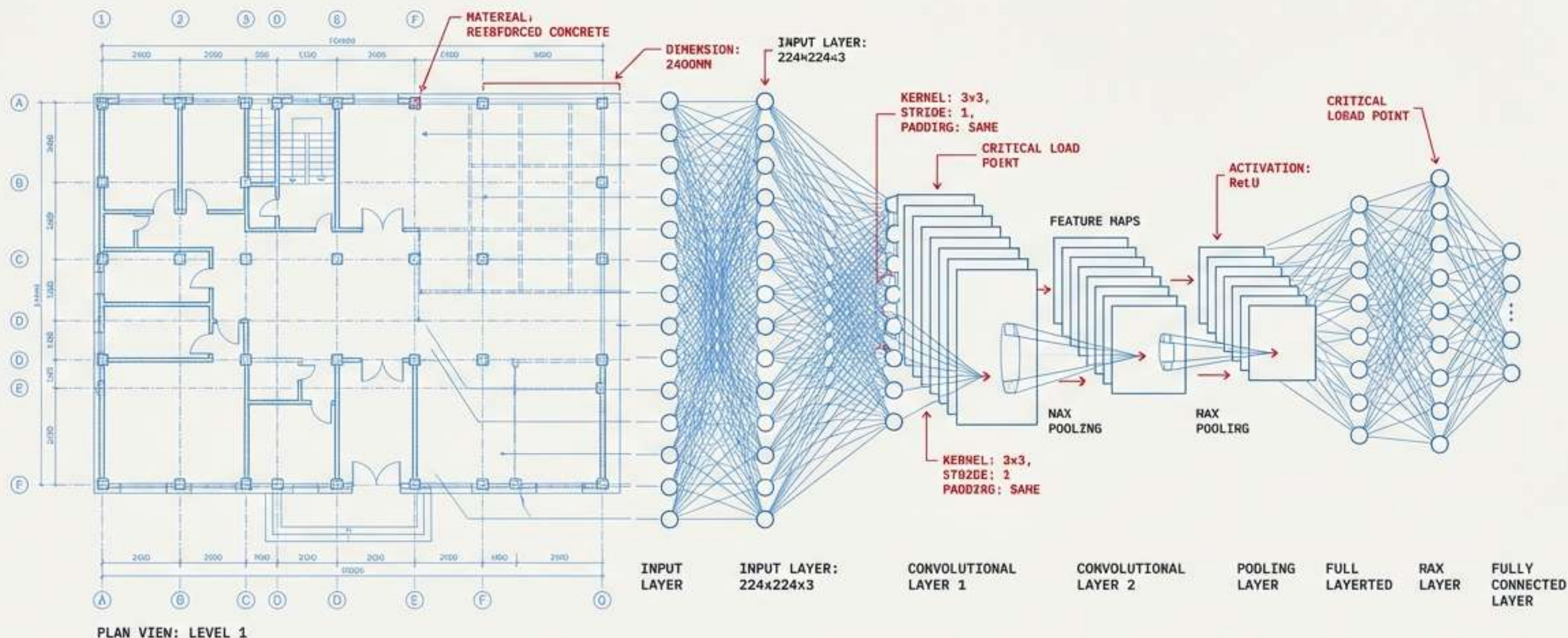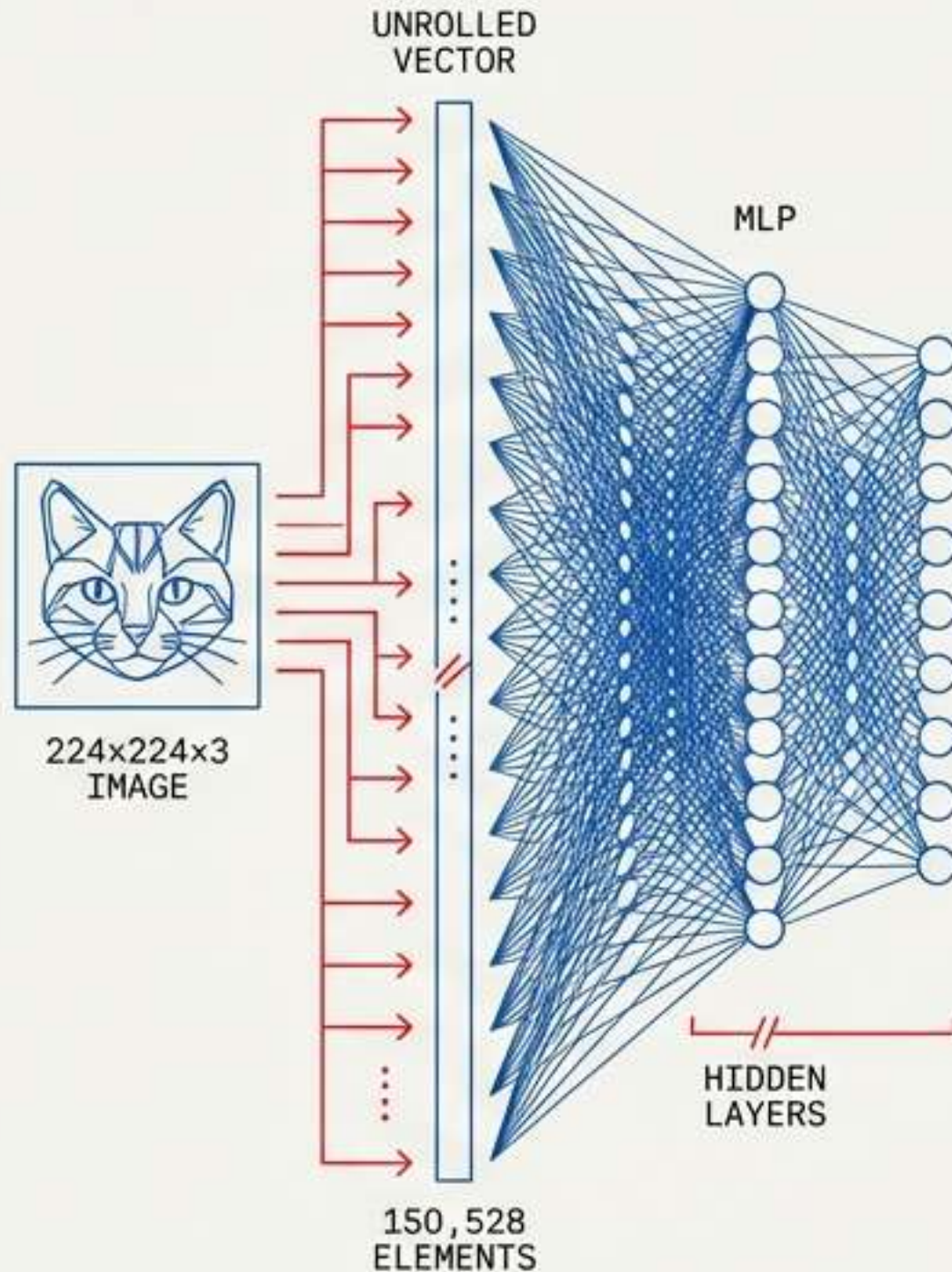# The Architect's Blueprint for Convolutional Neural Networks

