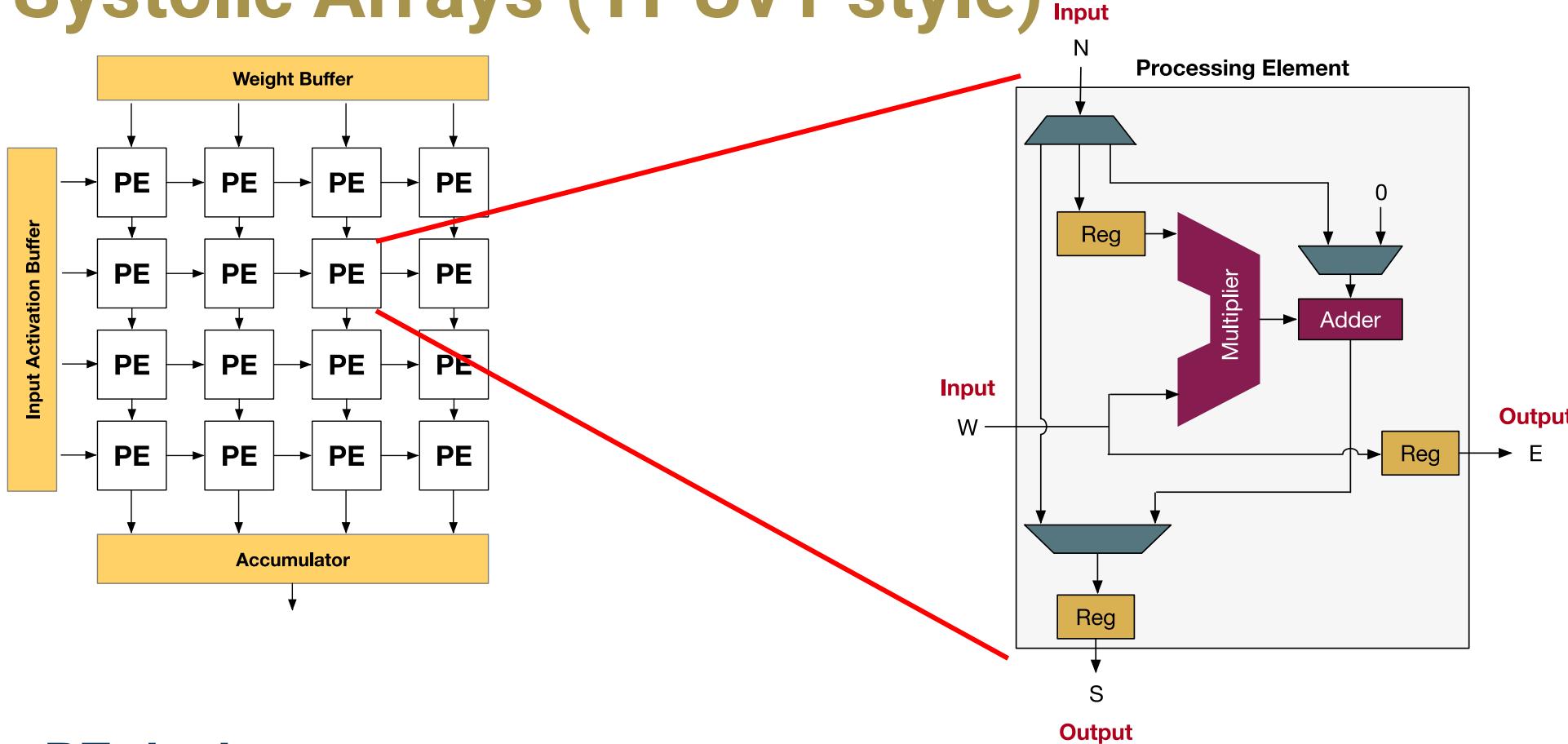
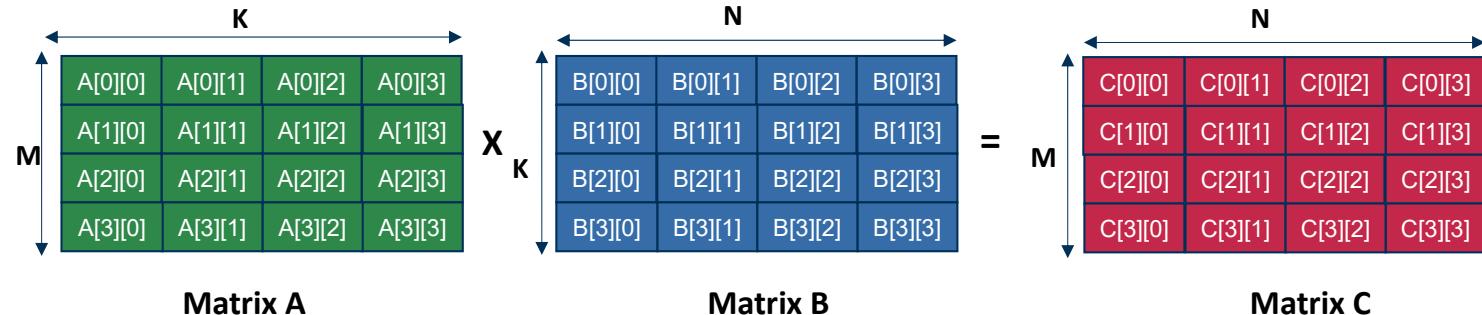
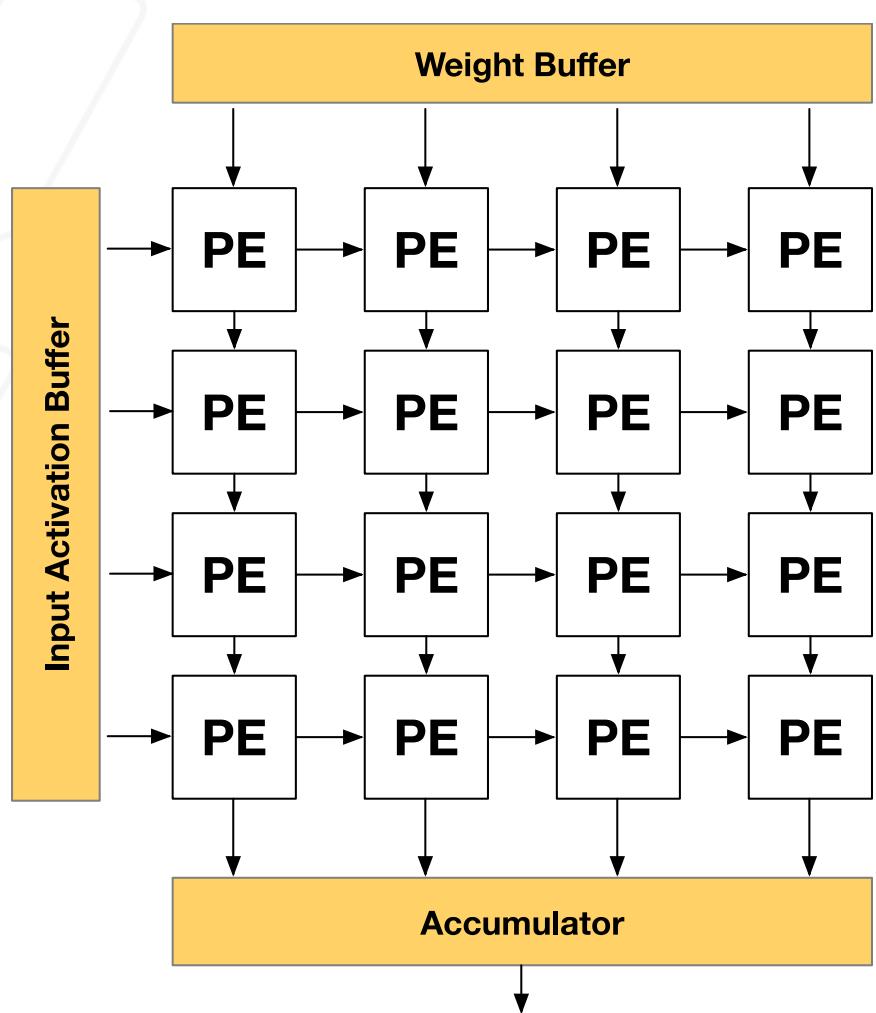


# PEs in Systolic Arrays (TPUv1 style)



- **Simple PE design**
  - No local scratchpad; several one-slot registers
  - Enables to increase the number of PEs (small area for each PE -> more PEs within the same area)

# Matrix Multiplication on a Systolic Array



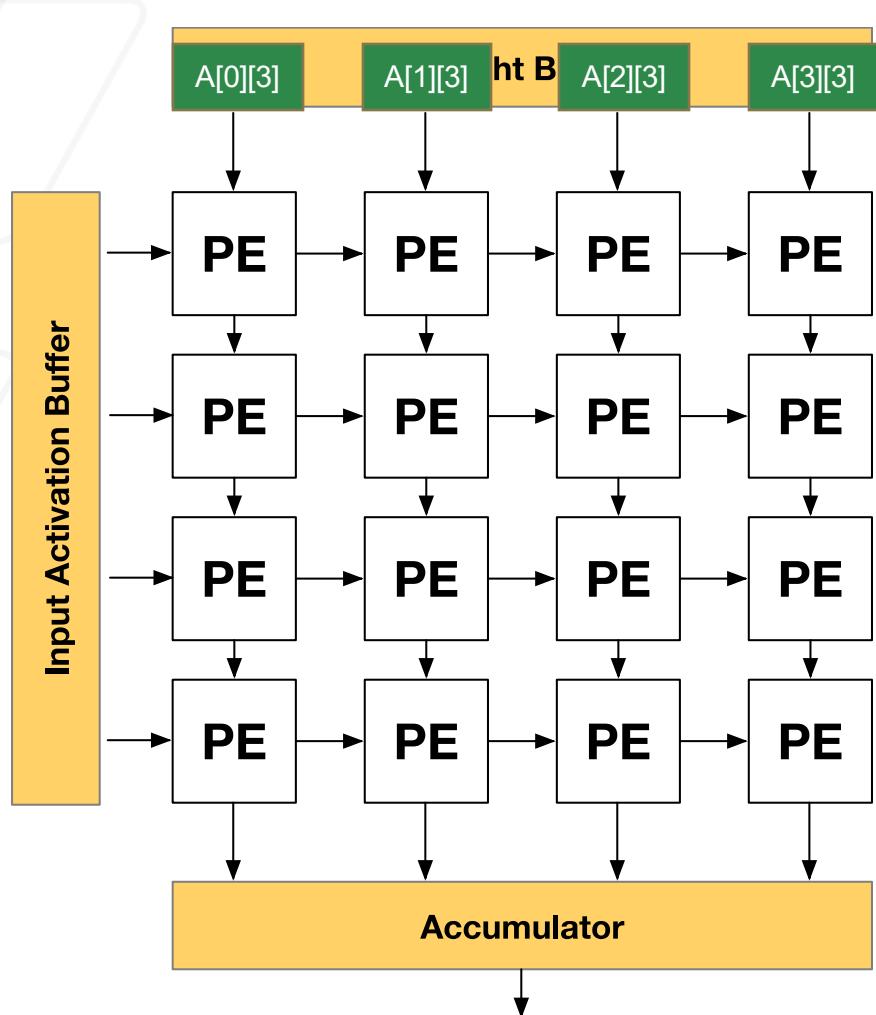
```
for m in range(4):  
    for n in range(4):  
        for k in range(4):  
            C[m][n] += A[m][k] * B[k][n]
```

Assume following “mapping”

- Column: Reduced (aka contracted) Dimension (K)
- Row: Spatial Dimension (either M or N)

Note: alternate mapping styles possible

# Matrix Multiplication on a Systolic Array



Matrix A			
M	K		
A[0][0]	A[0][1]	A[0][2]	A[0][3]
A[1][0]	A[1][1]	A[1][2]	A[1][3]
A[2][0]	A[2][1]	A[2][2]	A[2][3]
A[3][0]	A[3][1]	A[3][2]	A[3][3]

