

High Performance CNN Inference

CS 217 Stanford
Ardavan Pedram

Analysis of Convolution Approaches

1. GEMM (default)

Pros:

- Generic and stable
- Easy to implement (problem mapped into a BLAS call)
- Optimized solution if good BLAS is provided

Cons:

- Additional memory to store the intermediate data
- Rely heavily on optimized BLAS

2. Spatial domain

Pros:

- Avoids additional memory copy
- Speedy with optimized code

Cons:

- Rely on individually optimized kernels according to given params, or even given HW architecture

3. FFT domain

Pros:

- Lower computational complexity

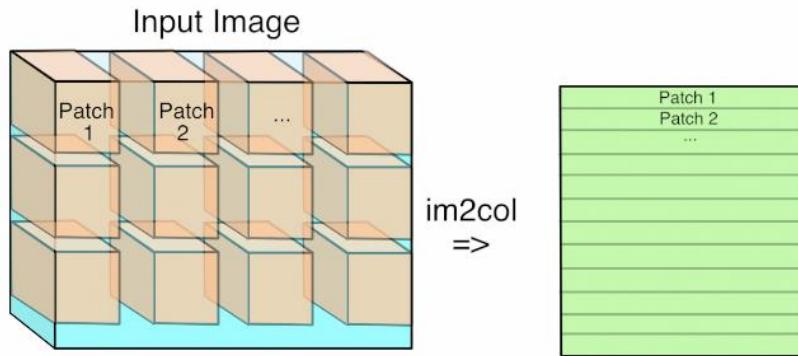
Cons:

- Additional memory to save FFT data
- Overhead is big for small kernel size, or large stride

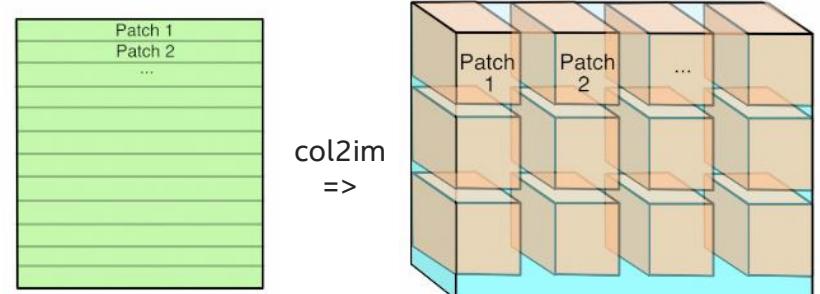
GEMM Based Convolution

Flatten input data and kernels, solve the convolution as a matrix multiplication problem

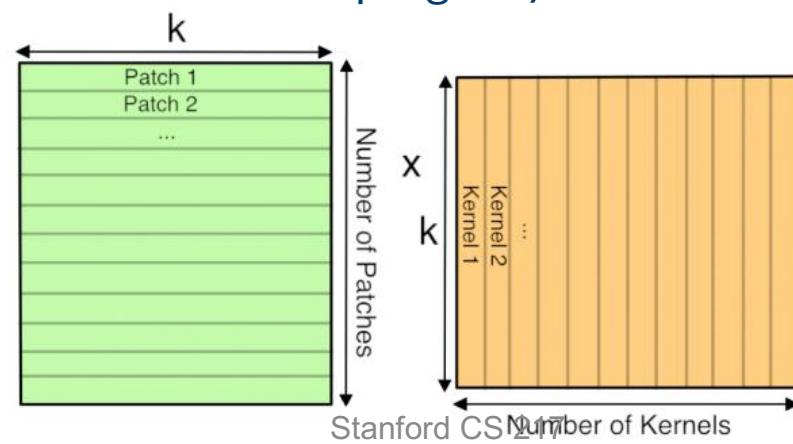
Step1: data flattening



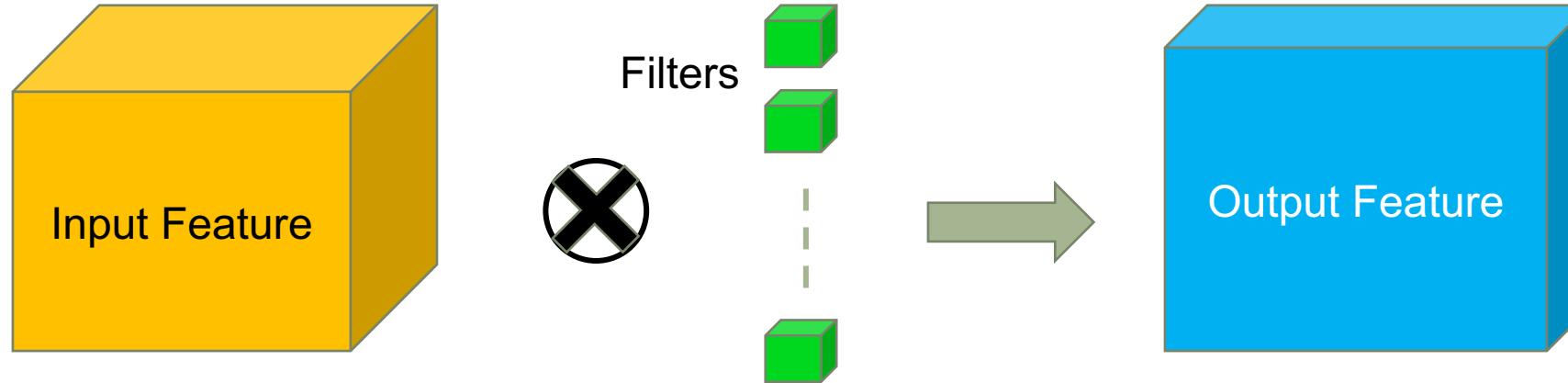
Step3: data unflattening



Step2: matrix multiply
Usually mapped into a BLAS
(Basic Linear Algebra
Subprogram) call

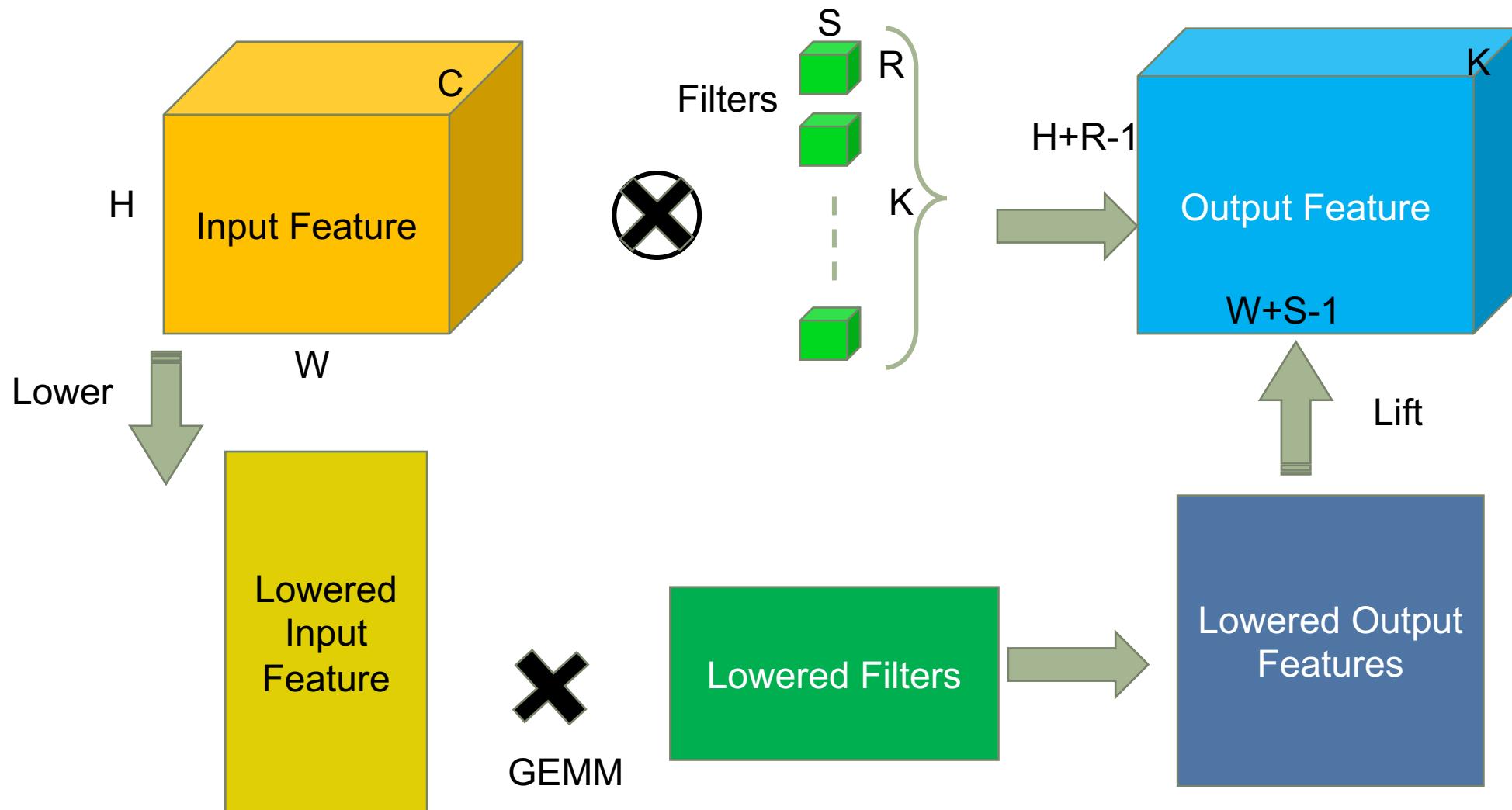


Caffe Con Troll [Hadjis et. al.]

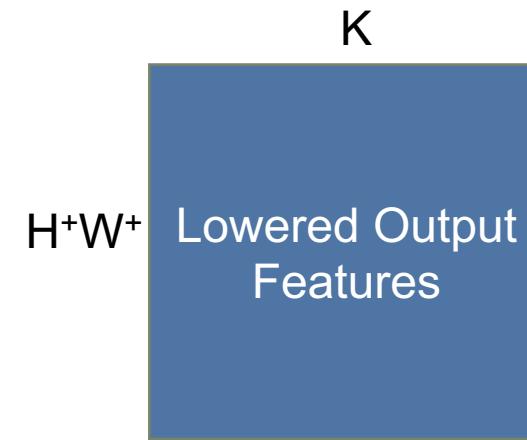
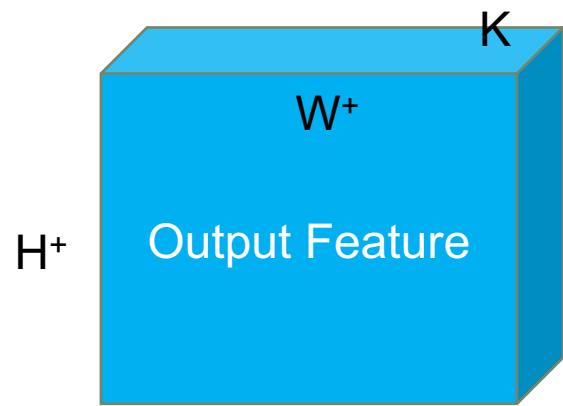


Can We turn the problem into GEMM?

