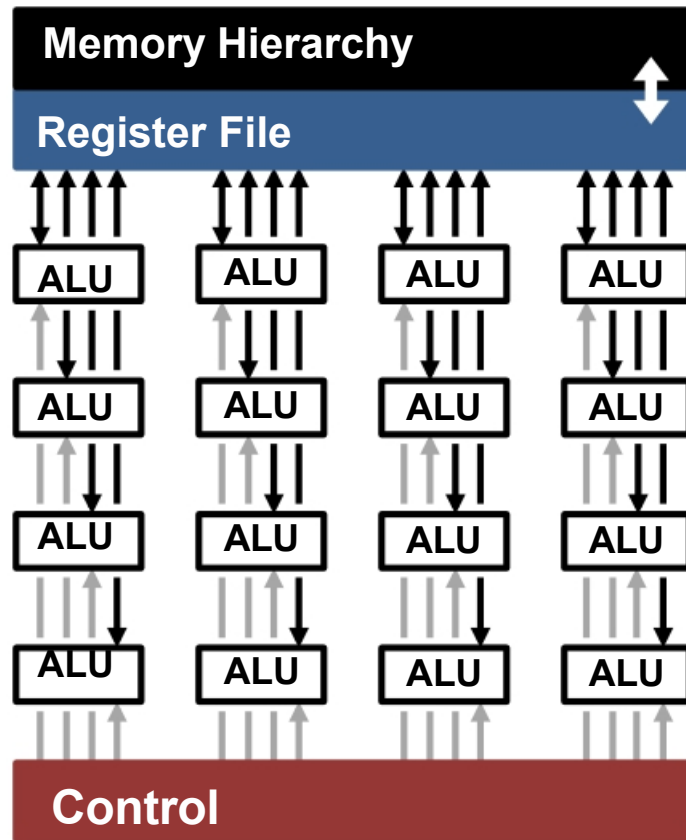


---

# DNN Accelerator Architectures

# Highly-Parallel Compute Paradigms

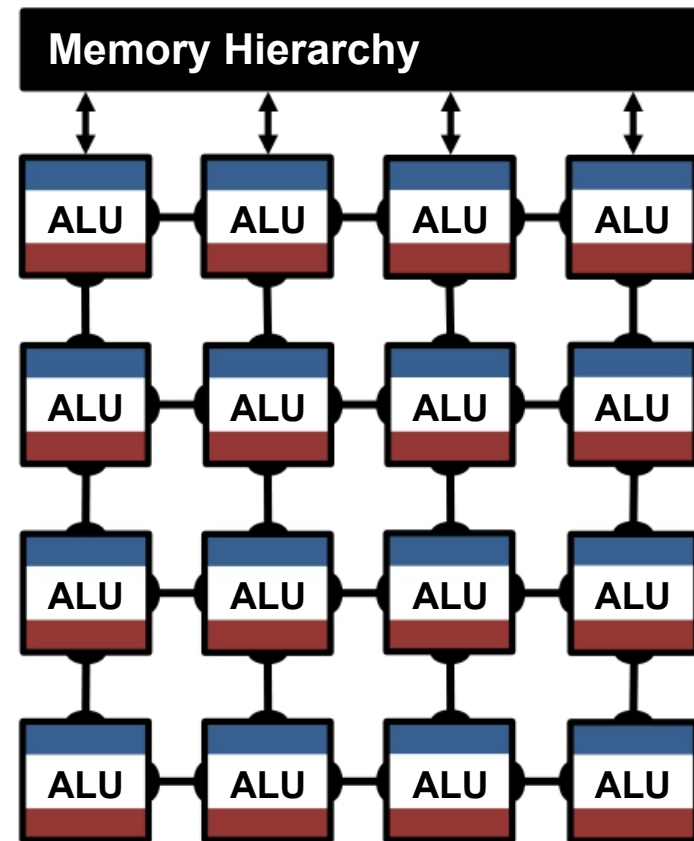
## Temporal Architecture (SIMD/SIMT)



Mostly in **CPU & GPU**

Centralized control for many ALUs (cannot communicate directly with each other directly)

## Spatial Architecture (Dataflow Processing)

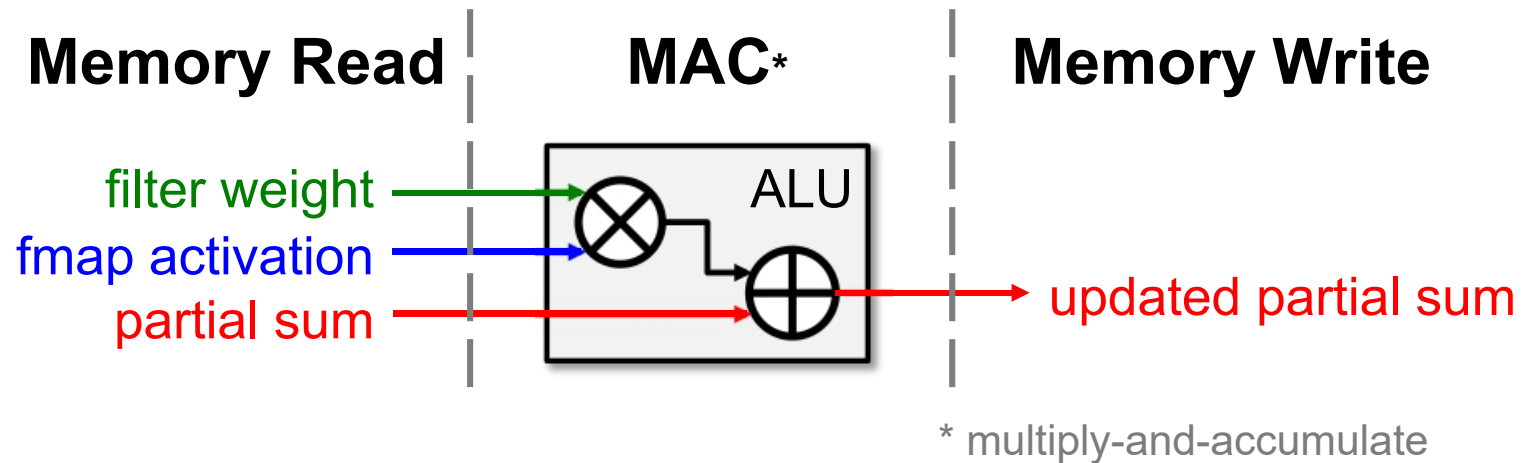


Data are passed from one ALU to another ALU

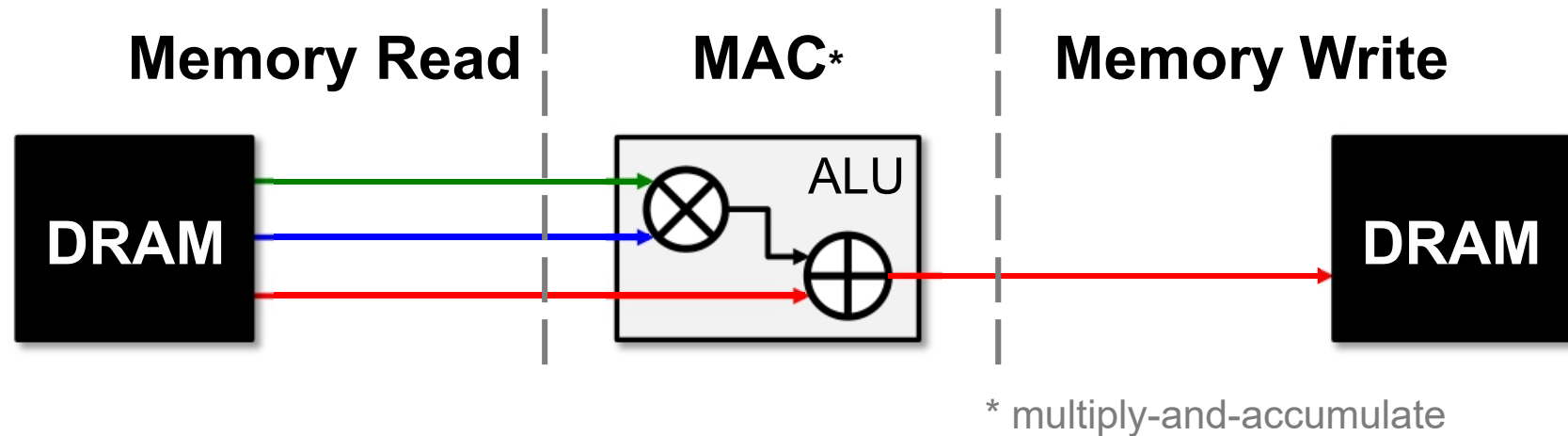
Mostly in ASIC or FPGA

# Memory Access is the Bottleneck

---



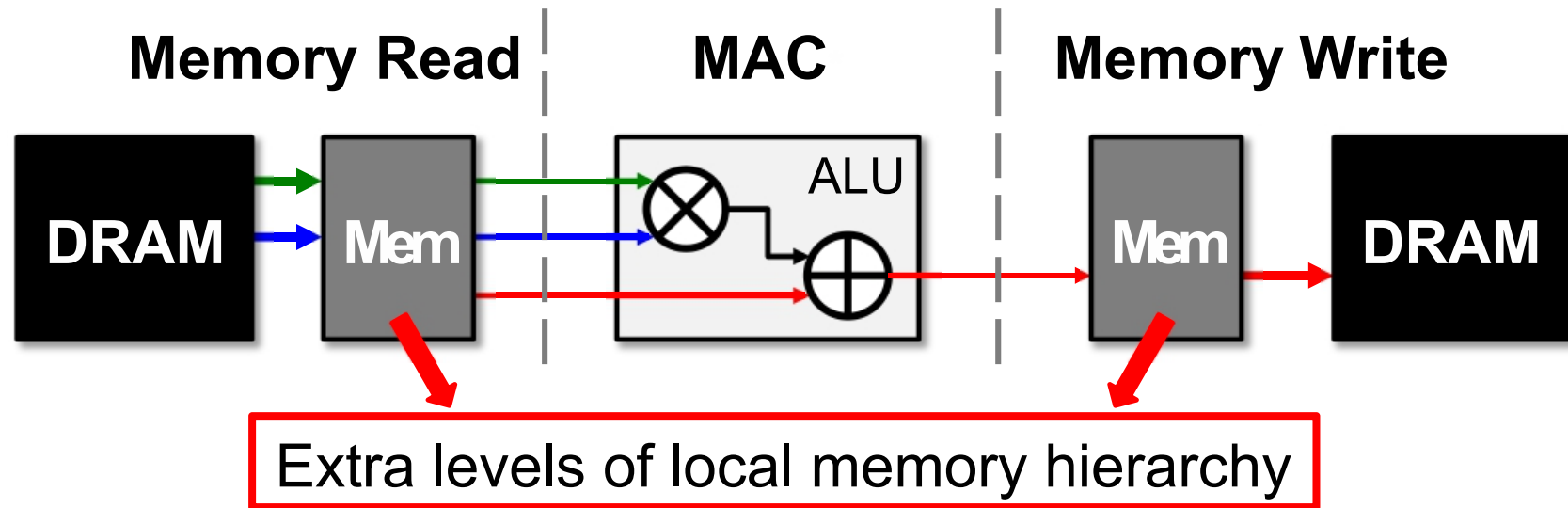
# Memory Access is the Bottleneck



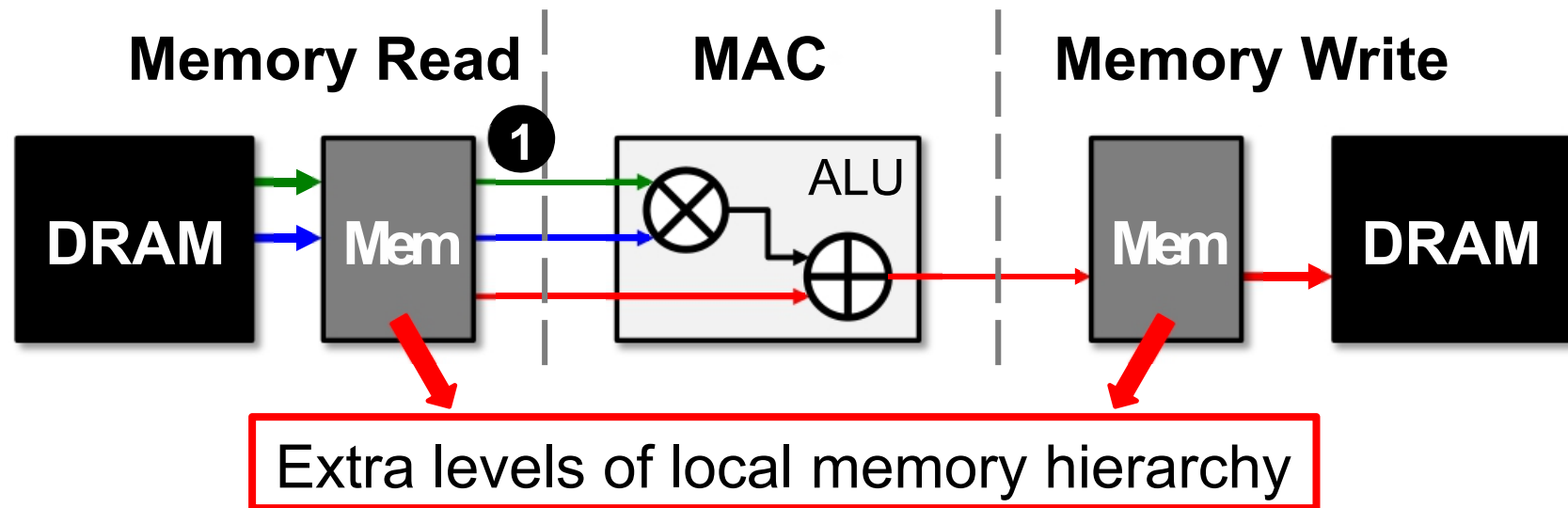
Worst Case: all memory R/W are **DRAM** accesses

