

Neural Network Basics

João Campagnolo

Computational
Neuroscience class

Spring 2024

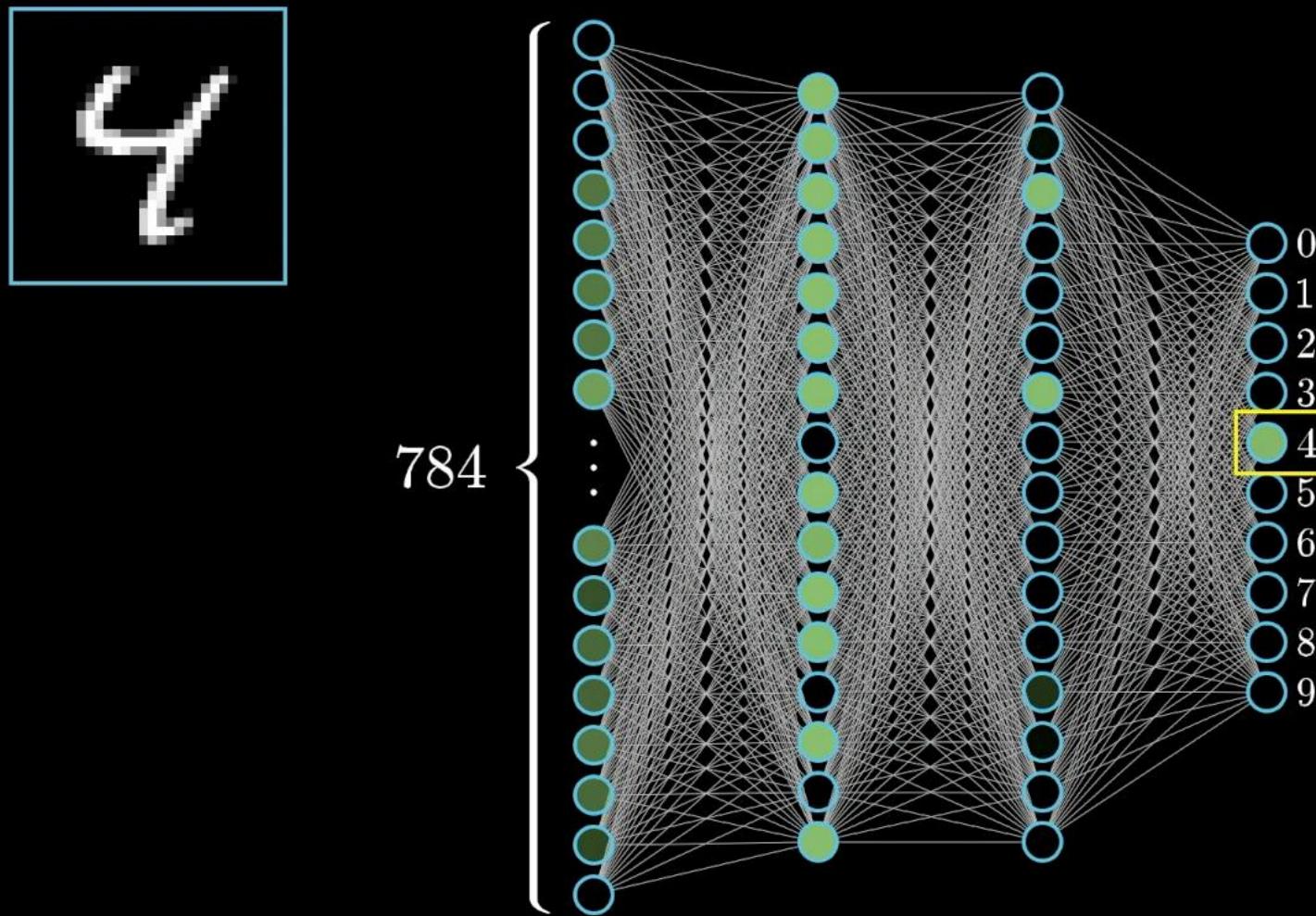
UNIVERSITY OF COPENHAGEN



Let's start with the example from [3b1b](#)

Plain vanilla

(aka “multilayer perceptron”)

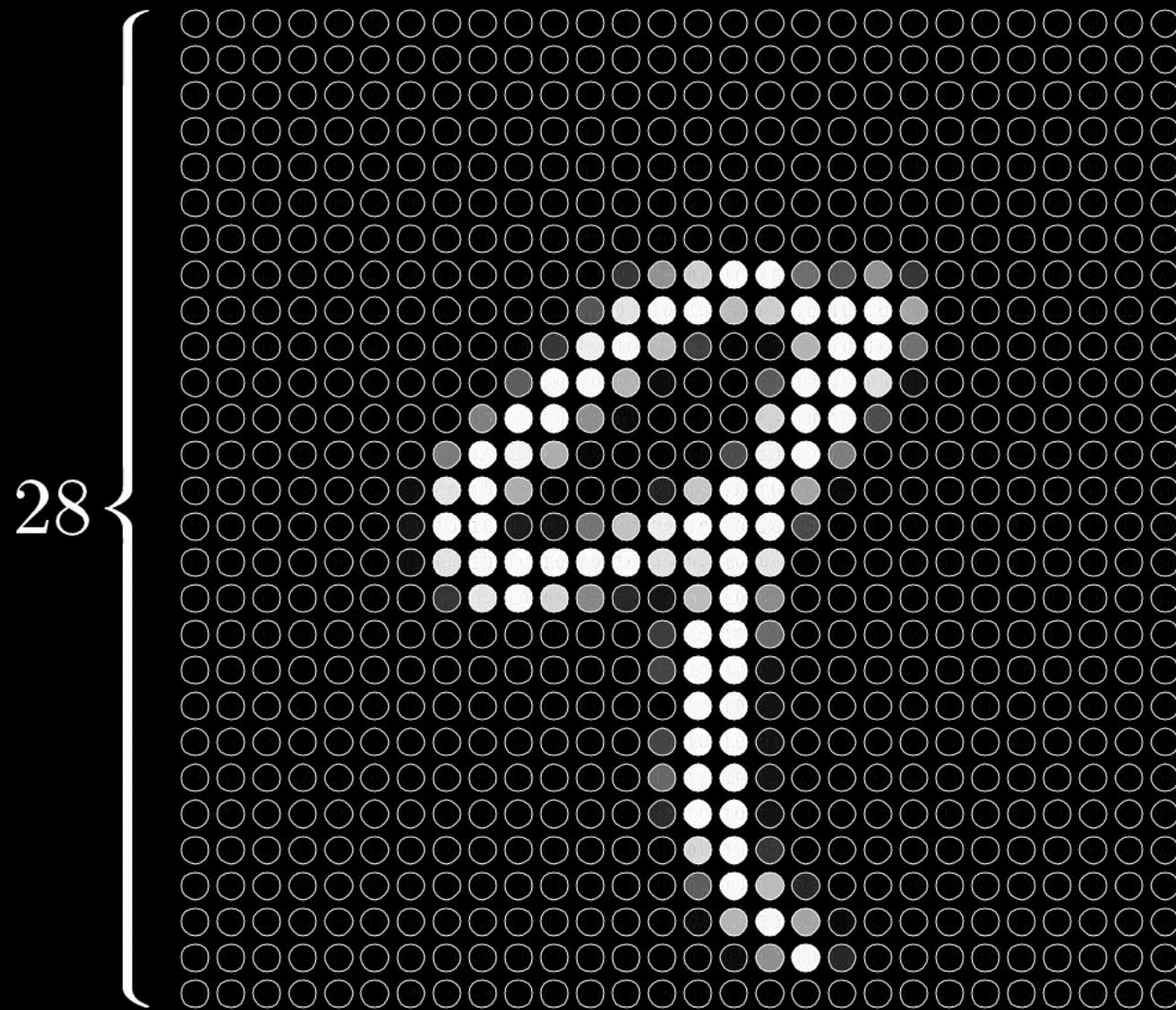


A horizontal color bar consisting of six circles of increasing size from left to right. Each circle contains a numerical value: 0.0, 0.2, 0.4, 0.6, 0.8, and 1.0. The circles are arranged in a row, with the first three having black outlines and the last three having white outlines.

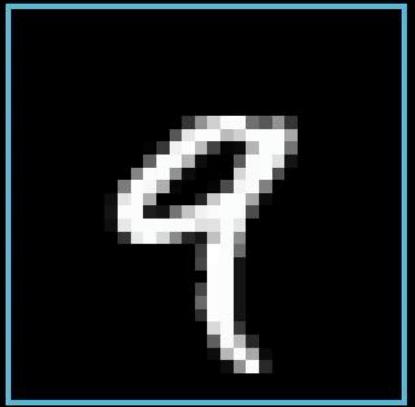
Each pixel in the original image has a value between 0.0 (black) and 1.0 (white).

Every neuron has an activation between 0.0 and 1.0, sort of analogous to how neurons in the brain can be active or inactive.

