

4-2 Convolution Neural Network I

Zhonglei Wang

WISE and SOE, XMU, 2025

Contents

1. Introduction

2. Basic structure

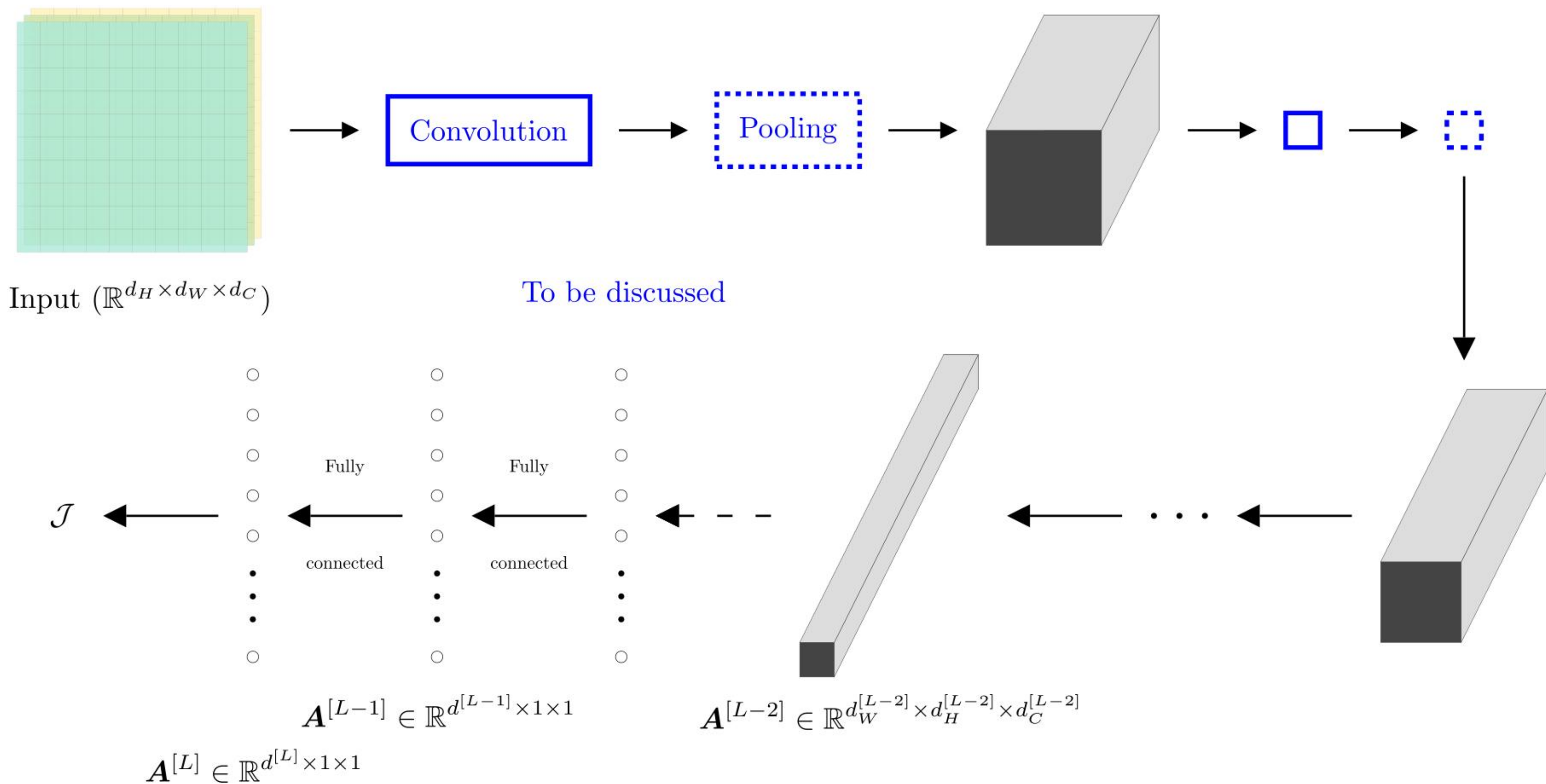
3. Forward propagation

4. Backpropagation

5. Understanding convolution networks

Introduction

1. Given an image, why not use a fully connected neural network?
 - Number of parameters is **overwhelmingly** large, easily leading to overfitting
 - Overlooks the spatial structure
 - Not location (translation) invariant
2. Consider a new network, convolution neural network (CNN), for image processing



Summary

1. Generally, $d_C^{[l]}$ increases as l increases, but $d_W^{[l]}$ and $d_H^{[l]}$ decrease
2. The vectorization step does not involve any model parameter
3. Typically, there are three “layers” for CNN
 - **Convolution** layers (CONV)
 - **Pooling** layers (POOL)
 - Fully connected layers (FC)

Notation

1. l : layer index
2. $d_H^{[l]}$: height of the “image” of the l th layer
3. $d_W^{[l]}$: width of the “image” of the l th layer
4. $d_C^{[l]}$: number of channels of the l th layer
5. $f^{[l]}$: kernel size associated with the l th layer
6. $p^{[l]}$: padding associated with the l th layer
7. $s^{[l]}$: stride associated with the l th layer

Vectorization

1. $\mathbf{A}^{[l-1]} \in \mathbb{R}^{n \times d_H^{[l-1]} \times d_W^{[l-1]} \times d_C^{[l-1]}}$: input for the l th layer

2. $\mathbf{Z}^{[l]} \in \mathbb{R}^{n \times d_H^{[l]} \times d_W^{[l]} \times d_C^{[l]}}$: linear transformed result

$$\mathbf{Z}^{[l]} = \mathbf{A}^{[l-1]} \text{ “} * \text{” } \mathbf{W}^{[l]} \text{ “} + \text{” } \mathbf{b}^{[l]}$$

$$\mathbf{A}^{[l]} = \sigma^{[l]}(\mathbf{Z}^{[l]})$$

- $\mathbf{W}^{[l]} \in \mathbb{R}^{d_C^{[l]} \times f^{[l]} \times f^{[l]} \times d_C^{[l-1]}}$: kernels
- $\mathbf{b}^{[l]} \in \mathbb{R}^{d_C^{[l]} \times 1 \times 1 \times 1}$: bias term
- “ $*$ ”: convolution for each sample and each channel
- “ $+$ ”: one common bias for each channel
- $\sigma^{[l]}(\cdot)$: activation function for the l th layer

