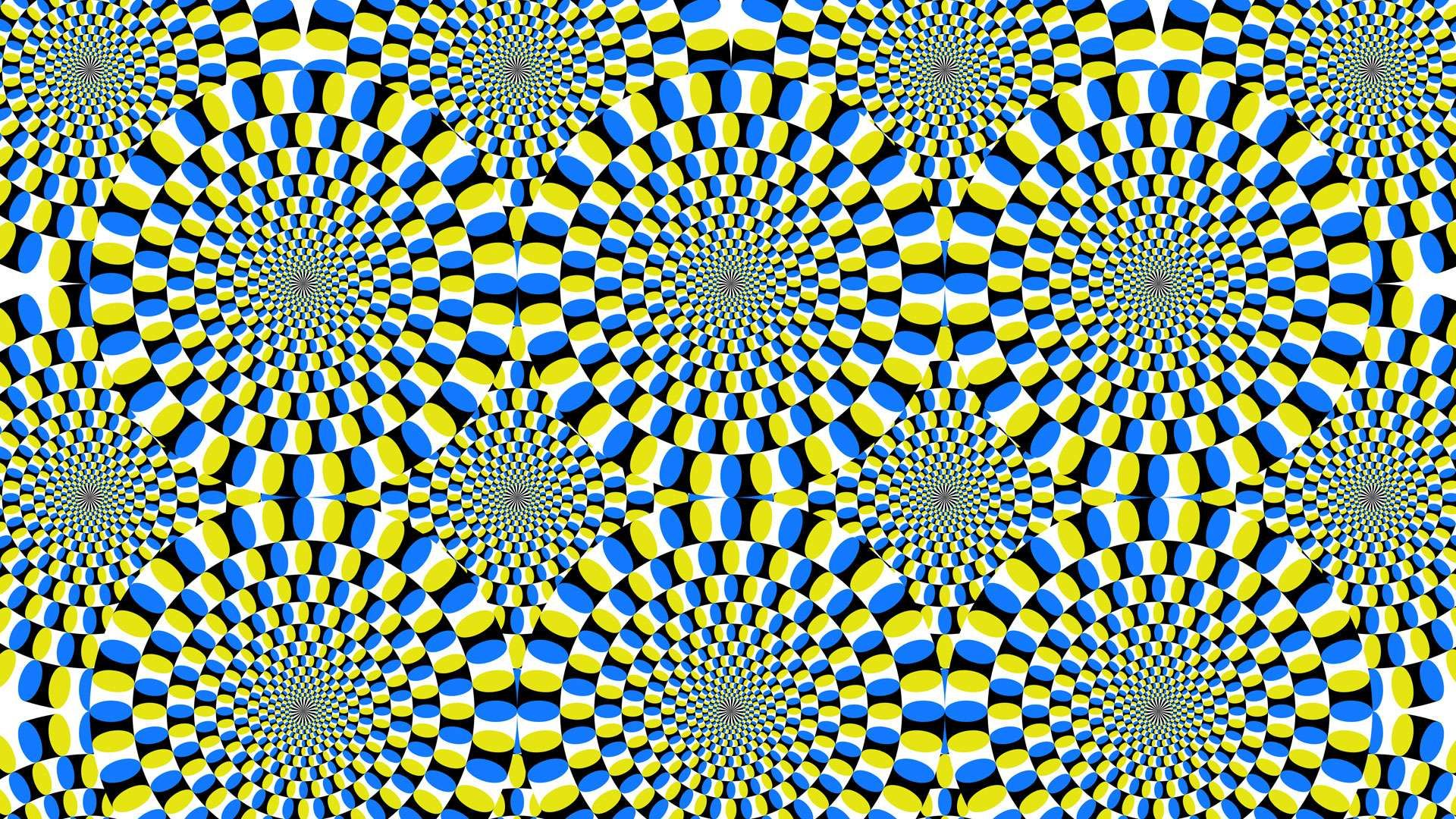


Adversarial Attacks in the Real World





Mathis Hammel

 @MathisHammel

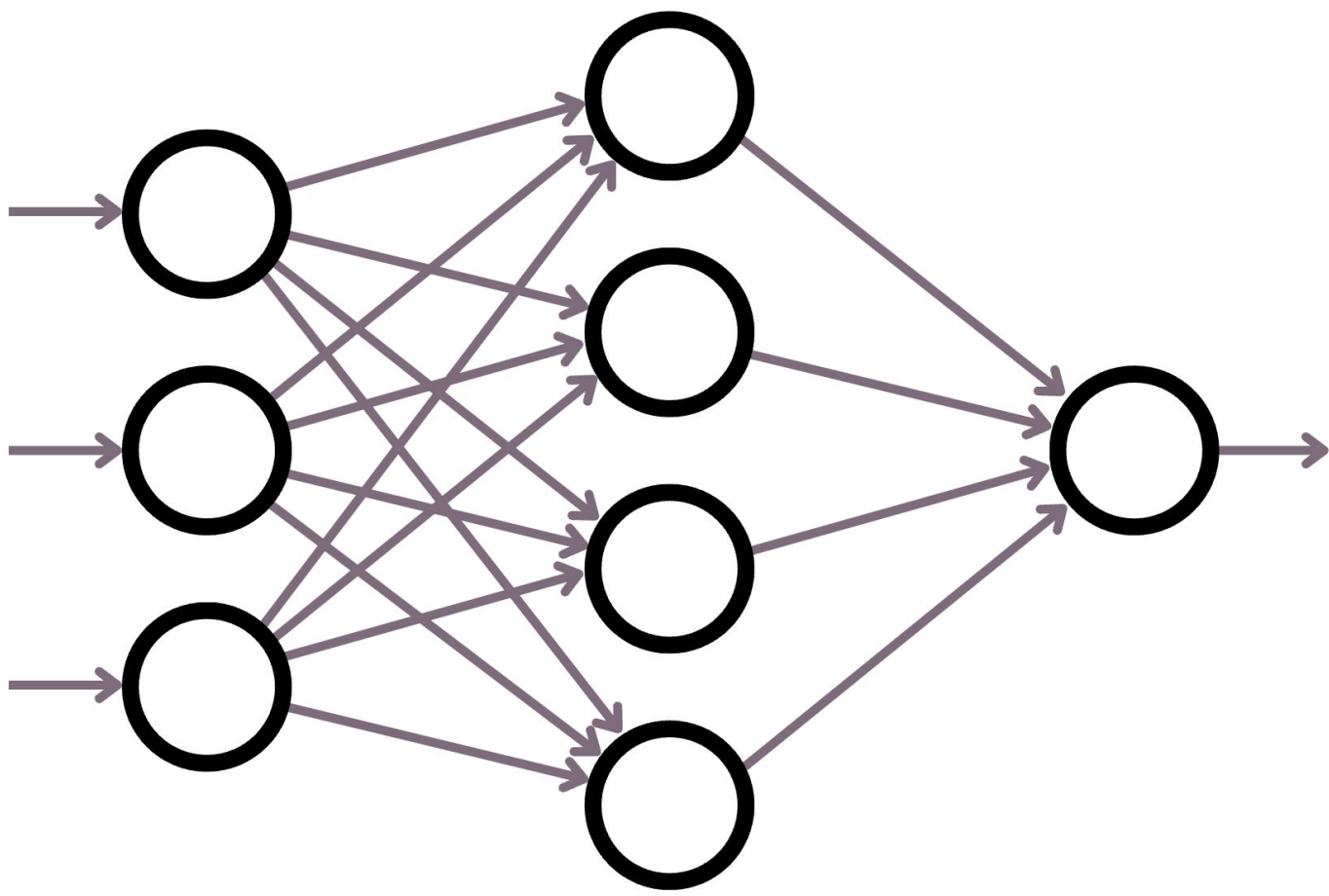
Independent consultant

- Technical competitions
- Algorithms
- Cybersecurity R&D

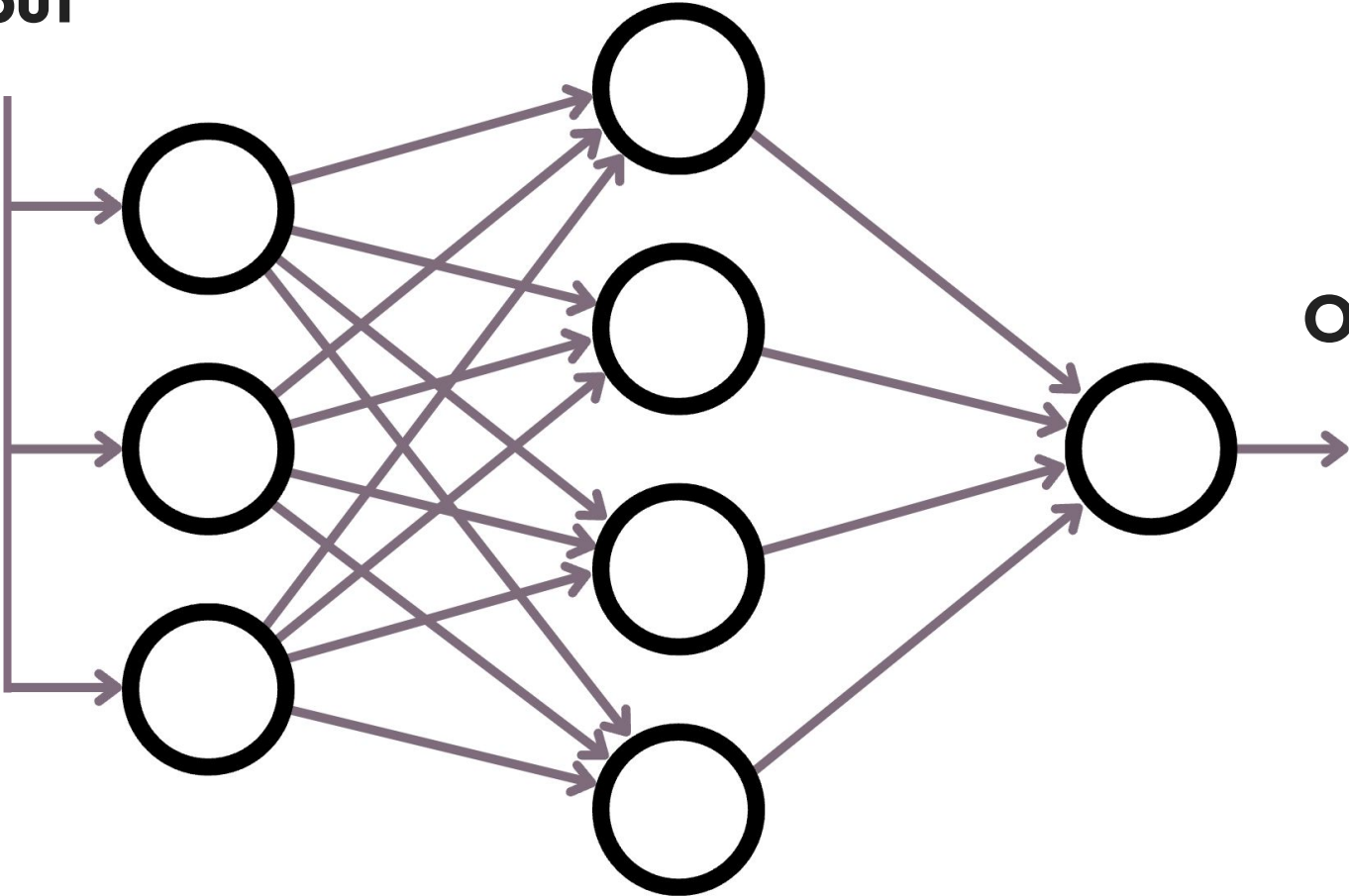


Let's start with the basics.

