

## JO Approach

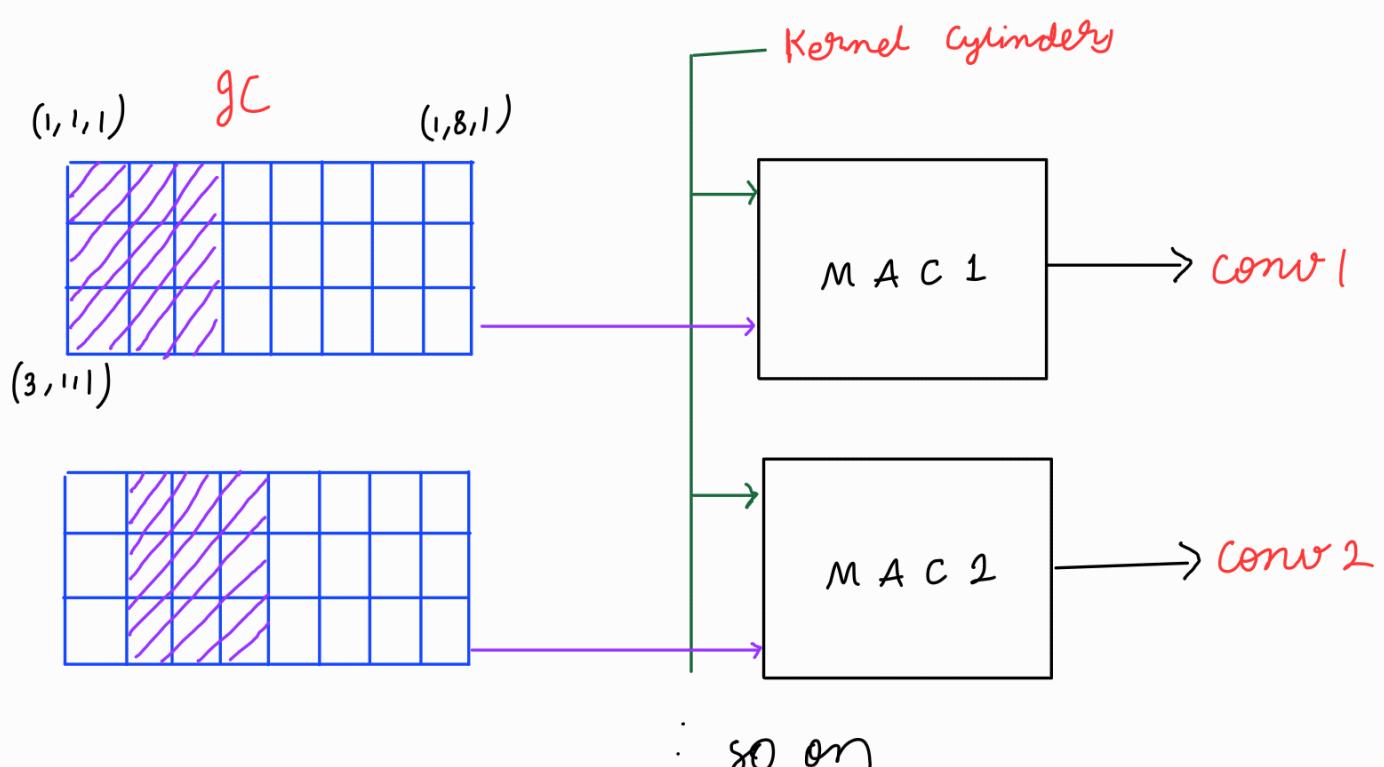
### MAC Architecture :-

Let us use some numbering system for input cylinders

	Col 1	Col 2	Col 3	Col 4	Col 5	Col 6	Col 7	Col 8
Row 1	Cy 1	Cy 4	Cy 7	Cy 10	Cy 13	Cy 16	Cy 19	Cy 22
Row 2	Cy 2	Cy 5	Cy 8	Cy 11	Cy 14	Cy 17	Cy 20	Cy 23
Row 3	Cy 3	Cy 6	Cy 9	Cy 12	Cy 15	Cy 18	Cy 21	Cy 24

\* Cy 1 = Input Cylinder 1 = Input[ row = 1, col = 1, chn = 1:8 ]

→ Recall from our "Convolution Approach"



So, MAC 1 requires Cy 1 to Cy 9

MAC 2 requires Cy 3 to Cy 12

... so on

→ Let us assume we need 1 clock cycle to get a cylinder.

