

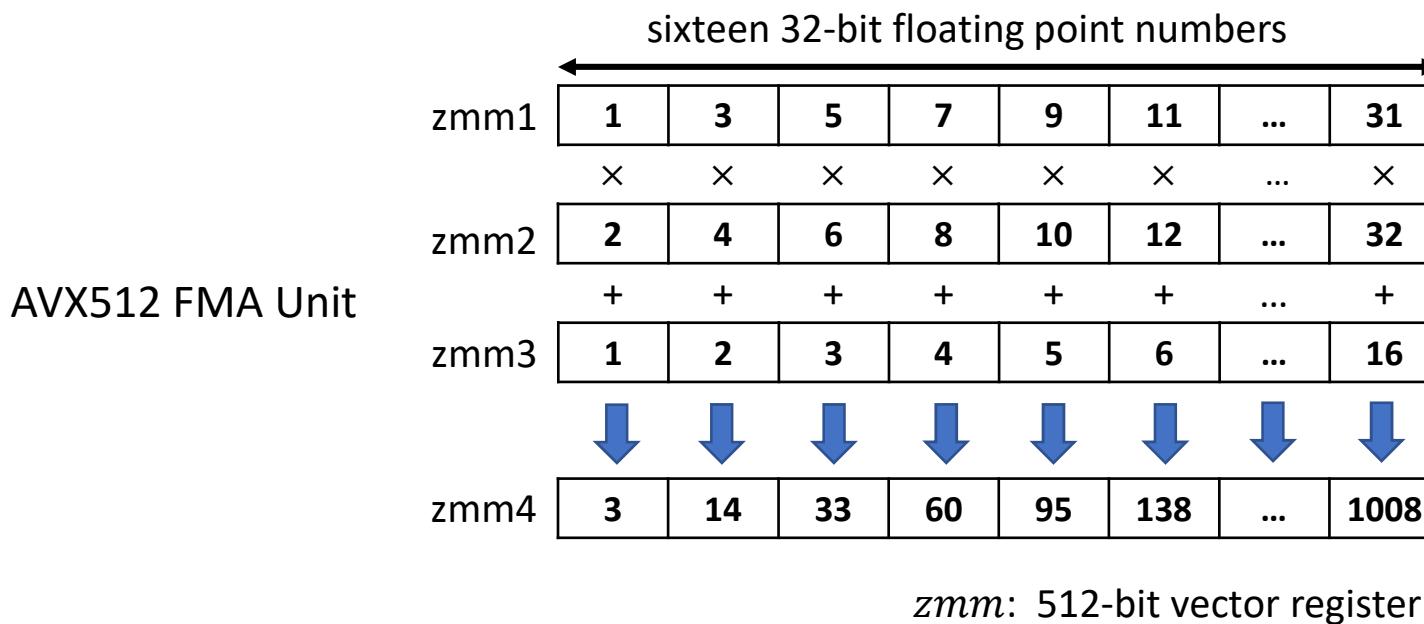
# An Efficient Multicore CPU Implementation for Convolution- Pooling Computation in CNN

Hiroki Kataoka, Kohei Yamashita, Yasuaki Ito, Koji Nakano  
Hiroshima University

Akihiko Kasagi, Tsuguchika Tabaru  
Fujitsu Laboratories Ltd.

# Contribution

- We present an efficient implementation of the convolutional-pooling in the multicore CPU with AVX-512 FMA units.
  - We use convolution interchange to reduce the computational cost.
  - Our proposed method is up to 2.82 times faster than the multiple convolution and then pooling by DNNL.
  - We incorporate the proposed implementation into TensorFlow to perform them as a TensorFlow operation.

