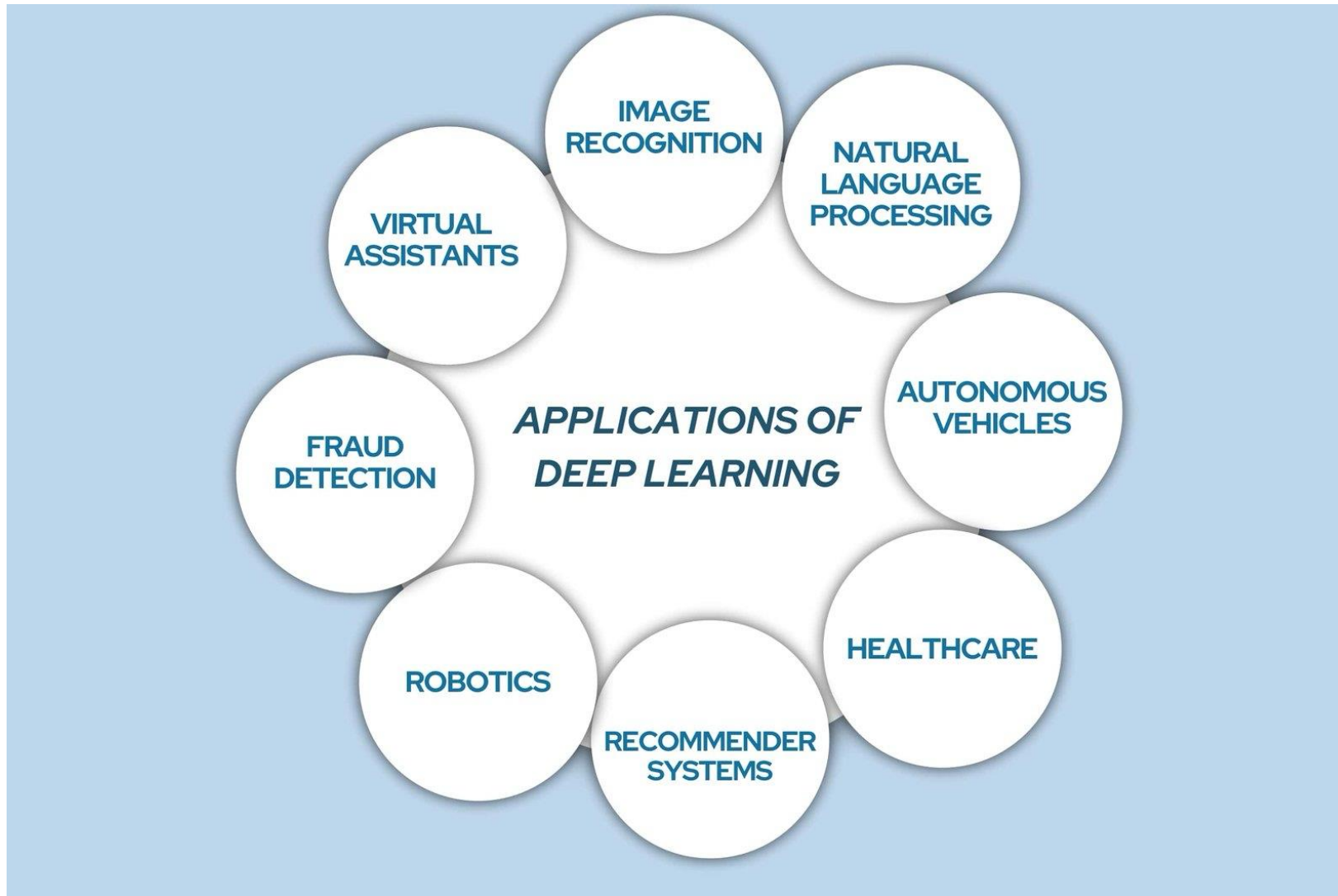


# Performance and Power Modeling of ML Accelerators

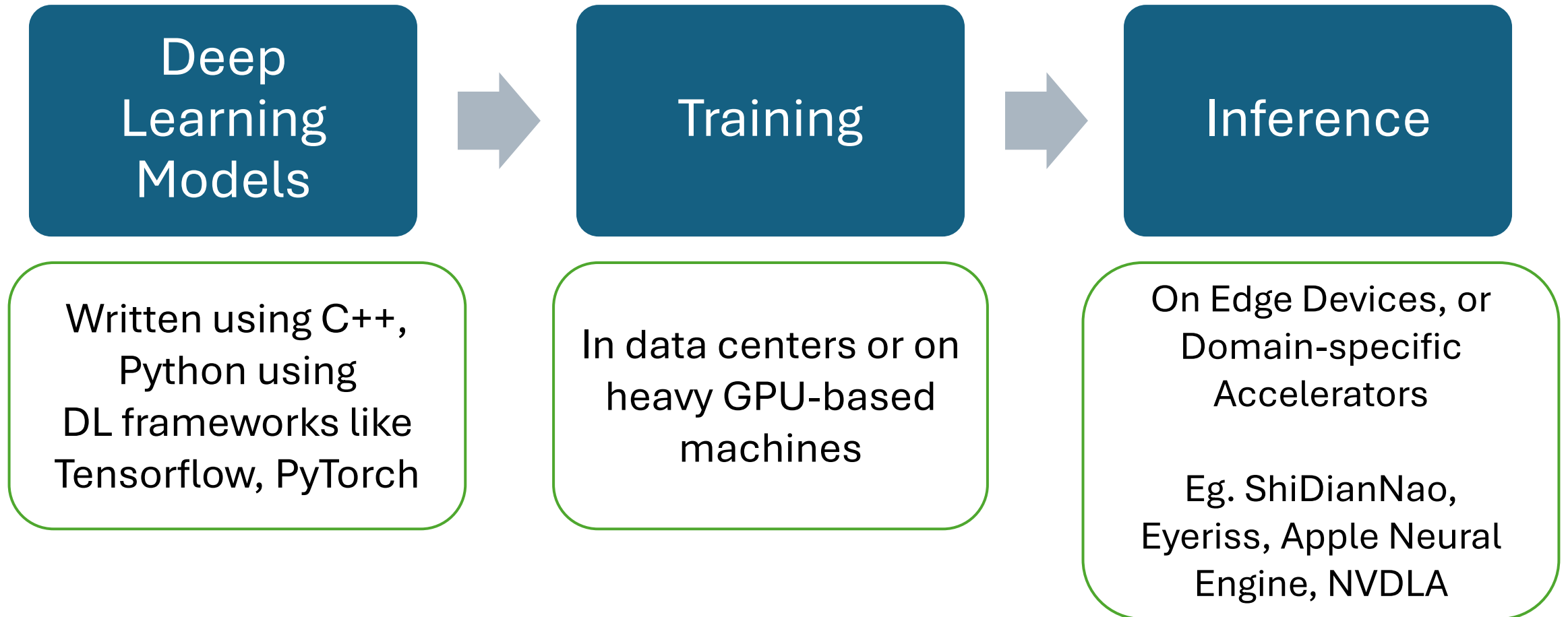
Ayushi Agarwal

Course: Synthesis of Digital Systems (COL719)

