# Implementation 01: GEMM

Bei Yu
CSE Department, CUHK
byu@cse.cuhk.edu.hk

(Latest update: September 2, 2024)

2024 Fall

# Convolution Basis

| | Input Feature Map | Filter | 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

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

Input Feature Map · Filter · 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

Input Feature Map

Filter

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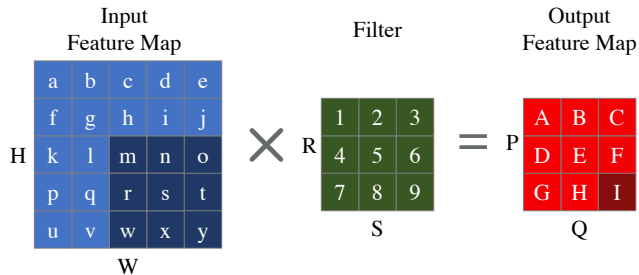
- H: Height of input feature map
- W: Width of input feature map
- R: Height of filter
- S: Width of filter
- P: Height of output feature map
- Q: Width of output feature map
- stride: # of rows/columns traversed per step

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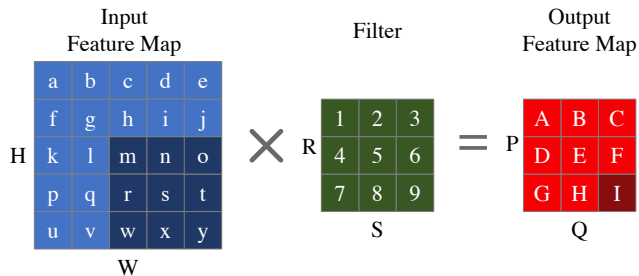- P: Height of output feature map
- Q: Width of output feature map
- stride: # of rows/columns traversed per step

- 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

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

- 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
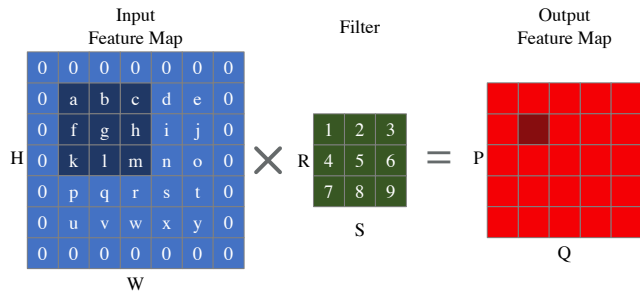- stride: # of rows/columns traversed per step
- padding: # of zero rows/columns added

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

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

- 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
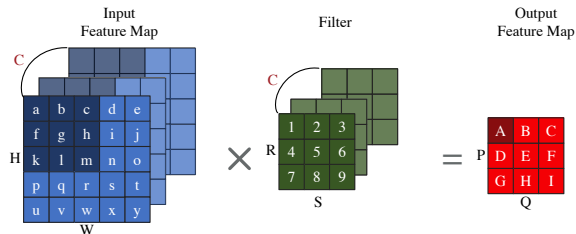
- stride: # of rows/columns traversed per step

- padding: # of zero rows/columns added

- C: # of input channels
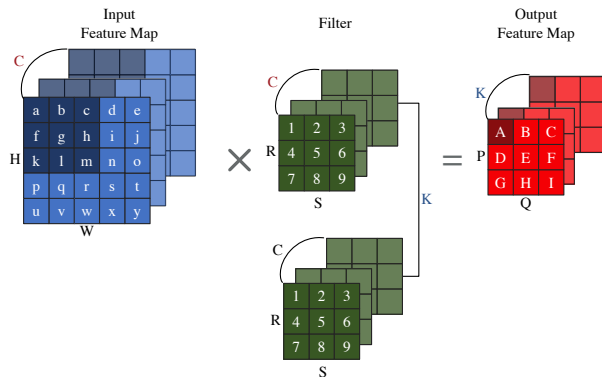
- H: Height of input feature map
- W: Width of input feature map
- R: Height of filter
- S: Width of filter
- P: Height of output feature map
- Q: Width of output feature map
- stride: # of rows/columns traversed per step
- padding: # of zero rows/columns added
- C: # of input channels
- K: # of output 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
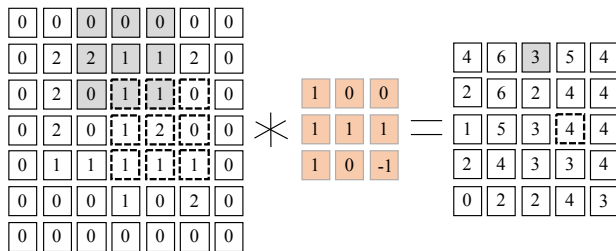- padding: # of zero rows/columns added
- C: # of input channels
- K: # of output channels
- N: Batch size
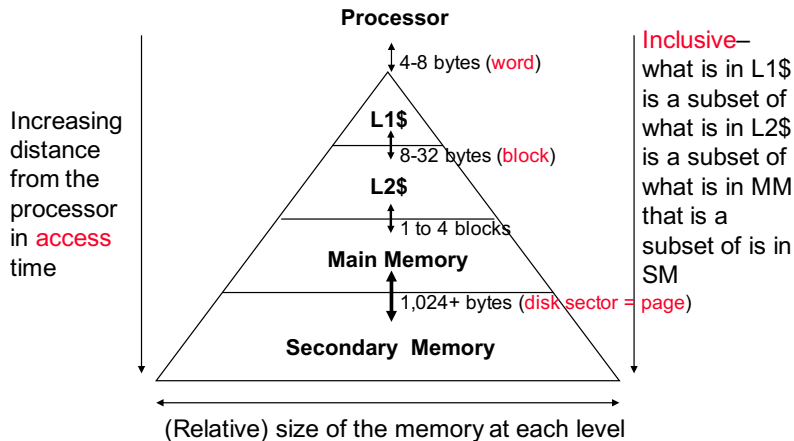
Direct convolution: No extra memory overhead

