

1 2 3 4 5 6 7 8 9 10 11 12 13 14 0

# VISION IN CODE

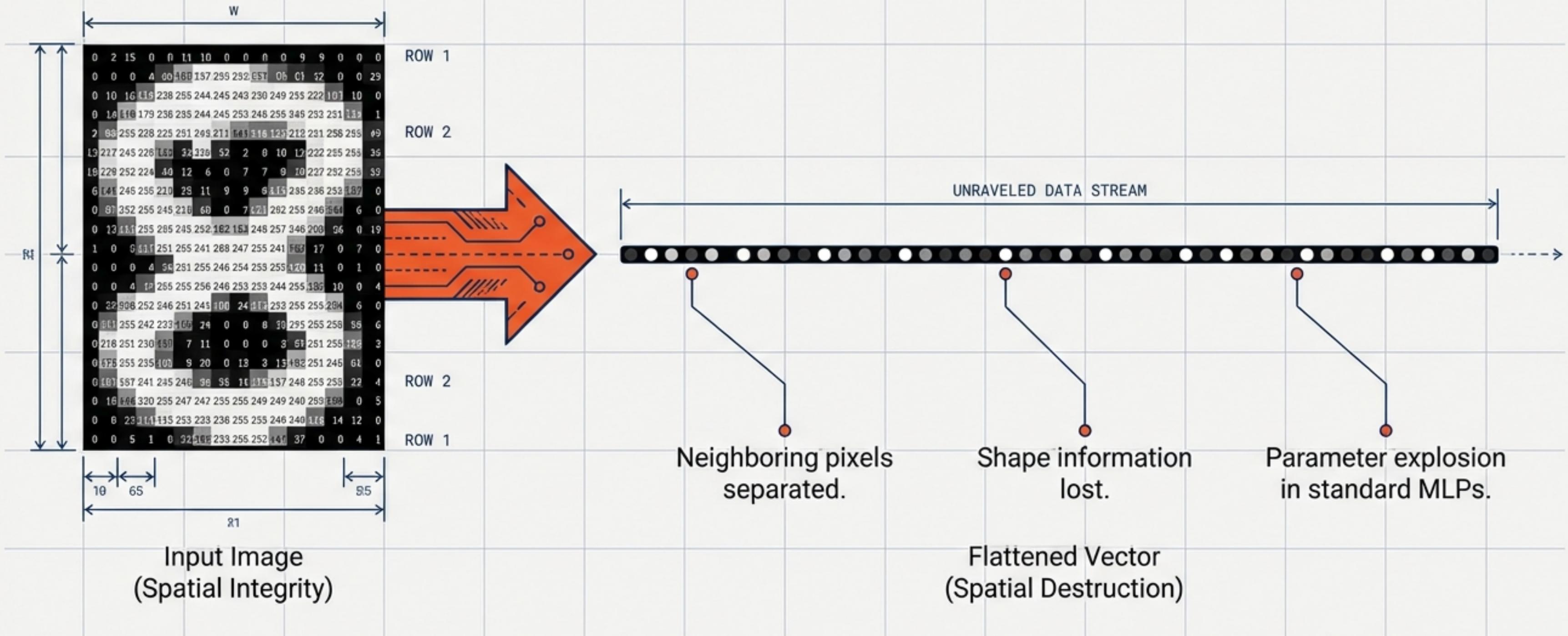
The Mechanics of Convolutional Neural Networks



0 10 20 30 40 50 60 70 80 90 100 110 120

0 25 45 120

# THE PROBLEM WITH FLATTENING.



# THE SOLUTION: LOCAL RECEPTIVE FIELDS.

125	100	50	200	23	45	112	67	89	32
180	140	210							
8	22	45	55	10	78	110	10	78	99
240		15	160	8	22	45	11	45	69
22	45	11	67	87	88	90	13	250	30
120	88	90	64	51	92	23	47	0	
63	140	55	48	160	9	78	240	160	15
42	180	160	15	8	22	45	11	66	11
8	22	45	11	67	88	90	13	250	30
15	8	250	13	250	19	23	45	45	0

## 01. SPATIAL AWARENESS

Pixels are only related to their immediate neighbors. Distant pixels are irrelevant.

## 02. PARAMETER SHARING

The same filter (weights) scans the entire image. A vertical edge in the top-left is the same feature as a vertical edge in the bottom-right.

## 03. THE FLASHLIGHT METAPHOR

We scan the dark room with a beam, processing one patch at a time.

