

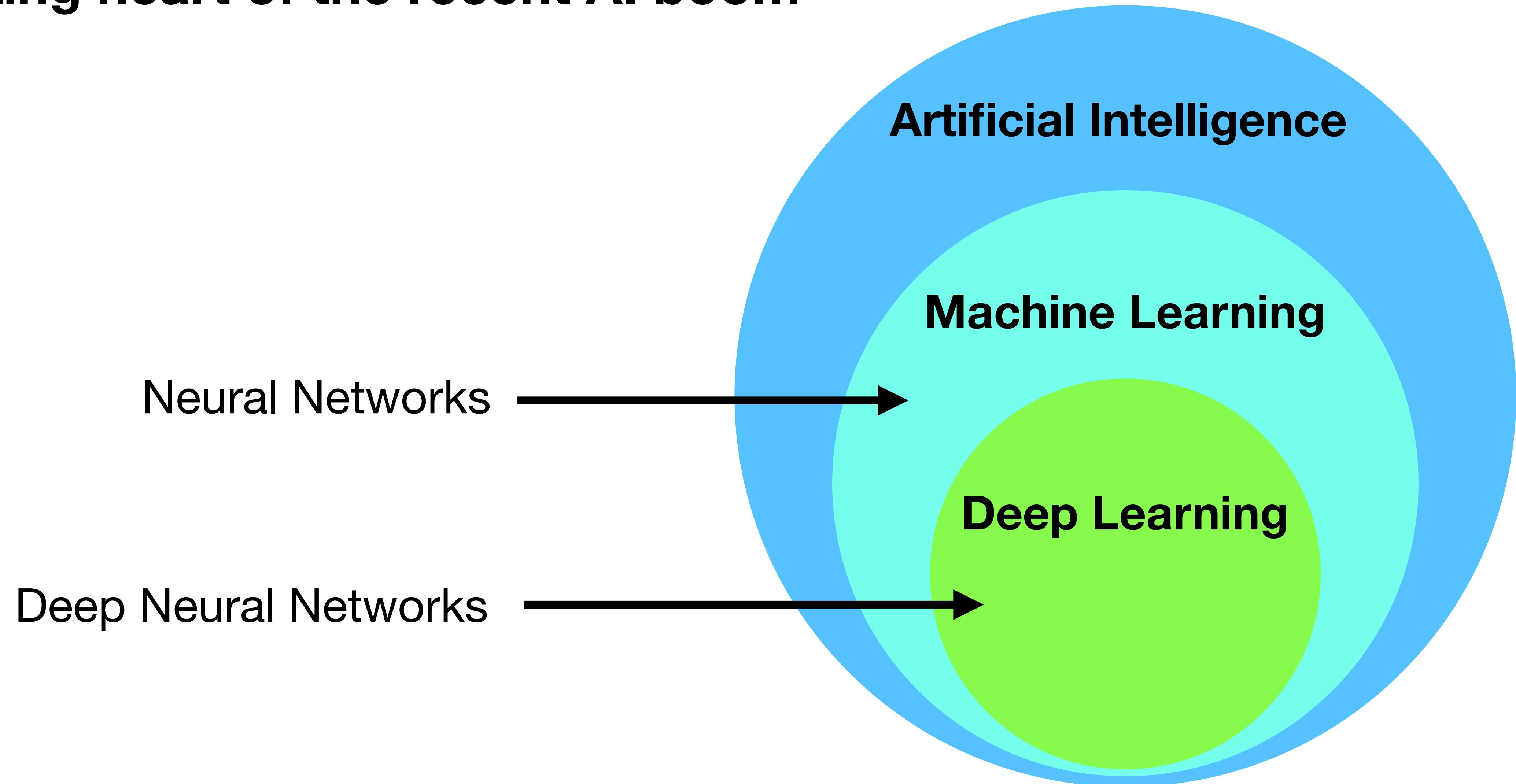
# **Accelerating Neural Networks with CUDA**

## **Multi-Core Programming**

**Parsa Toopchinezhad - Dr. Jahangiry - Spring 2024**

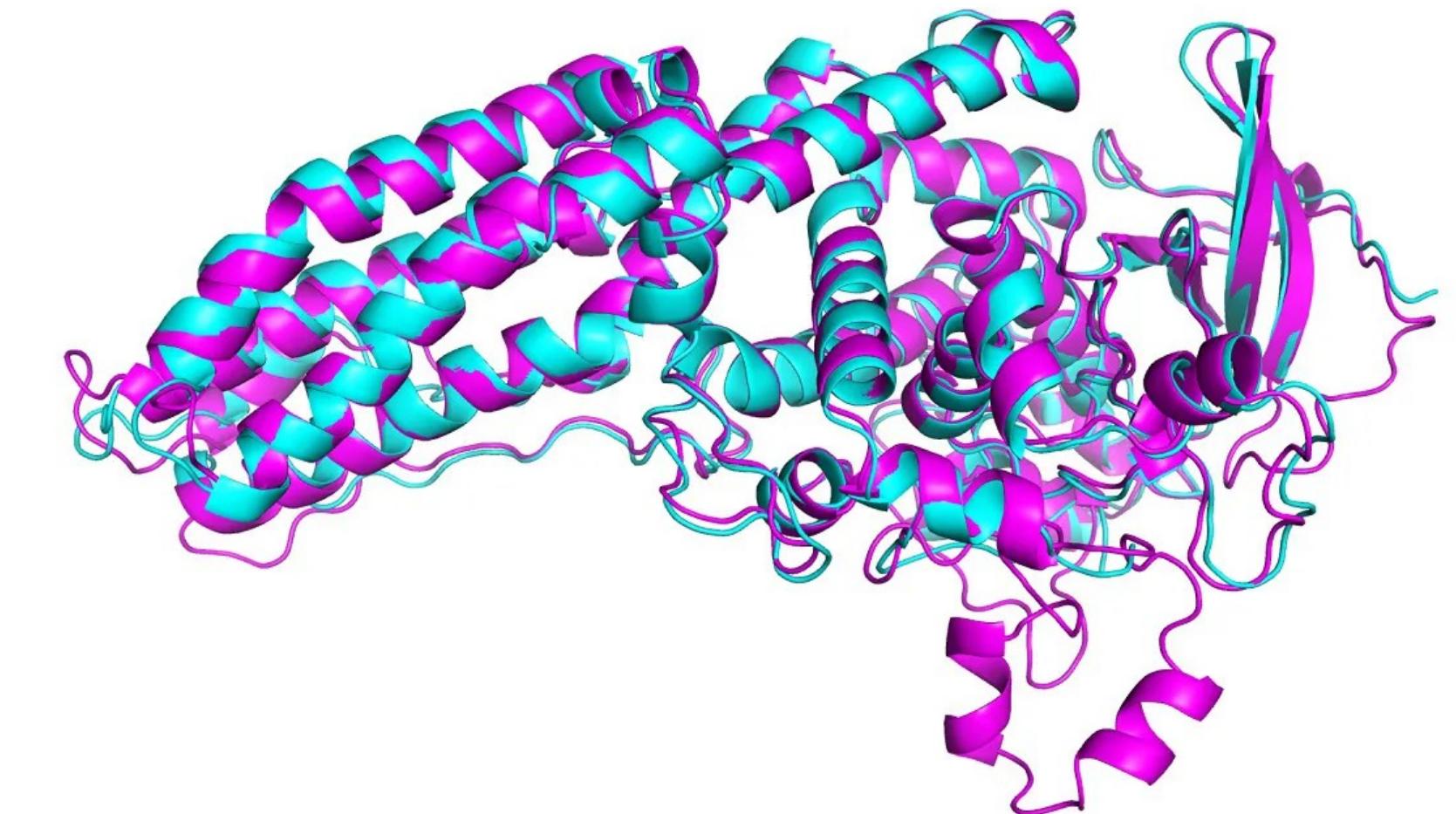
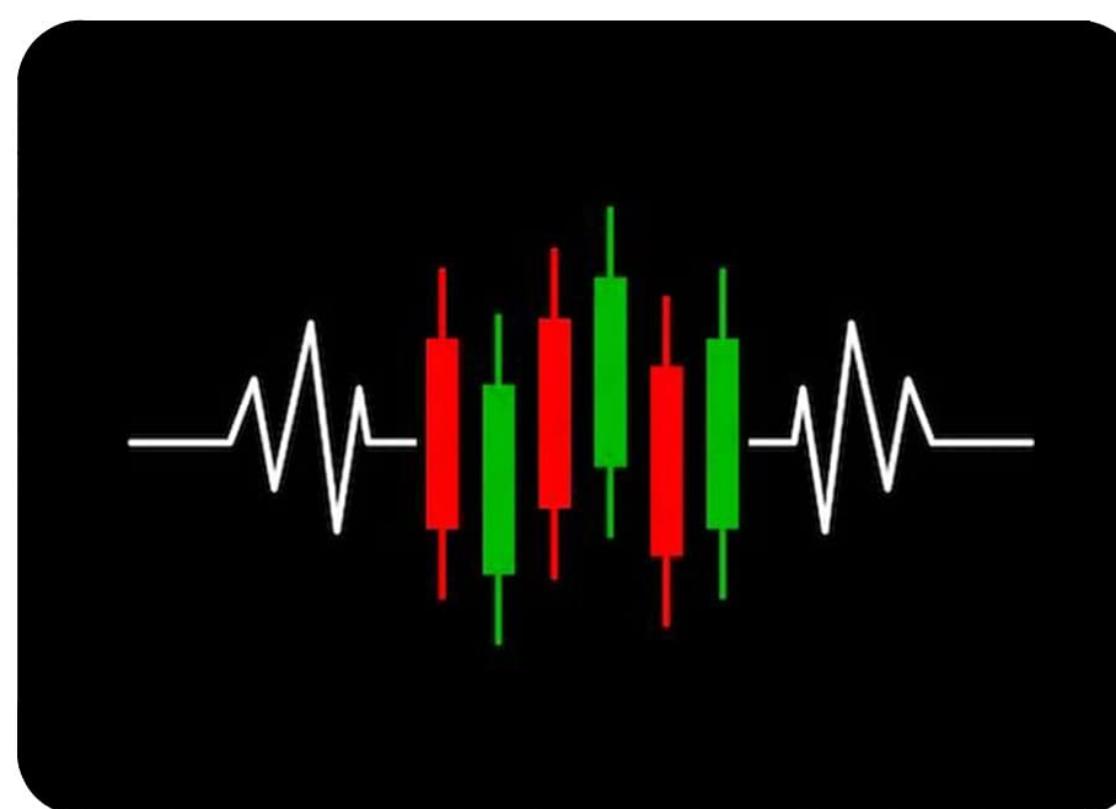
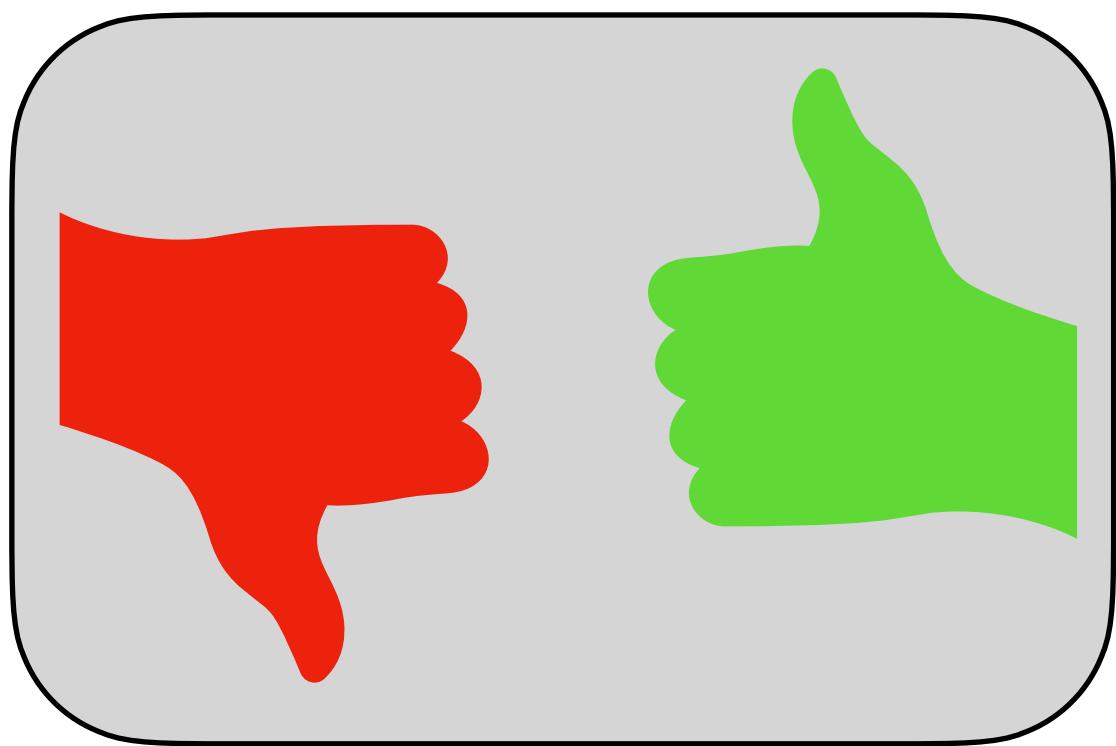
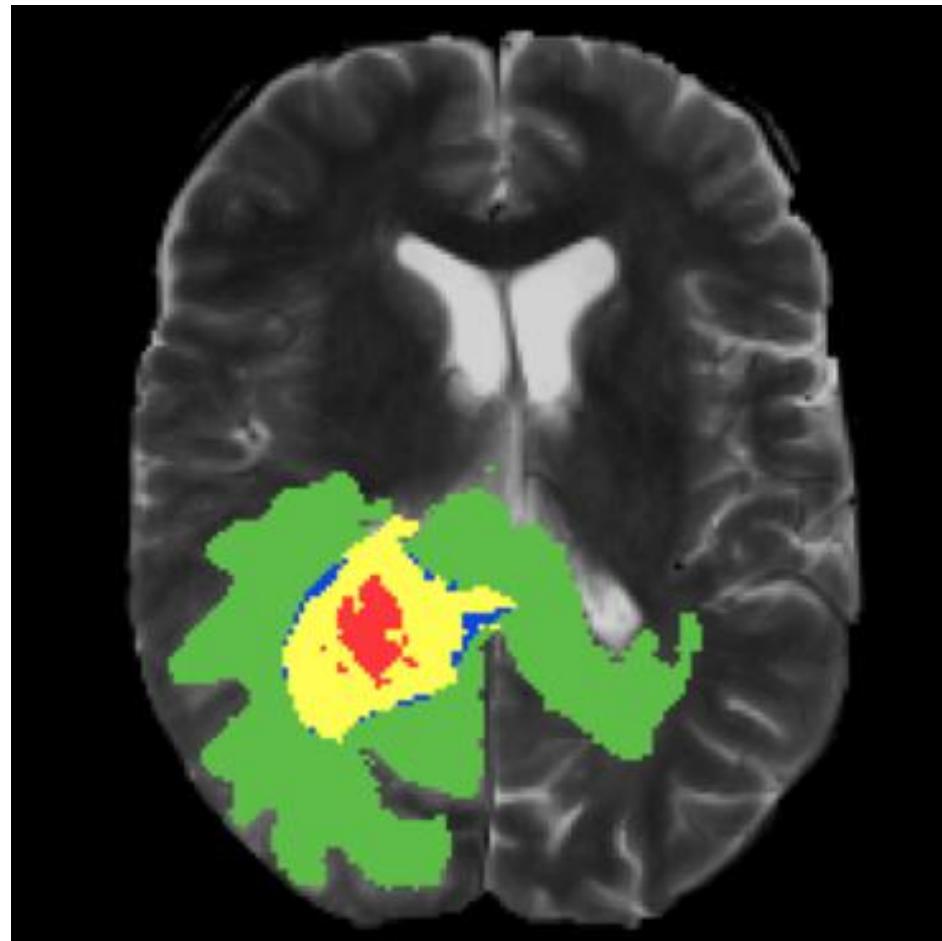
# What Are Neural Networks?

The beating heart of the recent AI boom



# Applications of Neural Networks

Dawn of an AI powered world



# Novel Applications of NNs

## Discovering faster matrix multiplication (AlphaTensor)

$$\begin{bmatrix} a_{1,1} & a_{1,2} & a_{1,3} & a_{1,4} & a_{1,5} \\ a_{2,1} & a_{2,2} & a_{2,3} & a_{2,4} & a_{2,5} \\ a_{3,1} & a_{3,2} & a_{3,3} & a_{3,4} & a_{3,5} \\ a_{4,1} & a_{4,2} & a_{4,3} & a_{4,4} & a_{4,5} \end{bmatrix} \times \begin{bmatrix} b_{1,1} & b_{1,2} & b_{1,3} & b_{1,4} & b_{1,5} \\ b_{2,1} & b_{2,2} & b_{2,3} & b_{2,4} & b_{2,5} \\ b_{3,1} & b_{3,2} & b_{3,3} & b_{3,4} & b_{3,5} \\ b_{4,1} & b_{4,2} & b_{4,3} & b_{4,4} & b_{4,5} \\ b_{5,1} & b_{5,2} & b_{5,3} & b_{5,4} & b_{5,5} \end{bmatrix} = \begin{bmatrix} c_{1,1} & c_{1,2} & c_{1,3} & c_{1,4} & c_{1,5} \\ c_{2,1} & c_{2,2} & c_{2,3} & c_{2,4} & c_{2,5} \\ c_{3,1} & c_{3,2} & c_{3,3} & c_{3,4} & c_{3,5} \\ c_{4,1} & c_{4,2} & c_{4,3} & c_{4,4} & c_{4,5} \end{bmatrix}$$

