



CMSC 5743

Efficient Computing of Deep Neural Networks

Implementation 01: GEMM-1

Bei Yu

CSE Department, CUHK

byu@cse.cuhk.edu.hk

(Latest update: February 7, 2023)

2023 Spring



① Convolution Basis

② Im2Col

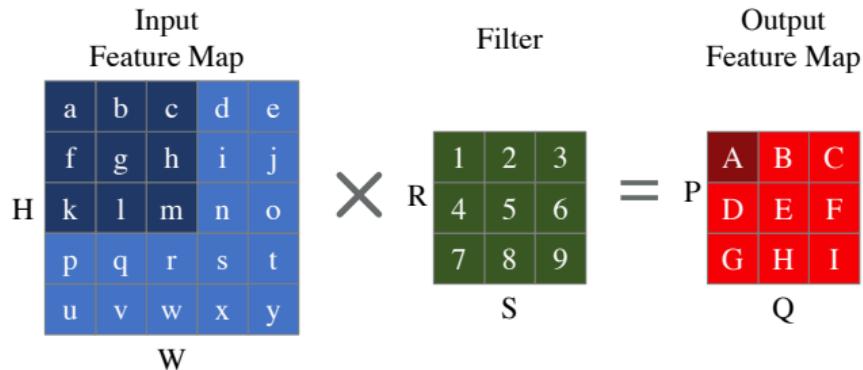
③ Direct Convolution

④ Memory Layout



Convolution Basis

2D-Convolution



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

- 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

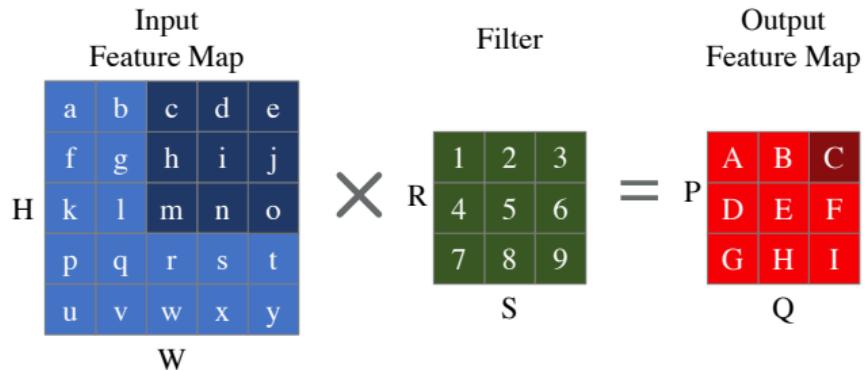
2D-Convolution



$$\begin{array}{c} \text{Input} \\ \text{Feature Map} \\ \hline \begin{matrix} a & b & c & d & e \\ f & g & h & i & j \\ k & l & m & n & o \\ p & q & r & s & t \\ u & v & w & x & y \end{matrix} \\ H \quad W \end{array} \times \begin{array}{c} \text{Filter} \\ \hline \begin{matrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{matrix} \\ R \quad S \end{array} = \begin{array}{c} \text{Output} \\ \text{Feature Map} \\ \hline \begin{matrix} A & B & C \\ D & E & F \\ G & H & I \end{matrix} \\ P \quad Q \end{array}$$

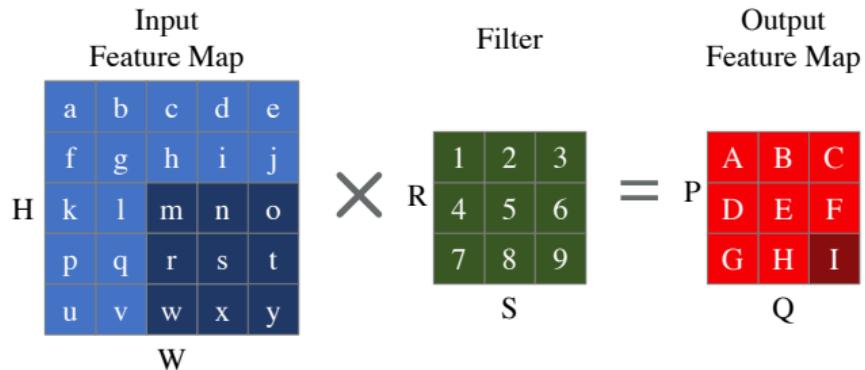
- H : Height of input feature map
- W : Width of input feature map
- R : Height of filter
- S : Width of filter
- P : Height of output feature map
- Q : Width of output feature map
- **stride**: # of rows/columns traversed per step

2D-Convolution



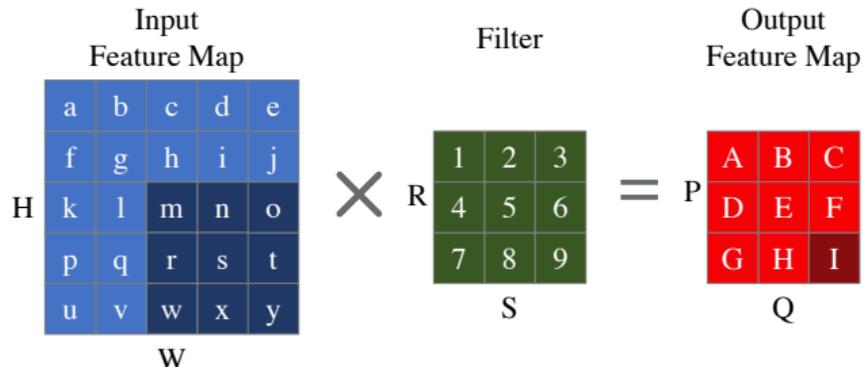
- **H**: Height of input feature map
- **W**: Width of input feature map
- **R**: Height of filter
- **S**: Width of filter
- **P**: Height of output feature map
- **Q**: Width of output feature map
- **stride**: # of rows/columns traversed per step