Matrix B			
N	K		
B[0][0]	B[0][1]	B[0][2]	B[0][3]
B[1][0]	B[1][1]	B[1][2]	B[1][3]
B[2][0]	B[2][1]	B[2][2]	B[2][3]
B[3][0]	B[3][1]	B[3][2]	B[3][3]

Matrix C			
N	M		
C[0][0]	C[0][1]	C[0][2]	C[0][3]
C[1][0]	C[1][1]	C[1][2]	C[1][3]
C[2][0]	C[2][1]	C[2][2]	C[2][3]
C[3][0]	C[3][1]	C[3][2]	C[3][3]

Phase 1:  
Load Weight Tensor

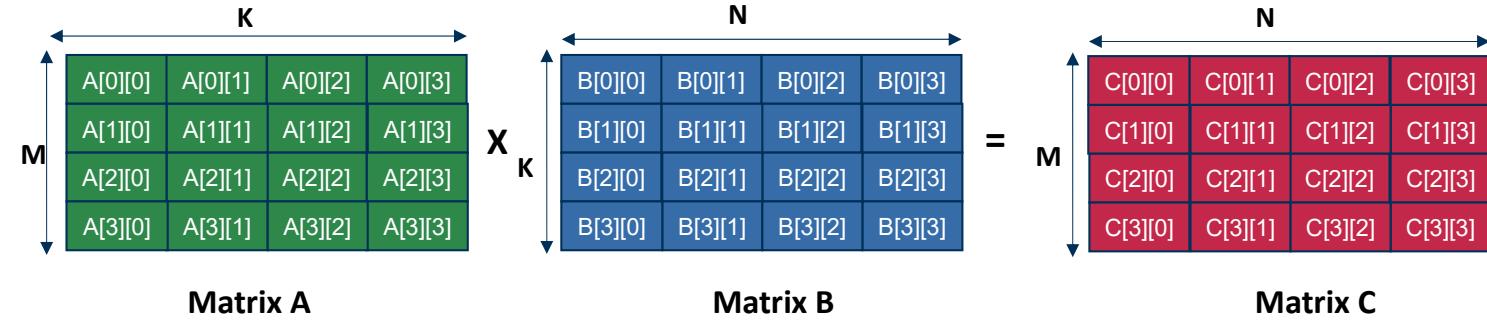
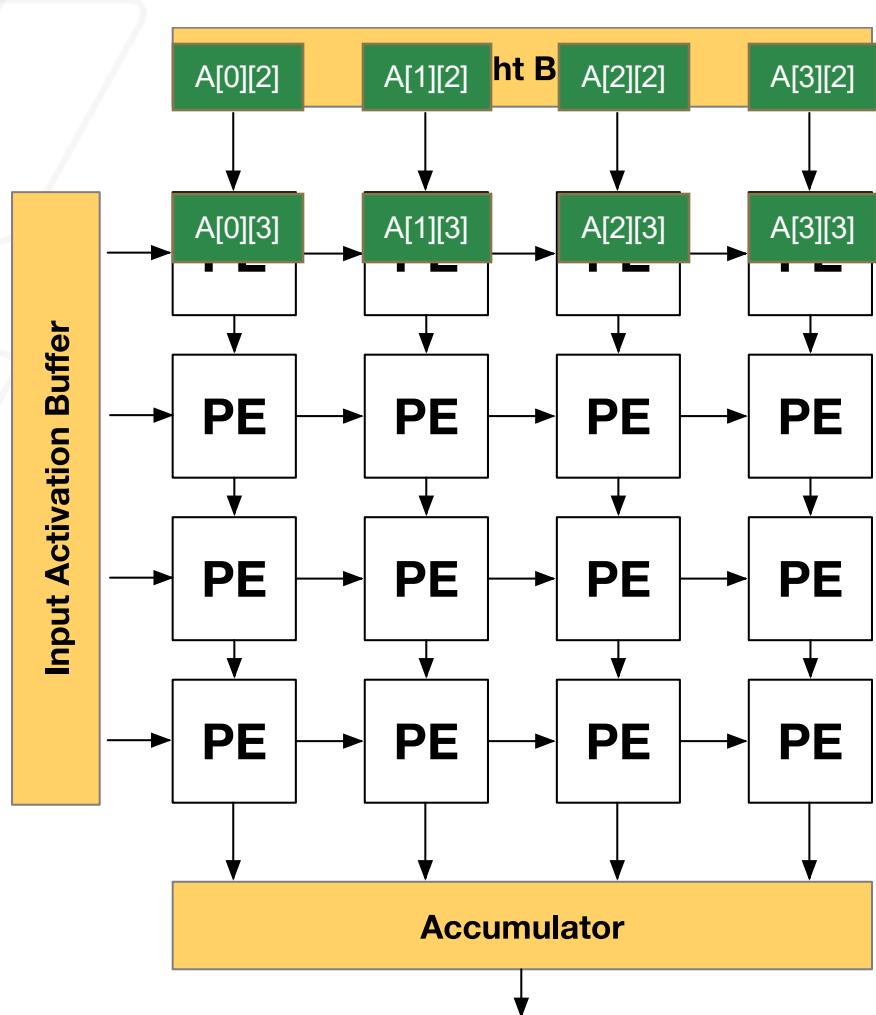
```
for m in range(4):  
    for n in range(4):  
        for k in range(4):  
            C[m][n] += A[m][k] * B[k][n]
```

Assume following “mapping”

- Column: Reduced (aka contracted) Dimension (K)
- Row: Spatial Dimension (either M or N)

Note: alternate mapping styles possible

# Matrix Multiplication on a Systolic Array



```

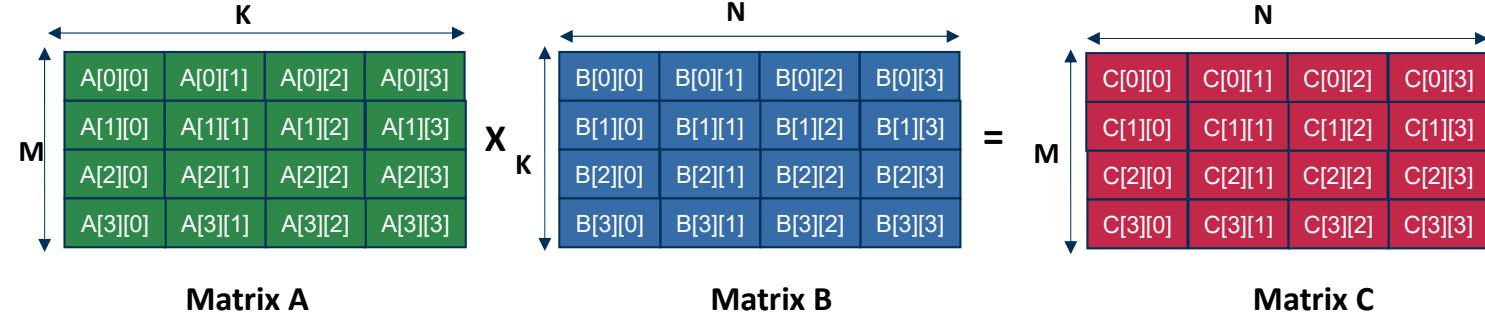
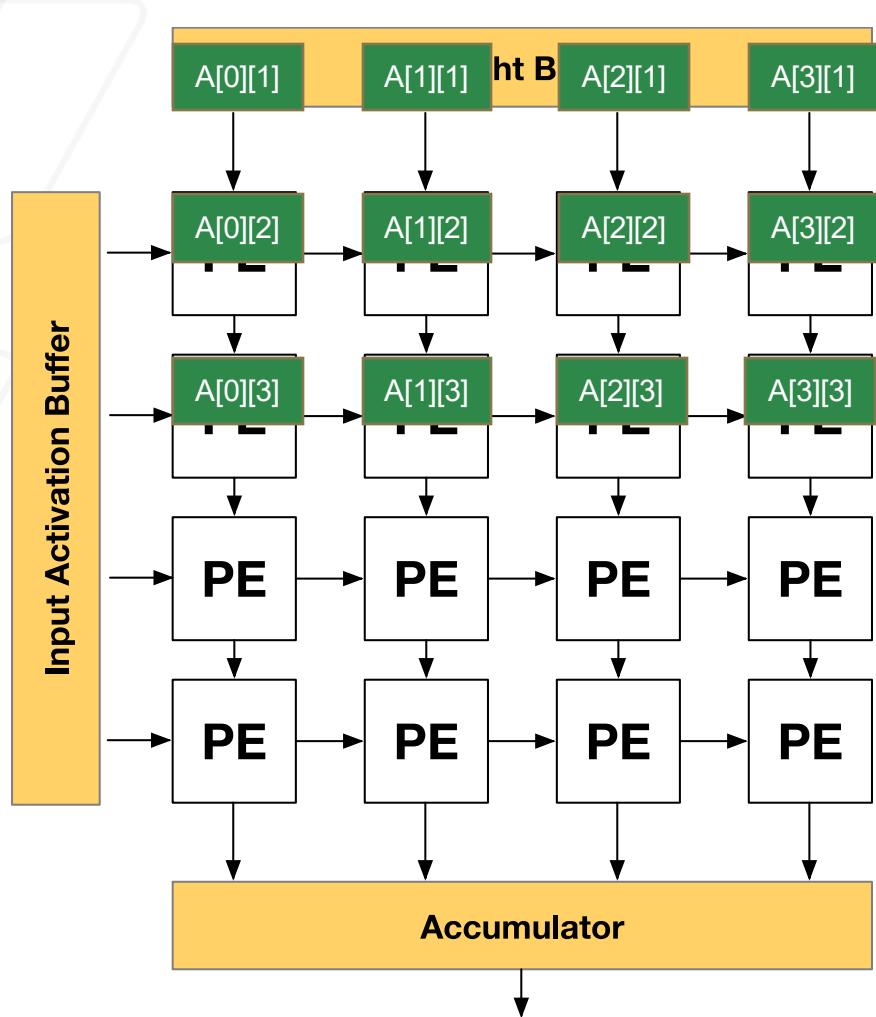
for m in range(4):
    for n in range(4):
        for k in range(4):
            C[m][n] += A[m][k] * B[k][n]
    
```

**Assume following “mapping”**

- Column: Reduced (aka contracted) Dimension (K)
- Row: Spatial Dimension (either M or N)

