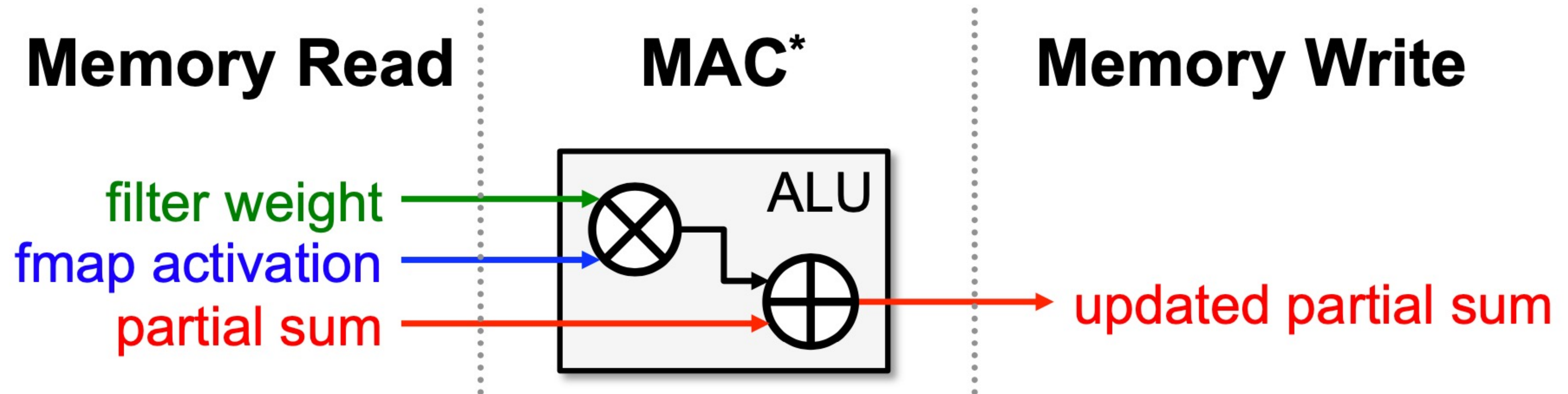


# Eyeriss

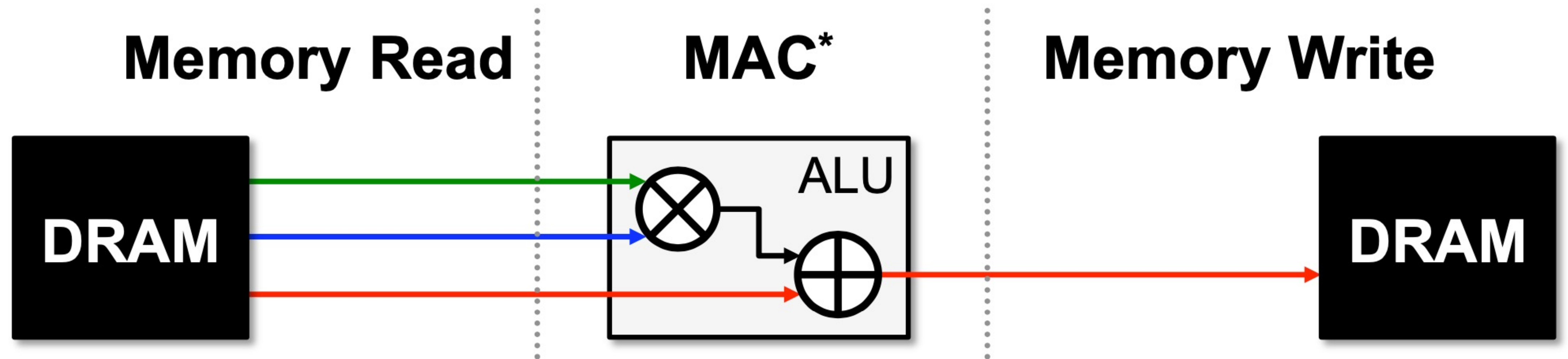
An Energy-Efficient Reconfigurable Accelerator for Deep Convolutional Neural Networks