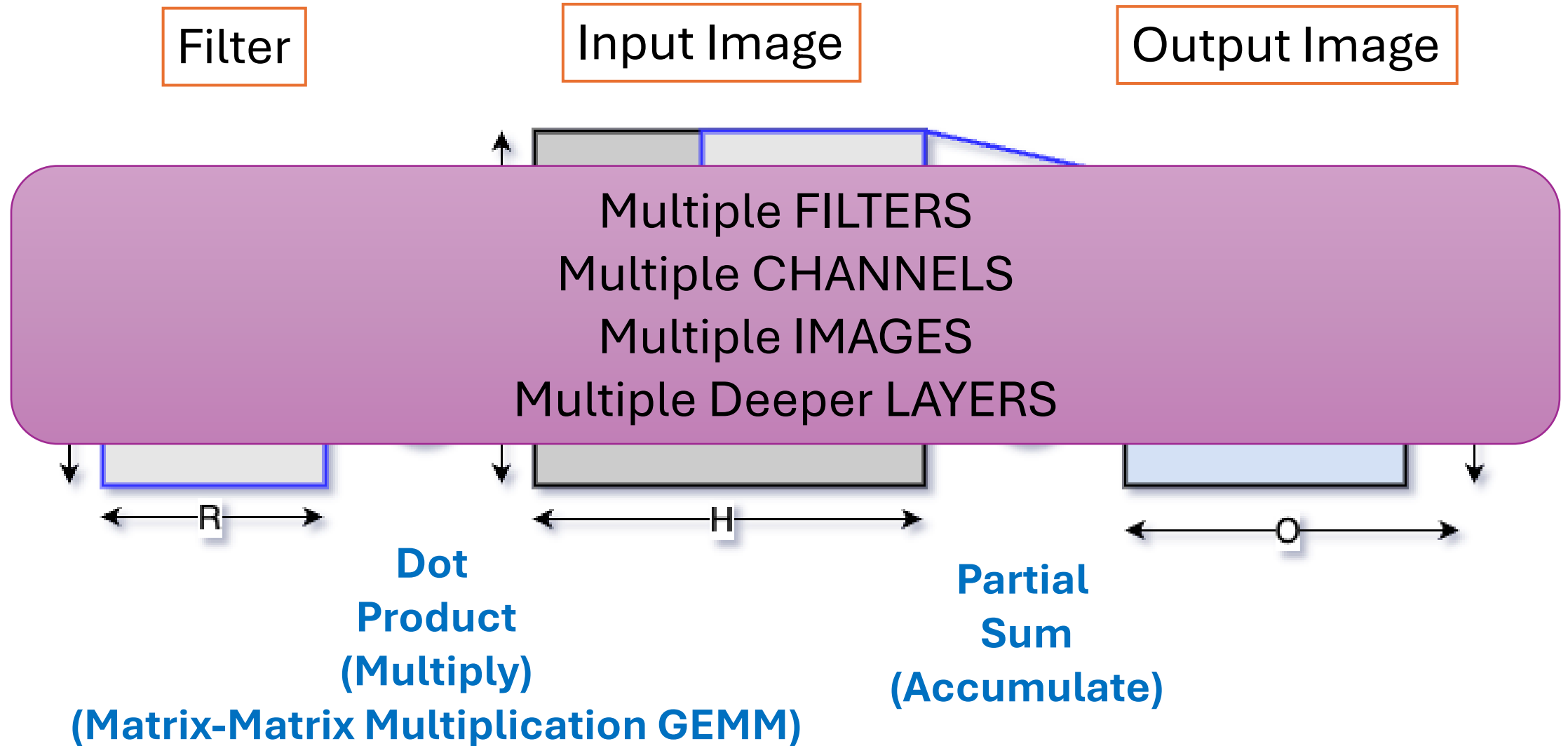
# Deep Learning Applications



# Deep Learning Application Pipeline



# Deep Convolutional Neural Networks



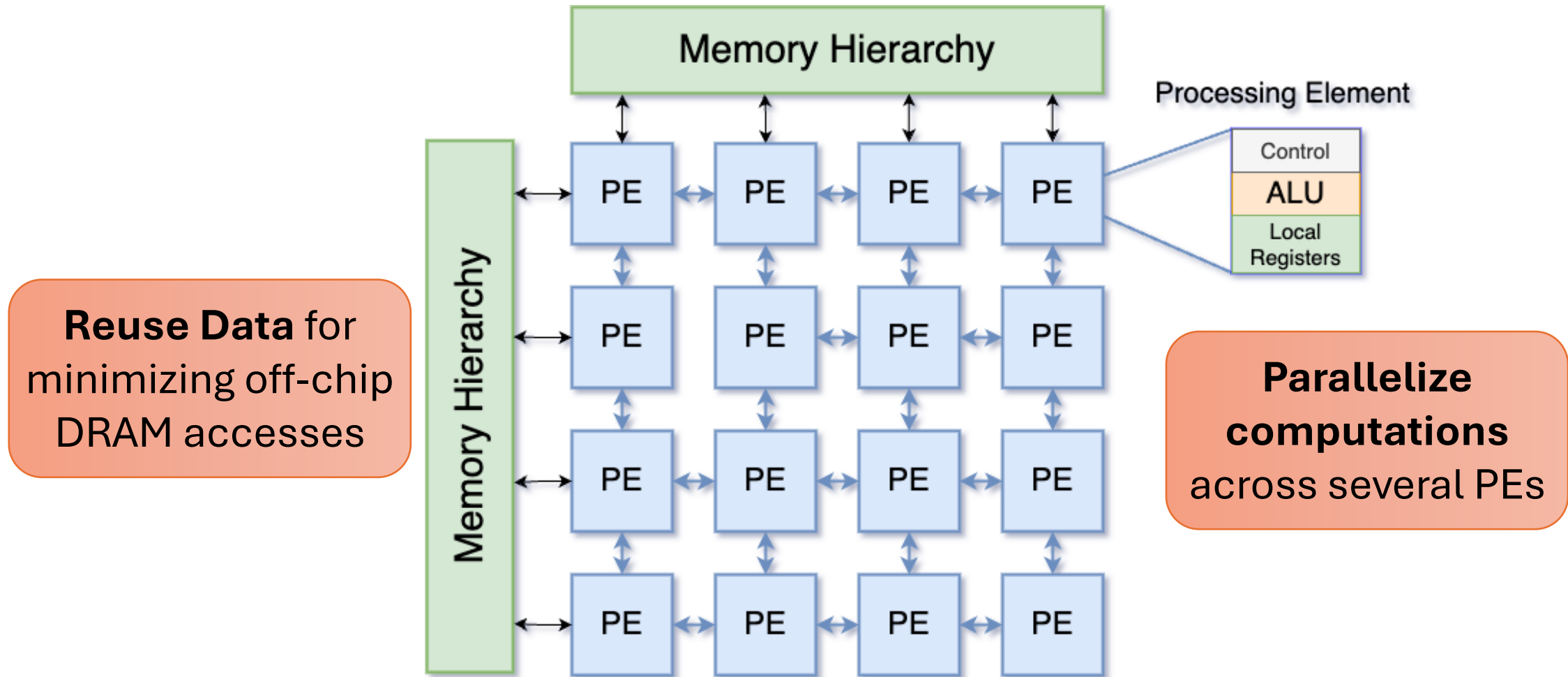
# Challenges with DNN computations

- Large number of parameters in modern deep networks
  - Require many parallel computations → Not suitable for CPU execution
  - Require many data transfers from the memory to compute units → GPU computation is energy-inefficient