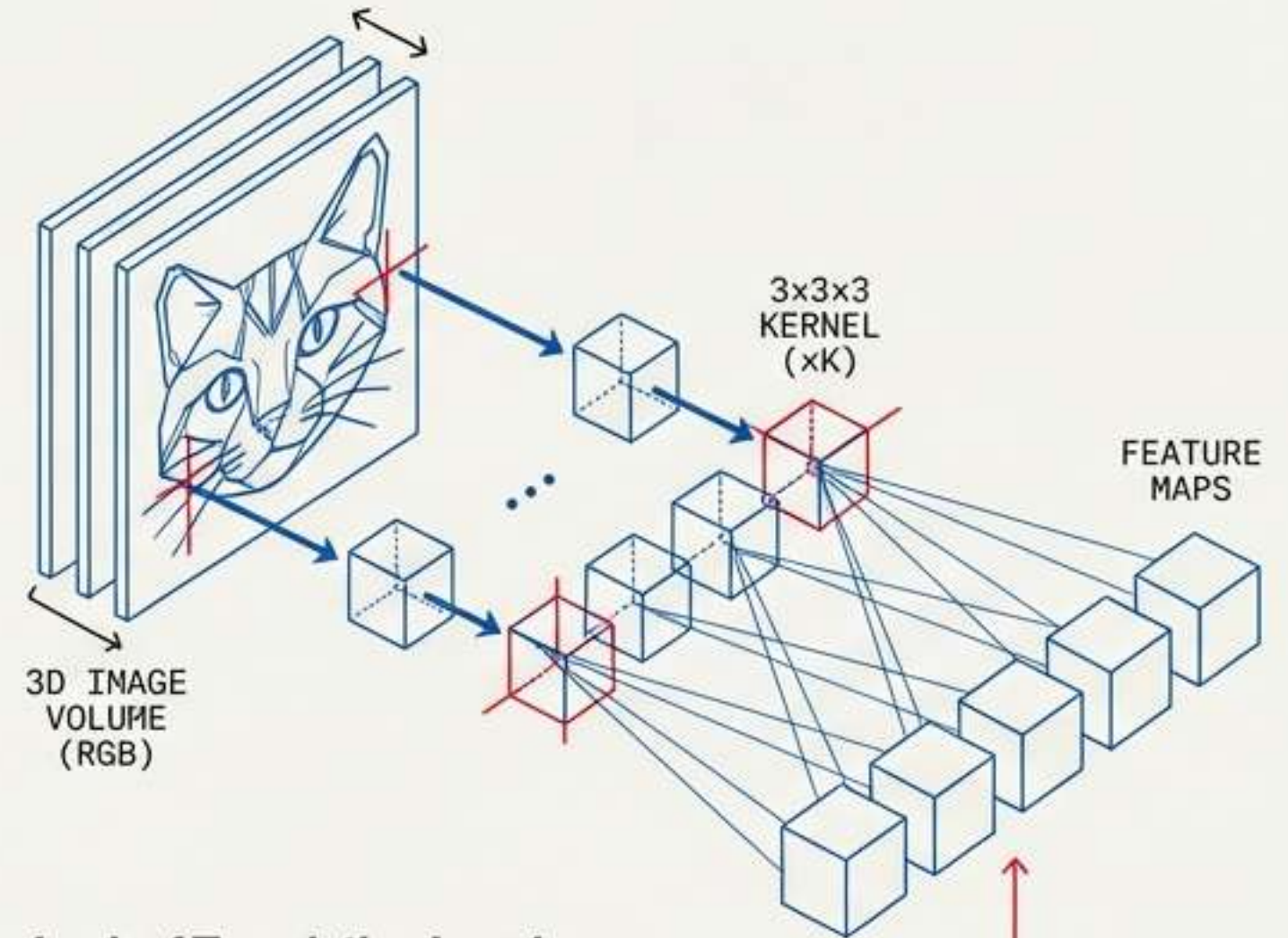A Foundational Guide to Building and Understanding Image Recognition Models

# Why Traditional Blueprints Fail for Image Data

Key idea: Standard Multilayer Perceptrons (MLPs) are ill-suited for images due to two critical flaws: parameter explosion and a lack of translation invariance.



UNROLLED VECTOR

MLP

224x224x3 IMAGE

150,528 ELEMENTS

HIDDEN LAYERS

## Parameter Explosion & Overfitting

An MLP requires a unique neuron for every single pixel. For a standard 224×224×3 color image, this means an input layer with 150,528 neurons. r neurons. A modest network with a few hidden layers would exceed 300 Billion trainable parameters, making it incredibly slow to train and highly prone to overfitting.

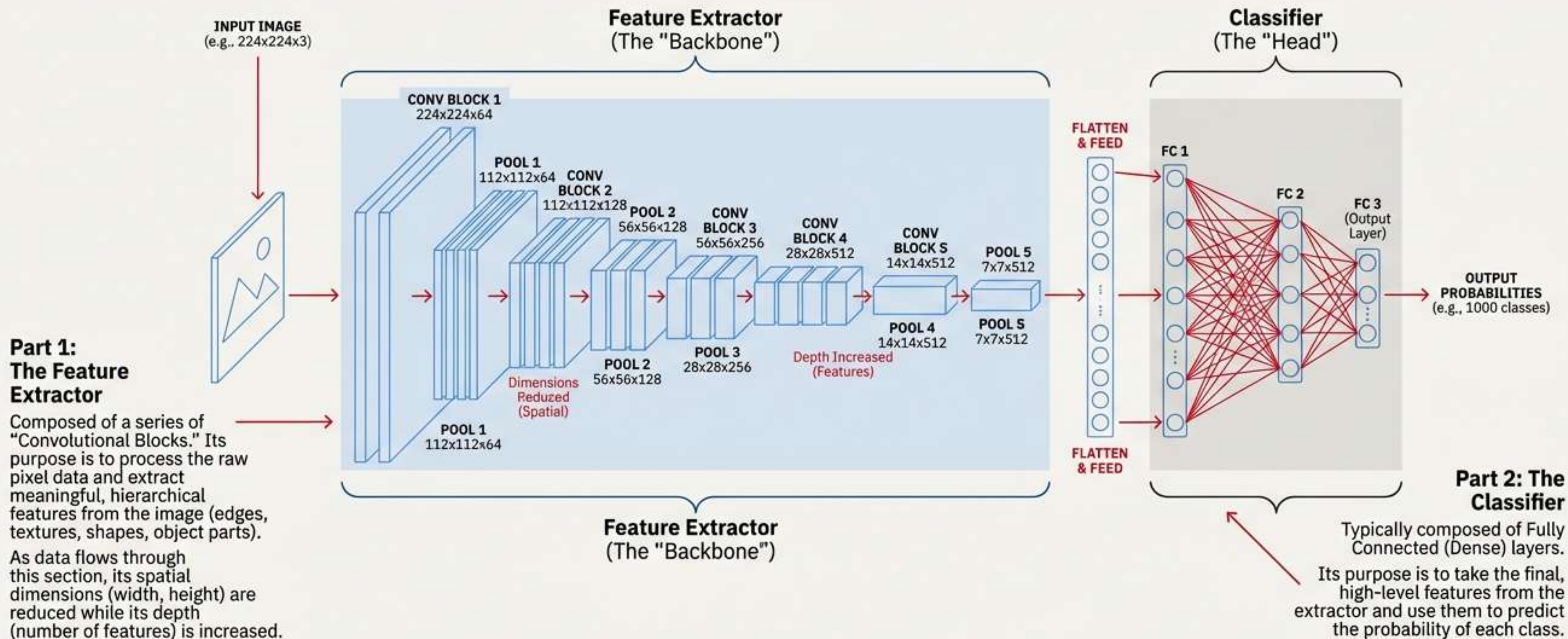3D IMAGE VOLUME (RGB)

3x3x3 KERNEL (×K)

