

# CMSC 5743



## Efficient Computing of Deep Neural Networks

### Im01: Convolution Speedup

Bei Yu

CSE Department, CUHK

[byu@cse.cuhk.edu.hk](mailto:byu@cse.cuhk.edu.hk)

(Latest update: November 28, 2022)

2023 Spring



## These slides contain/adapt materials developed by

- Minsik Cho and Daniel Brand (2017). “MEC: memory-efficient convolution for deep neural network”. In: *Proc. ICML*
- Asit K. Mishra et al. (2017). “Fine-grained accelerators for sparse machine learning workloads”. In: *Proc. ASPDAC*, pp. 635–640
- Jongsoo Park et al. (2017). “Faster CNNs with direct sparse convolutions and guided pruning”. In: *Proc. ICLR*
- UC Berkeley EE290: “Hardware for Machine Learning”  
<https://inst.eecs.berkeley.edu/~ee290-2/sp20/>



① Convolution Basis

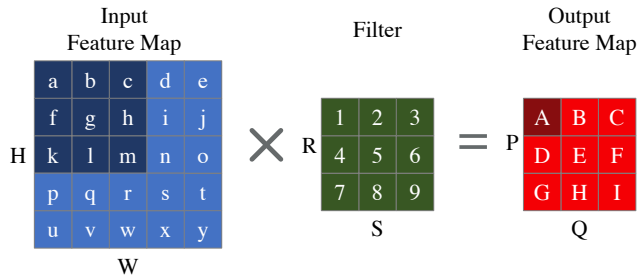
② GEMM

③ Direct Convolution

④ Sparse Convolution

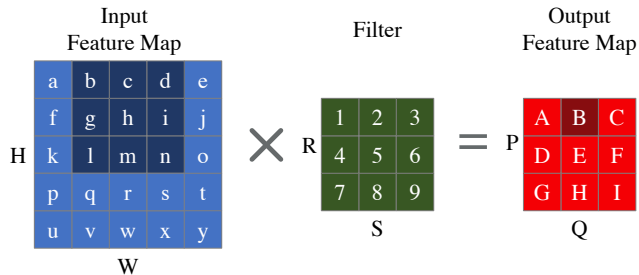


# Convolution Basis

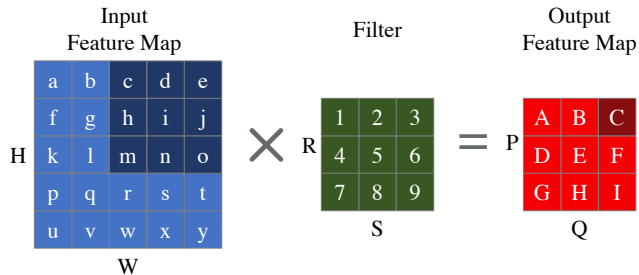


$$\begin{aligned}
 A &= a \cdot 1 + b \cdot 2 + c \cdot 3 \\
 &\quad + f \cdot 4 + g \cdot 5 + h \cdot 6 \\
 &\quad + k \cdot 7 + l \cdot 8 + m \cdot 9
 \end{aligned}$$

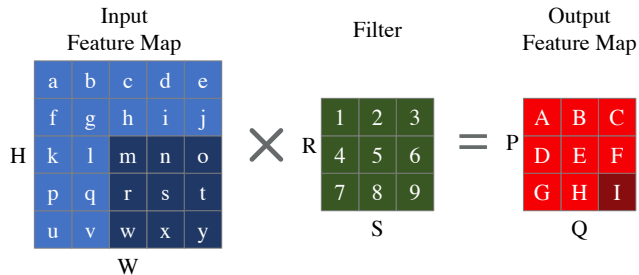
- H: Height of input feature map
- W: Width of input feature map
- R: Height of filter
- S: Width of filter
- P: Height of output feature map
- Q: Width of output feature map



- H: Height of input feature map
- W: Width of input feature map
- R: Height of filter
- S: Width of filter
- P: Height of output feature map
- Q: Width of output feature map
- **stride**: # of rows/columns traversed per step

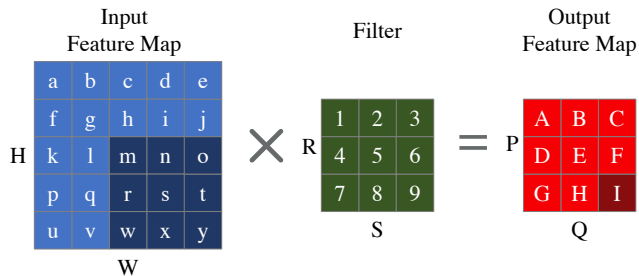


- H: Height of input feature map
- W: Width of input feature map
- R: Height of filter
- S: Width of filter
- P: Height of output feature map
- Q: Width of output feature map
- **stride**: # of rows/columns traversed per step