Vectorization

1. For simplicity, we show forward- and back-propagation for $n = 1$

- $\mathbf{A}^{[l-1]} \in \mathbb{R}^{d_H^{[l-1]} \times d_W^{[l-1]} \times d_C^{[l-1]}}$
- $\mathbf{Z}^{[l]} \in \mathbb{R}^{d_H^{[l]} \times d_W^{[l]} \times d_C^{[l]}}$

2. Forward propagation:

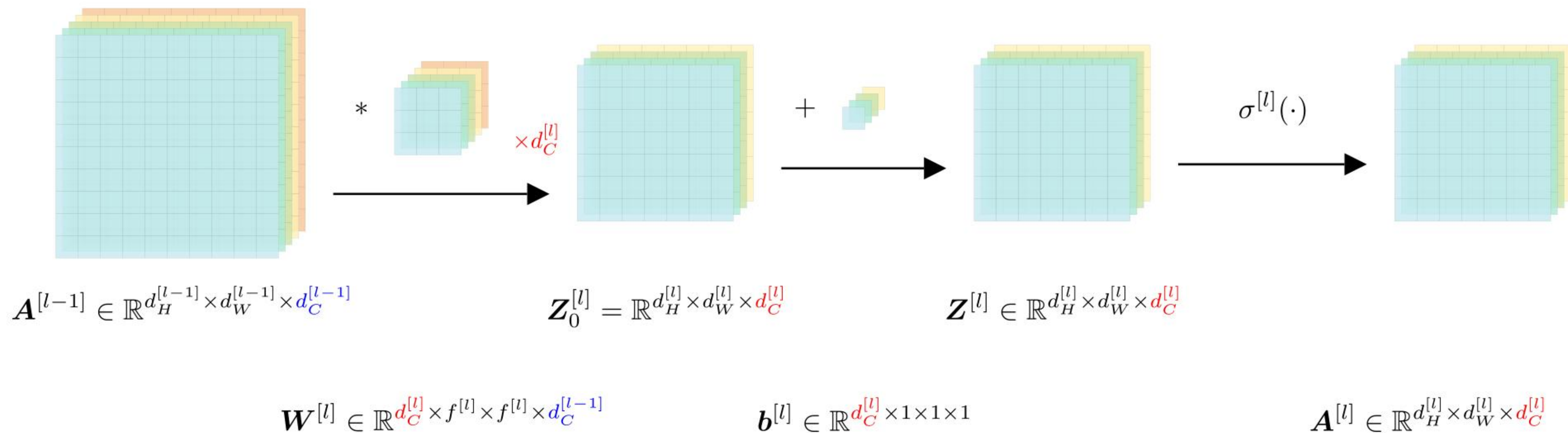
$$\begin{aligned}\mathbf{Z}^{[l]} &= \mathbf{A}^{[l-1]} \text{ “ } * \text{ ” } \mathbf{W}^{[l]} \text{ “ } + \text{ ” } \mathbf{b}^{[l]} \\ \mathbf{A}^{[l]} &= \sigma^{[l]}(\mathbf{Z}^{[l]})\end{aligned}$$

Convolution

1. Consider:

- $p^{[l]} = 0$
- $s^{[l]} = 1$

Convolution



Pooling

1. Pooling is typically used after convolution layers,
 - Reduce height and width for each **channel**
 - Achieves robustness by neglecting useless or repeated information
 - Increase the receptive field
 - Achieve robustness small translations, rotations, and other distortions in the input
 - Prevent overfitting, especially when combined with dropout

Pooling

1. Two types of pooling:

- Average pooling
- Max pooling

2. Notice

- Pooling is conducted for **each channel**
- Pooling does not involve new model parameters

Average pooling

1. Input size: 6×6
2. Kernel size: 2×2
3. Stride: 2

Average pooling

Input

3	6	5	4	8	9
1	7	9	6	8	0
5	0	9	6	2	0
5	2	6	3	7	0
9	0	3	2	3	1
3	1	3	7	1	7

Kernel

0.25	0.25
0.25	0.25

Result

4.25	6.0	6.25
3.0	6.0	2.25
3.25	3.75	3.0

Max pooling

1. Input size: 6×6
2. Kernel size: 2×2
3. Stride: 2

Max pooling

Input

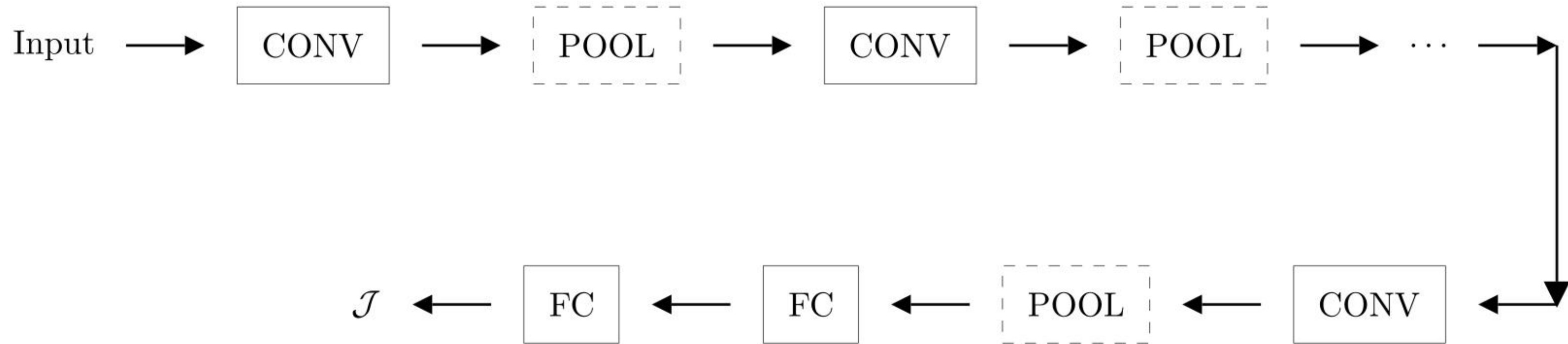
3	6	5	4	8	9
1	7	9	6	8	0
5	0	9	6	2	0
5	2	6	3	7	0
9	0	3	2	3	1
3	1	3	7	1	7