*Note: alternate mapping styles possible*

# Matrix Multiplication on a Systolic Array



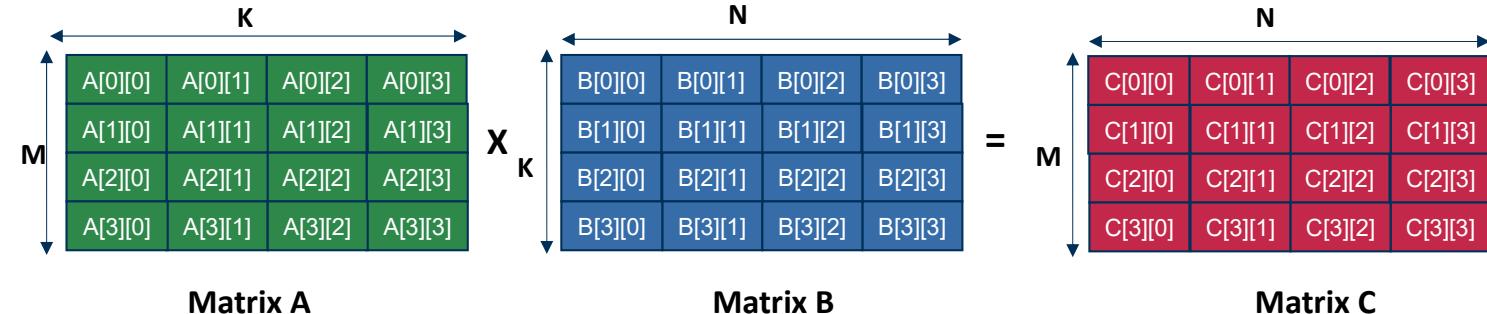
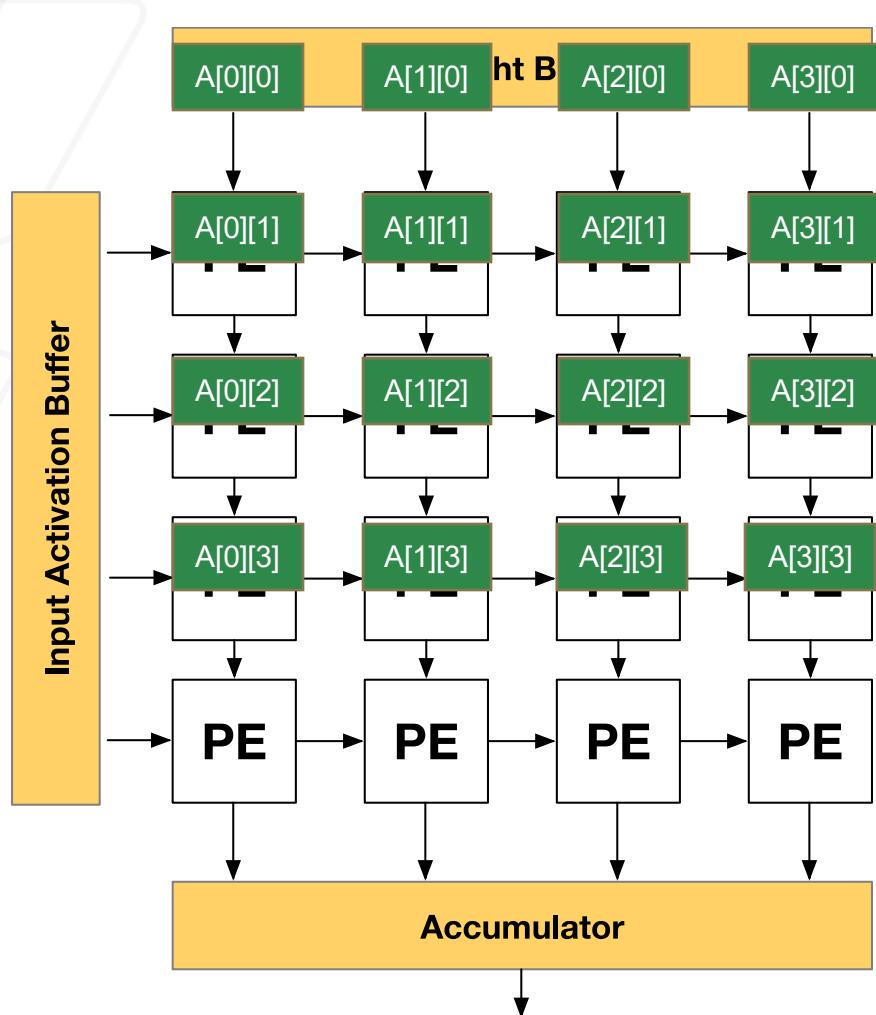
```
for m in range(4):
    for n in range(4):
        for k in range(4):
            C[m][n] += A[m][k] * B[k][n]
```

Assume following “mapping”

- Column: Reduced (aka contracted) Dimension (K)
- Row: Spatial Dimension (either M or N)

Note: alternate mapping styles possible

# Matrix Multiplication on a Systolic Array



**Phase 1:**  
Load Weight Tensor

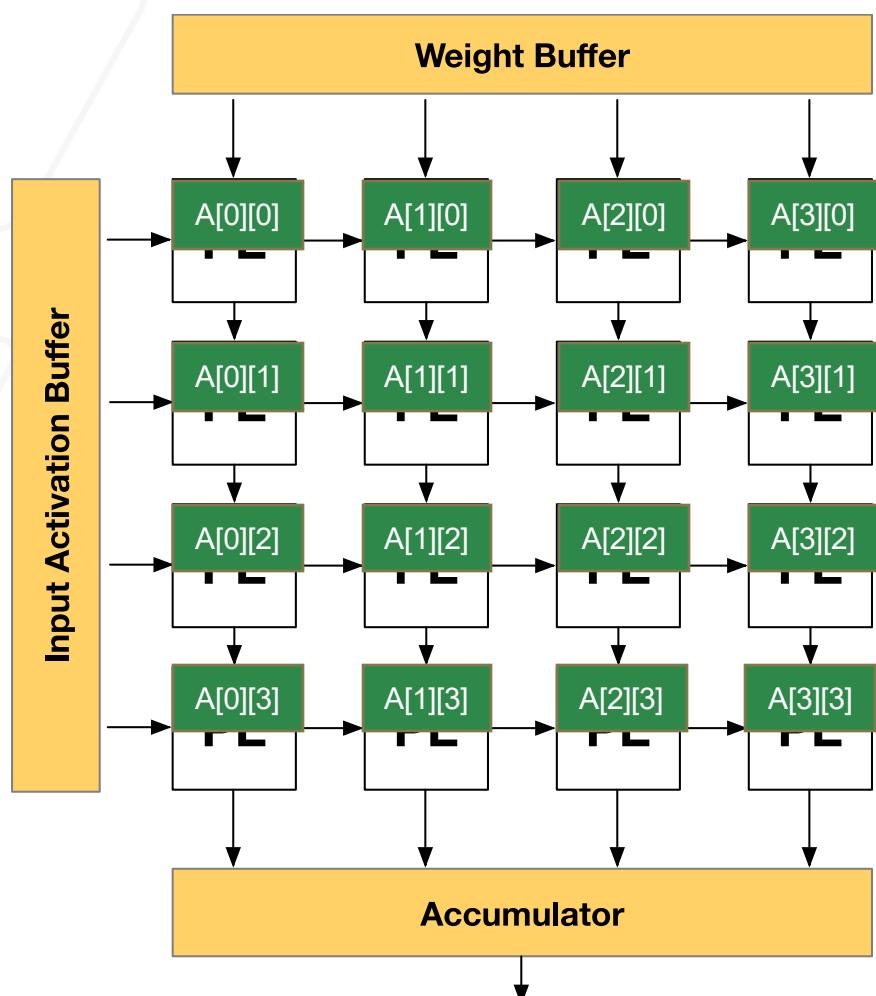
```
for m in range(4):
    for n in range(4):
        for k in range(4):
            C[m][n] += A[m][k] * B[k][n]
```

**Assume following “mapping”**

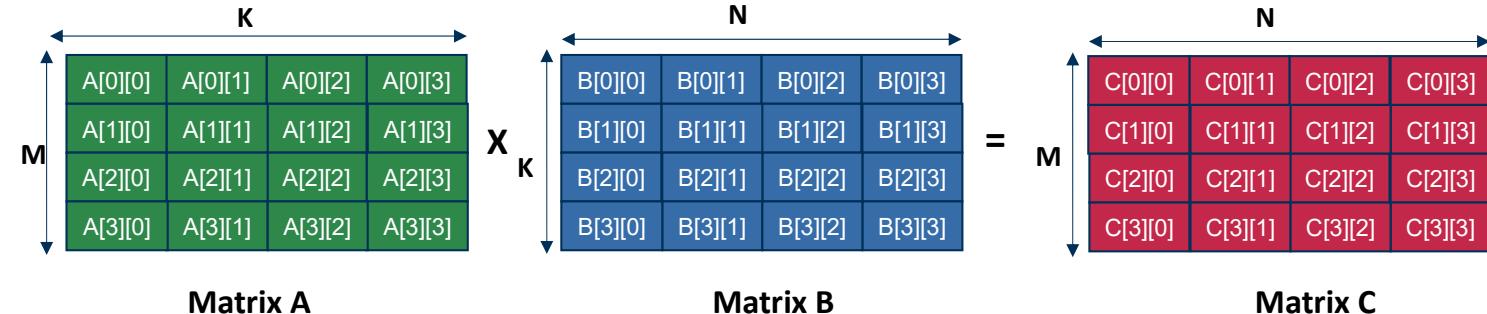
- Column: Reduced (aka contracted) Dimension (K)
- Row: Spatial Dimension (either M or N)

*Note: alternate mapping styles possible*

# Matrix Multiplication on a Systolic Array



Matrix A will now stay resident  
(i.e., “stationary”) in the PEs



```
for m in range(4):  
    for n in range(4):  
        for k in range(4):  
            C[m][n] += A[m][k] * B[k][n]
```

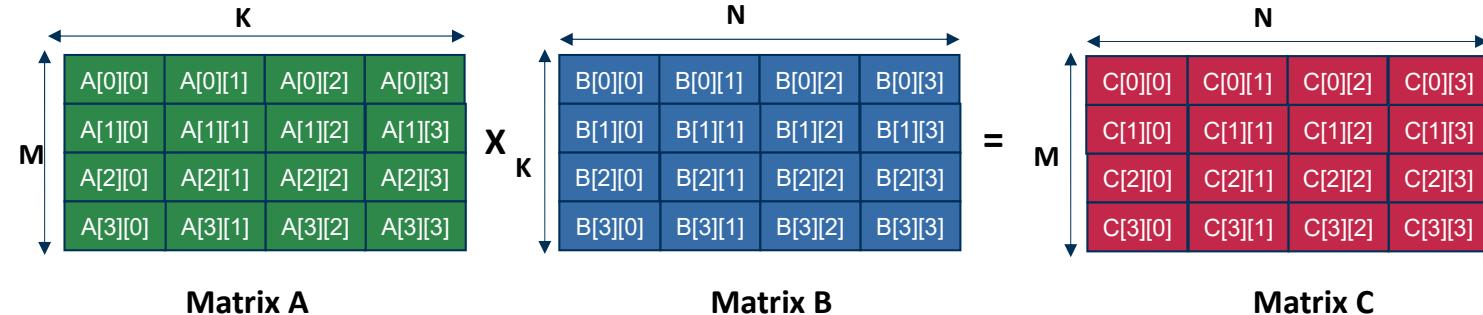
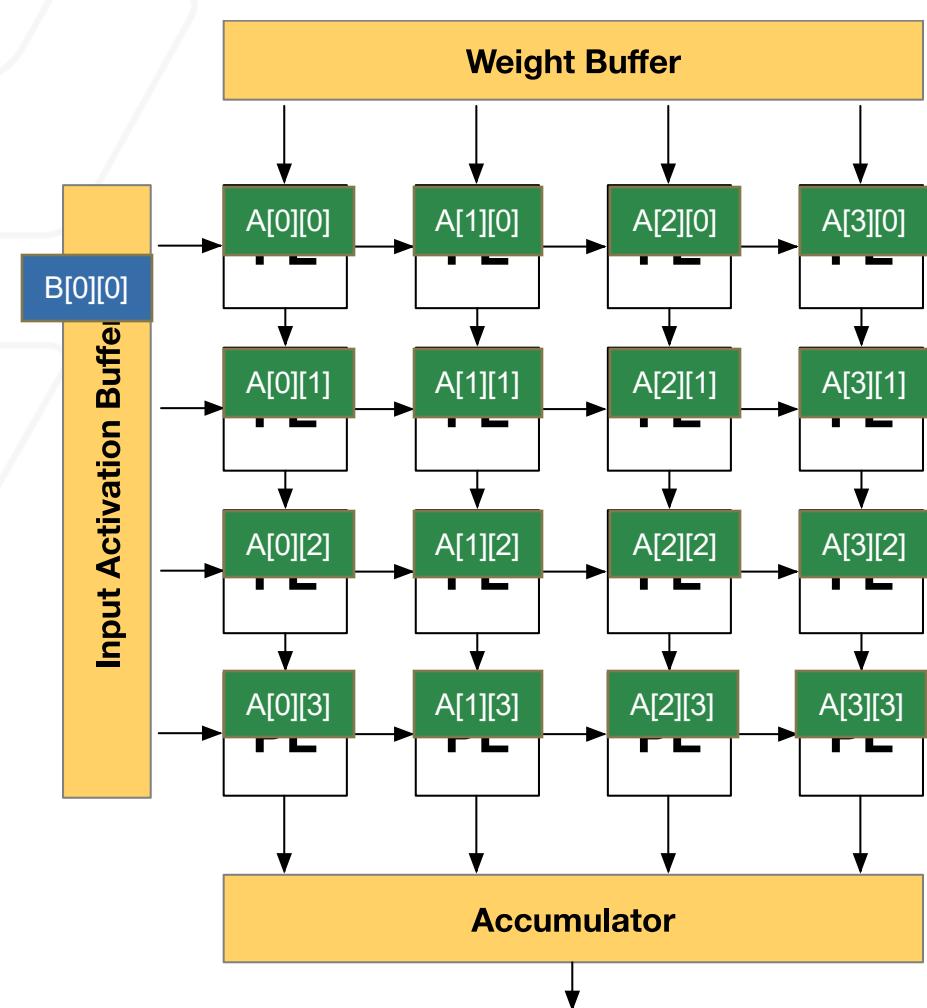
Assume following “mapping”

