

# Deep Learning



# Convolutional Neural Networks



# Learning Objectives

By the end of this lesson, you will be able to:

- 👁️ Analyze how to use image data for training
- 👁️ Interpret the structure and functionality of various convolutional neural network (CNN) architectures
- 👁️ Apply filters and pooling technique in CNN
- 👁️ Employ CNN for image classification
- 👁️ Analyze models using TensorBoard



## Business Scenario

A startup is working on an image recognition system designed to aid in medical diagnoses through medical imaging.

The company explores the application of convolutional neural network (CNN) algorithms, training their system to identify various medical conditions from X-rays and CT scans. They use TensorBoard to visualize the performance of these models, adjusting as needed. In their pursuit of obtaining optimal results, they experiment with various CNN filters, including horizontal and vertical sobel filters, blur filters, and outline filters, to identify the most effective ones for their specific medical imaging use case.

In a bid to enhance system accuracy, they're contemplating the integration of a residual neural network (ResNet) architecture.





# Introduction to CNN



## Discussion

# Discussion: CNN

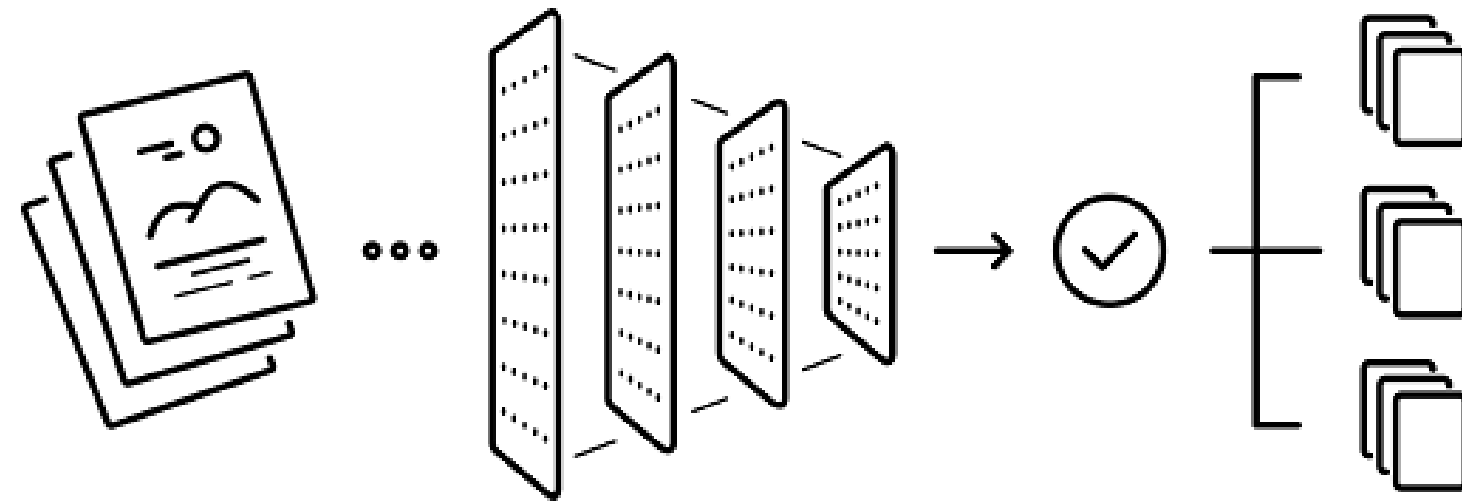
Duration: 10 minutes

- What is a CNN?
- What are some real-world applications of CNNs?



# Convolutional Neural Network (CNN)

A convolutional neural network (CNN) is a deep learning model specifically designed for visual data analysis, extracting features through convolutional and pooling layers to achieve high-level pattern recognition and classification.



Its architecture is modeled after the visual cortex and has transformed computer vision tasks.

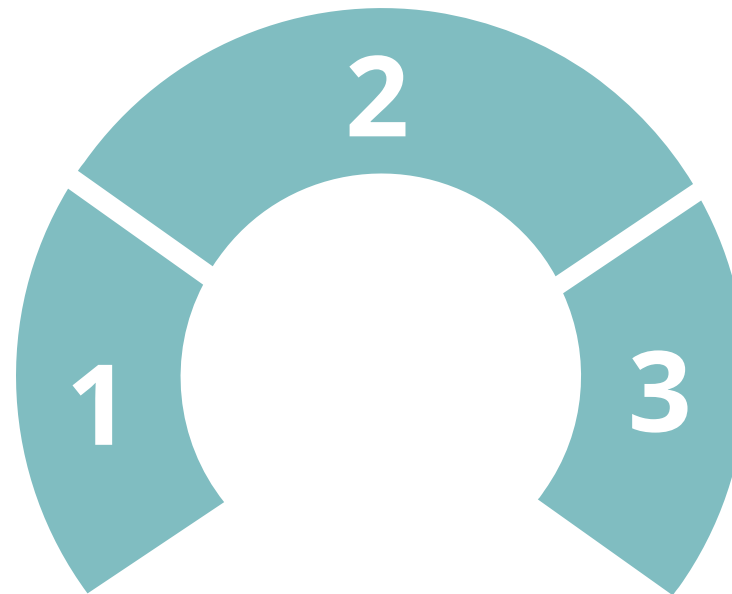


# Advantages of CNN

Some of the advantages of CNN are:

They automatically detect essential features without any human supervision.

Compared to feed-forward neural networks, CNNs have higher accuracy in resolving image recognition problems.



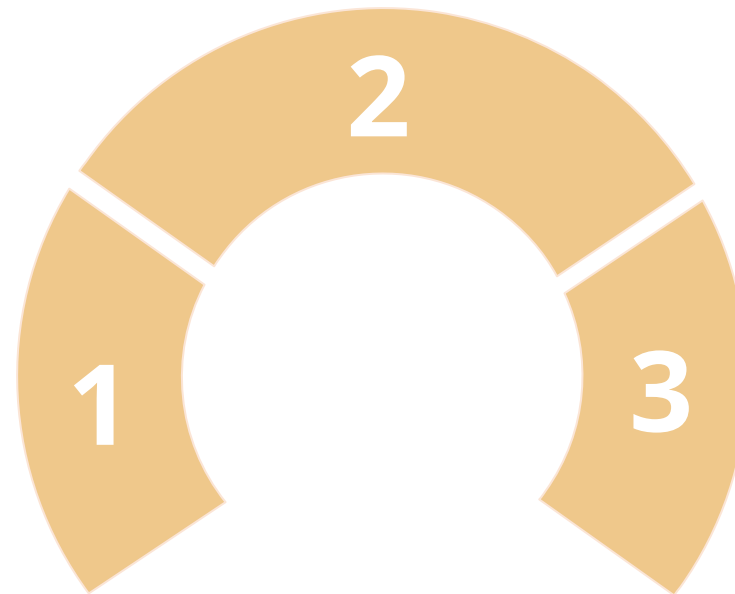
The results are more accurate than generic machine-learning techniques.

# Disadvantages of CNN

Some of the disadvantages of CNN are:

It requires more computing power and huge amounts of training data.

It does not encode the position and orientation of the object.



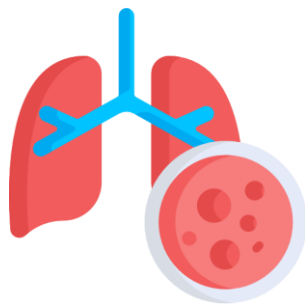
It cannot be spatially invariant with the input data.

# CNN Applications

CNN is largely used in computer vision applications such as:



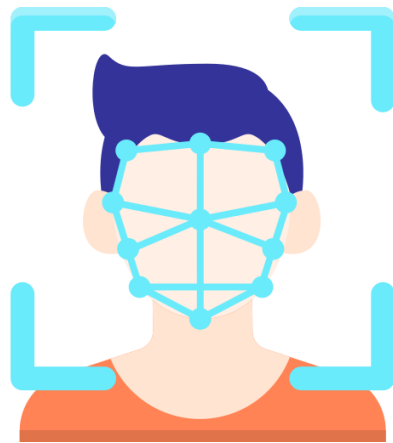
Detection of tools in a factory, where a CNN model can be trained to detect misplaced tools by factory workers



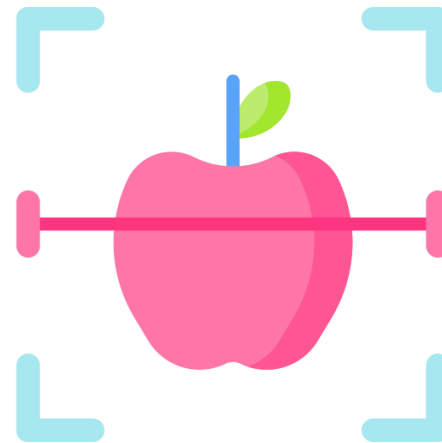
Detection of pulmonary fibrosis, where a CNN model can be fed a large dataset of a patient's lung images to identify any scarring in lung tissues

# CNN Applications

CNN is a popular algorithm used widely in the field of computer vision for the following applications:



Face recognition



Object detection

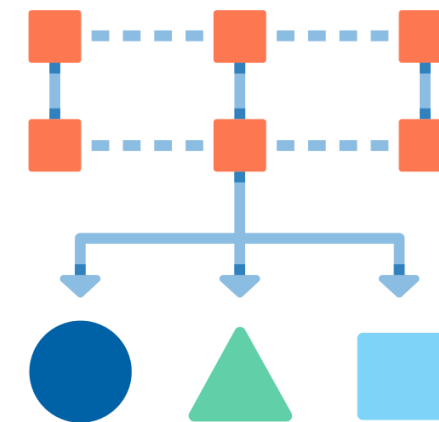


Image classification tasks

It uses a set of filters to extract features from the image.

## Uses of CNN

Consider an example of a dataset containing 60,000 images, each with dimensions (28,28,3) representing height, width, and color channels, respectively

If these images are to be processed by a Feed-Forward Neural Network (FFN), each image must first be flattened.

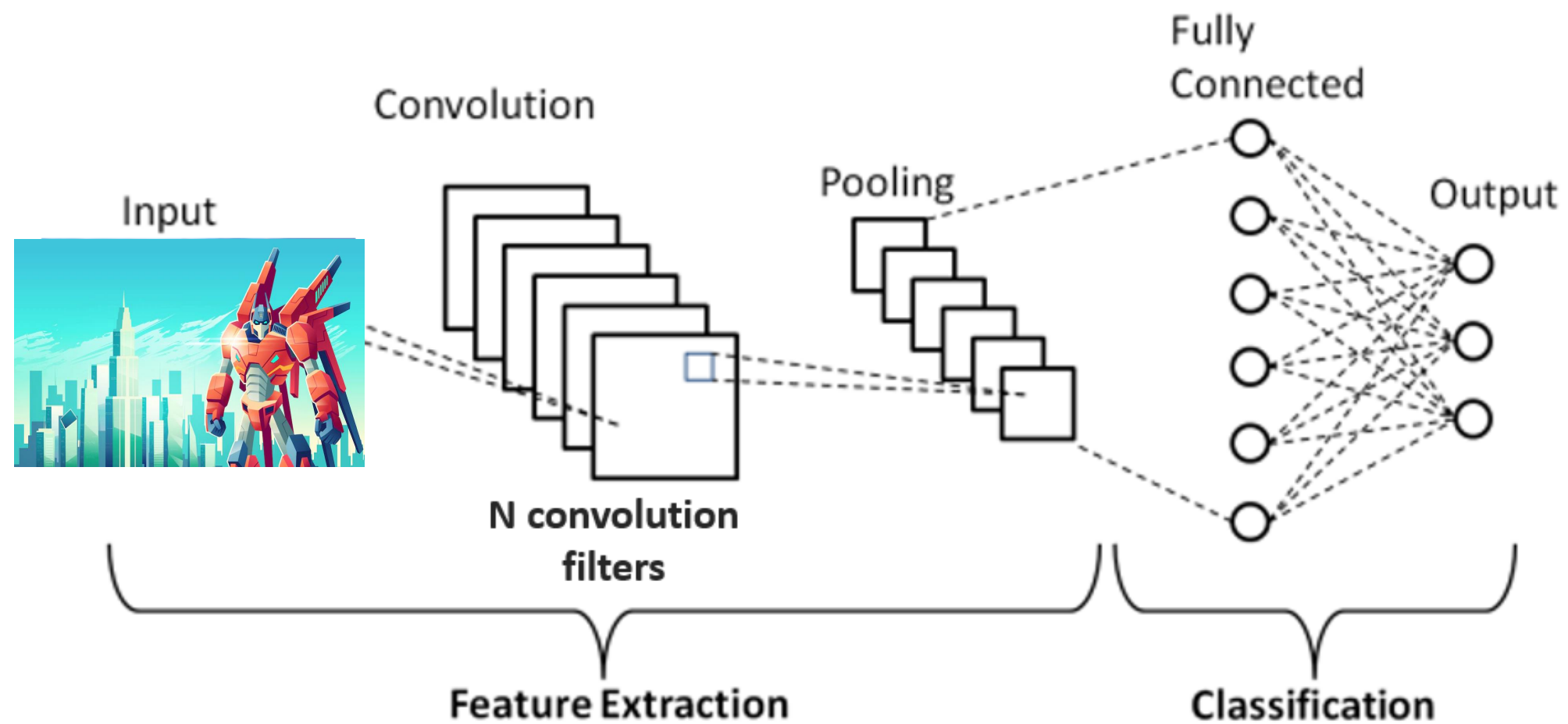
The shape of each image after flattening will be (2352).

When images have dimensions (1000, 1000, 3) representing height, width, and color channels (RGB), a feed-forward neural network (FFN) requires additional processing power for computations.

CNN was introduced to avoid such issues.

# Uses of CNN

It extracts features from the images and converts them into lower dimensions without losing their characteristics.



# Uses of CNN

Considering the image:

The initial image size is (400, 400, 3), and without convolution,  $400 \times 400 \times 3 = 4,80,000$  neurons will be needed in the input layer.

After applying convolution, the input tensor dimension is reduced to  $1 \times 1 \times 3$ . Therefore, only three neurons are needed in the first layer of the FFN.

## Discussion: CNN

Duration: 10 minutes



- What is a CNN?

**Answer:** A CNN, or Convolutional Neural Network, is a type of deep learning model that is widely used for processing and analyzing visual data, such as images and videos.

- What are some real-world applications of CNNs?

**Answer:** CNNs have numerous real-world applications, including image classification, object detection, facial recognition, self-driving cars, medical image analysis, and video analysis. They excel in tasks that require understanding and interpretation of visual data.





## Getting Started with Image Data

# Image Data

It consists of three dimensions. Example:



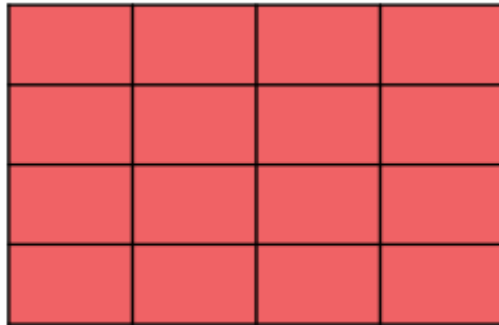
$(400, 400, 3)$  denotes the shape of the image.

