

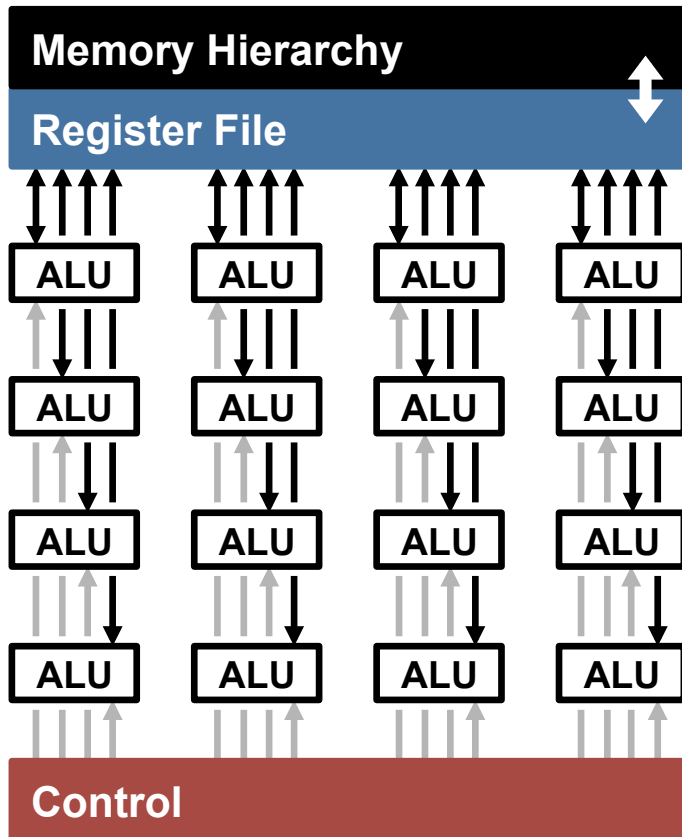
# DNN Accelerator Architectures

## CICS/MTL Tutorial (2017)

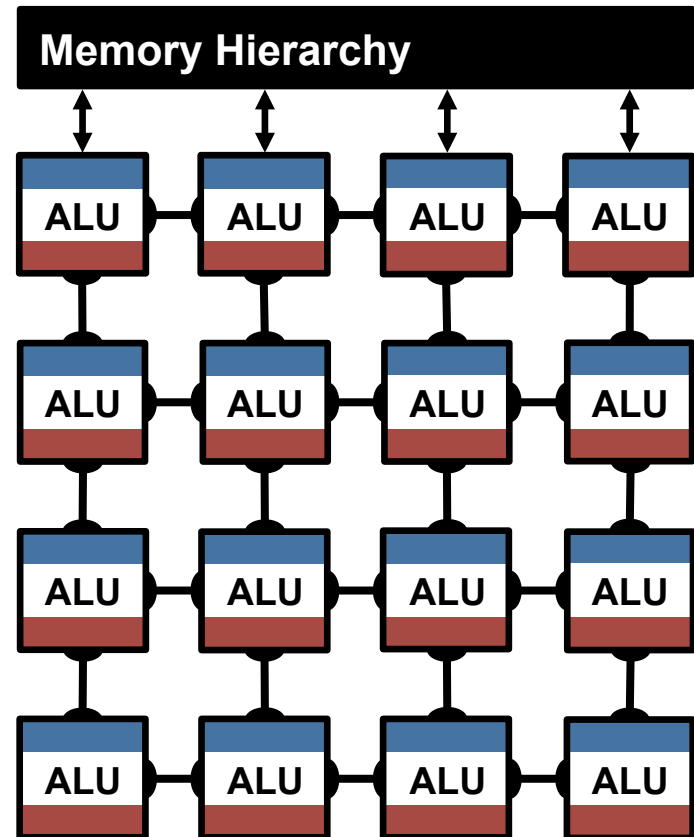
Website: <http://eyeriss.mit.edu/tutorial.html>

# Highly-Parallel Compute Paradigms

## Temporal Architecture (SIMD/SIMT)

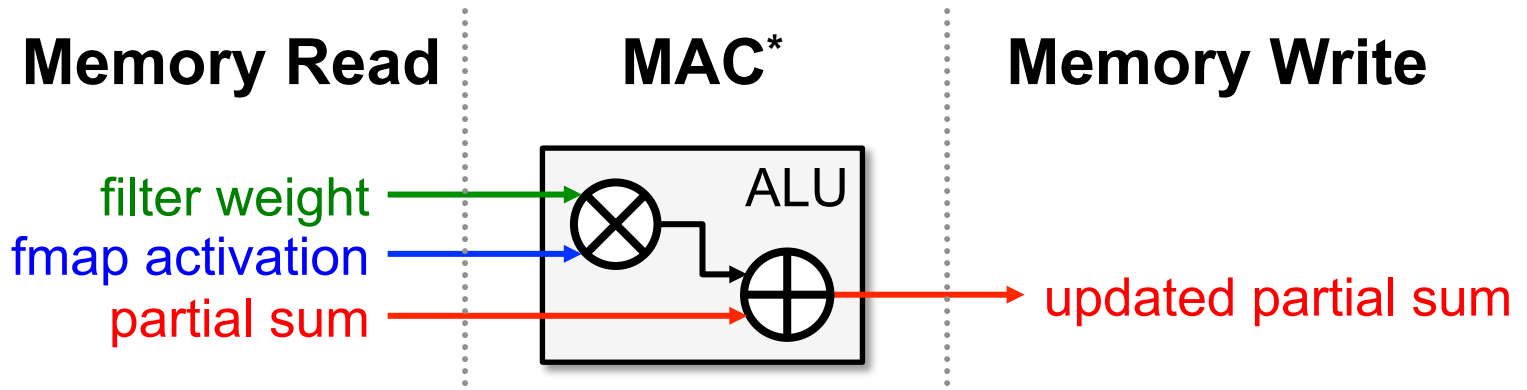


## Spatial Architecture (Dataflow Processing)



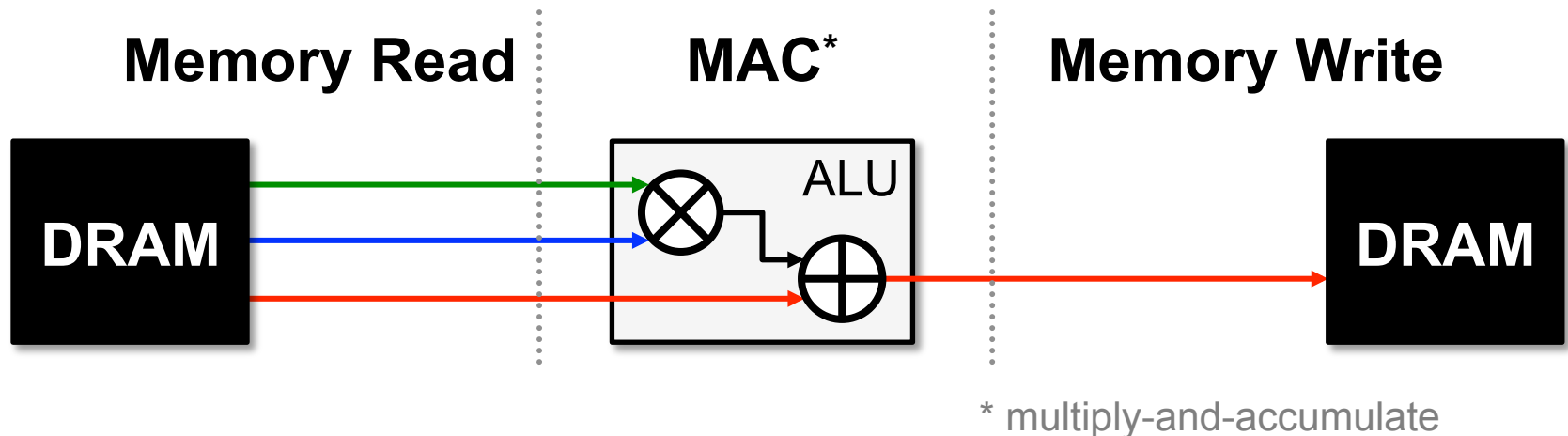
# Memory Access is the Bottleneck

---



\* multiply-and-accumulate

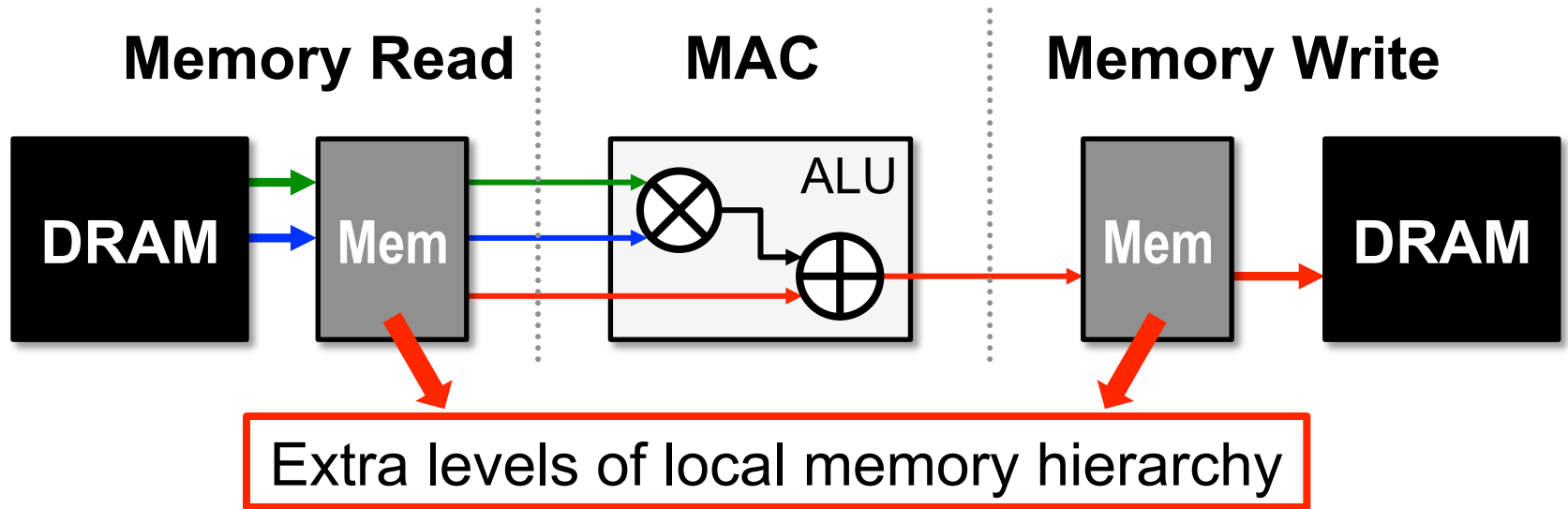
# Memory Access is the Bottleneck



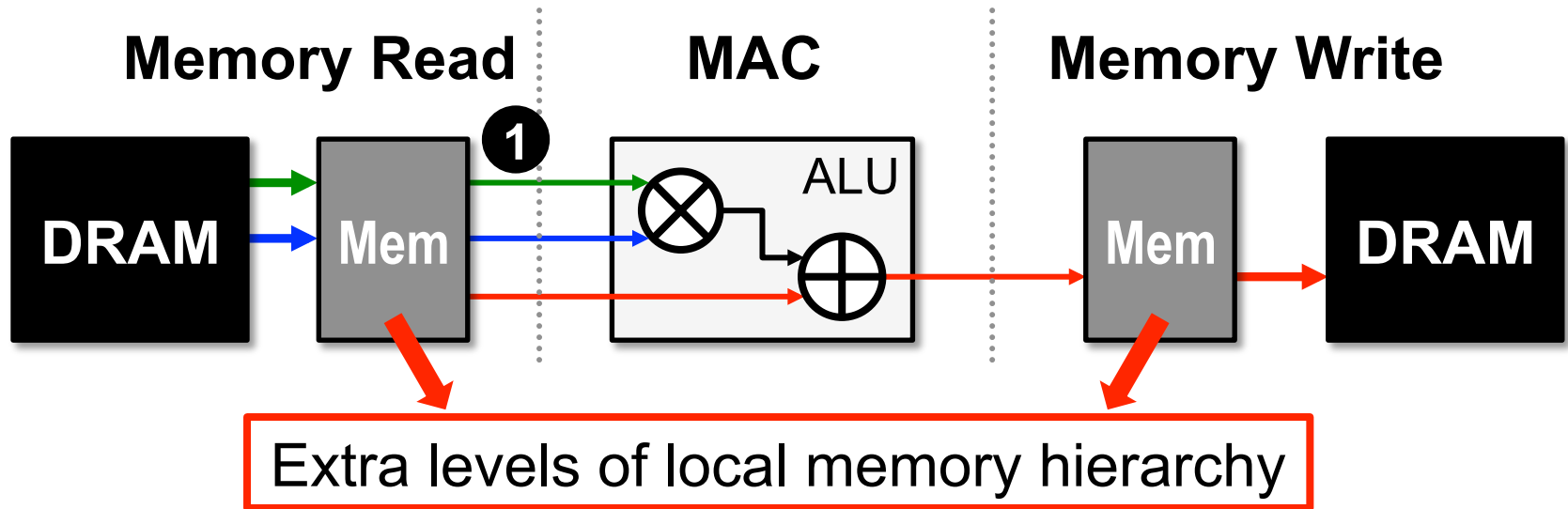
Worst Case: all memory R/W are **DRAM** accesses

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

# Memory Access is the Bottleneck



# Memory Access is the Bottleneck



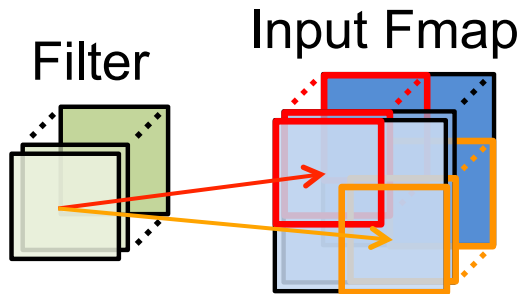
Opportunities: **1** data reuse

# Types of Data Reuse in DNN

---

## Convolutional Reuse

CONV layers only  
(sliding window)

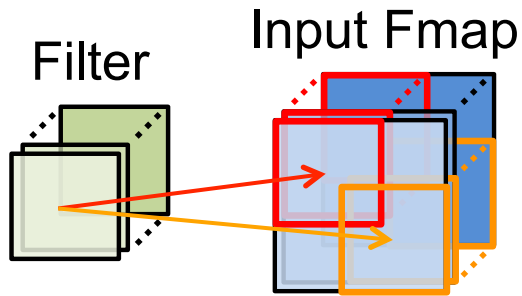


Reuse: **Activations**  
**Filter weights**

# Types of Data Reuse in DNN

## Convolutional Reuse

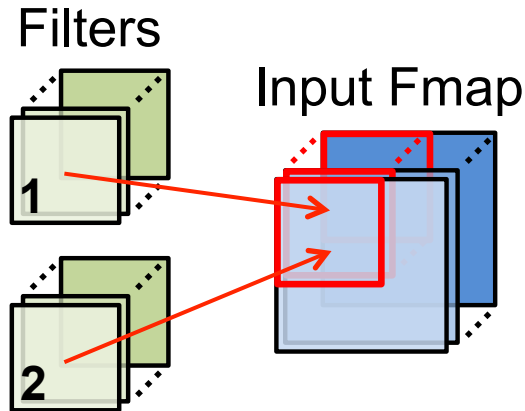
CONV layers only  
(sliding window)



Reuse: **Activations**  
**Filter weights**

## Fmap Reuse

CONV and FC layers



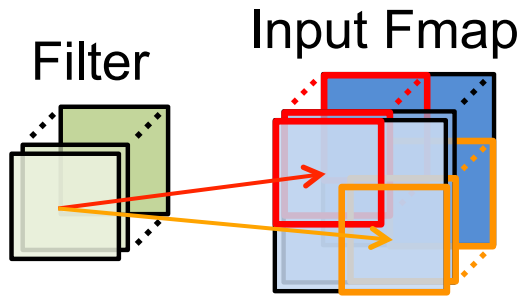
Reuse: **Activations**



# Types of Data Reuse in DNN

## Convolutional Reuse

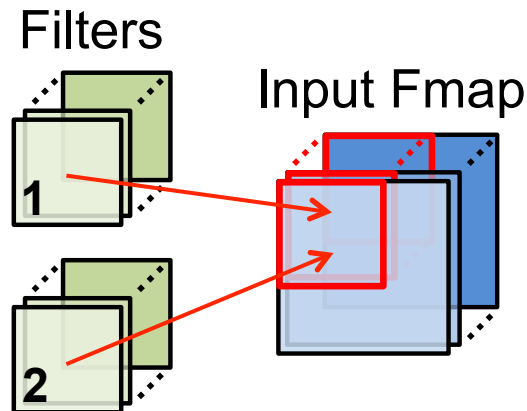
CONV layers only  
(sliding window)



Reuse: **Activations**  
**Filter weights**

## Fmap Reuse

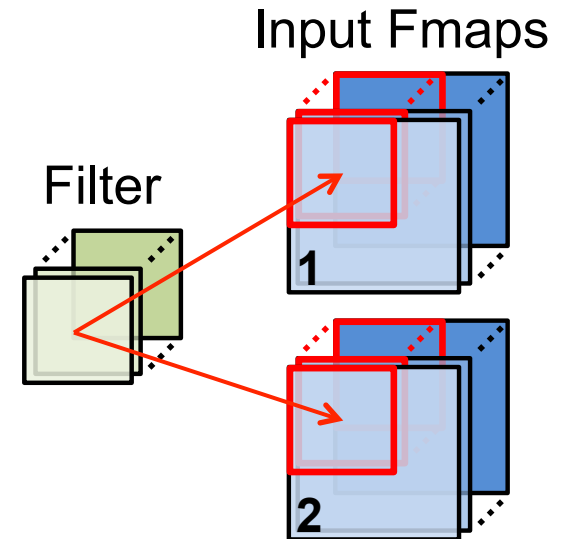
CONV and FC layers



Reuse: **Activations**

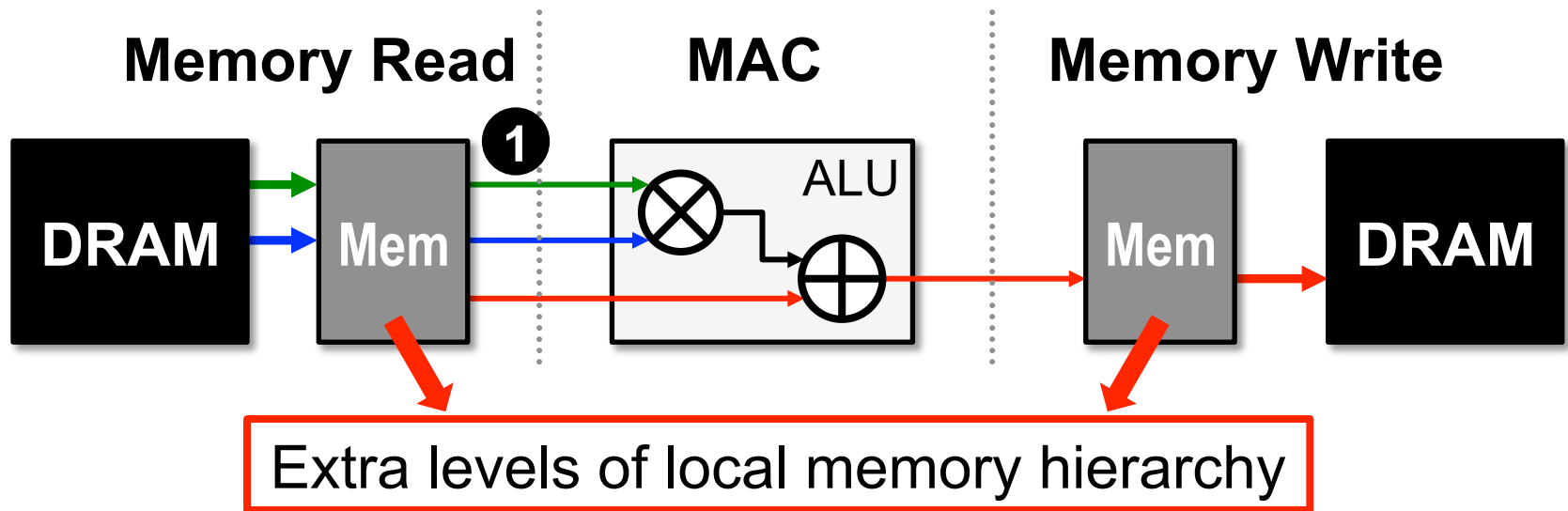
## Filter Reuse

CONV and FC layers  
(batch size > 1)



