# IT5005 Artificial Intelligence
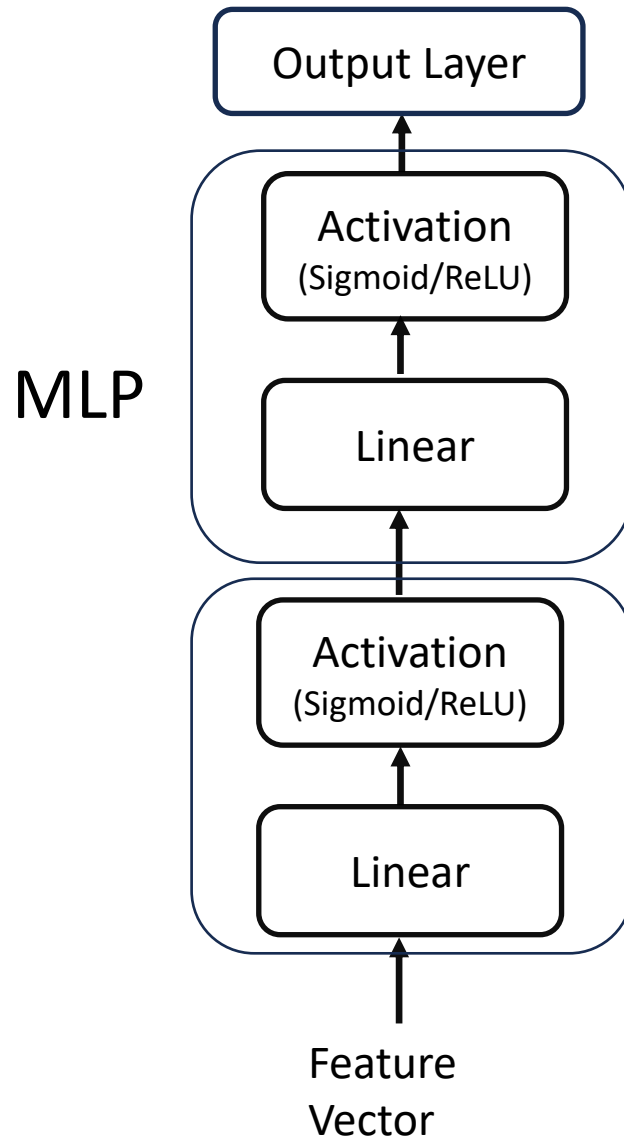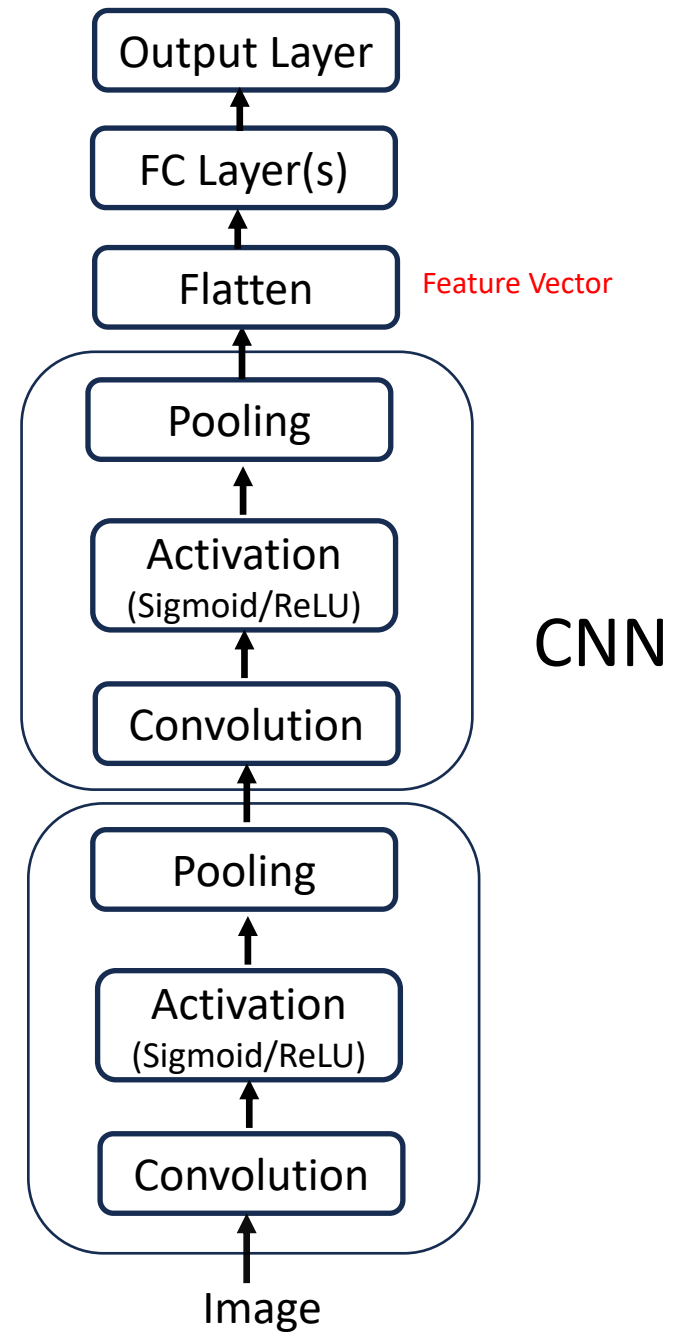
Sirigina Rajendra Prasad
AY2025/2026: Semester 1

## CNN Contd.
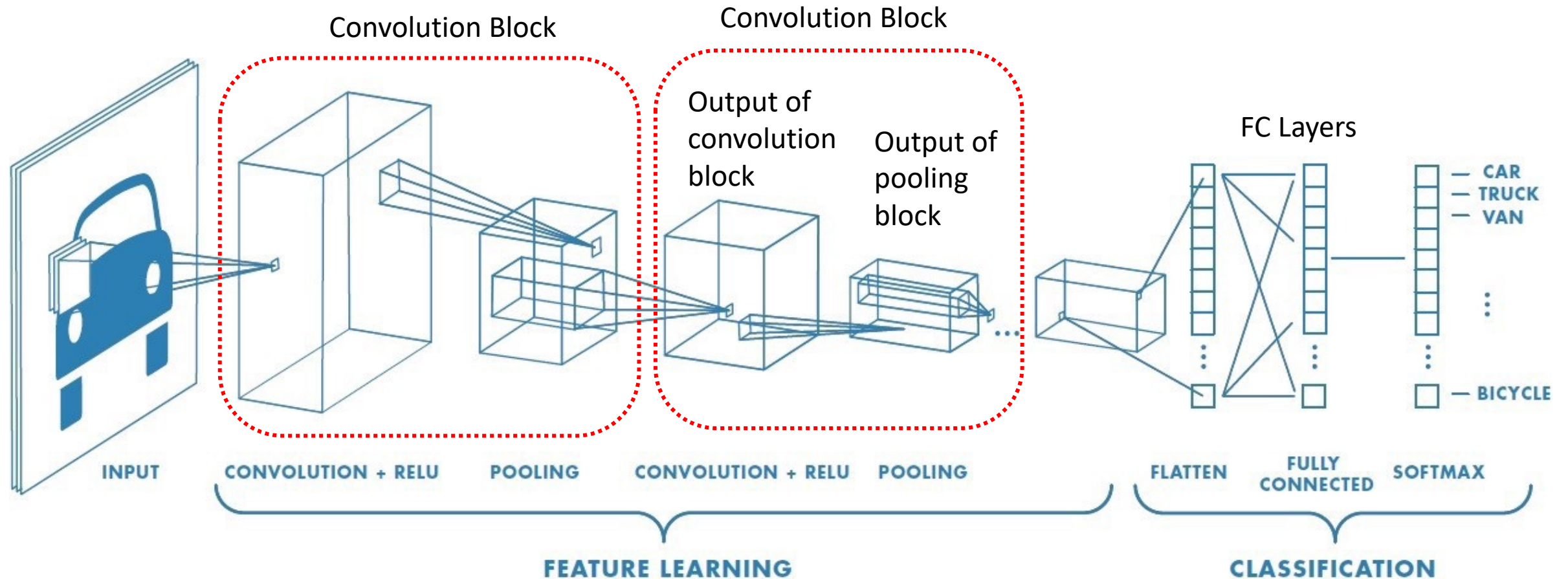
# MLP Vs CNN
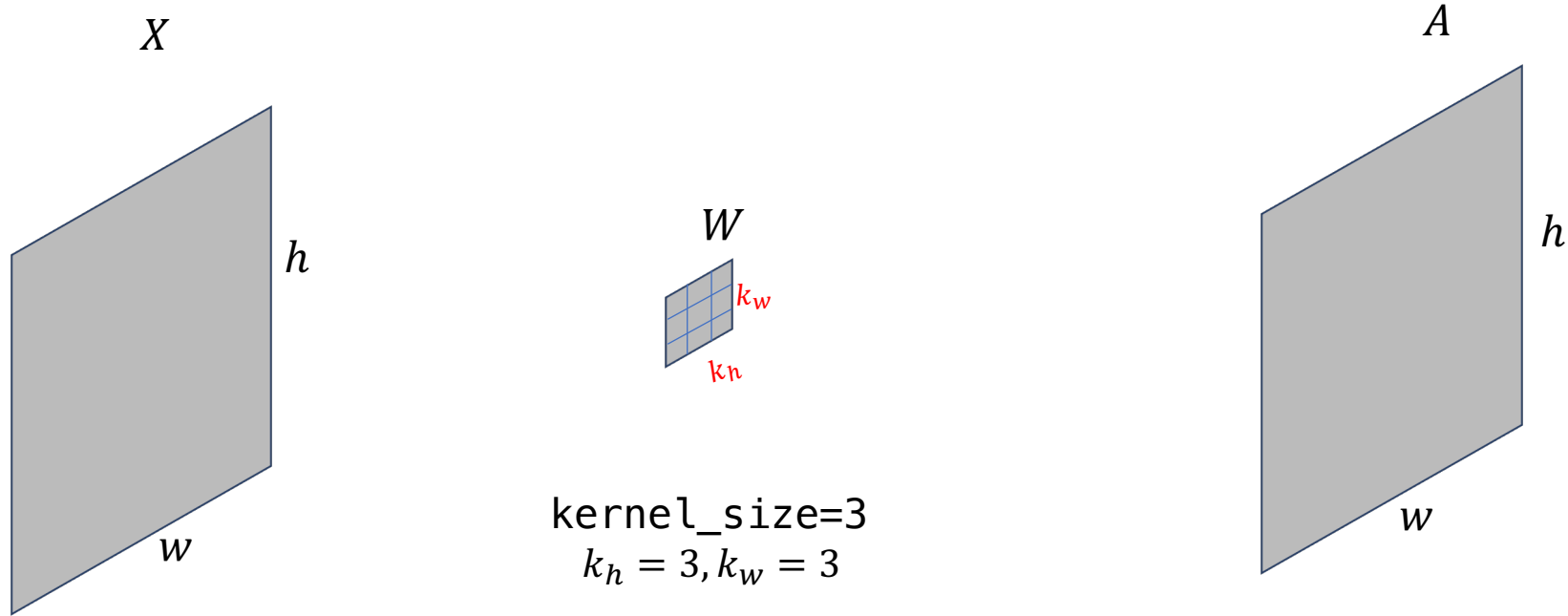


MLP

- Output Layer
- Activation (Sigmoid/ReLU)
- Linear
- Activation (Sigmoid/ReLU)
- Linear
- Feature Vector

Vs

CNN

- Output Layer
- FC Layer(s)
- Flatten — Feature Vector
- Pooling
- Activation (Sigmoid/ReLU)
- Convolution
- Pooling
- Activation (Sigmoid/ReLU)
- Convolution
- Image

# A Closer Look at CNN Architecture

Image credit: https://towardsdatascience.com/a-comprehensive-guide-to-convolutional-neural-networks-the-eli5-way-3bd2b1164a53

