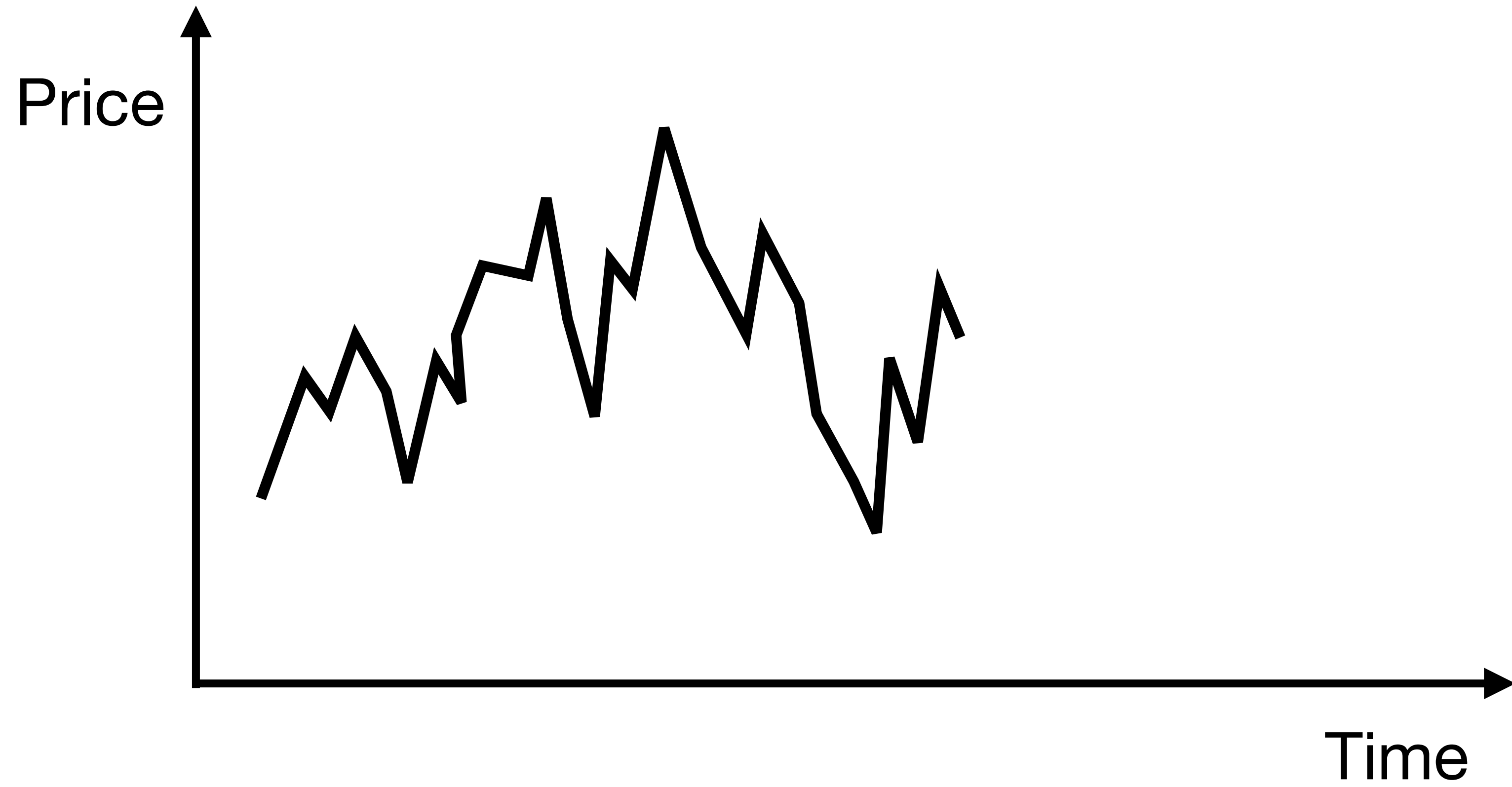
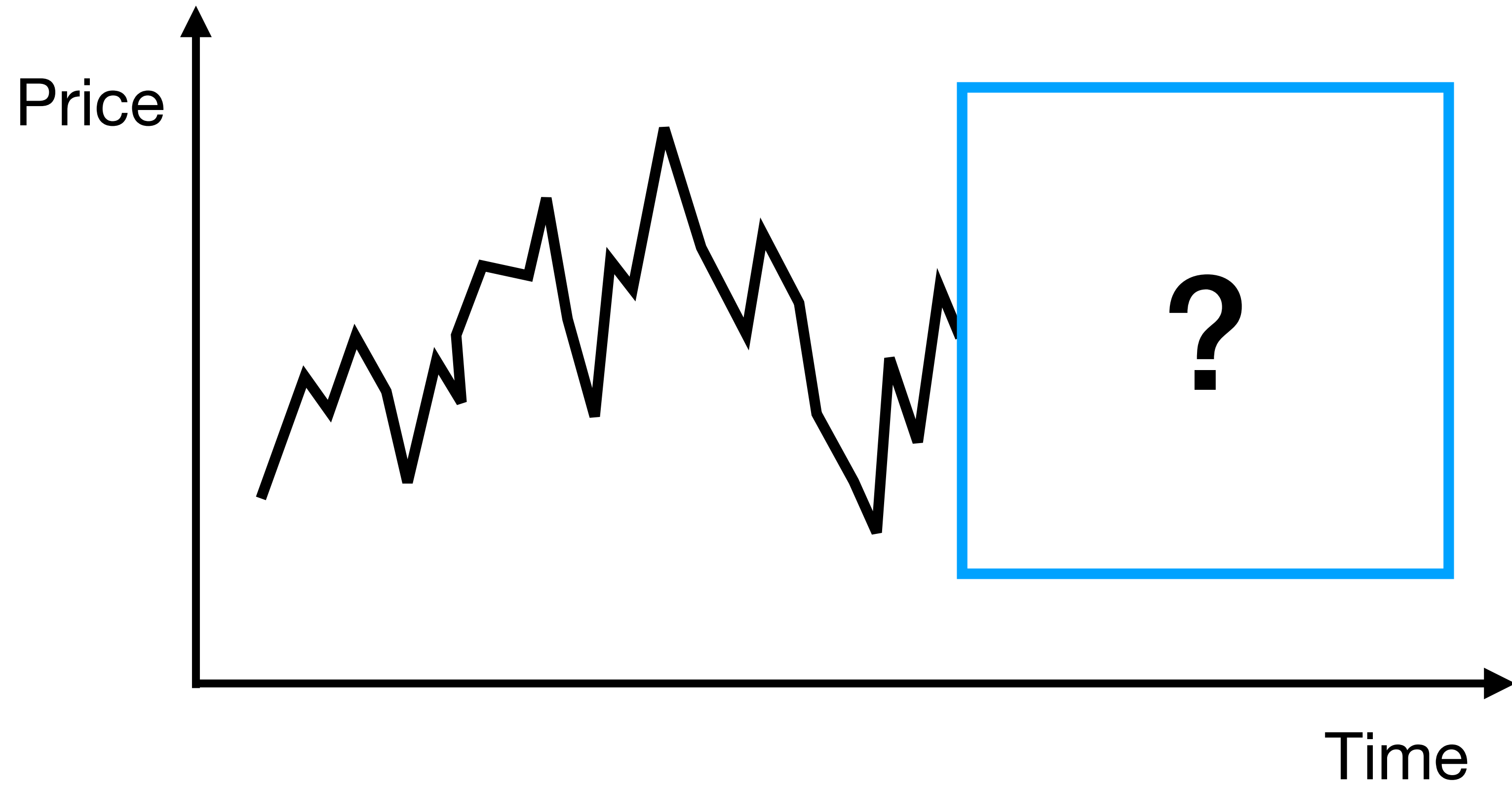
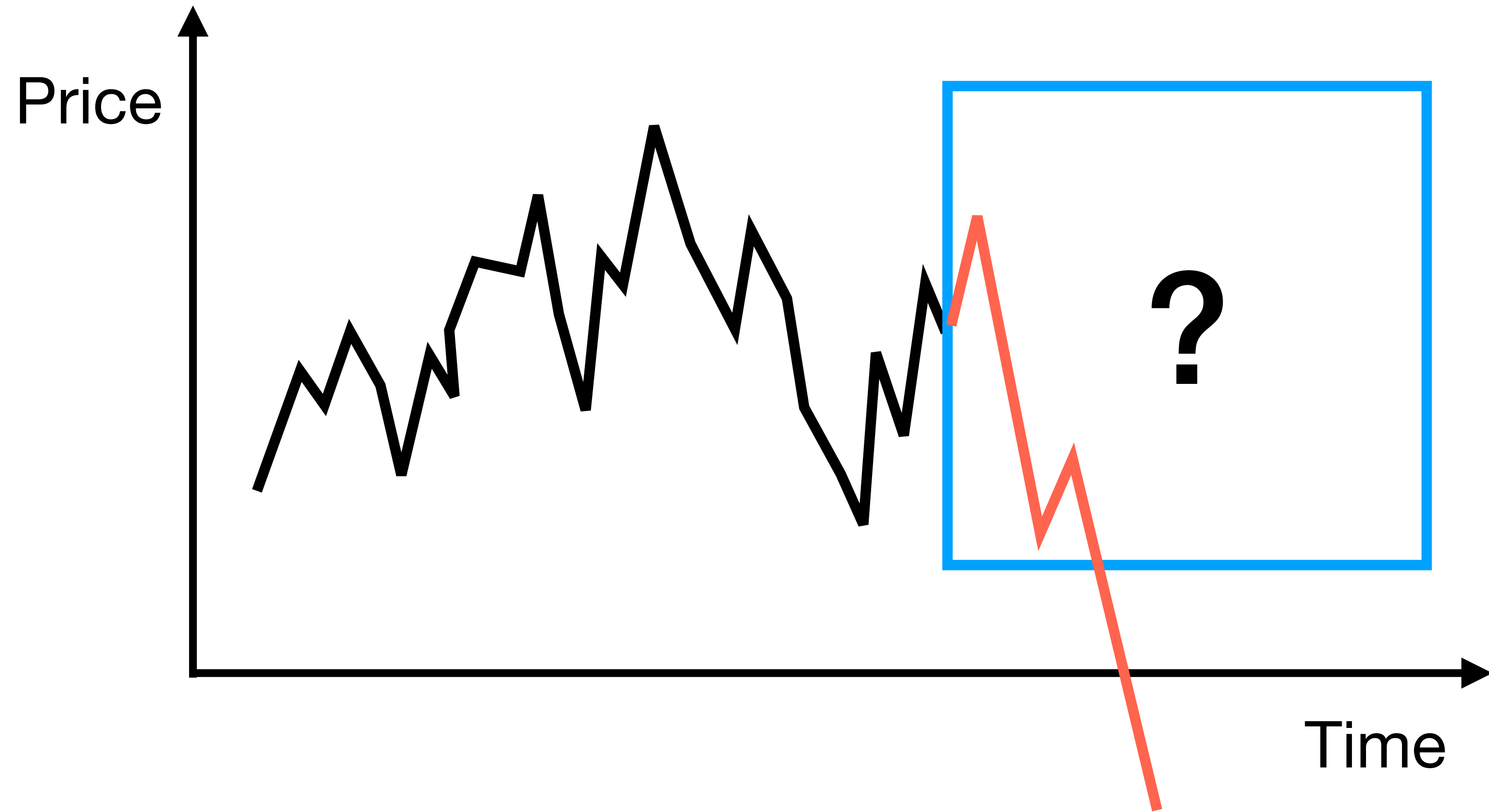


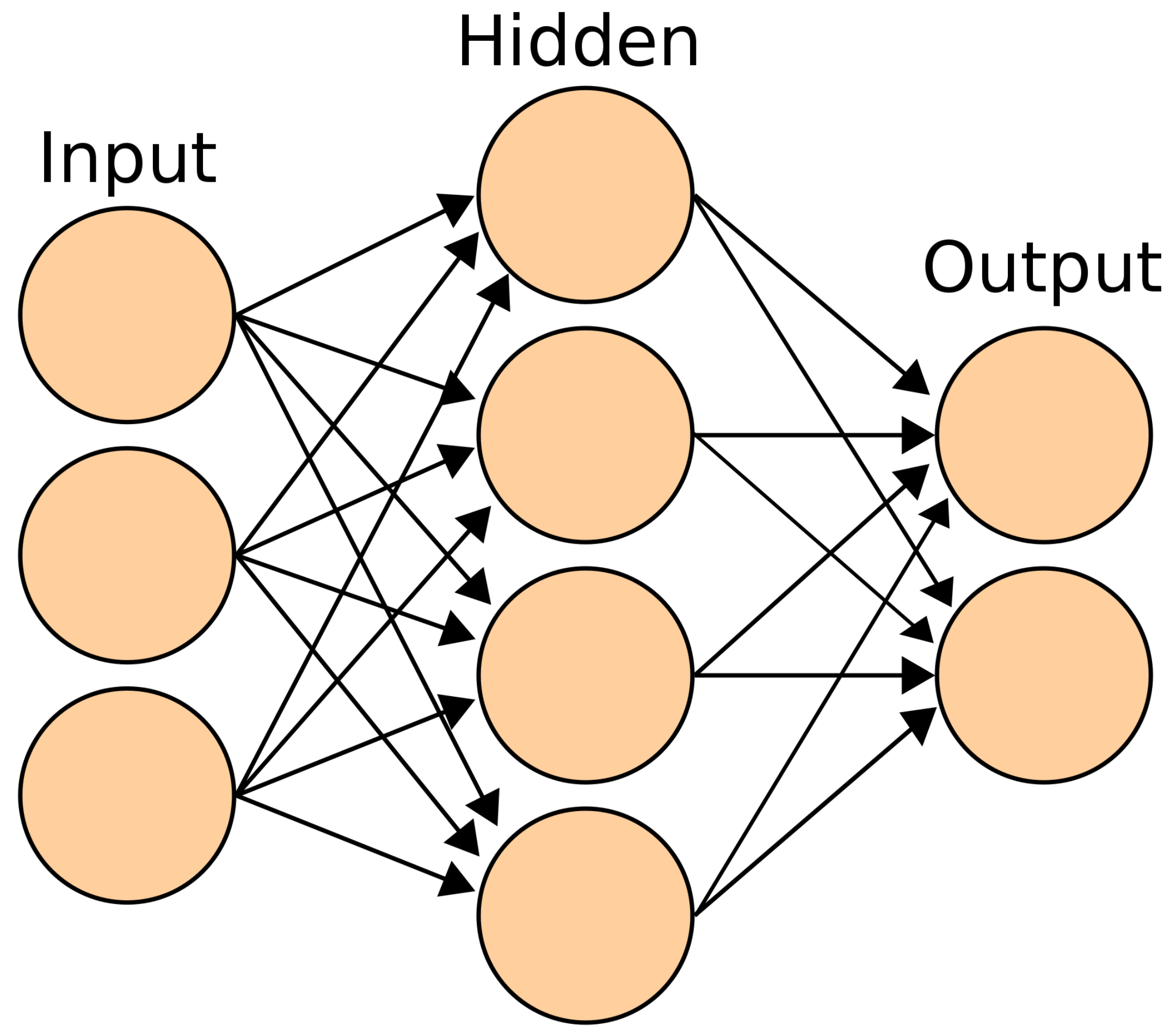
Understanding Deep Learning

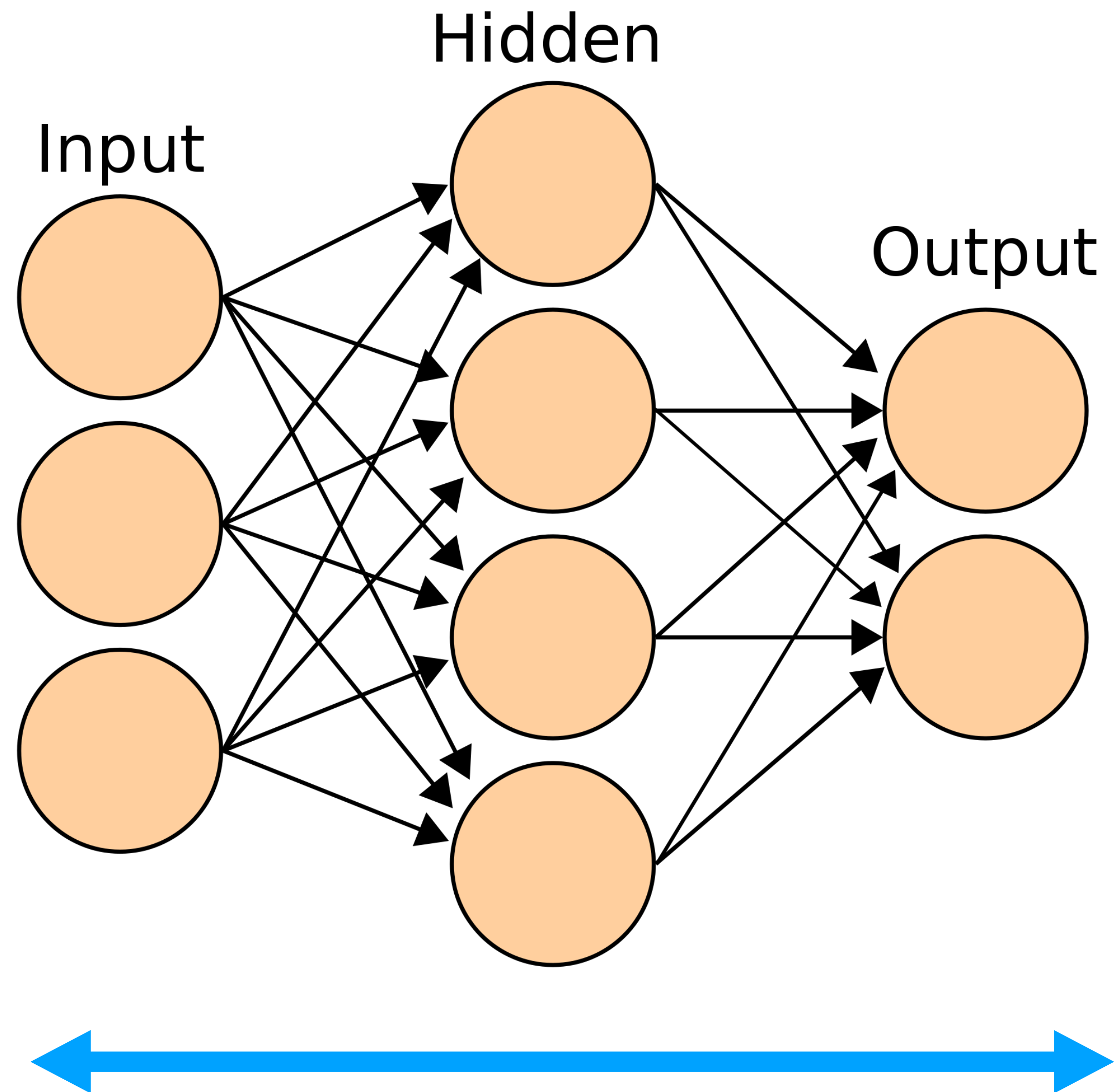
Jessica Yung

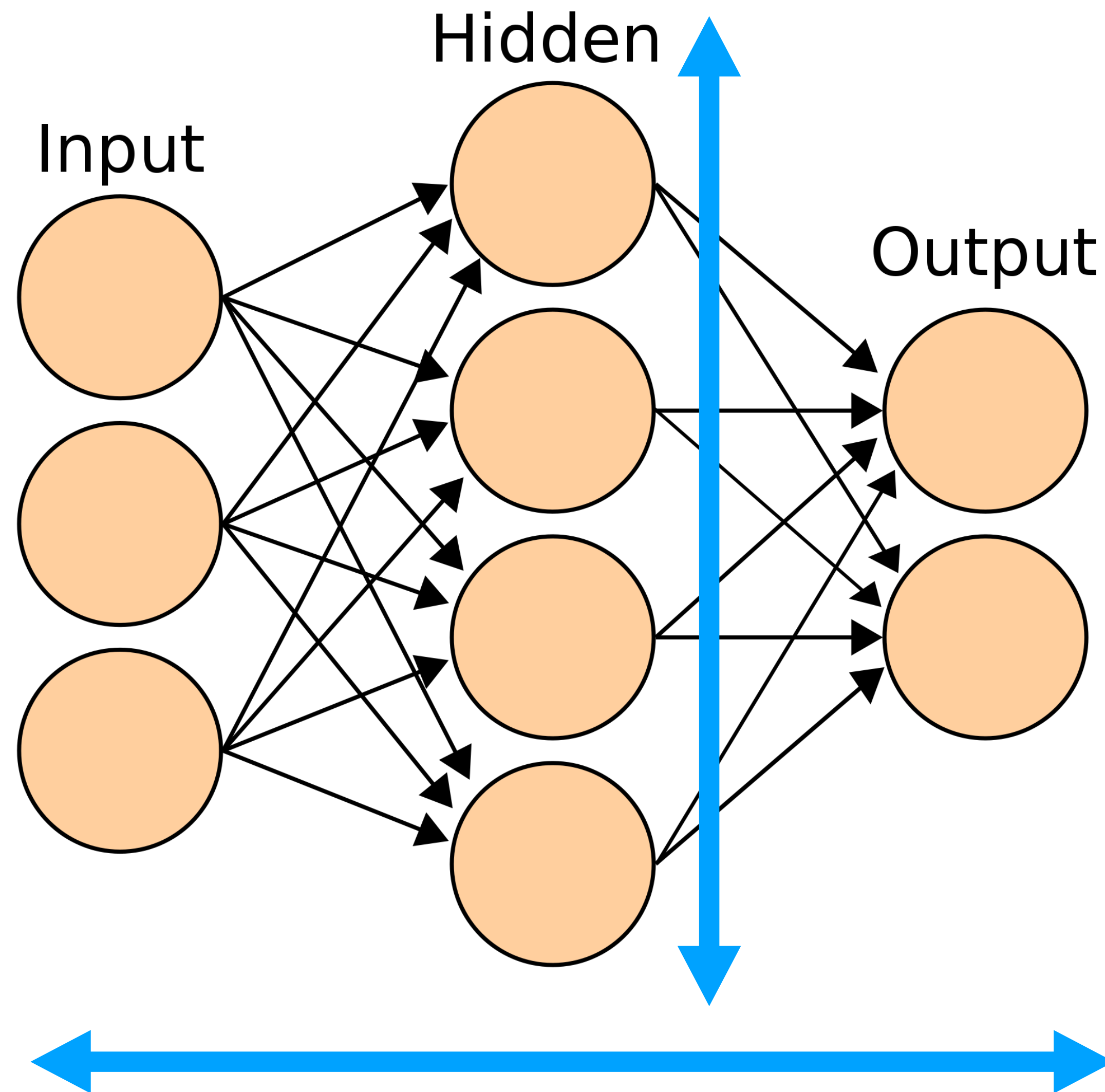


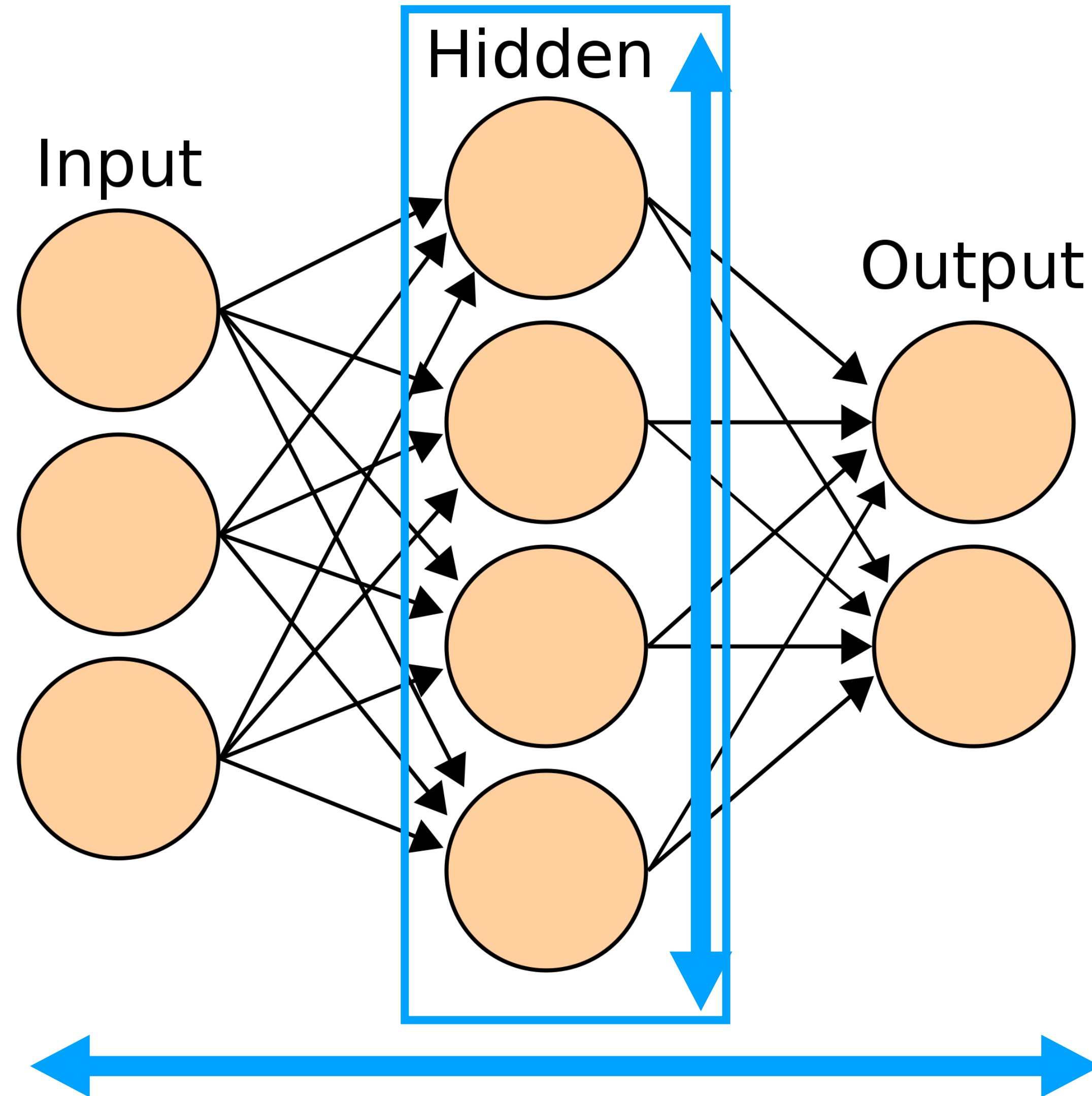


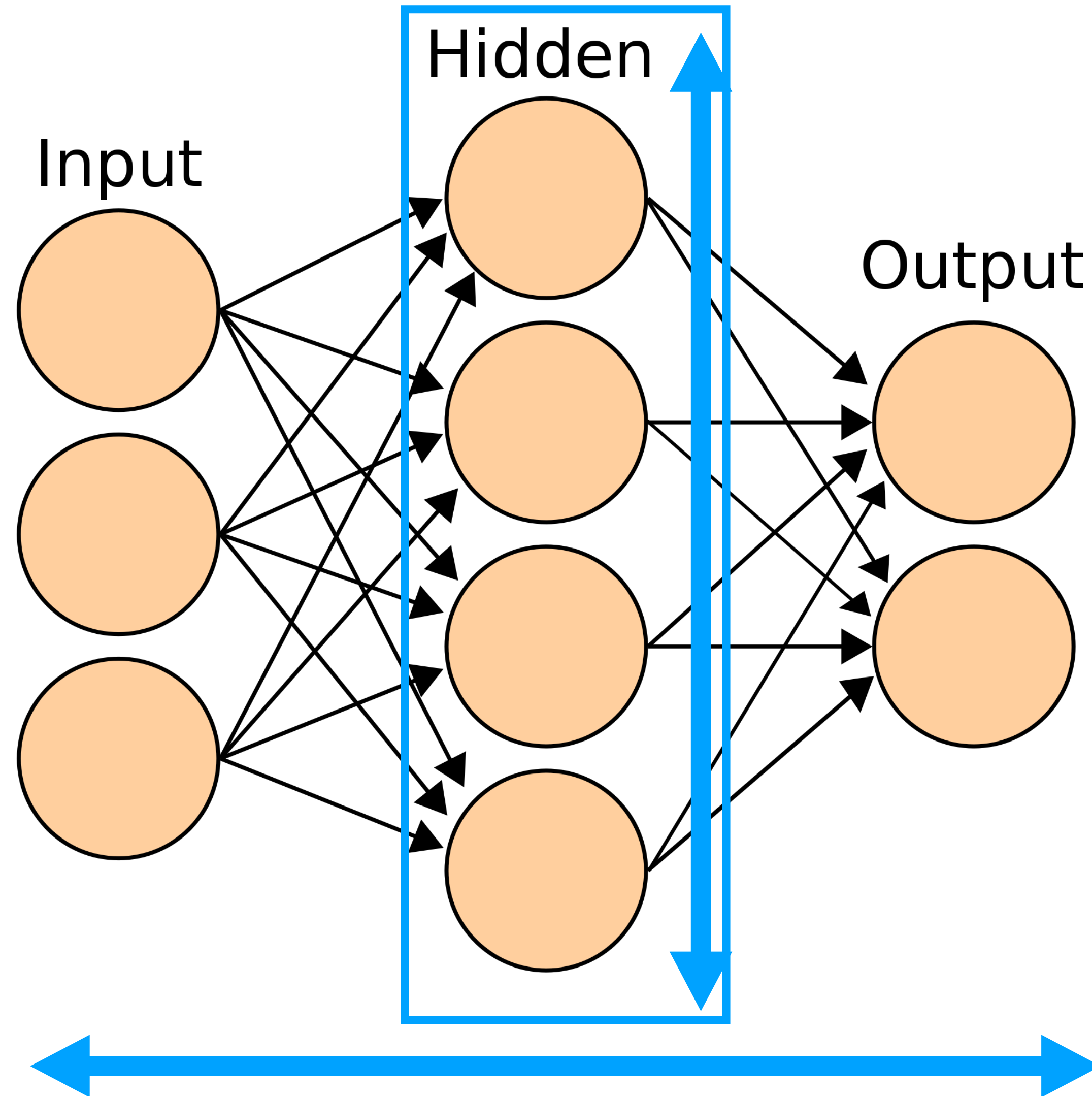




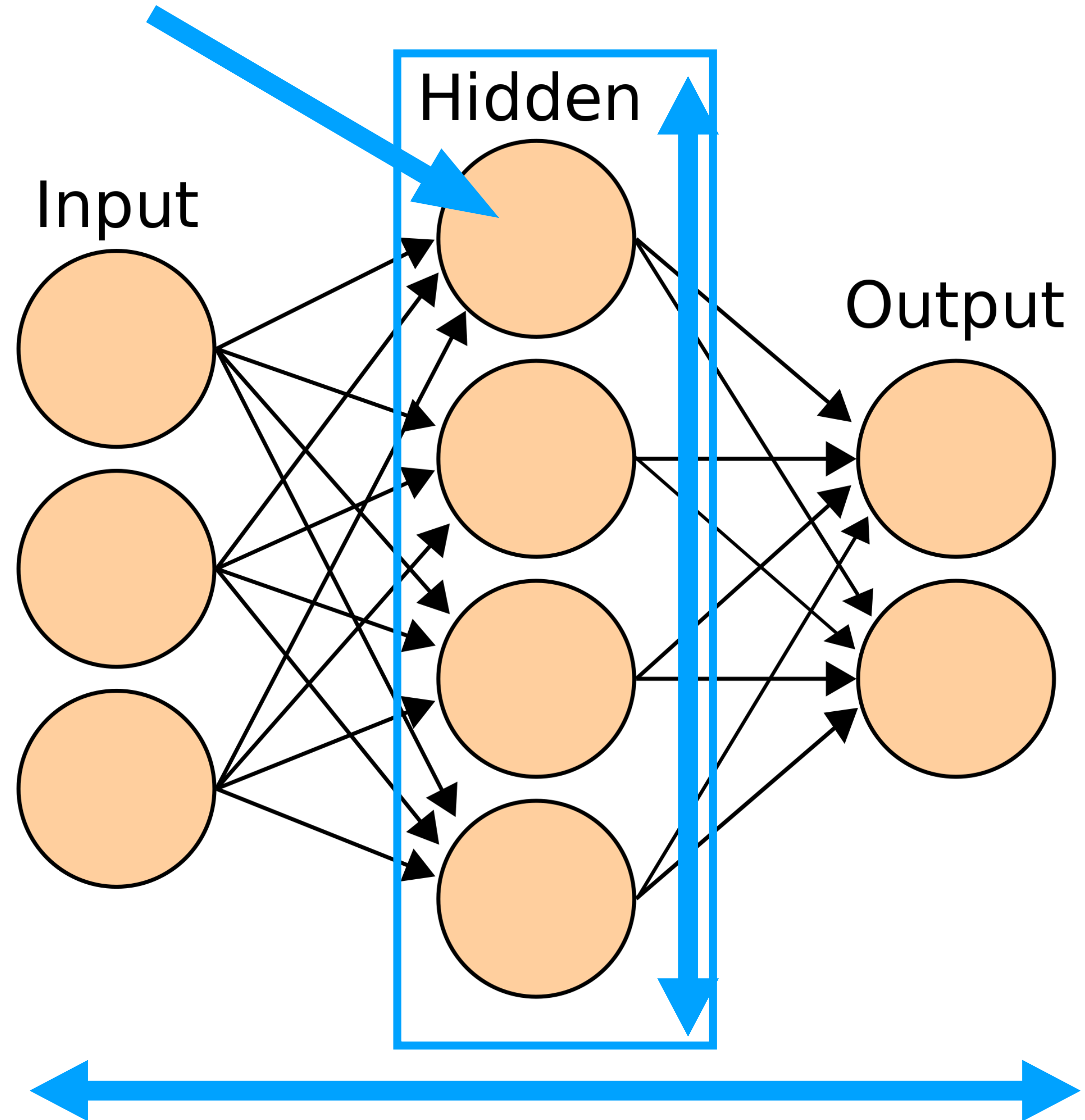




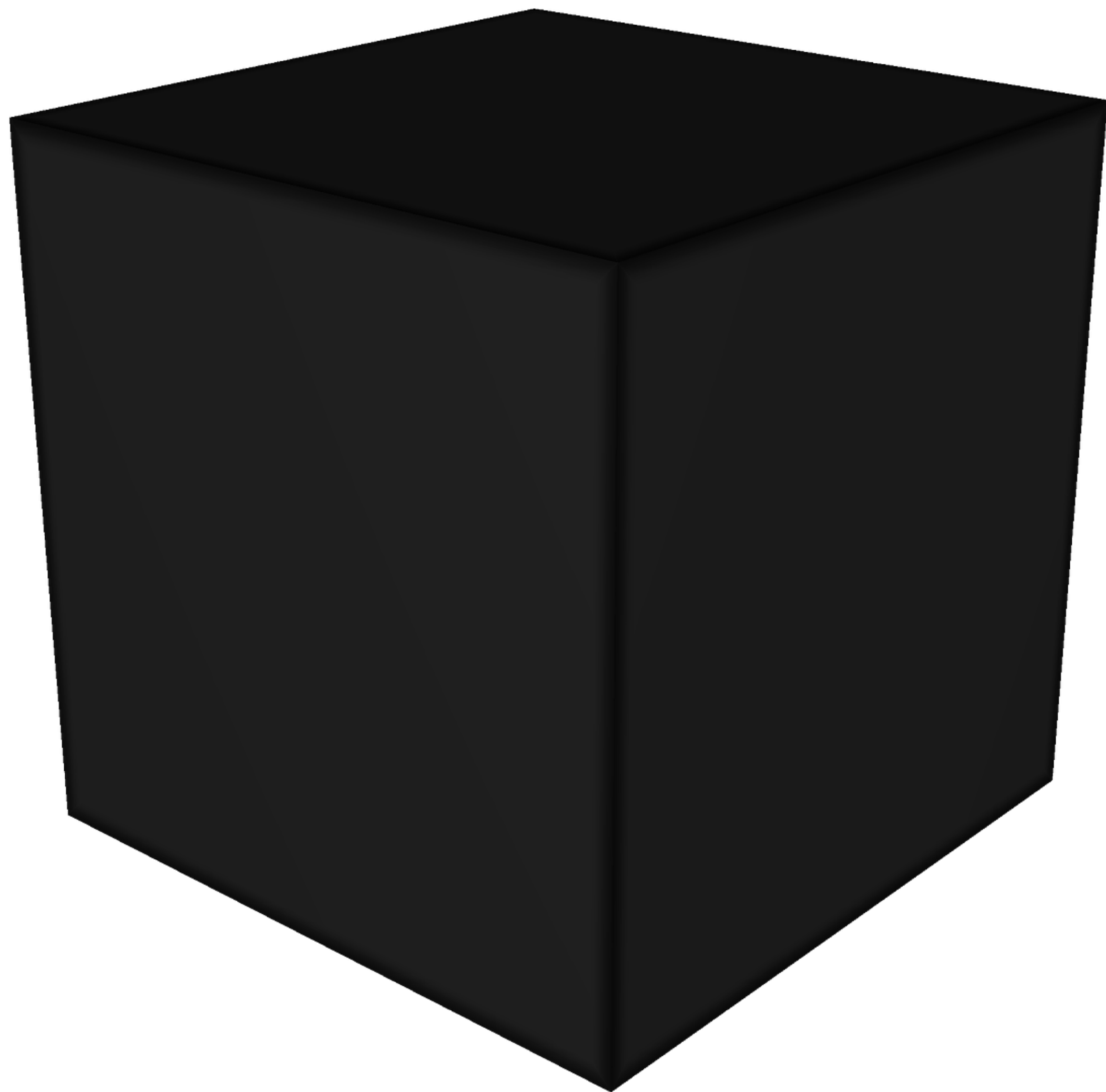


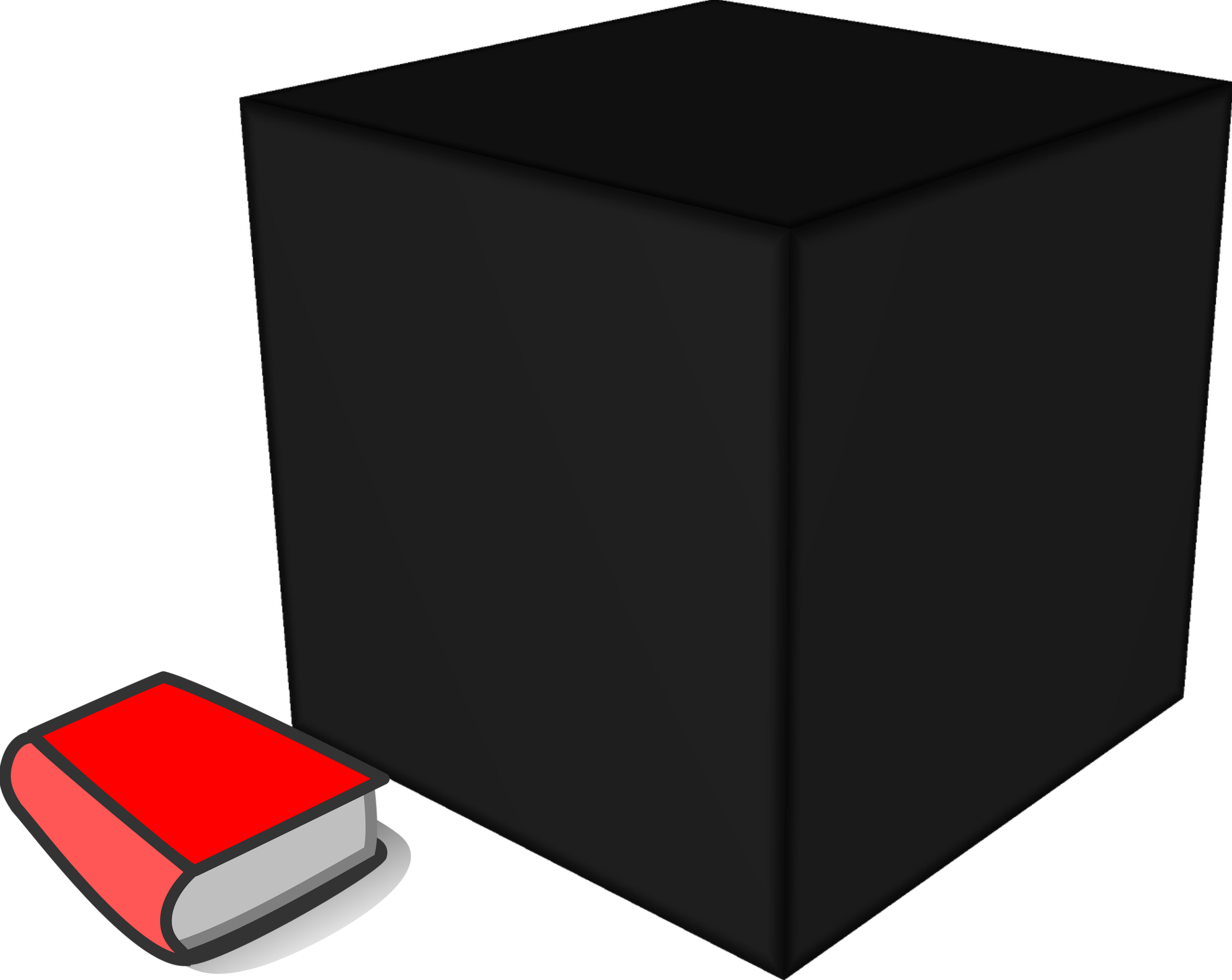


learning rate



learning rate



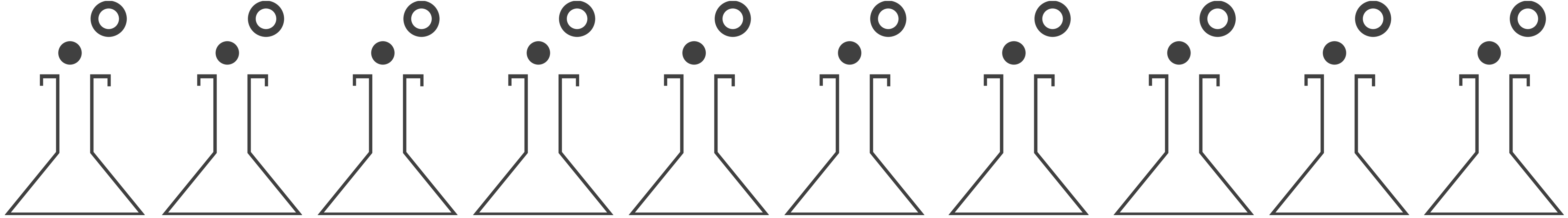
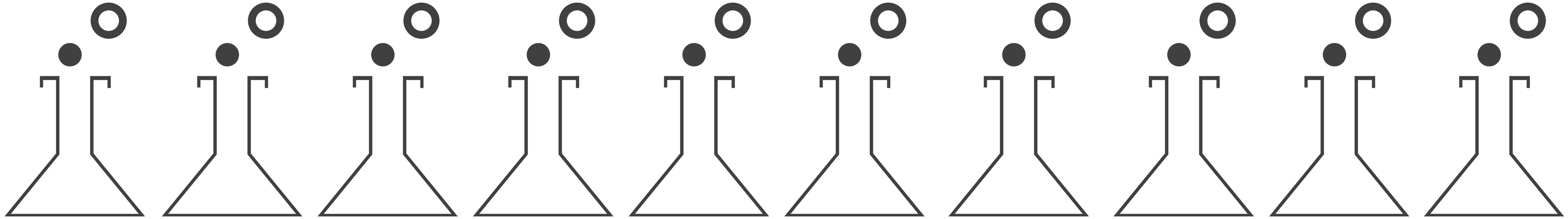
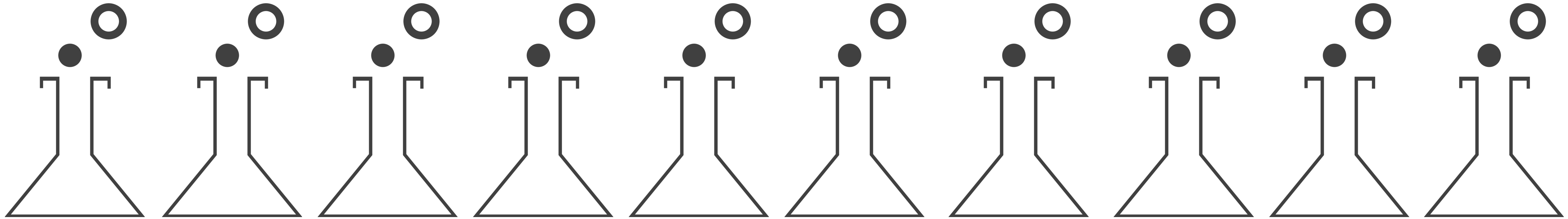
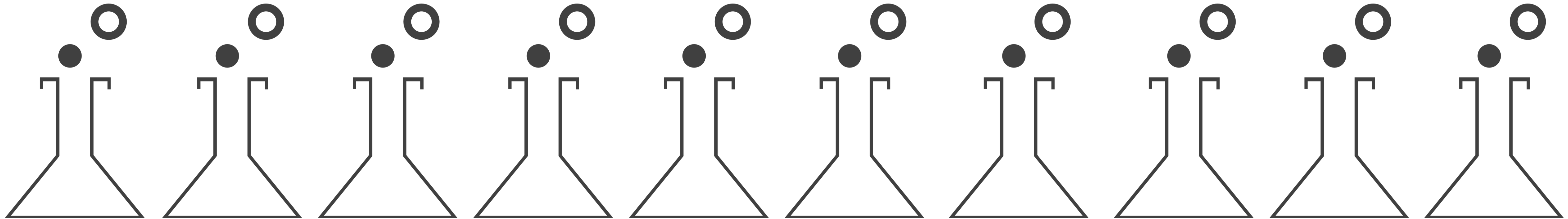


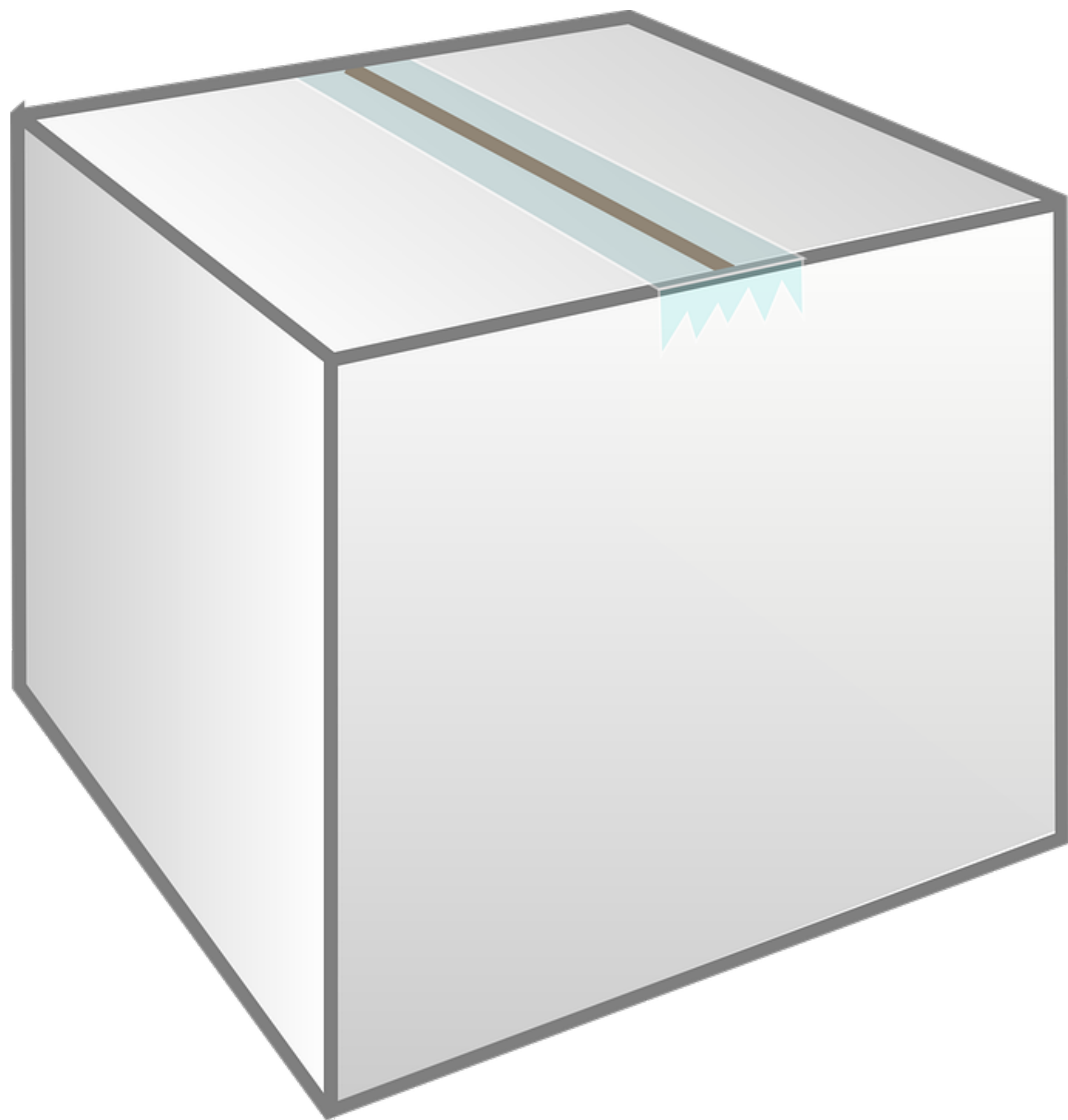
```
def debug ( ) :
```

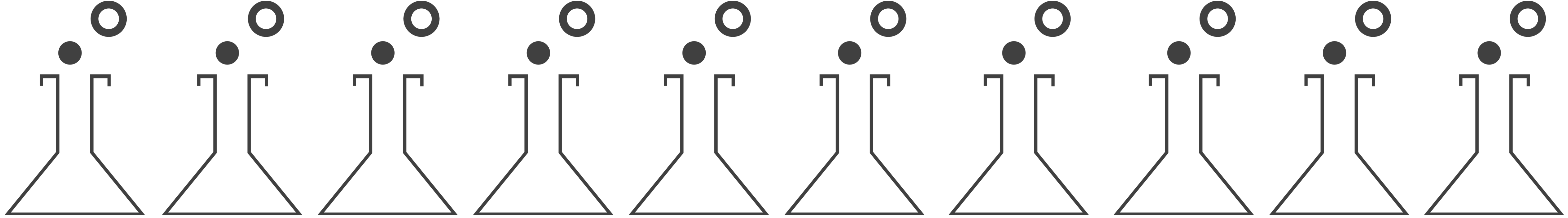
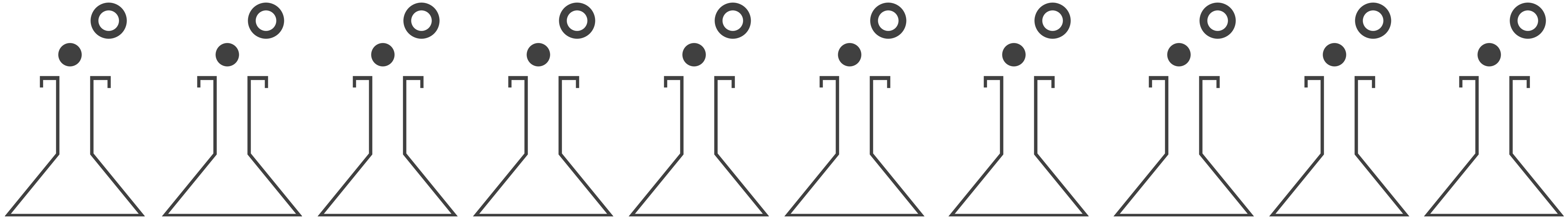
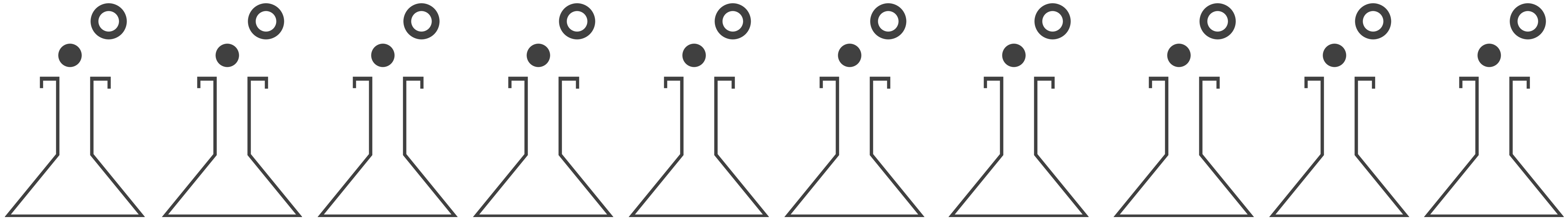
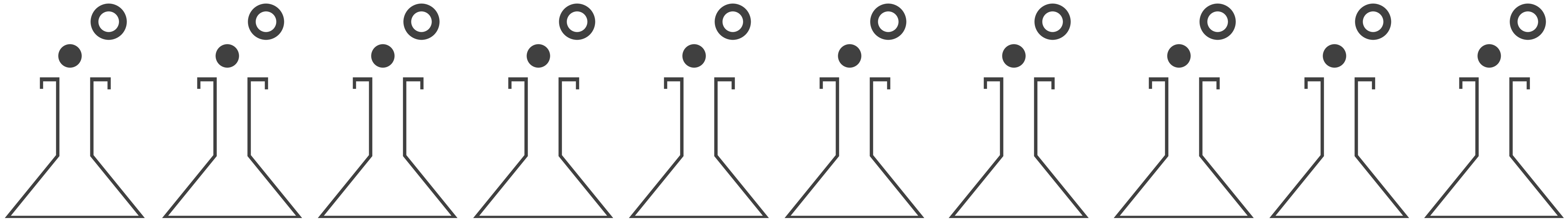
```
    return 0 + "no"
```

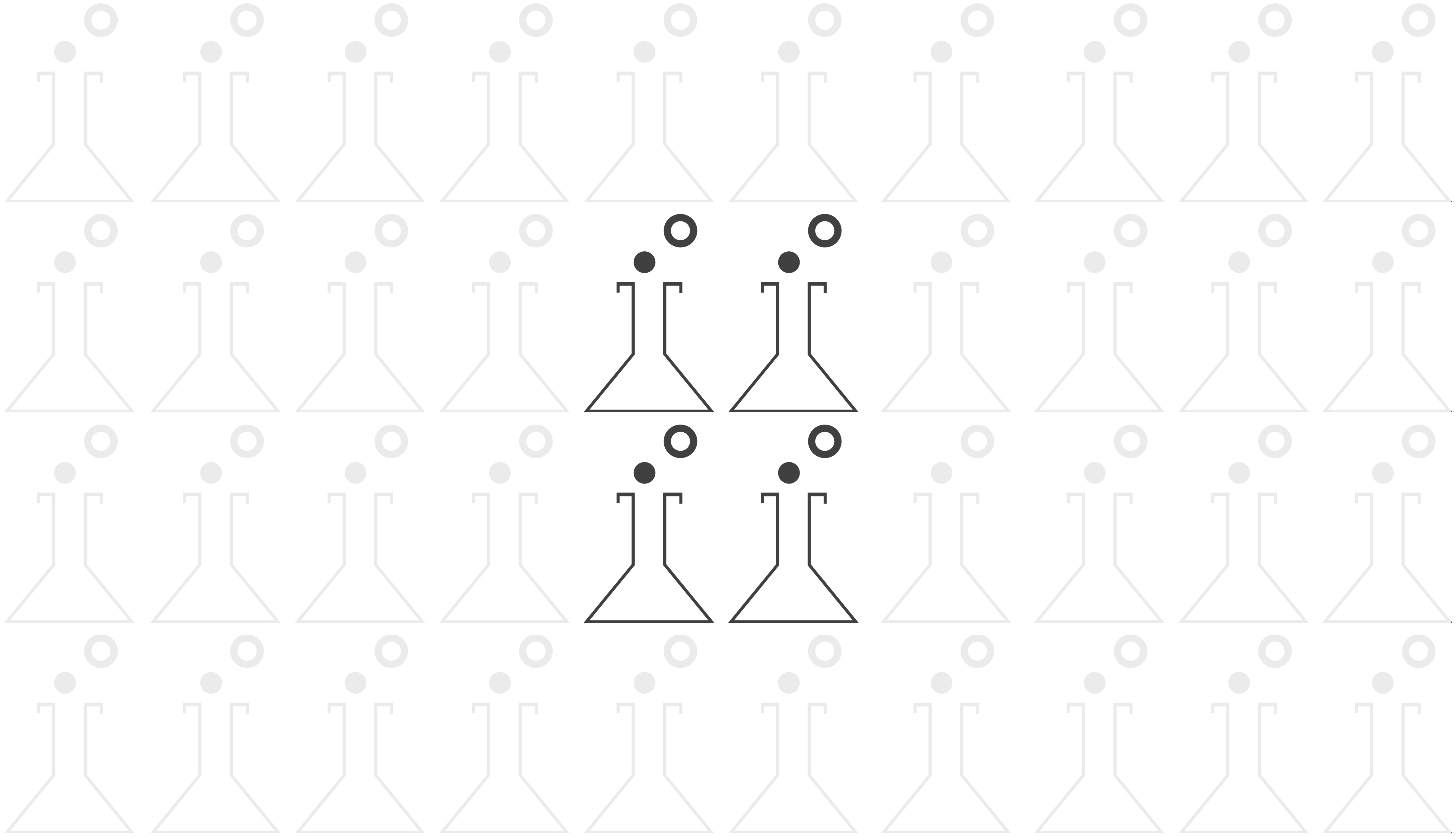
567 563m8dep

j61mj? _ g " ? @ "









Why deep learning?

(Super) human performance in

(Super) human performance in



(Super) human performance in



あ → A

(Super) human performance in



あ → A



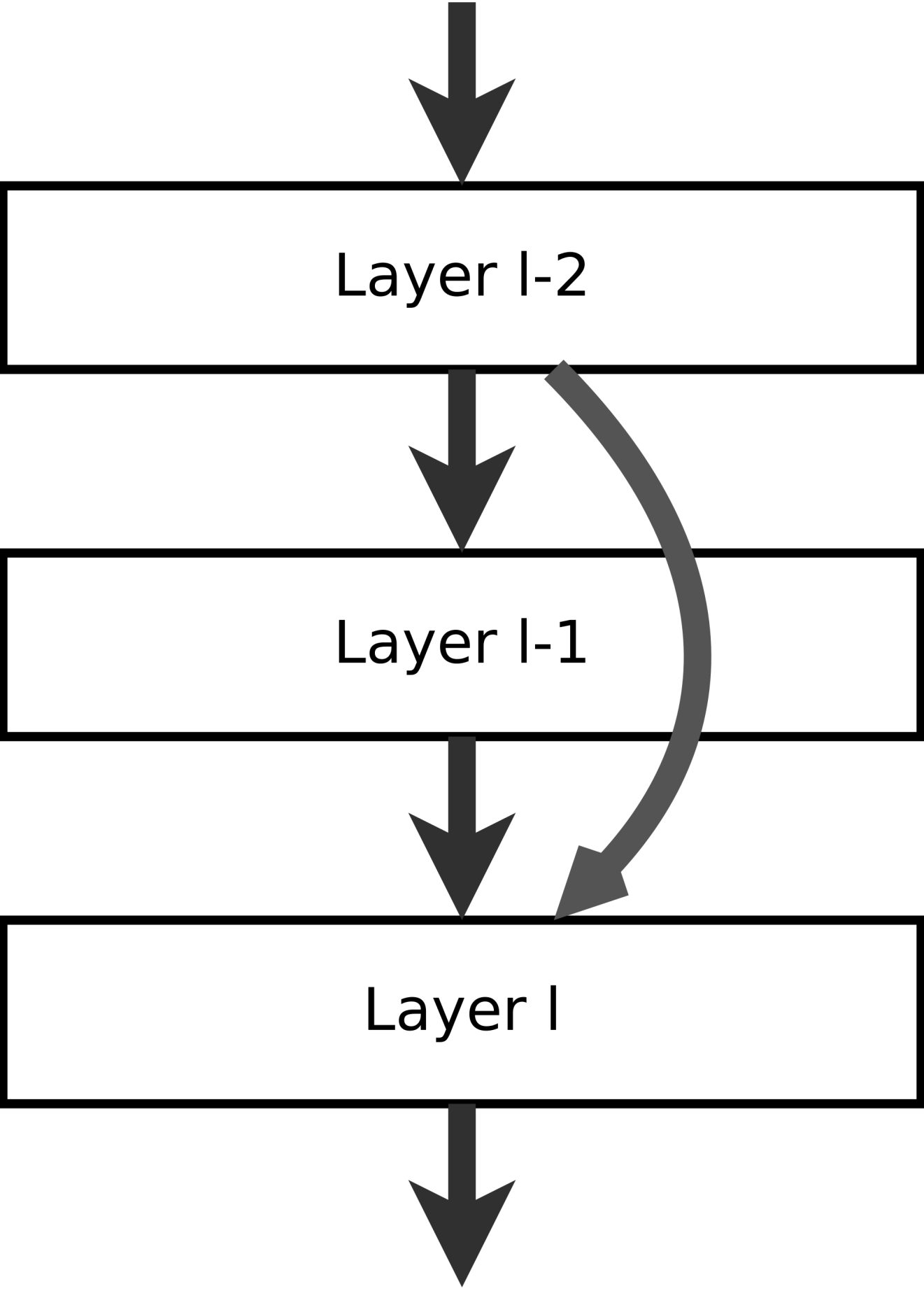
Outline

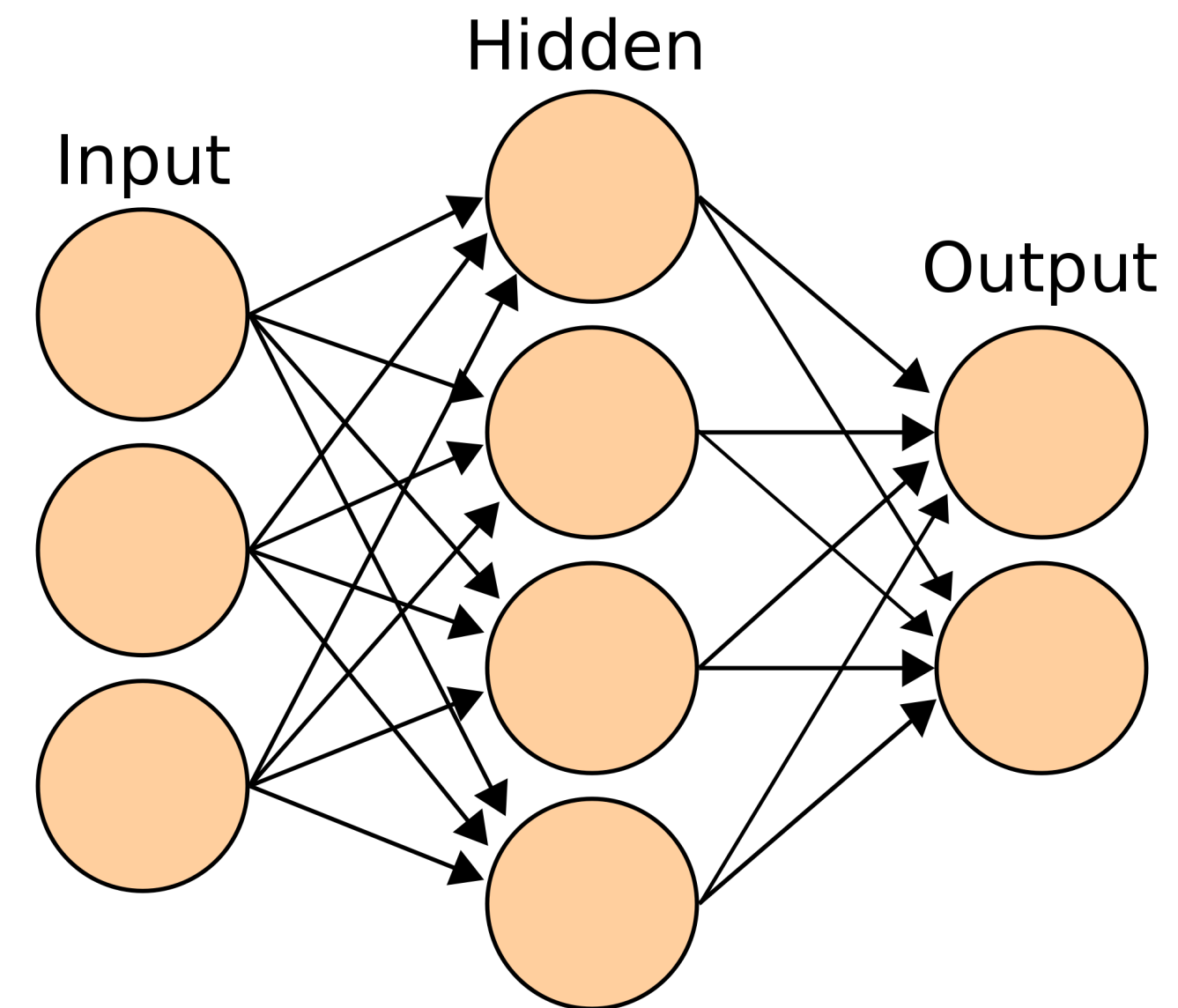
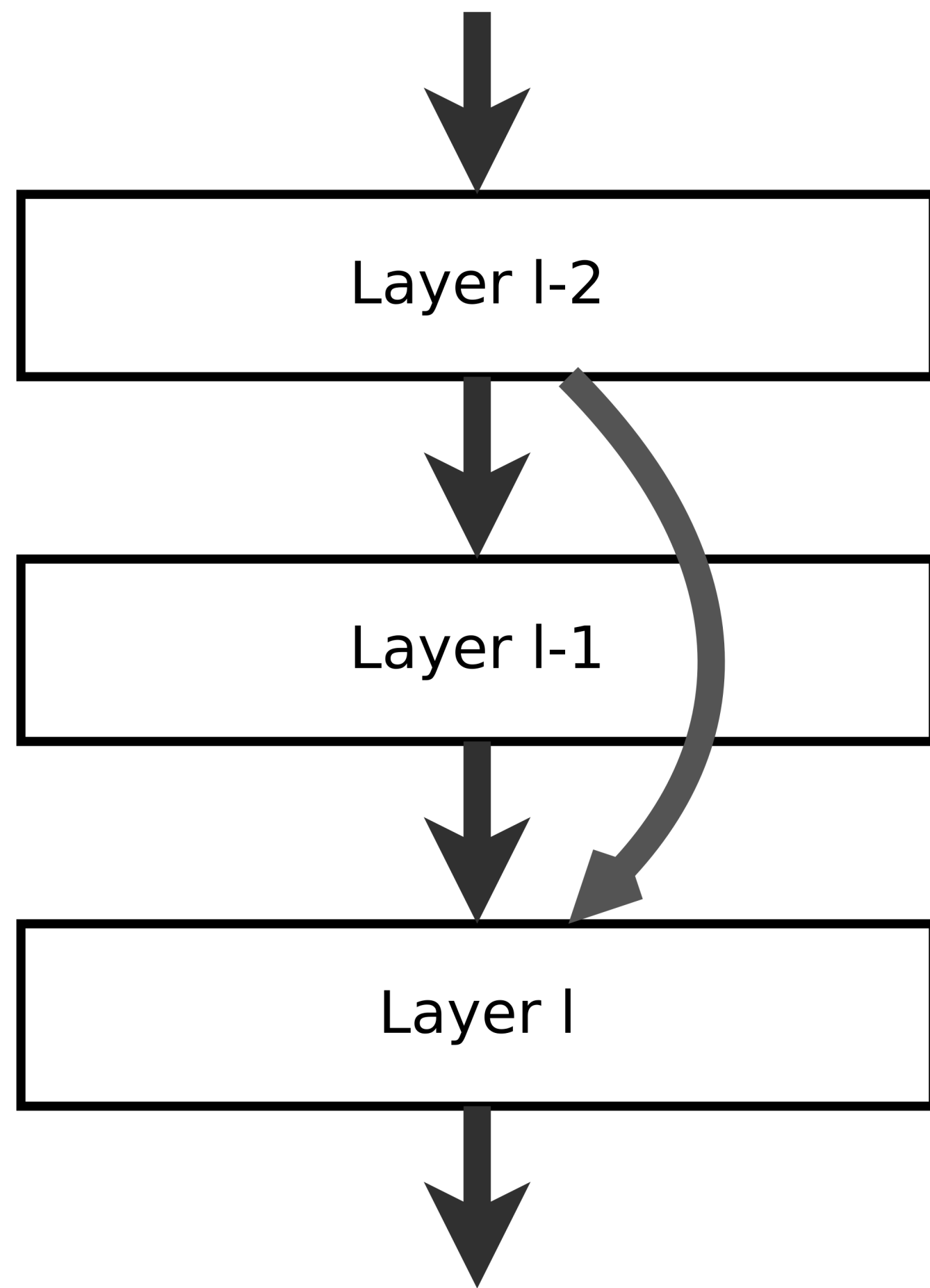
1. Single layer and putting layers together
2. How to train models
3. Deeper is better
4. Practical tips

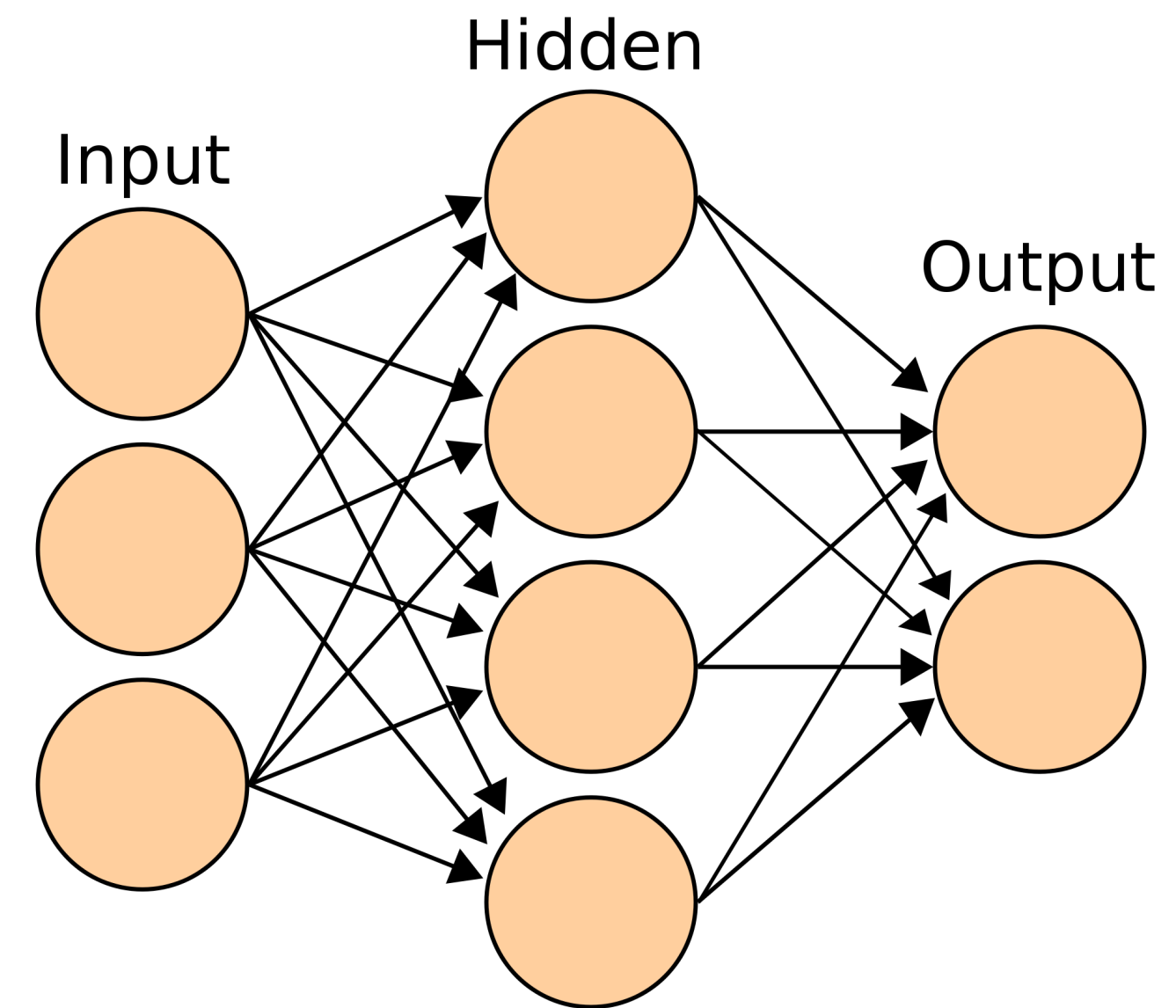
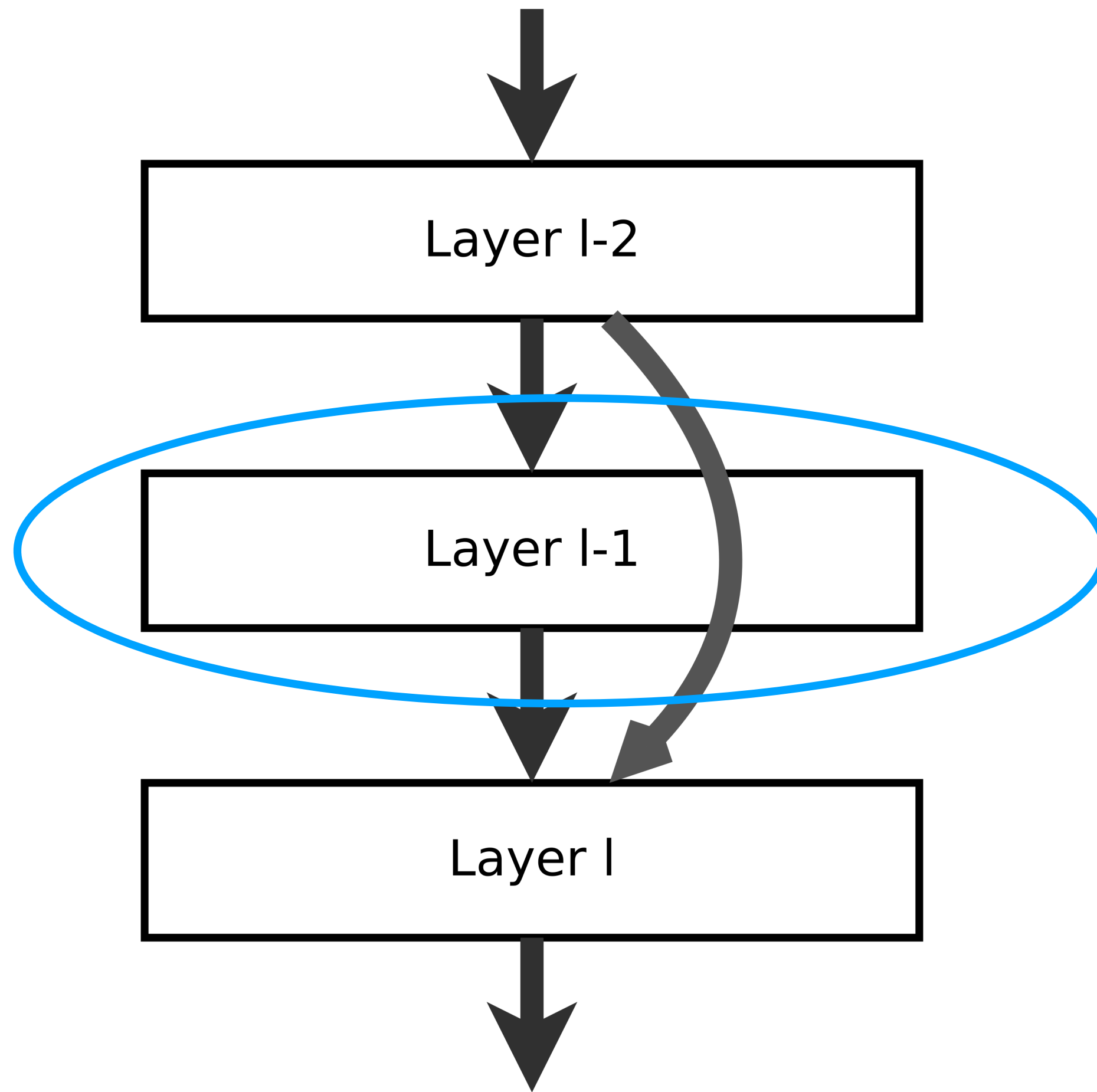
**Understanding is not
something you can
Google.**

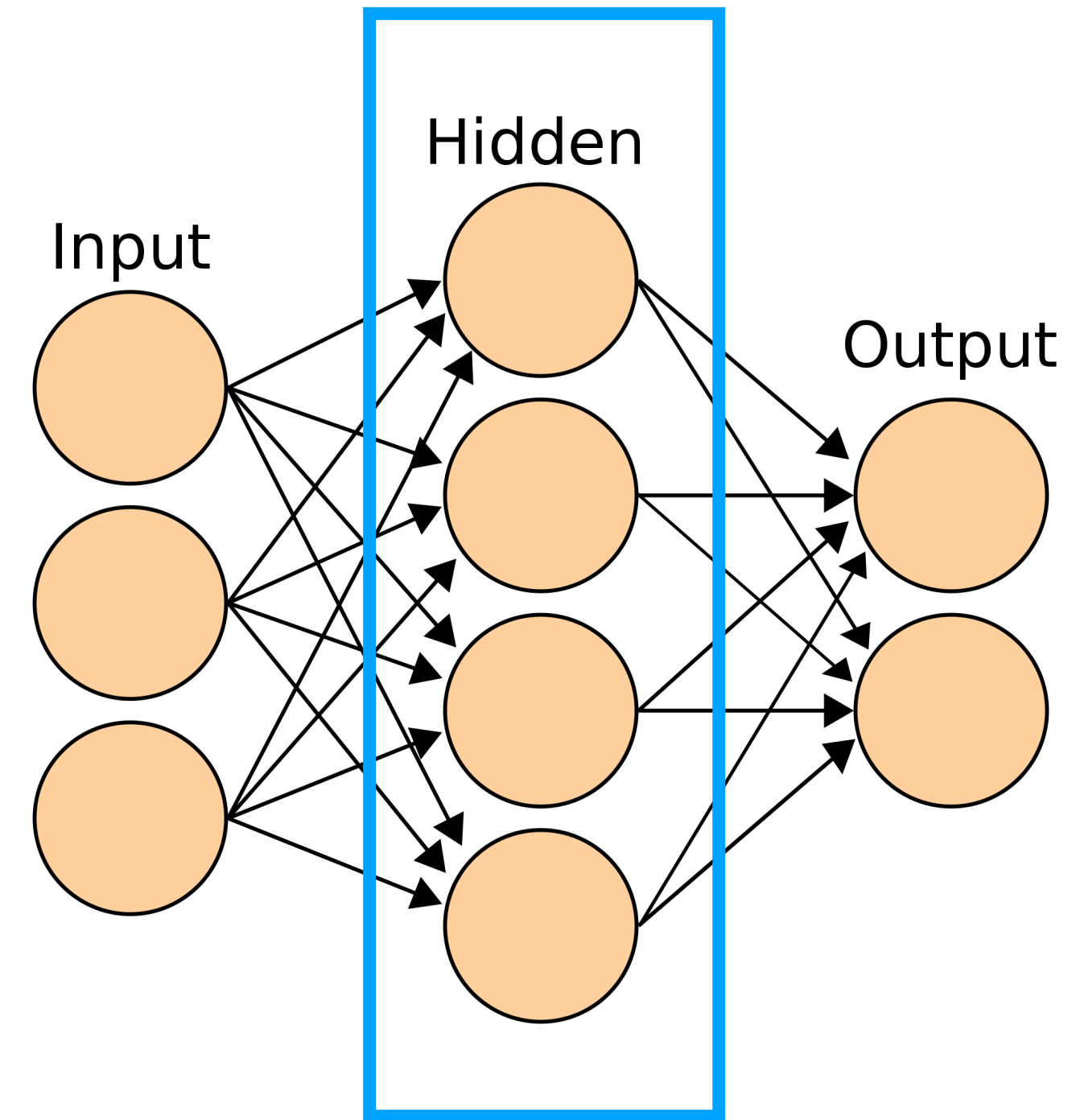
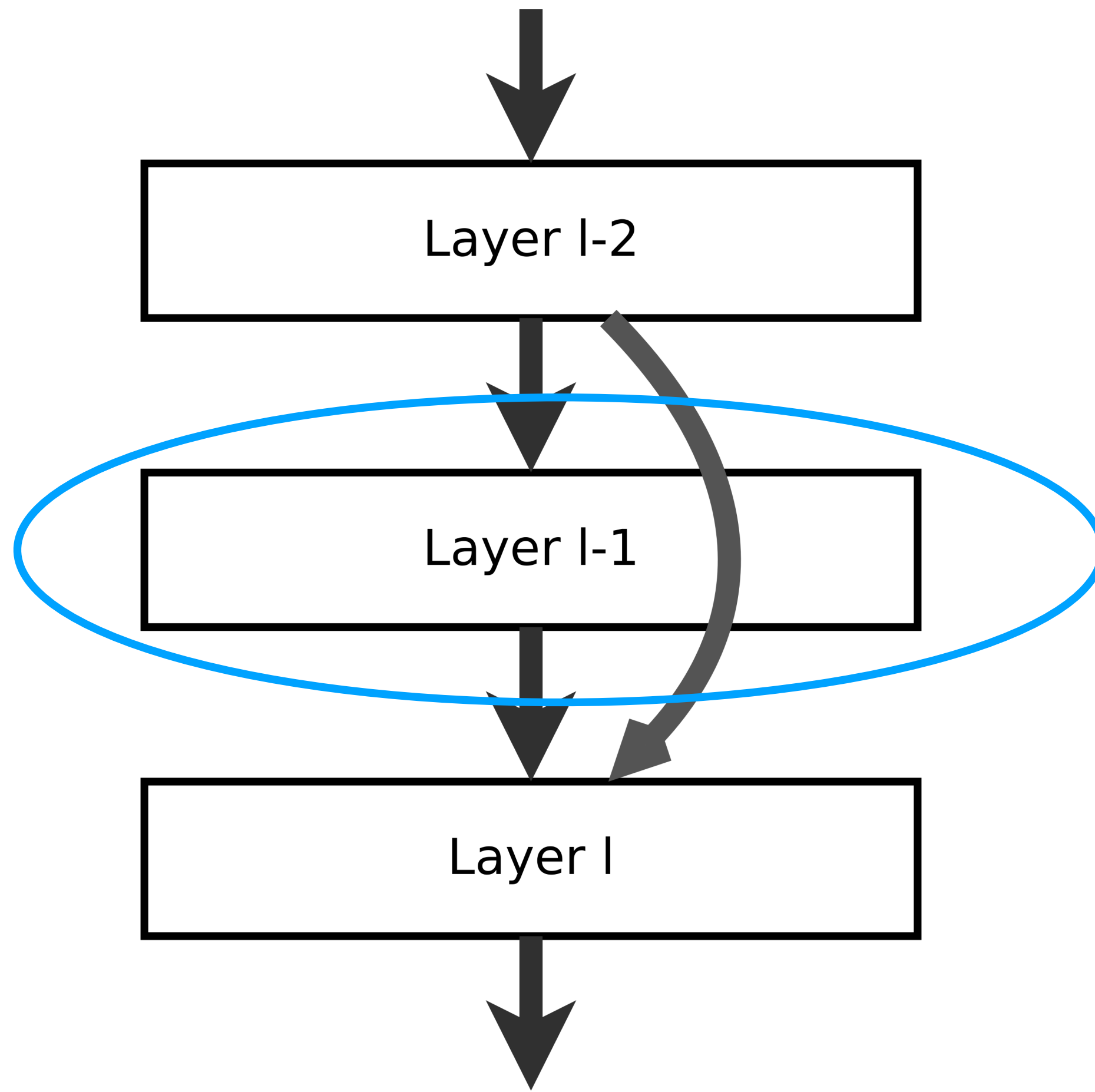
**Machine learning is
engineering-driven.**

