

CSDS 451: Designing High Performant Systems for AI

Lecture 9

9/23/2025

Sanmukh Kuppannagari

sanmukh.kuppannagari@case.edu

<https://sanmukh.research.st/>

Case Western Reserve University

Outline

- Accelerating Convolutional Neural Network Models: Basics

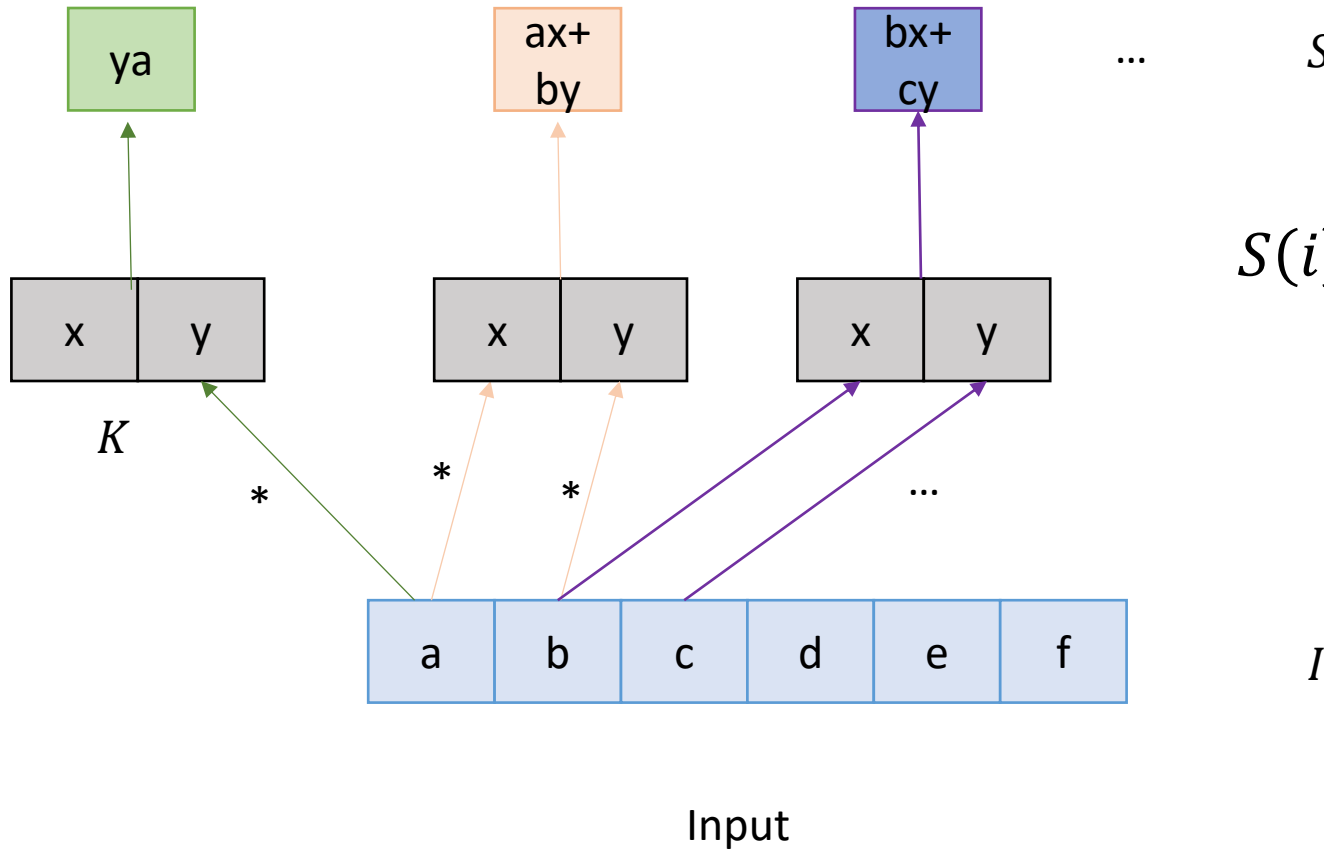
Reading Materials

- Deep Learning, by Ian Goodfellow, Yoshua Bengio, Aaron Courville
 - Available Online: <https://www.deeplearningbook.org/>
 - CNNs: Chapter 9
- A Guide to Convolution Arithmetic for Deep Learning
 - Available Online: <https://arxiv.org/pdf/1603.07285.pdf>

Convolutional Neural Networks

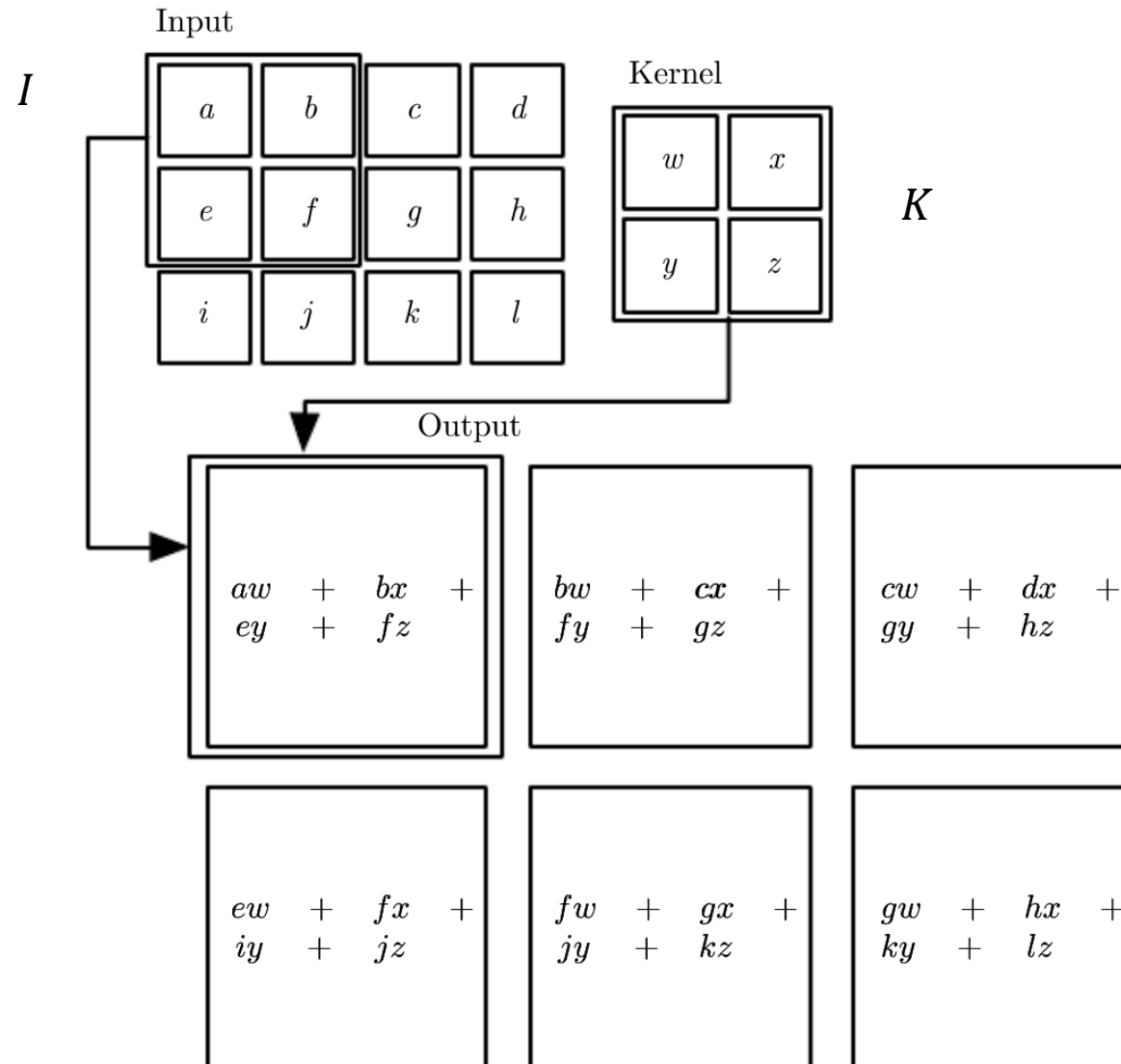
- Neural Network specialized for a grid-like topology
 - 1D: timeseries data, e.g., (12 pm, 63), (1pm, 65), (2 pm, 68),...
 - 2D: Images
 - 3D: Images with multiple channels, e.g., RGB images
- Name derives from “Convolution” operations performed in the “Convolution” Layers

1D Convolution



$$S(i) = \sum_m I(i + m) * K(m)$$

2D Convolution



$$S(i, j) = \sum_m \sum_n I(i + m, j + n) K(m, n)$$

S

Source: Deep Learning Book

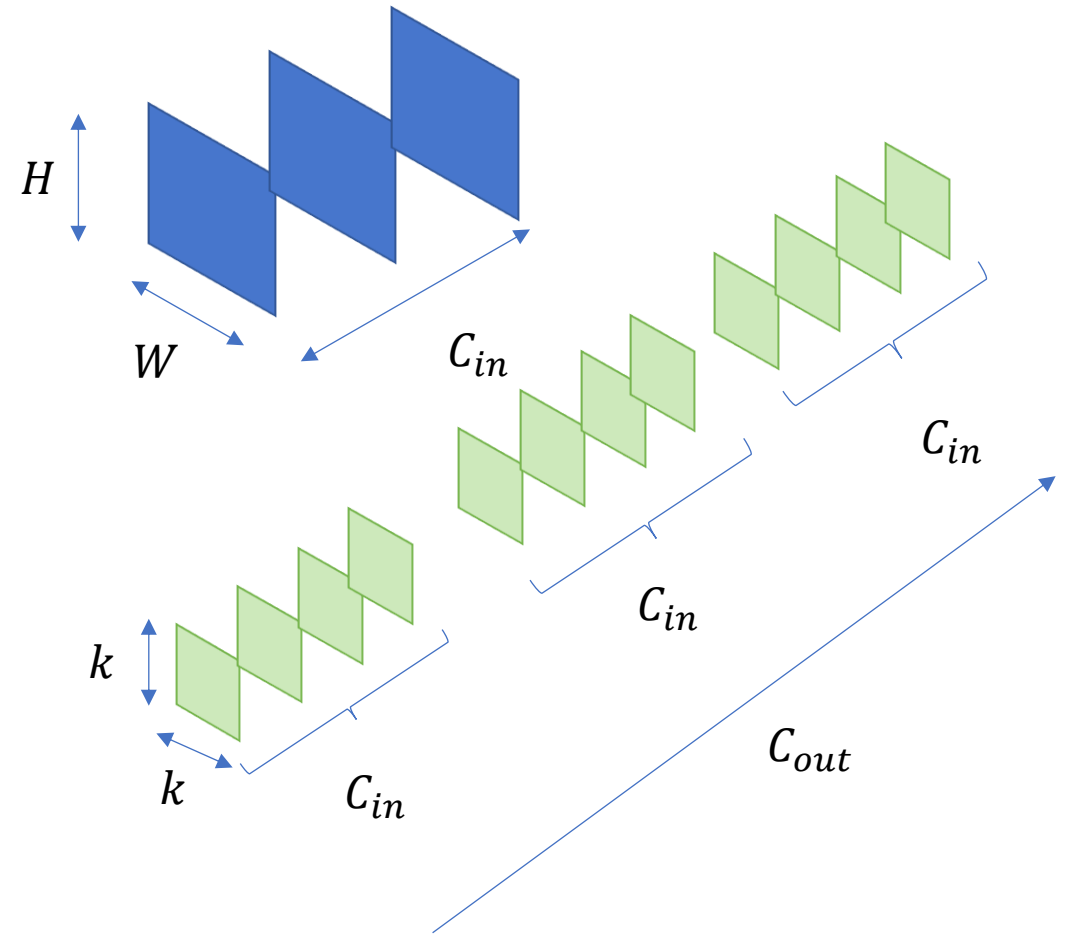
What happens for corner cases???

Convolutional Neural Network

- A Neural Network with:
 - Convolution Layers
 - Activations
 - Pooling Layers
 - ...

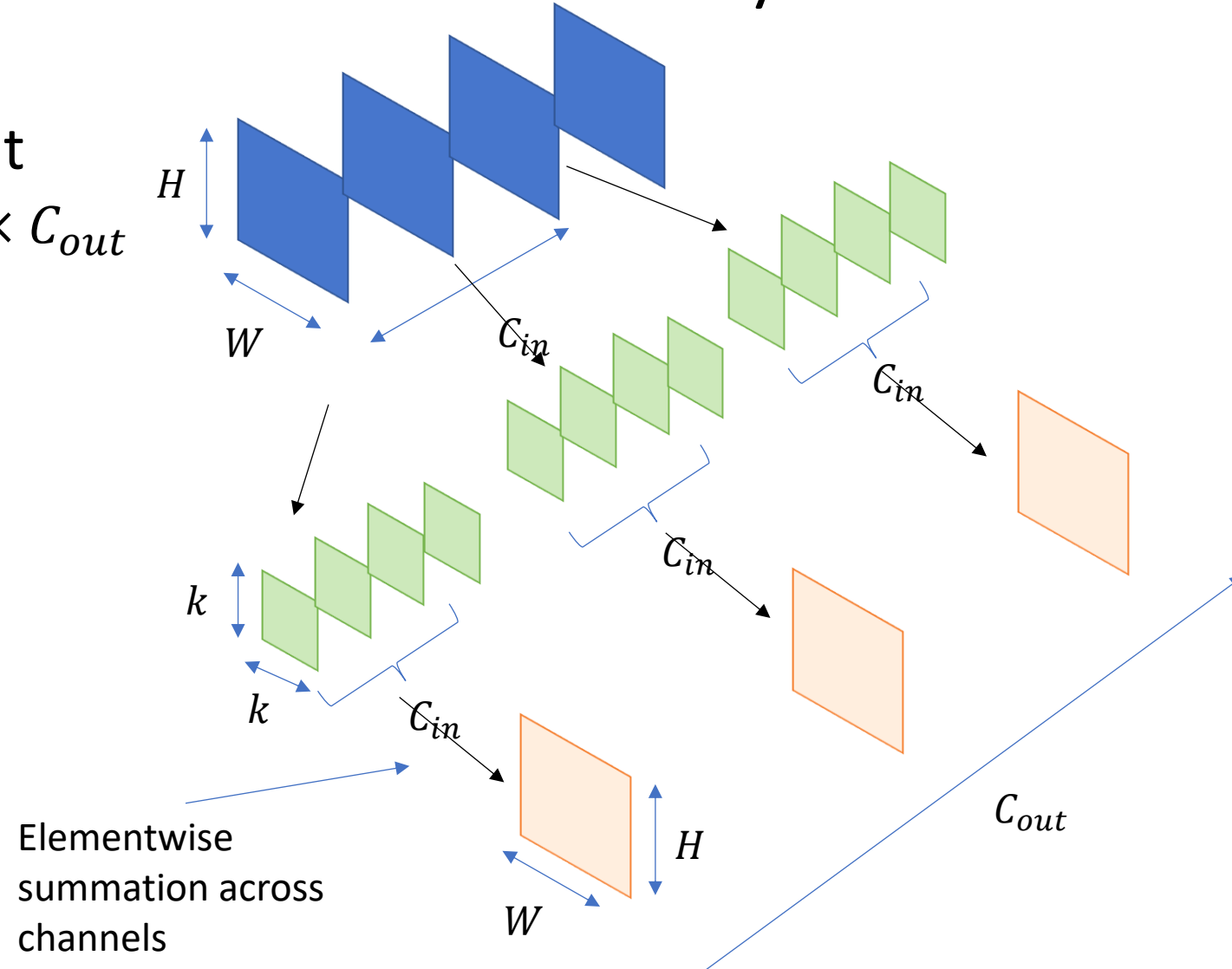
Convolution Layer

- Given an Input Feature Map
 - Size: $H \times W \times C_{in}$
 - H : Height of the input feature map
 - W : Width of the input feature map
 - C_{in} : Number of Channels in the feature map
- $C_{in} \times C_{out}$ Kernels, each of size $k \times k$



