

Accelerating Inference for DNNs

Yu-Hsin Chen

PhD, MIT
NVIDIA Research

CS217
October 18, 2018

Contributors to this Talk



Joel Emer

*Senior Distinguished
Research Scientist*

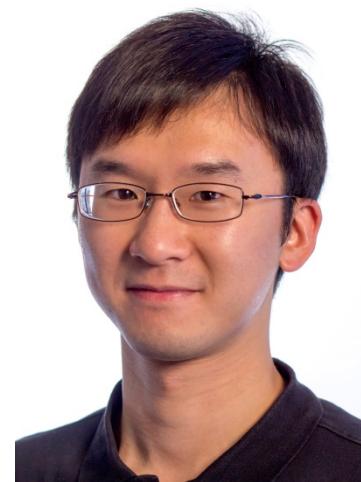
NVIDIA

Professor
MIT



Vivienne Sze

Professor
MIT



Yu-Hsin Chen

Research Scientist
NVIDIA



Tien-Ju Yang

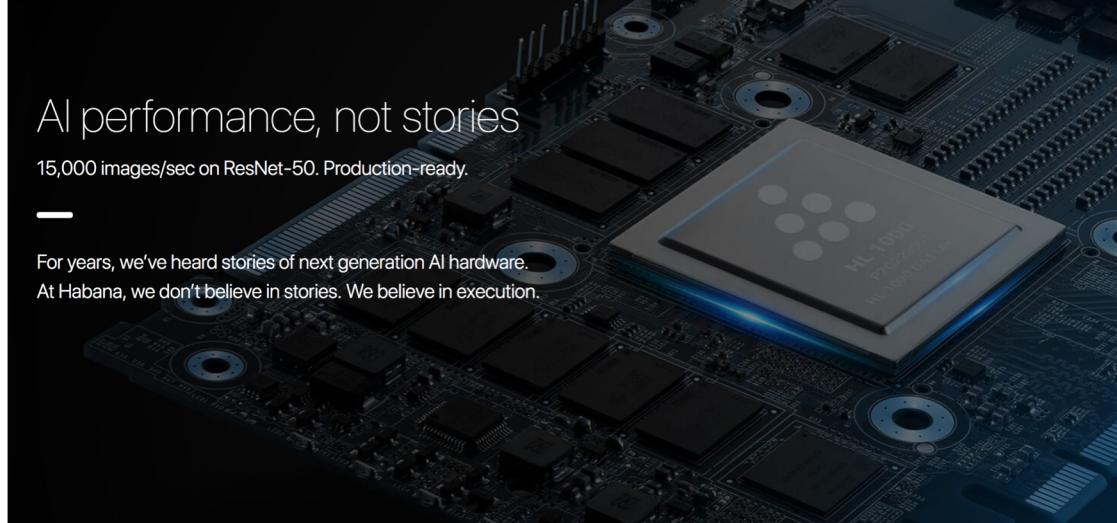
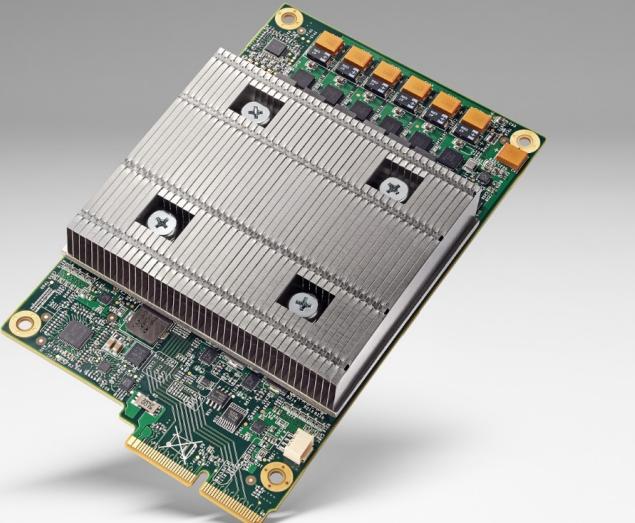
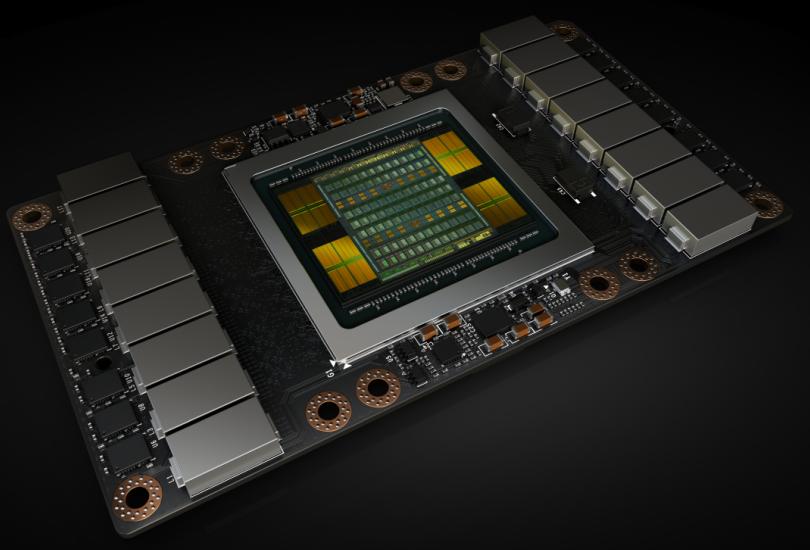
PhD Candidate
MIT

DNNs are Everywhere



[image sources: Amazon, nature, sciencenews.org]

Hardware Designed for DNNs



AI performance, not stories

15,000 images/sec on ResNet-50. Production-ready.

For years, we've heard stories of next generation AI hardware.
At Habana, we don't believe in stories. We believe in execution.

[image sources: NVIDIA, Huawei, Google, Habana Labs]

Key Metrics to Evaluate DNN Hardware

- **Accuracy**
 - Evaluate using the appropriate DNN model and dataset
- **Programmability**
 - Support multiple applications
- **Throughput / Latency**
 - GOPS, frame rate, delay
- **Energy / Power**
 - Energy per operation
 - DRAM Bandwidth
- **Cost**
 - Area (memory size, # of cores)

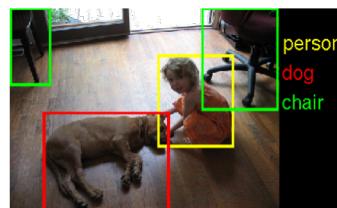
MNIST

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
1 2 2 2 2 3 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



Computer Vision



Speech Recognition



CHIP

BW

DRAM

Outlines

- **Architecture Design for DNN Accelerators**
- **DNN Model and Hardware Co-Design**

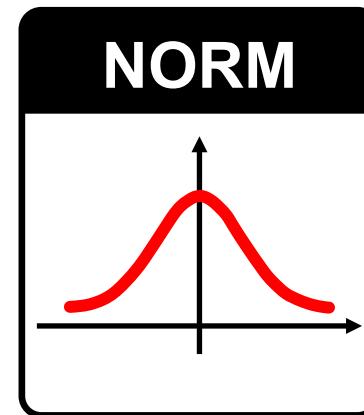
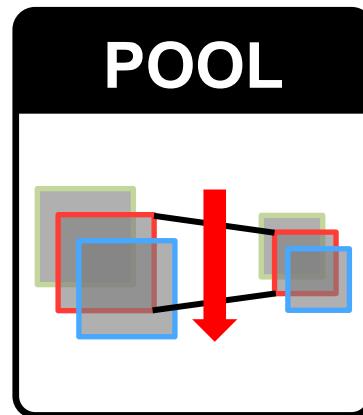
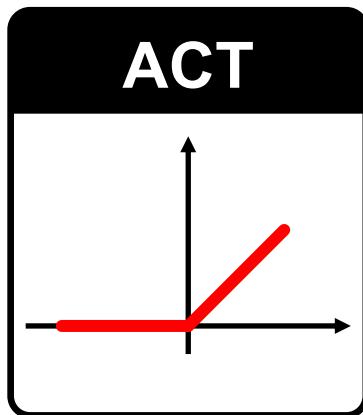
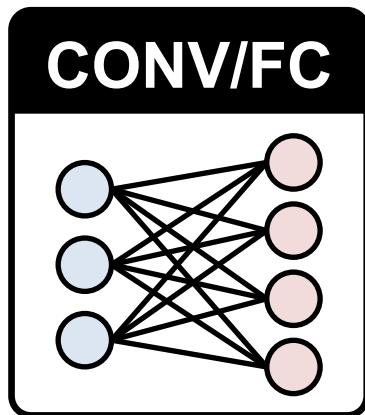
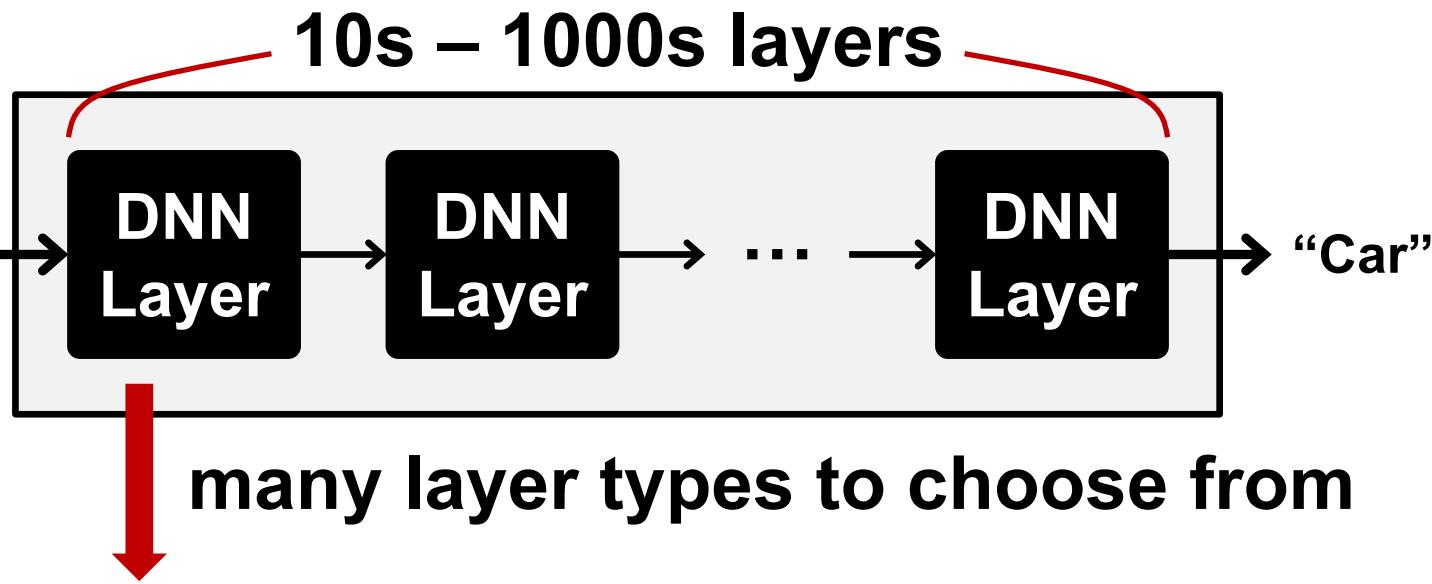
For more details:

Efficient Processing of Deep Neural Networks: A Tutorial and Survey
V.Sze, Y.-H. Chen, T.-J. Yang, J. Emer
Proceedings of IEEE
<https://arxiv.org/abs/1703.09039>

Your Takeaways

- Understand the trade-offs between various DNN accelerator architectures
- Assess the utility of various optimization approaches
- Understand recent implementation trends and opportunities

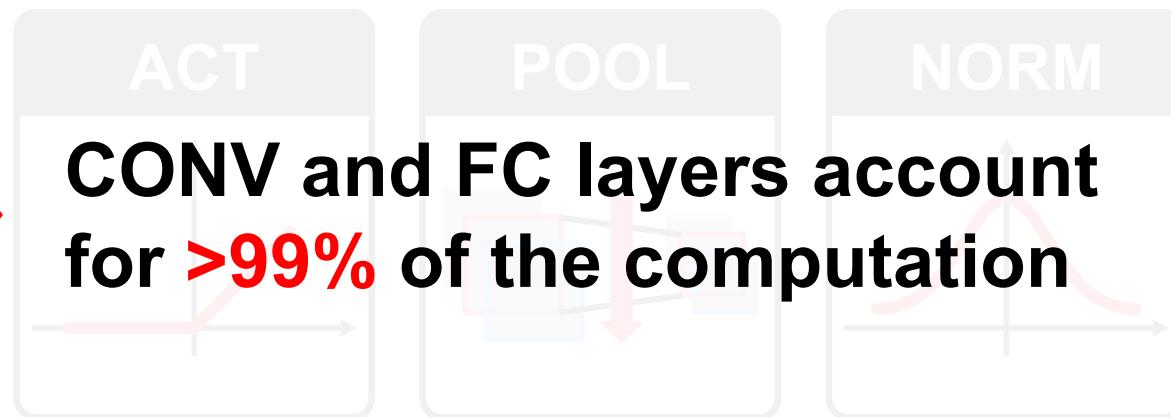
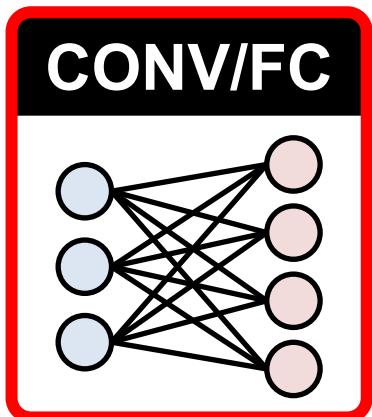
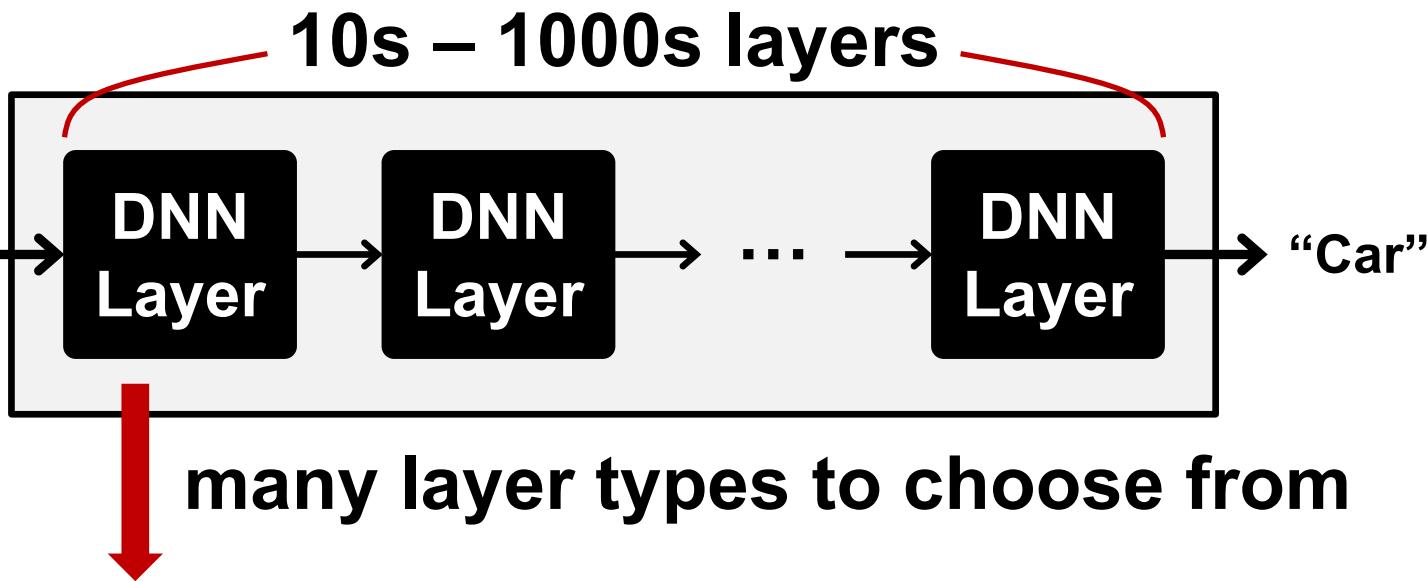
Primer on DNNs



...

[1] Tesla

Primer on DNNs

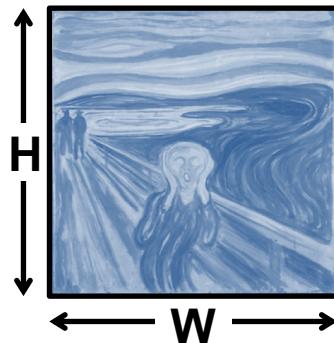
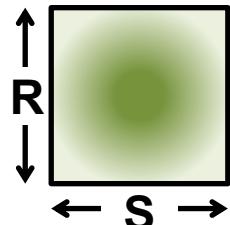


[1] Tesla

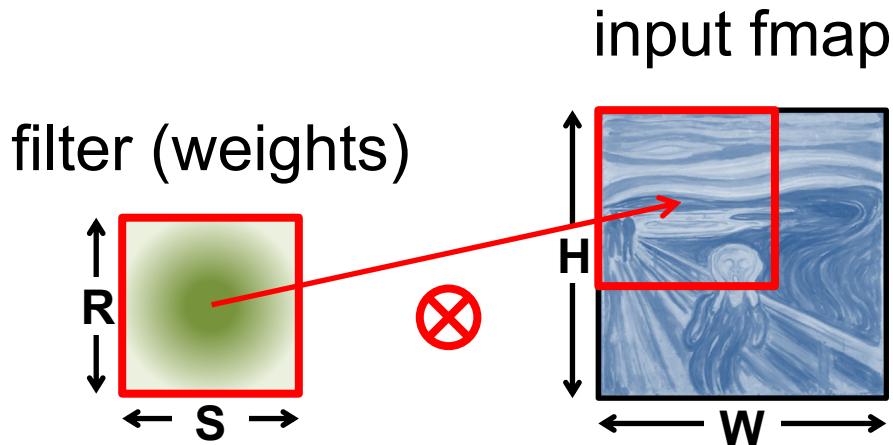
Convolutional (CONV) Layer

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

filter (weights)

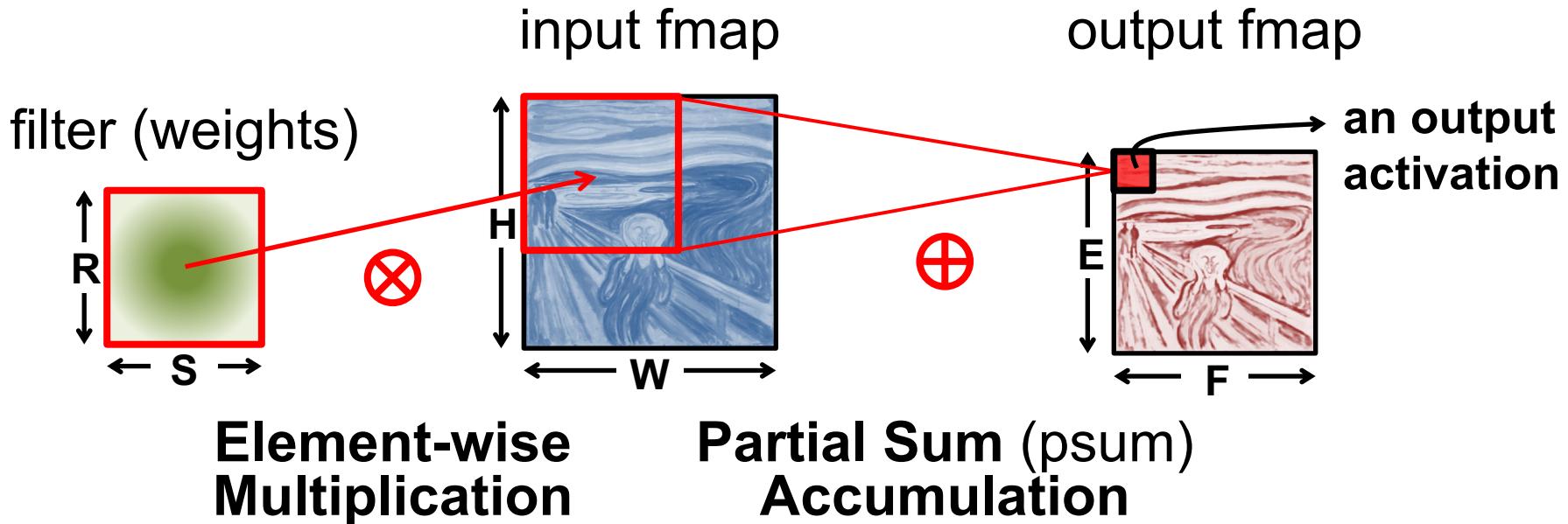


Convolutional (CONV) Layer

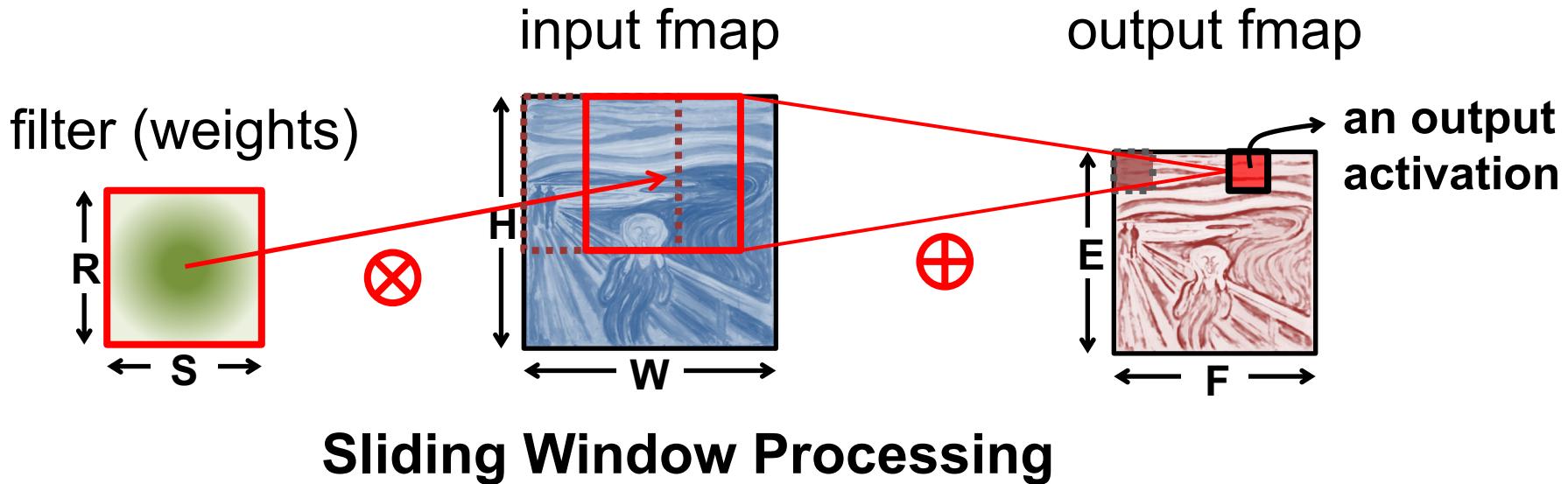


**Element-wise
Multiplication**

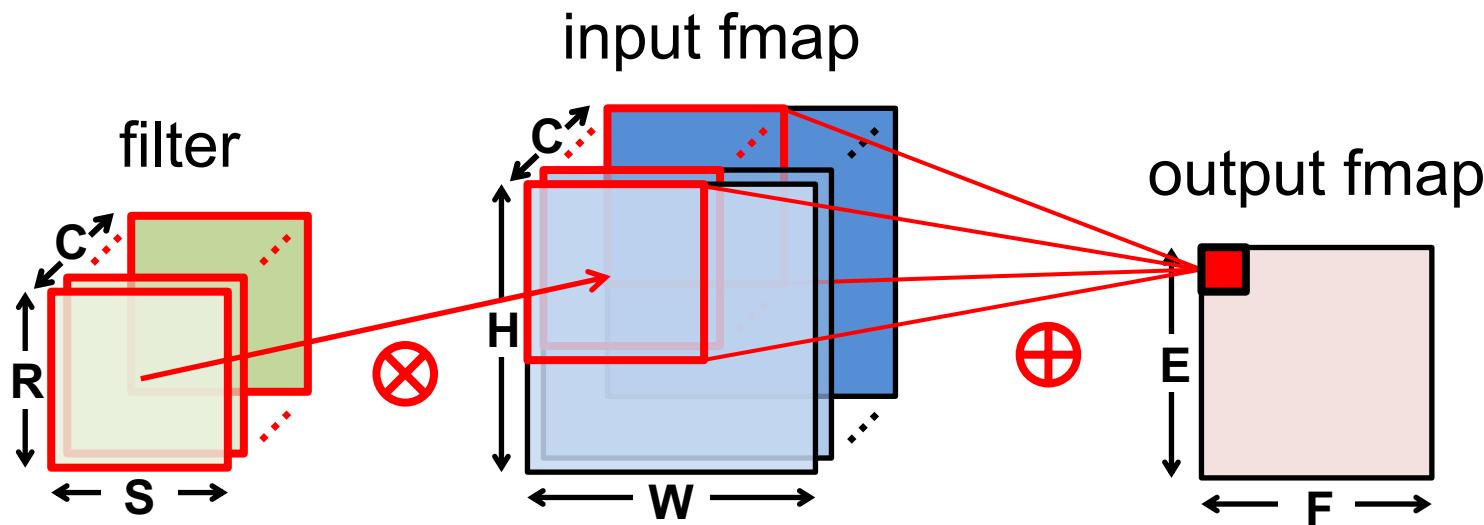
Convolutional (CONV) Layer



Convolutional (CONV) Layer

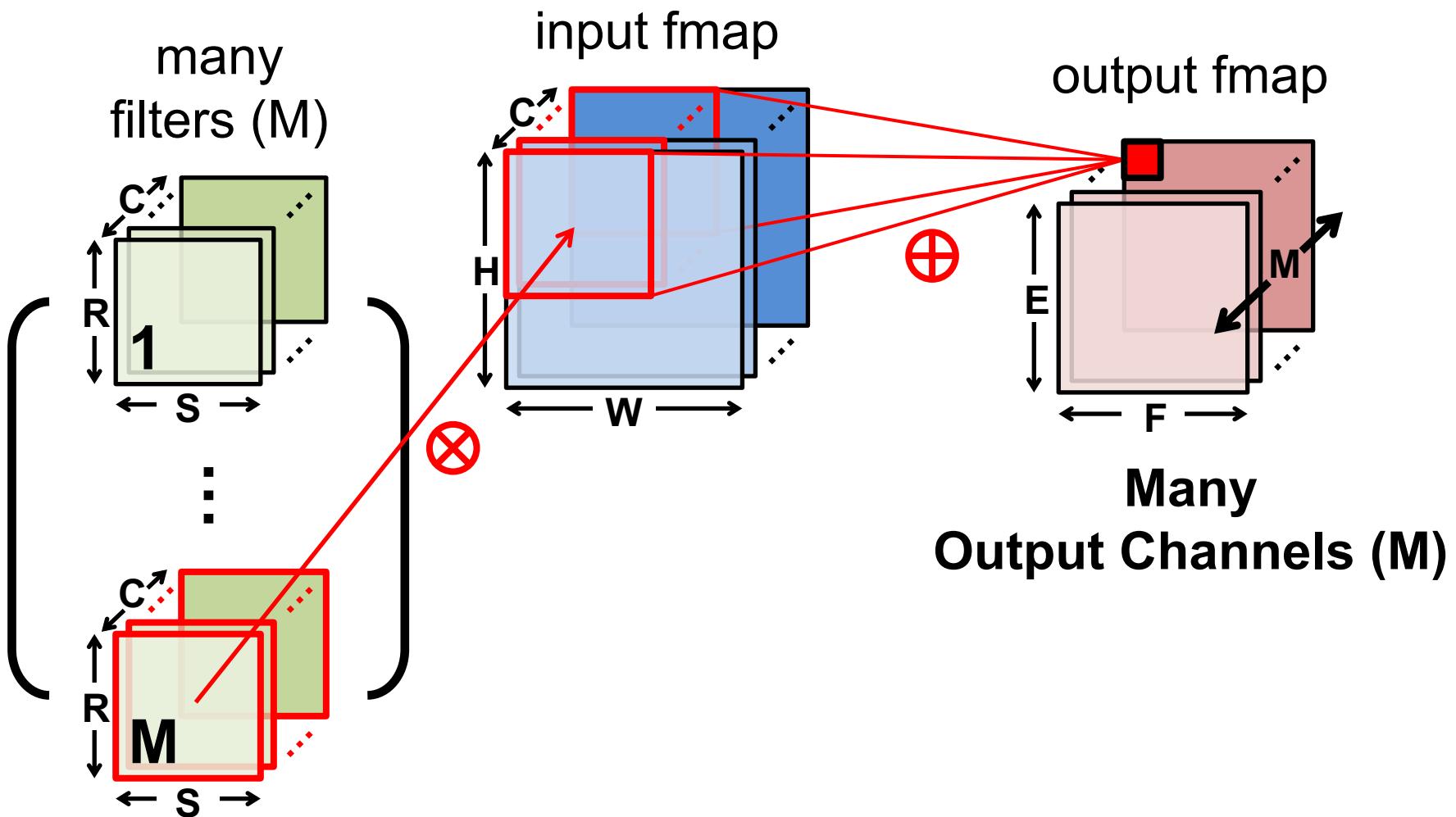


Convolutional (CONV) Layer

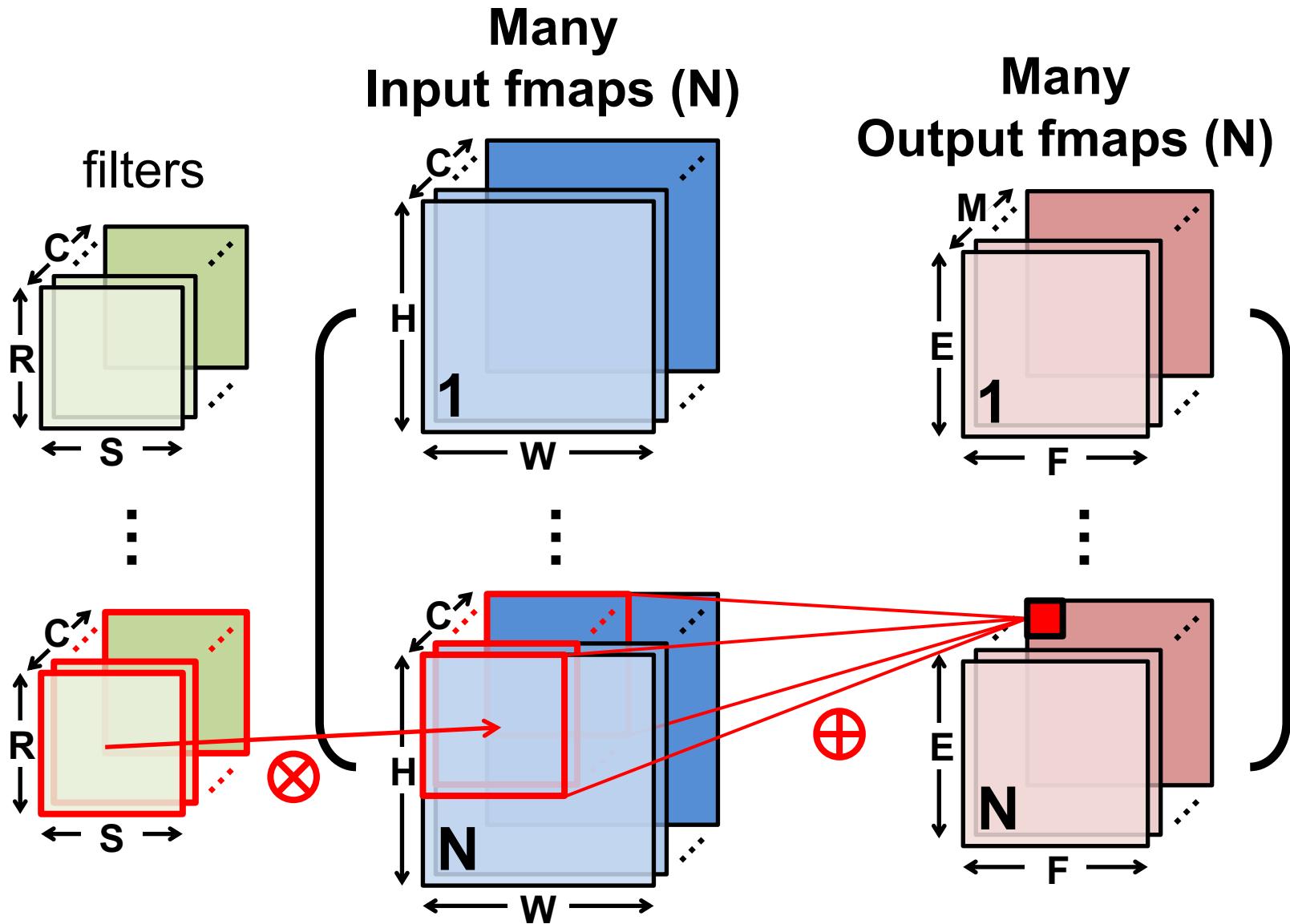


Many Input Channels (C)

