

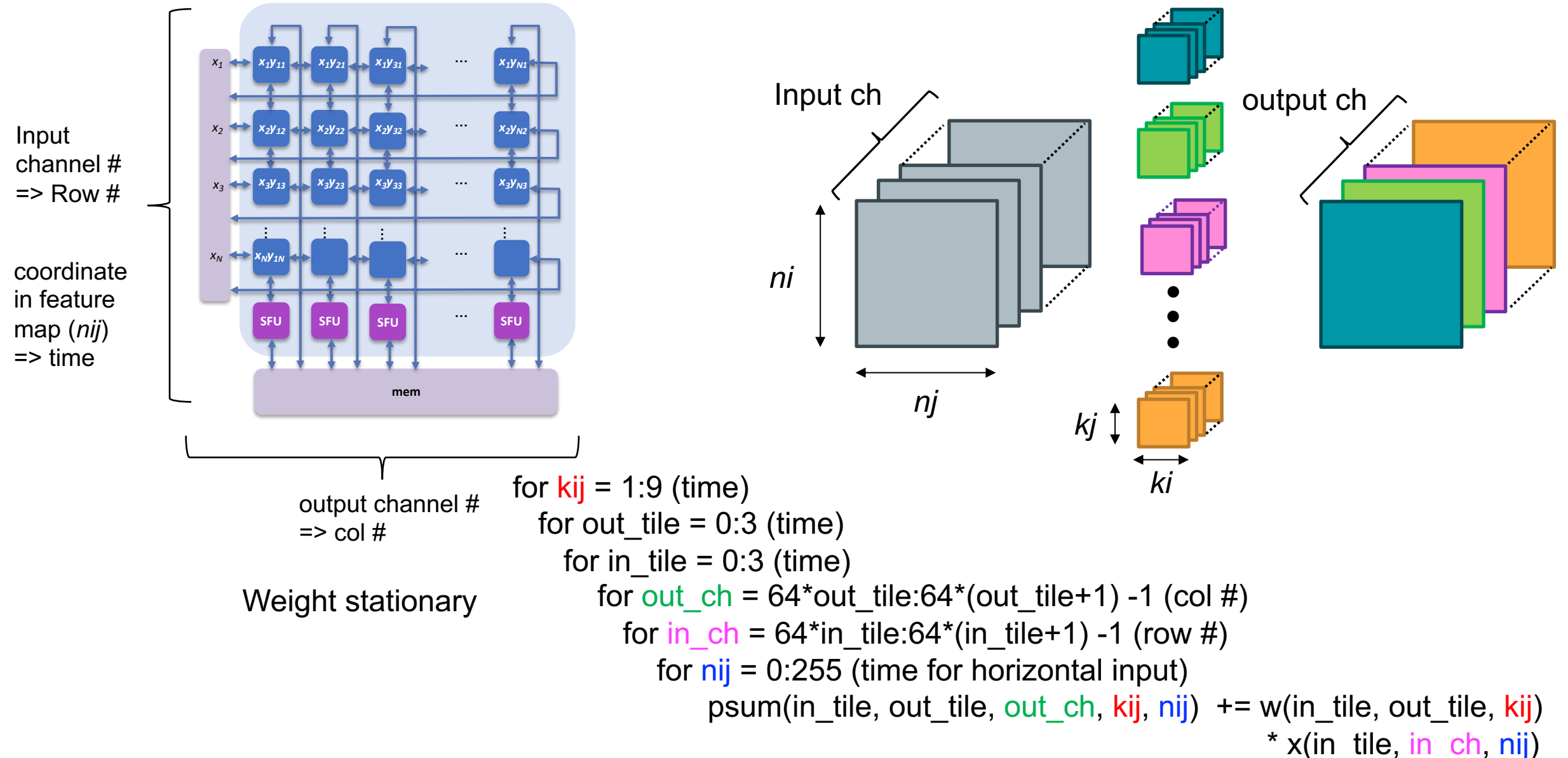
ECE284 Fall 21 W4S2

Low-power VLSI Implementation for Machine Learning

Prof. Mingu Kang

UCSD Computer Engineering

Tiling of the Workload (e.g., in/output channel # = 256)



[Example1] VGGNet_Hardware_Mapping

- VGGNet 3rd layer
- Weight stationary with an array size = 64 X 64
- Input and output channels are both 64
- Goals is to format input, psum, and output in the following form:
- Input (a_pad) format: [ic index (row #), nij (time step)]
- psum (psum) format: [oc index (col #), nij (time step), kij]
- Output (out) format: [oc index, o_nij]

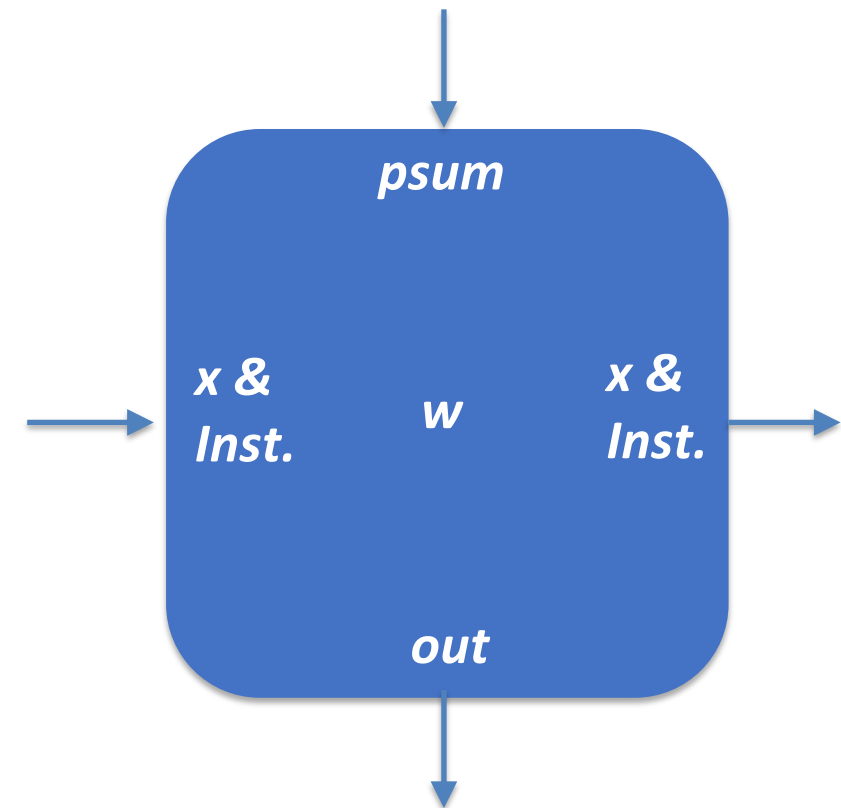
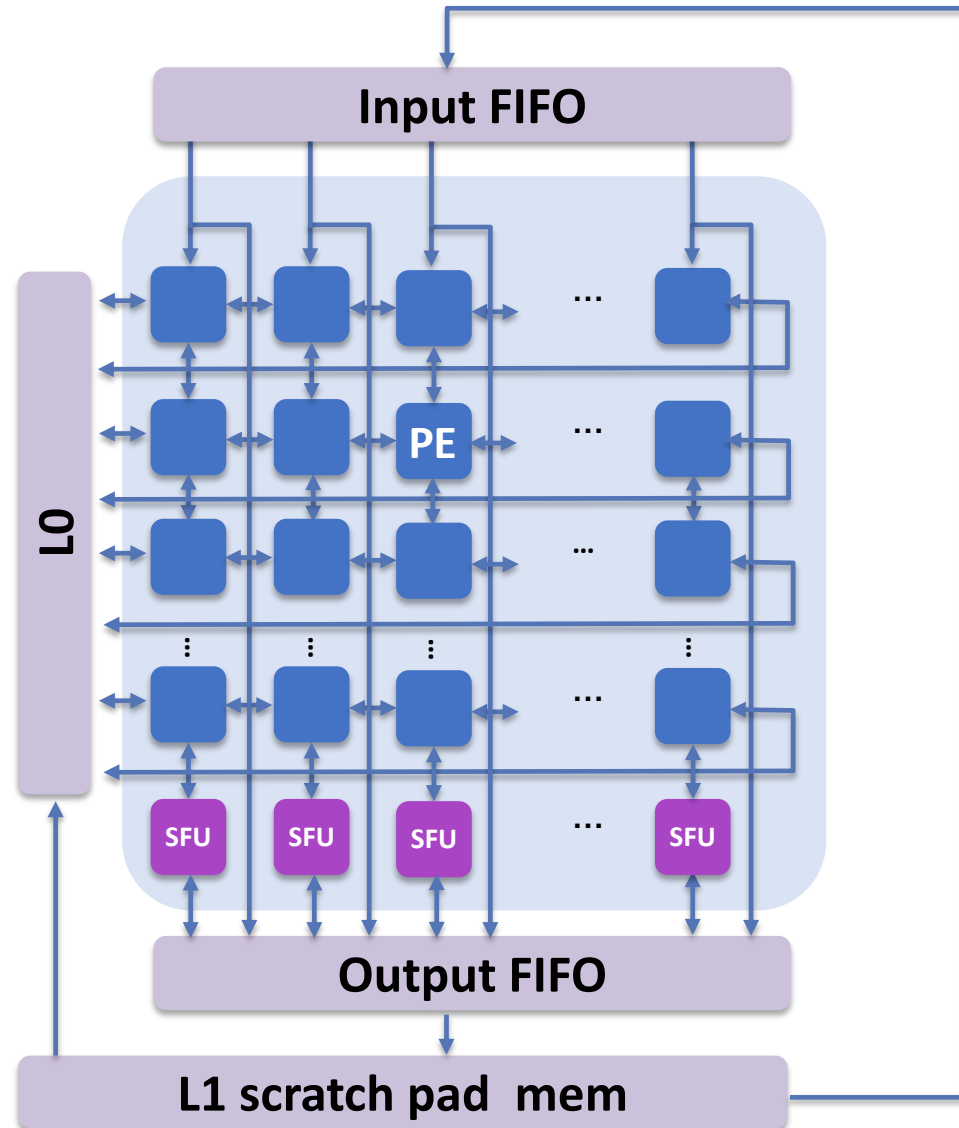
Abbreviation - ic: input channel, oc: output channel, o_nij: nij index for output

[HW_prob1] VGGNet_Hardware_Mapping_Tiling

- VGGNet 3rd layer
- Weight stationary with an array size = 16 X 16
- Input and output channels are both 64
- Needs to be tiled into 16 pieces
- Goals is to format input, psum, and output in the following form:
- Input (a_tile) format: [ic_tile, oc_tile, ic index (row #), nij (time step)]
- psum (psum) format: [ic_tile, oc_tile, oc index (col #), nij (time step), kij]
- Output (out) format: [oc index, o_nij]