FEATURE MAPS

## Lack of Translation Invariance

MLPs learn features based on absolute pixel position. The network reacts differently if the subject of the image is shifted. This forces the model to re-learn the same object at every possible location, making it inefficient and unreliable.
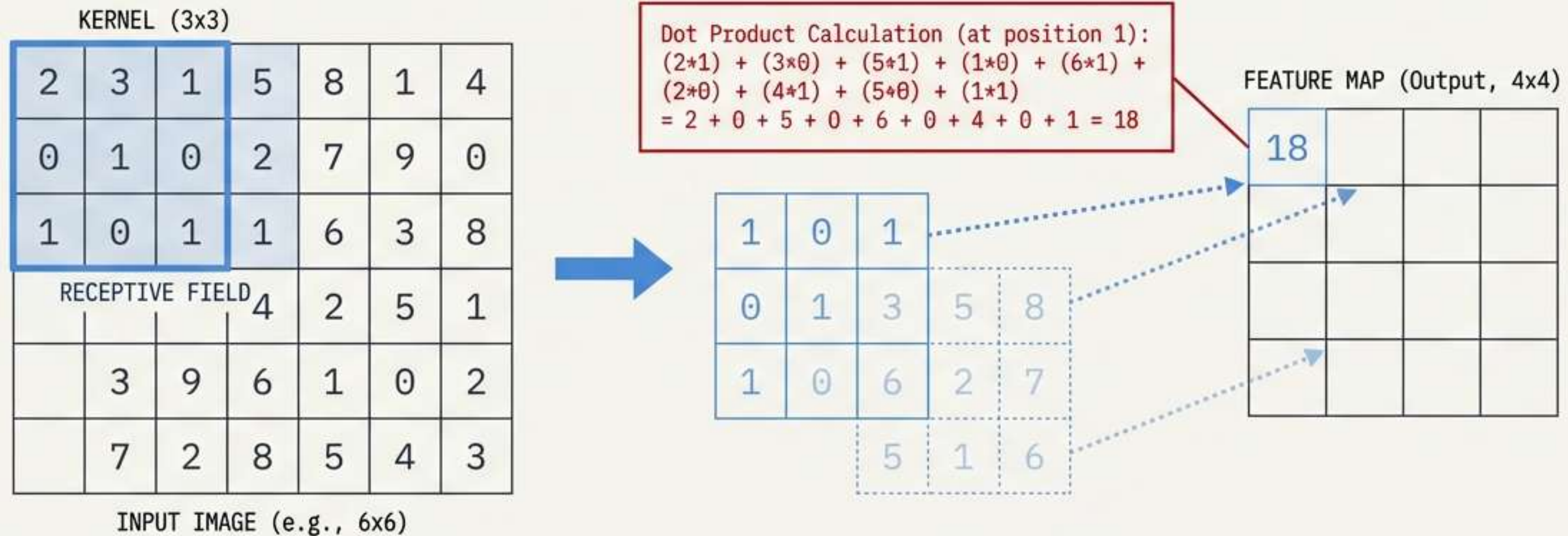
NotebookLM

# The CNN Blueprint: A Two-Part Structure

A CNN is composed of two distinct segments: an upstream **Feature Extractor** that learns to see, and a downstream **Classifier** that learns to decide.



**INPUT IMAGE** (e.g.. 224x224x3)

**Feature Extractor** (The "Backbone")

**Classifier** (The "Head")

CONV BLOCK 1 224x224x64

POOL 1 112x112x64

CONV BLOCK 2 112x112x128

POOL 2 56x56x128

CONV BLOCK 3 56x56x256

CONV BLOCK 4 28x28x512

CONV BLOCK S 14x14x512

POOL 5 7x7x512

POOL 2 56x56x128

POOL 3 28x28x256

POOL 4 14x14x512

POOL S 7x7x512

Dimensions Reduzed (Spatial)

Depth Increased (Features)

POOL 1 112x112x64

FLATTEN & FEED

FC 1

FC 2

FC 3 (Output Layer)

OUTPUT PROBABILITIES (e.g., 1000 classes)

FLATTEN & FEED

Feature Extractor (The "Backbone")

**Part 1: The Feature Extractor**

Composed of a series of "Convolutional Blocks." Its purpose is to process the raw pixel data and extract meaningful, hierarchical features from the image (edges, textures, shapes, object parts).

As data flows through this section, its spatial dimensions (width, height) are reduced while its depth (number of features) is increased.

**Part 2: The Classifier**

Typically composed of Fully Connected (Dense) layers.

Its purpose is to take the final, high-level features from the extractor and use them to predict the probability of each class.

NotebookLM

# Core Material I: Convolution for Extracting Local Patterns

The convolution operation is the heart of a CNN. It involves sliding a small, learnable matrix (a kernel or filter) over the input to detect specific local features like edges, textures, or shapes.

**KERNEL (3x3)**

| 2 | 3 | 1 | 5 | 8 | 1 | 4 |
|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 2 | 7 | 9 | 0 |
| 1 | 0 | 1 | 1 | 6 | 3 | 8 |

**RECEPTIVE FIELD**

| | 4 | 2 | 5 | 1 |
| 3 | 9 | 6 | 1 | 0 | 2 |
| 7 | 2 | 8 | 5 | 4 | 3 |

**INPUT IMAGE (e.g., 6x6)**

Dot Product Calculation (at position 1):
$(2*1) + (3*0) + (5*1) + (1*0) + (6*1) + (2*0) + (4*1) + (5*0) + (1*1)$
$= 2 + 0 + 5 + 0 + 6 + 0 + 4 + 0 + 1 = 18$

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

**FEATURE MAP (Output, 4x4)**

| 18 | | | |
|----|--|--|--|
| | | | |
| | | | |
| | | | |

## THE OPERATION

- A small filter (e.g., 3x3) slides over the input image.
- At each position, a dot product is computed between the filter's values and the corresponding input pixels. This region of the input is called the RECEPTIVE FIELD.
- This process produces a new grid called a FEATURE MAP or ACTIVATION MAP, which highlights where the specific feature was detected.

## KEY BENEFIT: PARAMETER SHARING

The *same* filter (with the same set of weights) is used across the entire image. This dramatically reduces the number of parameters compared to an MLP and allows the network to detect a feature regardless of its position (translation invariance).
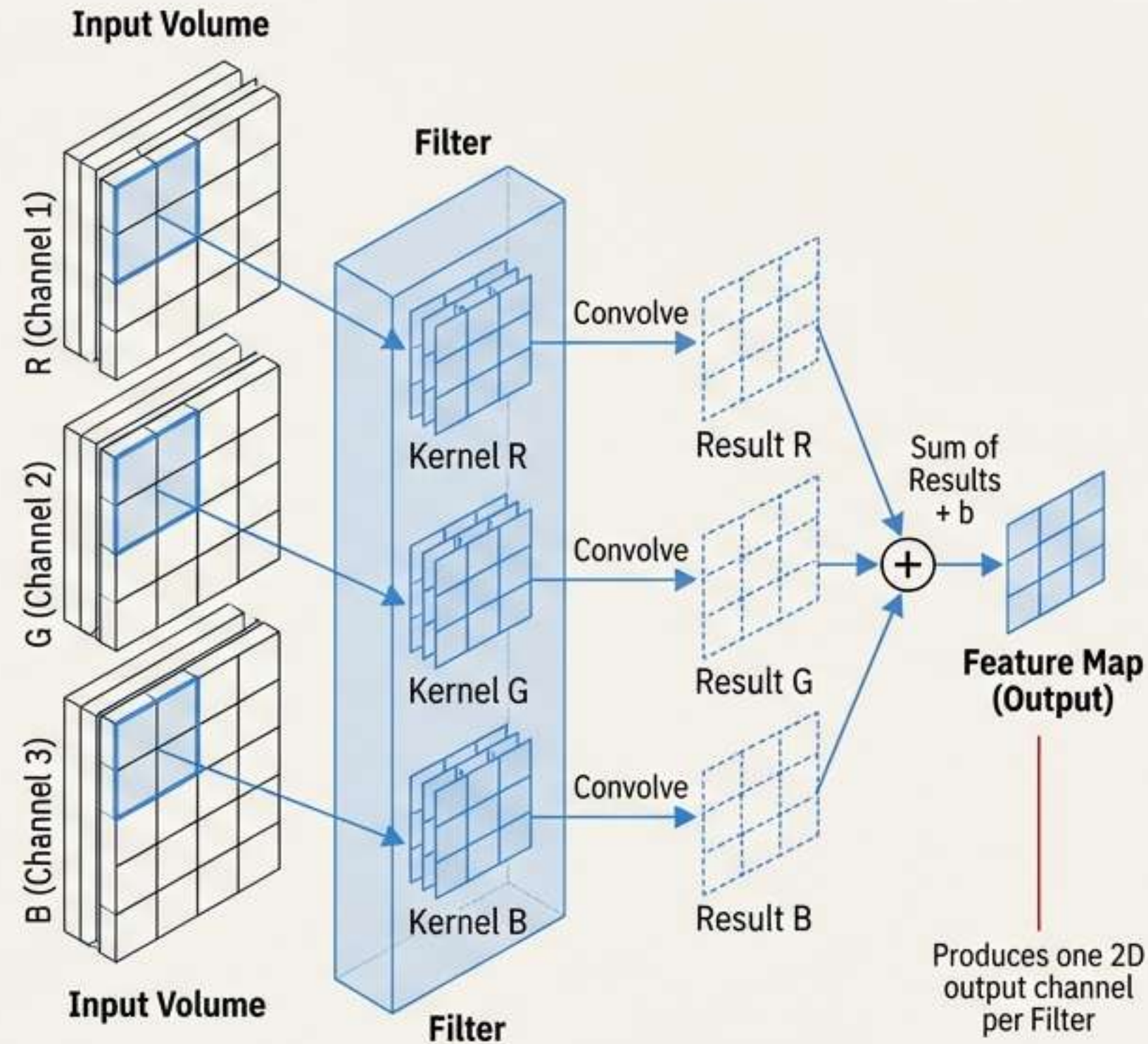
# Anatomy of a Convolution: Filters, Kernels, and Feature Detection