# Memory Access is the Bottleneck

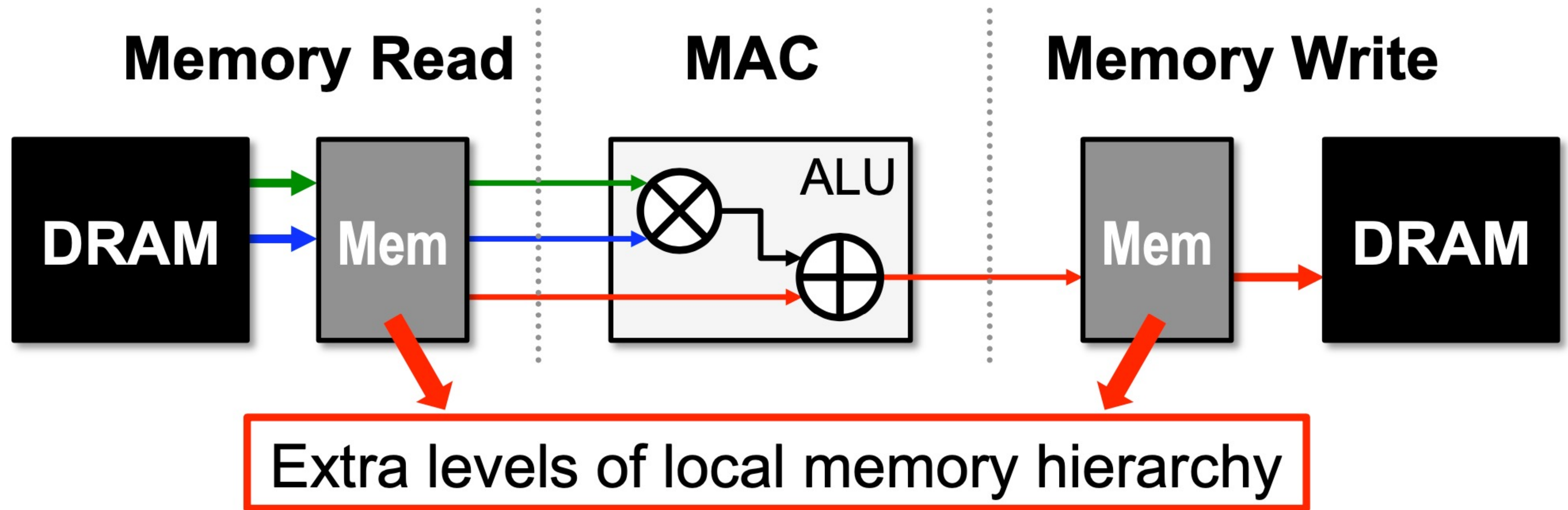






all memory R/W are **DRAM** accesses => **worst case**

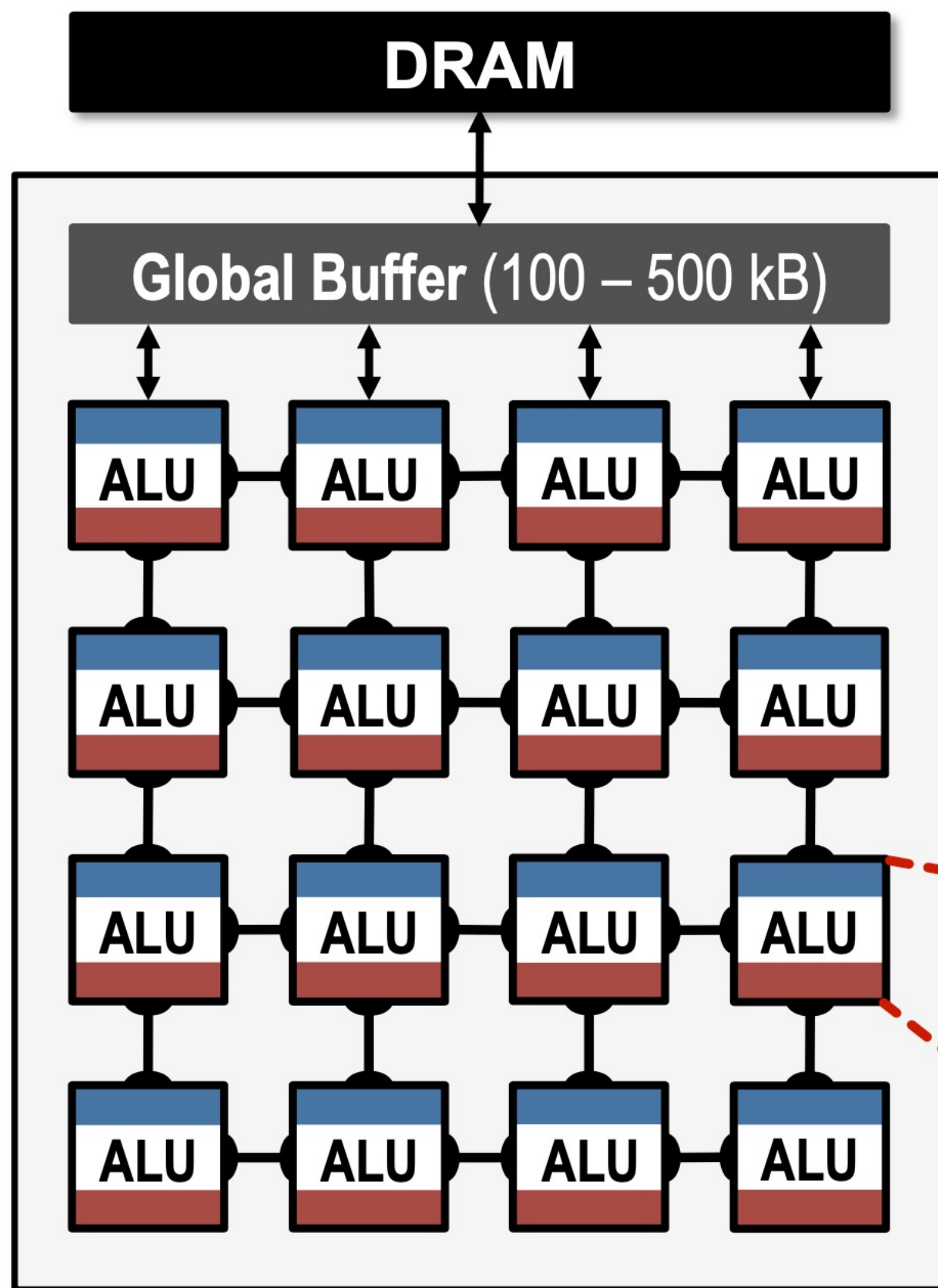




- data reuse
- local accumulation



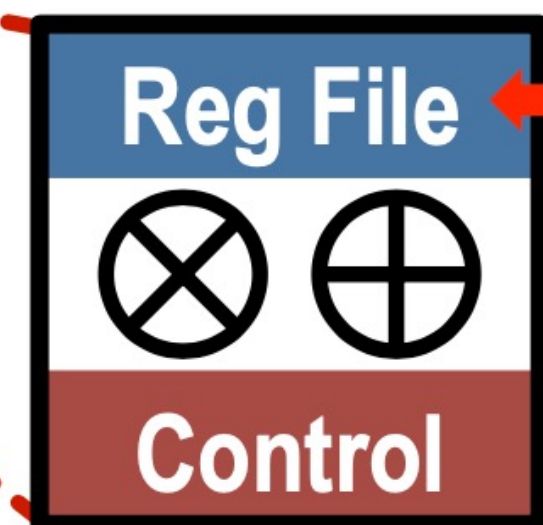
# Architecture



## Local Memory Hierarchy

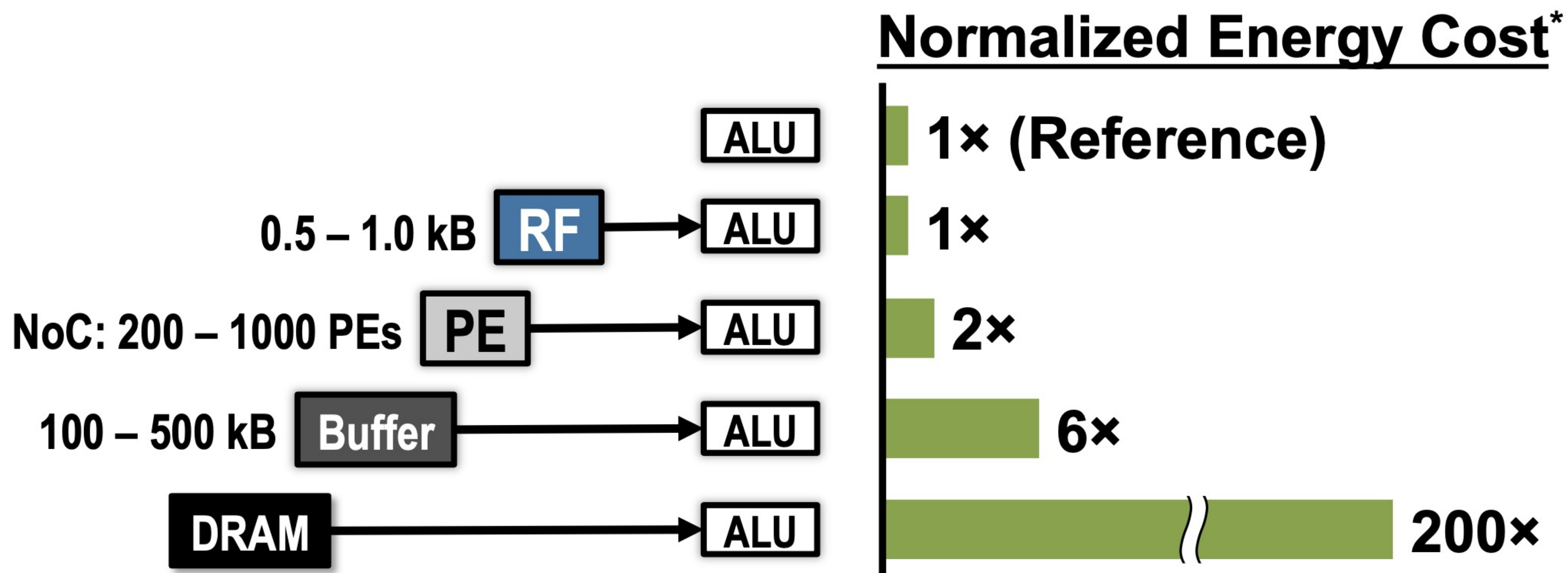
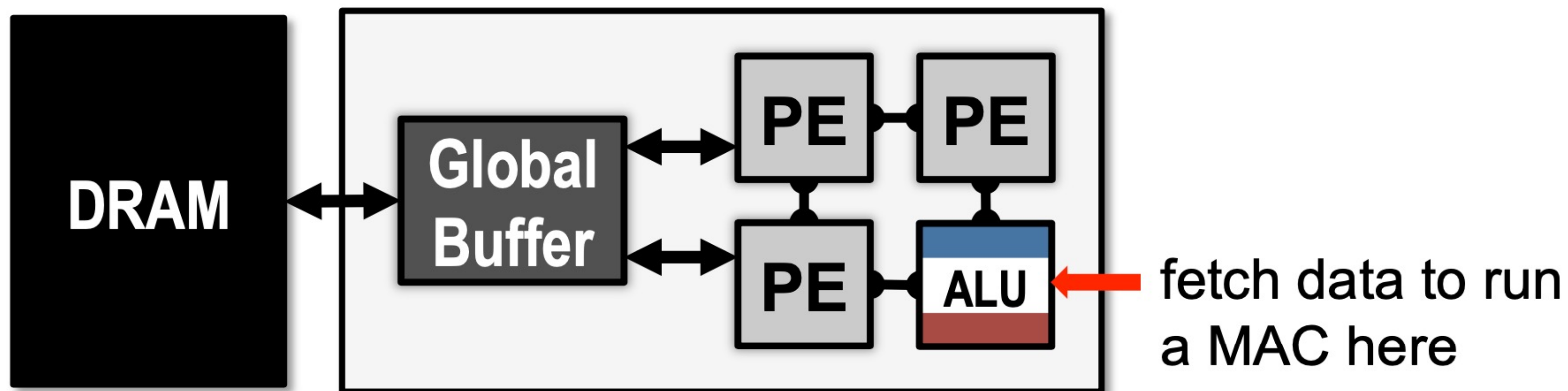
- Global Buffer
- Direct inter-PE network
- PE-local memory (RF)