Im2col

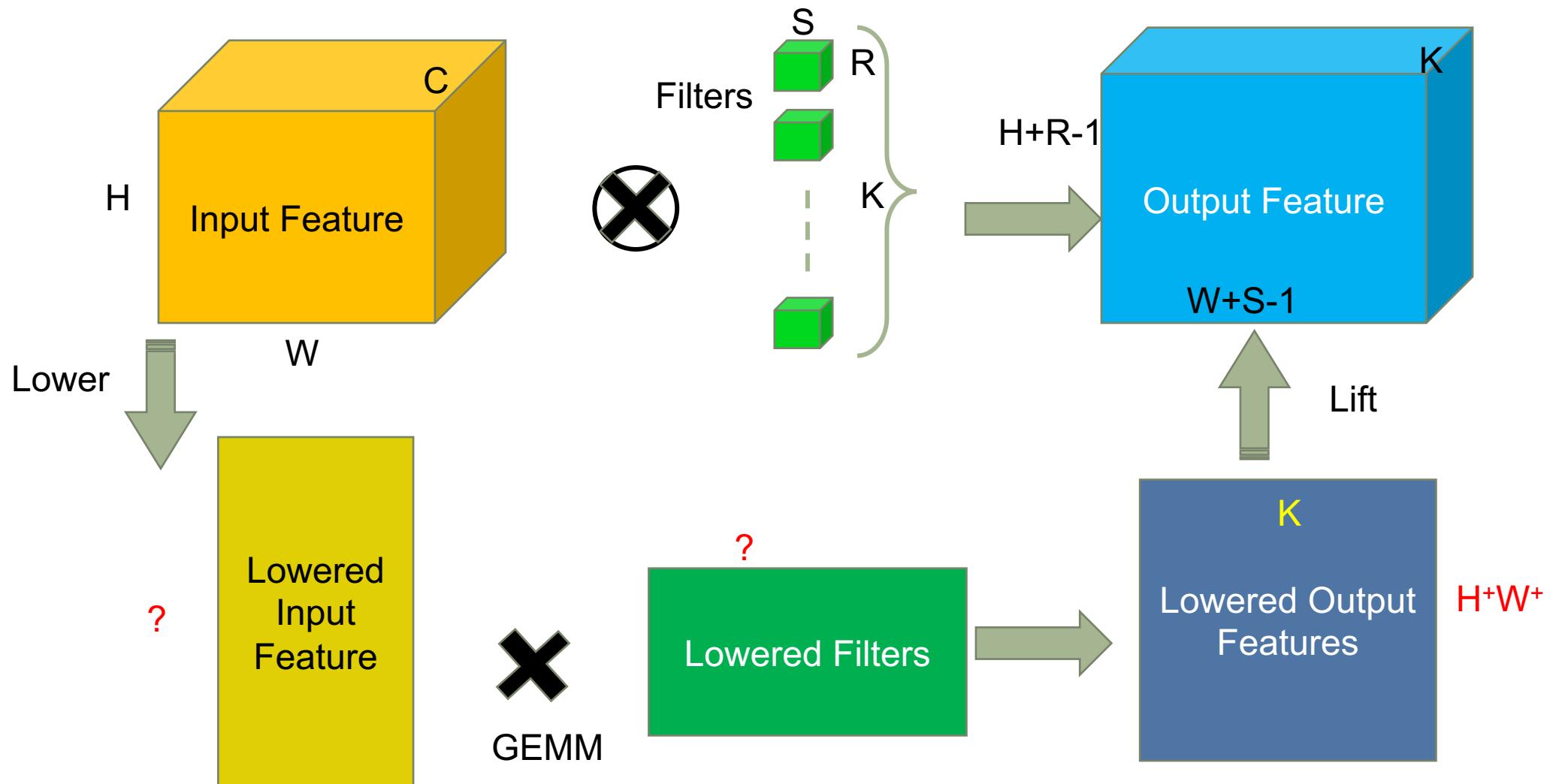


Turn a Tensor into Matrix

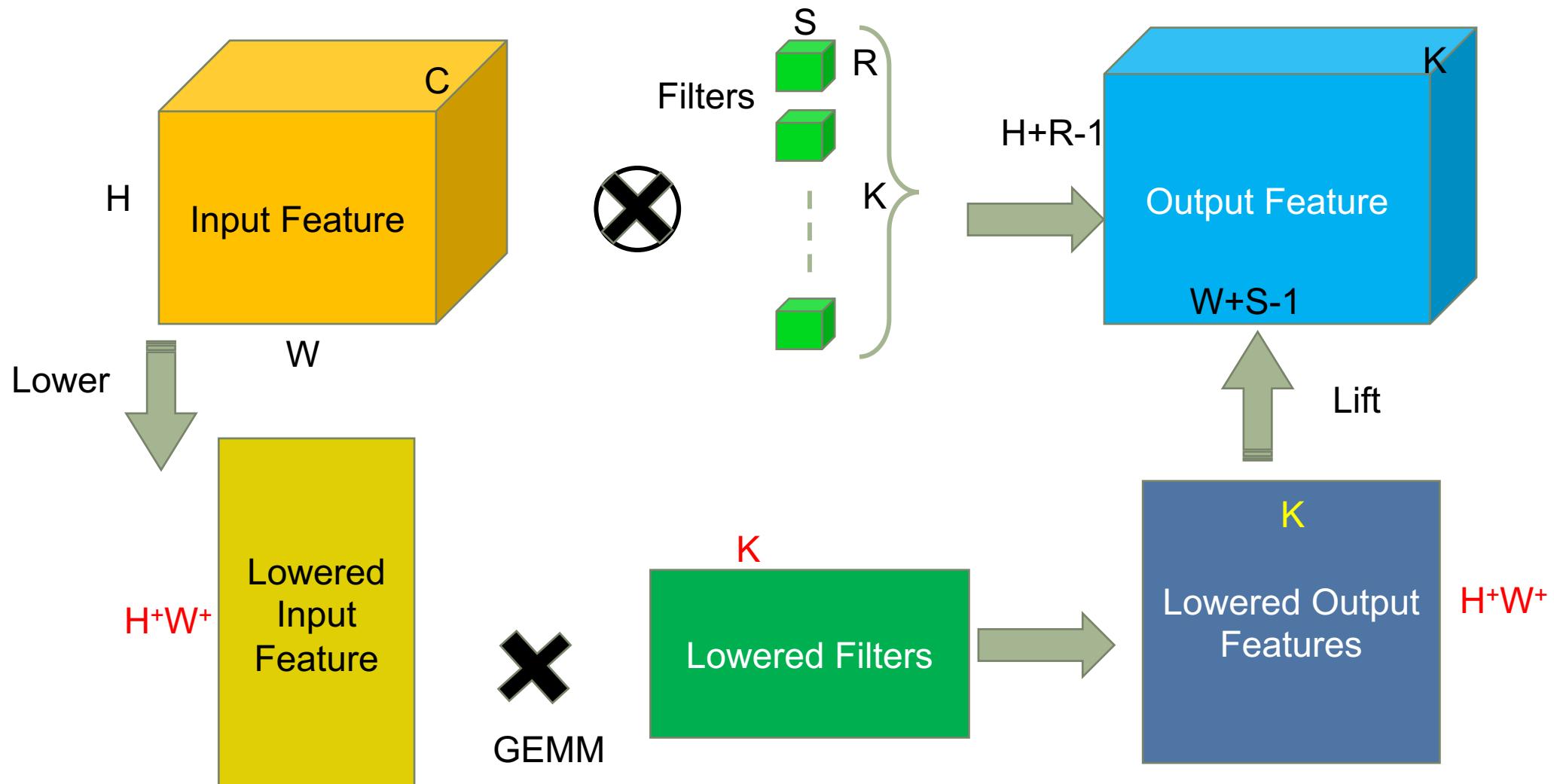


Remember What was a reduction dimension?

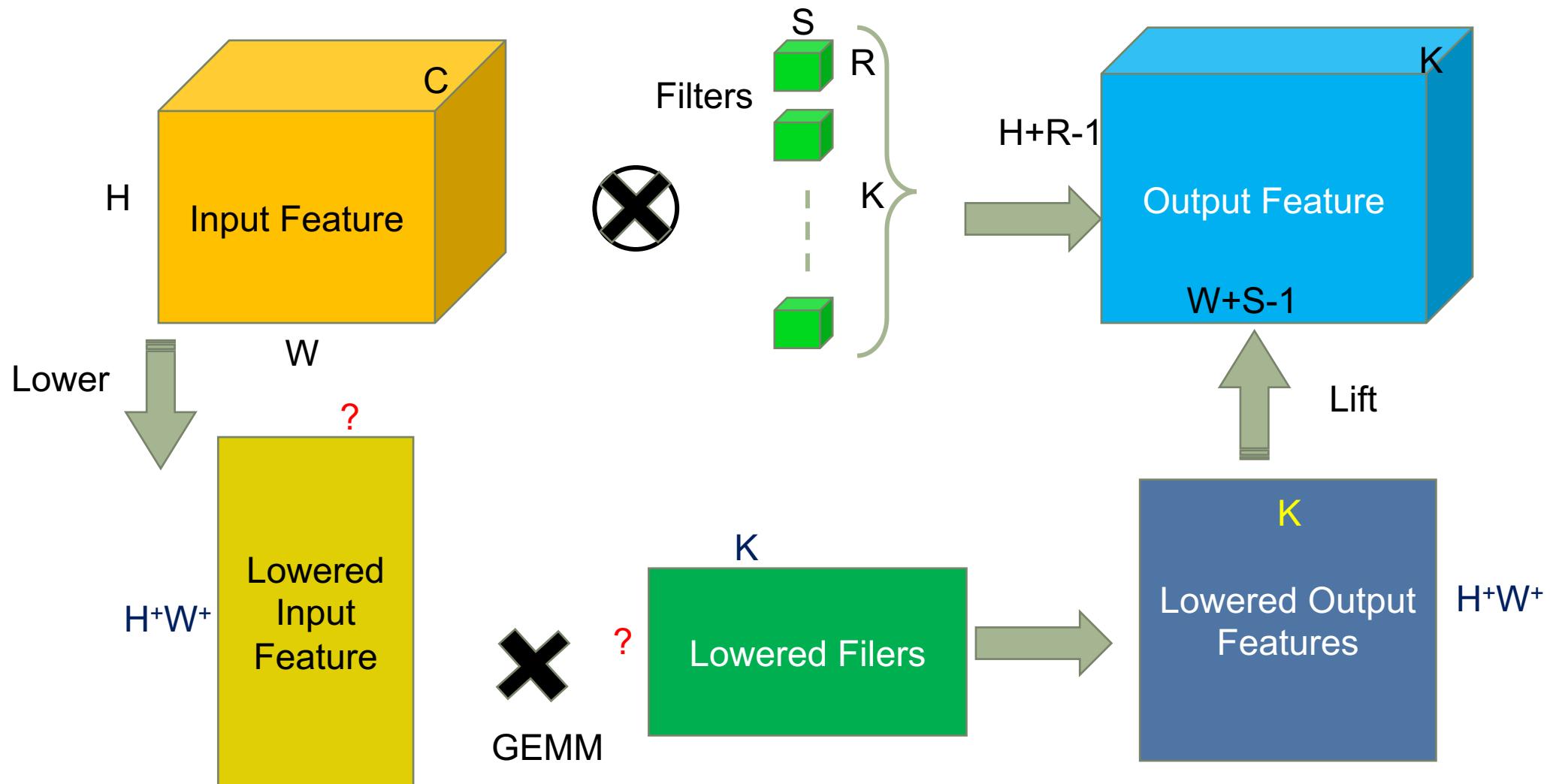
Im2col



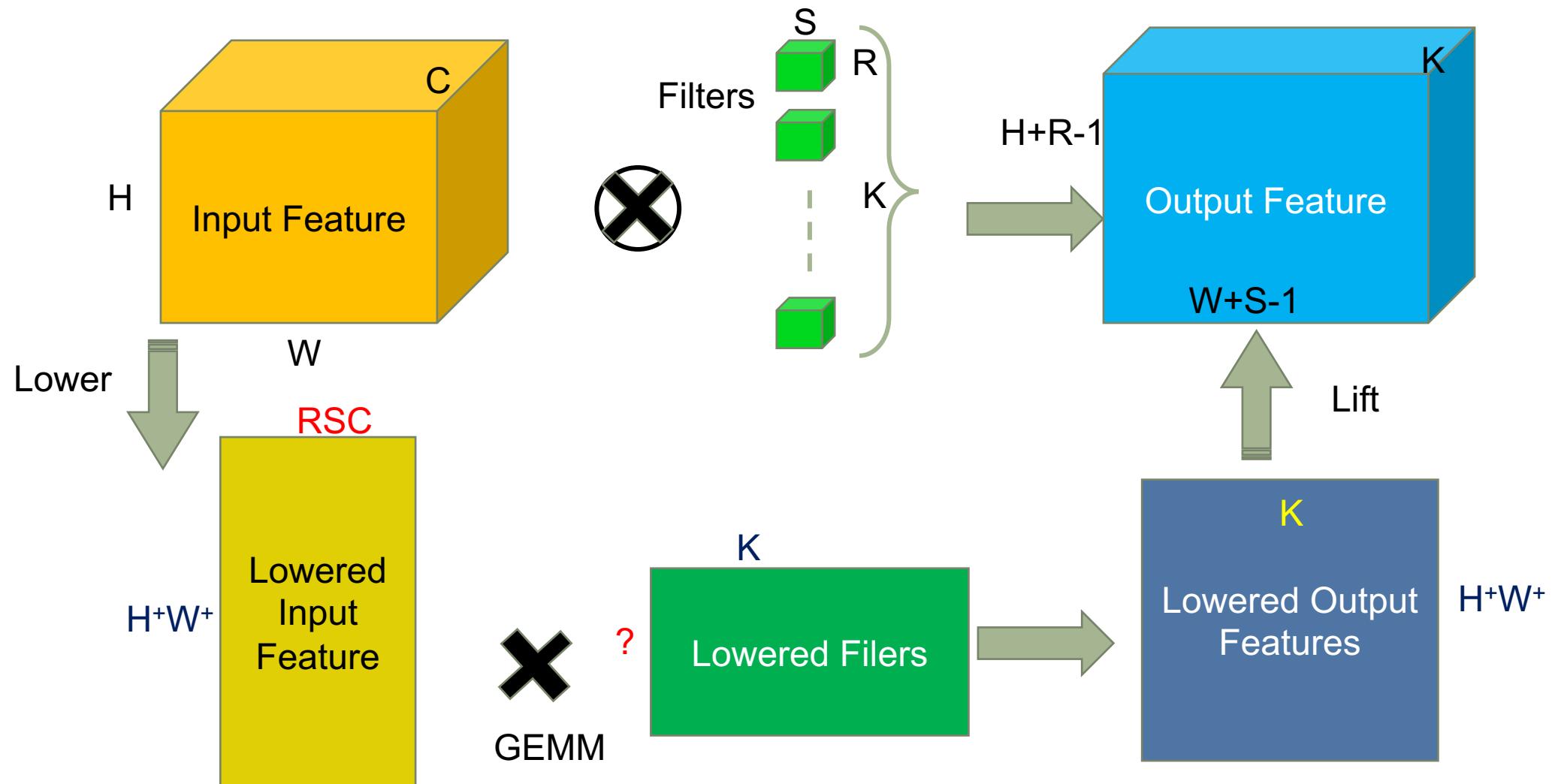
Im2col



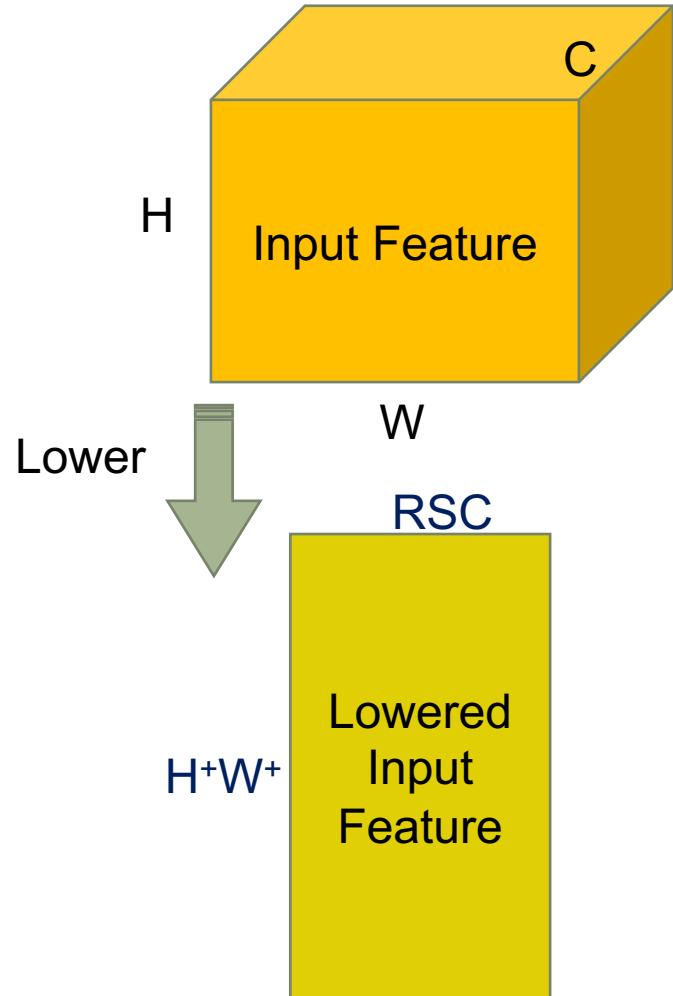
Im2col



Im2col



What Really Happened?