2D-Convolution



- **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

2D-Convolution

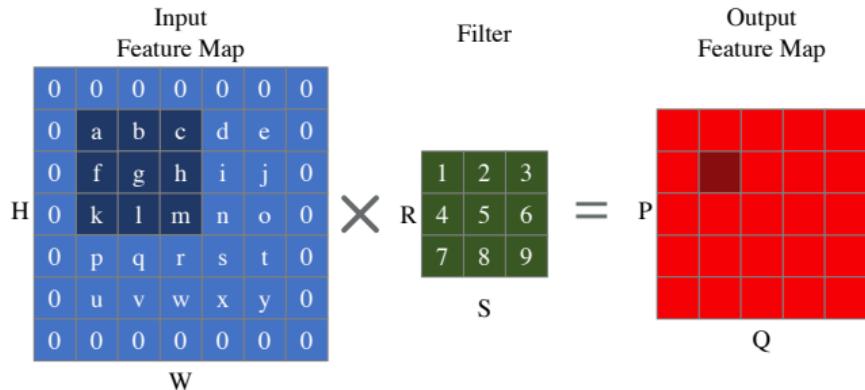


- H : Height of input feature map
- W : Width of input feature map
- R : Height of filter
- S : Width of filter
- P : Height of output feature map
- Q : Width of output feature map
- **stride**: # of rows/columns traversed per step

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

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

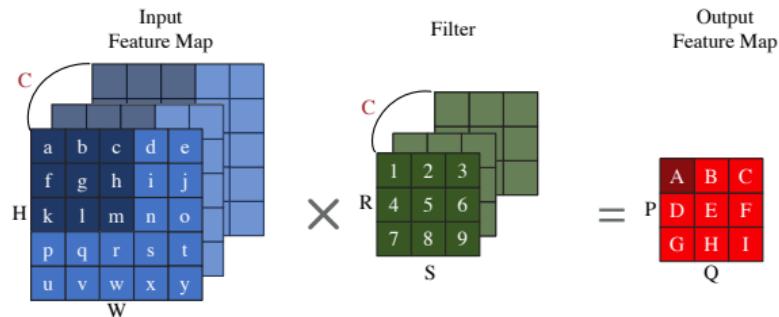
2D-Convolution



$$P = \frac{(H - R + 2 \cdot \text{pad})}{\text{stride}} + 1;$$
$$Q = \frac{(W - S + 2 \cdot \text{pad})}{\text{stride}} + 1.$$

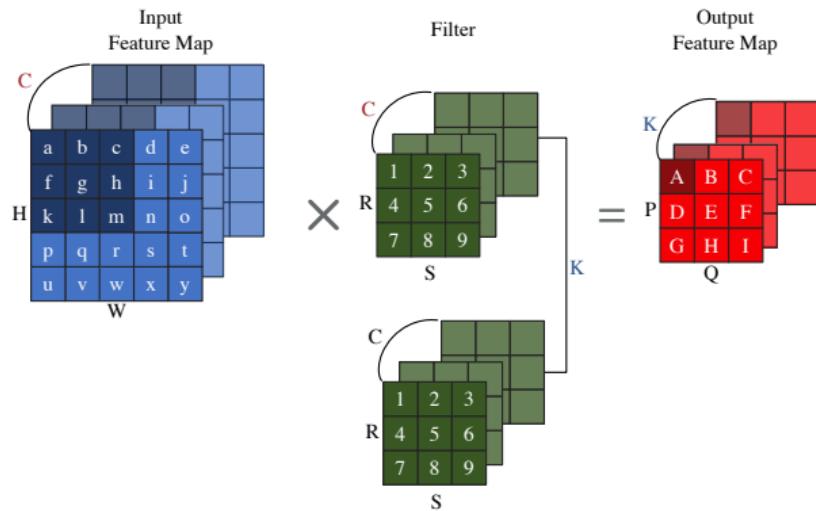
- H : Height of input feature map
- W : Width of input feature map
- R : Height of filter
- S : Width of filter
- P : Height of output feature map
- Q : Width of output feature map
- **stride**: # of rows/columns traversed per step
- **padding**: # of zero rows/columns added

3D-Convolution



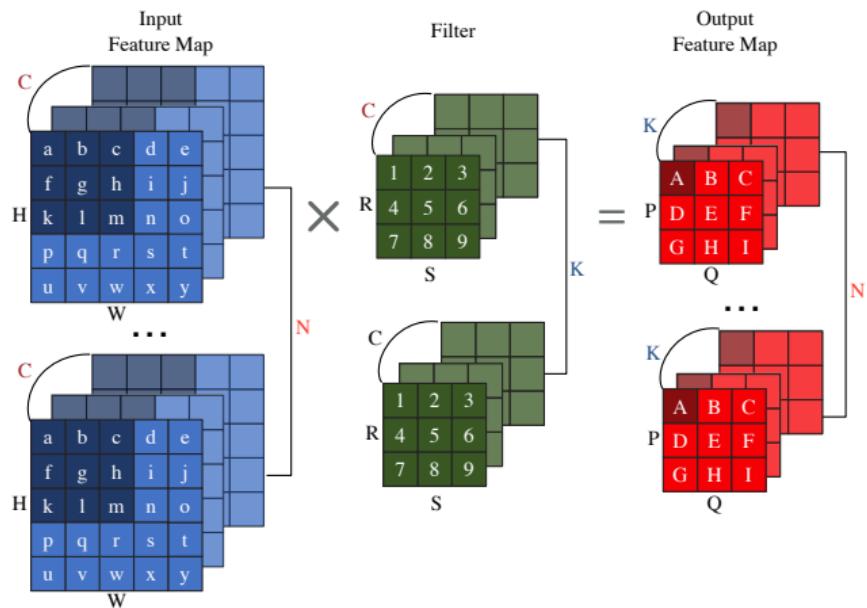
- H : Height of input feature map
- W : Width of input feature map
- R : Height of filter
- S : Width of filter
- P : Height of output feature map
- Q : Width of output feature map
- **stride**: # of rows/columns traversed per step
- **padding**: # of zero rows/columns added
- C : # of input channels

3D-Convolution



- 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

3D-Convolution



- **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



The diagram shows the computation of a 3x3 convolution kernel on a 7x7 input. The input is a 7x7 matrix of integers from 0 to 2. The kernel is a 3x3 matrix of integers from -1 to 4. The result is a 5x5 output matrix. The convolution is performed directly without padding or stride changes.