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

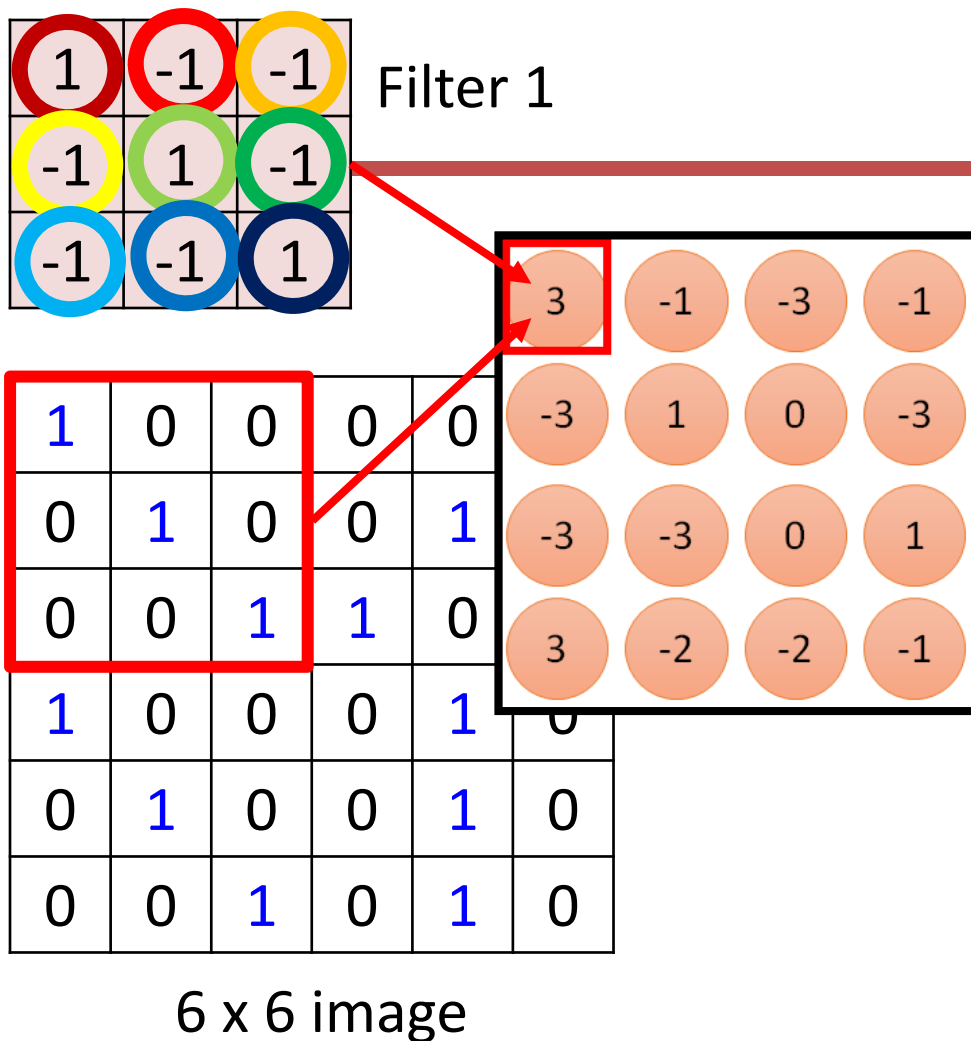
# Memory Access is the Bottleneck



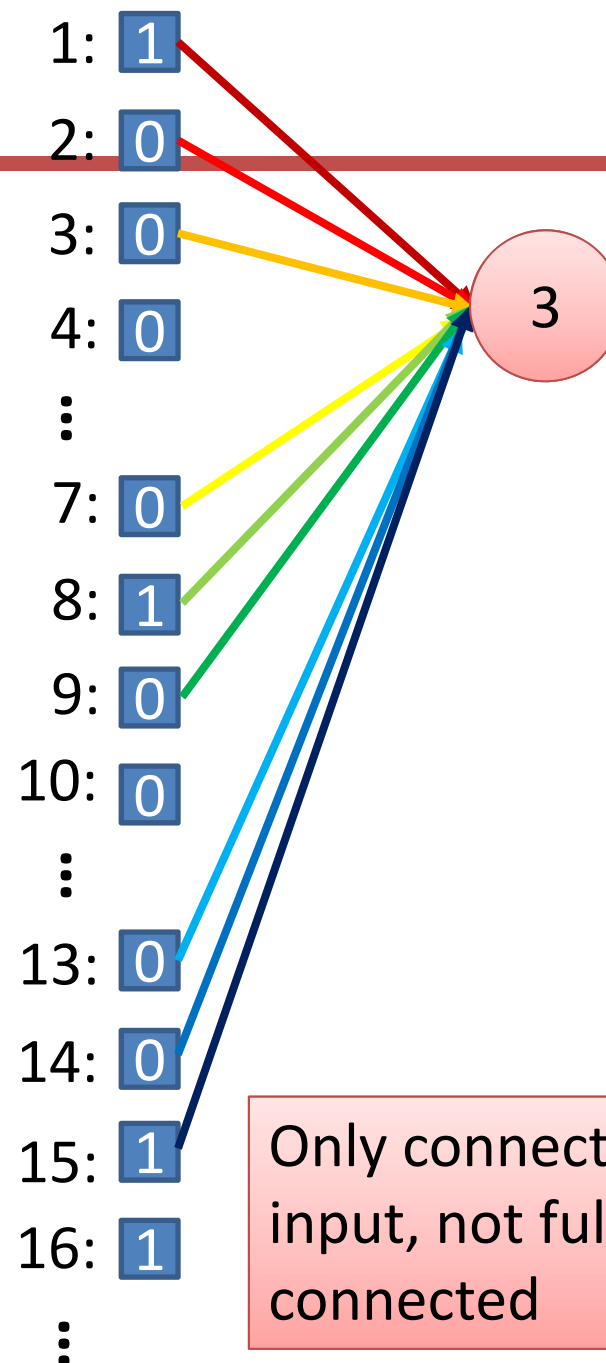
# Memory Access is the Bottleneck

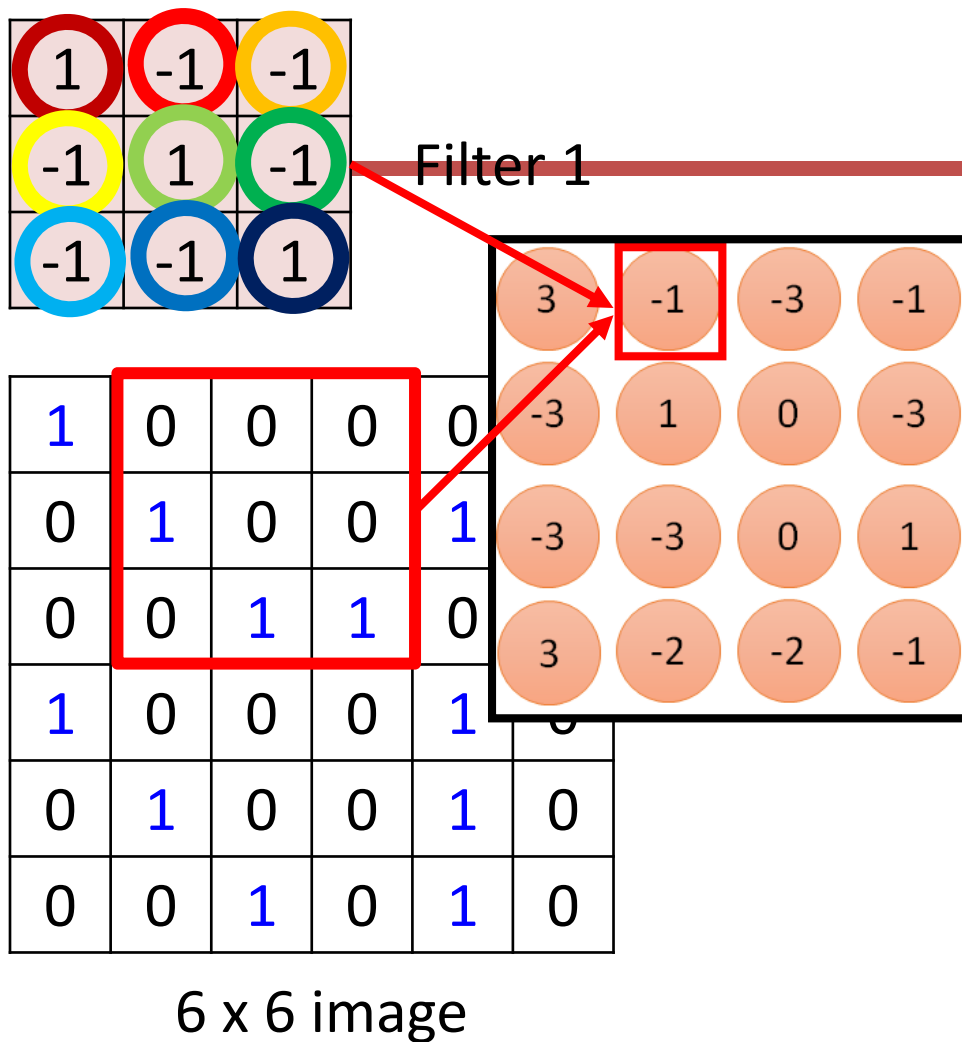


Opportunities: **1** data reuse



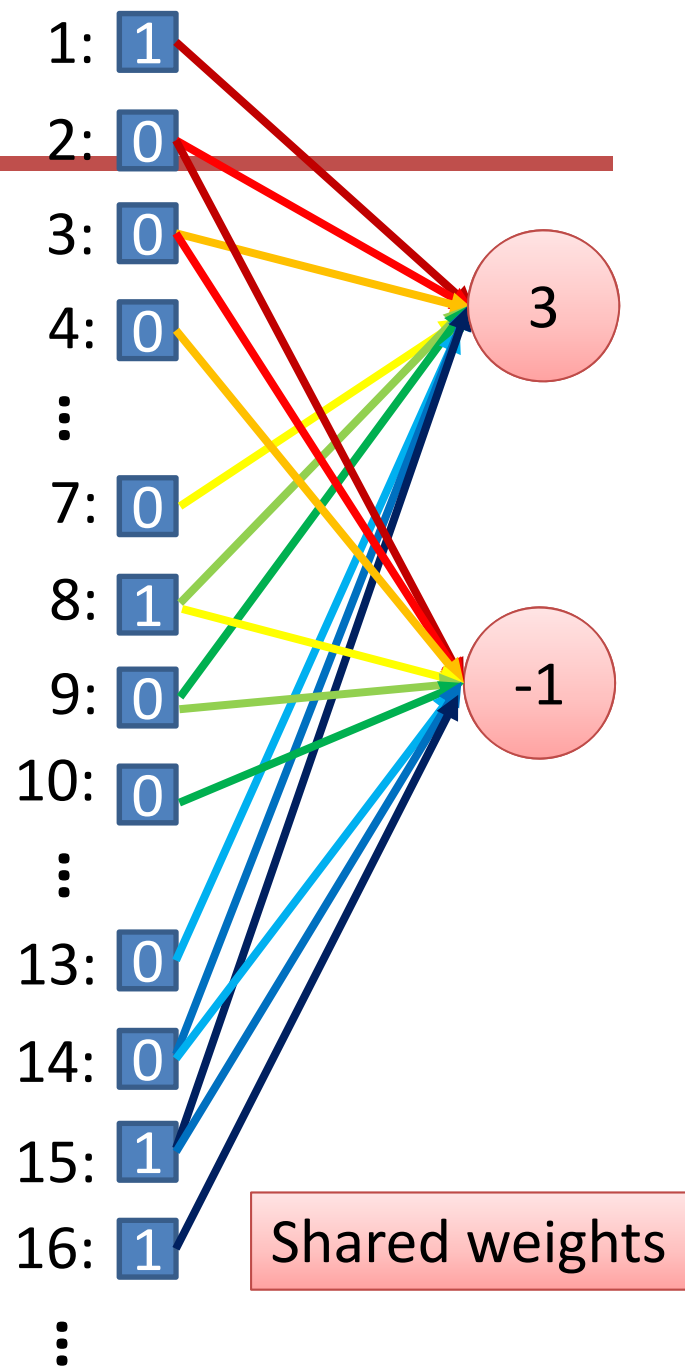
Less parameters!





Less parameters!

Even less parameters!

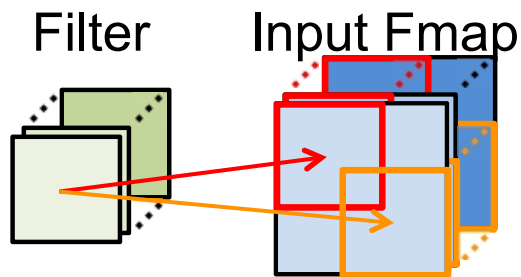




# Types of Data Reuse in DNN

## Convolutional Reuse

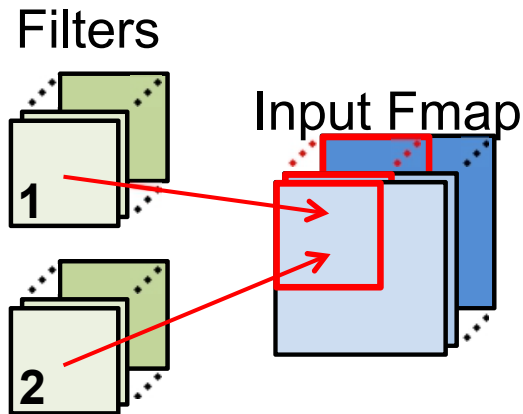
In **CONV** layers only  
(sliding window)



Reuse Both **Activations**  
and **Filter weights**

## Fmap Reuse

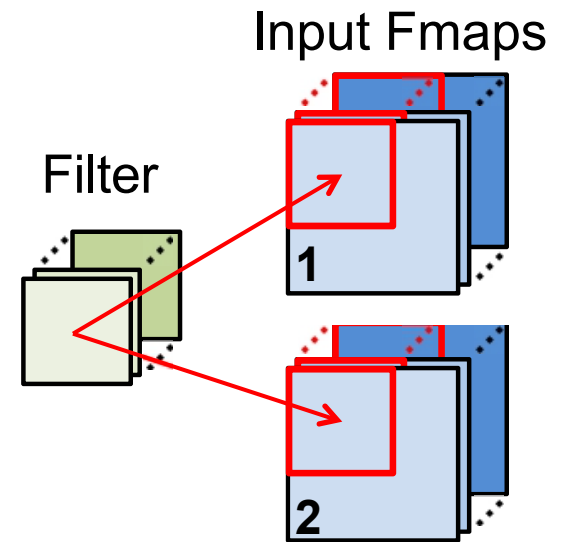
In both **CONV** and FC  
layers



Reuse **Activations**

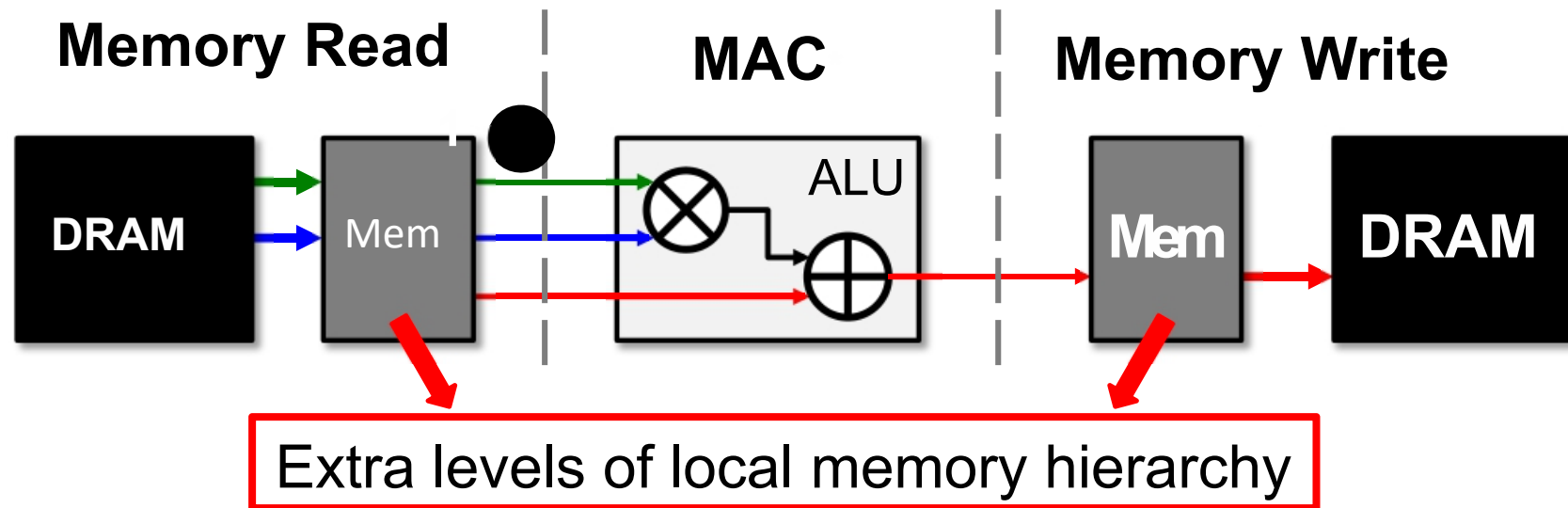
## Filter Reuse

In both **CONV** and **FC**  
layers (batch size > 1)