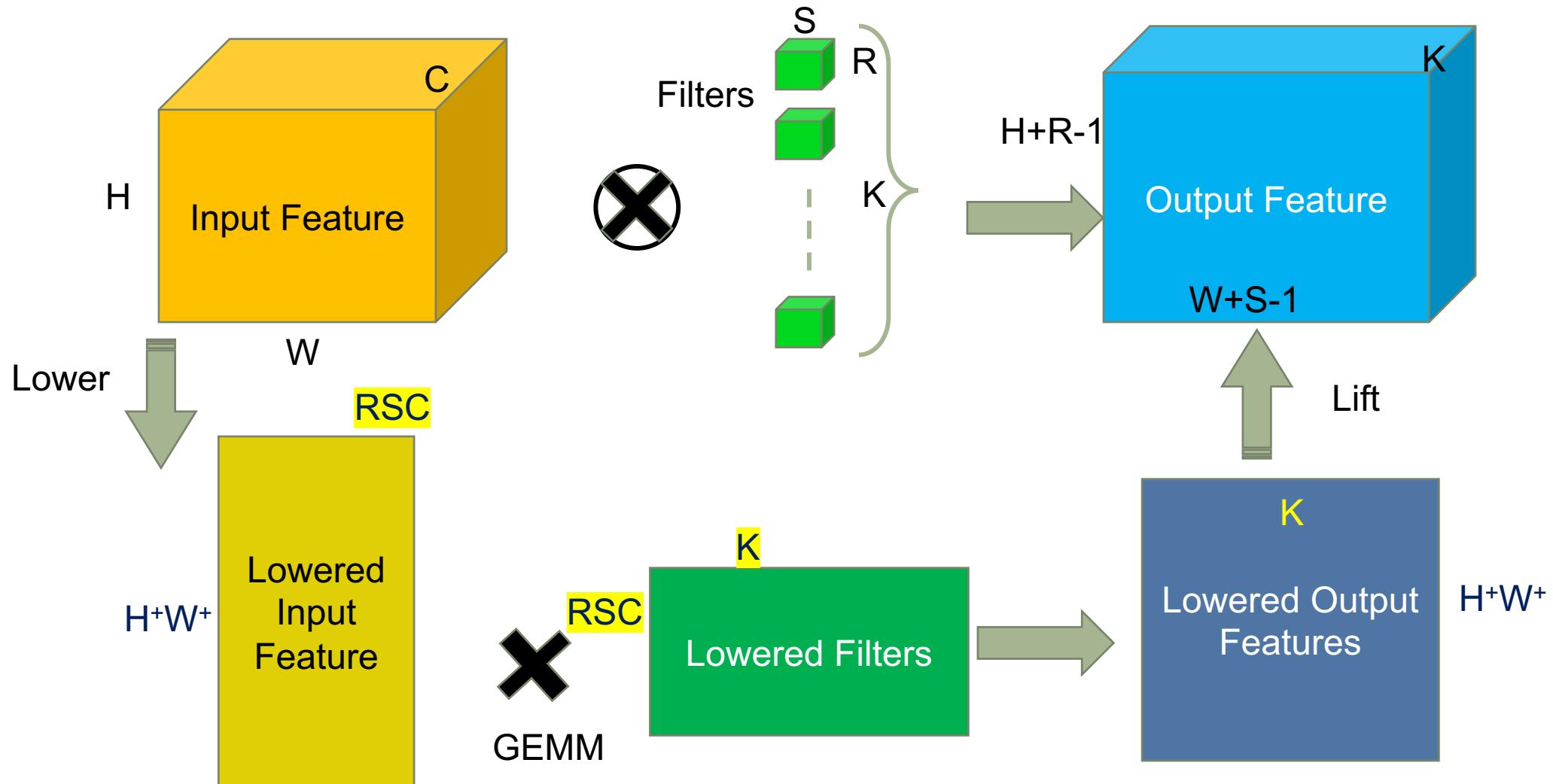
What was the tensor size?

What is the lowered matrix size?

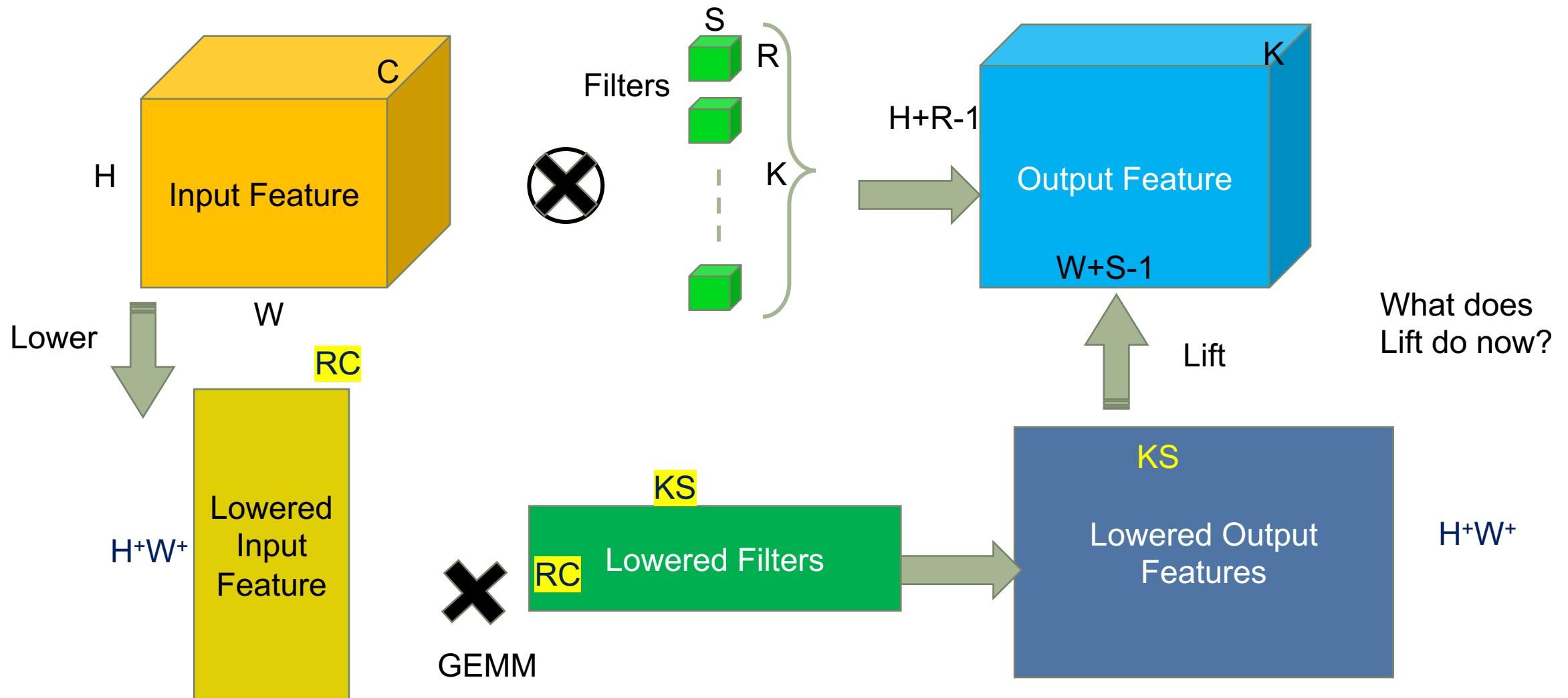
What does lowering mean?

What is the impact on Memory?

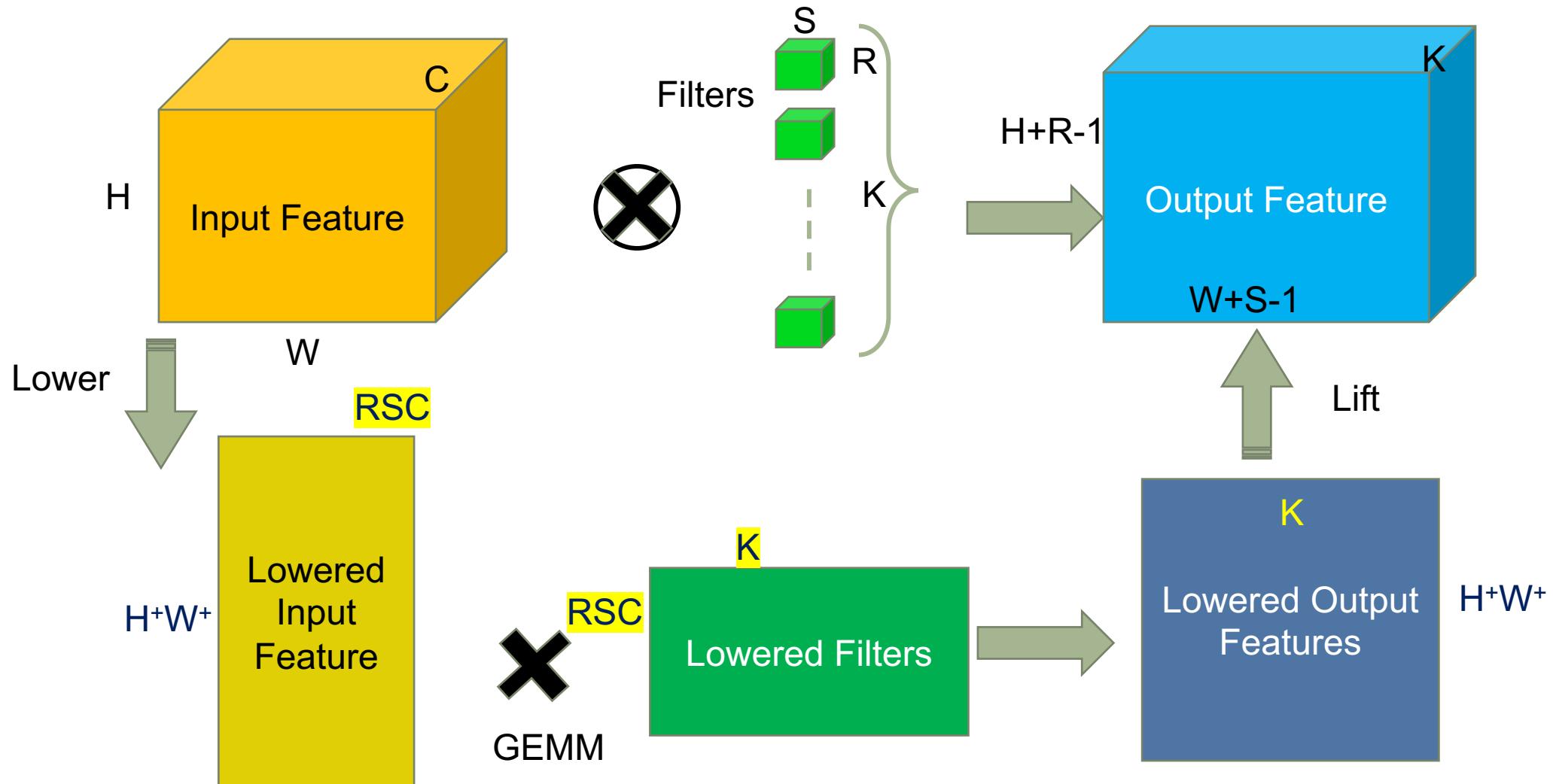
Are there other options?



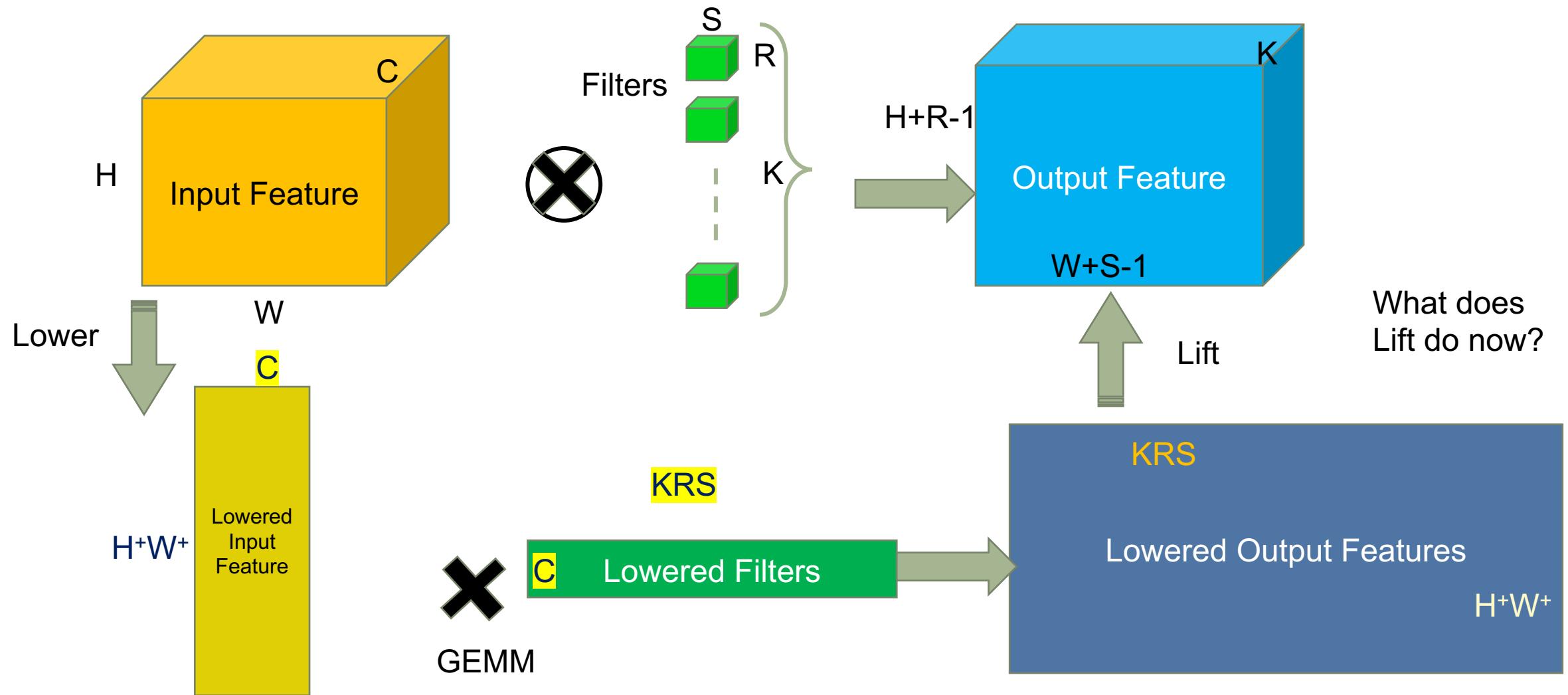
Play with the reduction dimension



Are there other options?

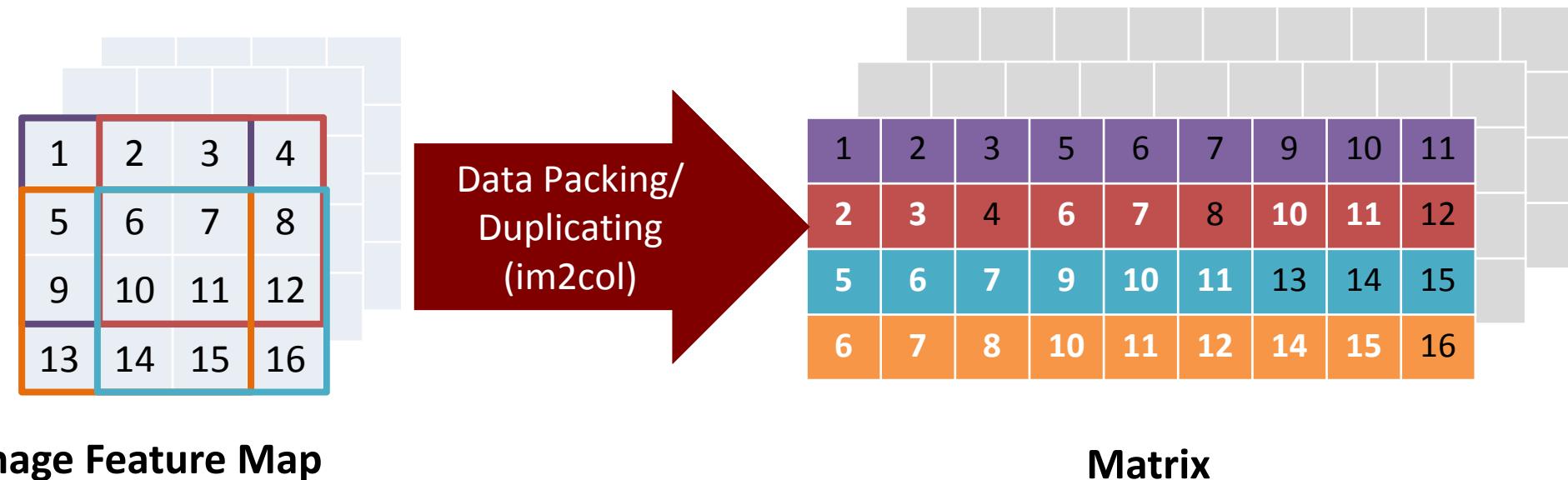


Is there another option?