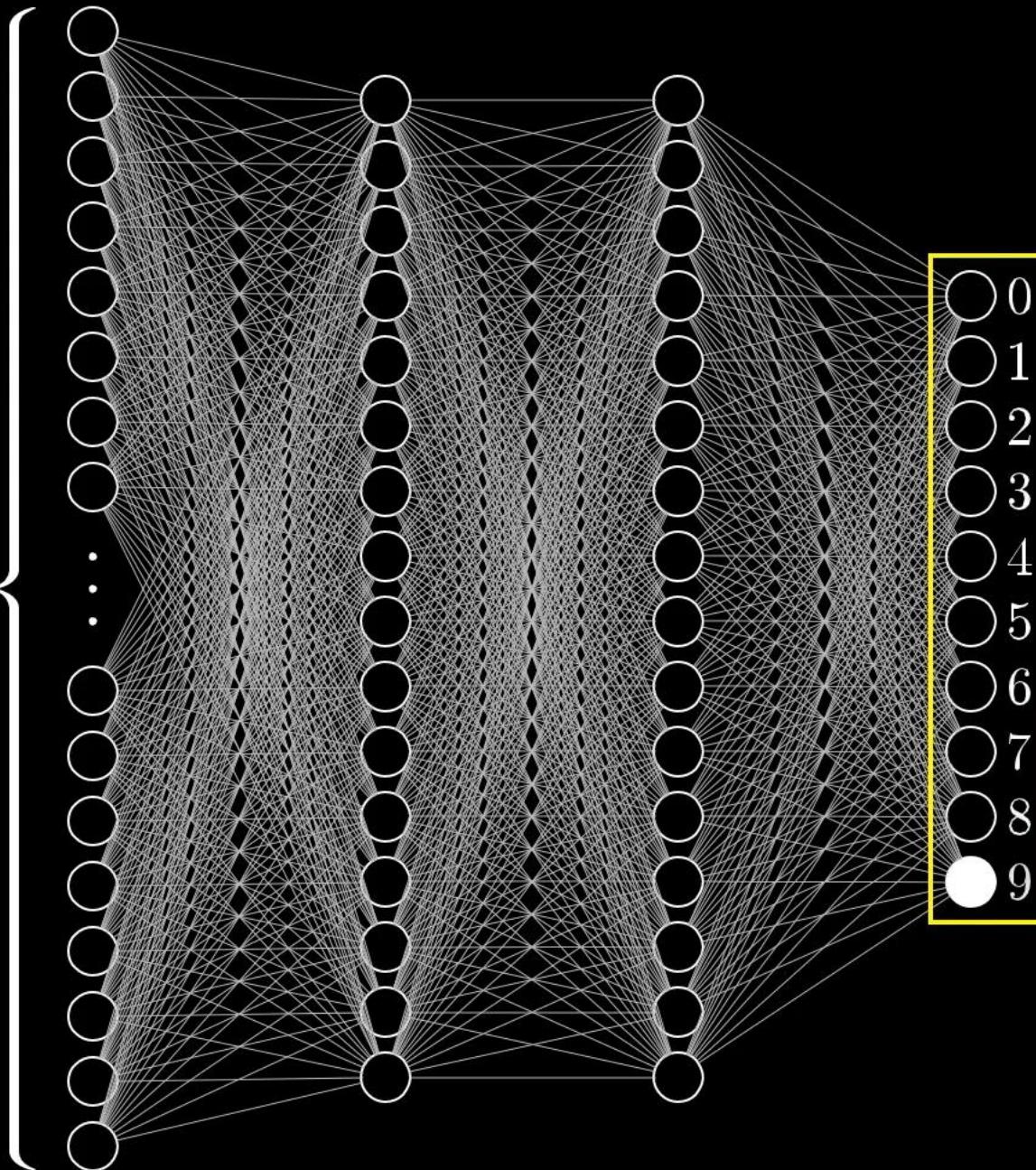
28



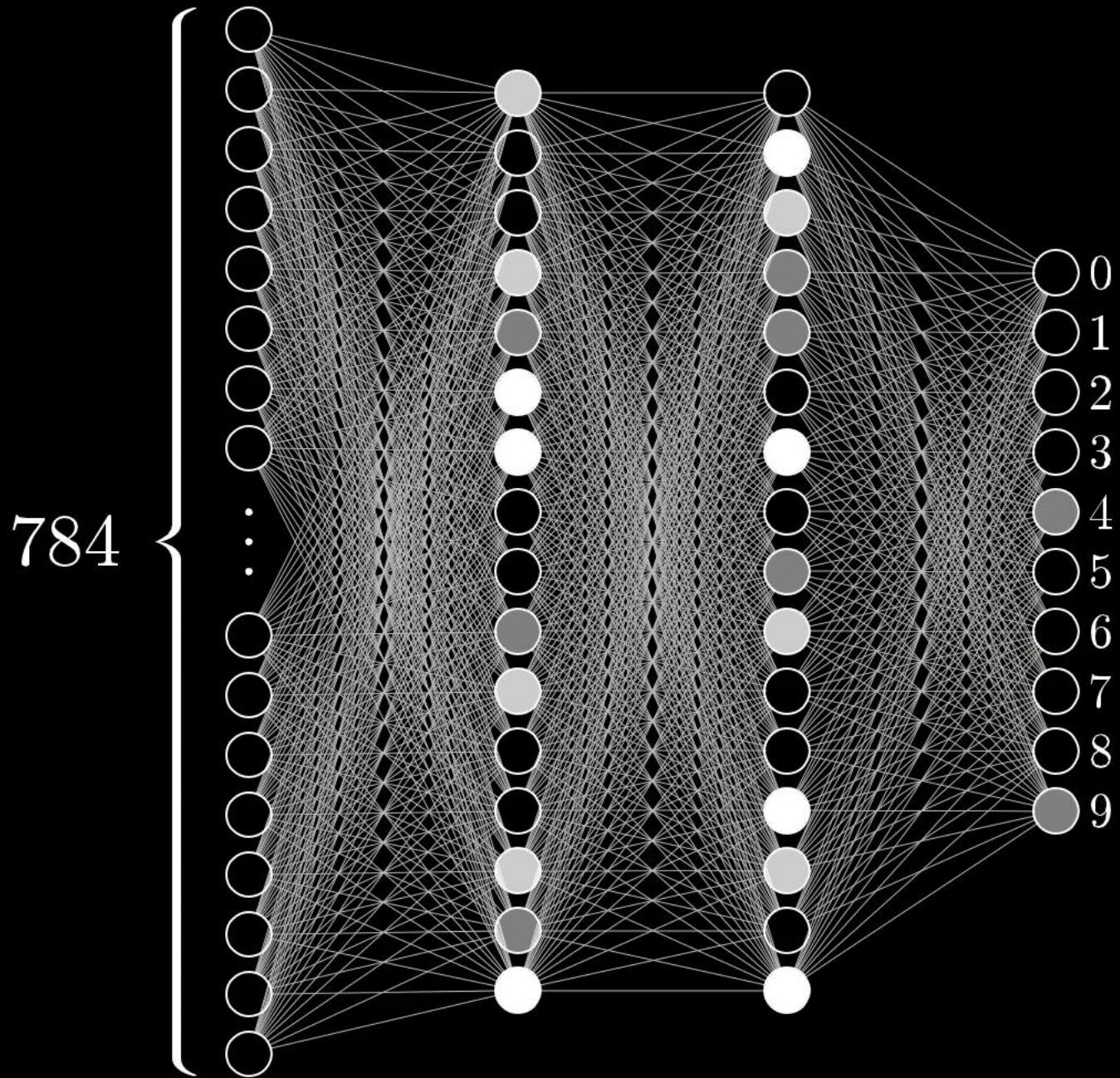
$$28 \times 28 = 784$$



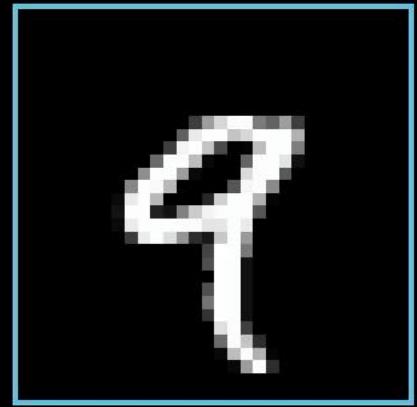
784



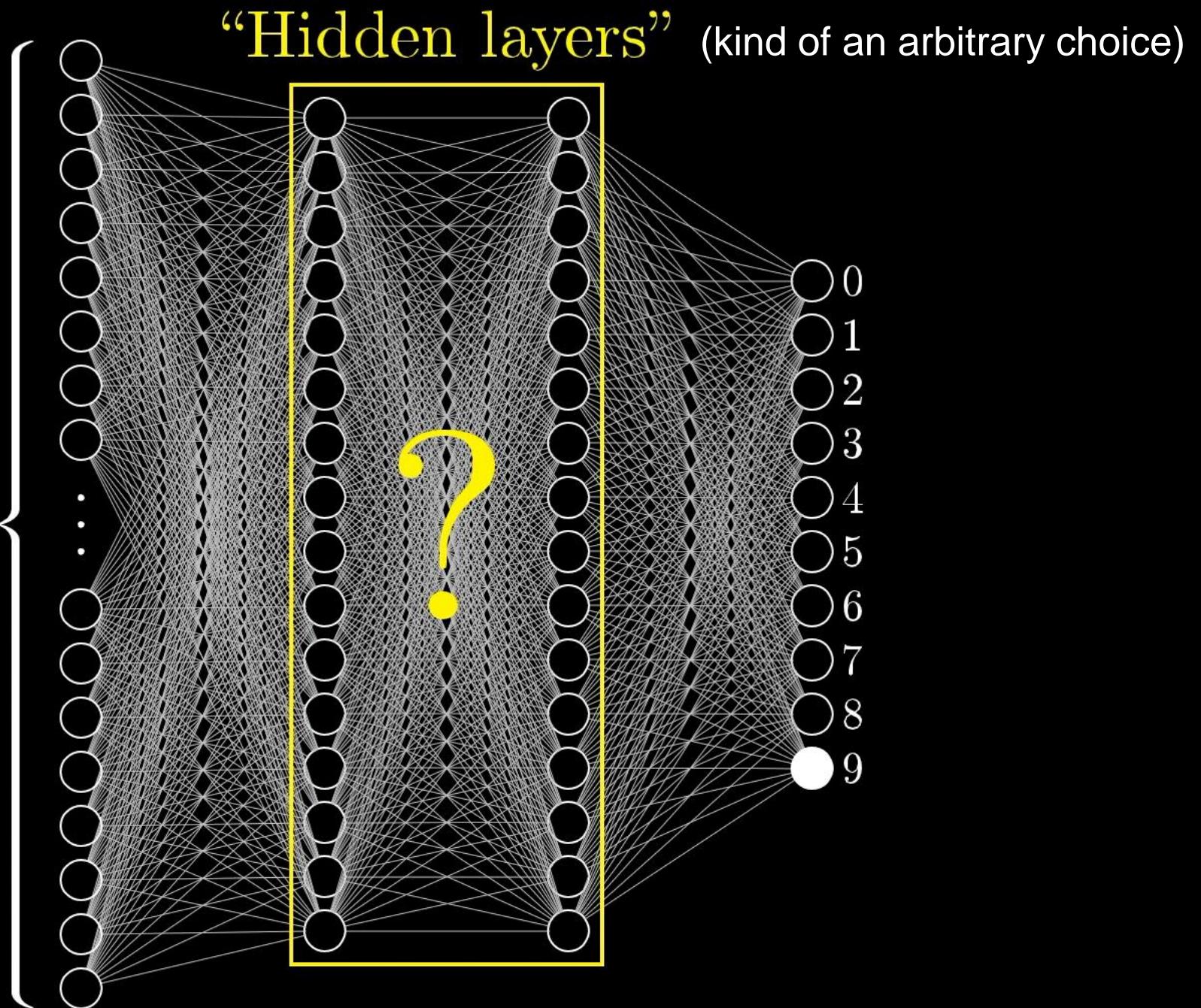
Set input and
output
dimensions



What is the
output the
network will
choose?

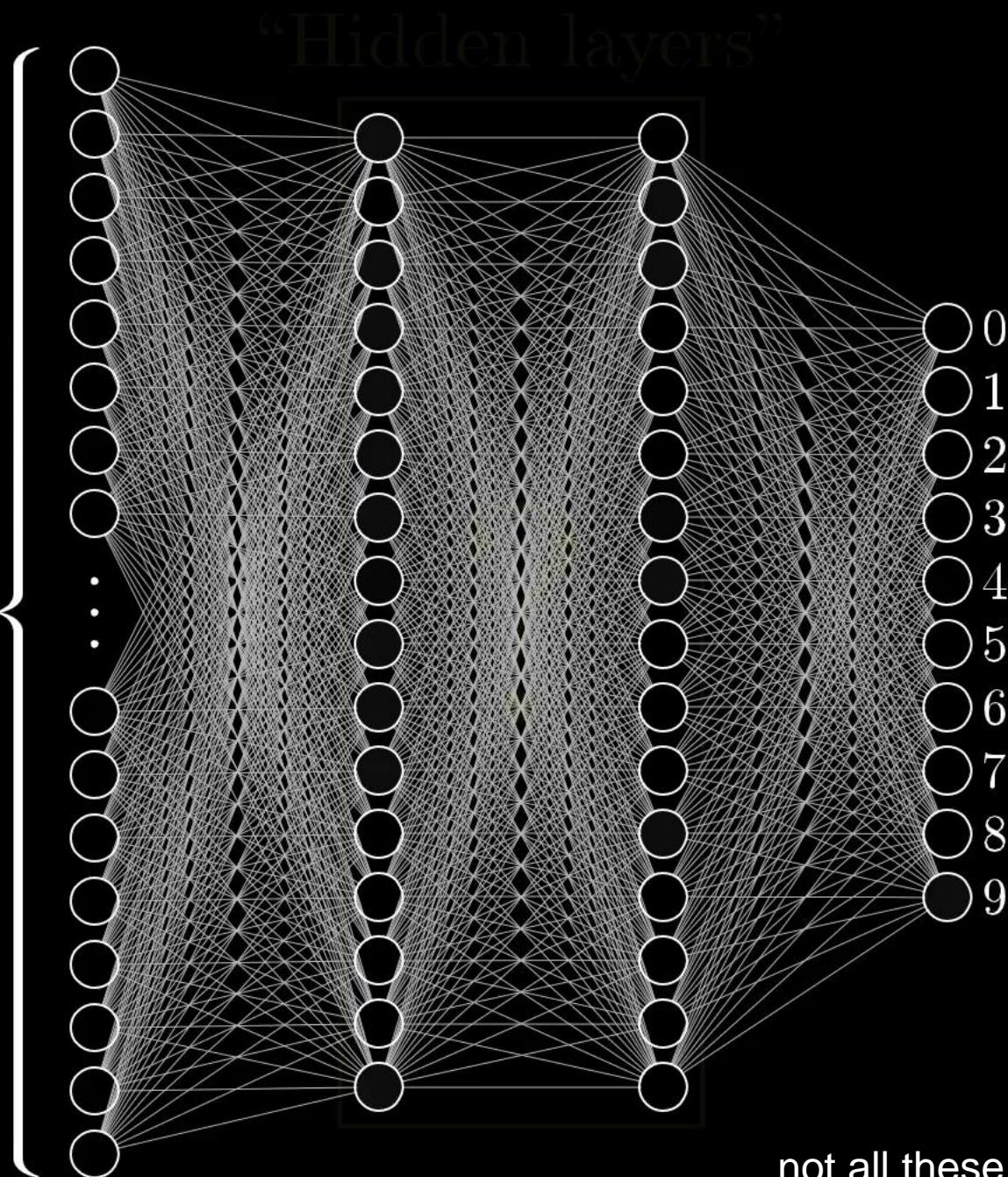


784





784

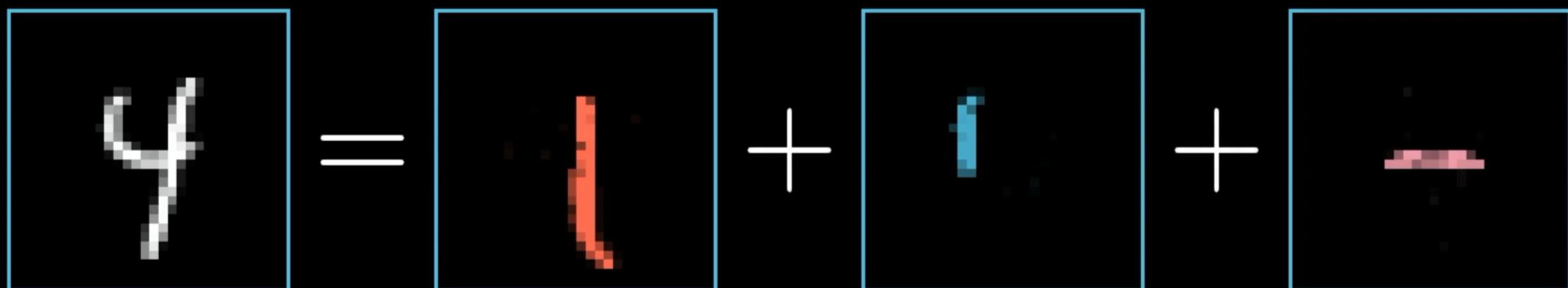


“Hidden layers”

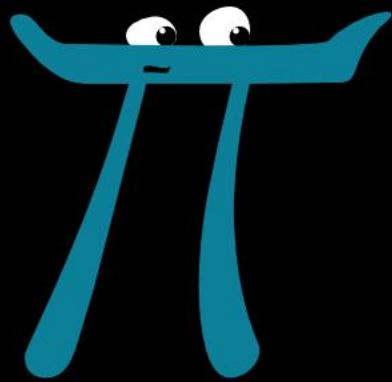
not all these connections are equal



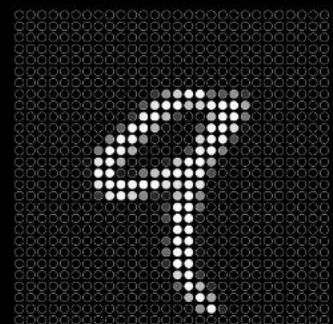
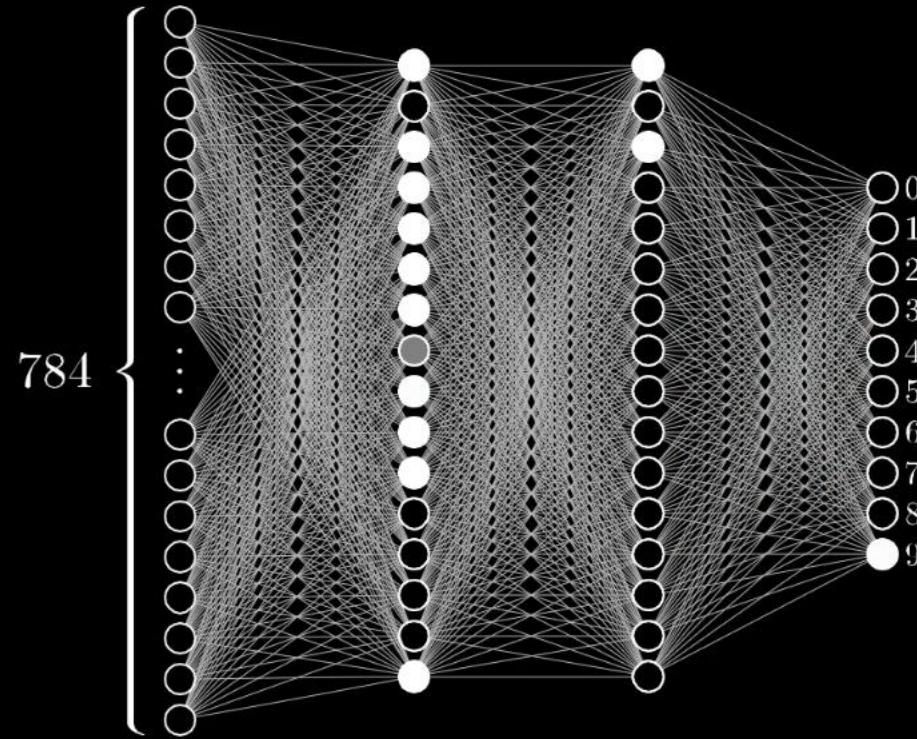
What's the best
hope for what
those middle
layers are
doing?