# THE RAW INPUT.

255	128	0	54	255
128	0	54	255	0
54	255	128	0	54
78	54	54	100	200
255	78	150	150	255
20	78	150	45	180
255	20	45	45	255

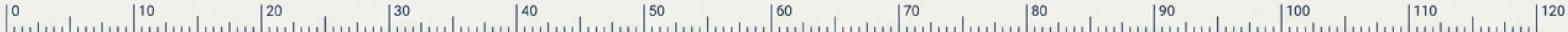
Raw Pixel Data [0 - 255]



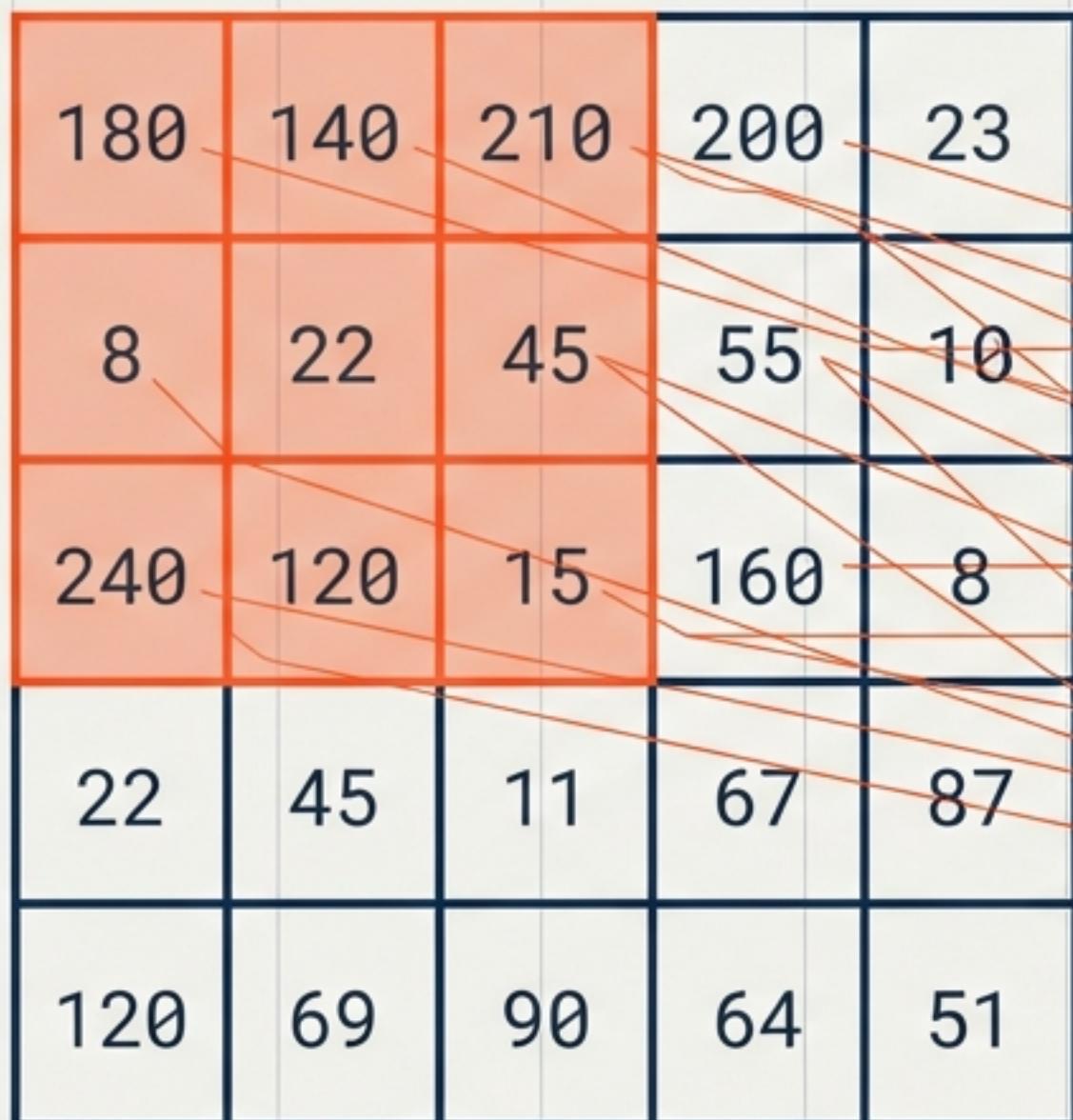
1.0	0.502	0.0	0.212	0.212	0.392
0.212	0.392	0.784	0.784	0.588	0.304
0.306	0.506	0.588	0.784	0.588	0.706
1.0	0.078	0.176	0.706	0.176	0.706
0.176	0.706	0.588	0.304	0.588	0.706