Conventional approach penalizes memory to optimize performance

casts convolution to **Matrix-Matrix-Multiplication**, in order to utilize high performance MMM implementations in Basic-Linear-Algebra-Subroutines(BLAS) libraries



Direct Convolution is Better

Higher performance, zero memory overheads



	Matrix-Matrix-Multiplication	Direct Convolution
Packing	Yes <ul style="list-style-type: none">• Over 10x additional memory for image data• Performance penalty	No <ul style="list-style-type: none">• Zero memory overheads• No performance penalty
Computation Performance	Less than expected theoretic peak of GEMM	Close to system's theoretic peak

A Deeper Dive Into Convolution Loops

- Look at GEMM Blocking
 - It is Simpler
- Directly Analyze CNN Blocking

Recap of GEMM Loops

- Rules for each new character

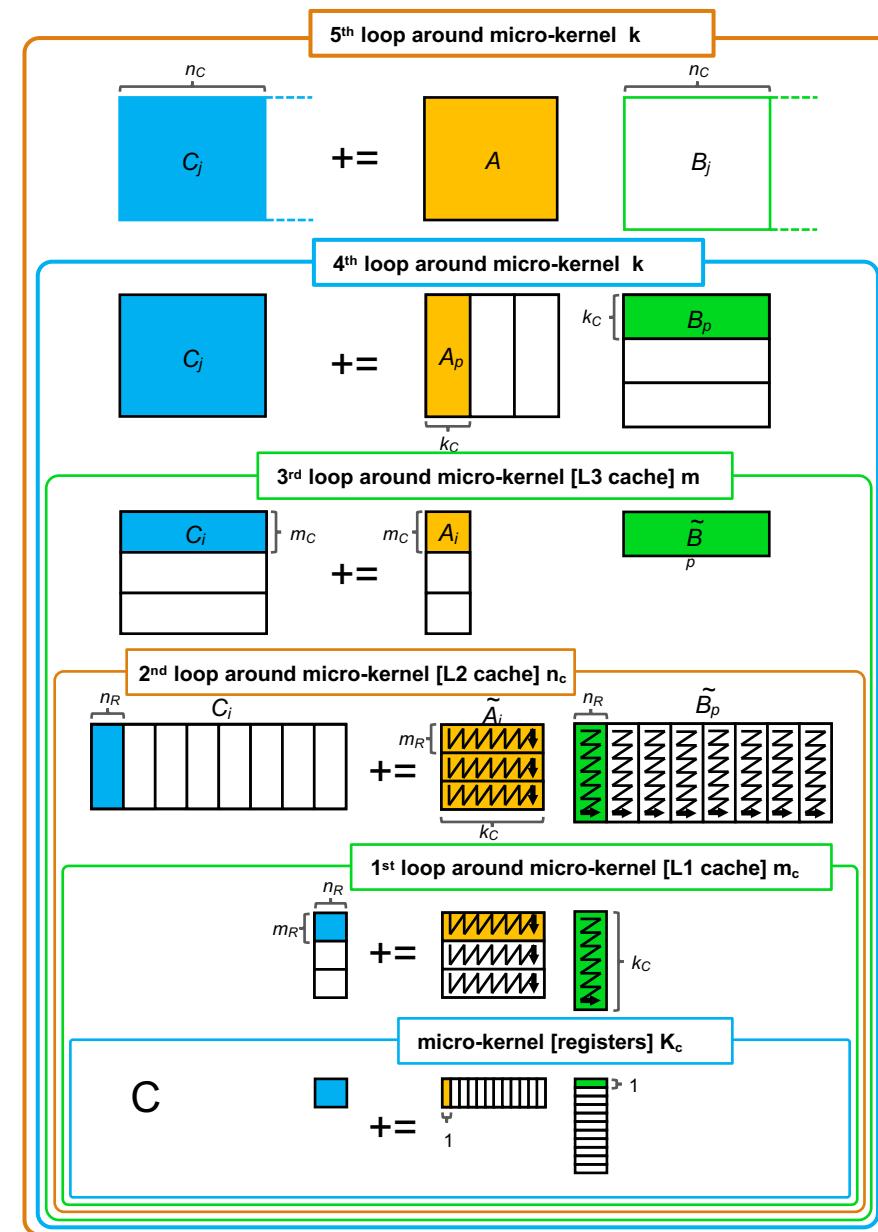
- Buffers
- Re-fetch rate

- $m_r n_r k_c m_c n_c m \ k \ n$
- $m_0 n_0 k_0 m_1 n_1 m_2 k_1 n_2$

```

for  $j_c = 0, \dots, n - 1$  in steps of  $n_c$ 
  for  $p_c = 0, \dots, k - 1$  in steps of  $k_c$ 
     $B(p_c : p_c + k_c - 1, j_c : j_c + n_c - 1) \rightarrow B_c$ 
    for  $i_c = 0, \dots, m - 1$  in steps of  $m_c$ 
       $A(i_c : i_c + m_c - 1, p_c : p_c + k_c - 1) \rightarrow A_c$ 
    for  $j_r = 0, \dots, n_c - 1$  in steps of  $n_r$ 
      for  $i_r = 0, \dots, m_c - 1$  in steps of  $m_r$ 
        for  $p_r = 0, \dots, k_c - 1$  in steps of 1
           $C_c(i_r : i_r + m_r - 1, j_r : j_r + n_r - 1)$ 
             $= A_c(i_r : i_r + m_r - 1, p_r)$ 
             $\cdot B_c(p_r, j_r : j_r + n_r - 1)$ 

```



Recap of GEMM Loops

- Rules for each new character

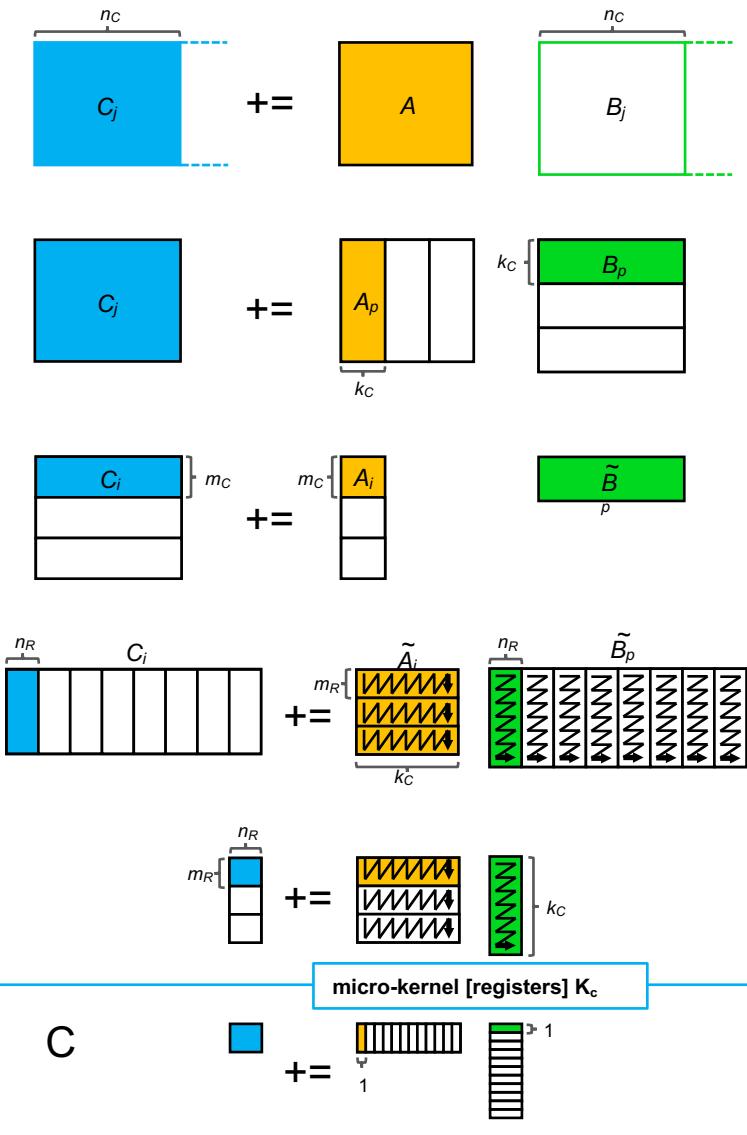
- Buffers
- Re-fetch rate

- $m_r n_r k_c$
- $m_0 n_0 k_0$

```

for  $j_c = 0, \dots, n - 1$  in steps of  $n_c$ 
  for  $p_c = 0, \dots, k - 1$  in steps of  $k_c$ 
     $B(p_c : p_c + k_c - 1, j_c : j_c + n_c - 1) \rightarrow B_c$ 
    for  $i_c = 0, \dots, m - 1$  in steps of  $m_c$ 
       $A(i_c : i_c + m_c - 1, p_c : p_c + k_c - 1) \rightarrow A_c$ 
    _____
    for  $j_r = 0, \dots, n_c - 1$  in steps of  $n_r$ 
      for  $i_r = 0, \dots, m_c - 1$  in steps of  $m_r$ 
        for  $p_r = 0, \dots, k_c - 1$  in steps of 1
           $C_c(i_r : i_r + m_r - 1, j_r : j_r + n_r - 1)$ 
             $+ = A_c(i_r : i_r + m_r - 1, p_r)$ 
             $\cdot B_c(p_r, j_r : j_r + n_r - 1)$ 

```



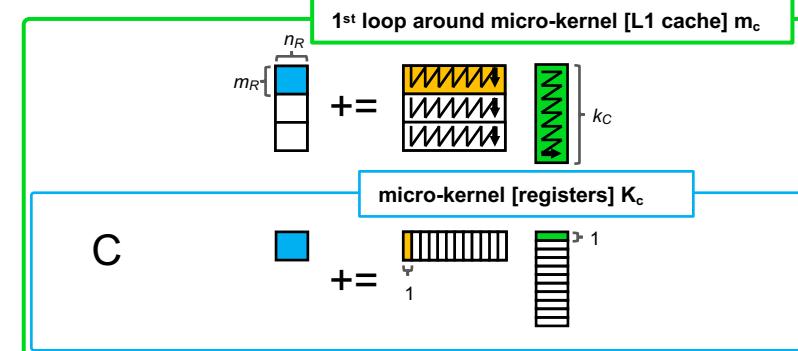
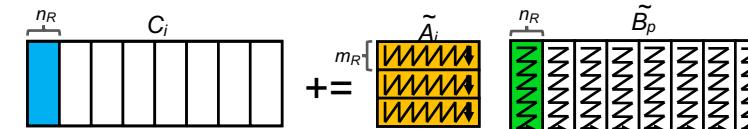
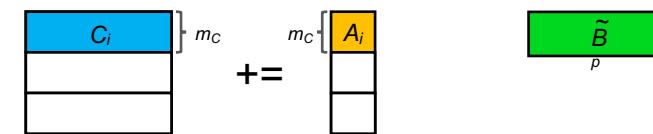
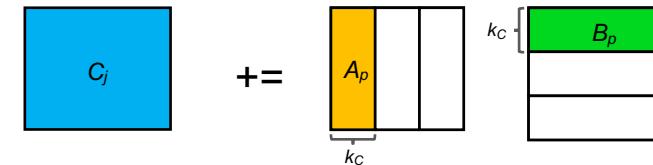
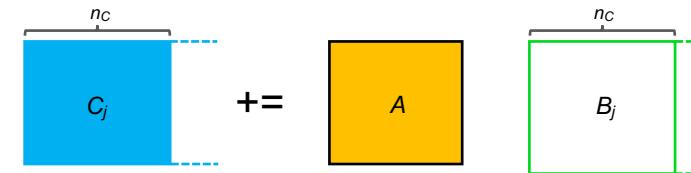
- Rules for each new character

- Buffers
- Re-fetch rate

- $m_r n_r k_c m_c$
- $m_0 n_0 k_0 m_1$

```

for  $j_c = 0, \dots, n - 1$  in steps of  $n_c$ 
  for  $p_c = 0, \dots, k - 1$  in steps of  $k_c$ 
     $B(p_c : p_c + k_c - 1, j_c : j_c + n_c - 1) \rightarrow B_c$ 
    for  $i_c = 0, \dots, m - 1$  in steps of  $m_c$ 
       $A(i_c : i_c + m_c - 1, p_c : p_c + k_c - 1) \rightarrow A_c$ 
    _____
    for  $j_r = 0, \dots, n_c - 1$  in steps of  $n_r$ 
      for  $i_r = 0, \dots, m_c - 1$  in steps of  $m_r$ 
      _____
      for  $p_r = 0, \dots, k_c - 1$  in steps of 1
         $C_c(i_r : i_r + m_r - 1, j_r : j_r + n_r - 1)$ 
         $+ = A_c(i_r : i_r + m_r - 1, p_r)$ 
         $\cdot B_c(p_r, j_r : j_r + n_r - 1)$ 
  
```



- Rules for each new character

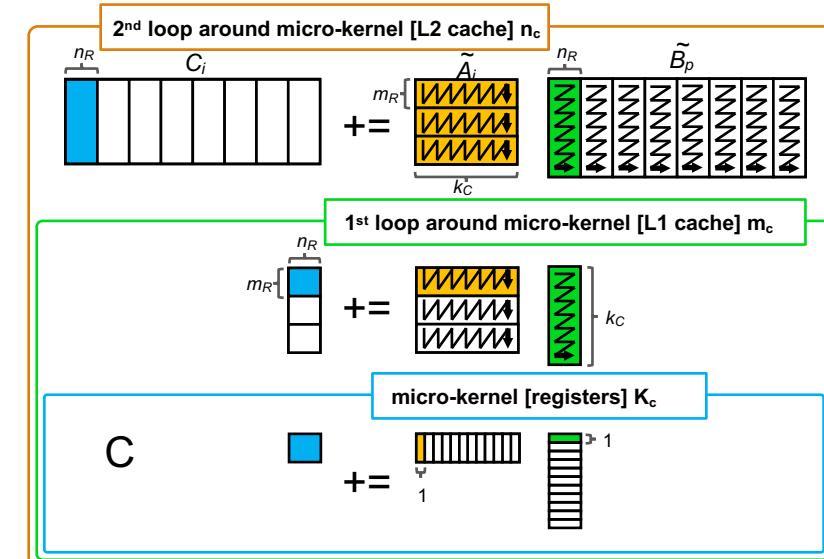
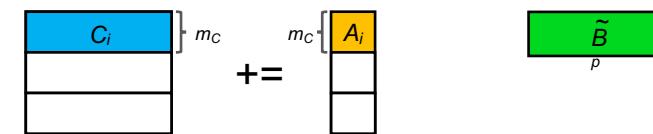
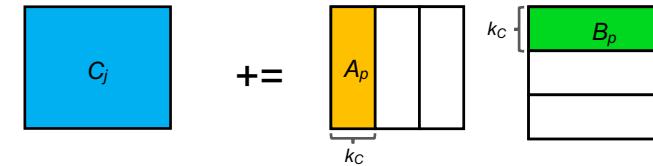
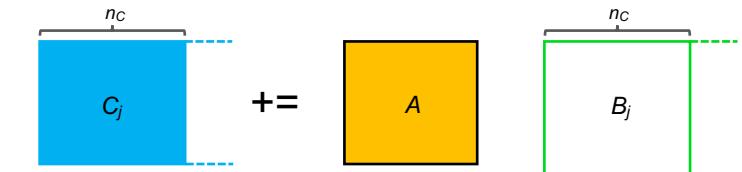
- Buffers
- Re-fetch rate

