# Convolutional Neural Networks (CNN)

**Prof. Seungchul Lee**

**Industrial AI Lab.**

# Machine Learning vs. Deep Learning

- Machine learning
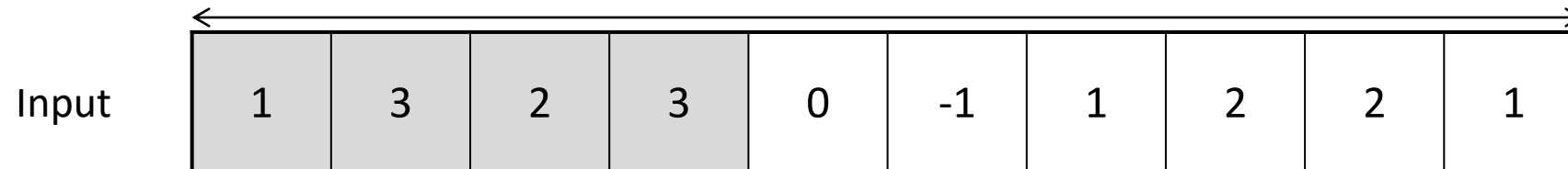


- Deep learning

# Convolution

# 1D Convolution

- (actually cross-correlation)

Input

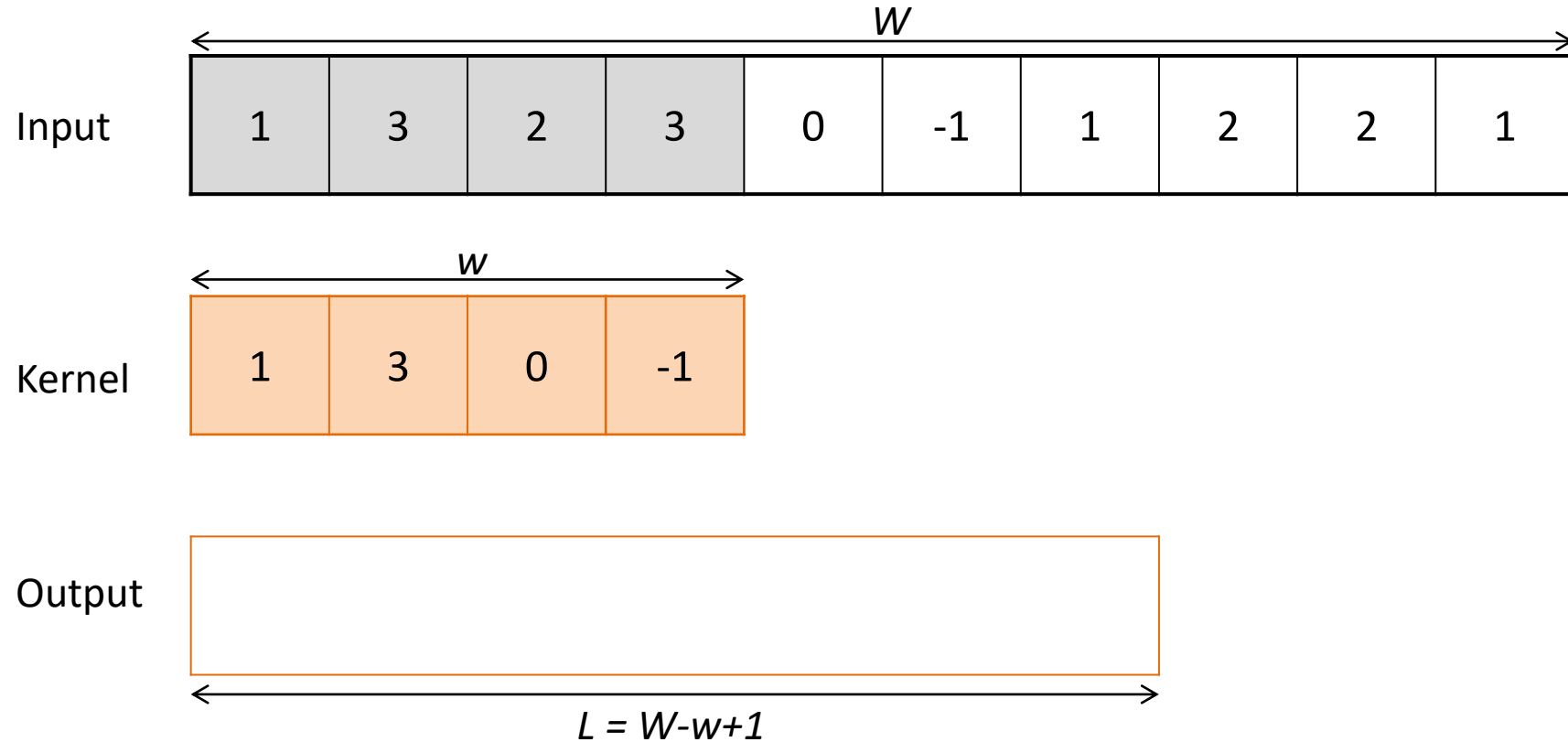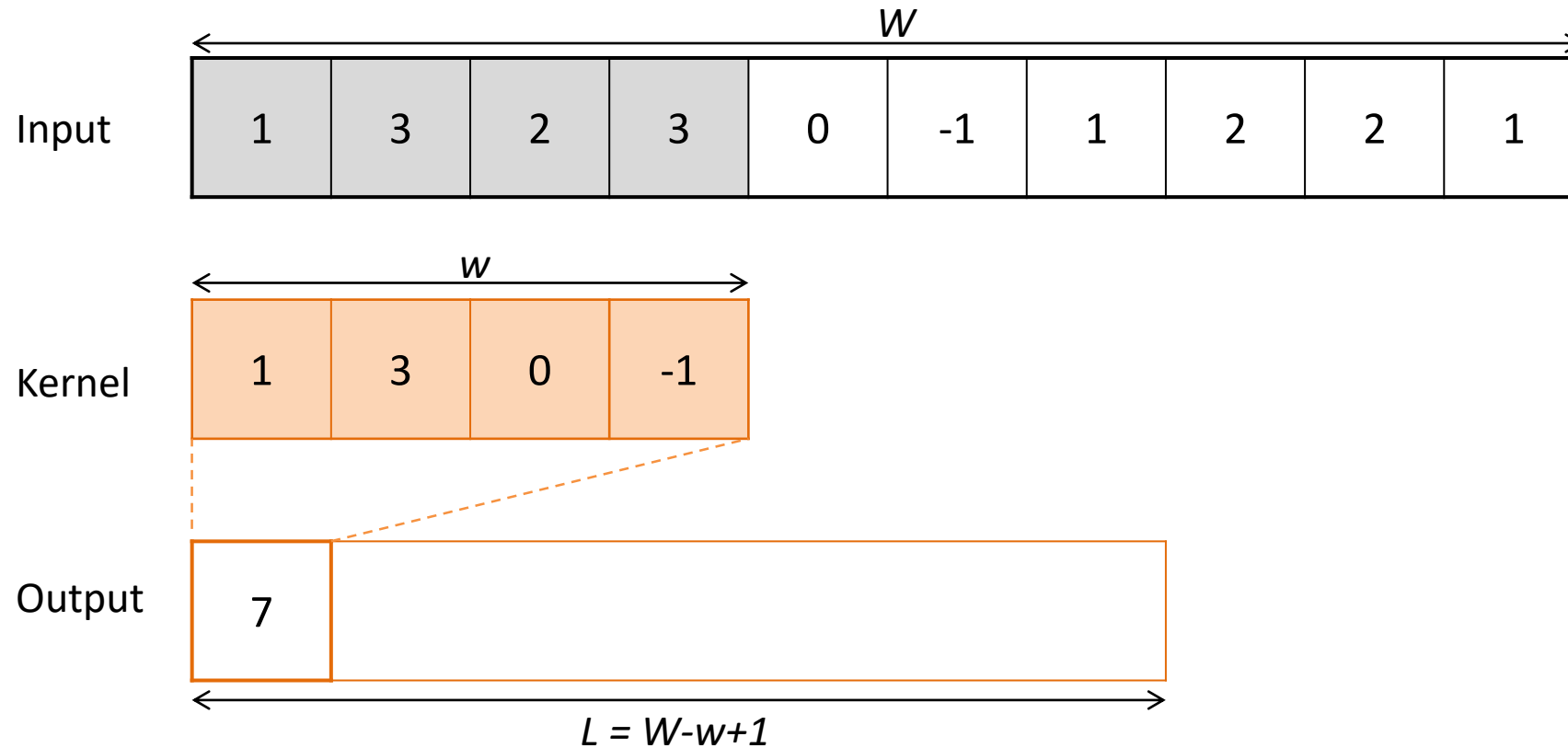| 1 | 3 | 2 | 3 | 0 | -1 | 1 | 2 | 2 | 1 |
|---|---|---|---|---|----|---|---|---|---|

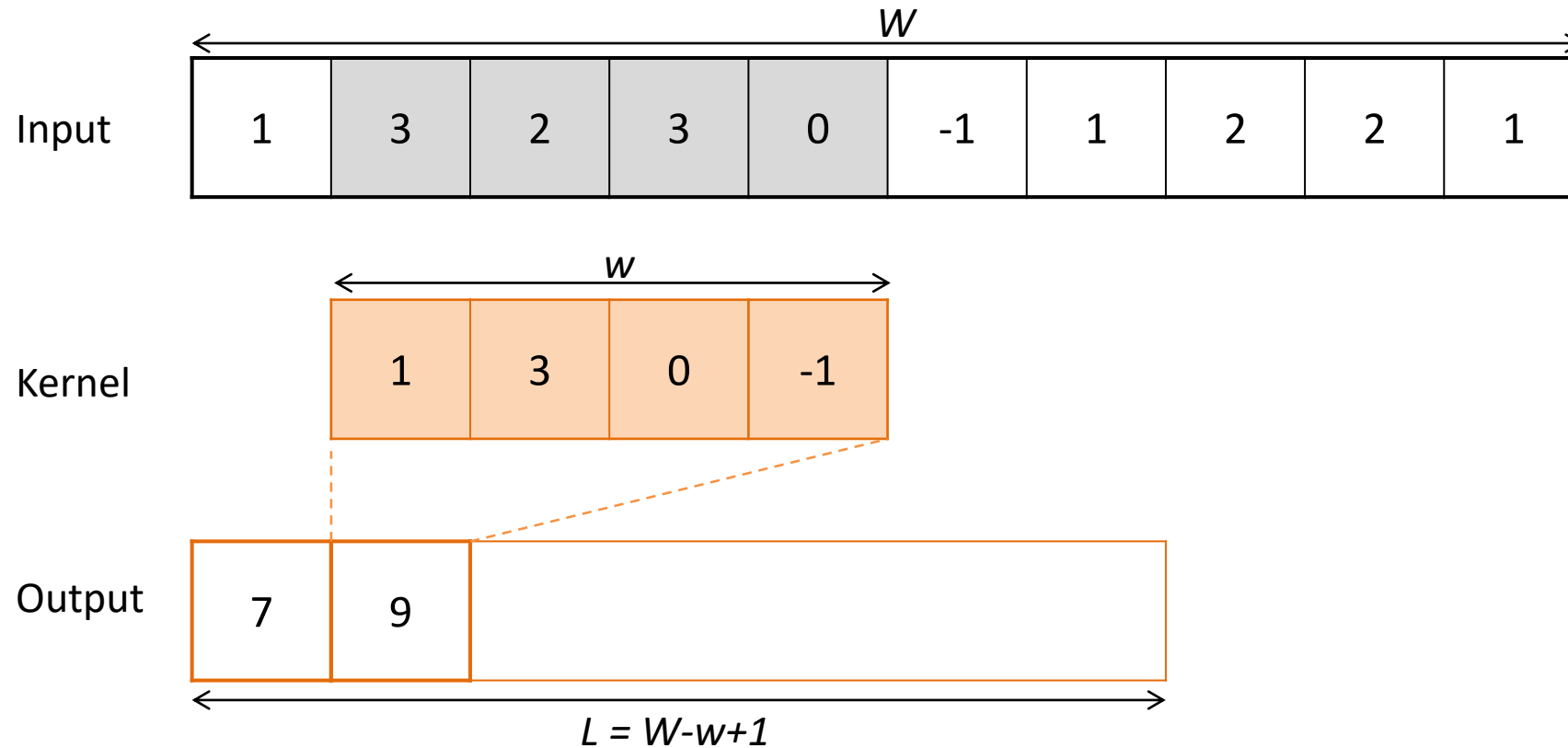# 1D Convolution

- (actually cross-correlation)

