

## 8 Systolic Arrays [30 points]

A systolic array consists of 3x4 Processing Elements (PEs), interconnected as shown in Figure 1. The inputs of the systolic array are labeled as H0, H1, H2 and V0,V1,V2,V3. Figure 2 shows the PE logic, which performs a multiply and accumulate operation (MAC), and it saves the result in an internal register (reg). Figure 2 also shows how each PE propagates its inputs. We make the following assumptions:

- The latency of each MAC 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 value of all registers is zero.
- You can input a value more than once in the systolic array.

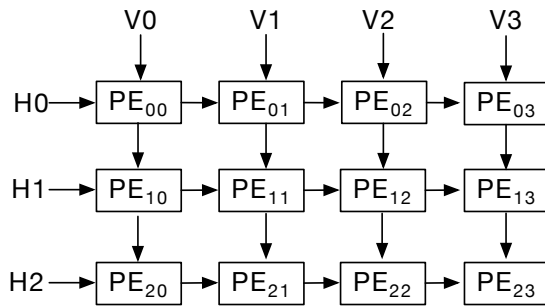


Figure 1: PE array

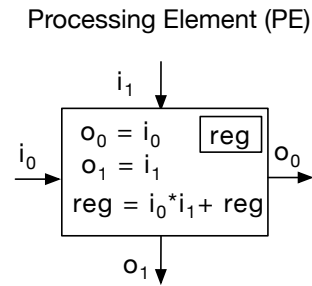


Figure 2: Processing Element (PE)

Your goal is to use this systolic array to perform the convolution of a 3x3 image (matrix I) with three 2x2 filters (matrices F, G, and H), to obtain three outputs (matrices O, U, and E):

$$\begin{matrix} I_{00} & I_{01} & I_{02} \\ I_{10} & I_{11} & I_{12} \\ I_{20} & I_{21} & I_{22} \end{matrix} \otimes \begin{matrix} F_{00} & F_{01} \\ F_{10} & F_{11} \end{matrix} = \begin{matrix} O_{00} & O_{01} \\ O_{10} & O_{11} \end{matrix}$$

$$\begin{matrix} I_{00} & I_{01} & I_{02} \\ I_{10} & I_{11} & I_{12} \\ I_{20} & I_{21} & I_{22} \end{matrix} \otimes \begin{matrix} G_{00} & G_{01} \\ G_{10} & G_{11} \end{matrix} = \begin{matrix} U_{00} & U_{01} \\ U_{10} & U_{11} \end{matrix}$$

$$\begin{matrix} I_{00} & I_{01} & I_{02} \\ I_{10} & I_{11} & I_{12} \\ I_{20} & I_{21} & I_{22} \end{matrix} \otimes \begin{matrix} H_{00} & H_{01} \\ H_{10} & H_{11} \end{matrix} = \begin{matrix} E_{00} & E_{01} \\ E_{10} & E_{11} \end{matrix}$$

As an example, the convolution of the matrix I with the filter F is computed as follows:

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

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

1. The input values (matrices I, F, G, and H) in the correct input ports of the systolic array (the values can be repeated).
2. The output values and the corresponding PE where the outputs (matrices O, U, and E) are generated.

Fill the gaps only with relevant information.

cycle	H0	H1	H2	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>
0	F <sub>00</sub>			I <sub>00</sub>															
1	F <sub>01</sub>	G <sub>00</sub>		I <sub>01</sub>	I <sub>01</sub>														
2	F <sub>10</sub>	G <sub>01</sub>	H <sub>00</sub>	I <sub>10</sub>	I <sub>02</sub>	I <sub>10</sub>													
3	F <sub>11</sub>	G <sub>10</sub>	H <sub>01</sub>	I <sub>11</sub>	I <sub>11</sub>	I <sub>11</sub>	I <sub>11</sub>	O <sub>00</sub>											
4		G <sub>11</sub>	H <sub>10</sub>		I <sub>12</sub>	I <sub>20</sub>	I <sub>12</sub>		O <sub>01</sub>			U <sub>00</sub>							
5			H <sub>11</sub>			I <sub>21</sub>	I <sub>21</sub>			O <sub>10</sub>			U <sub>01</sub>			E <sub>00</sub>			
6							I <sub>22</sub>				O <sub>11</sub>			U <sub>10</sub>			E <sub>01</sub>		
7															U <sub>11</sub>			E <sub>10</sub>	
8																			E <sub>11</sub>
9																			
10																			
11																			
12																			
13																			
14																			
15																			