DNN Model	Number of Parameters (weights)
AlexNet (2012)	3.98M
VGGNet-16 (2014)	28.25M
GoogLeNet (2015)	6.77M
ResNet-50 (2016)	23M
DLRM (2019)	540M

# Domain-Specific Accelerators

Example: Eyeriss, Google TPU, NVDLA, ShiDianNao



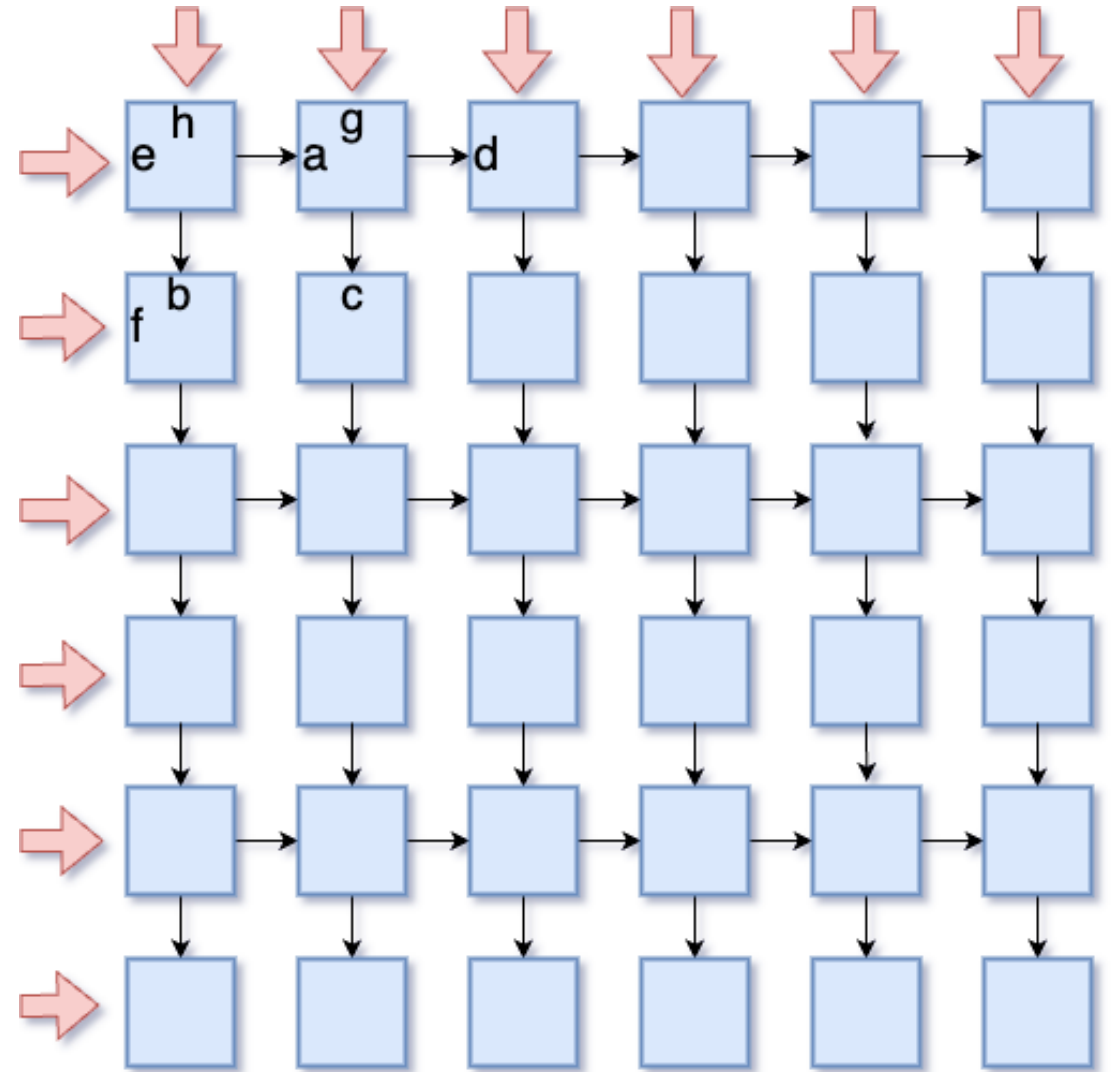
# Systolic Array Architecture

- Designed for Efficient parallel computation
- Composed of identical PEs
- Simpler data routing and communication btw PEs
- Pipelined Dataflow