# 1D Convolution

- (actually cross-correlation)

# 1D Convolution

- (actually cross-correlation)

# 2D Convolution

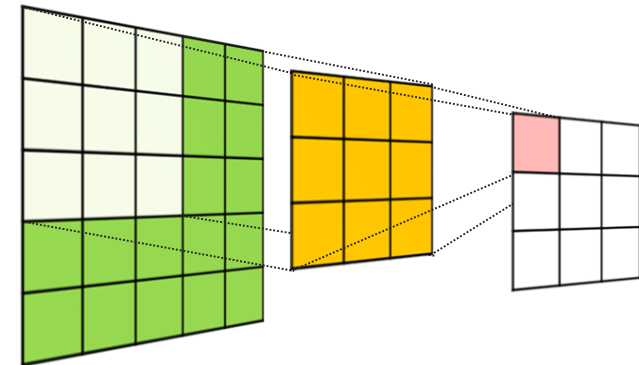# Convolution on Image (= Convolution in 2D)

- Filter (or Kernel)
  - Discrete convolution can be viewed as **element-wise multiplication** by a matrix
  - Modify or enhance an image by filtering
  - Filter images to emphasize certain features or remove other features
  - Filtering includes smoothing, sharpening and edge enhancement



Image

Convolved Feature

Image    Kernel    Output

# Convolution Mask + Neural Network



fully-connected unit

# Locality



- Locality: objects tend to have a local spatial support
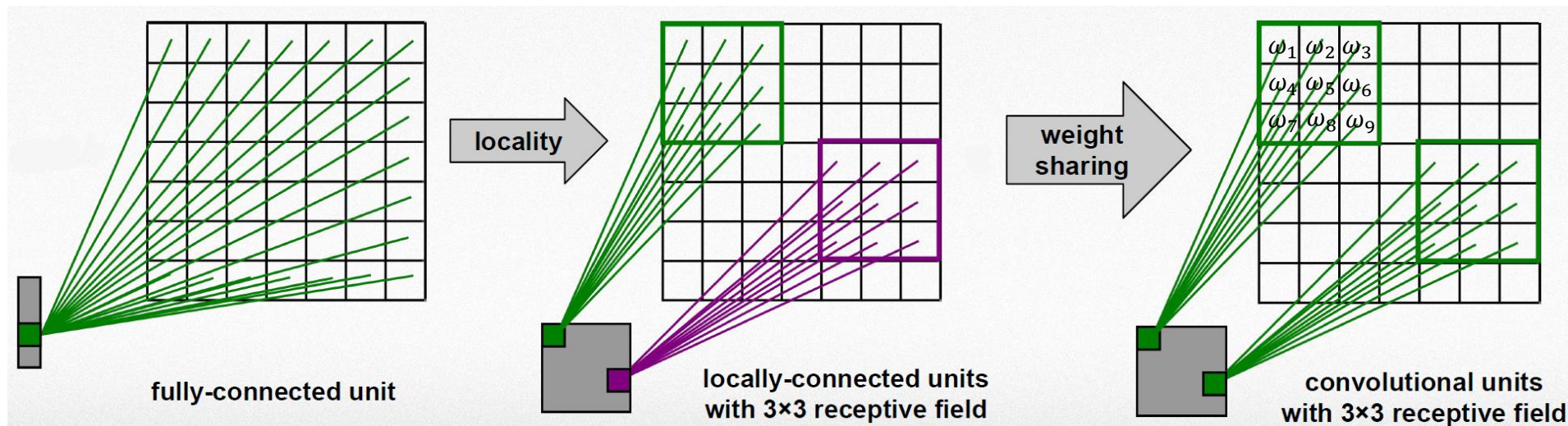  - fully-connected layer → locally-connected layer

# Locality
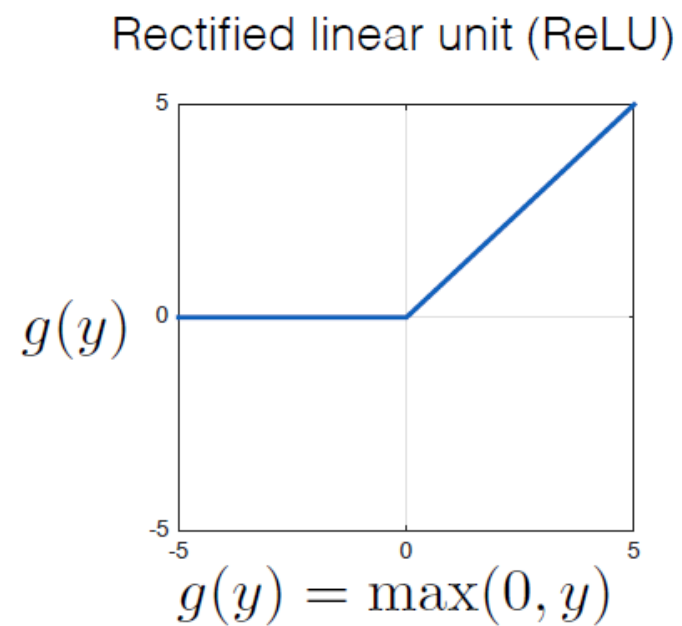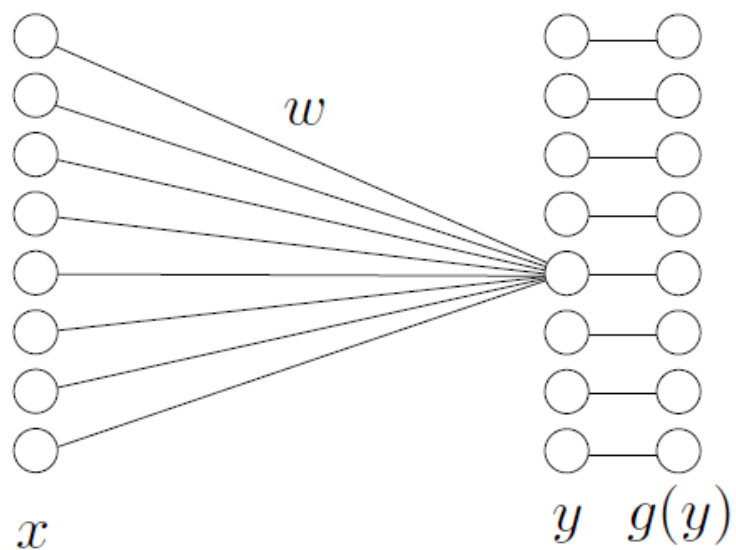


- Locality: objects tend to have a local spatial support
  - fully-connected layer → locally-connected layer

We are not designing the kernel, but are learning the kernel from data
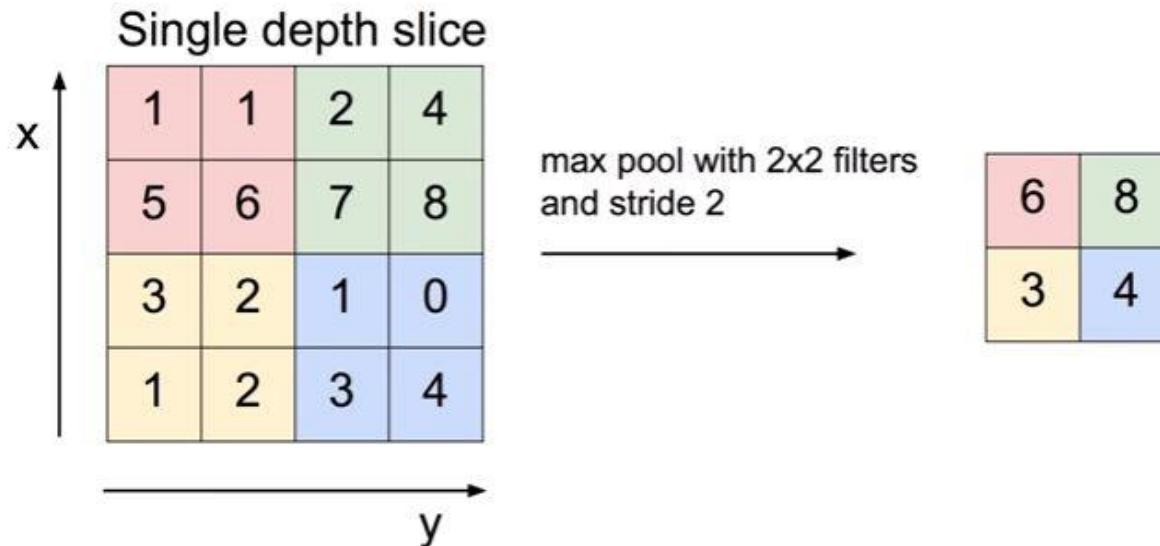→ Learning feature extractor from data



fully-connected unit

locality

locally-connected units
with 3×3 receptive field

weight
sharing

$\omega_1$ $\omega_2$ $\omega_3$
$\omega_4$ $\omega_5$ $\omega_6$
$\omega_7$ $\omega_8$ $\omega_9$

convolutional units
with 3×3 receptive field

# Nonlinear Activation Function



$x$ $\quad$ $w$ $\quad$ $y$ $\quad$ $g(y)$

Rectified linear unit (ReLU)

