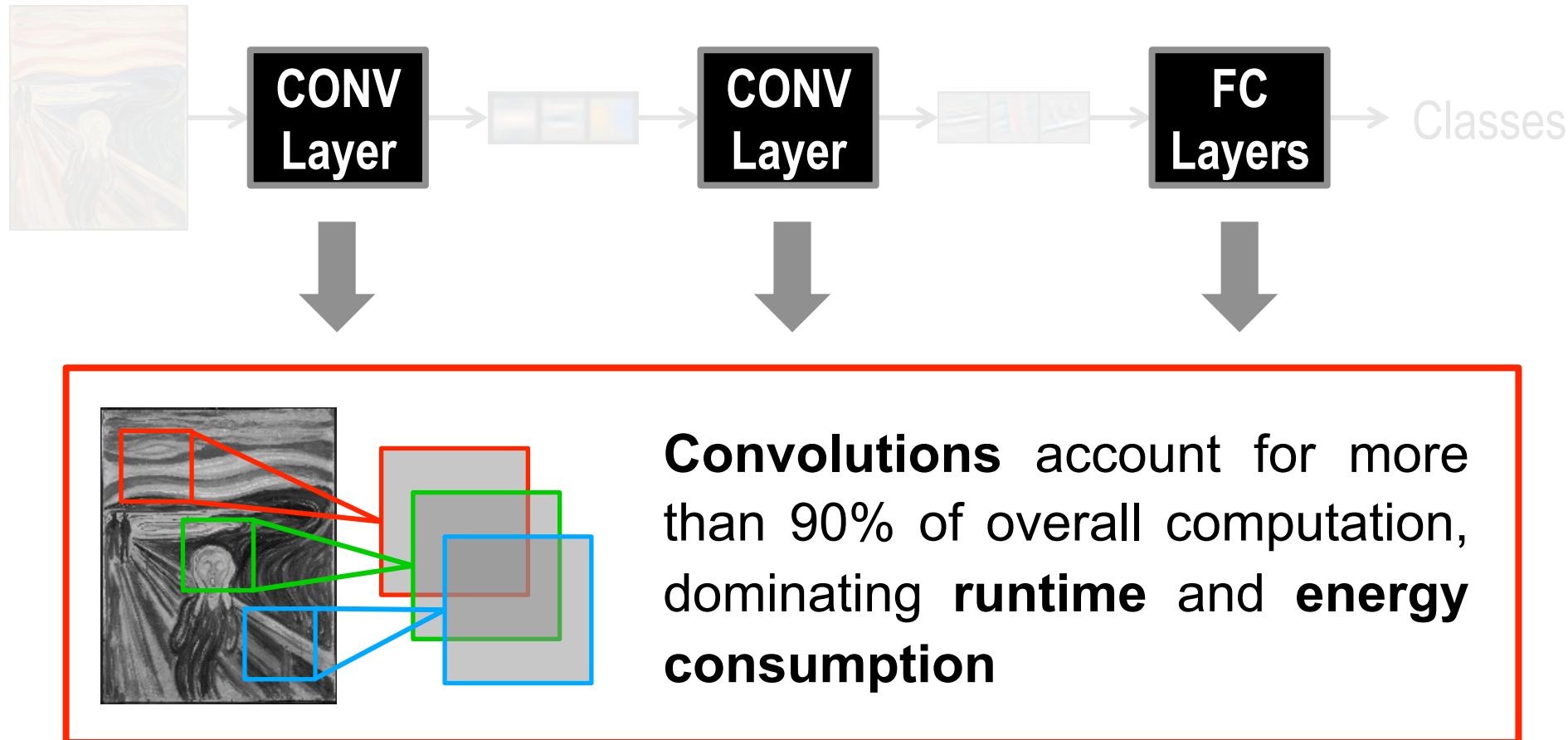
Modern *deep* CNN: up to 1000 CONV layers



Deep Convolutional Neural Networks

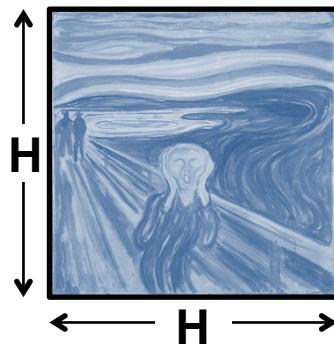
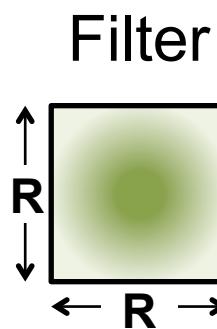


Deep Convolutional Neural Networks

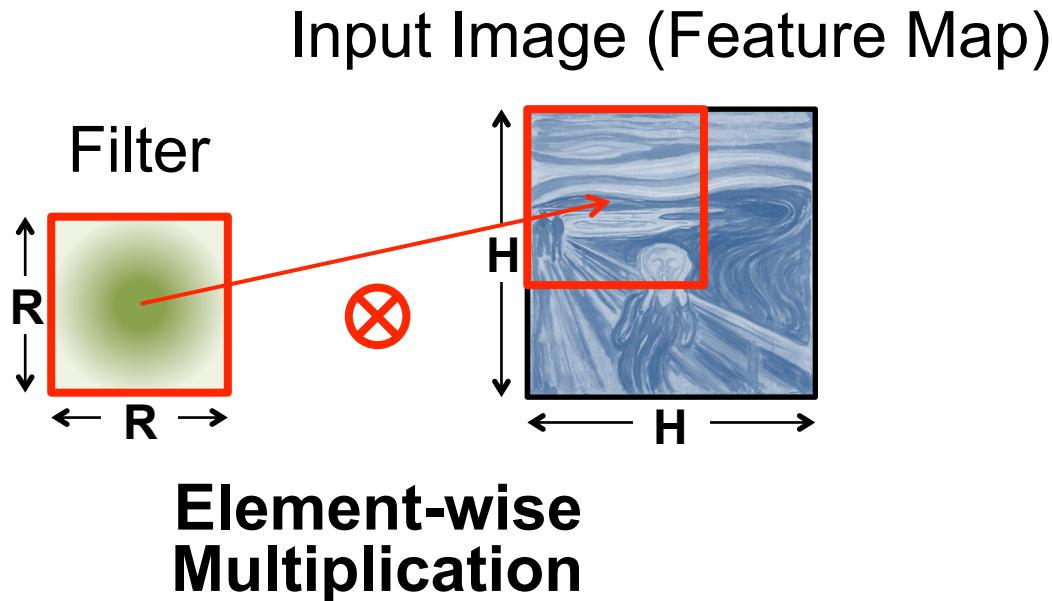


High-Dimensional CNN Convolution

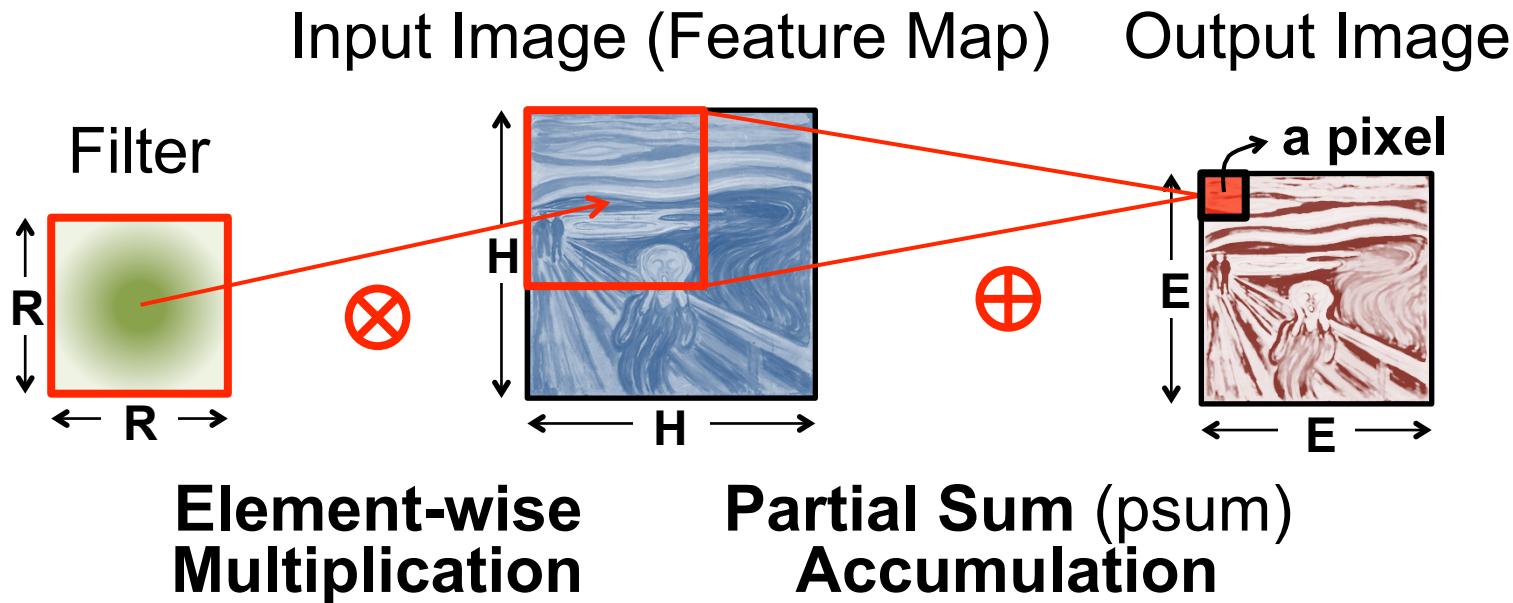
Input Image (Feature Map)



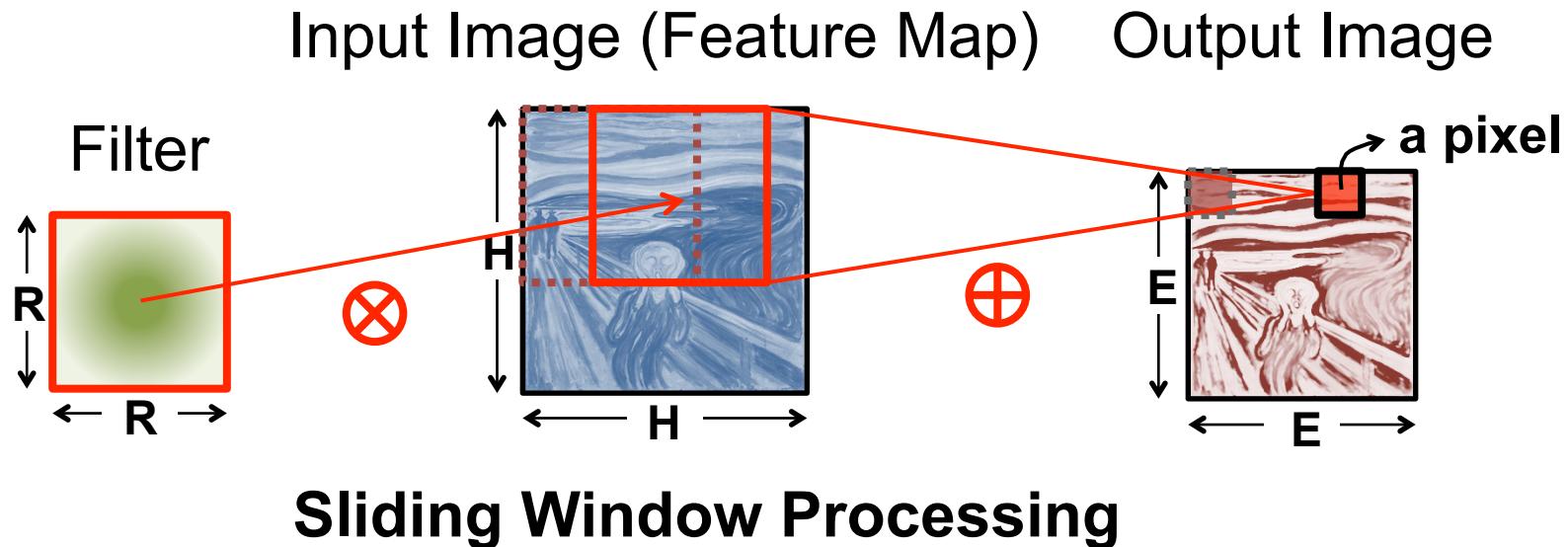
High-Dimensional CNN Convolution



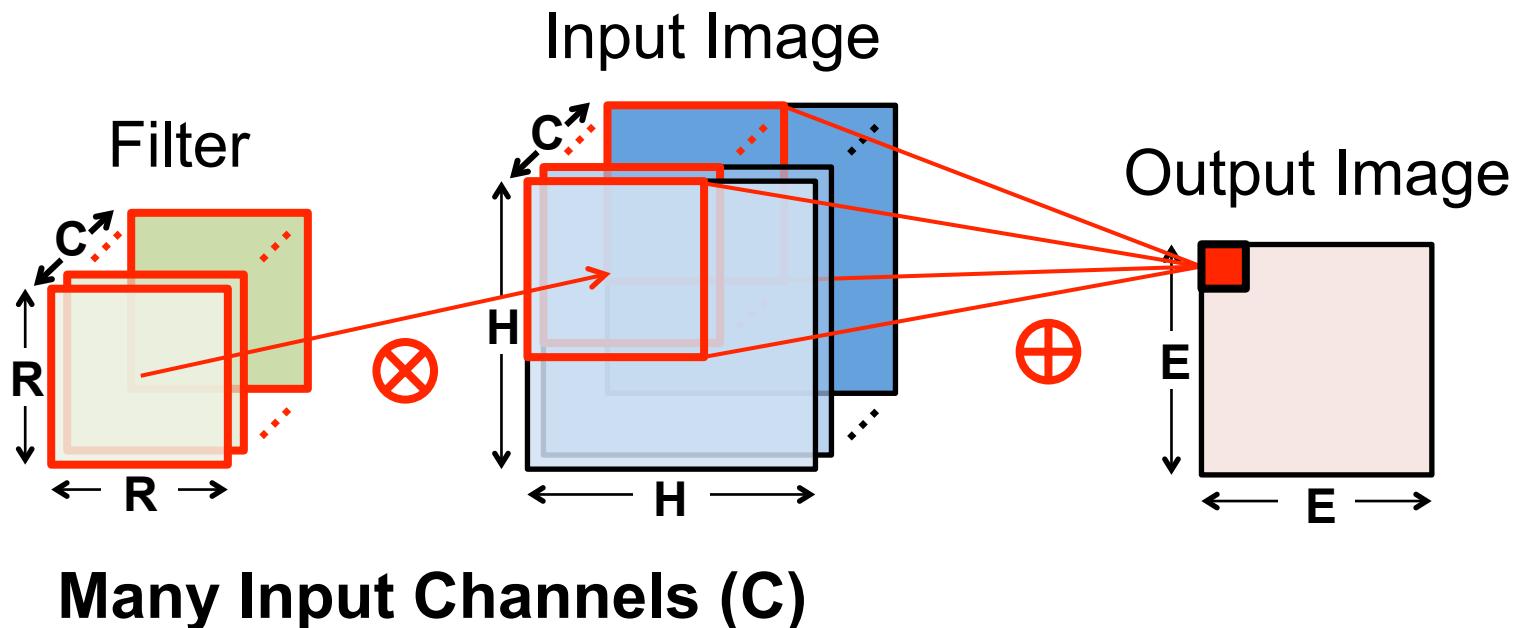
High-Dimensional CNN Convolution



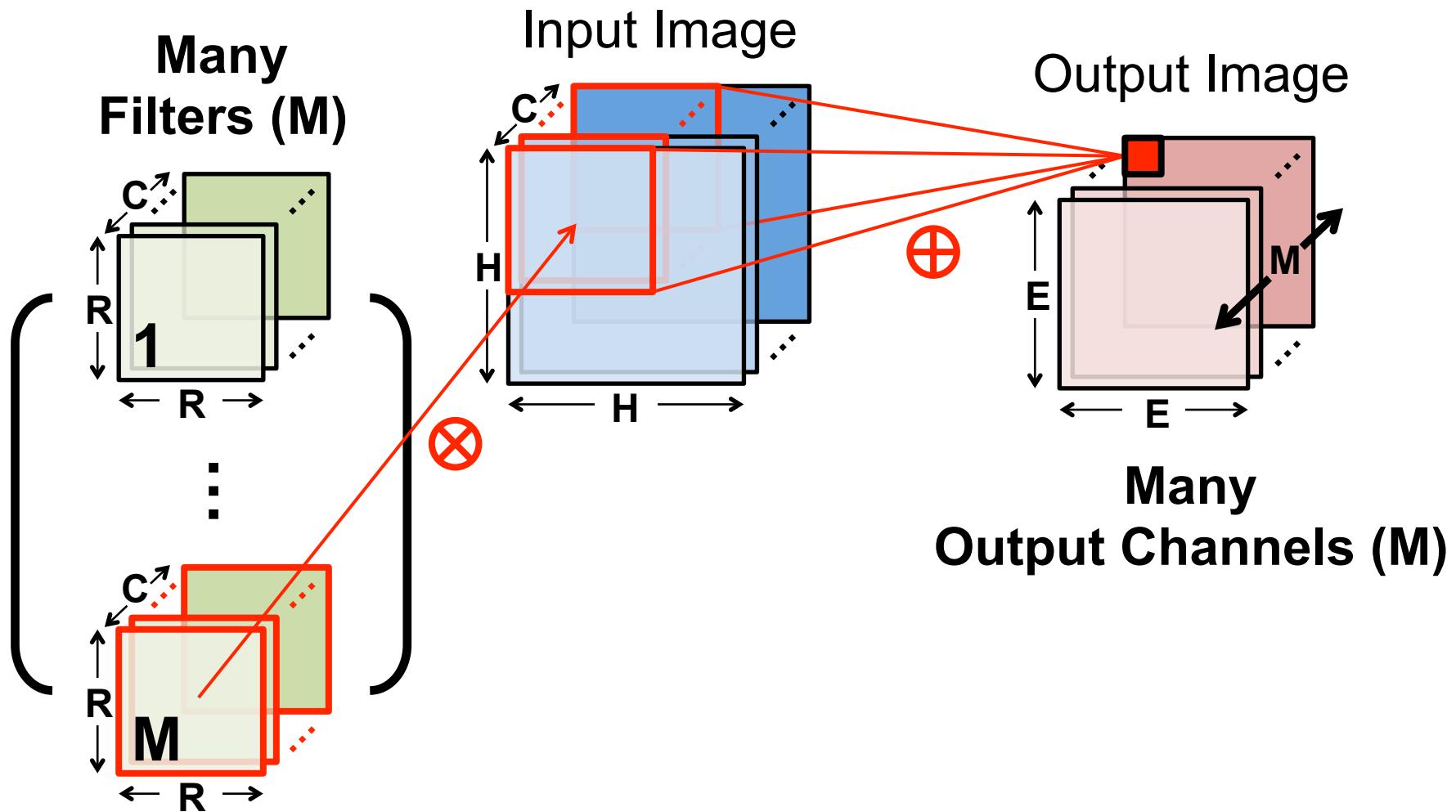
High-Dimensional CNN Convolution



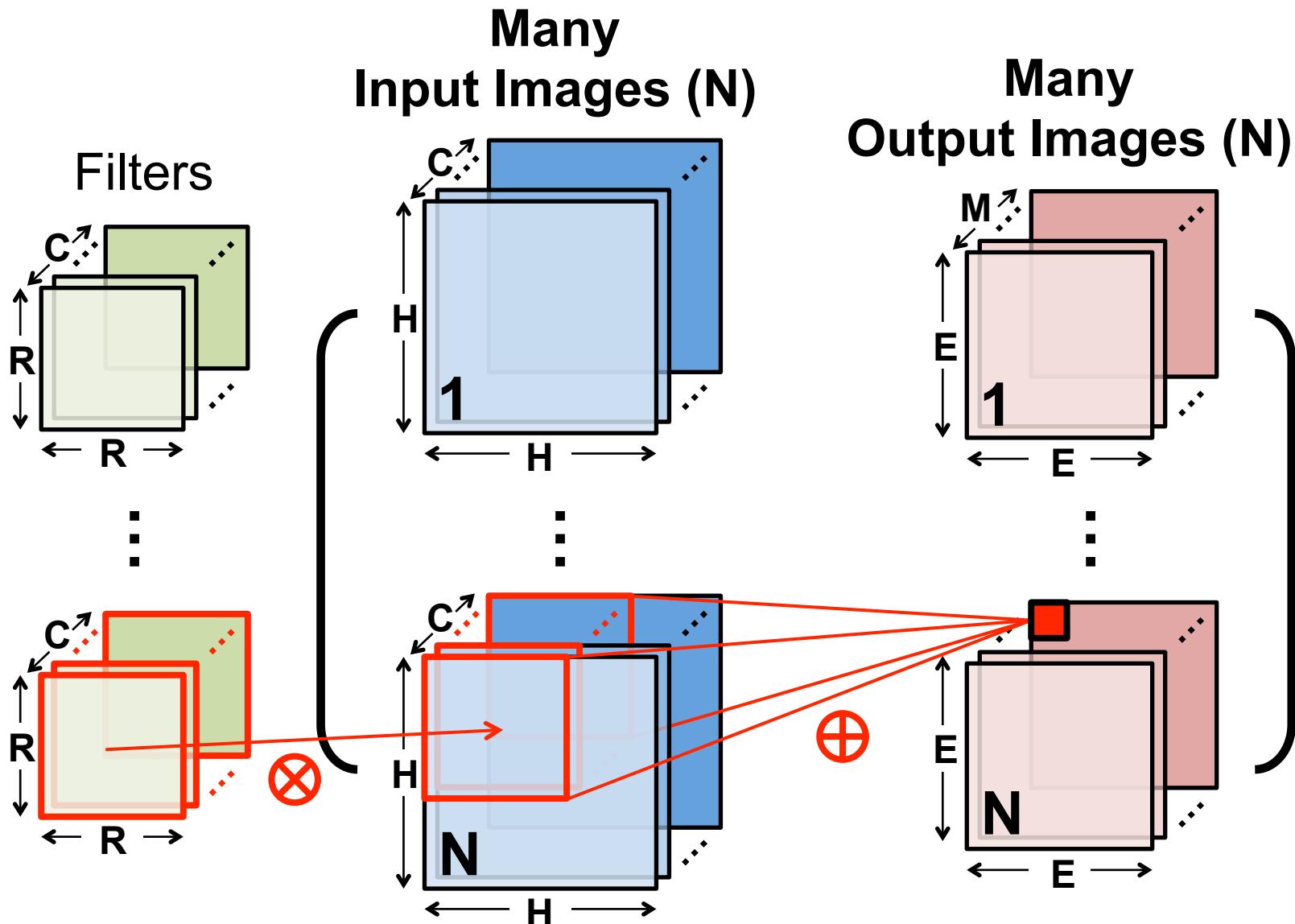
High-Dimensional CNN Convolution



High-Dimensional CNN Convolution



High-Dimensional CNN Convolution

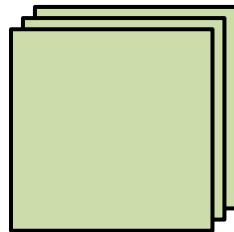


Large Sizes with Varying Shapes

AlexNet¹ Convolutional Layer Configurations

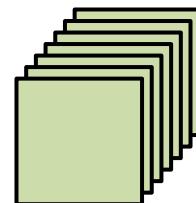
Layer	Filter Size (R)	# Filters (M)	# Channels (C)	Stride
1	11x11	96	3	4
2	5x5	256	48	1
3	3x3	384	256	1
4	3x3	384	192	1
5	3x3	256	192	1

Layer 1



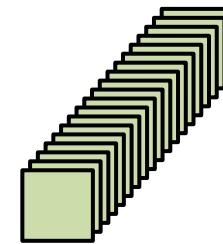
34k Params
105M MACs

Layer 2



307k Params
224M MACs

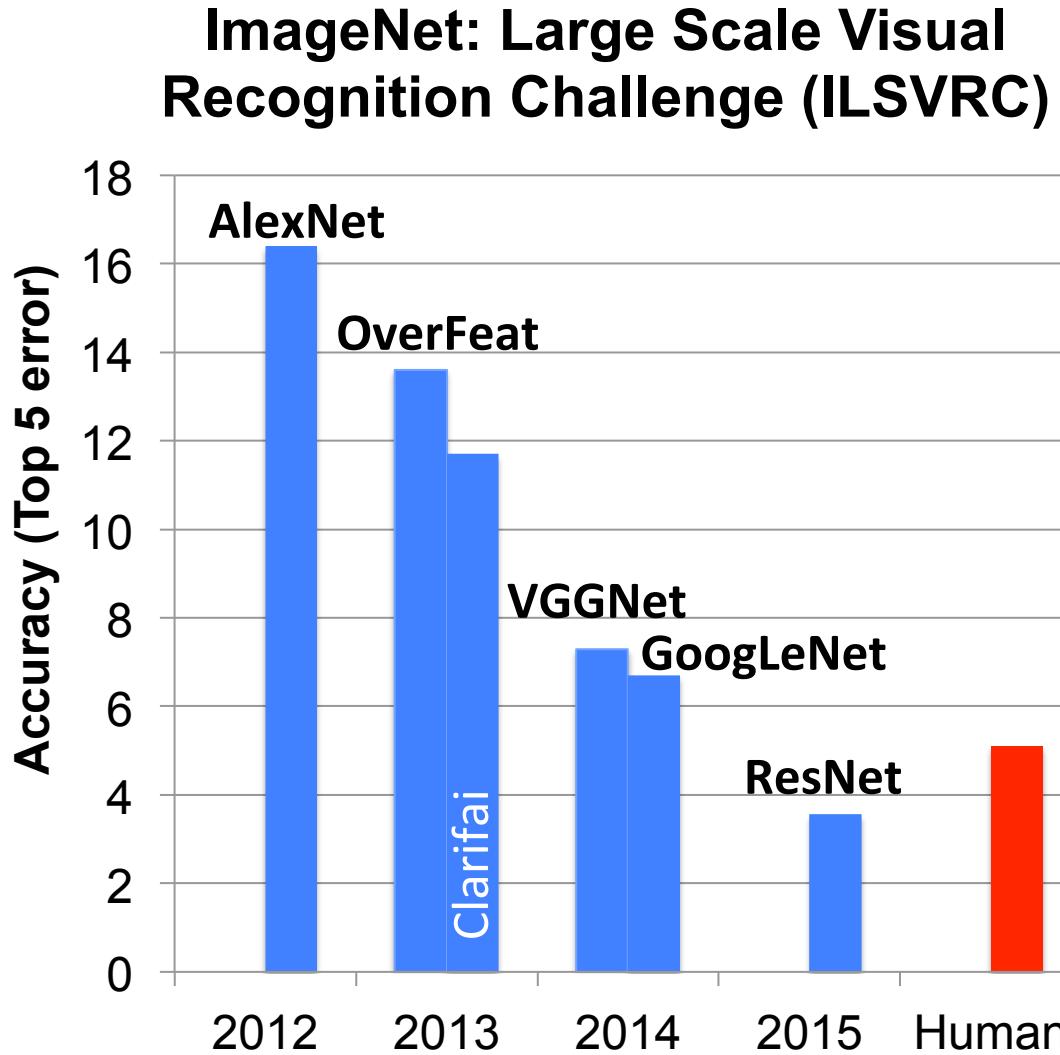
Layer 3



885k Params
150M MACs

Popular DNNs

- LeNet (1998)
- AlexNet (2012)
- OverFeat (2013)
- VGGNet (2014)
- GoogleNet (2014)
- ResNet (2015)



[O. Russakovsky et al., IJCV 2015]

Summary of Popular DNNs

Metrics	LeNet-5	AlexNet	VGG-16	GoogLeNet (v1)	ResNet-50
Top-5 error	n/a	16.4	7.4	6.7	5.3
Input Size	28x28	227x227	224x224	224x224	224x224
# of CONV Layers	2	5	16	21 (depth)	49
Filter Sizes	5	3, 5, 11	3	1, 3, 5, 7	1, 3, 7
# of Channels	1, 6	3 - 256	3 - 512	3 - 1024	3 - 2048
# of Filters	6, 16	96 - 384	64 - 512	64 - 384	64 - 2048
Stride	1	1, 4	1	1, 2	1, 2
# of Weights	2.6k	2.3M	14.7M	6.0M	23.5M
# of MACs	283k	666M	15.3G	1.43G	3.86G
# of FC layers	2	3	3	1	1
# of Weights	58k	58.6M	124M	1M	2M
# of MACs	58k	58.6M	124M	1M	2M
Total Weights	60k	61M	138M	7M	25.5M
Total MACs	341k	724M	15.5G	1.43G	3.9G

CONV Layers increasingly important!