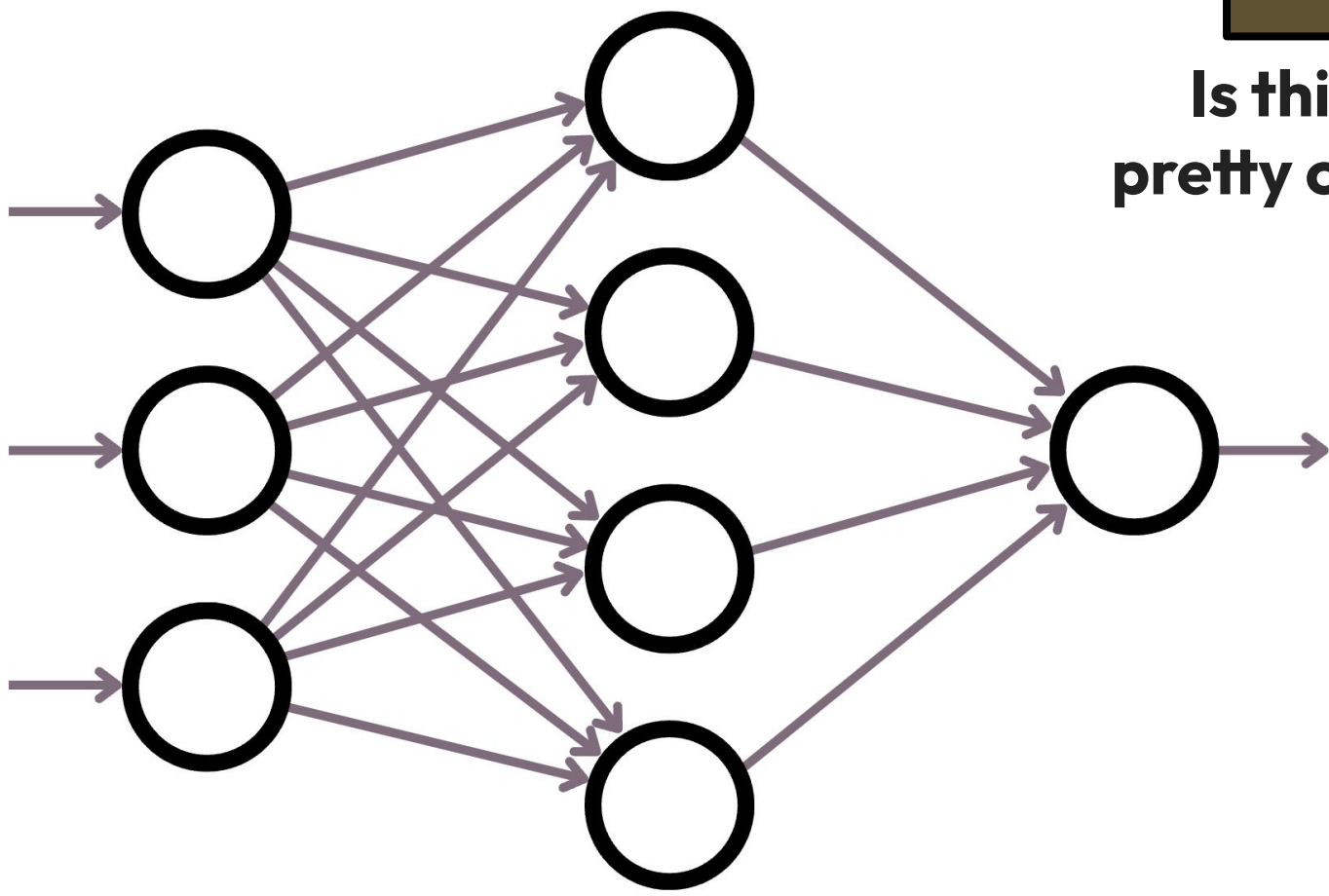
**What *really* is
a neural network?**



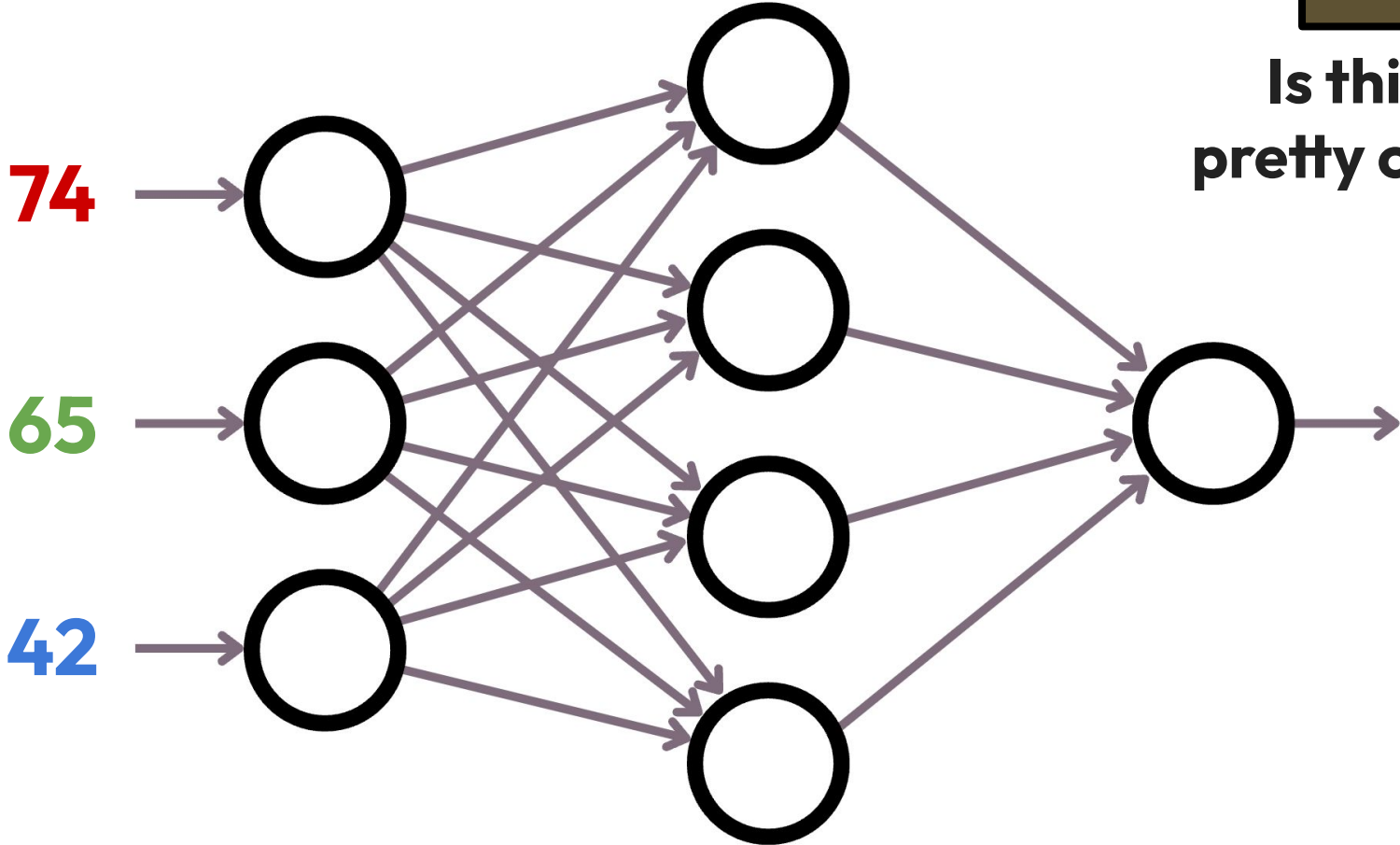
Input



Output



**Is this a
pretty color?**

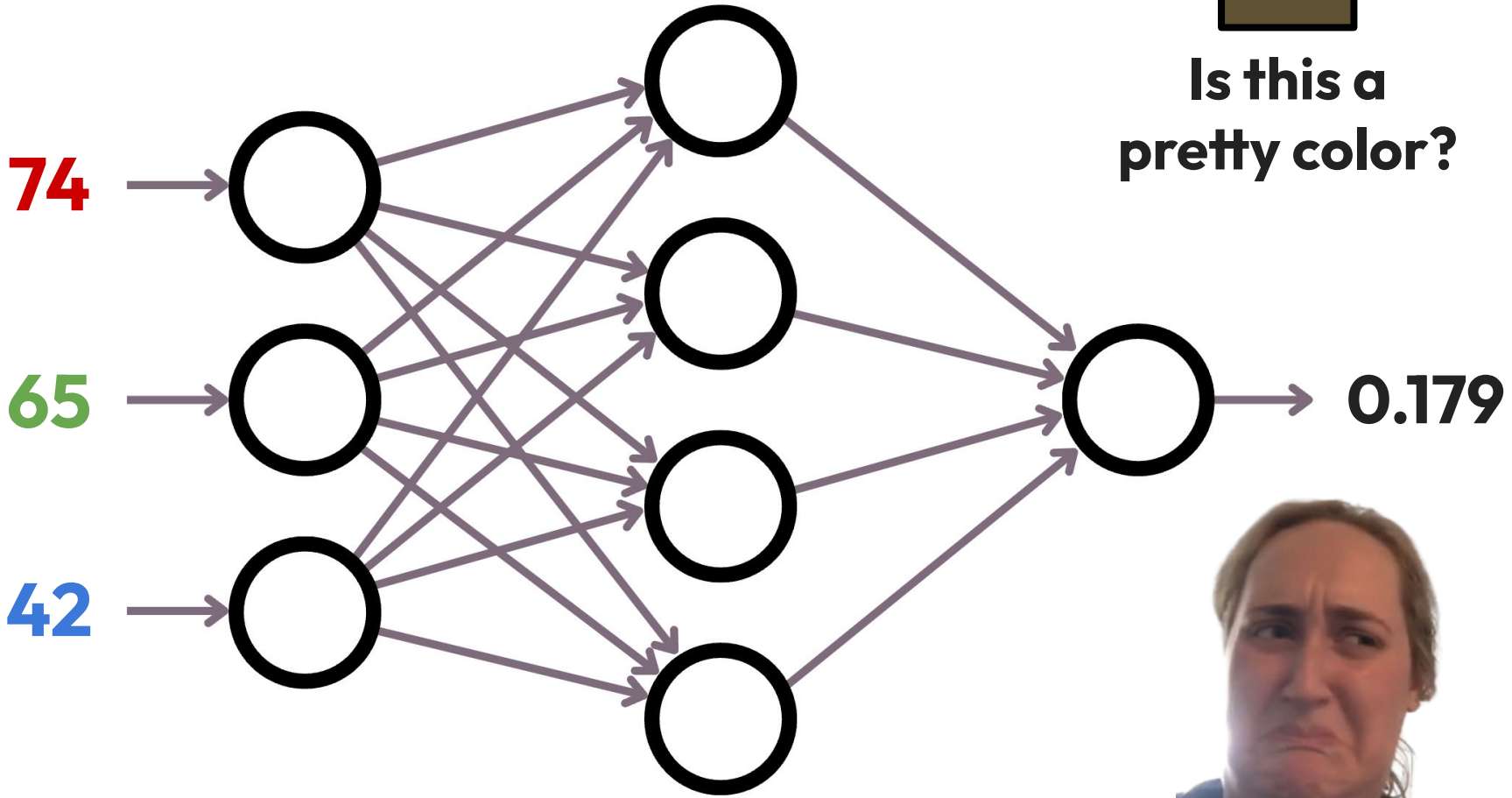


Is this a
pretty color?

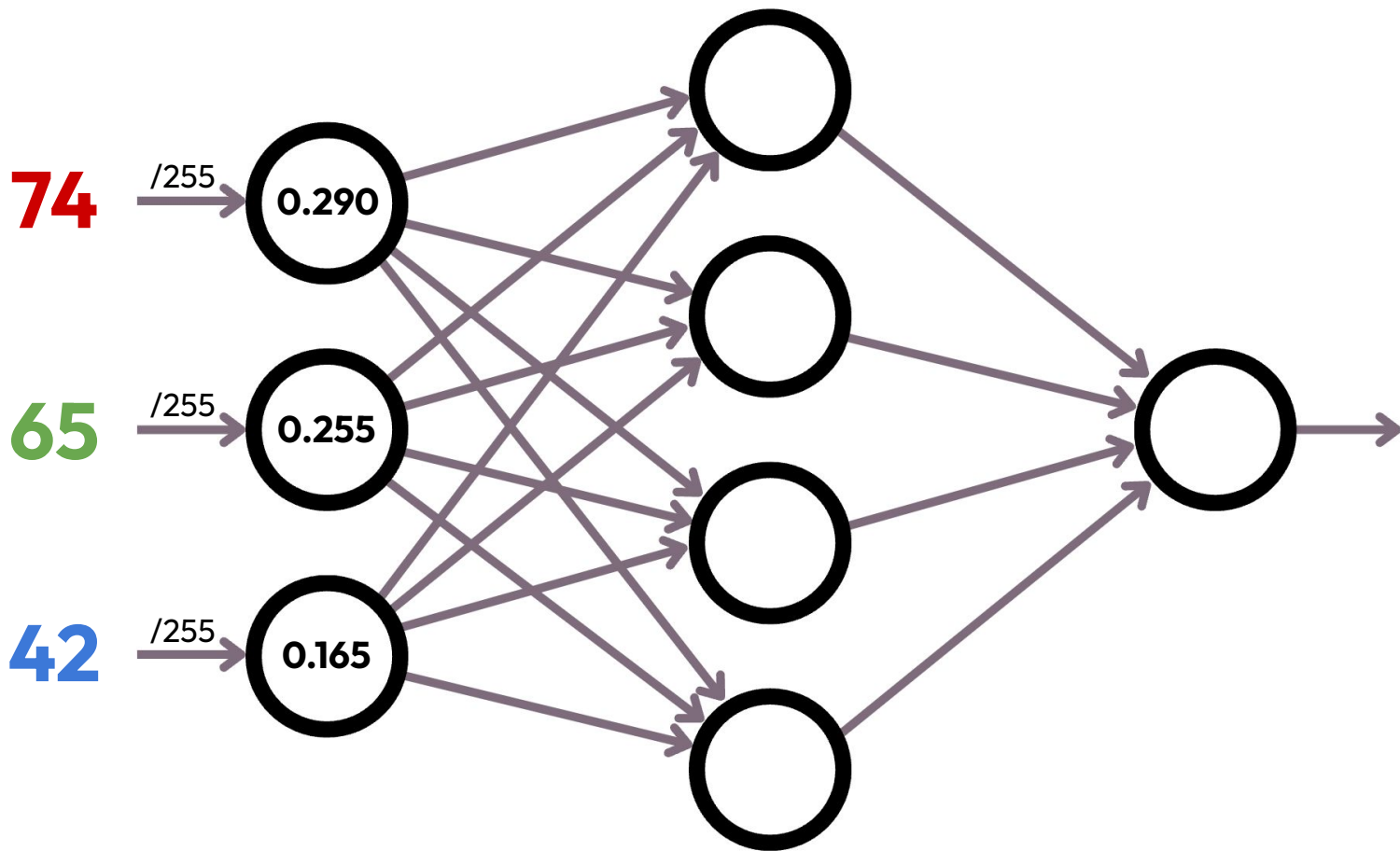
74

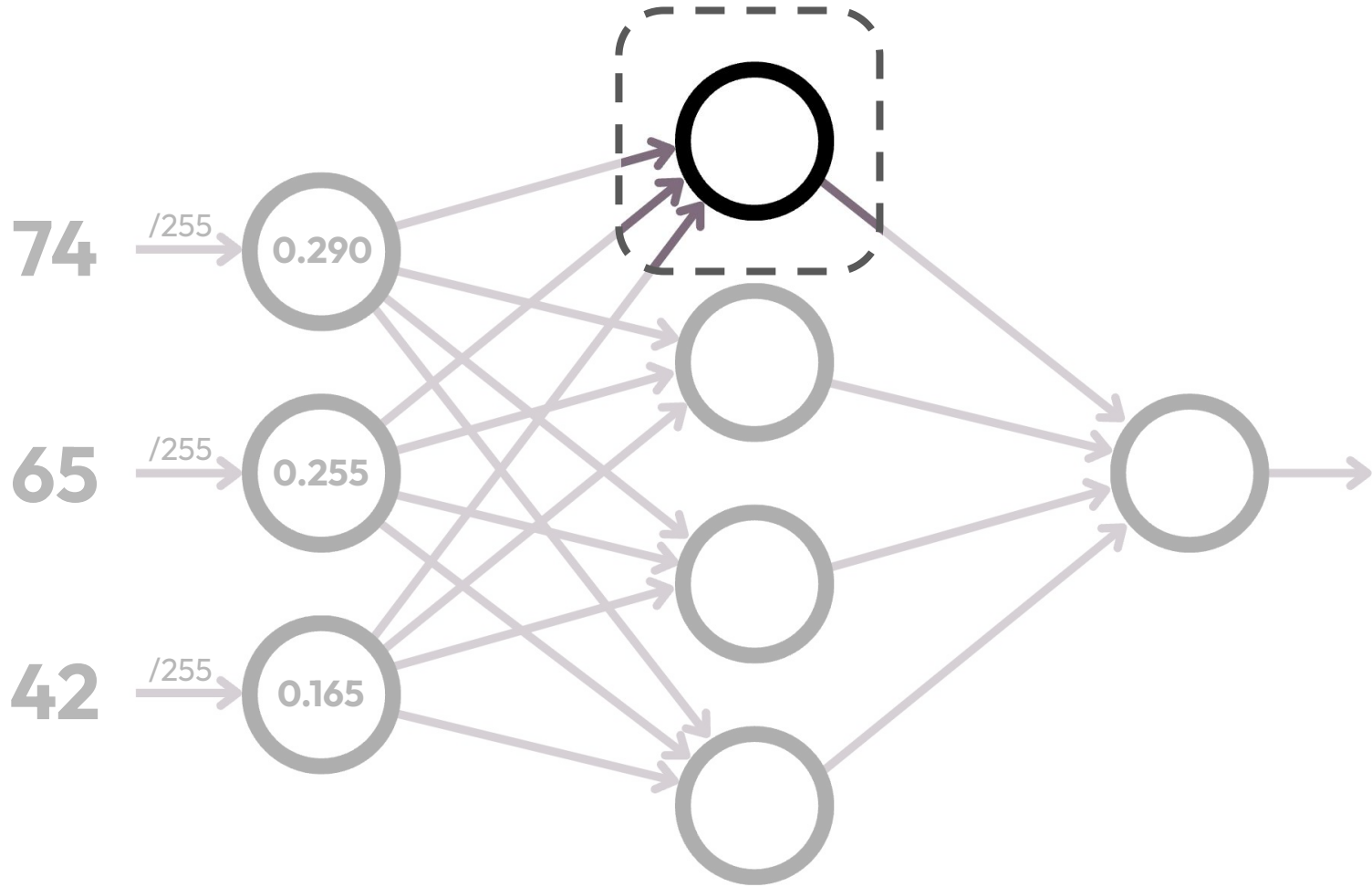
65

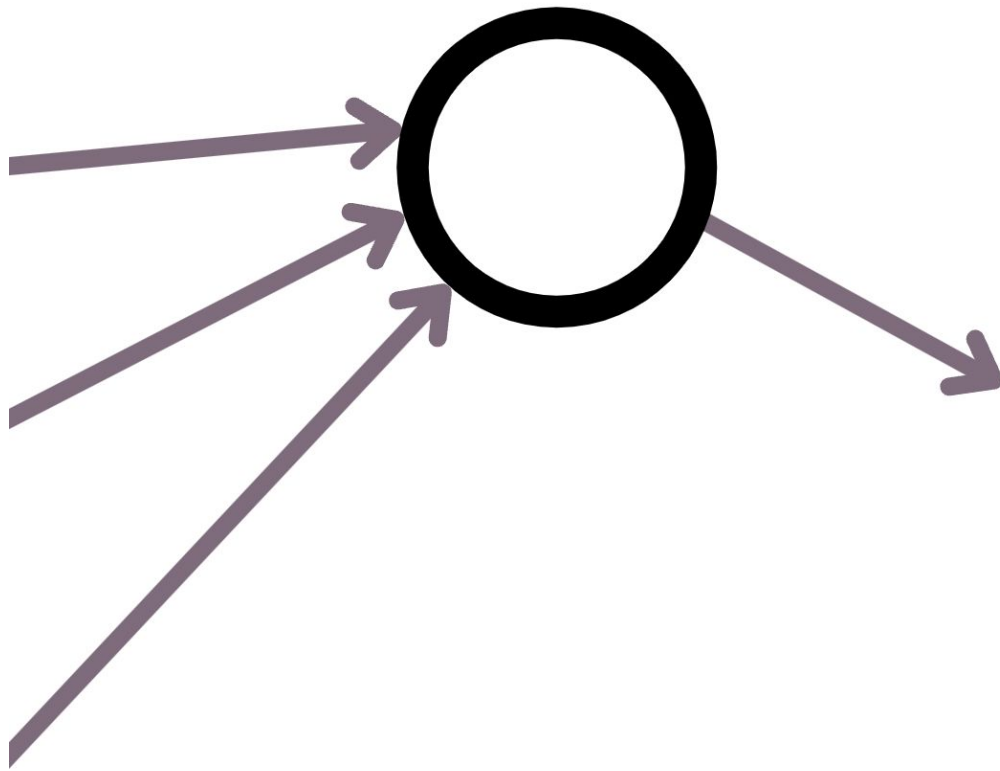
42



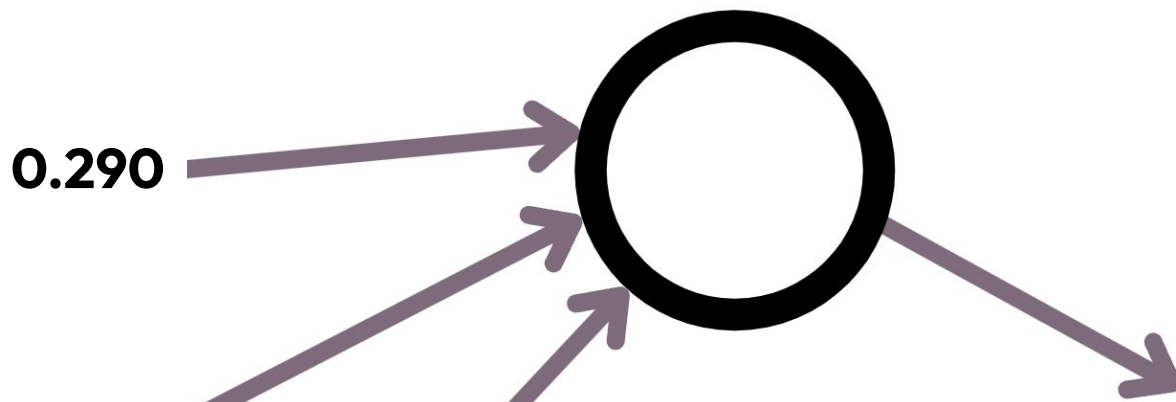
What really happens inside?