- Traditional: 100
- Strassen: 80
- AlphaTensor: 76



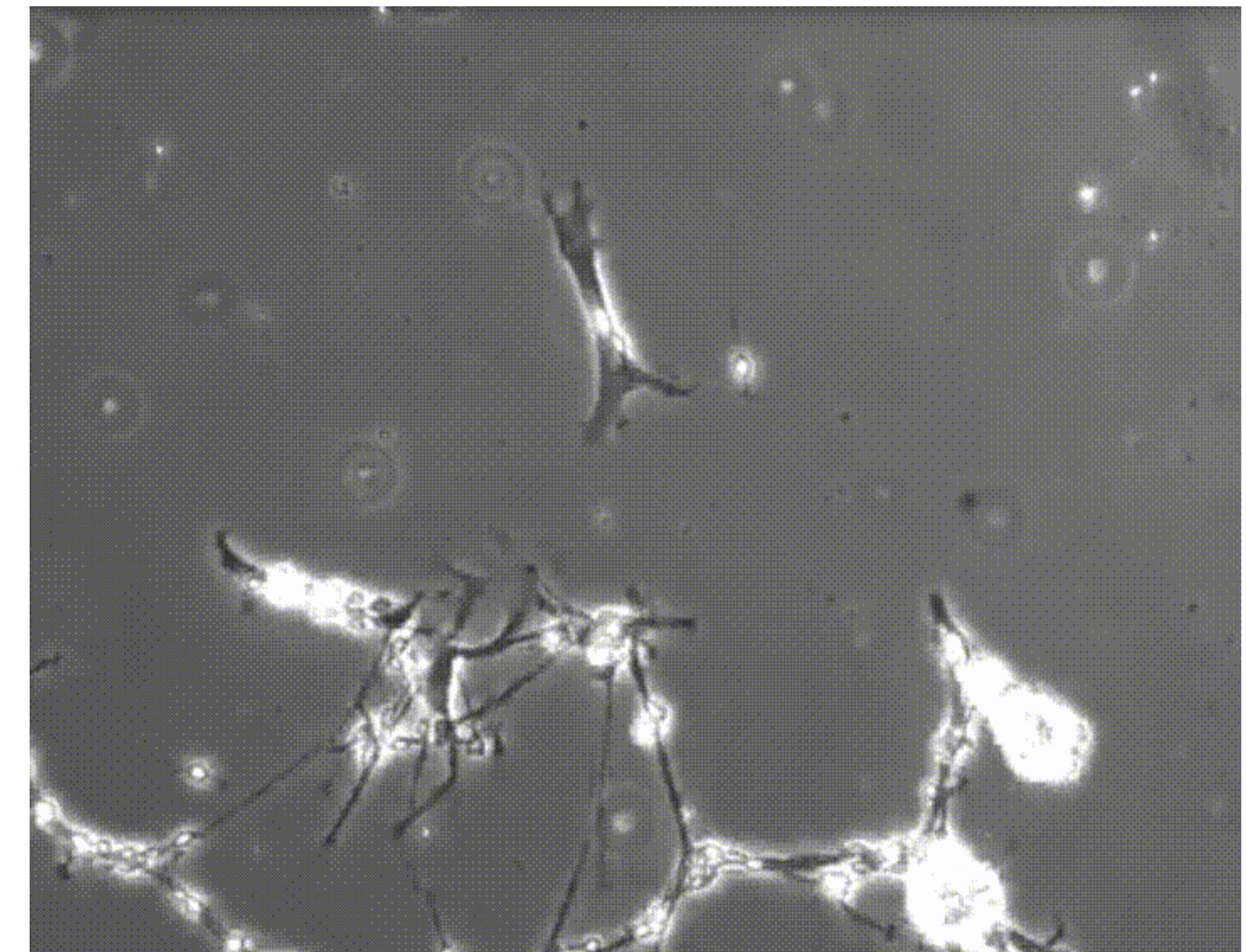
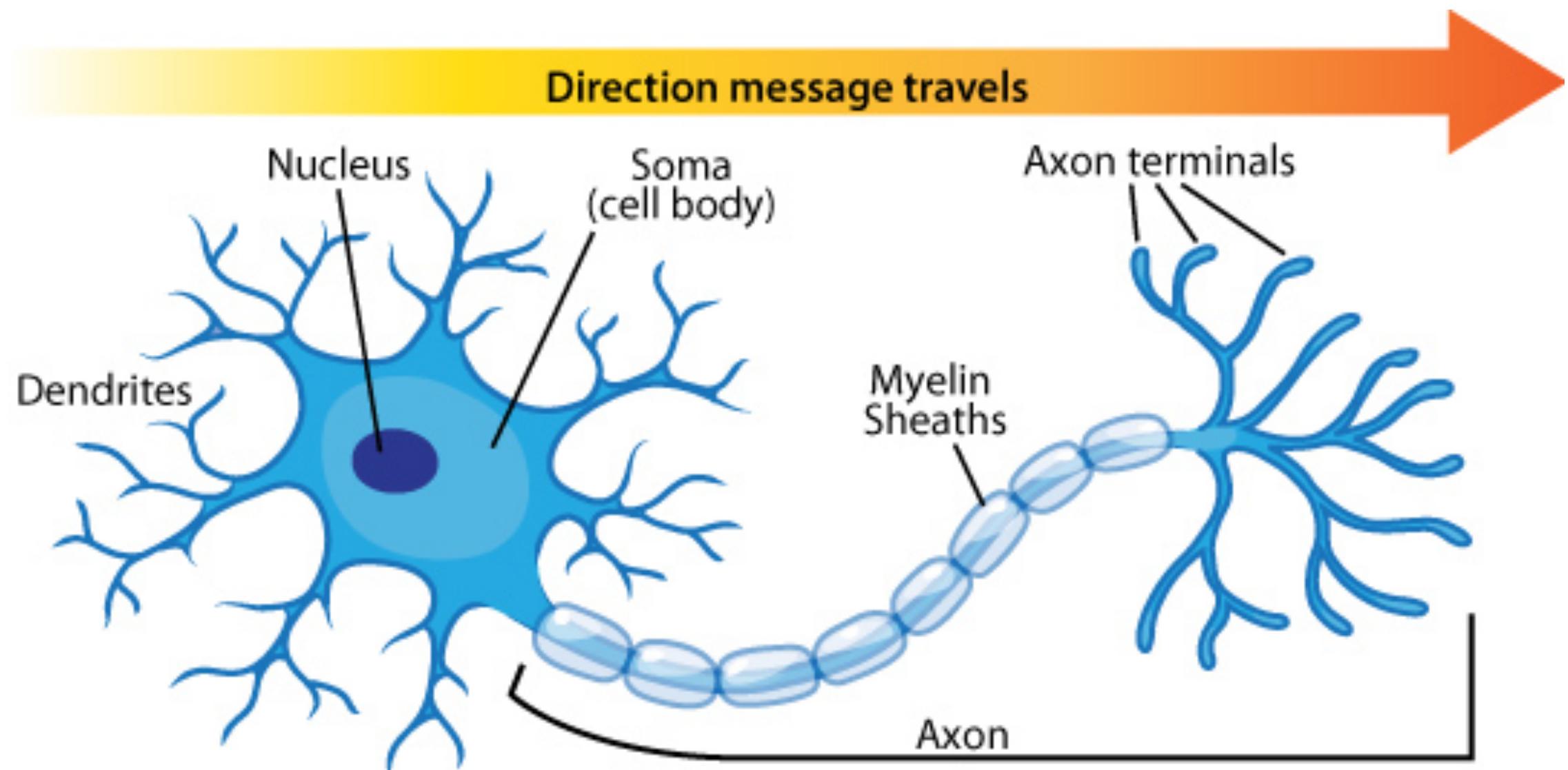
# Cool, but how do neural networks work?

(First we'll look at biological neural networks)

# How Does The Brain Work?

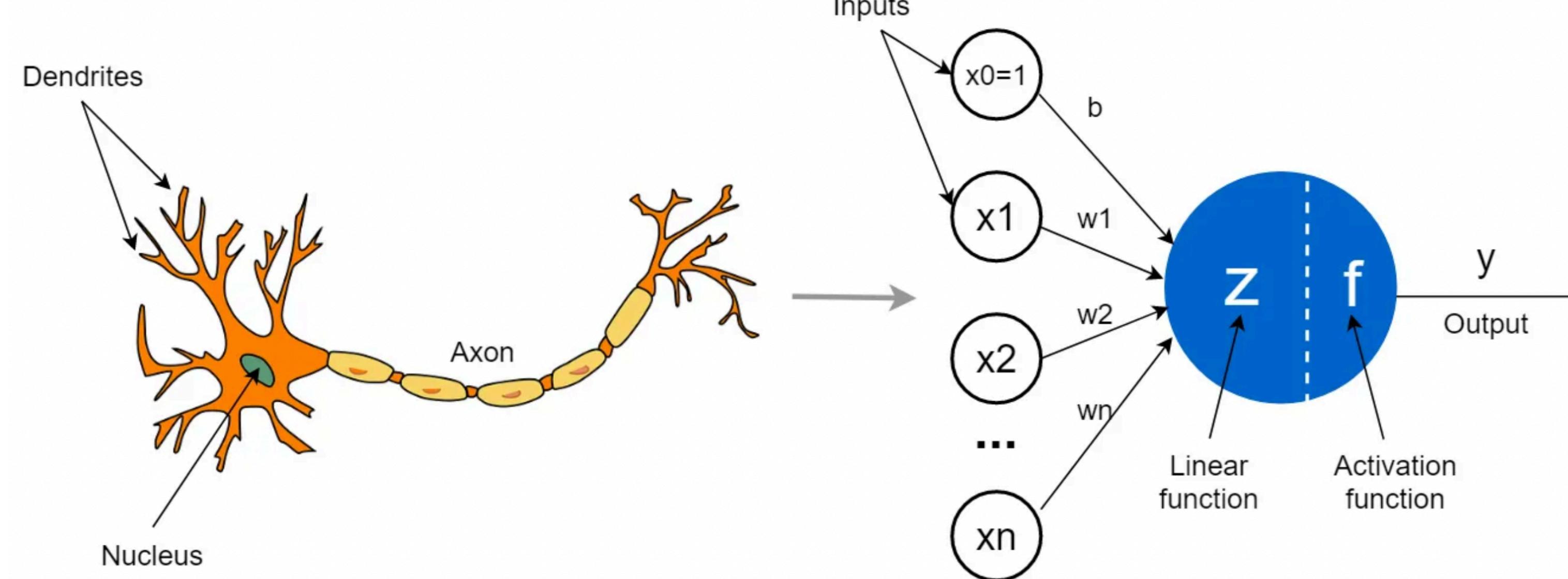
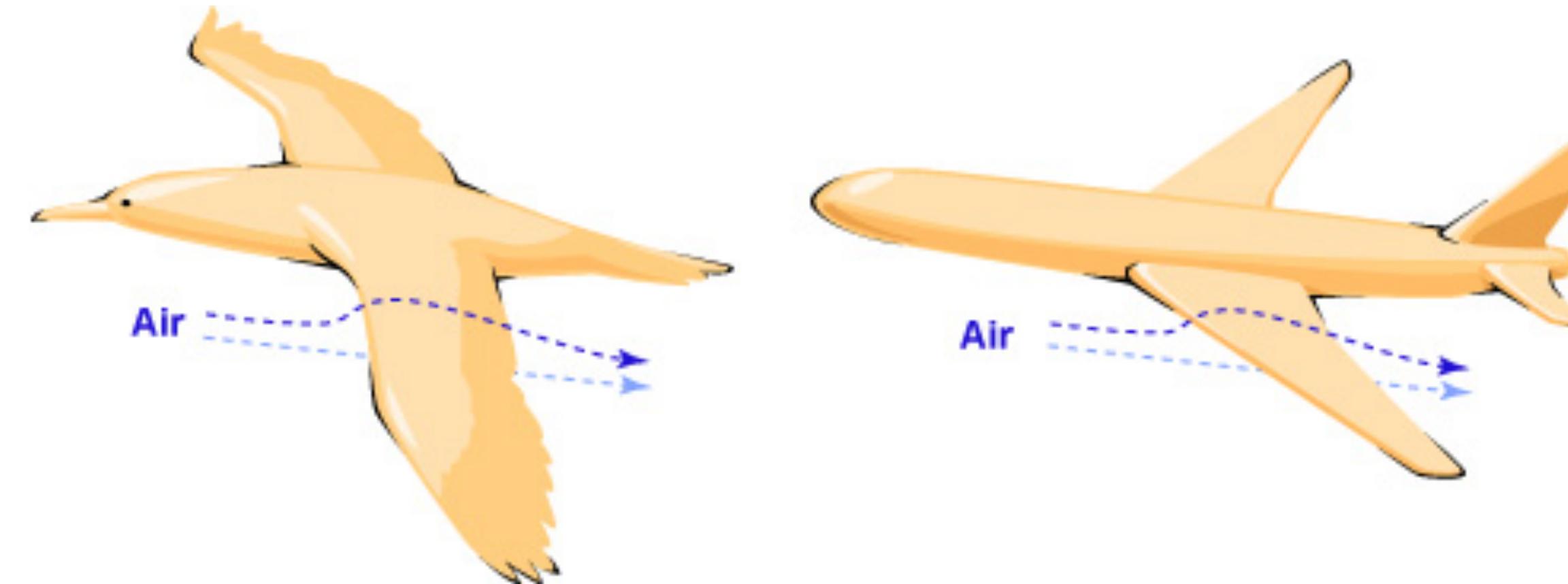
The least known part of our body

- “If the brain were so simple that we could understand it, we would be so simple that we couldn’t.” - Emerson M. Pugh
- Around 86 billion neurons



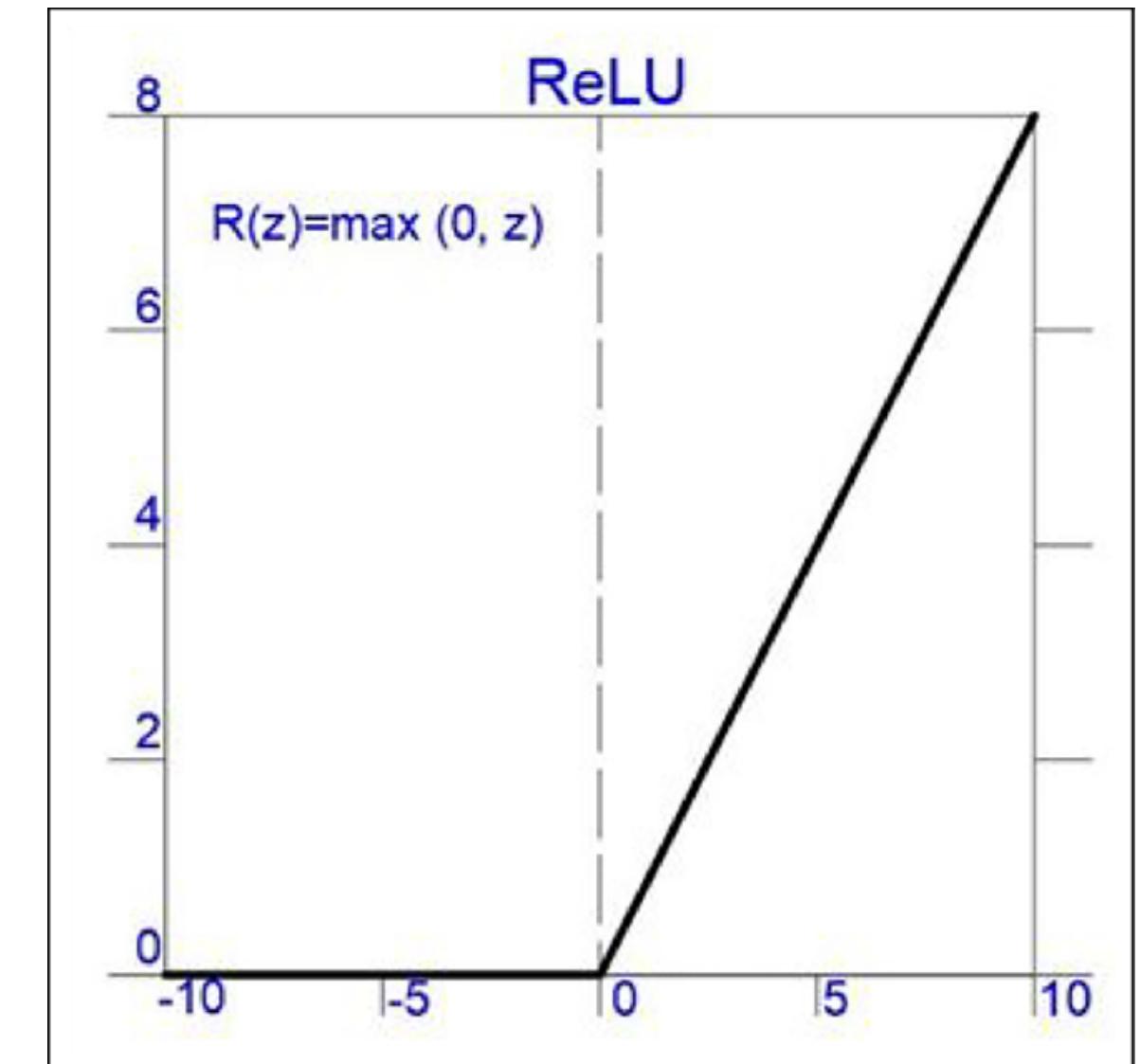
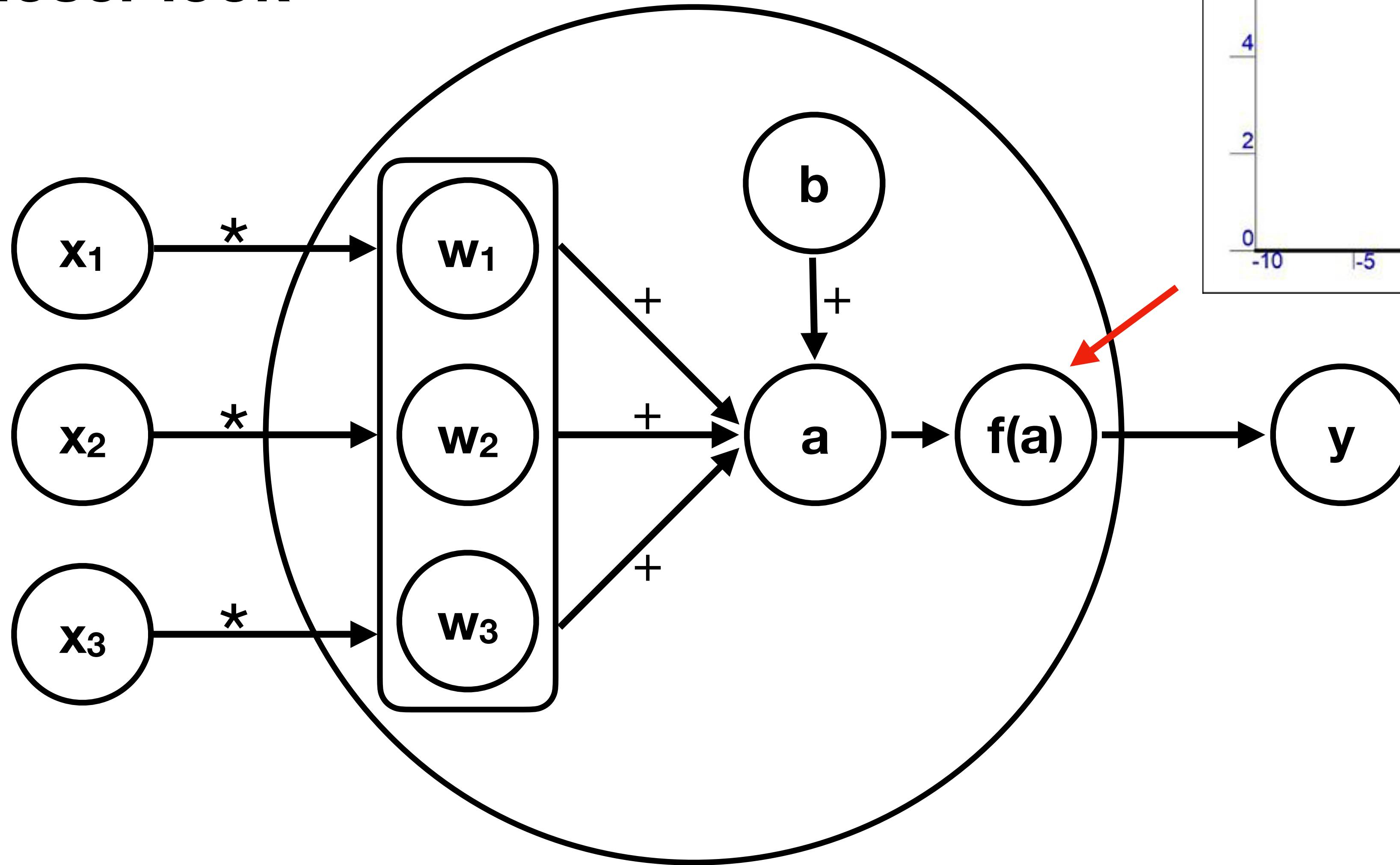
# Modelling Nature

