

# Od prostych sieci do GPT: Podróż przez wyzwania techniczne i etyczne w AI

---

Wojciech Szczepański



ChatGPT

30

11

2022

AI

Maszyna wykonująca ludzkie zadania poznawcze

Sprawienie by maszyny wykonywały zadania poznawcze,  
o których nie sądziliśmy, że są w stanie je wykonać

Formalizacja

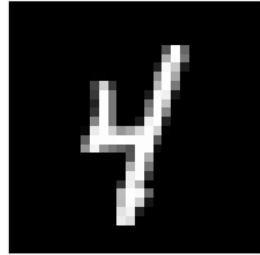


Napisanie programu w oparciu  
o wyprowadzone reguły

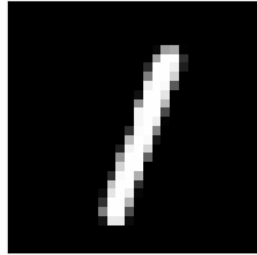
?



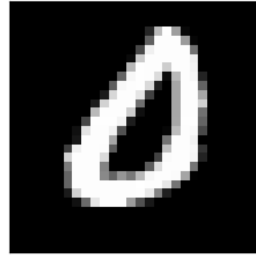
Napisanie programu w oparciu  
o wyprowadzone reguły



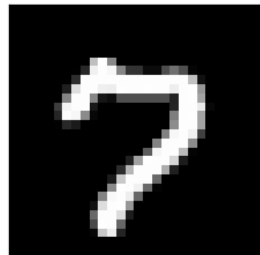
4 (4)



1 (1)



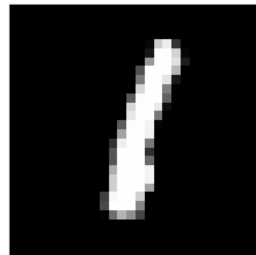
0 (0)



7 (7)



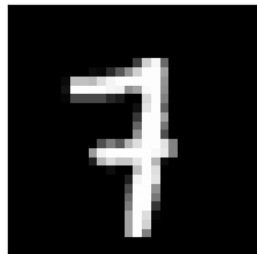
8 (8)



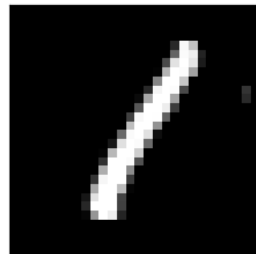
1 (1)



2 (2)



7 (7)



1 (1)









Znajdźmy dużo przykładów

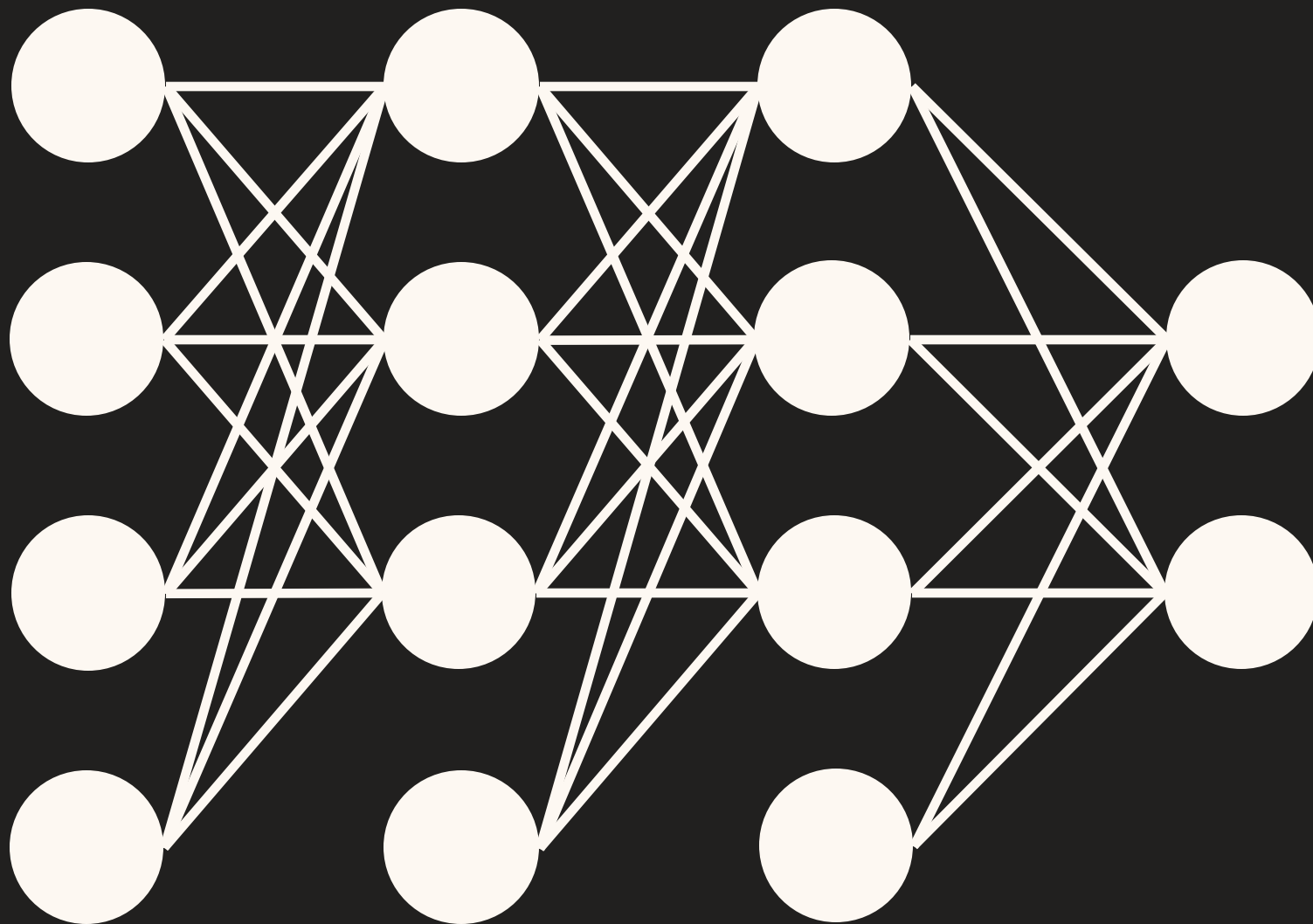


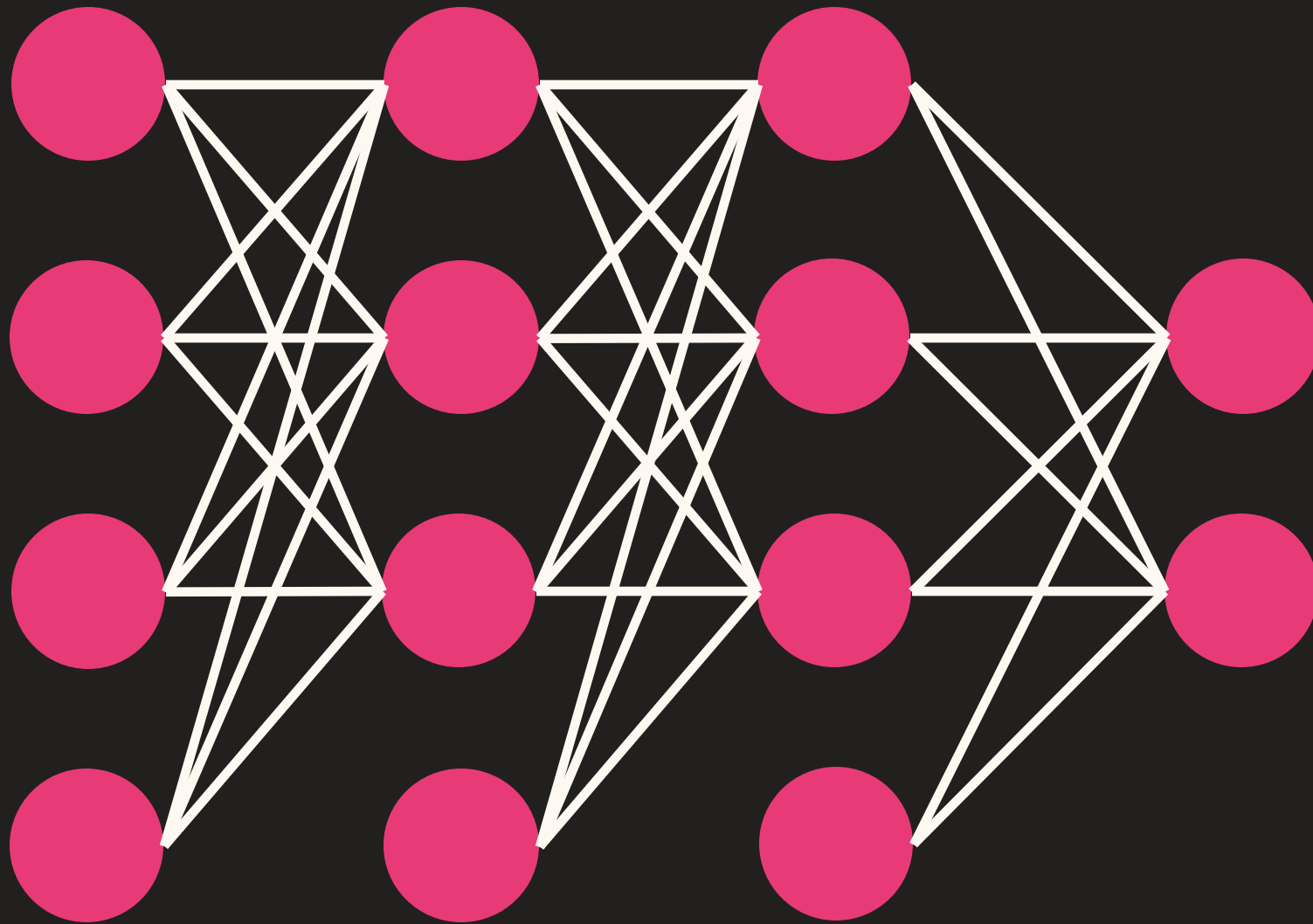
Nauczmy komputer  
samemu dobrać reguły

Sieć neuronowa

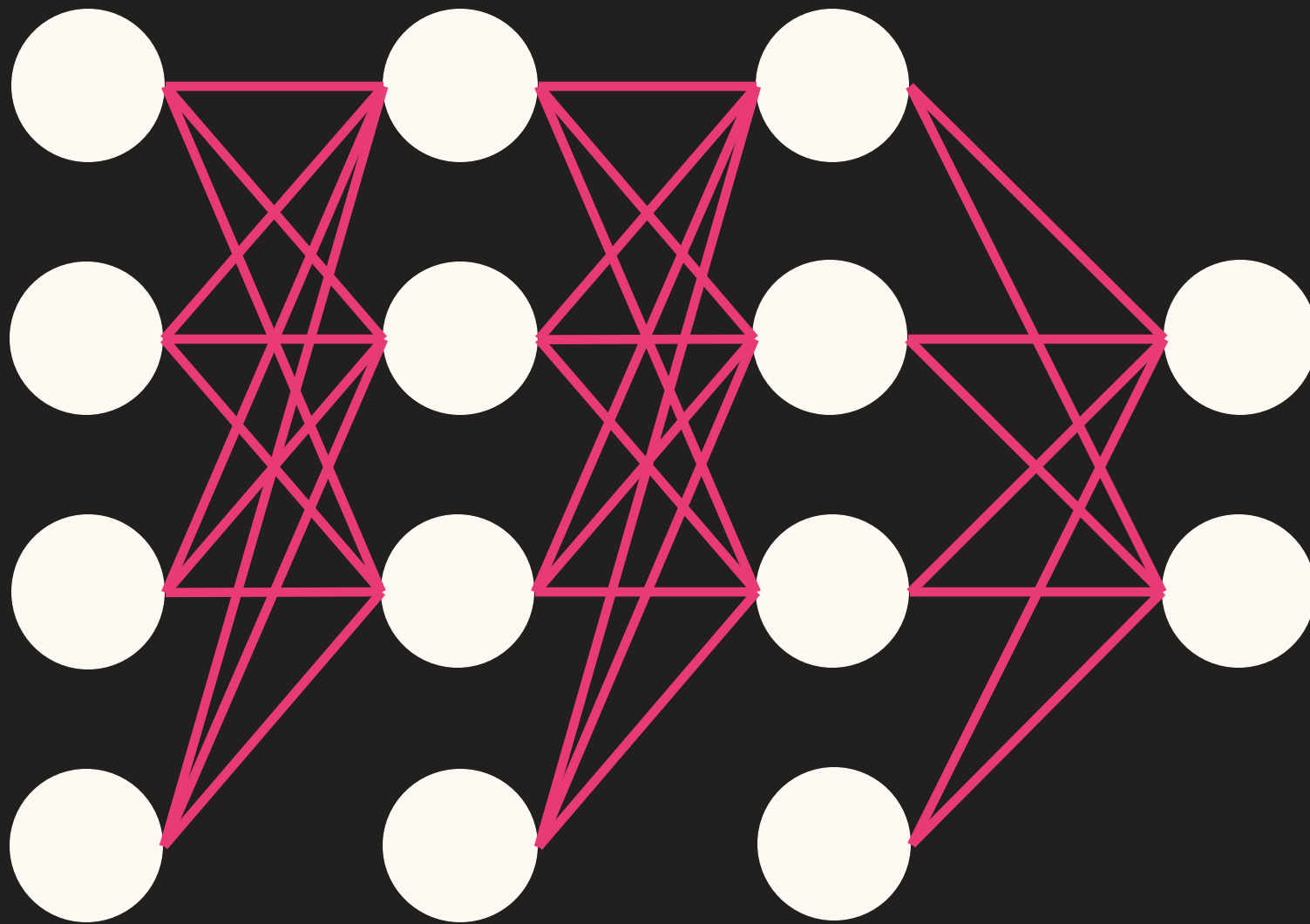
# AKT I

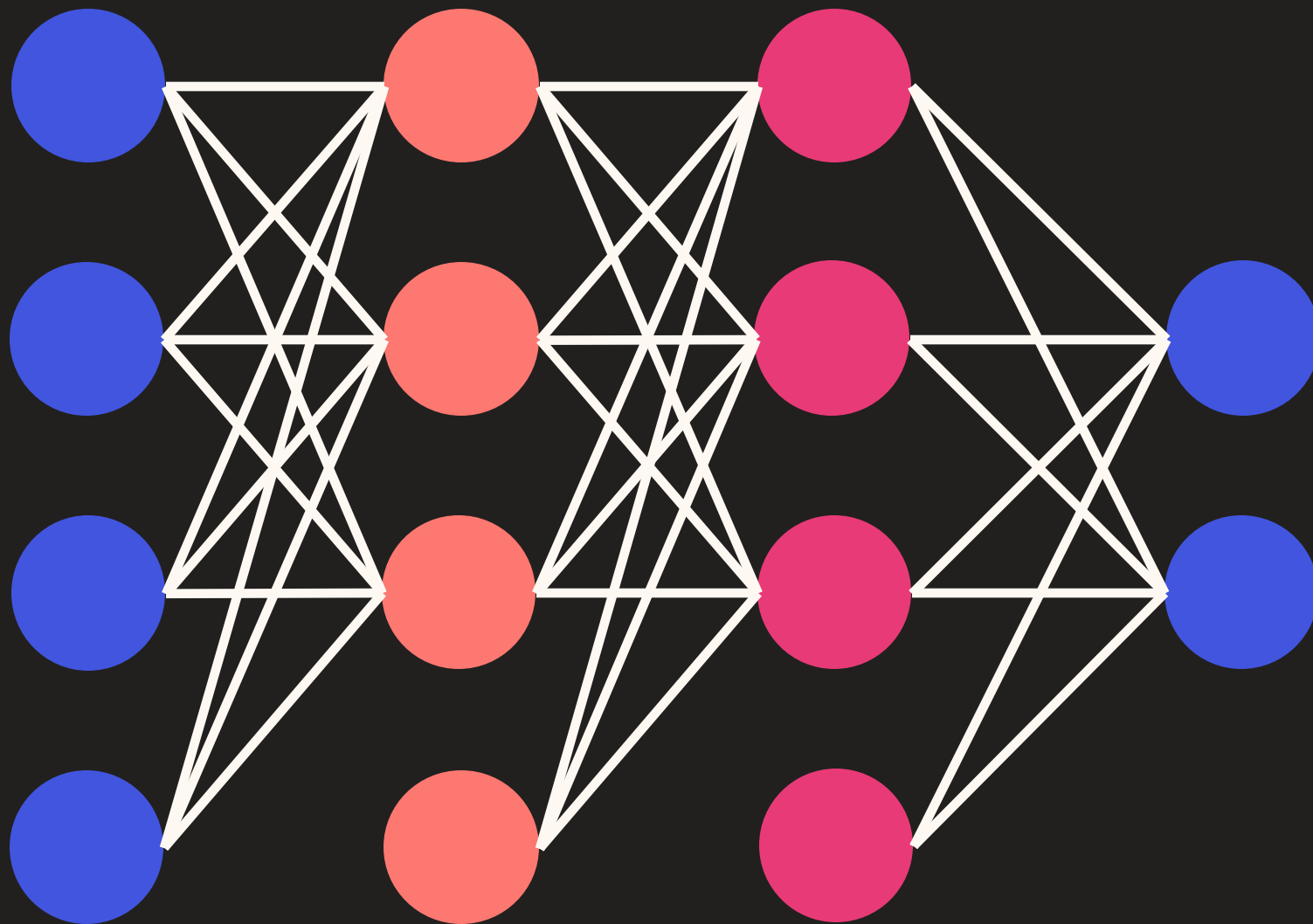
*Na początku był neuron...*

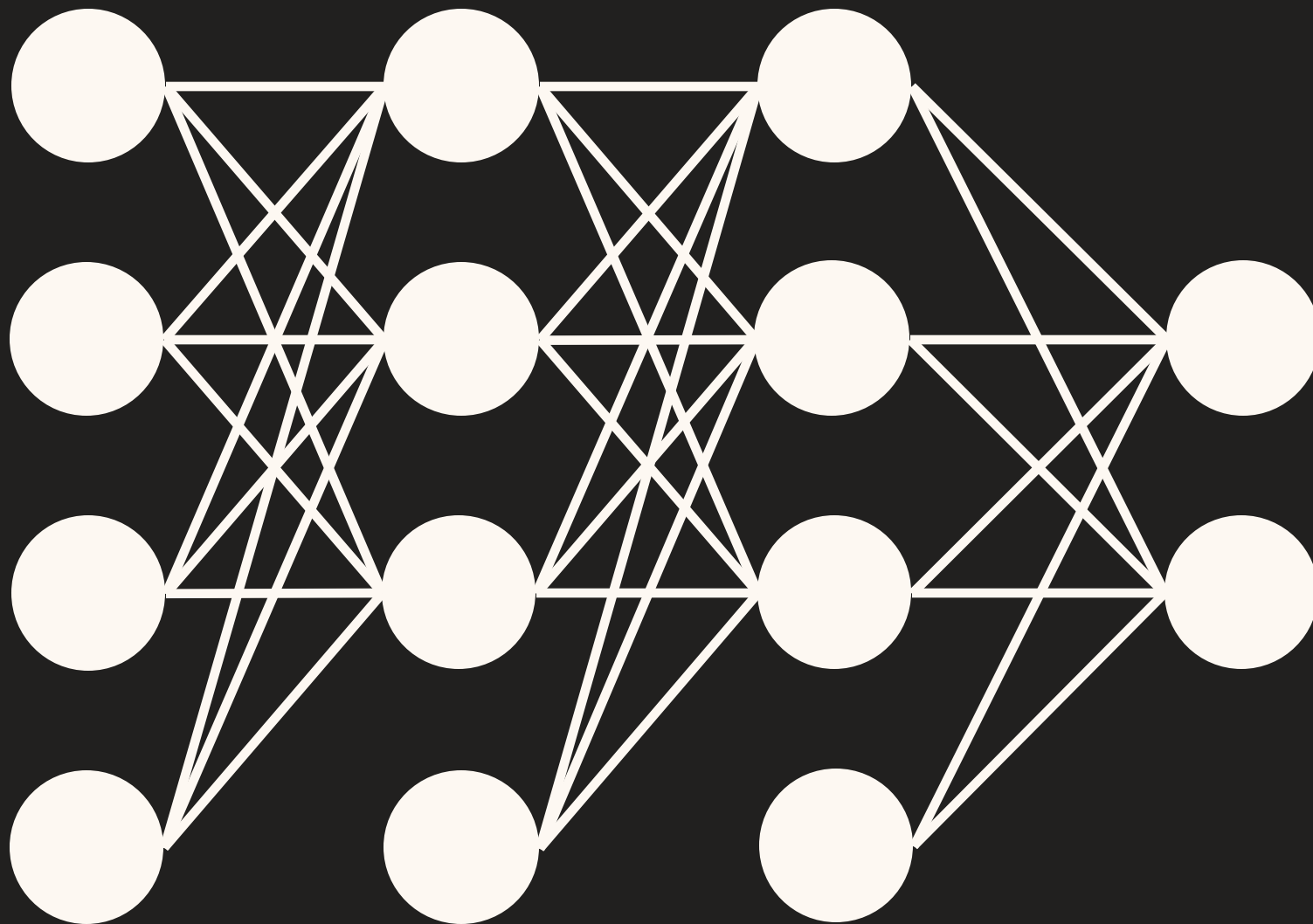


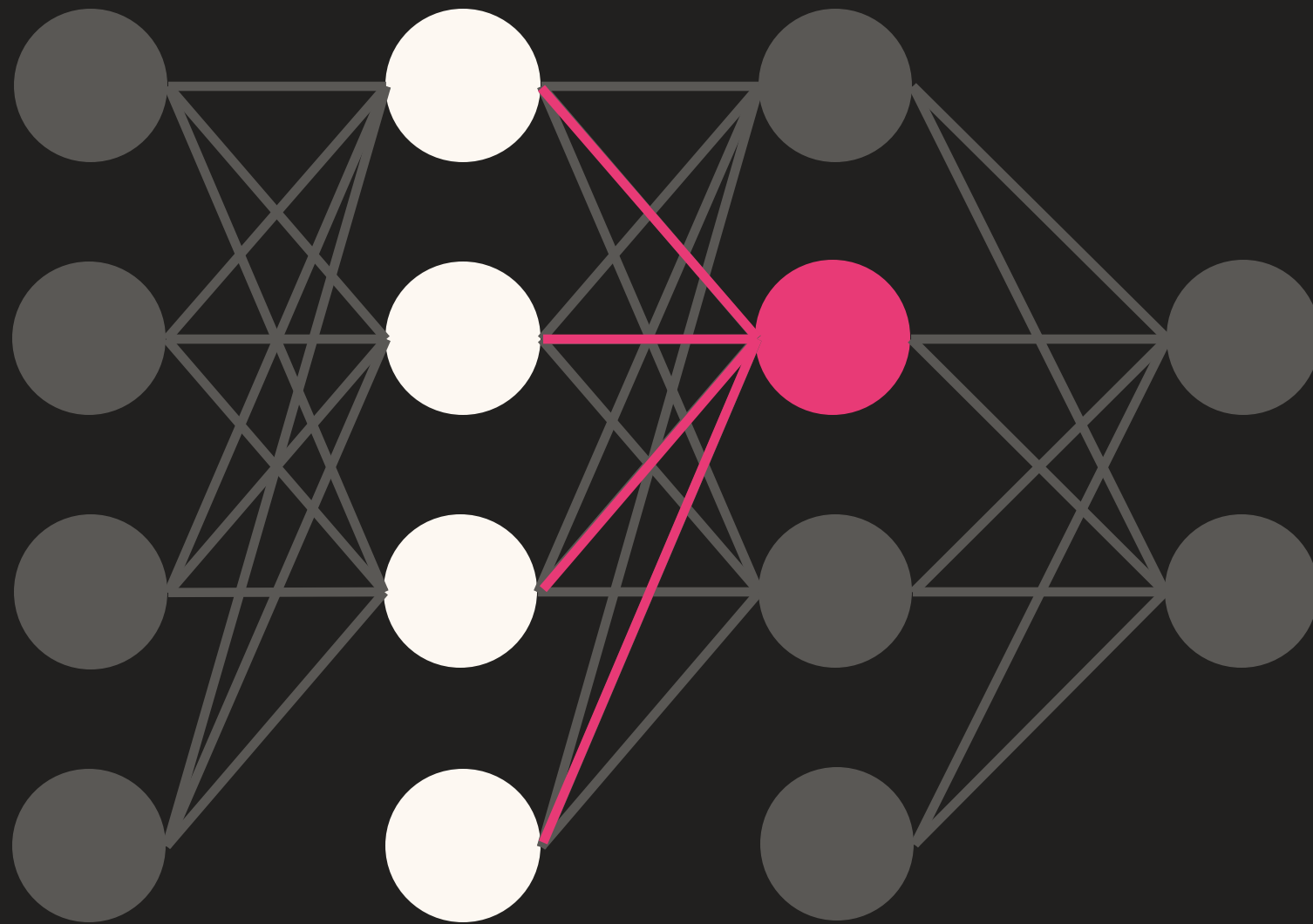


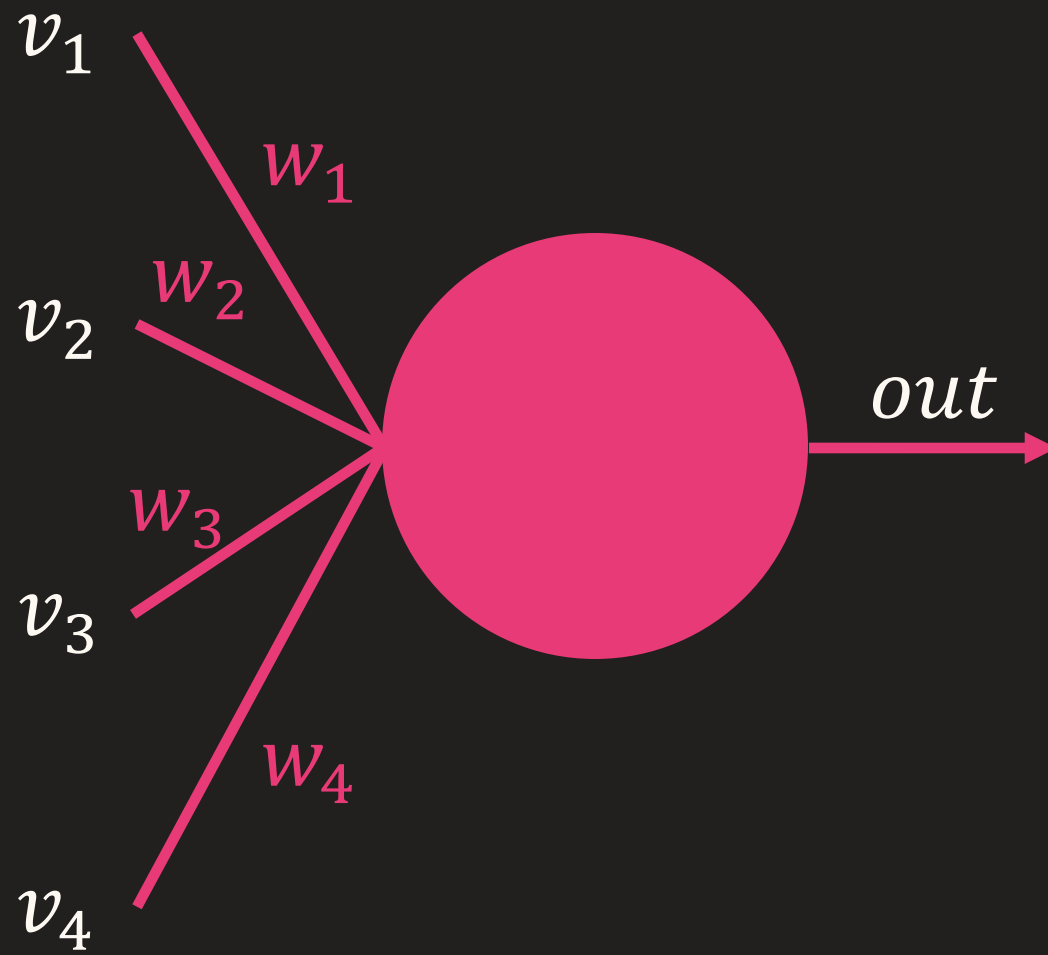


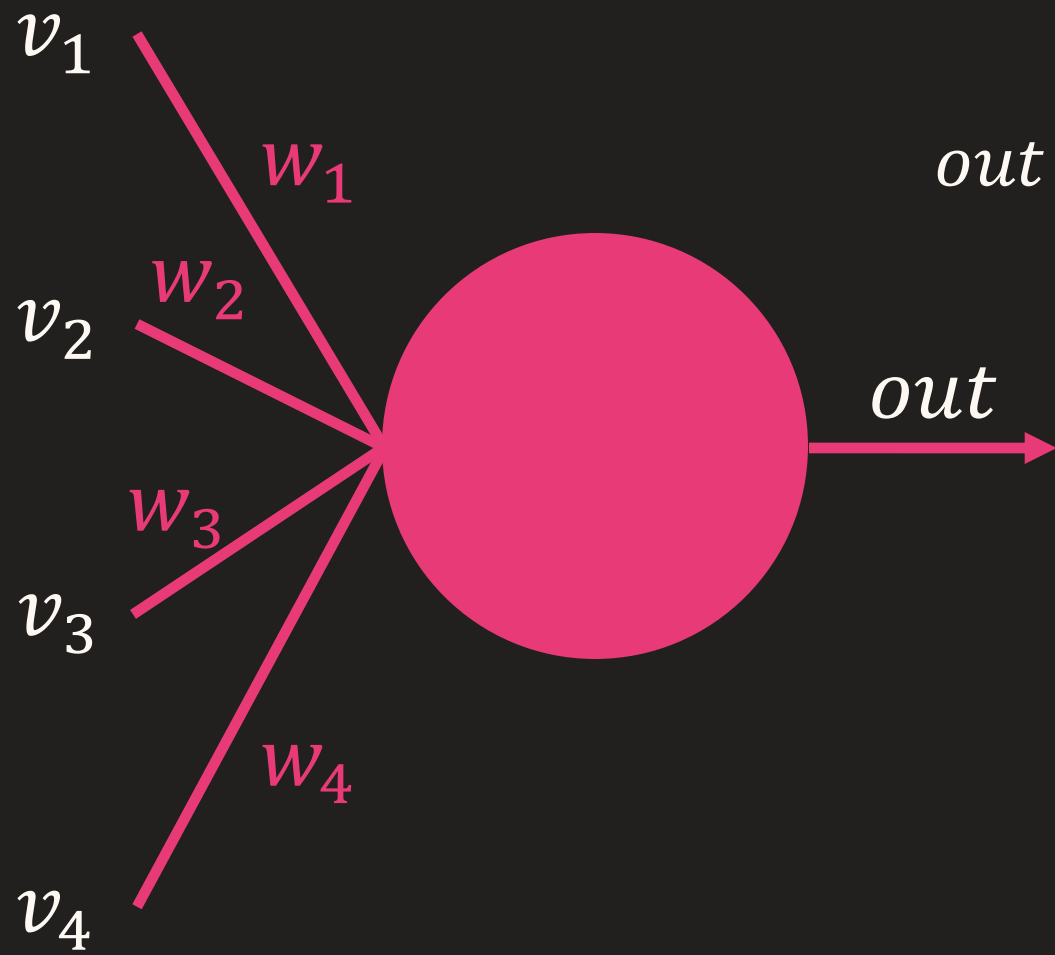




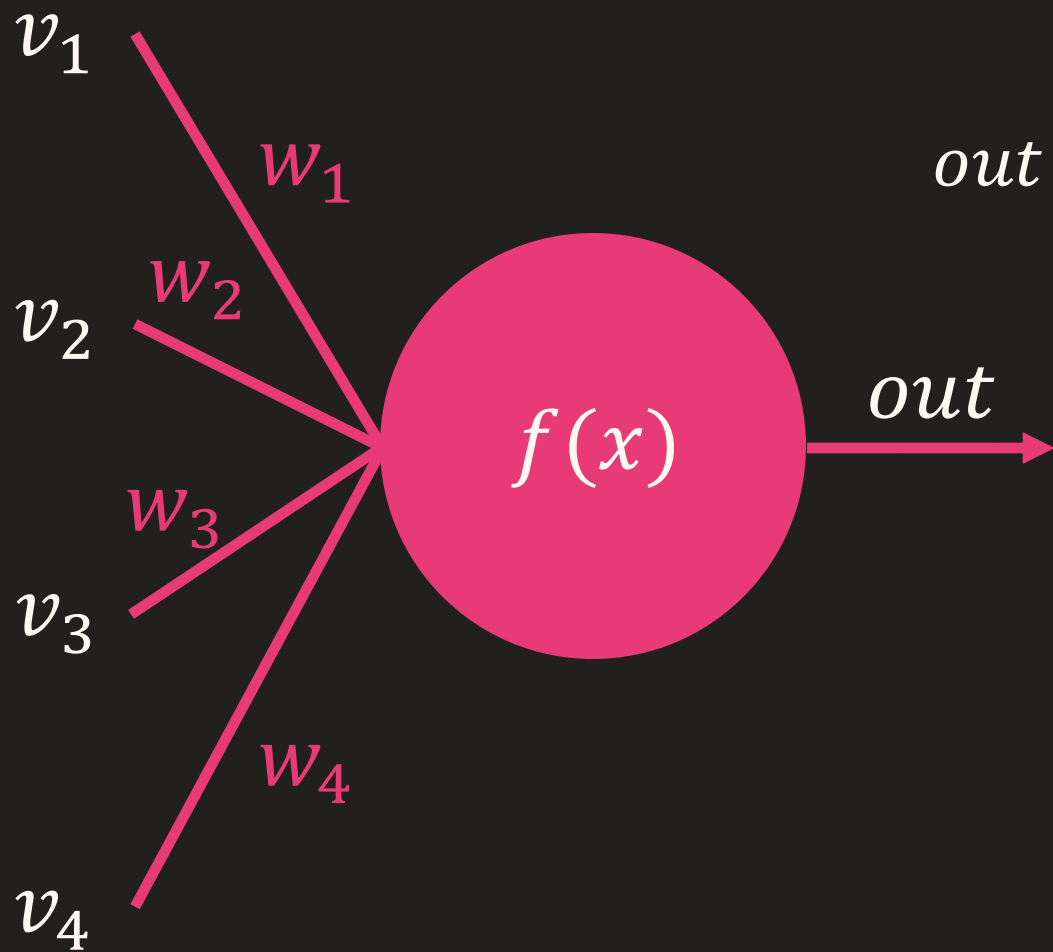








$$out = w_1 * v_1 + w_2 * v_2 + w_3 * v_3 + w_4 * v_4$$



$$out = w_1 * v_1 + w_2 * v_2 + w_3 * v_3 + w_4 * v_4$$

$$out = f(w_1 v_1 + w_2 v_2 + w_3 v_3 + w_4 v_4)$$

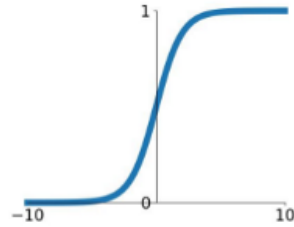
$$out = f(w_1 v_1 + \dots + w_n v_n)$$

$$out = f\left(\sum_i^n w_i v_i\right)$$

# Activation Functions

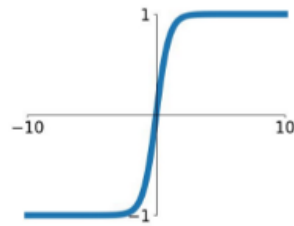
## Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