## Terminology

- **Kernel**: A single 2D matrix of learnable weights designed to detect one specific feature.

- **Filter**: A collection of kernels. For a multi-channel input (e.g., 3 RGB channels), a filter will have a corresponding kernel for each input channel. The number of filters in a layer is a design choice and determines the depth (number of channels) of the output.

## In CNNs

The network doesn't use pre-defined kernels. It *learns* the optimals use the optimal values for the kernels during training to detect features that are most useful for the classification task.

**Input Volume**

R (Channel 1)

G (Channel 2)

B (Channel 3)

**Input Volume**

**Filter**

Kernel R

Kernel G

Kernel B

**Filter**

Convolve → Result R

Convolve → Result G

Convolve → Result B

Sum of Results + b

$\oplus$

**Feature Map (Output)**

Produces one 2D output channel per Filter

## Example: The Sobel Kernel

A hand-crafted kernel used in traditional image processing to detect vertical edges. The values (+1s on the left, -1s on the right) numerically approximate the horizontal derivative, producing a high response where there are sharp vertical changes in intensity.

| +1 | 0 | -1 |
|----|---|----|
| +2 | 0 | -2 |
| +1 | 0 | -1 |

Positive/Negative weights create contrast

**Sobel Kernel (Vertical Edges)**

High response on vertical edges

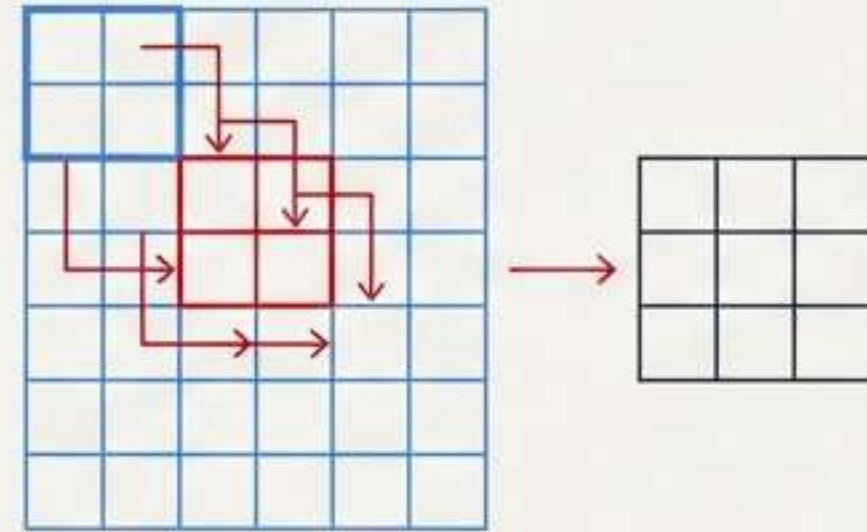**Input Image** → **Resulting Feature Map**

# Controlling the Blueprint: Stride and Padding

**Purpose:** These hyperparameters determine the spatial size of the output feature map, allowing us to preserve or reduce dimensions strategically.
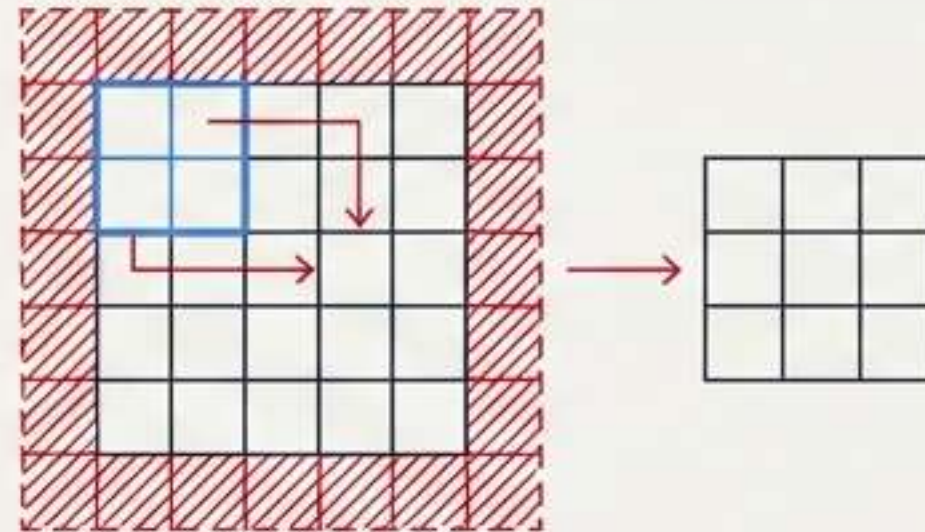
## Stride

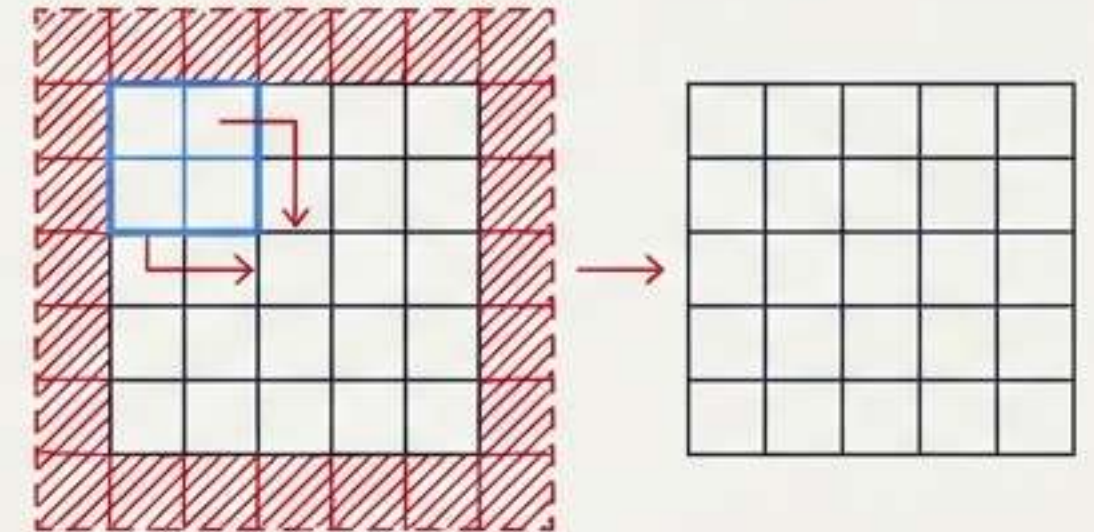**Definition:** The number of pixels the kernel moves at each step.