Getting inspiration from the world around us



# Artificial Neuron

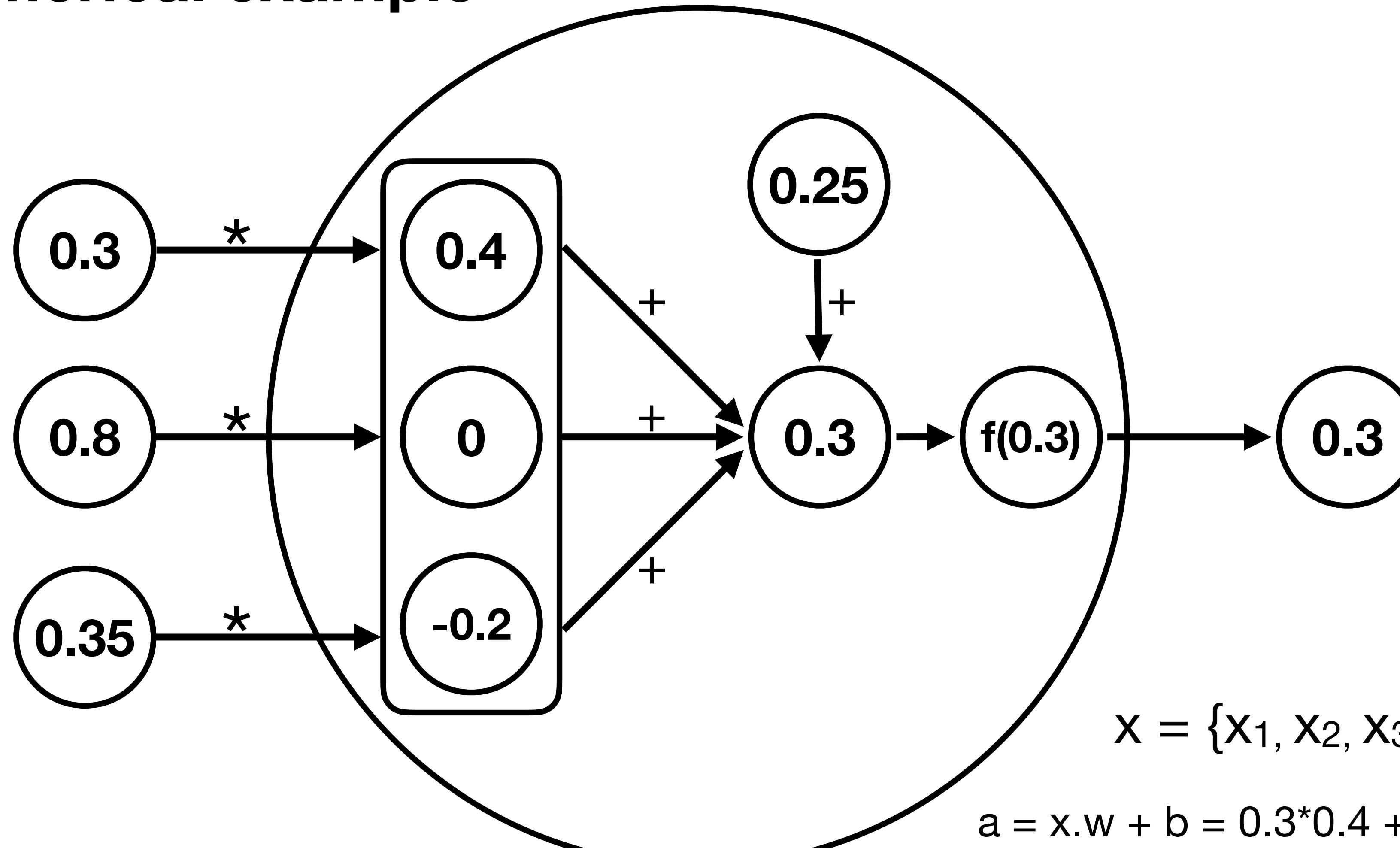
A closer look



Activation Function

# Artificial Neuron

## Numerical example



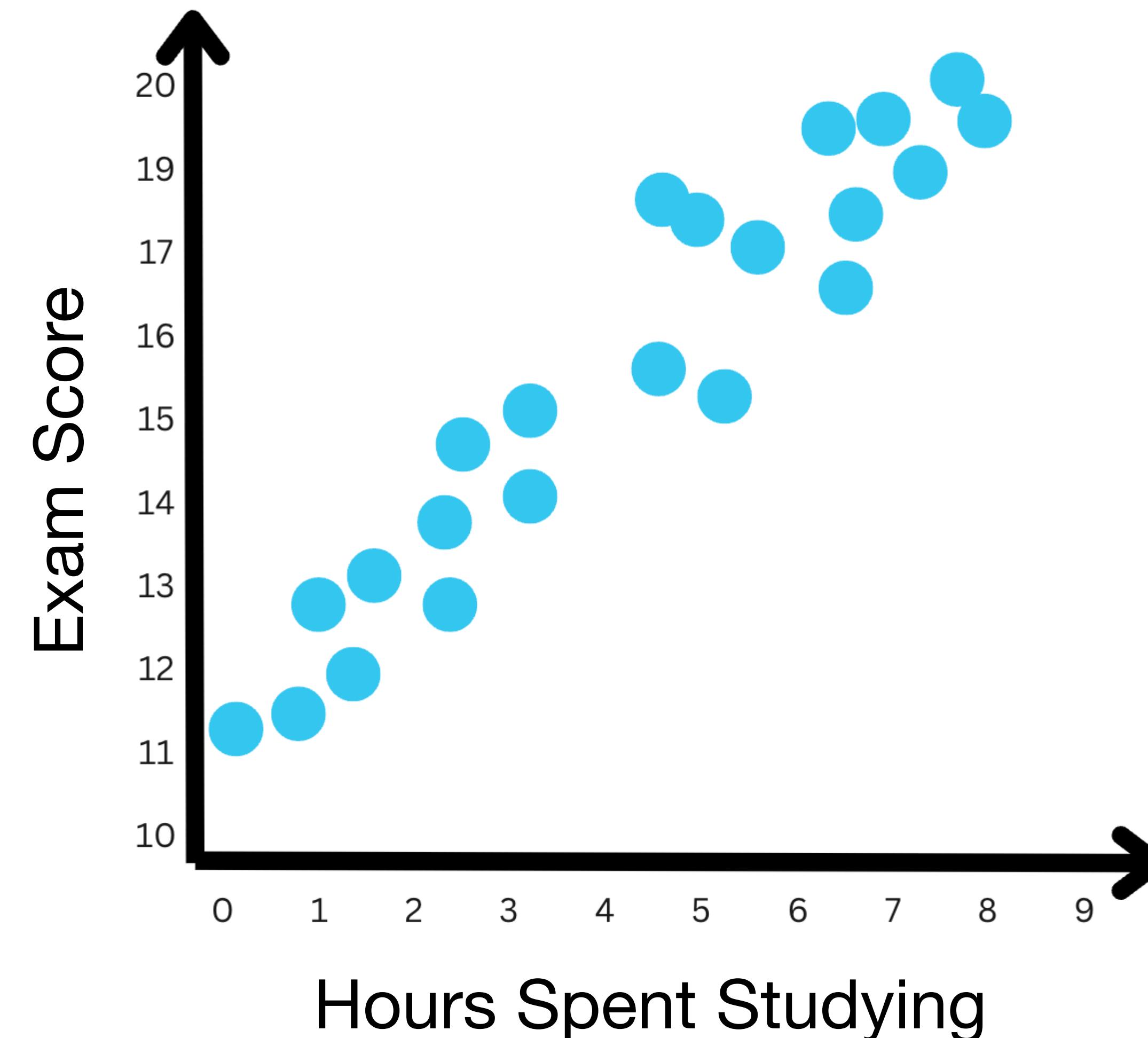
$$x = \{x_1, x_2, x_3\} \quad w = \{w_1, w_2, w_3\}$$

$$a = x \cdot w + b = 0.3 \cdot 0.4 + 0.8 \cdot 0 + 0.35 \cdot -0.2 + 0.25 = 0.3$$

$$y = f(a) = \text{ReLU}(0.3) = 0.3$$

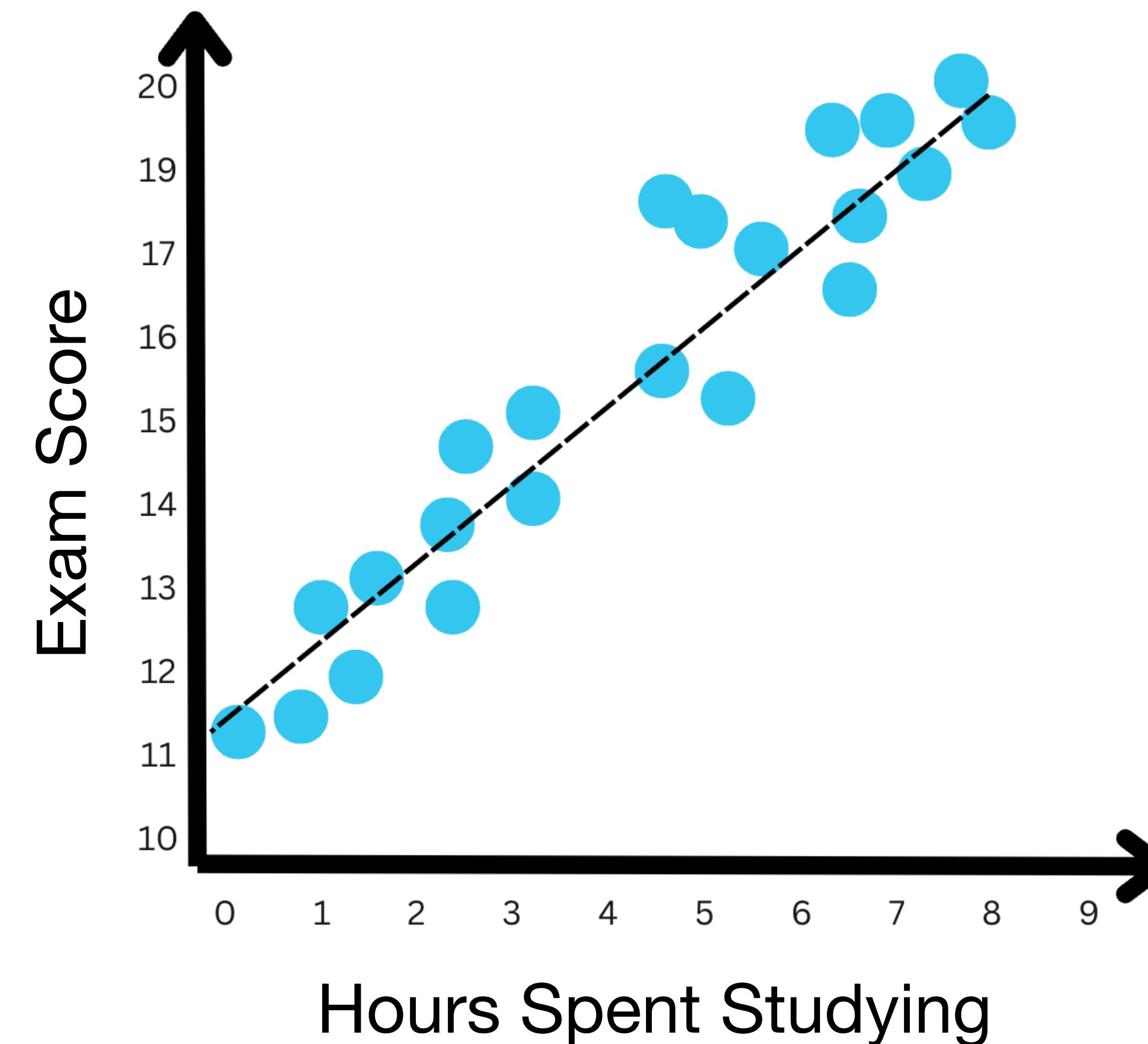
# What Can a Perceptron Do?

Simple yet powerful



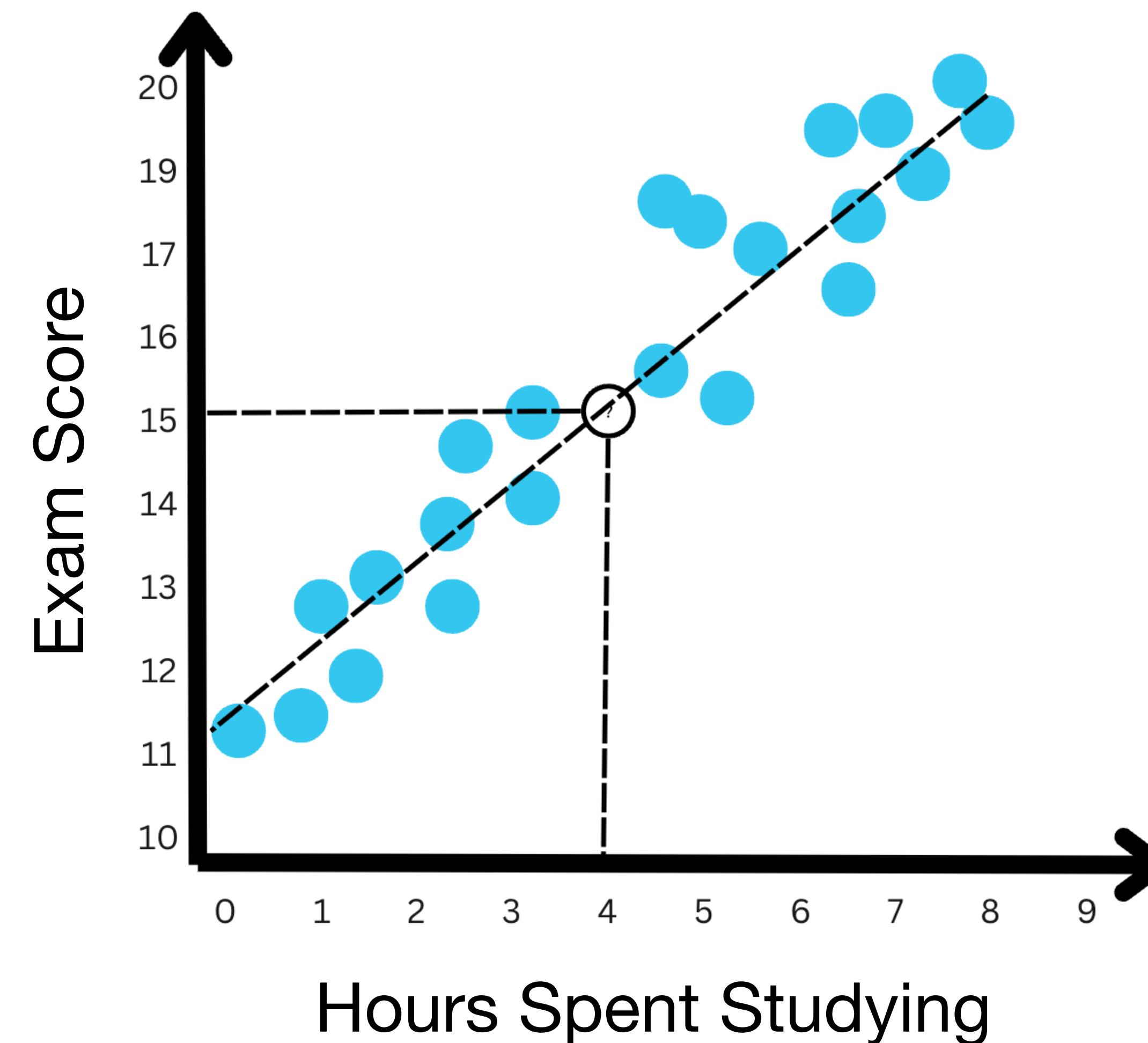
# What Can a Perceptron Do?

Simple yet powerful



# What Can a Perceptron Do?

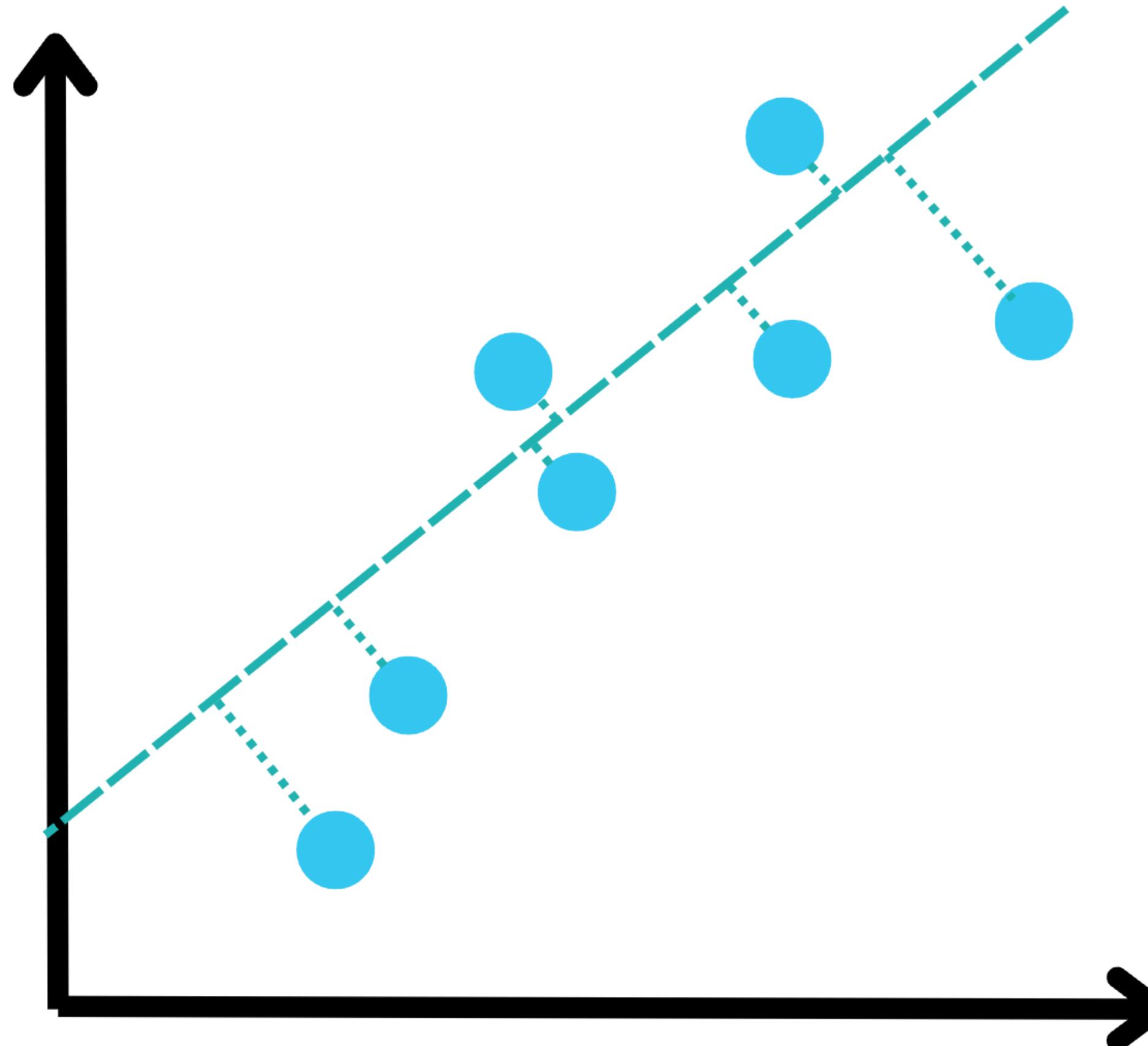
Simple yet powerful



**Goal:** find a line that  
is close to all points

# Calculating the Error

Formulating it to as an optimization problem



$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

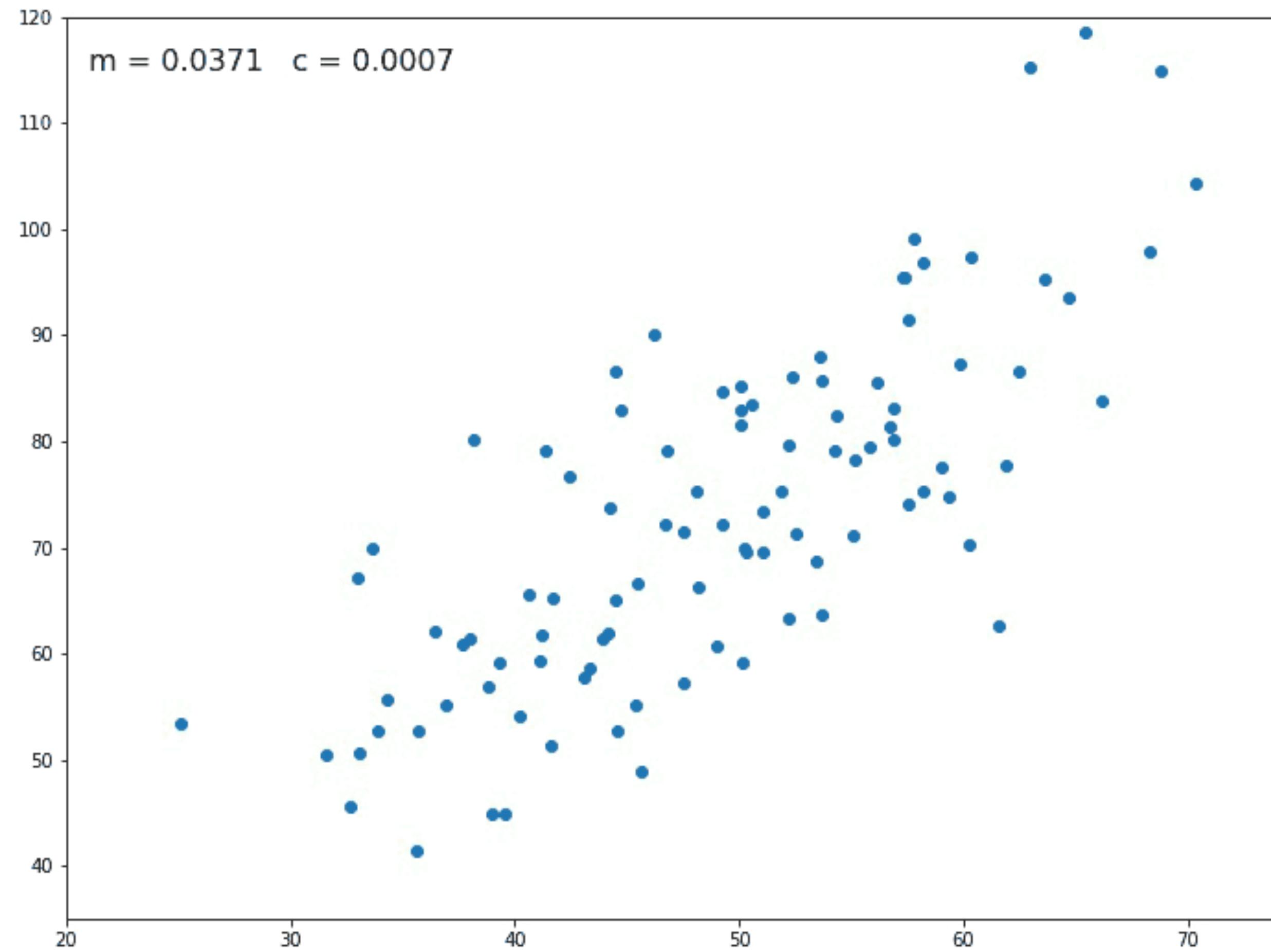
میانگین      فاصله نقطه تا خط

**Goal:** find  $w$  and  $b$  such that they minimize the MSE

(These values are found via calculating and stepping in the opposite side of the derivatives, we wont focus on this today)