Kernel

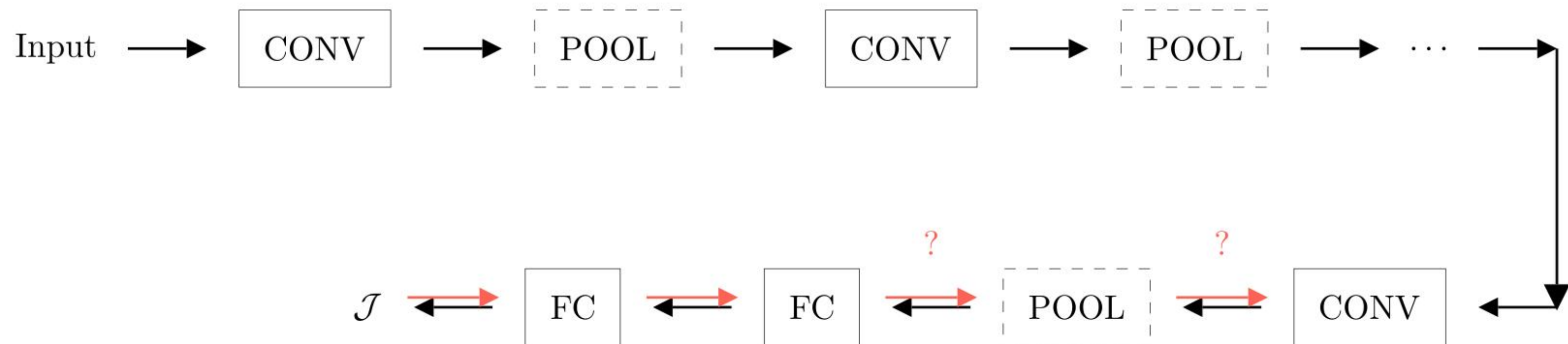
Result

7	9	9
5	9	7
9	7	7

Forward propagation



Backpropagation



Backpropagation

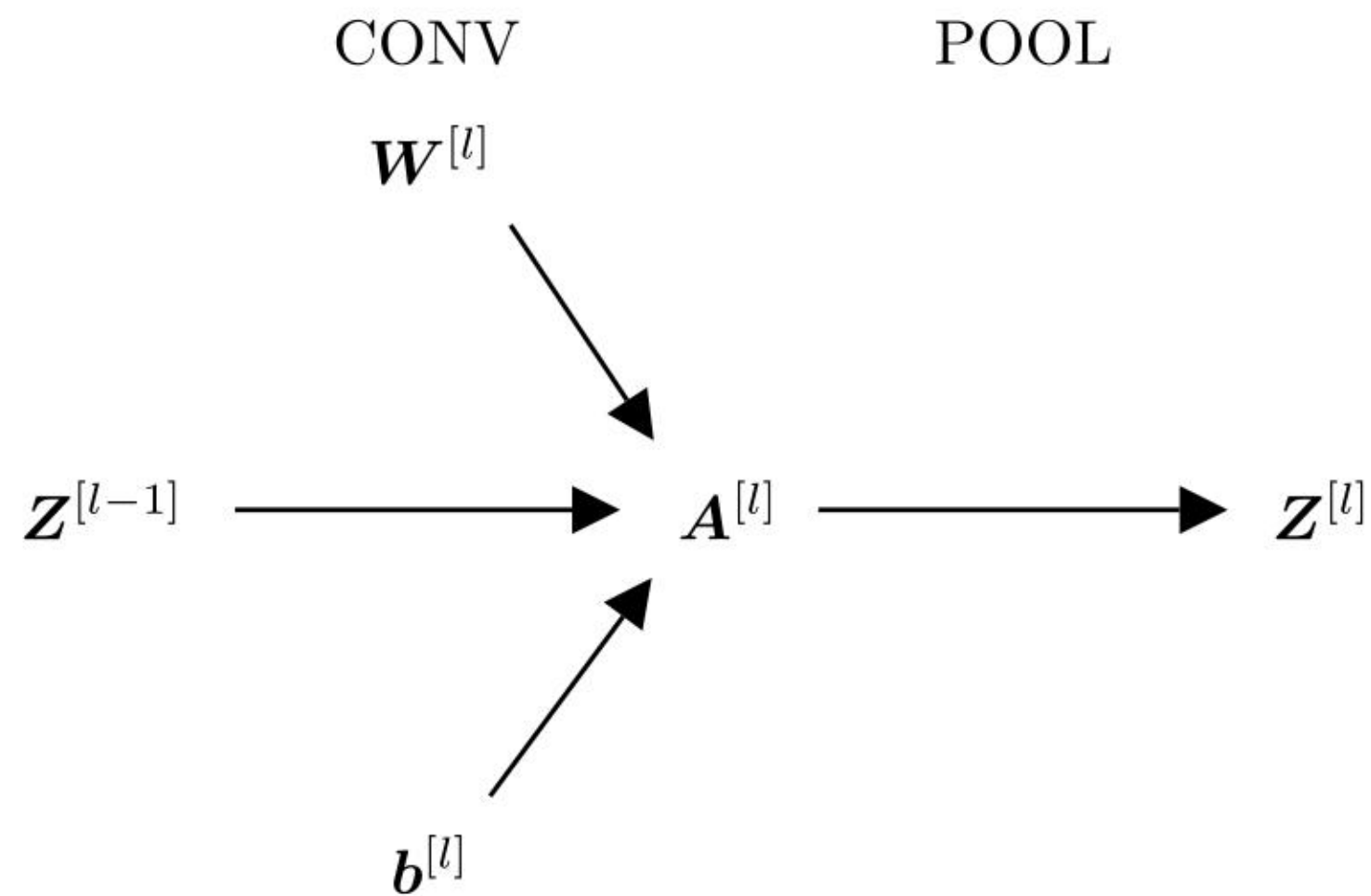
1. We have learnt backpropagation for fully connected layers
2. Thus, it remains to show the backpropagation for
 - The pooling layer (POOL)
 - The convolution layer (CONV)
3. In the following, we assume a POOL is conducted right after a CONV