Reuse: **Filter weights**

# Memory Access is the Bottleneck

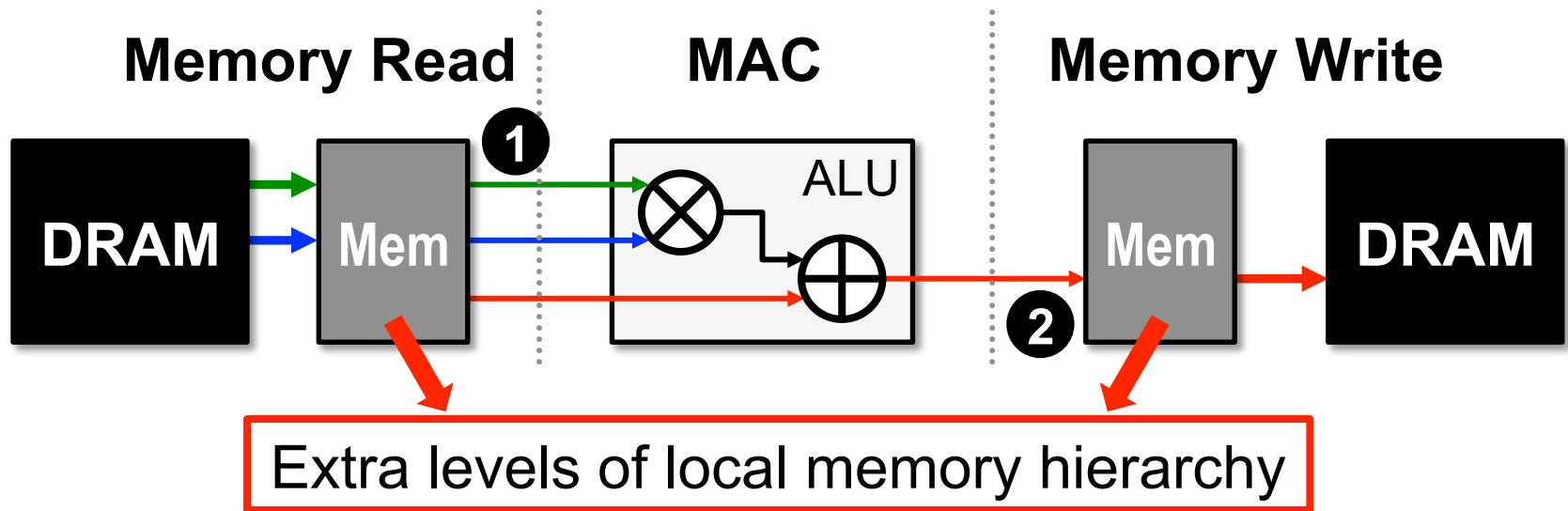


Opportunities: **1** data reuse

- 1** Can reduce DRAM reads of *filter/fmap* by up to **500x**\*\*

\*\* AlexNet CONV layers

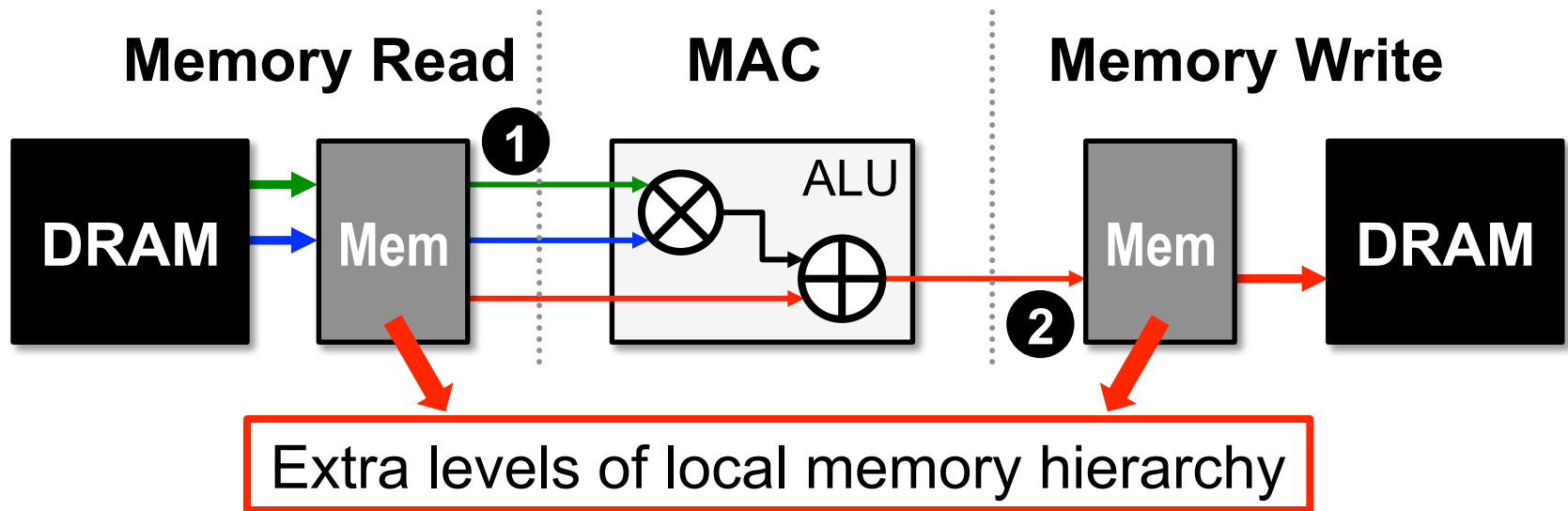
# Memory Access is the Bottleneck



Opportunities:    ① data reuse    ② local accumulation

- ① Can reduce DRAM reads of *filter/fmap* by up to 500×
- ② *Partial sum* accumulation does **NOT** have to access DRAM

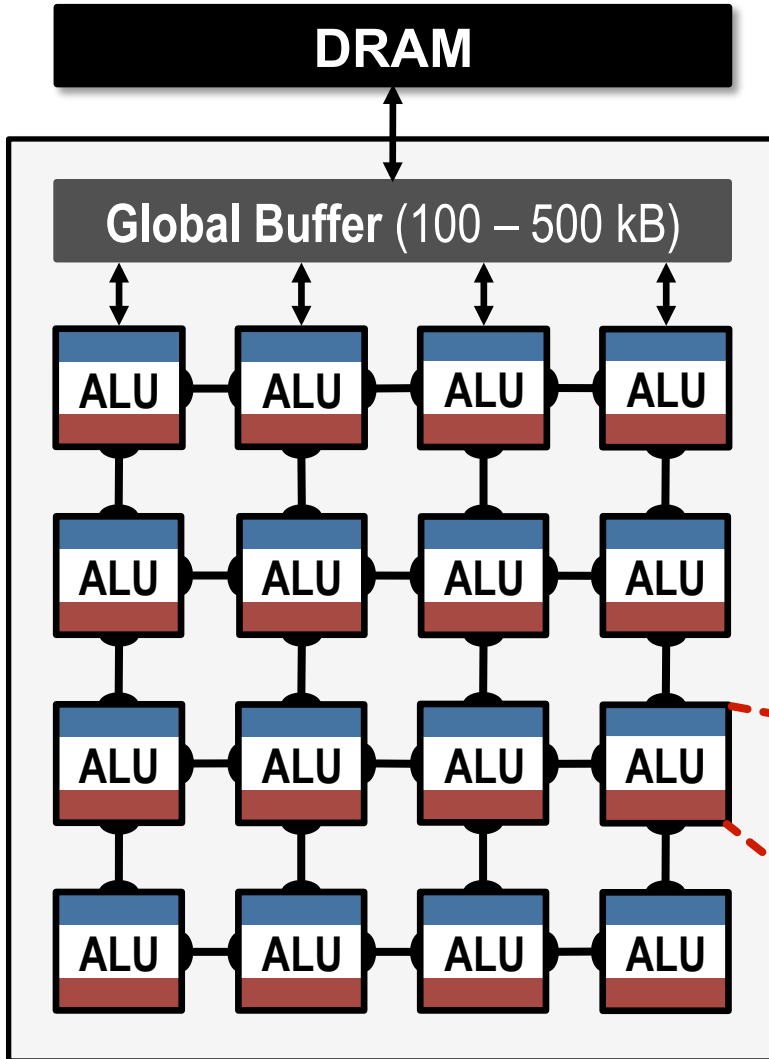
# Memory Access is the Bottleneck



Opportunities:    ① data reuse    ② local accumulation

- ① Can reduce DRAM reads of *filter/fmap* by up to **500×**
- ② **Partial sum** accumulation does **NOT** have to access DRAM
  - Example:      DRAM access in AlexNet can be reduced from **2896M** to **61M** (best case)

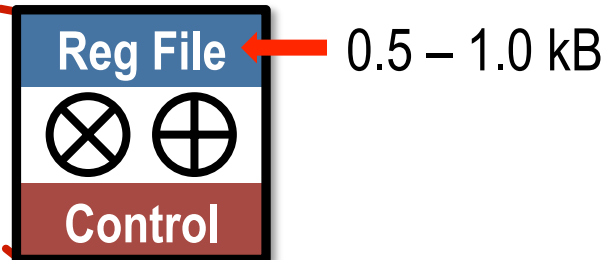
# Spatial Architecture for DNN



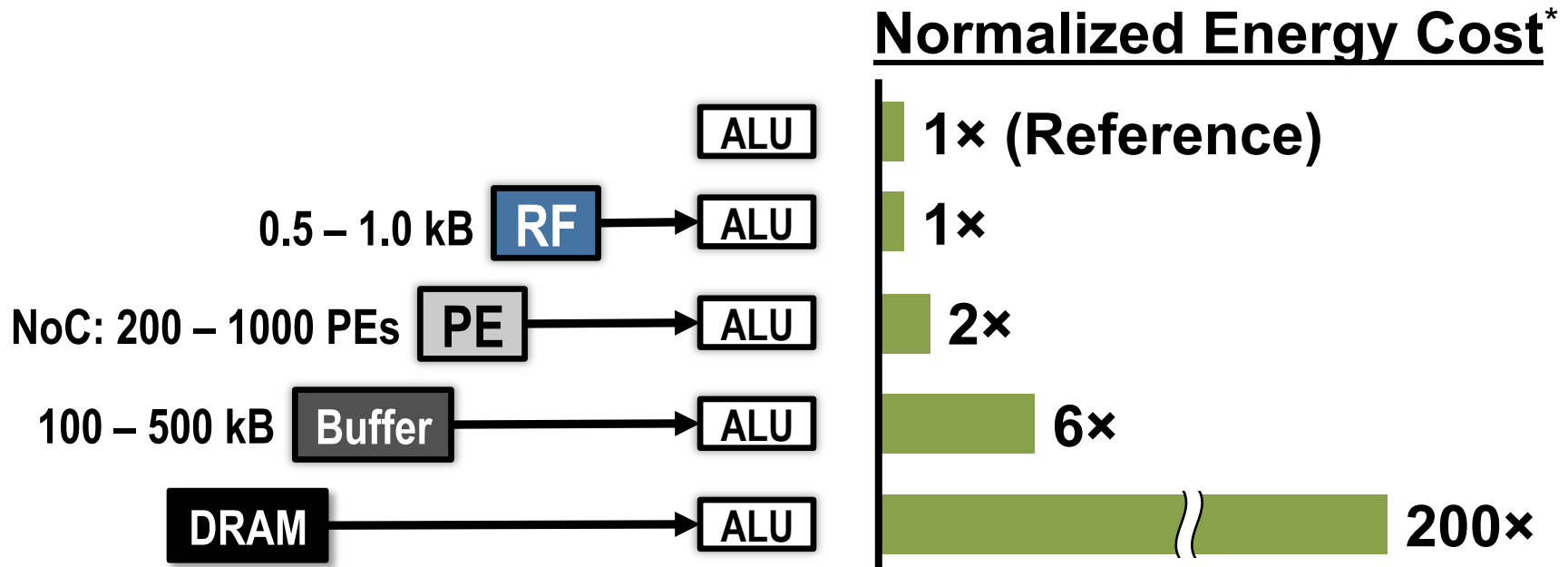
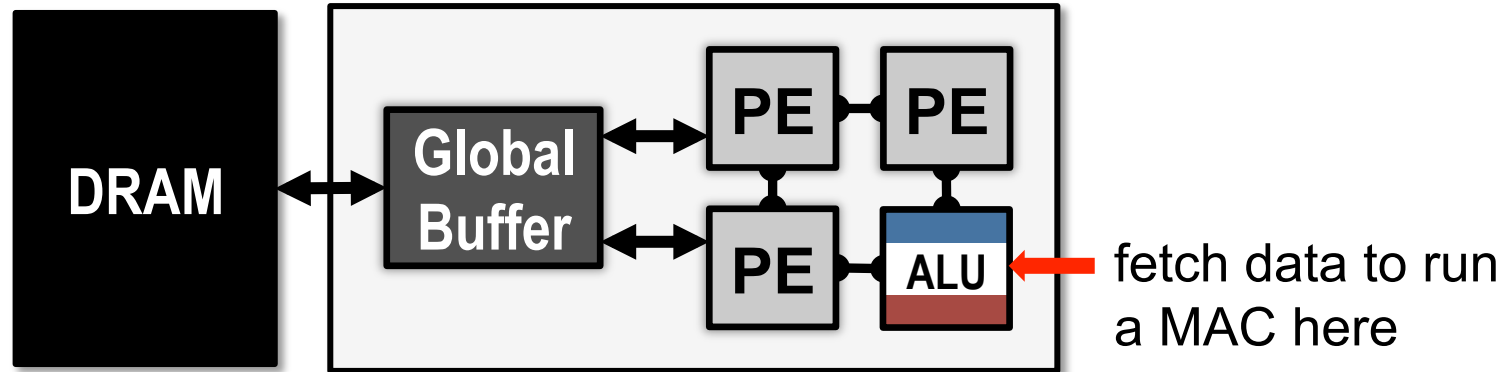
## Local Memory Hierarchy

- Global Buffer
- Direct inter-PE network
- PE-local memory (RF)

## Processing Element (PE)



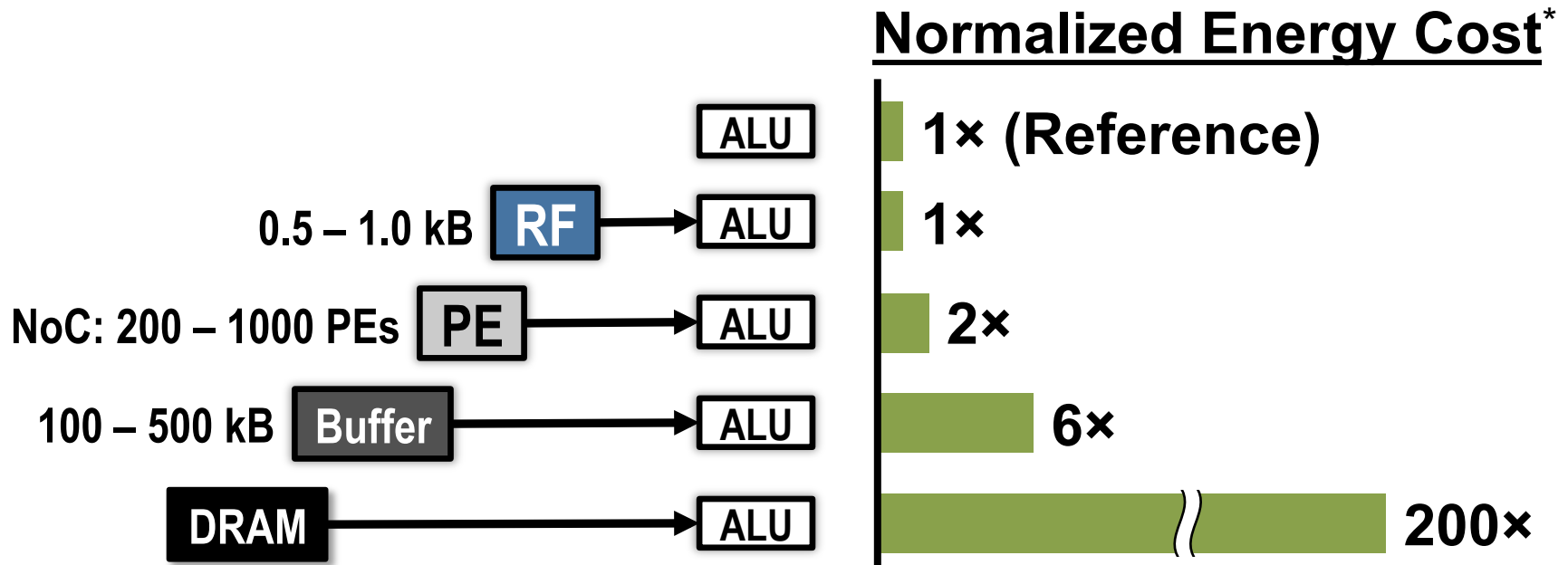
# Low-Cost Local Data Access



\* measured from a commercial 65nm process

# Low-Cost Local Data Access

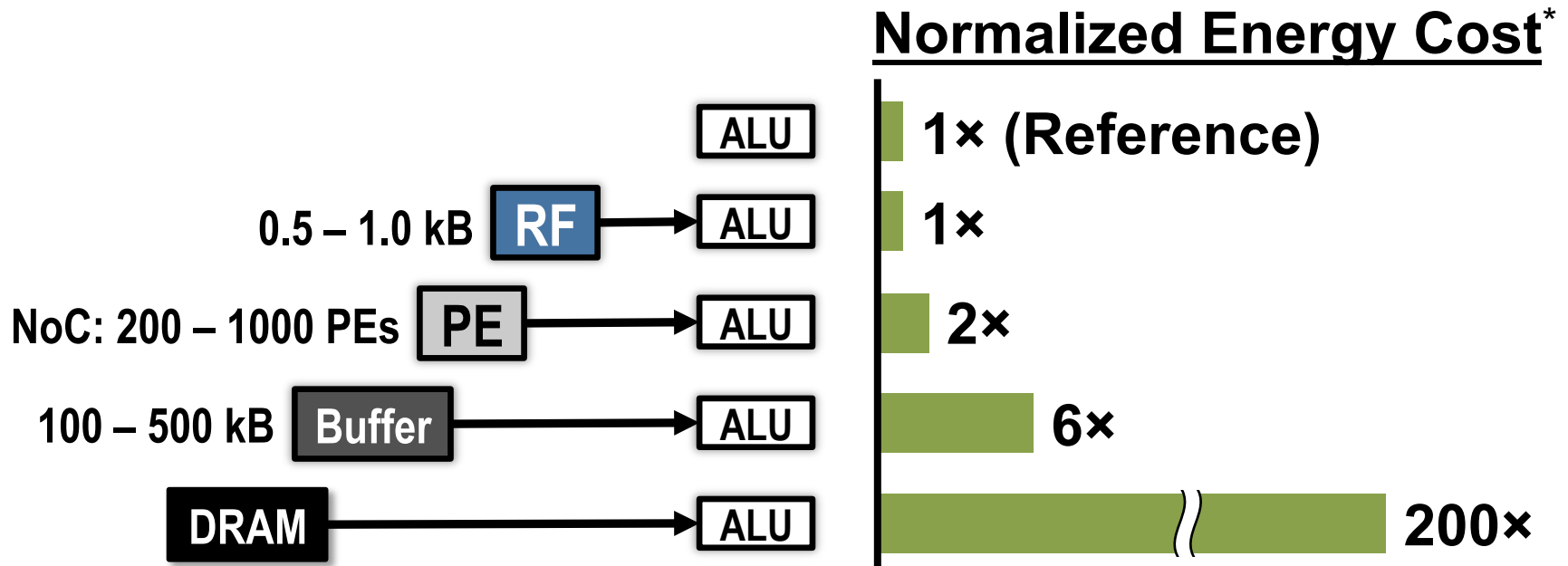
How to exploit **1** data reuse and **2** local accumulation with *limited* low-cost local storage?



# Low-Cost Local Data Access

How to exploit ❶ **data reuse** and ❷ **local accumulation** with *limited* low-cost local storage?

specialized **processing dataflow** required!

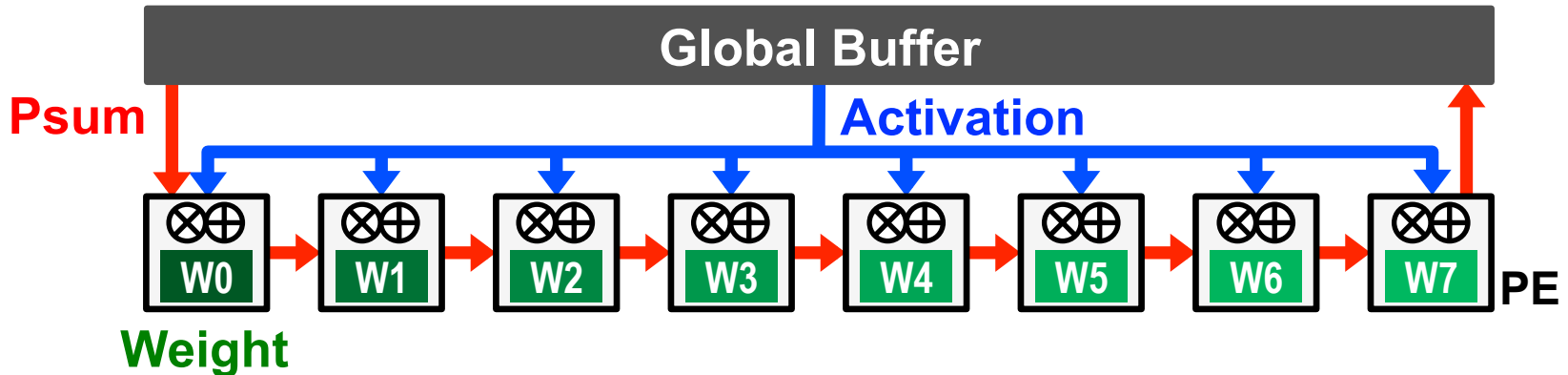




# Dataflow Taxonomy

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

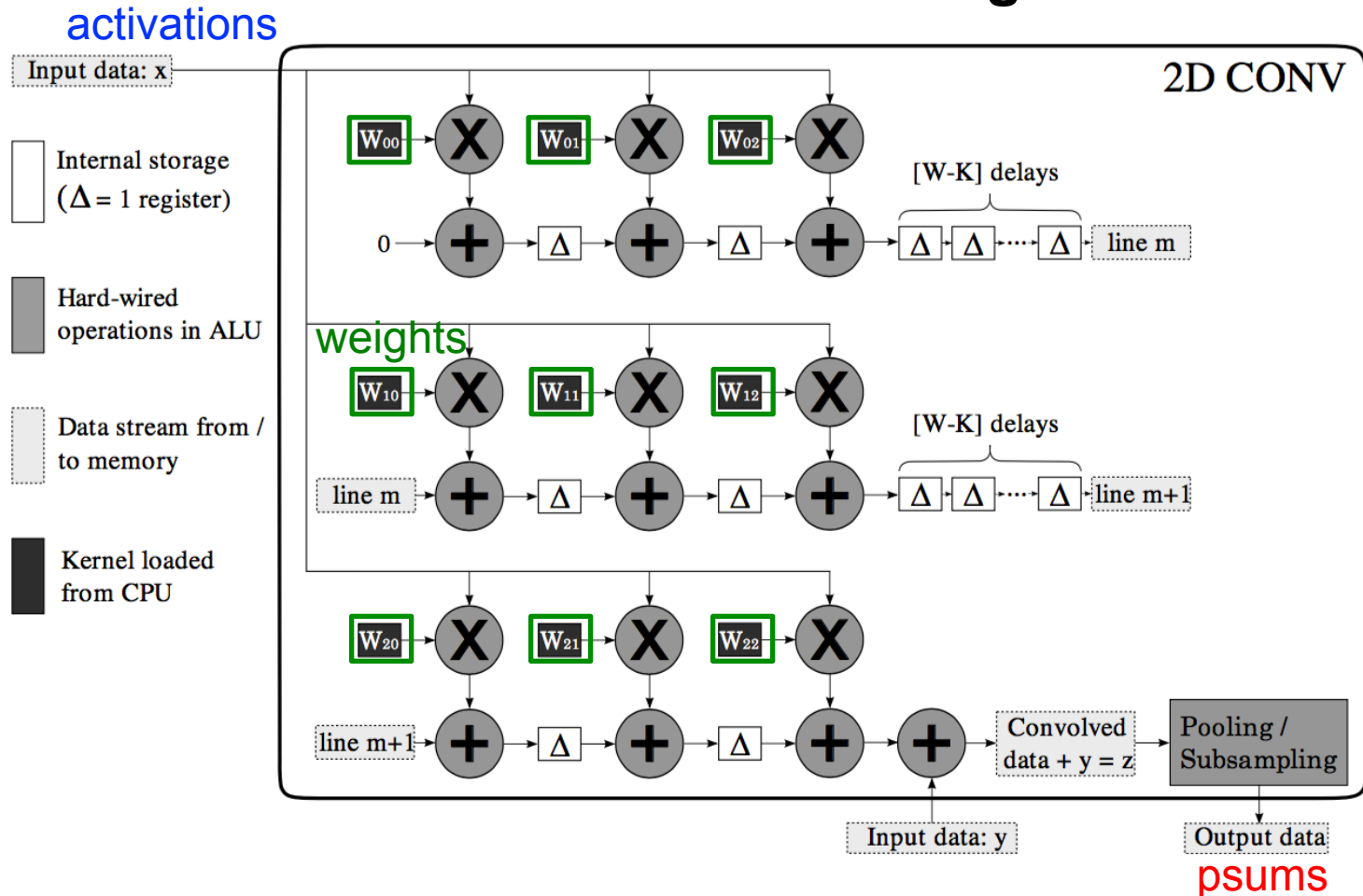
# Weight Stationary (WS)



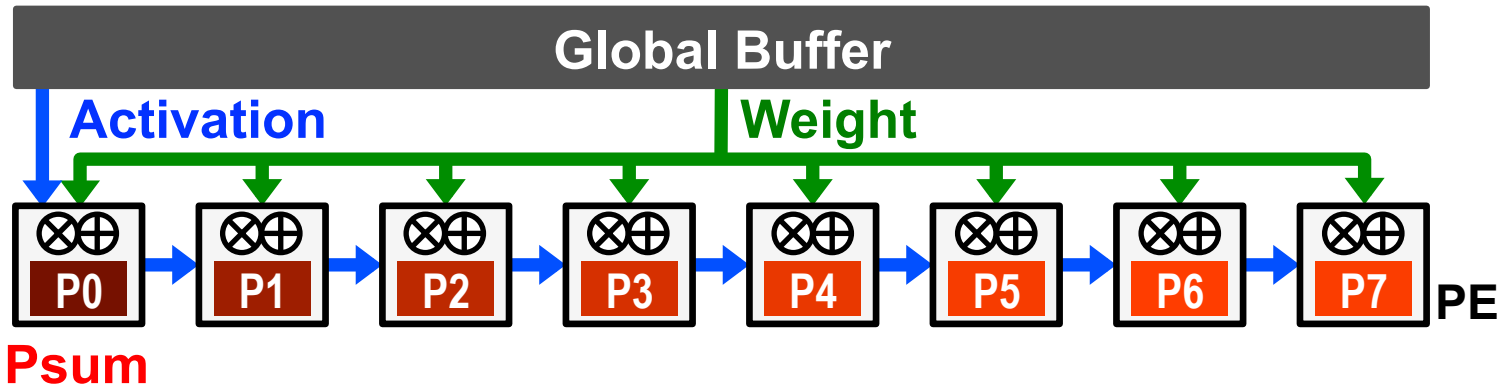
- **Minimize weight** read energy consumption
  - maximize convolutional and filter reuse of weights
- **Broadcast activations** and **accumulate psums** spatially across the PE array.