# Convolution Layer

```
nn.Conv2d(in_channels= , out_channels= , kernel_size=3, stride=1, padding=1)
```

$X$

$h$

$w$

$W$
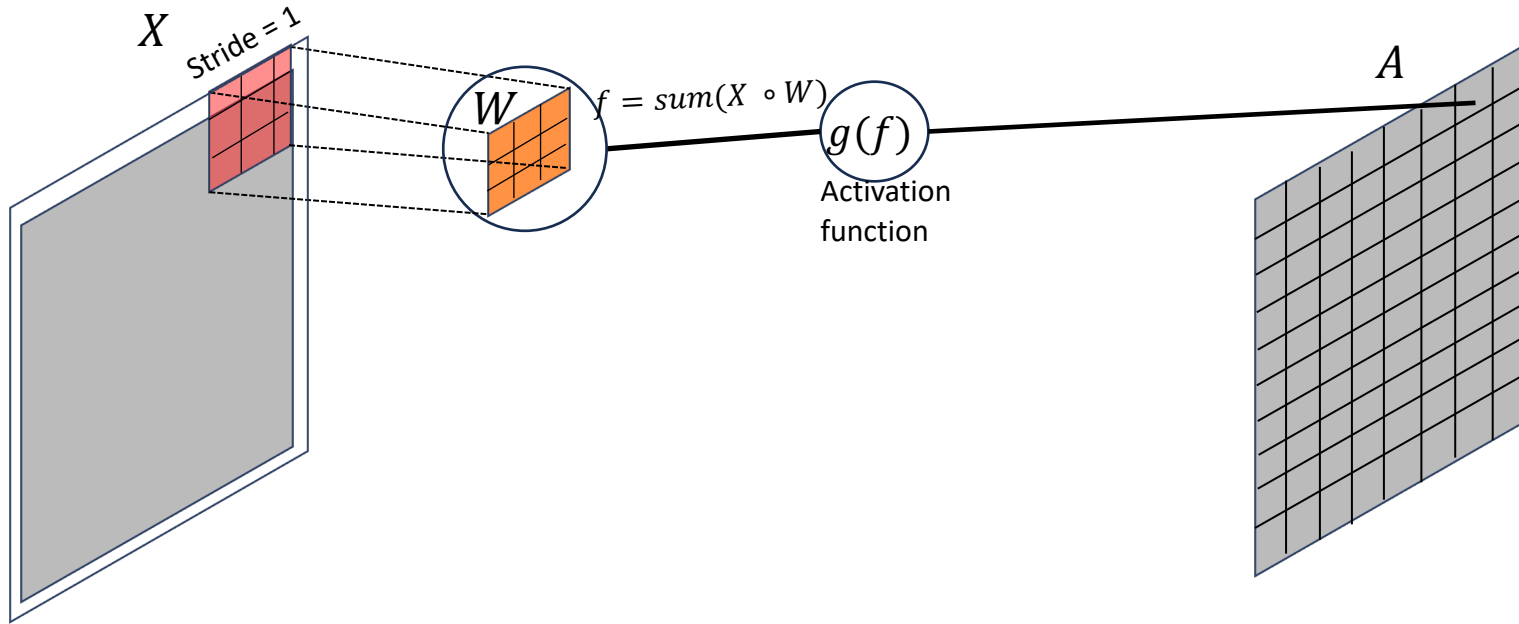
$k_w$

$k_h$

kernel_size=3
$k_h = 3, k_w = 3$

$A$

$h$

$w$

# Convolution Layer

$$f = sum(X \circ W)$$
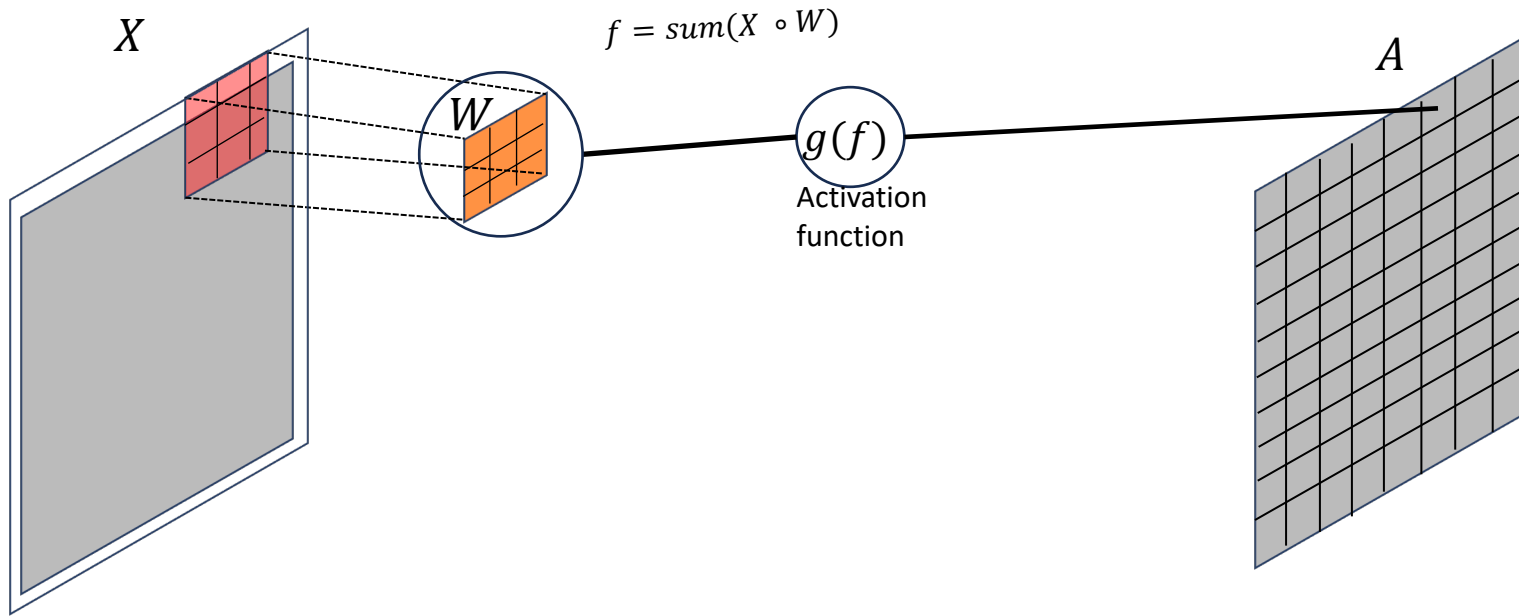
1. $X \circ W$ : Elementwise multiplication of $W$ and $X$
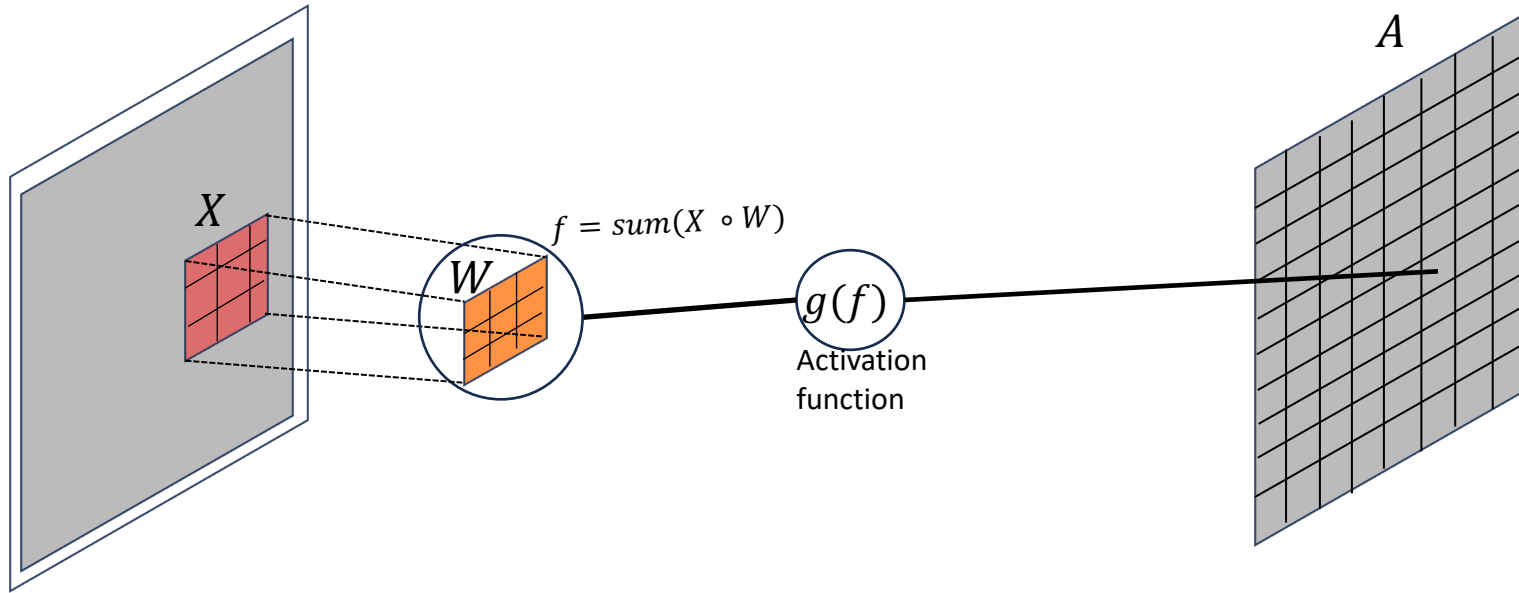2. $sum(X \circ W)$ sum of the elements of $X \circ W$

MAC: Multiply-Accumulate

$X$

Padding = 1

$f = sum(X \circ W)$

$W$

$g(f)$

Activation function

$A$

# Convolution + Pooling Layer: Gray Scale Image



$X$

Stride = 1

$W$

$f = sum(X \circ W)$

$g(f)$

Activation function

$A$

# Convolution + Pooling Layer: Gray Scale Image

$X$

$f = sum(X \circ W)$

$A$

$W$

$g(f)$

Activation function

# Convolution + Pooling Layer: Gray Scale Image

$X$

$W$

$f = sum(X \circ W)$

$g(f)$

Activation
function

$A$

# Convolution + Pooling Layer: Gray Scale Image



$X$

$W$

$f = sum(X \circ W)$

$g(f)$

Activation function

$A$

What is the shape of $A$?

How many parameters (weights)?

# Size of Each Kernel's Output

$$O = \left\lfloor \frac{I + 2P - K}{S} \right\rfloor + 1$$

where

$O$ is the output dimension (either height or width).

$I$ is the output dimension (either height or width).

$P = padding$

$K = kernel\_size$

$S = stride$

Usually, the convolution layer is designed to retain the height and width of input, i.e., $O = I$

Height and width of convolution layer's output are independent of $in\_channels$ and $out\_channels$

```
nn.Conv2d(in_channels= 1,out_channels= 1, kernel_size=3, stride=1, padding=1)
```

# Parameter Count for Convolution Layers

- # of Parameters = ($kernel\_size$ * $kernel\_size$ * $in\_channels$ * $out\_channels$) + # of bias terms

- Each kernel (neuron) has one bias term
  - # of kernels = $out\_channels$
  - # of bias terms = $out\_channels$

- Parameter count is independent of $stride$ and $padding$

```
nn.Conv2d(in_channels= , out_channels= , kernel_size=3, stride=1, padding=1)
```
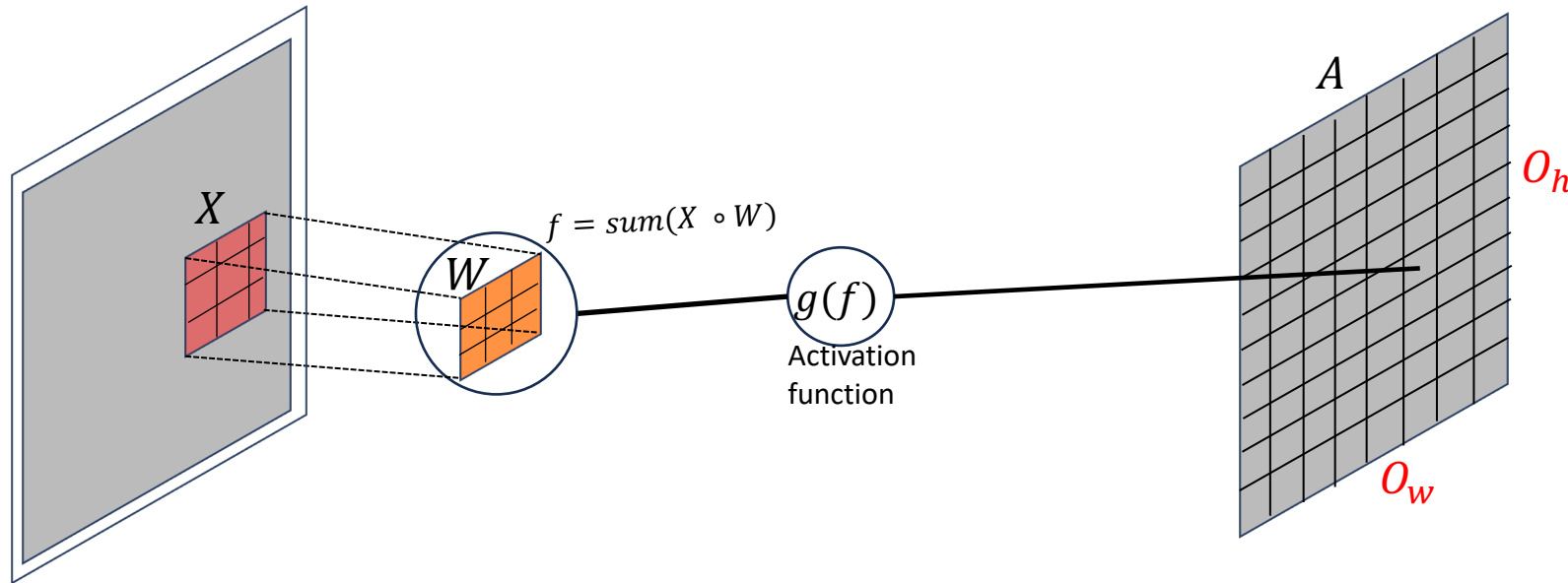
# FLOPS for Convolution Layer