Normalized Tensor [0.0 - 1.0]

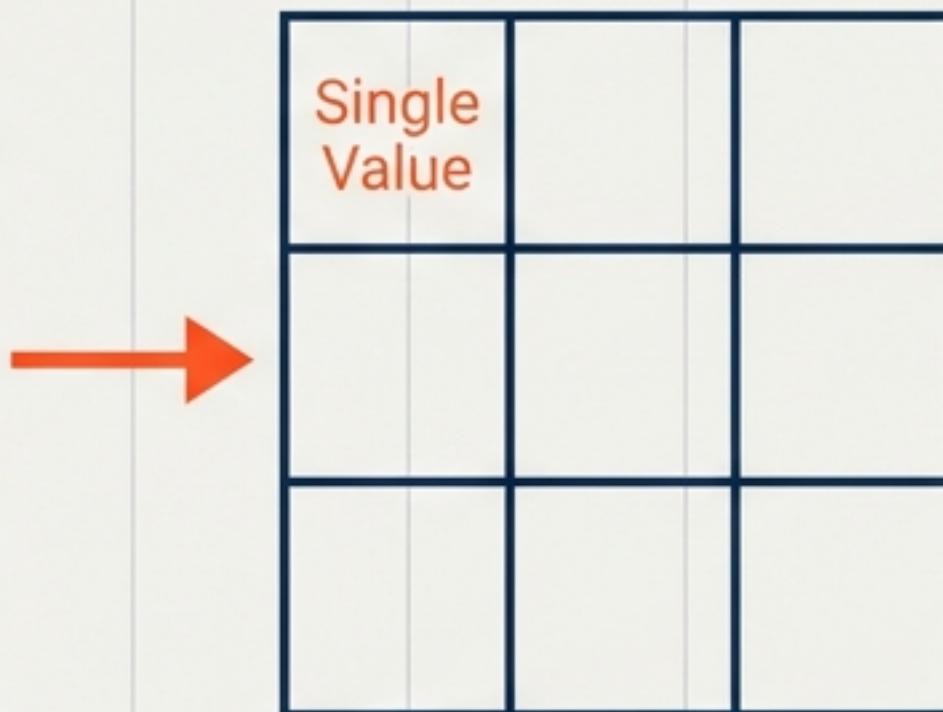
WHY NORMALIZE? Large integers destabilize gradients. Small floats (0-1) ensure efficient convergence.



# THE ENGINE: CONVOLUTION.



$$\begin{aligned} & (180 \times w_1) \\ & + (140 \times w_2) \\ & + (210 \times w_3) \\ & + (8 \times w_4) \\ & + (22 \times w_5) \\ & + (45 \times w_6) \\ & + (240 \times w_7) \\ & + (120 \times w_8) \\ & + (15 \times w_9) \\ = & \text{Single Value} \end{aligned}$$



## THE KERNEL (FILTER)

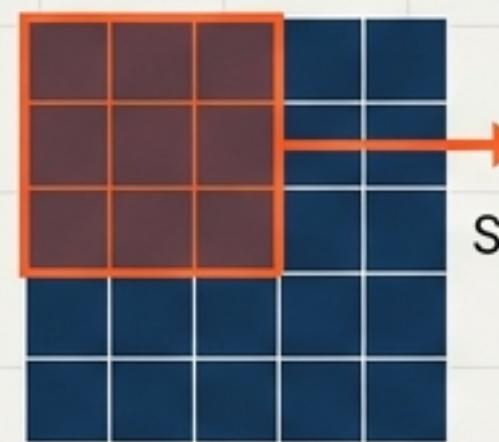
A tiny matrix of learnable weights that slides over the input.

## THE DOT PRODUCT

Element-wise multiplication followed by a sum. Measures how much the input patch matches the filter pattern.

# CONTROLLING THE SCAN

## Stride (s)



Stride 1



Stride 2

Larger stride = Smaller output.

## Padding (p)

0	0	0	0	0	0	0
0						0
0						0
0						0
0						0
0	0	0	0	0	0	0

Zero-padding preserves spatial dimensions.

