



香港中文大學
The Chinese University of Hong Kong

CMSC5743

L02: CNN Accurate Speedup I

Bei Yu

(Latest update: September 28, 2020)

Fall 2020



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/>

Overview



Convolution 101

GEMM

Sparse Convolution

Direct Convolution

Further Discussions

Overview



Convolution 101

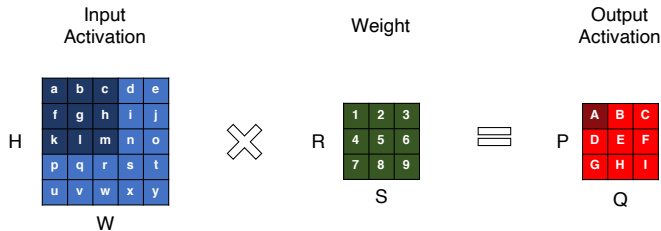
GEMM

Sparse Convolution

Direct Convolution

Further Discussions

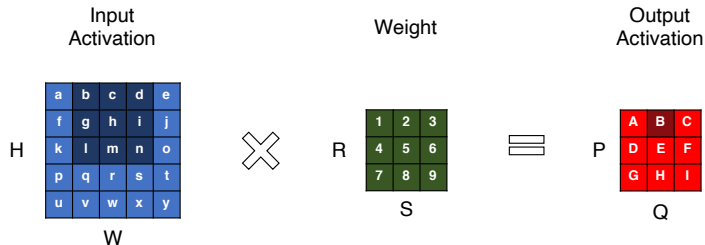
2D-Convolution



H: Height of Input Activation
W: Width of Input Activation
R: Height of Weight
S: Width of Weight
P: Height of Output Activation
Q: Width of Output Activation

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

2D-Convolution



H: Height of Input Activation
W: Width of Input Activation
R: Height of Weight
S: Width of Weight
P: Height of Output Activation
Q: Width of Output Activation
stride: # of rows/columns traversed per step

2D-Convolution



Input
Activation

	a	b	c	d	e
	f	g	h	i	j
H	k	l	m	n	o
	p	q	r	s	t
	u	v	w	x	y
					W



Weight

	1	2	3
	4	5	6
R	7	8	9
			S

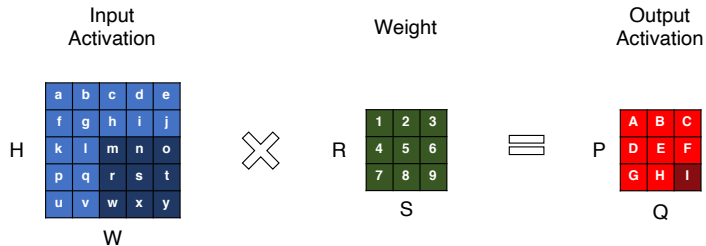


Output
Activation

	A	B	C
	D	E	F
P	G	H	I
			Q

H: Height of Input Activation
W: Width of Input Activation
R: Height of Weight
S: Width of Weight
P: Height of Output Activation
Q: Width of Output Activation
stride: # of rows/columns traversed per step

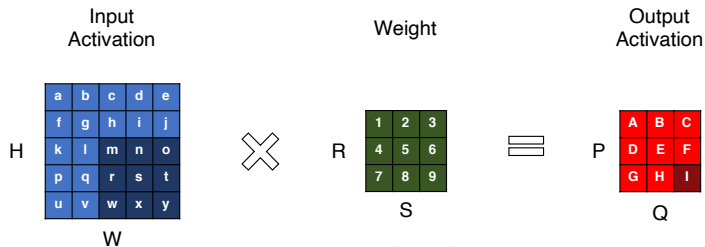
2D-Convolution



H: Height of Input Activation
W: Width of Input Activation
R: Height of Weight
S: Width of Weight
P: Height of Output Activation
Q: Width of Output Activation
stride: # of rows/columns traversed per step

$$l = m * 1 + n * 2 + o * 3 \\ + r * 4 + s * 5 + t * 6 \\ + w * 7 + x * 8 + y * 9$$

2D-Convolution

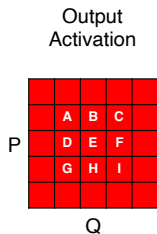
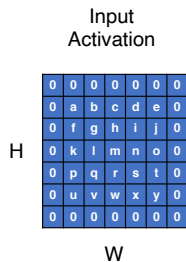


H: Height of Input Activation
W: Width of Input Activation
R: Height of Weight
S: Width of Weight
P: Height of Output Activation
Q: Width of Output Activation
stride: # of rows/columns traversed per step

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

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

2D-Convolution

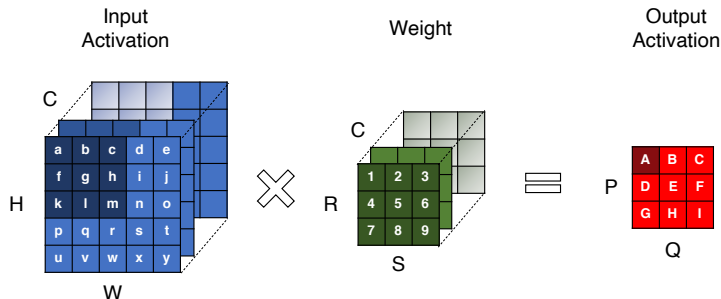


$$P = \frac{(H - R + 2 * pad)}{stride} + 1$$

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

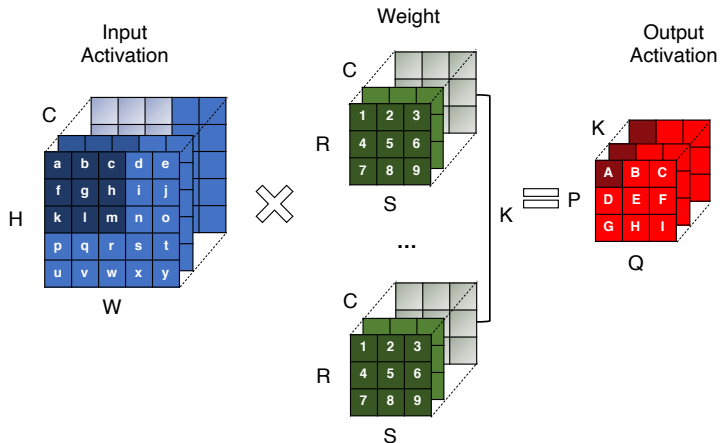
- H:** Height of Input Activation
- W:** Width of Input Activation
- R:** Height of Weight
- S:** Width of Weight
- P:** Height of Output Activation
- Q:** Width of Output Activation
- stride:** # of rows/columns traversed per step
- padding:** # of zero rows/columns added

