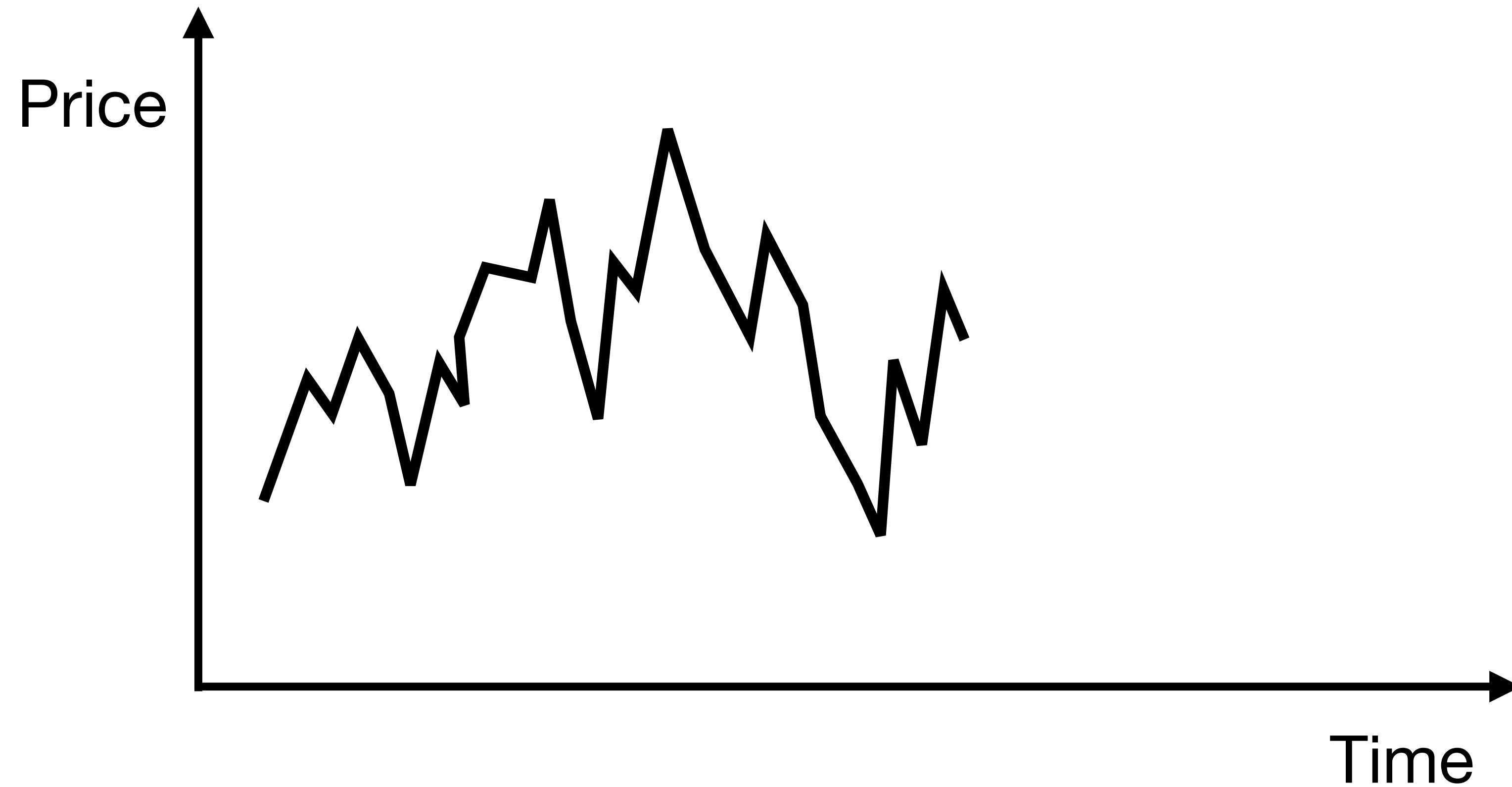
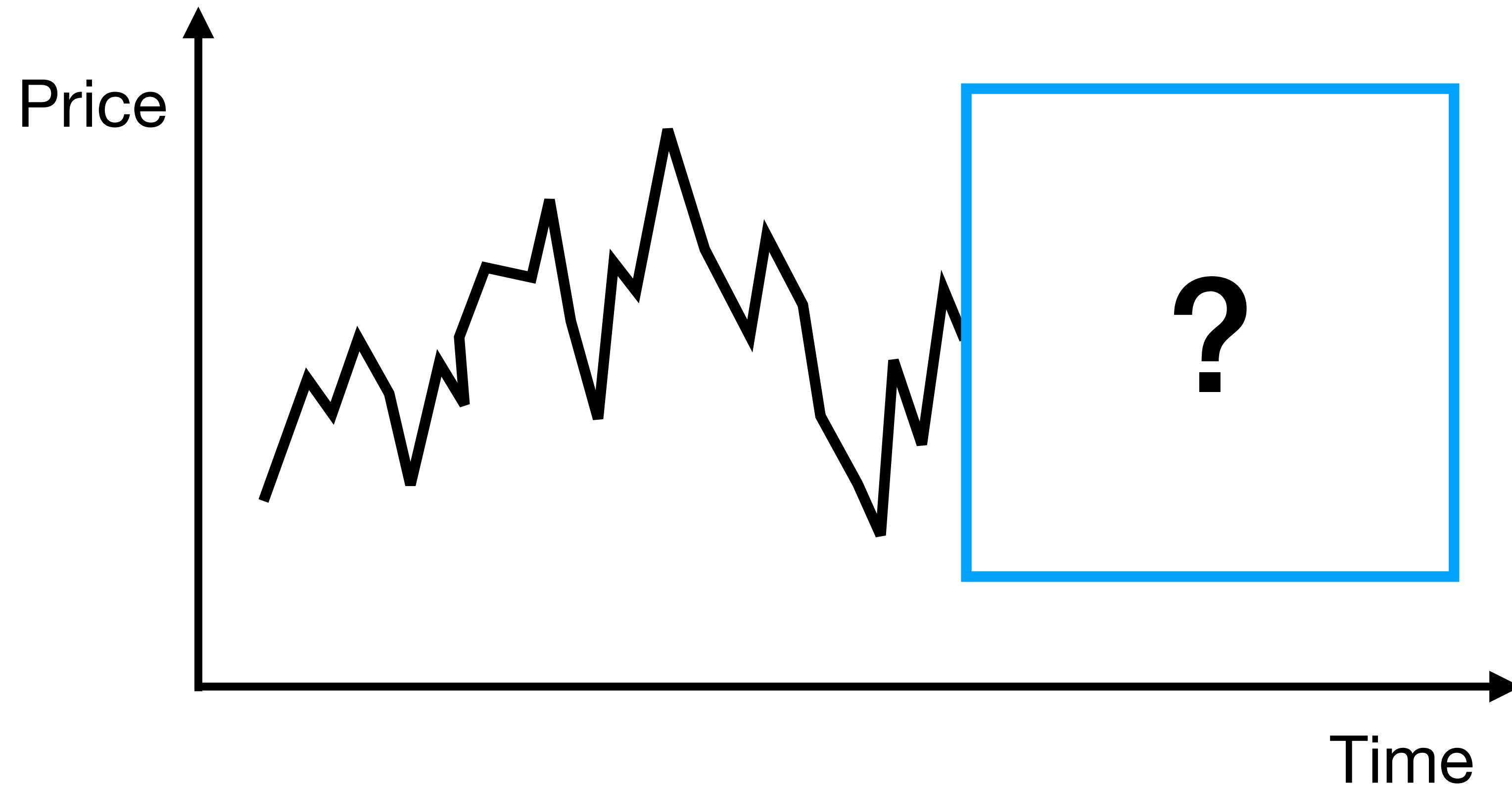
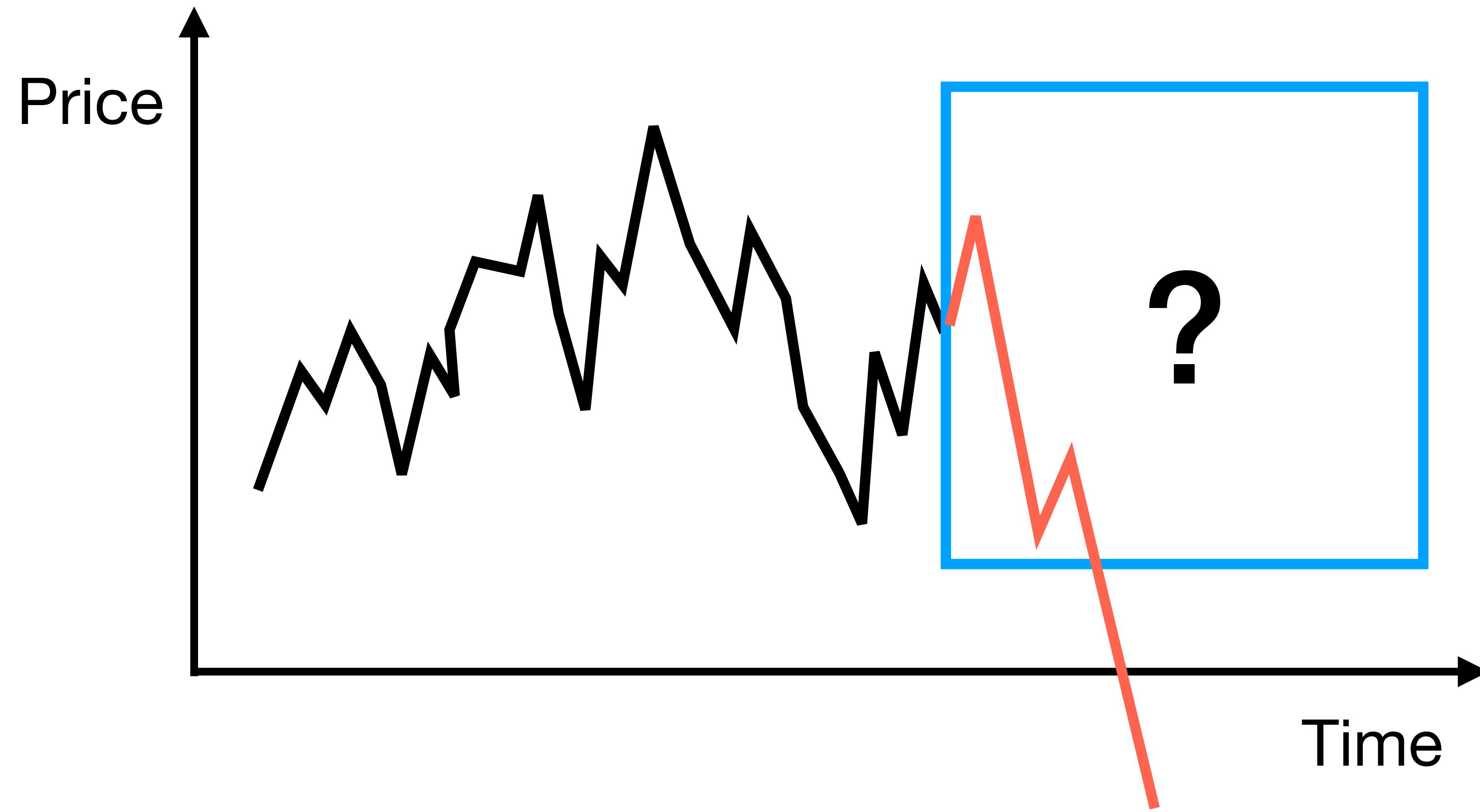


