

## 12 BONUS: Systolic Arrays [35 points]

A systolic array consists of 4x4 Processing Elements (PEs), interconnected as shown in Figure 1. The inputs of the systolic array are labeled as  $H_0, H_1, H_2, H_3$  and  $V_0, V_1, V_2, V_3$ . Figure 2 shows the PE logic, which performs a multiply and accumulate MAC operation and saves the result to an internal register (*reg*). Figure 2 also shows how each PE propagates its inputs. We make the following assumptions:

- The latency of each MAC operation is one cycle.
- The propagation of the values from  $i_0$  to  $o_0$ , and from  $i_1$  to  $o_1$ , takes one cycle.
- The initial values of all internal registers is zero.

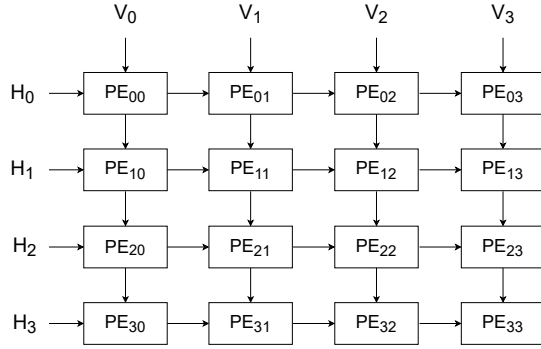


Figure 1: PE array

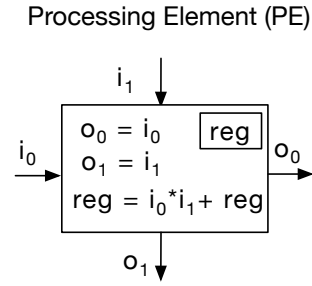


Figure 2: Processing Element (PE)

Your goal is to use the example systolic array shown in Figure 1 to perform the convolution ( $\otimes$ ) of a 3x3 image (matrix  $I_{3 \times 3}$ ) with four 2x2 filters (matrices  $A_{2 \times 2}$ ,  $B_{2 \times 2}$ ,  $C_{2 \times 2}$ , and  $D_{2 \times 2}$ ), to obtain four 2x2 outputs (matrices  $W_{2 \times 2}$ ,  $X_{2 \times 2}$ ,  $Y_{2 \times 2}$ , and  $Z_{2 \times 2}$ ):

$$\begin{matrix} I_{00} & I_{01} & I_{02} \\ I_{10} & I_{11} & I_{12} \\ I_{20} & I_{21} & I_{22} \end{matrix} \quad (\otimes) \quad \begin{matrix} A_{00} & A_{01} \\ A_{10} & A_{11} \end{matrix} = \begin{matrix} W_{00} & W_{01} \\ W_{10} & W_{11} \end{matrix}$$

$$\begin{matrix} I_{00} & I_{01} & I_{02} \\ I_{10} & I_{11} & I_{12} \\ I_{20} & I_{21} & I_{22} \end{matrix} \quad (\otimes) \quad \begin{matrix} B_{00} & B_{01} \\ B_{10} & B_{11} \end{matrix} = \begin{matrix} X_{00} & X_{01} \\ X_{10} & X_{11} \end{matrix}$$

$$\begin{matrix} I_{00} & I_{01} & I_{02} \\ I_{10} & I_{11} & I_{12} \\ I_{20} & I_{21} & I_{22} \end{matrix} \quad (\otimes) \quad \begin{matrix} C_{00} & C_{01} \\ C_{10} & C_{11} \end{matrix} = \begin{matrix} Y_{00} & Y_{01} \\ Y_{10} & Y_{11} \end{matrix}$$

$$\begin{matrix} I_{00} & I_{01} & I_{02} \\ I_{10} & I_{11} & I_{12} \\ I_{20} & I_{21} & I_{22} \end{matrix} \quad (\otimes) \quad \begin{matrix} D_{00} & D_{01} \\ D_{10} & D_{11} \end{matrix} = \begin{matrix} Z_{00} & Z_{01} \\ Z_{10} & Z_{11} \end{matrix}$$

As an example, the convolution of the matrix  $I_{3 \times 3}$  with the filter  $A_{2 \times 2}$  is computed as follows:

- $W_{00} = I_{00} * A_{00} + I_{01} * A_{01} + I_{10} * A_{10} + I_{11} * A_{11}$
- $W_{01} = I_{01} * A_{00} + I_{02} * A_{01} + I_{11} * A_{10} + I_{12} * A_{11}$
- $W_{10} = I_{10} * A_{00} + I_{11} * A_{01} + I_{20} * A_{10} + I_{21} * A_{11}$
- $W_{11} = I_{11} * A_{00} + I_{12} * A_{01} + I_{21} * A_{10} + I_{22} * A_{11}$

You should compute the four convolutions in the minimum possible number of cycles. Fill the following table with:

1. The input elements (from matrices  $I_{3 \times 3}$ ,  $A_{2 \times 2}$ ,  $B_{2 \times 2}$ ,  $C_{2 \times 2}$ , and  $D_{2 \times 2}$ ) in the correct input ports of the systolic array ( $H_0$ ,  $H_1$ ,  $H_2$ ,  $H_3$  and  $V_0$ ,  $V_1$ ,  $V_2$ ,  $V_3$ ). (Hint: If necessary, an input element can be concurrently streamed into several input ports of the array.)
2. The output values and the corresponding PE where the output elements (of matrices  $W_{2 \times 2}$ ,  $X_{2 \times 2}$ ,  $Y_{2 \times 2}$ , and  $Z_{2 \times 2}$ ) are generated.

Fill the blanks only with relevant information.

cycle	H0	H1	H2	H3	V0	V1	V2	V3	PE <sub>00</sub>	PE <sub>01</sub>	PE <sub>02</sub>	PE <sub>03</sub>	PE <sub>10</sub>	PE <sub>11</sub>	PE <sub>12</sub>	PE <sub>13</sub>	PE <sub>20</sub>	PE <sub>21</sub>	PE <sub>22</sub>	PE <sub>23</sub>	PE <sub>30</sub>	PE <sub>31</sub>	PE <sub>32</sub>	PE <sub>33</sub>
0	A <sub>00</sub>				I <sub>00</sub>																			
1	A <sub>01</sub>	B <sub>00</sub>			I <sub>01</sub>	I <sub>01</sub>																		
2	A <sub>10</sub>	B <sub>01</sub>	C <sub>00</sub>		I <sub>10</sub>	I <sub>02</sub>	I <sub>10</sub>																	
3	A <sub>11</sub>	B <sub>10</sub>	C <sub>01</sub>	D <sub>00</sub>	I <sub>11</sub>	I <sub>11</sub>	I <sub>11</sub>	I <sub>11</sub>	W <sub>00</sub>															
4		B <sub>11</sub>	C <sub>10</sub>	D <sub>01</sub>		I <sub>12</sub>	I <sub>20</sub>	I <sub>12</sub>		W <sub>01</sub>			X <sub>00</sub>											
5			C <sub>11</sub>	D <sub>10</sub>			I <sub>21</sub>	I <sub>21</sub>			W <sub>10</sub>			X <sub>01</sub>			Y <sub>00</sub>							
6				D <sub>11</sub>				I <sub>22</sub>				W <sub>11</sub>			X <sub>10</sub>			Y <sub>01</sub>			Z <sub>00</sub>			
7																X <sub>11</sub>			Y <sub>10</sub>			Z <sub>01</sub>		
8																				Y <sub>11</sub>			Z <sub>10</sub>	
9																								Z <sub>11</sub>
10																								
11																								
12																								
13																								
14																								
15																								