

Comprehensive Guide to Convolutional Neural Networks

September 28, 2025

1 CNN Structure and Operating Principles

Convolutional neural networks (CNNs) exploit spatial locality and translation invariance by sharing weights across spatial positions. An input image tensor $\mathbf{X} \in \mathbb{R}^{C_{\text{in}} \times H \times W}$ is convolved with kernels $\mathbf{K} \in \mathbb{R}^{C_{\text{out}} \times C_{\text{in}} \times k_h \times k_w}$ to produce feature maps:

$$Y_{c,i,j} = \sum_{c'=1}^{C_{\text{in}}} \sum_{u=0}^{k_h-1} \sum_{v=0}^{k_w-1} K_{c,c',u,v} X_{c',is+u-p,js+v-p}, \quad (1)$$

where s is the stride and p the padding. Stride controls downsampling, while padding preserves boundary information. Figure ?? illustrates how receptive fields grow across layers.

1.1 Feature Hierarchies

Early convolutional layers learn edge and texture detectors; middle layers capture motifs such as corners or object parts; deeper layers integrate global semantics. Weight sharing dramatically reduces parameter counts compared with fully connected layers.

1.2 Padding, Stride, and Dilation

Padding extends the input with zeros (or reflection/replication) so that convolutions can cover border pixels. Dilation introduces gaps between kernel elements, expanding the receptive field to $RF = (k-1)d + 1$ without increasing parameters. Careful combinations of stride and dilation maintain resolution while encoding global context.

1.3 Pooling and Downsampling

Pooling layers summarize local neighborhoods, providing robustness to small translations. Max pooling emphasizes the strongest activation, average pooling smooths representations, and strided convolutions can replace pooling entirely for learnable downsampling. Global average pooling reduces each feature map to a scalar, enabling fully convolutional classifiers.

1.4 Normalization and Activation

CNNs typically use rectified linear unit (ReLU) or its variants (LeakyReLU, GELU) to introduce nonlinearity. Batch normalization normalizes activations using mini-batch statistics, mitigating internal covariate shift and enabling higher learning rates. Layer normalization or group normalization are alternative choices for small batches.

Listing 1: PyTorch implementation of a convolutional stem with normalization and pooling.

```
1 import torch.nn as nn
2
3 class ConvStem(nn.Module):
4     def __init__(self, in_channels=3, hidden_channels=(64, 128, 256))
5     :
6         super().__init__()
7         layers = []
8         current = in_channels
9         for out_channels in hidden_channels:
10             layers += [
11                 nn.Conv2d(current, out_channels, kernel_size=3,
12                           stride=1, padding=1, bias=False),
13                 nn.BatchNorm2d(out_channels),
14                 nn.ReLU(inplace=True),
15                 nn.MaxPool2d(kernel_size=2, stride=2)
16             ]
17             current = out_channels
18         self.stem = nn.Sequential(*layers)
19
20     def forward(self, x):
21         return self.stem(x)
```

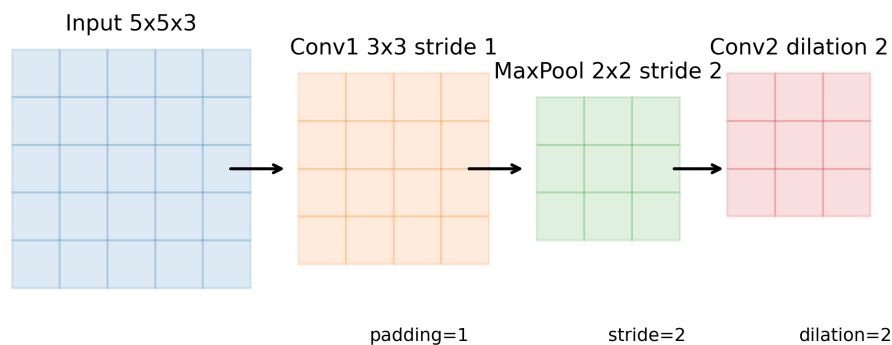


Figure 1: Illustration of convolutions, padding, stride, and pooling expanding receptive fields across layers.

2 Classical Architectures: A Comparative Overview

Over three decades, CNN architectures evolved to address optimization challenges, computational efficiency, and scalability. Figure ?? highlights key milestones, and Table ?? contrasts design choices.

2.1 LeNet-5 (1998)

Yann LeCun’s LeNet pioneered digit recognition by stacking convolution and subsampling layers before fully connected classifiers. Its use of weight sharing enabled practical training on limited hardware.

2.2 AlexNet (2012)

AlexNet conquered ImageNet using deeper networks, ReLU activations to avoid saturation, extensive data augmentation, and dropout regularization. Two GPUs processed disjoint halves of feature maps, showcasing the importance of hardware acceleration.

2.3 VGG (2014)

VGG-16/19 advocated uniform 3×3 convolutions stacked deep, demonstrating that depth and small kernels can approximate larger receptive fields with fewer parameters. The model’s simplicity made it a go-to backbone for transfer learning.

2.4 ResNet (2015)

Residual connections allowed ResNet to train over 150 layers by reformulating layers as residual functions $\mathcal{F}(\mathbf{x}) + \mathbf{x}$. Identity shortcuts preserve gradient flow and enable easier optimization.

2.5 Inception Family (2014–2016)

GoogLeNet (Inception-v1) introduced multi-branch modules combining various kernel sizes. Later versions factorized 5×5 into stacked 3×3 convolutions and added batch normalization. Inception-v3/v4 further refined factorization and regularization.

2.6 EfficientNet (2019)

EfficientNet scales depth, width, and resolution via a compound coefficient, starting from a mobile inverted bottleneck block with squeeze-and-excitation attention. Balanced scaling achieved state-of-the-art accuracy with fewer FLOPs.