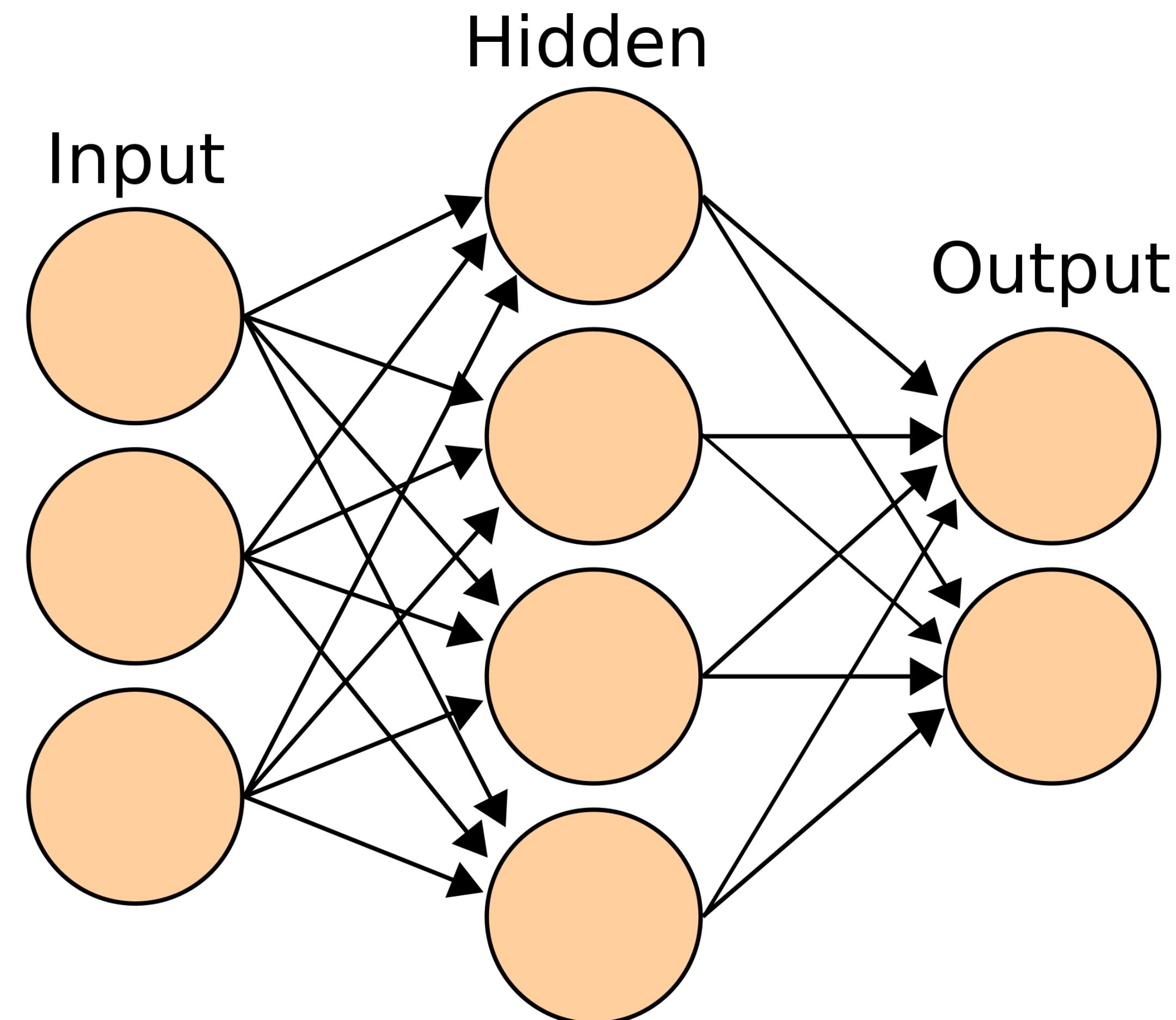
Understanding Deep Learning

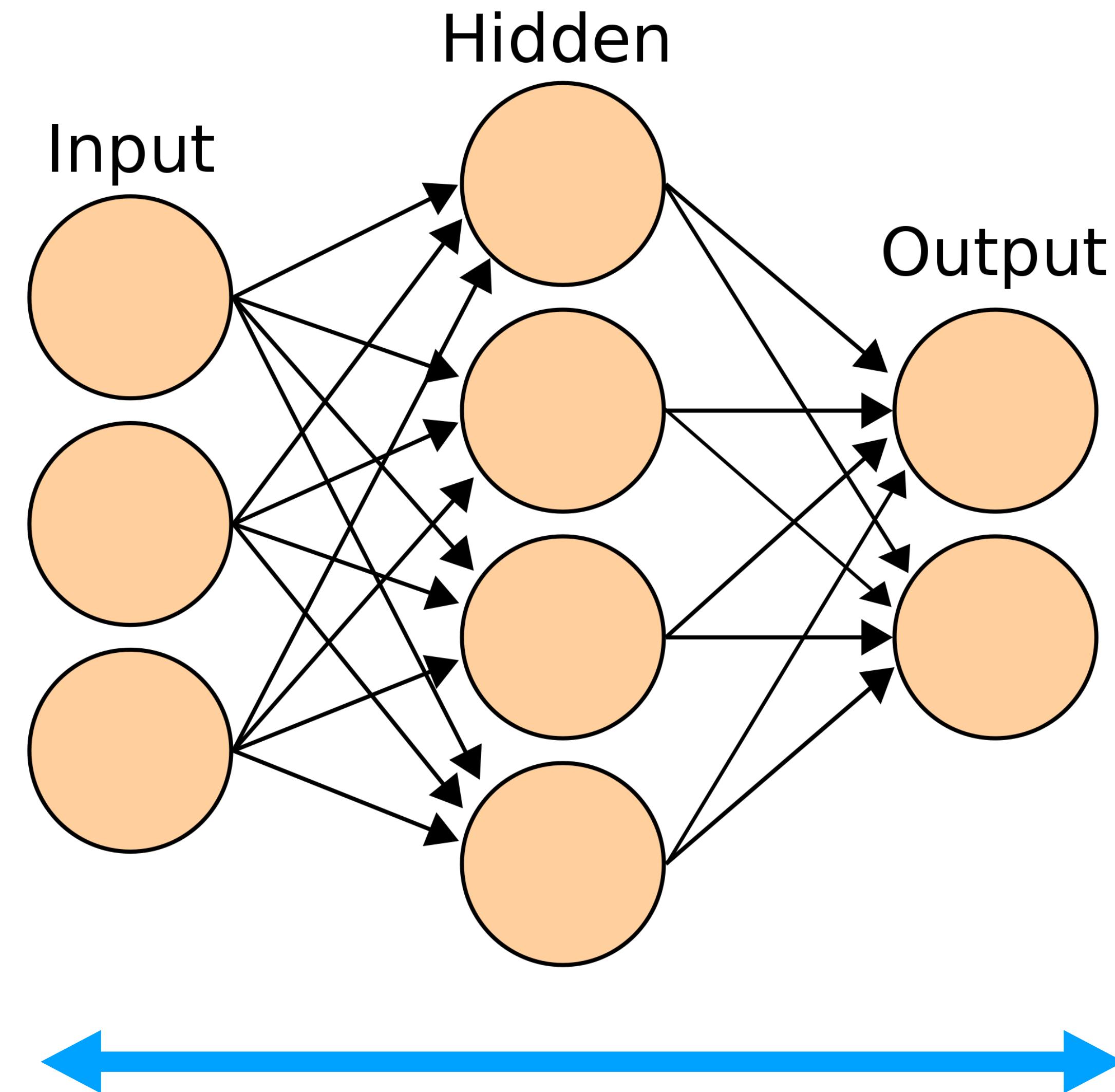