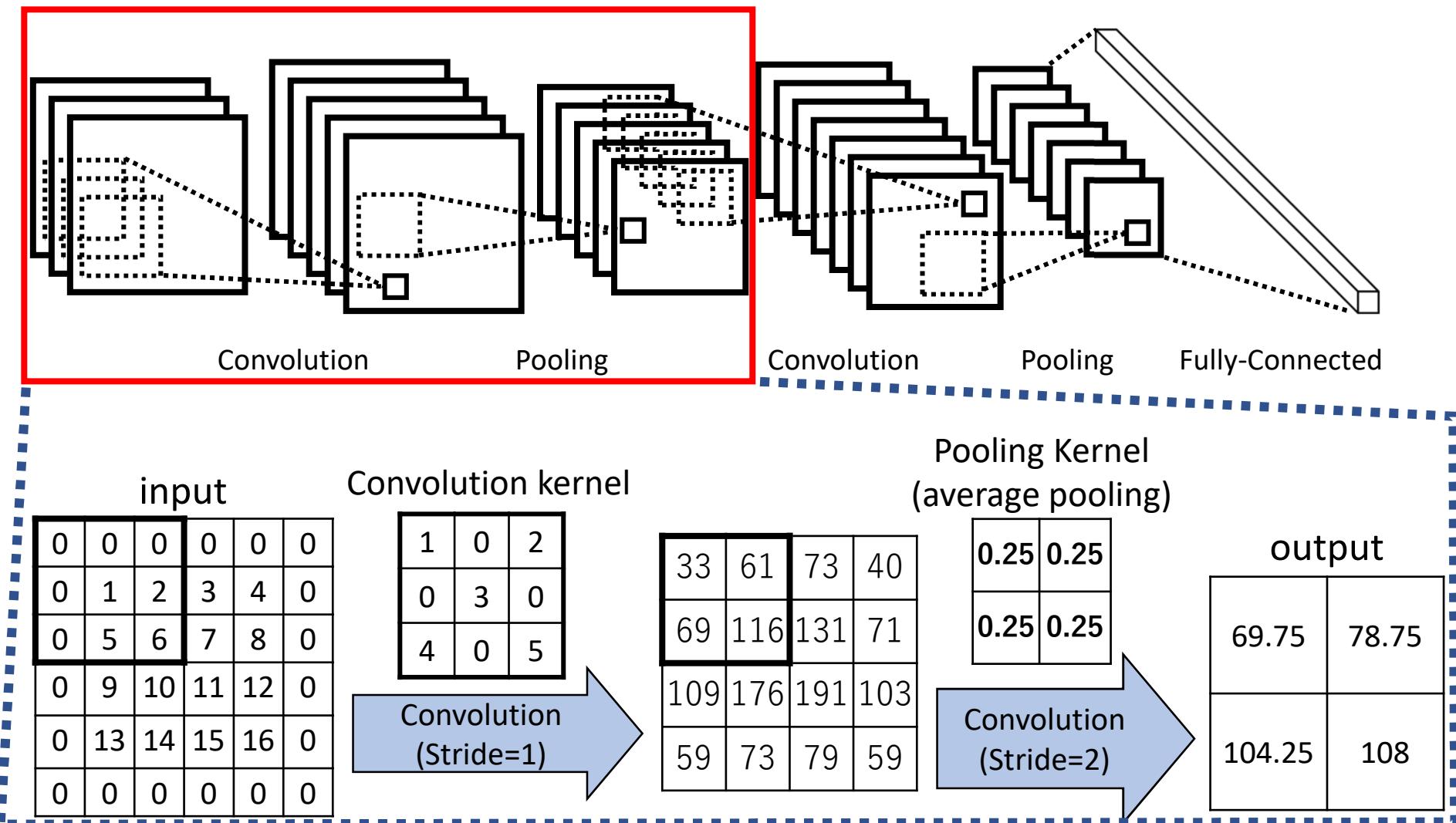


# Convolutional Neural Network

Convolutional Neural Network (CNN) is commonly used for deep learning

CNN consists of some convolution operations and pooling operations

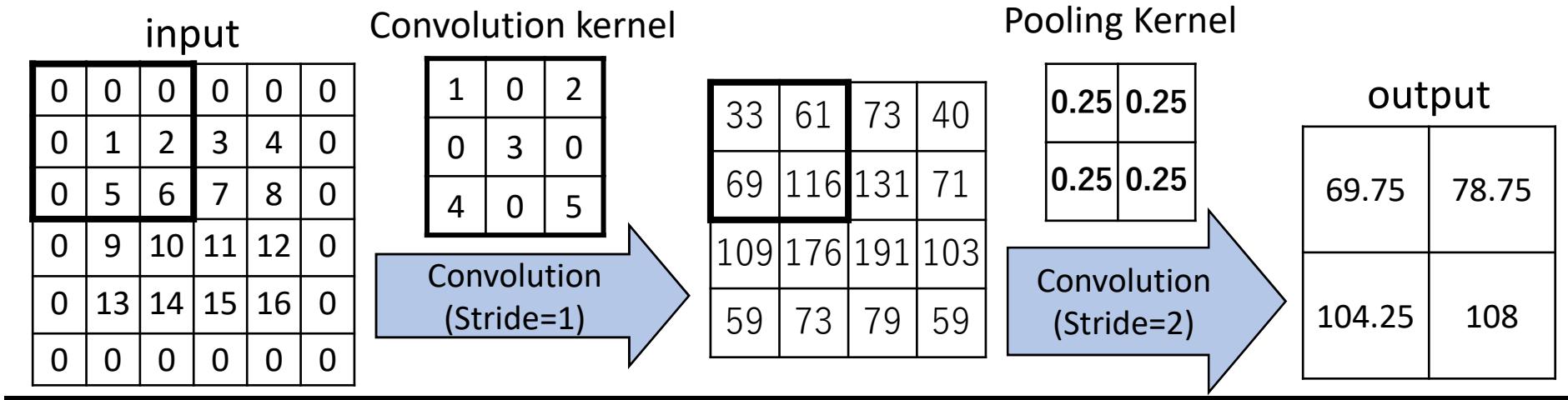