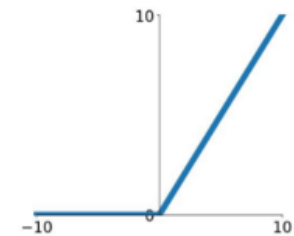
## tanh

$$\tanh(x)$$



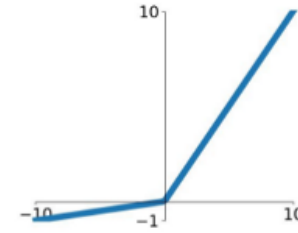
## ReLU

$$\max(0, x)$$



## Leaky ReLU

$$\max(0.1x, x)$$

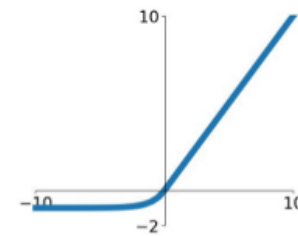


## Maxout

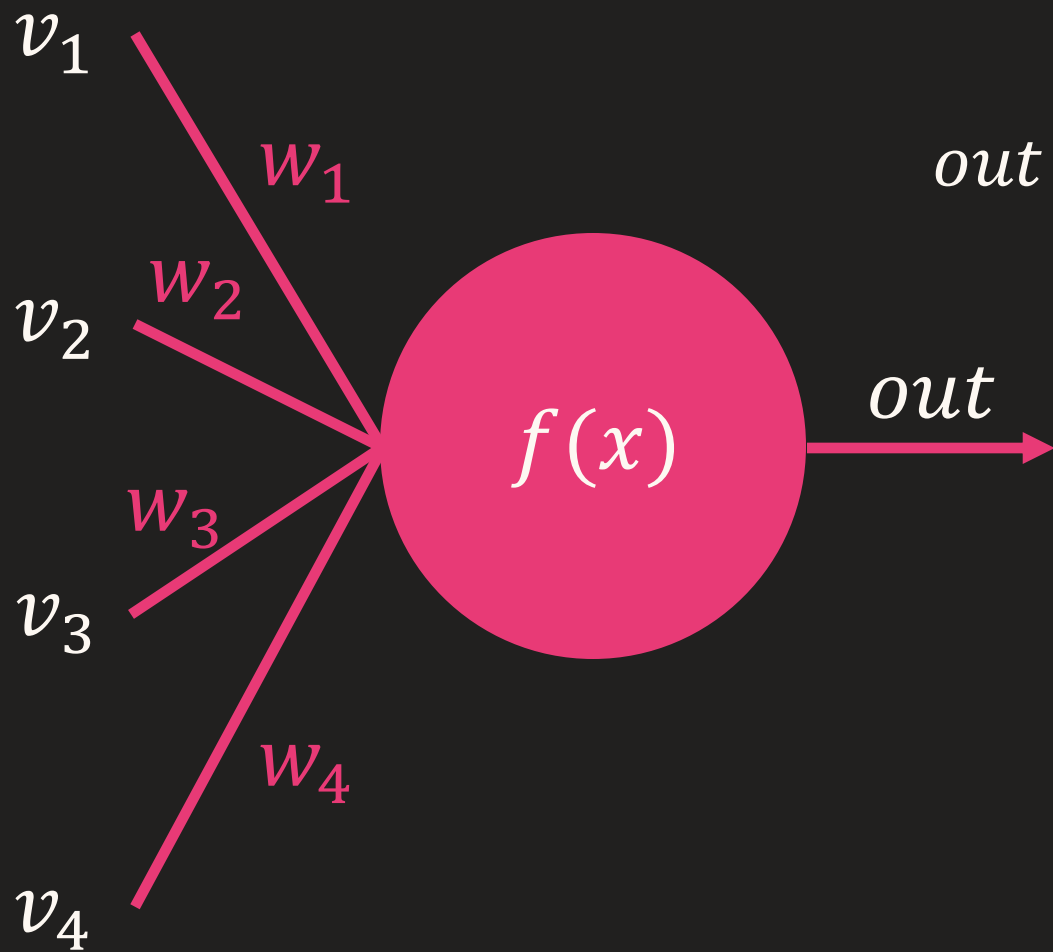
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

## ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$





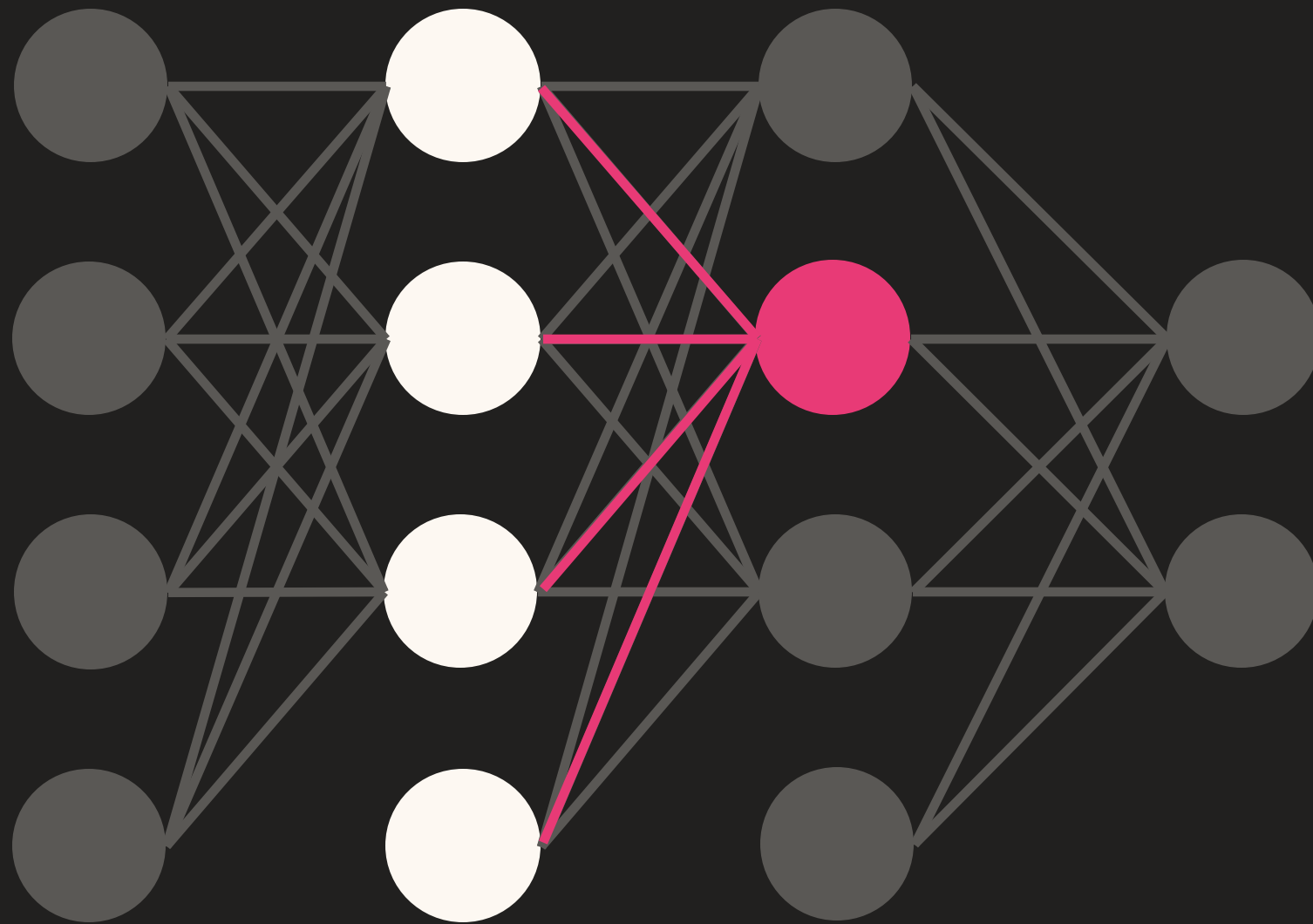


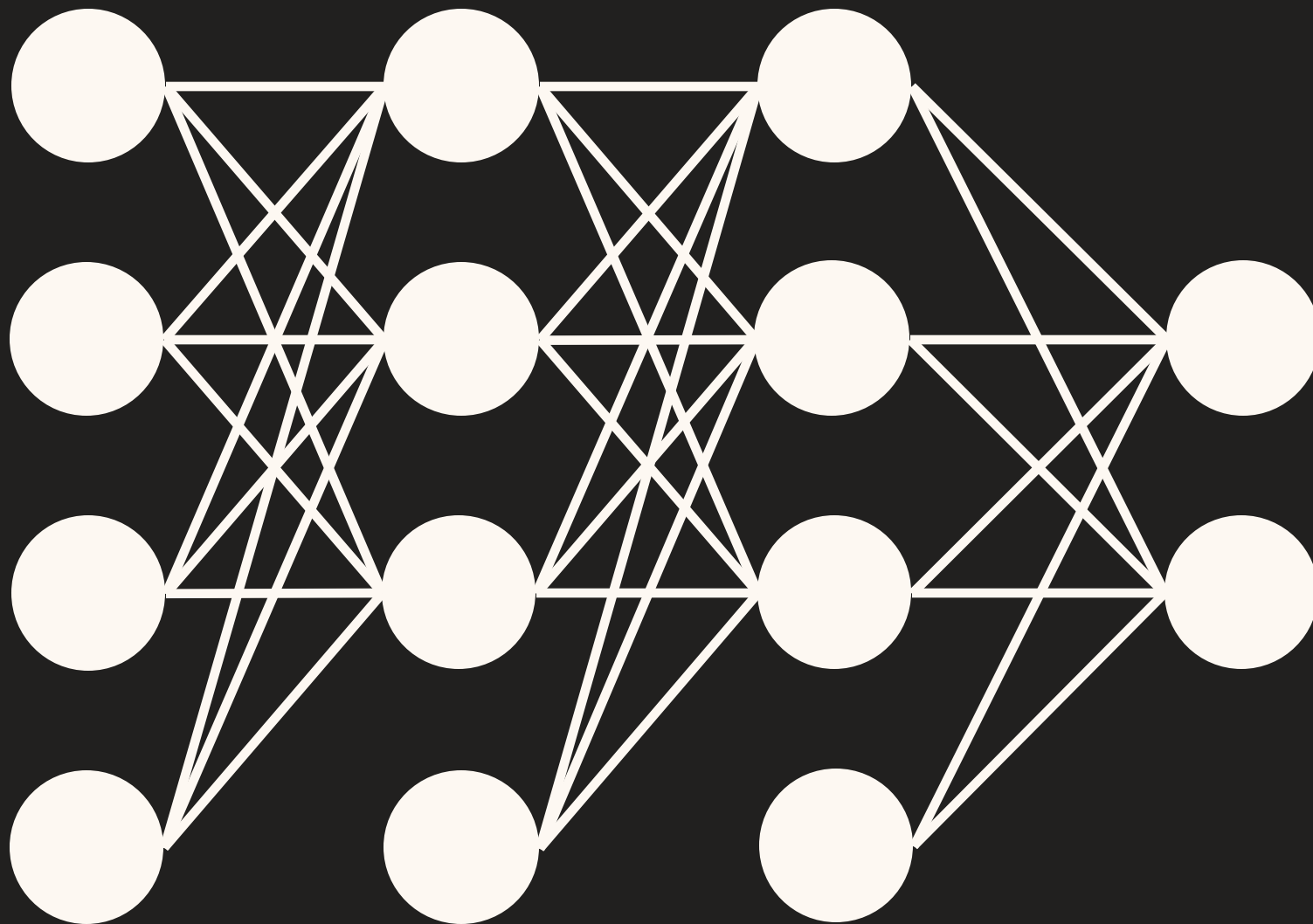
$$out = w_1 * v_1 + w_2 * v_2 + w_3 * v_3 + w_4 * v_4$$

$$out = f(w_1 v_1 + w_2 v_2 + w_3 v_3 + w_4 v_4)$$

$$out = f(w_1 v_1 + \dots + w_n v_n)$$

$$out = f\left(\sum_i^n w_i v_i\right)$$





Jak sieć neuronowa się uczy?



$$\begin{bmatrix} 132 & \dots & 67 \\ \vdots & \ddots & \vdots \\ 15 & \dots & 42 \end{bmatrix}$$

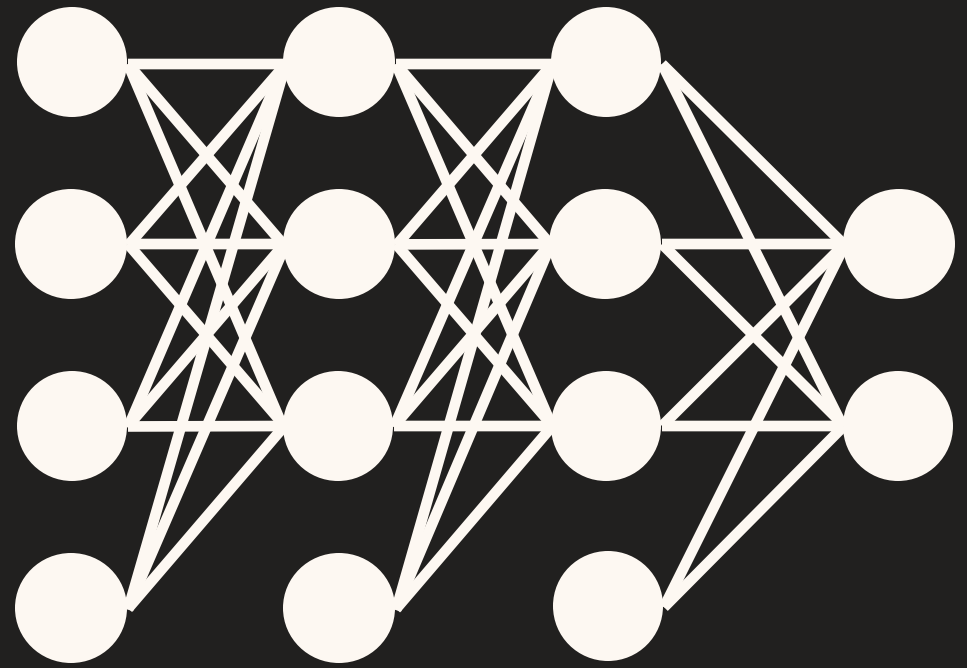


$$\begin{bmatrix} 0.52 & \dots & 0.26 \\ \vdots & \ddots & \vdots \\ 0.06 & \dots & 0.16 \end{bmatrix}$$

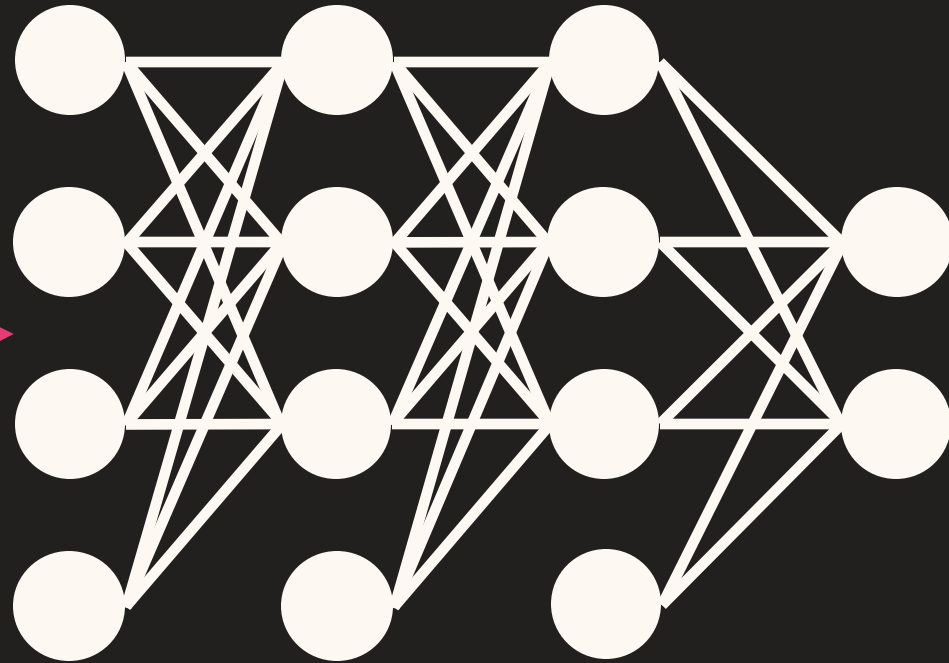
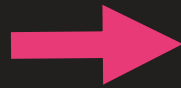
$\{Kot, Pies\}$  

$\begin{bmatrix} 1 \\ 0 \end{bmatrix}$  - Kot

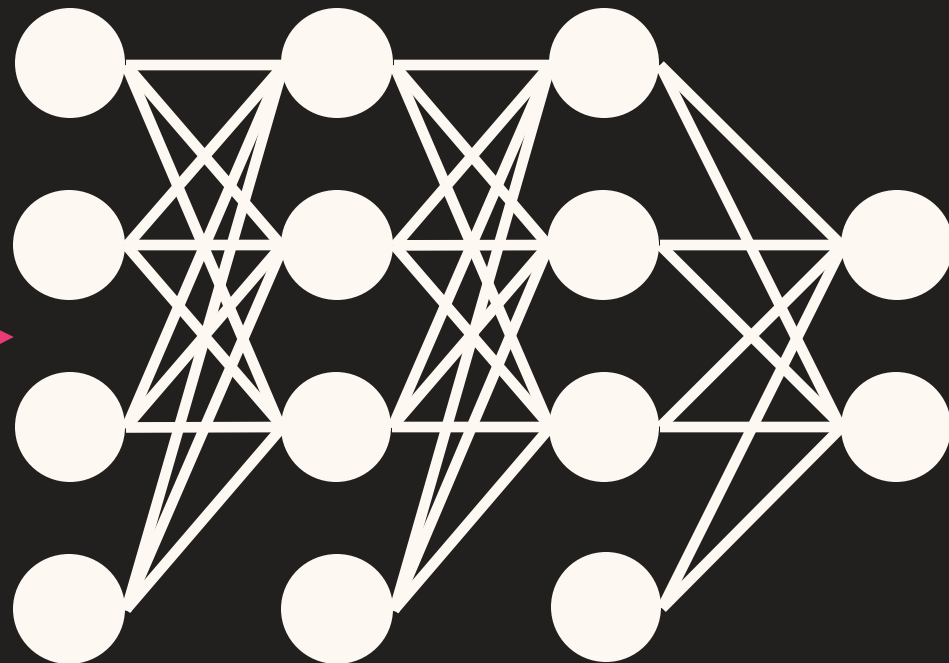
$\begin{bmatrix} 0 \\ 1 \end{bmatrix}$  - Pies



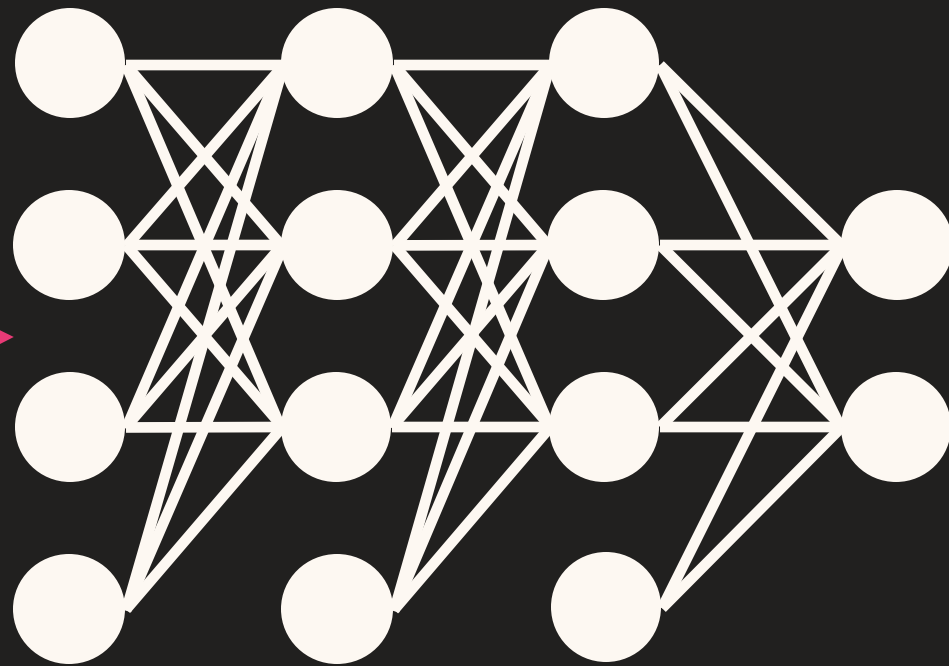




$$\begin{bmatrix} 0.52 & \dots & 0.26 \\ \vdots & \ddots & \vdots \\ 0.06 & \dots & 0.16 \end{bmatrix}$$

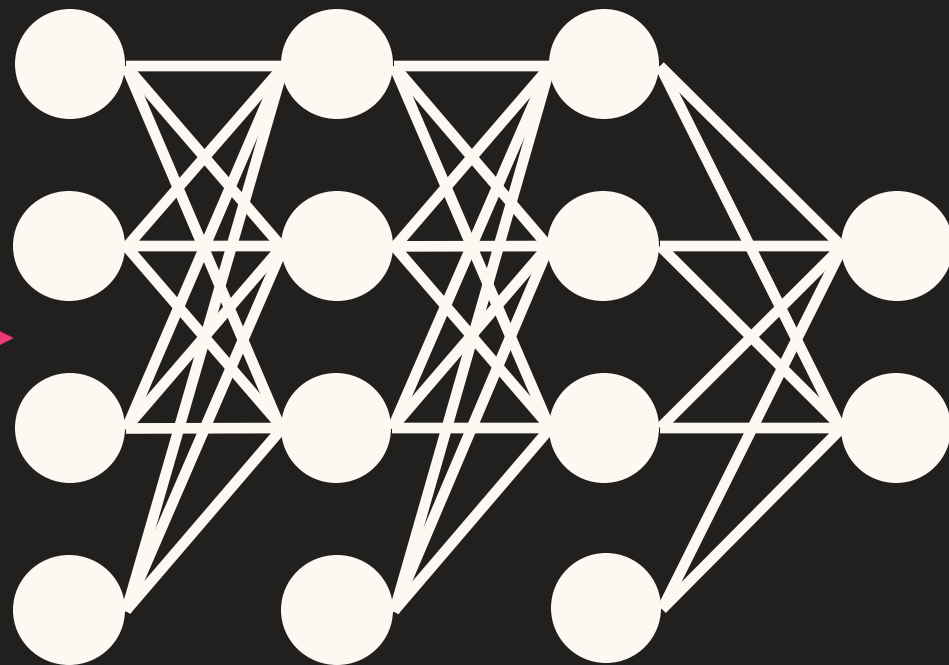


$$\begin{bmatrix} 0.52 & \dots & 0.26 \\ \vdots & \ddots & \vdots \\ 0.06 & \dots & 0.16 \end{bmatrix}$$



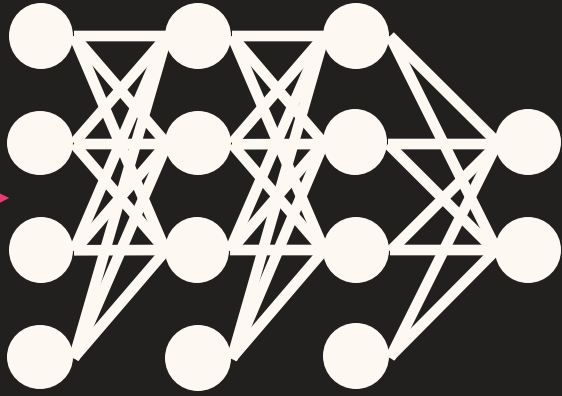
$$\begin{bmatrix} 0.3 \\ 0.7 \end{bmatrix}$$

$$\begin{bmatrix} 0.52 & \dots & 0.26 \\ \vdots & \ddots & \vdots \\ 0.06 & \dots & 0.16 \end{bmatrix}$$

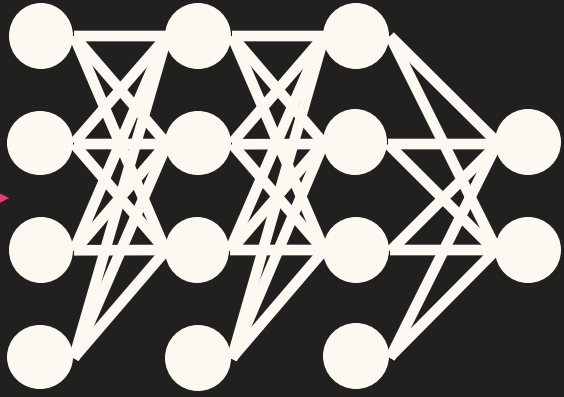


$$\begin{bmatrix} 0.3 \\ 0.7 \end{bmatrix} - \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

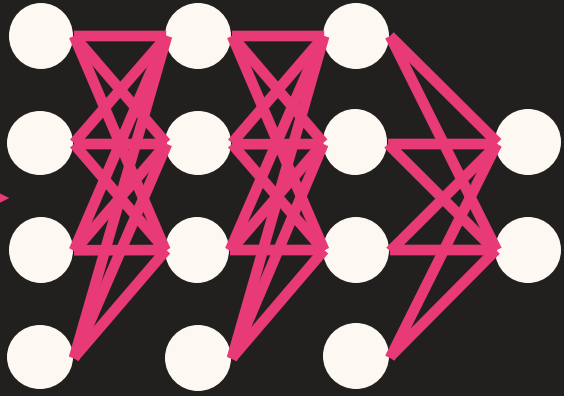
$$\begin{bmatrix} 0.52 & \dots & 0.26 \\ \vdots & \ddots & \vdots \\ 0.06 & \dots & 0.16 \end{bmatrix}$$



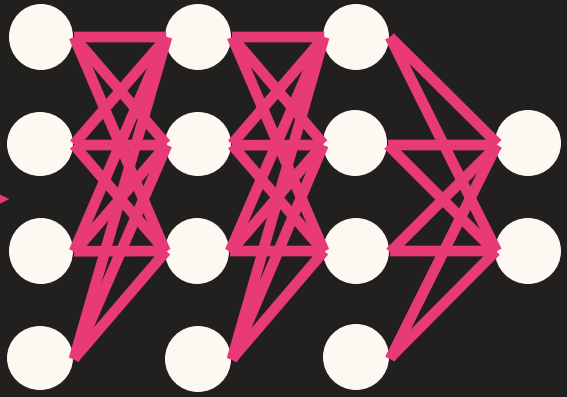
$$L(p, y) = \begin{bmatrix} 0.3 \\ 0.7 \end{bmatrix} - \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} 0.52 & \dots & 0.26 \\ \vdots & \ddots & \vdots \\ 0.06 & \dots & 0.16 \end{bmatrix}$$


$$L(p, y) = p - y$$

$$\begin{bmatrix} 0.52 & \dots & 0.26 \\ \vdots & \ddots & \vdots \\ 0.06 & \dots & 0.16 \end{bmatrix}$$


$$\downarrow L(p, y) = p - y$$

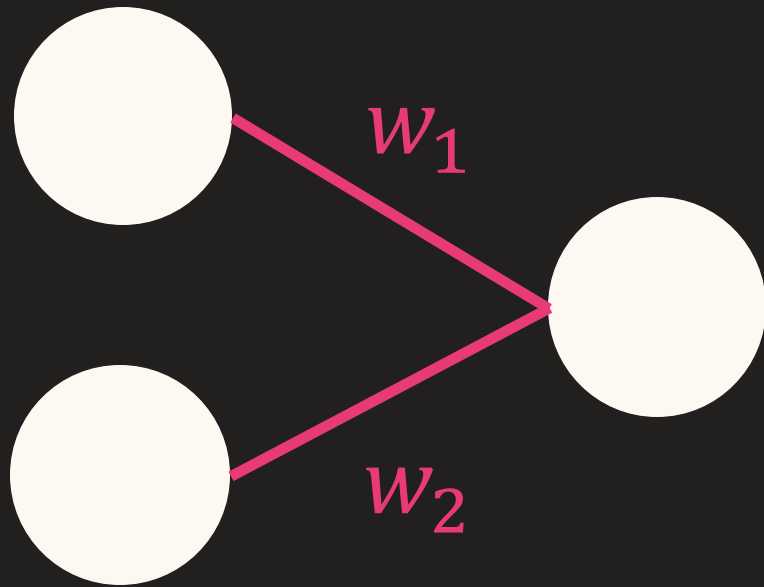
$$\begin{bmatrix} 0.52 & \dots & 0.26 \\ \vdots & \ddots & \vdots \\ 0.06 & \dots & 0.16 \end{bmatrix}$$