3D-Convolution



H: Height of Input Activation
W: Width of Input Activation
R: Height of Weight
S: Width of Weight
P: Height of Output Activation
Q: Width of Output Activation
stride: # of rows/columns traversed per step
padding: # of zero rows/columns added
C: # of Input Channels

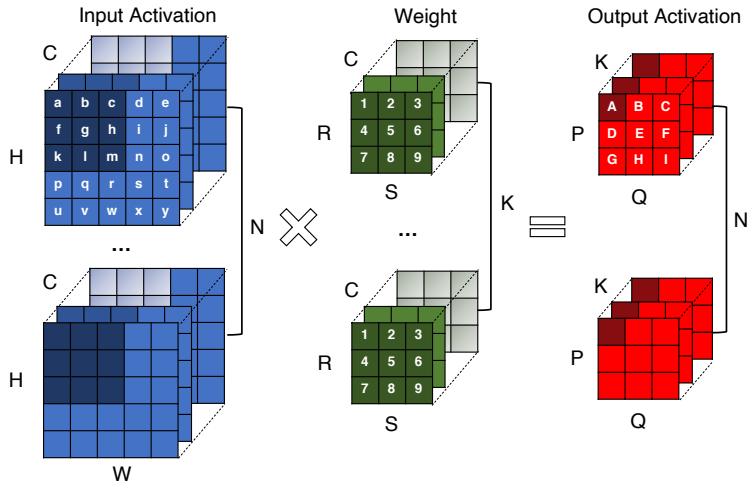
3D-Convolution



H: Height of Input Activation
W: Width of Input Activation
R: Height of Weight
S: Width of Weight
P: Height of Output Activation
Q: Width of Output Activation
stride: # of rows/columns traversed per step
padding: # of zero rows/columns added

C: # of Input Channels
K: # of Output Channels

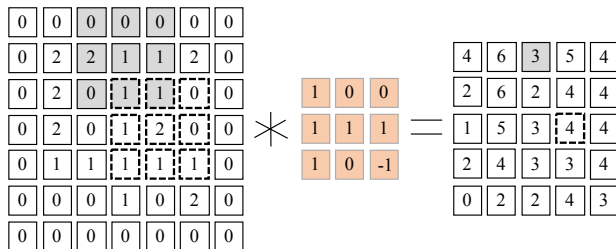
3D-Convolution



H: Height of Input Activation
W: Width of Input Activation
R: Height of Weight
S: Width of Weight
P: Height of Output Activation
Q: Width of Output Activation
stride: # of rows/columns traversed per step
padding: # of zero rows/columns added

C: # of Input Channels
K: # of Output Channels
N: Batch size

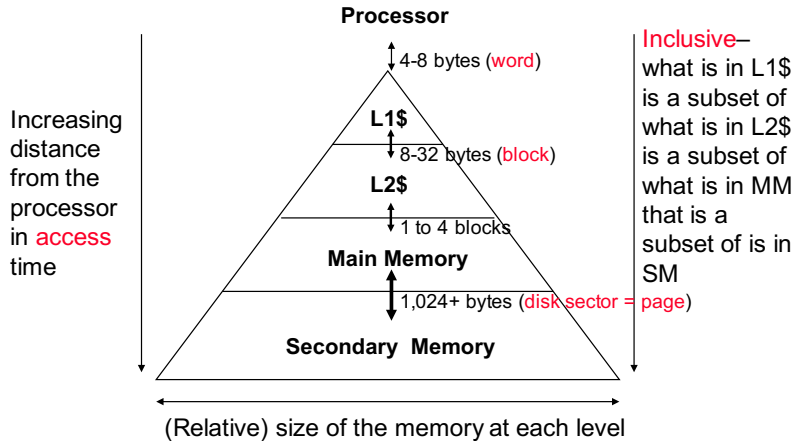
Convolution 101



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

Overview



Convolution 101

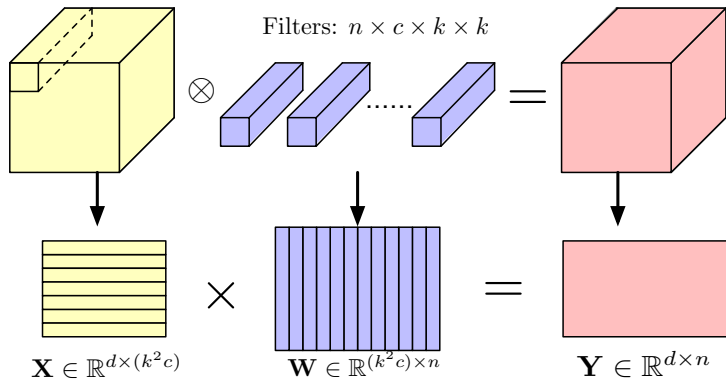
GEMM

Sparse Convolution

Direct Convolution

Further Discussions

Im2col (Image2Column) Convolution



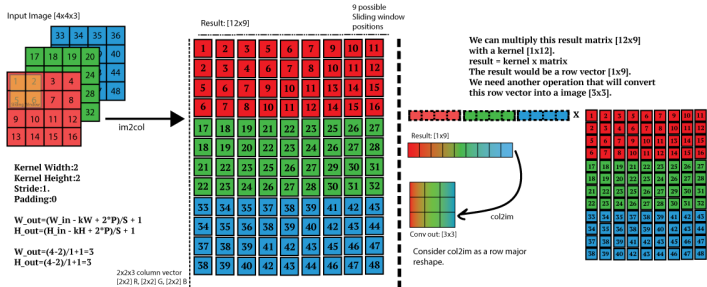
- ▶ Transform convolution to **matrix multiplication**
- ▶ **Unified** calculation for both convolution and fully-connected layers



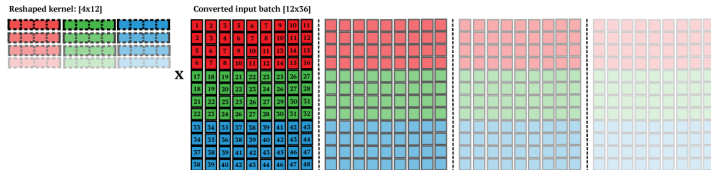
Im2col (Image2Column): Another View

Image to column operation (im2col)