Our techniques accelerate



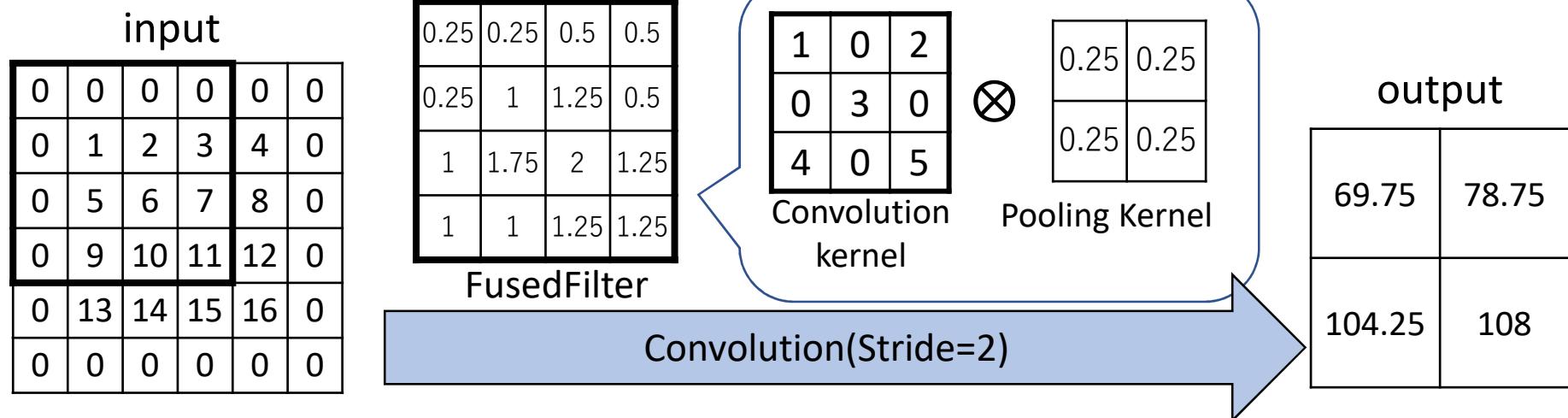
# Fused Filter algorithm (proposed)

The idea of the fused-filter method is to perform only one convolution to compute the convolution-pooling.