Abbreviation - ic: input channel, oc: output channel, o_nij: nij index for output

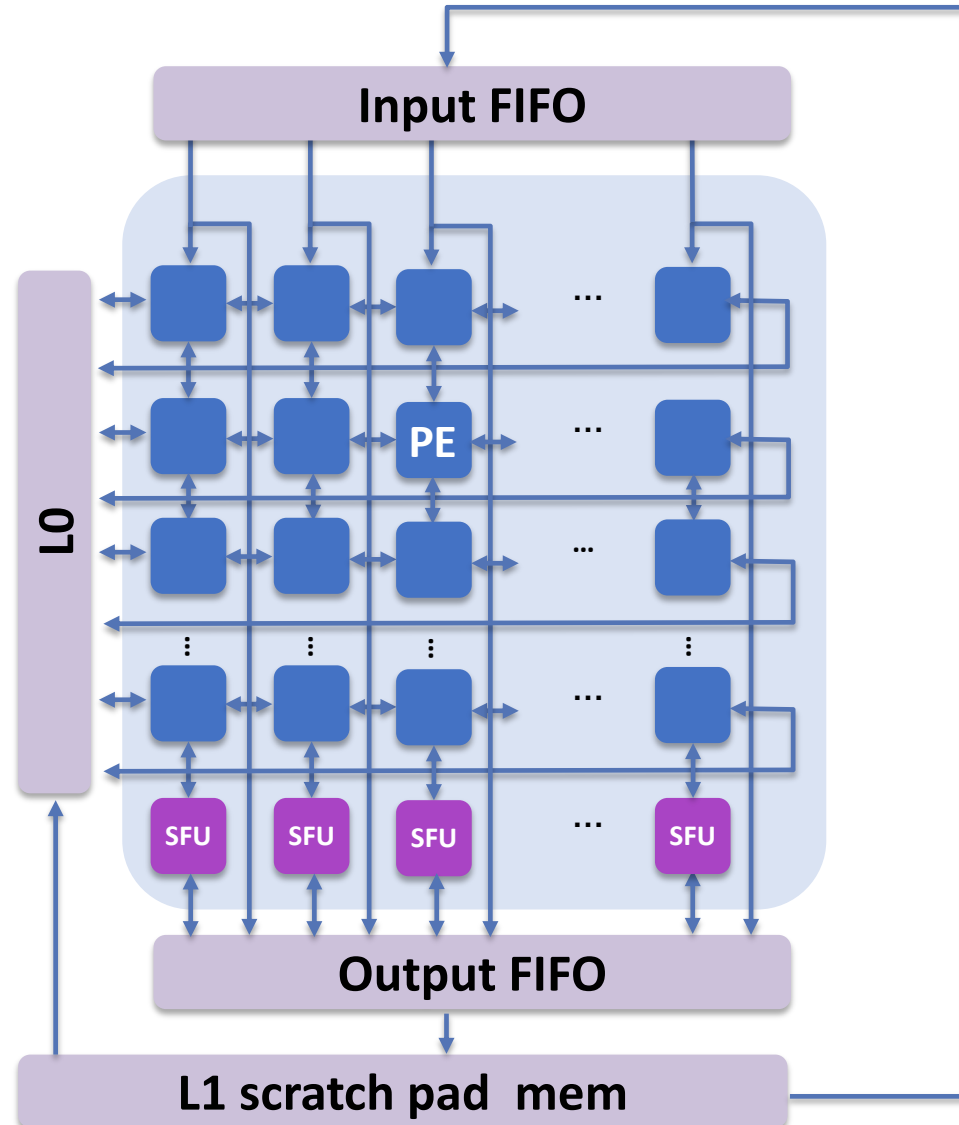
Complete Architecture and Instruction Flow



Data and Instruction flow

J.Oh, S. Lee, M. Kang, et al., "A 3.0 TFLOPS 0.62V Scalable Processor Core for High Compute Utilization AI Training and Inference", VLSI symp 2020

Processing Stage (Convolution)

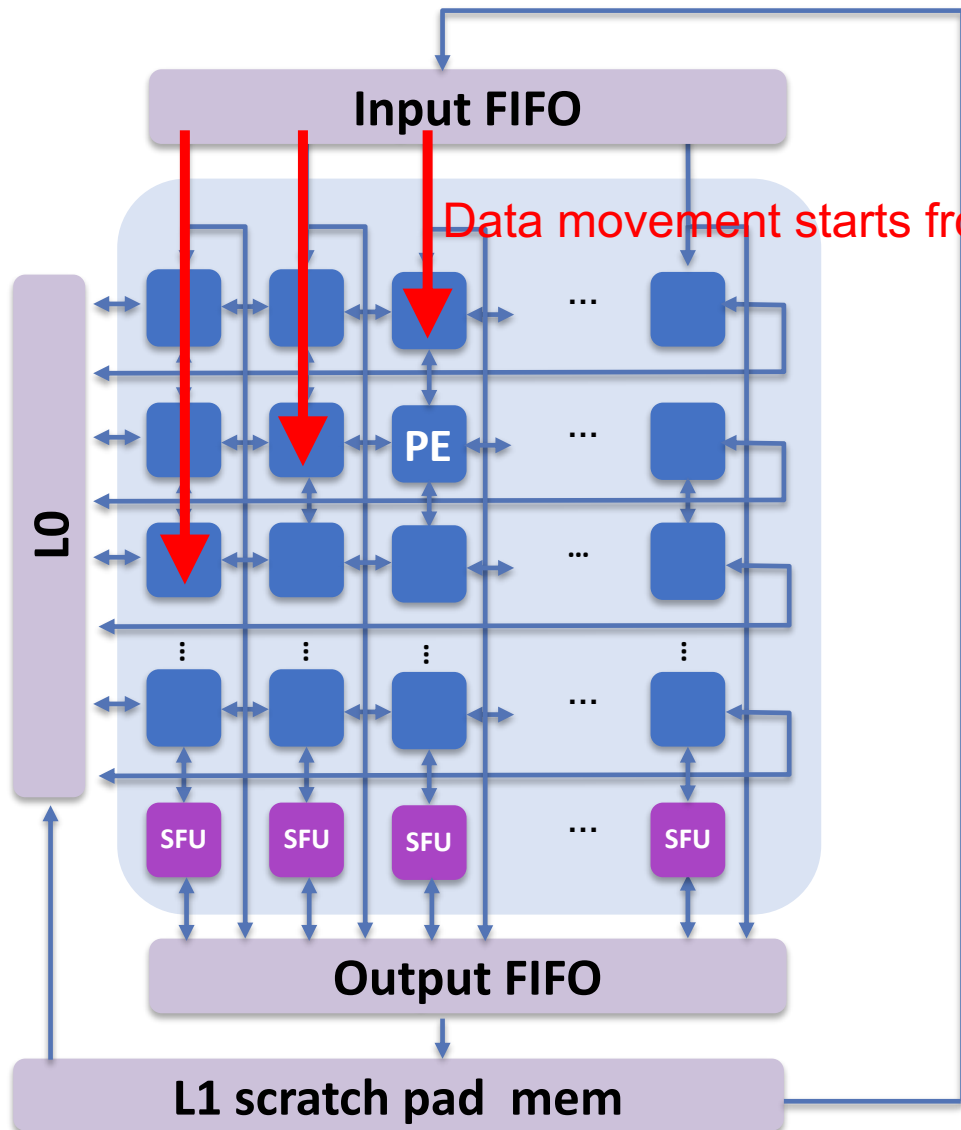


Processing stage

1. L1 scratch pad loading (e.g., from DRAM)
2. Kernel data loading to PE register (via IFIFO)
3. L0 data loading
(simultaneously with Execution)
4. Execution with PEs (psum computation)
5. psum movement to L1 scratch pad (via OFIFO)
(simultaneously with Execution)
6. Accumulation in SFU
(psum brought back from L1 scratch pad)
Then, output store in L1 scratch pad

Kernel Data Loading (with 16 X 16 array)

16 words given at a time



Data movement starts from left column

Note: all data packets (word) received at a time from L1 scratch pad, but outputs at different timing.
-> IFIFO is needed

1st cyc $w_{ic0,oc0}, w_{ic0,oc1}, w_{ic0,oc2}, \dots w_{ic0,oc15}$

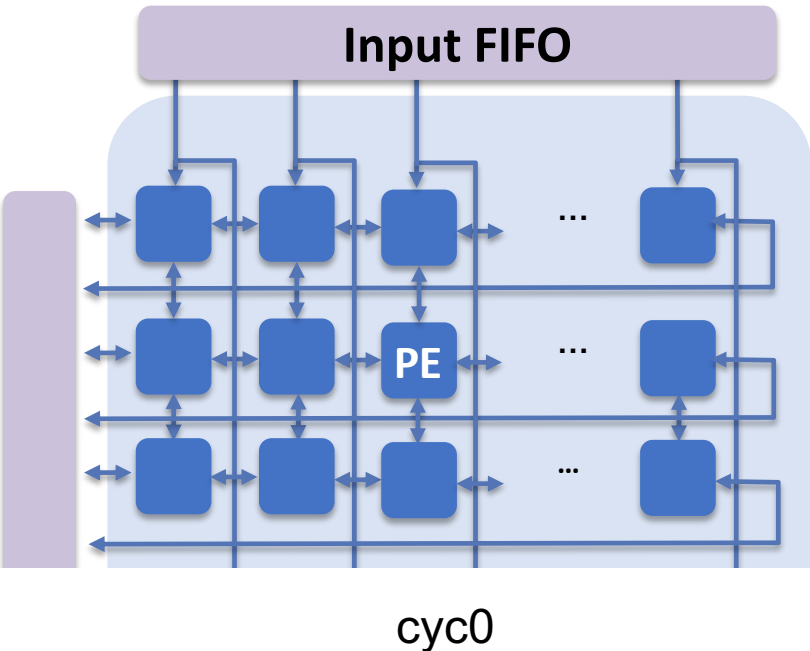
2nd cyc $w_{ic1,oc0}, w_{ic1,oc1}, w_{ic1,oc2}, \dots w_{ic1,oc15}$

...