Backpropagation

1. Denote

- $\mathbf{A}^{[l-1]} \in \mathbb{R}^{? \times ? \times d_C^{[l-1]}}$: output of the $(l-1)$ th CONV
- $\mathbf{Z}^{[l-1]} \in \mathbb{R}^{d_H^{[l-1]} \times d_W^{[l-1]} \times d_C^{[l-1]}}$: output of the $(l-1)$ th POOL, taking $\mathbf{A}^{[l-1]}$ as input
- $\mathbf{W}^{[l]} \in \mathbb{R}^{d_C^{[l]} \times f^{[l]} \times f^{[l]} \times d_C^{[l-1]}}$: kernel for the l th layer, taking $\mathbf{Z}^{[l-1]}$ as input
- $\mathbf{b}^{[l]} \in \mathbb{R}^{d_C^{[l]} \times 1 \times 1 \times 1}$: bias for the l th layer, taking $\mathbf{Z}^{[l-1]}$ as input
- $\mathbf{A}^{[l]} \in \mathbb{R}^{? \times ? \times d_C^{[l]}}$: output of the l th CONV
- $\mathbf{Z}^{[l]} \in \mathbb{R}^{d_H^{[l]} \times d_W^{[l]} \times d_C^{[l]}}$: output of the l th POOL

Structure



1. Assume that $d\mathbf{Z}^{[l]}$ is available

Backpropagation for average POOL

$A^{[l]}$

3	6	5	4	8	9
1	7	9	6	8	0
5	0	9	6	2	0
5	2	6	3	7	0
9	0	3	2	3	1
3	1	3	7	1	7

Kernel M

0.25	0.25
0.25	0.25

$Z^{[l]}$

4.25	6.0	6.25
3.0	6.0	2.25
3.25	3.75	3.0

Backpropagation for average POOL

1. For average pooling, we have

$$d\mathbf{A}^{[l]} = d\mathbf{Z}^{[l]} \otimes \mathbf{M}$$

- $\mathbf{A} \otimes \mathbf{B}$: Kronecker production of two matrices \mathbf{A} and \mathbf{B}

Backpropagation for Max POOL

$\mathbf{A}^{[l]}$

3	6	5	4	8	9
1	7	9	6	8	0
5	0	9	6	2	0
5	2	6	3	7	0
9	0	3	2	3	1
3	1	3	7	1	7

Kernel

$\mathbf{Z}^{[l]}$

7	9	9
5	9	7
9	7	7

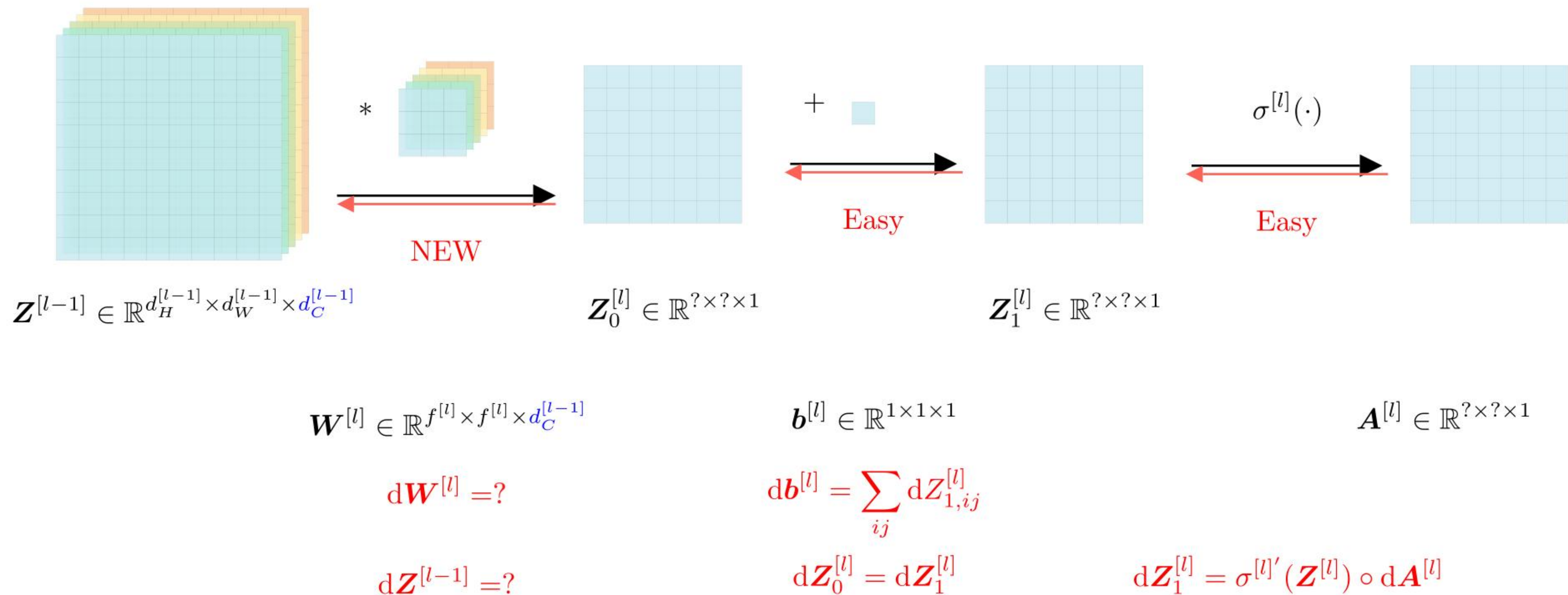
Backpropagation for Max POOL

1. For Max pooling, $d\mathbf{A}^{[l]}$ is obtained by
 - stacking small matrices associated with each elements in $\mathbf{Z}^{[l]}$ times a 2×2 mask matrix
 - Those 2×2 matrices consist of 0s and only one 1
2. We have finished the backpropagation from $d\mathbf{Z}^{[l]}$ to $d\mathbf{A}^{[l]}$
 - POOL does not involve any model parameters
3. Assume the availability of $d\mathbf{A}^{[l]}$ for the following analysis

Backpropagation for CONV

1. For simplicity, let $d_C^{[l]} = 1$.
2. Assume $d\mathbf{A}^{[l]}$ to be available

Backpropagation for CONV



Backpropagation for CONV

1. To find derivative with respect to something, we need to find where its information is contained
2. Initialize $d\mathbf{Z}^{[l-1]}$ as a zero matrix of the same dimension as $\mathbf{Z}^{[l-1]}$
3. Informally, obtain

$$d\mathbf{Z}_{\text{slice},ij}^{[l-1]} + = dZ_{0,ij}^{[l]} \times \mathbf{W}^{[l]}$$