I. Single Layer

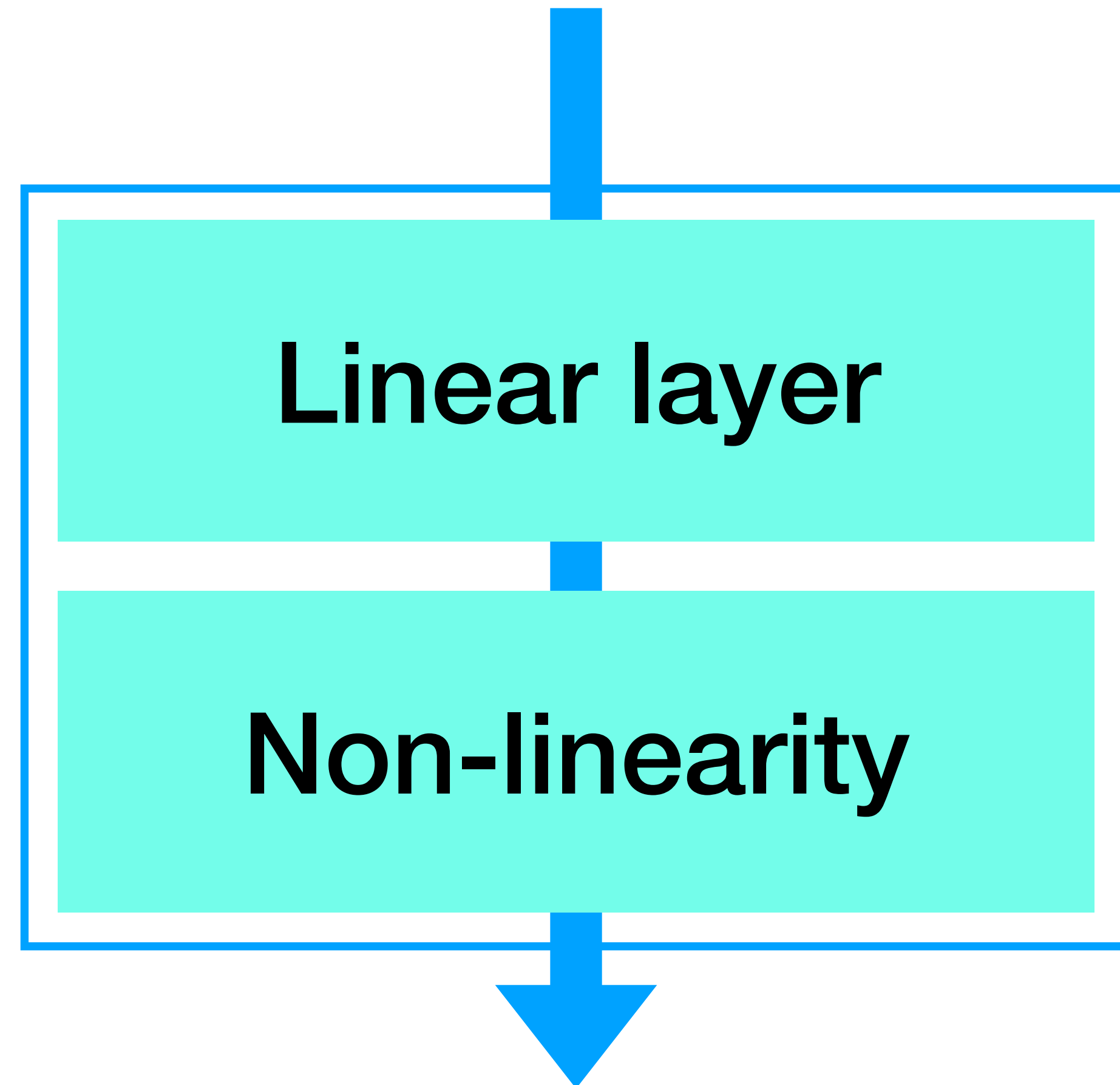






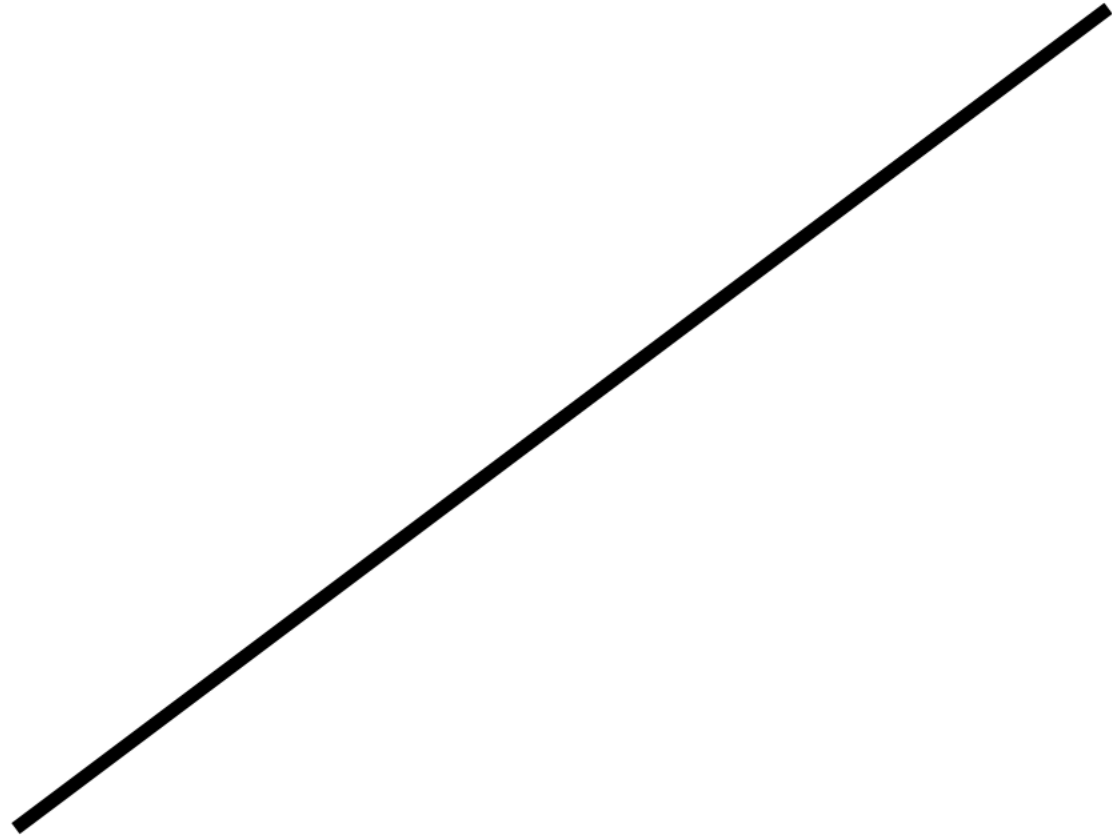


Fully connected layer

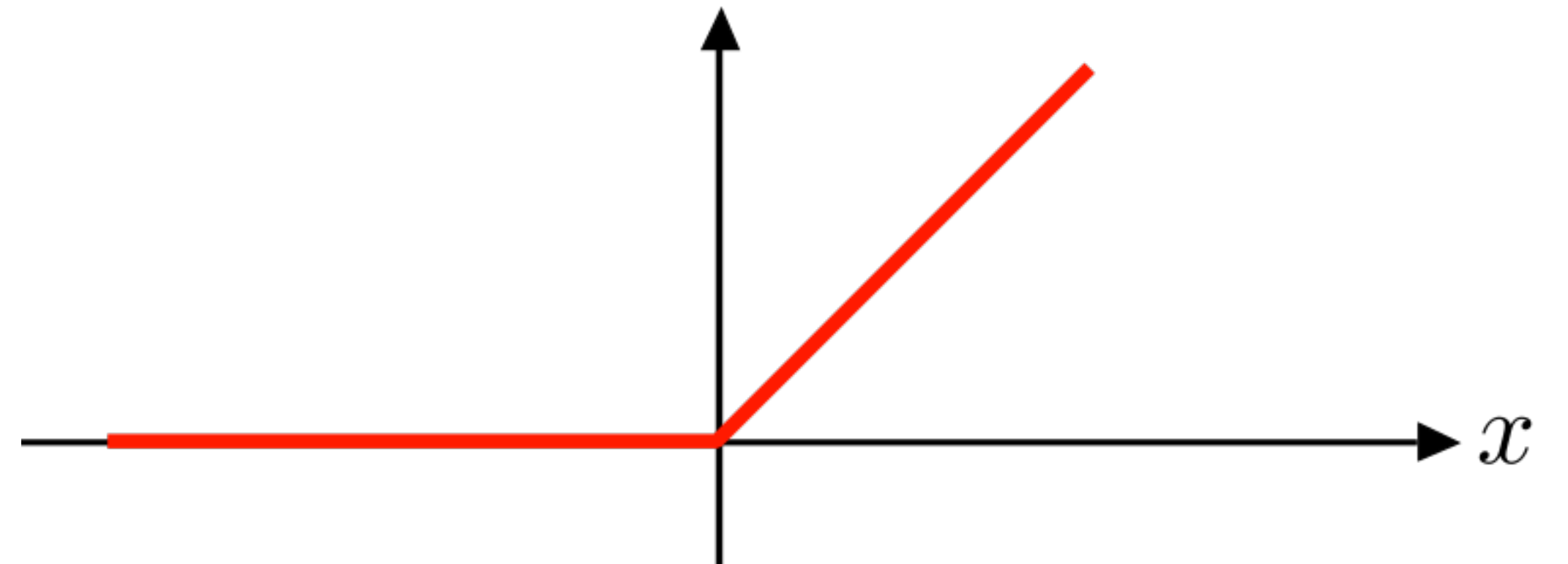


A single layer

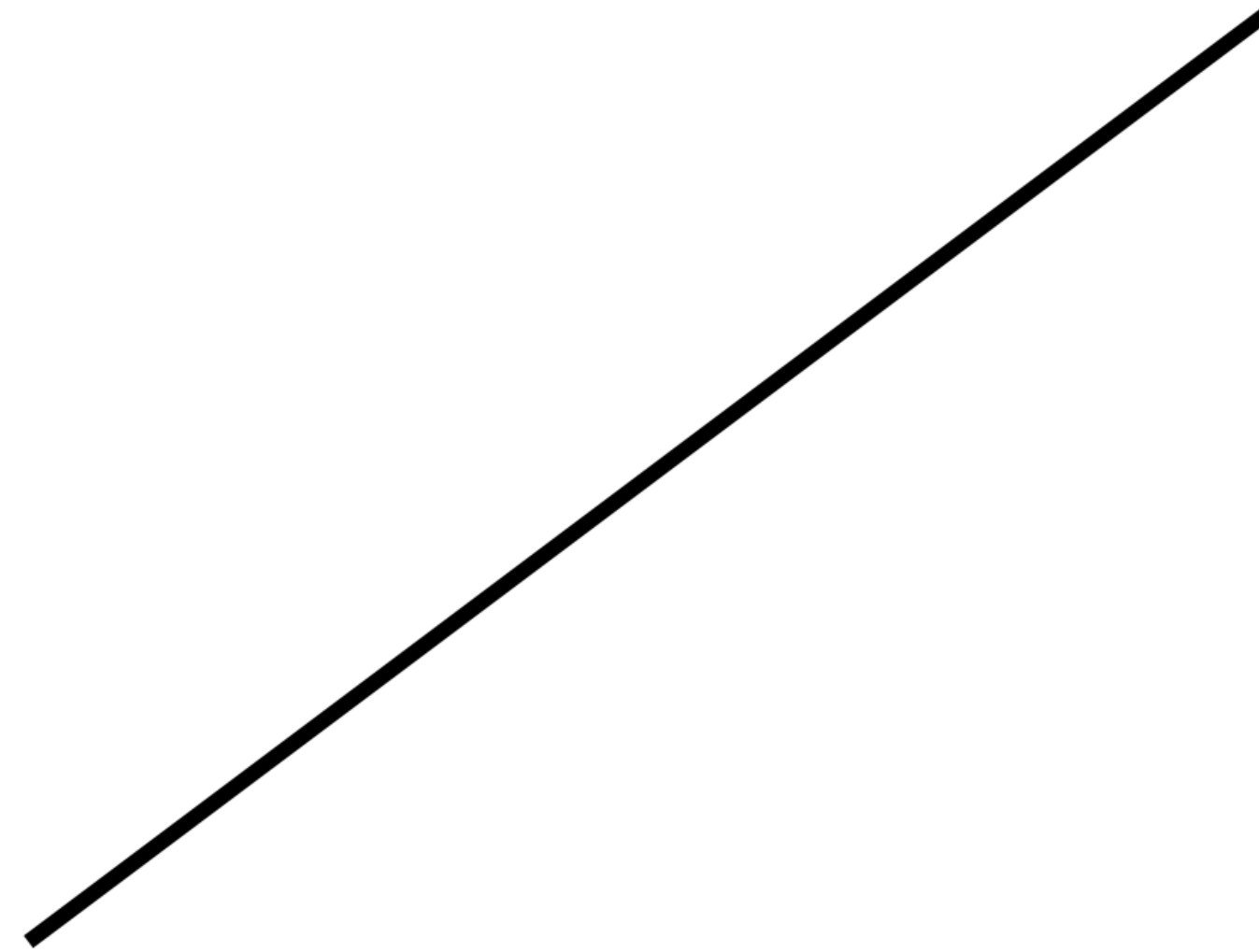
Linear layer



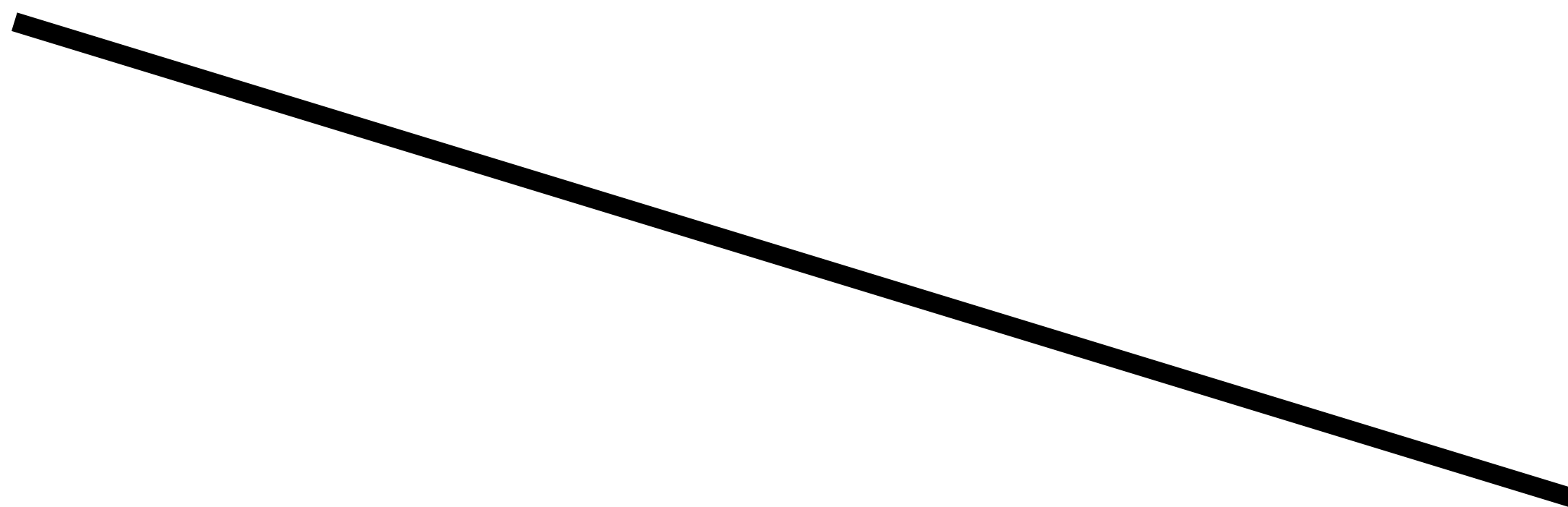
Non-linearity

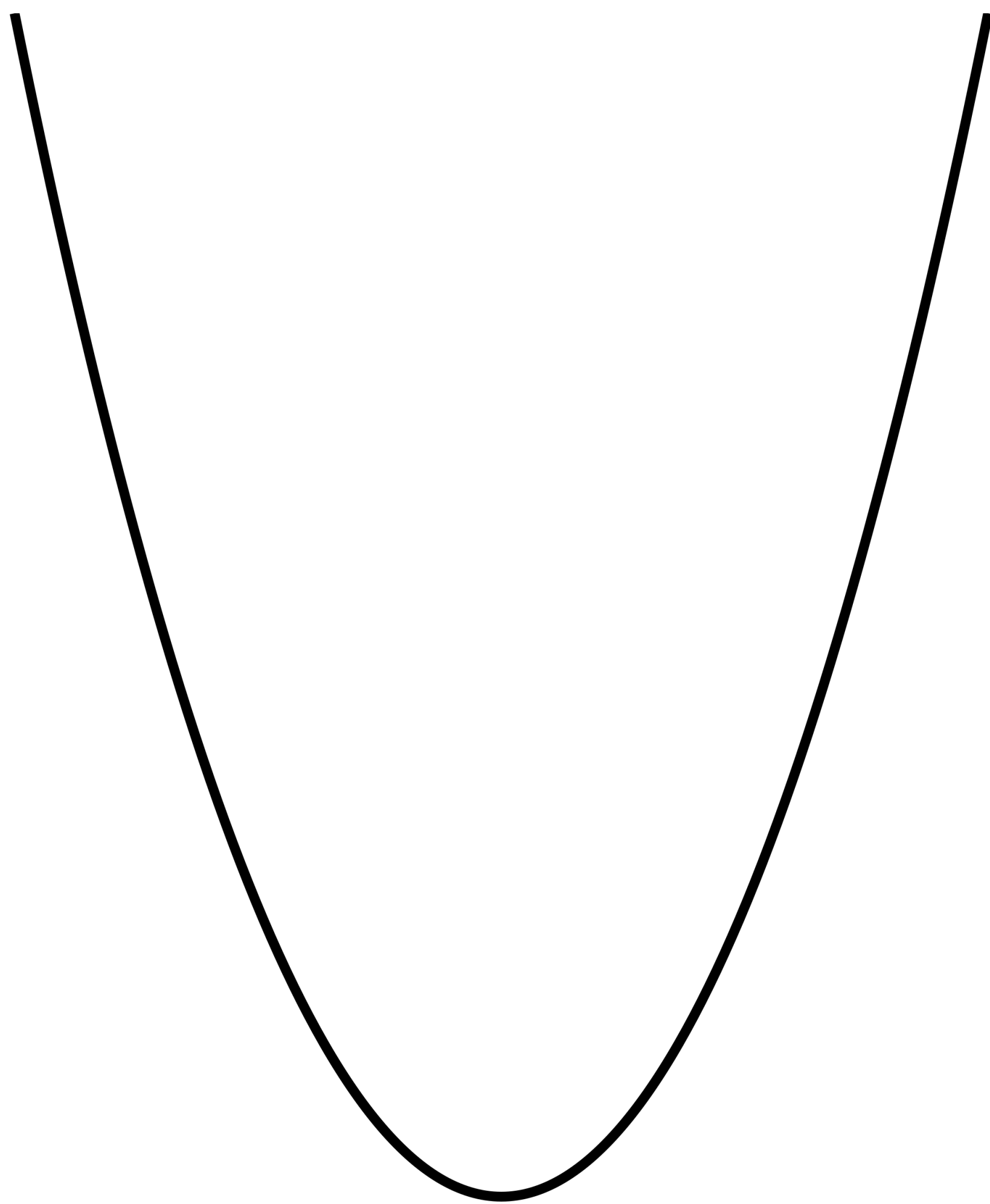


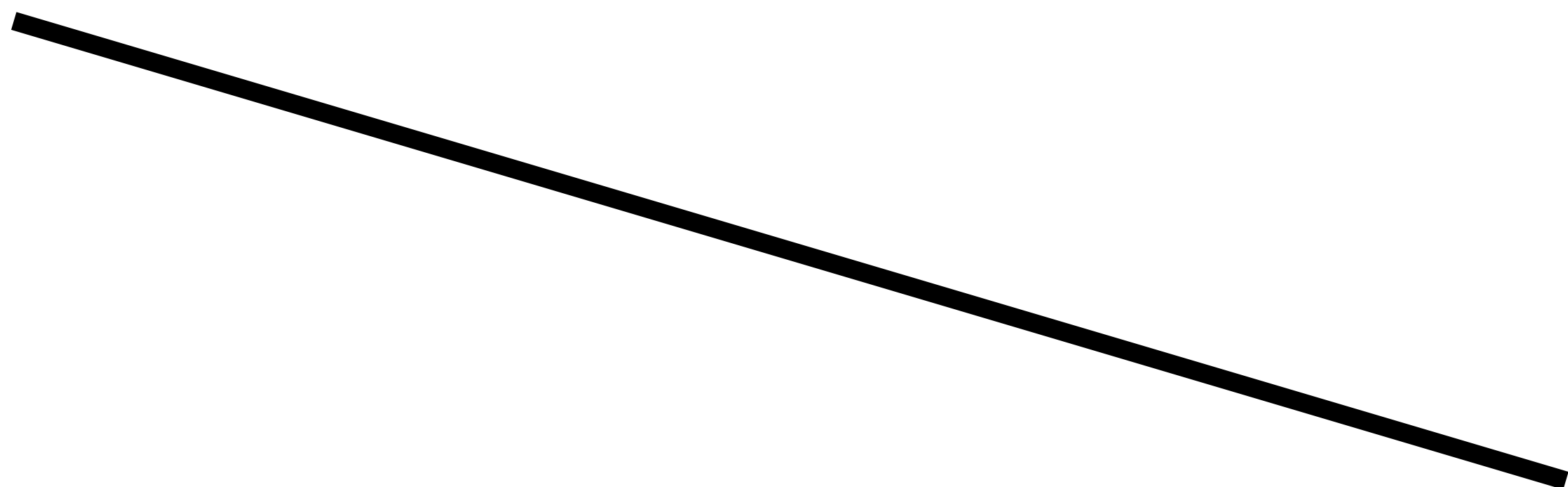
Linear layer

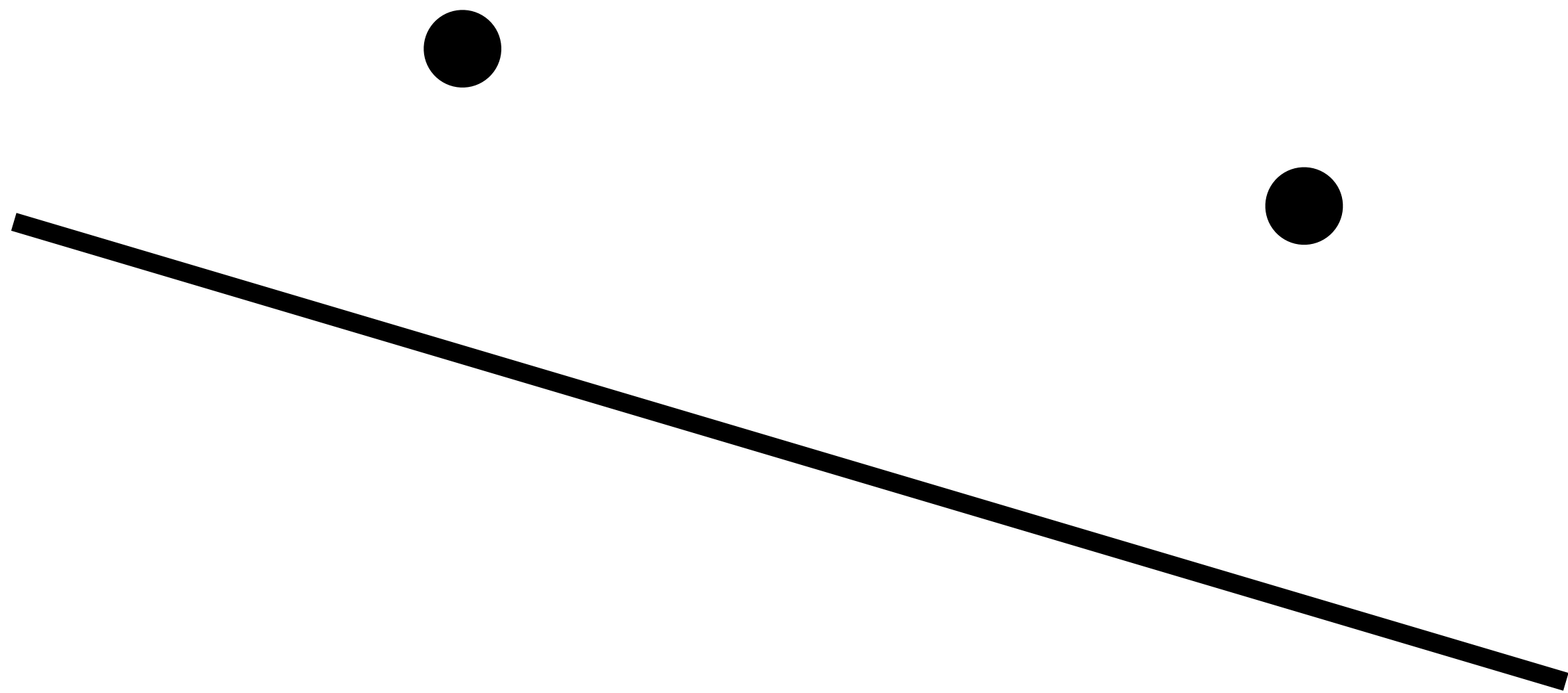


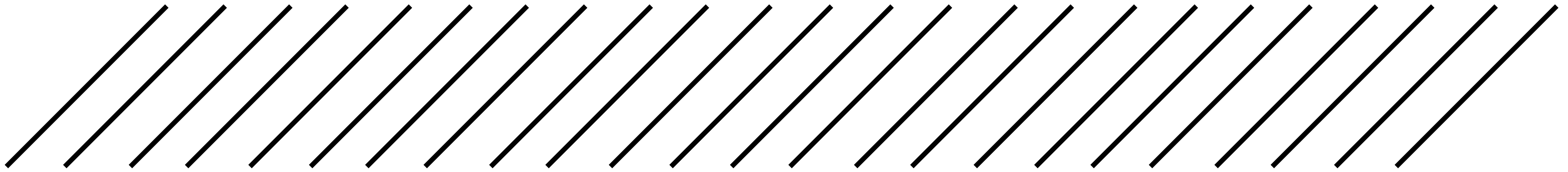
$$\mathbf{W}\mathbf{x} + \mathbf{b}$$

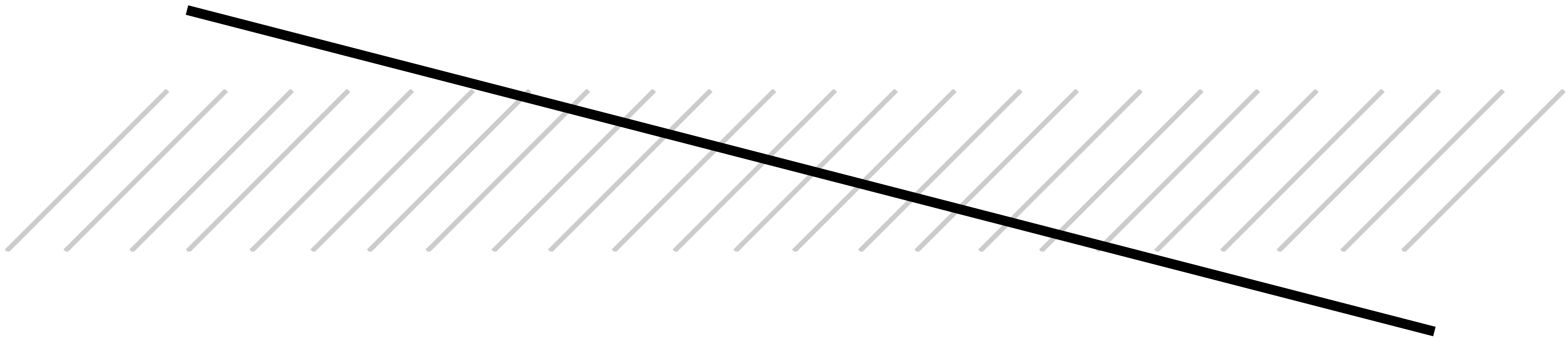


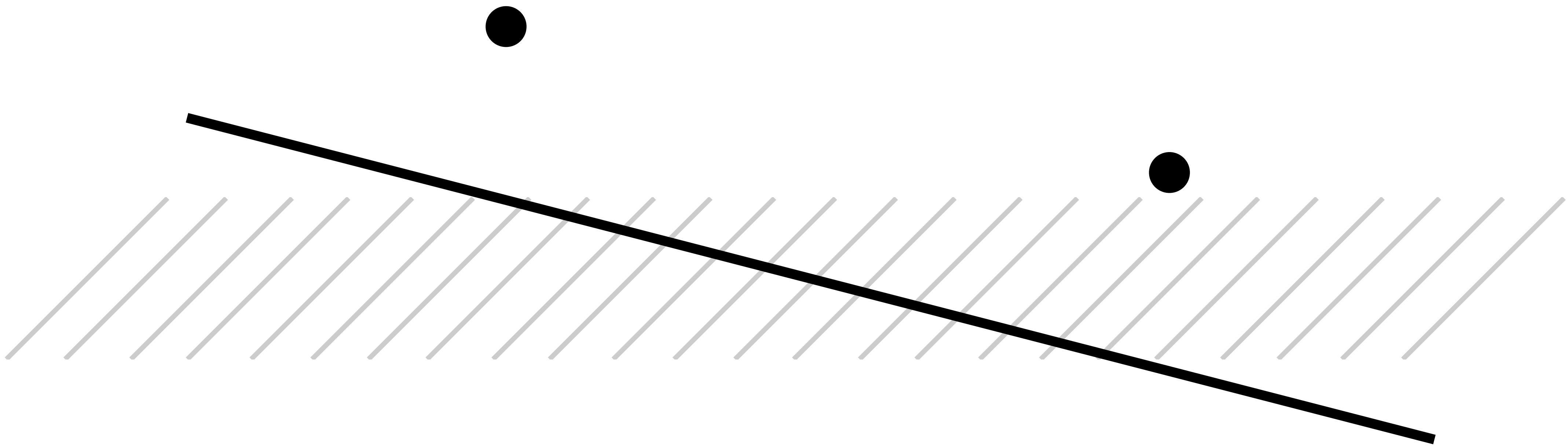


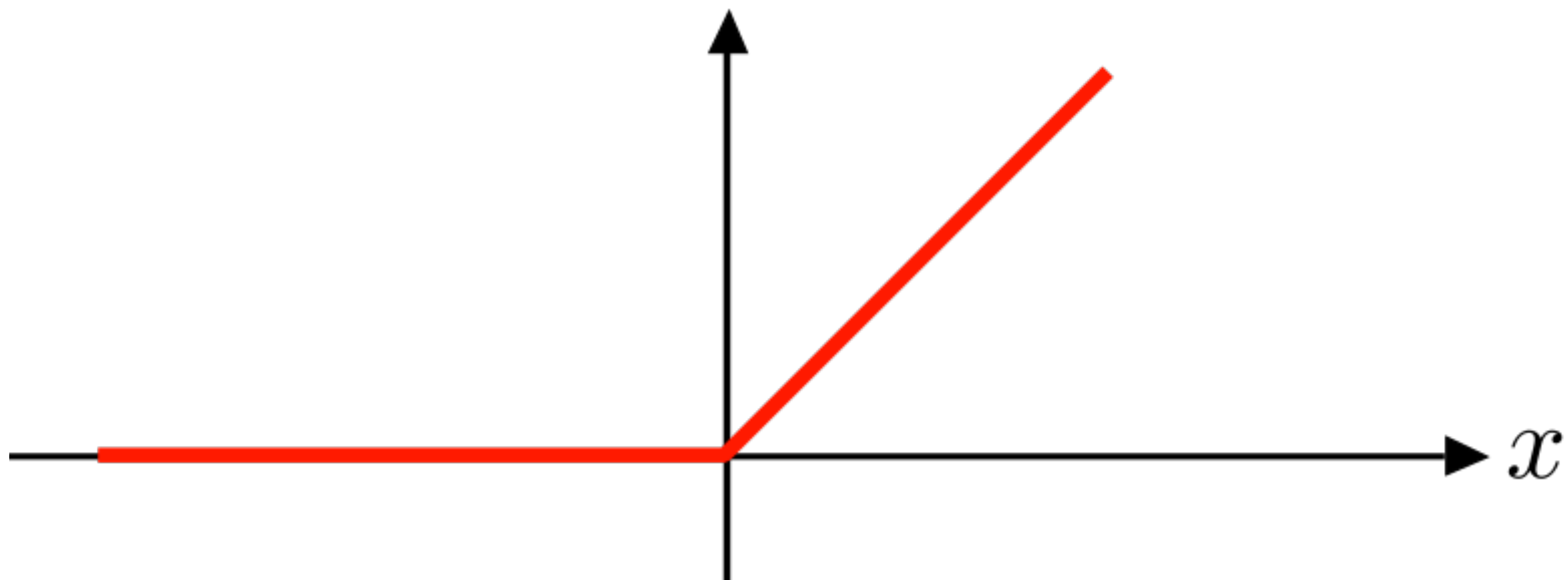


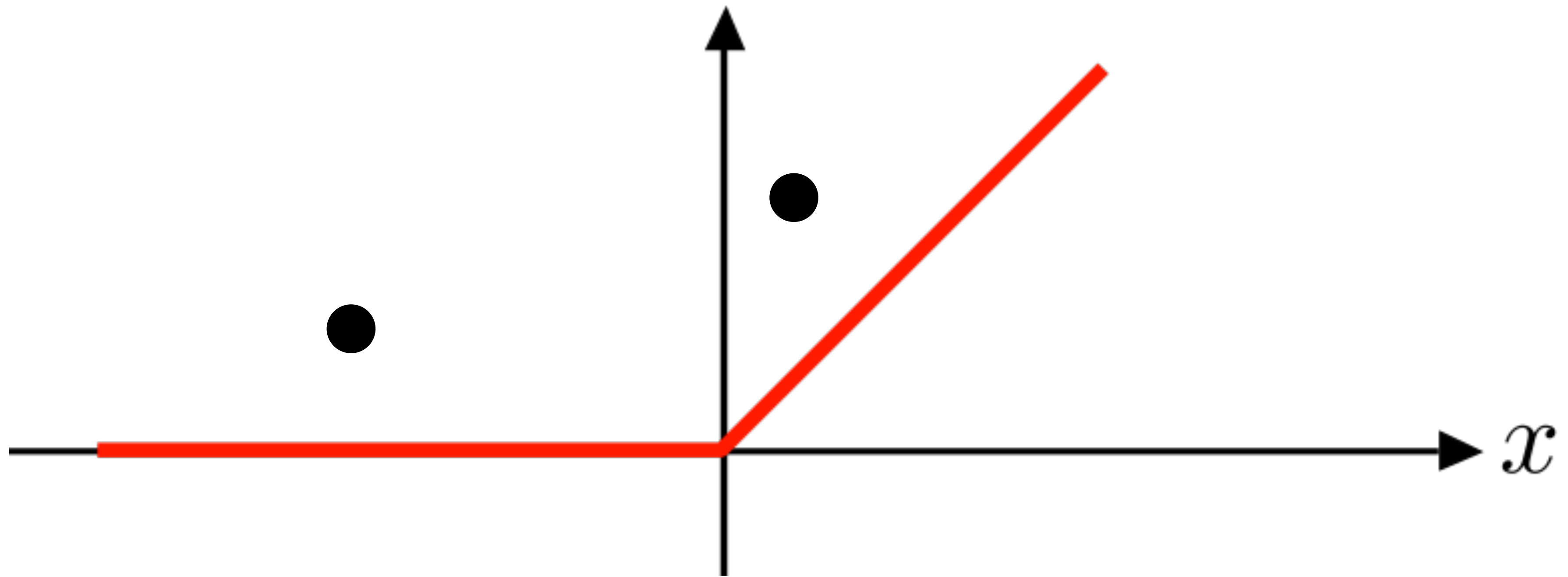






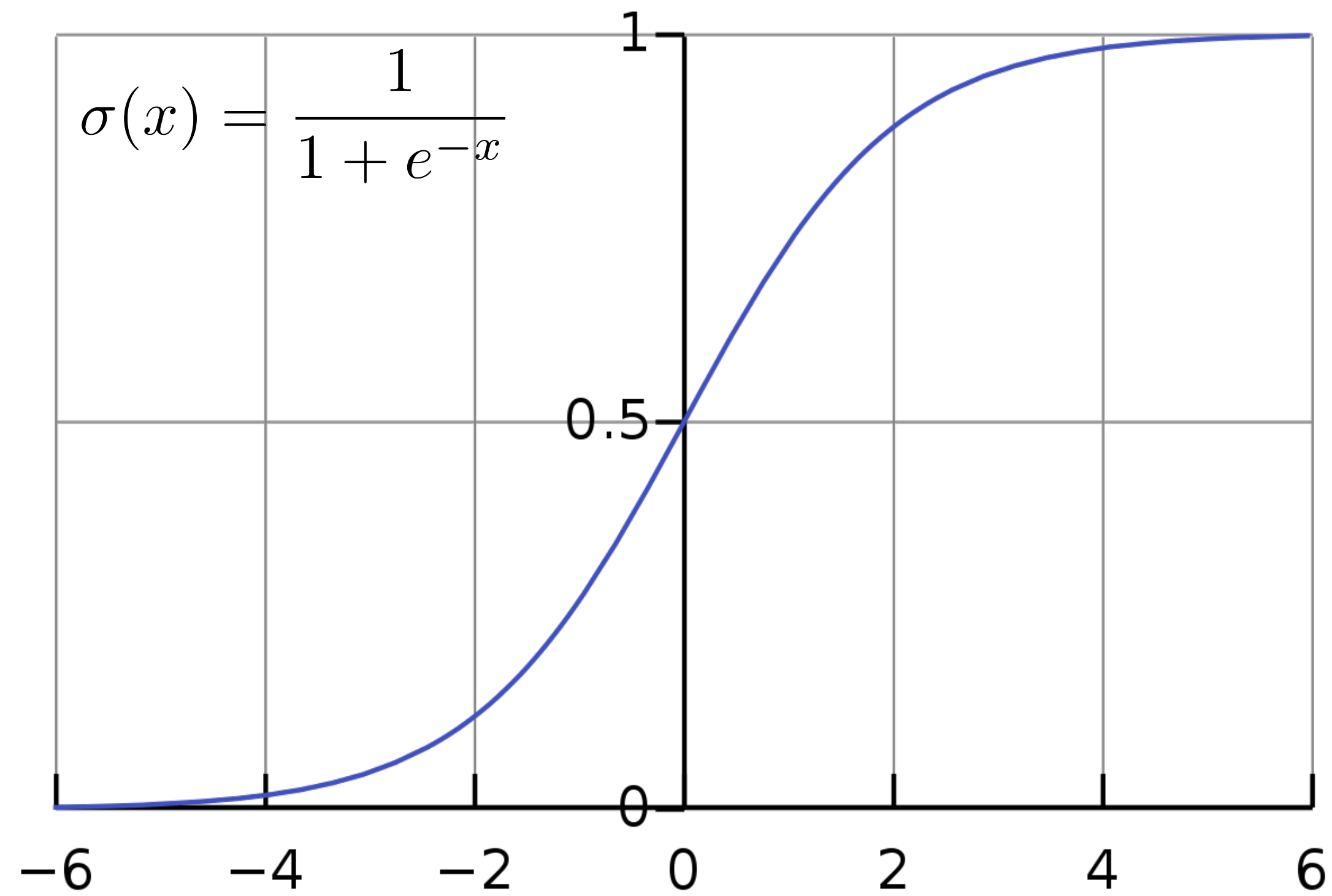




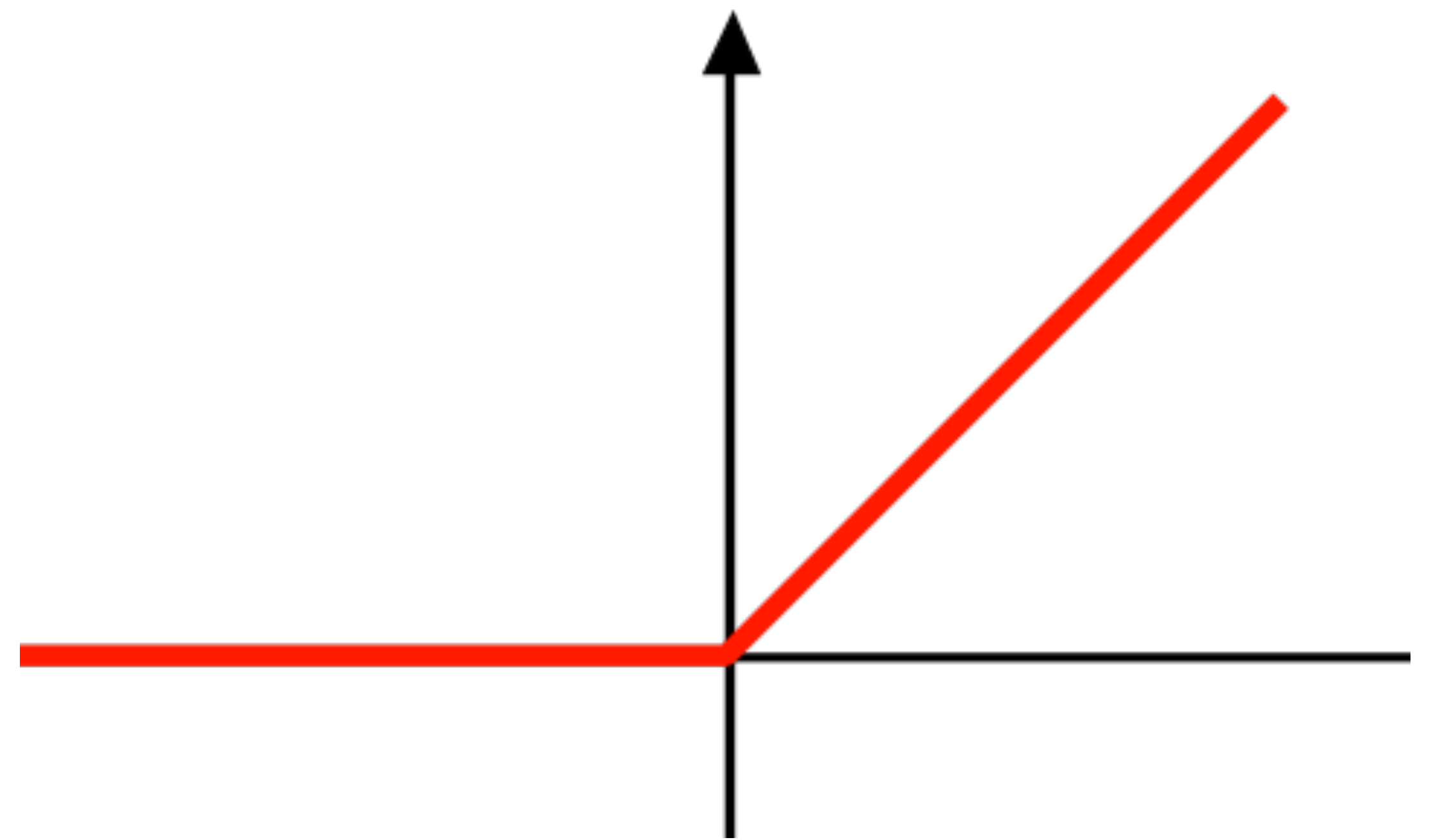


Non-linearities

Sigmoid

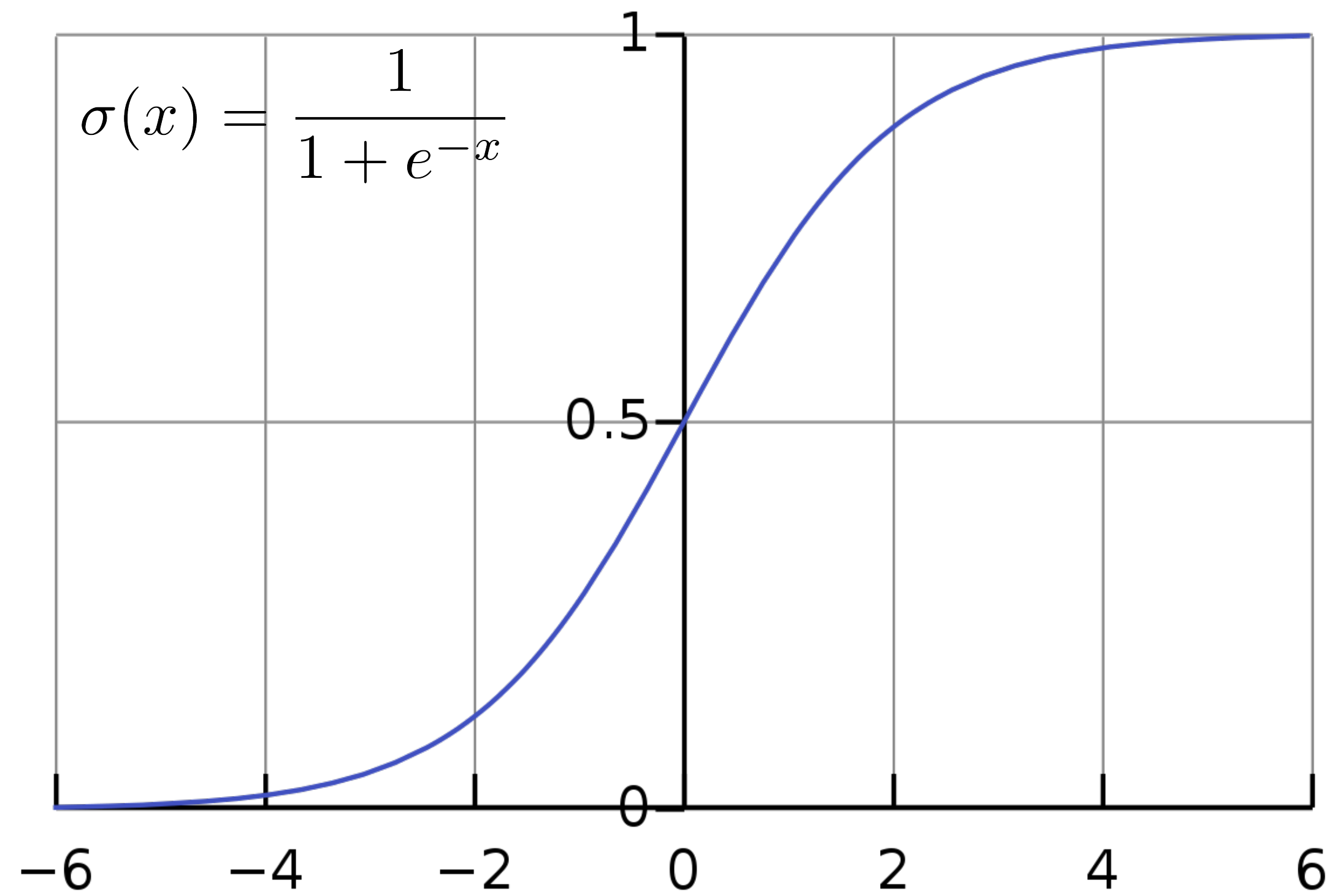


ReLU

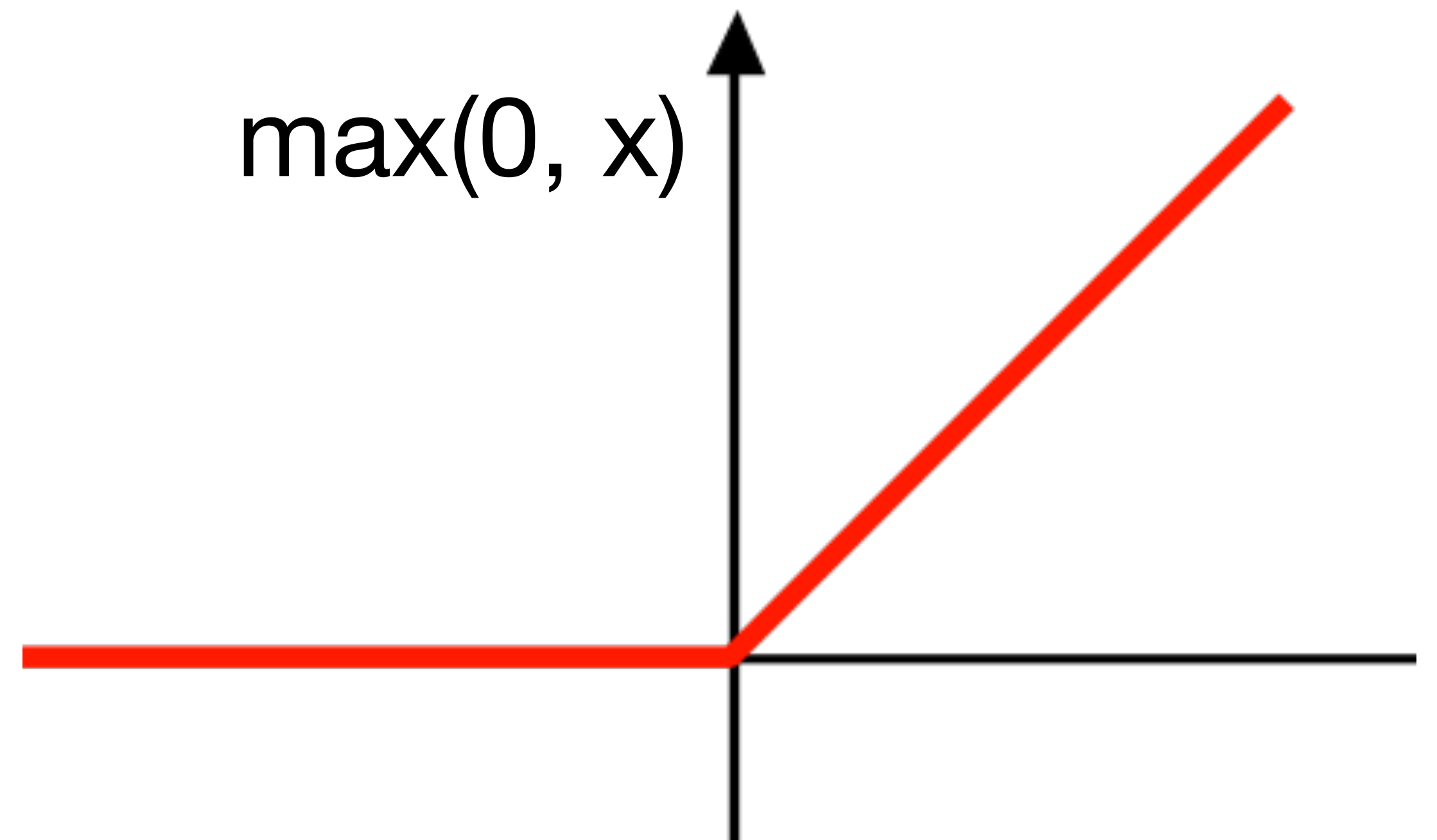


Non-linearities

Sigmoid



ReLU

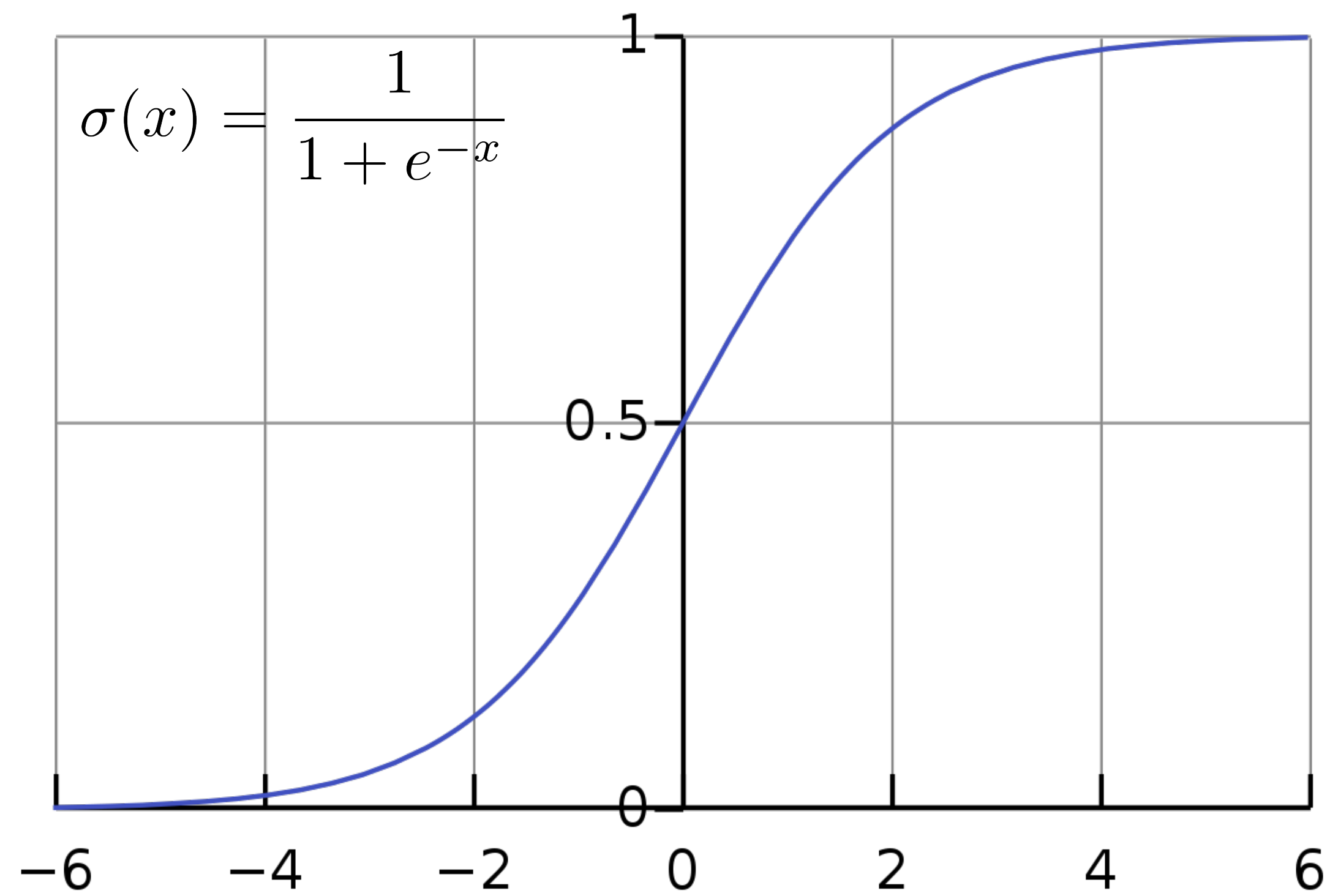


Output layer

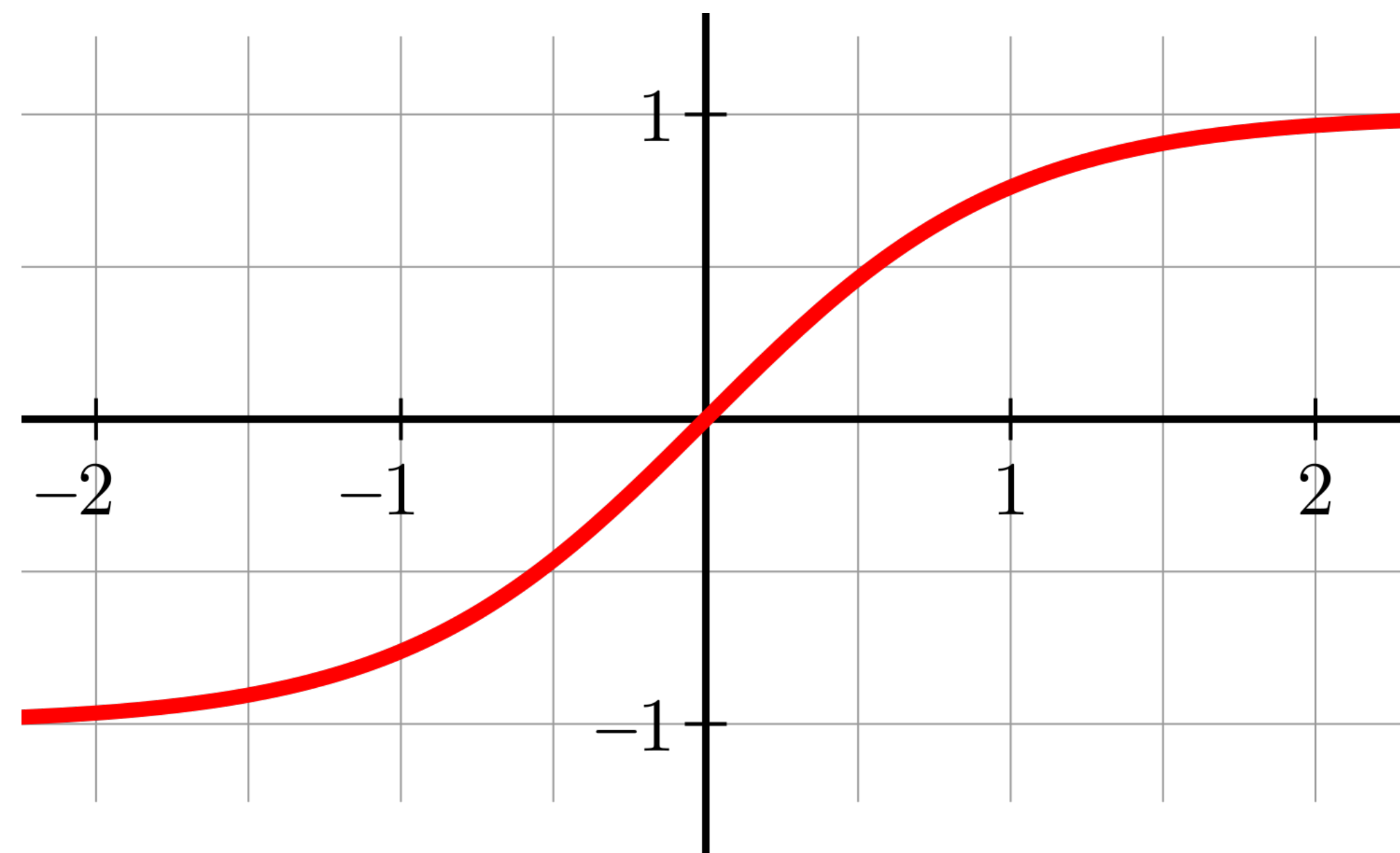




Sigmoid

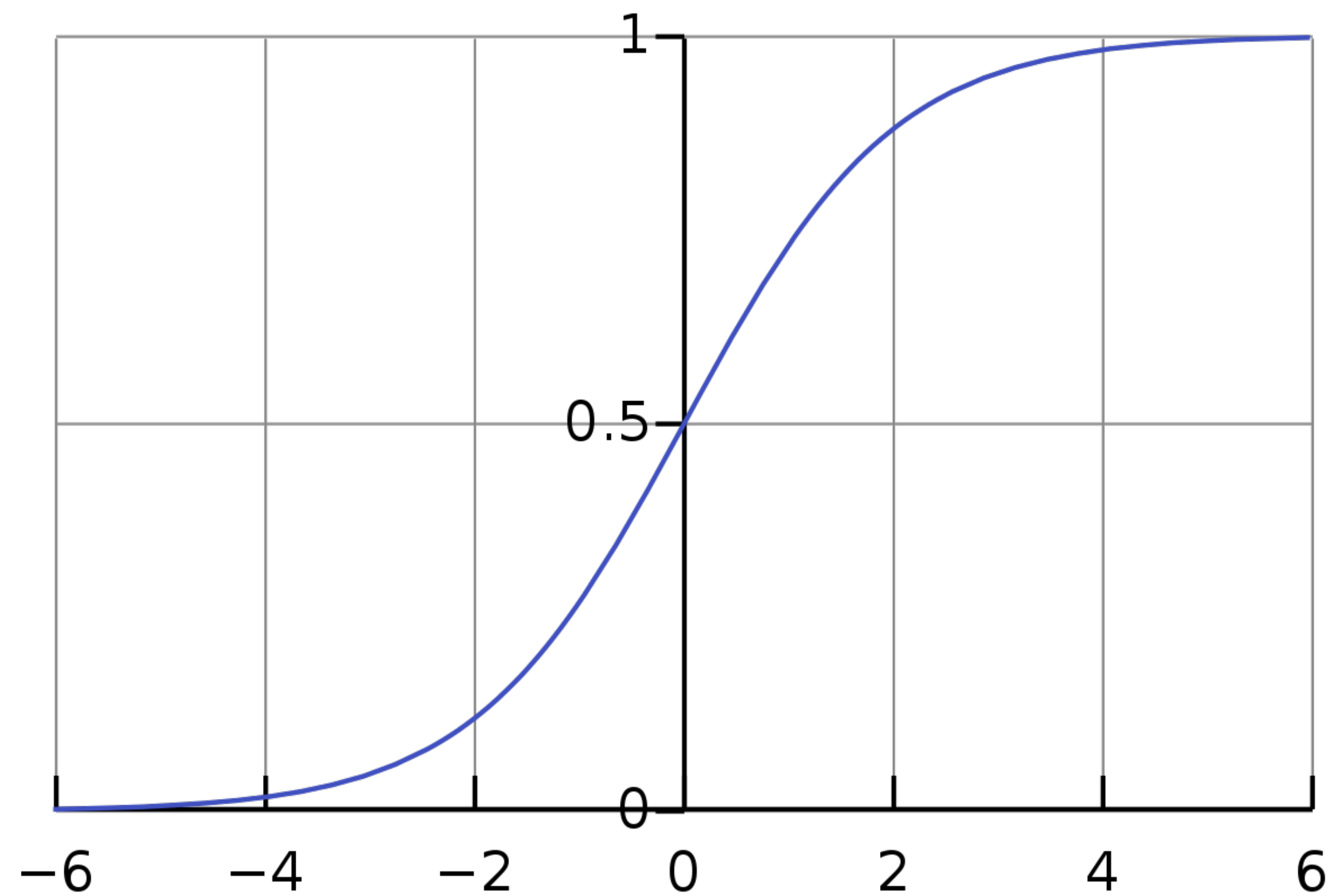


tanh





Sigmoid



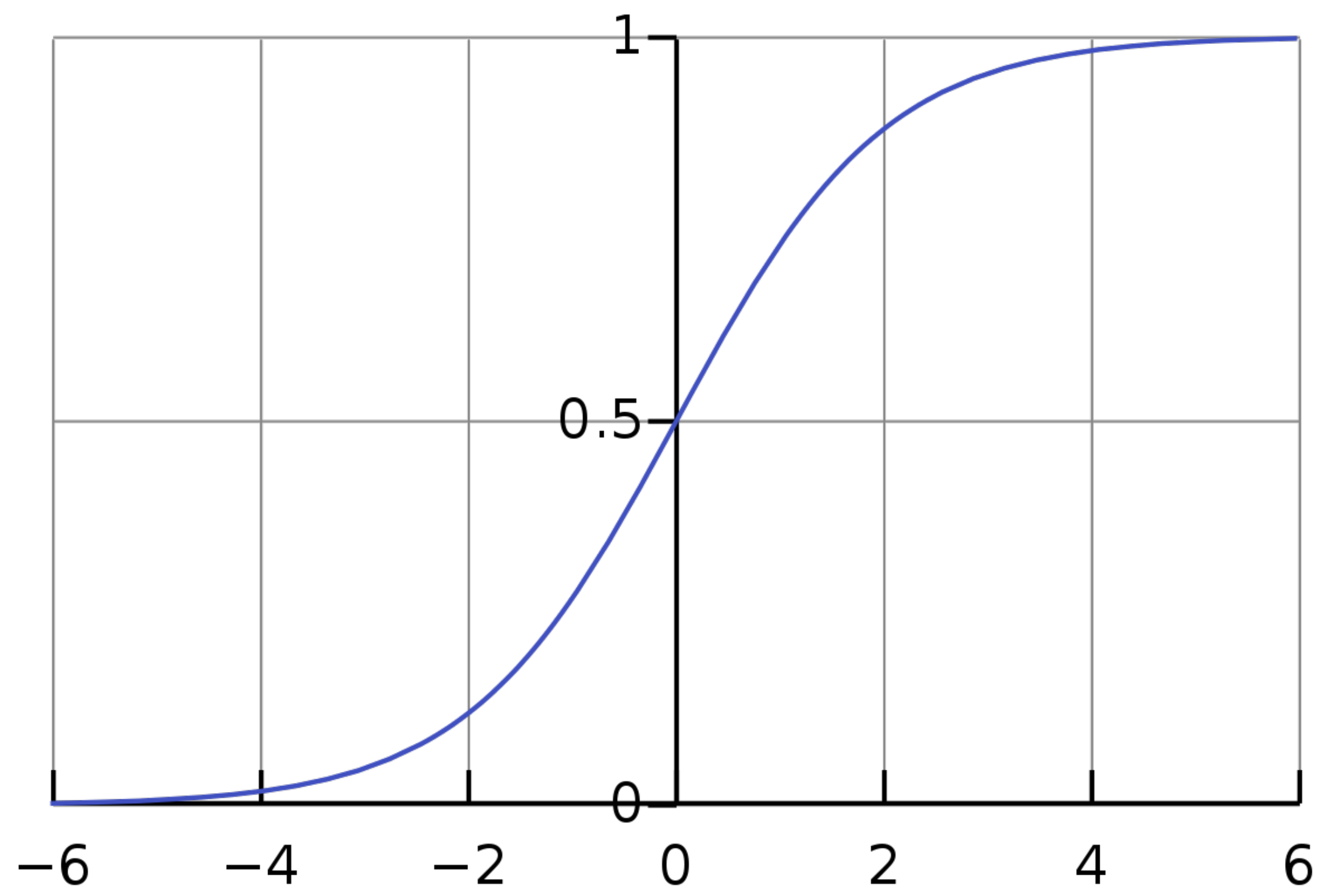
Linear layer

A graph of a linear function, represented by a straight black line with a positive slope. The line starts at a low value on the left and increases steadily towards the right. Below the line is the mathematical expression $Wx + b$.
$$Wx + b$$

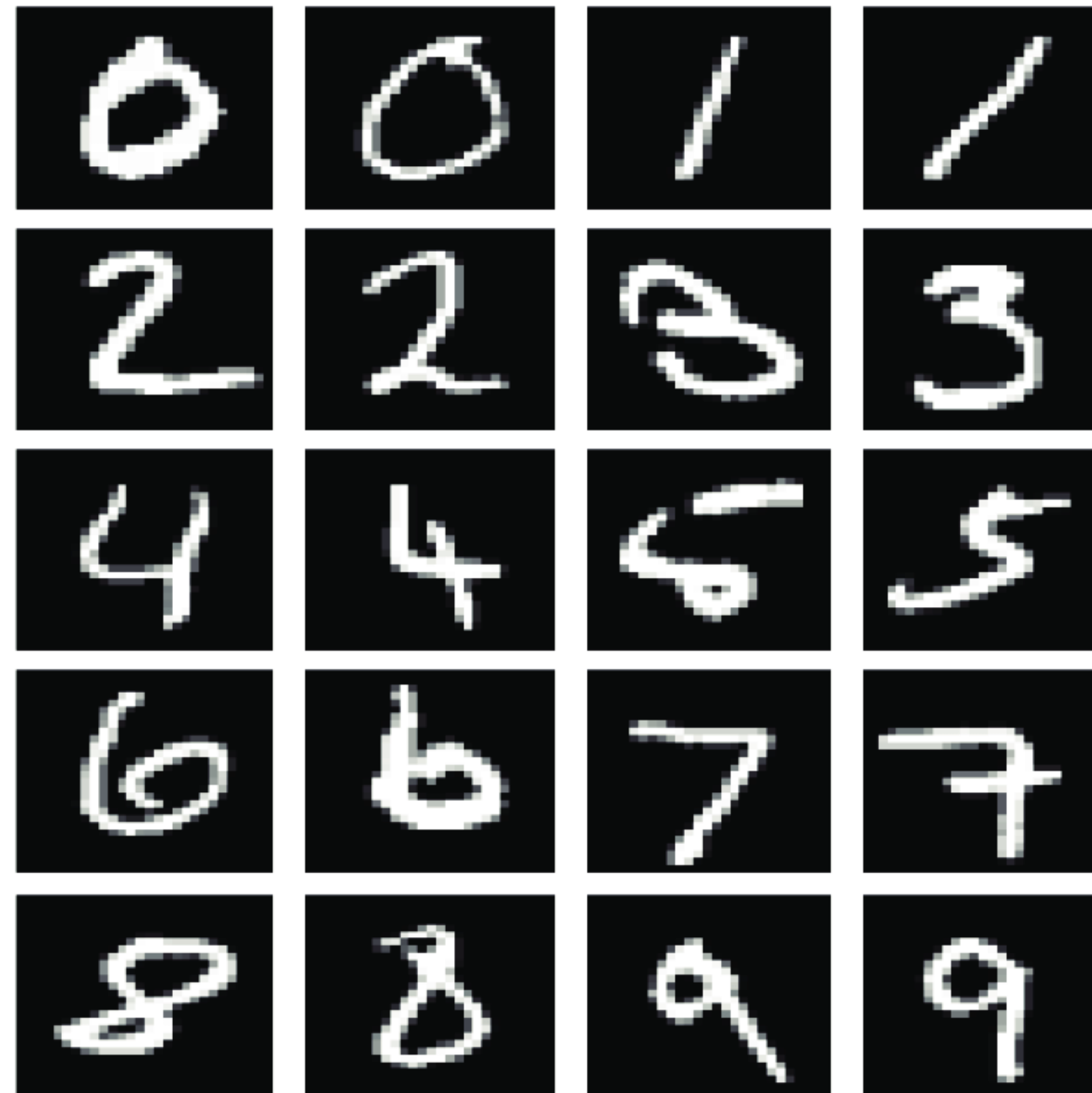




Sigmoid



MNIST dataset

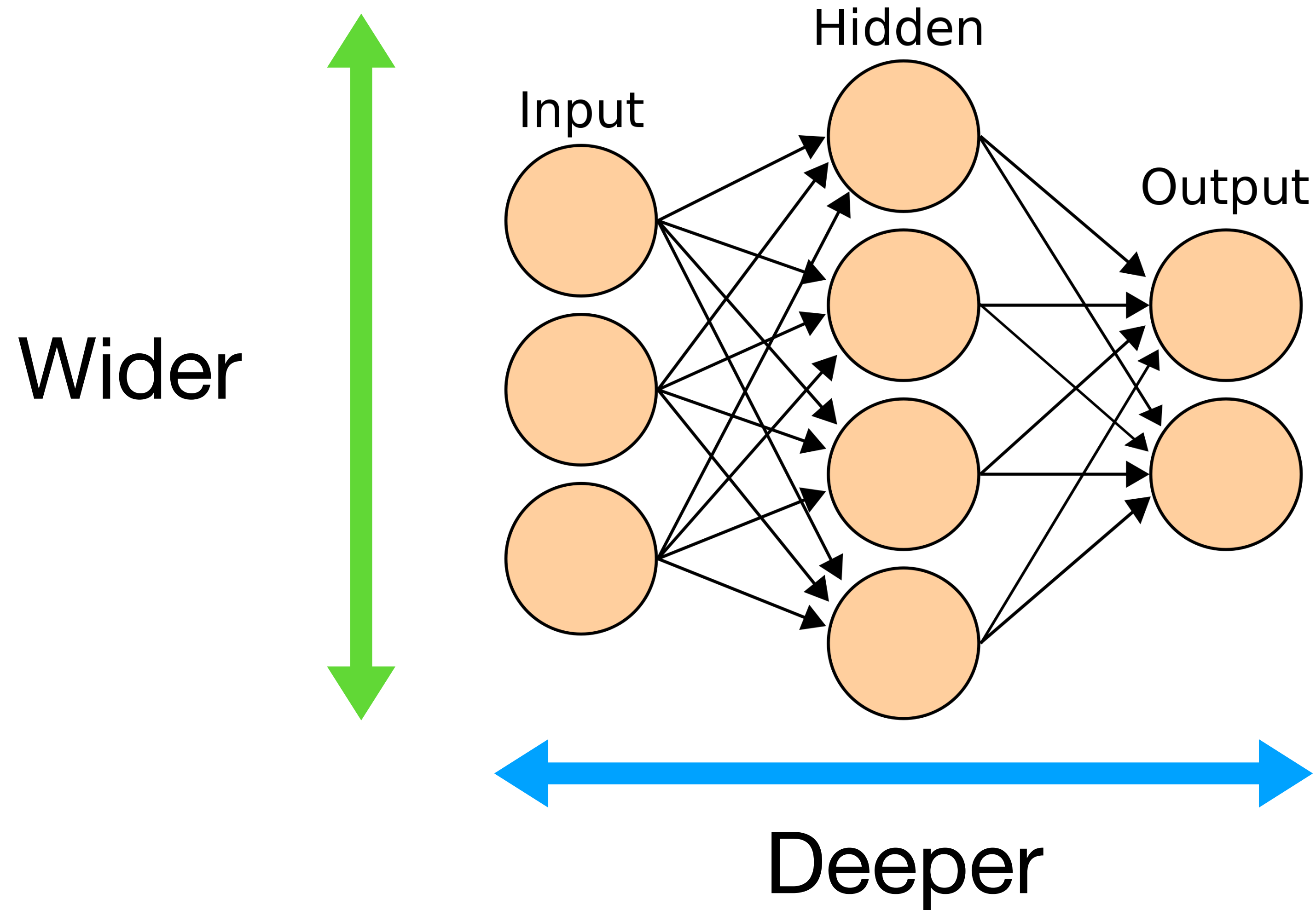


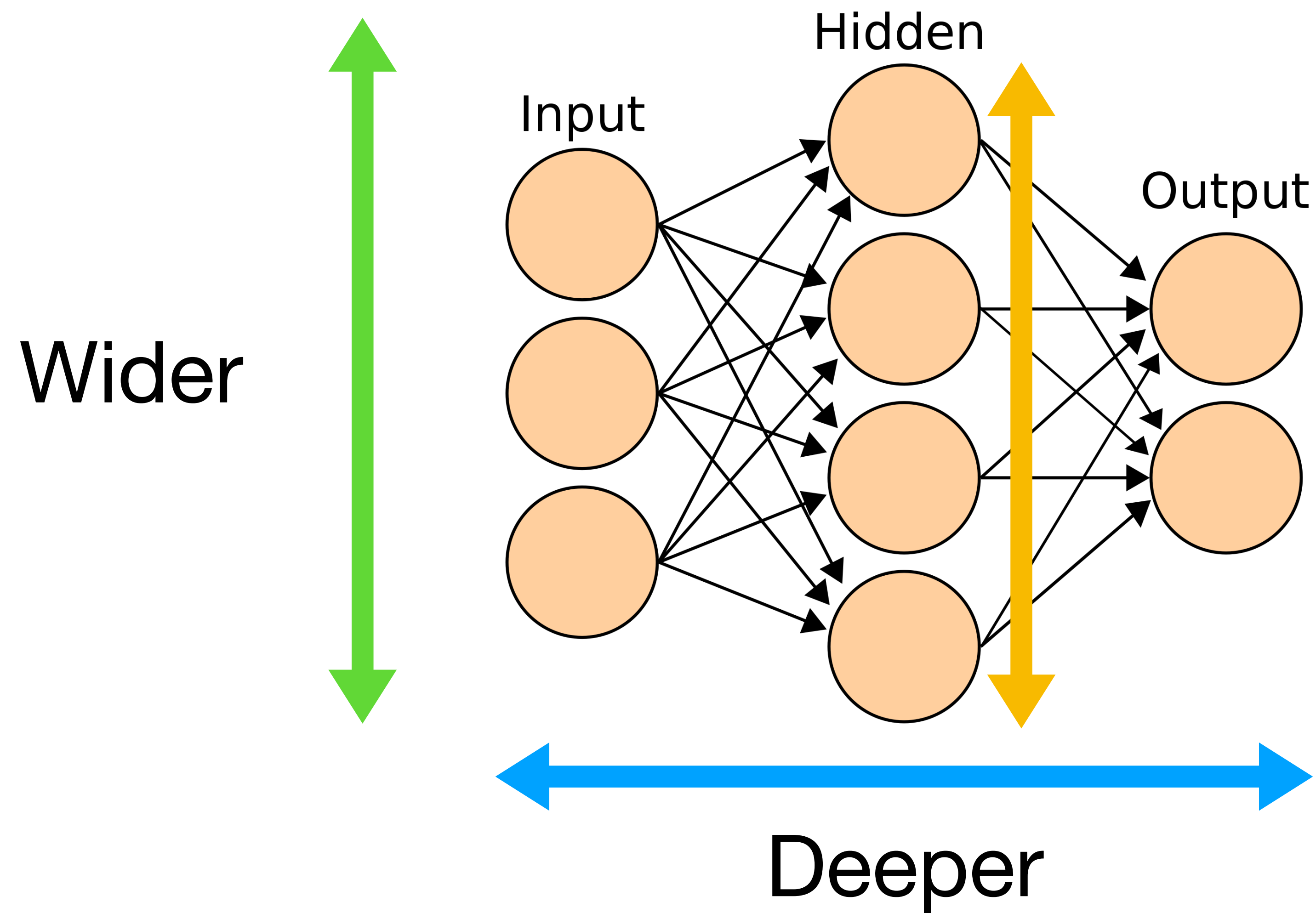
`[-1, 4, 7]`



Softmax(3)

`[0.0003 0.047 0.95]`

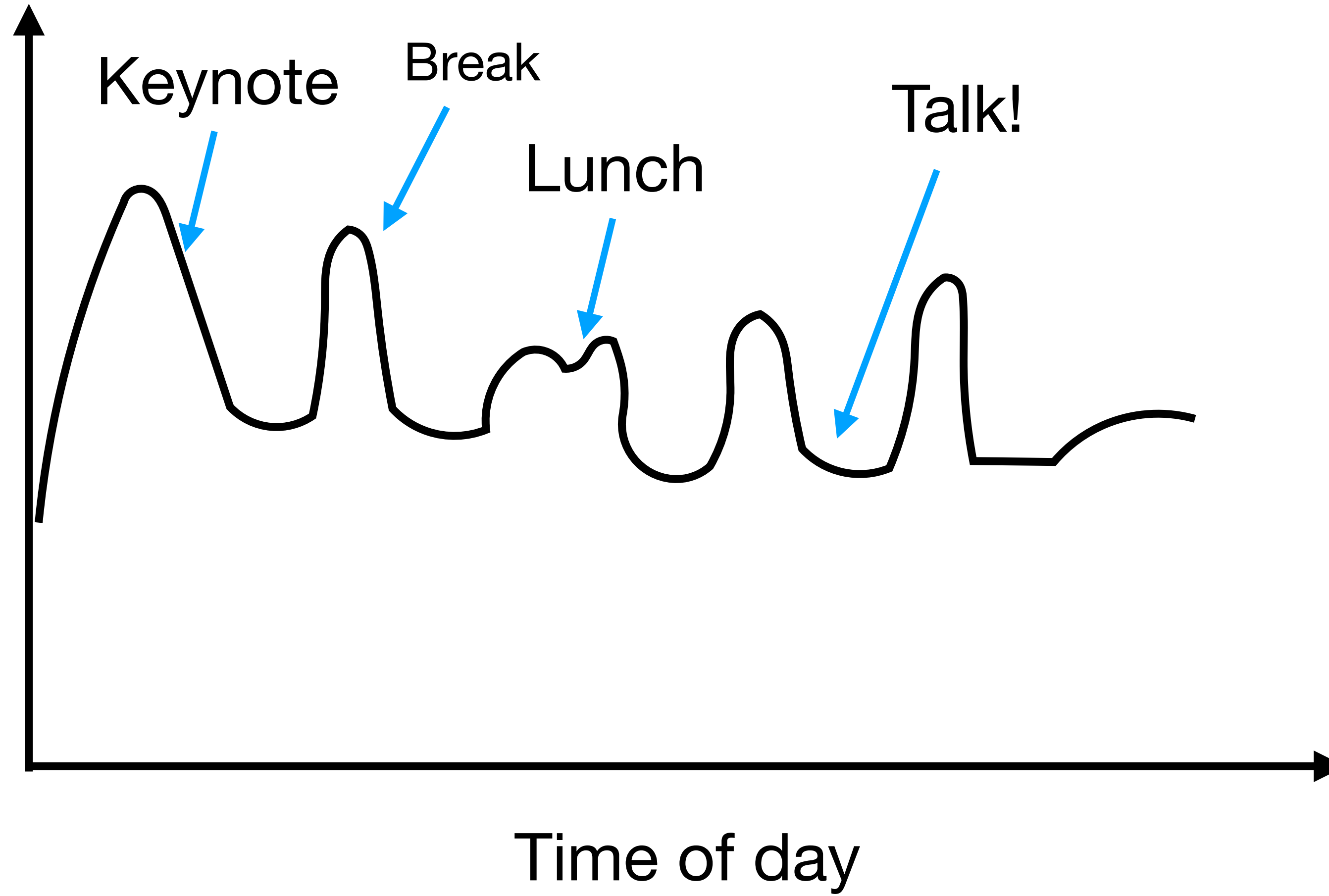




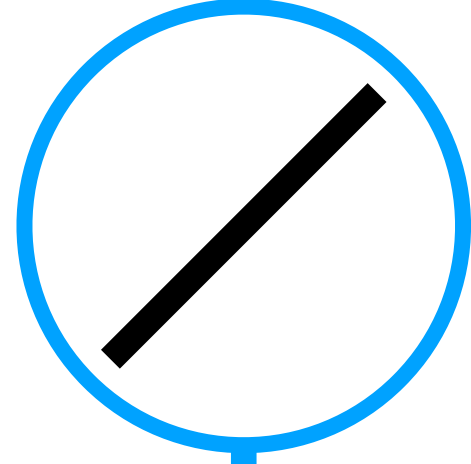
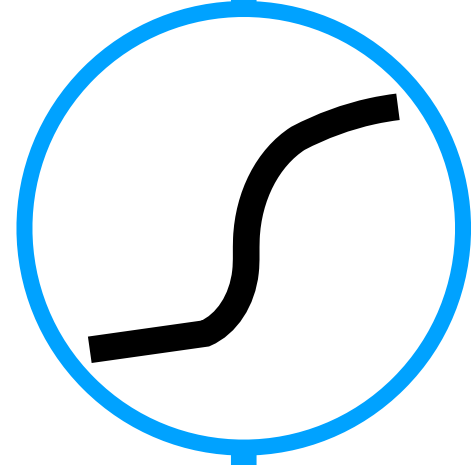
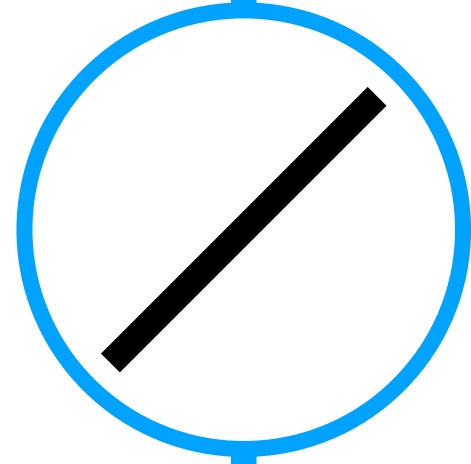




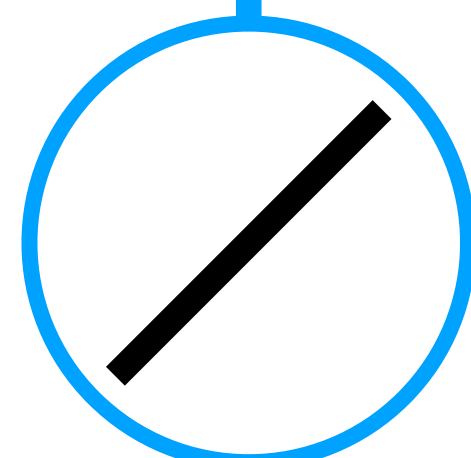
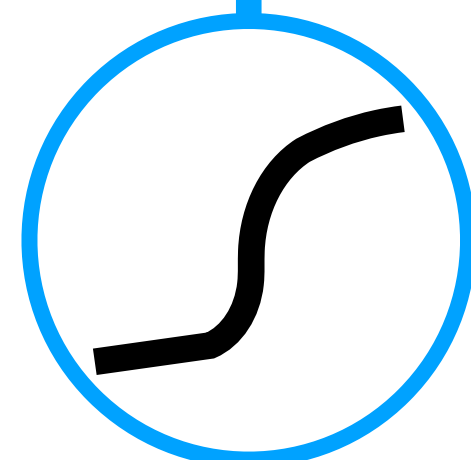
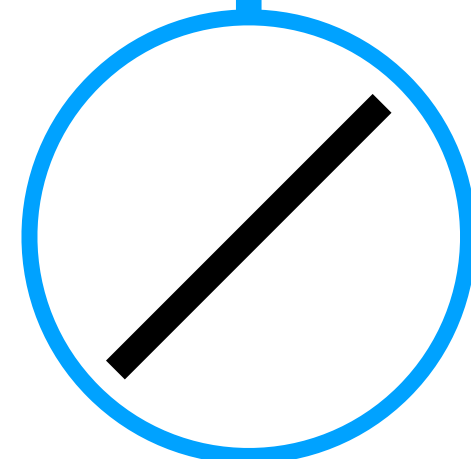
Coffee
needed



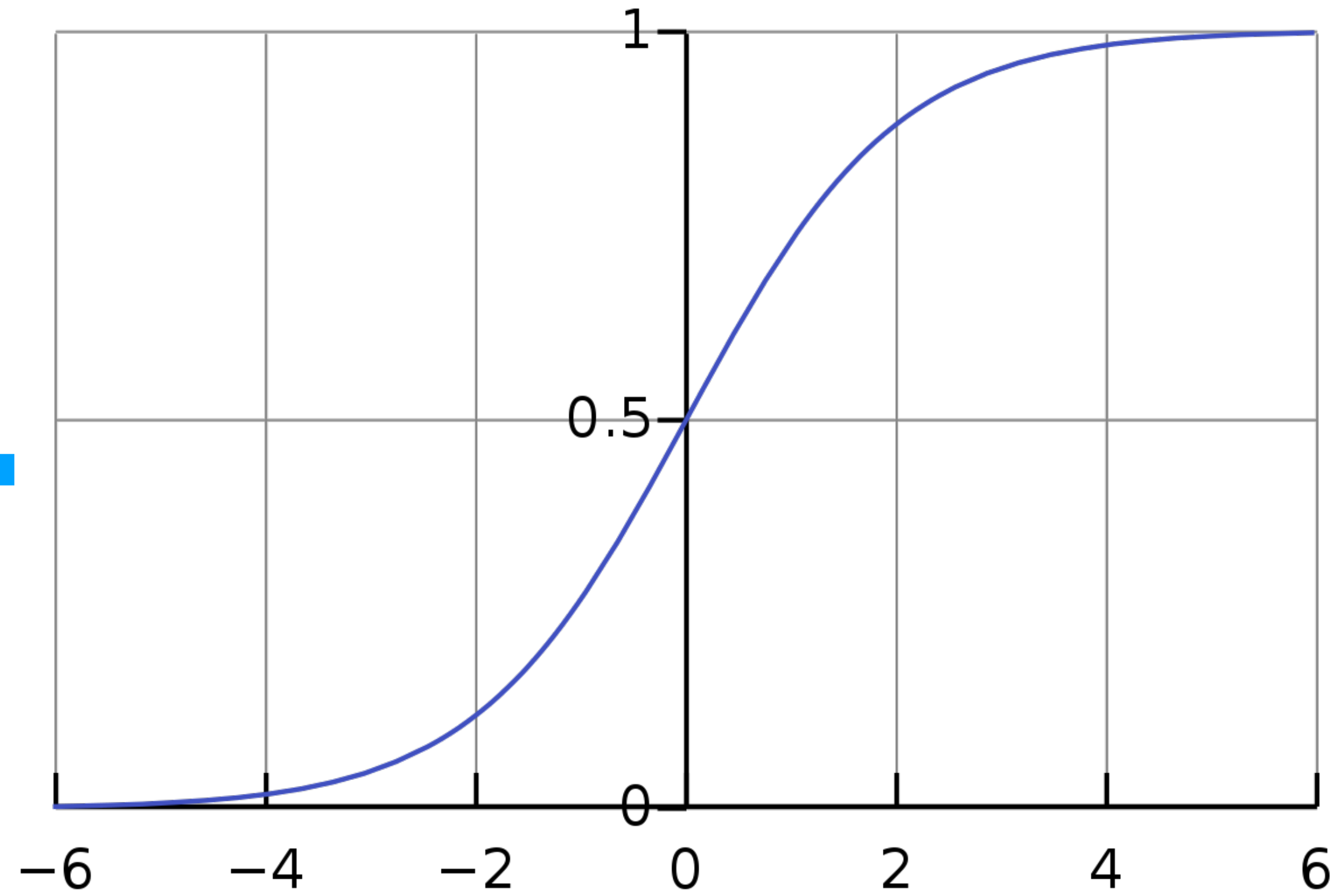
In



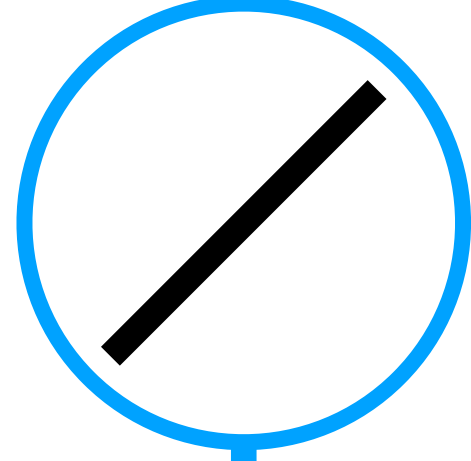
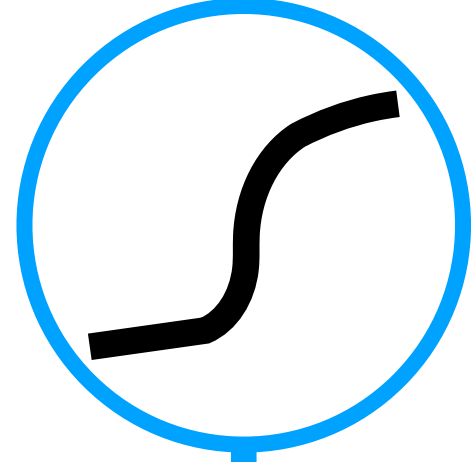
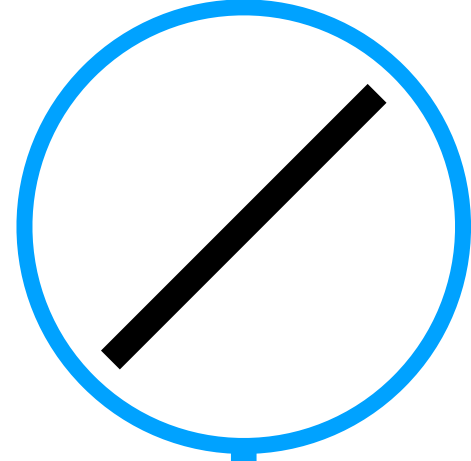
In



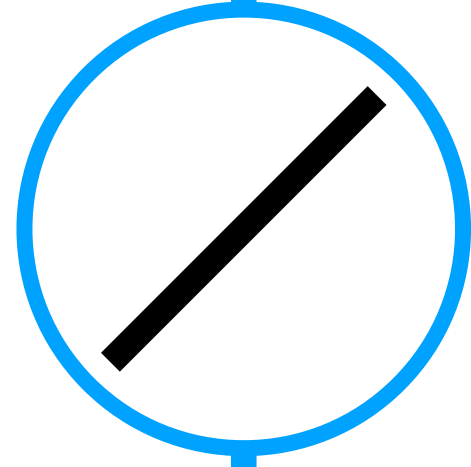
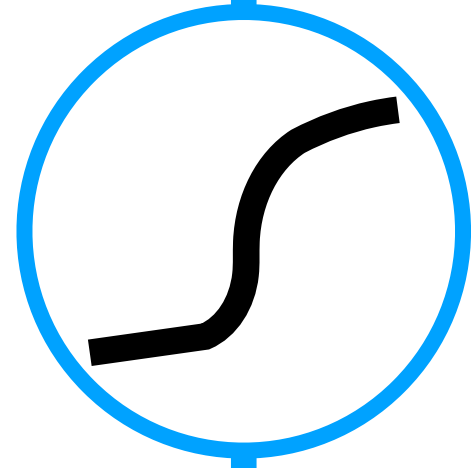
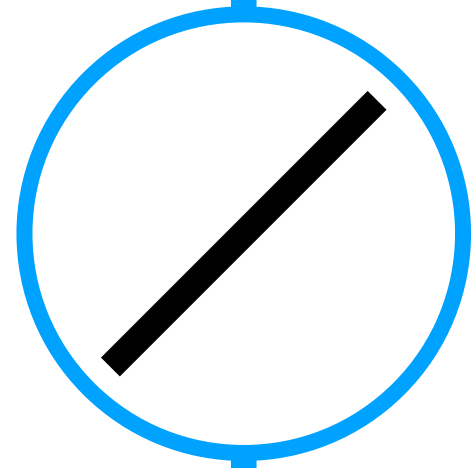
Sigmoid



In



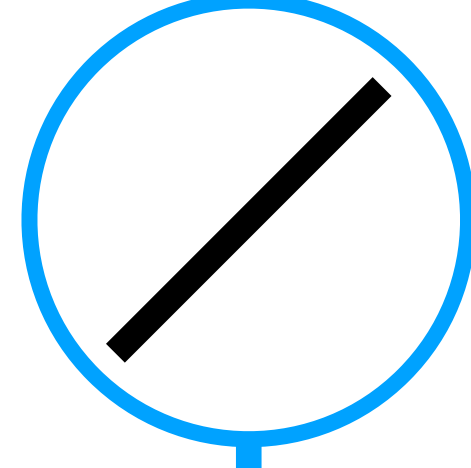
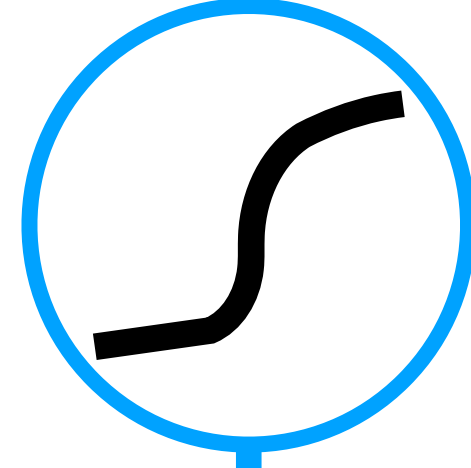
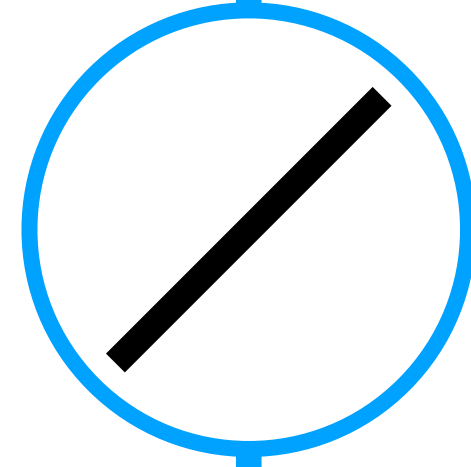
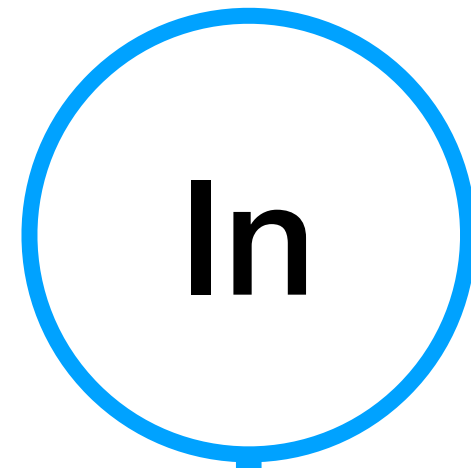
In



Coffee
needed



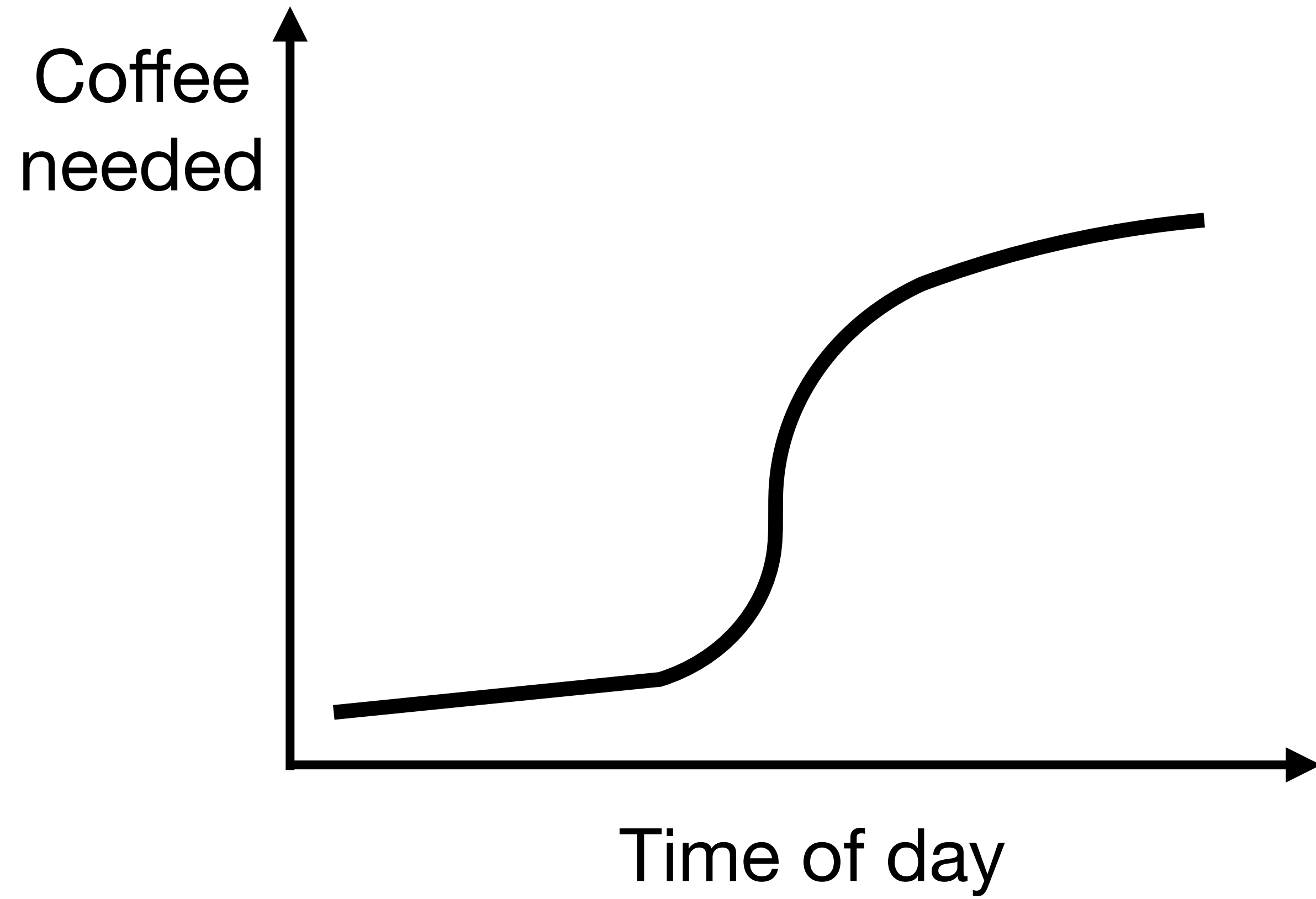
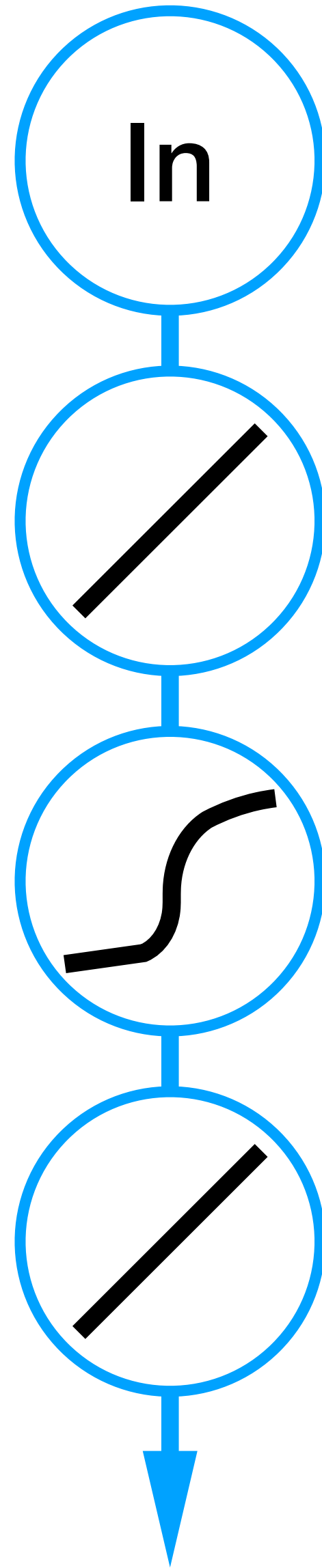
Sophisticated feature