- FLOPS: Floating Point Operations

Output height and width

Kernel height and width

- # of FLOPS = $2 * out\_channels * O_h * O_w * in\_channels * k_h * k_w$

1 Multiplication and 1 Addition

```
nn.Conv2d(in_channels=1, out_channels= 1, kernel_size= 3 , stride=1, padding=1)
```
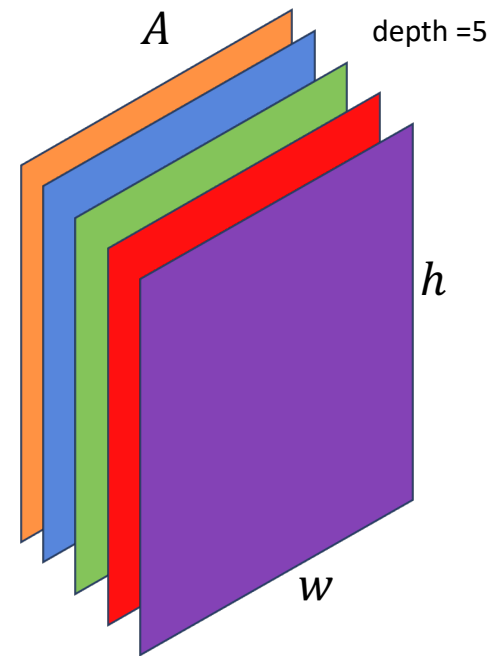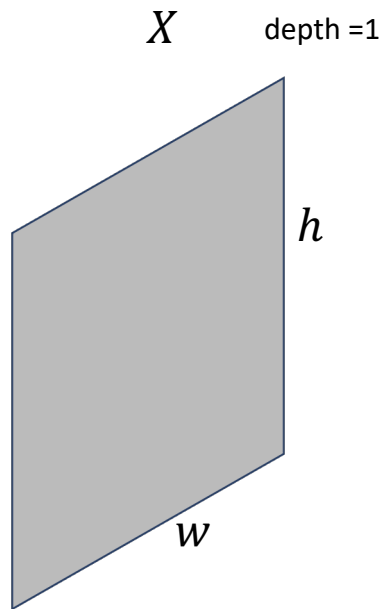


$X$

$W$

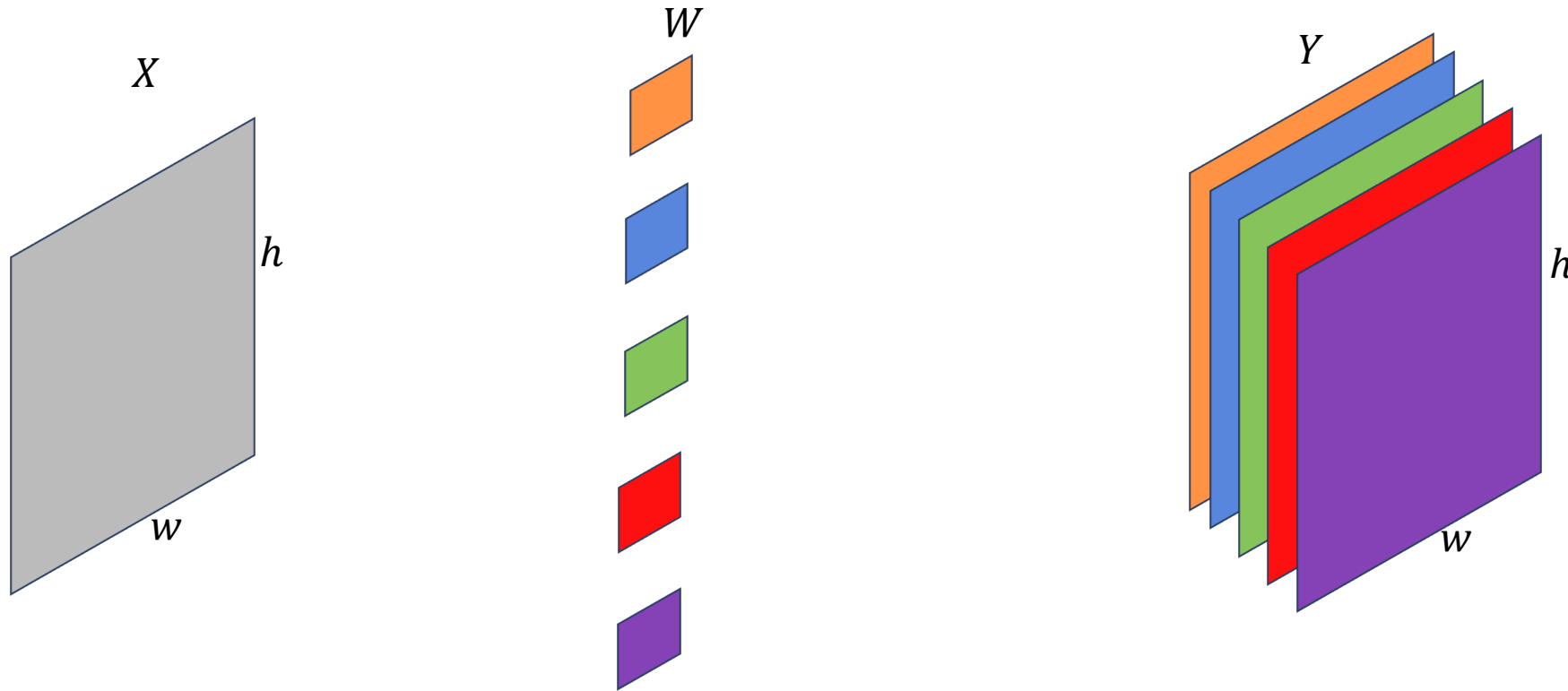$f = sum(X \circ W)$

$g(f)$
Activation
function

$A$

$O_h$

$O_w$

# How to increase depth?

# Convolution Layer: Gray Scale Image

How to increase the depth?

Multiple Kernels!!!

$X$     depth =1

$h$

$w$

$W$

$A$     depth =5

$h$

$w$

# Convolution Layer

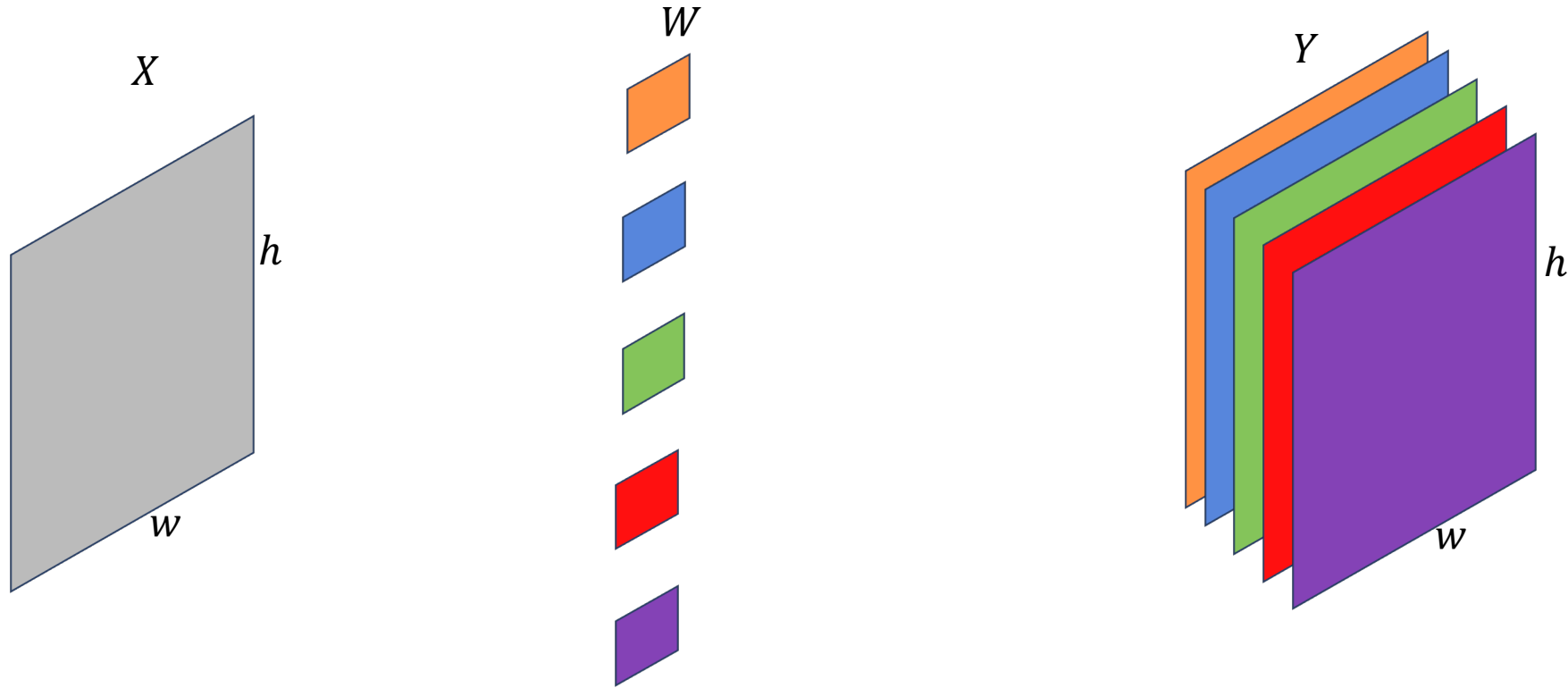

$X$     $h$     $w$

$W$

$Y$     $h$     $w$

```
nn.Conv2d(in_channels=__, out_channels= __, kernel_size= 3 , stride=1, padding=1)
```
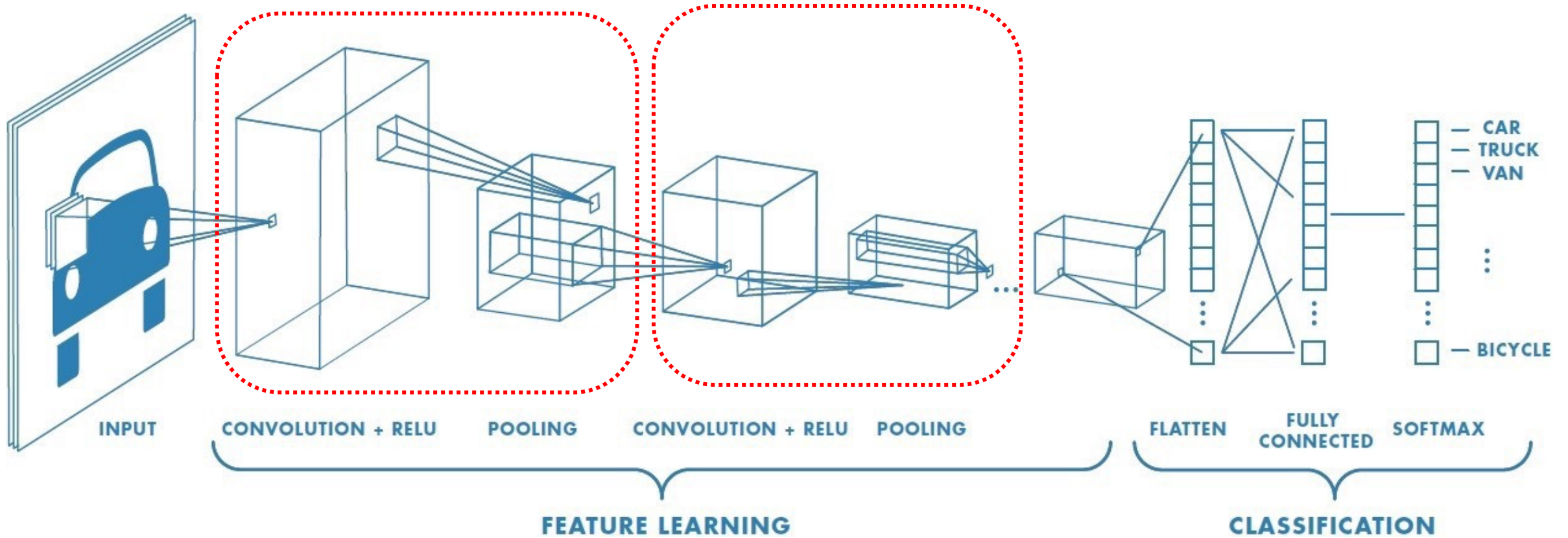
**What is the dimension of each kernel**?

# FLOPS for Convolution Layer

```
nn.Conv2d(in_channels=1, out_channels= 5, kernel_size= 3 , stride=1, padding=1)
```

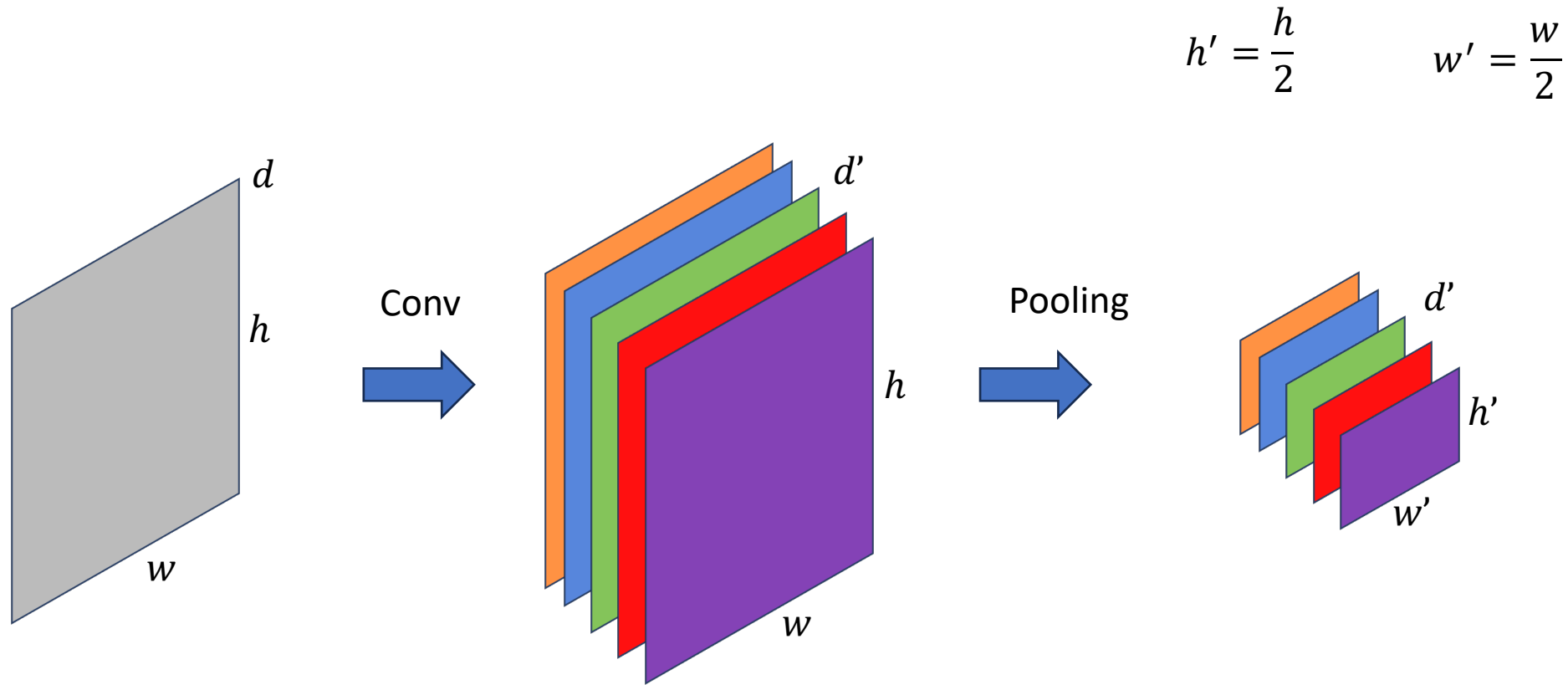# of FLOPS $= 2 * out\_channels * O_h * O_w * in\_channels * k_h * k_w$