Complexity versus Difficulty of Task

- Evaluate hardware using the appropriate DNN model and dataset
 - Difficult tasks typically require larger models
 - Different datasets for different tasks

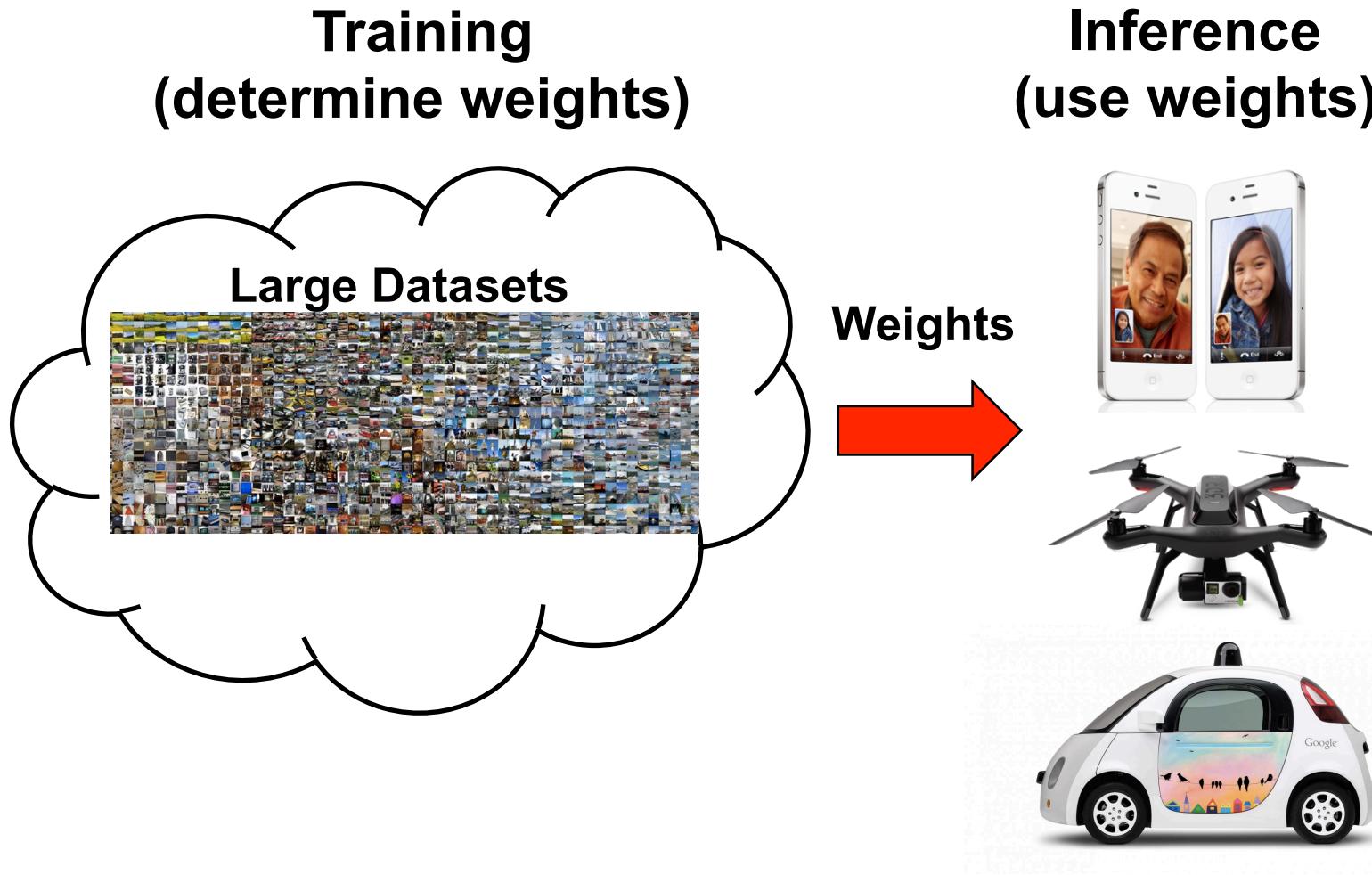
MNIST

3 6 8 1 7 9 6 6 9 1
6 7 5 7 8 6 3 4 8 5
2 1 7 9 7 1 2 8 4 6
4 8 1 9 0 1 8 8 9 4
7 6 1 8 6 4 1 5 6 0
7 5 9 2 6 5 8 1 9 7
2 2 2 2 2 3 4 4 8 0
0 2 3 8 0 7 3 8 5 7
0 1 4 6 4 6 0 2 4 3
7 1 2 8 7 6 9 8 6 1

ImageNet



Training vs. Inference



Challenges

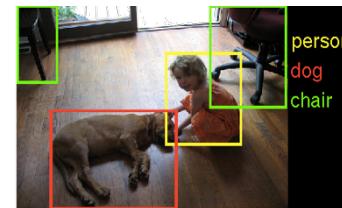
Key Metrics

- **Accuracy**
 - Measured on a publicly available dataset
 - Popular DNN Models
- **Programmability**
 - Support multiple applications
 - Different weights
- **Energy/Power**
 - Energy per operation
 - DRAM Bandwidth
- **Throughput/Latency**
 - GOPS, frame rate, delay
- **Cost**
 - Area (memory and logic size)

ImageNet



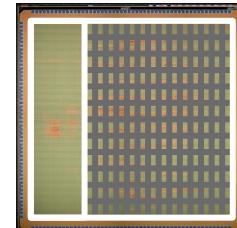
Computer
Vision



Speech
Recognition



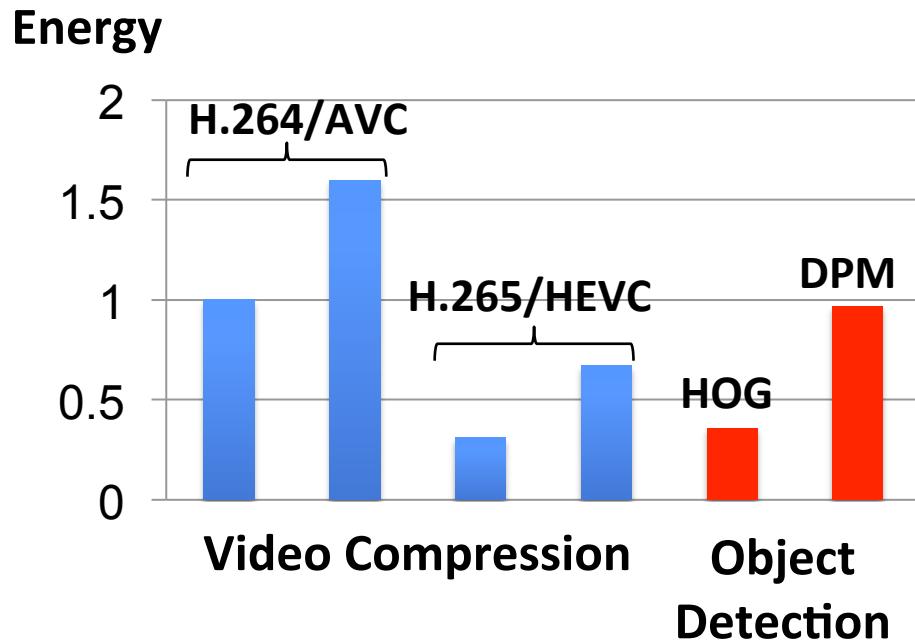
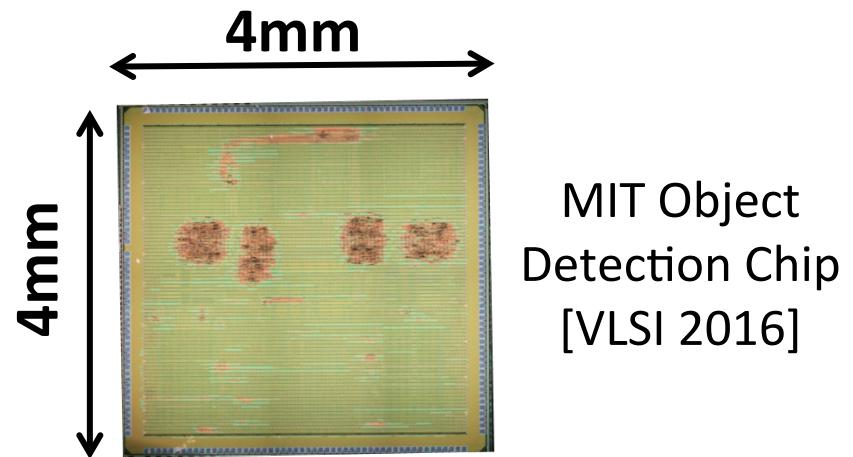
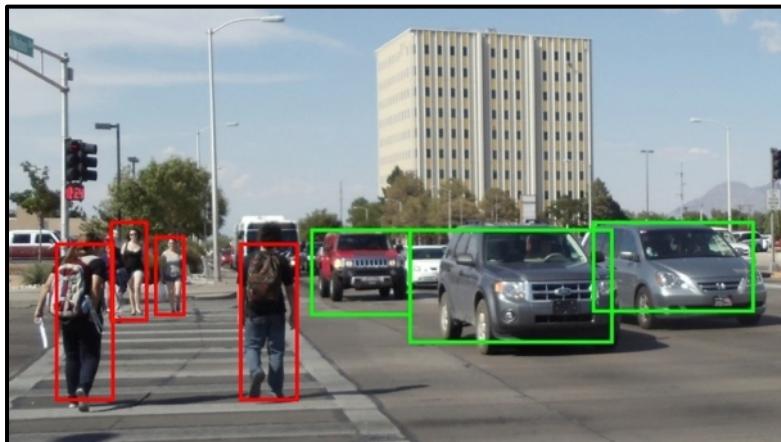
Chip



DRAM

Typical Constraints on Video Coding

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



Opportunities in Architecture

GPUs and CPUs Targeting Deep Learning

Intel Knights Landing (2016)



Nvidia PASCAL GP100 (2016)



Knights Mill: next gen Xeon
Phi “optimized for deep
learning”

Use matrix multiplication libraries on CPUs and GPUs

Accelerate Matrix Multiplication

- Implementation: **Matrix Multiplication (GEMM)**
 - **CPU:** OpenBLAS, Intel MKL, etc
 - **GPU:** cuBLAS, cuDNN, etc
- Optimized by tiling to storage hierarchy

Map DNN to a Matrix Multiplication

- Convert to matrix mult. using the **Toeplitz Matrix**

Convolution:



**Toeplitz Matrix
(w/ redundant data)**

Matrix Mult:

Map DNN to a Matrix Multiplication

- Convert to matrix mult. using the **Toeplitz Matrix**

Convolution:



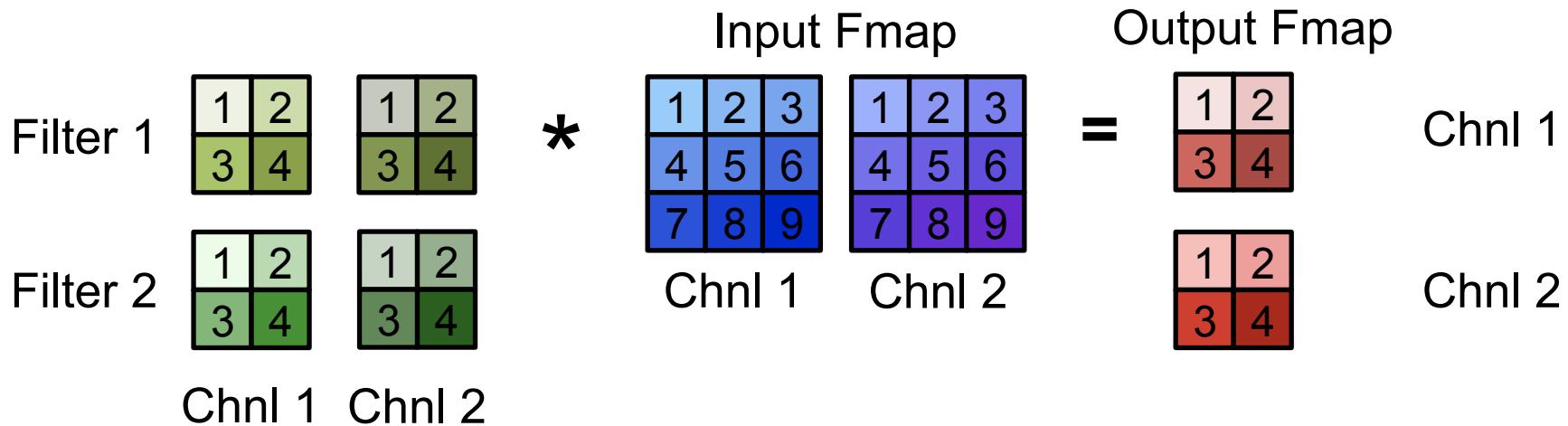
**Toeplitz Matrix
(w/ redundant data)**

Matrix Mult:

Data is repeated

Map DNN to a Matrix Multiplication

- Multiple Channels and Filters



Map DNN to a Matrix Multiplication

- Multiple Channels and Filters

**Toeplitz Matrix
(w/ redundant data)**

	Chnl 1				Chnl 2			
Filter 1	1	2	3	4	1	2	3	4
Filter 2	1	2	3	4	1	2	3	4

$$\times \begin{array}{|c|c|c|c|} \hline 1 & 2 & 4 & 5 \\ \hline 2 & 3 & 5 & 6 \\ \hline 4 & 5 & 7 & 8 \\ \hline 5 & 6 & 8 & 9 \\ \hline 1 & 2 & 4 & 5 \\ \hline 2 & 3 & 5 & 6 \\ \hline 4 & 5 & 7 & 8 \\ \hline 5 & 6 & 8 & 9 \\ \hline \end{array} = \begin{array}{|c|c|c|c|} \hline 1 & 2 & 3 & 4 \\ \hline 1 & 2 & 3 & 4 \\ \hline \end{array}$$

Chnl 1

Chnl 2

Chnl 1

Chnl 2

Goal: Reduced number of operations to **increase throughput**

Computation Transformations

- **Goal:** Bitwise same result, but reduce number of operations
- Focuses mostly on compute

Analogy: Gauss's Multiplication Algorithm

$$(a + bi)(c + di) = (ac - bd) + (bc + ad)i.$$

4 multiplications + 3 additions

$$k_1 = c \cdot (a + b)$$

$$k_2 = a \cdot (d - c)$$

$$k_3 = b \cdot (c + d)$$

$$\text{Real part} = k_1 - k_3$$

$$\text{Imaginary part} = k_1 + k_2.$$

3 multiplications + 5 additions

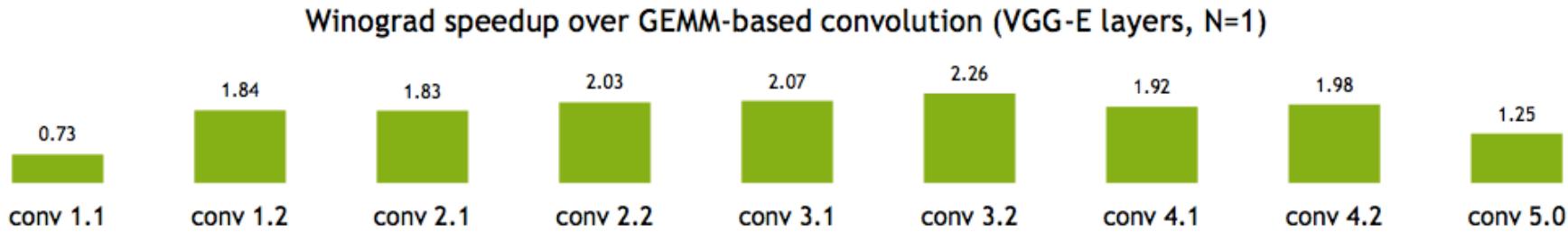
Reduce number of multiplications,
but increase number of additions

Reduce Operations in Matrix Multiplication

- **Winograd** [Lavin, CVPR 2016]
 - **Pro:** 2.25x speed up for 3x3 filter
 - **Con:** Specialized processing depending on filter size
- **Fast Fourier Transform** [Mathieu, ICLR 2014]
 - **Pro:** Direct convolution $O(N_o^2 N_f^2)$ to $O(N_o^2 \log_2 N_o)$
 - **Con:** Increase storage requirements
- **Strassen** [Cong, ICANN 2014]
 - **Pro:** $O(N^3)$ to $(N^{2.807})$
 - **Con:** Numerical stability

Winograd Performance Varies

Optimal convolution algorithm depends on convolution layer dimensions



Meta-parameters (data layouts, texture memory) afford higher performance

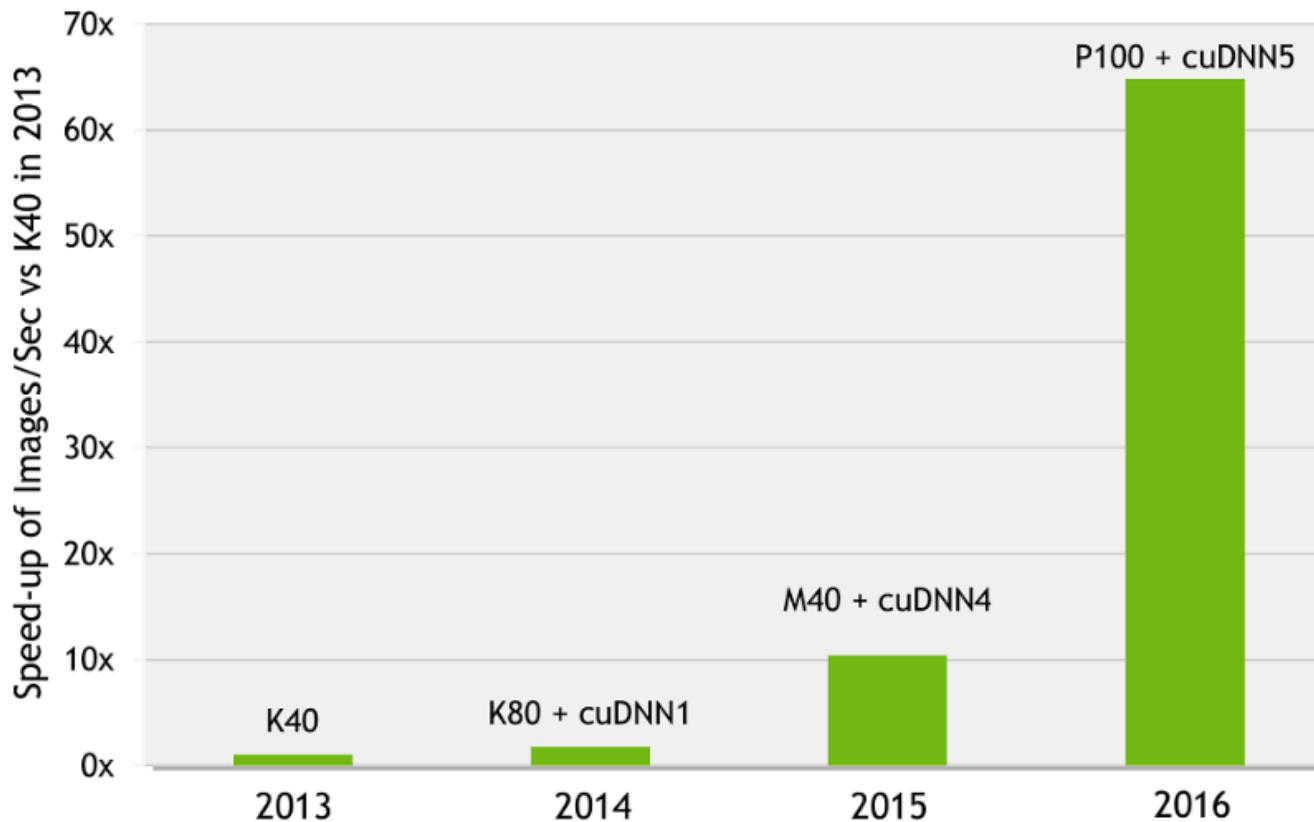
Using texture memory for convolutions: **13% inference speedup**

(GoogLeNet, batch size 1)

Source: Nvidia

cuDNN: Speed up with Transformations

60x Faster Training in 3 Years



AlexNet training throughput on:

CPU: 1x E5-2680v3 12 Core 2.5GHz. 128GB System Memory, Ubuntu 14.04

M40 bar: 8x M40 GPUs in a node, P100: 8x P100 NVLink-enabled

Source: Nvidia

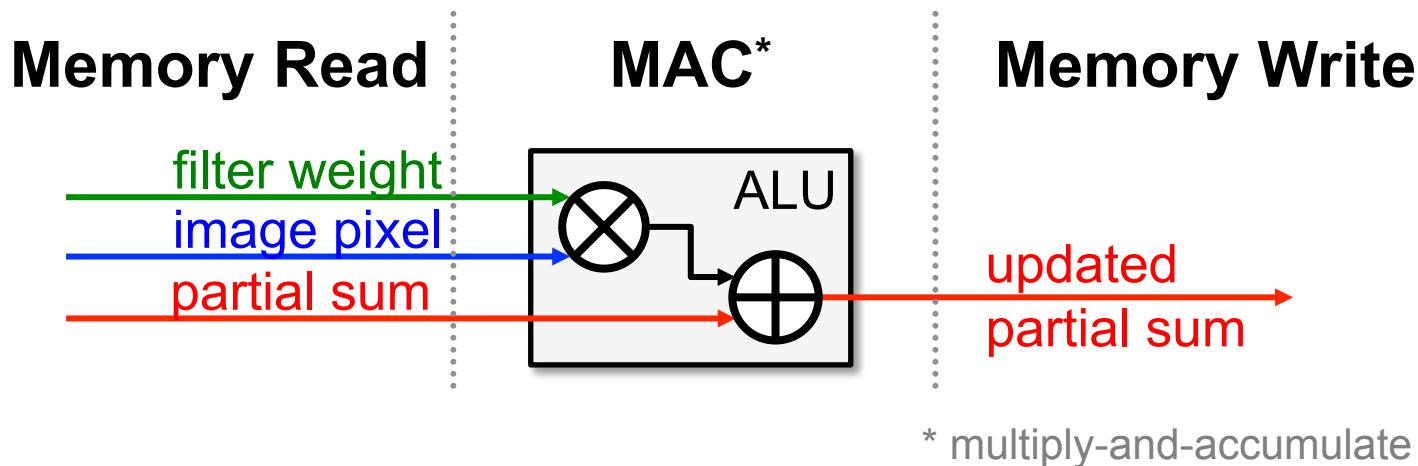
Specialized Hardware (Accelerators)

Properties We Can Leverage

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

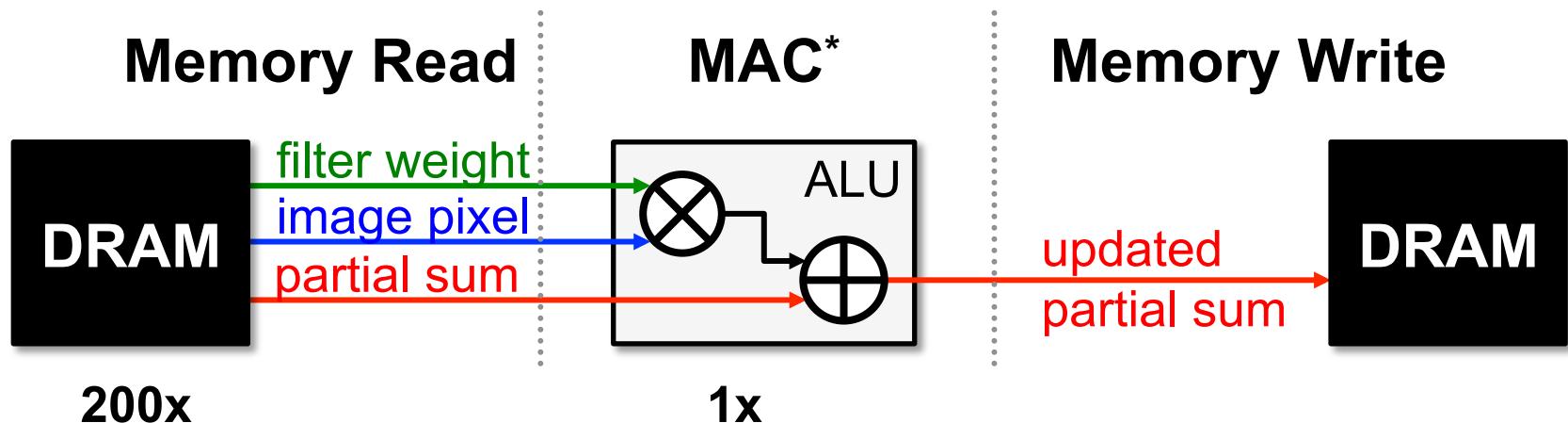
Properties We Can Leverage

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



Properties We Can Leverage

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

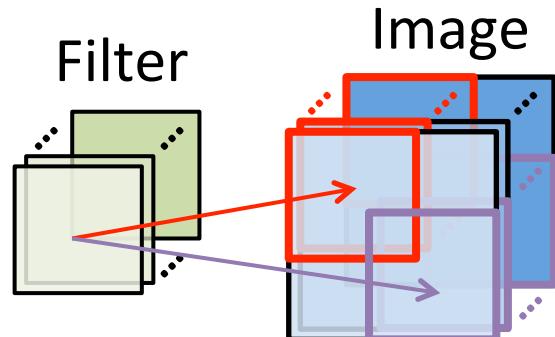


Worst Case: all memory R/W are **DRAM** accesses

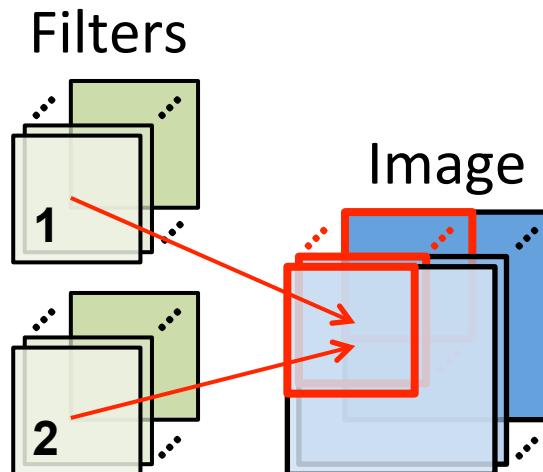
- Example: AlexNet [NIPS 2012] has **724M** MACs
→ **2896M** DRAM accesses required

Properties We Can Leverage

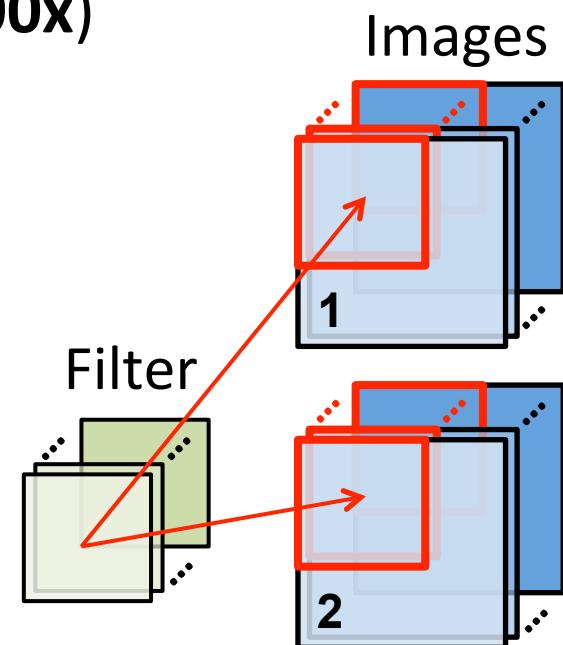
- Operations exhibit **high parallelism**
→ high throughput possible
- Input data reuse** opportunities (up to 500x)
→ exploit low-cost memory