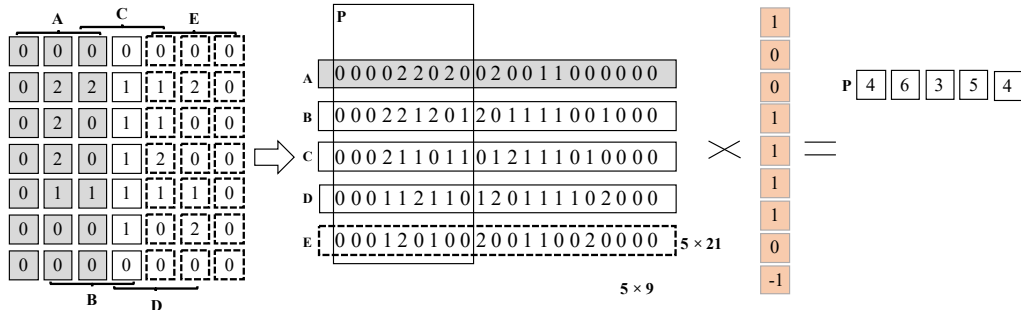
Slide the input image like a convolution but each patch become a column vector.



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



SOTA 1: Memory-efficient Convolution

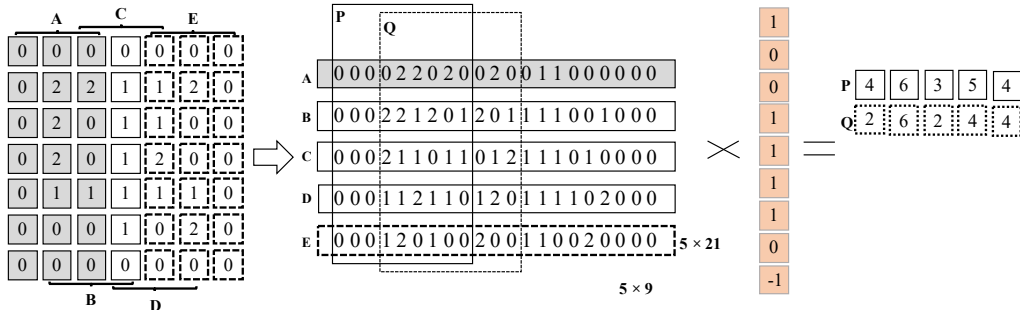


- ▶ Sub matrices in the lowered matrix will be “sgemm” ed in parallel
- ▶ Smaller memory foot print, cache locality, and explicit parallelism

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



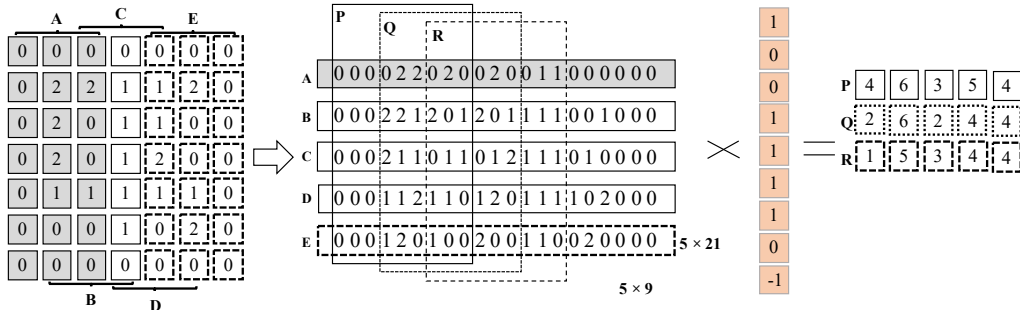
SOTA 1: Memory-efficient Convolution



- ▶ 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*.

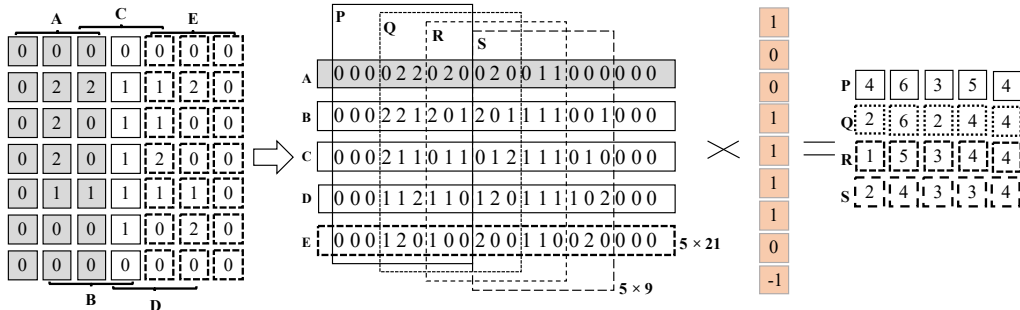
SOTA 1: Memory-efficient Convolution



- ▶ 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



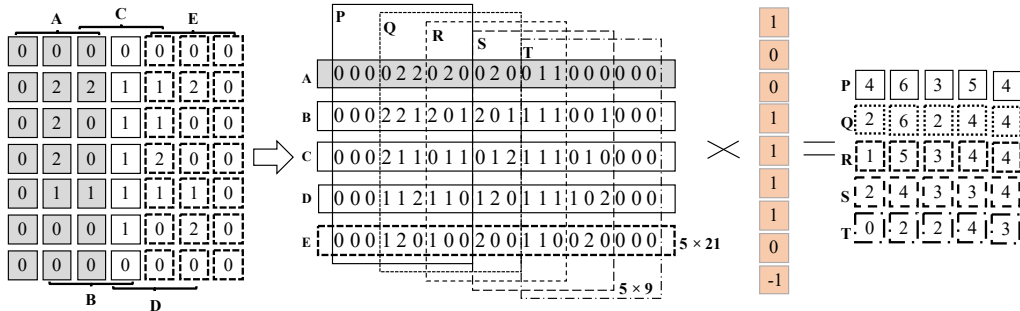
- ▶ Sub matrices in the lowered matrix will be “sgemm” ed in parallel
- ▶ Smaller memory foot print, cache locality, and explicit parallelism