- $m_r n_r k_c m_c n_c$
- $m_0 n_0 k_0 m_1 n_1$

```

for  $j_c = 0, \dots, n - 1$  in steps of  $n_c$ 
  for  $p_c = 0, \dots, k - 1$  in steps of  $k_c$ 
     $B(p_c : p_c + k_c - 1, j_c : j_c + n_c - 1) \rightarrow B_c$ 
    for  $i_c = 0, \dots, m - 1$  in steps of  $m_c$ 
       $A(i_c : i_c + m_c - 1, p_c : p_c + k_c - 1) \rightarrow A_c$ 
      for  $j_r = 0, \dots, n_c - 1$  in steps of  $n_r$ 
        for  $i_r = 0, \dots, m_c - 1$  in steps of  $m_r$ 
          for  $p_r = 0, \dots, k_c - 1$  in steps of 1
             $C_c(i_r : i_r + m_r - 1, j_r : j_r + n_r - 1)$ 
             $+ = A_c(i_r : i_r + m_r - 1, p_r)$ 
             $\cdot B_c(p_r, j_r : j_r + n_r - 1)$ 

```



- Rules for each new character

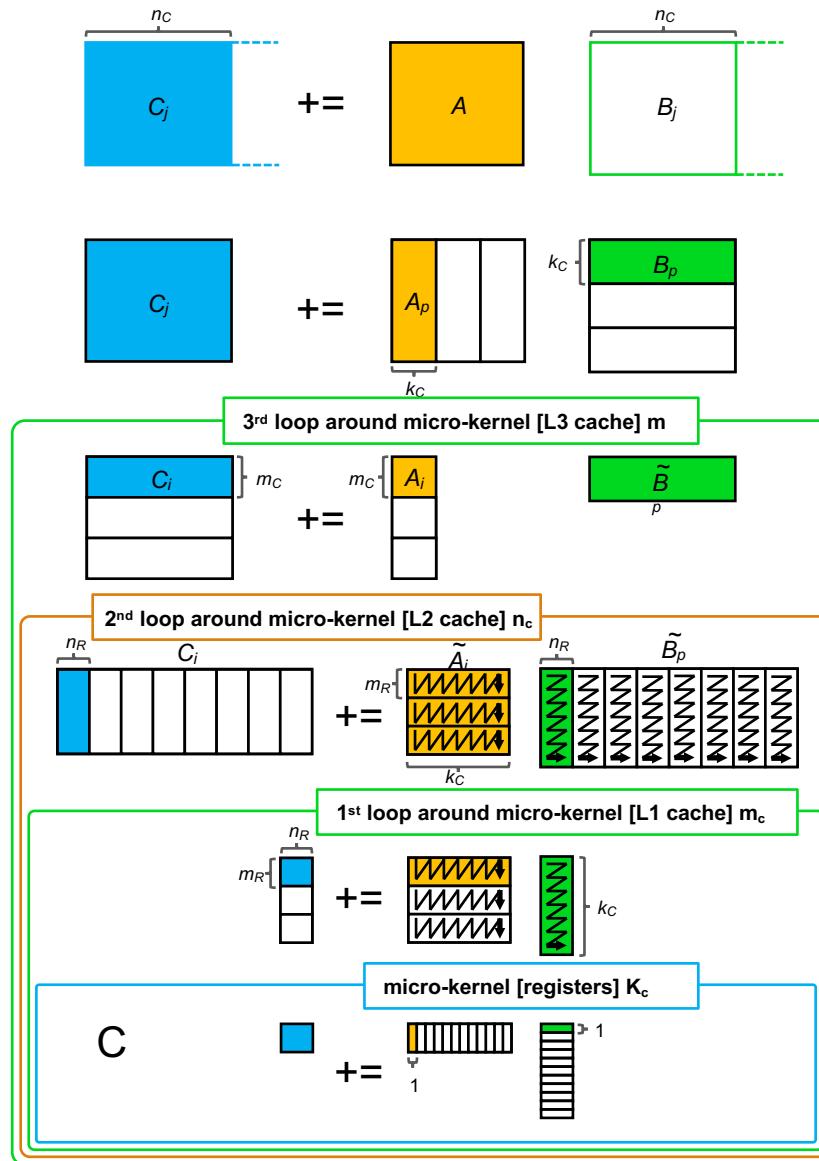
- Buffers
- Re-fetch rate

- $m_r n_r k_c m_c n_c m$
- $m_0 n_0 k_0 m_1 n_1 m_2$

```

for  $j_c = 0, \dots, n - 1$  in steps of  $n_c$ 
  for  $p_c = 0, \dots, k - 1$  in steps of  $k_c$ 
     $B(p_c : p_c + k_c - 1, j_c : j_c + n_c - 1) \rightarrow B_c$ 
    for  $i_c = 0, \dots, m - 1$  in steps of  $m_c$ 
       $A(i_c : i_c + m_c - 1, p_c : p_c + k_c - 1) \rightarrow A_c$ 
    for  $j_r = 0, \dots, n_r - 1$  in steps of  $n_r$ 
      for  $i_r = 0, \dots, m_c - 1$  in steps of  $m_r$ 
        for  $p_r = 0, \dots, k_c - 1$  in steps of 1
           $C_c(i_r : i_r + m_r - 1, j_r : j_r + n_r - 1)$ 
             $+ = A_c(i_r : i_r + m_r - 1, p_r)$ 
             $\cdot B_c(p_r, j_r : j_r + n_r - 1)$ 

```



- Rules for each new character

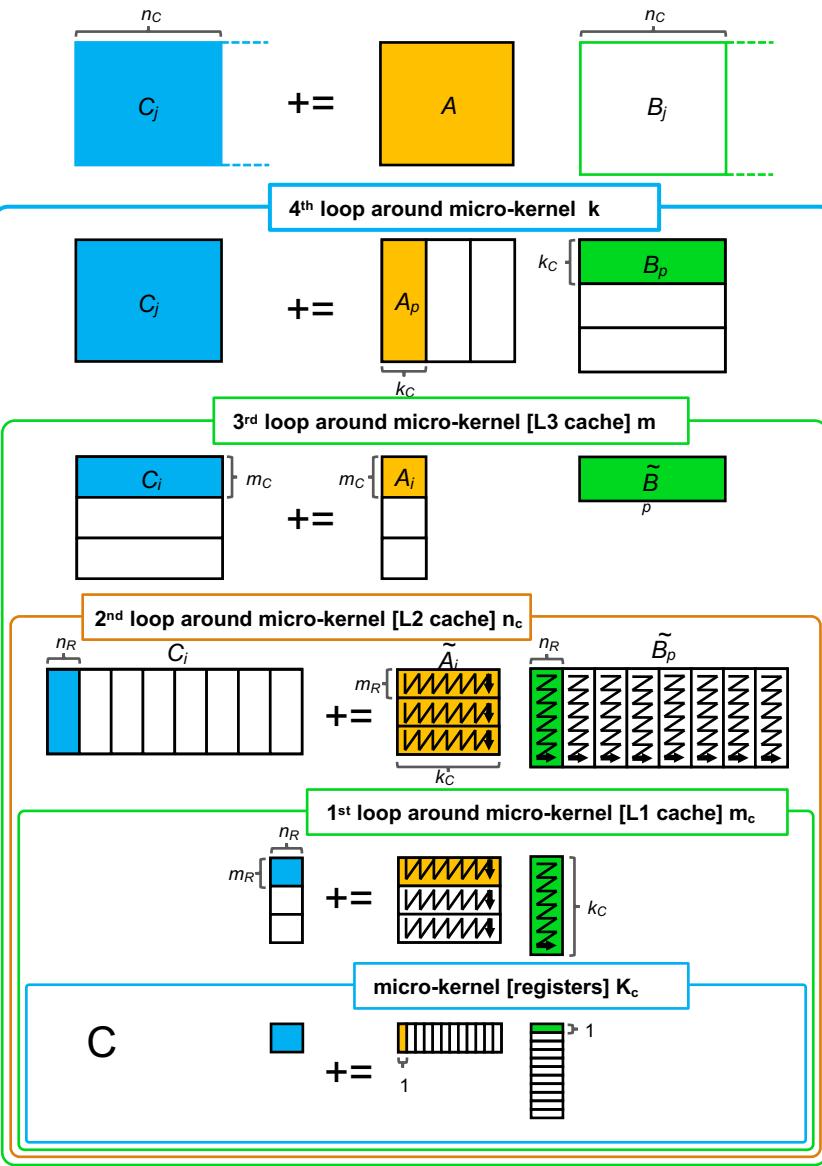
- Buffers
- Re-fetch rate

- $m_r n_r k_c m_c n_c m \ k$
- $m_0 n_0 k_0 m_1 n_1 m_2 k_1$

```

for  $j_c = 0, \dots, n - 1$  in steps of  $n_c$ 
  for  $p_c = 0, \dots, k - 1$  in steps of  $k_c$ 
     $B(p_c : p_c + k_c - 1, j_c : j_c + n_c - 1) \rightarrow B_c$ 
    for  $i_c = 0, \dots, m - 1$  in steps of  $m_c$ 
       $A(i_c : i_c + m_c - 1, p_c : p_c + k_c - 1) \rightarrow A_c$ 
    for  $j_r = 0, \dots, n_r - 1$  in steps of  $n_r$ 
      for  $i_r = 0, \dots, m_c - 1$  in steps of  $m_r$ 
        for  $p_r = 0, \dots, k_c - 1$  in steps of 1
           $C_c(i_r : i_r + m_r - 1, j_r : j_r + n_r - 1)$ 
             $+ = A_c(i_r : i_r + m_r - 1, p_r)$ 
             $\cdot B_c(p_r, j_r : j_r + n_r - 1)$ 

```



- Rules for each new character

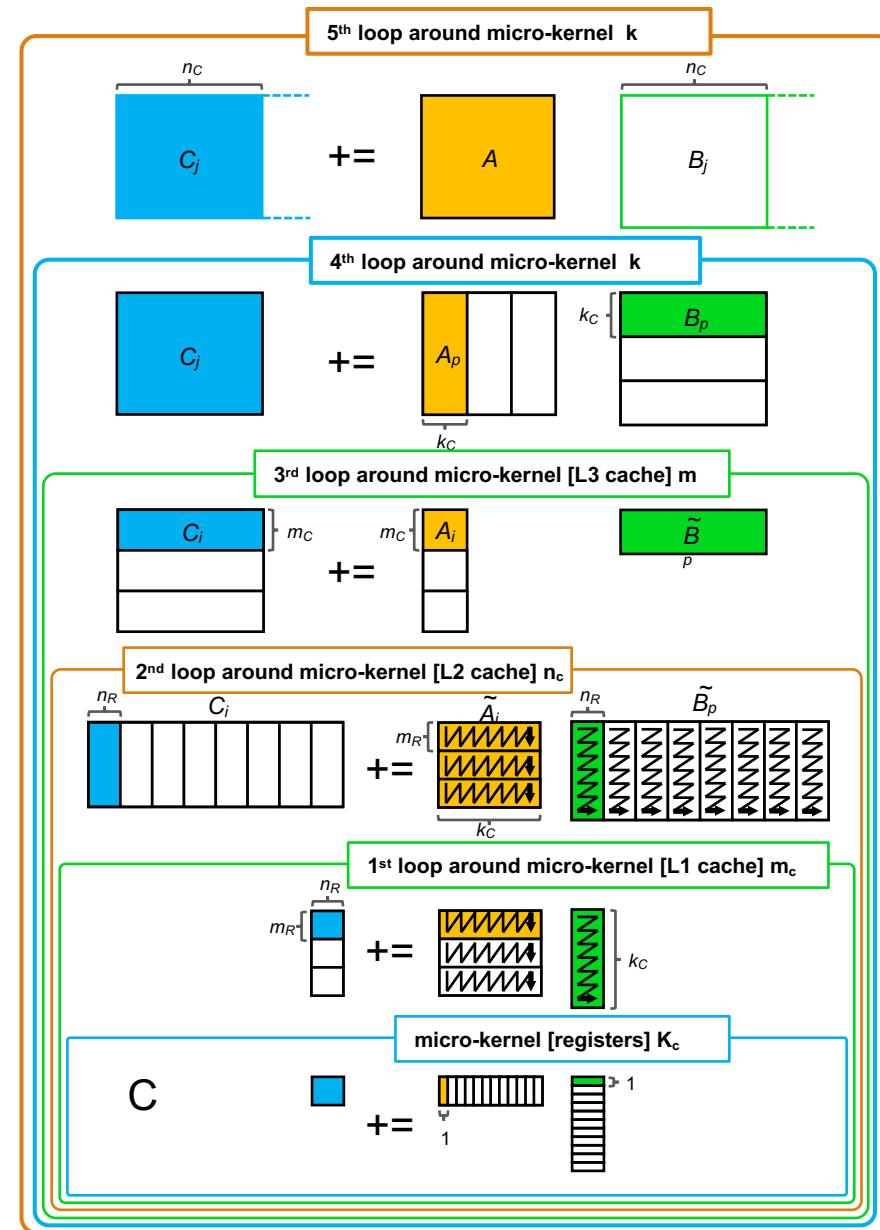
- Buffers
- Re-fetch rate

- $m_r n_r k_c m_c n_c m \ k \ n$
- $m_0 n_0 k_0 m_1 n_1 m_2 k_1 n_2$

```

for  $j_c = 0, \dots, n - 1$  in steps of  $n_c$ 
  for  $p_c = 0, \dots, k - 1$  in steps of  $k_c$ 
     $B(p_c : p_c + k_c - 1, j_c : j_c + n_c - 1) \rightarrow B_c$ 
    for  $i_c = 0, \dots, m - 1$  in steps of  $m_c$ 
       $A(i_c : i_c + m_c - 1, p_c : p_c + k_c - 1) \rightarrow A_c$ 
    for  $j_r = 0, \dots, n_r - 1$  in steps of  $n_r$ 
      for  $i_r = 0, \dots, m_c - 1$  in steps of  $m_r$ 
        for  $p_r = 0, \dots, k_c - 1$  in steps of 1
           $C_c(i_r : i_r + m_r - 1, j_r : j_r + n_r - 1)$ 
           $+ = A_c(i_r : i_r + m_r - 1, p_r)$ 
           $\cdot B_c(p_r, j_r : j_r + n_r - 1)$ 

```



GEMM Blocking

- Loop representation string
 - Rules for each new character
 - Buffers
 - Re-fetch rate