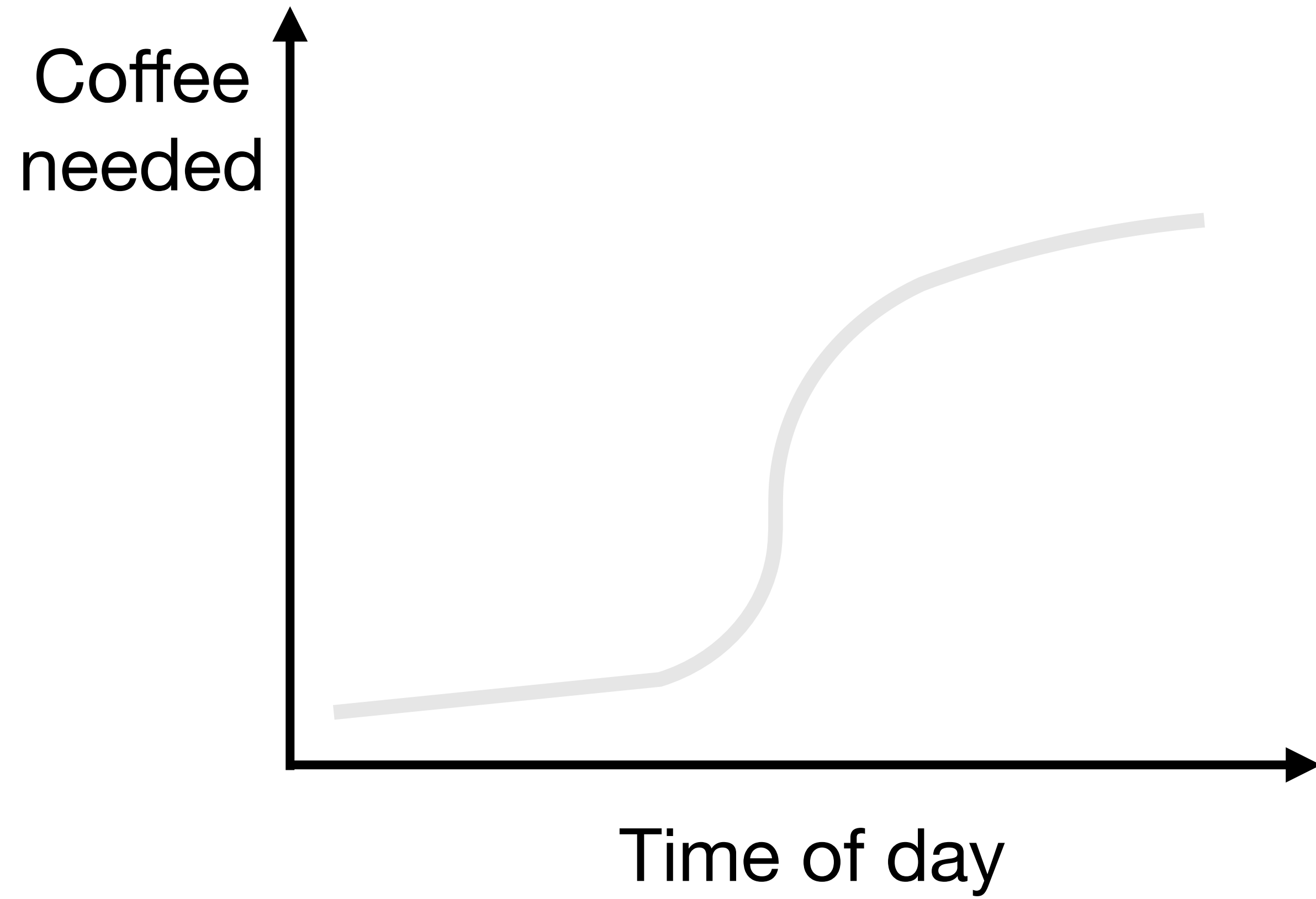
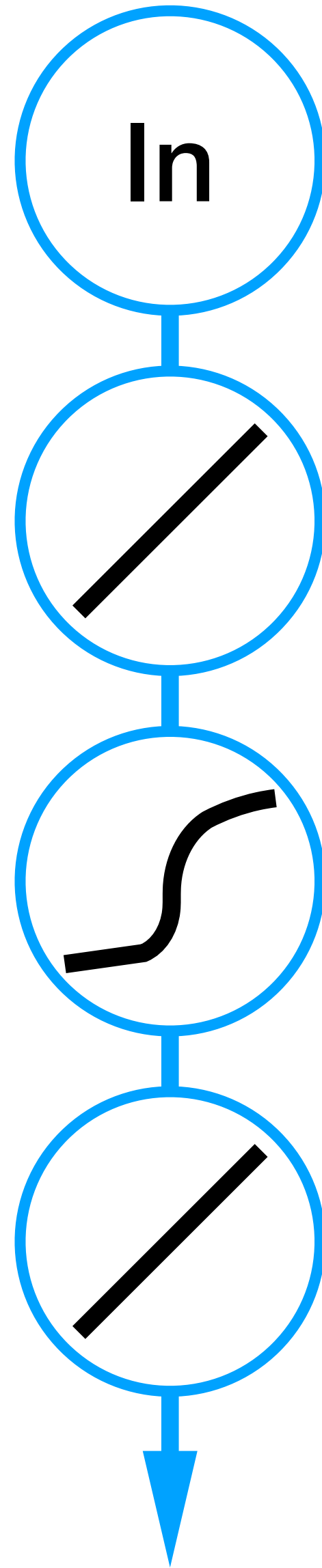


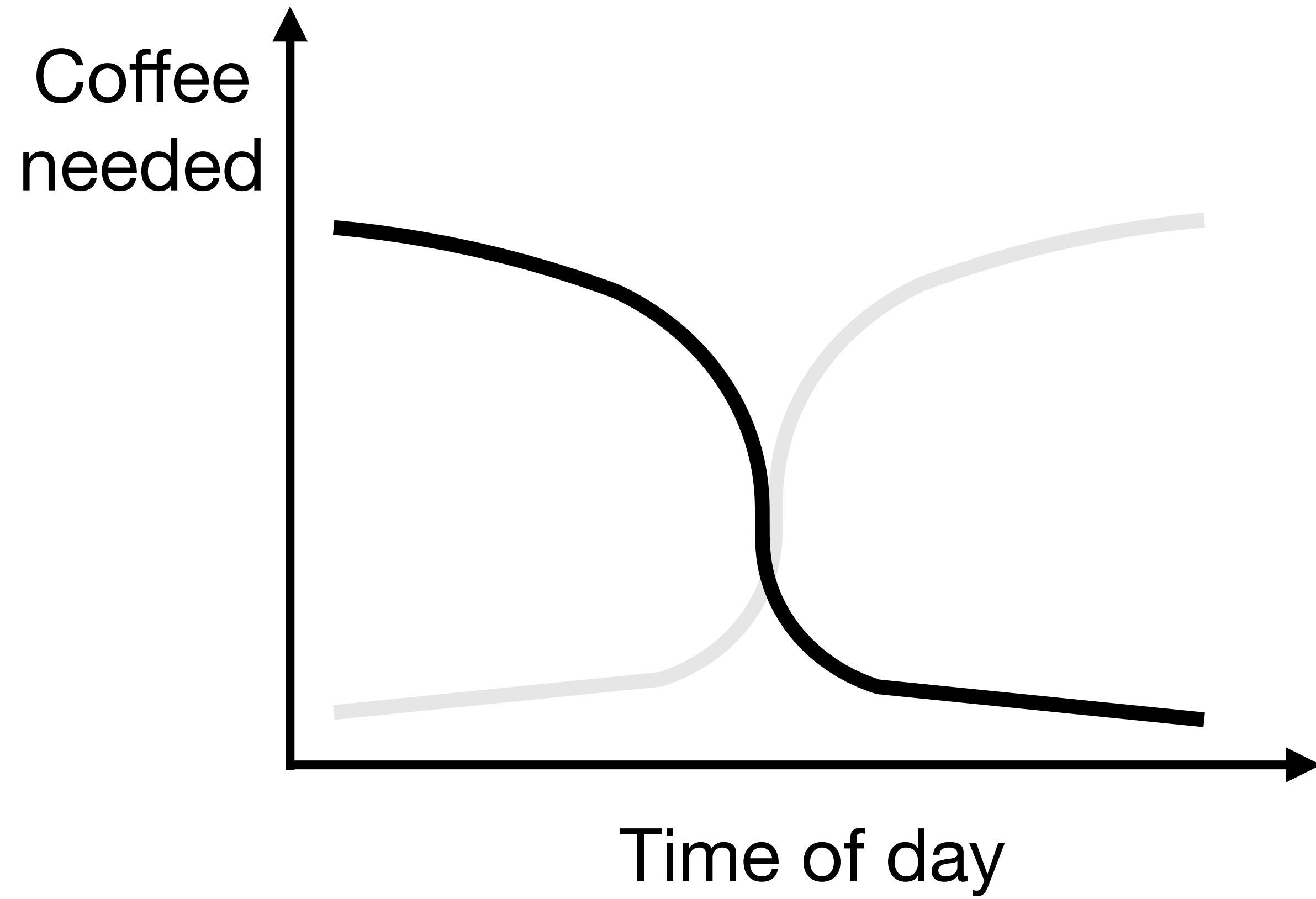
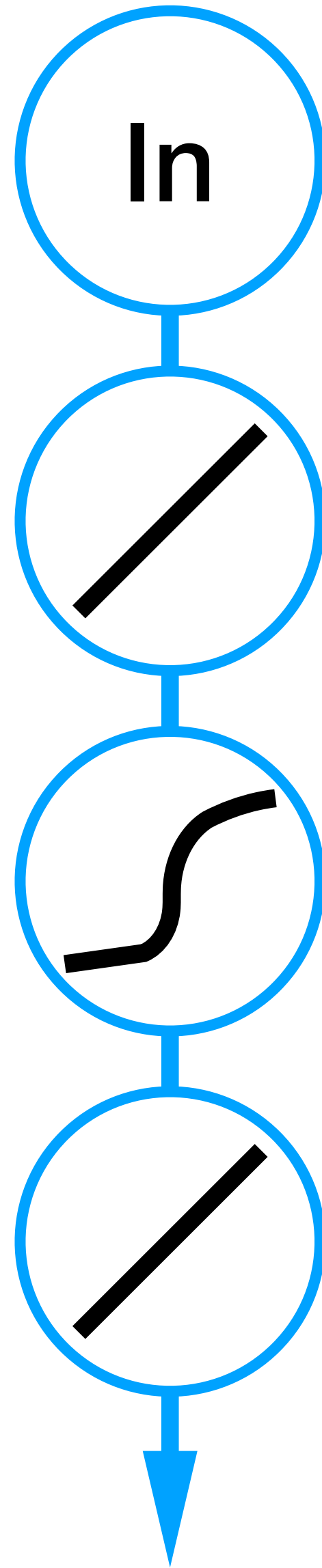
Coffee
needed

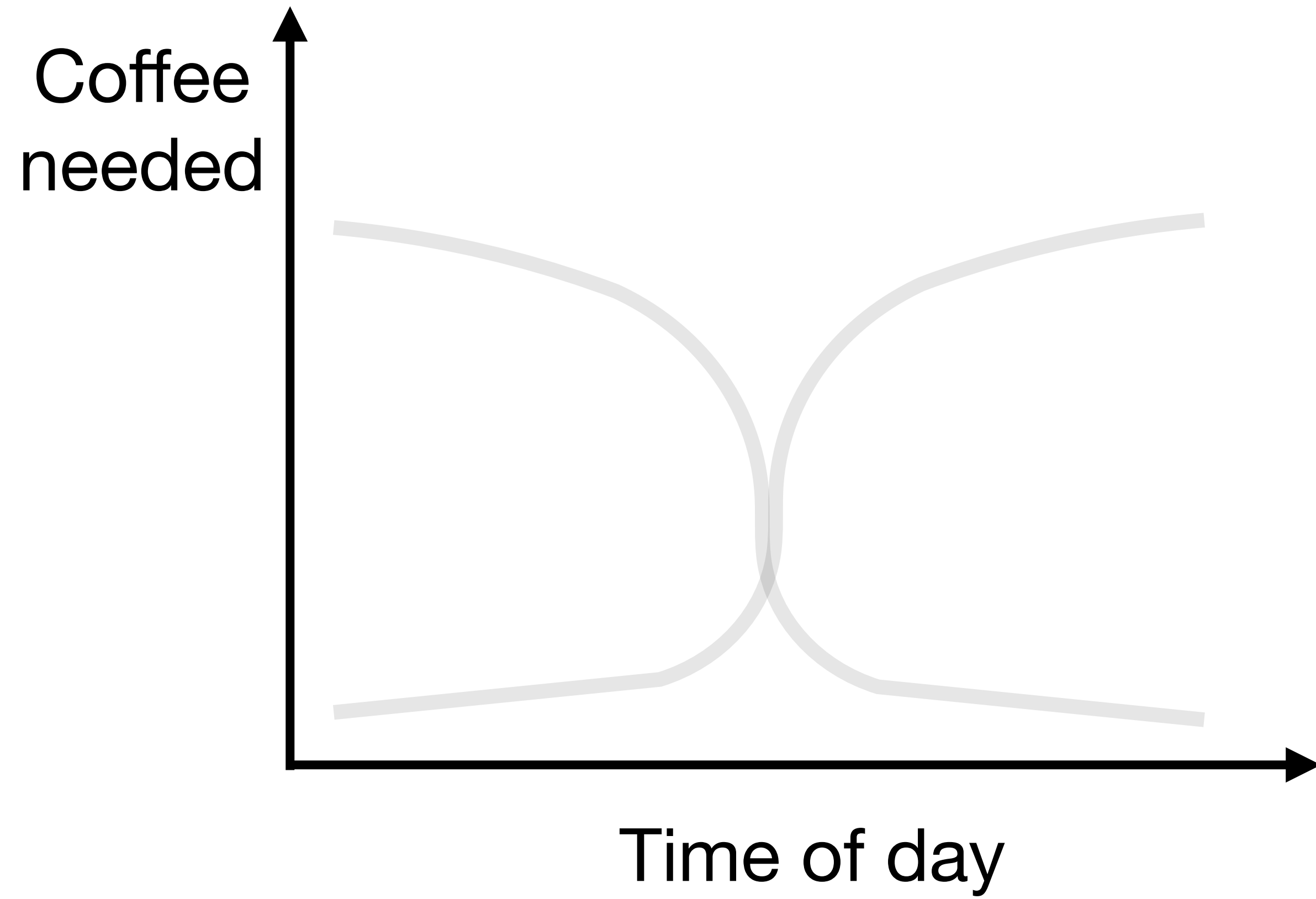
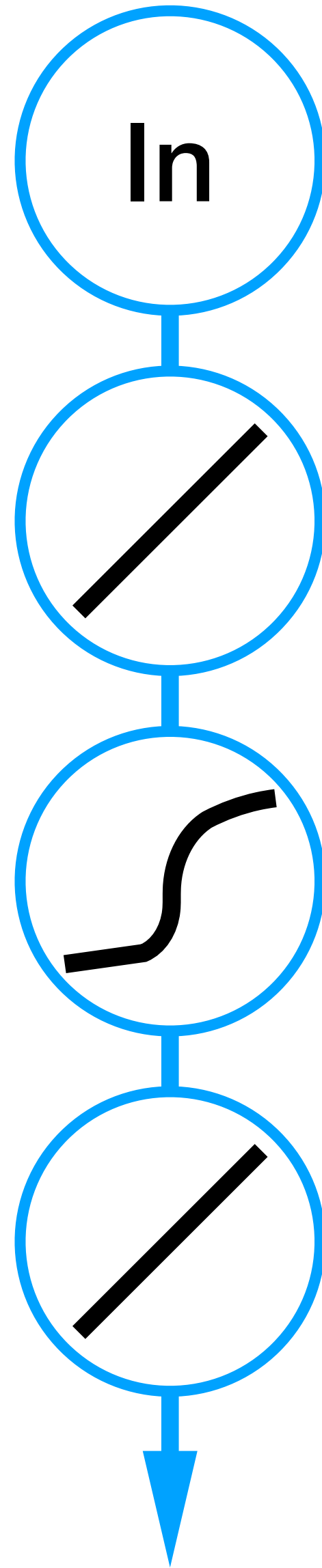


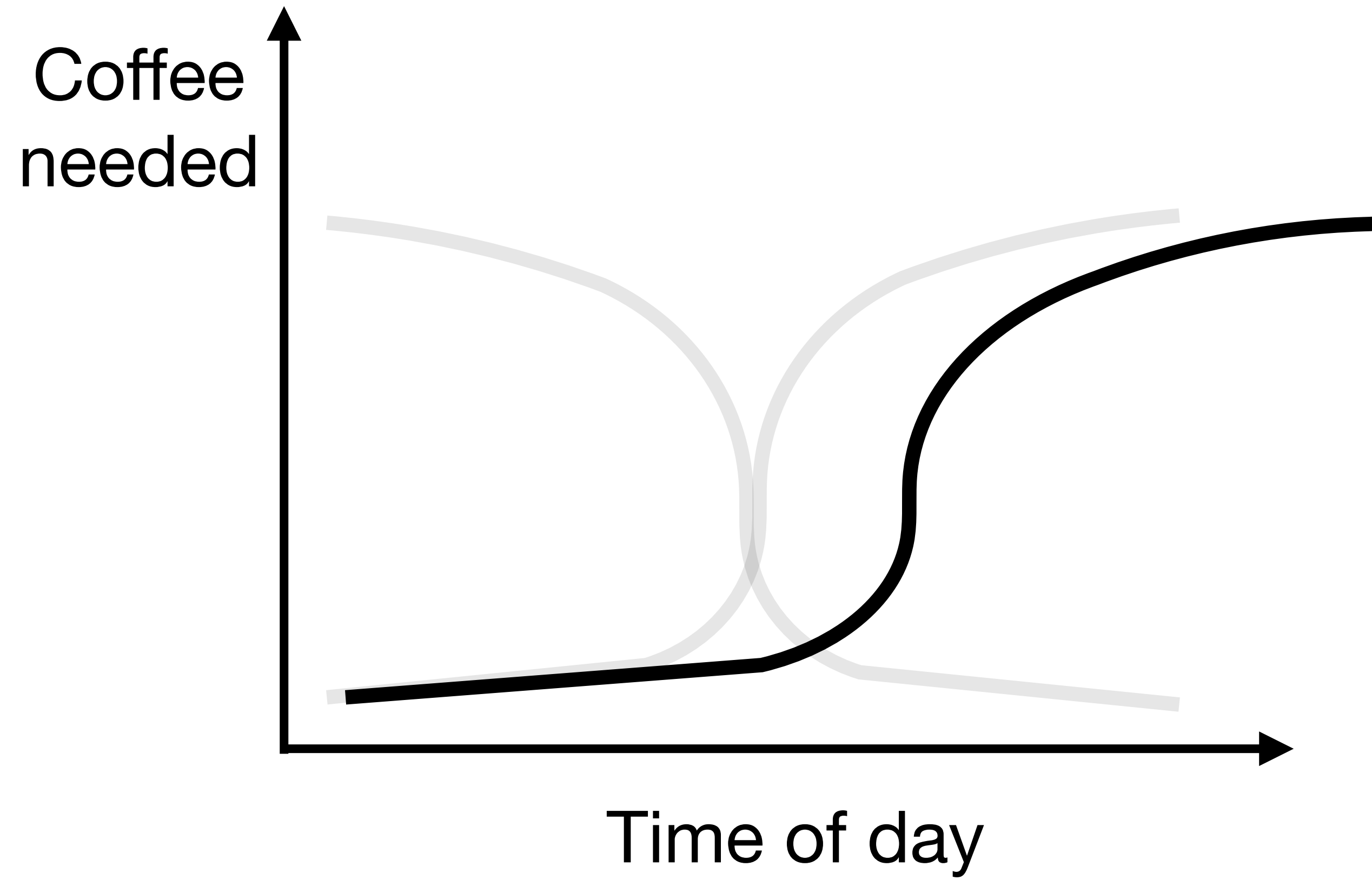
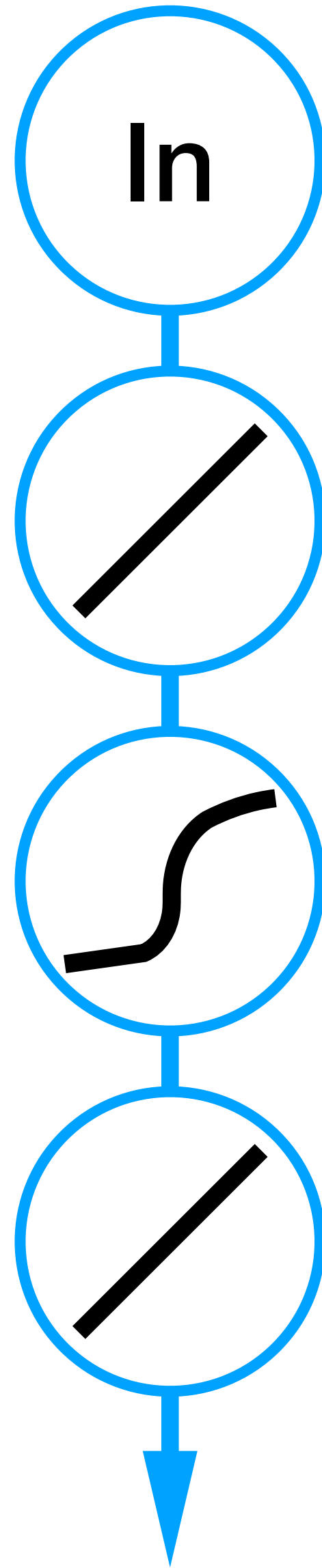
Time of day

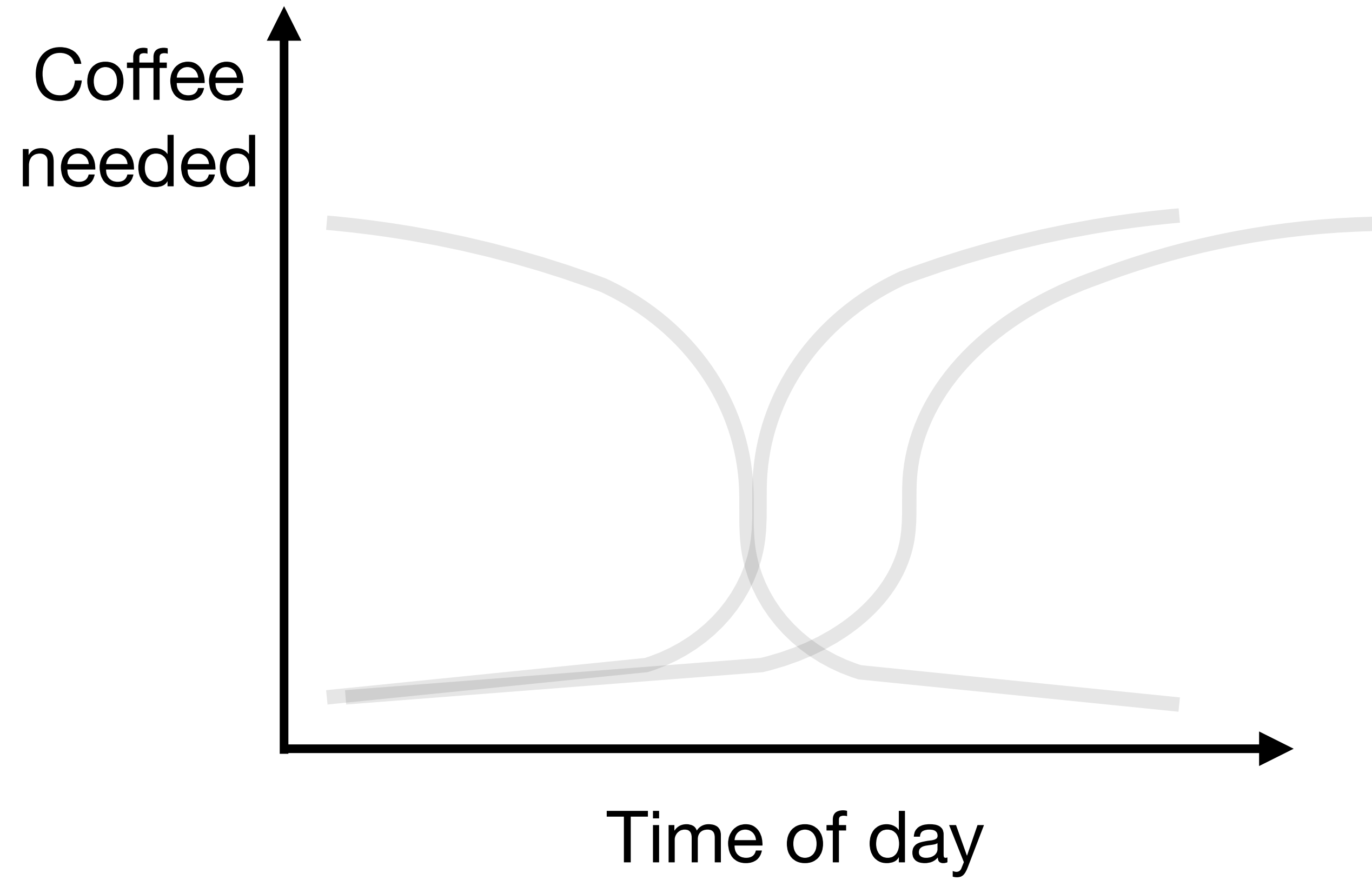
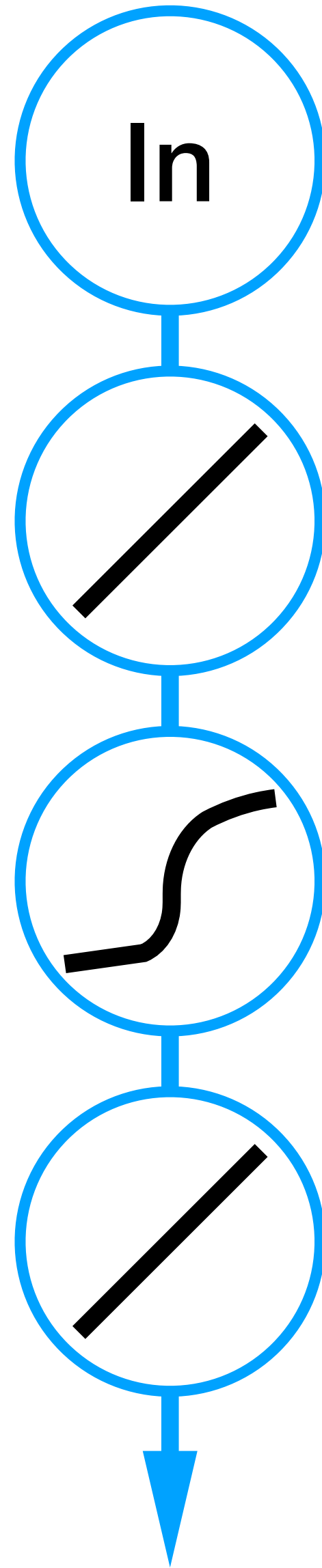


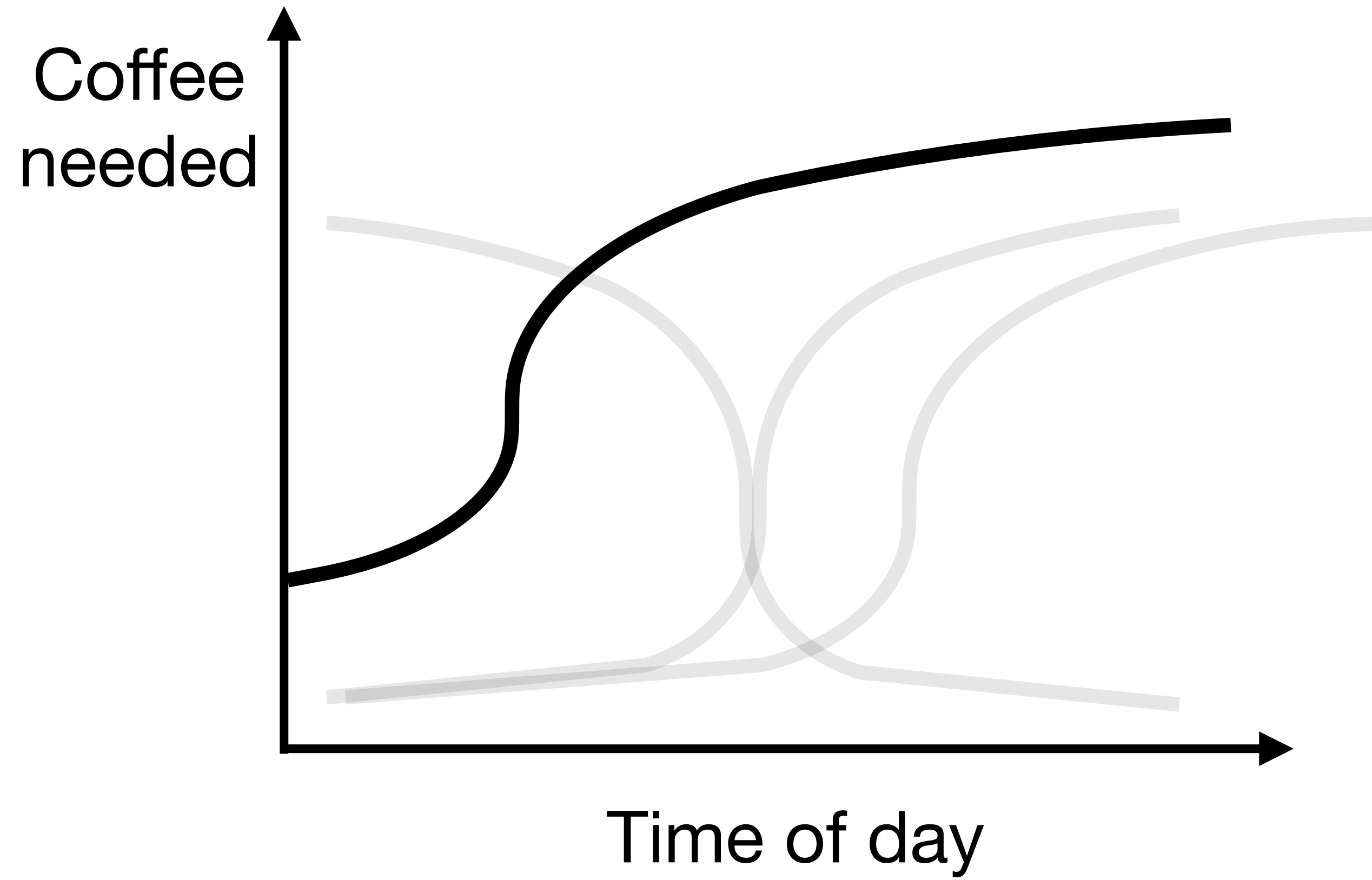
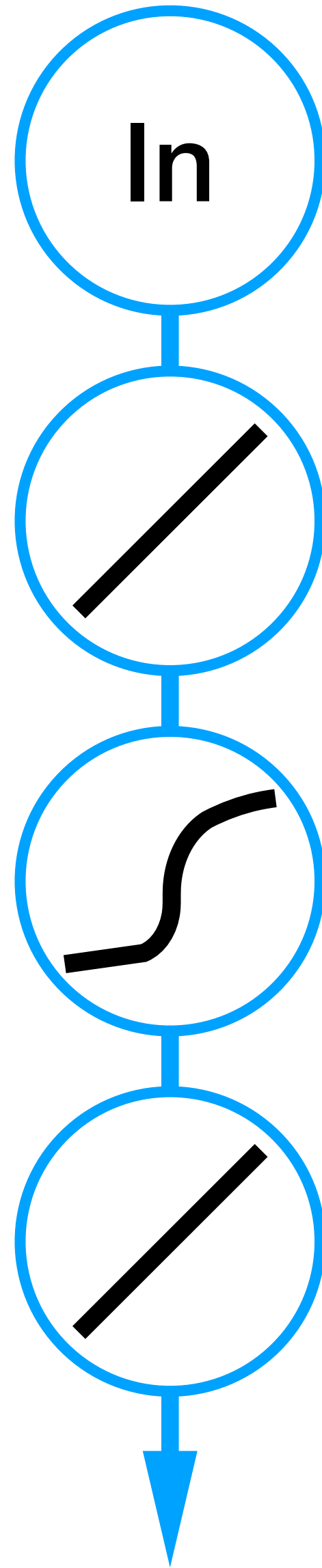


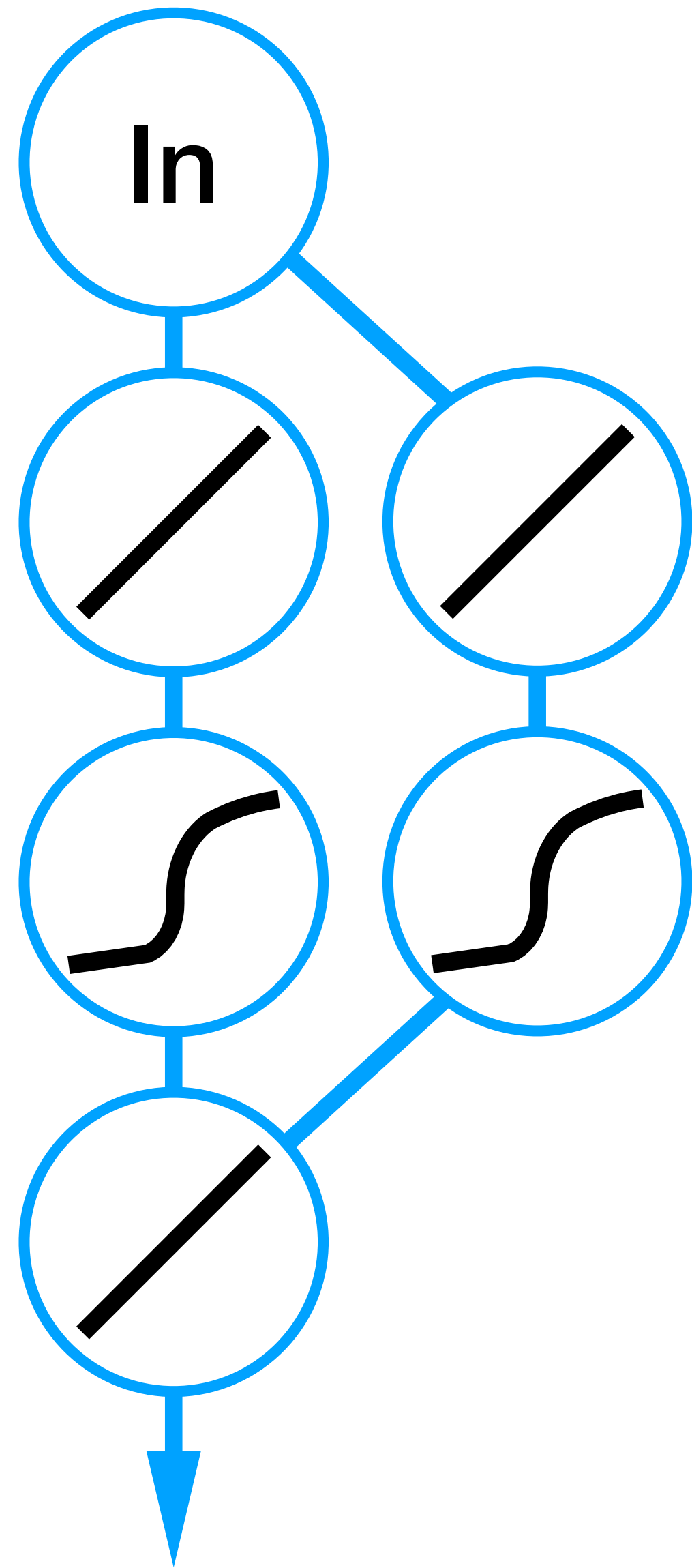


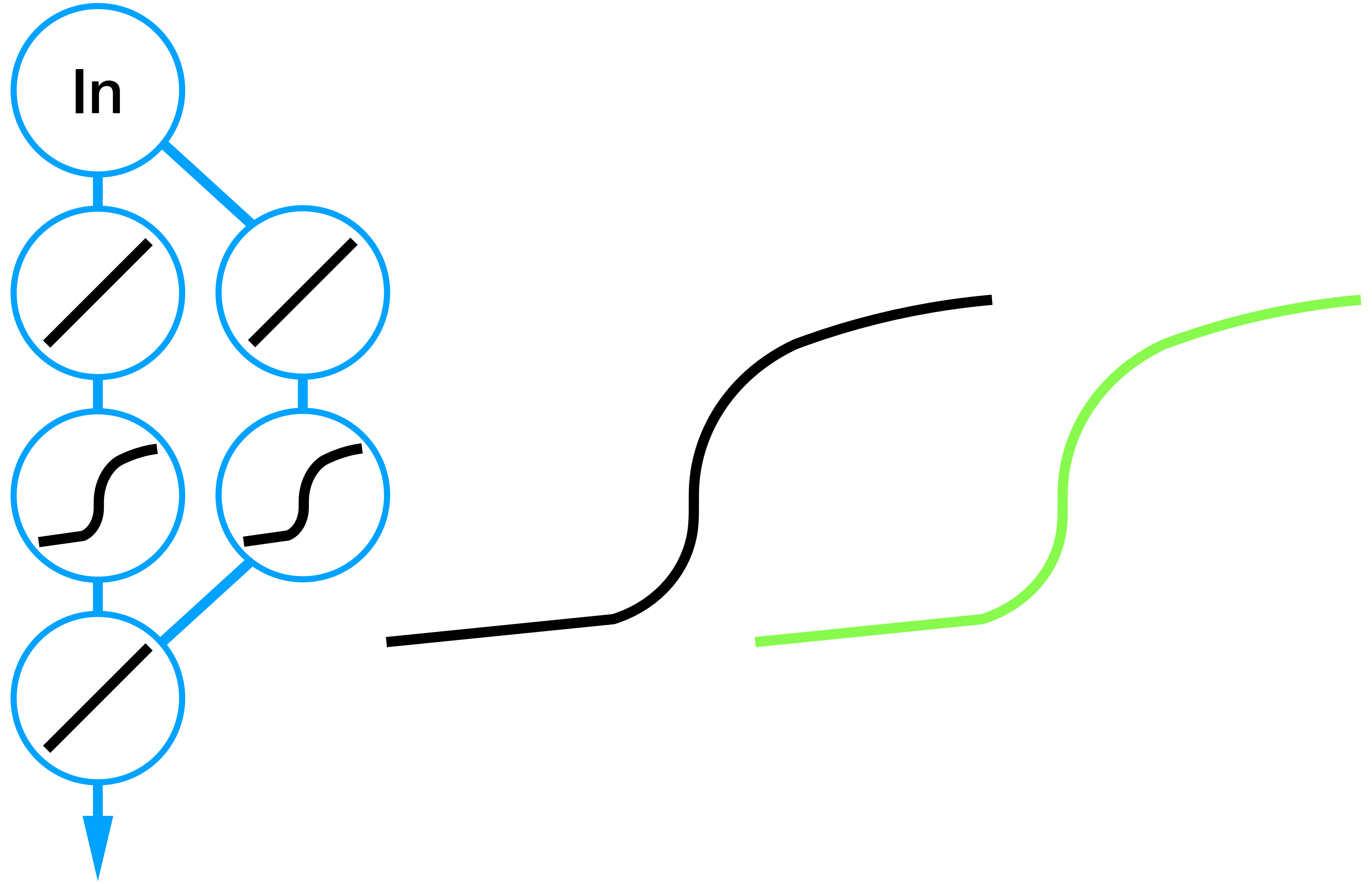


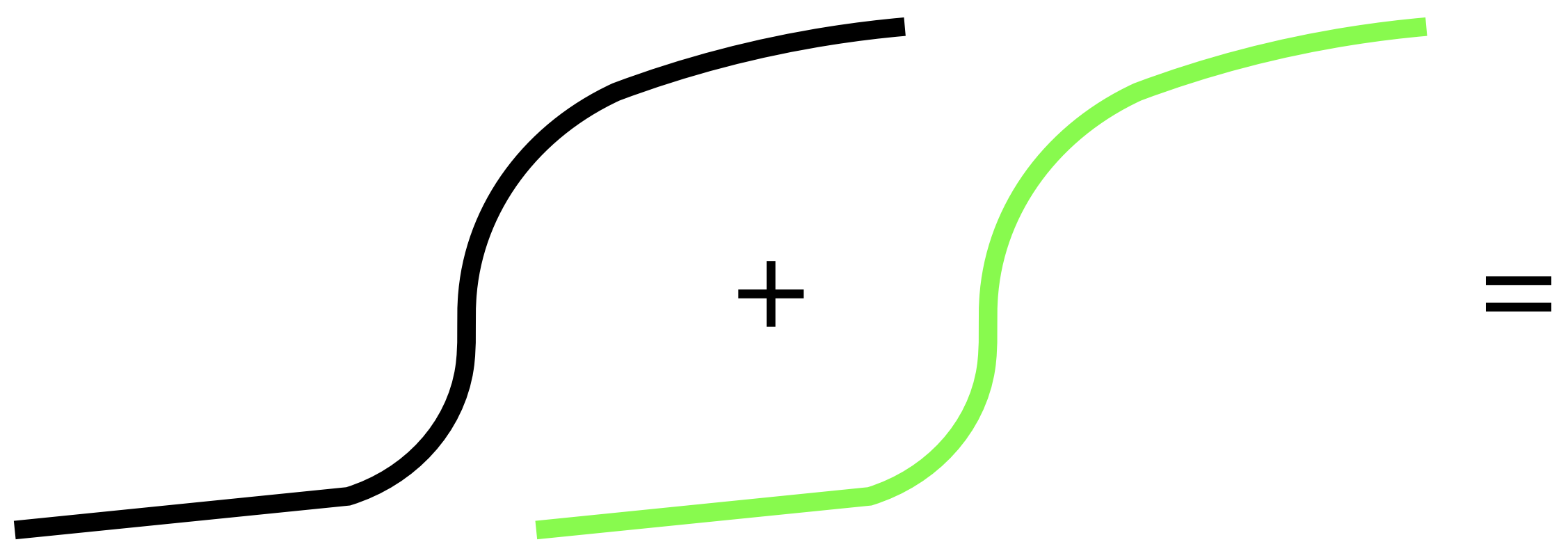


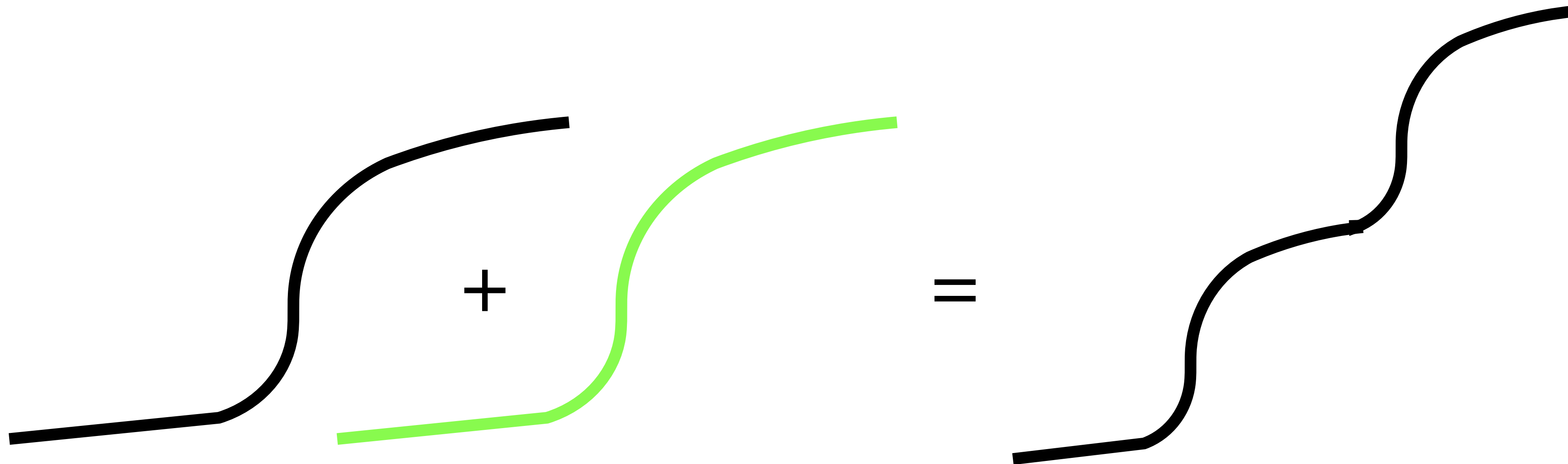


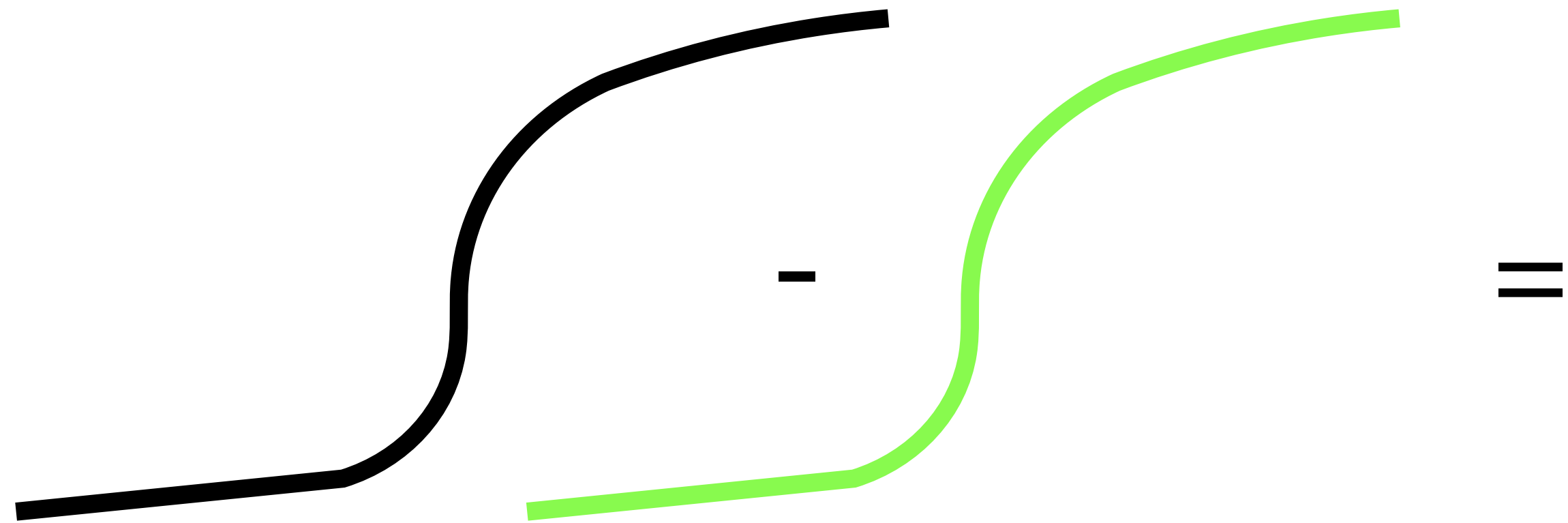
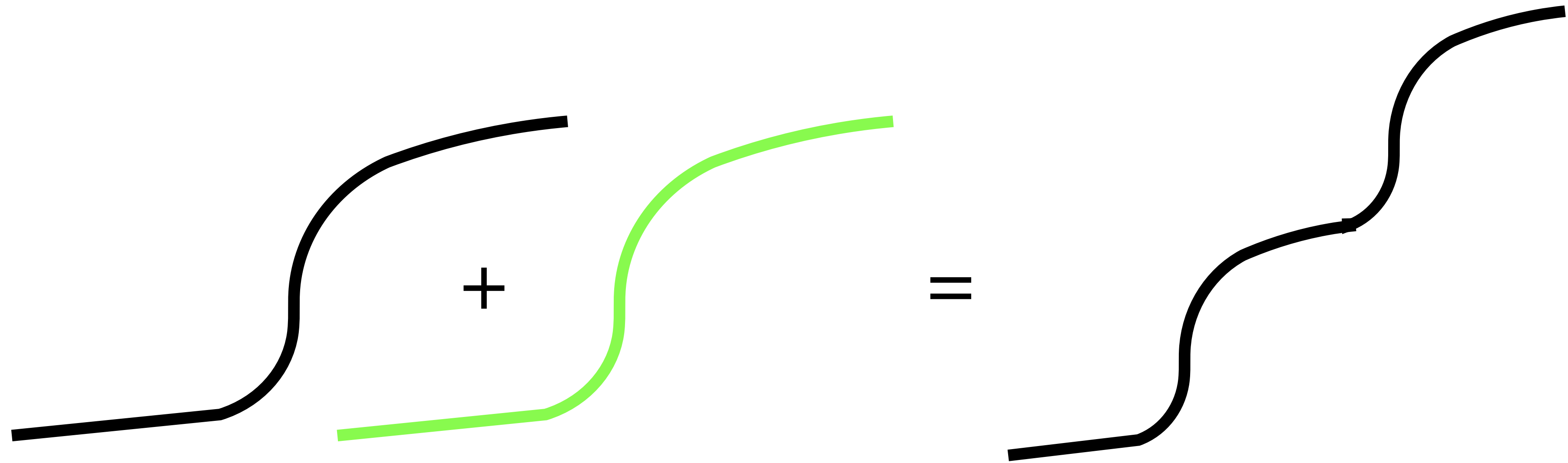


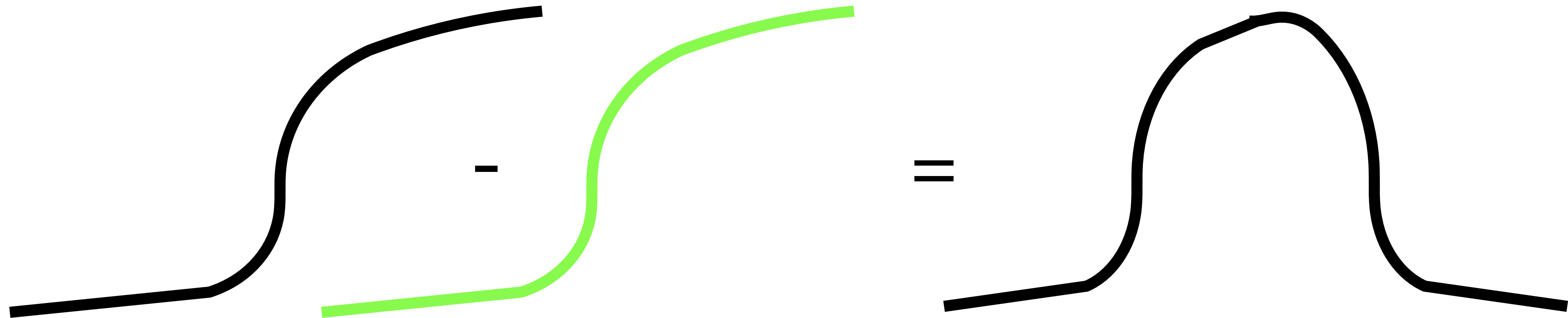
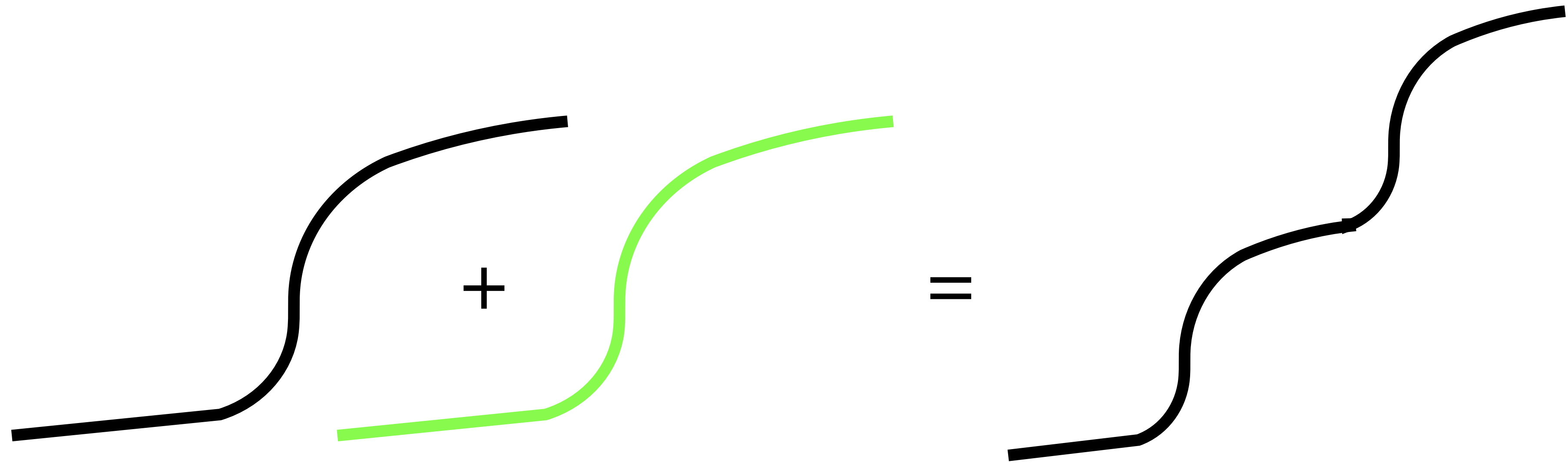