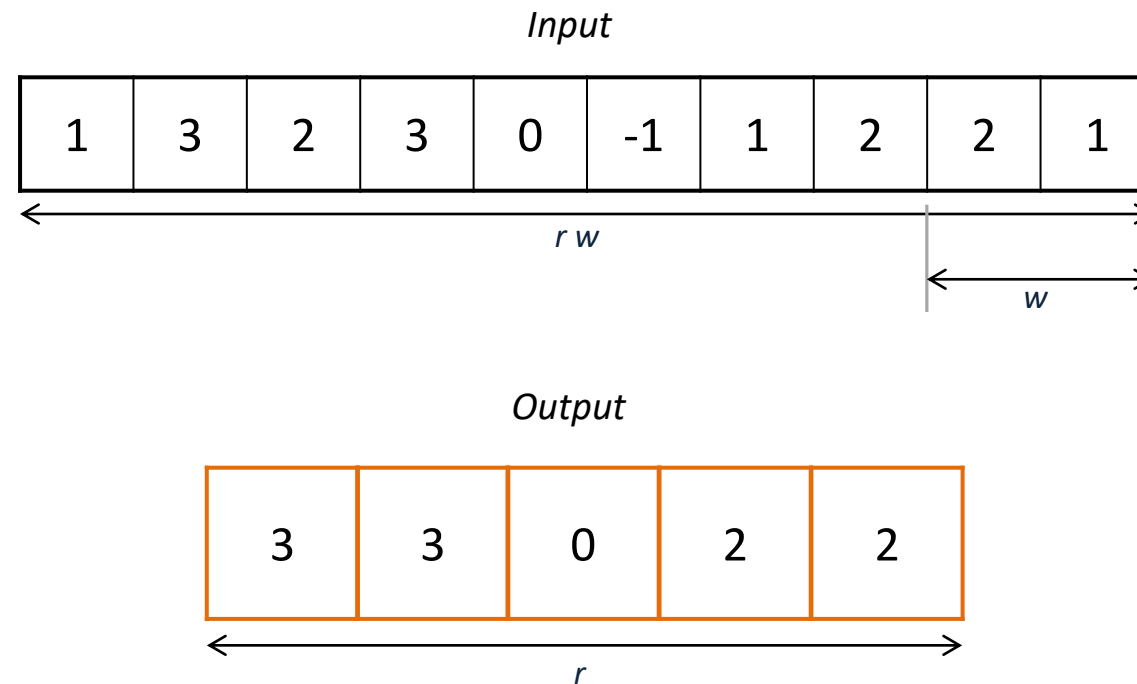$g(y)$

$$g(y) = \max(0, y)$$

# Pooling

# Pooling

- Compute a maximum value in a sliding window (max pooling)
- Reduce spatial resolution for faster computation
- Achieve invariance to local translation
- Max pooling introduces invariances
  - Pooling size : 2×2
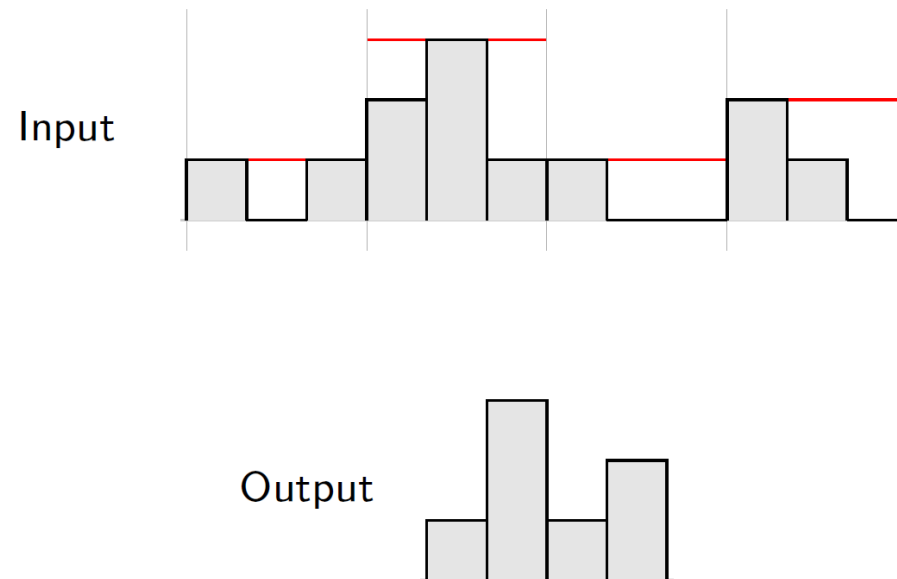  - No parameters: max or average of 2x2 units

Single depth slice

x

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

max pool with 2x2 filters
and stride 2

→

| 6 | 8 |
| 3 | 4 |

y

# Pooling

- Such an operation aims at grouping several activations into a single "more meaningful" one.

*Input*

| 1 | 3 | 2 | 3 | 0 | -1 | 1 | 2 | 2 | 1 |
|---|---|---|---|---|----|---|---|---|---|

$r\,w$

$w$

*Output*

| 3 | 3 | 0 | 2 | 2 |
|---|---|---|---|---|

$r$

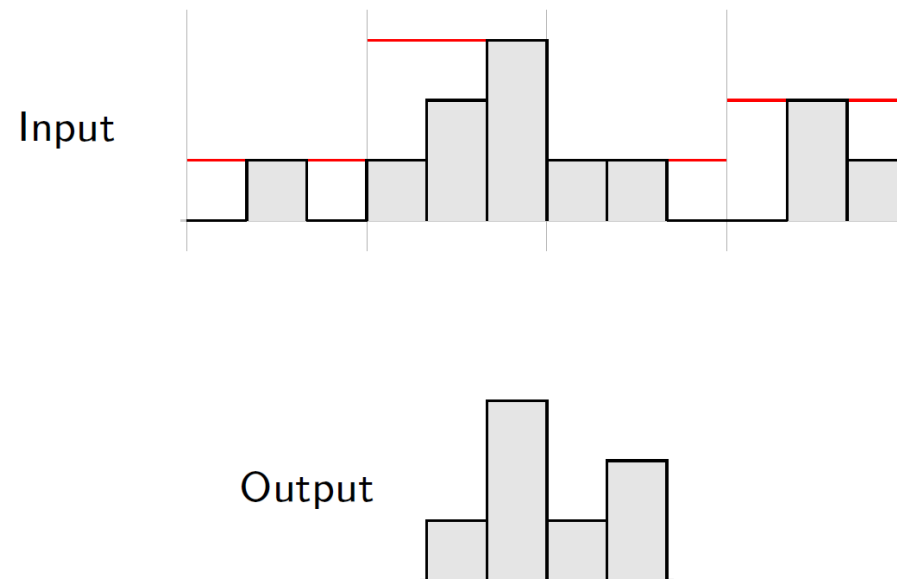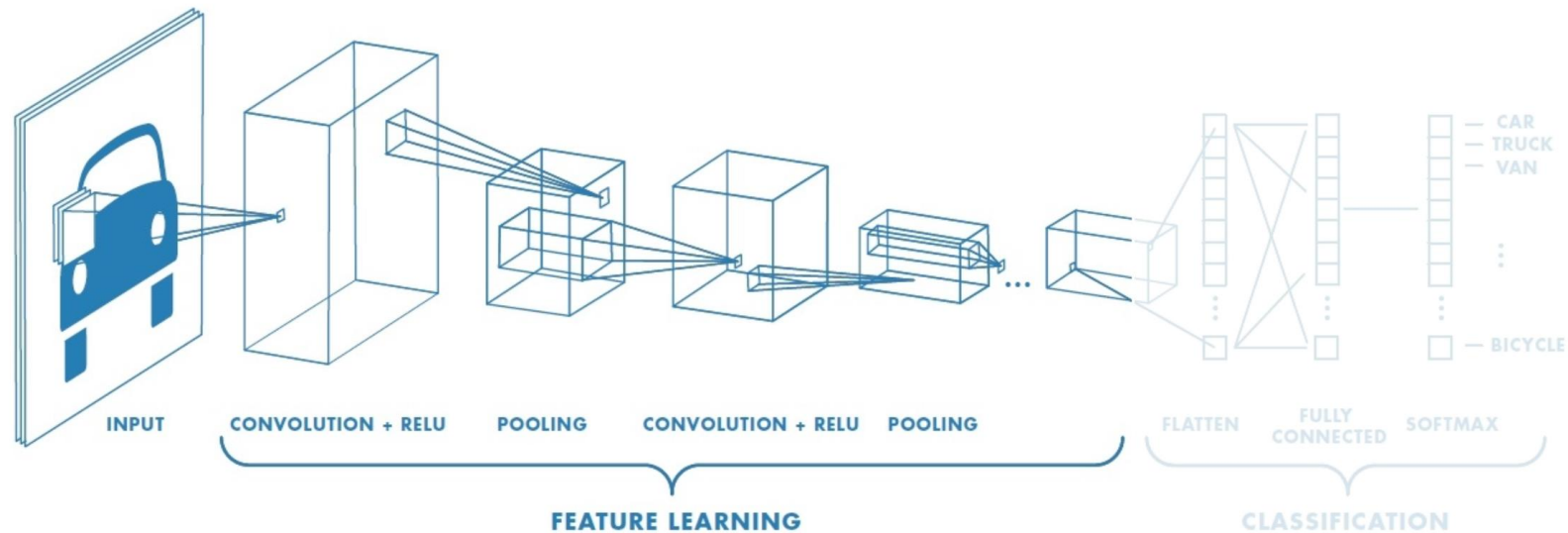- The average pooling computes average values per block instead of max values

# Pooling: Invariance

- Pooling provides invariance to any permutation inside one of the cell
- More practically, it provides a pseudo-invariance to deformations that result into local translations
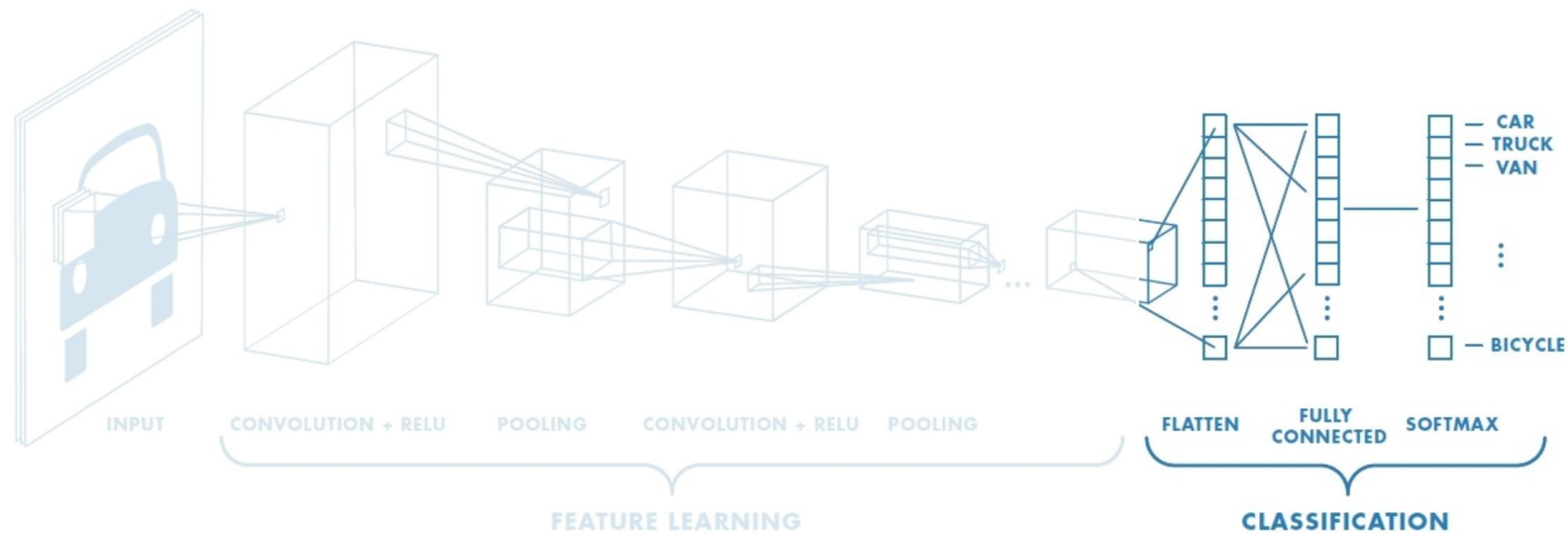
Input