## ◆ Naïve algorithm



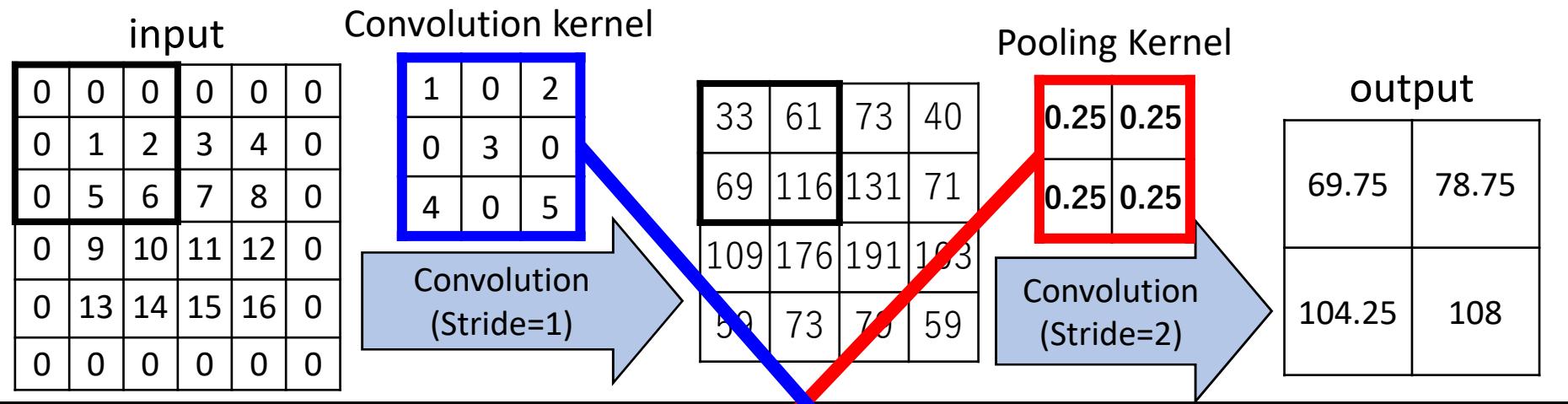
## ◆ Fused Filter algorithm



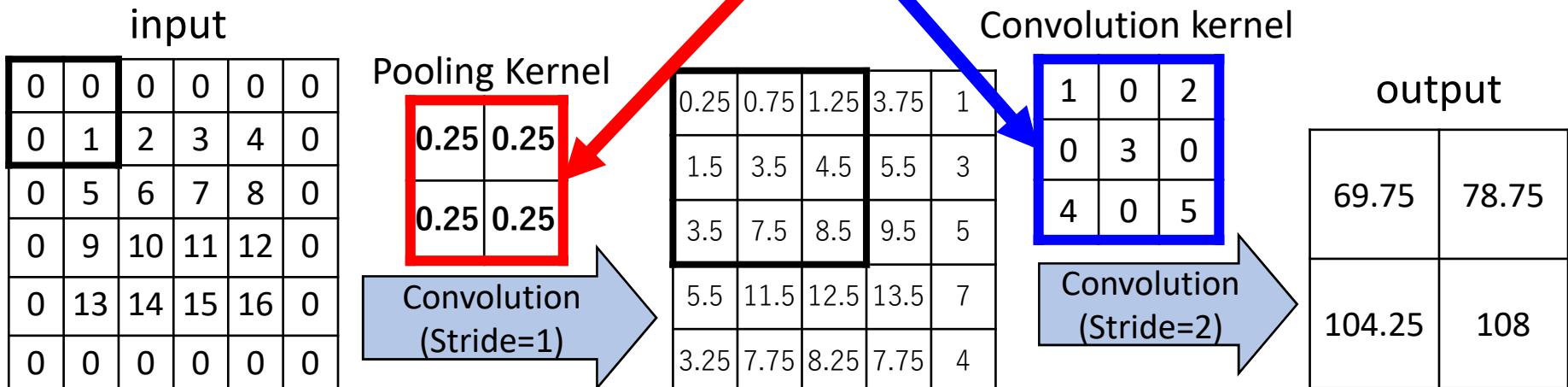
# Direct Sum algorithm (proposed)

The idea of the direct-sum method is to use the convolution interchange technique to implement the convolution-pooling.