$(400, 400)$  denotes the height and width of the image.

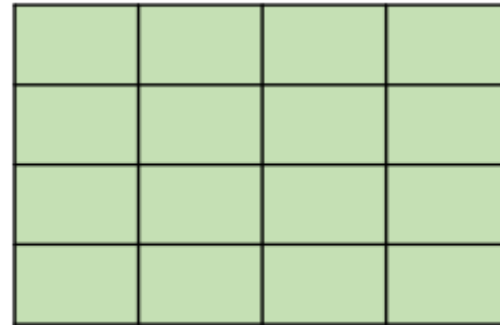
3 denotes the three channels in the image.

# Image Data

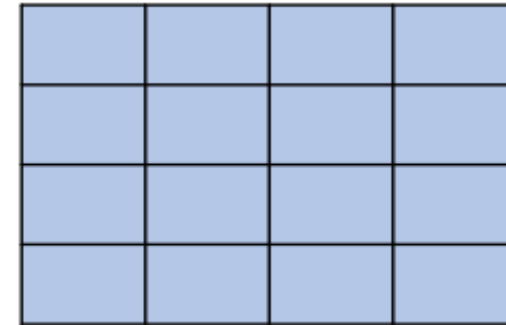
The three channels, R, G, and B, represent the color of the image in a combination of red, green, and blue as shown below:



R



G



B

For each channel, the values range from 0 to 255, where 0 represents the absence of that color and 255 represents the maximum intensity of that color.

# CNN with Image Data

In a CNN operation on the image data, convolution performs robust feature extraction.

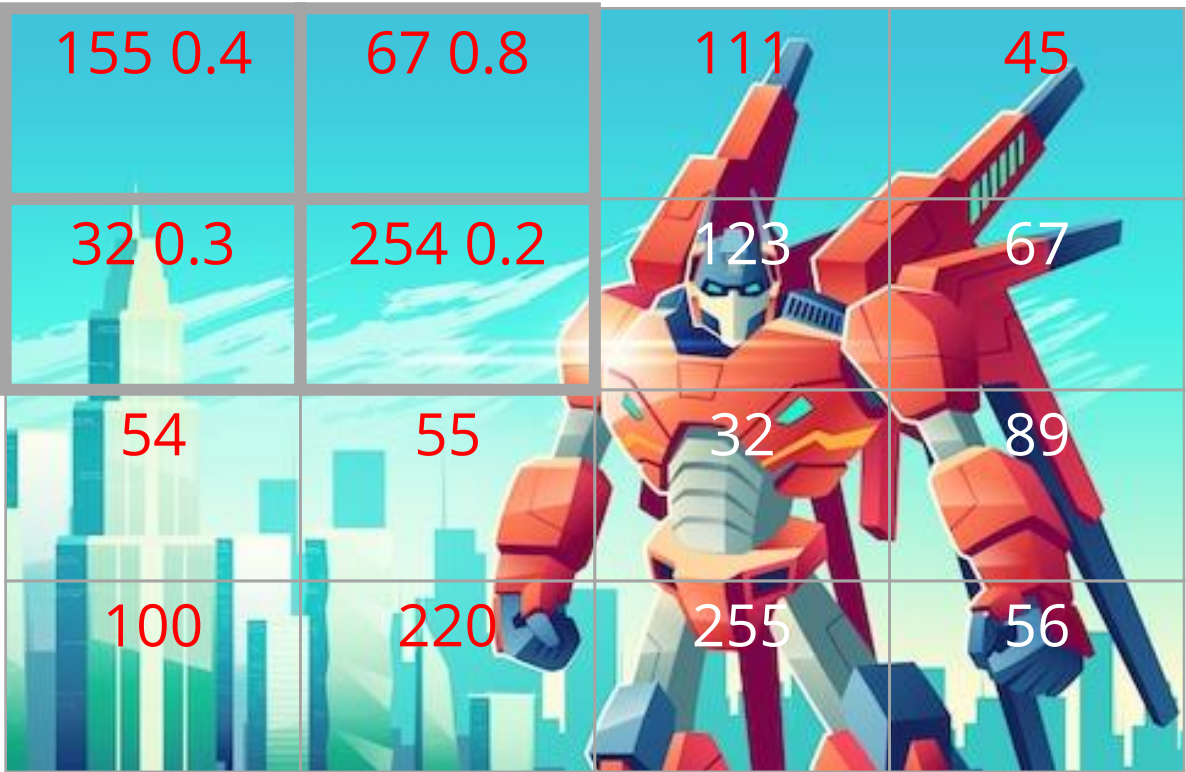


CNN extracts all the necessary information that is helpful to train models.

# Convolution Operation

It is the sum of the product of the filter and pixel values.

Consider the following image:



# Convolution Operation

The following can be observed:

There is an image with shape (4,4) and a filter with shape (2,2).

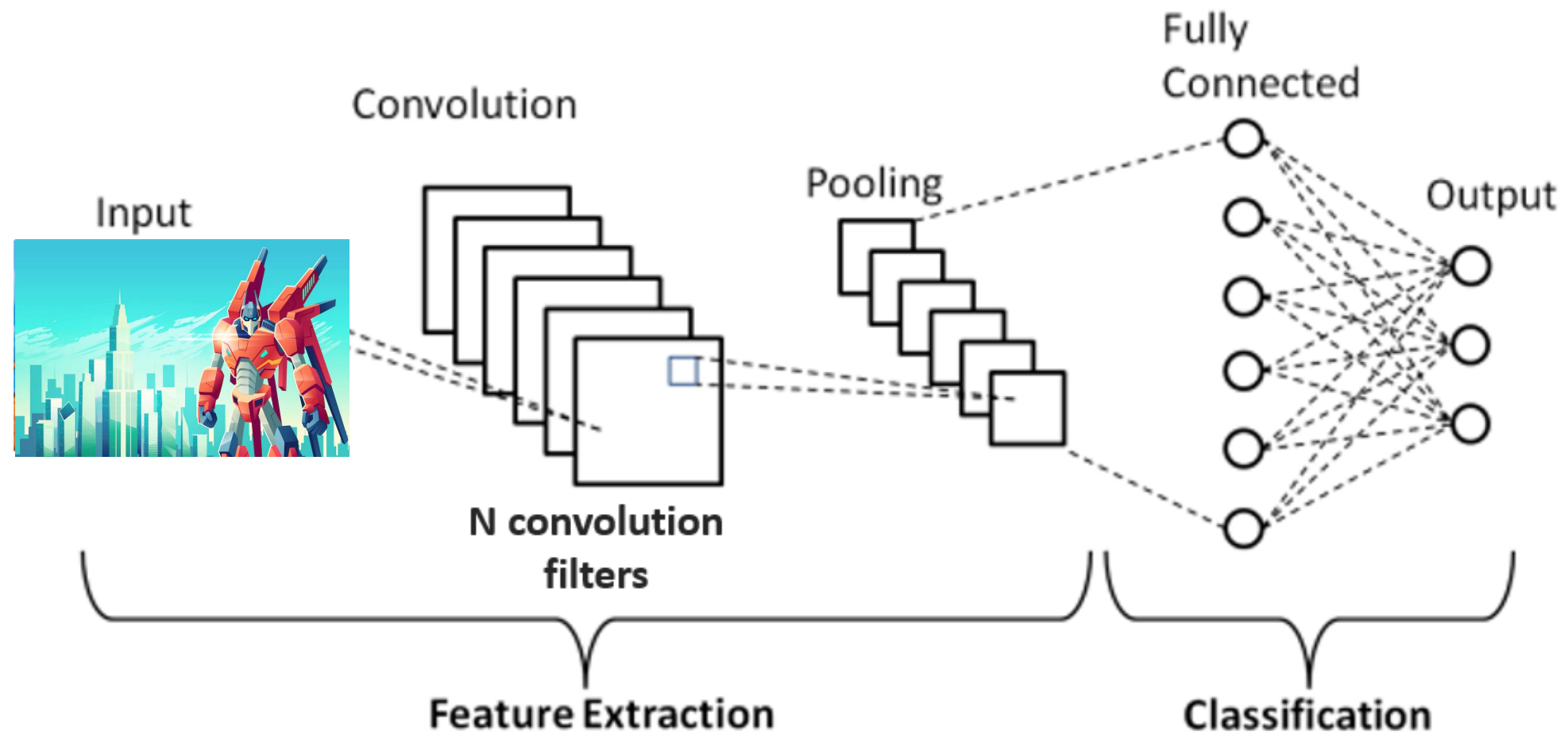
The filter values are convoluted with the image pixel values, and the filter is moved across the image by one step.

After convolution, one gets the information about specific features that the filter is designed to detect.

176	216.4	113.2
144	112	167
156	132	118

# Convolution Operation

Consider the following image:



# Convolution Operation

From the image, it can be observed that:

The input image is passed through  $N$  convolution filters. On valid padding, the output shape is  $(400, 400, n)$ .

After further pooling operations, the output matrix is flattened to form a feed-forward neural network for more processing.



## Assisted Practices



Let's understand the concept of CNN with image data using Jupyter Notebooks.

- 7.04\_Working with Image Data

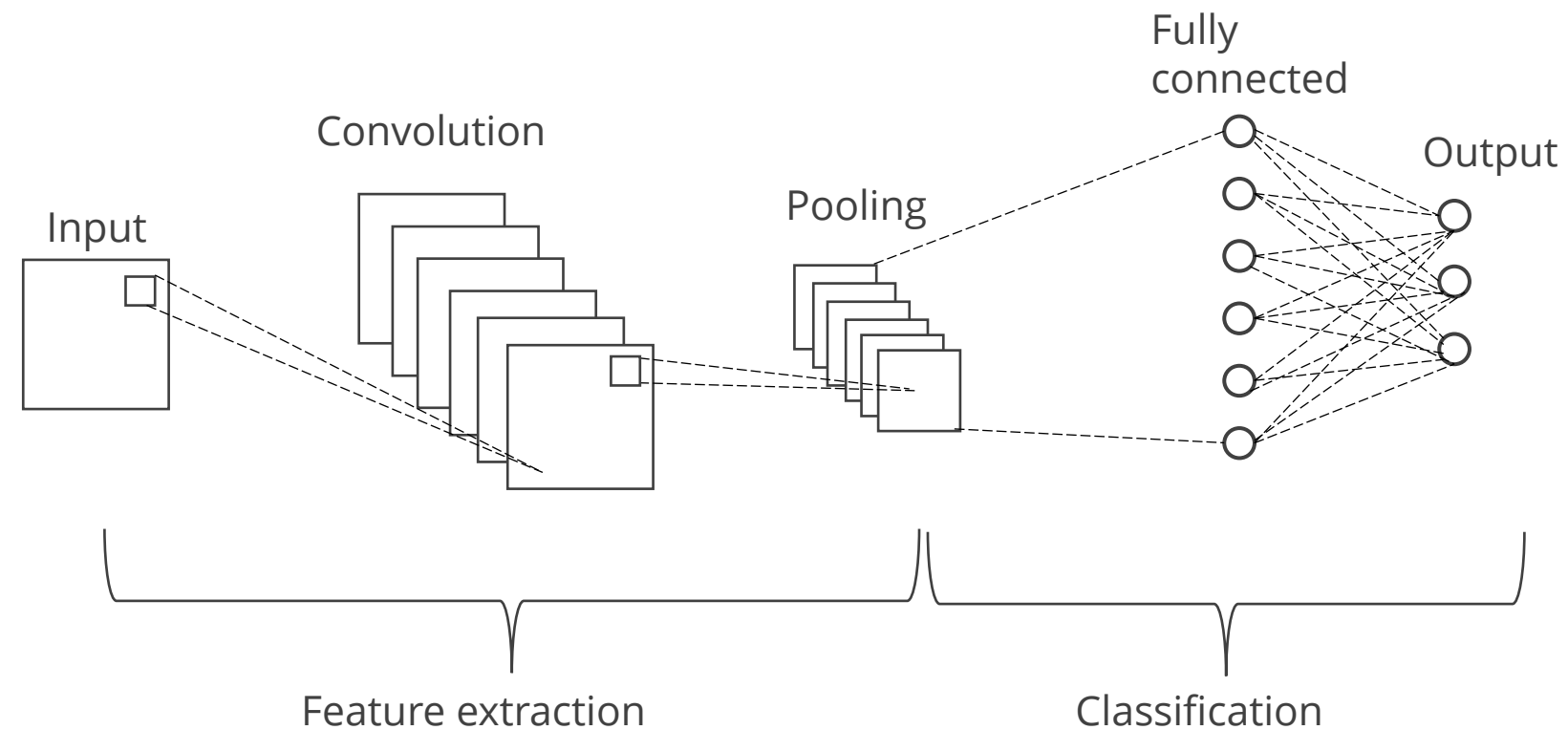
**Note:** Please refer to the Reference Material section to download the notebook files corresponding to each mentioned topic



# CNN Architecture

# CNN Architecture

A CNN combines a backpropagation algorithm with multiple layers, including convolution, pooling, and fully connected layers.



It aims to automatically and adaptively learn spatial hierarchies of input data.

# CNN Architecture

The basic architecture layers in a typical convolutional neural network are:

Convolution layer

Pooling layer

Fully connected layer

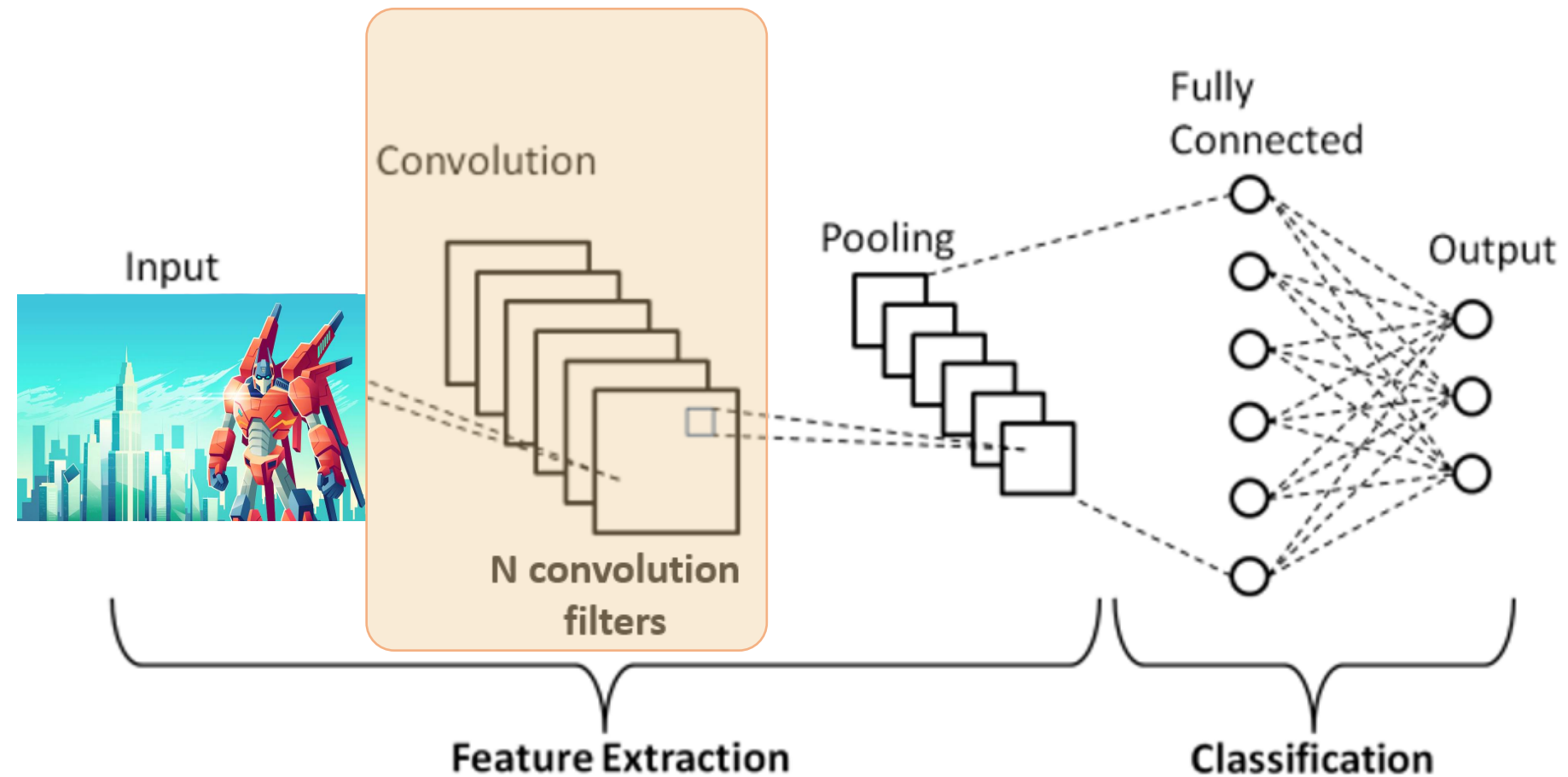
Activation function

Output layer

The parameters used are padding and strides.