Jessica Yung

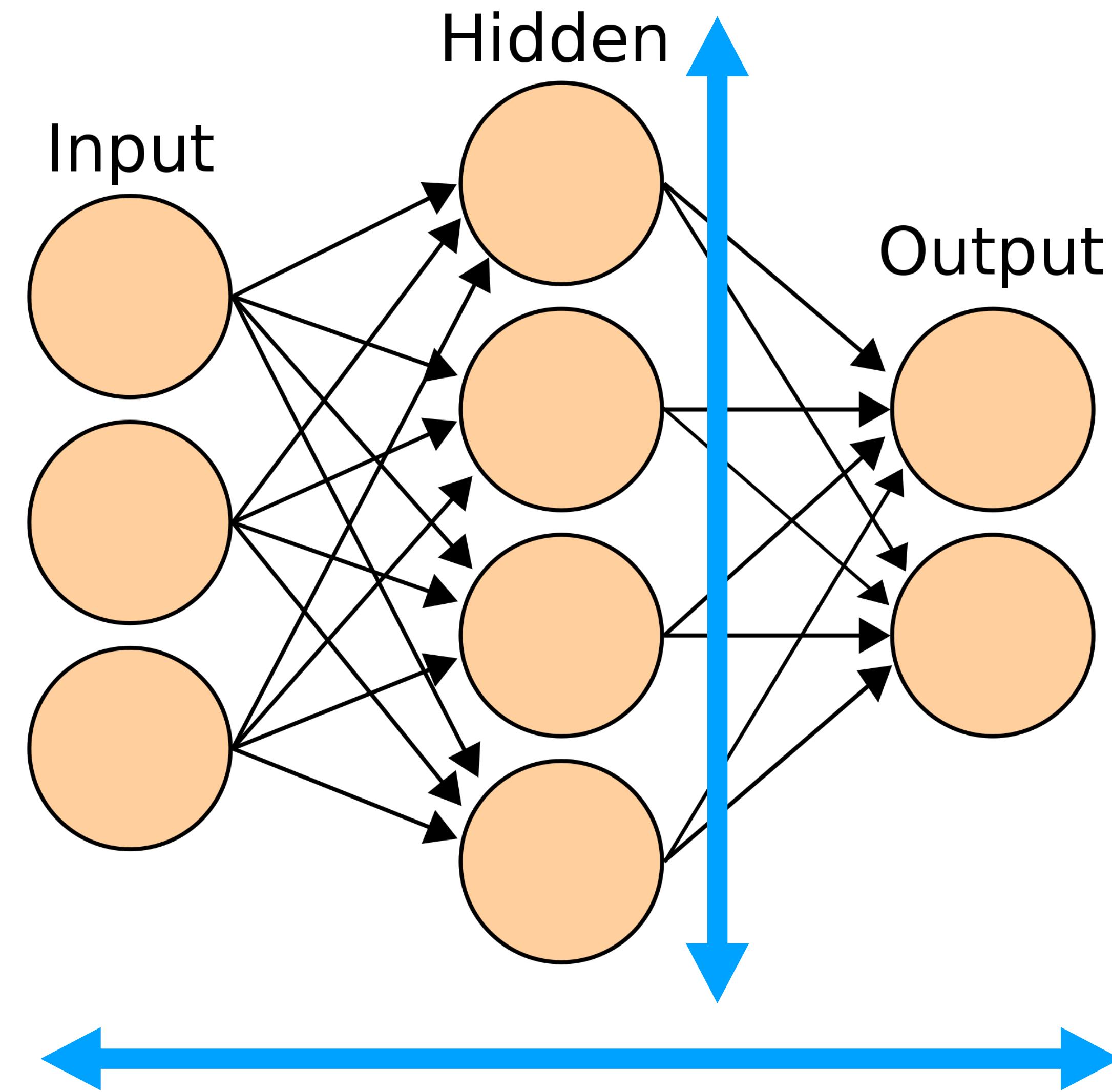


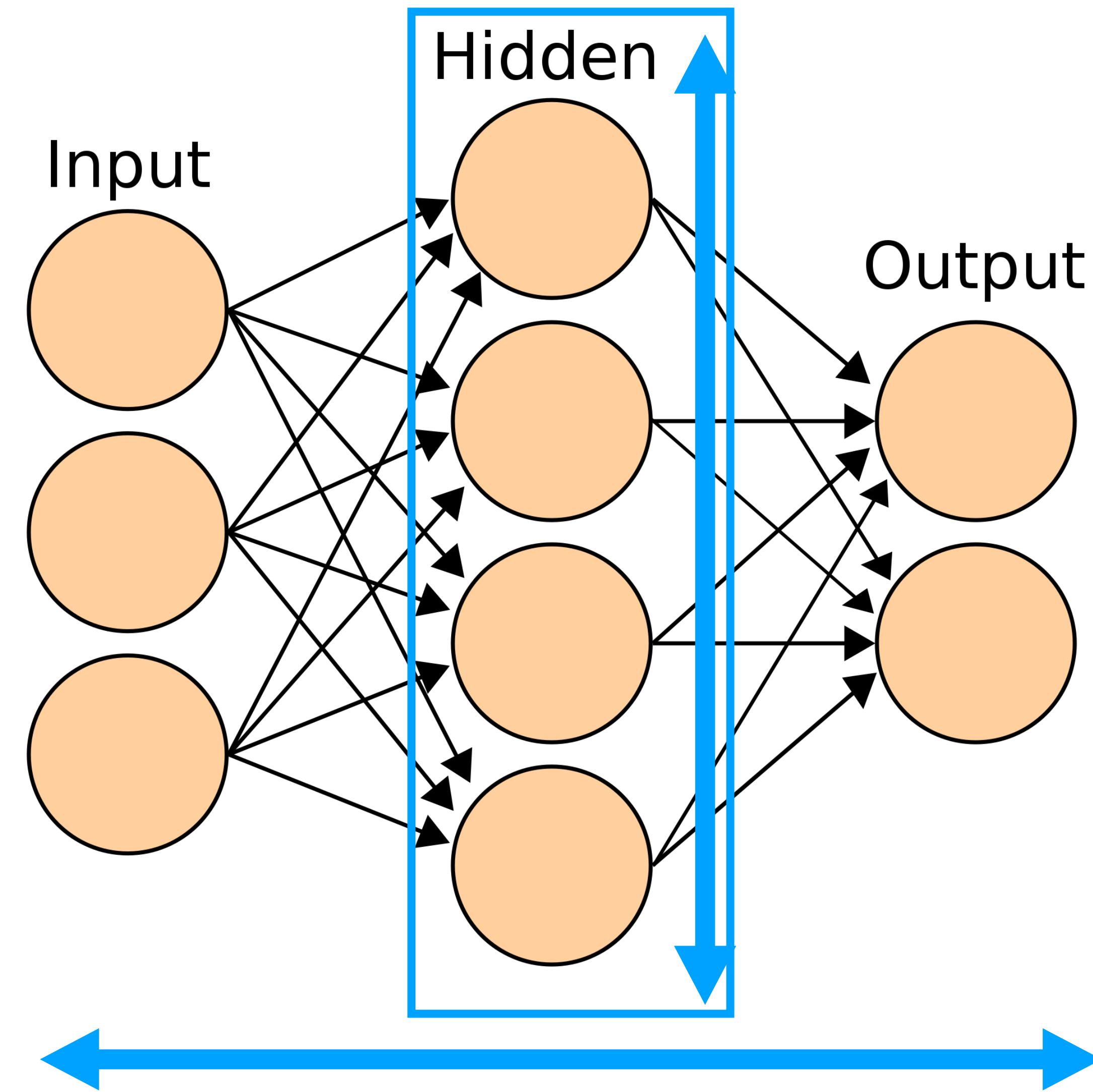