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

**Processor**

↓ 4-8 bytes (word)

Increasing distance from the processor in access time

**L1$**

↕ 8-32 bytes (block)

**L2$**

↕ 1 to 4 blocks

**Main Memory**

↕ 1,024+ bytes (disk sector = page)

**Secondary Memory**

(Relative) size of the memory at each level

Inclusive– what is in L1$ is a subset of what is in L2$ is a subset of what is in MM that is a subset of is in SM

- Spatial locality
- Temporal Locality

# Im2Col

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

- Input channel #: 2
- Output channel #: 1

Filters: $n \times c \times k \times k$

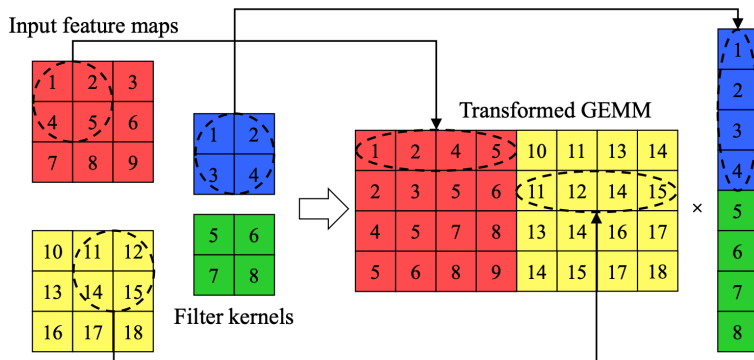$\mathbf{X} \in \mathbb{R}^{d \times (k^2 c)}$     $\mathbf{W} \in \mathbb{R}^{(k^2 c) \times n}$     $\mathbf{Y} \in \mathbb{R}^{d \times n}$

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

# Memory-efficient Convolution

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

---

[1]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

---

[1]Minsik Cho and Daniel Brand (2017). "MEC: memory-efficient convolution for deep neural network". In: *Proc. ICML*.
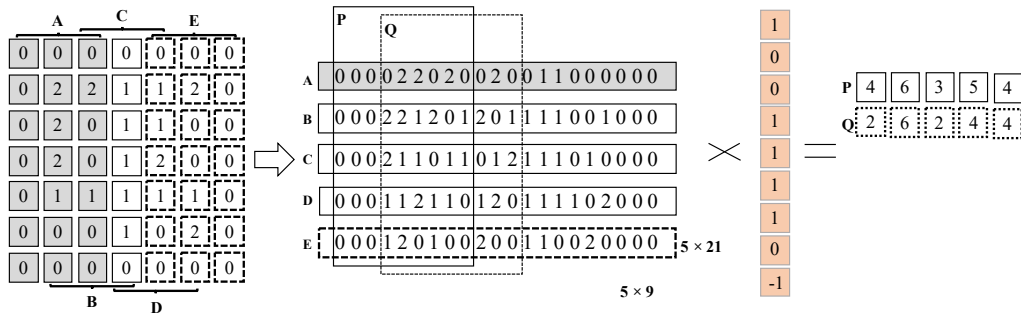
- Sub matrices in the lowered matrix will be "sgemm" ed in parallel
- Smaller memory foot print, cache locality, and explicit parallelism