**Effect:** A stride of 1 moves one pixel at a time, resulting in more detailed feature maps and larger outputs. A larger stride results in a smaller output size and less overlap between receptive fields.

## Padding

**Definition:** Adding extra pixels (usually zeros) around the input image border.

**Effect:** "Valid" padding (no padding) shrinks the output size. "Same" padding adds just enough pixels so the output spatial dimensions equal the input spatial dimensions. This preserves border information.

"Valid" Padding, Stride = 1     With Padding, Stride = 2     "Same" Padding, Stride = 1

$$O = (n - f + 2p) / s + 1$$

$O$: Output size       $p$: Padding
$n$: Input size        $s$: Stride
$f$: Kernel size

# Core Material II: Pooling to Reduce Complexity and Provide Invariance

**Purpose:** Pooling layers progressively reduce the spatial size of the feature maps to decrease the number of parameters and computation in the network. This also helps make the learned feature representations more robust to small translations in the input.
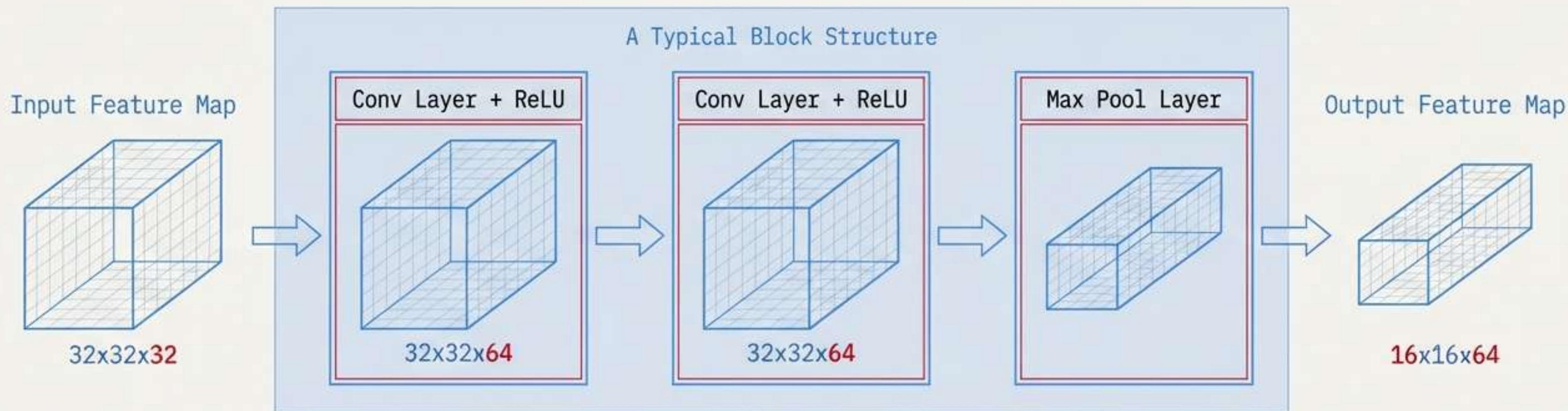


Input

Output
Pooled Feature Map

## Operation: Max Pooling

The most common type of pooling. A small window (e.g., 2x2) slides over the feature map. At each position, it outputs only the *maximum* value within the window. Crucially, the pooling filter has **no trainable parameters**. It's a fixed operation.

## Effect

It summarizes the features in a region, keeping the most prominent one. This downsampling reduces computational load and helps mitigate overfitting.

# The Assembly: Stacking Layers into Convolutional Blocks

**Key idea:** In practice, layers are not used in isolation. They are grouped into repeating patterns called "**Convolutional Blocks**" that form the backbone of the feature extractor.

A Typical Block Structure

| Input Feature Map | Conv Layer + ReLU | Conv Layer + ReLU | Max Pool Layer | Output Feature Map |
|---|---|---|---|---|
| 32x32x32 | 32x32x64 | 32x32x64 | | 16x16x64 |

A typical block consists of one or more convolutional layers followed by a pooling layer.

**Convolutional Layers:** Apply multiple filters to extract a rich set of features from the input. These layers increase the *depth* of the data (e.g., from 32 channels to 64 channels).

**Activation Function (ReLU):** After each convolution, a non-linear activation function is applied to enable the network to learn complex patterns.

**Max Pooling Layer:** Downsamples the resulting feature maps, reducing their spatial dimensions (e.g., from 32x32 to 16x16) before passing them to the next block.

# The Power of Depth: Learning Hierarchical Features

**Key Insight:** CNNs learn features in a hierarchical manner. Deeper layers combine features from earlier layers to detect increasingly complex structures.

Early Layers                Mid Layers                Deep Layers



## The Hierarchy

- **Layer 1:** Learns basic elements like edges, color blobs, and gradients.
- → **Intermediate Layers:** Combine edges to form textures, corners, and simple shapes.
- → **Deep Layers:** Assemble shapes into complex object parts (an eye, a wheel).

## Growing Receptive Fields

The **receptive field** is the region of the original input image that affects a single neuron's output. As we go deeper, each neuron's receptive field grows, allowing it to "see" and understand more global context from the original image.

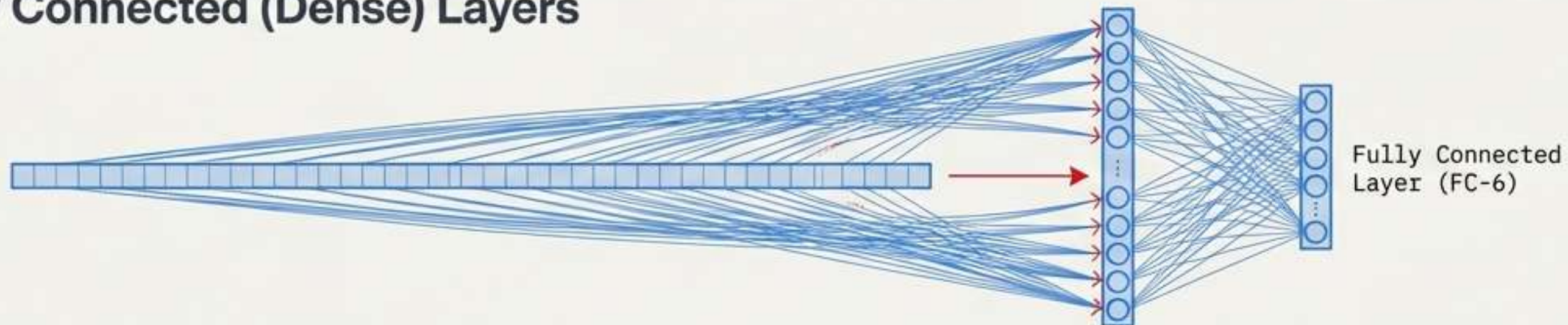# The Classifier Head: From Feature Maps to Predictions

The classifier's job is to take the high-level, spatially rich feature maps from the final convolutional block and transform them into class probabilities.