# Convolution Layer

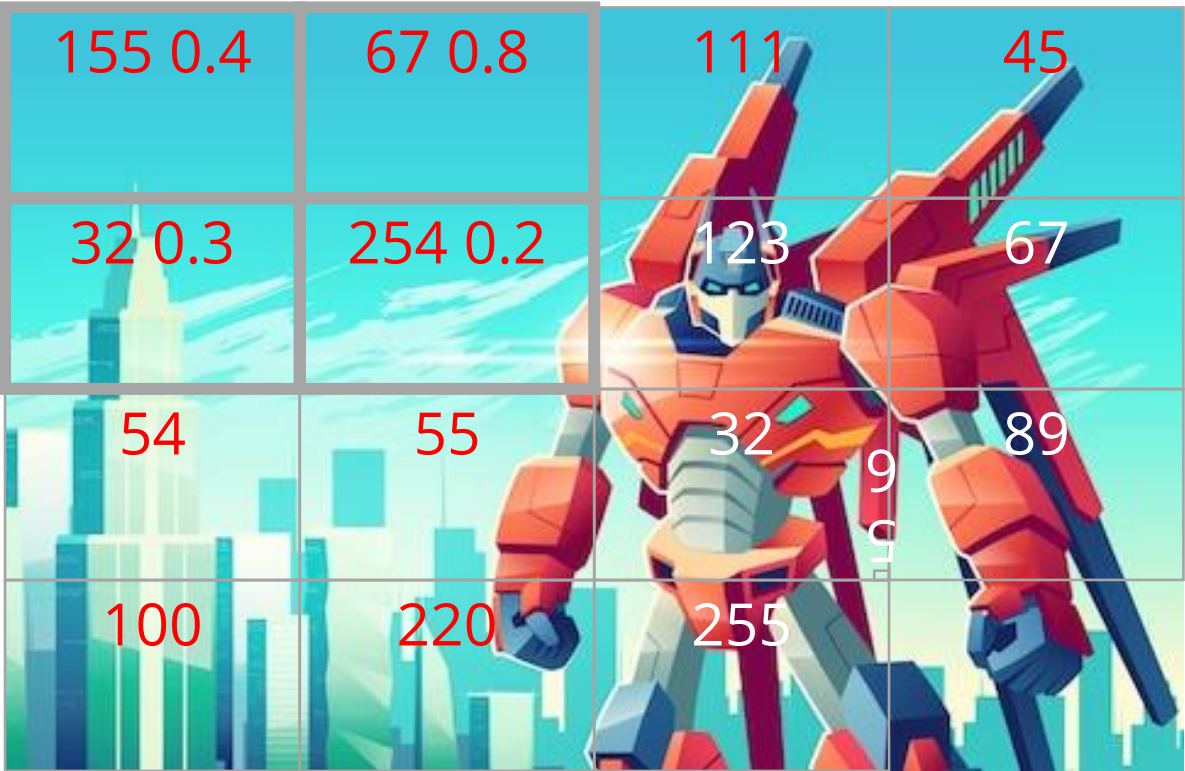
This layer is responsible for performing convolution operations on the given image and filter.



The convolution operation is the sum of the product of the filter values with pixel values.

# Convolution Layer

The convolution operation is performed on each patch of the image using a filter, and a feature matrix is created.



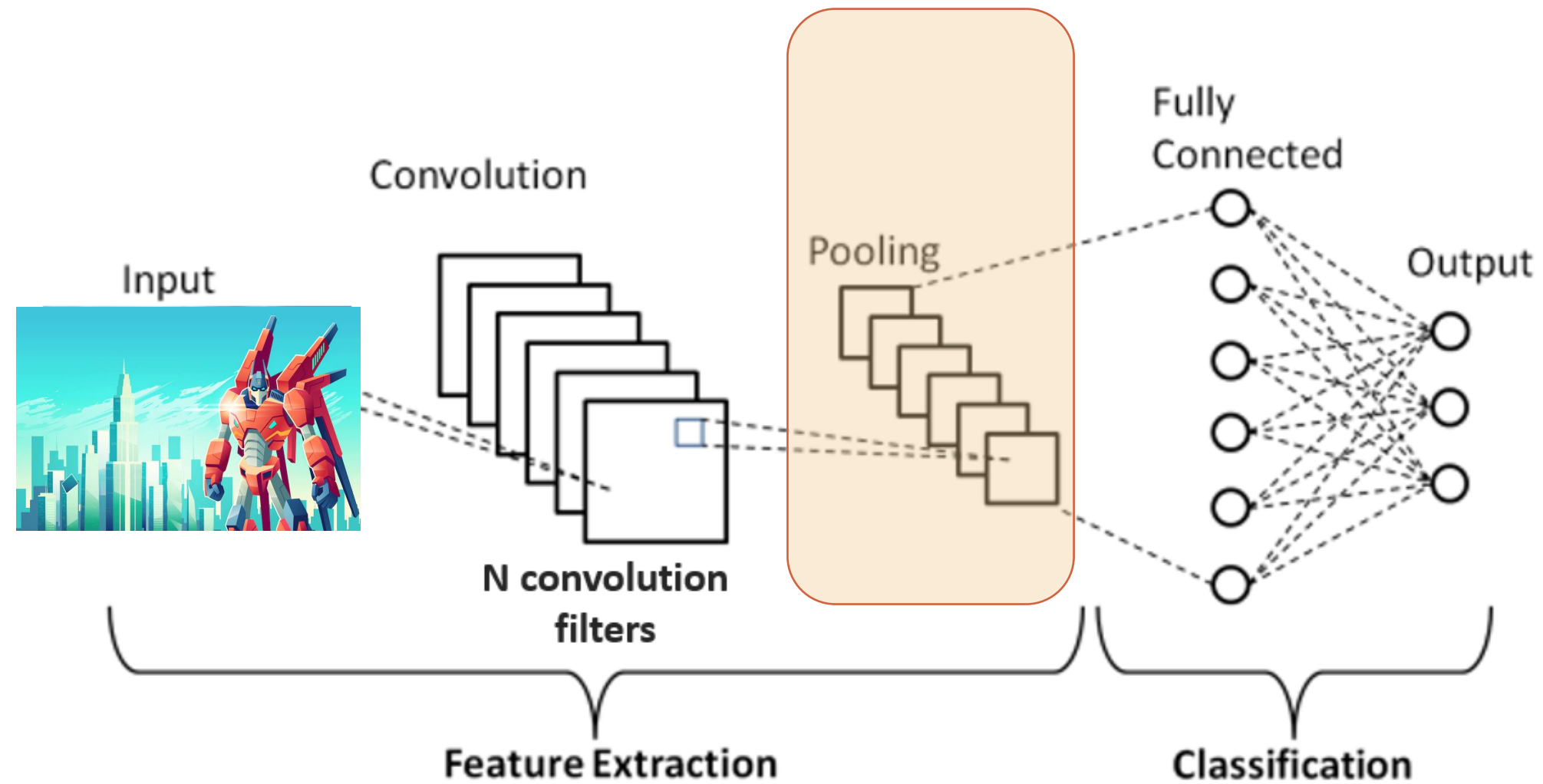
Convolution operation on the input image

176	216.4	113.2
144	112	167
156	132	118

Convolution output:  
Feature matrix

# Pooling Layer

This layer is responsible for calculating the largest value in each patch of the feature map.



# Pooling Layer

Some common functions used in the pooling operation are:

Maximum pooling

Calculates the maximum value for each patch of the feature map

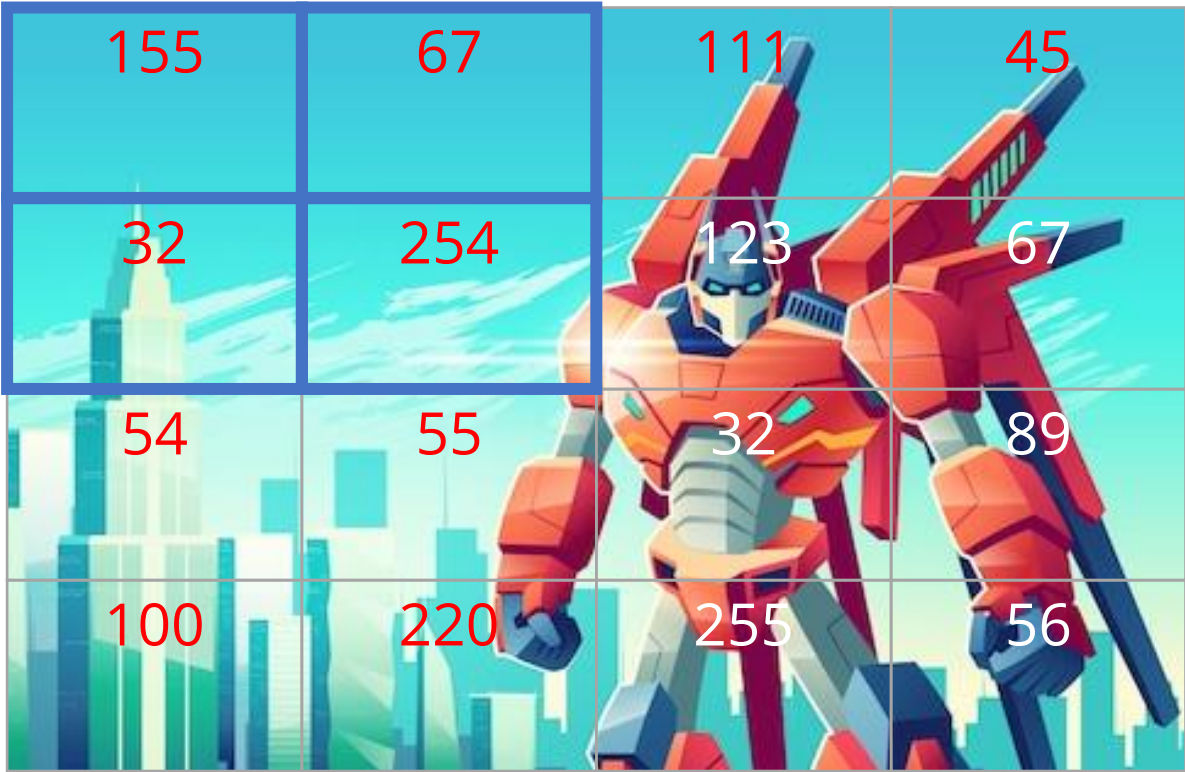
Average pooling

Calculates the average value for each patch on the feature map



# Pooling Layer

In the following image, a max-pooling window of shape (2,2) is used to calculate the maximum value in each (2,2) patch of the image.