---

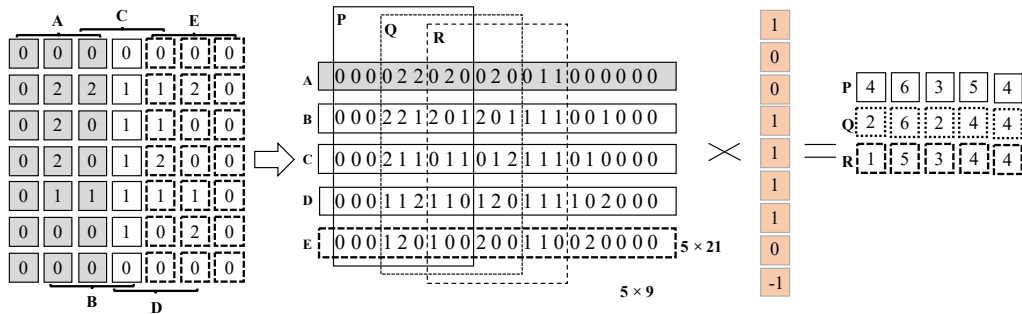[1]Minsik Cho and Daniel Brand (2017). "MEC: memory-efficient convolution for deep neural network". In: *Proc. ICML*.

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

---

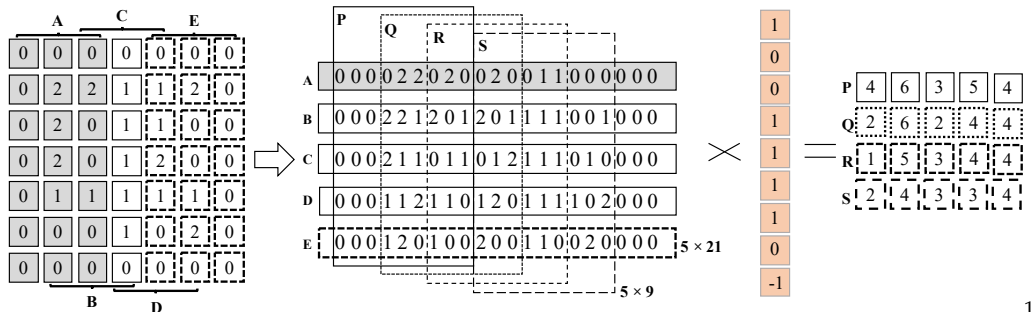[1] Minsik Cho and Daniel Brand (2017). "MEC: memory-efficient convolution for deep neural network". In: *Proc. ICML*.

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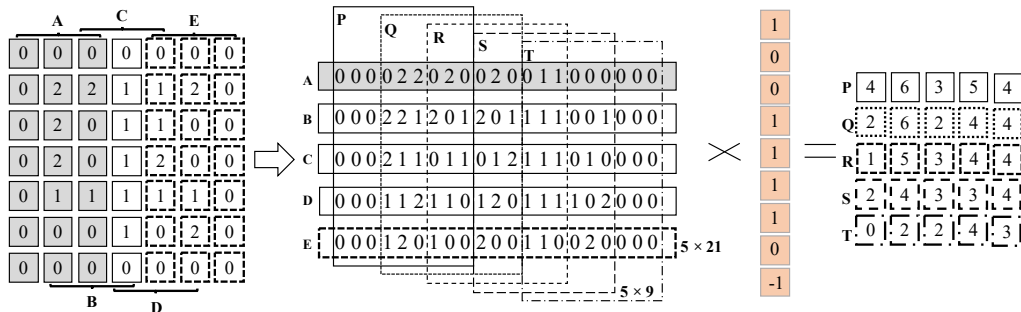- Smaller memory foot print, cache locality, and explicit parallelism

---

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

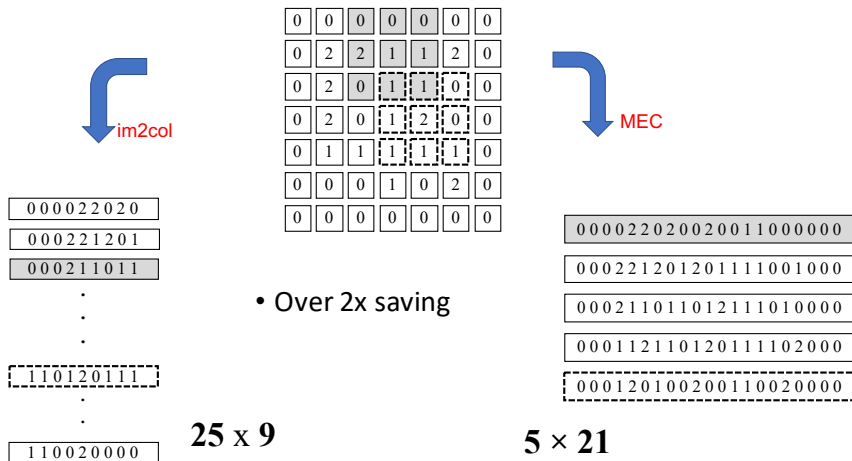Over $2\times$ memory saving[2]:



• Over 2x saving

**25** x **9**          **5** × **21**

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

# Memory Layout

# Data Layout Formats[3]

- 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

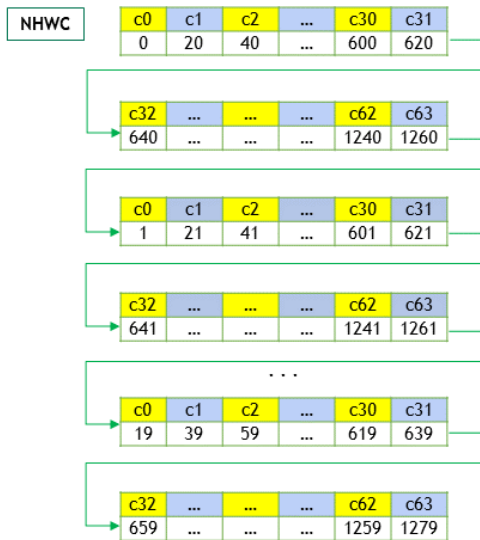| 1260 | 1261 | 1262 | 1263 |
|---|---|---|---|
| 1264 | 1265 | 1266 | 1267 |
| 1268 | 1269 | 1270 | 1271 |
| 1272 | 1273 | 1274 | 1275 |
| 1276 | 1277 | 1278 | 1279 |

. . .

# 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

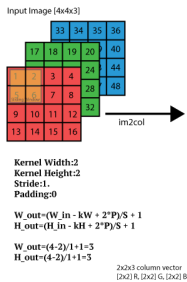- Proceed to the next batch (if $N$ is > 1)
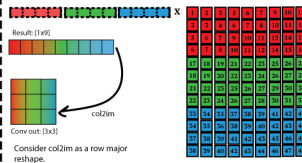
# Memory Layout in `Im2col`



Image to column operation (im2col)
Slide the input image like a convolution but each patch become a column vector.
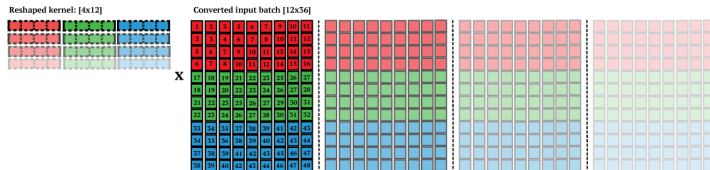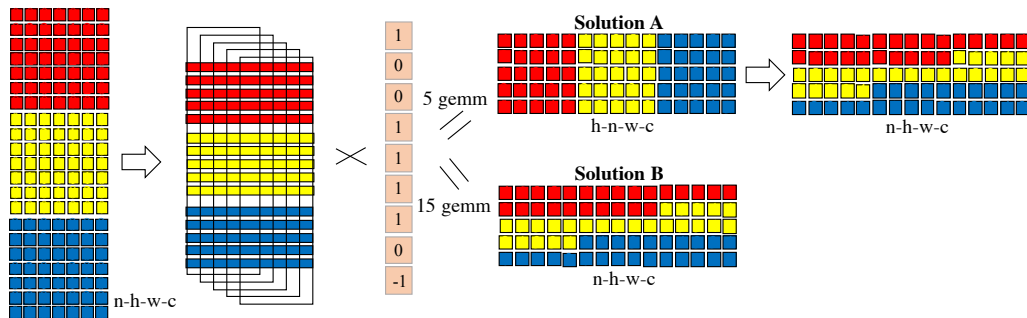
We can multiply this result matrix [12x9] with a kernel [1x12].
result = kernel x matrix
The result would be a row vector [1x9].
We need another operation that will convert this row vector into a image [3x3].

Kernel Width:2
Kernel Height:2
Stride:1.
Padding:0

$W\_out=(W\_in - kW + 2*P)/S + 1$
$H\_out=(H\_in - kH + 2*P)/S + 1$

$W\_out=(4-2)/1+1=3$
$H\_out=(4-2)/1+1=3$

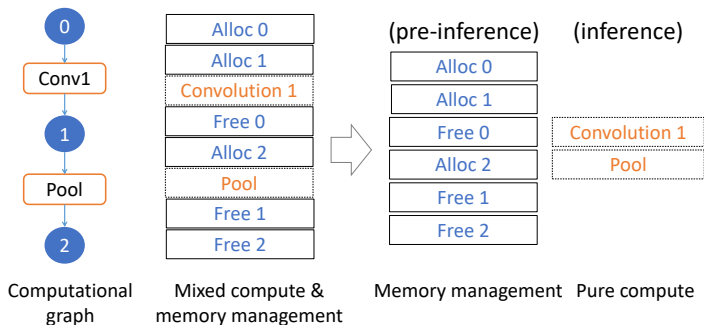Consider col2im as a row major reshape.

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

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

Computational graph — Mixed compute & memory management — Memory management — Pure compute

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