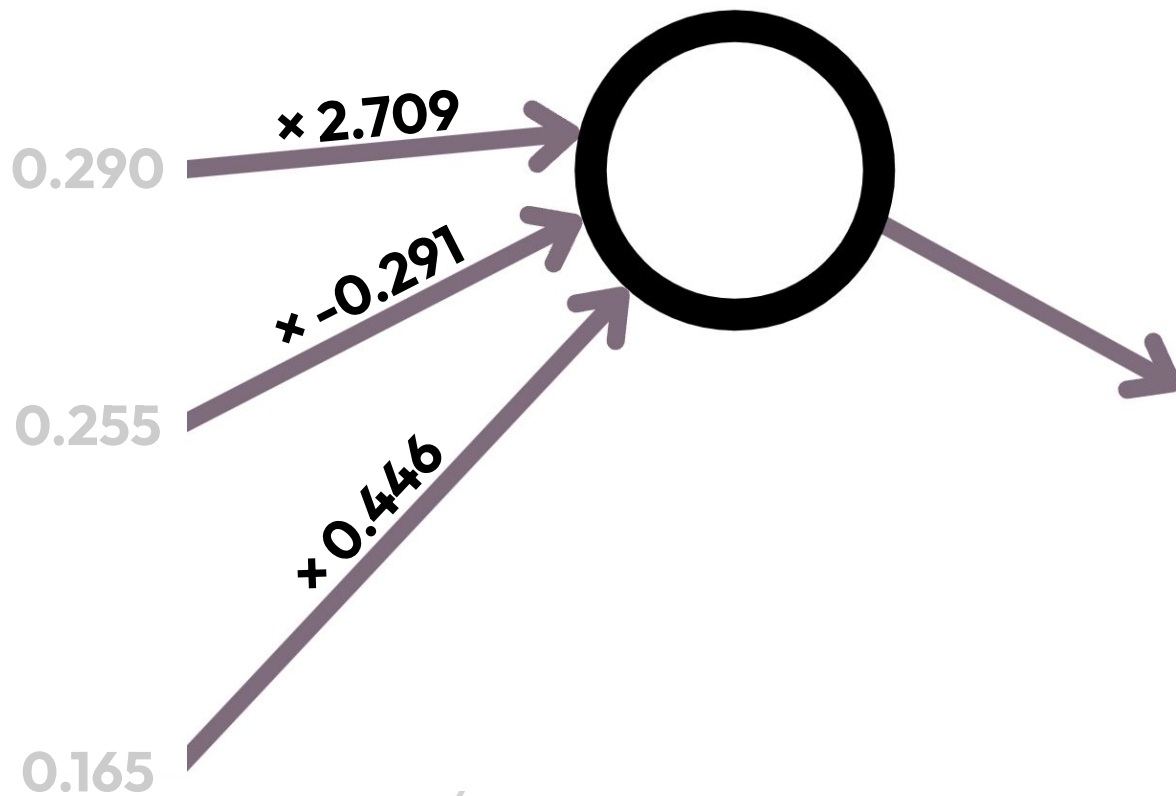




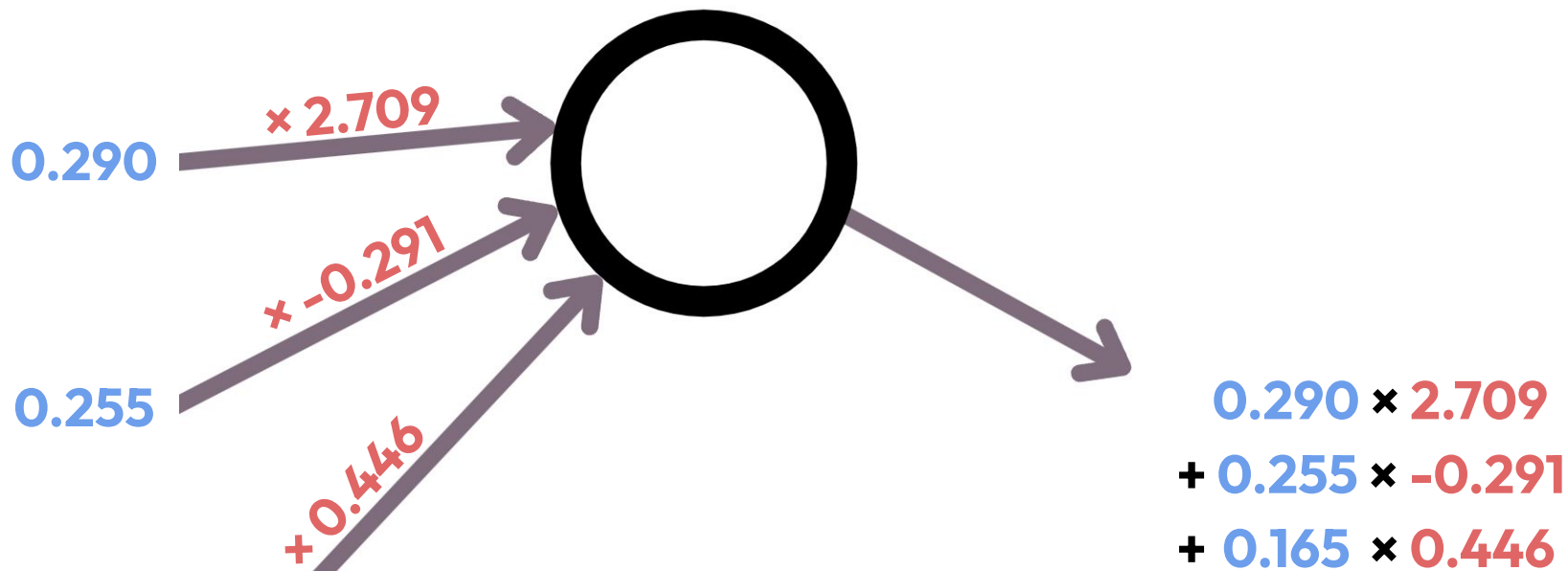
$$output = f \left(bias + \sum_i input_i \cdot weight_i \right)$$



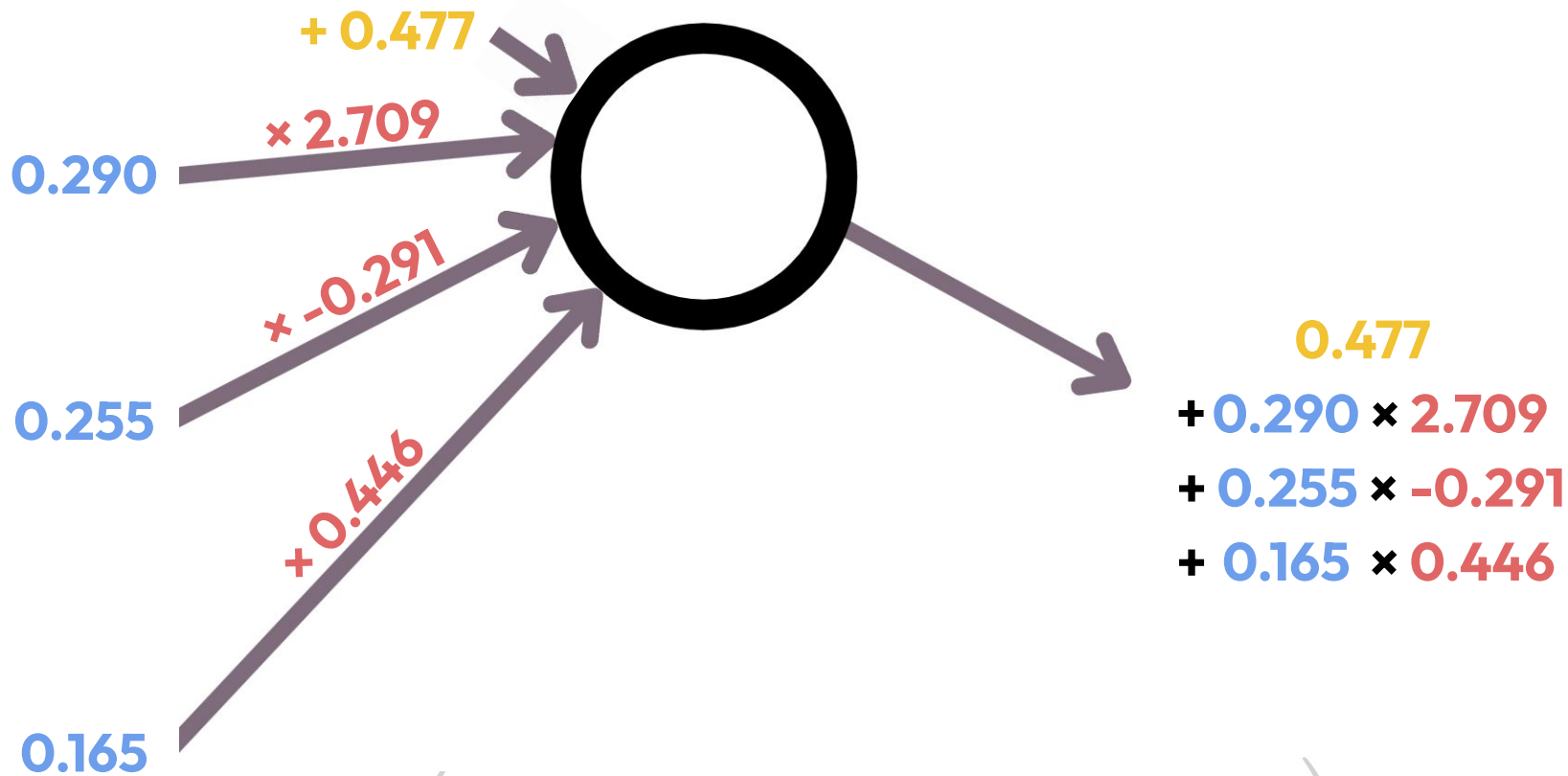
$$output = f \left(bias + \sum_i input_i \cdot weight_i \right)$$



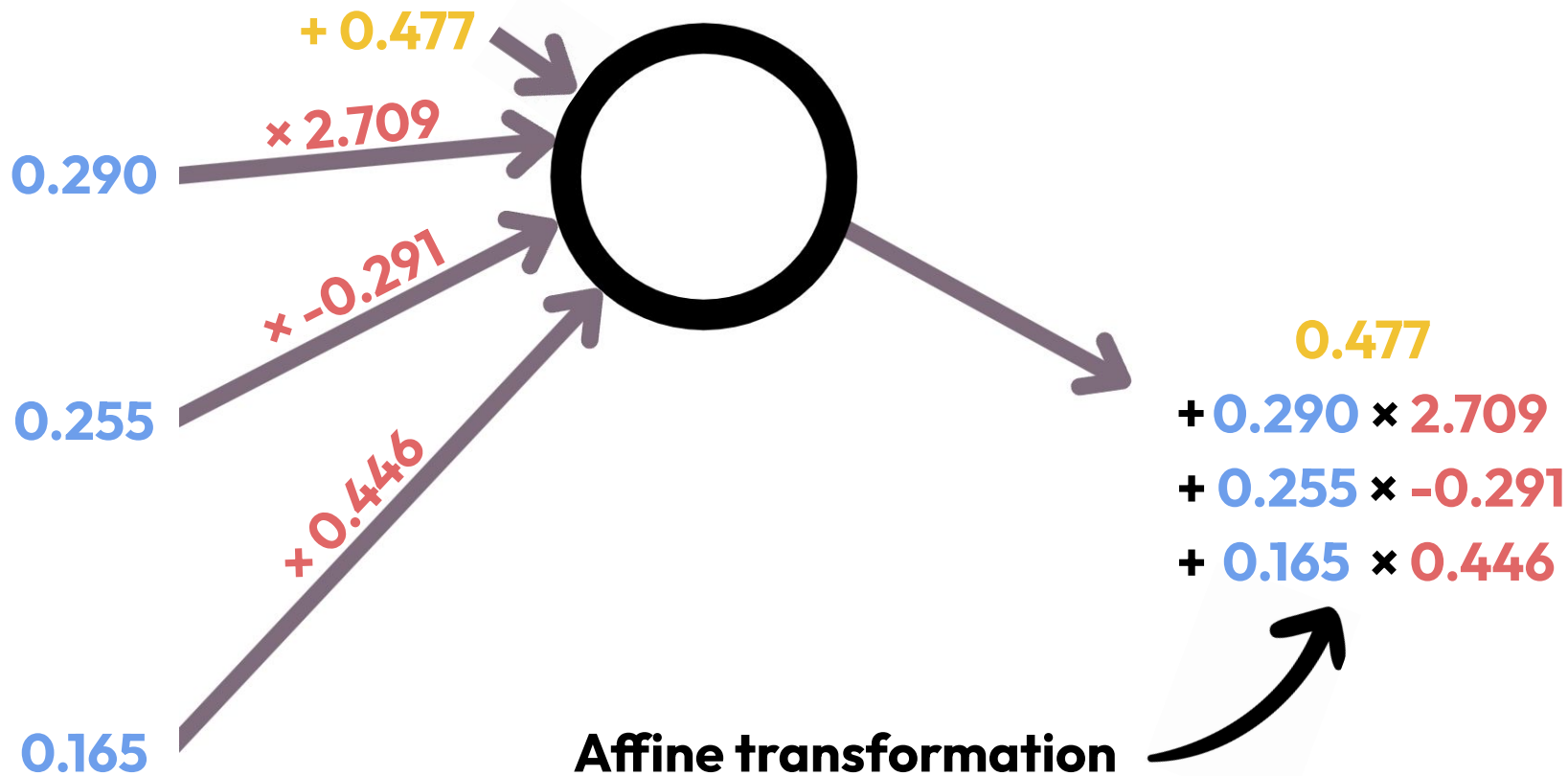
$$output = f \left(bias + \sum_i input_i \cdot weight_i \right)$$



$$output = f \left(bias + \sum_i \underline{input_i} \cdot \underline{weight_i} \right)$$

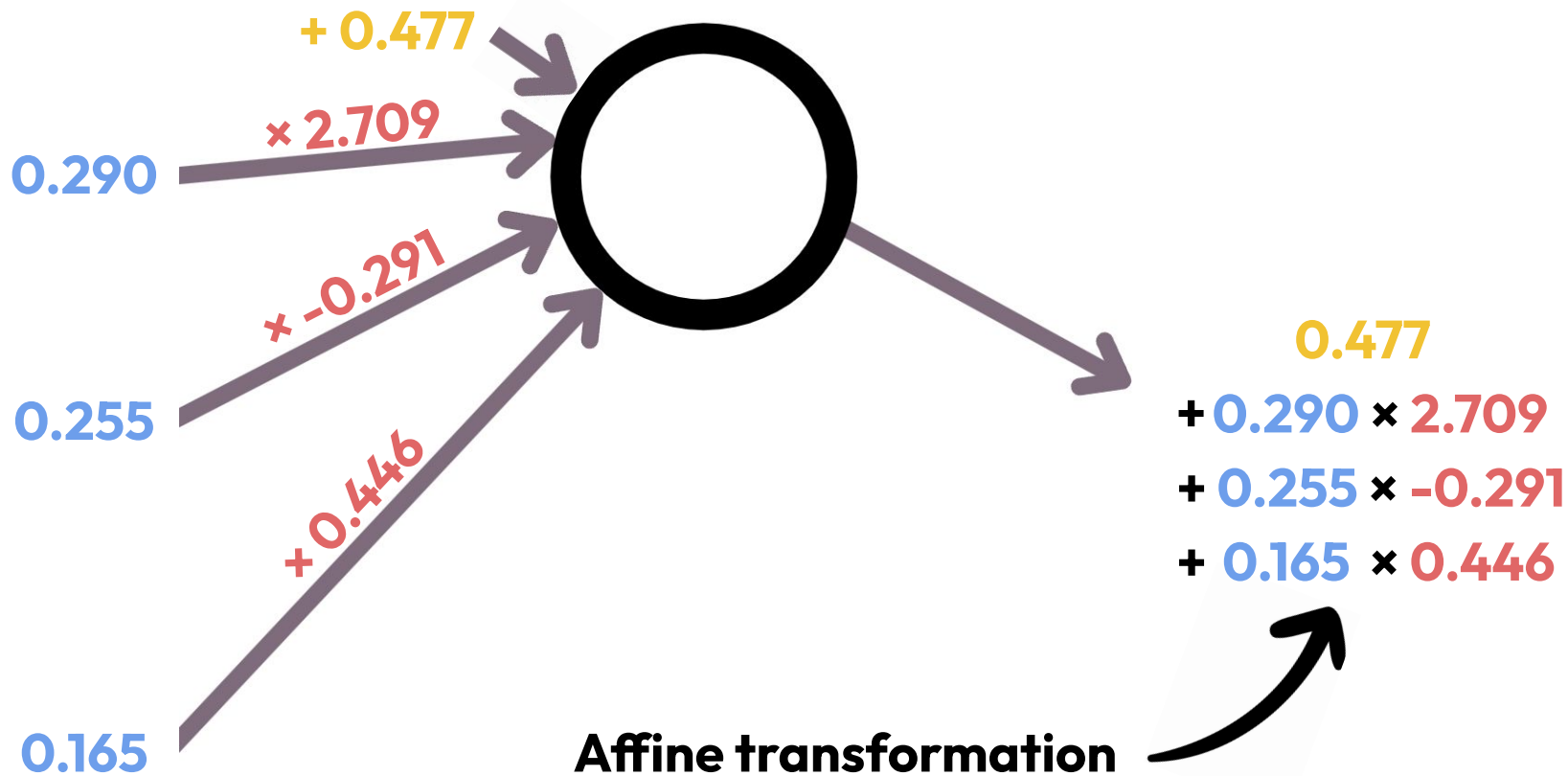


$$output = f \left(\underbrace{bias}_{0.477} + \sum_i \underbrace{input_i}_{0.290} \cdot \underbrace{weight_i}_{2.709} \right)$$