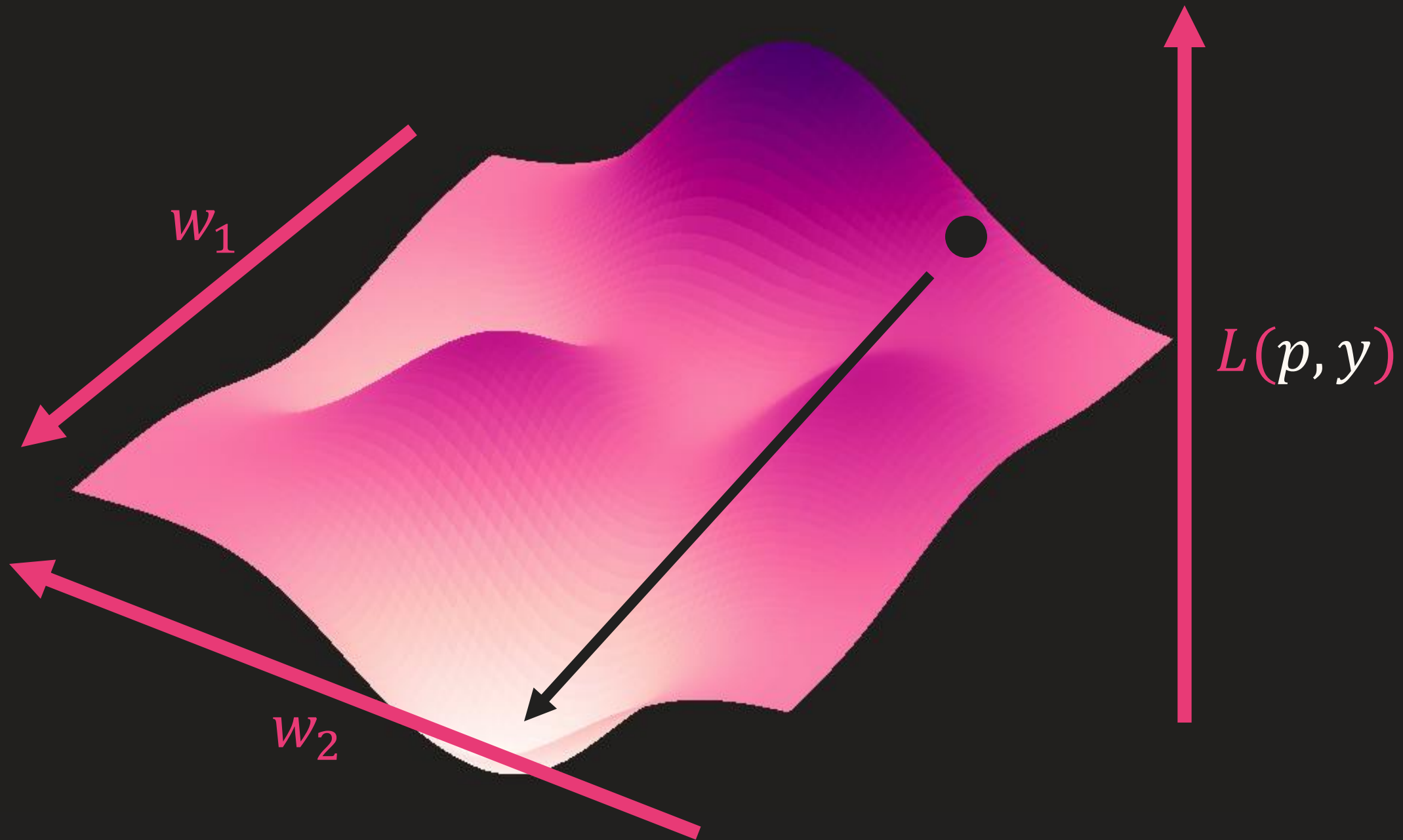
Parametry modelu



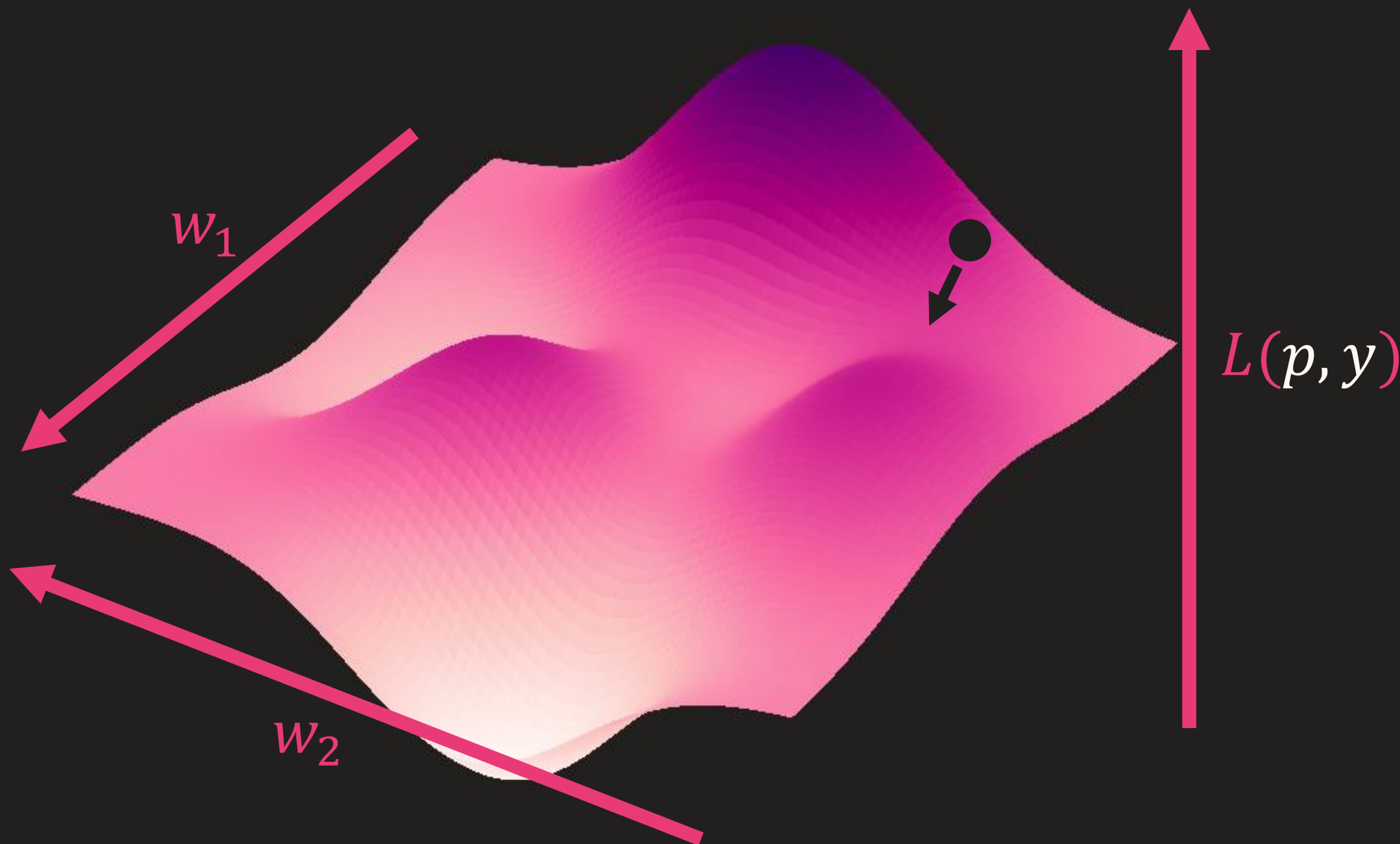
$$L(p, y) = p - y$$

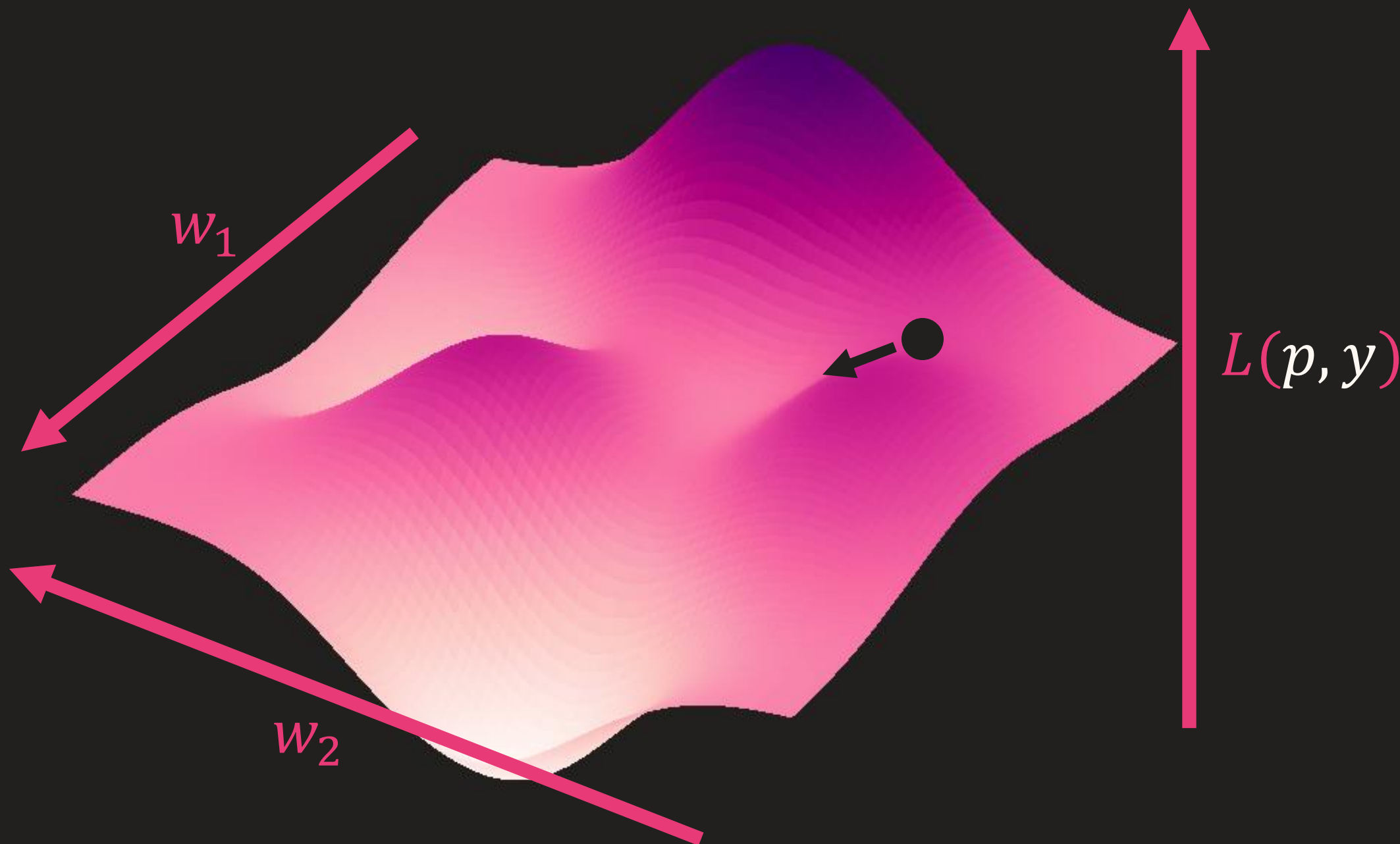


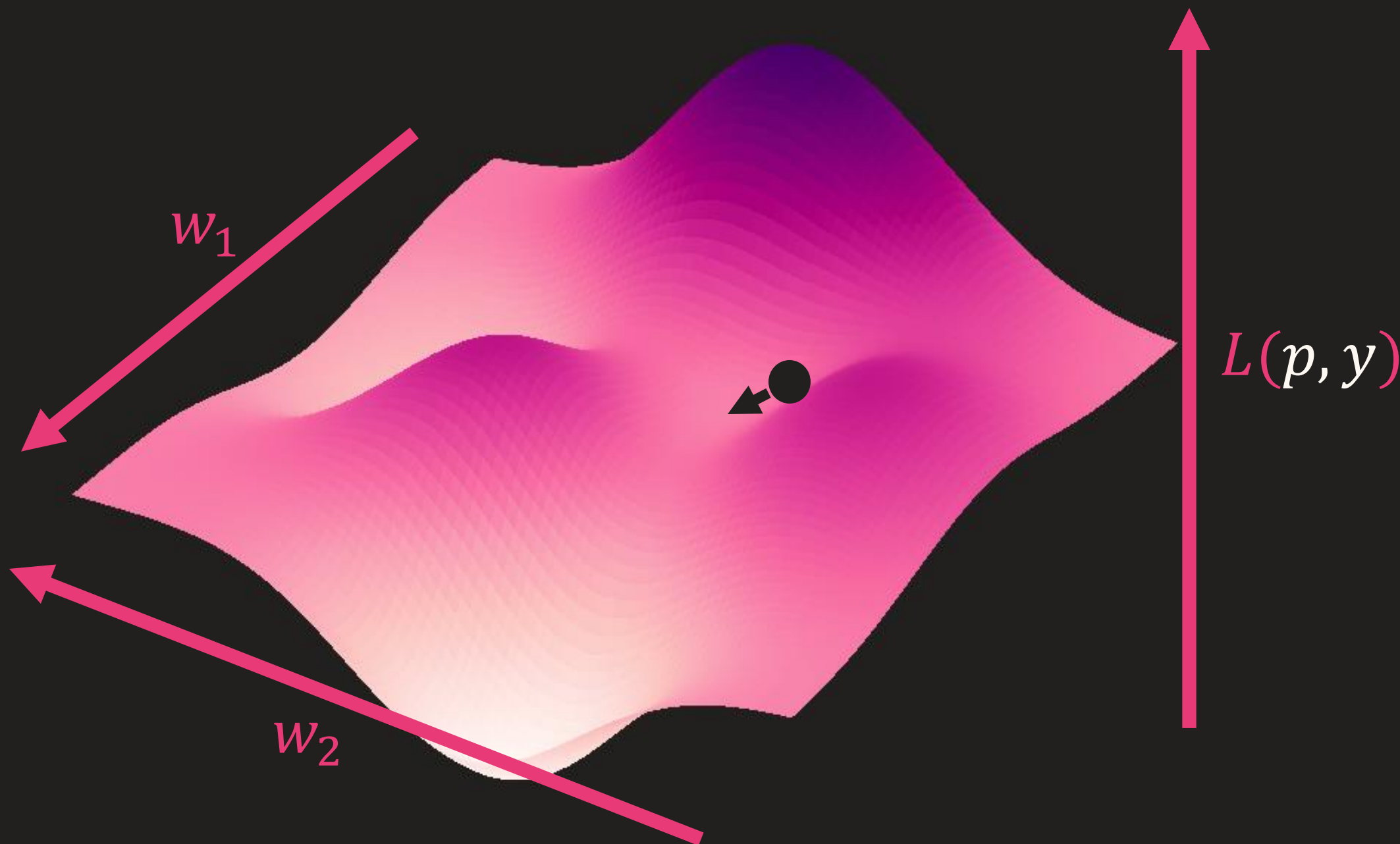


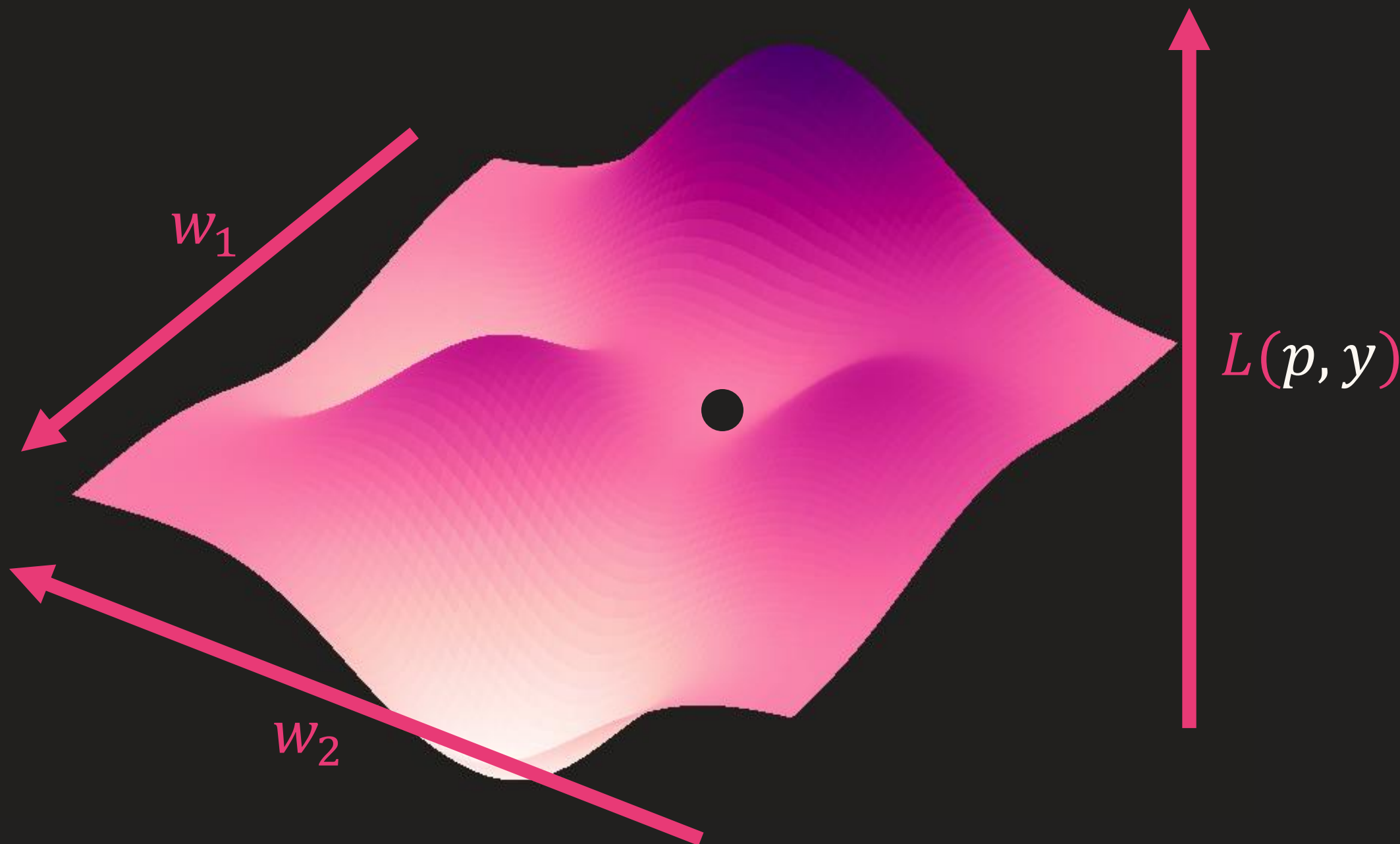








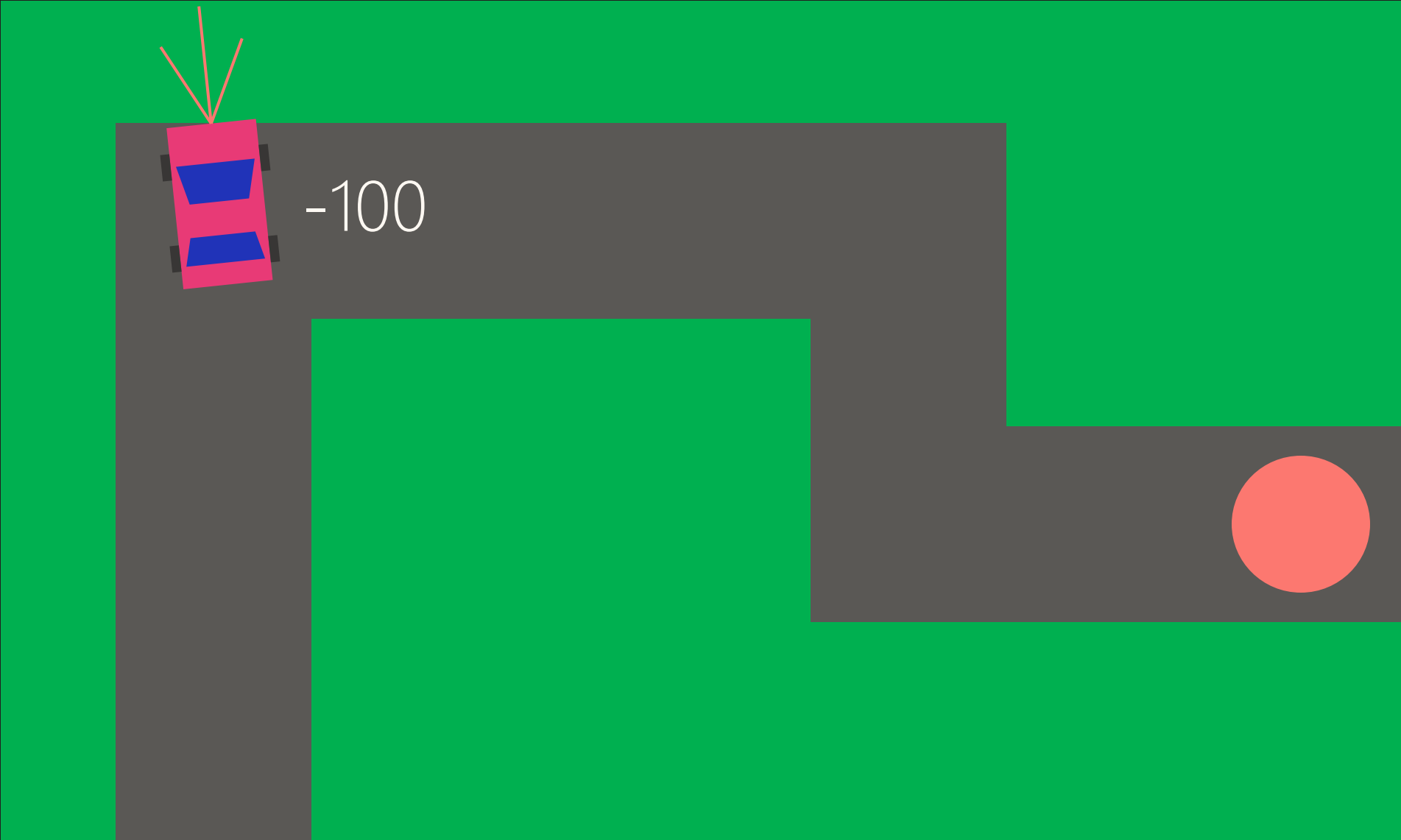


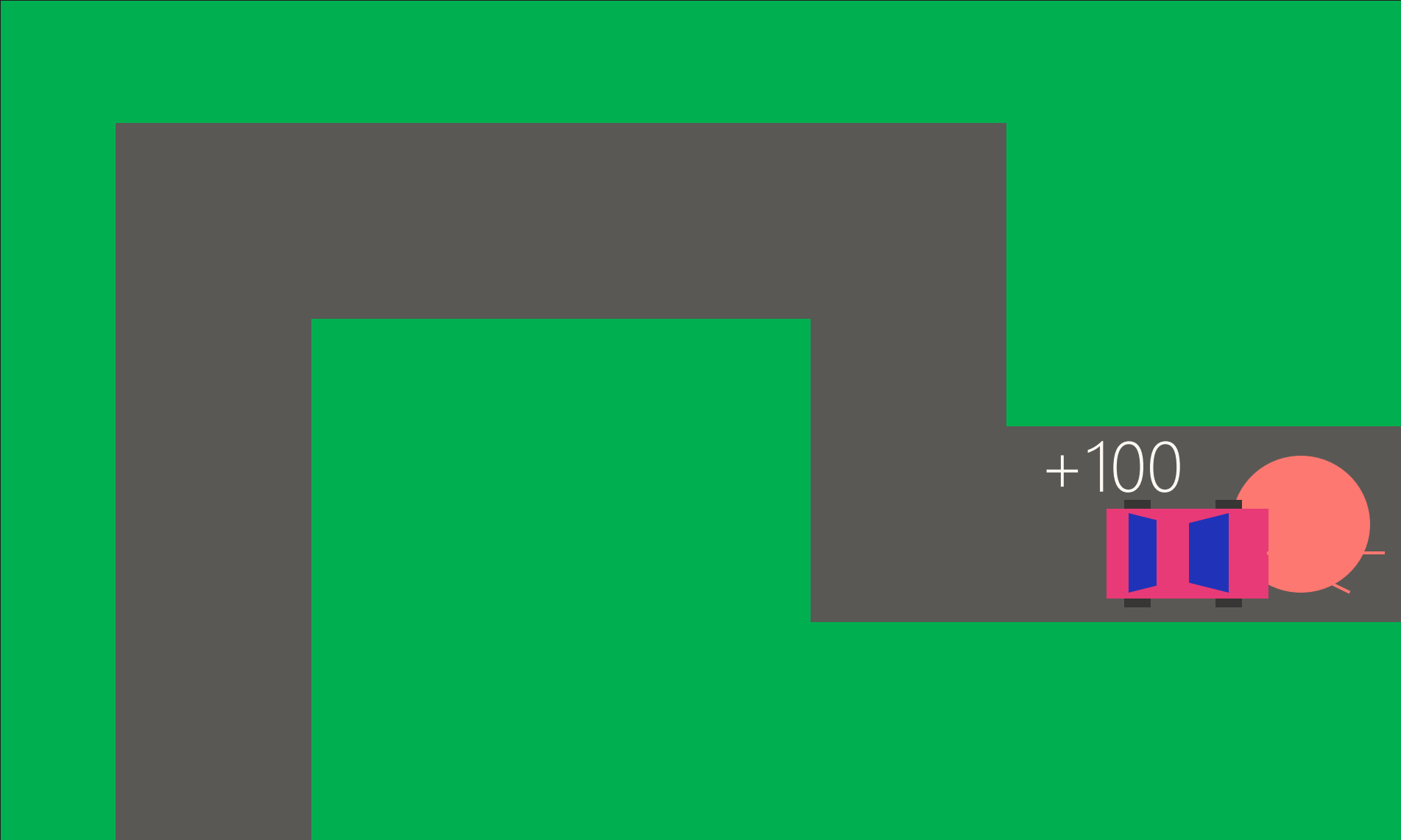




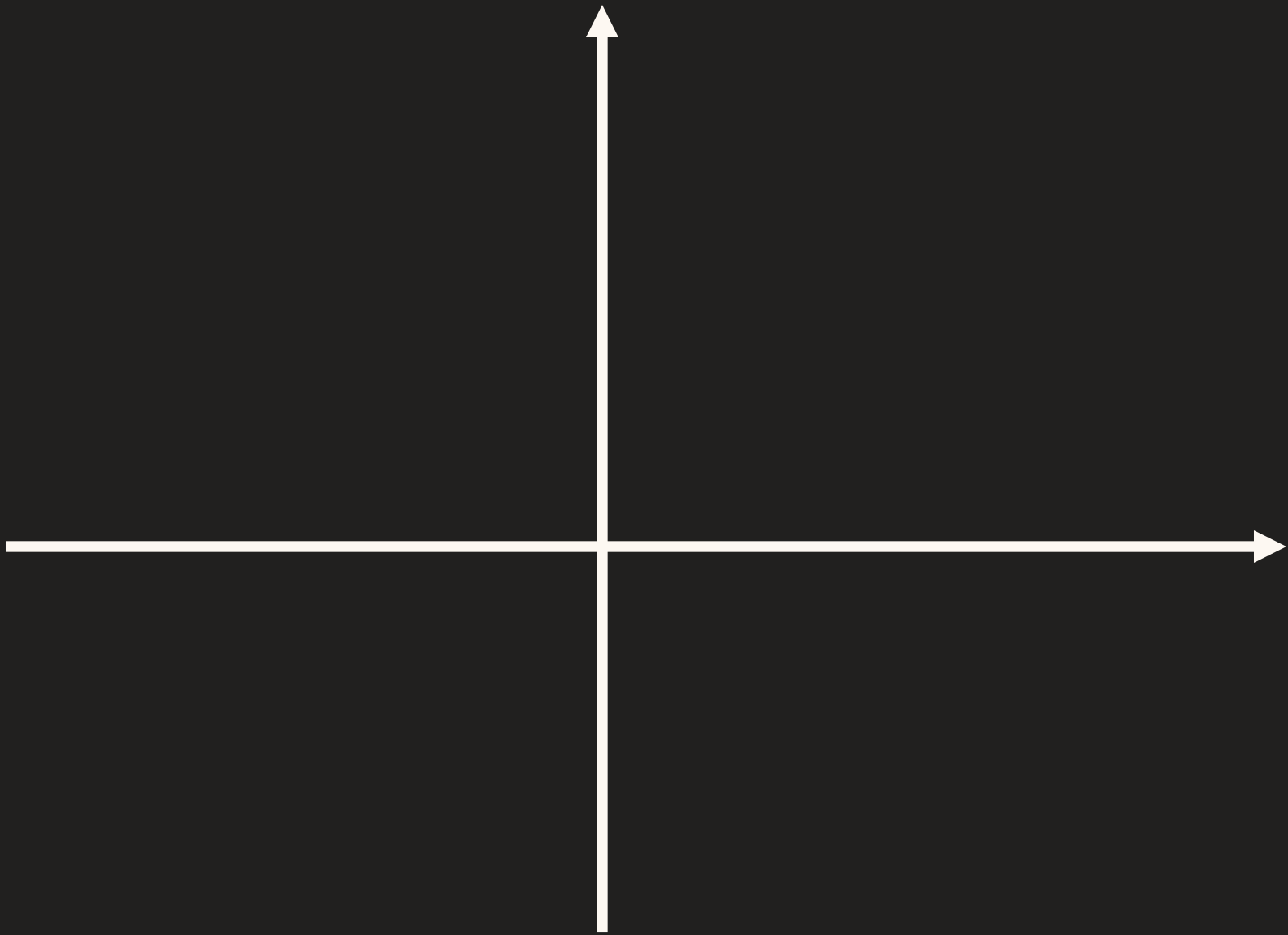




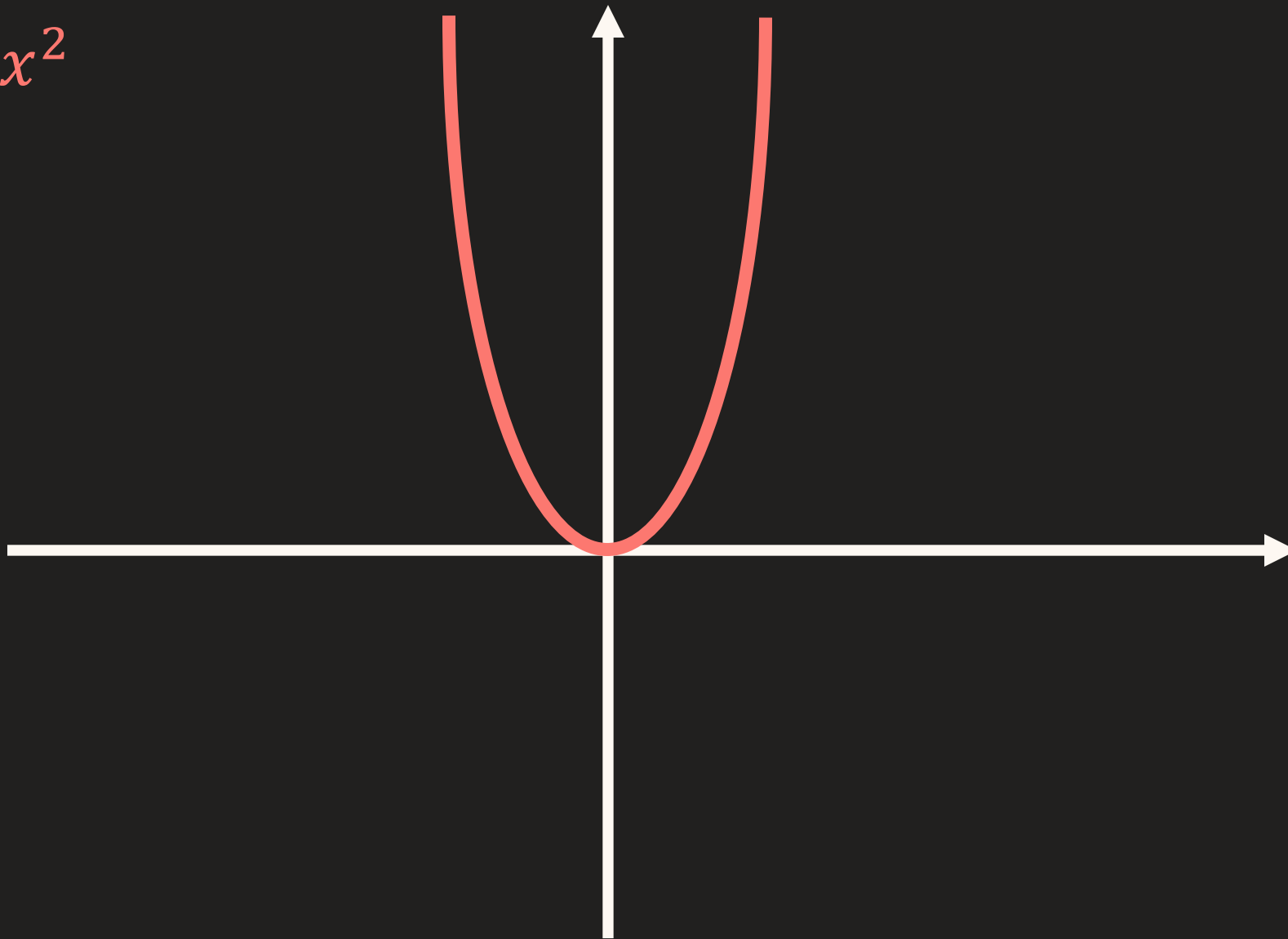




Generalizacja

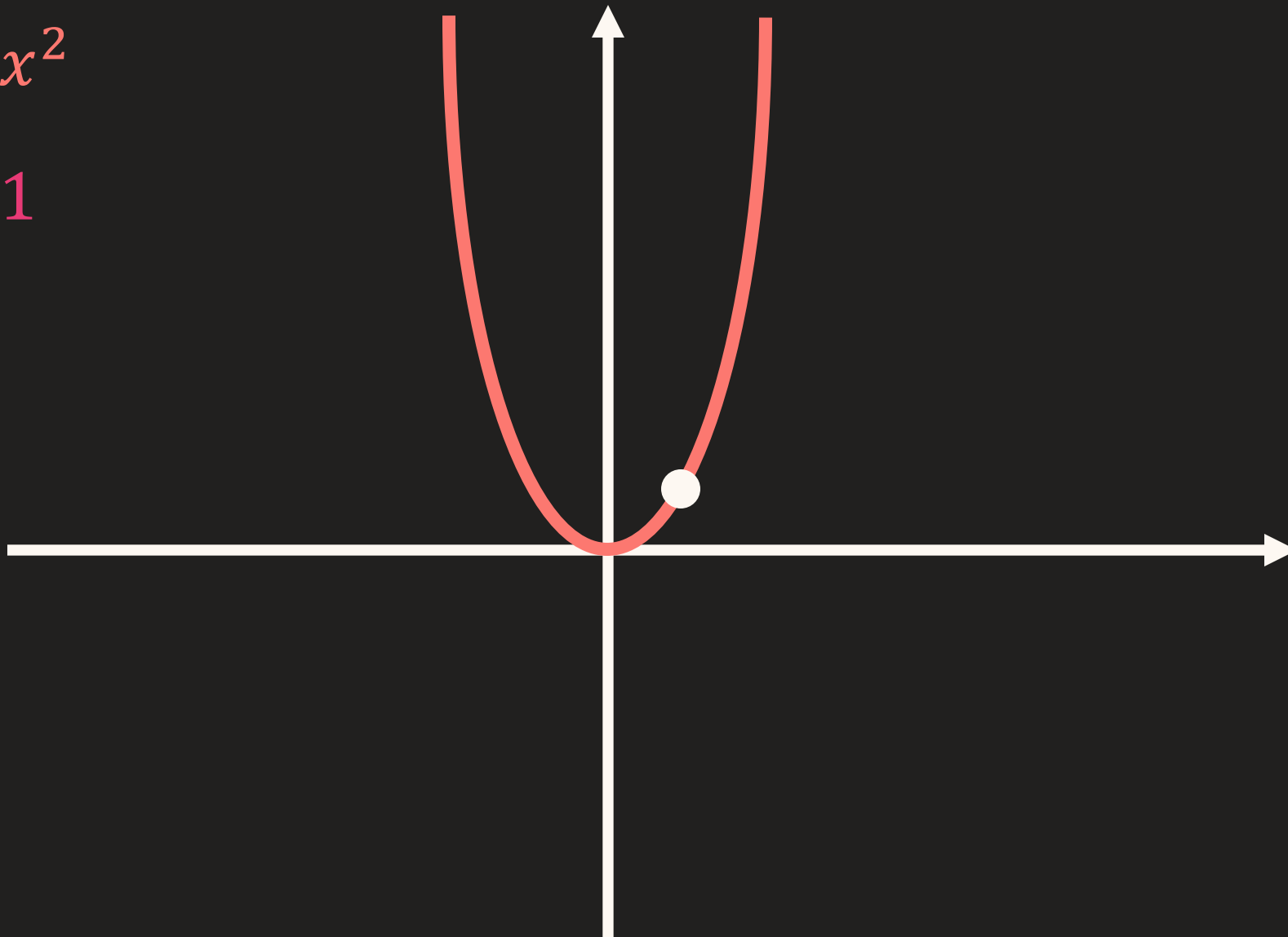


$$f(x) = x^2$$



$$f(x) = x^2$$

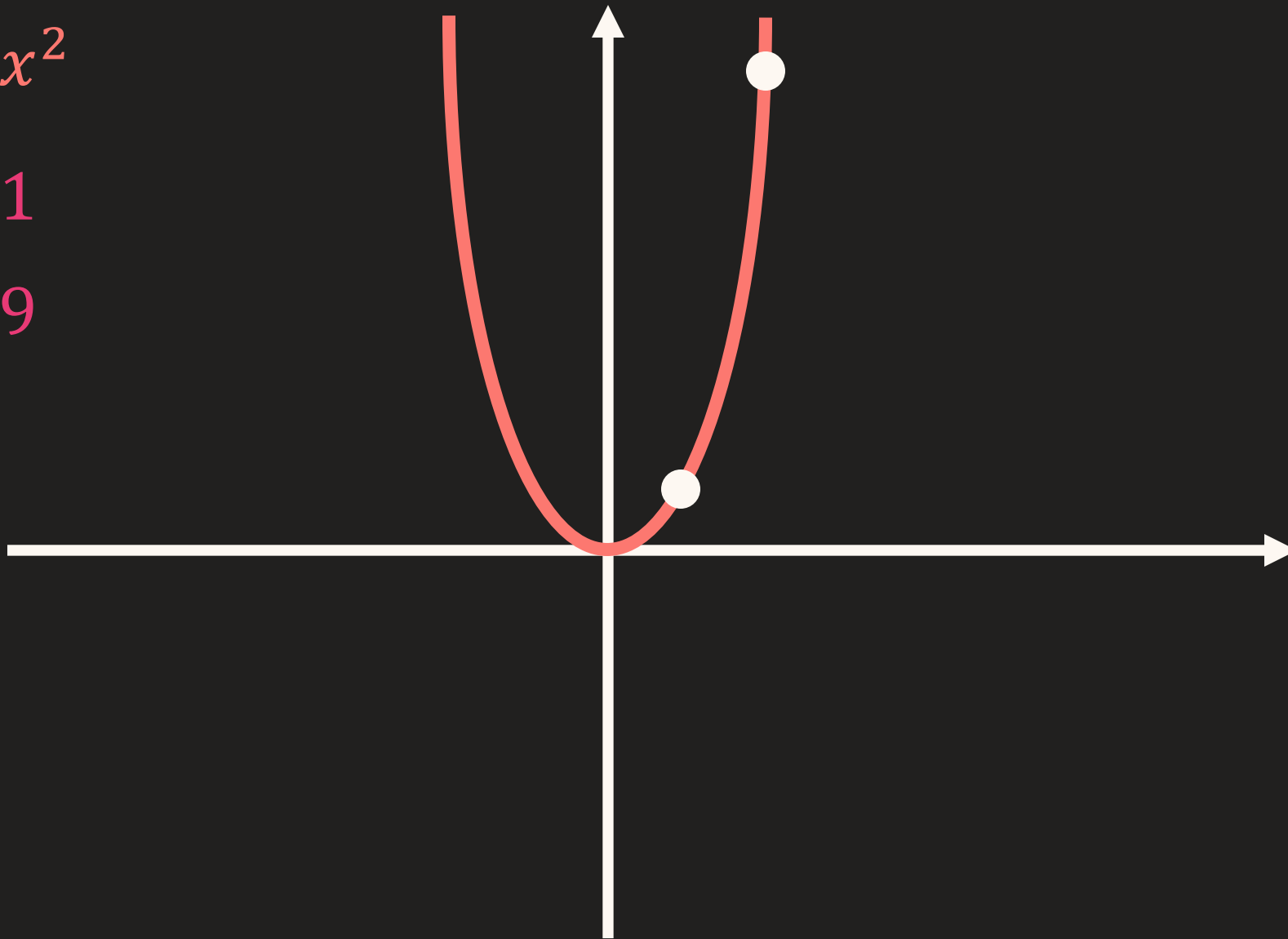
$$f(1) = 1$$



$$f(x) = x^2$$

$$f(1) = 1$$

$$f(3) = 9$$



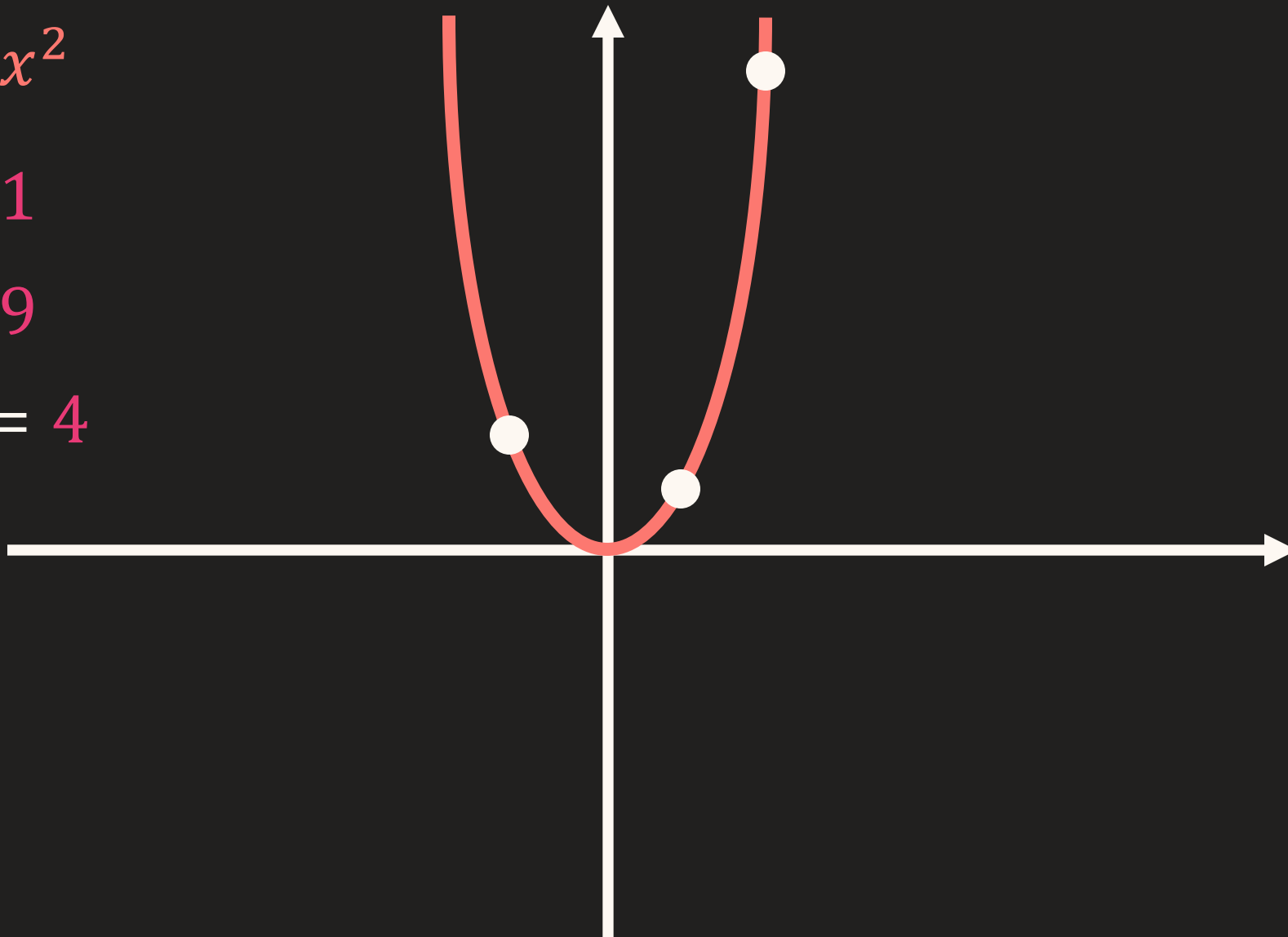


$$f(x) = x^2$$

$$f(1) = 1$$

$$f(3) = 9$$

$$f(-2) = 4$$

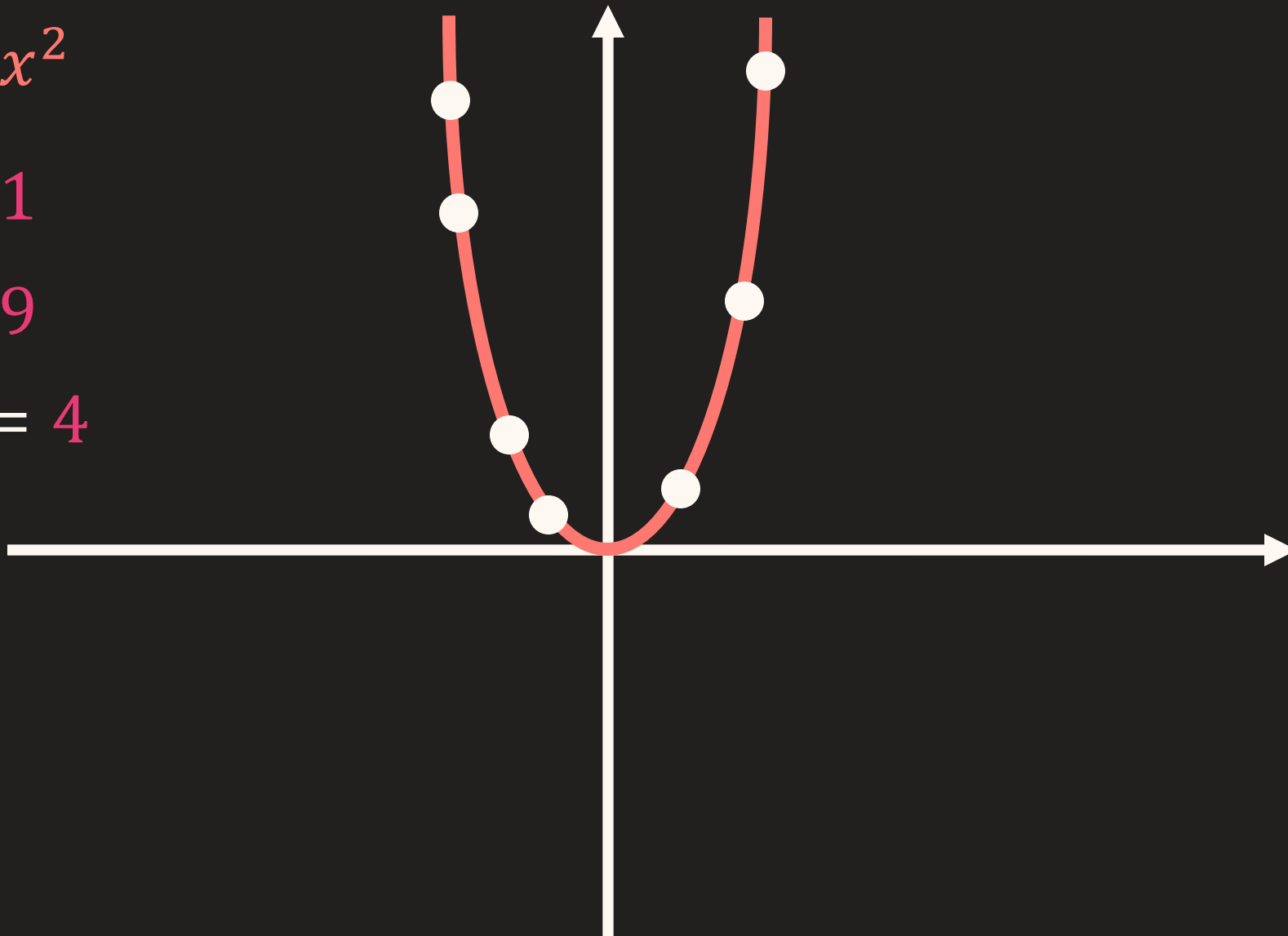


$$f(x) = x^2$$

$$f(1) = 1$$

$$f(3) = 9$$

$$f(-2) = 4$$

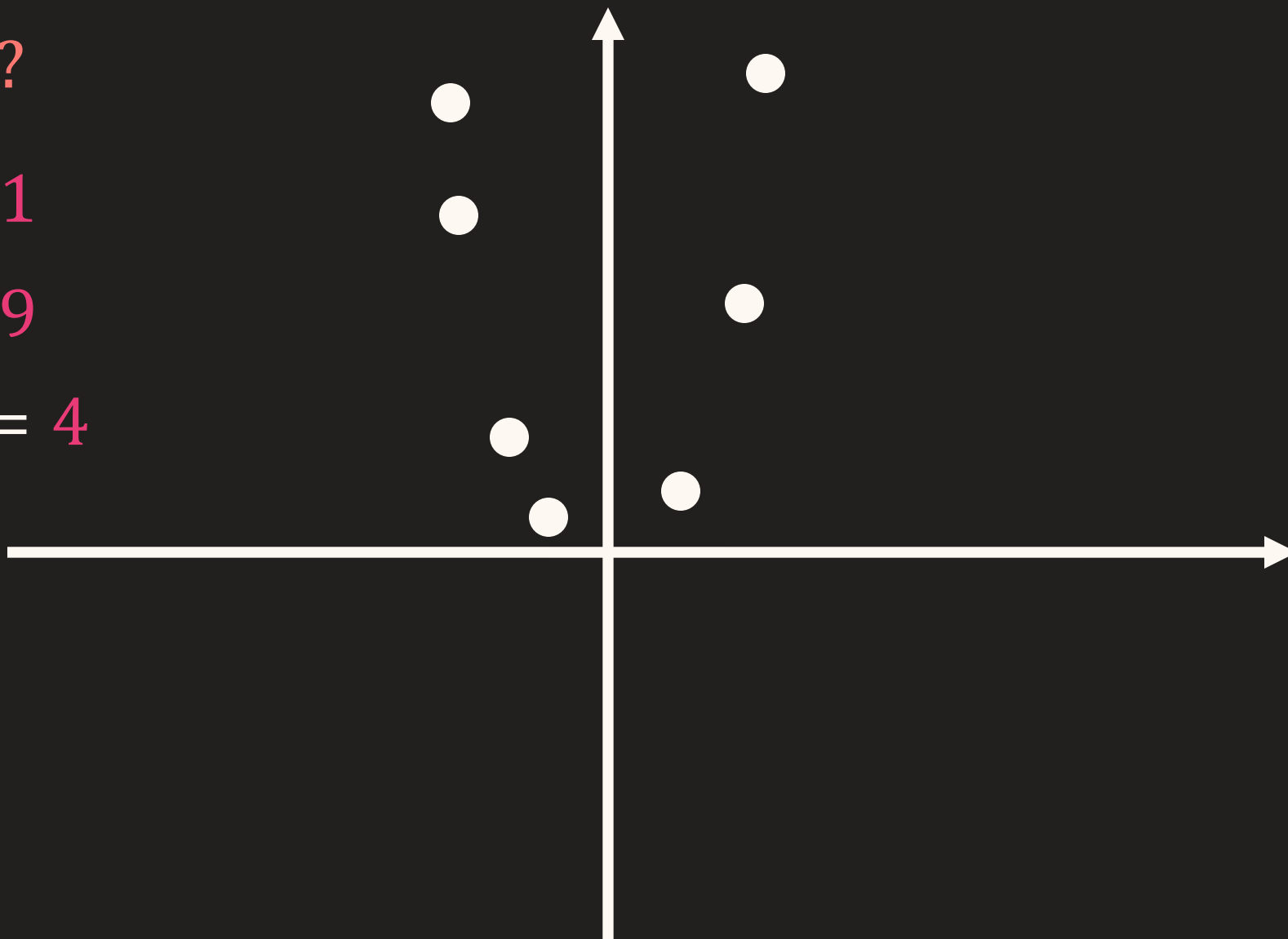


$$f(x) = ?$$

$$f(1) = 1$$

$$f(3) = 9$$

$$f(-2) = 4$$