## Step 1: Flattening



Final Feature Map
7x7x512

Flattened Vector
25,088 values

The 3D output of the feature extractor (e.g., 7×7×512) must be converted into a 1D vector. This "flattening" operation simply unrolls the multi-dimensional data into a long list of numbers (e.g., 25,088 values).

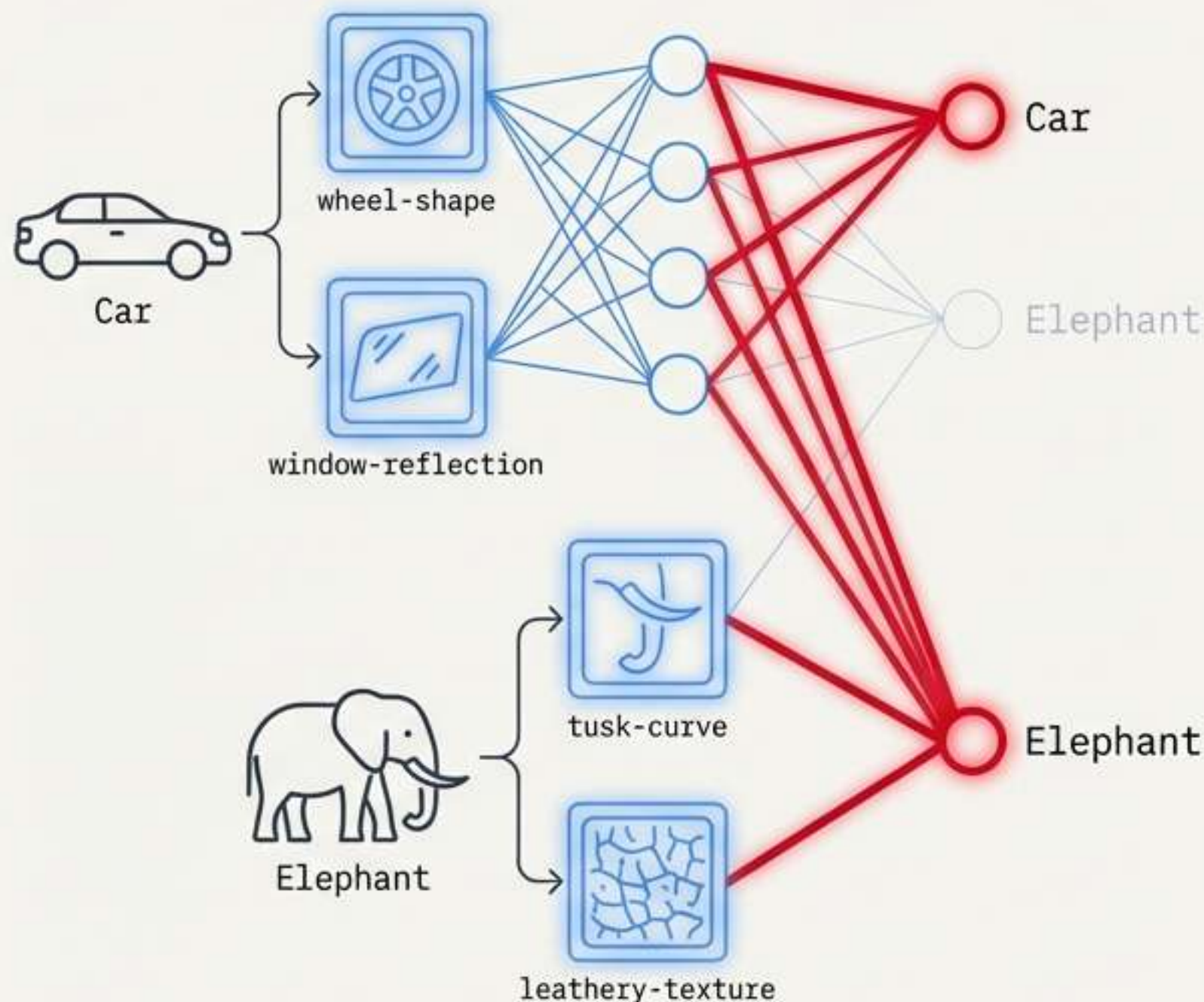## Step 2: Fully Connected (Dense) Layers



Fully Connected
Layer (FC-6)

This 1D vector is then fed into one or more standard fully connected layers. In a fully connected layer, every neuron is connected to *all* neurons in the previous layer. This allows the network to combine features from all spatial locations of the final feature map to make a holistic decision.

# How the Classifier Makes a Decision

## Intuition

The final feature maps from a trained CNN contain meaningful information about the image content. By connecting these features in alts in a fully connected manner, the classifier learns which combinations of features are indicative of which class.



Car

wheel-shape

window-reflection

Car

Elephant

tusk-curve

leathery-texture

Elephant

Car

Elephant

## Example

For an image of a **car**, feature maps corresponding to 'wheels,' 'windows,' and 'headlights' will have high activations. The trained weights in the fully connected layers will create strong pathways from these specific feature activations to the output neuron for the 'car' class.

## The Final Step: Softmax

The final fully connected layer outputs raw scores (logits) for each class. A **Softmax** function is applied to these scores to convert them into a probability distribution. It assigns decimal probabilities to each class, and all probabilities sum to 1. The class with the highest probability is the model's final prediction.

# Finishing Touches I: Building a Robust Network

**Key Idea:** Training deep networks is challenging. Regularization techniques are essential to prevent overfitting and improve generalization to unseen data.
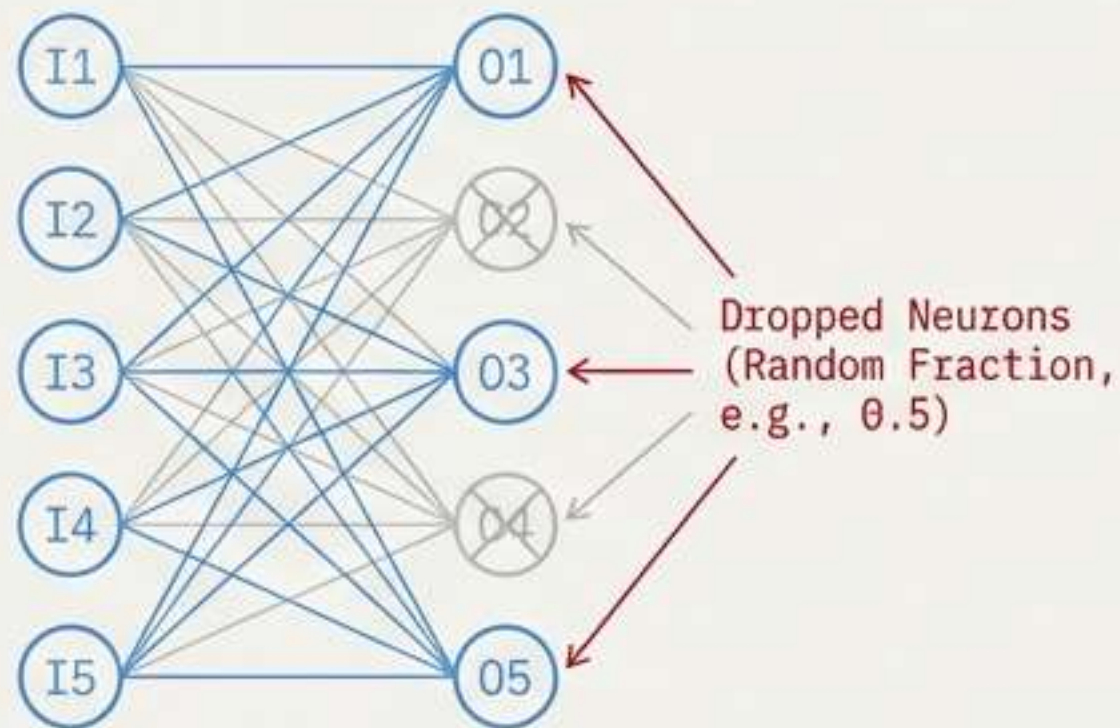
## Technique 1: Dropout