- $\mathbf{Z}_{\text{slice},ij}^{[l-1]}$: the part of $\mathbf{Z}^{[l-1]}$ used to obtain $Z_{0,ij}^{[l]}$

Backpropagation for CONV

1. It is obvious that every element in $\mathbf{Z}_0^{[l]}$ contains information of $\mathbf{W}^{[l]}$
2. Thus, we have

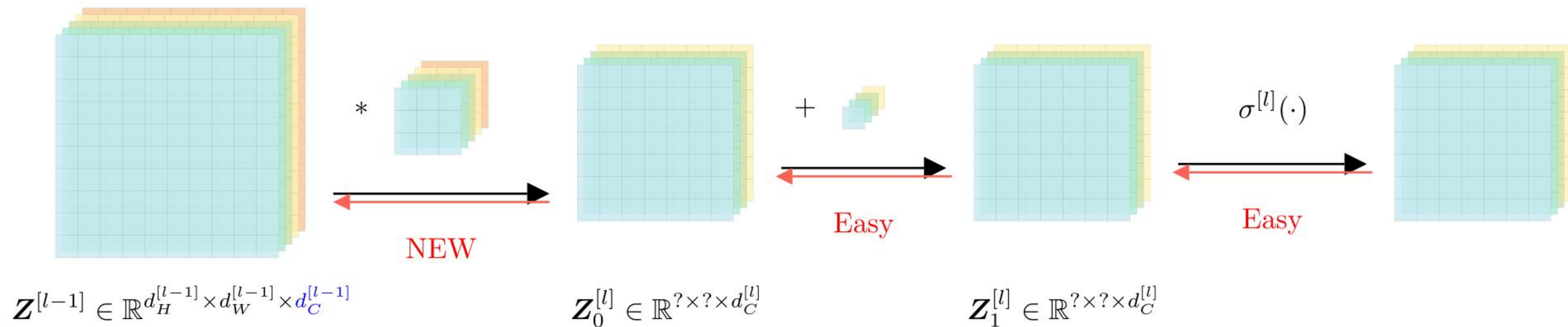
$$d\mathbf{W}^{[l]} = \sum_{ij} dZ_{0,ij}^{[l]} \times \mathbf{B}_{ij}^{[l-1]}$$

- $\mathbf{B}_{ij}^{[l-1]}$: the part of $\mathbf{Z}^{[l-1]}$ used to obtain $Z_{0,ij}^{[l]}$

Backpropagation for CONV

1. Consider the general case where there exist $d_C^{[l]}$ channels in the l th layer.

Backpropagation for CONV



$$\mathbf{W}^{[l]} \in \mathbb{R}^{d_C^{[l]} \times f^{[l]} \times f^{[l]} \times d_C^{[l-1]}}$$

$$\mathbf{b}^{[l]} \in \mathbb{R}^{d_C^{[l]} \times 1 \times 1 \times 1}$$

$$\mathbf{A}^{[l]} \in \mathbb{R}^{? \times ? \times d_C^{[l]}}$$

$$d\mathbf{W}^{[l]} = ?$$

$$d\mathbf{b}_c^{[l]} = \sum_{ij} dZ_{1,ijc}^{[l]} \quad (c = 1, \dots, d_C^{[l]})$$

$$d\mathbf{Z}^{[l-1]} = ?$$

$$d\mathbf{Z}_0^{[l]} = d\mathbf{Z}_1^{[l]}$$

$$d\mathbf{Z}_1^{[l]} = \sigma^{[l]'}(\mathbf{Z}^{[l]}) \circ d\mathbf{A}^{[l]}$$

Backpropagation for CONV

1. Just think about where the information is contained when calculating derivatives.
2. The information about $\mathbf{Z}^{[l-1]}$ is contained in every channel of $\mathbf{Z}_0^{[l]}$
3. Initialize $d\mathbf{Z}^{[l-1]}$ as a zero matrix of the same dimension as $\mathbf{Z}^{[l-1]}$
4. Informally, obtain

$$d\mathbf{Z}_{\text{slice},ij}^{[l-1]} + = \sum_{c=1}^{d_C^{[l]}} dZ_{0,ijc}^{[l]} \times \mathbf{W}_c^{[l]}$$

- $\mathbf{Z}_{\text{slice},ij}^{[l-1]}$: the part of $\mathbf{Z}^{[l-1]}$ used to obtain $Z_{0,ijc}^{[l]}$ for $c = 1, \dots, d_C^{[l]}$
- $\mathbf{W}_c^{[l]}$: the c th kernel in the l th layer

