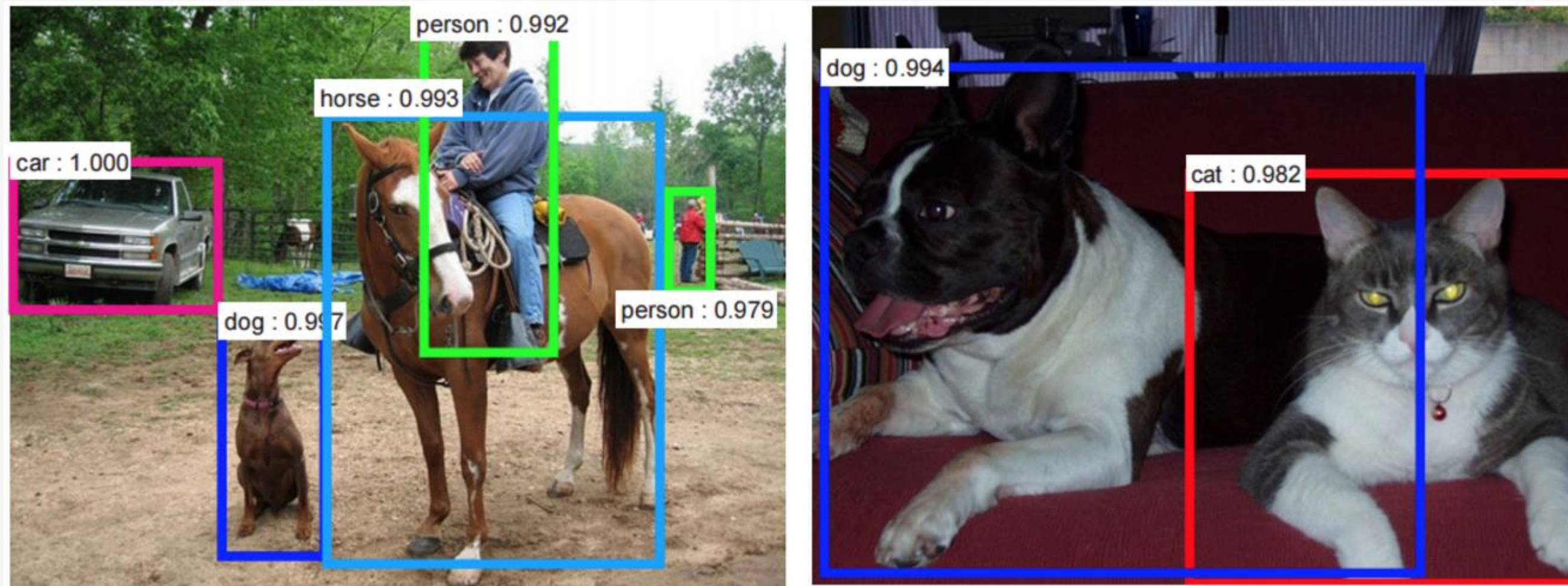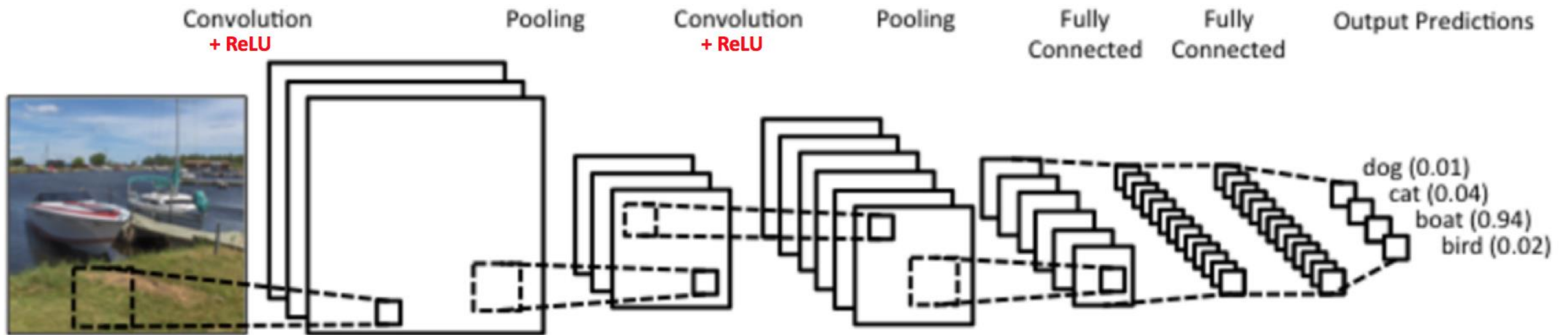# Convolutional Neural Networks
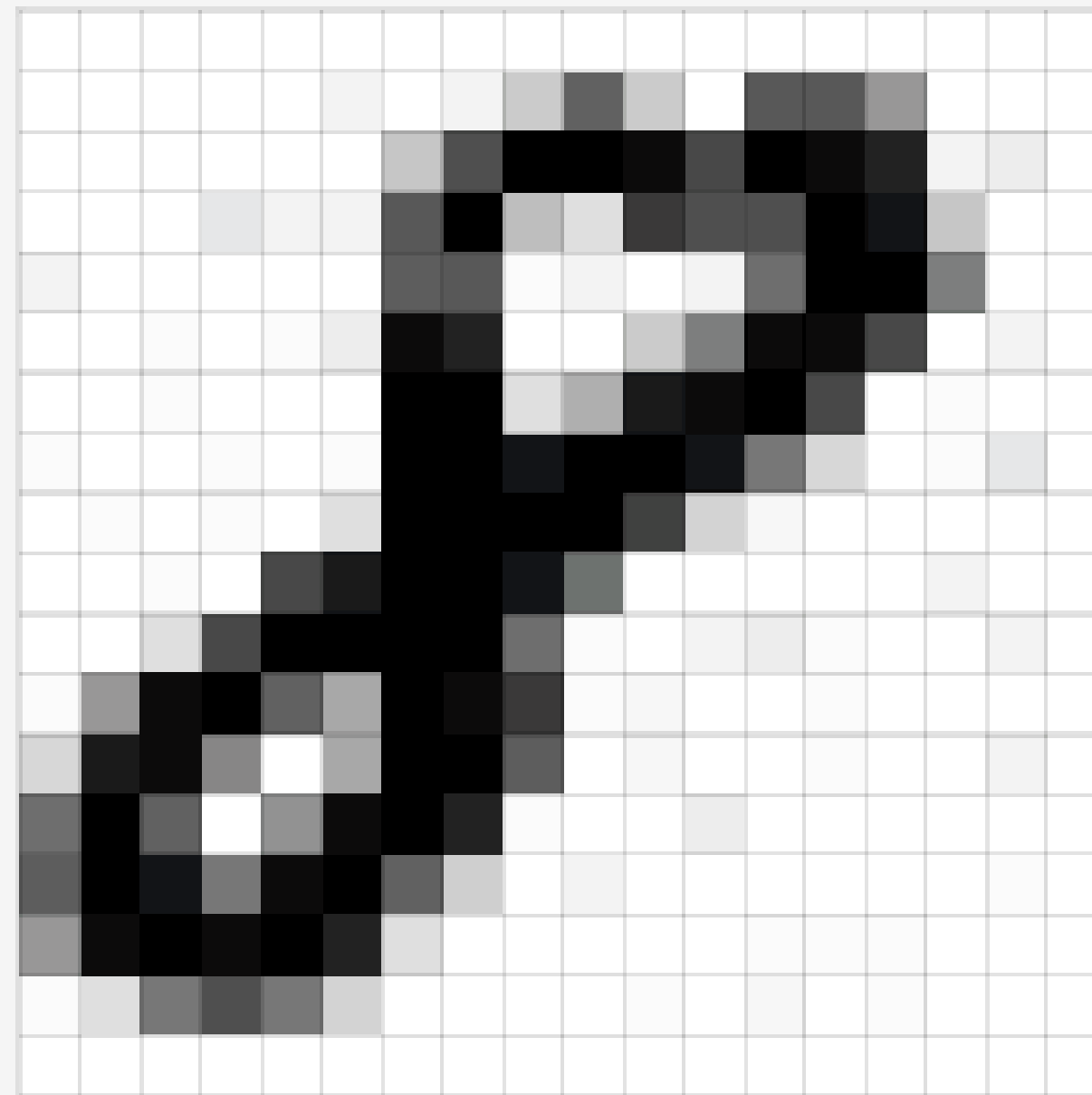
# Convolution Layer