- $m_r n_r k_c m_c n_c m \ k \ n$
- $m_0 n_0 k_0 m_1 n_1 m_2 k_1 n_2$

```

for  $j_c = 0, \dots, n - 1$  in steps of  $n_c$ 
  for  $p_c = 0, \dots, k - 1$  in steps of  $k_c$ 
     $B(p_c : p_c + k_c - 1, j_c : j_c + n_c - 1) \rightarrow B_c$  // Pack into  $B_c$ 
    for  $i_c = 0, \dots, m - 1$  in steps of  $m_c$ 
       $A(i_c : i_c + m_c - 1, p_c : p_c + k_c - 1) \rightarrow A_c$  // Pack into  $A_c$ 
    for  $j_r = 0, \dots, n_c - 1$  in steps of  $n_r$  // Macro-kernel
      for  $i_r = 0, \dots, m_c - 1$  in steps of  $m_r$ 
        for  $p_r = 0, \dots, k_c - 1$  in steps of 1 // Micro-kernel
           $C_c(i_r : i_r + m_r - 1, j_r : j_r + n_r - 1)$ 
           $+ = A_c(i_r : i_r + m_r - 1, p_r)$ 
           $\cdot B_c(p_r, j_r : j_r + n_r - 1)$ 

```

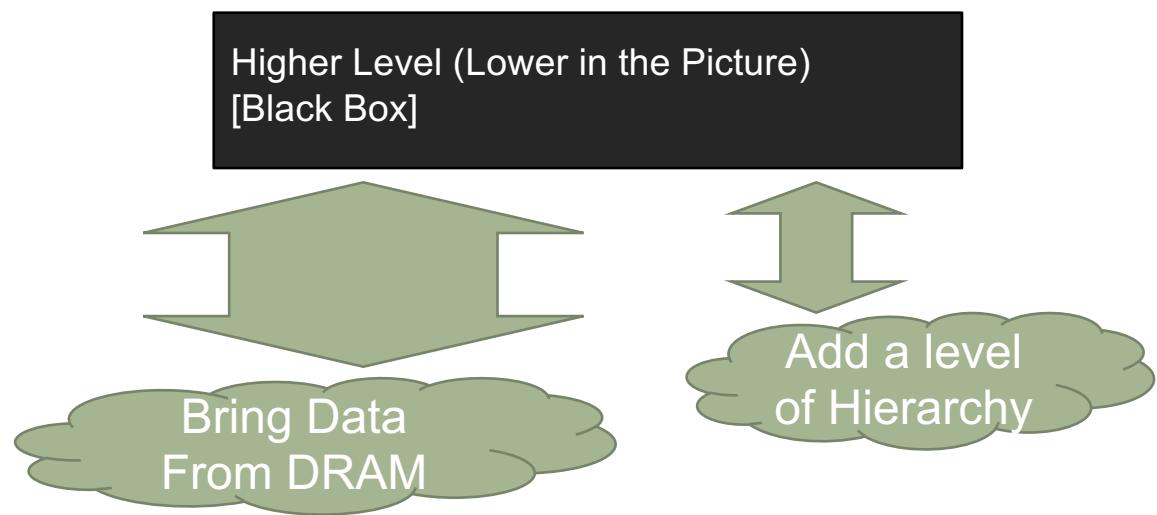
Recursive Formulation

- Fixed order as follows:

_____ $N_{i-1}M_{i-1}K_{i-1}$ 

- From left to right observing
 - N: New A

Buffer Type	Buffer Size	Buffer Refetch Rate



Recursive Formulation

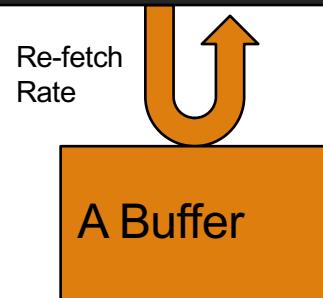
- Fixed order as follows:

$N_{i-1} M_{i-1} K_{i-1} N_i$ 

- From left to right observing
 - N : New A

Buffer Type	Buffer Size	Buffer Refetch Rate
A_i	$M_{i-1} K_{i-1}$	N_i / N_{i-1}

Higher Level (Lower in the Picture)
[Black Box]



Recursive Formulation

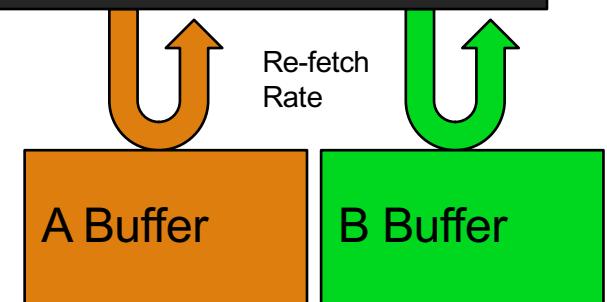
- Fixed order as follows:

$N_{i-1} M_{i-1} K_{i-1} N_i M_i$

Buffer Type	Buffer Size	Buffer Refetch Rate
A_i	$M_{i-1} K_{i-1}$	N_i / N_{i-1}
B_i	$N_{i-1} K_{i-1}$	M_i / M_{i-1}

- From left to right observing
 - N: New A
 - M: New B

Higher Level (Lower in the Picture)
[Black Box]



Recursive Formulation

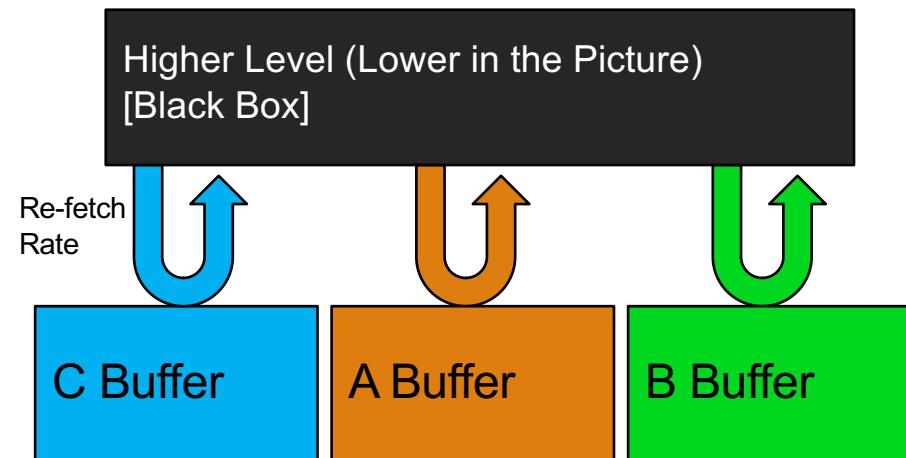
- Fixed order as follows:

_____ $N_{i-1} M_{i-1} K_{i-1}$ $N_i M_i K_i$ 

Buffer Type	Buffer Size	Buffer Refetch Rate
A_i	$M_{i-1} K_{i-1}$	N_i / N_{i-1}
B_i	$N_{i-1} K_{i-1}$	M_i / M_{i-1}
C_i	$M_{i-1} N_{i-1}$	$2K_i / K_{i-1}$

- From left to right observing

- N: New A
- M: New B
- K: New C

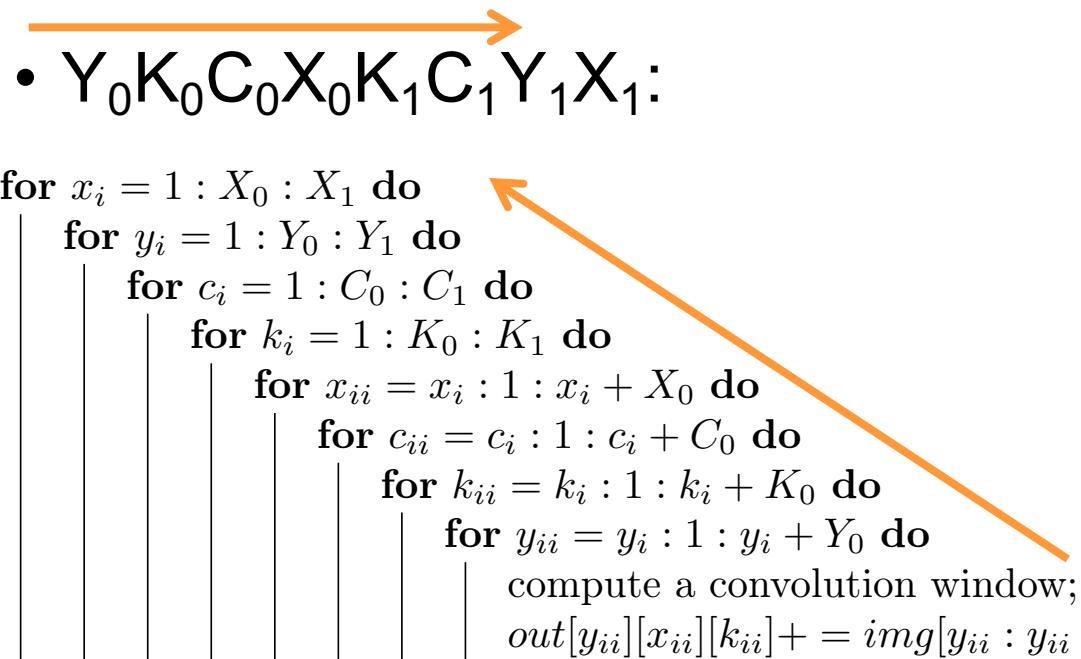


Formal Derivation

- Loop representation string
- Rules for each new character
 - Buffers
 - Re-fetch rate

• $Y_0 K_0 C_0 X_0 K_1 C_1 Y_1 X_1 :$

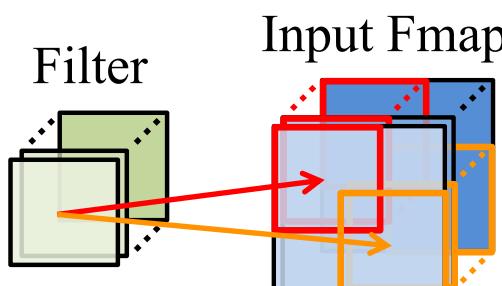
```
for  $x_i = 1 : X_0 : X_1$  do
  for  $y_i = 1 : Y_0 : Y_1$  do
    for  $c_i = 1 : C_0 : C_1$  do
      for  $k_i = 1 : K_0 : K_1$  do
        for  $x_{ii} = x_i : 1 : x_i + X_0$  do
          for  $c_{ii} = c_i : 1 : c_i + C_0$  do
            for  $k_{ii} = k_i : 1 : k_i + K_0$  do
              for  $y_{ii} = y_i : 1 : y_i + Y_0$  do
                compute a convolution window;
                 $out[y_{ii}][x_{ii}][k_{ii}] += img[y_{ii} : y_{ii} + F_h][x_{ii} : x_{ii} + F_w][c_{ii}] \times kernel[c_{ii}][k_{ii}]$ 
```



Types of Data Reuse in DNN

Convolutional Reuse

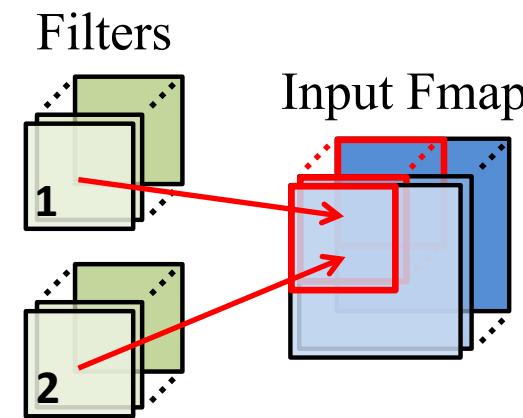
CONV layers only
(sliding window)



Reuse: **Activations**
Filter weights

Fmap Reuse

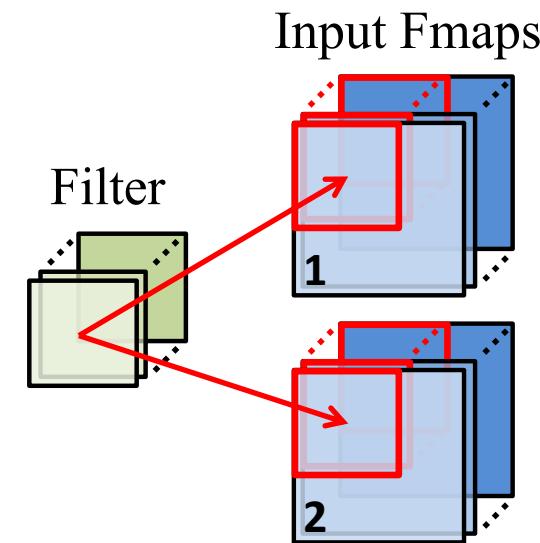
CONV and FC layers



Reuse: **Activations**

Filter Reuse

CONV and FC layers
(batch size > 1)



Reuse: **Filter weights**

Dataflow Taxonomy Summary

- **Weight Stationary**
 - minimize movement of **weights**
- **Output Stationary**
 - minimize movement of **partial sums**
- **Input Stationary**
 - minimize movement of **input activations**
- **No Local Reuse**
 - don't use any local PE storage. Maximize **global buffer size**

Dataflow Taxonomy Summary

- **Weight Stationary**
 - minimize movement of **weights**
- **Output Stationary**
 - minimize movement of **partial sums**
- **Input Stationary**
 - minimize movement of **input activations**
- **No Local Reuse**
 - don't use any local PE storage. Maximize **global buffer size**

Does this define the space of all possible dataflows?

Dataflow Taxonomy Summary

- **Weight Stationary**
 - minimize movement of **weights**
- **Output Stationary**
 - minimize movement of **partial sums**
- **Input Stationary**
 - minimize movement of **input activations**
- **No Local Reuse**
 - don't use any local PE storage. Maximize **global buffer size**

Does this define the space of all possible dataflows?

Not at all!

Can we have multiple **stationariness** in an architecture?