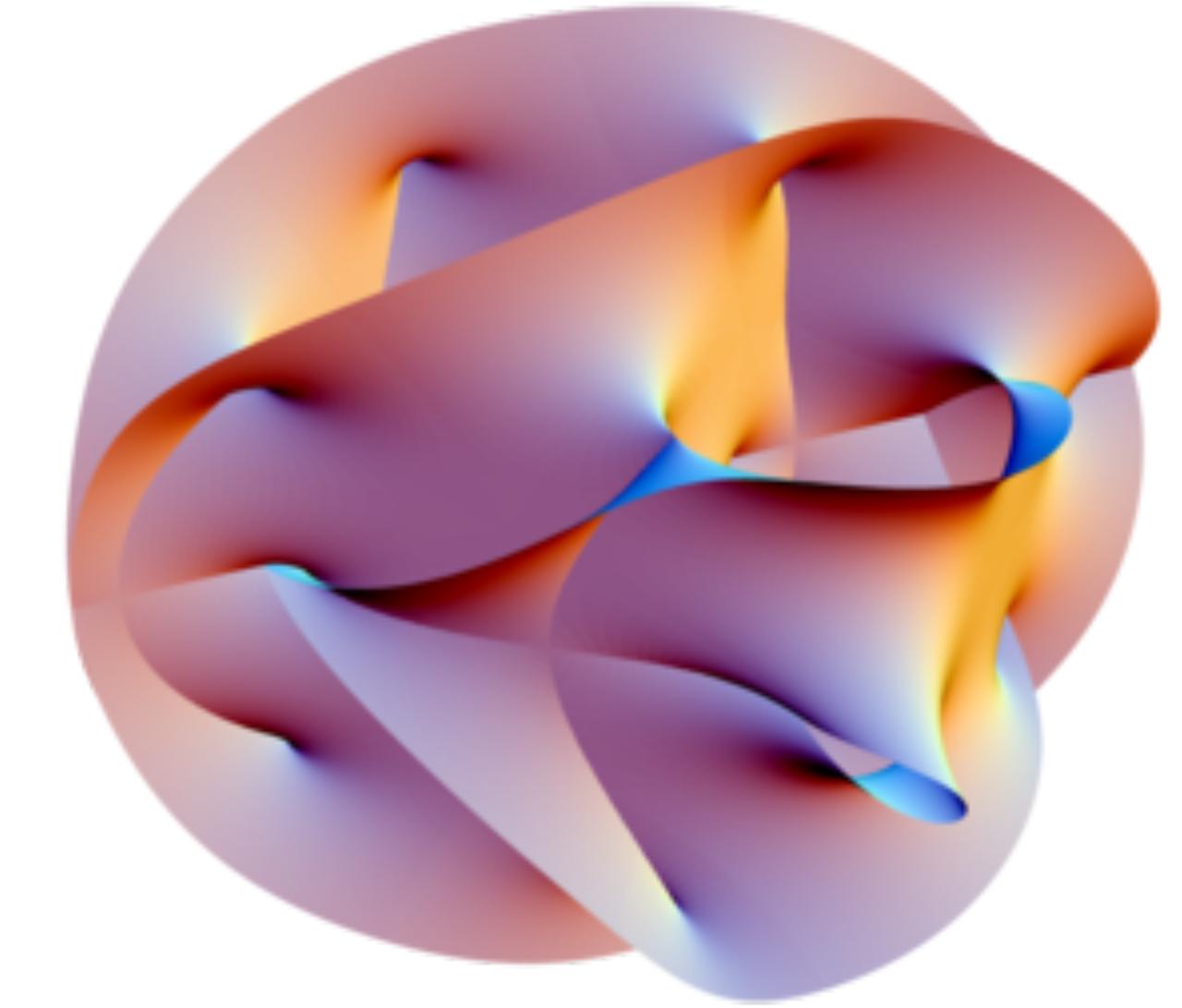
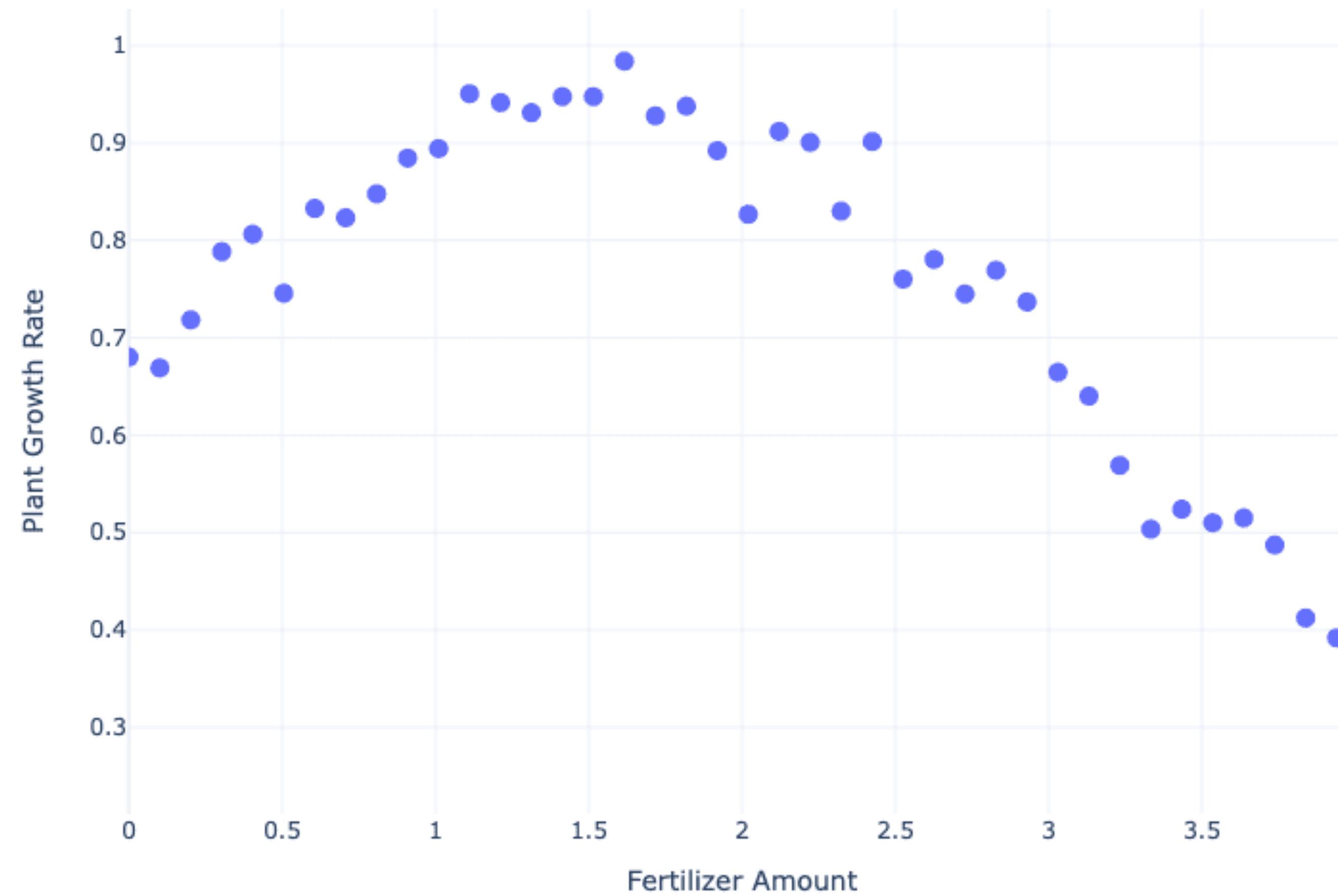
# Visualizing the Learning Process

## A perceptron in action



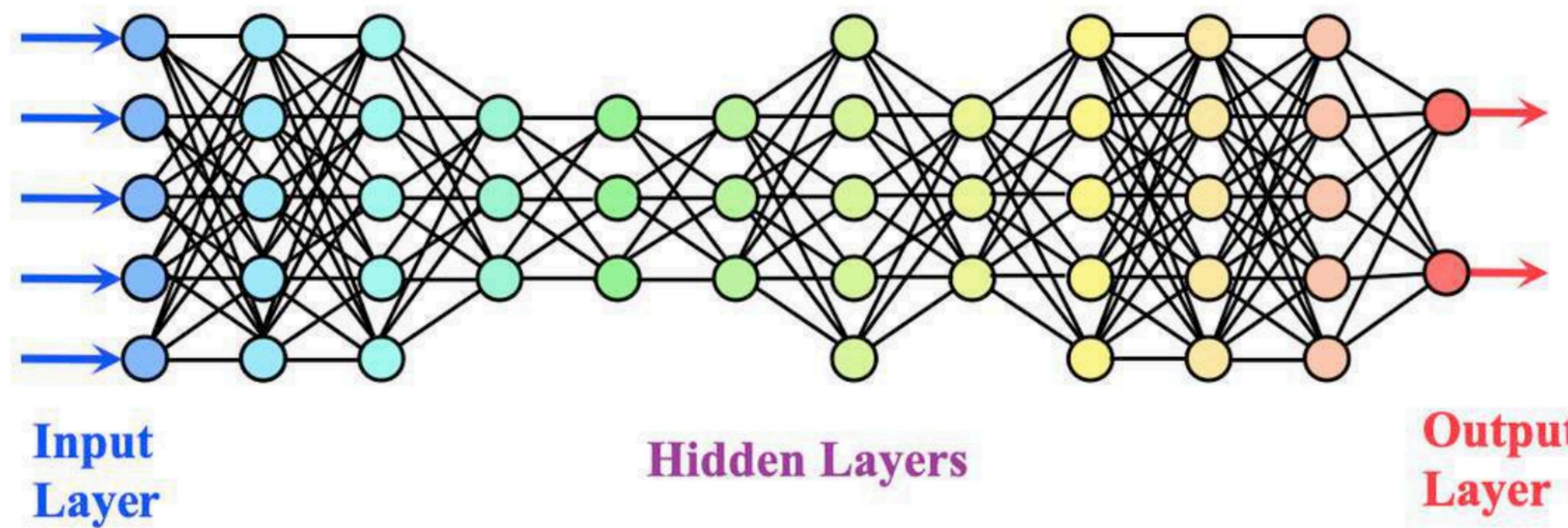
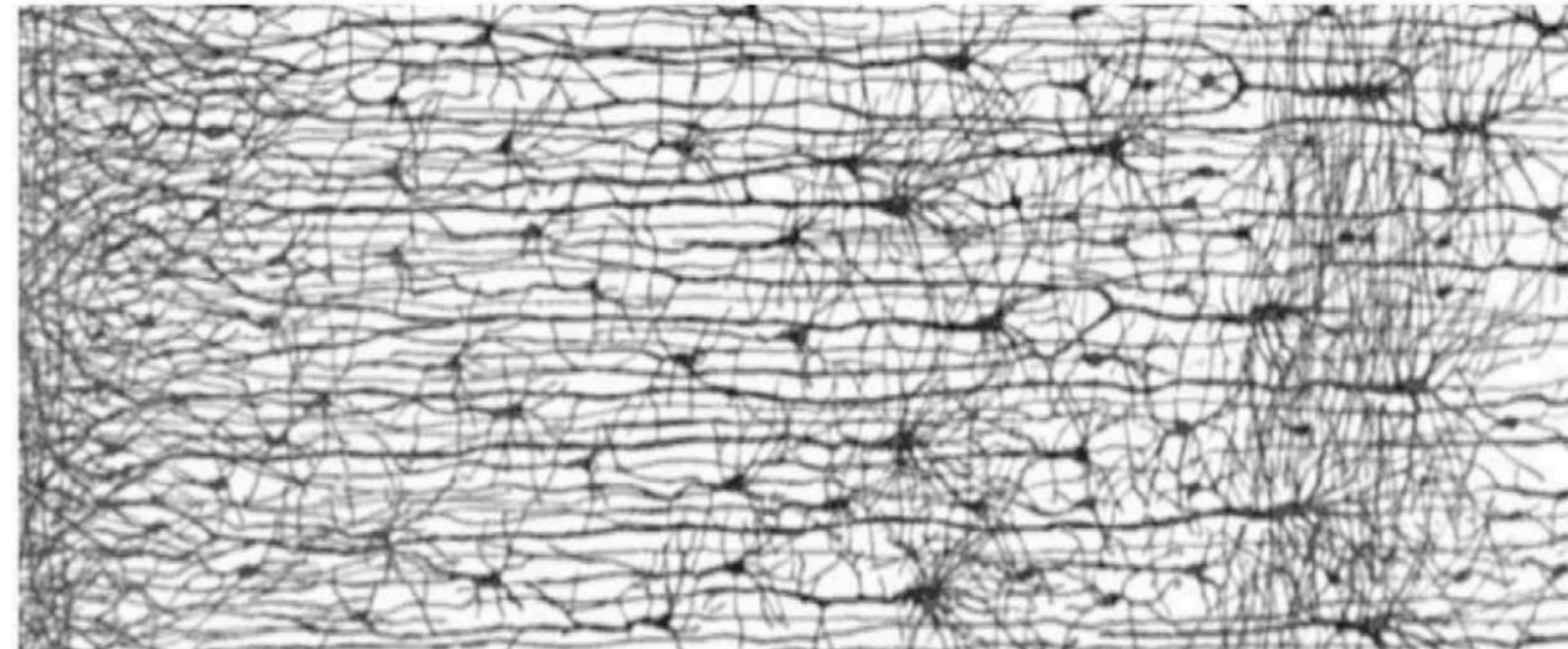
# Limits of Single Perceptrons