0	0	0	0	0	0	0
0	2	2	1	1	2	0
0	2	0	1	1	0	0
0	2	0	1	2	0	0
0	1	1	1	1	1	0
0	0	0	1	0	2	0
0	0	0	0	0	0	0

\times

1	0	0
1	1	1
1	0	-1

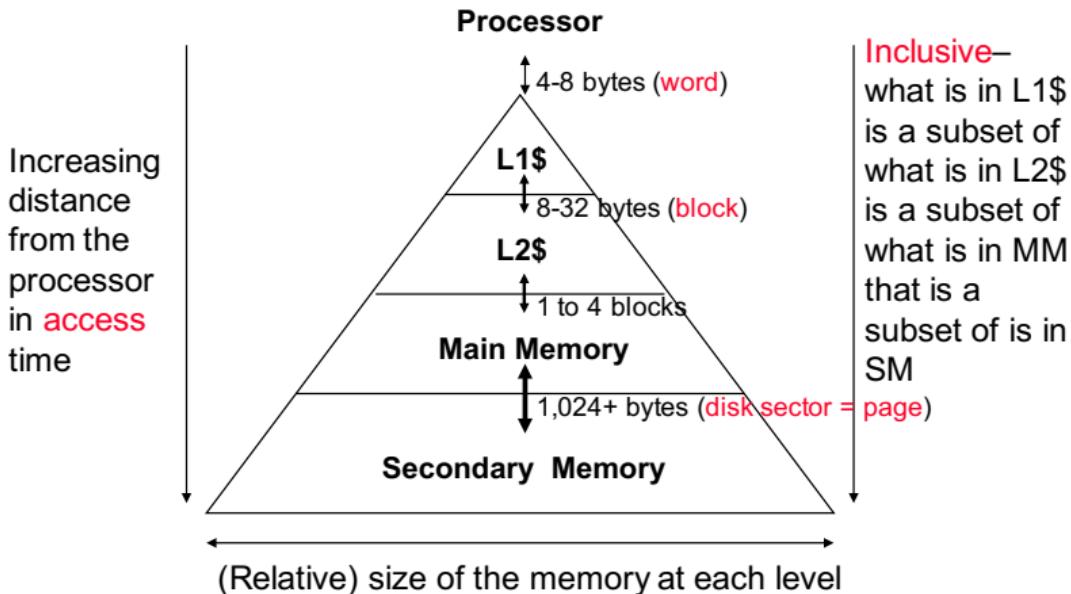
$=$

4	6	3	5	4
2	6	2	4	4
1	5	3	4	4
2	4	3	3	4
0	2	2	4	3

Direct convolution: No extra memory overhead

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

Background: Memory System



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



Im2Col

Im2col (Image2Column) Convolution

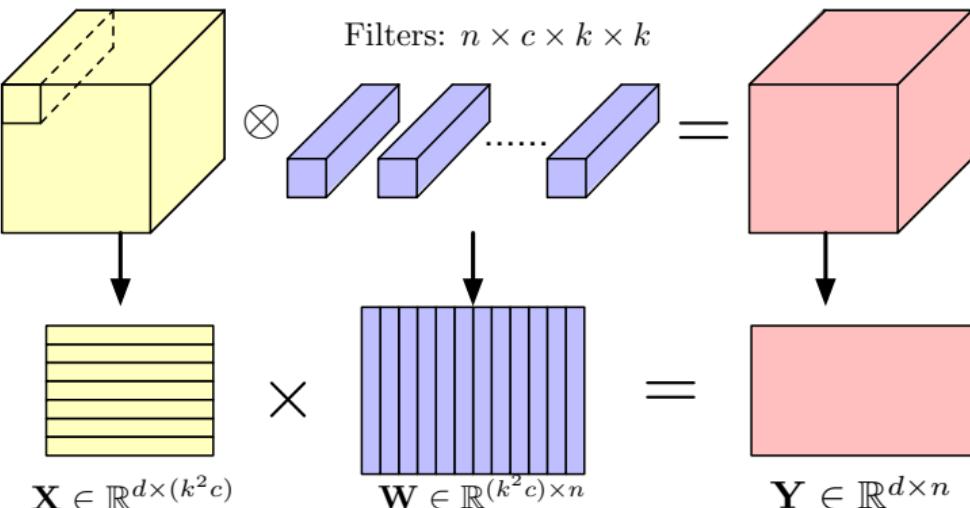


0	0	0	0	0	0	0	0
0	2	2	1	1	2	0	
0	2	0	1	1	0	0	
0	2	0	1	2	0	0	
0	1	1	1	1	1	0	
0	0	0	1	0	2	0	
0	0	0	0	0	0	0	

$$\begin{array}{c} \text{Input Image: } \\ \begin{matrix} 0 & 0 & 0 & 2 & 2 & 0 & 2 & 0 \\ 0 & 0 & 2 & 2 & 1 & 2 & 0 & 1 \\ 0 & 0 & 0 & 2 & 1 & 1 & 0 & 1 & 1 \\ \vdots & & & & & & & \\ 1 & 1 & 0 & 1 & 2 & 0 & 1 & 1 \\ \vdots & & & & & & & \\ 1 & 1 & 0 & 0 & 2 & 0 & 0 & 0 & 0 \end{matrix} \quad 25 \times 9 \\ \xrightarrow{\hspace{1cm}} \quad \times \quad = \\ \text{Filter: } \\ \begin{matrix} 1 \\ 0 \\ 0 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ -1 \end{matrix} \quad 9 \times 1 \end{array}$$

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

Im2col (Image2Column) Convolution



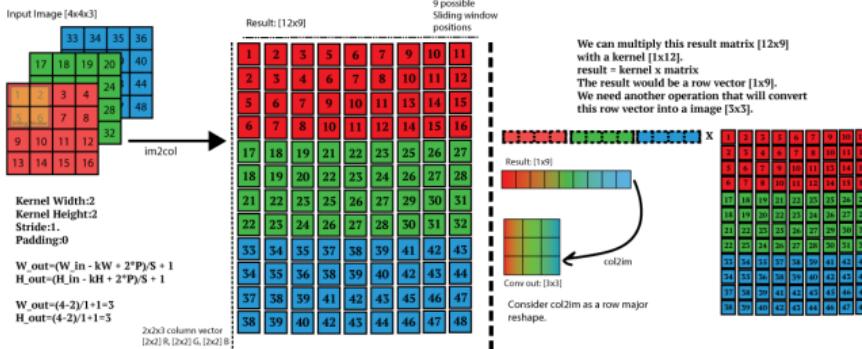
- Transform convolution to **matrix multiplication**
- **Unified** calculation for both convolution and fully-connected layers



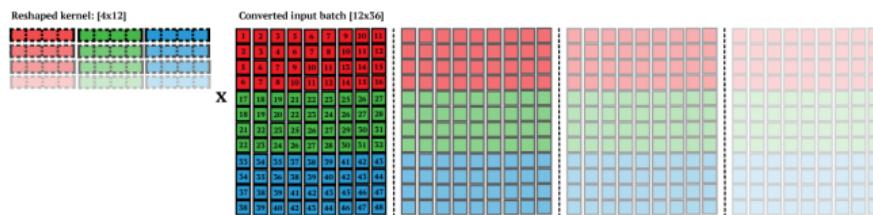
Im2col (Image2Column): Another View

Image to column operation (im2col)

Slide the input image like a convolution but each patch become a column vector.