**Purpose**
To prevent overfitting and improve generalization.

**How it Works**
During training, randomly sets a fraction of neuron activations to zero at each update step. This forces the network to learn redundant representations and prevents it from becoming too reliant on any single neuron.
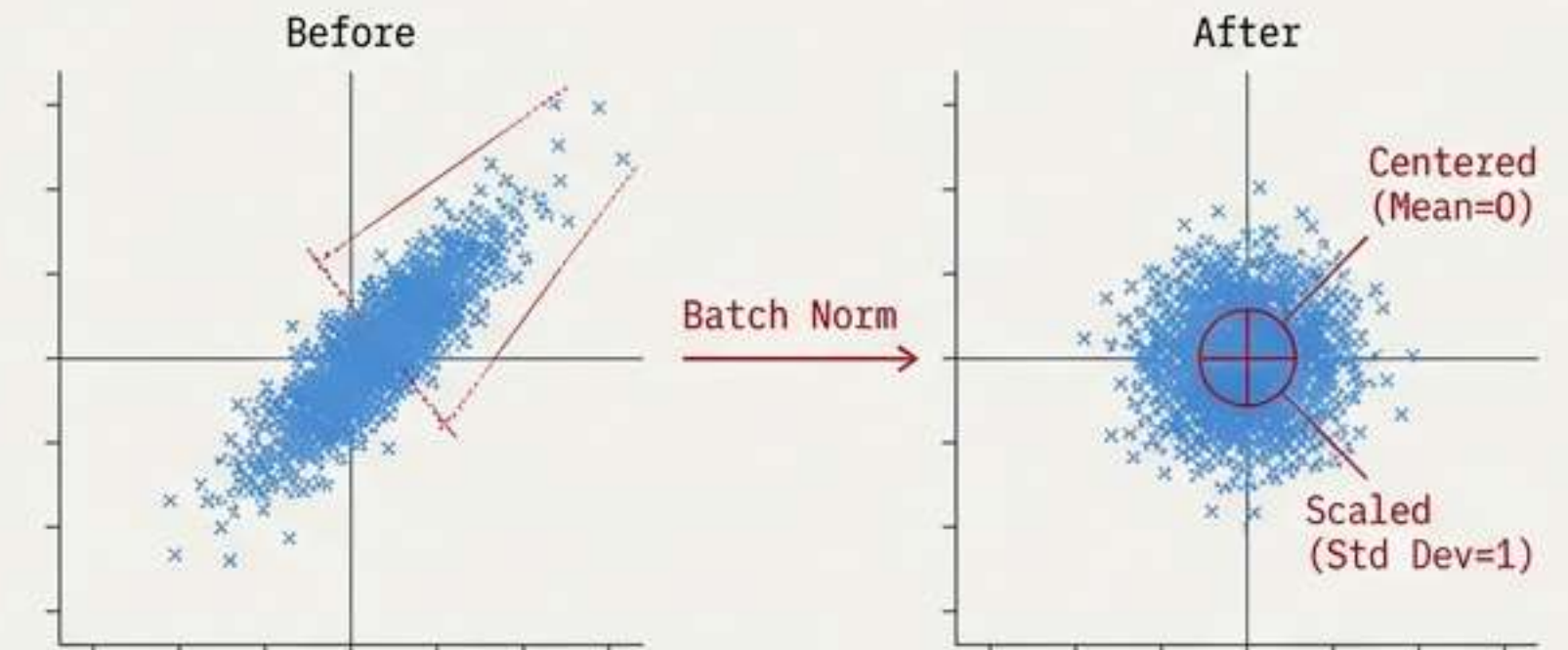
## Technique 2: Batch Normalization

**Purpose**
To stabilize training, accelerate convergence, and reduce sensitivity to weight initialization.
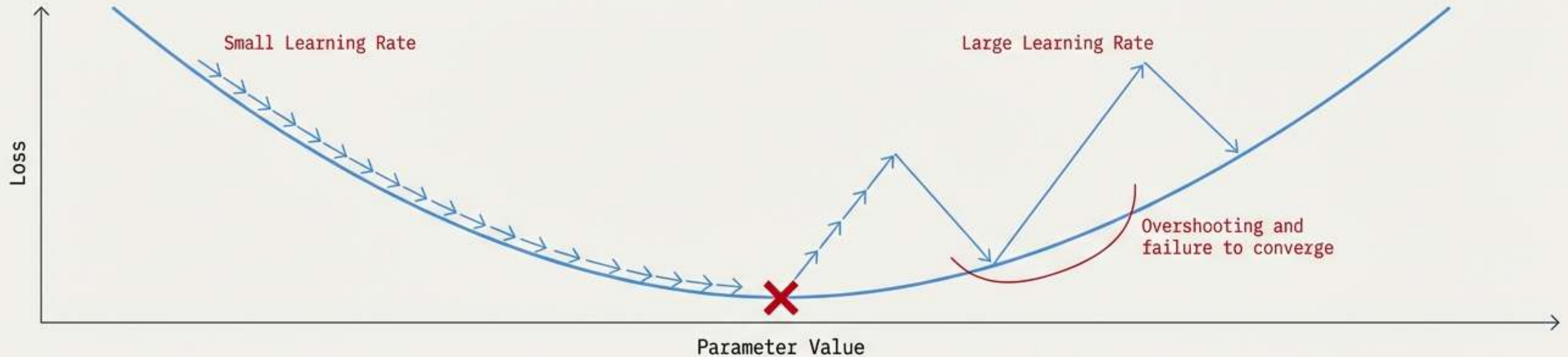
**How it Works**
Adds a normalization "layer" that standardizes the outputs of a previous layer by re-centering and re-scaling them across the current batch of data. This mitigates "internal covariate shift," where the distribution of layer inputs changes during training.

**Before and After**

# Finishing Touches II: Optimizing the Learning Process

**Key Idea:** The network learns by minimizing a Loss Function (e.g., cross-entropy loss) via an optimization algorithm that updates the weights.



## Hyperparameter 1: Learning Rate (LR)

Determines the size of the steps the optimizer takes to reach a local minimum of the loss function.

**The Trade-off**

A small LR leads to slow but stable convergence. A large LR can speed things up but risks overshooting the minimum and failing to converge.

## Optimizers: Adam vs. SGD

**SGD**
**Stochastic Gradient Descent (SGD):** The classic optimizer. Updates parameters for each training example or a small batch, following the direction of the slope downhill.