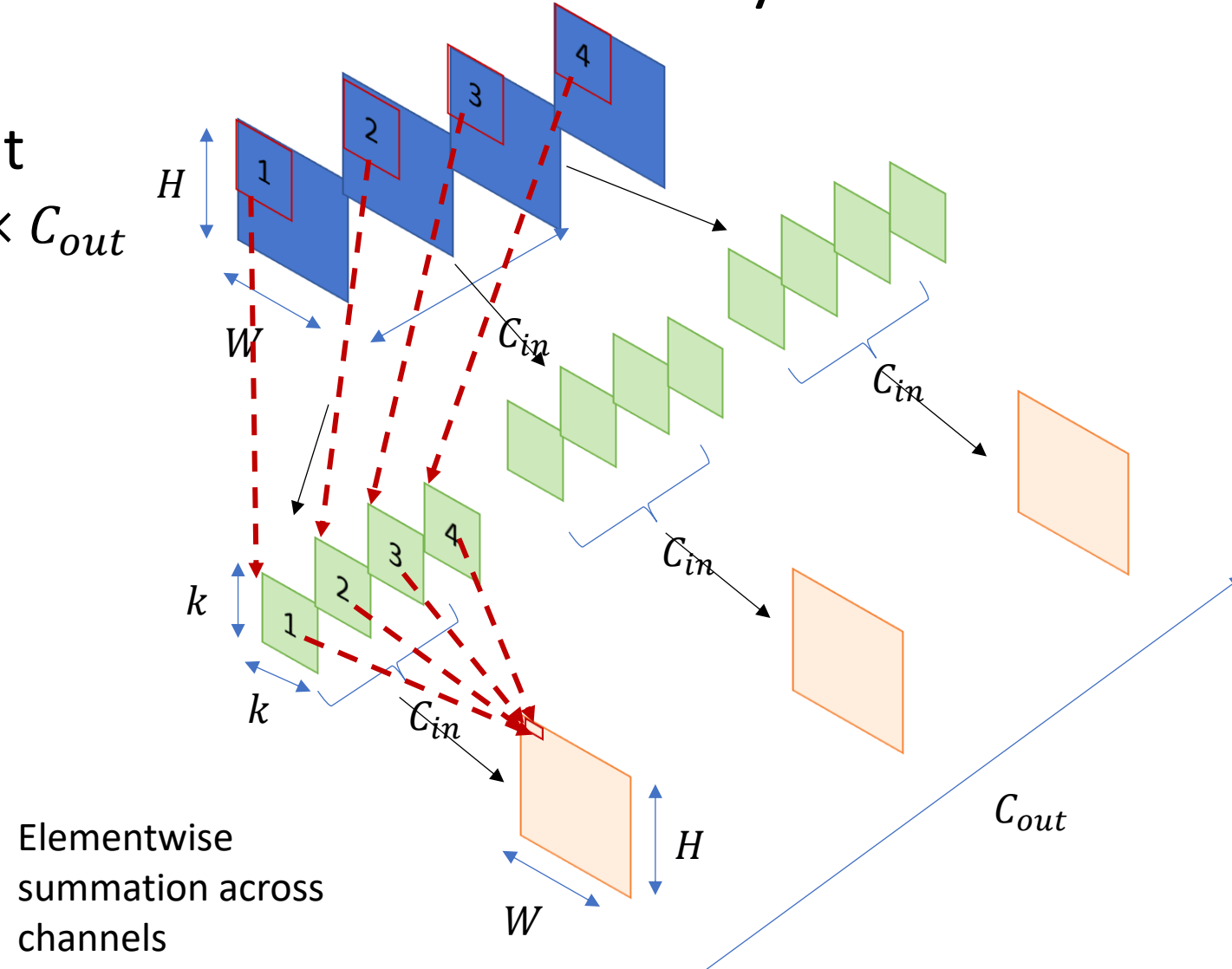
Convolution Layer

- Produce Output
 - Size: $H \times W \times C_{out}$



Convolution Layer

- Produce Output
 - Size: $H \times W \times C_{out}$



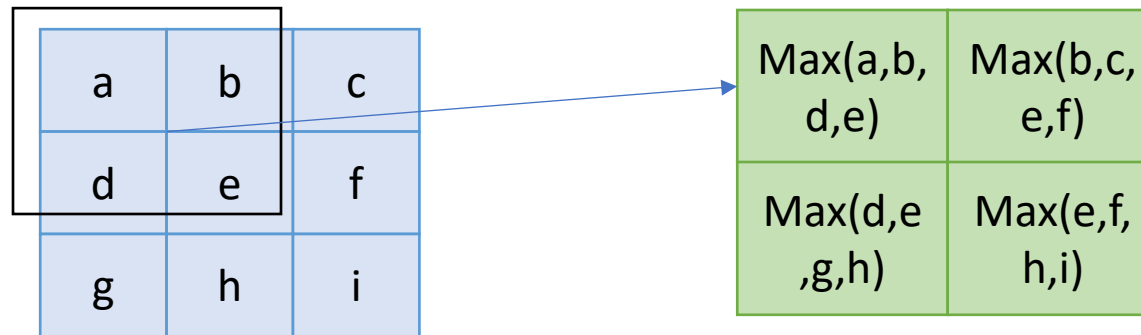
Activation Layer



σ : Activation Function, Sigmoid, ReLU, Softmax

Pooling Layer

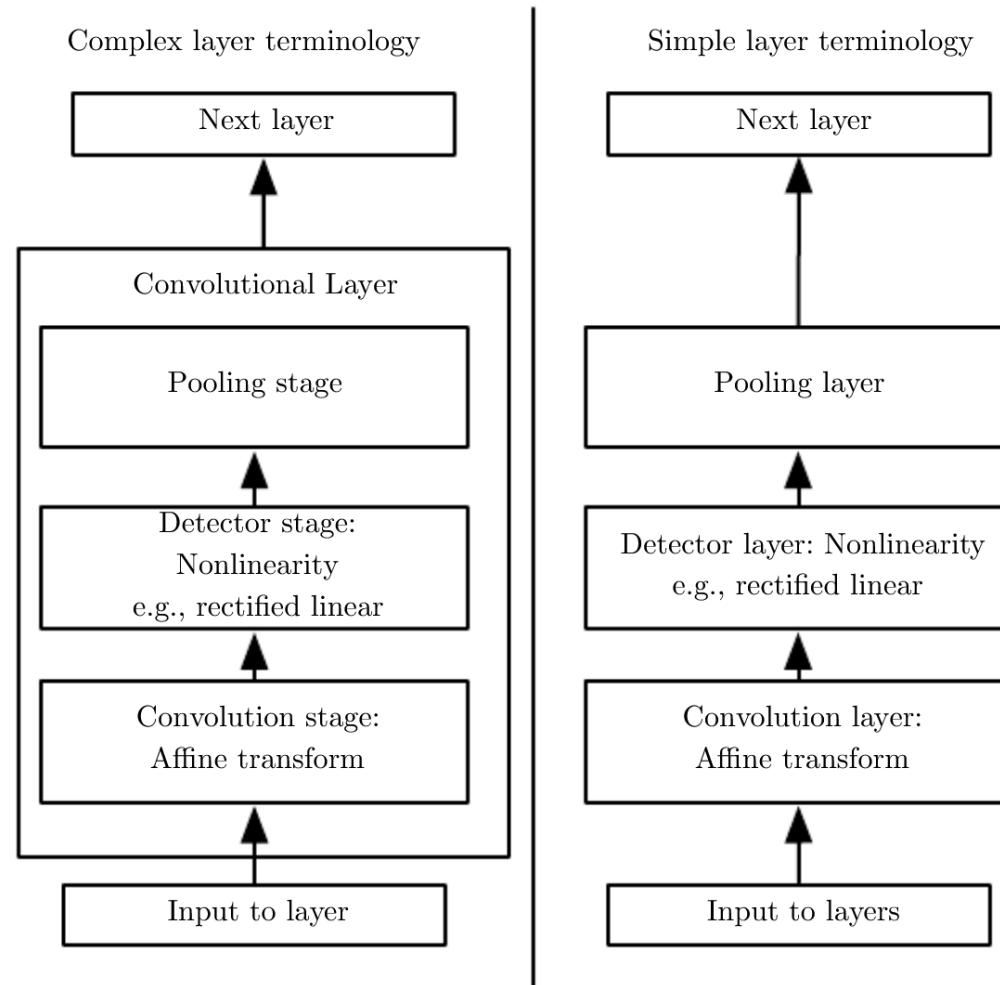
- Replaces the output of a grid with a summary of the grid



Max Pooling: Size
2 X 2

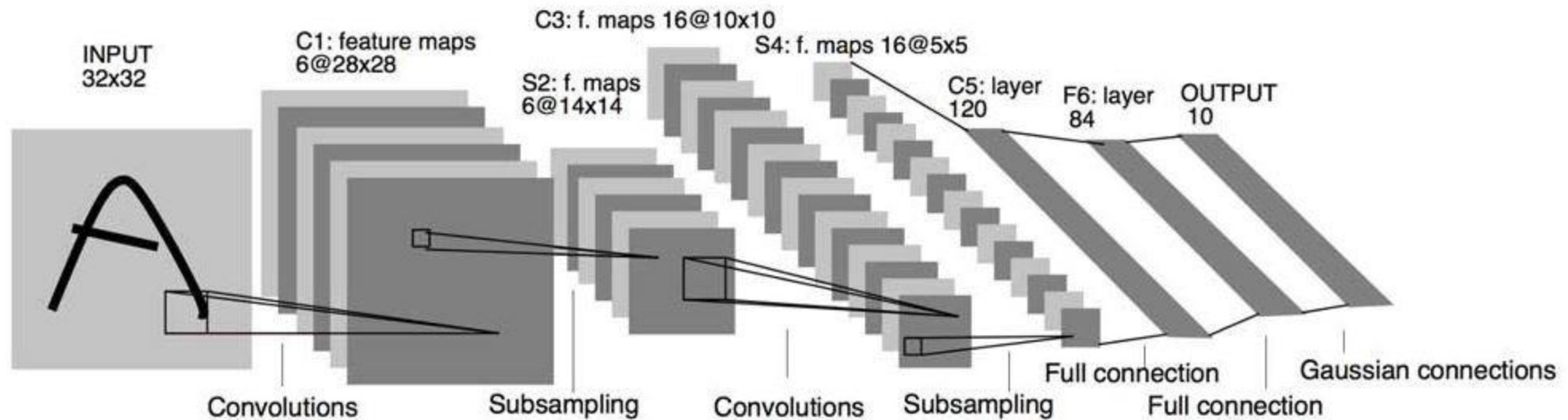
Pooling Types:
Average Pooling
Max Pooling

Convolution Neural Network



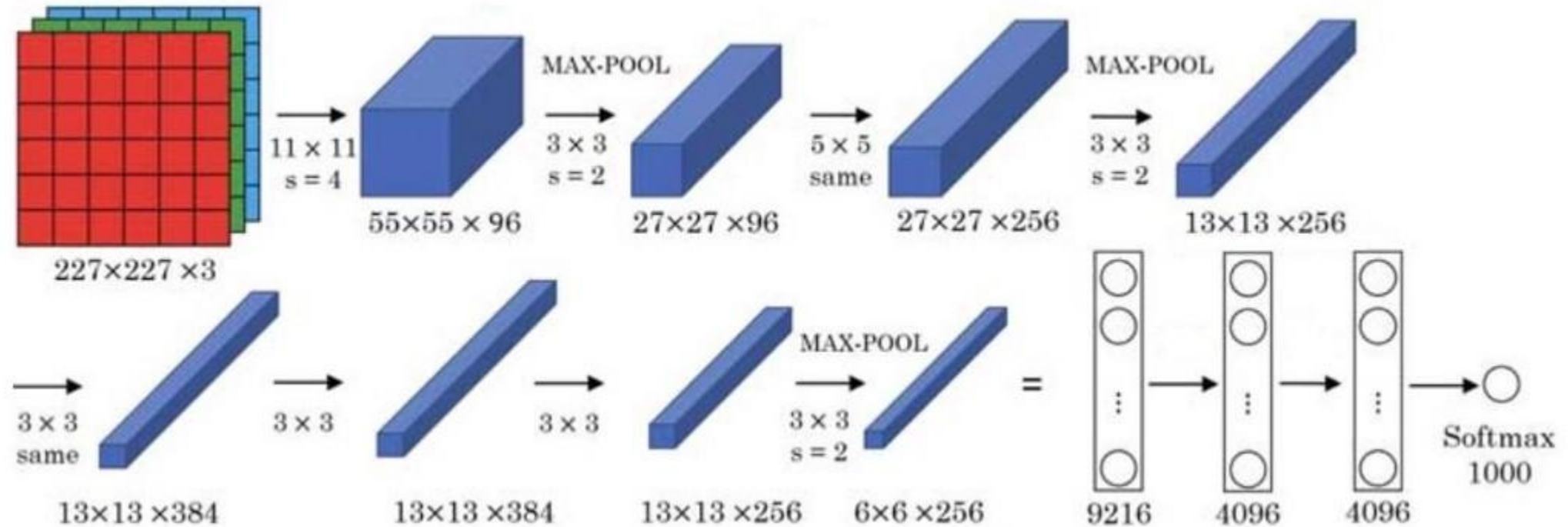
Source: Deep Learning Book

Example CNN Networks



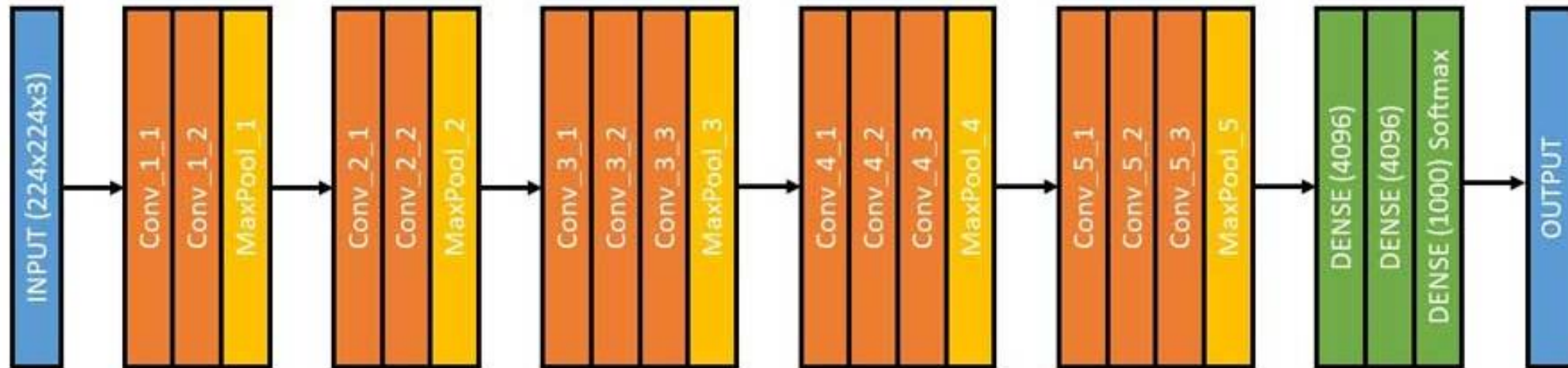
LeNet (1998): Gradient Based Learning Applied to Document Recognitions,
<http://yann.lecun.com/exdb/publis/pdf/lecun-01a.pdf>

Example CNN Networks



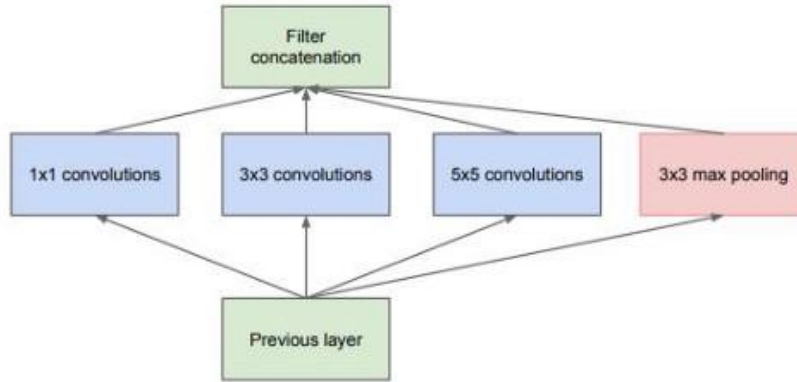
AlexNet (2012): ImageNet Classification with Deep Convolutional Neural Networks,
<https://proceedings.neurips.cc/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf>

Example CNN Networks

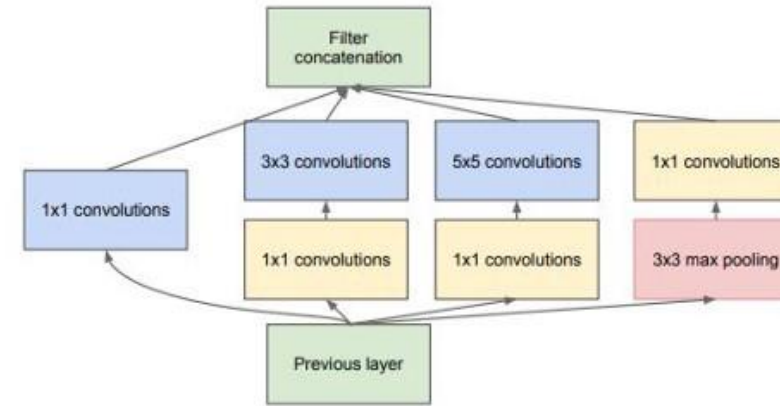


VGG16 (2014): VERY DEEP CONVOLUTIONAL NETWORKS
FOR LARGE-SCALE IMAGE RECOGNITION,
<https://arxiv.org/pdf/1409.1556.pdf>