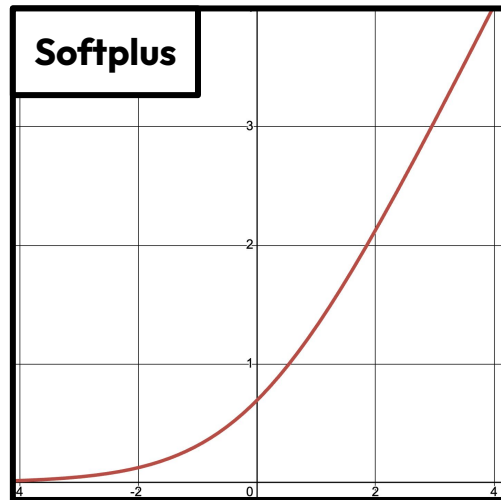
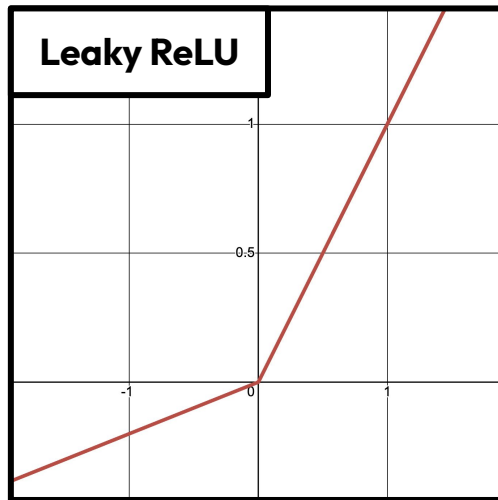
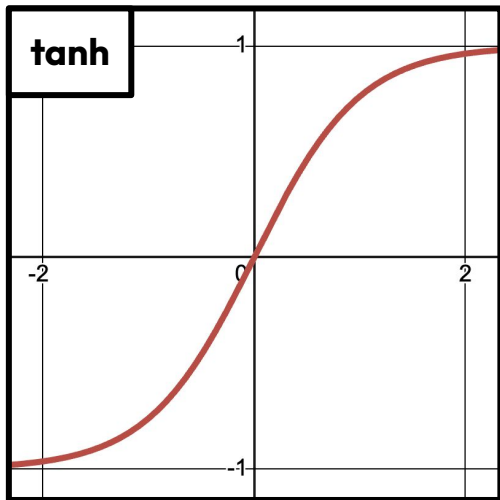
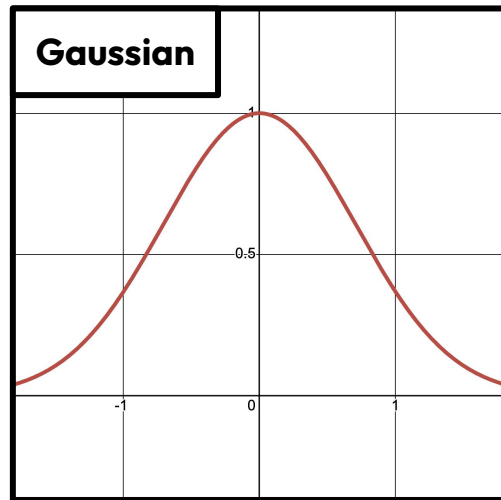
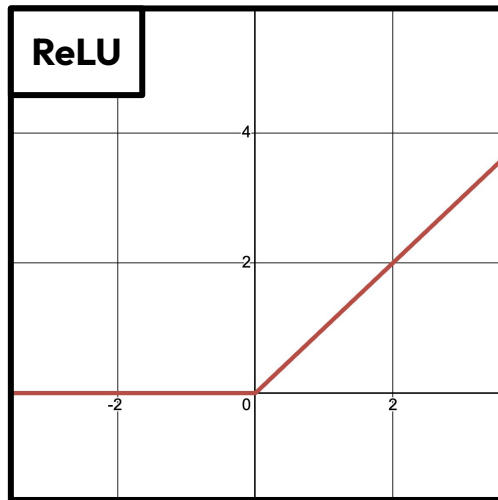
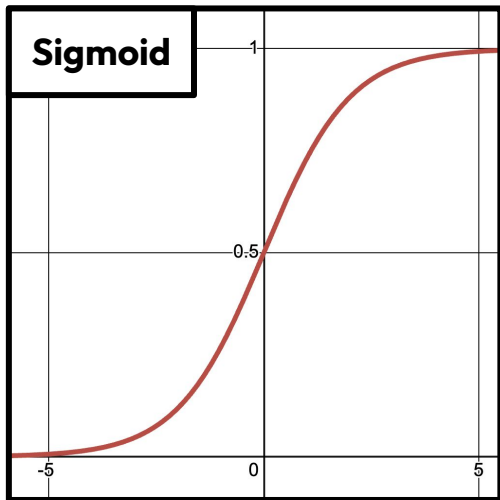


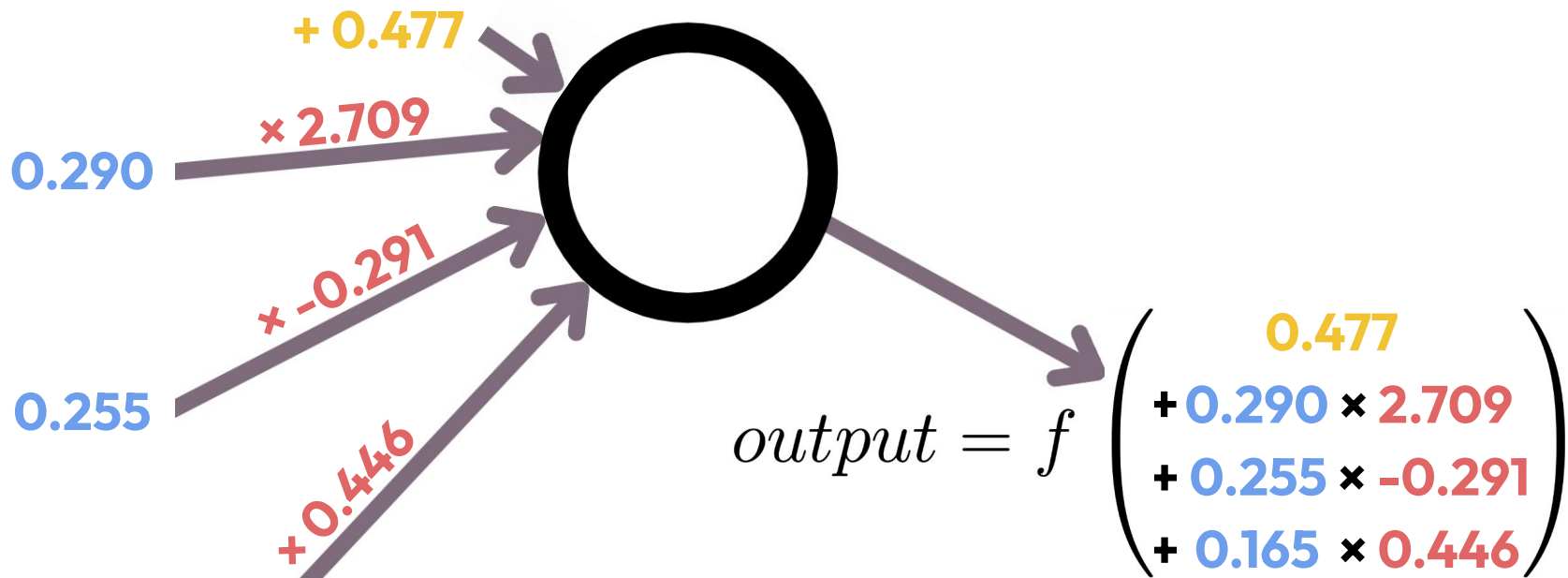
Affine transformation
is not a universal approximator

$$\begin{aligned} & 0.477 \\ & + 0.290 \times 2.709 \\ & + 0.255 \times -0.291 \\ & + 0.165 \times 0.446 \end{aligned}$$

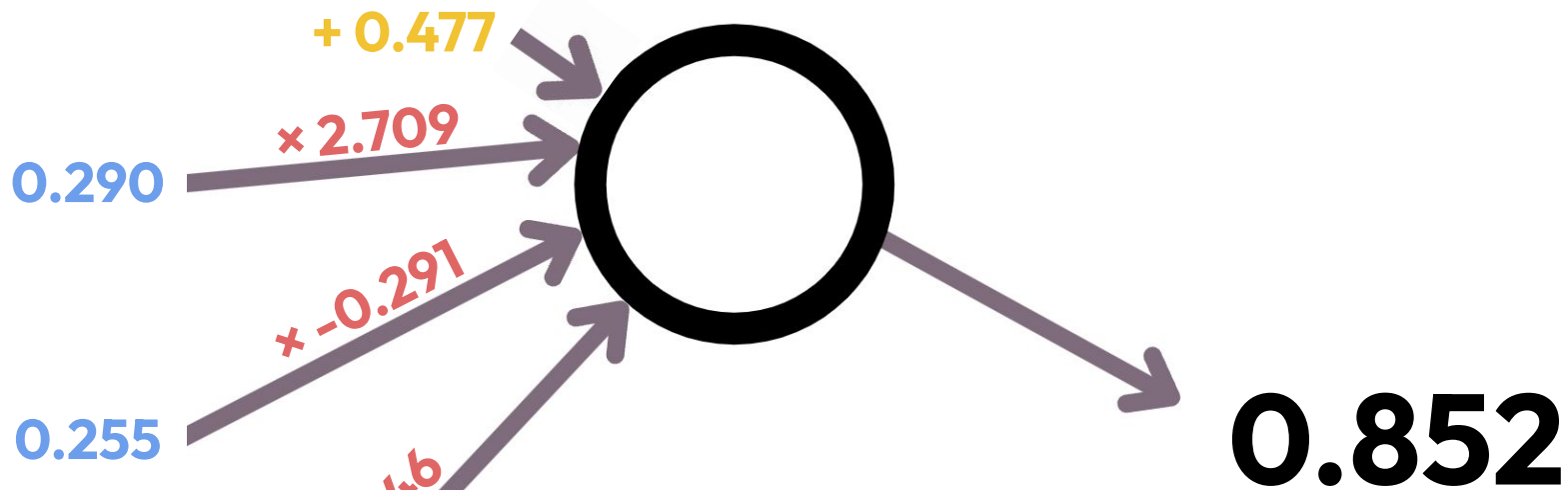


Affine transformation
is not a universal approximator
Let's add a non-linearity!