Convolutional (CONV) Layer

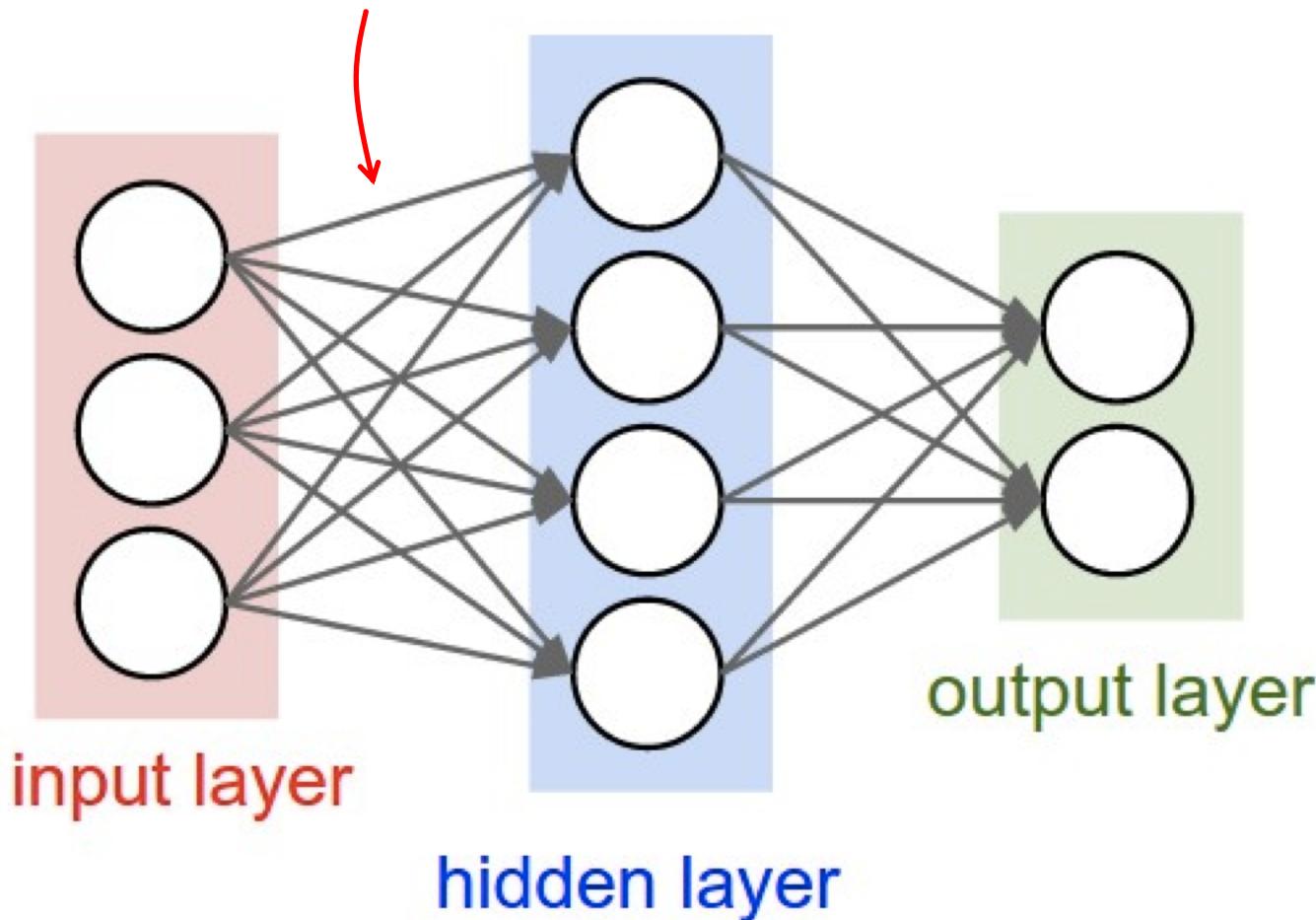


Convolutional (CONV) Layer



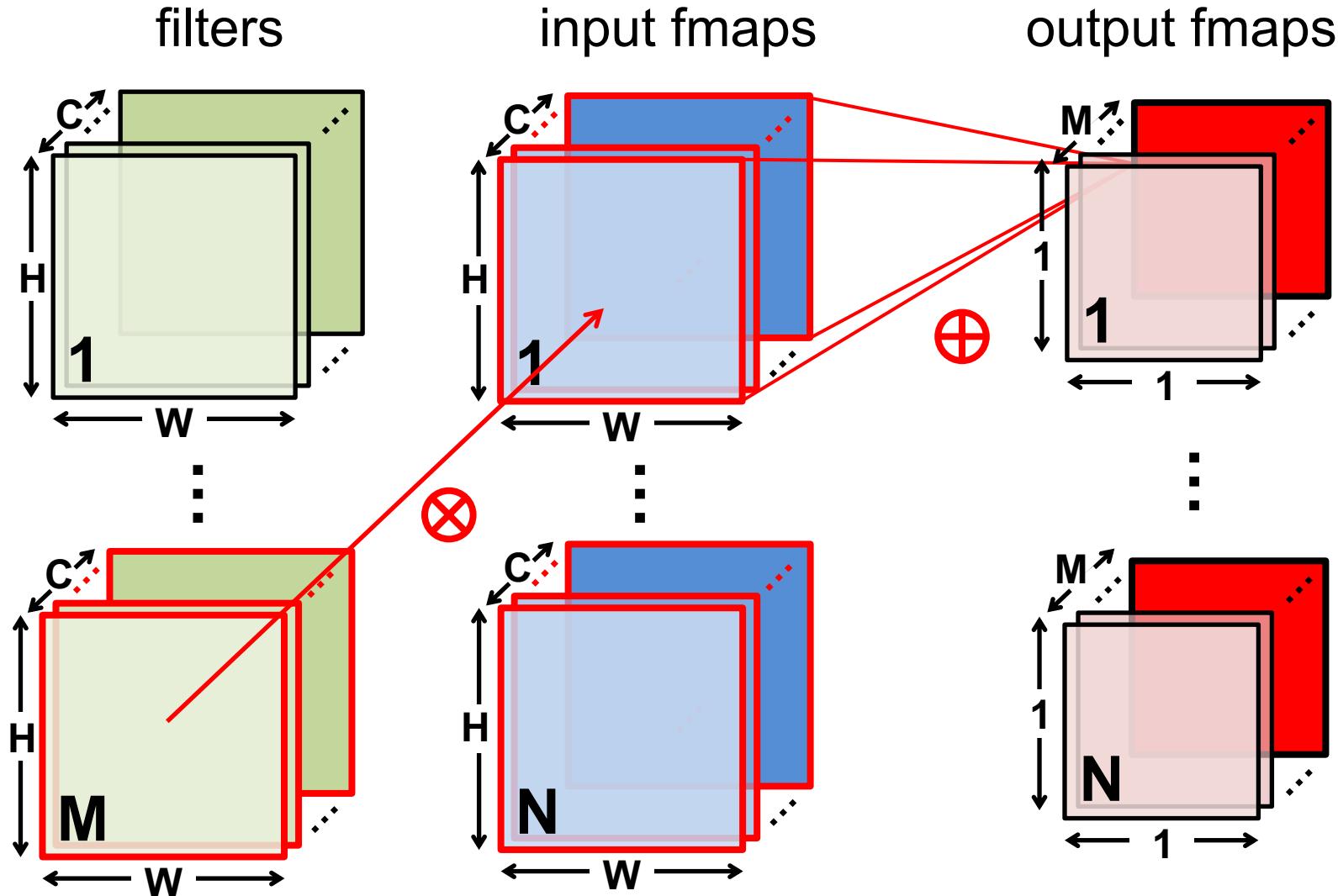
Fully-Connected (FC) Layer

All-to-all Connections



[images source: Stanford cs231n]

FC Layer – from CONV Layer POV

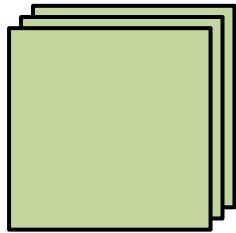


Widely Varying Layer Shapes

AlexNet^[1] Layer Shape Configurations

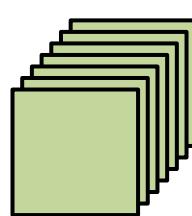
Layer	Filter Size (R)	# Filters (M)	# Channels (C)	Stride
1	11×11	96	3	4
2	5×5	256	48	1
3	3×3	384	256	1
8	1×1	1000	4096	1

Layer 1



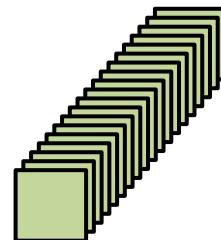
34k weights
105M MACs*

Layer 2



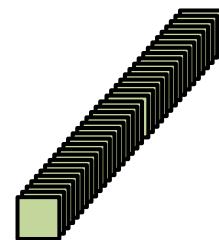
307k weights
224M MACs

Layer 3



885k weights
150M MACs

Layer 8



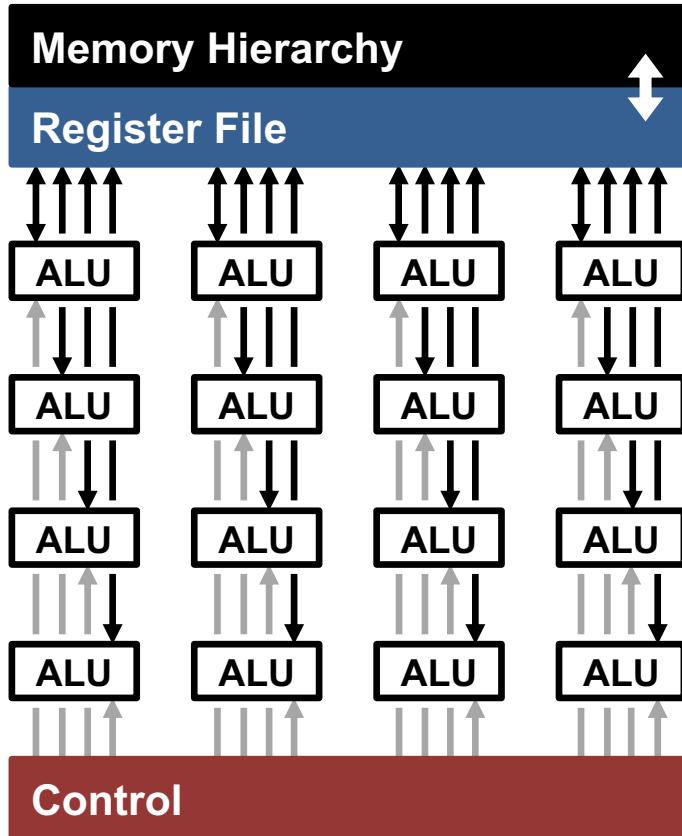
4096k weights
4M MACs

* multiply-accumulate

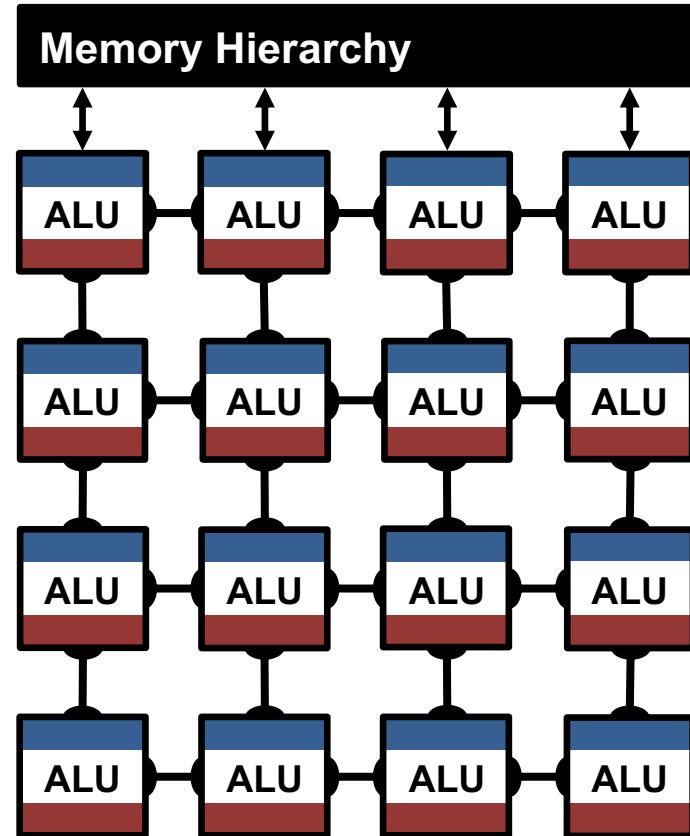
[1] Krizhevsky, NIPS 2012

Highly-Parallel Compute Paradigms

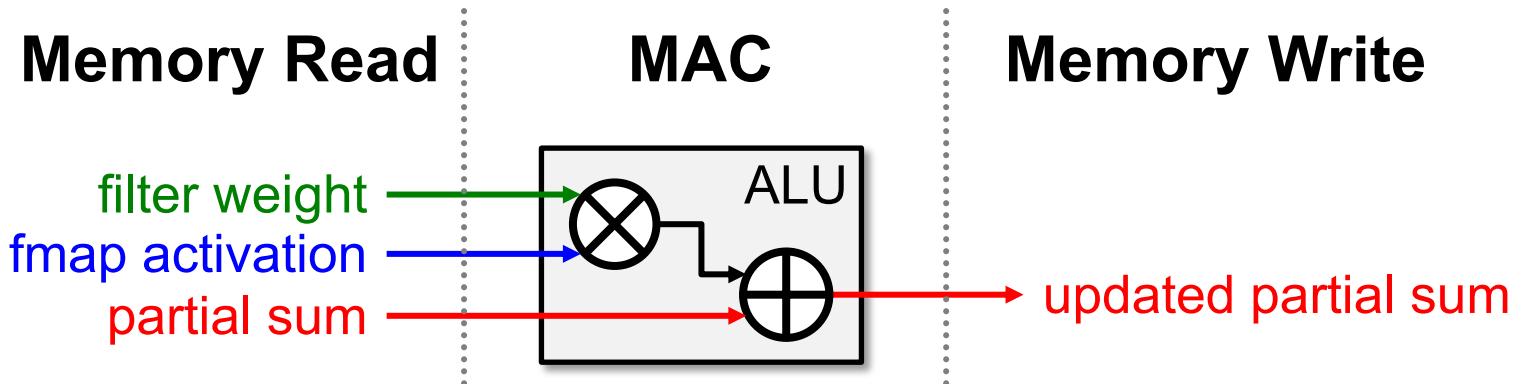
Temporal Architecture
(SIMD/SIMT)



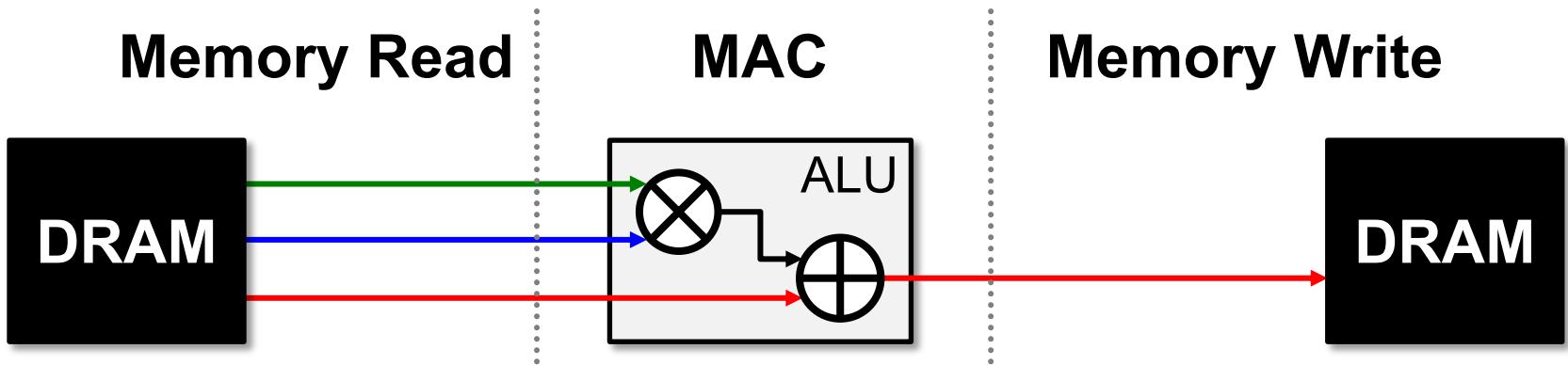
Spatial Architecture
(Dataflow Processing)



Memory Access is the Bottleneck



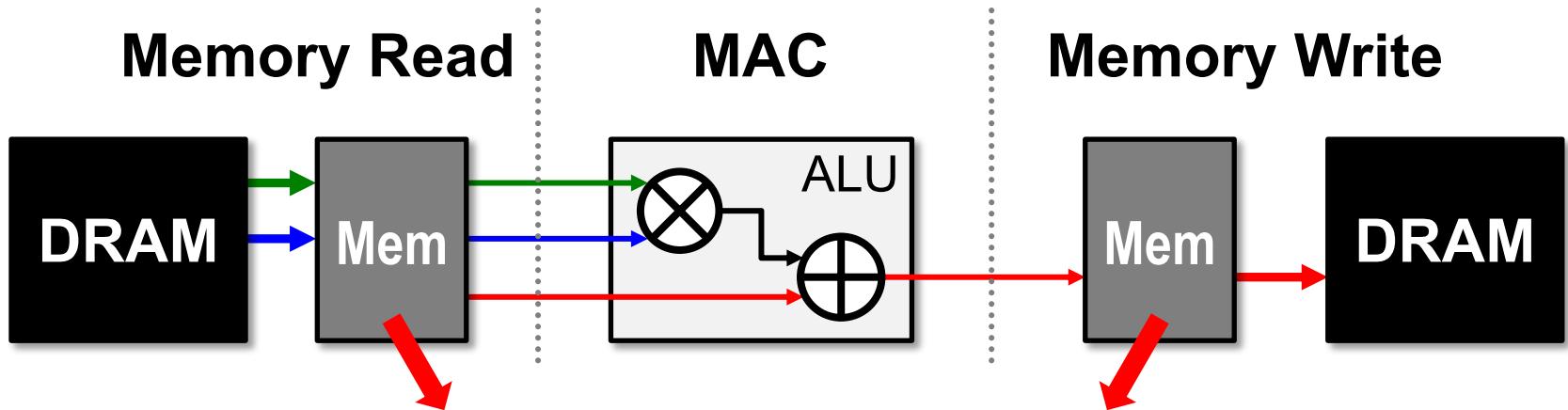
Memory Access is the Bottleneck



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

- Example: AlexNet has **724M** MACs
→ **2896M** DRAM accesses required

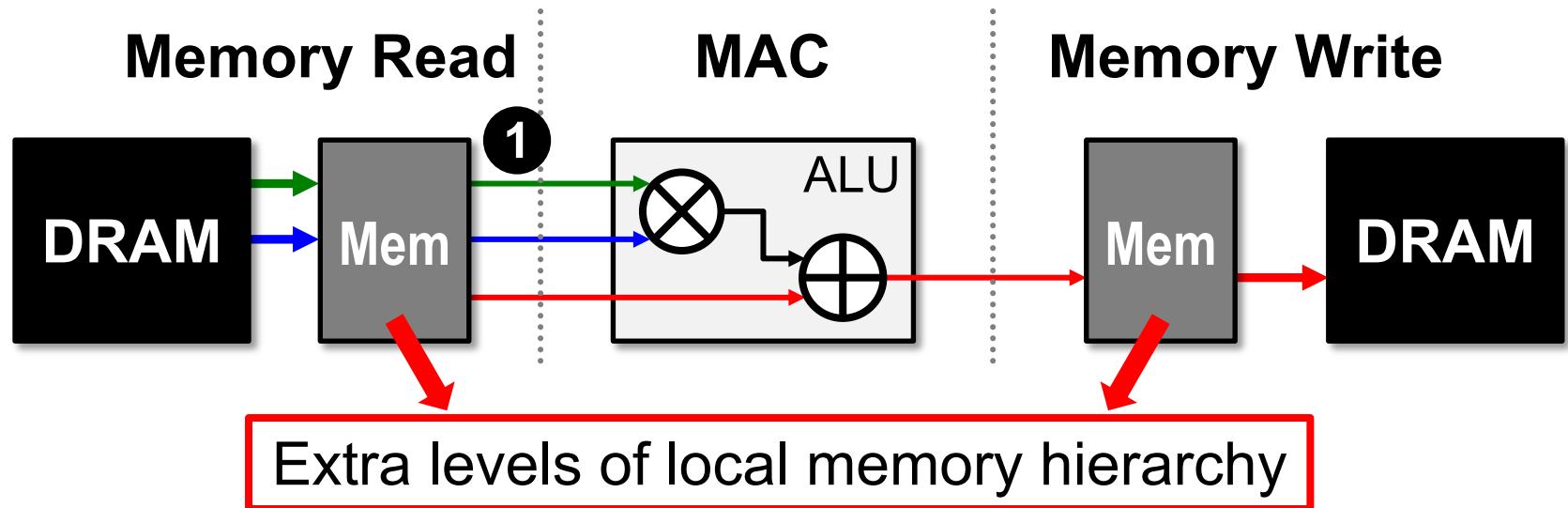
Memory Access is the Bottleneck



Extra levels of local memory hierarchy

Smaller, but Faster and more Energy-Efficient

Memory Access is the Bottleneck

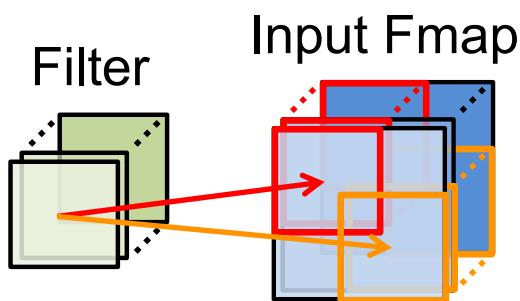


Opportunities: **① data reuse**

Types of Data Reuse in DNN

Convolutional Reuse

CONV layers only
(sliding window)

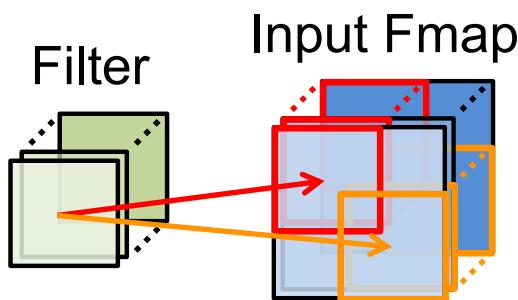


Reuse: Activations
Filter weights

Types of Data Reuse in DNN

Convolutional Reuse

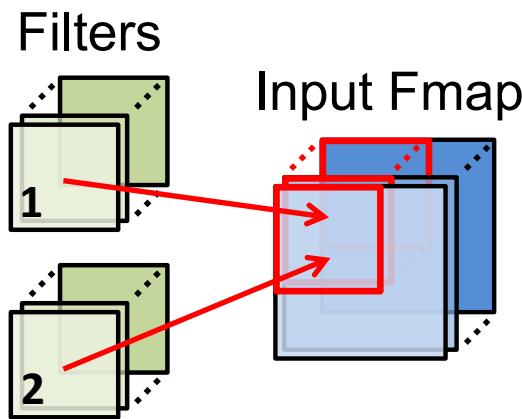
CONV layers only
(sliding window)



Reuse: Activations
Filter weights

Fmap Reuse

CONV and FC layers

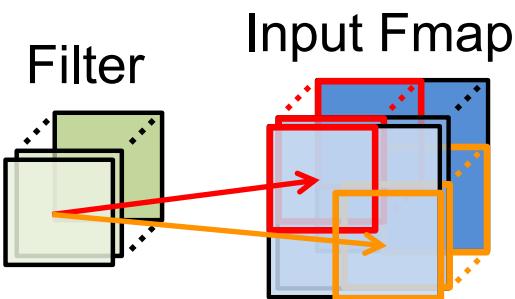


Reuse: Activations

Types of Data Reuse in DNN

Convolutional Reuse

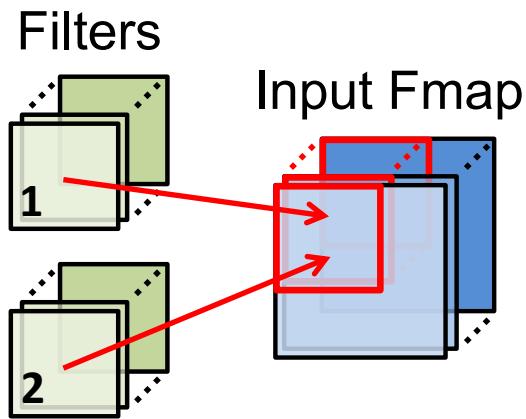
CONV layers only
(sliding window)



Reuse: Activations
Filter weights

Fmap Reuse

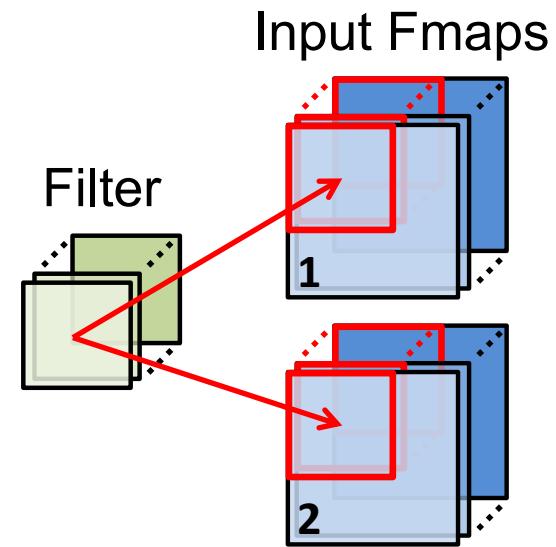
CONV and FC layers



Reuse: Activations

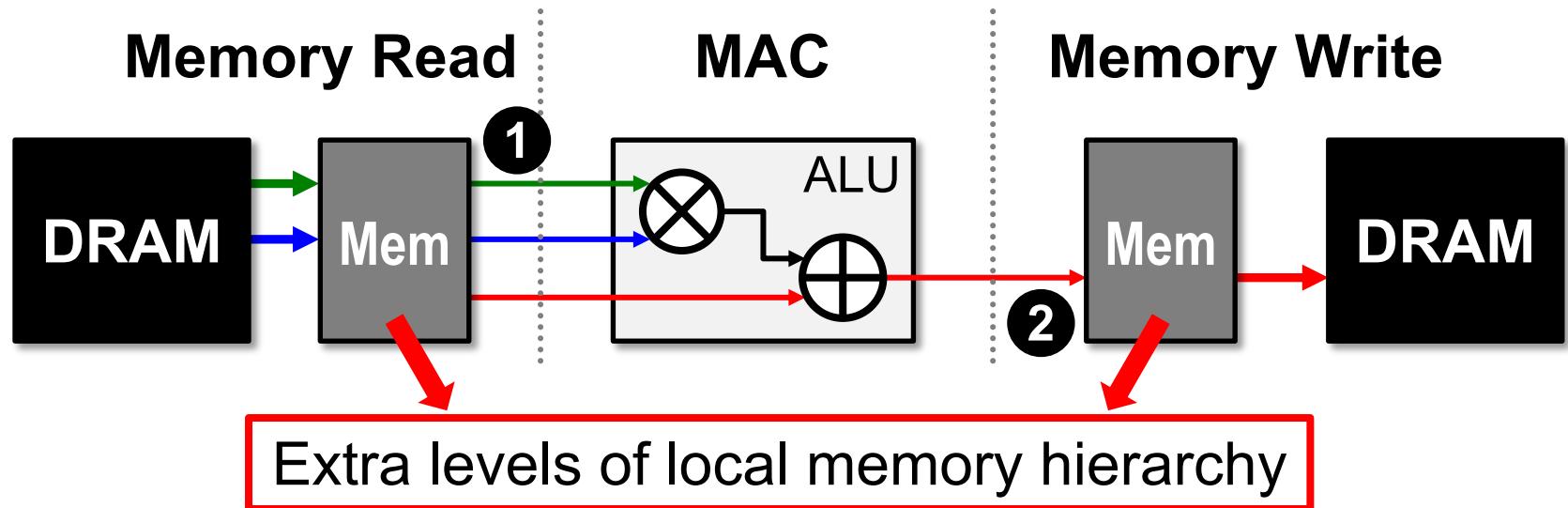
Filter Reuse

CONV and FC layers
(batch size > 1)



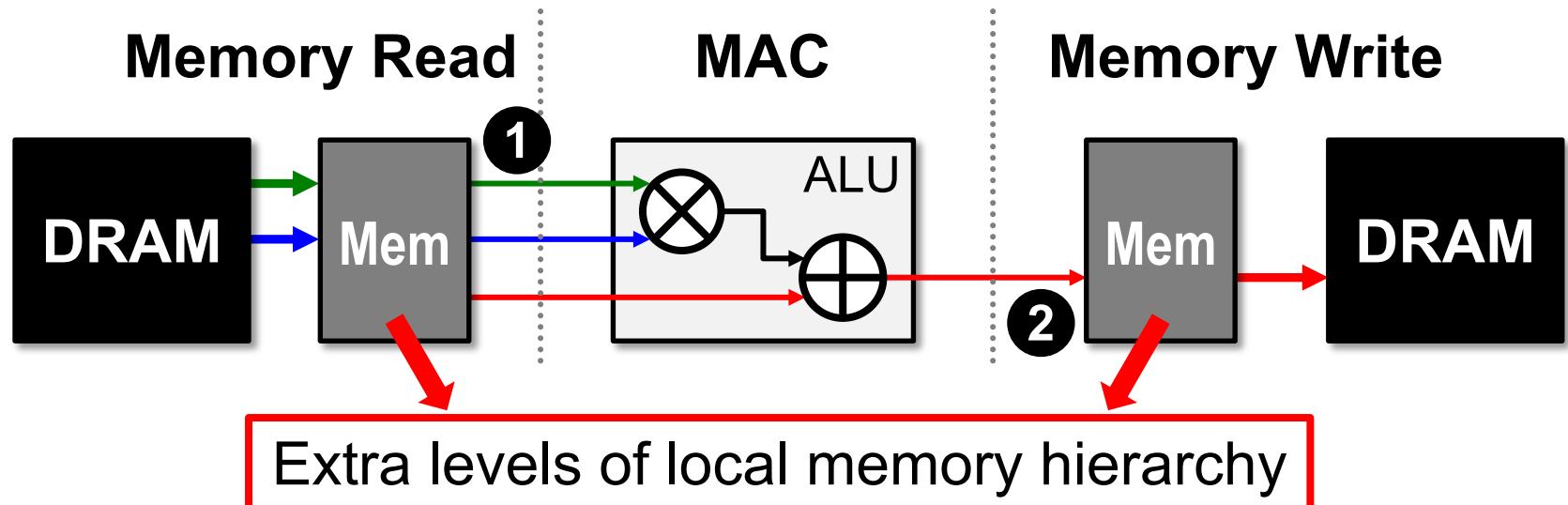
Reuse: Filter weights

Memory Access is the Bottleneck



Opportunities: **① data reuse ② local accumulation**

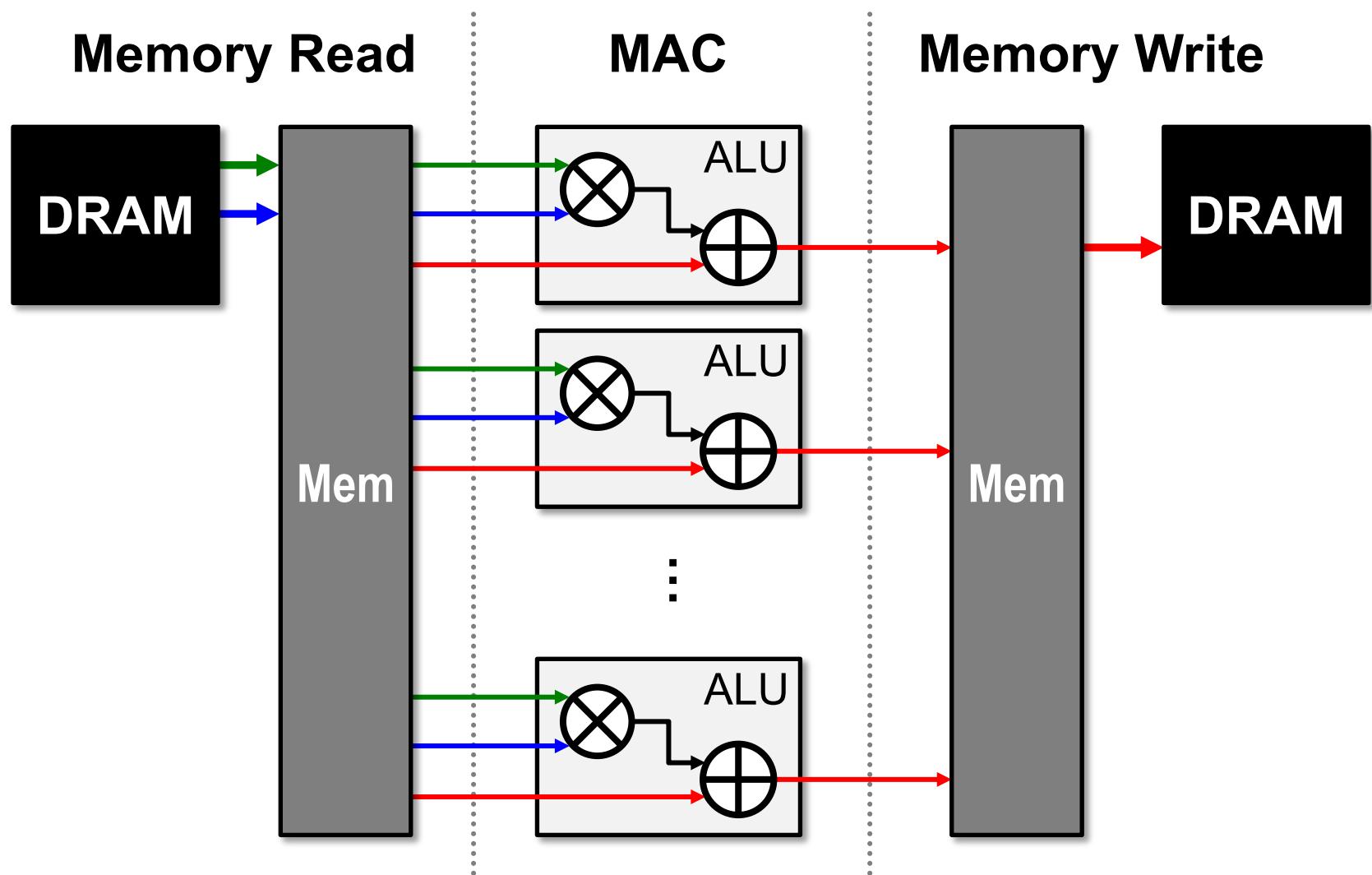
Memory Access is the Bottleneck



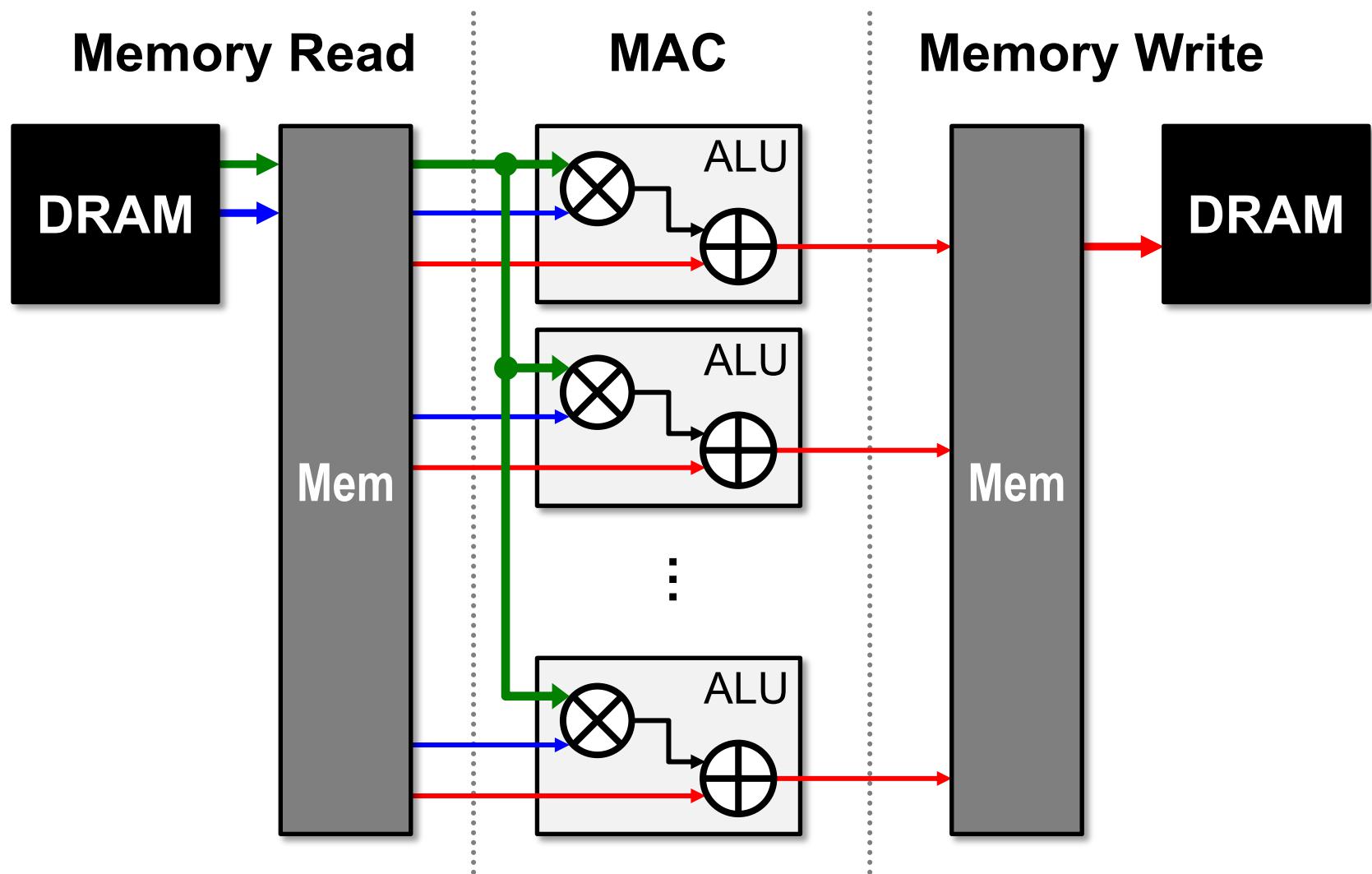
Opportunities: **① data reuse** **② local accumulation**