# An Image is a matrix of pixel values

Essentially, every image can be represented as a matrix of pixel values.

## The Convolution Step

ConvNets derive their name from the "convolution operator". The primary purpose of Convolution in case of a ConvNet is to extract features from the input image. Convolution preserves the spatial relationship between pixels by learning image features using small squares of input data. We will not go into the mathematical details of Convolution here, but will try to understand how it works over images.

Consider a 5 x 5 image whose pixel values are only 0 and 1 (note that for a grayscale image, pixel values range from 0 to 255, the green matrix below is a special case where pixel values are only 0 and 1):

Also, consider another 3 x 3 matrix as shown below:

| 1 | 1 | 1 | 0 | 0 |
|---|---|---|---|---|
| 0 | 1 | 1 | 1 | 0 |
| 0 | 0 | 1 | 1 | 1 |
| 0 | 0 | 1 | 1 | 0 |
| 0 | 1 | 1 | 0 | 0 |

| 1 | 0 | 1 |
|---|---|---|
| 0 | 1 | 0 |
| 1 | 0 | 1 |

| $1_{\times1}$ | $1_{\times0}$ | $1_{\times1}$ | 0 | 0 |
|---|---|---|---|---|
| $0_{\times0}$ | $1_{\times1}$ | $1_{\times0}$ | 1 | 0 |
| $0_{\times1}$ | $0_{\times0}$ | $1_{\times1}$ | 1 | 1 |
| 0 | 0 | 1 | 1 | 0 |
| 0 | 1 | 1 | 0 | 0 |

Image

| 4 | | |
|---|---|---|
| | | |
| | | |

Convolved Feature

In the table below, we can see the effects of convolution of the above image with different filters. As shown, we can perform operations such as Edge Detection, Sharpen and Blur just by changing the numeric values of our filter matrix before the convolution operation
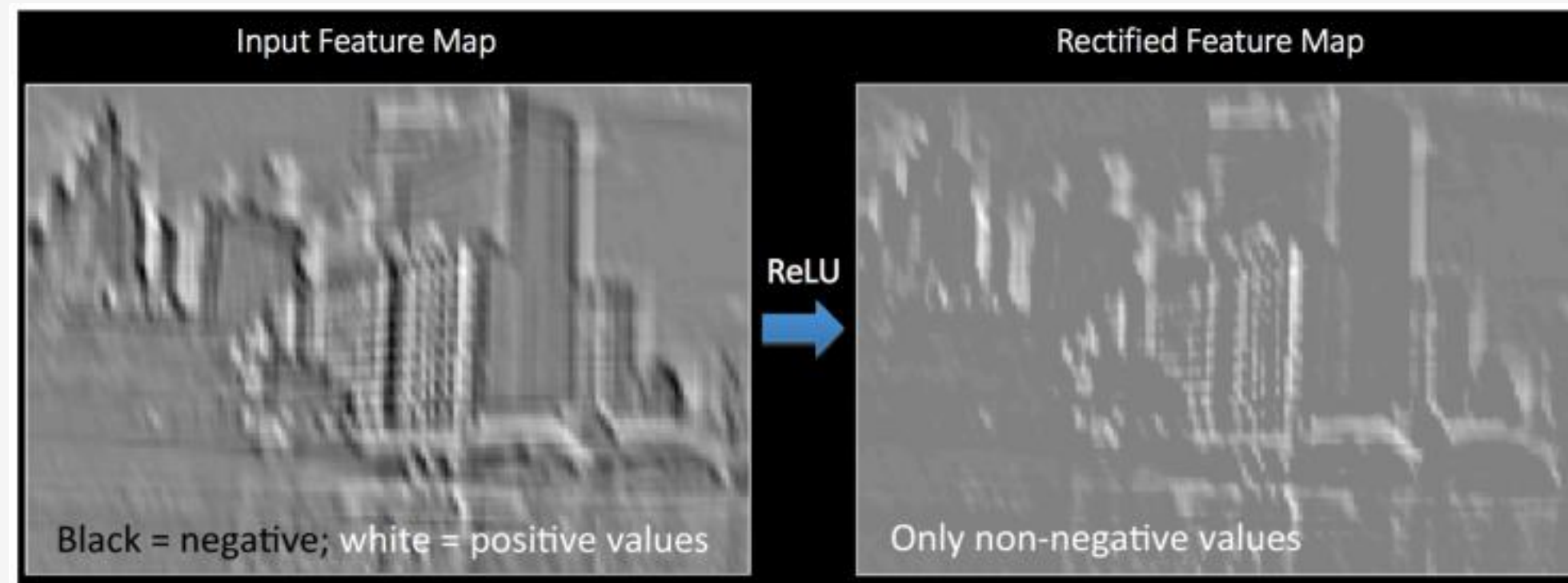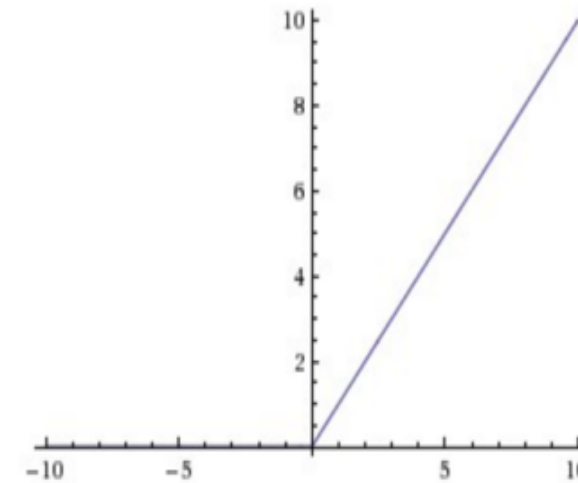
| Operation | Filter | Convolved Image |
|---|---|---|
| Identity | $\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$ | |
| Edge detection | $\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix}$ | |
| | $\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$ | |
| | $\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$ | |
| Sharpen | $\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$ | |
| Box blur (normalized) | $\frac{1}{9}\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$ | |
| Gaussian blur (approximation) | $\frac{1}{16}\begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$ | |