14th cyc $w_{ic14,oc0}, w_{ic14,oc1}, w_{ic14,oc2}, \dots w_{ic14,oc15}$

15th cyc $w_{ic15,oc0}, w_{ic15,oc1}, w_{ic15,oc2}, \dots w_{ic15,oc15}$

Kernel Data Loading (with 16 X 16 array)

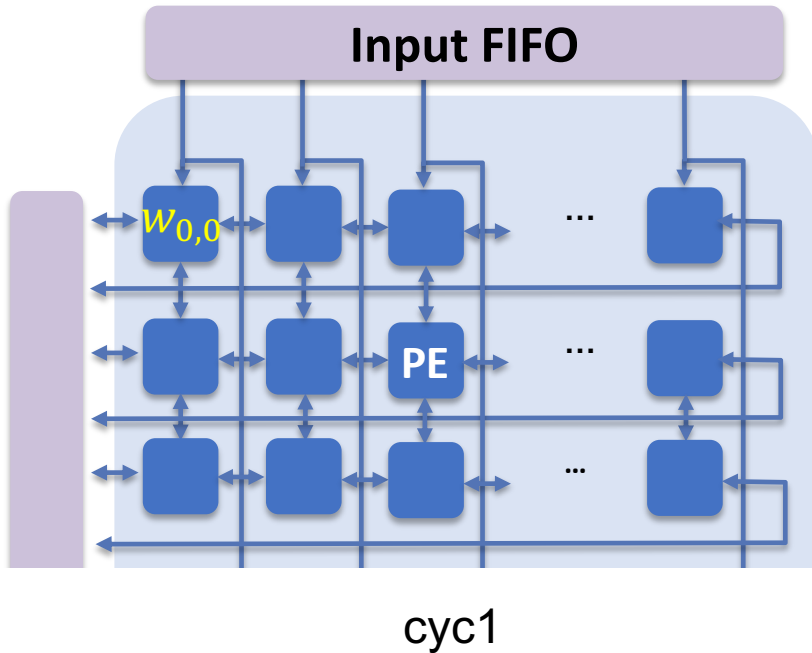


IFIFO Contents

$w_{ic0,oc0}$	$w_{ic0,oc1}$	$w_{ic0,oc2}$...	$w_{ic0,oc15}$

Read

Kernel Data Loading (with 16 X 16 array)

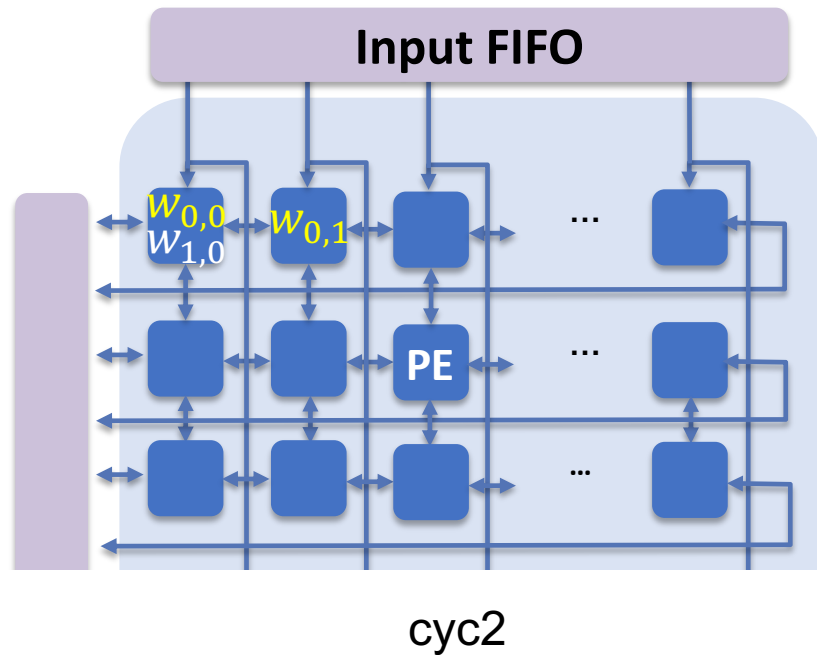


IFIFO Contents

cyc0	$w_{ic0,oc0}$	$w_{ic0,oc1}$	$w_{ic0,oc2}$		$w_{ic0,oc15}$
cyc1	$w_{ic1,oc0}$	$w_{ic1,oc1}$	$w_{ic1,oc2}$		$w_{ic1,oc15}$
		$w_{ic0,oc1}$	$w_{ic0,oc2}$		$w_{ic0,oc15}$
Read		➡	Read		

Yellow text: captured in the PE's local register

Kernel Data Loading (with 16 X 16 array)

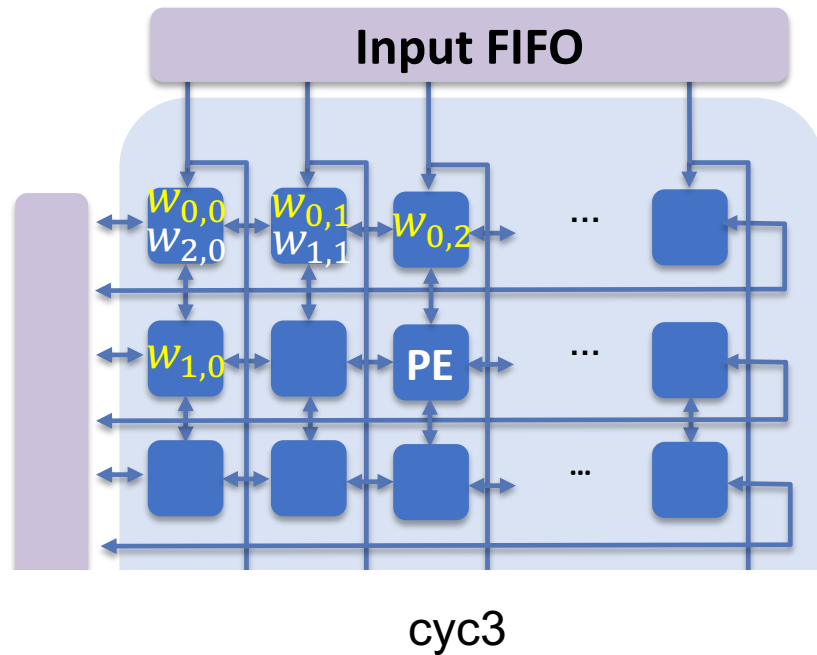


IFIFO Contents

cyc0	$w_{ic0,oc0}$	$w_{ic0,oc1}$	$w_{ic0,oc2}$		$w_{ic0,oc15}$
cyc1	$w_{ic1,oc0}$	$w_{ic1,oc1}$	$w_{ic1,oc2}$		$w_{ic1,oc15}$
		$w_{ic0,oc1}$	$w_{ic0,oc2}$		$w_{ic0,oc15}$
cyc2	$w_{ic2,oc0}$	$w_{ic2,oc1}$	$w_{ic2,oc2}$		$w_{ic2,oc15}$
		$w_{ic1,oc1}$	$w_{ic1,oc2}$		$w_{ic1,oc15}$
			$w_{ic0,oc2}$		$w_{ic0,oc15}$
Read → Read → Read					

Yellow text: captured in the PE's local register

Kernel Data Loading (with 16 X 16 array)

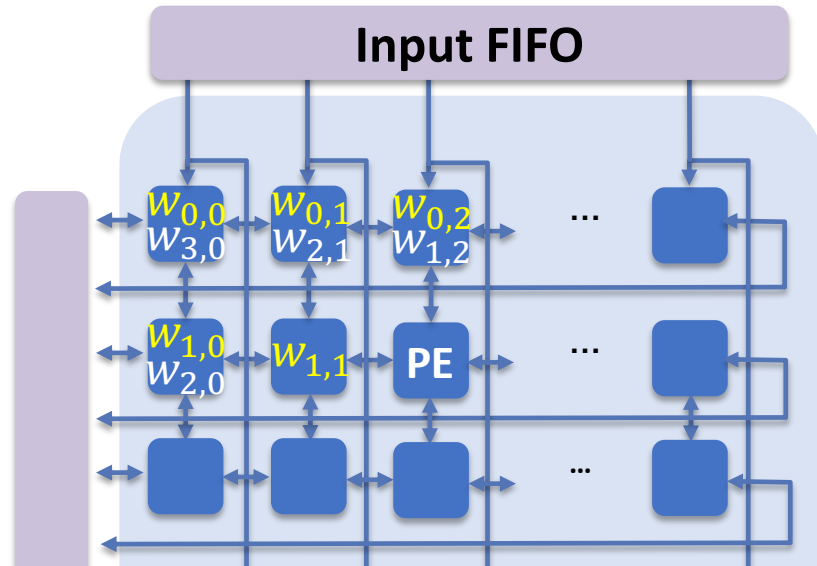


Yellow text: captured in the PE's local register

IFIFO Contents

cyc0	$W_{ic0,oc0}$	$W_{ic0,oc1}$	$W_{ic0,oc2}$		$W_{ic0,oc15}$
cyc1	$W_{ic1,oc0}$	$W_{ic1,oc1}$	$W_{ic1,oc2}$		$W_{ic1,oc15}$
		$W_{ic0,oc1}$	$W_{ic0,oc2}$		$W_{ic0,oc15}$
cyc2	$W_{ic2,oc0}$	$W_{ic2,oc1}$	$W_{ic2,oc2}$		$W_{ic2,oc15}$
		$W_{ic1,oc1}$	$W_{ic1,oc2}$		$W_{ic1,oc15}$
			$W_{ic0,oc2}$		$W_{ic0,oc15}$
cyc3	$W_{ic3,oc0}$	$W_{ic3,oc1}$	$W_{ic3,oc2}$		$W_{ic3,oc15}$
		$W_{ic2,oc1}$	$W_{ic2,oc2}$		$W_{ic2,oc15}$
			$W_{ic1,oc2}$		$W_{ic1,oc15}$
					$W_{ic0,oc15}$

Kernel Data Loading (with 16 X 16 array)



cyc4

Yellow text: captured in the PE's local register

IFIFO Contents

cyc2

$W_{ic2,oc0}$	$W_{ic2,oc1}$	$W_{ic2,oc2}$		$W_{ic2,oc15}$
	$W_{ic1,oc1}$	$W_{ic1,oc2}$		$W_{ic1,oc15}$
		$W_{ic0,oc2}$		$W_{ic0,oc15}$