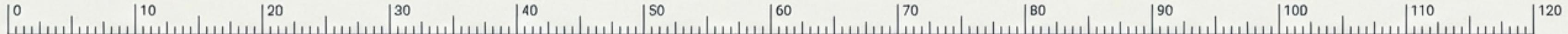
$$n_{\text{out}} = \text{floor}\left(\frac{n_{\text{in}} + 2p - k}{s}\right) + 1$$

n\_in: Input size

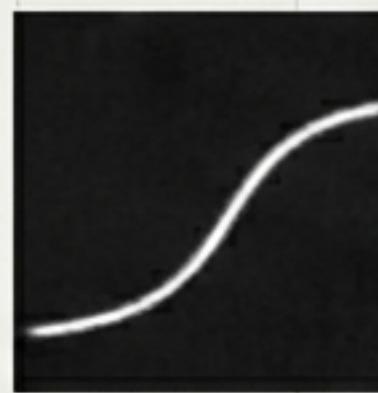
p: Padding

k: Kernel size

s: Stride



# THE OUTPUT: FEATURE MAPS & ReLU.



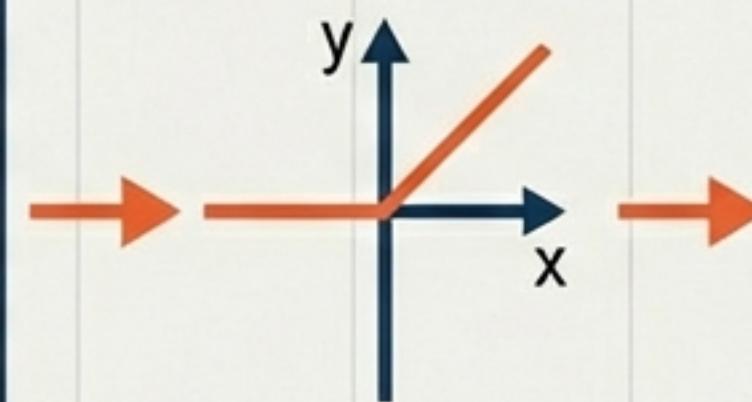
Input



Convolution

120	-15	-22	-45	-12
-15	150	180	-11	-33
-45	-11	200	180	-8
-33	-8	-67	150	-90
-10	-67	-90	-18	120

Feature Map  
(Pre-Activation)



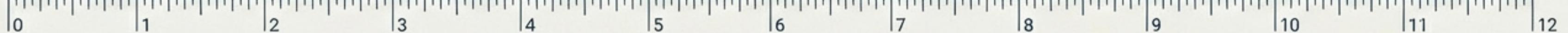
ReLU Function

120	0	0	0	0
0	150	180	0	0
0	0	200	180	0
0	0	0	150	0
0	0	0	0	120

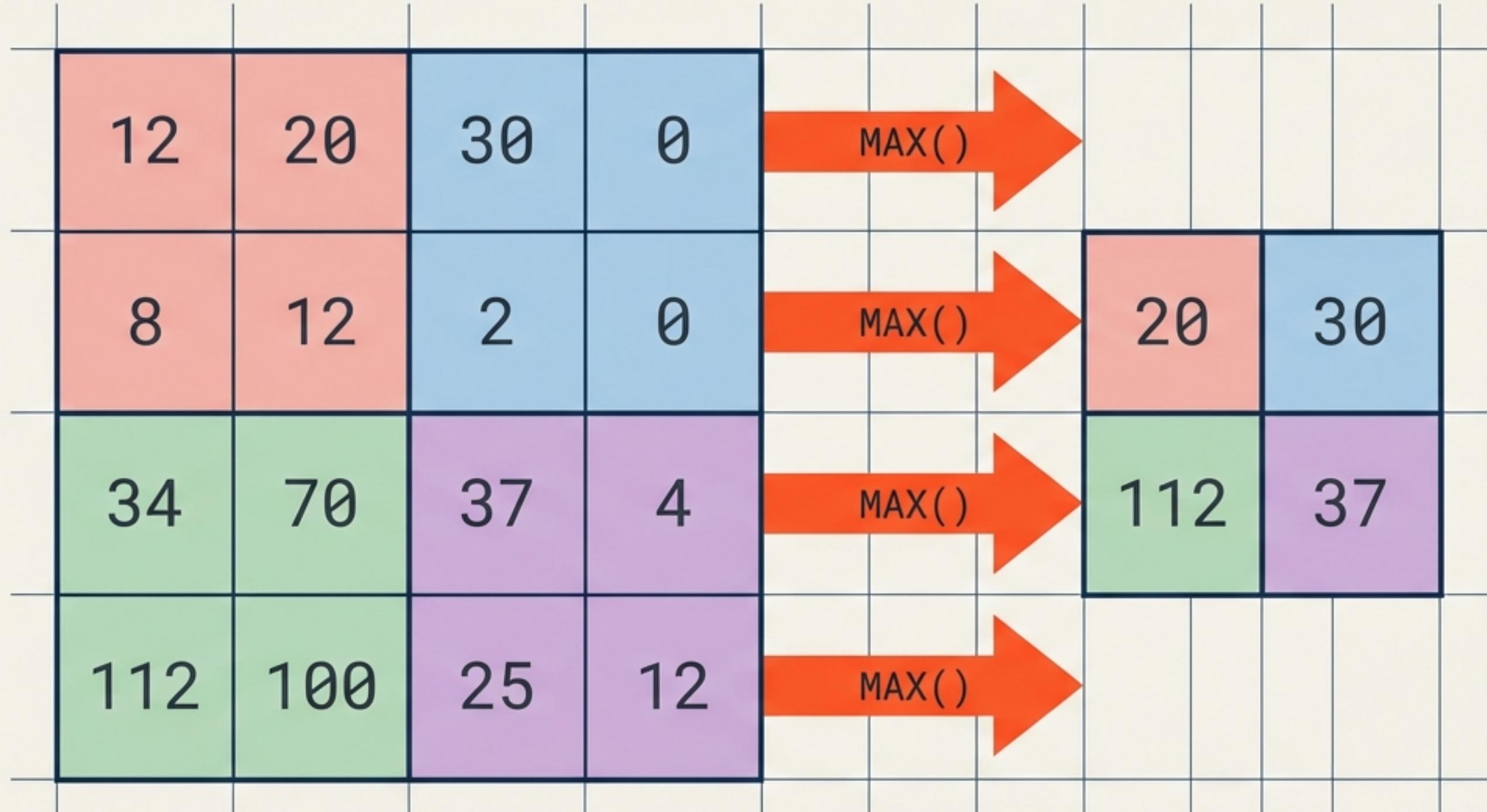
Feature Map  
(Post-Activation)