- Example: DRAM access in AlexNet can be reduced from **2896M** to **61M** (best case)

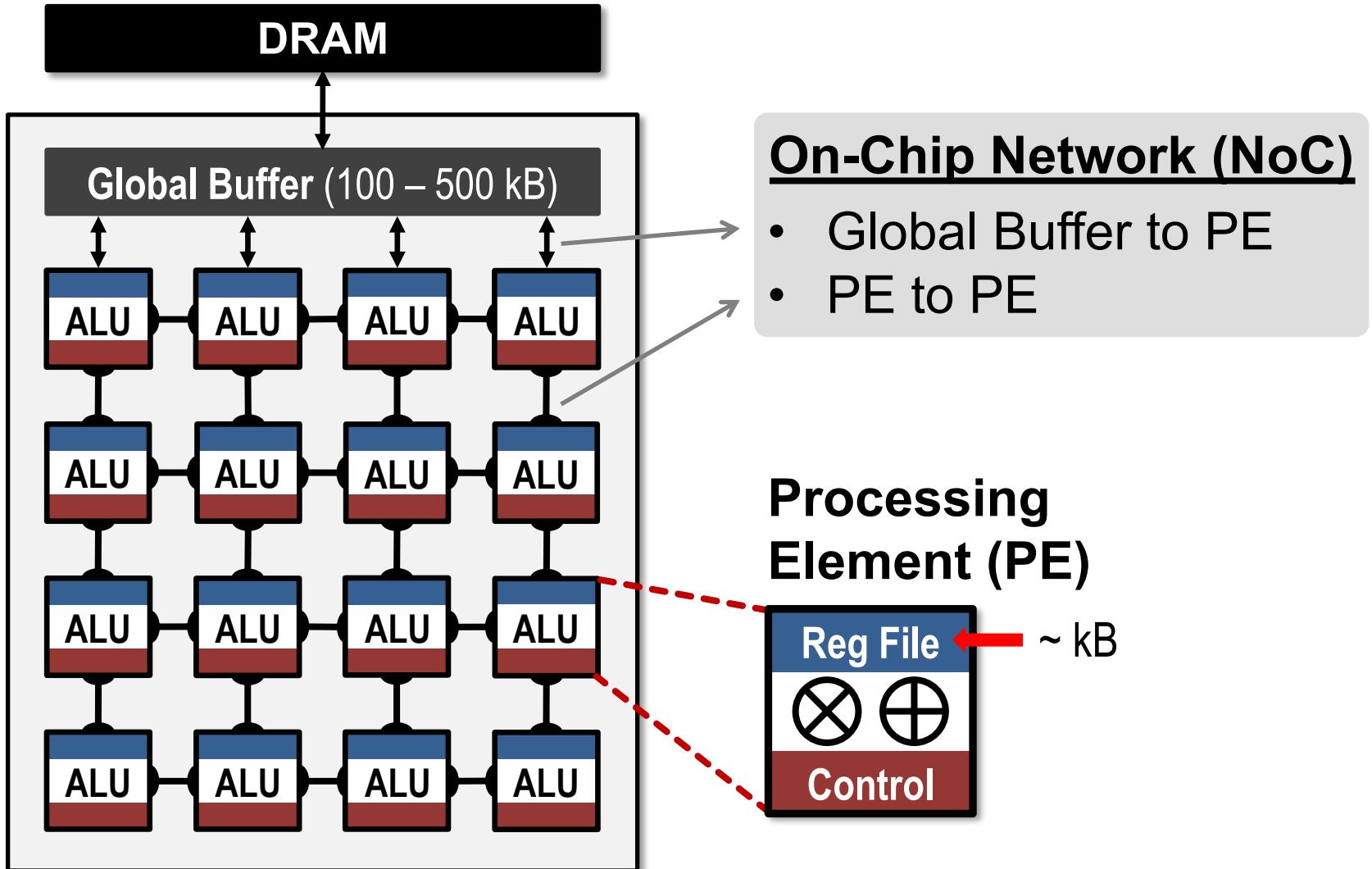
Leverage Parallelism for Higher Performance



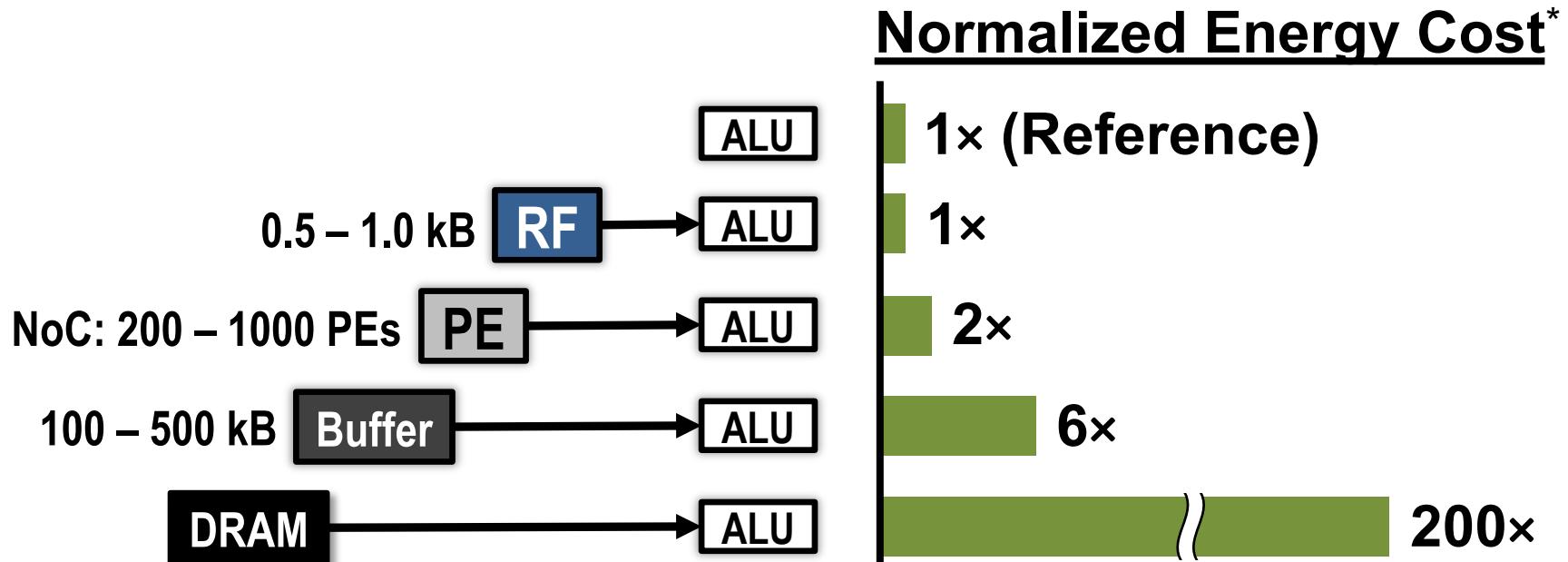
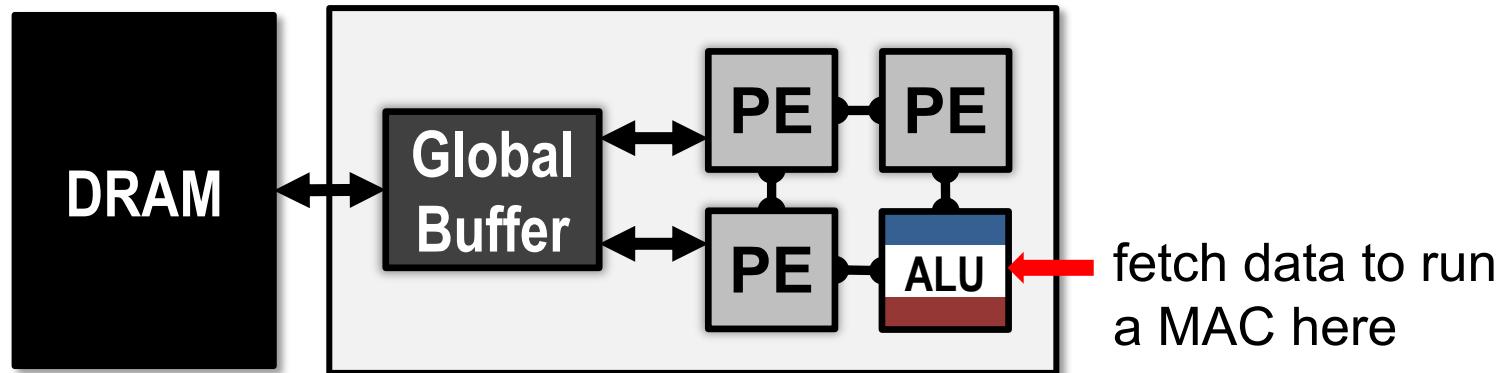
Leverage Parallelism for *Spatial* Data Reuse



Spatial Architecture



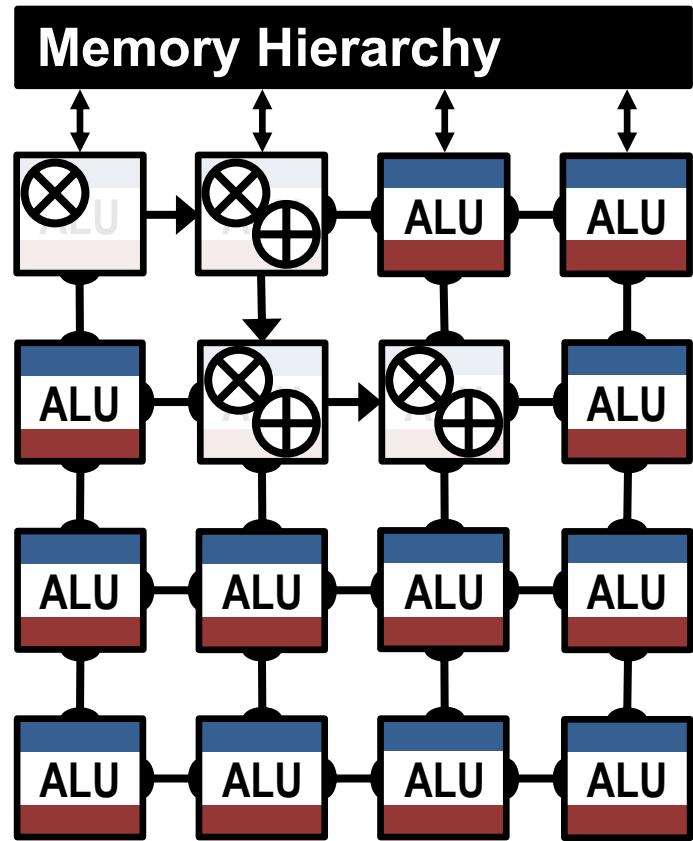
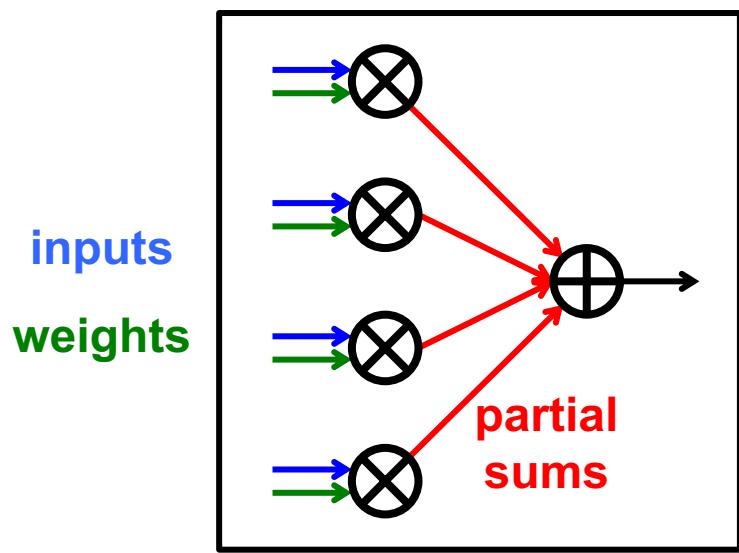
Low-Cost Local Data Access



* measured from a commercial 65nm process

How to Map the Dataflow?

CONV/FC Layer

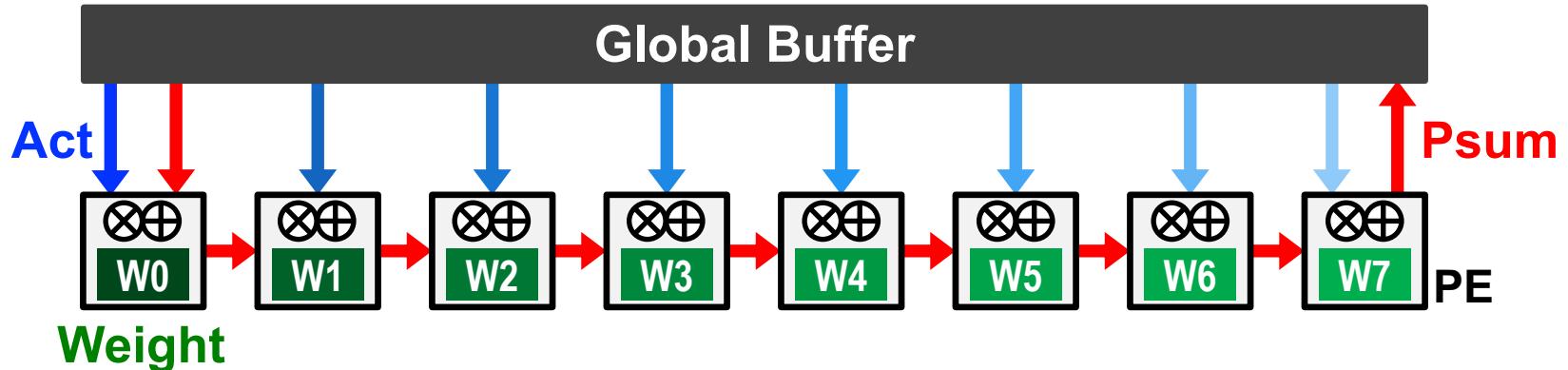


How to maximally exploit **data reuse** with the
low-cost memory hierarchy and **parallelism**?

Dataflow Taxonomy

- Weight Stationary (WS)
- Output Stationary (OS)
- Input Stationary (IS)
- No Local Reuse (NLR)

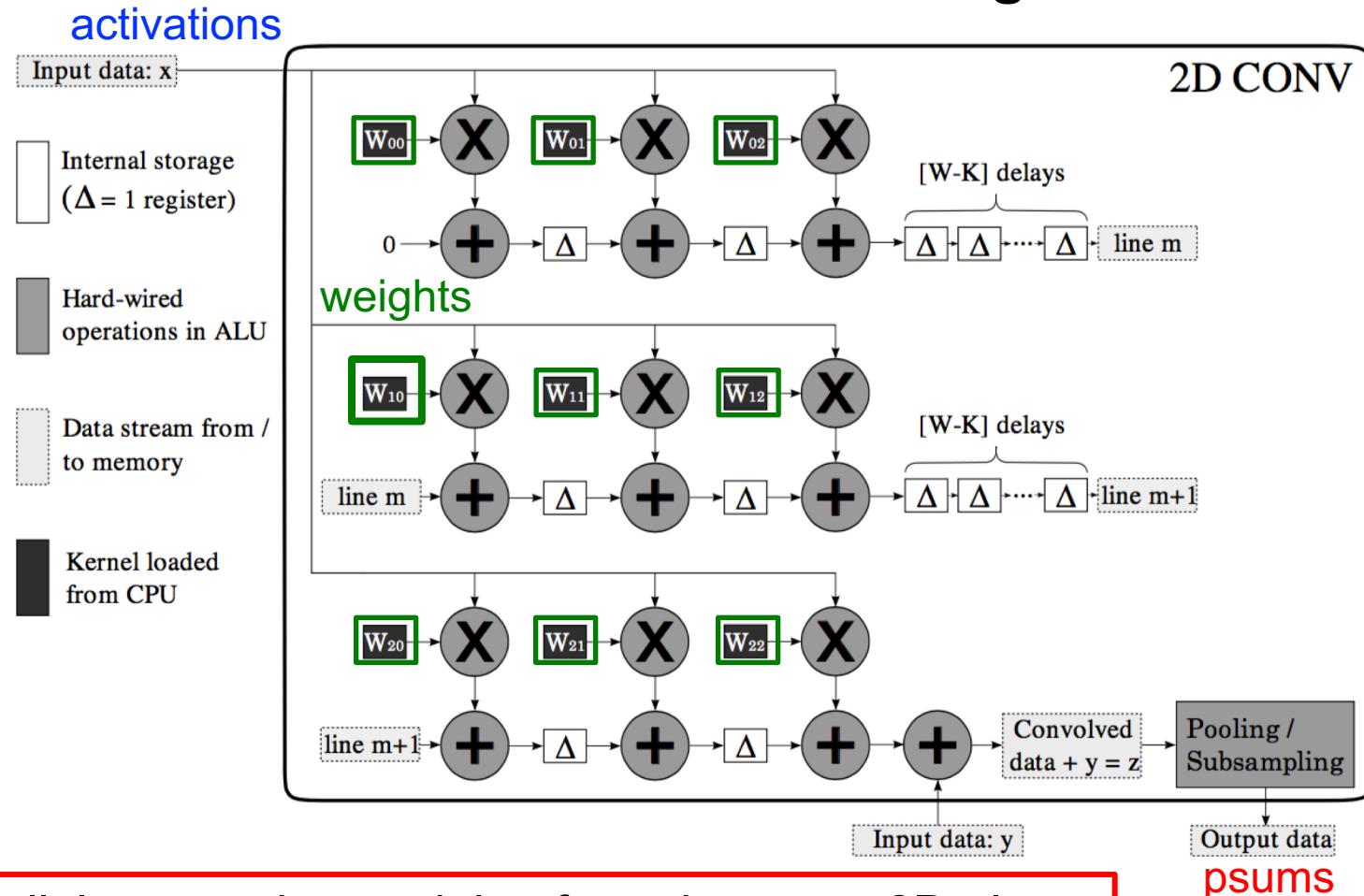
Weight Stationary (WS) Dataflow



- Minimize **weight** read energy consumption
 - maximize convolutional and filter reuse of weights
- Examples:
 - [nn-X (NeuFlow), CVPRW 2014] [Park, ISSCC 2015]
 - [Origami, GLSVLSI 2015] [Google TPU, ISCA 2017]

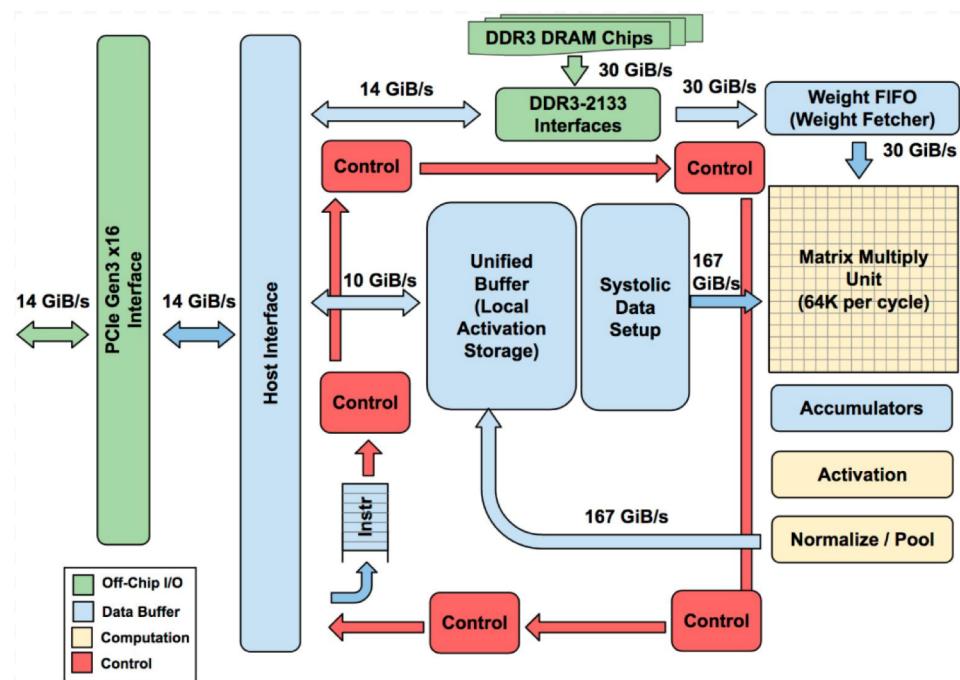
WS Example: nn-X (NeuFlow)

A 3×3 2D Convolution Engine

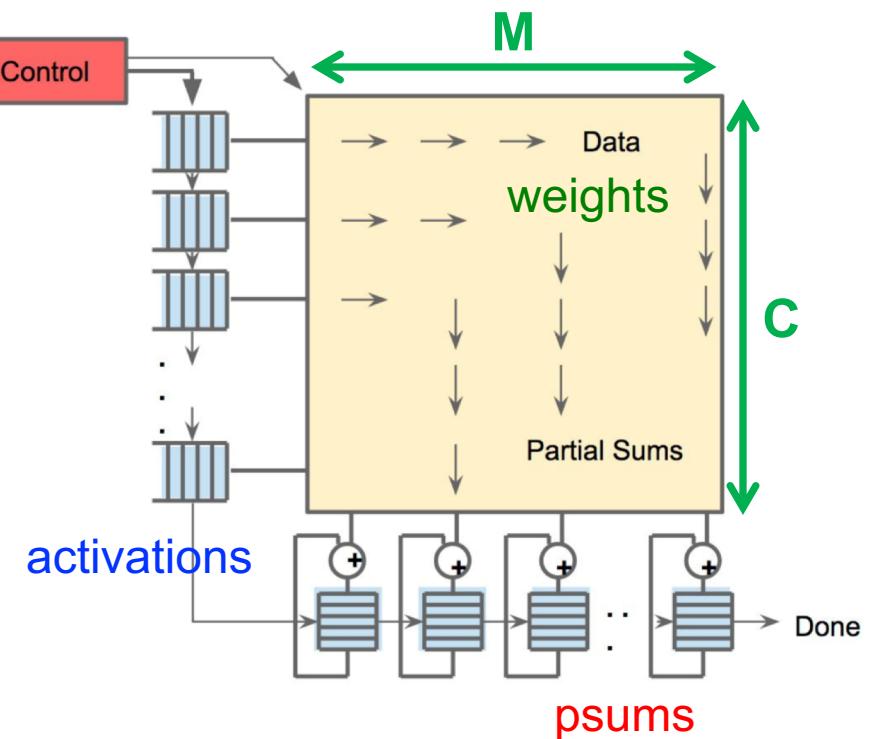


WS Example: Google TPU

Top-Level Architecture

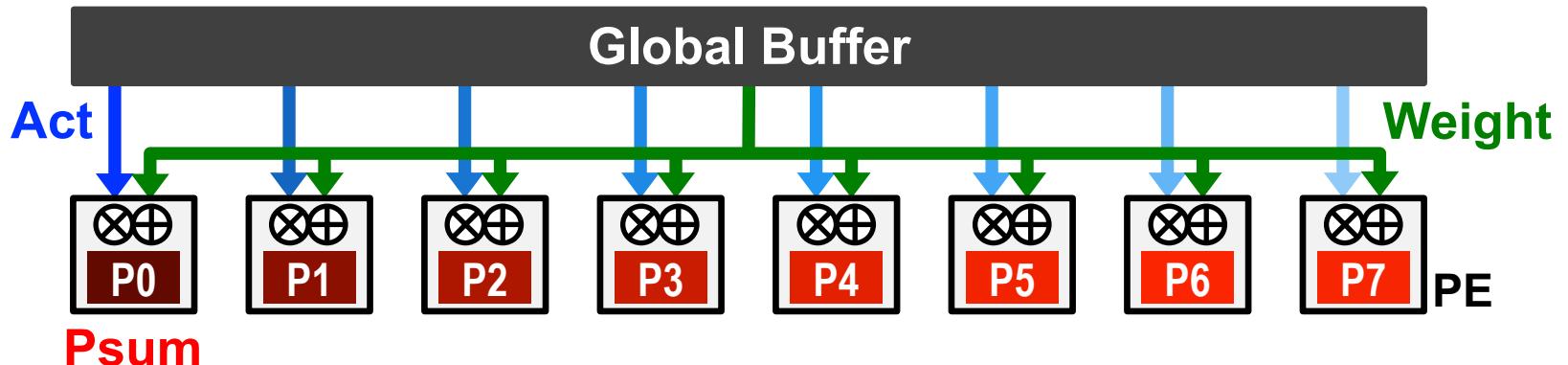


Matrix Multiply Unit



Parallel processing: weights from different 2D planes

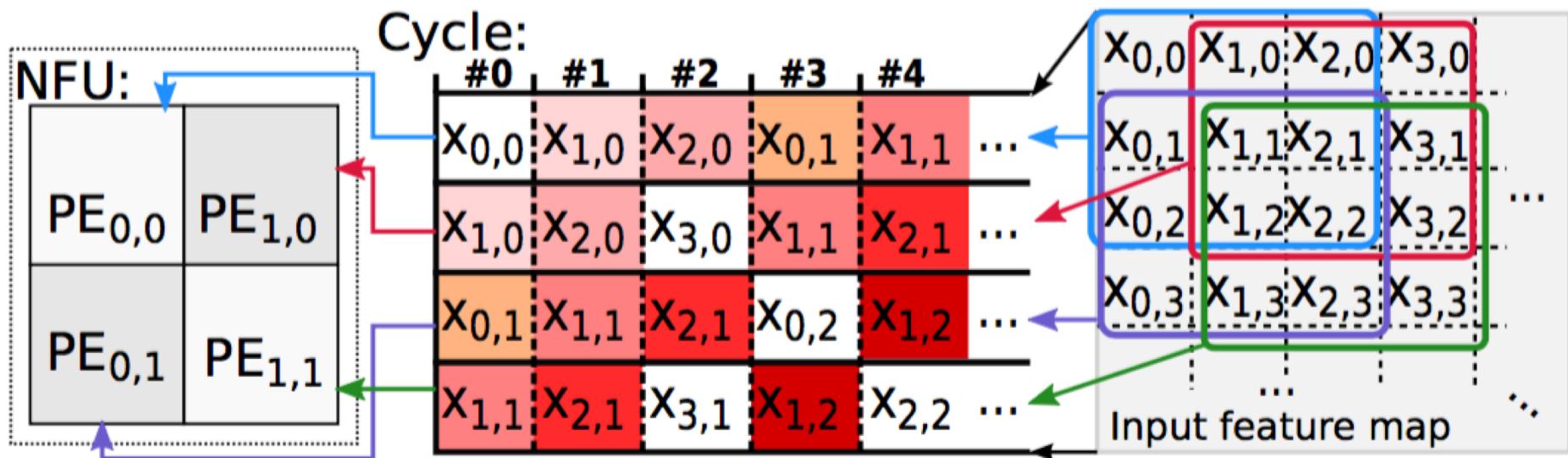
Output Stationary (OS) Dataflow



- **Minimize partial sum R/W energy consumption**
 - maximize local accumulation
- **Examples:**
 - [Gupta, ICML 2015]
 - [ShiDianNao, ISCA 2015]
 - [ENVISION, ISSCC 2017]
 - [Thinker, JSSC 2017]

OS Example: ShiDianNao

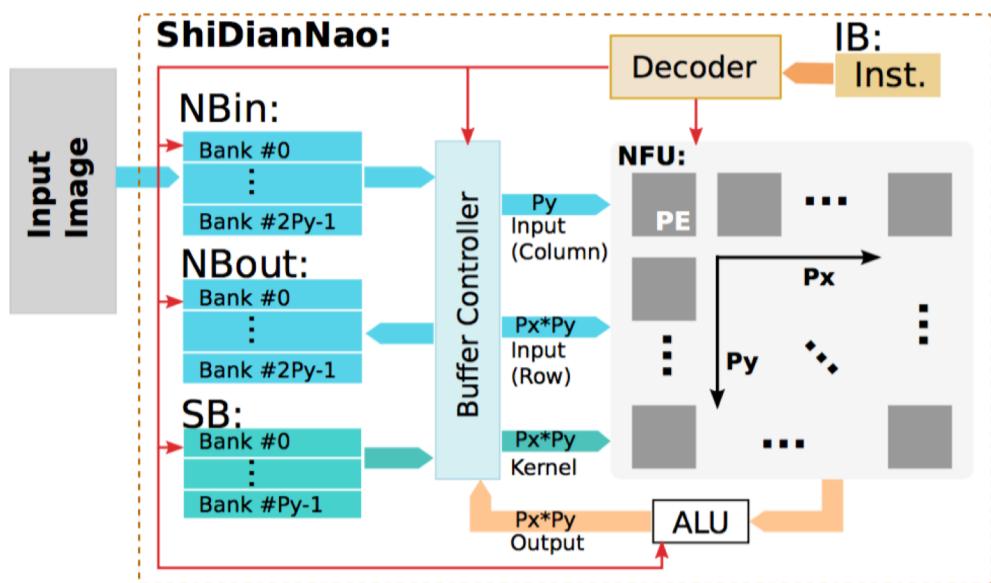
Dataflow in the PE Array



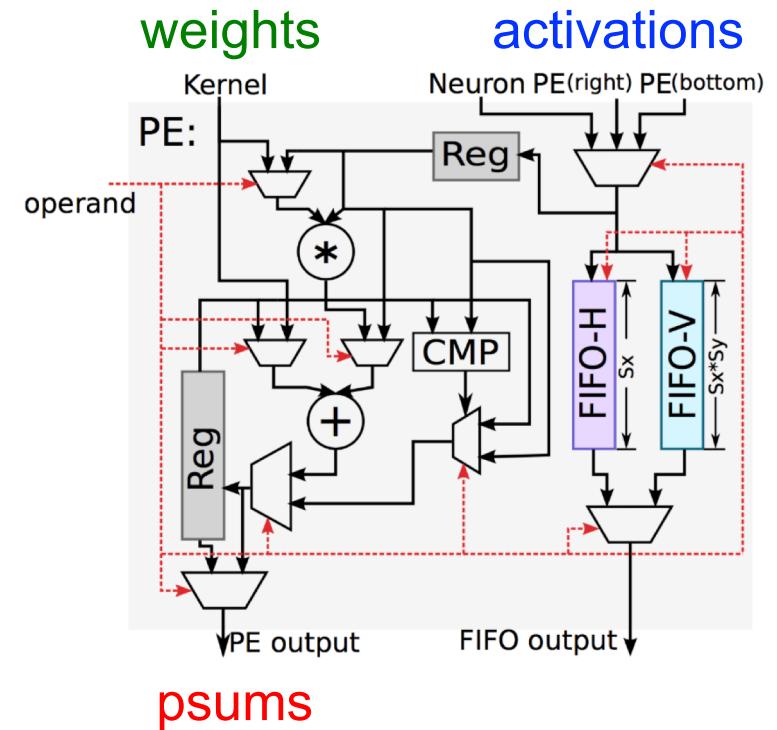
Parallel processing: many output activations from the same 2D plane

OS Example: ShiDianNao

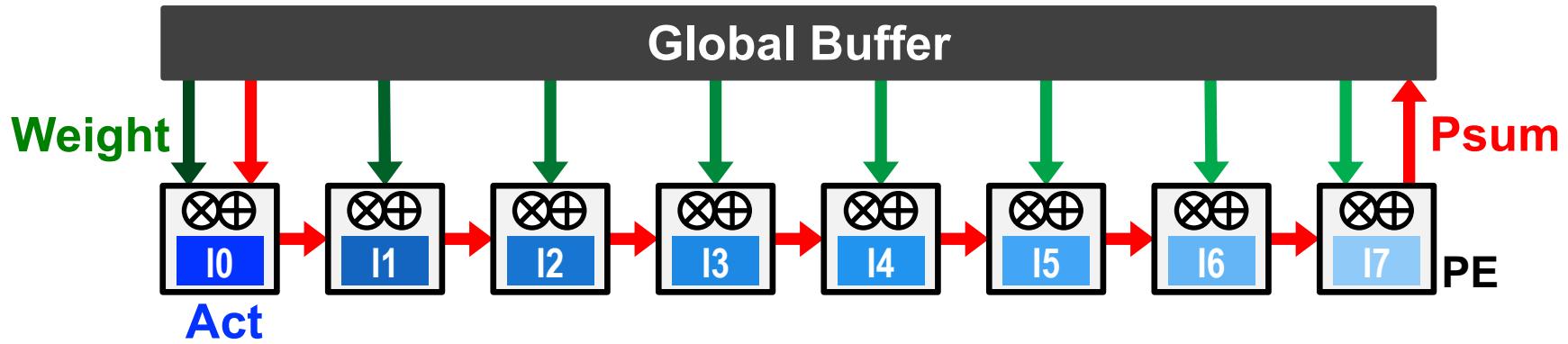
Top-Level Architecture



PE Architecture

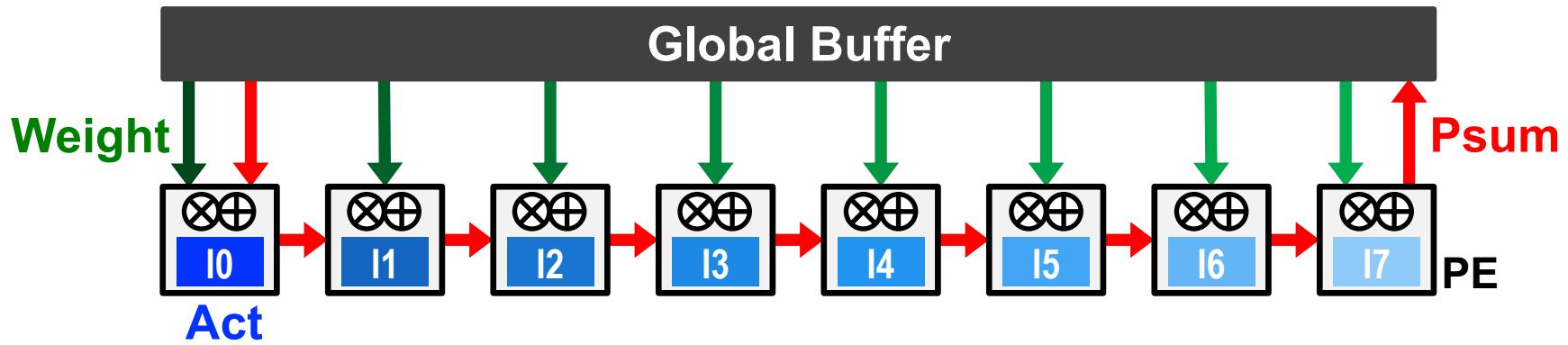


Input Stationary (IS) Dataflow



- **Minimize Input Act.** read energy consumption
 - maximize convolutional and fmap reuse of acts
- **Examples:**
 - [SCNN, /SCA 2017]

Input Stationary (IS) Dataflow

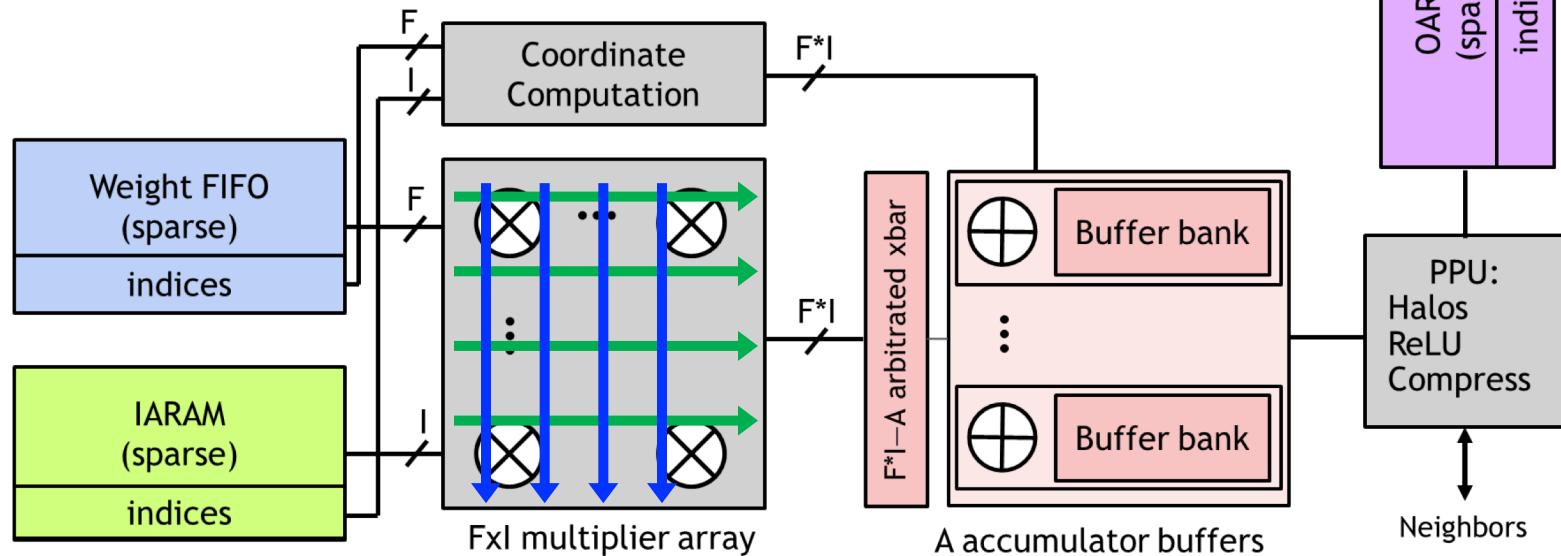


- **Minimize Input Act.** read energy consumption
 - maximize convolutional and fmap reuse of acts
- **Examples:**
 - [SCNN, /SCA 2017]

Why is it not as popular?