## Processing Element (PE)



0.5 – 1.0 kB

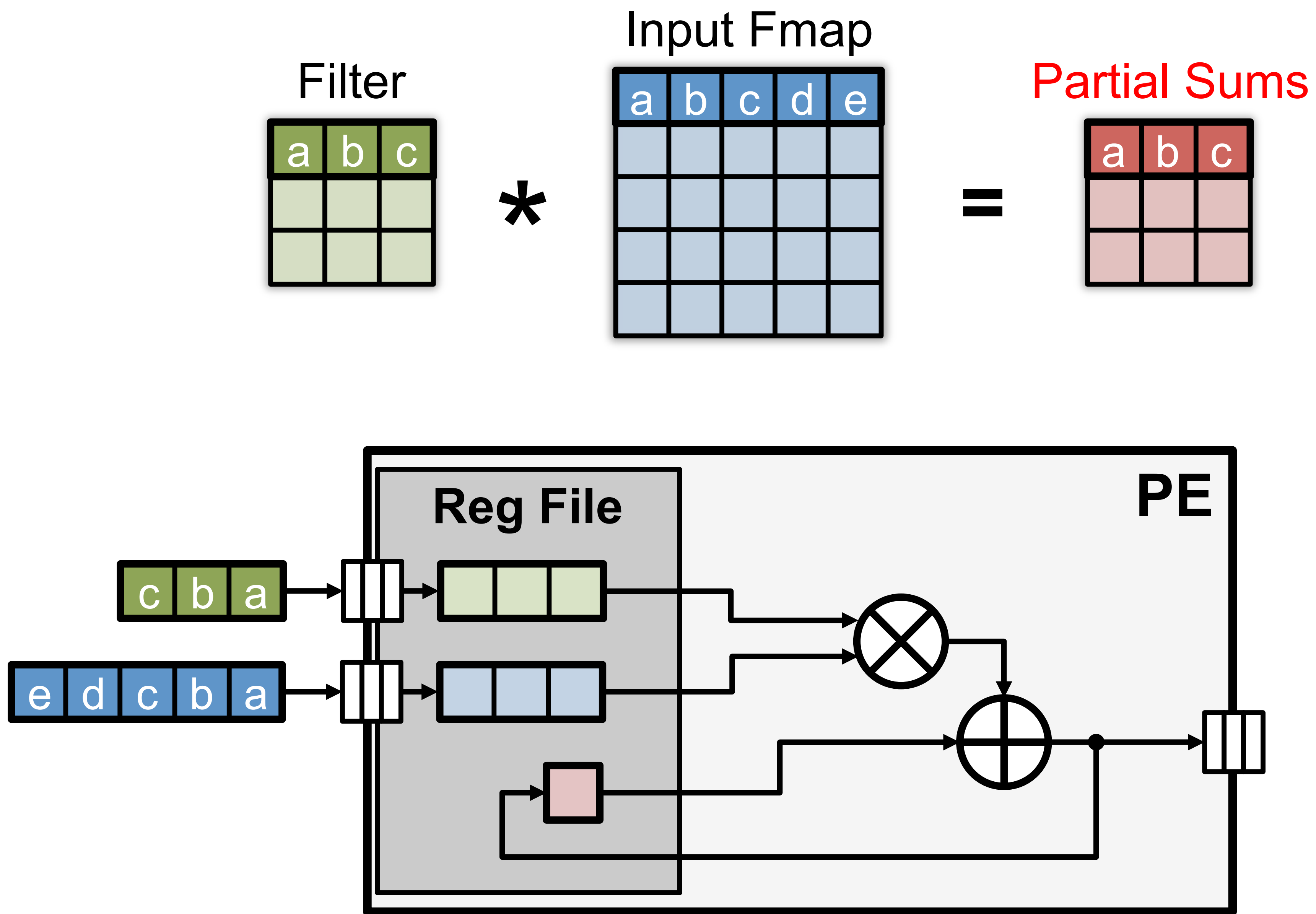




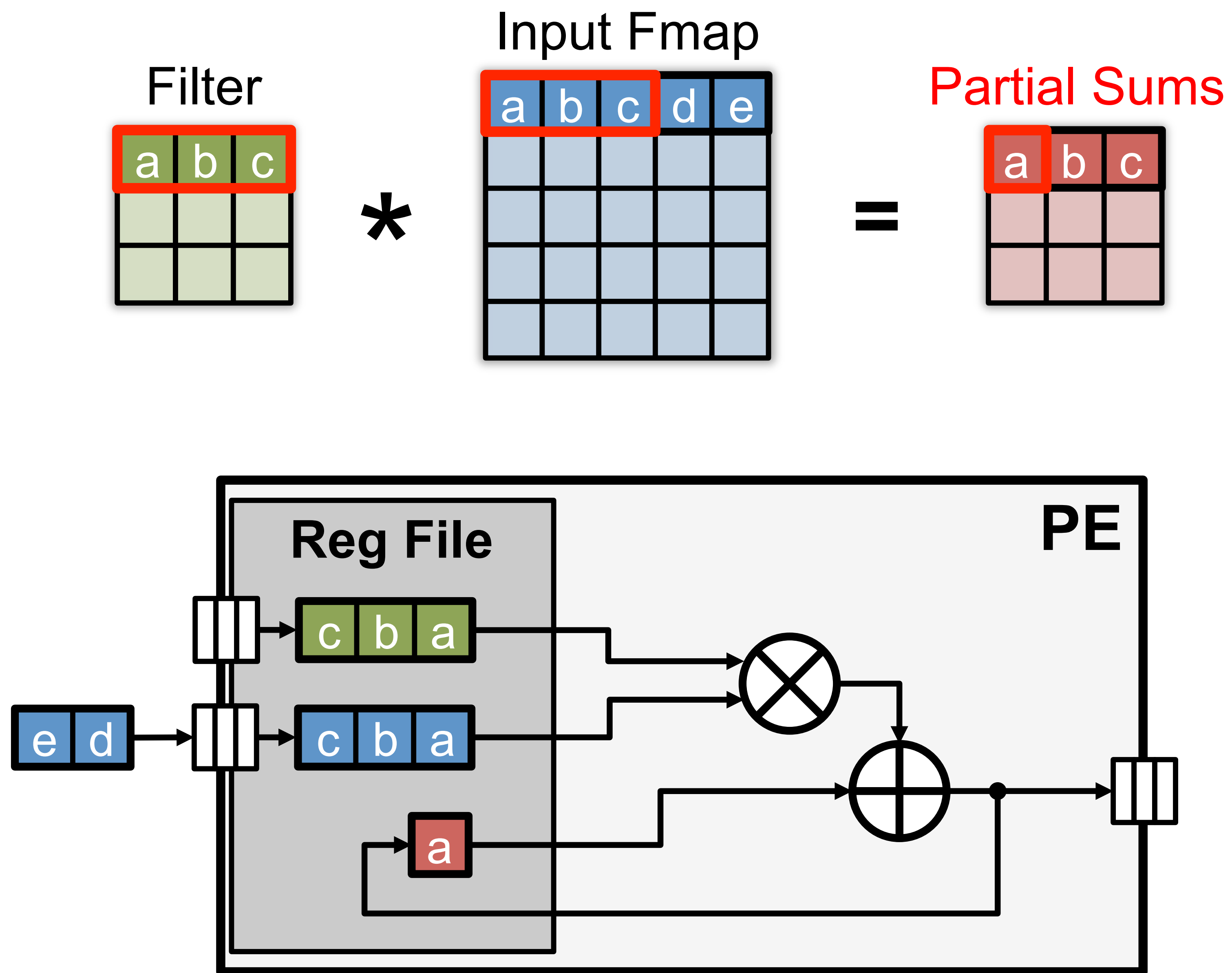


# Dataflow Taxonomy

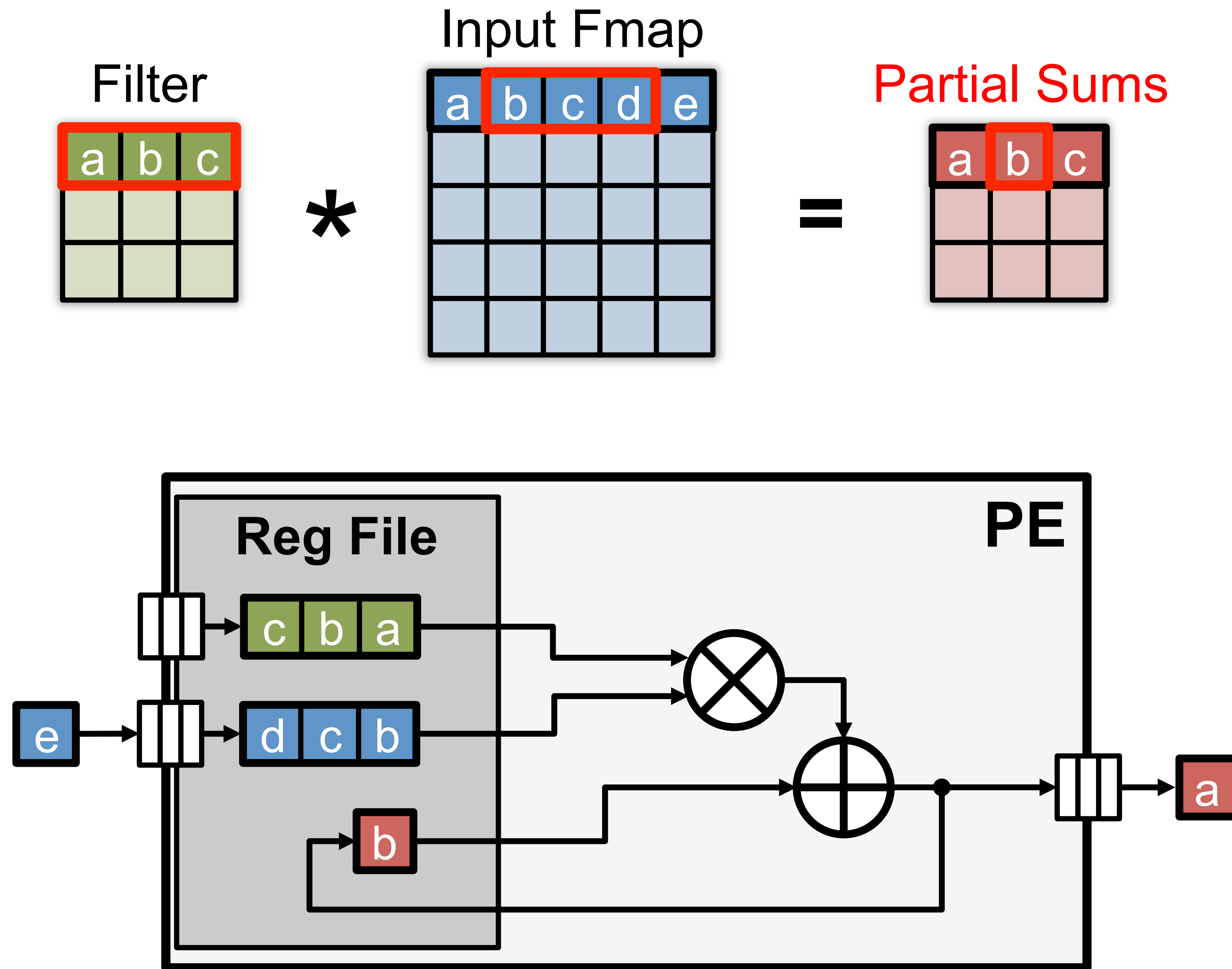
- Weight Stationary
- Output Stationary
- **Row Stationary**