## ACTIVATION (ReLU)

Rectified Linear Unit. It silences negative values (noise) and keeps positive matches. This non-linearity allows the network to learn complex shapes.



# DOWNSAMPLING: MAX POOLING



## THE PURPOSE

- 1. REDUCE COMPUTATION:** Smaller spatial dimensions = fewer calculations later.
- 2. INVARIANCE:** We care that a feature exists, not precisely where it is.
- 3. DENOISING:** Only the strongest signals survive.

# HIERARCHY OF ABSTRACTION

**High-Level Features**  
Roboto Mono



**Mid-Level Features**  
Roboto Mono

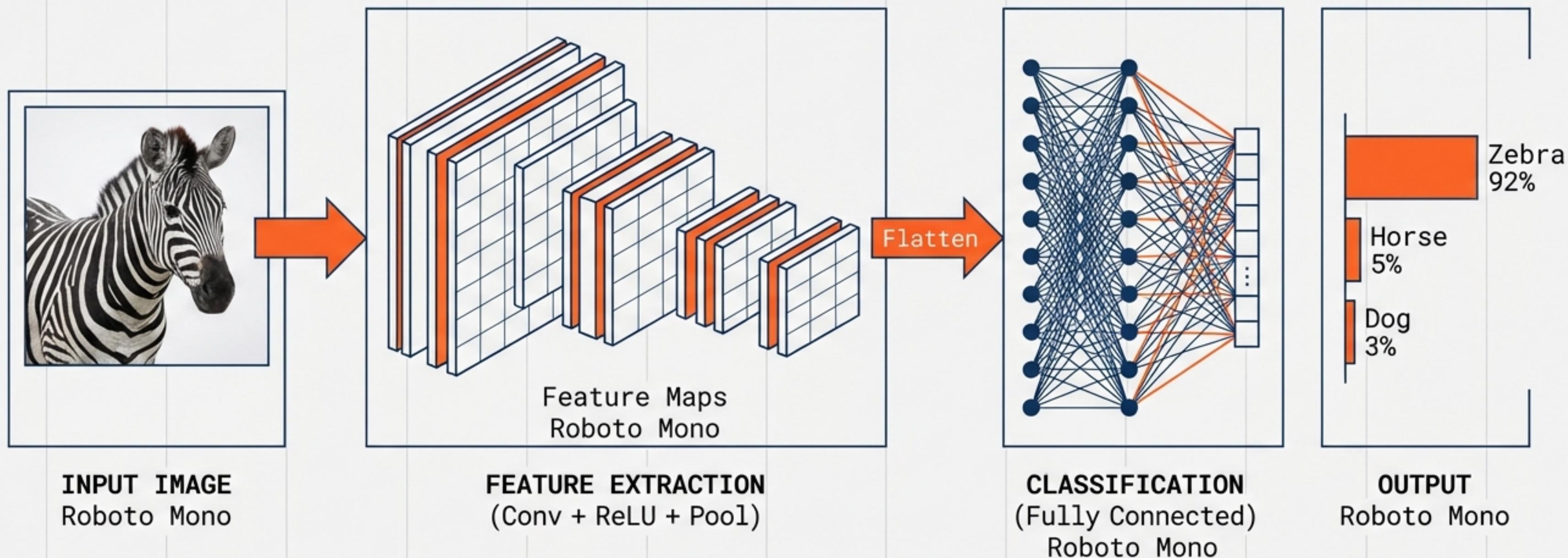


**Low-Level Features**  
Roboto Mono

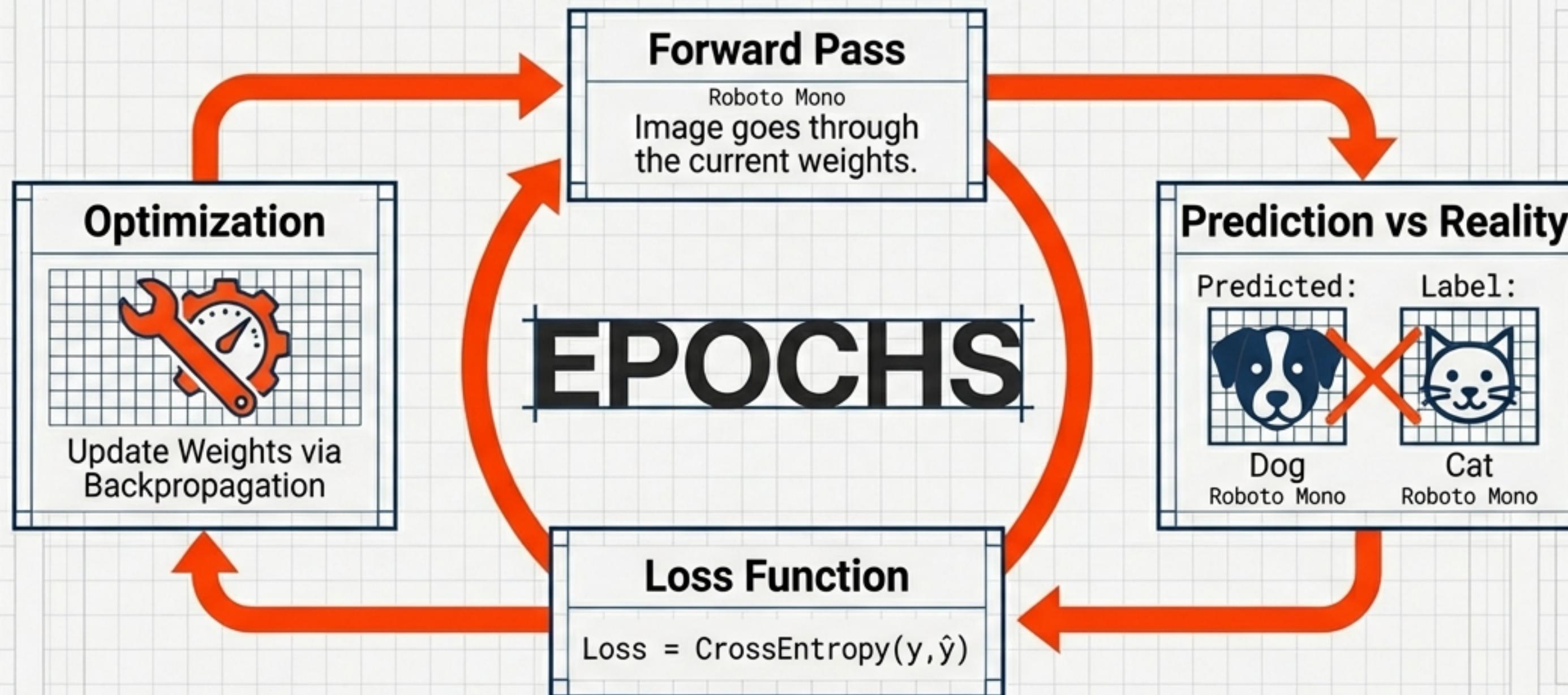


"Deep Learning is the process of building complex understanding from simple beginnings."  
Graphite Roboto Mono

# THE COMPLETE PIPELINE



# HOW IT LEARNS: TRAINING LOOP.



The kernels are not hard-coded. They start as random noise and are sculpted over thousands of iterations to minimize the error (Loss).

