

CNN

Summary of Topics Covered

1. Image Encoding & CNN Concepts

◆ What is Image Encoding?

- Process of converting an image into a numerical format that can be fed into a machine learning model.
- Output is typically a **tensor** (multi-dimensional array) representing pixel intensities.

◆ Why CNNs?

- Convolutional Neural Networks (CNNs) are **designed to process grid-like data** (like images and audio spectrograms).
- They are **translation-invariant** and capable of learning **hierarchical features** (from edges to complex shapes).

2. CNN Architecture Components

◆ Convolution Layer

- Applies **filters (kernels)** to extract local patterns.
- Each filter outputs a **feature map**.

◆ Activation Function

- Usually **ReLU** (Rectified Linear Unit) is used after each convolution to introduce non-linearity.

◆ Pooling Layer

- **Reduces spatial dimensions** while retaining important information.
- Common methods: **Max Pooling**, **Average Pooling**.

◆ Fully Connected Layers

- Flatten the feature maps and pass them through **dense layers** for classification or regression.
-

3. Image to Feature Conversion using CNNs

- Raw images → Convolutions → Feature Maps → Flatten → Feature Vectors
 - These **learned features** can then be:
 - Used for classification (e.g., digits, objects)
 - Passed into ML models like SVM, KNN, etc.
 - Visualized using PCA/t-SNE for clustering
-

4. Application to Audio via Spectrograms

◆ Spectrograms as Images

- Audio signals are **converted to 2D spectrograms** (time vs. frequency).
- These are treated like grayscale images by CNNs.

◆ Benefit

- Allows the **application of image-based deep learning** techniques to audio data (e.g., voice recognition, music classification).
-

Relevant Additional Insights

Transfer Learning

- You can use **pre-trained CNNs** (e.g., VGG, ResNet) to extract image/audio features without training from scratch.
- Useful when data is limited.

Feature Extraction for ML

- Instead of using CNNs for full classification, you can extract **intermediate layer outputs** (feature vectors) and:
 - Feed them into PCA/t-SNE for visualization.
 - Use clustering (KMeans, DBSCAN).
 - Calculate similarity metrics (Euclidean, Cosine distance).

🔄 Standardization

- Always **normalize image/audio inputs** (e.g., scale pixel values from 0–255 to 0–1).

✅ Overall Workflow

[Raw Image / Audio]



[Preprocessing (Resizing, Grayscale, Spectrogram)]



[CNN Layers: Convolution → ReLU → Pooling → FC]

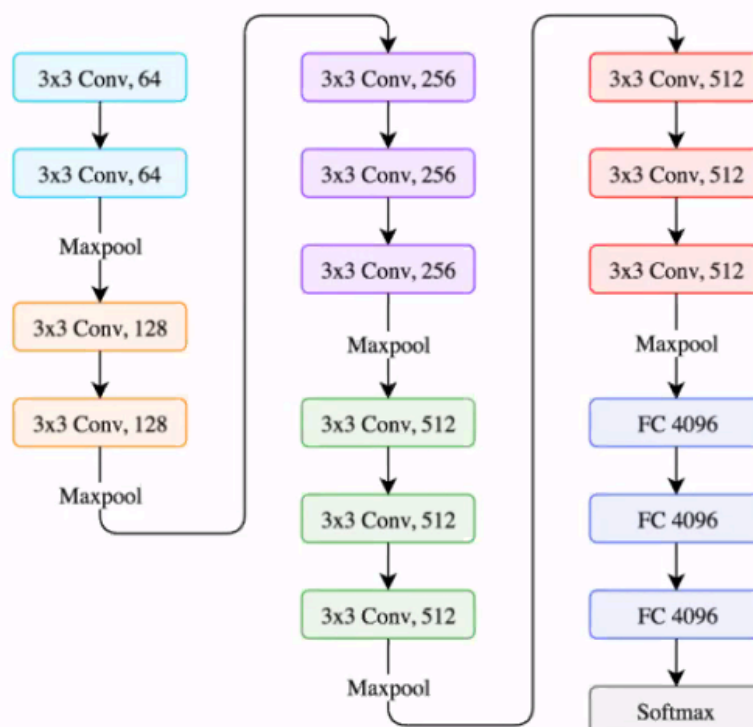


[Feature Vector / Class Label]