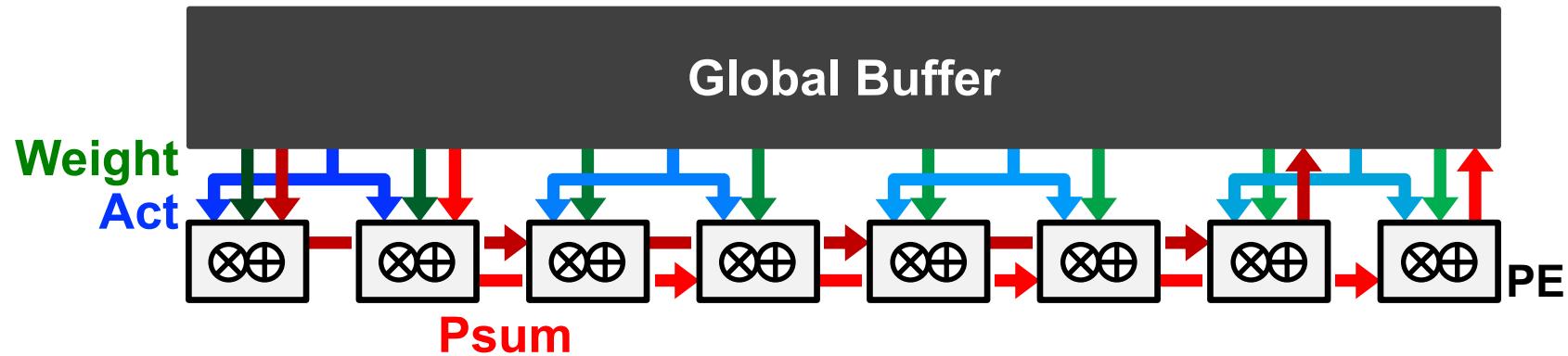
IS Example: SCNN

SCNN PE Architecture (64 PEs in SCNN)



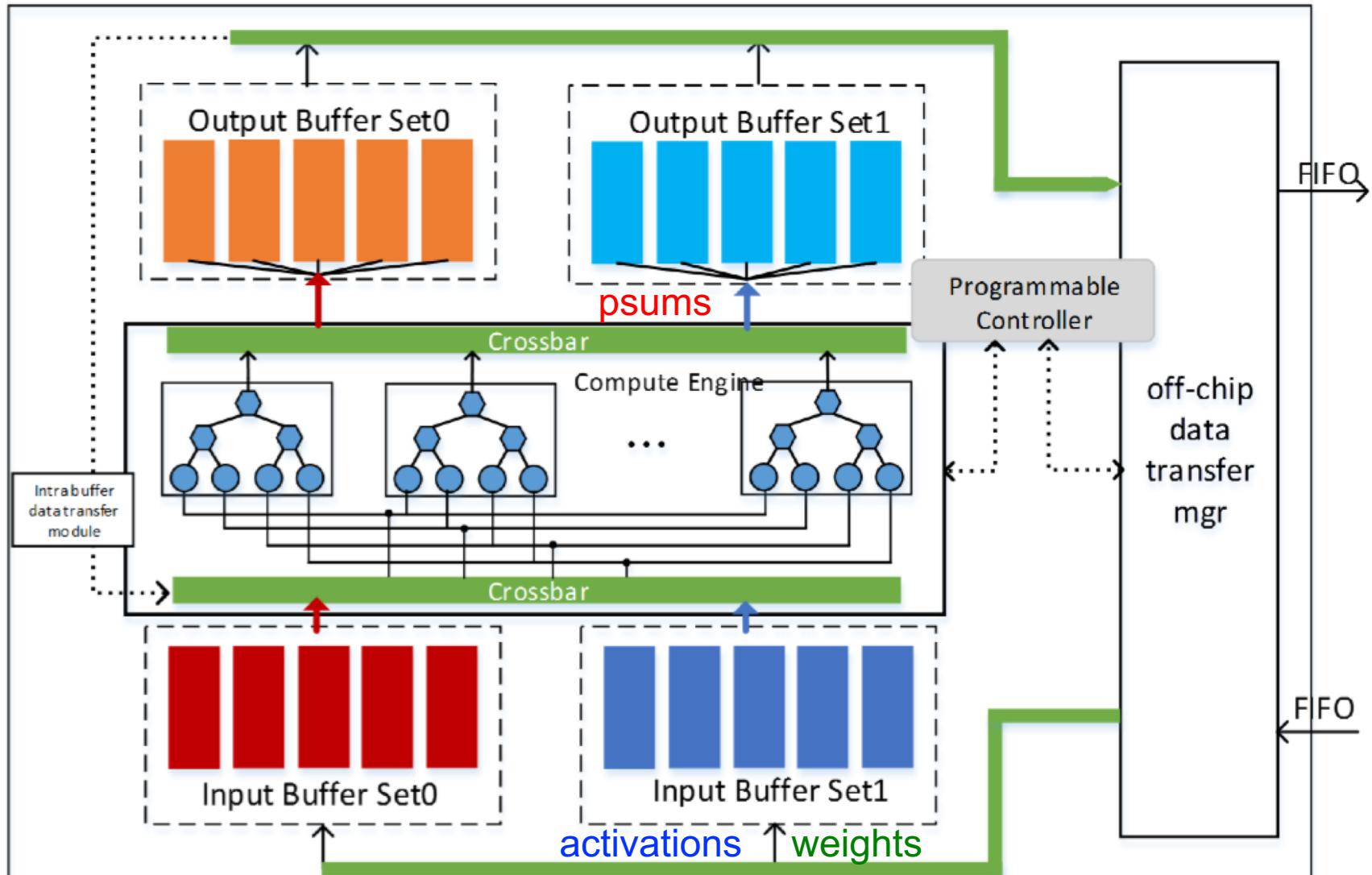
Parallel processing: many input activations from the same 2D plane

No Local Reuse (NLR) Dataflow



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

NLR Example: UCLA



Dataflow Taxonomy Summary

- **Weight Stationary**
 - minimize movement of **weights**
- **Output Stationary**
 - minimize movement of **partial sums**
- **Input Stationary**
 - minimize movement of **input activations**
- **No Local Reuse**
 - don't use any local PE storage. Maximize **global buffer size**

Dataflow Taxonomy Summary

- **Weight Stationary**
 - minimize movement of **weights**
- **Output Stationary**
 - minimize movement of **partial sums**
- **Input Stationary**
 - minimize movement of **input activations**
- **No Local Reuse**
 - don't use any local PE storage. Maximize **global buffer size**

Does this define the space of all possible dataflows?

Dataflow Taxonomy Summary

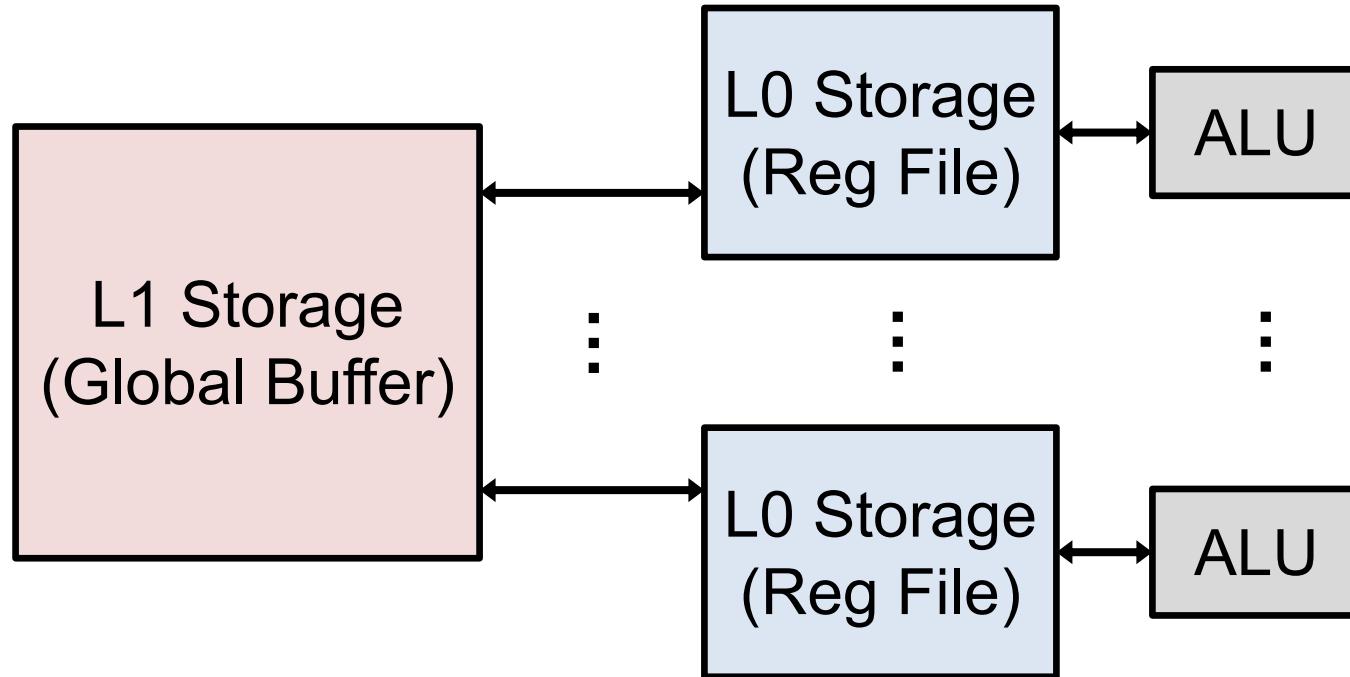
- **Weight Stationary**
 - minimize movement of **weights**
- **Output Stationary**
 - minimize movement of **partial sums**
- **Input Stationary**
 - minimize movement of **input activations**
- **No Local Reuse**
 - don't use any local PE storage. Maximize **global buffer size**

Does this define the space of all possible dataflows?

Not at all!

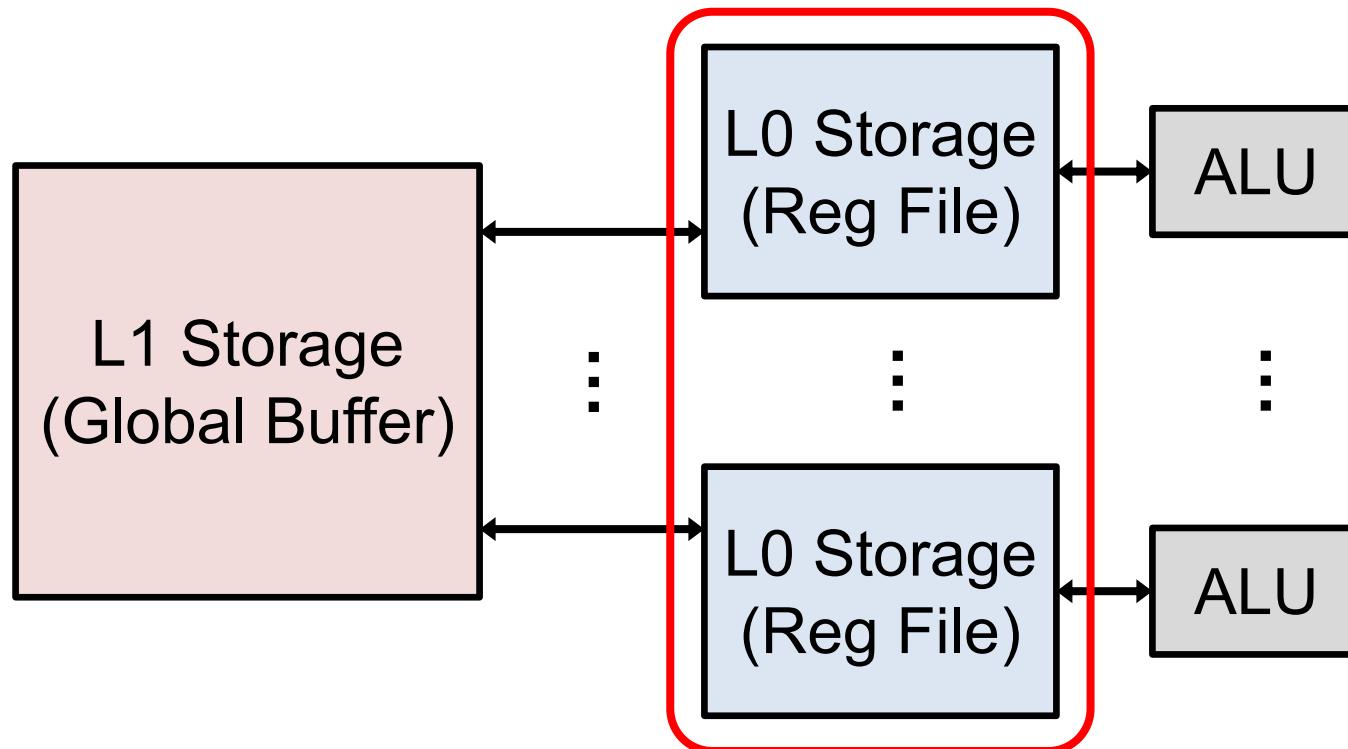
Can we have multiple **stationariness** in an architecture?

Multi-level Storage Hierarchy vs. Dataflow



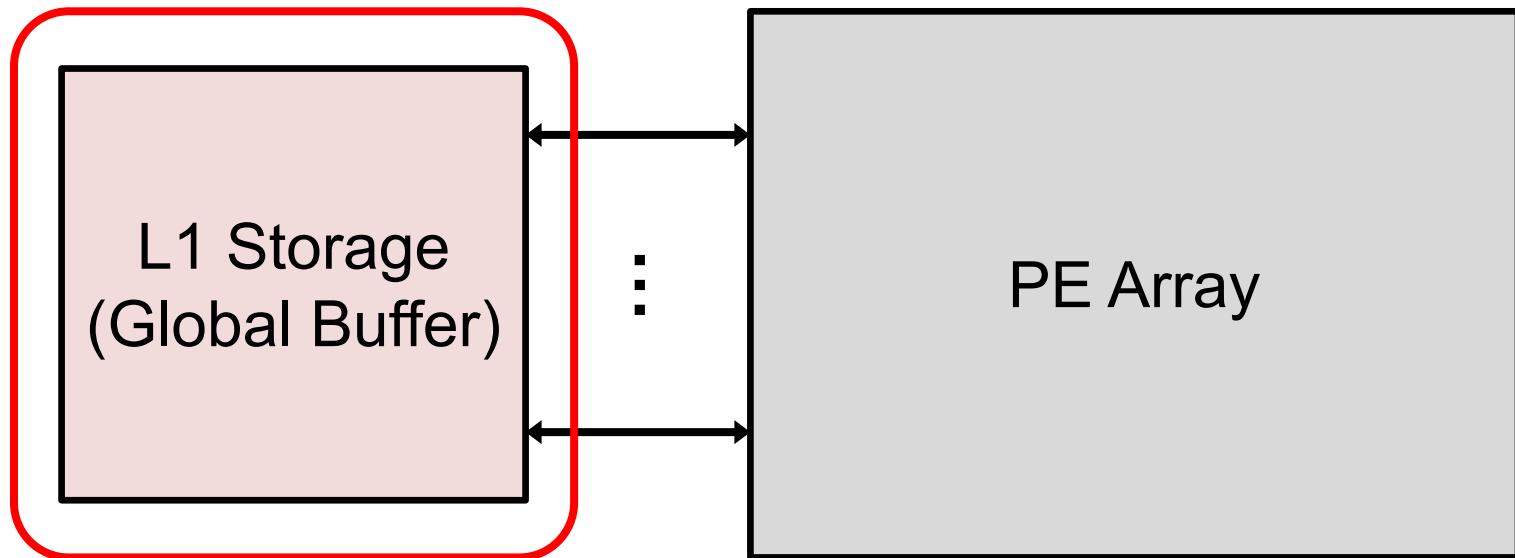
Multi-level Storage Hierarchy vs. Dataflow

Weight Stationary



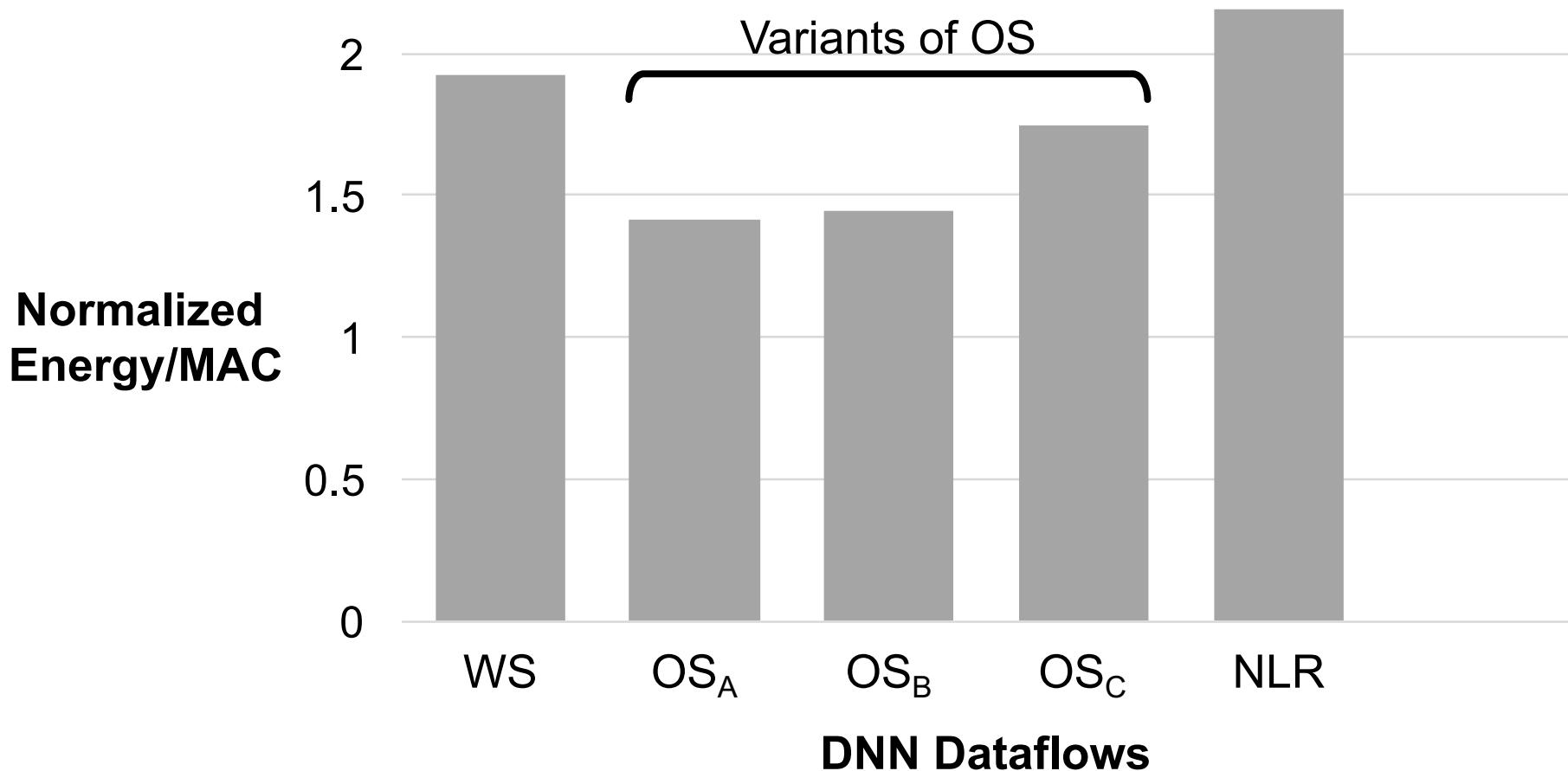
Multi-level Storage Hierarchy vs. Dataflow

Output Stationary



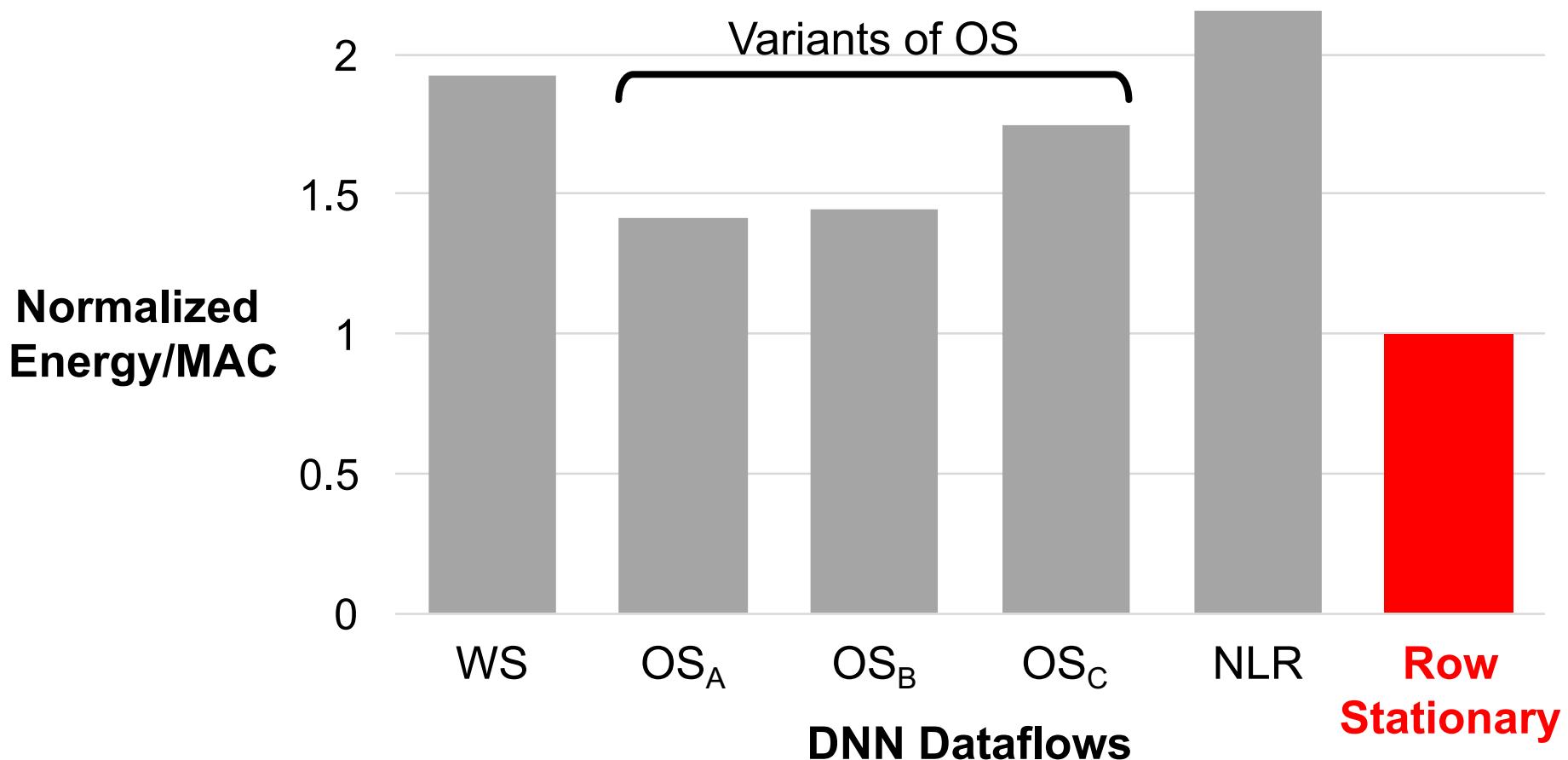
Energy Efficiency Comparison

- Same total area
- AlexNet CONV layers
- 256 PEs
- Batch size = 16



Energy Efficiency Comparison

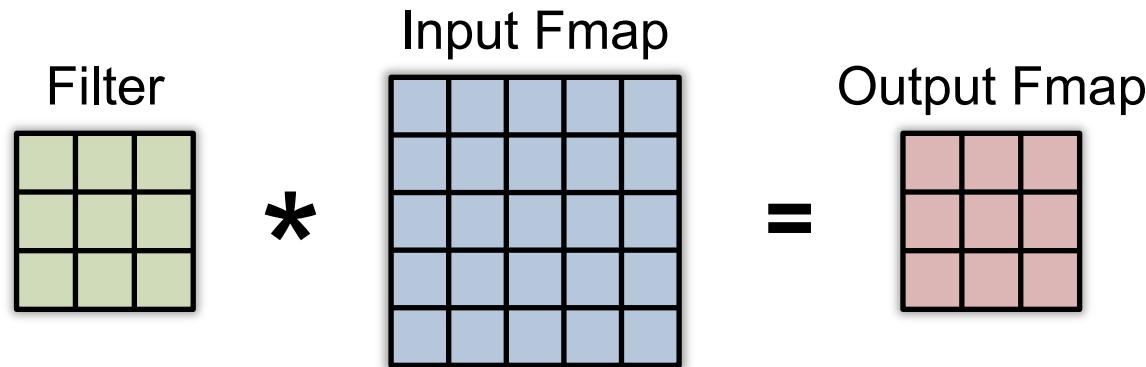
- Same total area
- AlexNet CONV layers
- 256 PEs
- Batch size = 16



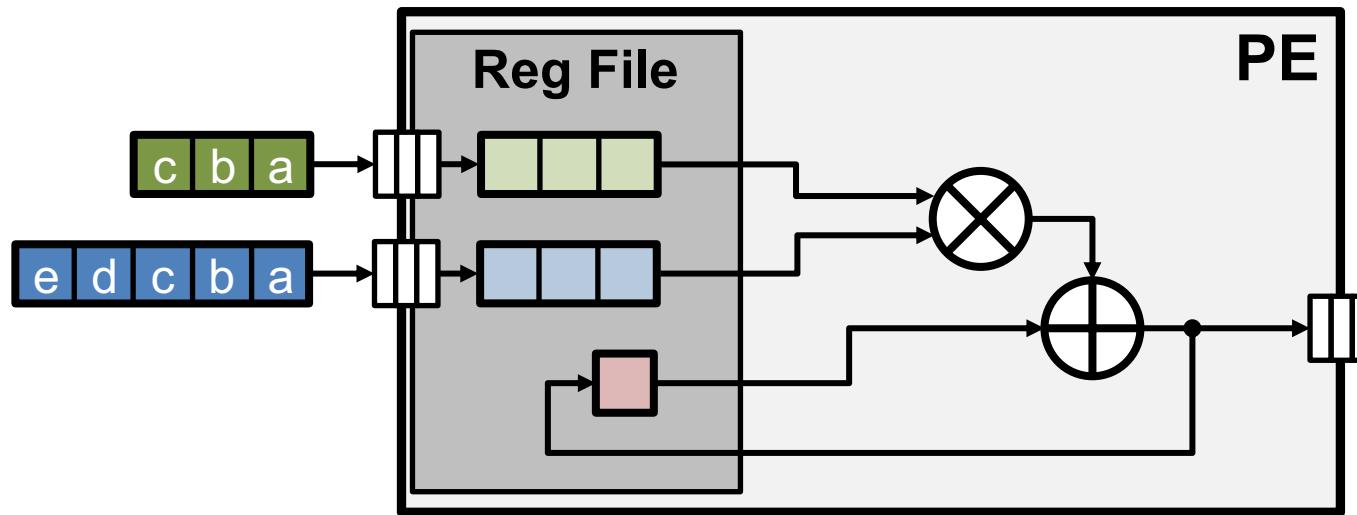
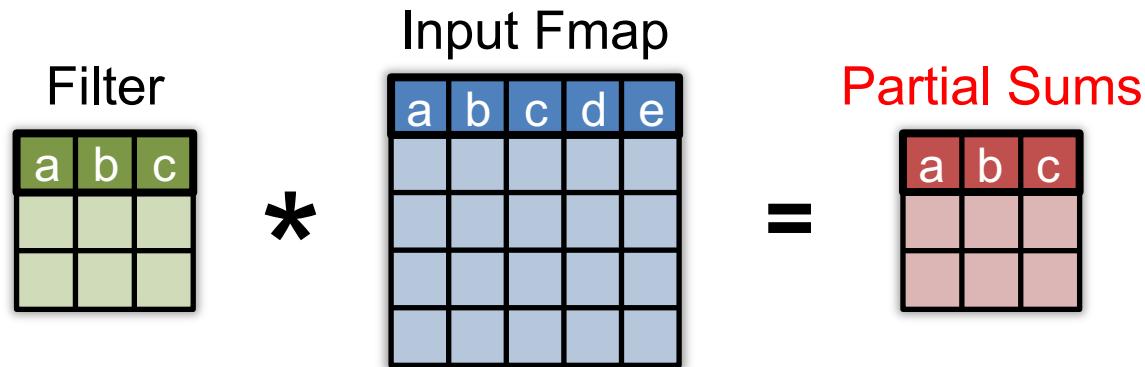
Energy-Efficient Dataflow: Row Stationary (RS)

- **Maximize** reuse and accumulation at **RF**
- Optimize for **overall** energy efficiency instead for *only* a certain data type

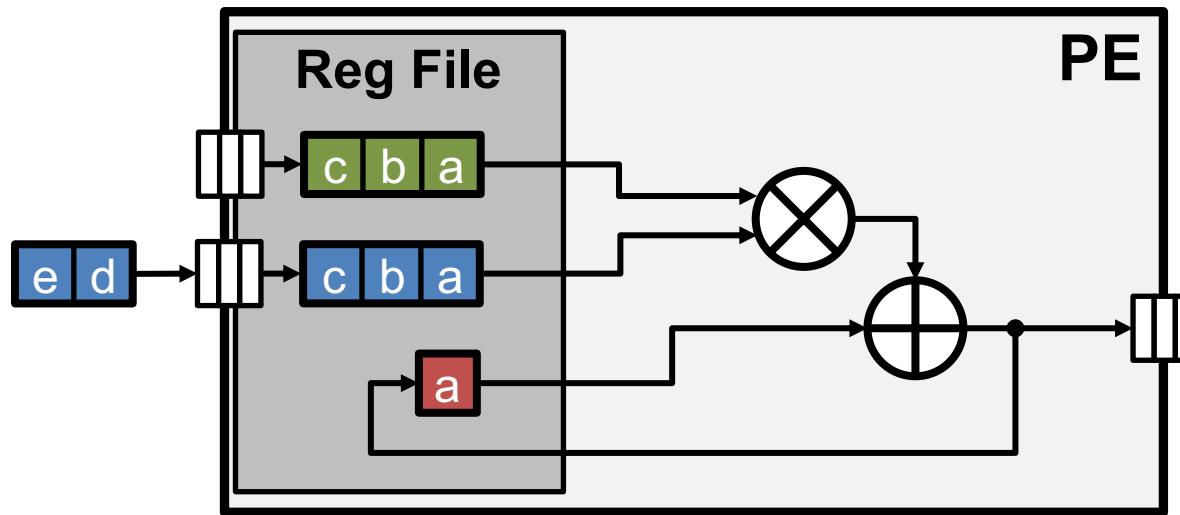
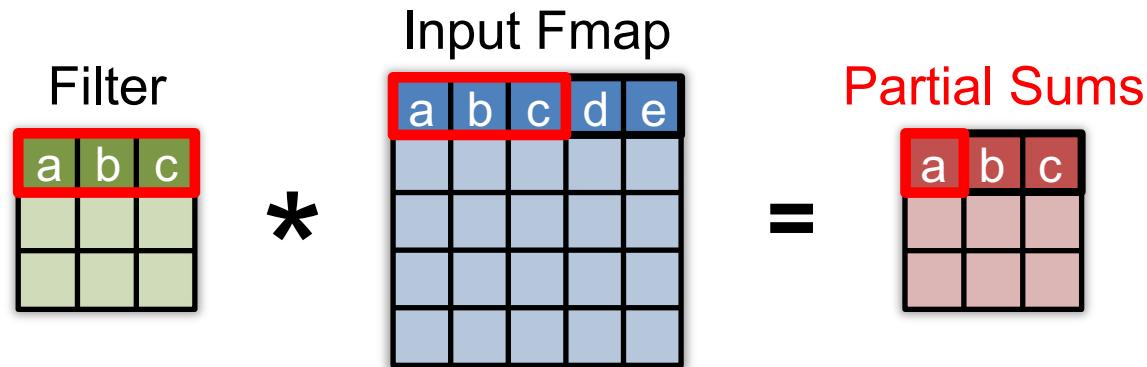
1D Row Convolution in PE