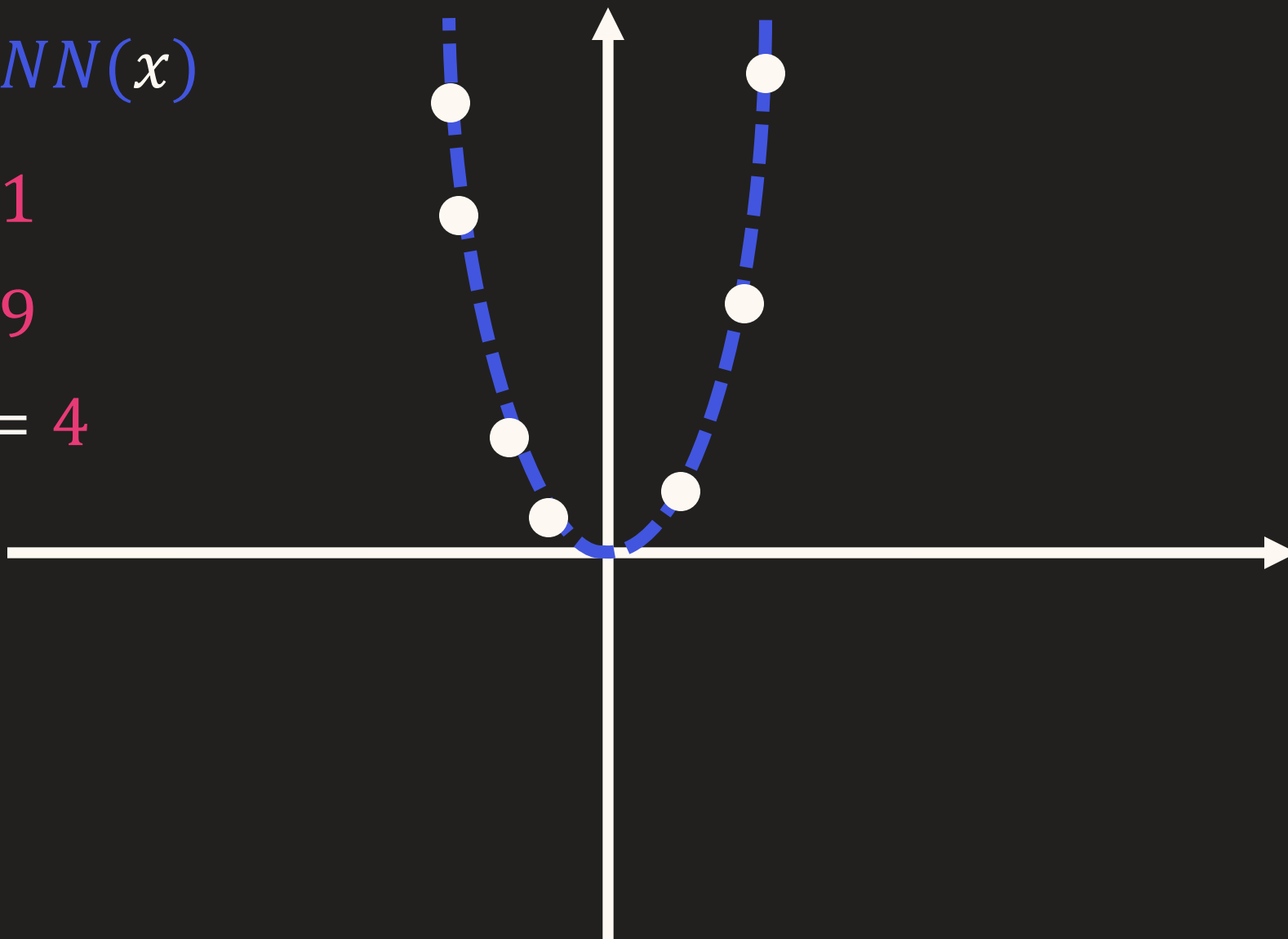


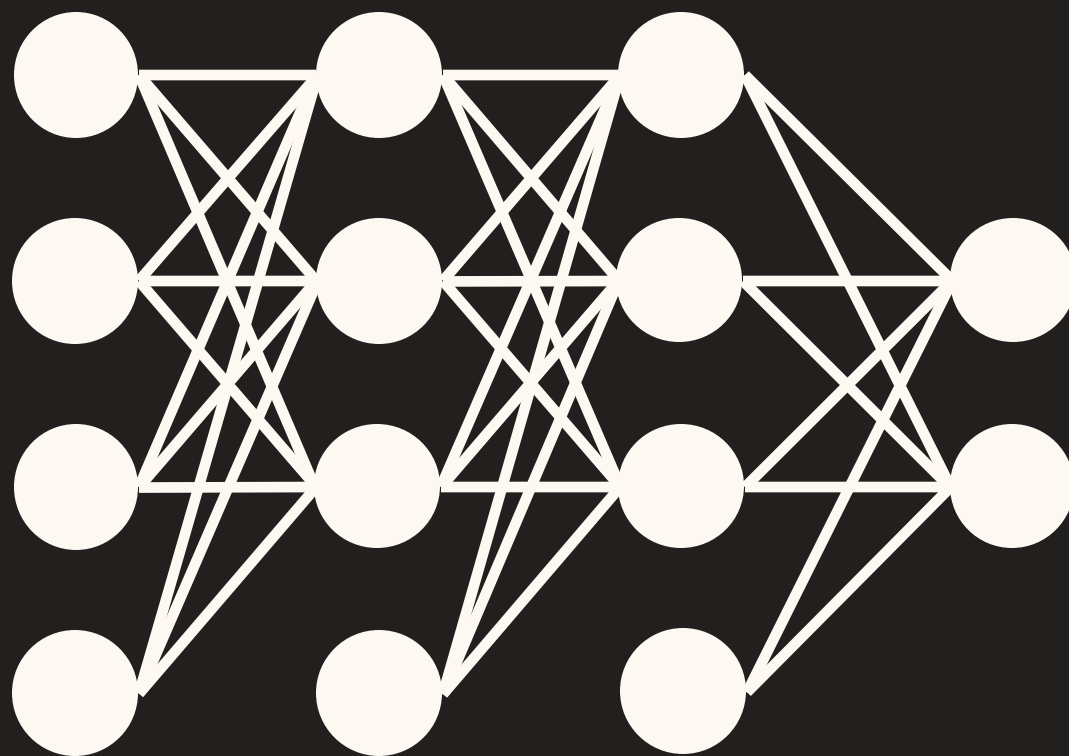
$$f(x) \approx NN(x)$$

$$f(1) = 1$$

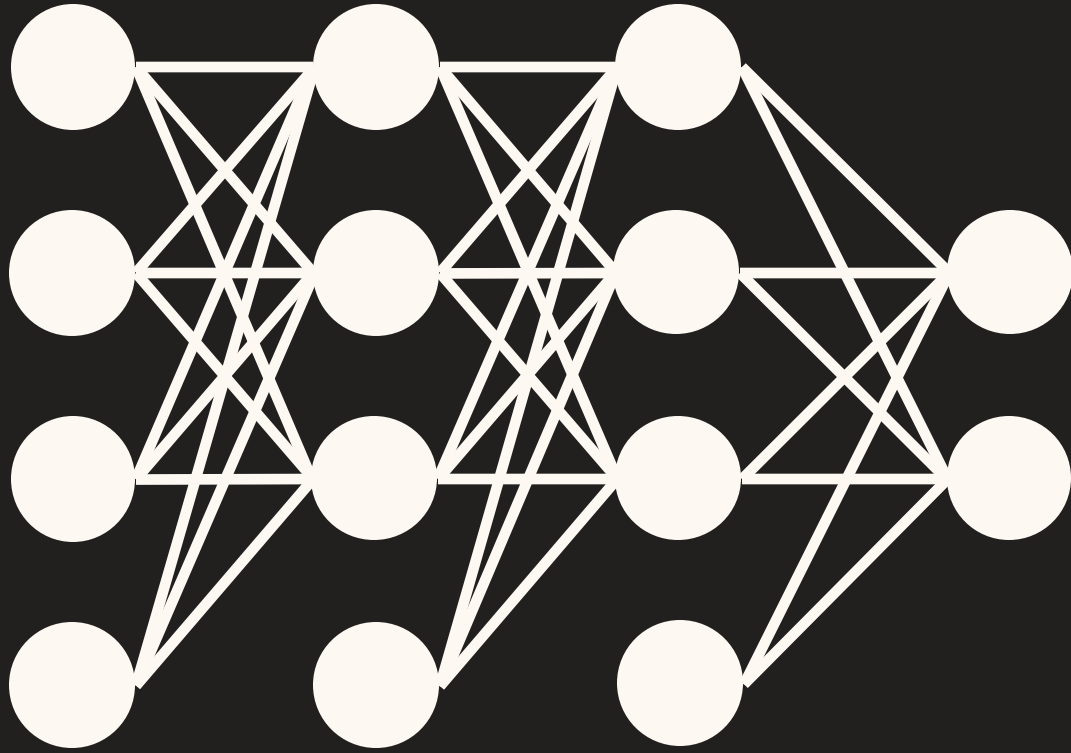
$$f(3) = 9$$

$$f(-2) = 4$$





MLP



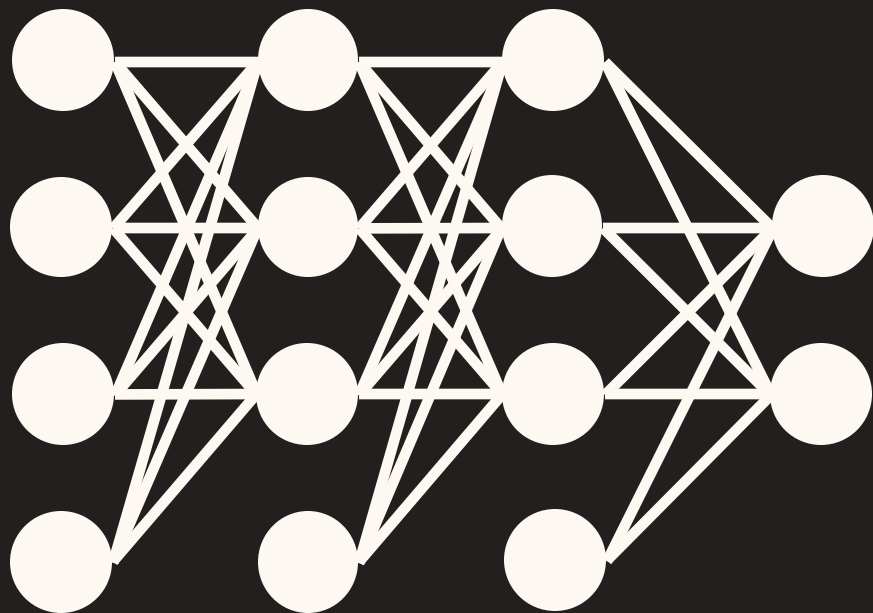
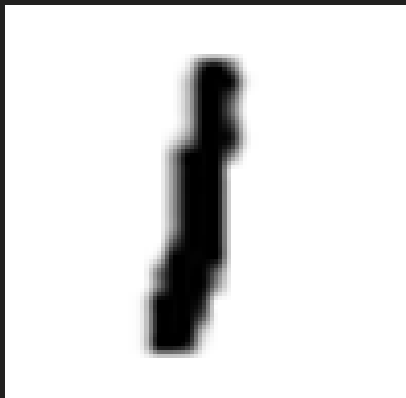
MLP

**50s**

⋮

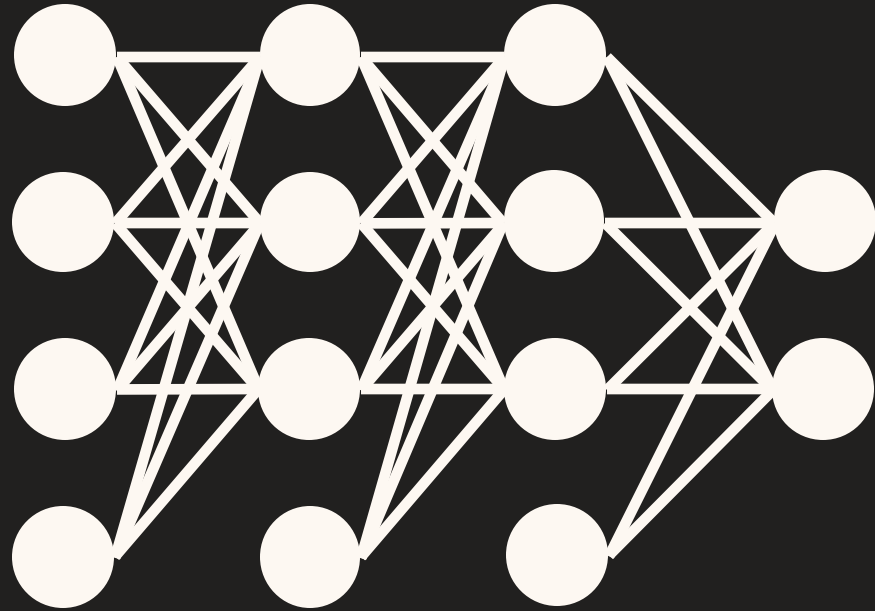
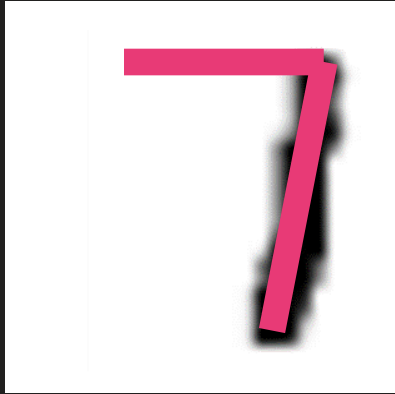
**90s**

Problemy?



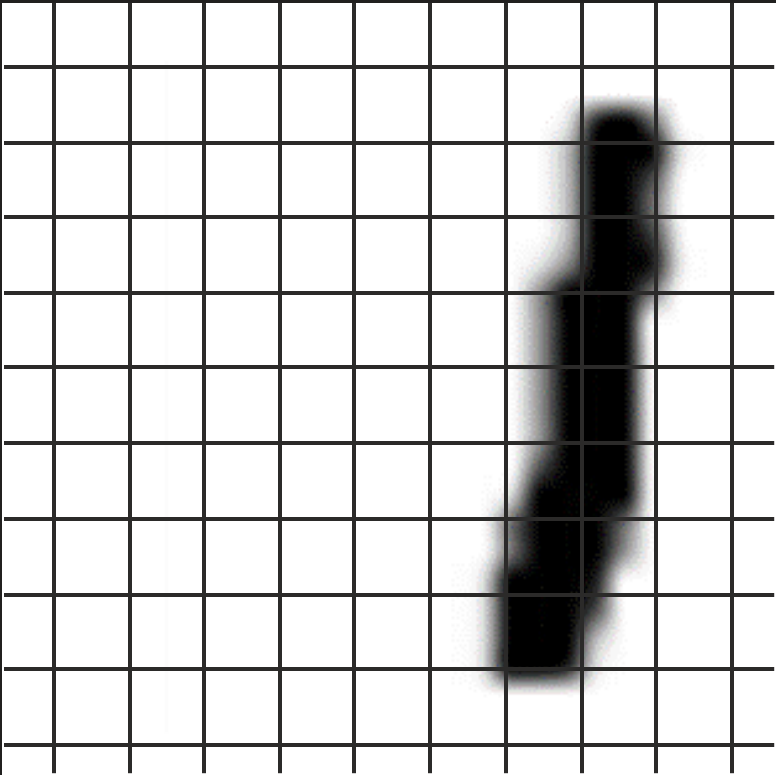
1

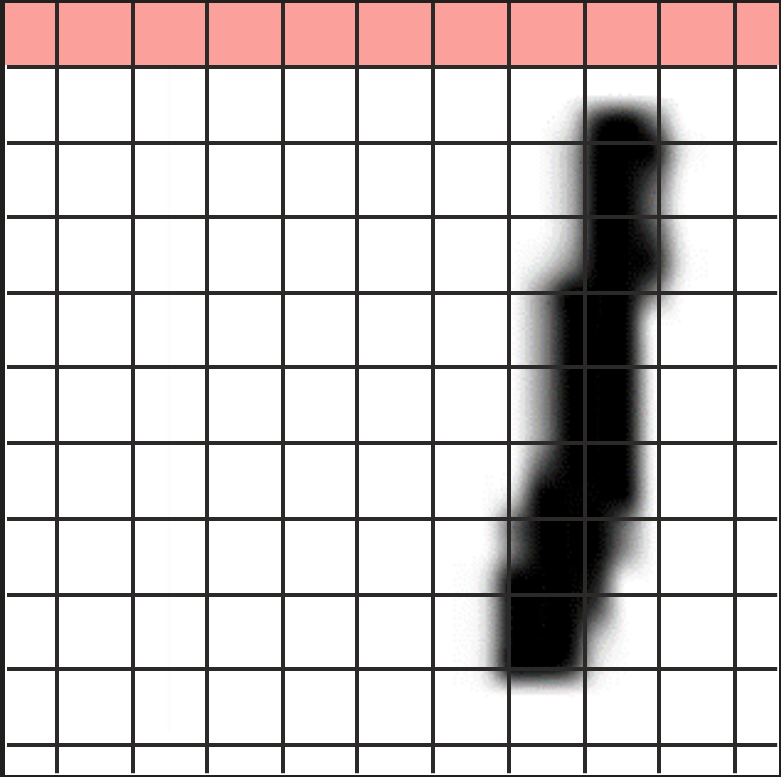


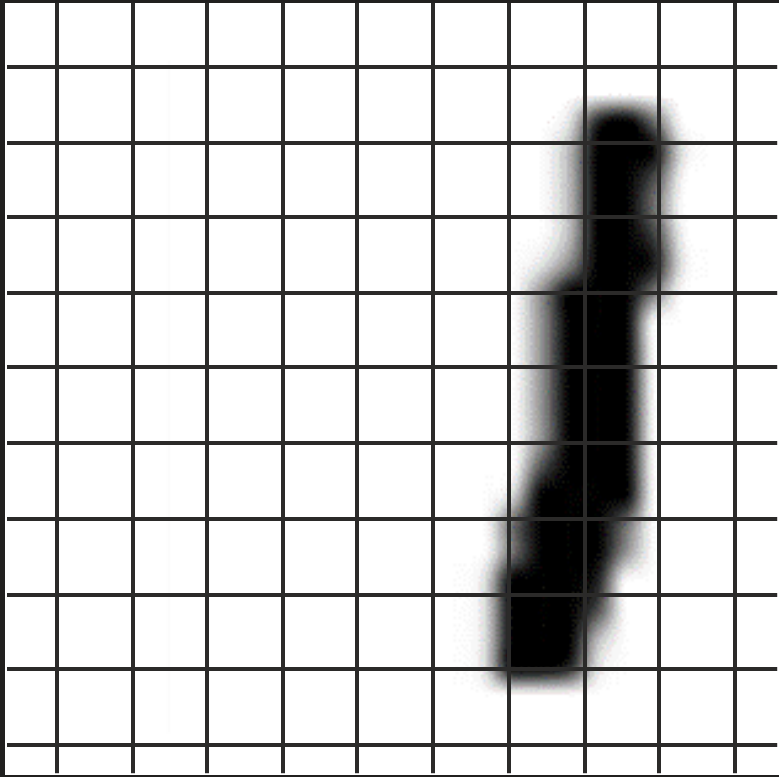


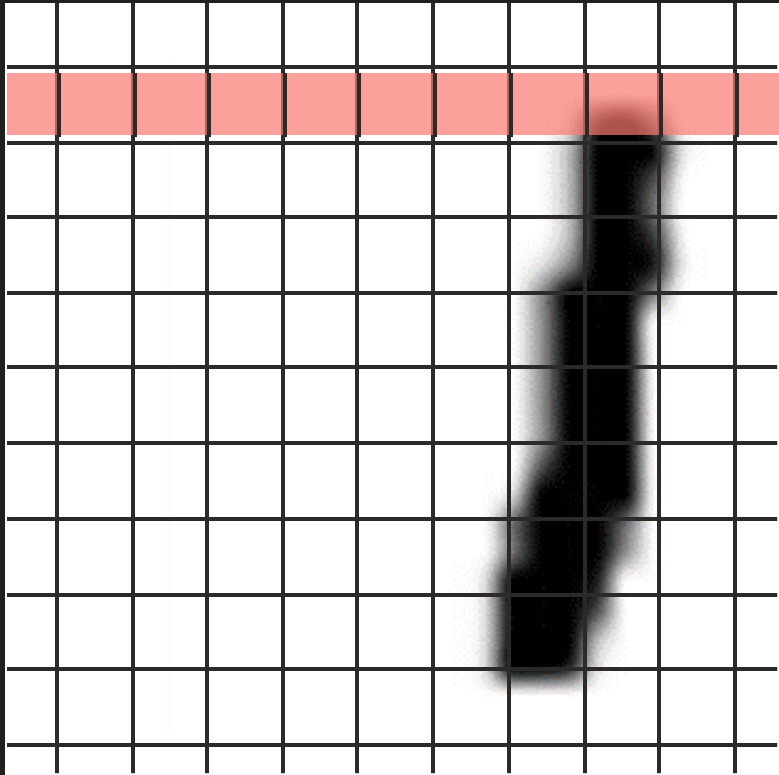
7

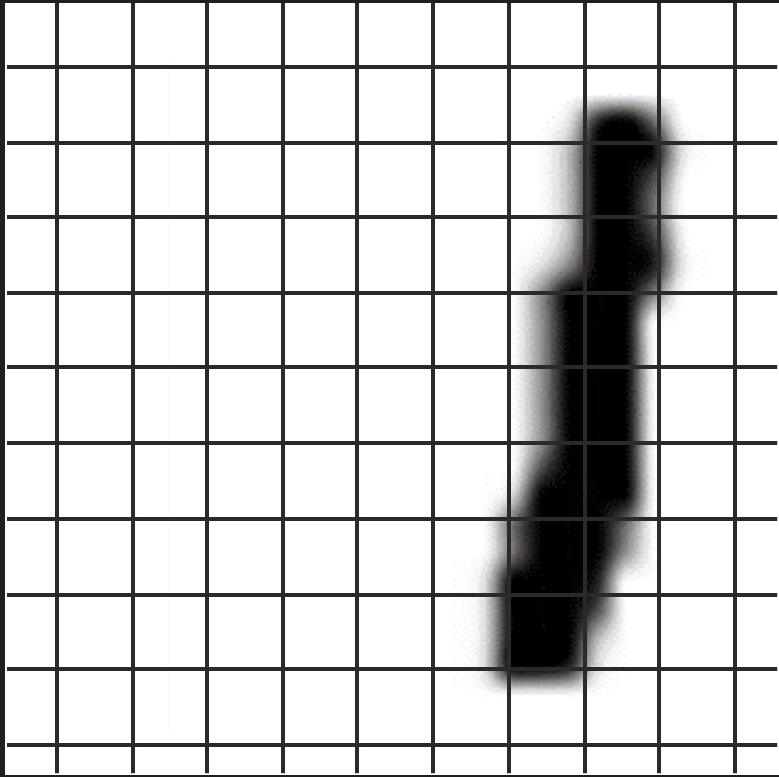


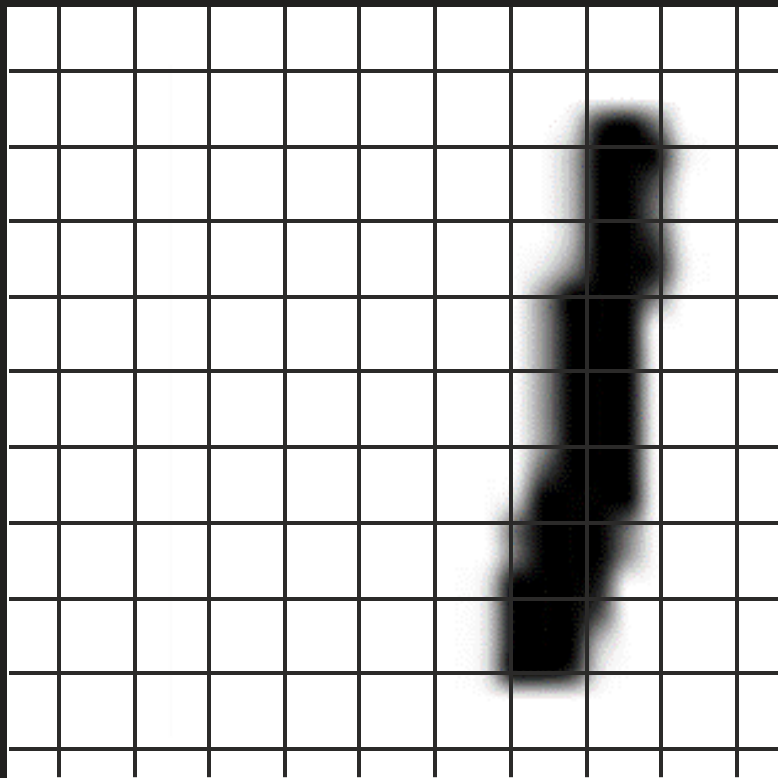




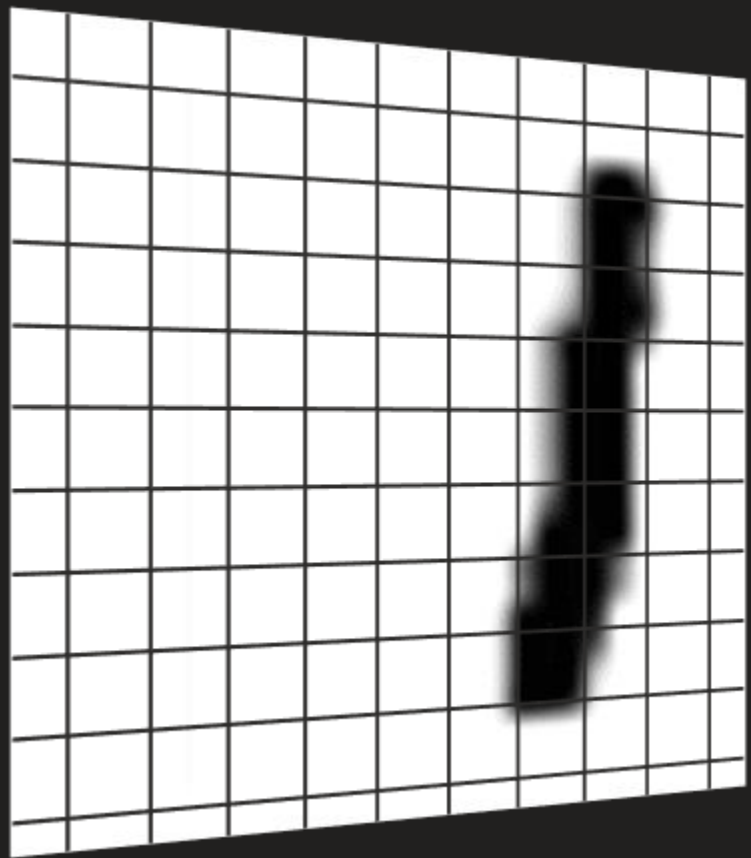


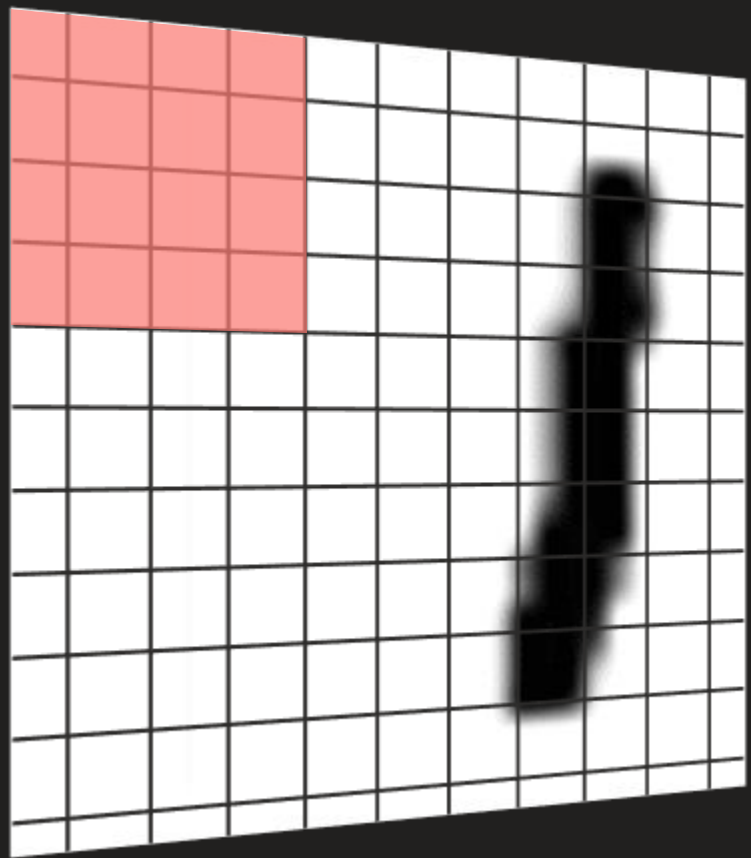


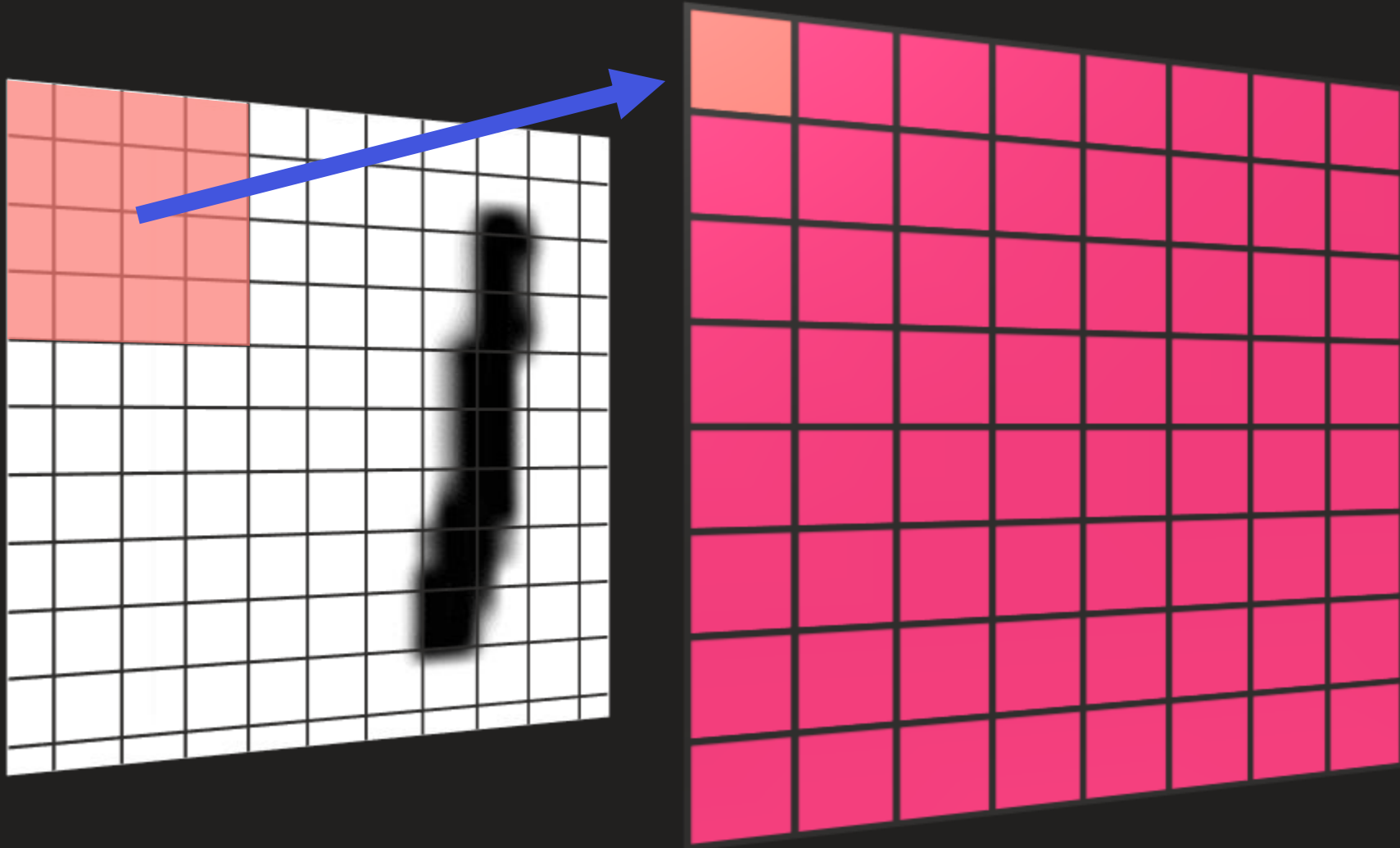


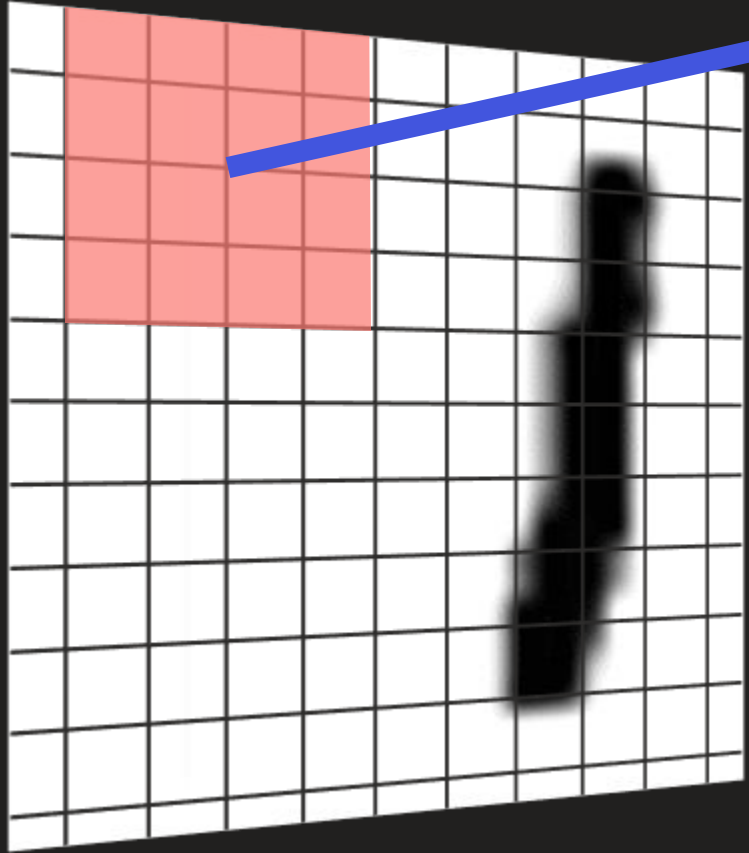

$$\begin{bmatrix} 132 \\ 12 \\ 21 \\ 231 \\ \vdots \\ 12 \\ 0 \\ 255 \\ 52 \end{bmatrix}$$