(Use for ML tasks: clustering, classification, similarity)

CNN: VGGNet : A Practical CNN Implementation



4

This image is a diagram of **VGGNet**, specifically the **VGG-16 architecture**, which is a widely used **Convolutional Neural Network (CNN)** model for image classification. Below is an explanation of each part of this architecture:

VGGNet: A Practical CNN Implementation

VGGNet is known for its **simplicity and depth**, using only:

- 3×3 convolution filters
 - Max pooling layers
 - Fully connected (FC) layers
 - Softmax for final classification
-

Structure Breakdown

Block 1:

- 3×3 Conv, 64
- 3×3 Conv, 64
- **Maxpool**

➡ Two convolution layers with 64 filters each, followed by a pooling layer to reduce spatial dimensions.

Block 2:

- 3×3 Conv, 128
- 3×3 Conv, 128
- **Maxpool**

➡ Increased number of filters (128), again two convolutions + pooling.


Block 3:

- 3×3 Conv, 256 × 3
- **Maxpool**

➡ Three convolution layers, now with 256 filters, adding more complexity and depth.


Block 4:

- `3×3 Conv, 512` × 3
- **Maxpool**

 Three layers with 512 filters each – feature maps are very rich at this point.


Block 5:

- `3×3 Conv, 512` × 3
- **Maxpool**

 Same as previous, final set of convolutional layers.

Fully Connected Layers

- `FC 4096`
- `FC 4096`
- `FC 4096`

 These layers act like a typical neural network (dense layers), processing the high-level features extracted by the convolutional part.

Softmax Layer

- Produces **class probabilities** for classification tasks (e.g., image categories).
-

General Flow:

Input Image →

Conv Layers (Feature Extraction) →

MaxPooling (Downsampling) →

Flatten →

Fully Connected Layers →

Softmax (Classification Output)

Key Highlights

- **All convolutions are 3×3:** Small receptive fields allow stacking more layers, which increases depth and expressiveness.
 - **Maxpooling** reduces spatial size and computations.
 - **Fully connected layers** integrate the abstracted features.
 - **Used in many transfer learning applications** (often pretrained on ImageNet).
-

What is Pooling?

Pooling is a **downsampling operation** that reduces the spatial dimensions (width and height) of feature maps while keeping the most important information.

Purpose:

- Reduce **computational complexity**
 - Control **overfitting**
 - Provide a form of **translation invariance** (i.e., the exact position of features matters less)
-

Types of Pooling

1. Max Pooling (most common)

- Takes the **maximum value** from each region (patch) of the feature map.

Example:

A **2×2 max pooling** with stride 2 over this matrix:

```
[1 3]
[2 4]
```

 Output: **4** (the max value in the 2×2 block)

Applied to a larger matrix:

```
Input:      Output (2×2 max pooling):
[[1, 3, 2, 1],  [[4, 3],
[4, 2, 1, 5],  →  [7, 8]]
```

[3, 6, 1, 2],
[7, 2, 8, 9]]

◆ 2. Average Pooling

- Takes the **average value** in each patch.
 - Less common, but sometimes used in place of max pooling in specific architectures.
-

🧠 How Pooling Works in CNNs

- Applied after convolution + activation (like ReLU)
 - Reduces the size of feature maps, e.g.:
 - From $32 \times 32 \rightarrow 16 \times 16$ with 2×2 pooling
 - Typically, **non-overlapping pooling** is used (stride = kernel size)
-

✅ Benefits of Pooling

Benefit	Explanation
Dimensionality reduction	Smaller feature maps = faster training & fewer parameters
Translation invariance	Detects features regardless of small shifts in position
Prevent overfitting	Less detailed spatial information reduces model complexity

🖼️ Visual Summary

Feature Map (before) → [Pooling] → Smaller, summarized Feature Map (after)