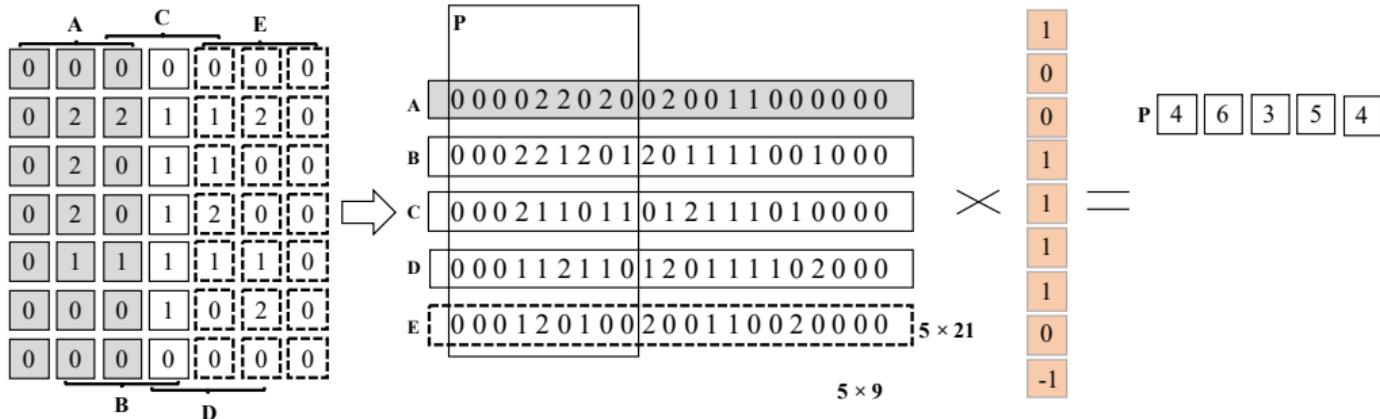
We get true performance gain when the kernel has a large number of filters, ie: F=4 and/or you have a batch of images (N=4). Example for the input batch [4x4x3x4], convolved with 4 filters [2x2x3x2]. The only problem with this approach is the amount of memory



1

¹https://leonardoaraujosantos.gitbook.io/artificial-intelligence/machine_learning/deep_learning/convolution_layer/making_faster