## ◆ Naïve algorithm



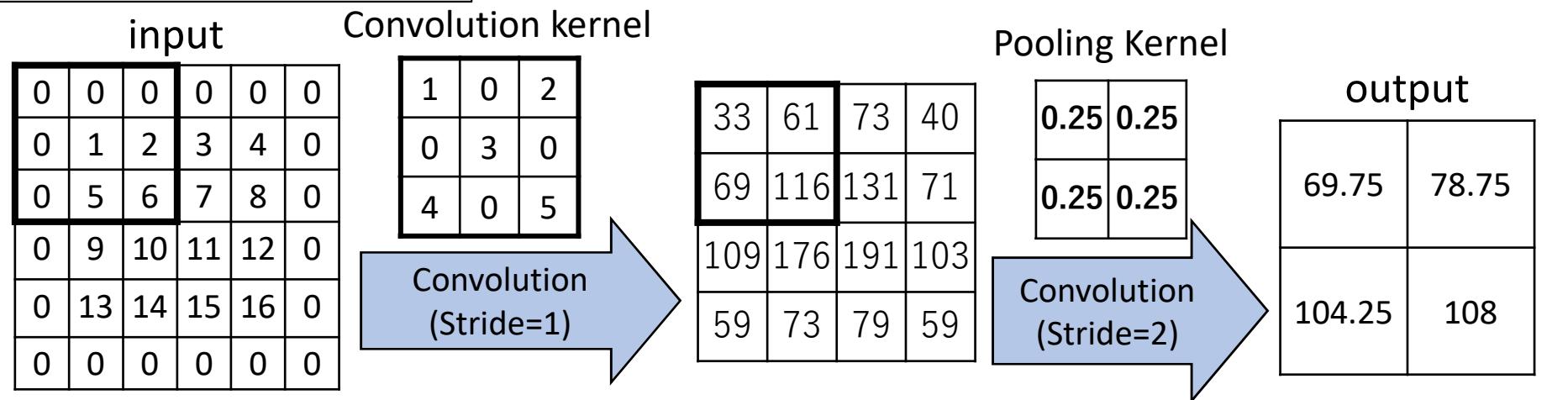
## ◆ Direct Sum algorithm



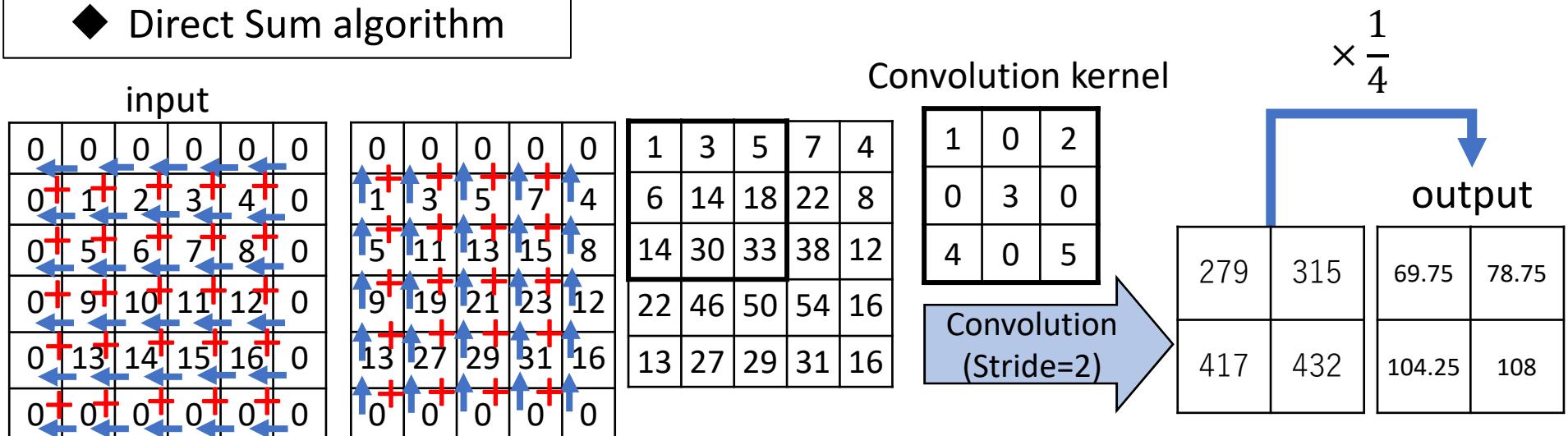
# Direct Sum algorithm (proposed)

Instead of the 2x2 pooling, we perform pairwise-sum twice.