High Parallelism

Ease of Implementation

Efficient Data Reuse



# Model GEMM using Systolic Array

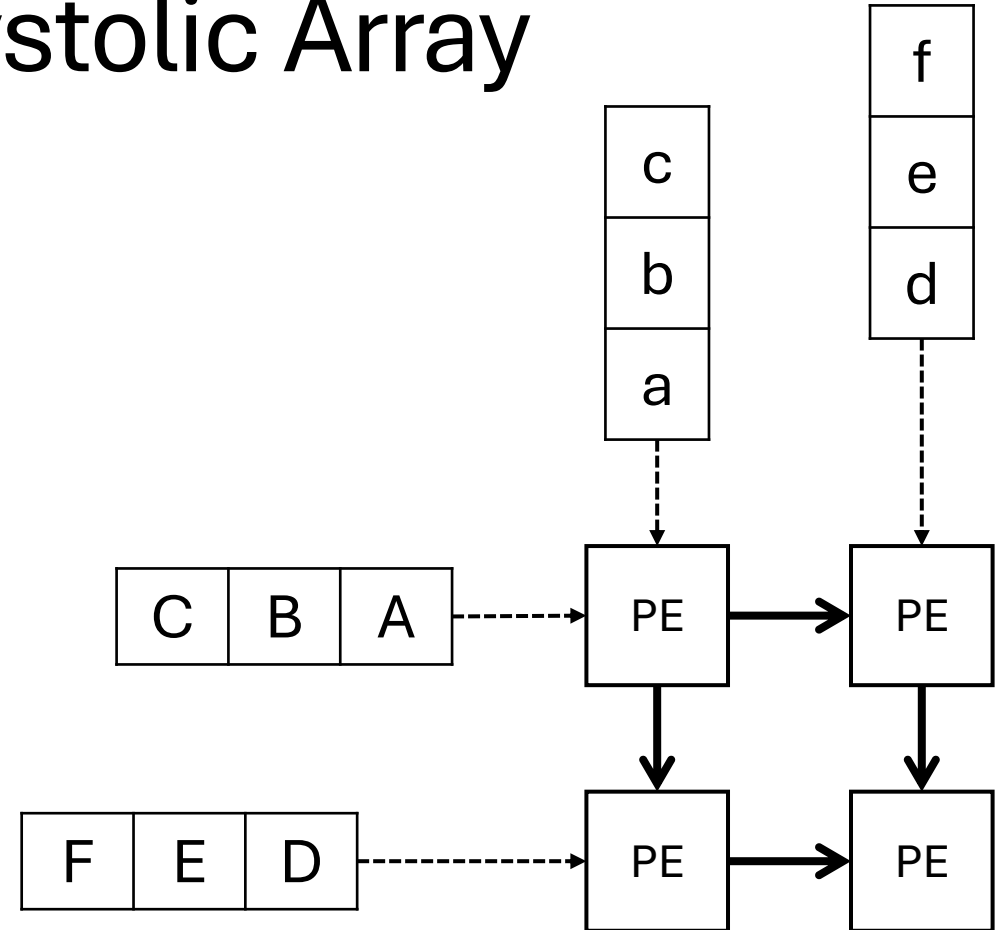
A	B	C
D	E	F

 $\times$ 

a	d
b	e
c	f

  
 $=$ 

$Aa + Bb + Cc$	$Ad + Be + Cf$
$Da + Eb + Fc$	$Dd + Ee + Ff$





# Model GEMM using Systolic Array

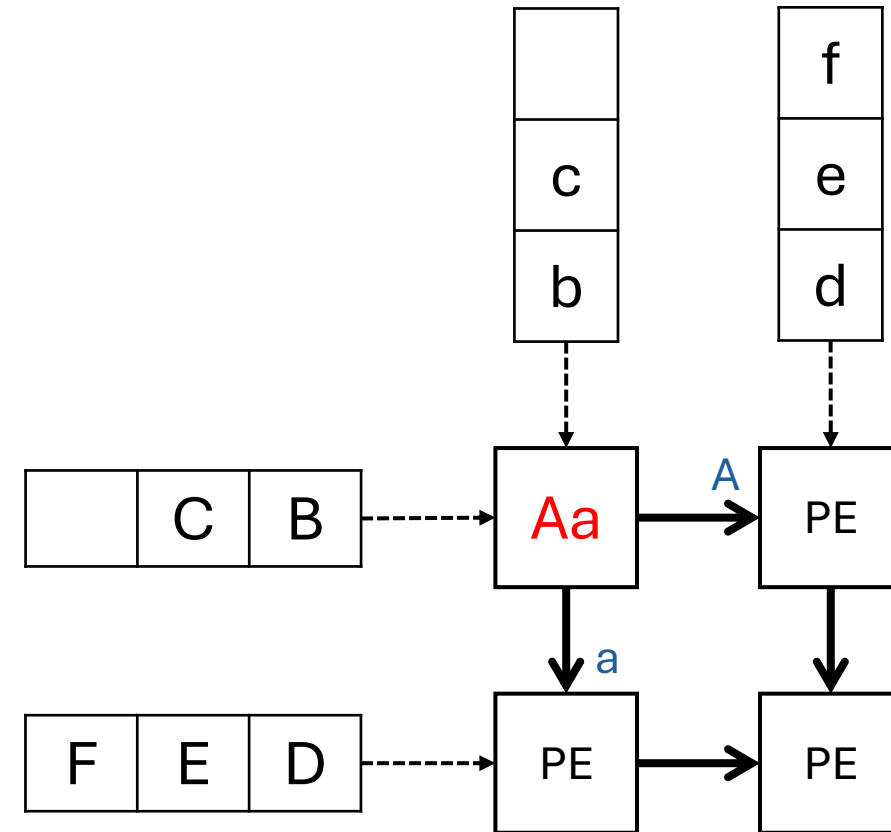
A	B	C
D	E	F

×

a	d
b	e
c	f

  
=

Aa + Bb + Cc	Ad + Be + Cf