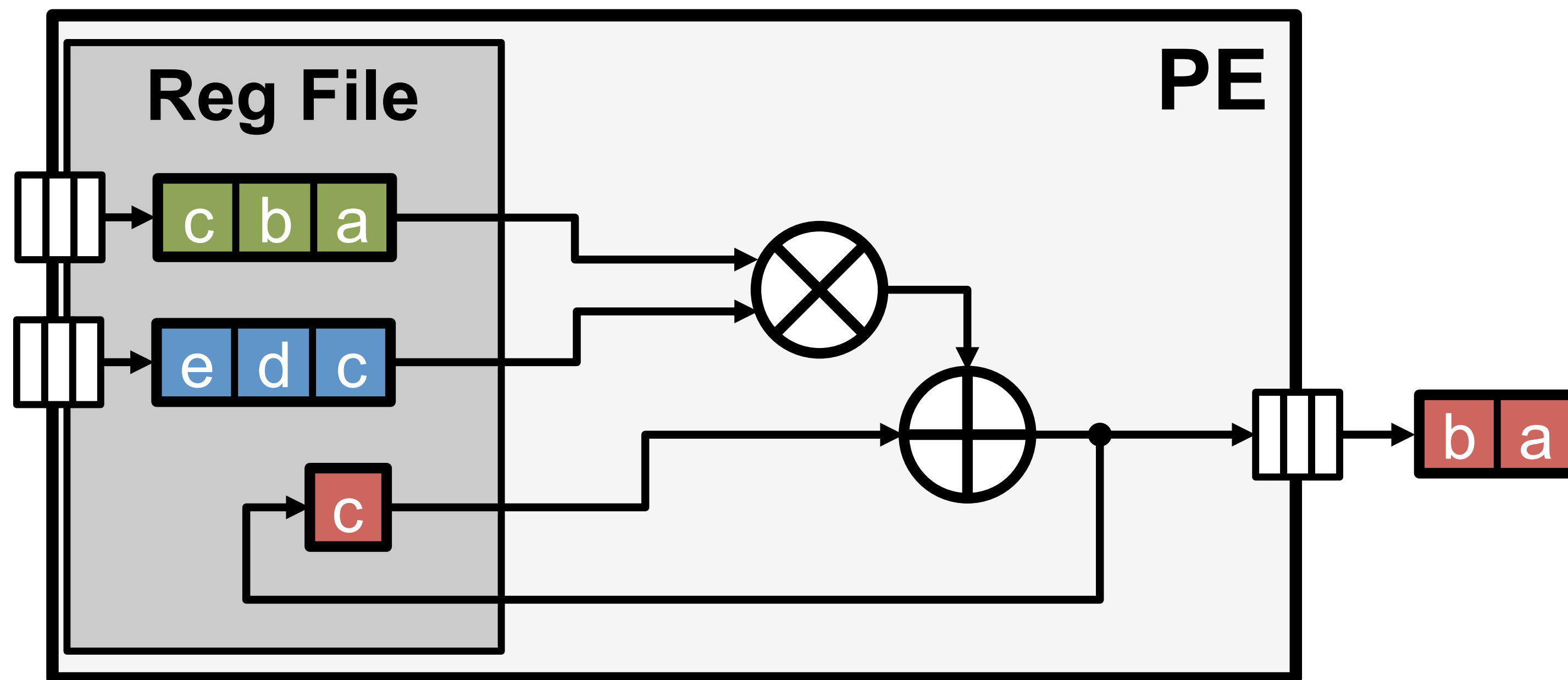
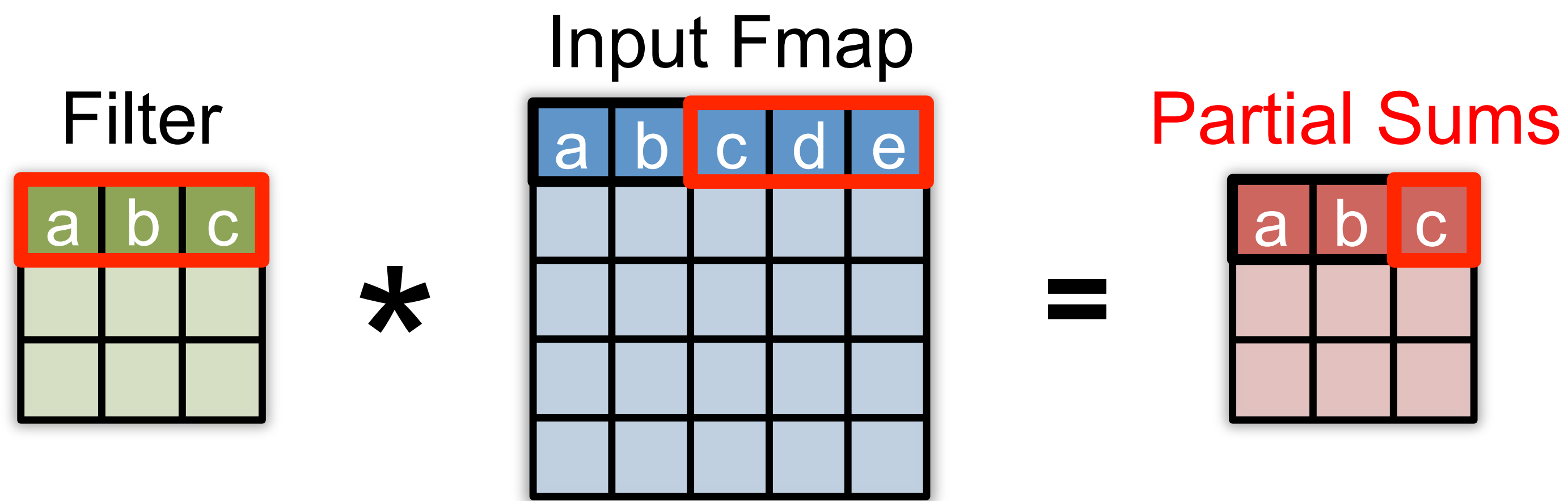