Reuse **Filter weights**

# Memory Access is the Bottleneck

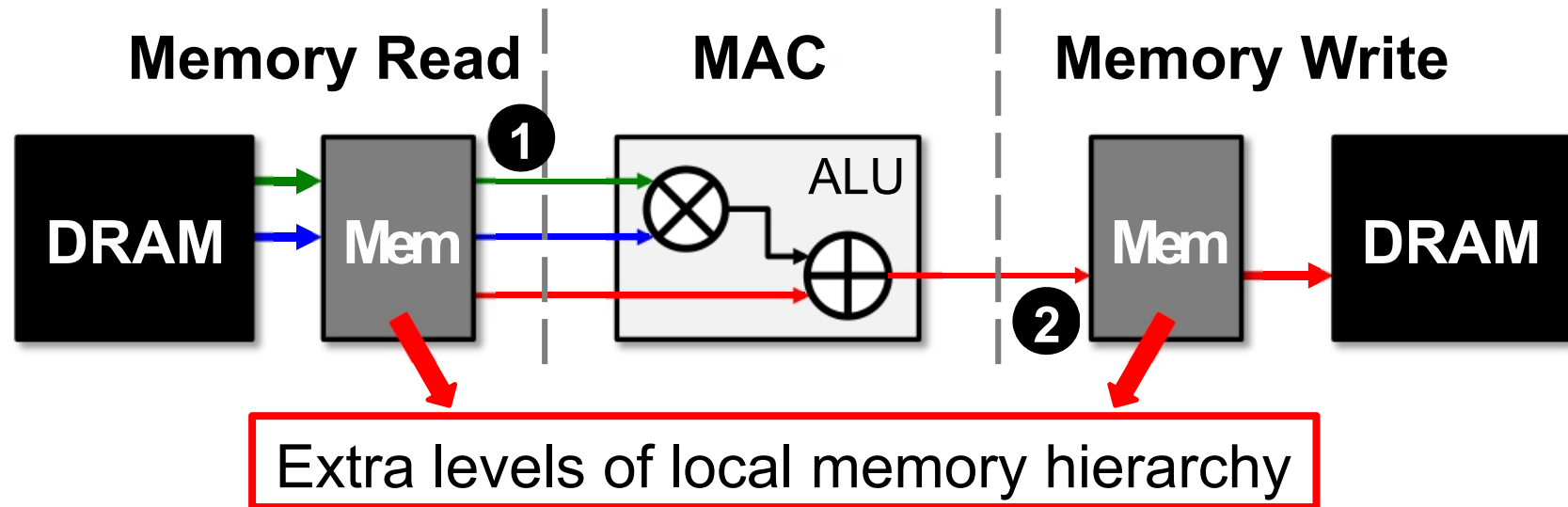


Opportunities: ① data reuse

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

\*\* AlexNet CONV layers

# Memory Access is the Bottleneck

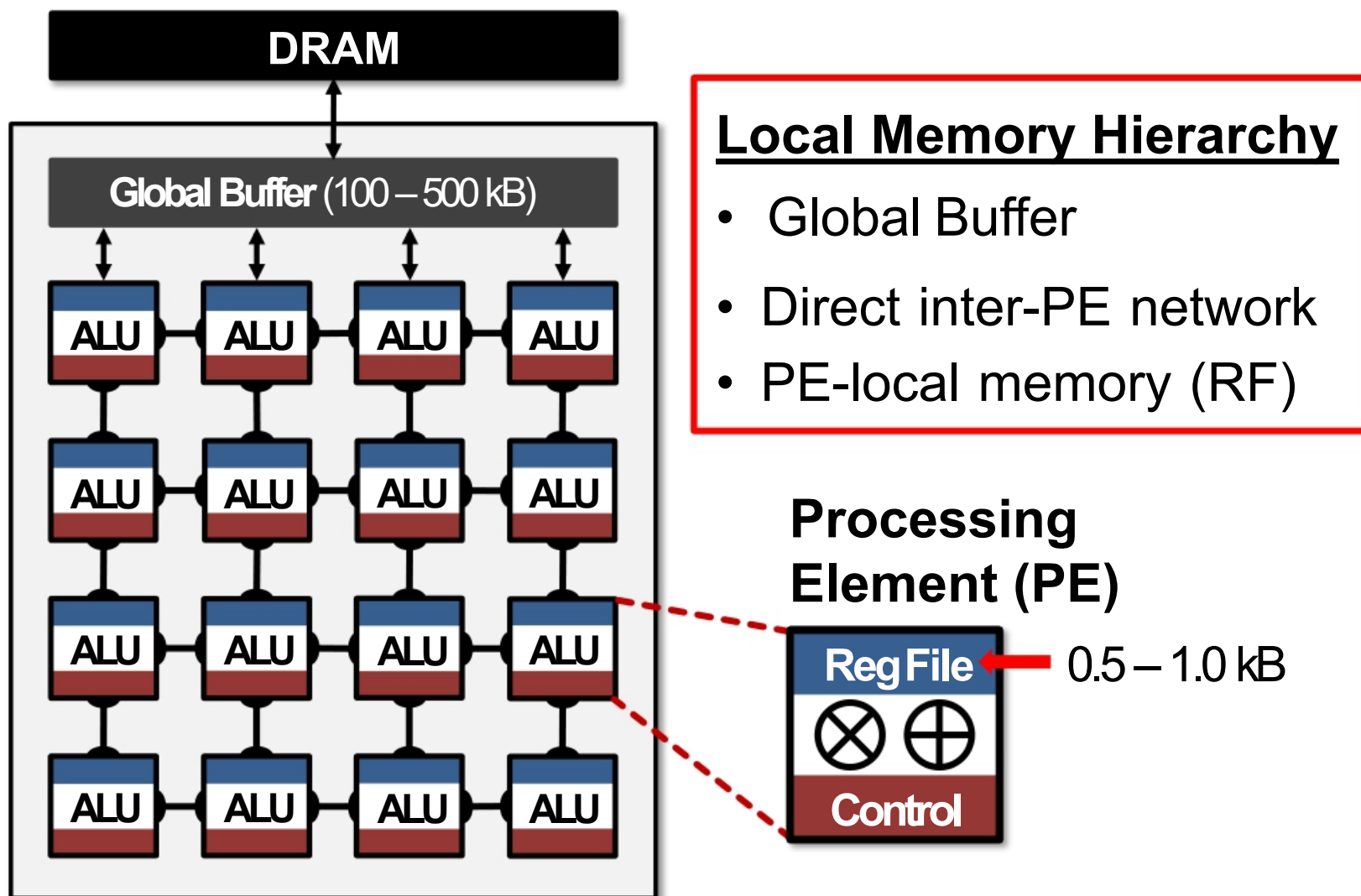


Opportunities: ① data reuse ② local accumulation

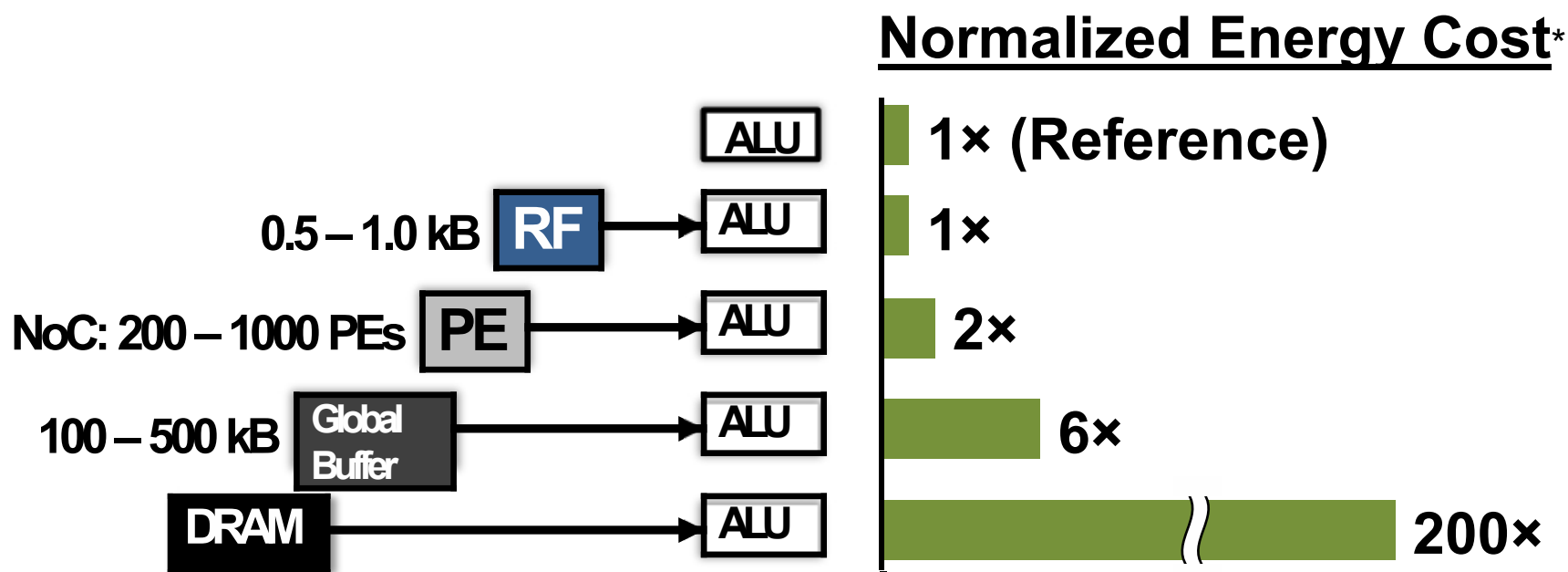
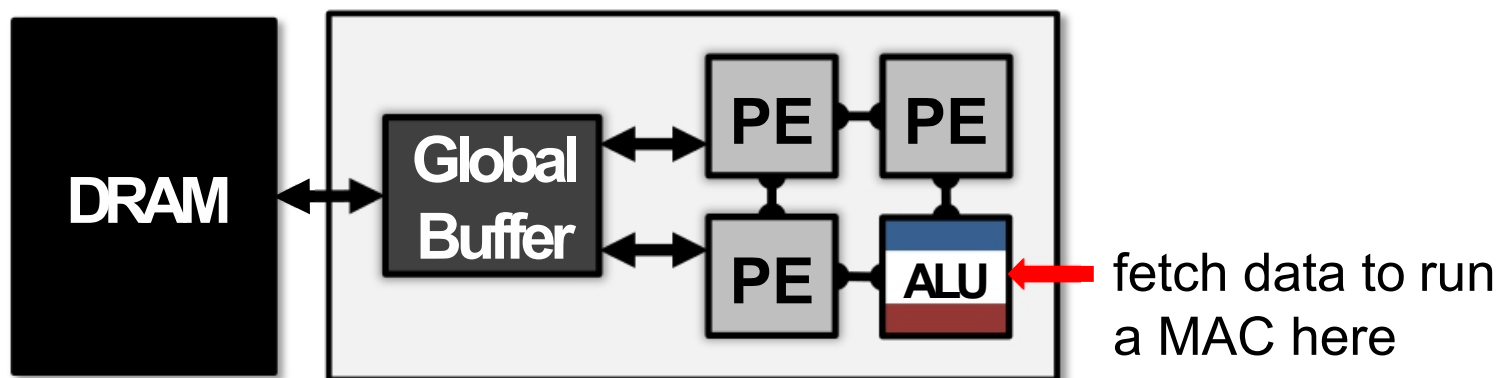
- 1) Can reduce DRAM reads of **filter/fmap** by up to **500×**
- 2) **Partial sum** accumulation does **NOT** have to access DRAM

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

# Spatial Architecture for DNN



# Low-Cost Local Data Access

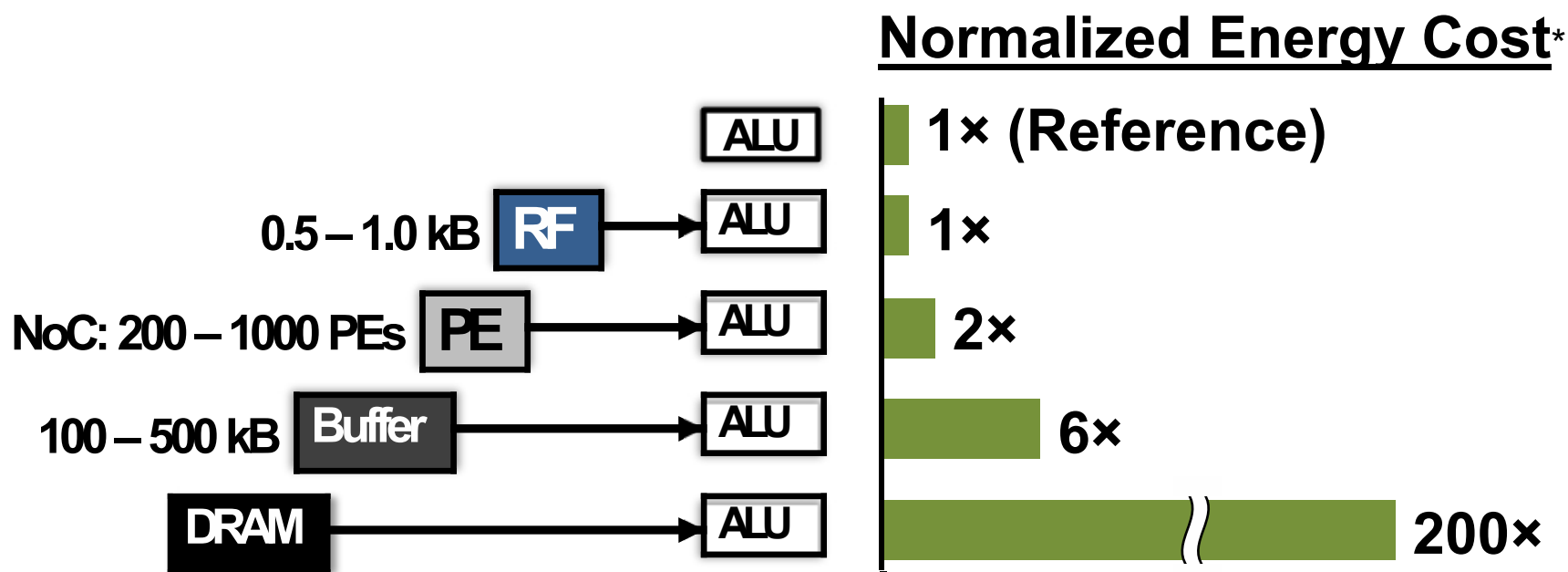


\* measured from a commercial 65nm process<sup>16</sup>

# Low-Cost Local Data Access

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

**specialized processing dataflow required!**



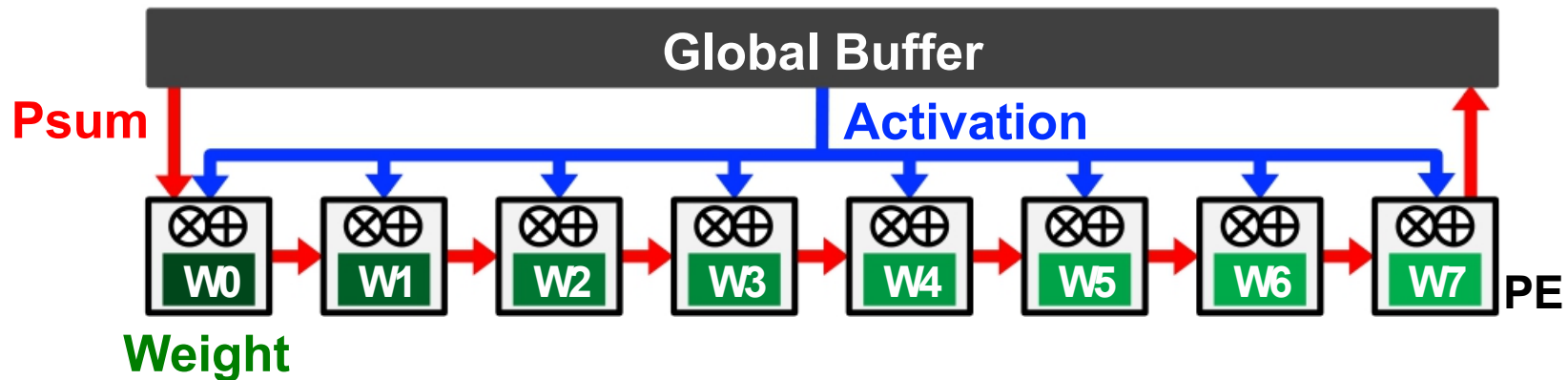
\* measured from a commercial 65nm process<sup>17</sup>