So, Cy 1 arrives at  $t = 1$

Cy 2 arrives at  $t = 2$

:

Cy  $\infty$  arrives at  $t = \infty$

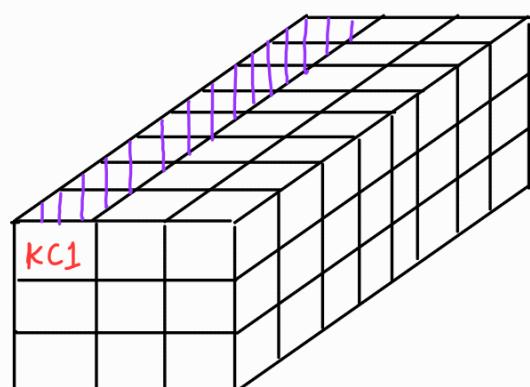
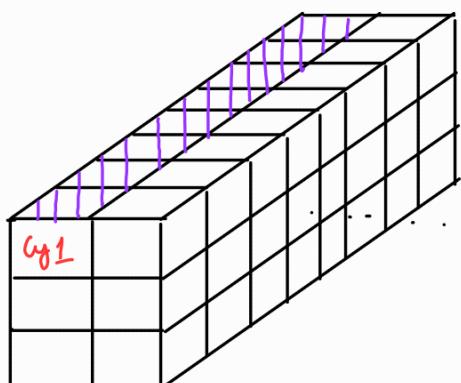
→ Also, no. of clock cycles required  
for half-precision multiplier } = 2 cycles

no. of clock cycles required  
for half-precision adder } = 1 cycle

### i) Multiplier and Adder structure

→ Multiplier and adder structure for convolution between "1" Input cylinder and "1" Kernel cylinder