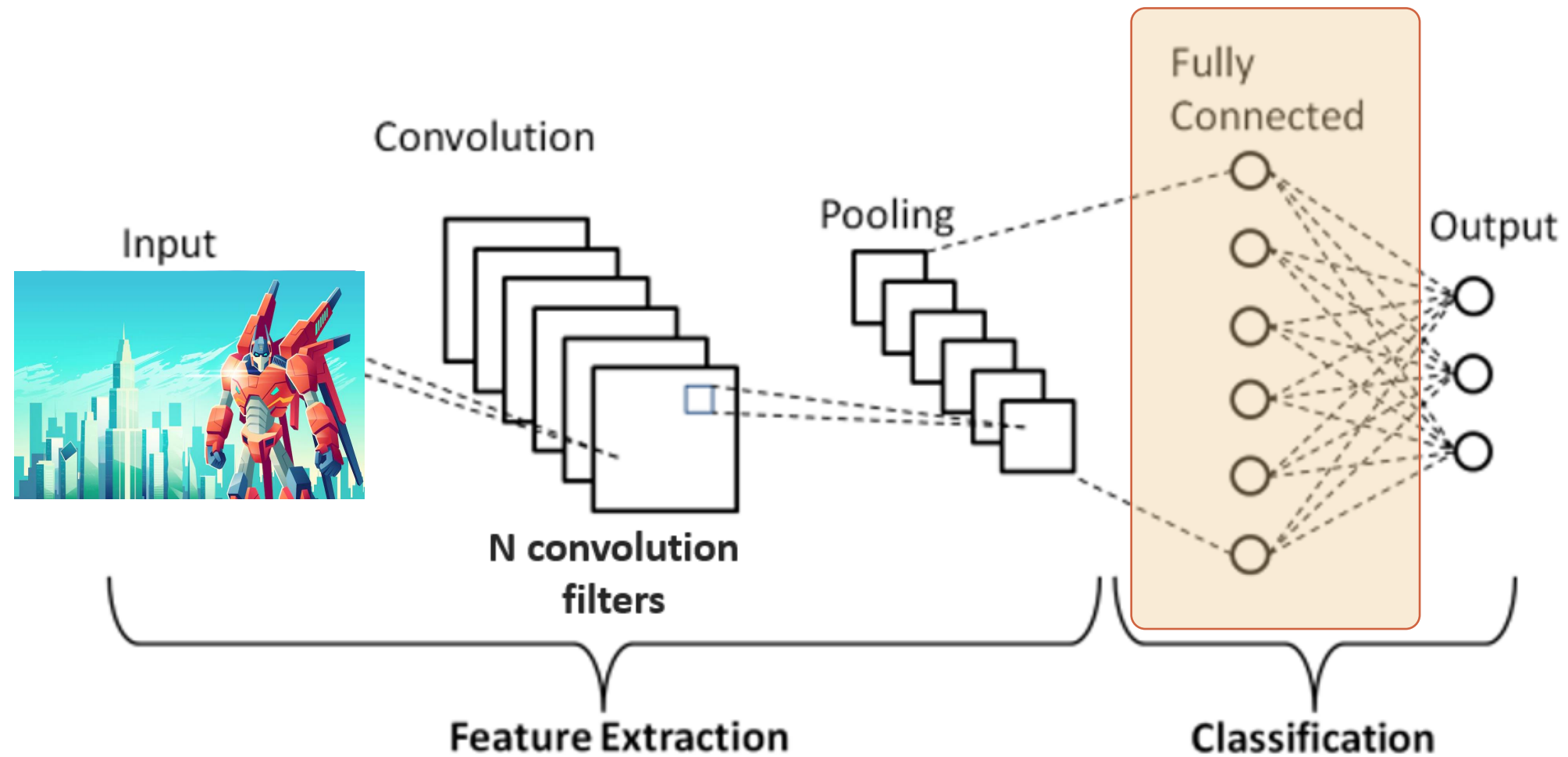
Max-pool operation on the input image

155	254	123
254	254	123
220	255	255

Max-pool output

# Fully Connected Layer

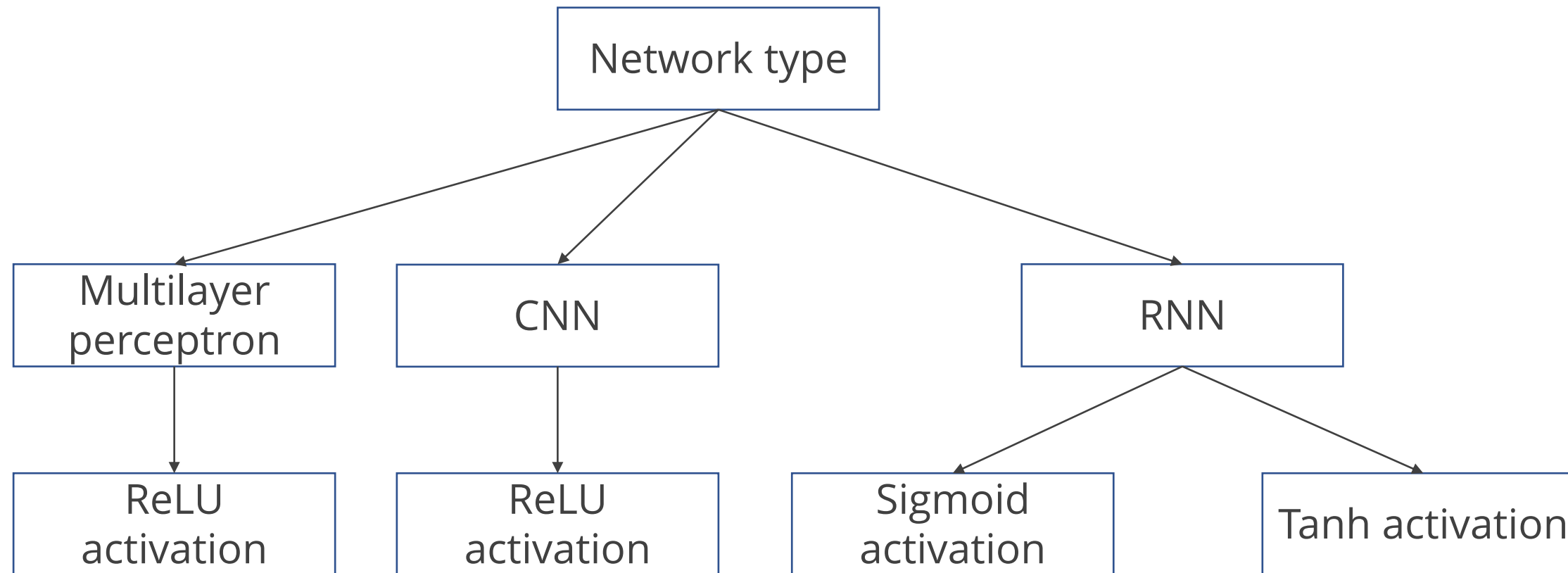
It is a feed-forward neural network and forms the last few layers in the network.



The input to this layer is the output from the final pooling or convolutional layer.

# Activation Function

The activation function calculates the weighted sum and adds bias to decide if a neuron should be activated.



It introduces nonlinearity into the neuron's output to perform more complex tasks.

# Output Layer

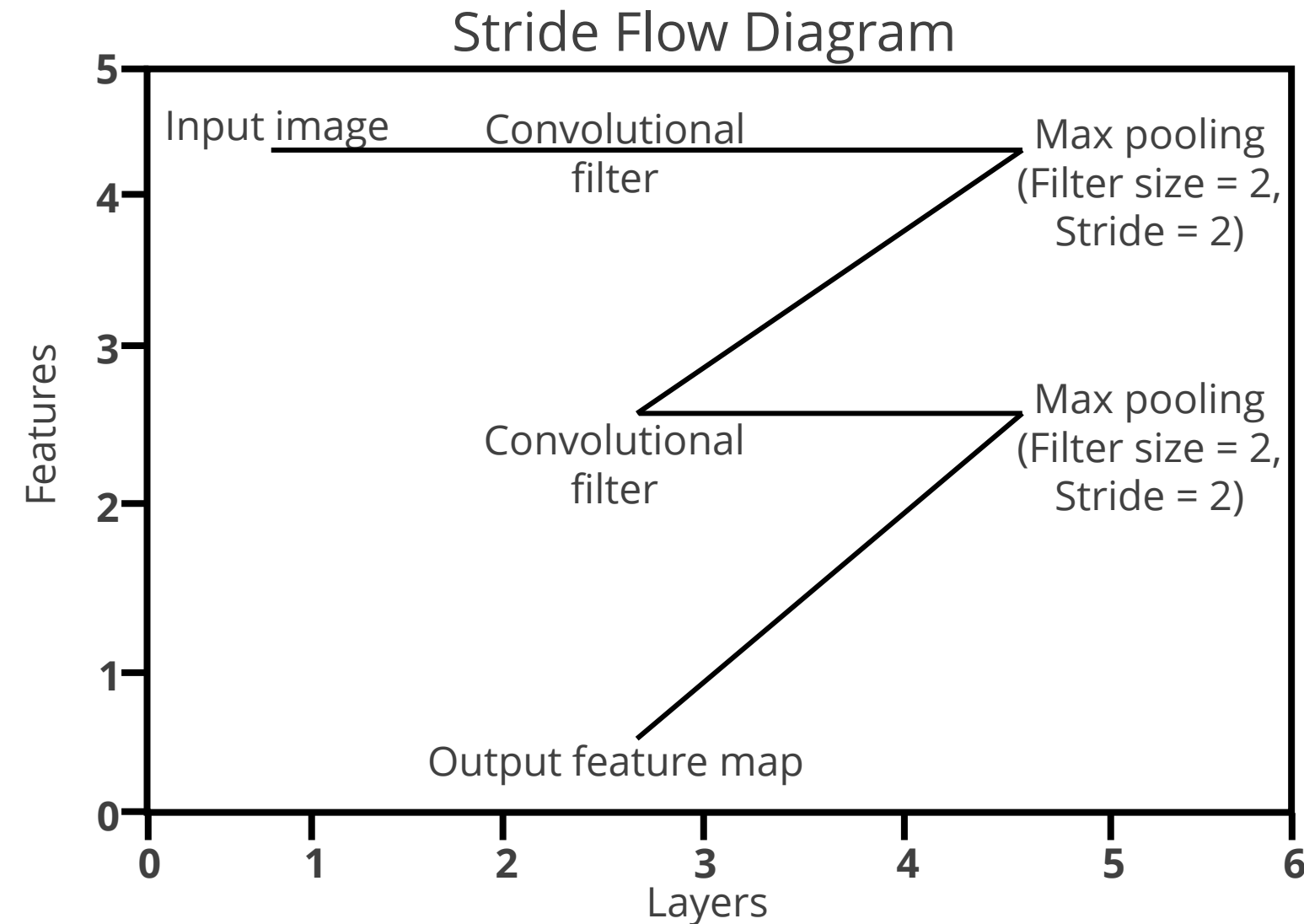
The output layer in a CNN is responsible for producing the final predictions or classifications based on the extracted features from the previous layers.

It consists of one or more fully connected layers, followed by an activation function such as Softmax for classification tasks.

The output layer's weights are learned through backpropagation during training to minimize loss and improve prediction accuracy.

# CNN Architecture Parameters: Strides

The movement of the processing window is controlled by the stride in CNN operations like convolution and max-pooling.

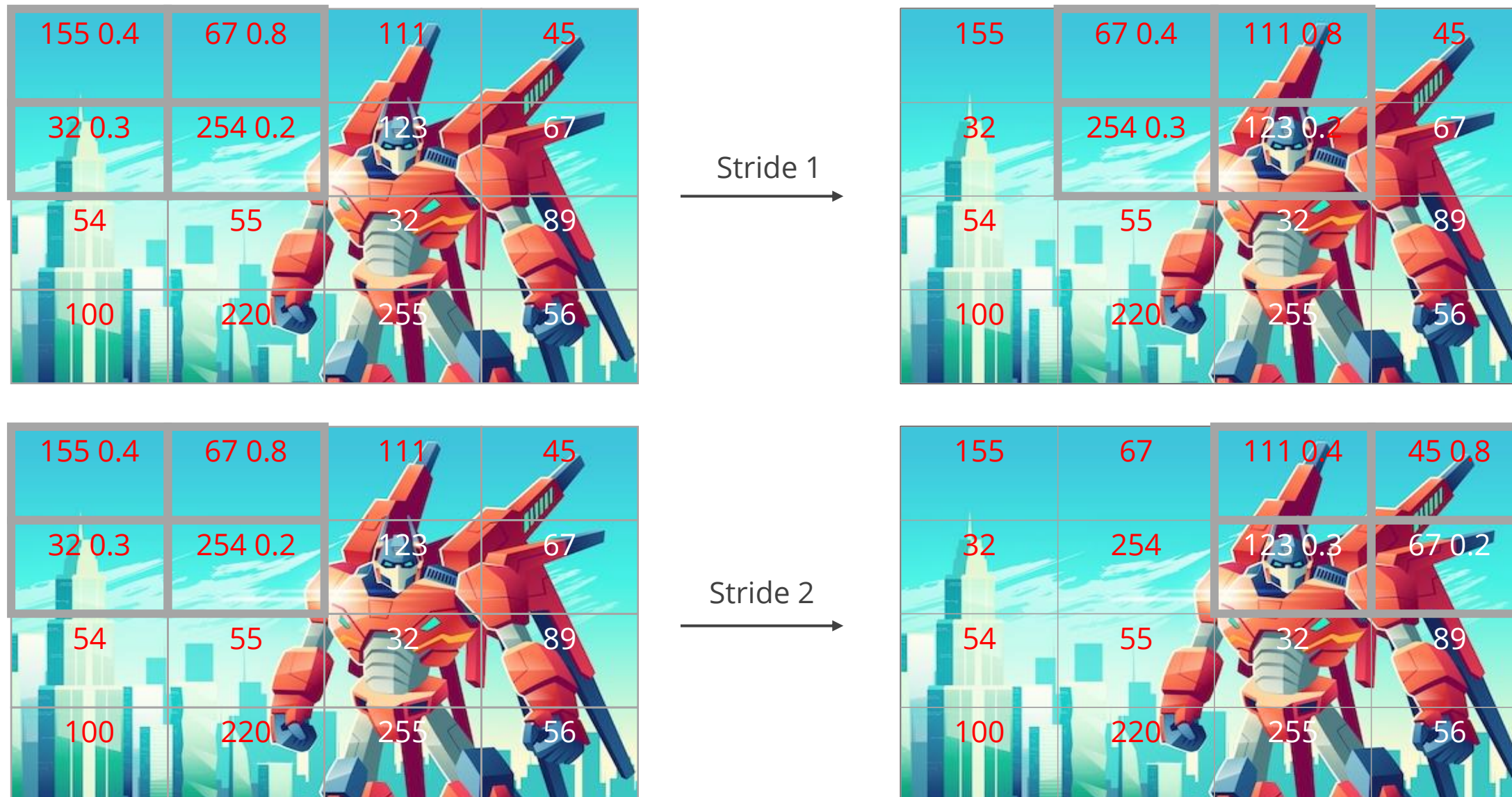


It denotes the number of pixel shifts across the input matrix. This affects the output size and may have an impact on the computational effectiveness and amount of detail in the generated feature maps.



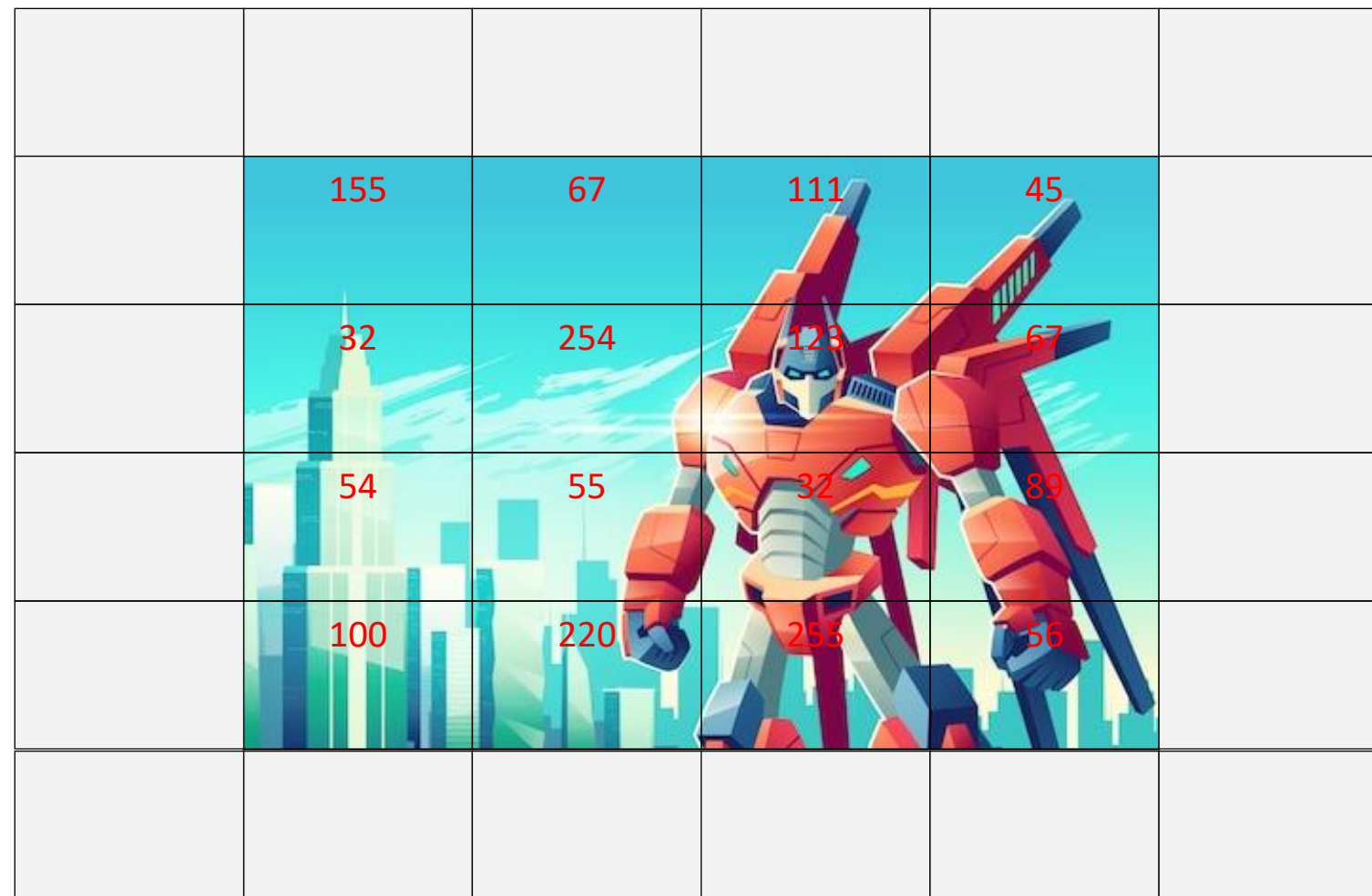
# CNN Architecture Parameters: Strides

For example, if the stride is set to 1, then the kernel moves horizontally and vertically by one pixel.