2

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



SOTA 1: Memory-efficient Convolution



2

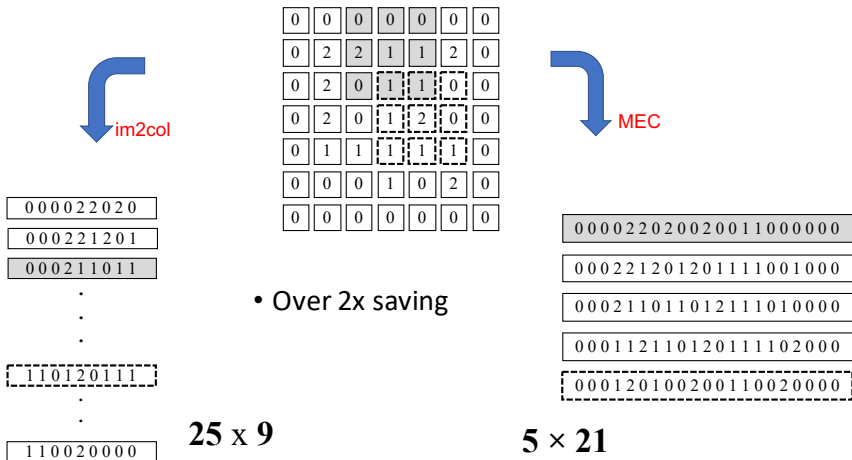
- ▶ Sub matrices in the lowered matrix will be “sgemm” ed in parallel
- ▶ Smaller memory foot print, cache locality, and explicit parallelism

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*.

Overview



Convolution 101

GEMM

Sparse Convolution

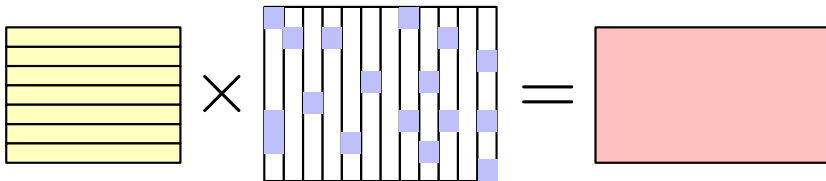
Direct Convolution

Further Discussions

Sparse Convolution



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



Sparse Convolution: Naive Implementation 1



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**
-

Sparse Convolution: Naive Implementation 1



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

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						
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
Clock Cycle	1	2	3	4	5	6	7

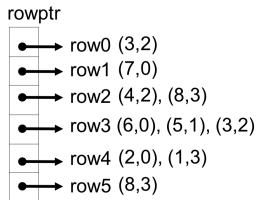
Sparse Matrix Representation



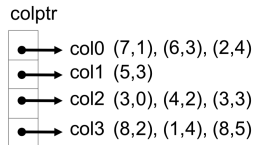
A

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

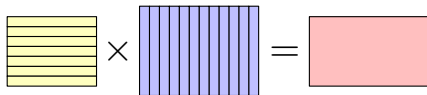


Compressed Sparse Row (CSR)



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.



SOTA 2: Sparse Convolution

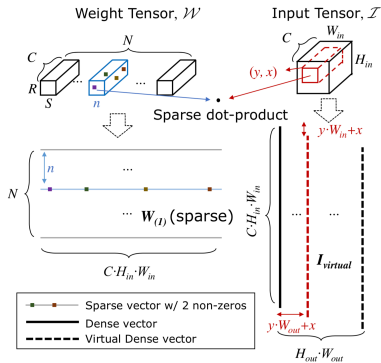


Figure 1: Conceptual view of the direct sparse convolution algorithm. Computation of output value at (y,x) th position of n th output channel is highlighted.

```

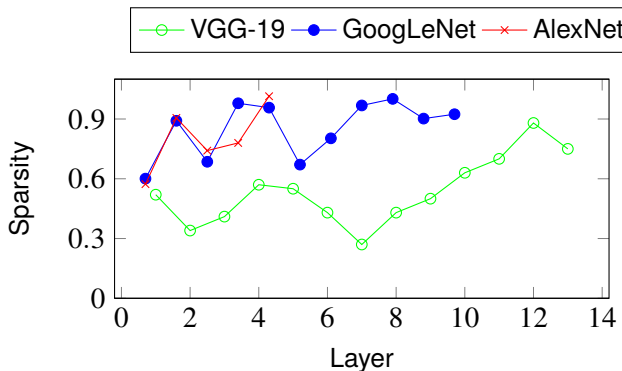
for each output channel n {
  for j in [W.rowptr[n], W.rowptr[n+1]] {
    off = W.colidx[j]; coeff = W.value[j]
    for (int y = 0; y < H_OUT; ++y) {
      for (int x = 0; x < W_OUT; ++x) {
        out[n][y][x] += coeff*in[off+f(0,y,x)]
      }
    }
  }
}
    
```

Figure 2: Sparse convolution pseudo code. Matrix \mathbf{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) , $\text{W.colidx}[j]$ contains the offset to (c,r,s) th element of tensor in, which is pre-computed by layout function as $f(c,r,s)$. If 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 in by $(0,y,x)$.

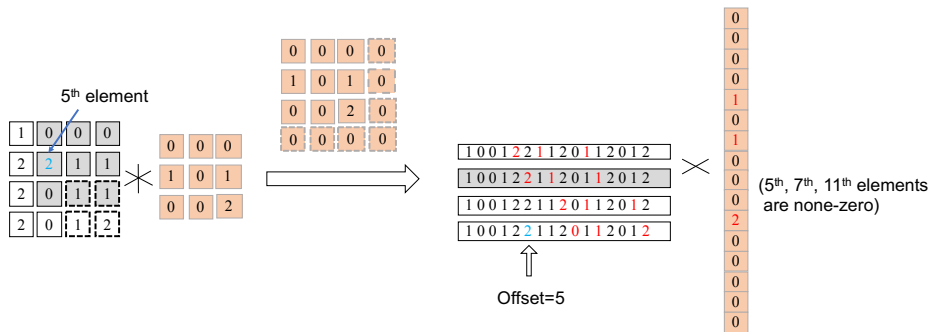
Discussion: Sparse-Sparse Convolution



- ▶ 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.