**Convolutional
Reuse**
(pixels, weights)



**Image
Reuse**
(pixels)



**Filter
Reuse**
(weights)

Images

Filters

Image

Filter

2

1

Filter

Image

2

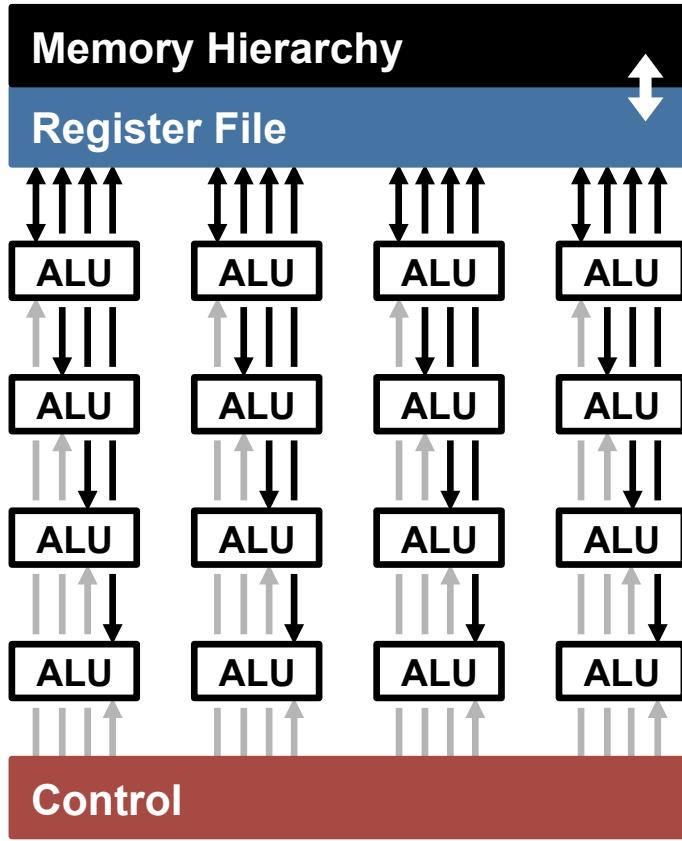
Filter

1

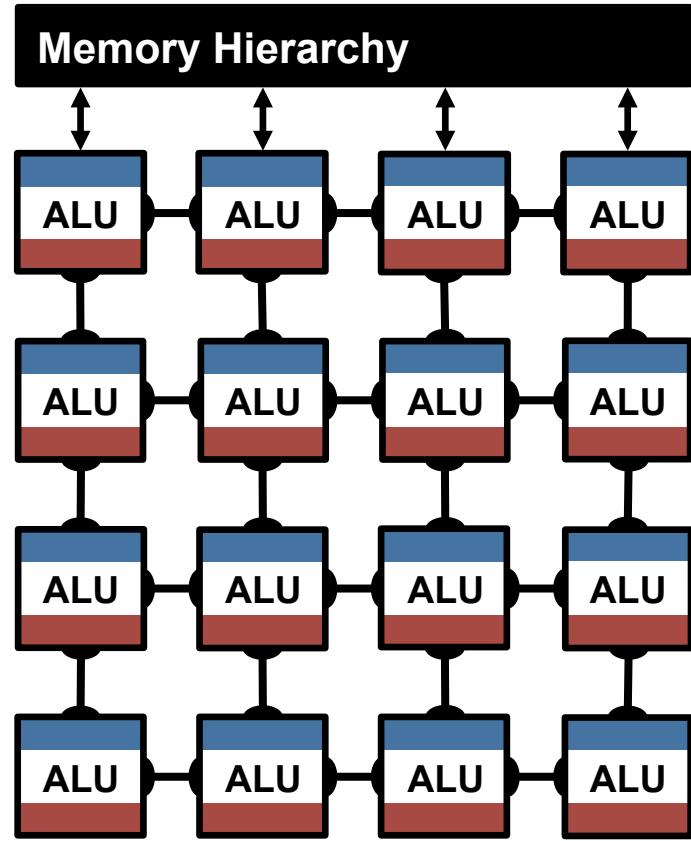
2

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)

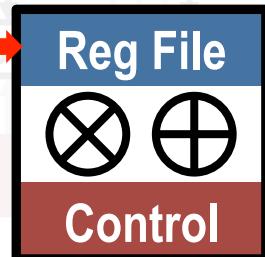
Inter-PE Communication

Sharing among regions of PEs

Processing
Element (PE)

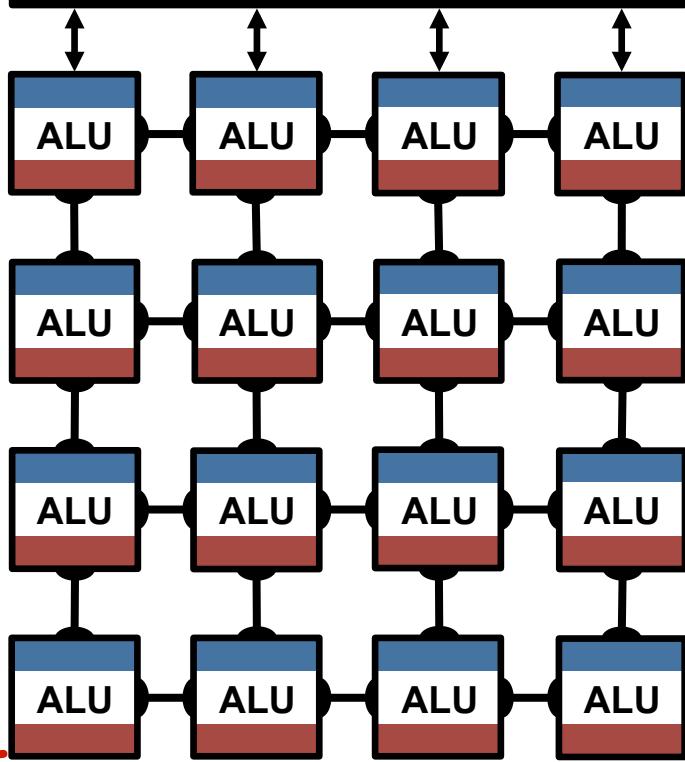
0.5 – 1.0 kB

Control



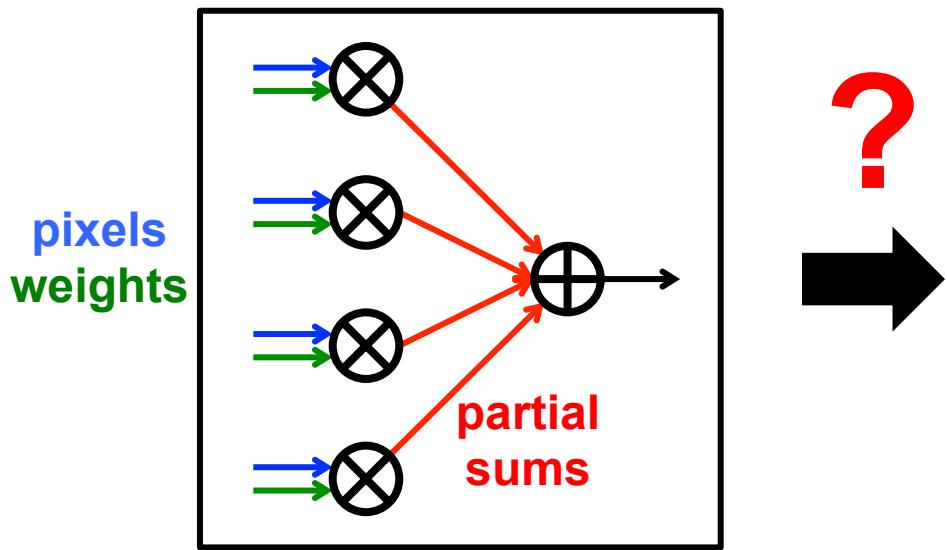
Spatial Architecture
(Dataflow Processing)

Memory Hierarchy



How to Map the Dataflow?

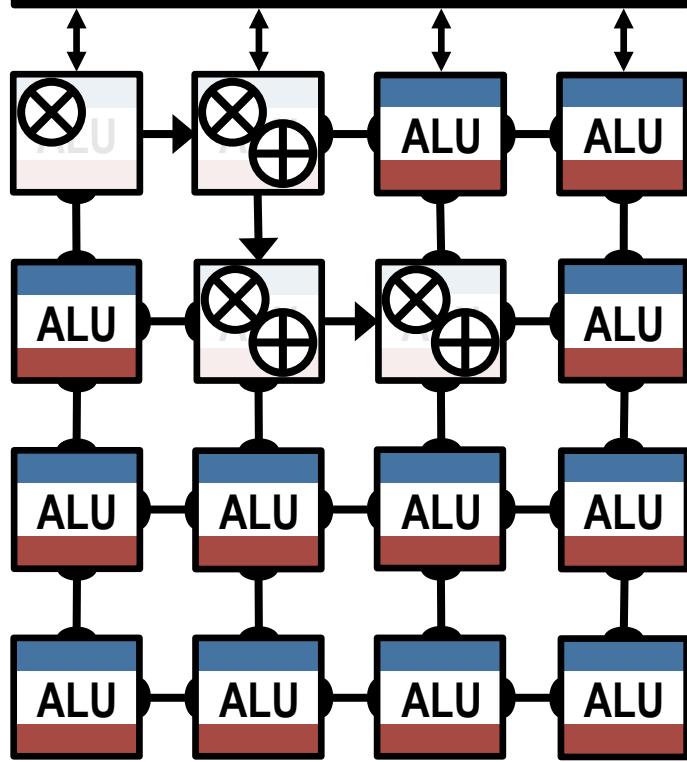
CNN Convolution



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

Spatial Architecture (Dataflow Processing)

Memory Hierarchy

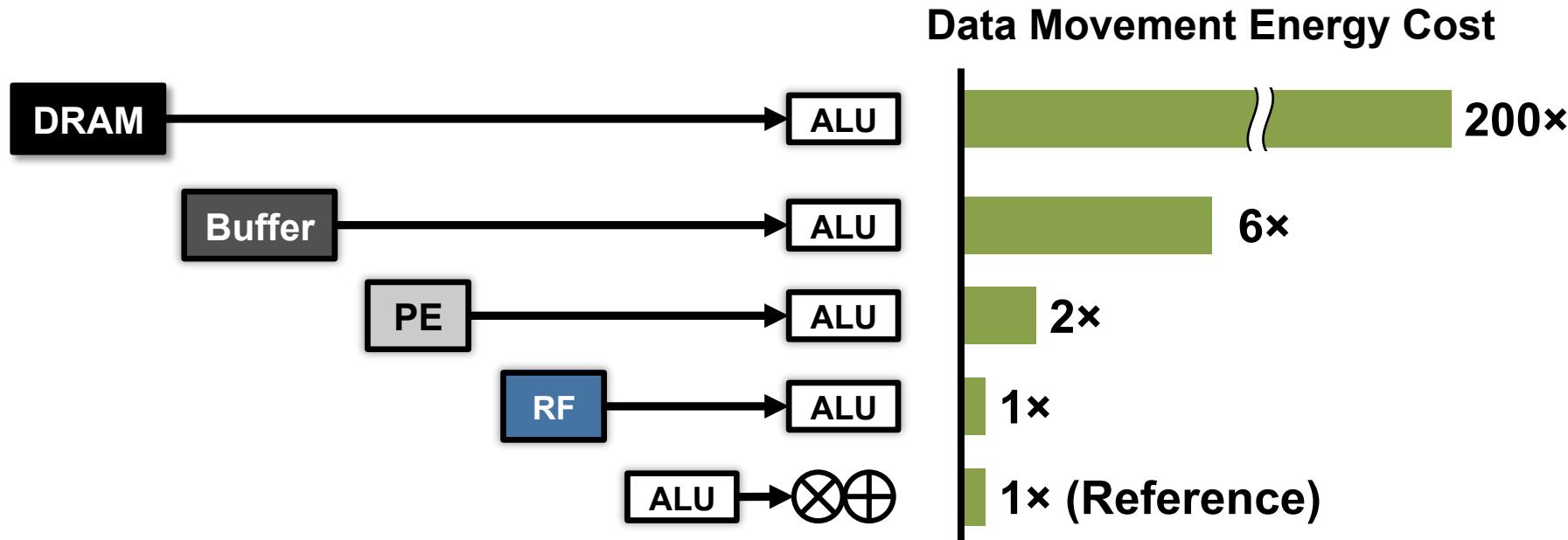
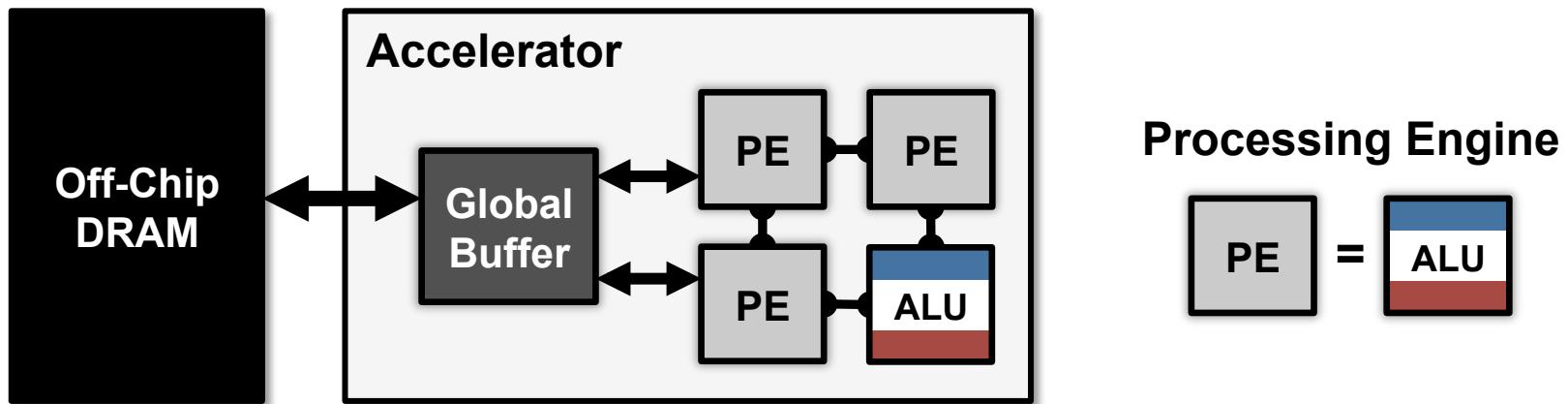


Energy-Efficient Dataflow

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

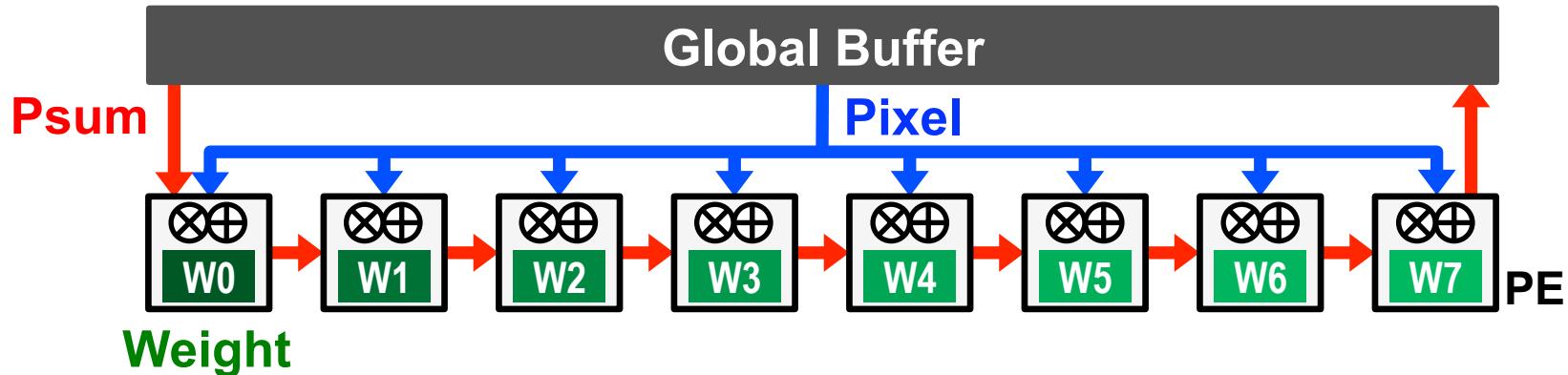
Maximize data reuse and accumulation at RF

Data Movement is Expensive



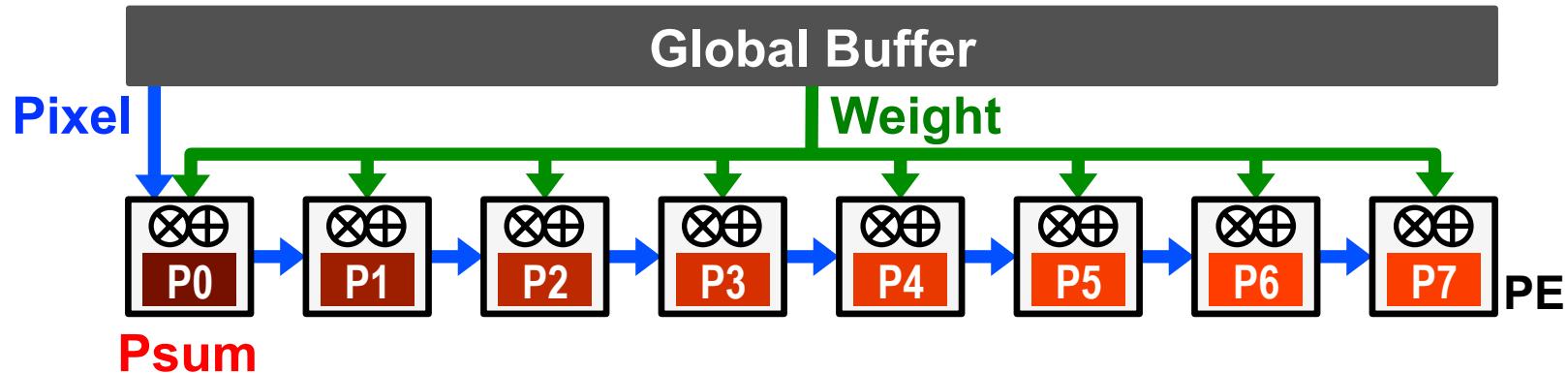
Maximize data reuse at lower levels of hierarchy

Weight Stationary (WS)



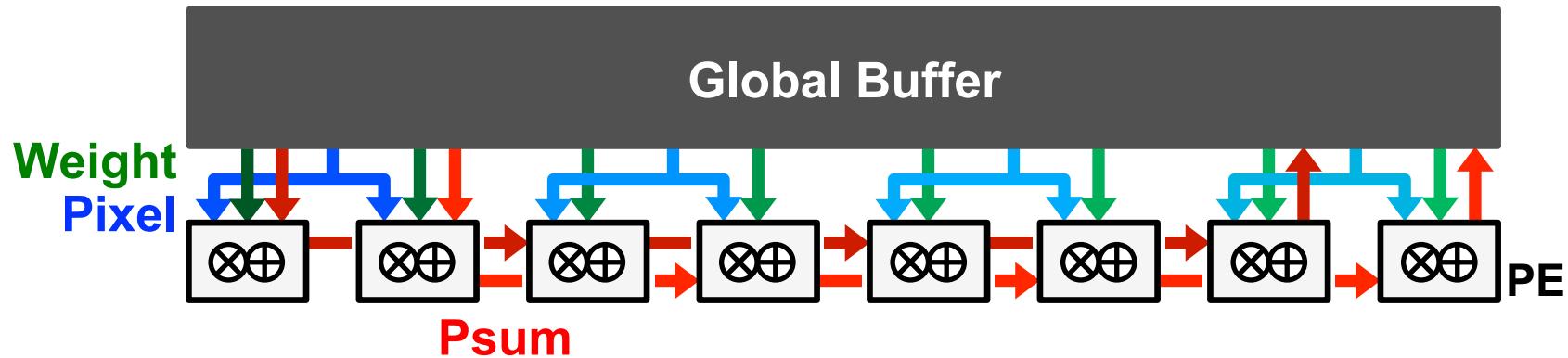
- **Minimize weight** read energy consumption
 - maximize convolutional and filter reuse of weights
- **Examples:**
 - [Chakradhar, /SCA 2010] [nn-X (NeuFlow), CVPRW 2014]
 - [Park, /SSCC 2015] [Origami, GLSVLSI 2015]

Output Stationary (OS)



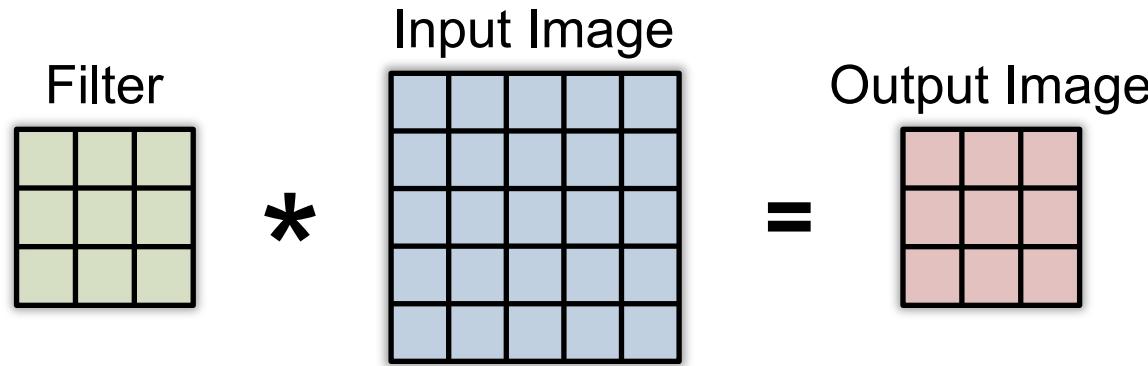
- **Minimize partial sum R/W energy consumption**
 - maximize local accumulation
- **Examples:**
 - [Gupta, ICML 2015]
 - [ShiDianNao, ISCA 2015]
 - [Peemen, ICCD 2013]

No Local Reuse (NLR)

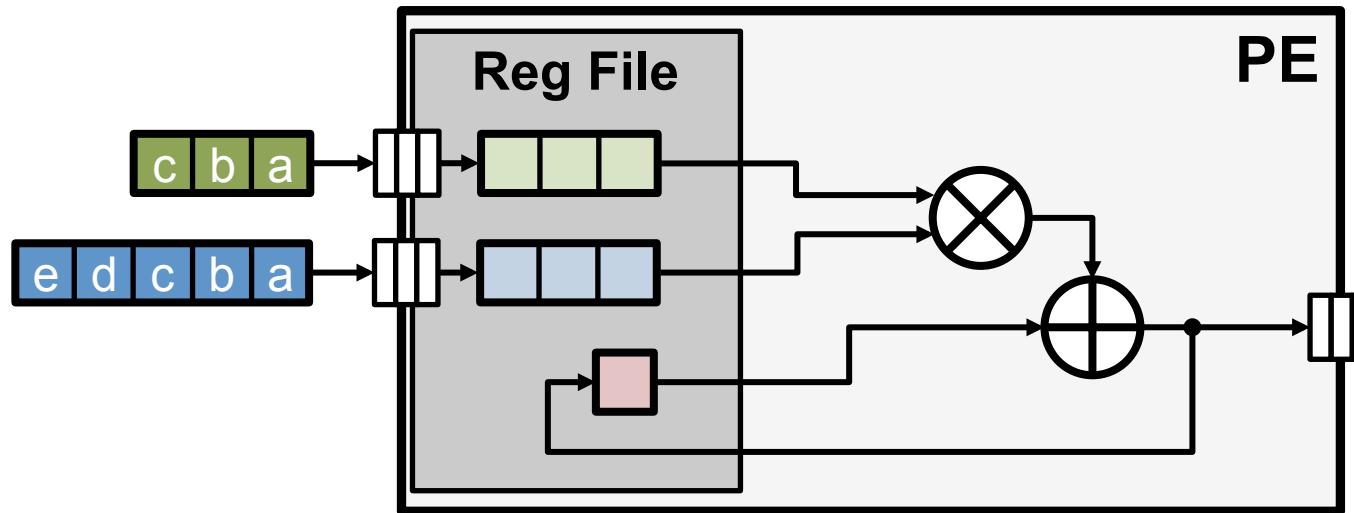
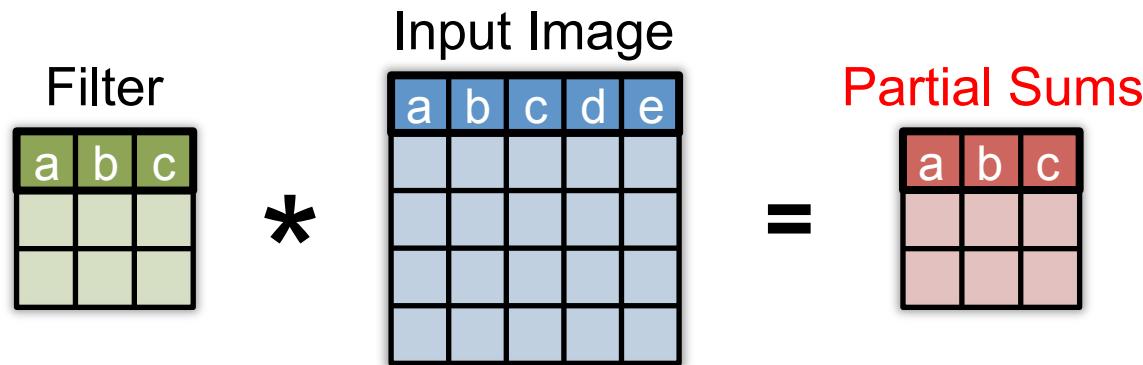


- Use a **large global buffer** as shared storage
 - Reduce **DRAM** access energy consumption
- **Examples:**
 - [DianNao, ASPLOS 2014]
 - [DaDianNao, MICRO 2014]
 - [Zhang, FPGA 2015]

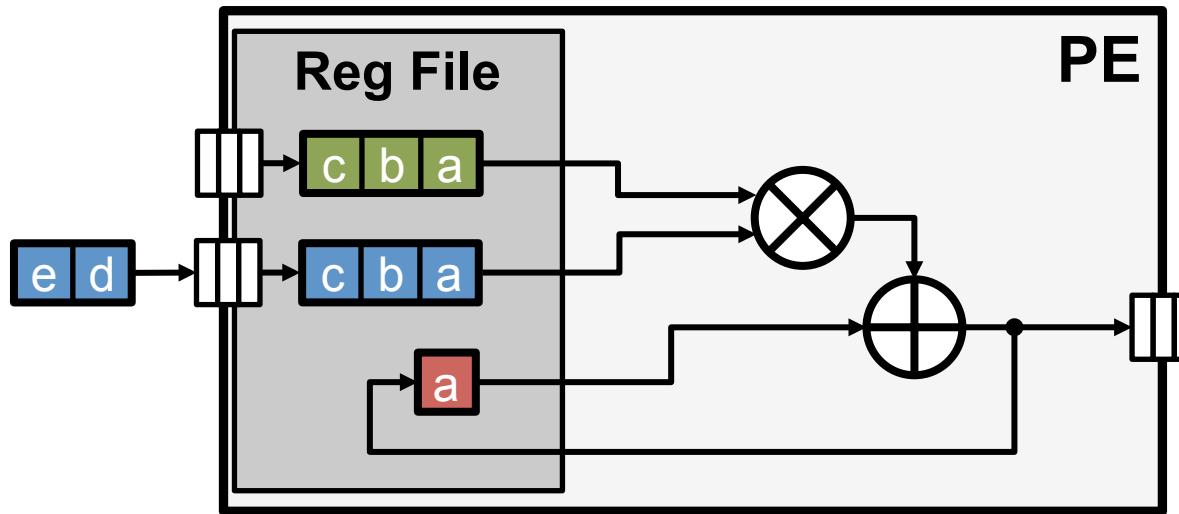
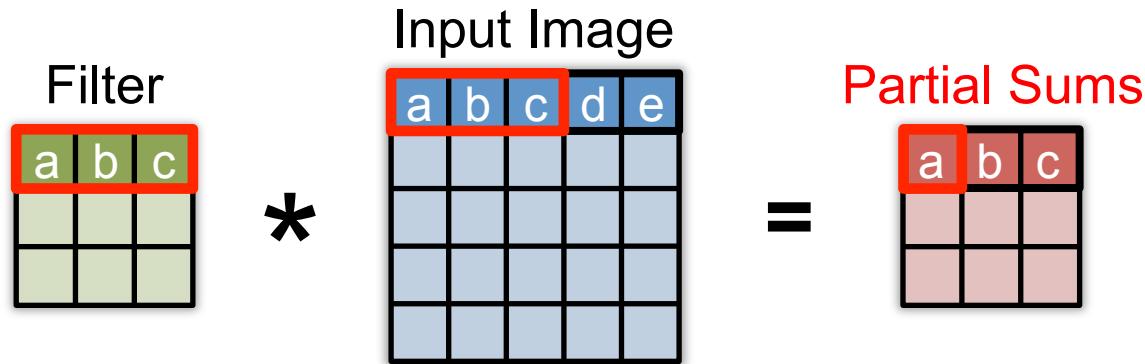
Row Stationary: Energy-efficient Dataflow