Example CNN Networks

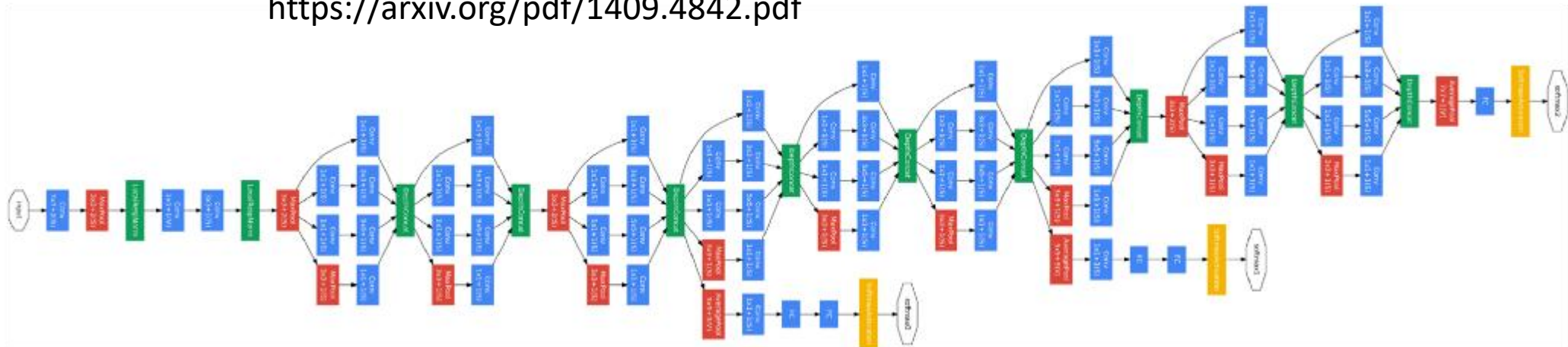


(a) Inception module, naïve version

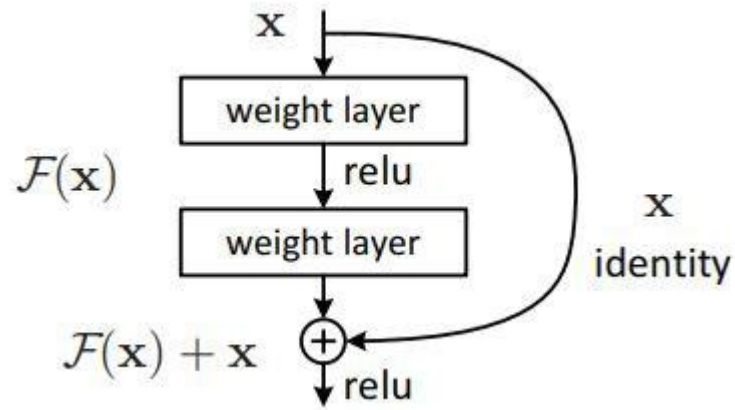


(b) Inception module with dimension reductions

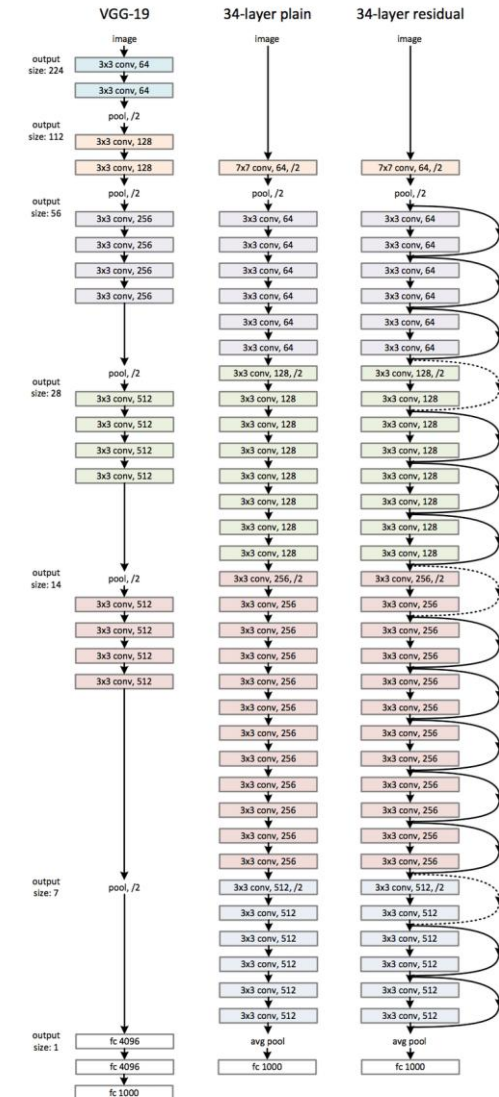
GoogLeNet/InceptionNet (2014): Going deeper with convolutions,
<https://arxiv.org/pdf/1409.4842.pdf>



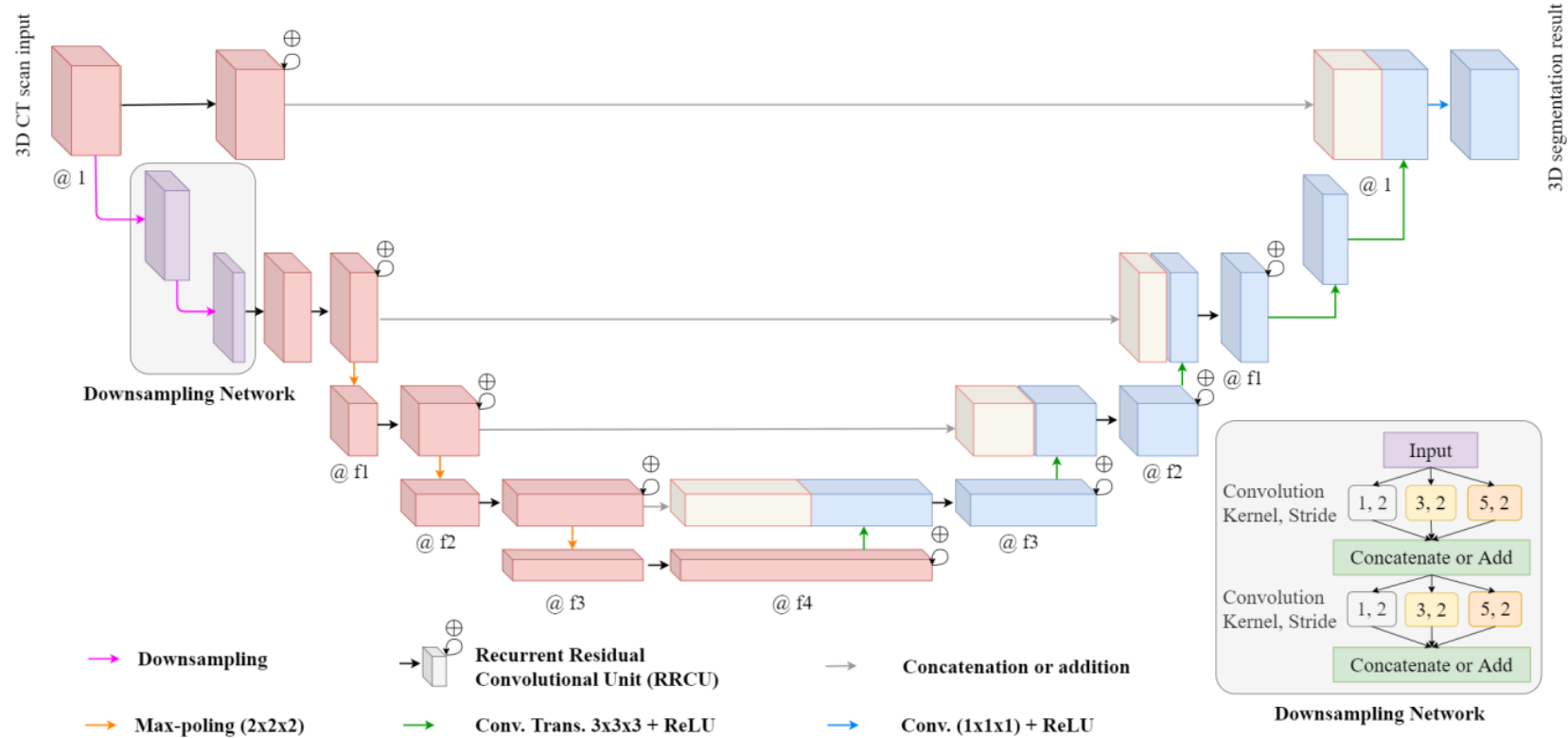
Example CNN Networks



ResNet (2015): Deep Residual Learning for Image Recognition,
<https://arxiv.org/pdf/1512.03385.pdf>



Example CNN Networks



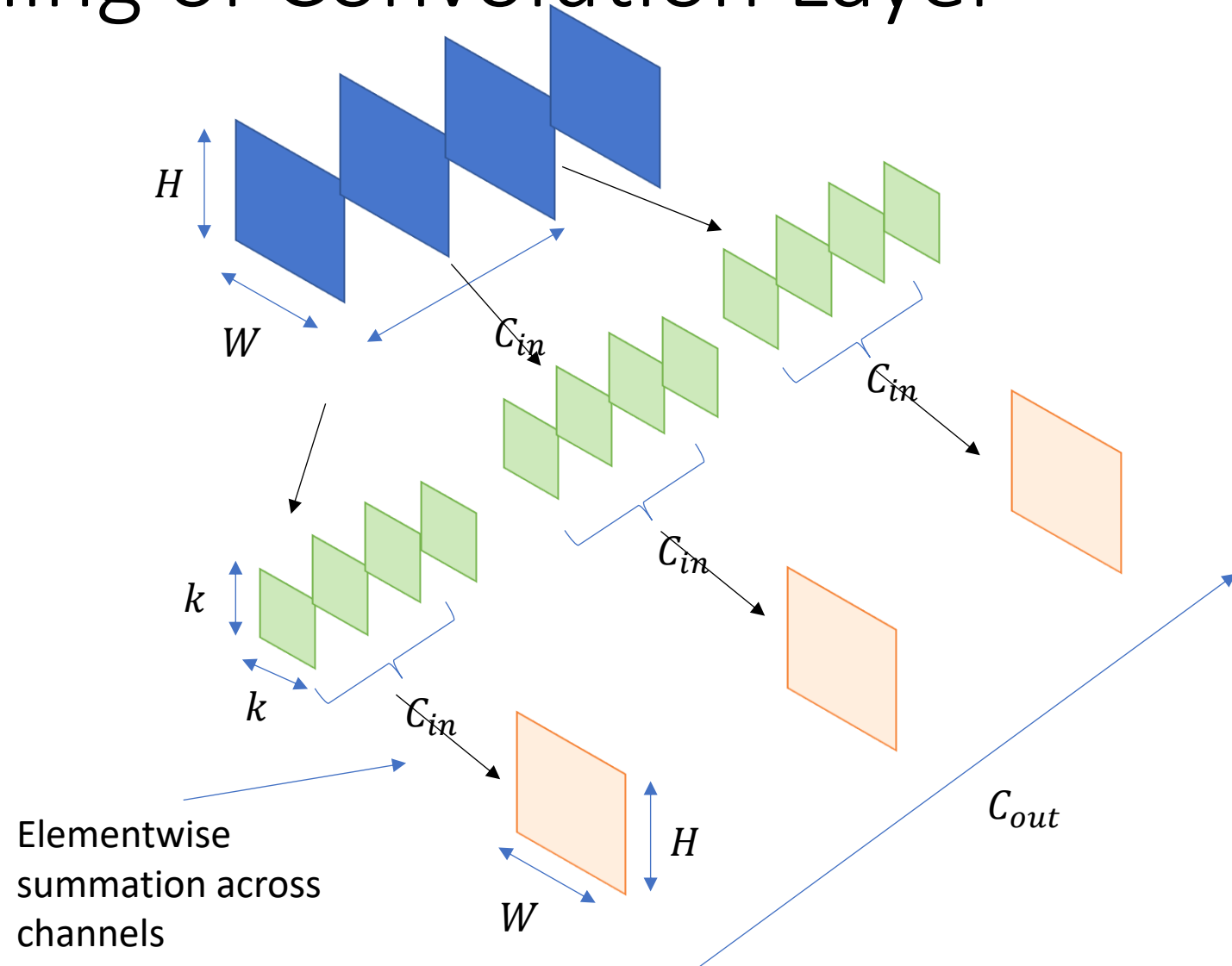
U-net: State of the art in image segmentation, especially in biomedical domain

Example CNN Networks

- EfficientNet (2019): Tan, M., & Le, Q. V. (2020). EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks. *ArXiv:1905.11946 [Cs, Stat]*. <http://arxiv.org/abs/1905.11946>
- FixEfficientNets (2022): Touvron, H., Vedaldi, A., Douze, M., & Jégou, H. (2020b). Fixing the train-test resolution discrepancy: FixEfficientNet. *ArXiv:2003.08237 [Cs]*. <http://arxiv.org/abs/2003.08237>
- AmoebaNets (2019): Regularized Evolution for Image Classifier Architecture Search, <https://ojs.aaai.org/index.php/AAAI/article/view/4405>

Performance Modeling of Convolution Layer

- Number of Operations: ??

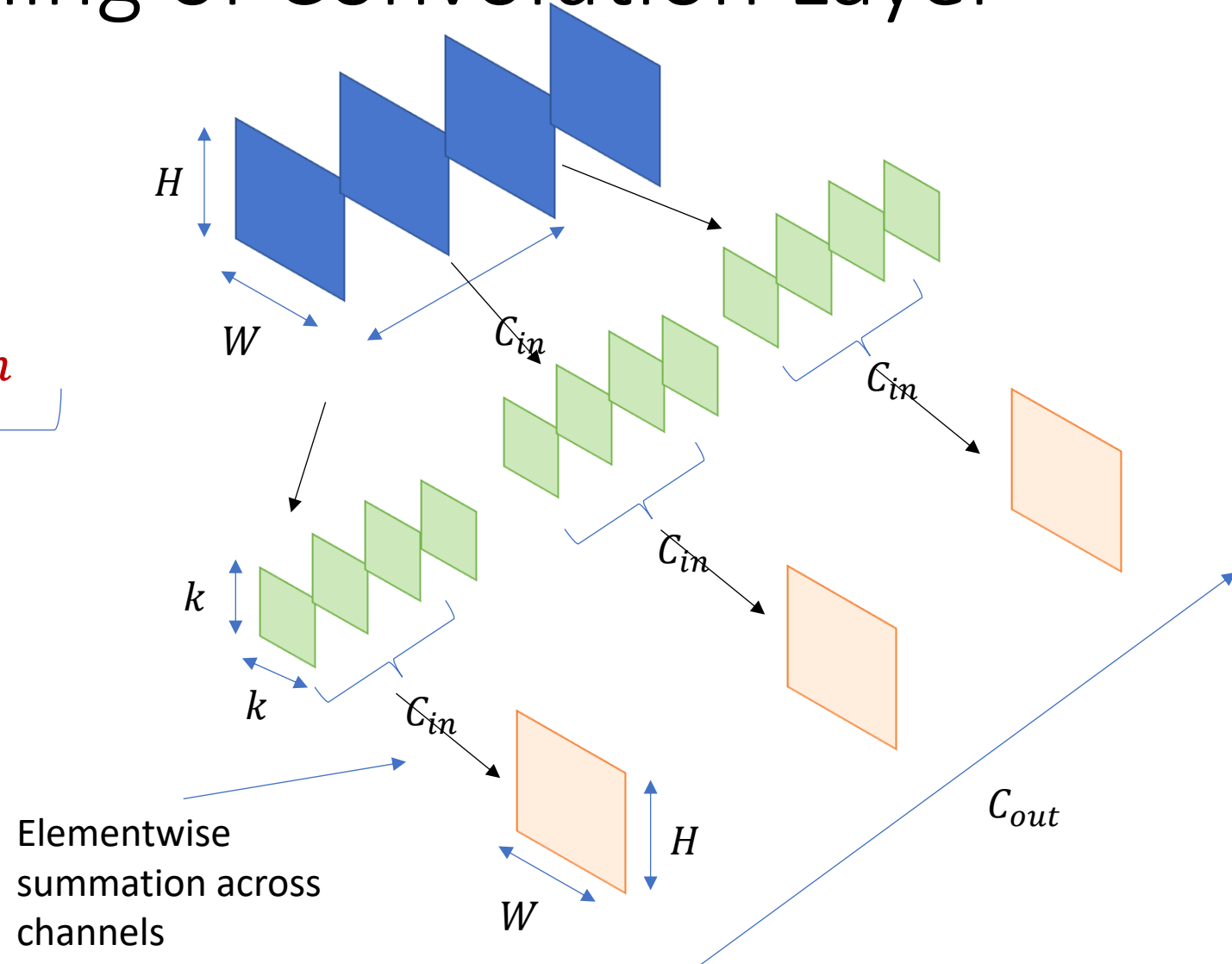


Performance Modeling of Convolution Layer

- Number of Operations:

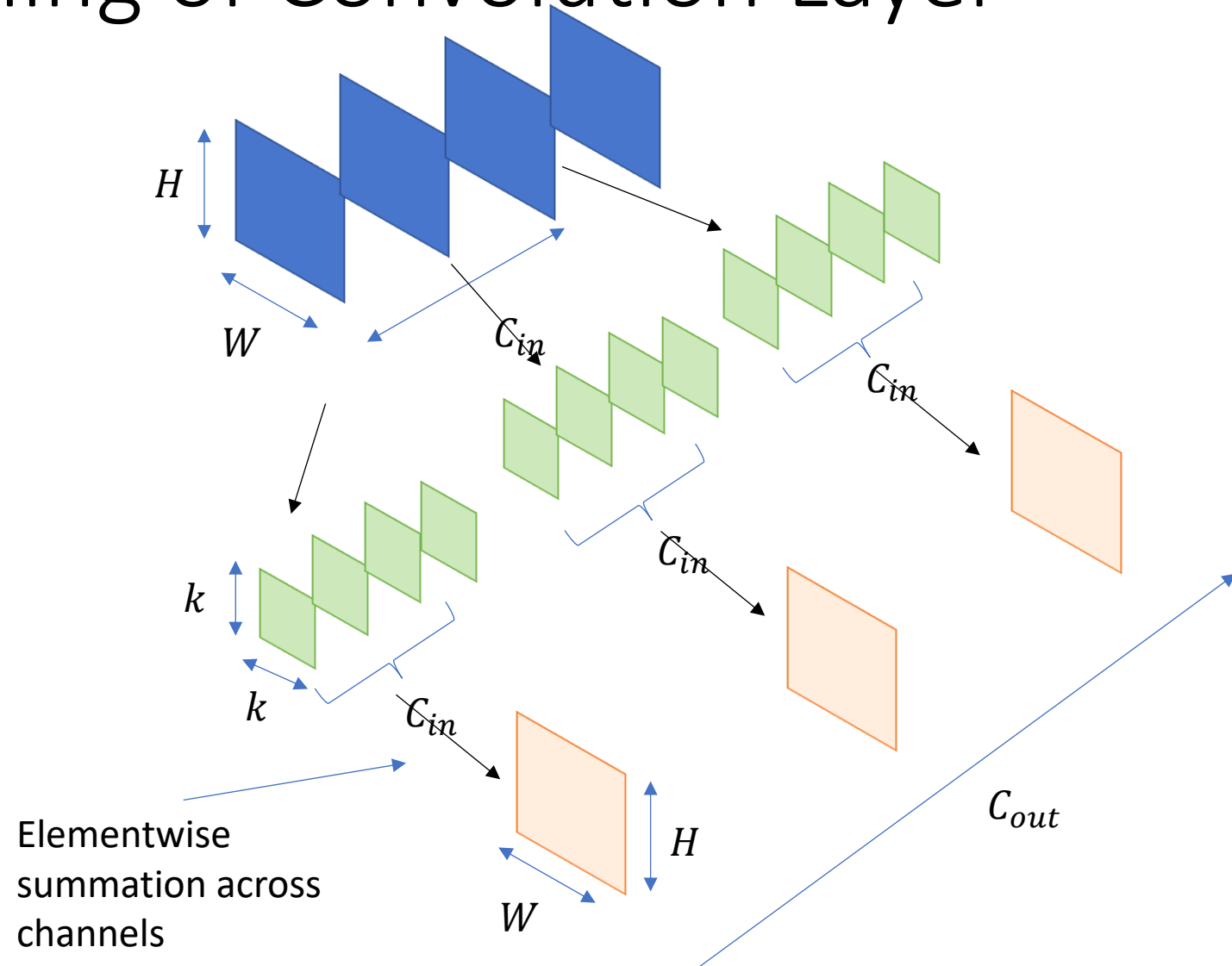
$$H \times W \times C_{out} \times k \times k \times C_{in}$$