Da + Eb + Fc	Dd + Ee + Ff



# Model GEMM using Systolic Array

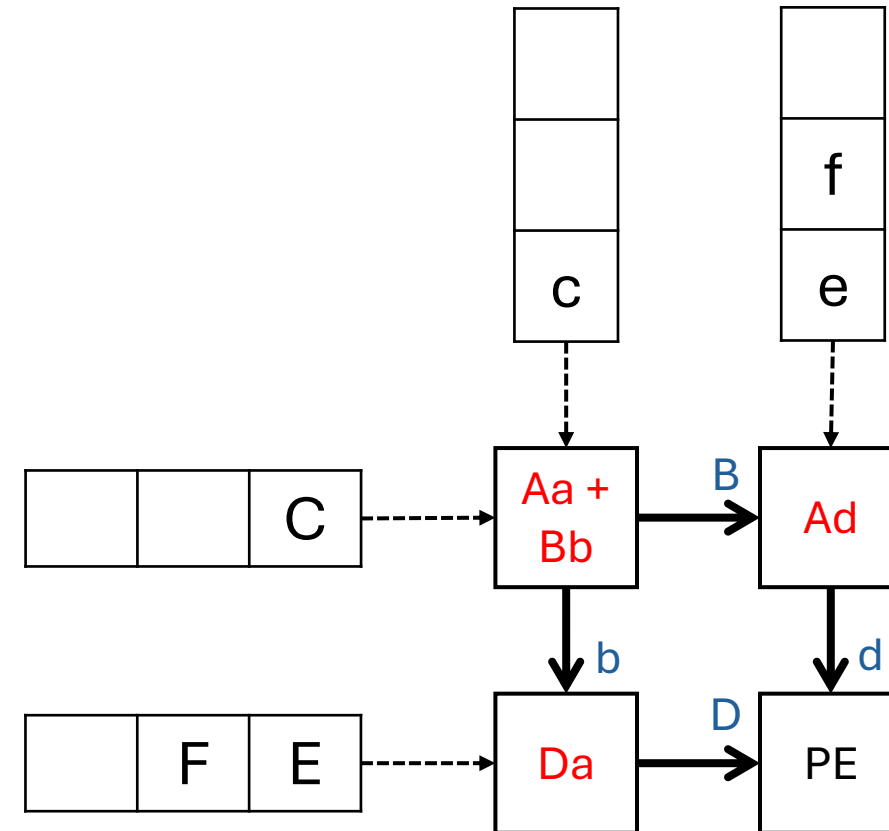
A	B	C
D	E	F

×

a	d
b	e
c	f

  
=

Aa + Bb + Cc	Ad + Be + Cf
Da + Eb + Fc	Dd + Ee + Ff



# Model GEMM using Systolic Array

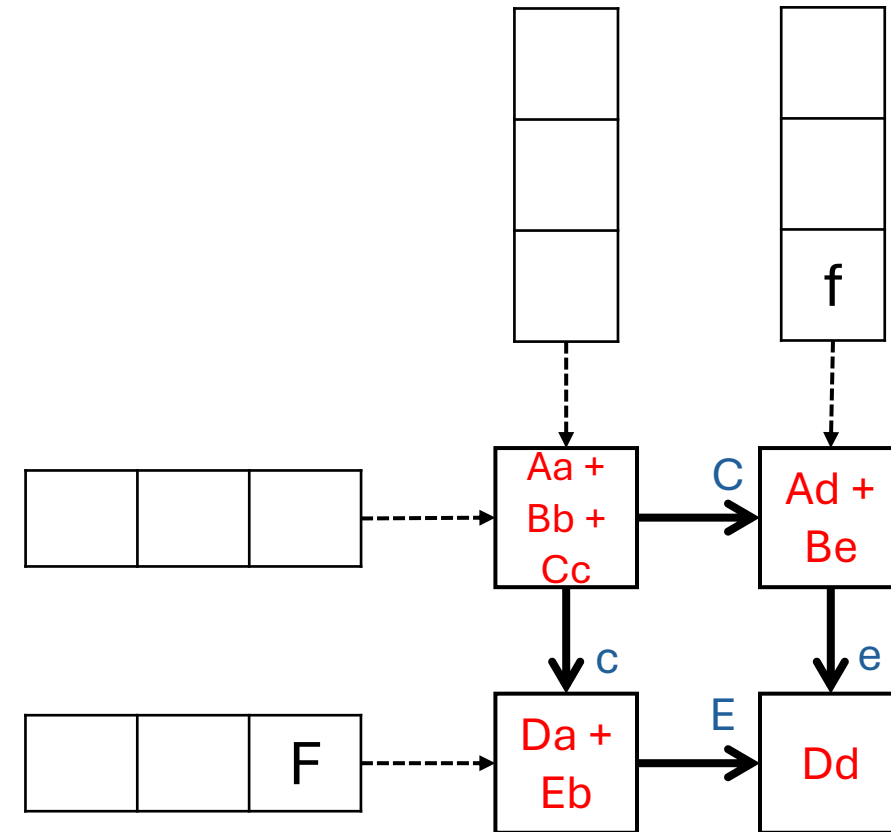
A	B	C
D	E	F

×

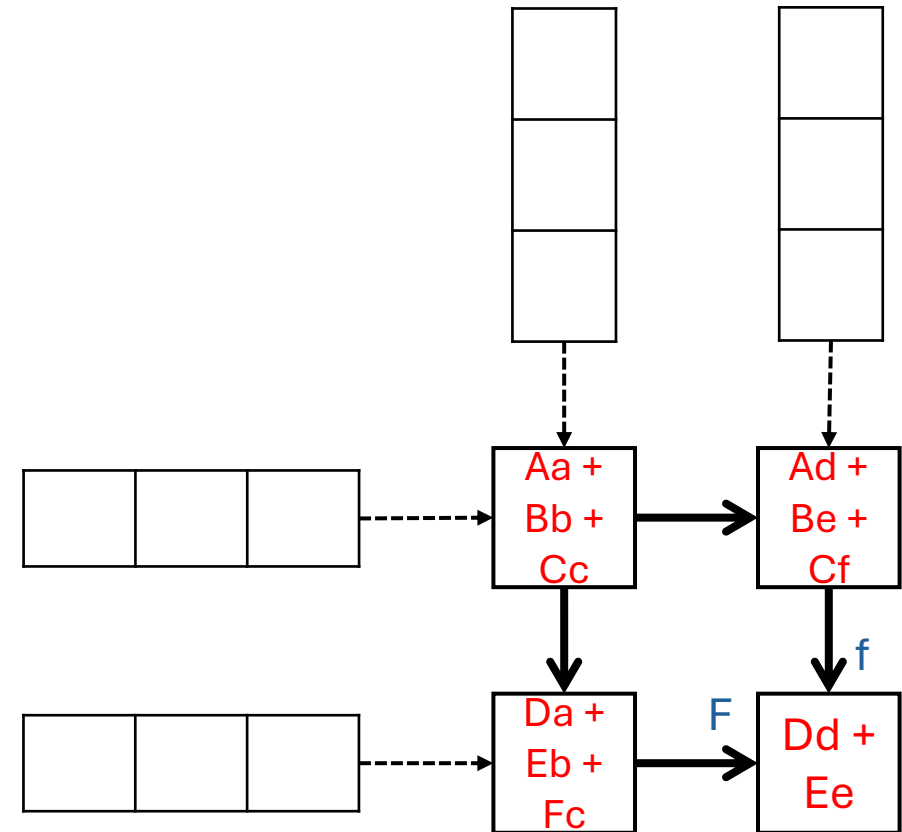
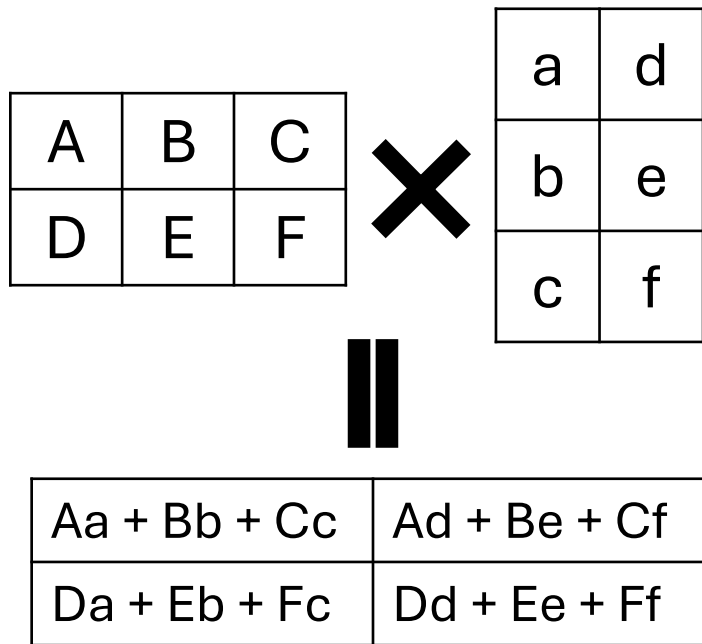
a	d
b	e
c	f