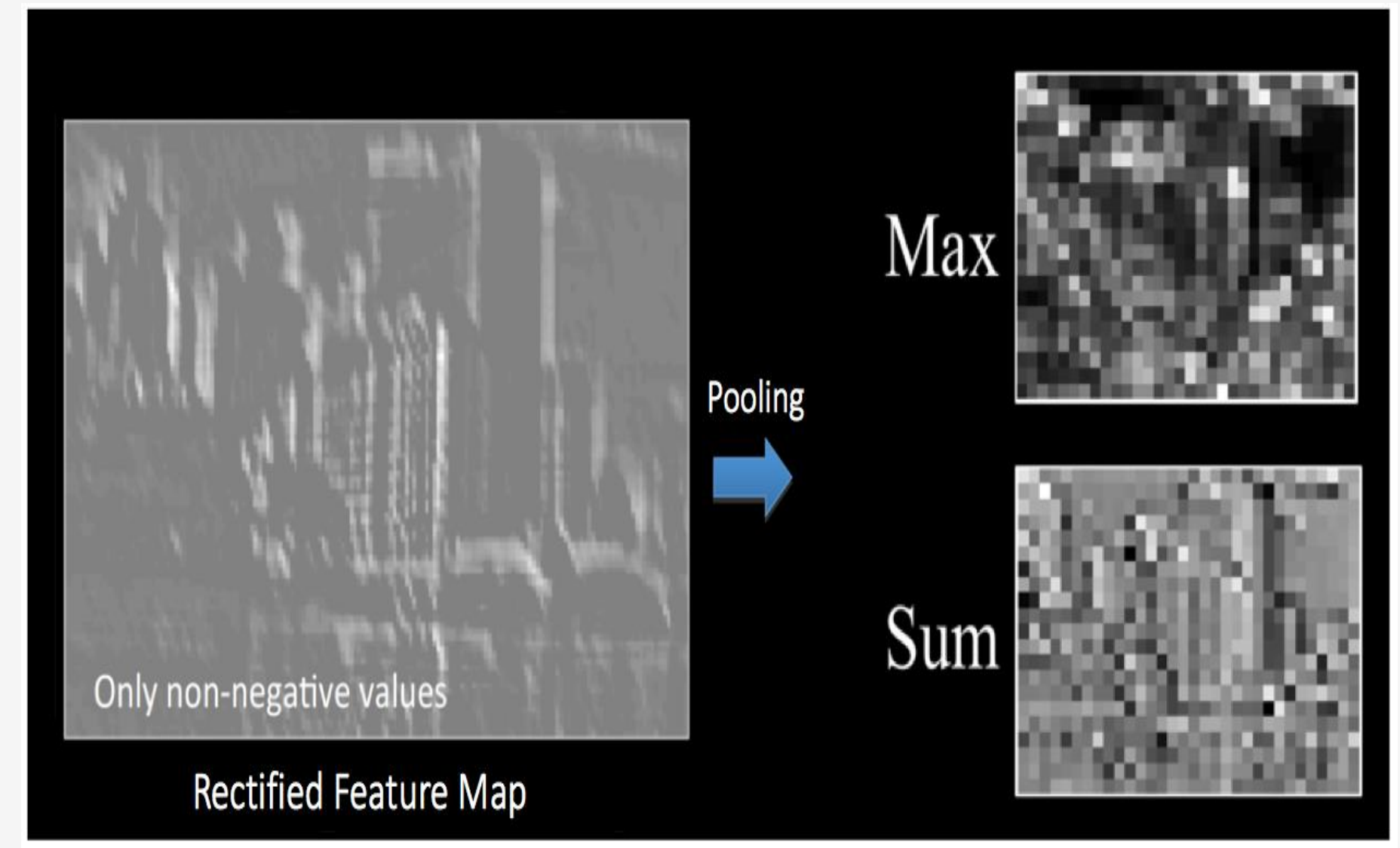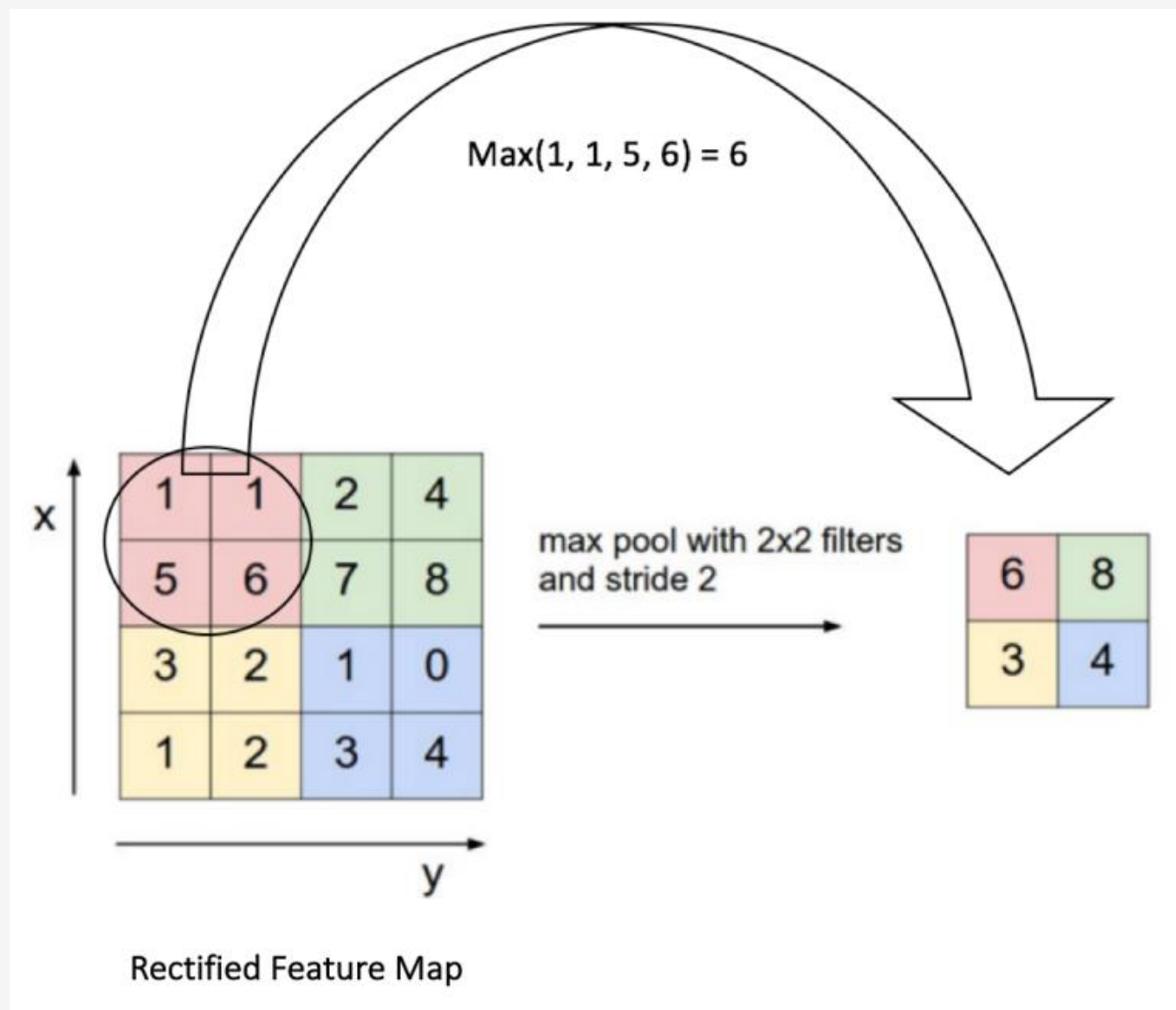Input