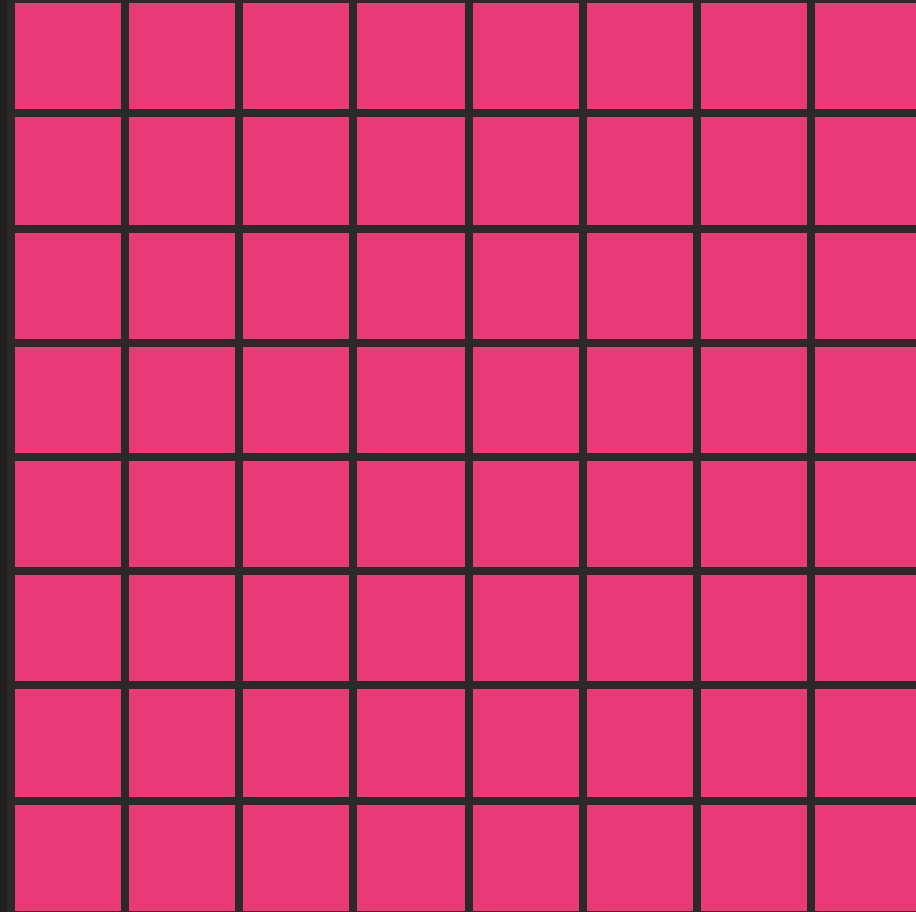


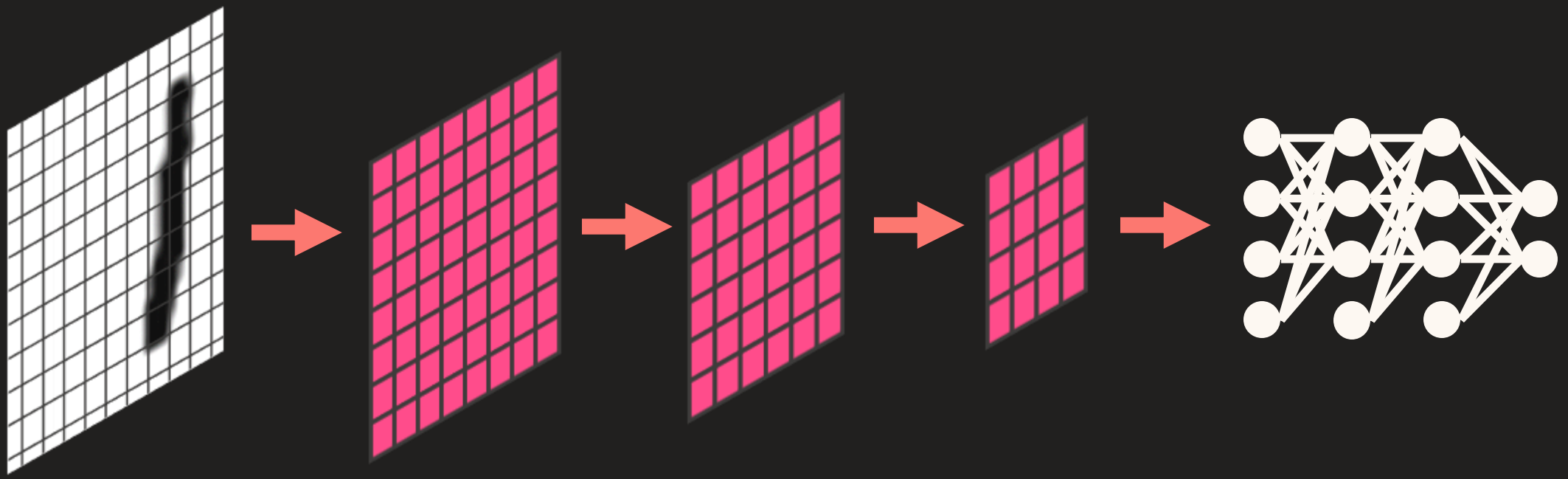




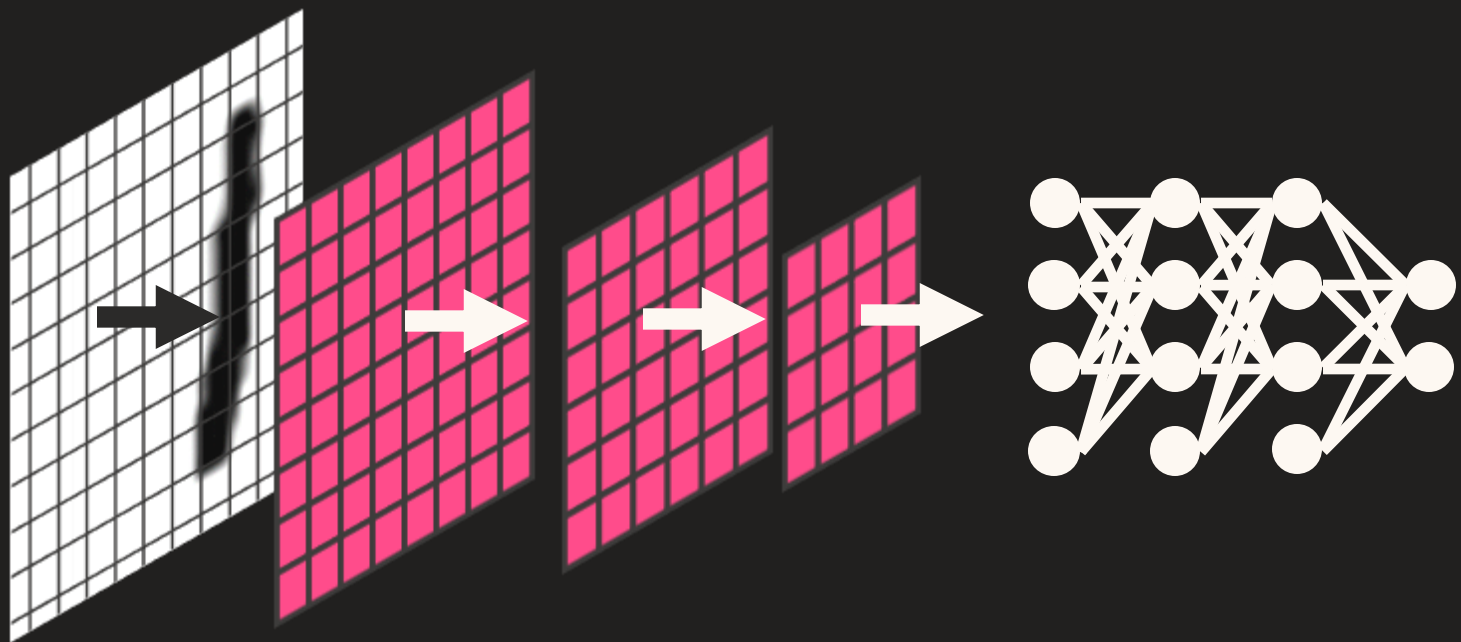








CNN



CNN

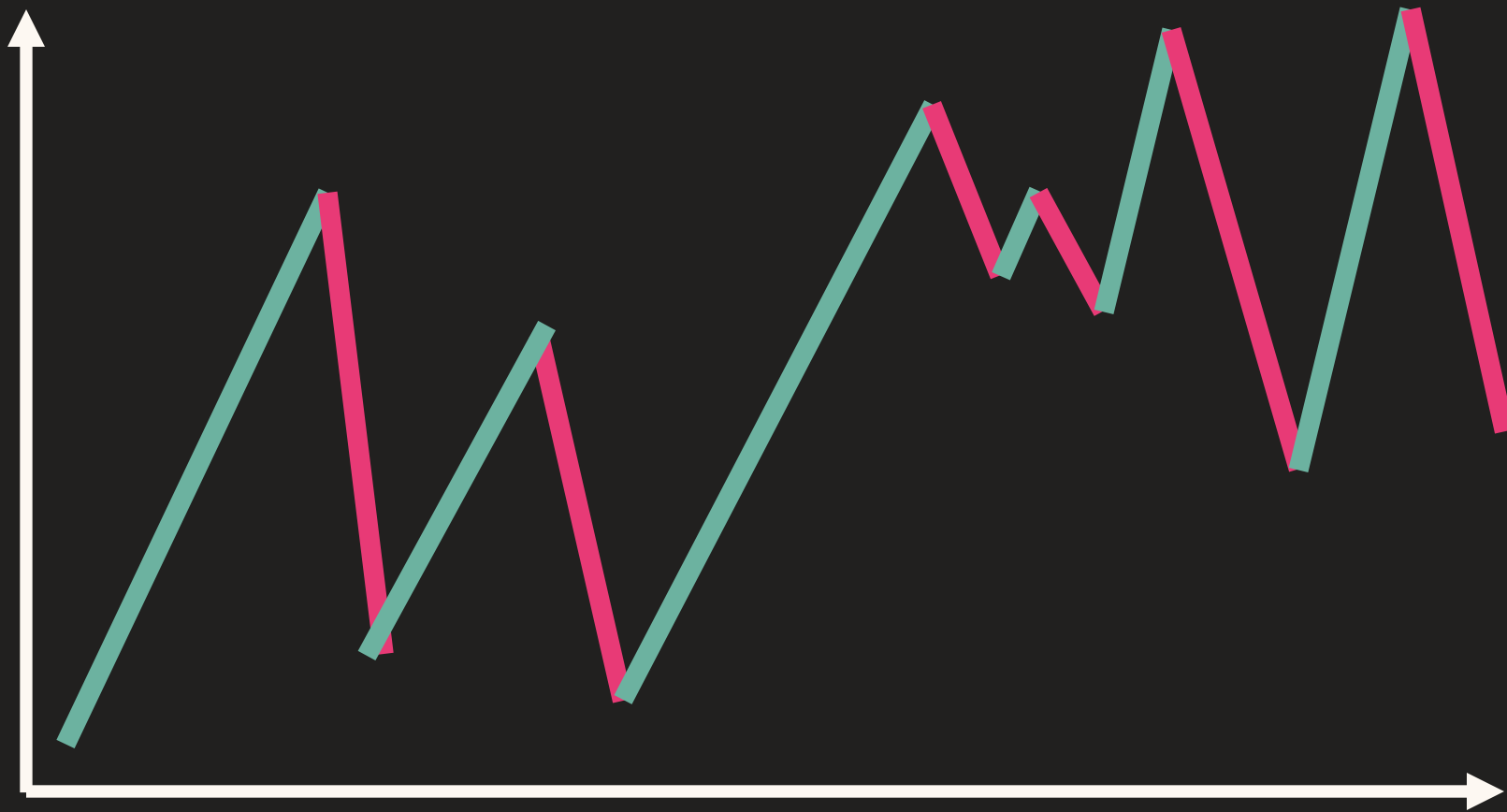
80s

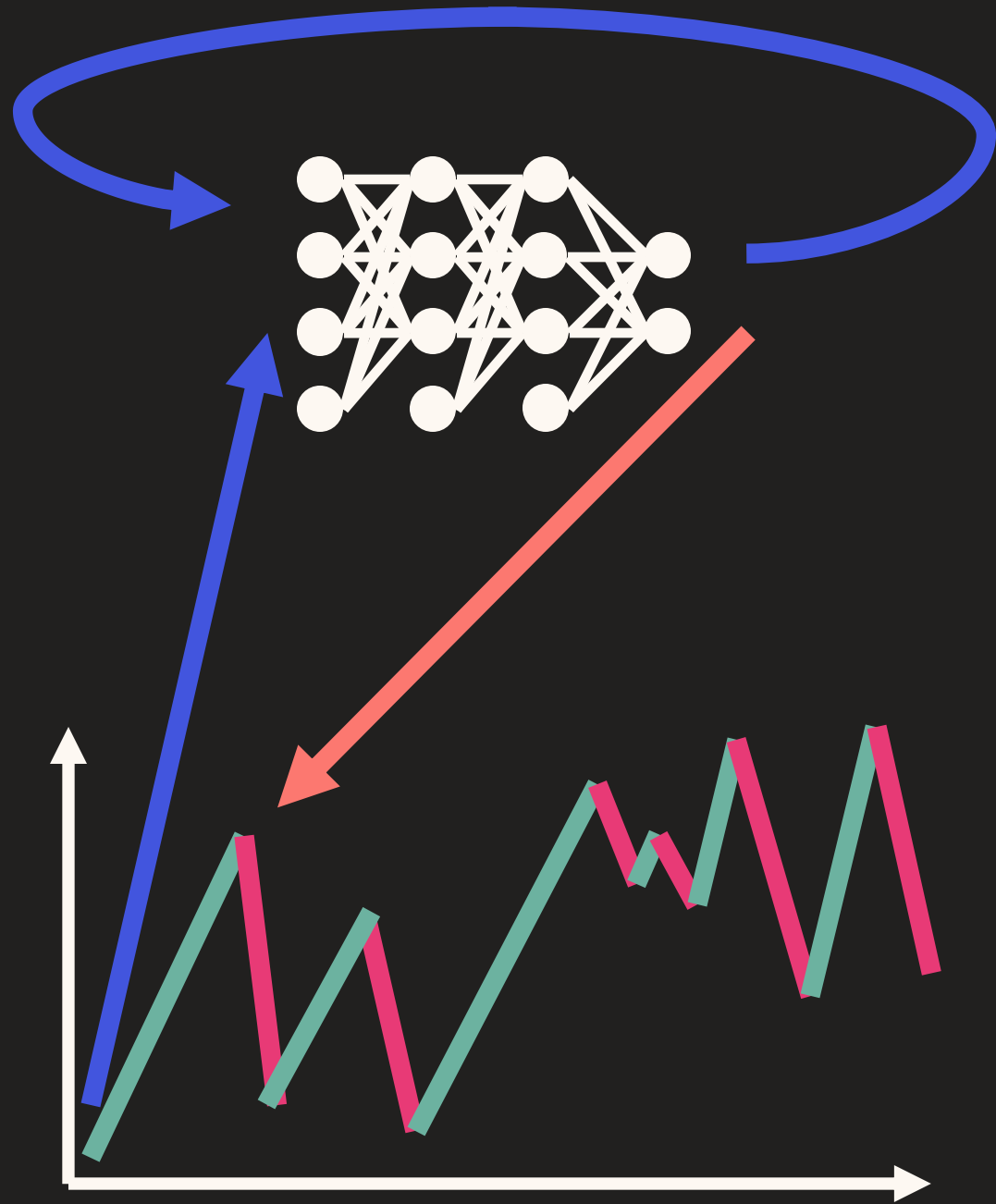
⋮

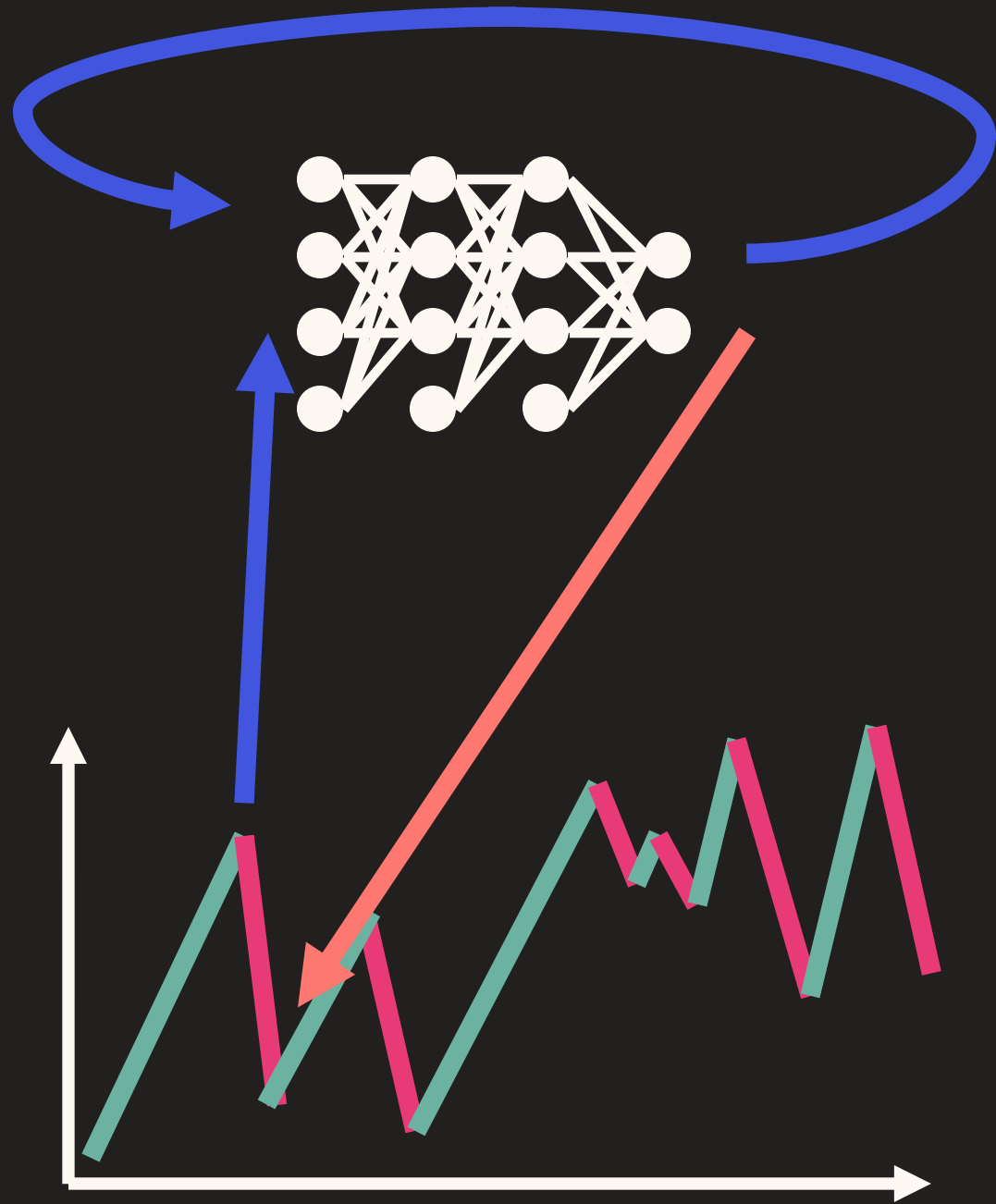
90s

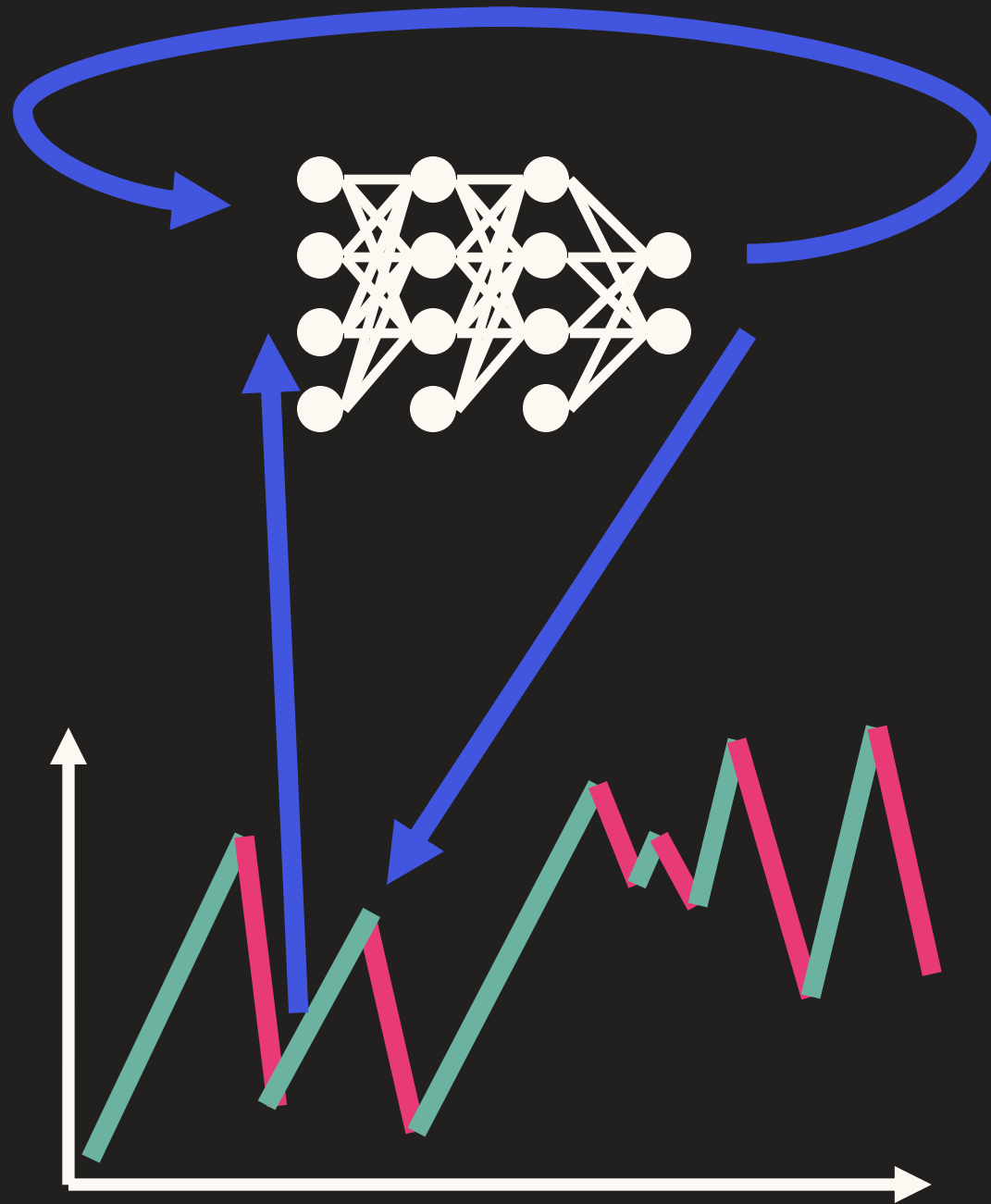


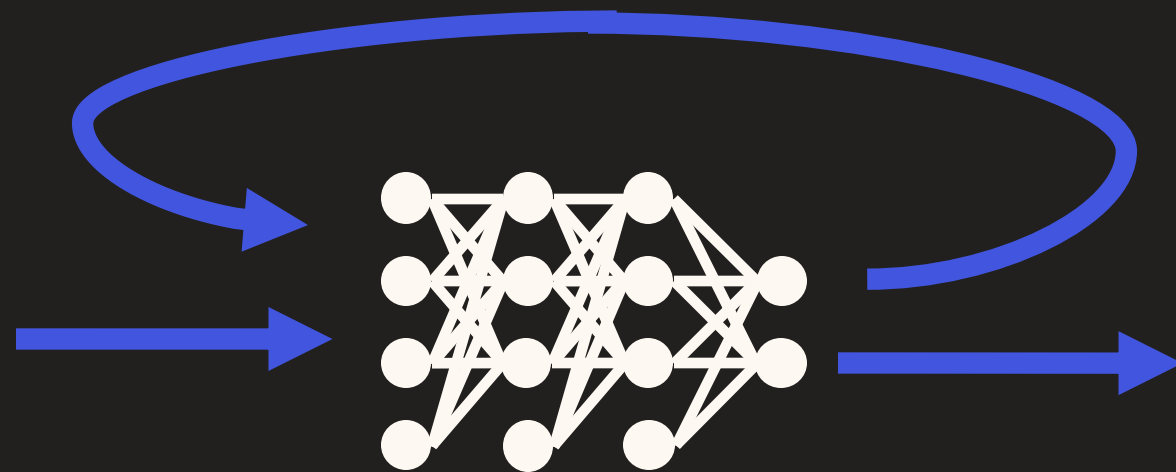
Dane czasowe i szeregowy



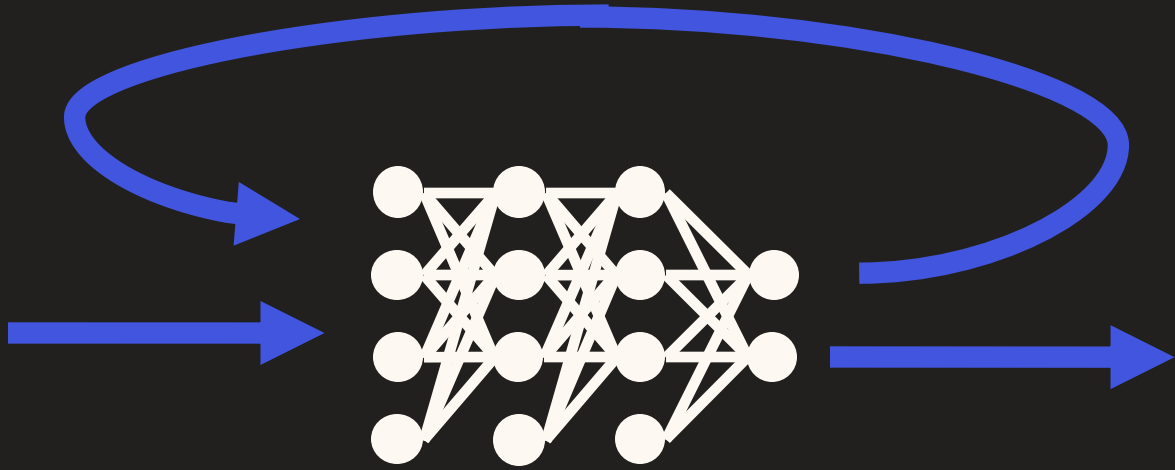








RNN



RNN

80s

⋮

90s

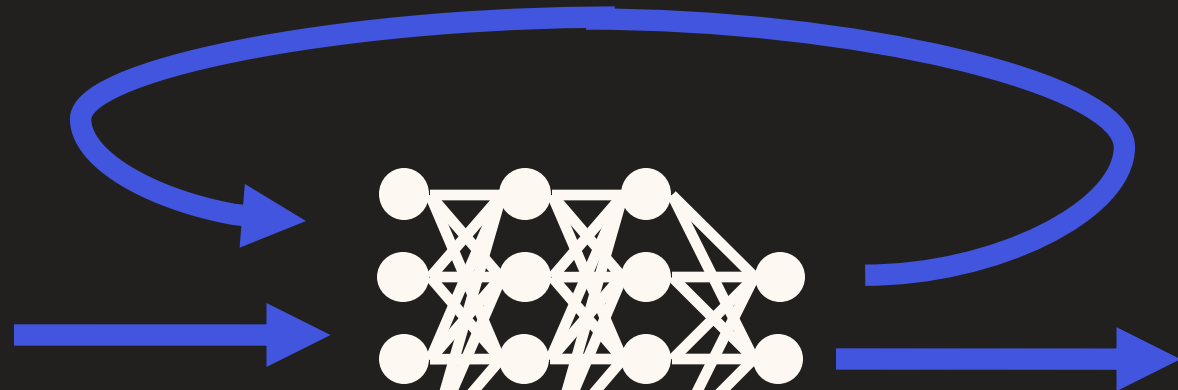
Problemy?



$$1.1^{100}$$

13780.612...

$x100$



1.1

$$0.9^{100}$$

0.00002656...

# AKT II

*Jak rozmawiać z maszyną?*

AI

Stolicą Francji jest Paryż.

Jakie miasto jest stolicą Francji?



Być, albo nie być?  
Oto jest pytanie.



Być, albo nie być? Oto

Być, albo nie być? Oto

Być, albo nie być? Oto

Być, albo nie być? Oto jest

Być, albo nie być? Oto jest

Być, albo nie być? Oto jest pytanie



Jakie miasto jest stolicą  
Francji?

Pytam, ponieważ mam  
jutro sprawdzian z  
Geografii



Pytanie:  
Jakie miasto jest stolicą  
Francji?  
Odpowiedź:

Stolicą

Pytanie:  
Jakie miasto jest stolicą  
Francji?  
Odpowiedź:

Stolicą Francji

Pytanie:

Jakie miasto jest stolicą  
Francji?

Odpowiedź:

Stolicą Francji jest

Pytanie:

Jakie miasto jest stolicą  
Francji?

Odpowiedź:

Stolicą Francji jest Paryż

ChatGPT

# GPT

Generative Pretrained Transformer

Generatywny

Przedtrenowany (?)

Transformer

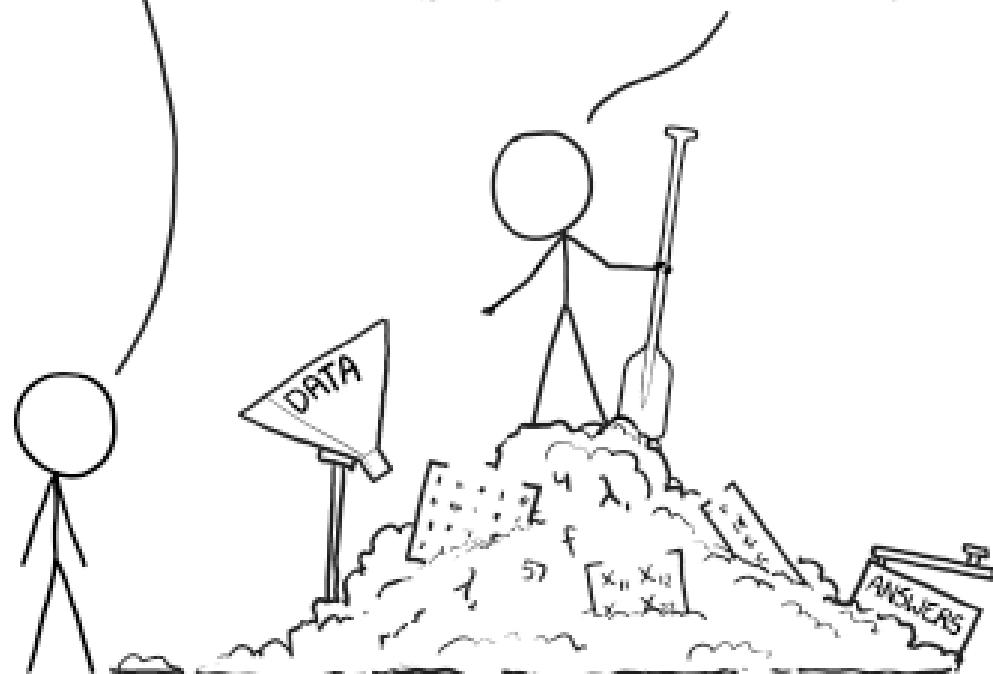
Transformer

THIS IS YOUR MACHINE LEARNING SYSTEM?