- Column: Reduced (aka contracted) Dimension (K)
- Row: Spatial Dimension (either M or N)

Note: alternate mapping styles possible

# Matrix Multiplication on a Systolic Array

**Phase 2:**  
Stream Activation Tensor



```
for m in range(4):
    for n in range(4):
        for k in range(4):
            C[m][n] += A[m][k] * B[k][n]
```

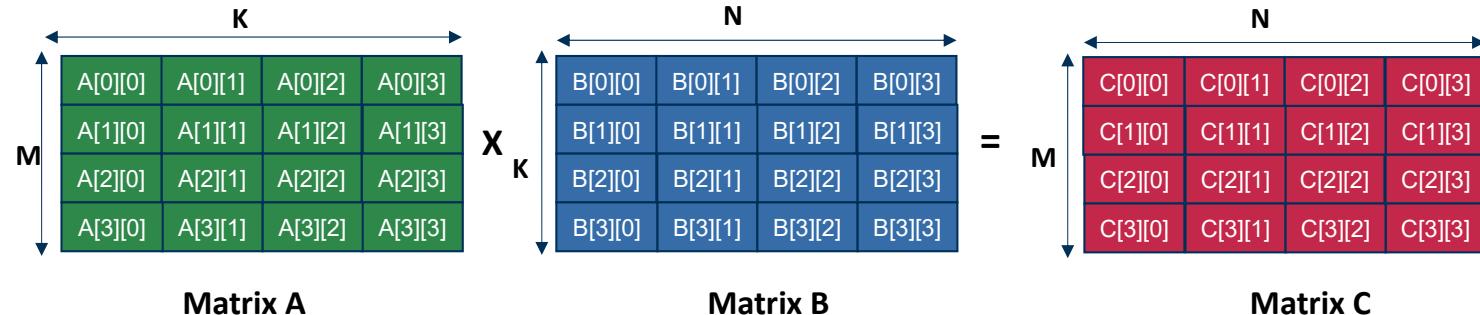
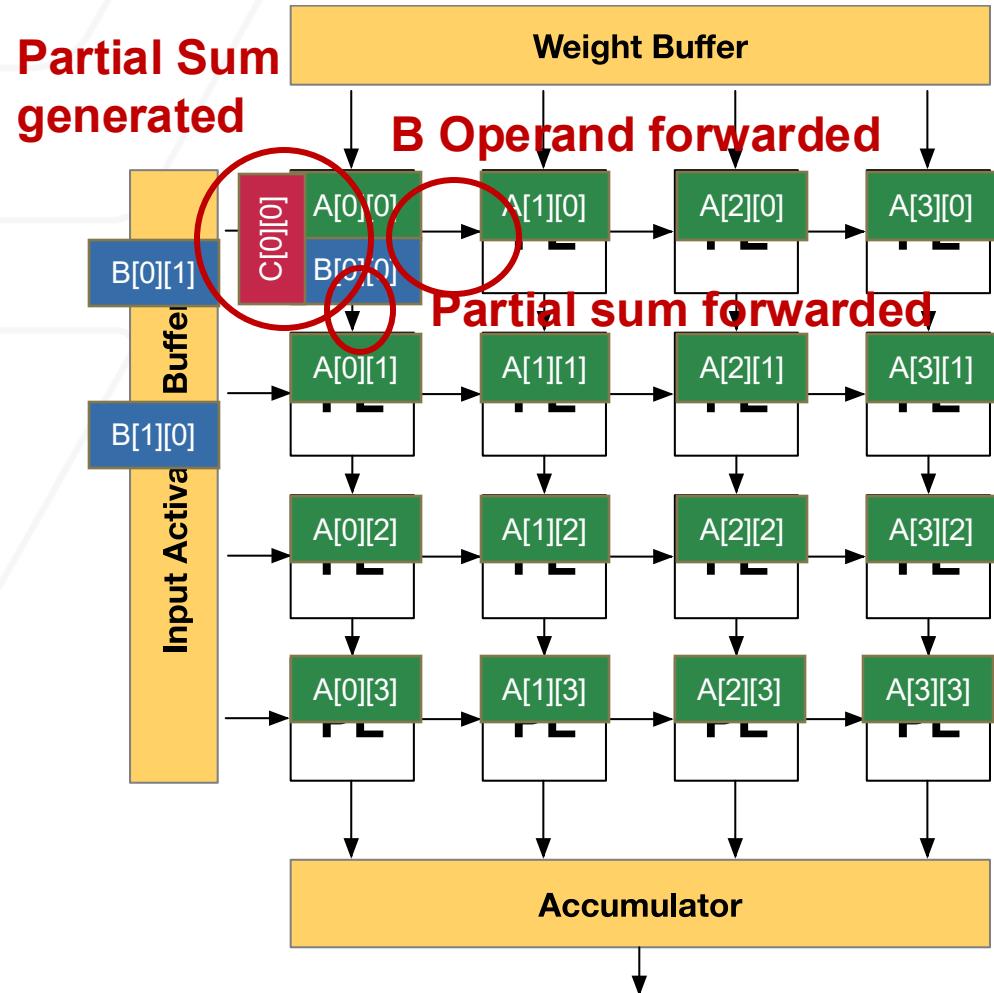
Assume following “mapping”

- Column: Reduced (aka contracted) Dimension ( $K$ )
- Row: Spatial Dimension (either  $M$  or  $N$ )

Note: alternate mapping styles possible

# Matrix Multiplication on a Systolic Array

**Phase 2:**  
Stream Activation Tensor



```
for m in range(4):
    for n in range(4):
        for k in range(4):
            C[m][n] += A[m][k] * B[k][n]
```

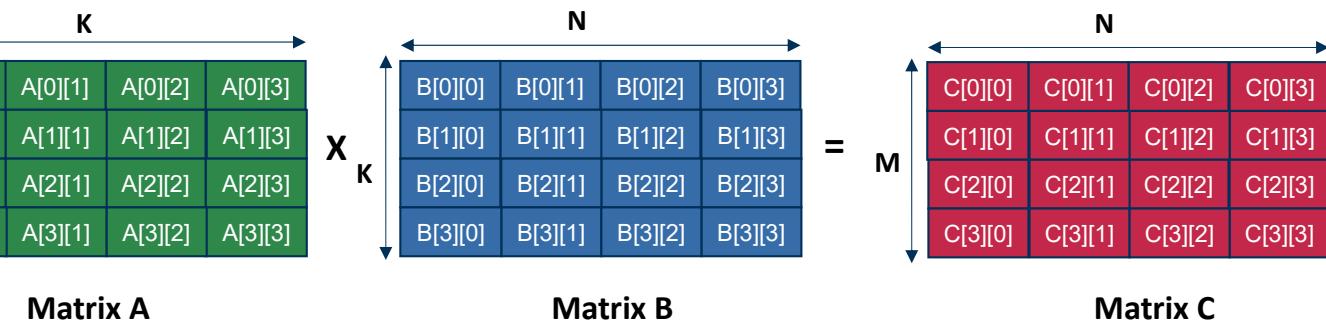
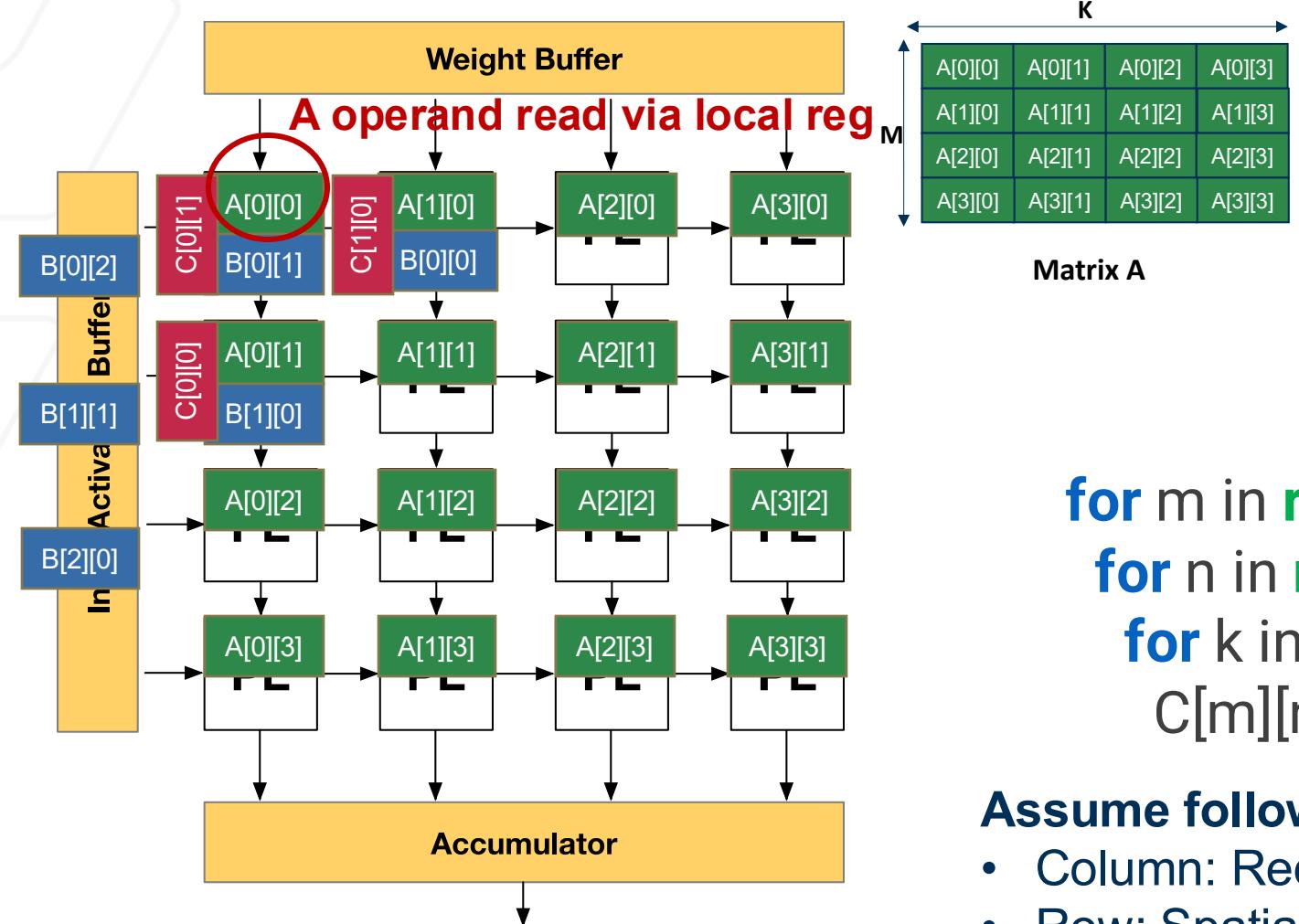
**Assume following “mapping”**