  
=

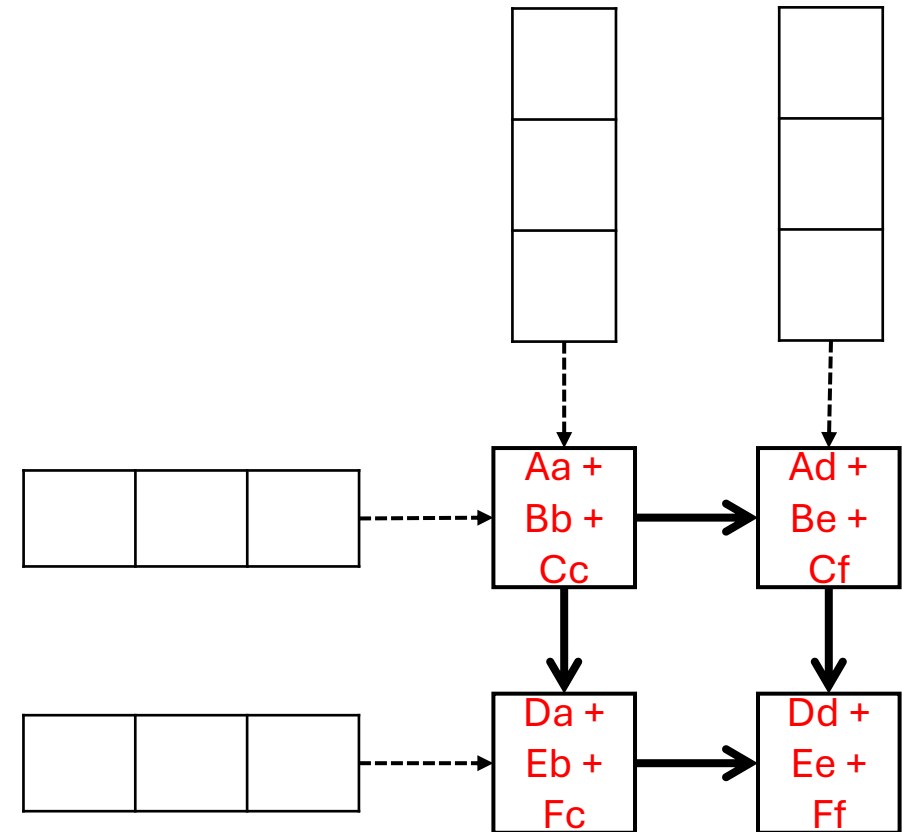
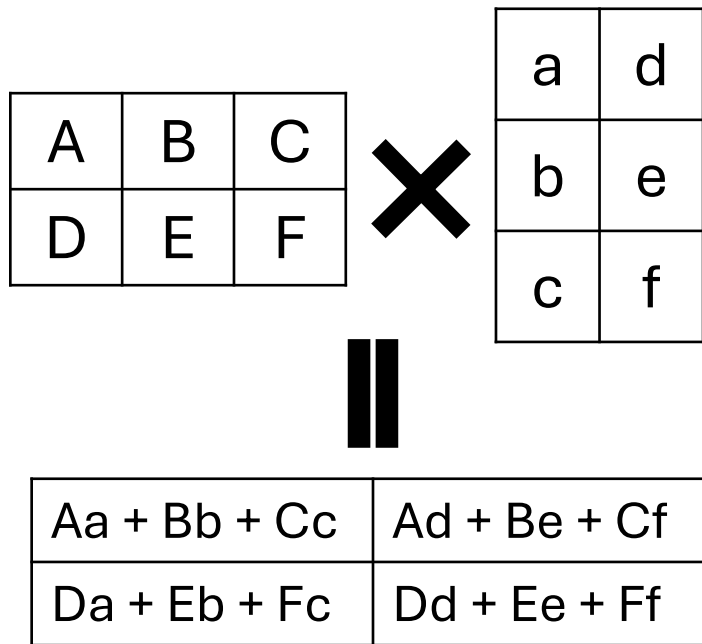
Aa + Bb + Cc	Ad + Be + Cf
Da + Eb + Fc	Dd + Ee + Ff