**Adam**
**Adam (Adaptive Moment Estimation):** A more advanced optimizer that computes adaptive learning rates for each parameter. It often converges faster and is more memory-efficient than standard SGD.
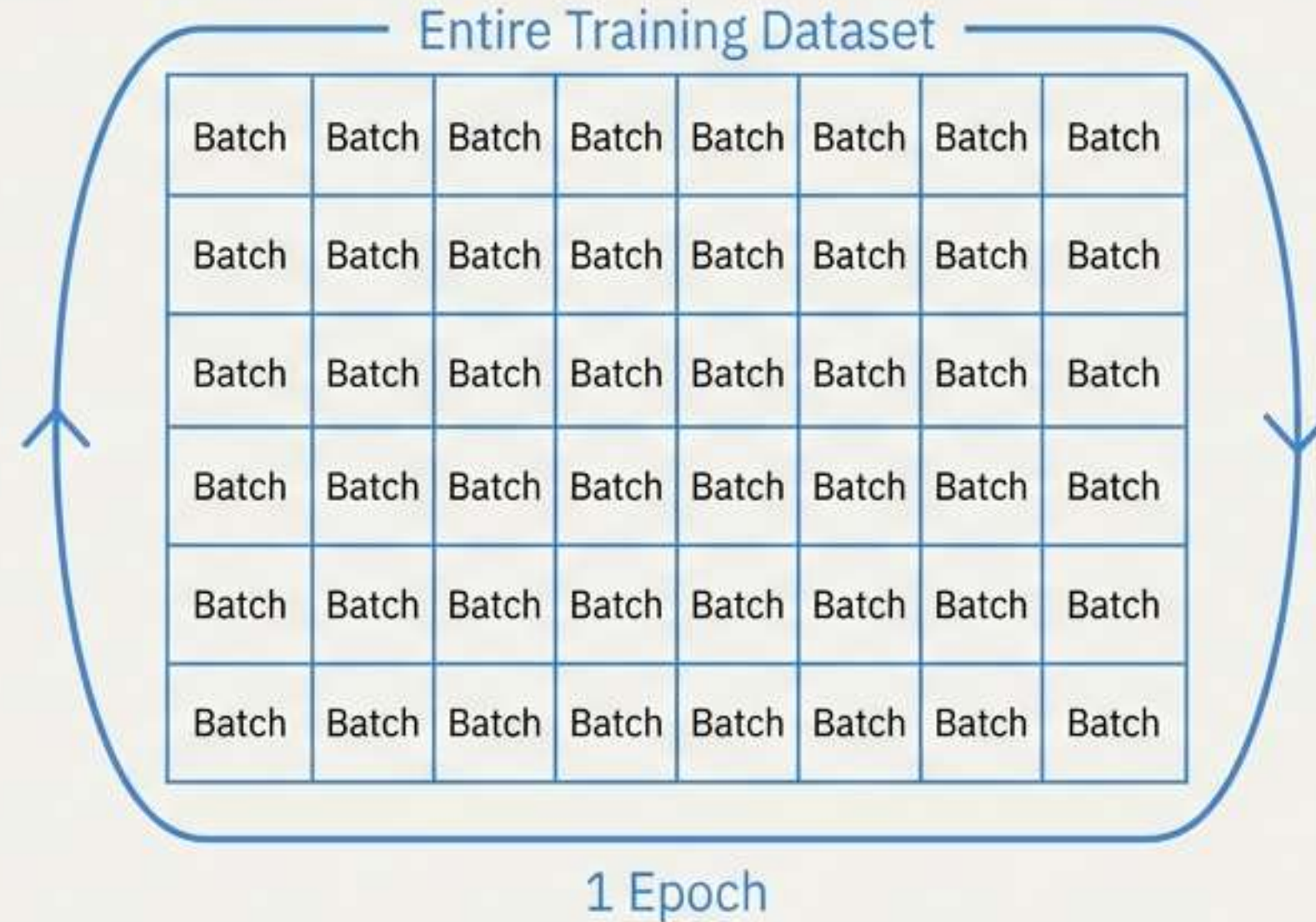
# Practical Considerations: Epochs and Batch Size

## Term 1: Epochs

One epoch is one full pass of the training algorithm over the *entire* training dataset.

## Usage

The number of epochs is the number of times the network will see the full dataset during training.

Entire Training Dataset

| Batch | Batch | Batch | Batch | Batch | Batch | Batch | Batch |
|-------|-------|-------|-------|-------|-------|-------|-------|
| Batch | Batch | Batch | Batch | Batch | Batch | Batch | Batch |
| Batch | Batch | Batch | Batch | Batch | Batch | Batch | Batch |
| Batch | Batch | Batch | Batch | Batch | Batch | Batch | Batch |
| Batch | Batch | Batch | Batch | Batch | Batch | Batch | Batch |
| Batch | Batch | Batch | Batch | Batch | Batch | Batch | Batch |

1 Epoch

## Term 2: Batch Size

The number of training examples utilized in one iteration (one weight update).
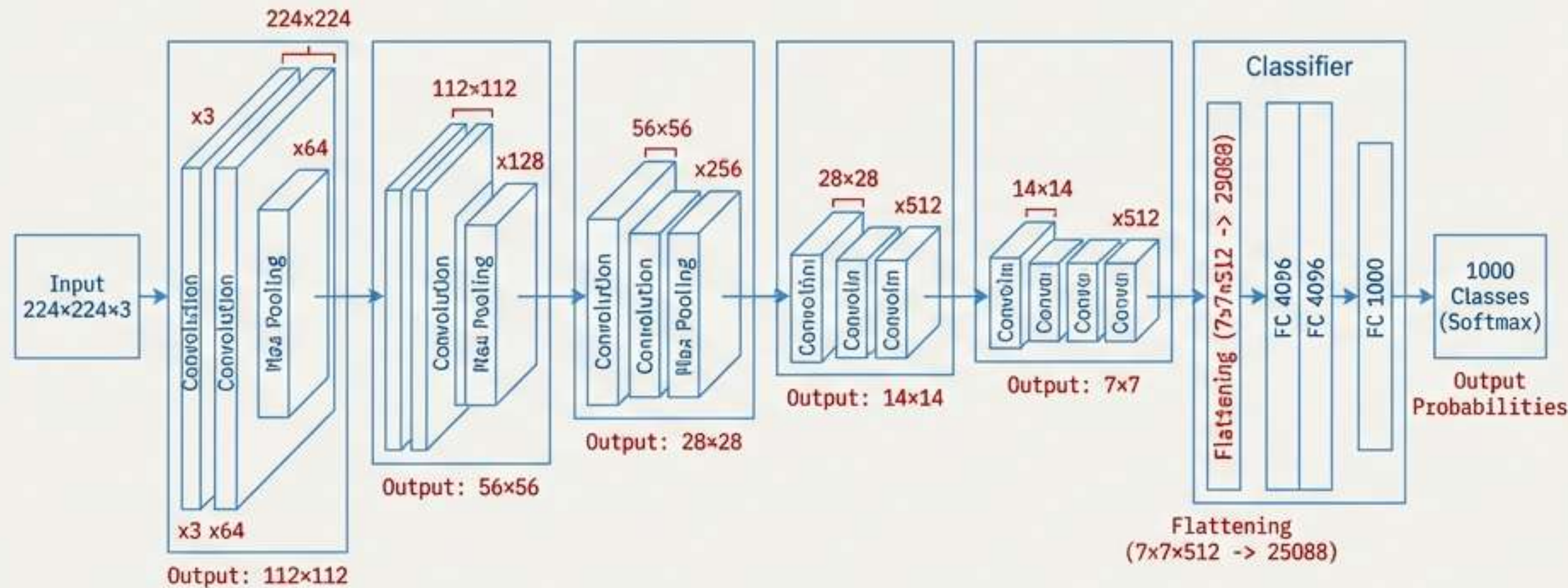
## Effect

Larger batch sizes require more memory but can decrease processing time. Smaller batch sizes require less memory but take longer to process the full dataset.

"The teacher said you will have 6 tests (epochs) of 1 hour during the whole year and each test will have 30 questions (batch size)."

# The Complete Blueprint: VGG-16 Deconstructed

**Key Idea:** We can now look at the full VGG-16 architecture and understand the purpose and flow of each component.



## The Journey of an Image

1. The 224×224×3 input enters the first convolutional block.

2. As it passes through the five blocks, we observe the spatial dimensions systematically decrease (224 -> 112 -> 56 -> 28 -> 14 -> 7).

3. Simultaneously, the depth (number of learned features) systematically increases (3 -> 64 -> 128 -> 256 -> 512 -> 512).

4. The final 7x7x512 feature map is flattened and passed to the classifier, which produces probabilities for 1,000 classes.

## Conclusion

This structure of feature extraction followed by classification is a foundational pattern in modern computer vision. By understanding these core building blocks, you can begin to analyze, use, and even design your own powerful neural networks.

NotebookLM