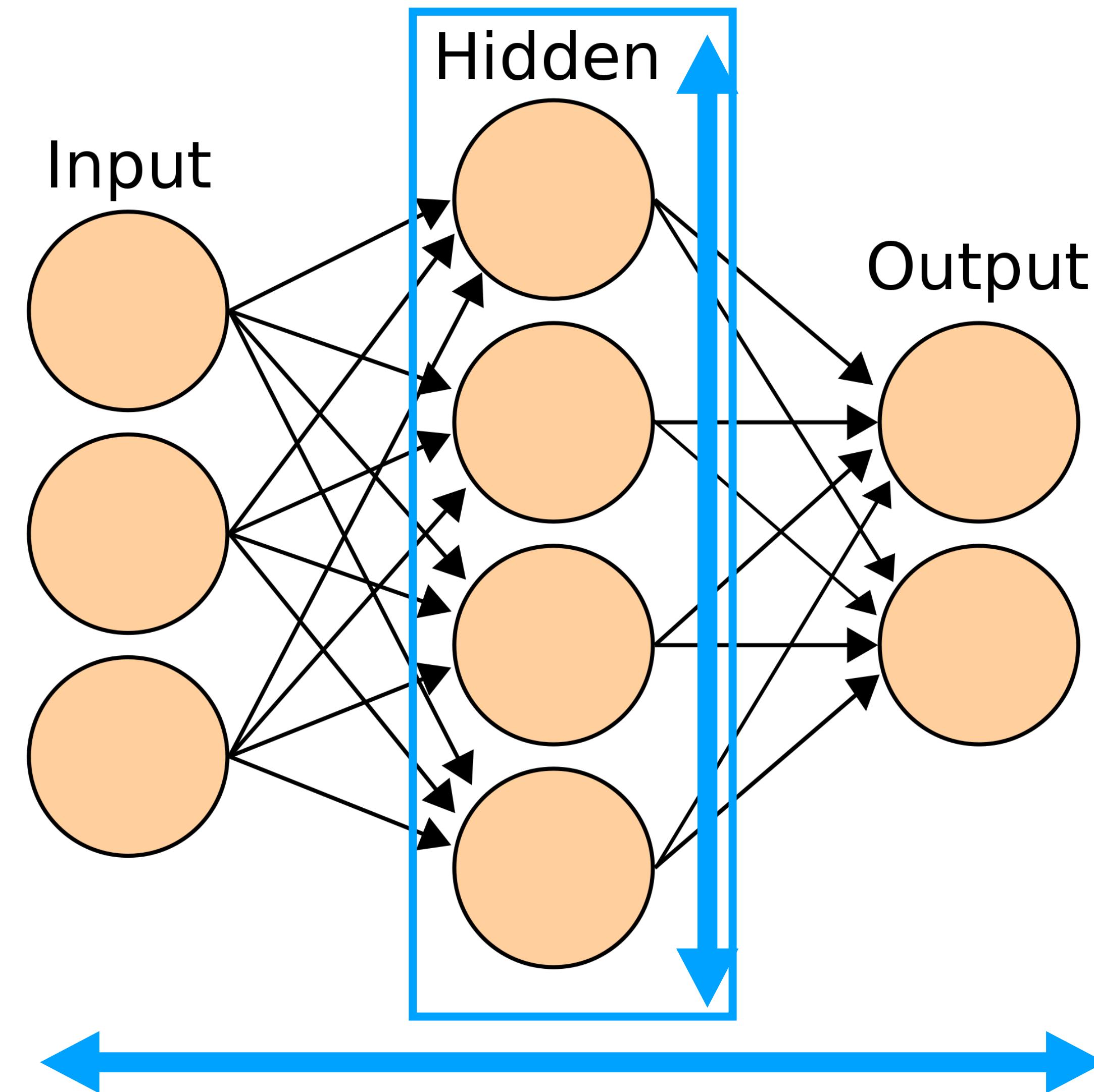




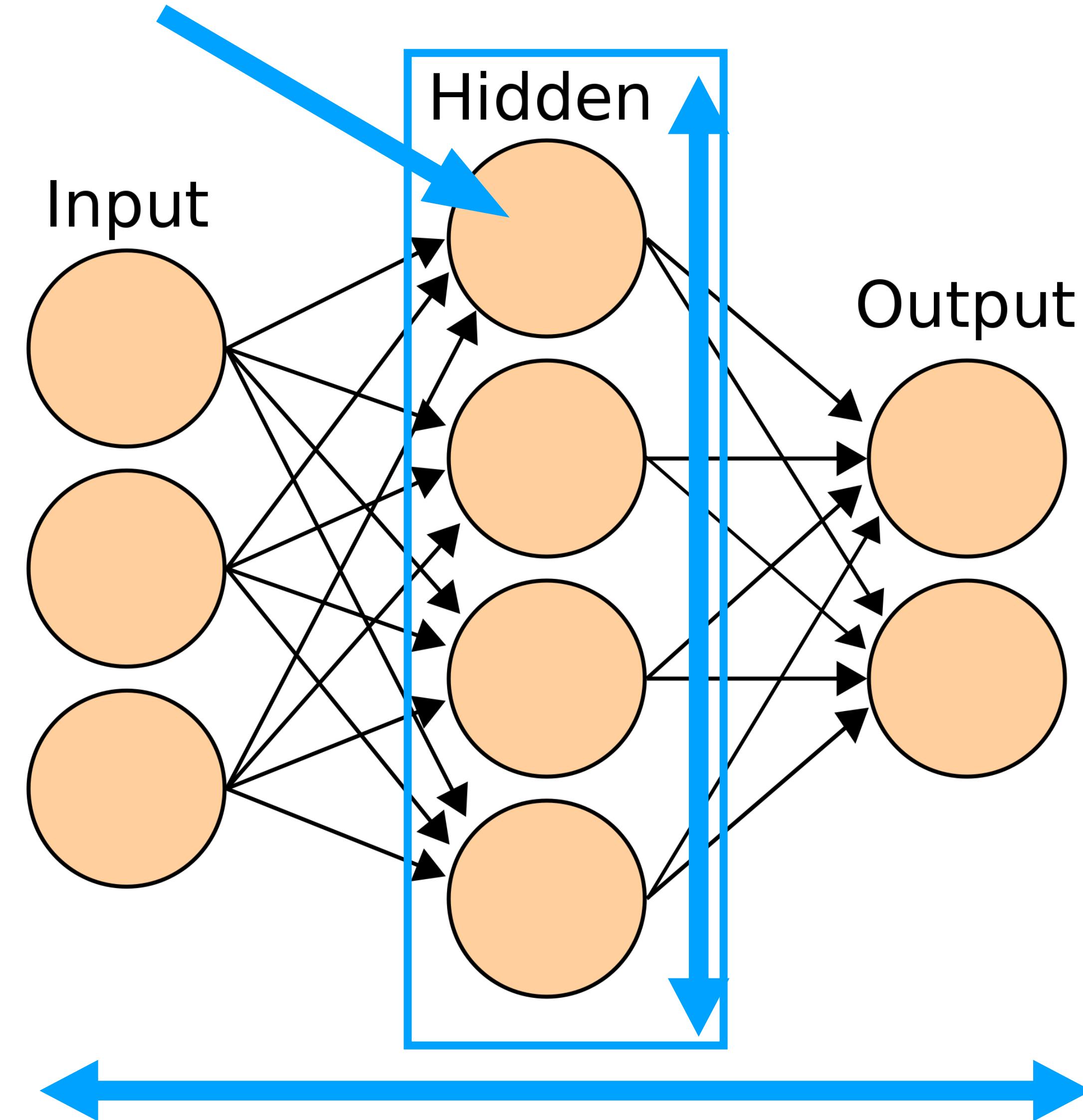