- Column: Reduced (aka contracted) Dimension (K)
- Row: Spatial Dimension (either M or N)

*Note: alternate mapping styles possible*

# Matrix Multiplication on a Systolic Array

**Phase 2:**  
Stream Activation Tensor



```

for m in range(4):
    for n in range(4):
        for k in range(4):
            C[m][n] += A[m][k] * B[k][n]
    
```

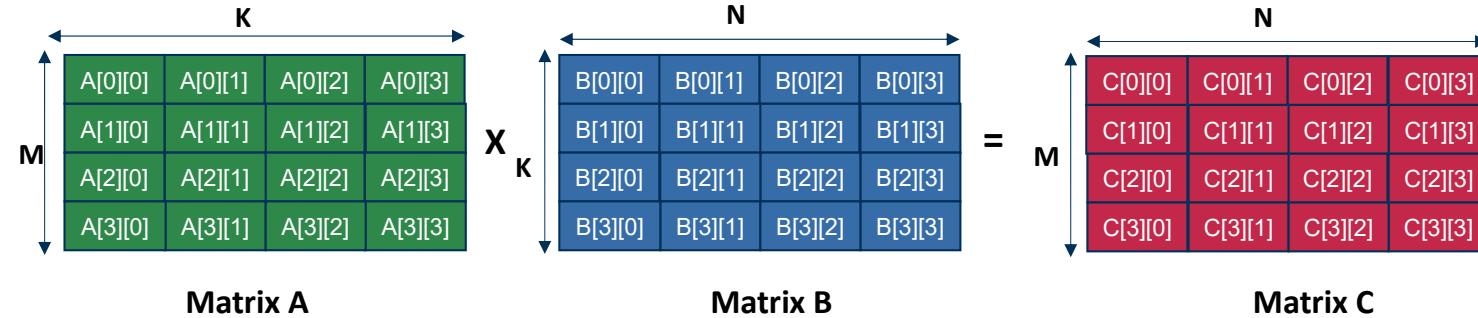
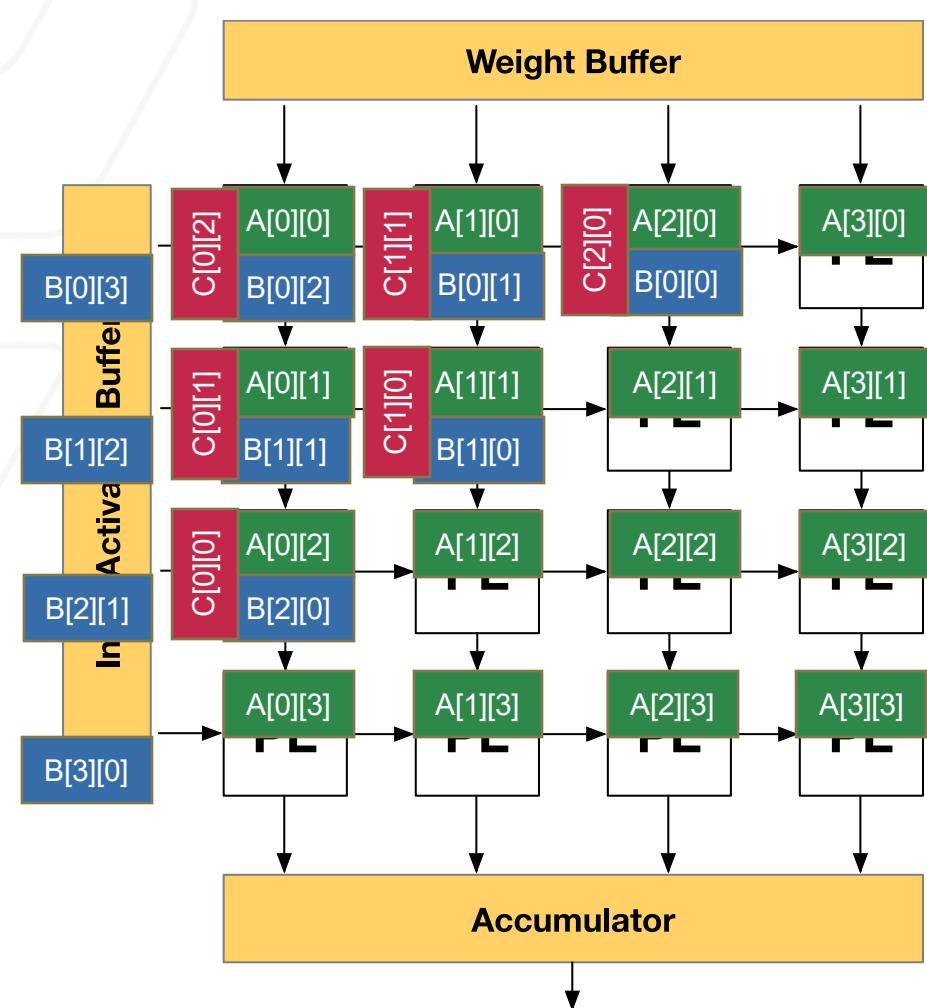
**Assume following “mapping”**

- Column: Reduced (aka contracted) Dimension (K)
- Row: Spatial Dimension (either M or N)

*Note: alternate mapping styles possible*

# Matrix Multiplication on a Systolic Array

**Phase 2:**  
Stream Activation Tensor



```

for m in range(4):
    for n in range(4):
        for k in range(4):
            C[m][n] += A[m][k] * B[k][n]
    
```

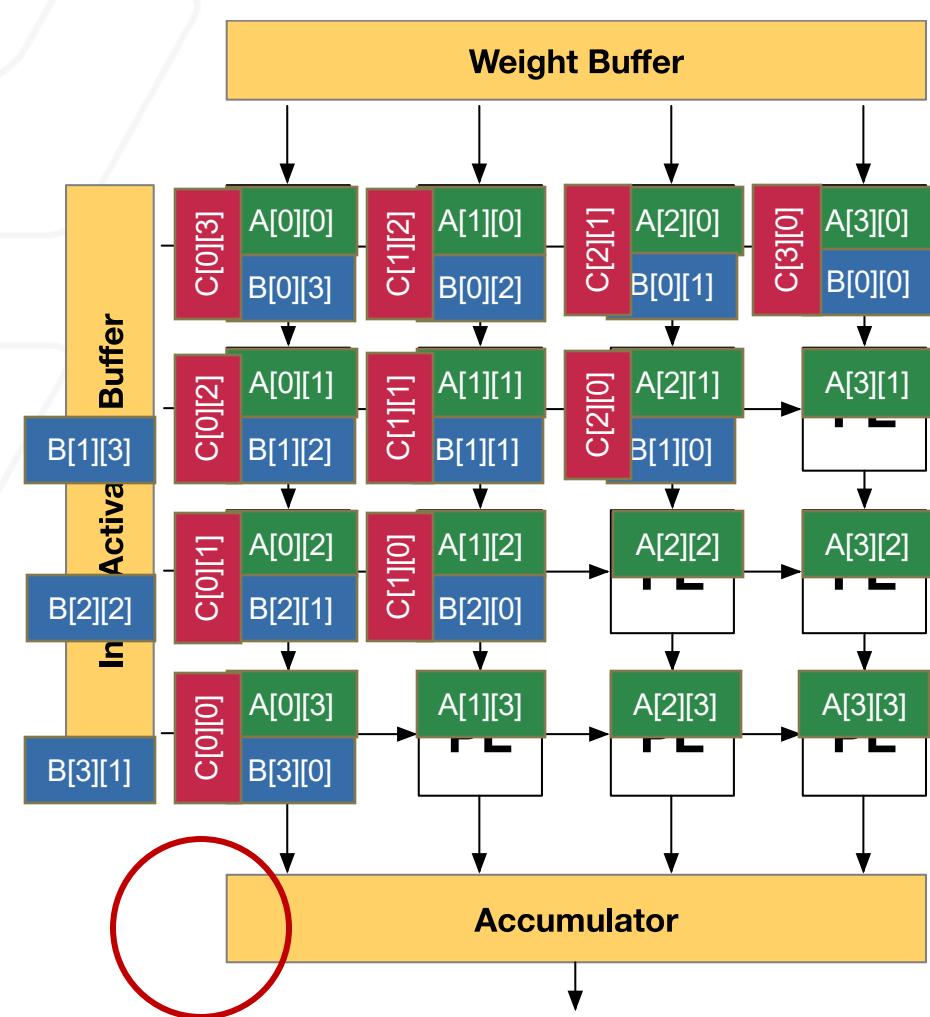
**Assume following “mapping”**

- Column: Reduced (aka contracted) Dimension (K)
- Row: Spatial Dimension (either M or N)

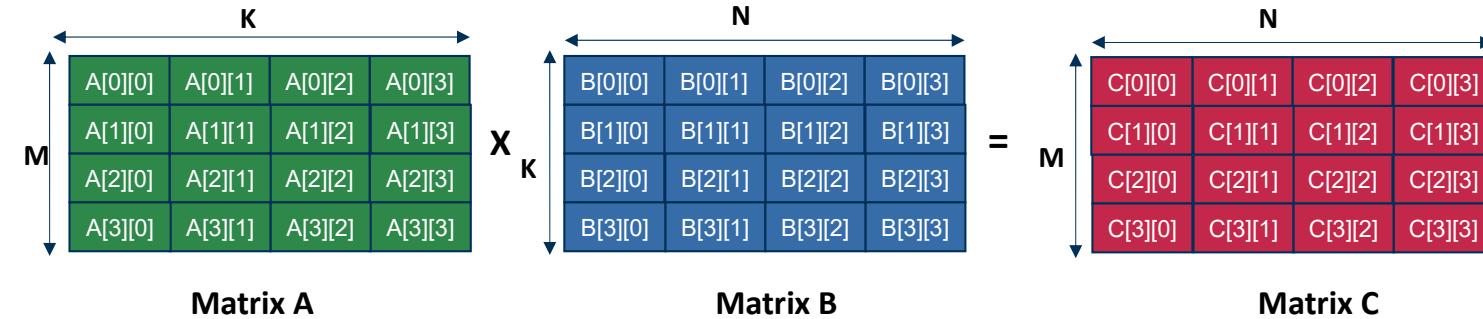
*Note: alternate mapping styles possible*

# Matrix Multiplication on a Systolic Array

**Phase 2:**  
Stream Activation Tensor



First output generated



```
for m in range(4):
    for n in range(4):
        for k in range(4):
            C[m][n] += A[m][k] * B[k][n]
```