Backpropagation for CONV

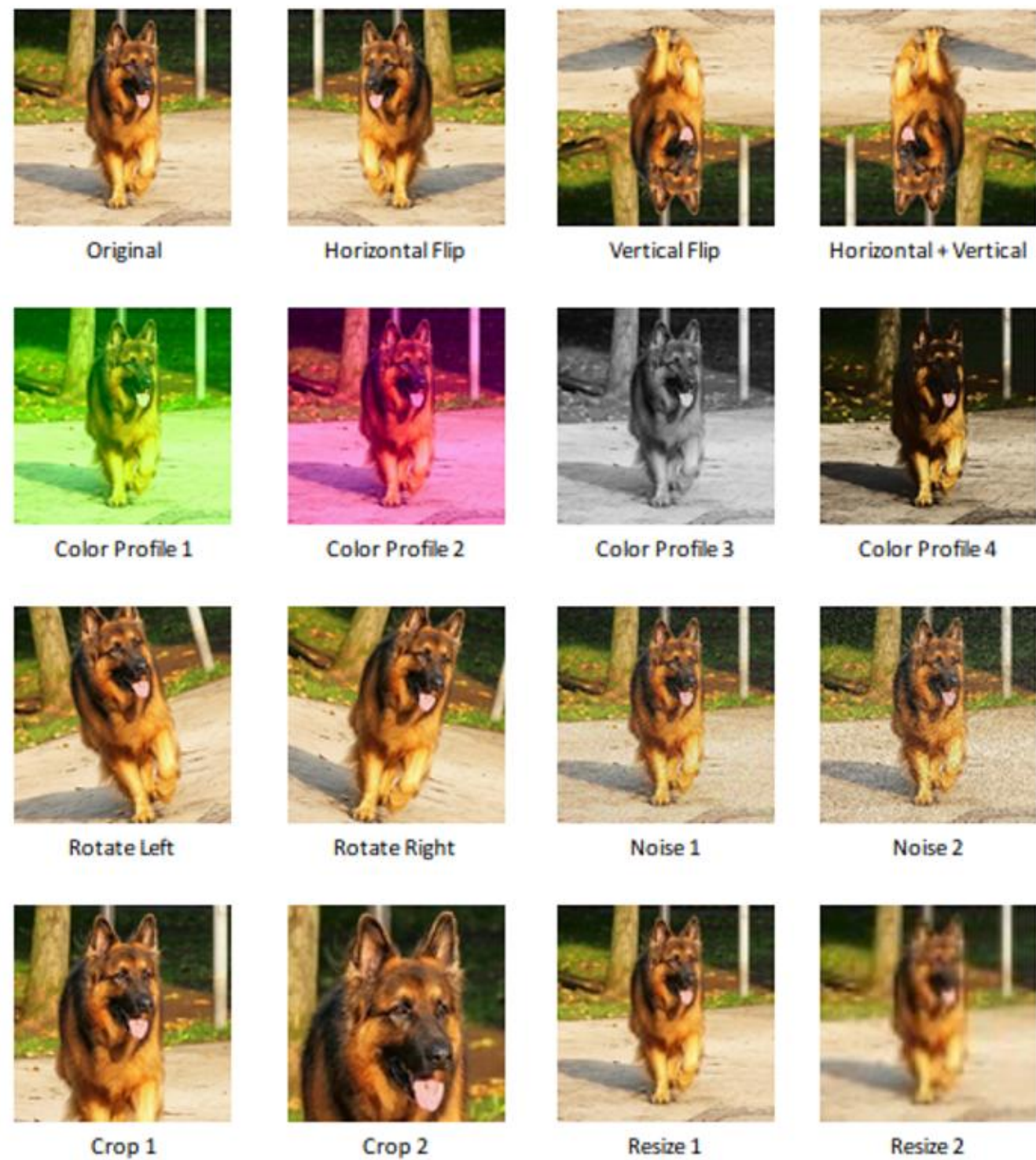
1. It is obvious that every element in $\mathbf{Z}_{0,c}^{[l]}$ contains information of $\mathbf{W}_c^{[l]}$
for $c = 1, \dots, d_C^{[l]}$

2. Thus, we have

$$d\mathbf{W}_{\mathbf{c}}^{[l]} = \sum_{ij} dZ_{0,ij\mathbf{c}}^{[l]} \times \mathbf{B}_{ij}^{[l-1]} \quad (c = 1, \dots, d_C^{[l]})$$

- $\mathbf{B}_{ij}^{[l-1]}$: the part of $\mathbf{Z}^{[l-1]}$ used to obtain $Z_{0,ij\mathbf{c}}^{[l]}$

Data augmentation



[<https://www.volansys.com/blog/data-augmentation-in-ml/>]