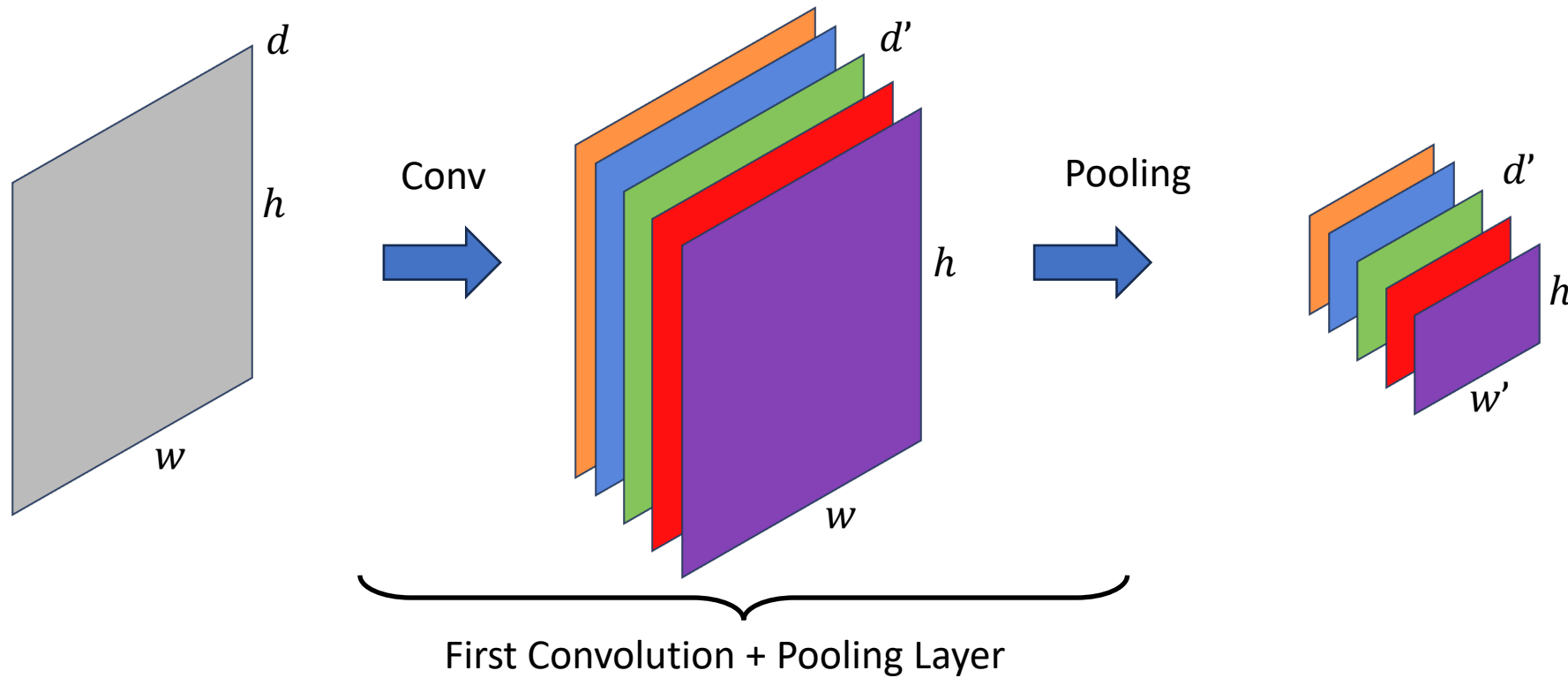How to reduce height and width?
Why should we reduce height and width?



# of FLOPS = $2 * out\_channels * O_h * O_w * in\_channels * k_h * k_w$

# How to reduce height and width?

$$h' = \frac{h}{2} \qquad w' = \frac{w}{2}$$

# How to reduce height and width?

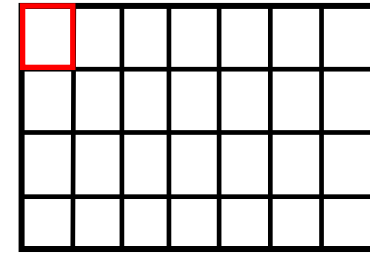$$h' = \frac{h}{2} \qquad w' = \frac{w}{2}$$



First Convolution + Pooling Layer
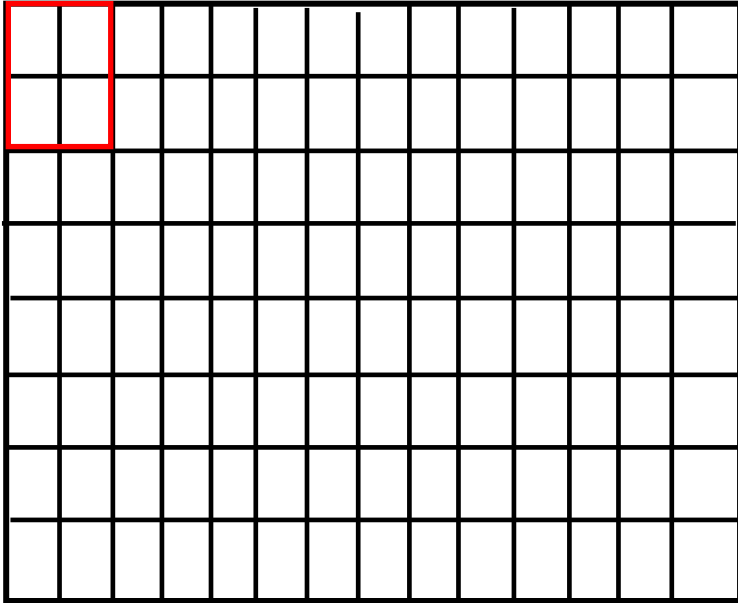
```
nn.Conv2d(in_channels=1, out_channels= 5, kernel_size= 3 , stride=1, padding=1)

max_pool_2D = nn.MaxPool2d(kernel_size = 2,stride = 2)
```

# Max Pooling
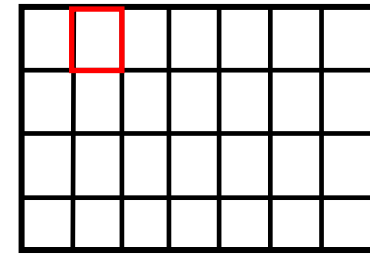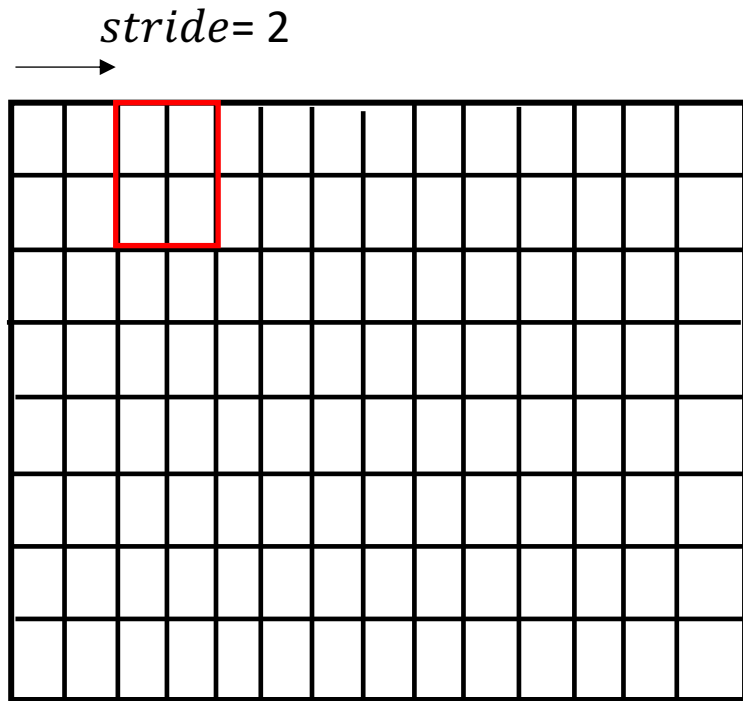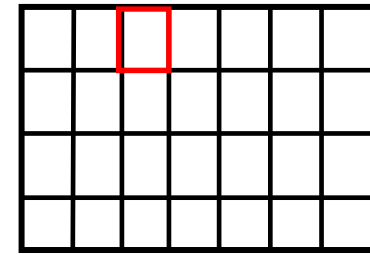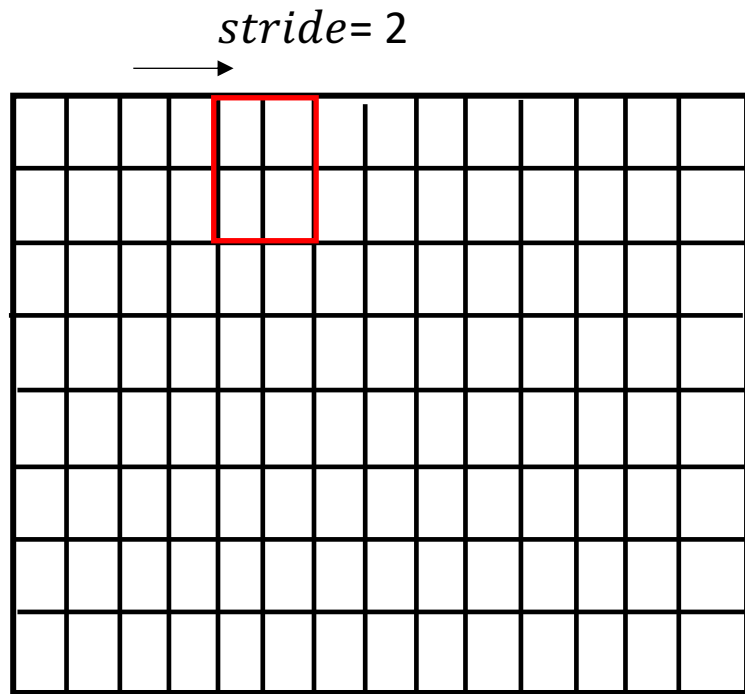
*kernel_size* = 2

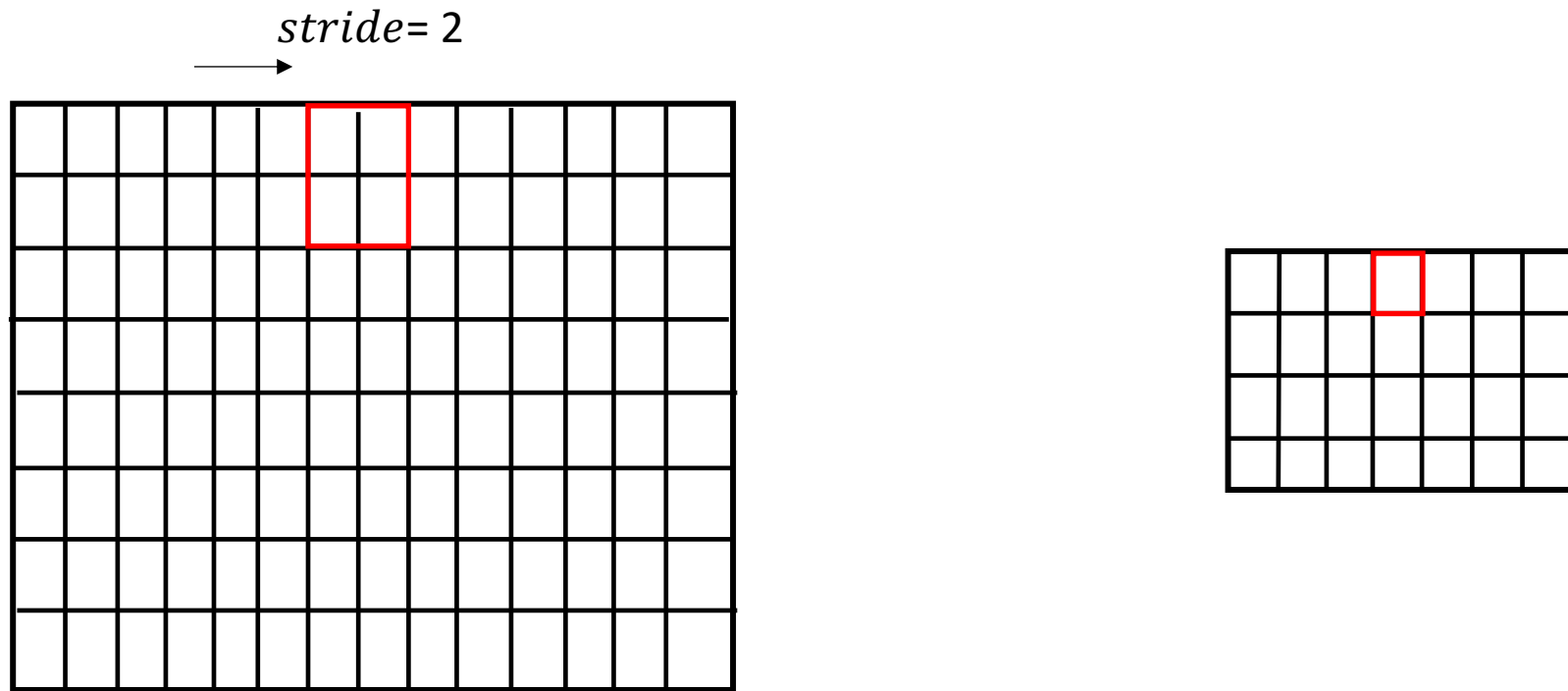max_pool_2D = nn.MaxPool2d(kernel_size = 2,stride = 2)