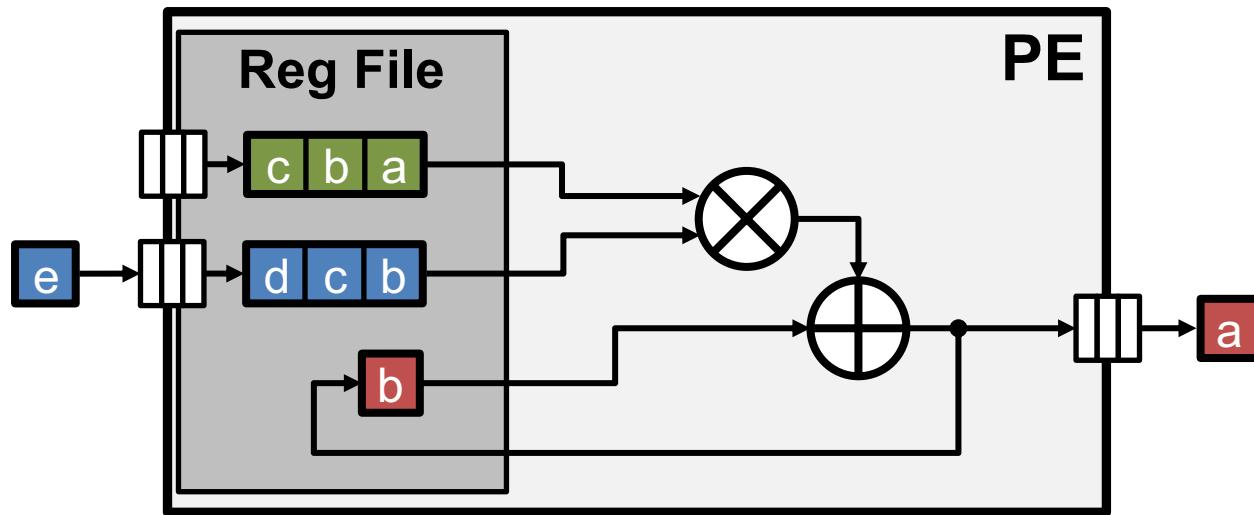
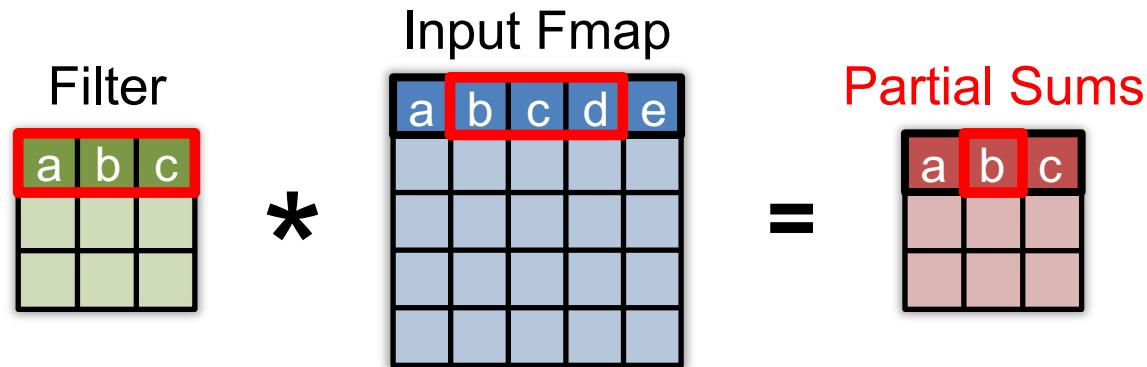
1D Row Convolution in PE



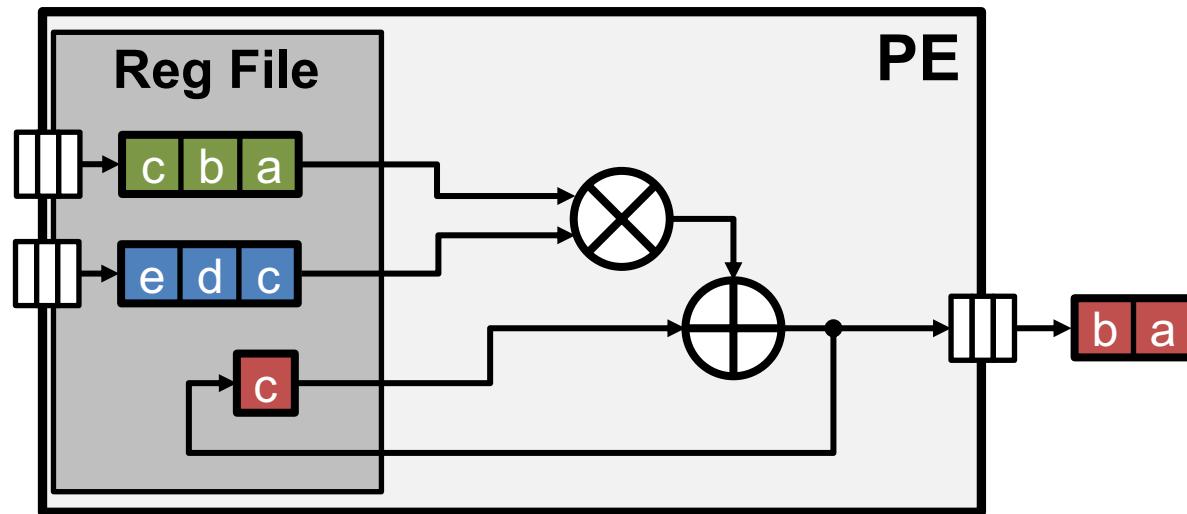
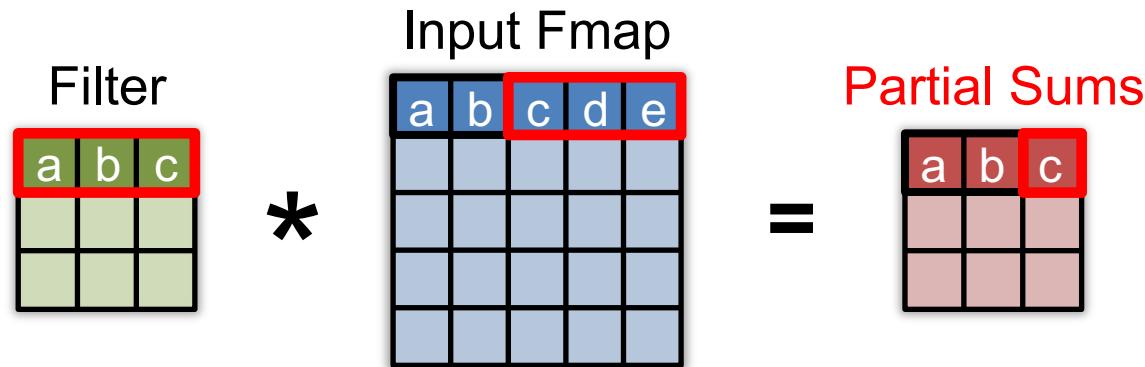
1D Row Convolution in PE



1D Row Convolution in PE

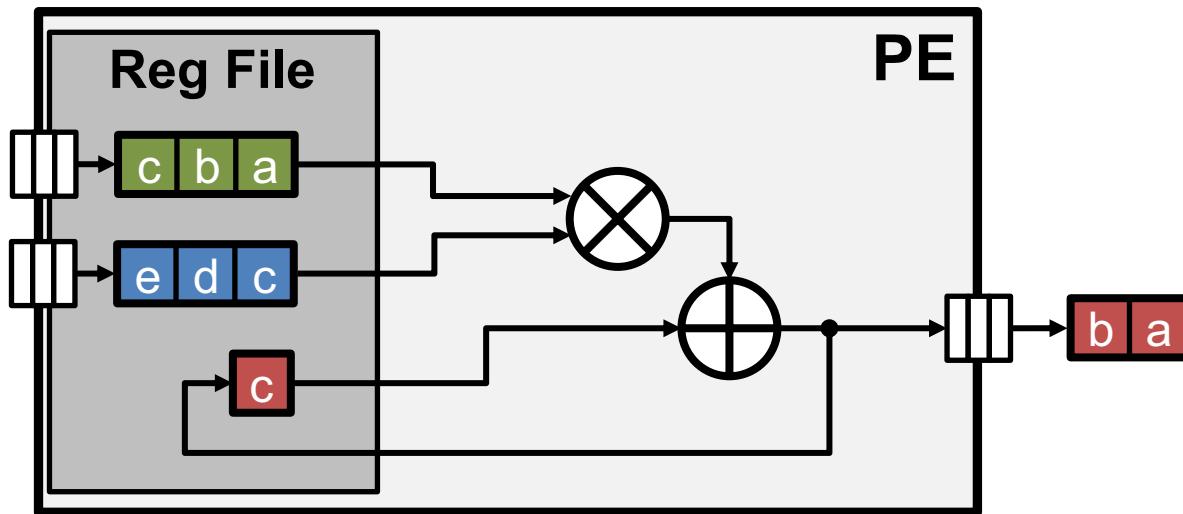


1D Row Convolution in PE



1D Row Convolution in PE

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

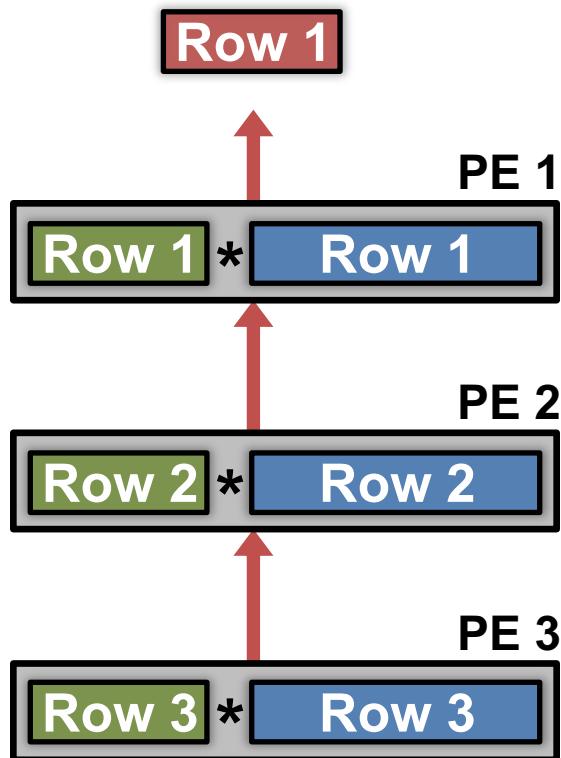


2D Convolution in PE Array



$$\begin{array}{c} \text{Icon of a 2x2 matrix with green and yellow cells} \\ * \end{array} = \begin{array}{c} \text{Icon of a 2x2 matrix with red and pink cells} \end{array}$$

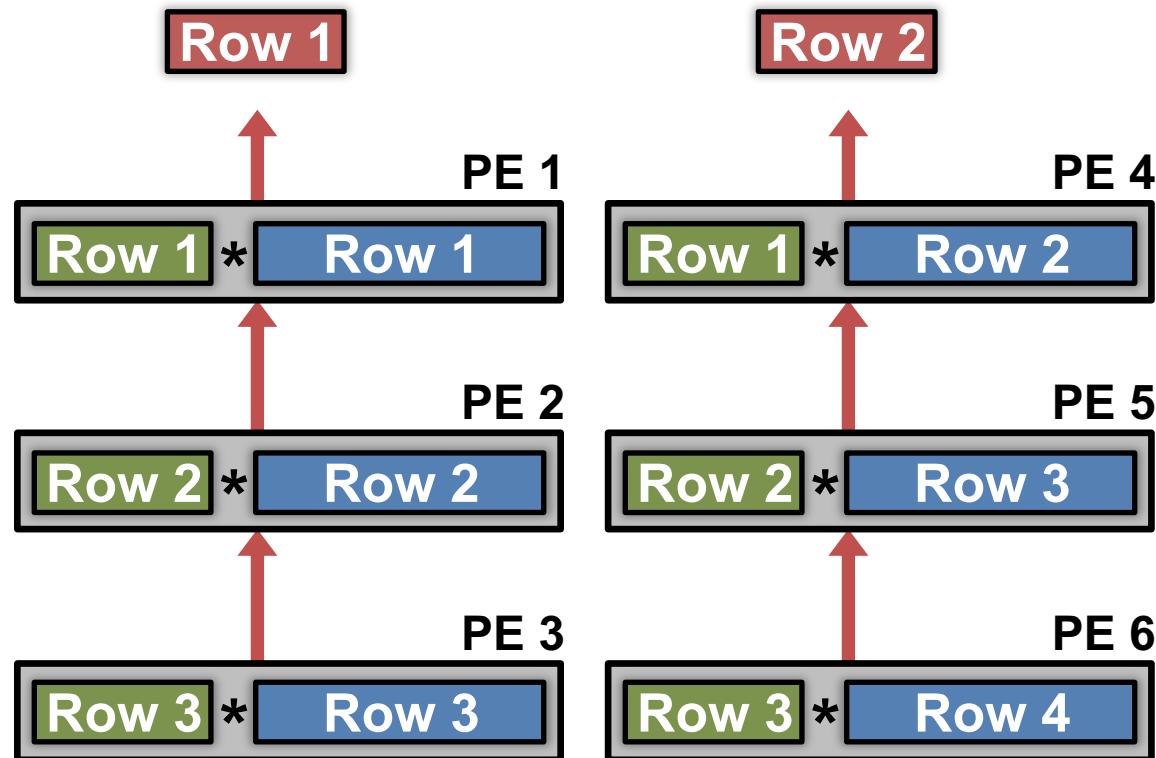
2D Convolution in PE Array



$$\begin{matrix} \text{Green Grid} \\ * \end{matrix} = \begin{matrix} \text{Red Grid} \end{matrix}$$

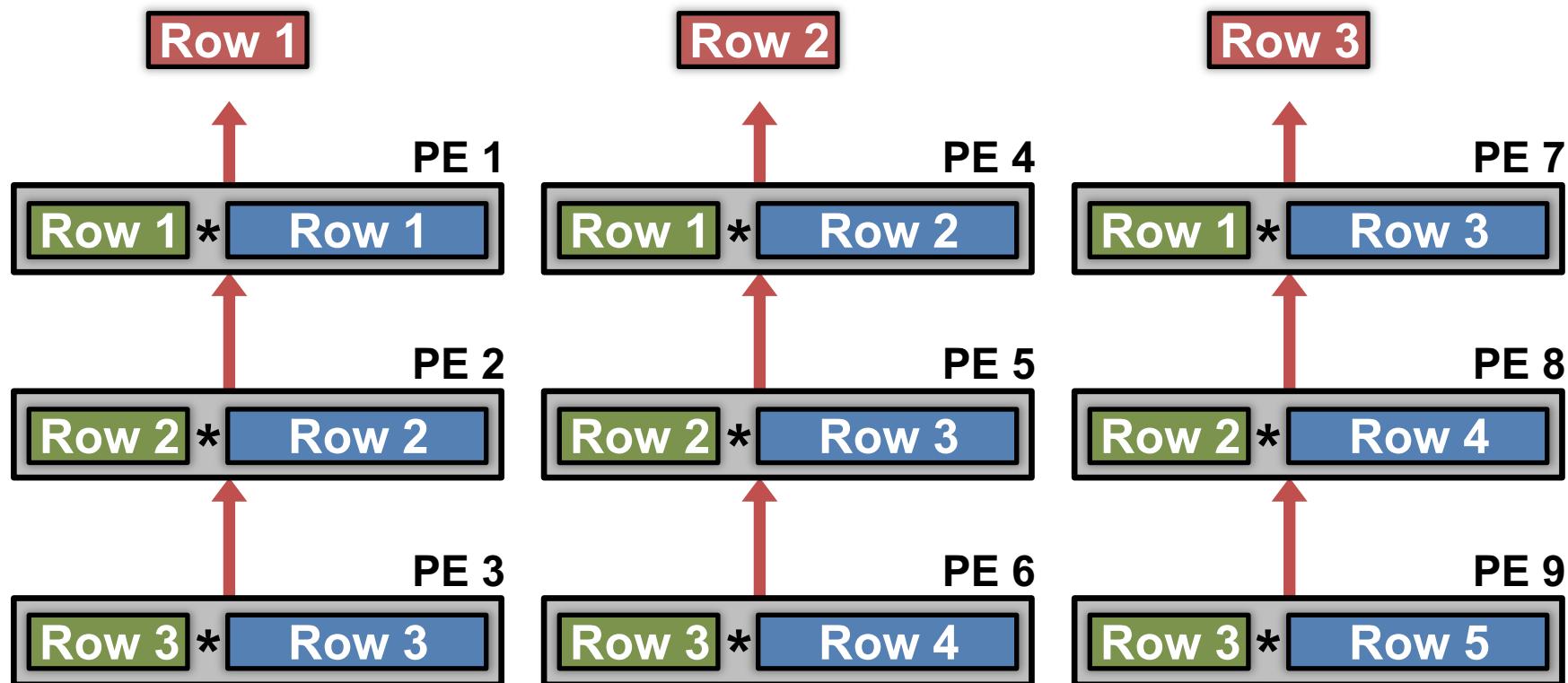
A diagram showing a 2D convolution operation. On the left is a green 3x3 grid. In the center is a multiplication symbol (*). To the right is an equals sign (=). On the far right is a red 3x3 grid.

2D Convolution in PE Array



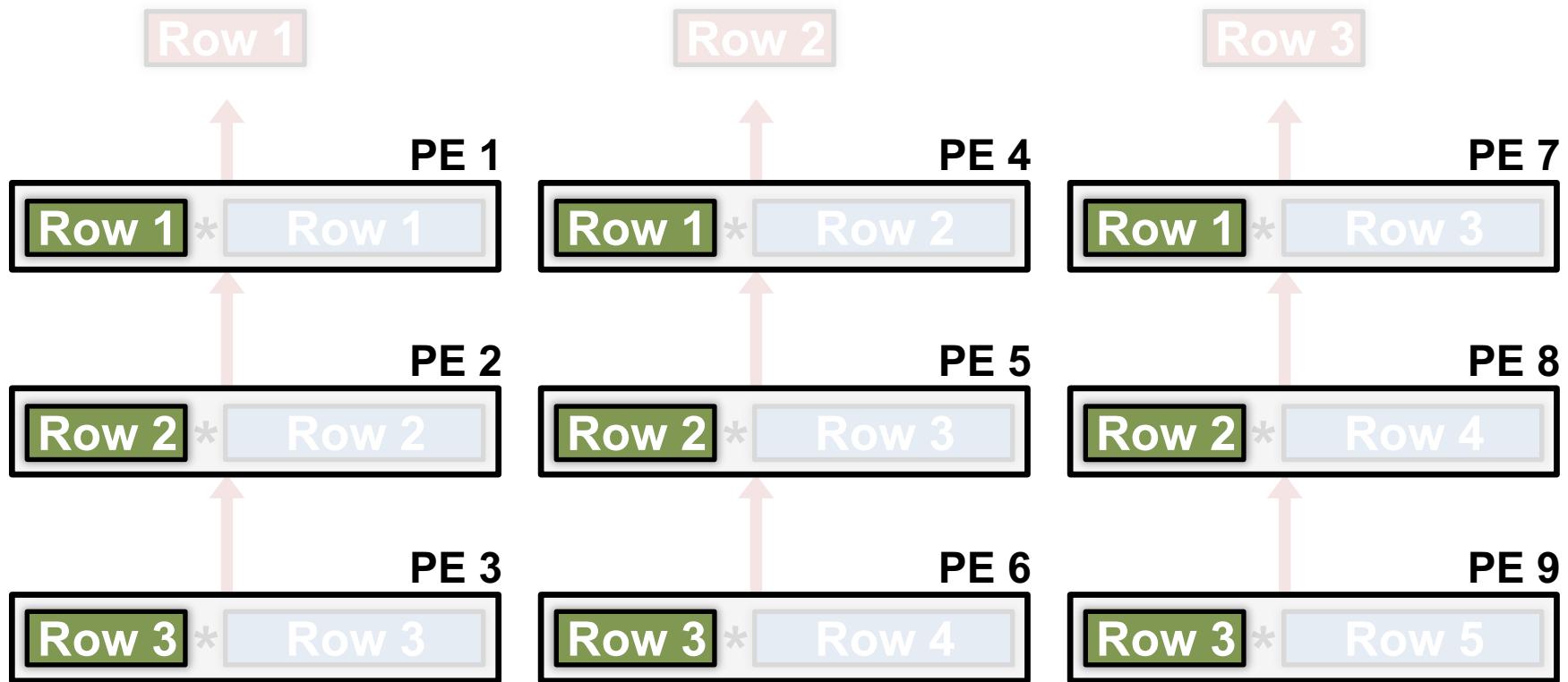
$$\begin{array}{c} \text{[Green]} \\ \times \end{array} \begin{array}{c} \text{[Blue]} \\ \times \end{array} = \begin{array}{c} \text{[Red]} \end{array} \quad \begin{array}{c} \text{[Green]} \\ \times \end{array} \begin{array}{c} \text{[Blue]} \\ \times \end{array} = \begin{array}{c} \text{[Red]} \end{array}$$

2D Convolution in PE Array



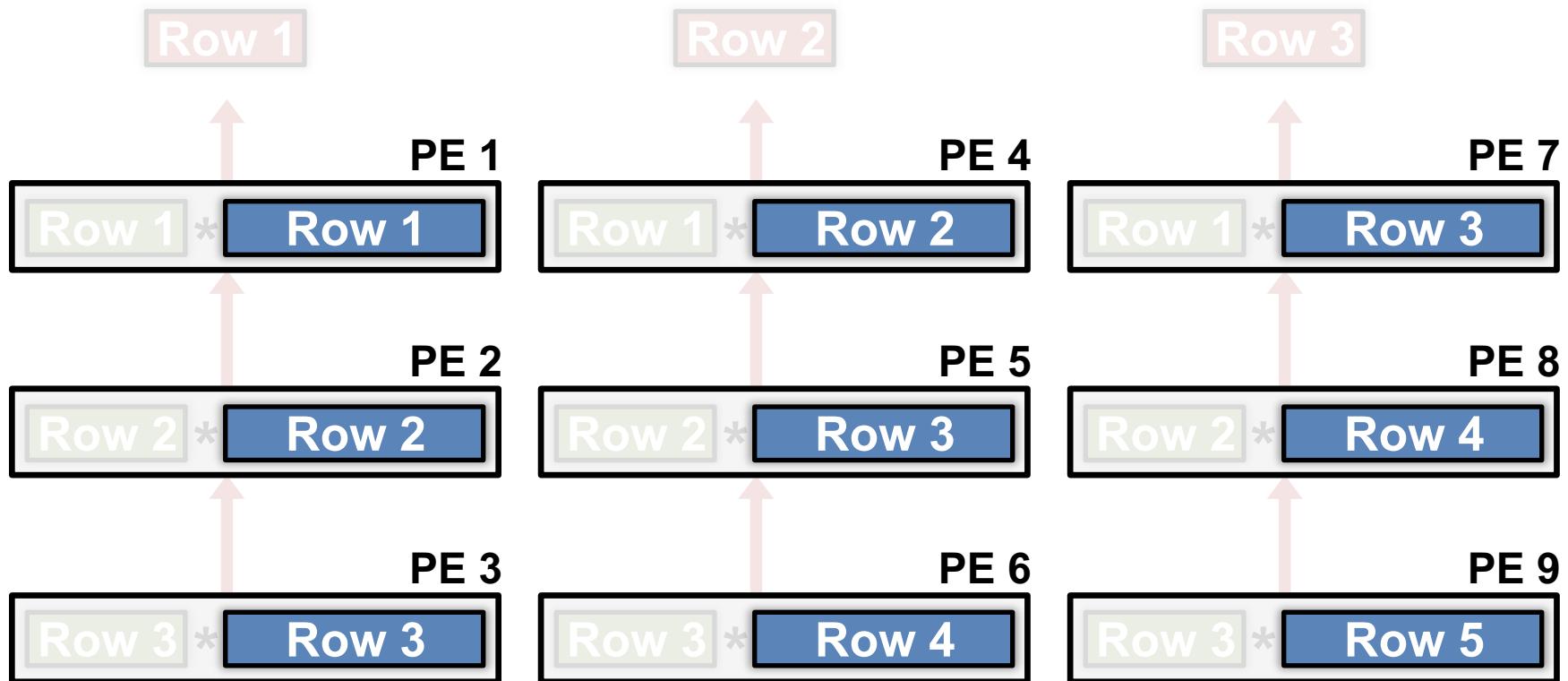
$$\begin{array}{c} \text{Green Grid} \\ \times \end{array} = \begin{array}{c} \text{Red Grid} \end{array}$$
$$\begin{array}{c} \text{Green Grid} \\ \times \end{array} = \begin{array}{c} \text{Pink Grid} \end{array}$$
$$\begin{array}{c} \text{Green Grid} \\ \times \end{array} = \begin{array}{c} \text{Pink Grid} \end{array}$$

Convolutional Reuse Maximized



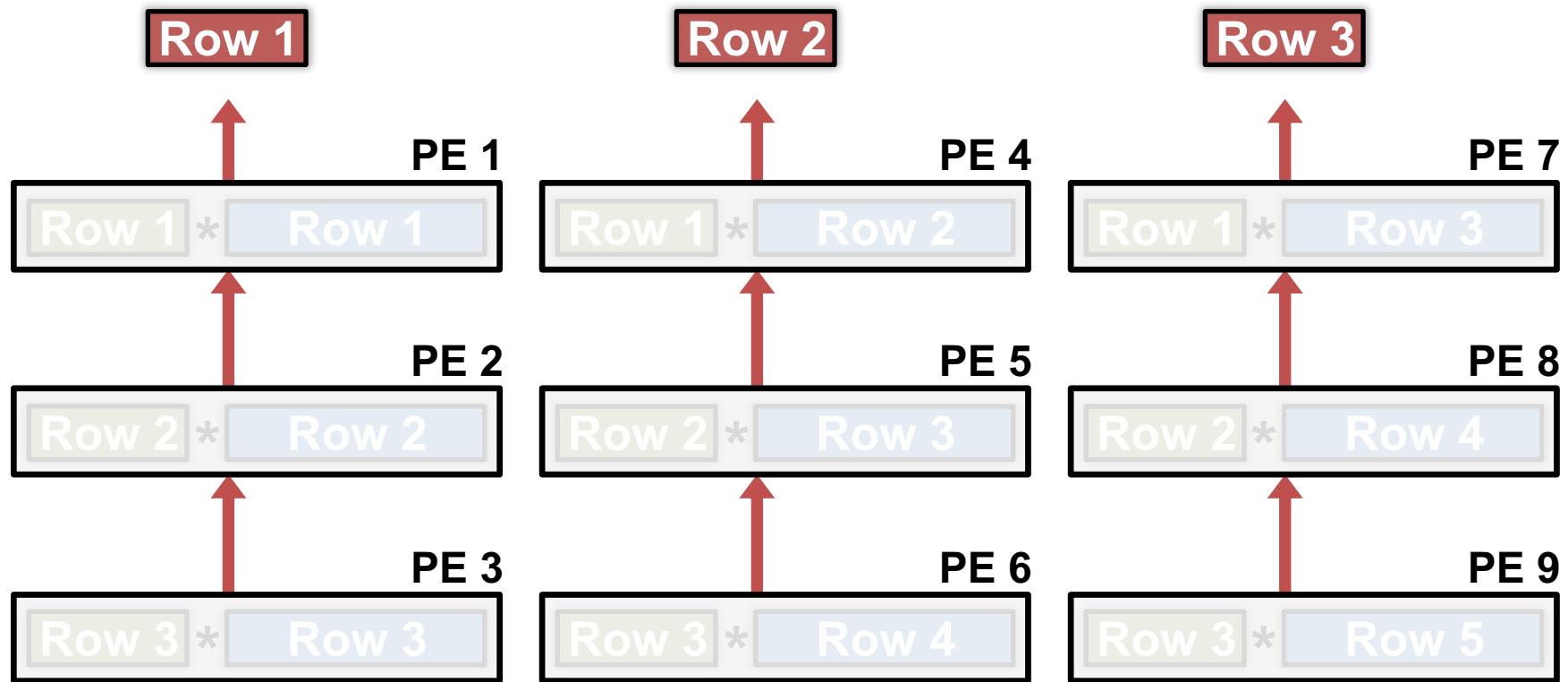
Filter rows are reused across PEs **horizontally**

Convolutional Reuse Maximized



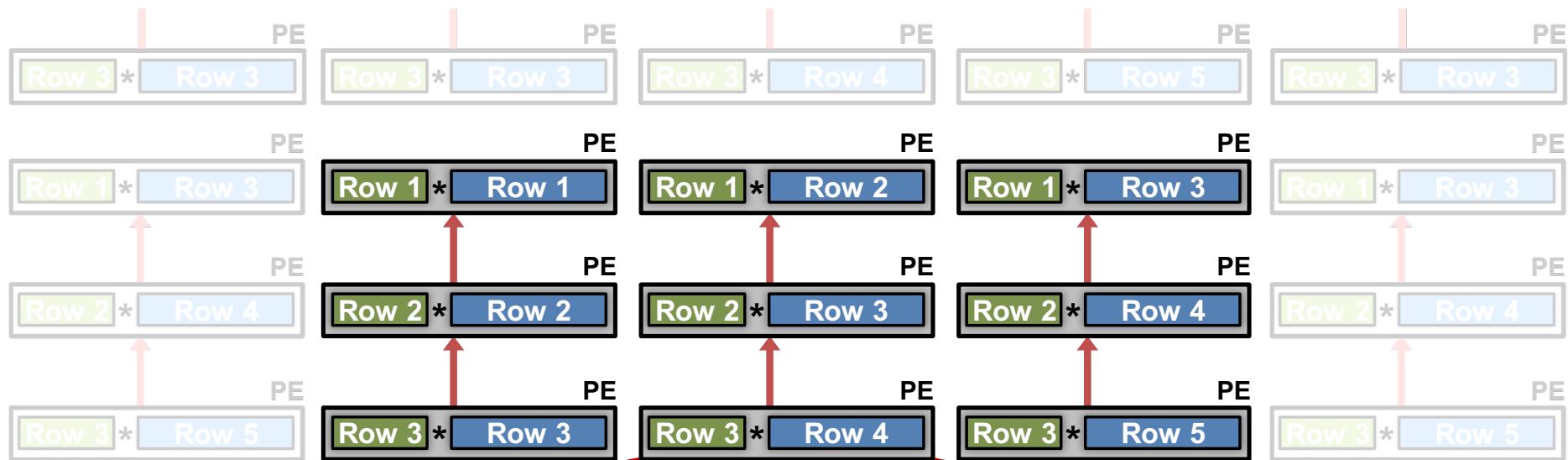
Fmap rows are reused across PEs **diagonally**

Maximize 2D Accumulation in PE Array



Partial sums accumulate across PEs **vertically**

Flexibility to Map Multiple Dimensions



Multiple **fmaps**:

$$\text{Filter 1} \quad \text{Fmap 1 \& 2} \quad \text{Psum 1 \& 2}$$
$$\text{---} * \quad \text{---} = \quad \text{---}$$

Multiple **filters**:

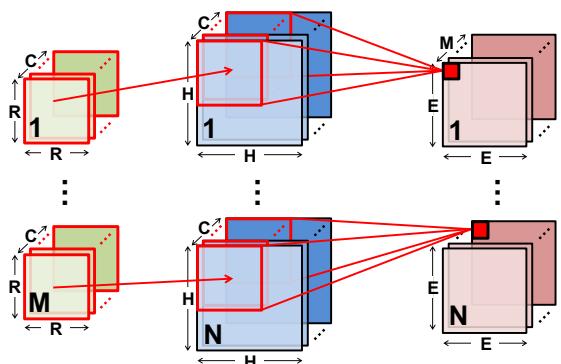
$$\text{Filter 1 \& 2} \quad \text{Fmap 1} \quad \text{Psum 1 \& 2}$$
$$\text{---} * \quad \text{---} = \quad \text{---}$$

Multiple **channels**:

$$\text{Filter 1} \quad \text{Image 1} \quad \text{Psum}$$
$$\text{---} * \quad \text{---} = \quad \text{---}$$

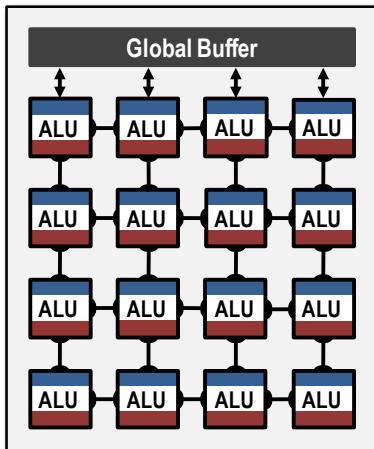
Finding the Optimal Mapping

DNN Configurations

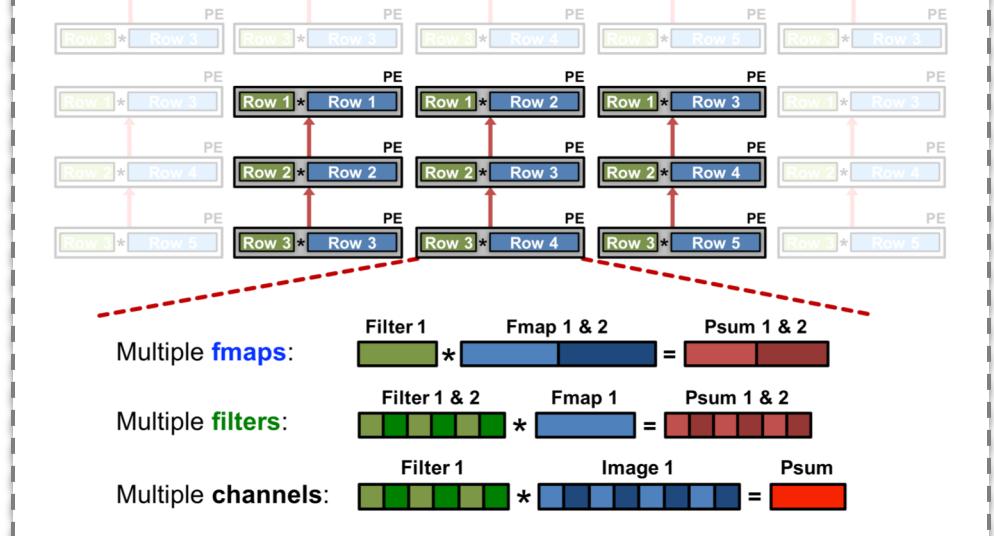


Optimization Compiler

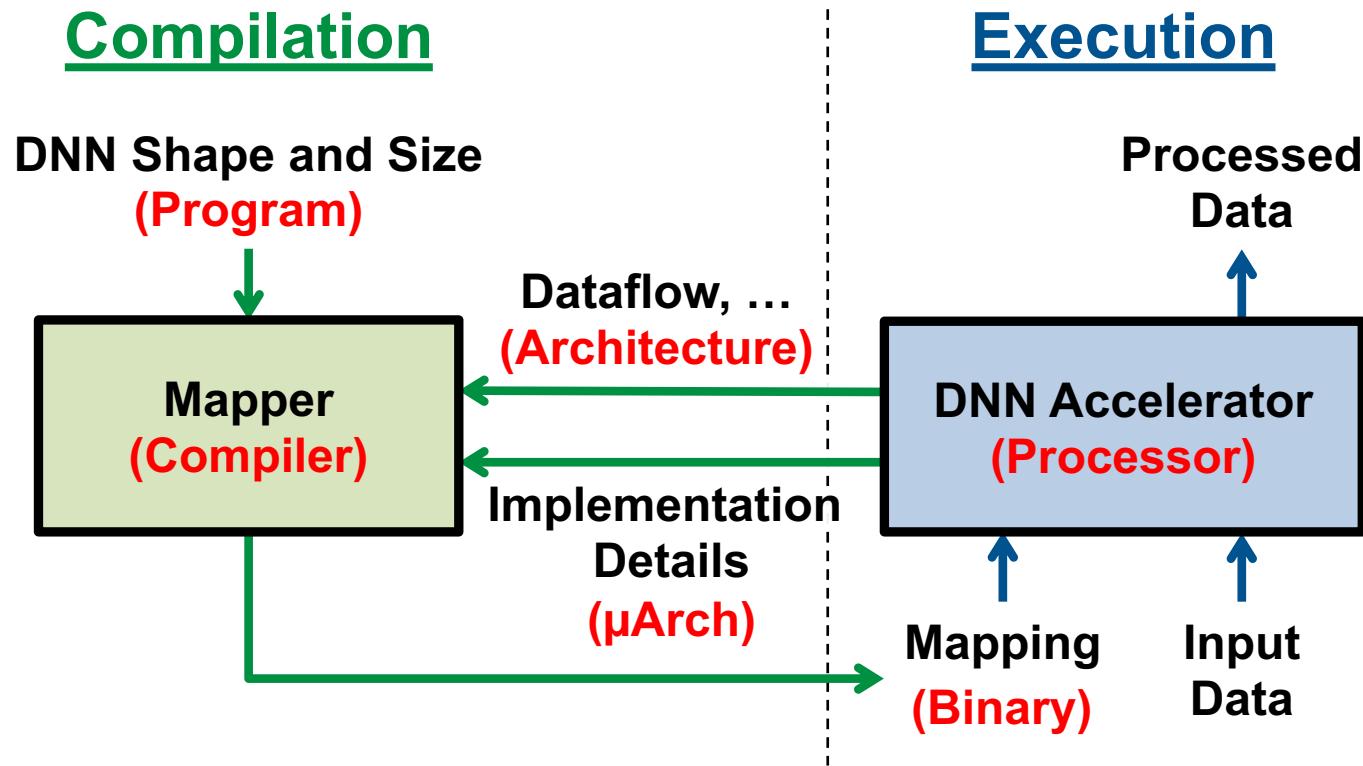
Hardware Resources



Row Stationary Mapping



Computer Architecture Analogy

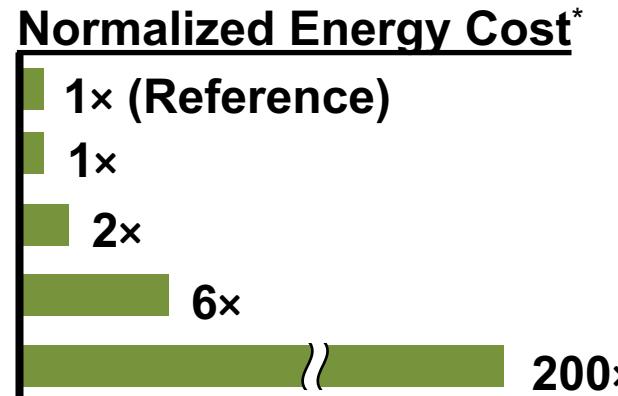
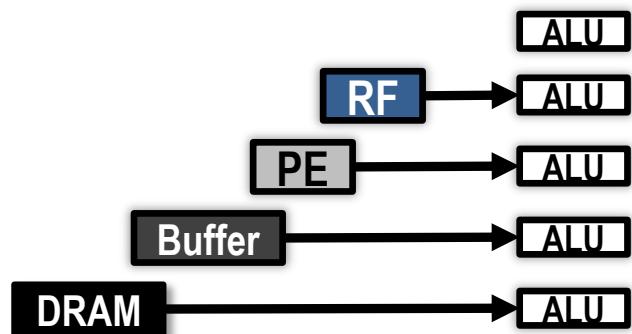