## ◆ Naïve algorithm



## ◆ Direct Sum algorithm

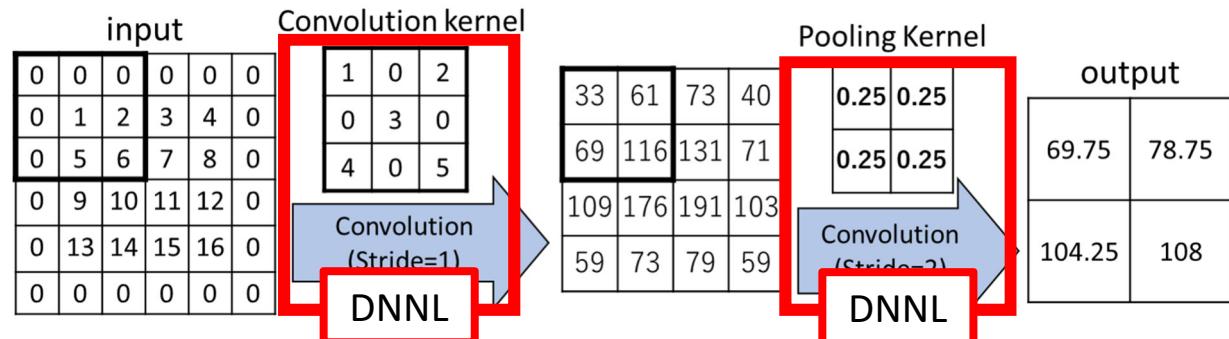


# Proposed implementation

We adopt DNNL functions to perform convolution and pooling since the functions are well-optimized. In our performance evaluation, we compare the following three algorithms.