- H: Height of input feature map
- W: Width of input feature map
- R: Height of filter
- S: Width of filter
- P: Height of output feature map
- Q: Width of output feature map
- **stride**: # of rows/columns traversed per step



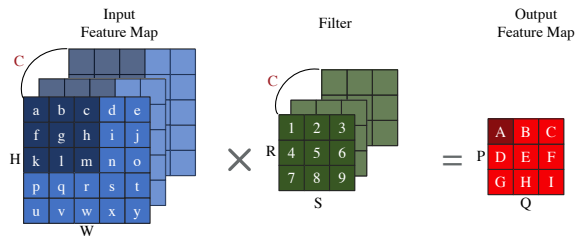


- H: Height of input feature map
- W: Width of input feature map
- R: Height of filter
- S: Width of filter
- P: Height of output feature map
- Q: Width of output feature map
- **stride**: # of rows/columns traversed per step

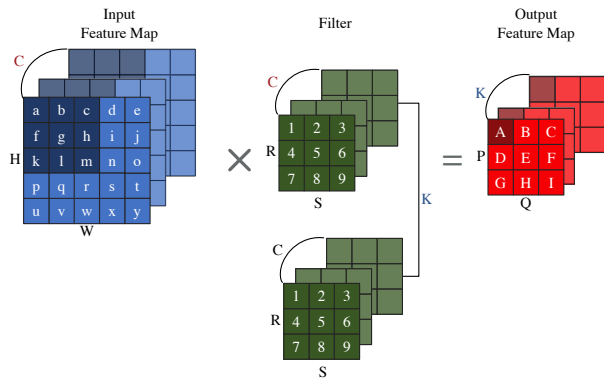
$$P = \frac{(H - R)}{\text{stride}} + 1;$$

$$Q = \frac{(W - S)}{\text{stride}} + 1.$$

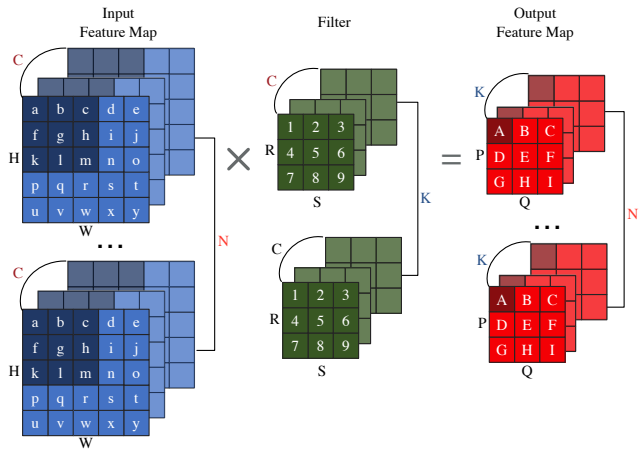




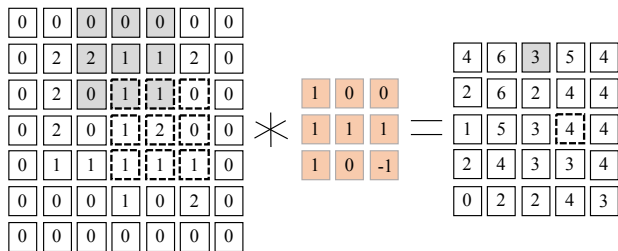
- H: Height of input feature map
- W: Width of input feature map
- R: Height of filter
- S: Width of filter
- P: Height of output feature map
- Q: Width of output feature map
- stride: # of rows/columns traversed per step
- padding: # of zero rows/columns added
- C: # of input channels



- $H$ : Height of input feature map
- $W$ : Width of input feature map
- $R$ : Height of filter
- $S$ : Width of filter
- $P$ : Height of output feature map
- $Q$ : Width of output feature map
- stride: # of rows/columns traversed per step
- padding: # of zero rows/columns added
- $C$ : # of input channels
- $K$ : # of output channels

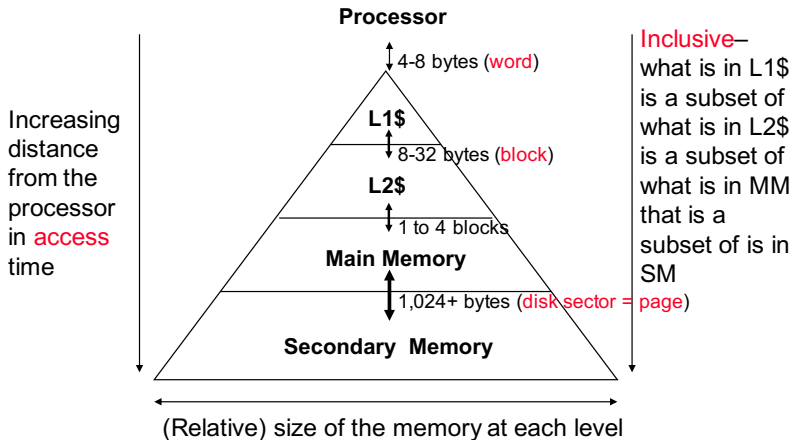


- $H$ : Height of input feature map
- $W$ : Width of input feature map
- $R$ : Height of filter
- $S$ : Width of filter
- $P$ : Height of output feature map
- $Q$ : Width of output feature map
- stride: # of rows/columns traversed per step
- padding: # of zero rows/columns added
- $C$ : # of input channels
- $K$ : # of output channels
- $N$ : Batch size