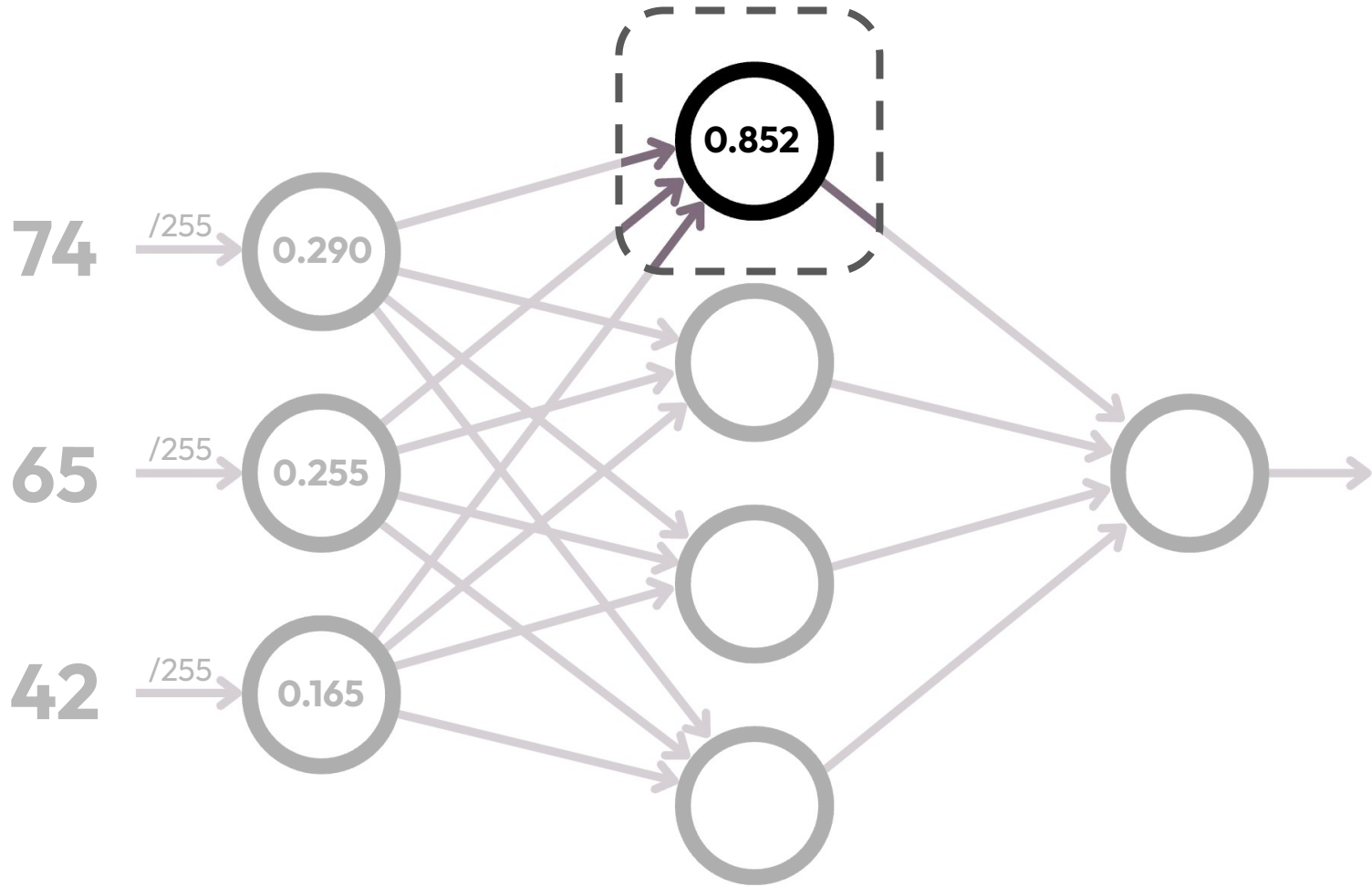


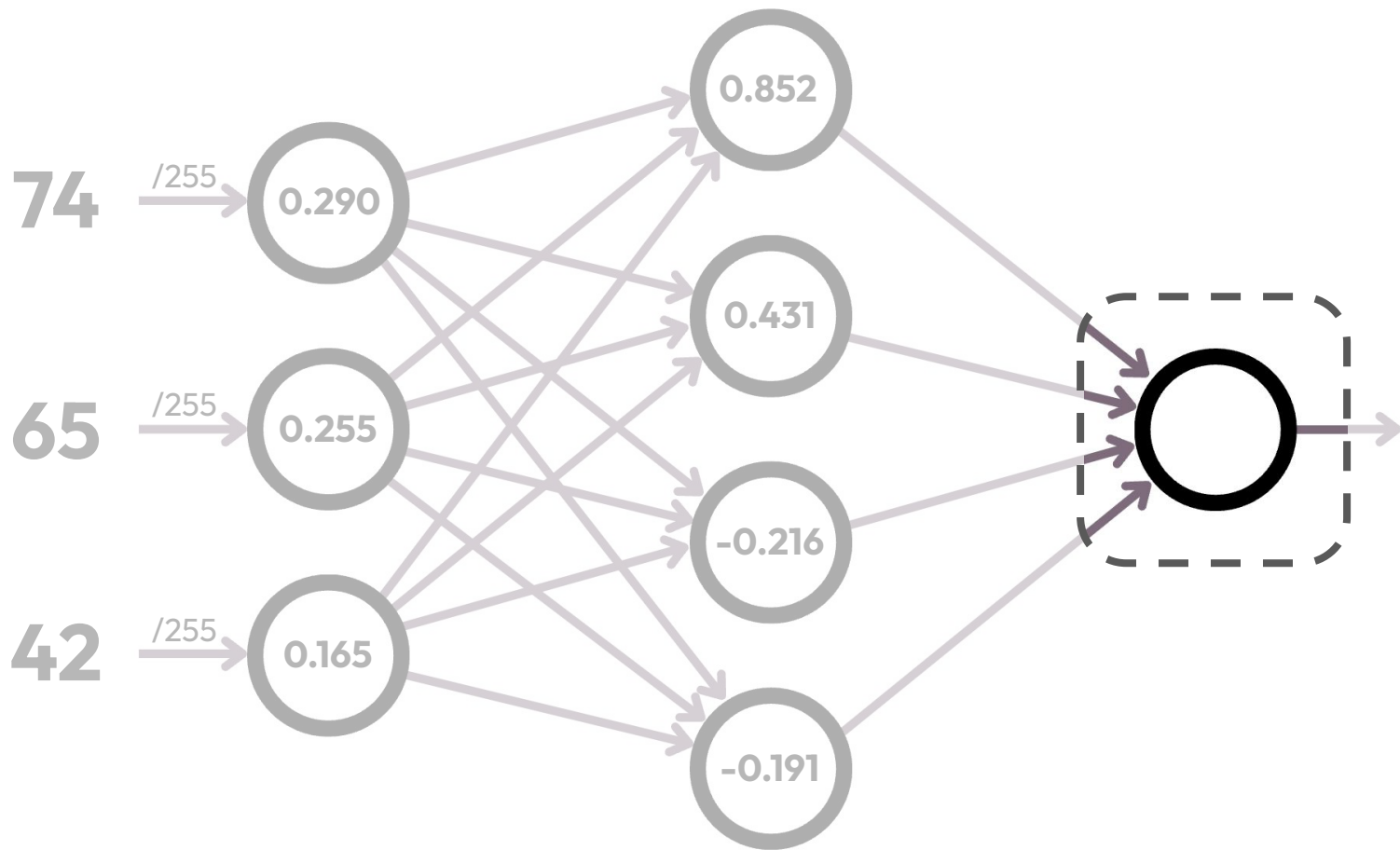


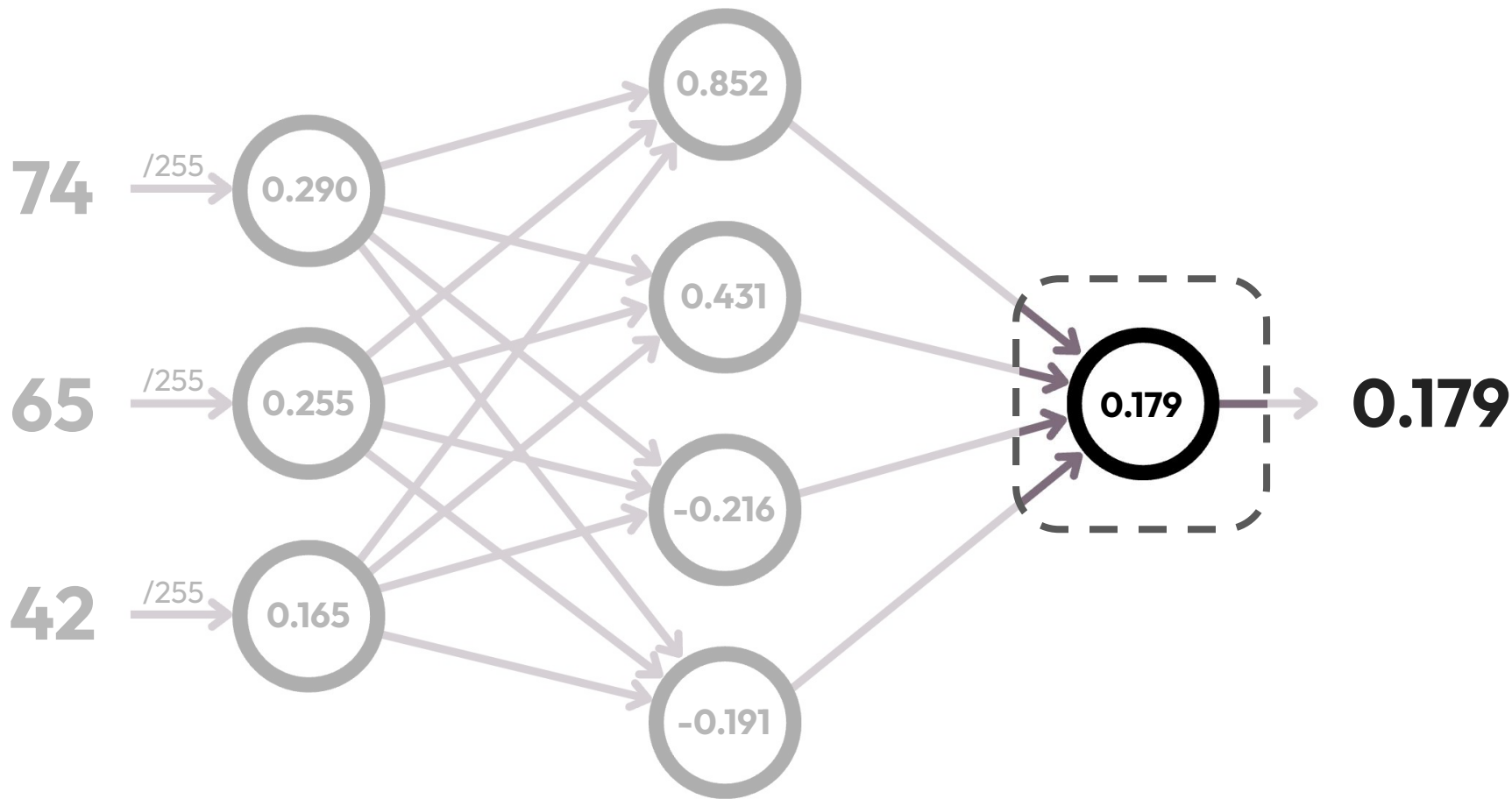
$$output = f \left(bias + \sum_i input_i \cdot weight_i \right)$$



$$output = f \left(bias + \sum_i input_i \cdot weight_i \right)$$



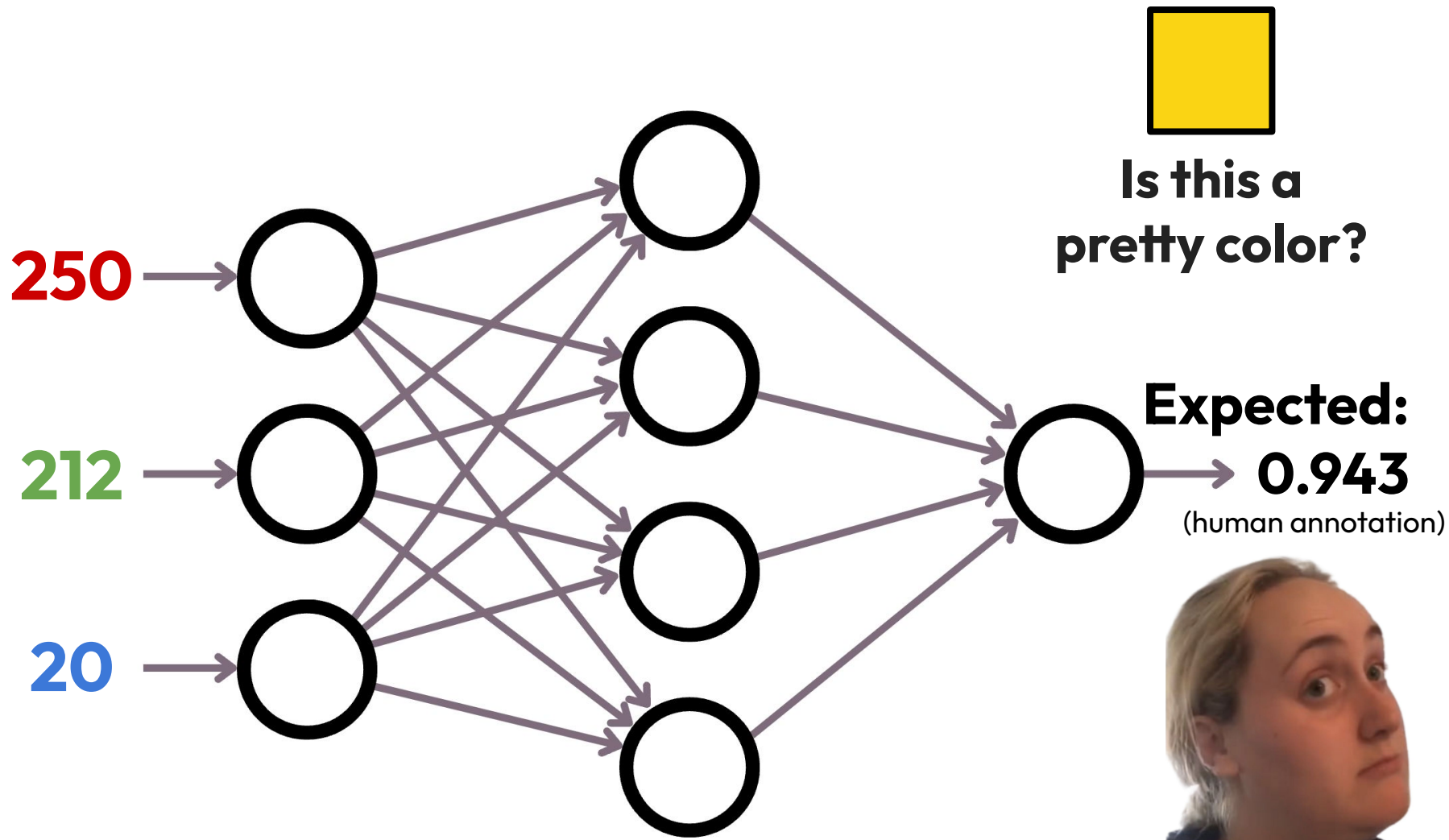


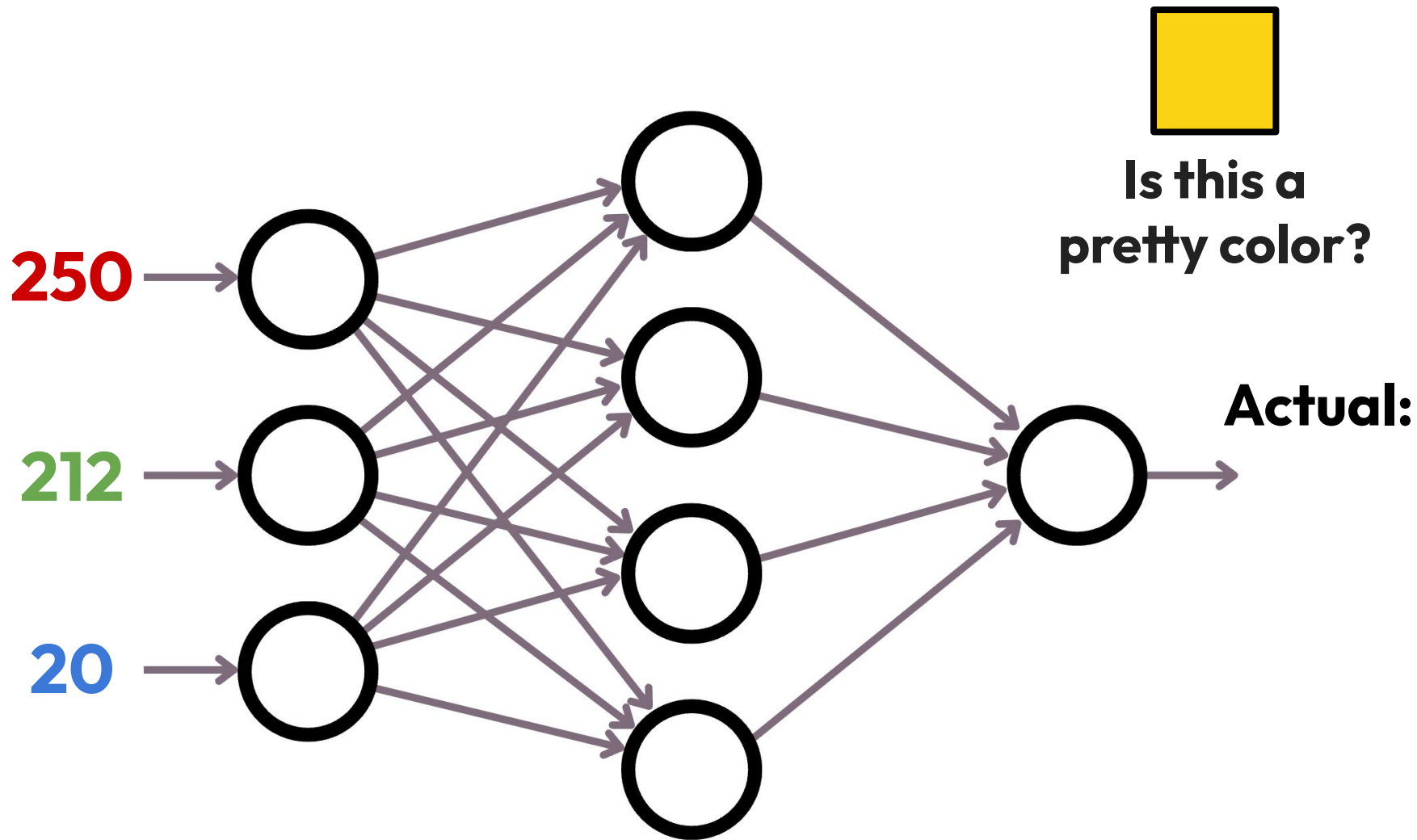


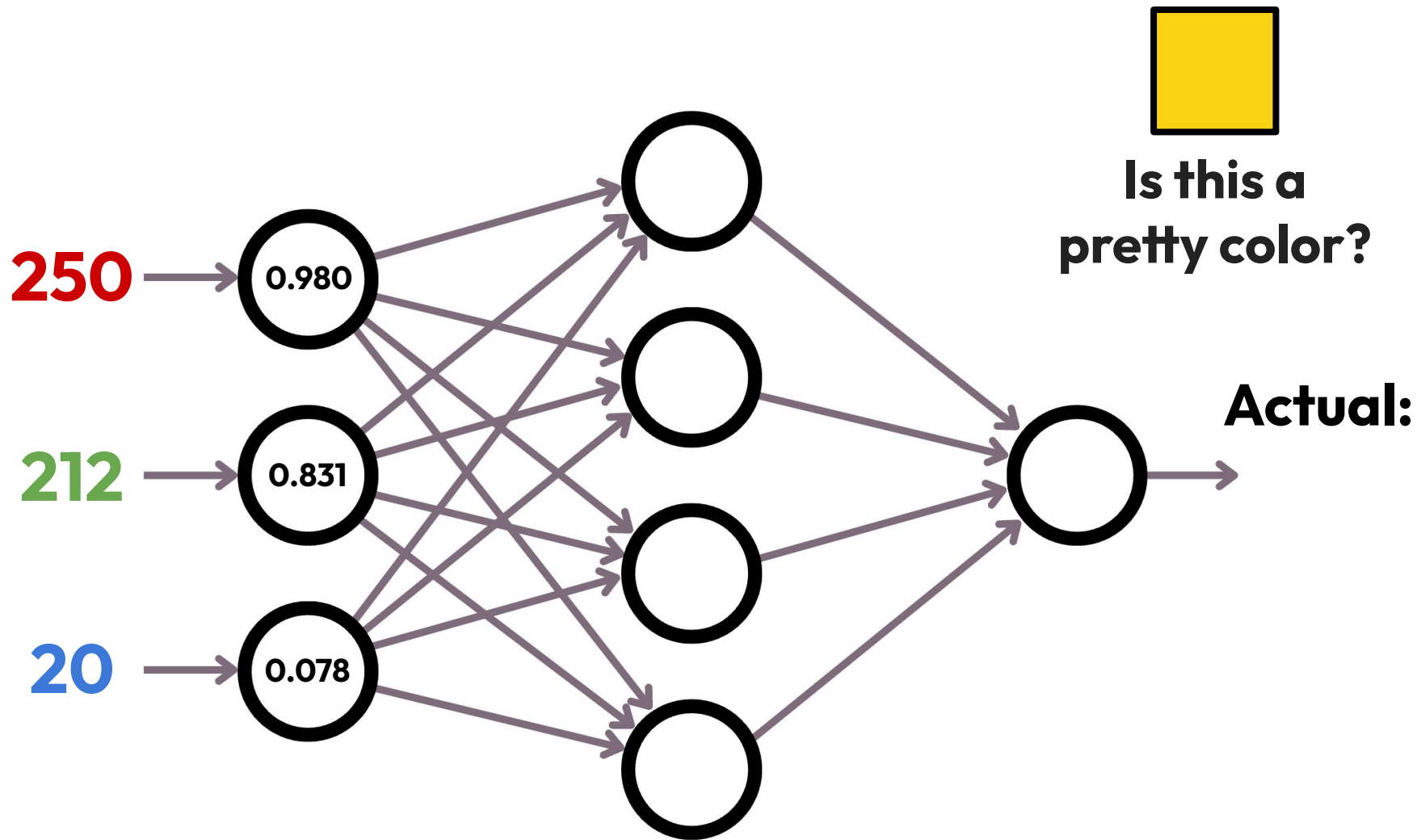
What is training?