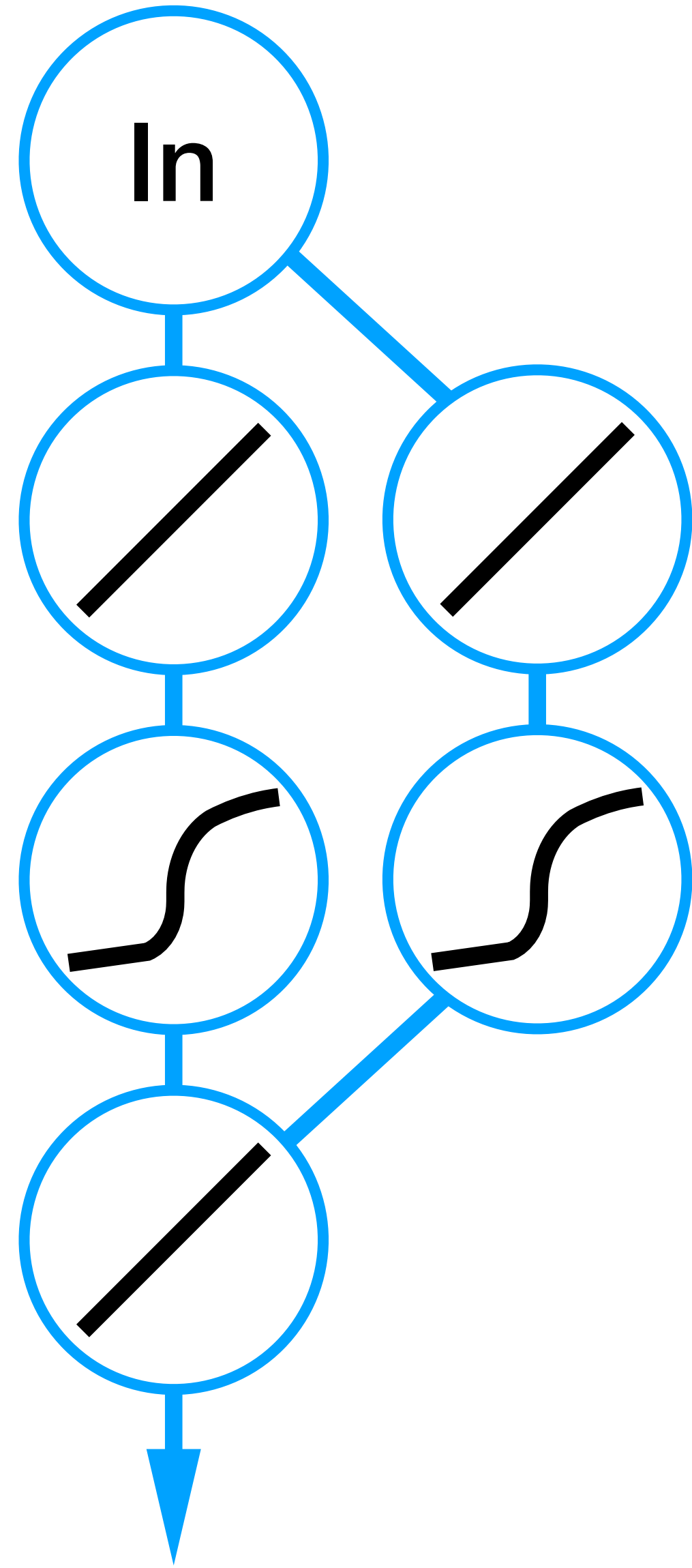


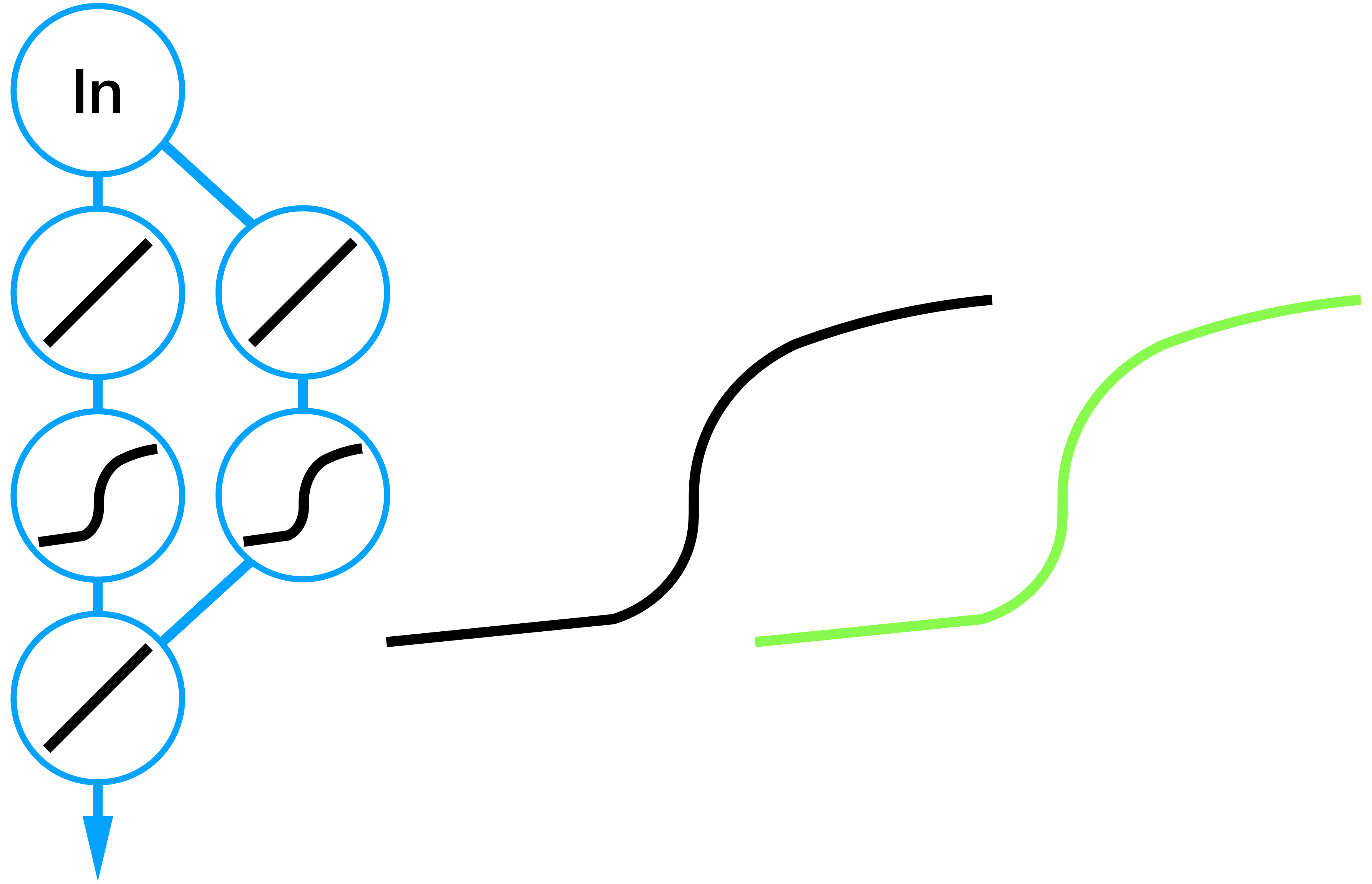


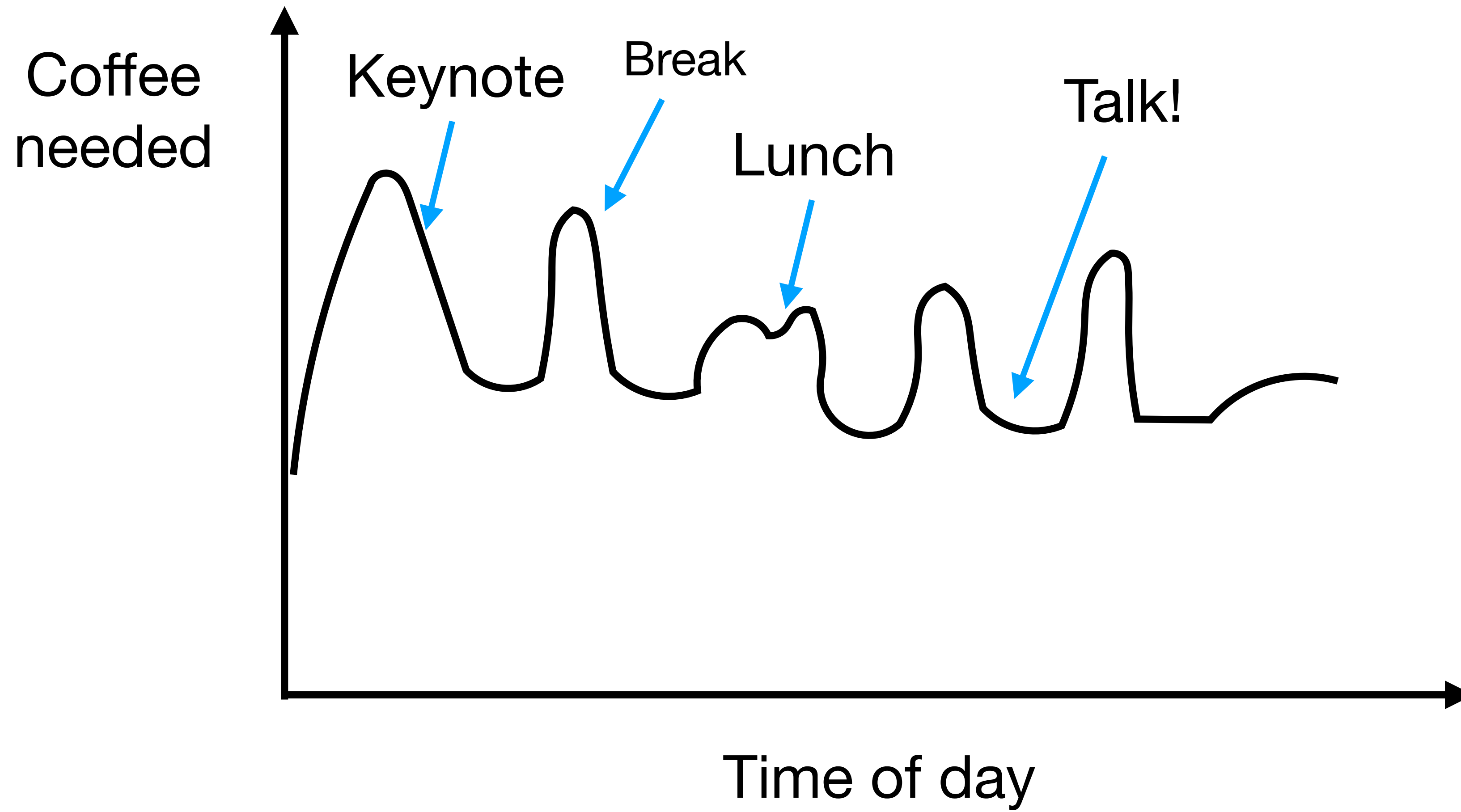


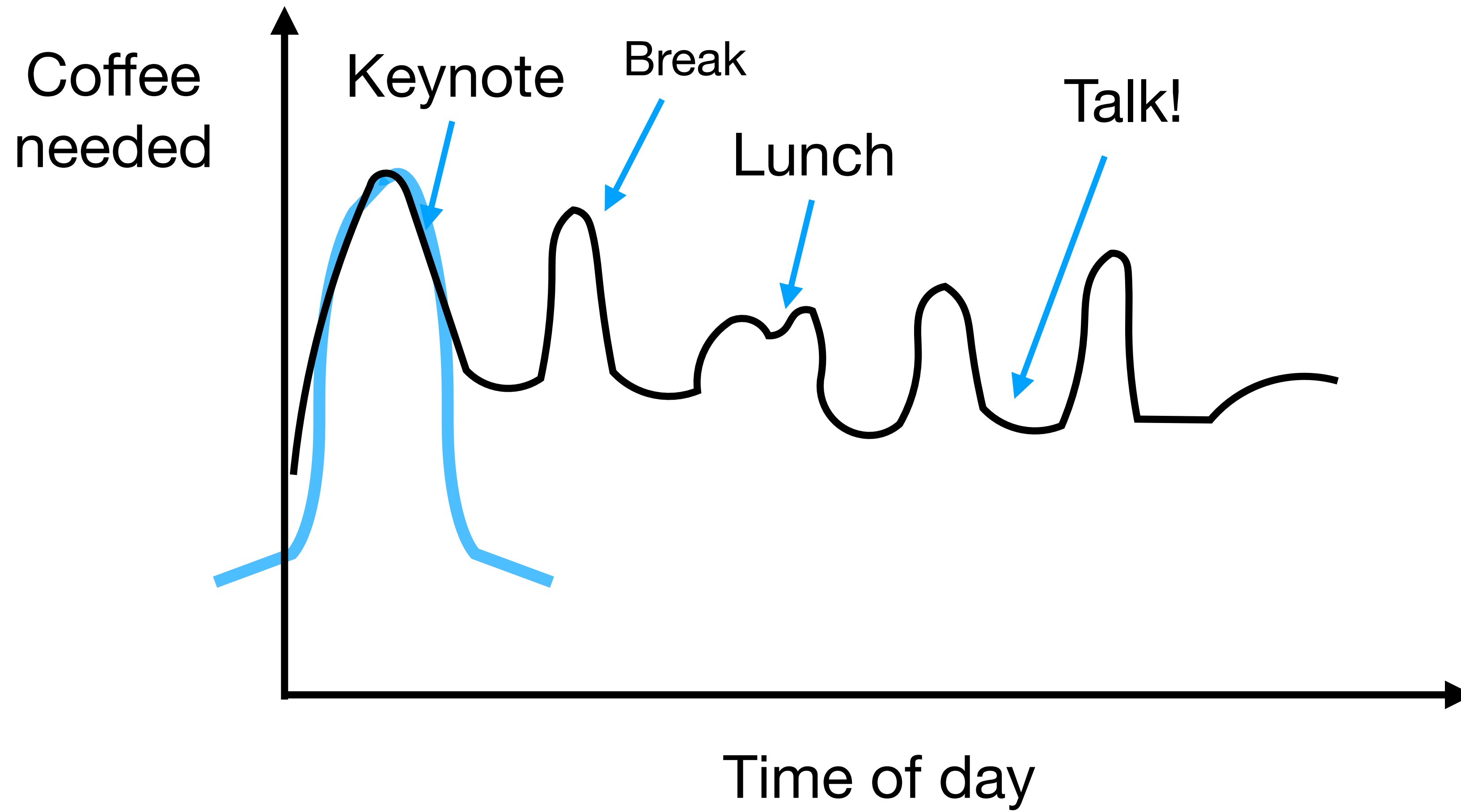


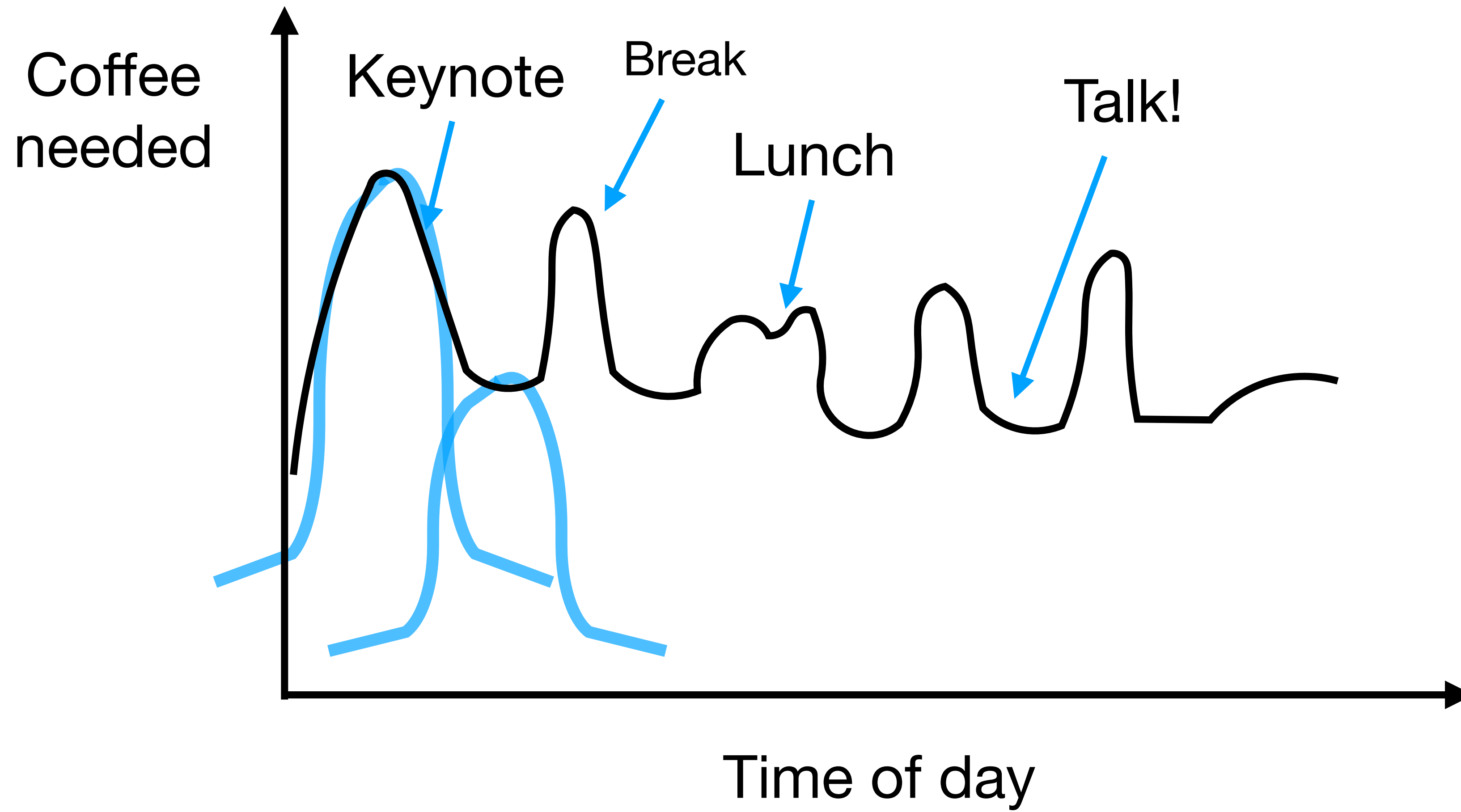


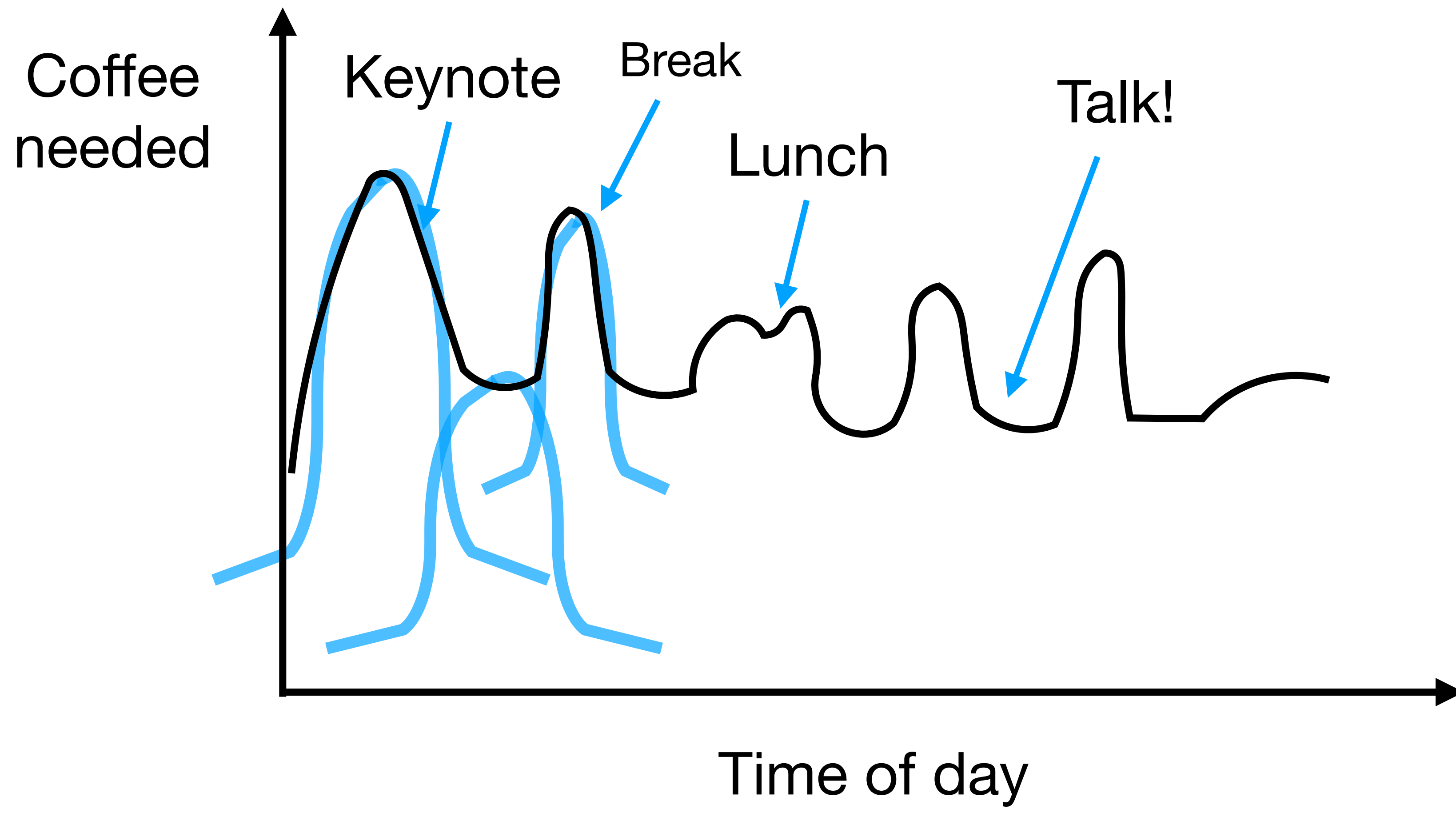


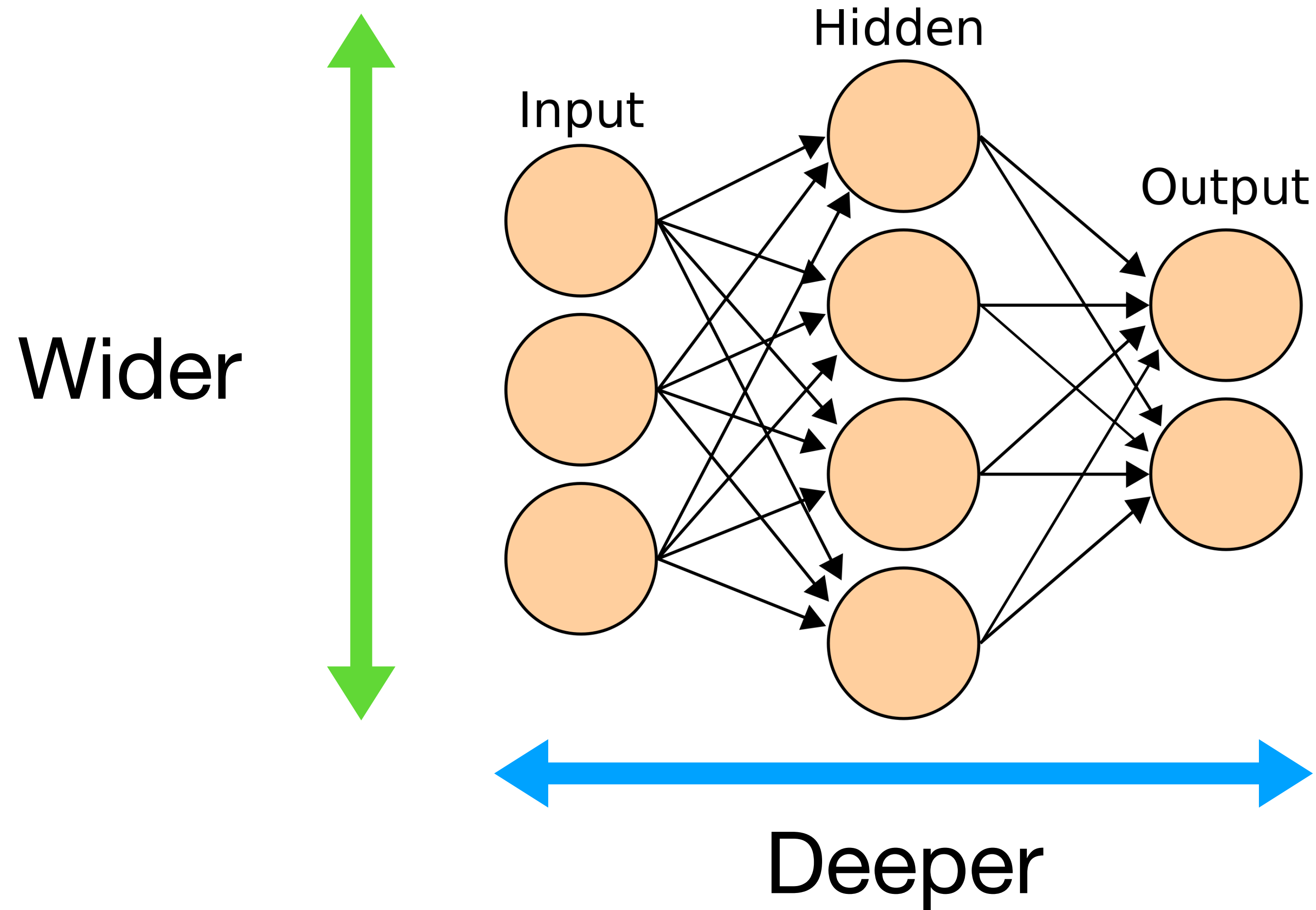








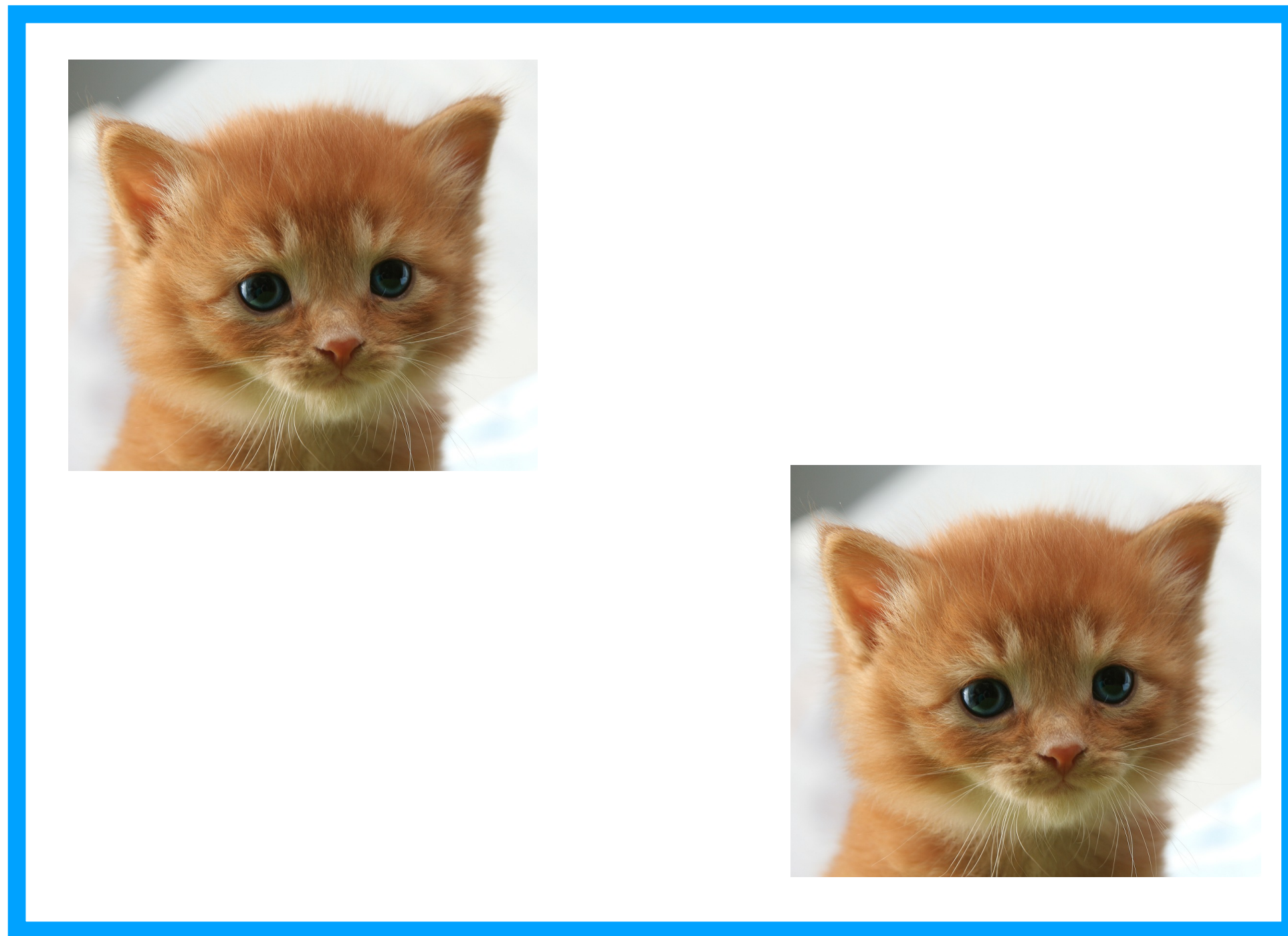




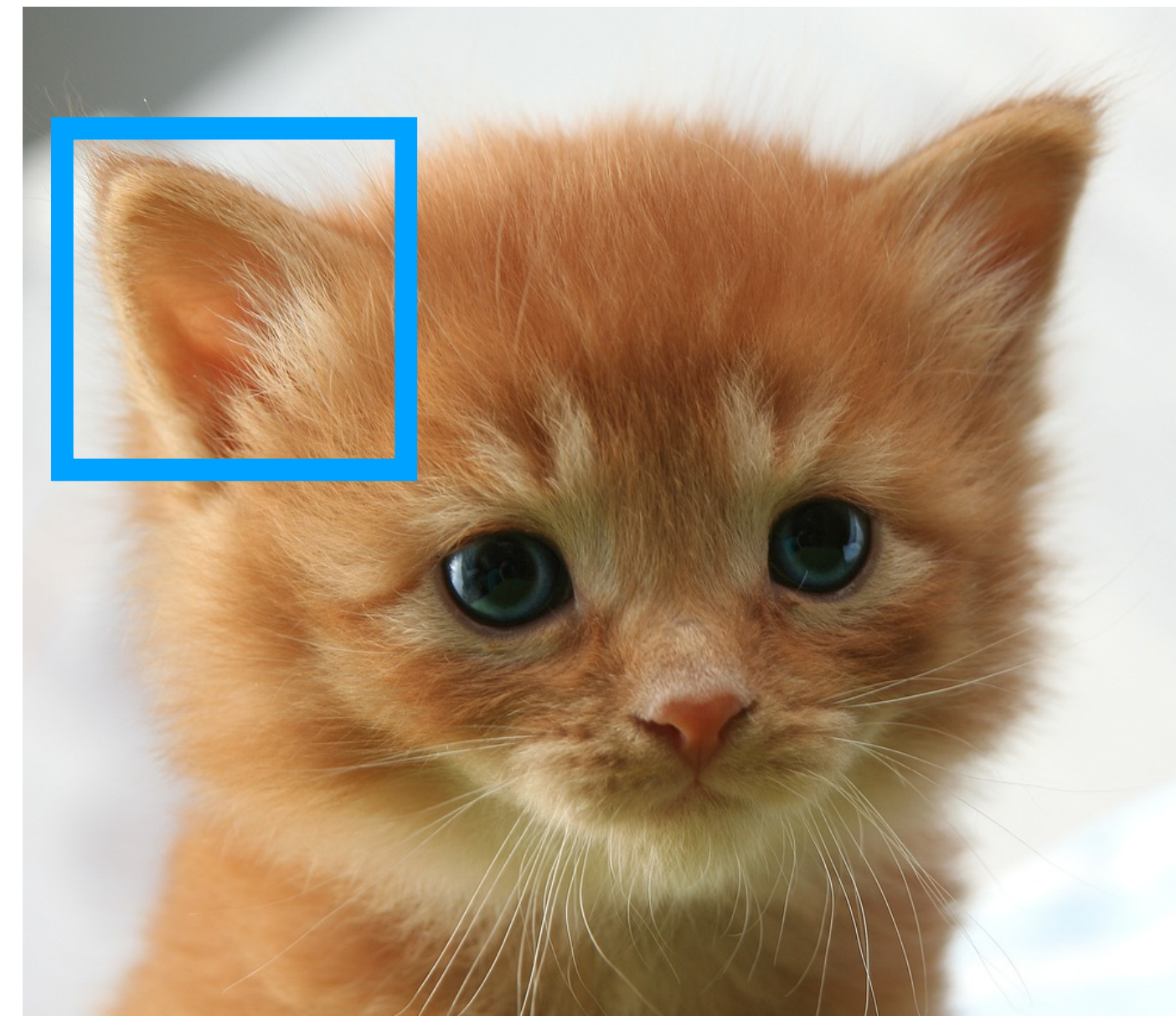
Foundational for CNNs and RNNs

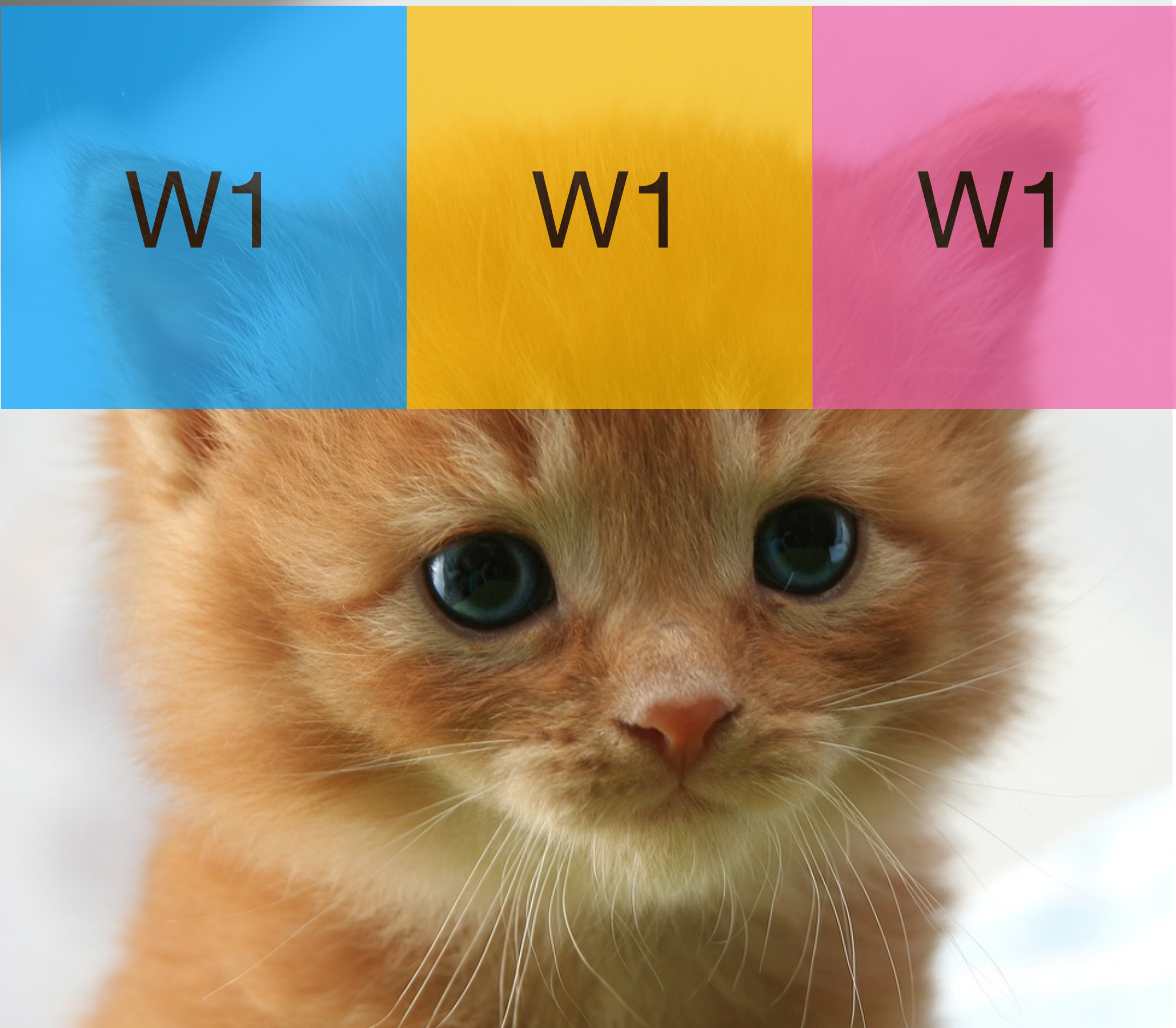
Convolutional layers

1. Translation invariance



2. Locality

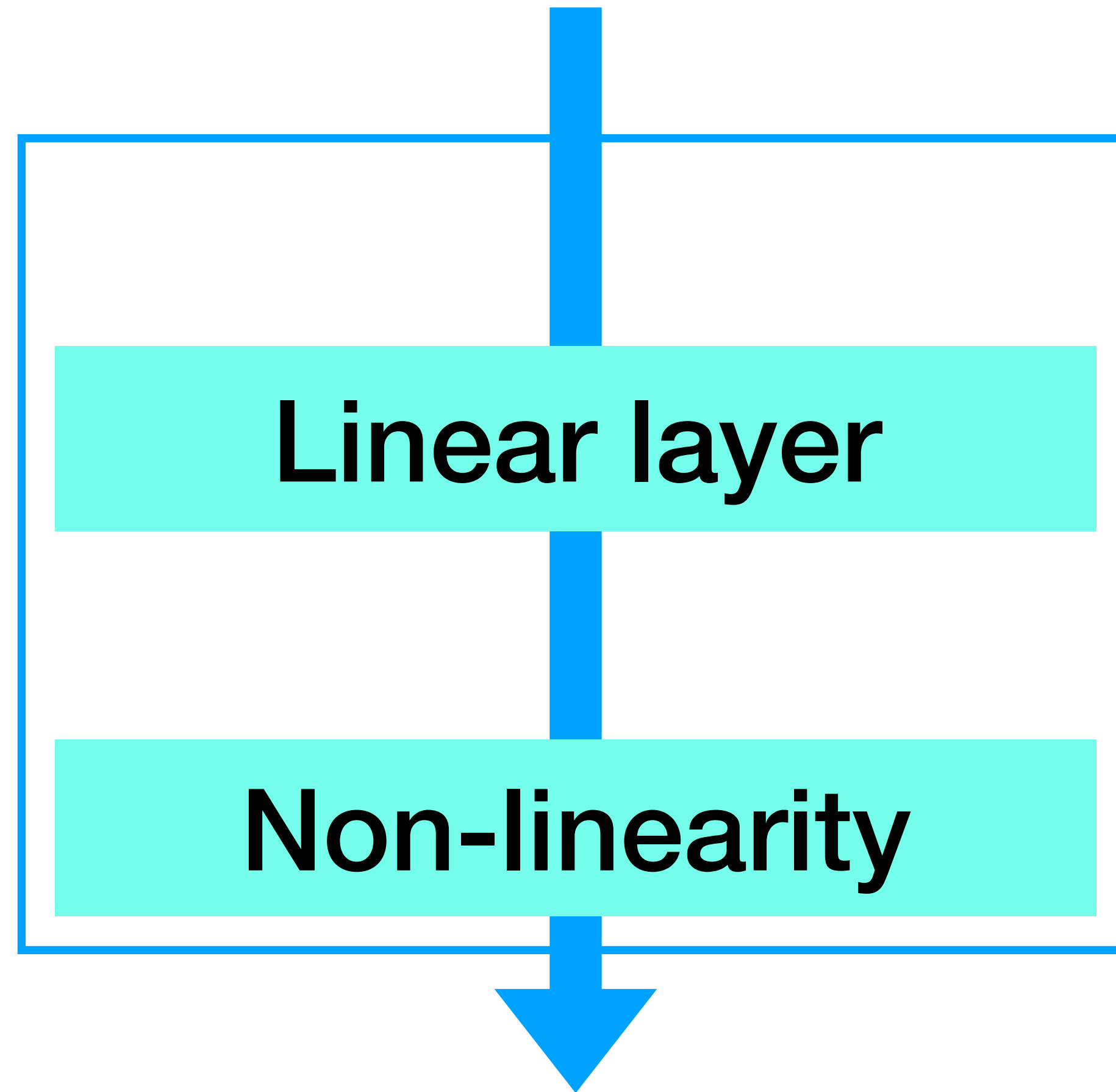




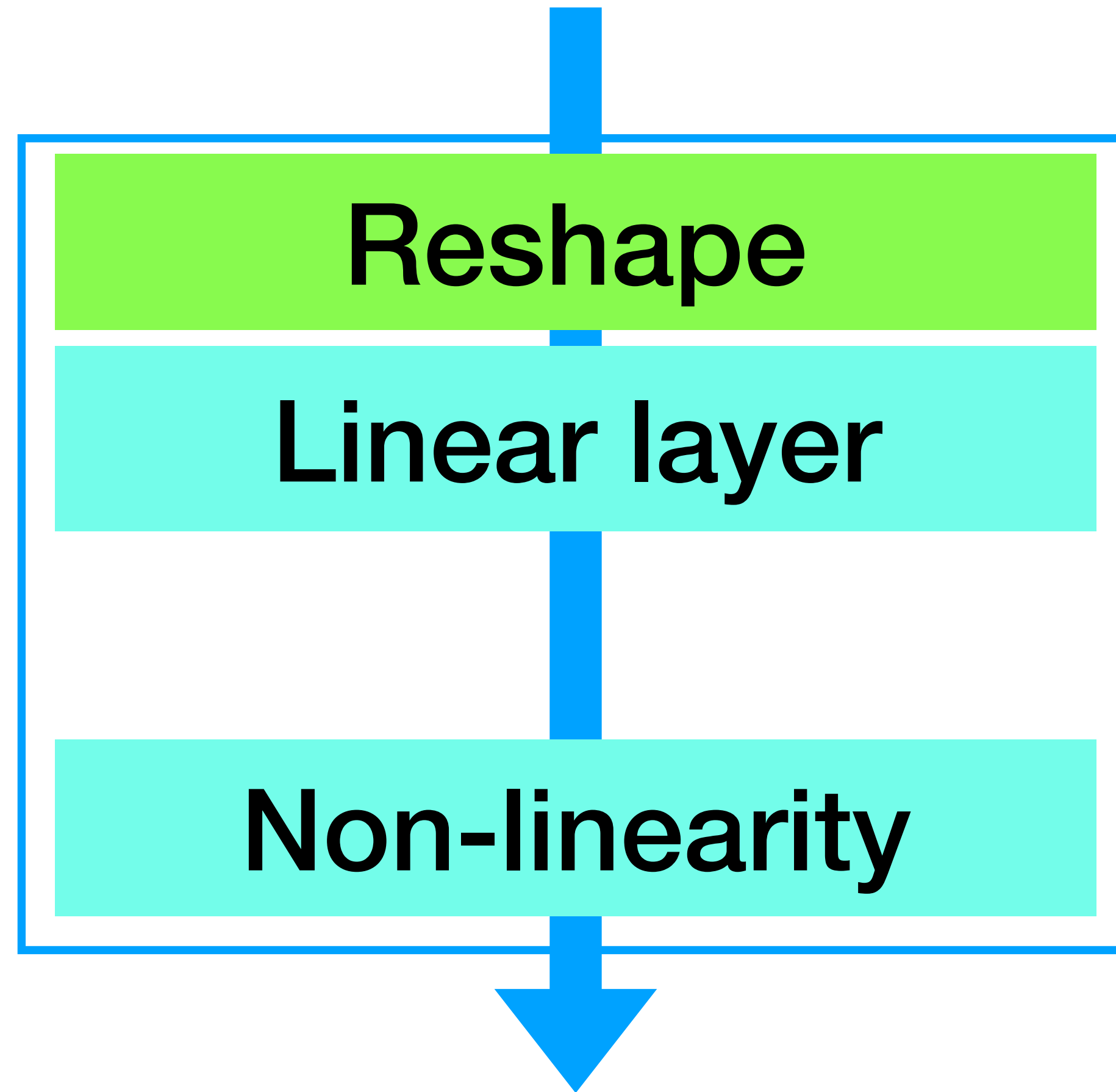
Output



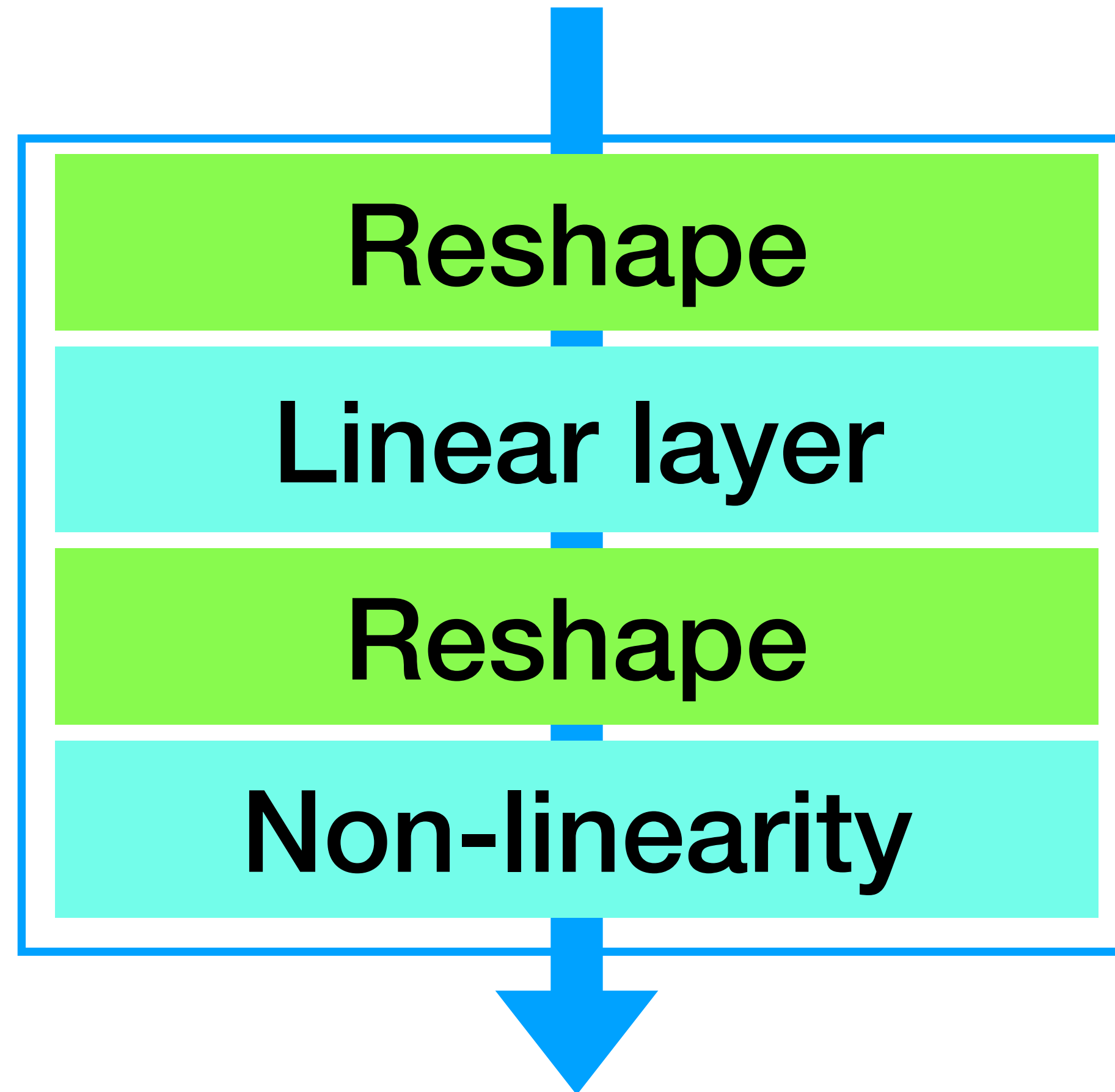
Convolutional layers



Convolutional layers

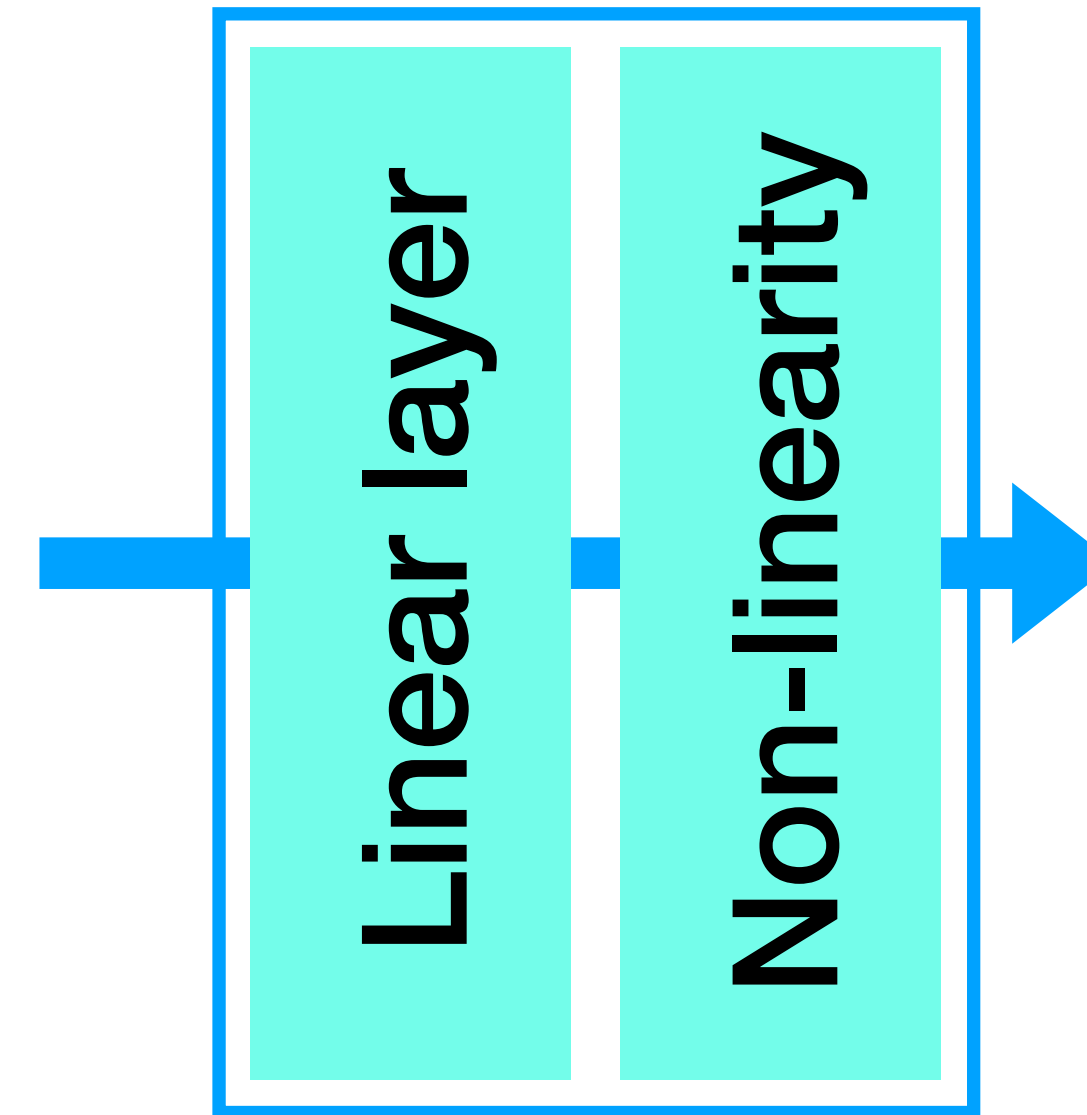


Convolutional layers

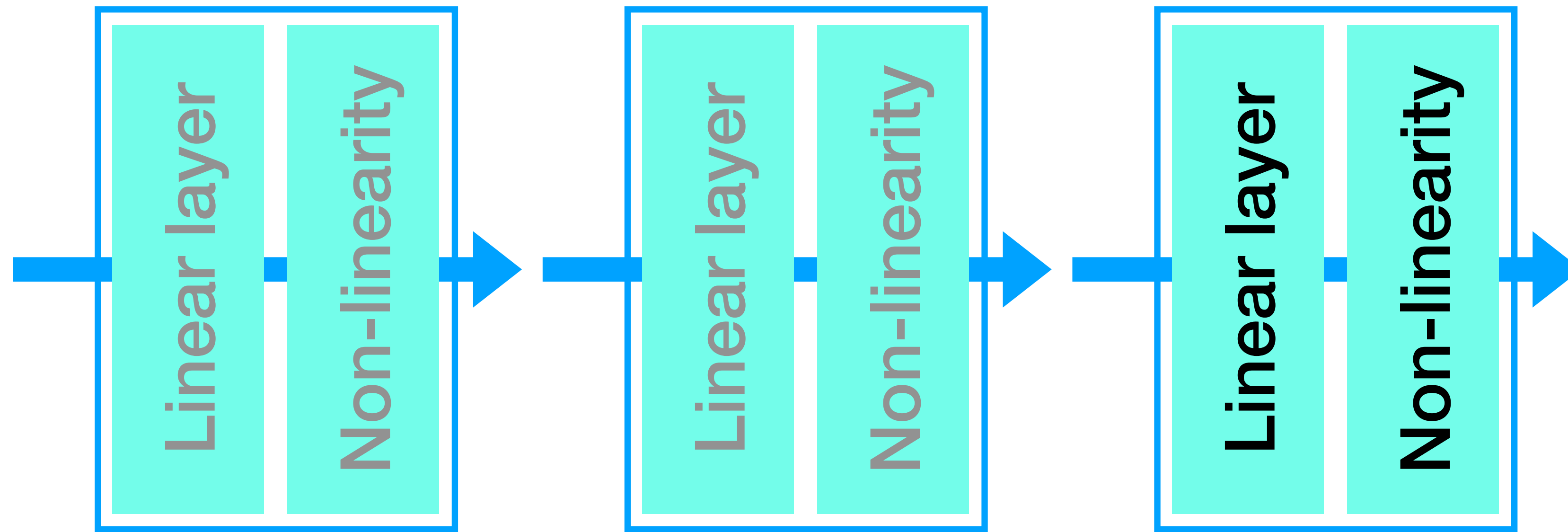


Recurrent layers

Recurrent layers

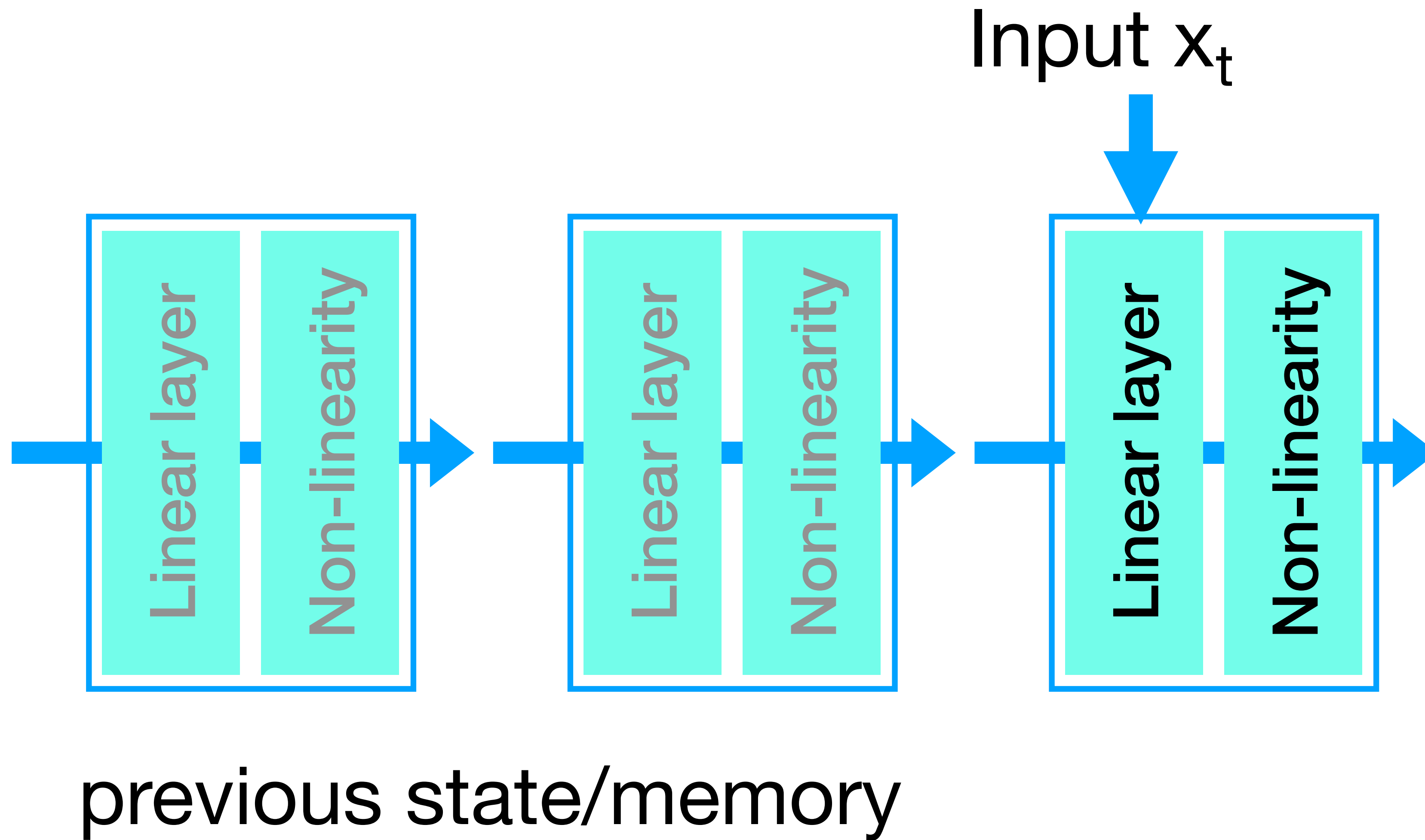


Recurrent layers

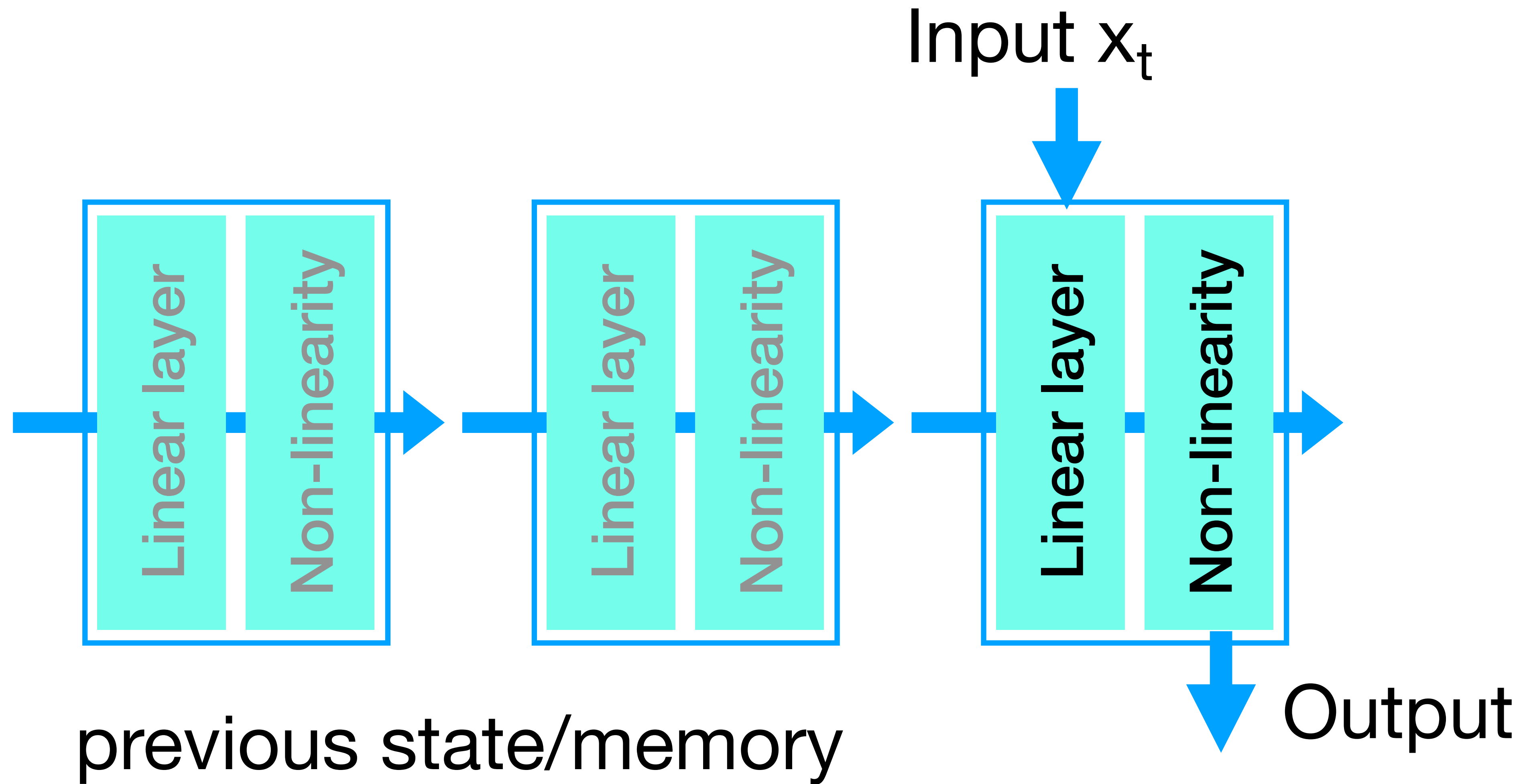


previous state/memory

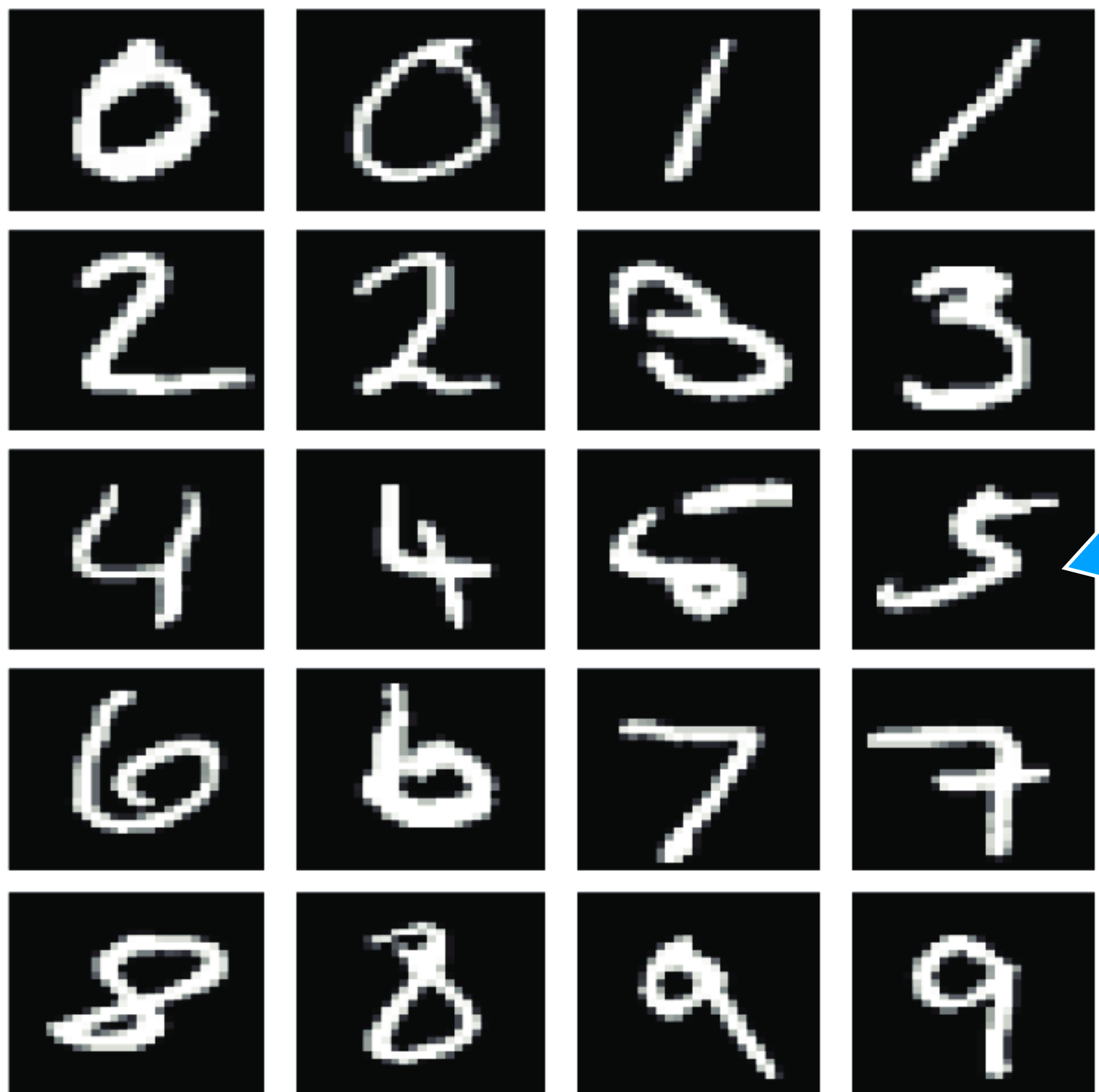
Recurrent layers



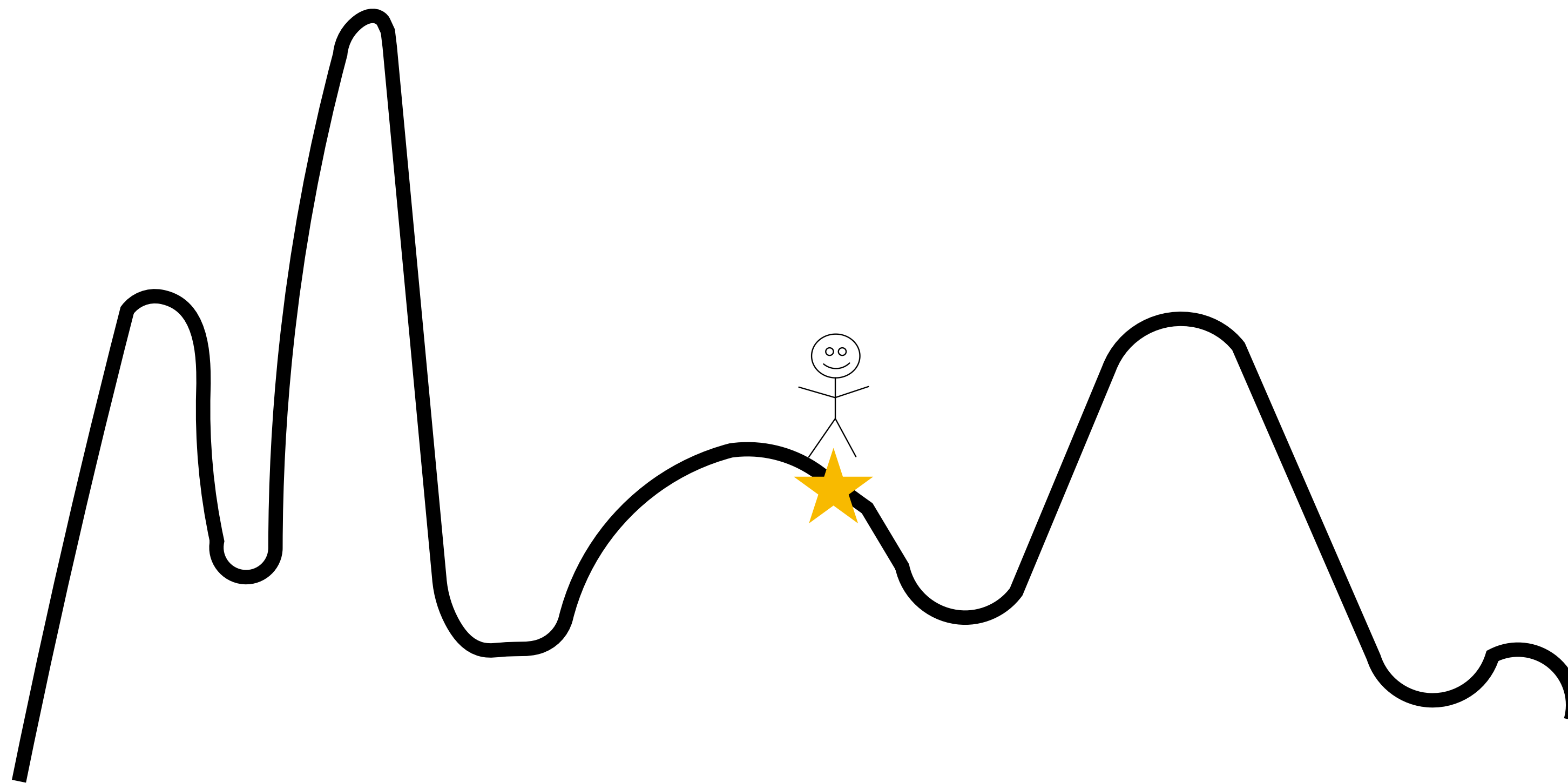
Recurrent layers

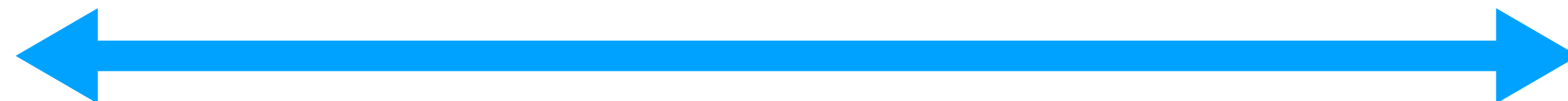
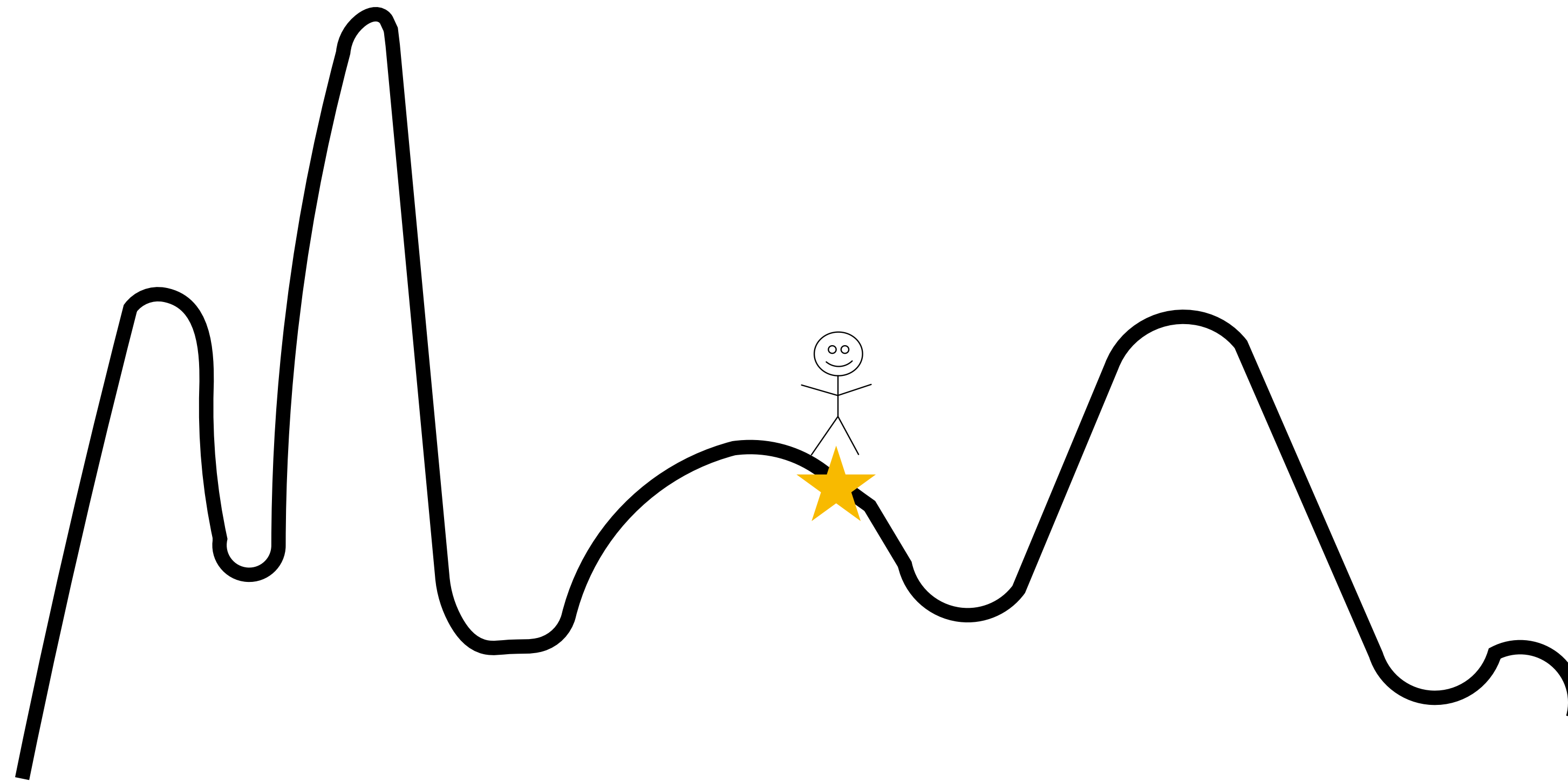