Direct convolution: No extra memory overhead

- Low performance
- Poor memory access pattern due to geometry-specific constraint
- Relatively short dot product

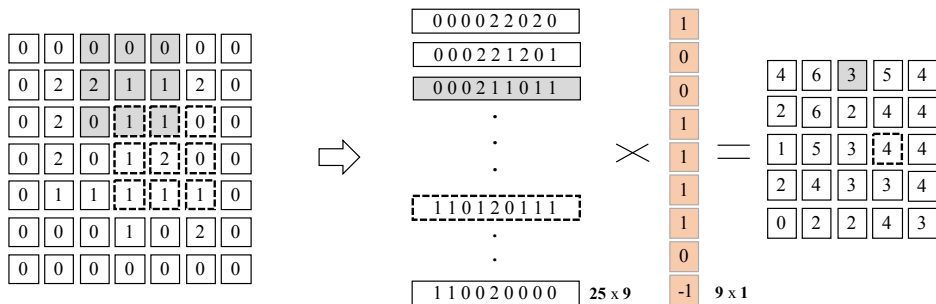


- **Spatial** locality
- **Temporal** Locality

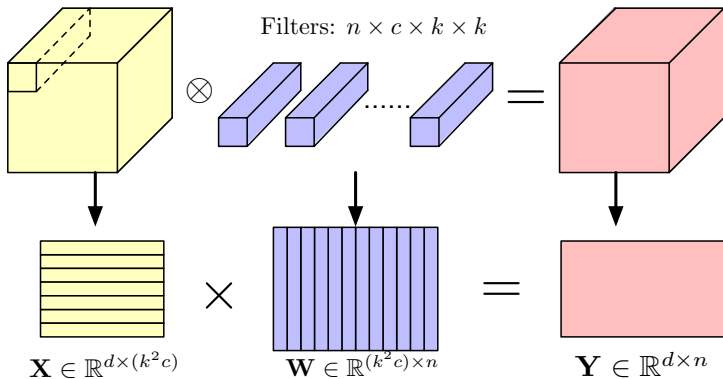


GEMM



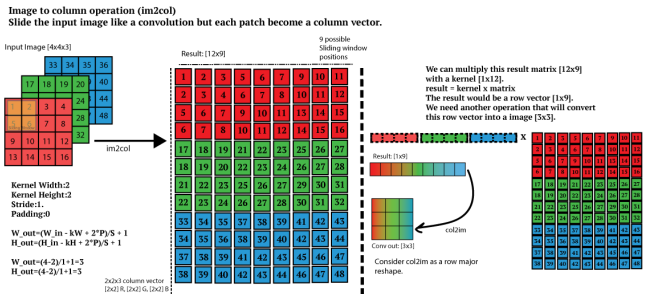


- Large extra memory overhead
- **Good** performance
- BLAS-friendly memory layout to enjoy SIMD/locality/parallelism
- Applicable for any convolution configuration on any platform

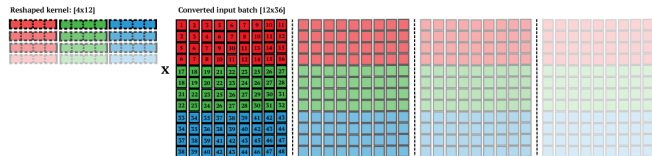


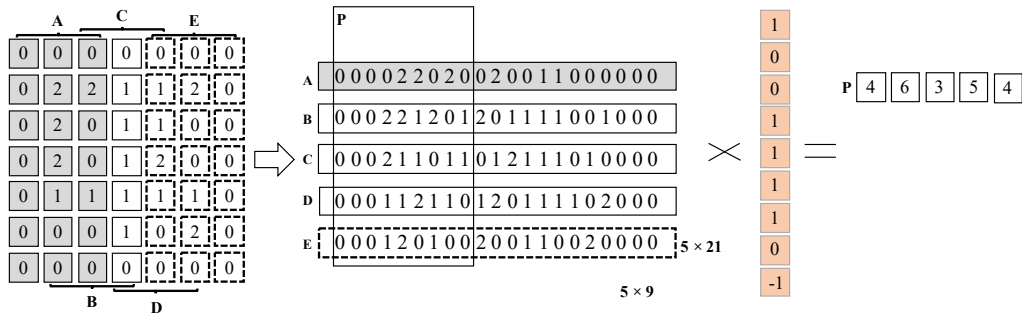
- Transform convolution to matrix multiplication
- Unified calculation for both convolution and fully-connected layers