# CALCULATING DIMENSIONS: AN EXAMPLE.

<p>Input Image: <math>32 \times 32 \times 3</math></p>	<p>The diagram illustrates the dimension changes through two layers. It starts with an input image of size <math>32 \times 32 \times 3</math>. After Layer 1 (Convolution), the dimensions become <math>14 \times 14 \times 6</math>. A red arrow points from the input to the output, with a vertical line at <math>32 \times 32</math> and another at <math>14 \times 14</math>. A downward arrow labeled "3 Channels" points from the input to the output, and an upward arrow labeled "6 Channels" points from the output back to the input.</p>
<p><b>Layer 1: Convolution</b></p> <p>(6 Kernels, Size 5x5, Stride 1)</p>	<p><b>Layer 2: Max Pooling</b></p> <p>(Size 2x2, Stride 2)</p>
<p>Calculation:</p> $(32 - 5 + 2(0)) / 1 + 1$ $= 28$	<p>Calculation:</p> $(28 - 2) / 2 + 1$ $= 14$
<p>Resulting Shape:</p> $28 \times 28 \times 6$	<p>Resulting Shape:</p> $14 \times 14 \times 6$

# FROM THEORY TO CODE (PyTorch)

```
class CNN_Model(nn.Module):
    def __init__(self):
        super().__init__()
        # Input: 3 channels (RGB), Output: 6 feature maps
        self.conv1 = nn.Conv2d(3, 6, 5)
        self.pool = nn.MaxPool2d(2, 2)
        self.conv2 = nn.Conv2d(6, 16, 5)

    def forward(self, x):
        # Convolution -> ReLU -> Pooling
        x = self.pool(F.relu(self.conv1(x)))
        x = self.pool(F.relu(self.conv2(x)))
        return x
```

Defines Kernels  
Downsampling

# CRITICAL ANALYSIS.

## PROS

 **Spatial Invariance:** Recognizes objects regardless of position in the frame.

 **Parameter Efficiency:** Shared weights reduce model size compared to dense networks.

 **State of the Art:** The gold standard for image and video tasks.

## CONS

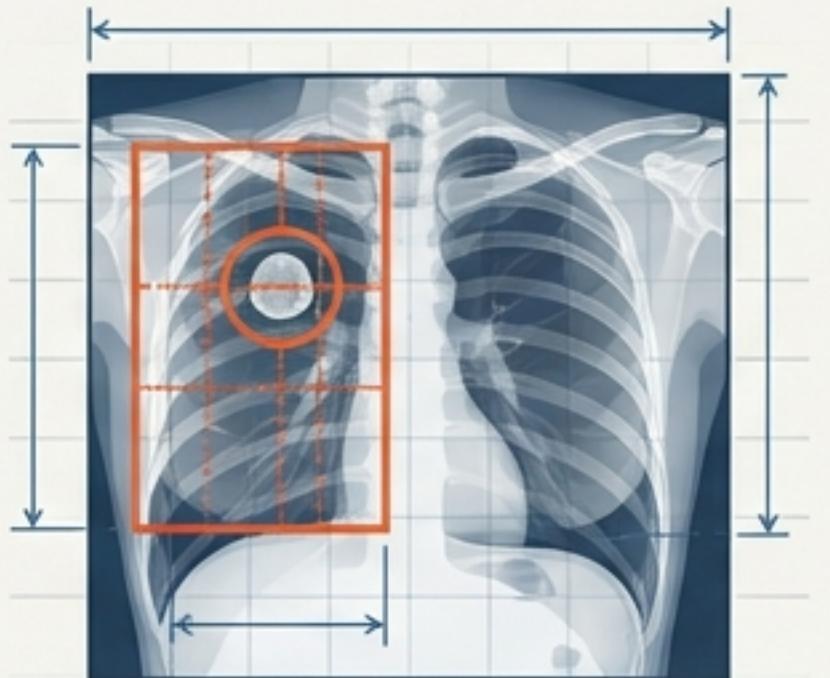
 **Compute Cost:** Training requires significant GPU acceleration.