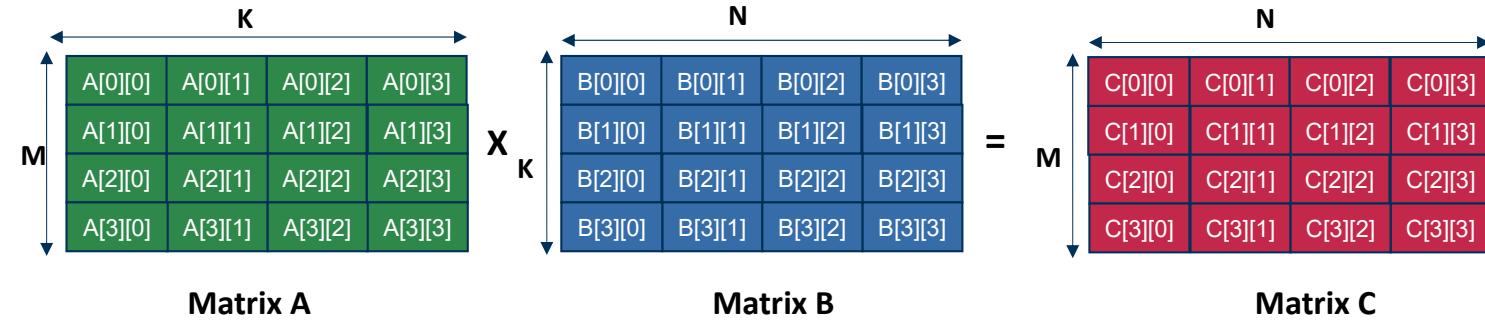
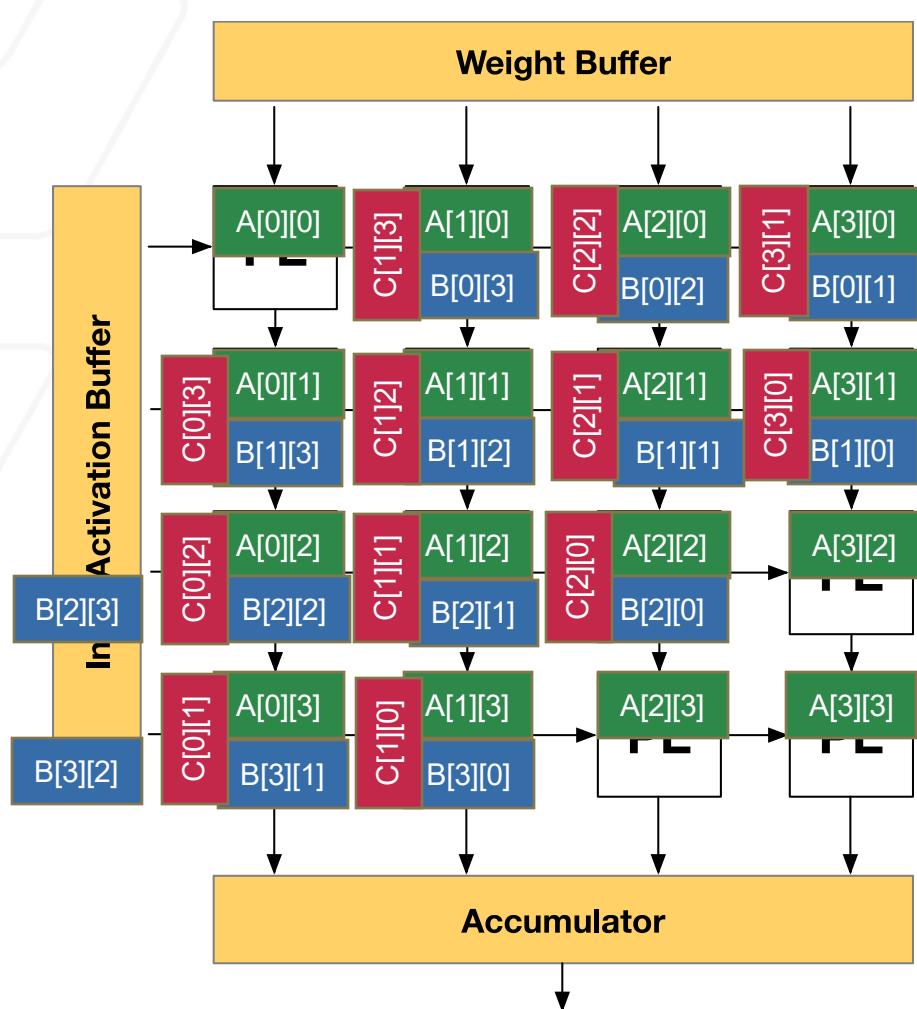
Assume following “mapping”

- Column: Reduced (aka contracted) Dimension (K)
- Row: Spatial Dimension (either M or N)

Note: alternate mapping styles possible

# Matrix Multiplication on a Systolic Array

**Phase 2:**  
Stream Activation Tensor



```

for m in range(4):
    for n in range(4):
        for k in range(4):
            C[m][n] += A[m][k] * B[k][n]
    
```

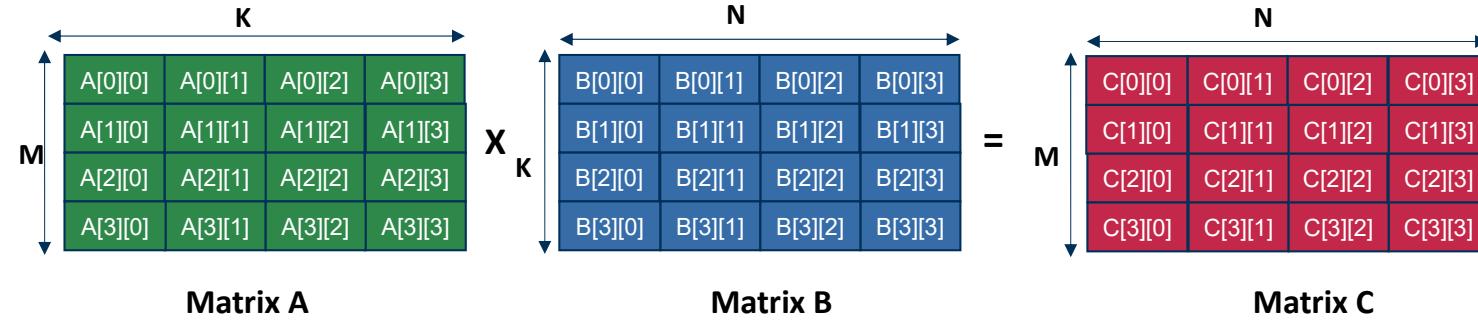
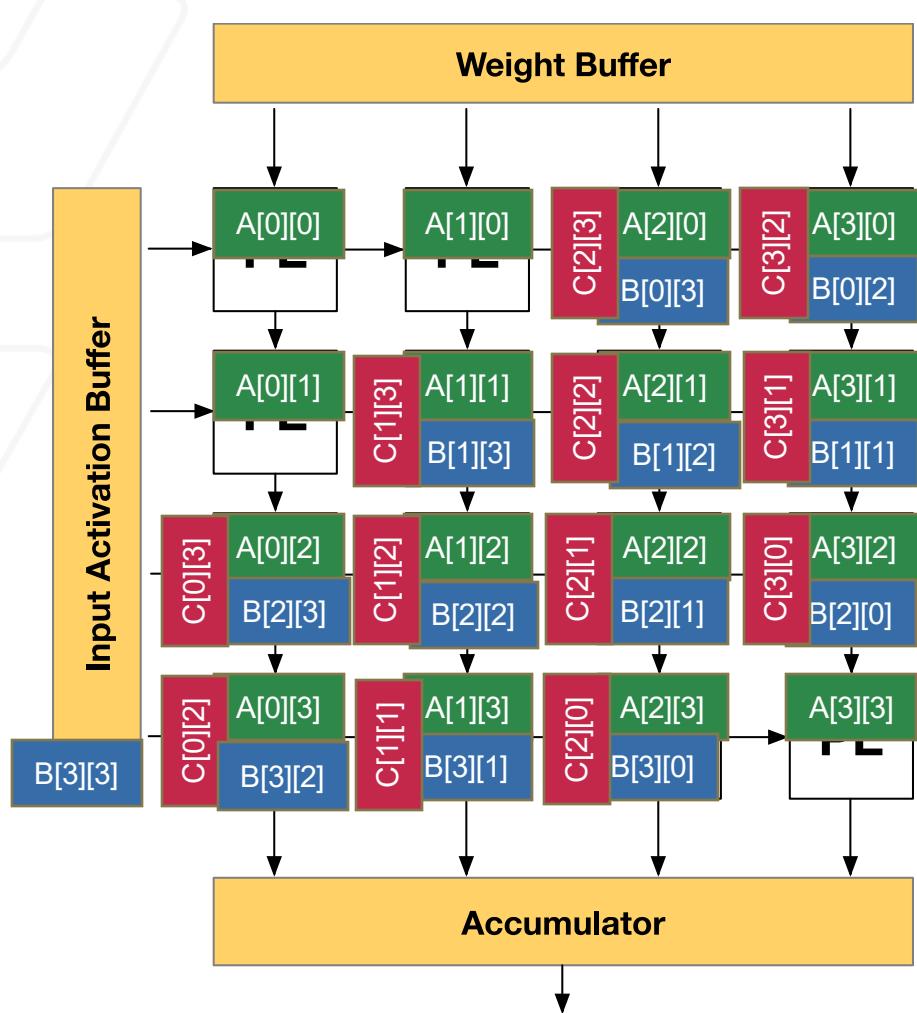
Assume following “mapping”

- Column: Reduced (aka contracted) Dimension ( $K$ )
- Row: Spatial Dimension (either  $M$  or  $N$ )

Note: alternate mapping styles possible

# Matrix Multiplication on a Systolic Array

**Phase 2:**  
Stream Activation Tensor



```

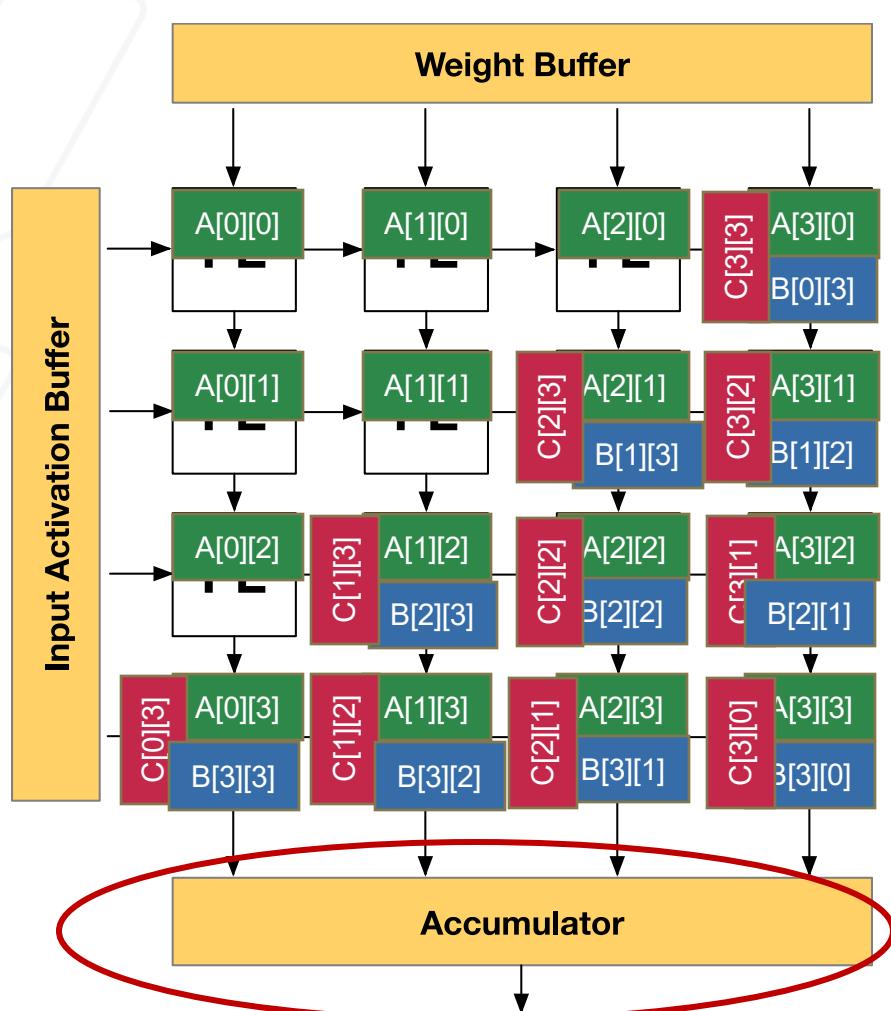
for m in range(4):
    for n in range(4):
        for k in range(4):
            C[m][n] += A[m][k] * B[k][n]
    
```