Output Pixels Operations per Pixels



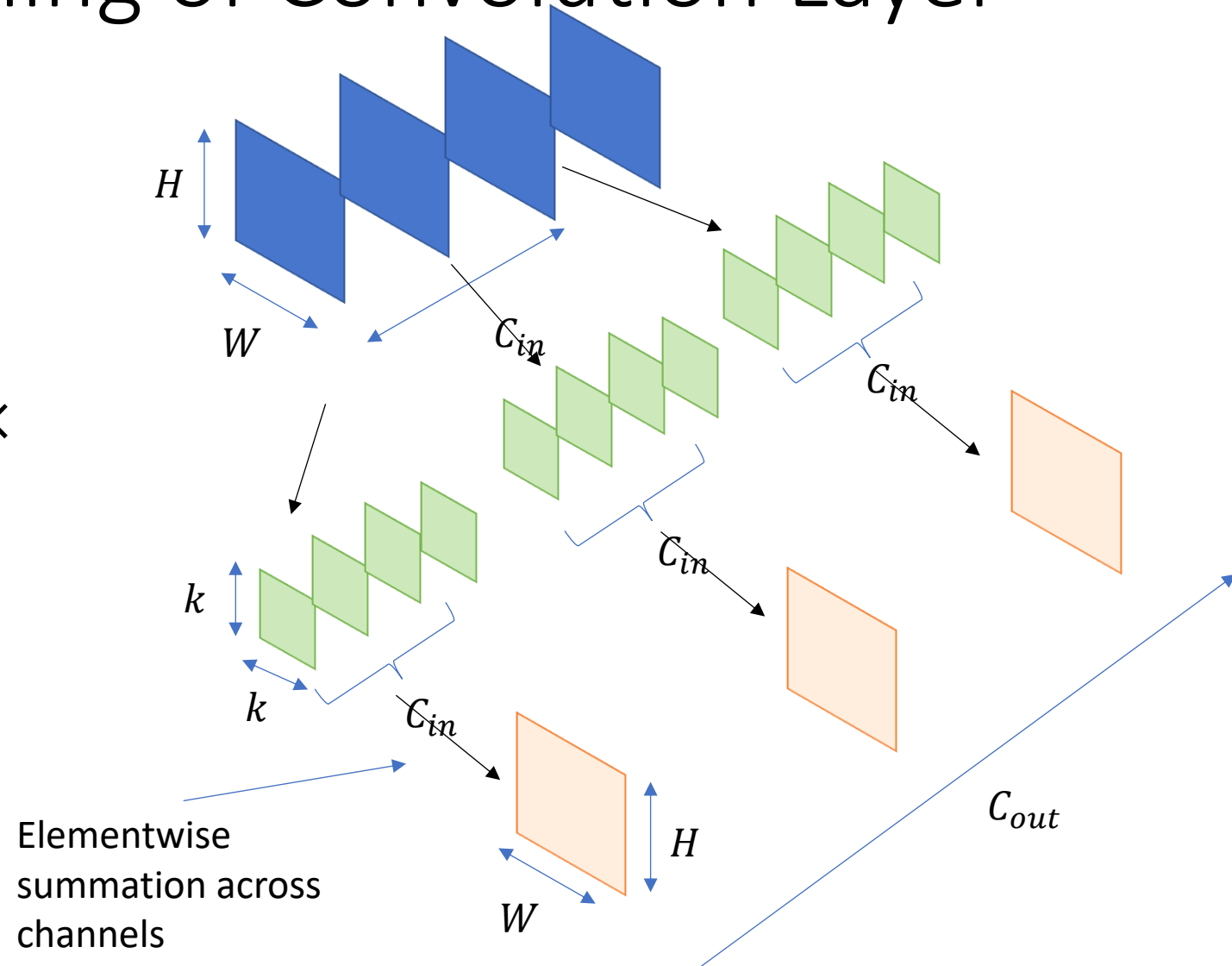
Performance Modeling of Convolution Layer

- Memory Transfers (infinite cache): ??



Performance Modeling of Convolution Layer

- Memory Transfers (infinite cache): ??
- $H \times W \times C_{in} + C_{in} \times C_{out} \times k \times k + H \times W \times C_{out}$
- $\approx H \times W \times (C_{in} + C_{out})$



Performance Modeling of Convolution Layer

- Theoretical Arithmetic Intensity

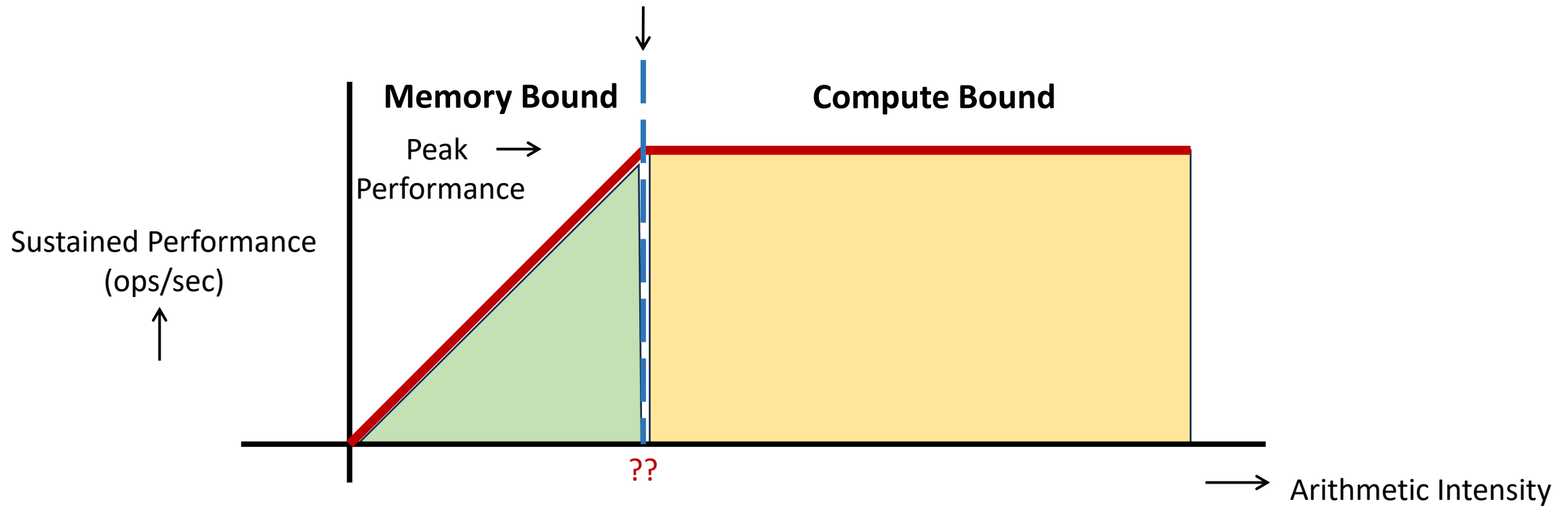
$$\frac{k^2 C_{in} C_{out}}{C_{in} + C_{out}}$$

Example value for $k = 9, C_{in} = 3, C_{out} = 64$

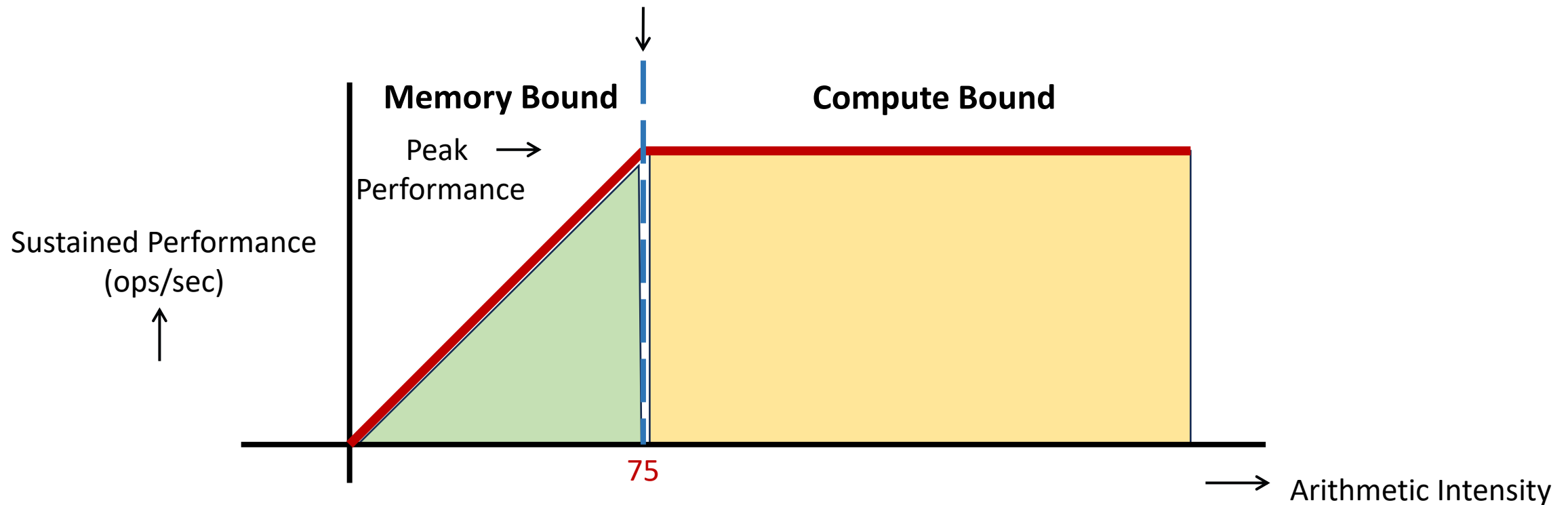
232

Example Machine Intensity = 75 (Nvidia A100*) how???

CNN on A100 Roofline Plot (assuming infinite cache)

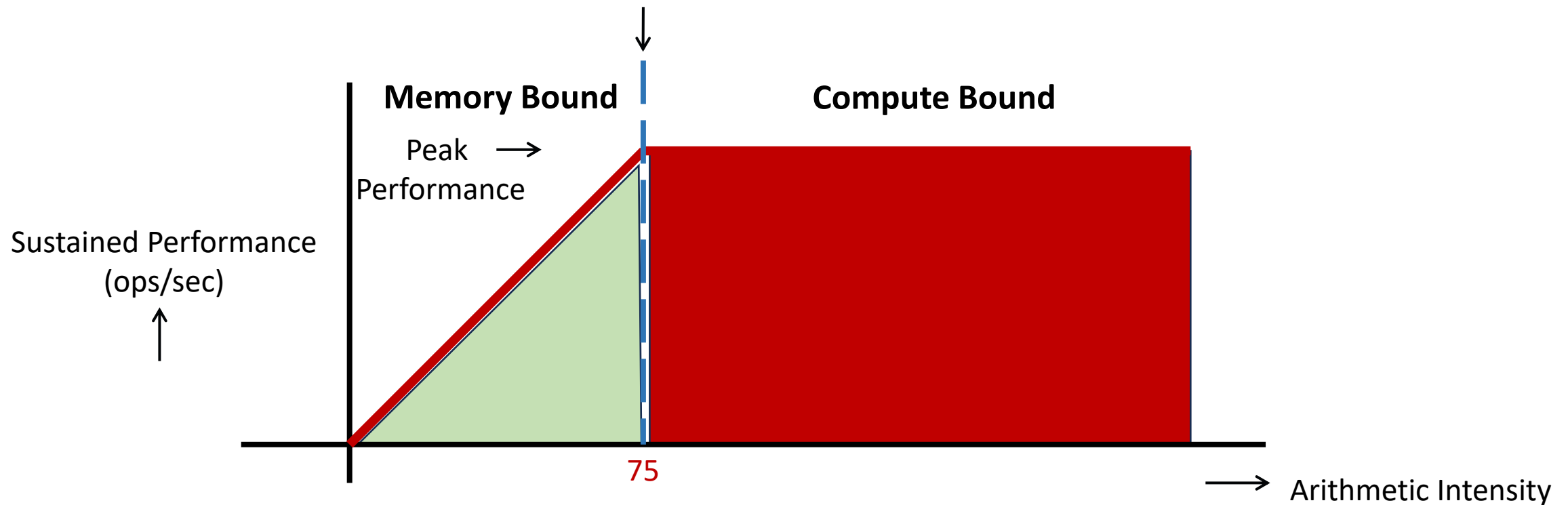


CNN on A100 Roofline Plot (assuming infinite cache)



Which area will the CNN Performance lie in?

CNN on A100 Roofline Plot (assuming infinite cache)



Which area will the CNN Performance lie in? **Compute Bound**

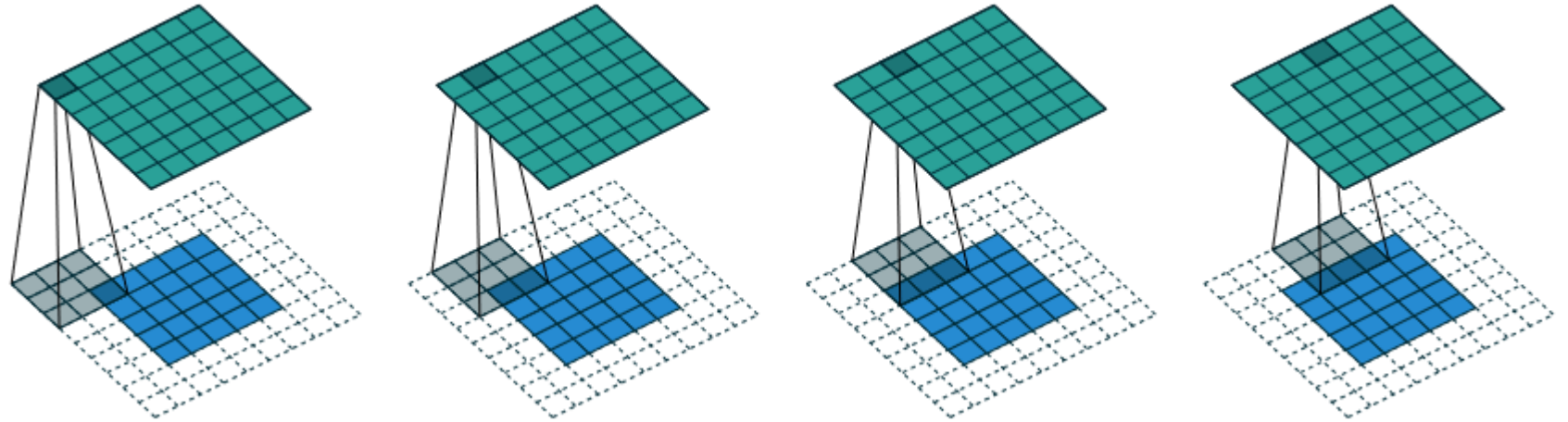
Implications

- Essentially a Compute Bound Problem, assuming large cache
- Objective: How to ensure efficiency is $O(1)$, i.e., all compute units are fully utilized
- If cache is limited
 - Objective 1: How to improve data reuse so that compute units have enough to work on?
 - Objective 2: How to ensure efficiency is $O(1)$, i.e., all compute units are fully utilized
- Another approach: increase Number of Devices
 - + Larger On-chip memory
 - + Higher bandwidth
 - - Communication/Synchronization overheads