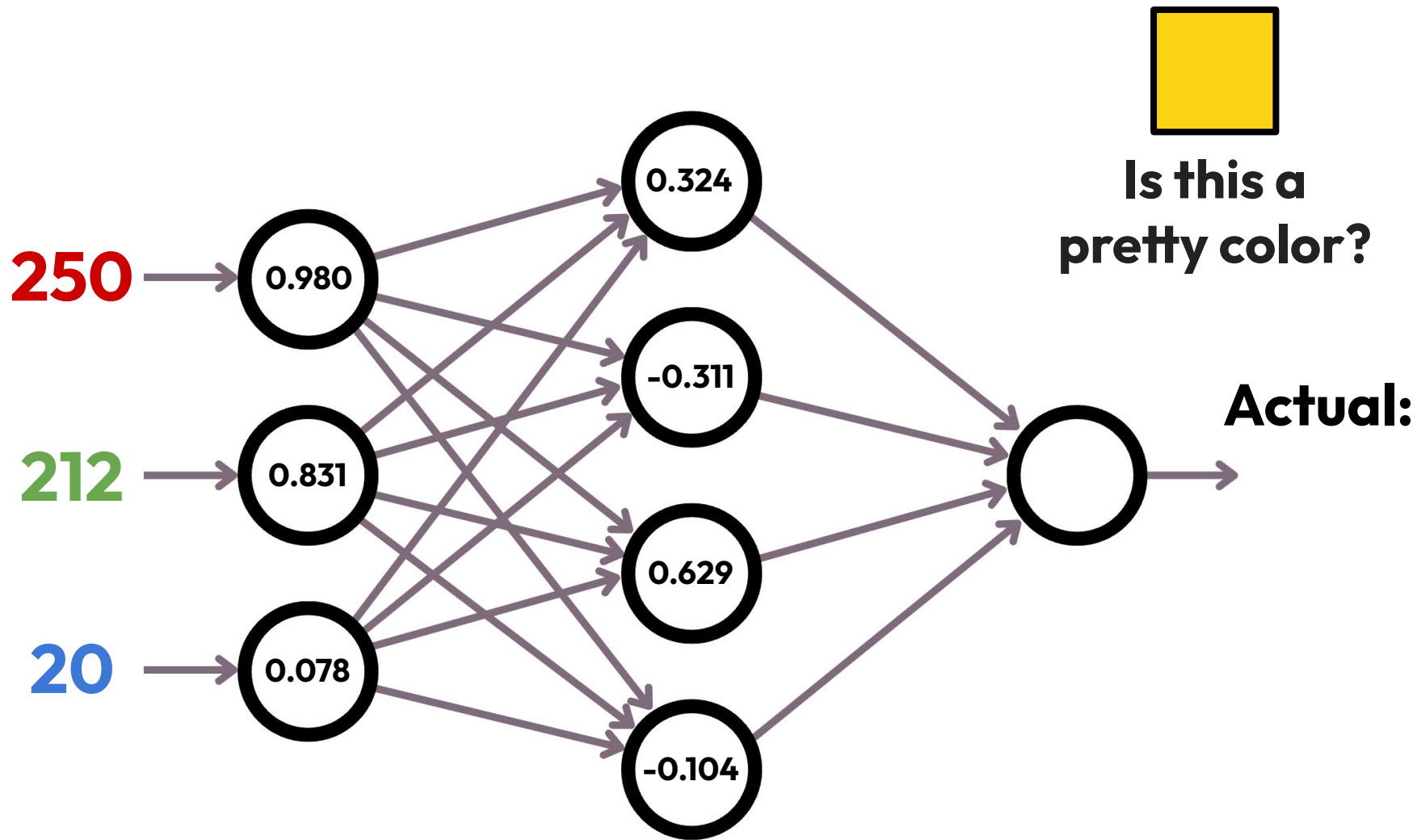
Initially, the neural net does not know what to do

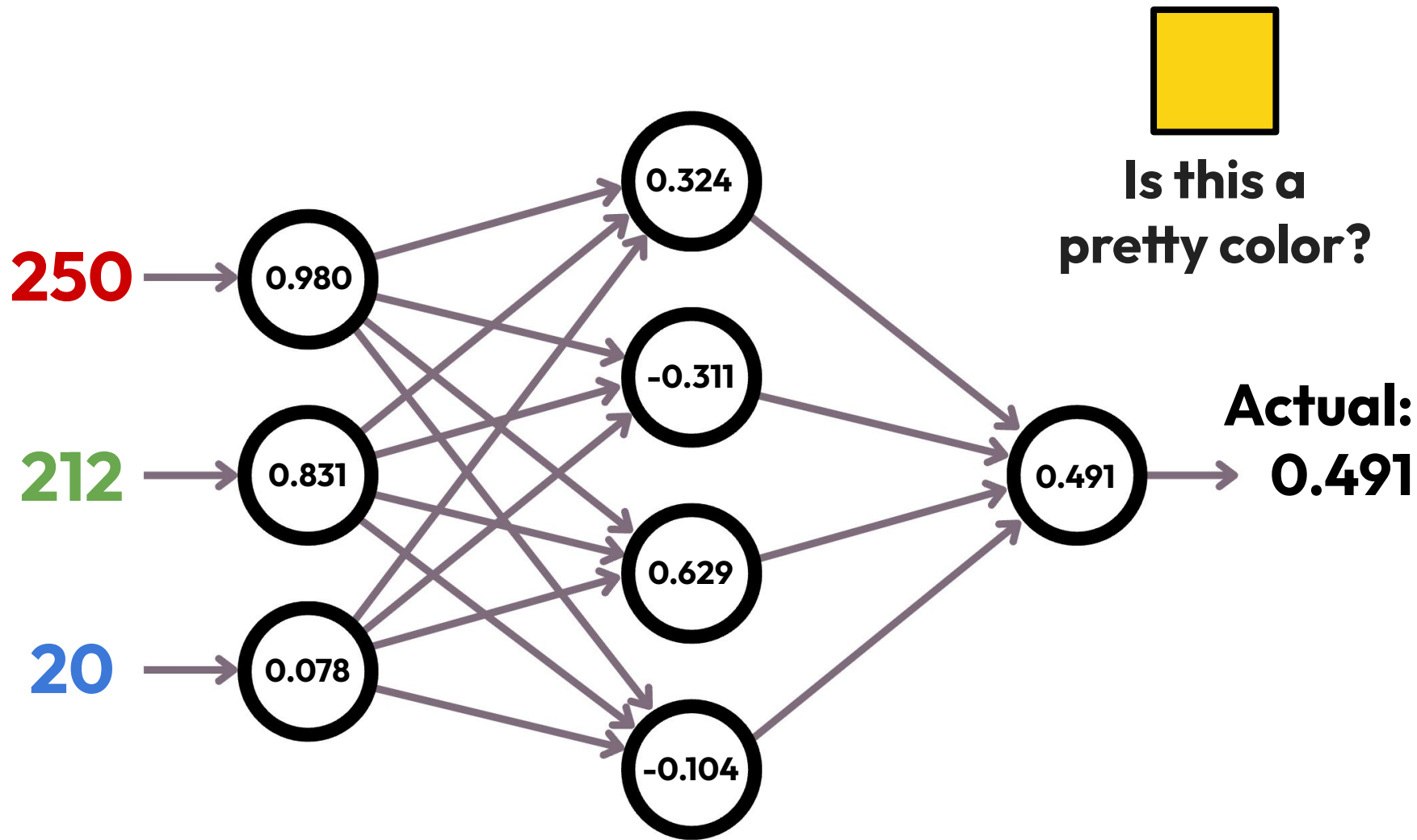
- It can approximate any function
- Adjust parameters (**weights** and **biases**)
- Learn their correct values through training