SOTA 1: Memory-efficient Convolution

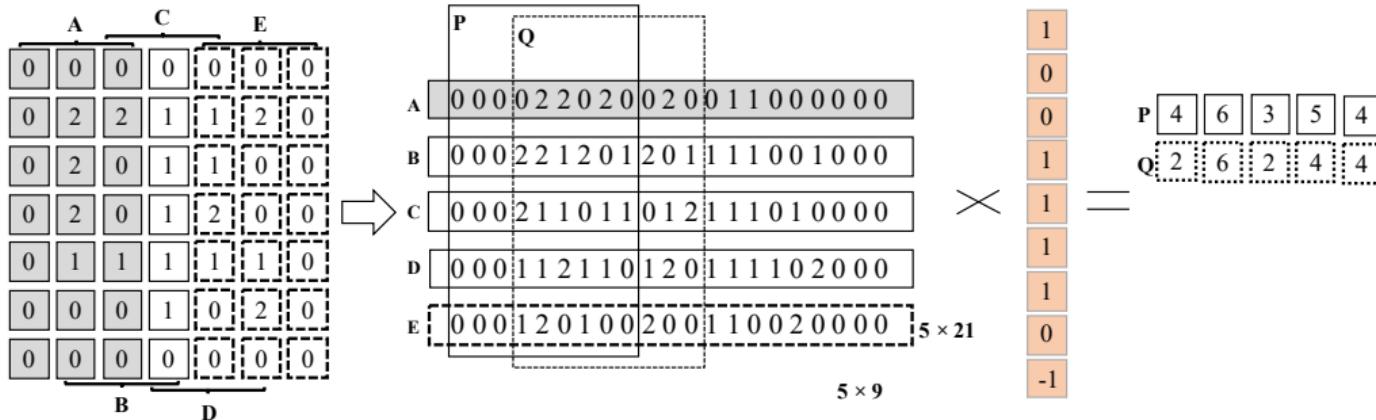


2

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

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

SOTA 1: Memory-efficient Convolution

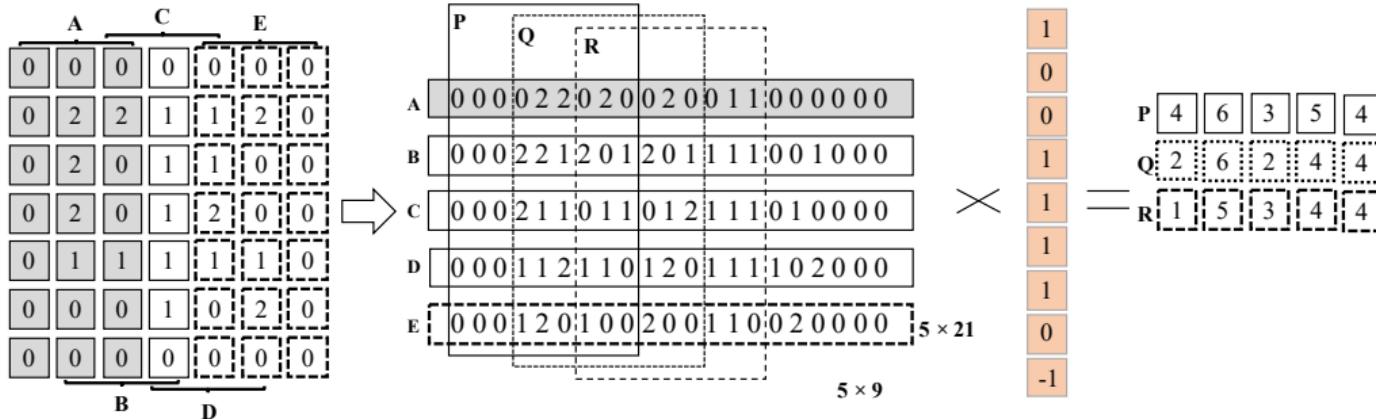


2

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

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

SOTA 1: Memory-efficient Convolution

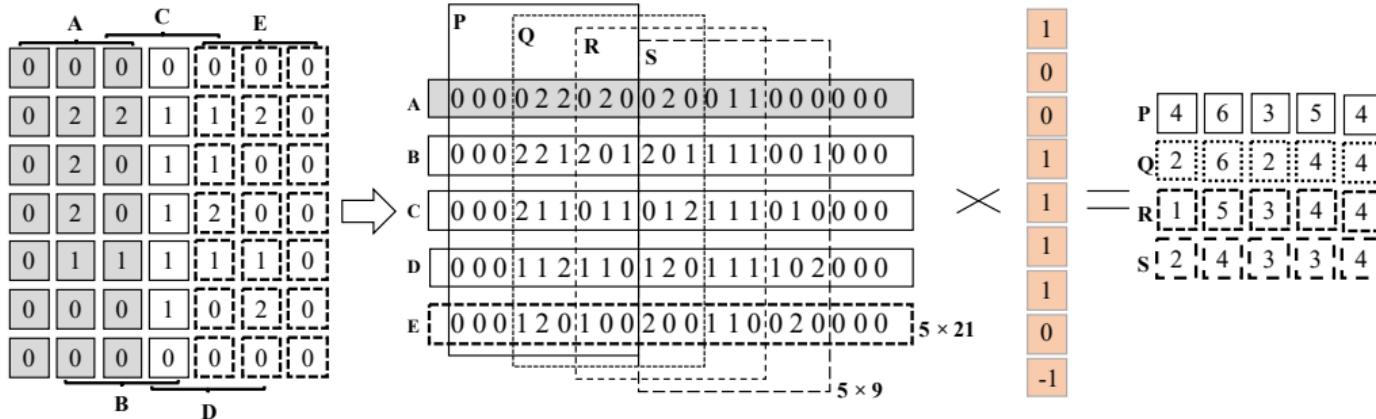


2

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

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

SOTA 1: Memory-efficient Convolution

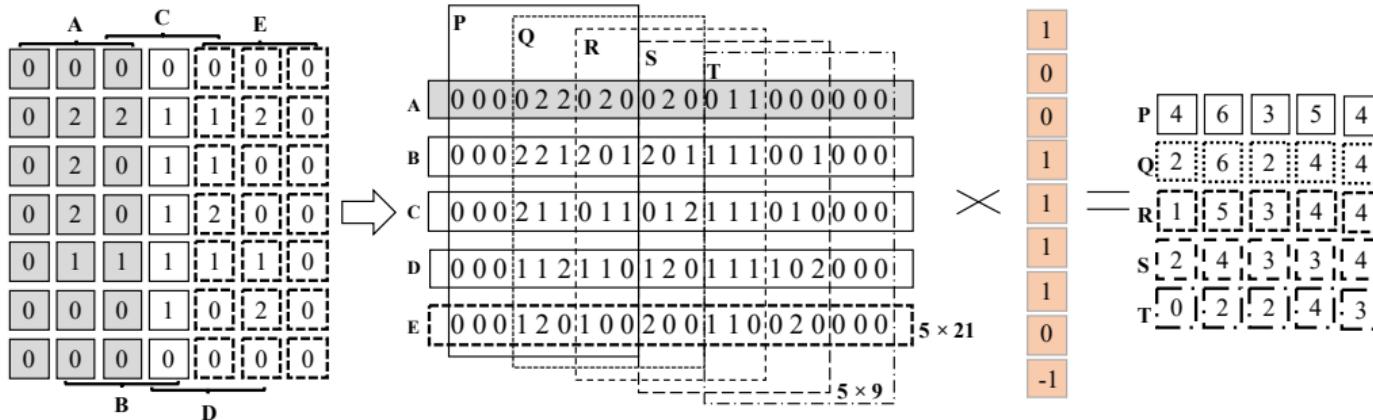


2

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

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

SOTA 1: Memory-efficient Convolution



2

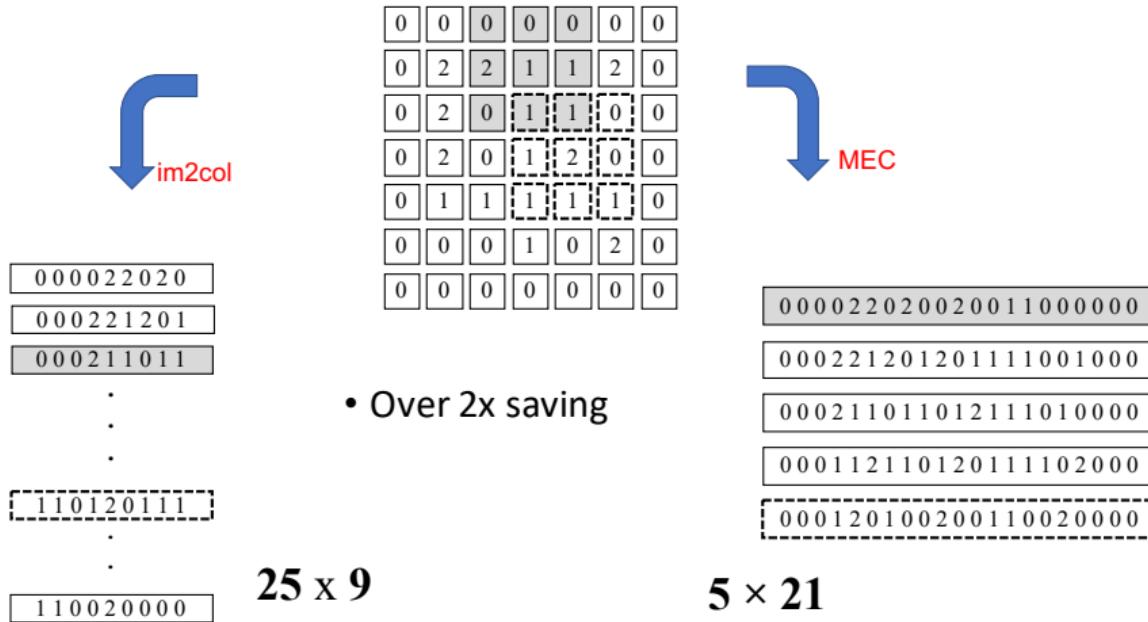
- Sub matrices in the lowered matrix will be “**sgemm**” ed in parallel
- Smaller memory foot print, cache locality, and explicit parallelism