**Assume following “mapping”**

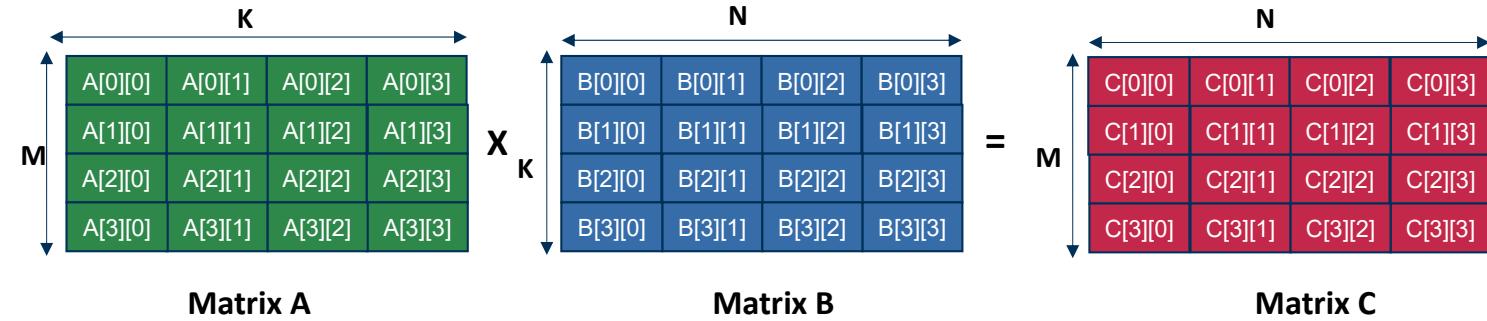
- Column: Reduced (aka contracted) Dimension (K)
- Row: Spatial Dimension (either M or N)

*Note: alternate mapping styles possible*

# Matrix Multiplication on a Systolic Array



**Four outputs every cycle in steady state**



**Phase 2:**  
Stream Activation Tensor

```

for m in range(4):
    for n in range(4):
        for k in range(4):
            C[m][n] += A[m][k] * B[k][n]
    
```

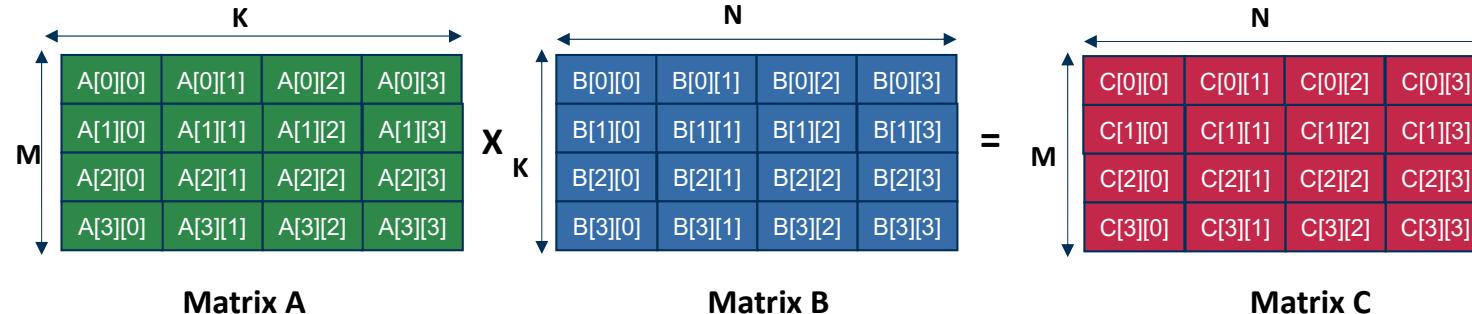
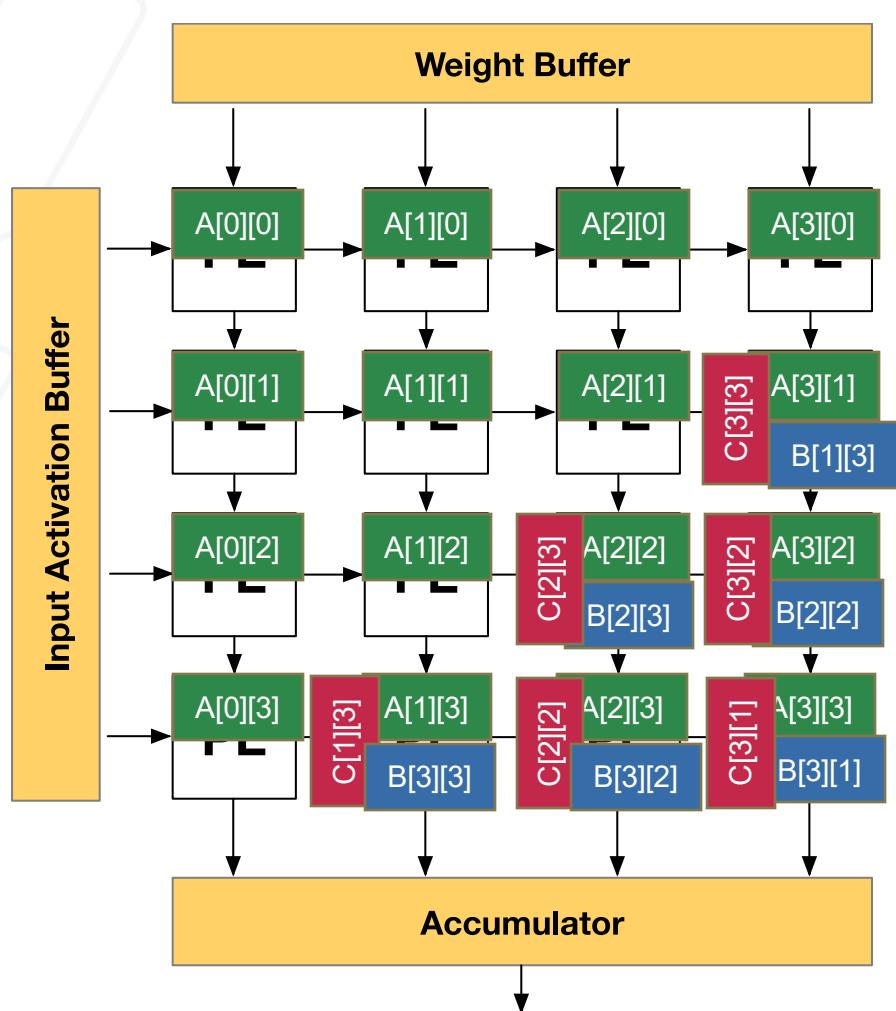
**Assume following “mapping”**

- Column: Reduced (aka contracted) Dimension (K)
- Row: Spatial Dimension (either M or N)

*Note: alternate mapping styles possible*

# Matrix Multiplication on a Systolic Array

**Phase 2:**  
Stream Activation Tensor



```
for m in range(4):
    for n in range(4):
        for k in range(4):
            C[m][n] += A[m][k] * B[k][n]
```

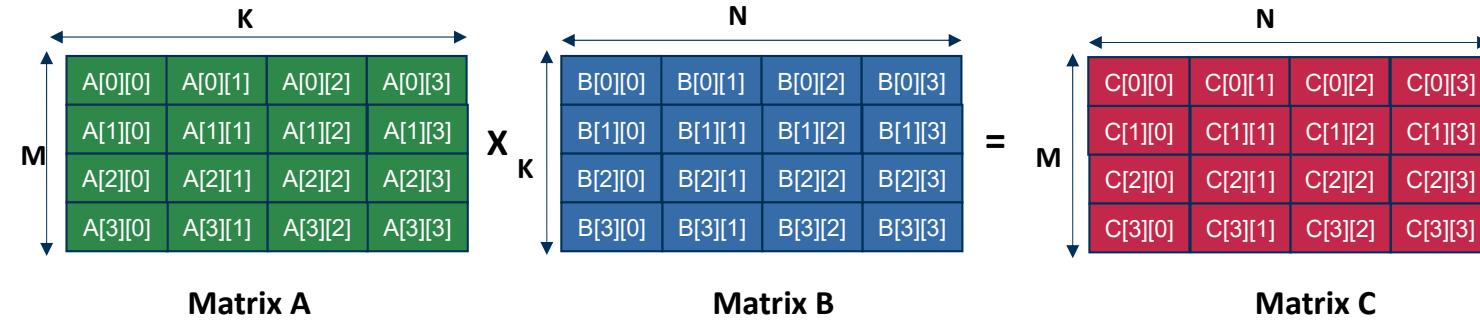
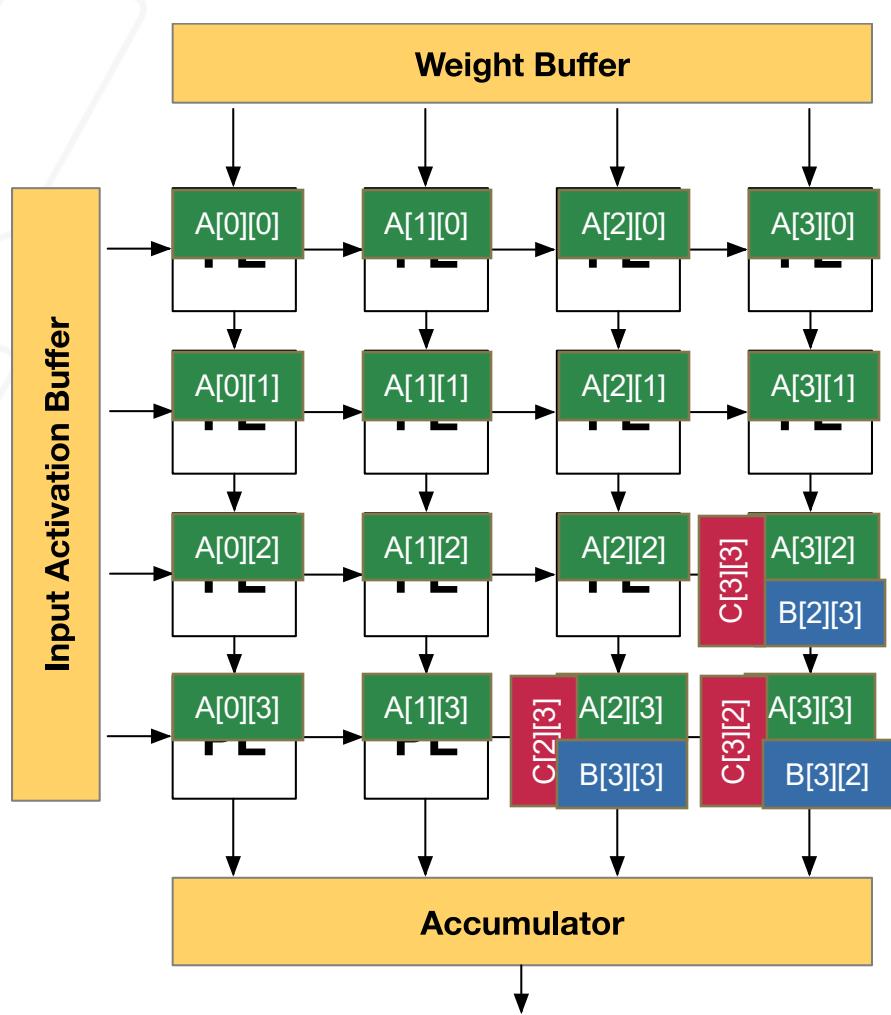
**Assume following “mapping”**

- Column: Reduced (aka contracted) Dimension (K)
- Row: Spatial Dimension (either M or N)

*Note: alternate mapping styles possible*

# Matrix Multiplication on a Systolic Array

**Phase 2:**  
Stream Activation Tensor



```
for m in range(4):
    for n in range(4):
        for k in range(4):
            C[m][n] += A[m][k] * B[k][n]
```