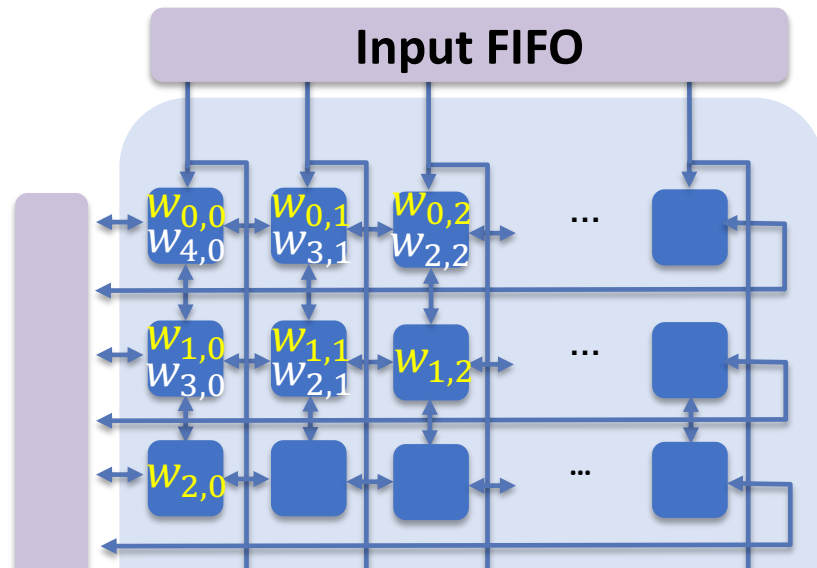
cyc3

$W_{ic3,oc0}$	$W_{ic3,oc1}$	$W_{ic3,oc2}$		$W_{ic3,oc15}$
	$W_{ic2,oc1}$	$W_{ic2,oc2}$		$W_{ic2,oc15}$
		$W_{ic1,oc2}$		$W_{ic1,oc15}$
				$W_{ic0,oc15}$

cyc4

$W_{ic4,oc0}$	$W_{ic4,oc1}$	$W_{ic4,oc2}$		$W_{ic4,oc15}$
	$W_{ic3,oc1}$	$W_{ic3,oc2}$		$W_{ic3,oc15}$
		$W_{ic2,oc2}$		$W_{ic2,oc15}$
				$W_{ic1,oc15}$
				$W_{ic0,oc15}$

Kernel Data Loading (with 16 X 16 array)



cyc5

Yellow text: captured in the PE's local register

IFIFO Contents

cyc3

$W_{ic3,oc0}$	$W_{ic3,oc1}$	$W_{ic3,oc2}$		$W_{ic3,oc15}$
	$W_{ic2,oc1}$	$W_{ic2,oc2}$		$W_{ic2,oc15}$
		$W_{ic1,oc2}$		$W_{ic1,oc15}$
				$W_{ic0,oc15}$

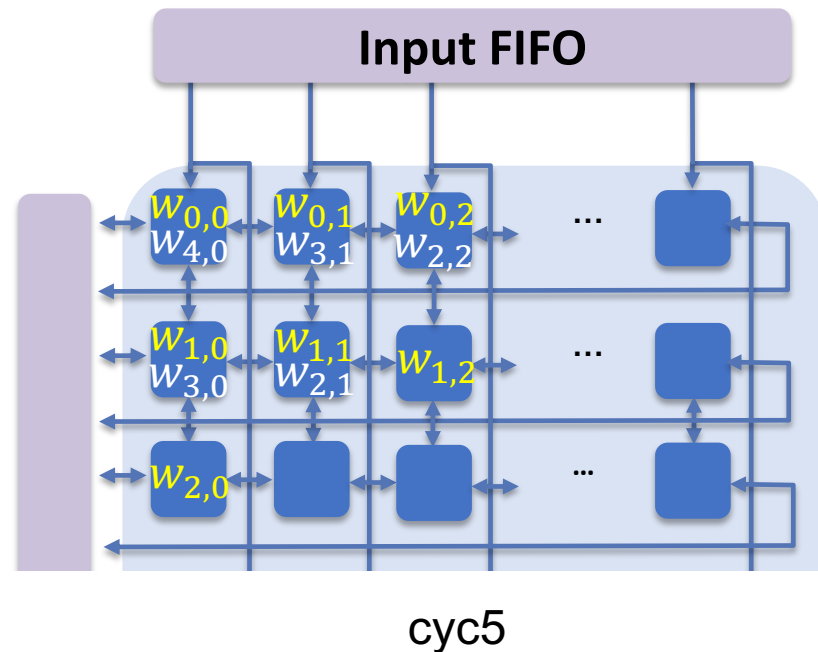
cyc4

$W_{ic4,oc0}$	$W_{ic4,oc1}$	$W_{ic4,oc2}$		$W_{ic4,oc15}$
	$W_{ic3,oc1}$	$W_{ic3,oc2}$		$W_{ic3,oc15}$
		$W_{ic2,oc2}$		$W_{ic2,oc15}$
				$W_{ic1,oc15}$
				$W_{ic0,oc15}$

cyc5

$W_{ic5,oc0}$	$W_{ic5,oc1}$	$W_{ic5,oc2}$		$W_{ic5,oc15}$
	$W_{ic4,oc1}$	$W_{ic4,oc2}$		$W_{ic4,oc15}$
		$W_{ic3,oc2}$		$W_{ic3,oc15}$
				$W_{ic2,oc15}$
				$W_{ic1,oc15}$
				$W_{ic0,oc15}$

Kernel Data Loading (with 16 X 16 array)



IFIFO Contents

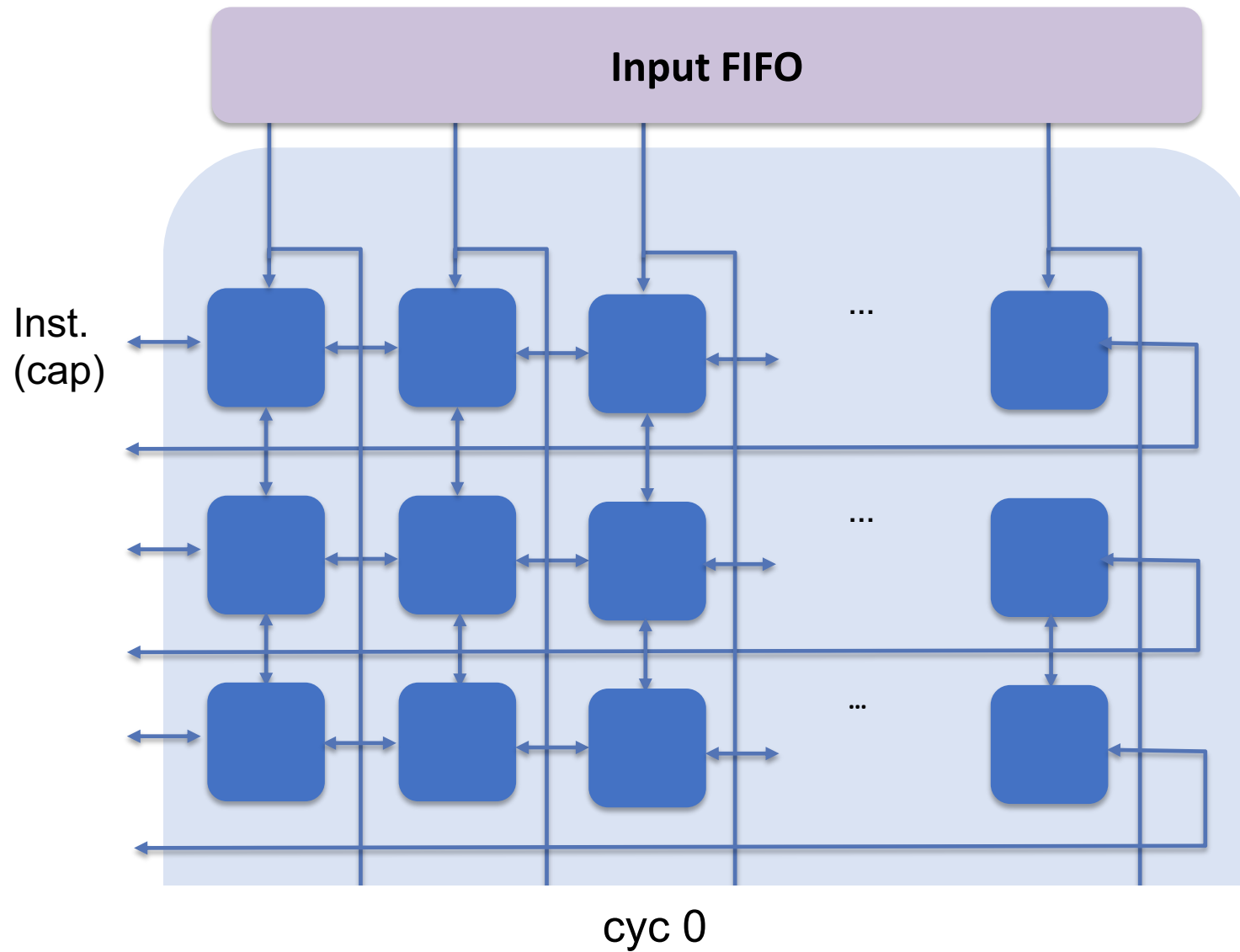
cyc5

$W_{ic5,oc0}$	$W_{ic5,oc1}$	$W_{ic5,oc2}$		$W_{ic5,oc15}$
	$W_{ic4,oc1}$	$W_{ic4,oc2}$		$W_{ic4,oc15}$
		$W_{ic3,oc2}$		$W_{ic3,oc15}$
				$W_{ic2,oc15}$
				$W_{ic1,oc15}$
				$W_{ic0,oc15}$