II. How train?

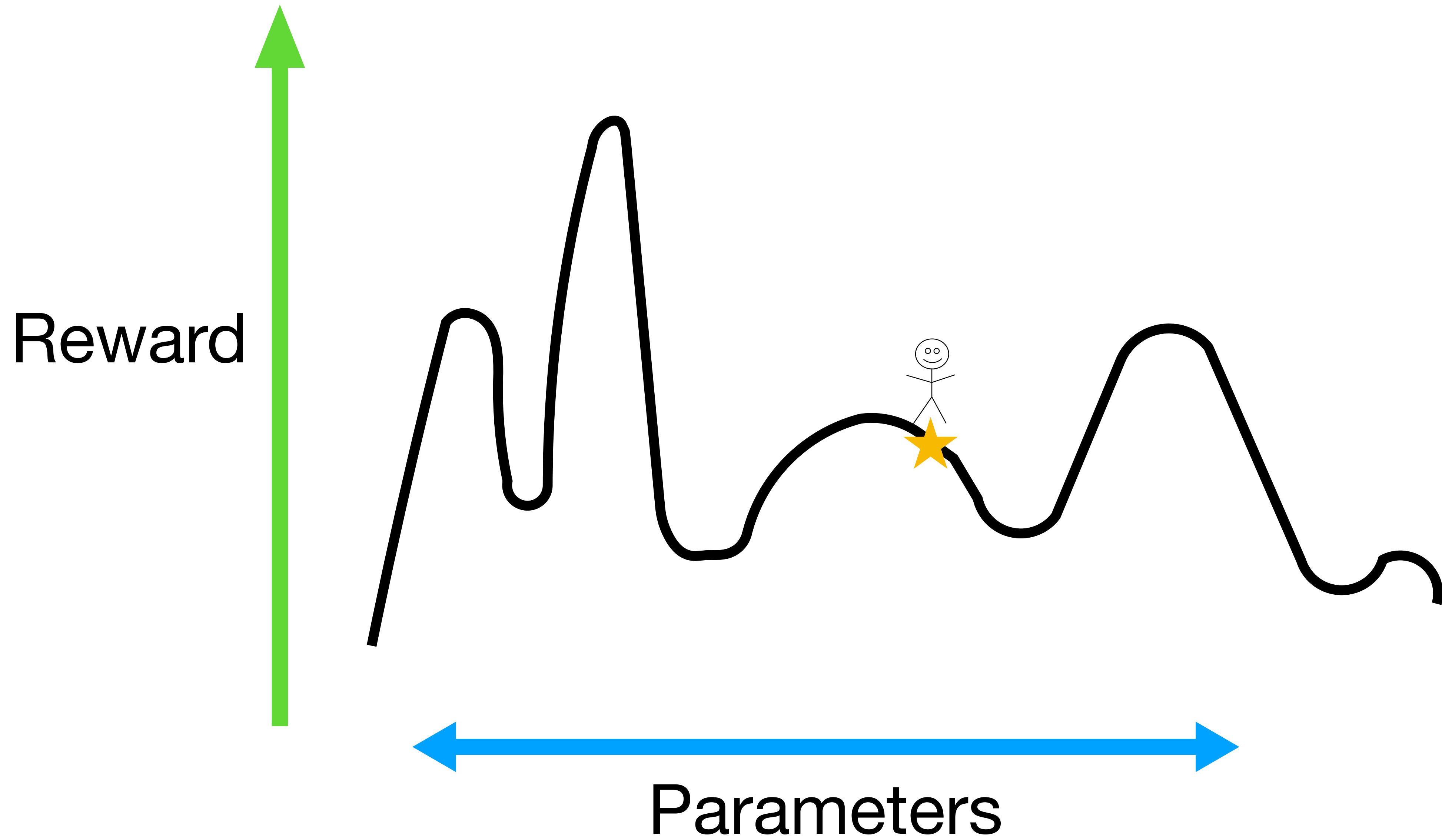


Metric:
~Prediction error

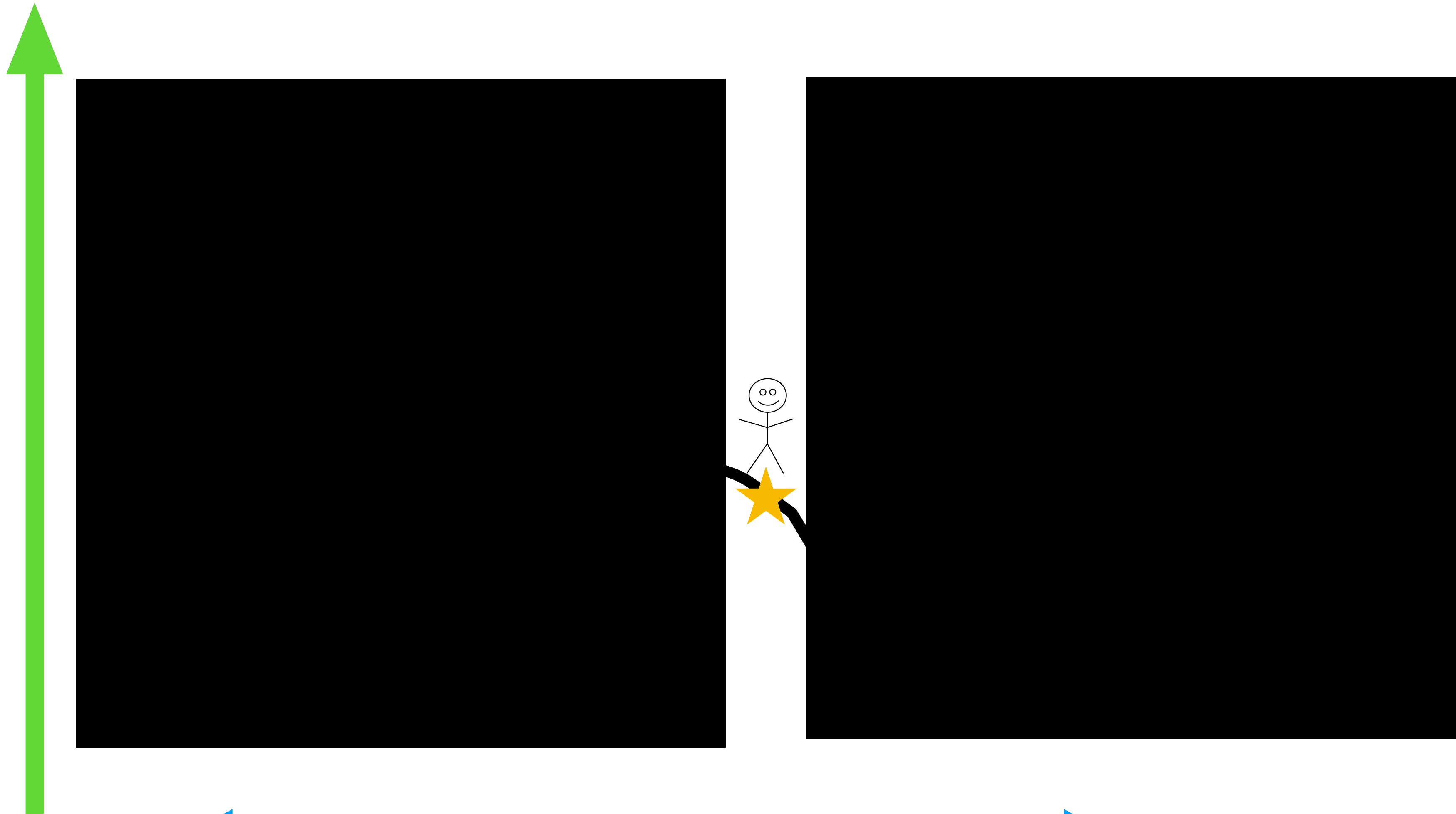




Parameters

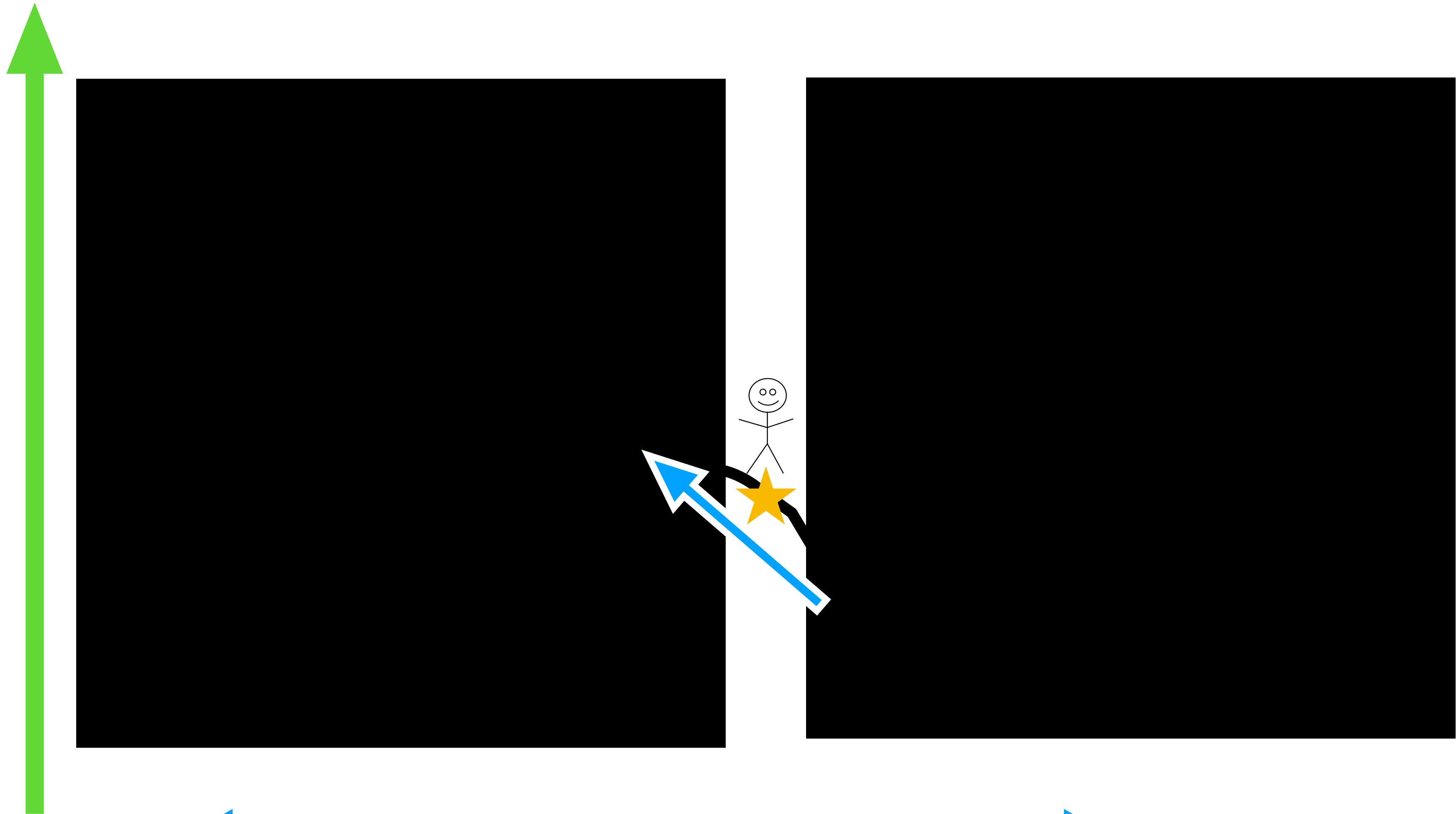


Reward

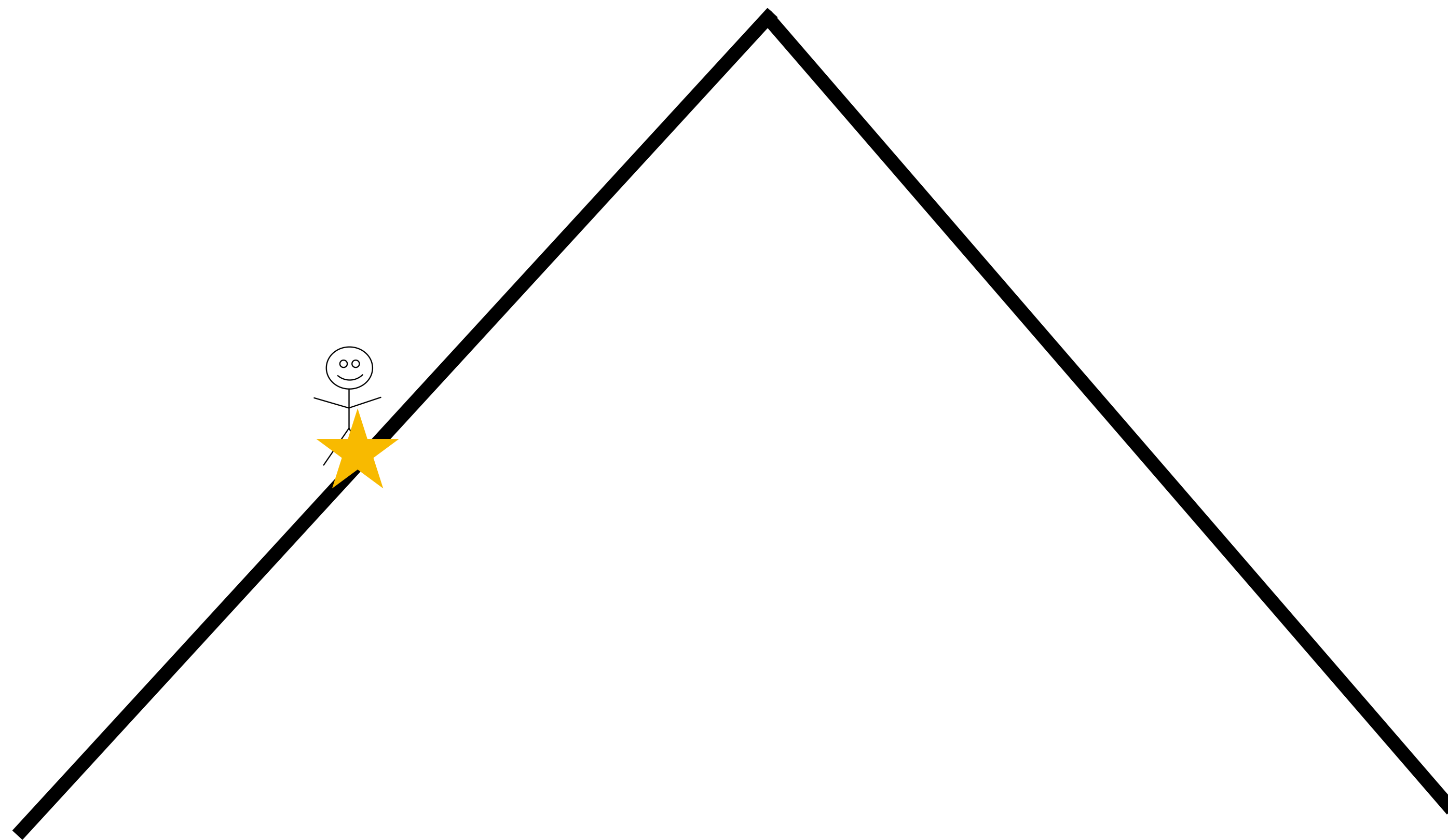


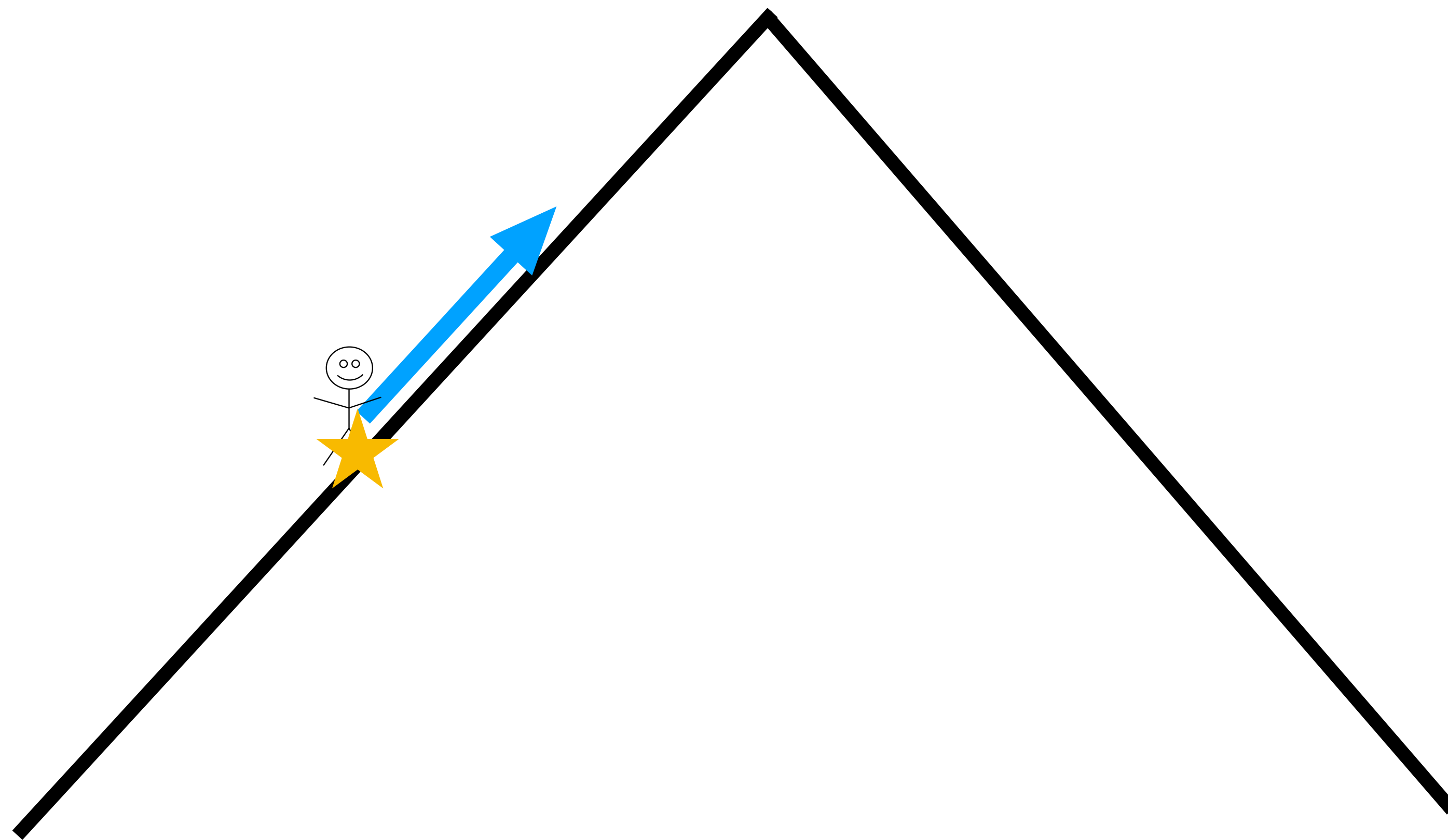
Parameters

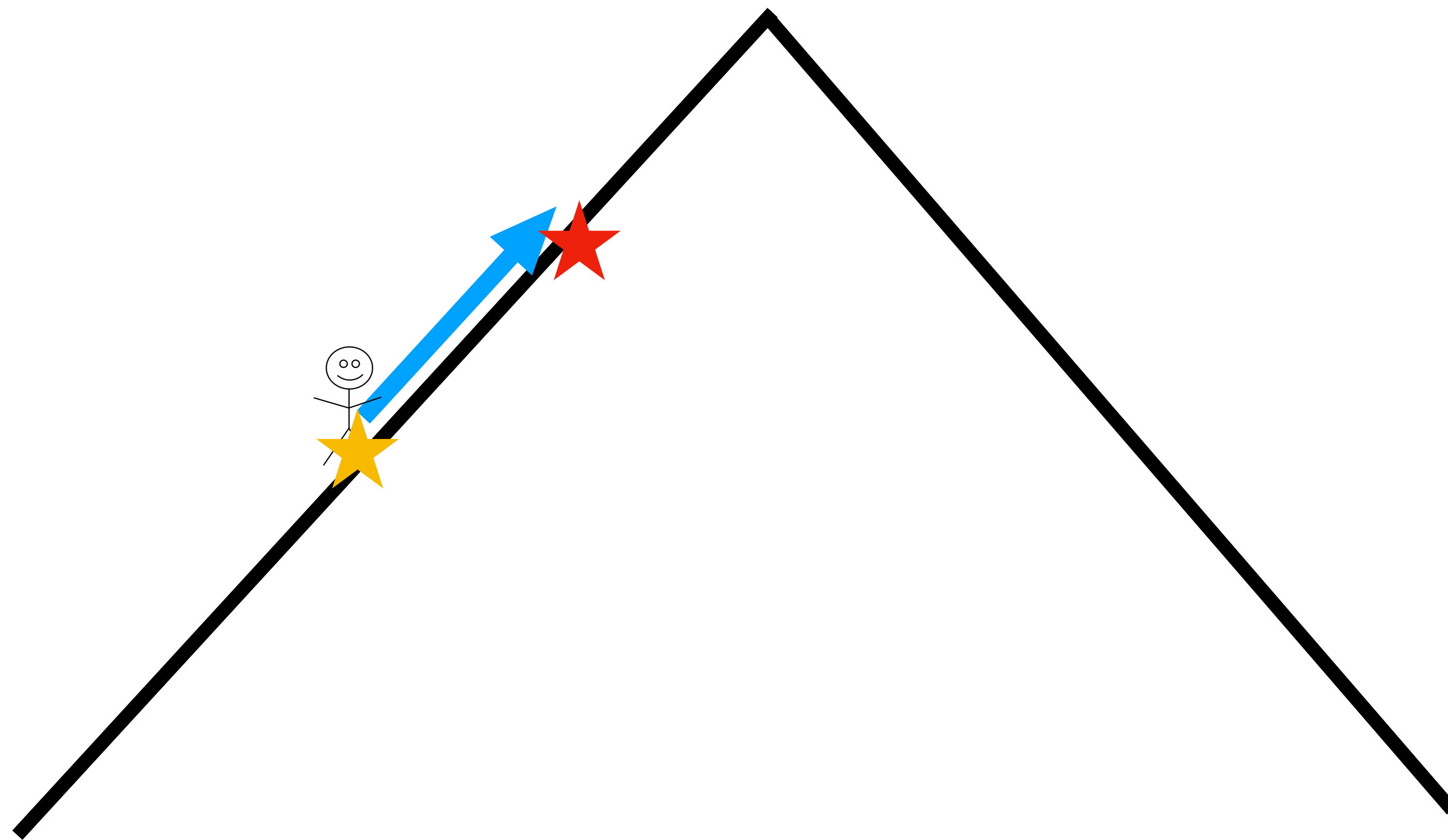
Reward

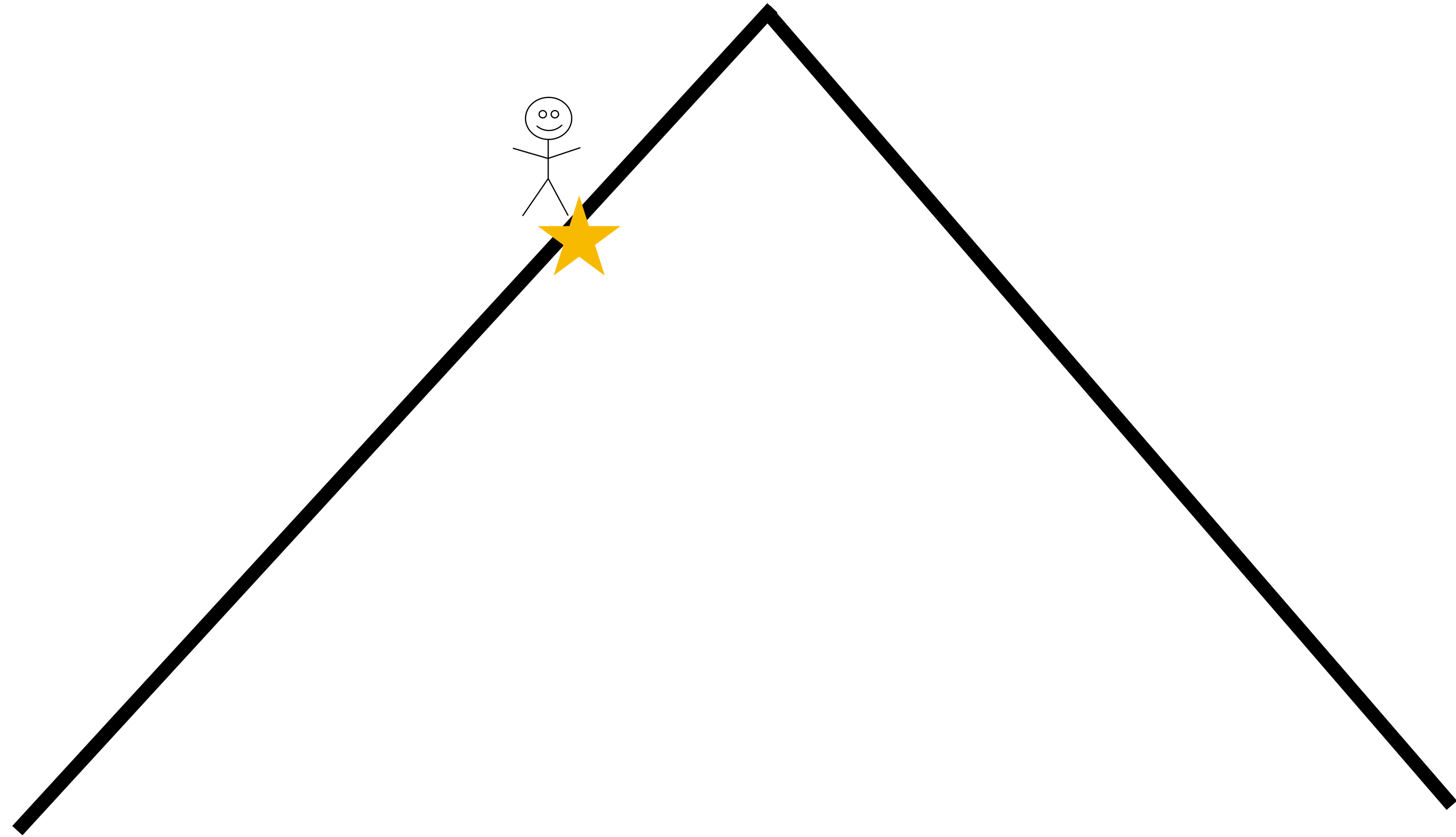


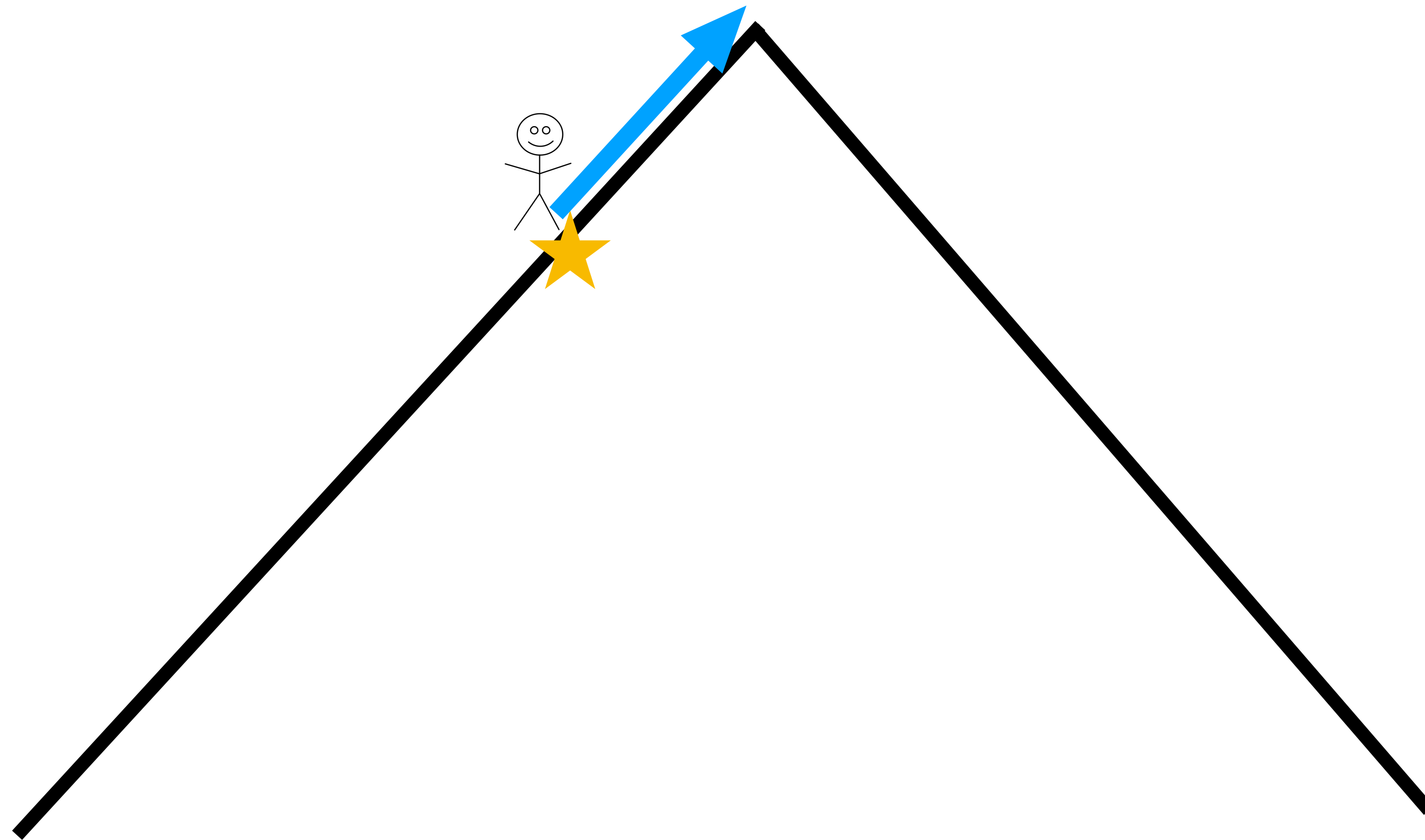
Parameters

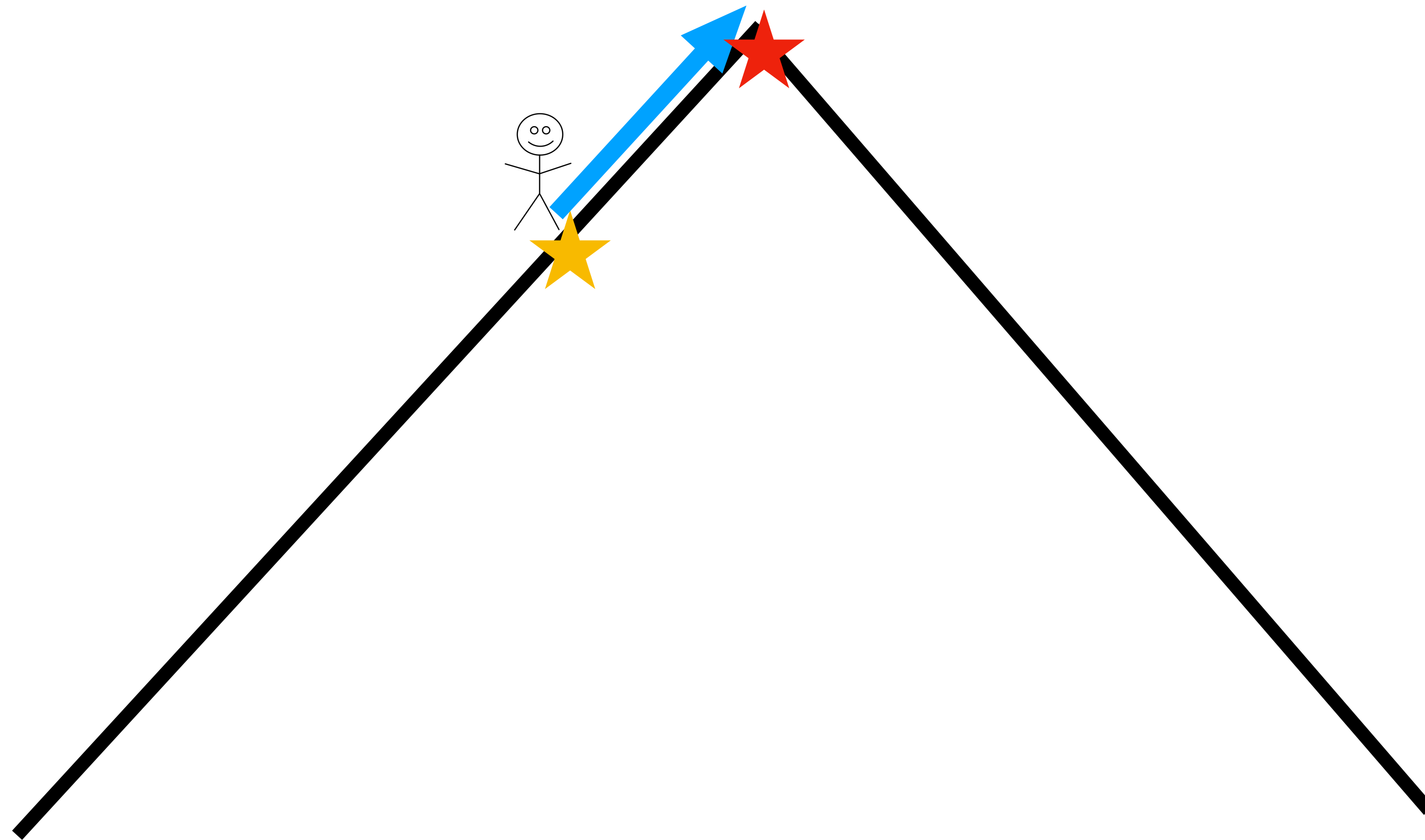


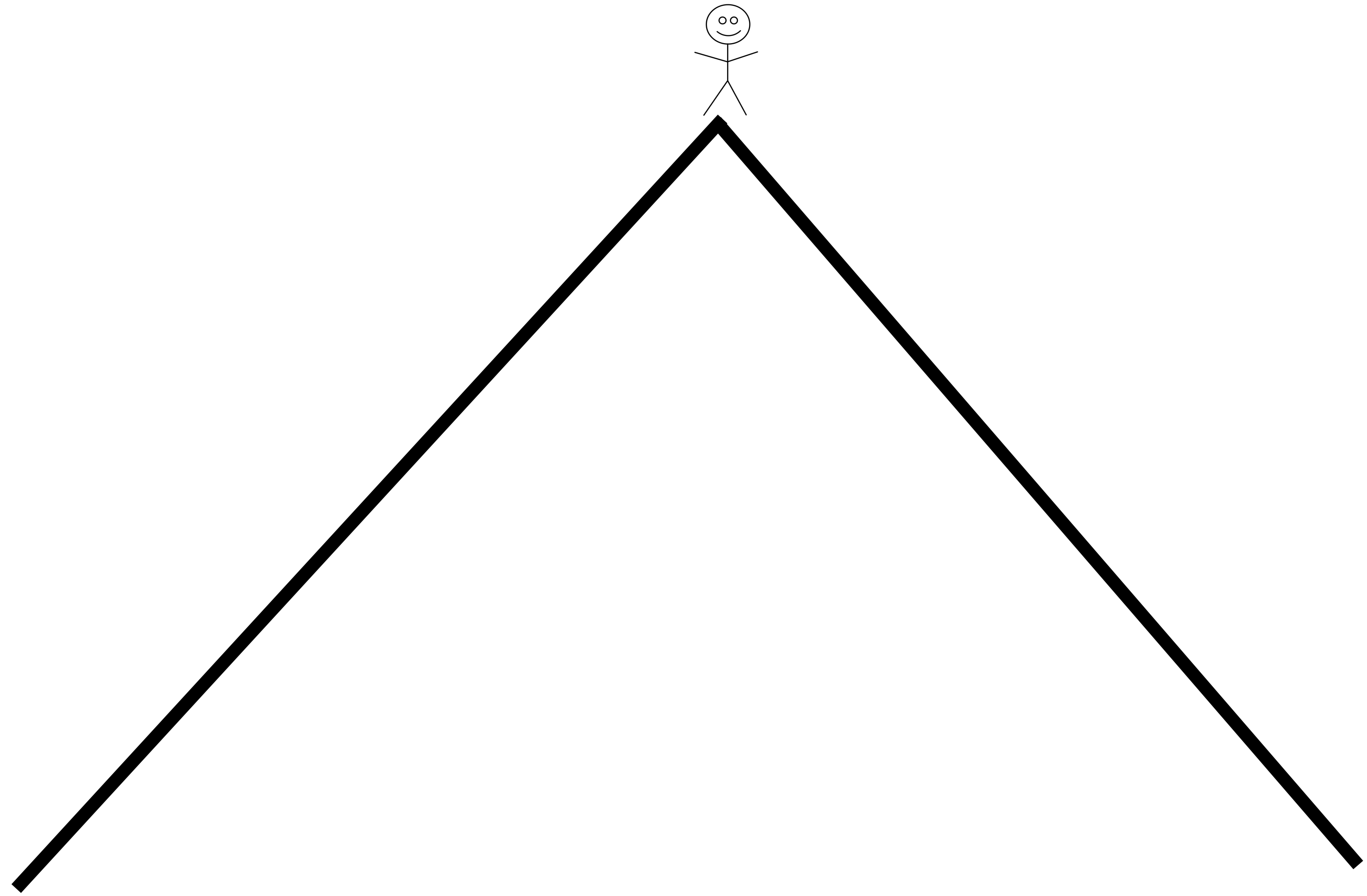


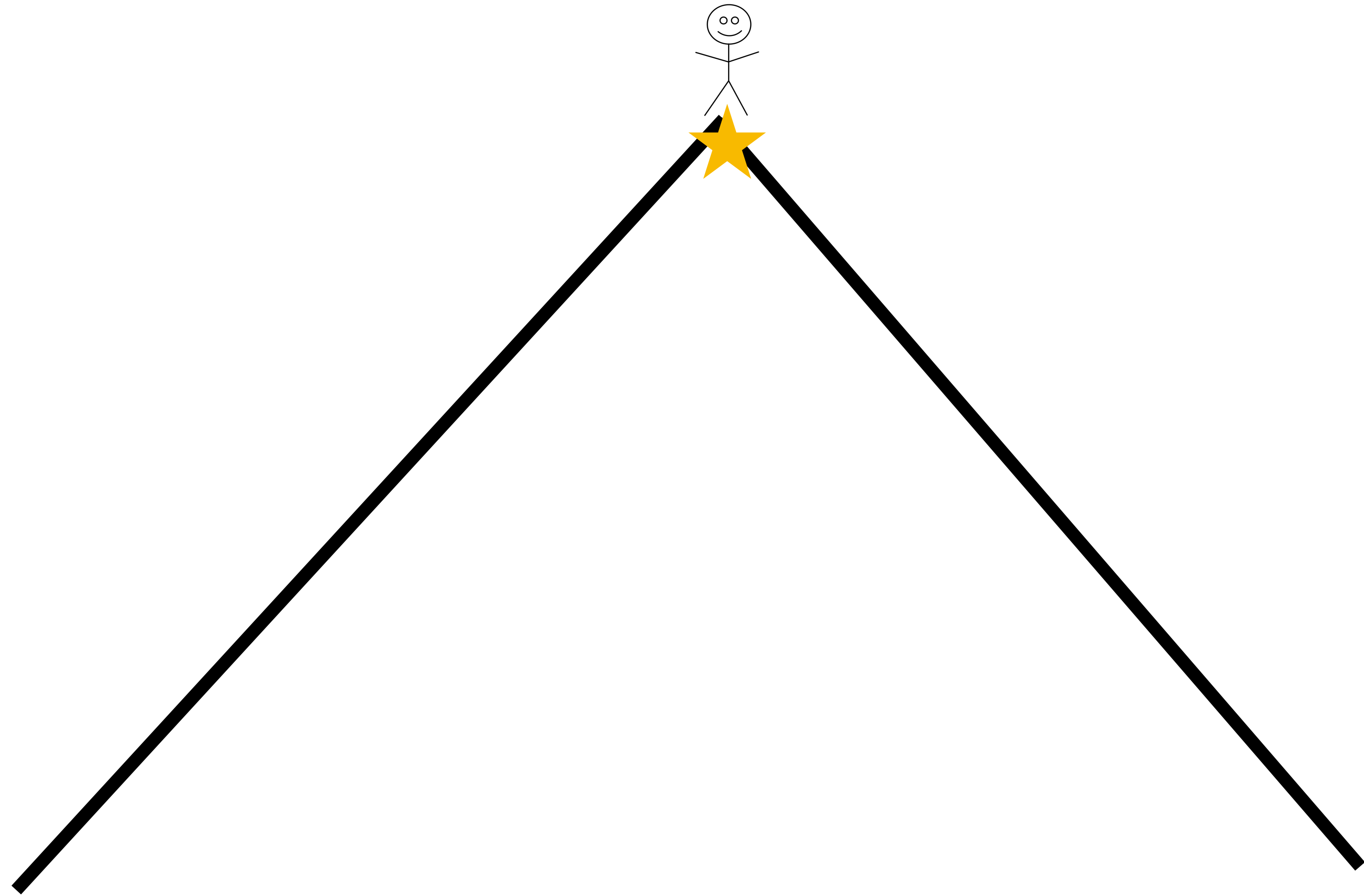


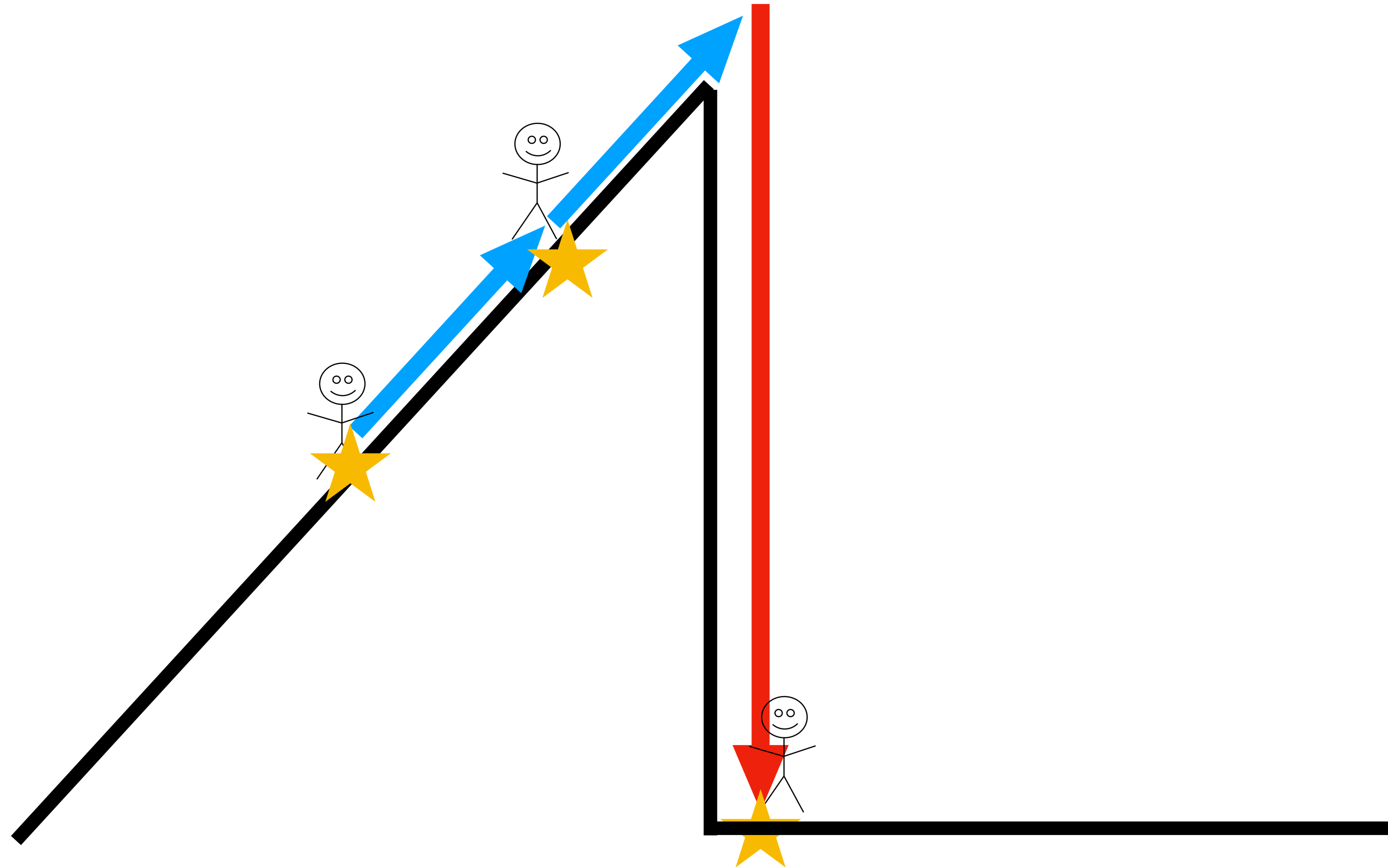


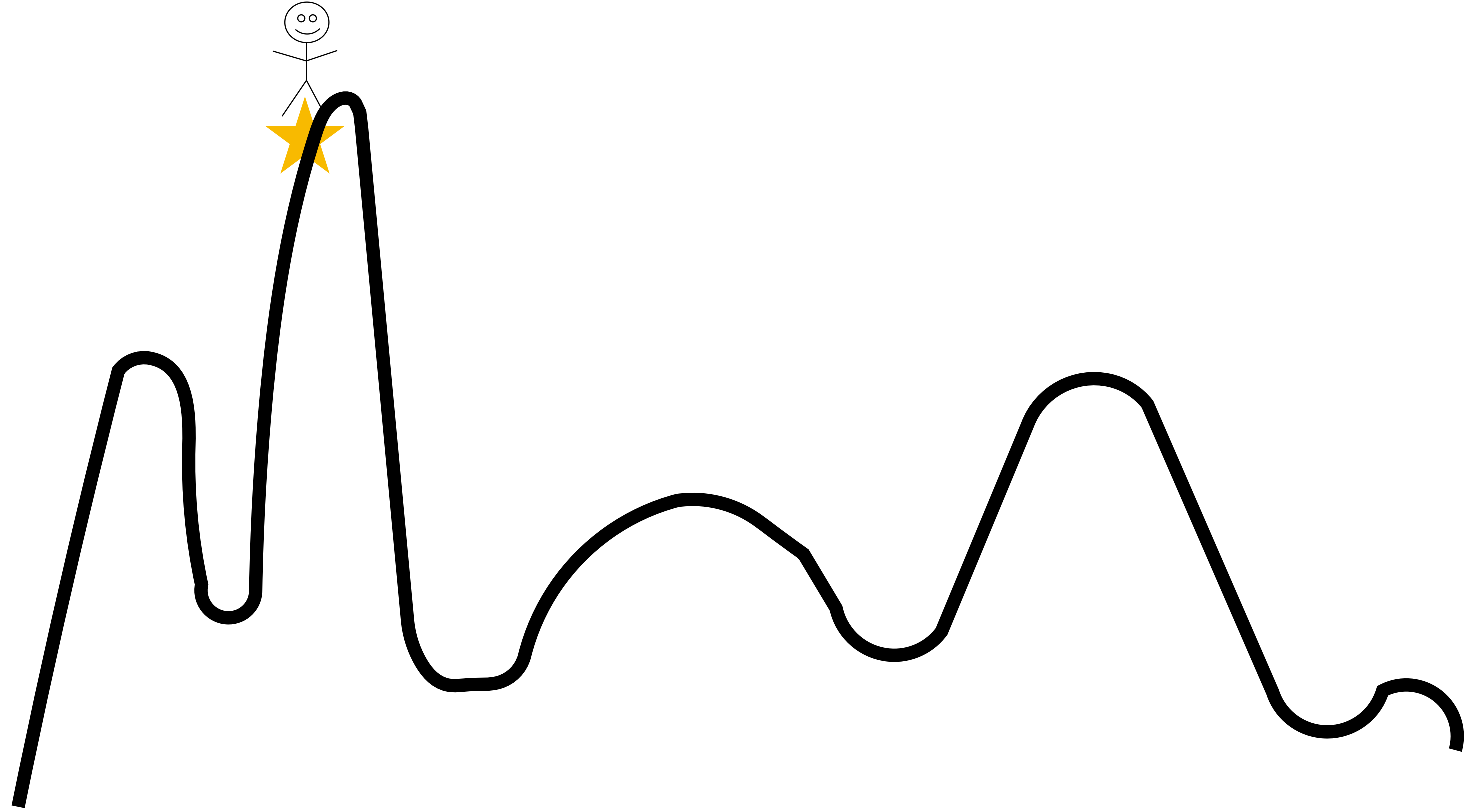


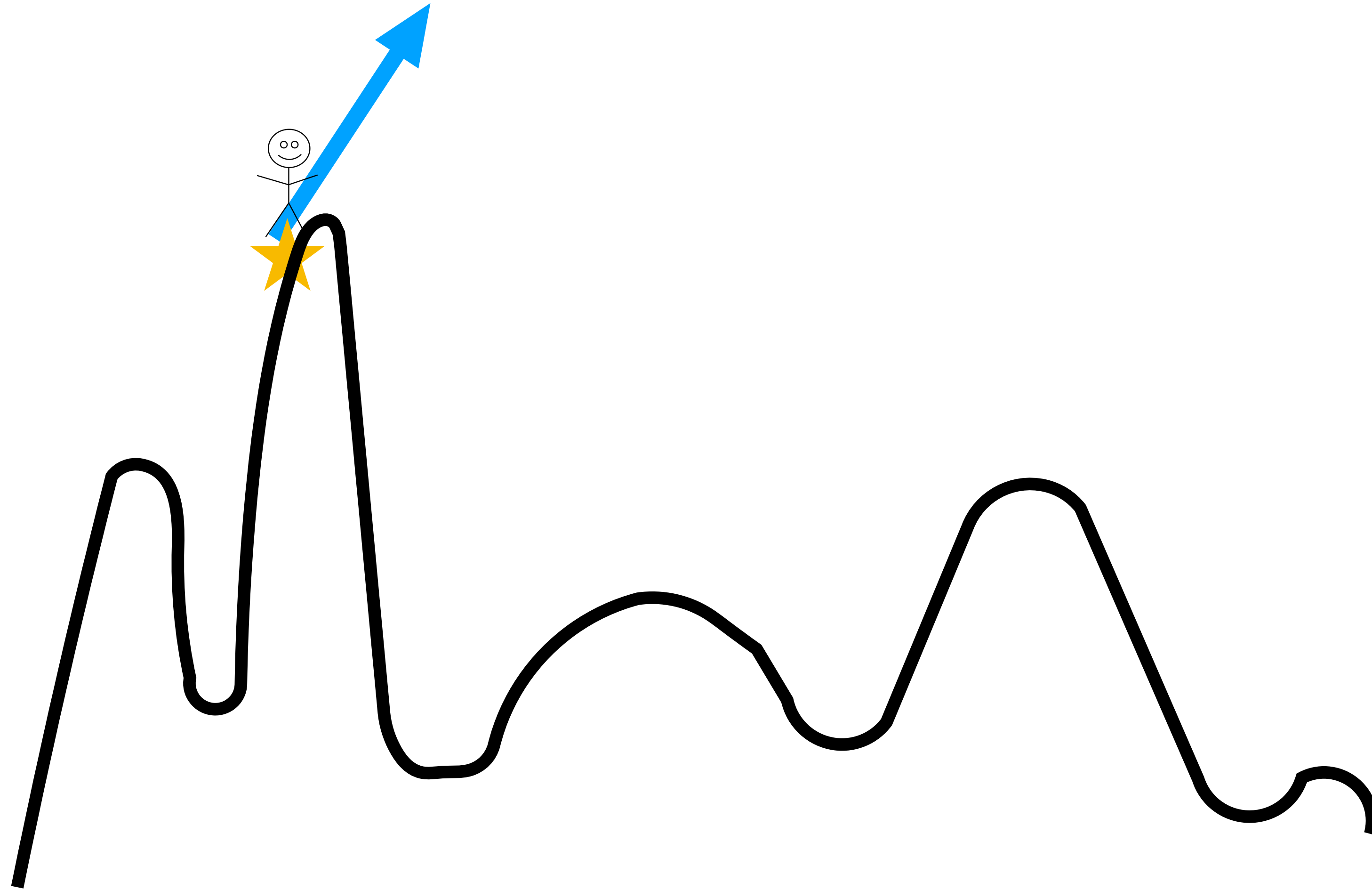


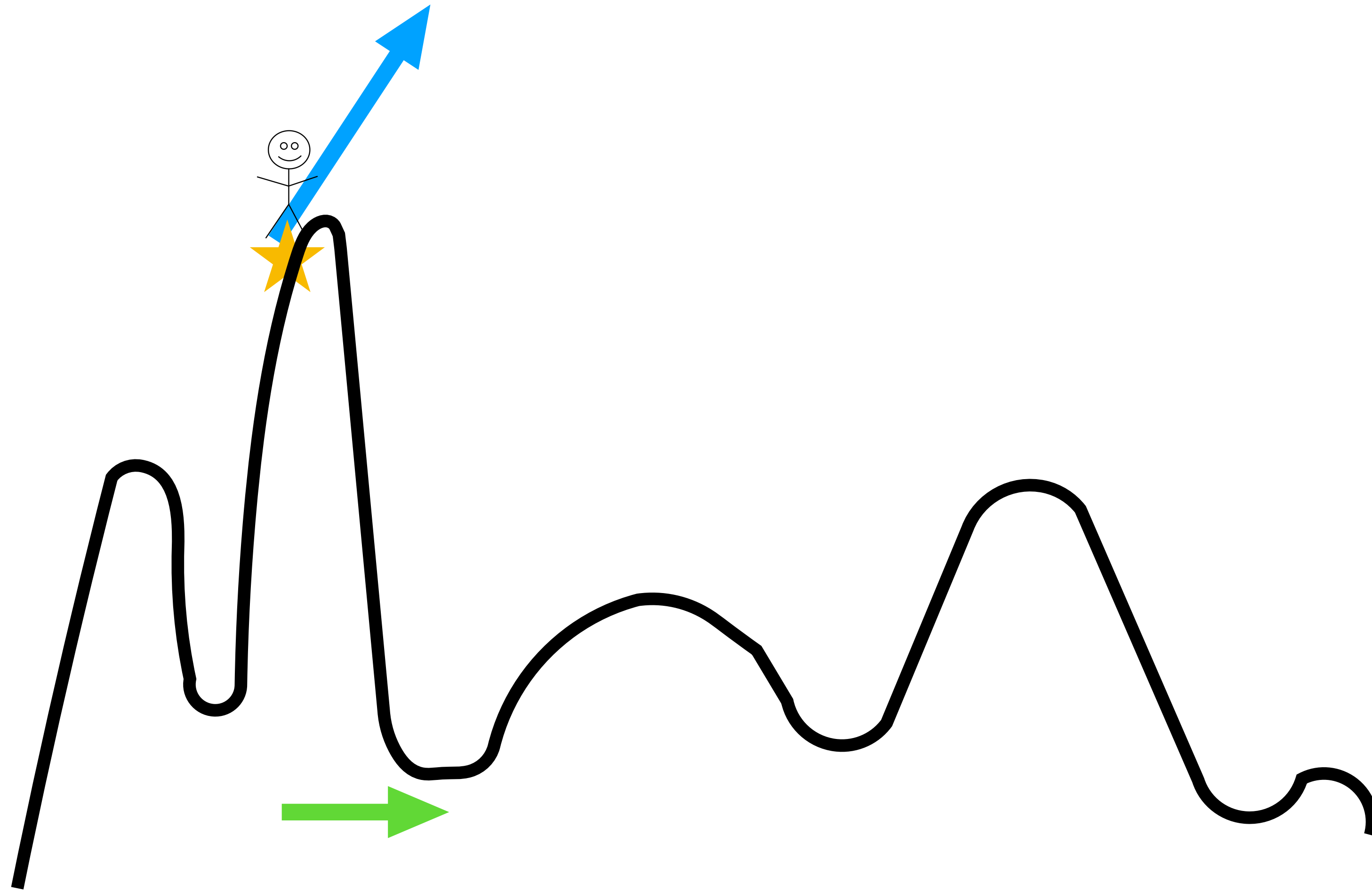


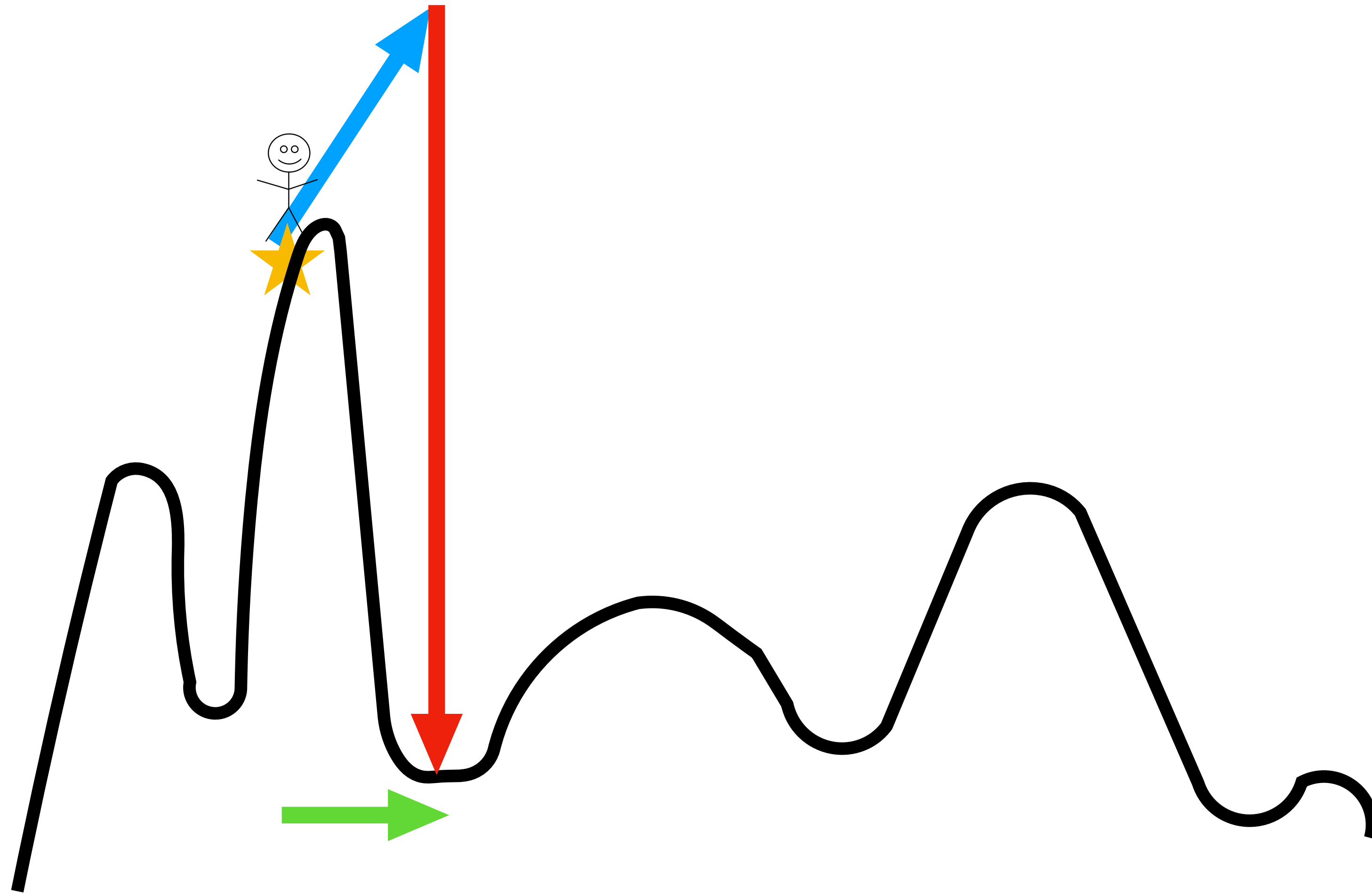


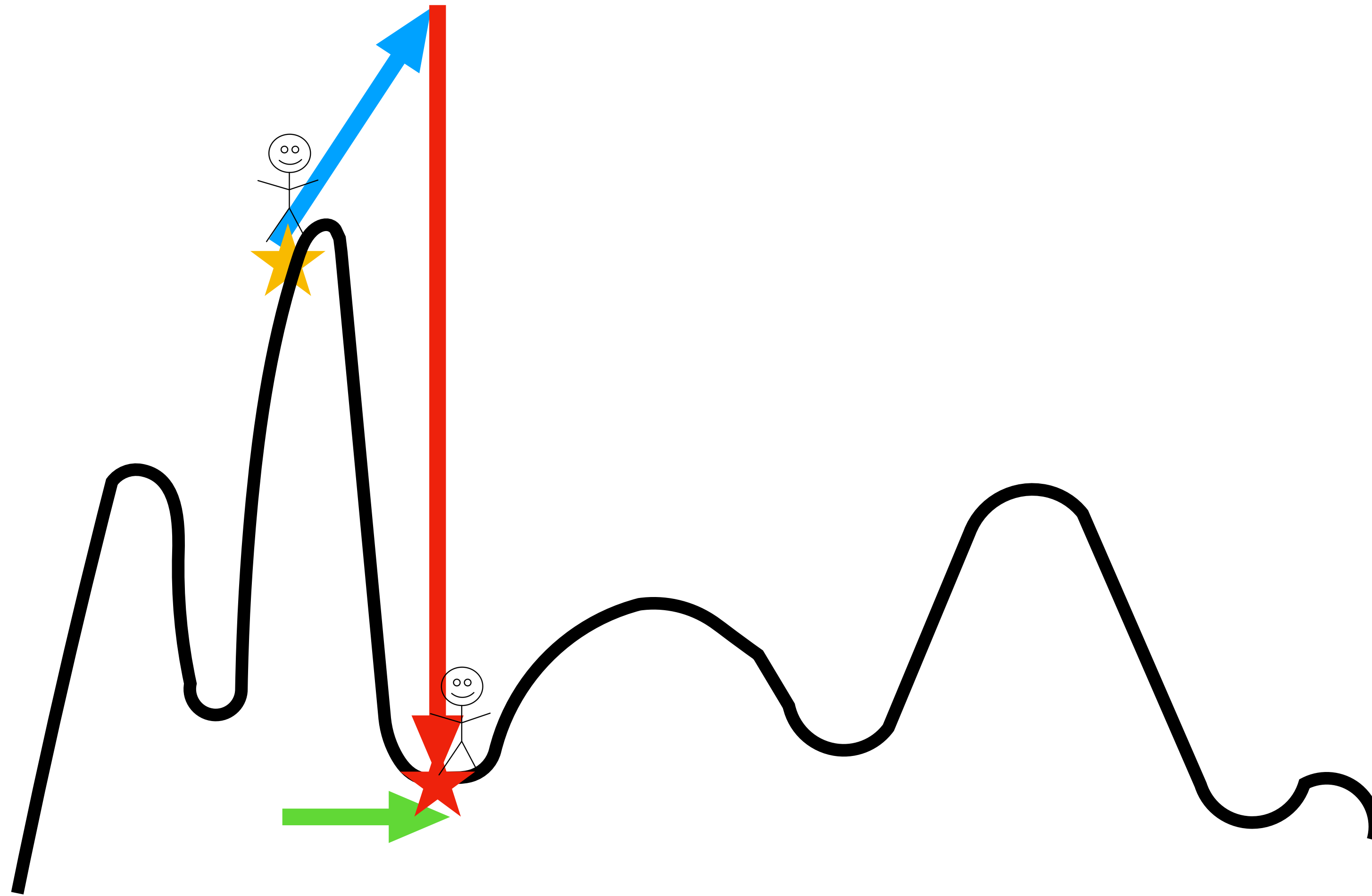


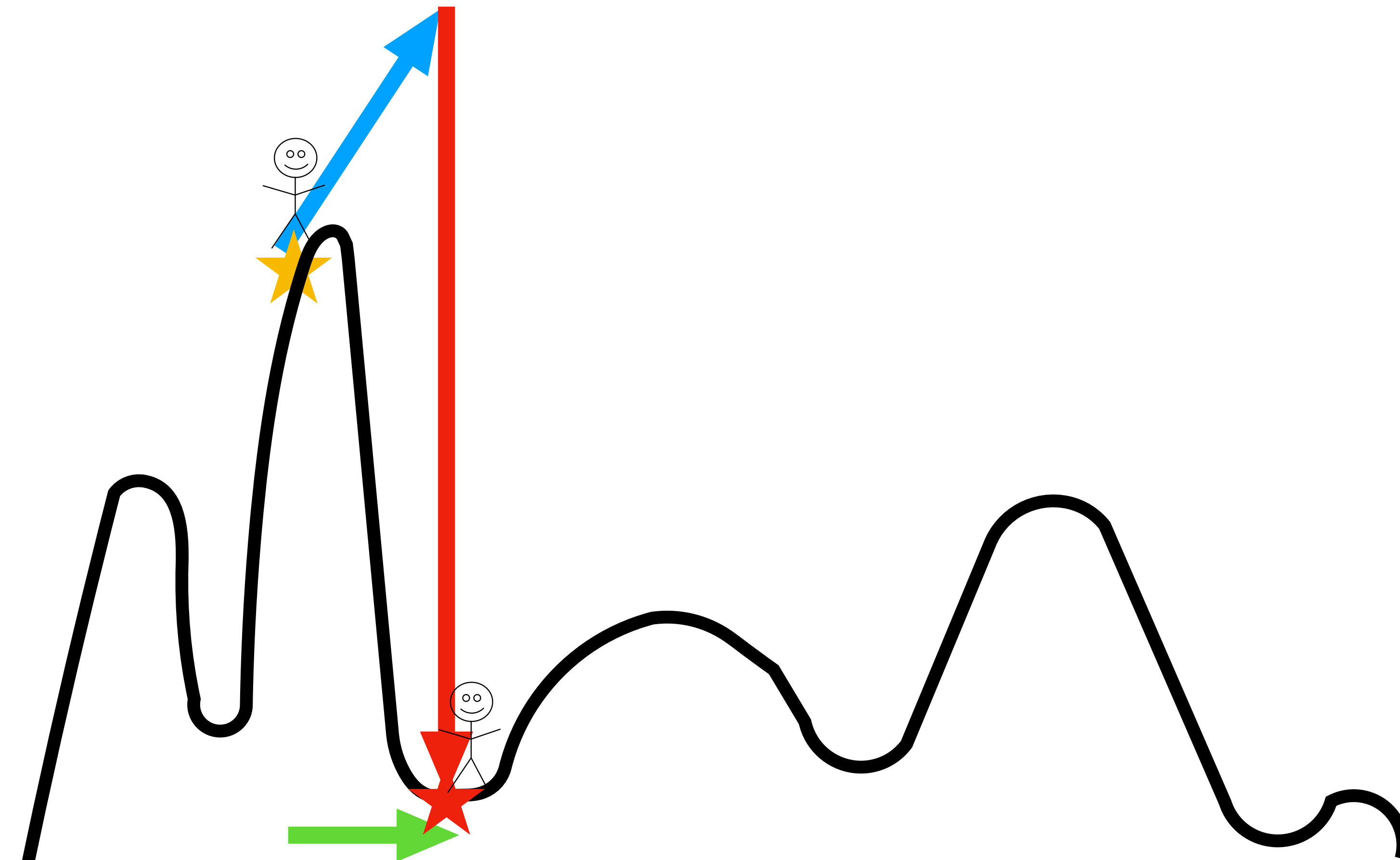












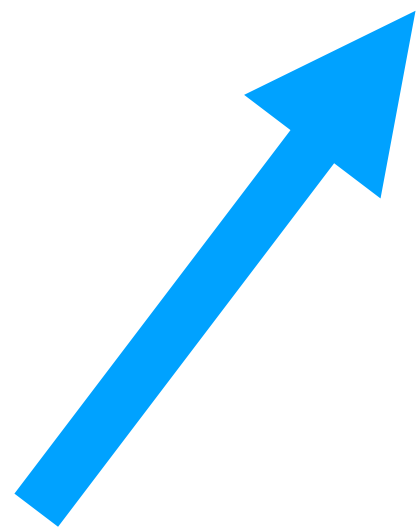
step size x gradient

Gradient descent

$$W = W - \alpha \nabla_w f$$

Gradient descent

$$W = W - \alpha \nabla_w f$$



Parameters

Gradient descent

$$W = W - \alpha \nabla_w f$$

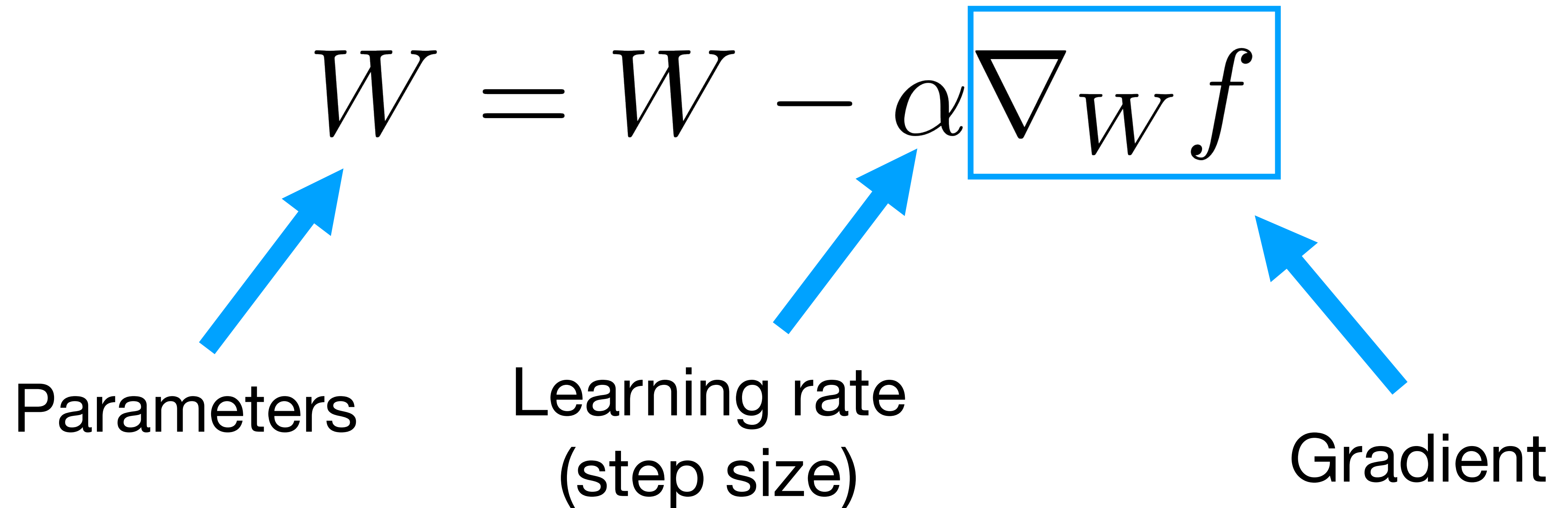
Parameters



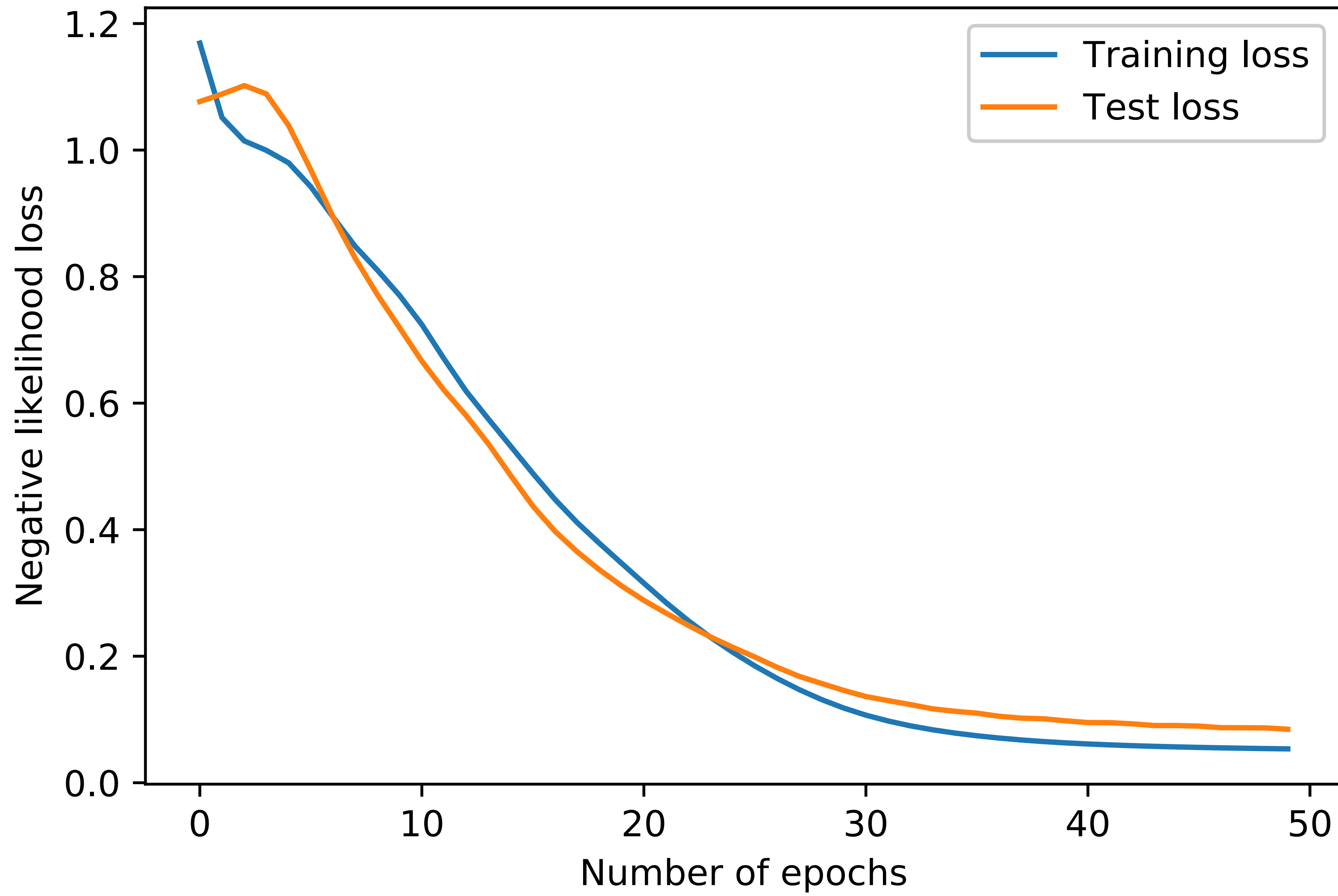
Gradient

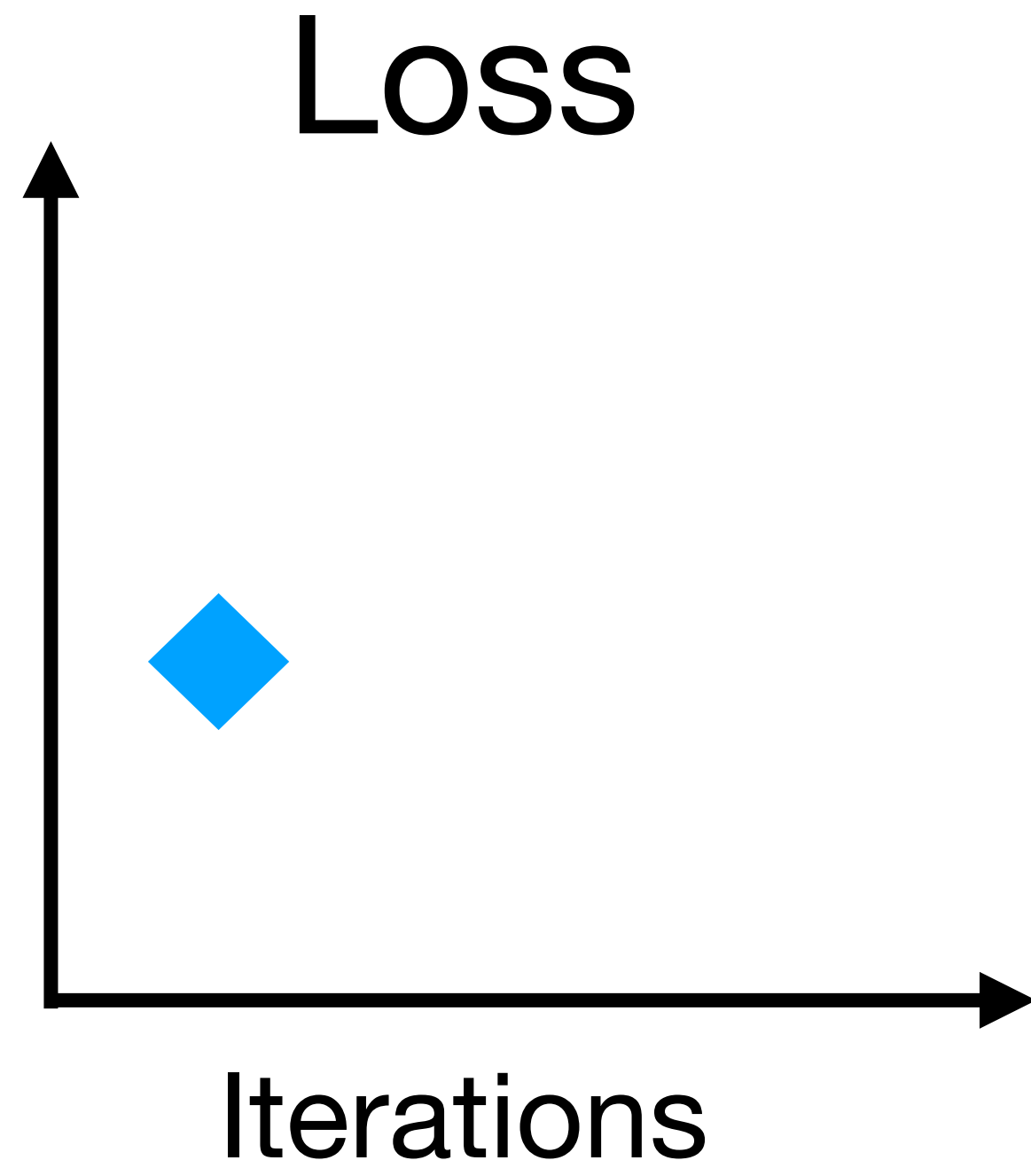
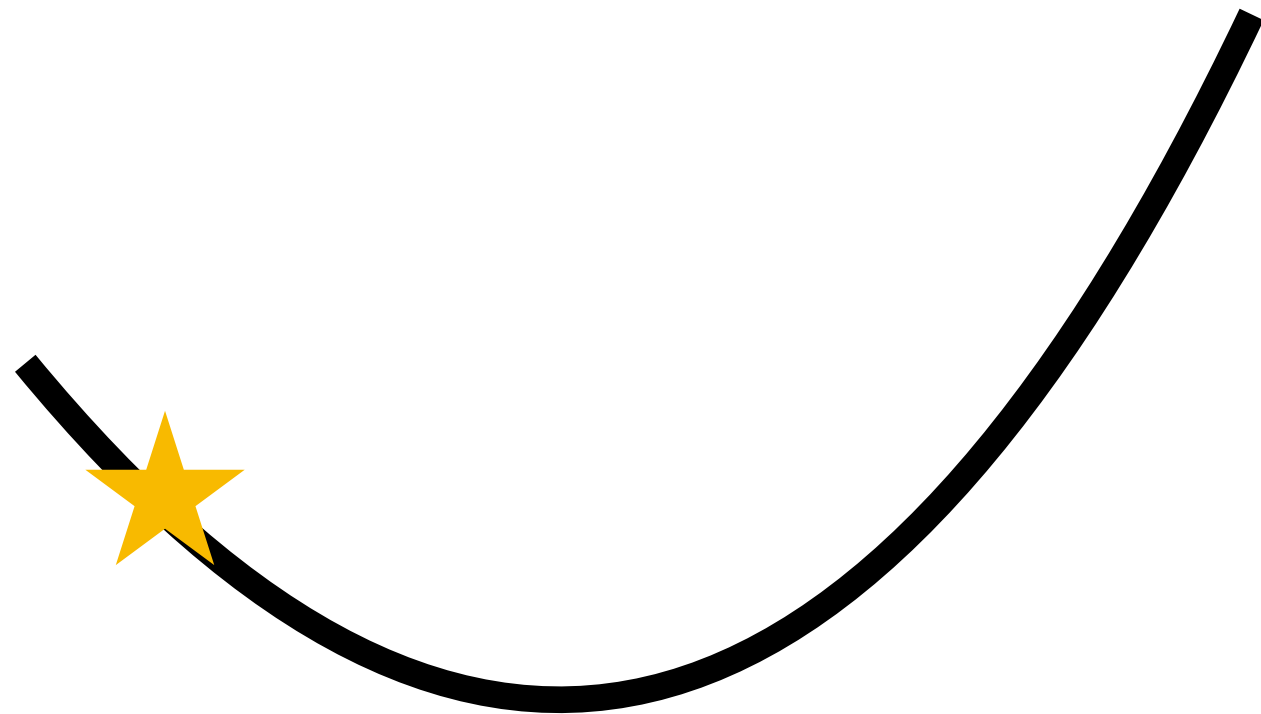


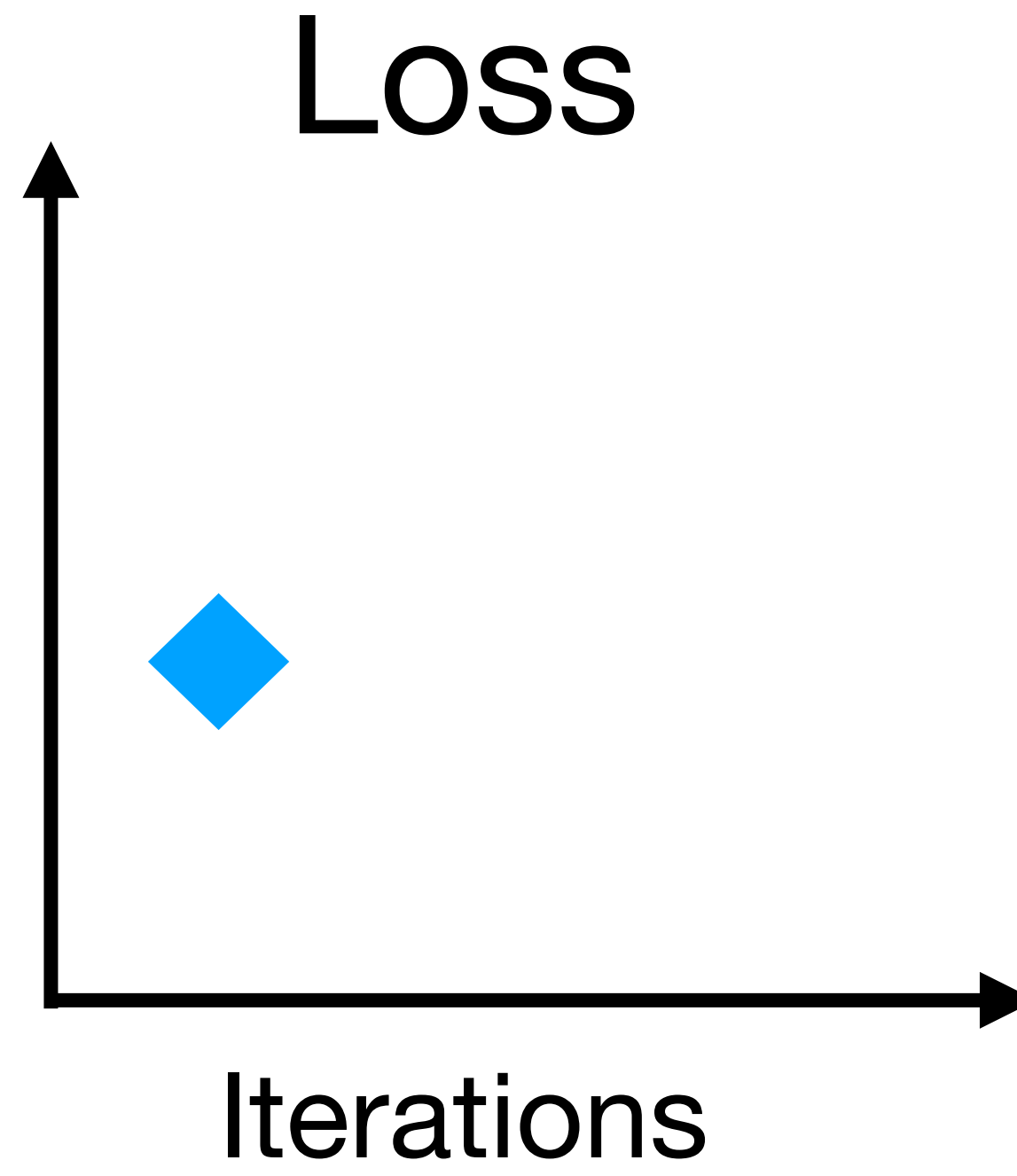
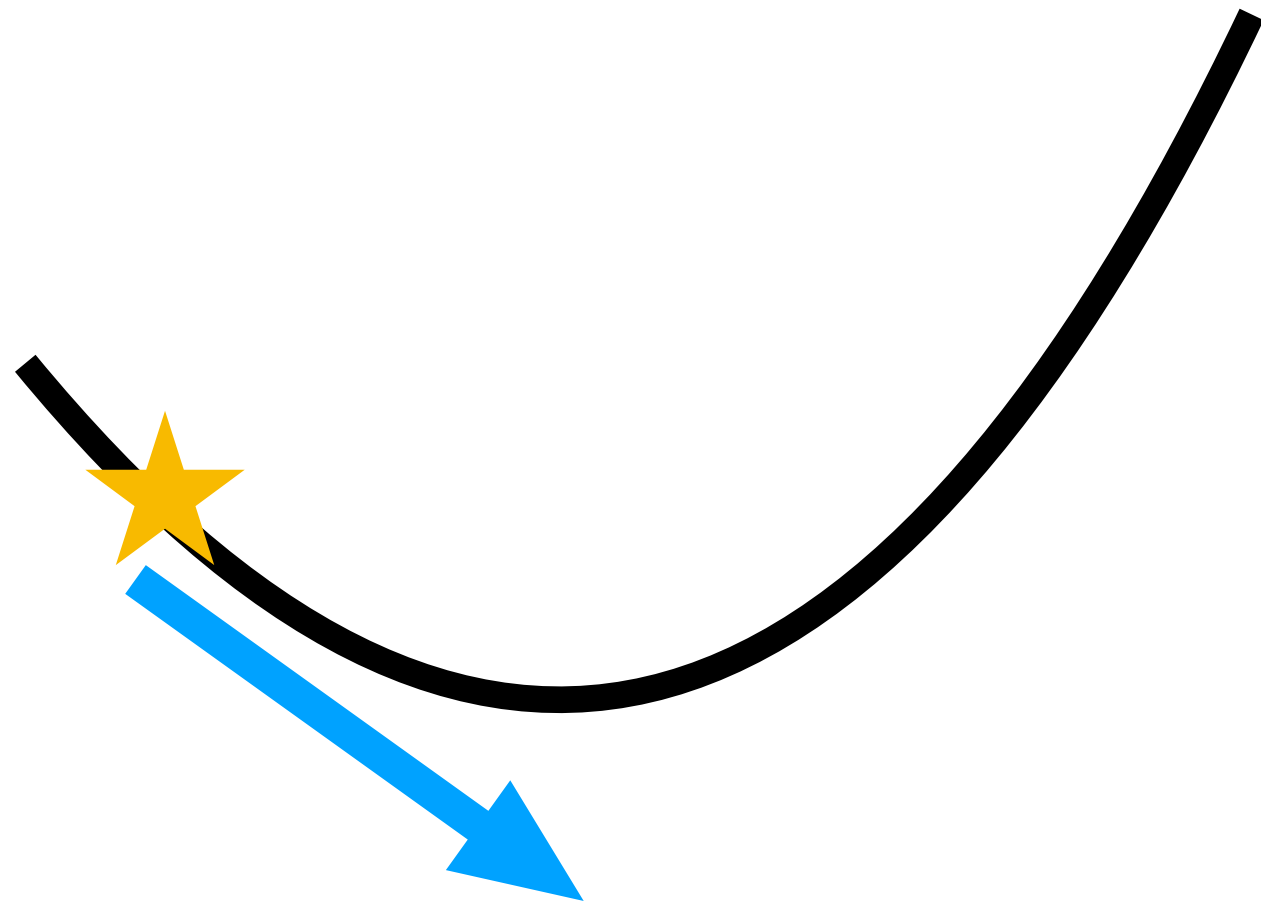
Gradient descent

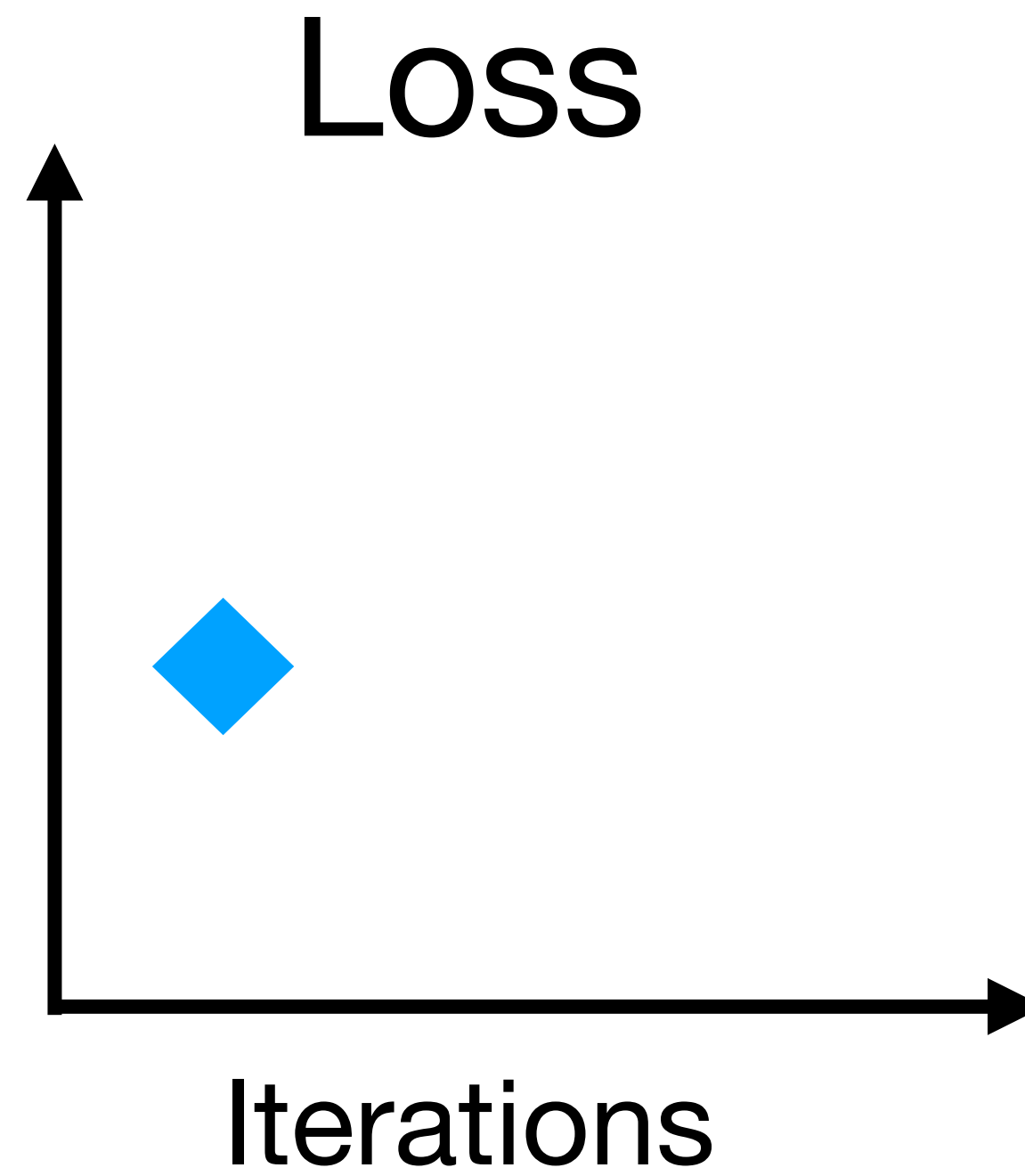
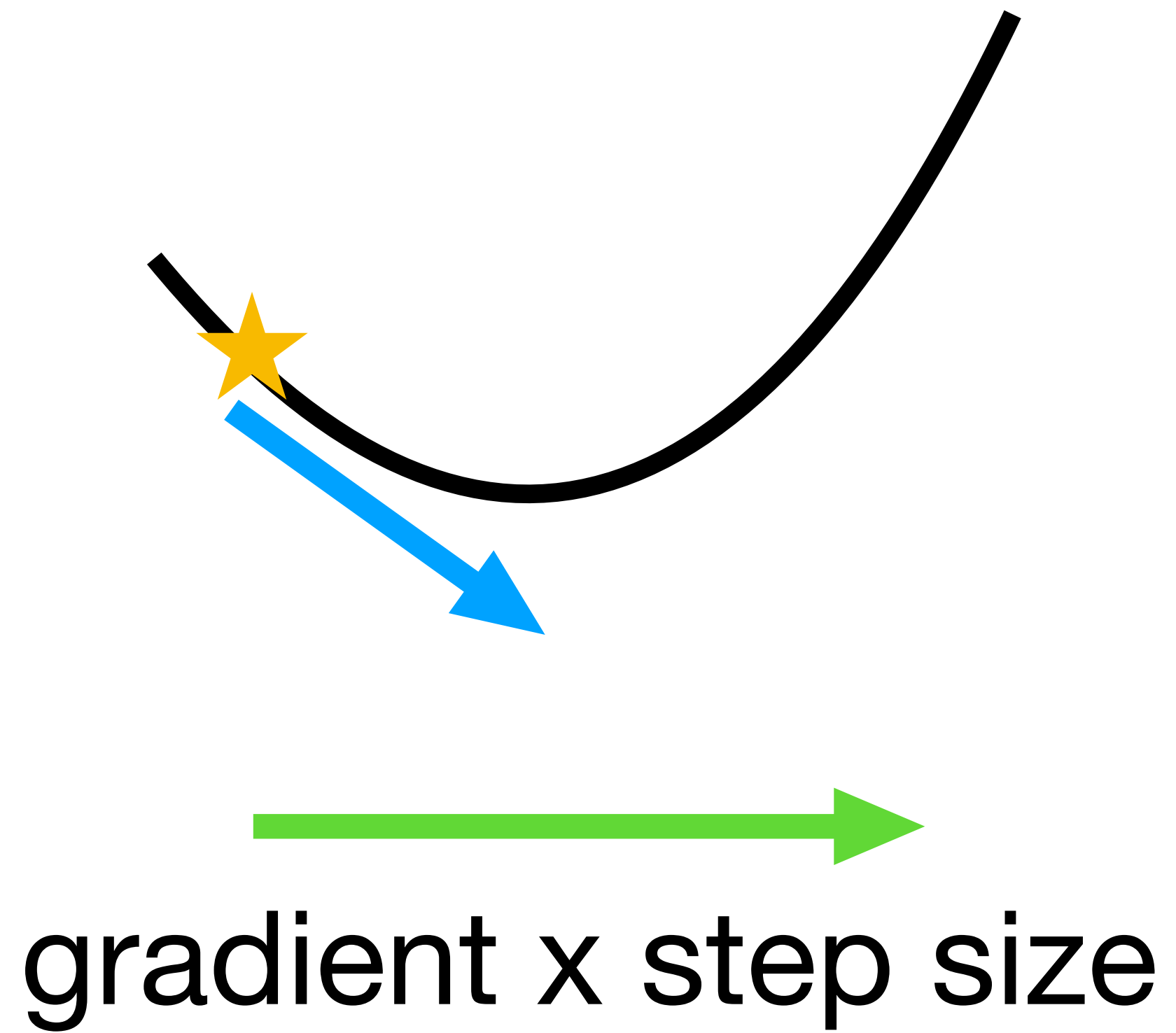


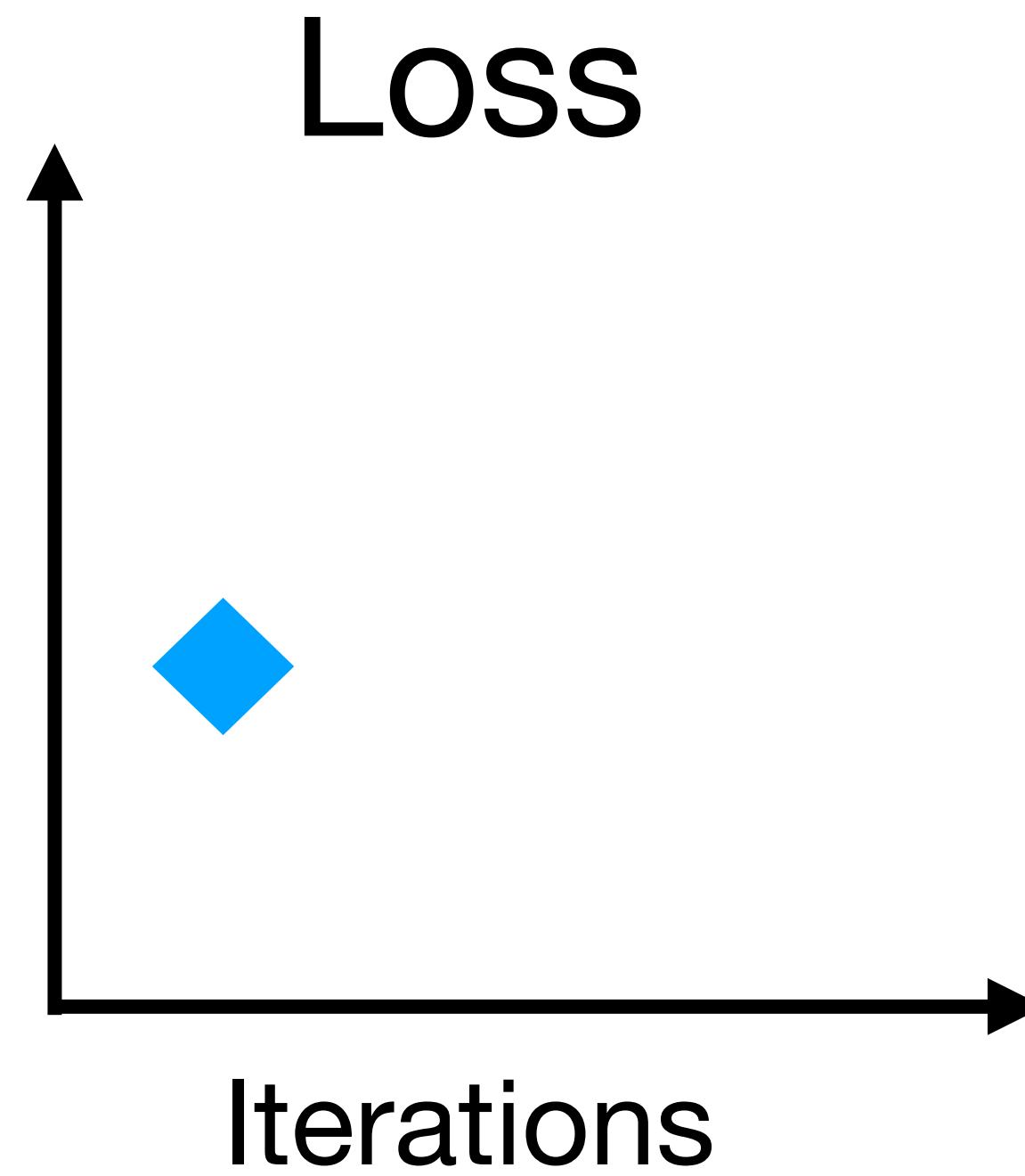
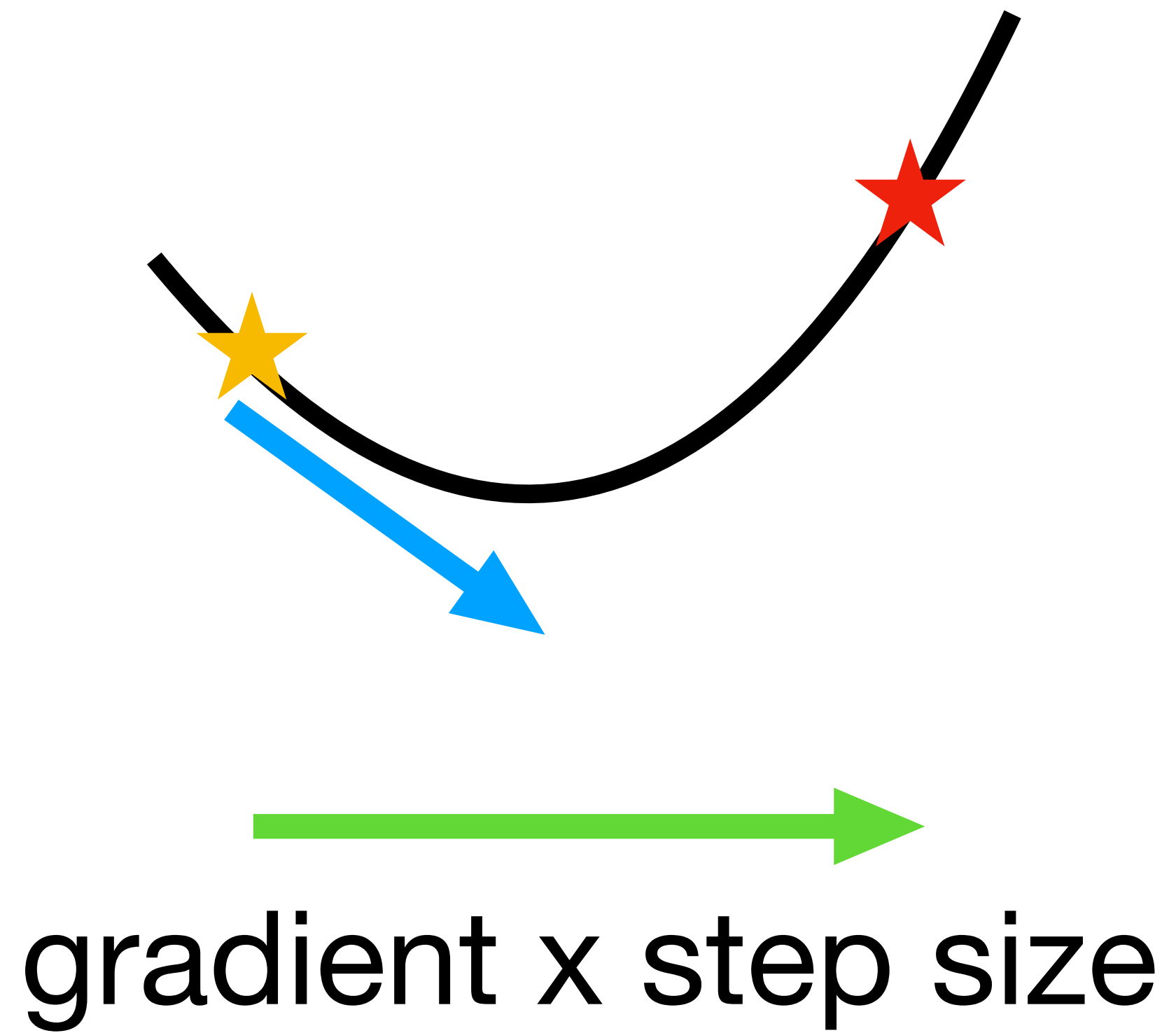
Good loss curves

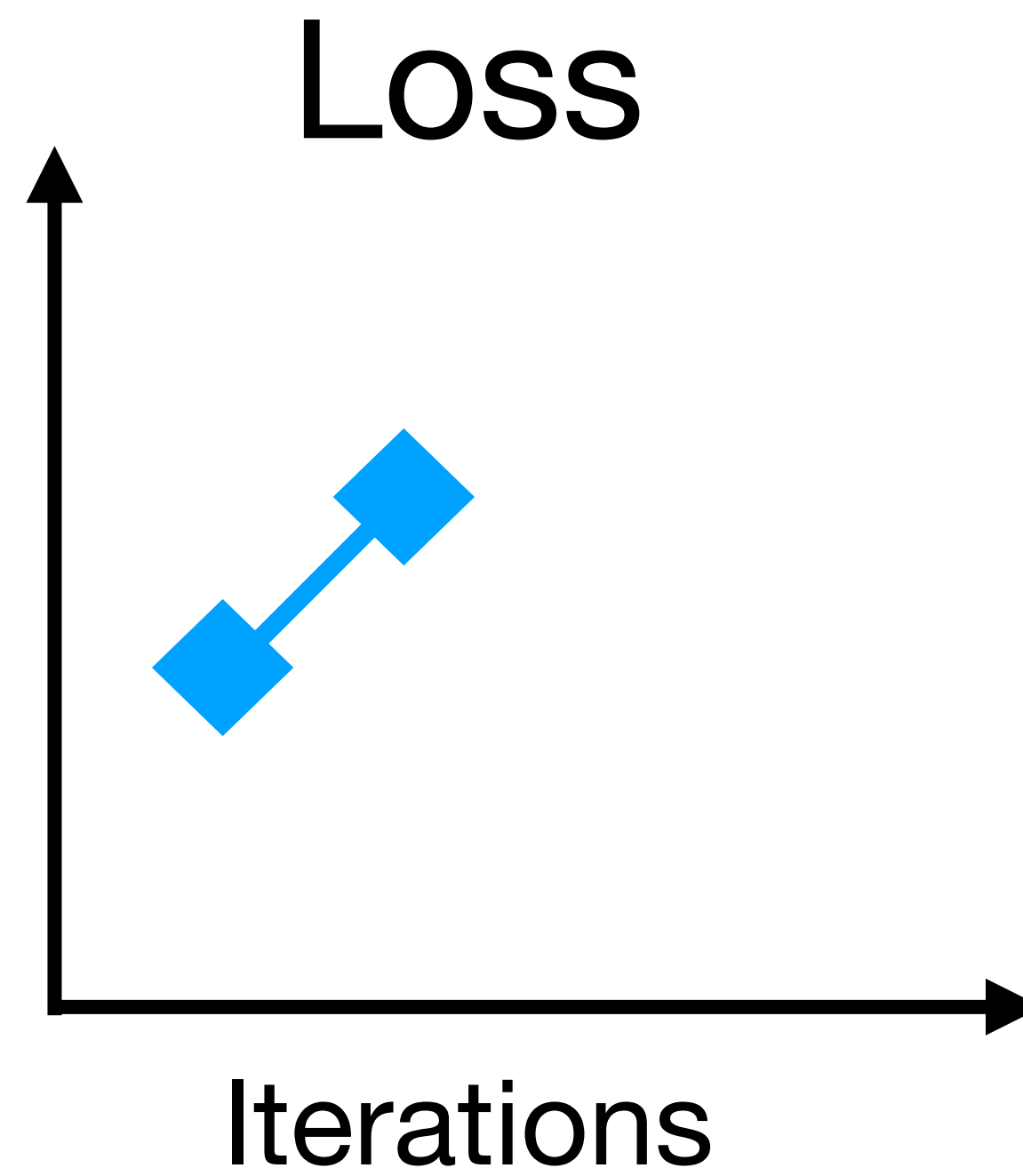
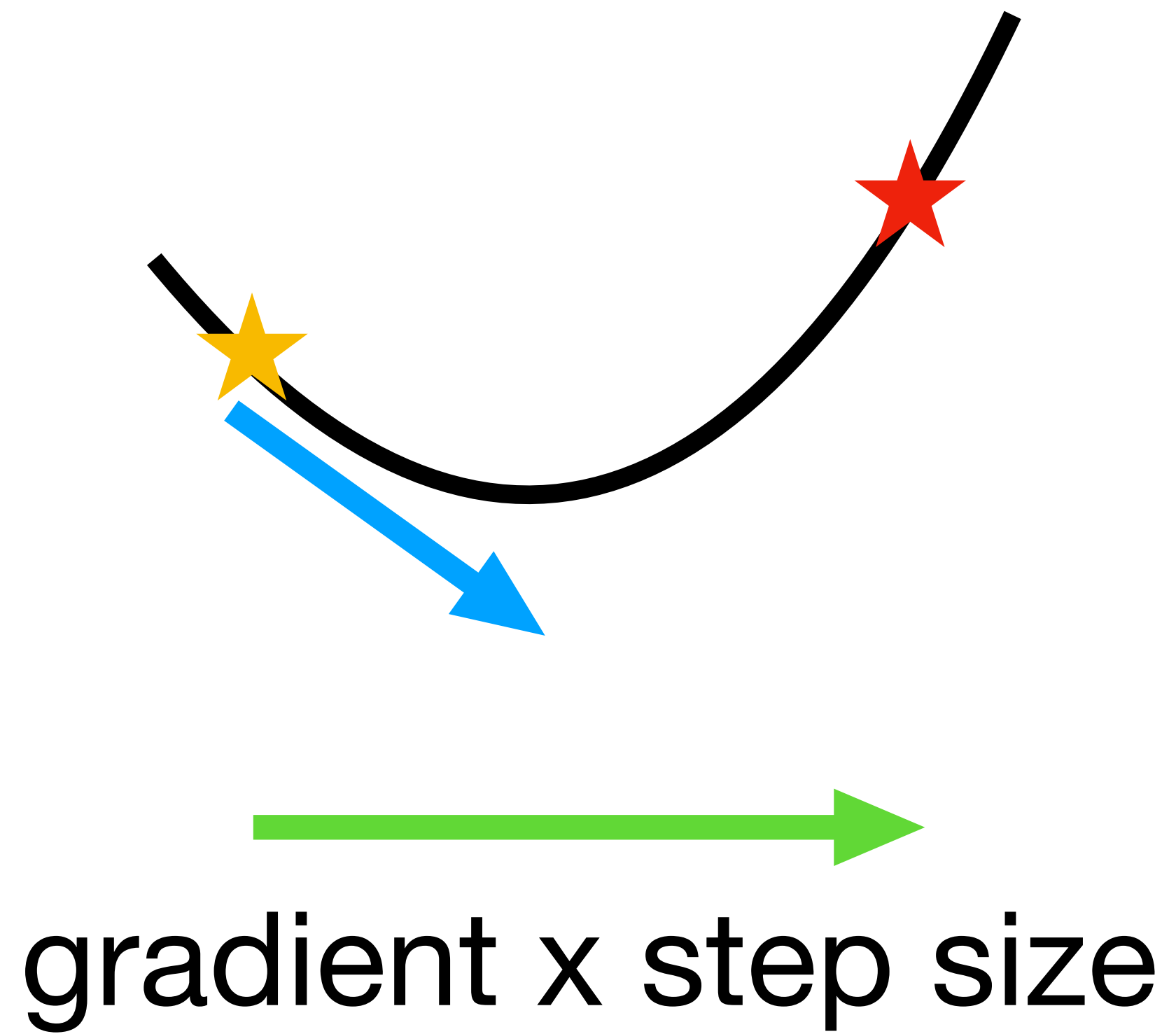


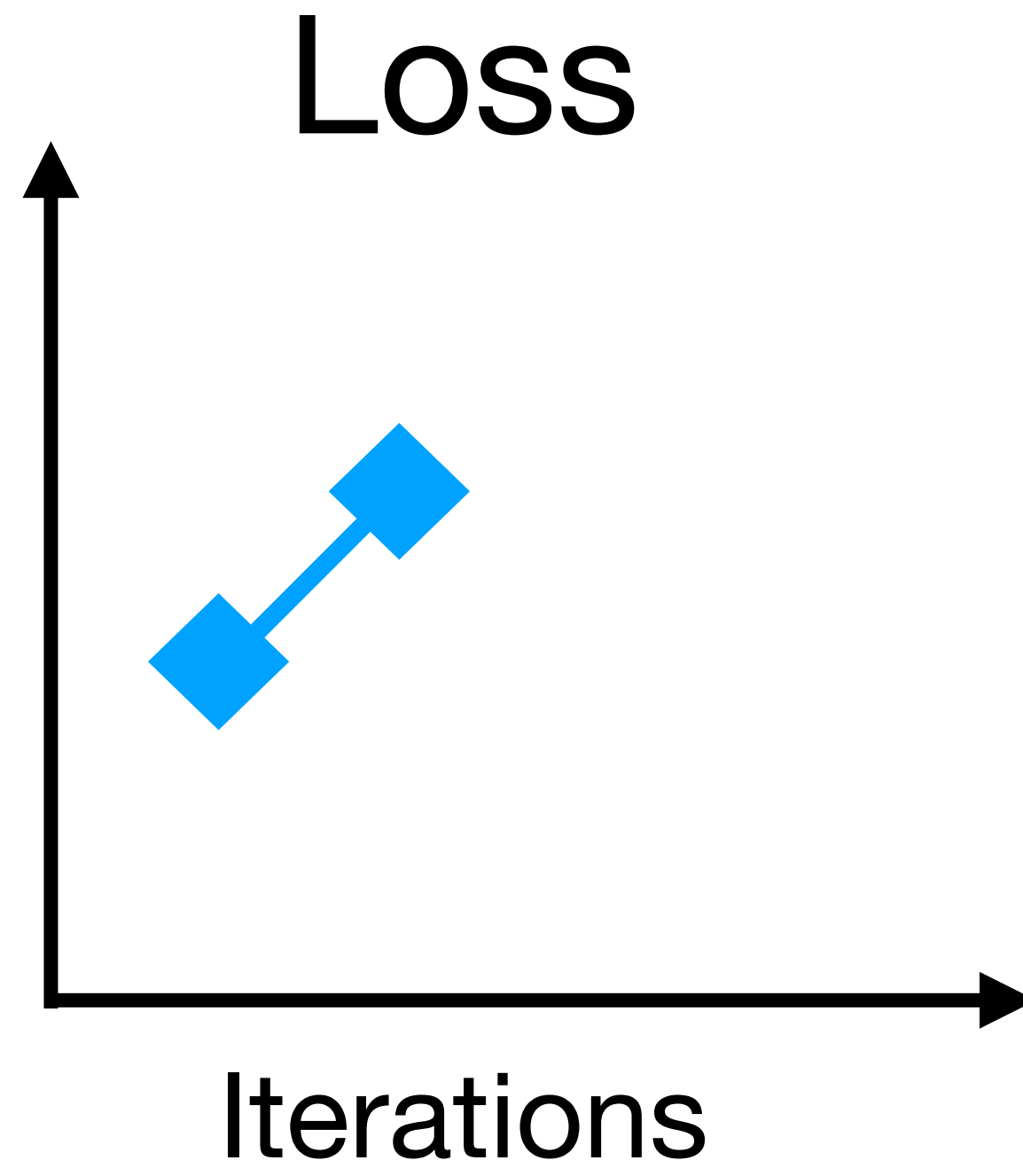
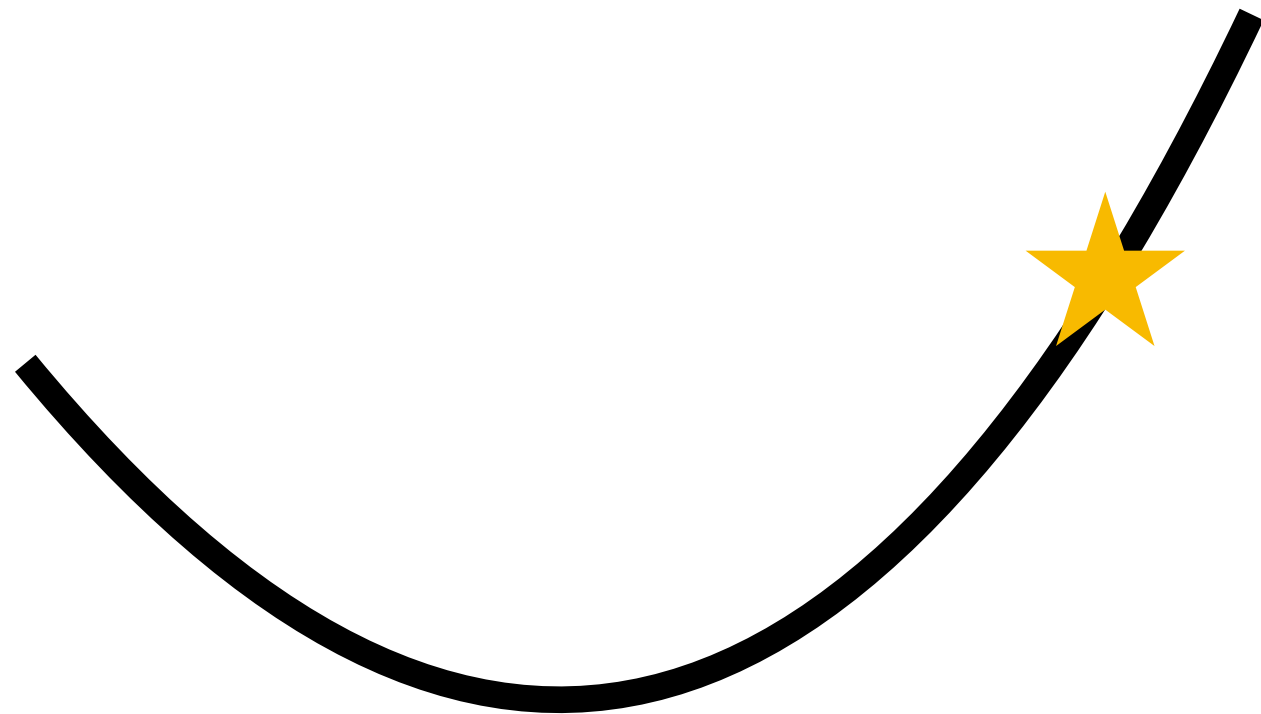


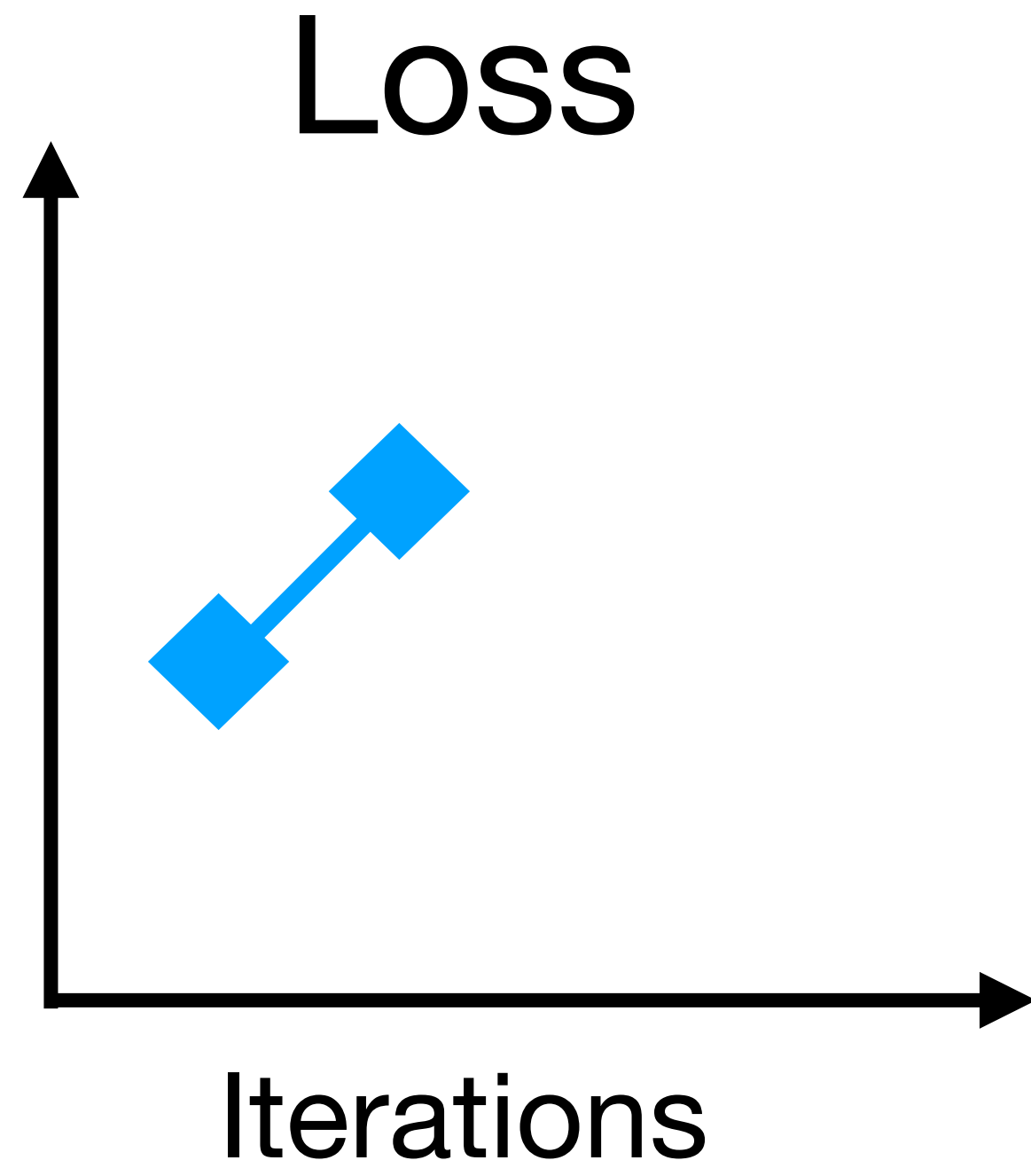
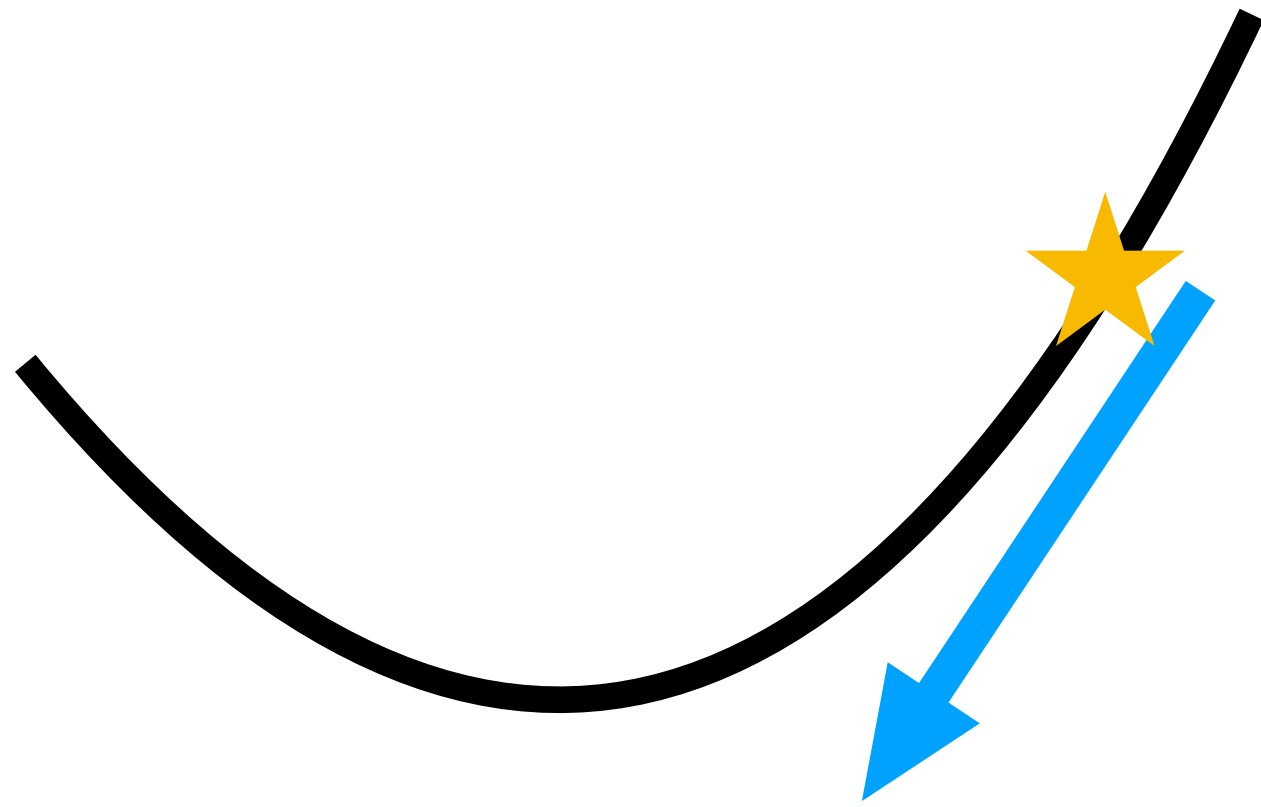


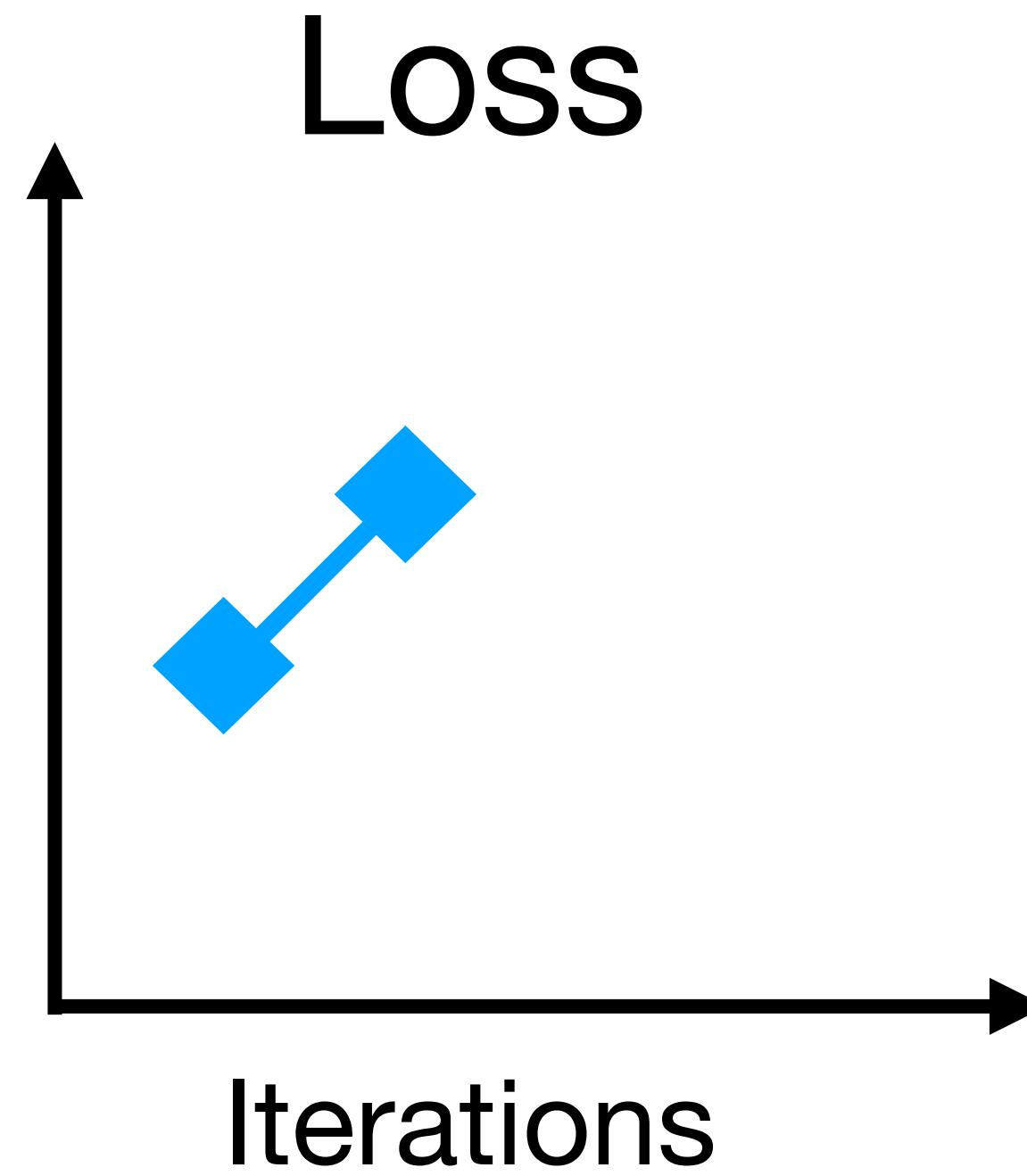
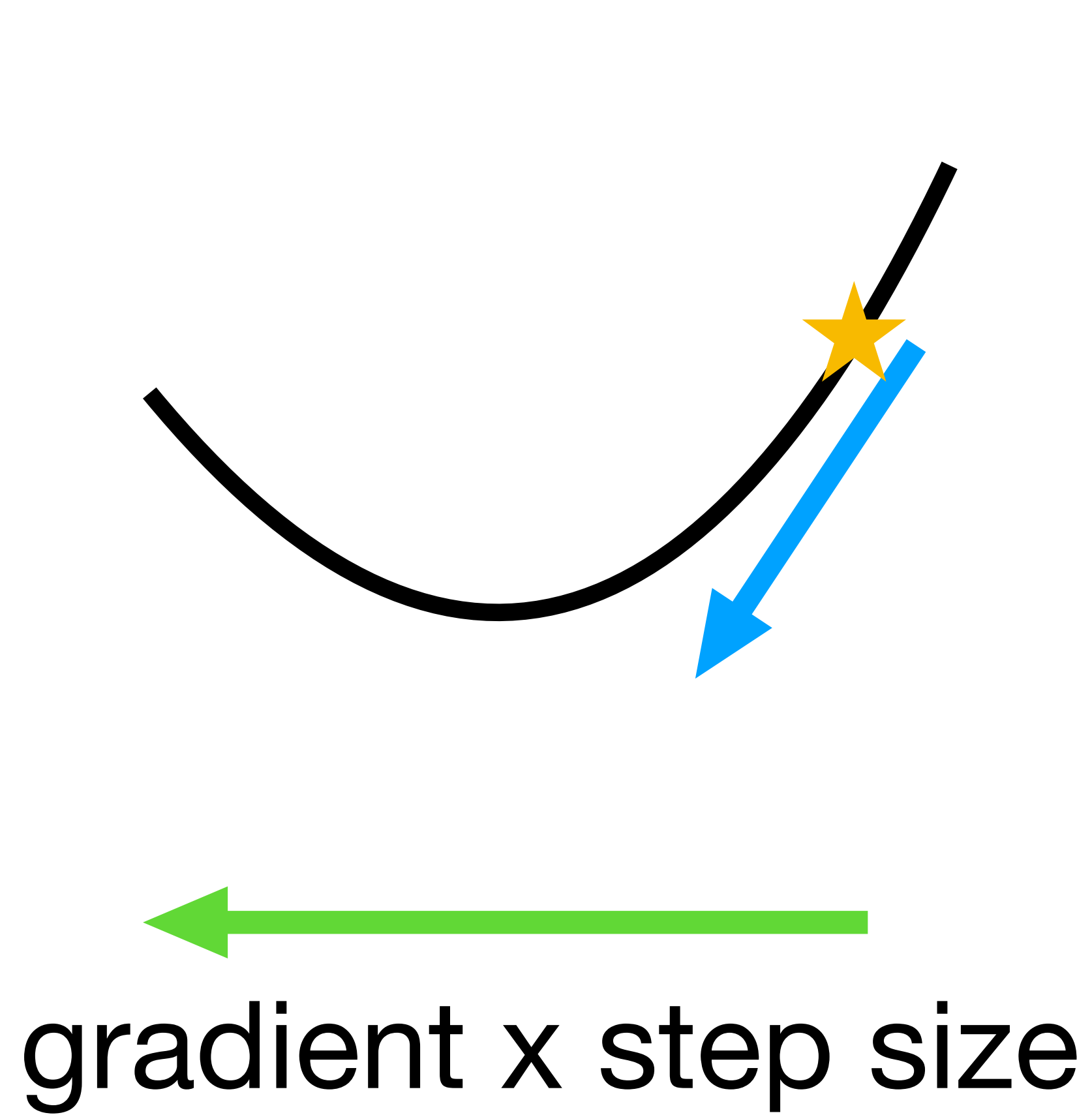


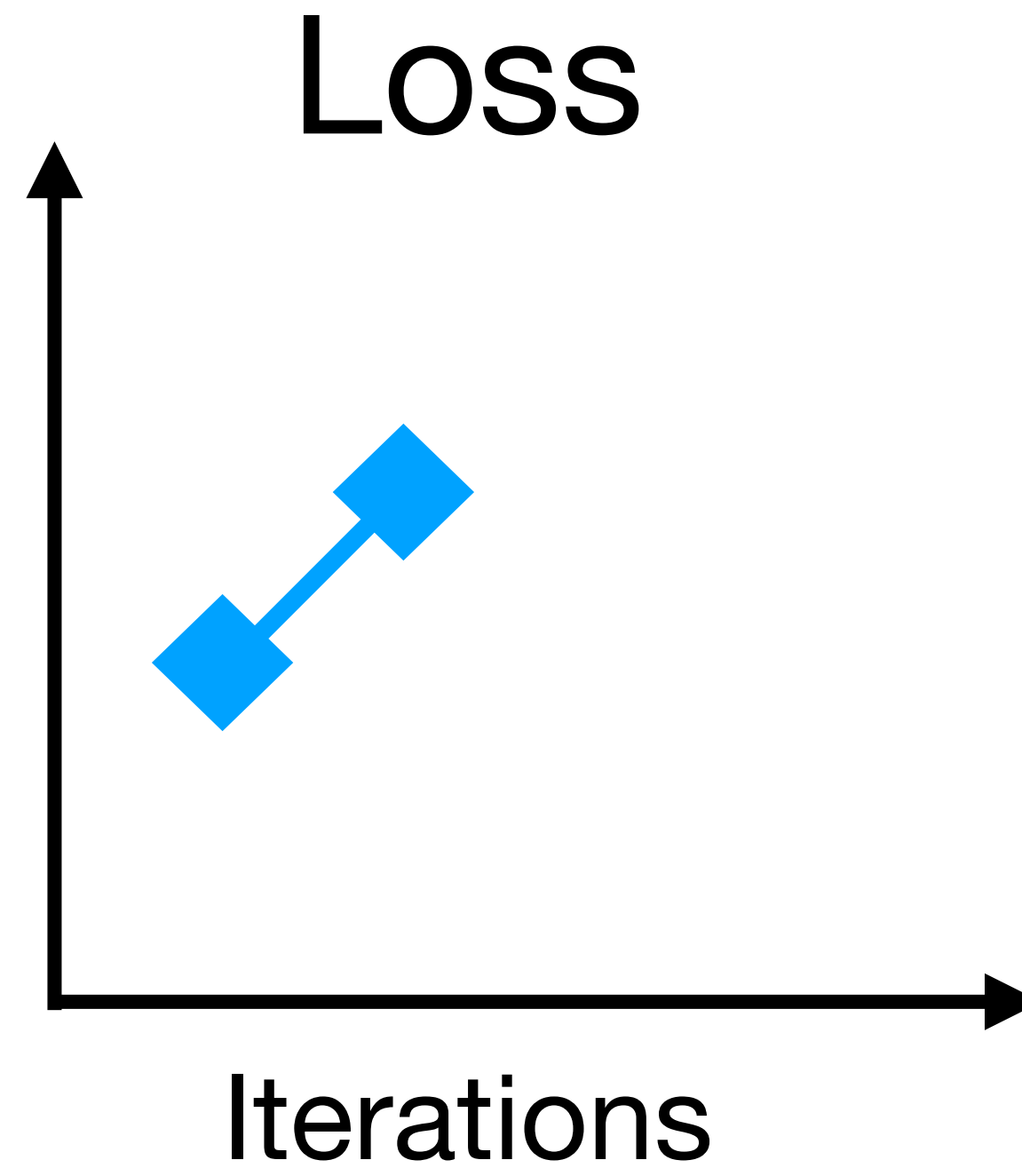
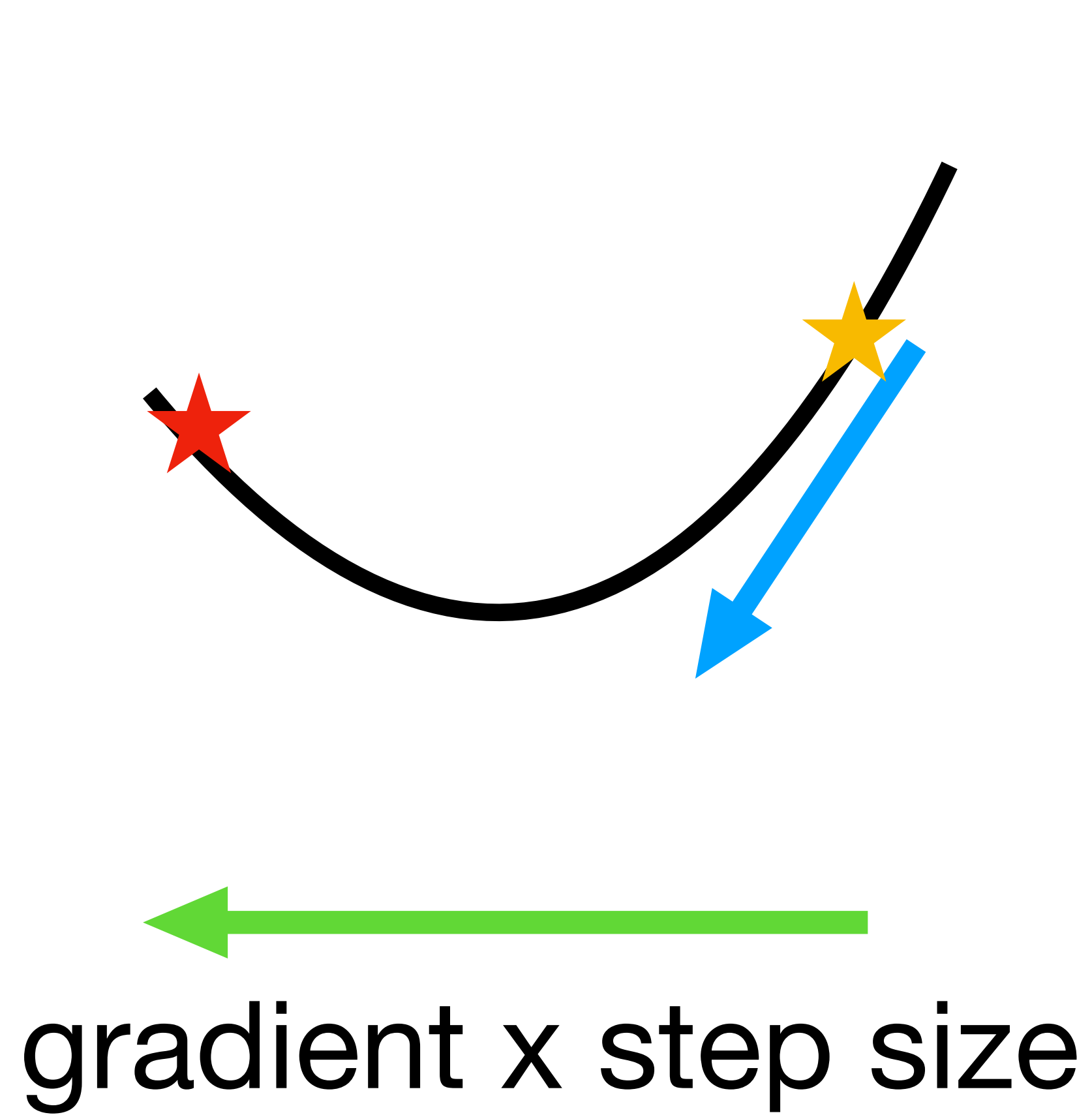


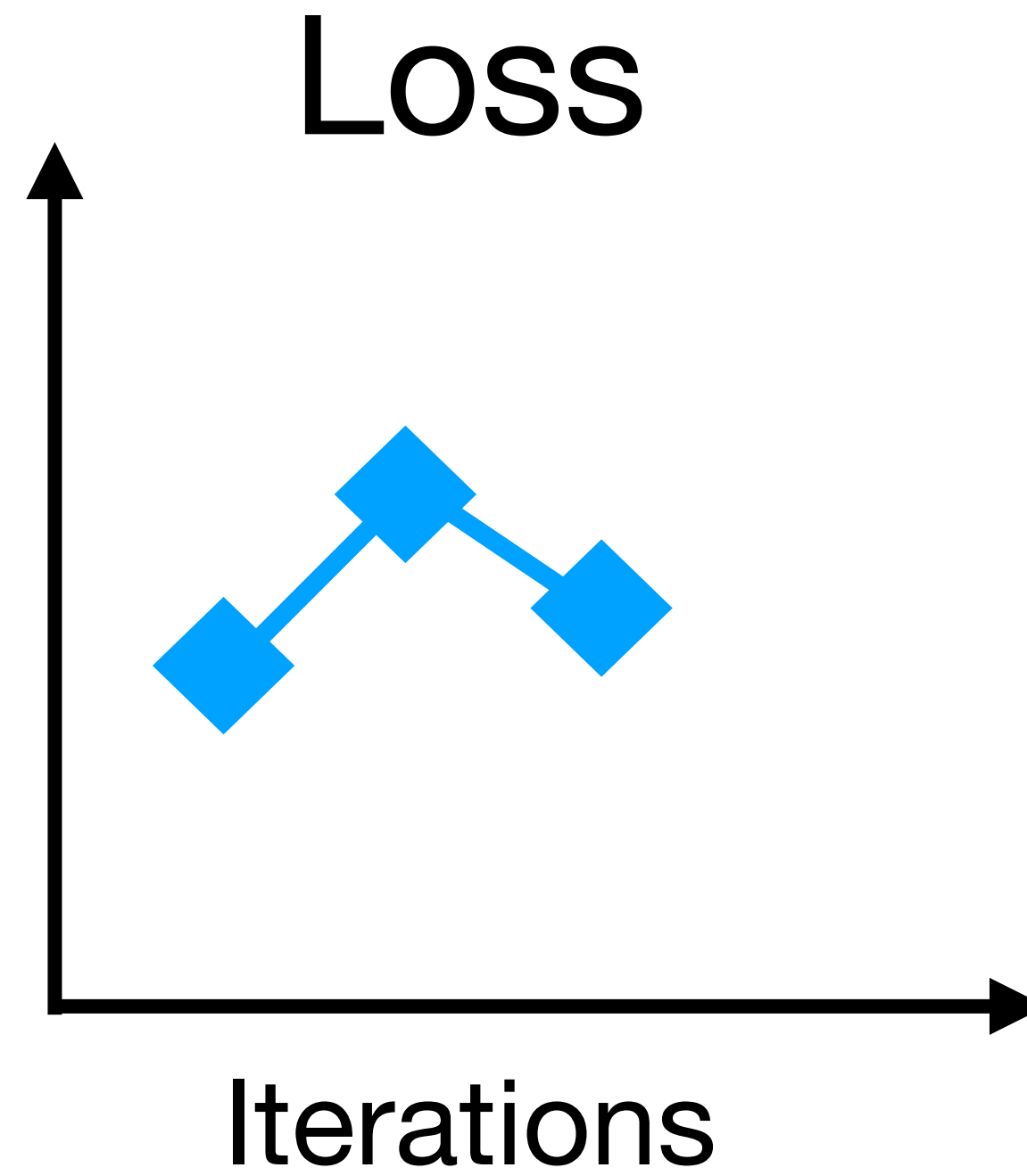
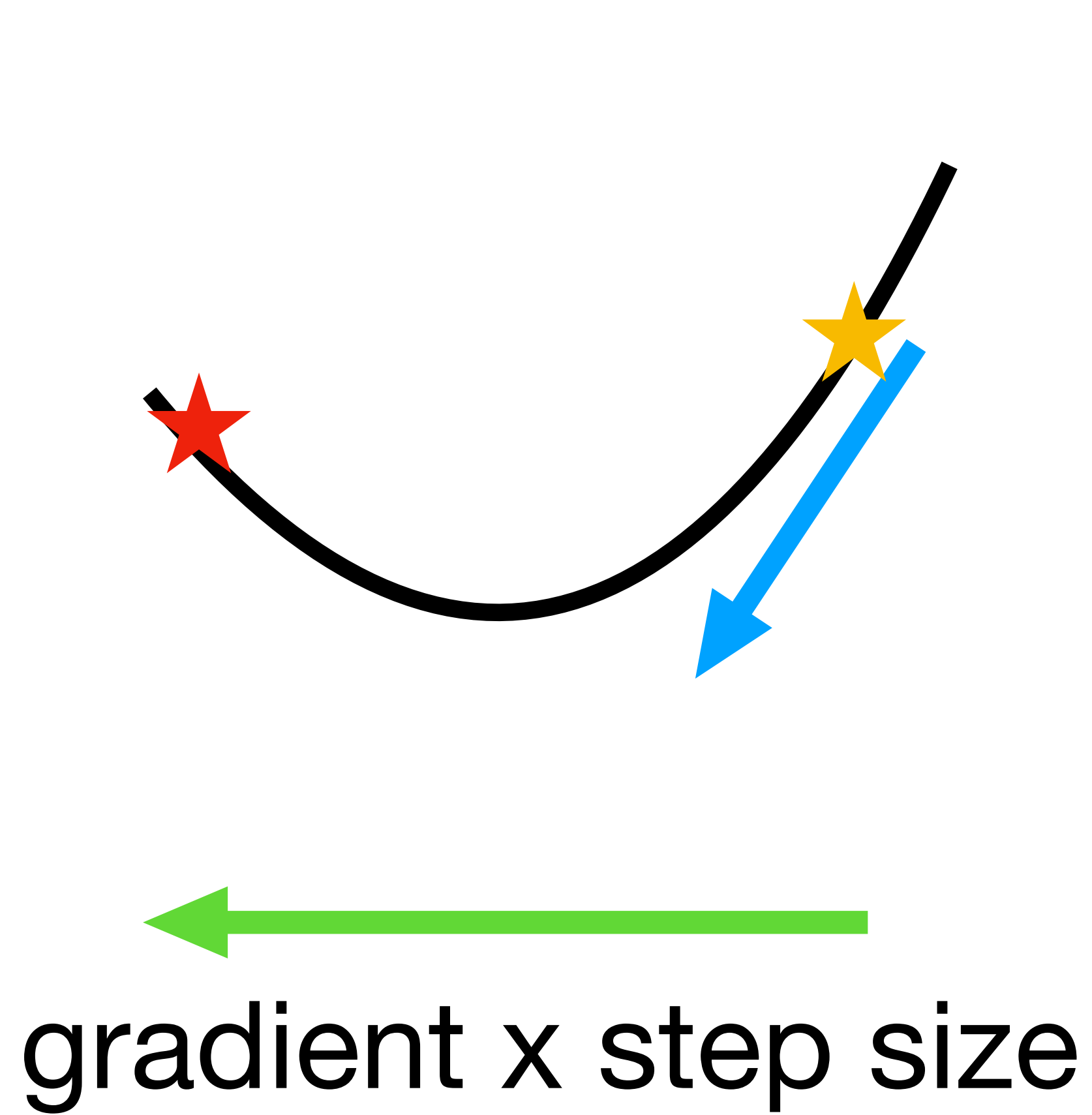