Evaluate Different Dataflows

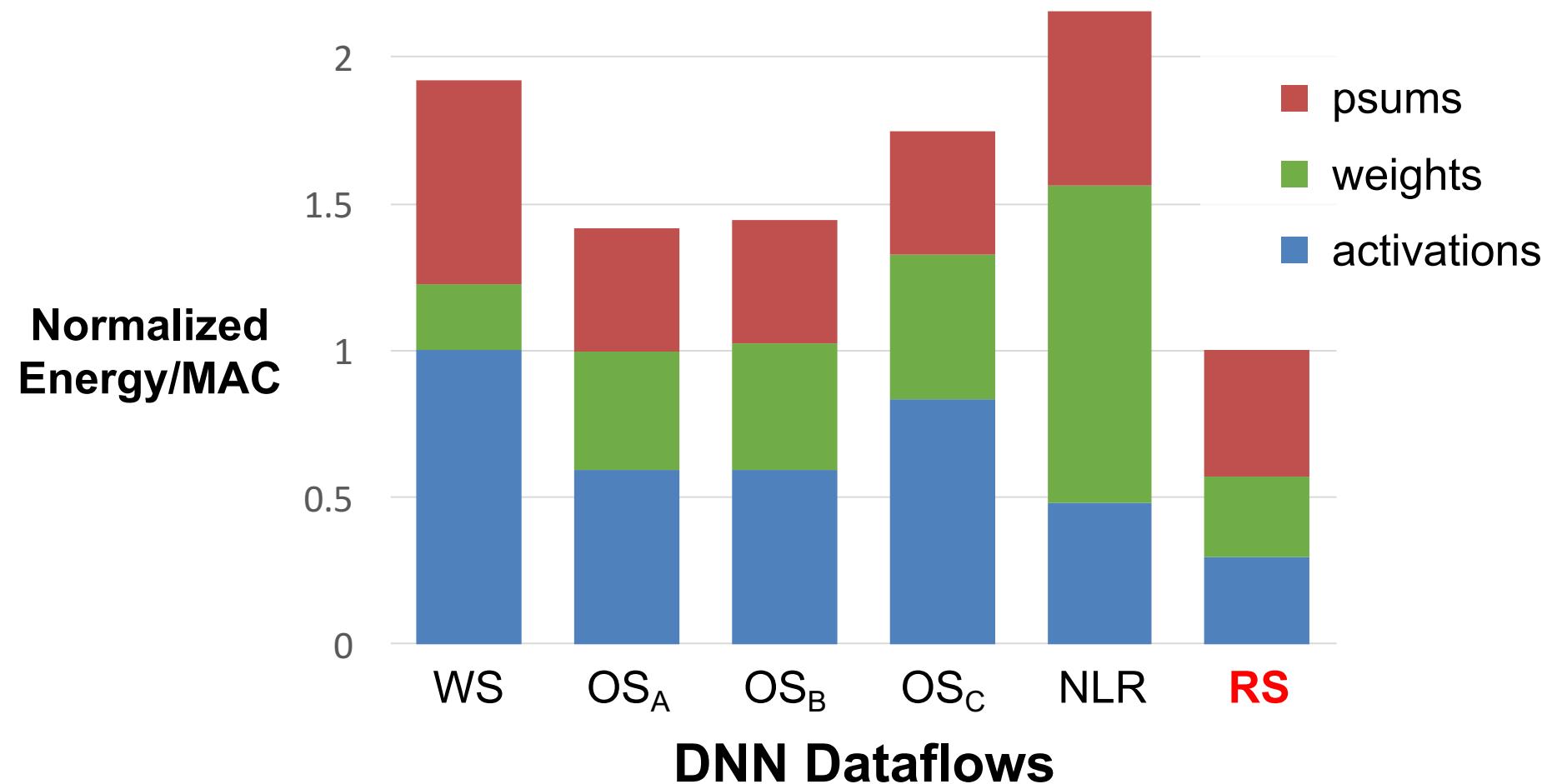
- **Weight Stationary**
 - minimize movement of **weights**
- **Output Stationary**
 - minimize movement of **partial sums**
- **No Local Reuse**
 - don't use any local PE storage. Maximize **global buffer size**
- **Row Stationary**

Evaluation Setup

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

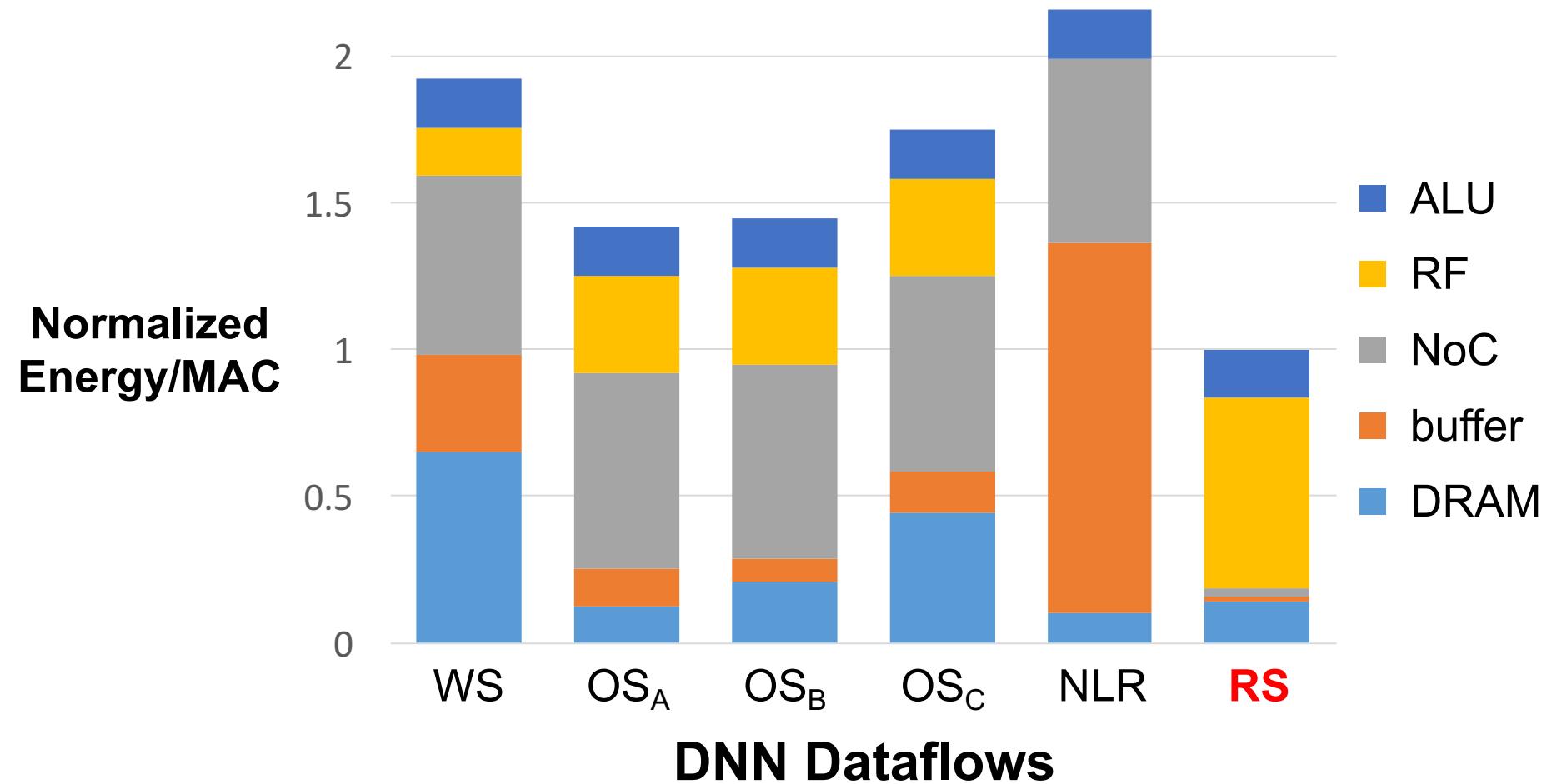


Dataflow Comparison: CONV Layers



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

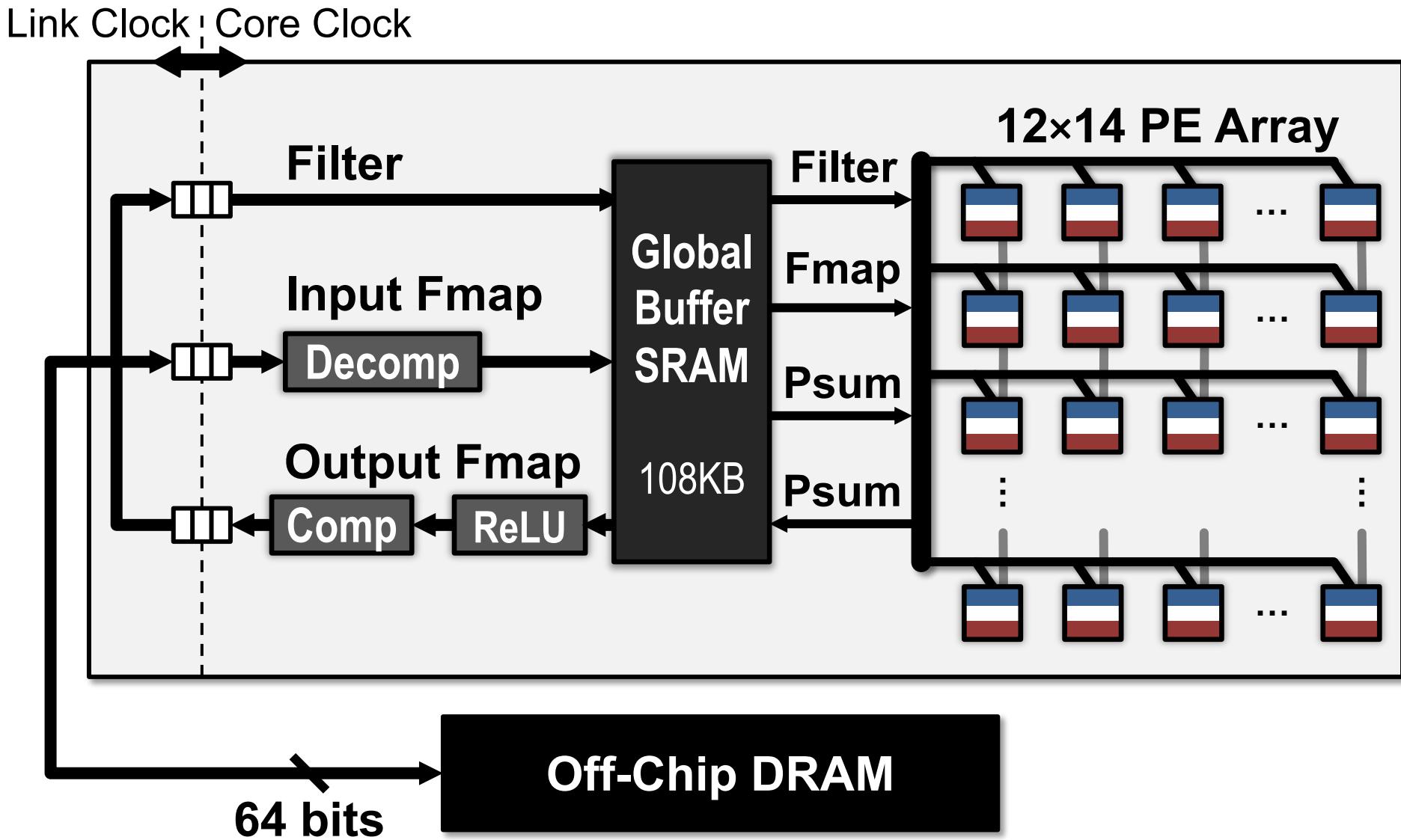
Dataflow Comparison: CONV Layers



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

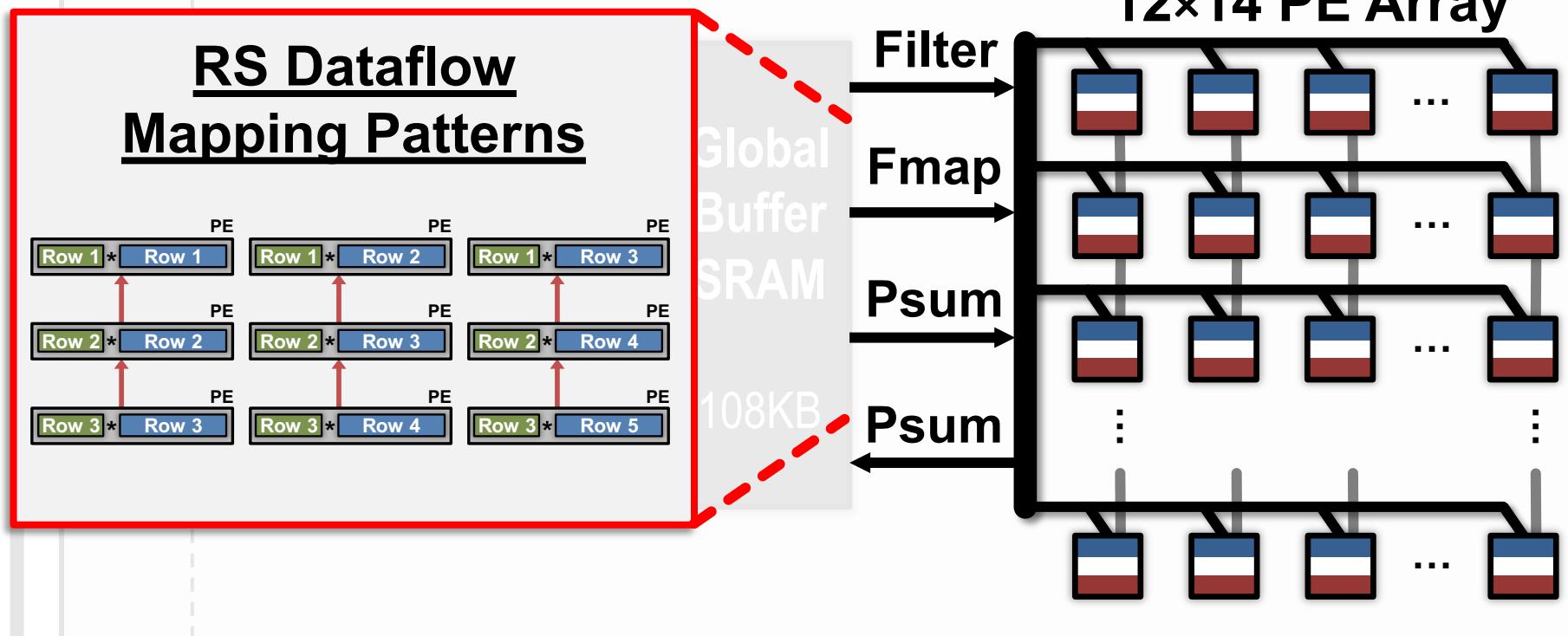
Hardware Architecture for RS Dataflow

Eyeriss: DNN Accelerator for RS Dataflow



Flexibility Required for Mapping

Link Clock | Core Clock

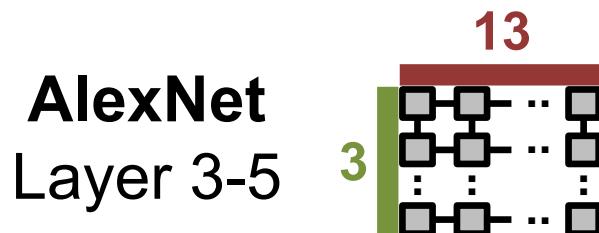


How to map **different shapes** with a **fixed-size PE array**?

64 bits

Flexible Mapping Strategy

Replication



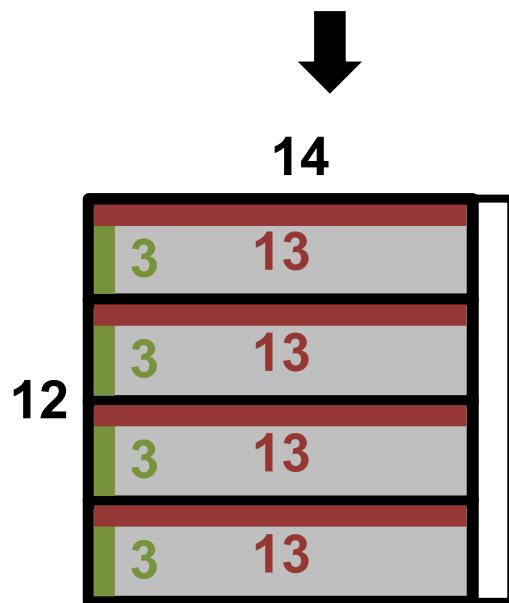
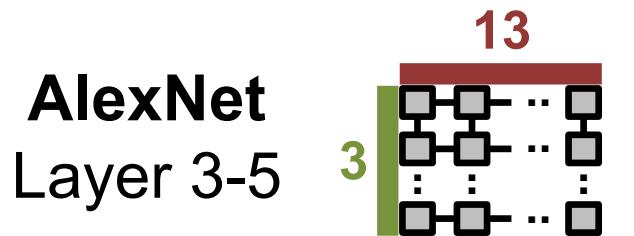
14



Physical PE Array

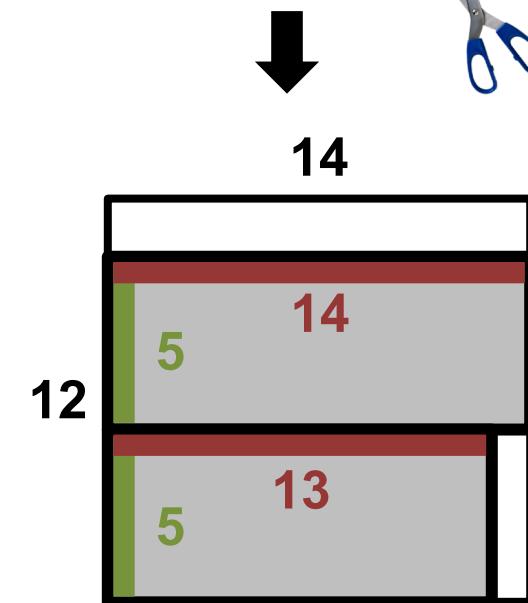
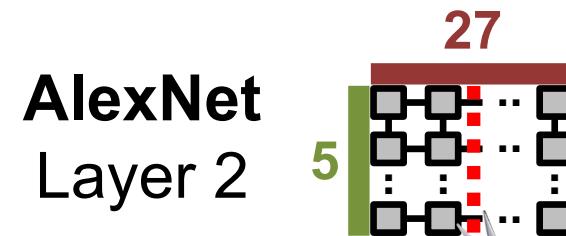
Flexible Mapping Strategy

Replication



Physical PE Array

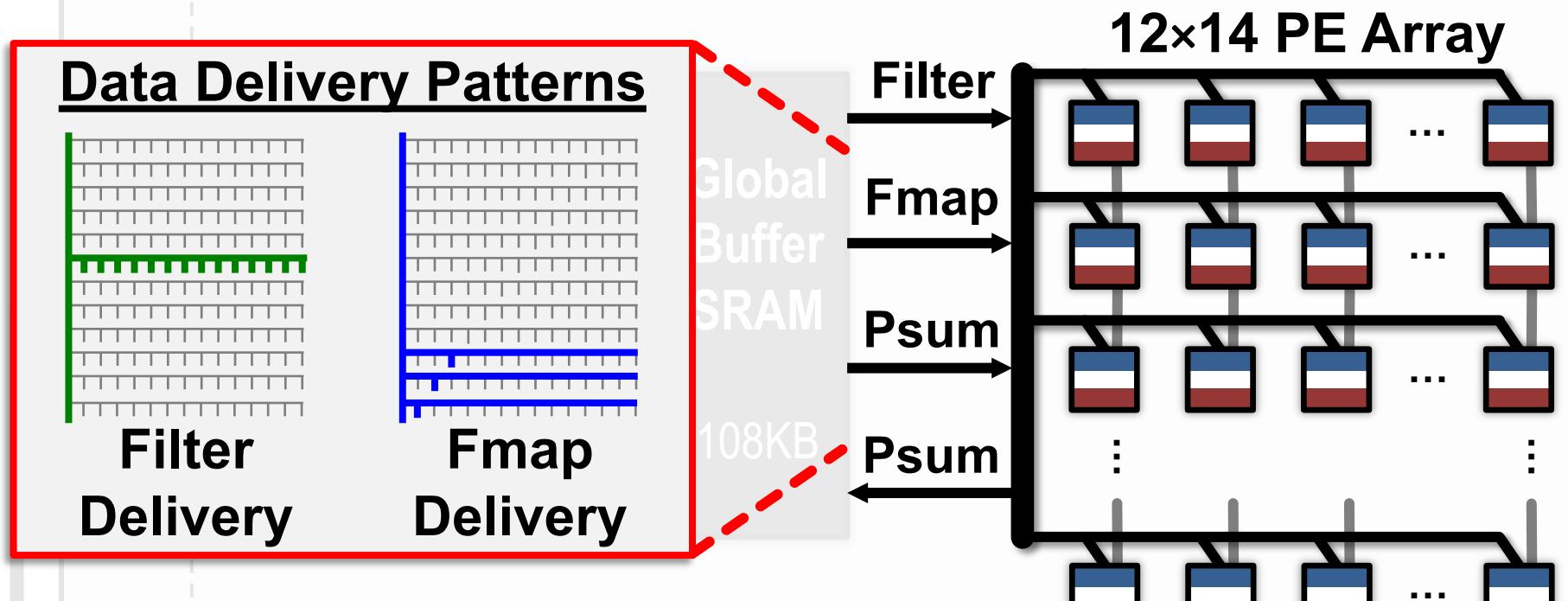
Folding



Physical PE Array

Data Delivery with On-Chip Network

Link Clock | Core Clock

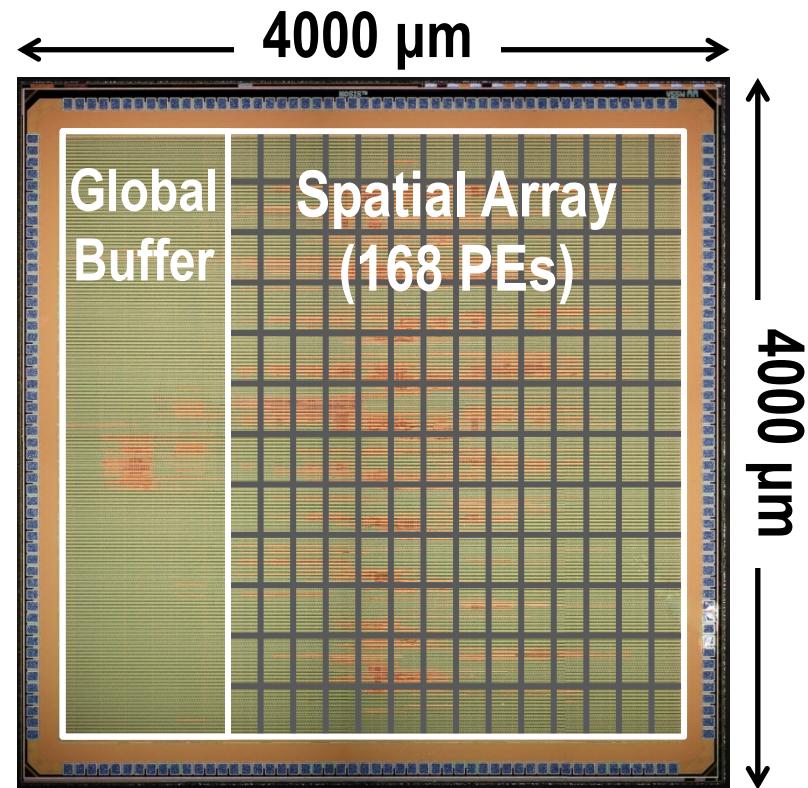


Supports any delivery pattern & saves over 80% energy compared to broadcast

64 bits

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 DNN 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



To support 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)

DNN Model & Hardware Co-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

Approaches

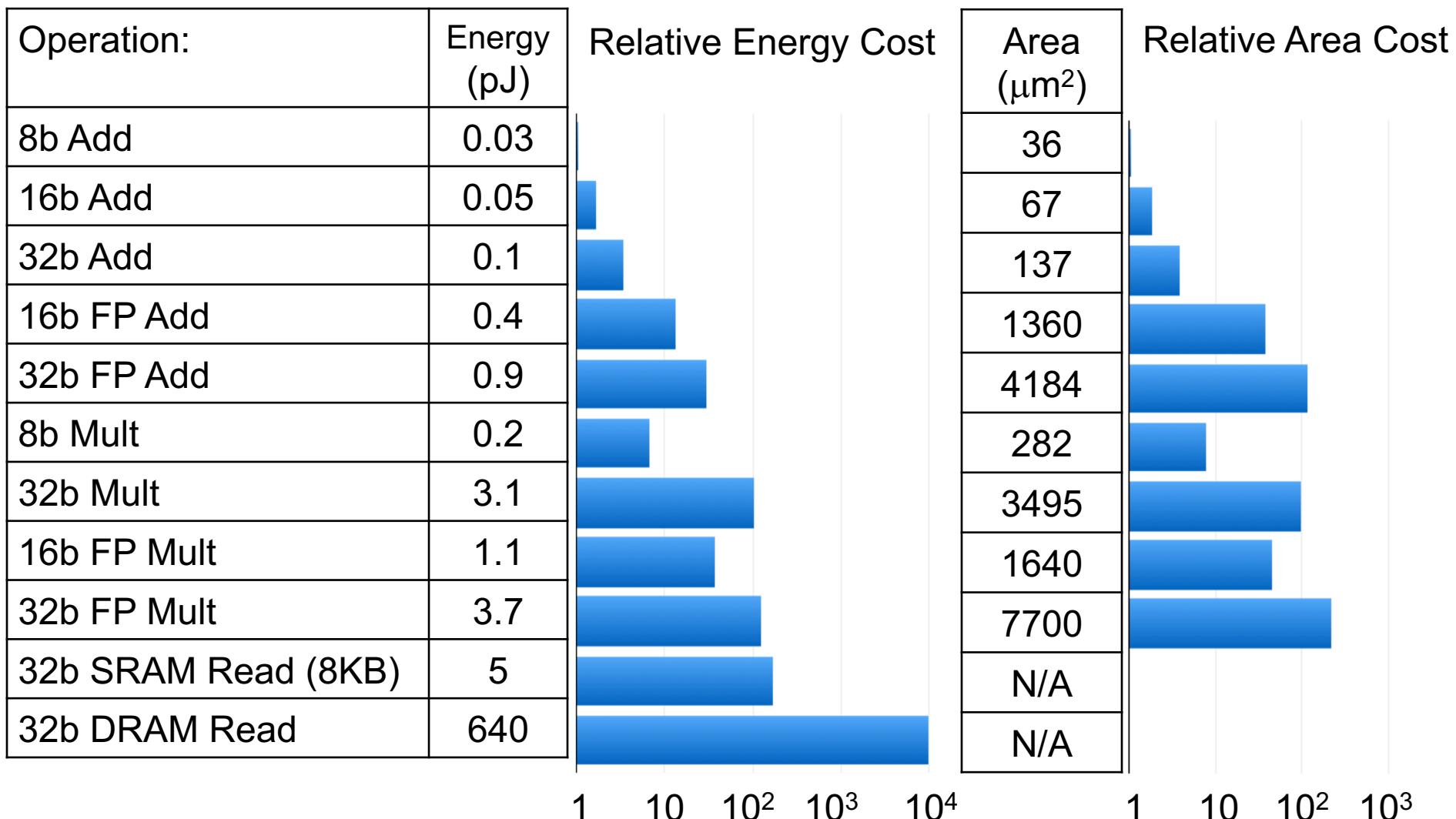
- **Reduce size of operands for storage/compute**
 - Floating point → Fixed point
 - Bit-width reduction
 - Non-linear quantization
- **Reduce number of operations for storage/compute**
 - Exploit Activation Statistics (Compression)
 - Network Pruning
 - Compact Network Architectures

Accuracy needs to be taken into account!

Approaches

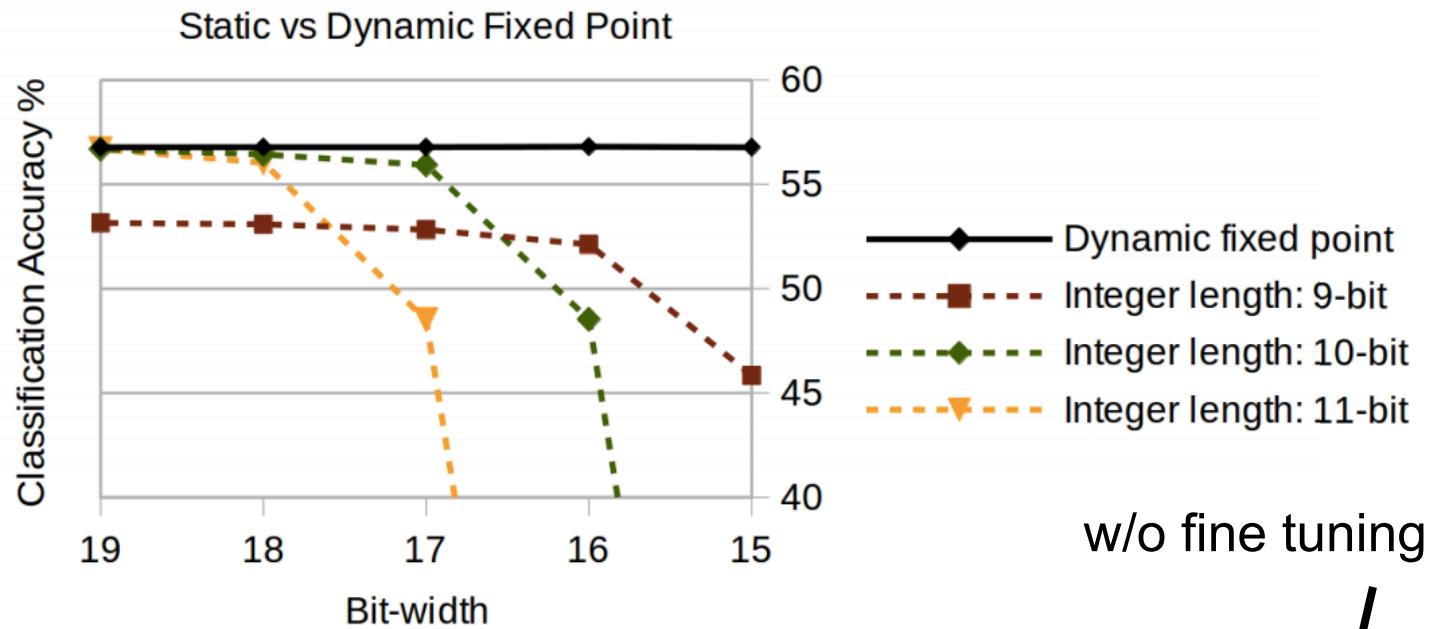
- **Reduce size of operands for storage/compute**
 - Floating point → Fixed point
 - Bit-width reduction
 - Non-linear quantization
- **Reduce number of operations for storage/compute**
 - Exploit Activation Statistics (Compression)
 - Network Pruning
 - Compact Network Architectures

Cost of Operations



Impact on Accuracy

Top-1 accuracy
on of CaffeNet
on ImageNet



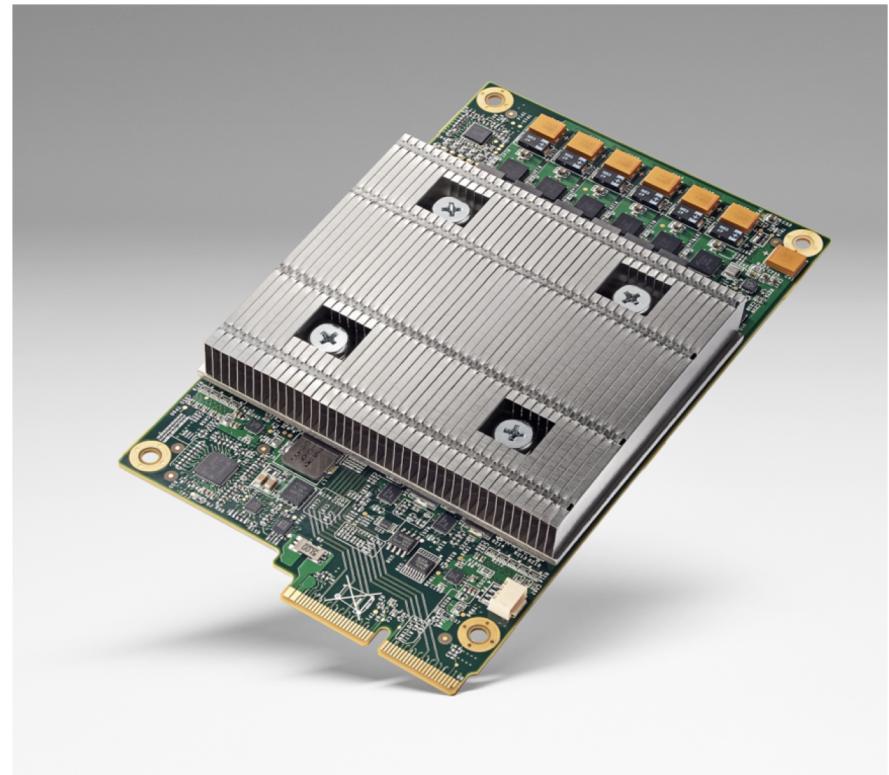
	Layer outputs	CONV parameters	FC parameters	32-bit floating point baseline	Fixed point accuracy
LeNet (Exp 1)	4-bit	4-bit	4-bit	99.1%	99.0% (98.7%)
LeNet (Exp 2)	4-bit	2-bit	2-bit	99.1%	98.8% (98.0%)
Full CIFAR-10	8-bit	8-bit	8-bit	81.7%	81.4% (80.6%)
SqueezeNet top-1	8-bit	8-bit	8-bit	57.7%	57.1% (55.2%)
CaffeNet top-1	8-bit	8-bit	8-bit	56.9%	56.0% (55.8%)
GoogLeNet top-1	8-bit	8-bit	8-bit	68.9%	66.6% (66.1%)