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

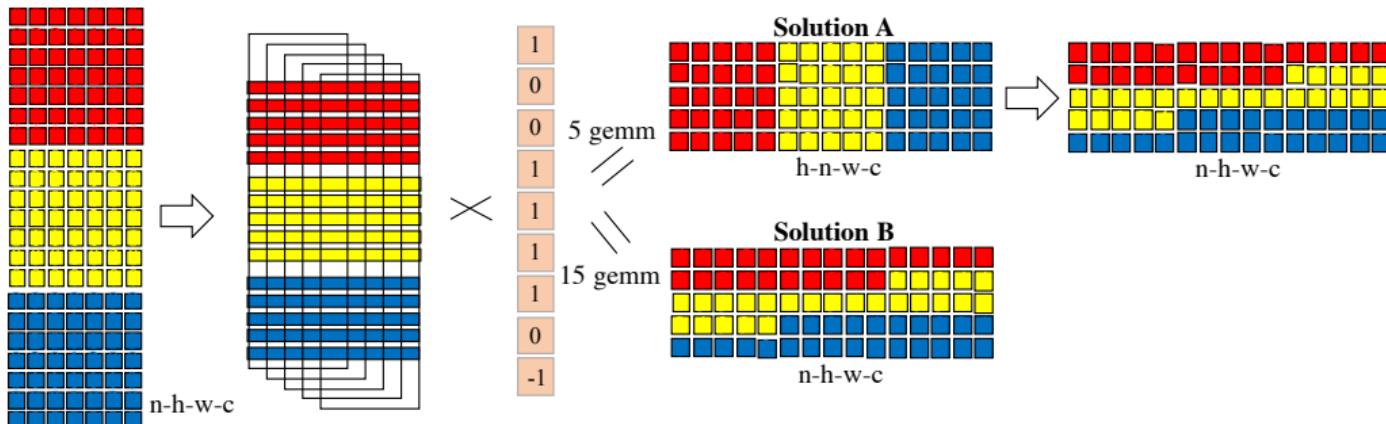
SOTA 1: Memory-efficient Convolution



Over $2\times$ memory saving³:



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

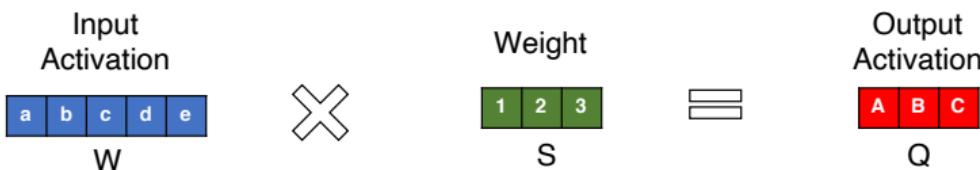


- MEC w. mini-batch: can use $n-h-w-c$ format
- Fusing convolution+pooling can be another solution



Direct Convolution

1D Convolution Example



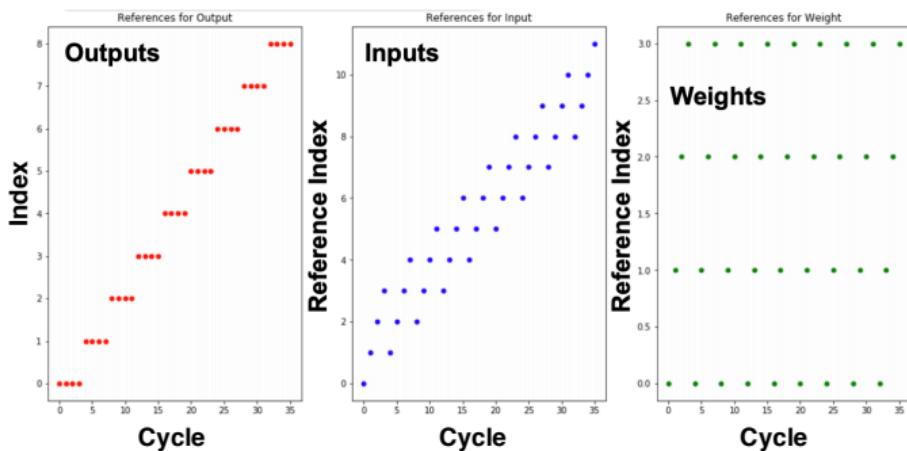
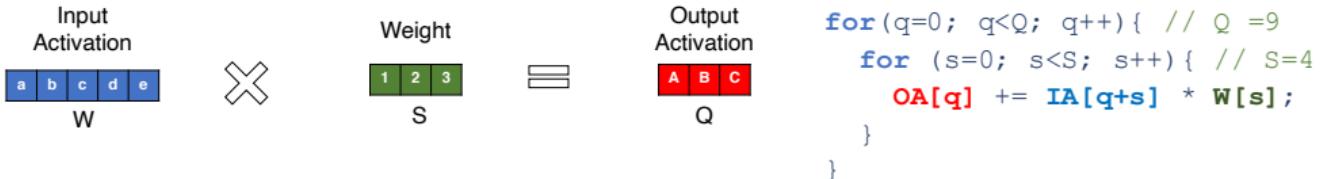
```
for(q=0; q<Q; q++) {  
    for (s=0; s<S; s++) {  
        OA[q] += IA[q+s] * W[s];  
    }  
}
```