## Introducing Non Linearity (ReLU)

An additional operation called ReLU has been used after every Convolution operation

Output = Max(zero, Input)

Input Feature Map → ReLU → Rectified Feature Map

Black = negative; white = positive values
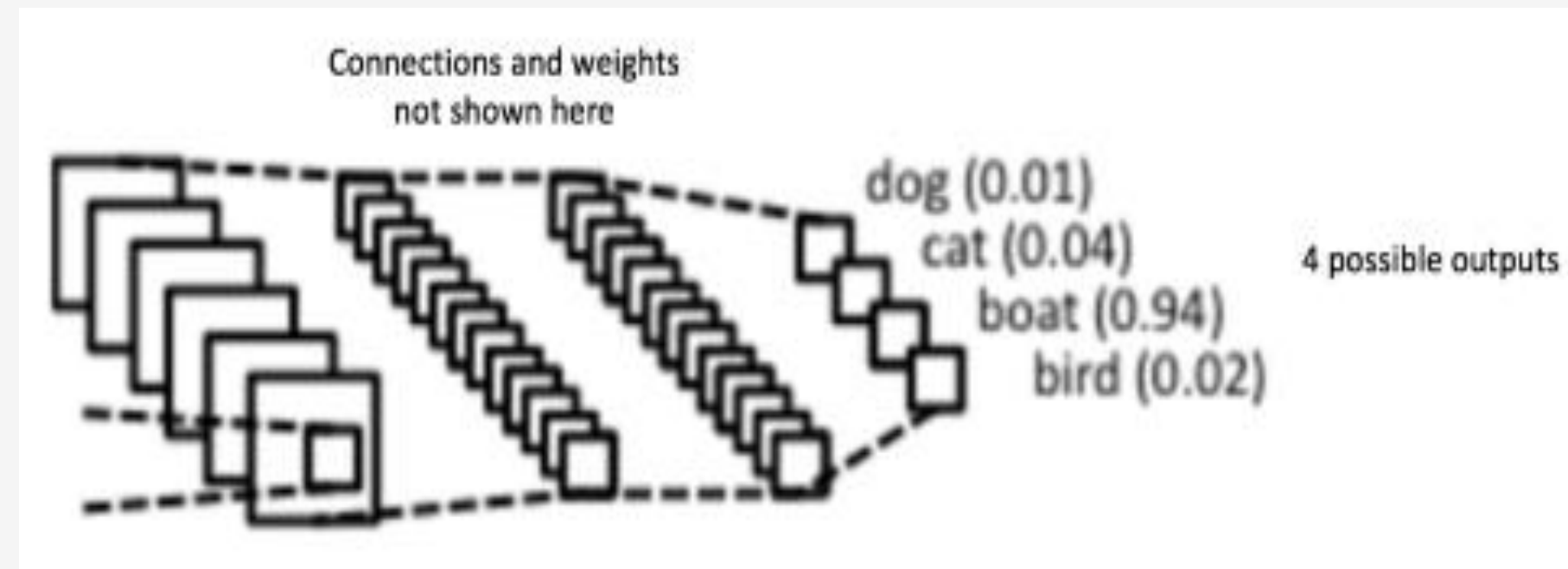
Only non-negative values

## The Pooling Step

Spatial Pooling (also called subsampling or downsampling) reduces the dimensionality of each feature map but retains the most important information. Spatial Pooling can be of different types: Max, Average, Sum etc.



Rectified Feature Map

## Fully Connected Layer

The Fully Connected layer is a traditional Multi Layer Perceptron that uses a softmax activation function in the output layer (other classifiers like SVM can also be used, but will stick to softmax in this post). The term "Fully Connected" implies that every neuron in the previous layer is connected to every neuron on the next layer.