Discussion: Sparse-Sparse Convolution



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

Overview



Convolution 101

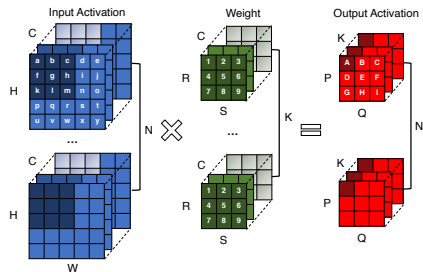
GEMM

Sparse Convolution

Direct Convolution

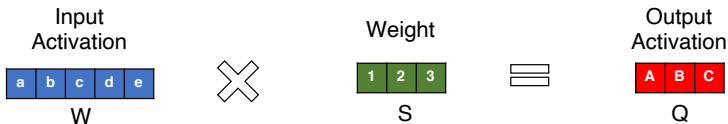
Further Discussions

Direct Convolution



```
for (n=0; n<N; n++) {
  for (k=0; k<K; k++) {
    for (p=0; p<P; p++) {
      for (q=0; q<Q; q++) {
        for (r=0; r<R; r++) {
          for (s=0; s<S; s++) {
            for (c=0; c<C; c++) {
              h = p * stride - pad + r;
              w = q * stride - pad + s;
              OA[n][k][p][q] +=
                IA[n][c][h][w]
                * W[k][c][r][s];
            }
          }
        }
      }
      OA[n][k][p][q] = Activation(OA[n][k][p][q]);
    }
  }
}
```

1D Convolution Example



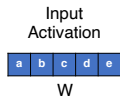
```
for(q=0; q<Q; q++){
  for (s=0; s<S; s++){
    OA[q] += IA[q+s] * W[s];
  }
}
```

**Output Stationary (OS)
Dataflow**

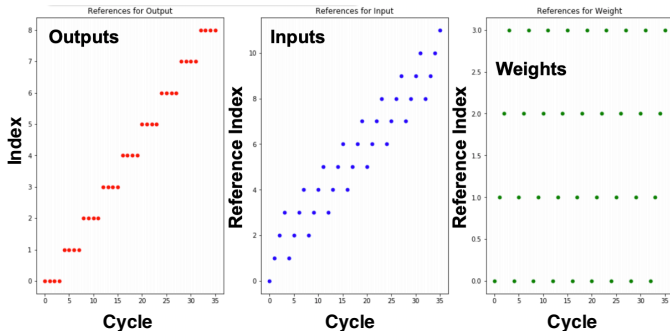
```
for (s=0; s<S; s++){
  for(q=0; q<Q; q++){
    OA[q] += IA[q+s] * W[s];
  }
}
```

**Weight Stationary (WS)
Dataflow**

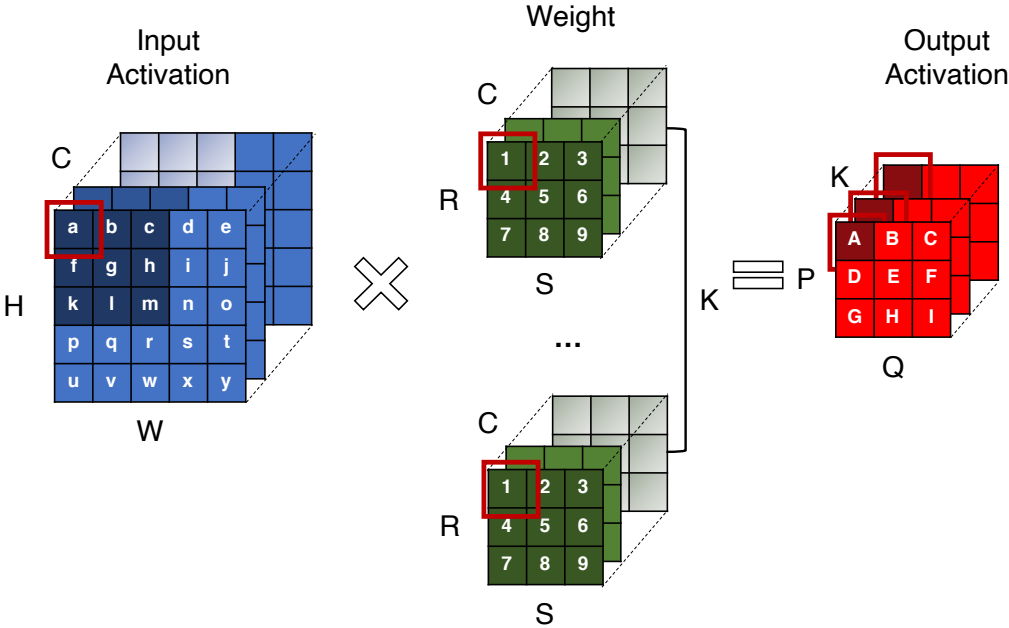
Buffer Access Pattern 1: Output Stationary



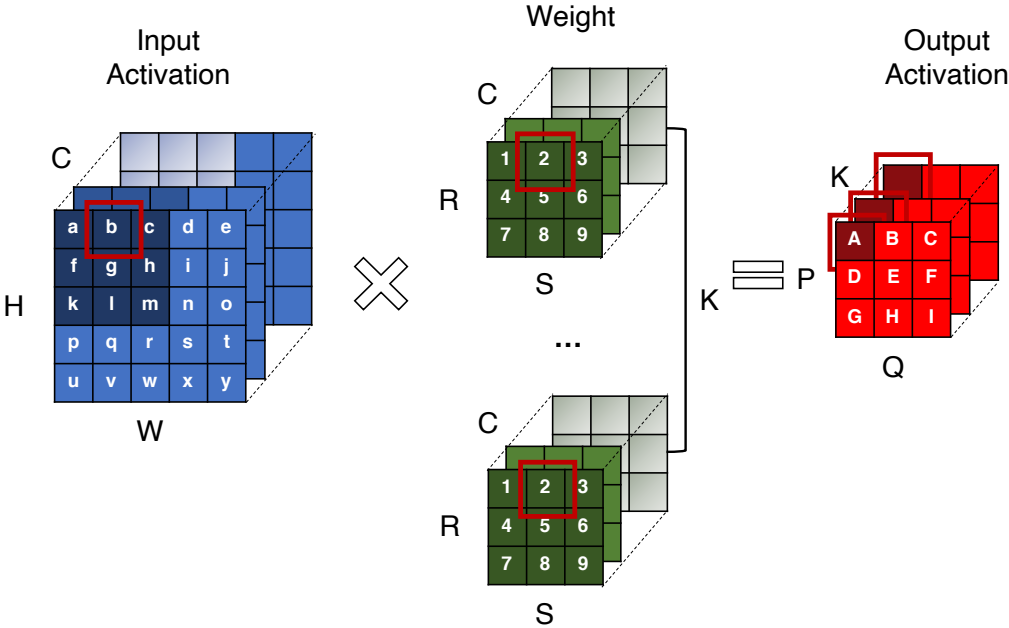
```
for (q=0; q<Q; q++){ // Q =9
  for (s=0; s<S; s++){ // S=4
    OA[q] += IA[q+s] * W[s];
  }
}
```