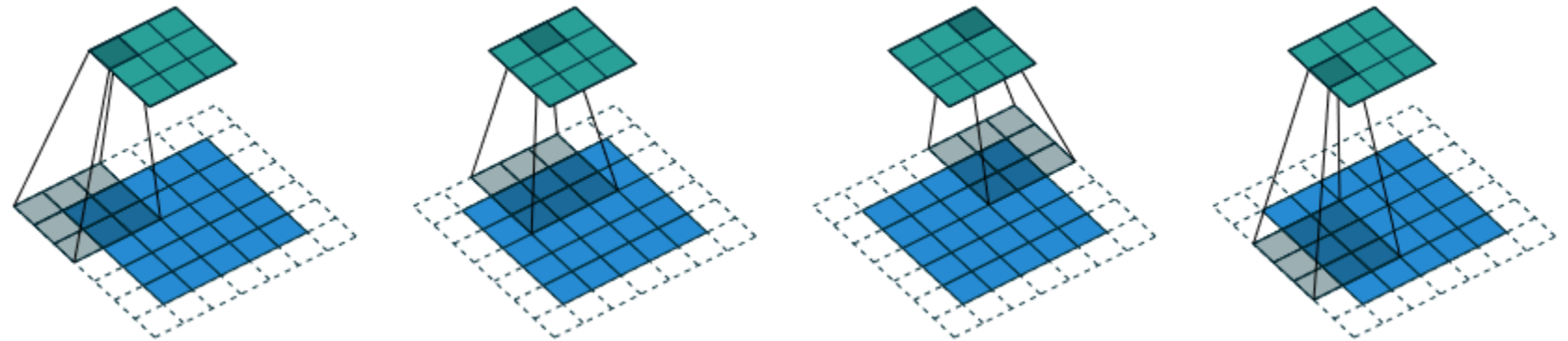
Variations

Reading Homework:

A guide to convolutional
arithmetic for deep learning



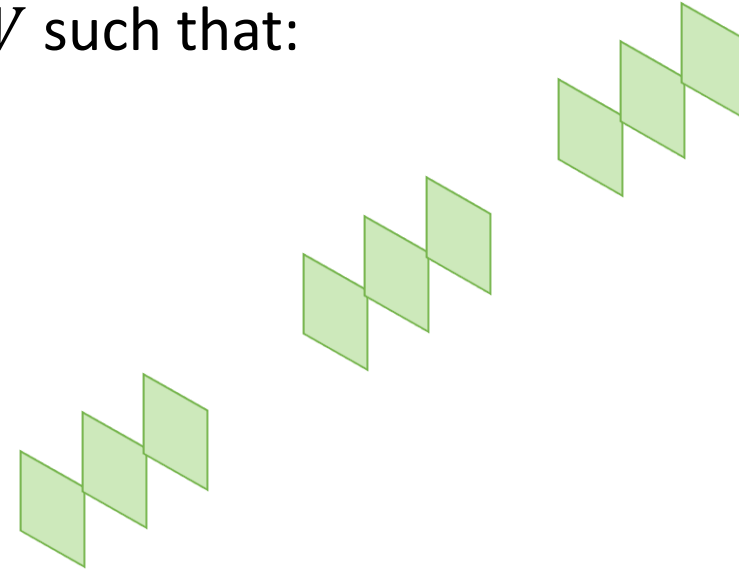
- Padding



- Strides

Convolutional Neural Network

- Training
 - Given $x_i \in X$ (image from a database), y_i : Label for the image
 - “Learn” a model F , defined by the filters W such that:
 - $\sum (y_i - F(x_i; W))^2$ is minimized
- Inference
 - Given a new image x'
 - Predict a label $y' = W(x')$
 - No change in the filters



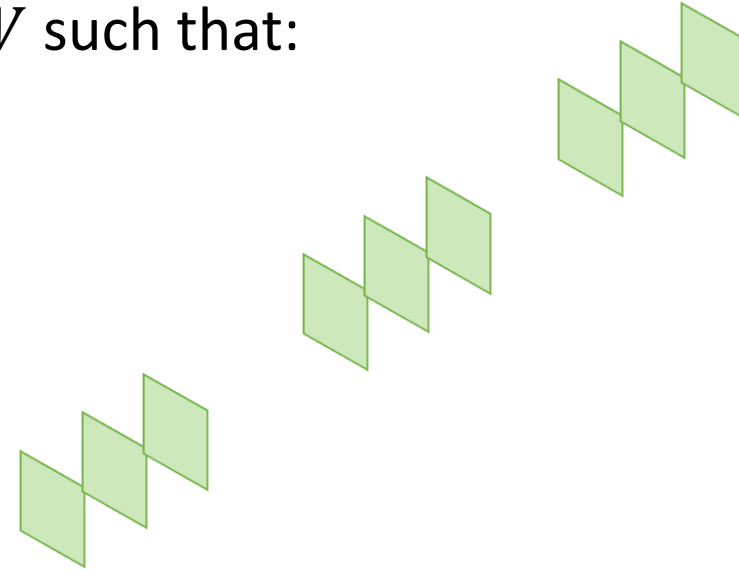
Convolutional Neural Network

- **Training**

- Given $x_i \in X$ (image from a database), y_i : Label for the image
- “Learn” a model F , defined by the filters W such that:
 - $\sum (y_i - F(x_i; W))^2$ is minimized

- **Inference**

- Given a new image x'
- Predict a label $y' = W(x')$
- No change in the filters



Training

$$E(W): \min_W \sum (y_i - F(x_i: W))^2$$

$$W_{t+1} \leftarrow W_t - \alpha \frac{\delta E(W)}{\delta W} \quad \leftarrow \text{Iteratively}$$

$$\frac{\delta E(W)}{\delta W} = 2 \sum_i (y_i - F(x_i: W)) \times \frac{\delta F(x_i: W)}{\delta W} \quad \leftarrow \text{For all samples, in each iteration}$$

Error

Training

L layers in CNN

$$W_1, W_2, W_3, \dots, W_L$$

$$O_1, O_2, O_3, \dots, O_L$$

Weights (Filters/Kernels)

Output of layer

$$2 * \text{Error} * \frac{\delta F(x_i : W_{1:L})}{\delta W_1}, \frac{\delta F(x_i : W_{1:L})}{\delta W_2}, \dots, \frac{\delta F(x_i : W_{1:L})}{\delta W_L}$$