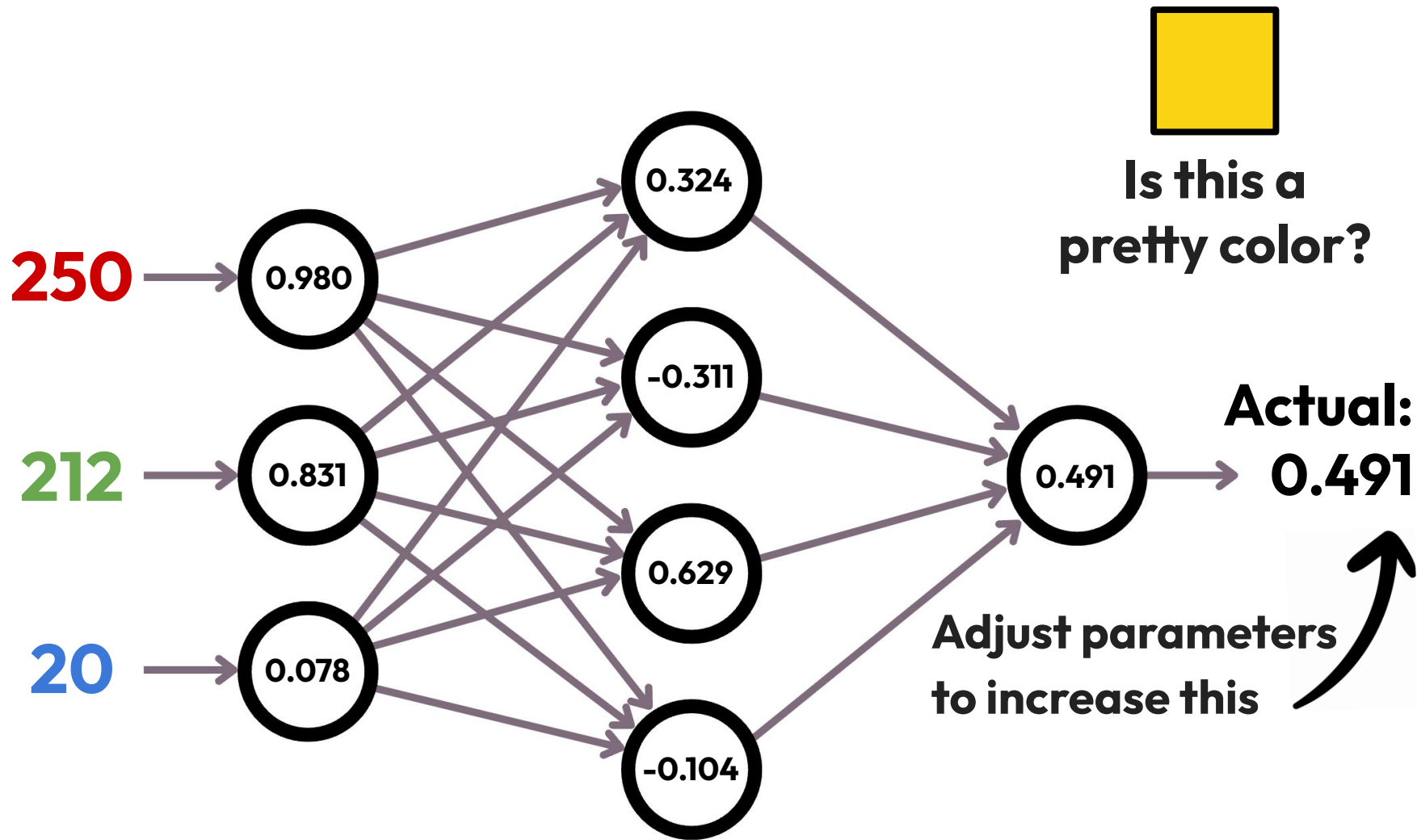


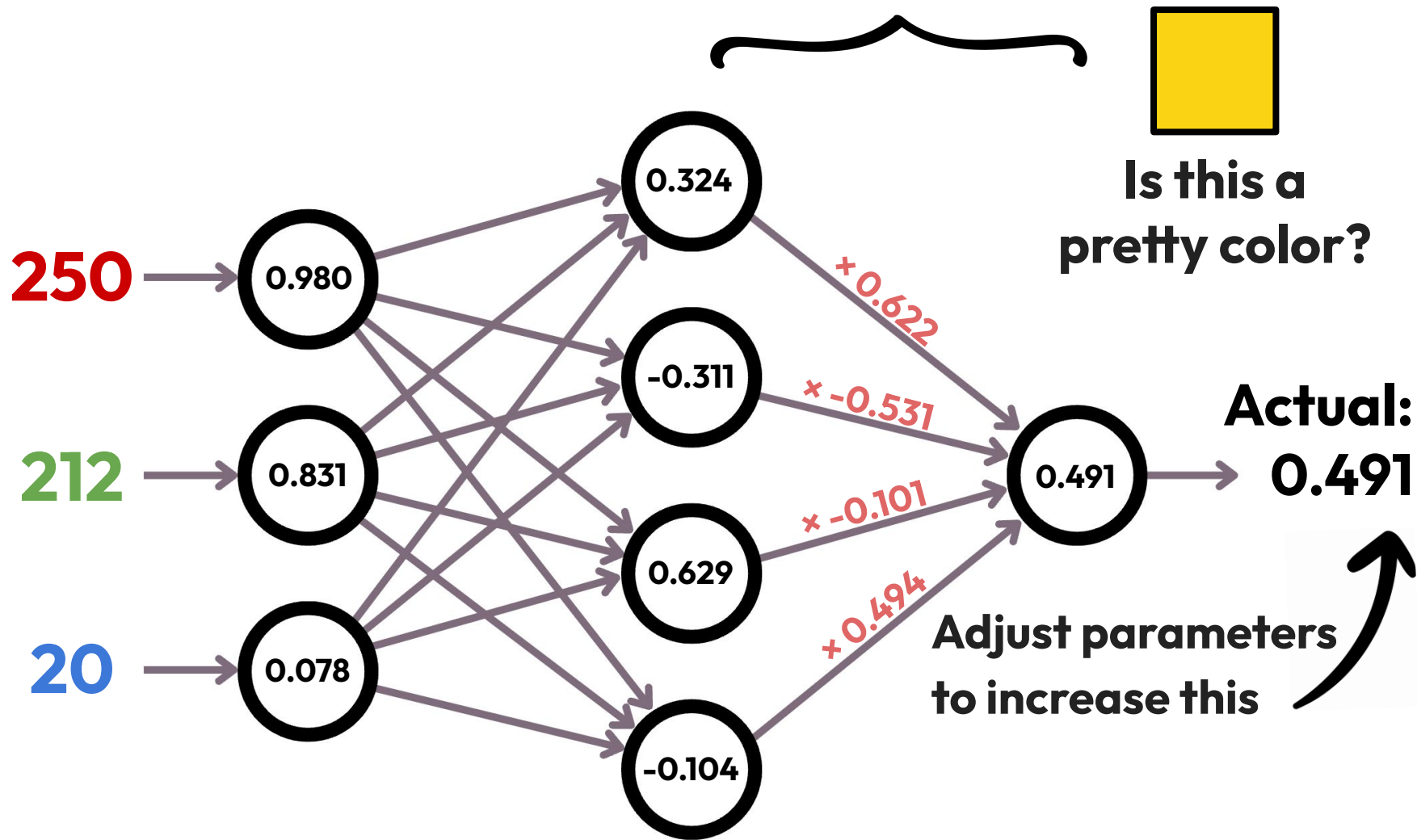


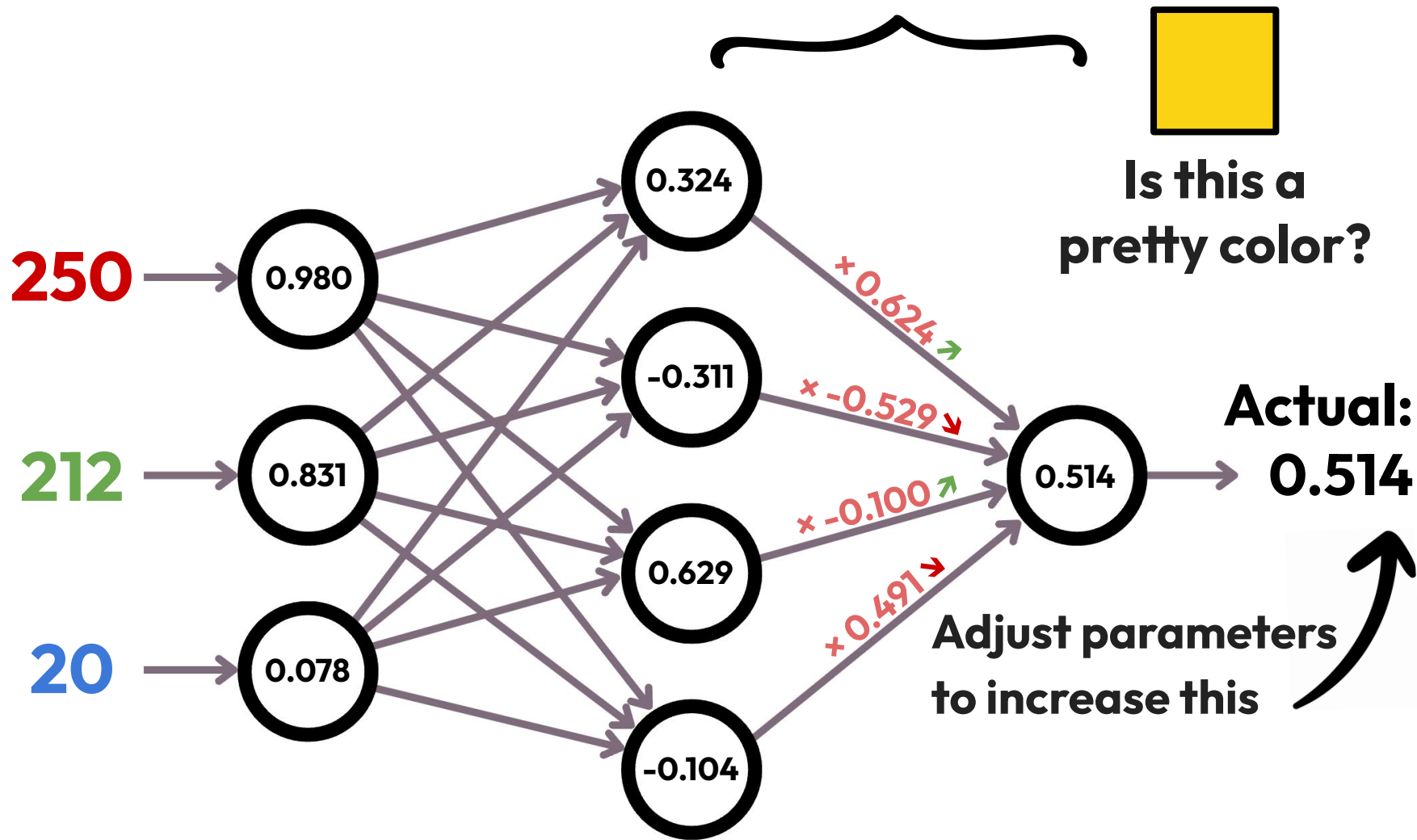


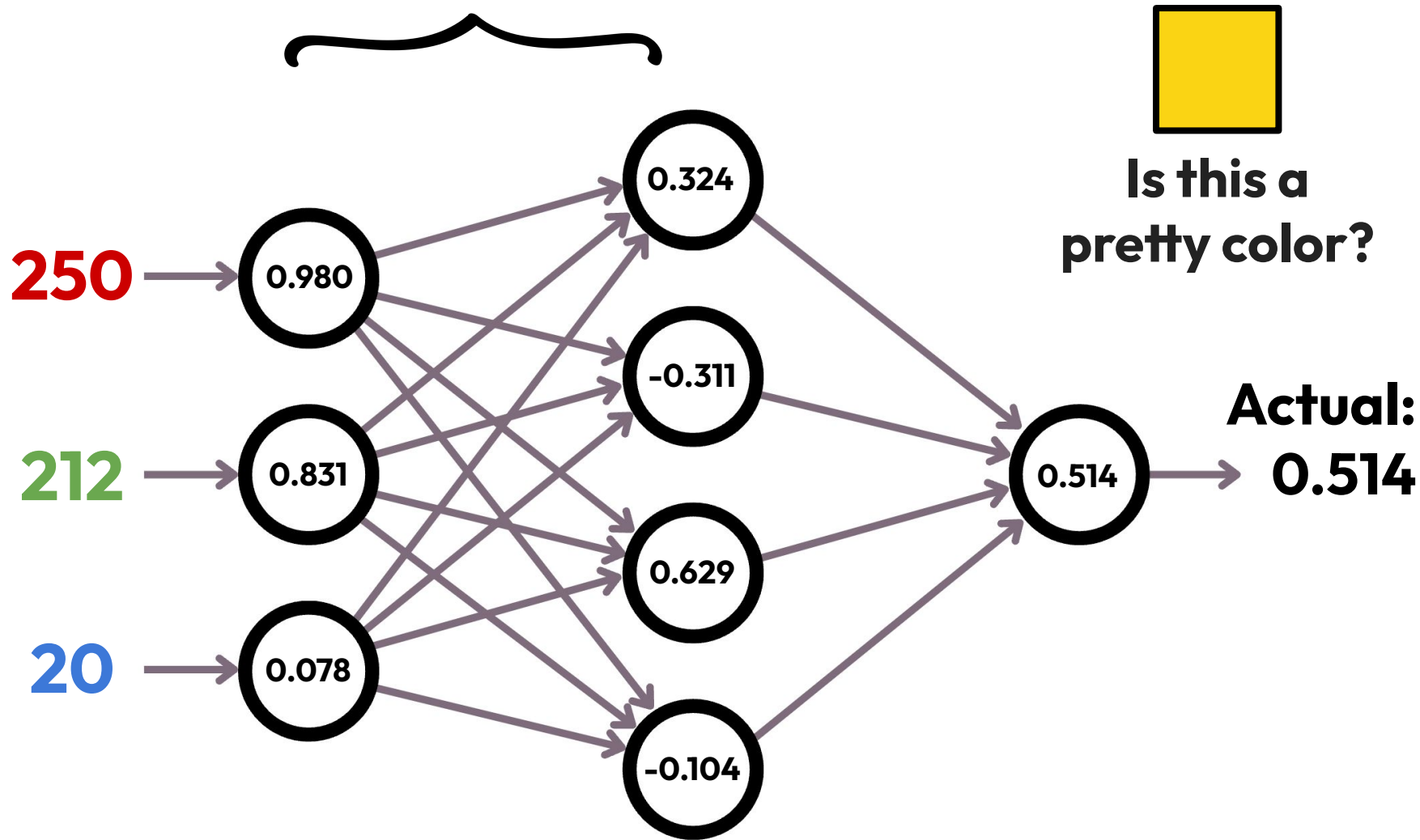


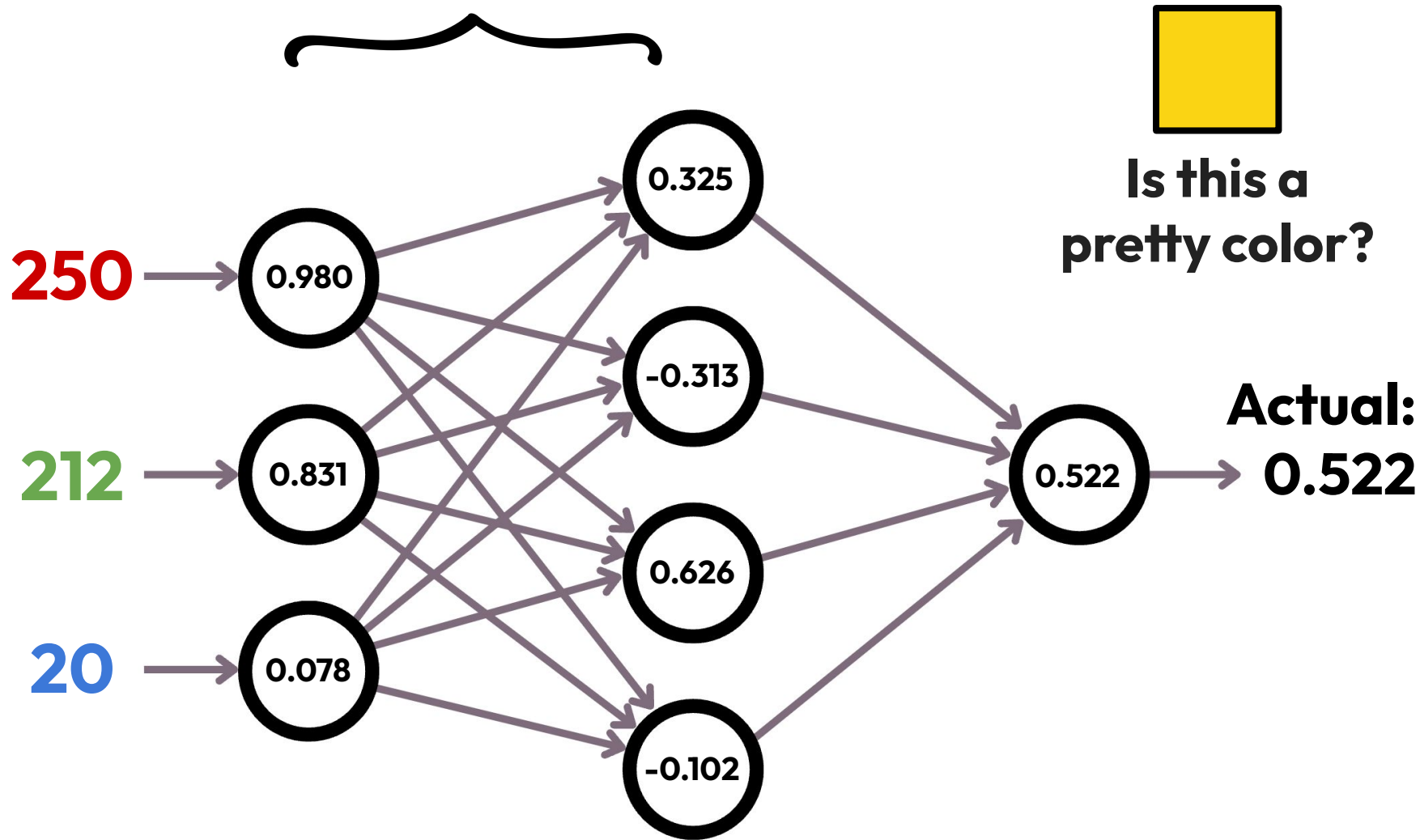


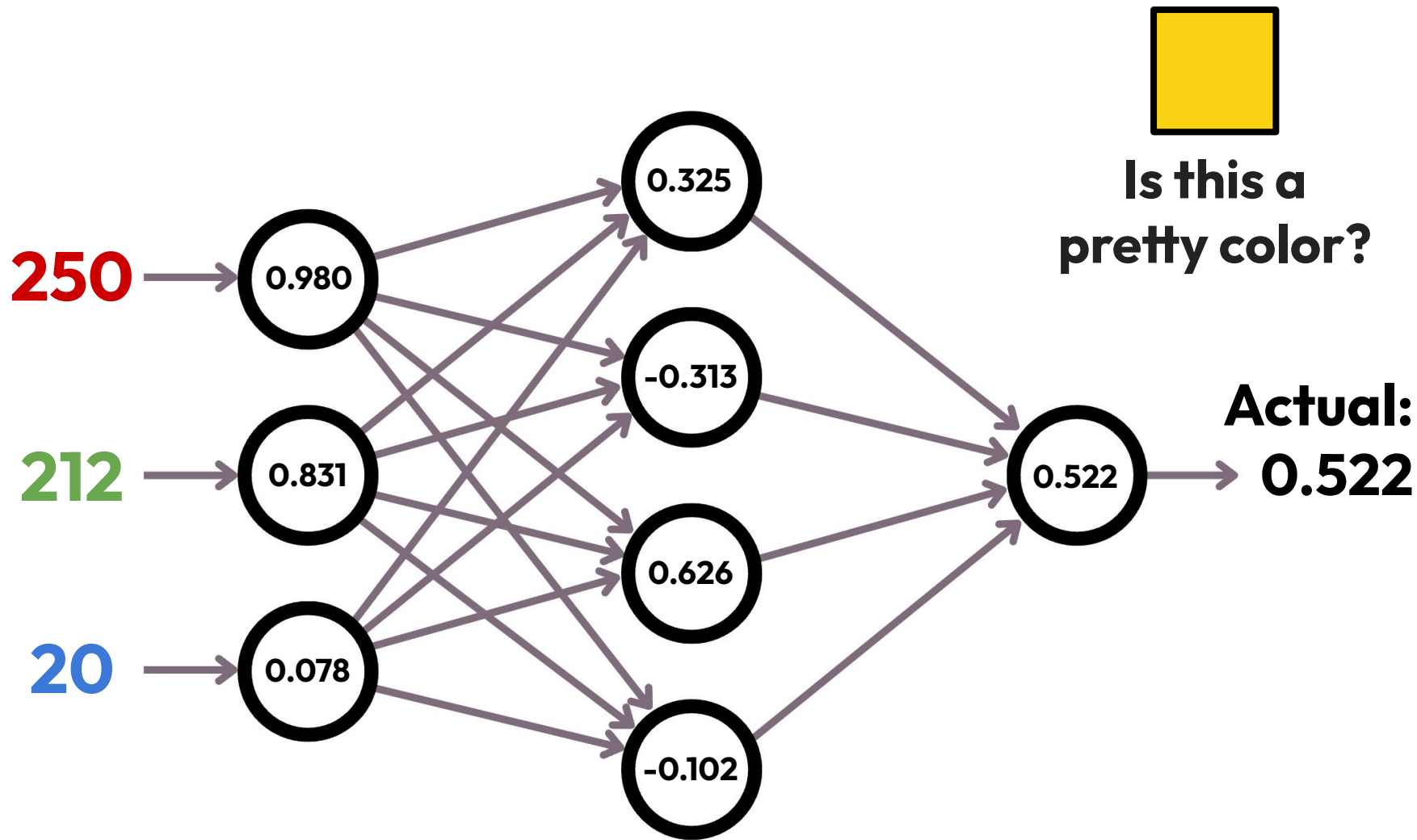












Repeat many times, over many examples

After enough training steps, the neural network
will approximate human evaluation

Parameter adjustments use the gradient

Neural nets must be **continuous** and **differentiable**

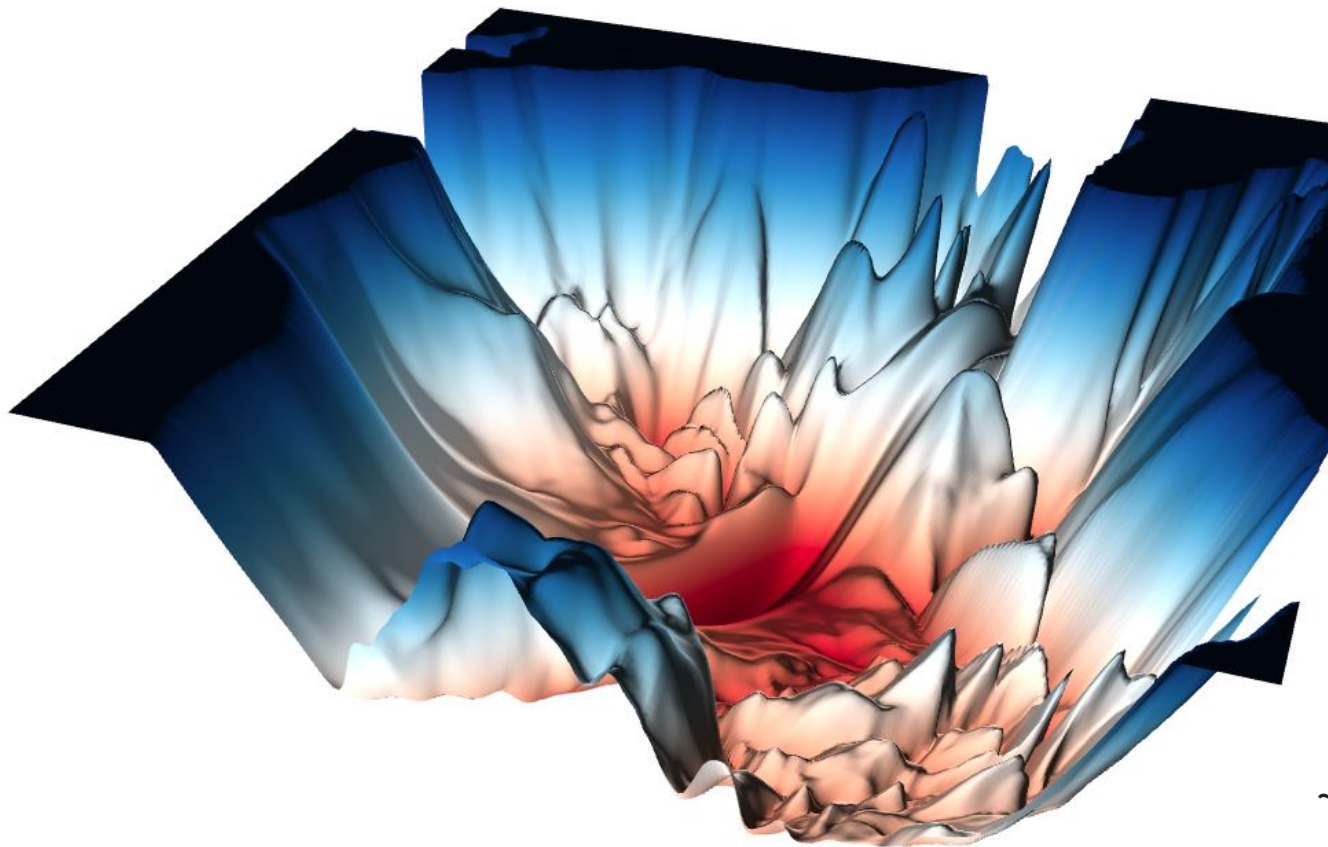
Parameter adjustments use the gradient

Neural nets must be **continuous** and **differentiable**

We can exploit these properties!