# Model GEMM using Systolic Array



# Model GEMM using Systolic Array



# Modeling Convolutions

Input Image

a	b	c	d	e
f	g	h	i	j
k	l	m	n	o
p	q	r	s	t
u	v	w	x	y

Filter 1

A	B	C
D	E	F
G	H	I

Filter 2

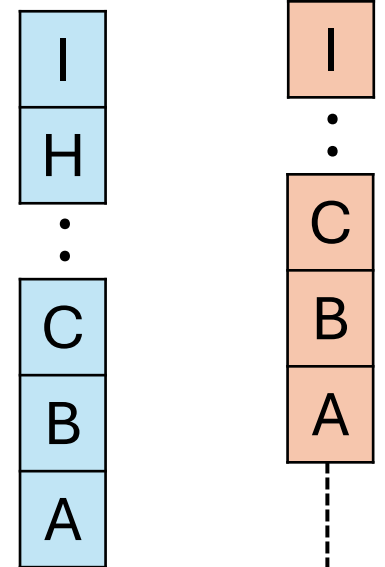
A	B	C
D	E	F
G	H	I

Vectorization

a	b	c	f	g	h	k	l	m
b	c	d	g	h	i	l	m	n

m l ... c b a

n m ... d c b



Data Flow Mapping

Output  
Channel 1

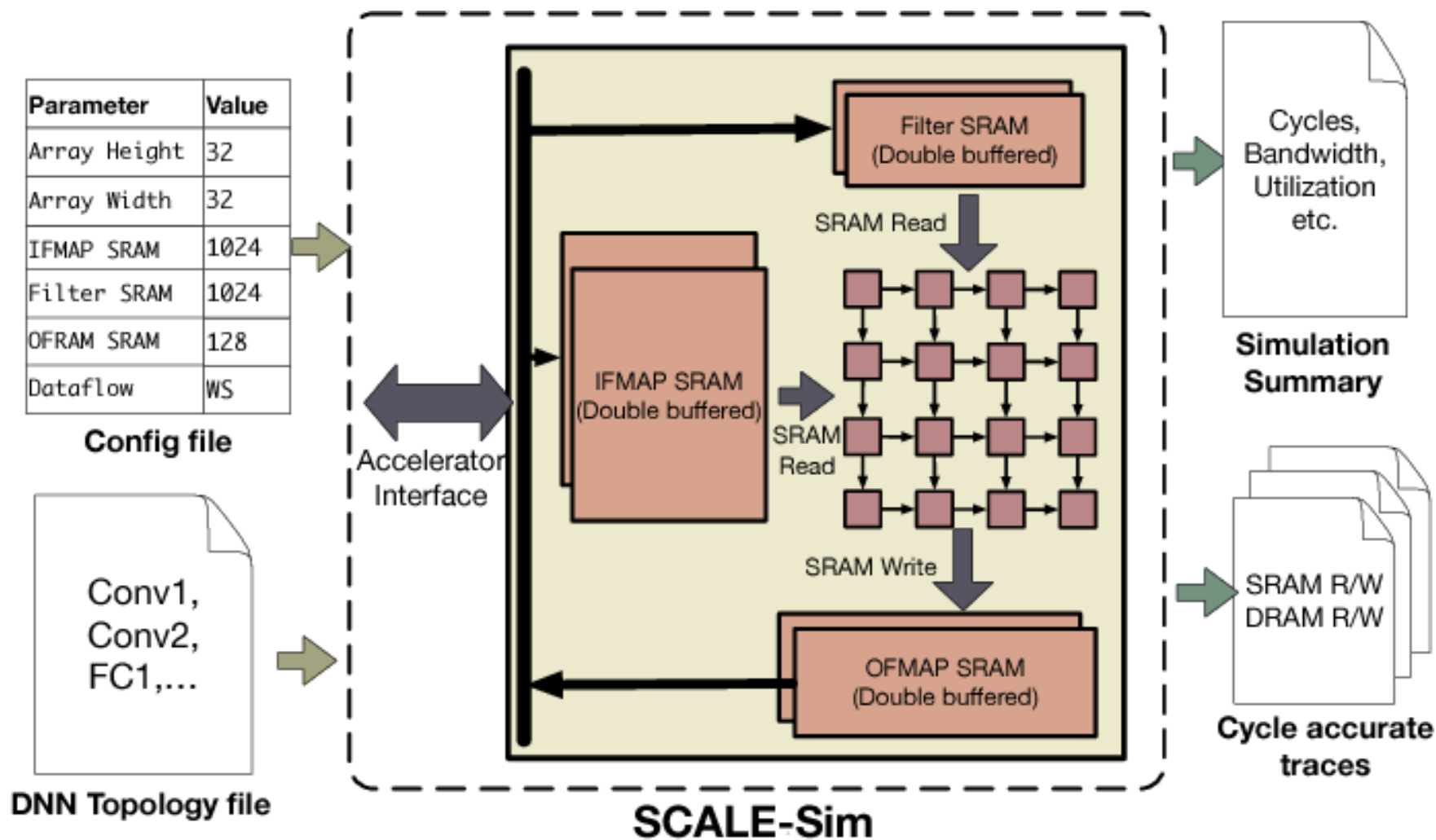
Output  
Channel 2

# Machine Learning Accelerator Trace-Driven Simulator: Scale-Sim: **S**ystolic **C**NN **A**cce**L**erator **S**imulator

- Analytical model for cycle-accurate computation timing, power/energy
  - This is for a specified accelerator hardware configuration and a neural network.
- Provides cycle-level traces for SRAM (on-chip accelerator memory) and off-chip DRAM.
- This can be used to analyze the access characteristics of several machine-learning models.
- Example Models: CNNs, RNNs, Transformers, Recommendation Models, etc.

# Scale-Sim: Analytical Model of Accelerator

Inputs



Outputs



# Scale-Sim: Analytical Model of Accelerator

- On-chip Memory Model
  - Double-buffered Memory
  - Model three sets of double-buffered memory for IFMAP, OFMAP, and FIL
- Compute Model
  - Using Systolic Array
- System Interface Model
  - Cycle-accurate DRAM traces that can be integrated into SoC
- Network Supported
  - Any network that can be mapped to Systolic Array as GEMM