Compute these for
all the samples

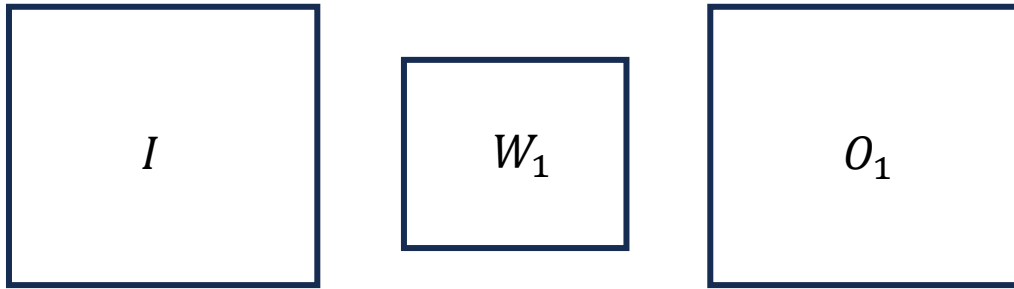
CNN Training: Forward Propagation



I

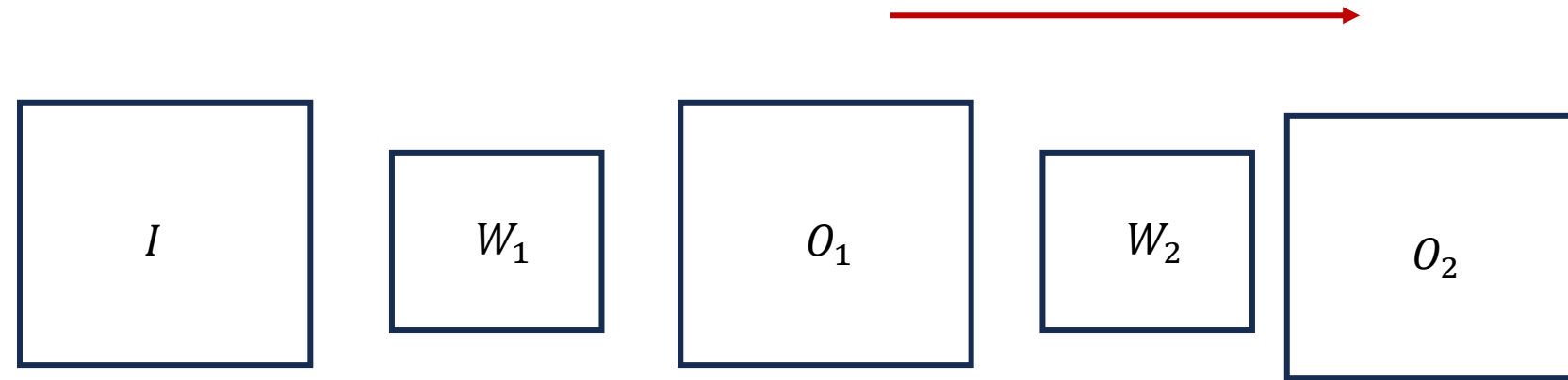
For a single image

CNN Training: Forward Propagation



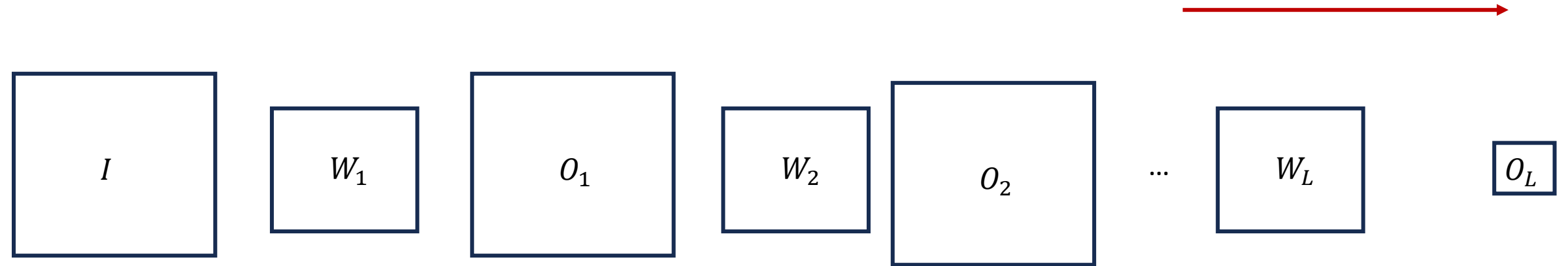
For a single image

CNN Training: Forward Propagation



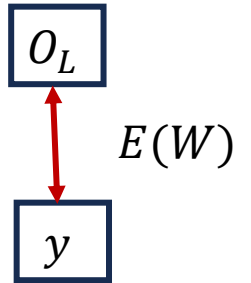
For a single image

CNN Training: Forward Propagation



For a single image

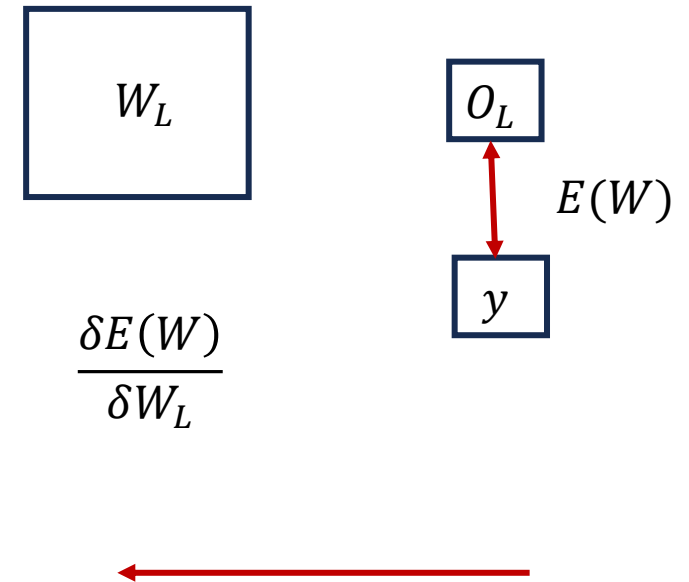
CNN Training: Error Calculation



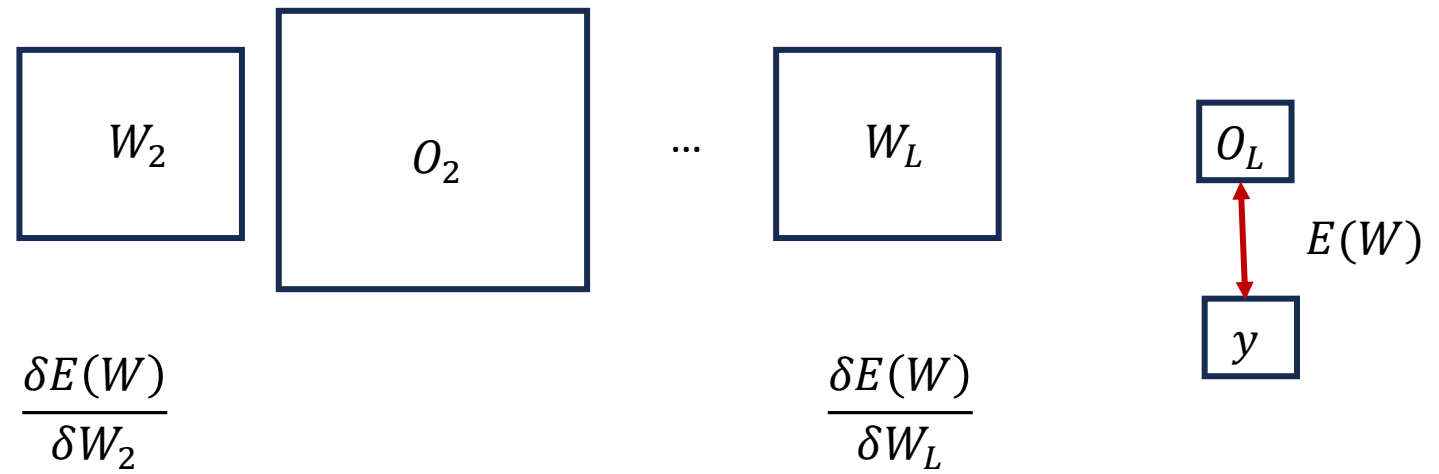
For a single image

CNN Training: Backward Propagation

For a single image

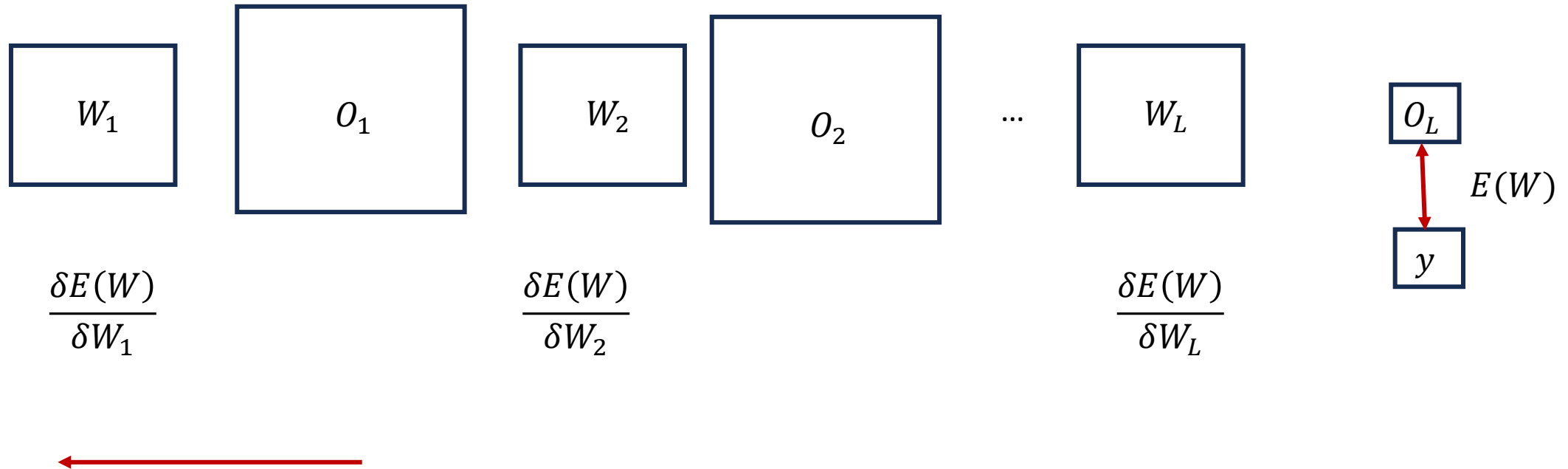


CNN Training: Backward Propagation



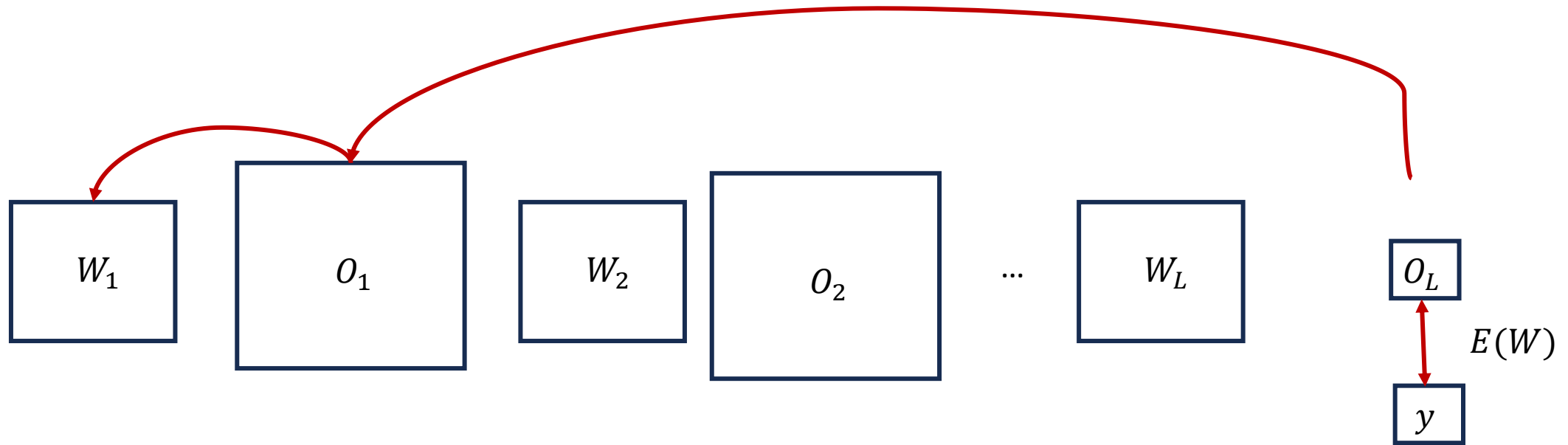
For a single image

CNN Training: Backward Propagation



For a single image

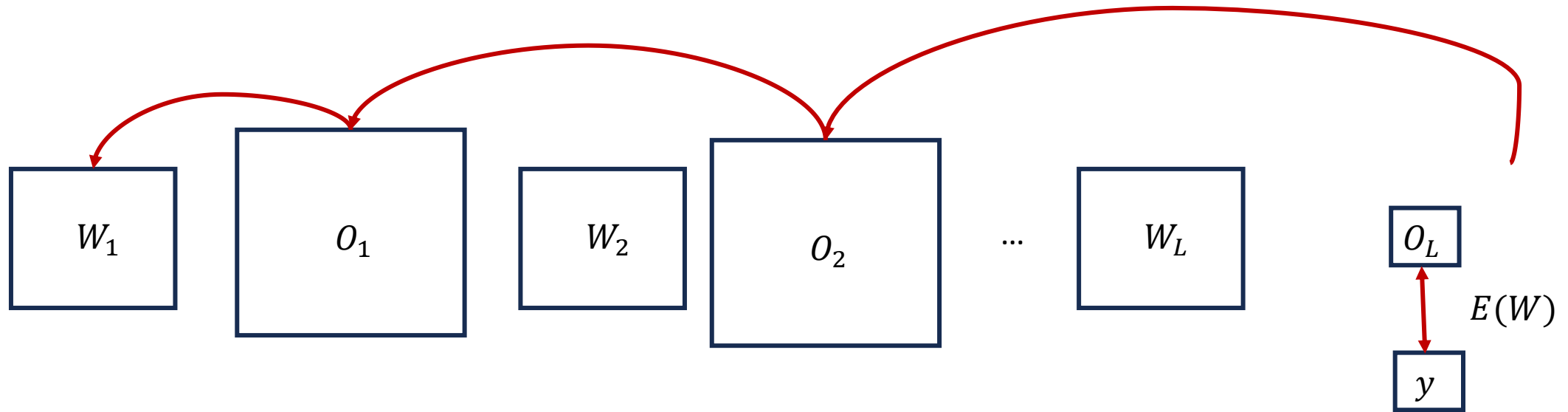
CNN Training: Backward Propagation



$$\frac{\delta E(W)}{\delta W_1} = \frac{\delta O_1}{\delta W_1} \times \frac{\delta E(W)}{\delta O_1}$$

For a single image

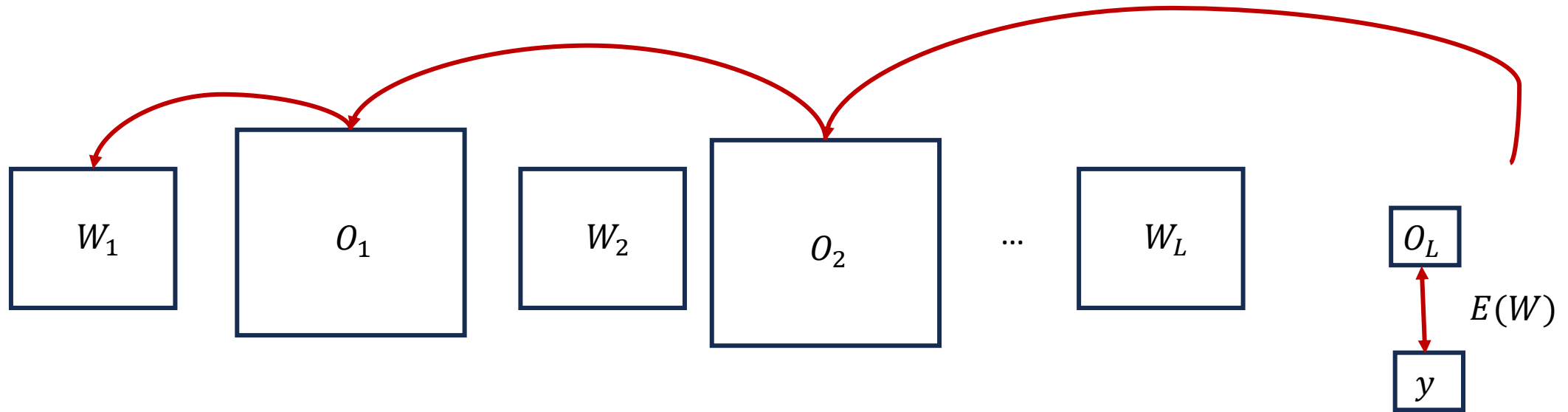
CNN Training: Backward Propagation



$$\frac{\delta E(W)}{\delta W_1} = \frac{\delta O_1}{\delta W_1} \times \frac{\delta O_2}{\delta O_1} \times \frac{\delta E(W)}{\delta O_2}$$

For a single image

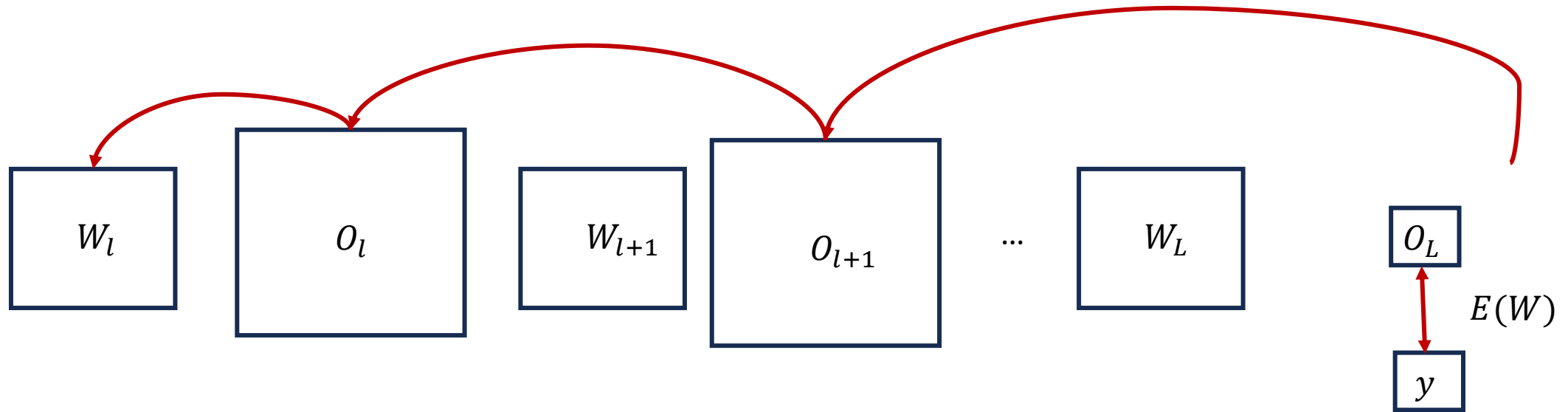
CNN Training: Backward Propagation



$$\frac{\delta E(W)}{\delta W_1} = \frac{\delta O_1}{\delta W_1} \times \frac{\delta O_2}{\delta O_1} \times \frac{\delta O_3}{\delta O_2} \times \dots \times \frac{\delta O_L}{\delta O_{L-1}} \times \frac{\delta E(W)}{\delta O_L}$$

For a single image

CNN Training: Backward Propagation



$$\frac{\delta E(W)}{\delta W_l} = \frac{\delta O_l}{\delta W_l} \times \frac{\delta O_{l+1}}{\delta O_l} \times \frac{\delta O_{l+2}}{\delta O_{l+1}} \times \dots \times \frac{\delta O_L}{\delta O_{L-1}} \times \frac{\delta E(W)}{\delta O_L}$$

For a single image

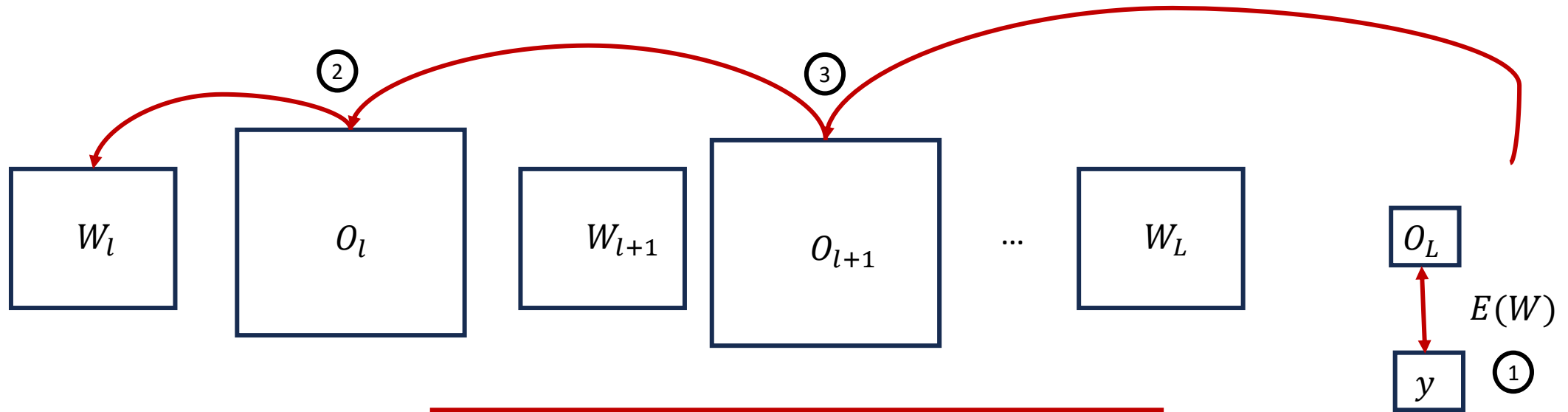
CNN Training: Backward Propagation

$$\frac{\delta E(W)}{\delta W_1} = \frac{\delta O_1}{\delta W_1} \times \frac{\delta O_2}{\delta O_1} \times \frac{\delta O_3}{\delta O_2} \times \cdots \times \frac{\delta O_L}{\delta O_{L-1}} \times \frac{\delta E(W)}{\delta O_L}$$

Reuse calculations performed for later layer
In earlier layer

$$\frac{\delta E(W)}{\delta W_l} = \frac{\delta O_l}{\delta W_l} \times \frac{\delta O_{l+1}}{\delta O_l} \times \frac{\delta O_{l+2}}{\delta O_{l+1}} \times \cdots \times \frac{\delta O_L}{\delta O_{L-1}} \times \frac{\delta E(W)}{\delta O_L}$$

CNN Training: Backward Propagation



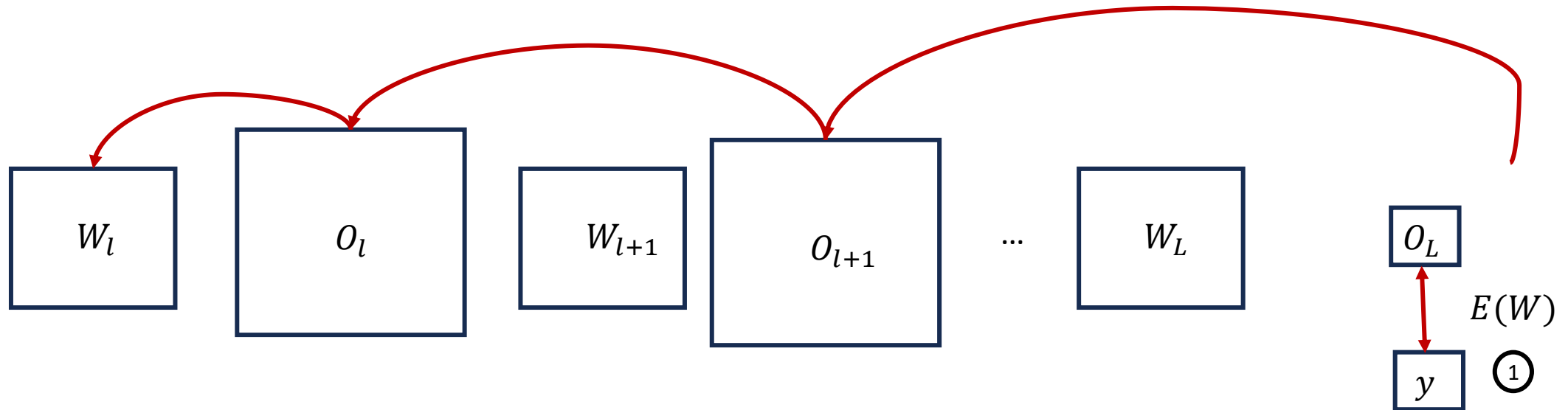
$$\frac{\delta E(W)}{\delta W_l} = \frac{\delta O_l}{\delta W_l} \times \frac{\delta O_{l+1}}{\delta O_l} \times \frac{\delta O_{l+2}}{\delta O_{l+1}} \times \dots \times \frac{\delta O_L}{\delta O_{L-1}} \times \frac{\delta E(W)}{\delta O_L}$$

The equation is annotated with circled numbers: (2) for $\frac{\delta O_l}{\delta W_l}$, (3) for $\frac{\delta O_{l+1}}{\delta O_l}$, and (1) for $\frac{\delta E(W)}{\delta O_L}$. The terms are grouped into nested dashed red boxes: the first two terms are in the outermost box, the next two are in a middle box, and the last two are in the innermost box.

We need to compute 3
types of terms

For a single image

CNN Training: Backward Propagation



$$\frac{\delta E(W)}{\delta W_l} = \frac{\delta O_l}{\delta W_l} \times \frac{\delta O_{l+1}}{\delta O_l} \times \frac{\delta O_{l+2}}{\delta O_{l+1}} \times \dots \times \frac{\delta O_L}{\delta O_{L-1}} \times \boxed{\frac{\delta E(W)}{\delta O_L} \textcircled{1}}$$

We need to compute 3
types of terms

For a single image

CNN Training: Backward Propagation

- $\frac{\delta E(W)}{\delta O_L}$???

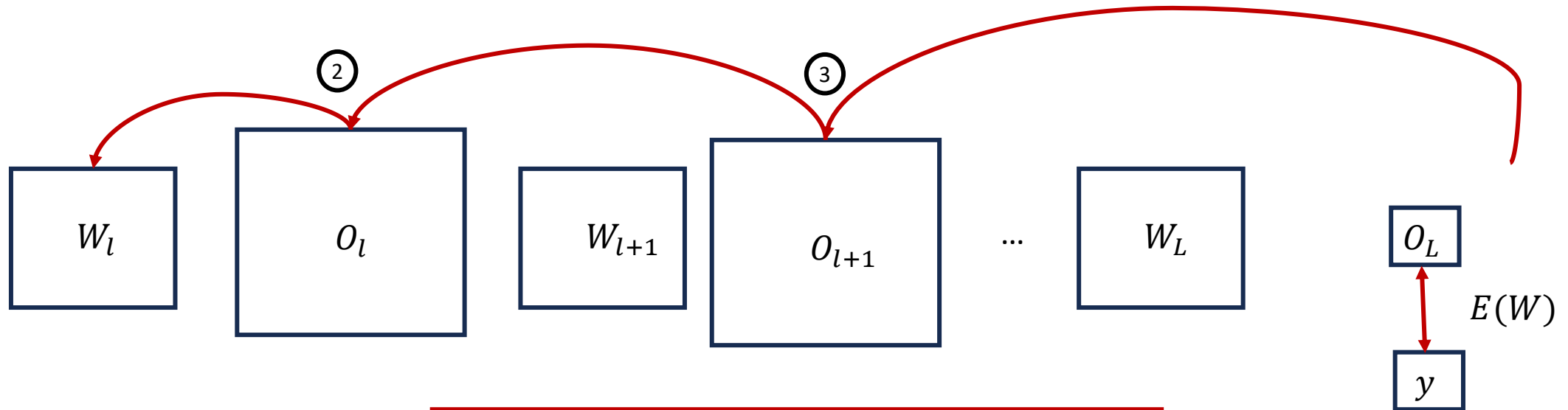
CNN Training: Backward Propagation

- $\frac{\delta E(W)}{\delta O_L} = \frac{\delta}{\delta O_L} (y_i - O_L)^2 = ??$

CNN Training: Backward Propagation

- $\frac{\delta E(W)}{\delta O_L} = \frac{\delta}{\delta O_L} (y_i - O_L)^2 = -2 * \textit{Error}_i$
- For i image

CNN Training: Backward Propagation



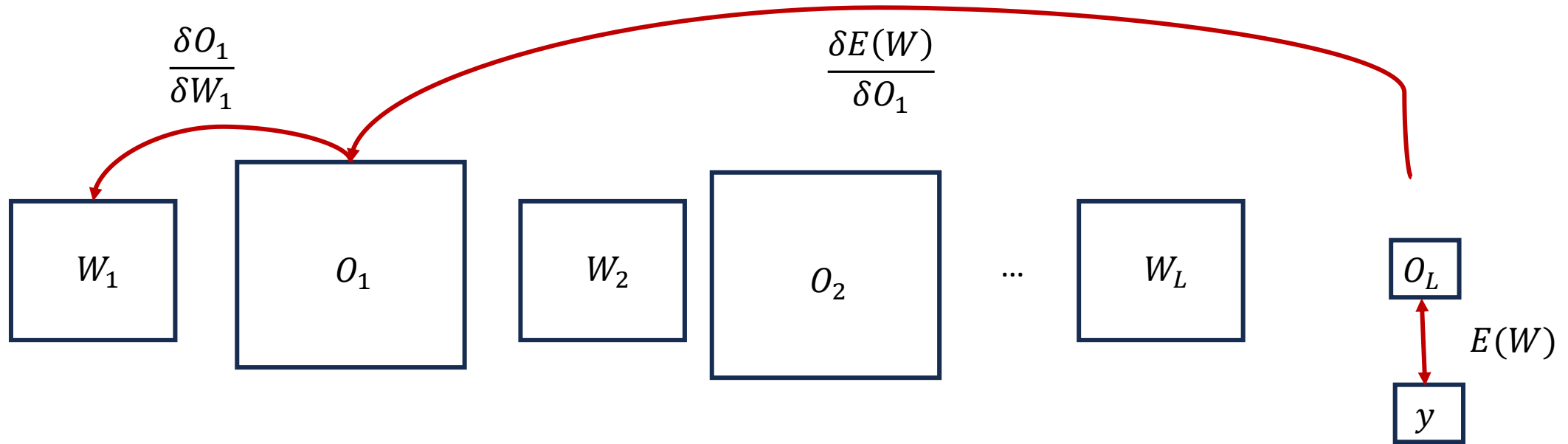
$$\frac{\delta E(W)}{\delta W_l} = \frac{\delta O_l}{\delta W_l} \times \frac{\delta O_{l+1}}{\delta O_l} \times \frac{\delta O_{l+2}}{\delta O_{l+1}} \times \dots \times \frac{\delta O_L}{\delta O_{L-1}} \times \frac{\delta E(W)}{\delta O_L}$$

The equation is enclosed in a red dashed box. The terms $\frac{\delta O_{l+1}}{\delta O_l}$ and $\frac{\delta O_L}{\delta O_{L-1}}$ are highlighted with red dashed boxes and labeled (3) and (2) respectively, corresponding to the arrows in the diagram above.