Output

# Pooling: Invariance

- Pooling provides invariance to any permutation inside one of the cell
- More practically, it provides a pseudo-invariance to deformations that result into local translations

# CNNs for Classification: Feature Learning

- Learn features in input image through convolution
- Introduce non-linearity through activation function (real-world data is non-linear!)
- Reduce dimensionality and preserve spatial invariance with pooling

# CNNs for Classification: Class Probabilities

- CONV and POOL layers output high-level features of input
- Fully connected layer uses these features for classifying input image
- Express output as probability of image belonging to a particular class



$$\text{softmax}(y_i) = \frac{e^{y_i}}{\sum_j e^{y_j}}$$

실시간 강의자료

# Images

# Images Are Numbers

# Images Are Numbers

# Images Are Numbers



What the computer sees

An image is just a matrix of numbers [0,255]!
i.e., 1080×1080×3 for an RGB image

Source: 6.S191 Intro. to Deep Learning at MIT

# Images

Original image

R

G

B



Gray image

# Multiple Filters (or Kernels)

input features      a bank of 2 filters      2-dimensional output features

F1

$\Sigma$

F2

$\Sigma$
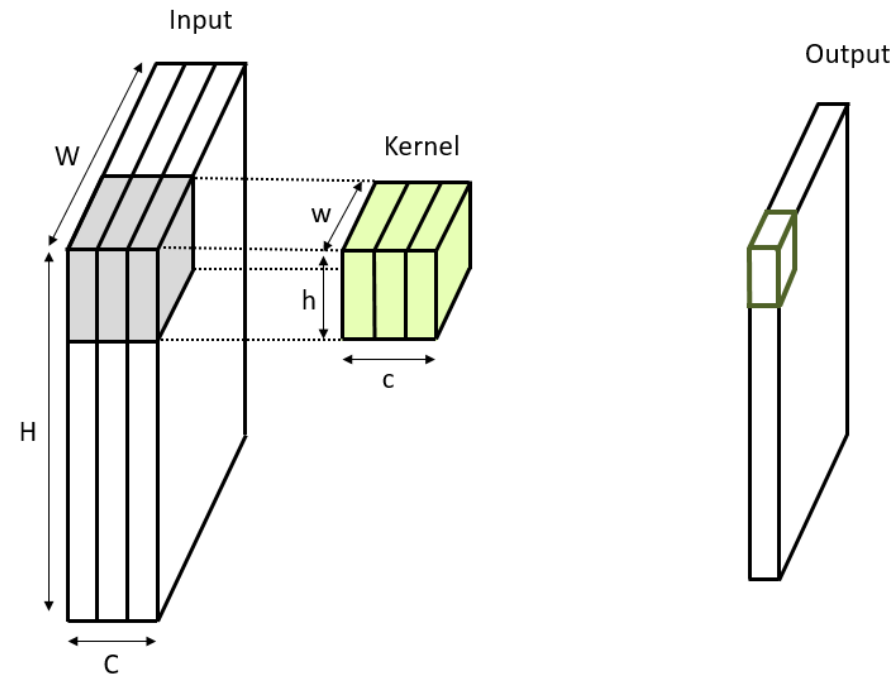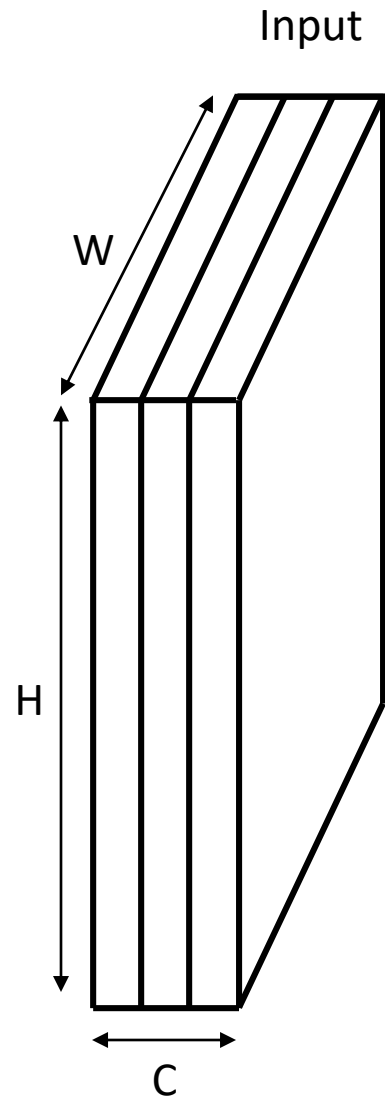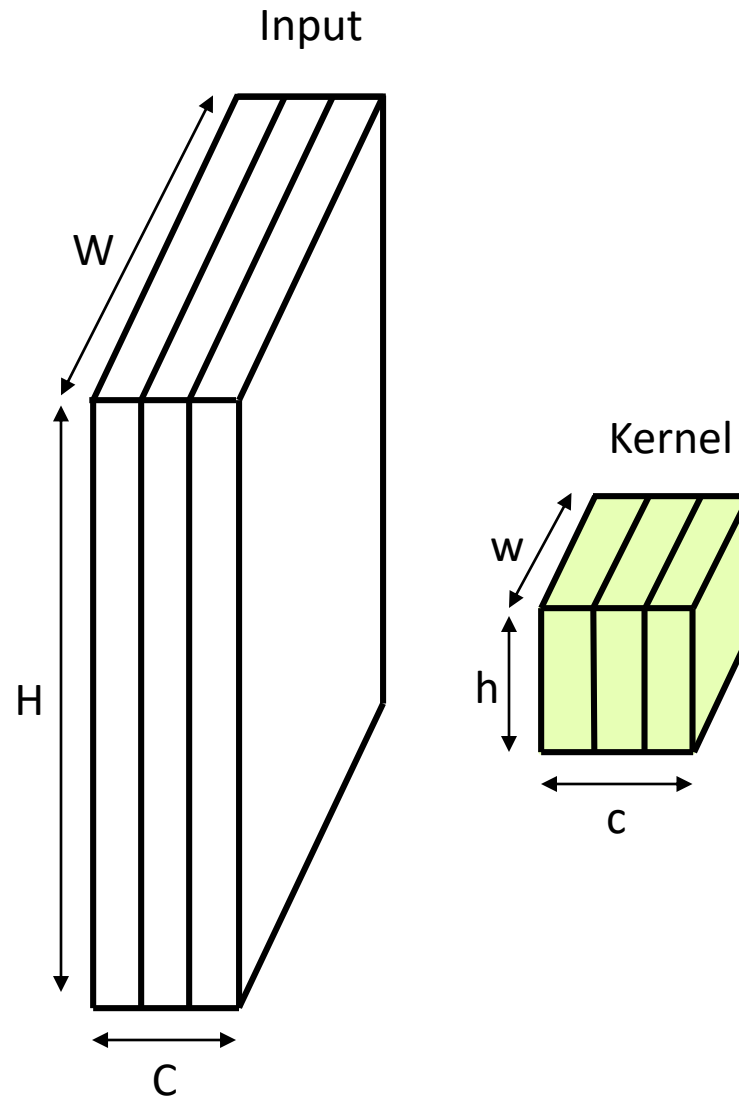
X