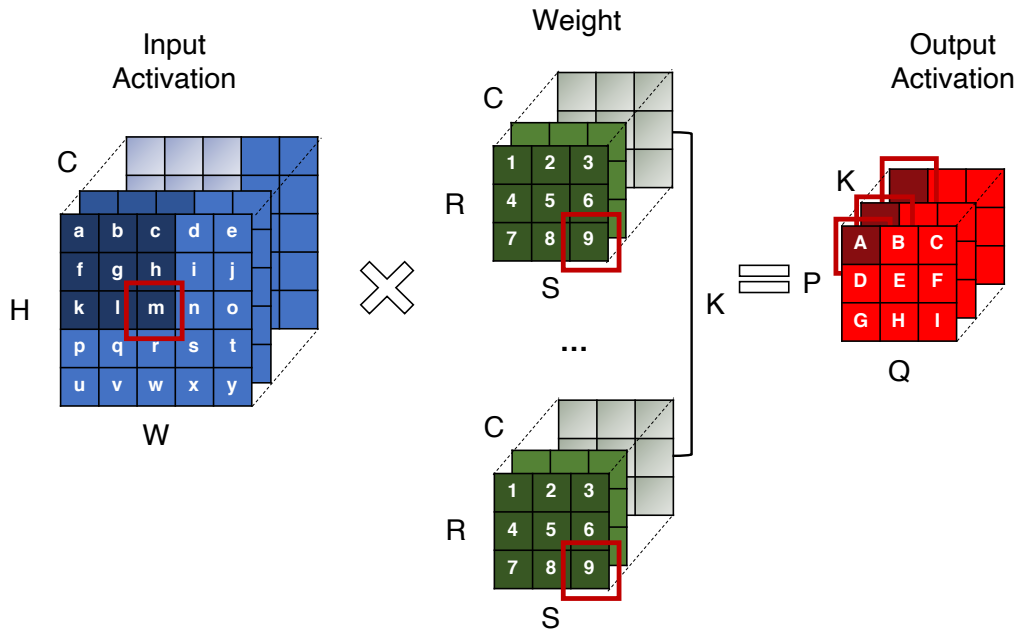
Output Stationary in 3D Convolution Scenario



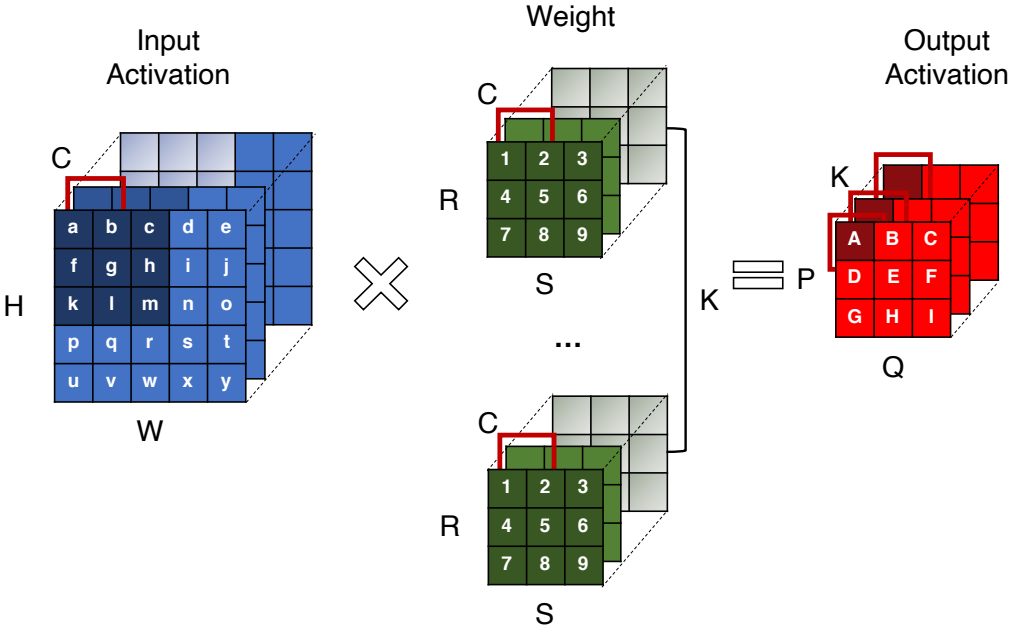
Output Stationary in 3D Convolution Scenario



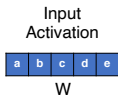
Output Stationary in 3D Convolution Scenario



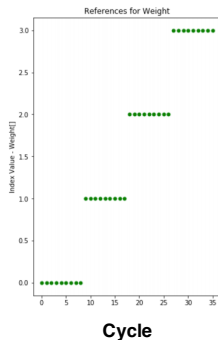
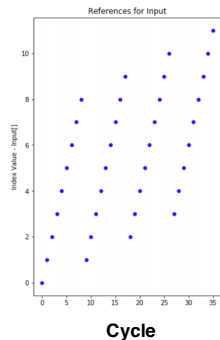
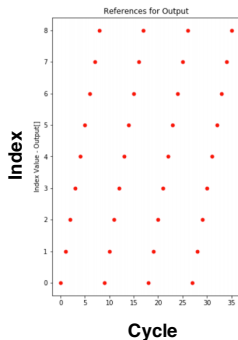
Output Stationary in 3D Convolution Scenario



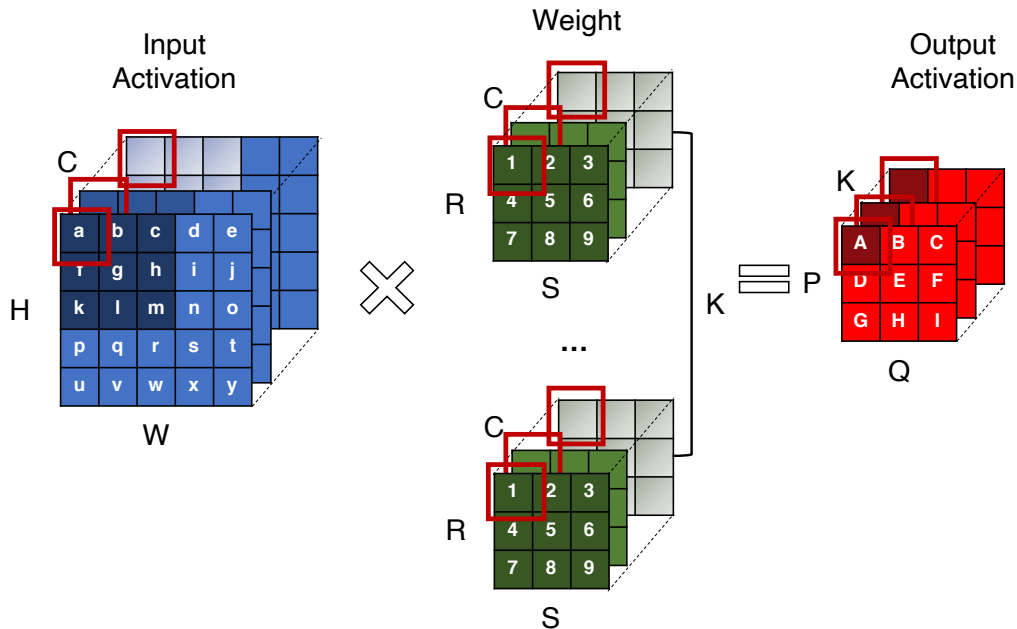
Buffer Access Pattern 2: Weight Stationary