# Im2col (Image2Column): Another View



We get true performance gain when the kernel has a large number of filters, i.e: F=4 and/or you have a batch of images (N=4). Example for the input batch [4x4x3x4], convolved with 4 filters [2x2x3x2]. The only problem with this approach is the amount of memory

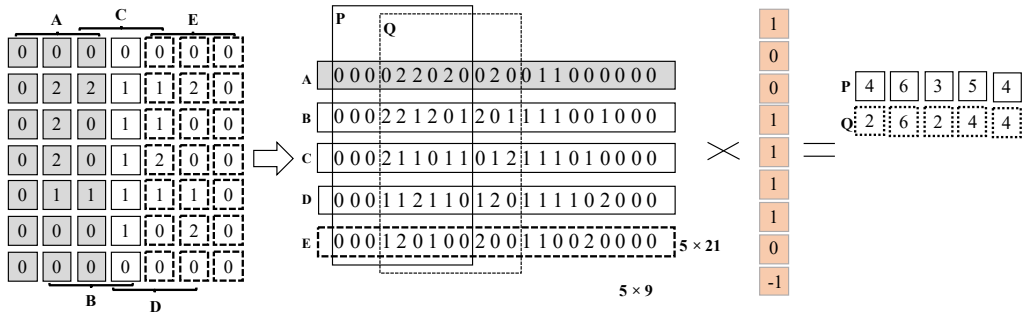




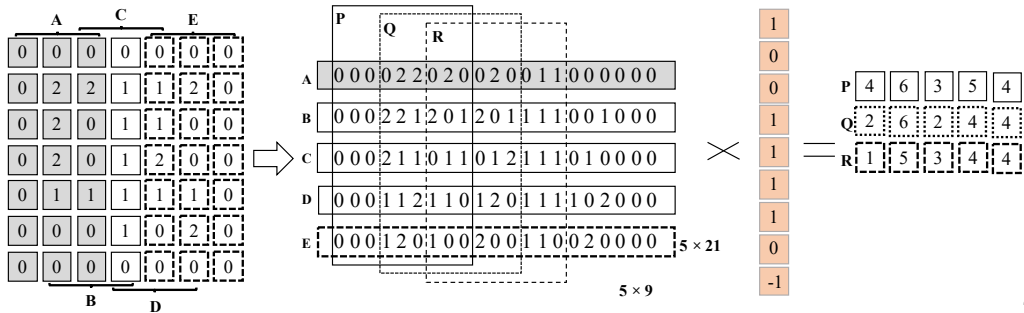
2

- Sub matrices in the lowered matrix will be “sgemm” ed in parallel
- Smaller memory foot print, cache locality, and explicit parallelism

<sup>2</sup>Minsik Cho and Daniel Brand (2017). “MEC: memory-efficient convolution for deep neural network”. In: *Proc. ICML*.



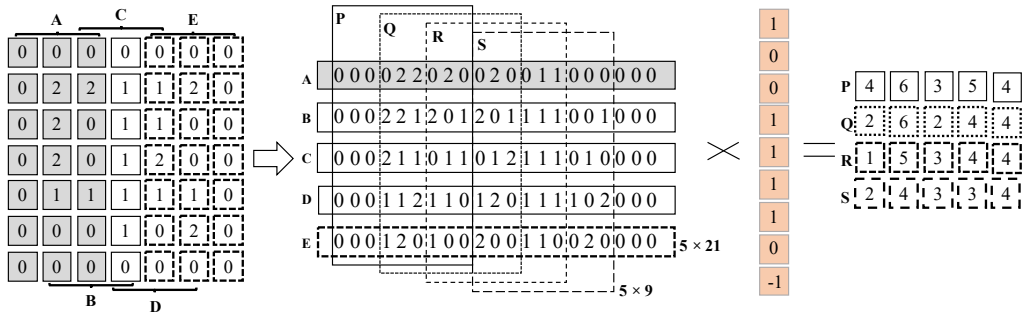
- Sub matrices in the lowered matrix will be “sgemm” ed in parallel
- Smaller memory foot print, cache locality, and explicit parallelism



2

- Sub matrices in the lowered matrix will be “sgemm” ed in parallel
- Smaller memory foot print, cache locality, and explicit parallelism

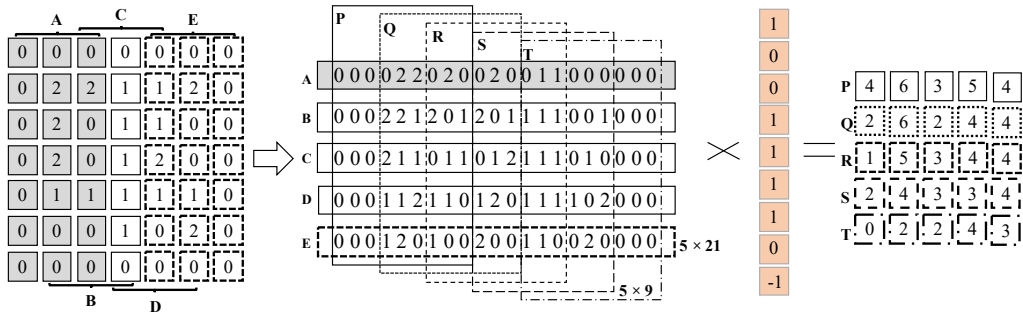
<sup>2</sup>Minsik Cho and Daniel Brand (2017). “MEC: memory-efficient convolution for deep neural network”. In: *Proc. ICML*.