Y

# Channels



- Colored image = tensor of shape (height, width, channels)
- Convolutions are usually computed for each channel and summed:
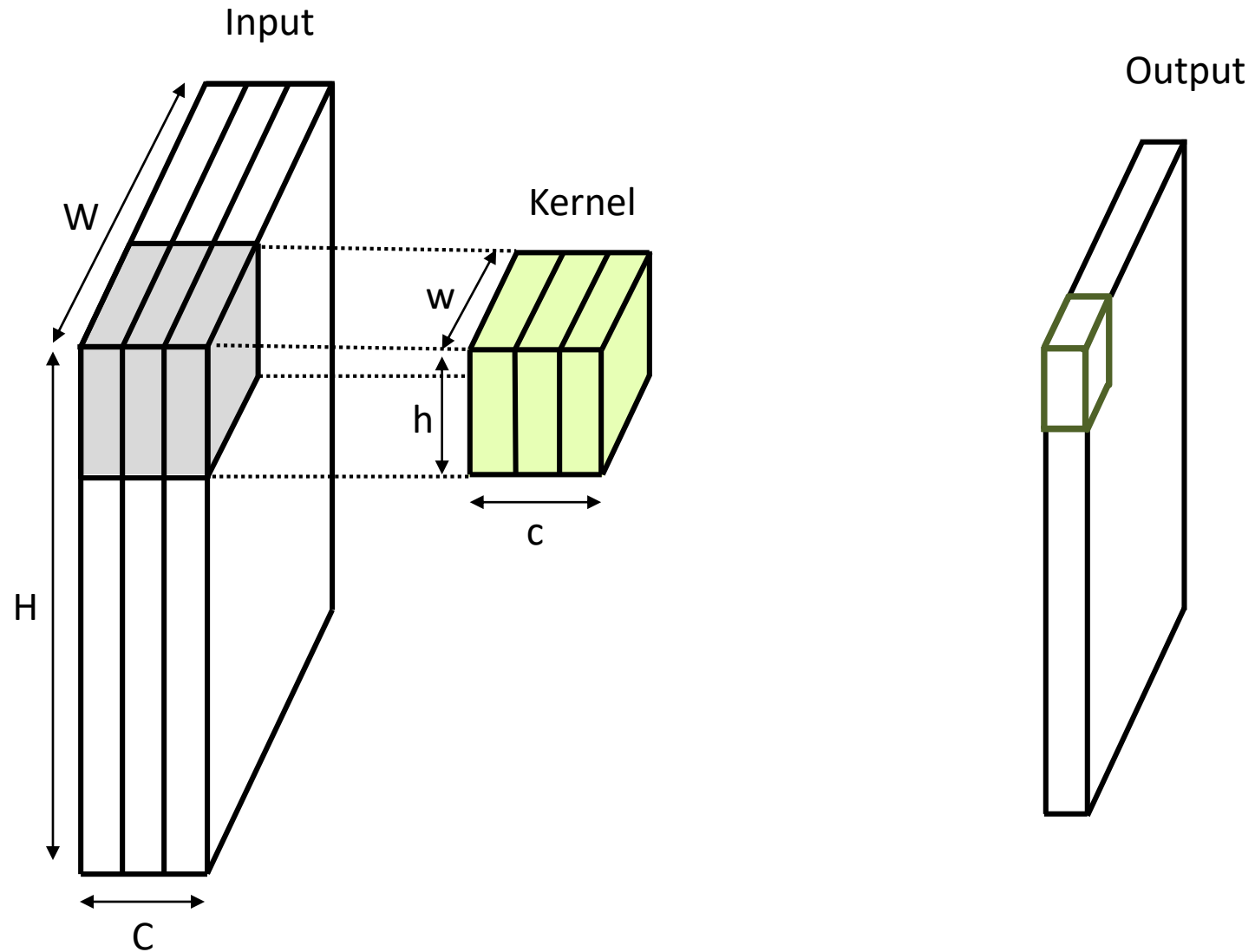- Kernel size aka receptive field (usually 1, 3, 5, 7, 11)
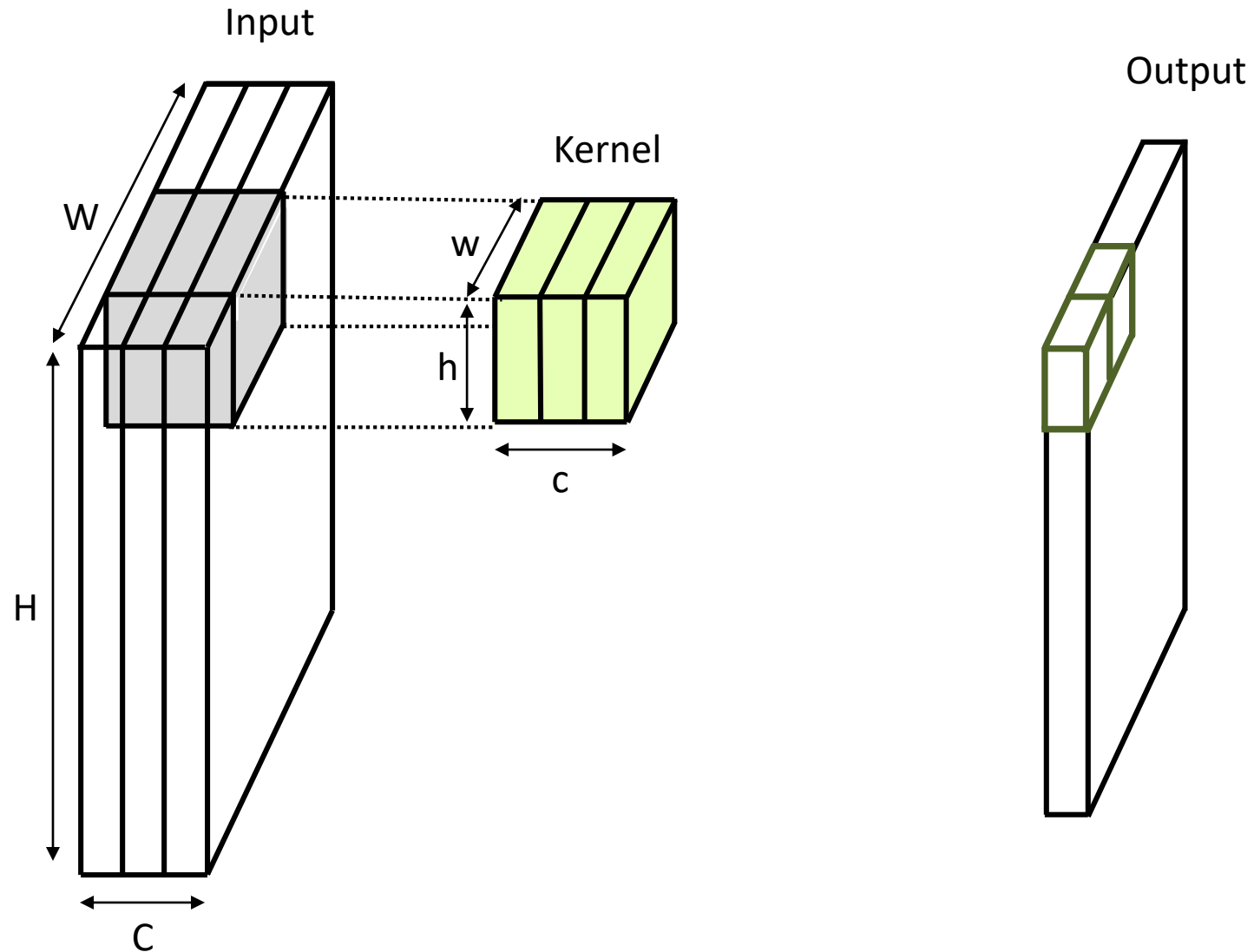
# Multi-channel 2D Convolution