We live in a high dimensional, non-linear world



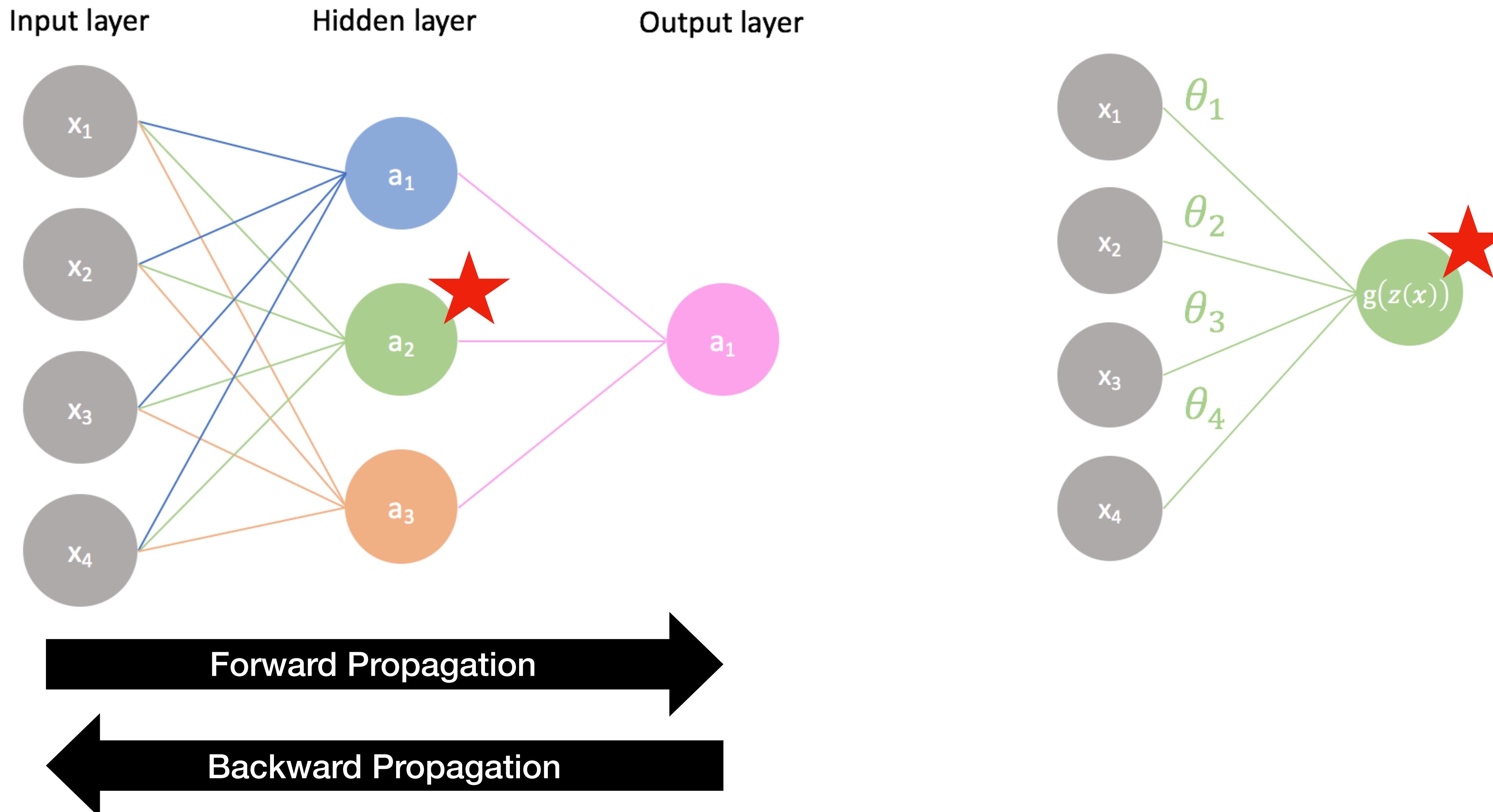
# Neural Networks

From perceptrons to layers to networks



# Neural Networks

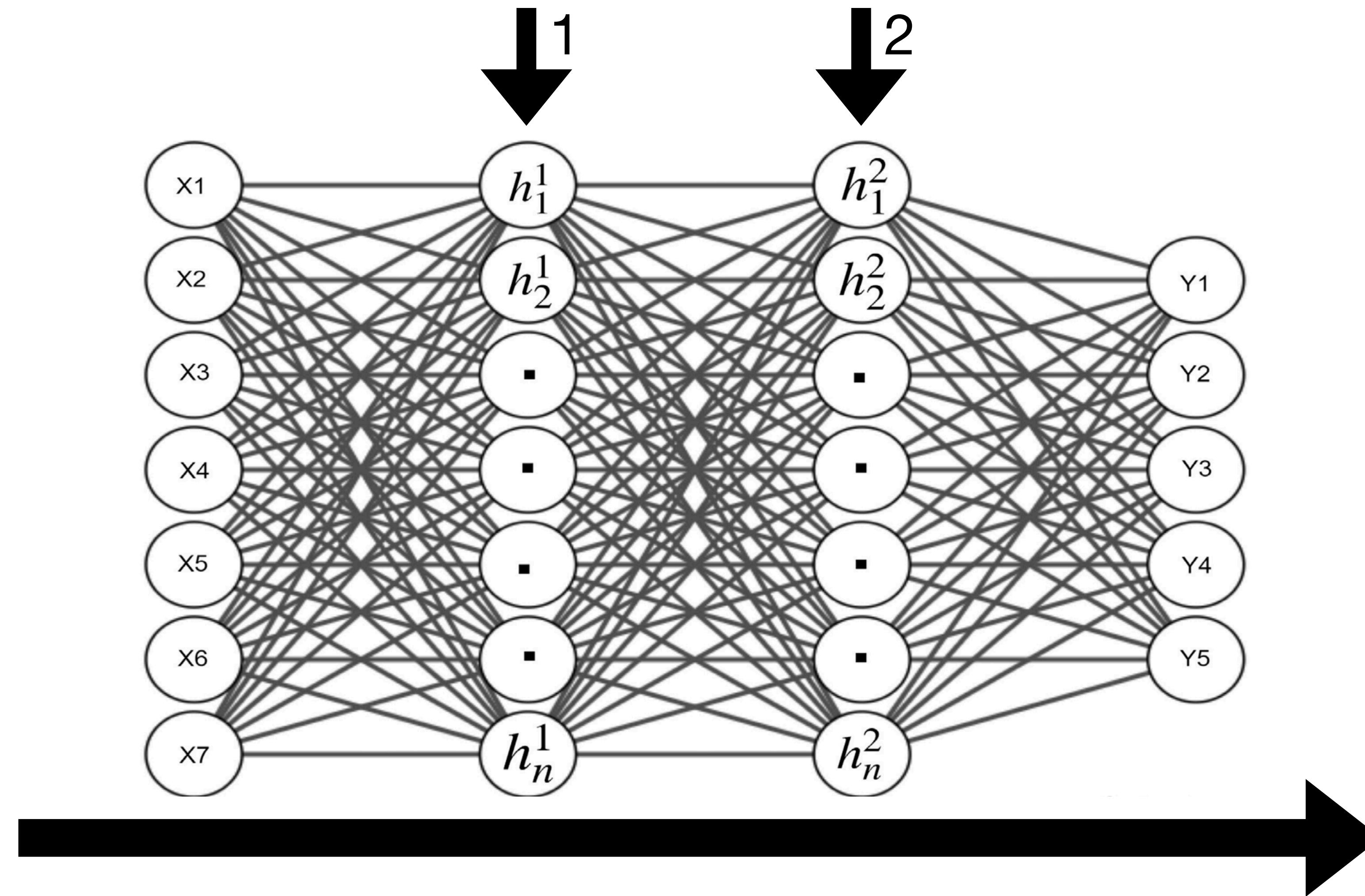
## A closer look



**How can we parallelize neural  
networks?**

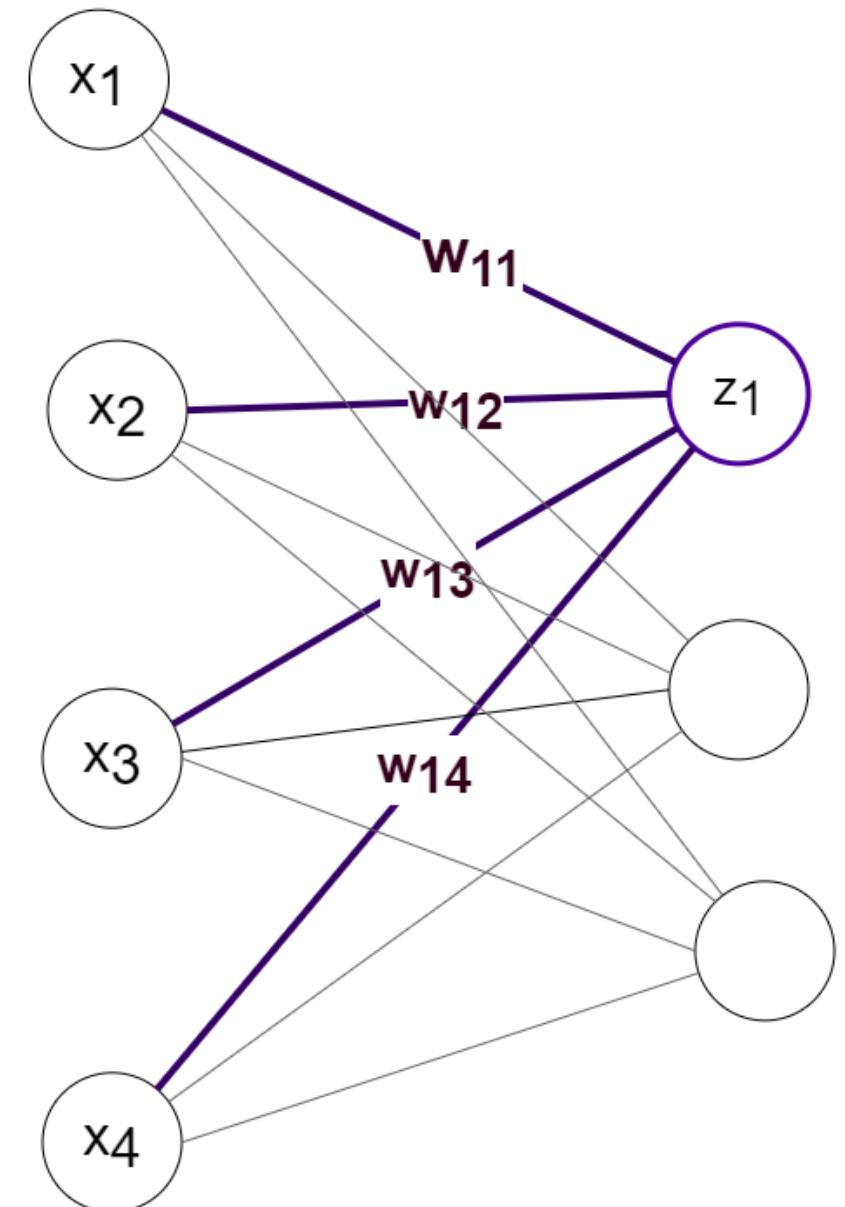
# What to Parallelize?

A combination of sequential and parallel operations



# A Naive Parallel Implementation

Thread 1



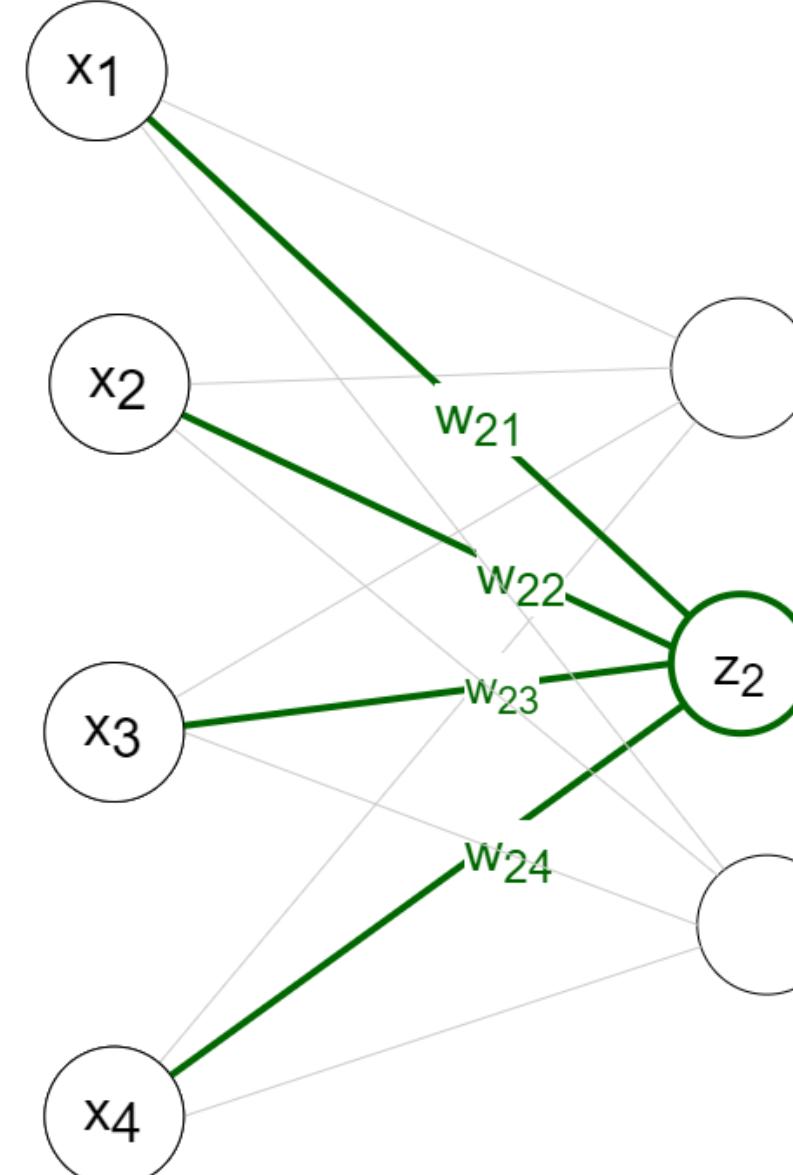
$$z_1 = w_{11} * x_1 + w_{12} * x_2 + w_{13} * x_3 + w_{14} * x_4 + b$$

Calculate Activation Function



$$a_1 = \text{sigmoid}(z_1)$$

Thread 2



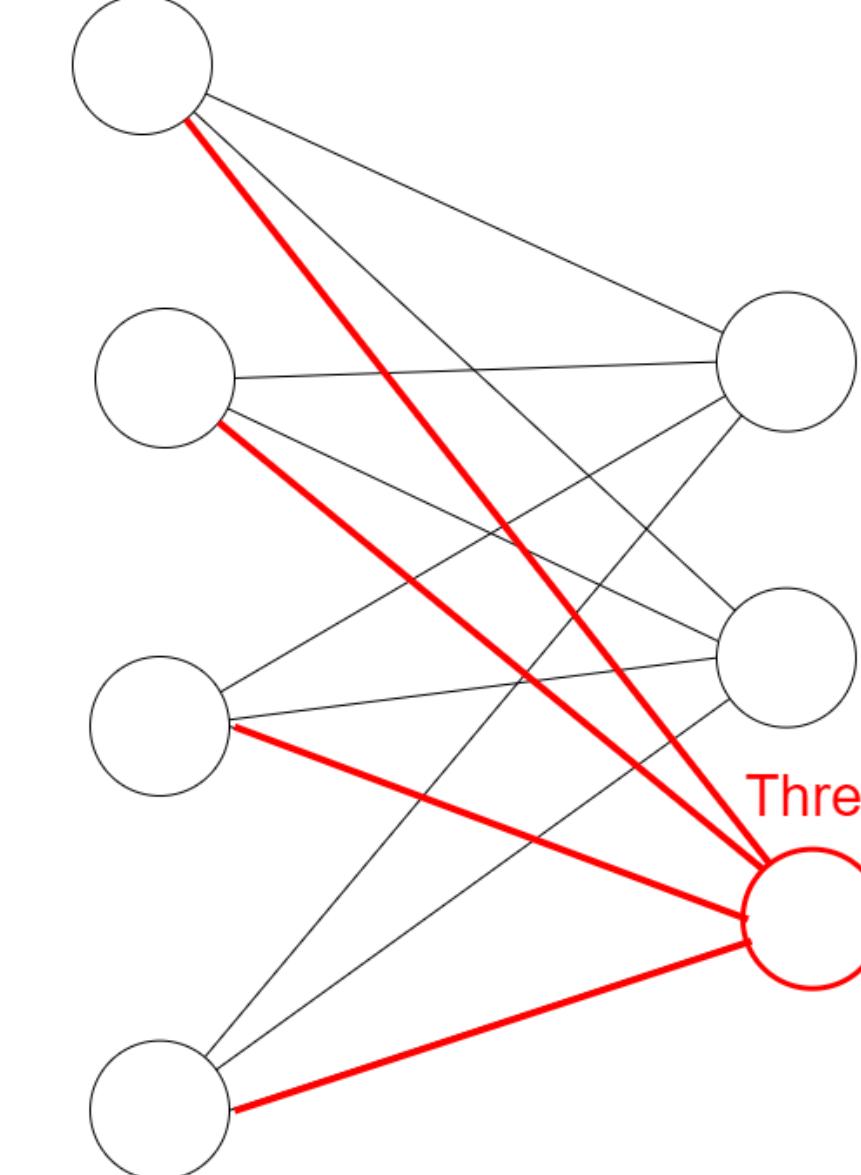
$$z_2 = w_{21} * x_1 + w_{22} * x_2 + w_{23} * x_3 + w_{24} * x_4 + b$$

Calculate Activation Function



$$a_2 = \text{sigmoid}(z_2)$$

Thread 3



$$z_3 = w_{31} * x_1 + w_{32} * x_2 + w_{33} * x_3 + w_{34} * x_4 + b$$

Calculate Activation Function



$$a_3 = \text{sigmoid}(z_3)$$

# A Naive Parallel Implementation

## Kernel code

```
__global__ void linear_layer_and_activation(float *weight_matrix, float *biases, float *x_inputs,
                                             float *z_values, float *activation_values,
                                             int nr_output_neurons, int nr_input_neurons)
{
    // We use the thread id so we can index different weights/neurons in each thread
    int id = threadIdx.x;

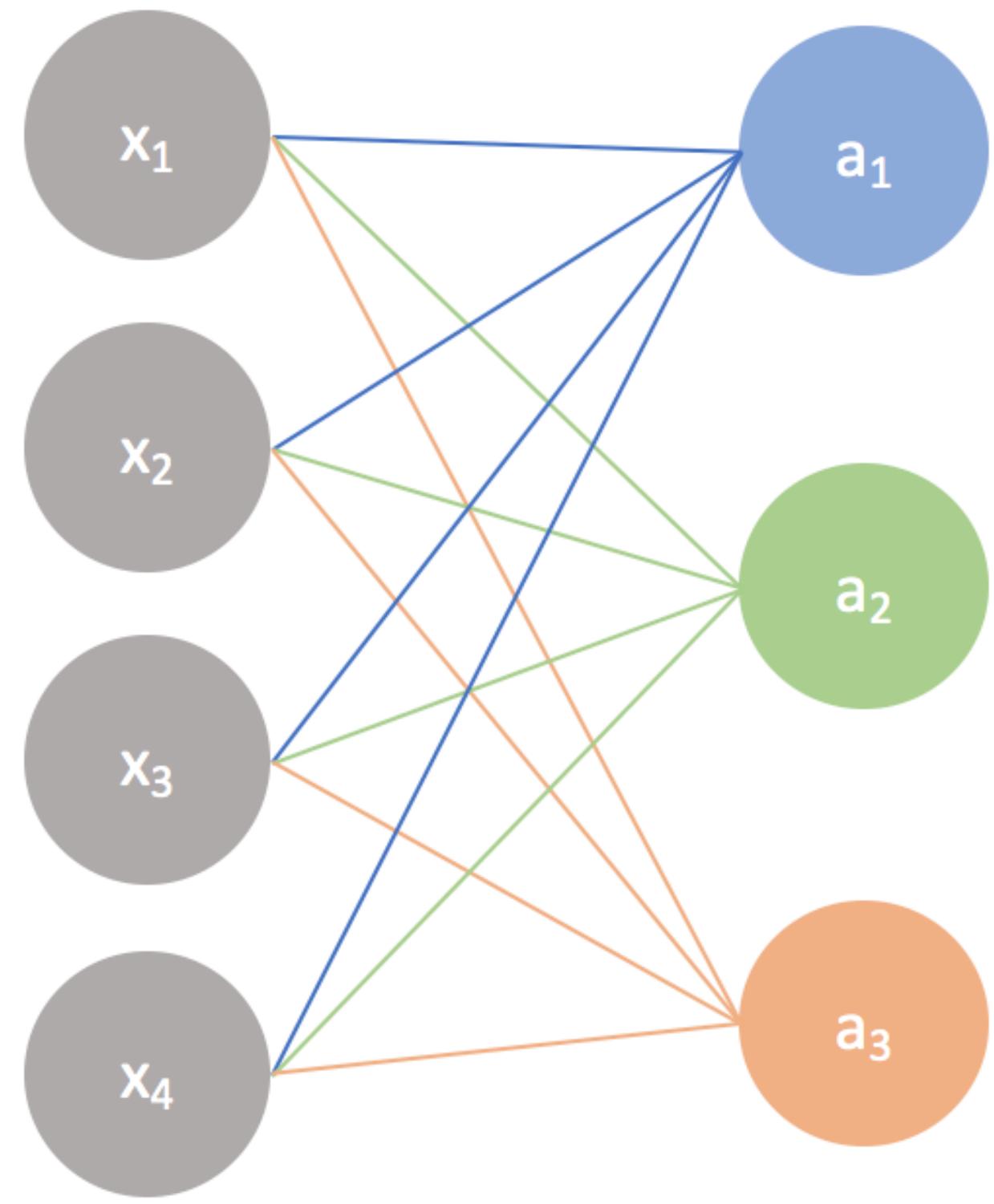
    // w*x
    // We loop over every incoming neuron and multiply it with it's weight and sum it to the current z_value
    for (int neuron_nr = 0; neuron_nr < nr_input_neurons; neuron_nr++)
    {
        z_values[id] += weight_matrix[(nr_input_neurons)* id + neuron_nr] * x_inputs[neuron_nr];
    }

    // w*x + b
    // Don't forget to add the bias
    z_values[id] += biases[id];

    // Then we compute the activation value just like in part1 of the tutorial
    // sig(w*x + b)
    activation_values[id] = 1.0 / (1.0 + exp(-z_values[id]));
}
```

# Neural Networks

## Forward propagation as matrix multiplication

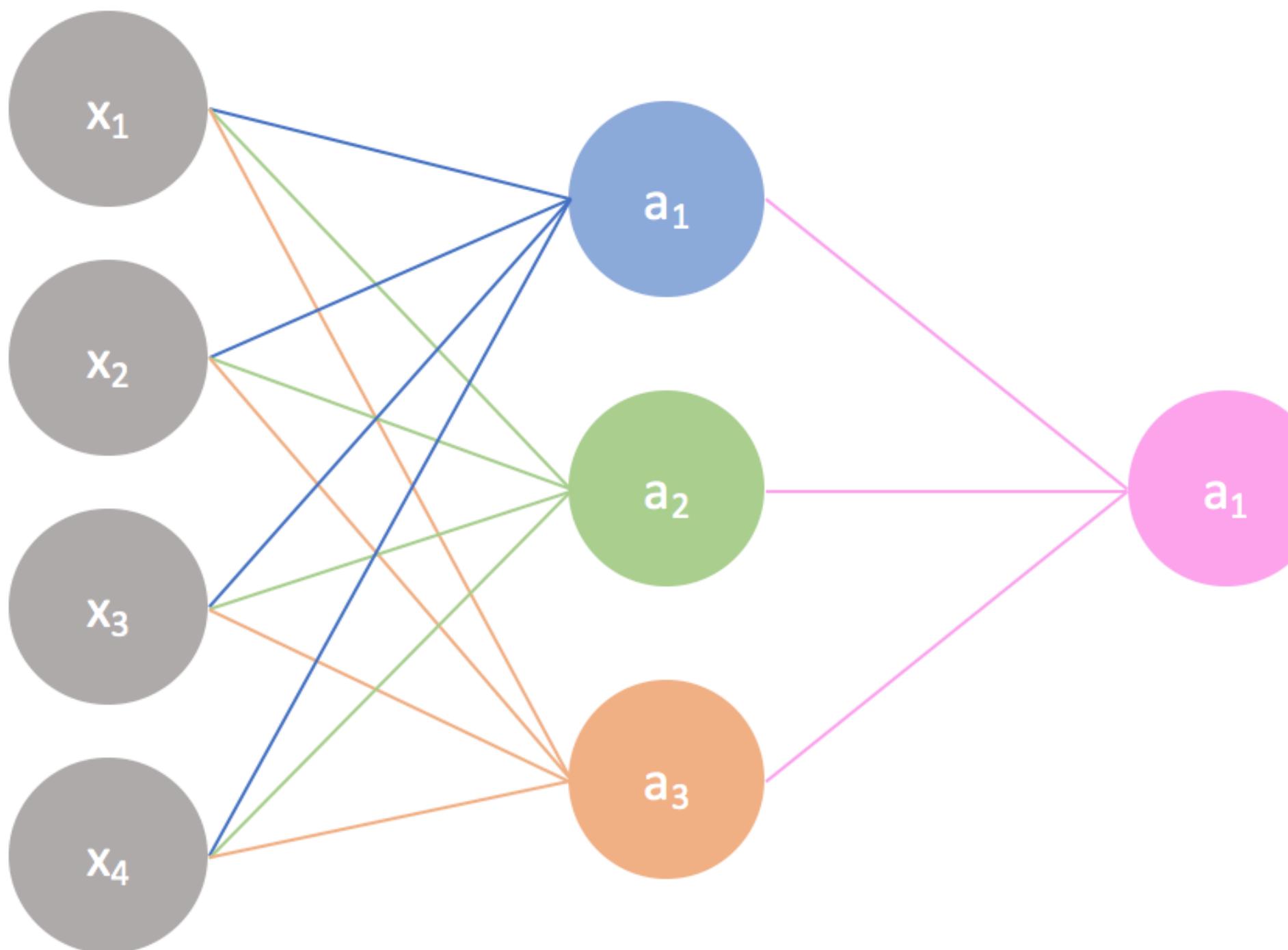


$$\begin{bmatrix} w_1 & w_2 & w_3 & w_4 \\ w_1 & w_2 & w_3 & w_4 \\ w_1 & w_2 & w_3 & w_4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} + \begin{bmatrix} b \\ b \\ b \end{bmatrix} = \begin{bmatrix} w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + b \\ w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + b \\ w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + b \end{bmatrix} \xrightarrow{\text{activation}} \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix}$$

The forward pass can be abstracted as a series of matrix multiplications, For which we have already studied efficient methods. We will come back to this.