2

- Sub matrices in the lowered matrix will be “sgemm” ed in parallel
- Smaller memory foot print, cache locality, and explicit parallelism

<sup>2</sup>Minsik Cho and Daniel Brand (2017). “MEC: memory-efficient convolution for deep neural network”. In: *Proc. ICML*.



2

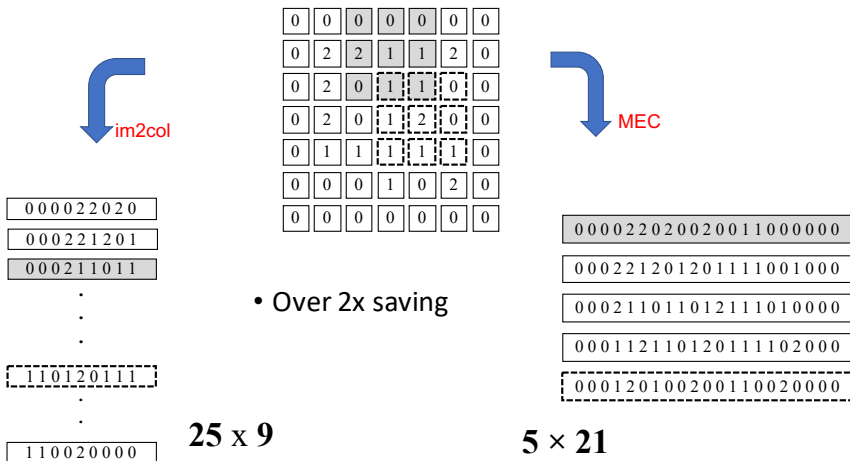
- Sub matrices in the lowered matrix will be “sgemm” ed in parallel
- Smaller memory foot print, cache locality, and explicit parallelism

<sup>2</sup>Minsik Cho and Daniel Brand (2017). “MEC: memory-efficient convolution for deep neural network”. In: *Proc. ICML*.





Over  $2\times$  memory saving<sup>3</sup>:

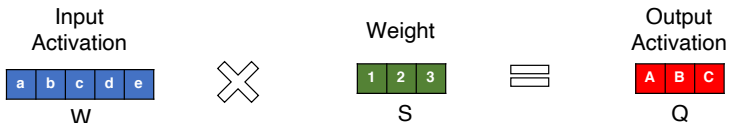


<sup>3</sup>Minsik Cho and Daniel Brand (2017). "MEC: memory-efficient convolution for deep neural network". In: *Proc. ICML*.



# Direct Convolution

# 1D Convolution Example



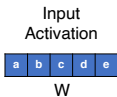
```
for(q=0; q<Q; q++){
  for (s=0; s<S; s++){
    OA[q] += IA[q+s] * W[s];
  }
}
```

**Output Stationary (OS)  
Dataflow**

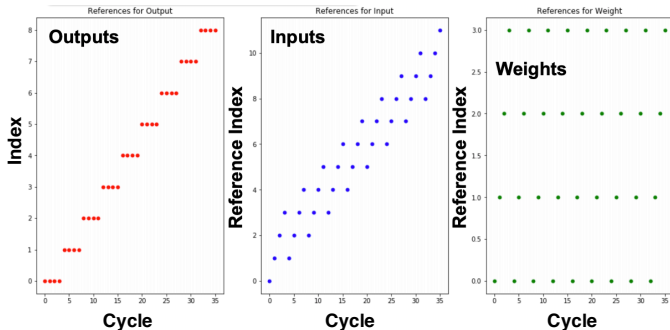
```
for (s=0; s<S; s++){
  for(q=0; q<Q; q++){
    OA[q] += IA[q+s] * W[s];
  }
}
```

**Weight Stationary (WS)  
Dataflow**

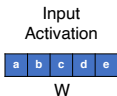
# Buffer Access Pattern 1: Output Stationary