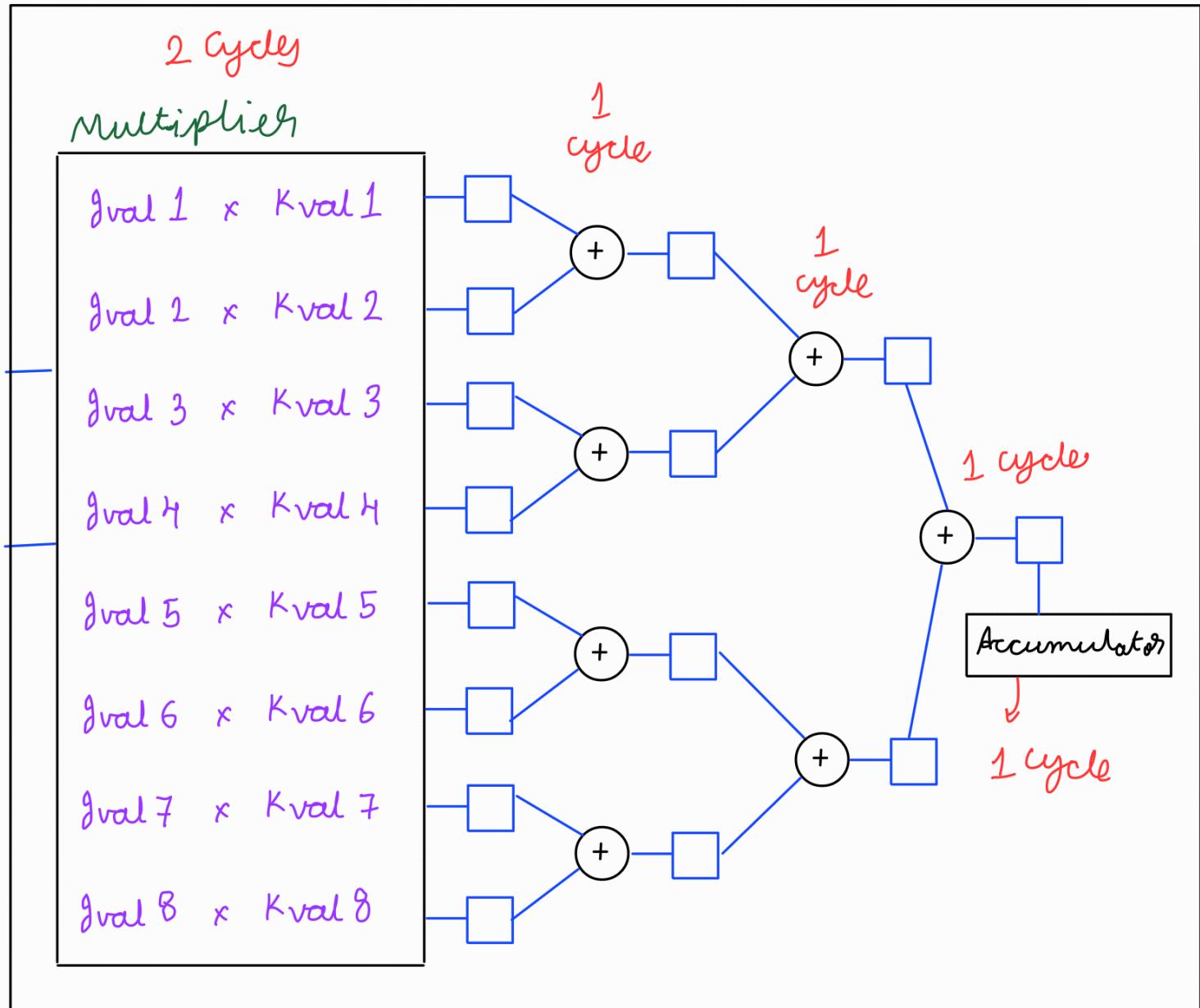
$$E_1 = Cy_1 * KC_1$$



Let  $Cy1 = [gval1, gval2, gval3, gval4, gval5, gval6, gval7, gval8]$

$KC1 = [Kval1, Kval2, Kval3, Kval4, Kval5, Kval6, Kval7, Kval8]$

MAC



*Mul*

*Add Acc*

*Mul* = 2 cycles

*Add Acc* = 4 cycles

→ For MAC 1, the cylinders Cy 1 to Cy 9 will arrive one after another, so instead of starting MAC 1 after waiting for 9 cylinders, we can pipeline the convolution of 9 cylinders. This will reduce no. of multipliers per MAC to 8 and no. of adders per MAC to 8

$t = 0$  Start

$t = 1$  Cy 1

$t = 2$  Cy 2

$t = 3$  Cy 3 Mul 1

$t = 4$  Cy 4 Mul 2

$t = 5$  Cy 5 Mul 3

$t = 6$  Cy 6 Mul 4

$t = 7$  Cy 7 Mul 5 Add Acc 1

$t = 8$  Cy 8 Mul 6 Add Acc 2

$t = 9$  Cy 9 Mul 7 Add Acc 3

$t = 10$  Mul 8 Add Acc 4

$t = 11$  Mul 9 Add Acc 5

$t = 12$  Add Acc 6

$t = 13$  Add Acc 7

$t = 14$  Add Acc 8

$t = 15$  Add Acc 9

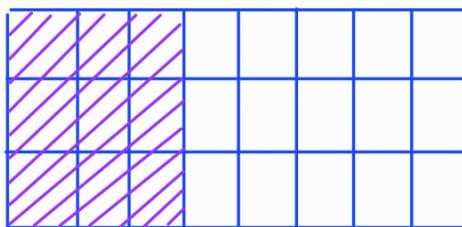
Hence , it takes 15 cycles for a MAC to calculate output once first cylinders arrives .