Input



W

H

C

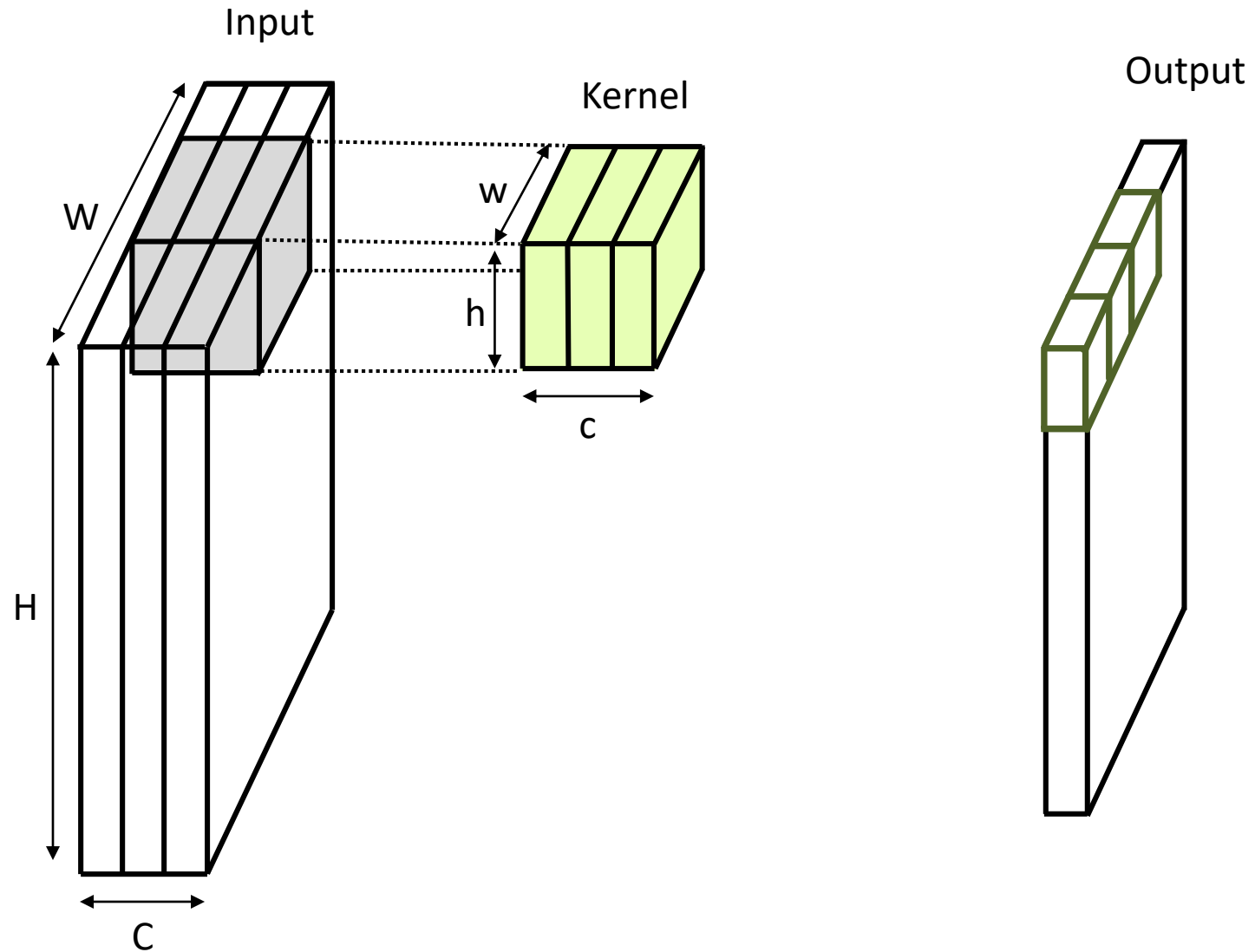# Multi-channel 2D Convolution



Input

Kernel

# Multi-channel 2D Convolution

Input

Kernel

Output



W

H

C

w

h

c

# Multi-channel 2D Convolution



Input

Kernel

Output

W

H

C

w

h

c

# Multi-channel 2D Convolution

# Multi-channel 2D Convolution



Input

Kernel

Output

W

H

C

w

h

c

# Multi-channel 2D Convolution

# Multi-channel 2D Convolution



Input

Kernel

Output

W

H

C

w

h

c

# Multi-channel 2D Convolution



Input

Kernel

Output

W

w

h

c

H

C

# Multi-channel 2D Convolution

# Multi-channel 2D Convolution



Input

Kernel

Output

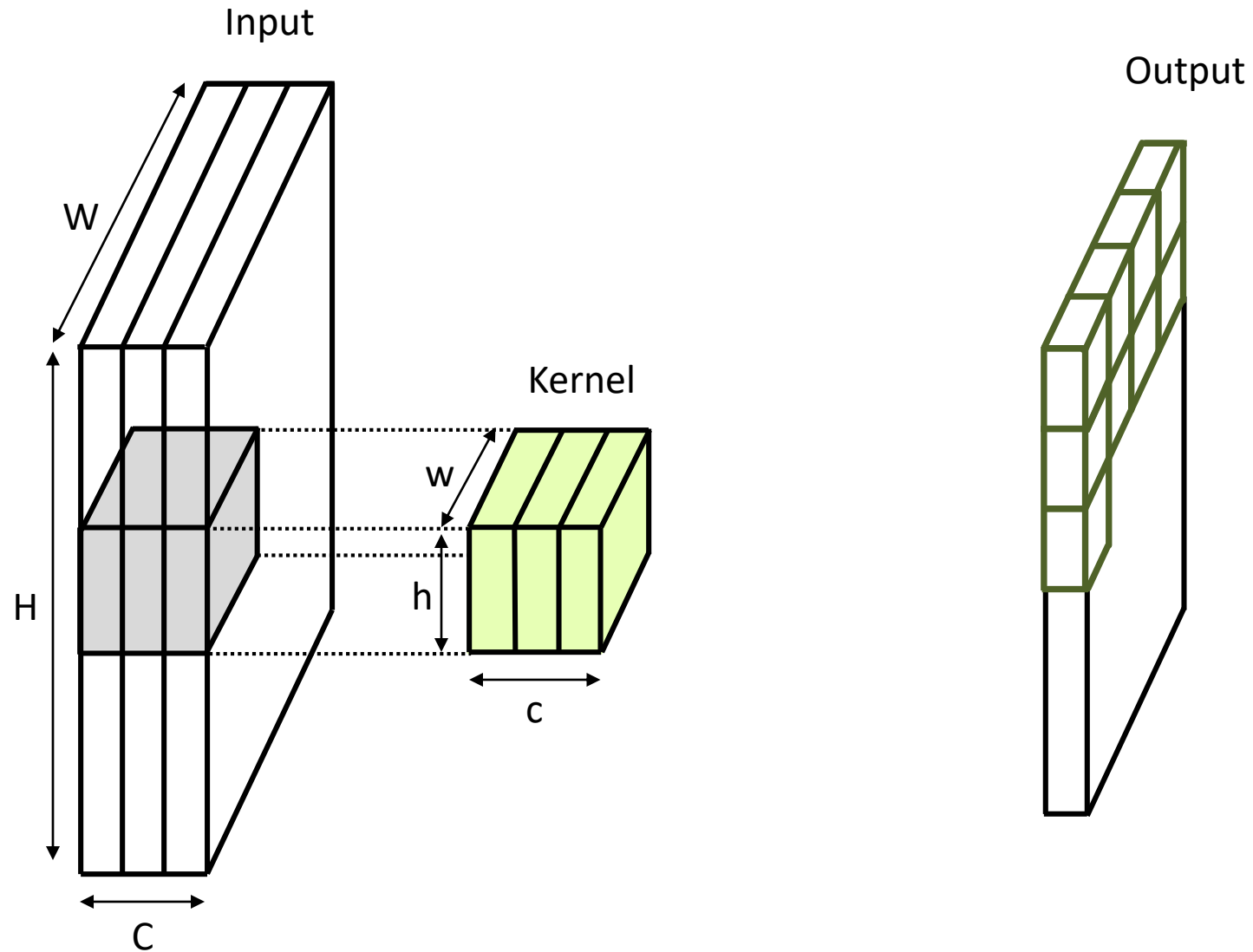# Multi-channel 2D Convolution



Input

W

H

C

Kernel

w

h

c

Output

$W - w + 1$

$H - h + 1$

1

# Multi-channel and Multi-kernel 2D Convolution

Input

Kernels

Output

W

H

C

w

h

c

D

W − w + 1

H − h + 1
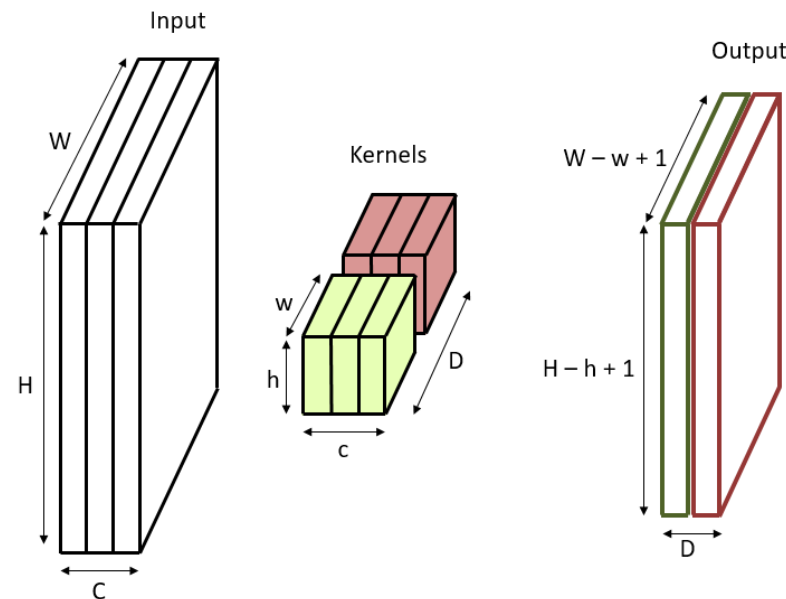
D

# Dealing with Shapes

- Activations or feature maps shape
  - Input $(W^i, H^i, C)$
  - Output $(W^o, H^o, D)$

- Kernel of Filter shape $(w, h, C, D)$
  - $w \times h$ Kernel size
  - $C$ Input channels
  - $D$ Output channels

- Numbers of parameters: $(w \times h \times C + 1) \times D$
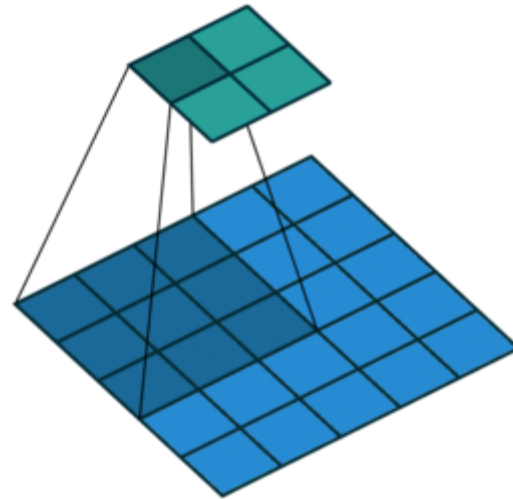  - bias

# Multi-channel 2D Convolution

- The kernel is not swiped across channels, just across rows and columns.

- Note that a convolution preserves the signal support structure.
  - A 1D signal is converted into a 1D signal, a 2D signal into a 2D, and neighboring parts of the input signal influence neighboring parts of the output signal.

- We usually refer to one of the channels generated by a convolution layer as an activation map.

- The sub-area of an input map that influences a component of the output as the receptive field of the latter.
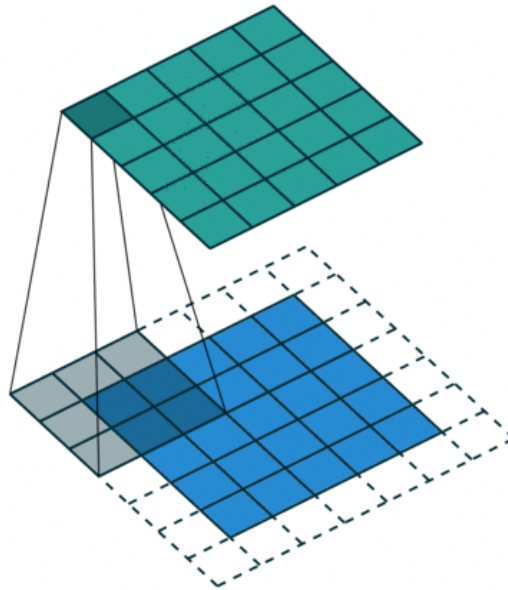
# Padding and Stride

# Strides

- Strides: increment step size for the convolution operator
- Reduces the size of the output map



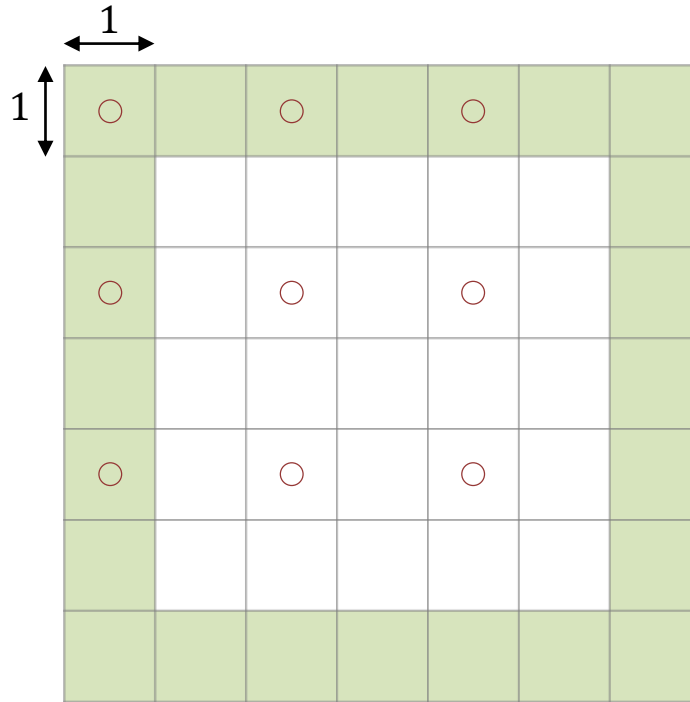Example with kernel size 3×3 and a stride of 2 (image in blue)

# Padding

- Padding: artificially fill borders of image
- Useful to <span style="color:red">keep spatial dimension constant</span> across filters
- Useful with strides and large receptive fields
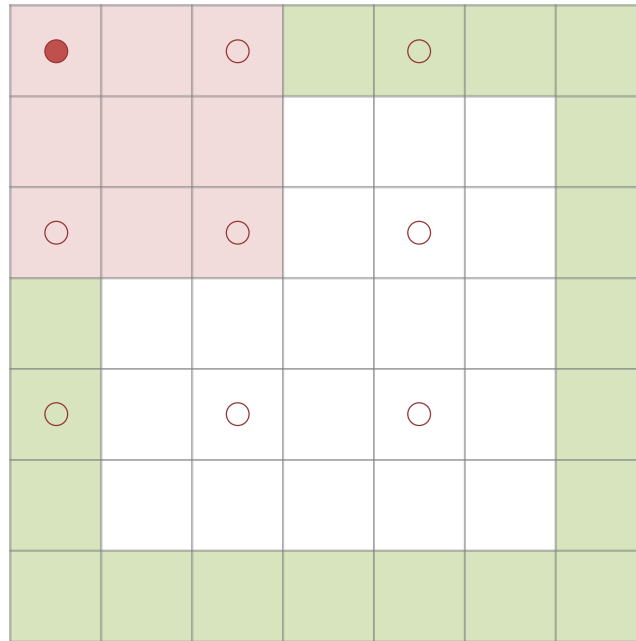- Usually fill with 0s

# Padding and Stride

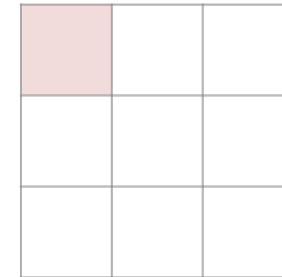- Here with $5 \times 5 \times C$ as input, a padding of $(1, 1)$, a stride of $(2, 2)$



Input

# Padding and Stride

- Here with $5 \times 5 \times C$ as input, a padding of $(1,1)$, a stride of $(2,2)$, and a kernel of size $3 \times 3 \times C$