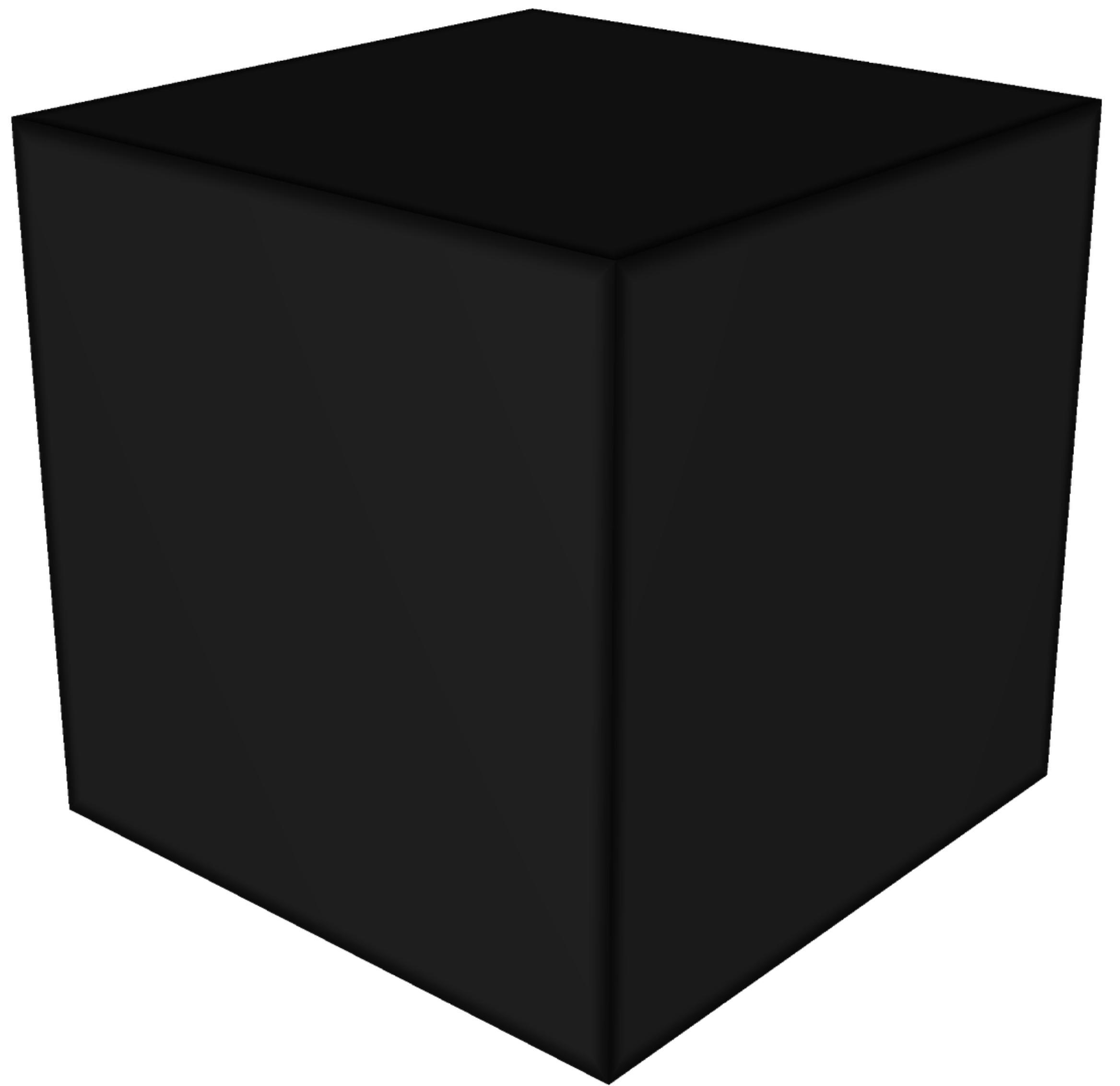


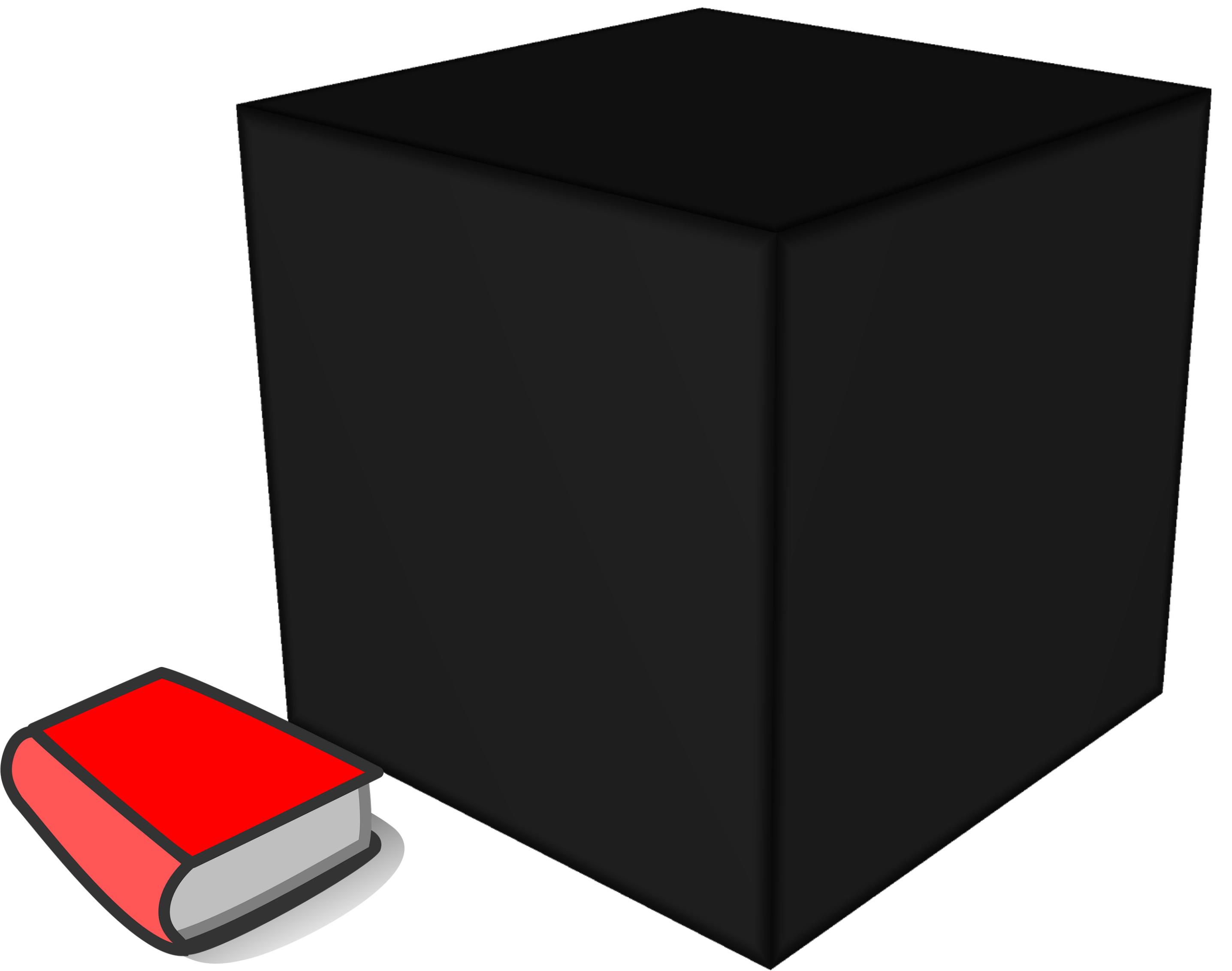