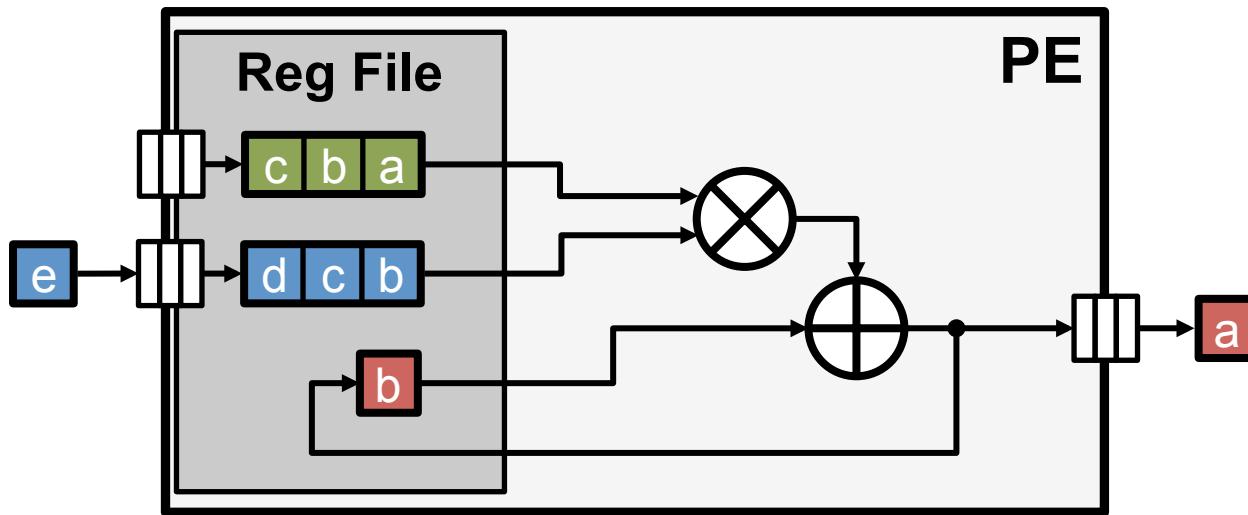
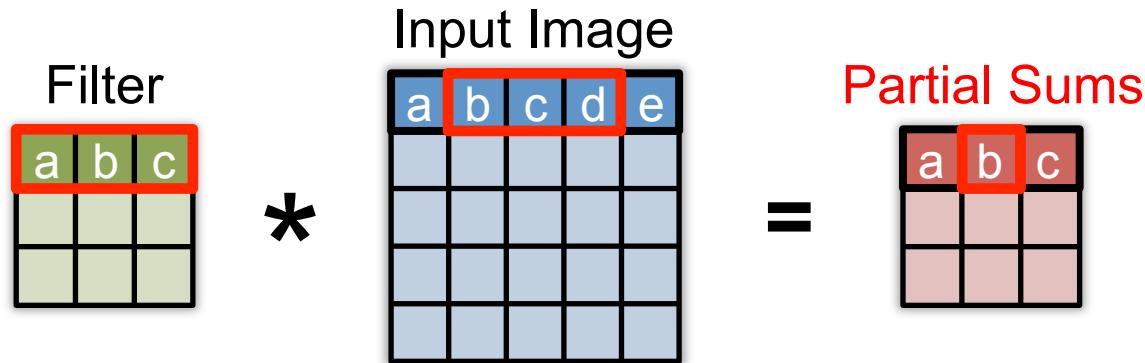
1D Row Convolution in PE



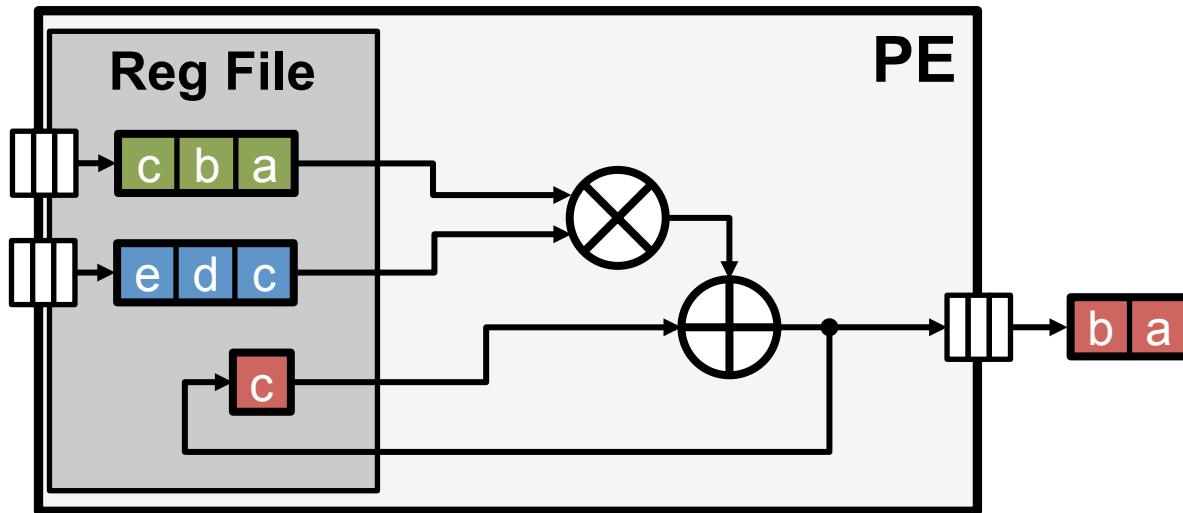
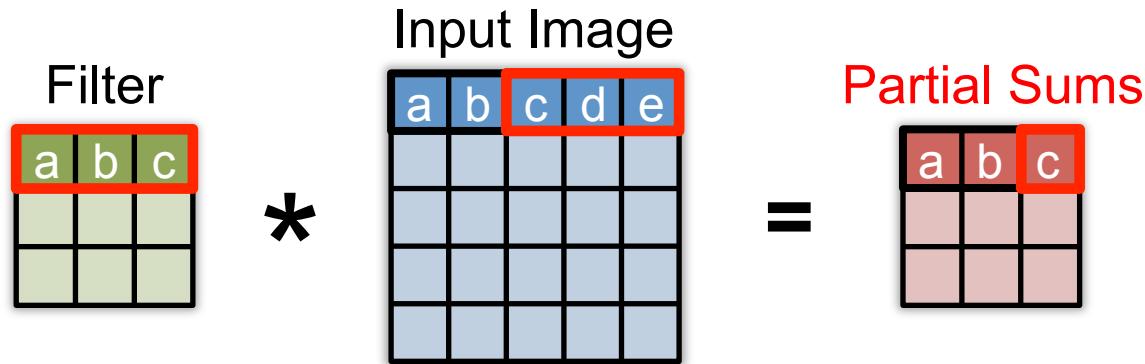
1D Row Convolution in PE



1D Row Convolution in PE

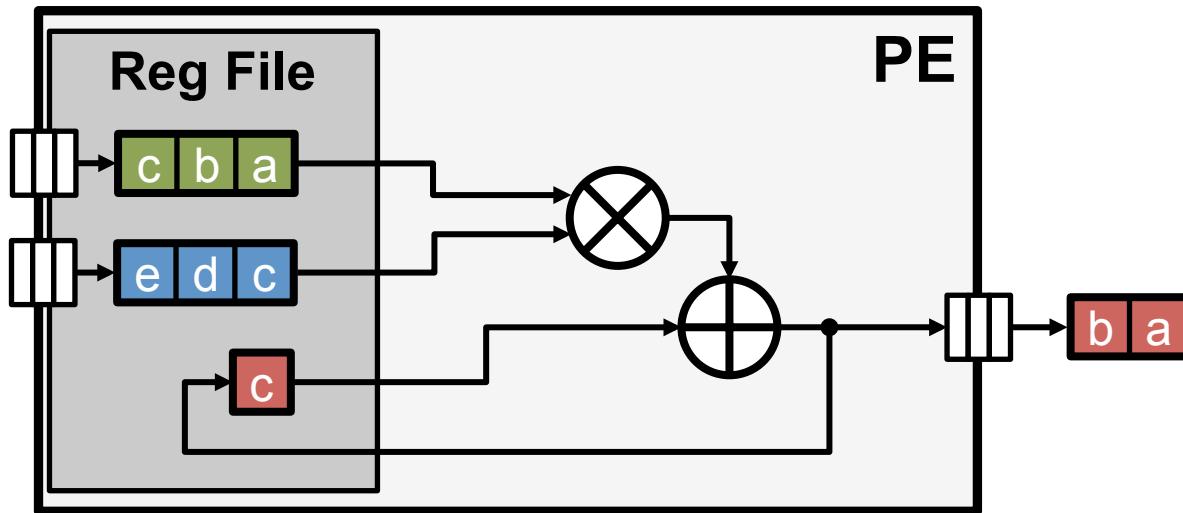


1D Row Convolution in PE



1D Row Convolution in PE

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



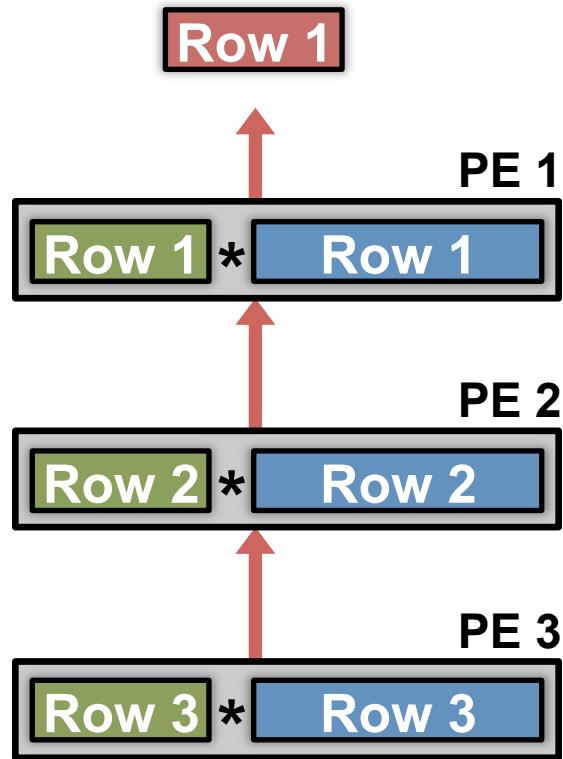
2D Convolution in PE Array



$$\begin{array}{c} \text{Input Row} \\ \begin{matrix} \text{Row 1} \\ \text{Row 2} \end{matrix} \end{array} * \begin{array}{c} \text{Filter} \\ \begin{matrix} \text{Filter Rows} \\ \begin{matrix} \text{Row 1} & \text{Row 2} \\ \text{Row 3} & \text{Row 4} \end{matrix} \end{matrix} \end{array} = \begin{array}{c} \text{Output Row} \\ \begin{matrix} \text{Row 1} \\ \text{Row 2} \end{matrix} \end{array}$$

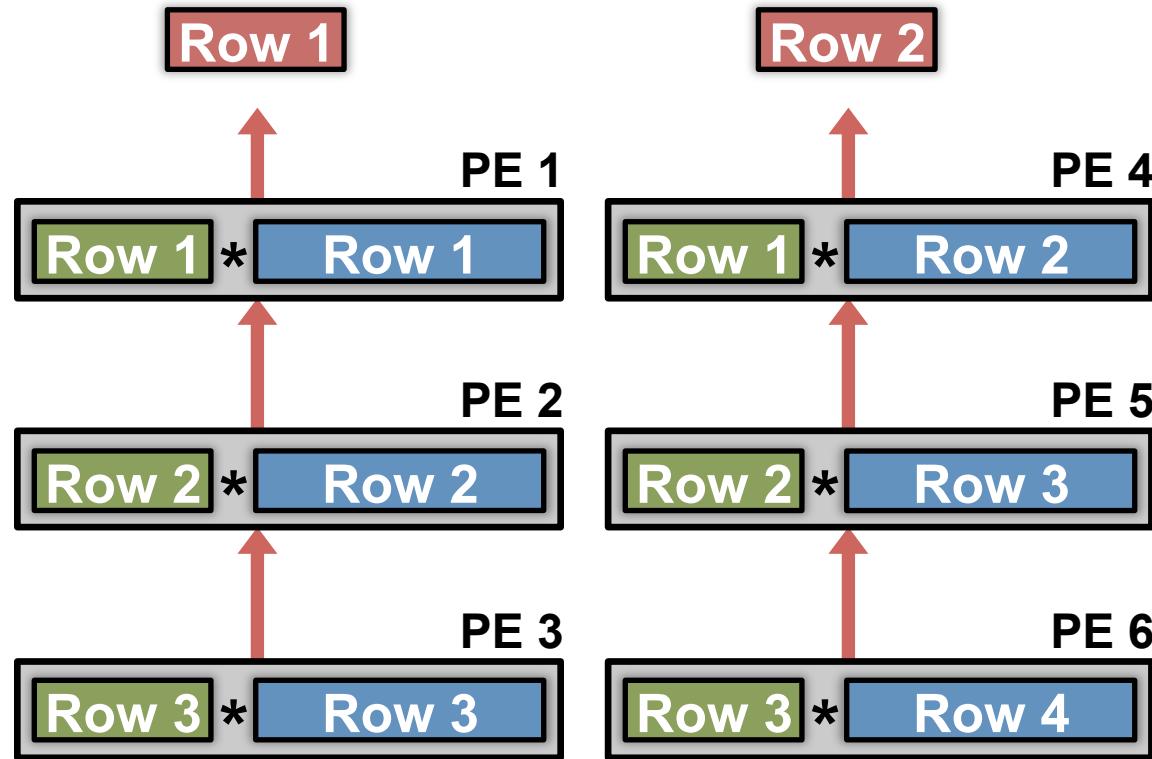
The diagram illustrates a 2D convolution operation. On the left, a green row labeled "Row 1" and a blue row labeled "Row 2" are multiplied by a filter represented as a 2x2 grid of blue squares. The result is an output row consisting of two red squares.

2D Convolution in PE Array



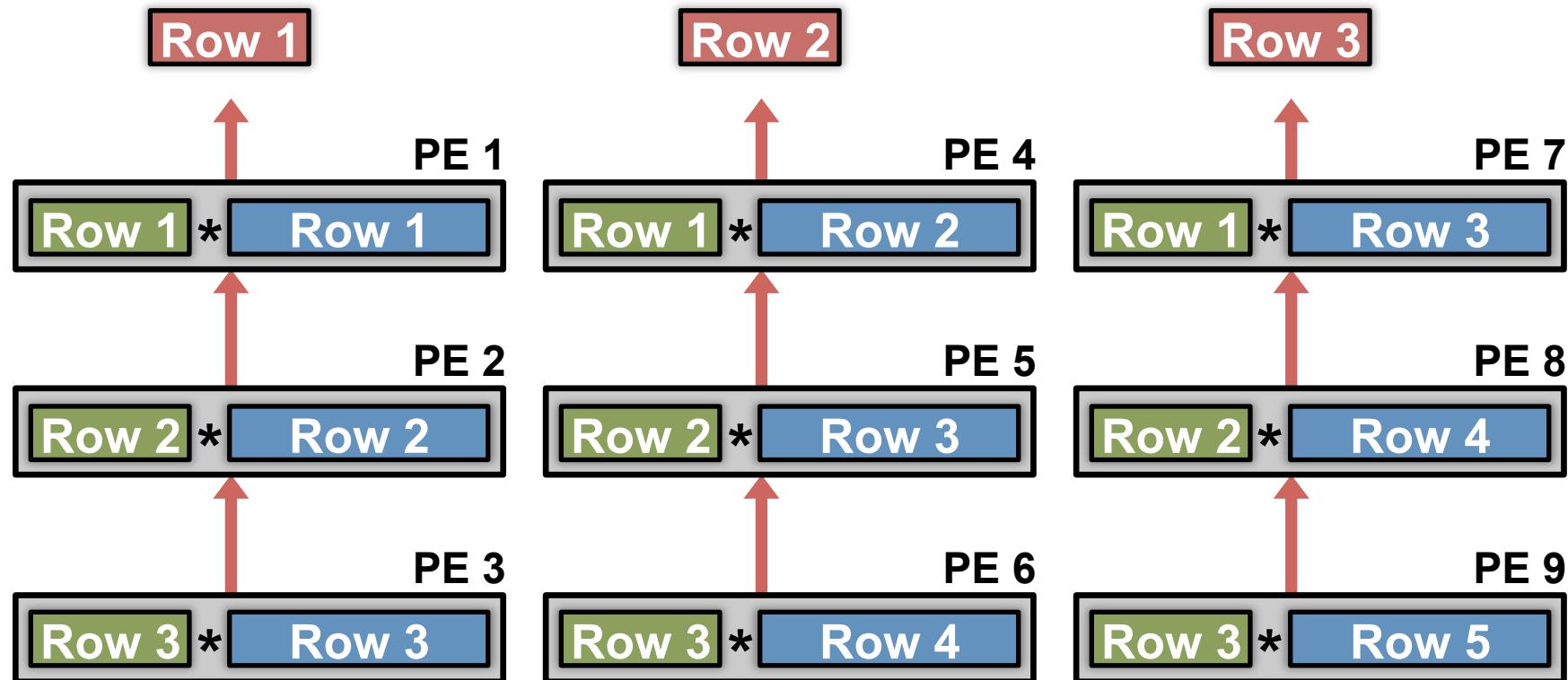
A diagram illustrating the result of a 3x3 convolution operation. It shows a green 3x3 input matrix multiplied by a blue 3x3 weight matrix, resulting in a red 3x3 output matrix. The multiplication is indicated by an asterisk (*) and the equals sign (=).

2D Convolution in PE Array



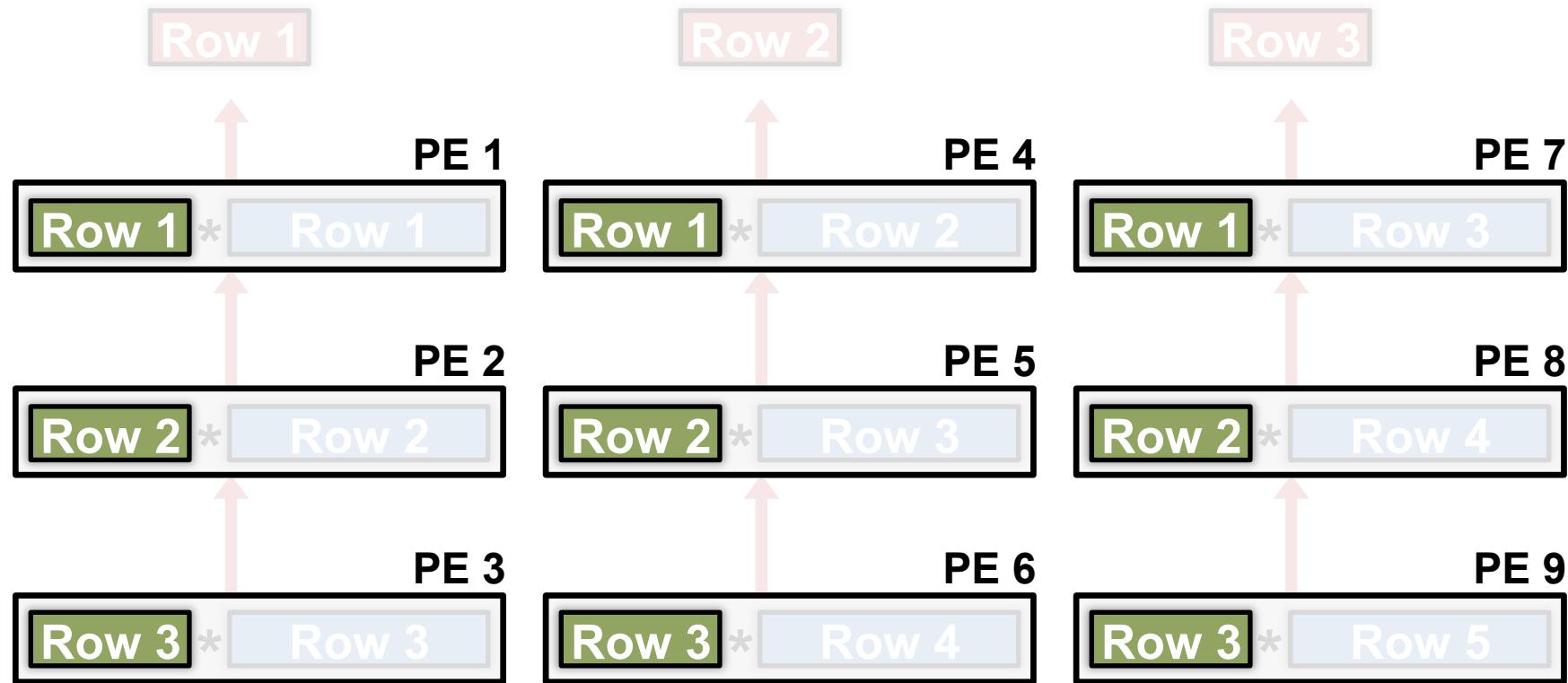
$$\begin{array}{c} \text{green 3x3 grid} \\ * \end{array} = \begin{array}{c} \text{red 3x3 grid} \end{array}$$
$$\begin{array}{c} \text{green 3x3 grid} \\ * \end{array} = \begin{array}{c} \text{red 3x3 grid} \end{array}$$

2D Convolution in PE Array



$$\begin{array}{c} \text{grid} \\ \times \\ \text{grid} \end{array} = \text{grid} \quad
 \begin{array}{c} \text{grid} \\ \times \\ \text{grid} \end{array} = \text{grid} \quad
 \begin{array}{c} \text{grid} \\ \times \\ \text{grid} \end{array} = \text{grid}$$

Convolutional Reuse Maximized



Filter rows are reused across PEs **horizontally**

Convolutional Reuse Maximized

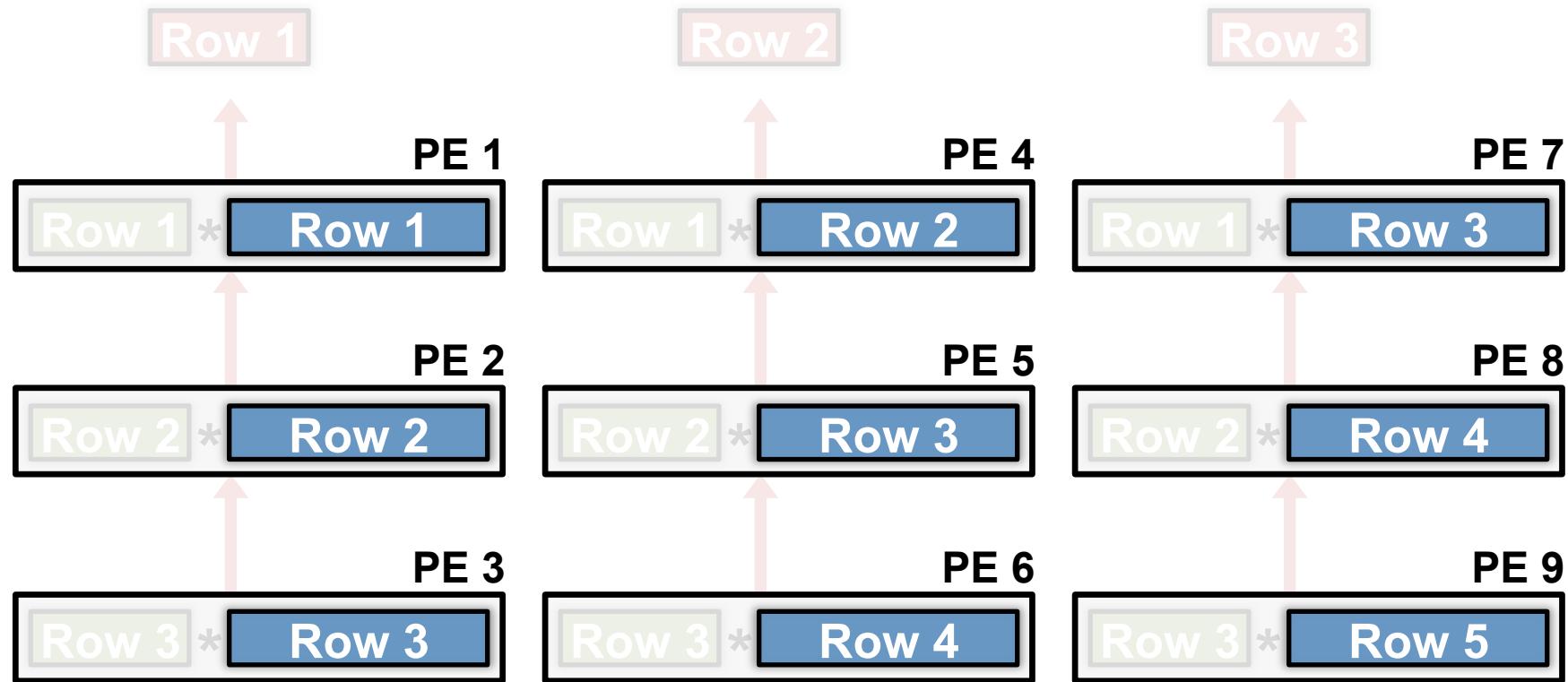
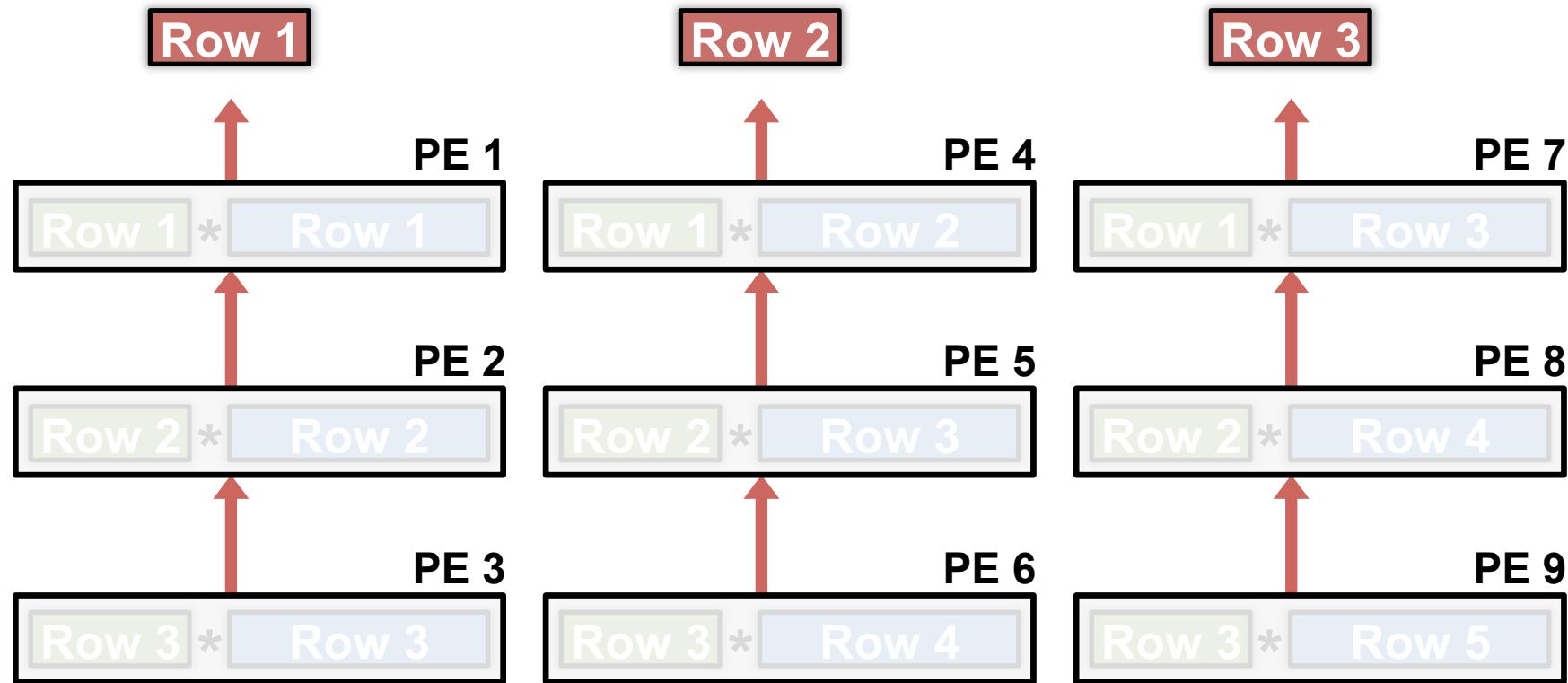


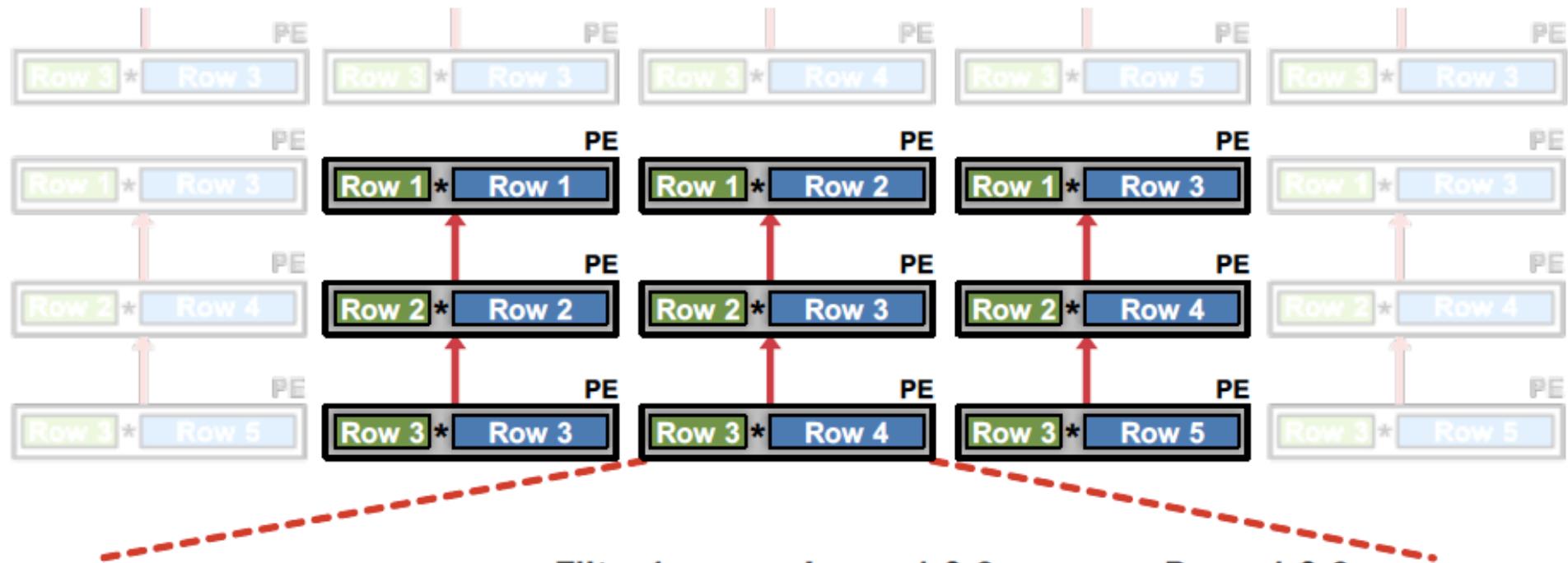
Image rows are reused across PEs **diagonally**