Reduced Precision in the Industry

8-bit INT is widely used for inference in commercial products



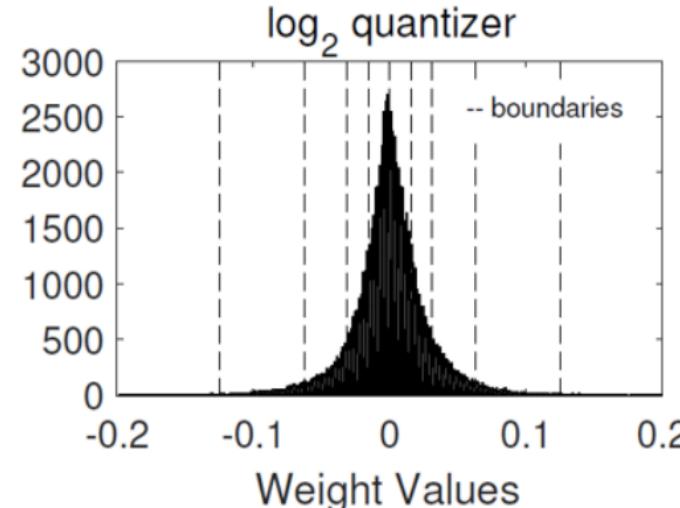
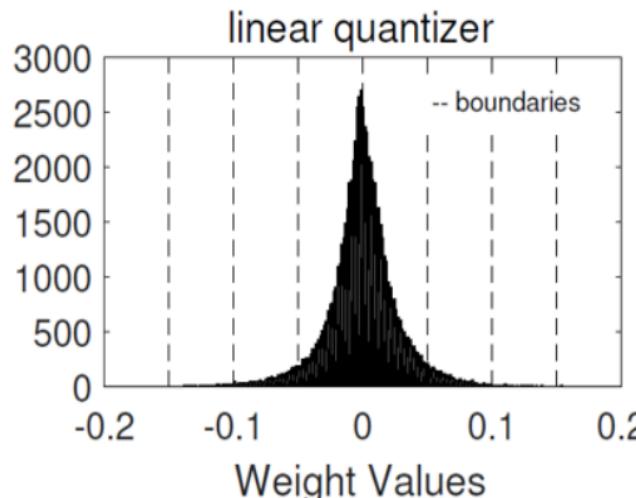
Nvidia Pascal (2016)



Google's TPU v1 (2016)

Reduced Precision in Research

Network	Precision		Accuracy Loss*
	Weights	Activations	
Binarized Neural Net (BNN)	1	1	29.8%
XNOR-Net	1 [†]	1 [†]	11.0%
HWGQ-Net	1 [†]	2 [†]	5.2%
Binary Weight Net (BWN)	1 [†]	32f	0.8%
Trained Ternary Quant (TTQ)	2 [†]	32f	0.6%
LogNet	4/5	4	3.2%



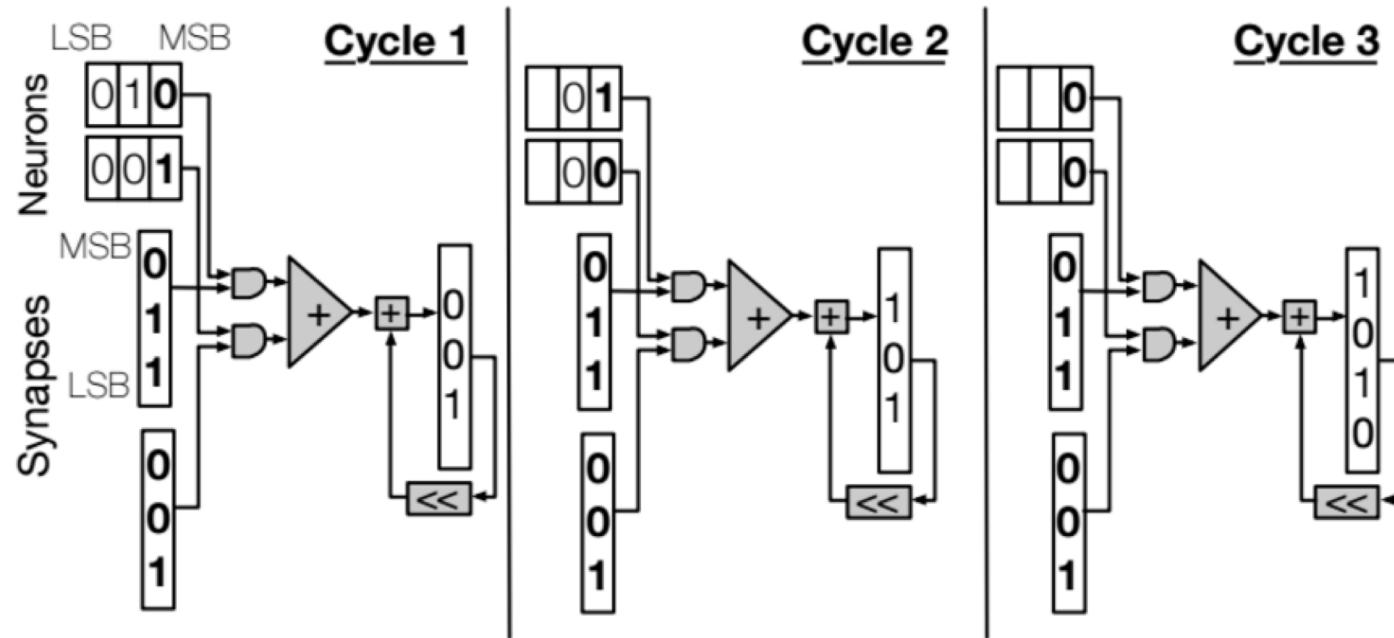
* AlexNet,
compared
to 32f

+ 32f in first
and last layers

Bitwidth Scaling for Higher Throughput

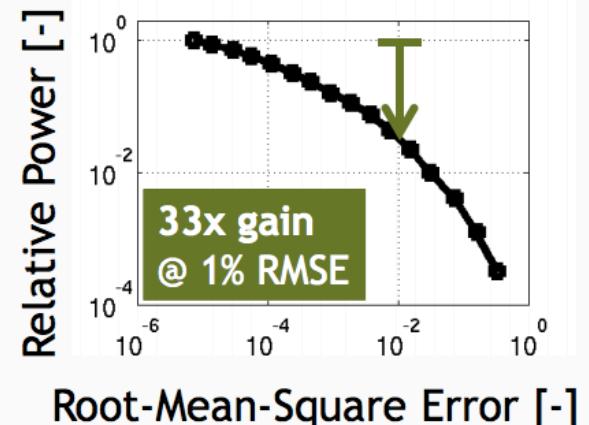
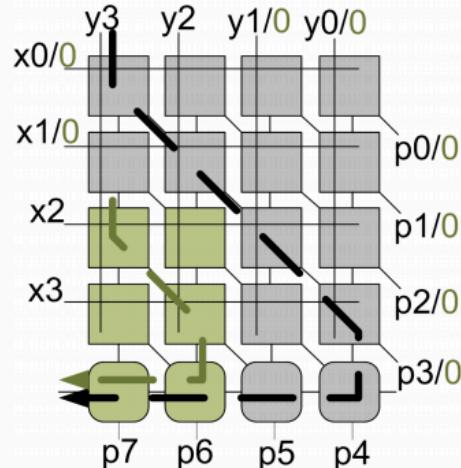
Bit-Serial Processing: Reduce Bitwidth → Skip Cycles
2.24x speed-up vs. 16-bit fixed

$$\sum_{i=0}^{N_i-1} s_i \times n_i = \sum_{i=0}^{N_i-1} s_i \times \sum_{b=0}^{P-1} n_i^b \times 2^b = \sum_{b=0}^{P-1} 2^b \times \sum_{i=0}^{N_i-1} n_i^b \times S_i$$



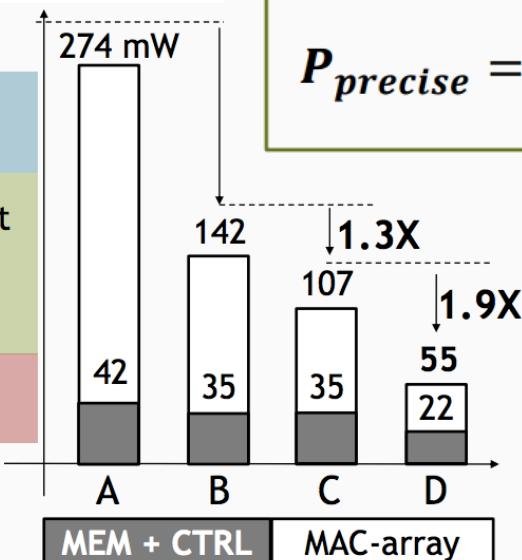
Bitwidth Scaling for Lower Power

Reduce Bitwidth
→ Shorter Critical Path
→ Reduce Voltage



AlexNet Layer 2 example:

- A. 2D-baseline @ 16 bit
- B. Precision-Scaling @ 7-7 bit
- C. Voltage-Scaling @ 0.9 V
- D. Sparse operation guarding



$$P_{precise} = \alpha CfV^2 \Rightarrow P_{imprecise} = \frac{\alpha}{k_1} Cf\left(\frac{V}{k_2}\right)^2$$

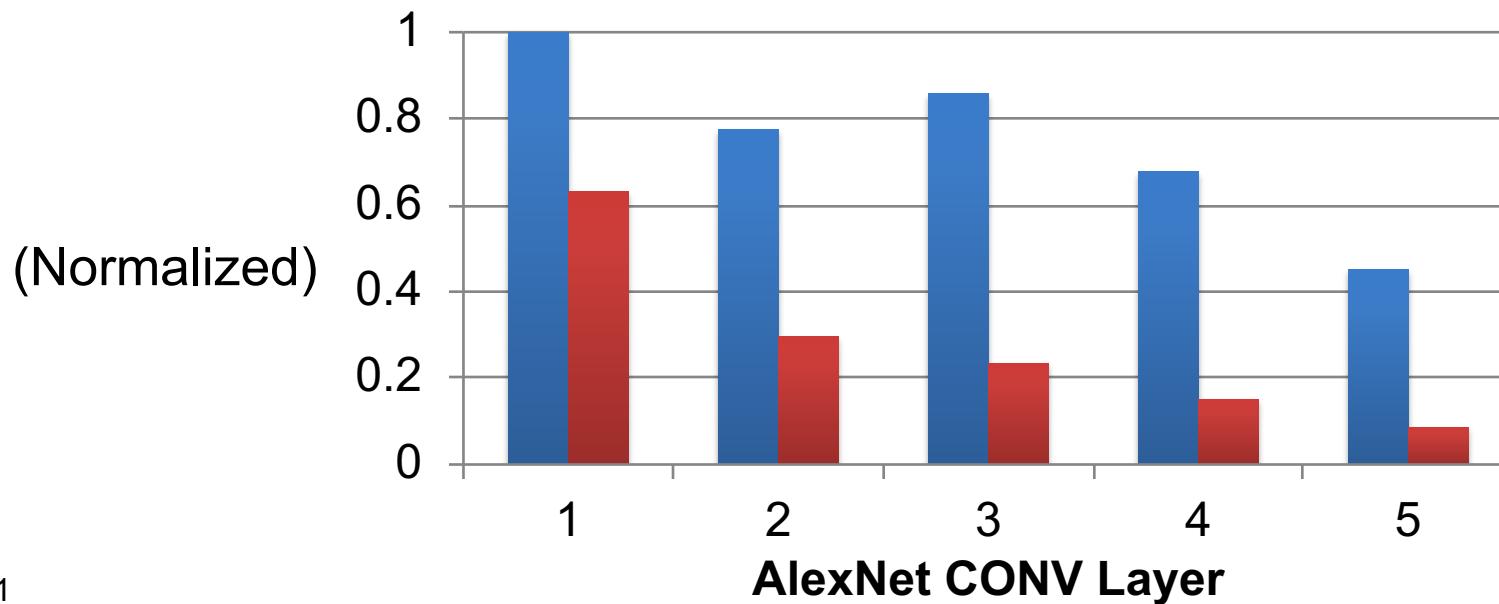
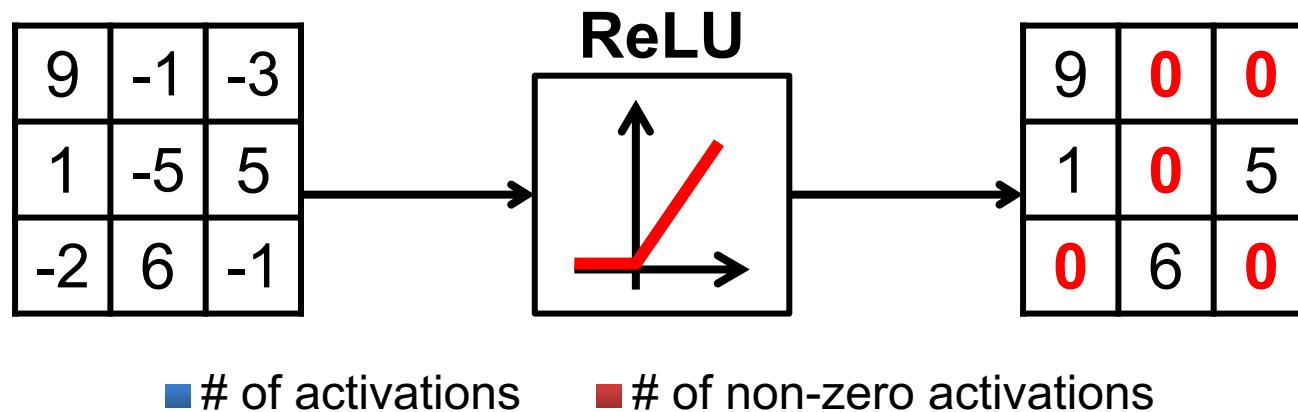
2.56x power reduction
vs. 16-bit fixed
(AlexNet layer 2)

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

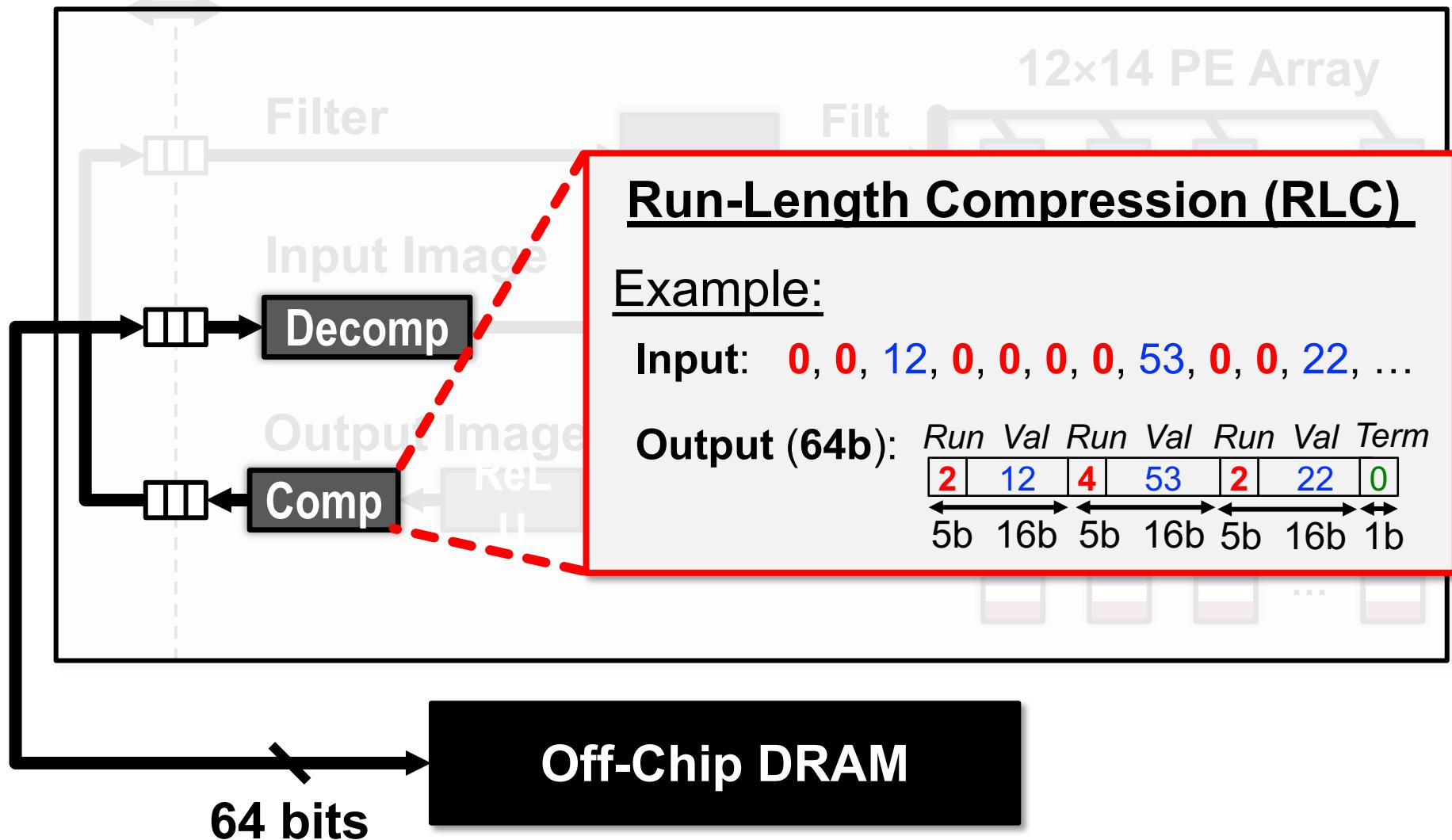
Sparsity in Feature Maps

Many **zeros** in feature maps after ReLU



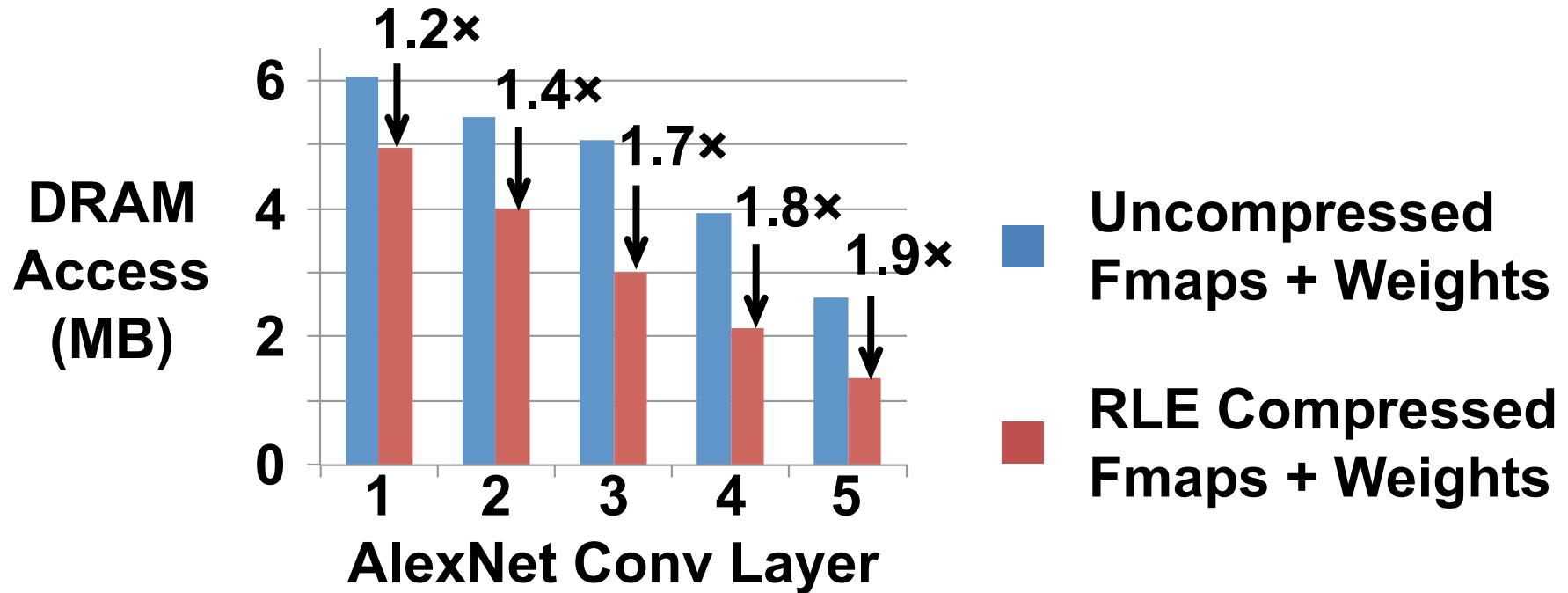
I/O Compression in Eyeriss

Link Clock | Core Clock



[Chen et al., ISSCC 2016]

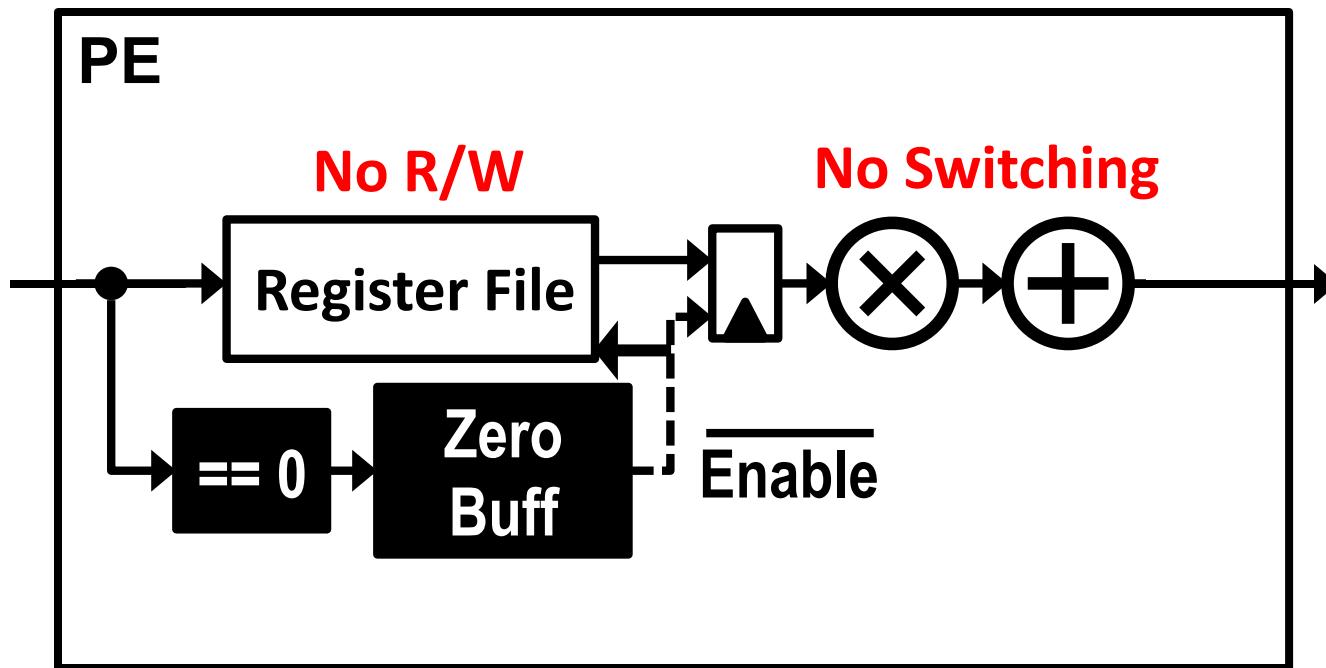
Compression Reduces DRAM BW



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

Data Gating / Zero Skipping in Eyeriss

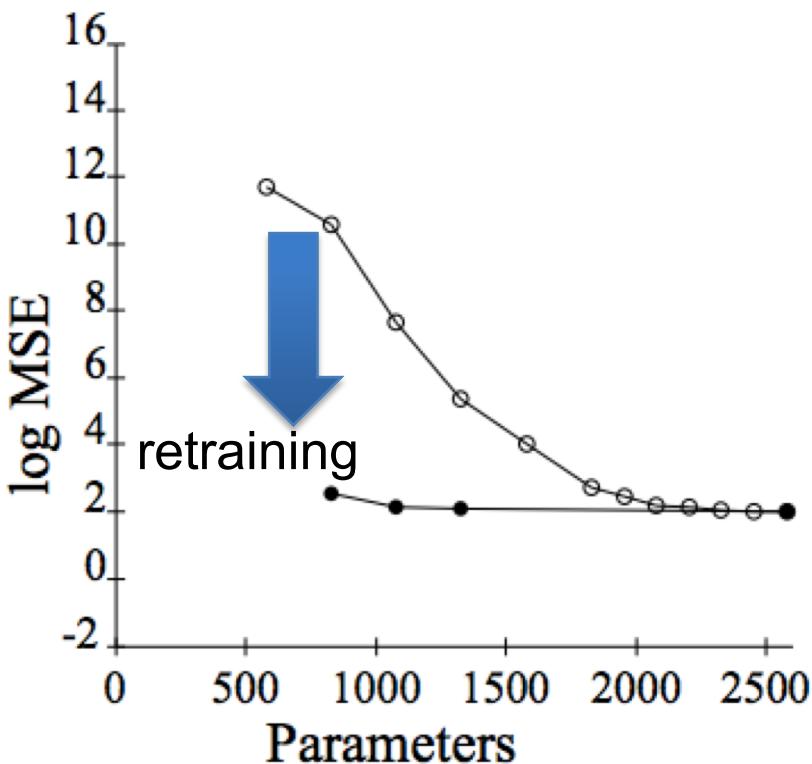
Reduce PE Power by **45%**



Pruning – Make Weights Sparse

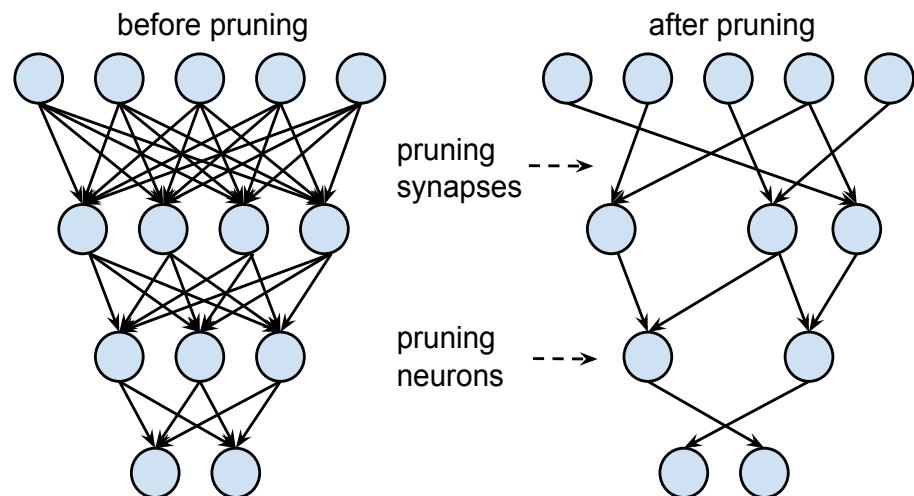
Optimal Brain Damage

[LeCun et al., NIPS 1989]



Magnitude-based Weight Pruning

[Han et al., NIPS 2015]

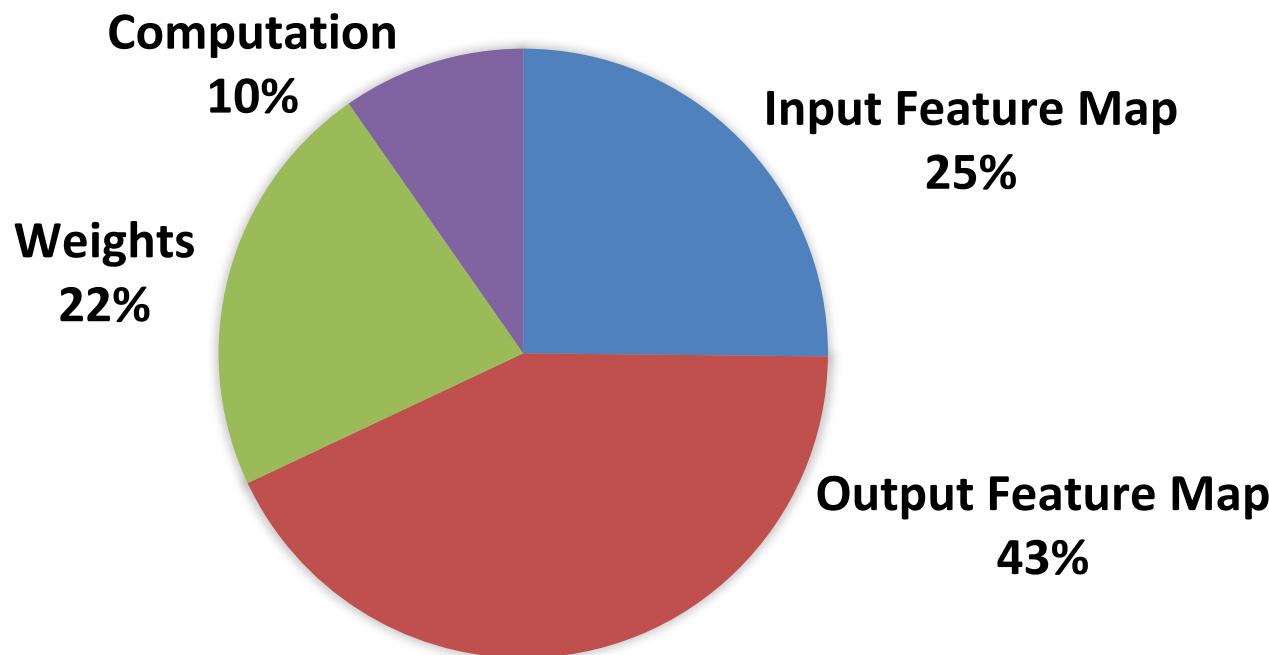


Example: AlexNet
Weight Reduction:
CONV layers 2.7x, **FC layers 9.9x**
Overall Reduction:
Weights 9x, MACs 3x