learning rate



learning rate

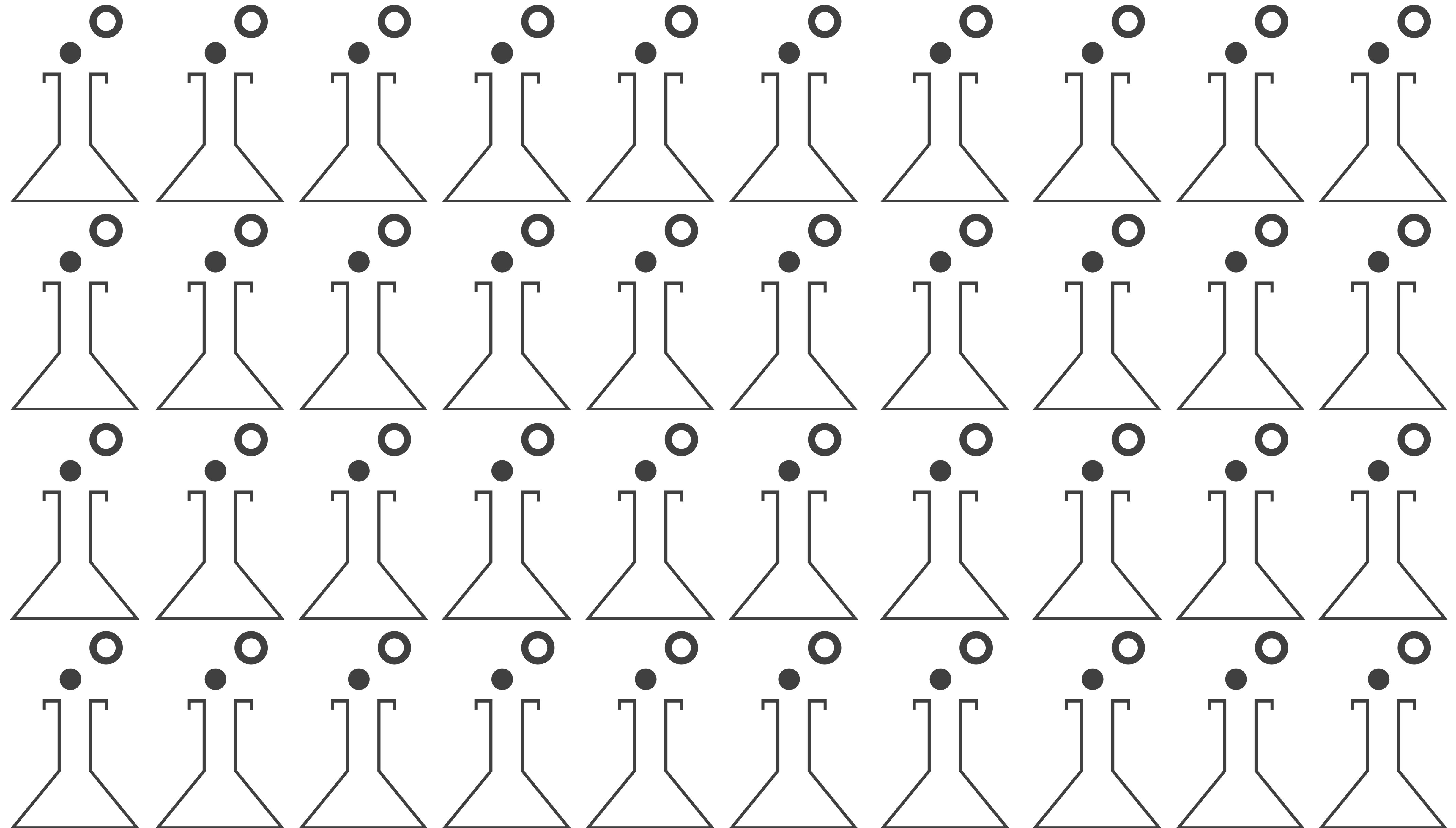


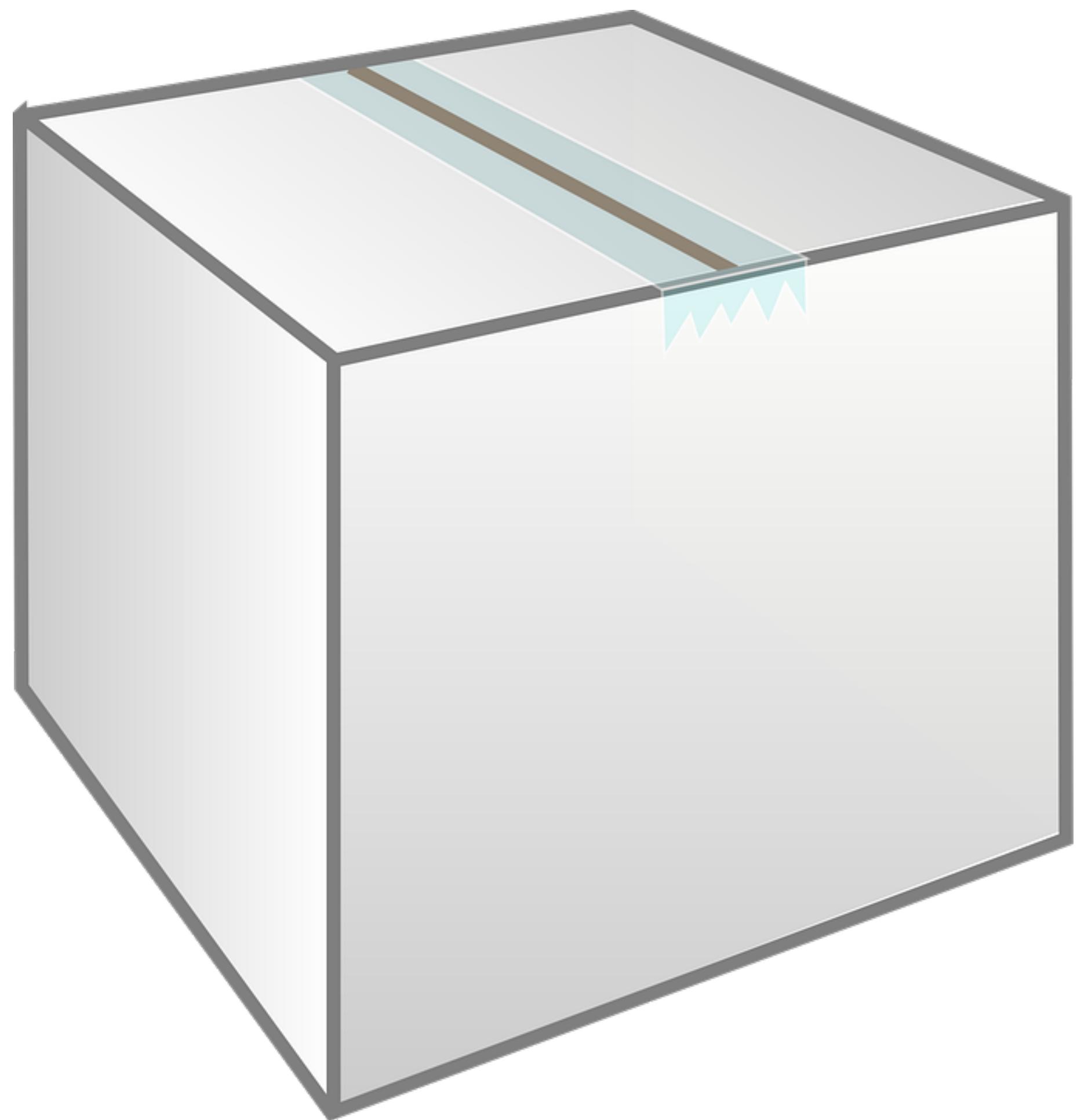