Maximize 2D Accumulation in PE Array



Partial sums accumulate across PEs **vertically**

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

Evaluate Reuse in Different Dataflows

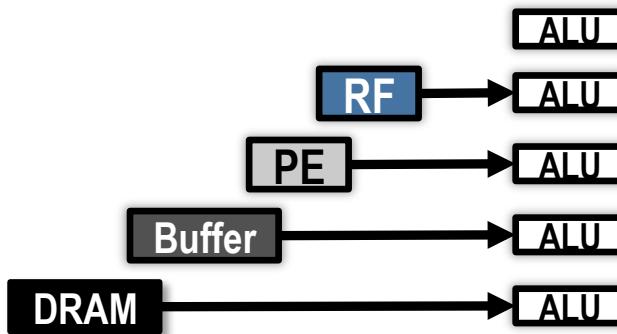
- **Weight Stationary**
 - Minimize movement of filter weights
- **Output Stationary**
 - Minimize movement of partial sums
- **No Local Reuse**
 - Don't use any local PE storage. Maximize global buffer size.
- **Row Stationary**

Evaluate Reuse in Different Dataflows

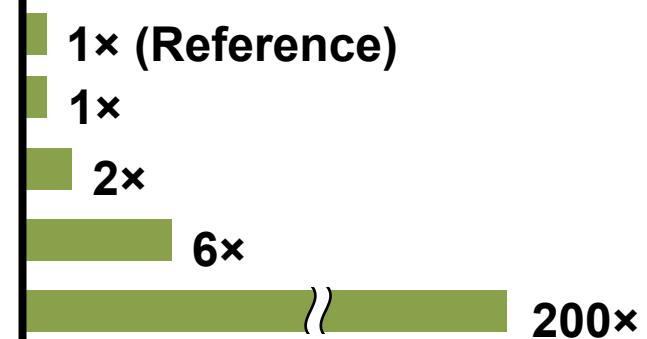
- **Weight Stationary**
 - Minimize movement of filter weights
- **Output Stationary**
 - Minimize movement of partial sums
- **No Local Reuse**
 - Don't use any local PE storage. Maximize global buffer size.
- **Row Stationary**

Evaluation Setup

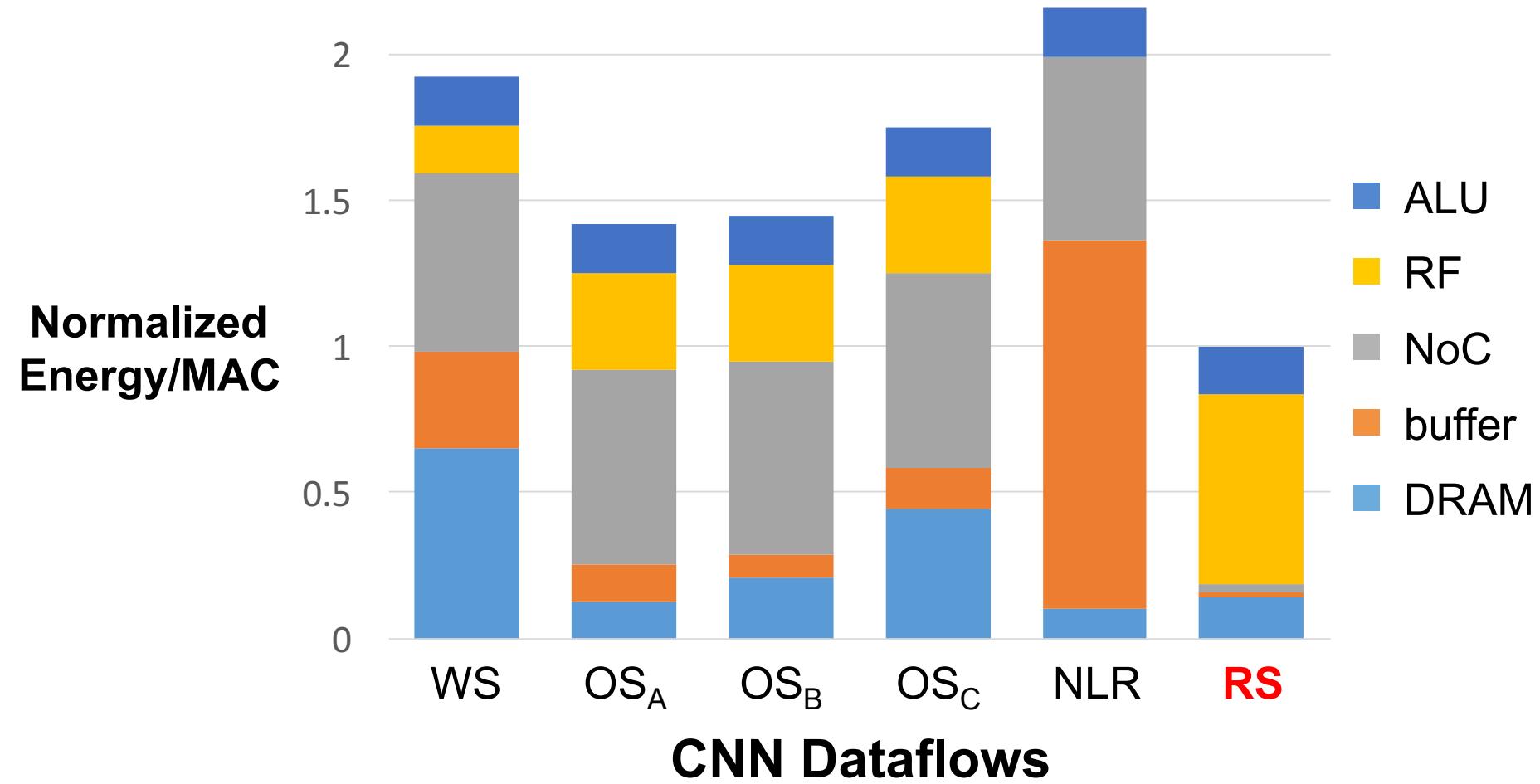
- Same Total Area
- AlexNet
- 256 PEs
- Batch size = 16



Normalized Energy Cost*

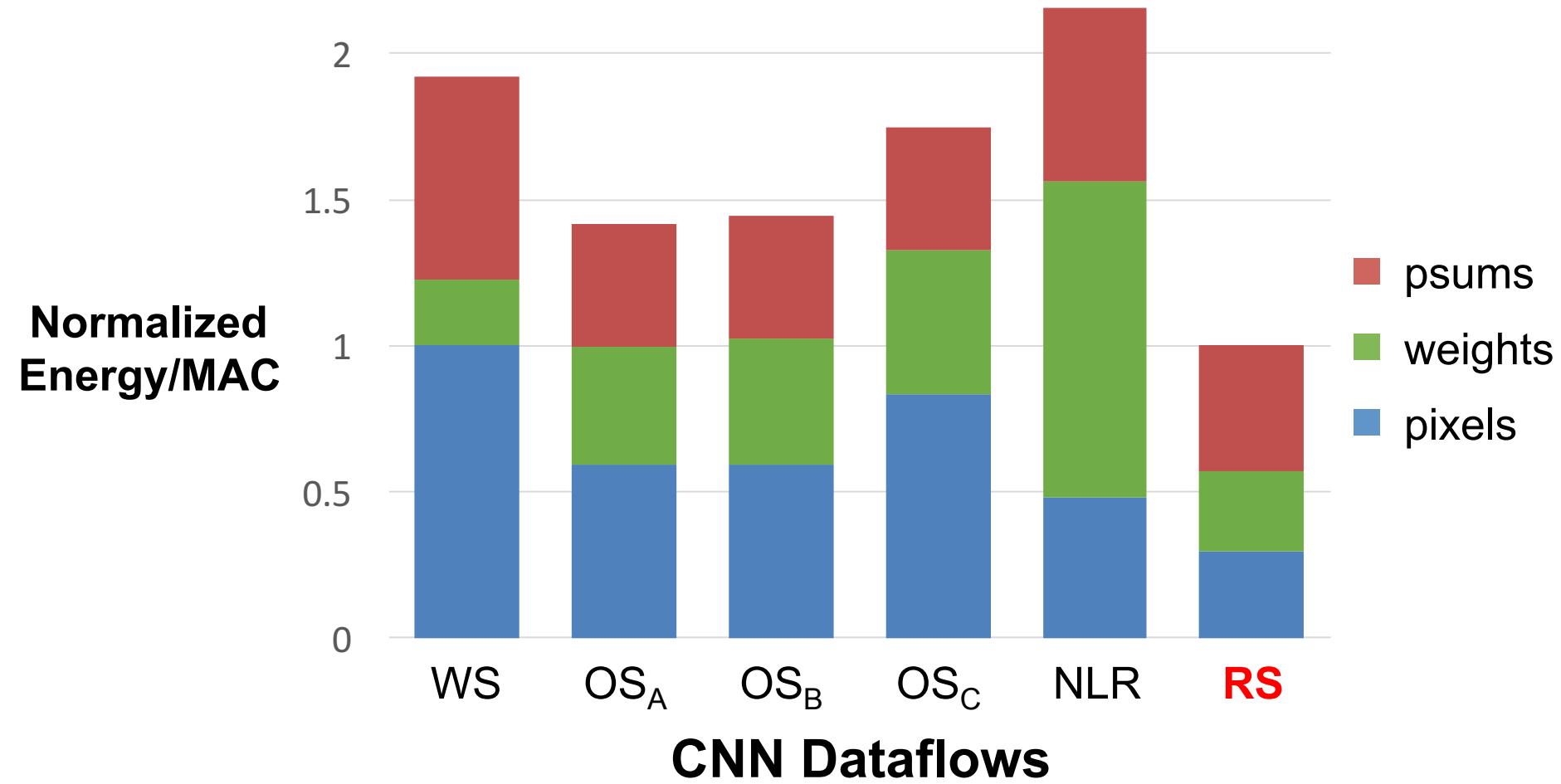


Dataflow Comparison: CONV Layers



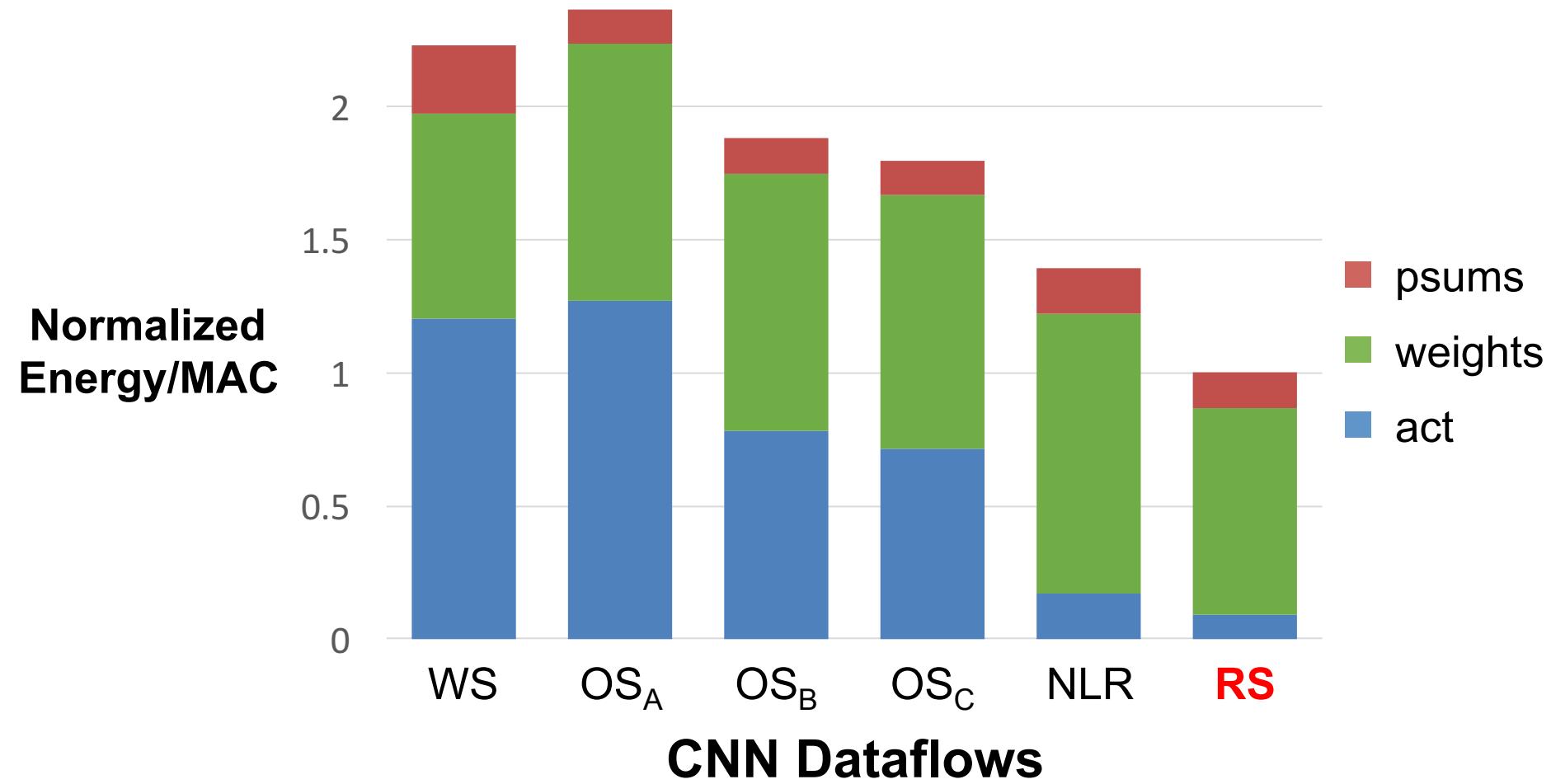
RS uses 1.4x – 2.5x lower energy than other dataflows

Dataflow Comparison: CONV Layers



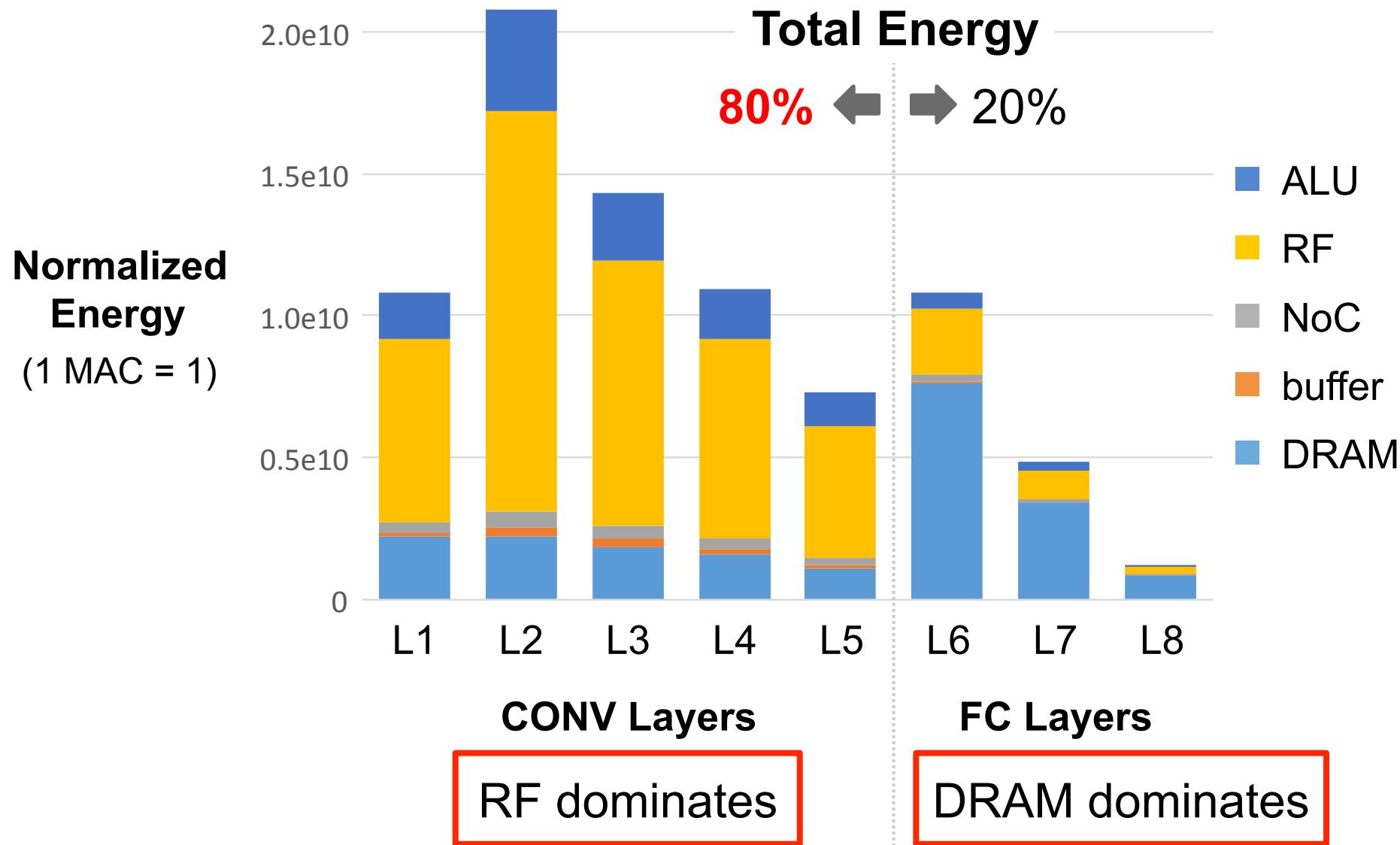
RS optimizes for the best **overall** energy efficiency

Dataflow Comparison: FC Layers

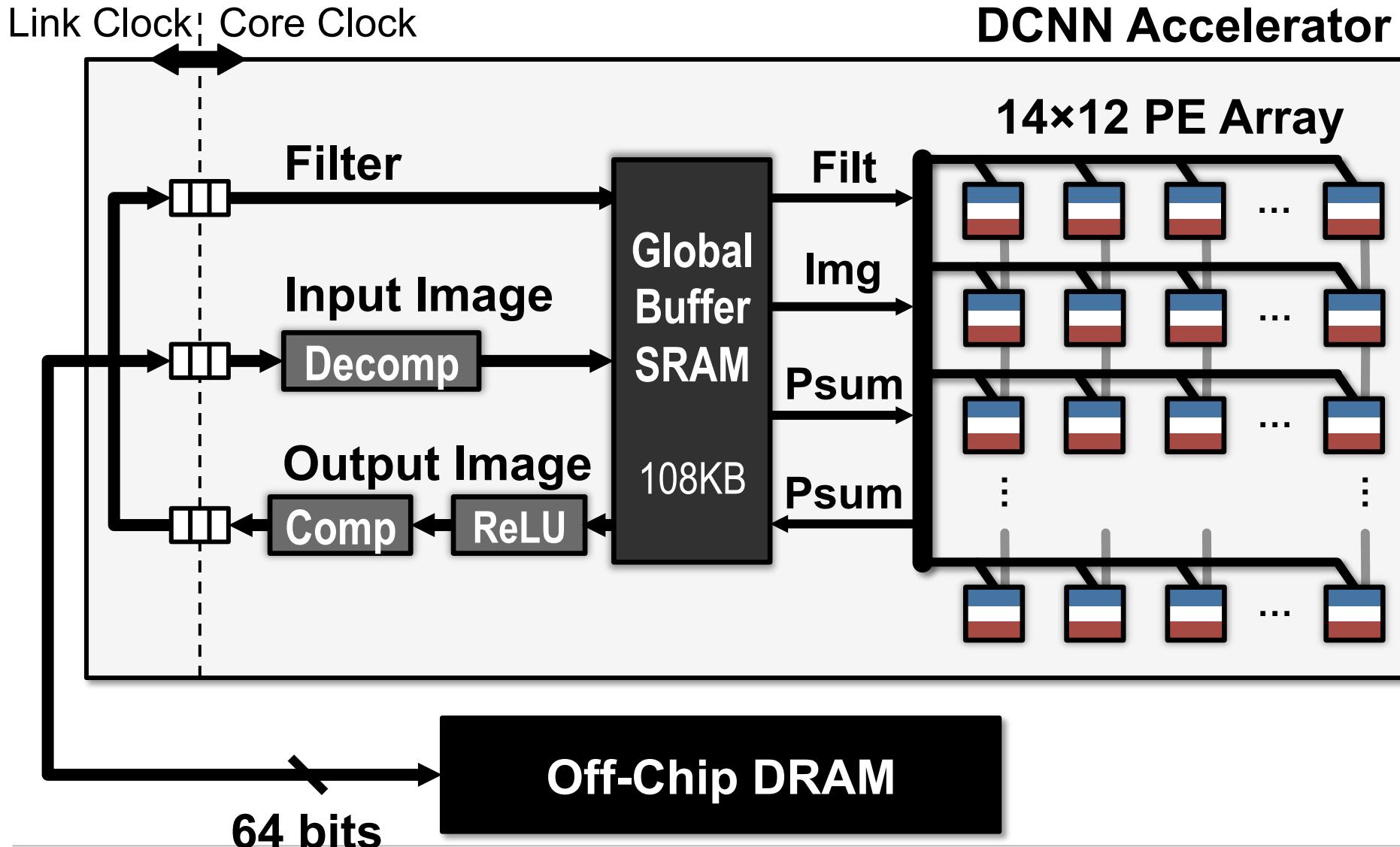


RS uses at least 1.3× lower energy than other dataflows

Row Stationary: Layer Breakdown



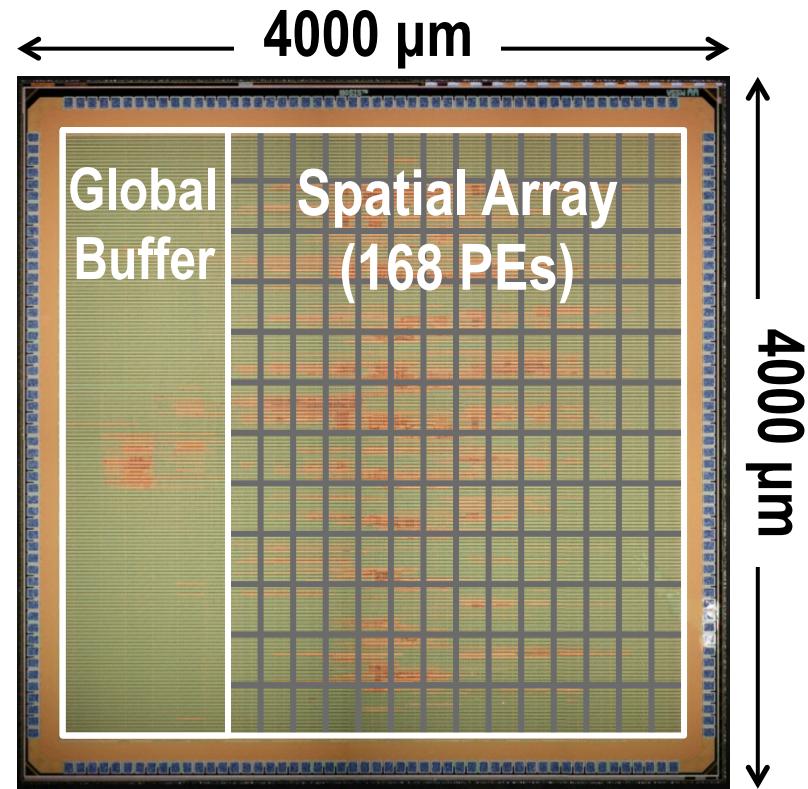
Eyeriss Deep CNN Accelerator



[Chen et al., ISSCC 2016]

Eyeriss Chip Spec & Measurement Results

Technology	TSMC 65nm LP 1P9M
On-Chip Buffer	108 KB
# of PEs	168
Scratch Pad / PE	0.5 KB
Core Frequency	100 – 250 MHz
Peak Performance	33.6 – 84.0 GOPS
Word Bit-width	16-bit Fixed-Point
Natively Supported CNN Shapes	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>	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
Throughput¹	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

The screenshot shows a web browser window titled "MIT Eyeriss Caffe Demo" at the URL "18.62.18.33:5000". The browser interface includes a toolbar with various icons and a sidebar with a vertical stack of application icons.

The main content of the page is titled "[ISSCC 2016] Paper 14.5: Eyeriss Caffe Demo". It features two main sections:

- 1. System Setup:** A photograph of a Jetson TK1 board and a VC707 board connected via a PCIe cable. The VC707 board has a red Eyeriss chip mounted on it. The text "Jetson TK1 VC707 + Eyeriss" is displayed below the image.
- 2. Eyeriss Die Photo:** A detailed photograph of the Eyeriss chip die. The die is square with a grid of processing elements (PEs). Labels indicate the "On-chip Buffer" and the "Spatial PE Array". Dimensions of 4000 μm by 4000 μm are shown.