Does the network
actually do this?



Could each neuron in the second layer of the network corresponds to some little edge?



Pixels



Edges

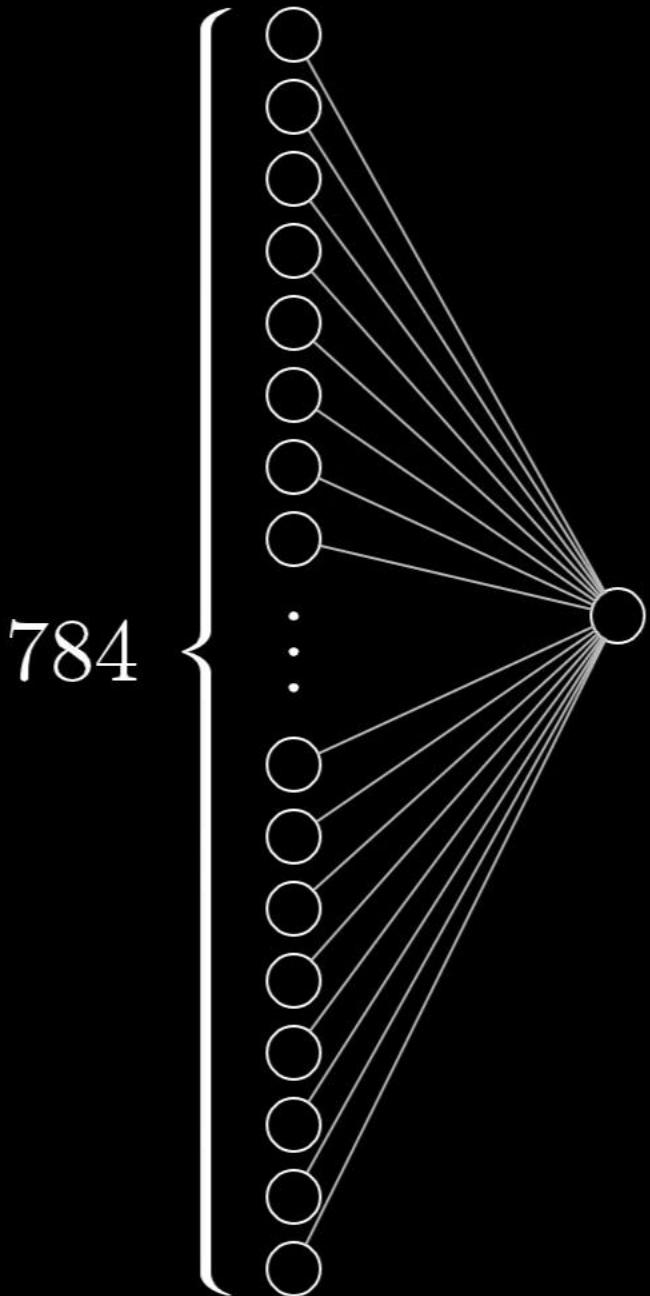


Shapes



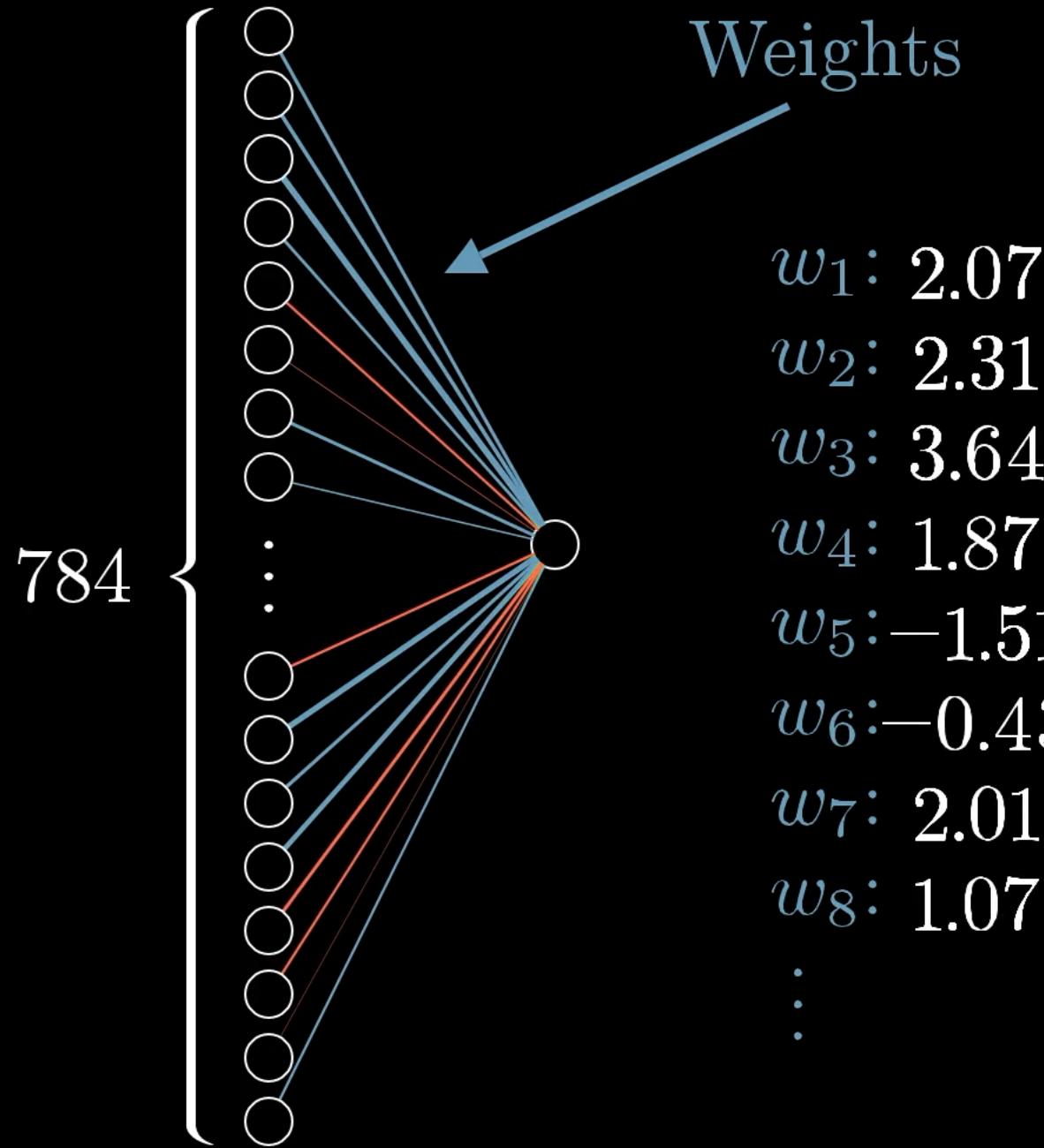
Digits

“The layered structure of the neural network is great because it allows you to break down difficult problems into bite-size steps, so that moving from one layer to the next is relatively straightforward.”

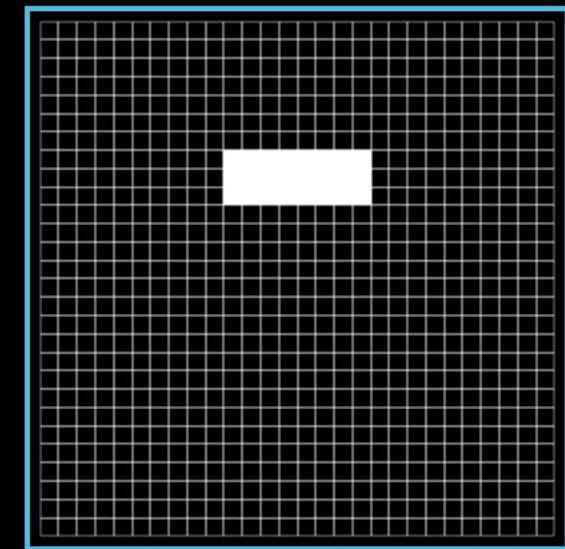


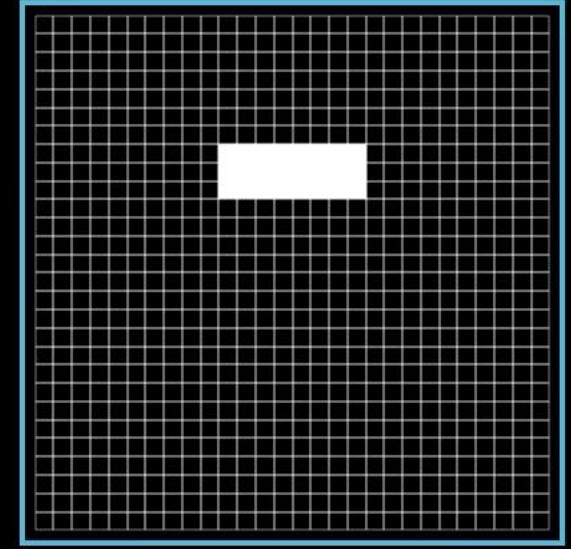
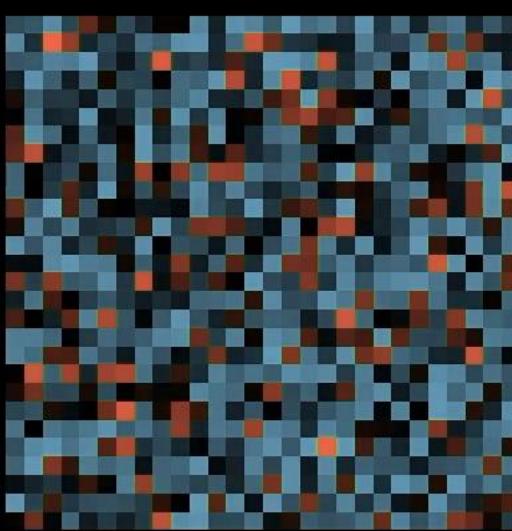
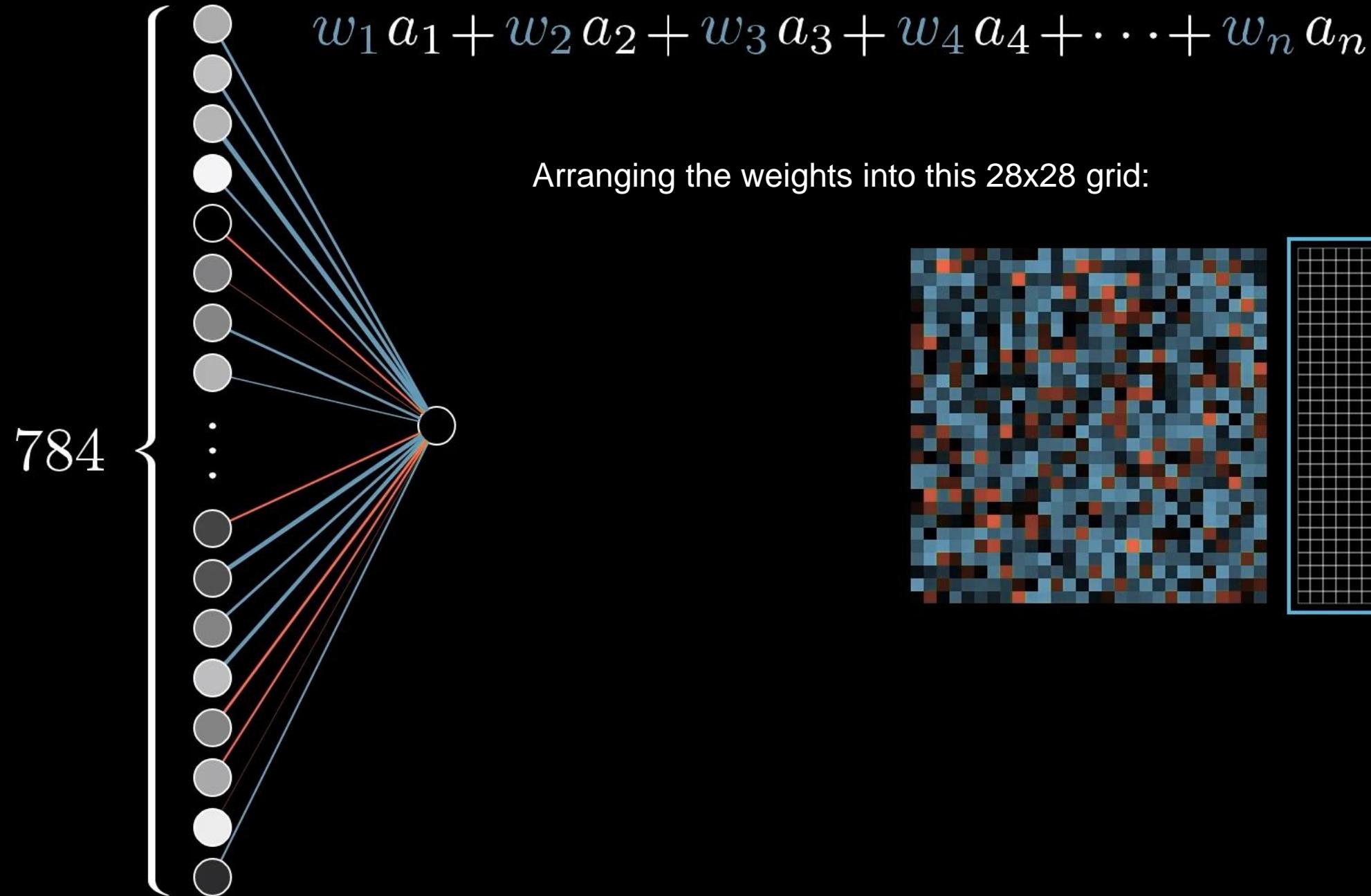
How Information Passes Between Layers