# Creating NNs with PyTorch

## Simplicity of high-level DL frameworks



```
import torch
import torch.nn as nn

class NeuralNetwork(nn.Module):
    def __init__(self):
        super(NeuralNetwork, self).__init__()
        self.model = nn.Sequential(
            nn.Linear(4, 3),
            nn.ReLU(),
            nn.Linear(3, 1)
        )

    def forward(self, x):
        return self.model(x)

# Create an instance of the NeuralNetwork class
model = NeuralNetwork()
print(model)
```

# **Convolutional NNs: Neural Networks on Steroids**

# Limits of Regular NNs

Importance of feature extraction



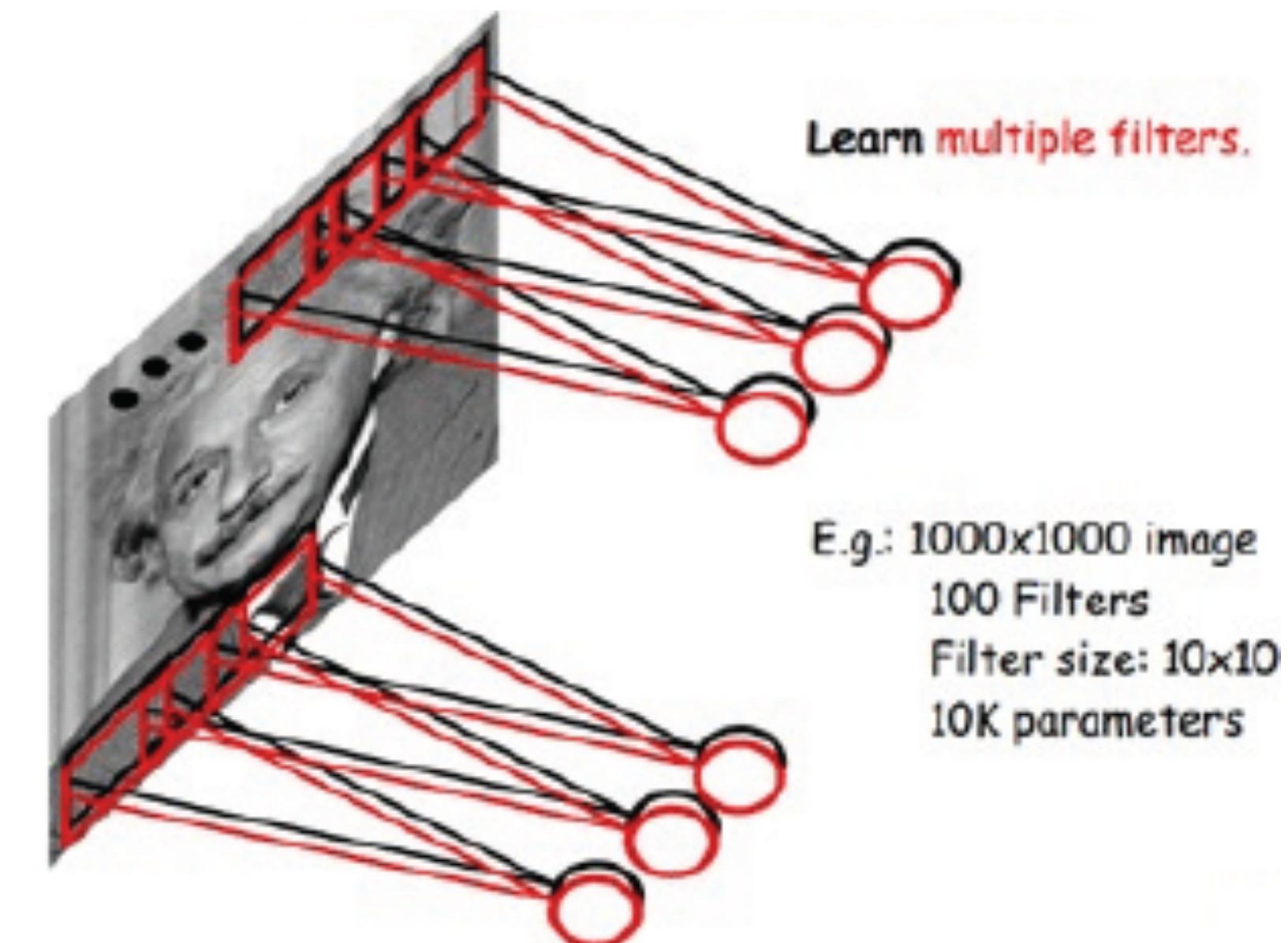
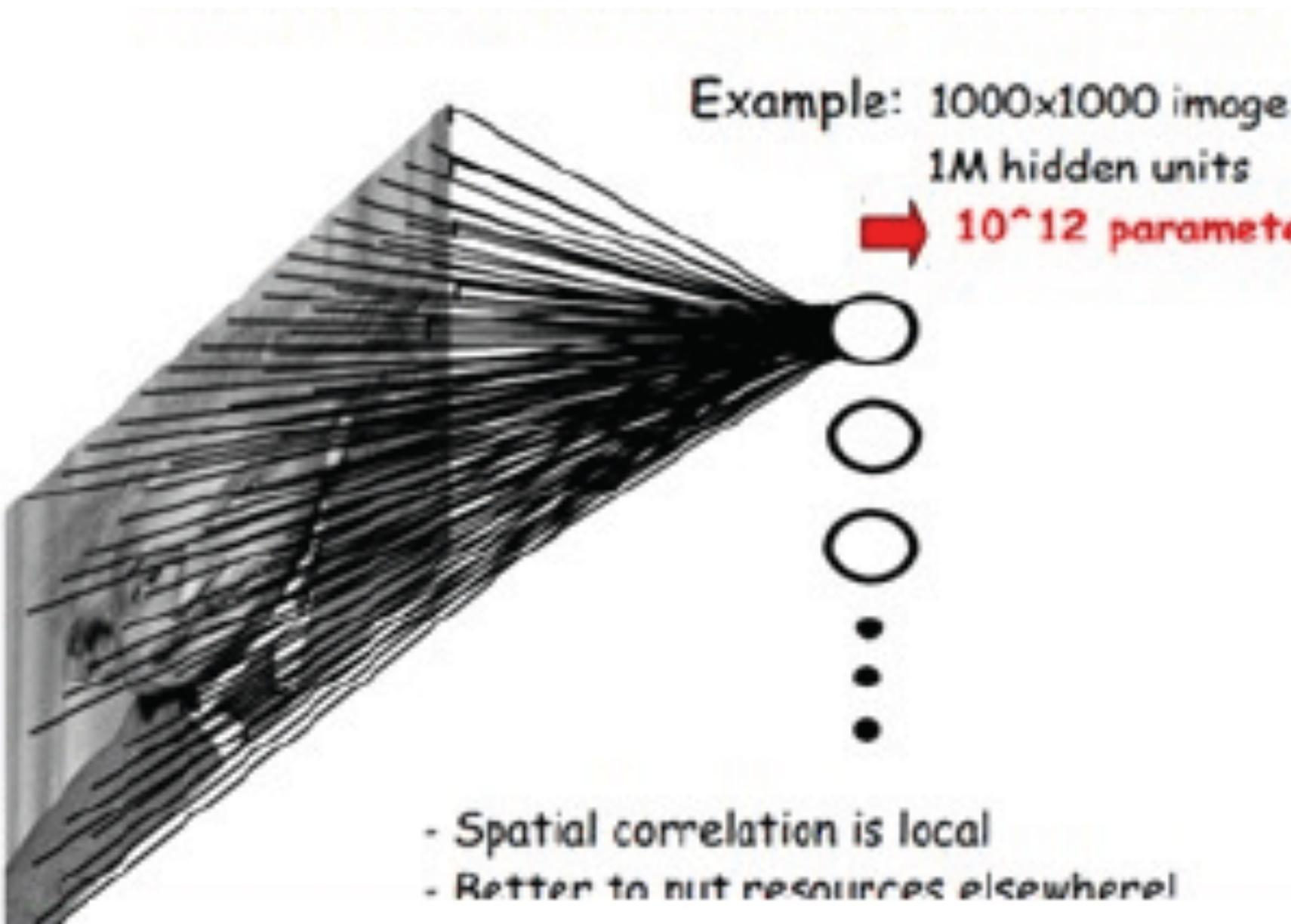
(a)



(b)

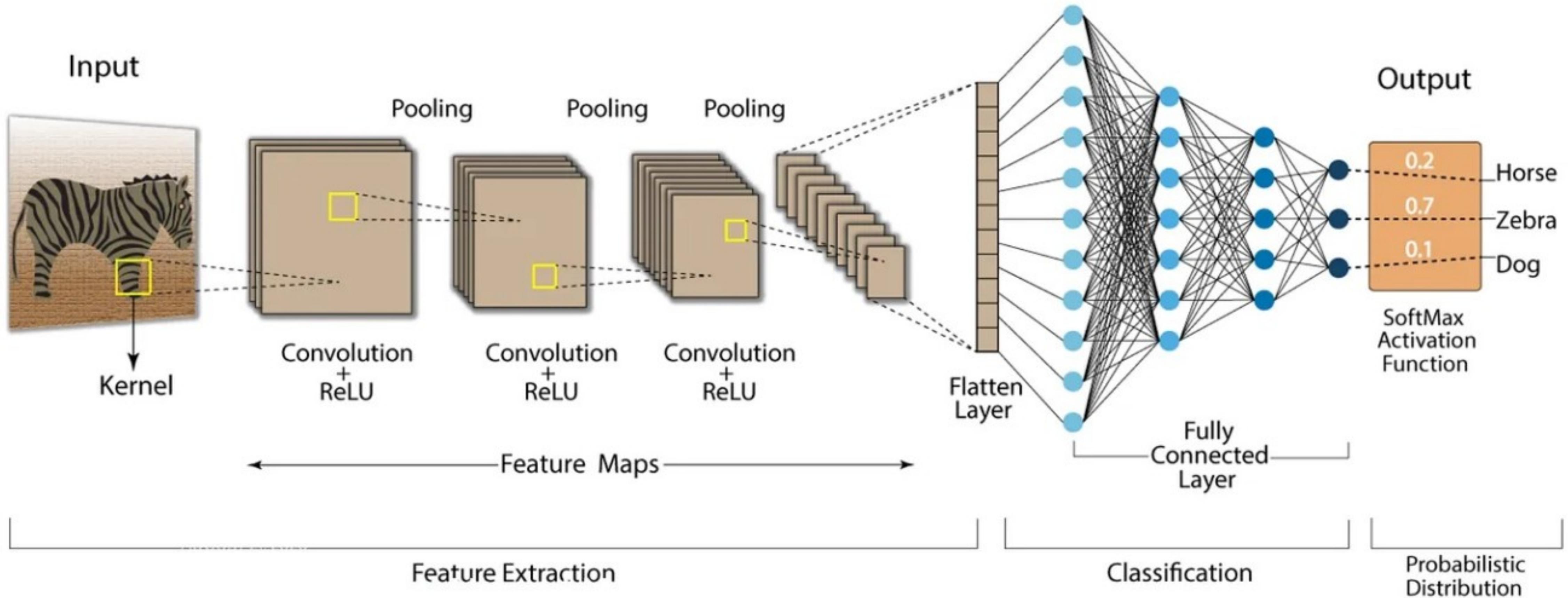
# Limits of Regular NNs

Unfeasible for large inputs (e.g. images)