# Max Pooling

$stride= 2$



```
max_pool_2D = nn.MaxPool2d(kernel_size = 2,stride = 2)
```

# Max Pooling



*stride*= 2

```
max_pool_2D = nn.MaxPool2d(kernel_size = 2,stride = 2)
```

# Max Pooling

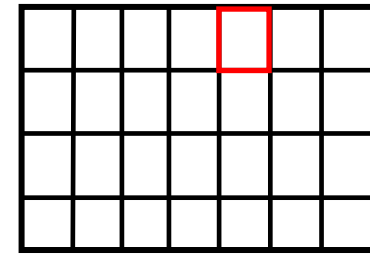$stride = 2$



```
max_pool_2D = nn.MaxPool2d(kernel_size = 2,stride = 2)
```

# Max Pooling

$stride = 2$



```
max_pool_2D = nn.MaxPool2d(kernel_size = 2,stride = 2)
```

# Max Pooling

$stride = 2$



```
max_pool_2D = nn.MaxPool2d(kernel_size = 2,stride = 2)
```
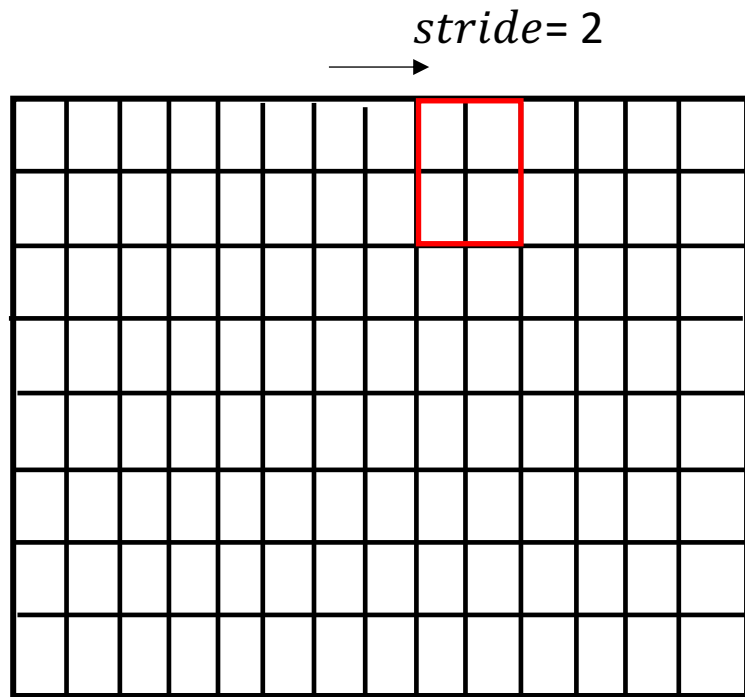
# Max Pooling

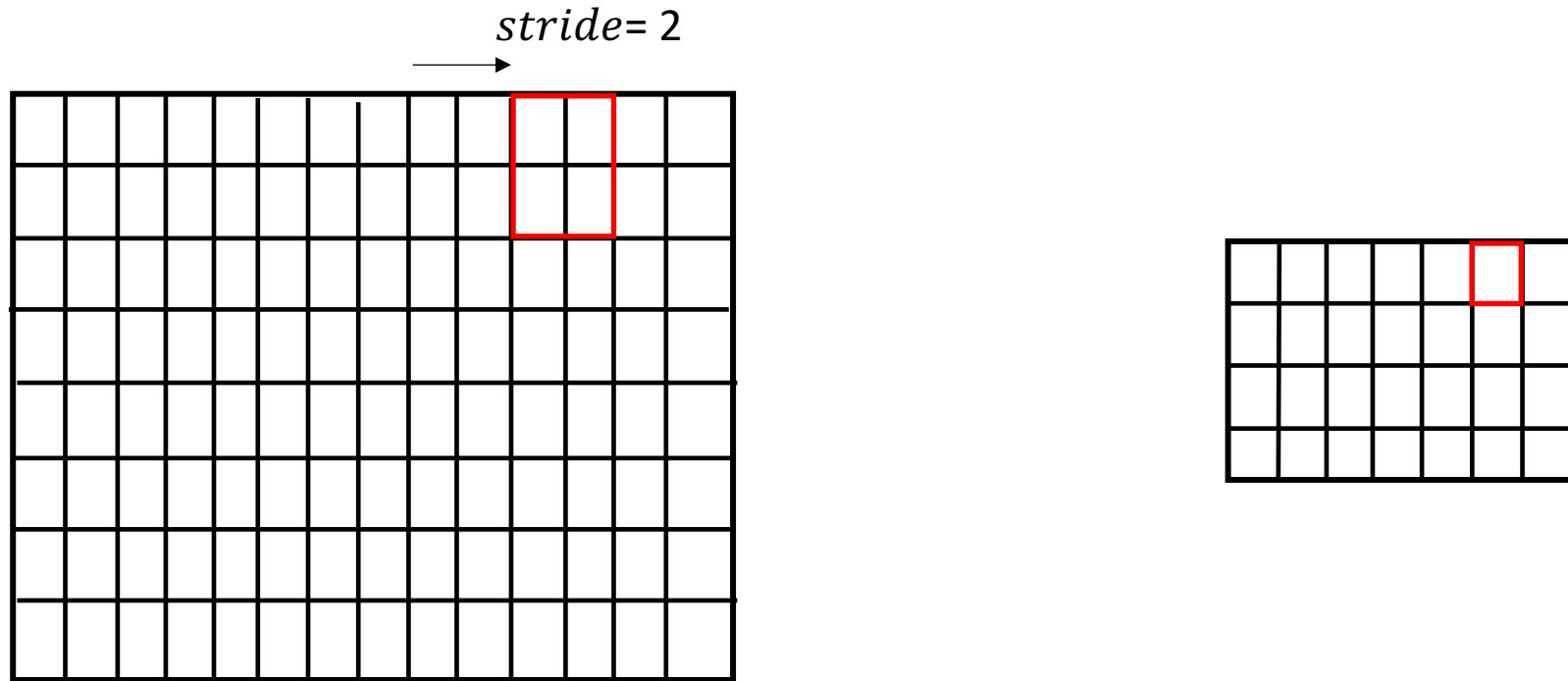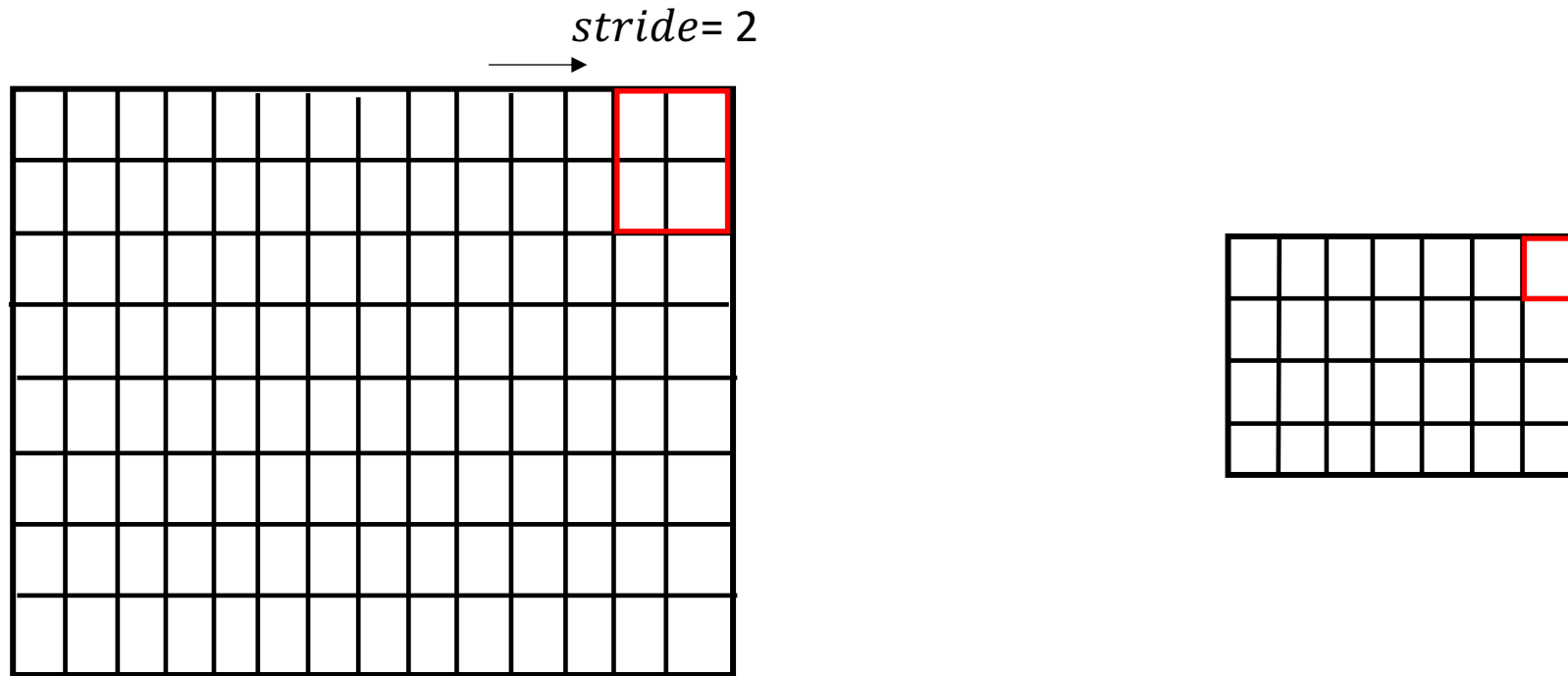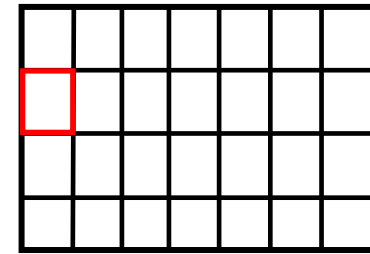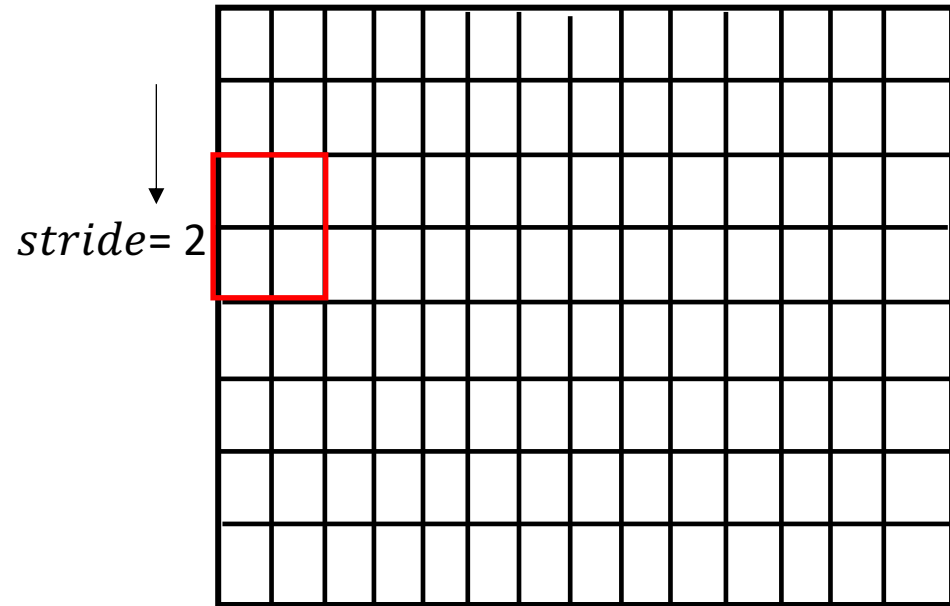$stride = 2$



```
max_pool_2D = nn.MaxPool2d(kernel_size = 2,stride = 2)
```