# CNN Architecture Parameters: Padding

It adds zeros to the input matrix symmetrically to make the shape of the output matrix the same as the input.



The gray area denotes the values padded with zero.

It reduces image shrinkage and increases image analysis accuracy.

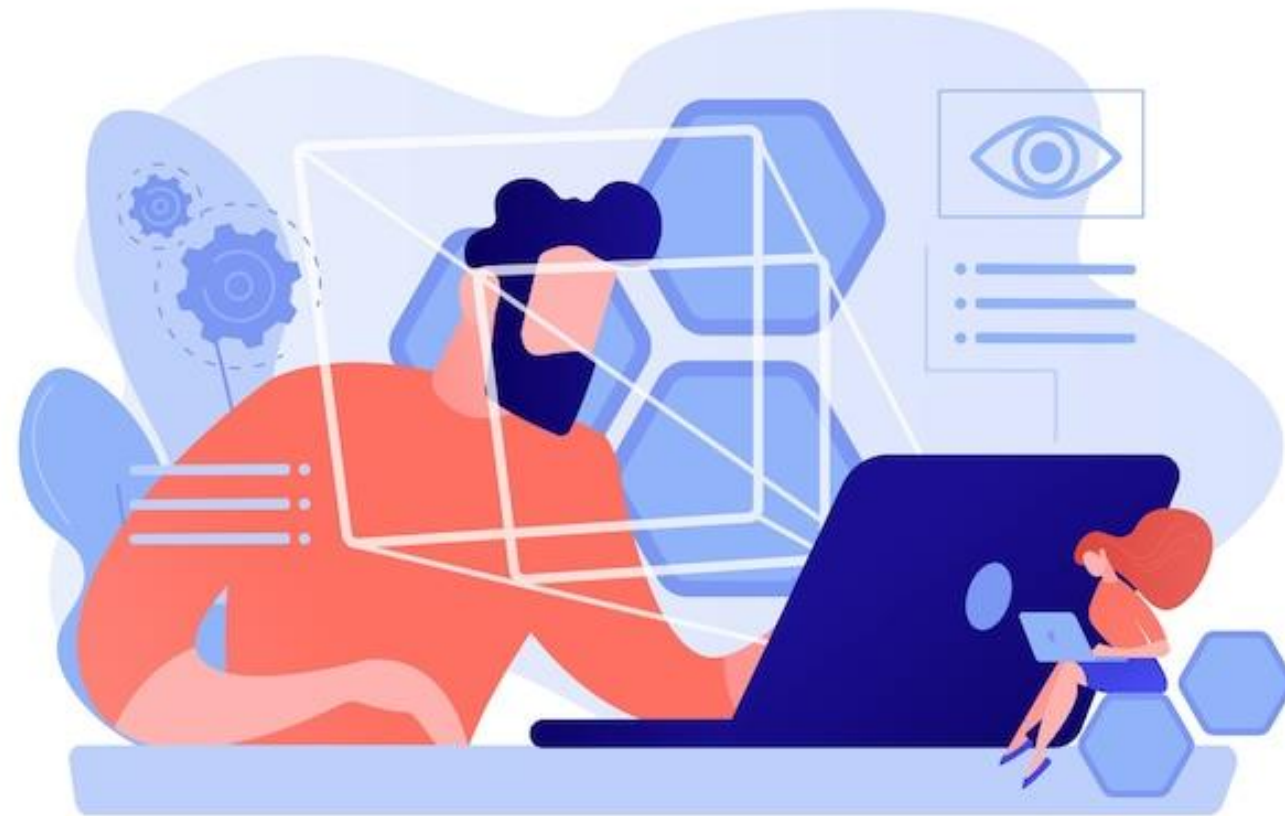


# ResNet



# ResNet

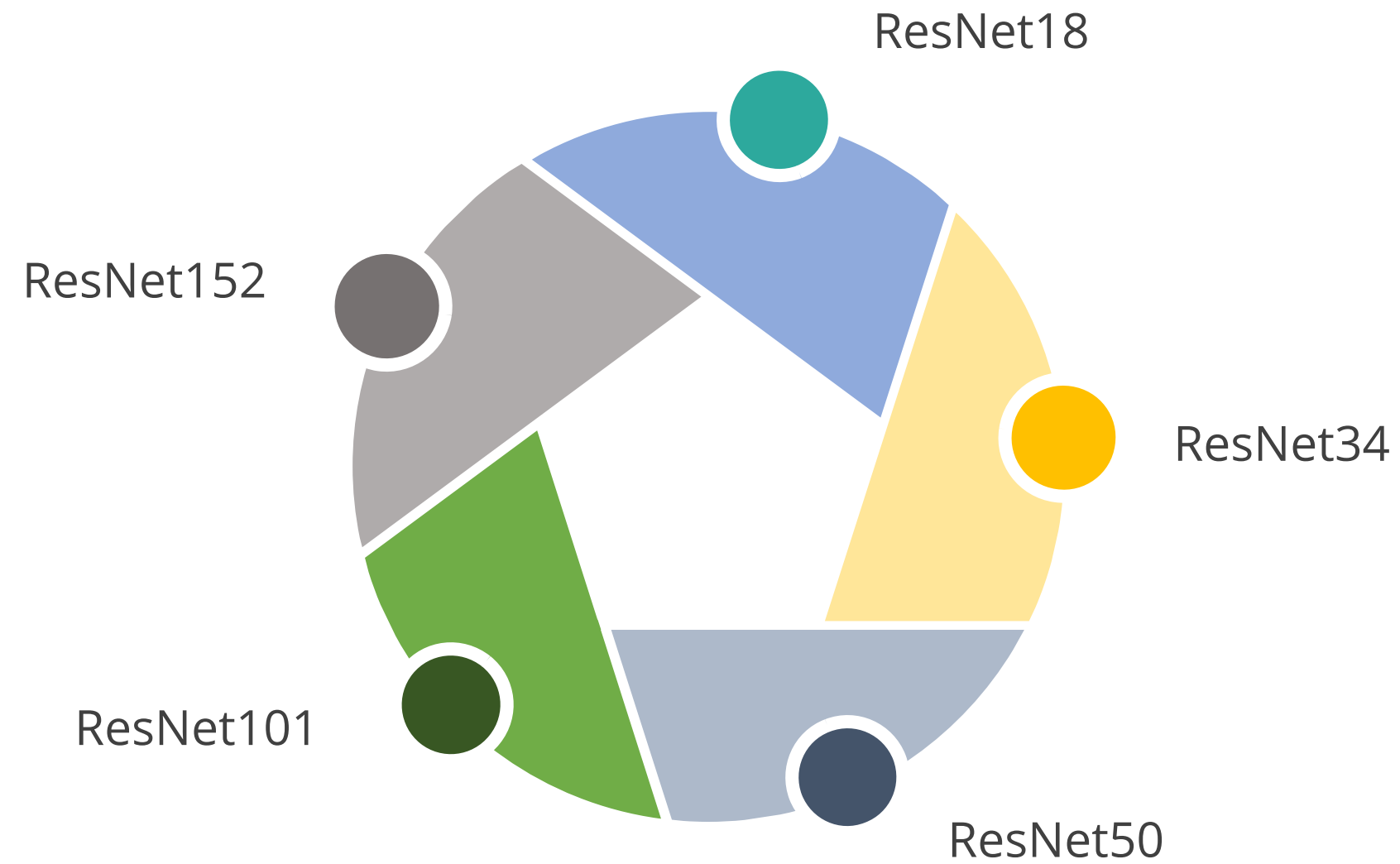
Residual neural network (ResNet) is a convolutional neural network architecture widely used in computer vision tasks.



It supports the construction of neural networks with thousands of convolutional layers.

# ResNet

ResNet has many variants such as:



The numbers represent the total number of layers in the neural network.

# ResNet50

It has 50 layers, which include 48 convolutional layers, one Max-pool layer, and one average pool layer.

It won the ILSVRC image classification challenge in 2015.

It is the only deep CNN architecture considered more efficient, with better performance.

# ResNet50

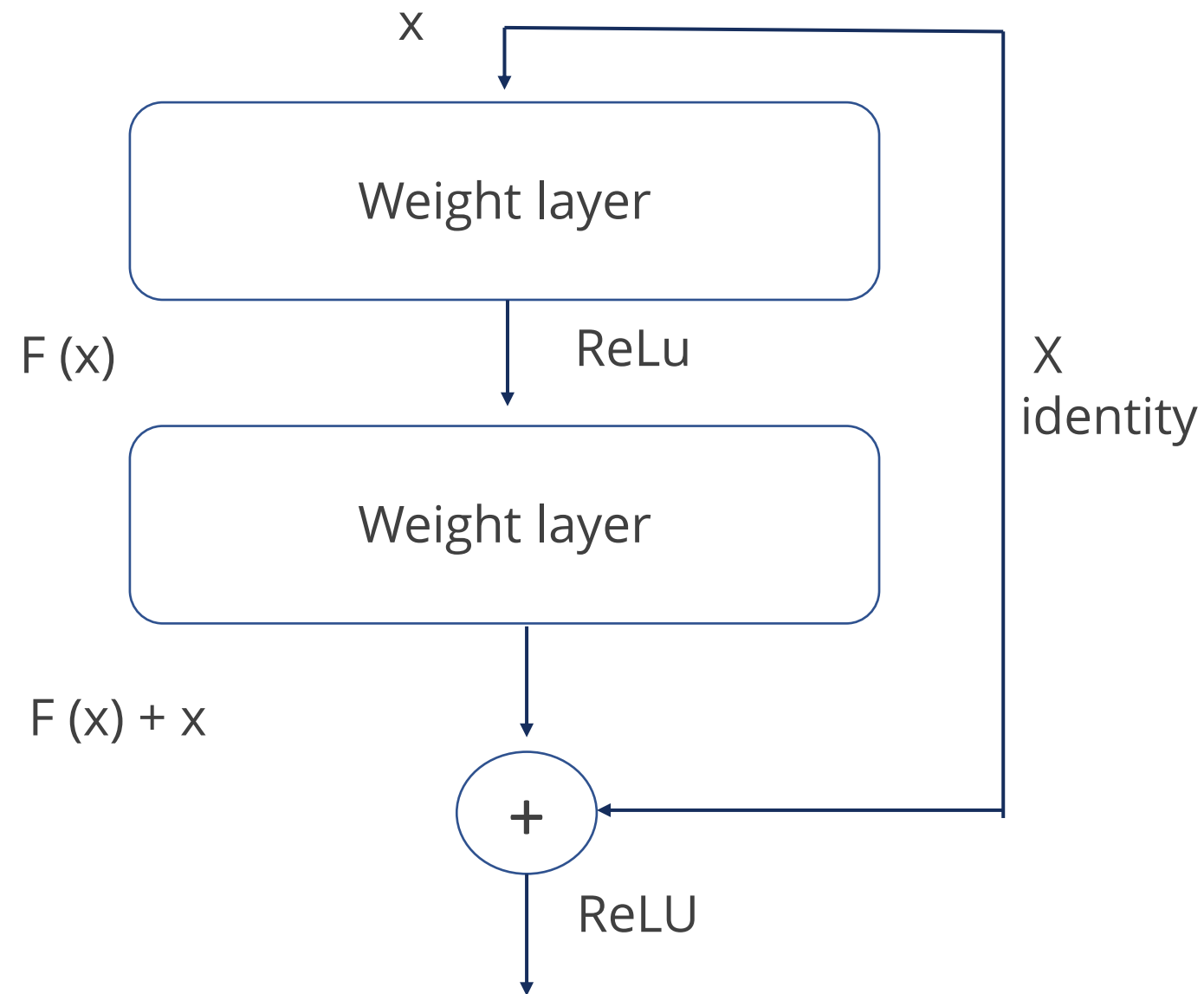
It's a deep network and uses residual connections.

Residual connections solve the problem of deeper model optimizations with representation power.

Residual connections copy the learned representations from a shallower model and add additional layers to establish identity mapping.

# ResNet50

An architectural feature called a skip connection enables direct feeding of the input to a later layer or set of layers.



This makes it easier for the model to learn complicated patterns and helps to solve the vanishing gradient problem in deep networks.

# ResNet50

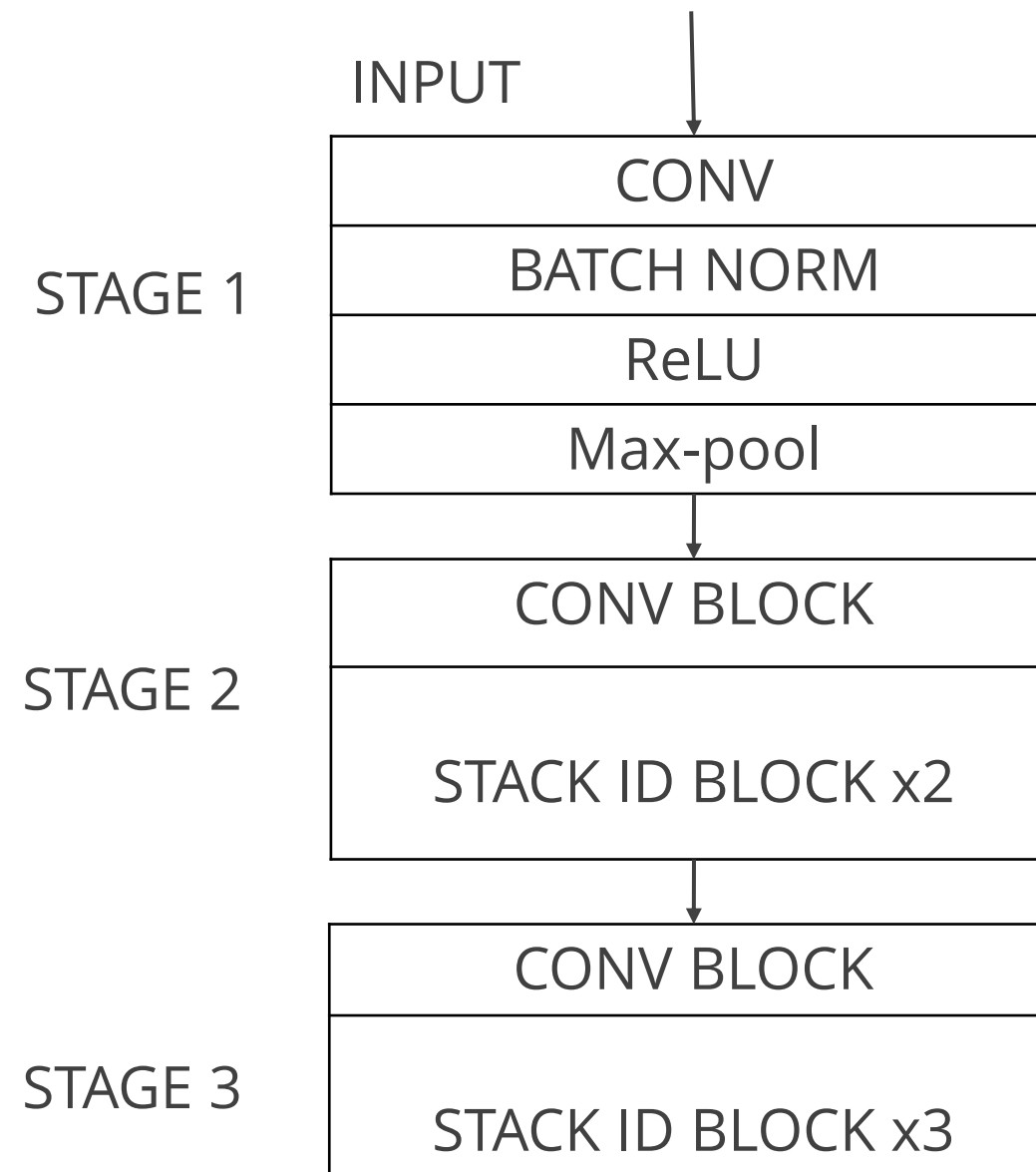
It takes the learnings from the earlier layers, passes their outputs further down, and sums them element-wise with the outputs from the skipped layers.