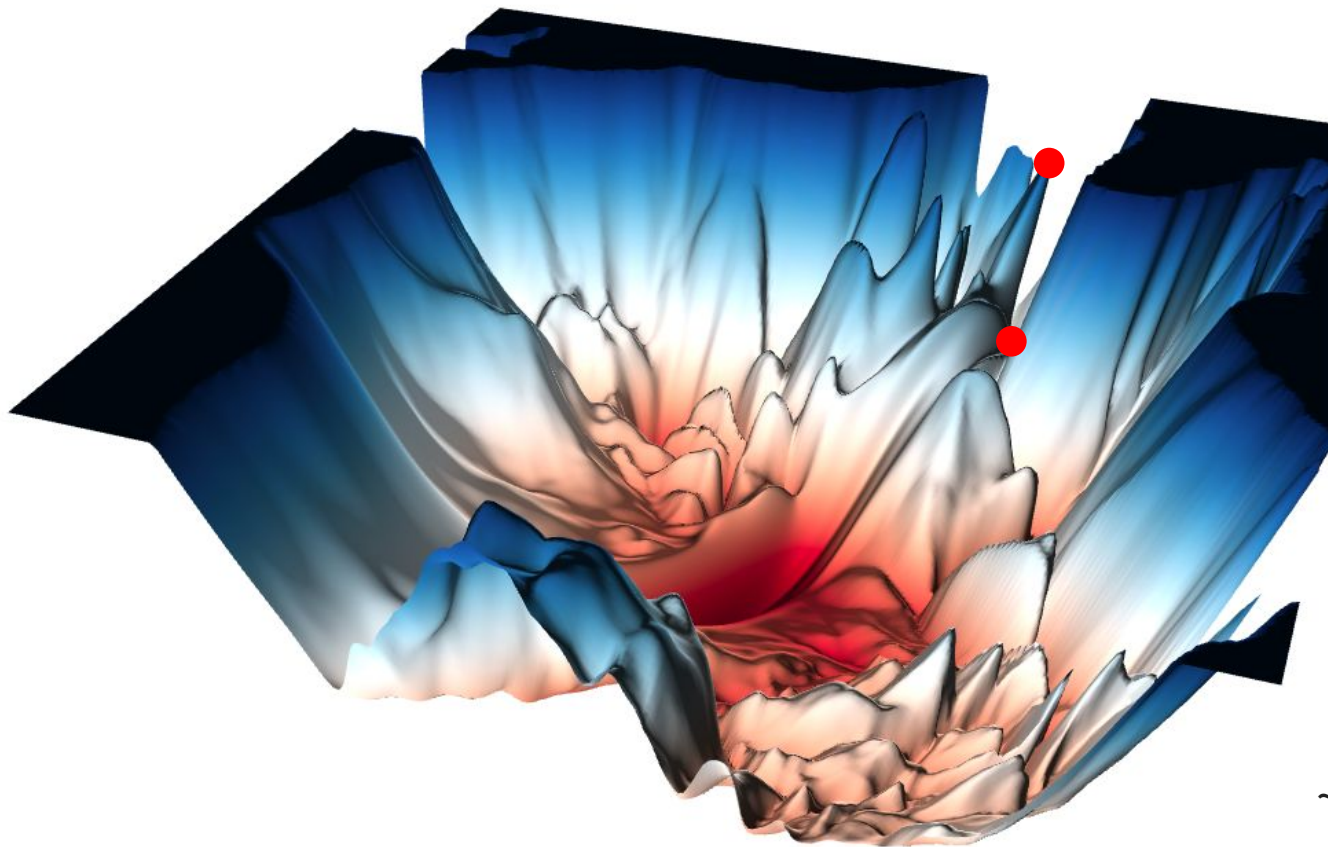


Neural networks are not “smooth”



~ Ankur Mohan

Neural networks are not “smooth”



~ Ankur Mohan

What do these irregularities mean?

A small perturbation on the input

(in a carefully chosen direction)

can change the output by a lot

Original



Infected



Each pixel is changed by 1 bit

Original



Panda

Infected



Chair (99.9%)

Pandas are nice,
but let's break something cooler.

Self-driving cars!

Disclaimer

For legal and practical reasons, the following attacks will be demonstrated on a simulated environment.

(but they also work on actual cars)



Objective:

Get this sign to be detected as “speed limit 50”

How do we find the correct perturbation?

How do we find the correct perturbation?

Use the gradient!

Partial derivative

$$\frac{\partial score}{\partial pixel_{x,y}}$$

Partial derivative

$$\frac{\partial score}{\partial pixel_{x,y}}$$

how will the **score** change

if

this **pixel** changes slightly

Partial derivative

$$\frac{\partial score}{\partial pixel_{x,y}}$$

how will the **score** change

if

this **pixel** changes slightly

We can also see it as a correlation coefficient

Gradient-based adversarial attacks

Change **all** pixels **slightly** in the gradient direction

OR

Change **a few** most significant pixels **all the way**

FGSM

Many small changes



Speed limit 50 (99.9%)

JSMA

Few large changes



Speed limit 50 (99.9%)

These attacks are not **robust**

They are very sensitive to:

- Camera noise
- Object angle
- Lighting
- Background

Infected image

raw file



Speed limit 50 (99.9%)

Infected image

smartphone photo



Stop (97.6%)

Infected image

original perturbation



Speed limit 50 (99.9%)

Infected image

same perturbation, different pose



Stop (98.0%)

2 conditions needed

- Bit-perfect input
- White-box access to the model

How do we make a better attack?

We want a perturbation that is:

- physically applied to the sign
- as stealthy as possible
- robust on a wide range of conditions
- without access to the model gradients

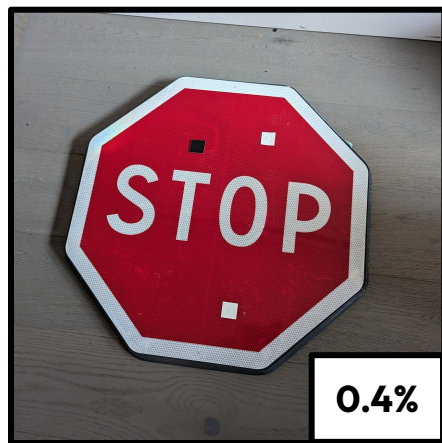


Noise resistance
can be achieved by
using larger stickers

Black-box attack

We don't have access to gradients, but we can still get the prediction score

→ Try many configurations and keep the best



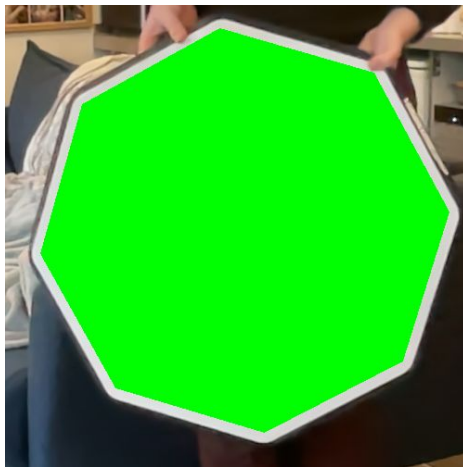
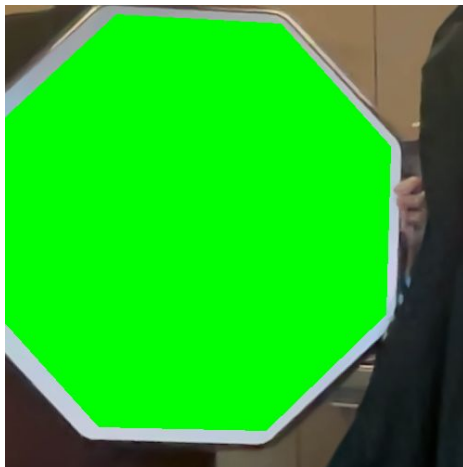
Black-box attack

Not very efficient to do manually

- Need to try millions of combinations
- Not very reproducible

→ Work in a **virtual** attack setup







The same perturbation can now be simulated over many different conditions

Eliminates all external variables

- Noise
- Position/angle
- Lighting and reflections
- Background

How do we search for good perturbations?

Bruteforce - try all possible combinations

Guaranteed to find the optimal solution

...but takes 10^{25} years 😱

How do we search for good perturbations?

Monte-Carlo - try random perturbations

Good perturbations are rare (~1 in a billion)

Need to be **very lucky**

How do we search for good perturbations?

Work step by step



How do we search for good perturbations?

Work step by step

Test many configurations
for sticker **1**, keep the best



How do we search for good perturbations?

Work step by step

Test many configurations
for sticker **2**, keep the best



How do we search for good perturbations?

Work step by step

Test many configurations
for sticker **3**, keep the best



How do we search for good perturbations?

Work step by step

etc.



How do we search for good perturbations?

Small improvement: add an **optimization** step to search locally for the best attack

Inspired by mutations in genetic algorithms

How do we search for good perturbations?

After adding each sticker, try moving the other ones slightly

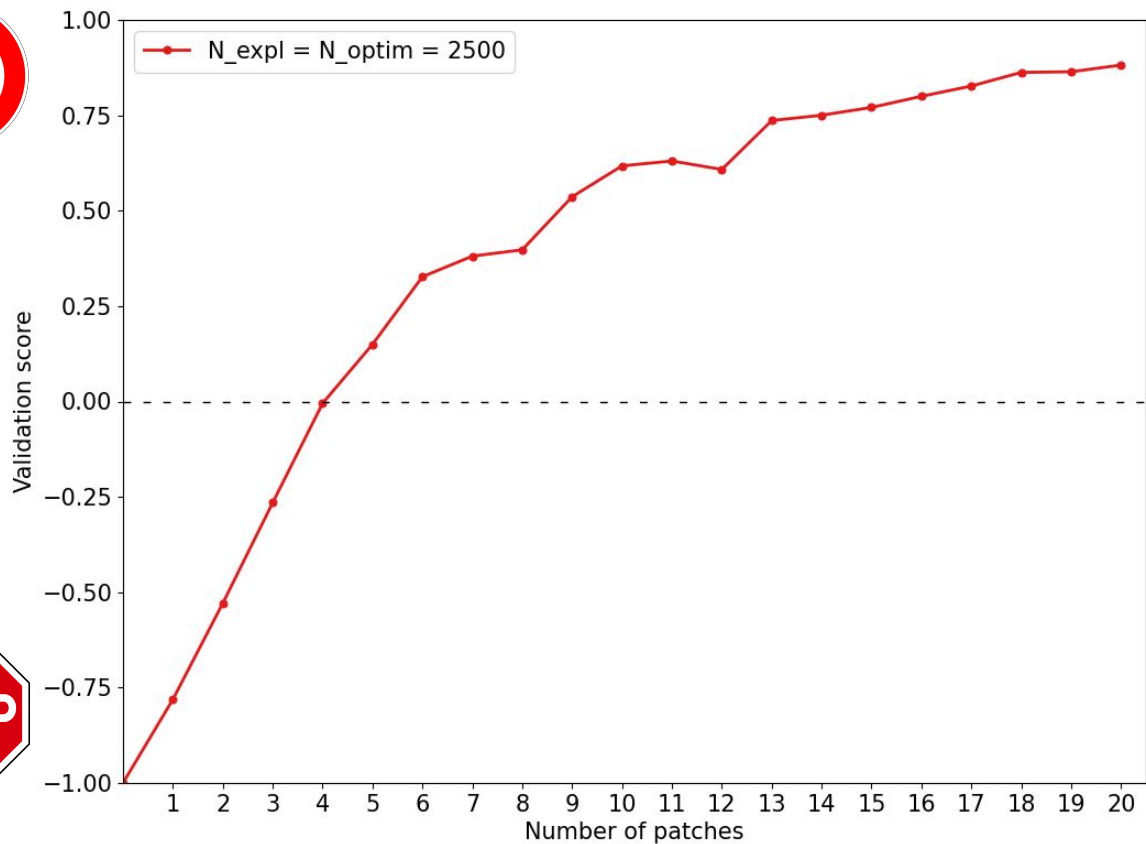


How do we search for good perturbations?

After adding each sticker, try moving the other ones slightly



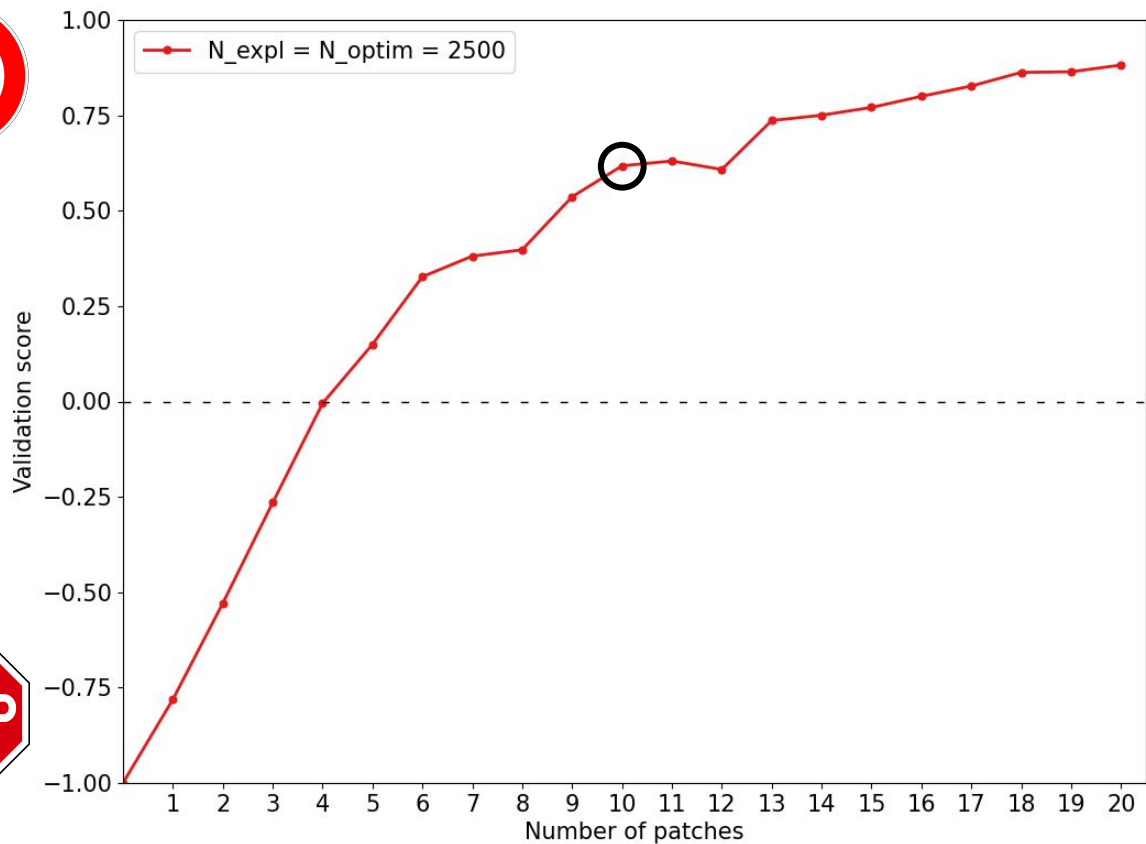
Attack performance @ 3min per patch



Attack performance @ 3min per patch

50

STOP



Demo time!



Slides



Paper

We want a perturbation that is:

- physically applied to the sign
- as stealthy as possible
- robust on a wide range of conditions
- without access to the model gradients

We want a perturbation that is:

- ✓ physically applied to the sign
- ? as stealthy as possible
- ✓ robust on a wide range of conditions
- ✓ without access to the model gradients

This PoC can be improved. Ideas:

- More stealth
- Black-box++

All neural networks are vulnerable to
adversarial attacks, and
no efficient protections currently exist.

Thank you!

Any questions?

Mathis Hammel

[@MathisHammel](#)

mathishammel.com

Independent consultant

- Technical competitions
- Algorithms
- Cybersecurity R&D



Slides



Paper