 **Data Hungry:** Needs massive labeled datasets (e.g., ImageNet) to generalize.

 **Sequence Struggle:** Less effective than Transformers for pure text/NLP tasks.

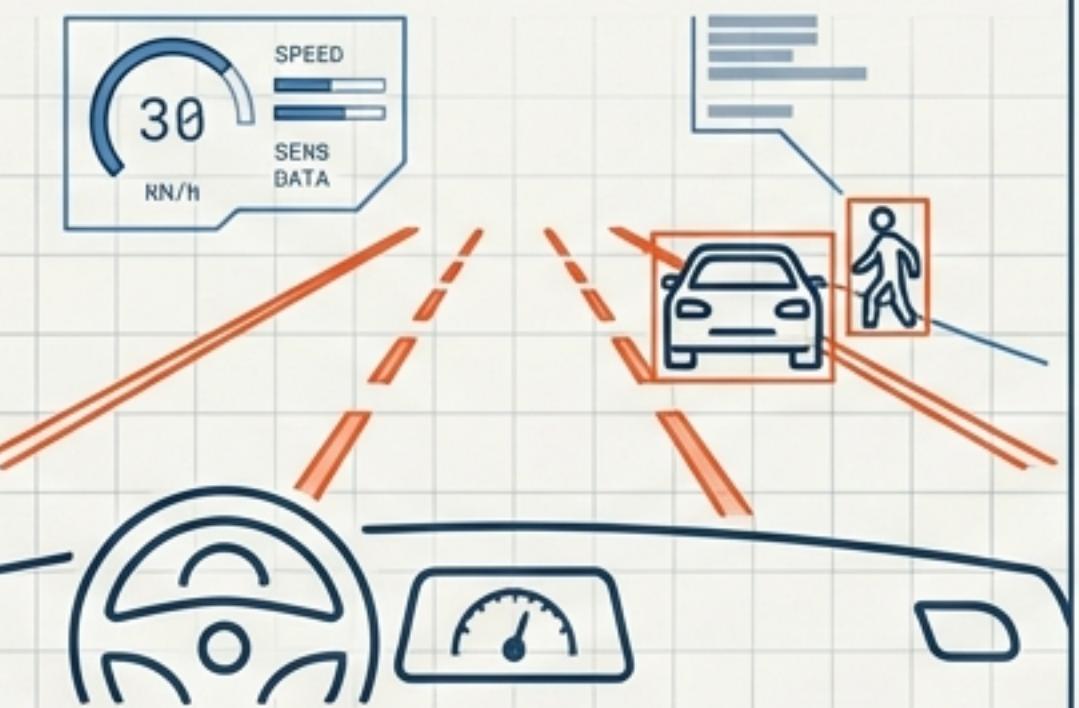
# REAL WORLD IMPACT.

Medical Imaging



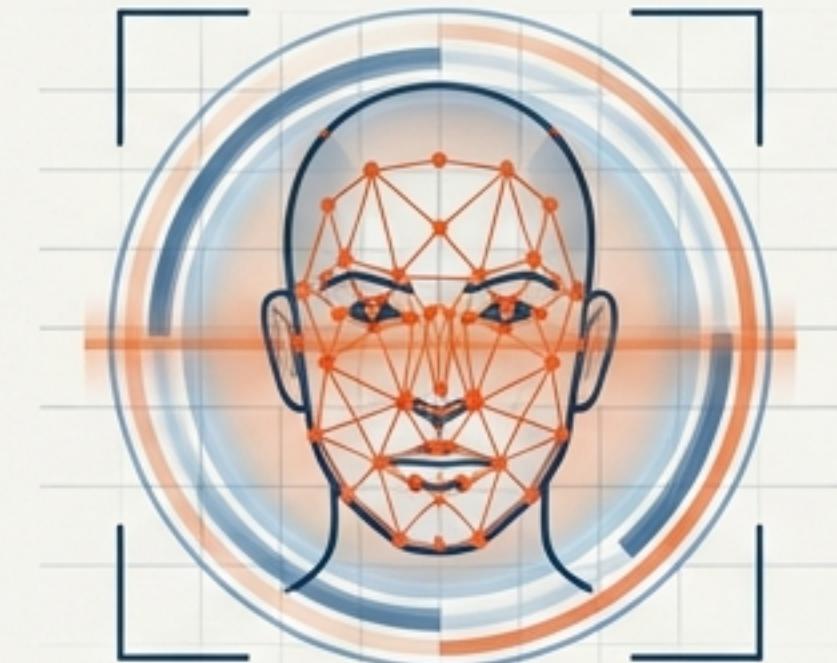
Tumor Detection & Diagnostics

Autonomous Systems



Self-Driving Perception

Security



Biometric Authentication

**While architecture evolves, the Convolutional Neural Network remains the bedrock of Computer Vision.**