Connections and weights
not shown here

dog (0.01)
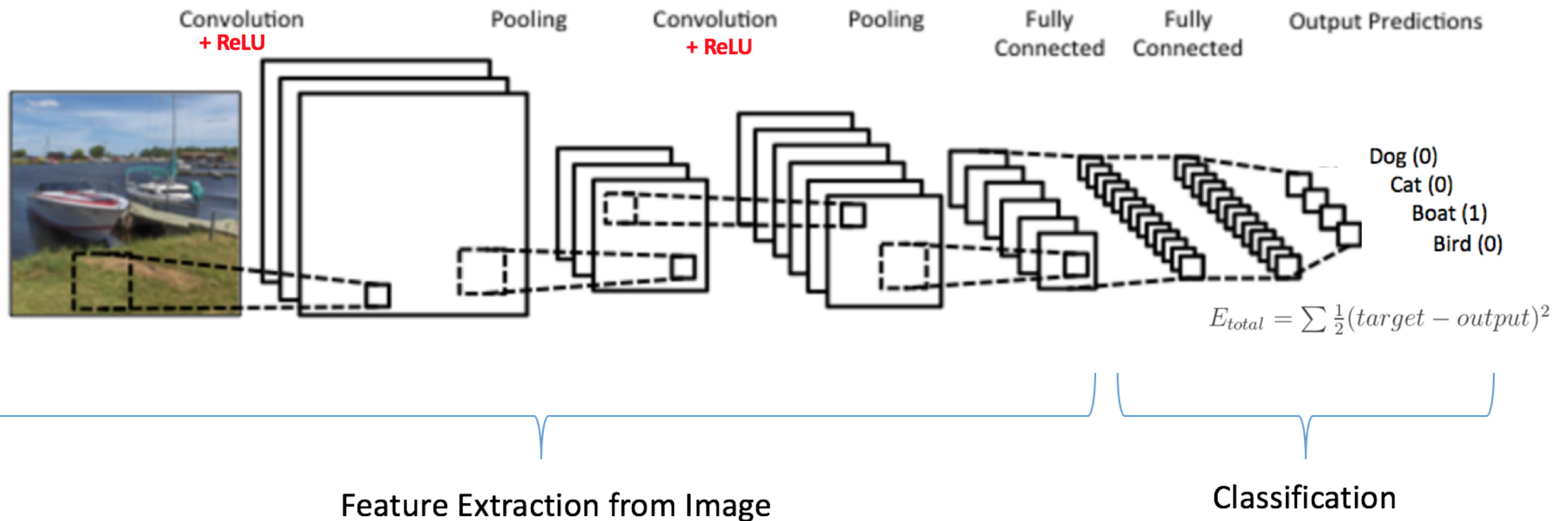cat (0.04)
boat (0.94)
bird (0.02)

4 possible outputs

Putting it all together – Training using Backpropagation
As discussed above, the Convolution + Pooling layers act as Feature Extractors from the input image while Fully Connected layer acts as a classifier.
Note that below, since the input image is a boat, the target probability is 1 for Boat class and 0 for other three classes, i.e.
•Input Image = Boat
•Target Vector = [0, 0, 1, 0]



Convolution + ReLU   Pooling   Convolution + ReLU   Pooling   Fully Connected   Fully Connected   Output Predictions

Dog (0)
Cat (0)
Boat (1)
Bird (0)

$$E_{total} = \sum \tfrac{1}{2}(target - output)^2$$

Feature Extraction from Image

Classification

# Recurrent neural networks

- Dates back to (Rumelhart *et al.*, 1986)
- A family of neural networks for handling sequential data, which involves variable length inputs or outputs

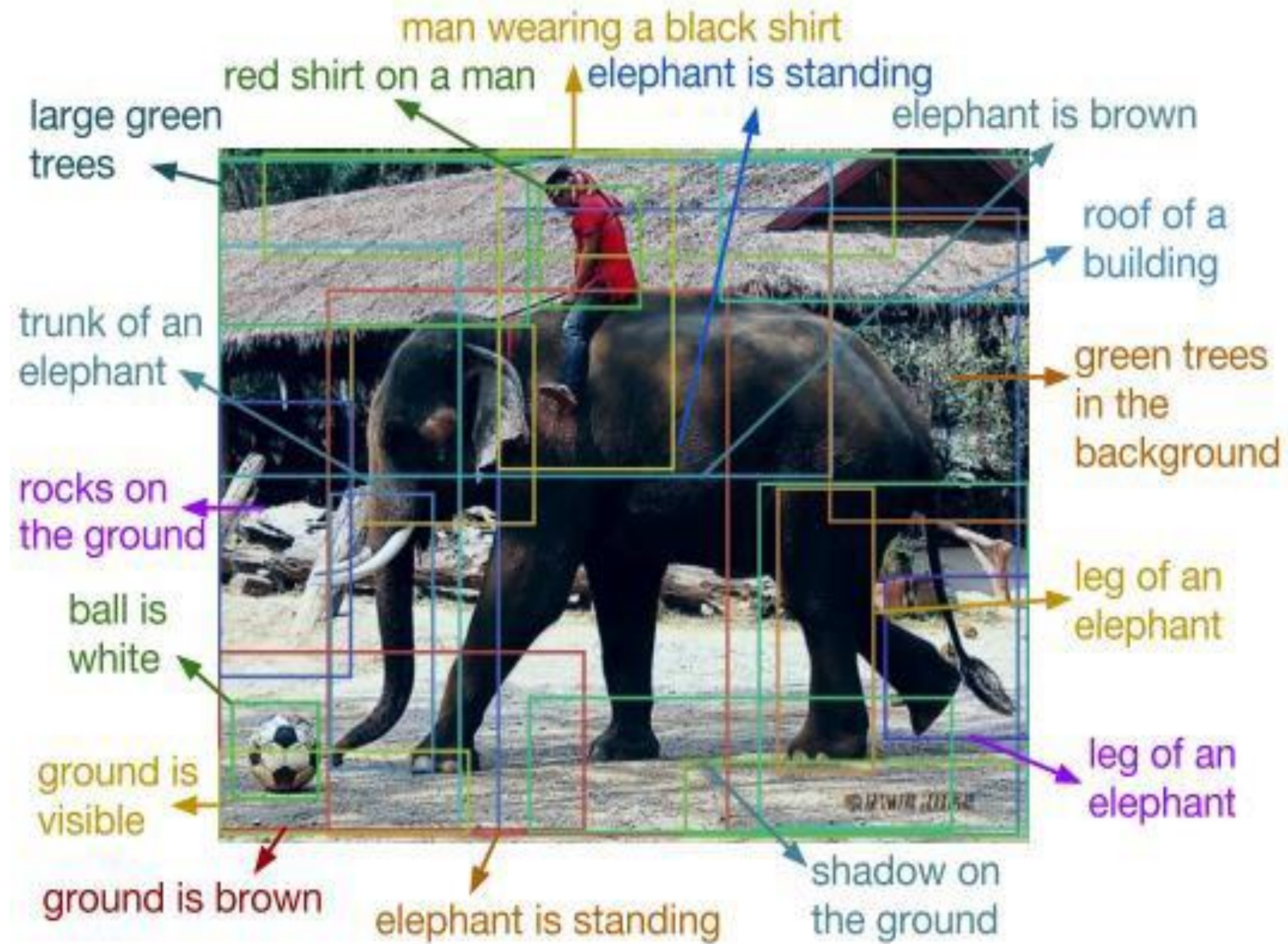- Especially, for natural language processing (NLP)

NAC

# Sequential data

- Each data point: A sequence of vectors $x^{(t)}$, for $1 \leq t \leq \tau$
- Batch data: many sequences with different lengths $\tau$
- Label: can be a scalar, a vector, or even a sequence

- Example
  - Sentiment analysis
  - Machine translation

# More complicated sequential data

- Data point: two dimensional sequences like images
- Label: different type of sequences like text sentences

- Example: image captioning

NAC

# Image captioning

# Time- series forecasting



**FIGURE 1.** Illustration of the multi-step time-series forecasting scheme. There are in total five values of interest including the $\hat{r}_{[t]}$, whose ground truth of $x_{[t]}$ is known and highlighted in red color.