Formal Derivation

- Loop representation string
- Rules for each new character
 - Buffers
 - Re-fetch rate

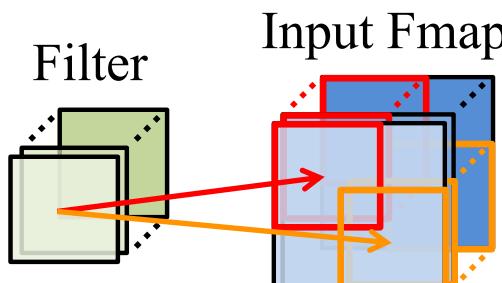
• $Y_0 K_0 C_0 X_0 K_1 C_1 Y_1 X_1 :$

```
for  $x_i = 1 : X_0 : X_1$  do
  for  $y_i = 1 : Y_0 : Y_1$  do
    for  $c_i = 1 : C_0 : C_1$  do
      for  $k_i = 1 : K_0 : K_1$  do
        for  $x_{ii} = x_i : 1 : x_i + X_0$  do
          for  $c_{ii} = c_i : 1 : c_i + C_0$  do
            for  $k_{ii} = k_i : 1 : k_i + K_0$  do
              for  $y_{ii} = y_i : 1 : y_i + Y_0$  do
                compute a convolution window;
                 $out[y_{ii}][x_{ii}][k_{ii}] += img[y_{ii} : y_{ii} + F_h][x_{ii} : x_{ii} + F_w][c_{ii}] \times kernel[c_{ii}][k_{ii}]$ 
```

Types of Data Reuse in DNN

Convolutional Reuse

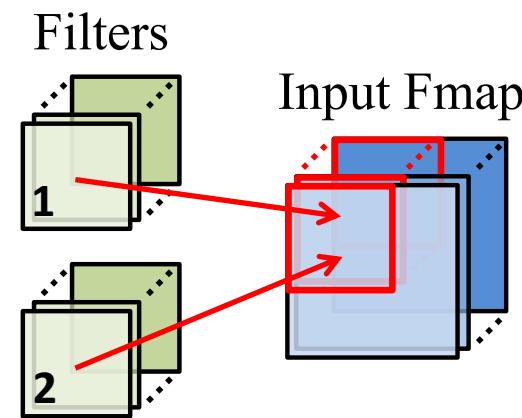
CONV layers only
(sliding window)



Reuse: **Activations**
Filter weights

Fmap Reuse

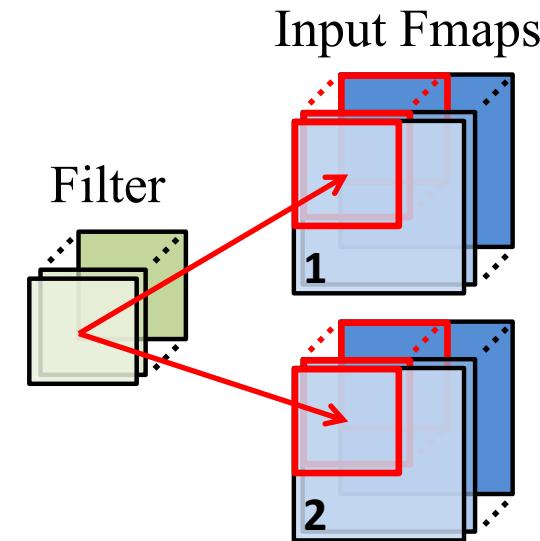
CONV and FC layers



Reuse: **Activations**

Filter Reuse

CONV and FC layers
(batch size > 1)



Reuse: **Filter weights**

Recursive Formulation

- Fixed order as follows:

_____ $Y_{i-1} X_{i-1} C_{i-1} K_{i-1} \dots \dots$



- From left to right observing

Higher Level
[Black Box]

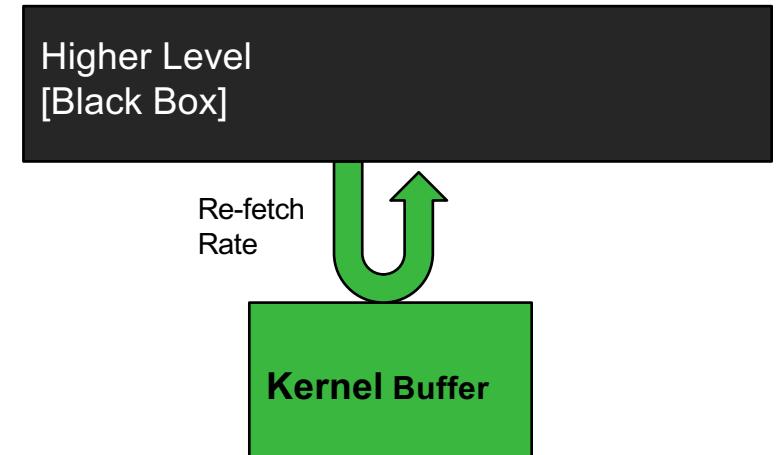
Recursive Formulation

Buffer Name	Buffer Size	Buffer Refetch Rate
KB_i	$C_{i-1}K_{i-1}F_hF_w$	$(X_iY_i)/(X_{i-1}Y_{i-1})$

- Fixed order as follows:

_____ $Y_{i-1}X_{i-1}C_{i-1}K_{i-1}Y_iX_i\dots\dots$ 

- From left to right observing
 - X or Y: New KB



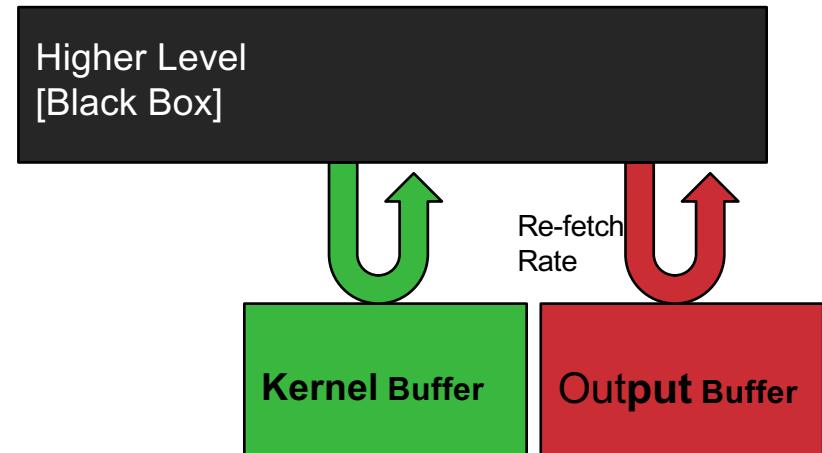
Recursive Formulation

Buffer Name	Buffer Size	Buffer Refetch Rate
KB_i OB_i	$C_{i-1}K_{i-1}F_hF_w$ $Y_iX_iK_{i-1}$	$(X_iY_i)/(X_{i-1}Y_{i-1})$ $2C_i/C_{i-1}$

- Fixed order as follows:

_____ $Y_{i-1}X_{i-1}C_{i-1}K_{i-1}Y_iX_iC_i\dots\dots$

- From left to right observing
 - X or Y: New KB
 - C: New OB



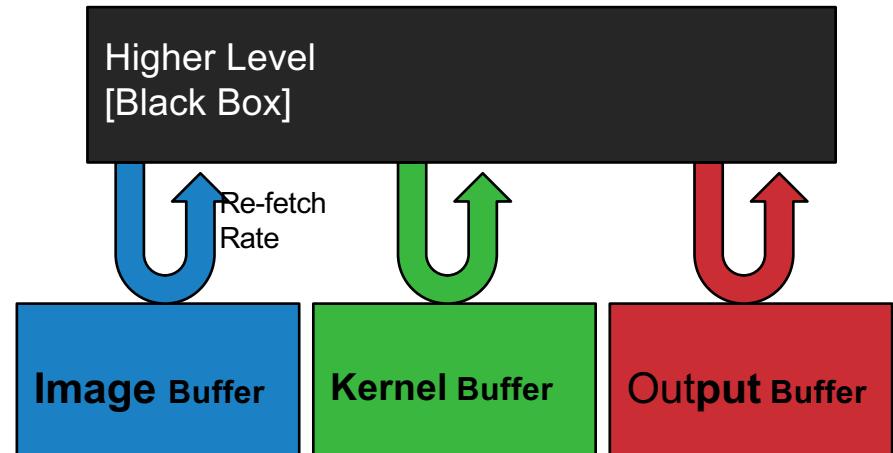
Recursive Formulation

Buffer Name	Buffer Size	Buffer Refetch Rate
KB_i	$C_{i-1}K_{i-1}F_hF_w$	$(X_iY_i)/(X_{i-1}Y_{i-1})$
OB_i	$Y_iX_iK_{i-1}$	$2C_i/C_{i-1}$
IB_i	$(Y_i + F_h - 1)(X_i + F_w - 1)C_i$	$(K_i(Y_i + F_h - 1)(X_i + F_w - 1))/(K_{i-1}Y_iX_i)$

- Fixed order as follows:

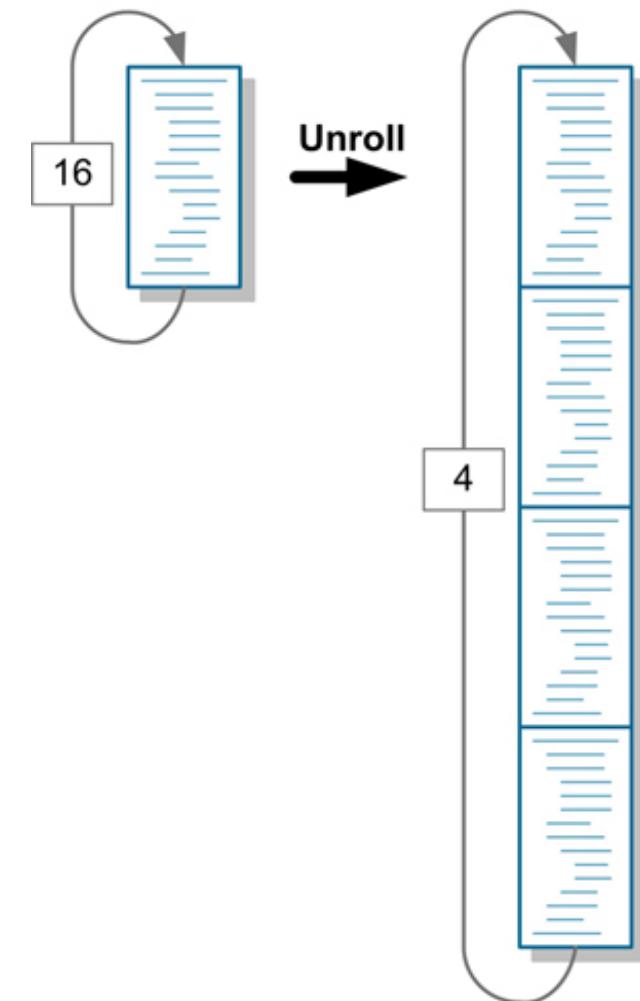
_____ $Y_{i-1}X_{i-1}C_{i-1}K_{i-1}Y_iX_iC_iK_i\dots\dots$

- From left to right observing
 - X or Y: New KB
 - C: New OB
 - K: New IB



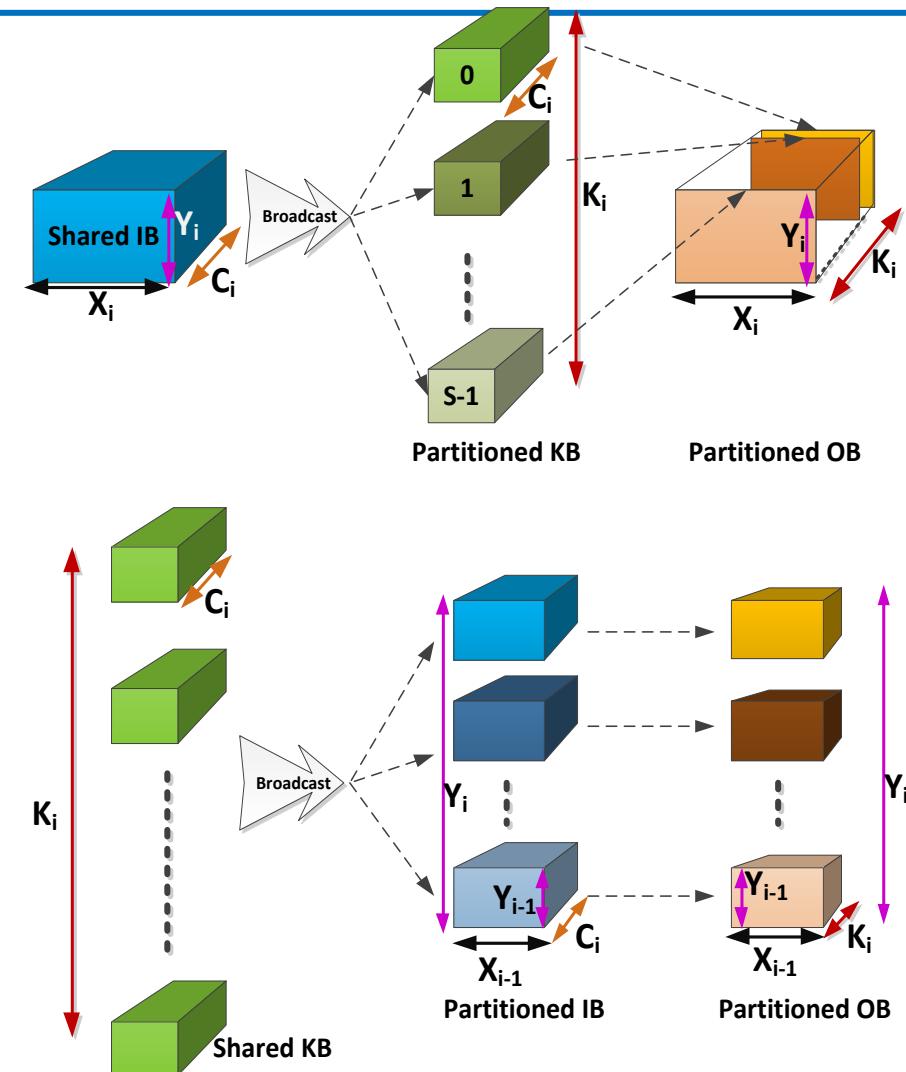
How Does Parallelism Work?

- Unroll Loops in Space
- Replicate Compute Units
- Perform Loops in parallel
 - SIMD
 - Multiple Cores
 - Multiple CPUs/GPUs



Multi-Core Parallelism and Partitioning

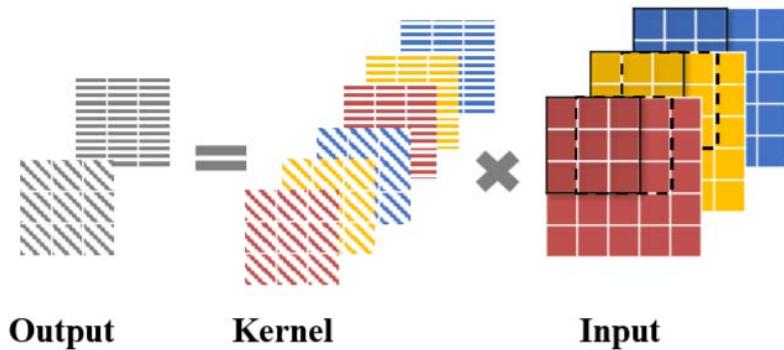
- Across channels (C)
 - Reduction
 - Synchronization
- Across kernels (K)
 - Synchronization
- Across image (X or Y)
 - Independent
- Broadcast operation
- Partitioned memory
- Nested parallelism



High Performance Direct Convolution

Map application to system architecture to fully exploit system capacity

High-dimensional data tensor



Multi-level nested loops and various data access patterns

```
for  $i = 1$  to  $C_i$  do
    for  $j = 1$  to  $C_o$  do
        for  $k = 1$  to  $W_o$  do
            for  $\ell = 1$  to  $H_o$  do
                for  $m = 1$  to  $W_f$  do
                    for  $n = 1$  to  $H_f$  do
                         $\mathcal{O}_{j,k,\ell} += \mathcal{I}_{i,k \times s+m, \ell \times s+n} \times \mathcal{F}_{i,j,m,n}$ 
```