---

# Dataflow Taxonomy

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

# Weight Stationary (WS)

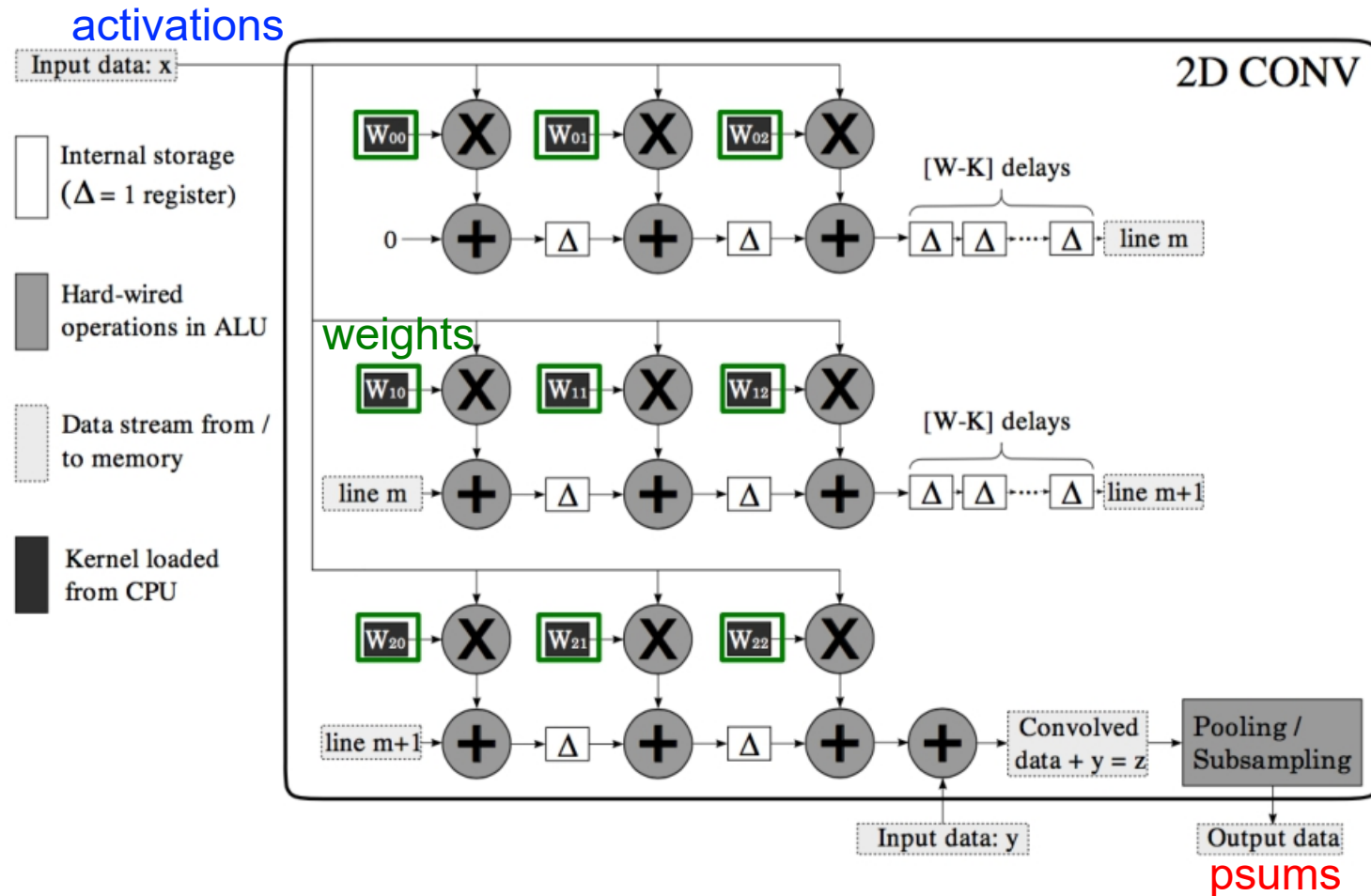


- **Minimize **weight**** read energy consumption
  - maximize convolutional and filter reuse of weights
- **Broadcast **activations**** and **accumulate **psums**** spatially across the PE array.

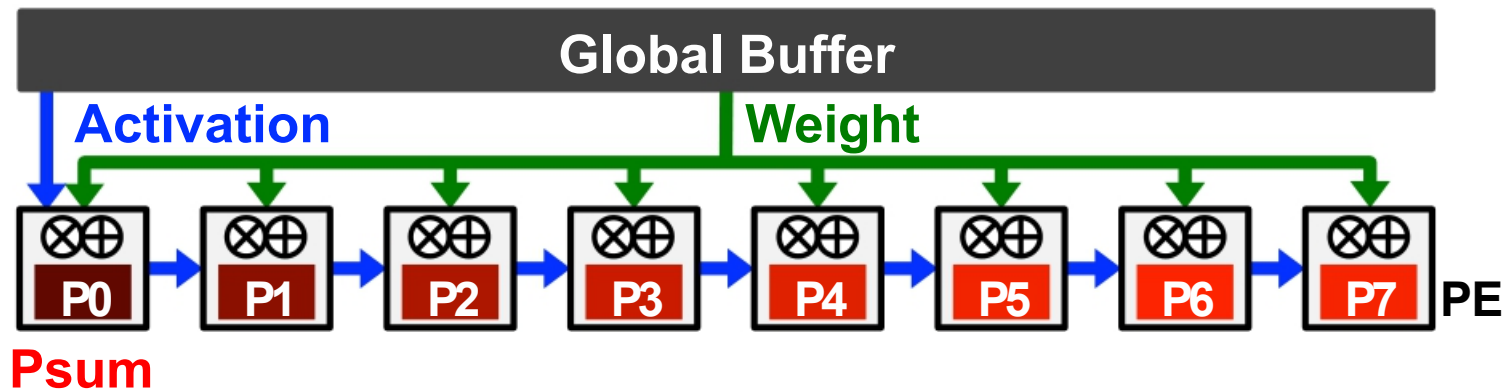


# WS Example: nn-X (NeuFlow)

## A 3×3 2D Convolution Engine



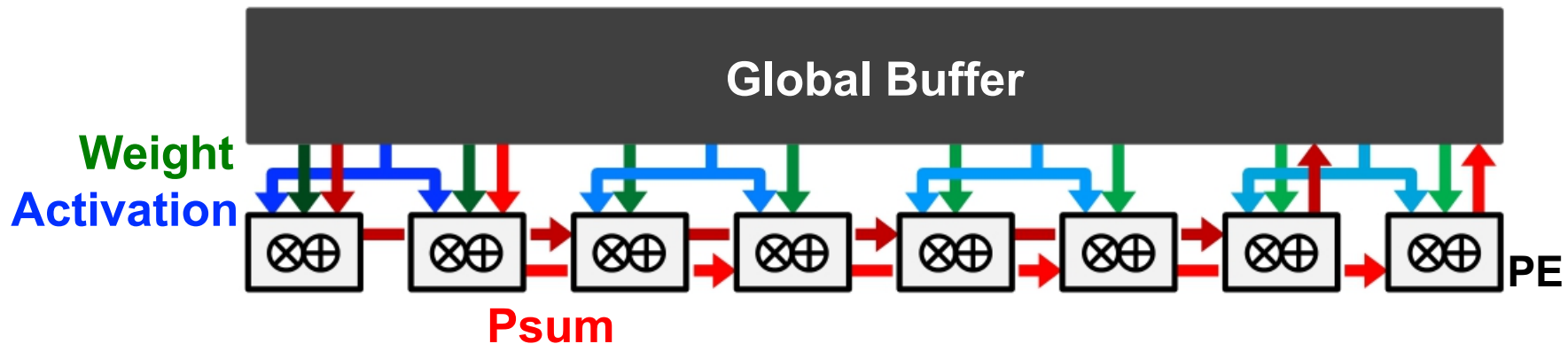
# Output Stationary (OS)



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

# No Local Reuse (NLR)

---



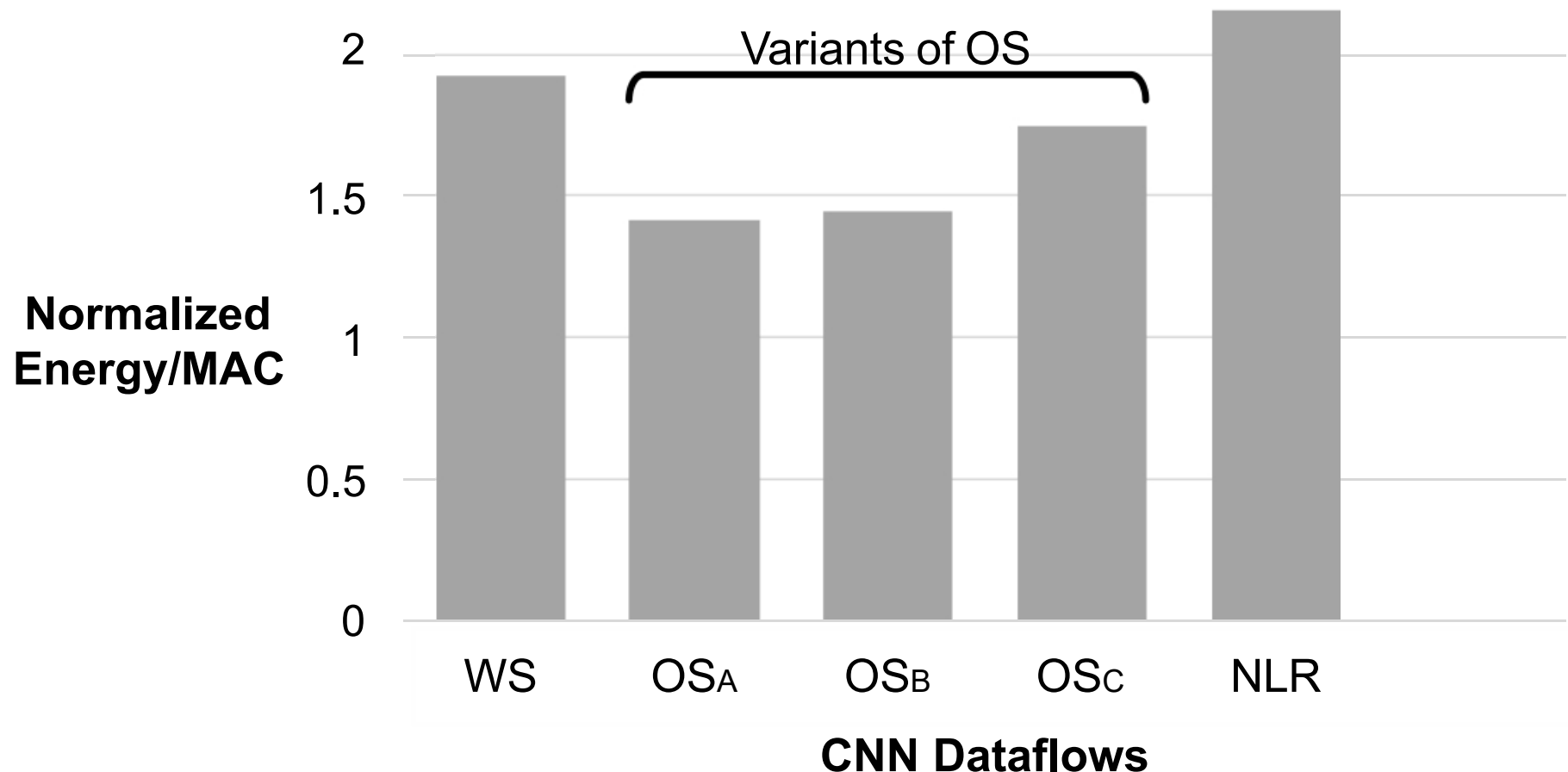
- Use a **large global buffer** as shared storage
  - Reduce **DRAM** access energy consumption
- Multicast **activations**, single-cast **weights**, and accumulate **psums** spatially across the PE array