# Max Pooling



*stride* = 2
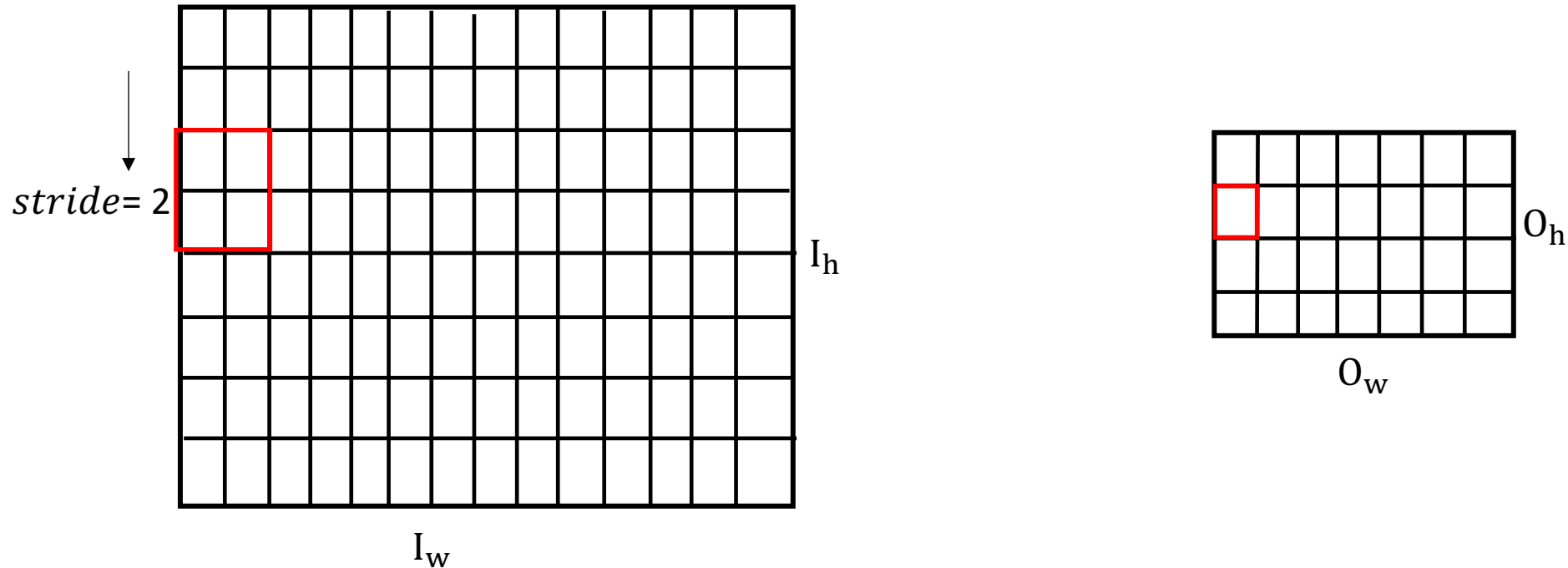
```
max_pool_2D = nn.MaxPool2d(kernel_size = 2,stride = 2)
```

# Output Size of Max Pool layer

- $O_h = \left\lfloor \dfrac{I_h + 2p - k}{s} \right\rfloor + 1$
- $O_w = \left\lfloor \dfrac{I_w + 2p - k}{s} \right\rfloor + 1$



*stride* = 2

$I_h$

$I_w$

$O_h$

$O_w$

```
max_pool_2D = nn.MaxPool2d(kernel_size = 2,stride = 2)
```

# FLOPS for MaxPool Layer

Output width

- FLOPS = $in\_channels * O_h * O_w * (k^2 - 1)$

  # of channels     Output height     # of comparisons

- $k$ is the kernel size



$O_h$

$O_w$

FLOPS: Floating Point Operations

Second convolution layer???



INPUT
CONVOLUTION + RELU
POOLING
CONVOLUTION + RELU
POOLING
FLATTEN
FULLY CONNECTED
SOFTMAX

— CAR
— TRUCK
— VAN

— BICYCLE

FEATURE LEARNING
CLASSIFICATION

# Adding another Conv+Pooling Layer

$$w' = \frac{w}{2} \qquad w'' = \frac{w'}{2}$$

$$h' = \frac{h}{2} \qquad h'' = \frac{w'}{2}$$



Second convolution layer???
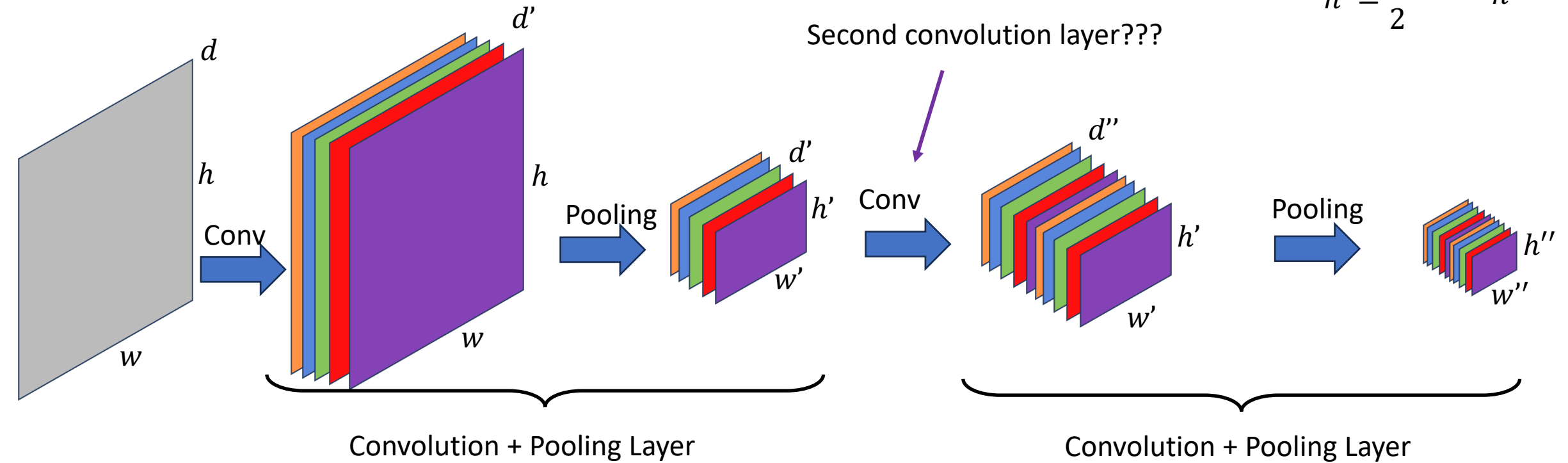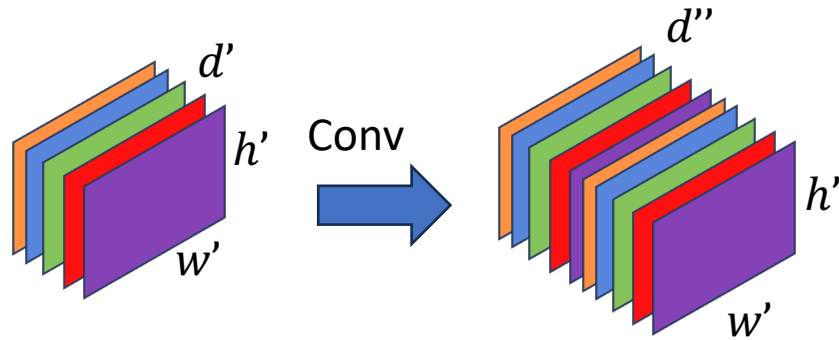
Conv

Pooling

Conv

Pooling

Convolution + Pooling Layer

Convolution + Pooling Layer

31

# Adding another Conv + Pooling Layer



$$w' = \frac{w}{2}$$
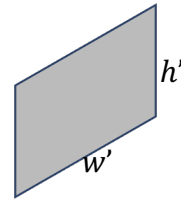
$$h' = \frac{h}{2}$$

**Questions**:
How many kernels?
What is the shape of each kernel?

# Adding another Conv+Pooling Layer

`nn.Conv2d(in_channels=__, out_channels= __, kernel_size= 3 , stride=1, padding=1)`

Shape of the kernel:_____

# Adding another Conv+Pooling Layer

```
nn.Conv2d(in_channels=__, out_channels= __, kernel_size= 3 , stride=1, padding=1)
```



$$f = sum(X \circ W)$$

# Adding another Conv+Pooling Layer

`nn.Conv2d(in_channels=__, out_channels= __, kernel_size= 3 , stride=1, padding=1)`



$$f = sum(X \circ W)$$

# Adding another Conv+Pooling Layer

```
nn.Conv2d(in_channels=__, out_channels= __, kernel_size= 3 , stride=1, padding=1)
```

# Adding another Conv+Pooling Layer

`nn.Conv2d(in_channels=__, out_channels= __, kernel_size= 3 , stride=1, padding=1)`

# Adding another Conv+Pooling Layer

```
nn.Conv2d(in_channels=__, out_channels= __, kernel_size= 3 , stride=1, padding=1)
```



$f = sum(X \circ W)$

# Adding another Conv+Pooling Layer

$$w' = \frac{w}{2}$$

$$h' = \frac{h}{2}$$



Conv

```
nn.Conv2d(in_channels=__, out_channels= __, kernel_size= 3 , stride=1, padding=1)
```
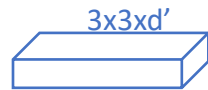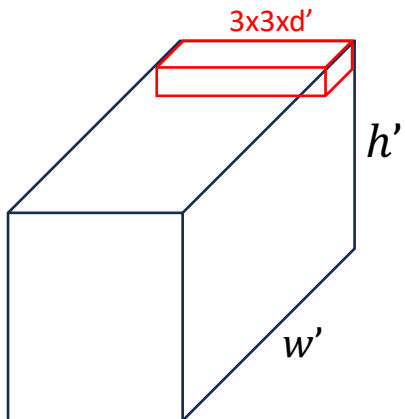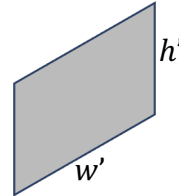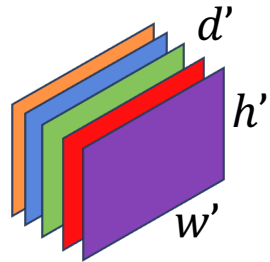
**Questions**:
How many kernels?
What is the shape of each kernel?