Direct Convolution

Comp Units
FP arithmetic units,
SIMD engine,
multicore, etc

L1 cache

L2 cache

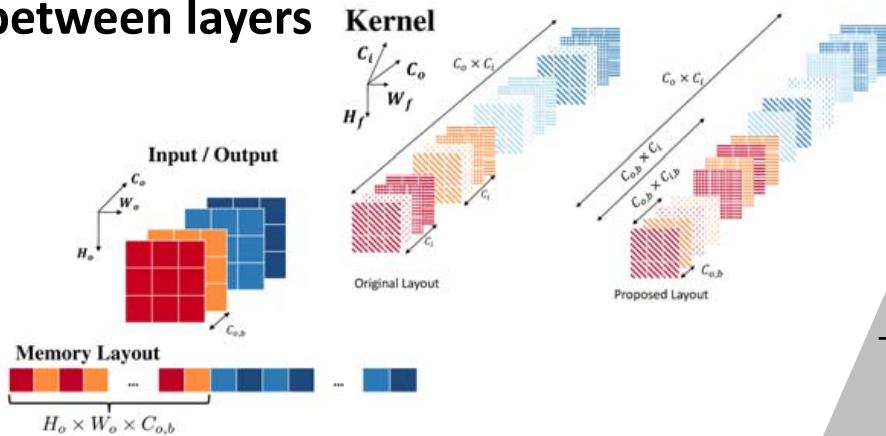
L3 cache

System Architecture

High Performance Direct Convolution

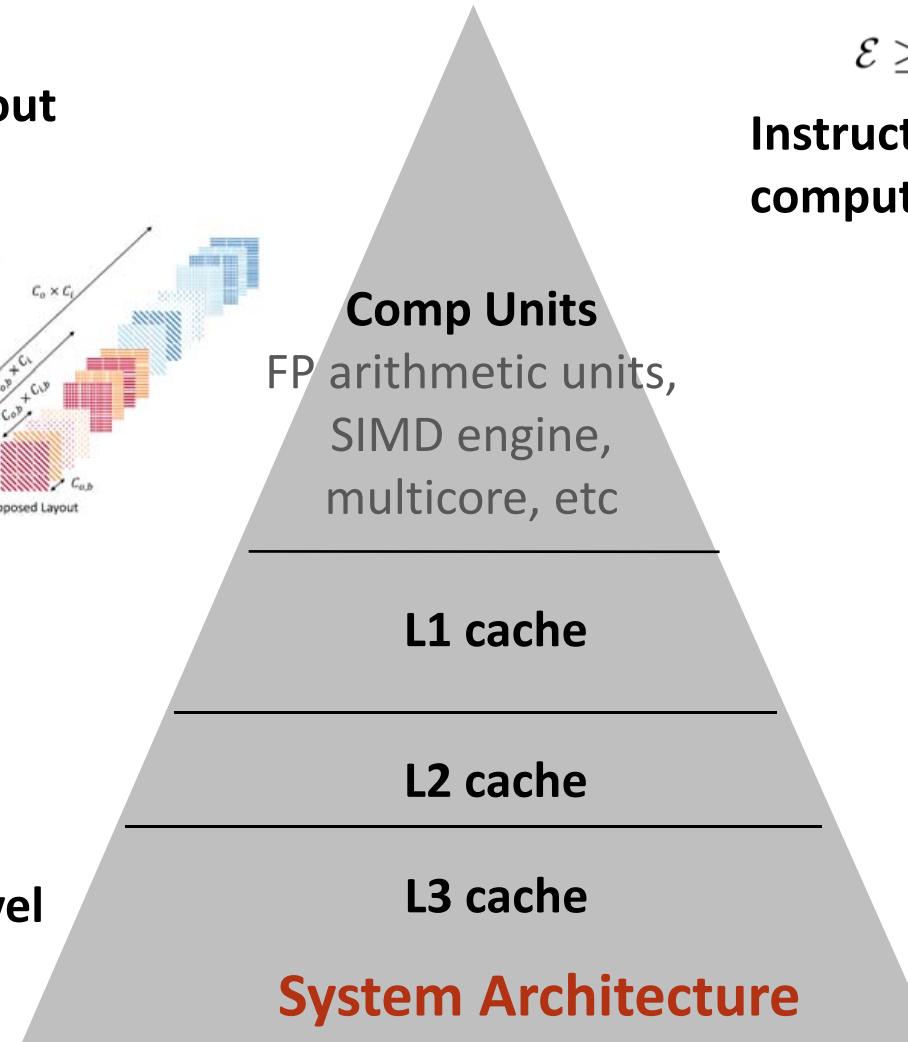
Analytical model to achieve system theoretical peak performance

Convolution-friendly data layout
that avoids format reshaping
between layers



```
for  $j' = 1$  to  $C_o/C_{o,b}$  in Parallel do  
  for  $i' = 1$  to  $C_i/C_{i,b}$  do
```

Parallelism to exploit core-level
concurrency



$$\mathcal{E} \geq N_{\text{vec}} N_{\text{fma}} L_{\text{fma}},$$

Instruction scheduling to saturate the
computation pipeline

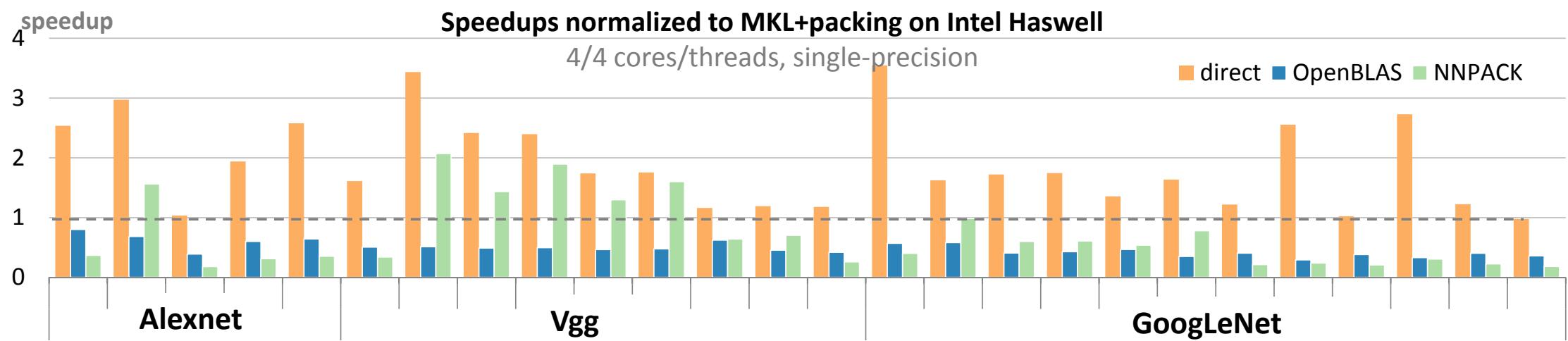
Algorithm 2 Reorder Convolution Algorithm

```
Input: Input  $\mathcal{I}$ , Kernel Weights  $\mathcal{F}$ , stride  $s$ ;  
Output: Output  $\mathcal{O}$   
for  $\ell = 1$  to  $H_o$  do  
  for  $n = 1$  to  $H_f$  do  
    for  $m = 1$  to  $W_f$  do  
      for  $i = 1$  to  $C_i$  do  
        for  $k = 1$  to  $W_o$  do  
          for  $j = 1$  to  $C_o$  do  
             $\mathcal{O}_{j,k,\ell} += \mathcal{I}_{i,k \times s + m, \ell \times s + n} \times \mathcal{F}_{i,j,m,n}$ 
```

Computation reordering that
maximizes data reuse

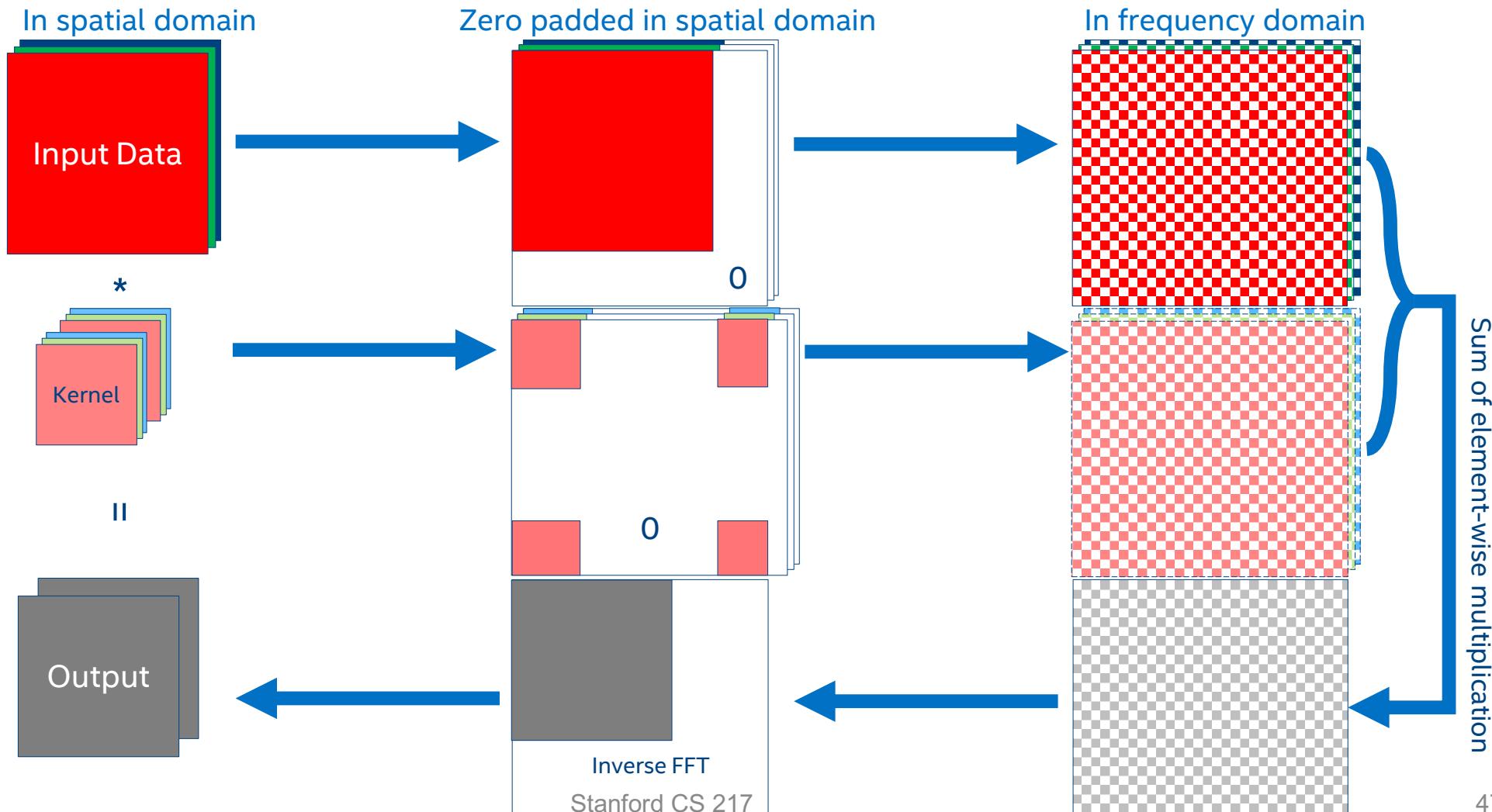
Loop tiling and
blocking to utilize the
memory hierarchies

Better performance



FFT based convolution

Convert input into Fourier Domain, apply element-wise multiplication to reduce complexity:
 $O(N^2K^2) \rightarrow O(N^2\log_2 N)$, where N is data size, K is kernel size



Winograd Method

- Fast algorithm for computing short convolutions
1. Polynomials
 - a) Chinese remainder theorem
 - b) Lagrange Interpolation
 2. Toom-Cook Algorithm
 3. Winograd method for FIR filters
 4. Gray and Levin paper on CNNs for Winograd

https://artofproblemsolving.com/community/c1157h990758_the_chinese_remainder_theorem_and_lagrange_interpolation

Interpolation

- To find the value of y for an x between different x - values x_0, x_1, \dots, x_n is called problem of interpolation.
- To find the value of y for an x which falls outside the range of x ($x < x_0$ or $x > x_n$) is called the problem of extrapolation.
- Theorem by Weierstrass in 1885, “Every continuous function in an interval (a,b) can be represented in that interval to any desired accuracy by a polynomial. ”
- Let us assign polynomial P_n of degree n (or less) that assumes the given data values
- $P_n(x_0) = y_0, P_n(x_1) = y_1, \dots, P_n(x_n) = y_n$
This polynomial P_n is called interpolation polynomial.
- x_0, x_1, \dots, x_n is called the nodes (tabular points, pivotal points or arguments).

Lagrange Interpolation

- Lagrange's interpolation formula with unequal intervals:
- Let $y = f(x)$ be continuous and differentiable in the interval (a, b) .
- Given the set of $n+1$ values $(x_0, y_0), (x_1, y_1), \dots, (x_n, y_n)$ of x and y , where the values of x need not necessarily be equally spaced.
- It is required to find $P_n(x)$, a polynomial of degree n such that y and $P_n(x)$ agree at the tabulated points.

Lagrange Interpolation

- This polynomial is given by the following formula:

$$\begin{aligned}y = f(x) \approx P_n(x) &= \frac{(x - x_1)(x - x_2) \dots (x - x_n)}{(x_0 - x_1)(x_0 - x_2) \dots (x_0 - x_n)} y_0 \\&+ \frac{(x - x_0)(x - x_2) \dots (x - x_n)}{(x_1 - x_0)(x_1 - x_2) \dots (x_1 - x_n)} y_1 + \dots \\&+ \frac{(x - x_0)(x - x_1) \dots (x - x_{n-1})}{(x_n - x_0)(x_n - x_1) \dots (x_n - x_{n-1})} y_n\end{aligned}$$

-
- In the Lagrange's interpolation formula y is treated as dependent variable and expressed as function of independent variable x .
 - Instead if x is treated as dependent variable and expressed as the function of independent variable y , then Lagrange's interpolation formula becomes

$$\begin{aligned}x = g(y) \approx P_n(y) &= \frac{(y - y_1)(y - y_2) \dots (y - y_n)}{(y_0 - y_1)(y_0 - y_2) \dots (y_0 - y_n)} x_0 \\&+ \frac{(y - y_0)(y - y_2) \dots (y - y_n)}{(y_1 - y_0)(y_1 - y_2) \dots (y_1 - y_n)} x_1 + \dots \\&+ \frac{(y - y_0)(y - y_1) \dots (y - y_{n-1})}{(y_n - y_0)(y_n - y_1) \dots (y_n - y_{n-1})} x_n\end{aligned}$$

- This relation is referred as Lagrange's inverse interpolation formula.

Chinese Remainder Theorem

- It is possible to uniquely determine a nonnegative integer given only its moduli with respect to each of several integers, provided that the integer is known to be smaller than the product of the moduli.

CRT Example

- Choose the moduli $m_0=3$, $m_1=4$, and $m_2=5$, and let $M= m_0m_1m_2= 60$.
- Given the integer c satisfying $0 \leq c < 60$, let $c_i = Rm_i [c]$.
- The Chinese remainder theorem says that there is a one-to-one map between the sixty values that c is allowed to take on and the sixty values that the vector of residues (c_0, c_1, c_2) can take on.
- Suppose that $c_0 = 2$, $c_1 = 1$, and $c_2 = 2$. The Unique Solution is ..

$$c \in \{2, 5, 8, 11, 14, 17, 20, 23, 26, 29, \dots\},$$

$$c \in \{1, 5, 9, 13, 17, 21, 25, 29, 33, \dots\},$$

$$c \in \{2, 7, 12, 17, 22, 27, 32, 37, \dots\}.$$

Relation between CRT and Lagrange

- Chinese Remaining Theorem for polynomials is related to Lagrange interpolation
- Lagrange Interpolation is used in Toom-Cook Fast Convolution
- CRT is used in Winograd Fast Convolution
- Read More about the relation:

Saints and Scoundrels and Two Theorems That Are Really the Same by Ezra Brown