Input

Output

# Padding and Stride

- Here with $5 \times 5 \times C$ as input, a padding of $(1,1)$, a stride of $(2,2)$, and a kernel of size $3 \times 3 \times C$
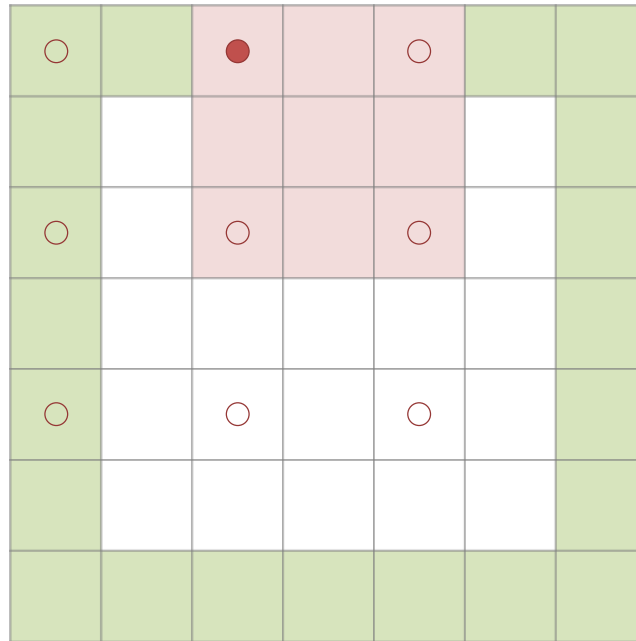


Input

Output

# Padding and Stride

- Here with $5 \times 5 \times C$ as input, a padding of $(1,1)$, a stride of $(2,2)$, and a kernel of size $3 \times 3 \times C$
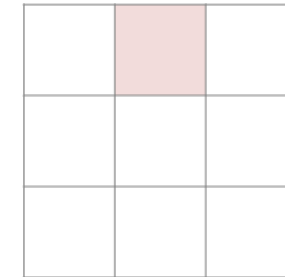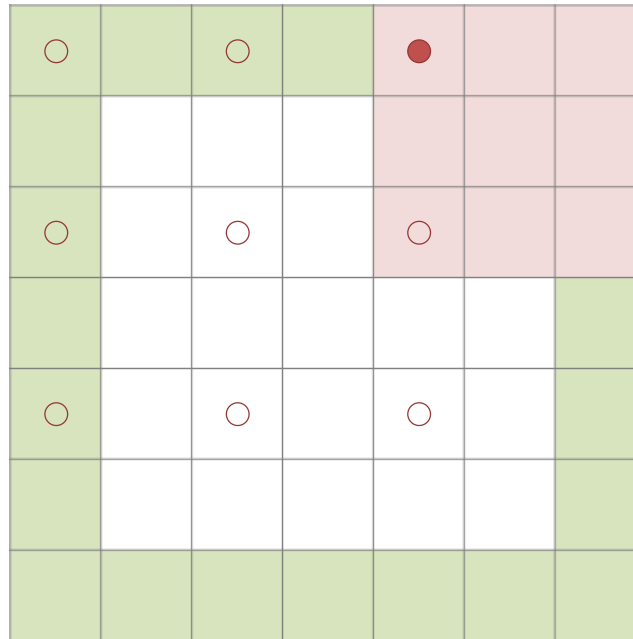


Input



Output

# Padding and Stride

- Here with $5 \times 5 \times C$ as input, a padding of (1,1), a stride of (2,2), and a kernel of size $3 \times 3 \times C$
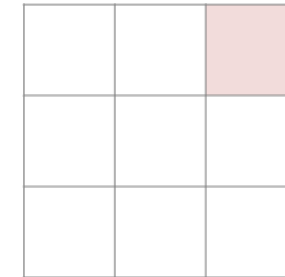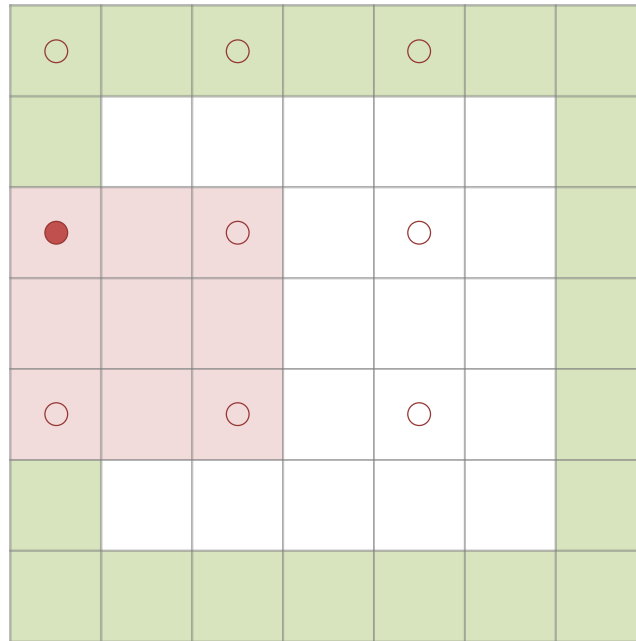


Input

Output

# Padding and Stride

- Here with $5 \times 5 \times C$ as input, a padding of $(1,1)$, a stride of $(2,2)$, and a kernel of size $3 \times 3 \times C$



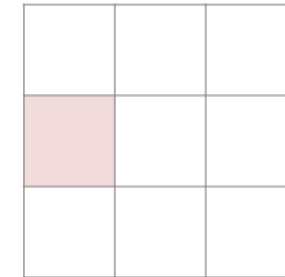Input

Output

# Padding and Stride

- Here with $5 \times 5 \times C$ as input, a padding of $(1,1)$, a stride of $(2,2)$, and a kernel of size $3 \times 3 \times C$

Input

Output

# Padding and Stride

- Here with $5 \times 5 \times C$ as input, a padding of $(1,1)$, a stride of $(2,2)$, and a kernel of size $3 \times 3 \times C$
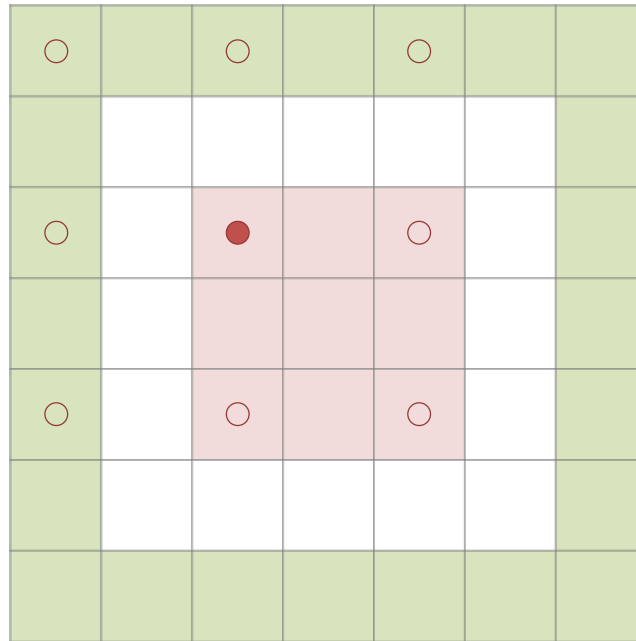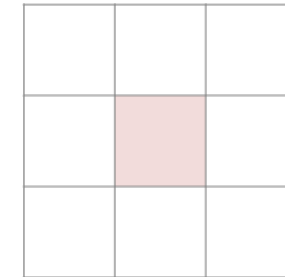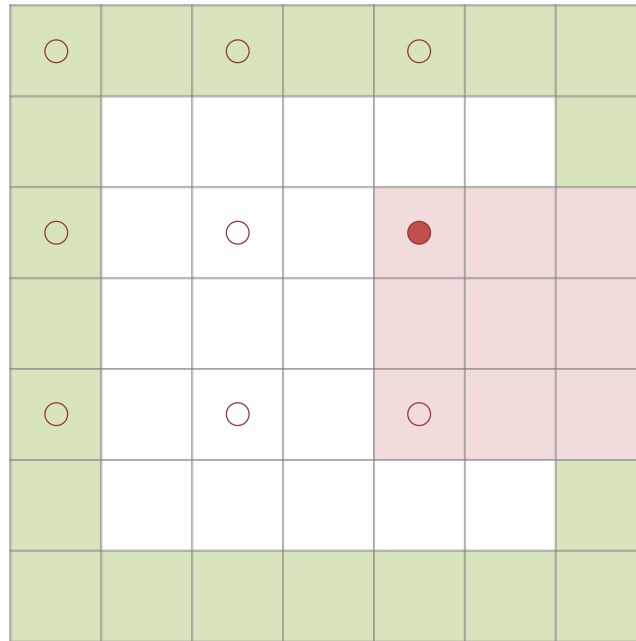


Input

Output

# Padding and Stride

- Here with $5 \times 5 \times C$ as input, a padding of $(1,1)$, a stride of $(2,2)$, and a kernel of size $3 \times 3 \times C$
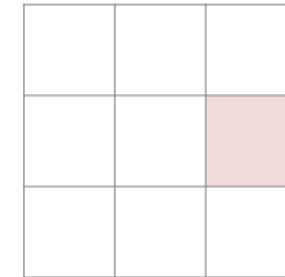