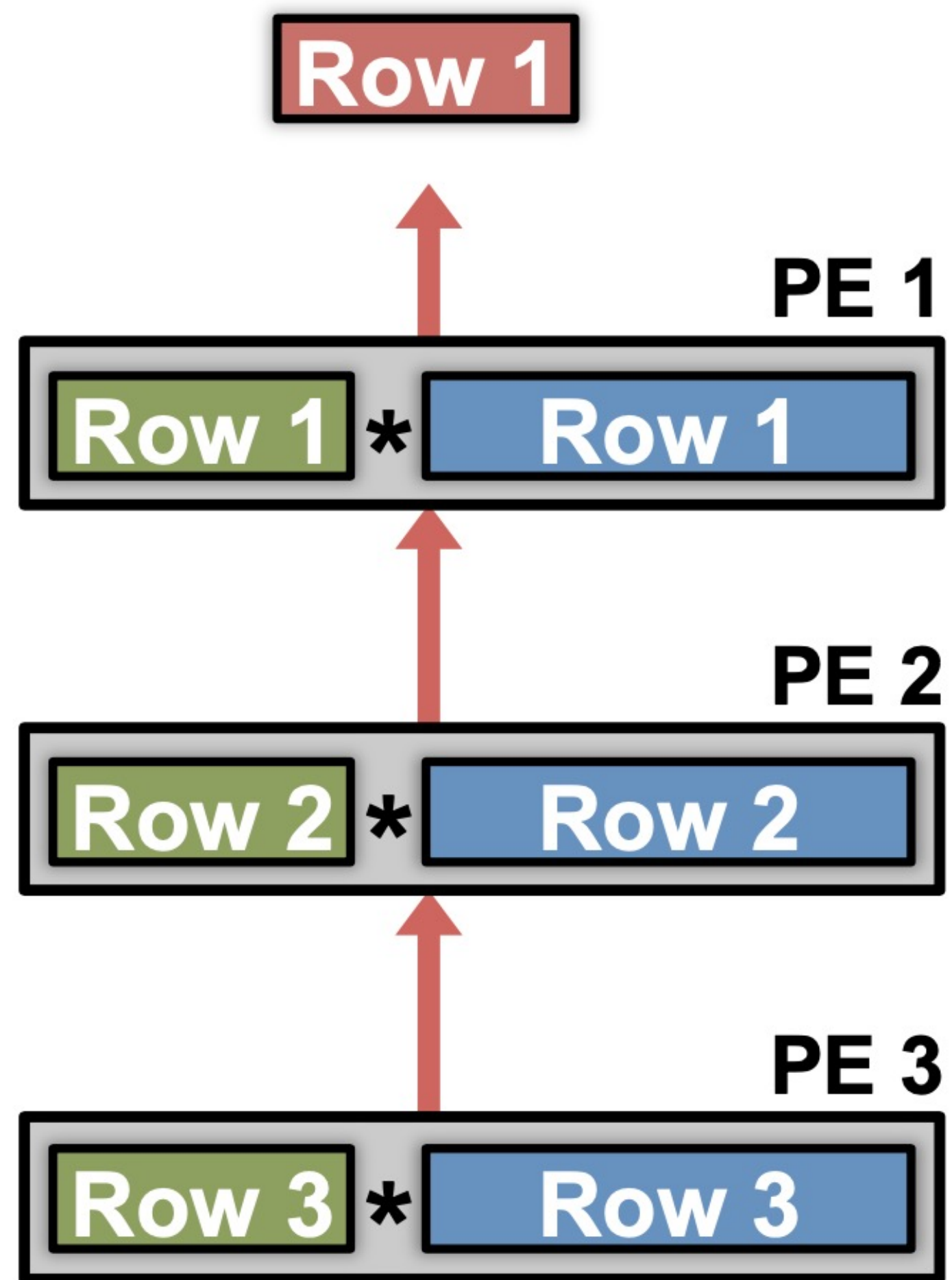






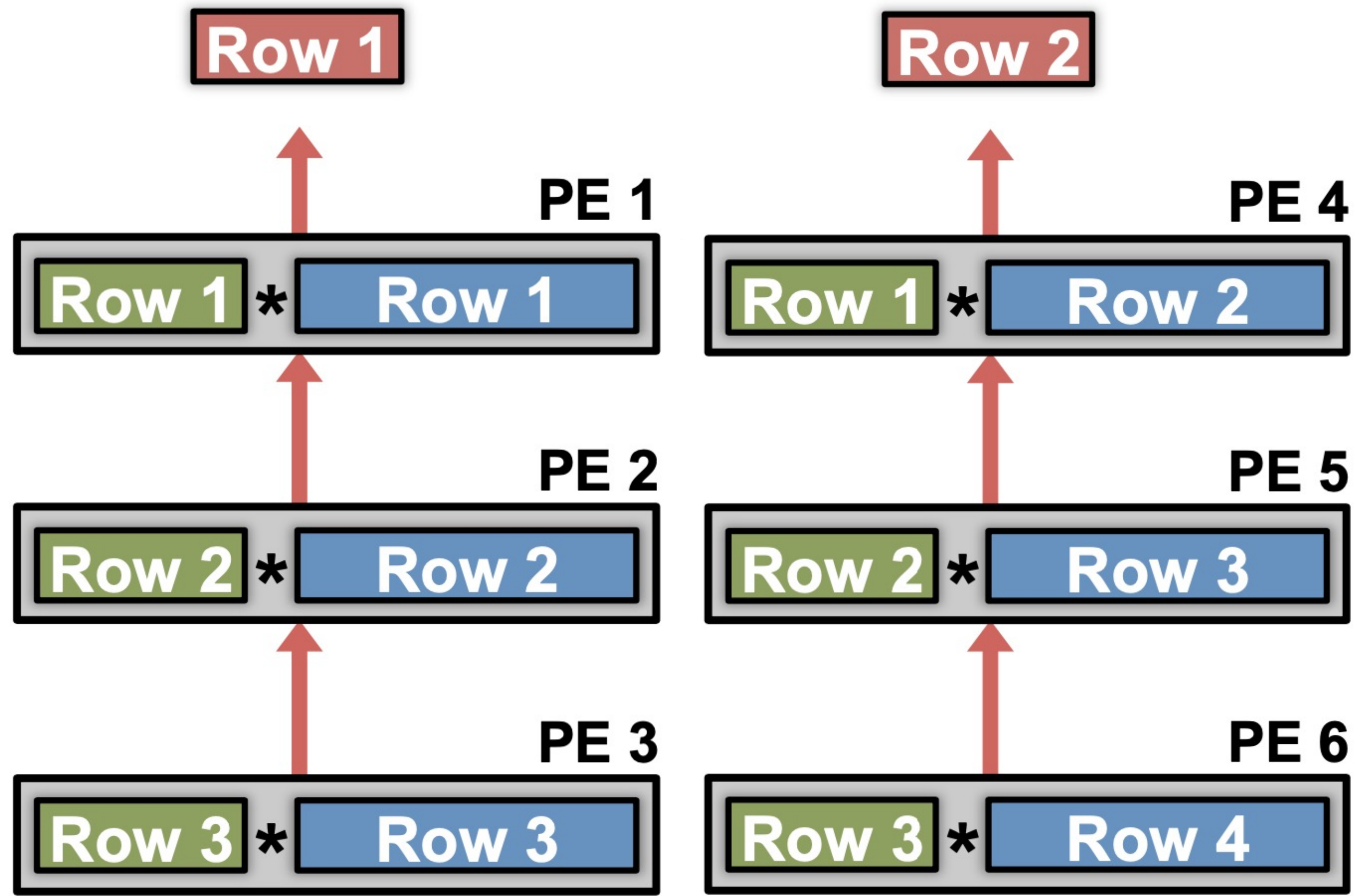






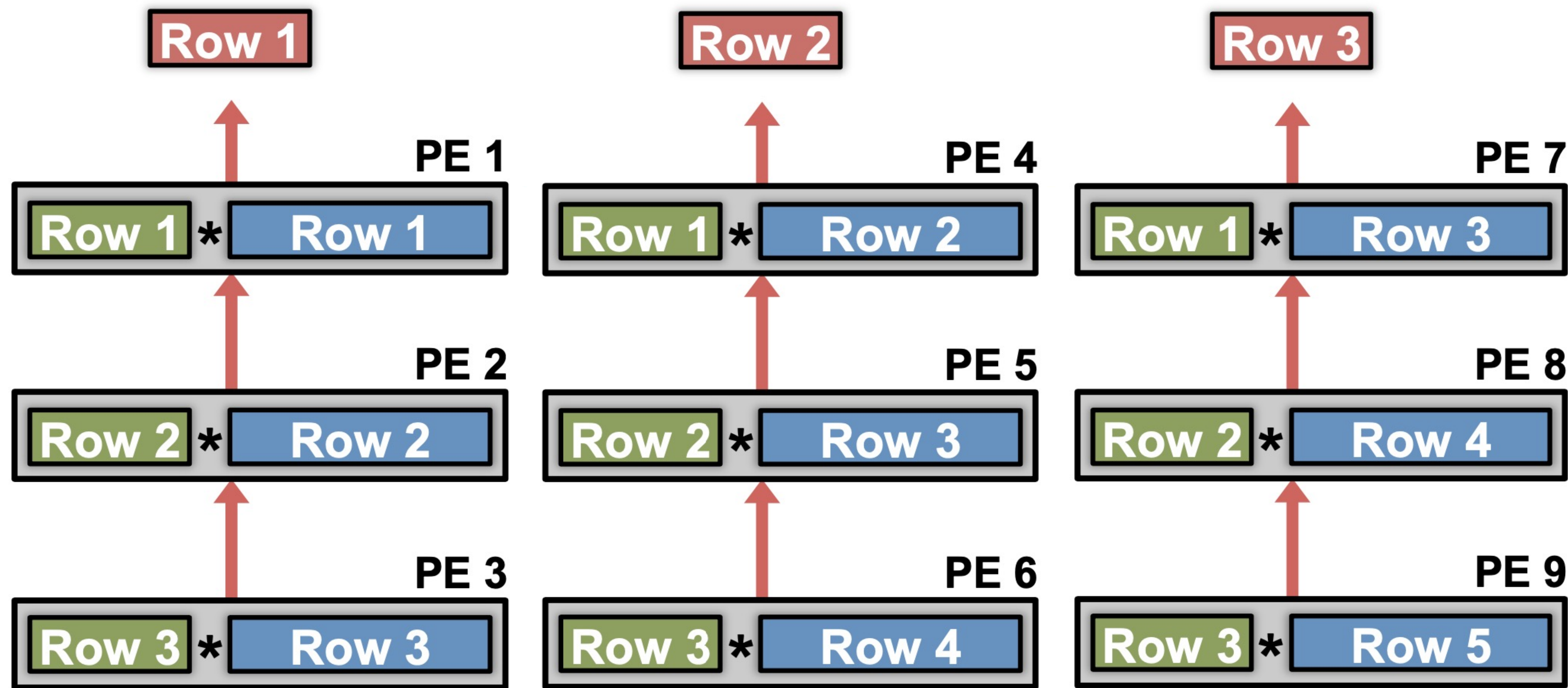
$$\begin{bmatrix} \text{Green} \\ \text{Green} \\ \text{Green} \end{bmatrix} * \begin{bmatrix} \text{Blue} & \text{Blue} & \text{Blue} & \text{Blue} \\ \text{Blue} & \text{Blue} & \text{Blue} & \text{Blue} \\ \text{Blue} & \text{Blue} & \text{Blue} & \text{Blue} \end{bmatrix} = \begin{bmatrix} \text{Pink} & \text{Pink} & \text{Pink} & \text{Pink} \\ \text{Pink} & \text{Pink} & \text{Pink} & \text{Pink} \\ \text{Pink} & \text{Pink} & \text{Pink} & \text{Pink} \end{bmatrix}$$