We need to compute 3
types of terms

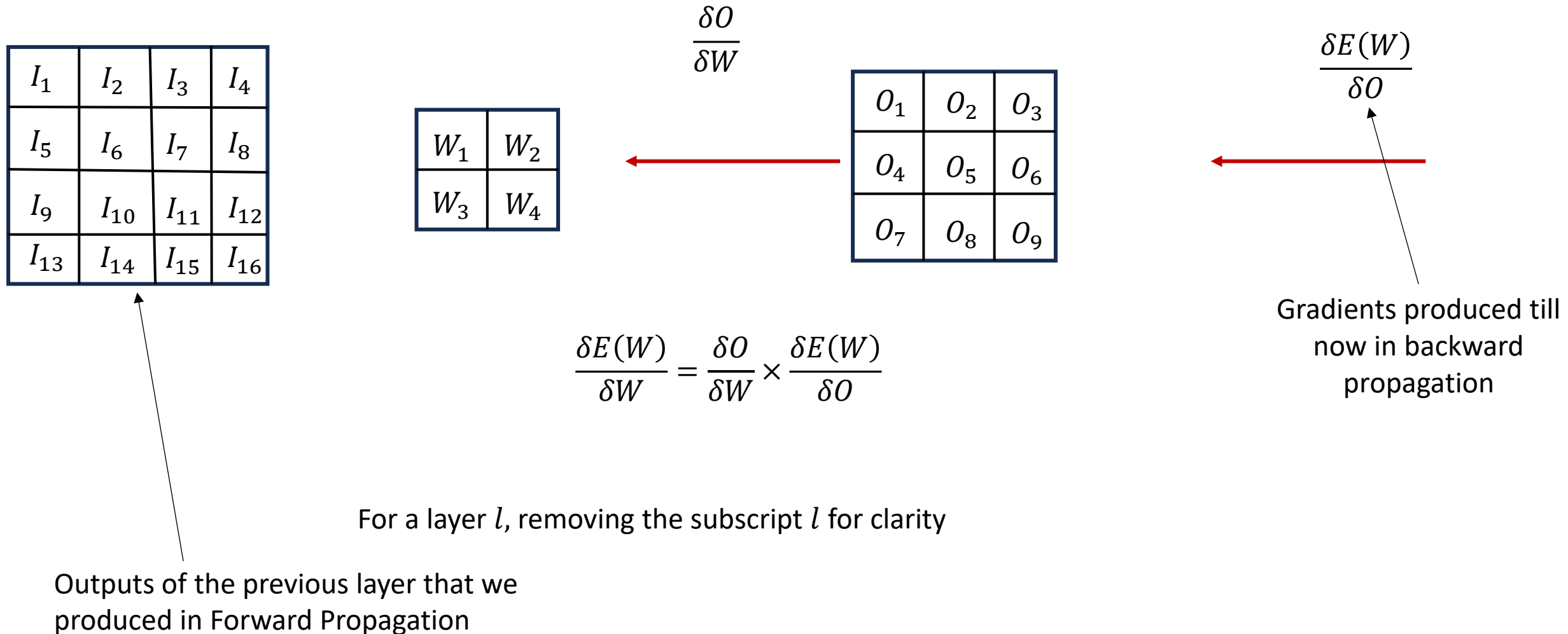
For a single image

CNN Training: Gradient Calculations



$$\frac{\delta E(W)}{\delta W_1} = \frac{\delta O_1}{\delta W_1} \times \frac{\delta E(W)}{\delta O_1}$$

CNN Training: Gradient Calculations



CNN Training: Gradient Calculations

I_1	I_2	I_3	I_4
I_5	I_6	I_7	I_8
I_9	I_{10}	I_{11}	I_{12}
I_{13}	I_{14}	I_{15}	I_{16}

W_1	W_2
W_3	W_4

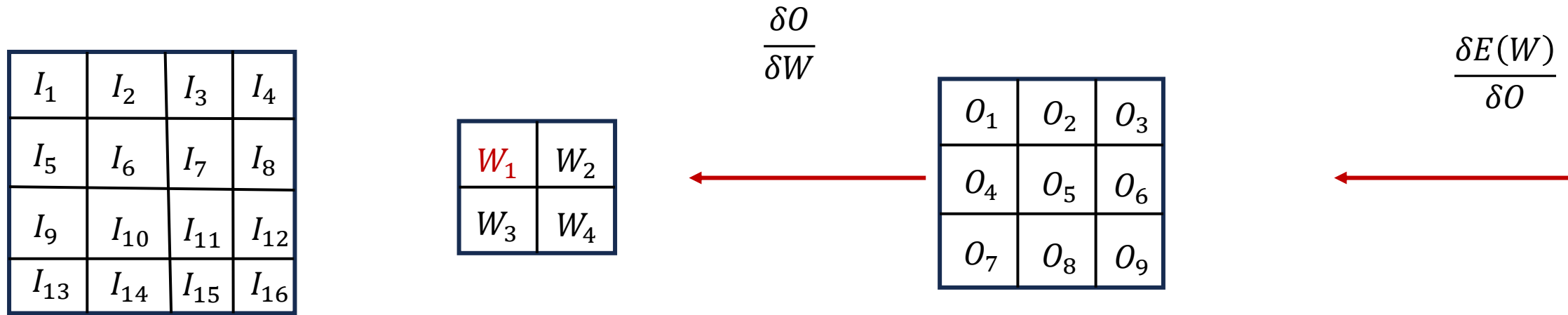
$$\frac{\delta O}{\delta W}$$

O_1	O_2	O_3
O_4	O_5	O_6
O_7	O_8	O_9

$$\frac{\delta E(W)}{\delta O}$$

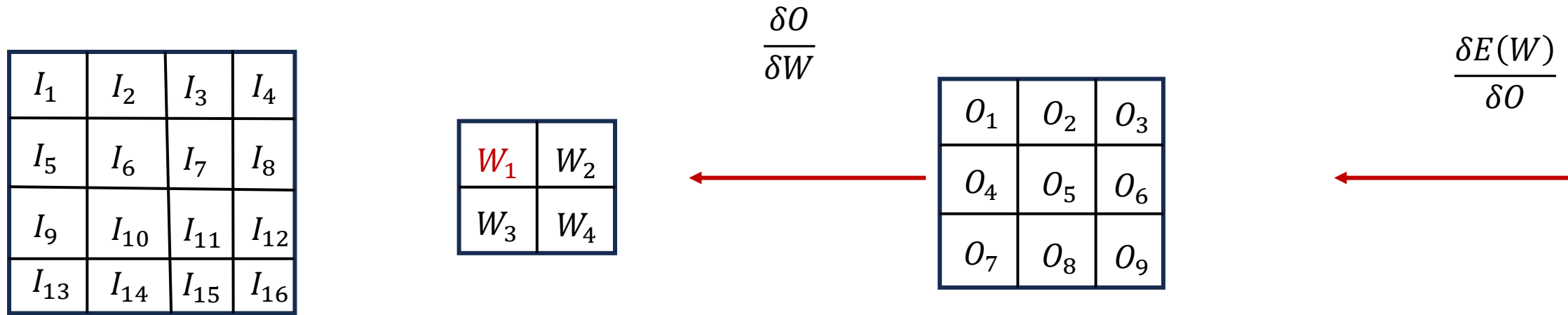
$$\frac{\delta E(W)}{\delta W_1} = \frac{\delta O}{\delta W_1} \times \frac{\delta E(W)}{\delta O}$$

CNN Training: Gradient Calculations



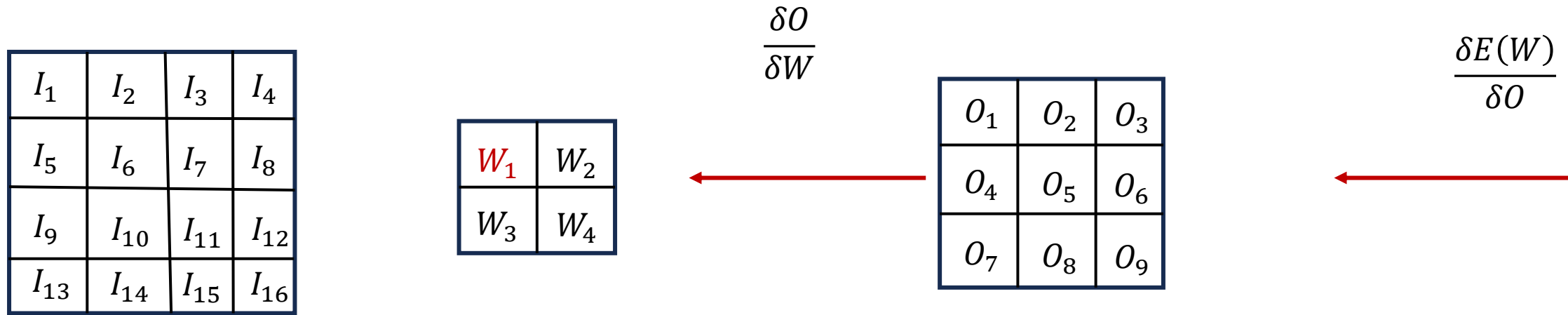
$$\frac{\delta E(W)}{\delta W_1} = \frac{\delta O_1}{\delta W_1} \times \frac{\delta E(W)}{\delta O_1} + \frac{\delta O_2}{\delta W_1} \times \frac{\delta E(W)}{\delta O_2} + \frac{\delta O_3}{\delta W_1} \times \frac{\delta E(W)}{\delta O_3} \dots + \frac{\delta O_9}{\delta W_1} \times \frac{\delta E(W)}{\delta O_9}$$

CNN Training: Gradient Calculations



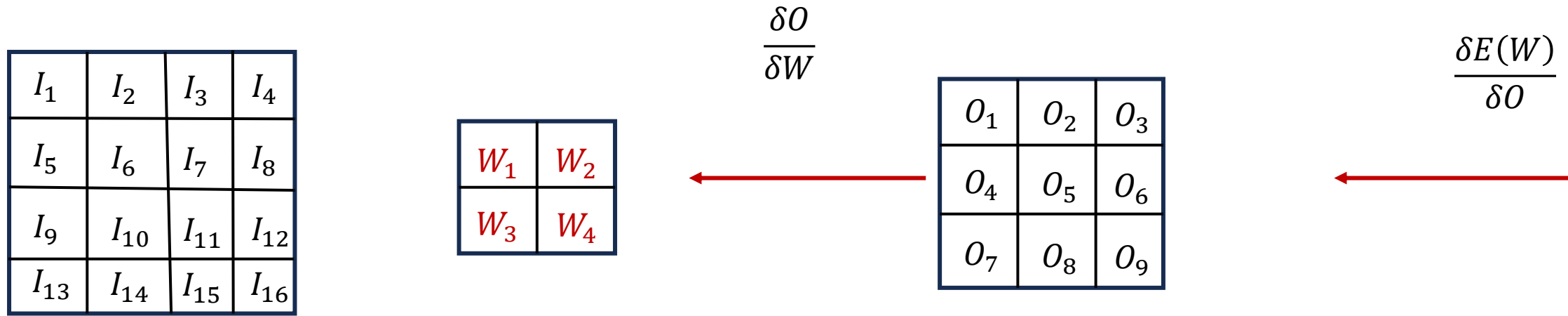
$$\frac{\delta E(W)}{\delta W_1} = \frac{\delta O_1}{\delta W_1} \times \frac{\delta E(W)}{\delta O_1} + \frac{\delta O_2}{\delta W_1} \times \frac{\delta E(W)}{\delta O_2} + \frac{\delta O_3}{\delta W_1} \times \frac{\delta E(W)}{\delta O_3} \dots + \frac{\delta O_9}{\delta W_1} \times \frac{\delta E(W)}{\delta O_9}$$

CNN Training: Gradient Calculations



$$\frac{\delta E(W)}{\delta W_1} = I_1 \times \frac{\delta E(W)}{\delta O_1} + I_2 \times \frac{\delta E(W)}{\delta O_2} + I_3 \times \frac{\delta E(W)}{\delta O_3} \dots + I_{11} \times \frac{\delta E(W)}{\delta O_9}$$

CNN Training: Gradient Calculations



$$\begin{aligned}\frac{\delta E(W)}{\delta W_1} &= I_1 \times \frac{\delta E(W)}{\delta O_1} + I_2 \times \frac{\delta E(W)}{\delta O_2} + I_3 \times \frac{\delta E(W)}{\delta O_3} \dots + I_{11} \times \frac{\delta E(W)}{\delta O_9} \\ \frac{\delta E(W)}{\delta W_2} &= I_2 \times \frac{\delta E(W)}{\delta O_1} + I_3 \times \frac{\delta E(W)}{\delta O_2} + I_6 \times \frac{\delta E(W)}{\delta O_3} \dots + I_{12} \times \frac{\delta E(W)}{\delta O_9} \\ \frac{\delta E(W)}{\delta W_3} &= I_5 \times \frac{\delta E(W)}{\delta O_1} + I_6 \times \frac{\delta E(W)}{\delta O_2} + I_7 \times \frac{\delta E(W)}{\delta O_3} \dots + I_{13} \times \frac{\delta E(W)}{\delta O_9}\end{aligned}$$

...

CNN Training: Gradient Calculations

I_1	I_2	I_3	I_4
I_5	I_6	I_7	I_8
I_9	I_{10}	I_{11}	I_{12}
I_{13}	I_{14}	I_{15}	I_{16}

W_1	W_2
W_3	W_4

$$\frac{\delta O}{\delta W}$$

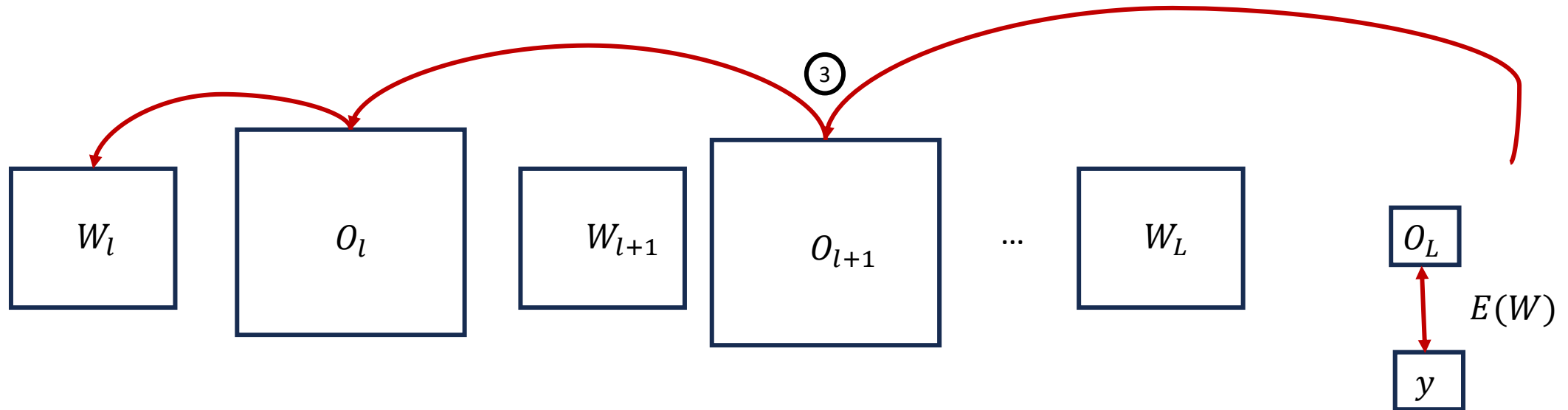
$\frac{\delta E(W)}{\delta O}$		
$\frac{\delta E}{\delta O_1}$	$\frac{\delta E}{\delta O_2}$	$\frac{\delta E}{\delta O_3}$
$\frac{\delta E}{\delta O_4}$	$\frac{\delta E}{\delta O_5}$	$\frac{\delta E}{\delta O_6}$
$\frac{\delta E}{\delta O_7}$	$\frac{\delta E}{\delta O_8}$	$\frac{\delta E}{\delta O_9}$

$$\frac{\delta E(W)}{\delta W} = \text{Conv}(I, \frac{\delta E(W)}{\delta O})$$

Note: Only inputs I and gradients $\frac{\delta E(W)}{\delta O}$ are needed. We don't need the output O for this layer. Why?