# WS Example: nn-X (NeuFlow)

## A 3x3 2D Convolution Engine



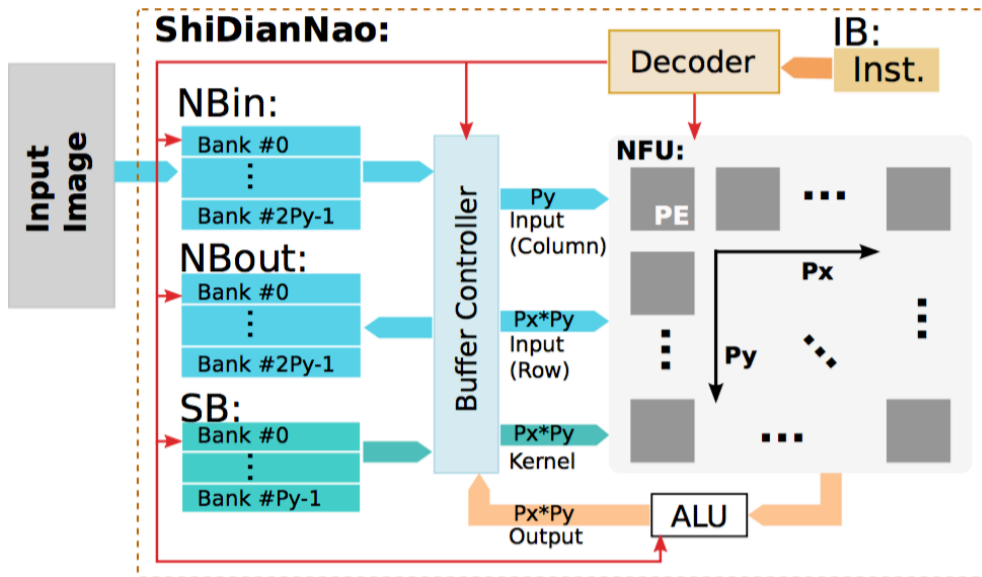
# Output Stationary (OS)



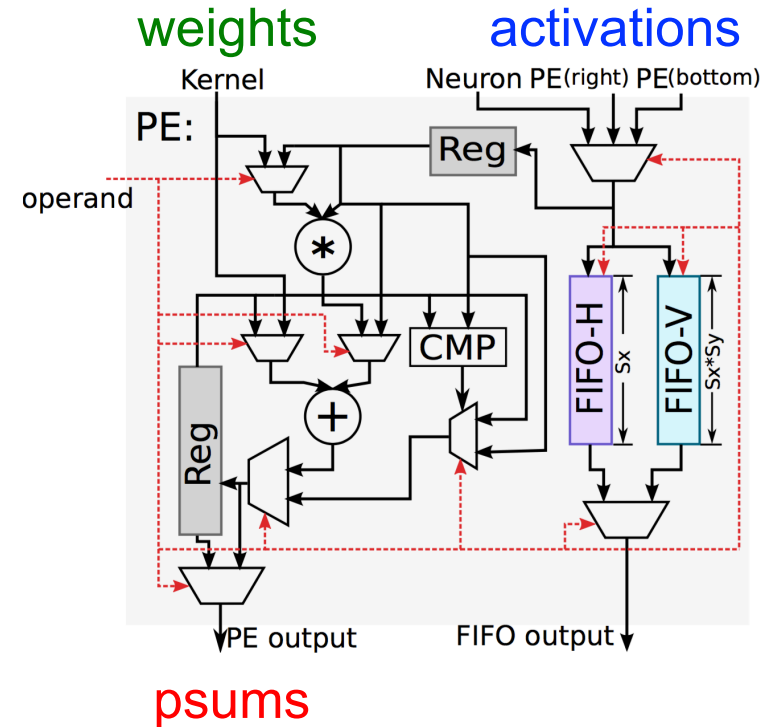
- **Minimize partial sum** R/W energy consumption
  - maximize local accumulation
- **Broadcast/Multicast filter weights** and reuse **activations** spatially across the PE array

# OS Example: ShiDianNao

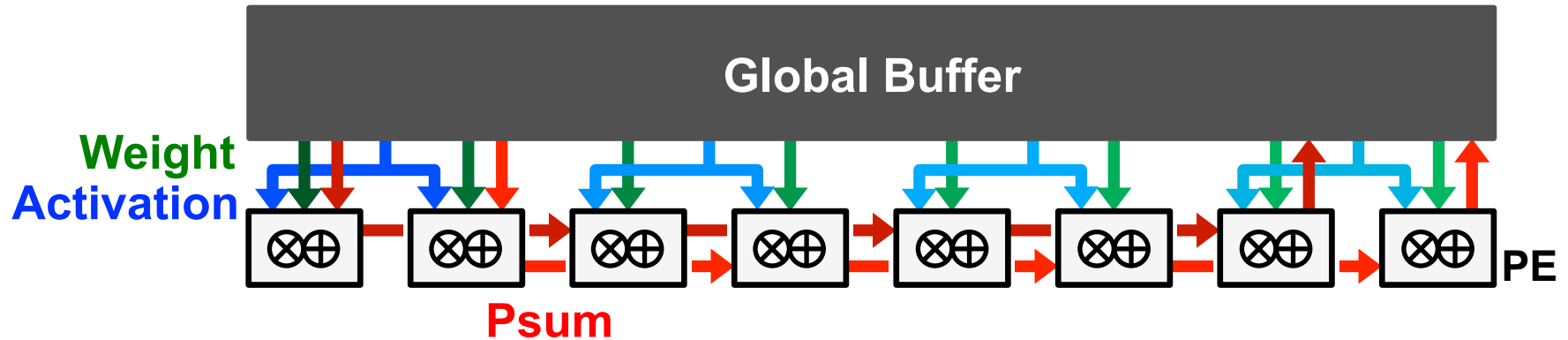
## Top-Level Architecture



## PE Architecture

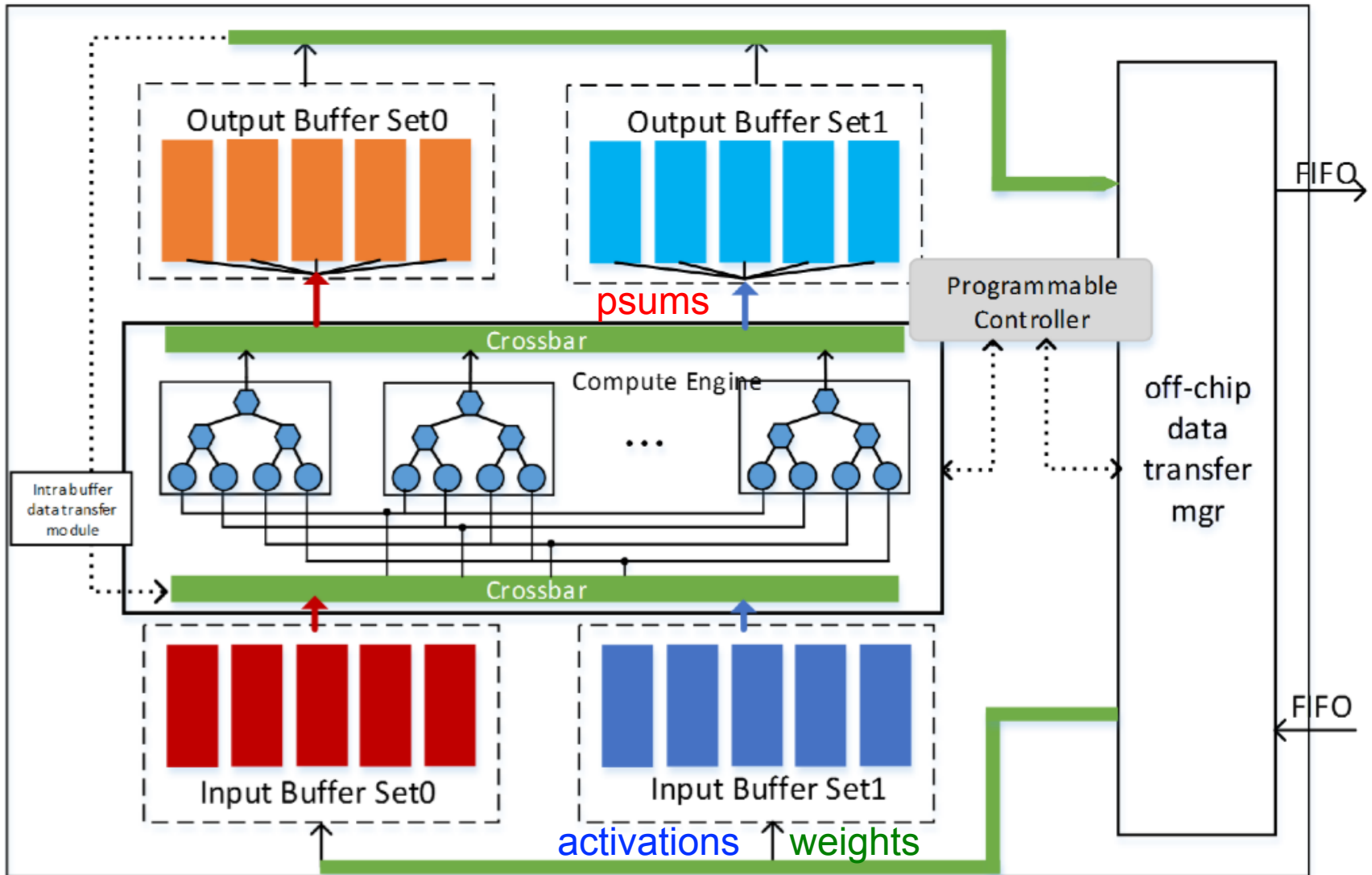


# No Local Reuse (NLR)



- Use a **large global buffer** as shared storage
  - Reduce **DRAM** access energy consumption
- **Multicast activations**, **single-cast weights**, and **accumulate psums** spatially across the PE array

# NLR Example: UCLA



# Taxonomy: More Examples

---

- **Weight Stationary (WS)**

[Chakradhar, *ISCA* 2010] [nn-X (NeuFlow), *CVPRW* 2014]

[Park, *ISSCC* 2015] [ISAAC, *ISCA* 2016] [PRIME, *ISCA* 2016]

- **Output Stationary (OS)**

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

[Gupta, *ICML* 2015] [Moons, *VLSI* 2016]

- **No Local Reuse (NLR)**

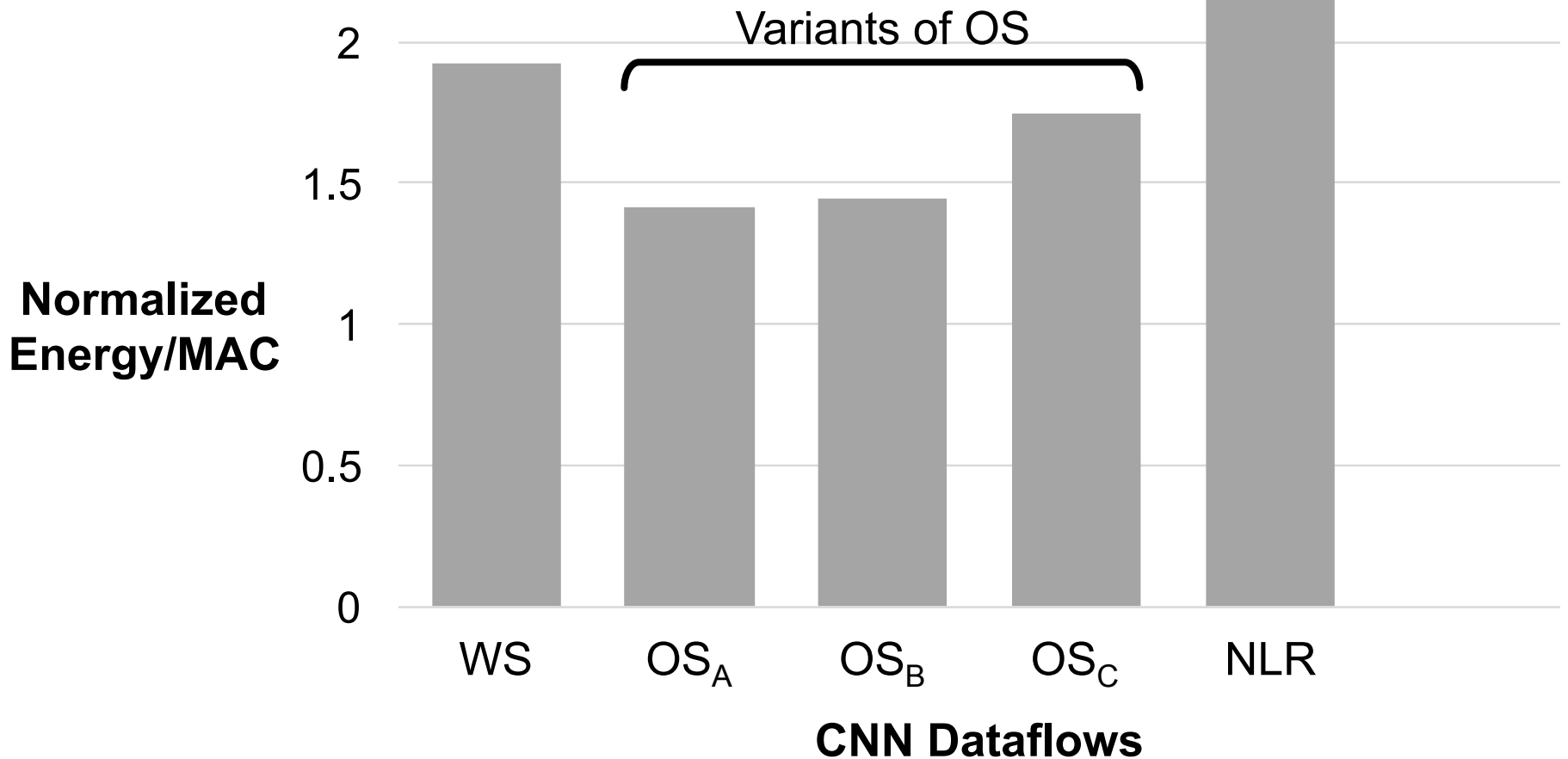
[DianNao, *ASPLOS* 2014] [DaDianNao, *MICRO* 2014]

[Zhang, *FPGA* 2015]



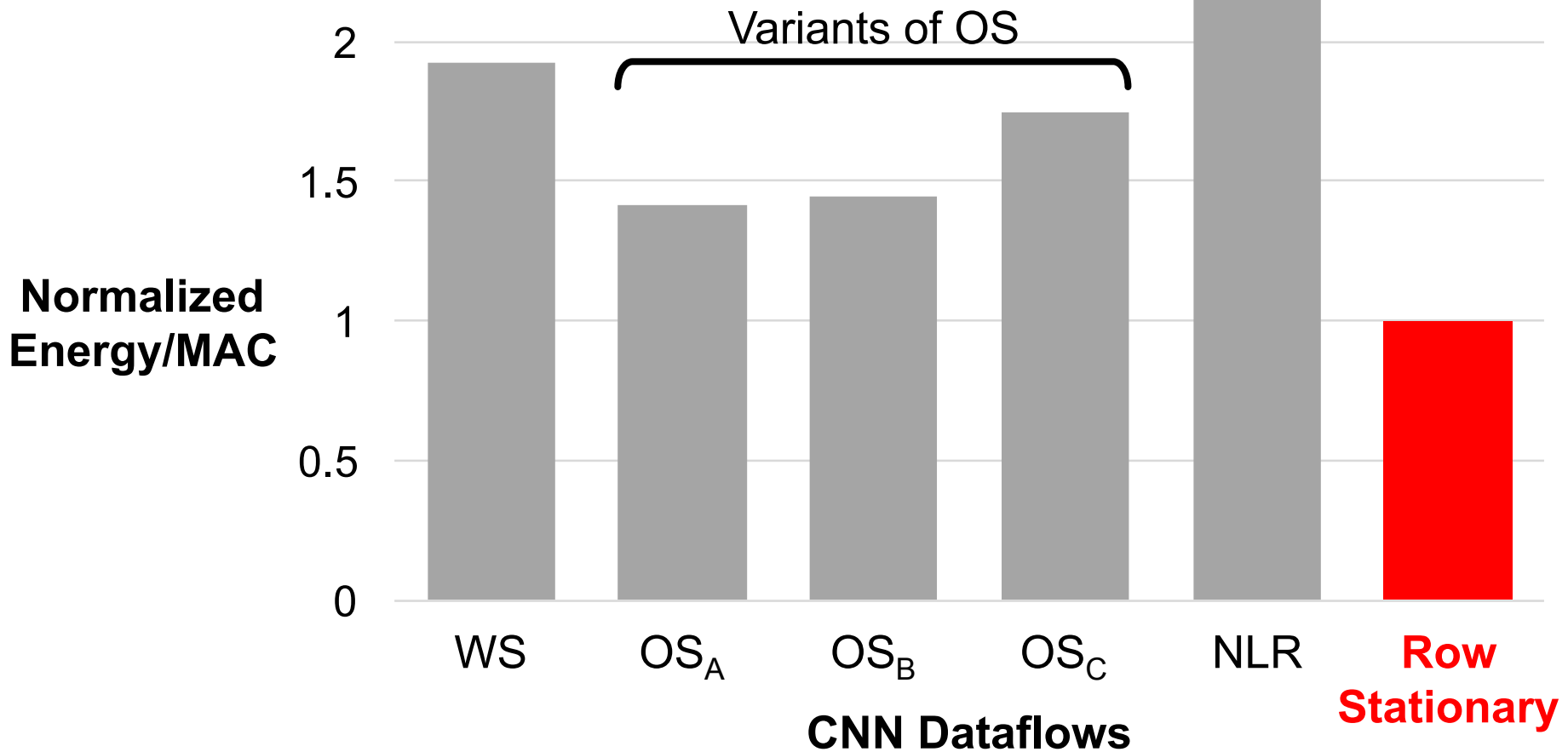
# Energy Efficiency Comparison

- 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

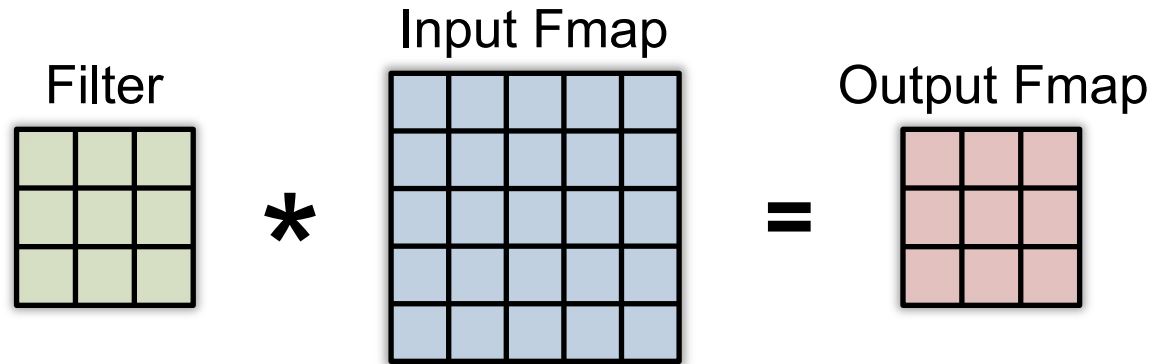


# Energy-Efficient Dataflow: Row Stationary (RS)

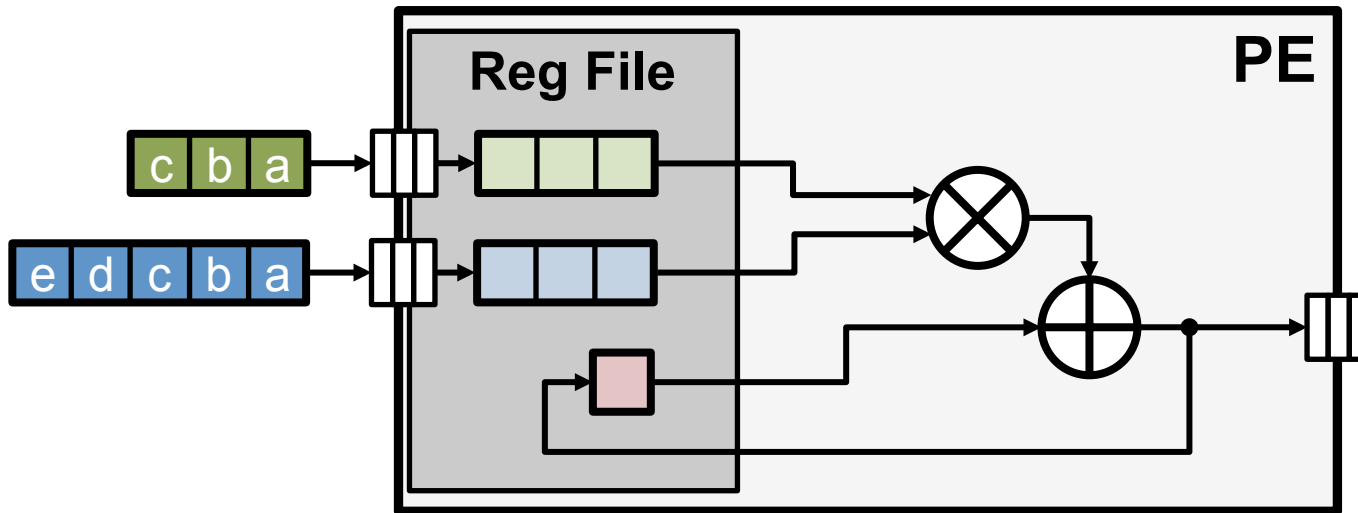
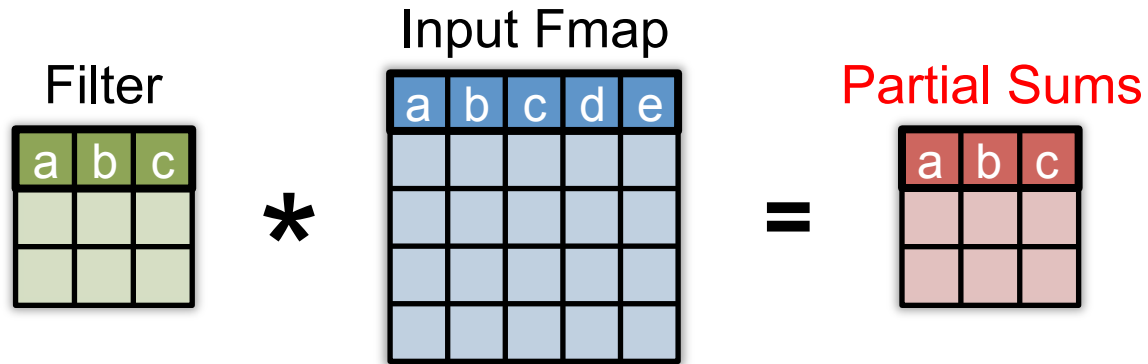
- **Maximize** reuse and accumulation at **RF**
- Optimize for **overall** energy efficiency instead for *only* a certain data type