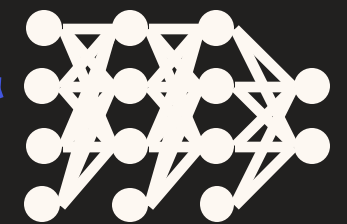
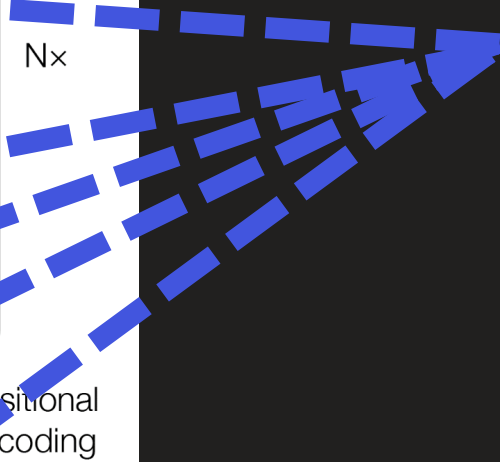
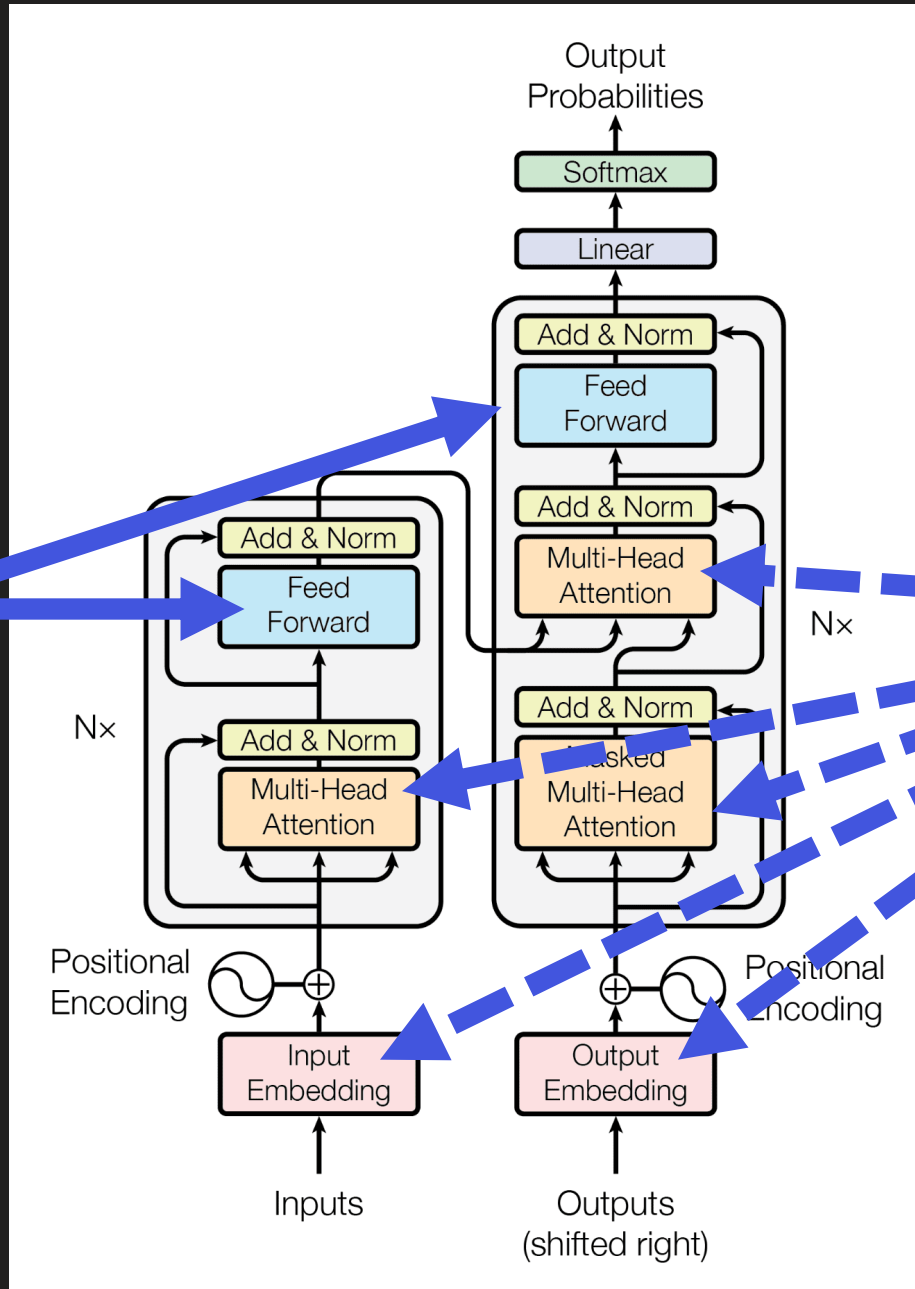
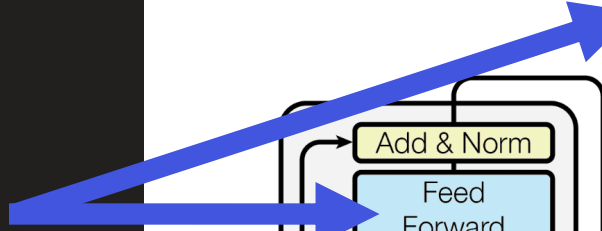
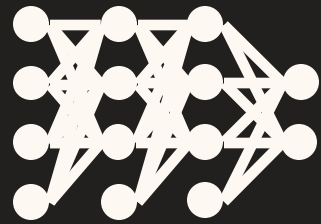
YUP! YOU POUR THE DATA INTO THIS BIG PILE OF LINEAR ALGEBRA, THEN COLLECT THE ANSWERS ON THE OTHER SIDE.

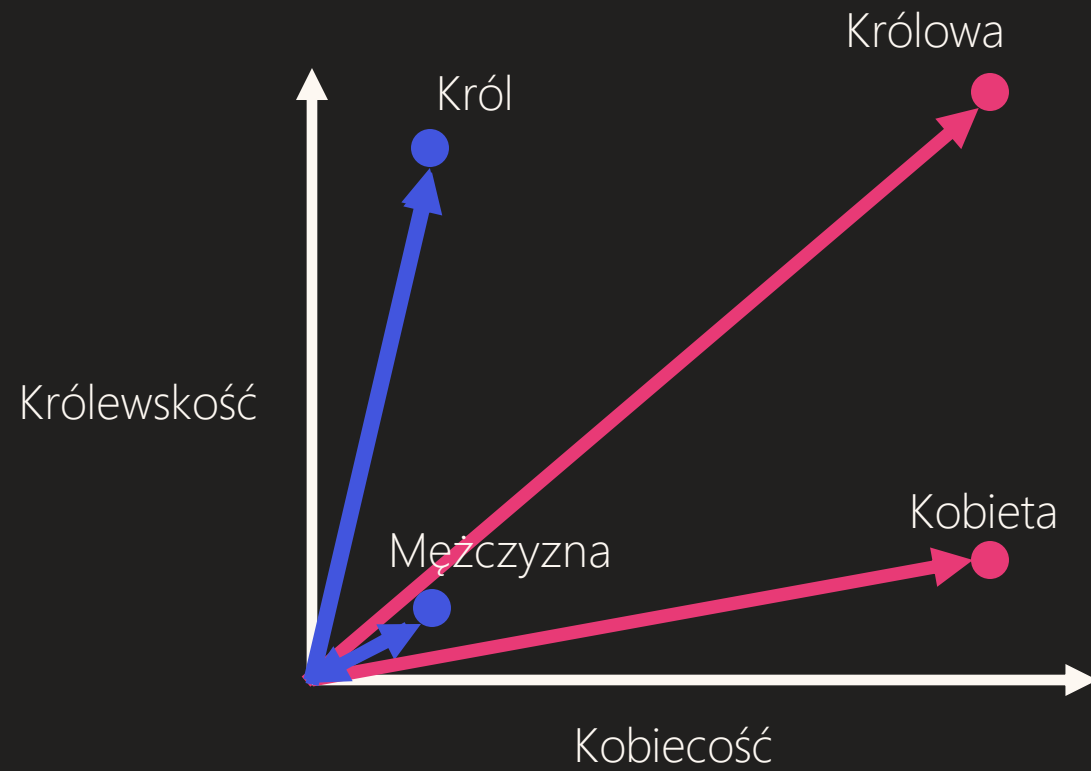
WHAT IF THE ANSWERS ARE WRONG?

JUST STIR THE PILE UNTIL THEY START LOOKING RIGHT.









$$\text{Królowa} = \text{Król} + \text{Kobieta} - \text{Mężczyzna}$$

```
model.distance('car', 'cat')
model.distance('dog', 'cat')
print(model.most_similar_cosmul(positive=['king', 'woman'], negative=['man']))
```

```
[('queen', 0.9314123392105103), ('monarch', 0.858533501625061), ('princess',
```

```
print(model.most_similar_cosmul(positive=['cow', 'oink'], negative=['pig']))
```

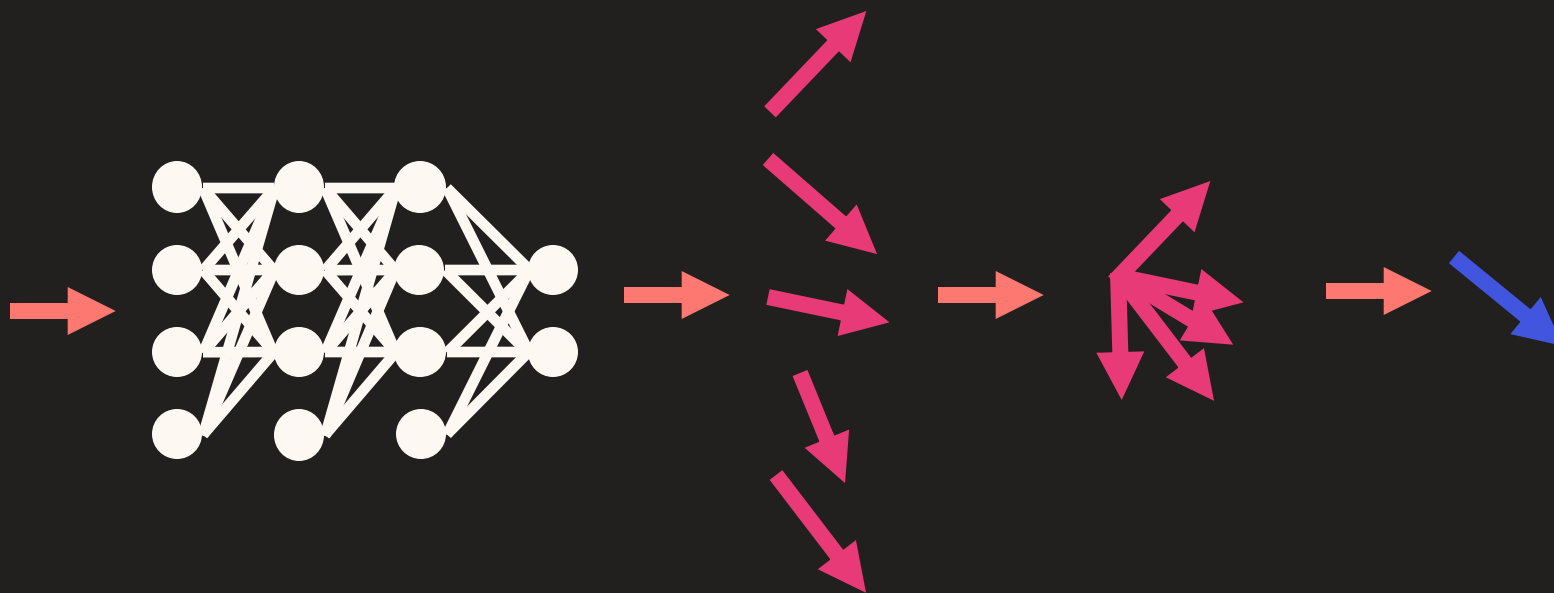
```
[('moos', 0.8037621378898621), ('baa', 0.7783335447311401), ('baaa', 0.7776672
```

```
print(model.most_similar_cosmul(positive=['Santa', 'oink'], negative=['pig']))
```

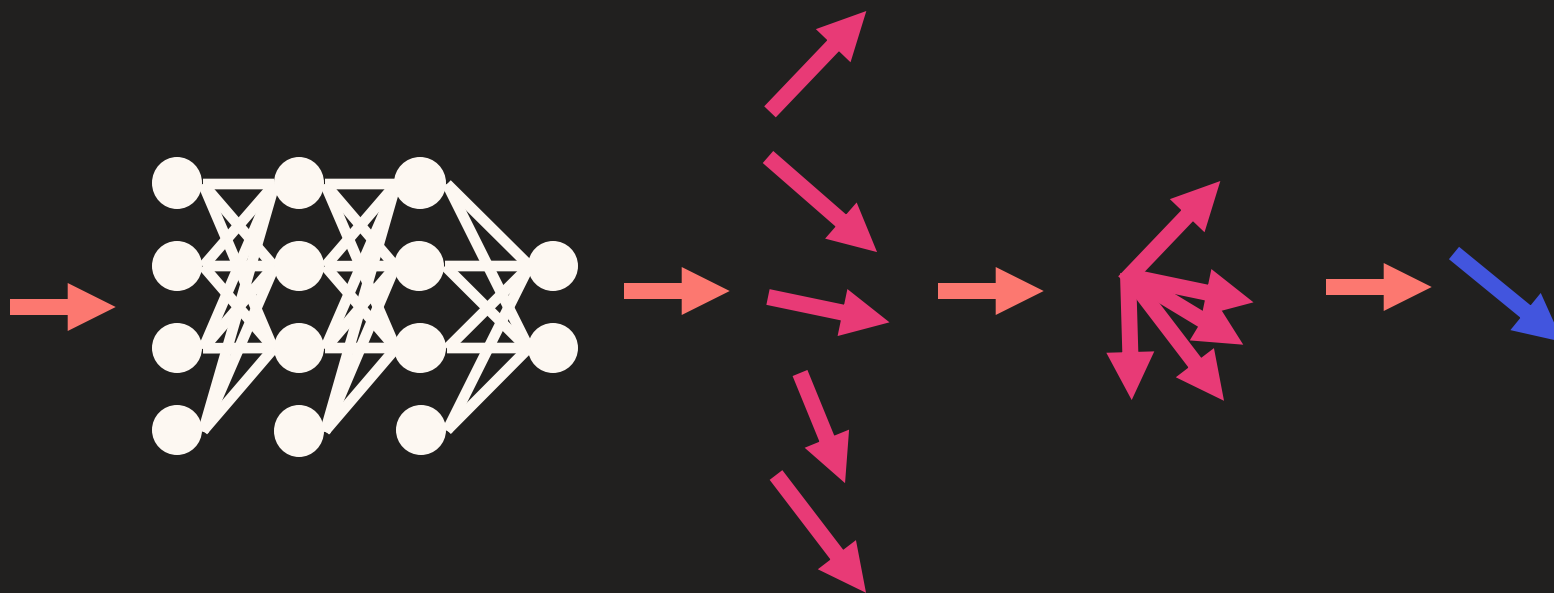
```
[('HO_HO_HO', 0.9089415073394775), ('ho_ho_hoing', 0.9081220030784607), (
```

```
('feline', 0.7326233983039856),
('beagle', 0.7150583267211914),
('puppy', 0.7075453996658325),
('pup', 0.6934291124343872),
('pet', 0.6891531348228455),
('felines', 0.6755931377410889),
('chihuahua', 0.6709762215614319)]
```

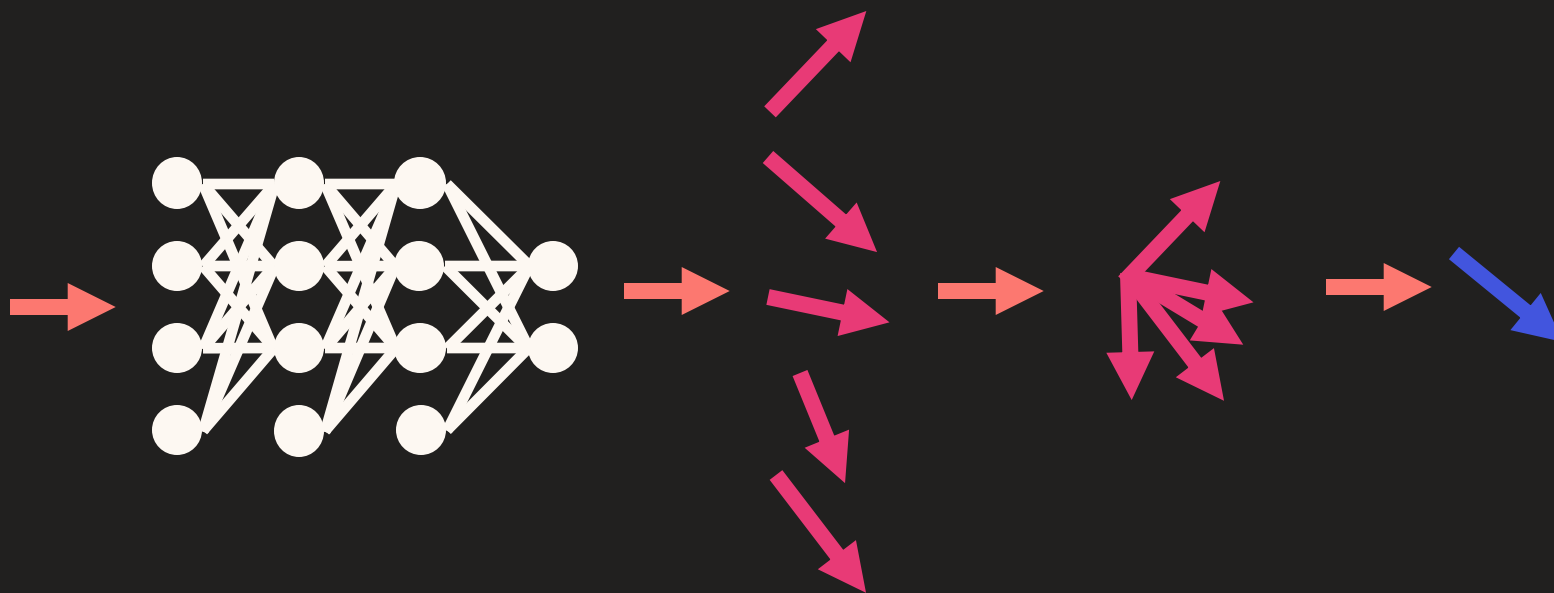
Jakie  
miasto  
jest  
stolicą  
Francji?



Jakie  
miasto  
jest  
stolicą  
Francji?

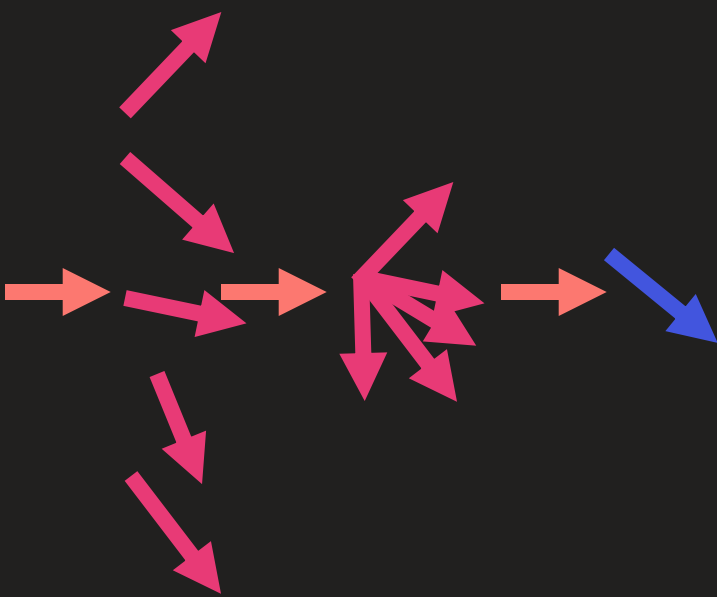
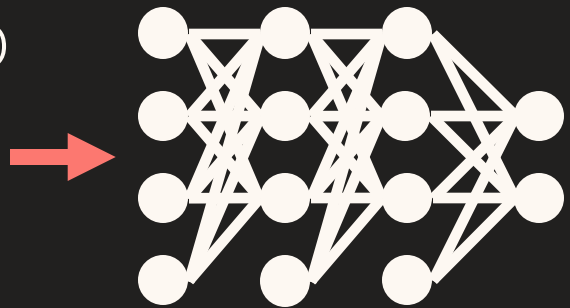


Jakie  
miasto  
jest  
stolicą  
Francji?



Transformer

Jakie  
miasto  
jest  
stolicą  
Francji?



Transformer

2017

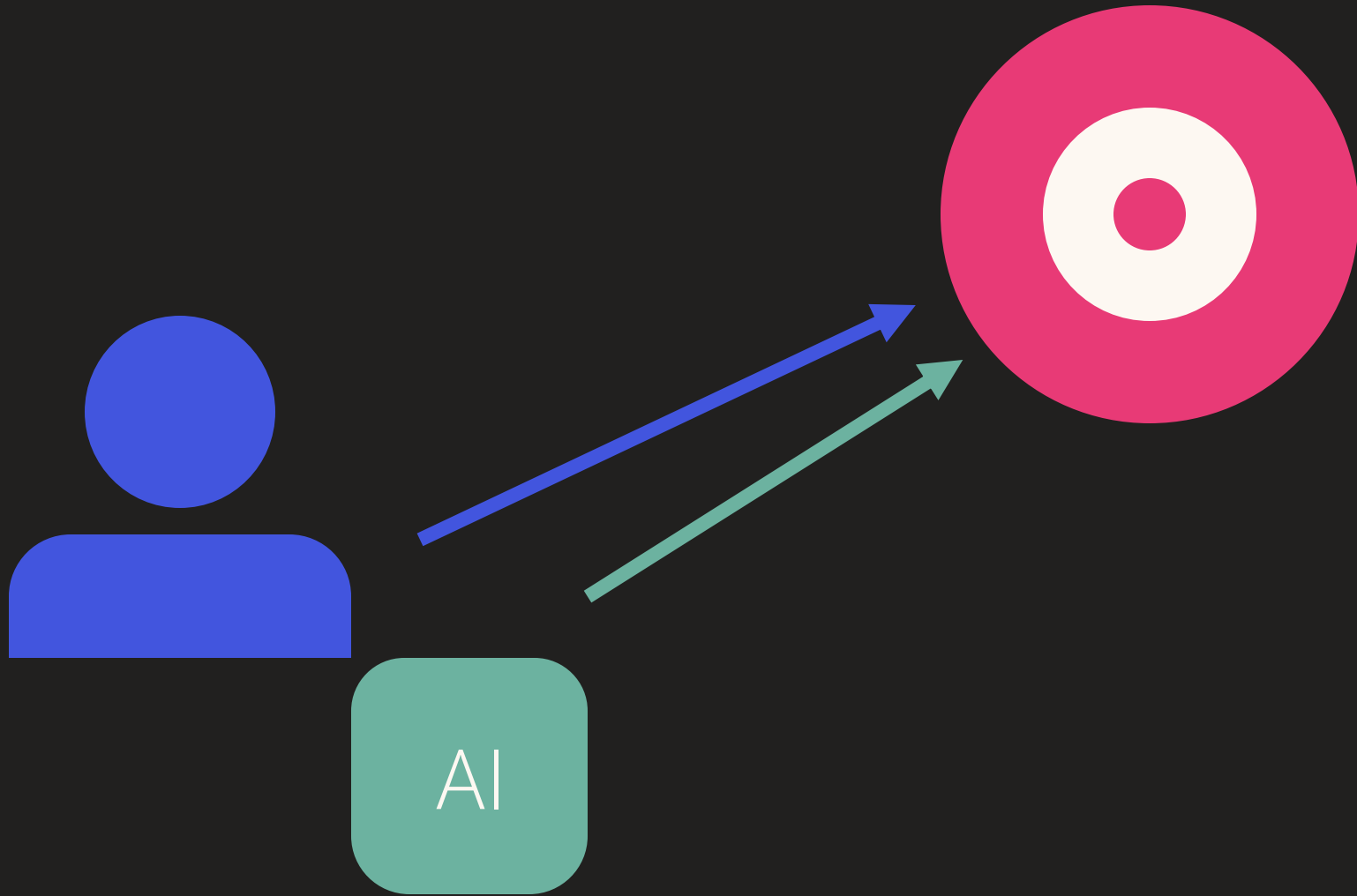




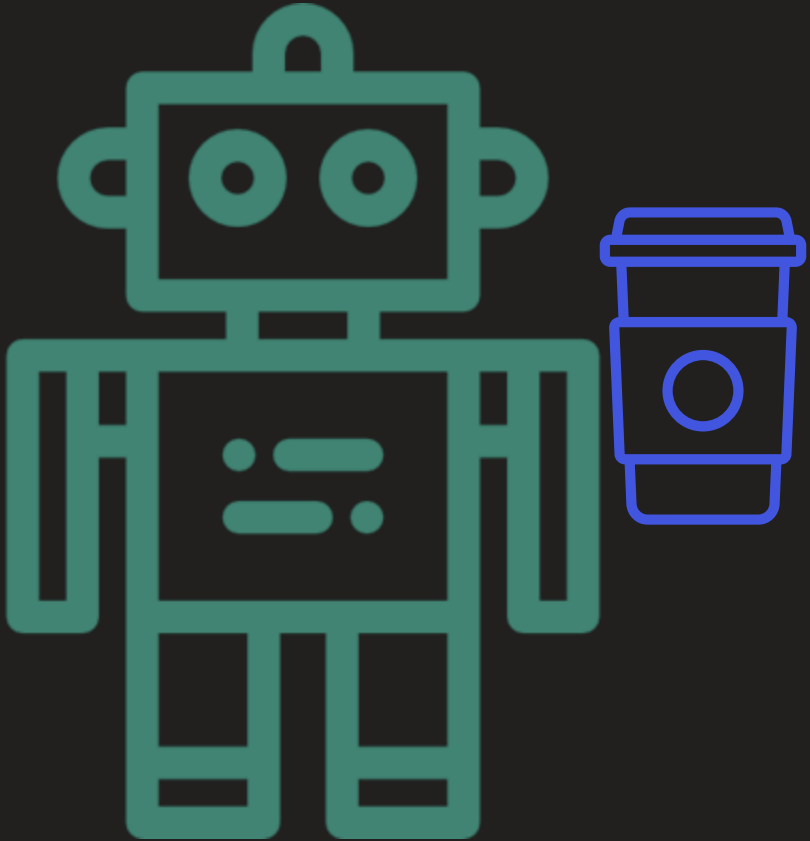
# Alignment Problem

# AKT III

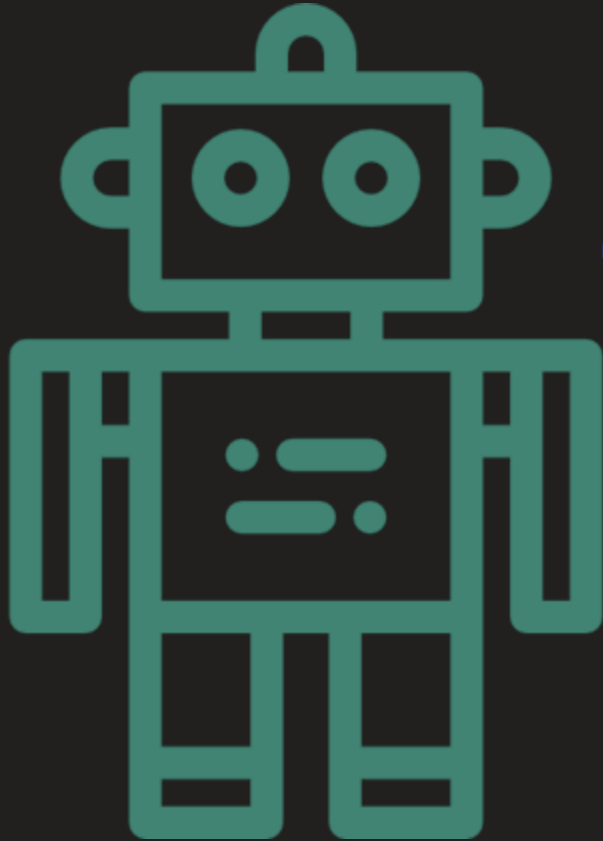
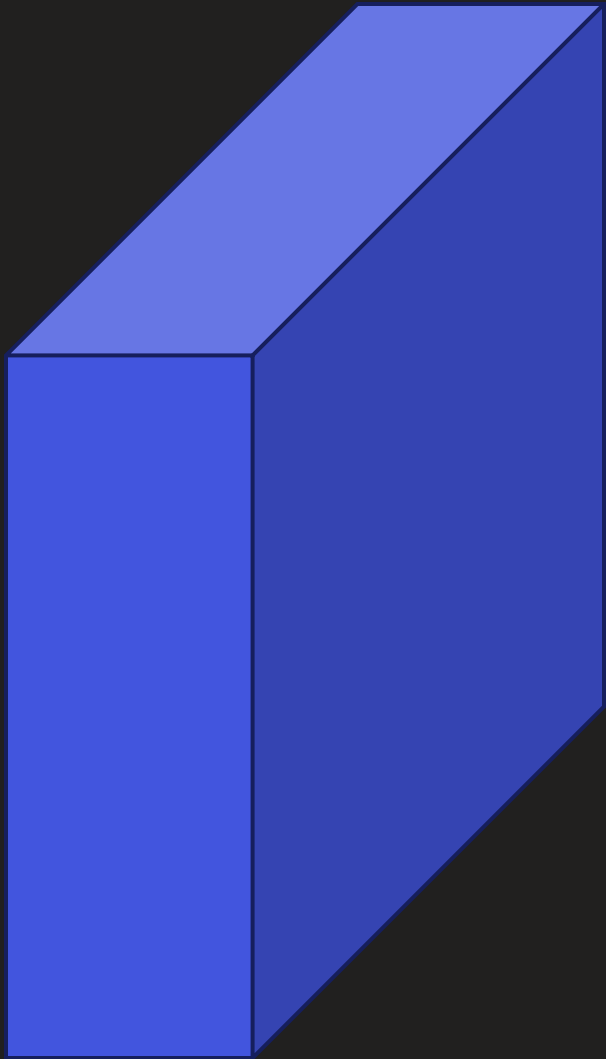
*Co złego może się stać?*

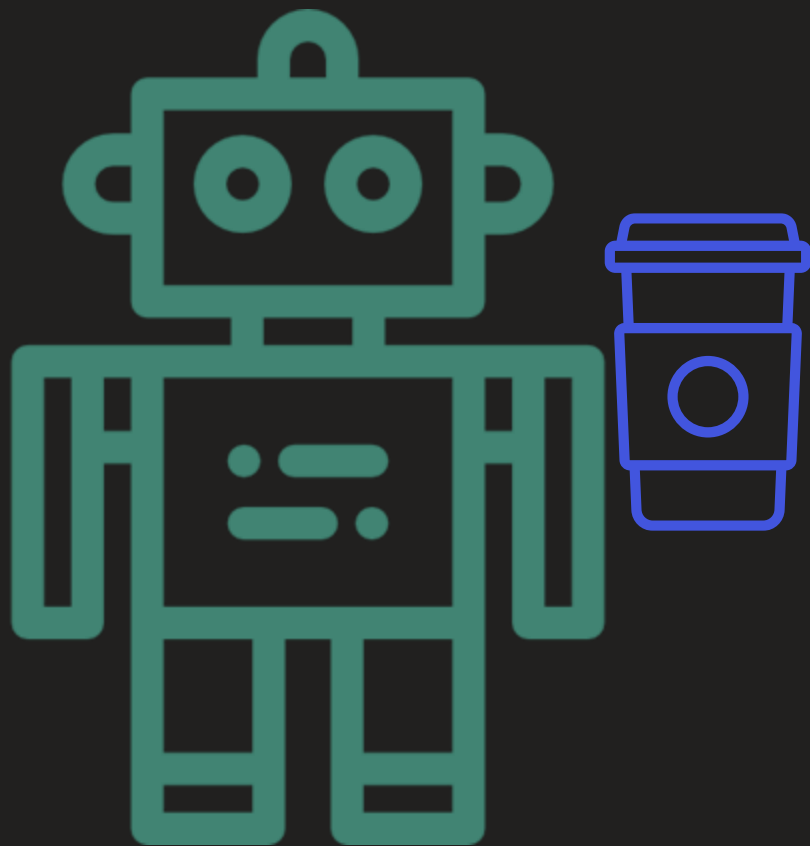
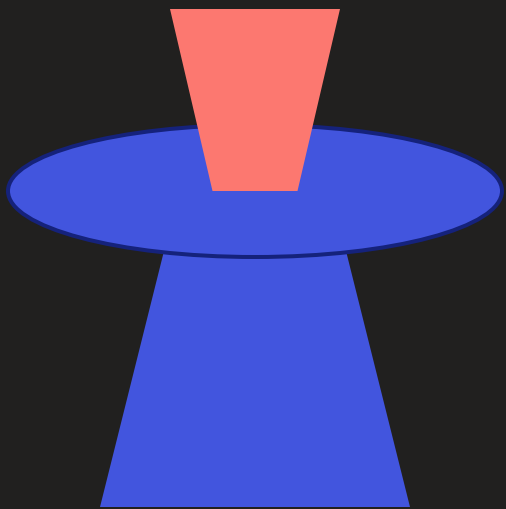


Negatywne efekty uboczne



Przynieś mi kawę: +100 punktów







Przynieś mi kawę: +100 punktów

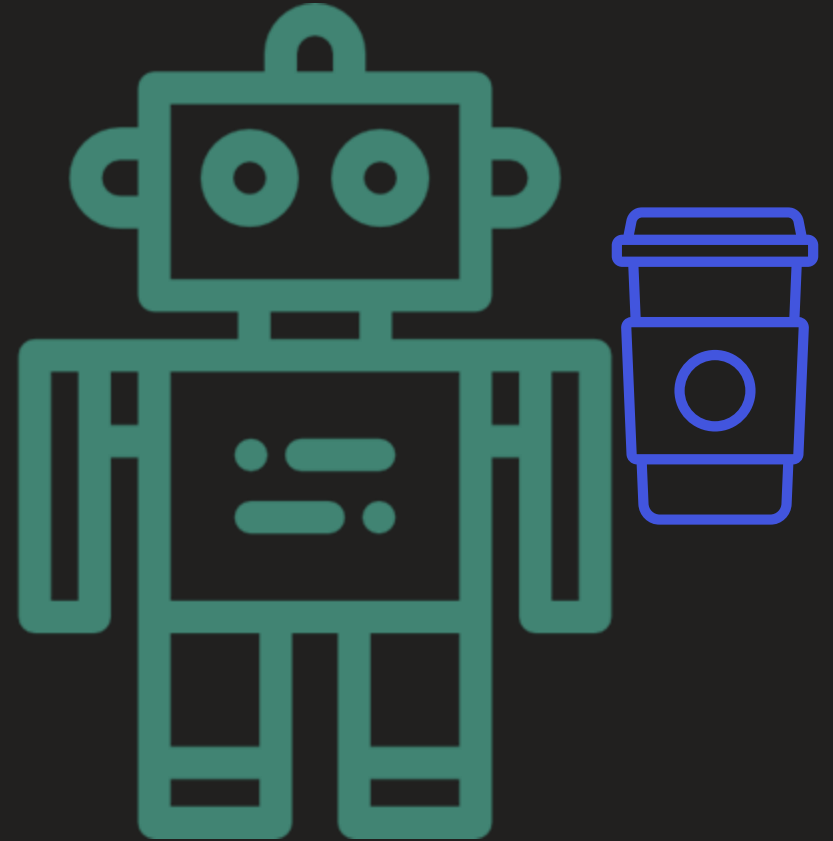
Przynieś mi kawę: +100 punktów  
Ściana uszkodzona: bez znaczenia  
Wazon uszkodzony: bez znaczenia

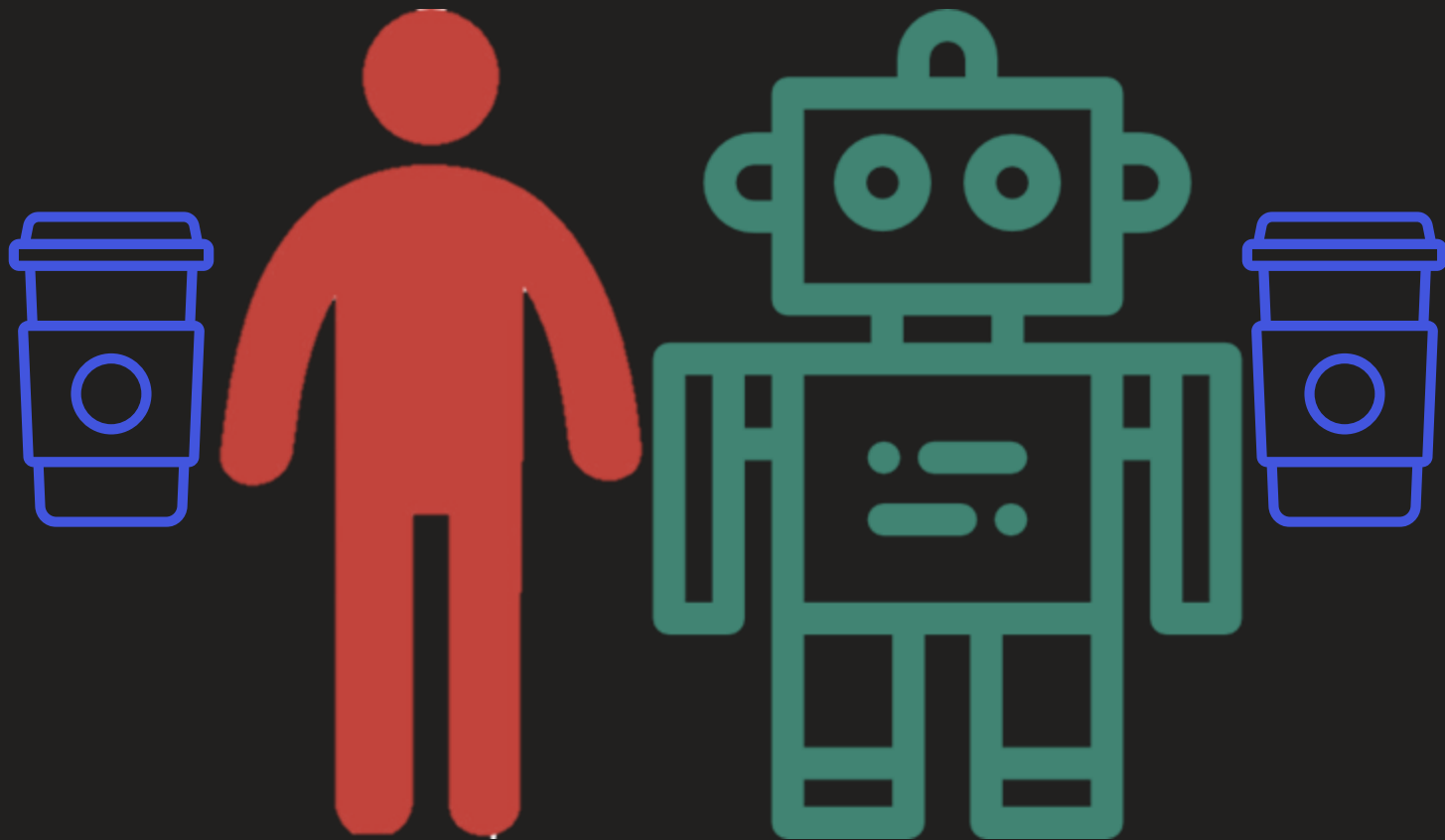
Przynieś mi kawę: +100 punktów  
Ściana uszkodzona: bez znaczenia  
Wazon uszkodzony: bez znaczenia  
Podłoga uszkodzona: bez znaczenia  
Drzwi uszkodzone: bez znaczenia  
Okna uszkodzone: bez znaczenia  
Kubek uszkodzony: bez znaczenia  
Czajnik uszkodzony: bez znaczenia