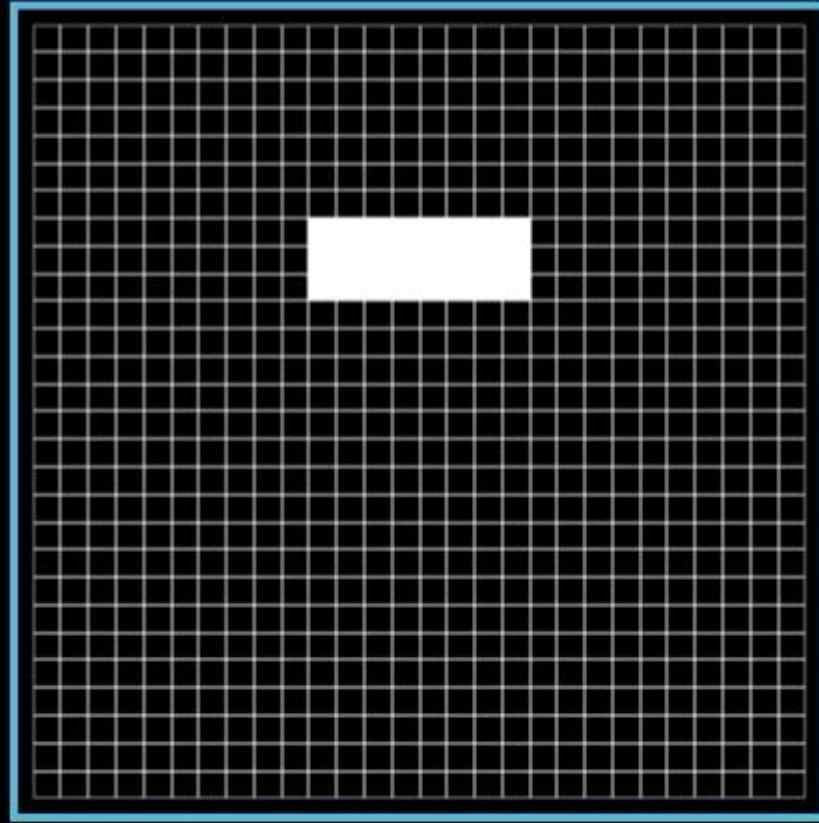
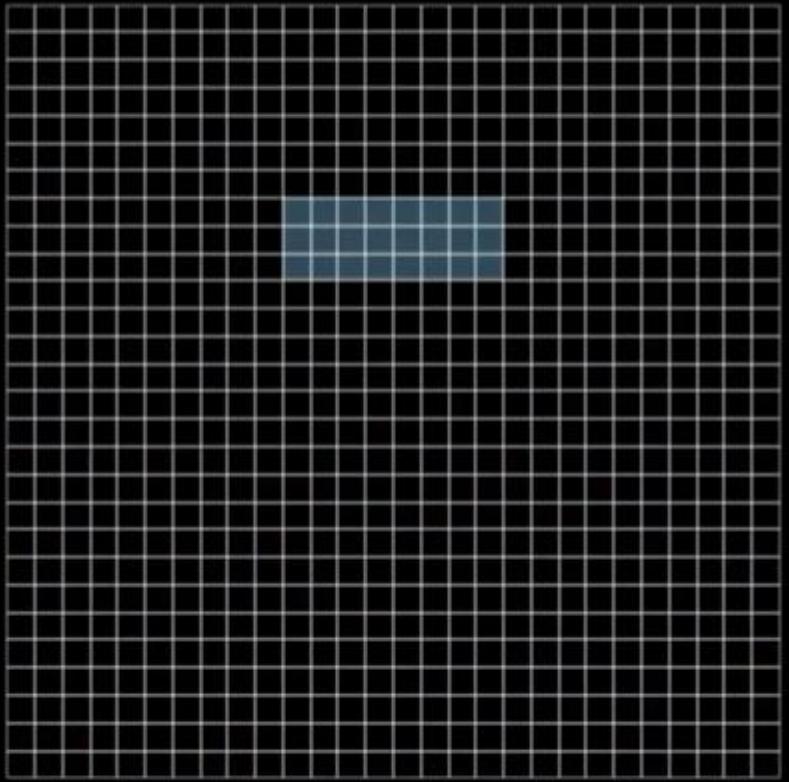
We want this one, specific neuron in the second layer to pick up on whether the image contains this one, specific edge.



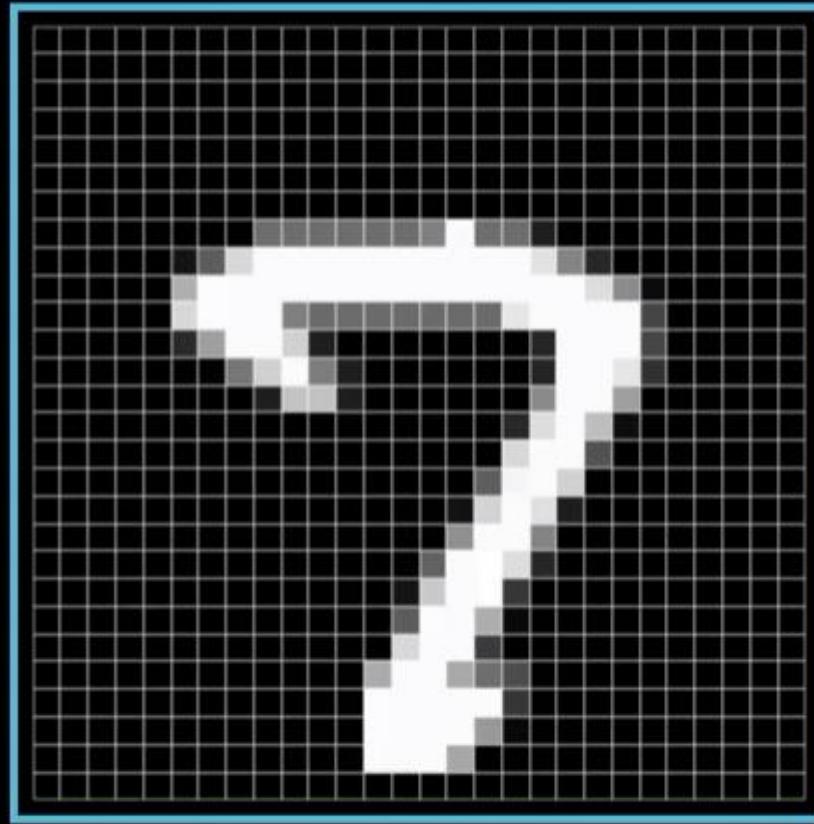
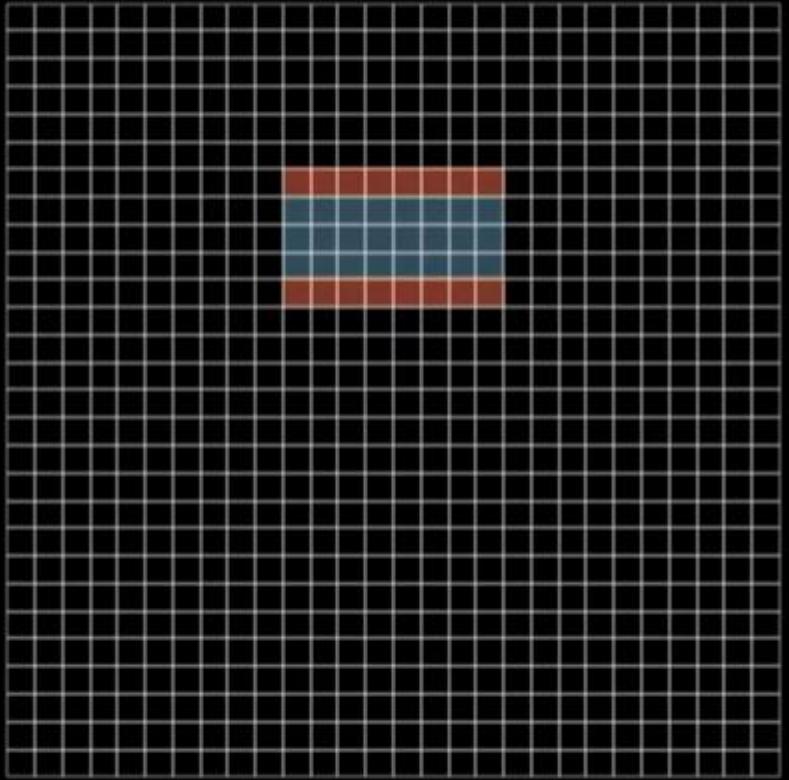
Each weight is an indication of how its neuron in the first layer is correlated with this new neuron in the second layer.



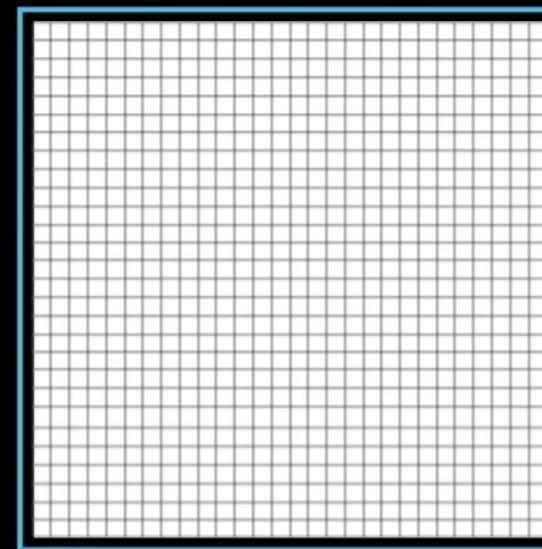
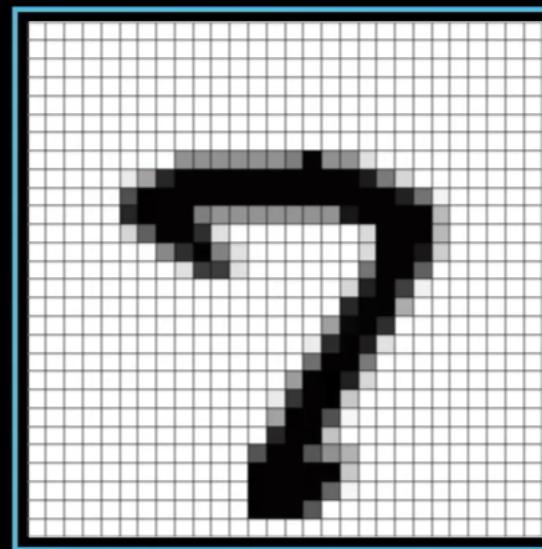
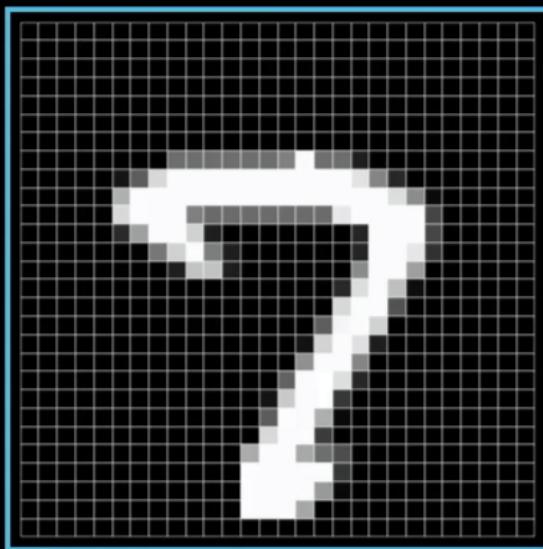
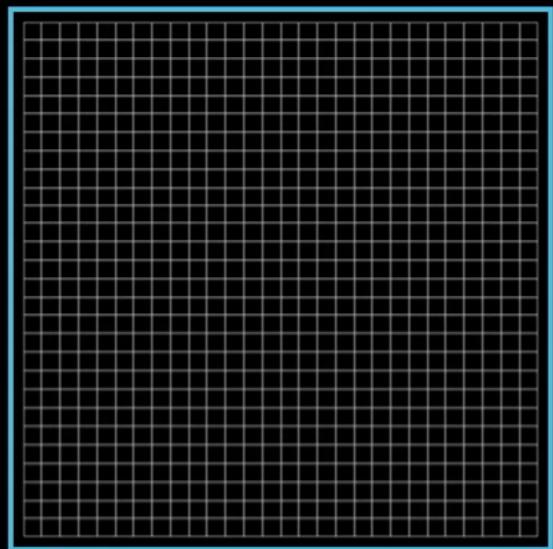
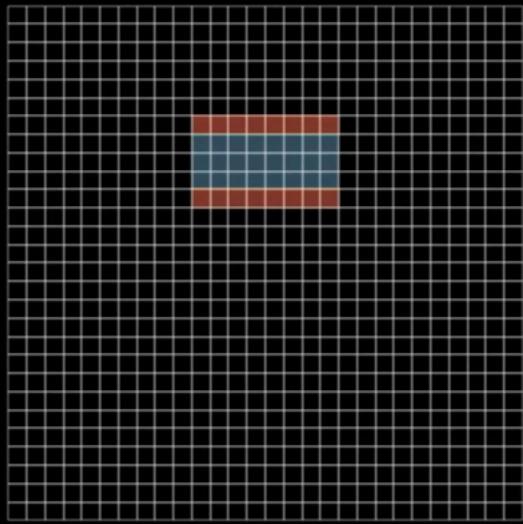




With these weights, the neuron in the second layer will be more activated when pixels in this region are more activated.



To really pick up on whether or not this is an edge, you might want to have some negative weights associated with the surrounding pixels.