→ Let us calculate effect of no. of MAC on total cycles required for 36 output values

Case 1 :- 4 MACs

After 1

MAC 1

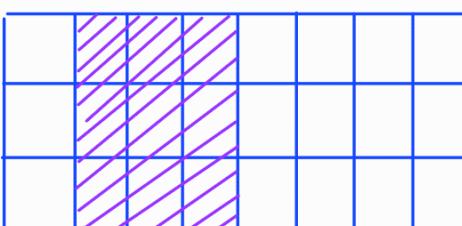


Start                    End

$t = 1$

$t = 15$

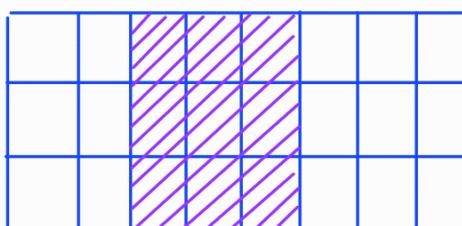
MAC 2



$t = 4$

$t = 18$

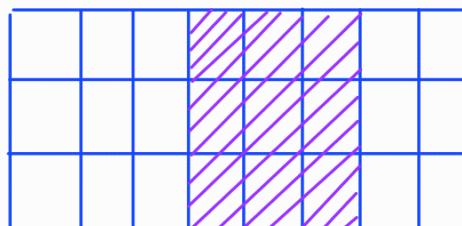
MAC 3



$t = 7$

$t = 21$

MAC 4

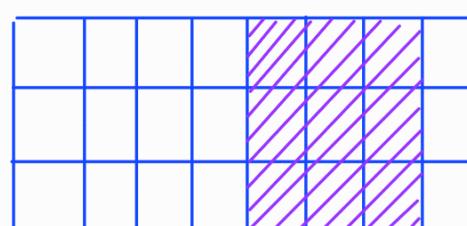


$t = 10$

$t = 24$

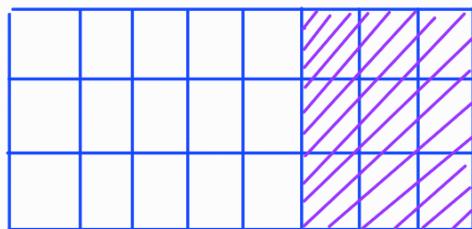
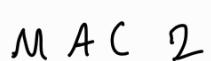
After 2

MAC 1



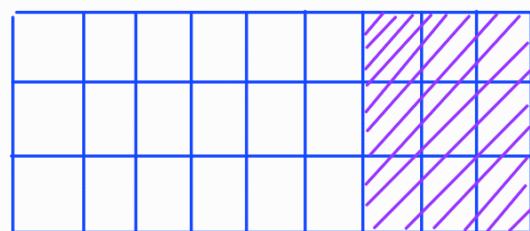
$t = 13$

$t = 27$



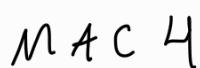
$$t = 16$$

$t = 30$



$$t = 19$$

$$t = 33$$



$$t = 22$$

$$t = 36$$

It is a



t = 106

$$t = 120$$

→ Hence for 4 MACs, it will take 120 cycles to get 36 output values.