Key Observations on Pruning

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

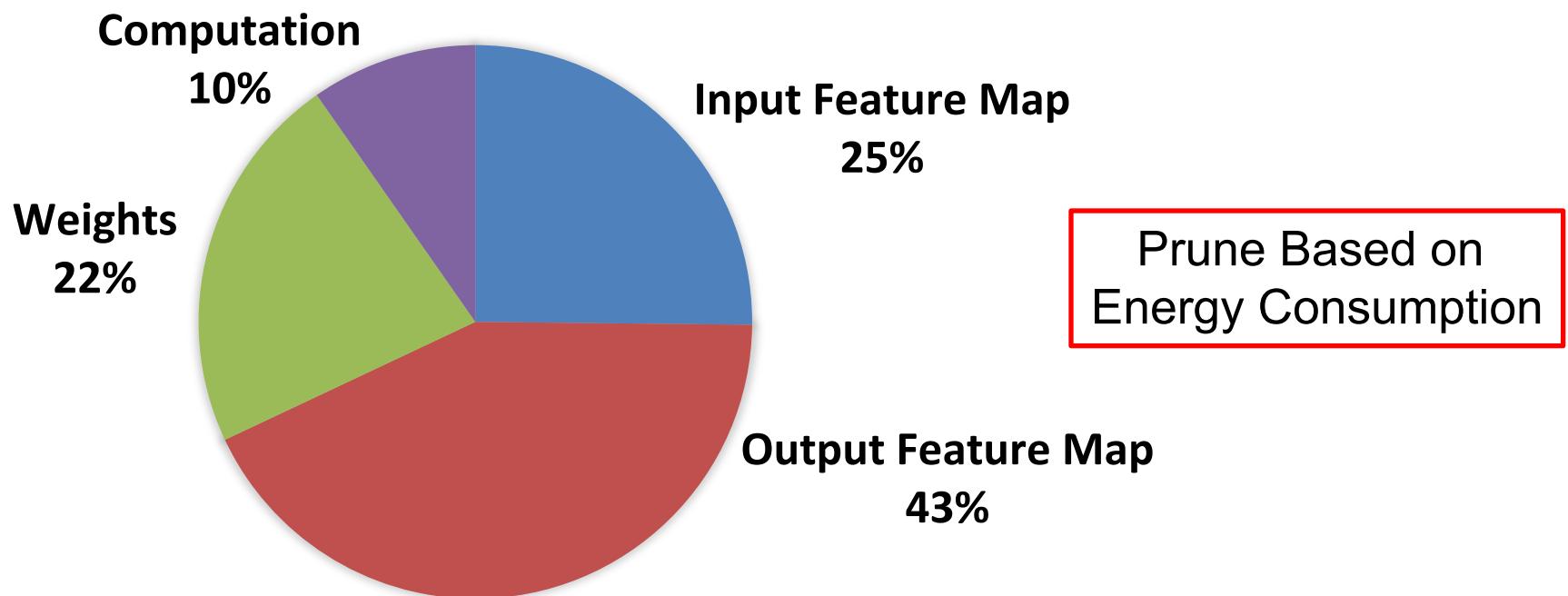
GoogLeNet Energy Consumption



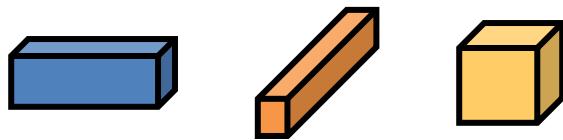
Key Observations on Pruning

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

GoogLeNet Energy Consumption



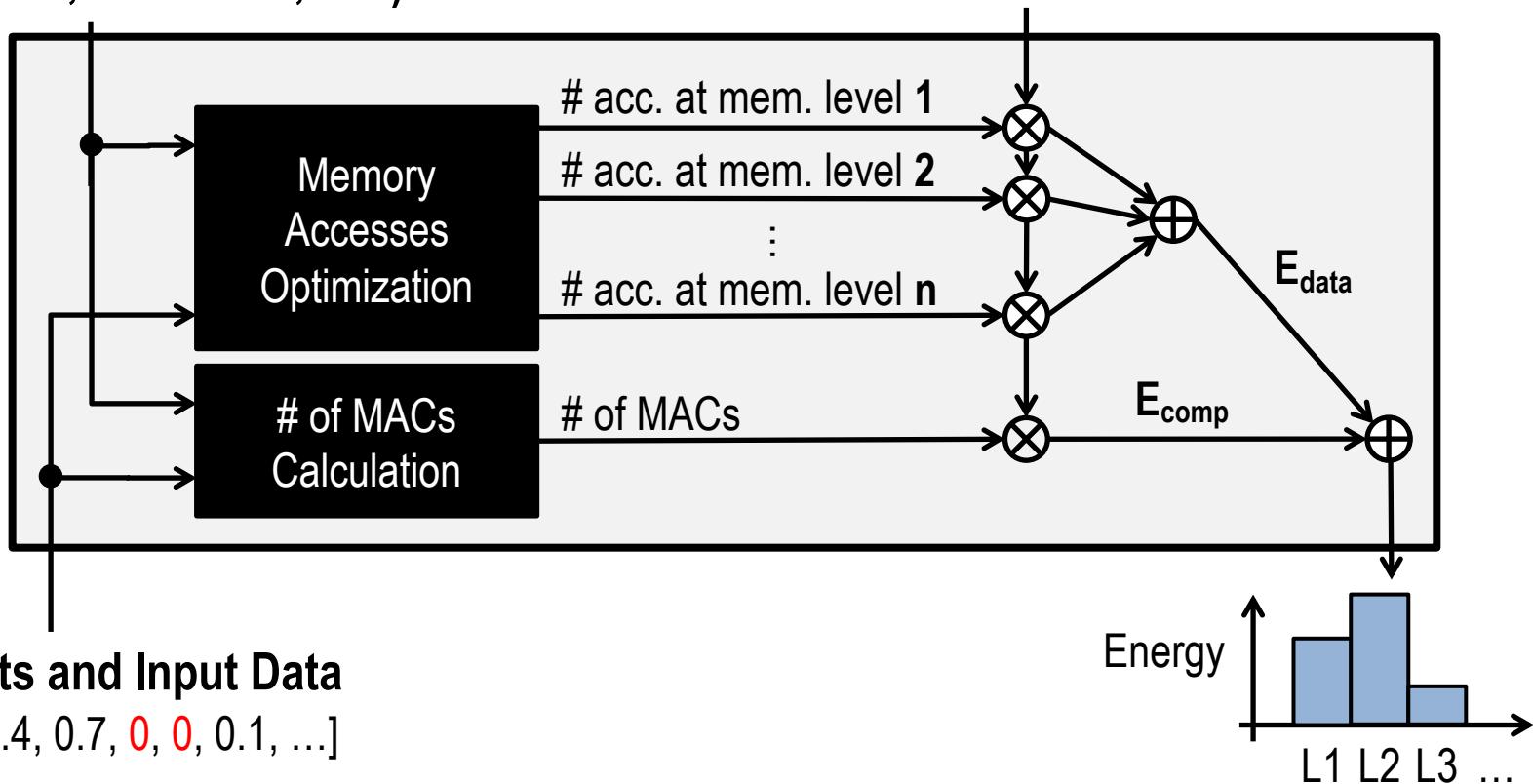
Energy-Evaluation Methodology



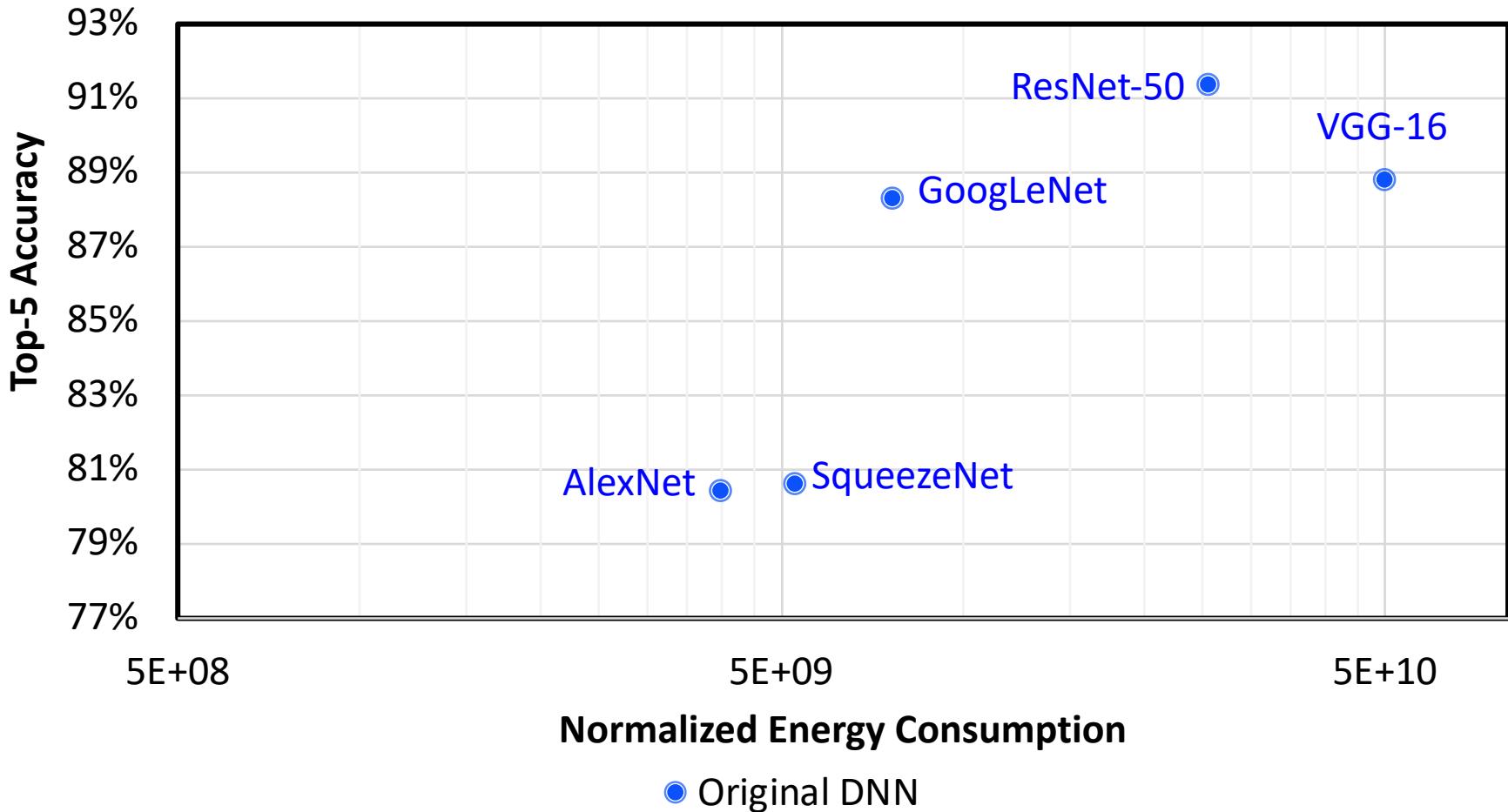
CNN Shape Configuration
(# of channels, # of filters, etc.)

Tool available at <http://eyeriss.mit.edu>

Hardware Energy Costs of each MAC and Memory Access

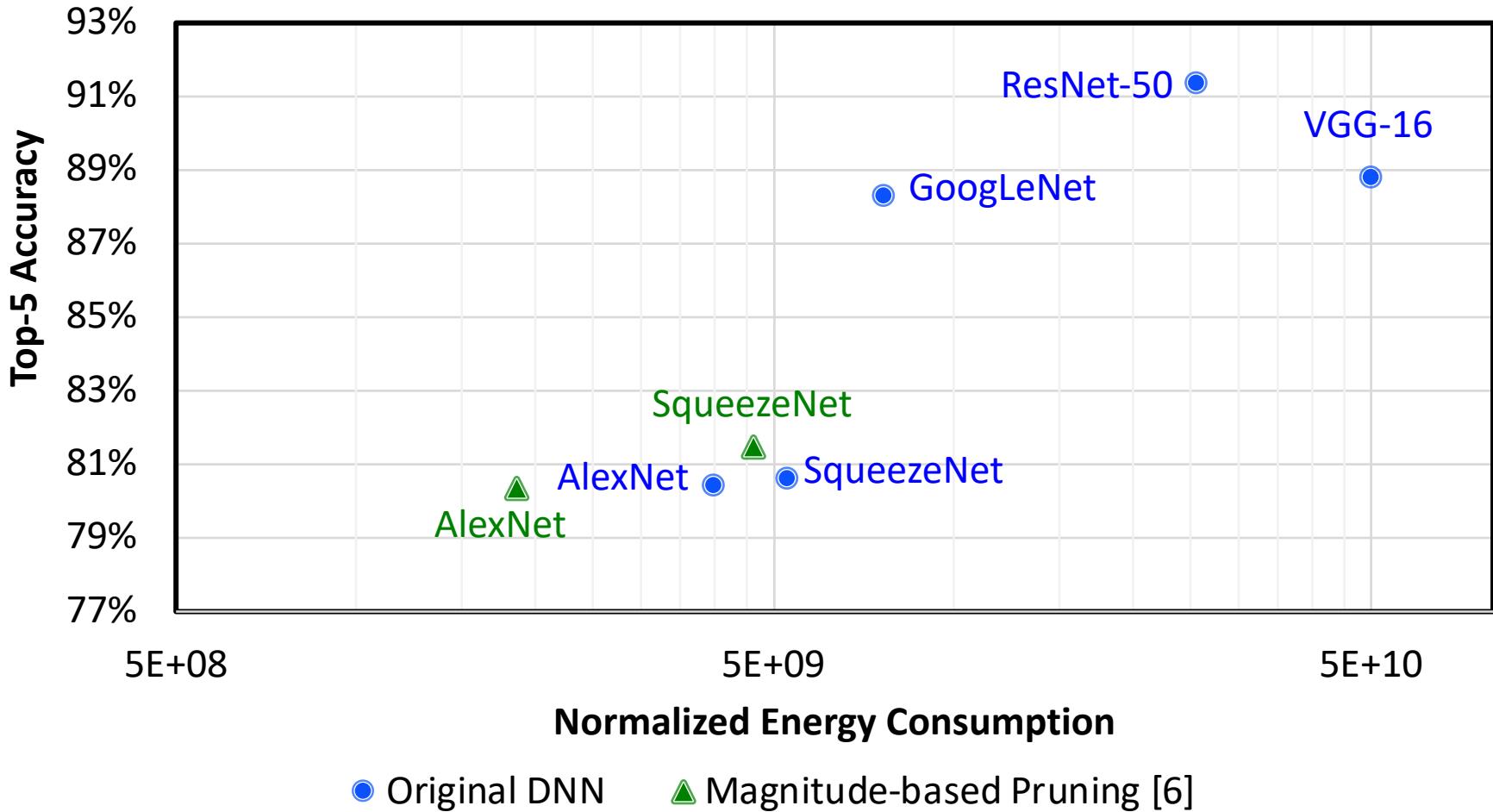


Energy Consumption of Existing DNNs



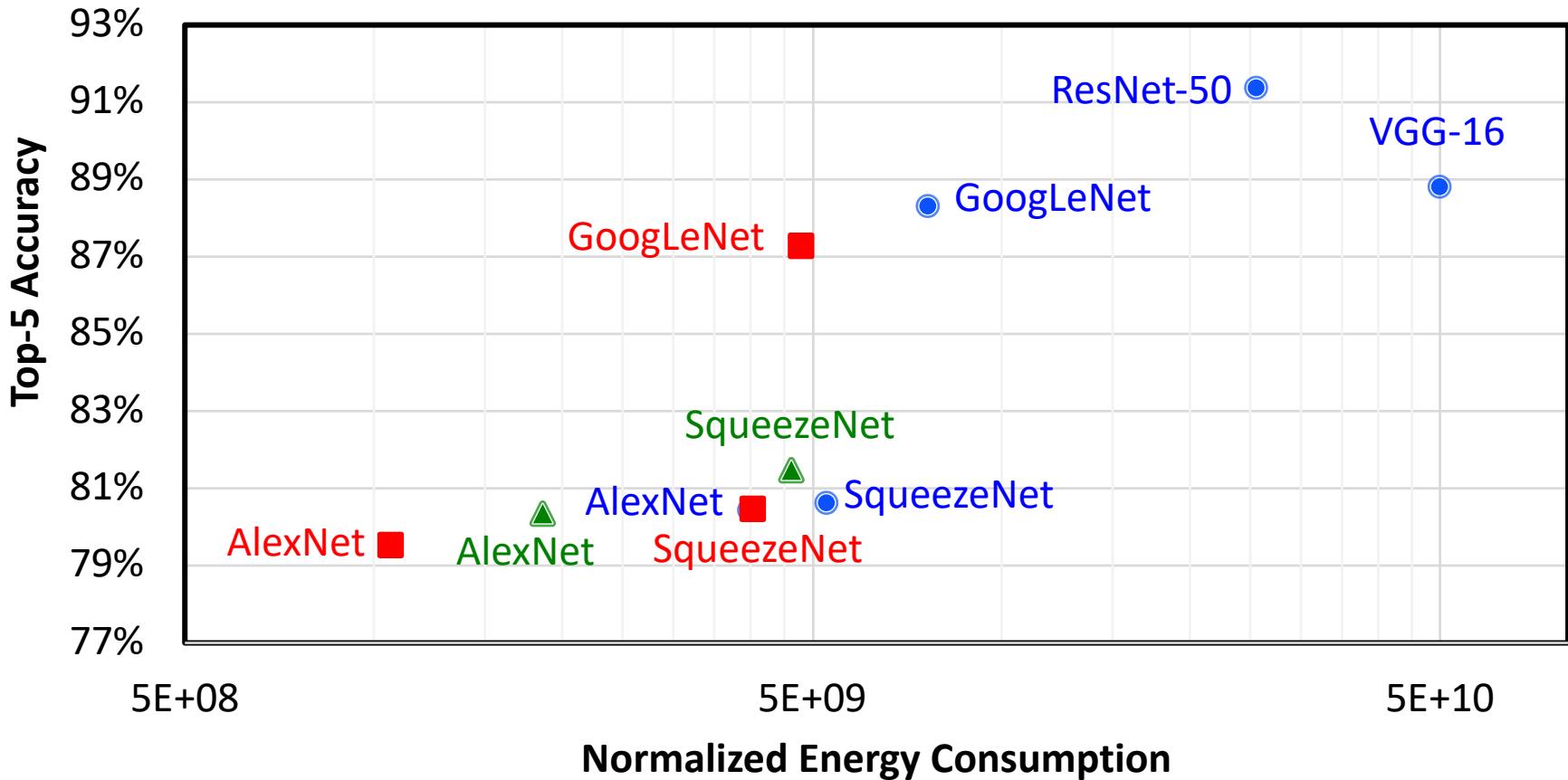
Deeper CNNs with fewer weights do not necessarily consume less energy than shallower CNNs with more weights

Magnitude-based Weight Pruning



Reduce number of weights by removing small magnitude weights

Energy-Aware Weight Pruning



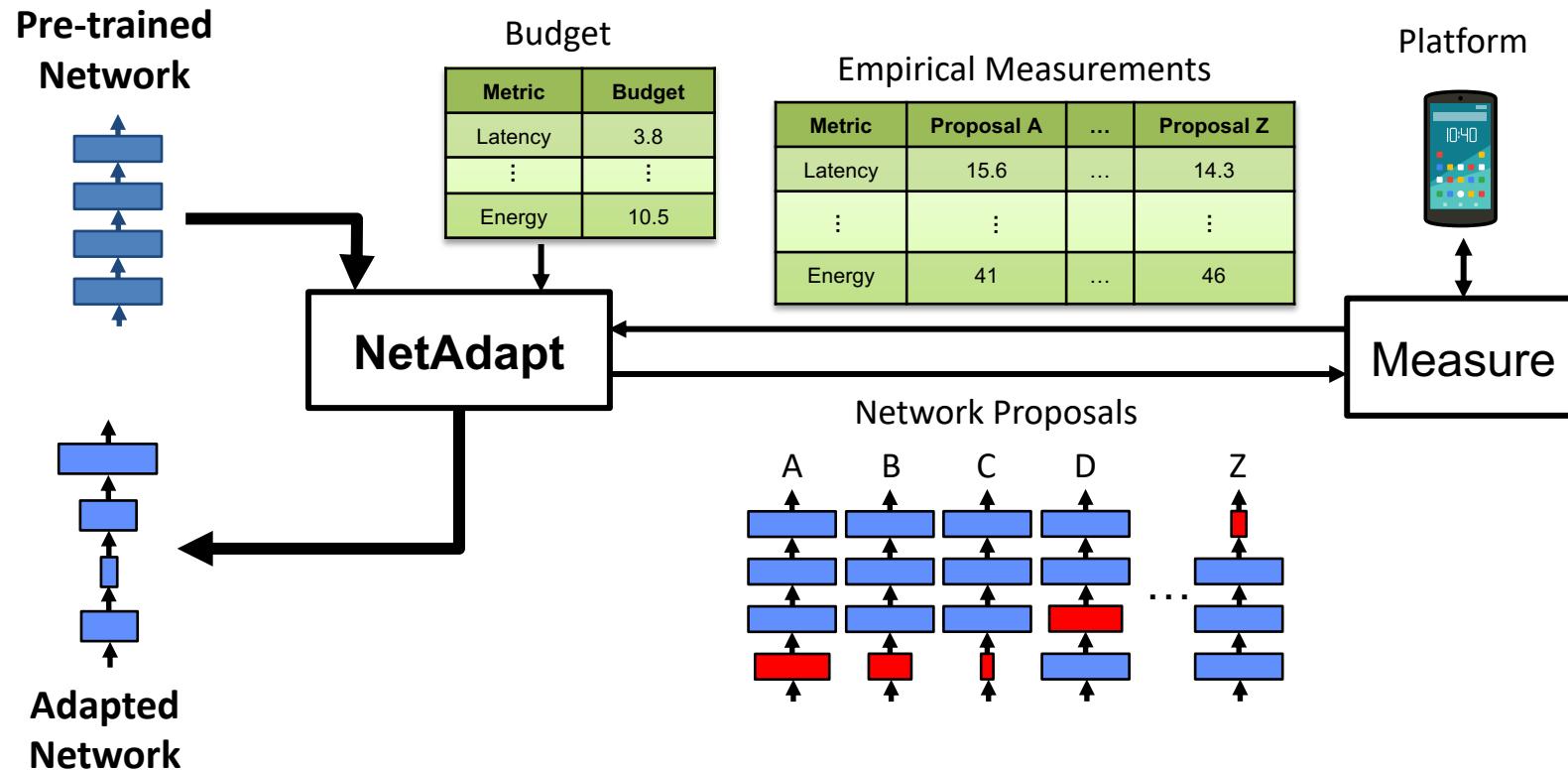
- Original DNN ▲ Magnitude-based Pruning [6] ■ Energy-aware Pruning (This Work)

Remove weights from layers in order of highest to lowest energy

3.7x reduction in AlexNet / **1.6x** reduction in GoogLeNet

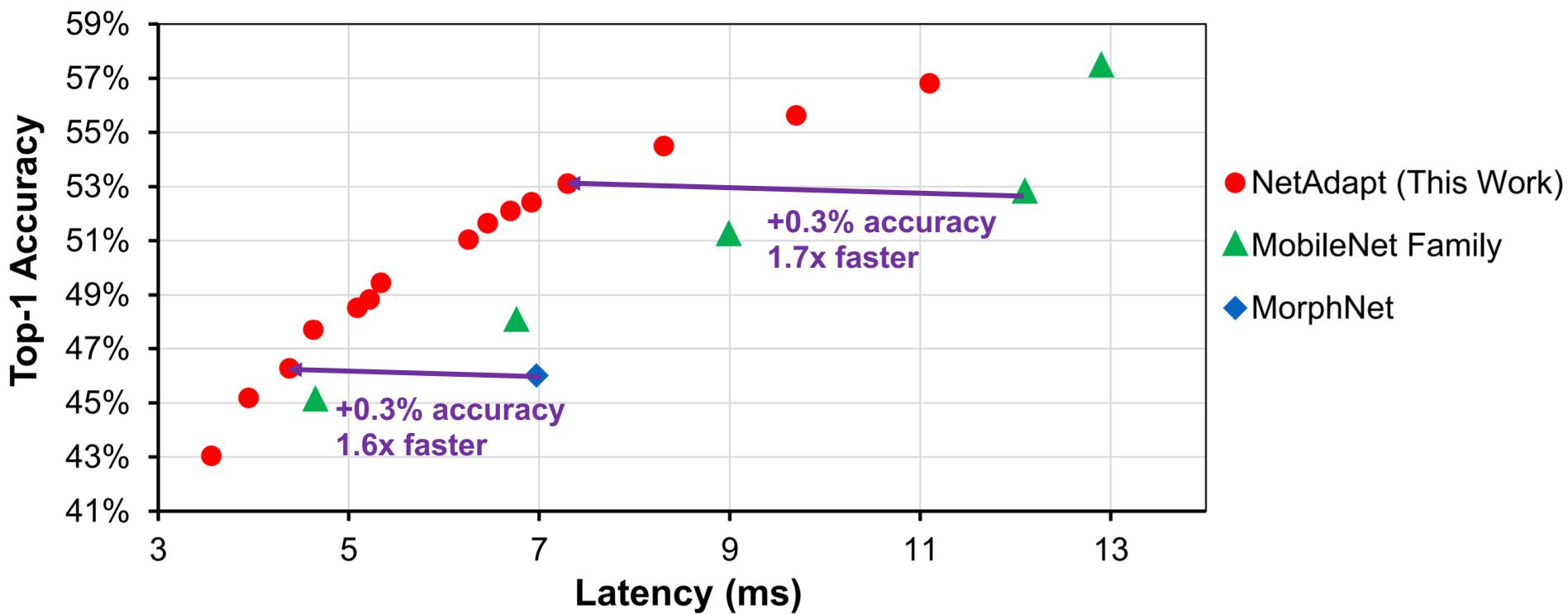
NetAdapt – Platform-Aware DNN Adaptation

- 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)



Latency vs. Accuracy Trade-off

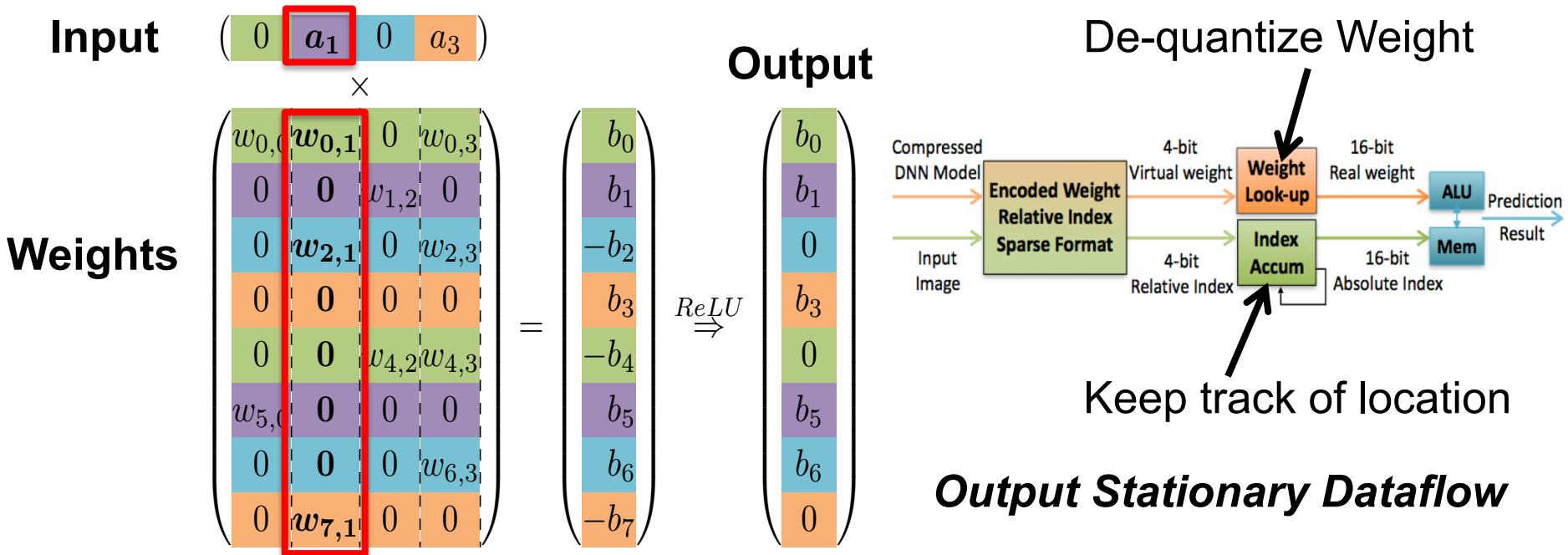
- NetAdapt boosts the real inference speed* of MobileNet by up to **1.7x** with higher accuracy



* Tested on the ImageNet dataset and a Google Pixel 1 CPU

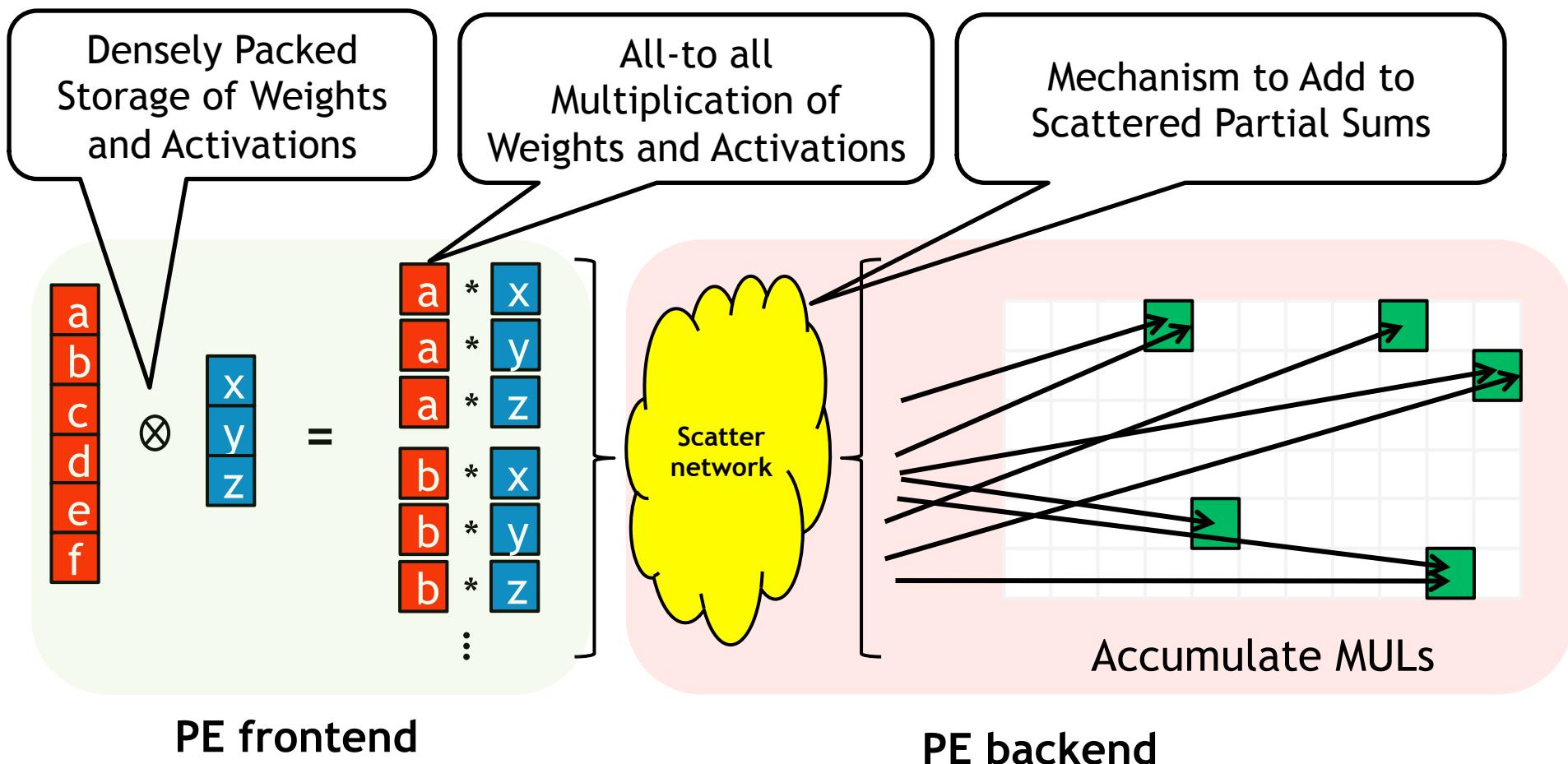
Sparse Architecture Example: EIE

- Supports pruned **FC Layers**
- Store weights in Compressed Sparse Column (CSC) format
- Read corresponding weight column when input is non-zero



Sparse Architecture Example: SCNN

- Supports pruned CONV Layers



PE frontend

PE backend

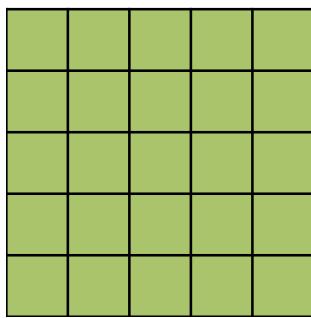
Input Stationary Dataflow

Compact Network Architecture Design

Build Network with series of Small Filters

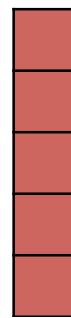
GoogLeNet / Inception v3

5x5 filter



decompose
→

5x1 filter

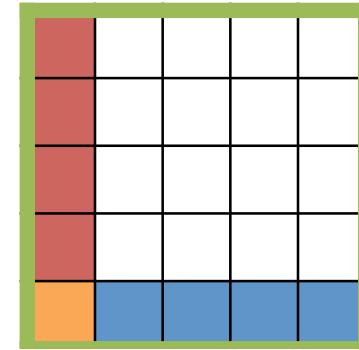


1x5 filter



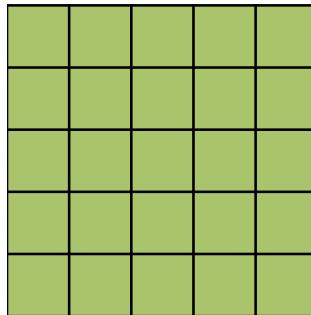
*separable
filters*

Apply sequentially



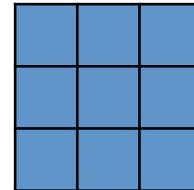
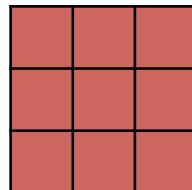
VGG-16

5x5 filter

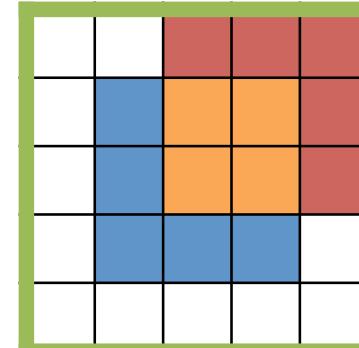


decompose
→

Two 3x3 filters

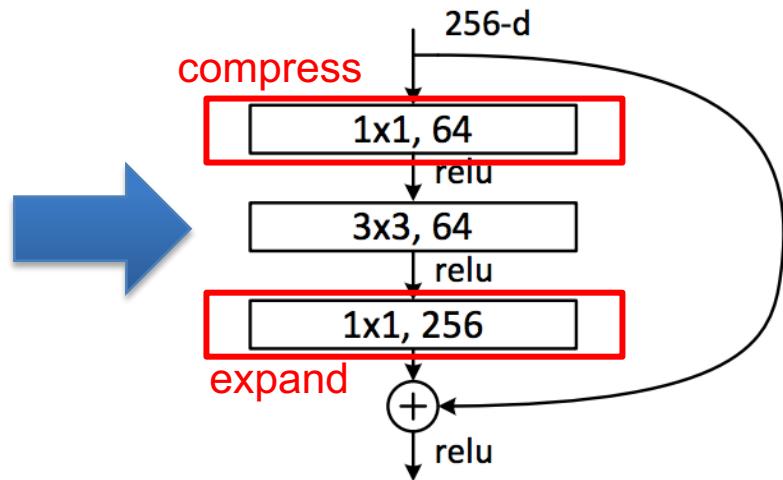
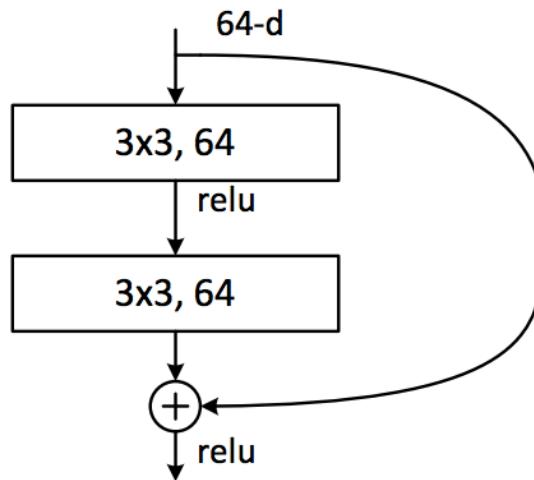


Apply sequentially

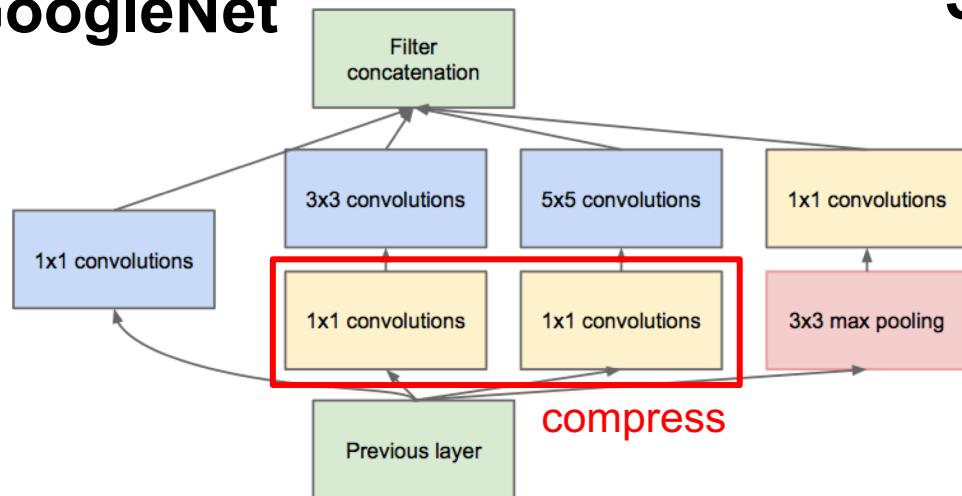


1×1 Bottleneck Layers are Popular

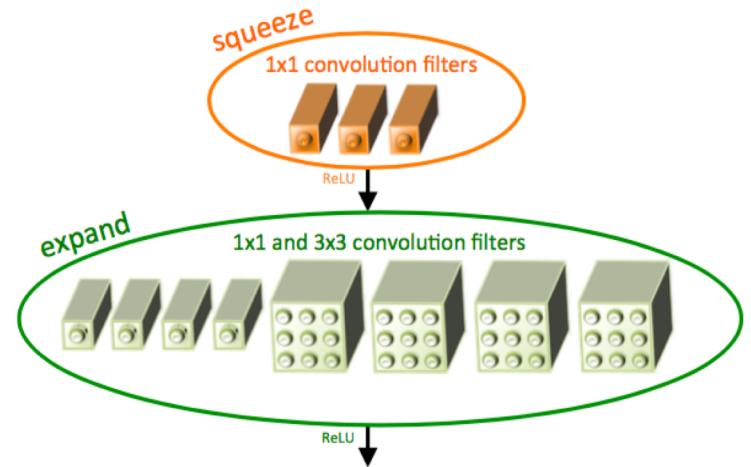
ResNet



GoogleNet



SqueezeNet



Summary

- **Dataflow** is one of the key attributes that define the architecture of DNN accelerators
- Co-design of **network model**, **dataflow**, and **hardware** is critical for the optimization of **performance**, **energy efficiency** and **flexibility** for DNN accelerators.
- It's important to consider a **comprehensive set of metrics** when evaluating different DNN solutions: **accuracy, flexibility, speed, energy, and cost**

Thank You

Resources:

Efficient Processing of Deep Neural Networks: A Tutorial and Survey

V.Sze, Y.-H. Chen, T.-J. Yang, J. Emer

Proceedings of IEEE

<https://arxiv.org/abs/1703.09039>

Eyeriss Website

More info on Eyeriss, Energy Modeling, Energy-Aware Pruning, NetAdapt, etc.

<http://eyeriss.mit.edu>