**Output Stationary (OS)
Dataflow**

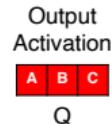
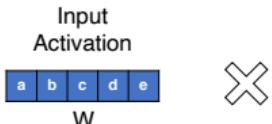
```
for (s=0; s<S; s++) {  
    for (q=0; q<Q; q++) {  
        OA[q] += IA[q+s] * W[s];  
    }  
}
```

**Weight Stationary (WS)
Dataflow**

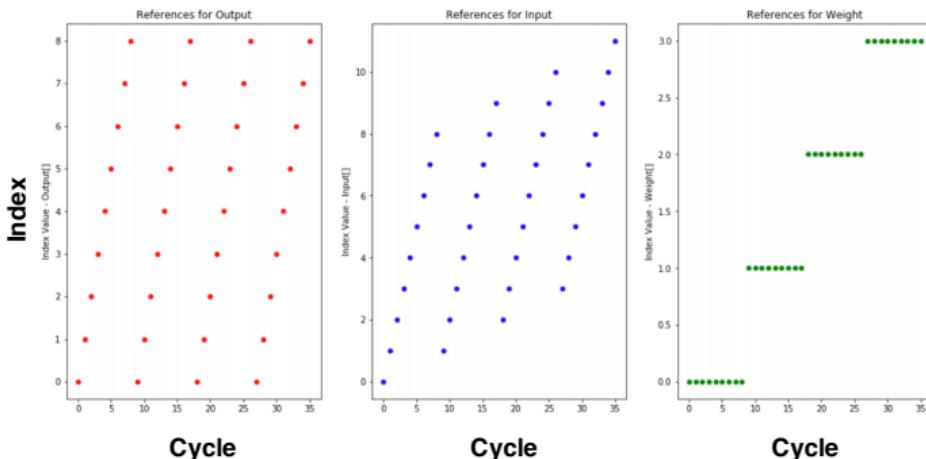
Buffer Access Pattern 1: Output Stationary



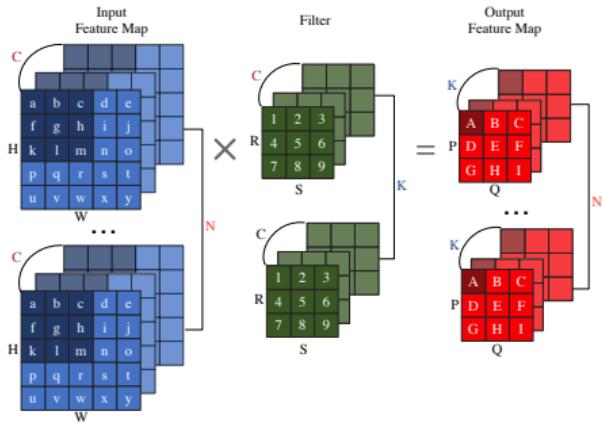
Buffer Access Pattern 2: Weight Stationary