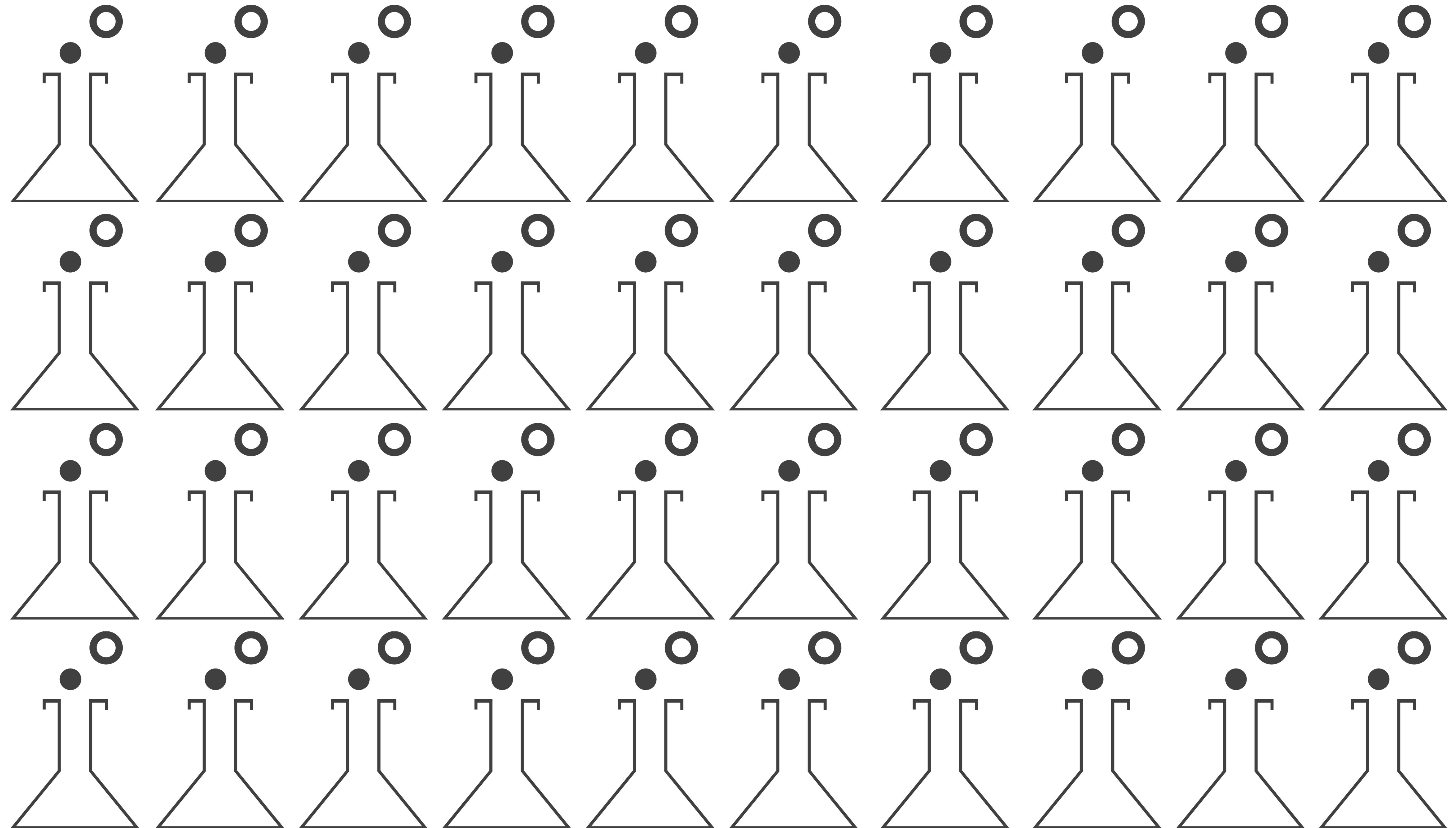
```
def debug():  
    return 0 + "no"
```