# Convolutional Neural Networks

Giving computers the power understand images



# The Convolution Operation



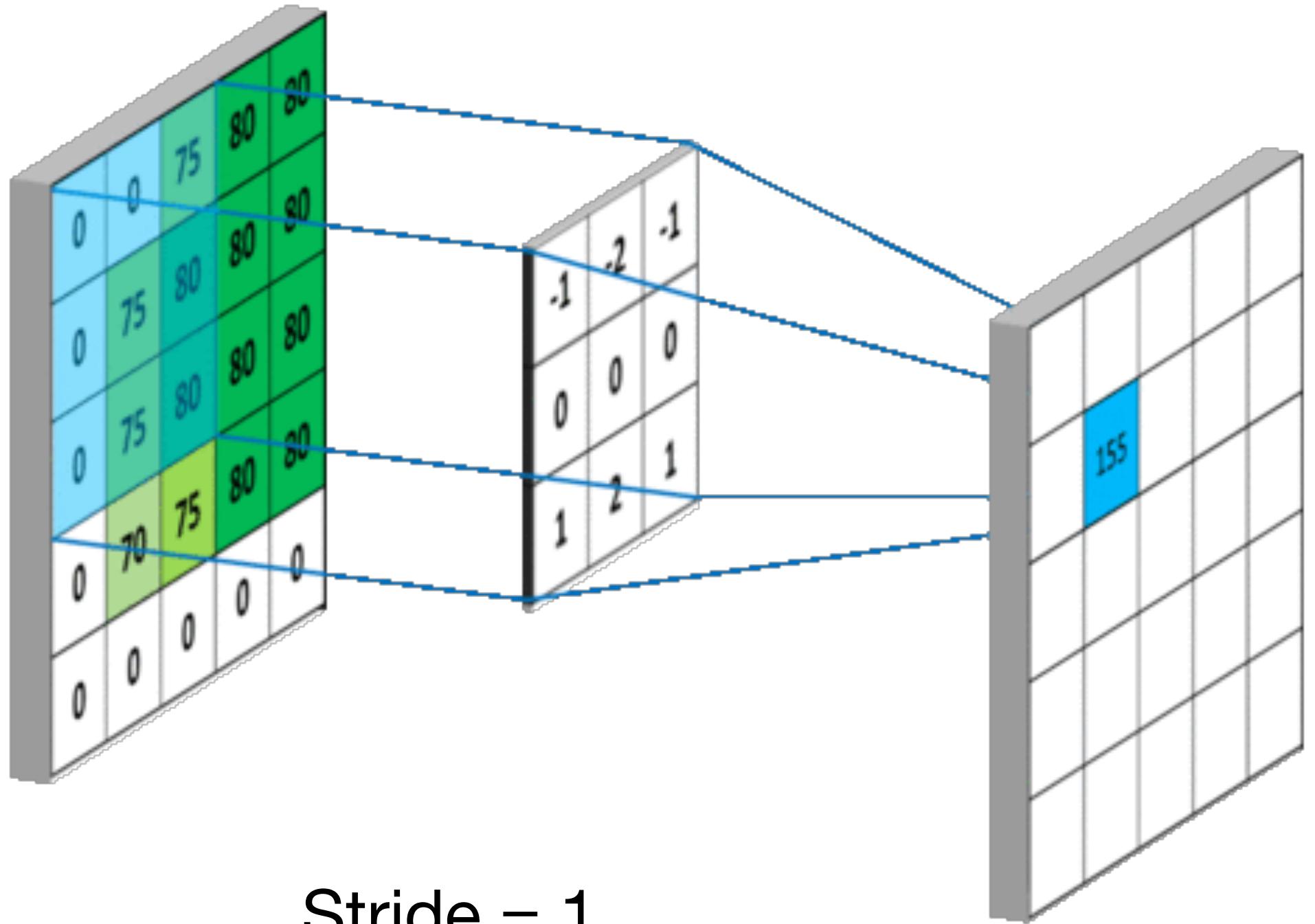
-1	-2	-1
0	0	0
1	2	1

(a) Vertical edge detector



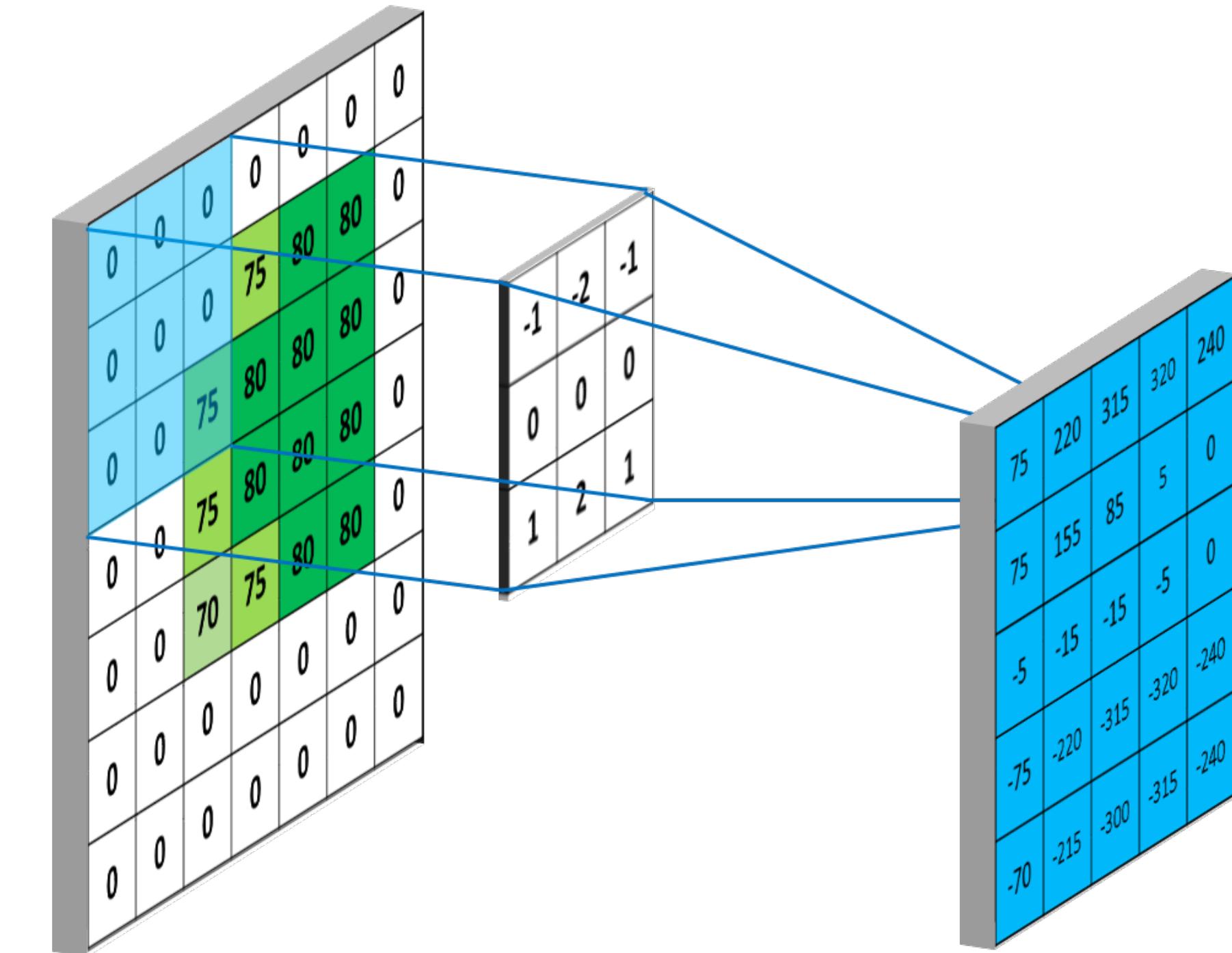
# Convolution Layers

The first half of CNNs



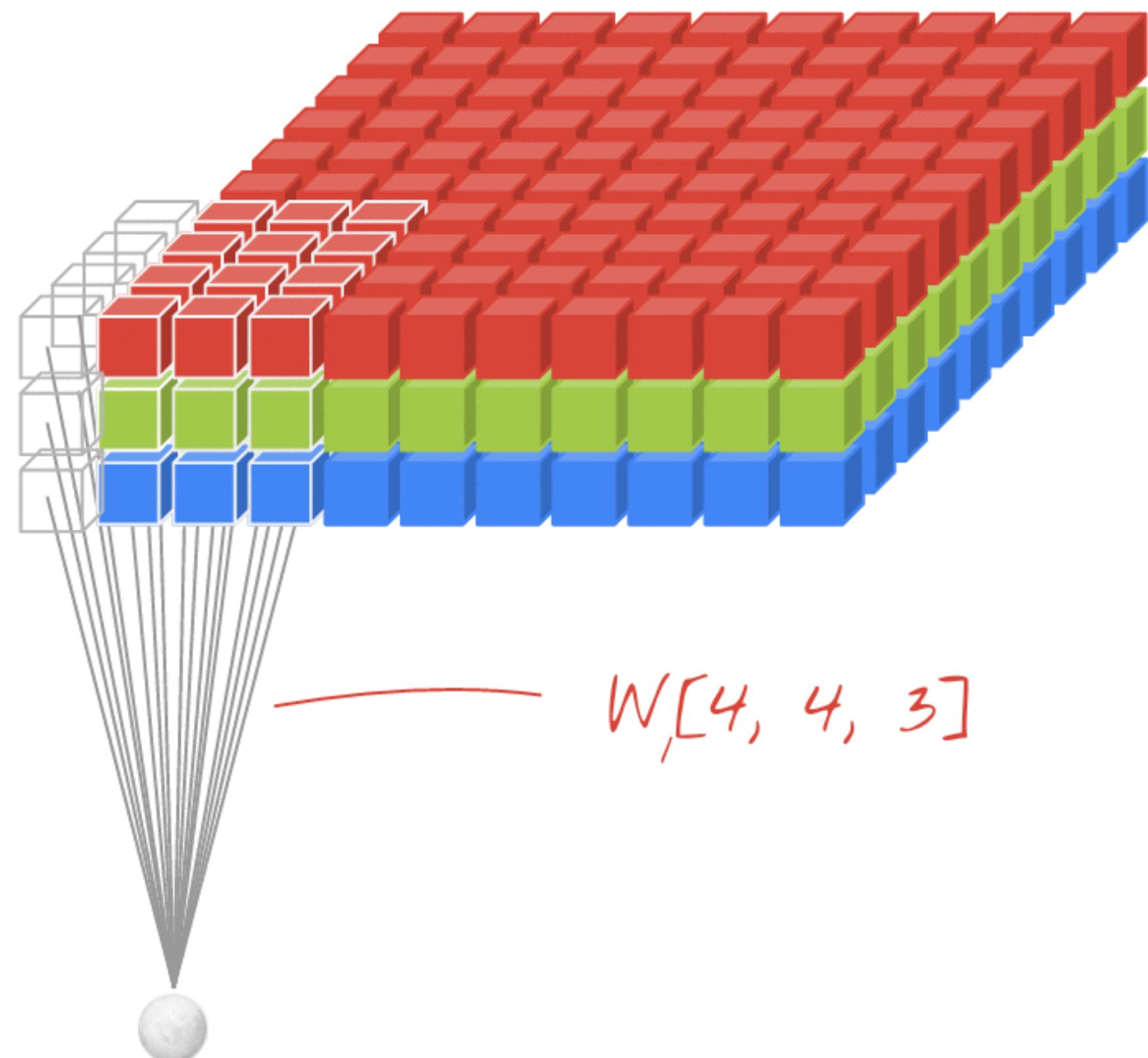
Stride = 1

Kernel Size = 3



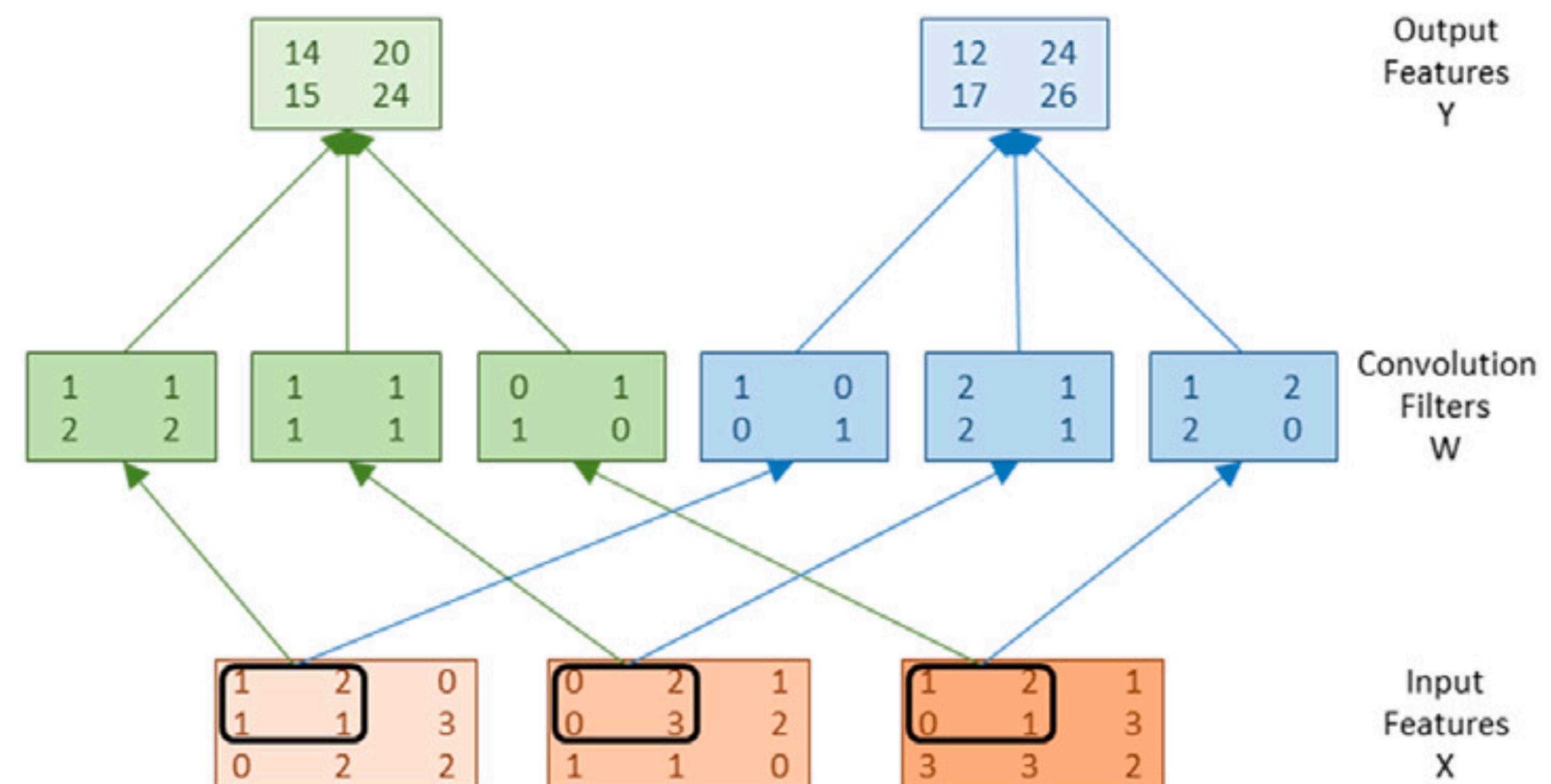
# Convolution Layers

## Multi-layer 3d layers



# Convolution Layers

Yet another example



# Coding Convolutional Layers

## A serialized implementation

The diagram illustrates the flow of data from input parameters to the convolution loop. It shows five input parameters: #filters, #img-channel, img-height, filter-size, and img-values. Arrows point from these parameters to the corresponding variables in the code: #filters to M, #img-channel to C, img-height to H, filter-size to K, and img-values to X. The output of the convolution loop is labeled 'output-img'.

```
01 void convLayer_forward(int M, int C, int H, int W, int K, float* X, float* W,
02                         float* Y) {
03     int H_out = H - K + 1;
04     int W_out = W - K + 1;           output-img
05
06     for(int m = 0; m < M; m++)          // for each output feature map
07         for(int h = 0; h < H_out; h++)    // for each output element
08             for(int w = 0; w < W_out; w++) {
09                 Y[m, h, w] = 0;
10                 for(int c = 0; c < C; c++) // sum over all input feature maps
11                     for(int p = 0; p < K; p++) // KxK filter
12                         for(int q = 0; q < K; q++)
13                             Y[m, h, w] += X[c, h + p, w + q] * W[m, c, p, q];
14 }
```

# Parallel Implementation; Kernel Launch

- each thread computes one element of one output feature map
- each block is 2d and computes a tile of  $\text{TILE\_WIDTH} \times \text{TILE\_WIDTH}$  pixels in one output feature map
- The grid is 2d, The first dimension (X) corresponds to the (M) output features maps covered by each block. The second dimension (Y) reflects the location of a block's output tile inside the output feature map.

```
01 # define TILE_WIDTH 16
02 W_grid = W_out/TILE_WIDTH; // number of horizontal tiles per output map
03 H_grid = H_out/TILE_WIDTH; // number of vertical tiles per output map
04 T = H_grid * W_grid;
05 dim3 blockDim(TILE_WIDTH, TILE_WIDTH, 1);
06 dim3 gridDim(M, T, N);
07 ConvLayerForward_Kernel<<< gridDim, blockDim>>>(...);
```

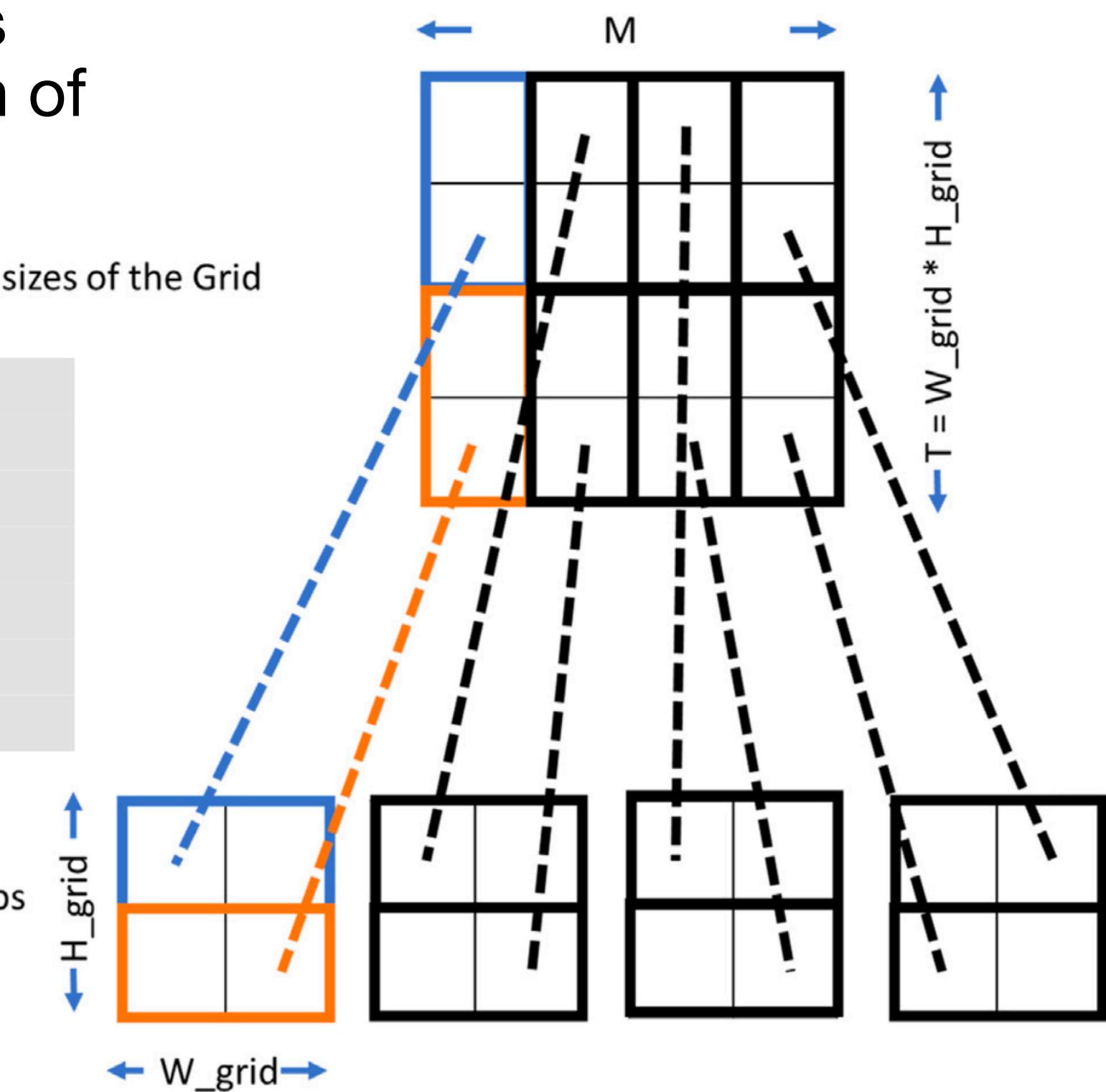
$$W_{\text{out}} = 128 \rightarrow W_{\text{grid}} = 128/16 = 8$$

$$H_{\text{out}} = 128 \rightarrow H_{\text{grid}} = 128/16 = 8$$

$$T = H_{\text{grid}} * W_{\text{grid}} = 8 * 8 = 64$$



X-Y dimension sizes of the Grid



# Coding Convolutional Layers

## A parallel implementation

```
01 __global__ void
02 ConvLayerForward_Kernel(int C, int W_grid, int K, float* X, float* W,
                           float* Y) {
03     int m = blockIdx.x;
04     int h = (blockIdx.y / W_grid)*TILE_WIDTH + threadIdx.y;
05     int w = (blockIdx.y % W_grid)*TILE_WIDTH + threadIdx.x;
06     int n = blockIdx.z;
07     float acc = 0.;
08     for (int c = 0; c < C; c++) { // sum over all input channels
09         for (int p = 0; p < K; p++) // loop over KxK filter
10             for (int q = 0; q < K; q++)
11                 acc += X[n, c, h + p, w + q] * W[m, c, p, q];
12     }
13     Y[n, m, h, w] = acc;
14 }
```

The diagram illustrates the memory access pattern for the convolution kernel. It shows five variables with arrows pointing to specific memory locations in the code:

- output-img**: Points to the variable `Y` at line 13.
- #img-channels**: Points to the variable `C` at line 2.
- filter-size**: Points to the variable `K` at line 2.
- img-values**: Points to the variable `X` at line 2.
- filter-values**: Points to the variable `W` at line 2.

Red annotations on the left side of the code provide context for the variables:

- Filter number**: Points to `m` at line 3.
- Start pixel height**: Points to `h` at line 4.
- Start pixel width**: Points to `w` at line 5.
- Image batch #**: Points to `n` at line 6.

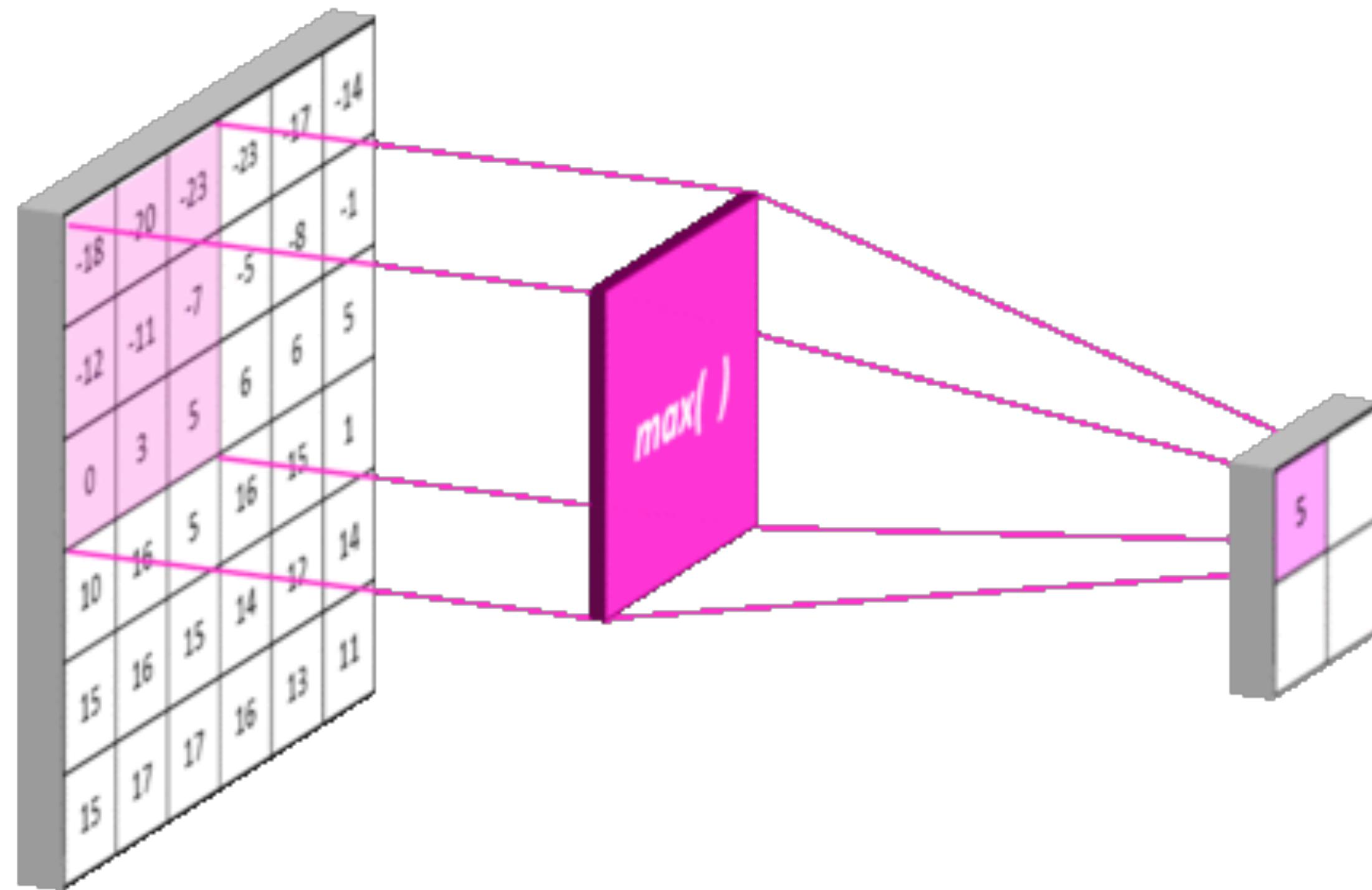
# Coding Convolution Layers

## Limitations and fixes

- High degree of parallelism
- Consumes too much global memory bandwidth
- Fix:
  - Using constant memory caching
  - Using shared memory tiling

# Pooling Layers

## The second half of CNNs



Stride = 3

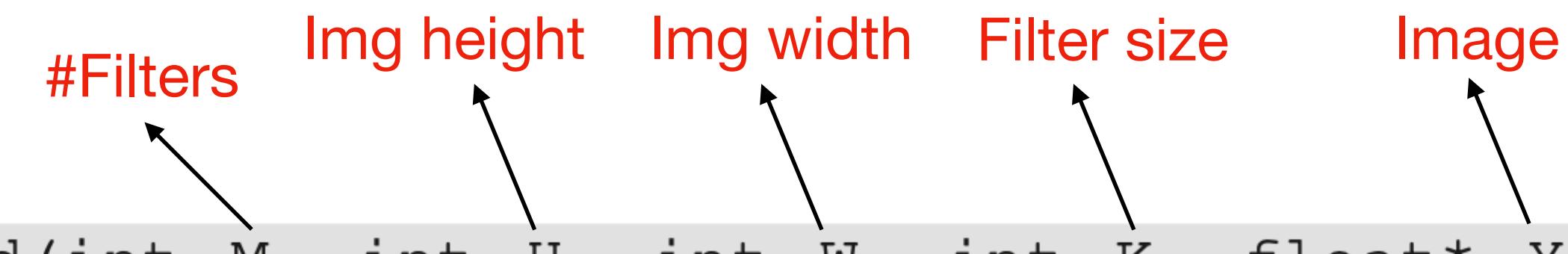
Kernel Size = 3

# Coding Pooling Layers

## A serialized implementation

```
01 void subsamplingLayer_forward(int M, int H, int W, int K, float* Y, float* S) {  
02     for(int m = 0; m < M; m++) // for each output feature map  
03         for(int h = 0; h < H/K; h++) // for each output element,  
04             for(int w = 0; w < W/K; w++) { // this code assumes that H and W  
05                 S[m, x, y] = 0.; // are multiples of K  
06                 for(int p = 0; p < K; p++) { // loop over KxK input samples  
07                     for(int q = 0; q < K; q++)  
08                         S[m, h, w] += Y[m, K*h + p, K*w+ q] / (K*K);  
09                 }  
10                     // add bias and apply non-linear activation  
11                     S[m, h, w] = sigmoid(S[m, h, w] + b[m]);  
12 }
```

#Filters      Img height      Img width      Filter size      Image



# Coding Pooling Layers

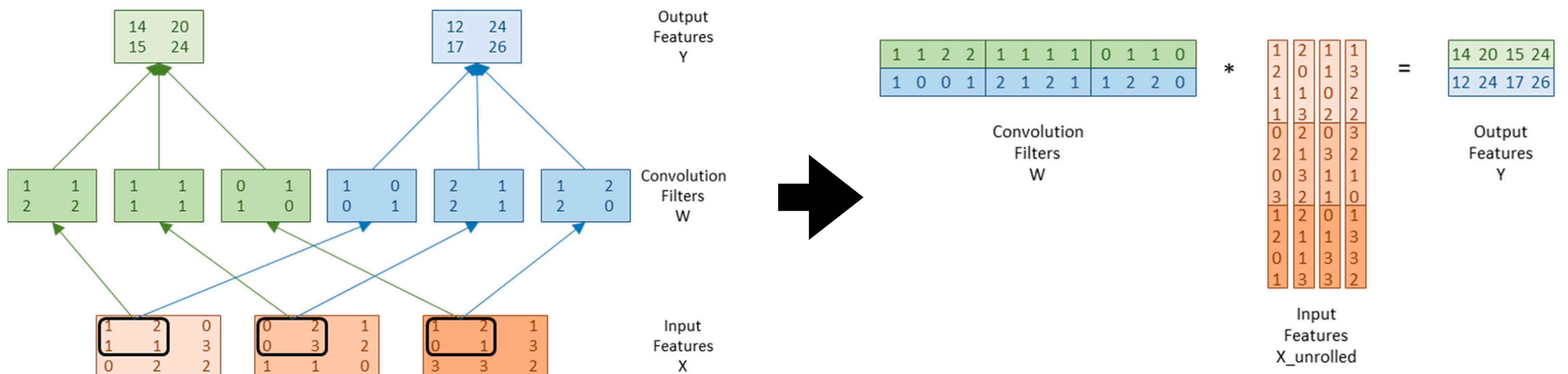
## A parallel implementation

**This might look good, but we can  
do much better...!**

Remember how we could represent a neural  
network as a series of matrix multiplications?