# Adding another Conv+Pooling Layer

```
nn.Conv2d(in_channels=__, out_channels= __, kernel_size= 3 , stride=1, padding=1)
```
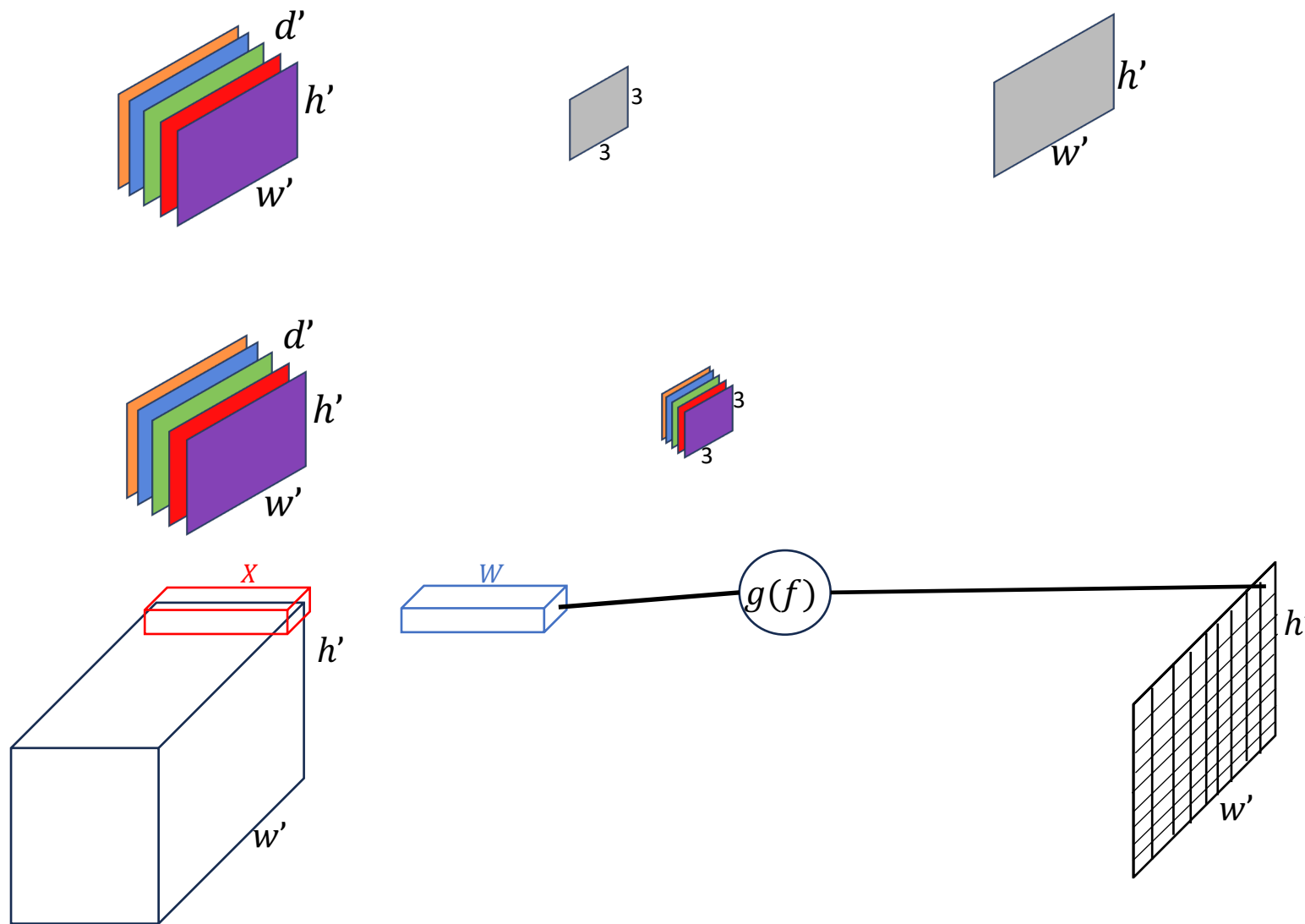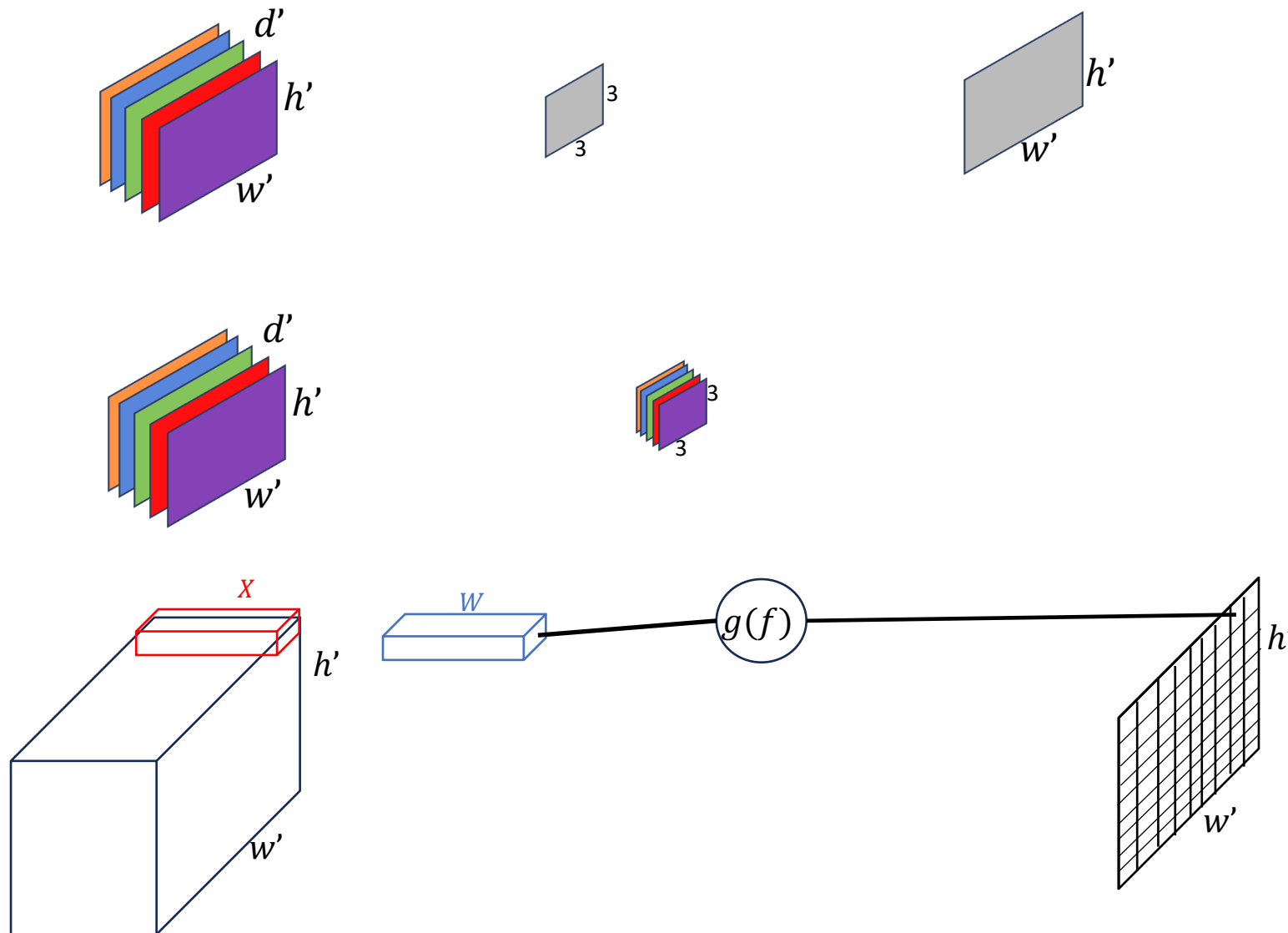
$d''$ kernels

$d''$ layers, one for each kernel

Convolution

$h'$

$w'$

$h'$

$w'$

# Adding another Conv+Pooling Layers

$$w' = \frac{w}{2}$$

$$h' = \frac{h}{2}$$



$d' = 5d$

Conv

Pooling

Conv

Convolution + Pooling Layer

**Questions**:
How many FLOPS?

# Adding another Conv+Pooling Layer

$$w' = \frac{w}{2} \qquad w'' = \frac{w'}{2}$$

$$h' = \frac{h}{2} \qquad h'' = \frac{w'}{2}$$



Convolution + Pooling Layer          Convolution + Pooling Layer

```
nn.MaxPool2d(kernel_size=__, stride=__)
```

# Adding another Conv+Pooling Layer



$d' = 5d$

$w'' = \dfrac{w'}{2}$

$h'' = \dfrac{w'}{2}$

Conv

Pooling

Conv

Pooling

Convolution + Pooling Layer

Convolution + Pooling Layer

```
nn.Conv2d(in_channels=__, out_channels= __, kernel_size= 3 , stride=1, padding=1)

nn.MaxPool2d(kernel_size=__, stride=__)
```
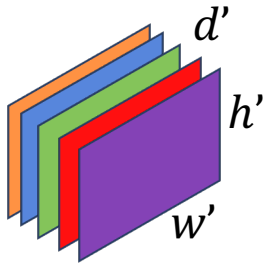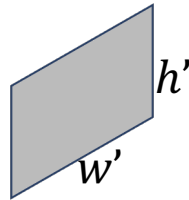
INPUT    CONVOLUTION + RELU    POOLING    CONVOLUTION + RELU    POOLING    FLATTEN    FULLY CONNECTED    SOFTMAX

CAR
TRUCK
VAN
BICYCLE

FEATURE LEARNING      CLASSIFICATION

# Flatten + FC Layers



What is the size of flattened vector?

Think of flatten vector as feature vector
Fully Connected (FC) layer is traditional neural network layer with predefined number of neurons

# Flattened Layer + FC Layer

Previous layer is flattened to vector

FC Layer

$d''$

$h''$

$w''$

$$Length = w'' * h'' * d''$$

Length of the flattened layer is equal to the size of the input for FC layer

Design of FC layers require the knowledge of size of final Conv+Pooling layer

# Flattened Layer + FC Layer(s)

Previous layer is flattened to vector

FC Layer (1)

FC Layer (2)

FC Layer (3)

Output Layer

$d''$

$h''$

$w''$

$Length = w'' * h'' * d''$

Deep Neural Network

Length of the flattened layer is equal to the size of the input for FC layer

Design of FC layers require the knowledge of size of final Conv+Pooling layer

# How many FLOPS per layer?

- FLOPS: Floating Point Operations

```
nn.Conv2d(in_channels=5, out_channels= 10, kernel_size= 3 , stride=1, padding=1)
```



10 kernels

10 layers, one for each kernel

$h' = w' = 16$

$d' = 5$

# How many FLOPS per layer?

`nn.Conv2d(in_channels=5, out_channels= 10, kernel_size= 3 , stride=1, padding=1)`

10 kernels

10 layers (or out_channels), one for each kernel

$d'$

$h'$

$w'$

Convolution

$h'$

$w'$

...

$h' = w' = 16 \qquad d' = 5$

$(out\_channels * out\_height * out\_width) * 2 * (in\_channels * kernel\_height * kernel\_width)$

# of pixels in output

(# of multiplications + # of additions ) per pixel

# Convolution Layer: RGB Input

`nn.Conv2d(in_channels=___, out_channels=__, kernel_size= 3 , stride=1, padding=1)`

$X$

$W$

$A$

Shape of kernel:_____

# Convolution Layer: RGB Input

```
nn.Conv2d(in_channels=__, out_channels=__, kernel_size= 3 , stride=1, padding=1)
```



$X$

$d$ kernels

$d$ layers, one for each kernel

Assuming kernel size of 3 for each kernel, how many parameters in the convolution layer?

# CNN Playground

- https://cs.stanford.edu/people/karpathy/convnetjs/demo/cifar10.html

# Types of Convolution

- Standard Convolution


- Depthwise Convolution


- Depthwise separable Convolution

# Standard Convolution



3x3xd'

$h'$

$w'$

Input

3x3xd'

Kernel

Output

# Recap: How many FLOPS per layer?

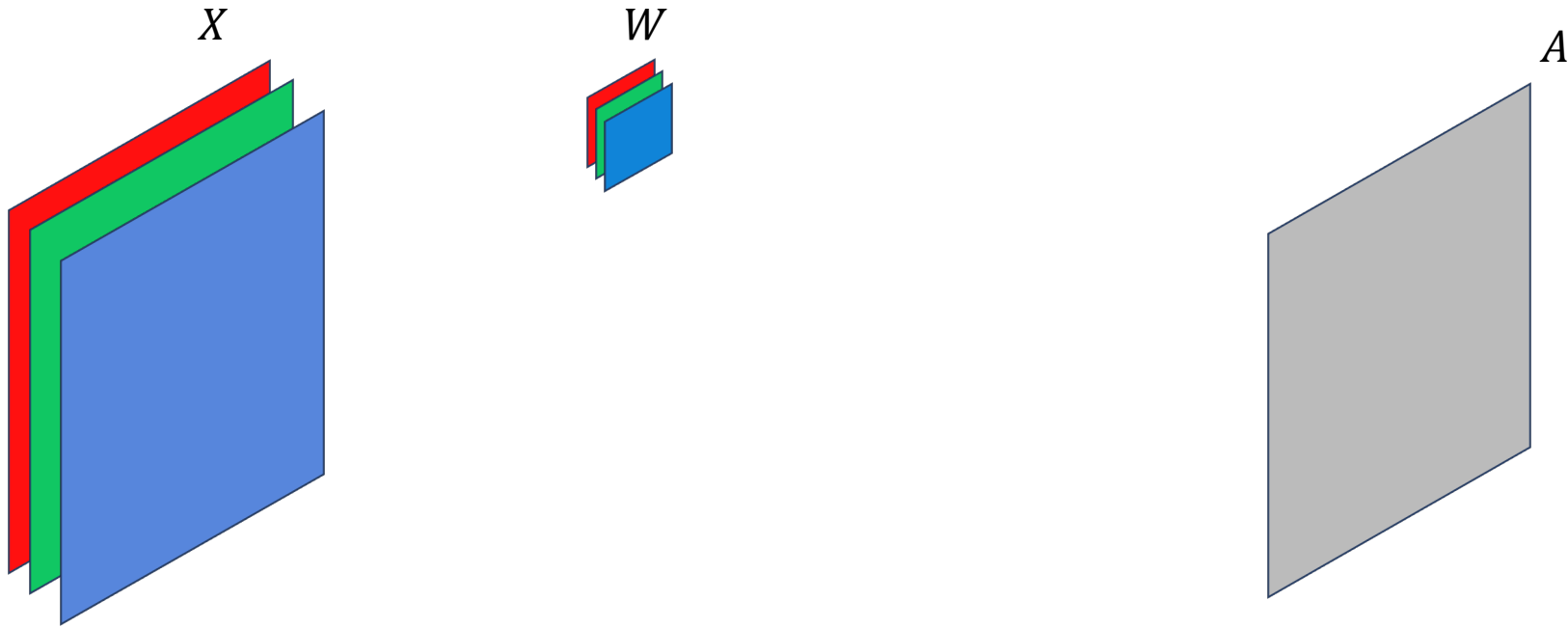`nn.Conv2d(in_channels=5, out_channels= 10, kernel_size= 3 , stride=1, padding=1)`

10 kernels

10 layers (or out_channels), one for each kernel

$d'$

$h'$

$w'$

Convolution

...

$h'$

$w'$

$h' = w' = 16$     $d' = 5$

$(out\_channels * out\_height * out\_width) * 2 * (in\_channels * kernel\_height * kernel\_width)$

\# of pixels in output

(\# of multiplications + \# of additions ) per pixel

# Depthwise Convolutions

d'

$h'$

$w'$

Input

3x3xd'

d' Kernels one for each input layer

d' output feature maps, one for each kernel

Output

# Depthwise Convolutions: # of FLOPs

Input

$3 \times 3 \ kernels \ of \ shape \ 3 \times 3 \times 1$

# of kernels: $d'$
one kernel for each input layer

$d'$ output feature maps, one for each kernel

Output

$(out\_channels * out\_height * out\_width) * 2 * (kernel\_height * kernel\_width)$

# Depthwise Separable Convolutions

No interaction between the layers with depthwise convolutions

$3\times3$ *kernels of shape* $3\times3\times1$

# of kernels: $d'$
one kernel for each input layer

d' output feature maps, one for each kernel

Input

Output

# Depthwise Separable Convolutions

depthwise convolutions

pointwise convolutions

$3\times3$ *kernels of shape* $3\times3\times1$

$d'$

$h'$

$w'$

# of kernels: $d'$

$1\times1\times d'$

$d'$ output feature maps, one for each kernel

Input

How to get multiple output feature maps at pointwise convolution layer?