Note: This operation where the output dimension is greater than input dimension is called deconvolution or transposed convolution

CNN Training: Backward Propagation

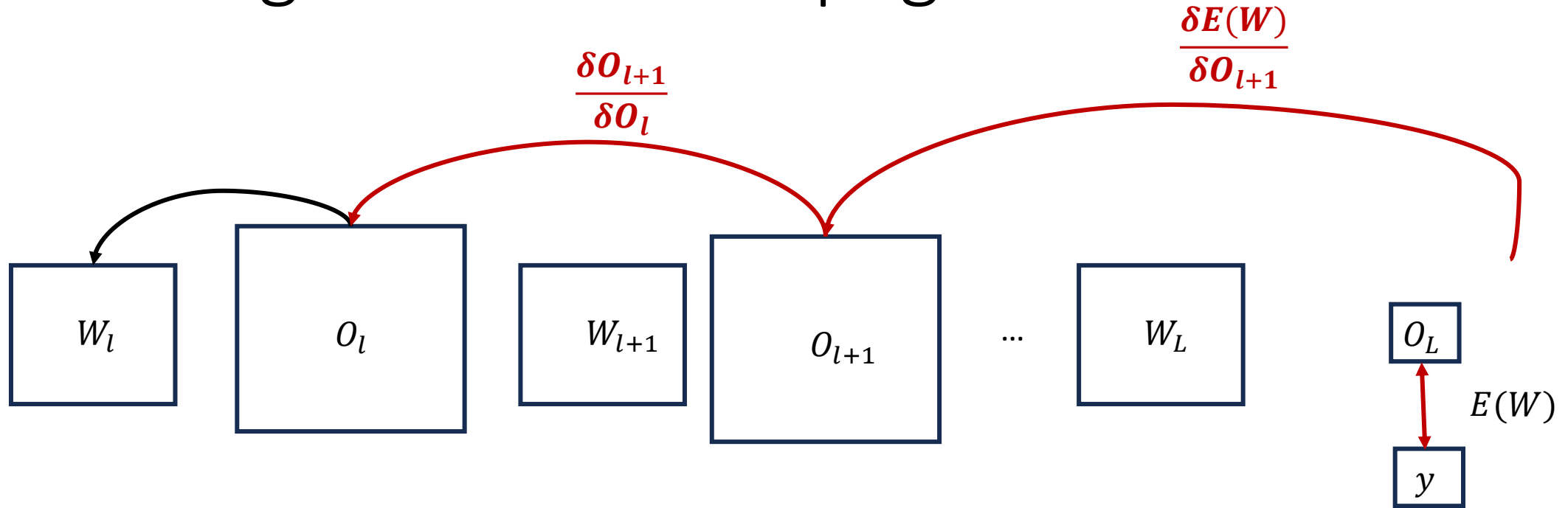


$$\frac{\delta E(W)}{\delta W_l} = \frac{\delta O_l}{\delta W_l} \times \frac{\delta O_{l+1}}{\delta O_l} \times \frac{\delta O_{l+2}}{\delta O_{l+1}} \times \dots \times \frac{\delta O_L}{\delta O_{L-1}} \times \frac{\delta E(W)}{\delta O_L}$$

We need to compute 3
types of terms

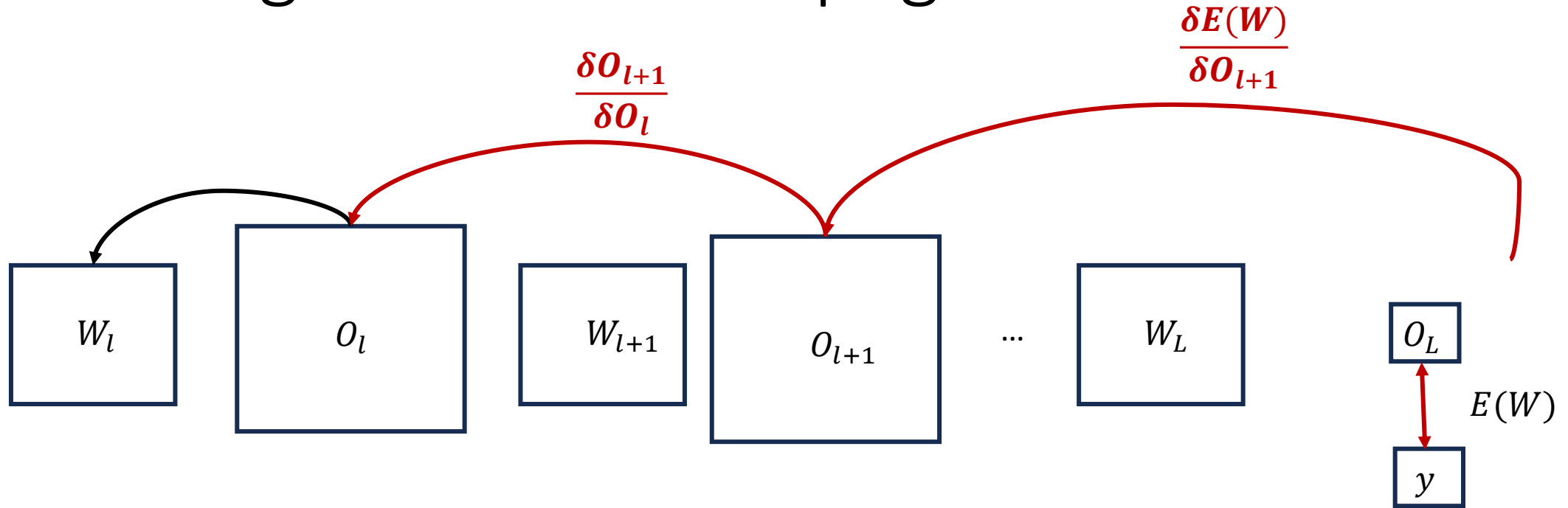
For a single image

CNN Training: Backward Propagation



$$\frac{\delta E(W)}{\delta O_l} = \frac{\delta O_{l+1}}{\delta O_l} \times \frac{\delta O_{l+2}}{\delta O_{l+1}} \times \dots \times \frac{\delta O_L}{\delta O_{L-1}} \times \frac{\delta E(W)}{\delta O_L}$$

CNN Training: Backward Propagation



$$\frac{\delta E(W)}{\delta O_l} = \frac{\delta O_{l+1}}{\delta O_l} \times \frac{\delta E(W)}{\delta O_{l+1}}$$

CNN Training: Backward Propagation

$$O_l$$

O_1	O_2	O_3	O_4
O_5	O_6	O_7	O_8
O_9	O_{10}	O_{11}	O_{12}
O_{13}	O_{14}	O_{15}	O_{16}

W_1	W_2
W_3	W_4

$$O_{l+1}$$

$\frac{\delta E}{\delta O_1}$	$\frac{\delta E}{\delta O_2}$	$\frac{\delta E}{\delta O_3}$
$\frac{\delta E}{\delta O_4}$	$\frac{\delta E}{\delta O_5}$	$\frac{\delta E}{\delta O_6}$
$\frac{\delta E}{\delta O_7}$	$\frac{\delta E}{\delta O_8}$	$\frac{\delta E}{\delta O_9}$

What all do we need for the computation? Do we need inputs or outputs to the layer?

$$\frac{\delta E(W)}{\delta O_l} = \text{Conv}(W_{180}, \frac{\delta E(W)}{\delta O_{l+1}})$$

HW: W_{180} – W rotated by 180 degrees

Key Operation in CNN

- Forward Propagation:

$$O_l = \textit{Conv}(W_l, I_l)$$

- Backward Propagation

$$\frac{\delta E(W)}{\delta W_l} = \textit{Conv}(I_l, \frac{\delta E(W)}{\delta O_l})$$

$$\frac{\delta E(W)}{\delta O_l} = \textit{Conv}(W_{180}, \frac{\delta E(W)}{\delta O_{l+1}})$$

Training

$$E(W): \min_W \sum (y_i - F(x_i: W))^2$$

$$W_{t+1} \leftarrow W_t - \alpha \frac{\delta E(W)}{\delta W} \quad \leftarrow \text{Iteratively}$$

$$\frac{\delta E(W)}{\delta W} = 2 \underbrace{\sum (y_i - F(x_i: W))}_{\text{Error}} \times \frac{\delta F(x_i: W)}{\delta W} \quad \leftarrow \text{For all samples, in each iteration}$$

Error

For multiple images in a batch: Simply sum up the gradients

$$\frac{\delta E(W)}{\delta W_l} = \frac{\delta O_l}{\delta W_l} \times \frac{\delta O_{l+1}}{\delta O_l} \times \frac{\delta O_{l+2}}{\delta O_{l+1}} \times \dots \times \frac{\delta O_L}{\delta O_{L-1}} \times \frac{\delta E(W)}{\delta O_L}$$

Training

L layers in CNN

$W_1, W_2, W_3, \dots, W_L$
 $O_1, O_2, O_3, \dots, O_L$

Weights (Filters/Kernels)

Output of layer

$$2 * \text{Error} * \frac{\delta F(x_i : W_{1:L})}{\delta W_1}, \frac{\delta F(x_i : W_{1:L})}{\delta W_2}, \dots, \frac{\delta F(x_i : W_{1:L})}{\delta W_L}$$

Training – Key Performance Objective

- Accuracy – We want to train a model that can predict with low error on unseen data
- Throughput – Given a large dataset, we want to train as quickly as possible

Training – Key Requirement

- For each image,
- Need to store the inputs of each layer
 - Note: Input I_l of layer l is output (activation) O_{l-1} of layer $l - 1$
- Need to have gradients $\frac{\delta E(W)}{\delta W_l}, \frac{\delta E(W)}{\delta O_l}, \frac{\delta E(W)}{\delta O_{l+1}}$ for each layer l
- Ungraded HW assignment: Can you calculate the difference in storage requirements between computing gradients from left to right vs right to left?

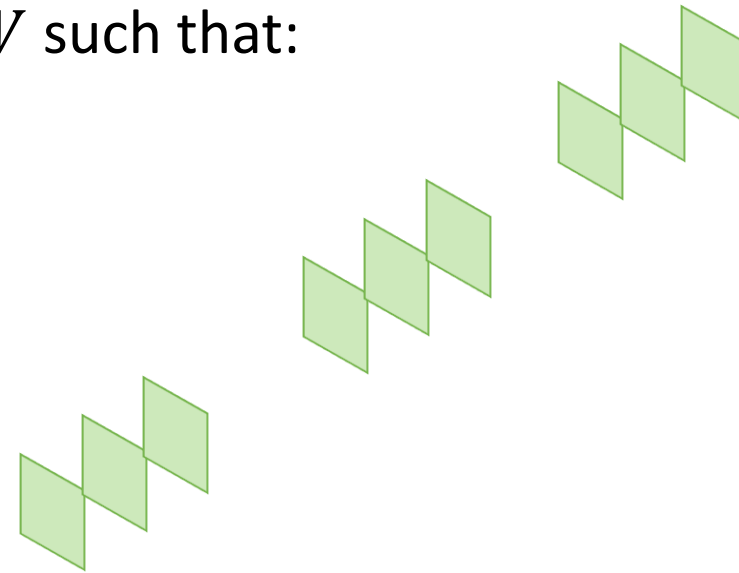
Convolutional Neural Network

- **Training**

- Given $x_i \in X$ (image from a database), y_i : Label for the image
- “Learn” a model F , defined by the filters W such that:
 - $\sum (y_i - F(x_i; W))^2$ is minimized

- **Inference**

- Given a new image x'
- Predict a label $y' = W(x')$
- No change in the filters



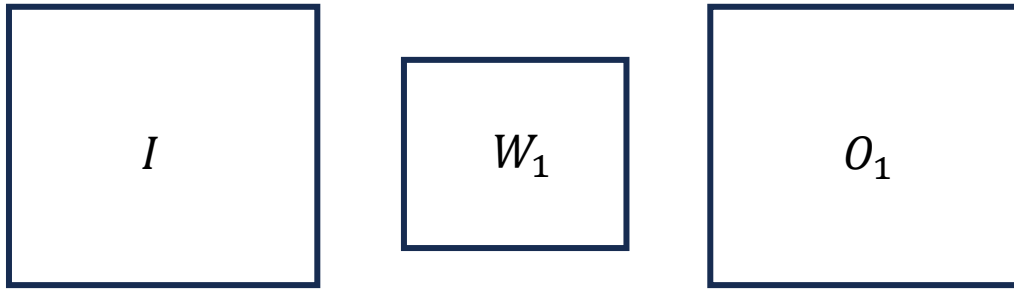
CNN Training: Forward Propagation



I

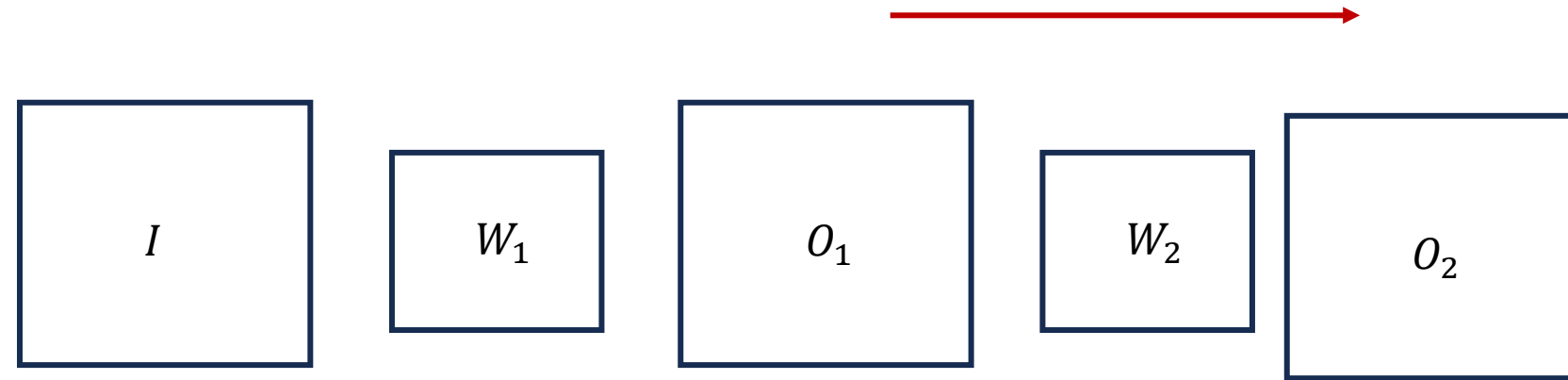
For a single image

CNN Training: Forward Propagation



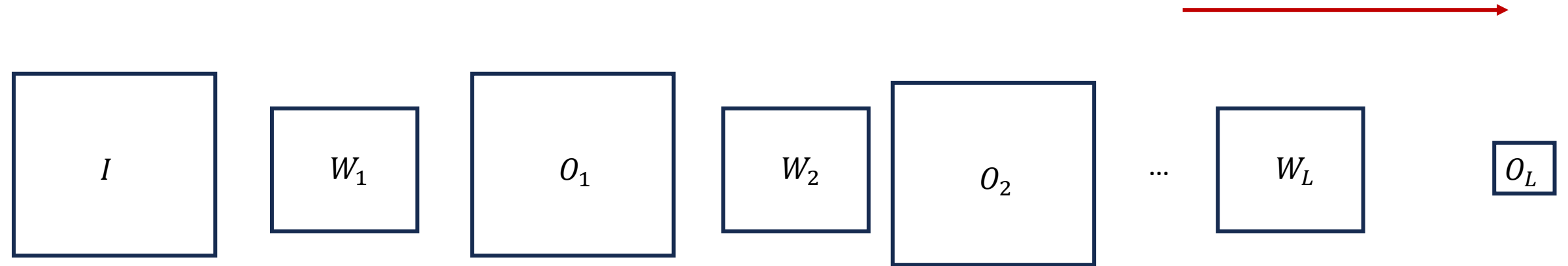
For a single image

CNN Training: Forward Propagation



For a single image

CNN Training: Forward Propagation



For a single image

Inference

- Simply the forward propagation portion
- No error calculation or backpropagation as ground truth doesn't exist
- Performance Objectives:
 - Accuracy – Although ground truth is not known
 - Latency – Perform inference as fast as possible on single or a batch of images

Next Class

- 9/25 Lecture 10
 - Accelerating Convolutional Neural Networks: Convolution as Matrix Multiplication

Thank You

- Questions?
- Email: sanmukh.kuppannagari@case.edu