With the addition of a skip connection, the output  $H(x)$  is defined as the sum of the original input  $x$  and the transformed output  $f(x)$ .

Skip connections do not introduce any extra parameters or computational complexity.

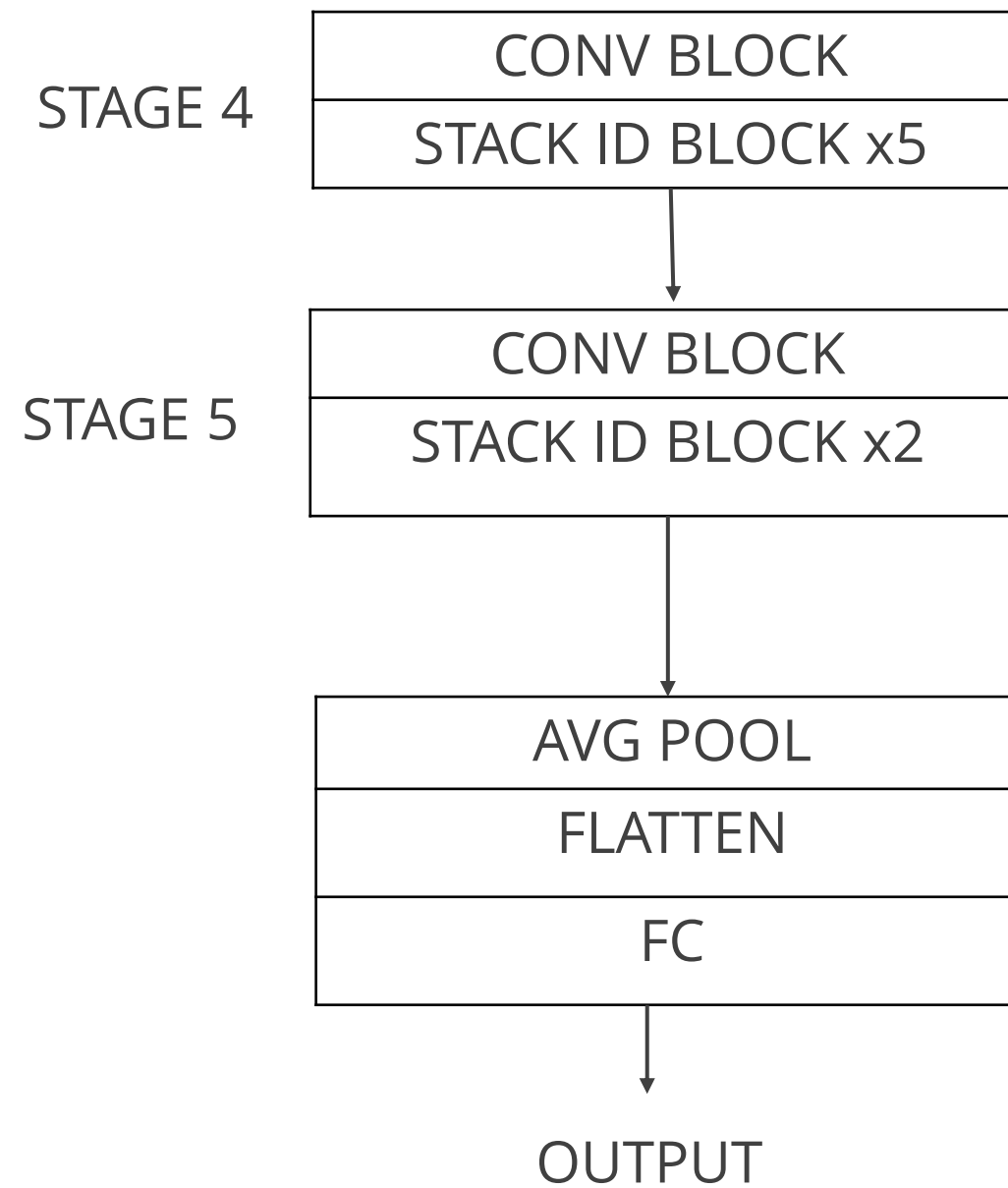
# ResNet50 Architecture

In the ResNet50 architecture, the residual blocks are stacked to improve representation power in further layers.



# ResNet50 Architecture

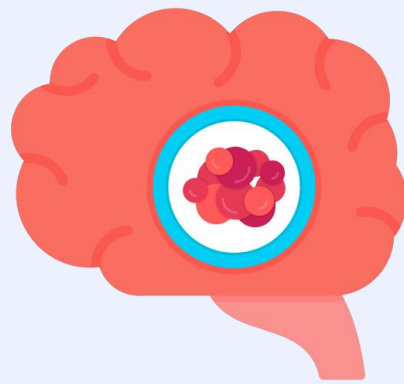
Every residual block has two 3x3 convolution layers.



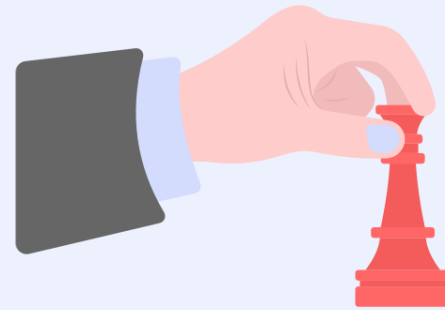


# ResNet Use Cases

Due to its large number of layers and residual connections, ResNet solves almost any computer vision problem with ease.



Detects brain tumors  
based on patients' brain  
MRI scan images



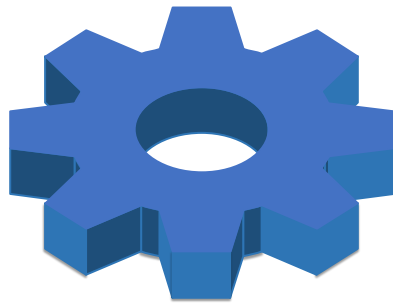
Recognizes player activities  
in games to produce  
equally challenging bots



Recognizes human  
emotions to understand  
their behaviour

# Commonly Used CNN Architectures

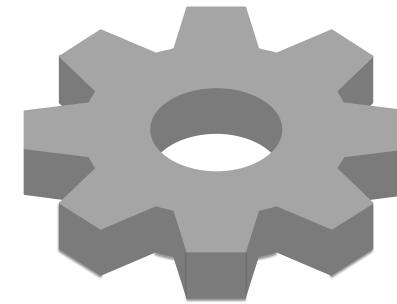
The well-known architectures of Convolutional Neural Networks (CNNs) are as follows:



VGG16



Alex Net



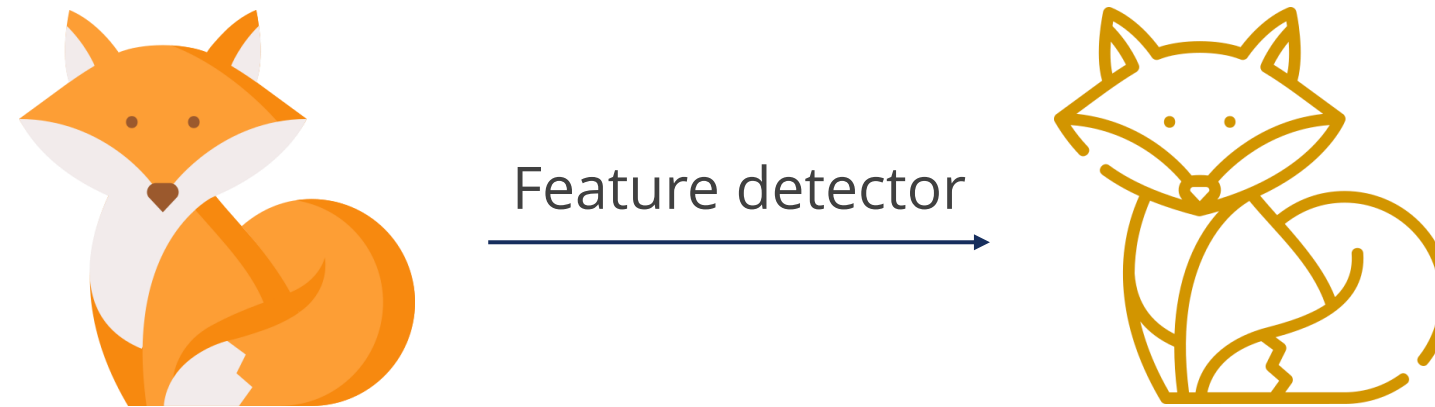
Google net or  
Inception



## Filters in CNN

## Filters in CNN

Filters detect spatial patterns, such as edges, by detecting changes in the intensity values of an image.

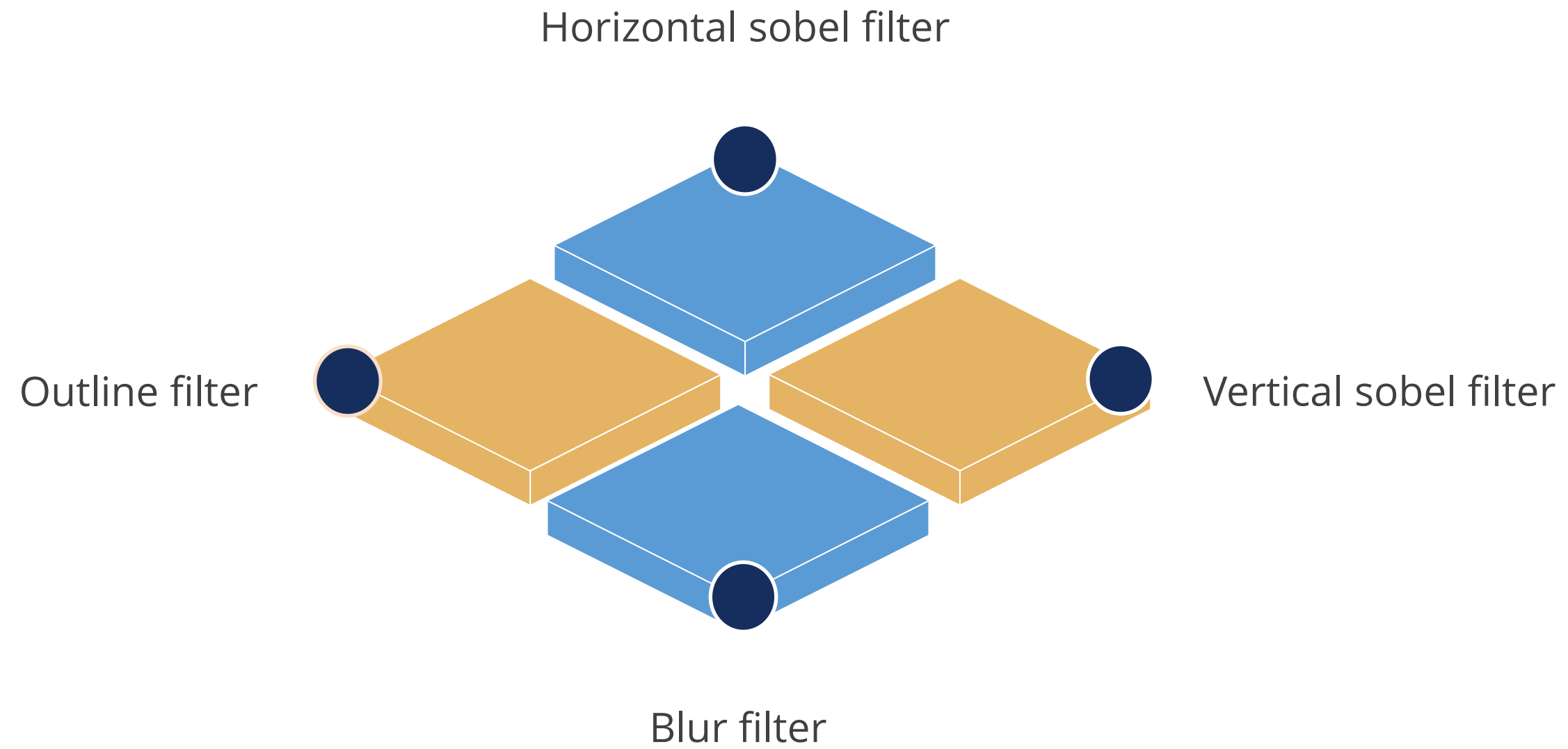


The convolution operation is the sum of the products of the filter matrix and the pixel values of the image.

The type of filter used affects the output produced after convolution.

# Filters in CNN

The different types of CNN filters are:



# Horizontal Sobel Filter

It is responsible for detecting edges horizontally in an image.

The filter matrix used here is:

$$\begin{bmatrix} -1 & -2 & -1 \\ 0 & 0 & 0 \\ 1 & 2 & 1 \end{bmatrix}$$

# Horizontal Sobel Filter

Consider the image and the output shown to understand its working



Sobel edge  
detector



# Vertical Sobel Filter

It is responsible for detecting edges vertically in an image.

The filter matrix used here is:

$$\begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix}$$



# Vertical Sobel Filter

Consider the image and the output shown to understand its working



Edge detection



# Blur Filter

It is responsible for blurring the image.