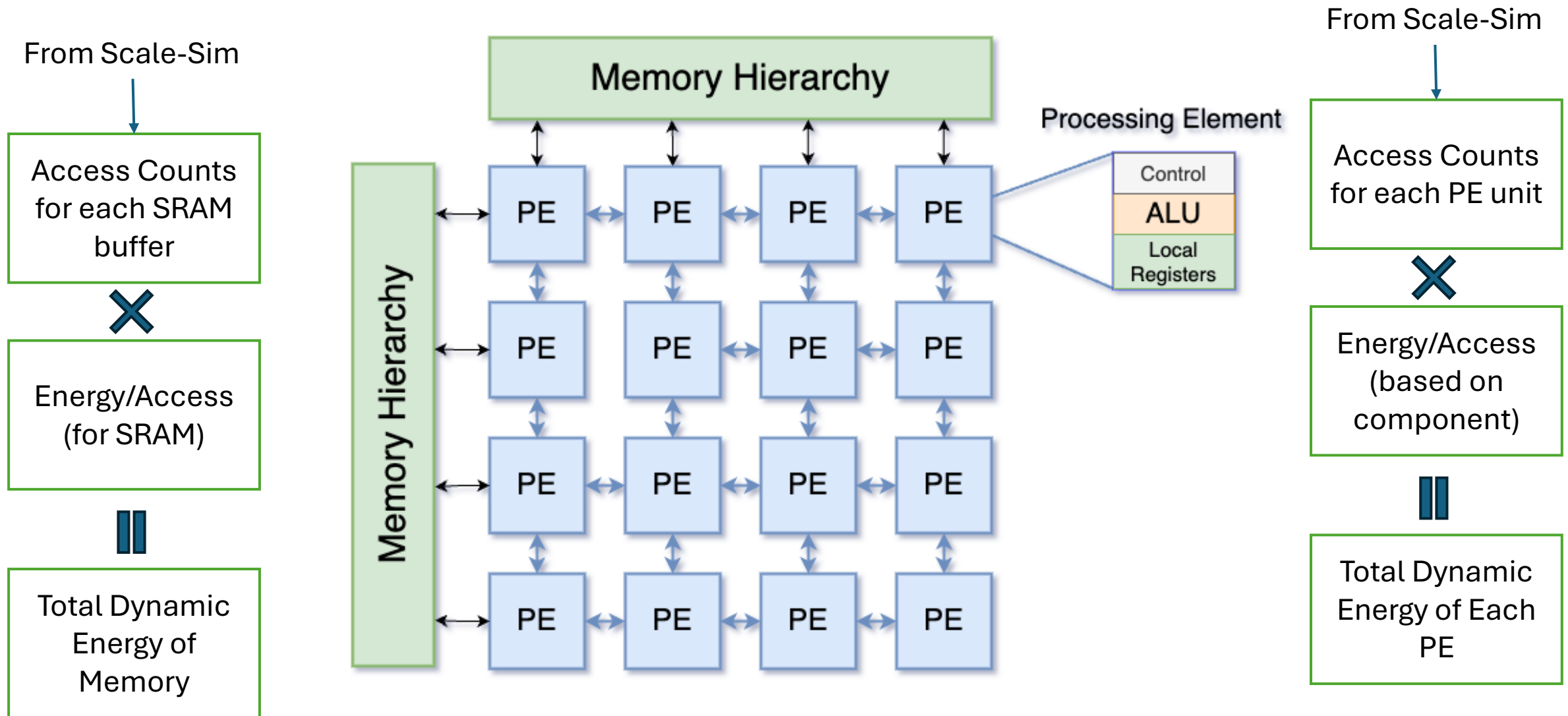
# Power Modeling of Accelerators

- $P = \alpha \cdot C \cdot V_{DD}^2 \cdot f$  ( $\alpha$  – Activity),  $E = P \cdot t$  ( $t$  – compute time)
- Design architecture-level power/energy estimation models.
- Provides accurate power/energy consumption analysis for accelerators.
- Predictions are based on factors like workload, architecture design, and operating conditions.
- Example of such models: Aladdin, Accelergy

# Power Modeling Advantages

- Allows designers to optimize for energy efficiency during the design phase.
- Power Optimization
  - Identify design bottlenecks
  - Optimize parameters like clock frequency or memory hierarchy to minimize power consumption.
- Thermal management:
  - Predict power dissipation to analyze the temperature/thermal map of the Accelerator
  - Design efficient cooling methods for the accelerator

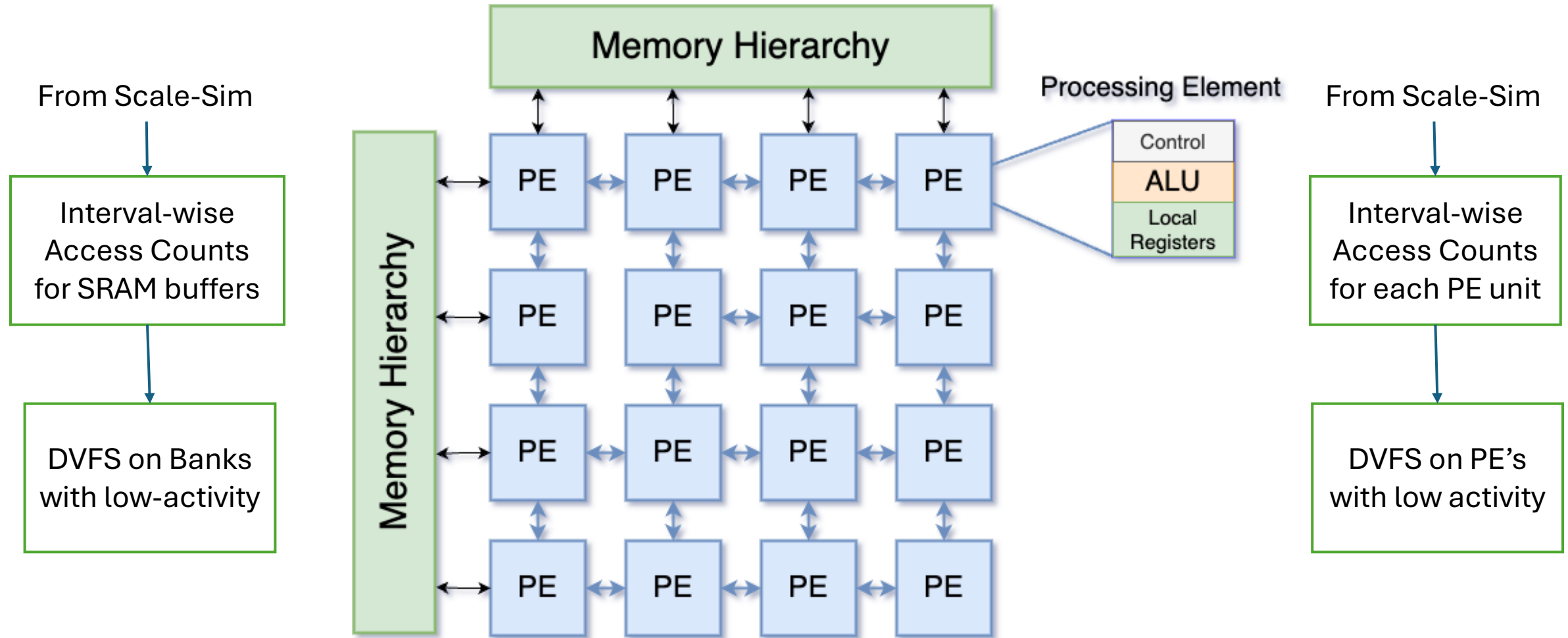
# Energy Model using Scale-Sim



# Dynamic Voltage and Frequency Scaling (DVFS)

- DVFS is used to reduce energy consumption.
- Dynamic Scaling of the supply voltage and frequency of the chip
  - Scaling based on workload
  - Scaling done during idle periods or low computational intensity periods
- Performance Trade-offs:
  - Accurate workload prediction – low and high compute intensity periods
    - Should not optimize high compute zones
  - Efficient voltage switching mechanisms
    - High switching latency – high overheads

# DVFS Policy using Scale-Sim



# Projects

- Design a DVFS policy for the accelerator, based on
  - Progress compared to a given deadline (time, in *ms*)
    - If AHEAD of the deadline, perform DVFS to optimize the accelerator's energy
  - or/and, Accelerator Activity during the computation
    - Based on activity in the accelerator, perform DVFS on components with low utilization or activity
- Design a Dynamic Thermal Management Policy for the Accelerator using Scale-Sim and Hotspot.
  - Collect interval-wise access counts for each component from Scale-Sim
  - Get a power trace from Accelergy
  - Use HotSpot to obtain a Thermal Map – identify the hot regions
  - Use DVFS to optimize temperature – maybe slow down the hot regions

# Thank You

Questions?





# Accelerator Energy Model – Accelergy tool

- Accelergy is an architecture-level energy estimator tool.
- It uses the action counts from the performance simulator to estimate the total energy.

From Scale-Sim or Accelerator  
Performance Simulator

Module	Action	Counts
Global SRAM	access()	x1
Buffer	access()	x2
MAC unit	compute()	x3

Module	Energy/action
Global SRAM	y1
Buffer	y2
MAC unit	y3

$$\begin{aligned} \text{Energy} \\ \text{Estimator} \\ = \\ x1*y1 + \\ x2*y2 + \\ x3*y3 \end{aligned}$$

Total Accelerator  
Energy Estimate

Sample Accelerator