Similarly

For 5 MAC's → 117 cycles for 35 output

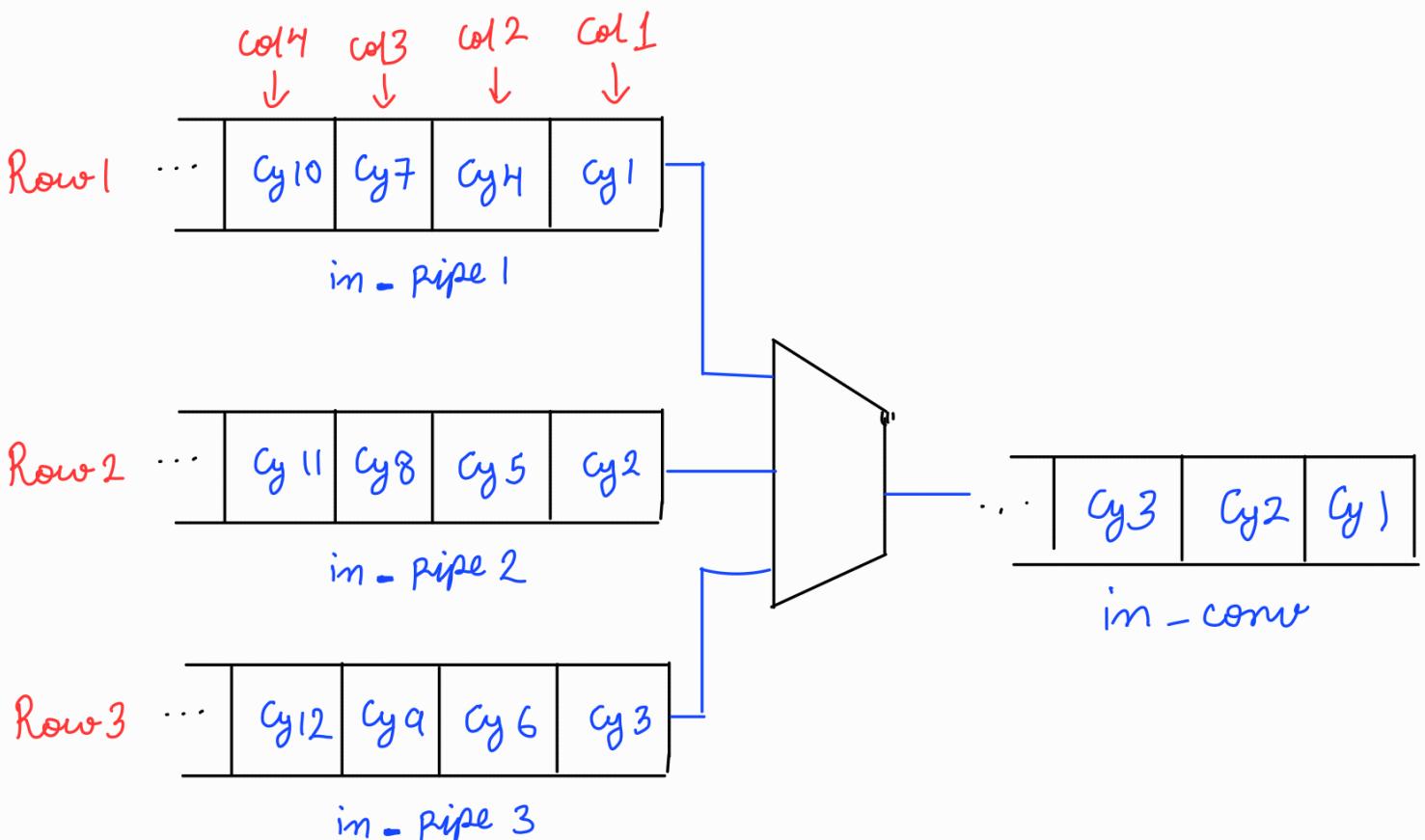
for 6 MAC's  $\rightarrow$  120 cycles for 36 output