The filter matrix used here is:

$$\begin{bmatrix} 0.0625 & 0.125 & 0.0625 \\ 0.125 & 0.25 & 0.125 \\ 0.0625 & 0.125 & 0.0625 \end{bmatrix}$$

# Blur Filter

The image shows a blurred output:

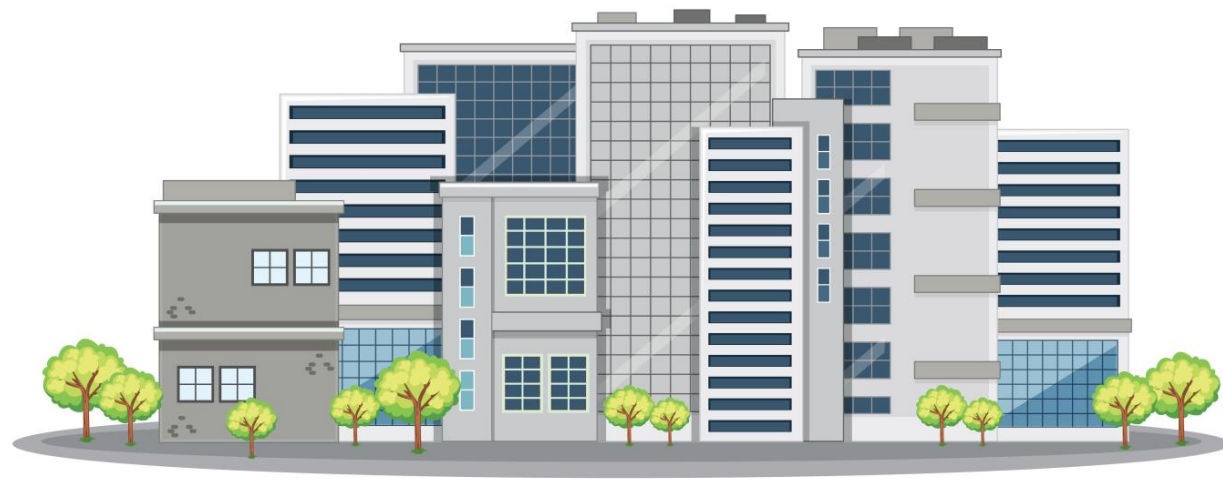


Image softening



# Outline Filter

It is responsible for detecting the outline of objects in an image.

The filter matrix used here is:

$$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$

# Outline Filter

The image shows the output with an outline filter.



Boundary  
extraction

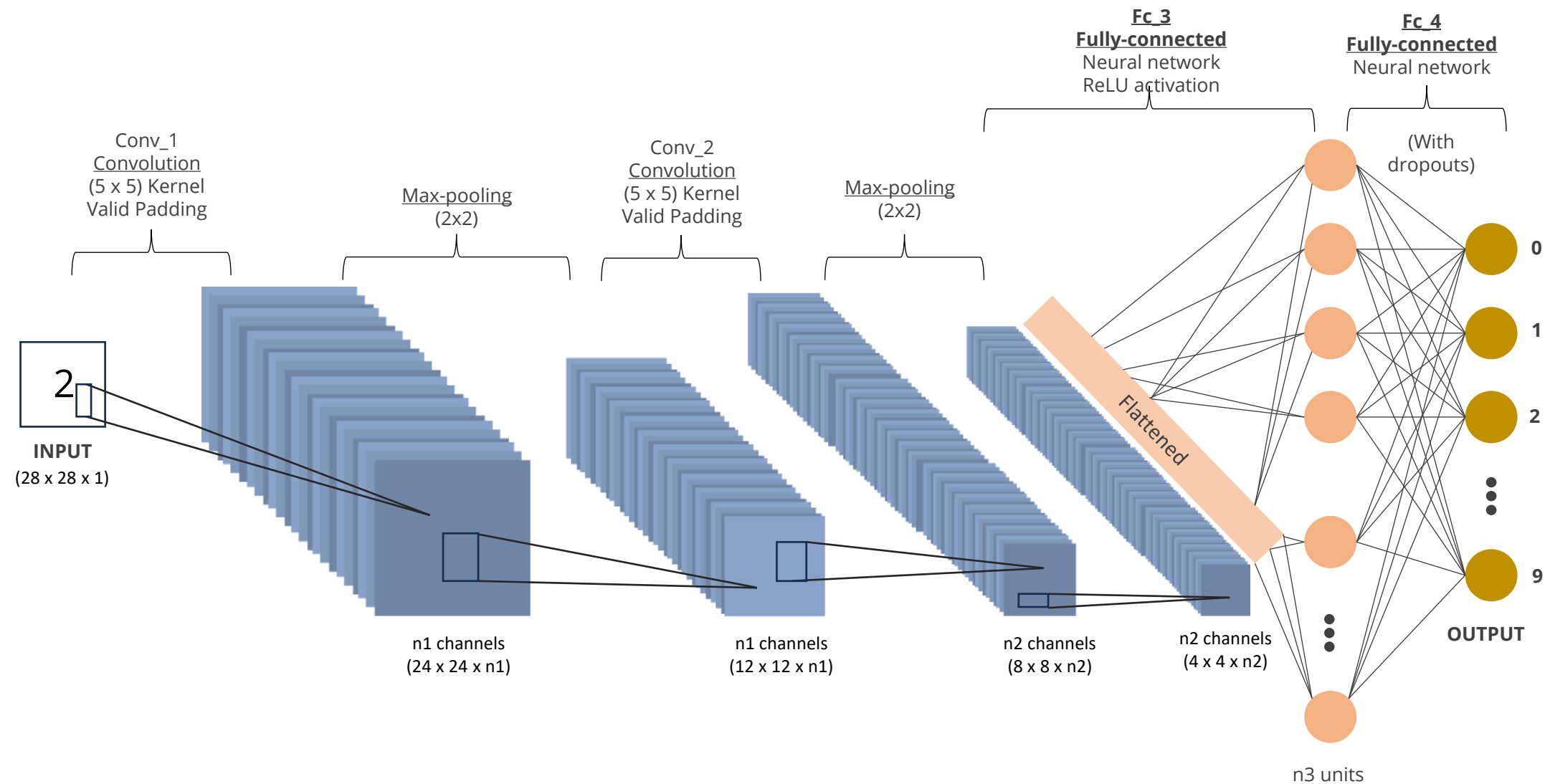




## Working of CNNs

# Working of CNNs

In a CNN, each filter's optimum value is acquired during training and not explicitly defined.

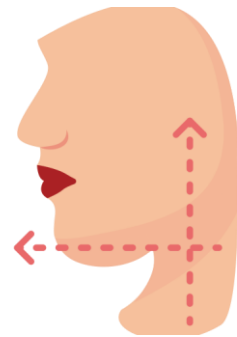


Humans can recognize and extract meaning from images by utilizing learned filters in the visual cortex.

# Working of CNNs

CNN filters learn detection of abstract concepts.

## Example



The first convolution may detect features like a person's shoulders or face contours.



The second convolution may detect the eye shape and the edges of the shoulder.



# Working of CNNs

CNNs can provide more abstract and comprehensive information by stacking convolutional layers atop each other.



They perform hierarchical feature learning, like the human brain's image recognition.

The problems to be solved and the features to be learned determine the convolutions to be applied.

## 2D Convolution Layer

The 2D convolution (conv2D) is the most common type of convolution applied.

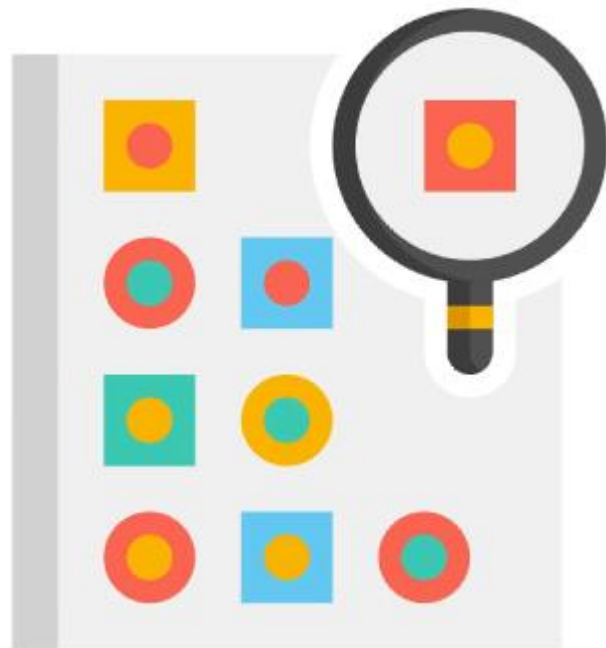


A filter or kernel has a height and a width, often smaller than the input image, and it slides over the entire image.

The receptive field is the image area where the filter is applied.

## 2D Convolution Layer

The conv2D filters extend across the three-color channels (RGB), with different filters for each channel.



- ◆ The individual channel convolutions are combined to produce the concluding image.
- ◆ The filters undergo random initialization and distribution to ensure they learn differently.
- ◆ The filters eventually learn to detect various image aspects.

## 2D Convolution Layer

Multiple conv2D filters are used in a single layer to recognize distinct features, as each filter learns a different feature.

Each filter serves as an input to the neural network's next layer.

### Example

If the first layer has 8 filters and the second has 32, each second layer filter sees 8 inputs, resulting in a 32 by 8 feature map.

A single output from each layer is obtained after the application of each of the eight feature maps of a single filter.

# Constraints

Though highly accurate, the conv2D layer has some drawbacks.

It is computationally expensive.

A large conv2D filter compilation is time-consuming, and stacking multiple filters in layers increases the number of calculations.

# Constraints

To overcome the constraints, the filter size can be reduced and the strides increased.



However, the filter's effective receptive field and the quantity of data it can capture are reduced.

## Assisted Practices



Let's understand the concept of CNN for image classification using Jupyter Notebooks.

- 7.08\_Image classification using CNN

**Note:** Please refer to the Reference Material section to download the notebook files corresponding to each mentioned topic

# Pooling in CNN





## Discussion

# Discussion: Pooling in CNN

Duration: 10 minutes

- What is a CNN?
- What are some real-world applications of CNNs?



# Pooling in CNN

The pooling operation involves sliding a two-dimensional filter over each channel of the feature map.



It summarizes the features that lie within the region covered by the filter.

# Pooling in CNN

For a feature map with dimensions  $nh \times nw \times nc$ , the output dimensions obtained after pooling are:

$$(nh - f + 1) / s \times (nw - f + 1) / s \times nc$$

Where,

Nh: Height of the feature map

Nw: Width of the feature map

nc: Number of channels in the feature map

F: Size of the filter

S: Stride length

A common CNN model architecture contains multiple stacked convolution and pooling layers.

# Use of Pooling in CNN

Pooling layers reduce the feature map dimensions.

It reduces the number of parameters to learn and the amount of computation in the network.

It summarizes the features present in a region of the feature map generated by a convolution layer.

It performs further operations on the summarized features instead of precisely positioned features.

It exhibits increased resilience to variations in the positions of features within the input image.

# Types of Pooling Layers

There are three types of pooling layers:



Max-pooling

Average pooling

Global pooling

# Max-pooling

It selects the maximum element from the region of the feature map covered by the filter.

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

Max-pool

Filter: (2x2)  
Stride: (2, 2)

9	7
8	6

The output of this layer is a feature map with the most prominent features of the previous feature map.

# Average Pooling

It computes the average of the elements present in the region of the feature map.

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

Average pool

Filter: (2x2)  
Stride: (2, 2)

4.25	4.25
4.25	3.5

It gives the average of the features present in a patch.



# Global Pooling

It reduces each channel in the feature map to a single value.

An  $n_h \times n_w \times n_c$  feature map is reduced to a  $1 \times 1 \times n_c$  feature map, equivalent to using a filter of dimensions  $n_h \times n_w$ .

It can either be global max-pooling or global average pooling.

# Discussion: Pooling in CNN

Duration: 10 minutes



- What is pooling in a CNN?

**Answer:** Pooling in a CNN is a technique used to downsample feature maps, reducing their spatial dimensions while retaining important information.

- What are the benefits of pooling in a CNN?

**Answer:** The benefits of pooling in a CNN include dimensionality reduction, translation invariance, feature extraction, and parameter reduction.



# Introduction to TensorBoard

# TensorBoard

It offers a web-based interface enabling visualization of diverse aspects related to model performance, data exploration, and real-time monitoring of training progress.



It can display image, text, and audio data and aims to decrease the complexity of neural networks.

# TensorBoard

TensorBoard visualizations can be graphs and histograms that are used to interpret the results of:

Loss

Accuracies

Other metrics from the model

# TensorBoard

To understand the concept better, consider the following images with two types of concrete:



Plain



Marred

Image Source: <https://data.mendeley.com/datasets/5y9wdsg2zt/2>

# TensorBoard

Construct a classifier to detect the two types of surfaces

The classifier classifies surfaces on construction sites to understand the withstanding capacity of buildings.

It can understand if surfaces are damaged due to earthquakes or natural calamities.

# TensorBoard

It can be initialized using the following command:

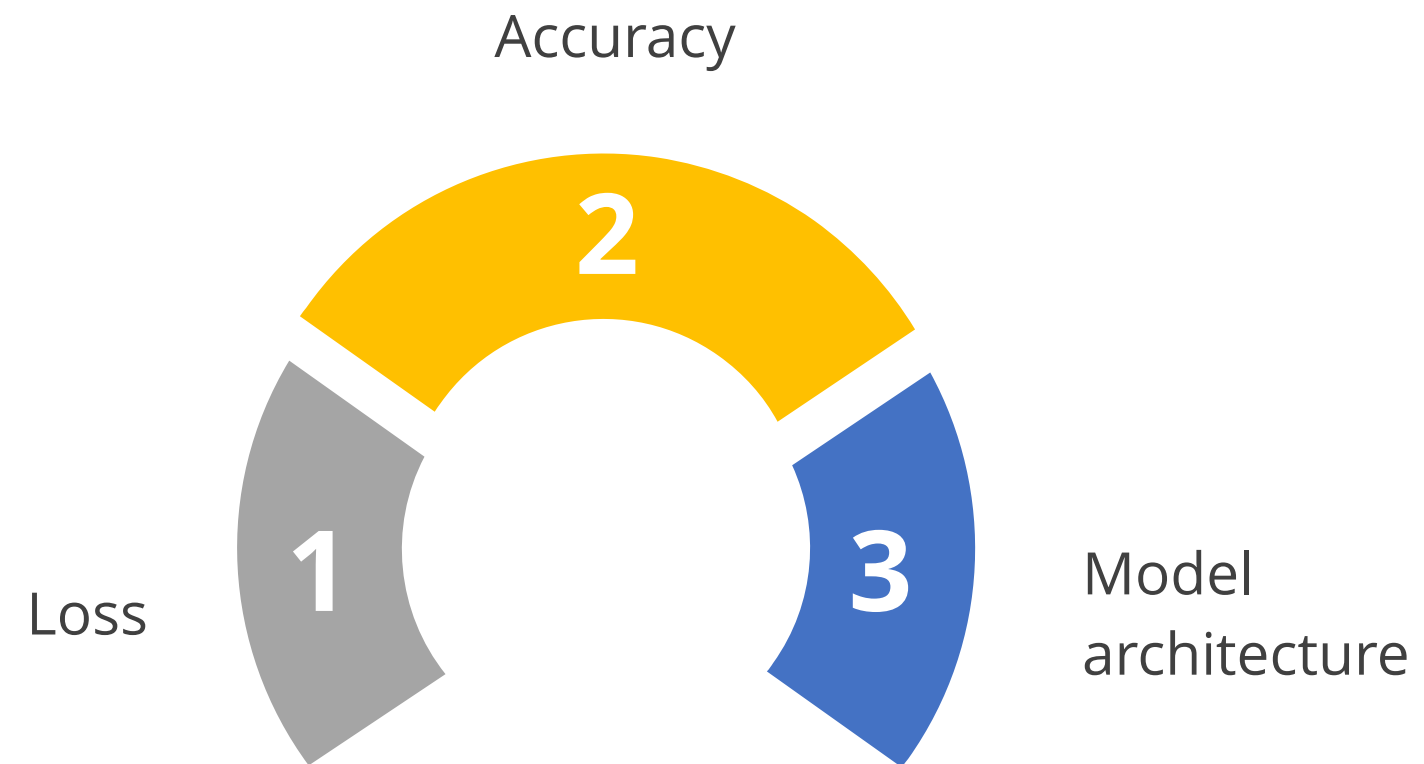
```
tensorboard --logdir path_to_logdir
```

logdir is where the logs of the model are saved during training.