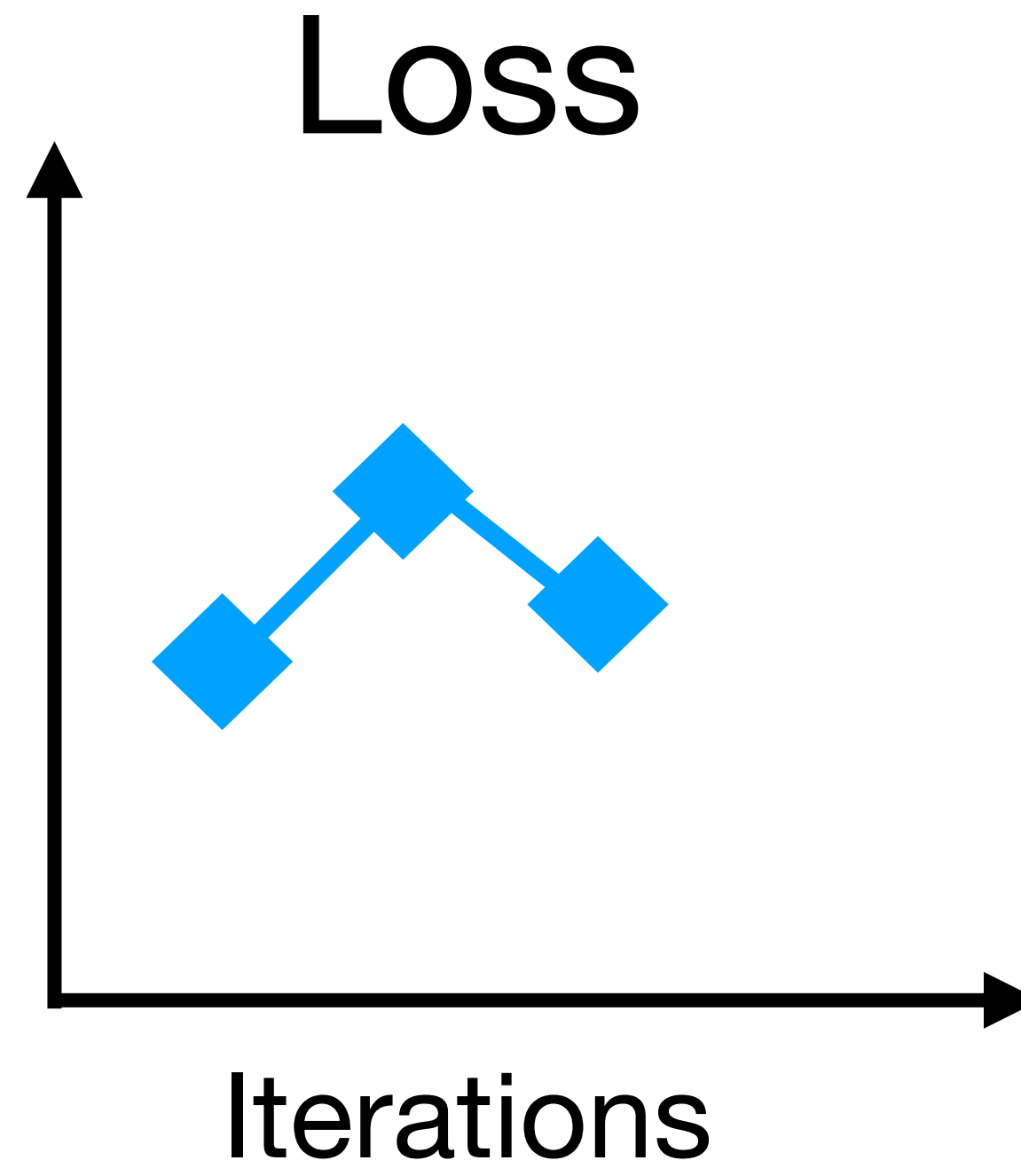
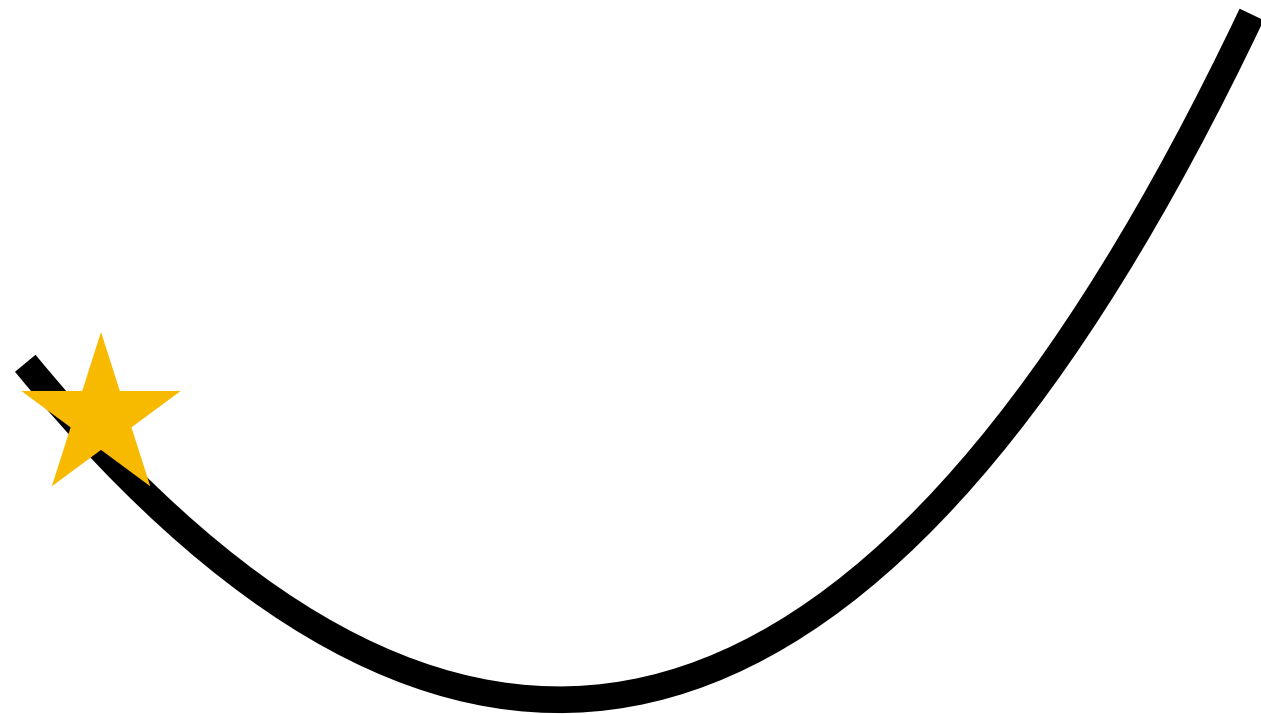


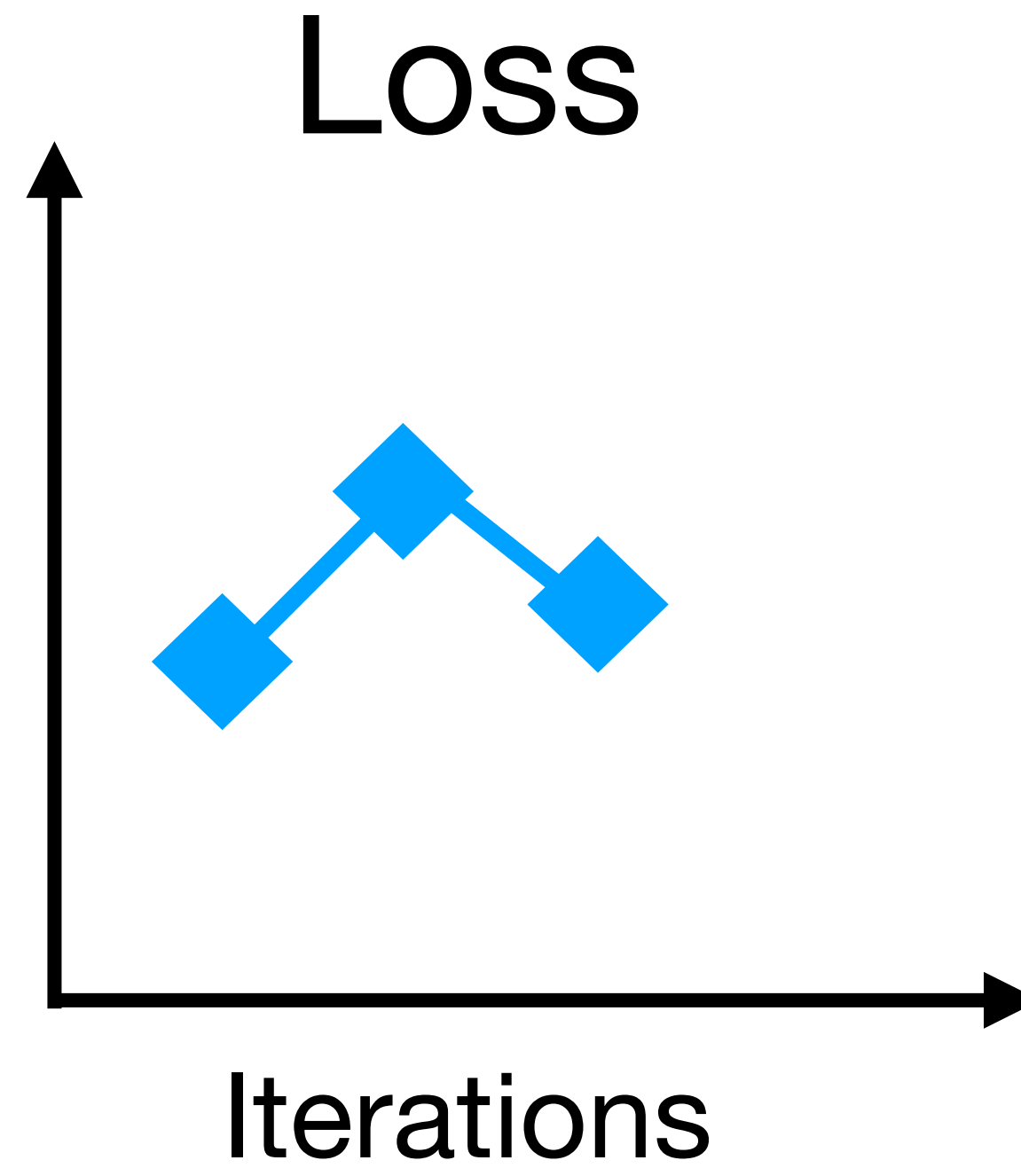
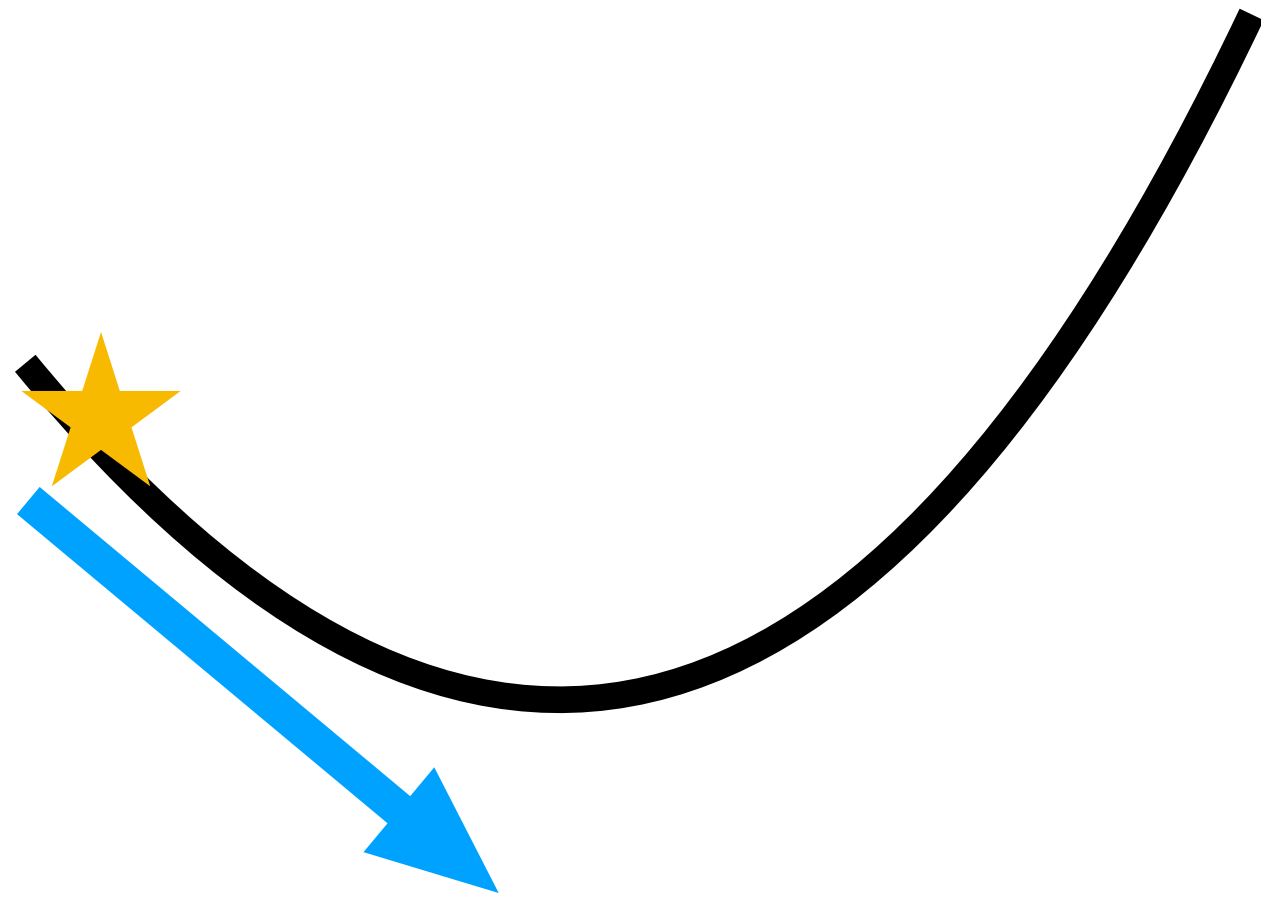


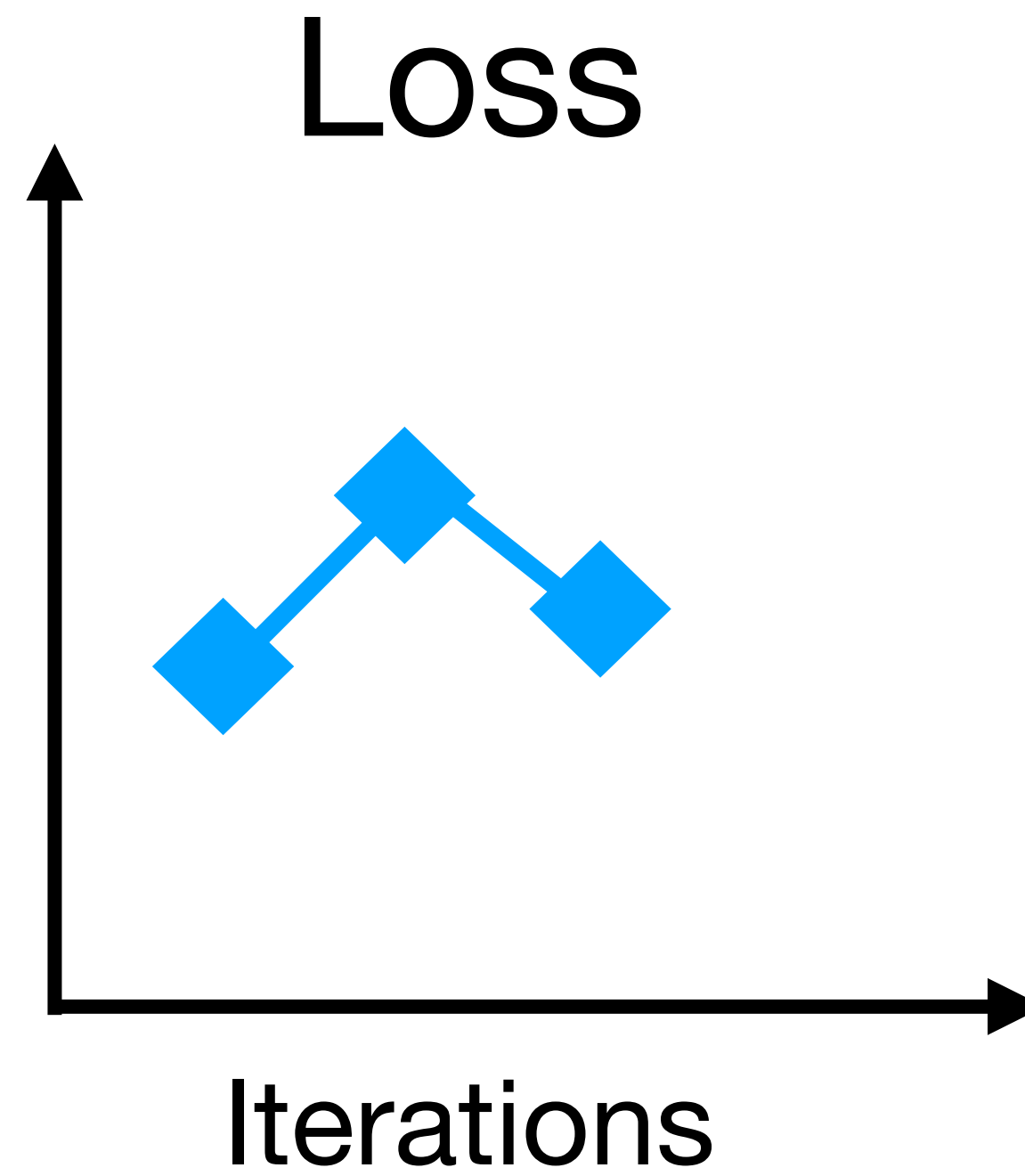
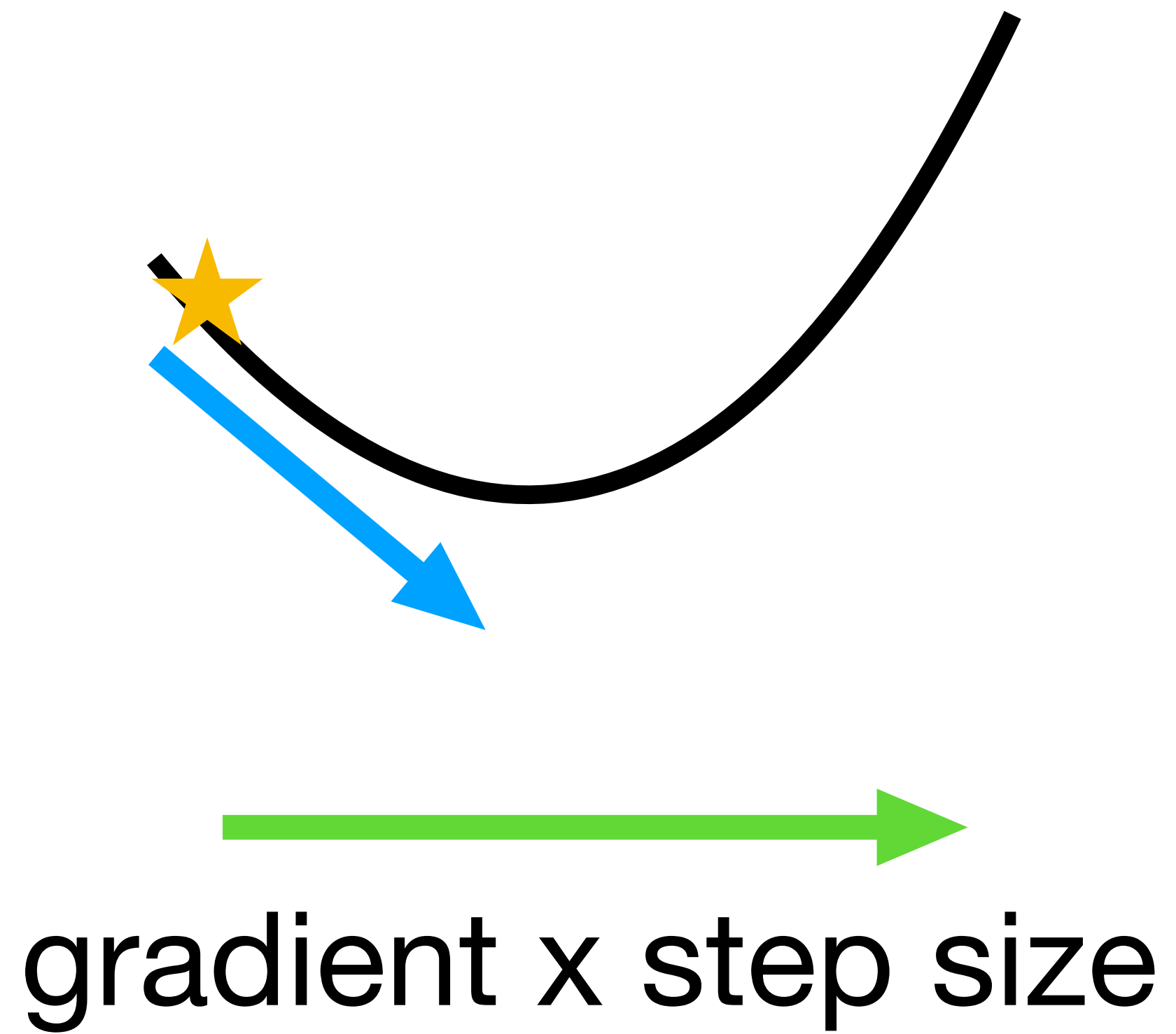


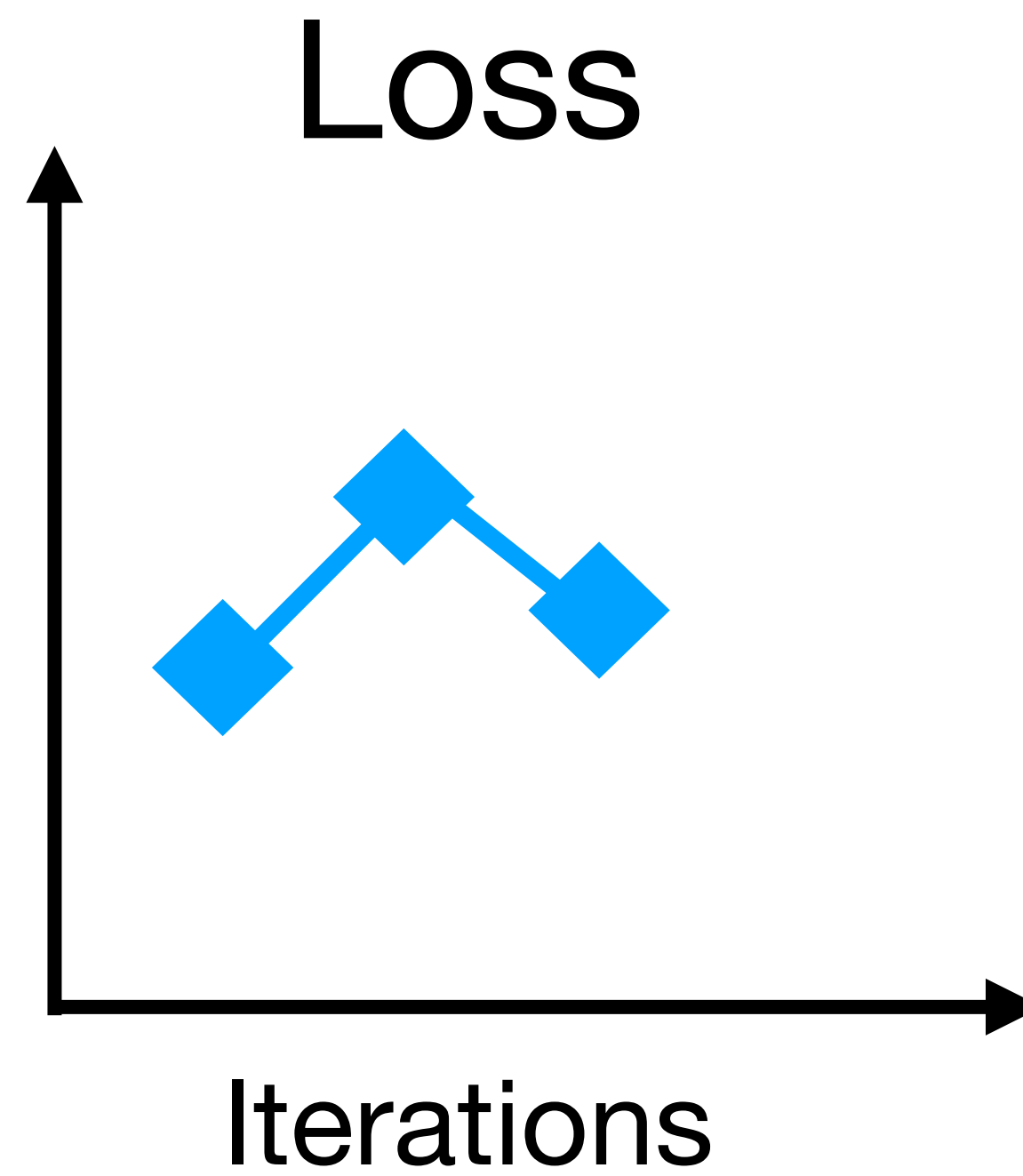
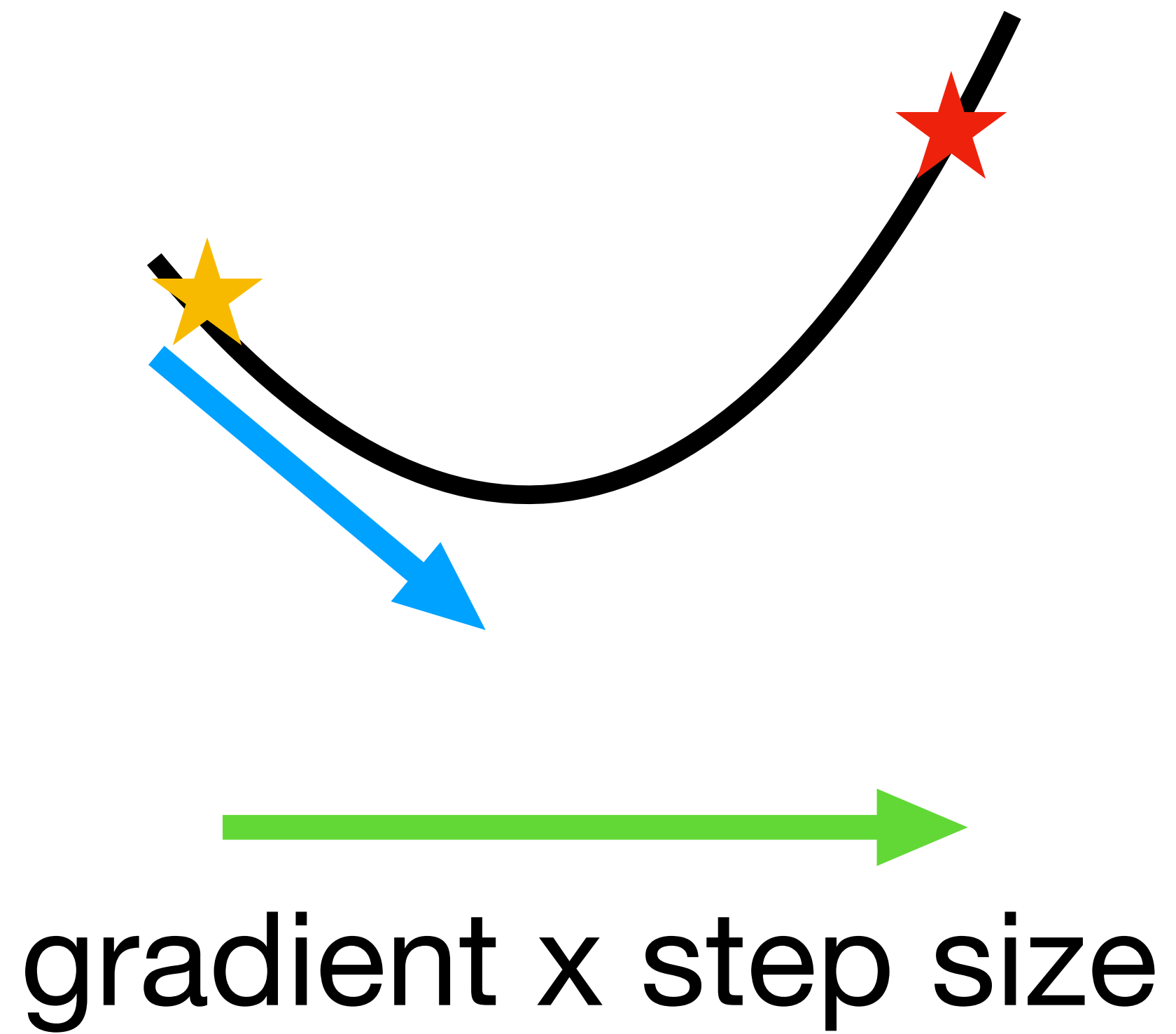


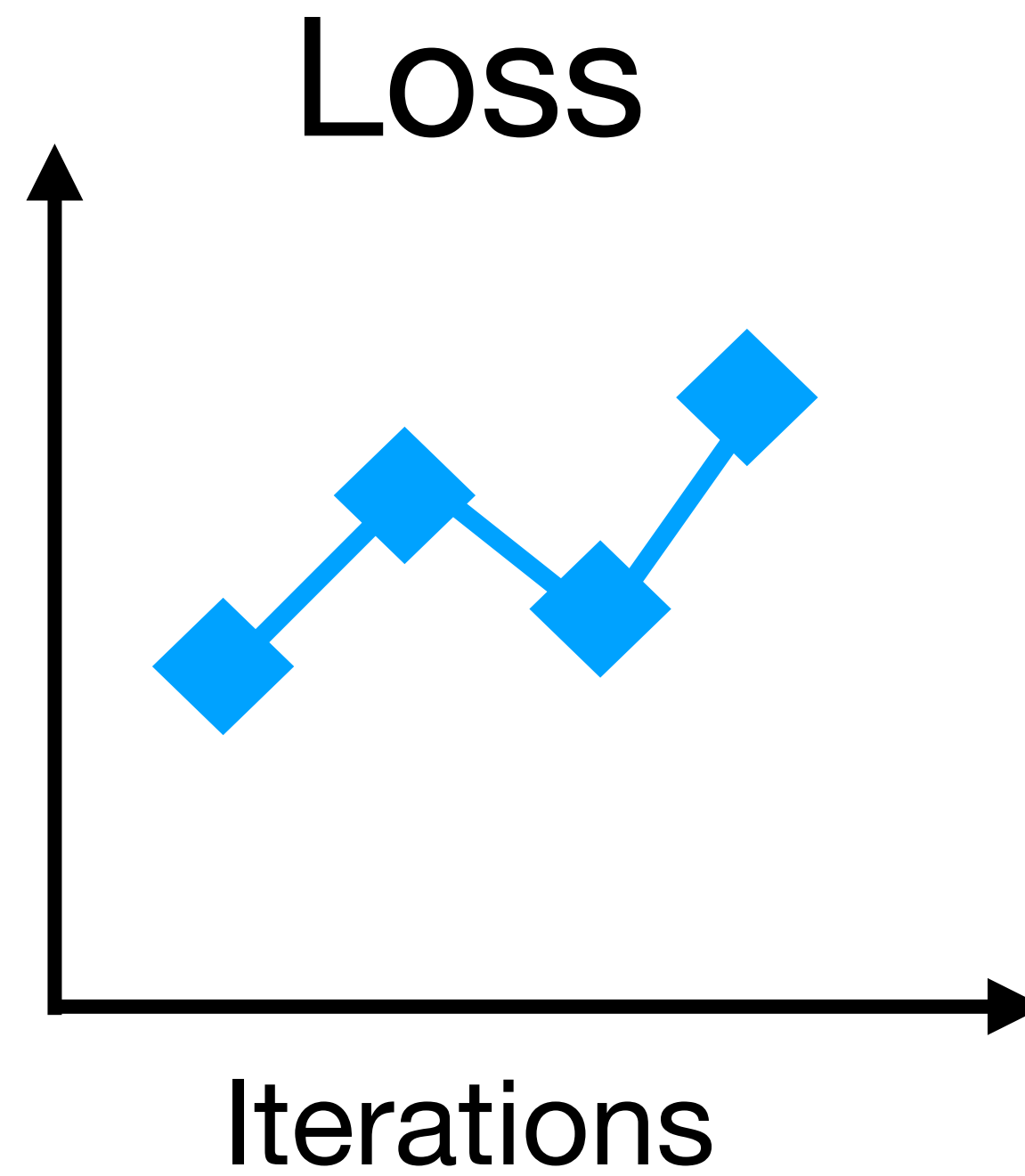
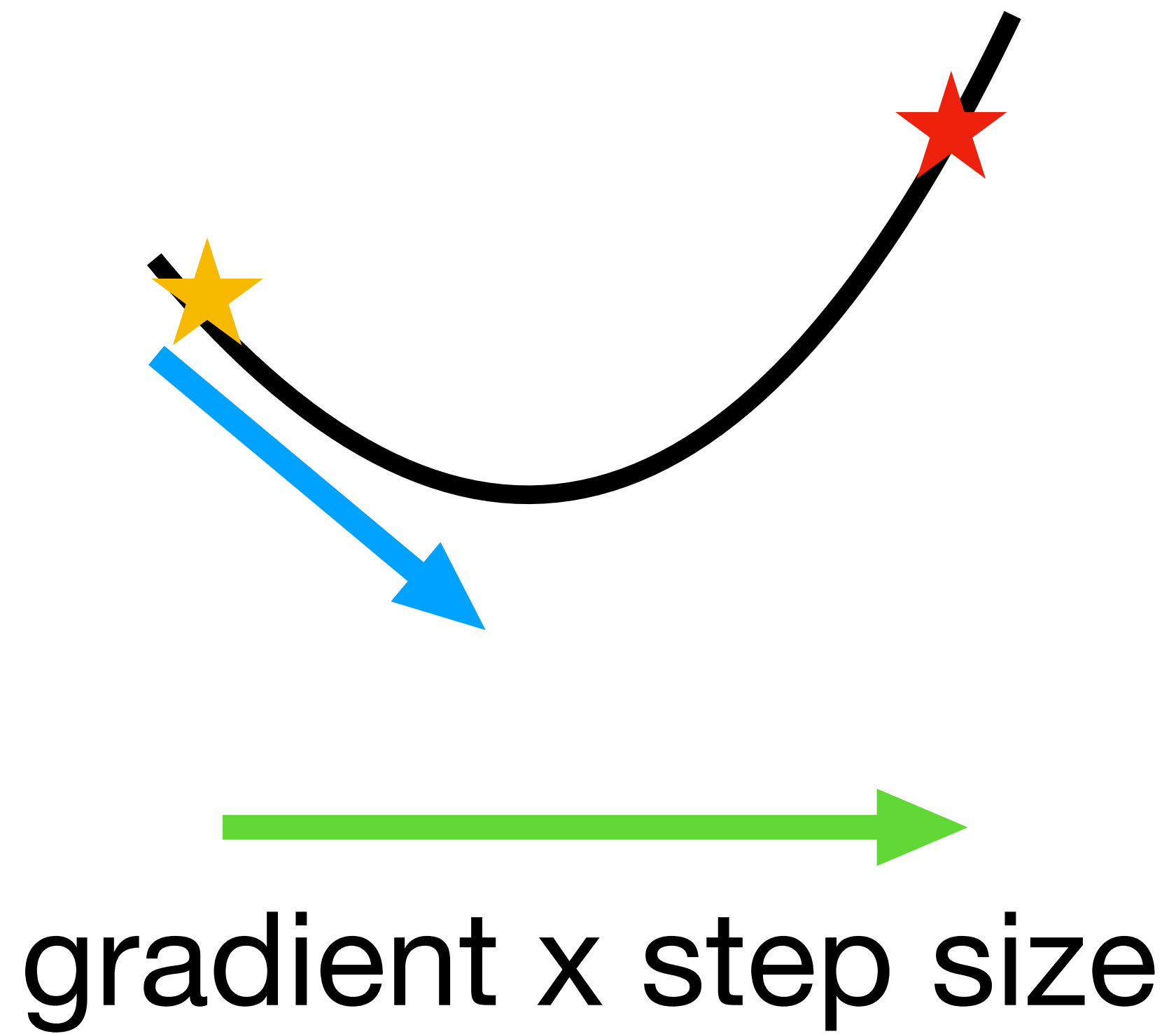


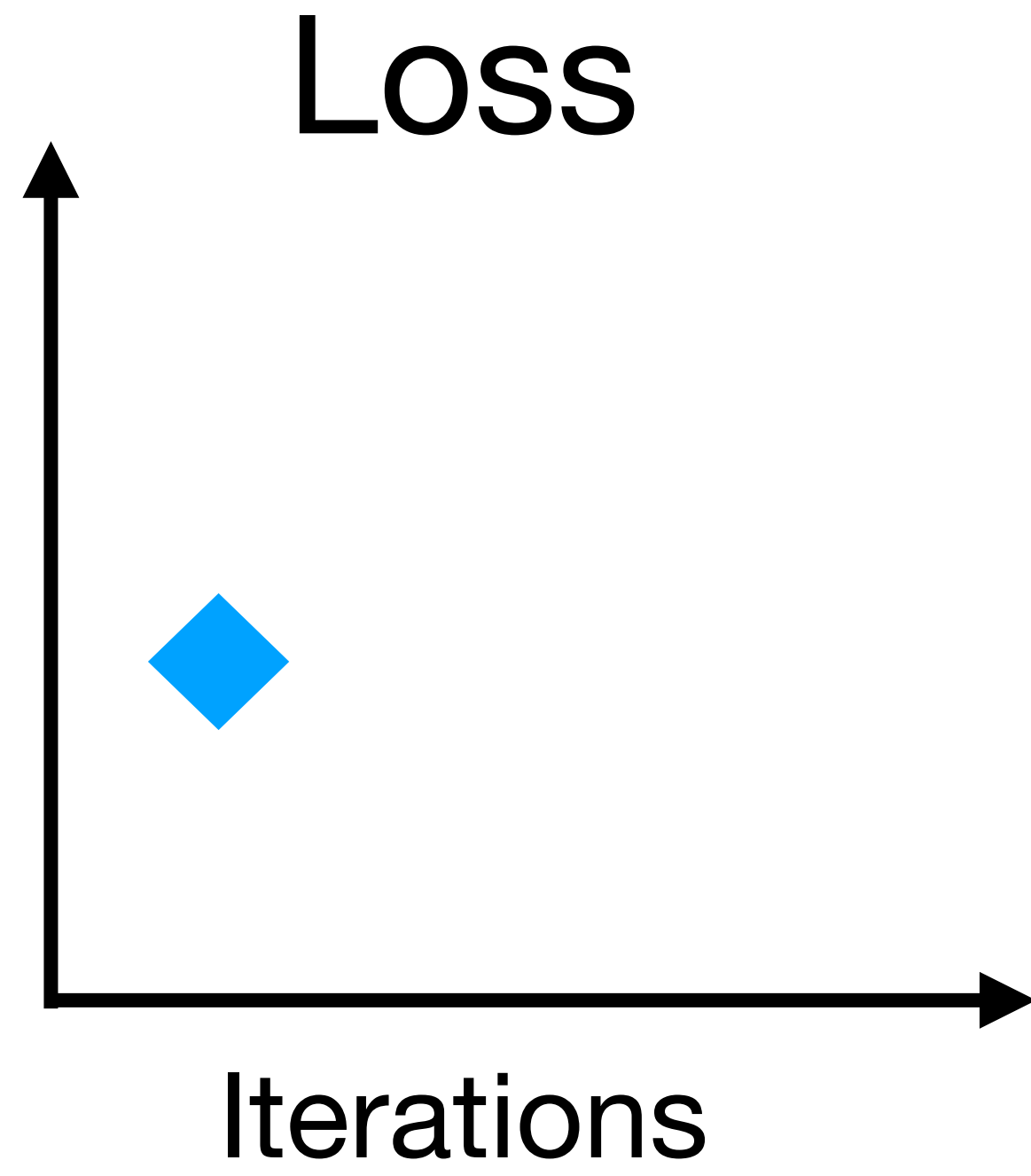
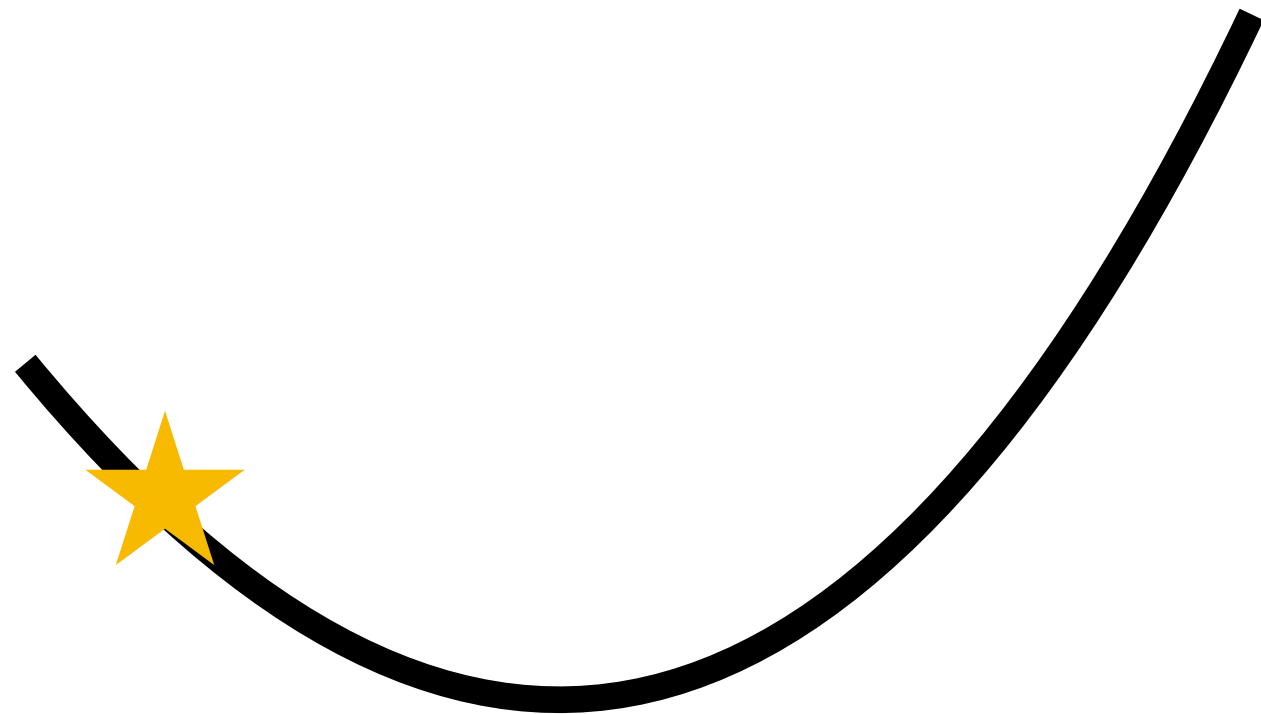


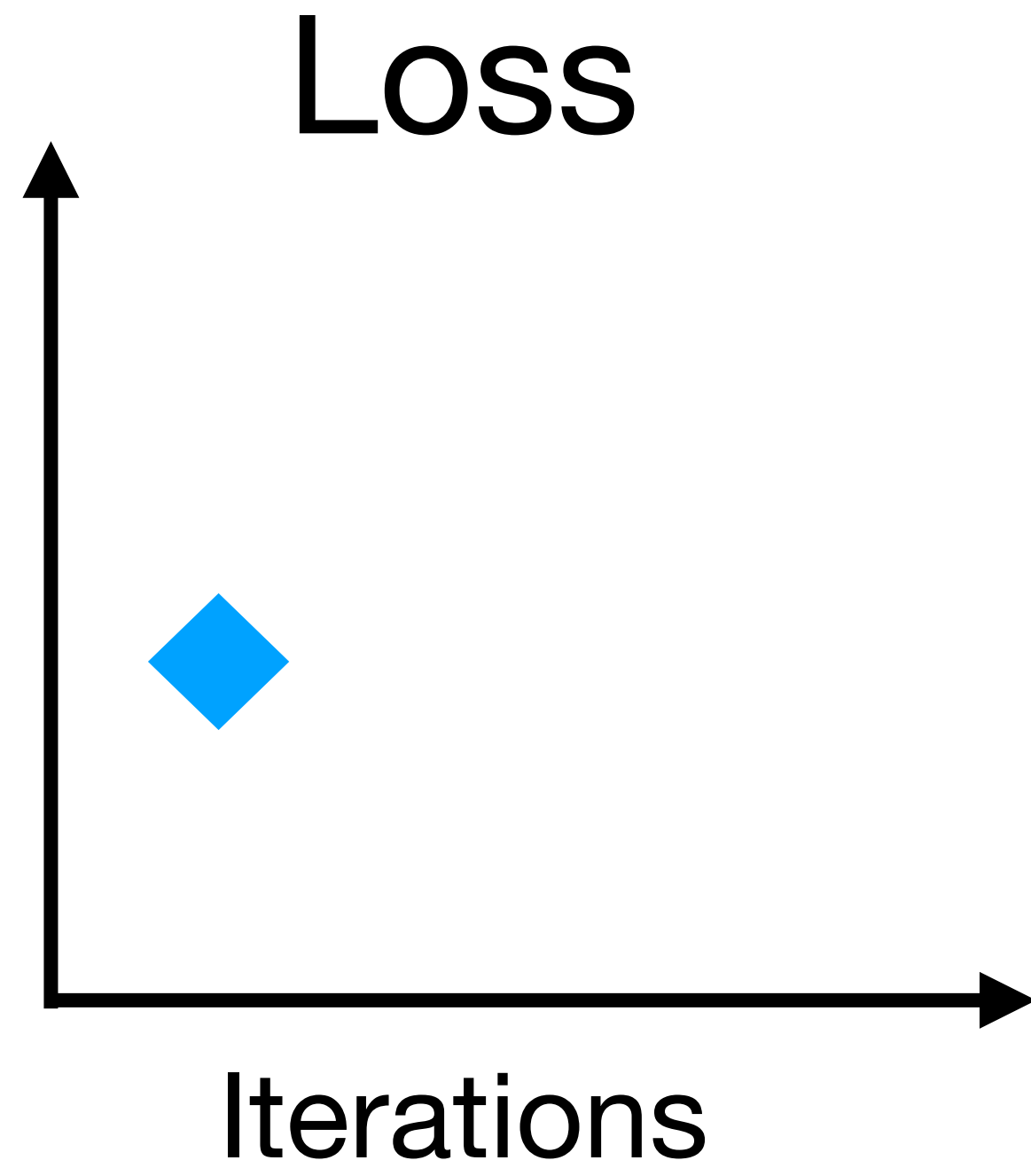
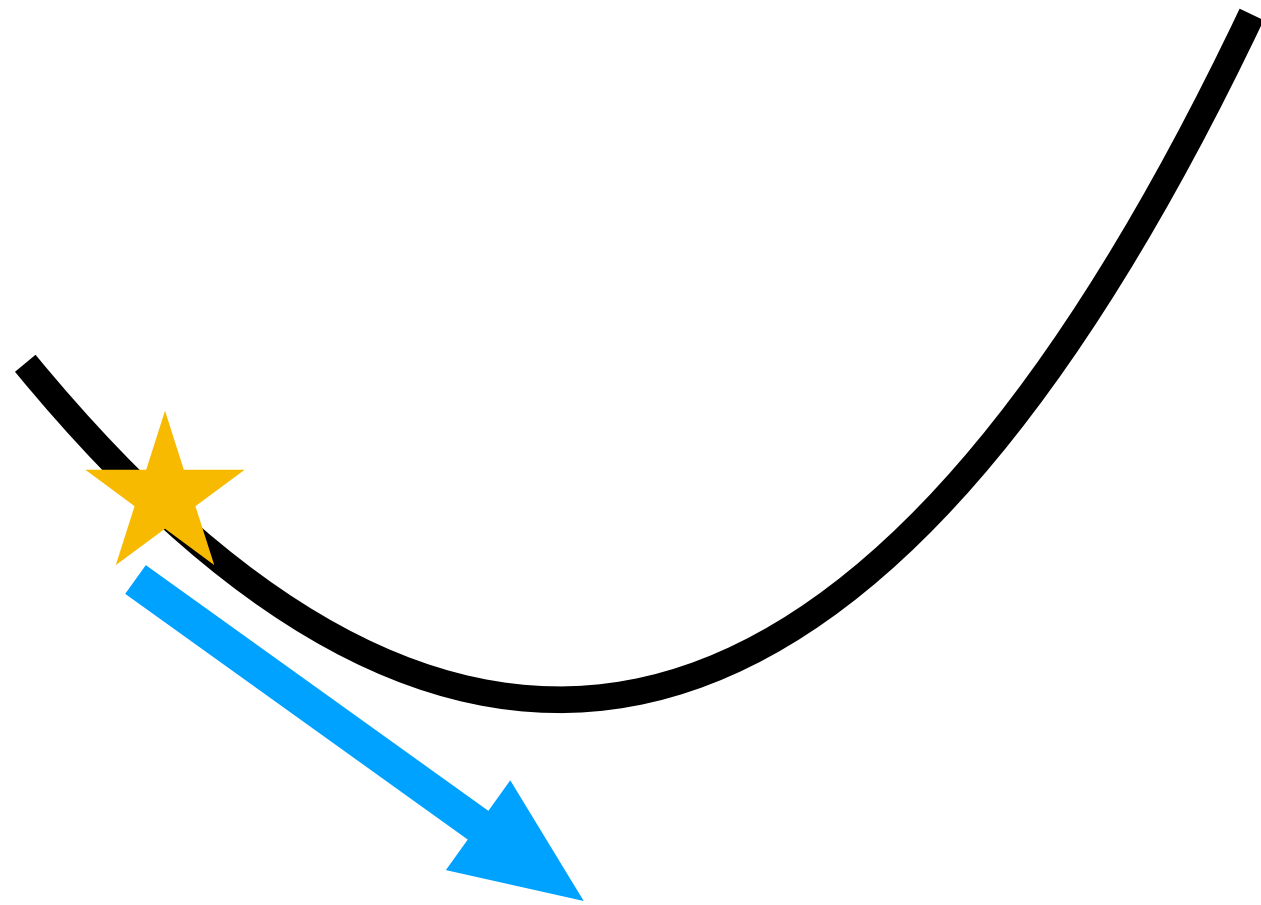


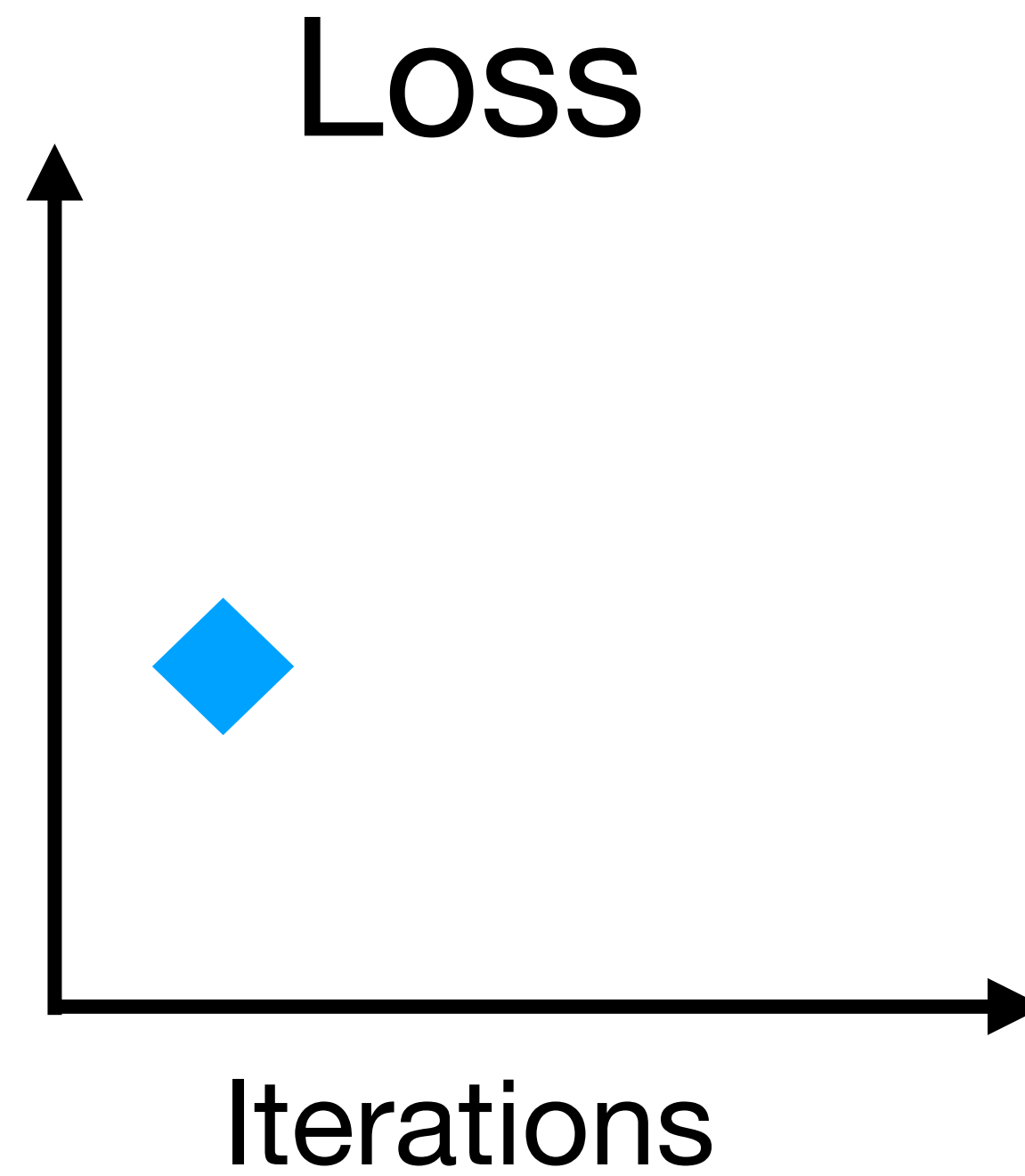
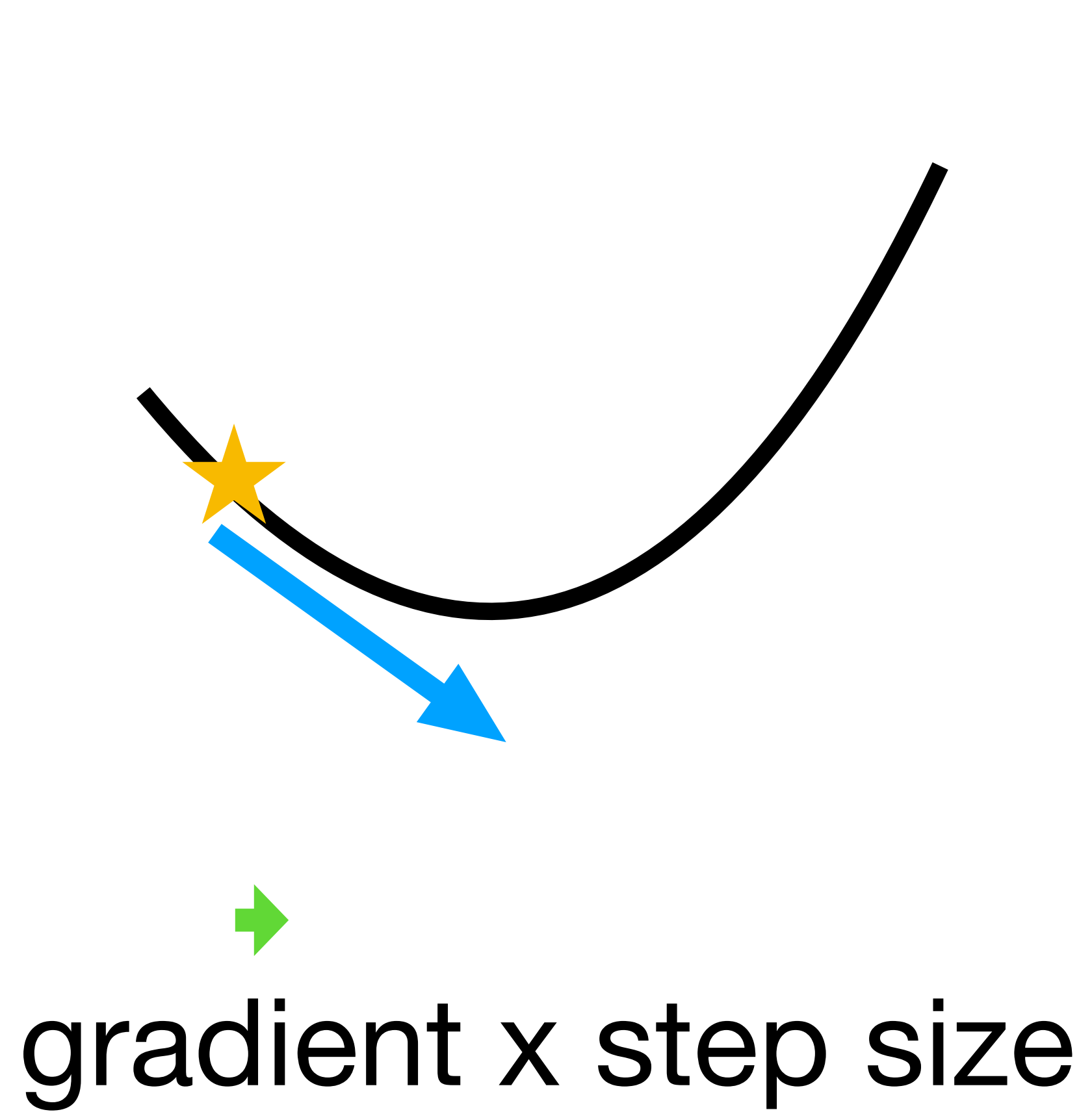


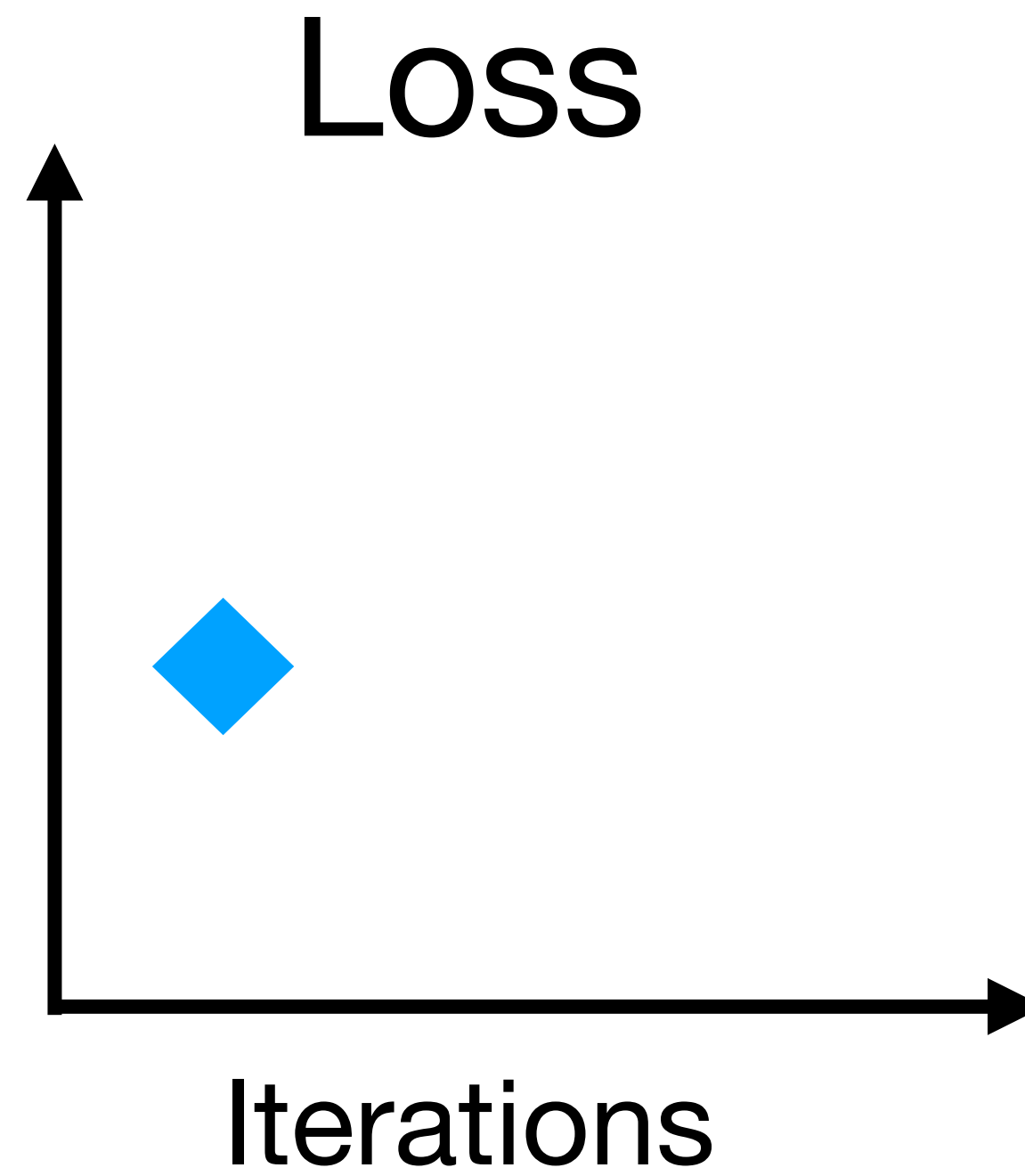
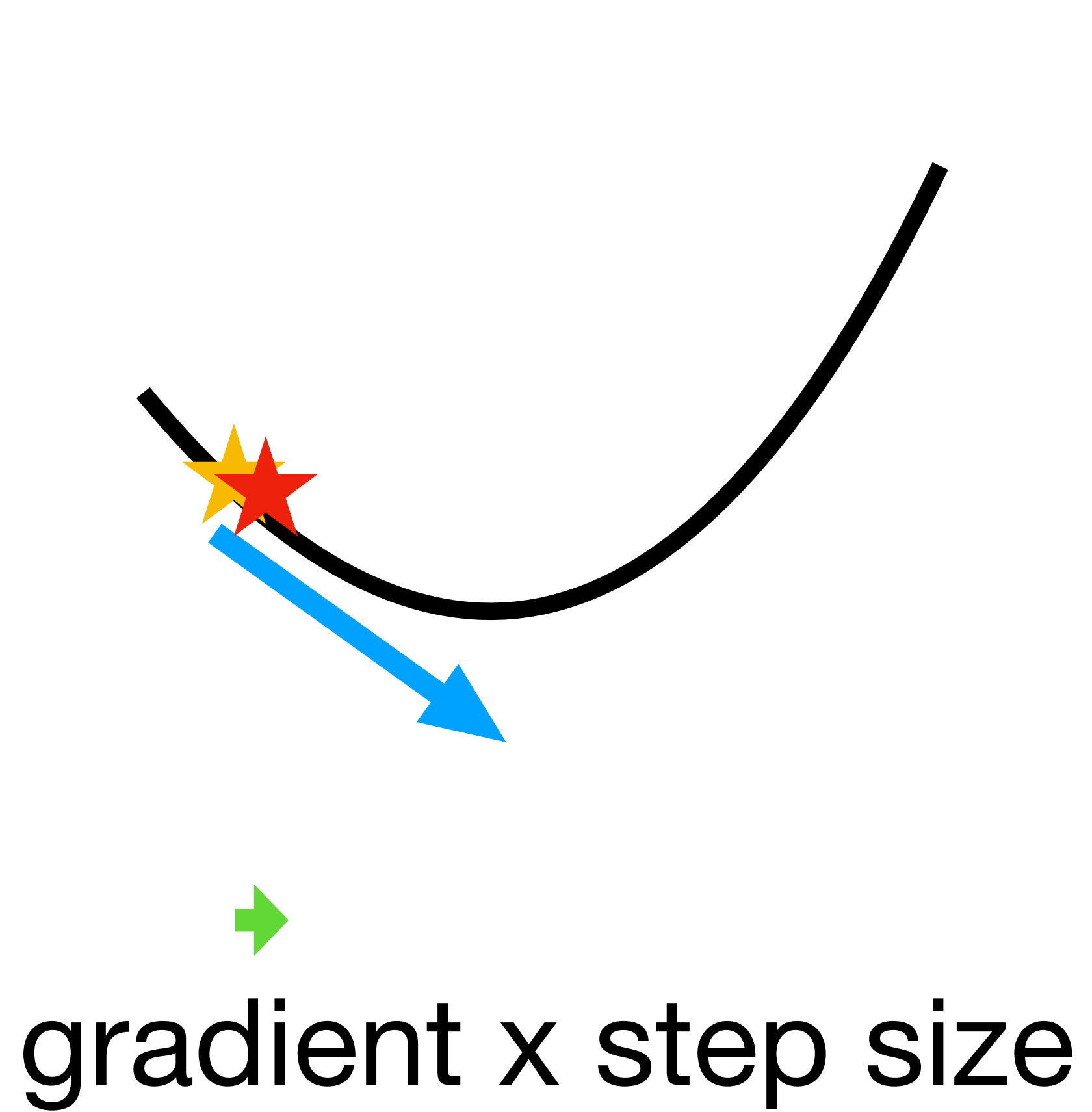


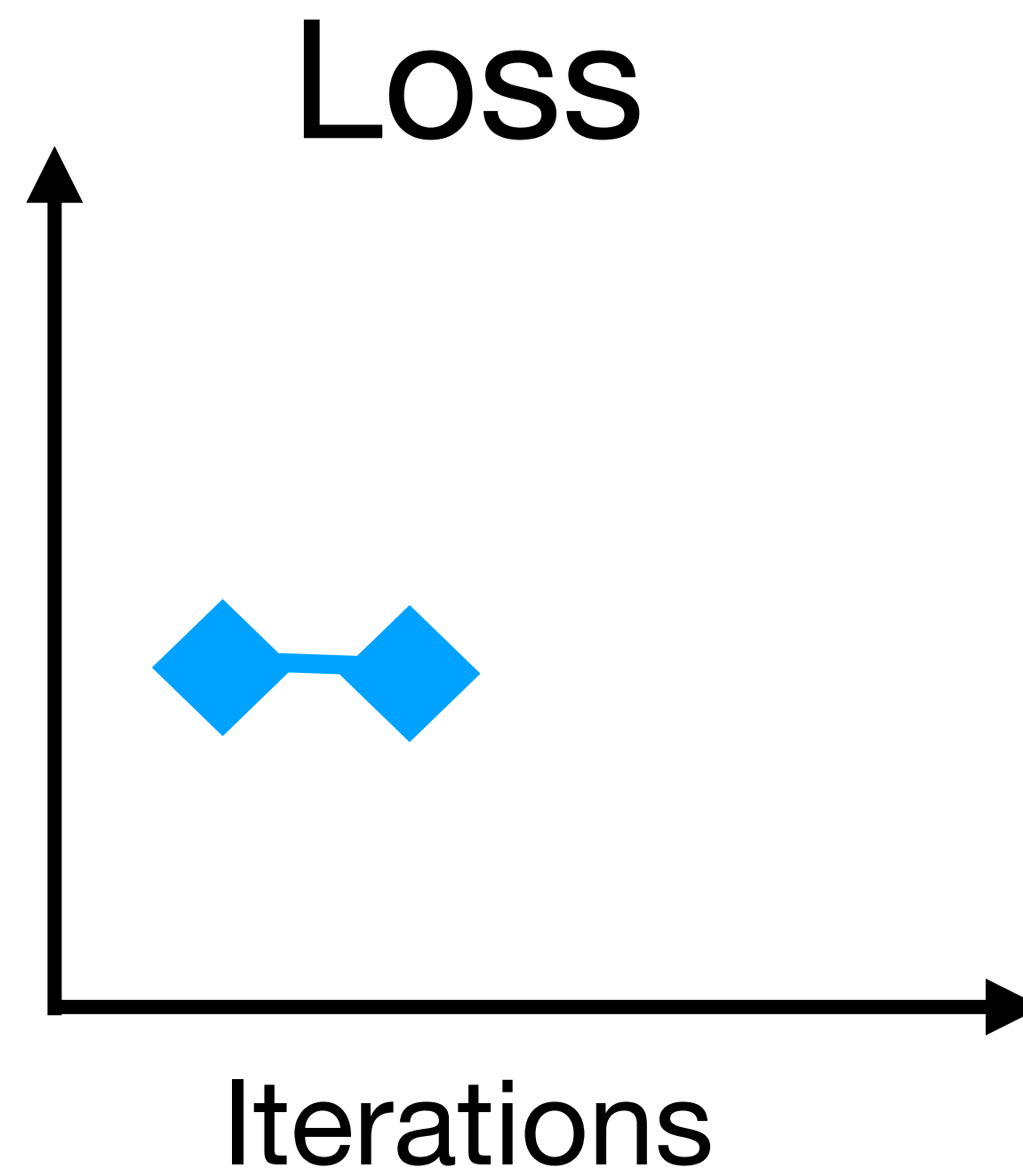
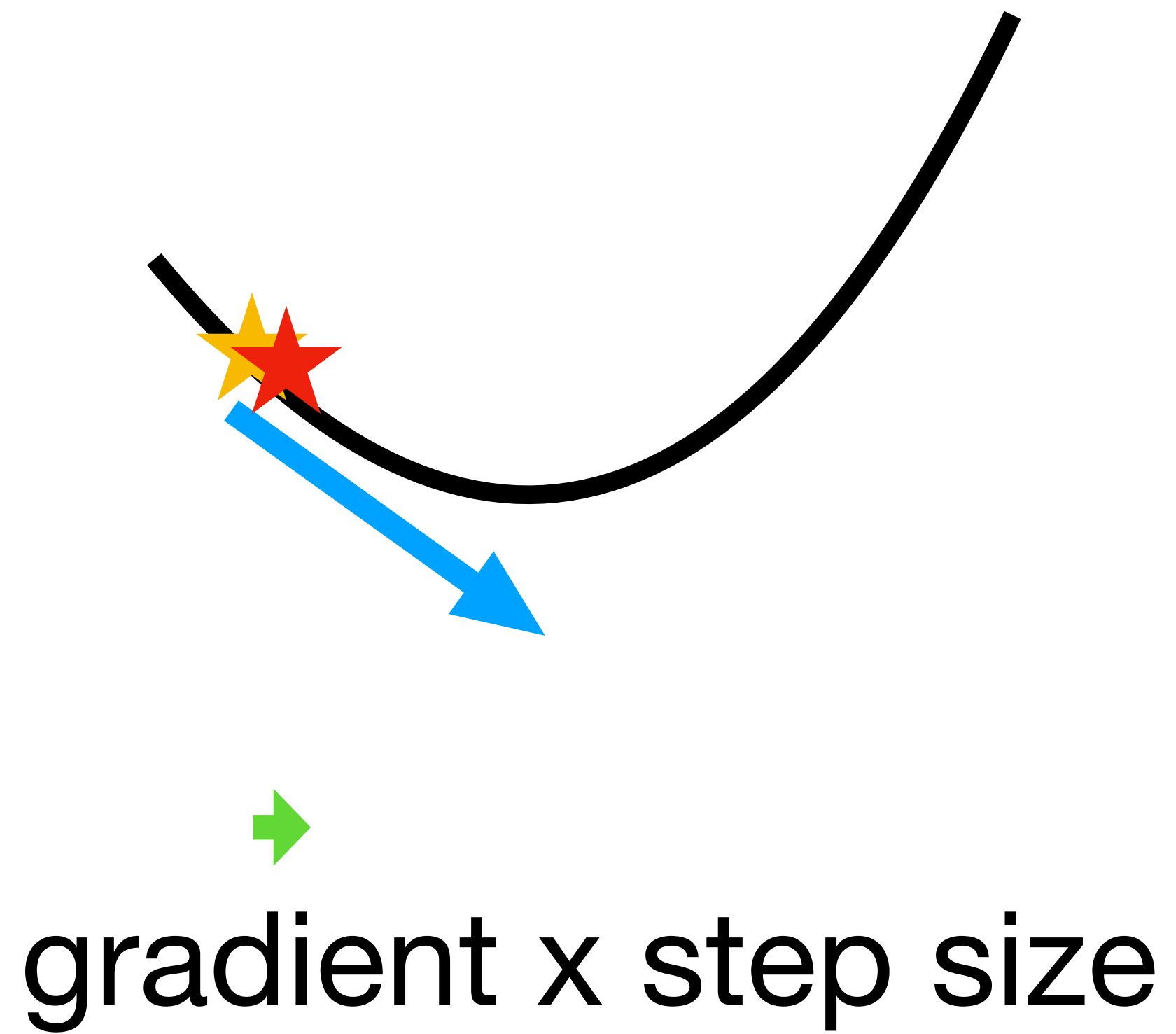




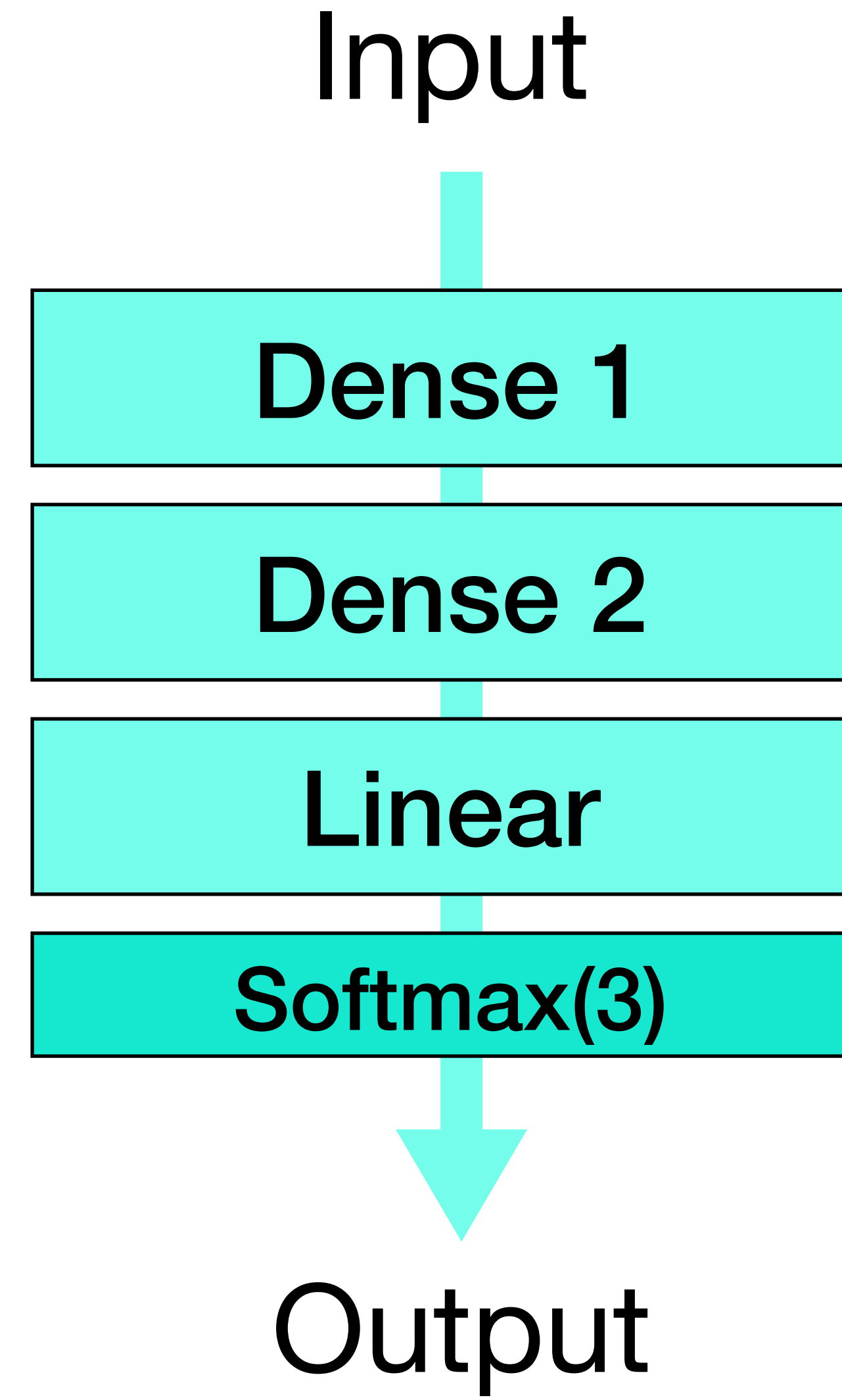


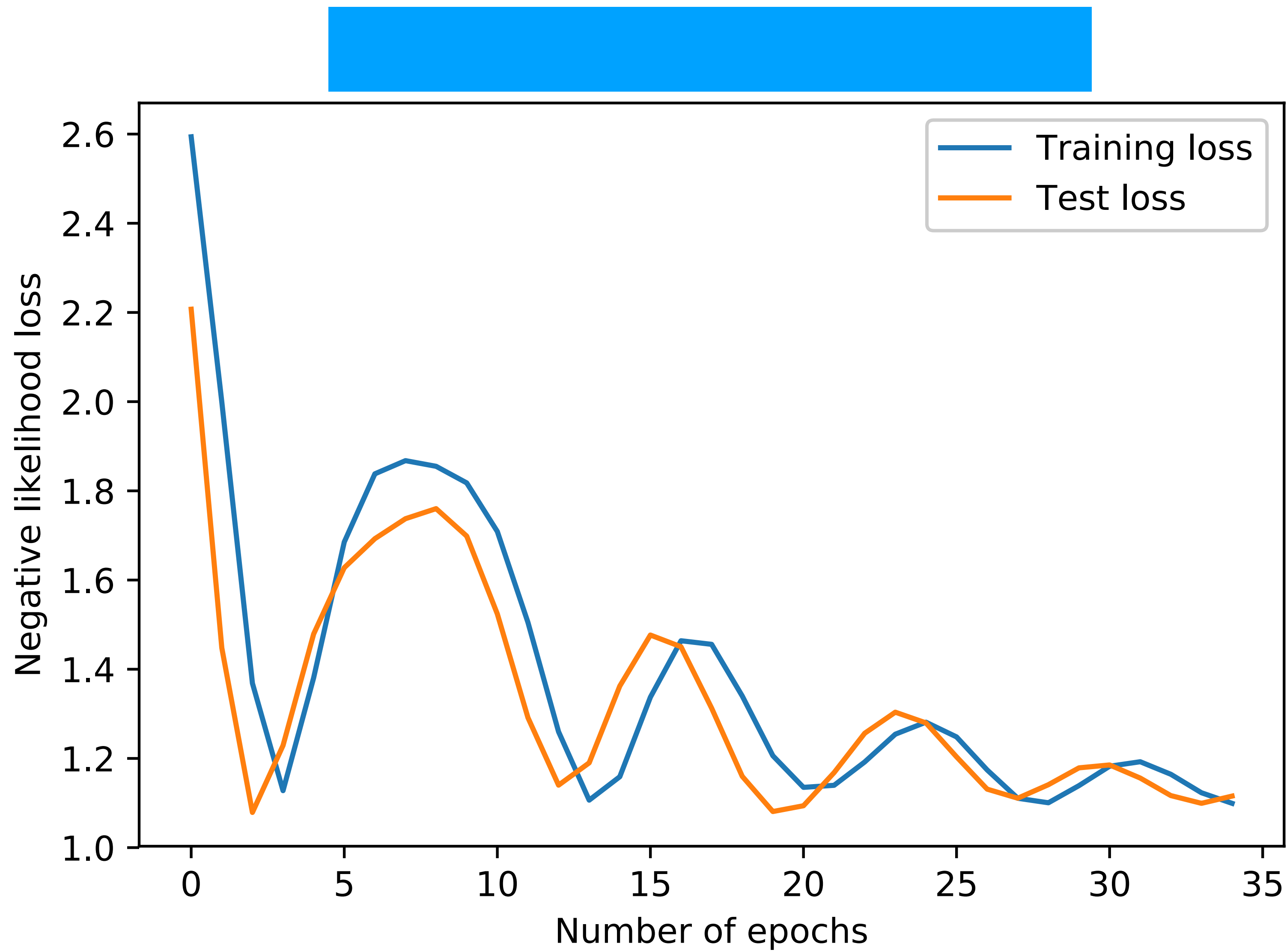




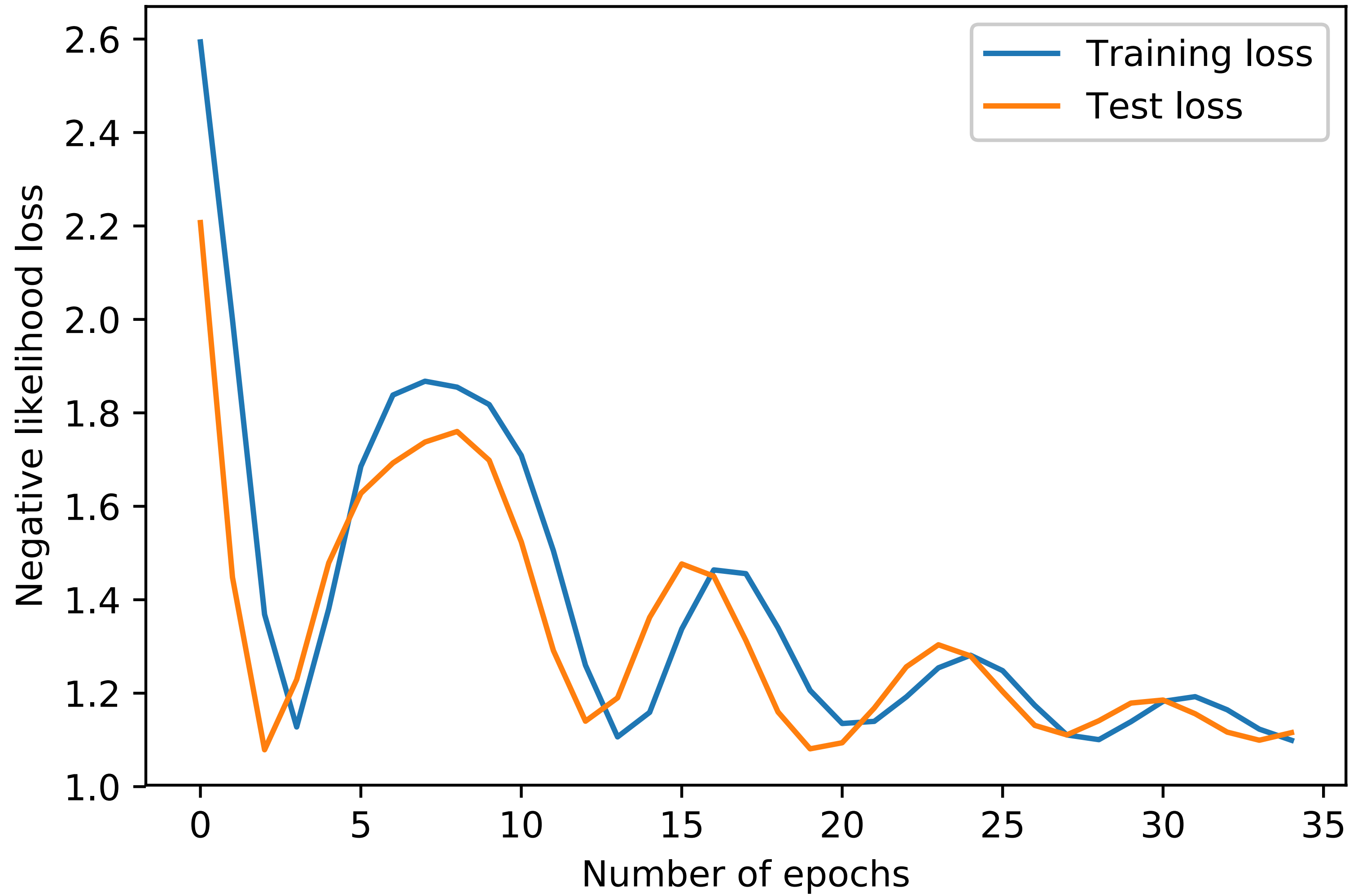


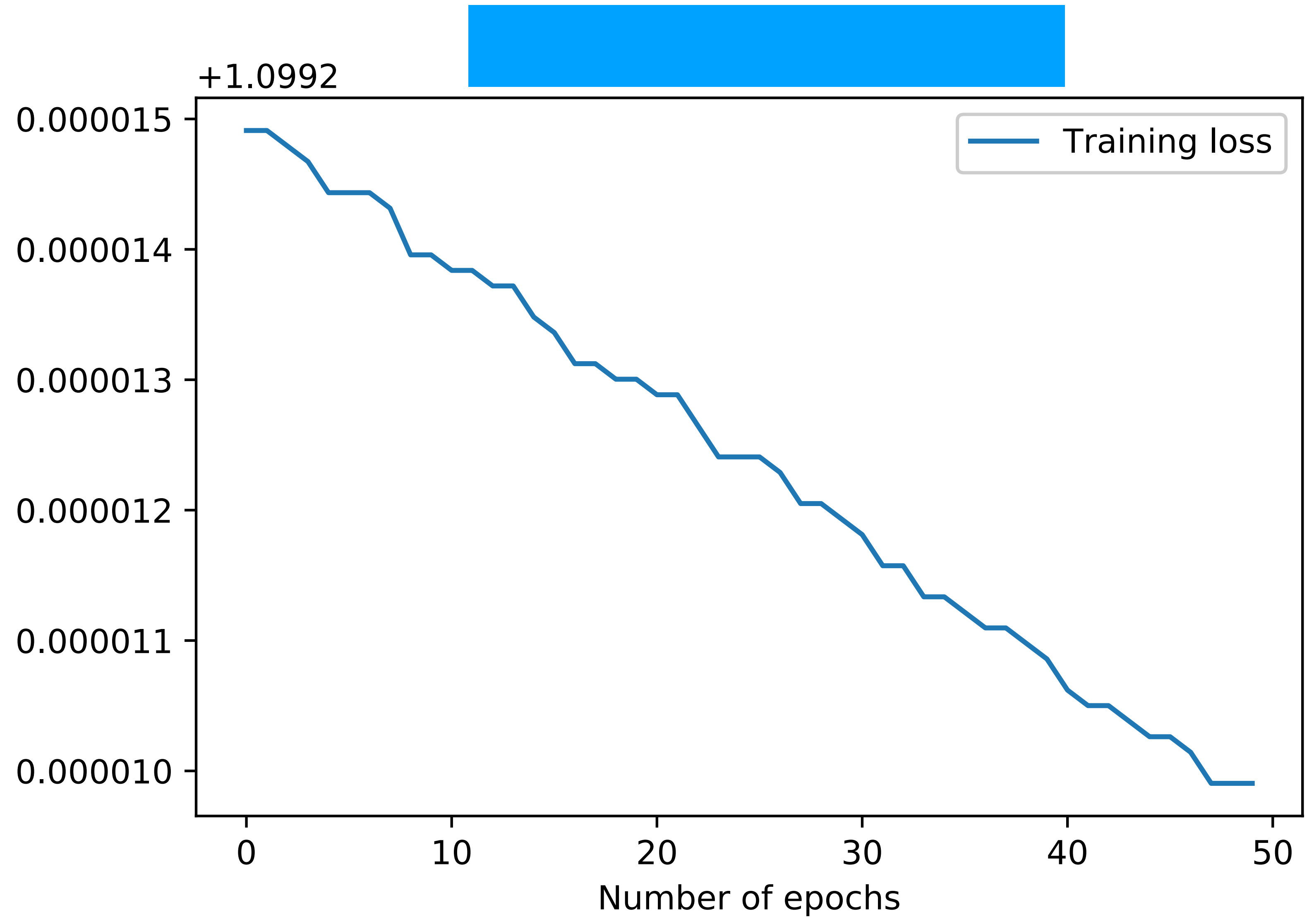
Iris dataset

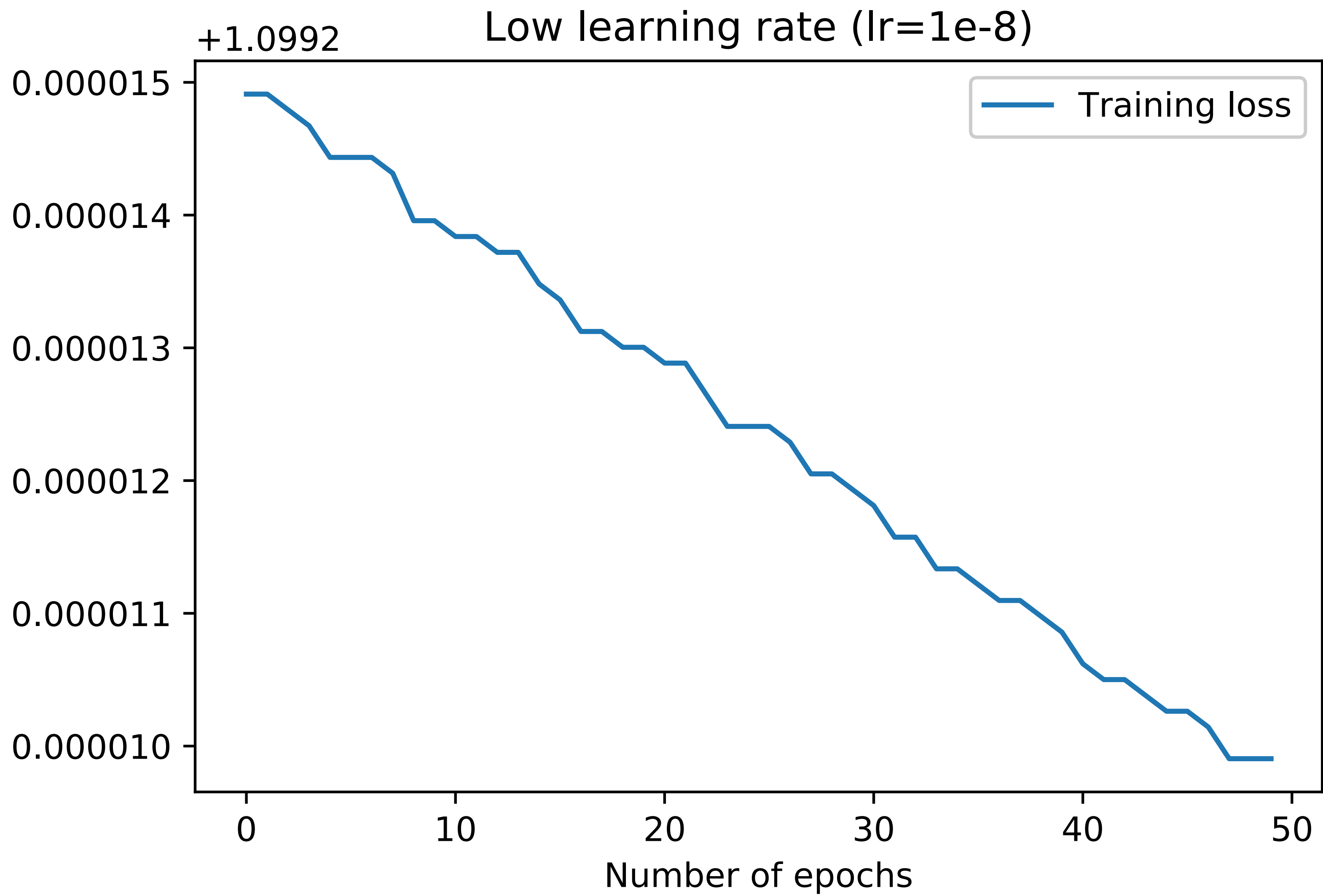




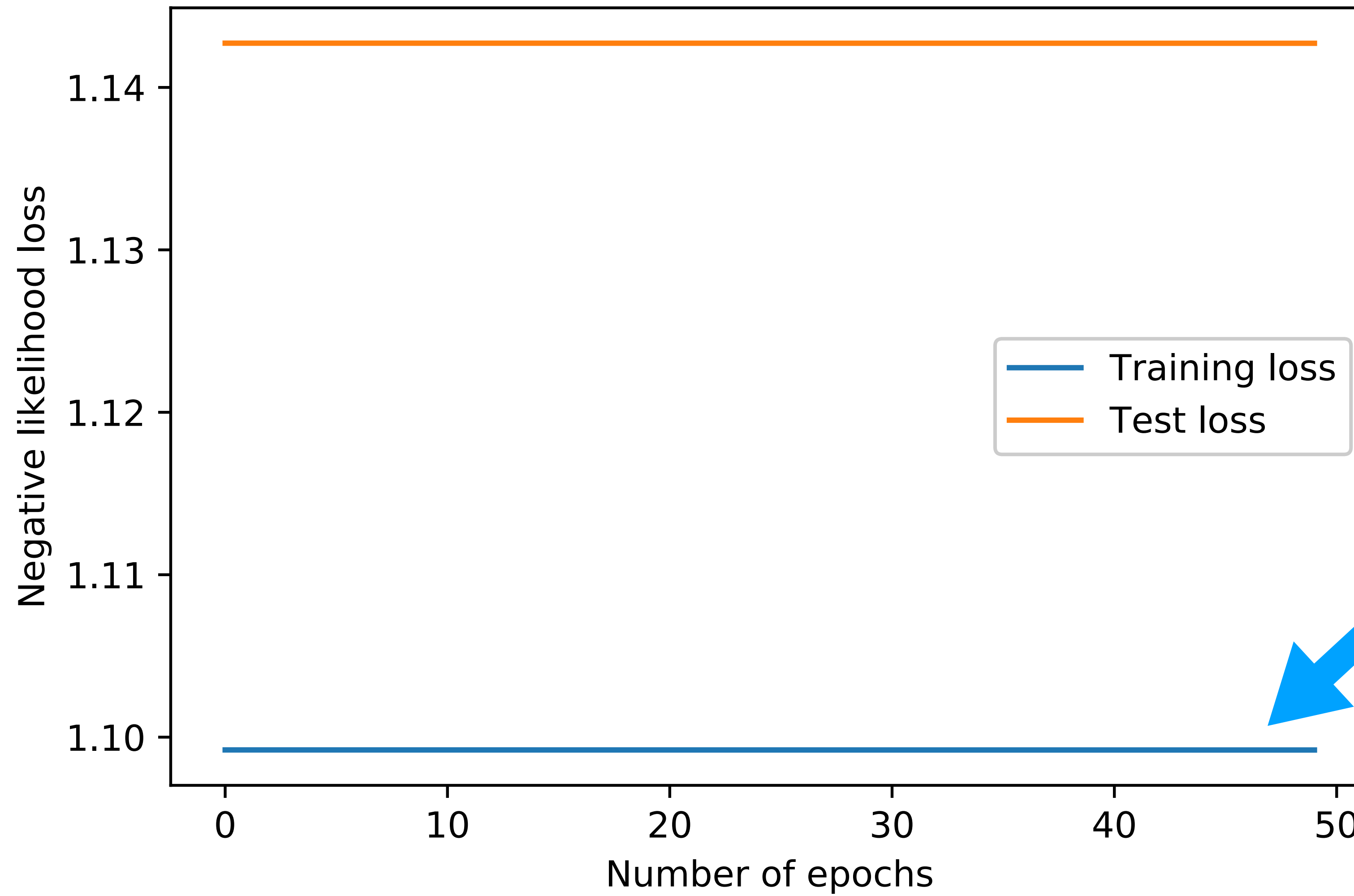
High learning rate loss curves(lr=2)