$$\begin{bmatrix} \text{Green} & \text{Green} & \text{Green} \\ \text{Green} & \text{Green} & \text{Green} \\ \text{Green} & \text{Green} & \text{Green} \end{bmatrix} * \begin{bmatrix} \text{Blue} & \text{Blue} & \text{Blue} & \text{Blue} \\ \text{Blue} & \text{Blue} & \text{Blue} & \text{Blue} \\ \text{Blue} & \text{Blue} & \text{Blue} & \text{Blue} \\ \text{Blue} & \text{Blue} & \text{Blue} & \text{Blue} \end{bmatrix} = \begin{bmatrix} \text{Red} & \text{Red} & \text{Red} \\ \text{Red} & \text{Red} & \text{Red} \\ \text{Red} & \text{Red} & \text{Red} \end{bmatrix}$$
$$\begin{bmatrix} \text{Green} & \text{Green} & \text{Green} \\ \text{Green} & \text{Green} & \text{Green} \\ \text{Green} & \text{Green} & \text{Green} \end{bmatrix} * \begin{bmatrix} \text{Blue} & \text{Blue} & \text{Blue} & \text{Blue} \\ \text{Blue} & \text{Blue} & \text{Blue} & \text{Blue} \\ \text{Blue} & \text{Blue} & \text{Blue} & \text{Blue} \\ \text{Blue} & \text{Blue} & \text{Blue} & \text{Blue} \end{bmatrix} = \begin{bmatrix} \text{Red} & \text{Red} & \text{Red} \\ \text{Red} & \text{Red} & \text{Red} \\ \text{Red} & \text{Red} & \text{Red} \end{bmatrix}$$



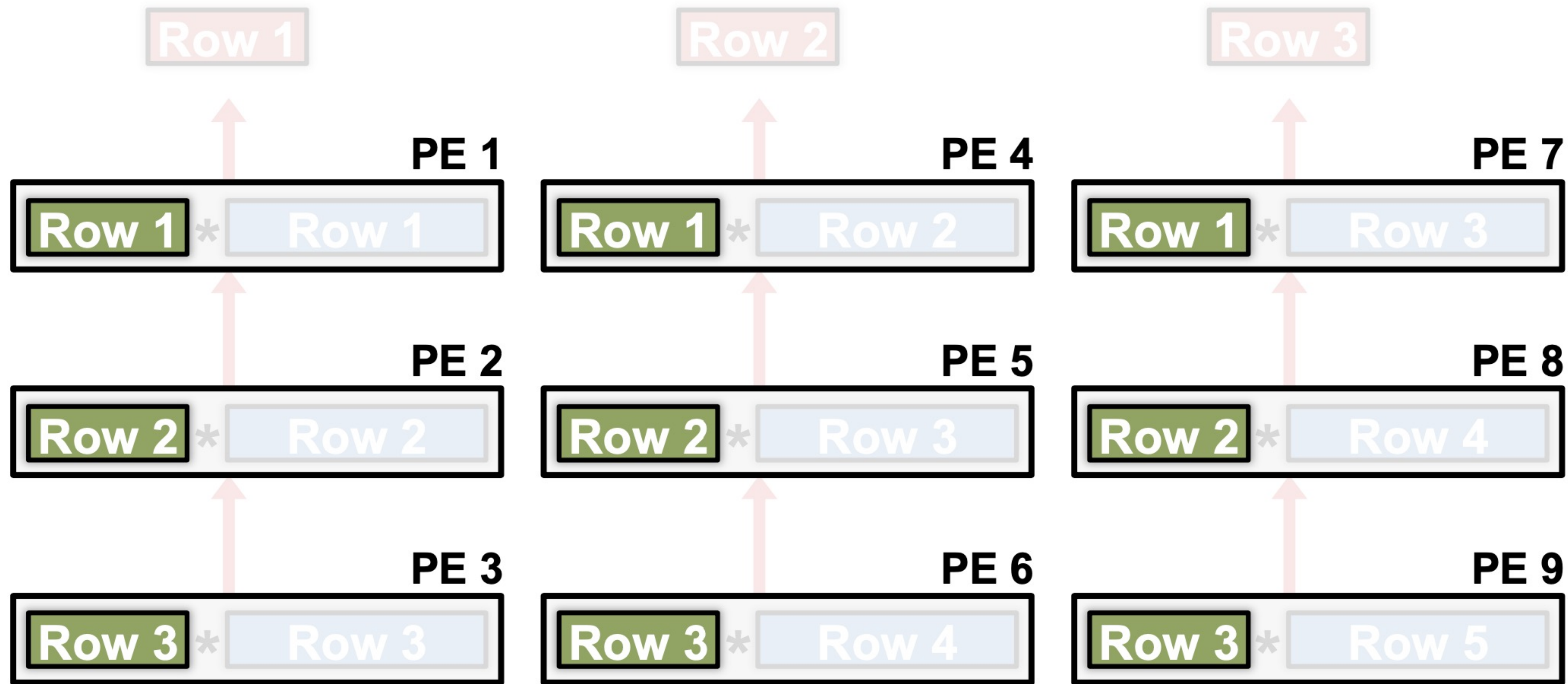


$$\begin{bmatrix} \text{Row 1} \\ \text{Row 2} \\ \text{Row 3} \end{bmatrix} * \begin{bmatrix} \text{Row 1} & \text{Row 2} & \text{Row 3} & \text{Row 4} & \text{Row 5} \end{bmatrix} = \begin{bmatrix} \text{Row 1} & \text{Row 2} & \text{Row 3} & \text{Row 4} & \text{Row 5} \end{bmatrix}$$

$$\begin{bmatrix} \text{Row 1} \\ \text{Row 2} \\ \text{Row 3} \end{bmatrix} * \begin{bmatrix} \text{Row 1} & \text{Row 2} & \text{Row 3} & \text{Row 4} & \text{Row 5} \end{bmatrix} = \begin{bmatrix} \text{Row 1} & \text{Row 2} & \text{Row 3} & \text{Row 4} & \text{Row 5} \end{bmatrix}$$

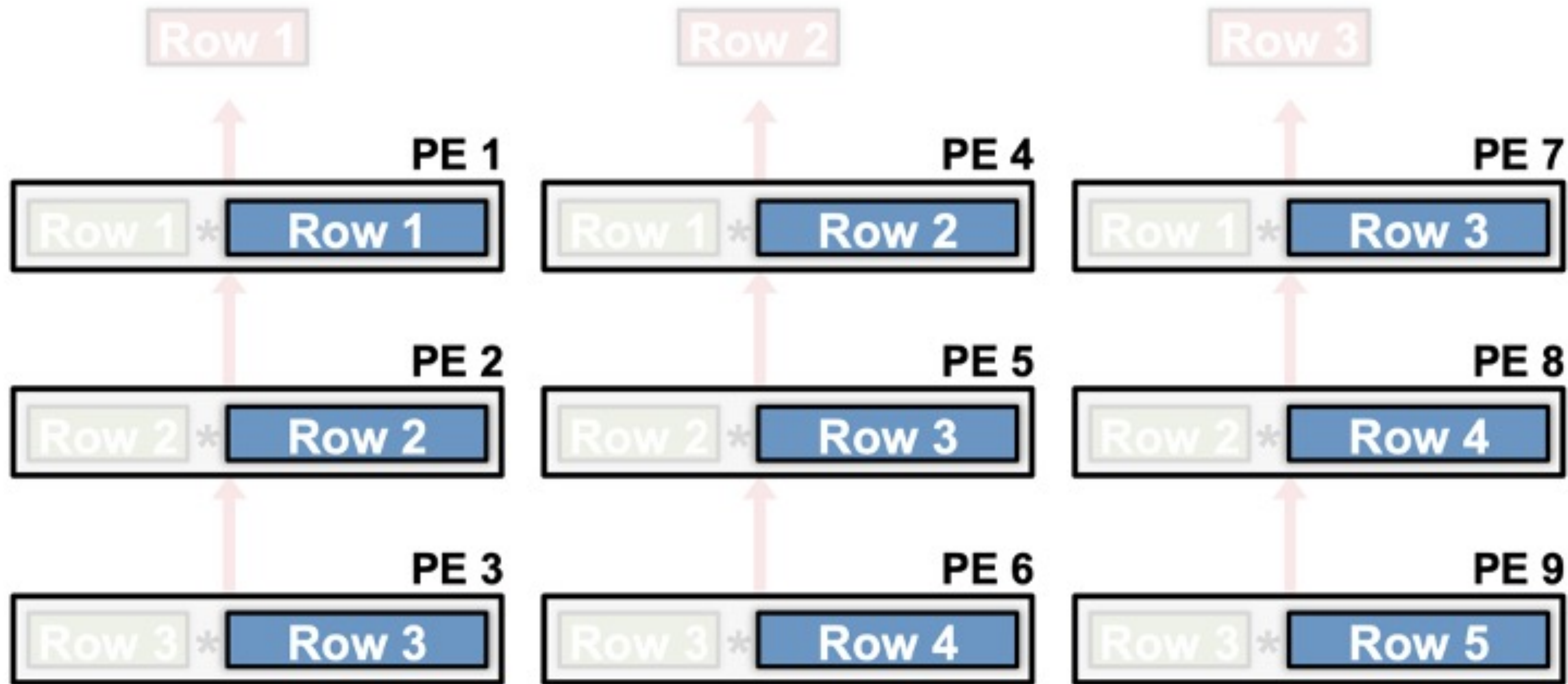
$$\begin{bmatrix} \text{Row 1} \\ \text{Row 2} \\ \text{Row 3} \end{bmatrix} * \begin{bmatrix} \text{Row 1} & \text{Row 2} & \text{Row 3} & \text{Row 4} & \text{Row 5} \end{bmatrix} = \begin{bmatrix} \text{Row 1} & \text{Row 2} & \text{Row 3} & \text{Row 4} & \text{Row 5} \end{bmatrix}$$