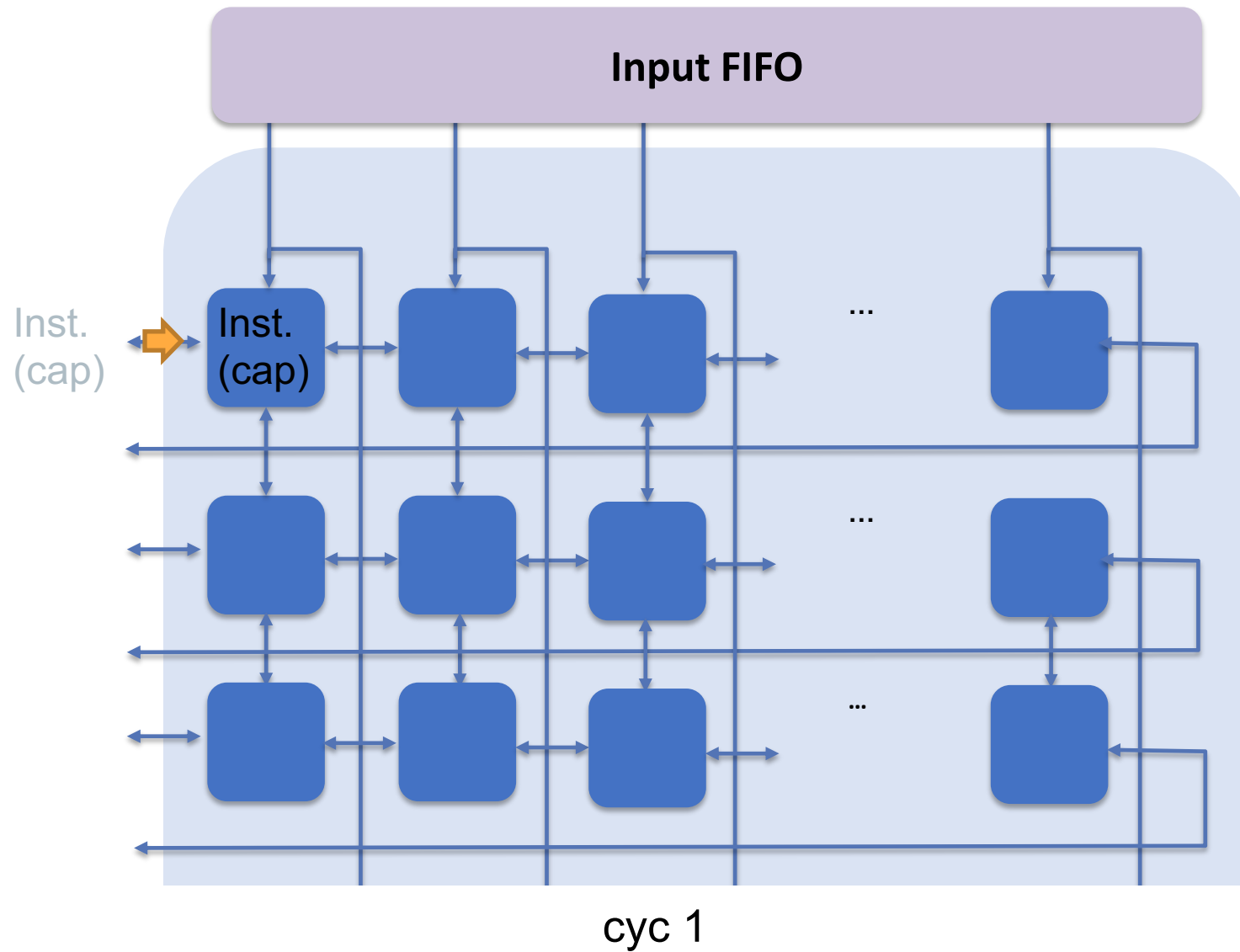
Observation

- IFIFO: requires up to (minimum) 16 depth at the last column
- PE: captures the value at the right timing by the instruction flow
- Roughly 16 (last column start) + 16 (last element start) + 16 (last element delivery)
= 48 cycles required to fill the entire array