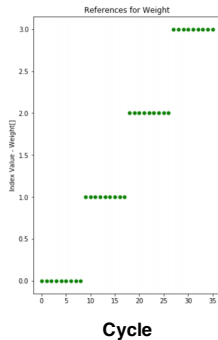
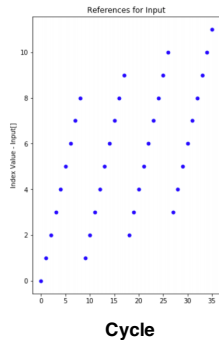
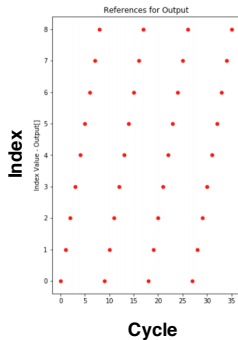
```
for (q=0; q<Q; q++){ // Q =9
  for (s=0; s<S; s++){ // S=4
    OA[q] += IA[q+s] * W[s];
  }
}
```

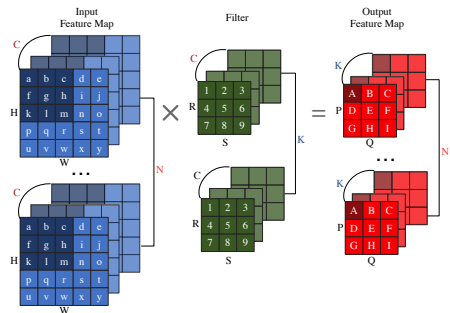


# Buffer Access Pattern 2: Weight Stationary



```
for (s=0; s<S; s++){ // S=4
  for (q=0; q<Q; q++){ // Q=9
    OA[q] += IA[q+s] * W[s];
  }
}
```

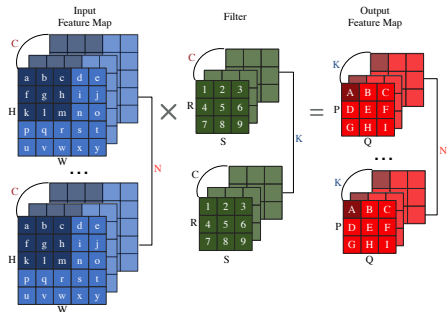




```

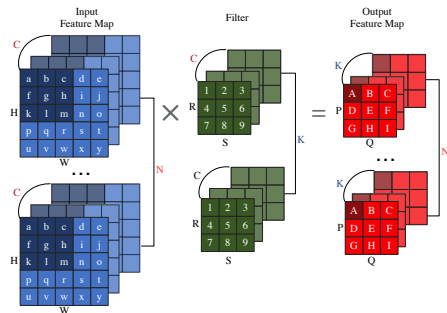
1  for (n=0; n<N; n++) {
2  for (k=0; k<K; k++) {
3  for (p=0; p<P; p++) {
4  for (q=0; q<Q; q++) {
5      OA[n][k][p][q] = 0;
6      for (r=0; r<R; r++) {
7      for (s=0; s<S; s++) {
8      for (c=0; c<C; c++) {
9          h = p * stride - pad + r;
10         w = q * stride - pad + s;
11         OA[n][k][p][q] += IA[n][c][h][w] * W[k][c][r][s];
12     } } } } } } } }
    
```

# Direct Convolution: Loop Ordering



```
1  for (n=0; n<N; n++) {
2  for (r=0; r<R; r++) {
3  for (s=0; s<S; s++) {
4  for (c=0; c<C; c++) {
5  for (k=0; k<K; k++) {
6      float curr_w = W[r][s][c][k];
7      for (p=0; p<P; p++) {
8      for (q=0; q<Q; q++) {
9          h = p * stride - pad + r;
10         w = q * stride - pad + s;
11         OA[n][k][p][q] += IA[n][c][h][w] * curr_w;
12     } } } } } }
```

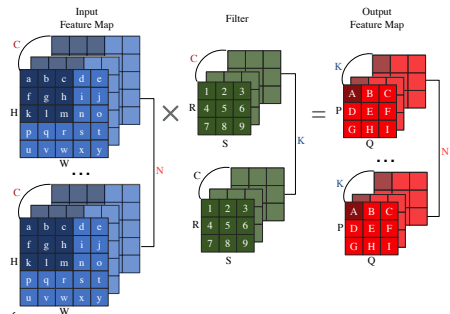
# Direct Convolution: Loop Ordering + Unrolling



```
1  for (n=0; n<N; n++) {
2  for (r=0; r<R; r++) {
3  for (s=0; s<S; s++) {
4  spatial_for (c=0; c<C; c++) {
5  spatial_for (k=0; k<K; k++) {
6  float curr_w = W[r][s][c][k];
7  for (p=0; p<P; p++) {
8  for (q=0; q<Q; q++) {
9      h = p * stride - pad + r;
10     w = q * stride - pad + s;
11     OA[n][k][p][q] += IA[n][c][h][w] * curr_w;
12 } } } } } }
```