# Row Stationary: Energy-efficient Dataflow

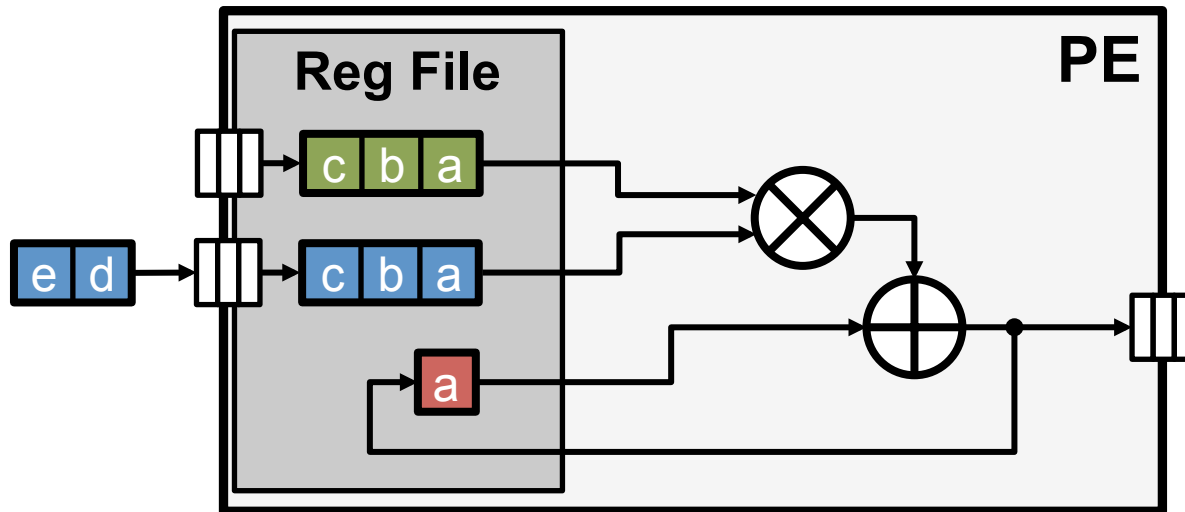
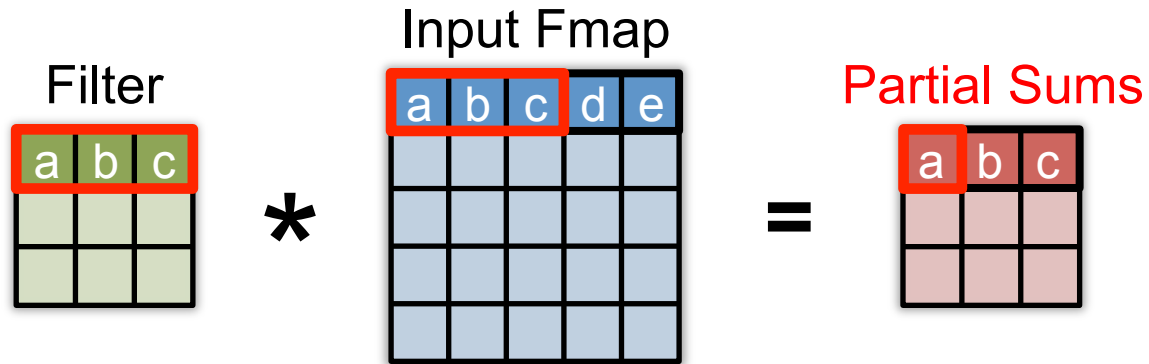
---



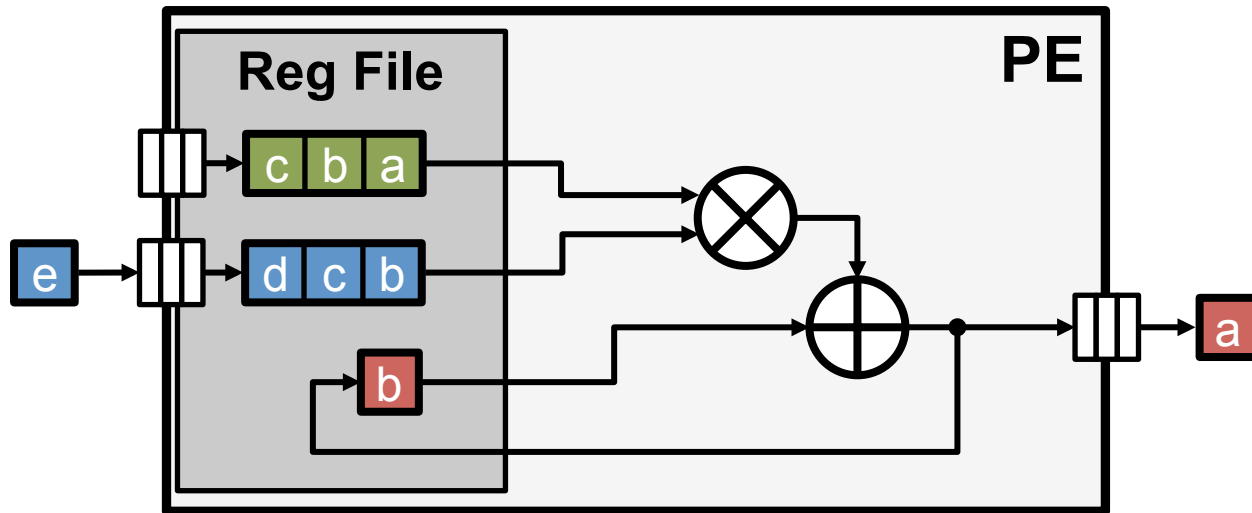
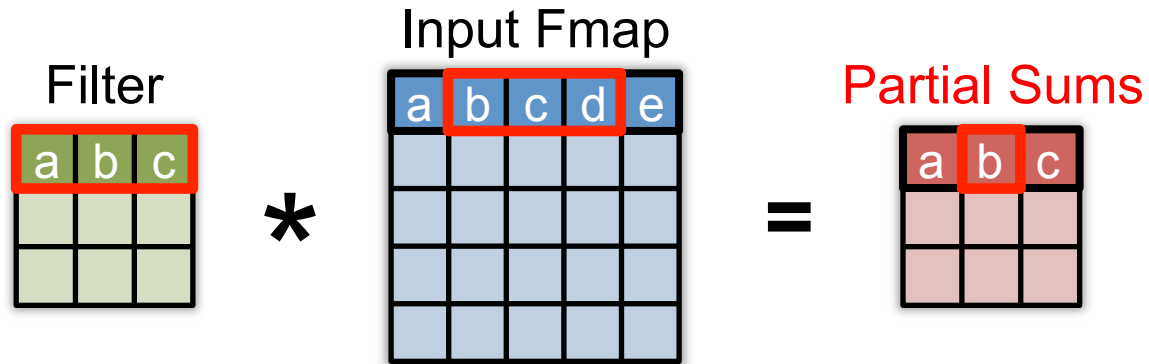
# 1D Row Convolution in PE



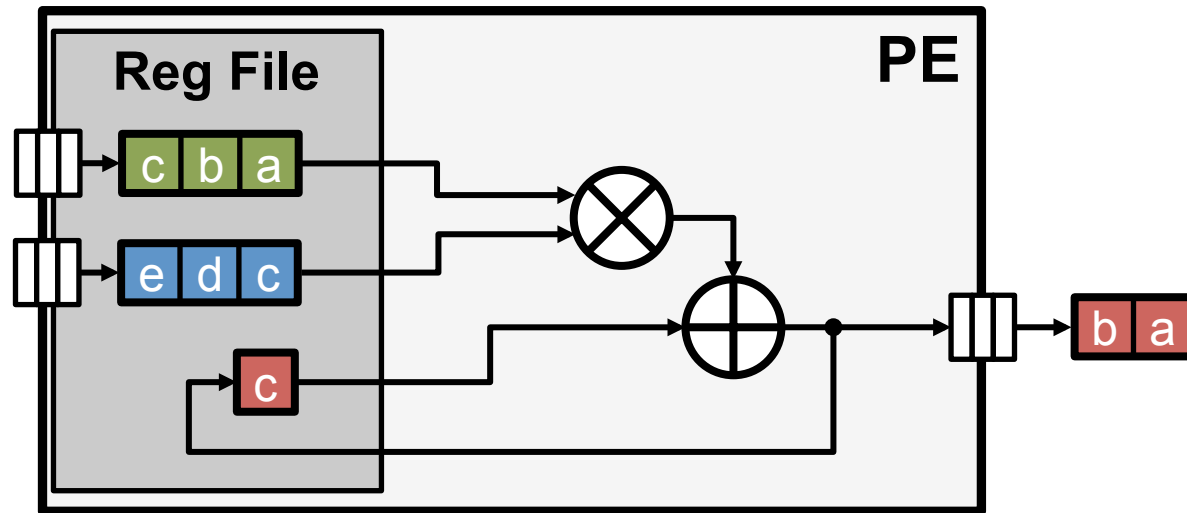
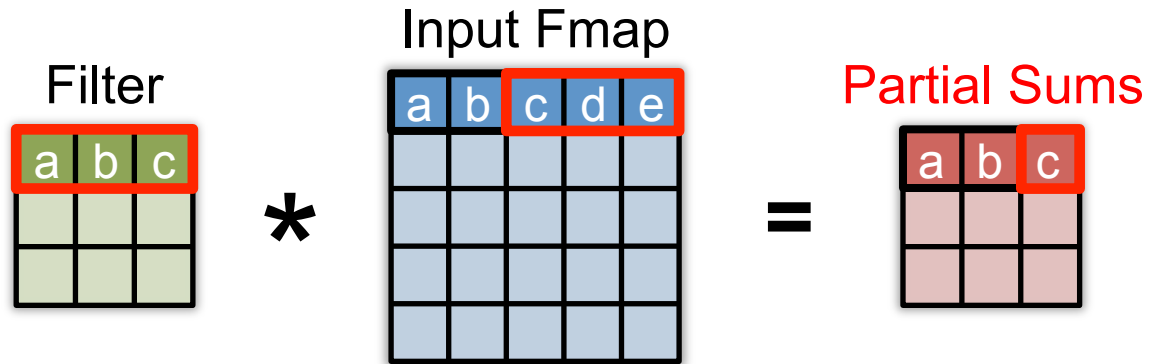
# 1D Row Convolution in PE



# 1D Row Convolution in PE



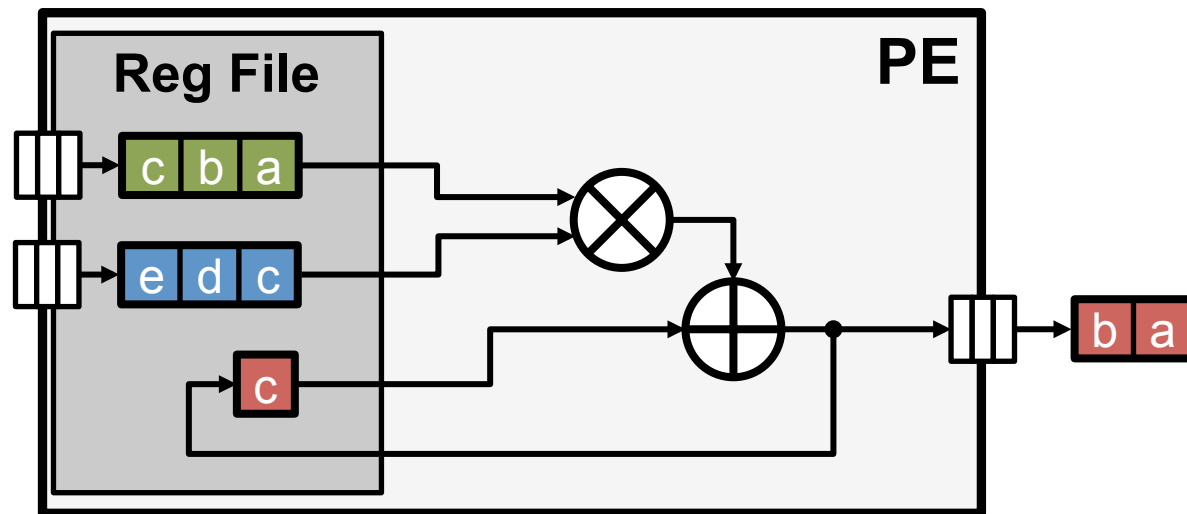
# 1D Row Convolution in PE





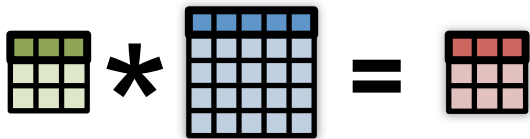
# 1D Row Convolution in PE

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



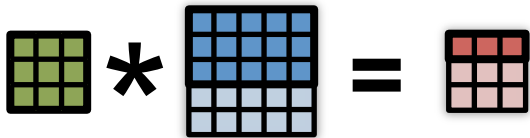
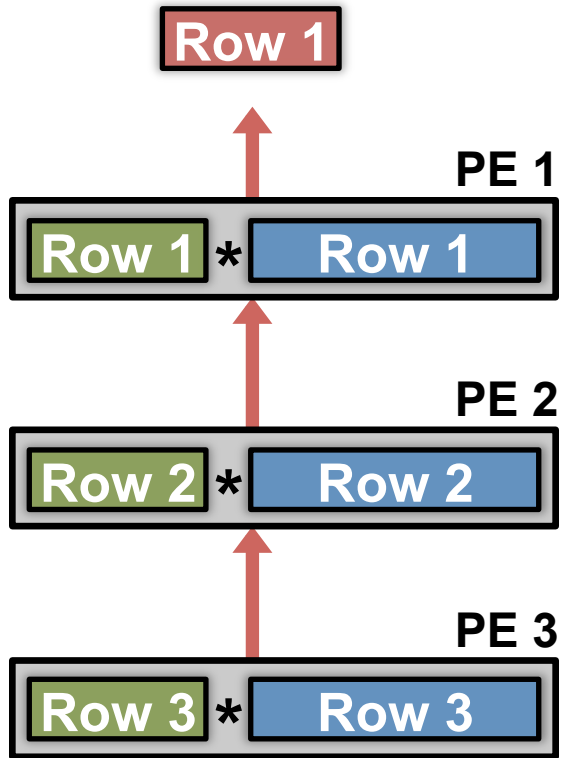
# 2D Convolution in PE Array

---

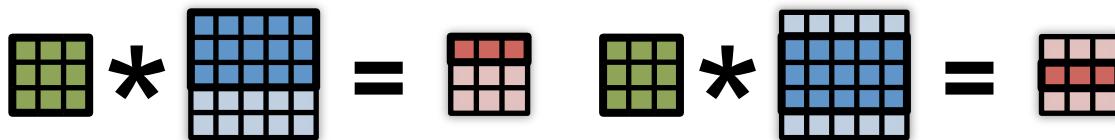
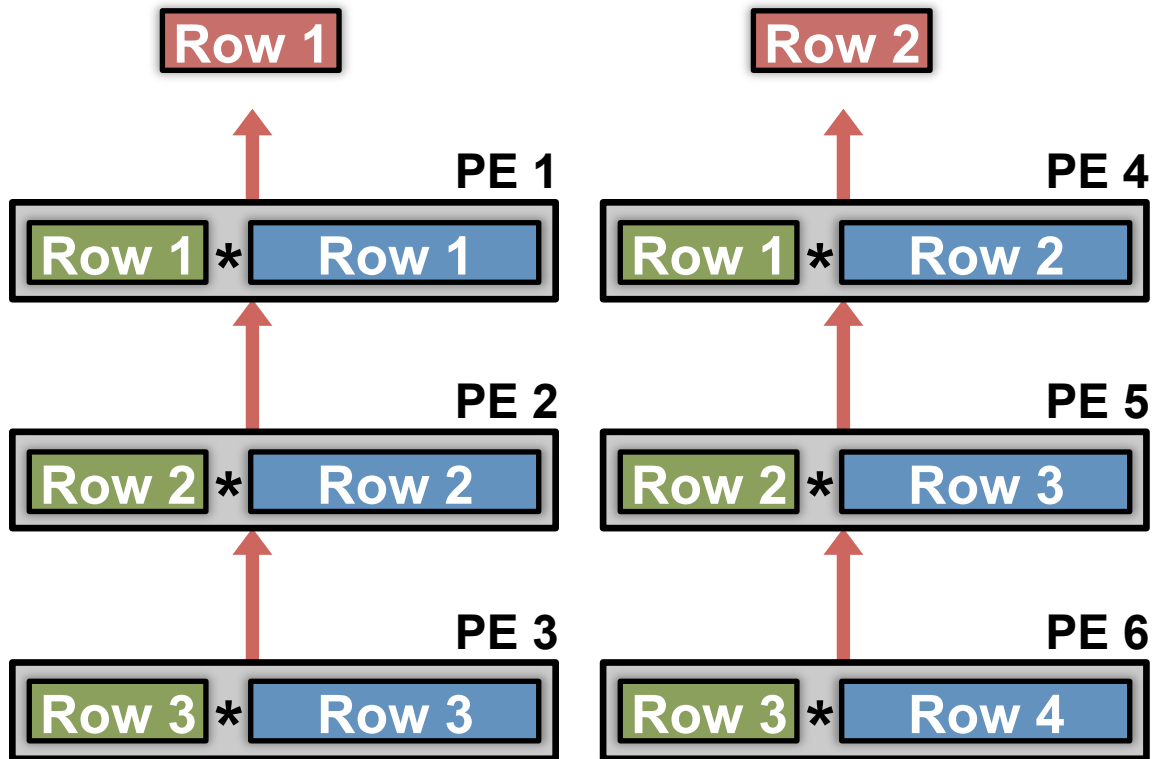


# 2D Convolution in PE Array

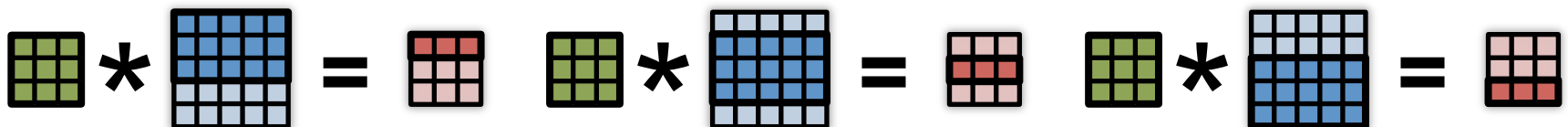
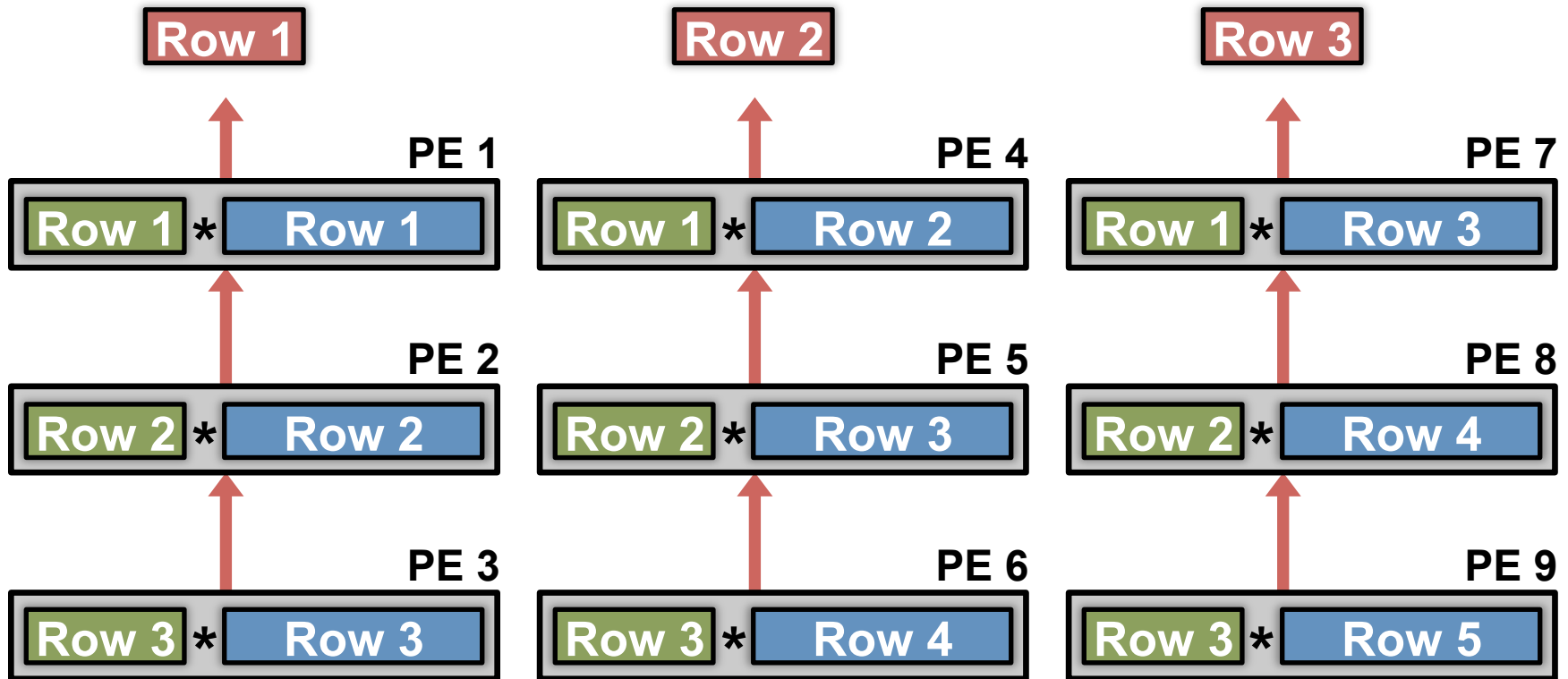
---