Understanding convolution networks

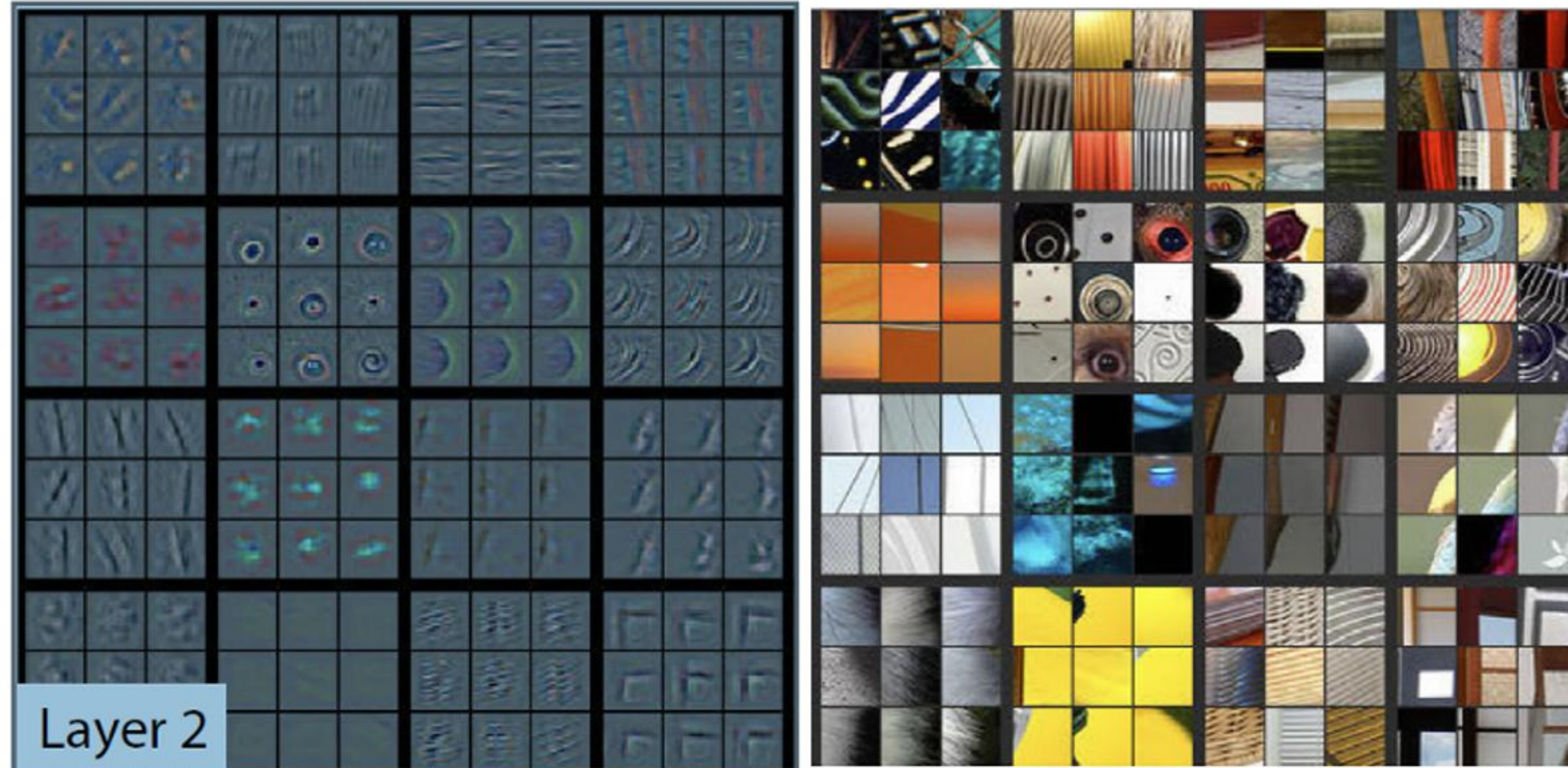


Layer 1



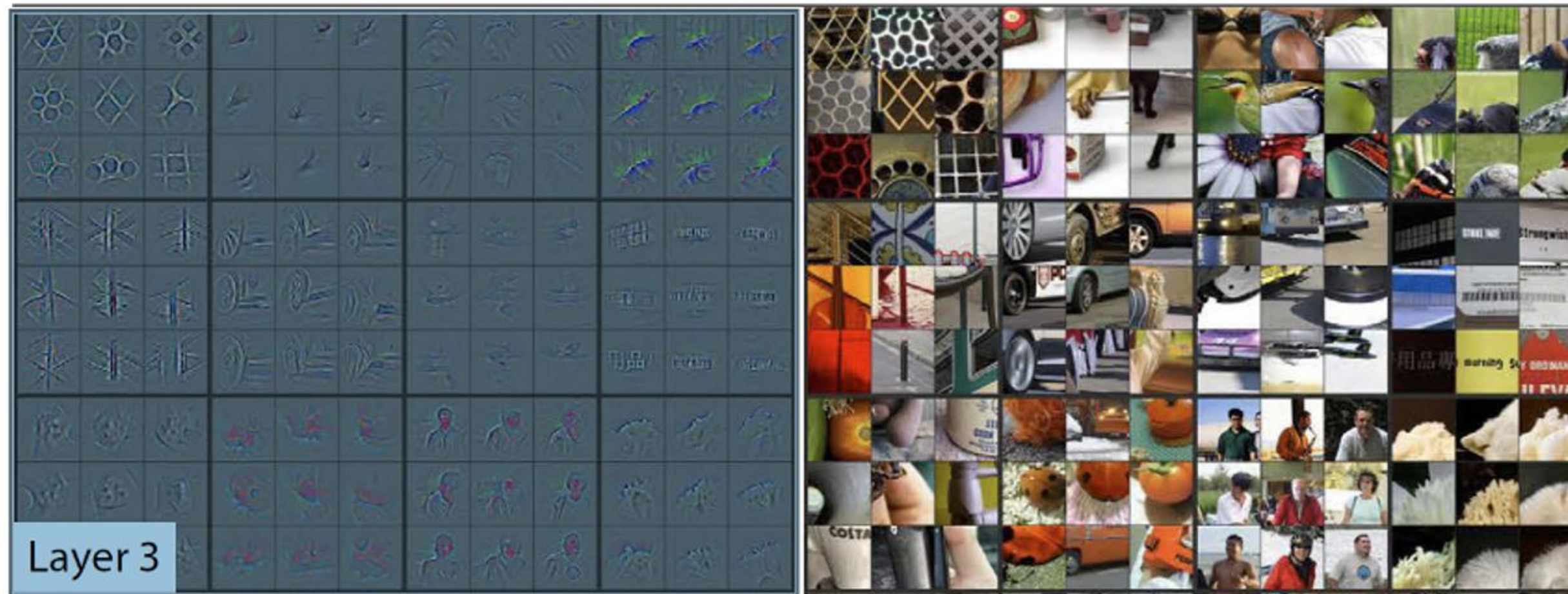
[Zeiler, M. D., & Fergus, R. (2014, September). Visualizing and understanding convolutional networks. In European conference on computer vision (pp. 818-833). Springer, Cham]

Understanding convolution networks



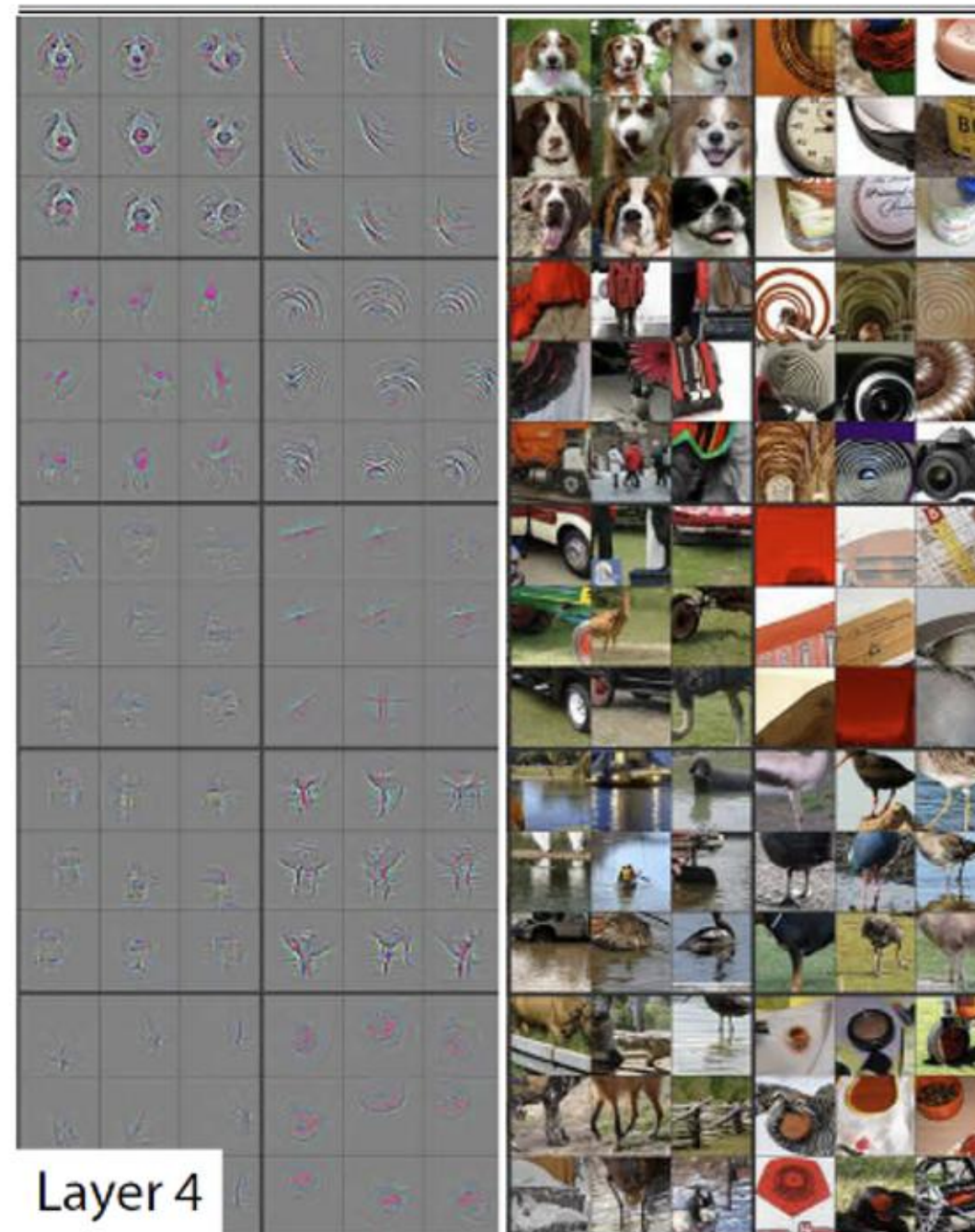
[Zeiler, M. D., & Fergus, R. (2014, September). Visualizing and understanding convolutional networks. In European conference on computer vision (pp. 818-833). Springer, Cham]