# Matrix Multiply Unit



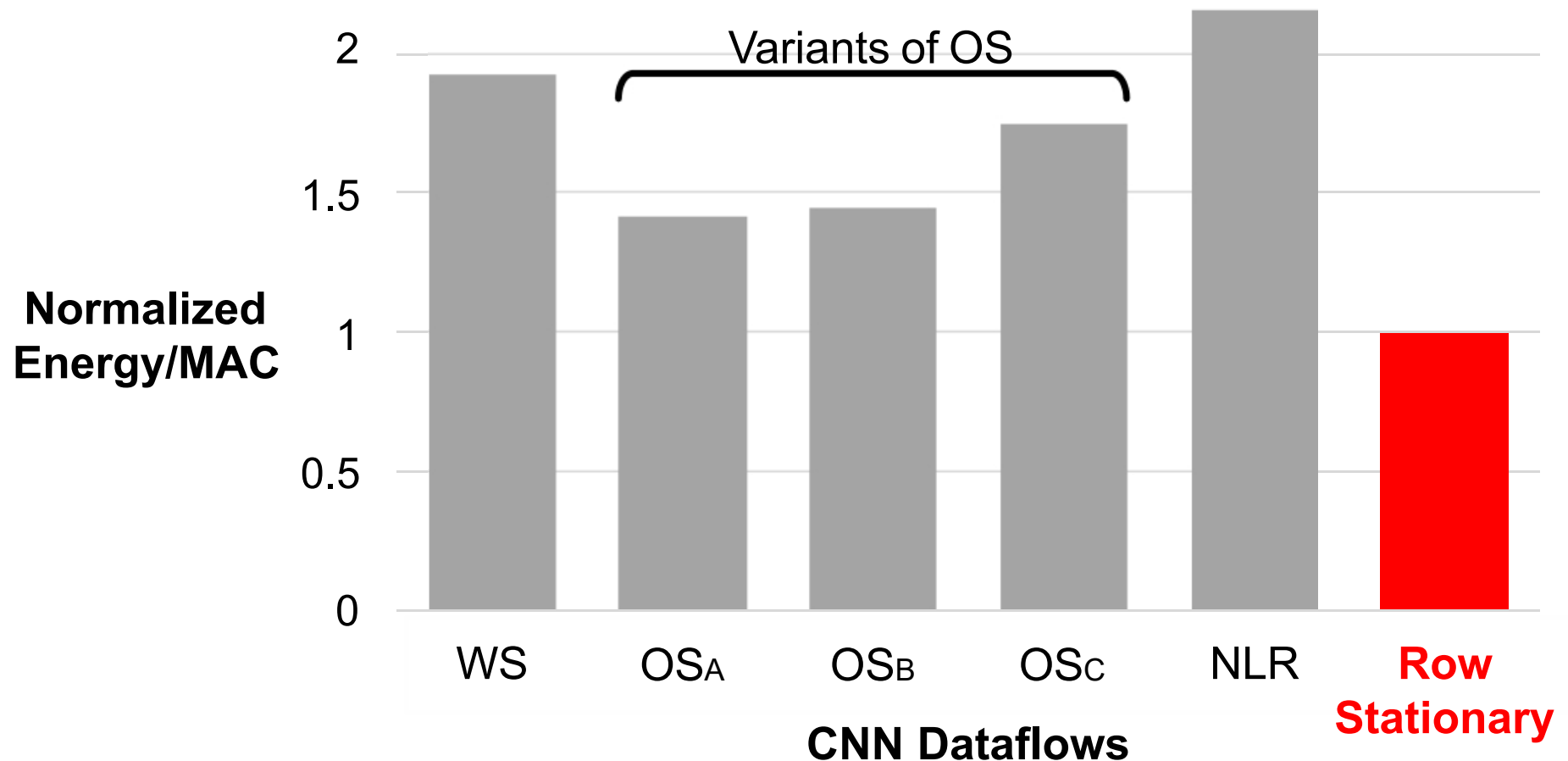
# Energy Efficiency Comparison

- 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



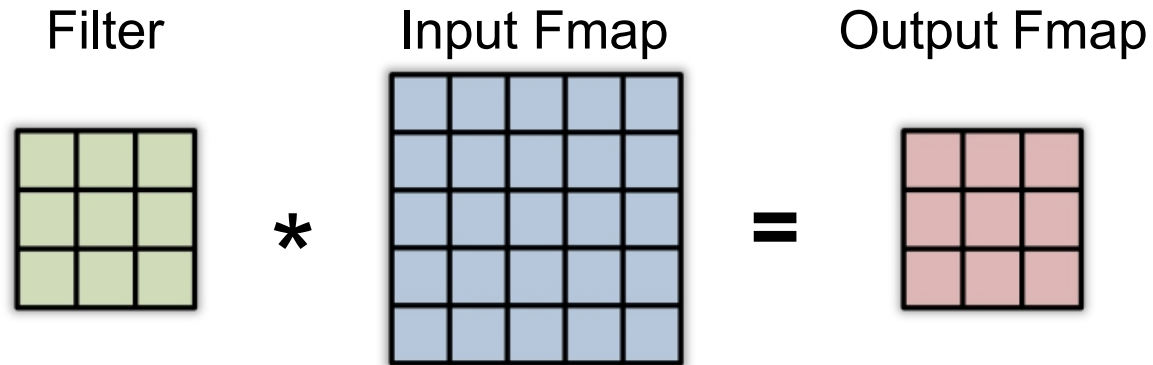
---

## 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

Filter

a	b	c

\*

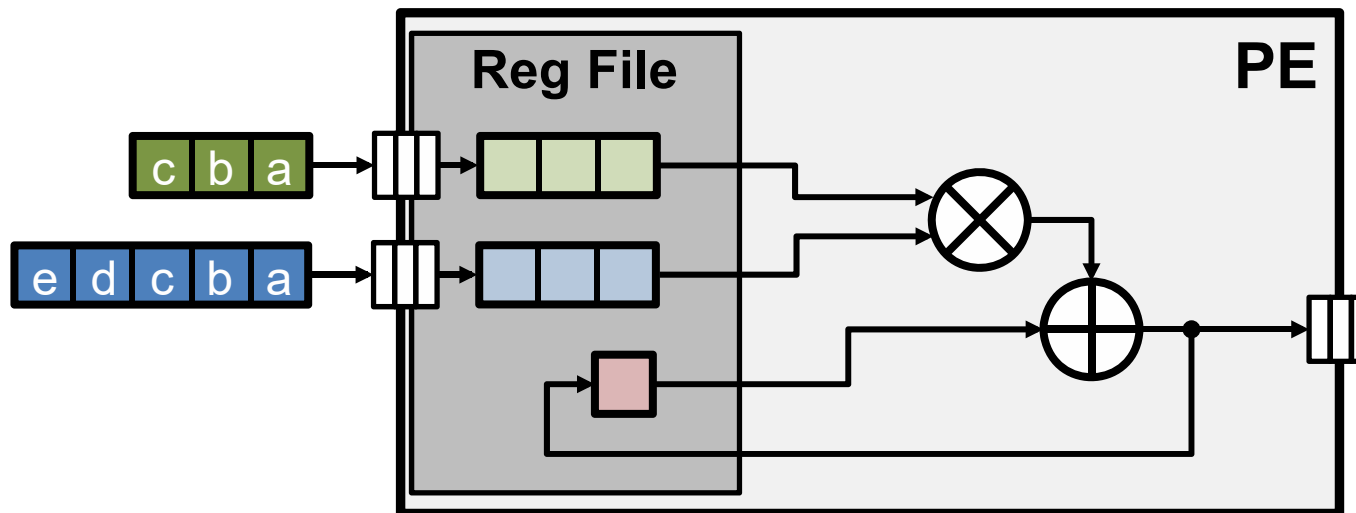
Input Fmap

a	b	c	d	e

=

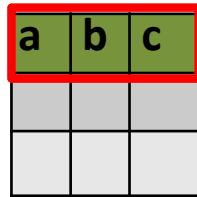
Partial Sums

a	b	c



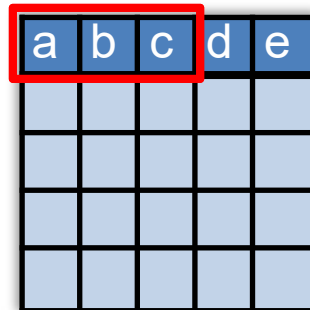
# 1D Row Convolution in PE

Filter



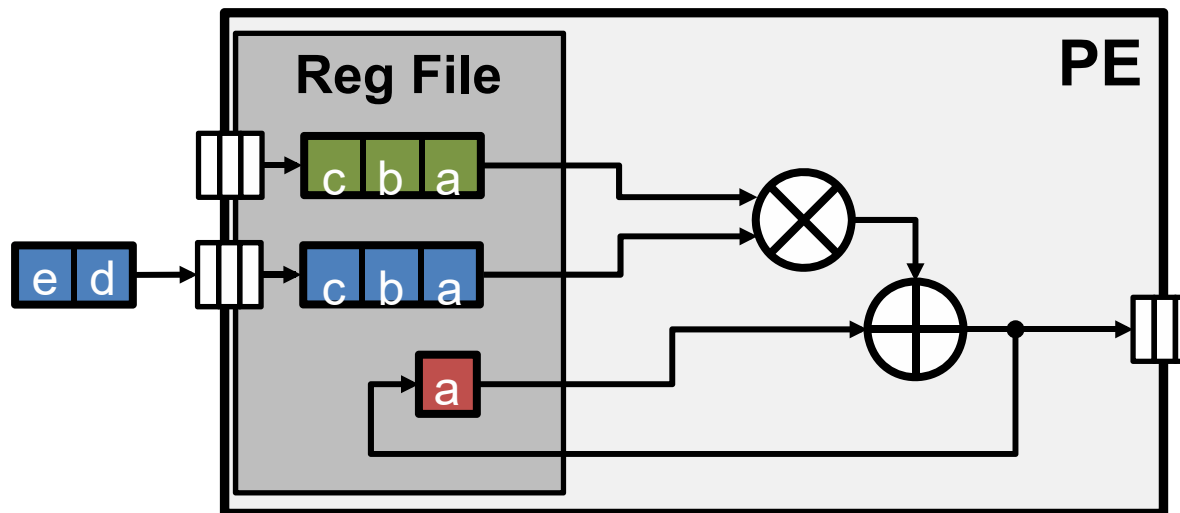
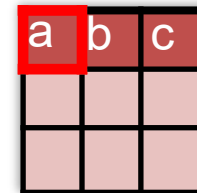
\*

Input Fmap

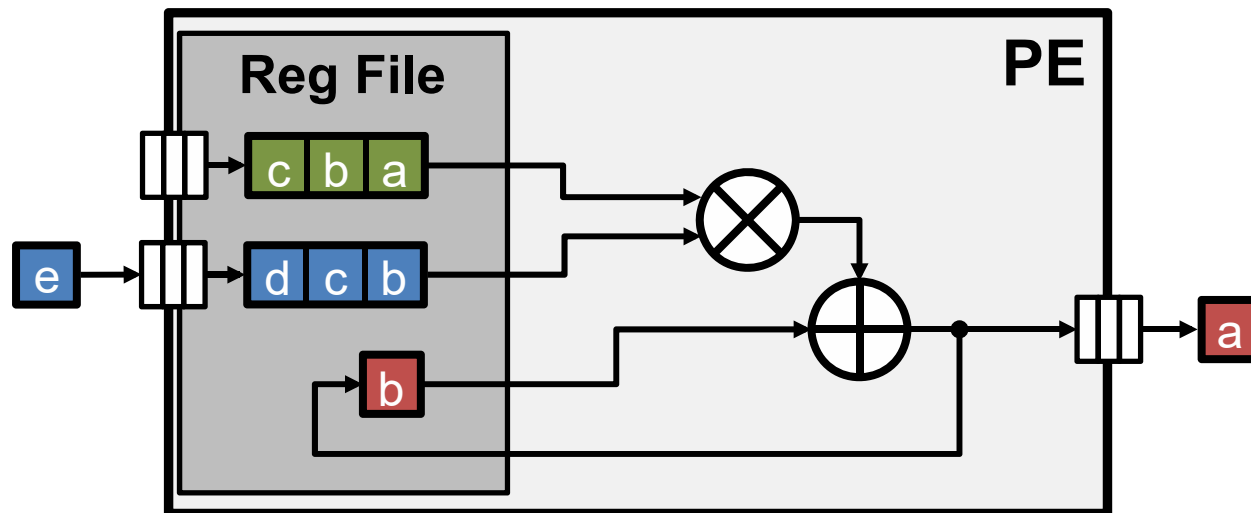
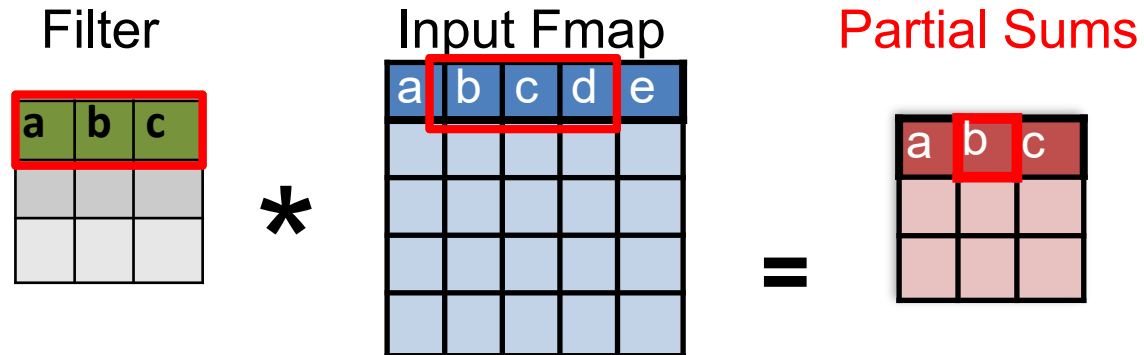


=

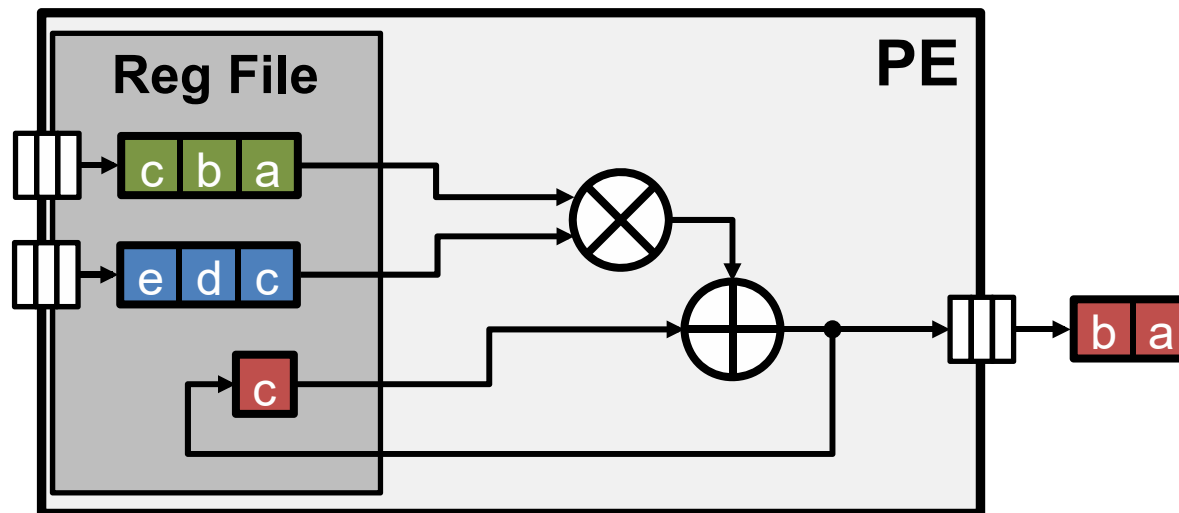
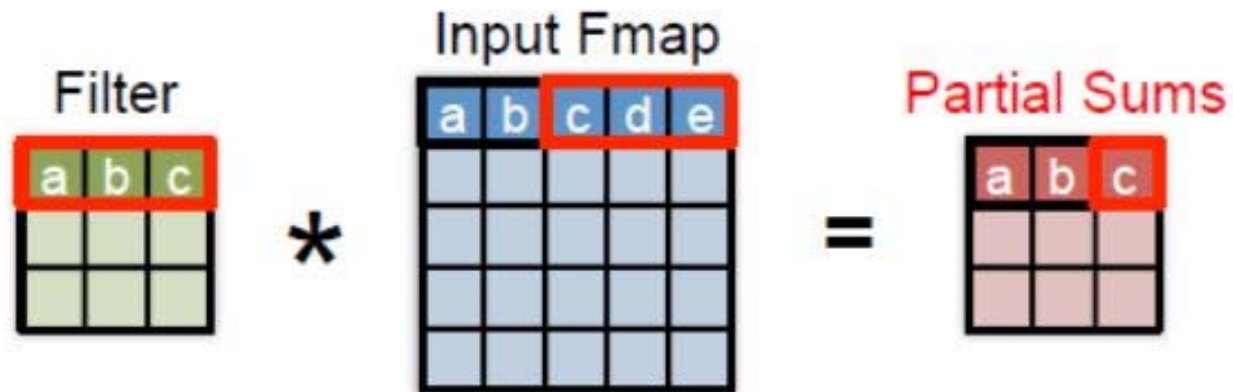
Partial Sums