Input

Output

# Padding and Stride

- Here with $5 \times 5 \times C$ as input, a padding of $(1,1)$, a stride of $(2,2)$, and a kernel of size $3 \times 3 \times C$
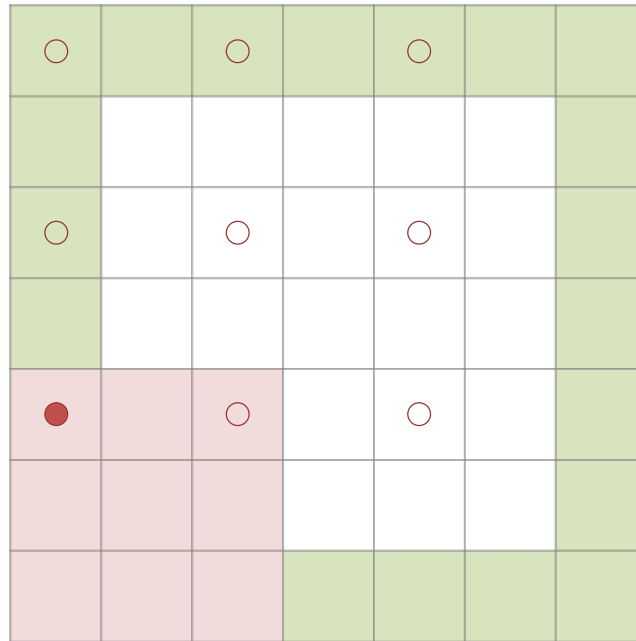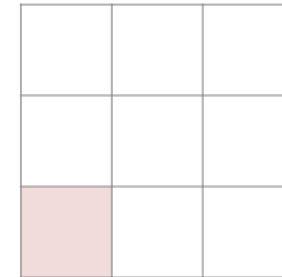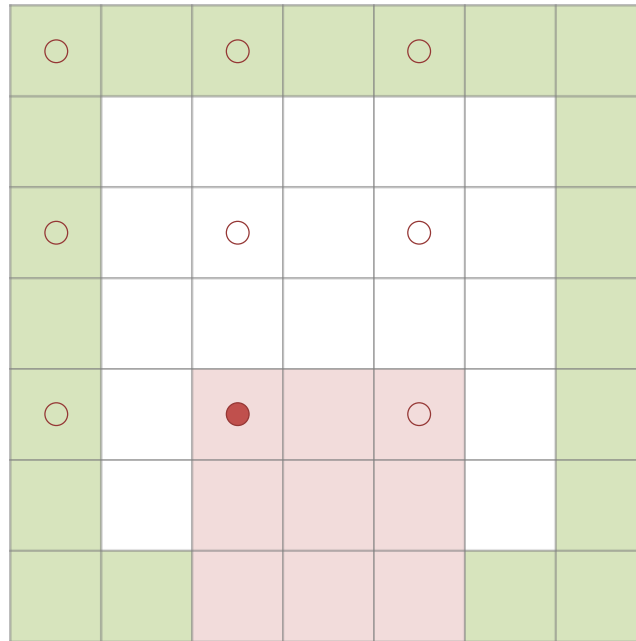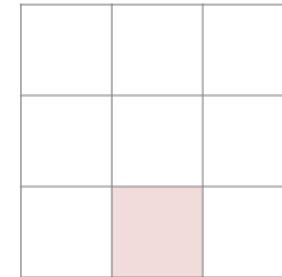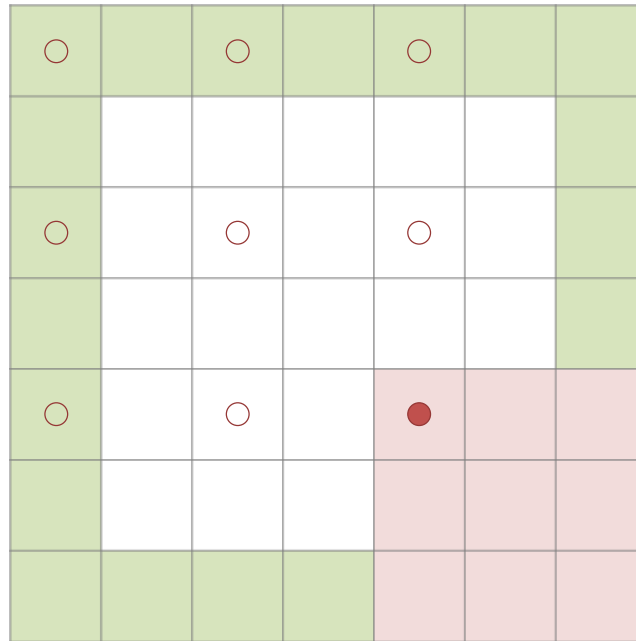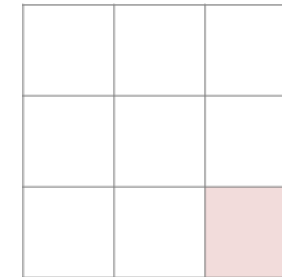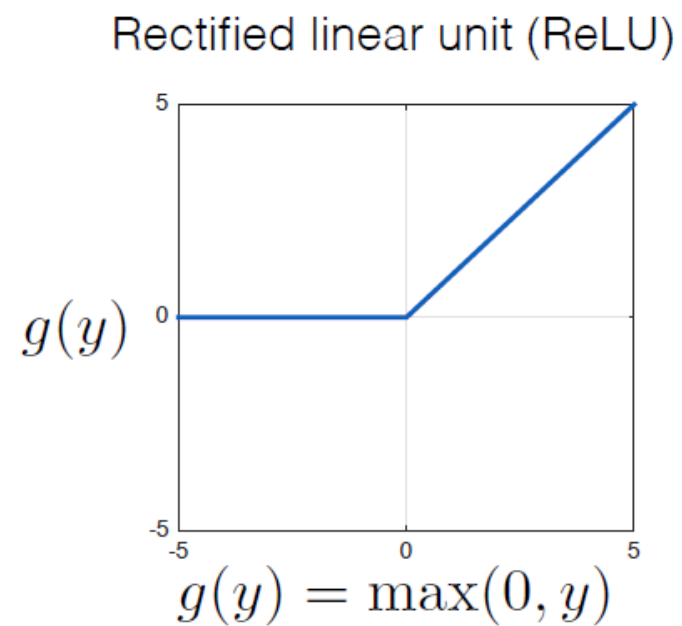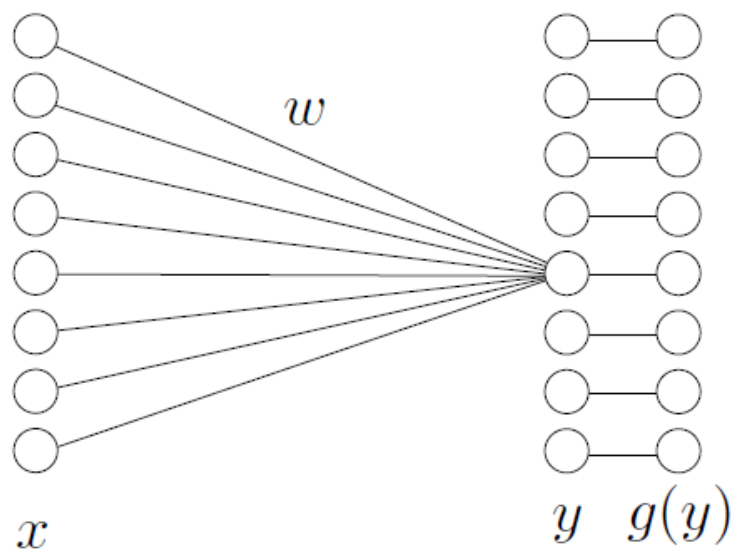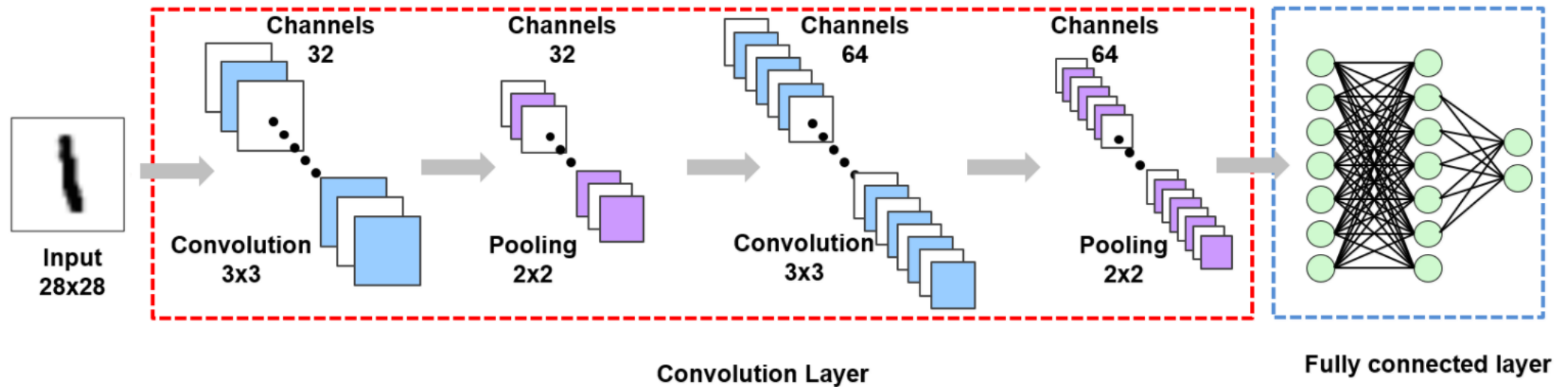
Input

Output

# Nonlinear Activation Function



$$g(y) = \max(0, y)$$

Rectified linear unit (ReLU)

# CNN in TensorFlow

# Lab: CNN with TensorFlow

- MNIST example
- To classify handwritten digits

# CNN Structure

```python
model = tf.keras.models.Sequential([
    tf.keras.layers.Conv2D(32,
                           (3,3),
                           activation = 'relu',
                           padding = 'SAME',
                           input_shape = (28, 28, 1)),

    tf.keras.layers.MaxPool2D((2,2)),

    tf.keras.layers.Conv2D(64,
                           (3,3),
                           activation = 'relu',
                           padding = 'SAME',
                           input_shape = (14, 14, 32)),

    tf.keras.layers.MaxPool2D((2,2)),

    tf.keras.layers.Flatten(),

    tf.keras.layers.Dense(128, activation = 'relu'),

    tf.keras.layers.Dense(10, activation = 'softmax')
])
```
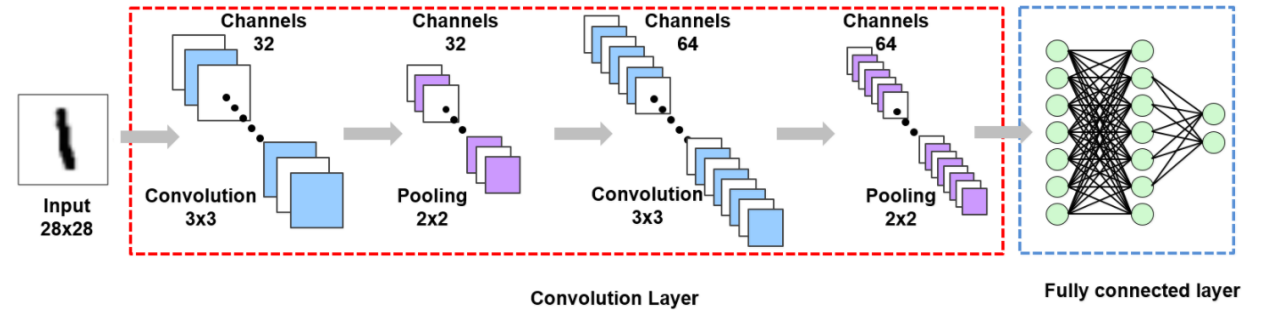
# Loss and Optimizer

- Loss
  - Classification: Cross entropy
  - Equivalent to applying logistic regression
- Optimizer
  - GradientDescentOptimizer
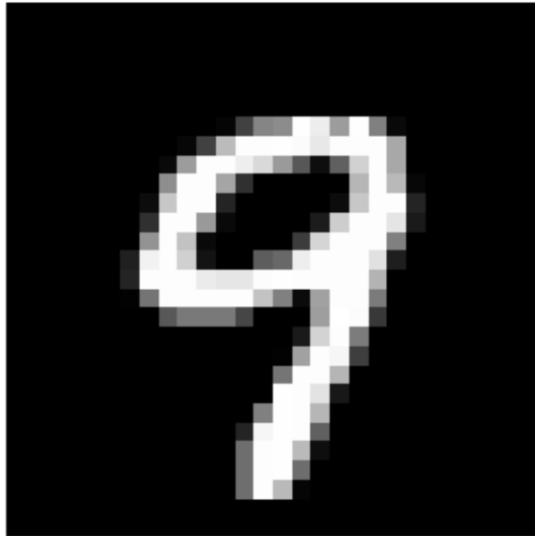  - AdamOptimizer: the most popular optimizer

```
model.compile(optimizer = 'adam',
              loss = 'sparse_categorical_crossentropy',
              metrics = ['accuracy'])
```

```
model.fit(train_x, train_y)
```

# Test or Evaluation

```
test_loss, test_acc = model.evaluate(test_x, test_y)
```

```
313/313 [==============================] - 1s 4ms/step - accuracy: 0.9838 - loss: 0.0466
loss = 0.05, Accuracy = 98 %
```



Prediction : 9