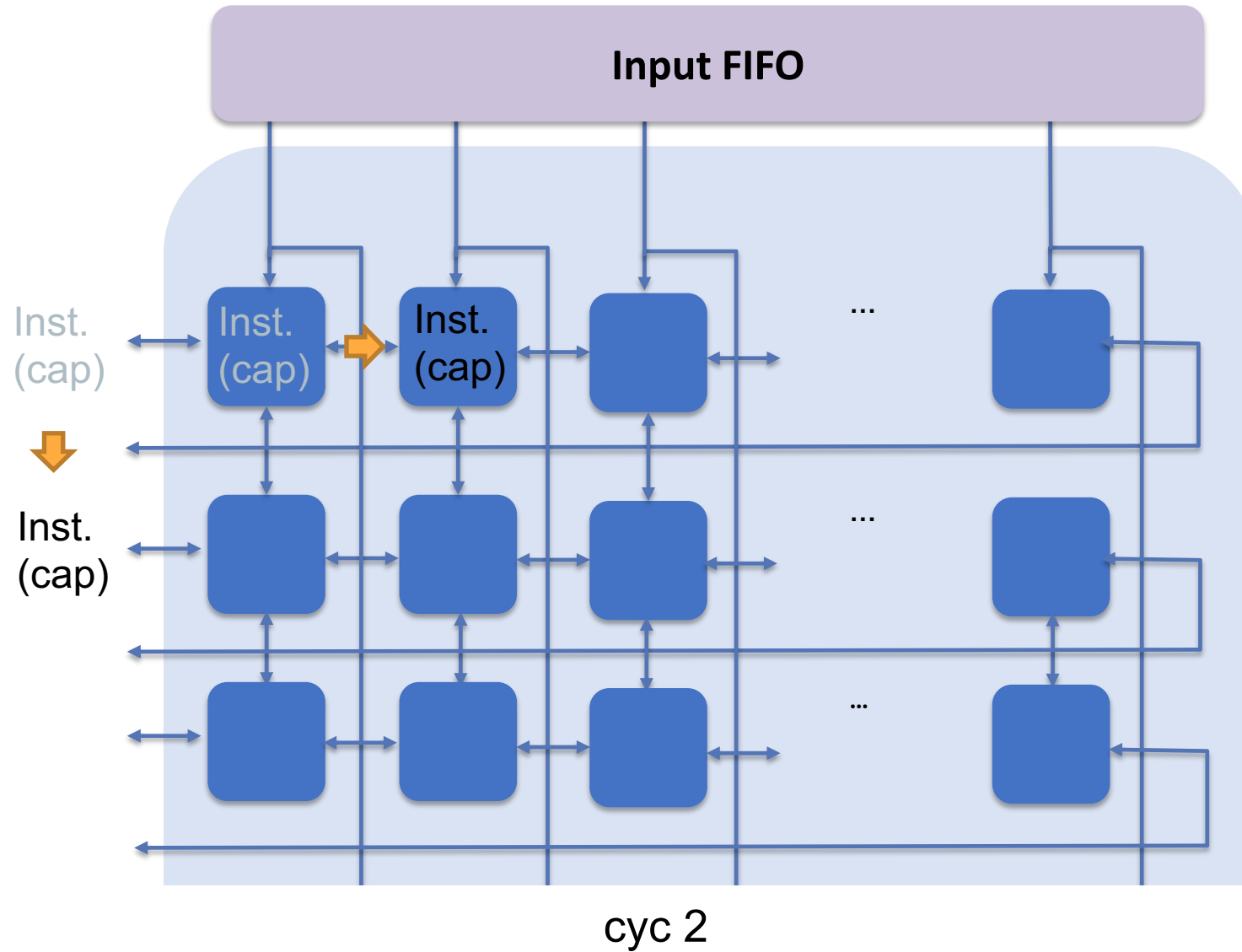
Instruction Flow



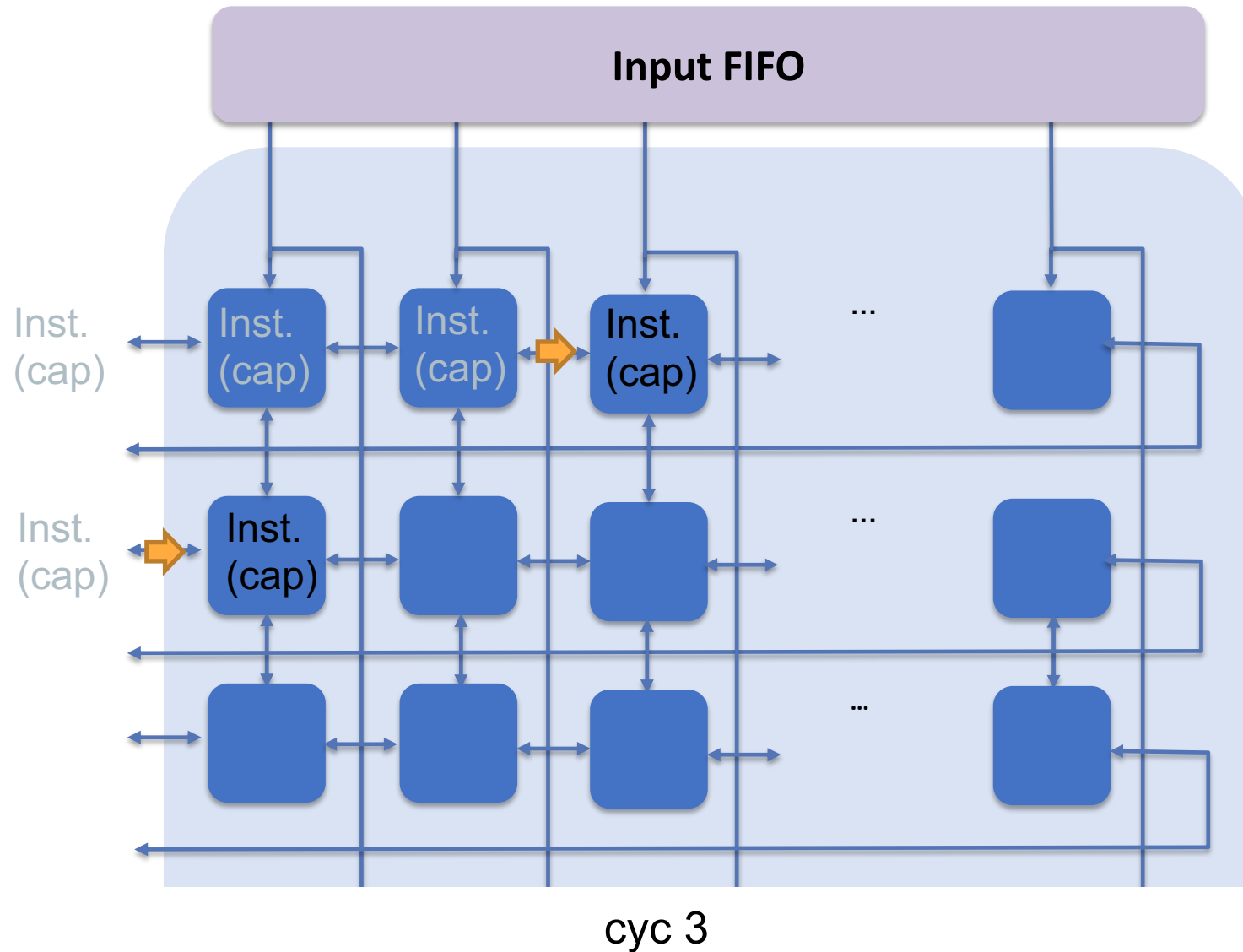
Instruction Flow



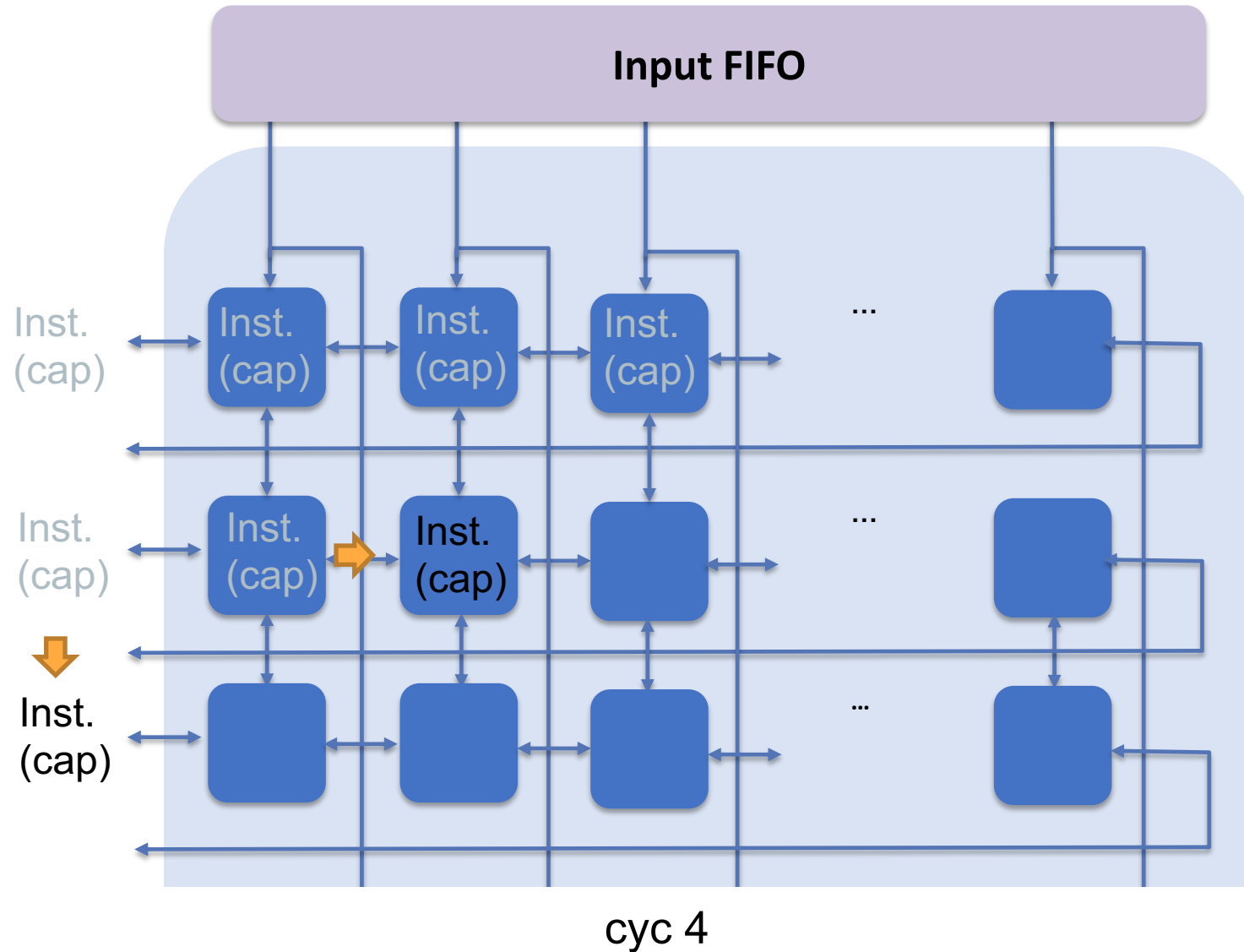
Instruction Flow



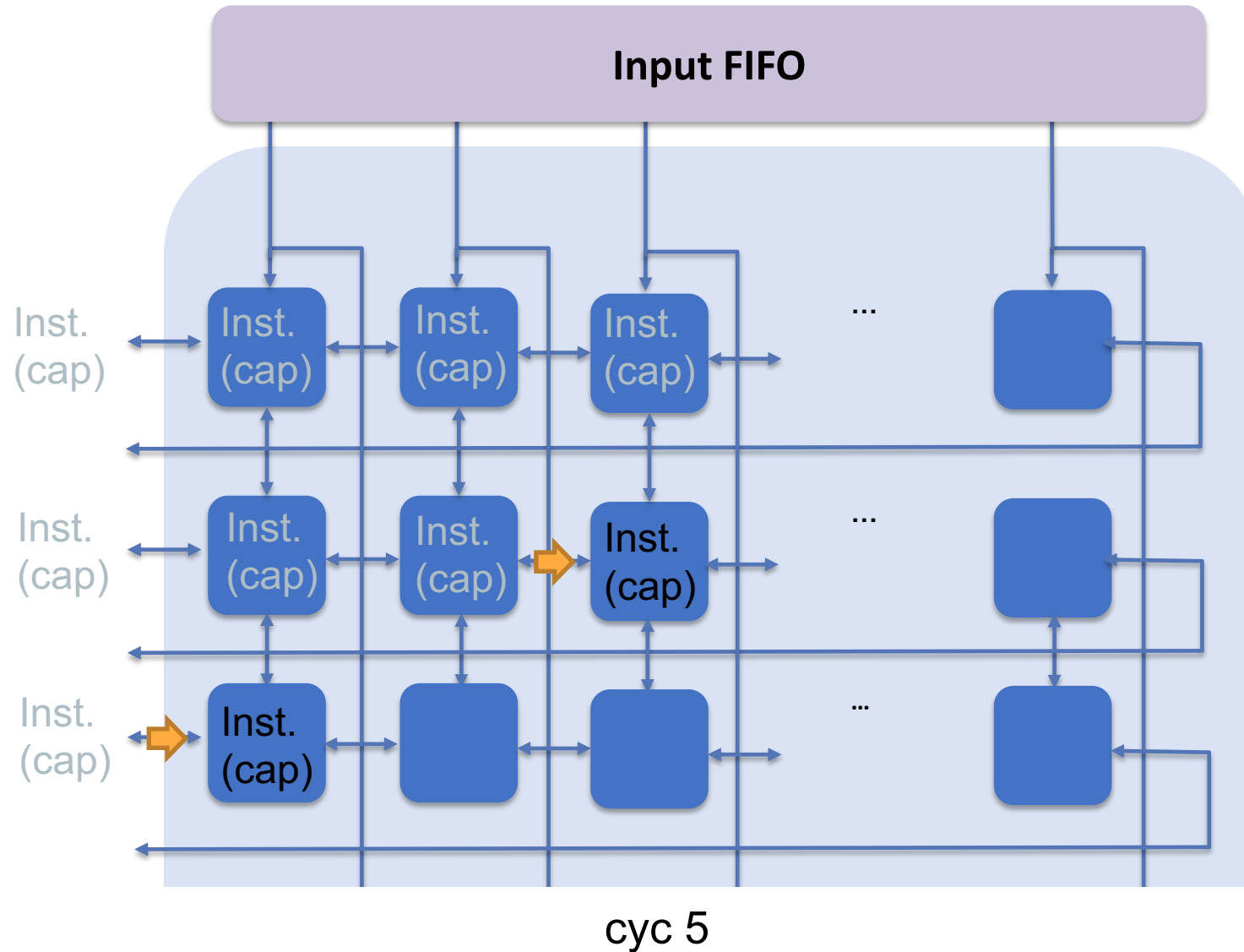
Instruction Flow



Instruction Flow

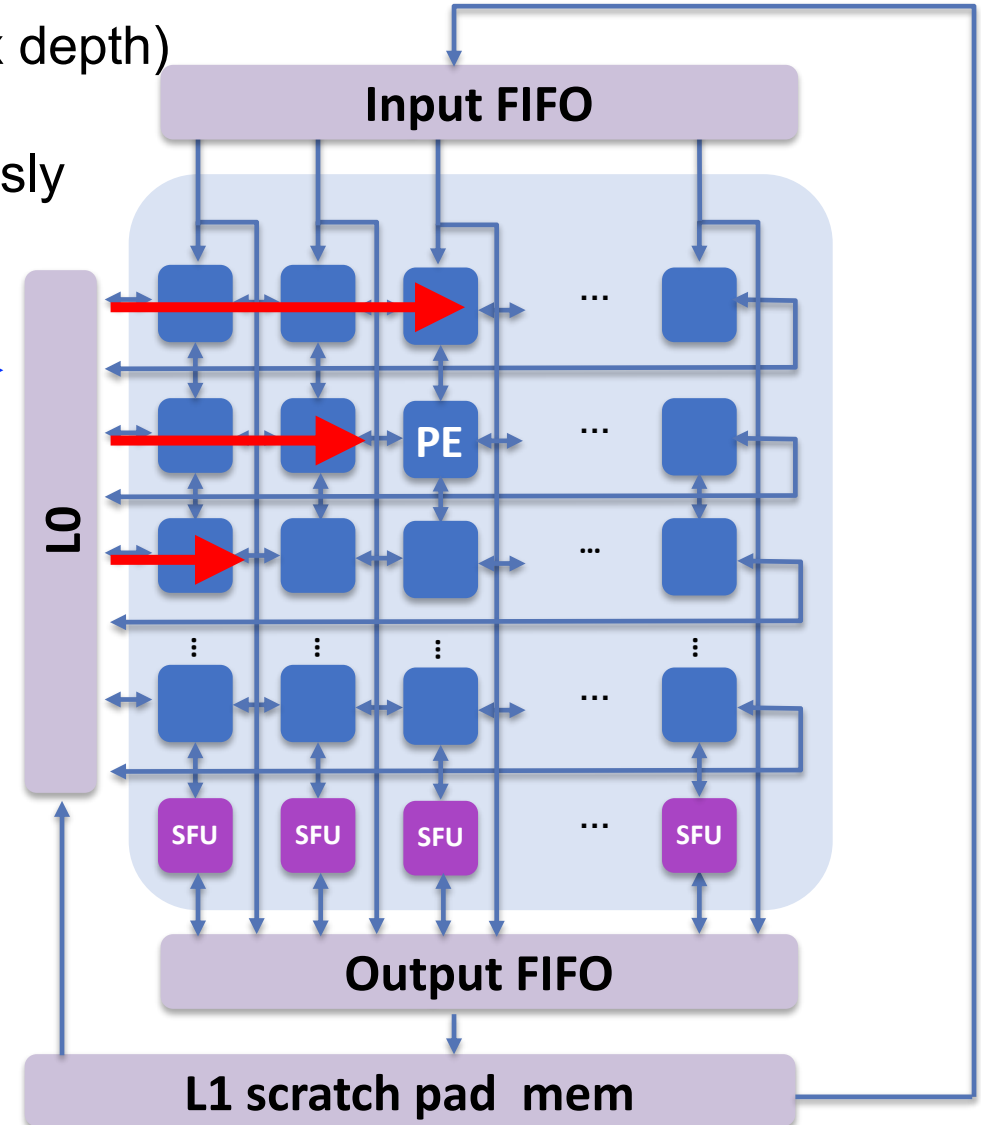
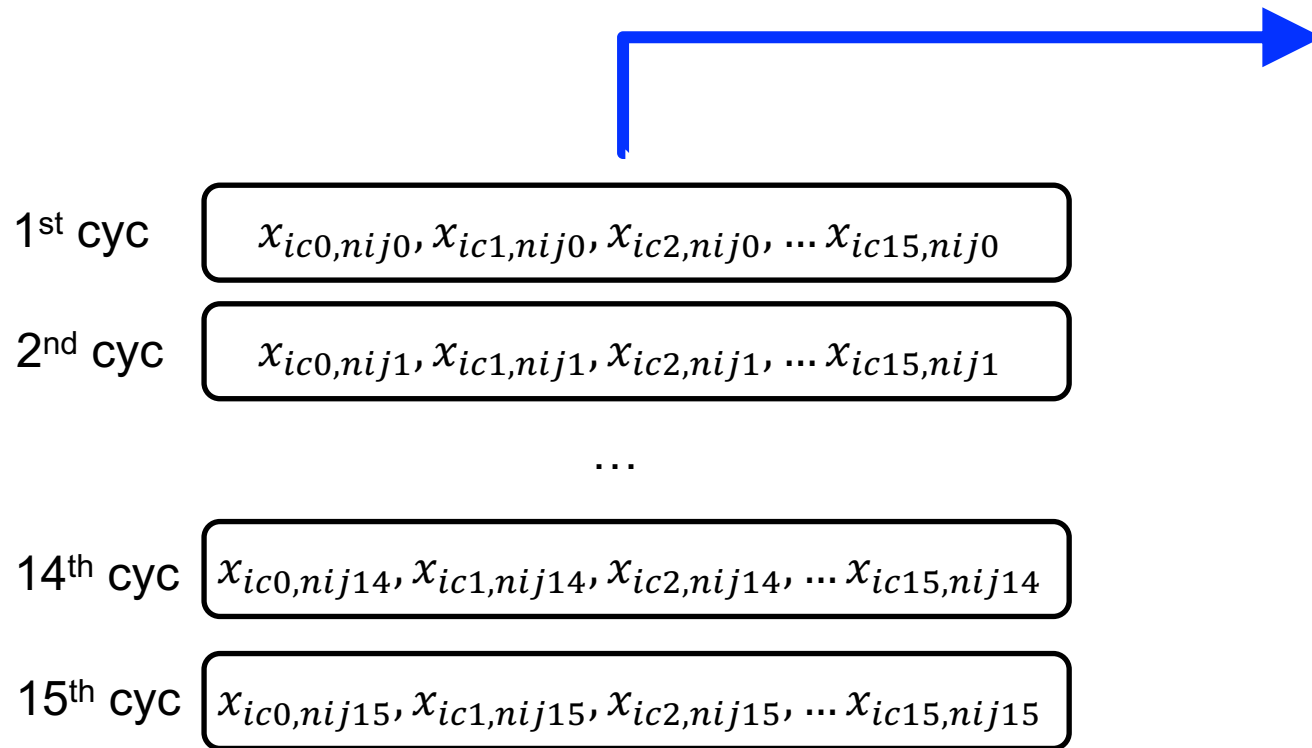