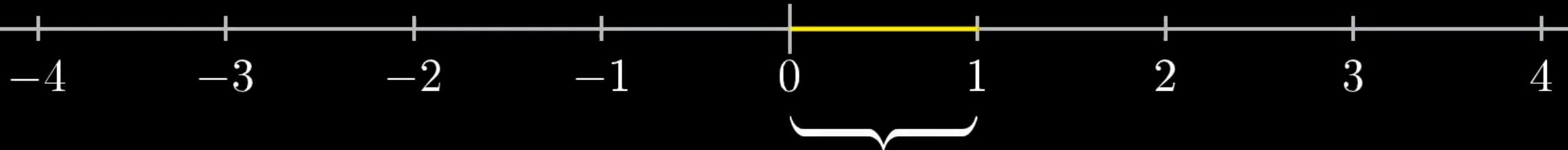
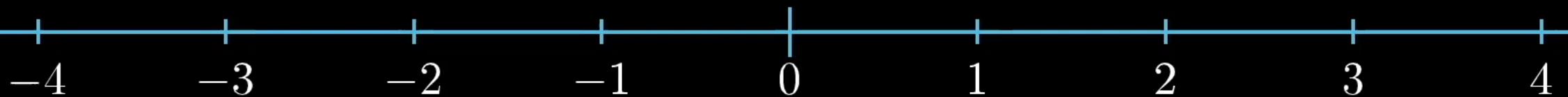
A

B

C

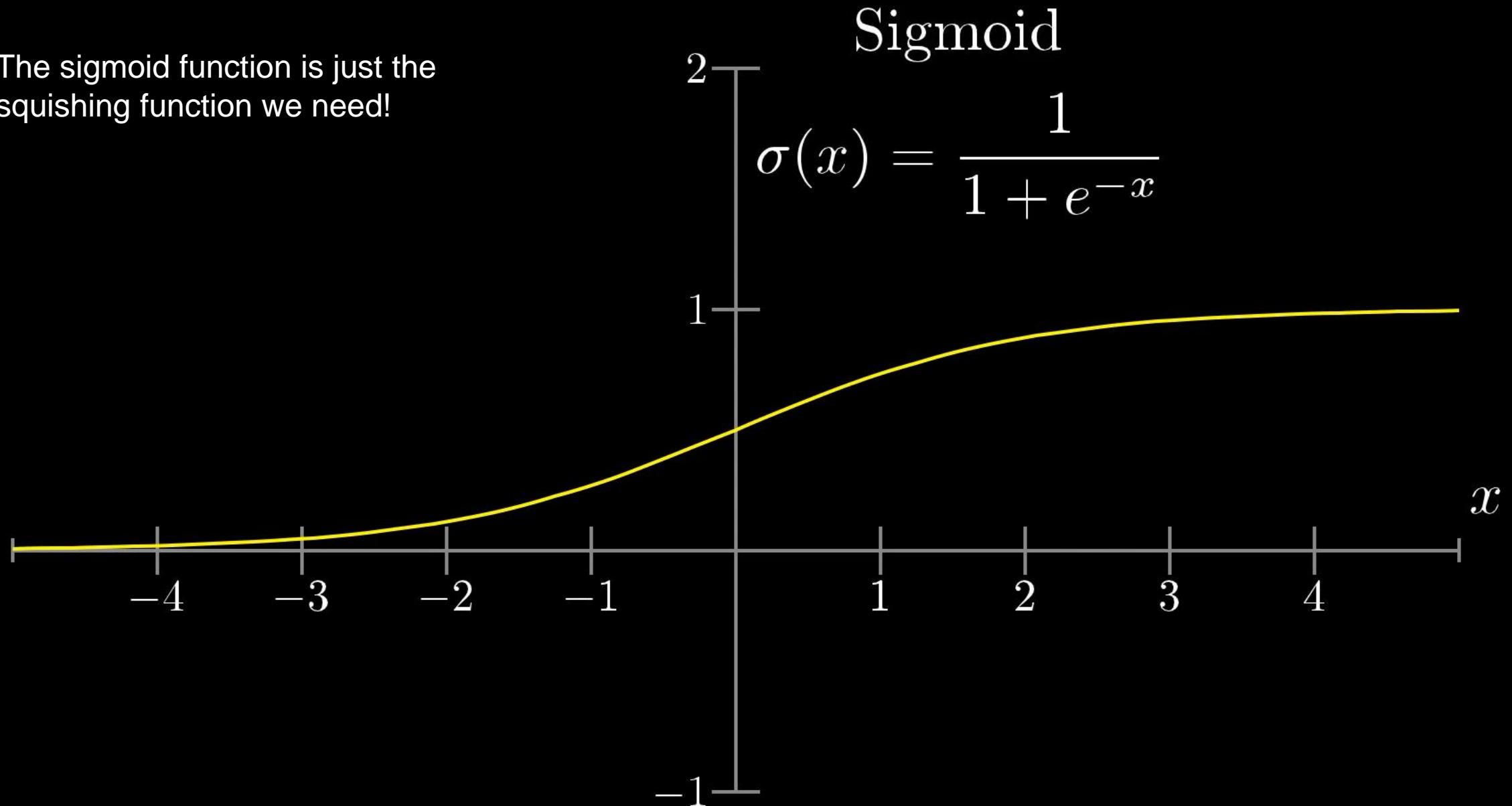
D

$$w_1a_1 + w_2a_2 + w_3a_3 + w_4a_4 + \cdots + w_na_n$$



Activations should be in this range

The sigmoid function is just the squishing function we need!

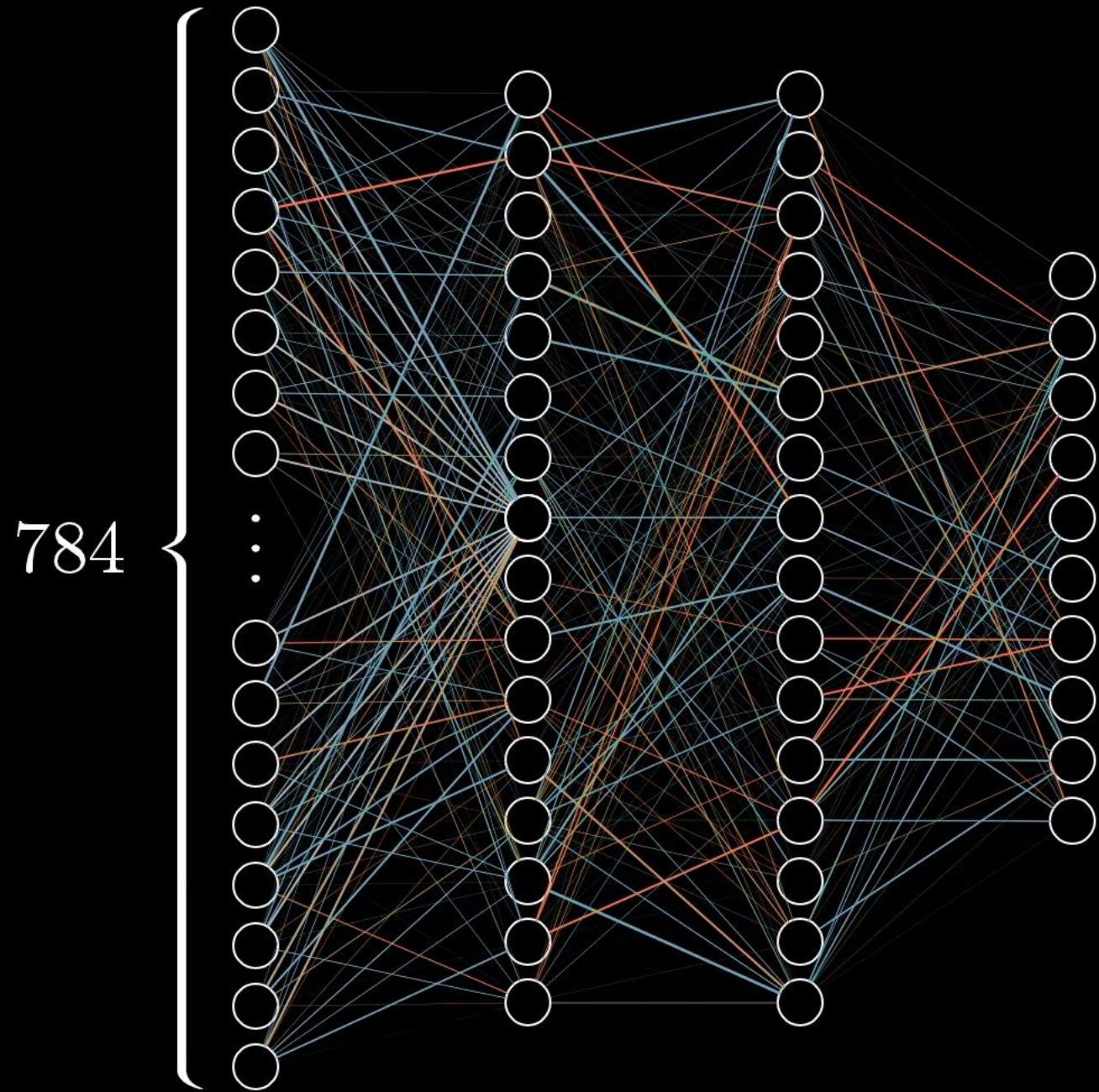


Maybe we only want it to be meaningfully active when that sum is bigger than, say, 10.

$$\sigma(w_1 a_1 + w_2 a_2 + w_3 a_3 + \cdots + w_n a_n \boxed{-10})$$

“bias”

Only activate meaningfully
when weighted sum > 10



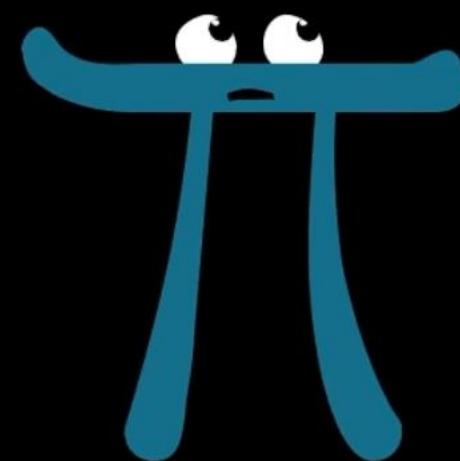
$784 \times 16 + 16 \times 16 + 16 \times 10$

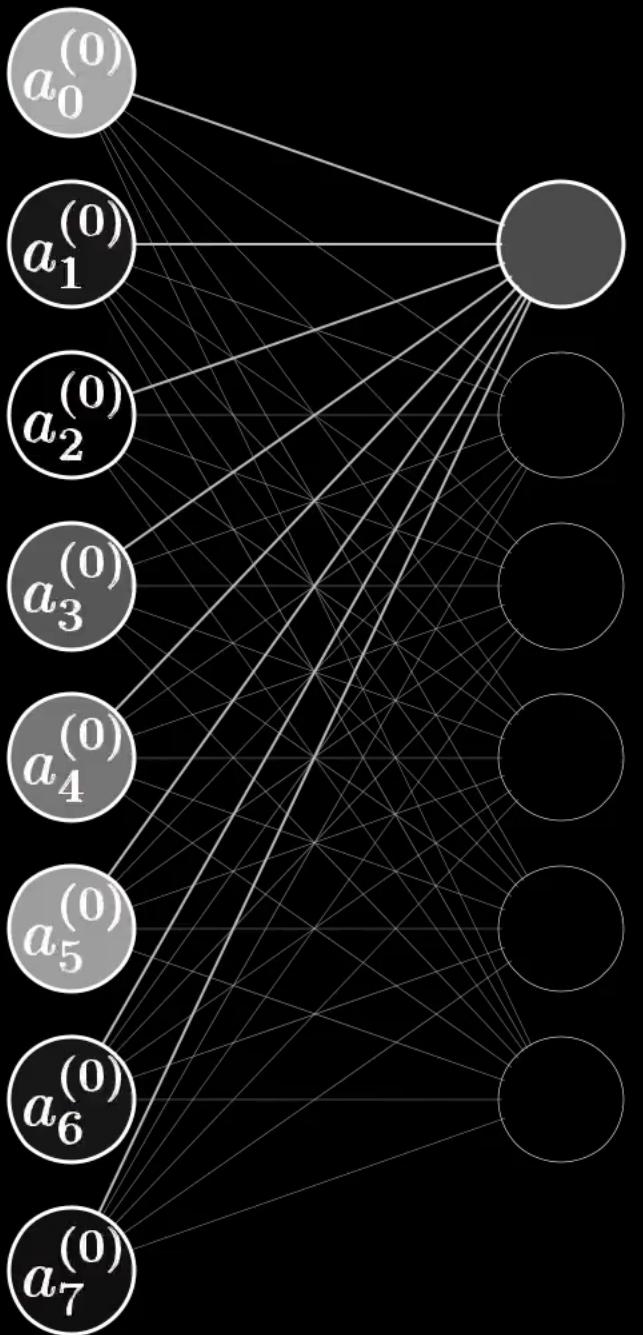
weights

16 + 16 + 10

biases

13,002

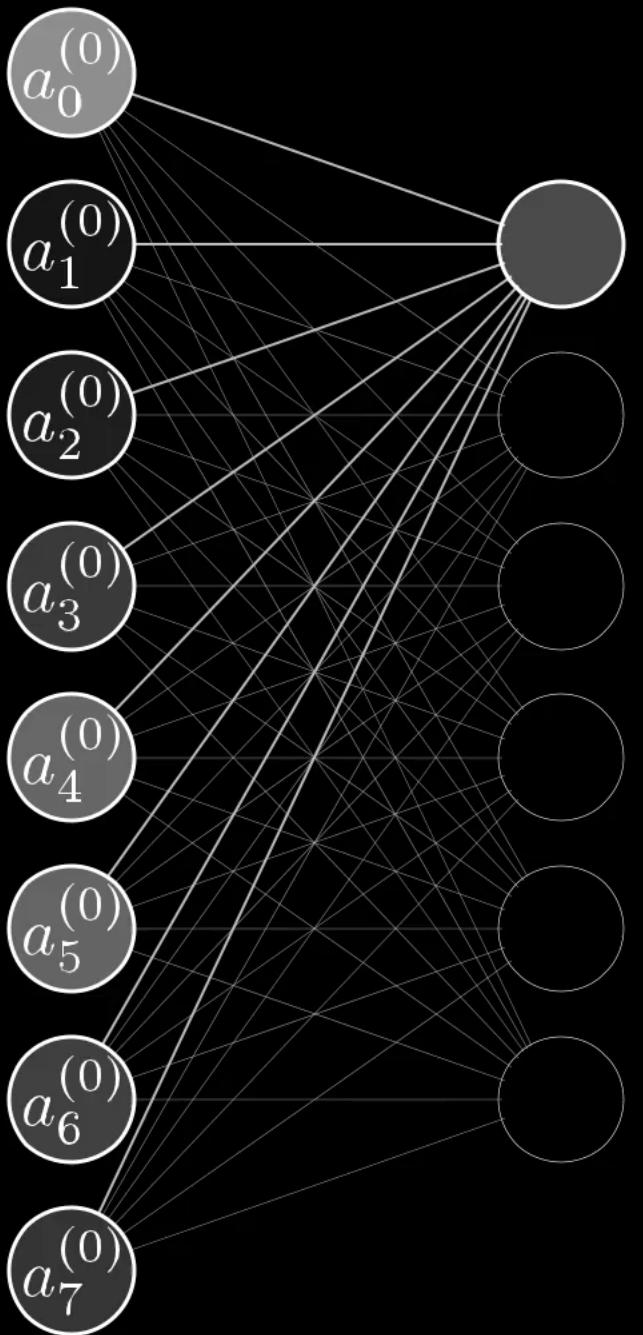




Sigmoid

$$a_0^{(1)} = \sigma \left(w_{0,0} a_0^{(0)} + w_{0,1} a_1^{(0)} + \cdots + w_{0,n} a_n^{(0)} + b_0 \right)$$