Below these sections is a "Classification" form:

- A text input field labeled "Provide an image URL..."
- A file upload input field labeled "Or upload an image: Choose File No file chosen"
- A blue "Classify URL" button

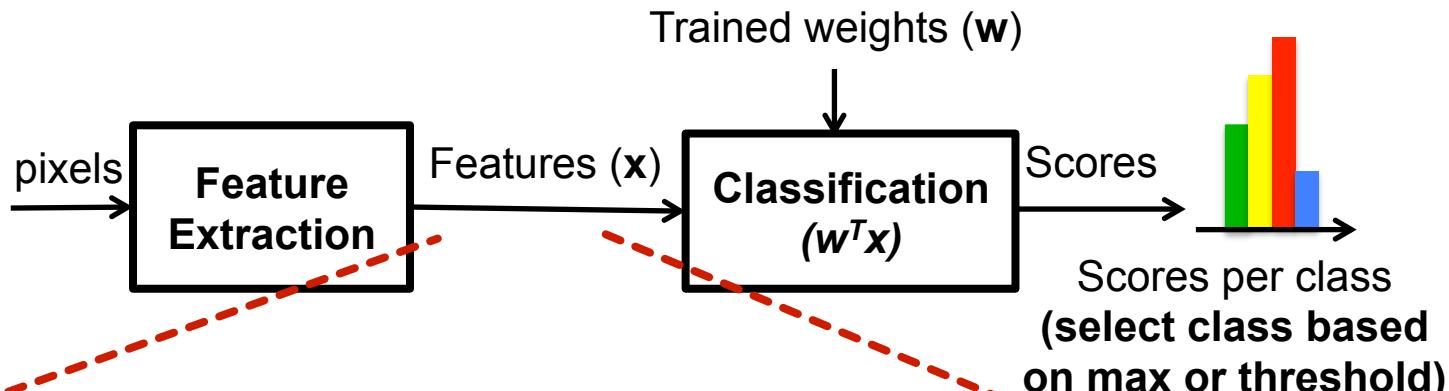
<https://vimeo.com/154012013>

Integrated with BVLC Caffe DL Framework

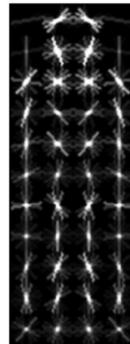
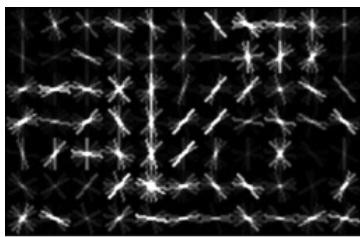
Hand-Crafted vs. Learned Features

Machine Learning Pipeline (Inference)

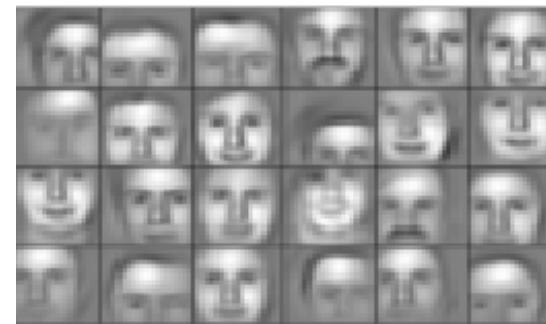
Image



Handcrafted Features
(e.g. HOG)

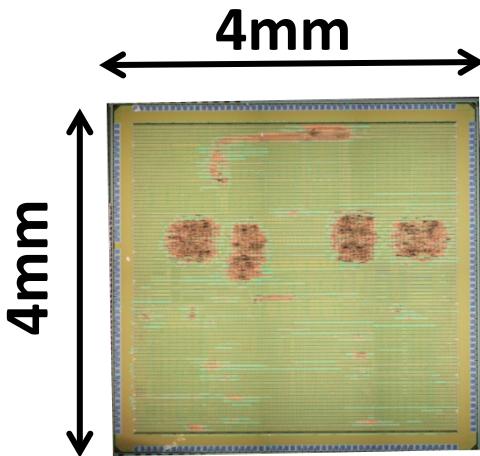


Learned Features
(e.g. DNN)

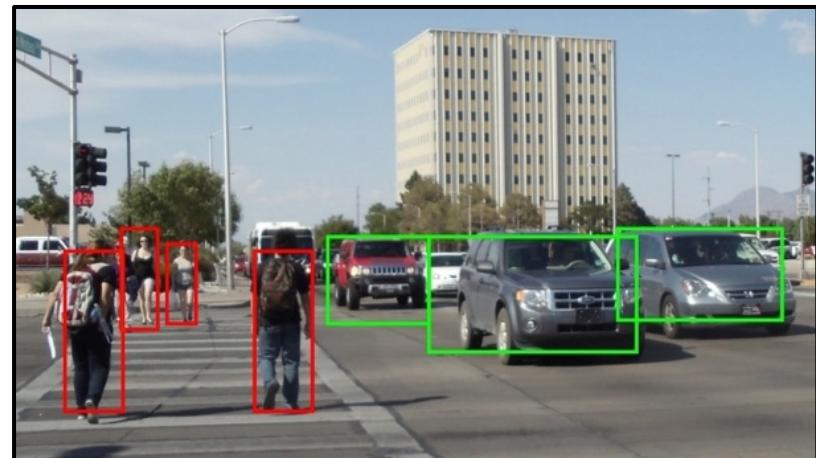


$$\text{Score} = \sum_n x_i w_i$$

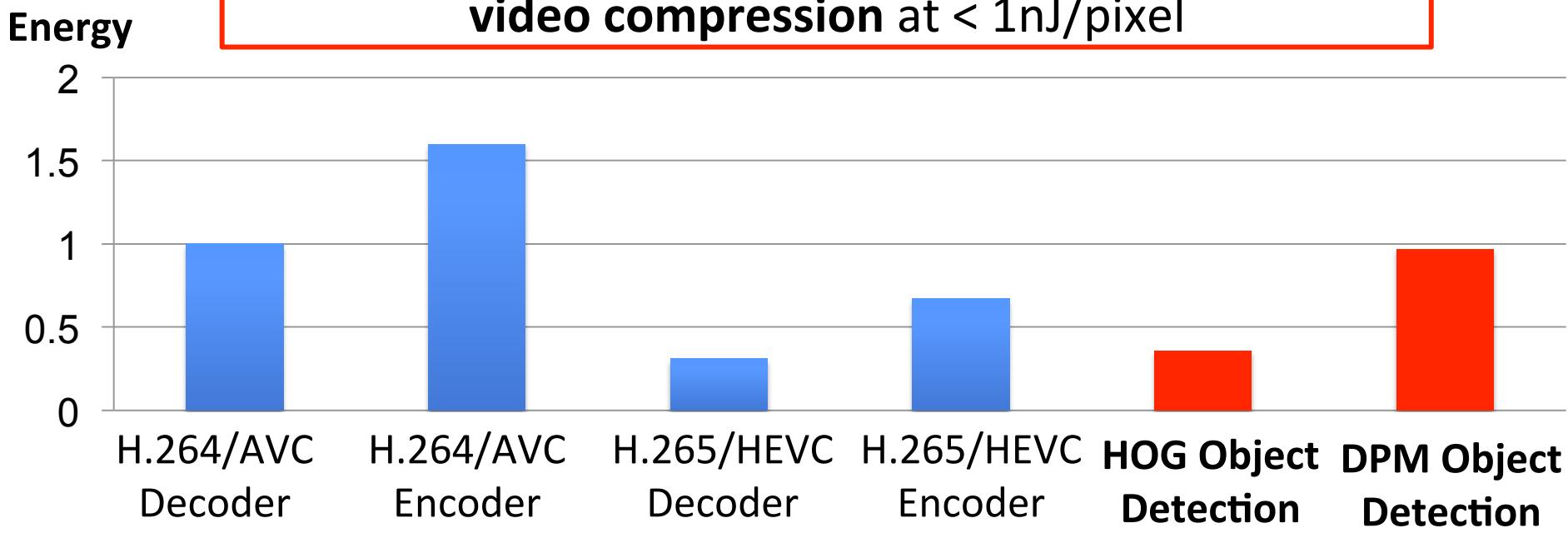
Energy-Efficient Object Detection



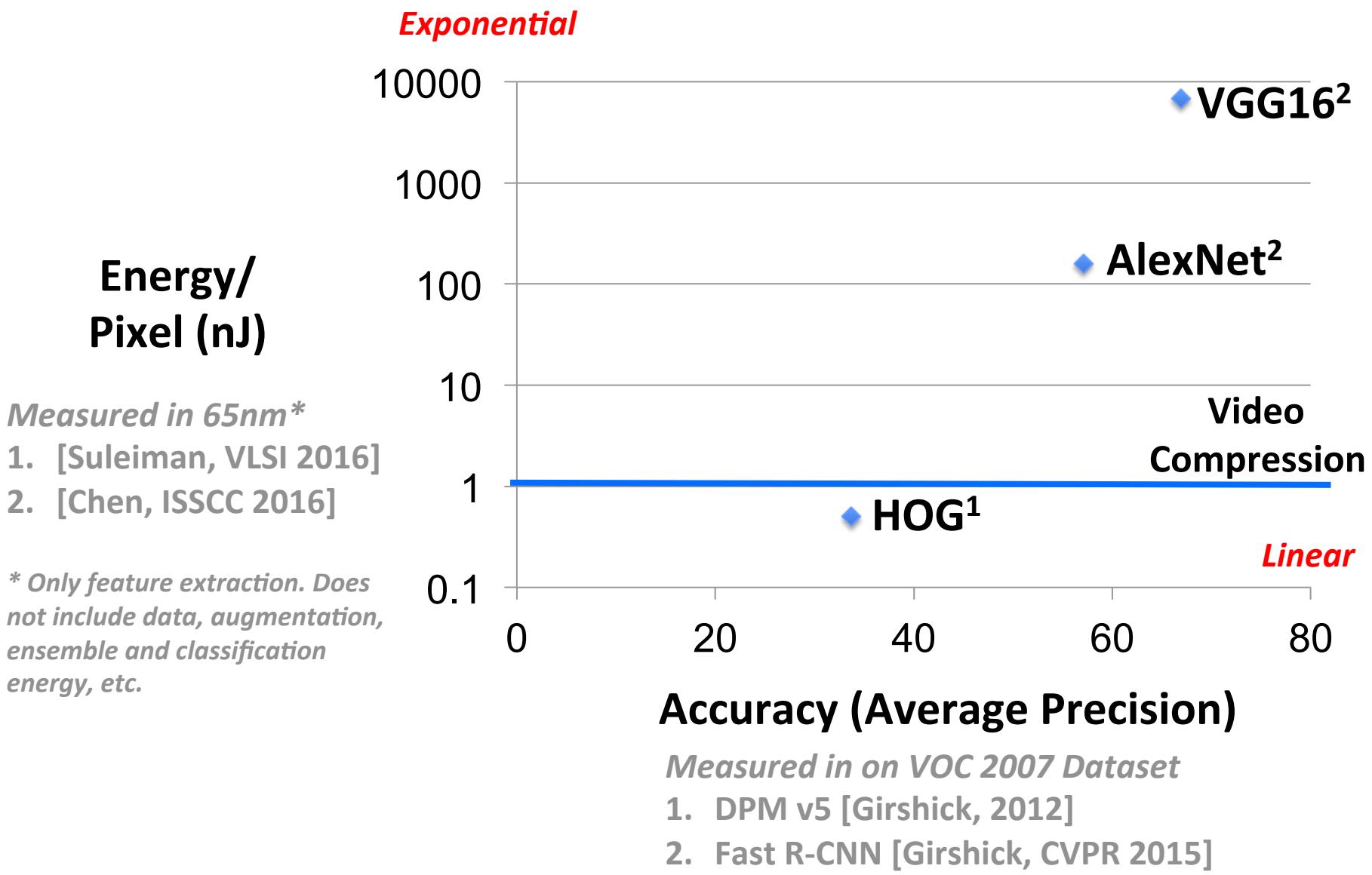
MIT Object
Detection Chip
[VLSI 2016]



Enable object detection to be as **energy-efficient** as
video compression at < 1nJ/pixel



Features: Energy vs. Accuracy

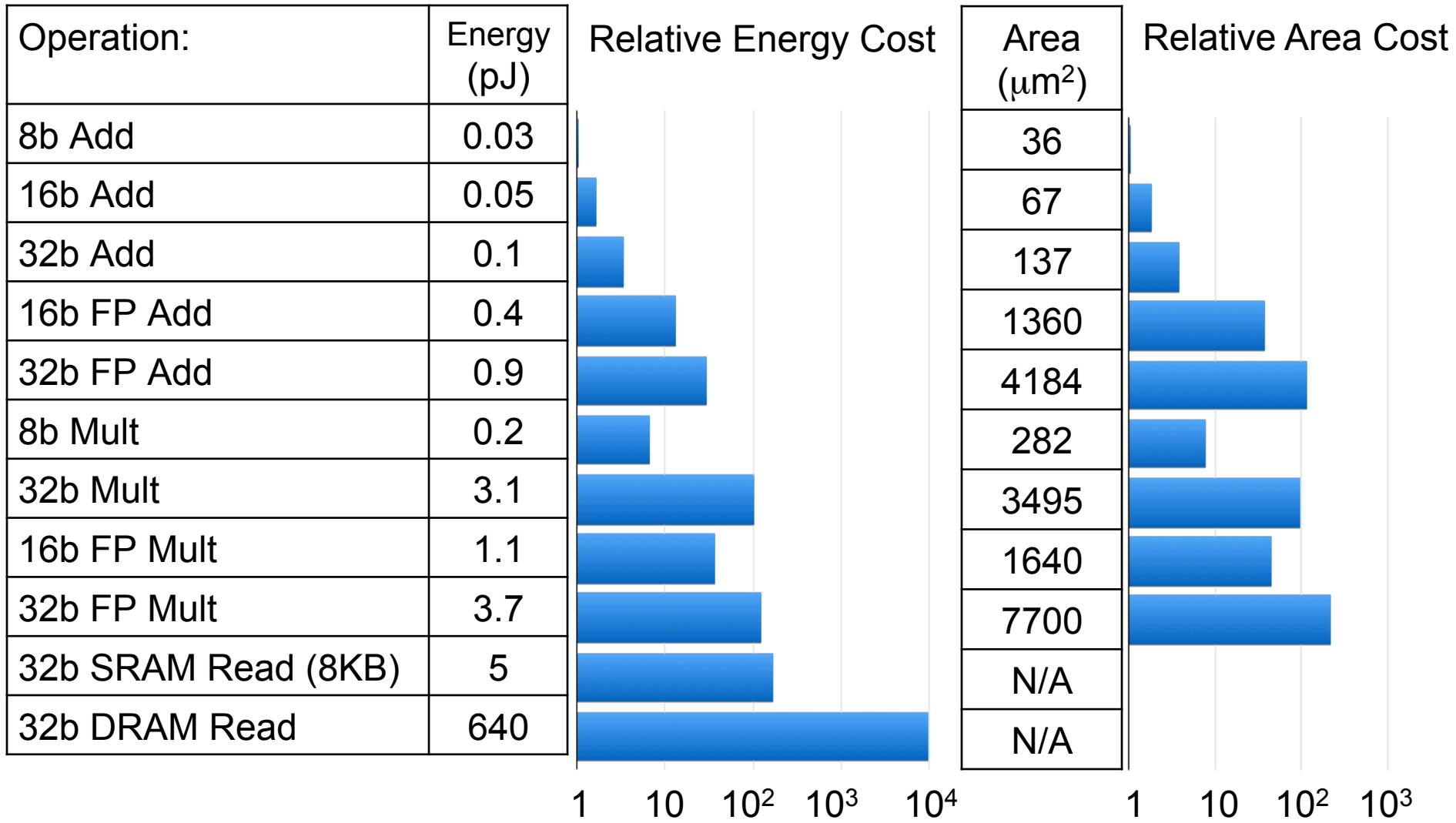


Opportunities in Joint Algorithm Hardware Design

Approaches

- Reduce size of operands for storage/compute
 - Floating point -> Fixed point
 - Bit-width reduction
 - Non-linear quantization
- Reduce number of operations for storage/compute
 - Exploit Activation Statistics (Compression)
 - Network Pruning
 - Compact Network Architectures

Cost of Operations

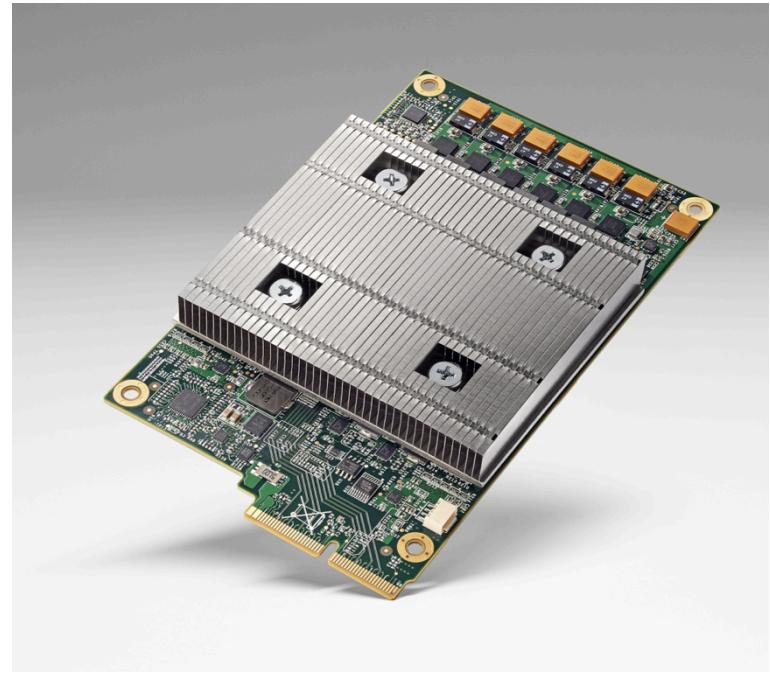


[Horowitz, "Computing's Energy Problem (and what we can do about it)", ISSCC 2014]

Commercial Products using 8-bit Integer



Nvidia's Pascal (2016)

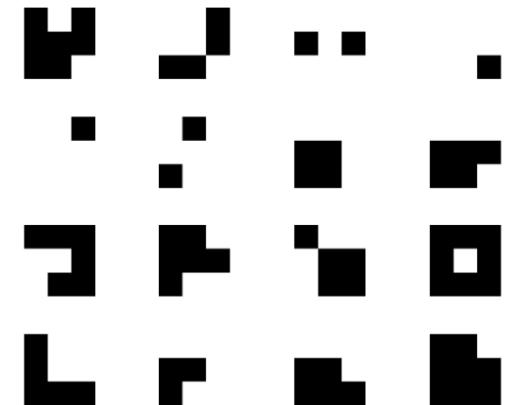


Google's TPU (2016)

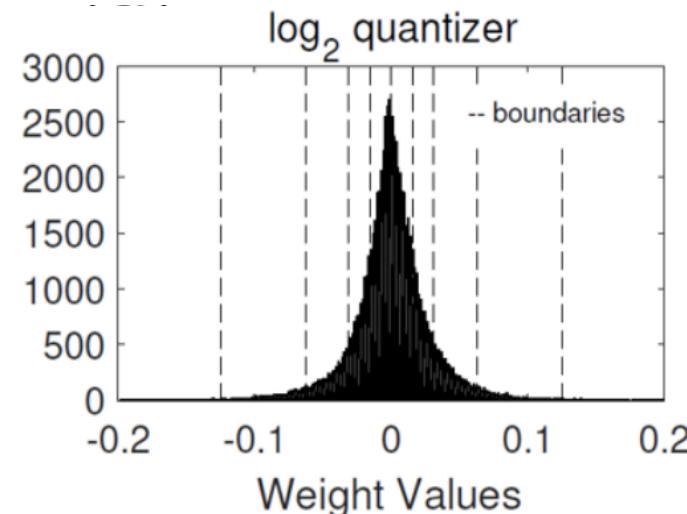
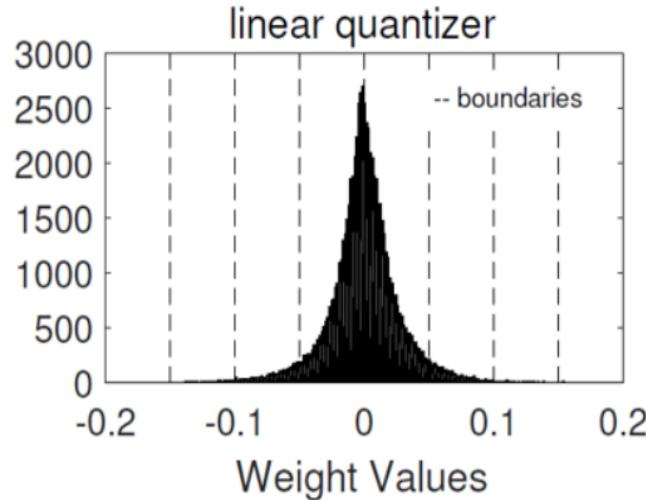
Reduced Precision in Research

- Reduce number of bits
 - Binary Nets [Courbariaux, NIPS 2015]
- Reduce number of unique weights
 - Ternary Weight Nets [Li, arXiv 2016]
 - XNOR-Net [Ratner, ECCV 2016]
- Non-Linear Quantization
 - LogNet [Lee, ICASSP 2017]

Binary Filters

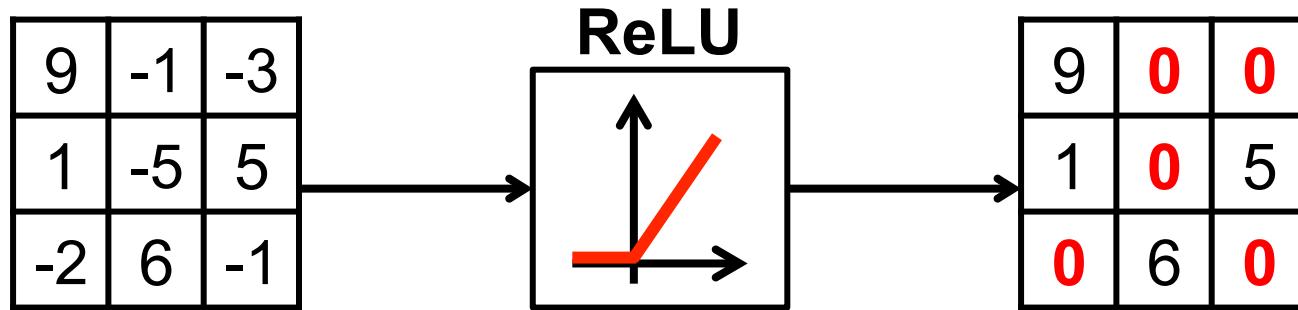


Log Domain Quantization



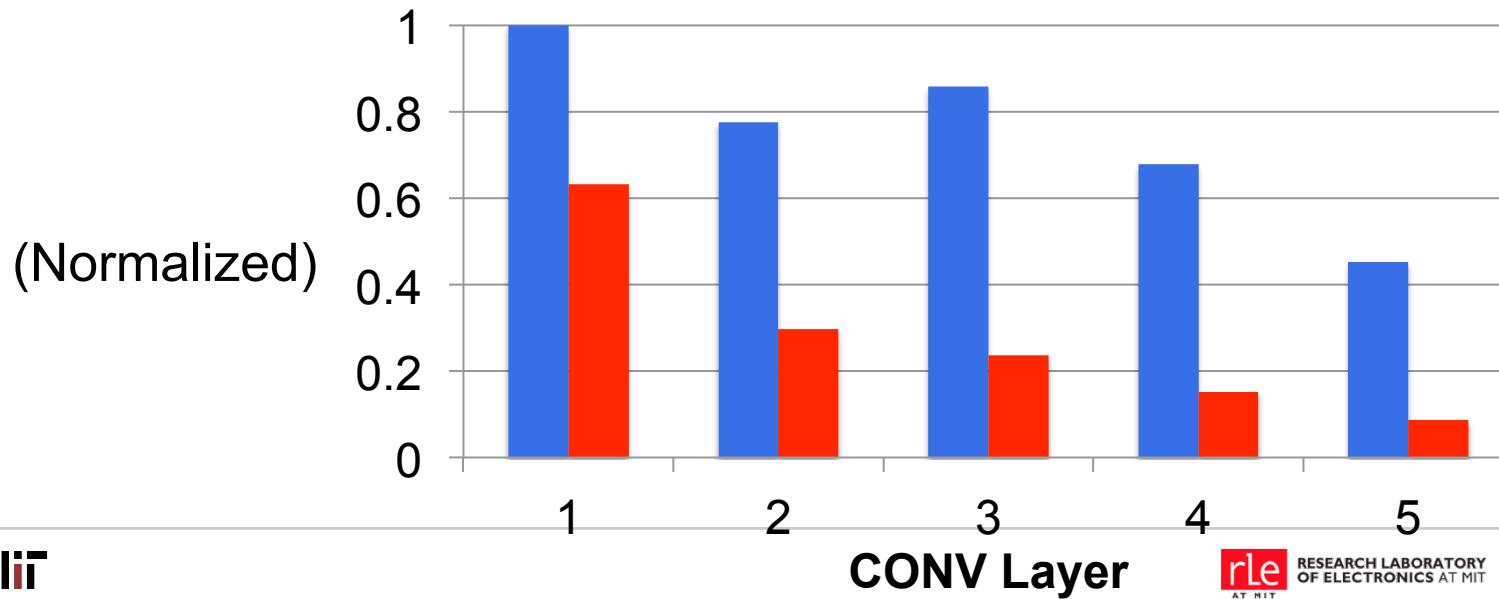
Sparsity in Data

Many **zeros** in output fmaps after ReLU



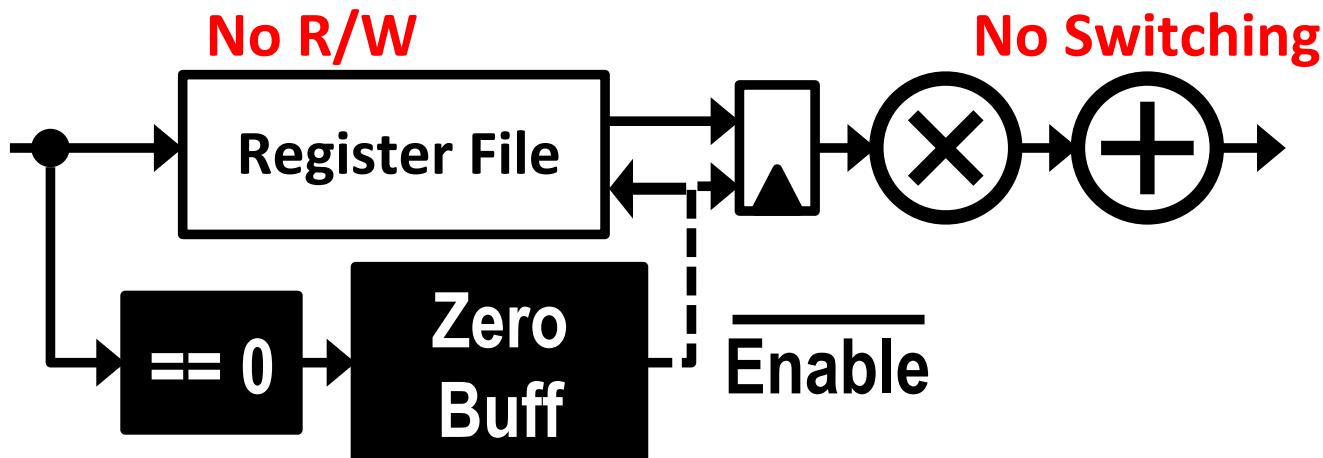
■ # of activations

■ # of non-zero activations

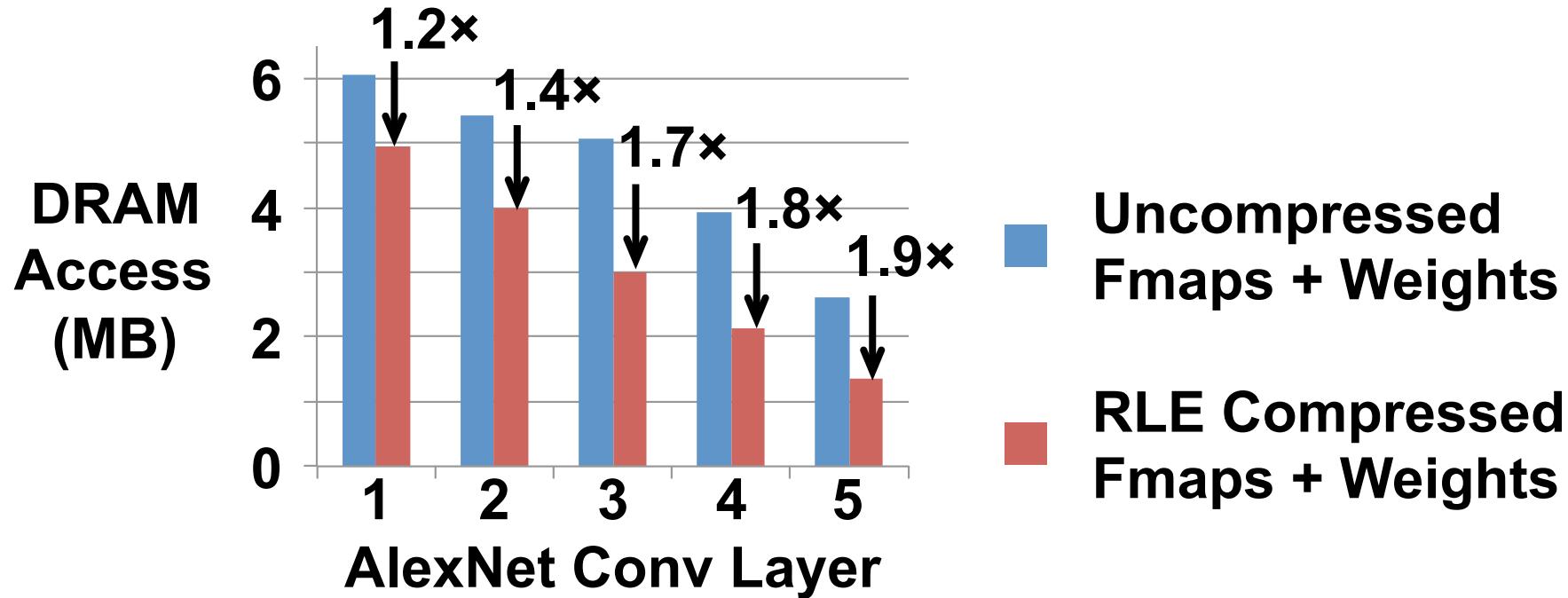


Zero Data Processing Gating

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



Compression Reduces DRAM BW



Simple RLC within 5% - 10% of theoretical entropy limit

Sparsity with Basis Projection

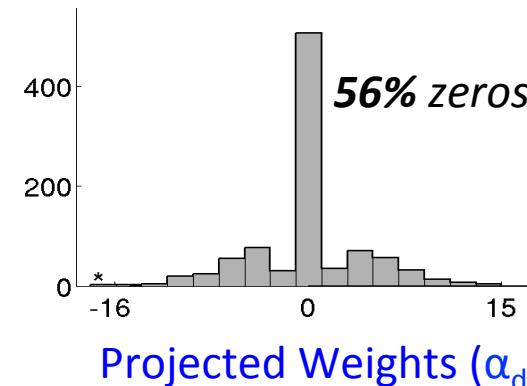
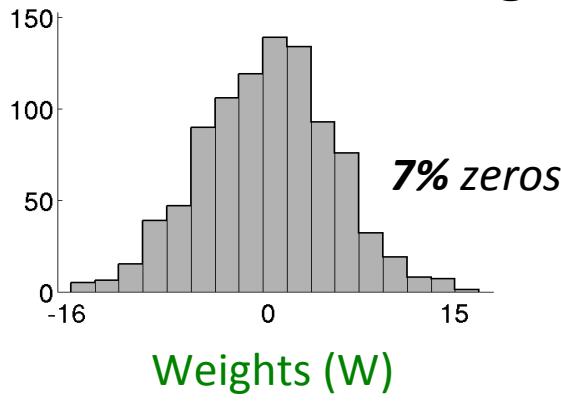
Reduce the number of multiplications by projecting onto a basis that increases sparsity (>1.8x power reduction)