Table 1: Comparison of influential CNN architectures. Parameter counts and FLOPs are approximate for ImageNet-scale models.

Architecture	Year	Depth	Params (M)	FLOPs (G)	Signature features
LeNet-5	1998	7	0.06	0.002	Conv + subsampling for digits
AlexNet	2012	8	61	1.5	ReLU, dropout, large kernels
VGG-16	2014	16	138	15.5	Stacked 3×3 convs
ResNet-50	2015	50	25	4.1	Residual bottlenecks
Inception-v3	2016	48	24	5.7	Multi-branch factorized convs
EfficientNet-B4	2019	82	19	4.2	Compound scaling, MBConv

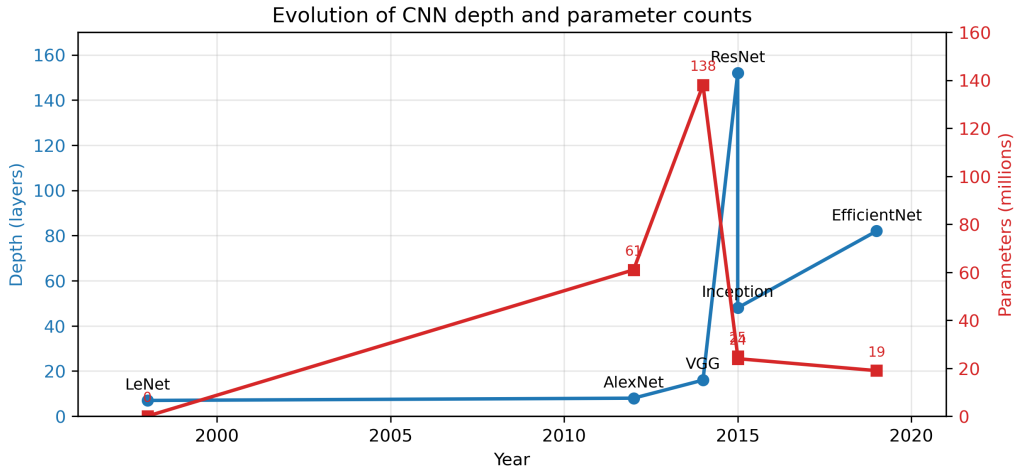


Figure 2: Timeline of CNN breakthroughs and their core innovations.

3 Applications Across Vision Tasks

CNNs underpin diverse computer vision applications. Figure ?? sketches high-level pipelines.

3.1 Image Classification

Classification networks map an input image to class probabilities. Modern classifiers leverage pretrained backbones, global pooling, and label smoothing to combat overconfidence. Mixup and CutMix augmentation further regularize training.

3.2 Object Detection

Detection augments CNN backbones with localization heads. Two-stage detectors (Faster R-CNN) generate region proposals before refinement, while one-stage models (YOLOv8, RetinaNet) predict boxes and classes densely. Feature pyramid networks (FPN) reuse multi-scale feature maps to detect objects spanning wide size ranges.

3.3 Semantic Segmentation

Segmentation assigns labels to individual pixels. Encoder-decoder architectures (U-Net, SegNet) recover spatial detail via skip connections. Dilated convolutions (DeepLab),

pyramid pooling modules (PSPNet), and attention enhance context reasoning. Metrics such as mean Intersection-over-Union (mIoU) measure performance.

3.4 Beyond 2D Vision

CNNs extend to video (3D convolutions, TimeSformer hybrids), medical imaging (multi-scale U-Nets), and multimodal perception (CNN backbones fused with transformers). Lightweight variants (MobileNet, ShuffleNet) enable deployment on edge devices.

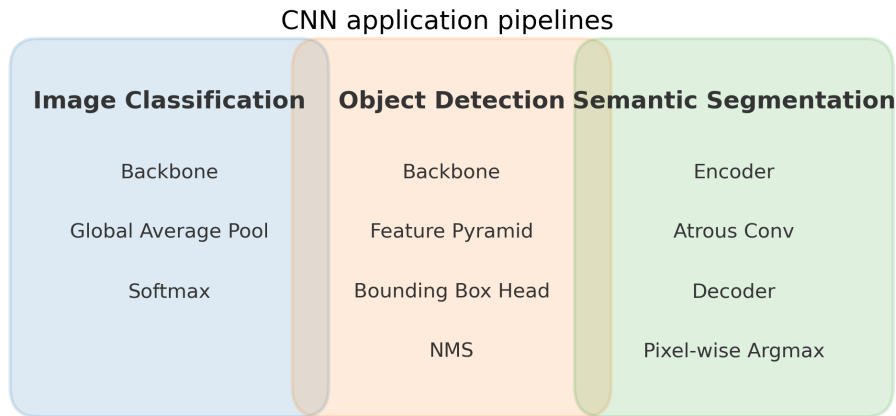


Figure 3: Representative CNN application pipelines: classification, detection, and semantic segmentation.

4 Training and Deployment Tips

- **Initialization**
normalization: Use He initialization with BatchNorm or GroupNorm for stable gradients.
- **Regularization:** Apply augmentation (random crop, color jitter, CutMix), dropout, and stochastic depth to combat overfitting.
- **Optimization:** Cosine learning-rate schedules with warmup and momentum SGD remain strong baselines; adaptive optimizers aid rapid prototyping.
- **Efficiency:** Prune filters, quantize weights, or use depthwise separable convolutions for mobile deployment.
- **Monitoring:** Track task-specific metrics (Top-1 accuracy, mAP, mIoU) and visualize feature maps to diagnose failures.