Instruction Flow



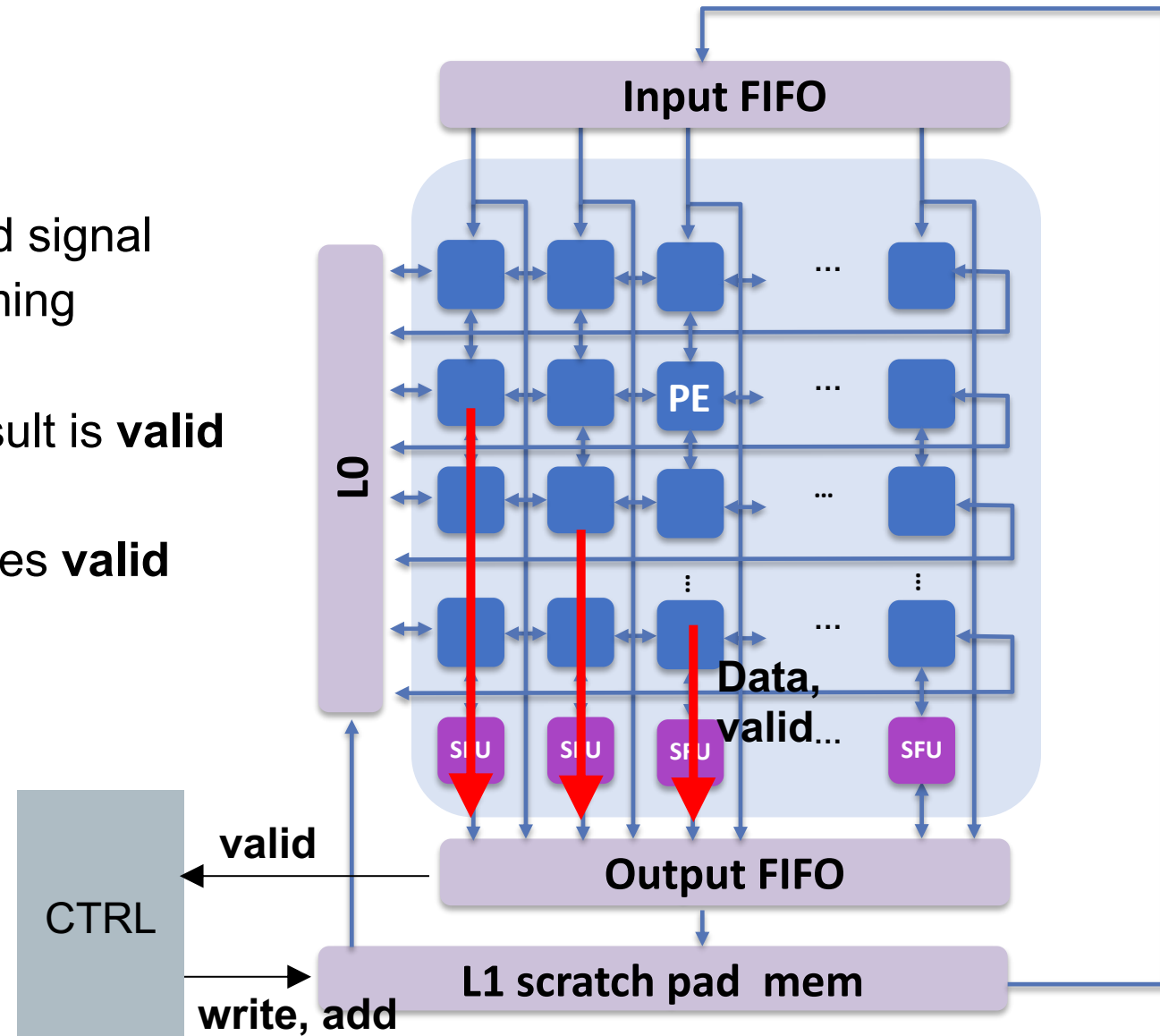
L0 Loading and Execution

- L0 is also a FIFO similar to Input FIFO (last row has max depth)
- Data is received at a time, but first row outputs first
- L0 loading and execution can be processed simultaneously
- Execution instruction also flows similar to kernel loading



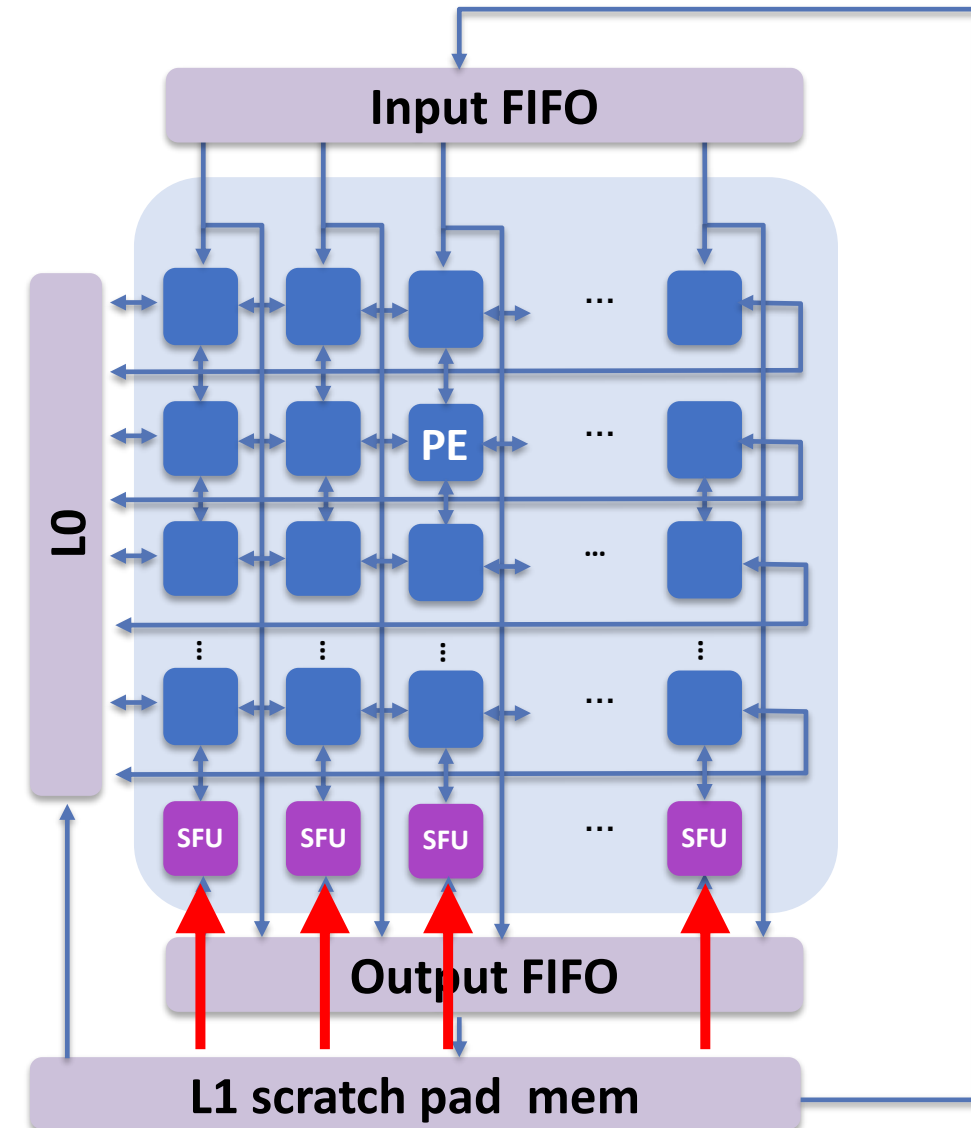
Output FIFO to L1 Scratch pad

- Data from each column sends the data and valid signal
- Output fifo (OFIFO) receives data at different timing from each column
- But, outputs at a time once the last column's result is **valid**
- Thus, OFIFO's first column requires more depth
- Once OFIFO has a complete row of data, it issues **valid** signal to CTRL, and write into L1 scratch pad