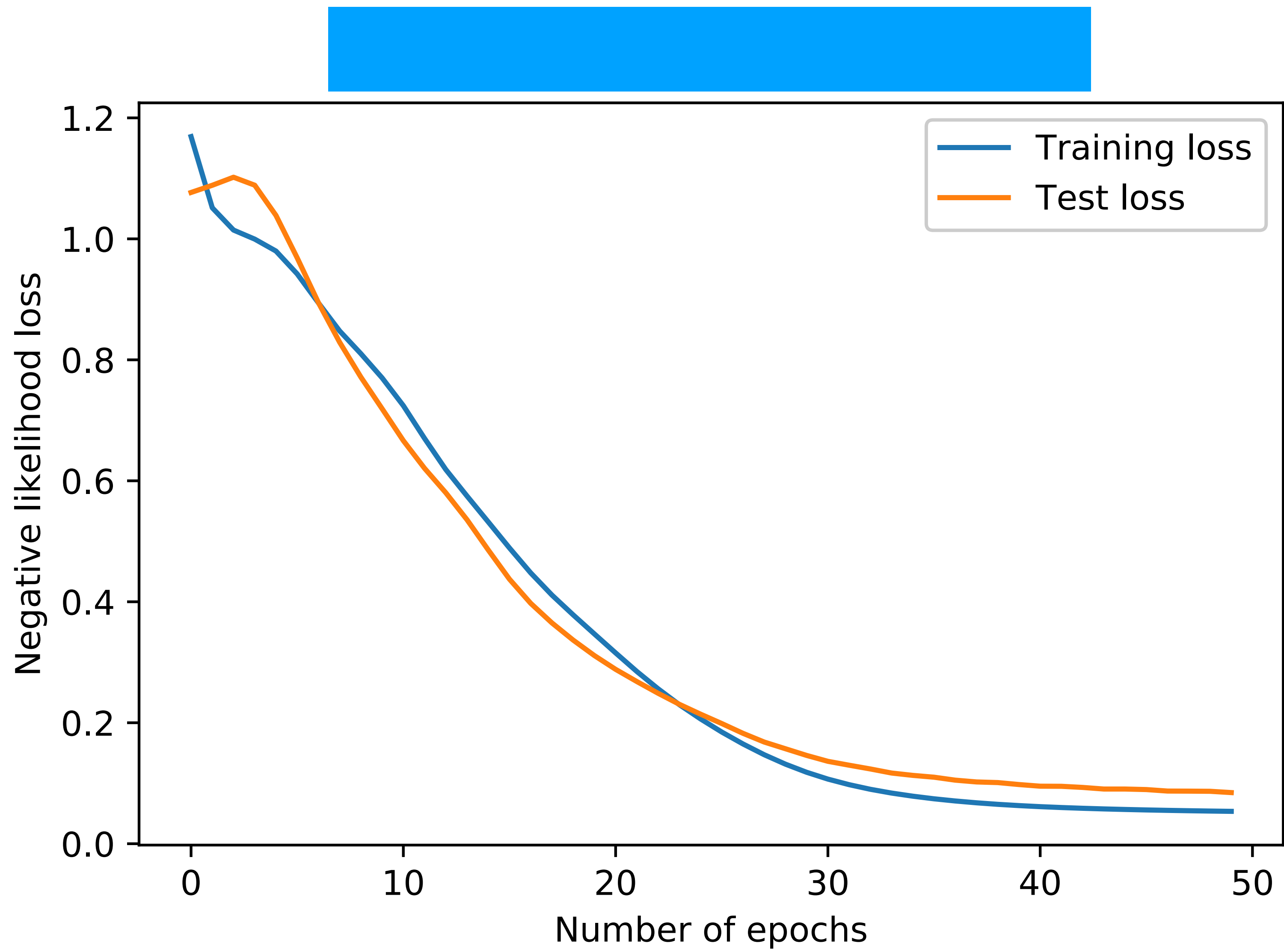




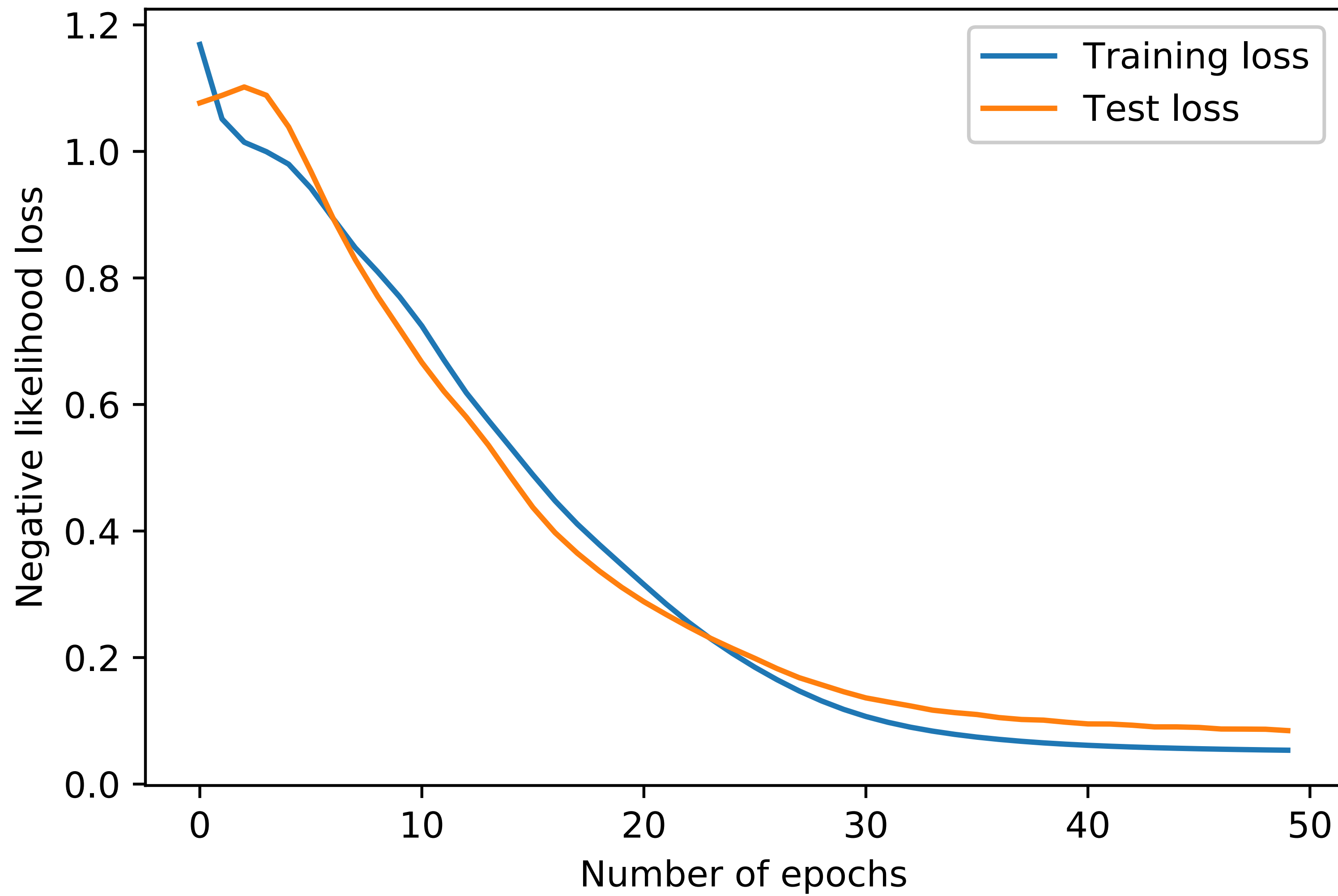


Low learning rate (lr=1e-8)

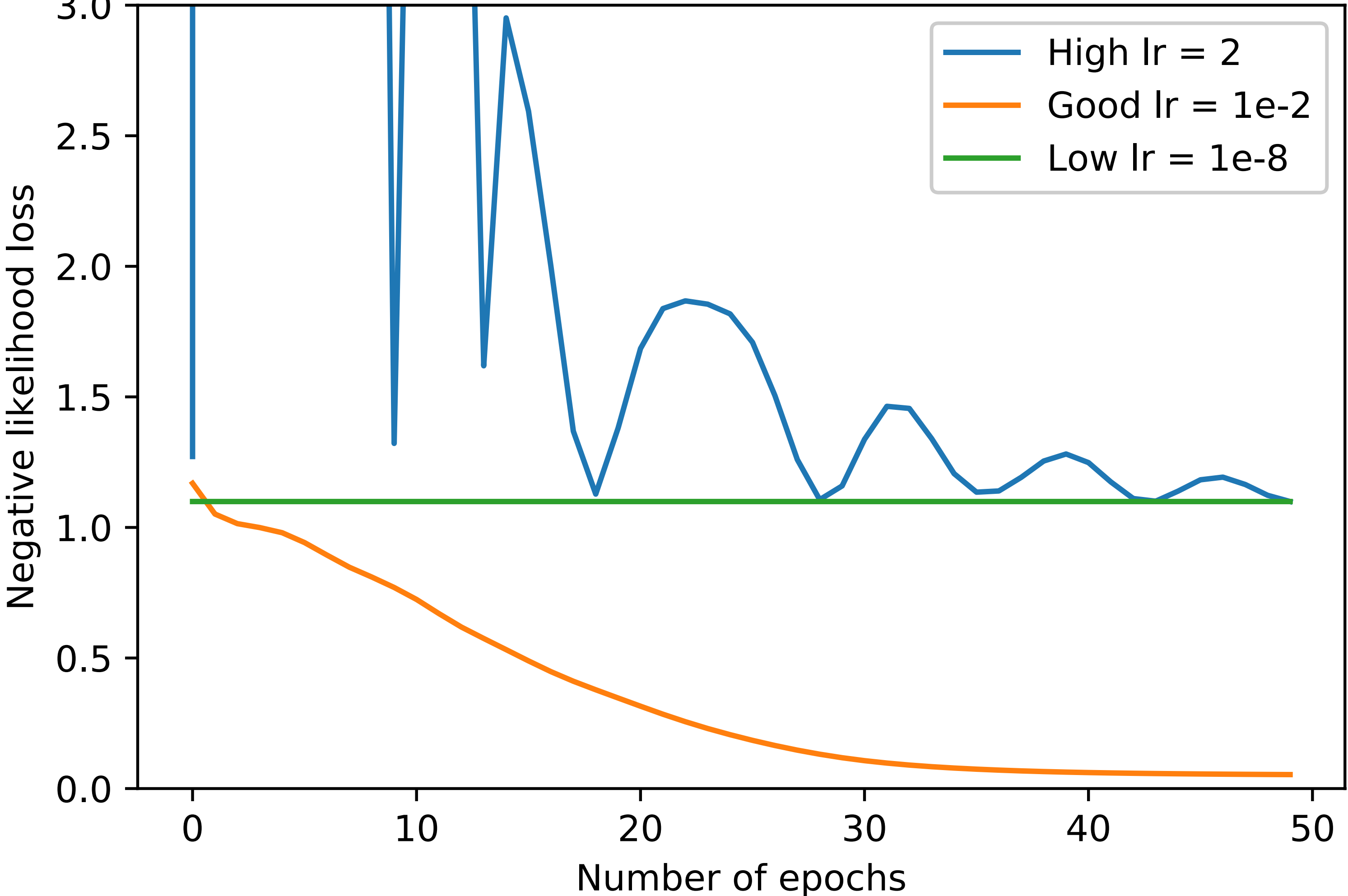




Good loss curves (lr=1e-2)



Training losses at different learning rates

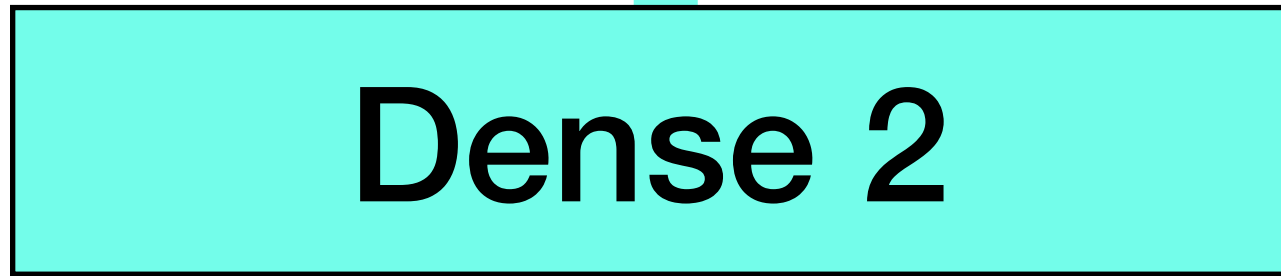
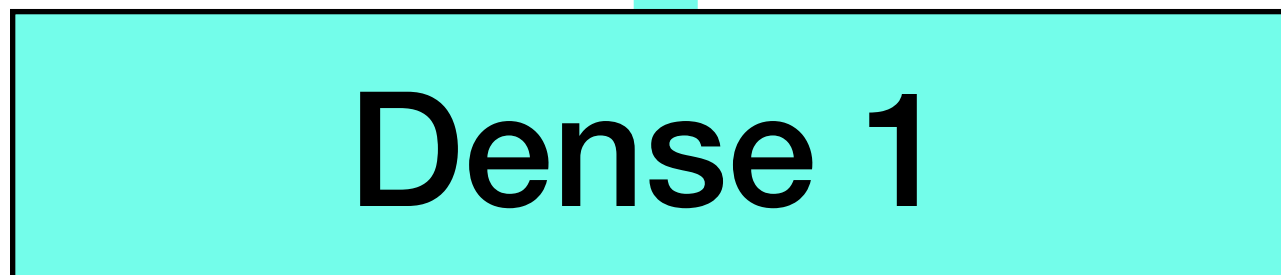


**How calculate
gradient?**

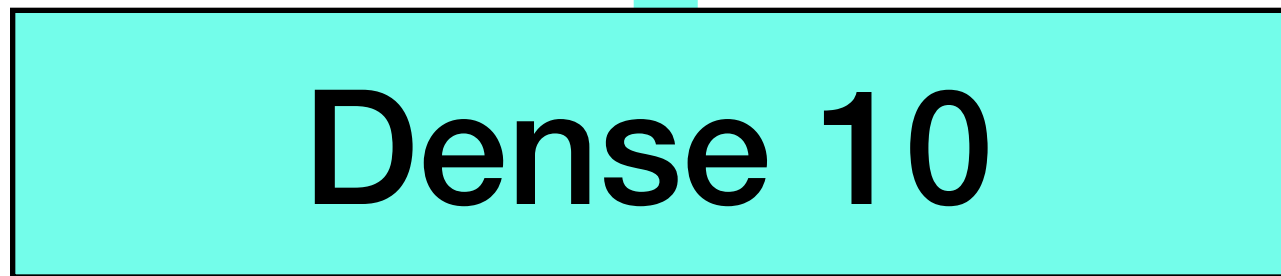
**How calculate
gradient?**



Input

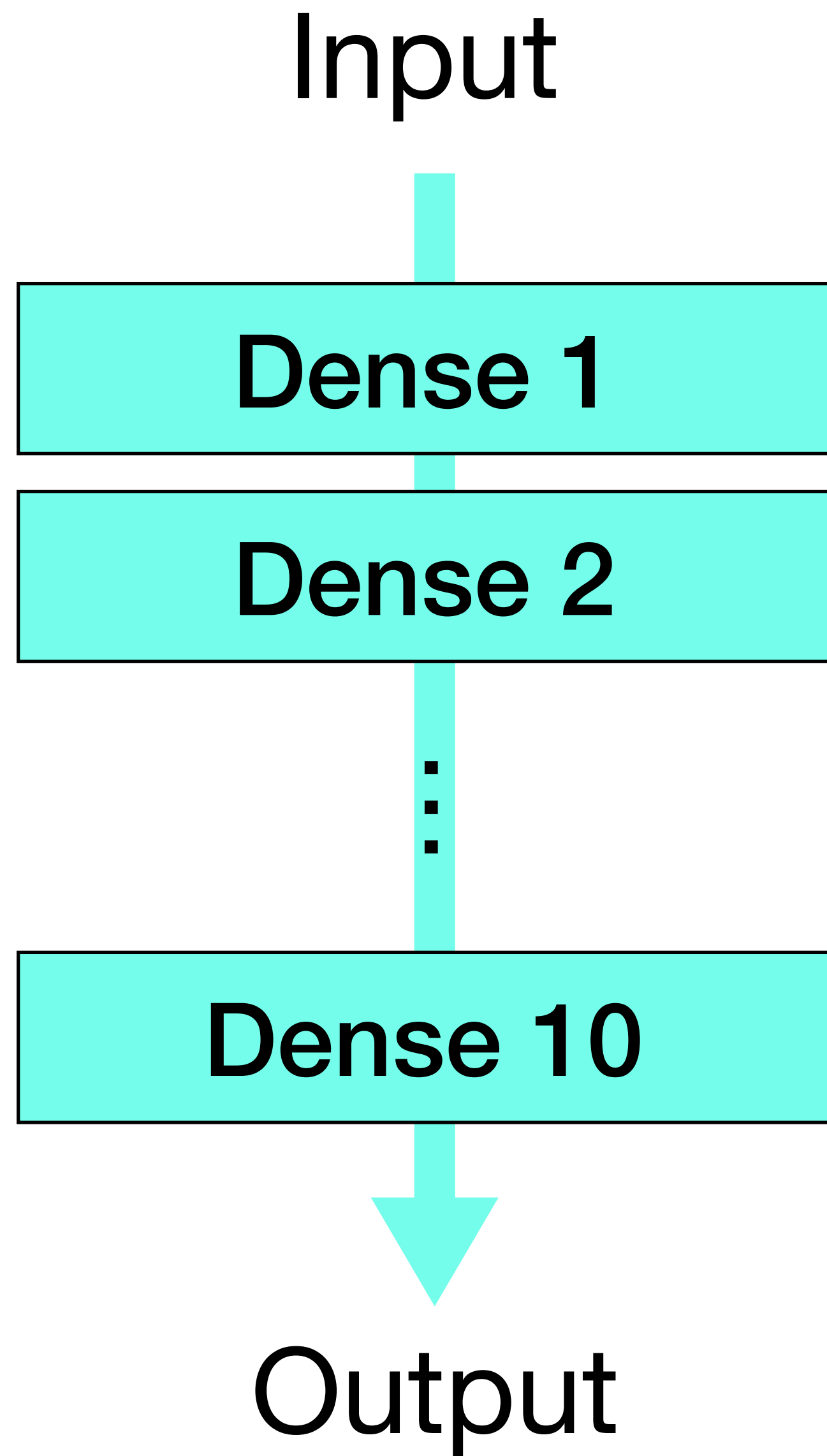


⋮

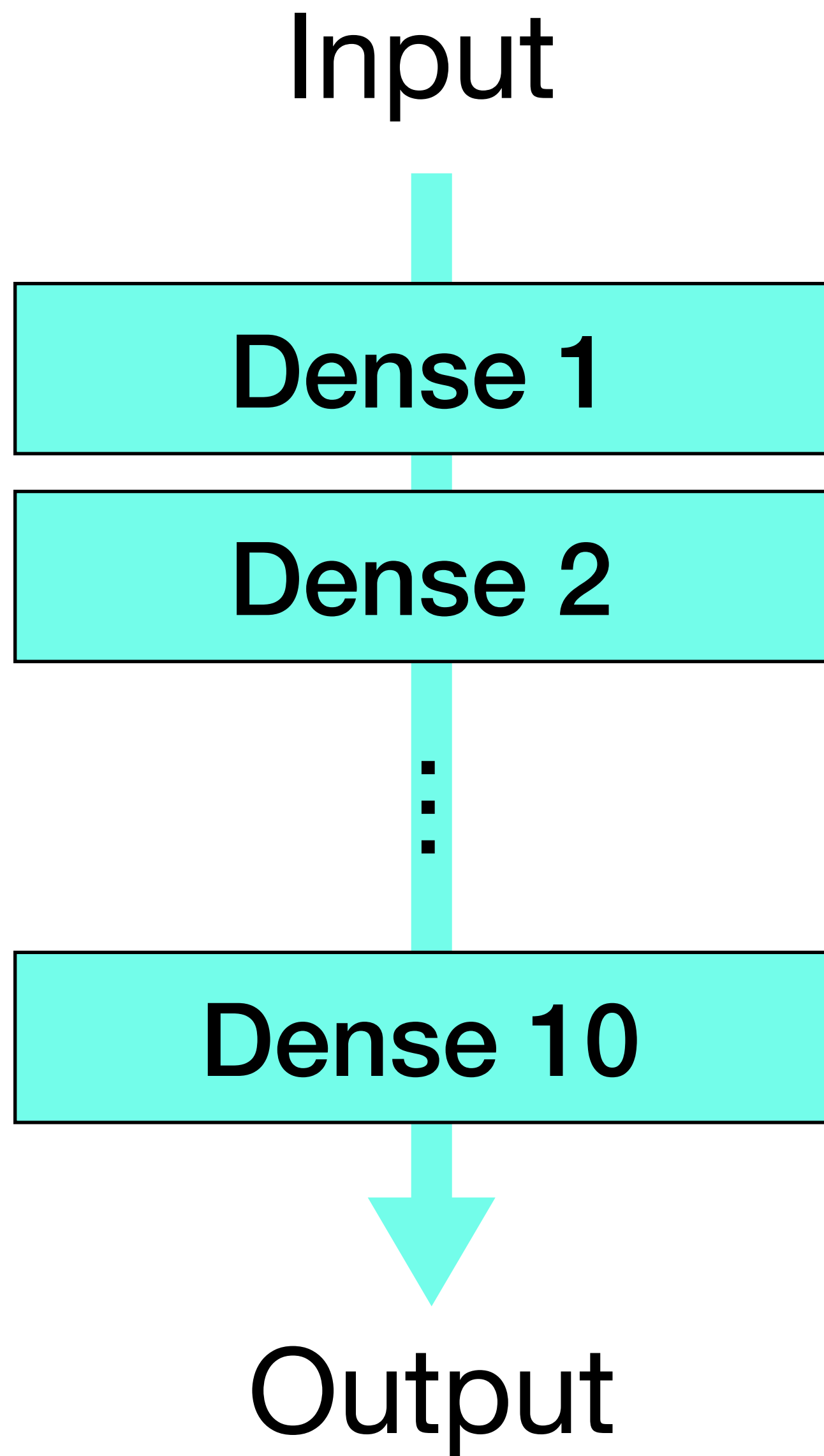


Output

$$\frac{\partial \text{Loss}}{\partial W_1}$$



$$\frac{\partial \text{Loss}}{\partial W_1}$$
$$= \frac{\partial \text{Loss}}{\partial O_{10}} \frac{\partial}{\partial W_1} (\sigma(W_{10} \sigma(W_9 \sigma(\dots))))$$

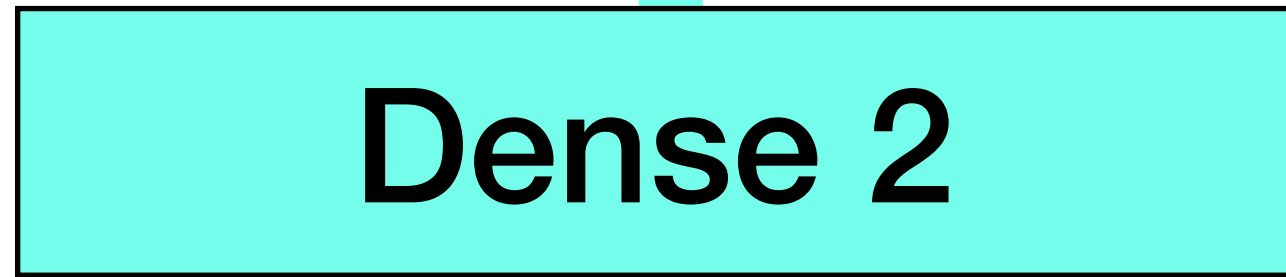
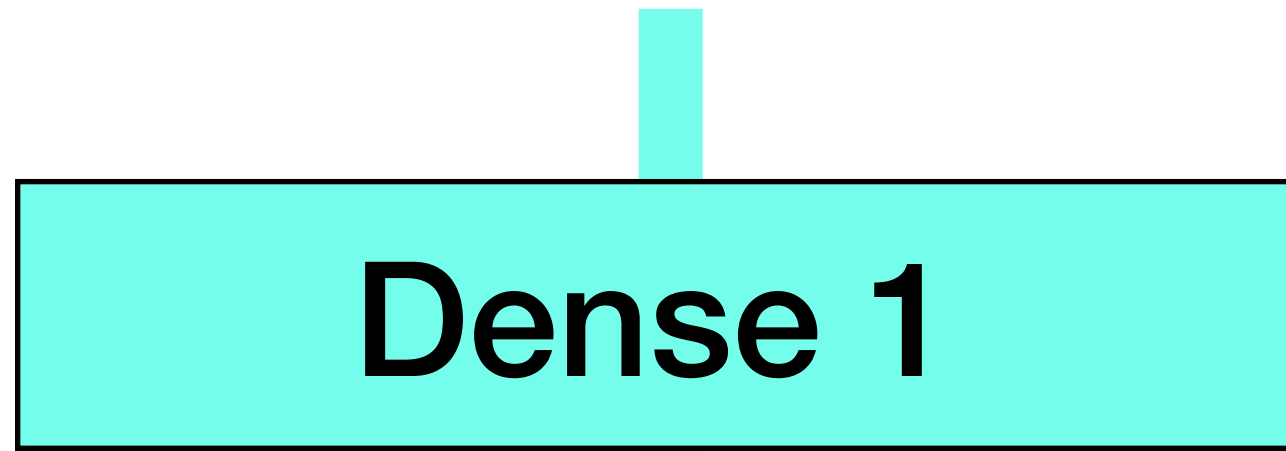


$$\frac{\partial \text{Loss}}{\partial W_1}$$

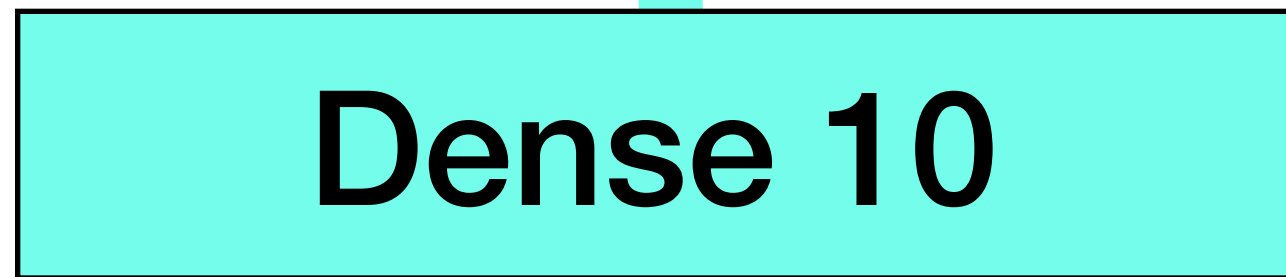
$$= \frac{\partial \text{Loss}}{\partial O_{10}} \frac{\partial}{\partial W_1} (\sigma(W_{10} \sigma(\dots \sigma(W_1 \sigma(\dots \sigma(x))))))$$

CENSORED

Input



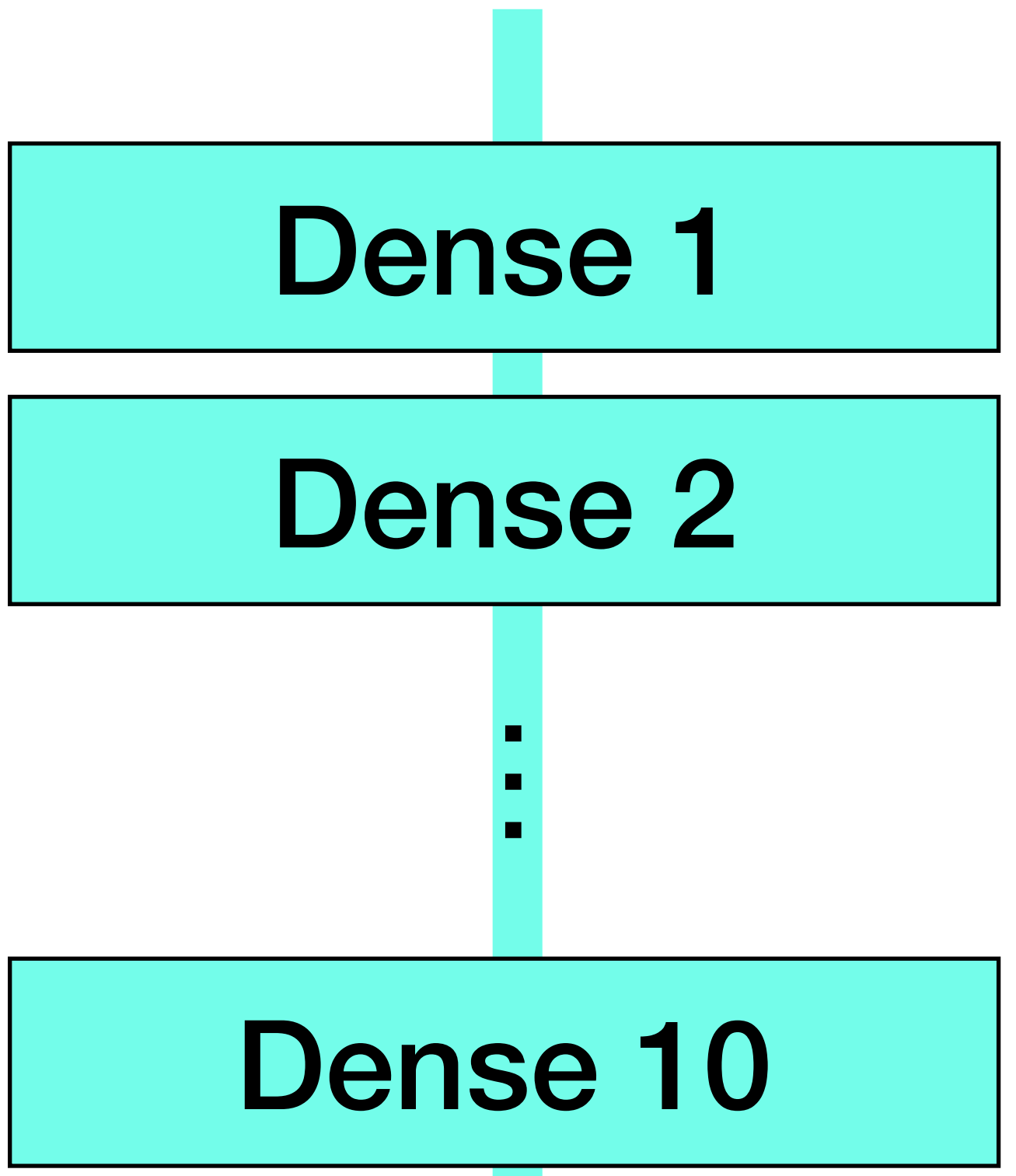
⋮



Output

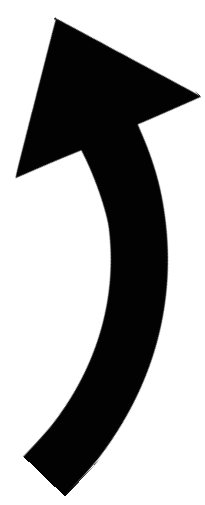
$$\frac{\partial \text{Loss}}{\partial W_1}$$

Input

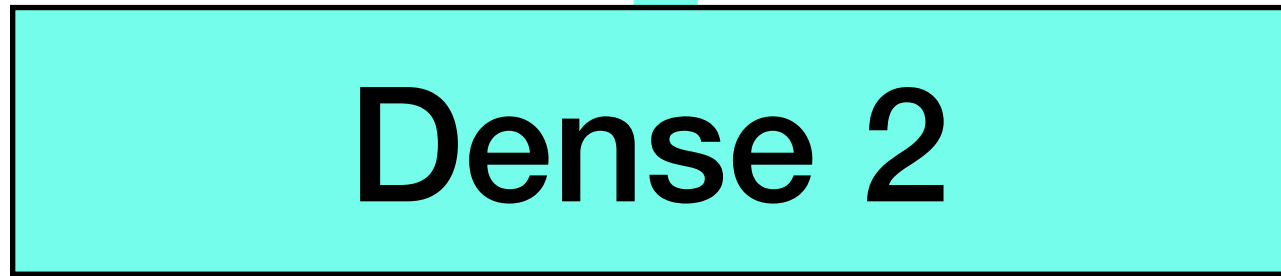
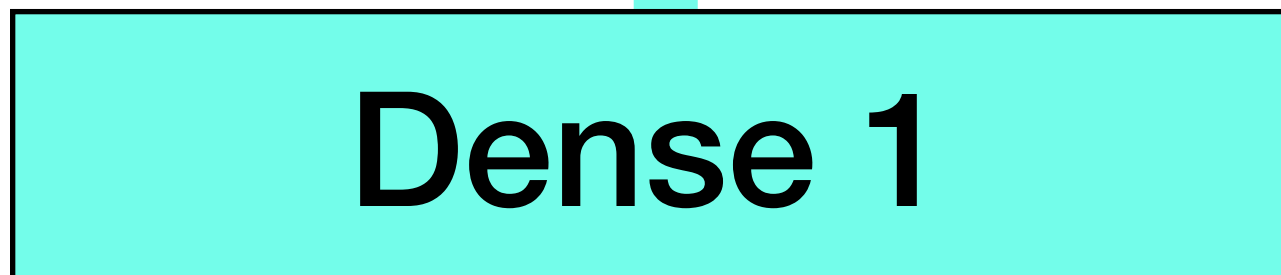


Output

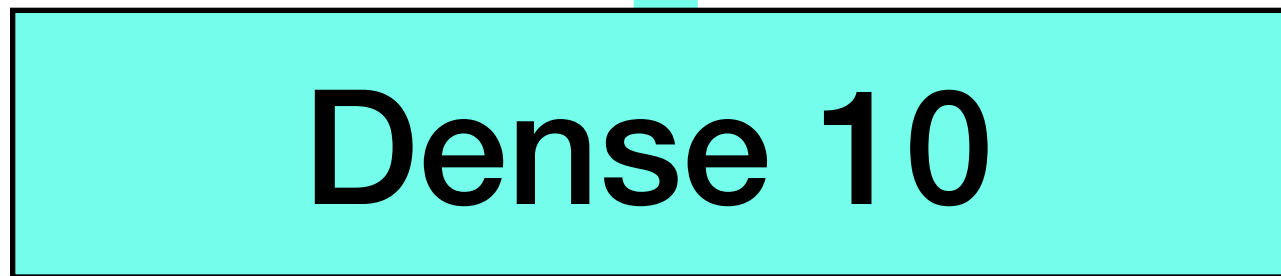
$$\frac{\partial \text{Loss}}{\partial W_1}$$



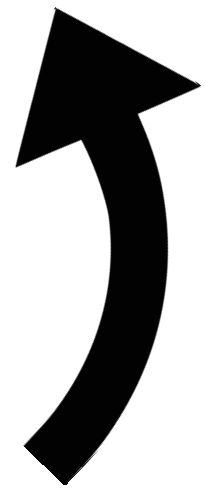
Input



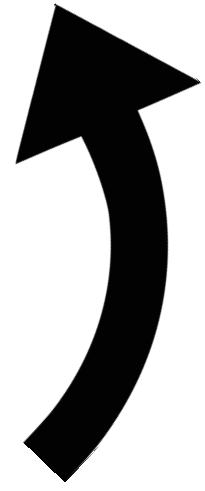
⋮

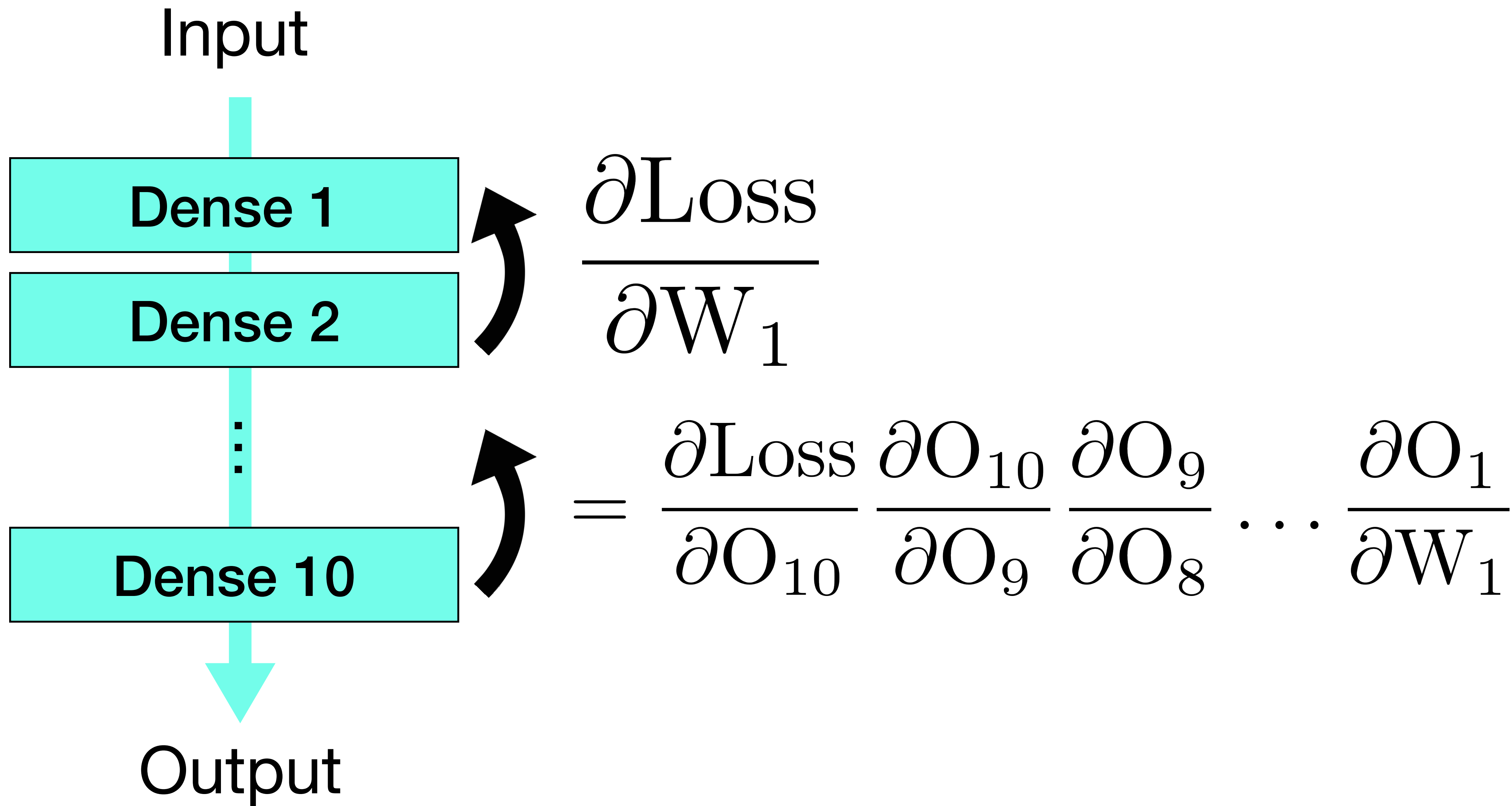


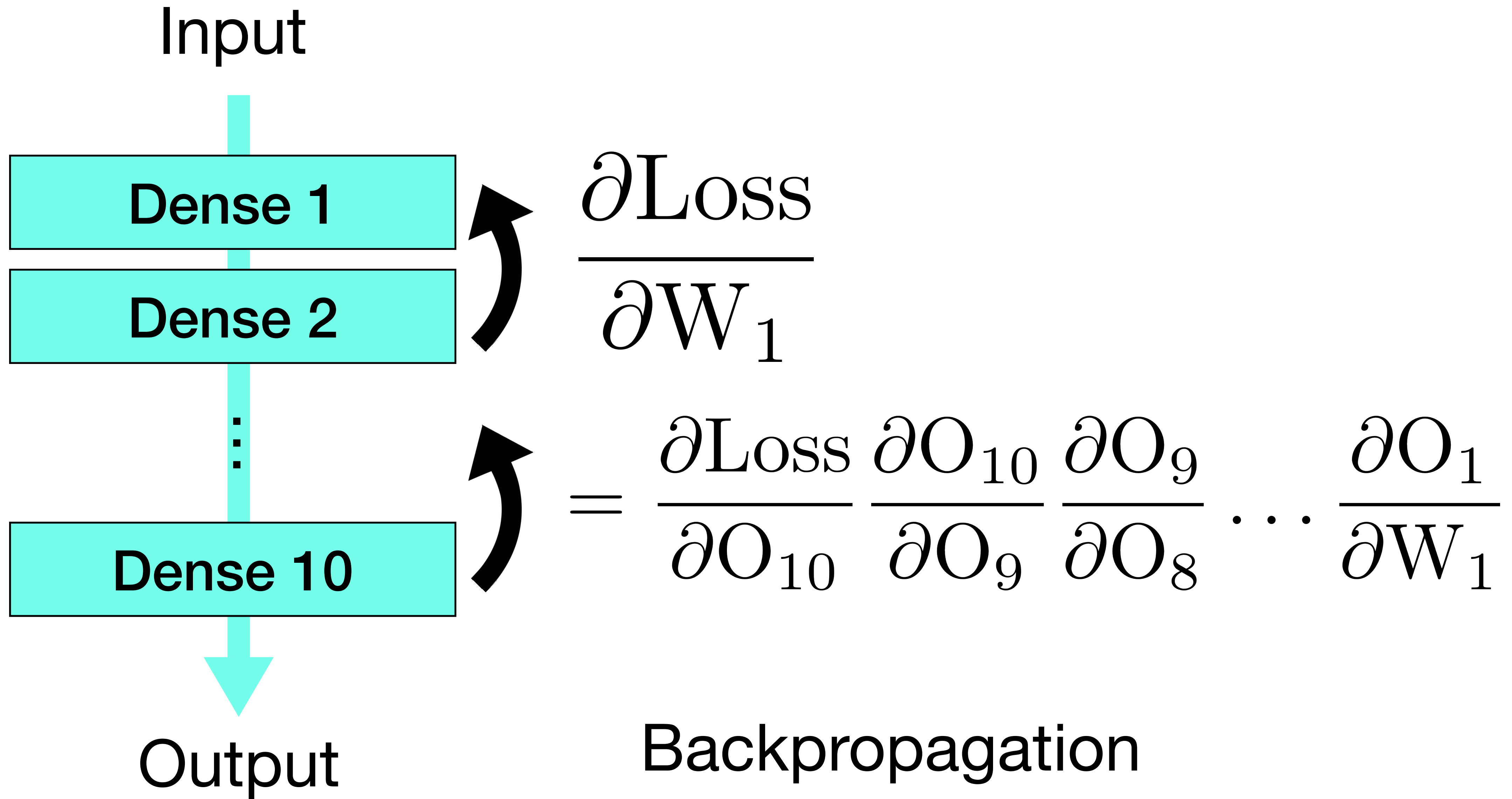
Output

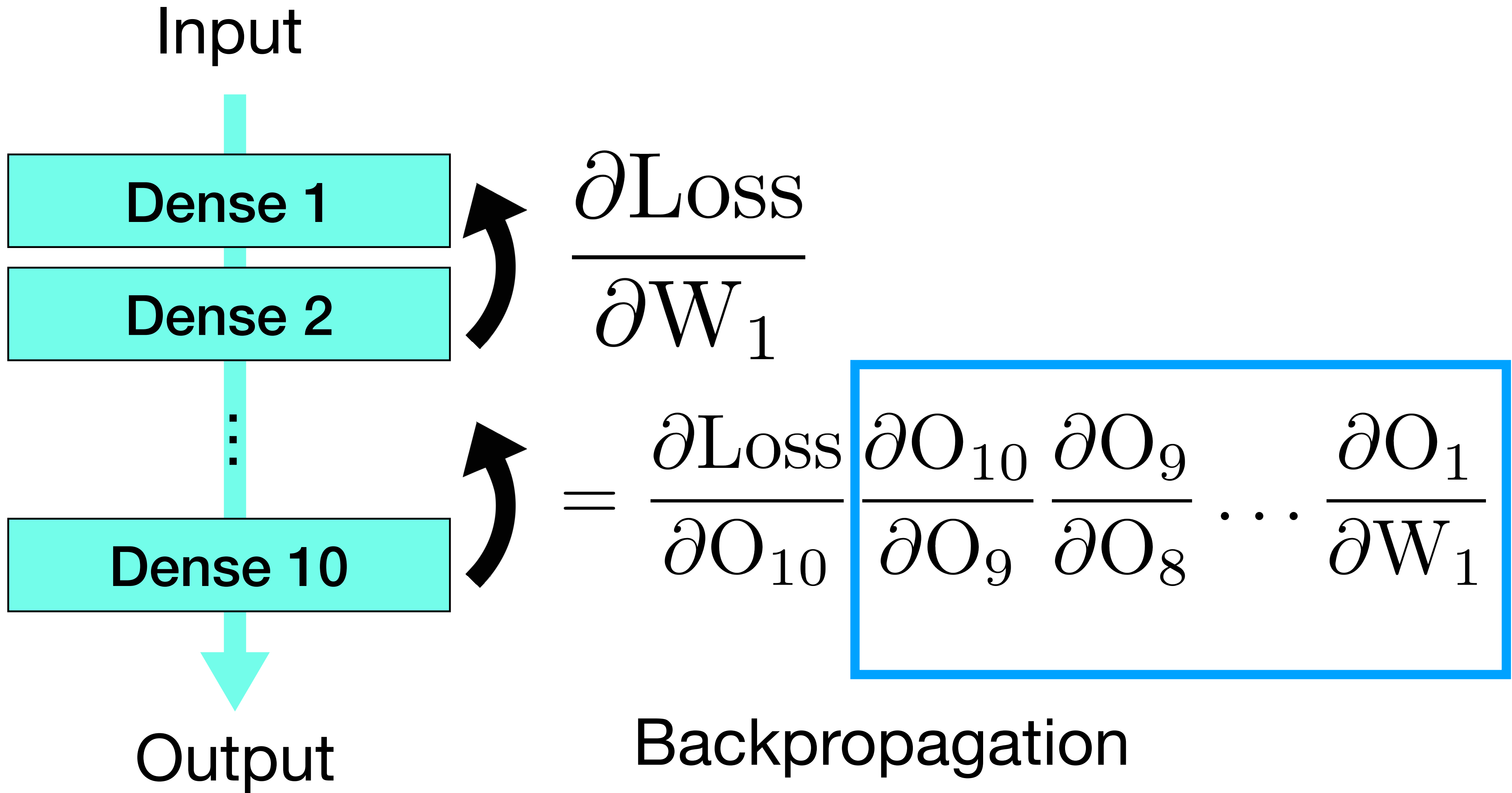


$$\frac{\partial \text{Loss}}{\partial W_1}$$

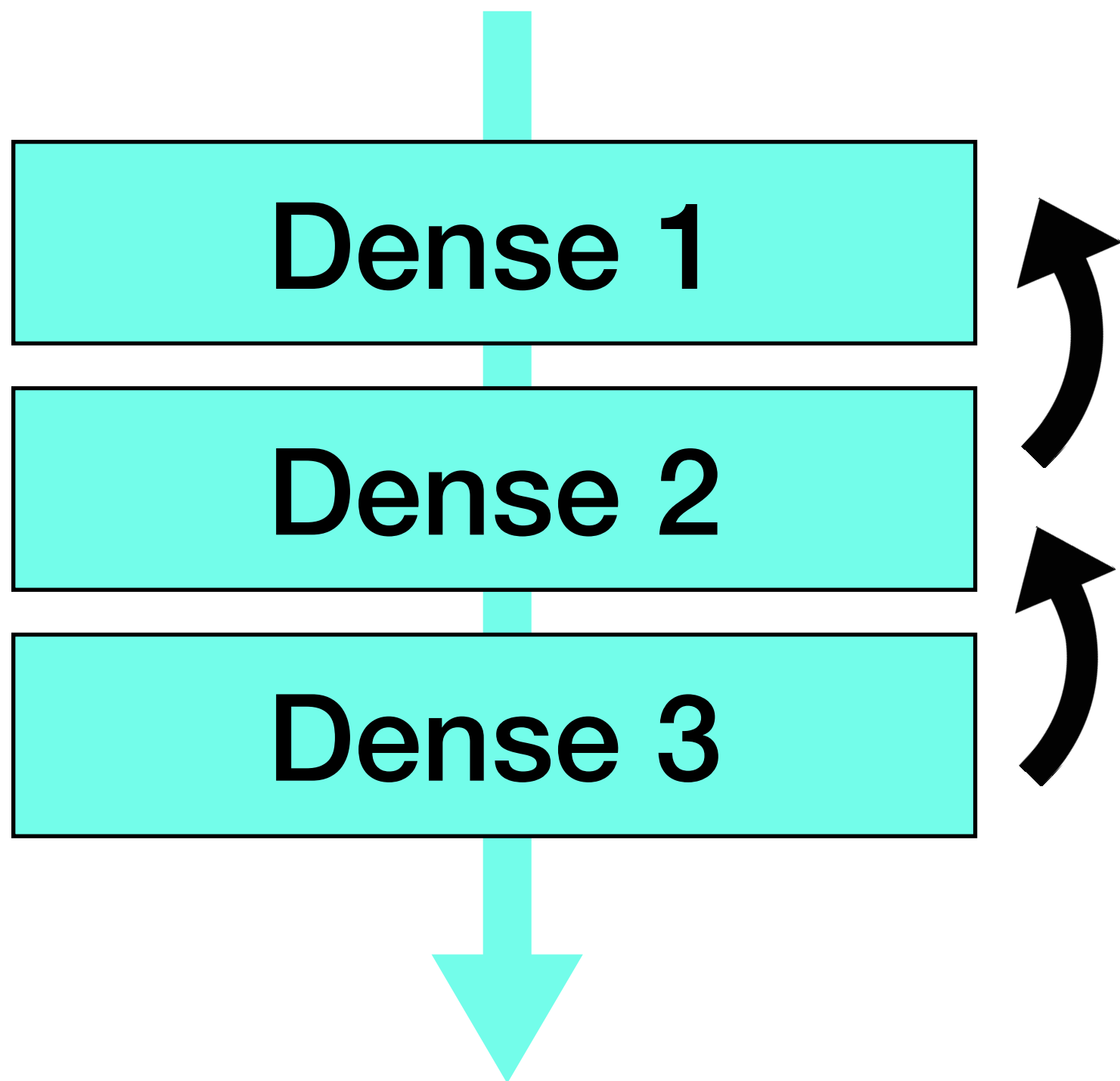








Input



Output

$$\frac{\partial O_3}{\partial O_2} \times \frac{\partial O_2}{\partial O_1} \times \frac{\partial O_1}{\partial W_1}$$

Input

Dense 1

Dense 2

Dense 3

Output

$$\frac{\partial O_3}{\partial O_2} \times \frac{\partial O_2}{\partial O_1} \times \frac{\partial O_1}{\partial W_1}$$

$$\frac{\partial O_3}{\partial O_2}$$

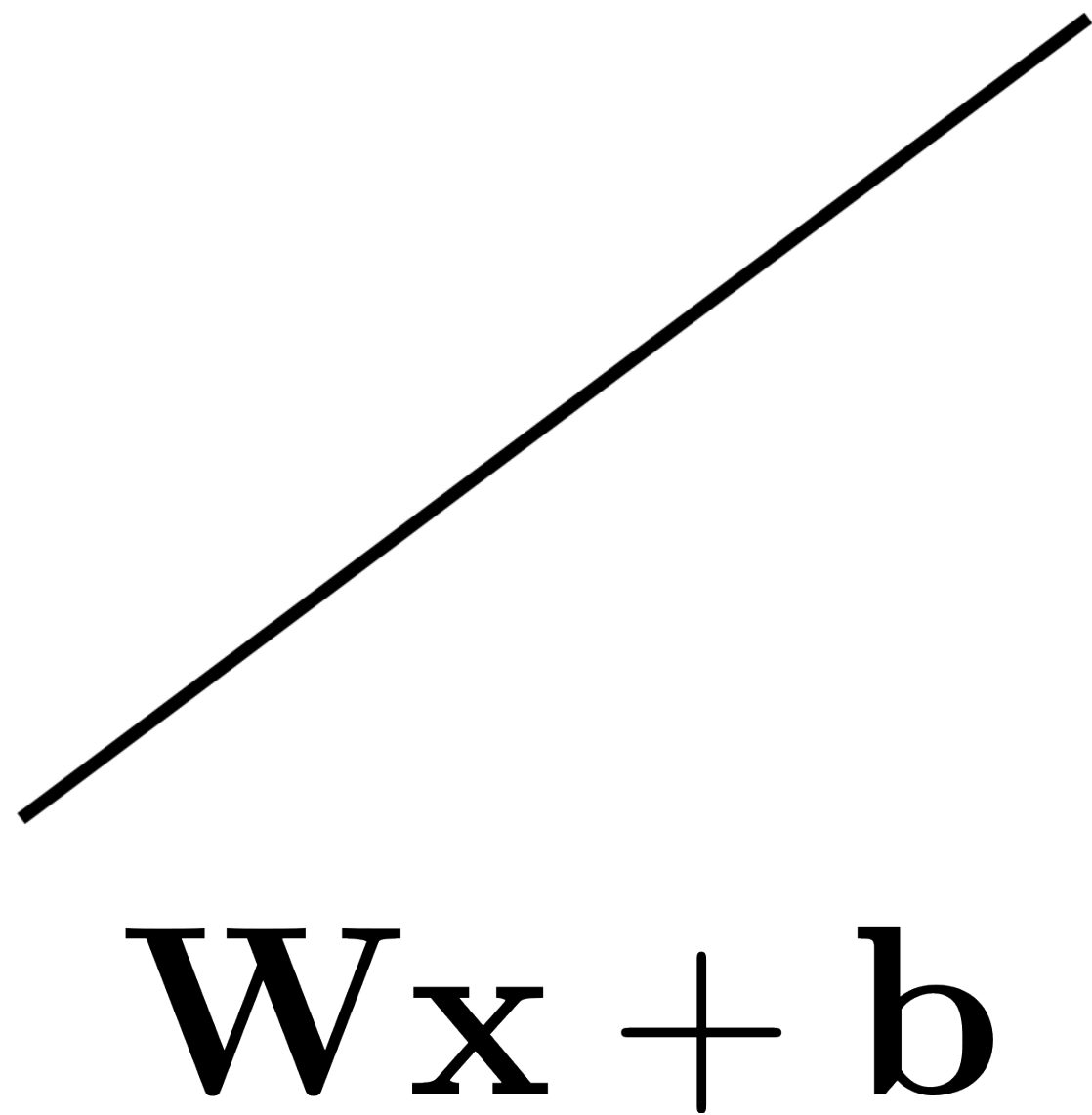
$$\frac{\partial O_3}{\partial O_2}$$

A single layer

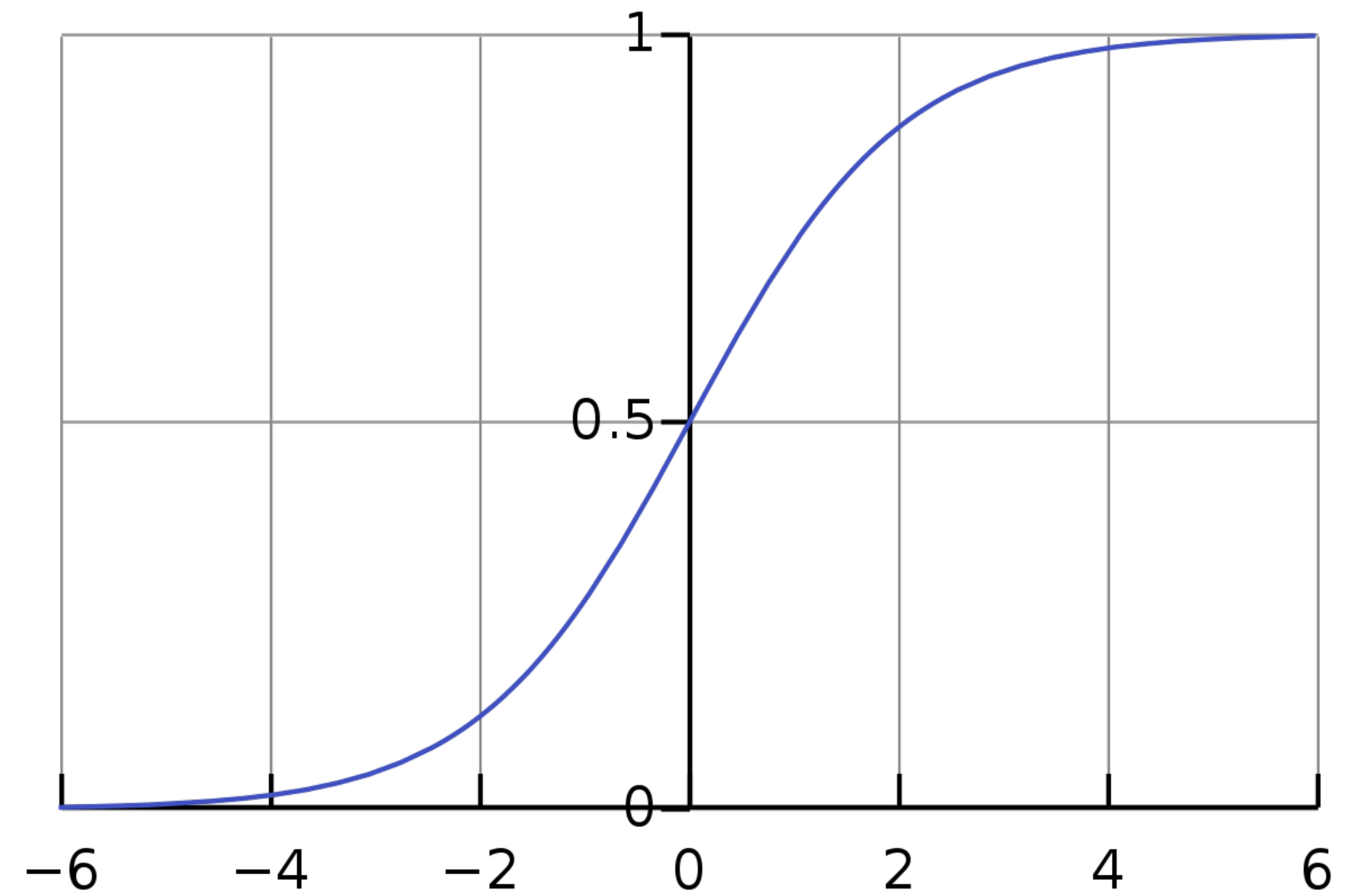
$$\frac{\partial O_3}{\partial O_2}$$

A single layer

Linear layer



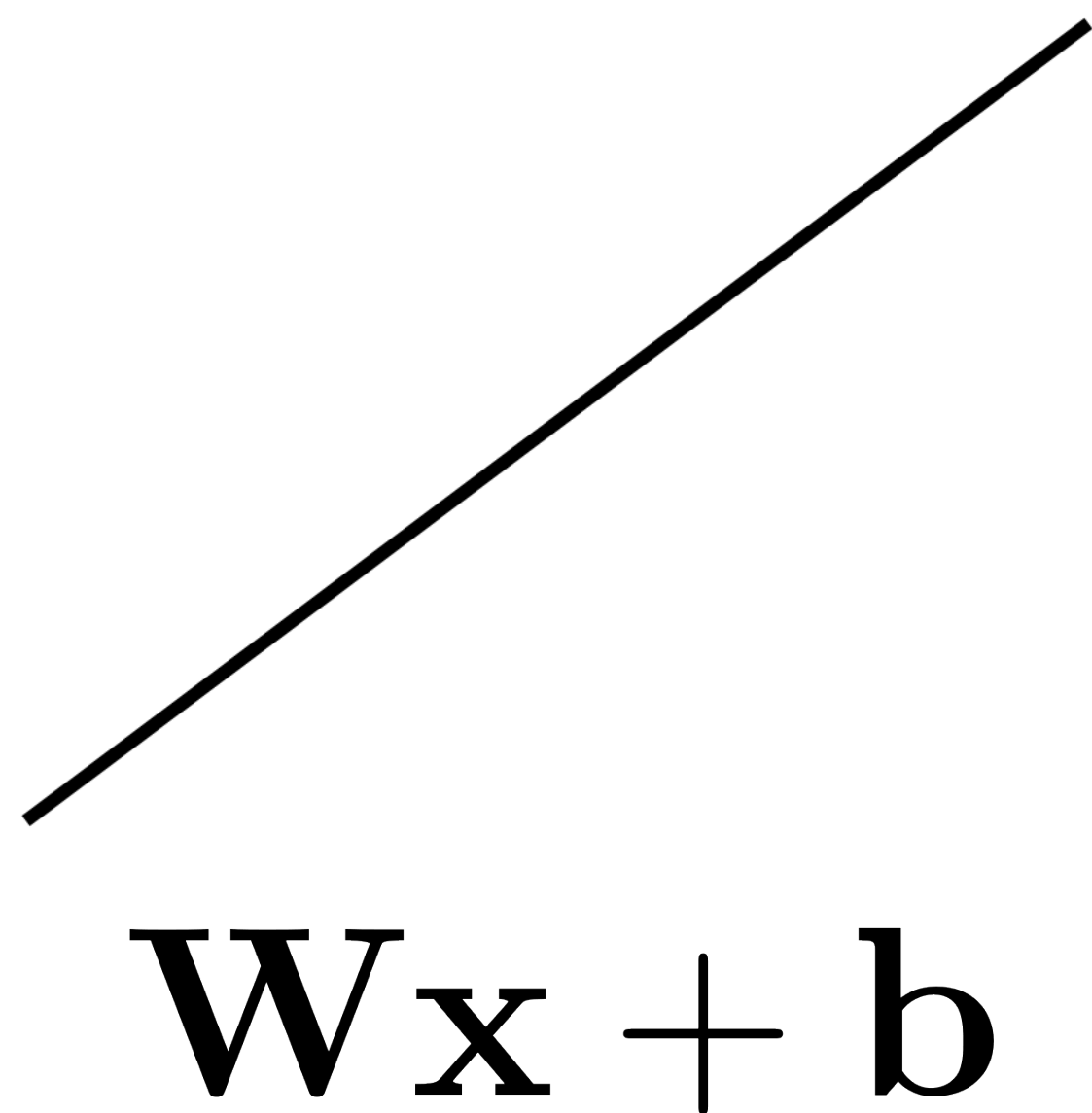
Non-linearity



$$\frac{\partial O_3}{\partial O_2}$$

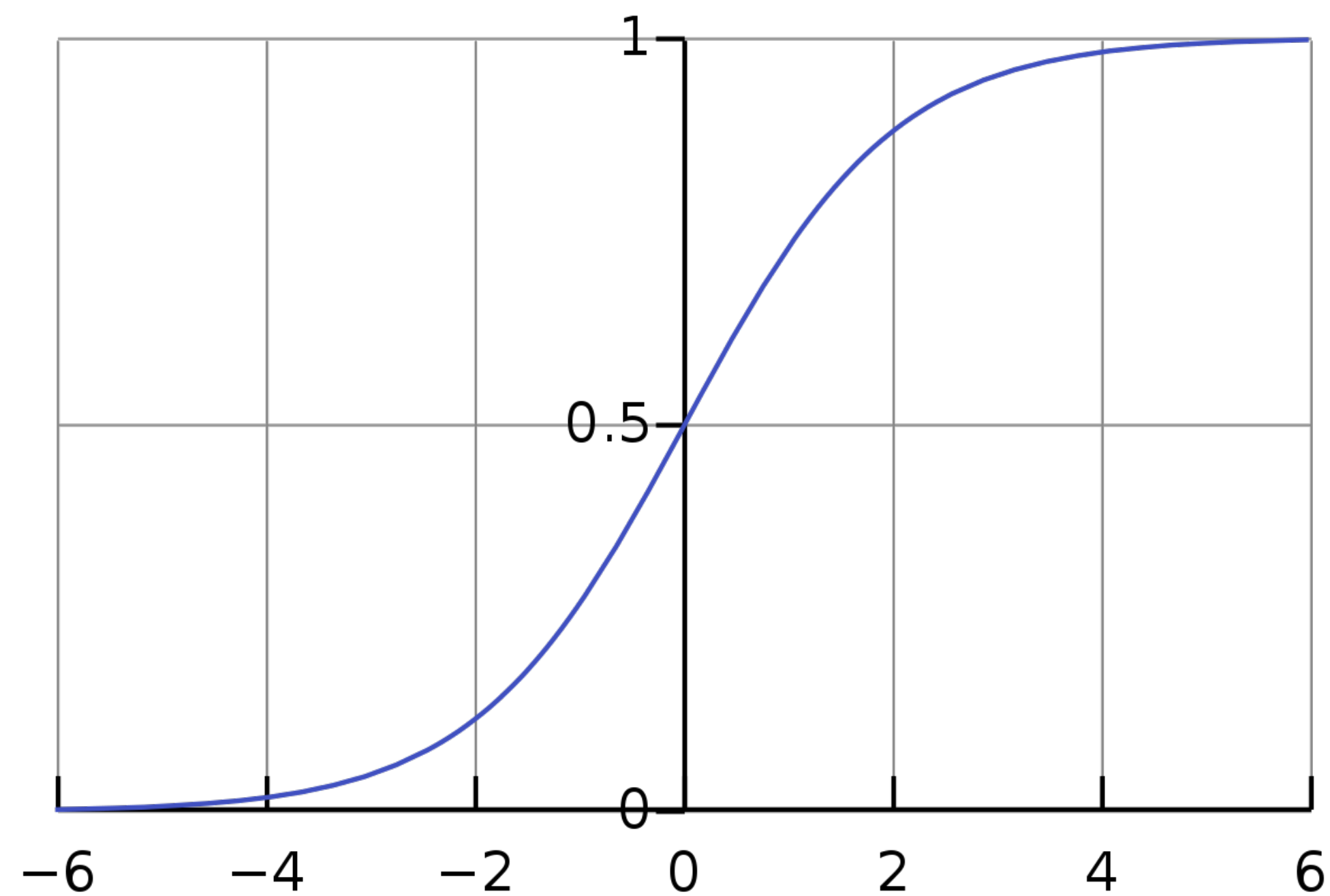
A single layer

Linear layer



Non-linearity

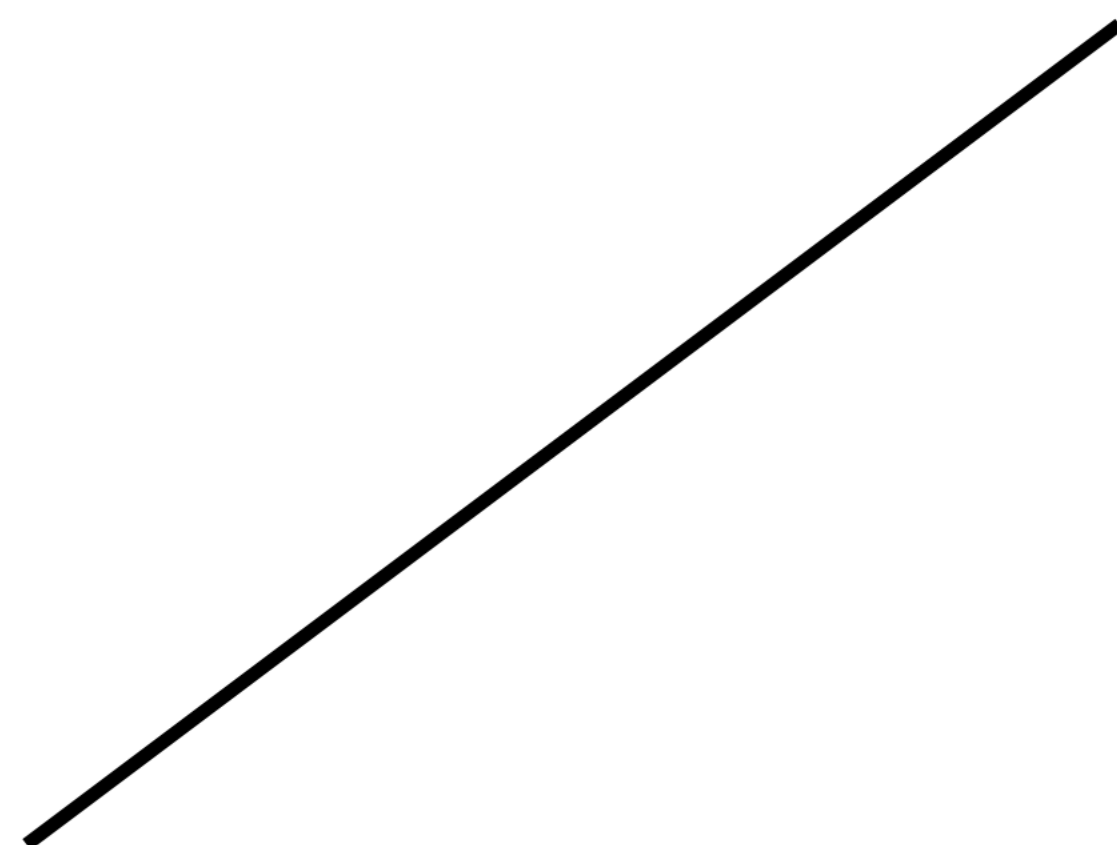
x



$$\frac{\partial O_3}{\partial O_2}$$

A single layer

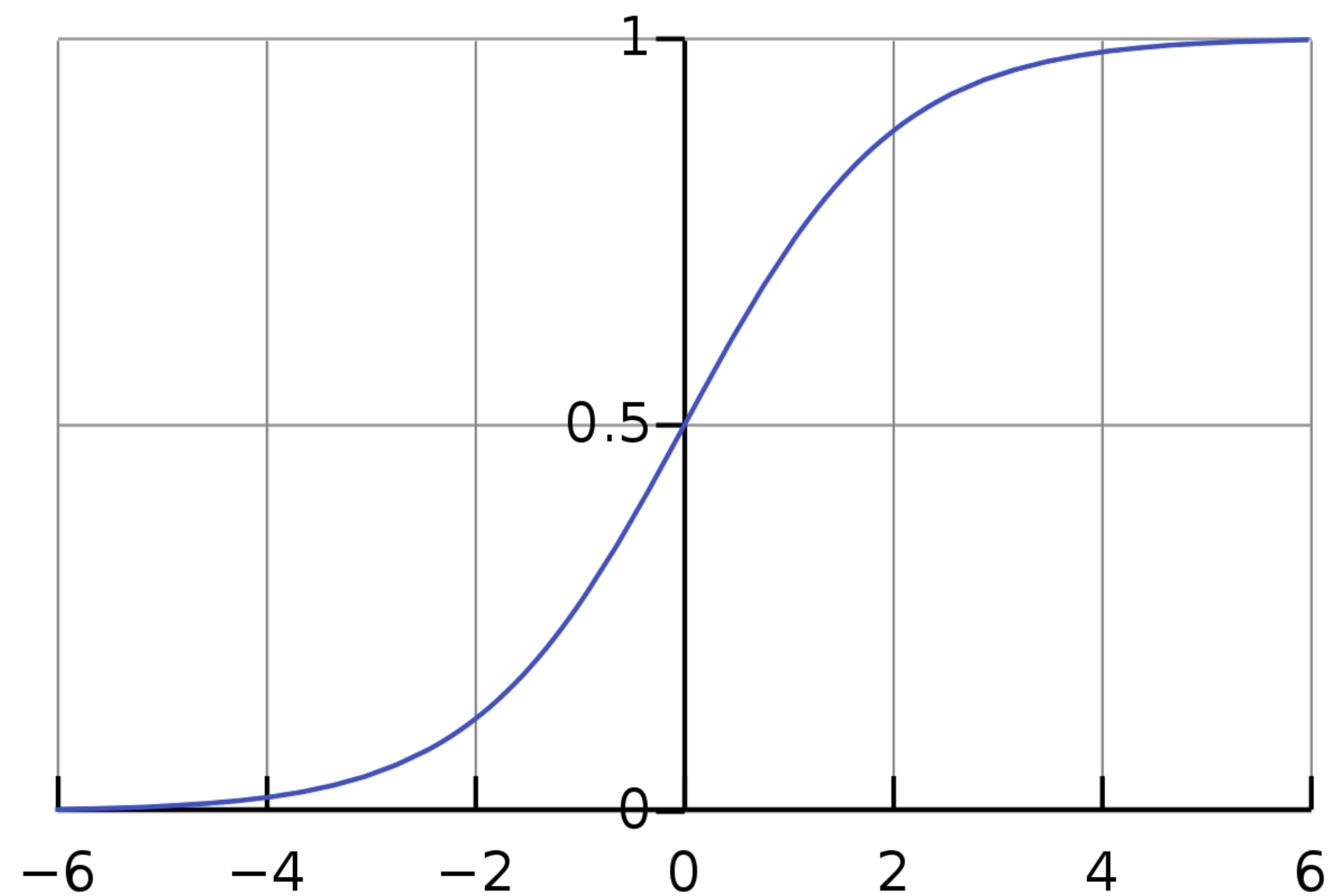
Linear layer

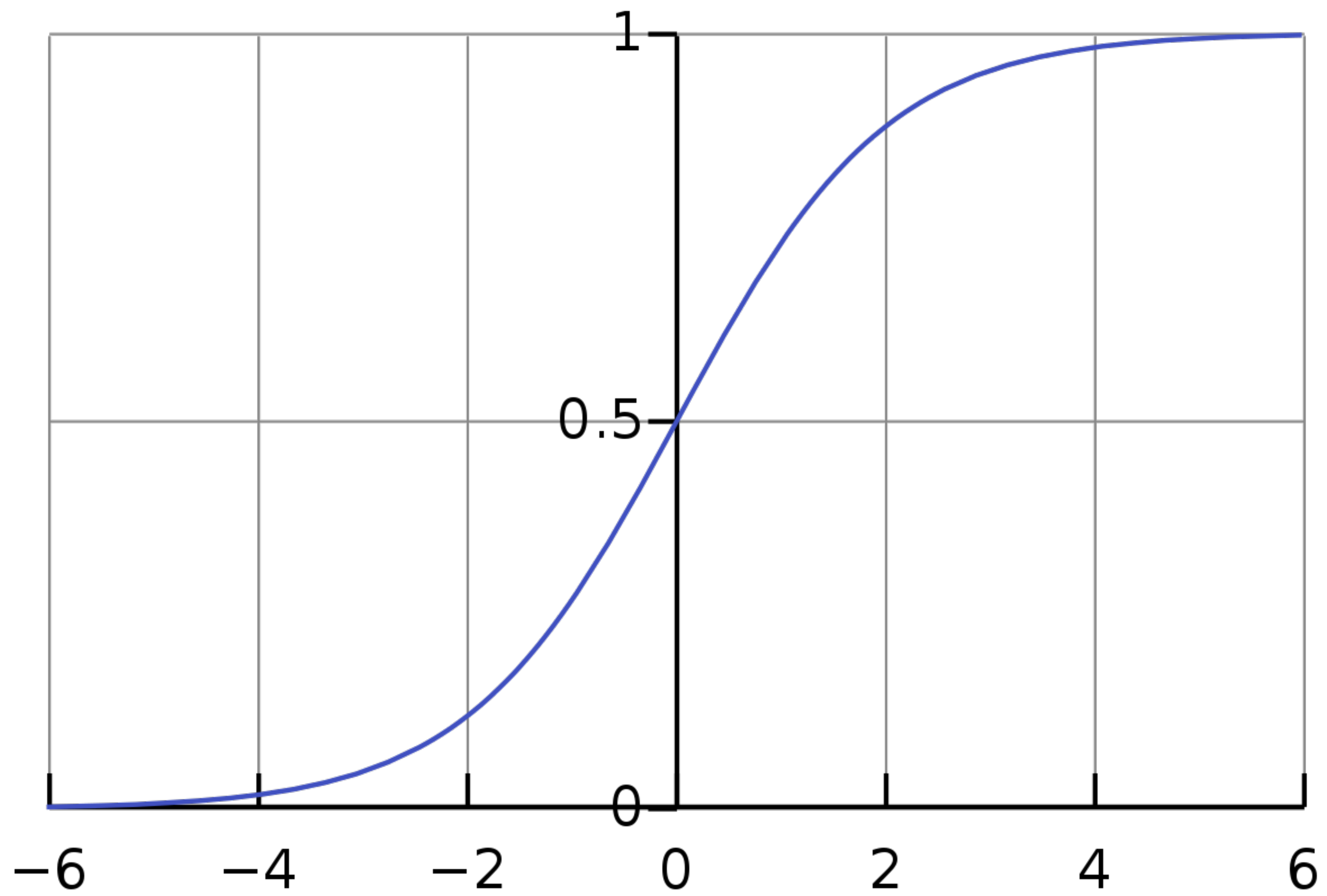


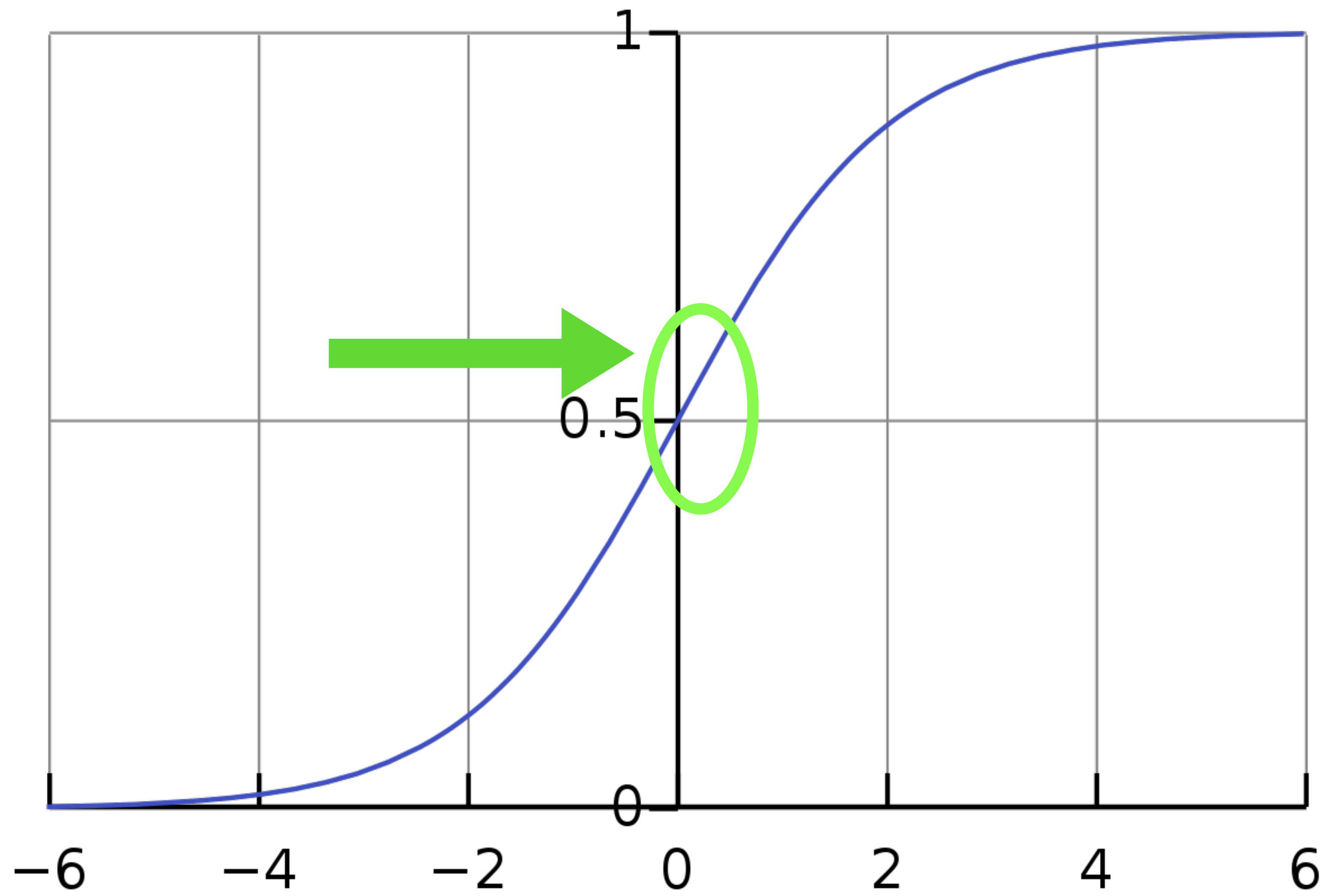
$$\mathbf{W}\mathbf{x} + \mathbf{b}$$

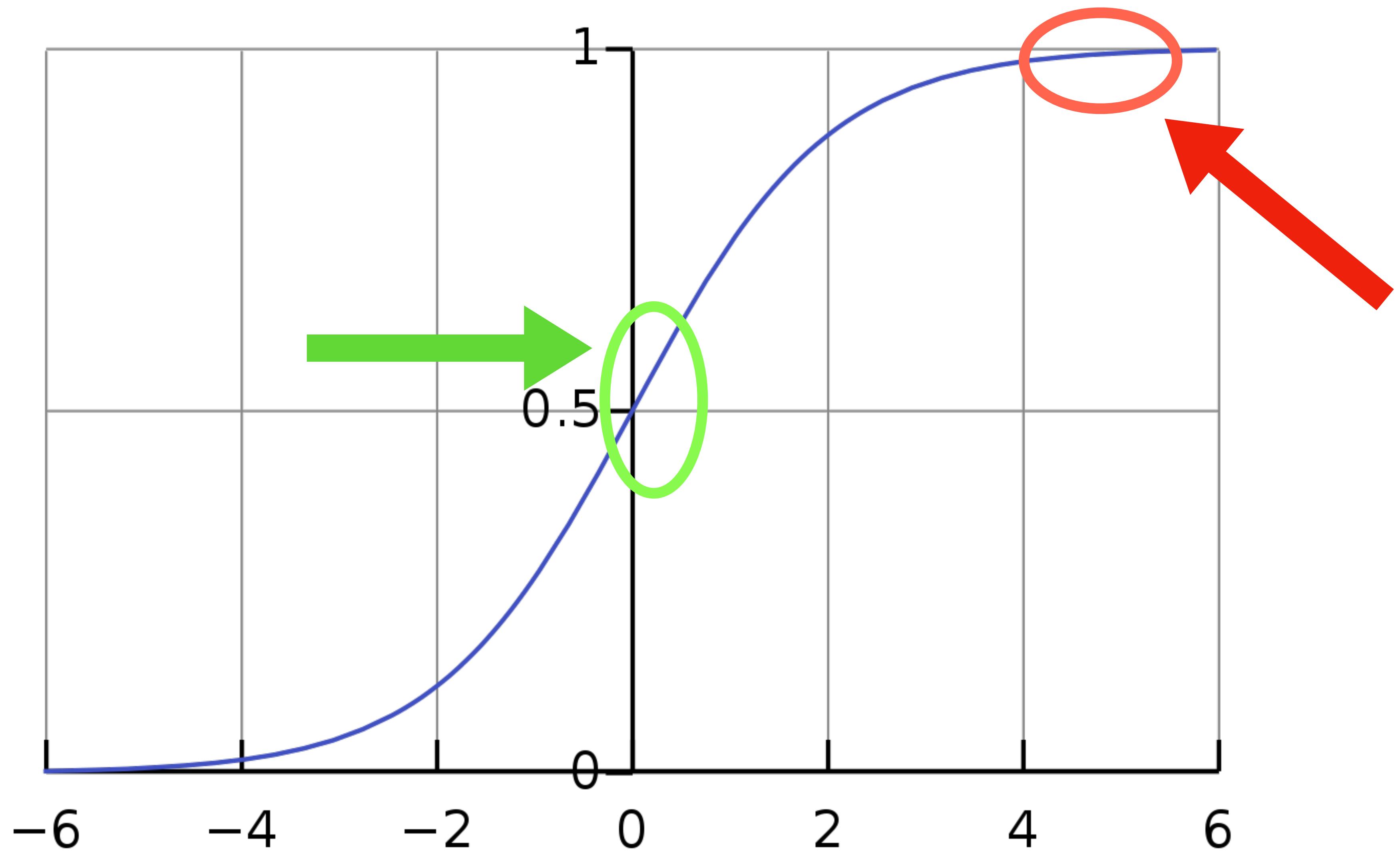
x

Non-linearity









**Initialise weights with
SMALL variance!**

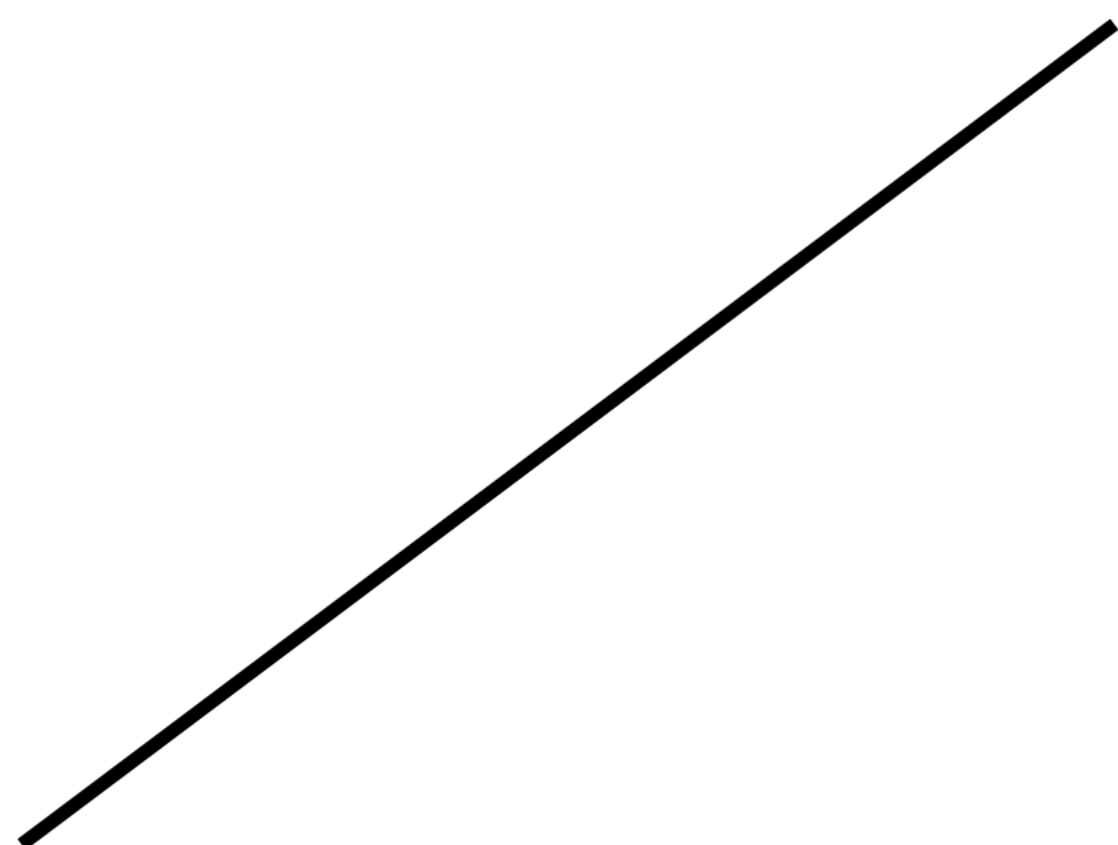
**Initialise weights with
SMALL variance!**