For 7 MAC's  $\rightarrow$  117 cycles for 35 output

## Input Module :-

Let us continue with 4 MAC blocks per engine

For storing input, let us use 3 pipes, each pipe store the data for 1 row

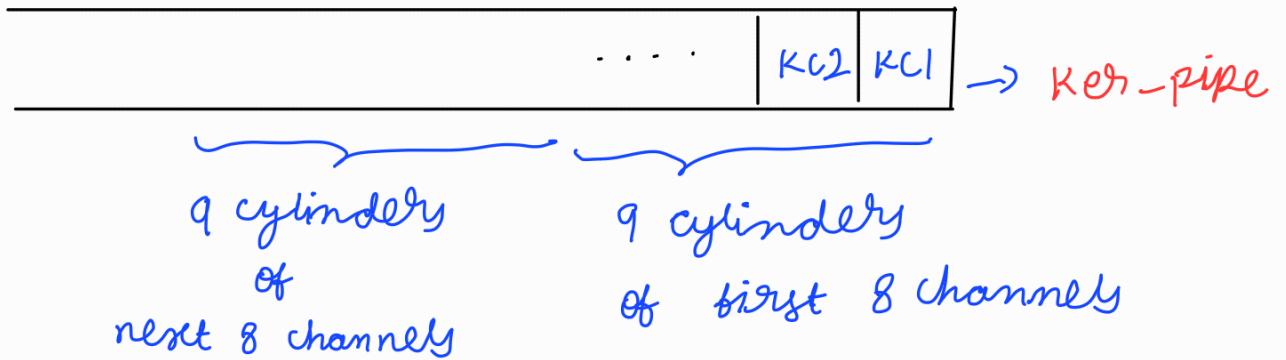


→ A diagrammatic representation for distribution  
of "in-comm" data to MAC buffers

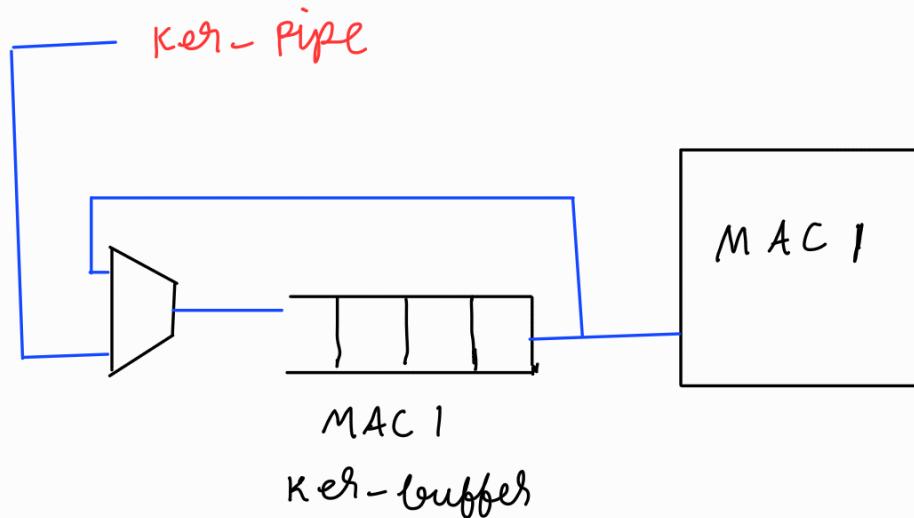
	iter 1	iter 2	iter 3	iter 4	iter 5	iter 6
MAC 1 in-buffer	Cy 1-3	Cy 4-6	Cy 7-9		Cy 13-15	Cy 16-18
MAC 2 in-buffer		Cy 4-6	Cy 7-9	Cy 10-12		Cy 16-18
MAC 3 in-buffer			Cy 7-9	Cy 10-12	Cy 13-15	
MAC 4 in-buffer				Cy 10-12	Cy 13-15	Cy 16-18

## Kernel Module :-

→ Similar to Arman's approach, we will store entire kernel in a pipe.



→ Storage of these kernels in MAC buffers.



## Output Module :-

→ We have output buffer of depth 36, we can send these values directly to the memory.