# TensorBoard

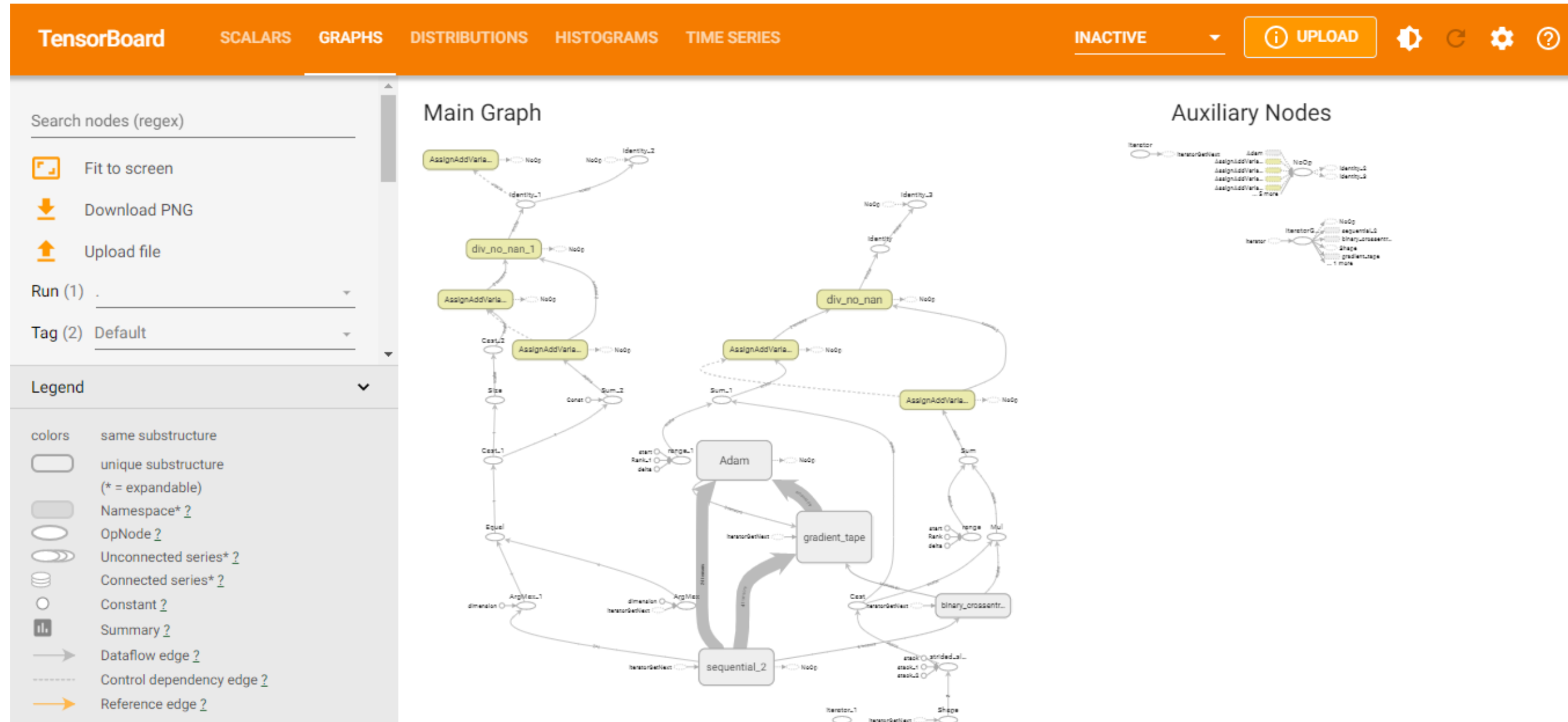
The logs of the model consist of:



It can also include custom metrics.

# TensorBoard

Consider the following image:



# TensorBoard

The image has five sections:

Scalars

The graphs of metrics, such as training loss and accuracy, are saved

Graphs

The model architecture is clearly portrayed

Distributions

The weight distribution along each layer is shown

Histograms

The histogram plots with a frequency of 1 along each layer are displayed

Time series

The distribution along each layer is checked, along with the time

These details help one understand and debug the machine learning model under study thoroughly.

# Assisted Practices



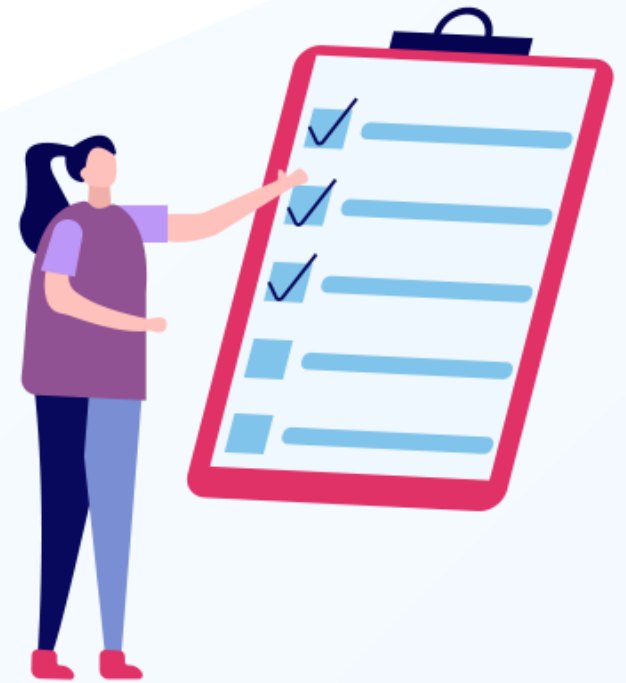
Let's understand the concept of introduction to TensorBoard using Jupyter Notebooks.

- 7.10\_Introduction to TensorBoard

**Note:** Please refer to the Reference Material section to download the notebook files corresponding to each mentioned topic

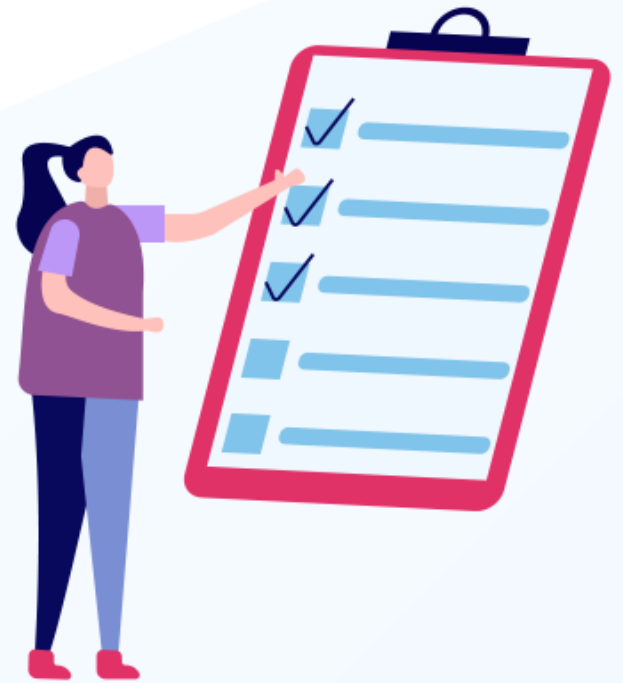
## Key Takeaways

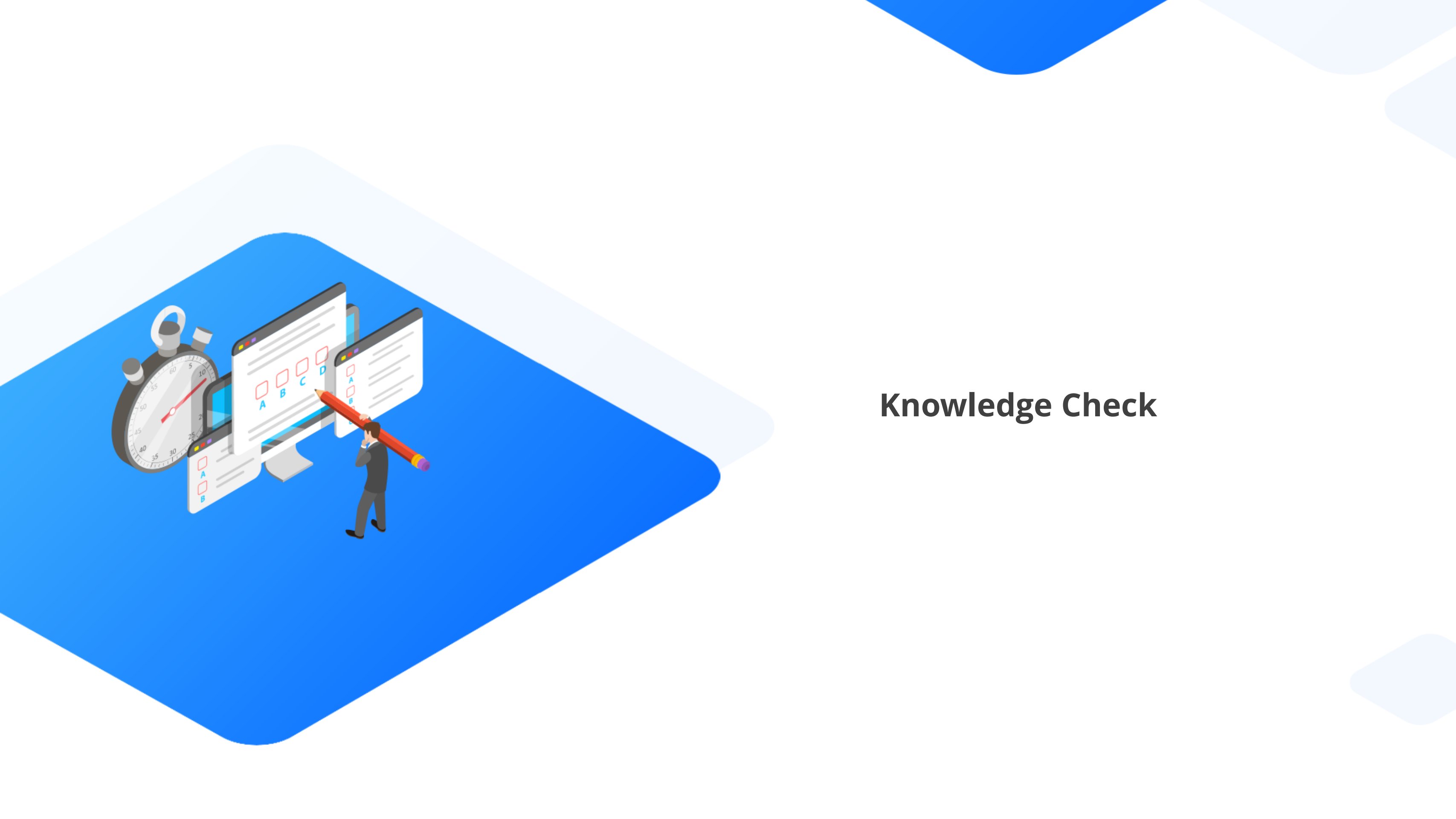
- CNN is a popular algorithm that is used widely in the field of computer vision.
- The convolution operation is the sum of the product of filter values to the pixel values.
- The three essential layers in a CNN are the convolution, pooling, and fully connected layers.
- Residual neural network (ResNet) is a convolutional neural network architecture widely used in computer vision tasks.



# Key Takeaways

- 👁 The different types of CNN filters are horizontal and vertical Sobel, blur, and outline filters.
- 👁 The Conv2D filters extend across the three-color channels (RGB), and the individual channel convolutions are combined to produce the concluding image.
- 👁 TensorBoard is an interface used to visualize, understand, and debug machine learning models.





## Knowledge Check

## Knowledge Check

1

**What does the shape of an image data represent?**

- A. Height and depth of the image
- B. Width, depth, and brightness of the image
- C. Height, width, and brightness of the image
- D. Height, width, and channels of the image





## Knowledge Check

1

What does the shape of an image data represent?

- A. Height and depth of the image
- B. Width, depth, and brightness of the image
- C. Height, width, and brightness of the image
- D. Height, width, and channels of the image

---

The correct answer is **D**

---

The shape of an image data represents the height, width, and channels of the image, denoted as (height, width, channels).



## Knowledge Check

2

**What does the convolution operation do in a CNN?**

- A. Flattens the image data
- B. Extracts all the necessary information from the image data
- C. Performs an addition operation on image data
- D. Subtracts the filter values from the pixel values in the image data



## Knowledge Check

2

**What does the convolution operation do in a CNN?**

- A. Flattens the image data
- B. Extracts all the necessary information from the image data
- C. Performs an addition operation on image data
- D. Subtracts the filter values from the pixel values in the image data

---

The correct answer is **B**

---

**In a CNN, the convolution operation extracts all the necessary information from the image data, which is helpful in training models.**



## Knowledge Check

3

**What is the purpose of pooling layers in CNNs?**

- A. To increase the number of parameters to learn and the amount of computation in the network
- B. To reduce the feature map dimensions and amount of computation in the network
- C. To increase the feature map dimensions and amount of computation in the network
- D. To reduce the feature map dimensions and number of layers in the network



## Knowledge Check

3

What is the purpose of pooling layers in CNNs?

- A. To increase the number of parameters to learn and the amount of computation in the network
- B. To reduce the feature map dimensions and amount of computation in the network
- C. To increase the feature map dimensions and amount of computation in the network
- D. To reduce the feature map dimensions and number of layers in the network



---

The correct answer is **B**

---

**Pooling layers reduce the feature map dimensions, which helps reduce the number of parameters to learn and the amount of computation in the network.**

# Lesson-End Project: Image Classifier with CIFAR10



**Problem statement:** Build a deep learning convolutional neural network to recognize characters using the Chars74k dataset

**Objective:** Build a neural network-based classification model to recognize characters using the following metrics:  
Use four convolution layers with a  $3 \times 3$  kernel and activation function as ReLU. Add maximum pooling layers after every other convolution layer and two hidden layers with dropout.

**Access:** Click on the **Lab** tab on the left side of the LMS panel. Copy the generated username and password. Click on the **Launch Lab** button. On the new page, enter the username and password you copied earlier into the respective fields. Click **Login** to start your lab session. A full-fledged Jupyter lab opens, which you can use for your hands-on practice and projects.



**Thank You!**