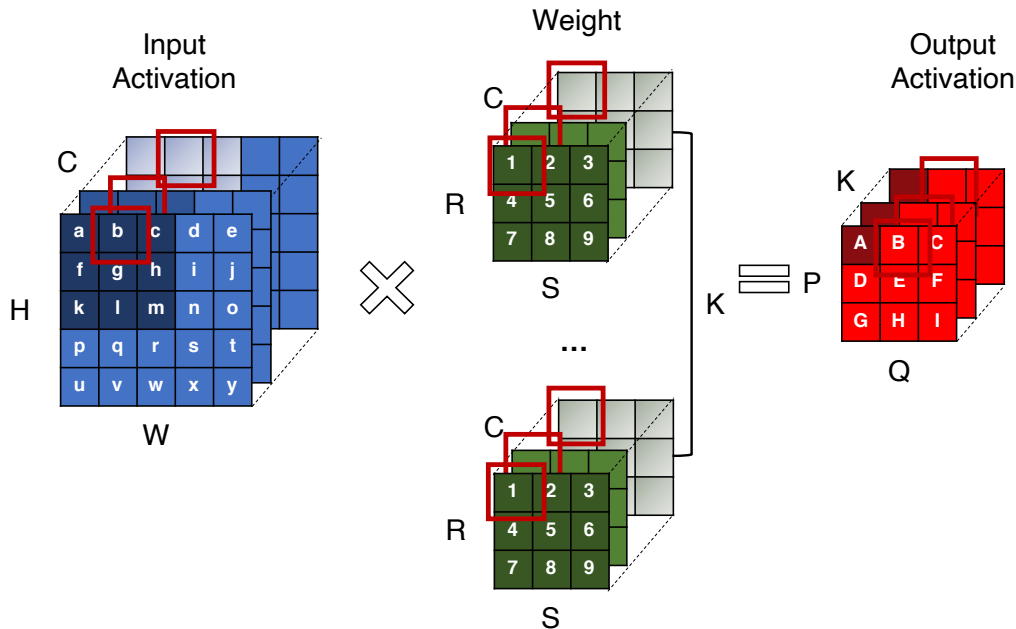
```
for (s=0; s<S; s++) { // S=4
  for (q=0; q<Q; q++) { // Q=9
    OA[q] += IA[q+s] * W[s];
  }
}
```



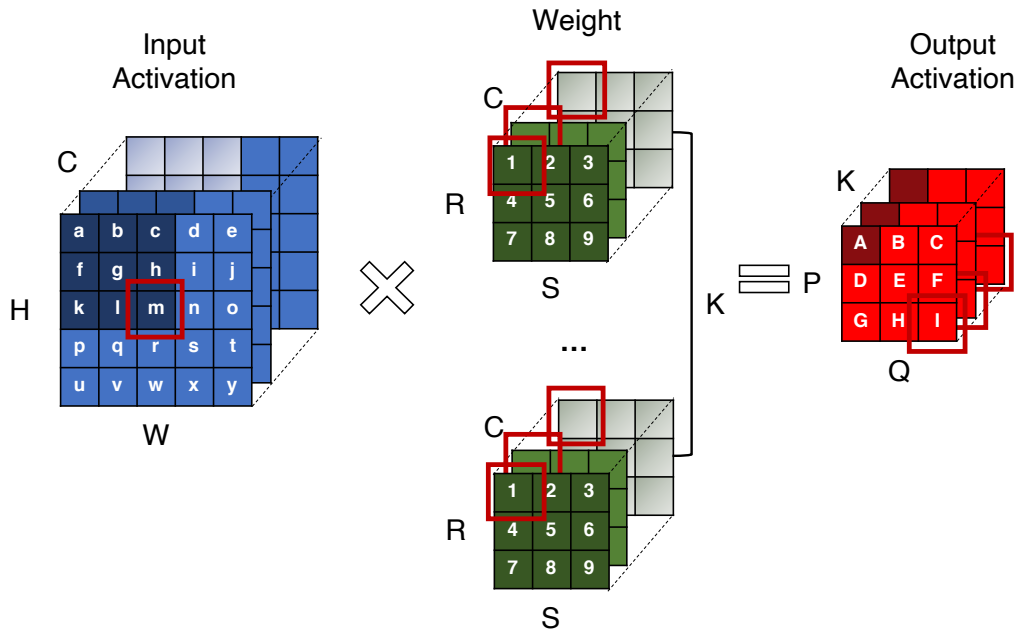
Weight Stationary in 3D Convolution Scenario



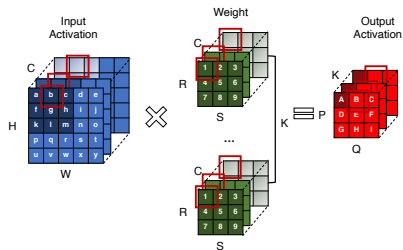
Weight Stationary in 3D Convolution Scenario



Weight Stationary in 3D Convolution Scenario



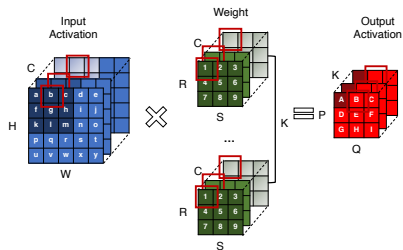
Weight Stationary Dataflow



- What we had before:

```
for (n=0; n<N; n++) {
  for (k=0; k<K; k++) {
    for (p=0; p<P; p++) {
      for (q=0; q<Q; q++) {
        OA[n][k][p][q] = 0;
        for (r=0; r<R; r++) {
          for (s=0; s<S; s++) {
            for (c=0; c<C; c++) {
              h = p * stride - pad + r;
              w = q * stride - pad + s;
              OA[n][k][p][q] +=
                IA[n][c][h][w]
                  * W[k][c][r][s];
            }
          }
        }
      }
    }
  }
}
```