Understanding convolution networks



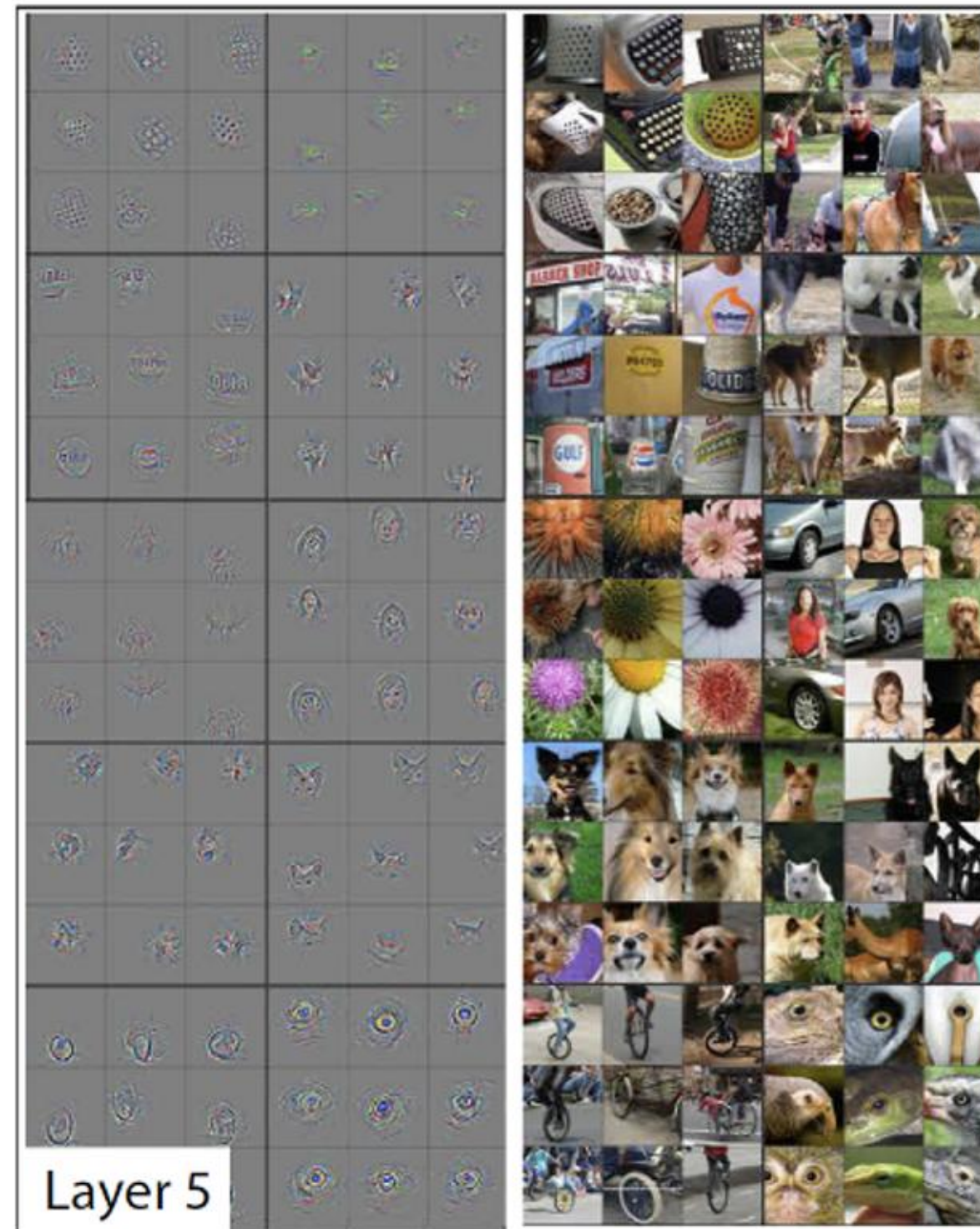
[Zeiler, M. D., & Fergus, R. (2014, September). Visualizing and understanding convolutional networks. In European conference on computer vision (pp. 818-833). Springer, Cham]

Understanding convolution networks



[Zeiler, M. D., & Fergus, R. (2014, September). Visualizing and understanding convolutional networks. In European conference on computer vision (pp. 818-833). Springer, Cham]

Understanding convolution networks



[Zeiler, M. D., & Fergus, R. (2014, September). Visualizing and understanding convolutional networks. In European conference on computer vision (pp. 818-833). Springer, Cham]