# 1D Row Convolution in PE

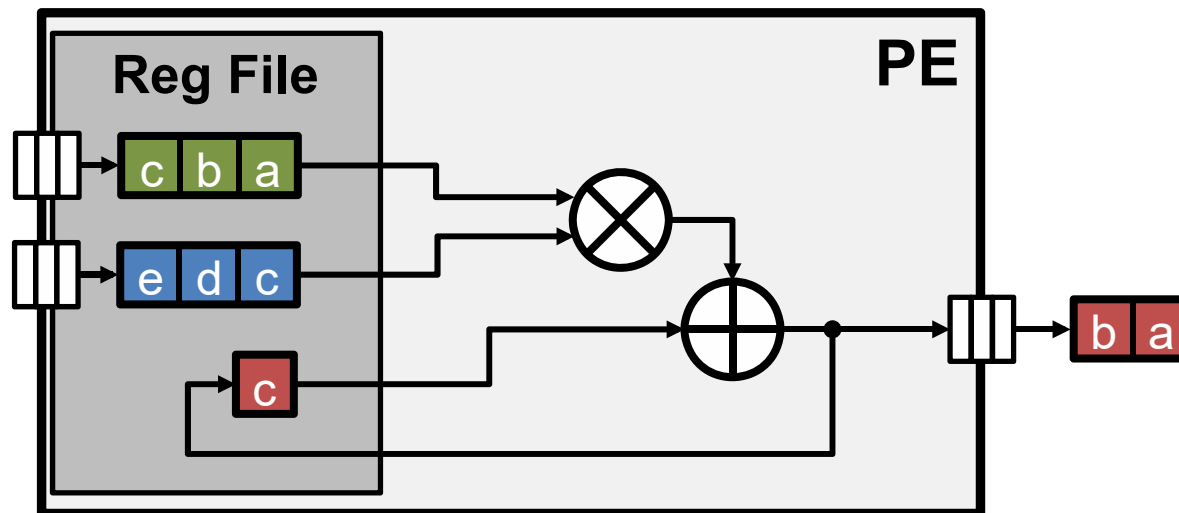


# 1D Row Convolution in PE



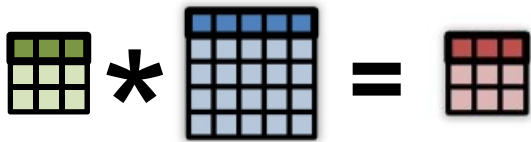
# 1D Row Convolution in PE

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



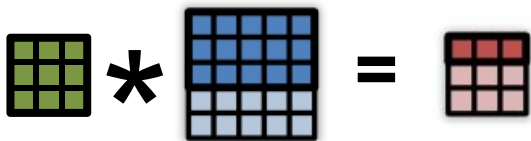
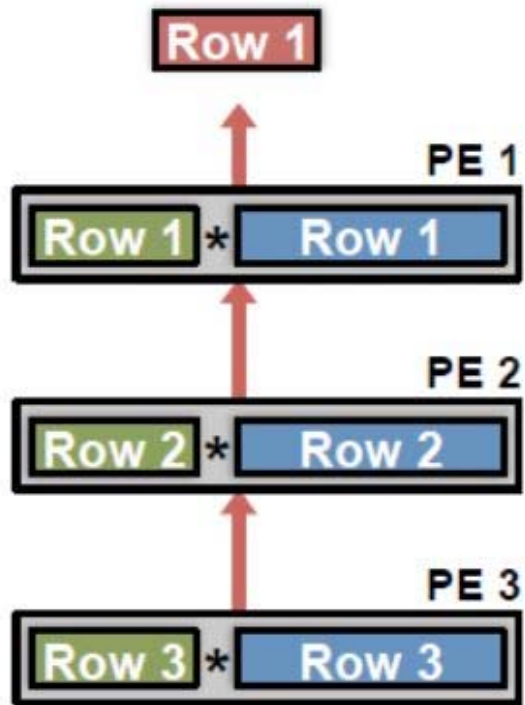
# 2D Convolution in PE Array

---

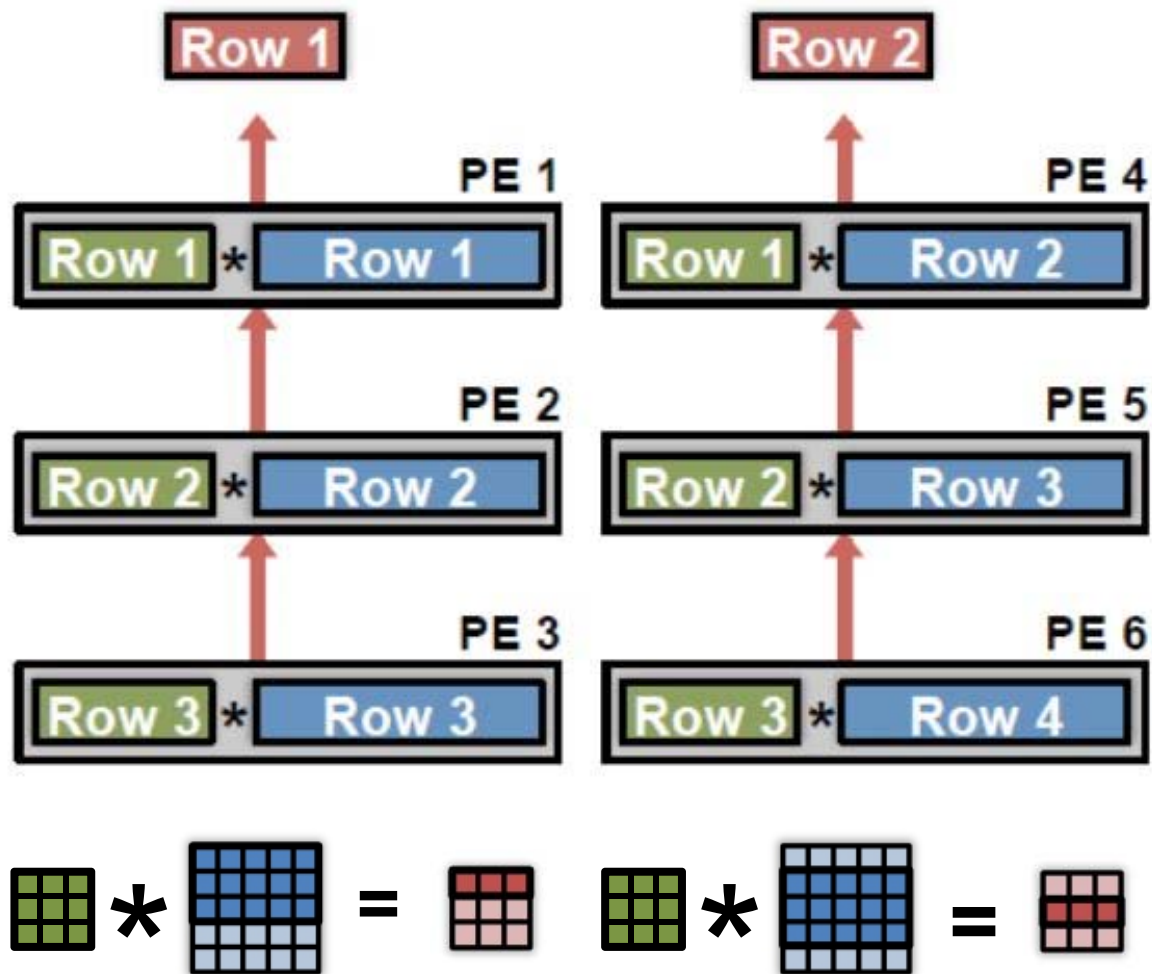


# 2D Convolution in PE Array

---

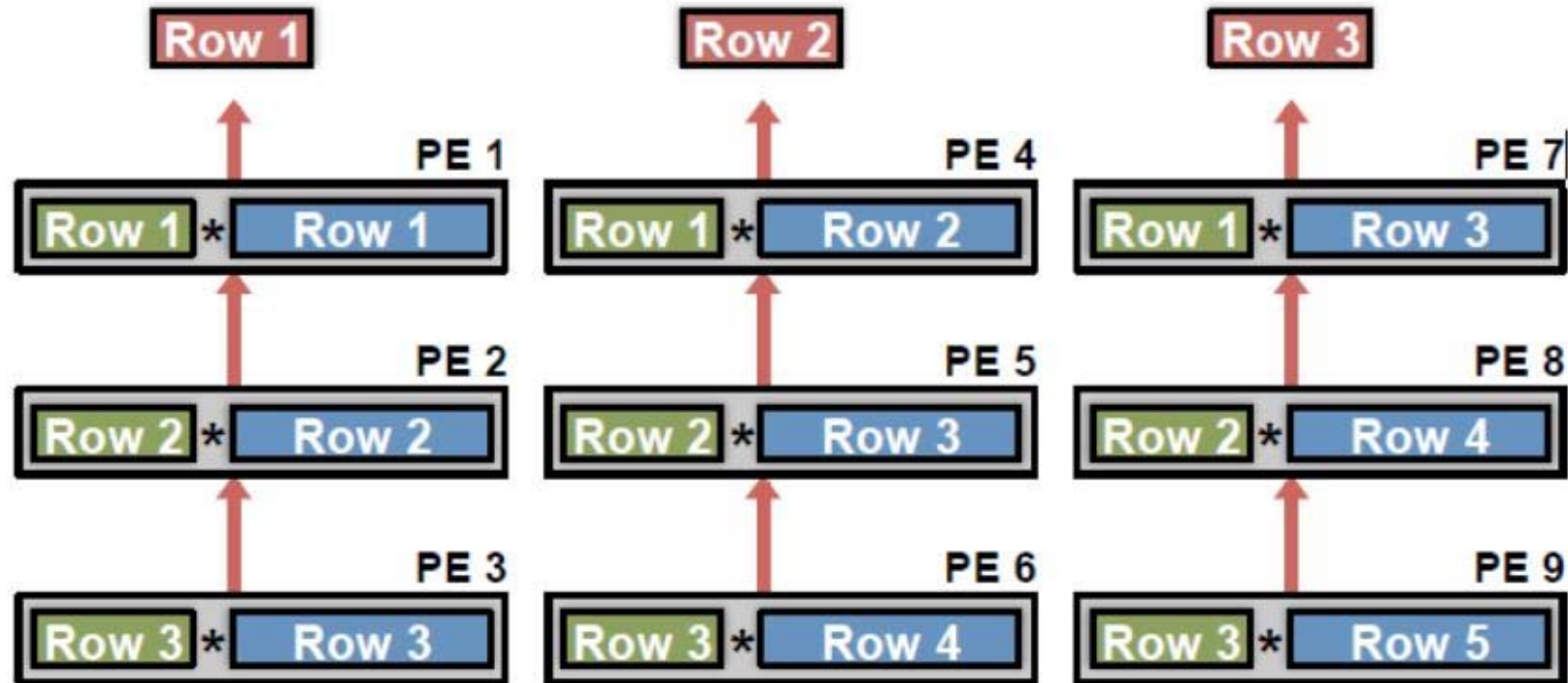


# 2D Convolution in PE Array





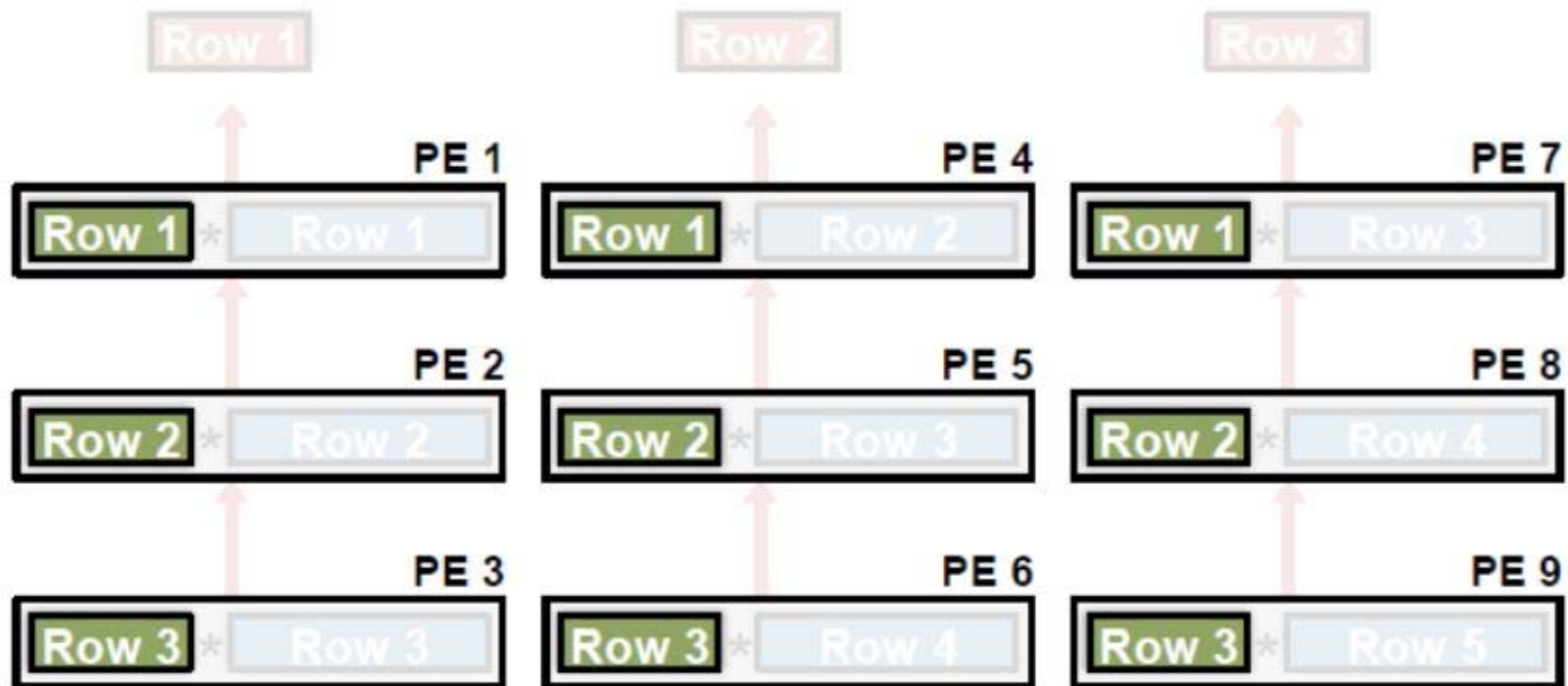
# 2D Convolution in PE Array



$$\begin{array}{c} \text{3x3 Green} \end{array} * \begin{array}{c} \text{5x5 Blue} \end{array} = \begin{array}{c} \text{3x3 Red} \\ \text{3x3 Pink} \end{array} \quad
 \begin{array}{c} \text{3x3 Green} \end{array} * \begin{array}{c} \text{5x5 Blue} \end{array} = \begin{array}{c} \text{3x3 Red} \\ \text{3x3 Pink} \end{array} \quad
 \begin{array}{c} \text{3x3 Green} \end{array} * \begin{array}{c} \text{5x5 Blue} \end{array} = \begin{array}{c} \text{3x3 Red} \\ \text{3x3 Pink} \end{array}$$