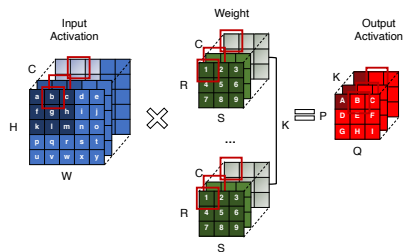

Weight Stationary Dataflow



- Change temporal ordering

```
for (n=0; n<N; n++) {  
  for (r=0; r<R; r++) {  
    for (s=0; s<S; s++) {  
      for (c=0; c<C; c++) {  
        float curr_w = W[r][s][c][k];  
        for (p=0; p<P; p++) {  
          for (q=0; q<Q; q++) {  
            h = p * stride - pad + r;  
            w = q * stride - pad + s;  
            OA[n][k][p][q] +=  
              IA[n][c][h][w]  
                * curr_w;  
          }  
        }  
      }  
    }  
  }  
}
```

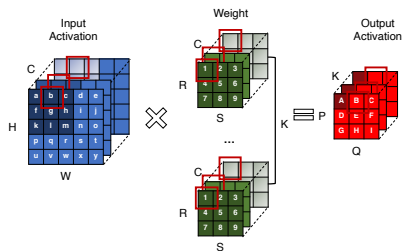
Weight Stationary Dataflow



- Apply spatial parallelism

```
for (n=0; n<N; n++) {
  for (r=0; r<R; r++) {
    for (s=0; s<S; s++) {
      spatial_for (c=0; c<C; c++) {
        spatial_for (k=0; k<K; k++) {
          float curr_w = W[r][s][c][k];
          for (p=0; p<P; p++) {
            for (q=0; q<Q; q++) {
              h = p * stride - pad + r;
              w = q * stride - pad + s;
              OA[n][k][p][q] +=
                IA[n][c][h][w]
                  * curr_w;
            }
          }
        }
      }
    }
  }
}
```

Weight Stationary Dataflow



- Apply temporal tiling

```
for (n=0; n<N; n++) {
  for (r=0; r<R; r++) {
    for (s=0; s<S; s++) {
      for (c_t=0; c_t<C/16; c_t++) {
        for (k_t=0; k_t<K/64; k_t++) {
          spatial_for (c_s=0; c_s<16; c_s++) {
            spatial_for (k_s=0; k_s<64; k_s++) {
              int curr_c = c_t * 16 + c_s;
              int curr_k = k_t * 64 + k_s;
              float curr_w = W[r][s][curr_c][curr_k];
              for (p=0; p<P; p++) {
                for (q=0; q<Q; q++) {
                  h = p * stride - pad + r;
                  w = q * stride - pad + s;
                  OA[n][curr_k][p][q] +=
                    IA[n][curr_c][h][w]
                    * curr_w;
                }
              }
            }
          }
        }
      }
    }
  }
}
```

Overview



Convolution 101

GEMM

Sparse Convolution

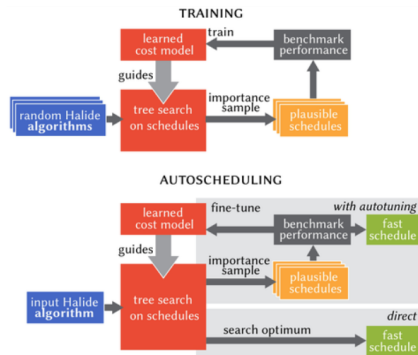
Direct Convolution

Further Discussions

Example: Halide, SIGGRAPH '2019



- ▶ <https://youtu.be/3uiEyEKji0M>
- ▶ “We generate schedules for Halide programs using tree search over the space of schedules guided by a learned cost model and optional autotuning. The cost model is trained by benchmarking thousands of randomly-generated Halide programs and schedules. The resulting code significantly outperforms prior work and human experts.”



5

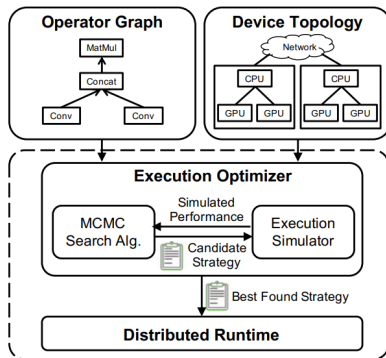
⁵Andrew Adams et al. (2019). “Learning to optimize halide with tree search and random programs”. In: *ACM Trans. Graph.* 38.4, 121:1–121:12. doi: 10.1145/3306346.3322967. URL: <https://doi.org/10.1145/3306346.3322967>.



Example: FlexFlow, SysML'2019

- ▶ “The optimizer uses a MCMC search algorithm to explore the space of possible parallelization strategies and iteratively proposes candidate strategies that are evaluated by an execution simulator.”

6



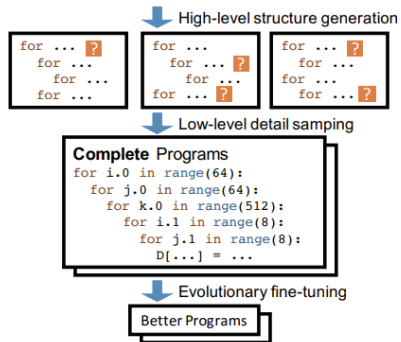
⁶Zhihao Jia, Matei Zaharia, and Alex Aiken (2019). “Beyond Data and Model Parallelism for Deep Neural Networks”. In: *Proceedings of Machine Learning and Systems 2019, MLSys 2019, Stanford, CA, USA, March 31 - April 2, 2019*. Ed. by Ameet Talwalkar, Virginia Smith, and Matei Zaharia. [mlsys.org. URL: https://proceedings.mlsys.org/book/265.pdf](https://proceedings.mlsys.org/book/265.pdf).



Example: Anso: AutoTVM v2.0, arXiv

- ▶ “We present Anso, a tensor program generation framework for deep learning applications. Compared with existing search strategies, Anso explores much more optimization combinations by sampling programs from a hierarchical representation of the search space.”

8



⁸Lianmin Zheng et al. (2020). “Anso: Generating High-Performance Tensor Programs for Deep Learning”. In: *CoRR* abs/2006.06762. arXiv: 2006.06762. URL: <https://arxiv.org/abs/2006.06762>.