# Direct Convolution: Loop Ordering + Unrolling + Tiling



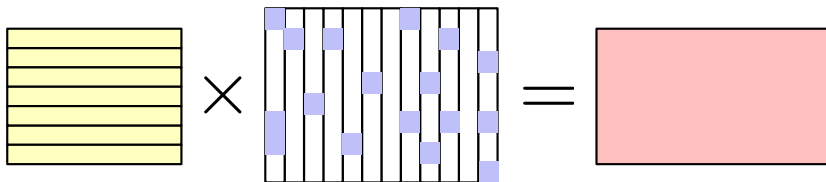
```
1  for (n=0; n<N; n++) {
2  for (r=0; r<R; r++) {
3  for (s=0; s<S; s++) {
4  for (c_t=0; c_t<C/16; c_t++) {
5  for (k_t=0; k_t<K/64; k_t++) {
6  spatial_for (c_s=0; c_s<16; c_s++) {
7  spatial_for (k_s=0; k_s<64; k_s++) {
8      int curr_c = c_t * 16 + c_s;
9      int curr_k = k_t * 64 + k_s;
10     float curr_w = W[r][s][curr_c][curr_k];
11     for (p=0; p<P; p++) for (q=0; q<Q; q++) {
12         h = p * stride - pad + r; w = q * stride - pad + s;
13         OA[n][curr_k][p][q] += IA[n][curr_c][h][w] * curr_w;
14     } } } } } }
```



# Sparse Convolution



- Our DNN may be **redundant**, and sometimes the filters may be **sparse**
- Sparsity can be helpful to overcome **over-fitting**





X			
0	0	3	0
7	0	0	0
0	0	4	8
6	5	3	0
2	0	0	1
0	0	0	8

 \* 

W
0
0
4
8

---

**Algorithm 1** Sparse Convolution Naive 1

---

- 1: **for all**  $w[i]$  **do**
  - 2:     **if**  $w[i] = 0$  **then**
  - 3:         Continue;
  - 4:     **end if**
  - 5:     output feature map  $Y \leftarrow X \times w[i]$ ;
  - 6: **end for**
-



$$\begin{array}{c}
 X \\
 \begin{array}{|c|c|c|c|}
 \hline
 0 & 0 & 3 & 0 \\
 \hline
 7 & 0 & 0 & 0 \\
 \hline
 0 & 0 & 4 & 8 \\
 \hline
 6 & 5 & 3 & 0 \\
 \hline
 2 & 0 & 0 & 1 \\
 \hline
 0 & 0 & 0 & 8 \\
 \hline
 \end{array}
 \end{array}
 *
 \begin{array}{c}
 W \\
 \begin{array}{|c|}
 \hline
 0 \\
 \hline
 0 \\
 \hline
 4 \\
 \hline
 8 \\
 \hline
 \end{array}
 \end{array}$$

---

## Algorithm 2 Sparse Convolution Naive 1

---

- 1: **for all**  $w[i]$  **do**
  - 2:     **if**  $w[i] = 0$  **then**
  - 3:         Continue;
  - 4:     **end if**
  - 5:     output feature map  $Y \leftarrow X \times w[i]$ ;
  - 6: **end for**
- 

**BAD** implementation for Pipeline!

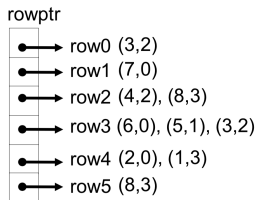
Instr. No.	Pipeline Stage						
	IF	ID	EX	MEM	WB		
1	IF	ID	EX	MEM	WB		
2		IF	ID	EX	MEM	WB	
3			IF	ID	EX	MEM	WB
4				IF	ID	EX	MEM
5					IF	ID	EX
<b>Clock Cycle</b>	1	2	3	4	5	6	7



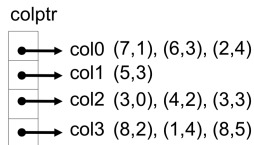
**A**

0	0	3	0
7	0	0	0
0	0	4	8
6	5	3	0
2	0	0	1
0	0	0	8

**A matrix example**

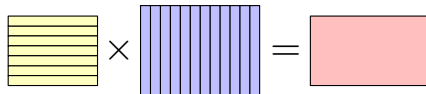


**Compressed Sparse Row (CSR)**



**Compressed Sparse Column (CSC)**

- CSR: Good for operation on **feature maps**
- CSC: Good for operation on **filters**
- We have better control on filters, thus usually CSC.





matrix \* sparse vector

$$\begin{array}{cccc} & \text{X} & & \\ \begin{array}{|c|c|c|c|} \hline 0 & 0 & 3 & 0 \\ \hline 7 & 0 & 0 & 0 \\ \hline 0 & 0 & 4 & 8 \\ \hline 6 & 5 & 3 & 0 \\ \hline 2 & 0 & 0 & 1 \\ \hline 0 & 0 & 0 & 8 \\ \hline \end{array} & * & \begin{array}{|c|} \hline \text{w} \\ \hline 0 \\ \hline 0 \\ \hline 4 \\ \hline 8 \\ \hline \end{array} & = & \begin{array}{|c|} \hline \text{Y} \\ \hline 12 \\ \hline 0 \\ \hline 16 \\ \hline 12 \\ \hline 0 \\ \hline 0 \\ \hline \end{array}\end{array}$$