Bias

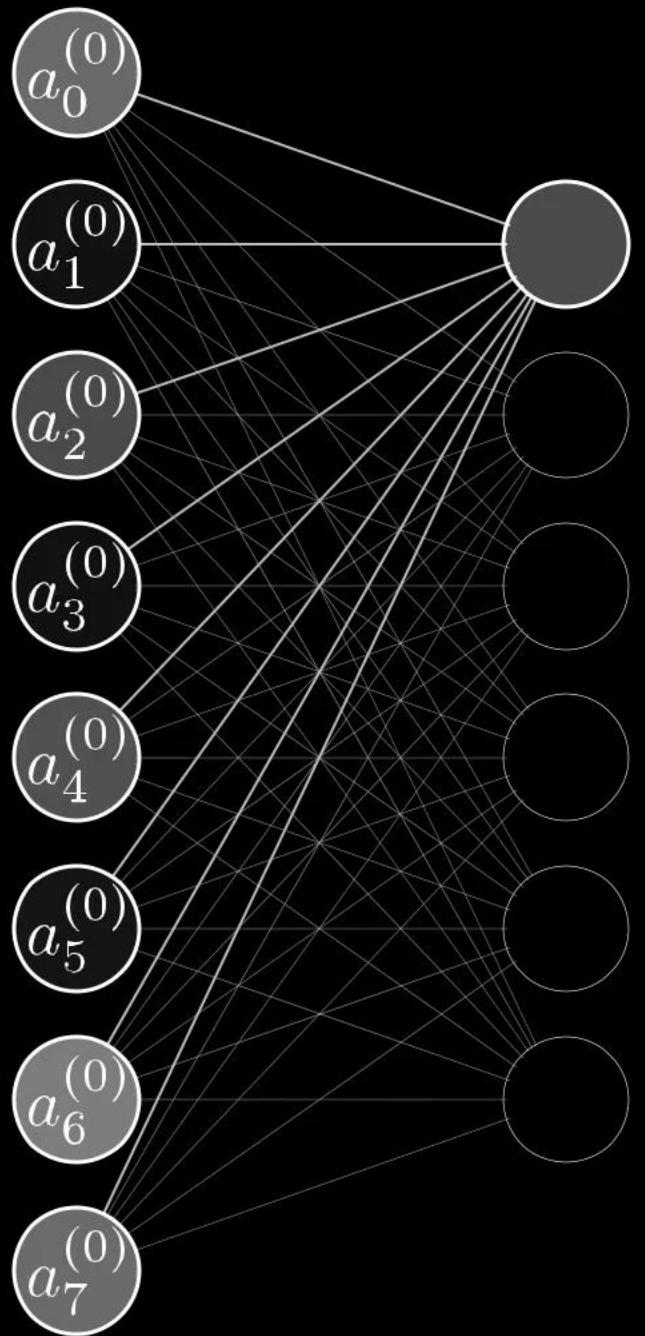


Sigmoid

$$a_0^{(1)} = \sigma \left(w_{0,0} a_0^{(0)} + w_{0,1} a_1^{(0)} + \cdots + w_{0,n} a_n^{(0)} + b_0 \right)$$

Bias

$$\begin{bmatrix} a_0^{(0)} \\ a_1^{(0)} \\ \vdots \\ a_n^{(0)} \end{bmatrix}$$

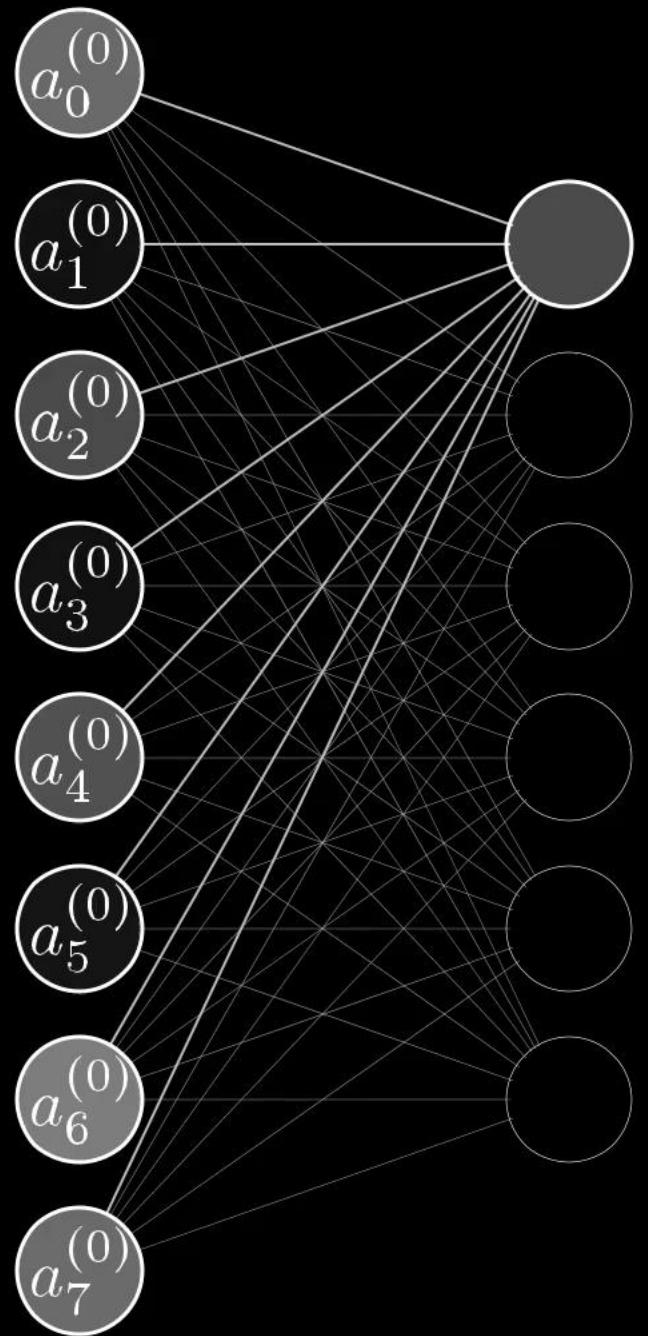


Sigmoid

$$a_0^{(1)} = \sigma \left(w_{0,0} a_0^{(0)} + w_{0,1} a_1^{(0)} + \cdots + w_{0,n} a_n^{(0)} + b_0 \right)$$

↑
Bias

$$\begin{bmatrix} w_{0,0} & w_{0,1} & \dots & w_{0,n} \\ w_{1,0} & w_{1,1} & \dots & w_{1,n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{k,0} & w_{k,1} & \dots & w_{k,n} \end{bmatrix} \begin{bmatrix} a_0^{(0)} \\ a_1^{(0)} \\ \vdots \\ a_n^{(0)} \end{bmatrix}$$

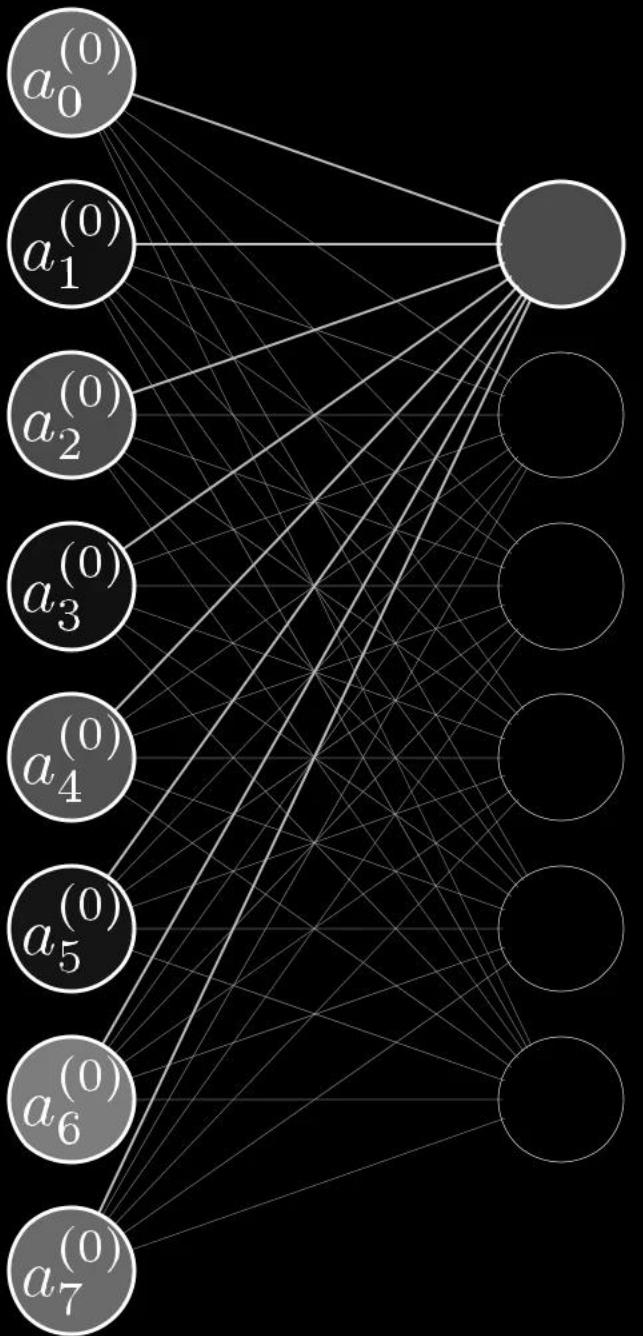


Sigmoid

$$a_0^{(1)} = \sigma \left(w_{0,0} a_0^{(0)} + w_{0,1} a_1^{(0)} + \cdots + w_{0,n} a_n^{(0)} + b_0 \right)$$

Bias

$$\begin{bmatrix} w_{0,0} & w_{0,1} & \cdots & w_{0,n} \\ w_{1,0} & w_{1,1} & \cdots & w_{1,n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{k,0} & w_{k,1} & \cdots & w_{k,n} \end{bmatrix} \begin{bmatrix} a_0^{(0)} \\ a_1^{(0)} \\ \vdots \\ a_n^{(0)} \end{bmatrix}$$

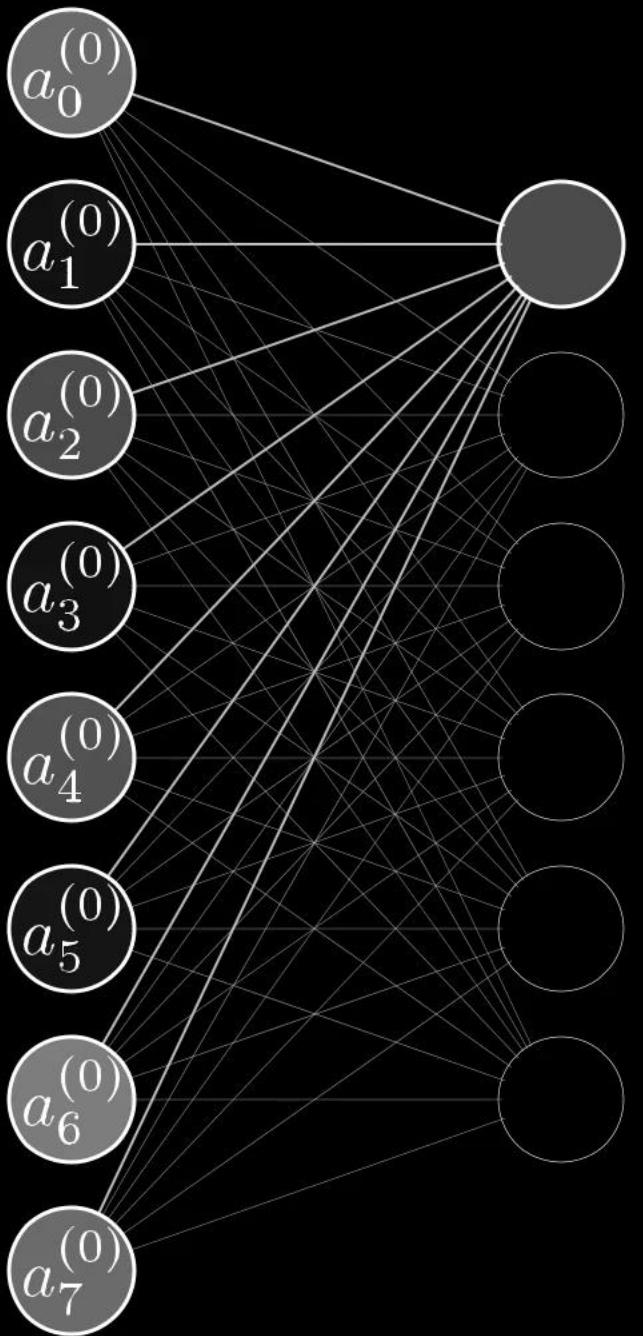


Sigmoid

$$a_0^{(1)} = \sigma \left(w_{0,0} a_0^{(0)} + w_{0,1} a_1^{(0)} + \cdots + w_{0,n} a_n^{(0)} + b_0 \right)$$

↑
Bias

$$\begin{bmatrix} w_{0,0} & w_{0,1} & \cdots & w_{0,n} \\ w_{1,0} & w_{1,1} & \cdots & w_{1,n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{k,0} & w_{k,1} & \cdots & w_{k,n} \end{bmatrix} \begin{bmatrix} a_0^{(0)} \\ a_1^{(0)} \\ \vdots \\ a_n^{(0)} \end{bmatrix} + \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_n \end{bmatrix}$$



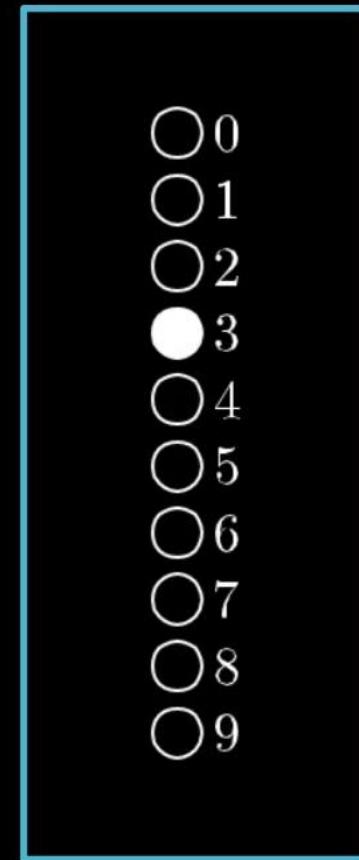
$$\sigma \left(\begin{bmatrix} w_{0,0} & w_{0,1} & \dots & w_{0,n} \\ w_{1,0} & w_{1,1} & \dots & w_{1,n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{k,0} & w_{k,1} & \dots & w_{k,n} \end{bmatrix} \begin{bmatrix} a_0^{(0)} \\ a_1^{(0)} \\ \vdots \\ a_n^{(0)} \end{bmatrix} + \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_n \end{bmatrix} \right)$$

Neural network function



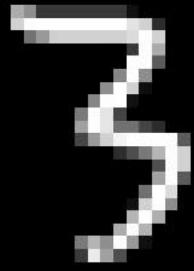
784 inputs

the entire network is just a **function!** It takes in 784 numbers as its input, and spits out 10 numbers as its output. It's an absurdly complicated function, because it takes over 13,000 parameters (weights and biases), and it involves iterating many matrix-vector products and sigmoid squishifications together. But it's just a function nonetheless.