Basis Projection Equation

$$\langle H, W \rangle = \left\langle H, \sum_d S_d \alpha_d \right\rangle = \sum_d \langle H, S_d \rangle \alpha_d = \sum_d P_d \alpha_d$$

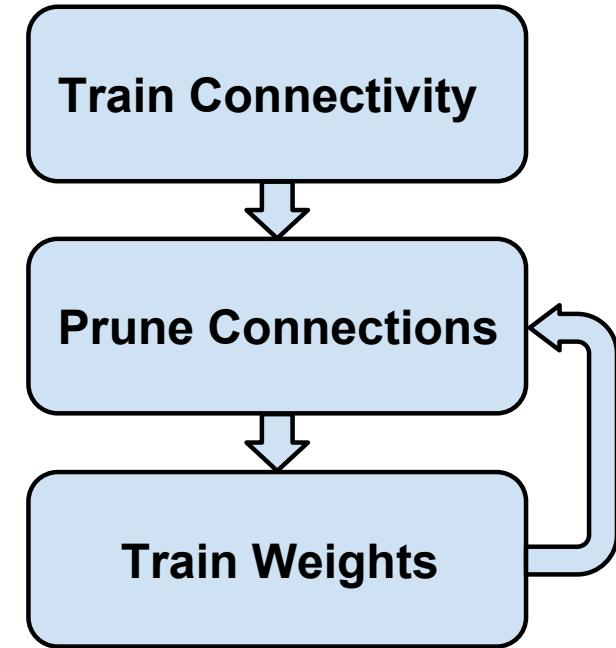
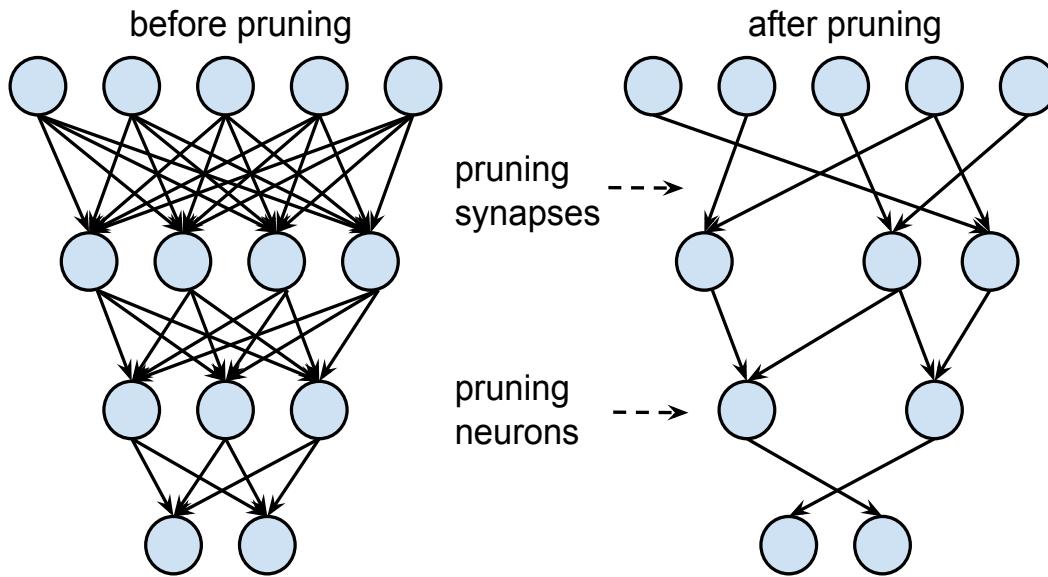
Features Weights Basis Projected Features Projected Weights

Histogram of Weights



Pruning – Make Weights Sparse

Prune based on *magnitude* of weights



Example: AlexNet

Weight Reduction: CONV layers 2.7x, FC layers 9.9x

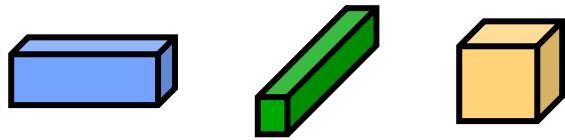
(*Most reduction on fully connected layers*)

Overall: 9x weight reduction, 3x MAC reduction

Key Metrics for Embedded DNN

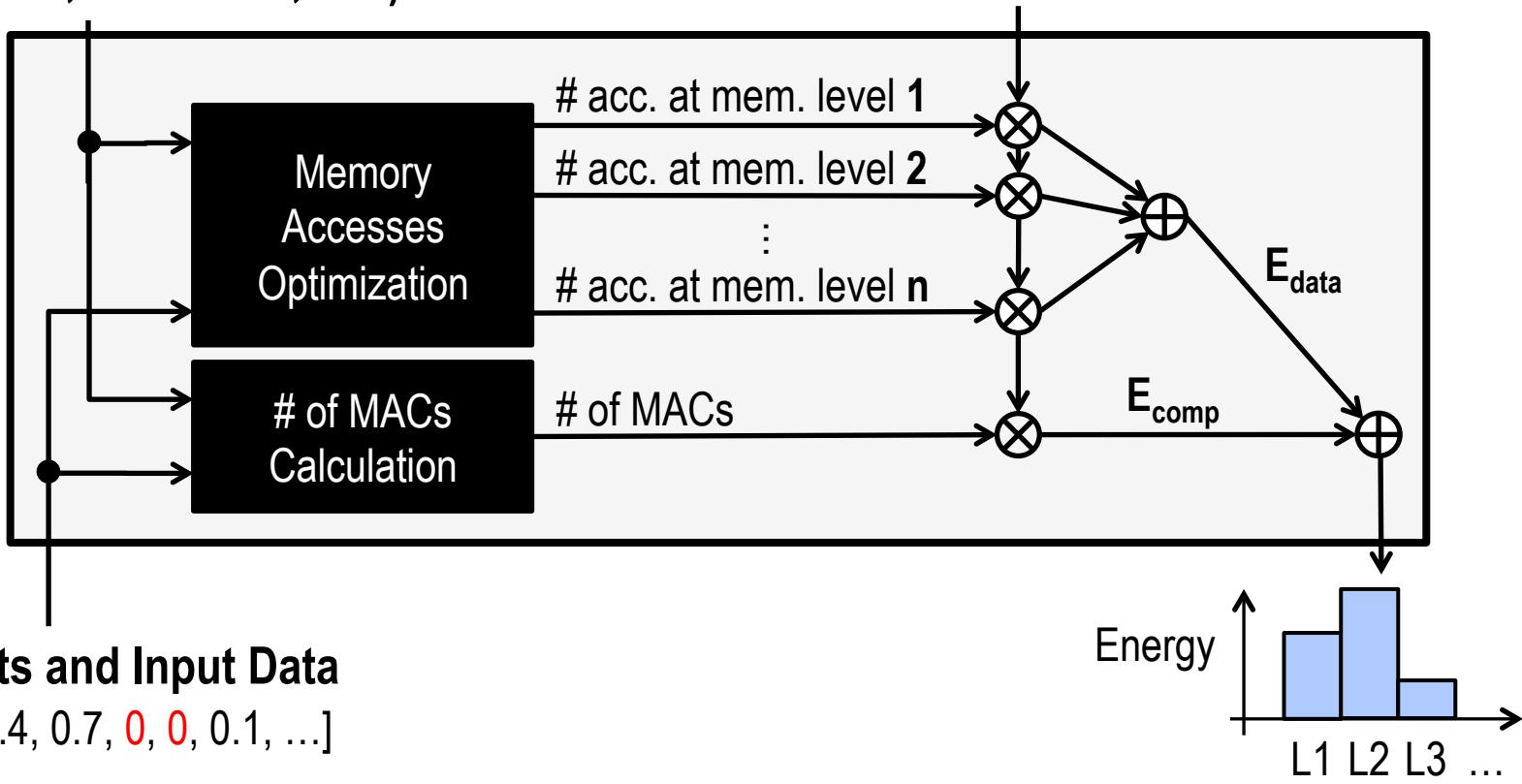
- Accuracy → Measured on Dataset
- Speed → Number of MACs
- Storage Footprint → Number of Weights
- Energy → ?

Energy-Evaluation Methodology



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

Hardware Energy Costs of each MAC and Memory Access



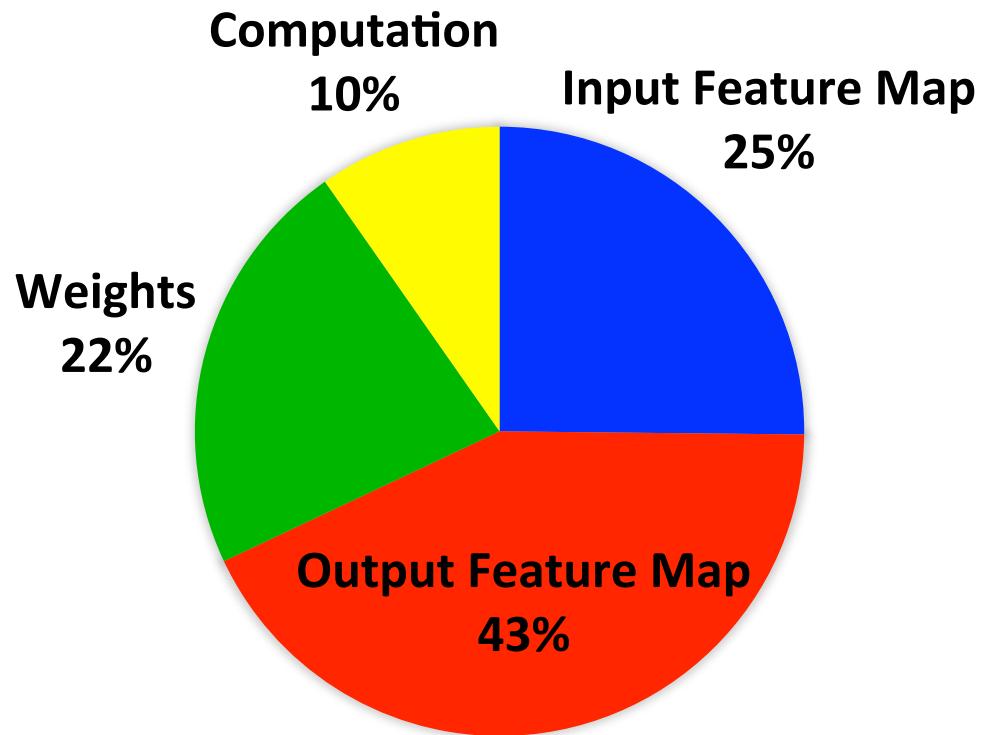
[Yang et al., CVPR 2017]

CNN Energy Consumption

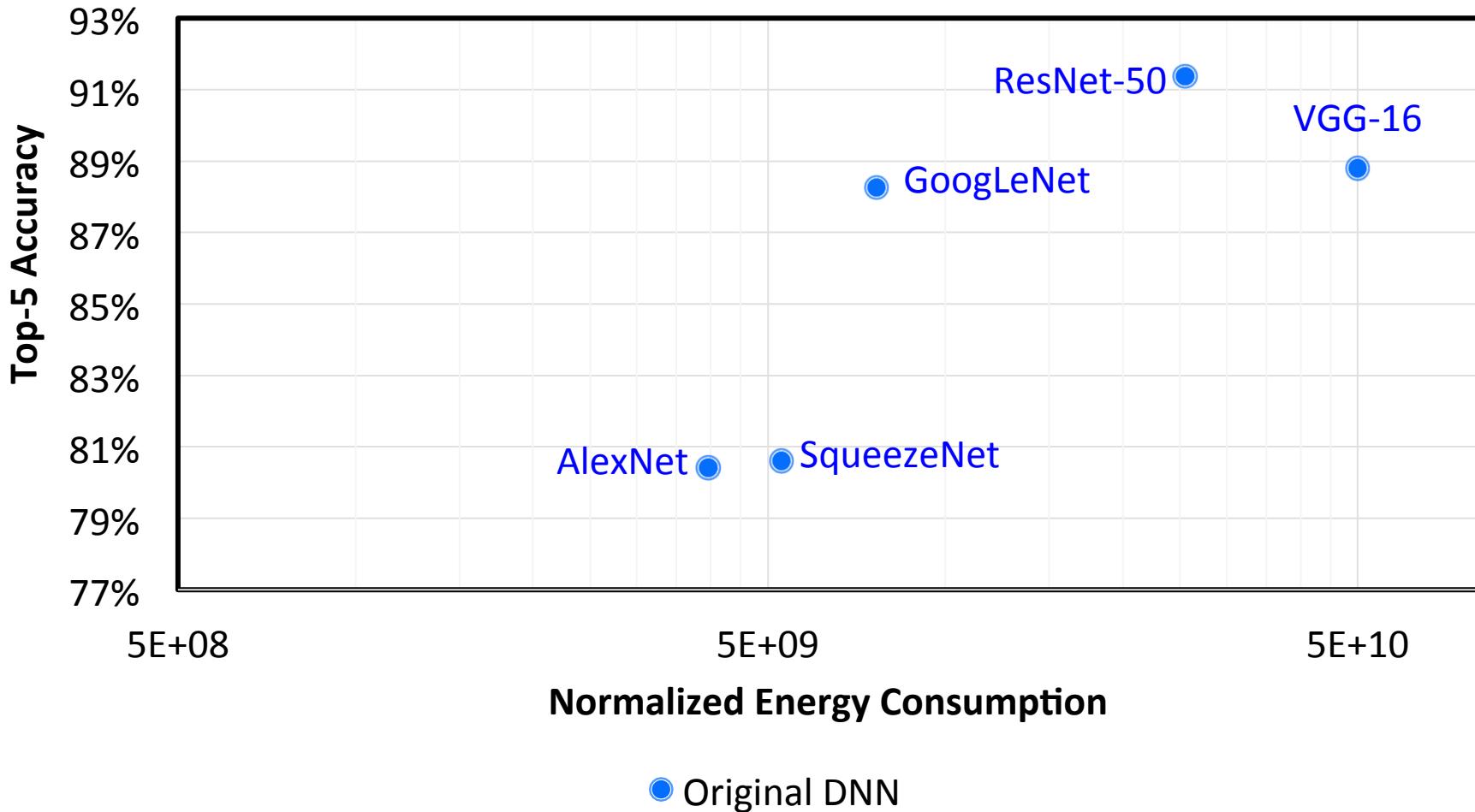
Key Observations

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

**Energy Consumption
of GoogLeNet**

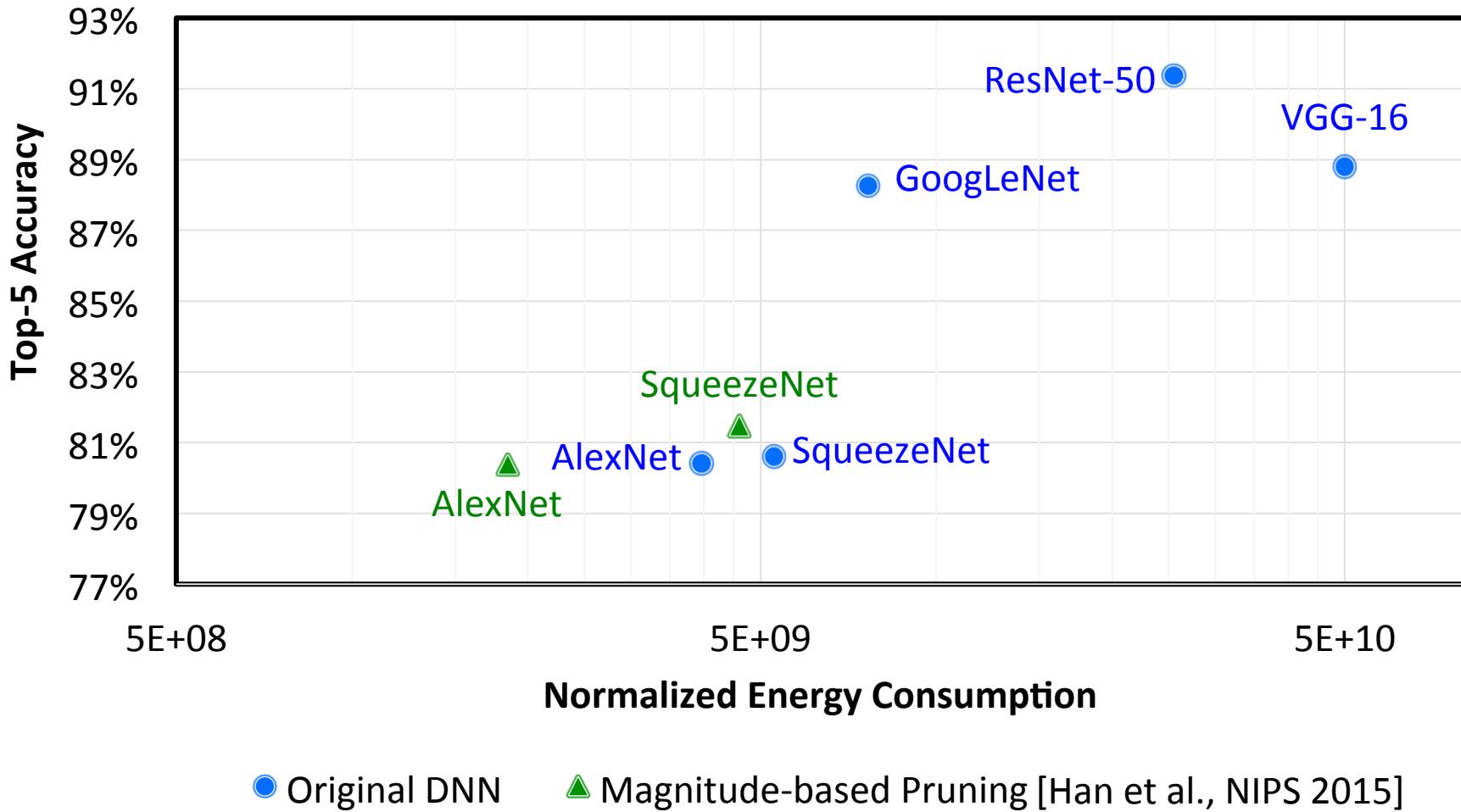


Energy Consumption of Existing DNNs



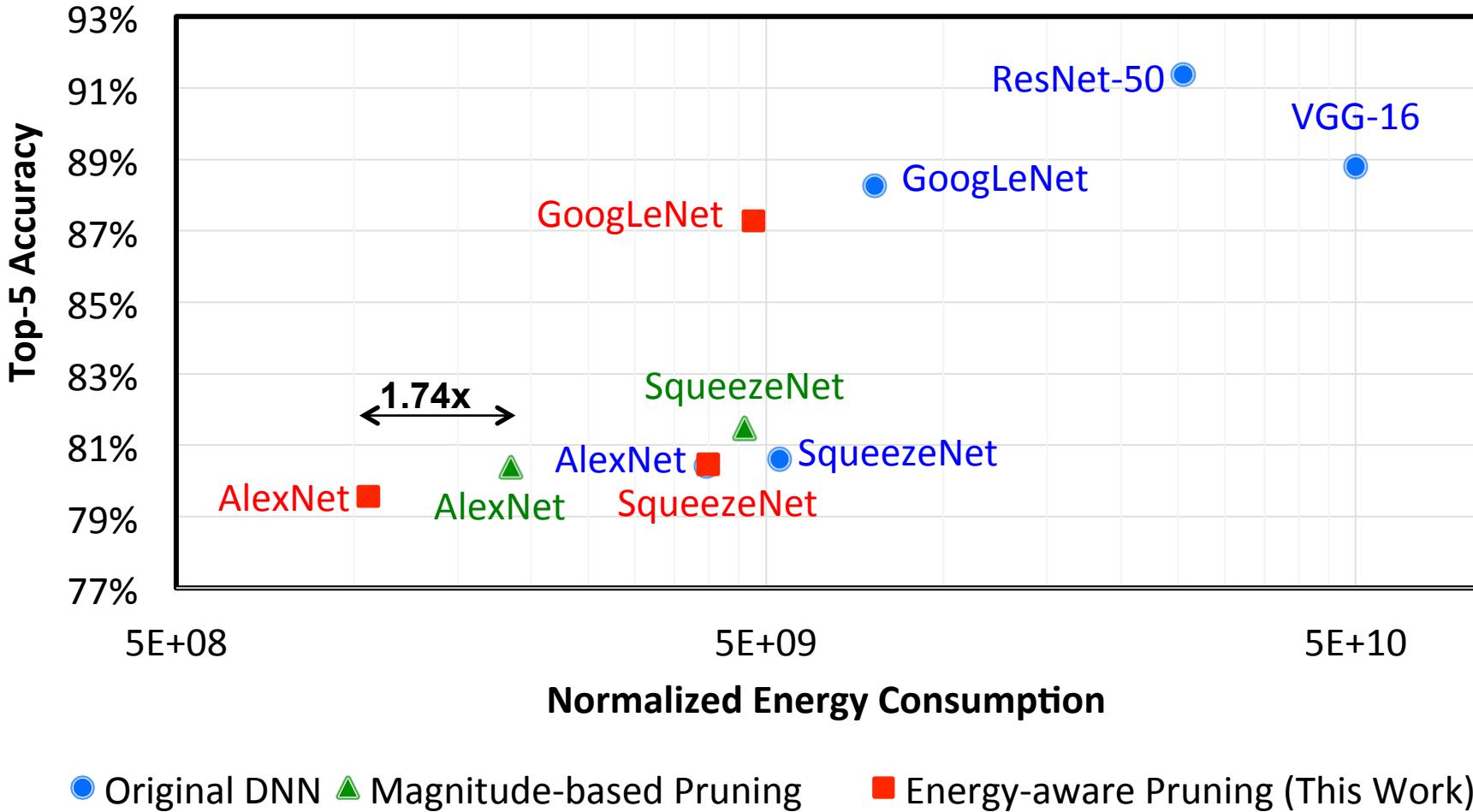
Deeper CNNs with fewer weights do not necessarily consume less energy than shallower CNNs with more weights

Magnitude-based Weight Pruning



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

Energy-Aware Pruning



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

Compact Network Architectures

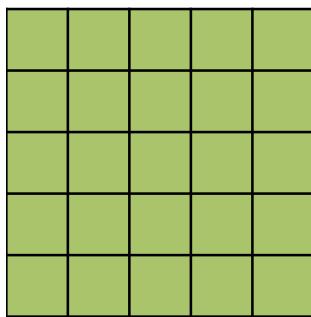
- Break large convolutional layers into a series of smaller convolutional layers
 - Fewer weights, but same effective receptive field
- Before Training: Network Architecture Design
- After Training: Decompose Trained Filters

Network Architecture Design

Build Network with series of Small Filters

GoogleNet/Inception v3

5x5 filter

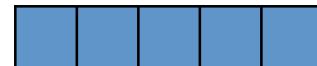


decompose
→

5x1 filter

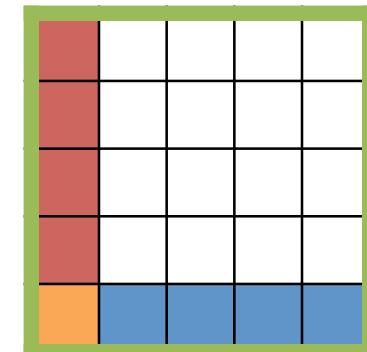


1x5 filter



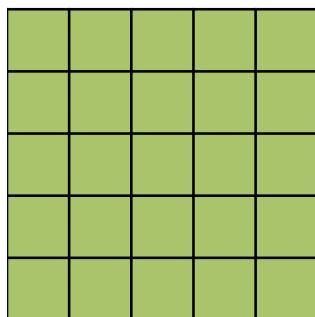
*separable
filters*

Apply sequentially



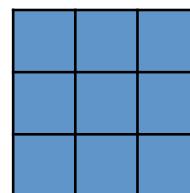
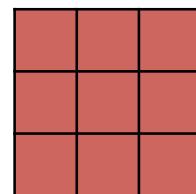
VGG-16

5x5 filter

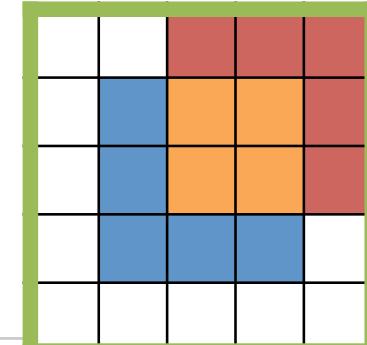


decompose
→

Two 3x3 filters



Apply sequentially



Network Architecture Design

Reduce size and computation with 1x1 Filter (**bottleneck**)

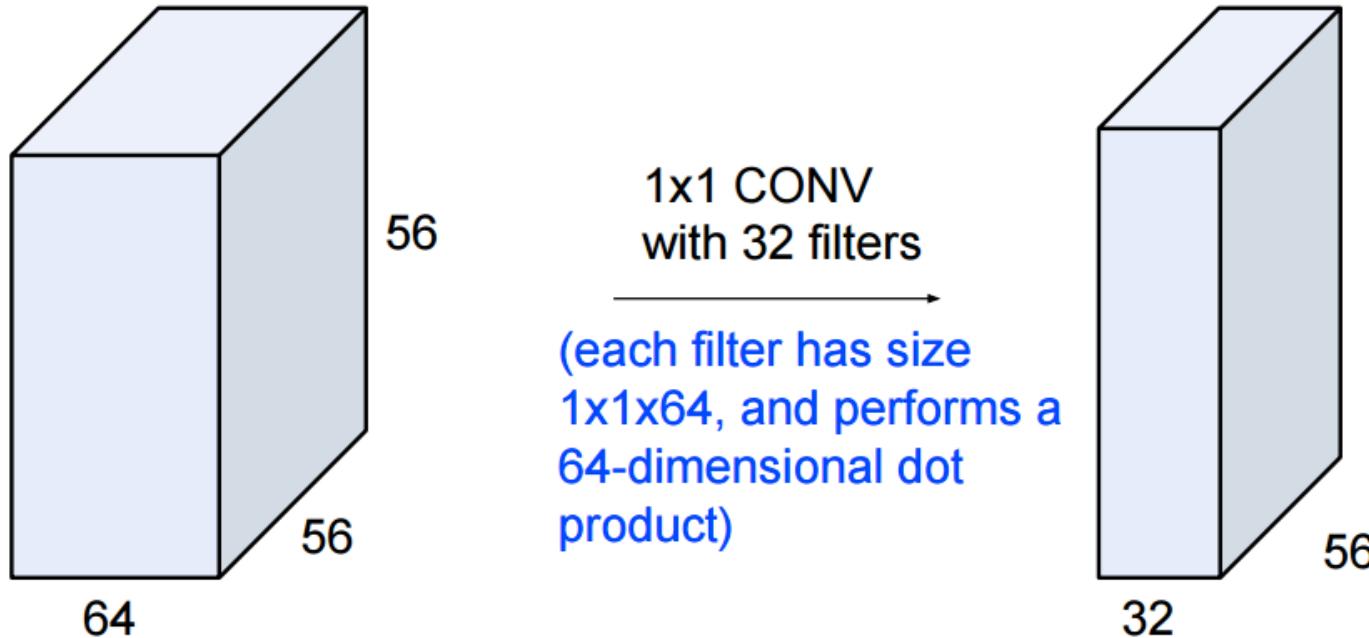


Figure Source:
Stanford cs231n

Used in Network In Network(NiN) and GoogLeNet

[Lin et al., ArXiV 2013 / ICLR 2014] [Szegedy et al., ArXiV 2014 / CVPR 2015]

Network Architecture Design

Reduce size and computation with 1x1 Filter (**bottleneck**)