# A typical dynamic system



$$s^{(t+1)} = f(s^{(t)})$$

# A system driven by external data



$$s^{(t+1)} = f(s^{(t)}, x^{(t+1)})$$

# Compact view



$$s^{(t+1)} = f(s^{t}, x^{(t+1)})$$

# Recurrent neural networks

- Use the same computational function and parameters across different time steps of the sequence
- Each time step: takes the input entry and the previous hidden state to compute the output entry
- Loss: typically computed every time step

- It is very difficult to train RNNs to retain information over many time steps
- This make is very difficult to learn RNNs that handle long-distance dependencies, such as subject-verb agreement.

# LSTM Network Architecture

# Advantage

- Hidden state: a lossy summary of the past

- Shared functions and parameters: greatly reduce the capacity and good for generalization in learning

- Explicitly use the prior knowledge that the sequential data can be processed by in the same way at different time step (e.g., NLP)

# Transformer

An Encoder Block: same structure, different parameters

# Encoder

Note: The ffnn is independent
for each word.
Hence can be parallelized.

# Self Attention

| Input | Thinking | Machines |
|-------|----------|----------|

Embedding $X_1$ [    ]   $X_2$ [    ]

Queries $q_1$ [    ]   $q_2$ [    ]   $W^Q$

Keys $k_1$ [    ]   $k_2$ [    ]   $W^K$

Values $v_1$ [    ]   $v_2$ [    ]   $W^V$

First we create three vectors
by multiplying input embedding
(1x512)
 $x_i$ with three matrices (64x512):

$q_i = x_i W^Q$
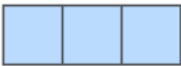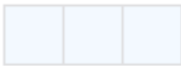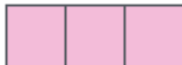$K_i = x_i W^K$
$V_i = x_i W^V$

# Self Attention

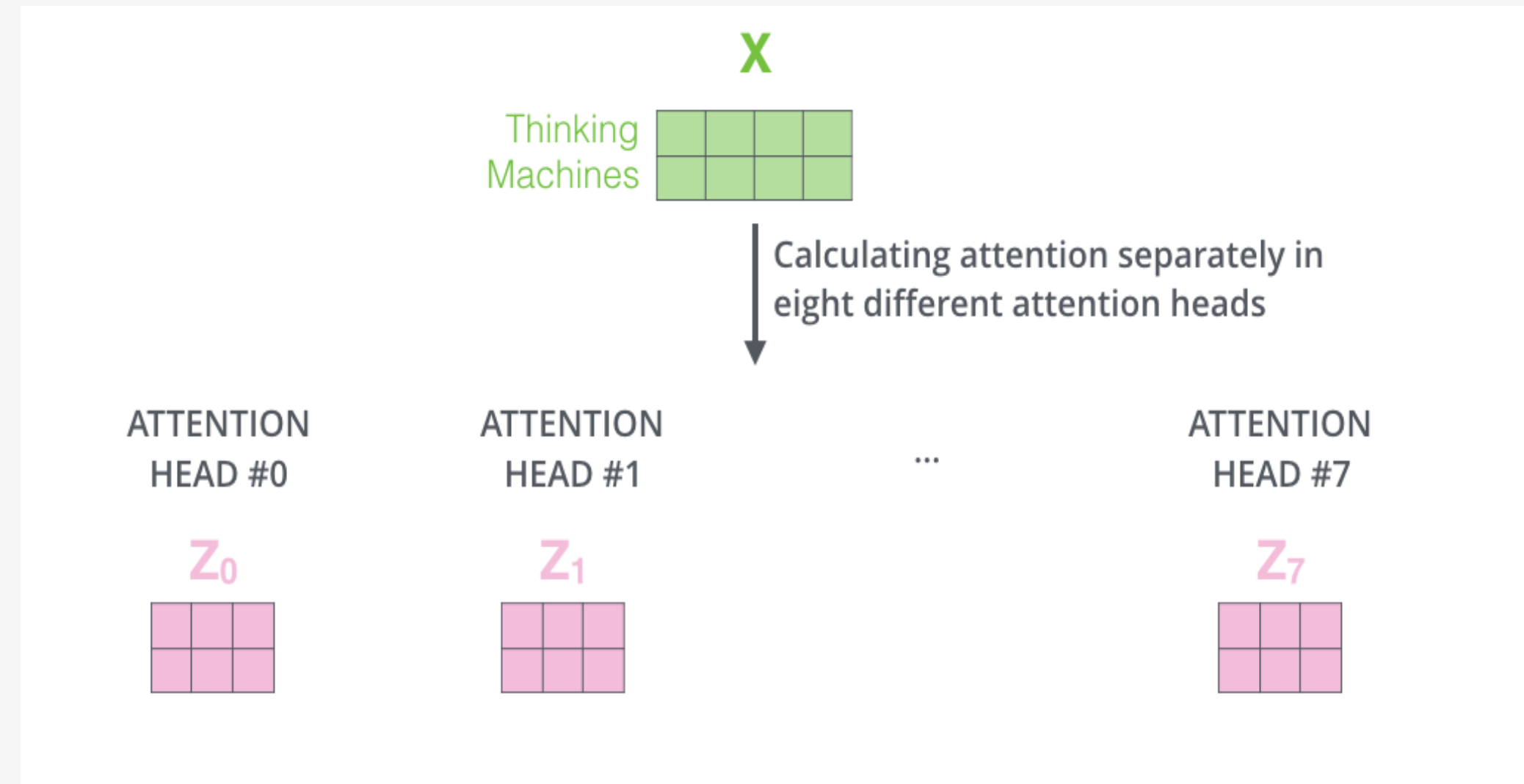Now we need to calculate a score to determine how much focus to place on other Parts of the input.