$$\begin{array}{cccc} \begin{array}{|c|c|c|c|} \hline 0 & 0 & 3 & 0 \\ \hline 7 & 0 & 0 & 0 \\ \hline 0 & 0 & 4 & 8 \\ \hline 6 & 5 & 3 & 0 \\ \hline 2 & 0 & 0 & 1 \\ \hline 0 & 0 & 0 & 8 \\ \hline \end{array} & * & \begin{array}{|c|} \hline 0 \\ \hline 0 \\ \hline 4 \\ \hline 8 \\ \hline \end{array} & = & \begin{array}{|c|} \hline 12 \\ \hline 0 \\ \hline 80 \\ \hline 12 \\ \hline 8 \\ \hline 64 \\ \hline \end{array}\end{array}$$

- **BAD** implementation for Spatial Locality!
- **Poor** memory access patterns

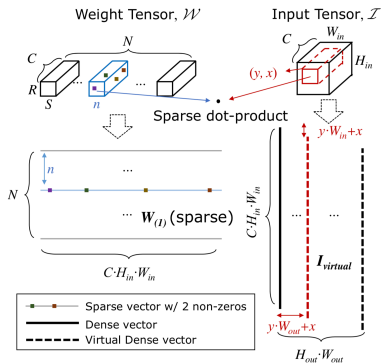


Figure 1: Conceptual view of the direct sparse convolution algorithm. Computation of output value at  $(y,x)$ th position of  $n$ th output channel is highlighted.

```

for each output channel n {
  for j in [W.rowptr[n], W.rowptr[n+1]] {
    off = W.colidx[j]; coeff = W.value[j]
    for (int y = 0; y < H_OUT; ++y) {
      for (int x = 0; x < W_OUT; ++x) {
        out[n][y][x] += coeff*in[off+f(0,y,x)]
      }
    }
  }
}
    
```

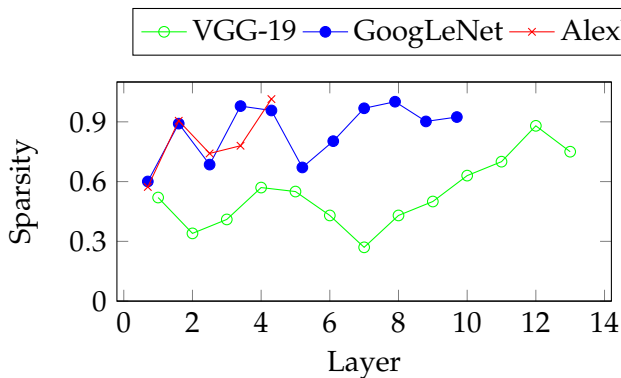
Figure 2: Sparse convolution pseudo code. Matrix  $W$  has *compressed sparse row* (CSR) format, where  $\text{rowptr}[n]$  points to the first non-zero weight of  $n$ th output channel. For the  $j$ th non-zero weight at  $(n,c,r,s)$ ,  $W.\text{colidx}[j]$  contains the offset to  $(c,r,s)$ th element of tensor  $\text{in}$ , which is pre-computed by layout function  $f(c,r,s)$ . If  $\text{in}$  has CHW format,  $f(c,r,s) = (cH_{in} + r)W_{in} + s$ . The “virtual” dense matrix is formed on-the-fly by shifting  $\text{in}$  by  $(0,y,x)$ .

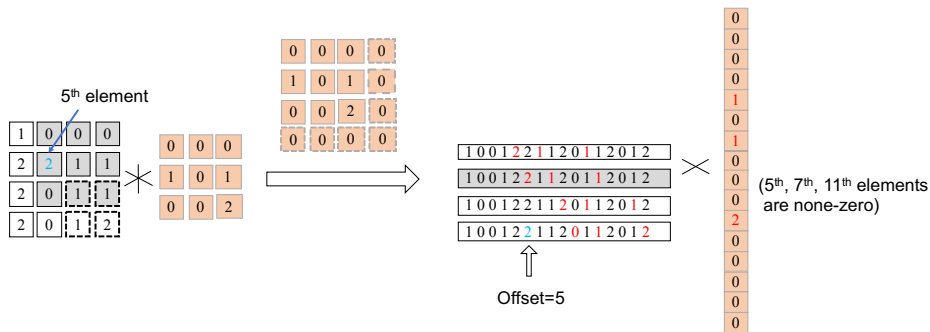
4





- Sparsity is a desired property for computation acceleration. (cuSPARSE library, direct sparse convolution, etc.)
- Sometimes not only the **filters** but also the **input feature maps** are sparse.





- Efficient programming implementation required; (Improve pipeline efficiency)
- When  $\text{sparsity}(\text{input}) = 0.9$ ,  $\text{sparsity}(\text{weight}) = 0.8$ , more than  $10\times$  speedup;
- Some other issues:
  - How to be compatible with pooling layer?
  - Transform between dense & sparse formats