$(out\_channels * out\_height * out\_width) * 2 * (kernel\_height * kernel\_width) + 2*(out\_height * out\_width * in\_channels)$

# Depthwise Separable Convolutions

depthwise convolutions

pointwise convolutions

$3\times3$ *kernels of shape* $3\times3\times1$

$d'$

$h'$

$w'$

# of kernels: $d'$

$d'$ output feature maps, one for each kernel

$1\times1\times d'$

$1\times1\times d'$

$1\times1\times d'$

Input

How to get multiple output feature maps at pointwise convolution layer?
Use multiple 1x1 kernels

$(out\_channels * out\_height * out\_width) * 2 * (kernel\_height * kernel\_width)$

# Mobilenets

**MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications**

Andrew G. Howard          Menglong Zhu          Bo Chen          Dmitry Kalenichenko
Weijun Wang          Tobias Weyand          Marco Andreetto          Hartwig Adam

Google Inc.

{howarda,menglong,bochen,dkalenichenko,weijunw,weyand,anm,hadam}@google.com

Key Idea: Depthwise Convolutions

https://arxiv.org/pdf/1704.04861

# CNN Architectures

- AlexNet

- VGG-16

- ResNet-50
  - Residual connections
  - Bottleneck layer to reduce parameter count

- MobileNetV2
  - Depthwise Convolution

# LetNet5 Architecture: Tutorial!!



FC (10) — `nn.Linear(in_features =__, out_features=__)`

FC (84) — `nn.Linear(in_features=__, out_features=__)`

FC (120) — `nn.Linear(in_features=__, out_features=__)`

2 × 2 AvgPool, stride 2 — `nn.AvgPool2d(kernel_size=__, stride=__)`

5 × 5 Conv (16) 16 kernels — `nn.Conv2d(in_channels=__, out_channels=__, kernel_size=__, stride=2, padding=2)`

2 × 2 AvgPool, stride 2 — `nn.AvgPool2d(kernel_size=__, stride=__)`

6 kernels
5 × 5 Conv (6), pad 2 — `nn.Conv2d(in_channels=__, out_channels=__, kernel_size=__, stride=1, padding=2)`

Image (28 × 28)

Image taken from: https://en.wikipedia.org/wiki/LeNet

# VGGNet: Homework

- Visual Geometry Group (VGG) at the University of Oxford



224×224×3  224×224×64
112×112×128
56×56×256
28×28×512
14×14×512
7×7×512
1×1×4096  1×1×1000

- convolution+ReLU
- max pooling
- fully connected+ReLU
- softmax

Table 1: **ConvNet configurations** (shown in columns). The depth of the configurations increases from the left (A) to the right (E), as more layers are added (the added layers are shown in bold). The convolutional layer parameters are denoted as "conv⟨receptive field size⟩-⟨number of channels⟩". The ReLU activation function is not shown for brevity.
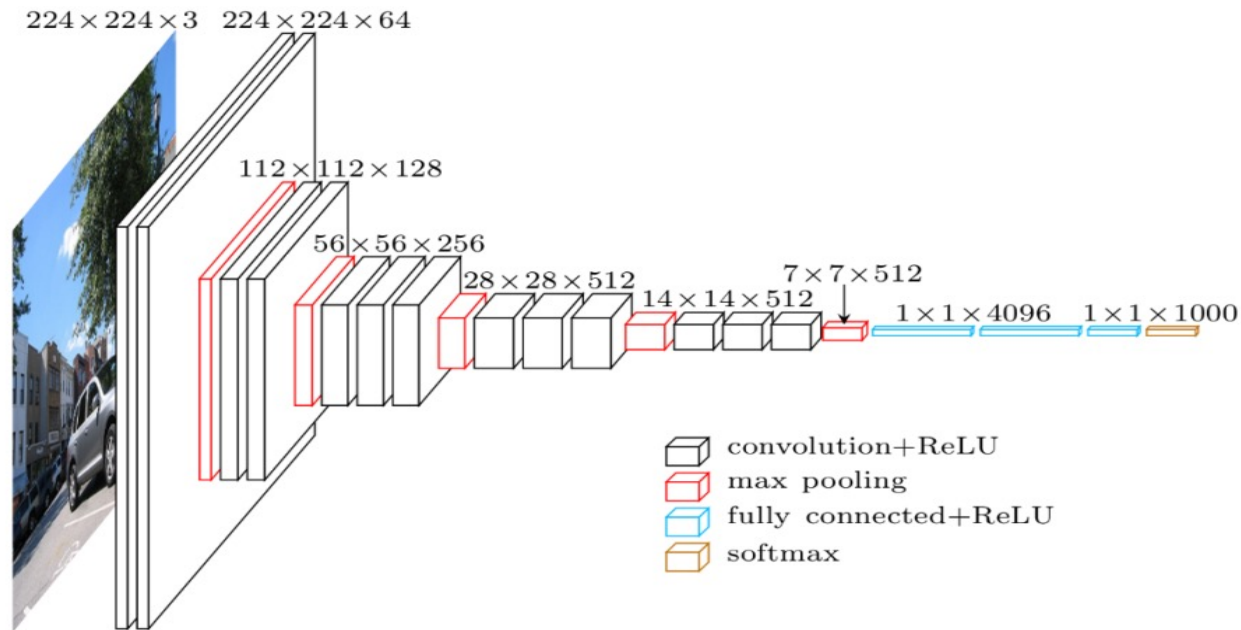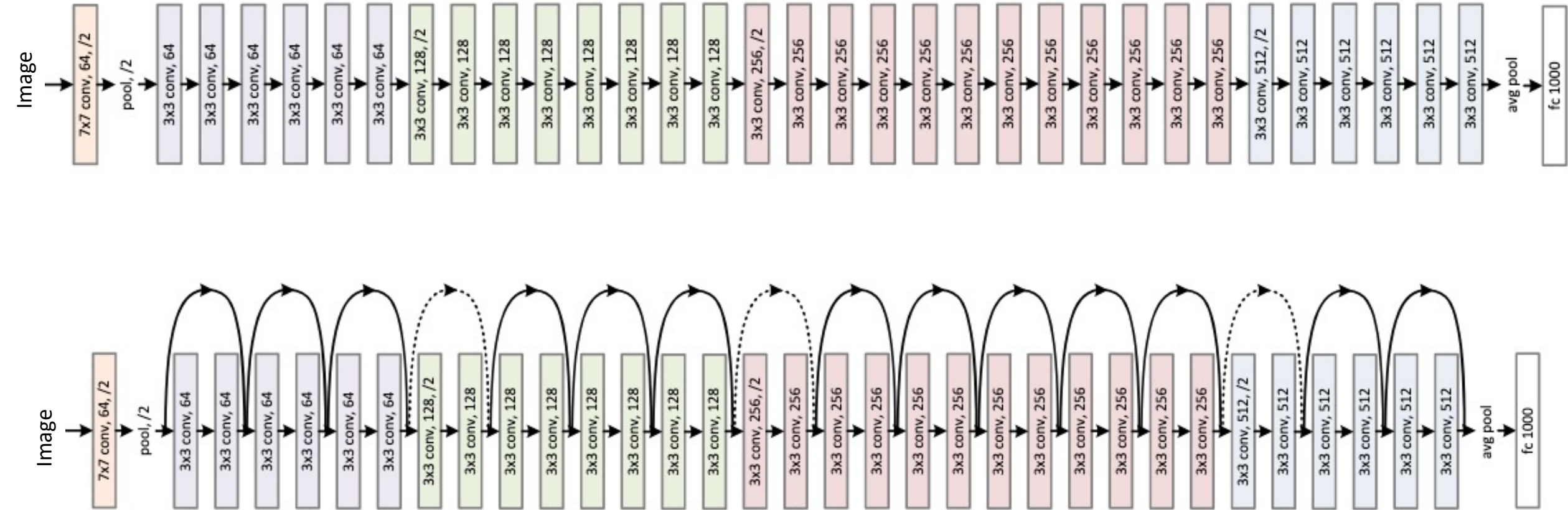
| ConvNet Configuration | | | | | |
|---|---|---|---|---|---|
| A | A-LRN | B | C | D | E |
| 11 weight layers | 11 weight layers | 13 weight layers | 16 weight layers | 16 weight layers | 19 weight layers |
| input (224 × 224 RGB image) | | | | | |
| conv3-64 | conv3-64 **LRN** | conv3-64 **conv3-64** | conv3-64 conv3-64 | conv3-64 conv3-64 | conv3-64 conv3-64 |
| maxpool | | | | | |
| conv3-128 | conv3-128 | conv3-128 **conv3-128** | conv3-128 conv3-128 | conv3-128 conv3-128 | conv3-128 conv3-128 |
| maxpool | | | | | |
| conv3-256 conv3-256 | conv3-256 conv3-256 | conv3-256 conv3-256 | conv3-256 conv3-256 **conv1-256** | conv3-256 conv3-256 **conv3-256** | conv3-256 conv3-256 conv3-256 **conv3-256** |
| maxpool | | | | | |
| conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 **conv1-512** | conv3-512 conv3-512 **conv3-512** | conv3-512 conv3-512 conv3-512 **conv3-512** |
| maxpool | | | | | |
| conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 **conv1-512** | conv3-512 conv3-512 **conv3-512** | conv3-512 conv3-512 conv3-512 **conv3-512** |
| maxpool | | | | | |
| FC-4096 | | | | | |
| FC-4096 | | | | | |
| FC-1000 | | | | | |
| soft-max | | | | | |

79

https://arxiv.org/pdf/1409.1556v6

# ResNet: Homework
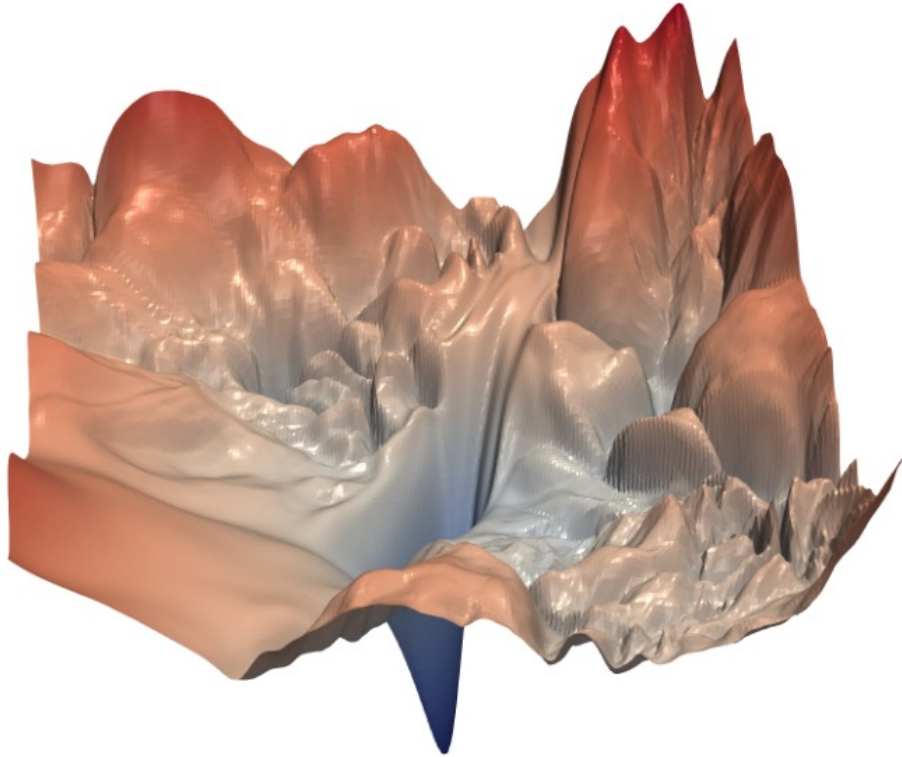
https://arxiv.org/abs/1512.03385

# Skip Connection: Tutorial!!

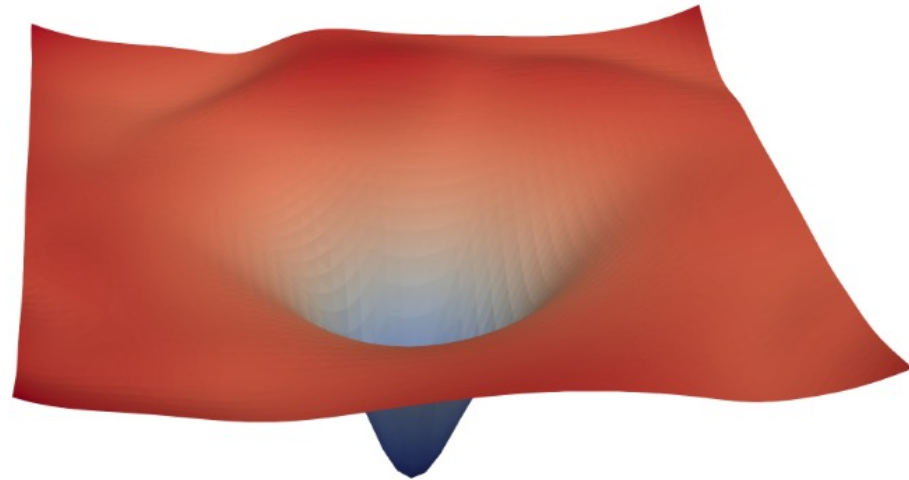aka Residual Connection



Basic building block of residual learning

# Skip Connections



(a) without skip connections

(b) with skip connections

82

Li, H., Xu, Z., Taylor, G., Studer, C. and Goldstein, T., 2018. Visualizing the loss landscape of neural nets. *Advances in neural information processing systems, 31*.