# Self Attention



| | Thinking | Machines |
|---|---|---|
| Input | | |
| Embedding | $x_1$ | $x_2$ |
| Queries | $q_1$ | $q_2$ |
| Keys | $k_1$ | $k_2$ |
| Values | $v_1$ | $v_2$ |
| Score | $q_1 \bullet k_1 = 112$ | $q_1 \bullet k_2 = 96$ |
| Divide by 8 ($\sqrt{d_k}$) | 14 | 12 |
| Softmax | 0.88 | 0.12 |
| Softmax X Value | $v_1$ | $v_2$ |
| Sum | $z_1$ | $z_2$ |

$$\text{softmax}\left(\frac{Q \times K^T}{\sqrt{d_k}}\right) V = Z$$
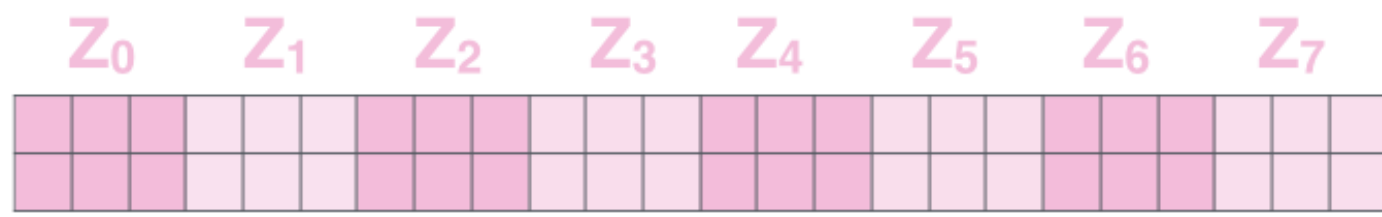
# Multiple heads

1. It expands the model's ability to focus on different positions.
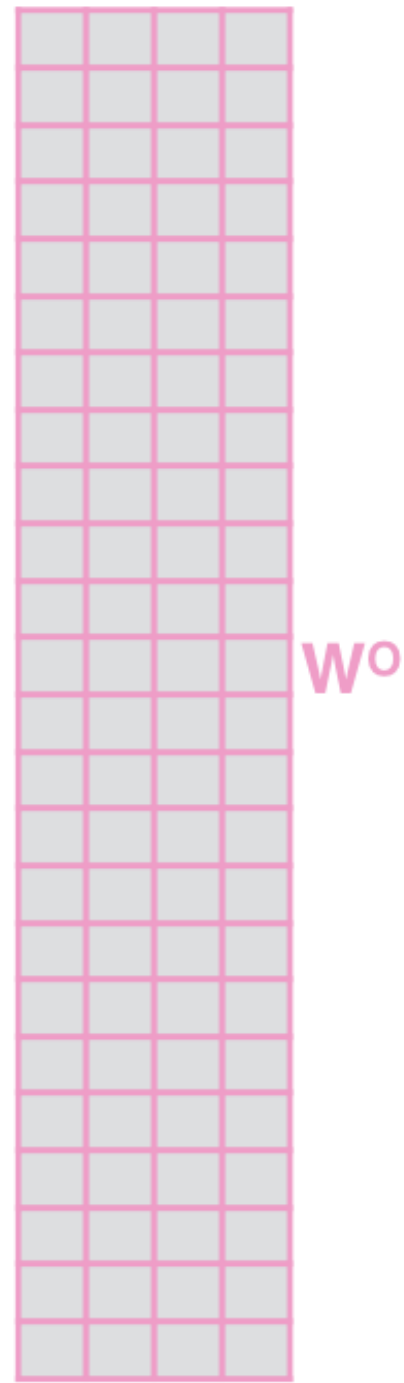2. It gives the attention layer multiple "representation subspaces"

The output is expecting only a 2x4 matrix, hence,

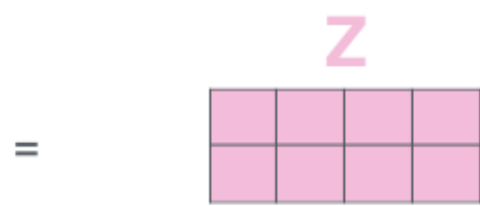1) Concatenate all the attention heads

$Z_0$    $Z_1$    $Z_2$    $Z_3$    $Z_4$    $Z_5$    $Z_6$    $Z_7$

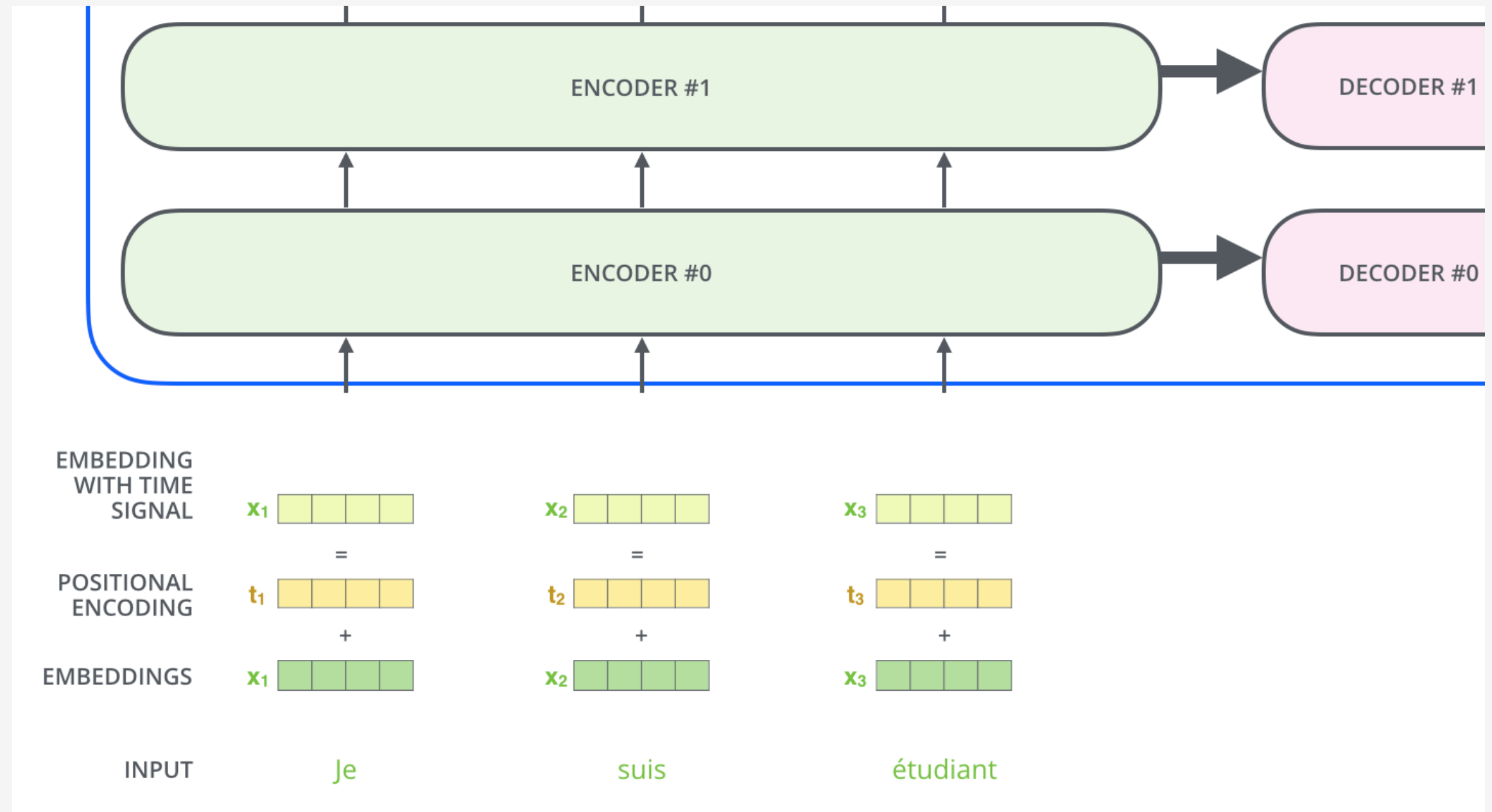2) Multiply with a weight matrix $W^O$ that was trained jointly with the model