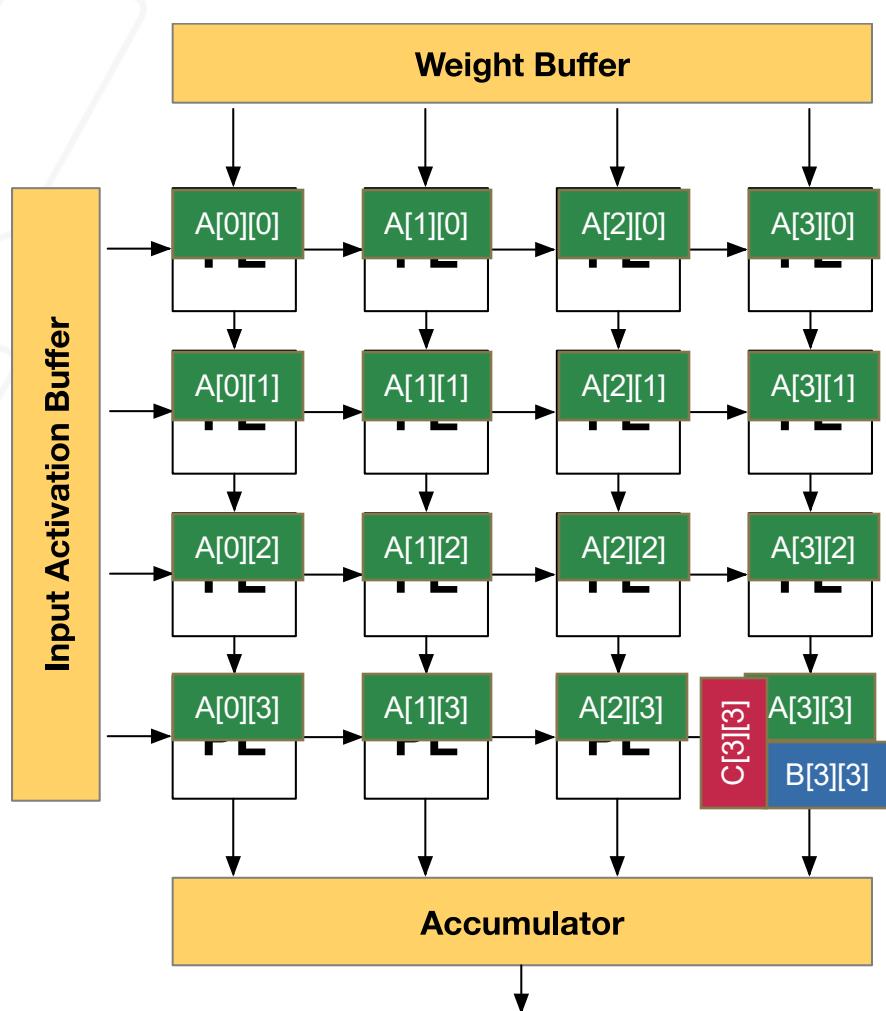
**Assume following “mapping”**

- Column: Reduced (aka contracted) Dimension (K)
- Row: Spatial Dimension (either M or N)

*Note: alternate mapping styles possible*

# Matrix Multiplication on a Systolic Array

Phase 2:  
Stream Activation Tensor



Matrix A			
M	K		
A[0][0]	A[0][1]	A[0][2]	A[0][3]
A[1][0]	A[1][1]	A[1][2]	A[1][3]
A[2][0]	A[2][1]	A[2][2]	A[2][3]
A[3][0]	A[3][1]	A[3][2]	A[3][3]

Matrix B			
N			
B[0][0]	B[0][1]	B[0][2]	B[0][3]
B[1][0]	B[1][1]	B[1][2]	B[1][3]
B[2][0]	B[2][1]	B[2][2]	B[2][3]
B[3][0]	B[3][1]	B[3][2]	B[3][3]

Matrix C			
N			
C[0][0]	C[0][1]	C[0][2]	C[0][3]
C[1][0]	C[1][1]	C[1][2]	C[1][3]
C[2][0]	C[2][1]	C[2][2]	C[2][3]
C[3][0]	C[3][1]	C[3][2]	C[3][3]

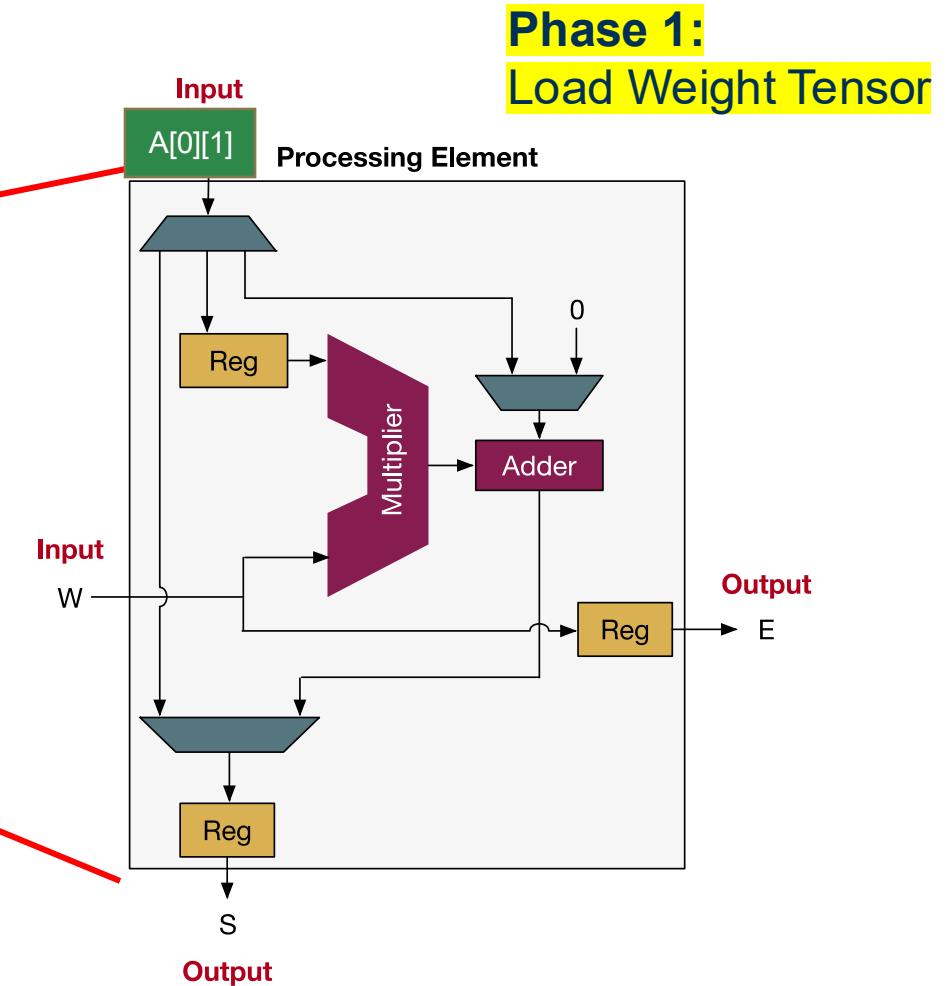
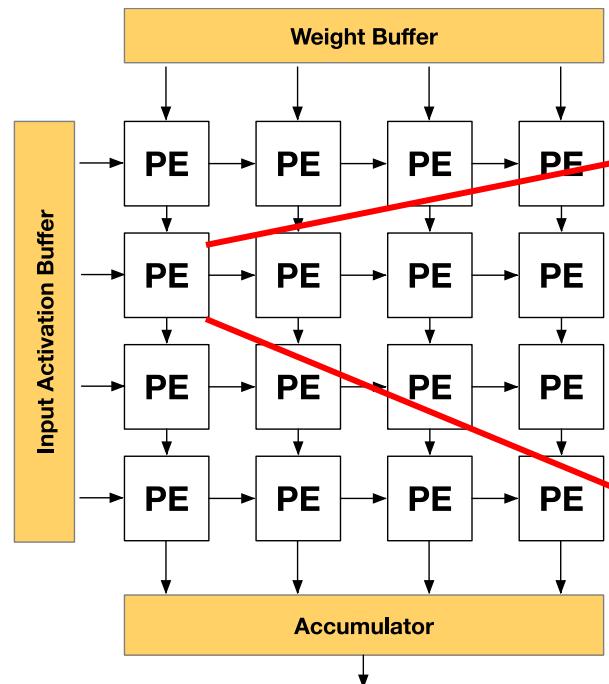
```
for m in range(4):  
    for n in range(4):  
        for k in range(4):  
            C[m][n] += A[m][k] * B[k][n]
```

Assume following “mapping”

- Column: Reduced (aka contracted) Dimension (K)
- Row: Spatial Dimension (either M or N)

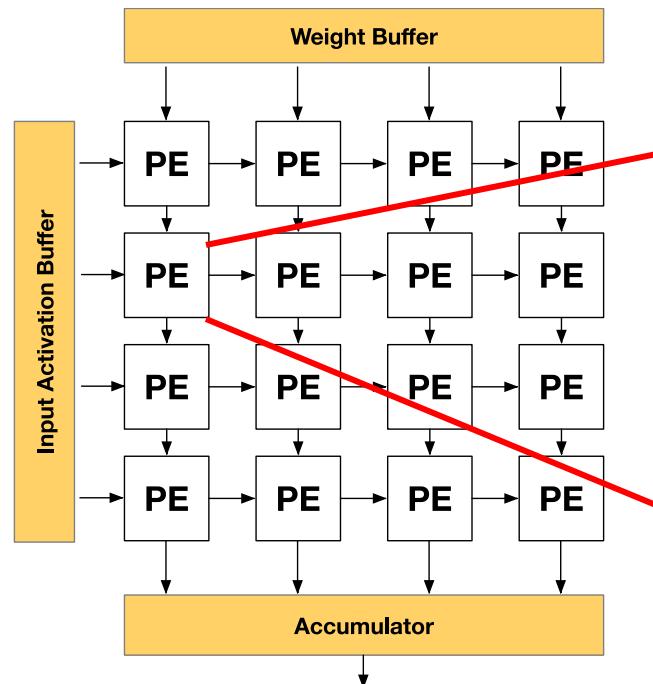
Note: alternate mapping styles possible

# Recap: Zooming into PE

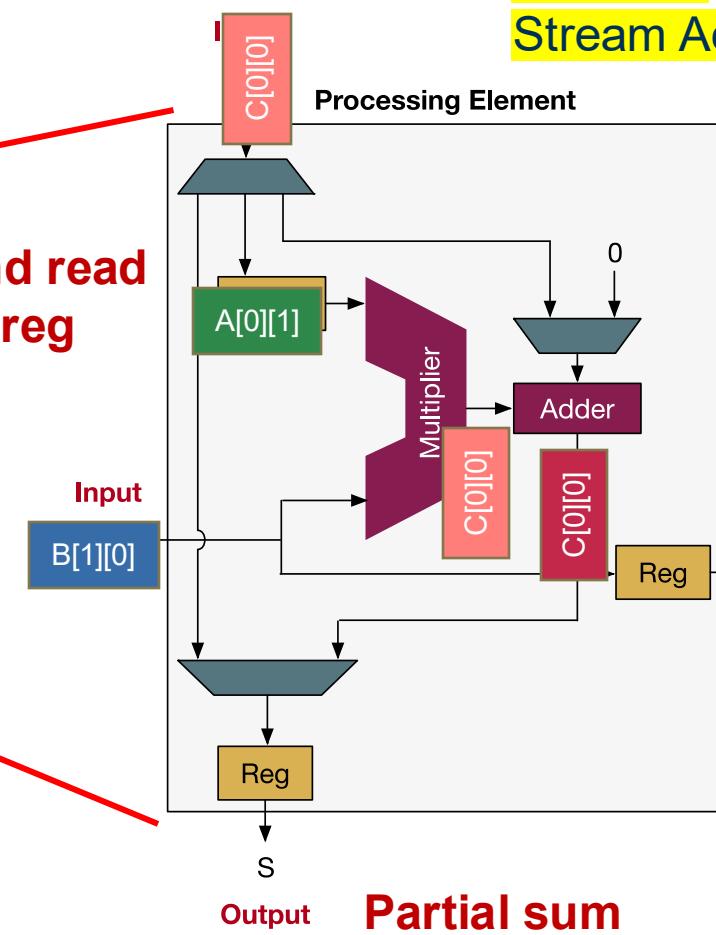


# Recap: Zooming into PE

Phase 2:  
Stream Activation Tensor



A operand  
read via local reg



B operand  
forwarded

Partial sum  
forwarded