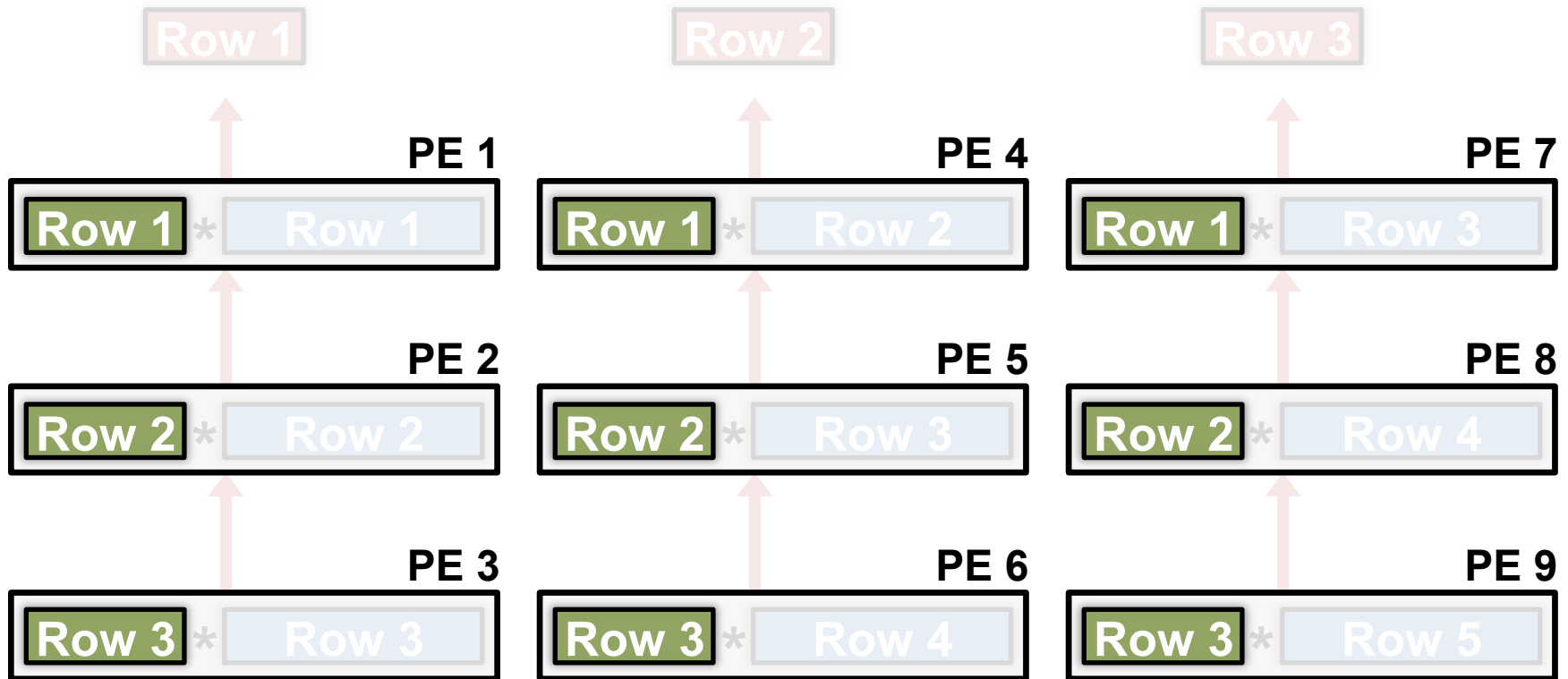
# 2D Convolution in PE Array



# 2D Convolution in PE Array

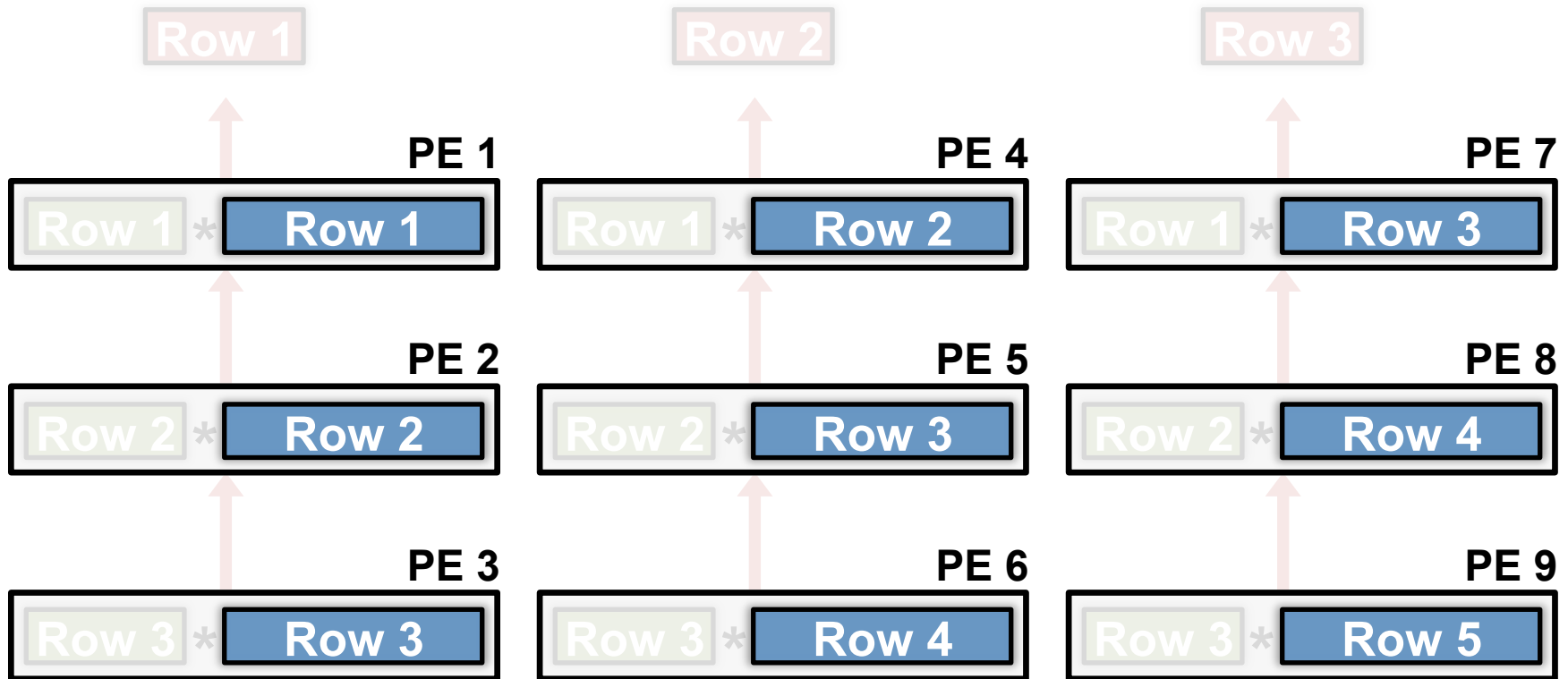


# Convolutional Reuse Maximized



**Filter rows** are reused across PEs **horizontally**

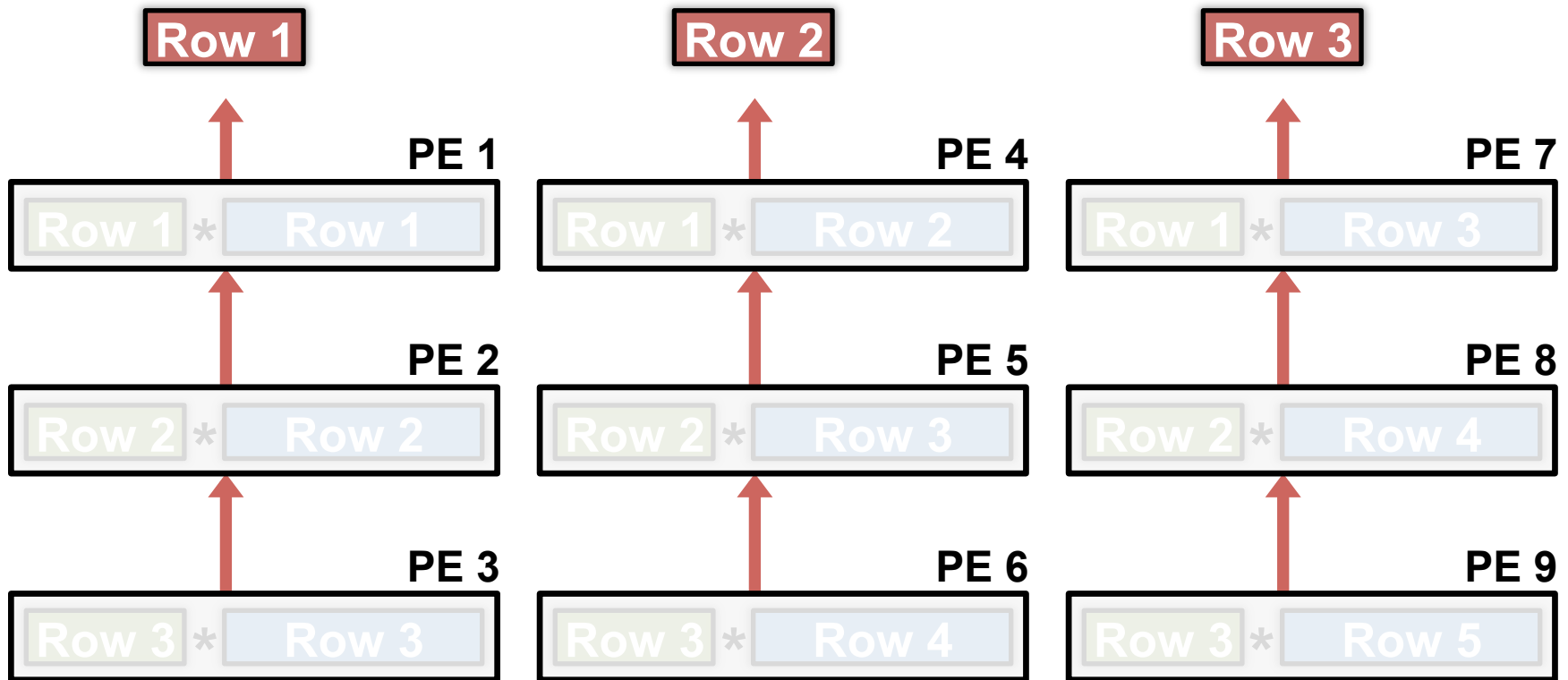
# Convolutional Reuse Maximized



**Fmap rows** are reused across PEs **diagonally**

# Maximize 2D Accumulation in PE Array

---



**Partial sums** accumulate across PEs **vertically**



# Dimensions Beyond 2D Convolution

---

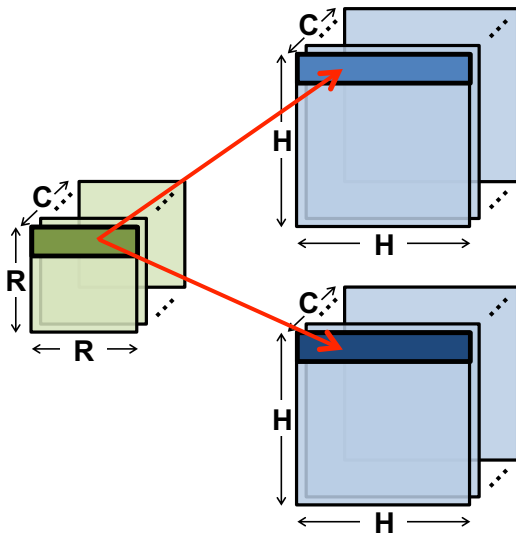
- 1 Multiple Fmaps
- 2 Multiple Filters
- 3 Multiple Channels

# Filter Reuse in PE

## 1 Multiple Fmaps

## 2 Multiple Filters

## 3 Multiple Channels



Channel 1 **Filter 1** **Row 1** \* **Fmap 1** **Row 1** = **Psum 1** **Row 1**

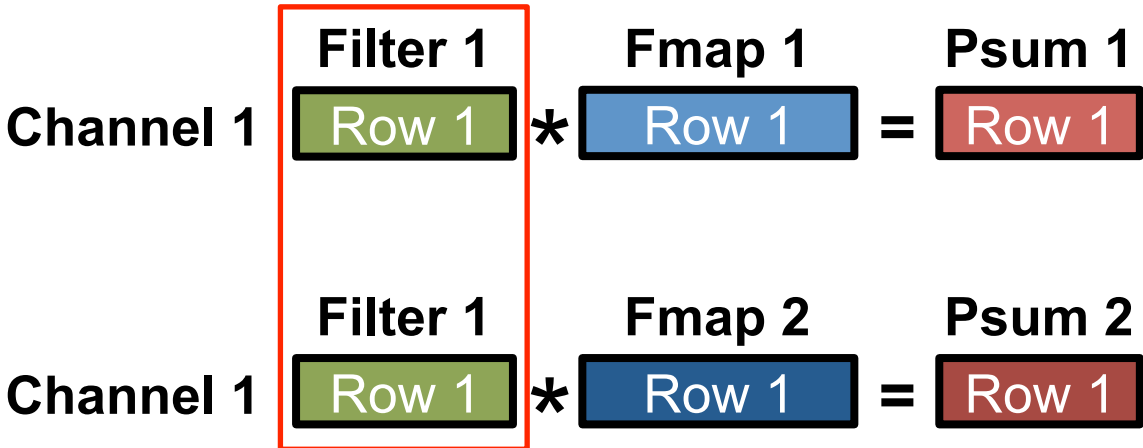
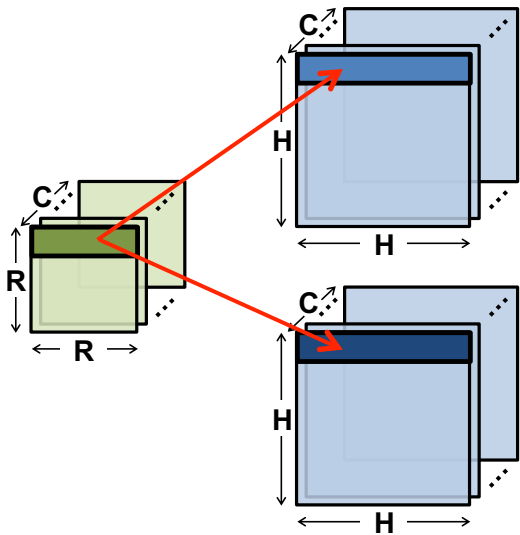
Channel 1 **Filter 1** **Row 1** \* **Fmap 2** **Row 1** = **Psum 2** **Row 1**

# Filter Reuse in PE

## 1 Multiple Fmaps

## 2 Multiple Filters

## 3 Multiple Channels



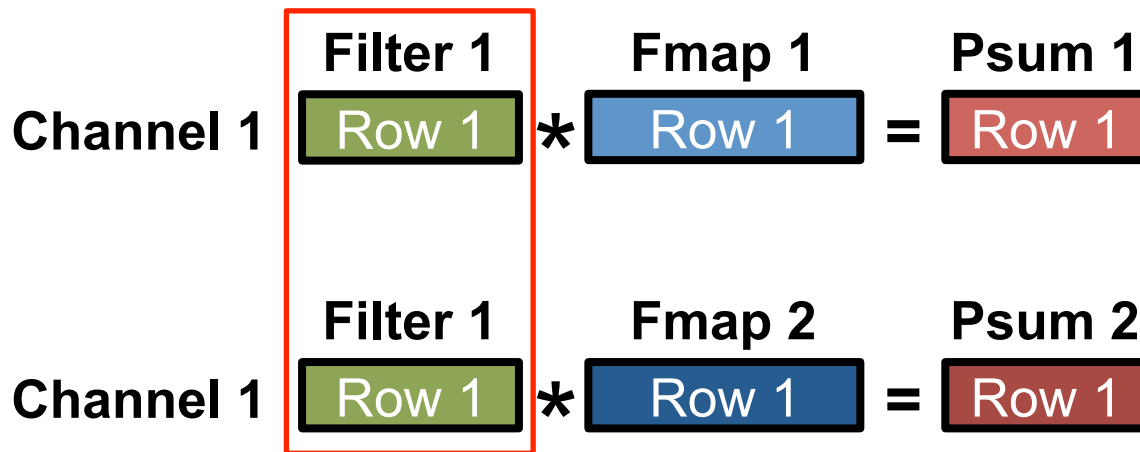
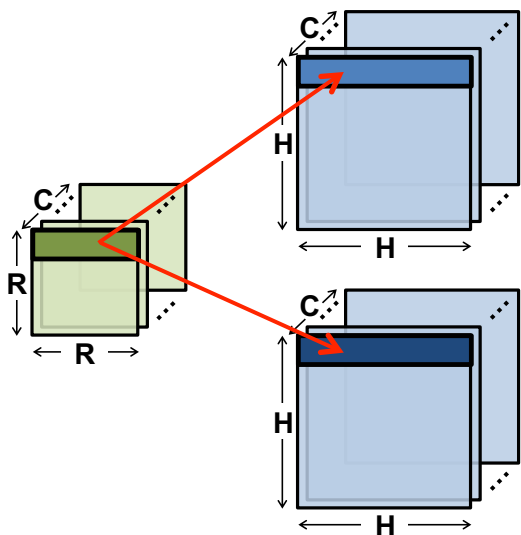
share the same filter row

# Filter Reuse in PE

## 1 Multiple Fmaps

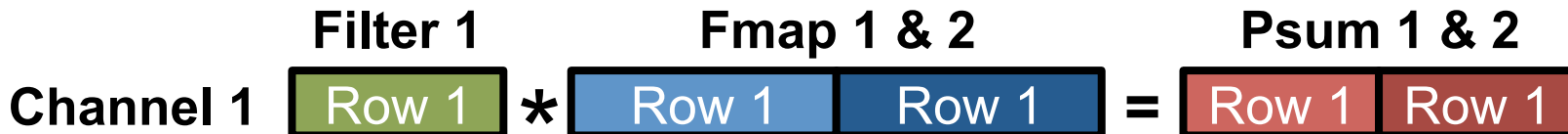
## 2 Multiple Filters

## 3 Multiple Channels



share the same filter row

Processing in PE: concatenate fmap rows

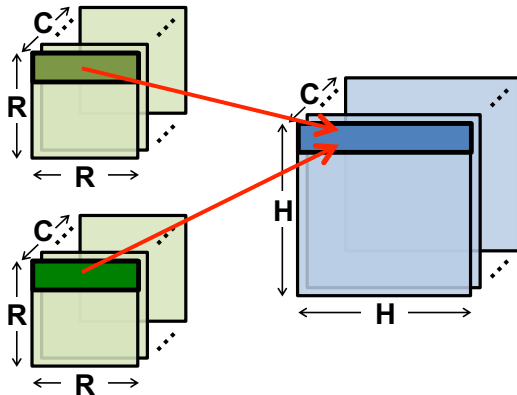


# Fmap Reuse in PE

1 Multiple Fmaps

2 Multiple Filters

3 Multiple Channels



$$\text{Channel 1} \quad \text{Filter 1} \quad \text{Fmap 1} \quad \text{Psum 1}$$
$$\text{Row 1} * \text{Row 1} = \text{Row 1}$$

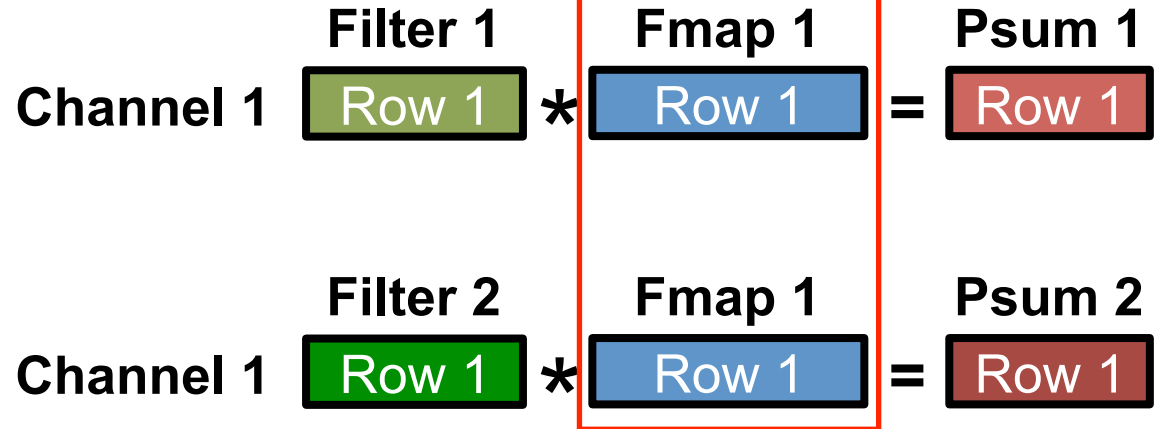
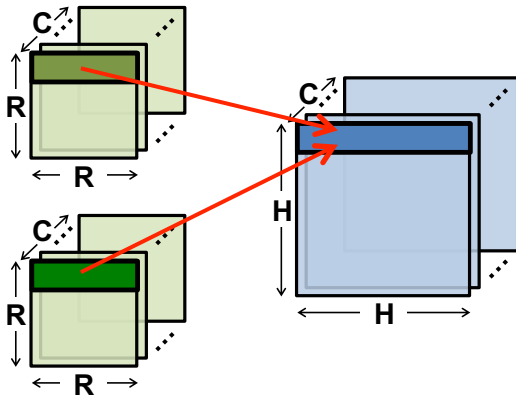
$$\text{Channel 1} \quad \text{Filter 2} \quad \text{Fmap 1} \quad \text{Psum 2}$$
$$\text{Row 1} * \text{Row 1} = \text{Row 1}$$

# Fmap Reuse in PE

1 Multiple Fmaps

2 Multiple Filters

3 Multiple Channels



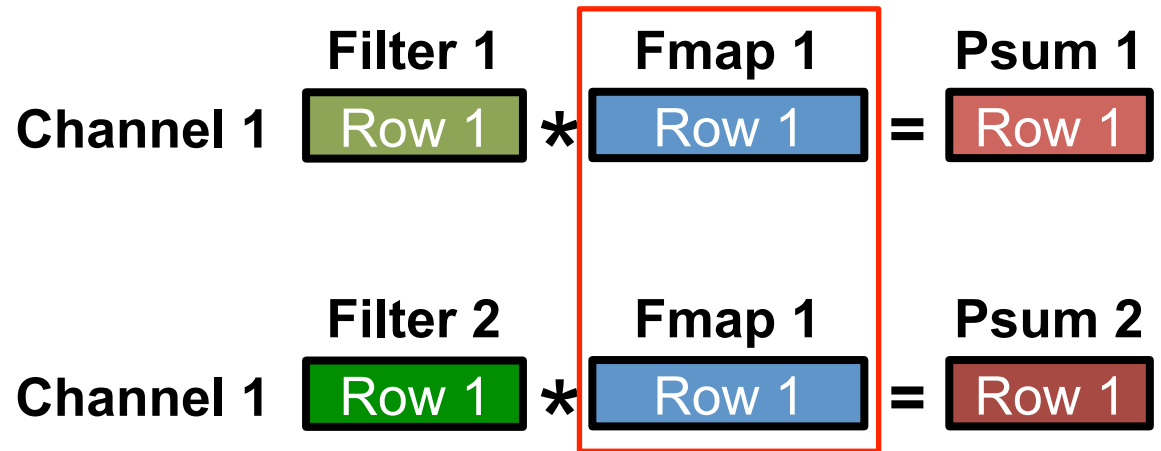
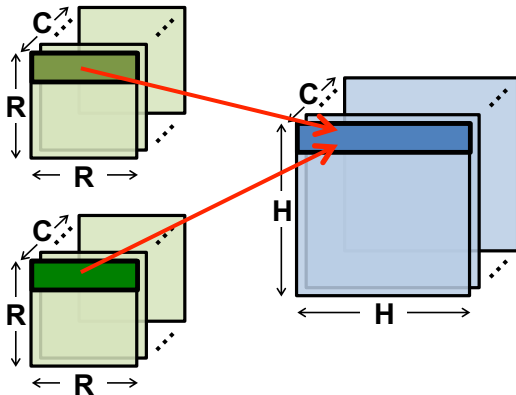
share the same fmap row