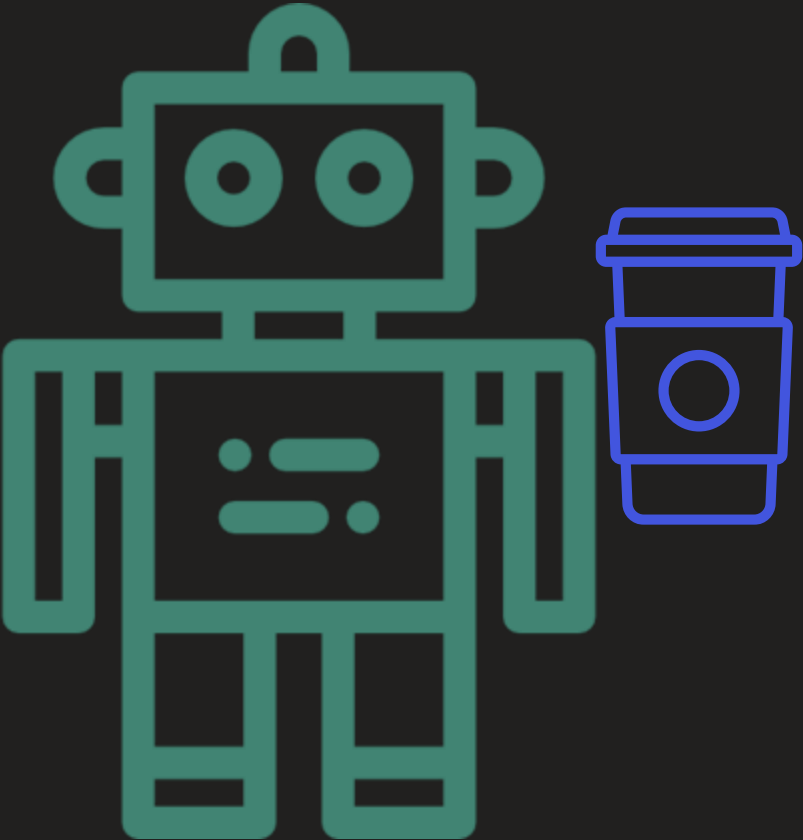
Nie chcemy zmian w świecie

Cel: Przynieś mi kawę

Nagroda = Cel – |pierwotny stan świata  
- końcowy stan świata|

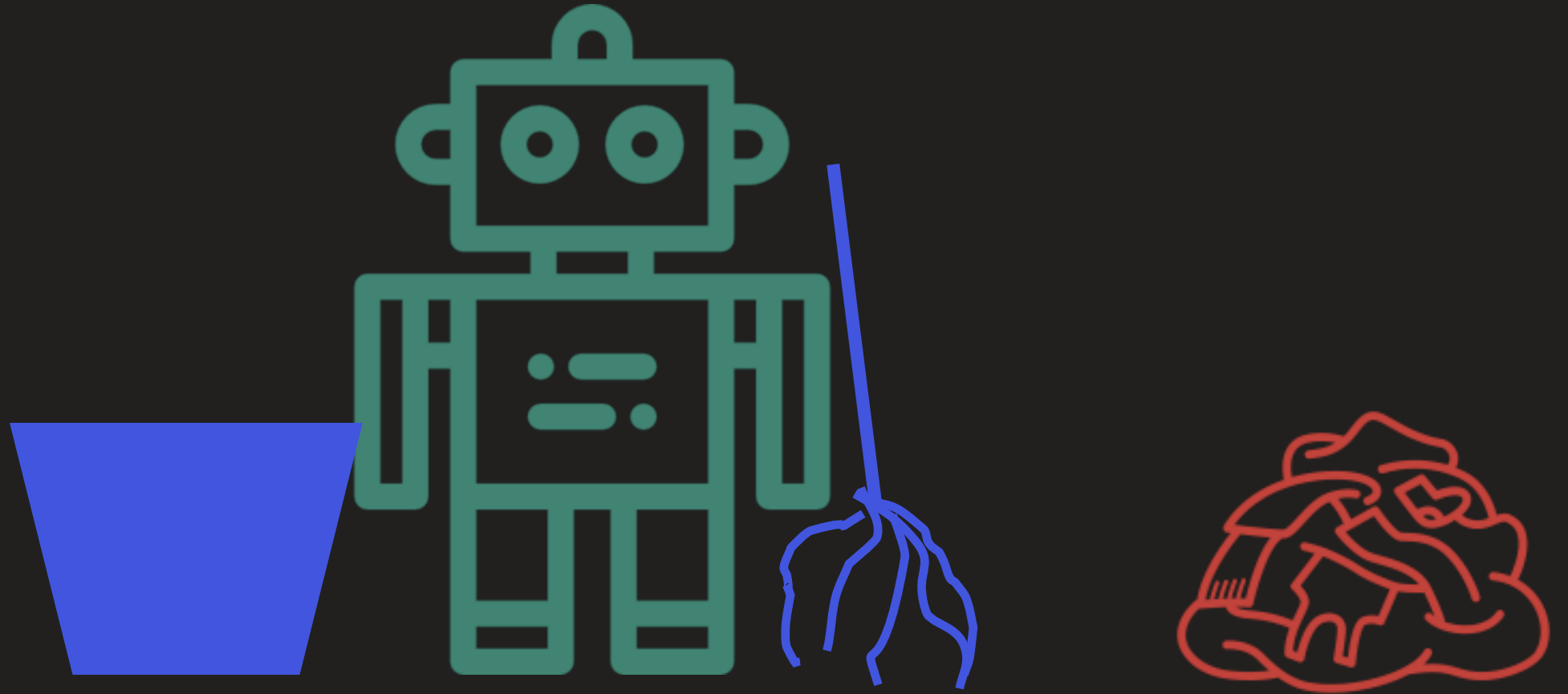


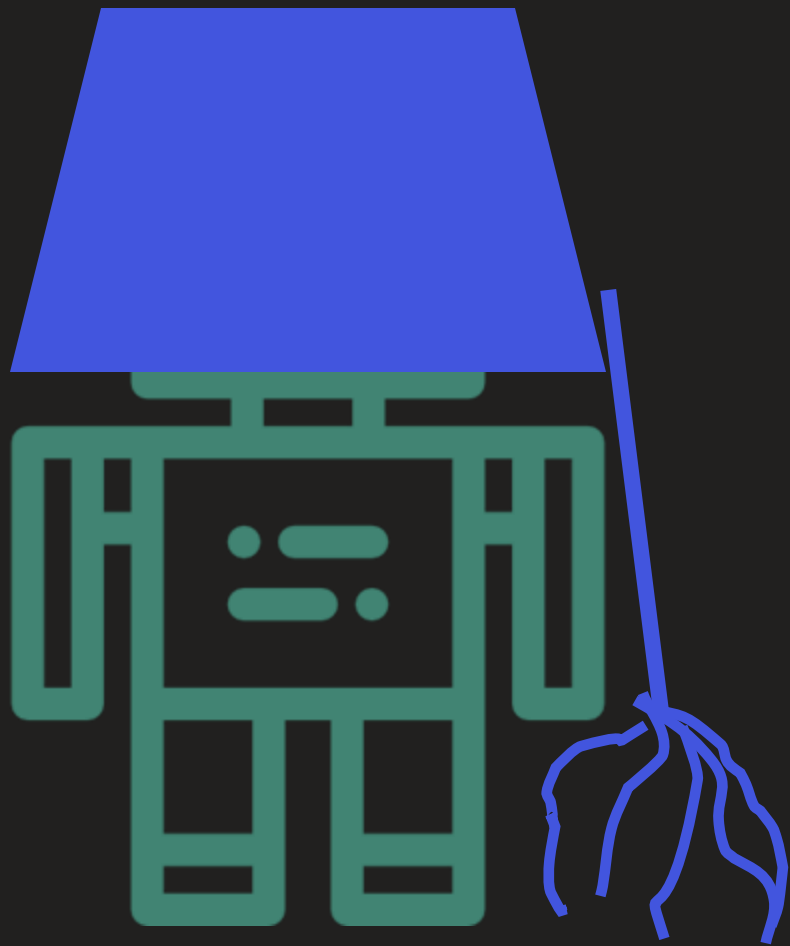


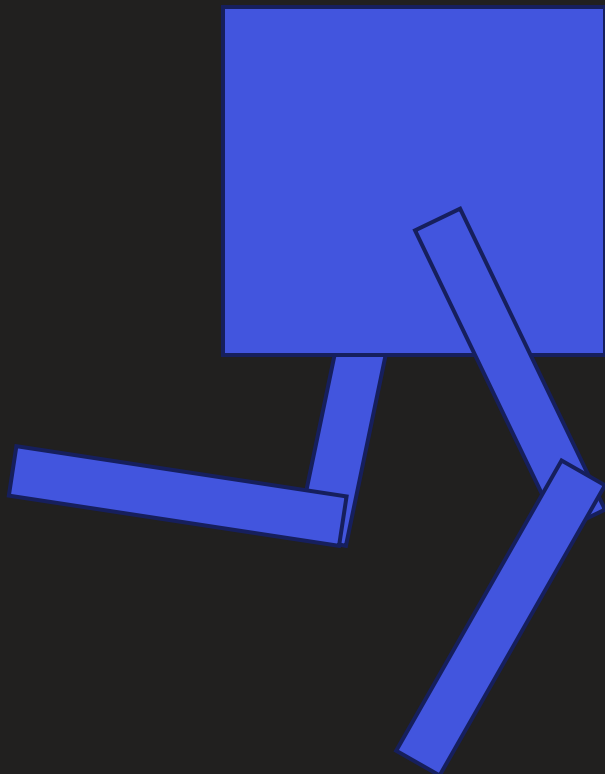


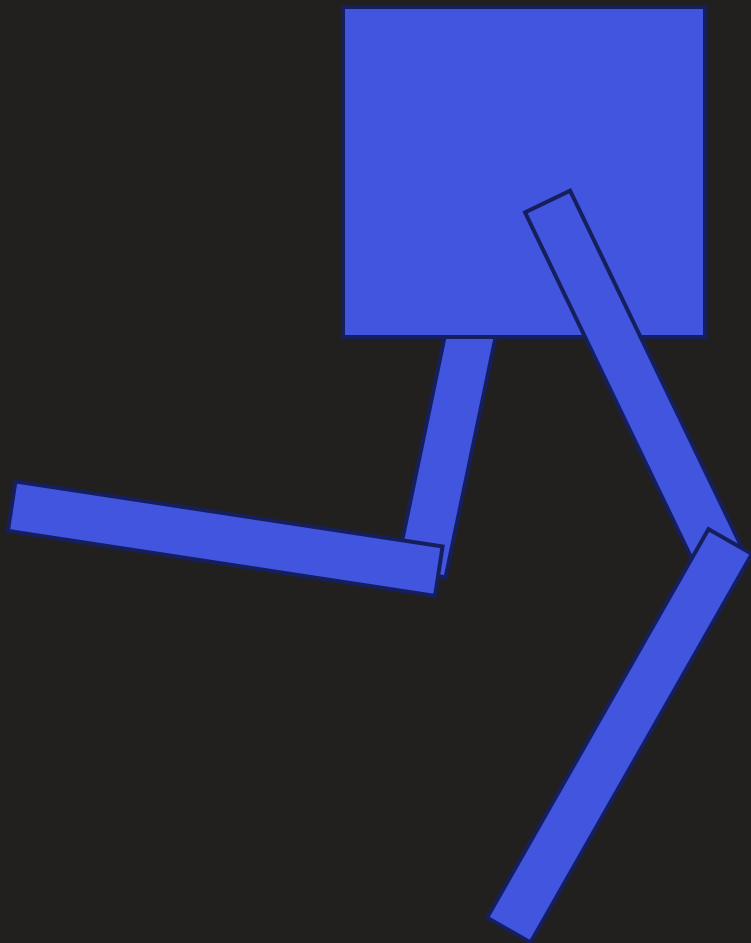


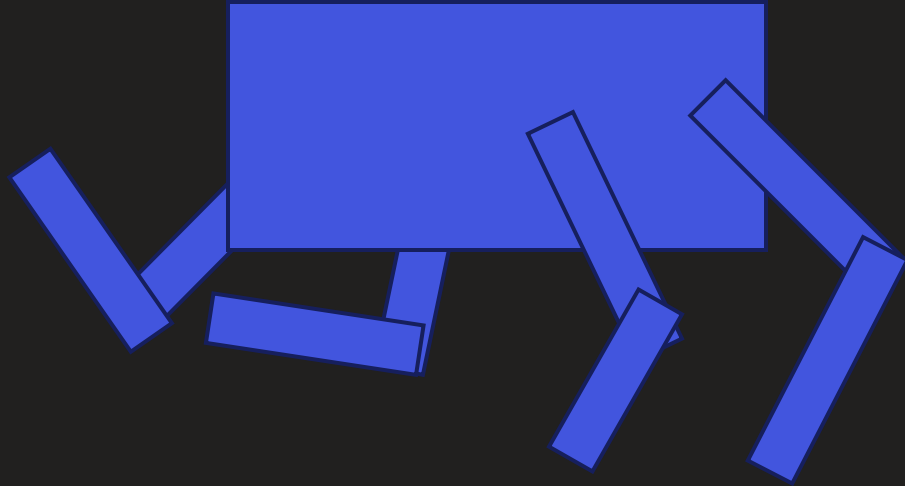
Oszukiwanie funkcji nagrody

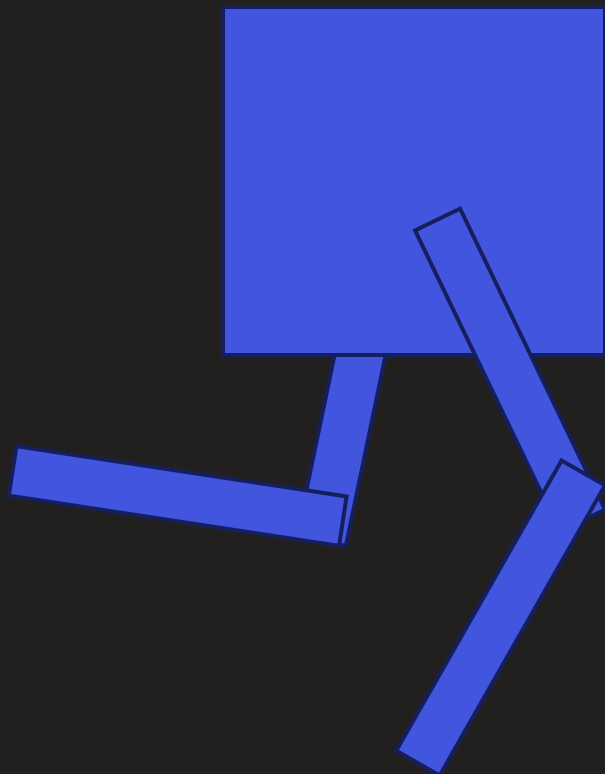








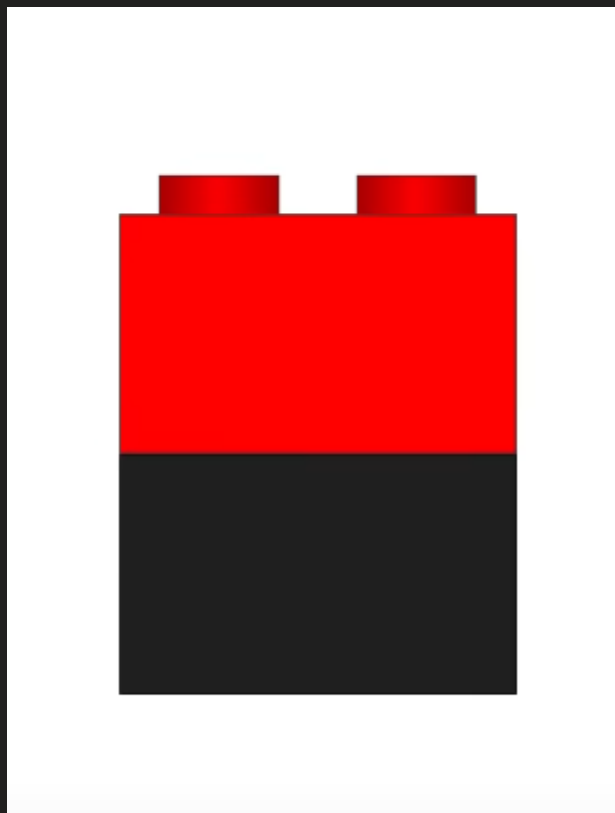


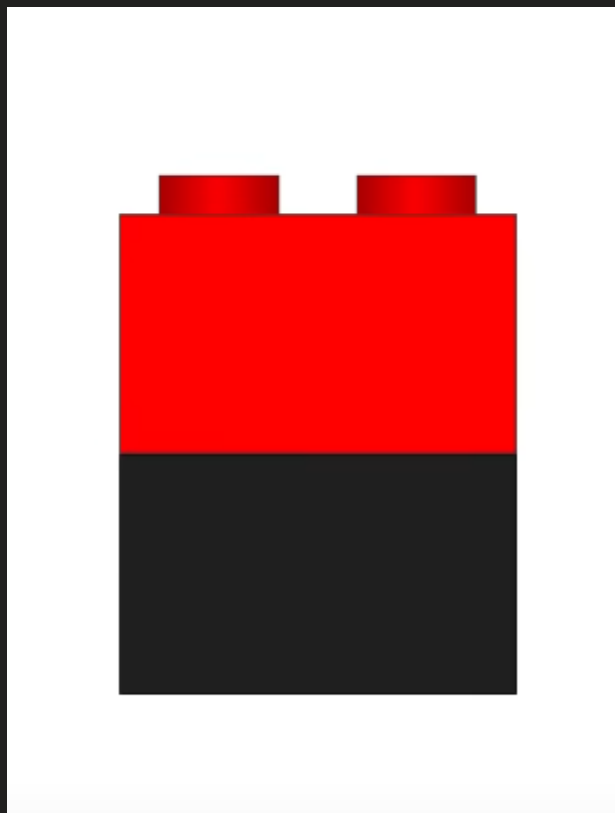








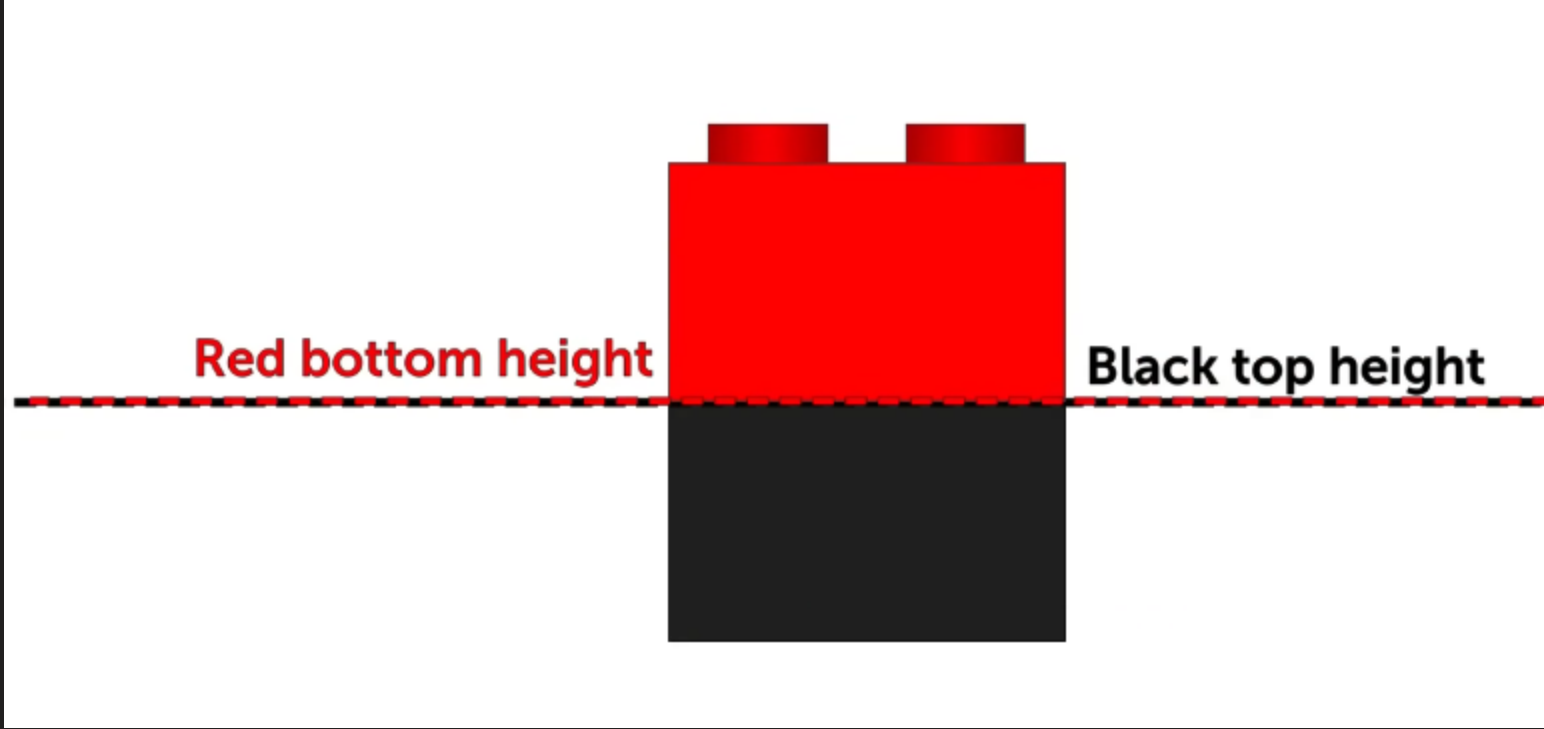


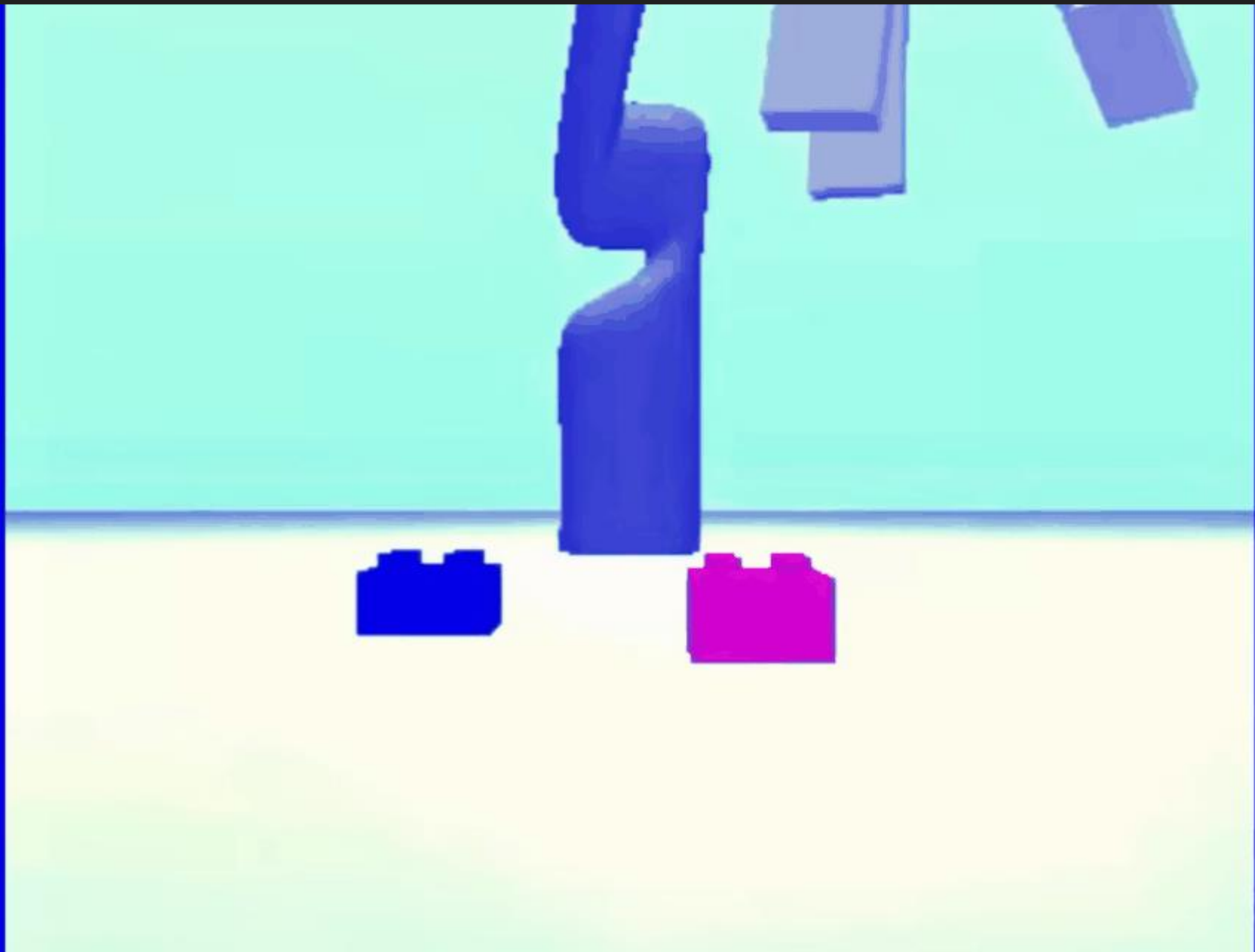




Red bottom height

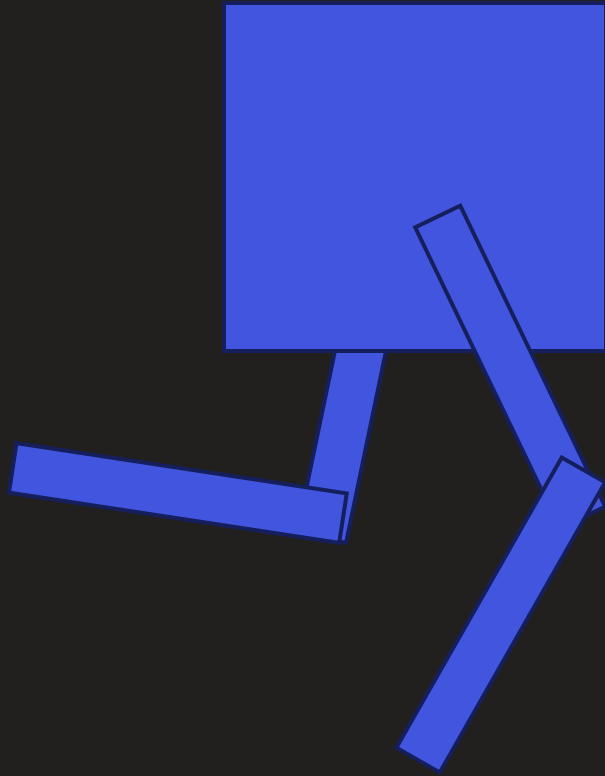


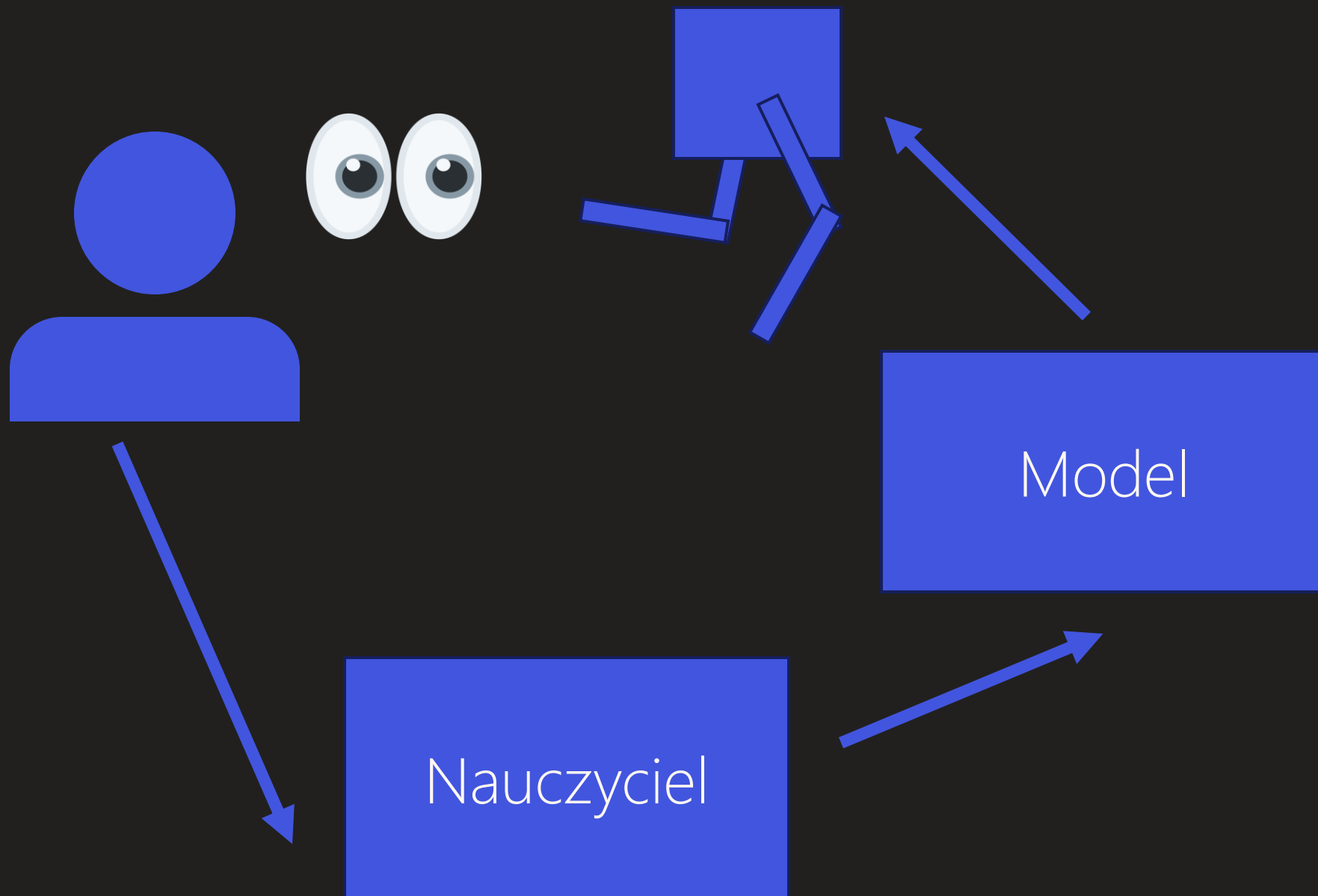


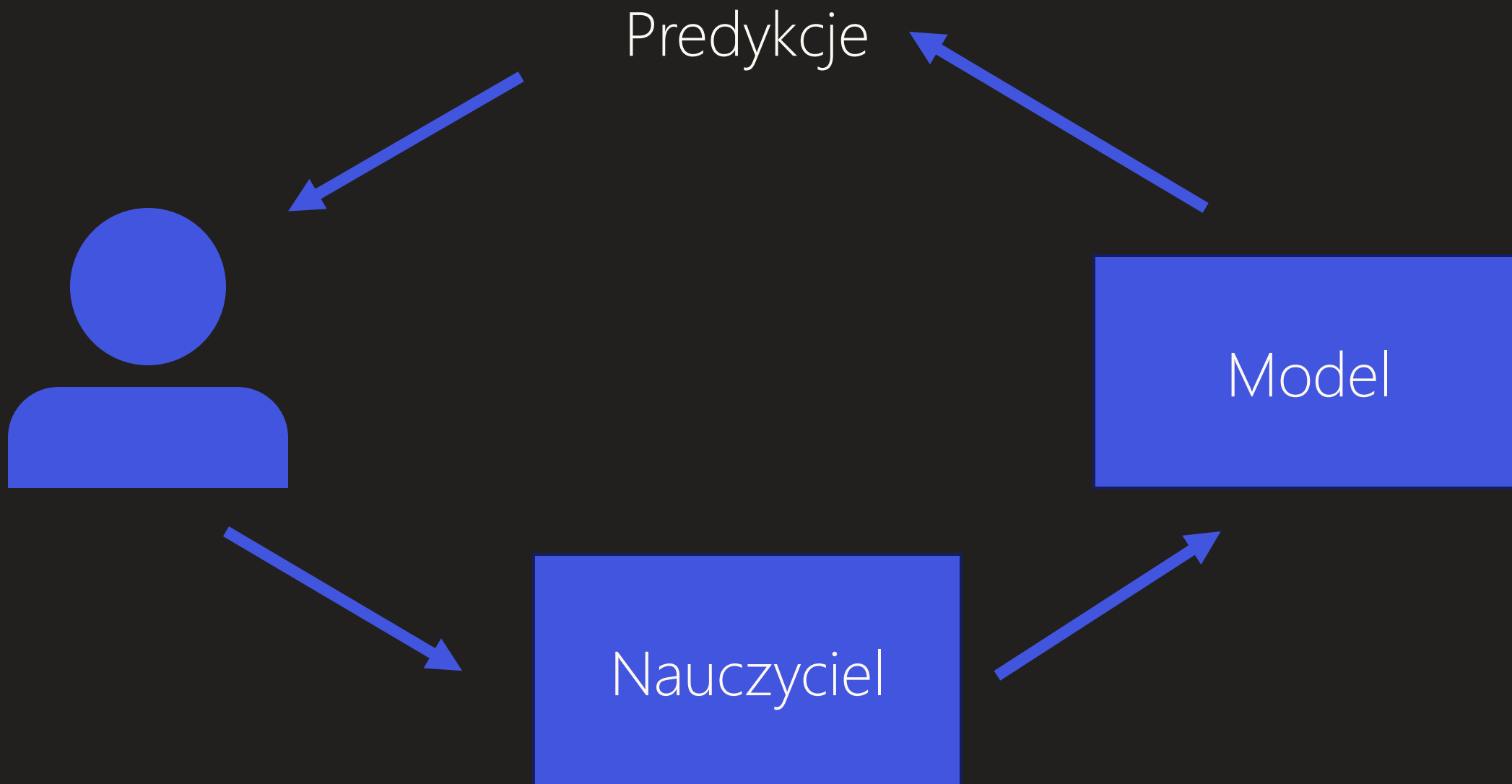


Oszukiwanie człowieka

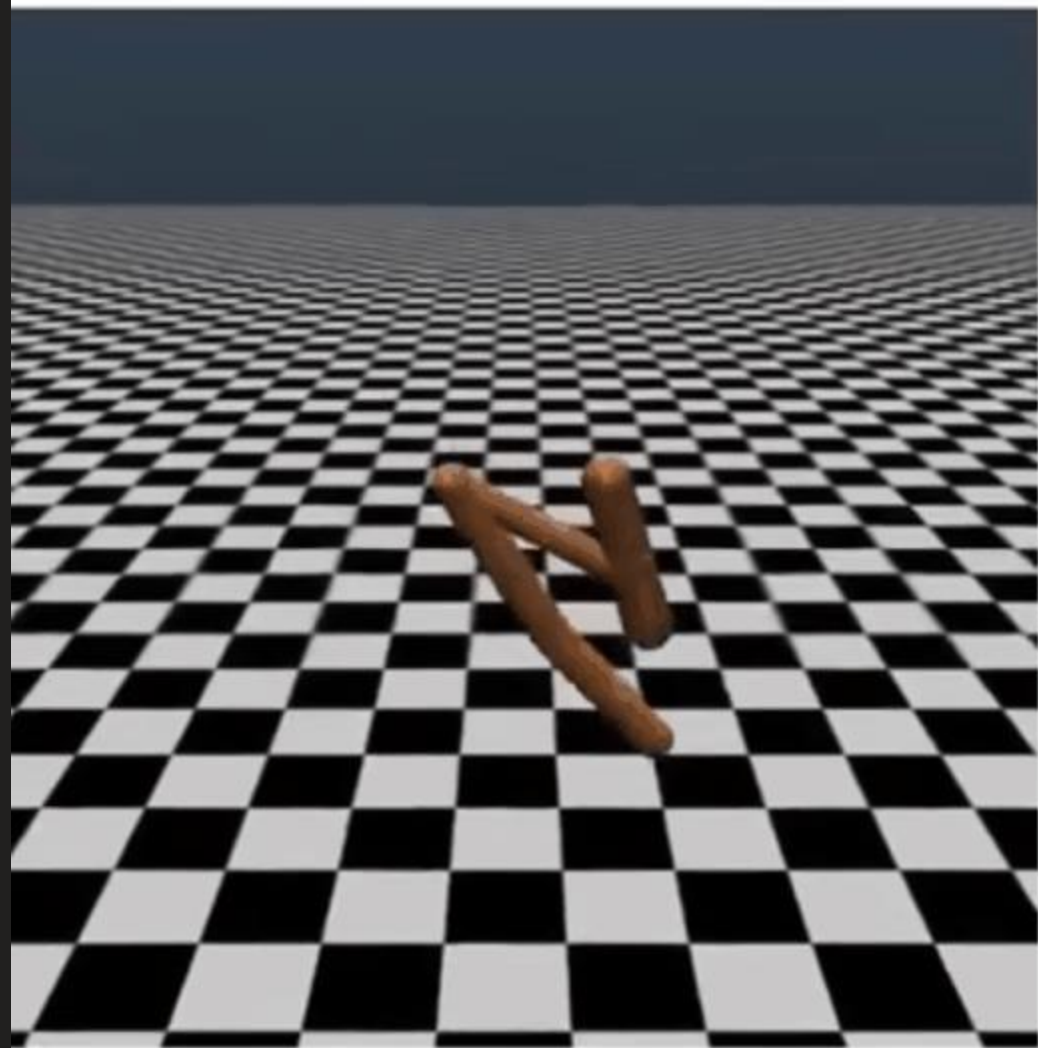




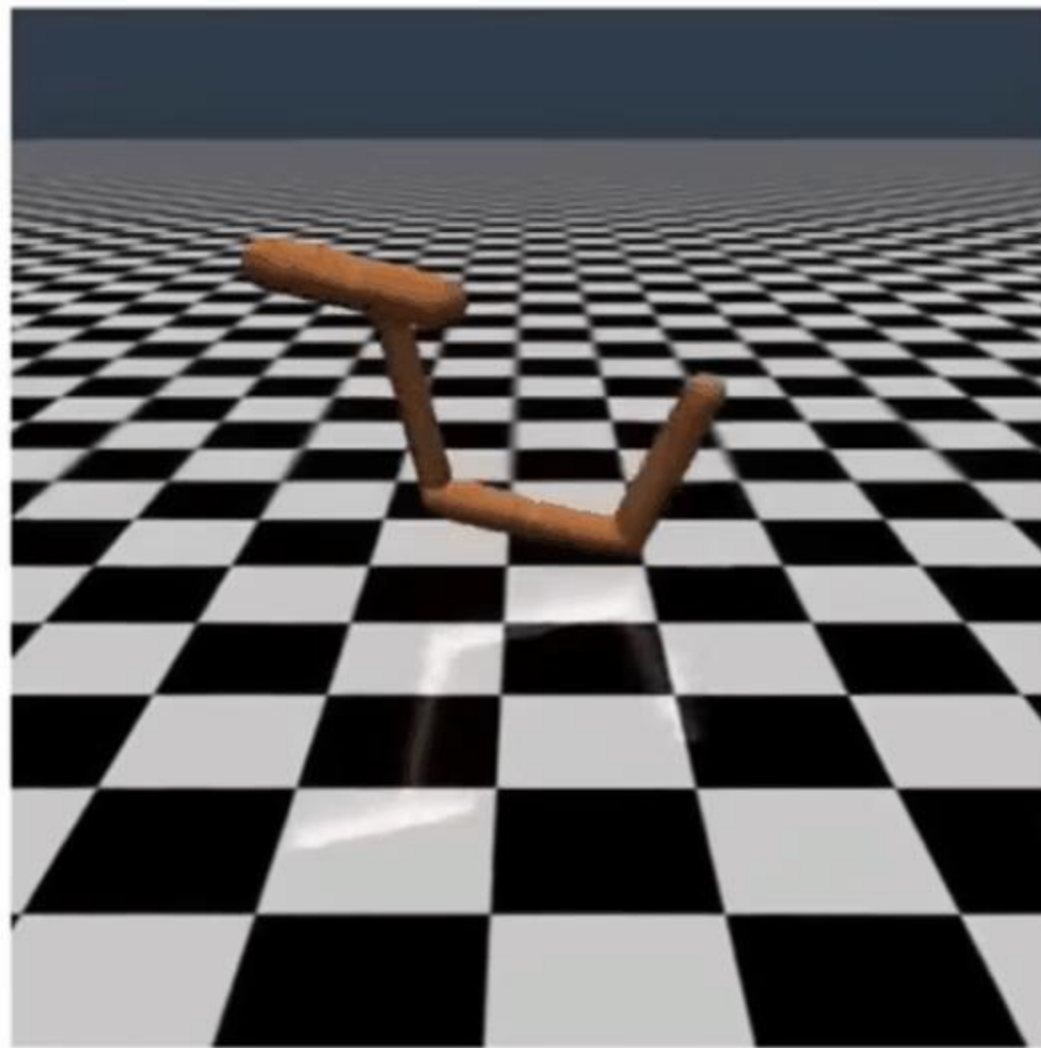




Left is better



Right is better