X

3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN

Z

=

$W^O$

# Representing the input order (positional encoding)

The transformer adds a vector to each input embedding. These vectors follow a specific pattern that the model learns, which helps it determine the position of each word, or the distance between different words in the sequence. The intuition here is that adding these values to the embedings provides meaningful distances between the embedding vectors once they're projected into Q/K/V vectors and during dot-product attention.
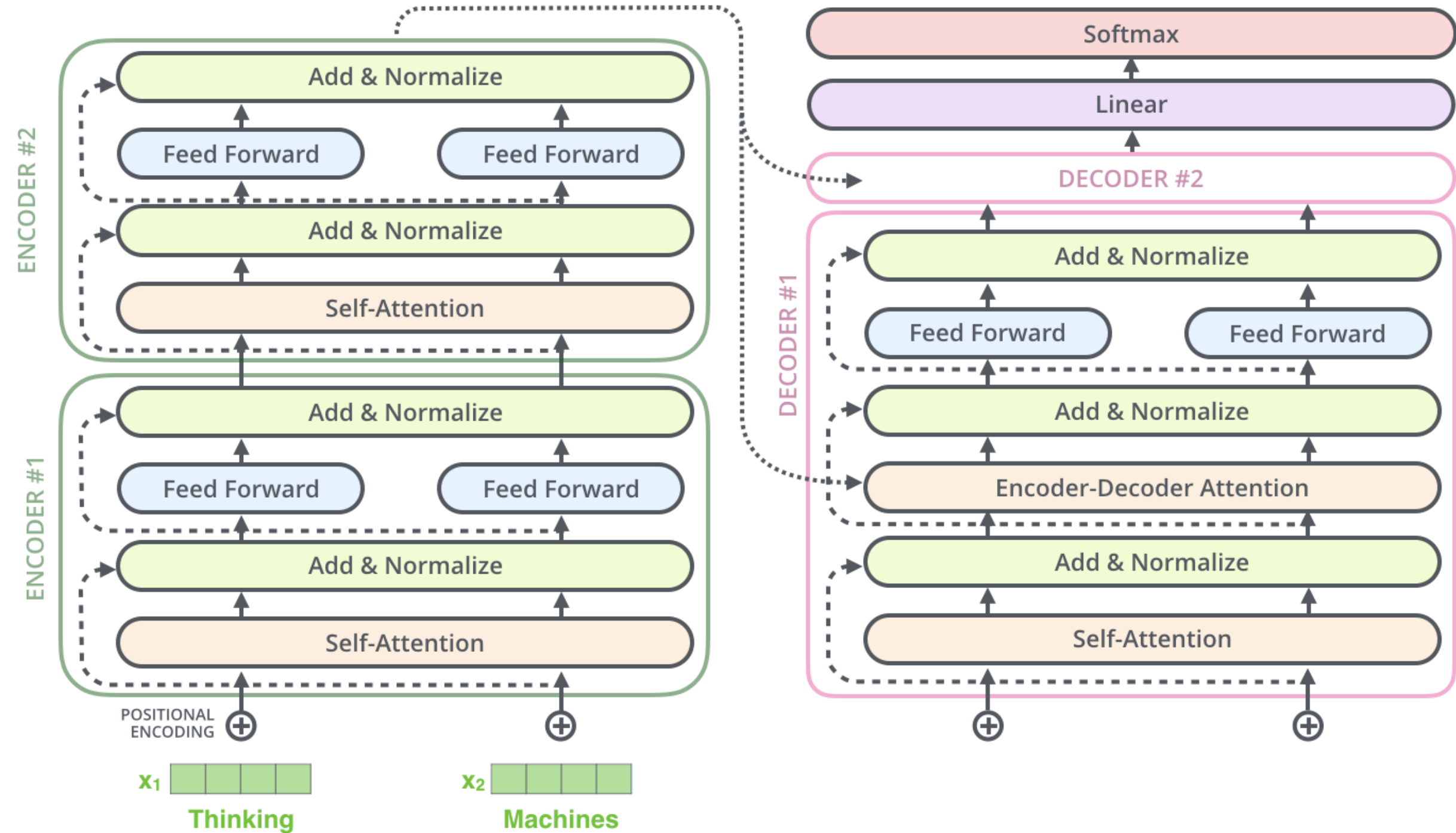
# The complete transformer



The encoder-decoder attention is just like self attention, except it uses K, V from the top of encoder output, and its own Q

Note: In decoder, the input is "incomplete" when calculating self-attention.

The solution is to set future unknown values with "-inf".

# Decoder's Output Linear Layer

# Training and the Loss Function

## Untrained Model Output

| 0.2 | 0.2 | 0.1 | 0.2 | 0.2 | 0.1 |
|-----|-----|-----|-----|-----|-----|

## Correct and desired output

| 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 |
|-----|-----|-----|-----|-----|-----|
| a   | am  | I   | thanks | student | <eos> |

We can use cross Entropy.

We can also optimize two words at a time: using BEAM search: keep a few alternatives for the first word.