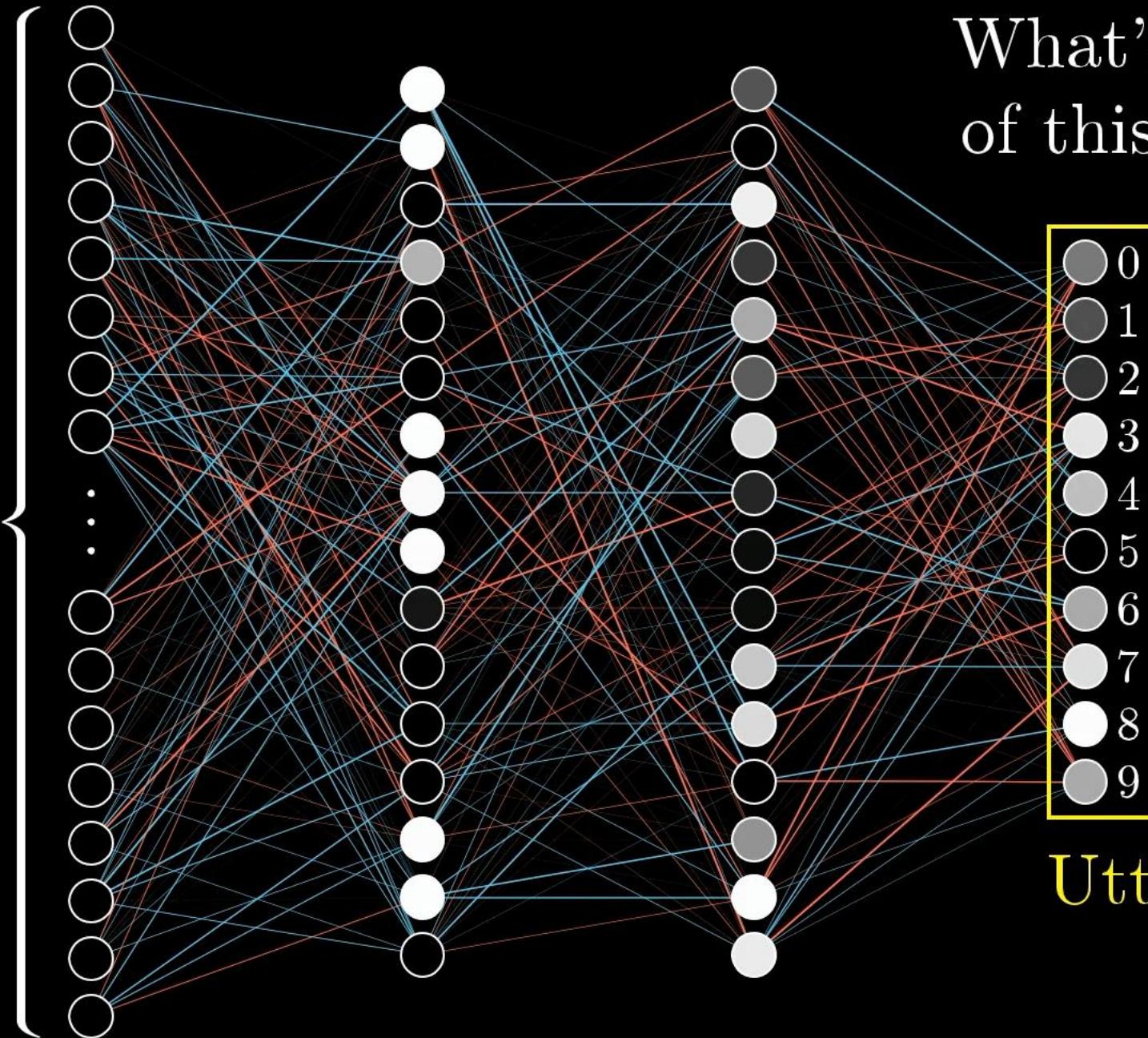
10 outputs

$(\boxed{0}, 0)(\boxed{6}, 6)(\boxed{3}, 3)(\boxed{6}, 6)(\boxed{7}, 7)(\boxed{8}, 8)(\boxed{0}, 0)(\boxed{9}, 9)$
 $(\boxed{5}, 5)(\boxed{4}, 4)(\boxed{3}, 3)(\boxed{4}, 6)(\boxed{5}, 5)(\boxed{8}, 8)(\boxed{9}, 9)(\boxed{5}, 5)$
 $(\boxed{4}, 4)(\boxed{4}, 4)(\boxed{7}, 7)(\boxed{2}, 2)(\boxed{0}, 0)(\boxed{3}, 3)(\boxed{2}, 2)(\boxed{8}, 8)$
 $(\boxed{9}, 9)(\boxed{1}, 1)(\boxed{9}, 9)(\boxed{2}, 2)(\boxed{2}, 2)(\boxed{7}, 7)(\boxed{9}, 9)(\boxed{4}, 4)$
 $(\boxed{8}, 8)(\boxed{7}, \text{The network learns by looking at examples of correctly identified digits.})(\boxed{5}, 5)(\boxed{3}, 3)$
 $(\boxed{2}, 2)(\boxed{3}, \text{The network learns by looking at examples of correctly identified digits.})(\boxed{1}, 1)(\boxed{5}, 5)$
 $(\boxed{8}, 8)(\boxed{4}, 4)(\boxed{1}, 7)(\boxed{7}, 7)(\boxed{4}, 4)(\boxed{4}, 4)(\boxed{4}, 4)(\boxed{4}, 4)(\boxed{2}, 2)$
 $(\boxed{0}, 0)(\boxed{7}, 7)(\boxed{2}, 2)(\boxed{4}, 4)(\boxed{8}, 8)(\boxed{2}, 2)(\boxed{6}, 6)(\boxed{9}, 9)$
 $(\boxed{1}, 9)(\boxed{2}, 2)(\boxed{8}, 8)(\boxed{7}, 7)(\boxed{6}, 6)(\boxed{1}, 1)(\boxed{1}, 1)(\boxed{2}, 2)$
 $(\boxed{3}, 3)(\boxed{9}, 9)(\boxed{1}, 1)(\boxed{6}, 6)(\boxed{5}, 5)(\boxed{1}, 1)(\boxed{1}, 1)(\boxed{0}, 0)$



784

Initialized with
totally random
weights and
biases, the
network is
terrible at
identifying digits.



What's the “cost”
of this difference?

0
1
2
3
4
5
6
7
8
9



0
1
2
3
4
5
6
7
8
9

Utter trash

What's the “cost”
of this difference?

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



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

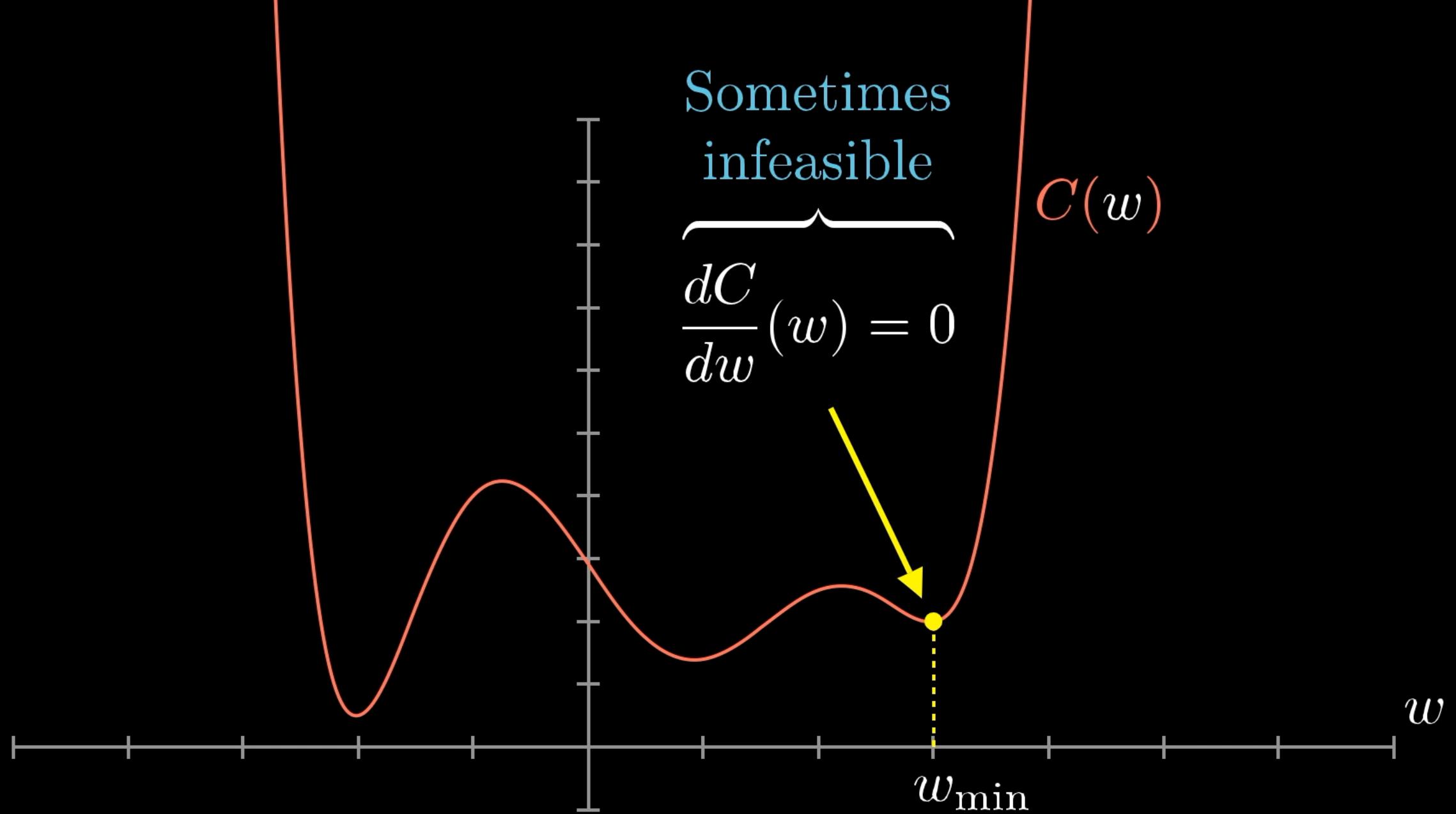
Utter trash

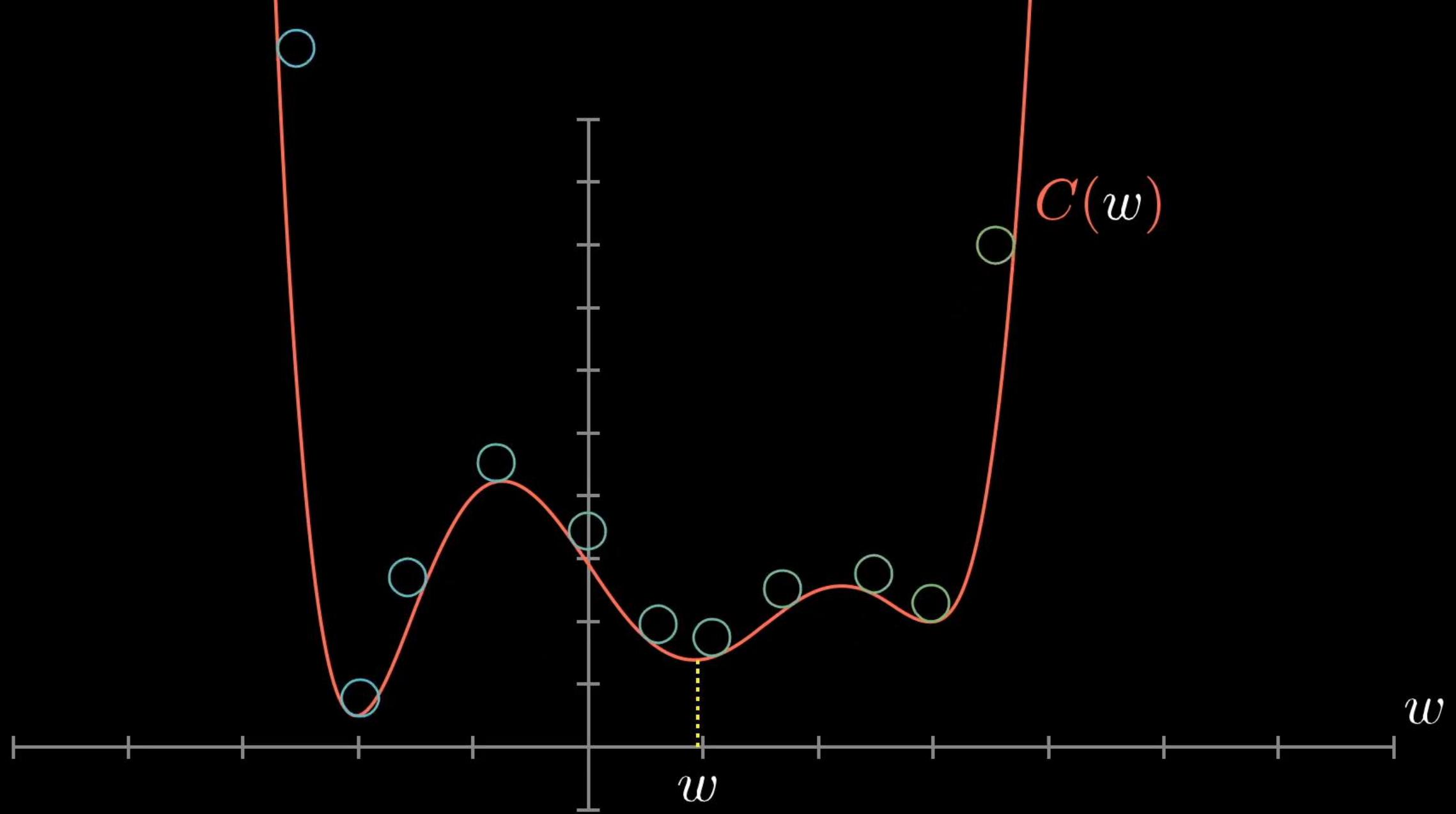
Cost of

3

The "cost" is calculated by
adding up the squares of the
differences between what we
got and what we want.