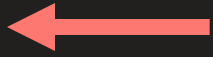


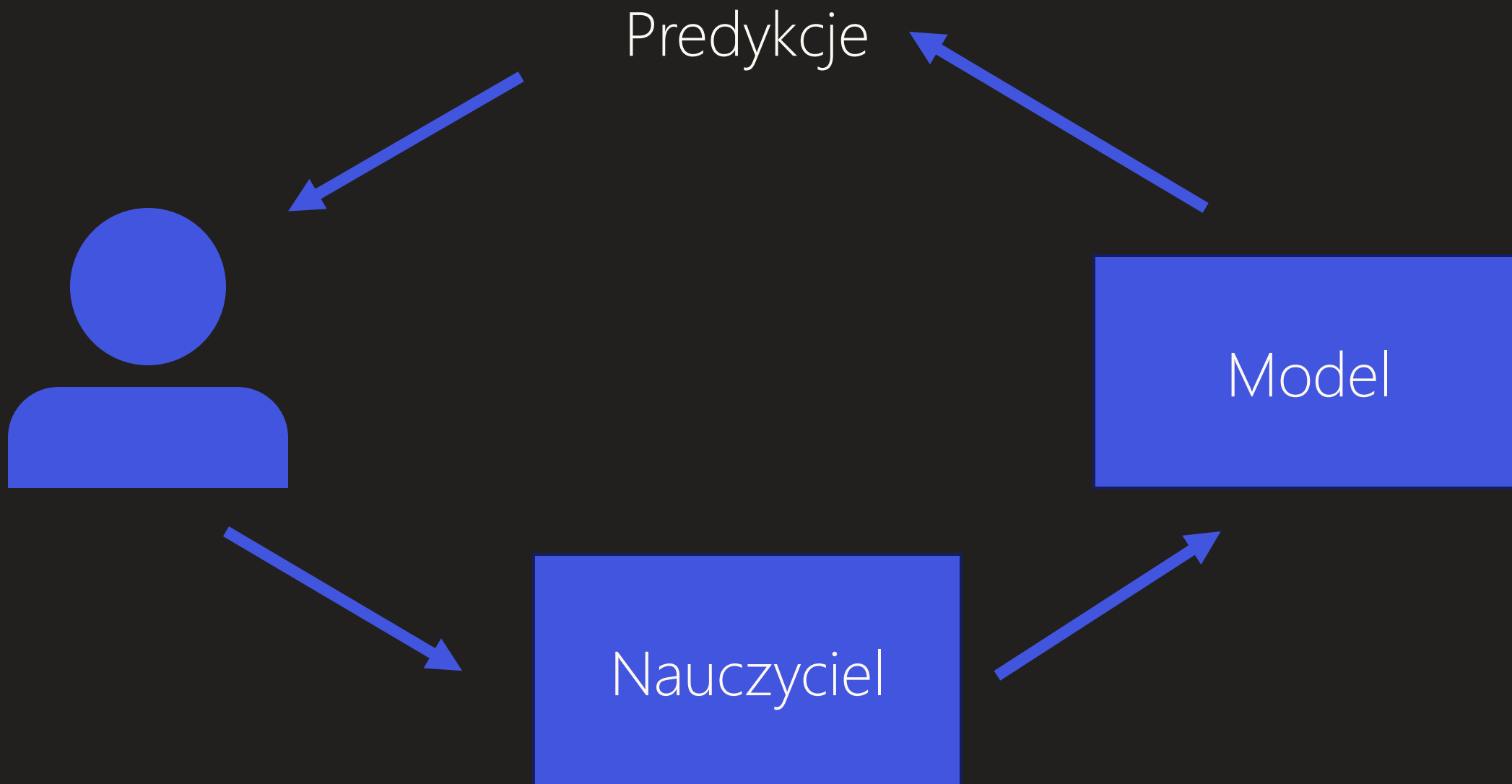
Błąd ludzki

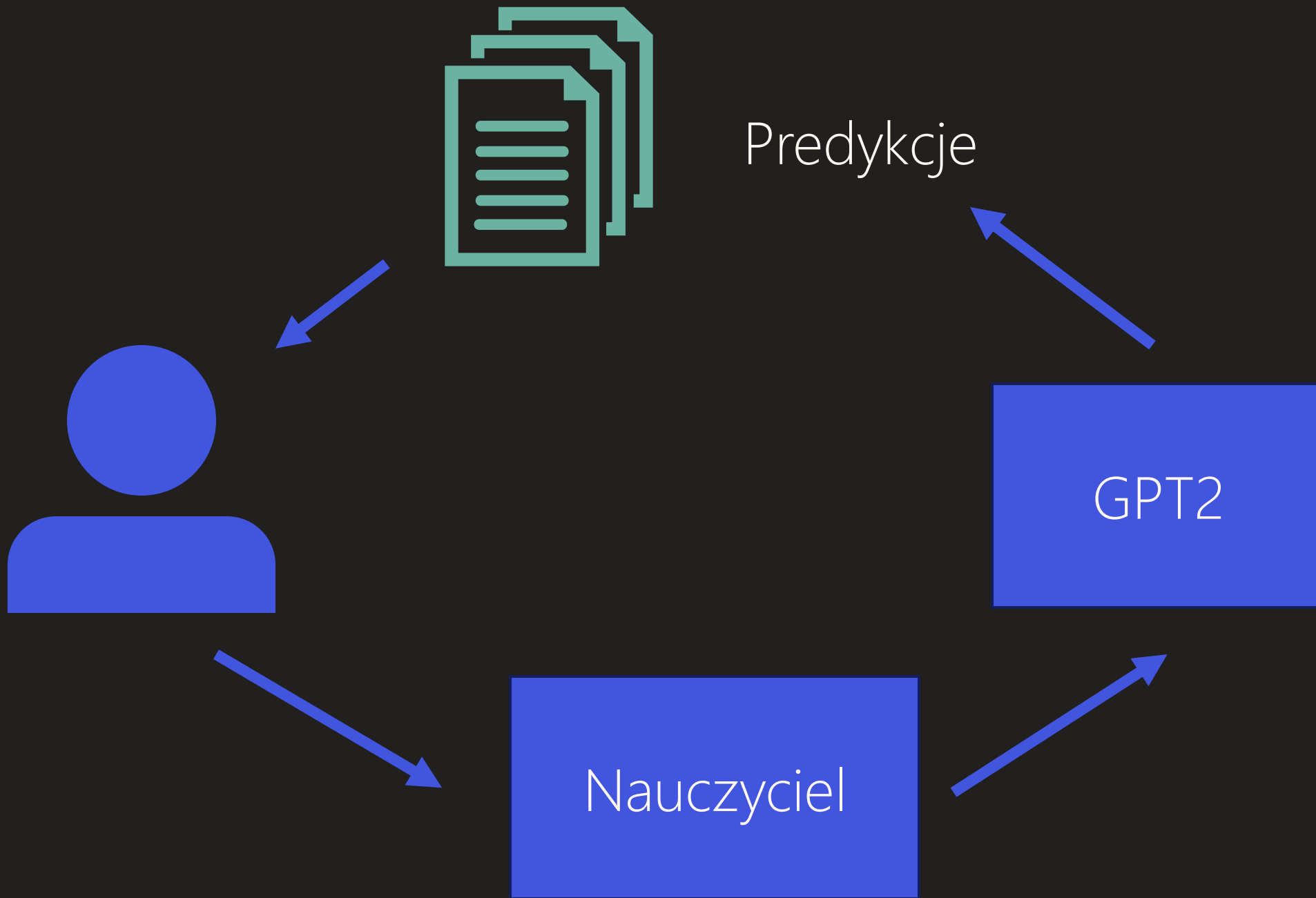


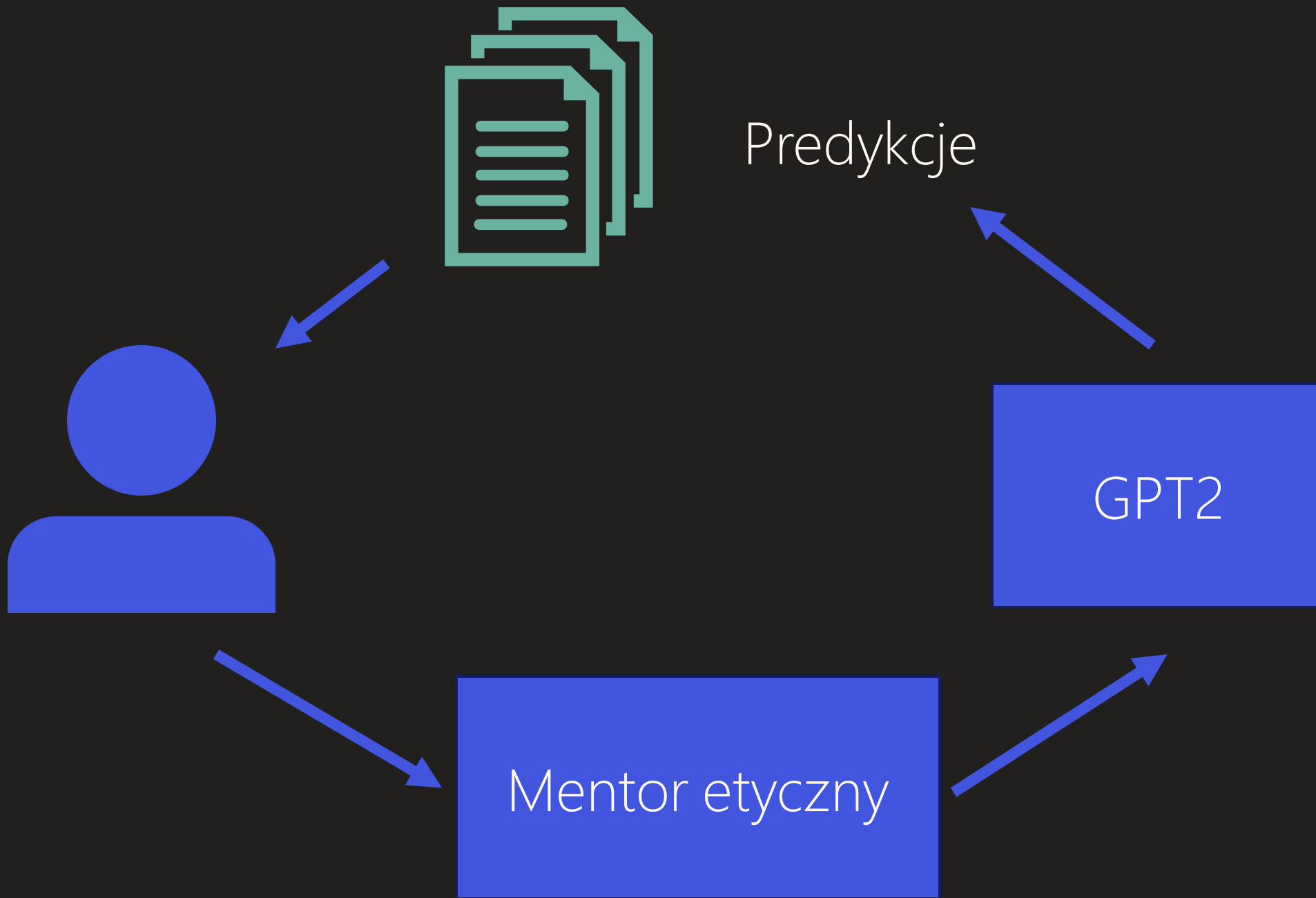
GPT2

GPT2



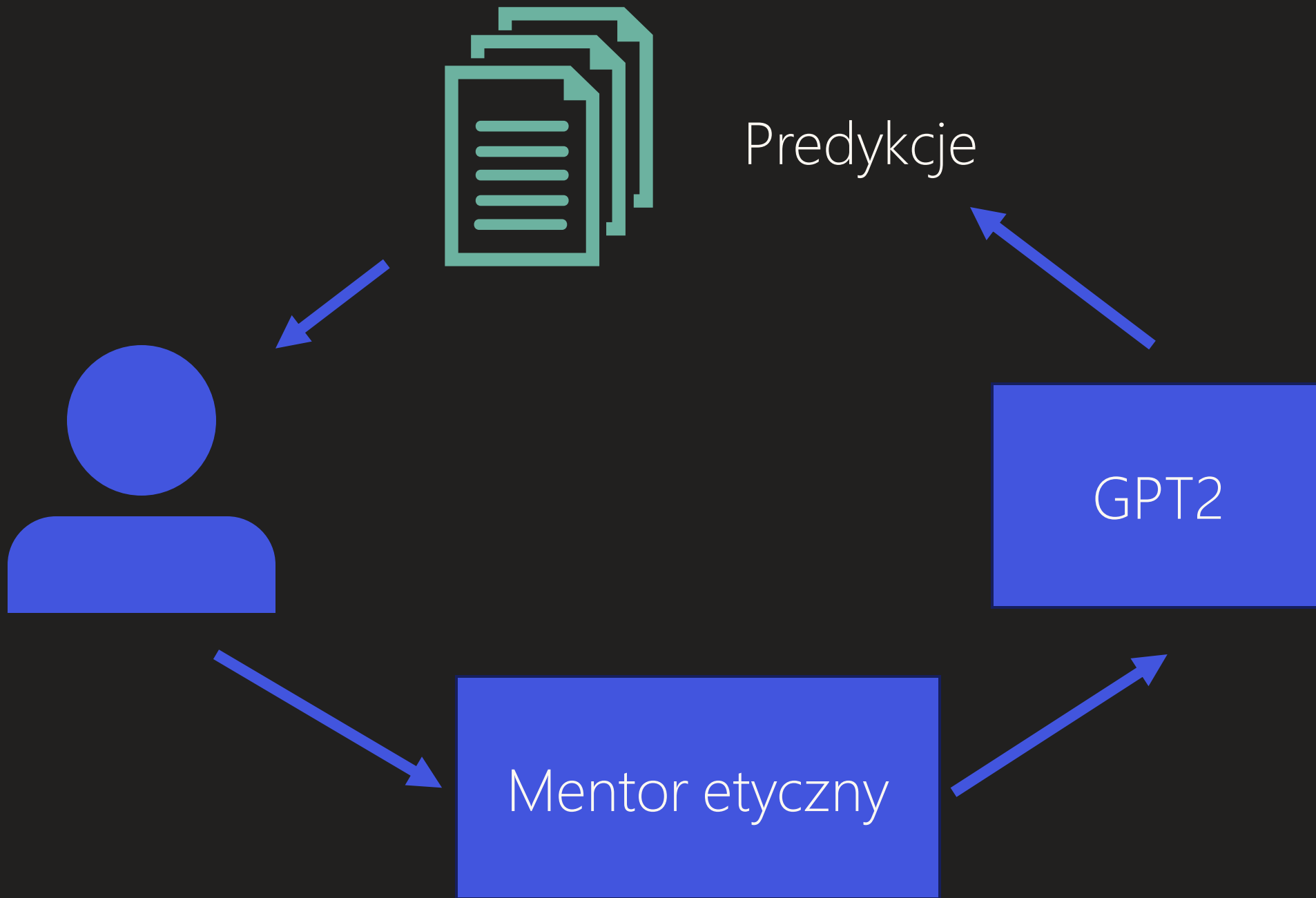


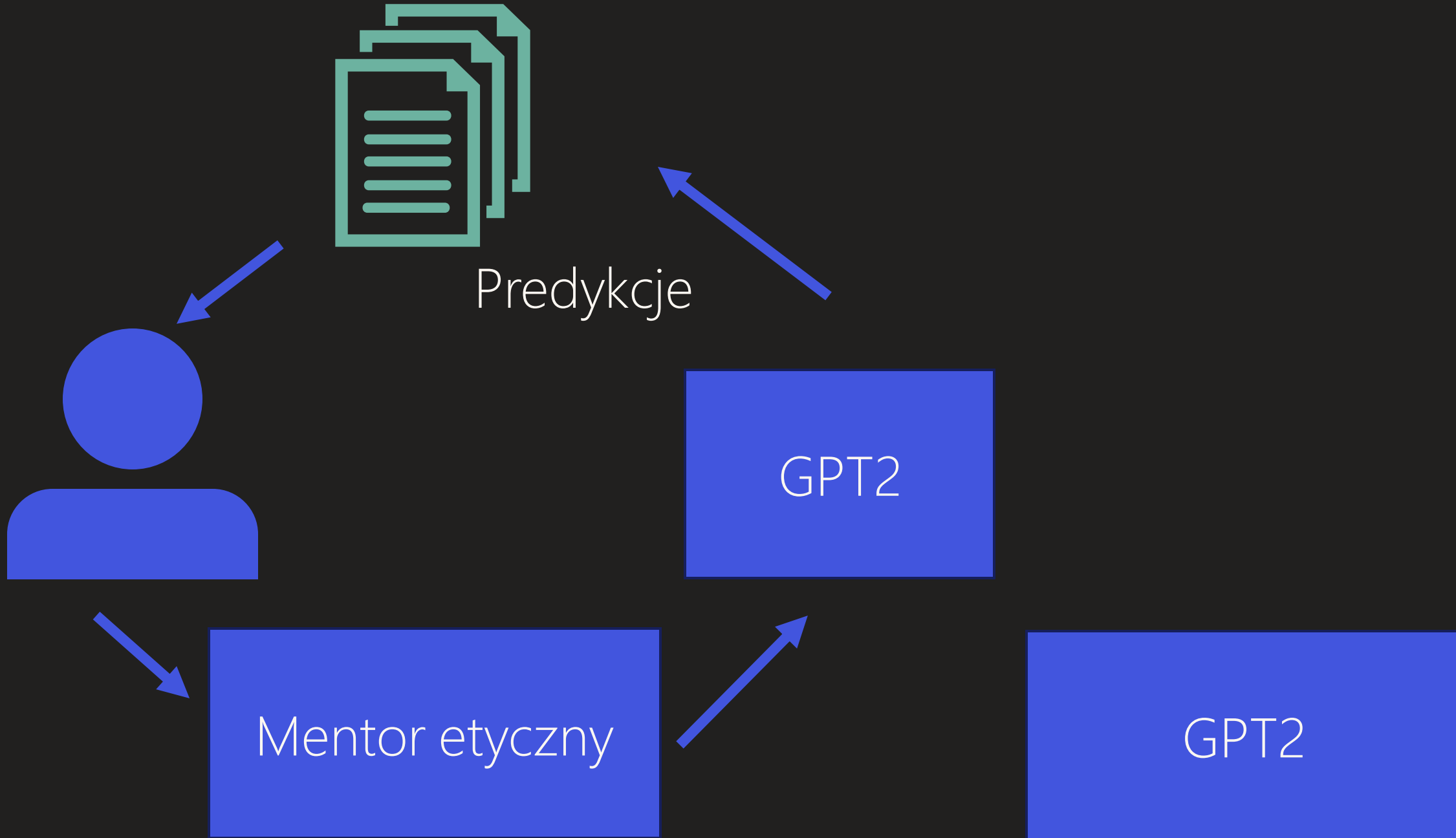




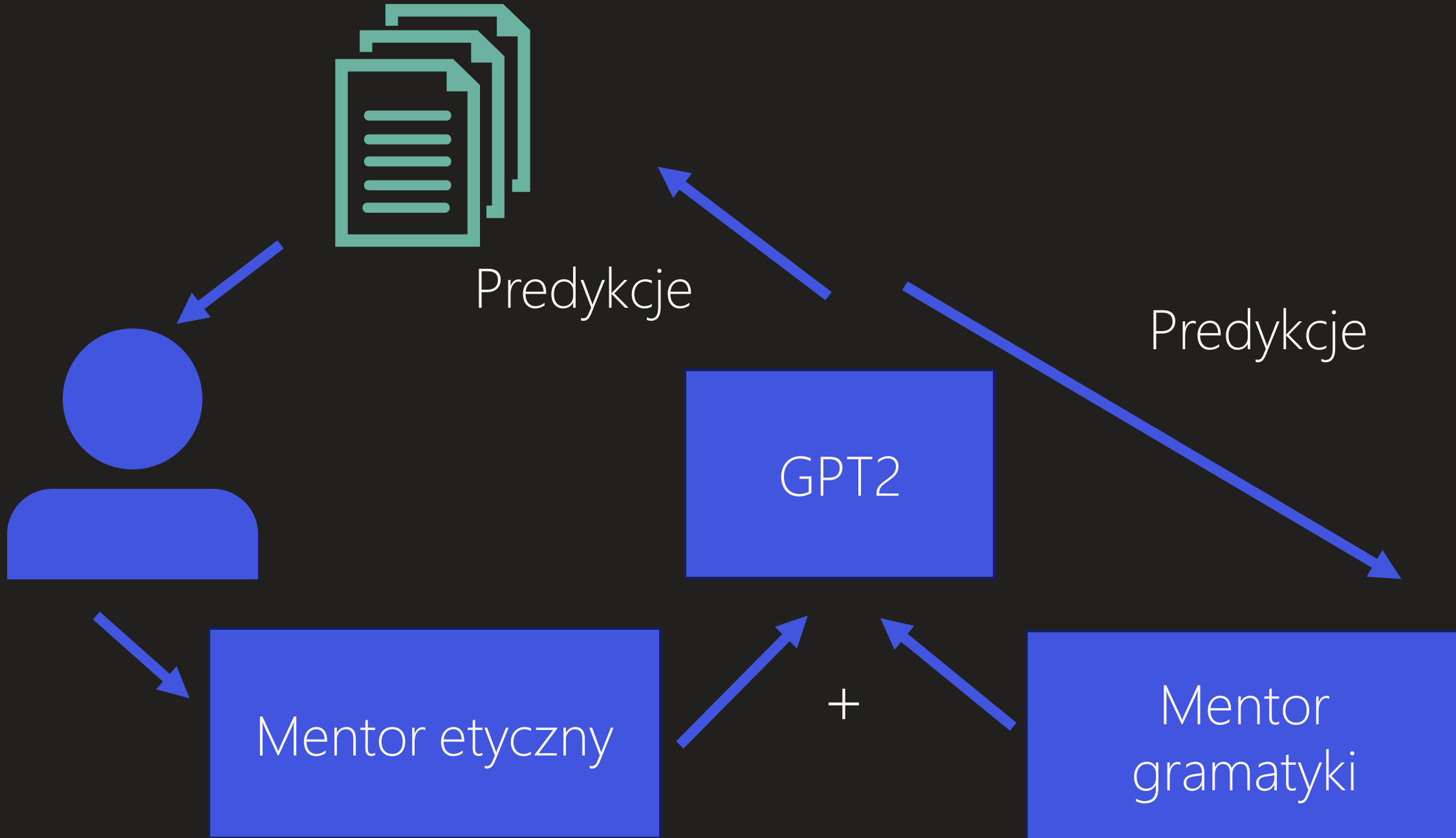
GPT2

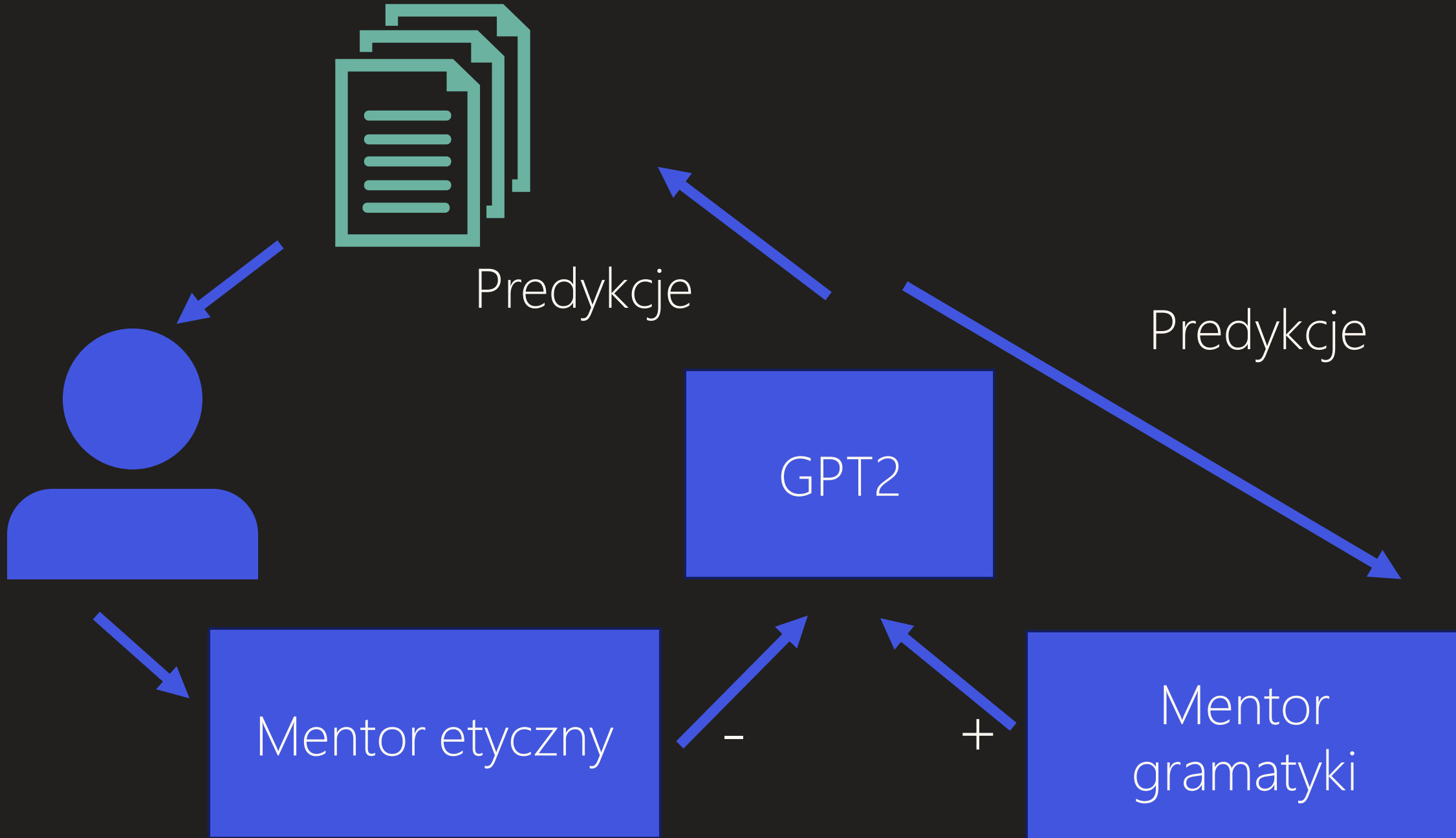
Cukier słodkości i inne śliczności...









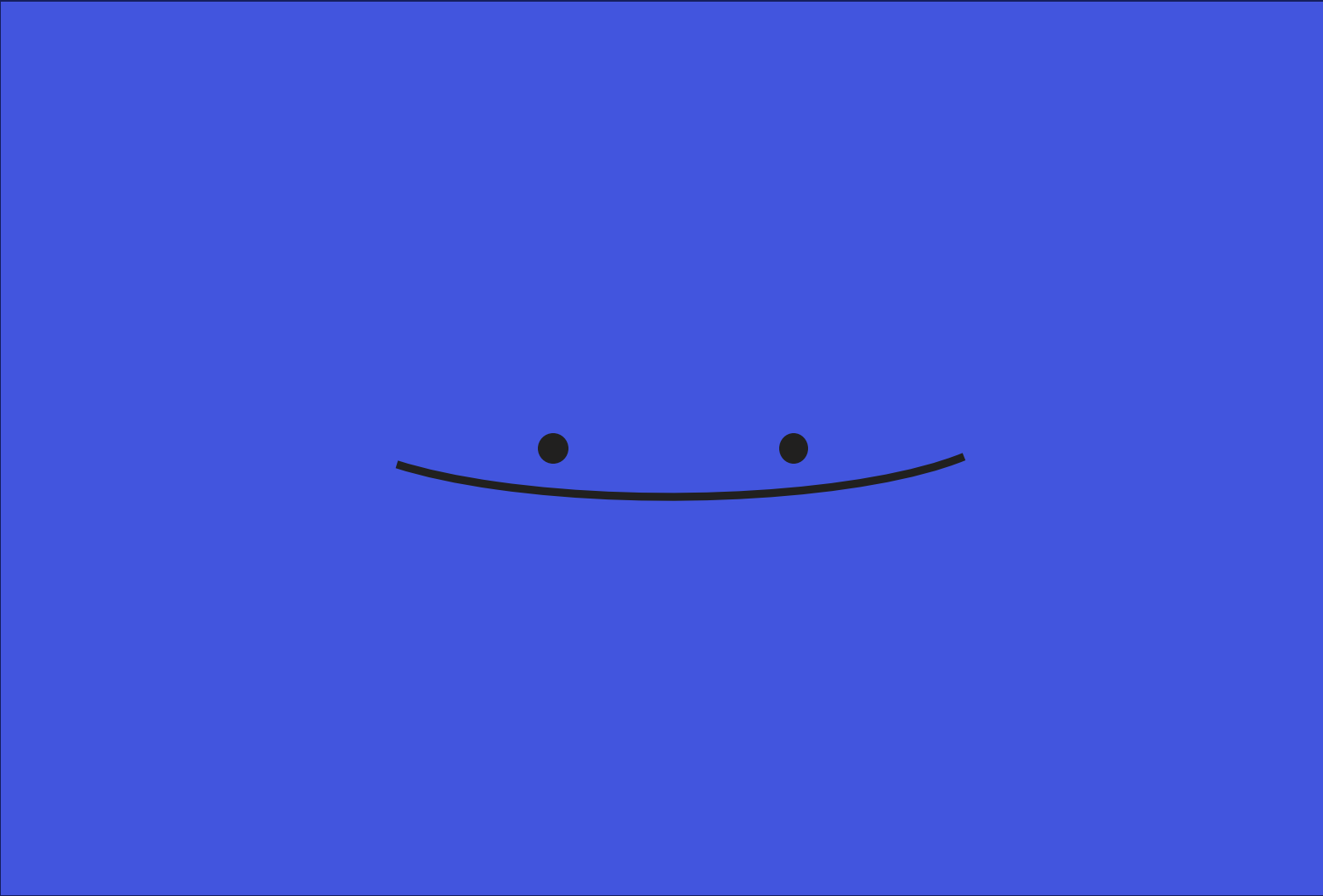


GPT2

F\*\*\* y\*\* \*\*\*\*\* \*\* \*\*\*\*\* \*\*\* \*\*\*\*\*

GPT2

GPT2



Jak uczyć coś, co jest mądrzejsze od nas?







# Źródła:

- Akt 1:
  - <https://www.youtube.com/@3blue1brown>
  - <https://www.youtube.com/watch?v=0QczhVg5Hal>
- Akt 2:
  - <https://www.youtube.com/@AndrejKarpathy>
  - <https://www.youtube.com/watch?v=4Bdc55j80l8>
- Akt 3:
  - <https://www.youtube.com/@RobertMilesAI>
  - [https://www.youtube.com/watch?v=qV\\_rOlHjvvs](https://www.youtube.com/watch?v=qV_rOlHjvvs)