e.g. Glorot initialisation

$$\frac{\partial O_3}{\partial O_2}$$

A single layer

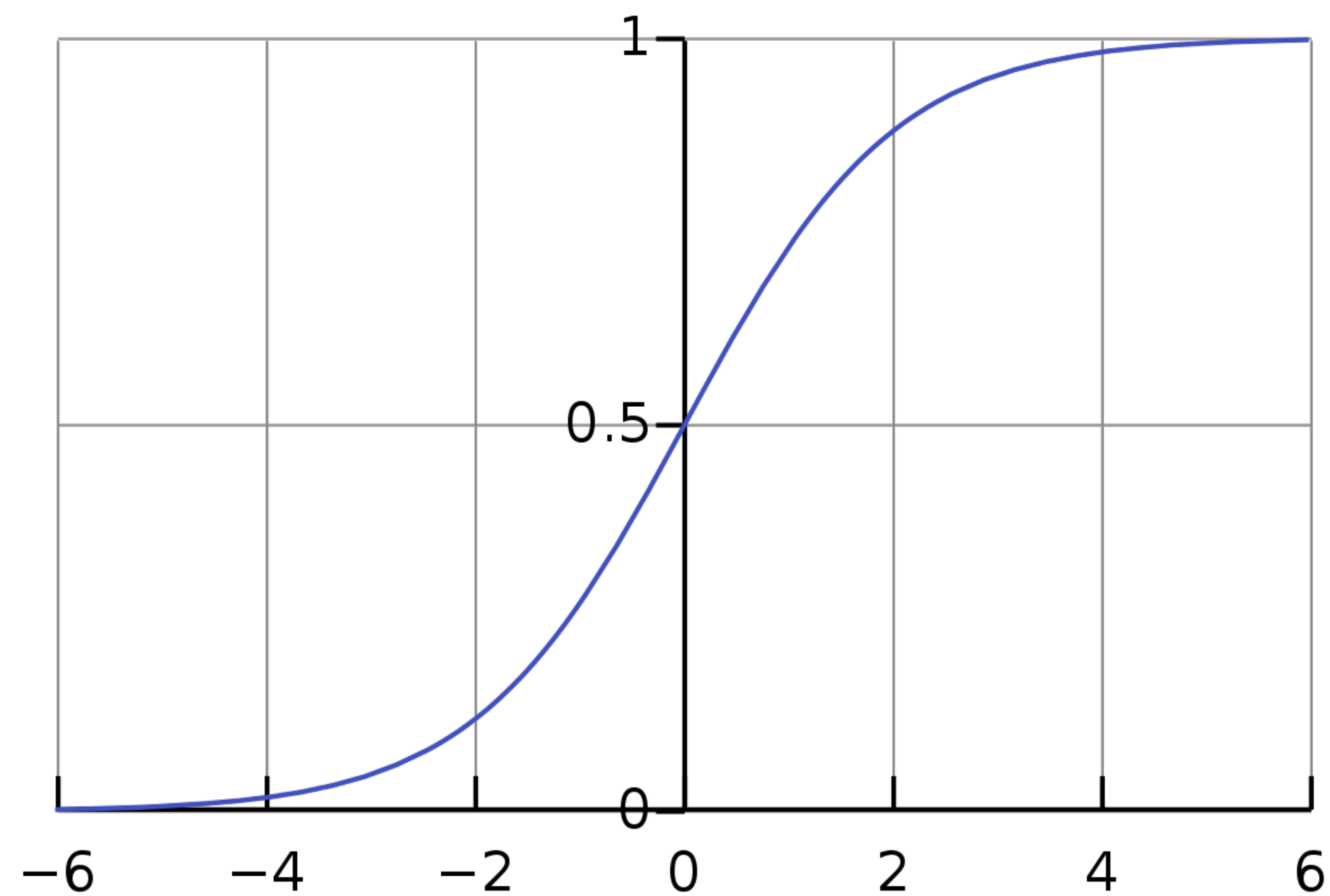
Linear layer



$$\mathbf{W}\mathbf{x} + \mathbf{b}$$

Non-linearity

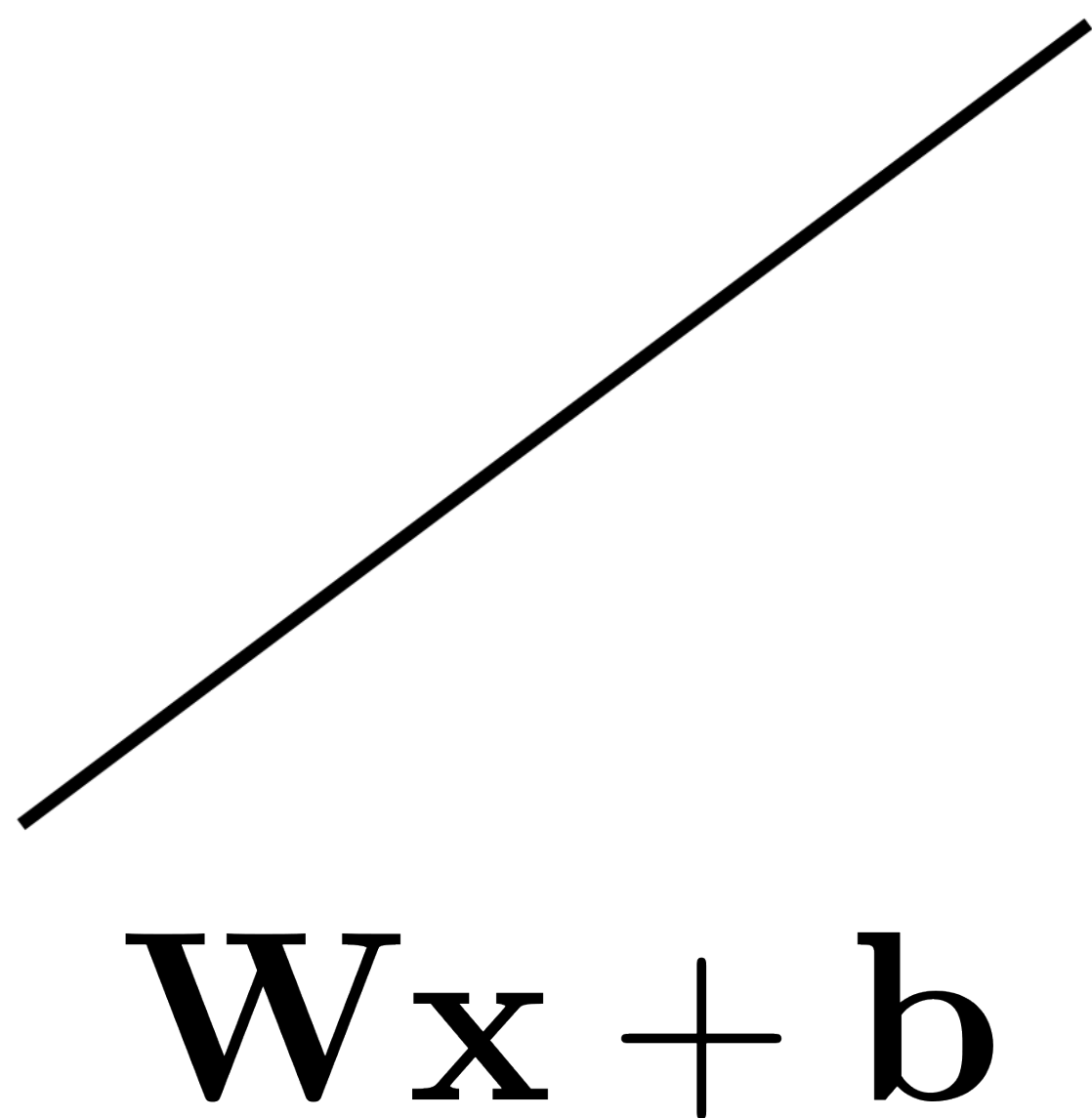
x



$$\frac{\partial O_3}{\partial O_2}$$

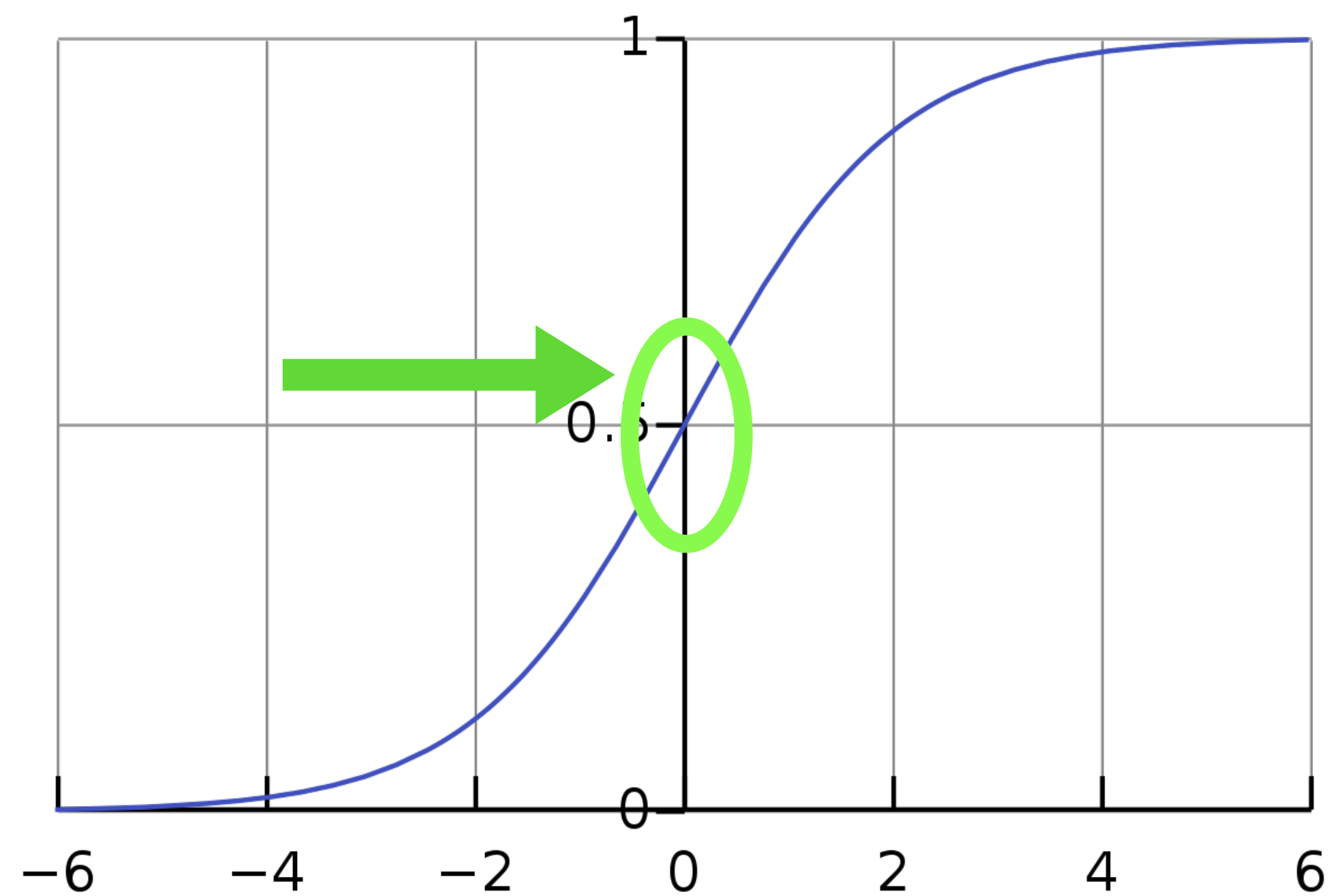
A single layer

Linear layer



Non-linearity

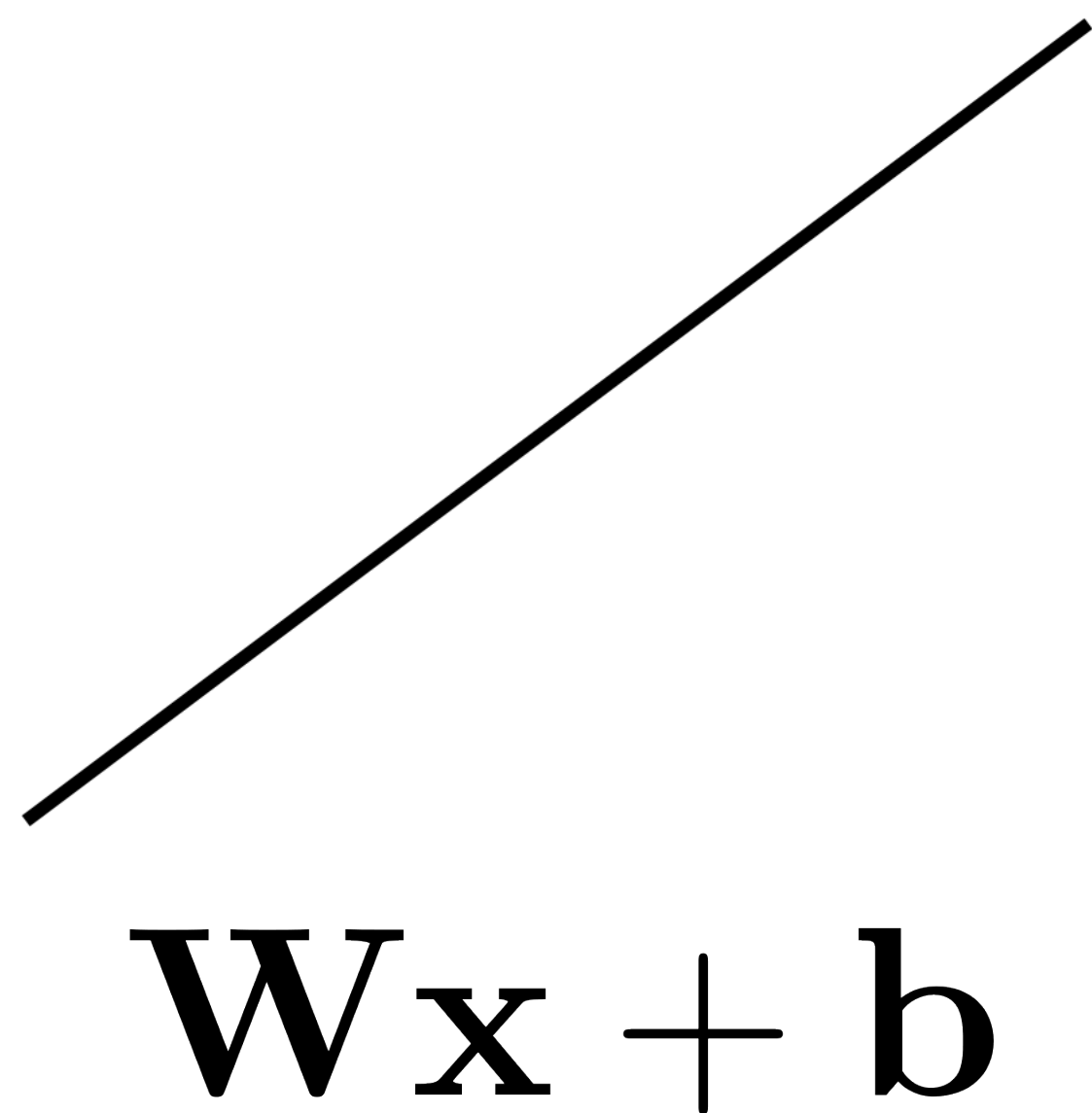
x



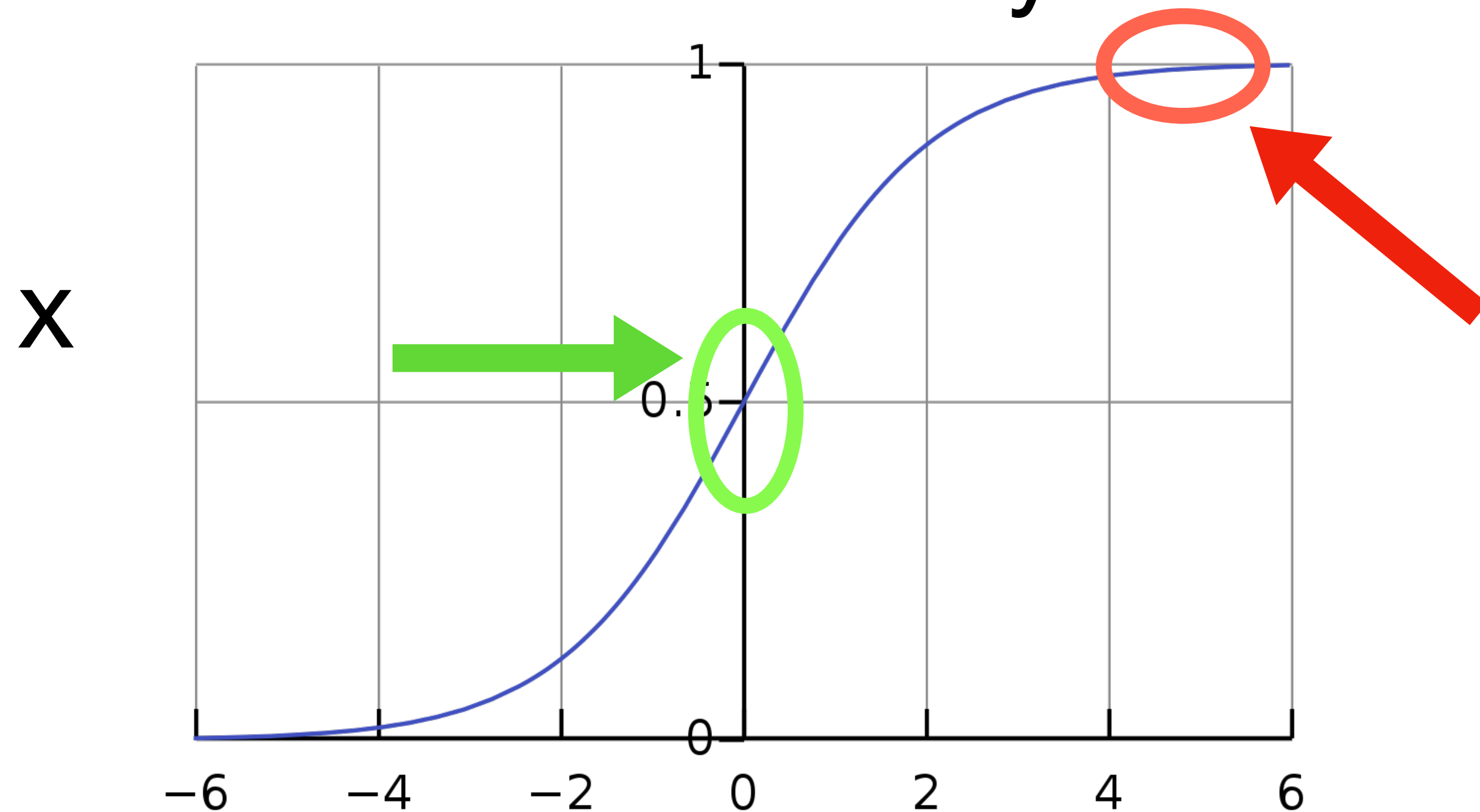
$$\frac{\partial O_3}{\partial O_2}$$

A single layer

Linear layer



Non-linearity



**Initialise weights with
SMALL variance!**

**Initialise weights with
SMALL variance!**

e.g. Glorot initialisation

III. Deeper is better

ImageNet



ImageNet



leopard

ImageNet



sabre-toothed tiger

leopard

cat

snow leopard

ImageNet



tree shrew

sabre-toothed tiger

leopard

cat

snow leopard

ImageNet



tree shrew

sabre-toothed tiger

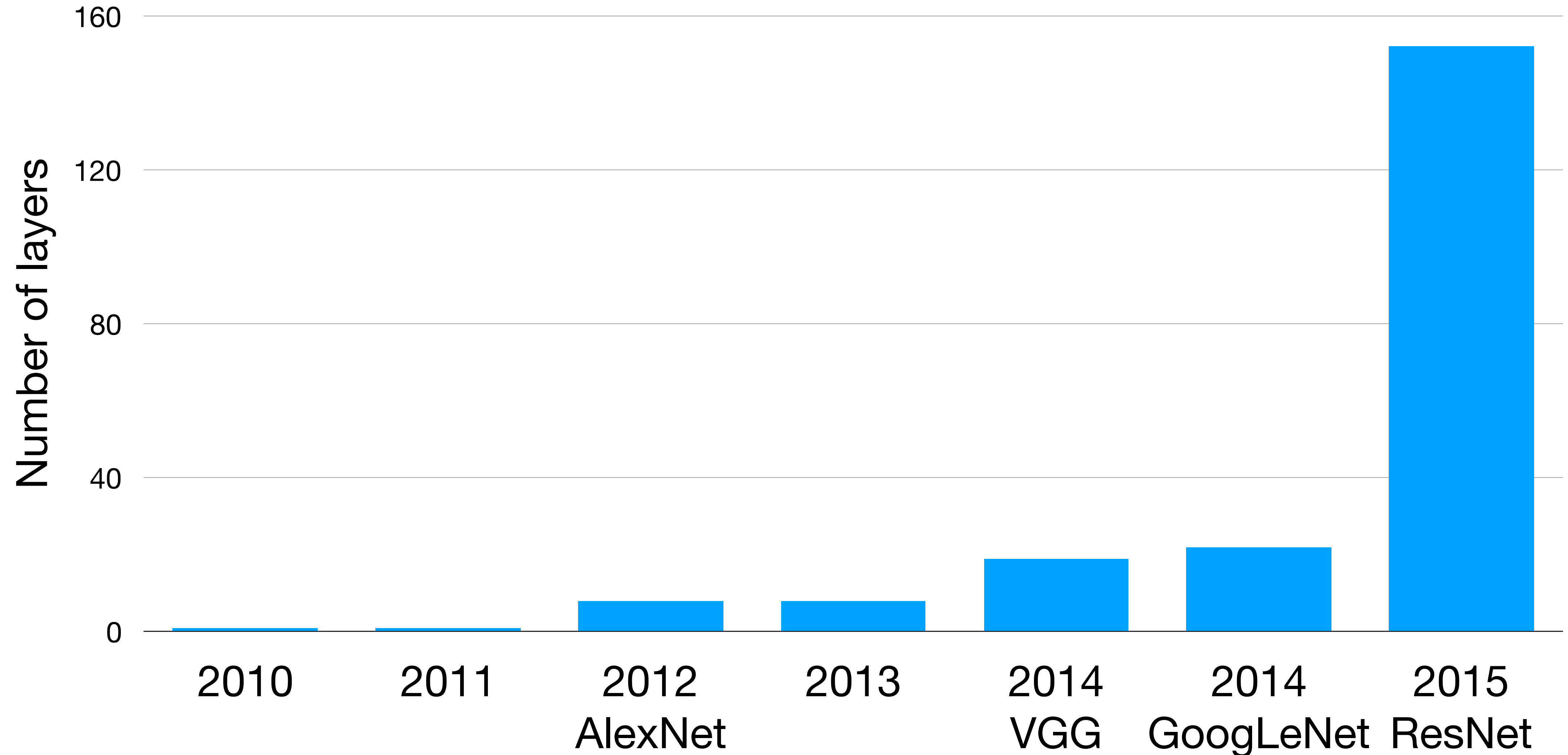
leopard

cat

snow leopard

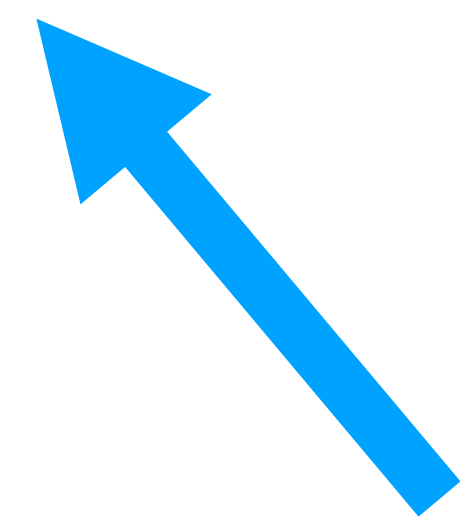
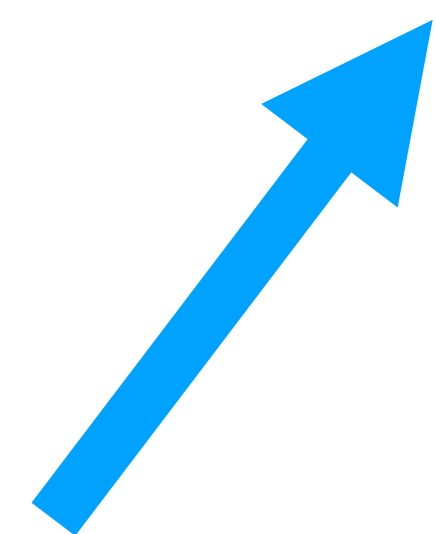
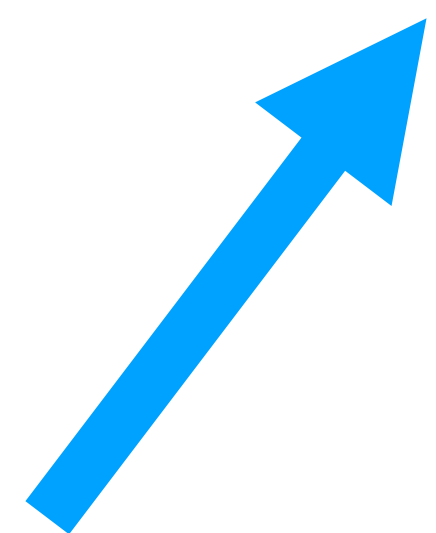
radiator

Revolution of Depth: Top models on ImageNet



Data from Kaiming He, Xiangyu Zhang, Shaoqing Ren and Jian Sun. 'Deep Residual Learning for Image Recognition', CVPR 2016.

$$W = W - \alpha \nabla_w f$$



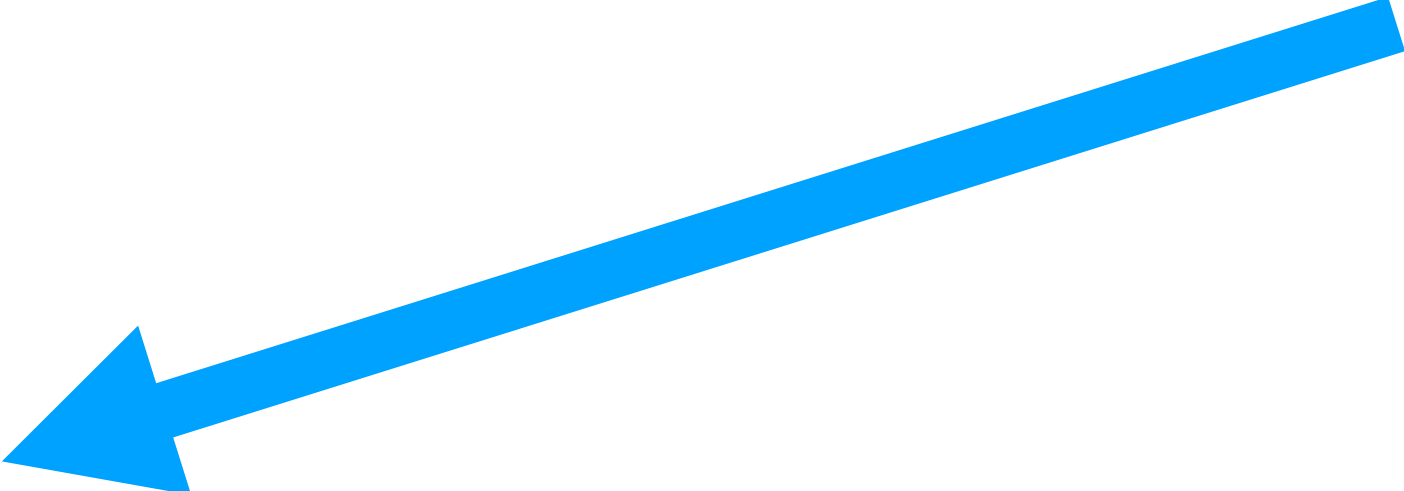
Parameters

Learning rate
(step size)

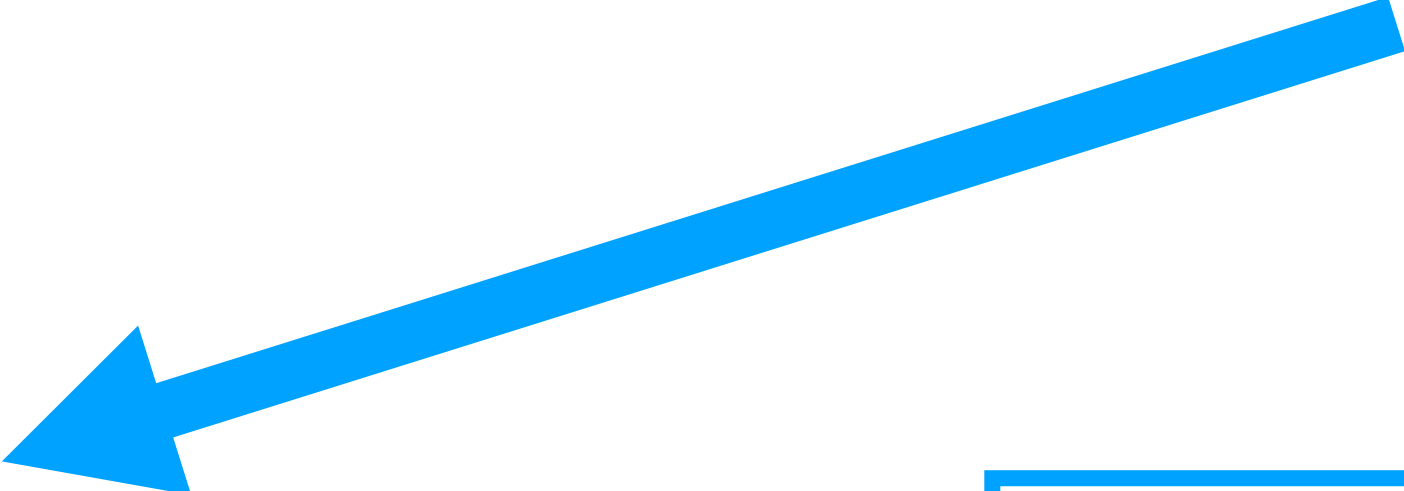
Gradient

$$W = W - \alpha \nabla_w f$$

$$W = W - \alpha \nabla_W f$$

$$\frac{\partial \text{Loss}}{\partial W_1} = \frac{\partial \text{Loss}}{\partial O_{10}} \frac{\partial O_{10}}{\partial O_9} \frac{\partial O_9}{\partial O_8} \cdots \frac{\partial O_1}{\partial W_1}$$


$$W = W - \alpha \nabla_w f$$


$$\frac{\partial \text{Loss}}{\partial W_1} = \frac{\partial \text{Loss}}{\partial O_{10}} \frac{\partial O_{10}}{\partial O_9} \frac{\partial O_9}{\partial O_8} \cdots \frac{\partial O_1}{\partial W_1}$$

$$\frac{\partial O_{100}}{\partial O_{99}} \times \frac{\partial O_{99}}{\partial O_{98}} \times \dots \times \frac{\partial O_2}{\partial O_1} \times \frac{\partial O_1}{\partial W_1}$$

$$\boxed{0.1} \times \frac{\partial O_{99}}{\partial O_{98}} \times \dots \times \frac{\partial O_2}{\partial O_1} \times \frac{\partial O_1}{\partial W_1}$$

$$\boxed{0.1} \times \boxed{0.1} \times \dots \times \frac{\partial O_2}{\partial O_1} \times \frac{\partial O_1}{\partial W_1}$$

$$\boxed{0.1} \times \boxed{0.1} \times \dots \times \boxed{0.1} \times \frac{\partial \mathcal{O}_1}{\partial W_1}$$

$$0.1 \times 0.1 \times \dots \times 0.1 \times 0.1$$

$$0.1 \times 0.1 \times \dots \times 0.1 \times 0.1$$

$$= 10^{-100}$$

Vanishing gradients

$$0.1 \times 0.1 \times \dots \times 0.1 \times 0.1$$

$$= 10^{-100}$$

$$\frac{\partial O_{100}}{\partial O_{99}} \times \frac{\partial O_{99}}{\partial O_{98}} \times \dots \times \frac{\partial O_2}{\partial O_1} \times \frac{\partial O_1}{\partial W_1}$$

$$\boxed{2} \times \frac{\partial O_{99}}{\partial O_{98}} \times \dots \times \frac{\partial O_2}{\partial O_1} \times \frac{\partial O_1}{\partial W_1}$$

$$\frac{\partial O_{100}}{\partial O_{99}} \times \frac{\partial O_{99}}{\partial O_{98}} \times \dots \times \frac{\partial O_2}{\partial O_1} \times \frac{\partial O_1}{\partial W_1}$$

$$\frac{\partial \mathcal{O}_{100}}{\partial \mathcal{O}_{99}} \times \frac{\partial \mathcal{O}_{99}}{\partial \mathcal{O}_{98}} \times \dots \times \frac{\partial \mathcal{O}_2}{\partial \mathcal{O}_1} \times \frac{\partial \mathcal{O}_1}{\partial W_1}$$

$$\frac{\partial \mathcal{O}_{100}}{\partial \mathcal{O}_{99}} \times \frac{\partial \mathcal{O}_{99}}{\partial \mathcal{O}_{98}} \times \dots \times \frac{\partial \mathcal{O}_2}{\partial \mathcal{O}_1} \times \frac{\partial \mathcal{O}_1}{\partial W_1}$$

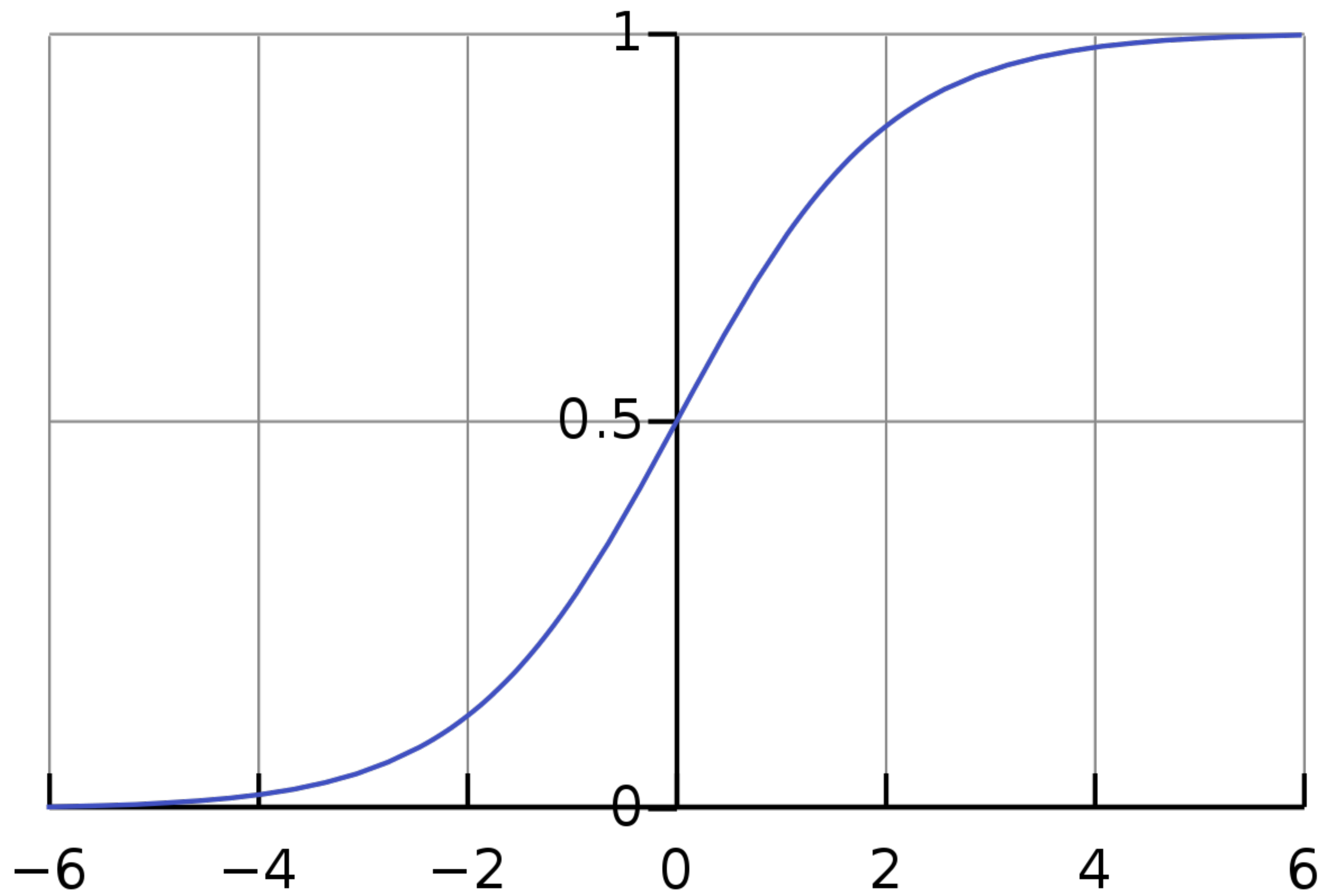
$$\frac{\partial \mathcal{O}_{100}}{\partial \mathcal{O}_{99}} 2 \times \frac{\partial \mathcal{O}_{99}}{\partial \mathcal{O}_{98}} 2 \times \dots \times \frac{\partial \mathcal{O}_2}{\partial \mathcal{O}_1} 2 \times \frac{\partial \mathcal{O}_1}{\partial W_1} 2$$

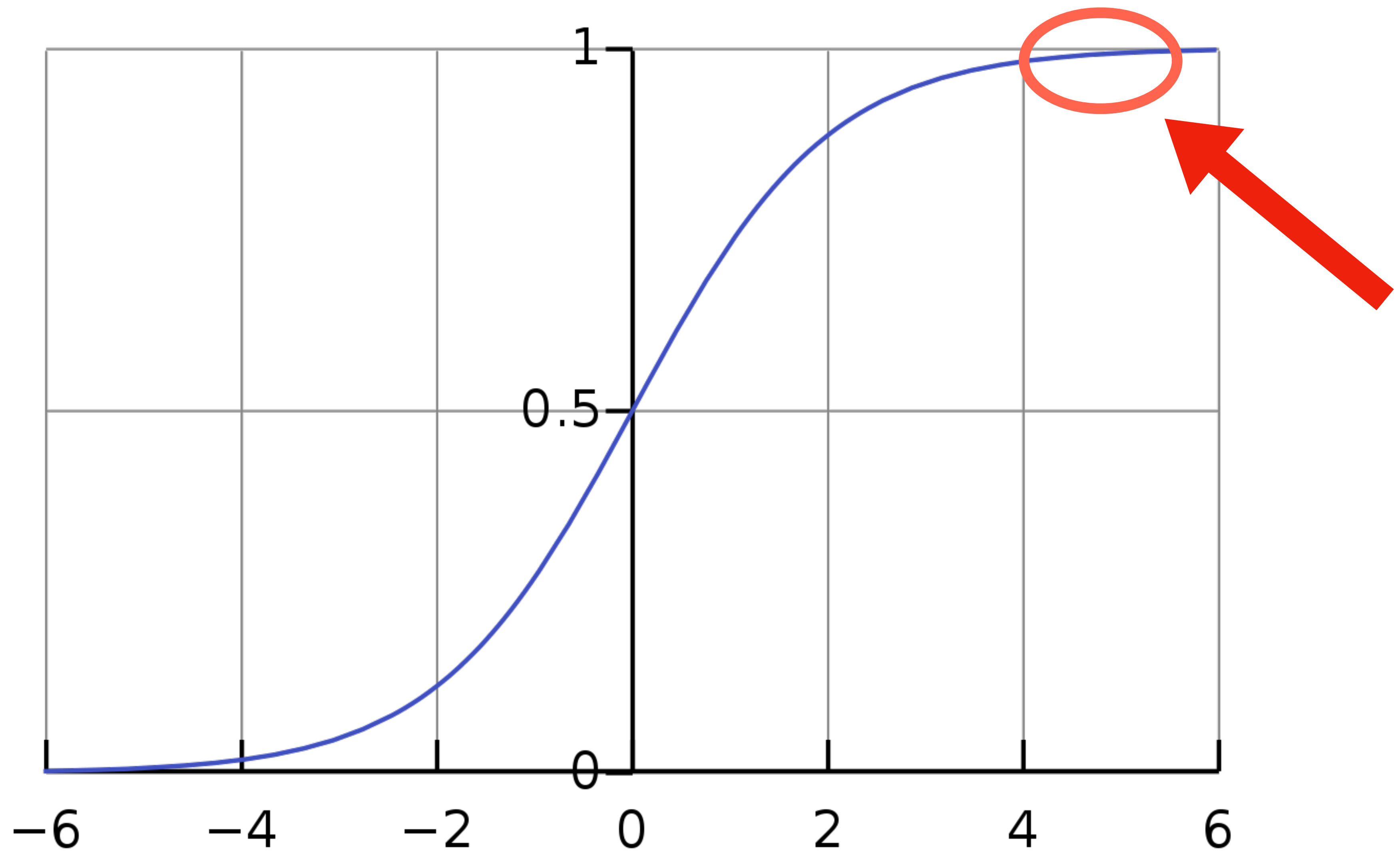
$$> 10^{30}$$

Exploding gradients

$$\frac{\partial \mathcal{O}_{100}}{\partial \mathcal{O}_{99}} \times \frac{\partial \mathcal{O}_{99}}{\partial \mathcal{O}_{98}} \times \dots \times \frac{\partial \mathcal{O}_2}{\partial \mathcal{O}_1} \times \frac{\partial \mathcal{O}_1}{\partial W_1}$$

$$> 10^{30}$$





Vanishing gradients

$$0.1 \times 0.1 \times \dots \times 0.1 \times 0.1$$

$$= 10^{-100}$$

Vanishing gradients

$$0.1 \times \text{[red X]} \times \dots \times 0.1 \times 0.1$$

$$= 10^{-100}$$

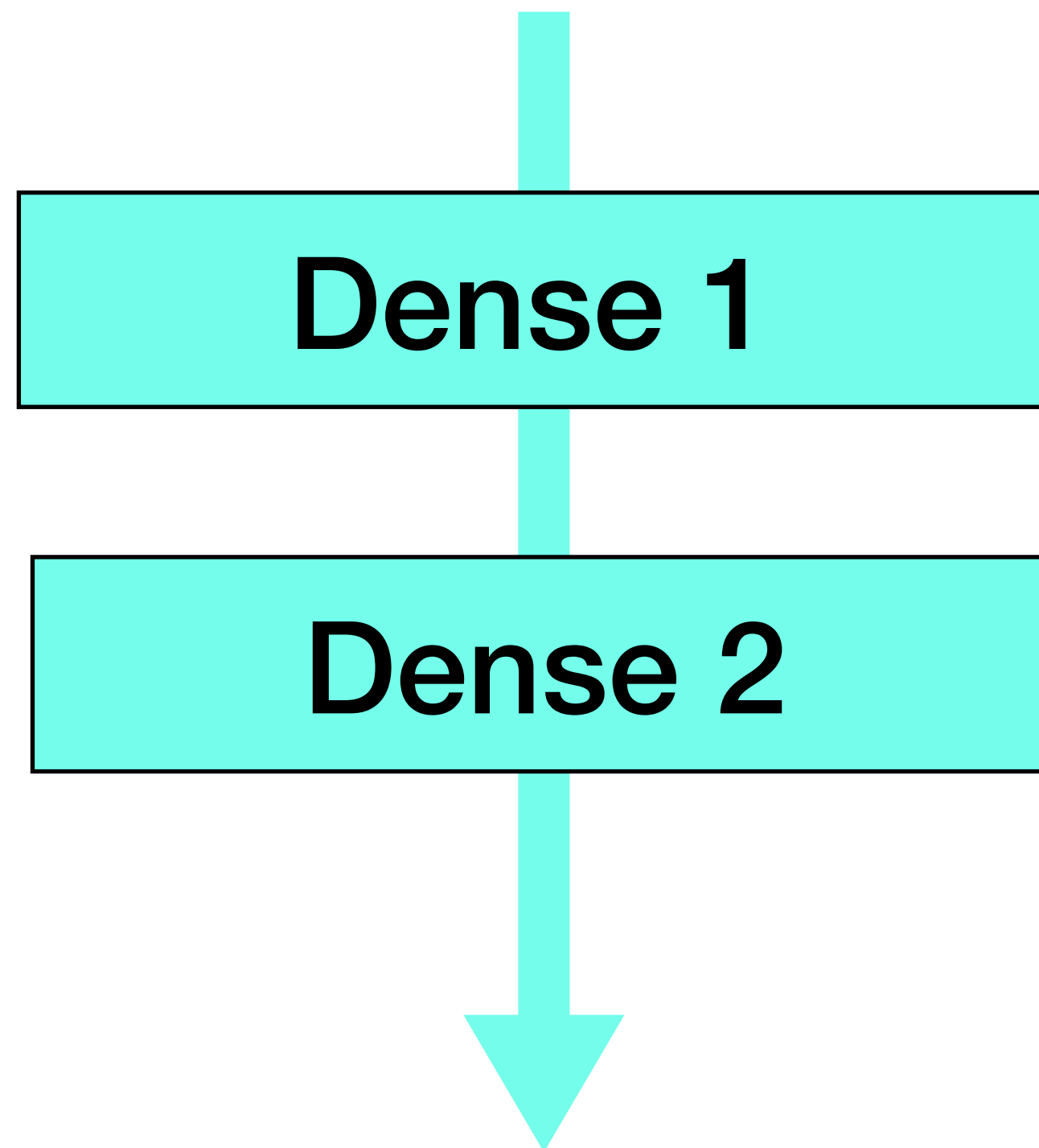
Vanishing gradients

$$\boxed{0.1} \times \boxed{1} \times \dots \times \boxed{0.1} \times \boxed{0.1}$$

$$= 10^{-100}$$

Skip connections

Input



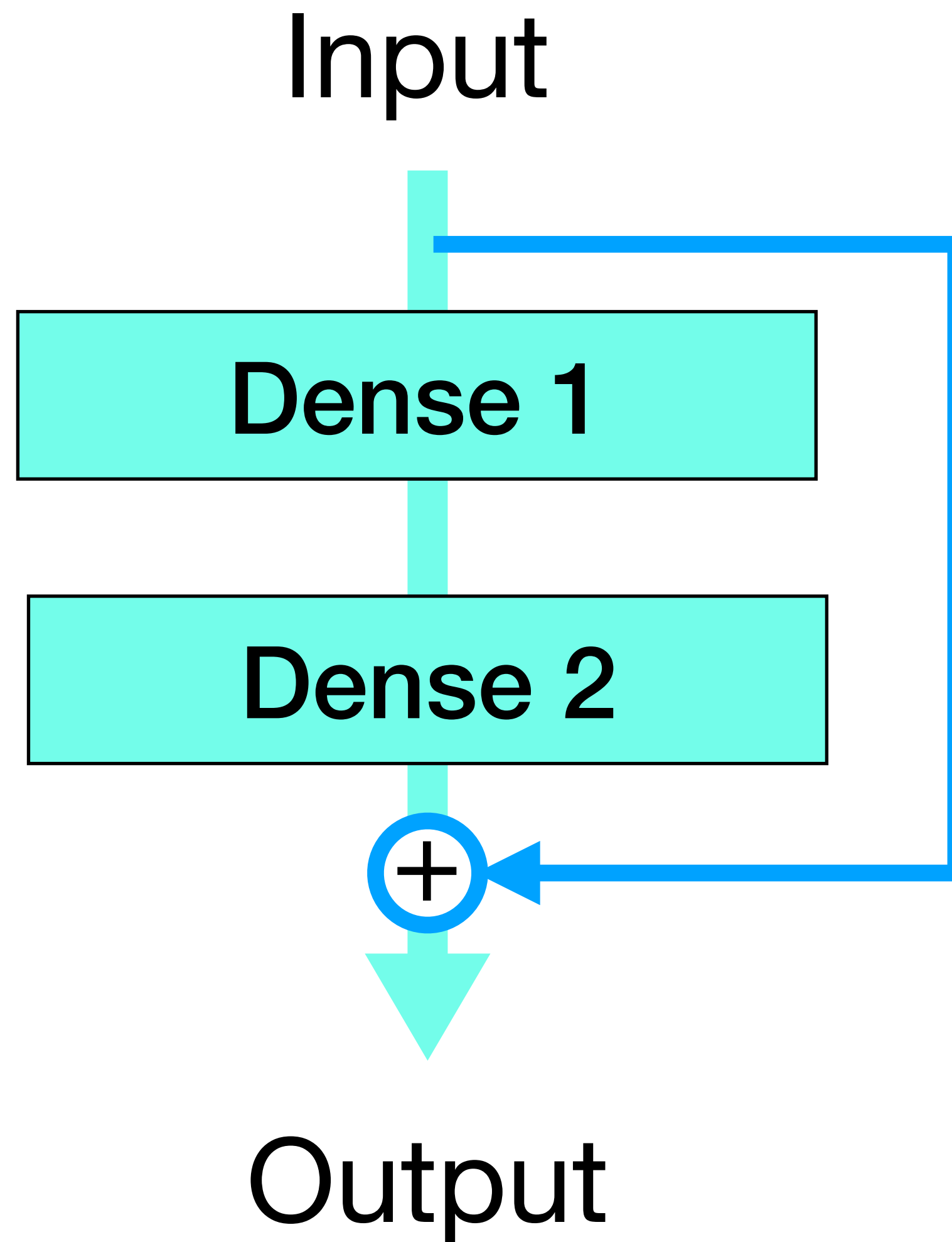
Dense 1

Dense 2

Output

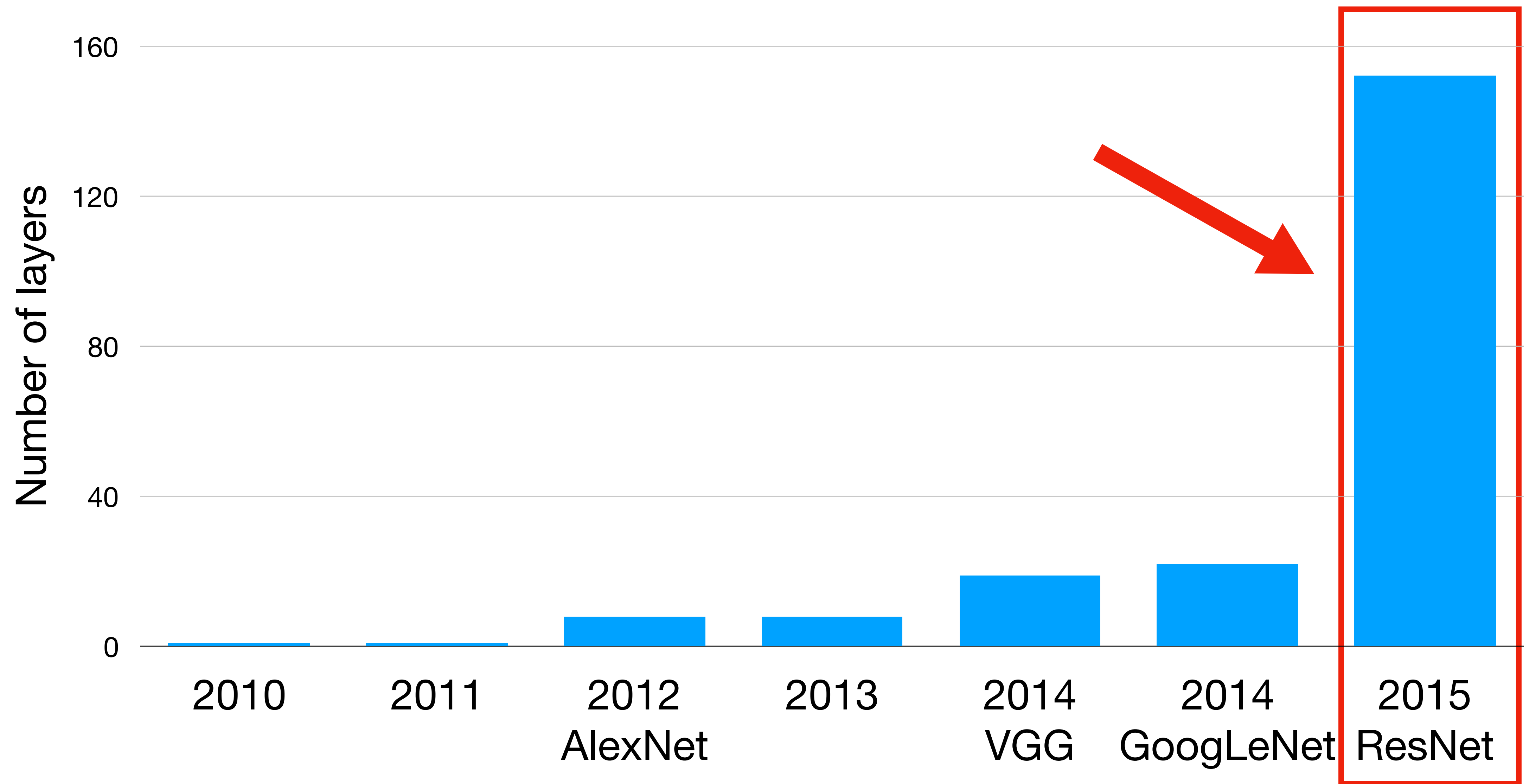
ResNets (He et. al., 2015)

Skip connections



ResNets (He et. al., 2015)

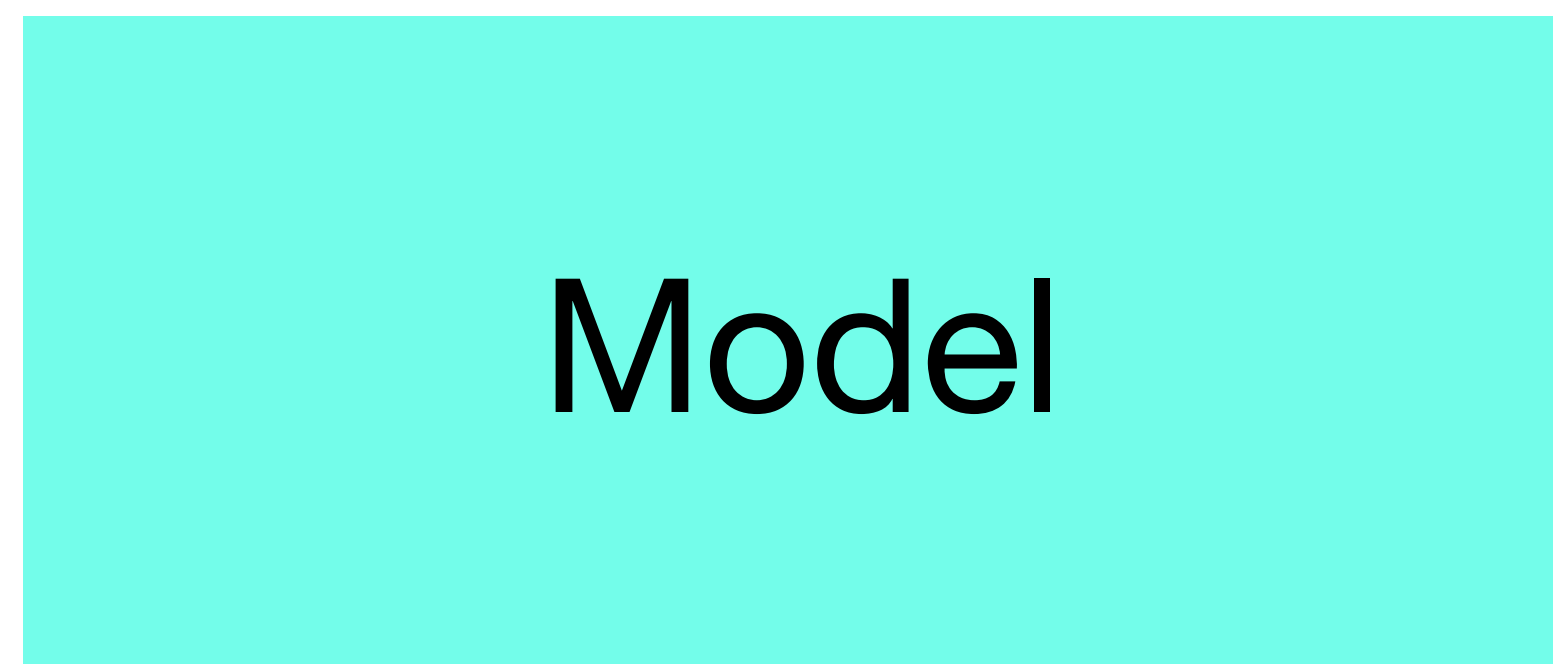
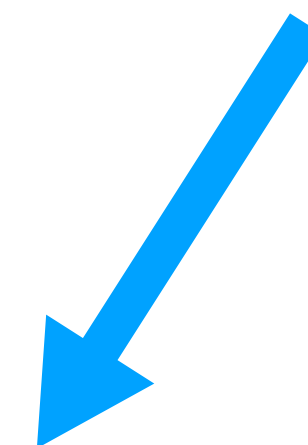
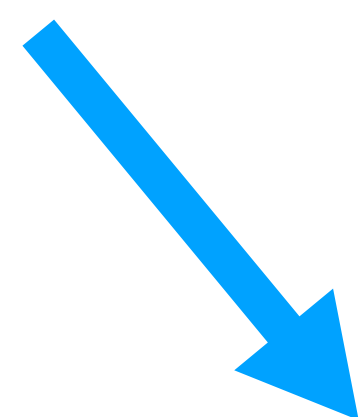
Revolution of Depth: Top models on ImageNet



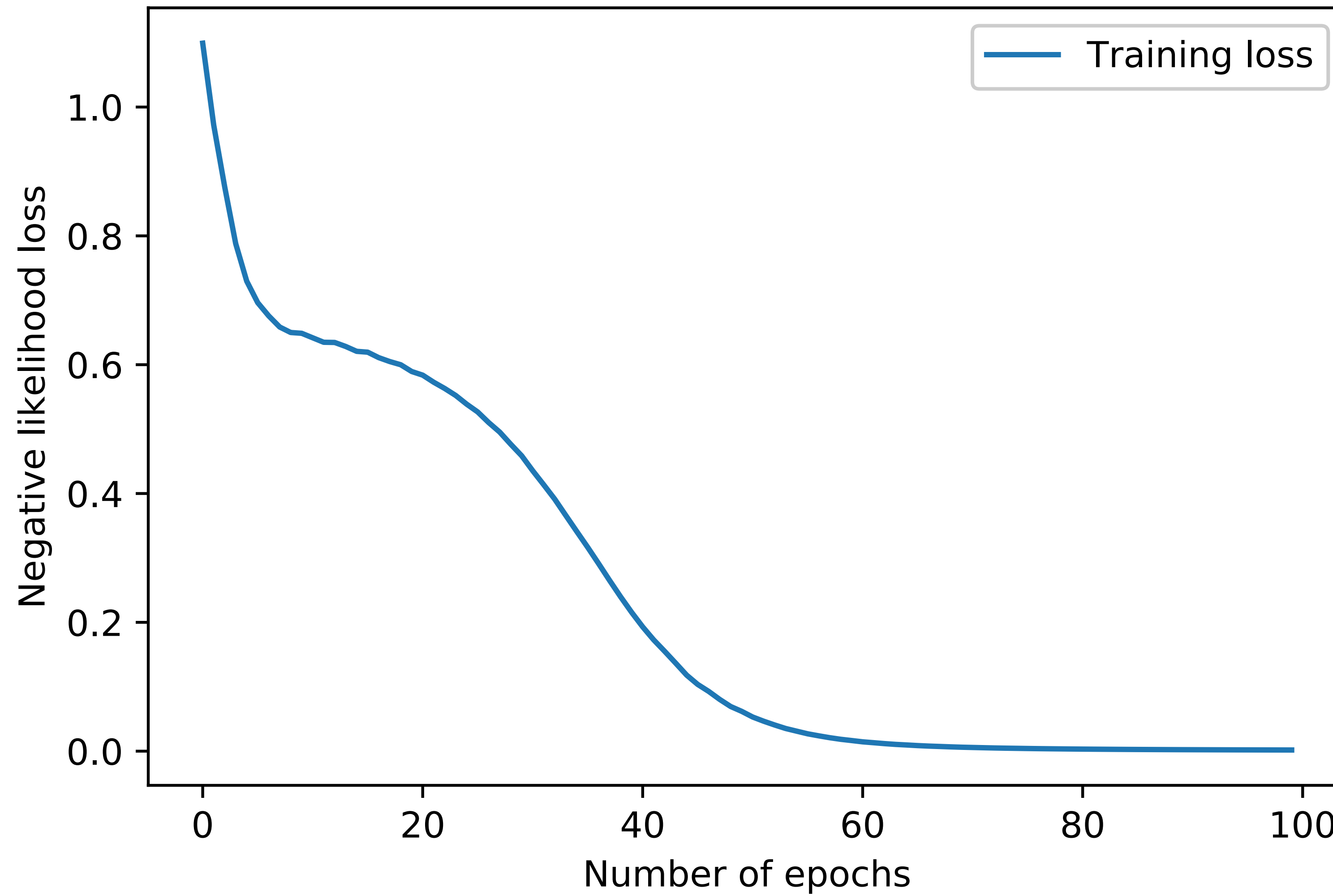
IV. Practical Tips

Overfit First

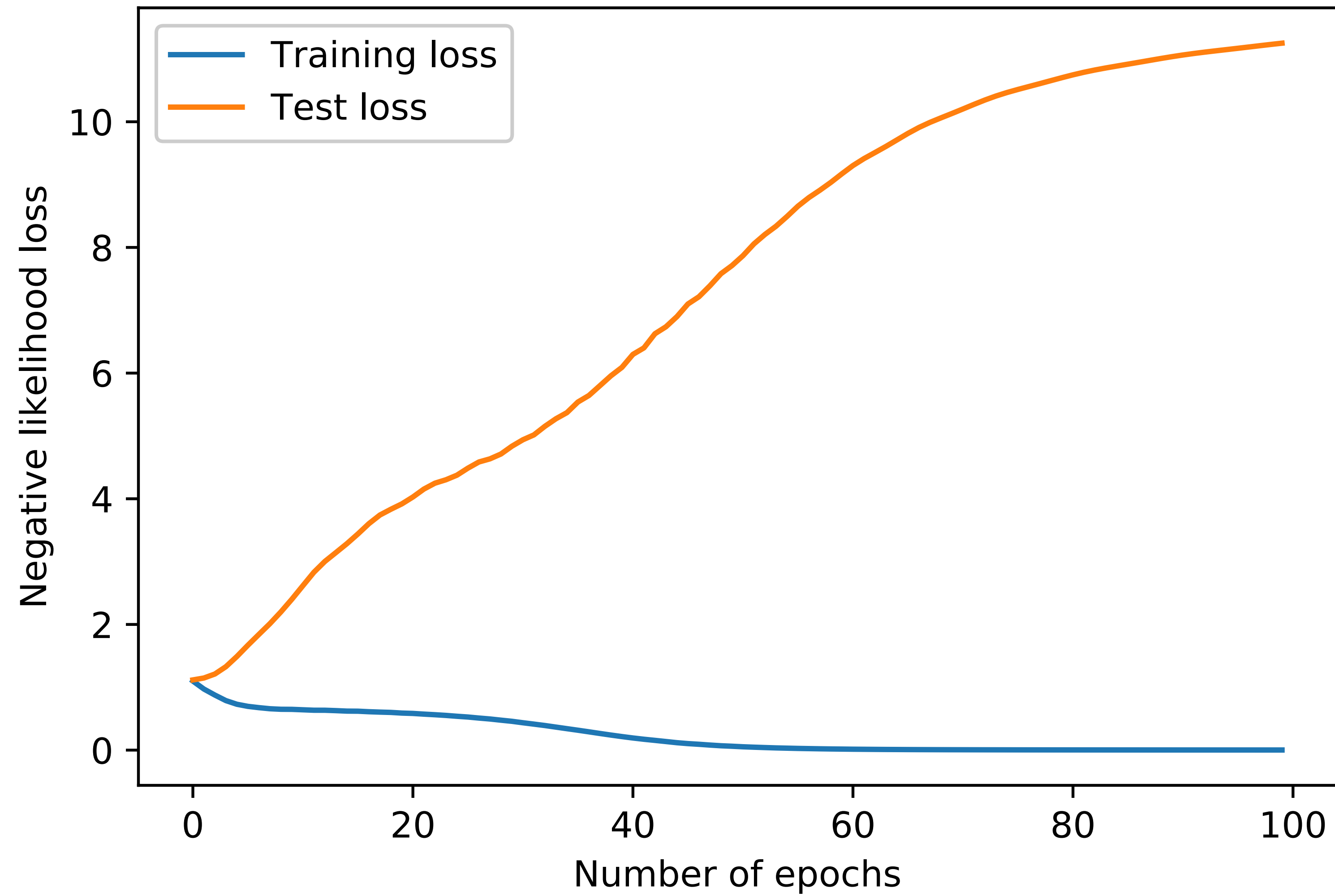
a few examples



Train with 3 examples

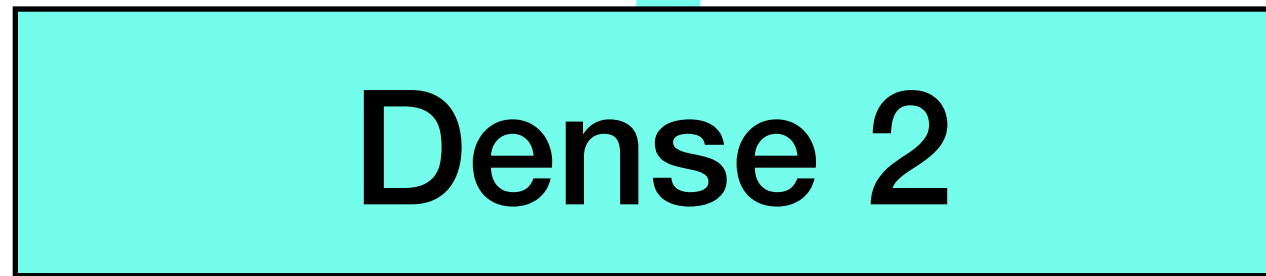
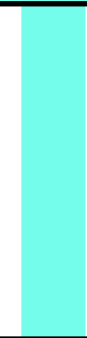
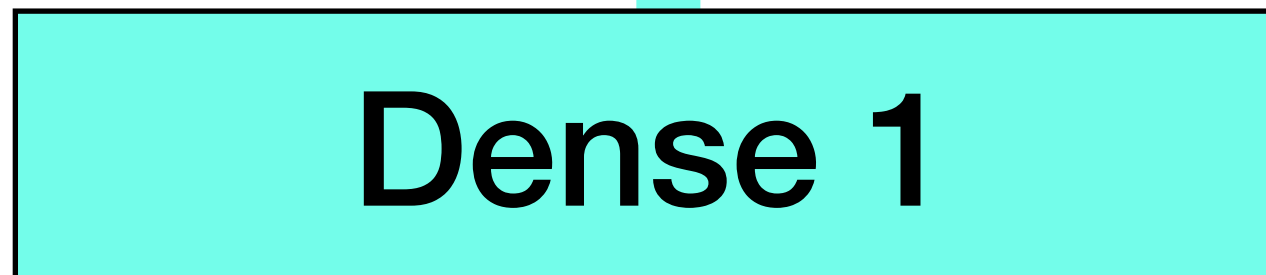
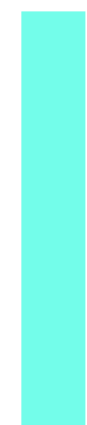


Train with 3 examples



Batchnorm

Input



Output

Input

Batchnorm

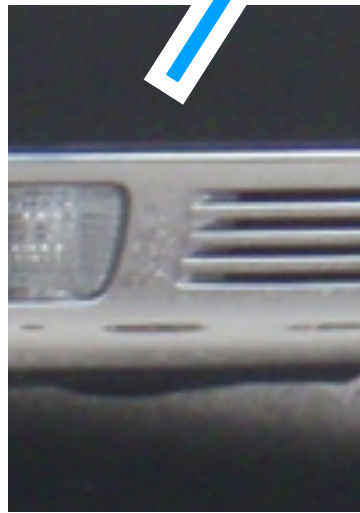
Dense 1

Batchnorm

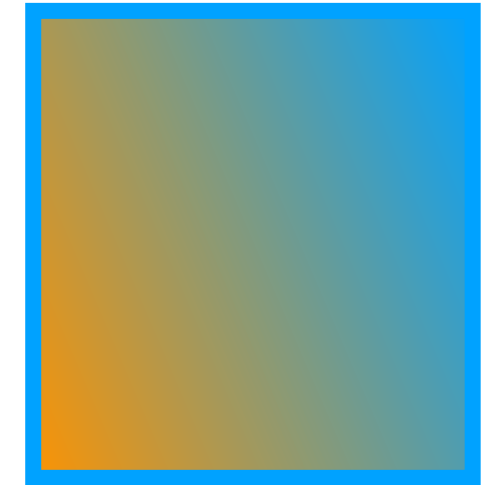
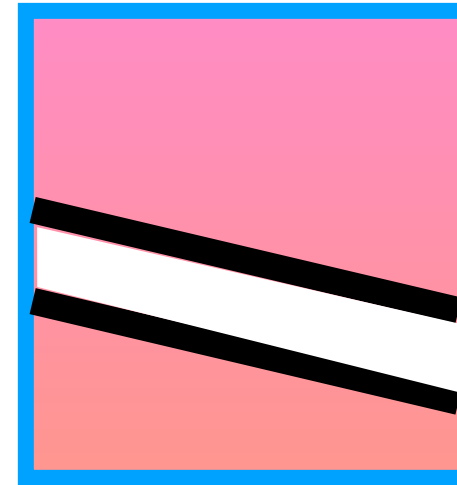
Dense 2

Output

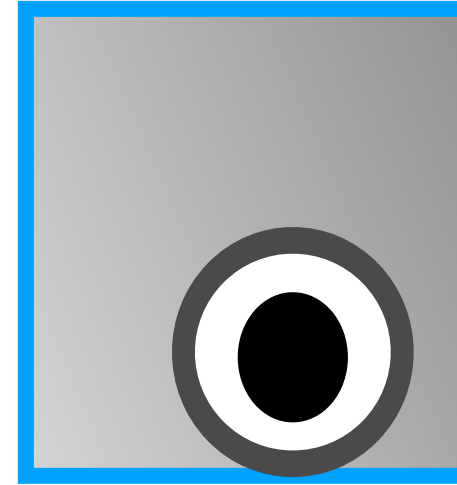
Transfer Learning



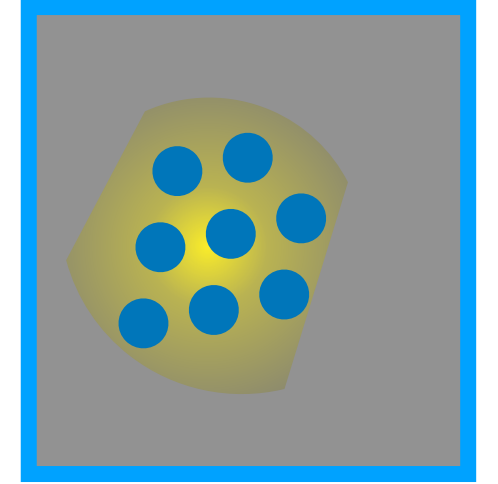
Low



Mid



High



**Trained
network**

e.g. ResNet

Last layer(s)

Mid/high level
features

**Trained
network**

e.g. ResNet

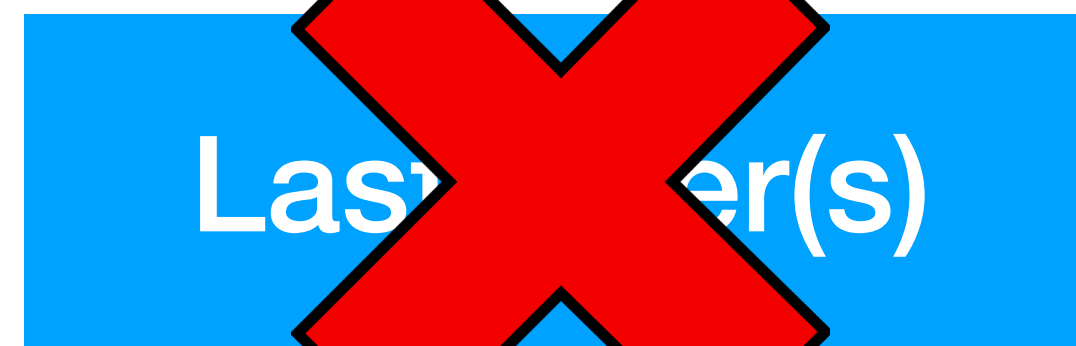
Last layer(s)

Mid/high level
features



**Trained
network**

e.g. ResNet

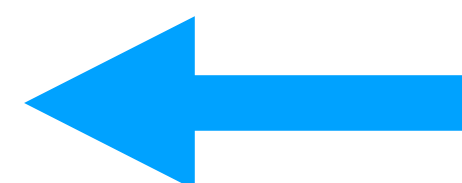
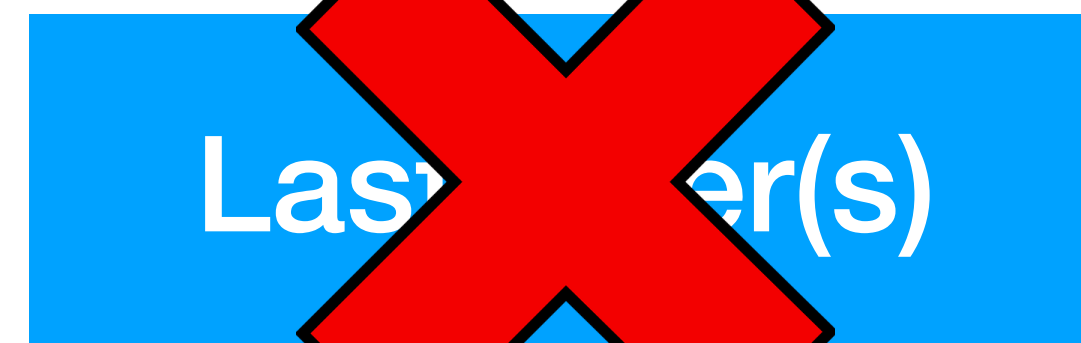


Last layer(s)

Mid/high level
features



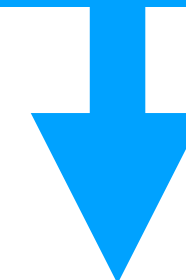
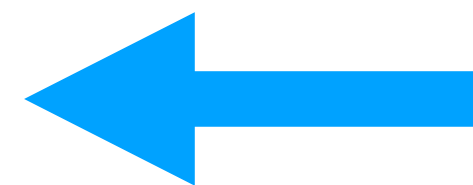
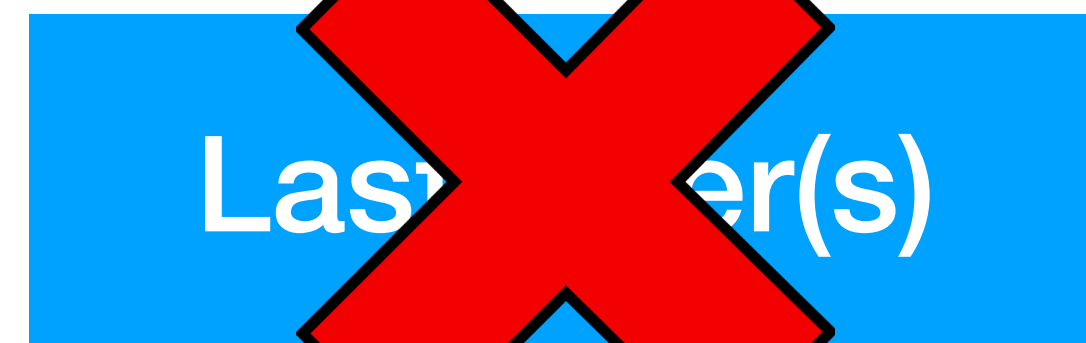
e.g. ResNet



Mid/high level features



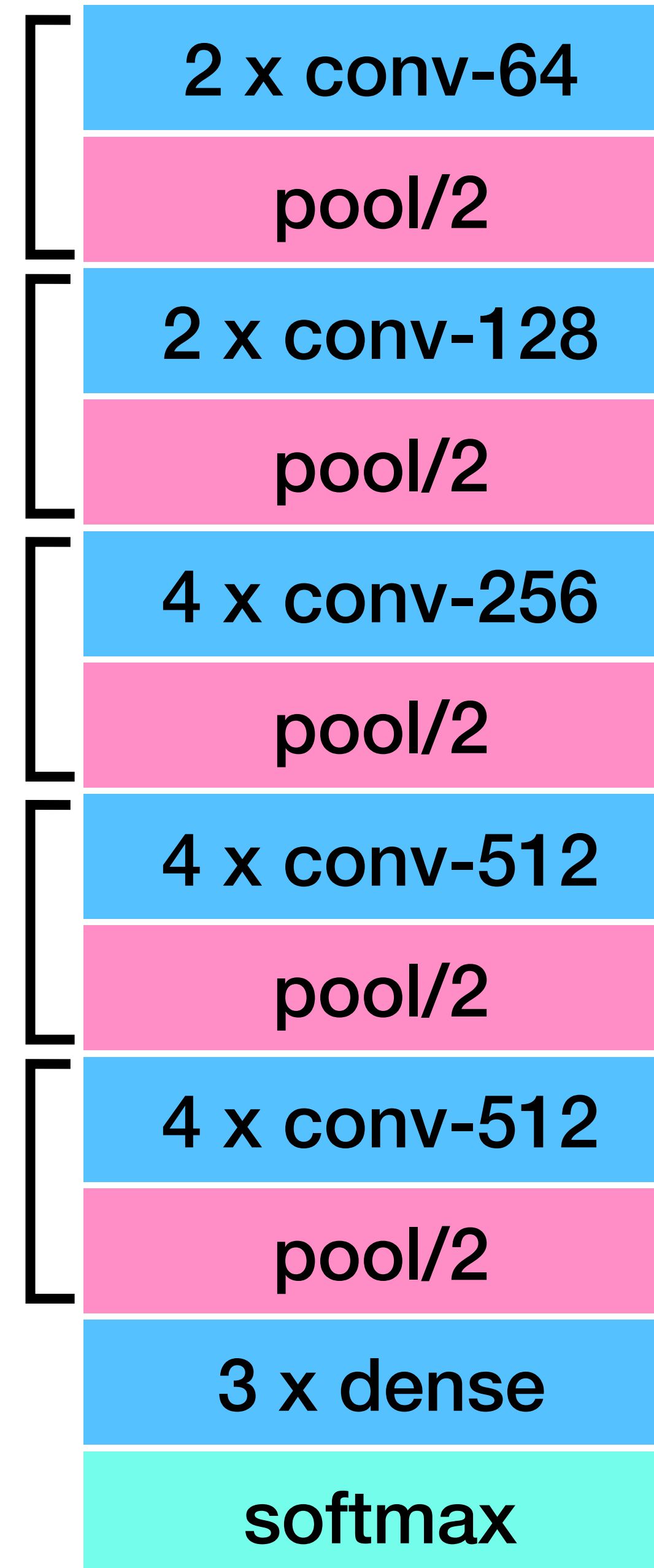
e.g. ResNet



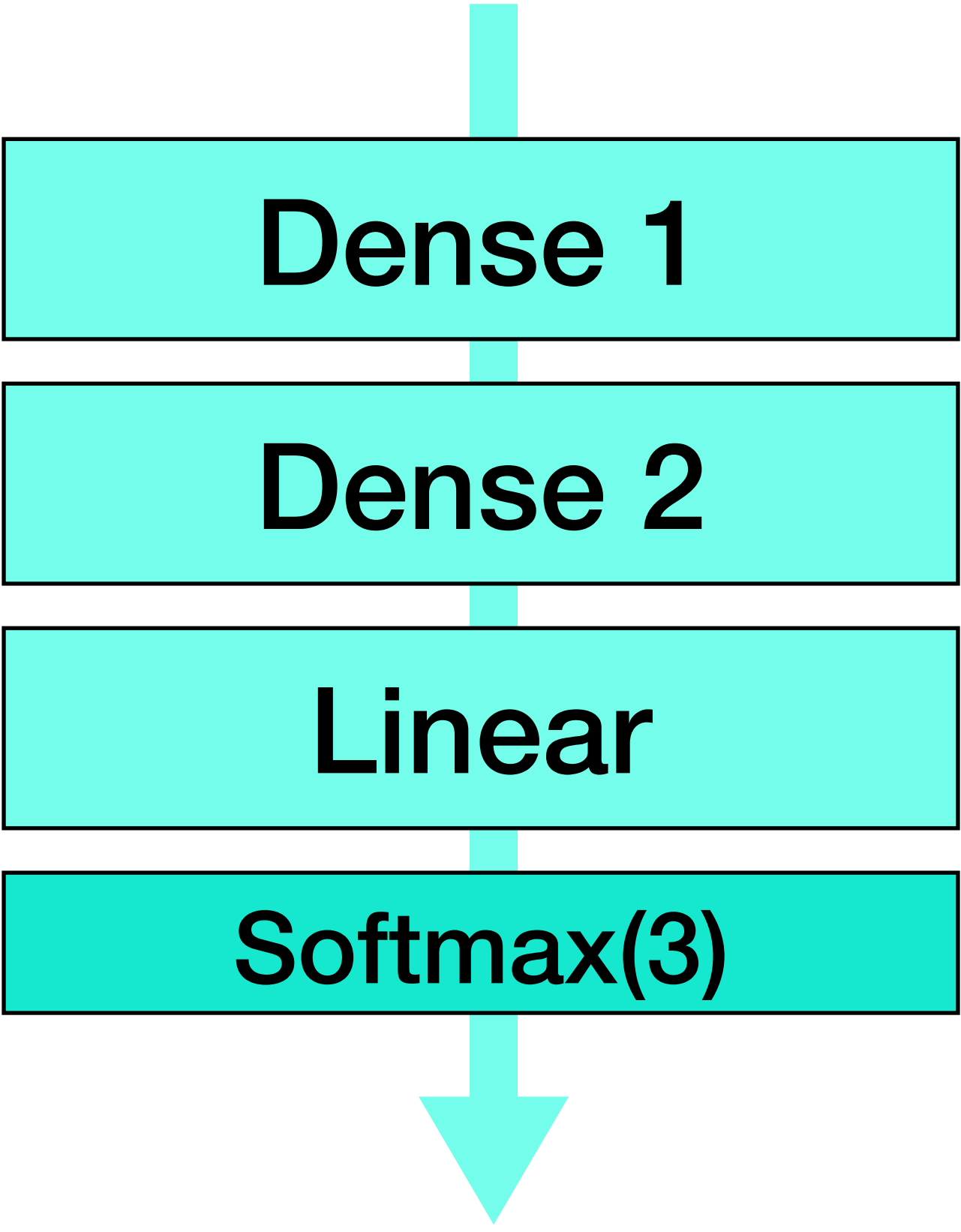
New output

Simplify

VGG-19



Dense 1



Input

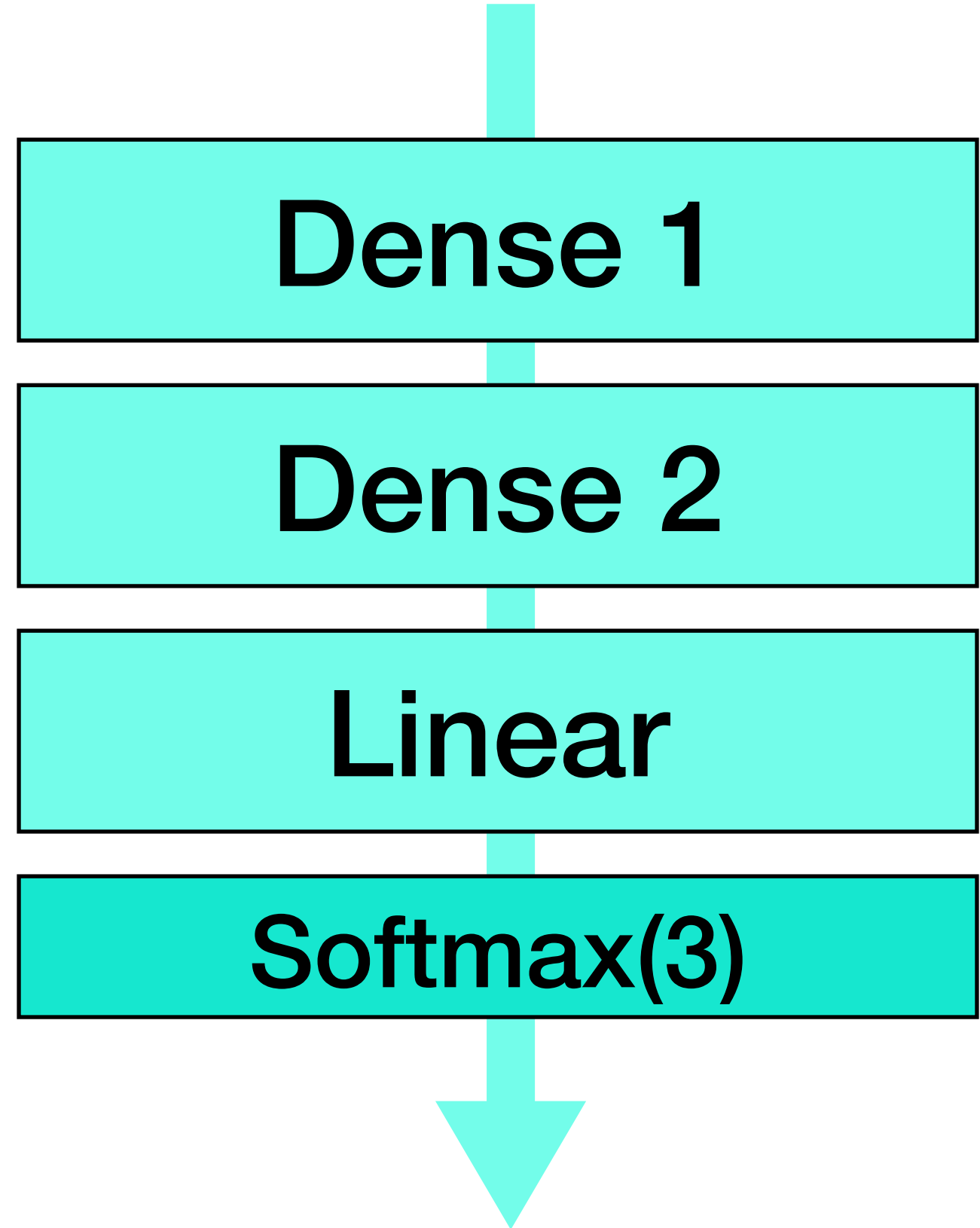
Dense 1

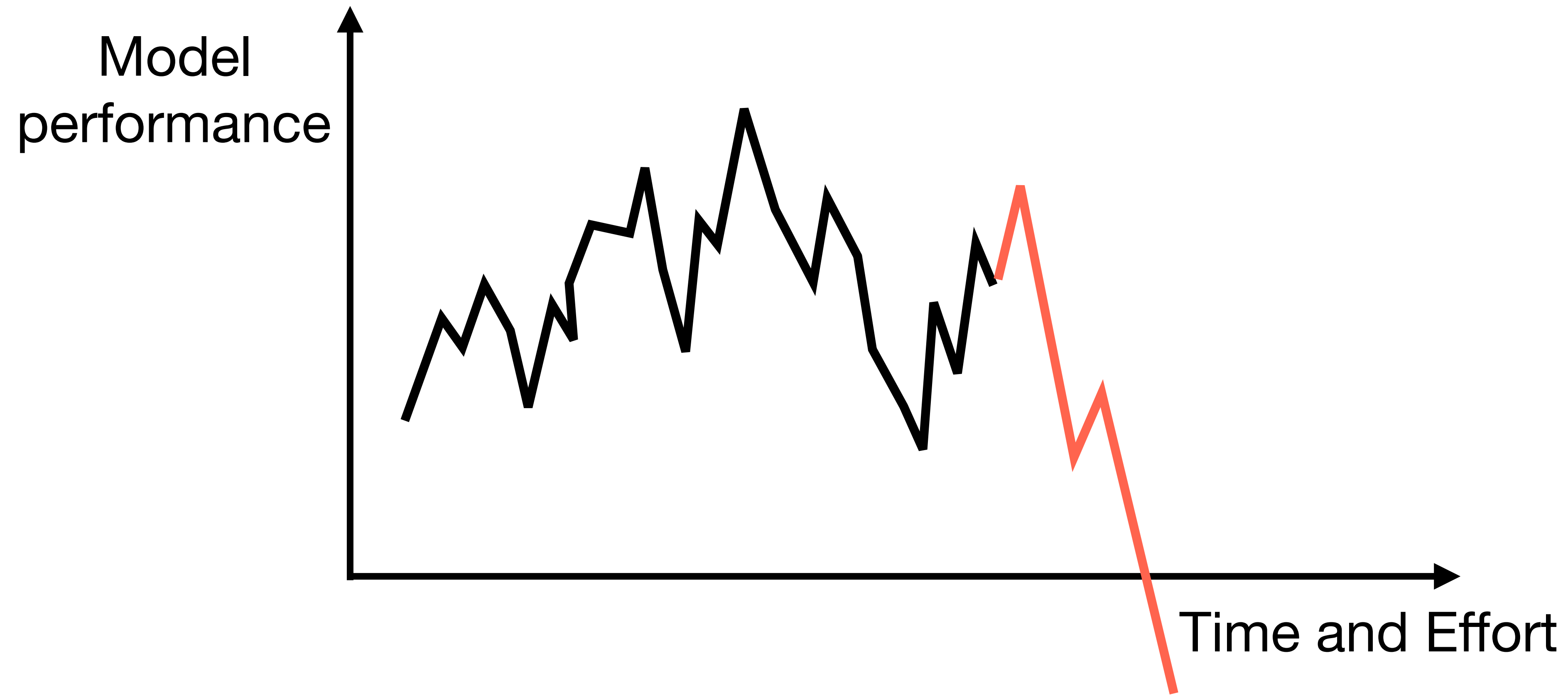
Dense 2

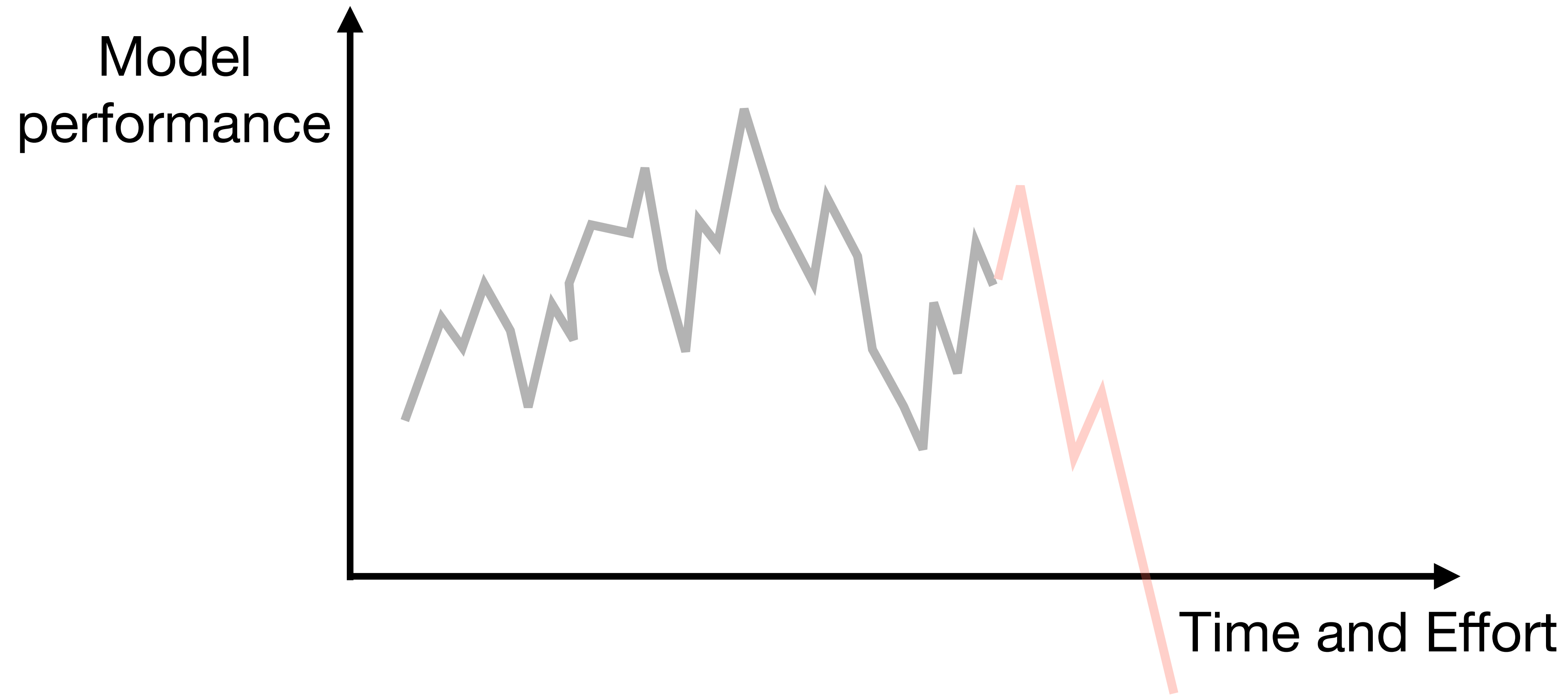
Linear

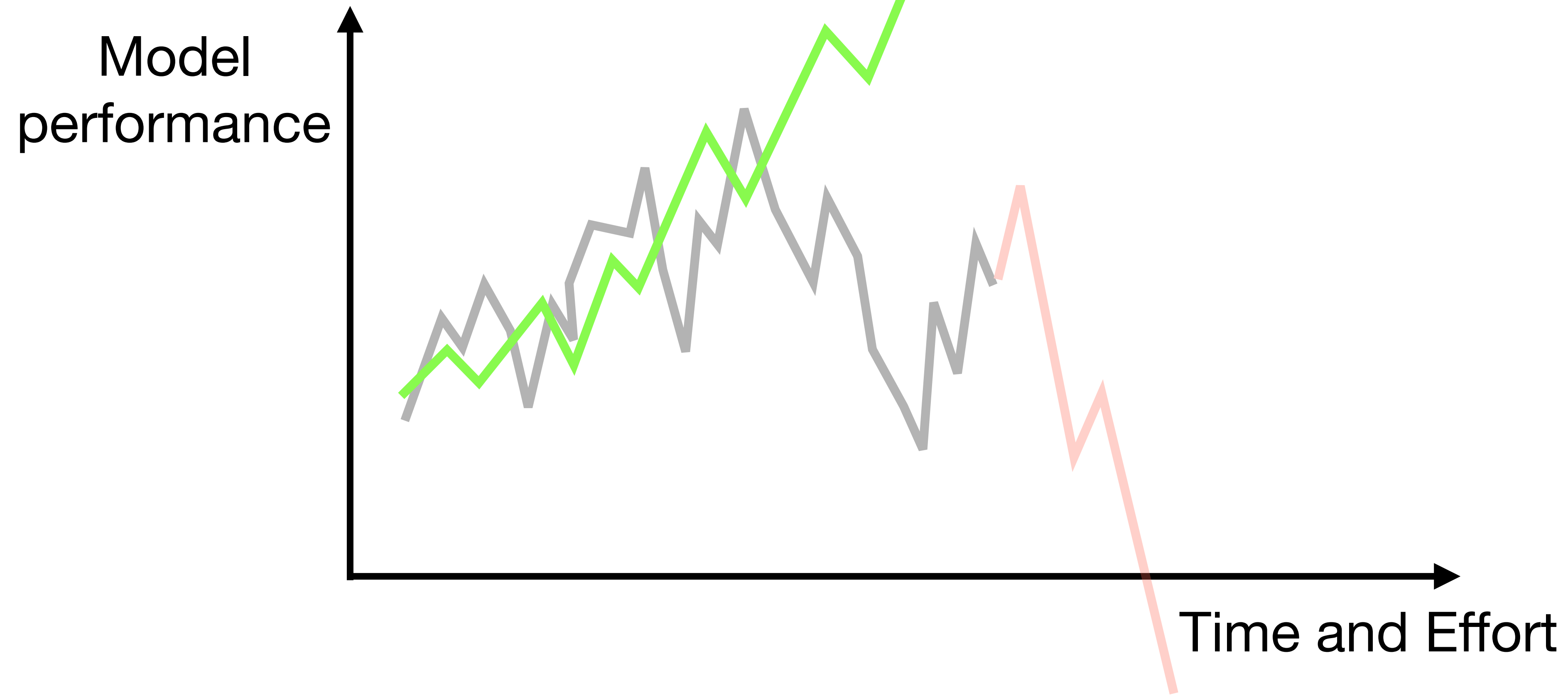
Softmax(3)

Output









Thank You!

Appendix

Resources on Coding Neural Networks

- Jason Brownlee's Machine Learning Mastery (lots of code)
- My blog posts on training models in TensorFlow (training a model, MLP, CNN)

Resources on Deep Learning

- [THE Deep Learning book \(Ian Goodfellow, Yoshua Bengio, Aaron Courville\)](#)
- [Stanford's course on Convolutional Neural Networks](#)
- [fast.ai courses on neural networks](#) (haven't tried much myself but I've heard they are good)