# Fmap Reuse in PE

1 Multiple Fmaps

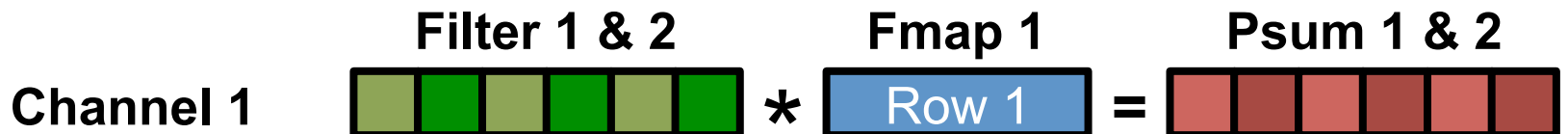
2 Multiple Filters

3 Multiple Channels



share the same fmap row

Processing in PE: interleave filter rows

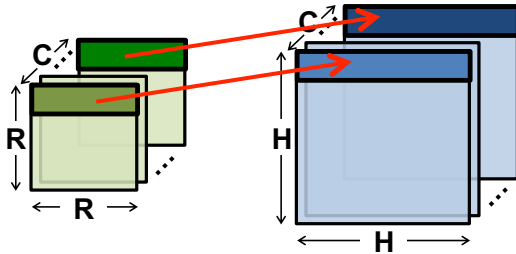


# Channel Accumulation in PE

1 Multiple Fmaps

2 Multiple Filters

3 Multiple Channels



Channel 1

Filter 1		Fmap 1	=	Psum 1
Row 1	*	Row 1	=	Row 1

Channel 2

Filter 1		Fmap 1	=	Psum 1
Row 1	*	Row 1	=	Row 1

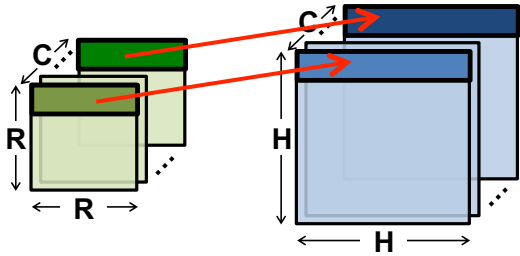


# Channel Accumulation in PE

1 Multiple Fmaps

2 Multiple Filters

3 Multiple Channels



$$\begin{array}{l} \text{Channel 1} \\ \text{Channel 2} \end{array} \begin{array}{l} \text{Filter 1} \\ \text{Filter 1} \end{array} \begin{array}{l} \text{Row 1} \\ \text{Row 1} \end{array} * \begin{array}{l} \text{Fmap 1} \\ \text{Fmap 1} \end{array} \begin{array}{l} \text{Row 1} \\ \text{Row 1} \end{array} = \begin{array}{l} \text{Psum 1} \\ \text{Psum 1} \end{array} \begin{array}{l} \text{Row 1} \\ \text{Row 1} \end{array}$$

accumulate psums

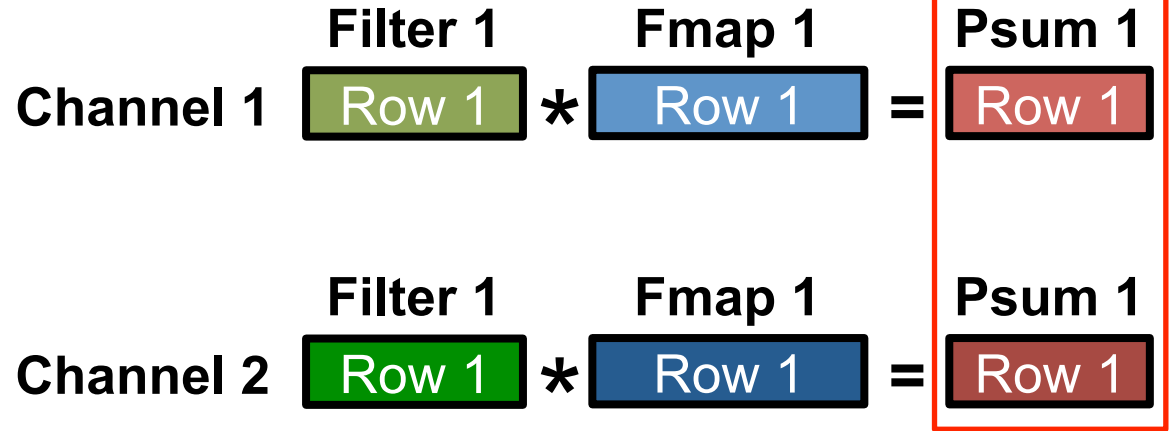
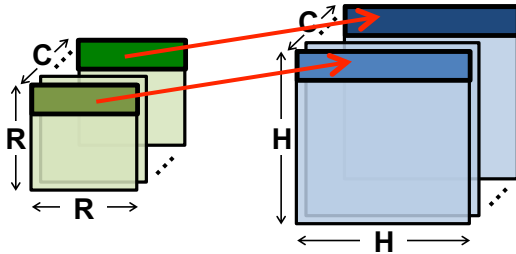
$$\begin{array}{l} \text{Row 1} \\ \text{Row 1} \end{array} + \begin{array}{l} \text{Row 1} \\ \text{Row 1} \end{array} = \begin{array}{l} \text{Row 1} \\ \text{Row 1} \end{array}$$

# Channel Accumulation in PE

1 Multiple Fmaps

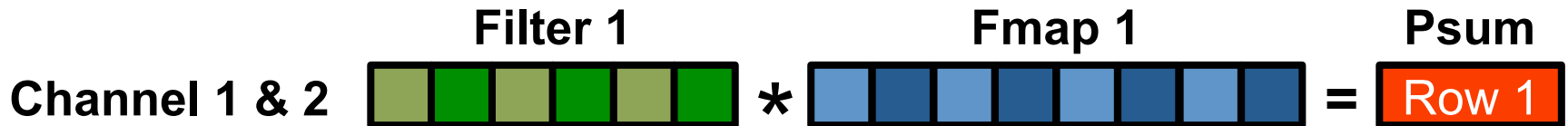
2 Multiple Filters

3 Multiple Channels

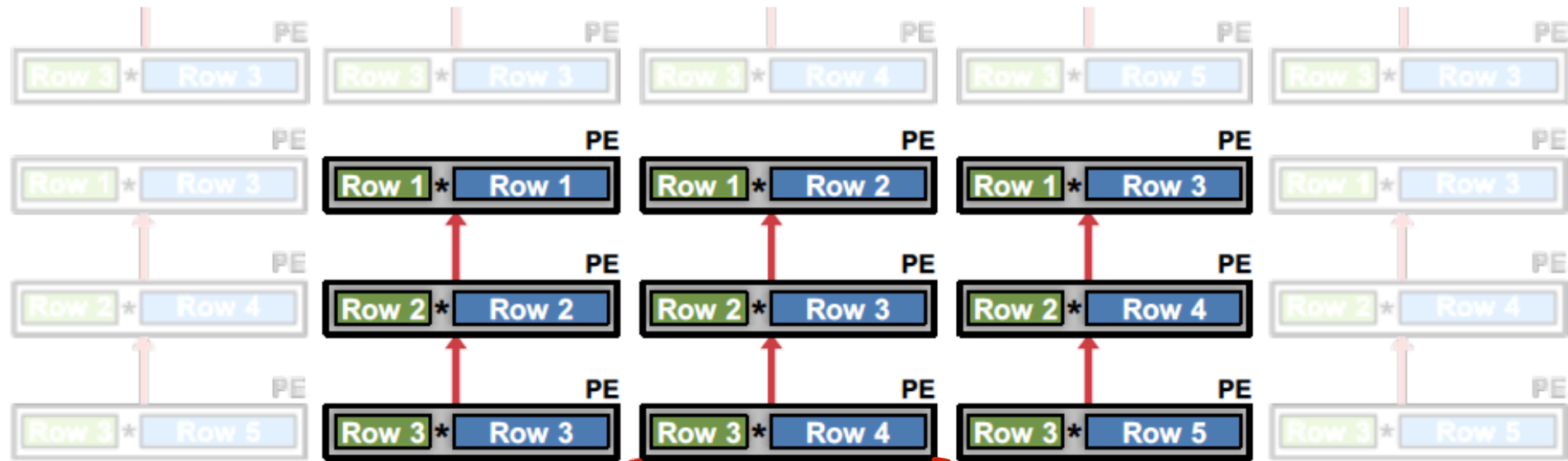


accumulate psums

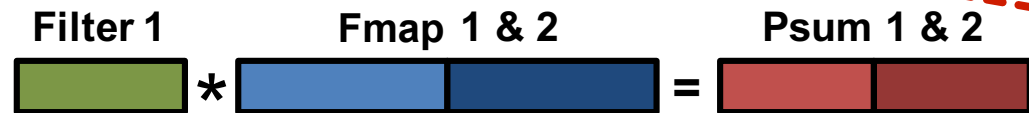
Processing in PE: interleave channels



# DNN Processing – The Full Picture



Multiple **fmaps**:



Multiple **filters**:



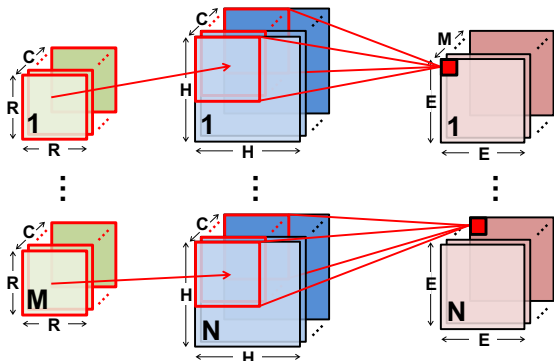
Multiple **channels**:



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

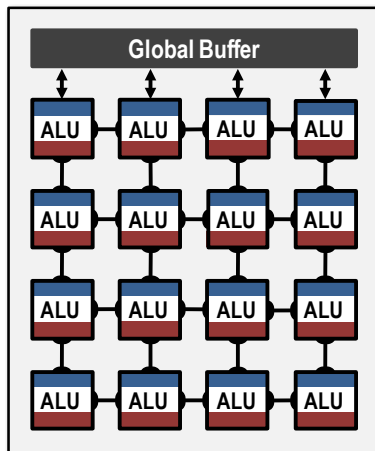
# Optimal Mapping in Row Stationary

## CNN Configurations

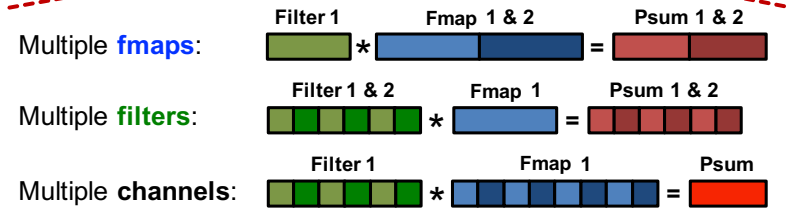
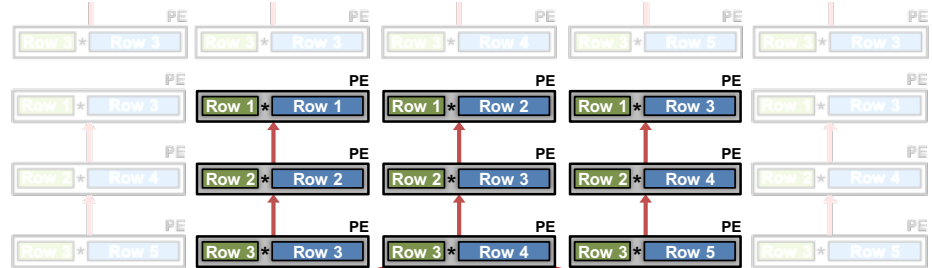


Optimization Compiler

## Hardware Resources



## Row Stationary Mapping



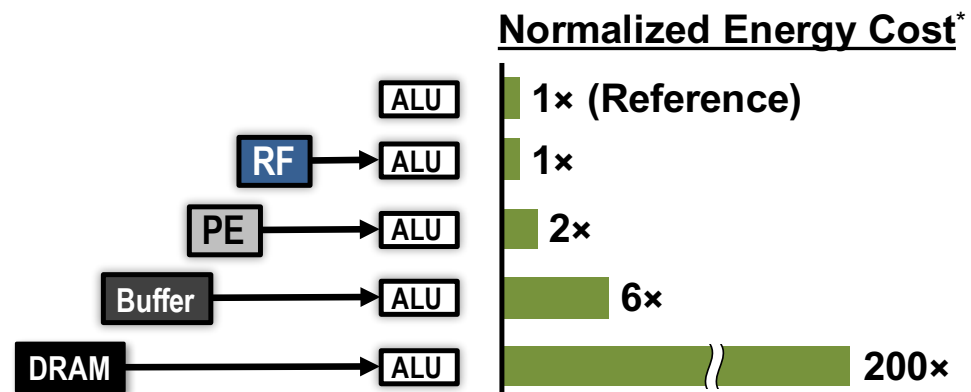
# Dataflow Simulation Results

# Evaluate Reuse in Different Dataflows

- **Weight Stationary**
  - Minimize movement of filter weights
- **Output Stationary**
  - Minimize movement of partial sums
- **No Local Reuse**
  - No PE local storage. Maximize global buffer size.
- **Row Stationary**

## Evaluation Setup

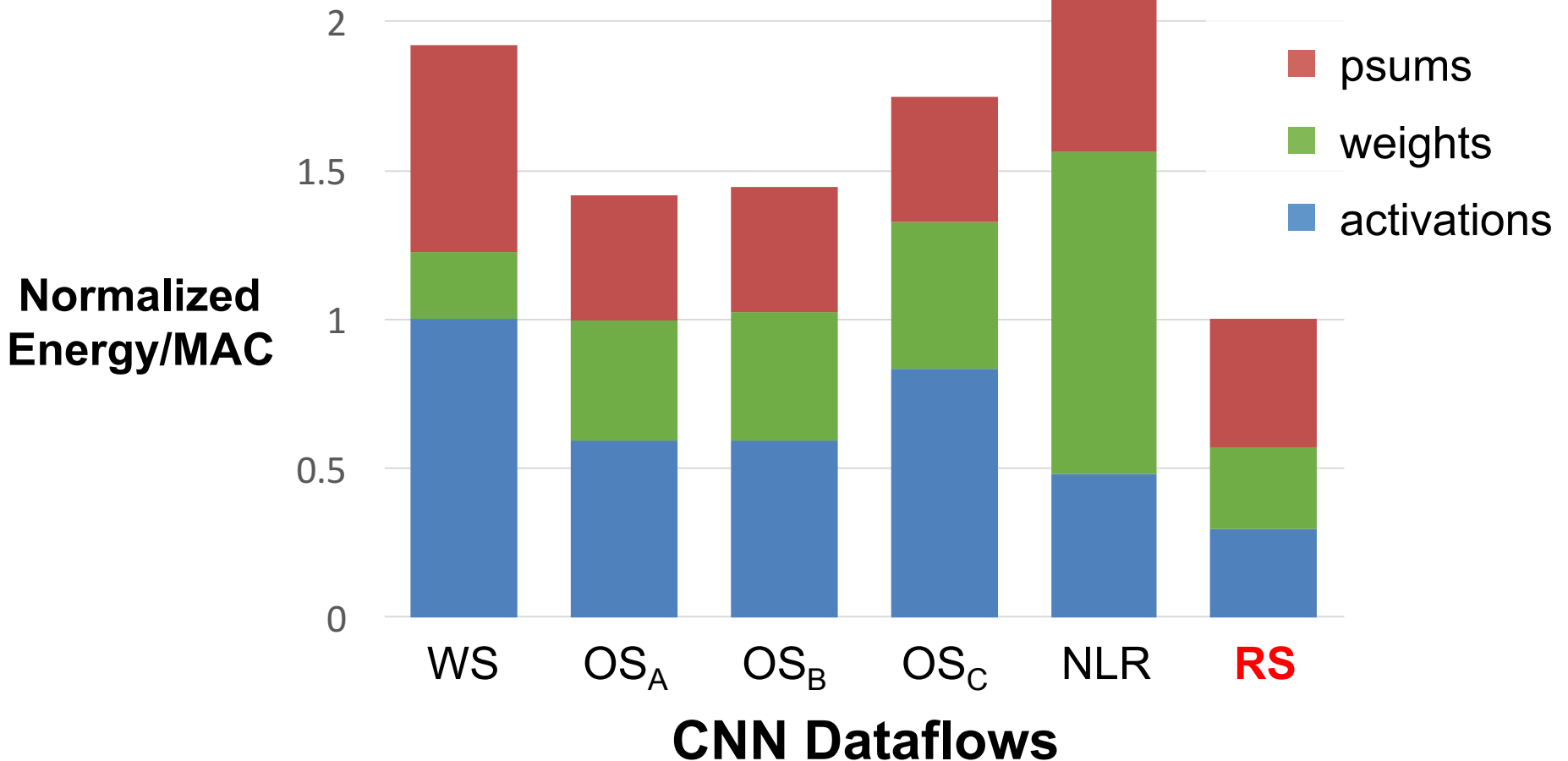
- same total area
- 256 PEs
- AlexNet
- batch size = 16



# Variants of Output Stationary

	$OS_A$	$OS_B$	$OS_C$
Parallel Output Region			
# Output Channels	Single	Multiple	Multiple
# Output Activations	Multiple	Multiple	Single
Notes	Targeting <b>CONV</b> layers		Targeting <b>FC</b> layers