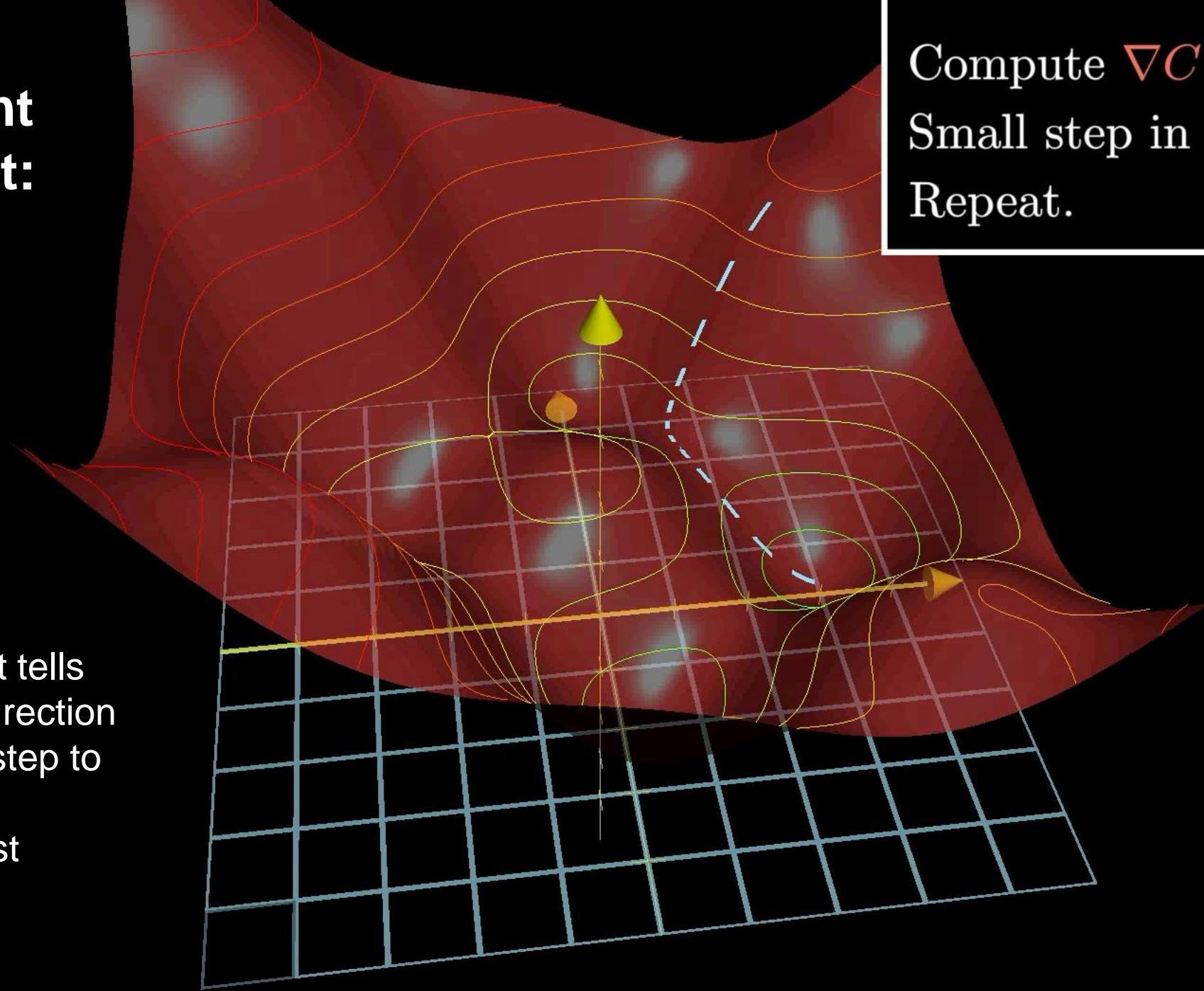
$$\left\{ \begin{array}{l} (0.43 - 0.00)^2 + \\ (0.28 - 0.00)^2 + \\ (0.19 - 0.00)^2 + \\ (0.88 - 1.00)^2 + \\ (0.72 - 0.00)^2 + \\ (0.01 - 0.00)^2 + \\ (0.64 - 0.00)^2 + \\ (0.86 - 0.00)^2 + \\ (0.99 - 0.00)^2 + \\ (0.63 - 0.00)^2 \end{array} \right\}$$





Gradient descent:

The gradient tells you which direction you should step to increase the function most quickly.



Compute ∇C
Small step in $-\nabla C$ direction
Repeat.

Each step: $-\eta \nabla C$; η =learning rate;
Take shorter steps closer to the minima

13,002 weights and biases

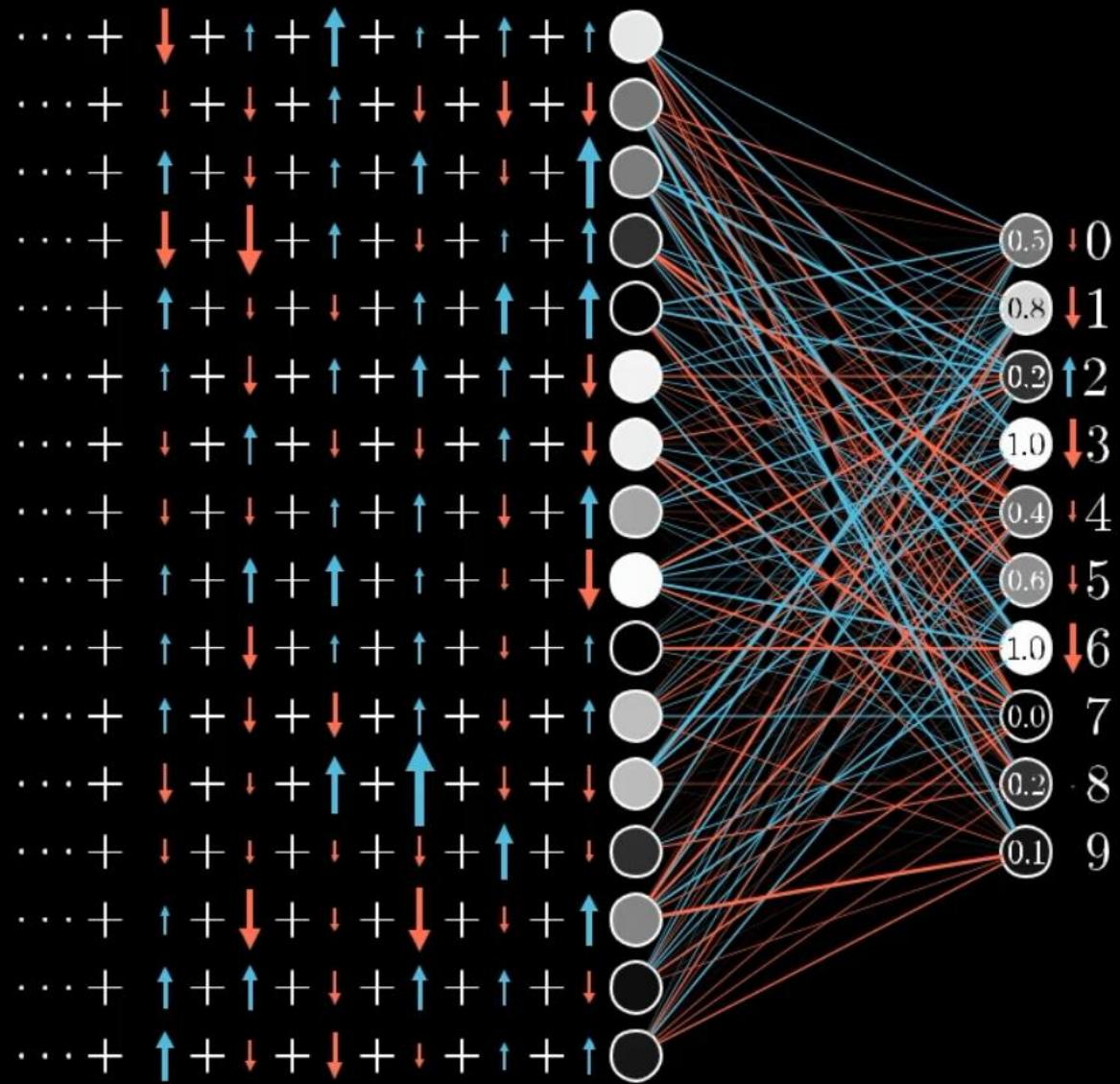
$$\vec{\mathbf{W}} = \begin{bmatrix} 2.25 \\ -1.57 \\ 1.98 \\ \vdots \\ -1.16 \\ 3.82 \\ 1.21 \end{bmatrix} \quad -\nabla C(\vec{\mathbf{W}}) = \begin{bmatrix} 0.18 \\ 0.45 \\ -0.51 \\ \vdots \\ 0.40 \\ -0.32 \\ 0.82 \end{bmatrix}$$

How to nudge all
weights and biases

$$0.2 = \sigma(w_0 a_0 + w_1 a_1 + \dots + w_{n-1} a_{n-1} + b)$$

There are three avenues that can team up together to increase this activation:

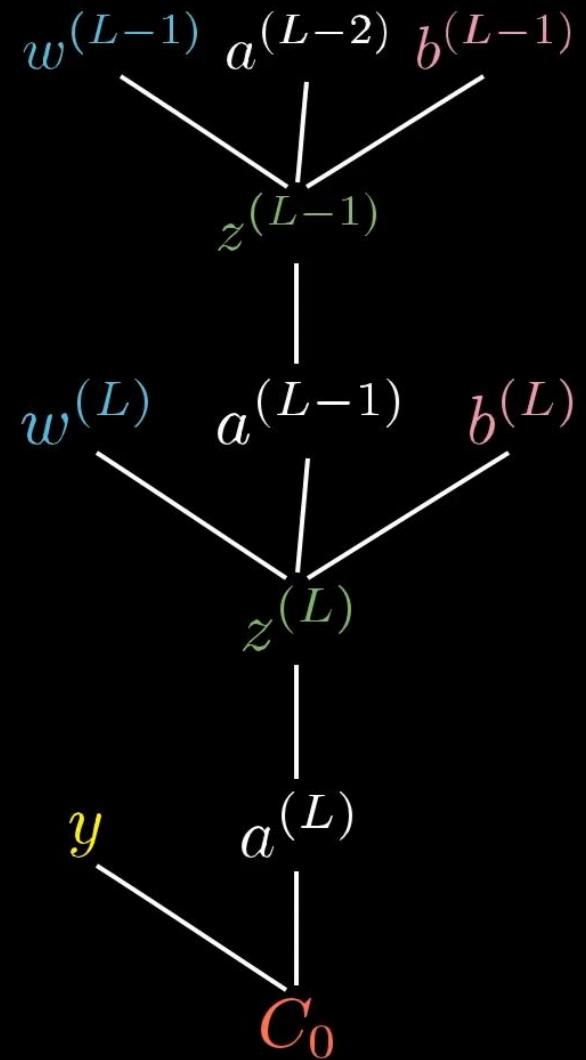
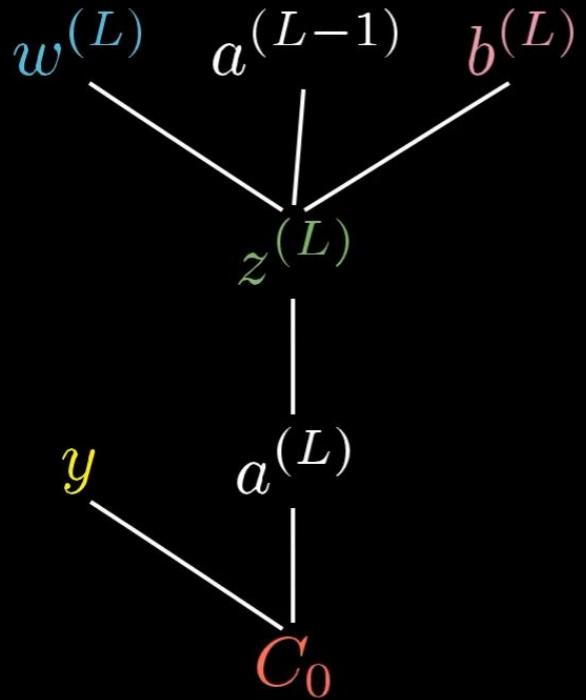
1. Increase the bias
2. Increase the weights
3. Change the activations from the previous layer



It's impossible to perfectly satisfy all these competing desires for activations in the second layer.

The best we can do is add up all the desired nudges to find the overall desired change.

Calculating the Gradient with Backpropagation:



$$\nabla C = \begin{bmatrix} \frac{\partial C}{\partial w^{(1)}} \\ \frac{\partial C}{\partial b^{(1)}} \\ \vdots \\ \frac{\partial C}{\partial w^{(L)}} \\ \frac{\partial C}{\partial b^{(L)}} \end{bmatrix}$$

$$\nabla C \leftarrow \begin{cases} \frac{\partial C}{\partial w_{jk}^{(l)}} = a_k^{(l-1)} \sigma'(z_j^{(l)}) \boxed{\frac{\partial C}{\partial a_j^{(l)}}} \\ \sum_{j=0}^{n_{l+1}-1} w_{jk}^{(l+1)} \sigma'(z_j^{(l+1)}) \frac{\partial C}{\partial a_j^{(l+1)}} \\ \text{or} \\ 2(a_j^{(L)} - y_j) \end{cases}$$