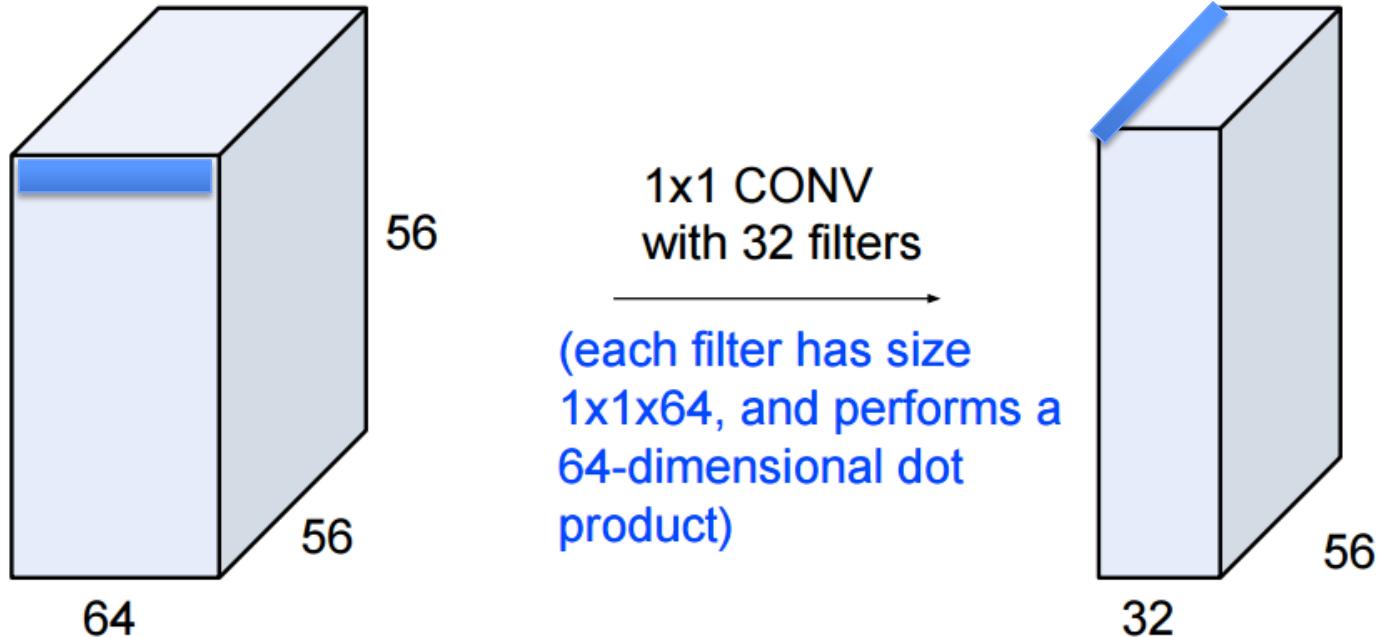


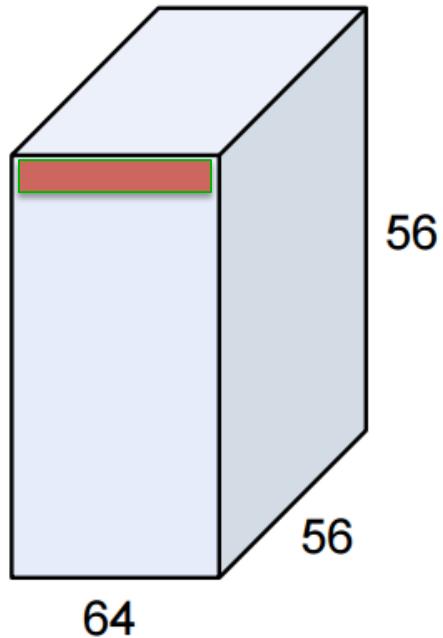
Figure Source:
Stanford cs231n

Used in Network In Network(NiN) and GoogLeNet

[Lin et al., ArXiV 2013 / ICLR 2014] [Szegedy et al., ArXiV 2014 / CVPR 2015]

Network Architecture Design

Reduce size and computation with 1x1 Filter (**bottleneck**)



1x1 CONV
with 32 filters
→
(each filter has size
1x1x64, and performs a
64-dimensional dot
product)

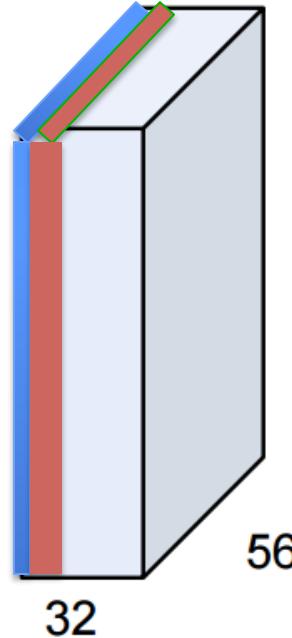


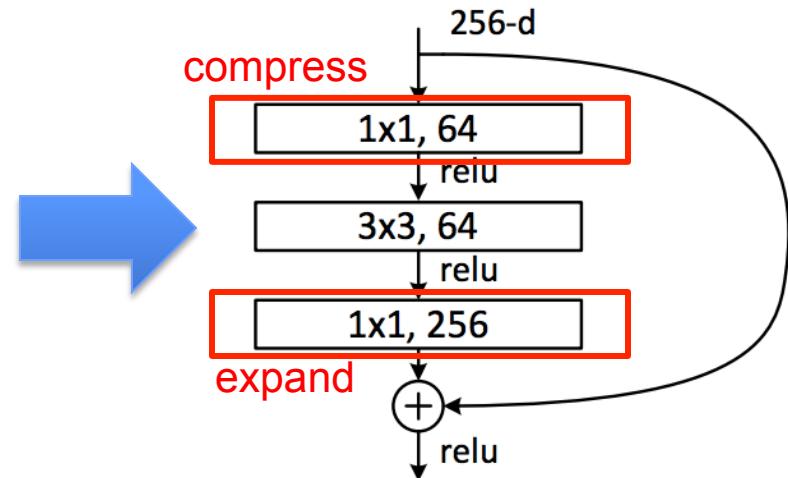
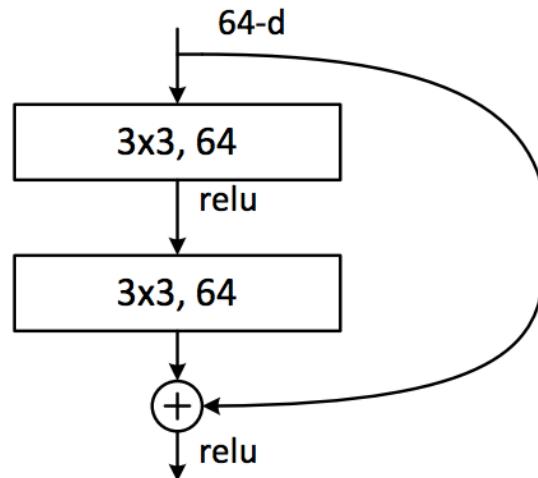
Figure Source:
Stanford cs231n

Used in Network In Network(NiN) and GoogLeNet

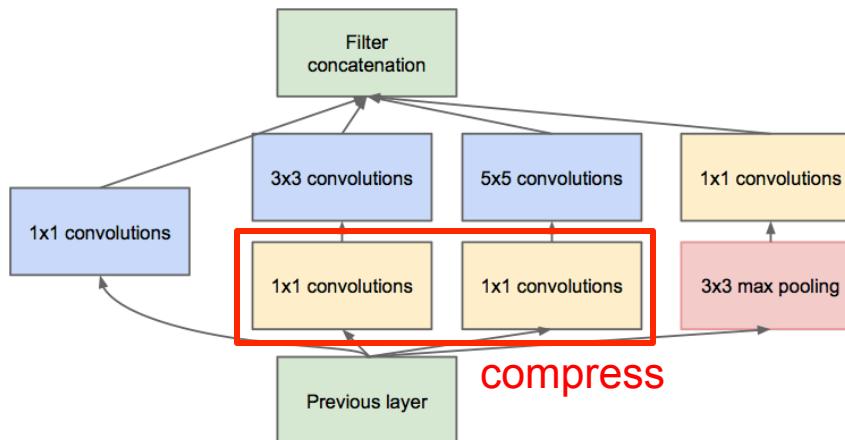
[Lin et al., ArXiV 2013 / ICLR 2014] [Szegedy et al., ArXiV 2014 / CVPR 2015]

Bottleneck in Popular DNN models

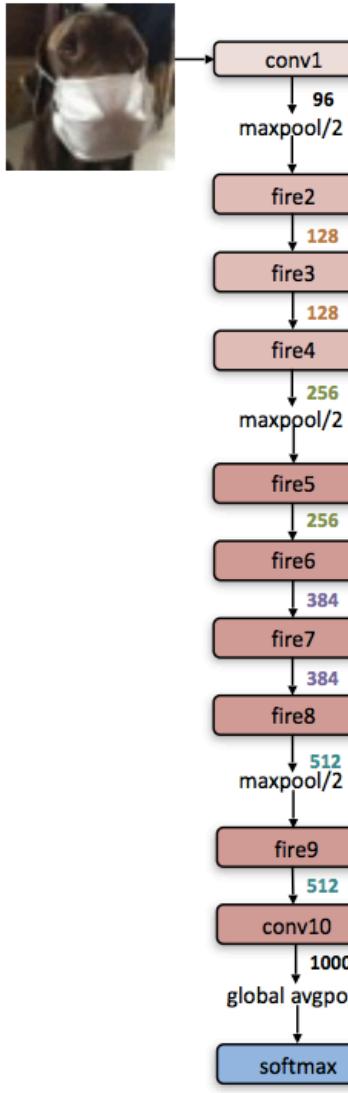
ResNet



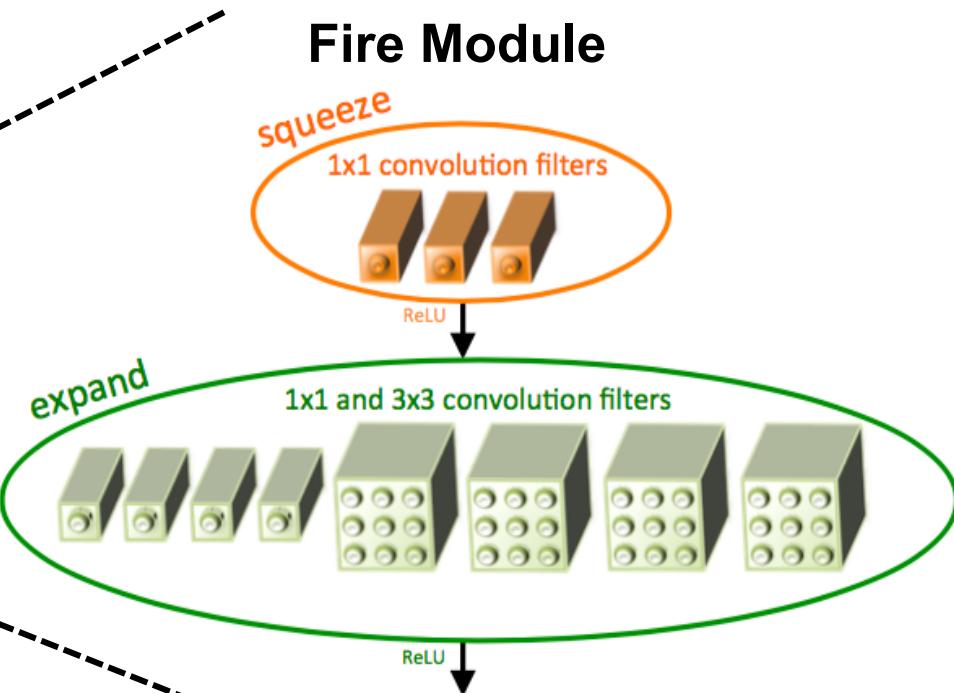
GoogleNet



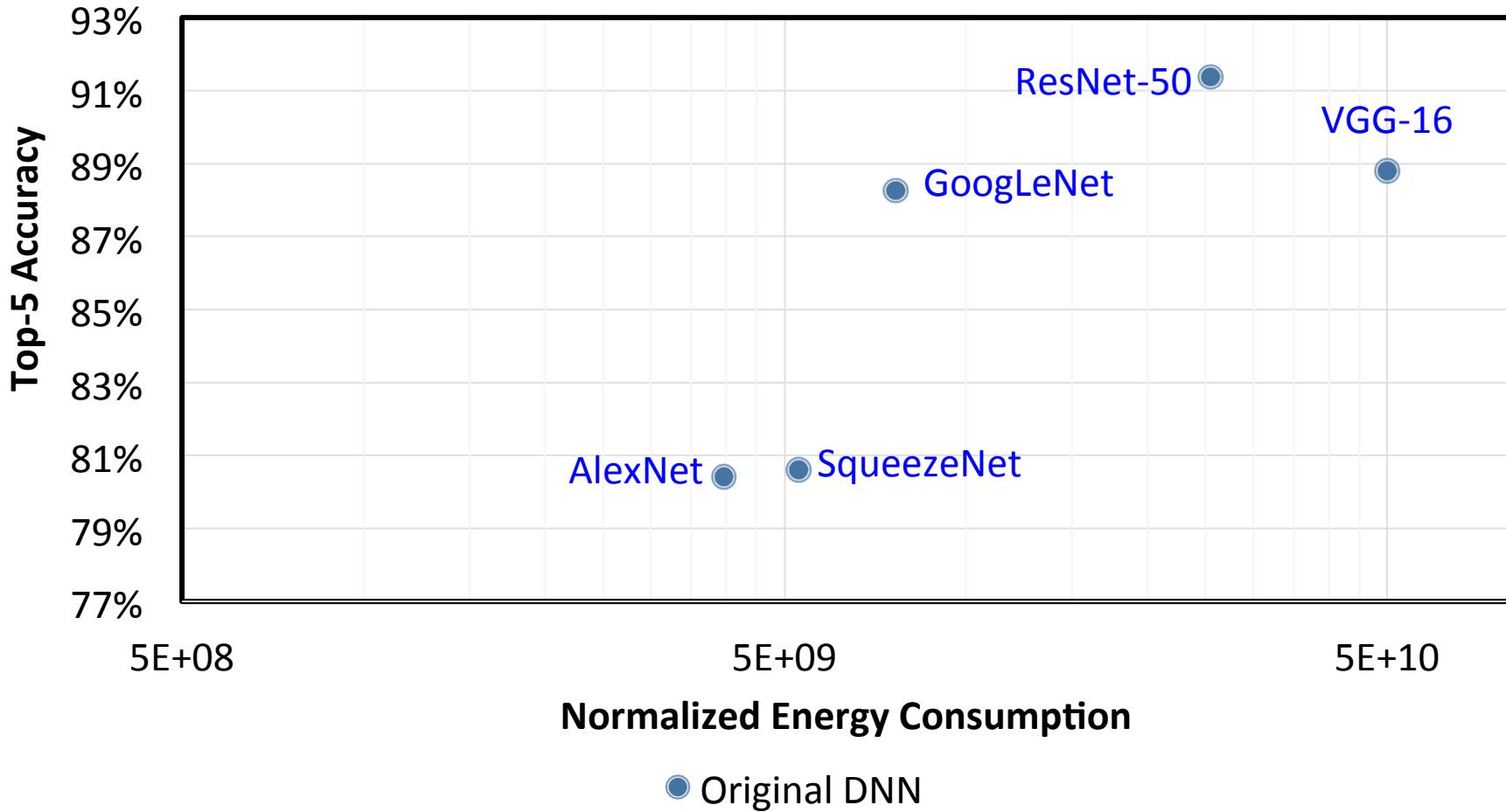
SqueezeNet



Reduce weights by reducing number of input channels by “squeezing” with 1x1
50x fewer weights than AlexNet (no accuracy loss)



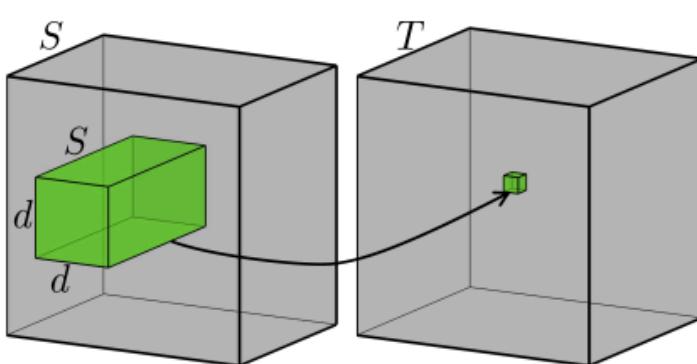
Energy Consumption of Existing DNNs



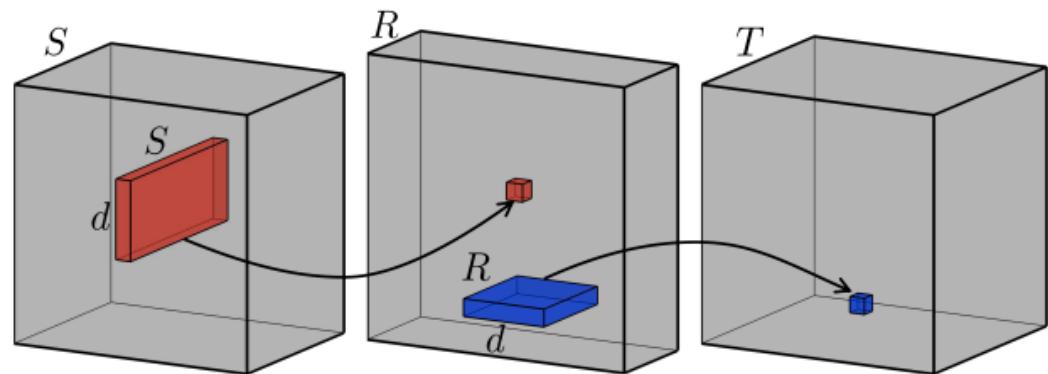
Deeper CNNs with fewer weights do not necessarily consume less energy than shallower CNNs with more weights

Decompose Trained Filters

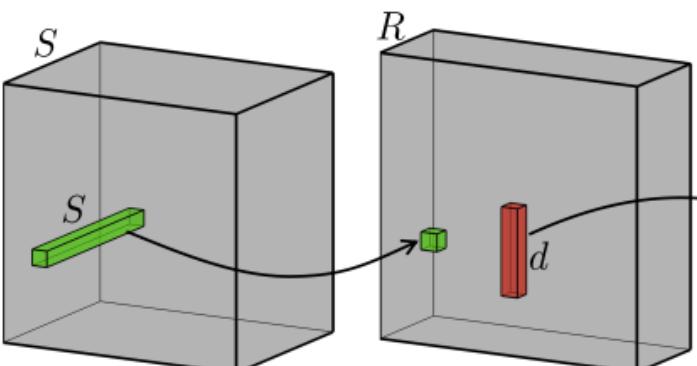
After training, perform **low-rank approximation** by applying tensor decomposition to weight kernel; then **fine-tune** weights for accuracy



(a) Full convolution

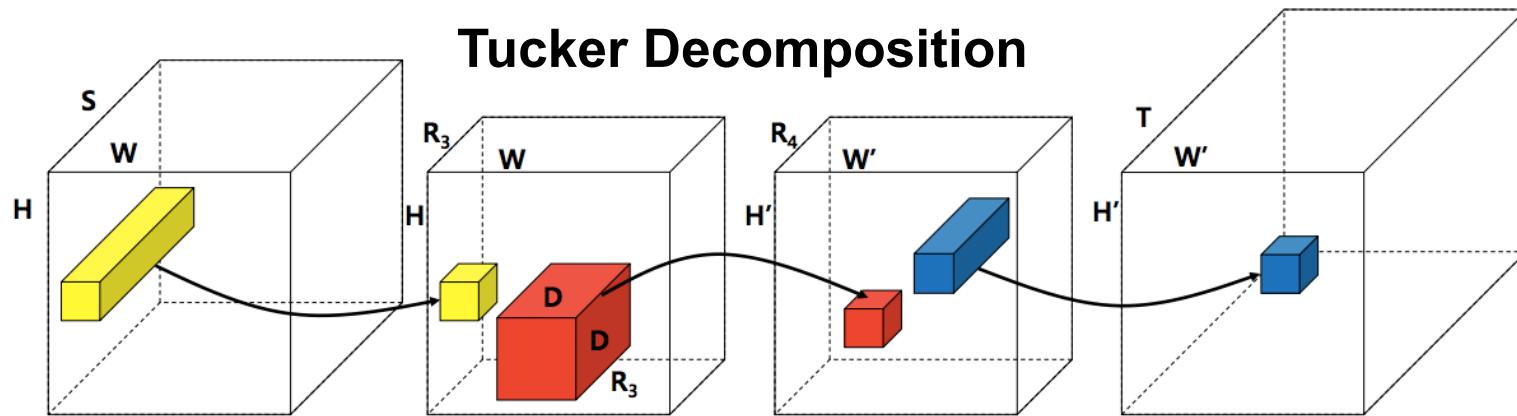


(b) Two-component decomposition (Jaderberg et al., 2014a)



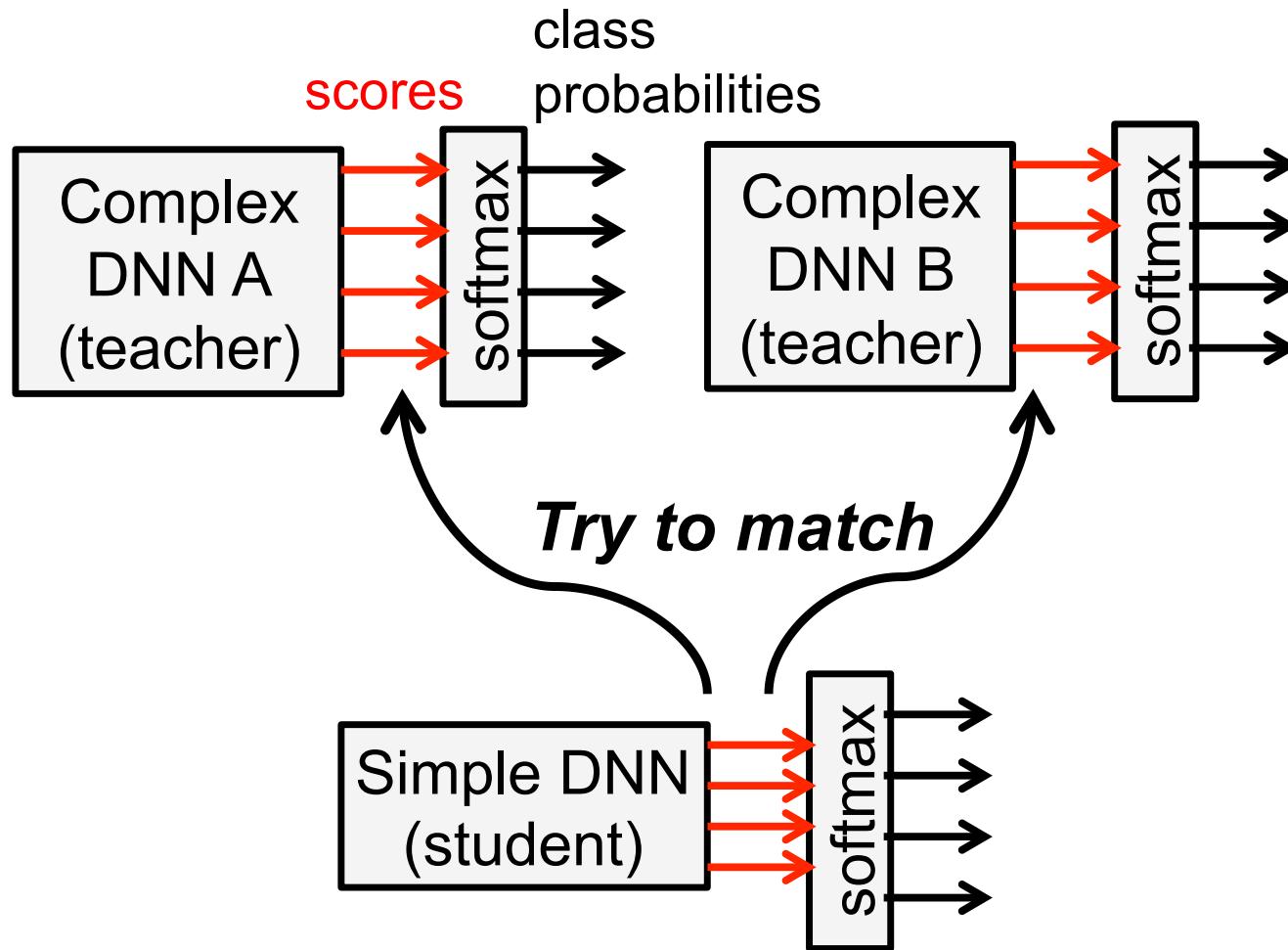
(c) CP-decomposition

Decompose Trained Filters on Phone



Model	Top-5	Weights	FLOPs	S6	Titan X
<i>AlexNet</i>	80.03	61M	725M	117ms	245mJ
<i>AlexNet*</i>	78.33	11M	272M	43ms	72mJ
(imp.)	(-1.70)	($\times 5.46$)	($\times 2.67$)	($\times 2.72$)	($\times 3.41$)
<i>VGG-S</i>	84.60	103M	2640M	357ms	825mJ
<i>VGG-S*</i>	84.05	14M	549M	97ms	193mJ
(imp.)	(-0.55)	($\times 7.40$)	($\times 4.80$)	($\times 3.68$)	($\times 4.26$)
<i>GoogLeNet</i>	88.90	6.9M	1566M	273ms	473mJ
<i>GoogLeNet*</i>	88.66	4.7M	760M	192ms	296mJ
(imp.)	(-0.24)	($\times 1.28$)	($\times 2.06$)	($\times 1.42$)	($\times 1.60$)
<i>VGG-16</i>	89.90	138M	15484M	1926ms	4757mJ
<i>VGG-16*</i>	89.40	127M	3139M	576ms	1346mJ
(imp.)	(-0.50)	($\times 1.09$)	($\times 4.93$)	($\times 3.34$)	($\times 3.53$)

Knowledge Distillation



[Bucilu et al., KDD 2006],[Hinton et al., arXiv 2015]

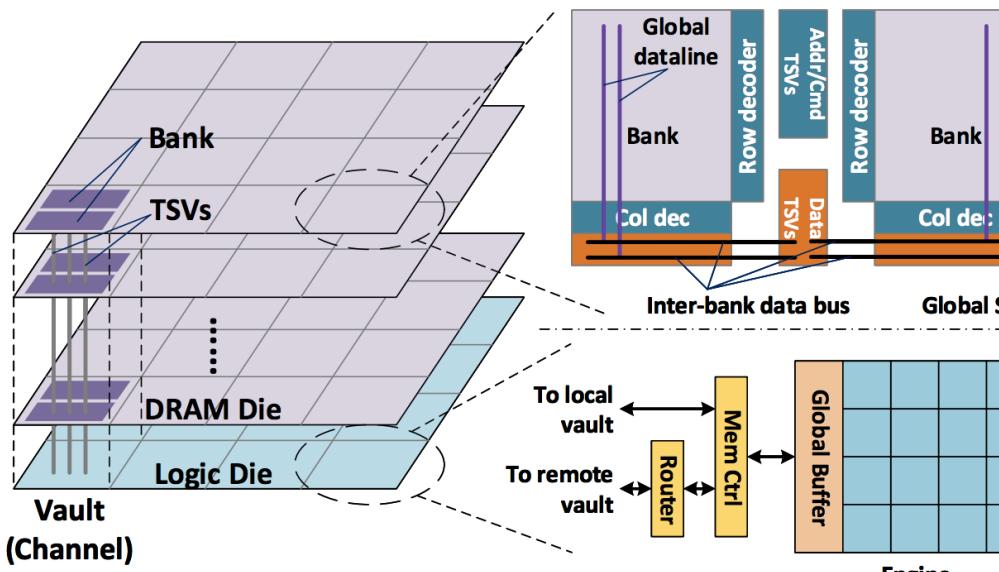
Opportunities in Advanced Technologies

Reduce data movement by embedding computation into memory and sensor

Advanced Memory Technologies

Many new memories and devices explored to reduce data movement

Stacked DRAM



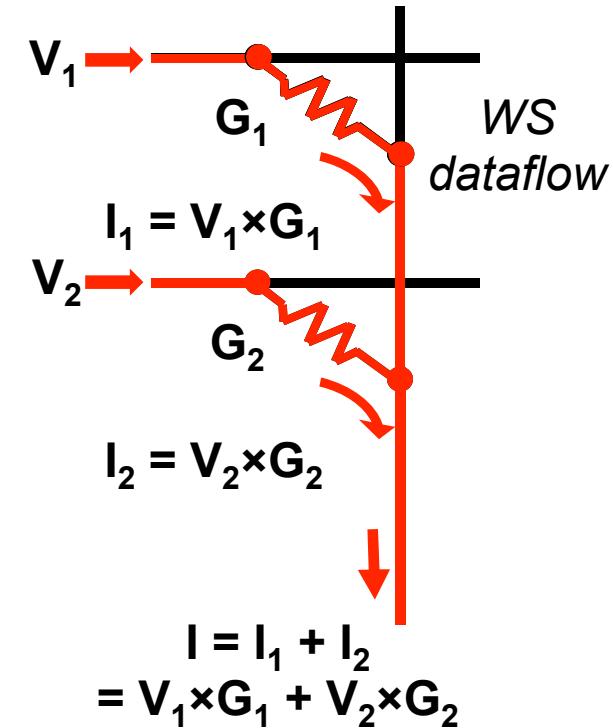
[Gao et al., Tetris, ASPLOS 2017]

[Kim et al., NeuroCube, ISCA 2016]

eDRAM

[Chen et al., DaDianNao, MICRO 2014]

Non-Volatile Resistive Memories



[Shafiee et al., ISCA 2016]

[Chi et al., PRIME, ISCA 2016]

Summary

- **Deep Learning is an important area of research**
 - Wide range of applications
 - Various methods to extract features (hand-crafted and learned)
- **Challenge is to balance the key metrics**
 - Accuracy, Energy, Throughput, Cost, etc.
- **Opportunities at various levels of hardware design**
 - Architecture, Joint Algorithm-Hardware, Mixed-Signal Circuits, Advanced Technologies
 - Important to consider interactions between levels to maximize impact

117

Acknowledgements



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



References

Overview Paper

V. Sze, Y.-H. Chen, T-J. Yang, J. Emer, “*Efficient Processing of Deep Neural Networks: A Tutorial and Survey*”, arXiv, 2017 <https://arxiv.org/pdf/1703.09039.pdf>

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

New Class Fall 2017

“6.S082/6.888 Hardware Architecture for Deep Learning”

For updates



Follow @eems_mit

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