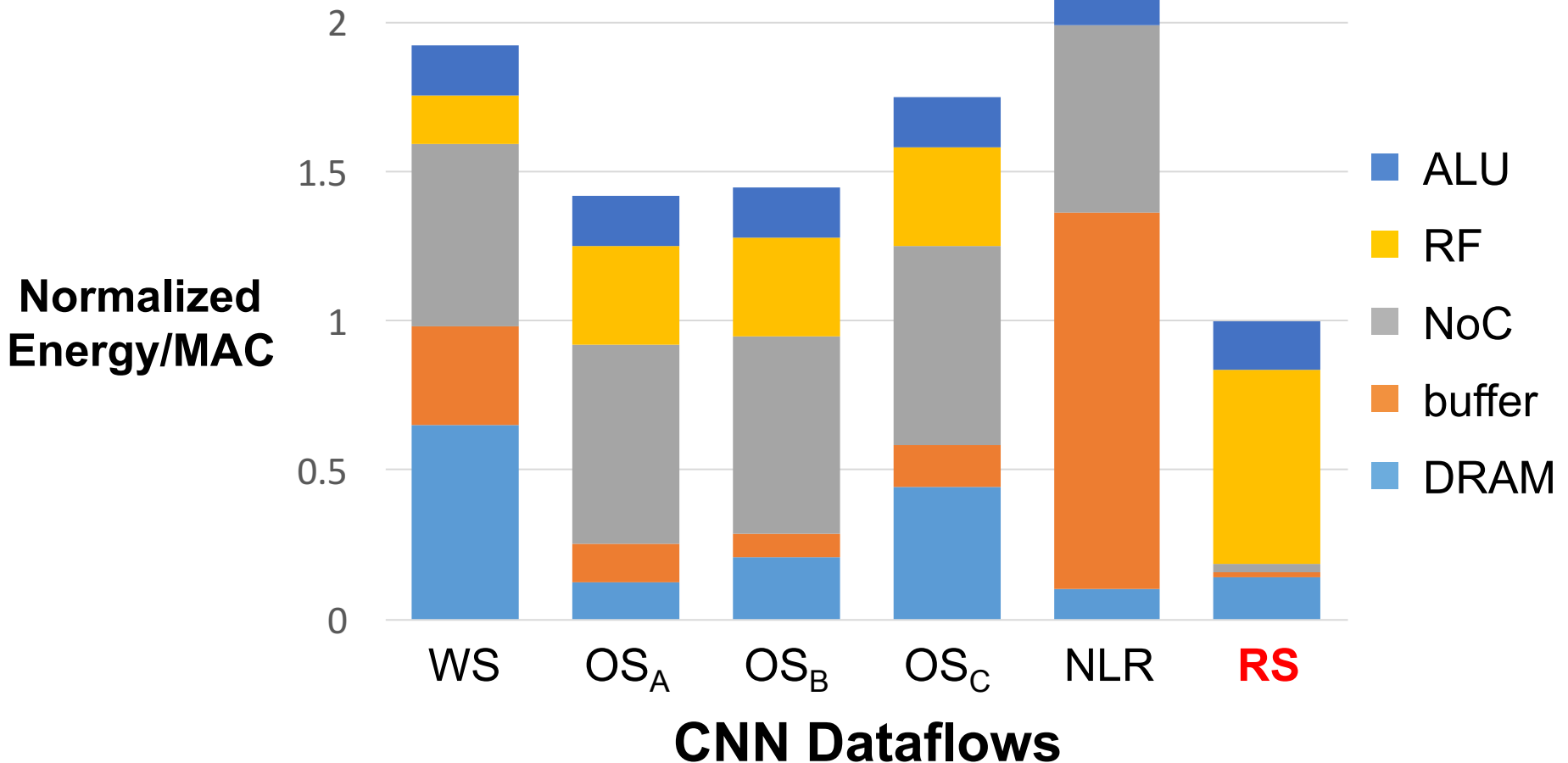
# Dataflow Comparison: CONV Layers



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

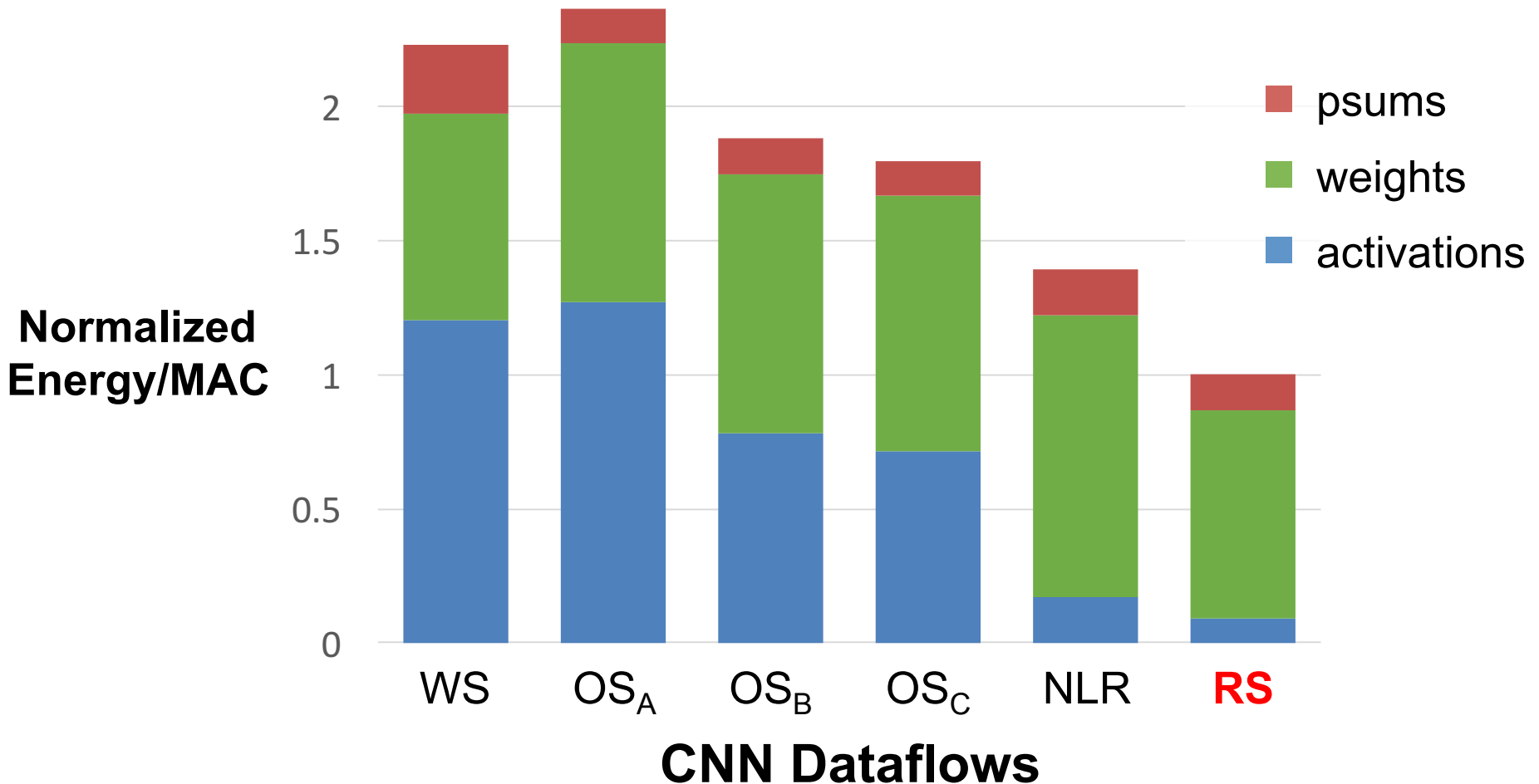


# Dataflow Comparison: CONV Layers



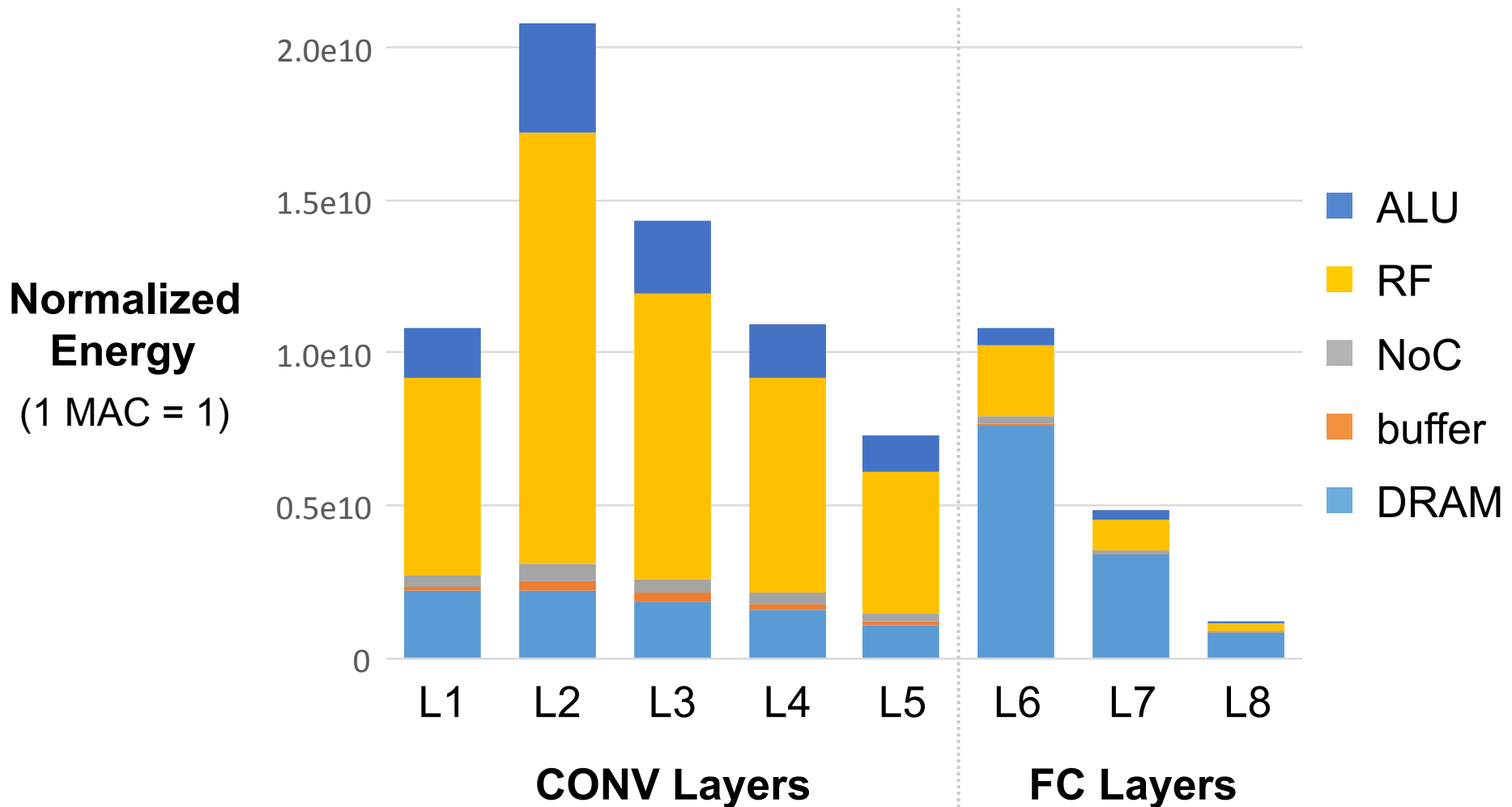
RS uses **1.4× – 2.5× lower** energy than other dataflows

# Dataflow Comparison: FC Layers

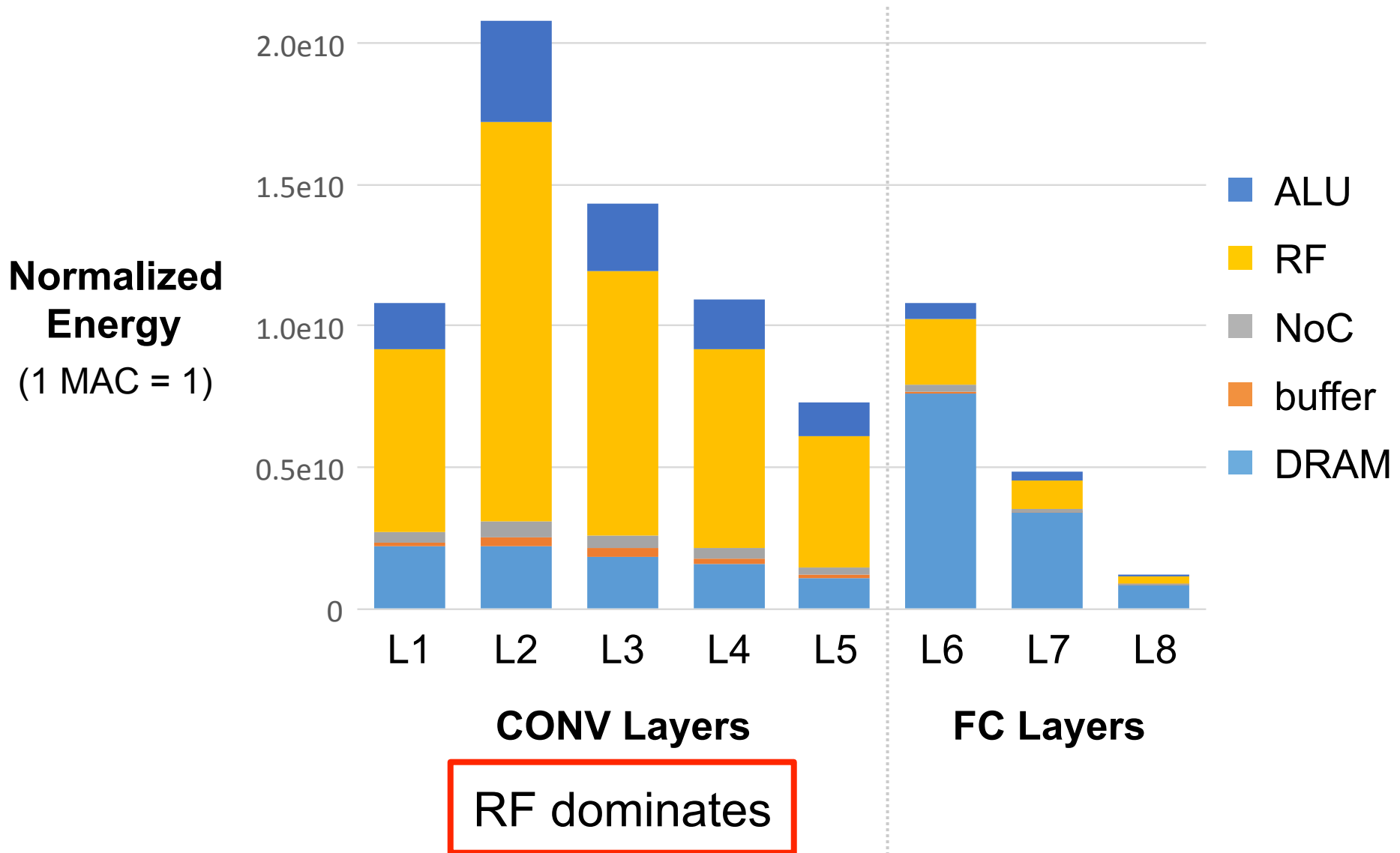


RS uses at least **1.3× lower** energy than other dataflows

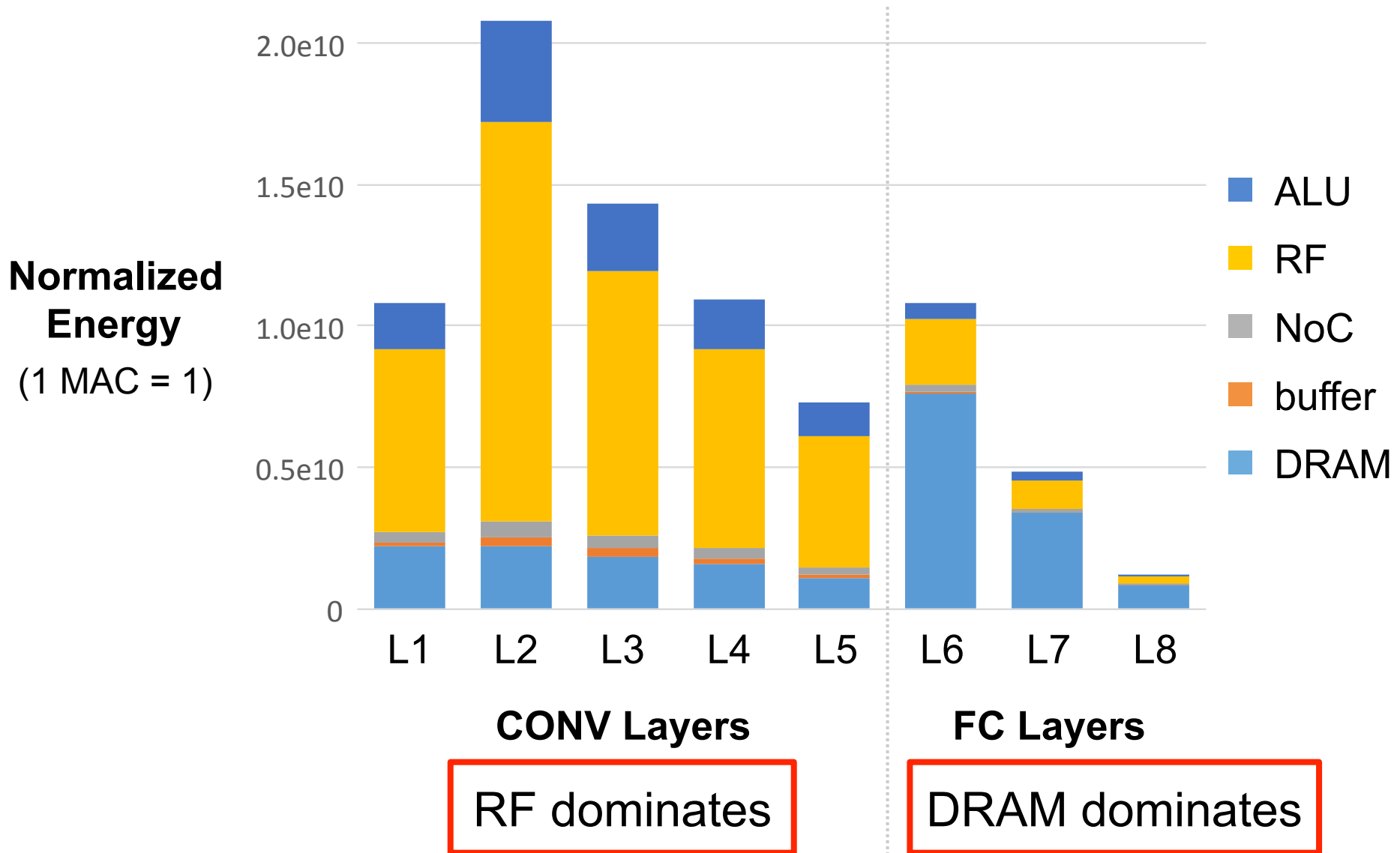
# Row Stationary: Layer Breakdown



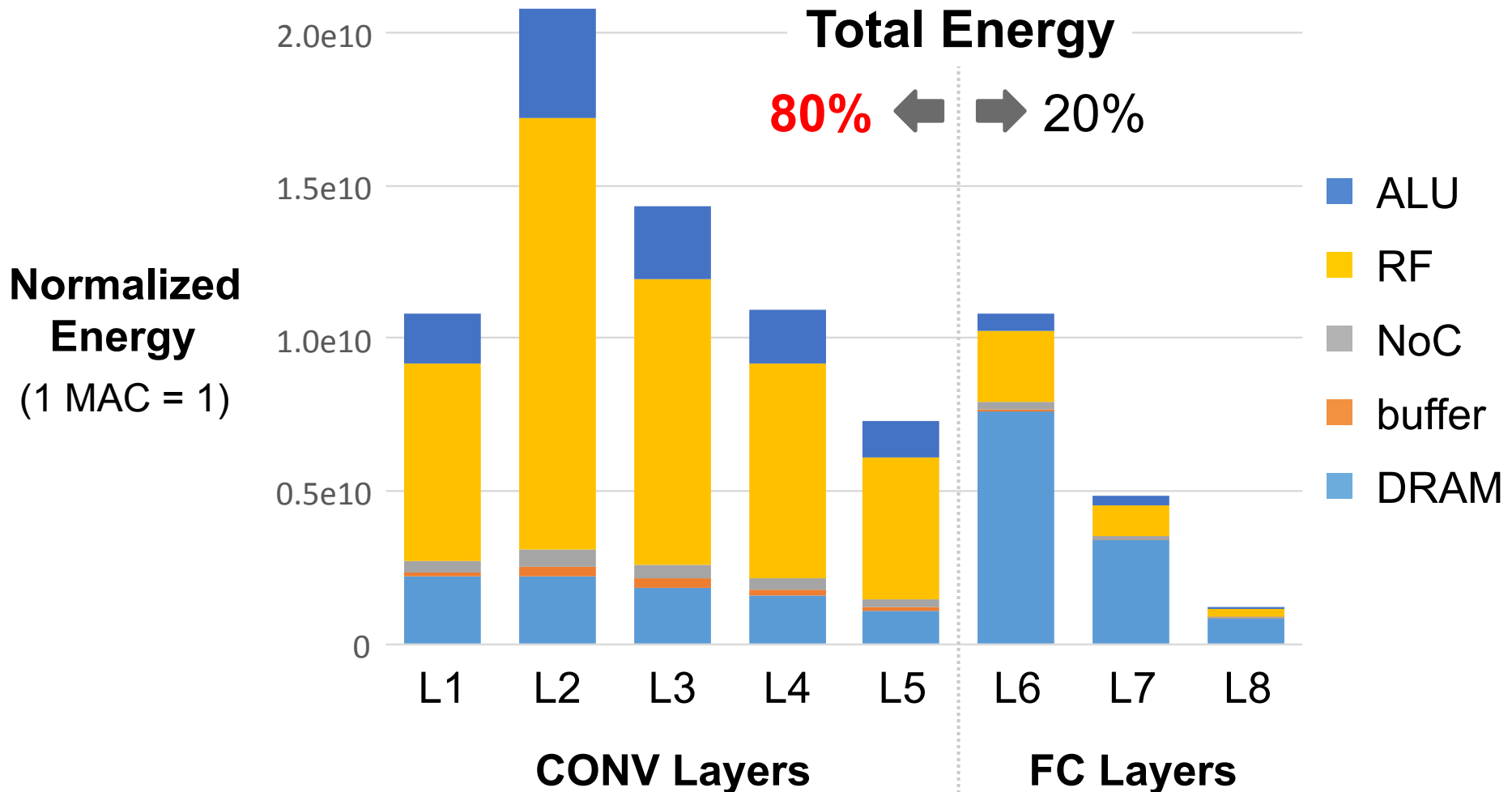
# Row Stationary: Layer Breakdown



# Row Stationary: Layer Breakdown



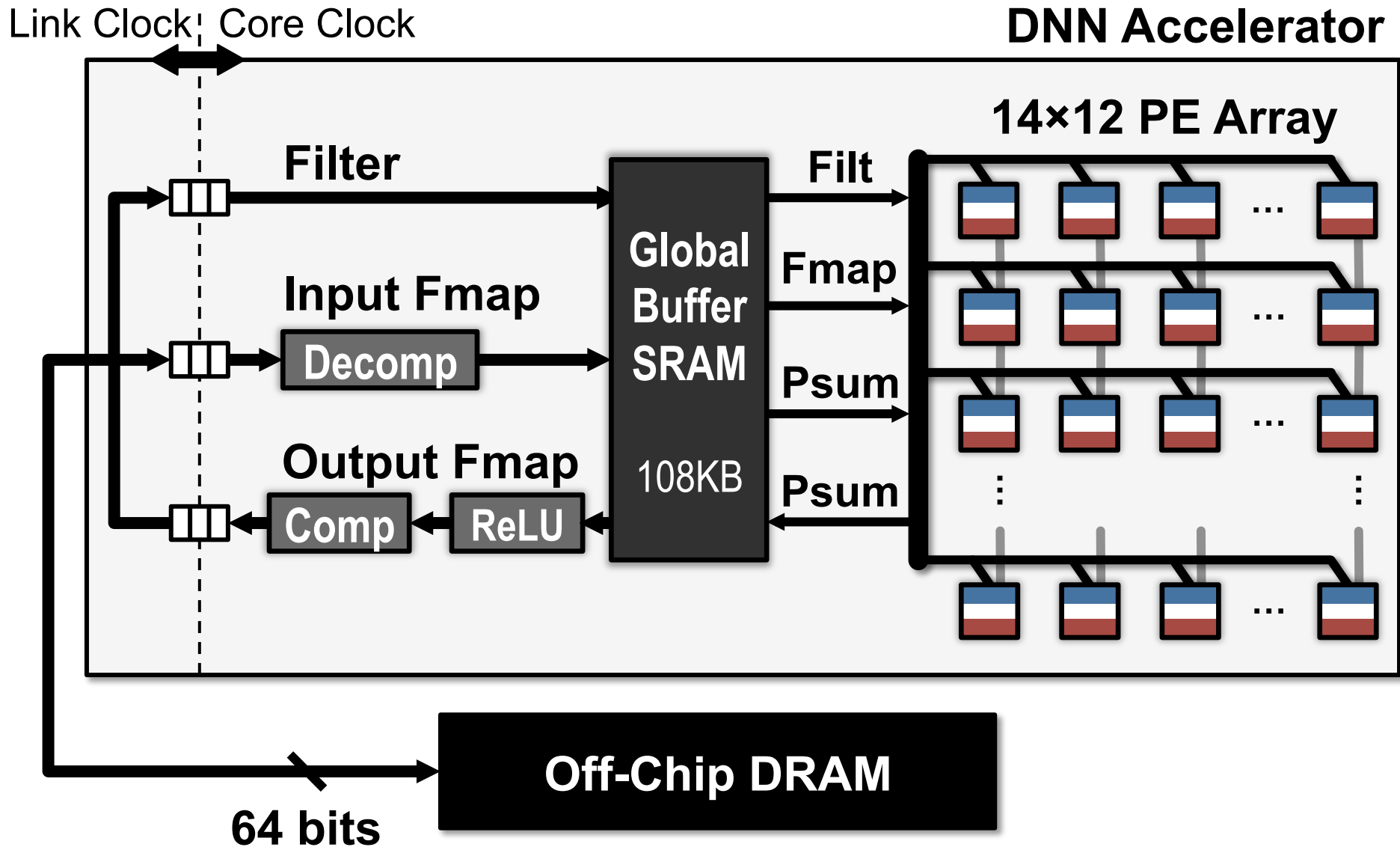
# Row Stationary: Layer Breakdown



CONV layers dominate energy consumption!