**Filter rows** are reused across PEs horizontally

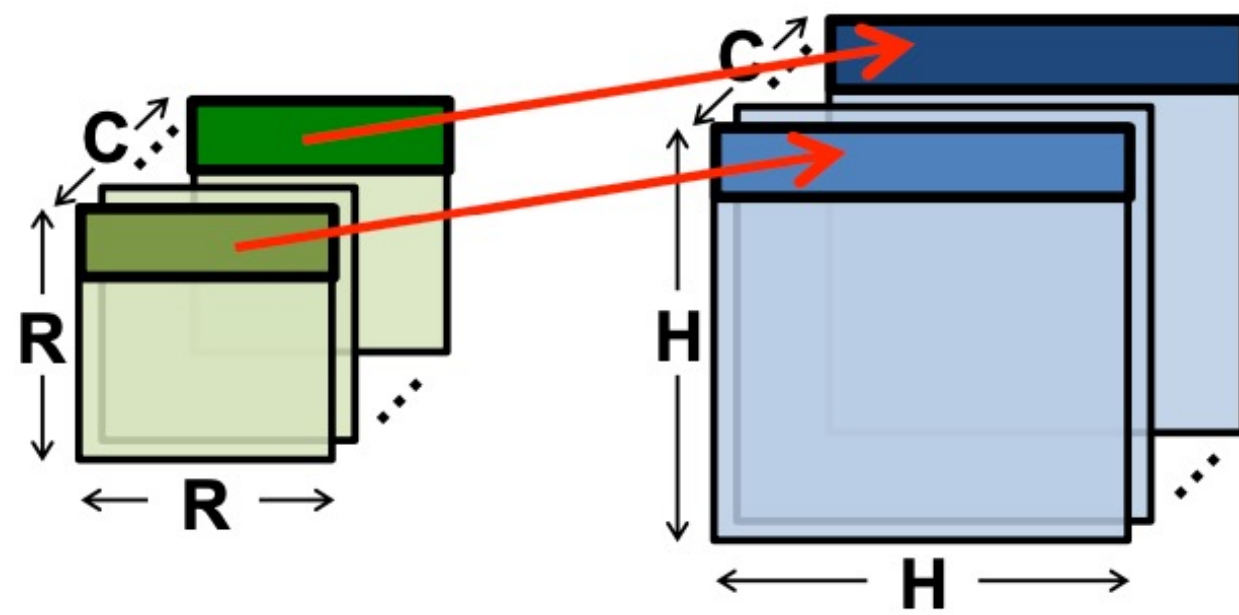




**Fmap rows** are reused across PEs **diagonally**



# multiple channels



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

accumulate psums

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



# Run-Length Coding (RLC)

Input: 0, 0, 12, 0, 0, 0, 0, 53, 0, 0, 22, ...

*Run Level Run Level Run Level Term*

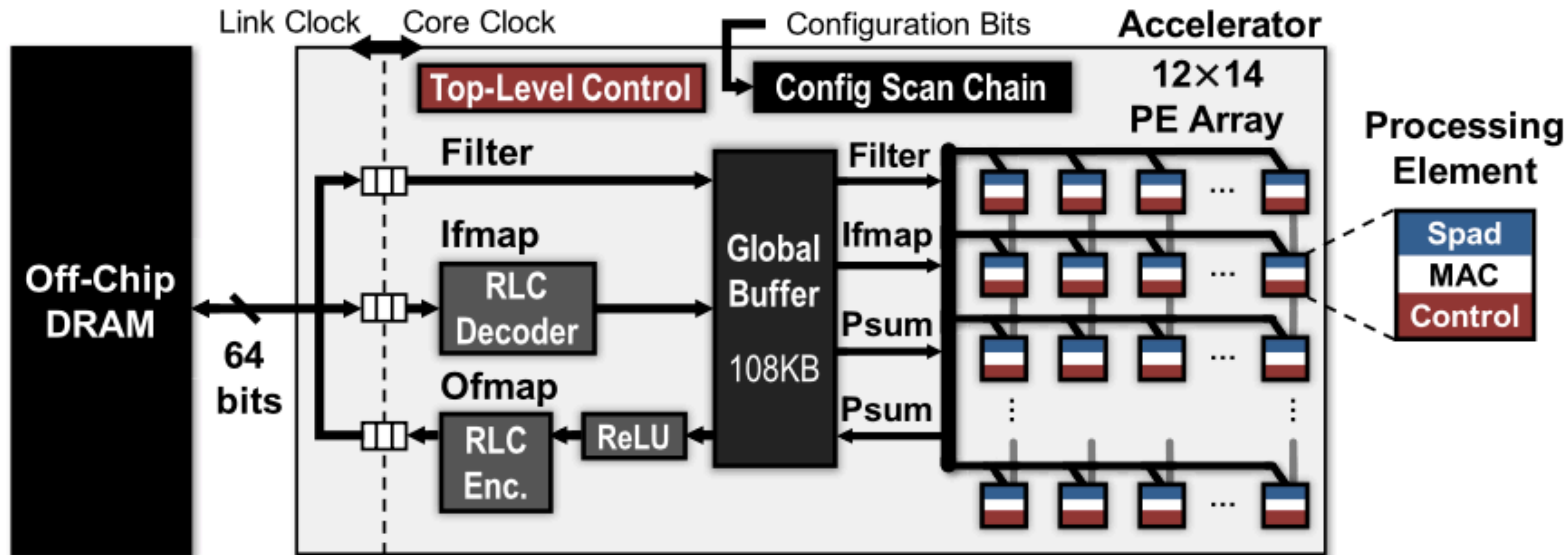
Output (64b):