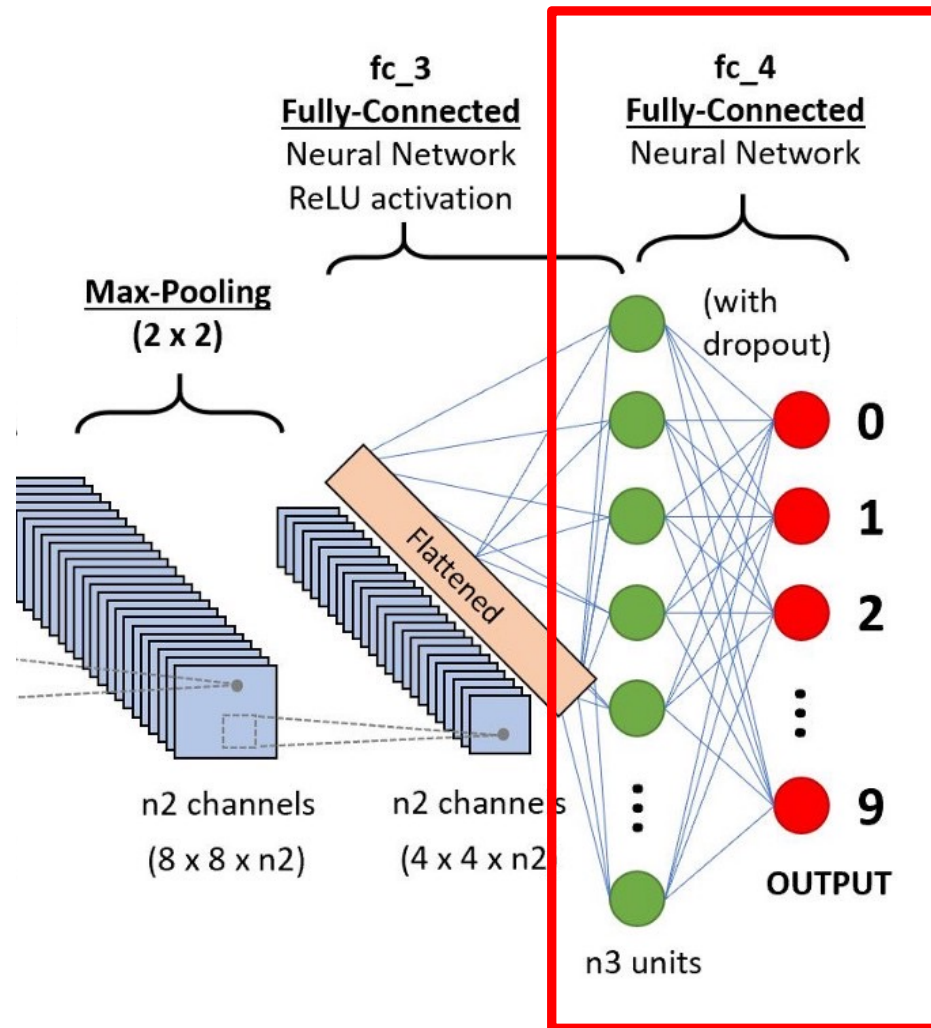


Accumulation (or Activation, e.g., ReLU)

- 1D special function unit (SFU) receives the data from L1 mem.
- Receives 9 consecutive vectors (for 3 X 3 kernels) and accumulate in SFP.
- Then, send back to L1 scratch pad mem.
- Similarly, activation functions (e.g., ReLU) can be processed right after the accumulation if needed



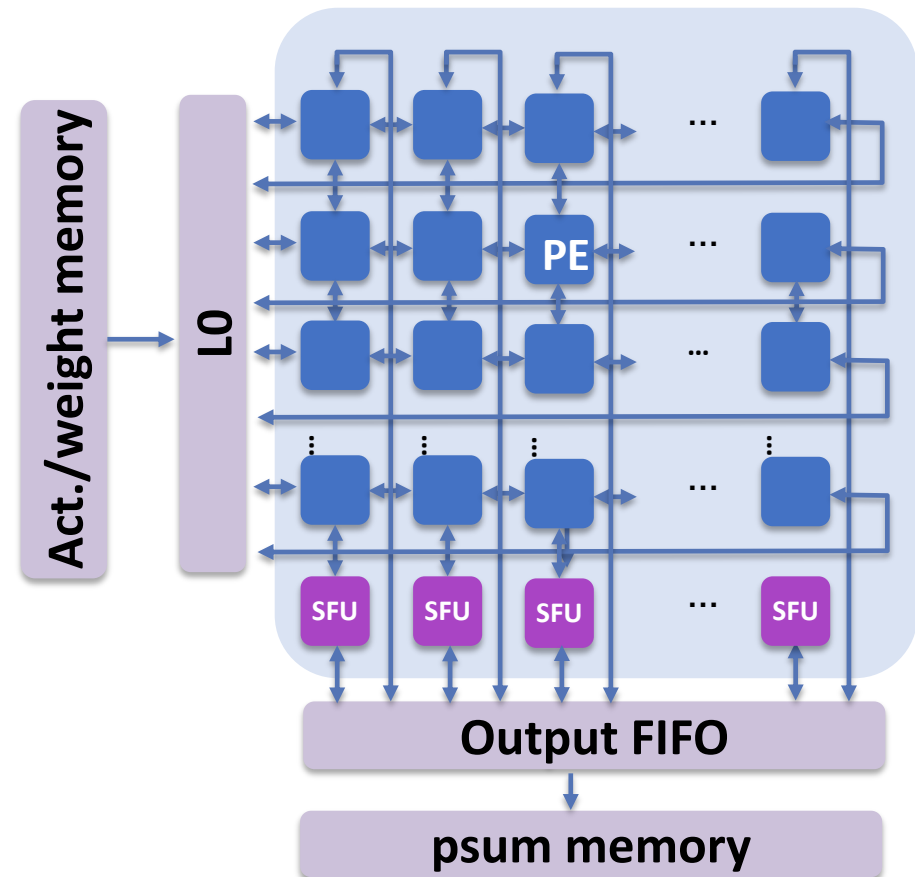
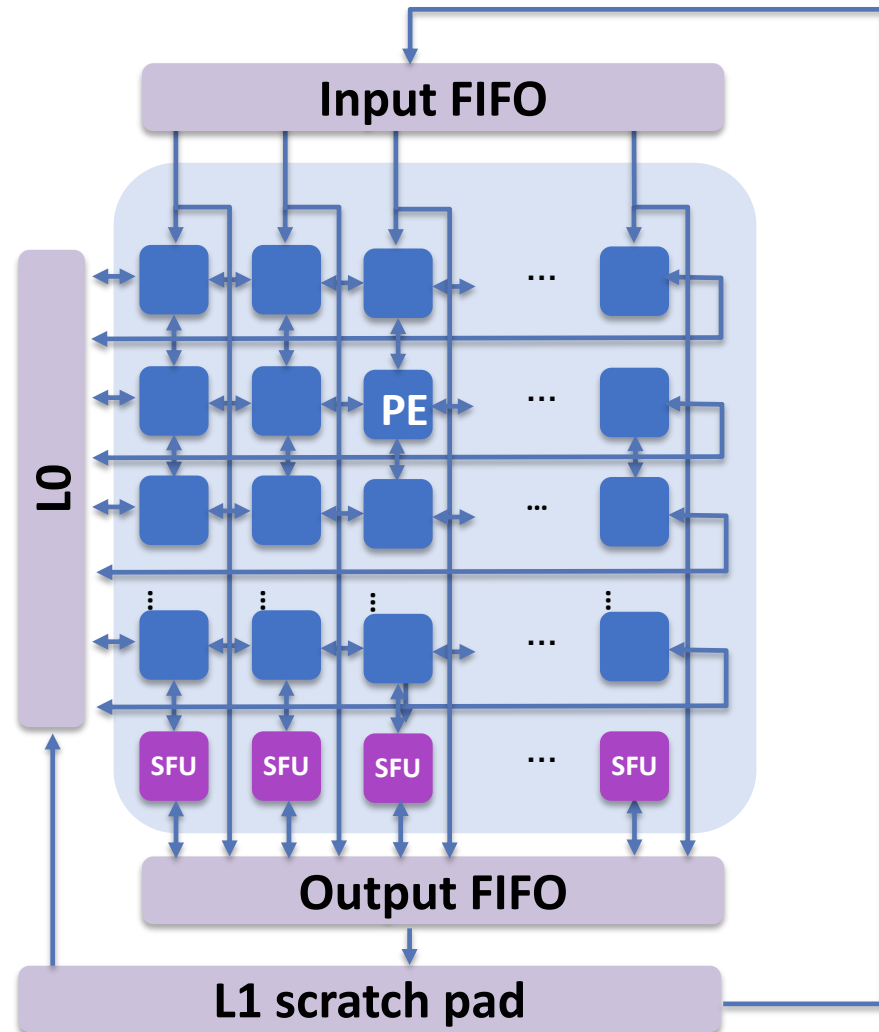
Fully-Connected Layer Computation



$$\begin{bmatrix} x_{ic0} & x_{ic1} & \dots & x_{ic255} \end{bmatrix} * \begin{bmatrix} W_{ic0,oc0} & W_{0,1} & \dots & W_{0,9} \\ W_{ic1,oc0} & W_{1,1} & \dots & W_{1,9} \\ \dots & \dots & \dots & \dots \\ W_{ic255,oc0} & W_{255,1} & \dots & W_{255,9} \end{bmatrix}$$

- Vector \times matrix computation (instead of mat \times mat)
- Weight is arranged in the array as above matrix
- No notion of n_{ij} , thus no data reuse
- Can be reused when batch-size > 1
- The kernel weight is used just once and discarded

Activation and PSUM Memory



L1 Scratch Pad Memory Capacity

- Assume that a layer should be computed without DRAM communication in the middle
- Weight / Activation: 4-bit, psum / output: 16-bit
- e.g., VGGNet 4-th layer
 $n_{ij} = 32 \times 32 = 1024$, $i_c = 64$, $o_c = 64$, $k_{ij} = 9$
- Padding not considered for simplicity
- Weight: $64 (i_c) \times 64 (o_c) \times 9 (k_{ij}) \times 4\text{bits} = 18 \text{ KB}$
- Activation: $64 (i_c) \times 1024 (n_{ij}) \times 4\text{bits} = 32 \text{ KB}$
- psum: $64 (o_c) \times 1024 \times 9 (k_{ij}) \times 16\text{bits} = 1152 \text{ KB}$ **Execute & acc parallel processing case**
=> psum: $64 (o_c) \times 1024 \times 1 \times 16\text{bits} = 128 \text{ KB}$
- fout: $64 (o_c) \times 1024 (n_{ij}) \times 16\text{bits} = 128 \text{ KB}$
- Roughly, < 1.2MB required assuming some data can be overwritten (e.g., psum, activation)