# Convolutional Reuse Maximized

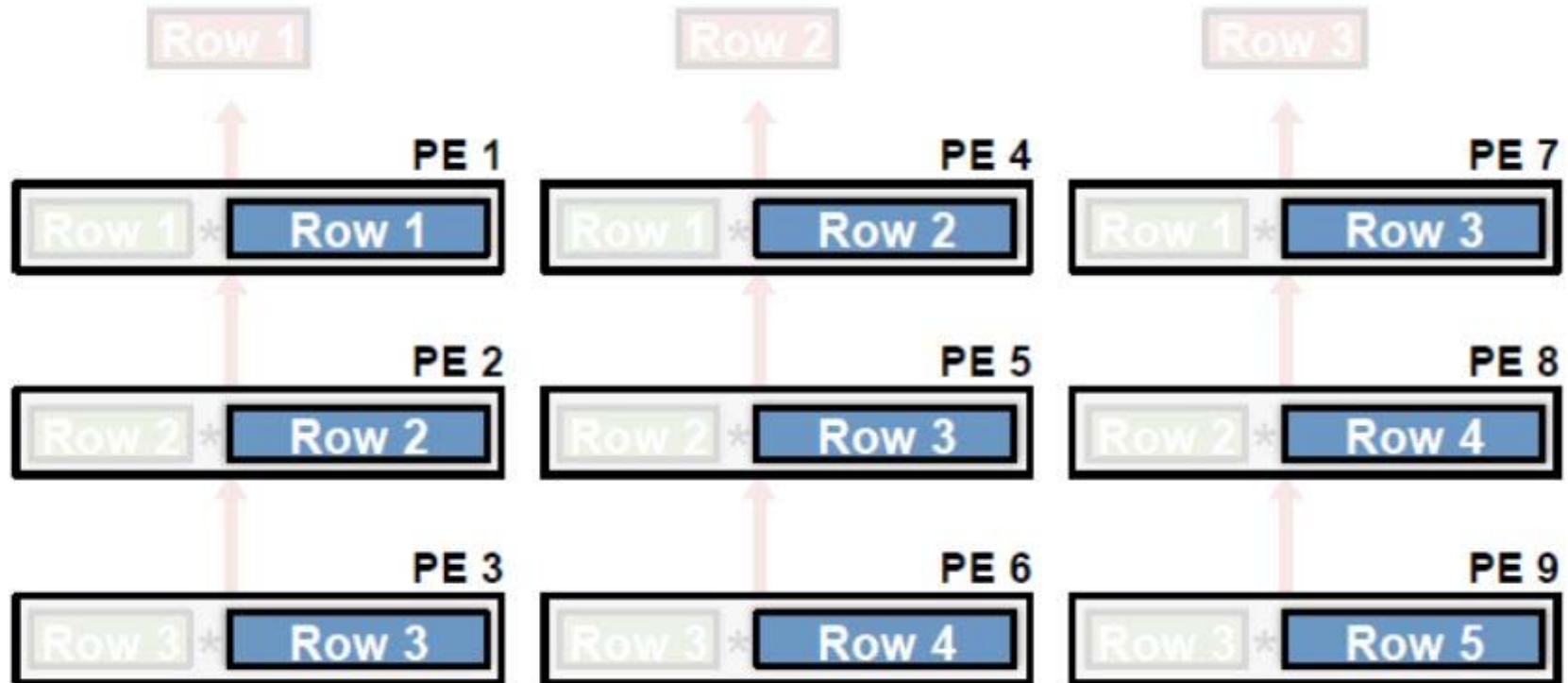
---



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

# Convolutional Reuse Maximized

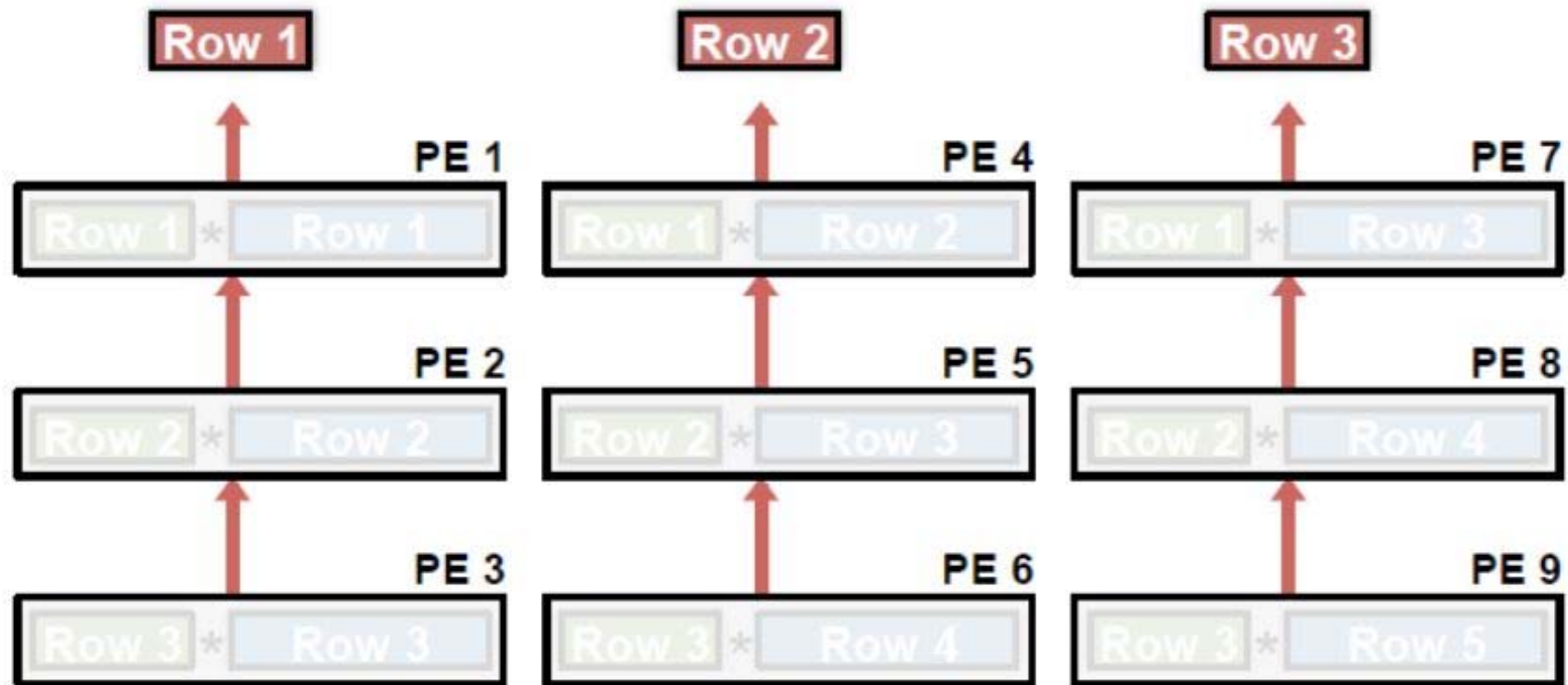
---



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

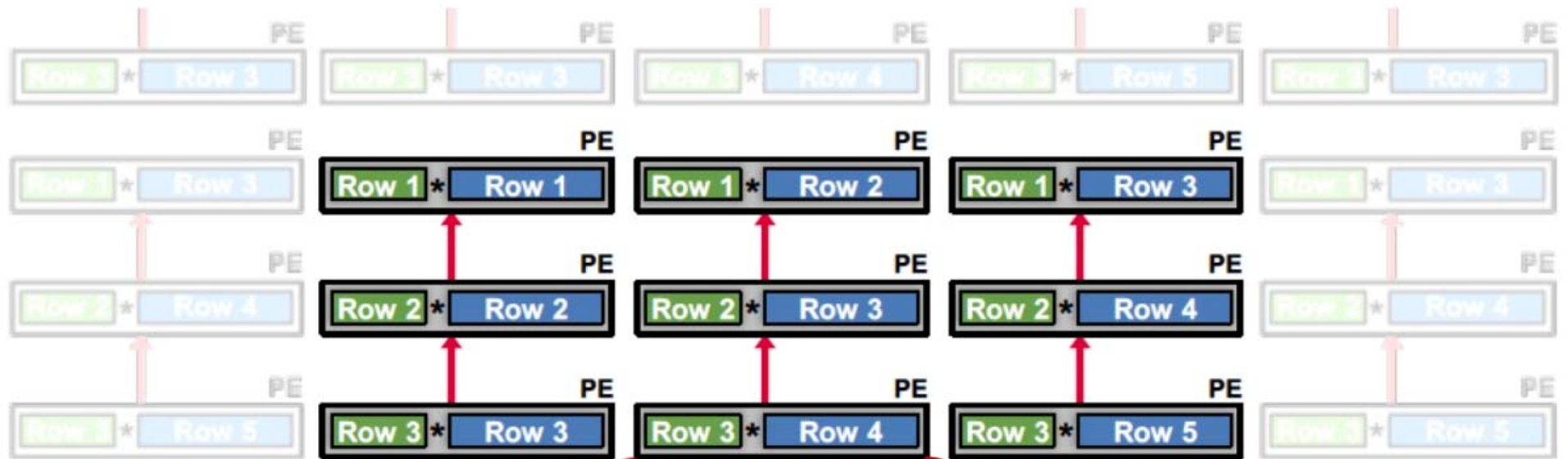
# Maximize 2D Accumulation in PE Array

---



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

# DNN Processing – The Full Picture



Multiple **fmaps**:



Multiple **filters**:

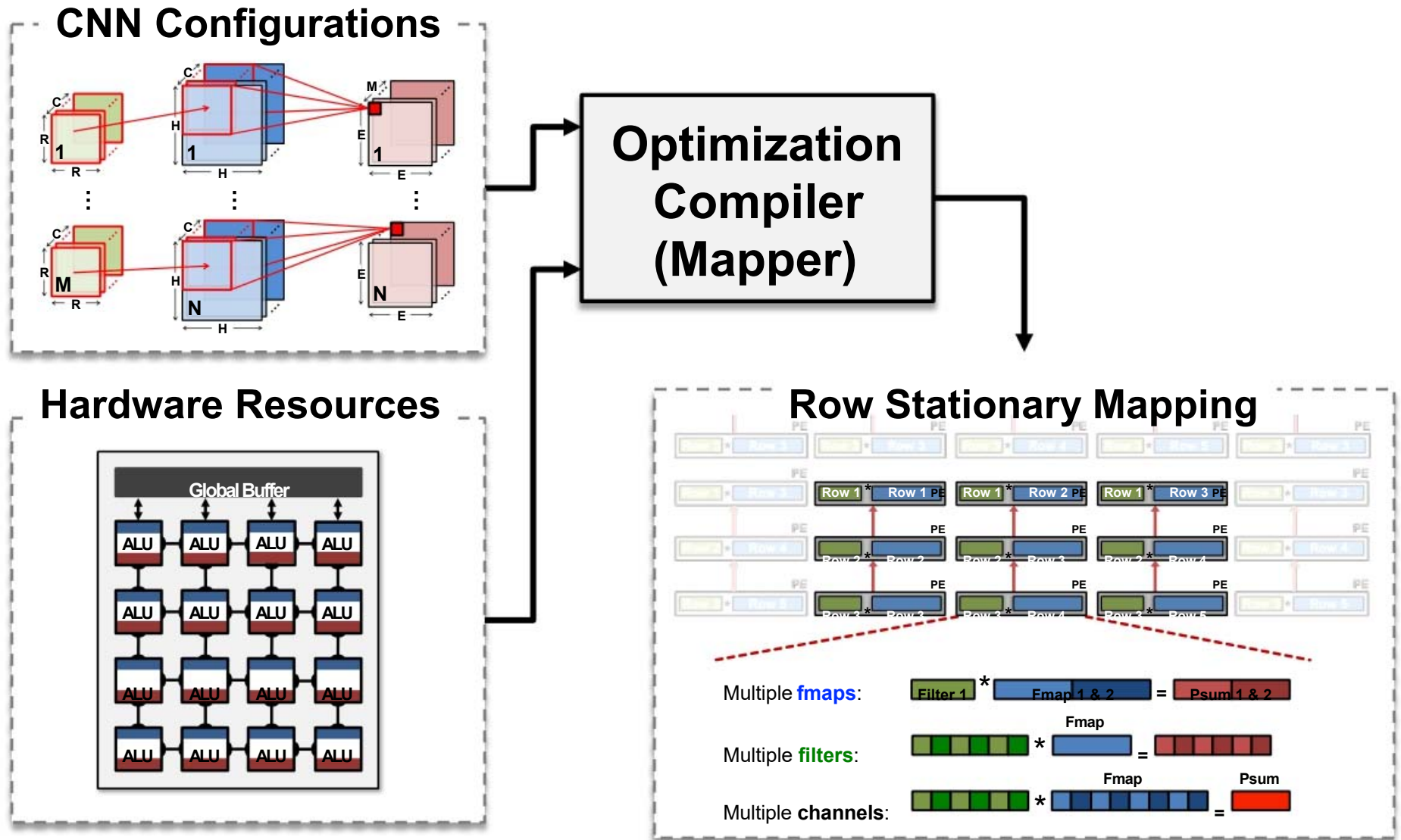


Multiple **channels**:

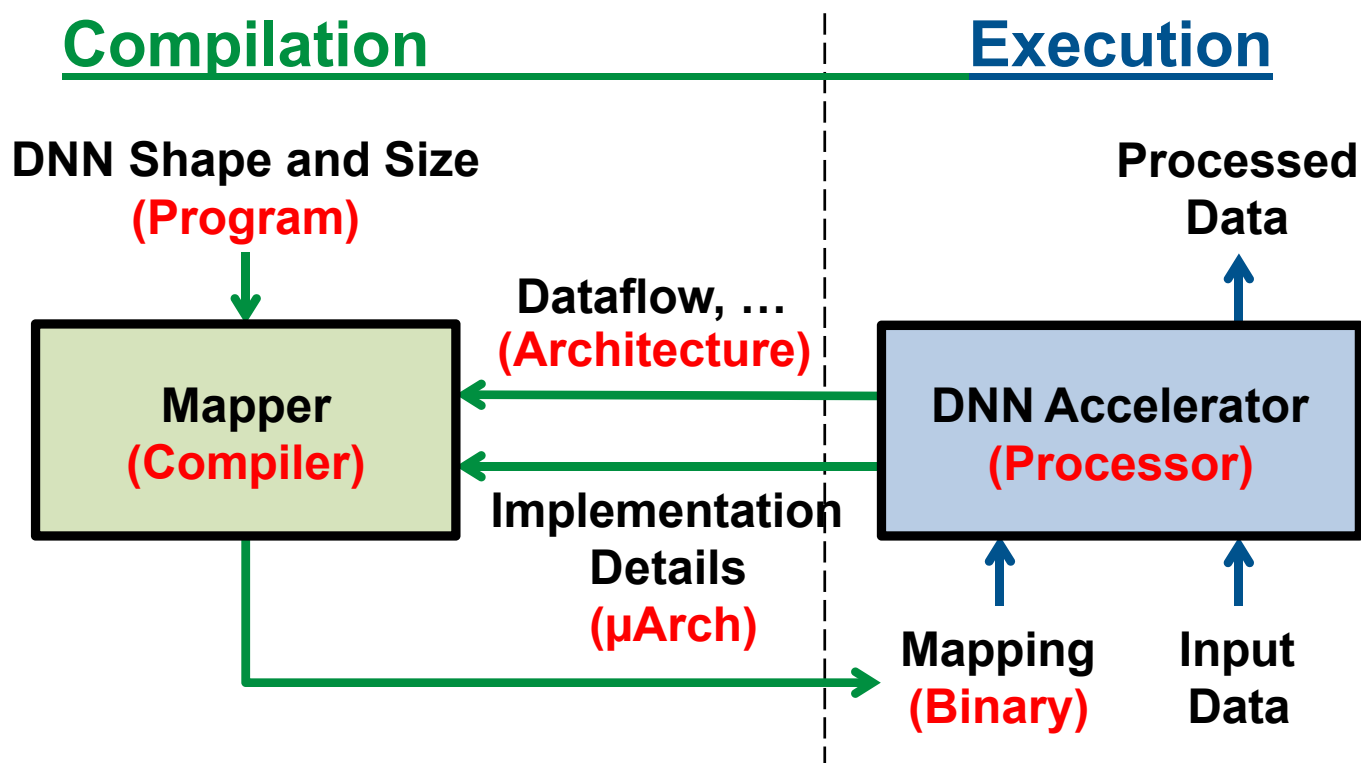


Map rows from **multiple fmaps**, **filters** and **channels** to same PE to exploit other forms of reuse and local accumulation