2	12	4	53	2	22	0
---	----	---	----	---	----	---

5b 16b 5b 16b 5b 16b 1b



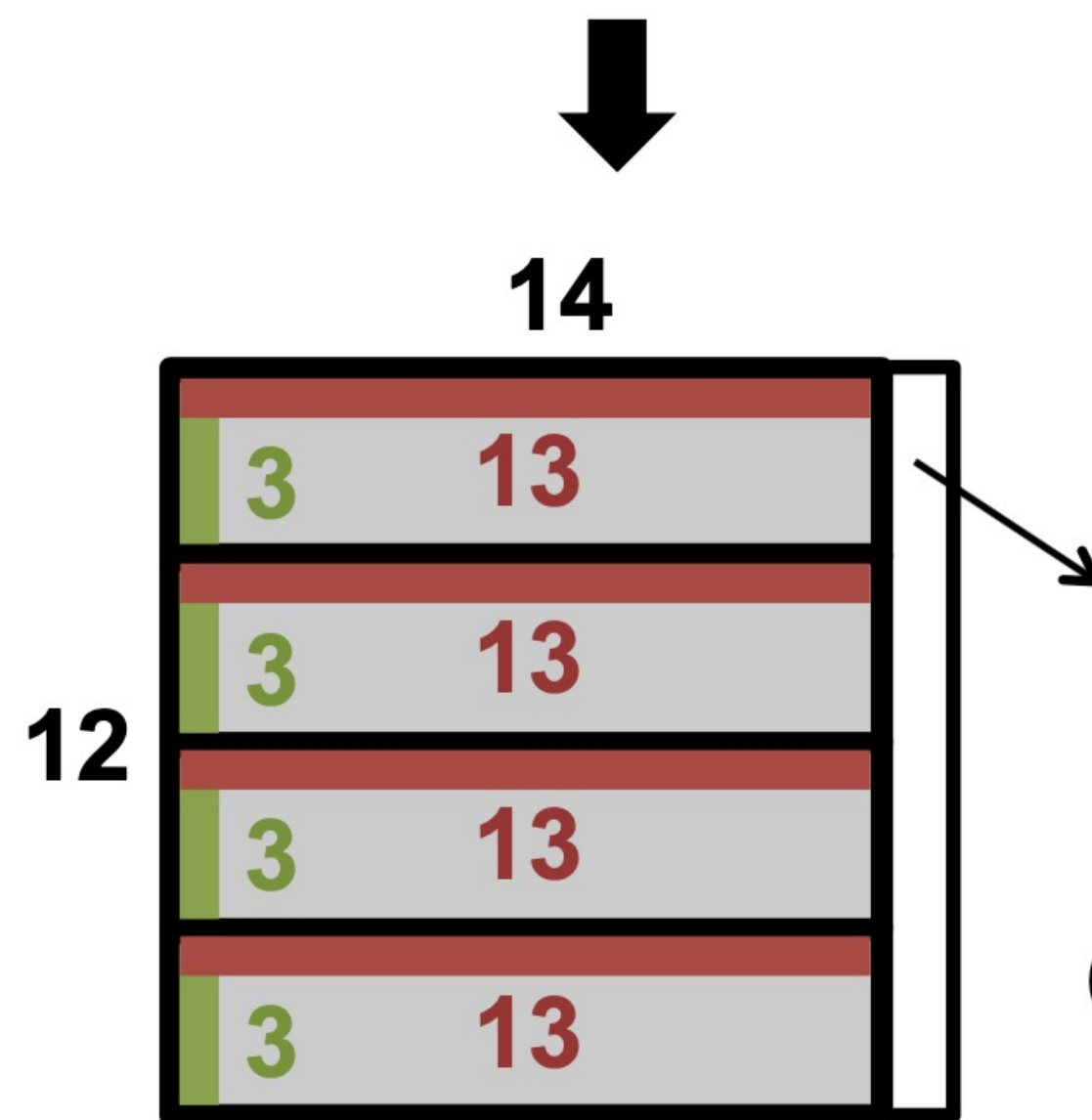
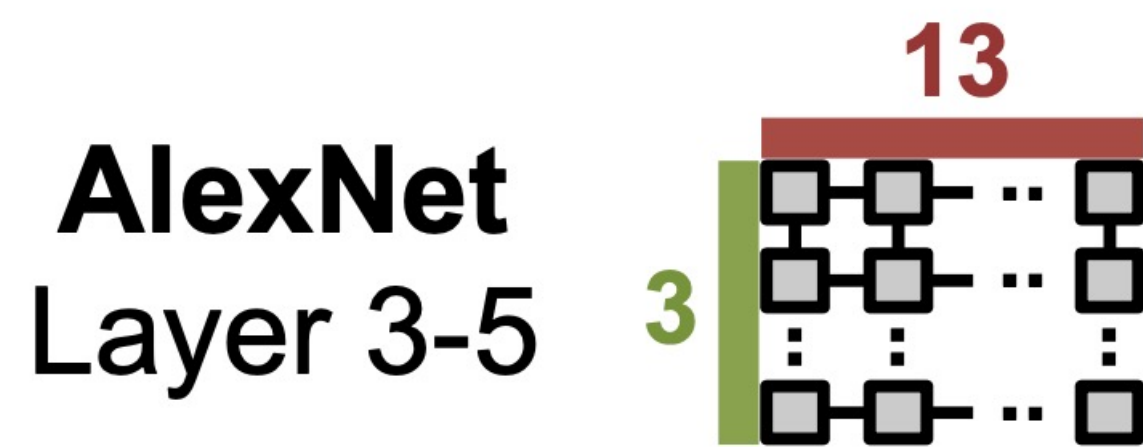
# Eyeriss system architecture





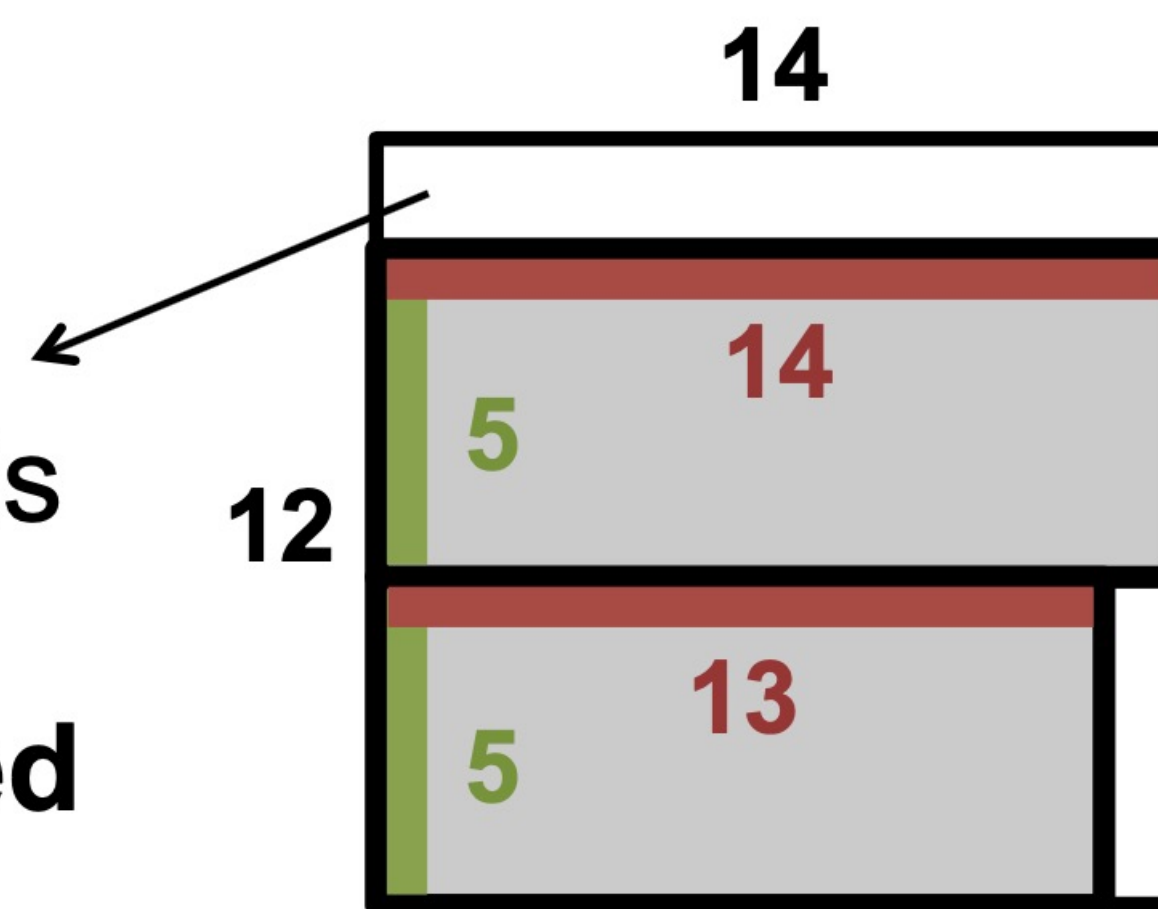
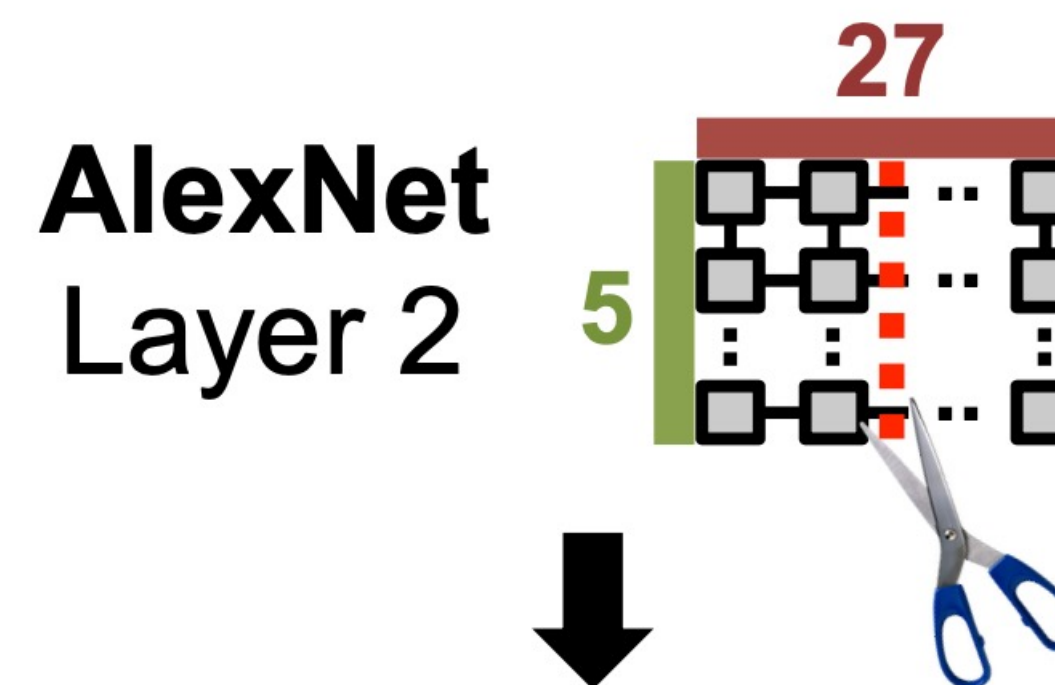
# Logical to Physical Mappings

## Replication



Physical PE Array

## Folding



Physical PE Array

Unused PEs  
are  
Clock Gated