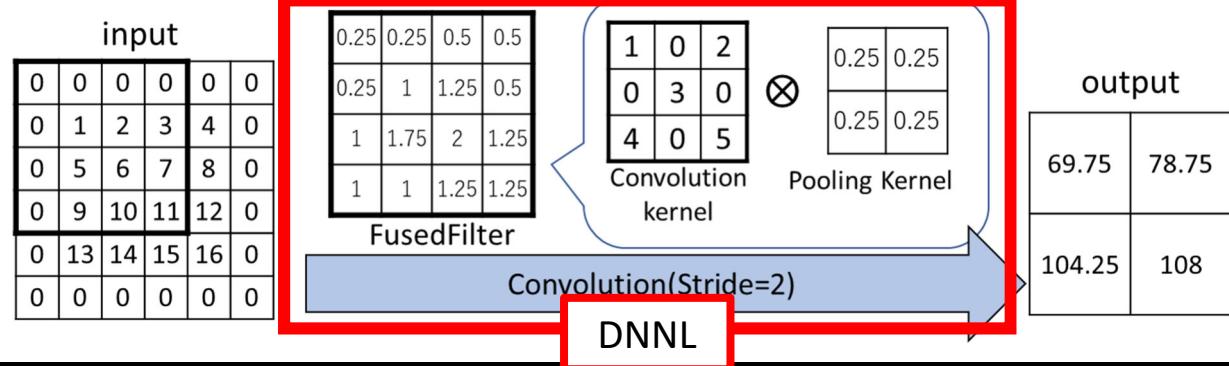
## Naïve algorithm

Convolution and then pooling



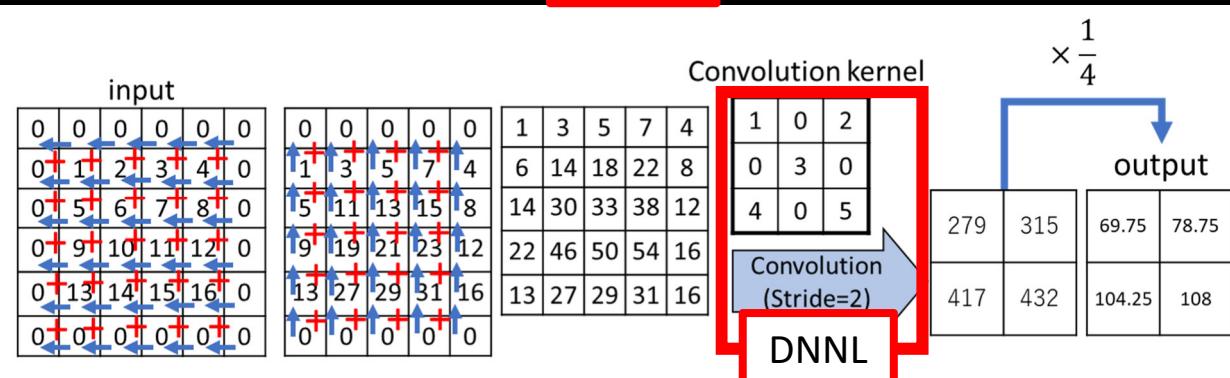
## Fused Filter algorithm

Convolution with the fused filter



## Direct Sum algorithm

Average pooling implemented by pairwise-sums in horizontal and vertical directions, and then convolution



# Performance Evaluation

**CPU:** Core i9-7980XE (18 physical cores)

**Maximum clock frequency:** 4.4GHz(single core), 3.4GHz(18 cores)

**OpenMP:** version 5.0, **DNNL:** version 1.1.0

Computation time [ms]

input:  $18 \times 18$ , kernel:  $3 \times 3$ , pooling size:  $2 \times 2$ , batch: 32

		input channel, output channel				
		32,32	64,64	128,128	256,256	512,512
naïve (DNNL)	convolution	0.133	0.272	0.816	3.750	11.007
	pooling	0.018	0.032	0.067	0.152	0.309
	total	<b>0.151</b>	<b>0.304</b>	<b>0.883</b>	<b>3.902</b>	<b>11.316</b>
proposed (fused-filter)	convolution	<b>0.100</b>	<b>0.210</b>	<b>0.572</b>	<b>2.055</b>	<b>8.849</b>
	speed-up	1.510	1.448	1.544	1.899	1.279
proposed (direct-sum)	direct-sum	0.022	0.030	0.066	0.222	0.707
	convolution	0.064	0.166	0.363	1.162	4.940
	total	<b>0.086</b>	<b>0.196</b>	<b>0.429</b>	<b>1.384</b>	<b>5.647</b>
	speed-up	1.756	1.551	2.058	2.820	2.004

# TensorFlow Incorporation

- We incorporate the proposed implementation into TensorFlow to perform them as a TensorFlow operation.



Computation time [ms]

input:  $18 \times 18$ , kernel:  $3 \times 3$ , pooling size:  $2 \times 2$ , batch: 32

input channel, output channel	32, 32	64, 64	128, 128	256, 256	512, 512
naïve	0.8296	1.6692	3.4670	9.1944	25.1610
incorporation	0.7047	1.1260	2.0238	4.0289	12.1624
speed-up	1.1772	1.4824	1.7131	2.2821	2.0687

**CPU:** Core i9-7980XE (18 physical cores)

**Maximum clock frequency:** 4.4GHz(single core), 3.4GHz(18 cores)

**OpenMP:** version 5.0, **DNNL:** version 1.1.0

# Conclusion

- We presented an efficient implementation of the convolutional-pooling in the multicore CPU with AVX-512 FMA units.
  - Convolution interchange to reduce the computational cost.
  - Our proposed method is up to 2.82 times faster than the multiple convolution and then pooling by DNNL.
  - We incorporate the proposed implementation into TensorFlow to perform them as a TensorFlow operation.

Contact: [yasuaki@cs.hiroshima-u.ac.jp](mailto:yasuaki@cs.hiroshima-u.ac.jp)