# Optimal Mapping in Row Stationary



# Computer Architecture Analogy



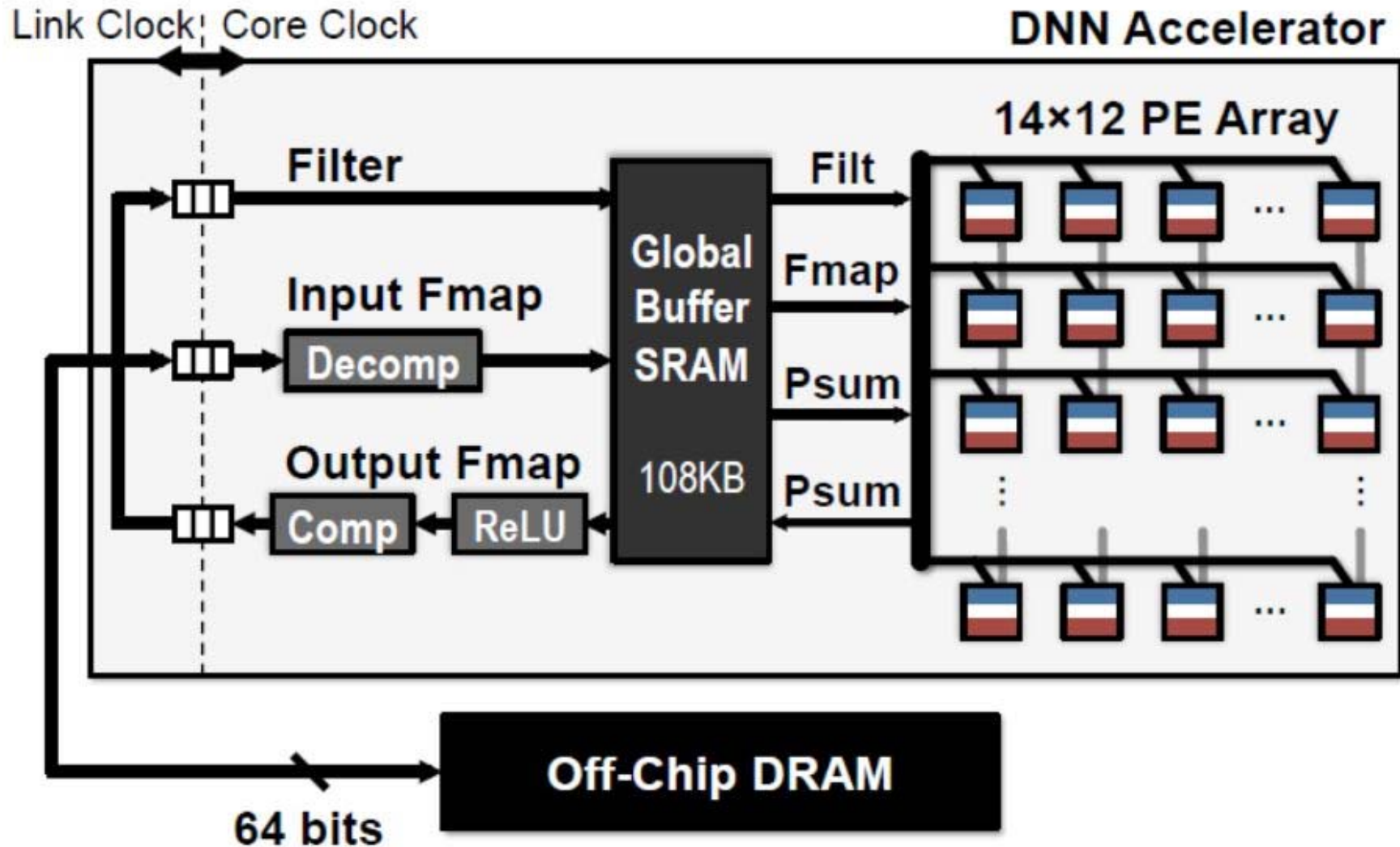


---

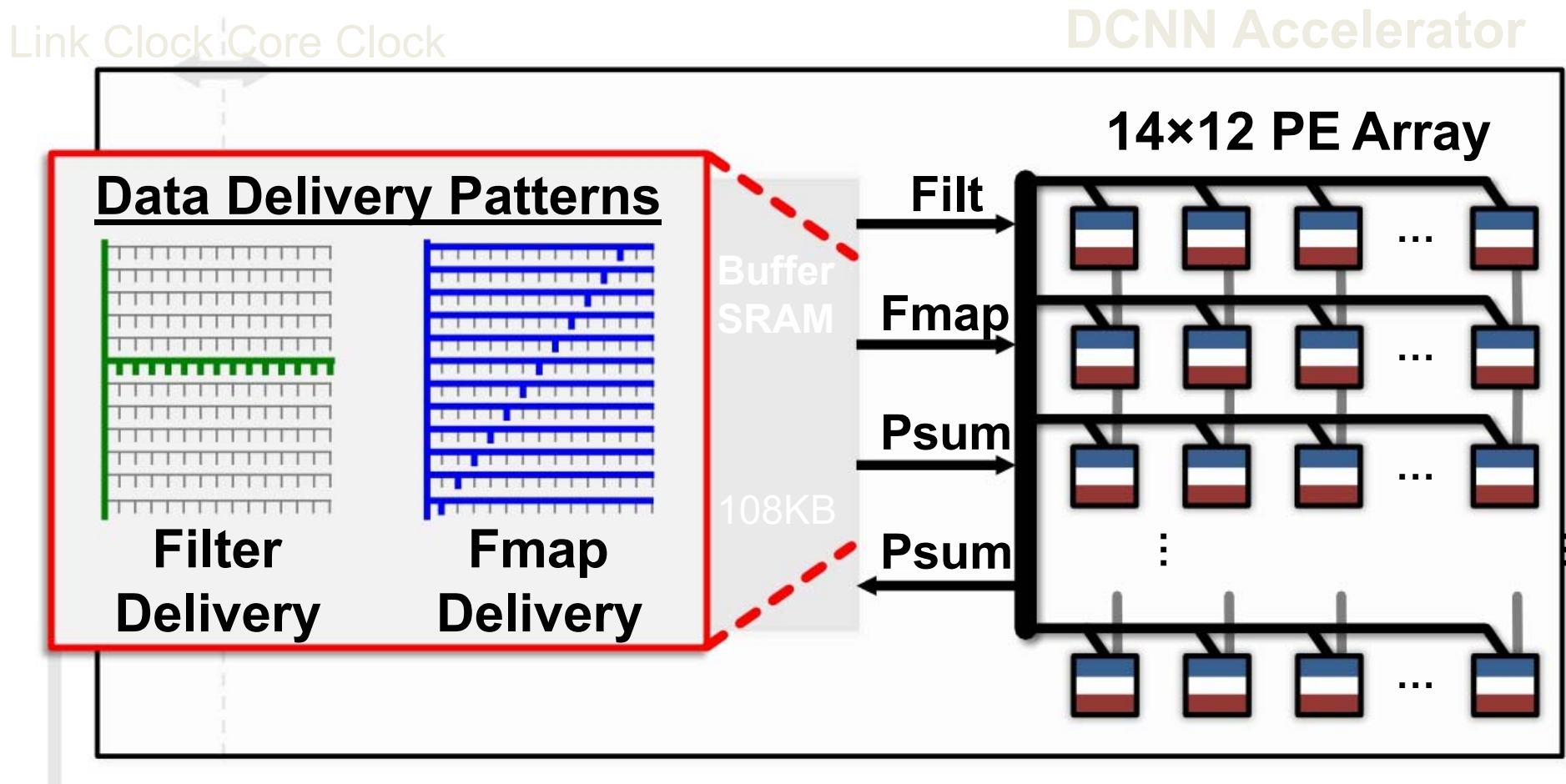
# Hardware Architecture for RS Dataflow



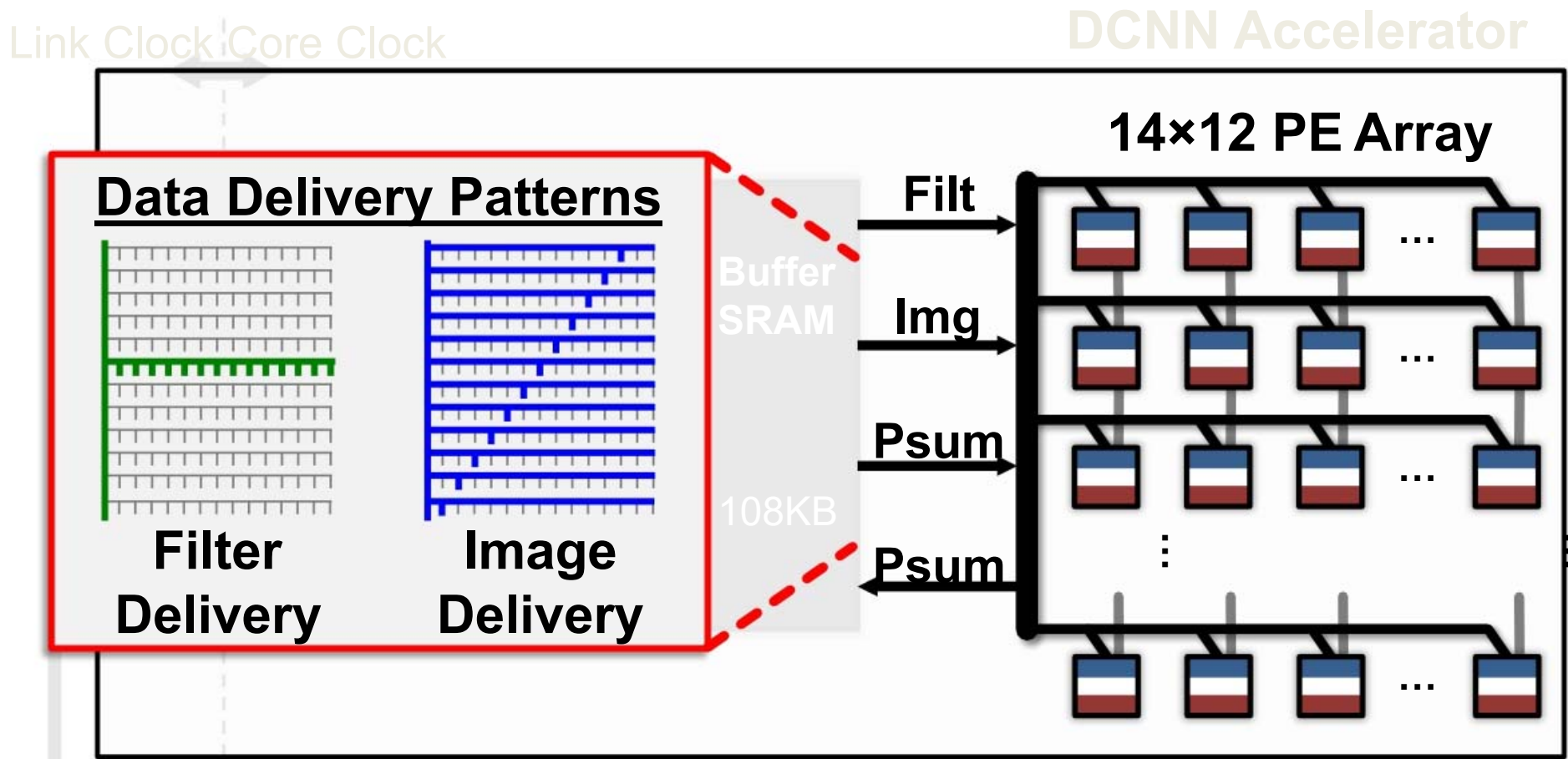
# Eyeriss DNN Accelerator



# Data Delivery with On-Chip Network



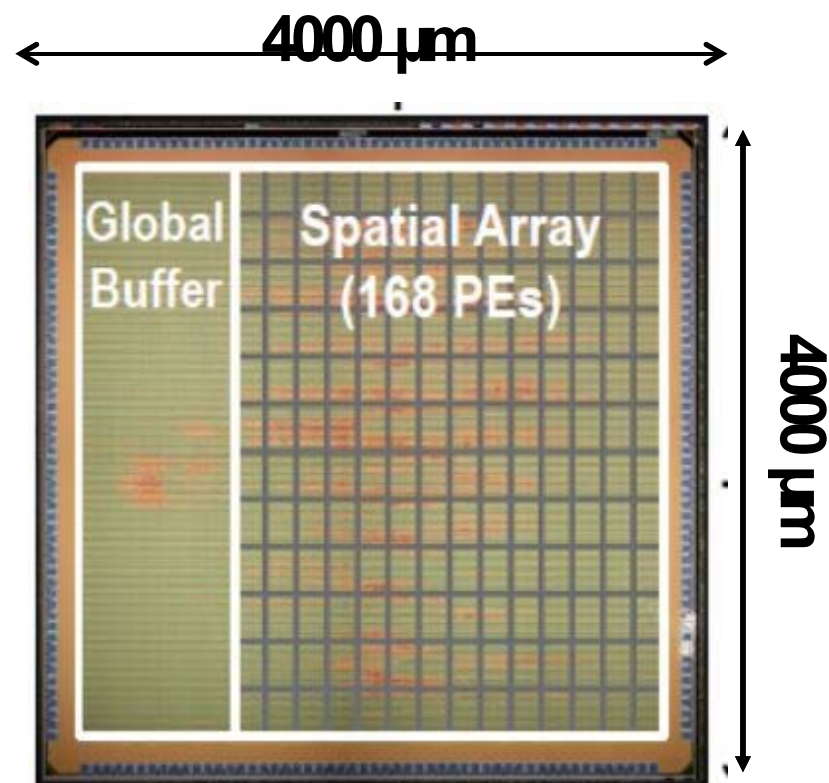
# Data Delivery with On-Chip Network



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

# 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**

---

# Backup Slides