# Hardware Architecture for RS Dataflow

# Eyeriss DNN Accelerator

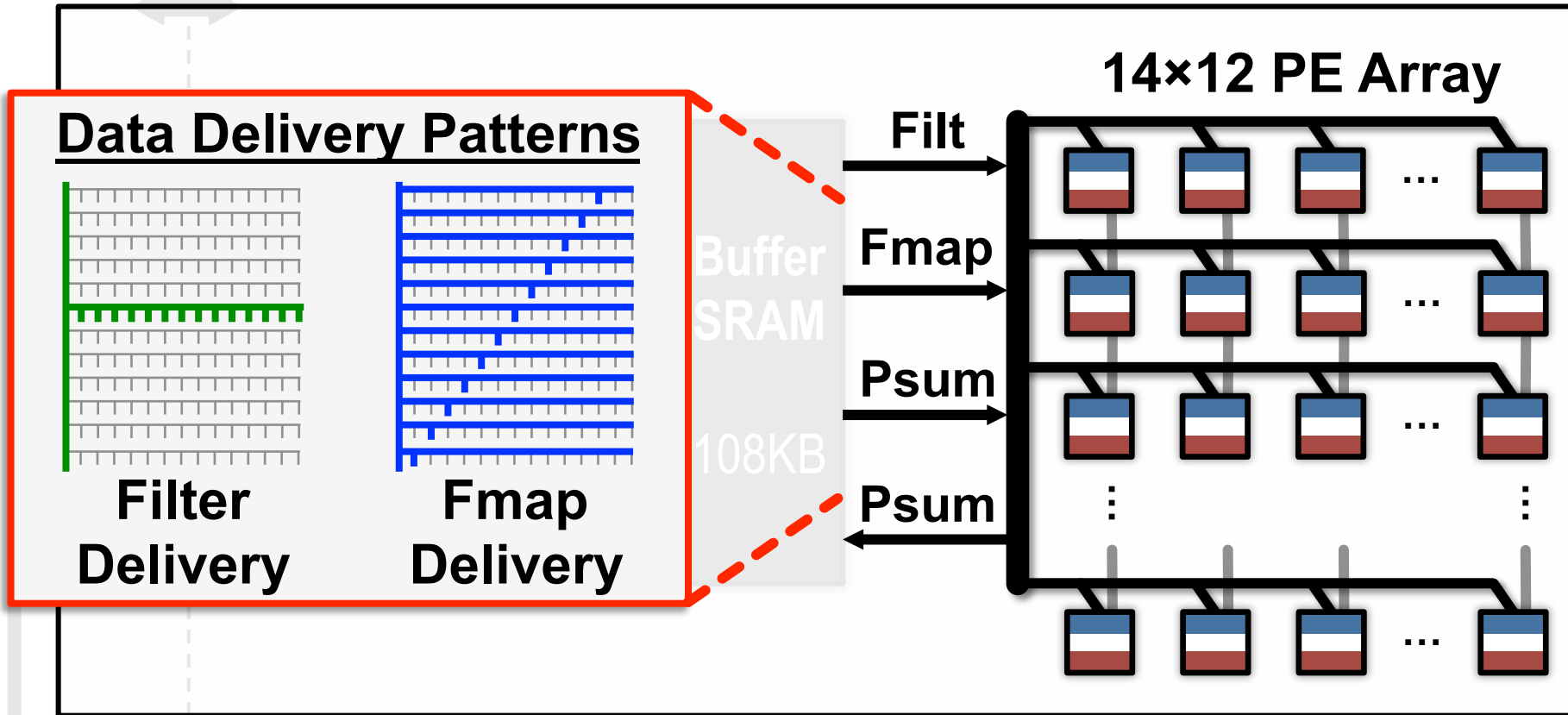




# Data Delivery with On-Chip Network

Link Clock | Core Clock

DCNN Accelerator

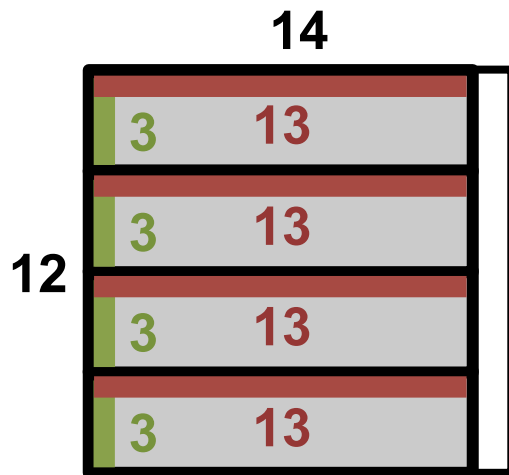
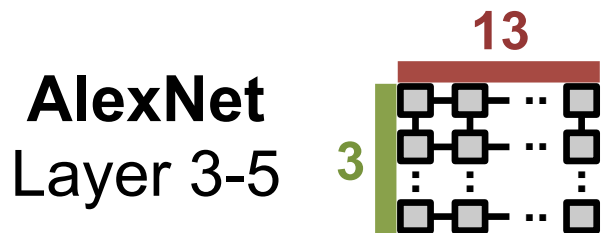


How to accommodate different shapes with fixed PE array?

64 bits

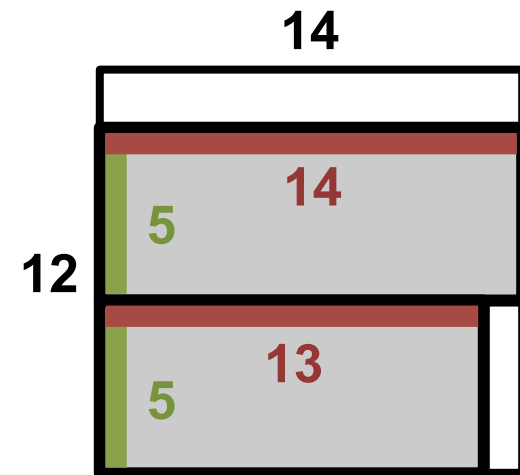
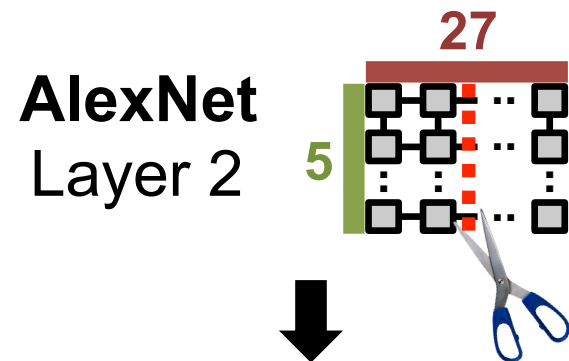
# Logical to Physical Mappings

## Replication



Physical PE Array

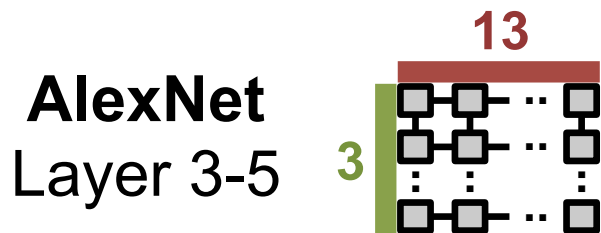
## Folding



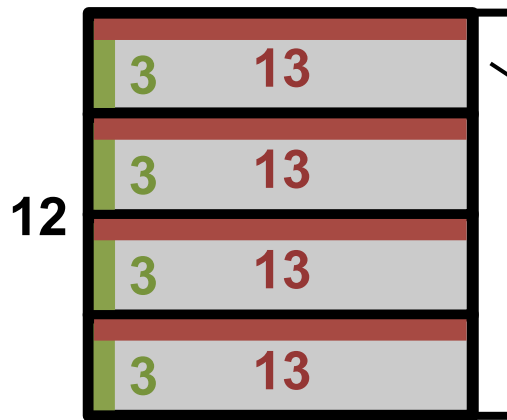
Physical PE Array

# Logical to Physical Mappings

## Replication

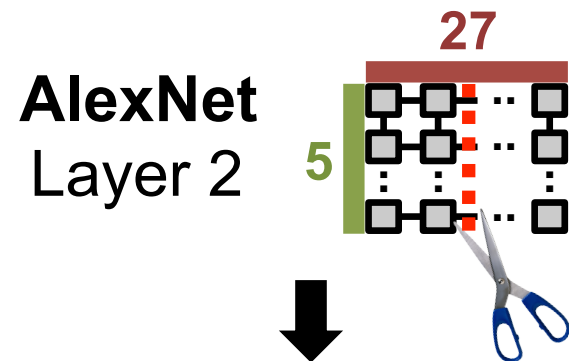


14

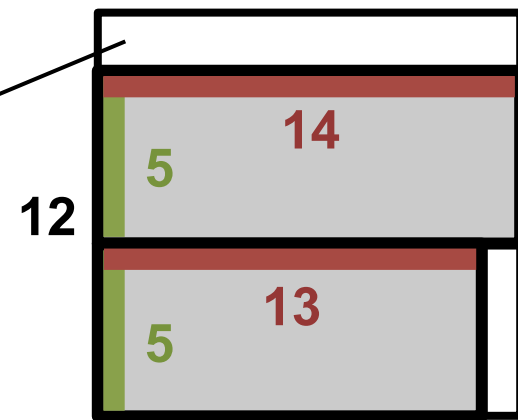


Physical PE Array

## Folding



14



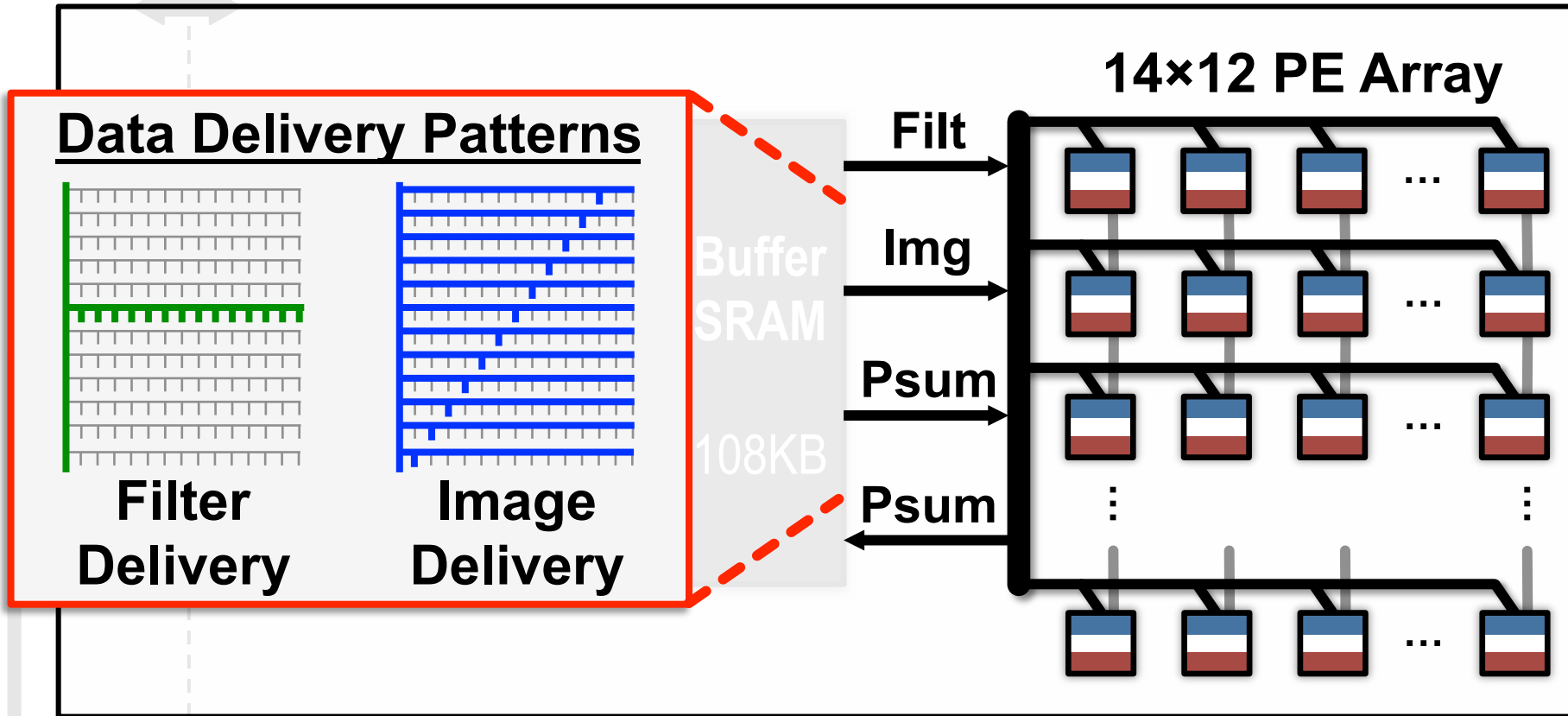
Physical PE Array

Unused PEs  
are  
Clock Gated

# Data Delivery with On-Chip Network

Link Clock | Core Clock

DCNN Accelerator

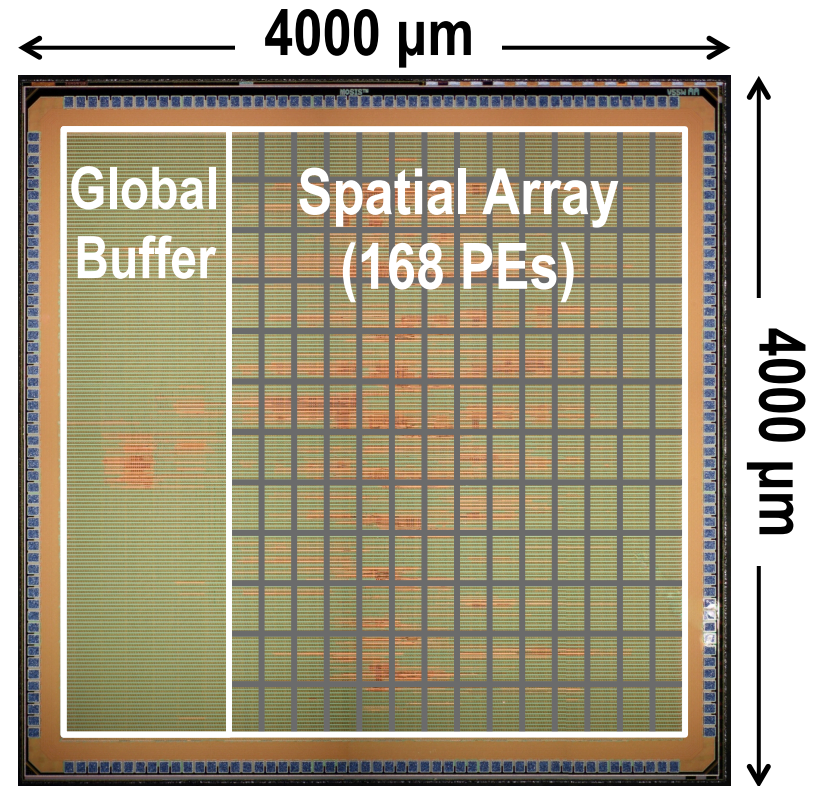


Compared to Broadcast, **Multicast** saves **>80%** of NoC energy

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)

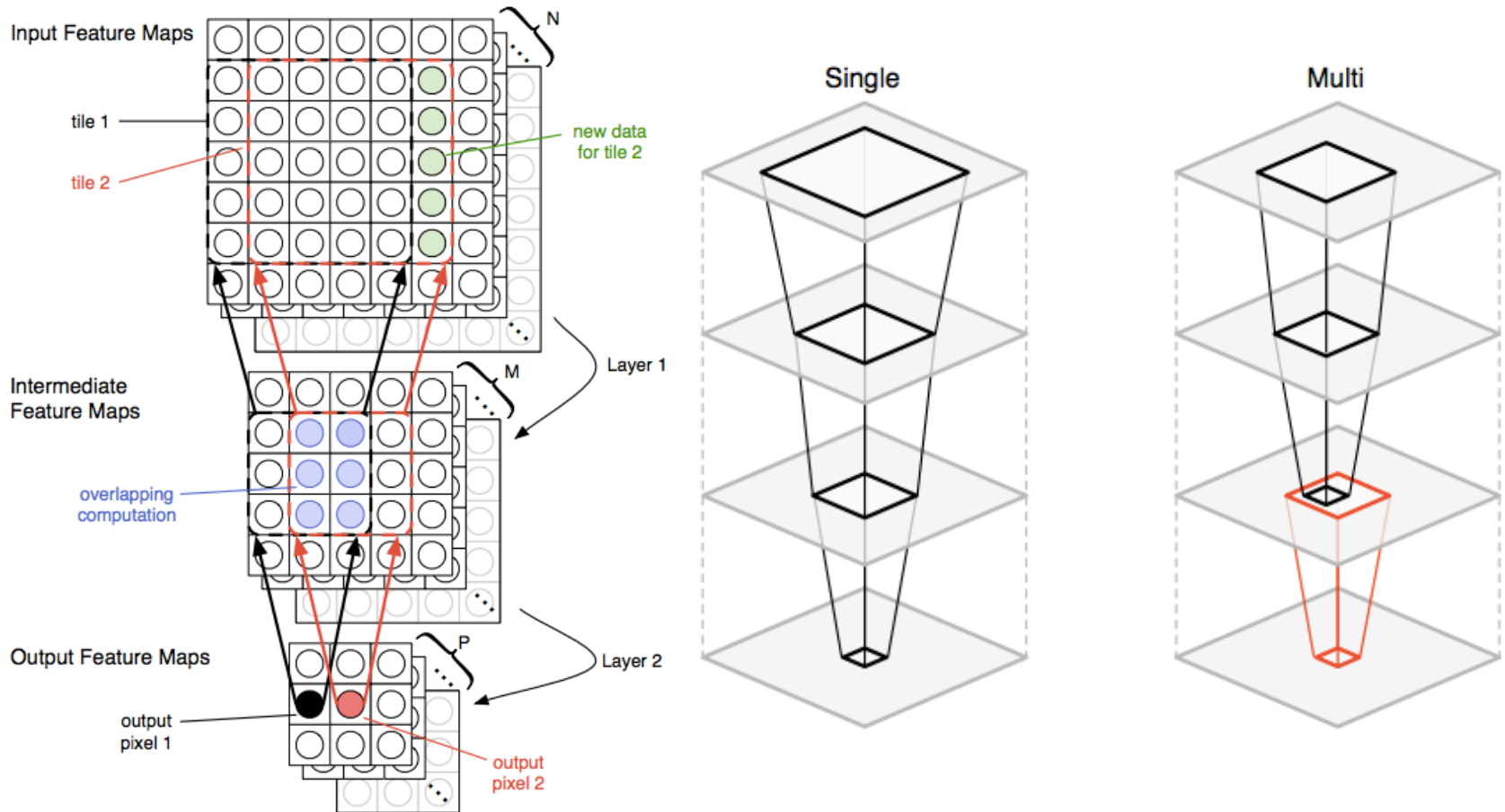
# Summary of DNN Dataflows

---

- **Weight Stationary**
  - Minimize movement of filter weights
  - Popular with processing-in-memory architectures
- **Output Stationary**
  - Minimize movement of partial sums
  - Different variants optimized for CONV or FC layers
- **No Local Reuse**
  - No PE local storage → maximize global buffer size
- **Row Stationary**
  - Adapt to the NN shape and hardware constraints
  - Optimized for overall **system energy efficiency**

# Fused Layer

- Dataflow across multiple layers



# Metrics for DNN Hardware

---

- **Measure energy and DRAM access relative to number of non-zero MACs and bit-width of MACs**
  - Account for impact of sparsity in weights and activations
  - Normalize DRAM access based on operand size
- **Energy Efficiency of Design**
  - $\text{pJ}/(\text{non-zero weight \& activation})$
- **External Memory Bandwidth**
  - $\text{DRAM operand access}/(\text{non-zero weight \& activation})$
- **Area Efficiency**
  - Total chip  $\text{mm}^2/\text{multi}$  (also include process technology)
  - Accounts for on-chip memory



# Website to Summarize Results

- <http://eyeriss.mit.edu/benchmarking.html>
- Send results or feedback to: [eyeriss@mit.edu](mailto:eyeriss@mit.edu)

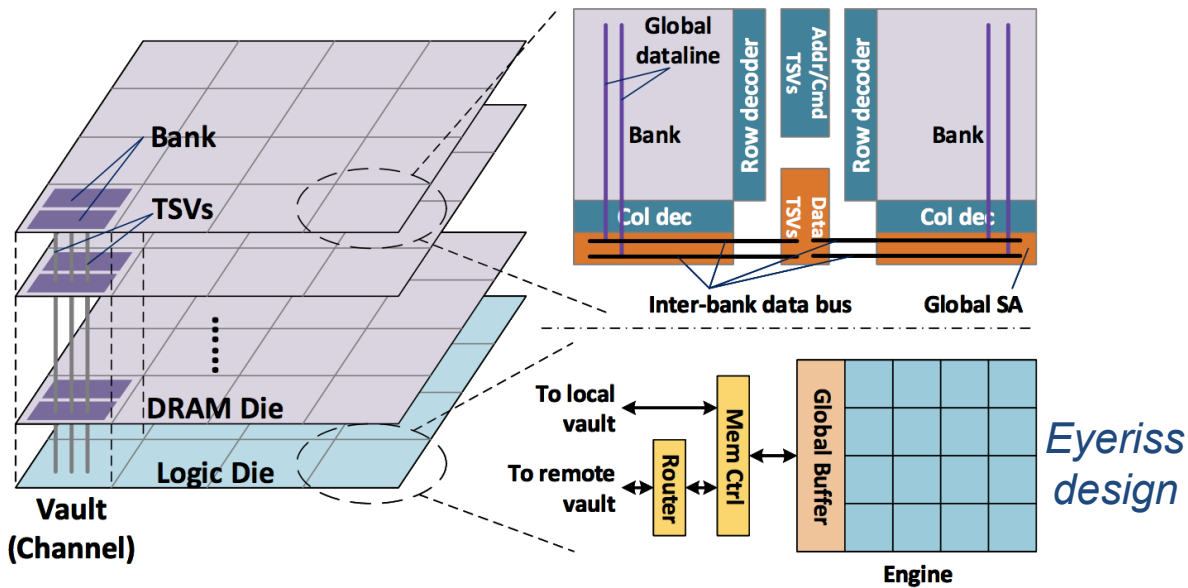
ASIC Specs	Input
Process Technology	65nm LP TSMC (1.0V)
Core area (mm <sup>2</sup> ) / multiplier	0.073
On-Chip memory (kB) / multiplier	1.14
Measured or Simulated	Measured
If Simulated, Syn or PnR? Which corner?	n/a

Metric	Units	Input
Name of CNN	Text	AlexNet
# of Images Tested	#	100
Bits per operand	#	16
Batch Size	#	4
# of Non Zero MACs	#	409M
Runtime	ms	115.3
Power	mW	278
<b>Energy/non-zero MACs</b>	<b>pJ/MAC</b>	<b>21.7</b>
<b>DRAM access/non-zero MACs</b>	<b>operands /MAC</b>	<b>0.005</b>

# Advanced Memory Technologies

Many new memories and devices explored to reduce data movement

## Stacked DRAM



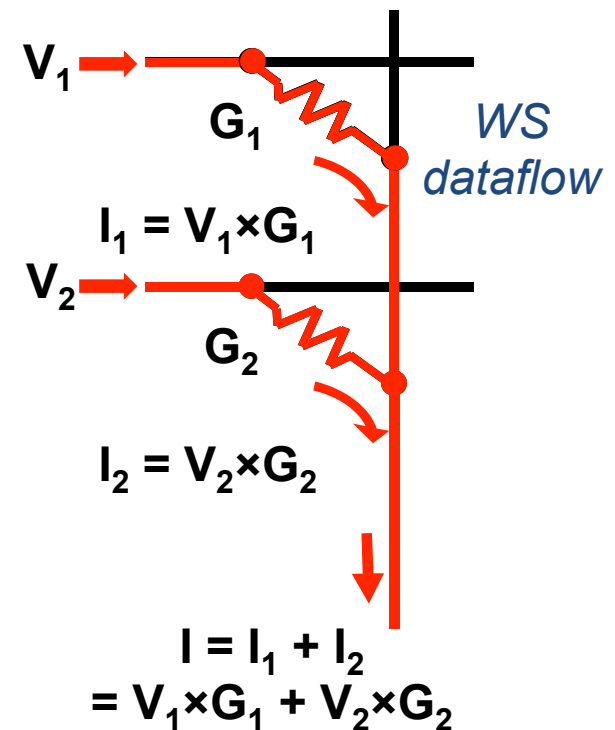
[Gao et al., Tetris, ASPLOS 2017]

[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]