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

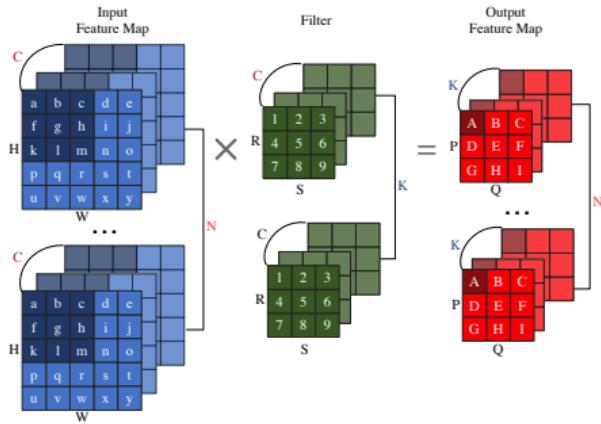


Direct Convolution



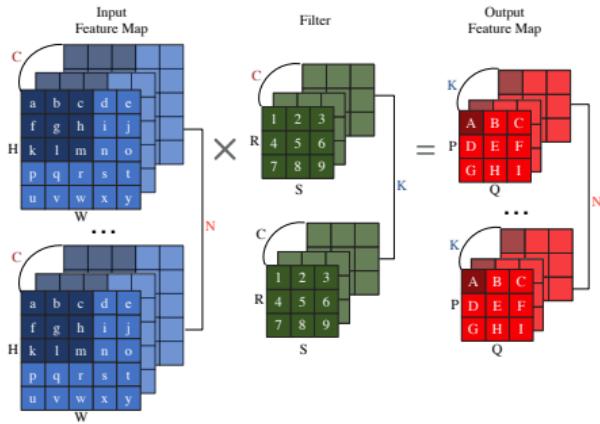
```
1   for (n=0; n<N; n++) {
2     for (k=0; k<K; k++) {
3       for (p=0; p<P; p++) {
4         for (q=0; q<Q; q++) {
5           OA[n][k][p][q] = 0;
6           for (r=0; r<R; r++) {
7             for (s=0; s<S; s++) {
8               for (c=0; c<C; c++) {
9                 h = p * stride - pad +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



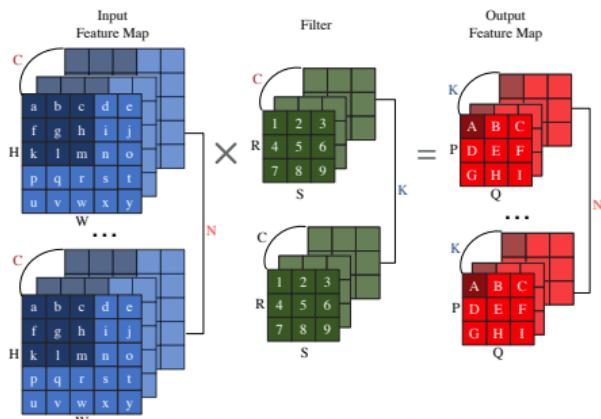
```
1  for (n=0; n<N; n++) {  
2      for (r=0; r<R; r++) {  
3          for (s=0; s<S; s++) {  
4              for (c=0; c<C; c++) {  
5                  for (k=0; k<K; k++) {  
6                      float curr_w = W[r][s][c][k];  
7                      for (p=0; p<P; p++) {  
8                          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                } } } } } }
```



Memory Layout



Data Layout Formats⁴

- **N** is the batch size
- **C** is the number of feature maps
- **H** is the image height
- **W** is the image width

EXAMPLE

N = 1

C = 64

H = 5

W = 4

c = 0

0	1	2	3
4	5	6	7
8	9	10	11
12	13	14	15
16	17	18	19

c = 1

20	21	22	23
24	25	26	27
28	29	30	31
32	33	34	35
36	37	38	39

c = 2

40	41	42	43
44	45	46	47
48	49	50	51
52	53	54	55
56	57	58	59

c = 30

600	601	602	603
604	605	606	607
608	609	610	611
612	613	614	615
616	617	618	619

c = 31

620	621	622	623
624	625	626	627
628	629	630	631
632	633	634	635
636	637	638	639

c = 32

640	641	642	643
644	645	646	647
648	649	650	651
652	653	654	655
656	657	658	659

...

c = 62

1240	1241	1242	1243
1244	1245	1246	1247
1248	1249	1250	1251
1252	1253	1254	1255
1256	1257	1258	1259

c = 63

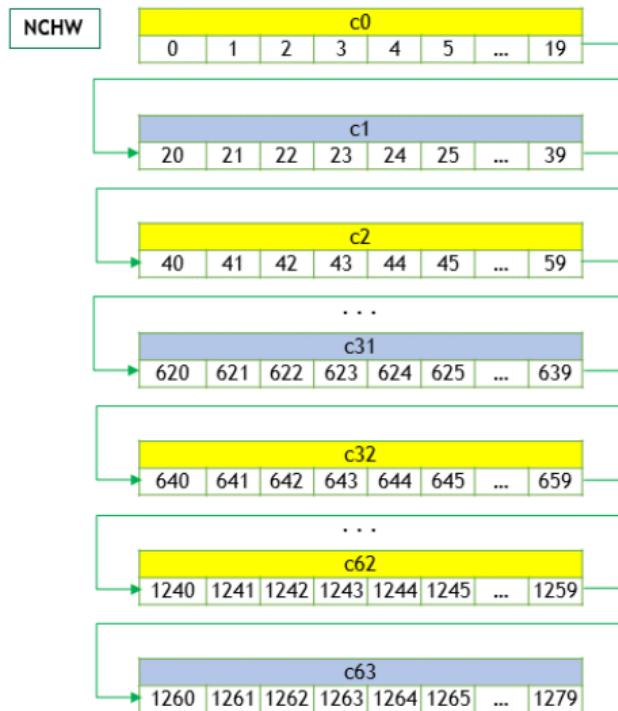
1260	1261	1262	1263
1264	1265	1266	1267
1268	1269	1270	1271
1272	1273	1274	1275
1276	1277	1278	1279

⁴<https://docs.nvidia.com/deeplearning/cudnn/developer-guide/index.html>

NCHW Memory Layout



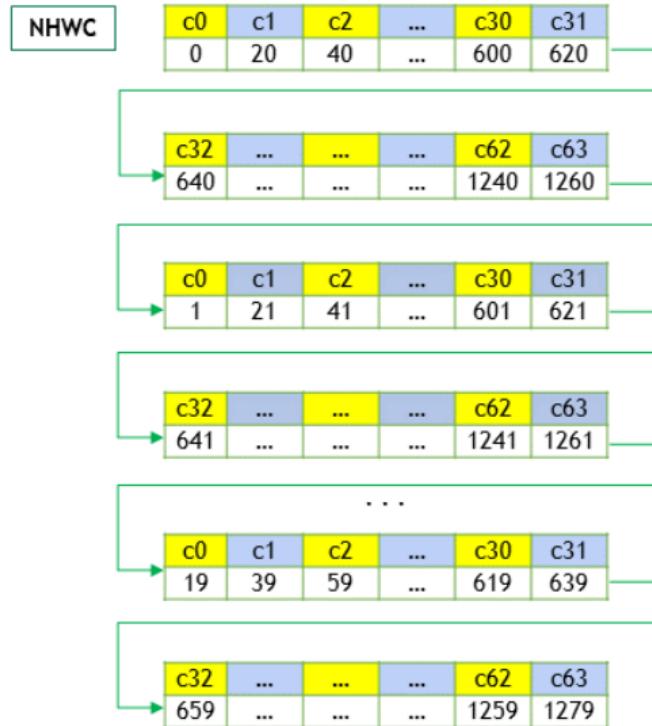
- Begin with first channel ($c=0$), elements arranged contiguously in row-major order
- Continue with second and subsequent channels until all channels are laid out



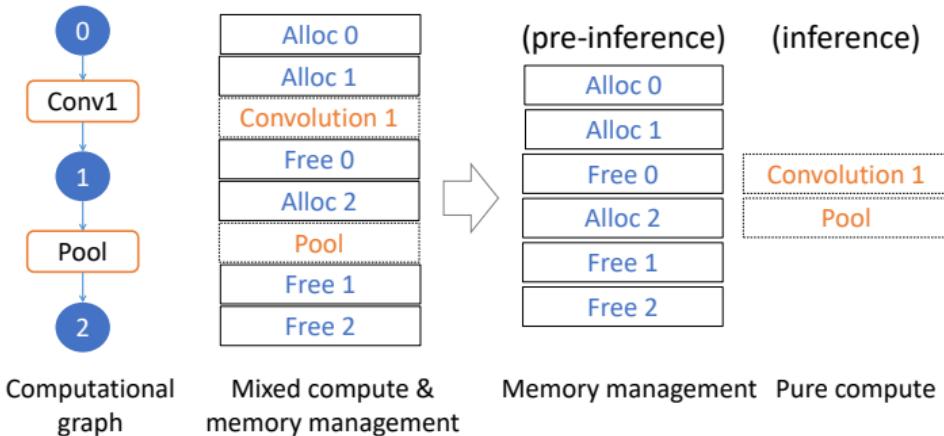


- Begin with the first element of channel 0, then proceed to the first element of channel 1, and so on, until the first elements of all the C channels are laid out
- Next, select the second element of channel 0, then proceed to the second element of channel 1, and so on, until the second element of all the channels are laid out
- Follow the row-major order of channel 0 and complete all the elements
- Proceed to the next batch (if N is > 1)

NHWC Memory Layout



Memory optimization of MNN



- MNN can infer the exact required memory for the entire graph:
 - virtually walking through all operations
 - summing up all allocation and freeing



- 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/>