# Unrolling Convolution Layers

Formulating the problem as matrix multiplication

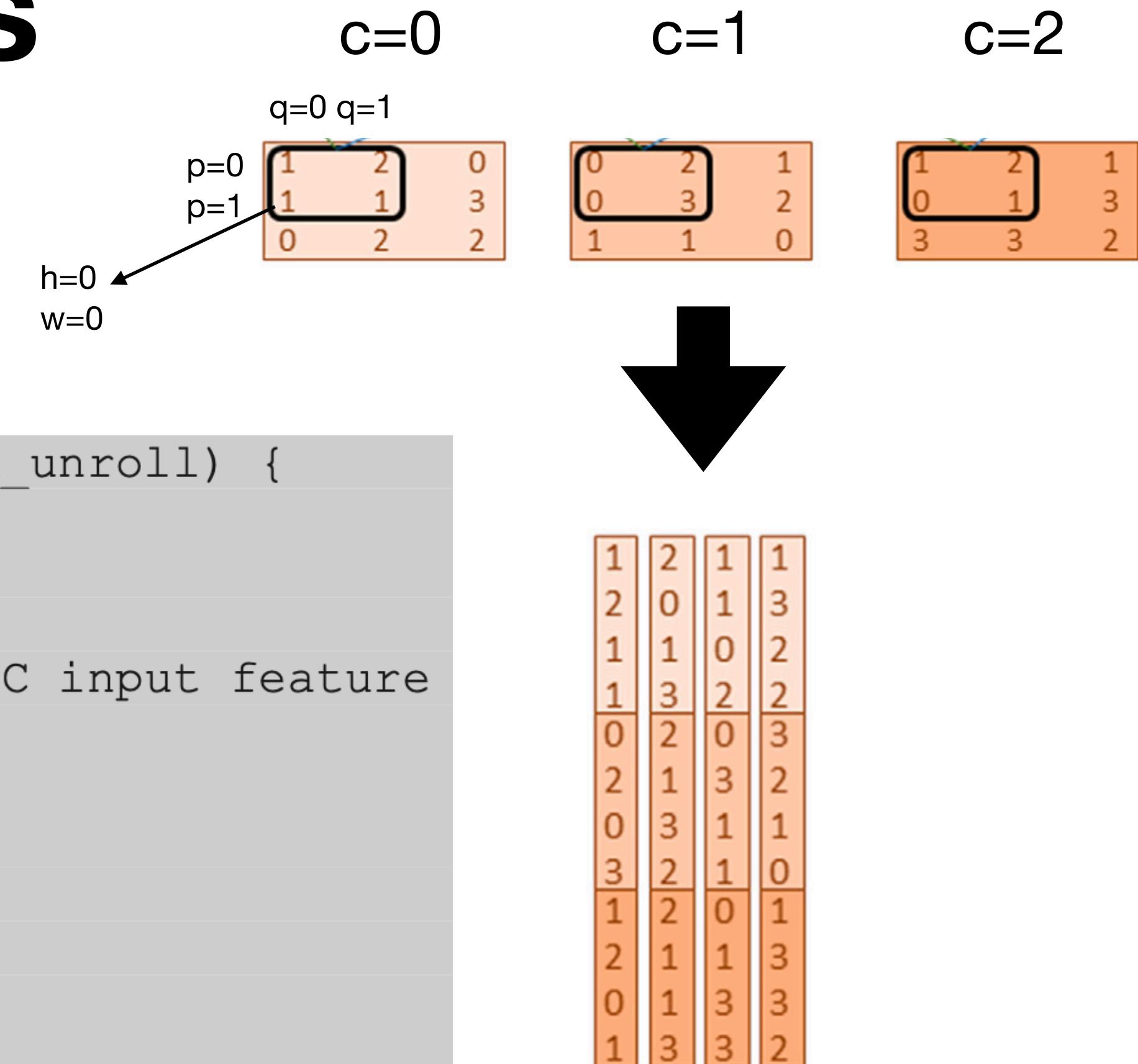


Once we unroll it, the resulting multiplication can be easily solved using warptiling or even cuBLAS

# Unrolling Convolution Layers

## Sequential implementation

```
01 void unroll(int C, int H, int W, int K, float* X, float* X_unroll) {  
02     int H_out = H - K + 1;  
03     int W_out = W - K + 1;  
04 #channels for(int c = 0; c < C; c++) {  
05     // Beginning row index of the section for channel C input feature  
06     // map in the unrolled matrix  
07     Start-index w_base = c * (K*K);  
08     for(int p = 0; p < K; p++) {  
09         for(int q = 0; q < K; q++) {  
10             for(int h = 0; h < H_out; h++) {  
11                 int h_unroll = w_base + p*K + q;  
12                 for(int w = 0; w < W_out; w++) {  
13                     int w_unroll = h * W_out + w;  
14                     X_unroll[h_unroll, w_unroll] = X(c, h + p, w + q);  
15                 }  
16             }  
17         }  
18     }
```



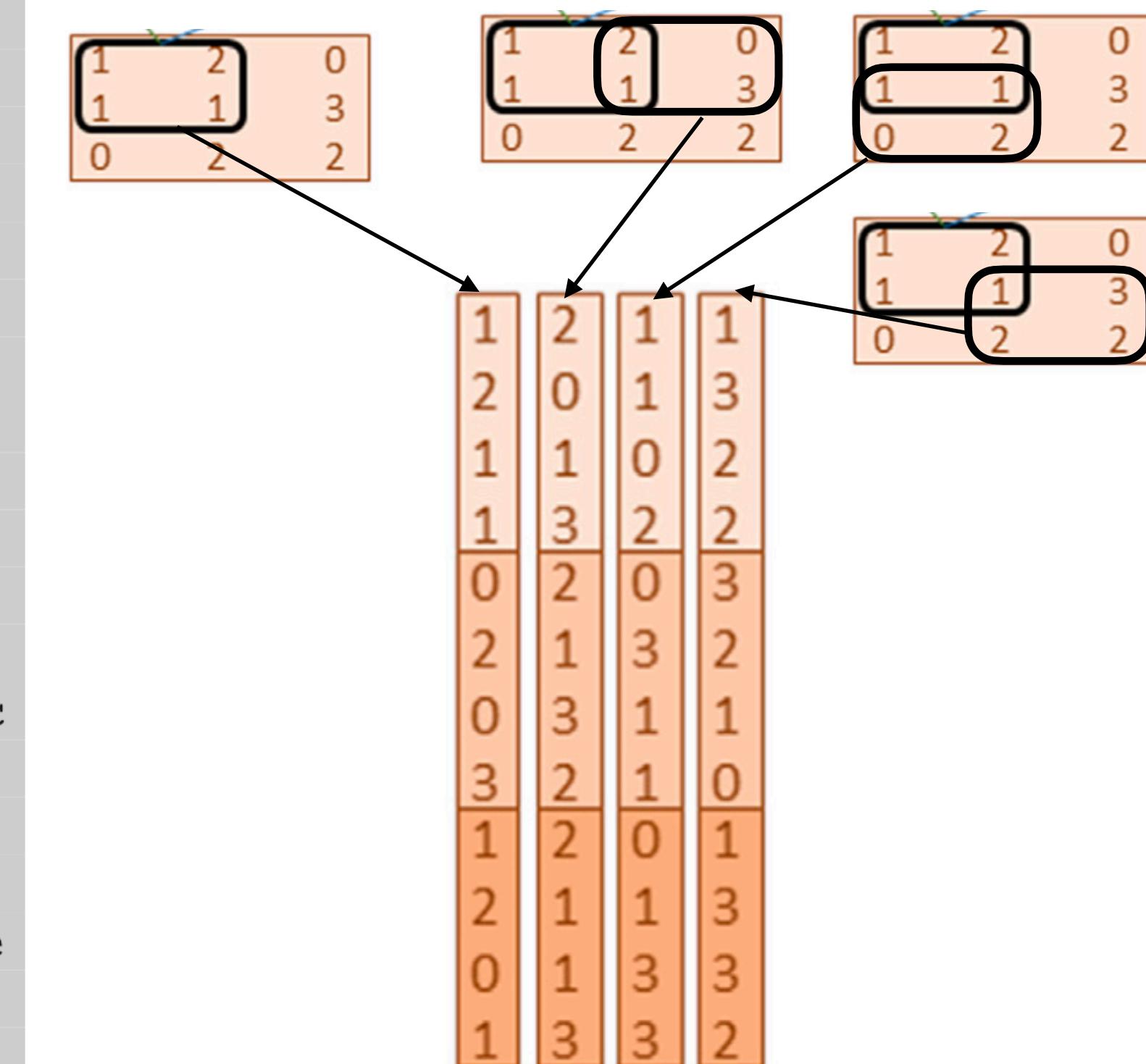
# Unrolling Convolution Layers

## Parallel implementation

#img-channel

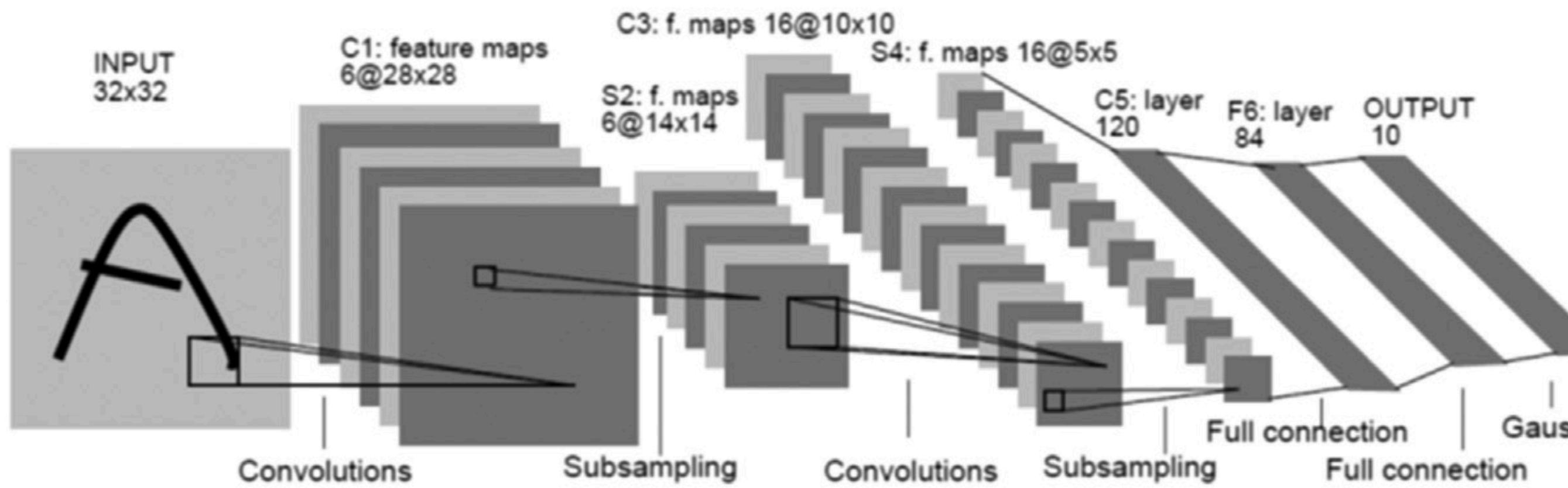
```
01 __global__ void
02 unroll_Kernel(int C, int H, int W, int K, float* X, float* X_unroll) {
03     int t = blockIdx.x * blockDim.x + threadIdx.x;
04     int H_out = H - K + 1;
05     int W_out = W - K + 1;
06     // Width of the unrolled input feature matrix
07     int W_unroll = H_out * W_out;
08     if (t < C * W_unroll) {
09         // Channel of the input feature map being collected by the thread
10         int c = t / W_unroll;
11         // Column index of the unrolled matrix to write a strip of
12         // input elements into (also, the linearized index of the output
13         // element for which the thread is collecting input elements)
14         int w_unroll = t % W_unroll;
15         // Horizontal and vertical indices of the output element
16         int h_out = w_unroll / W_out;
17         int w_out = w_unroll % W_out;
18         // Starting row index for the unrolled matrix section for channel c
19         int w_base = c * K * K;
20         for(int p = 0; p < K; p++) {
21             for(int q = 0; q < K; q++) {
22                 // Row index of the unrolled matrix for the thread to write
23                 // the input element into for the current iteration
24                 int h_unroll = w_base + p*K + q;
25                 X_unroll[h_unroll, w_unroll] = X[c, h_out + p, w_out + q];
26             }
27         }
28     }
29 }
```

- Each CUDA thread will be responsible for gathering ( $K \times K$ ) input elements from one input feature map
- The total number of threads will be  $(C \times H_{out} \times W_{out})$



# LeNet-5

## A revolution in image processing



Learn PyTorch in 1 hour:



Everything you need to know to get started

[www.aparat.com/v/MOhuQ](http://www.aparat.com/v/MOhuQ)



```
class LeNet5(nn.Module):
    def __init__(self, num_classes):
        super(ConvNeuralNet, self).__init__()
        self.layer1 = nn.Sequential(
            nn.Conv2d(1, 6, kernel_size=5, stride=1, padding=0),
            nn.BatchNorm2d(6),
            nn.ReLU(),
            nn.MaxPool2d(kernel_size = 2, stride = 2))
        self.layer2 = nn.Sequential(
            nn.Conv2d(6, 16, kernel_size=5, stride=1, padding=0),
            nn.BatchNorm2d(16),
            nn.ReLU(),
            nn.MaxPool2d(kernel_size = 2, stride = 2))
        self.fc = nn.Linear(400, 120)
        self.relu = nn.ReLU()
        self.fc1 = nn.Linear(120, 84)
        self.relu1 = nn.ReLU()
        self.fc2 = nn.Linear(84, num_classes)

    def forward(self, x):
        out = self.layer1(x)
        out = self.layer2(out)
        out = out.reshape(out.size(0), -1)
        out = self.fc(out)
        out = self.relu(out)
        out = self.fc1(out)
        out = self.relu1(out)
        out = self.fc2(out)
        return out
```

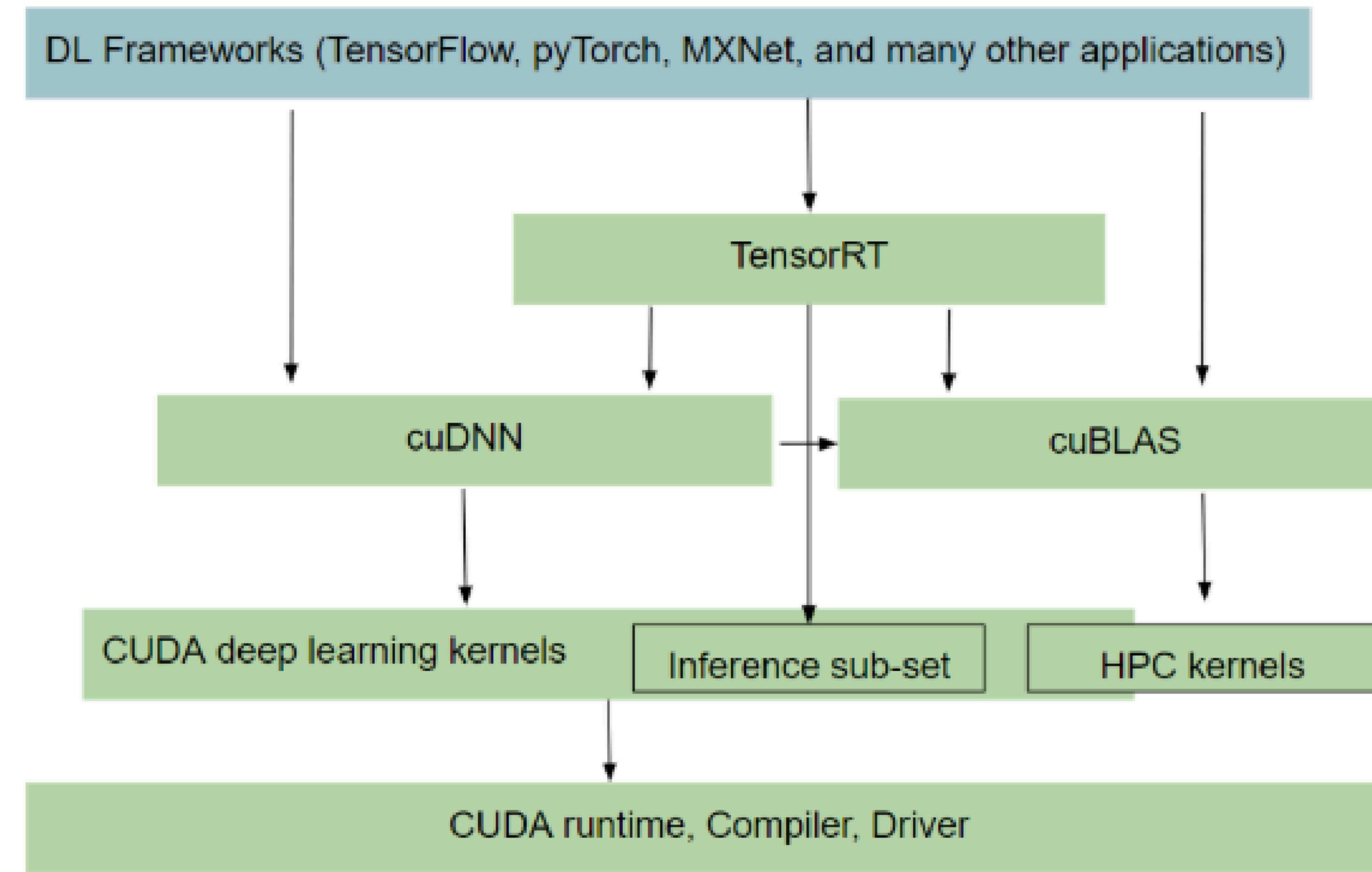
# CuDNN: A Higher Level of Abstraction

## NVIDIA's Deep Neural Network library

- A library of pre-written primitives CUDA kernels for deep neural networks
- Accelerated compute-bound operations like convolution and matmul
- Optimized memory-bound operations like pooling, softmax, normalization and activation
- Example APIs:
  - > `cudnnConvolutionBackwardBias()`
  - > `cudnnConvolutionBackwardData()`
  - > `cudnnConvolutionBackwardFilter()`
  - > `cudnnConvolutionBiasActivationForward()`
  - > `cudnnConvolutionForward()`
  - > `cudnnCreateAttnDescriptor()`
  - > `cudnnCreateConvolutionDescriptor()`
  - > `cudnnPoolingBackward()`
  - > `cudnnPoolingForward()`
  - > `cudnnReduceTensor()`
  - > `cudnnReorderFilterAndBias()`
  - > `cudnnRNNGetClip_v8()`
  - > `cudnnRNNSetClip_v8()`
  - > `cudnnScaleTensor()`
  - > `cudnnSetAttnDescriptor()`

# Where CuDNN Fits in the DL Ecosystem

## A backbone of modern DL Frameworks



# Summary and Conclusion

Things you have (hopefully) learned:

- How neural networks are changing the world
- Brief explanation of how neural networks work under the hood
- CNNs, convolutions and pooling layers
- How to both naively and professionally parallelize neural networks
- How libraries like cuBLAS and cuDNN can simplify the process

**The End**  
**Thanks for your attention**