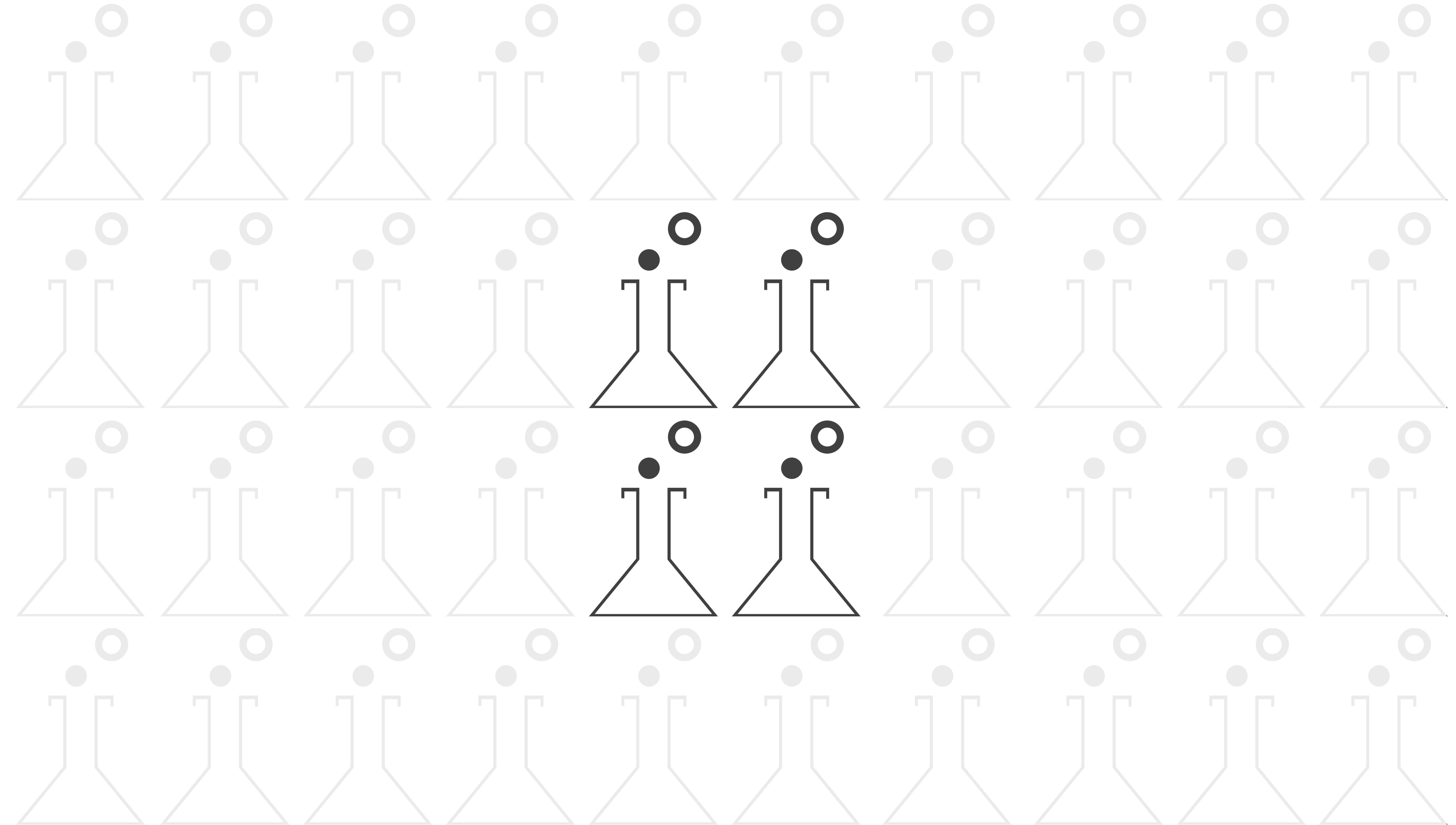
567 563m8dep

j61mj? _ g “?@”









Why deep learning?

(Super) human performance in

(Super) human performance in



(Super) human performance in

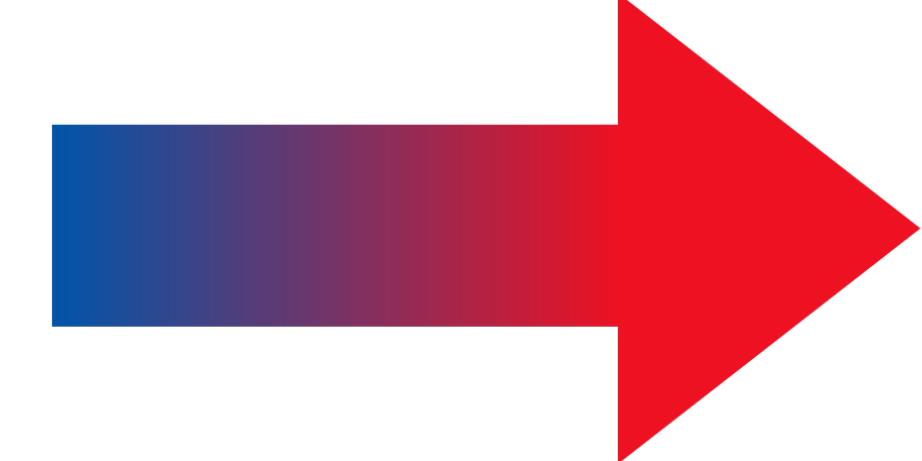


あ → A

A large black Japanese character 'あ' is positioned to the left of a horizontal arrow. The arrow is blue on the left side and transitions to red on the right side. To the right of the arrow is a large black English capital letter 'A'. The entire graphic is set against a white background.

(Super) human performance in



あ  A



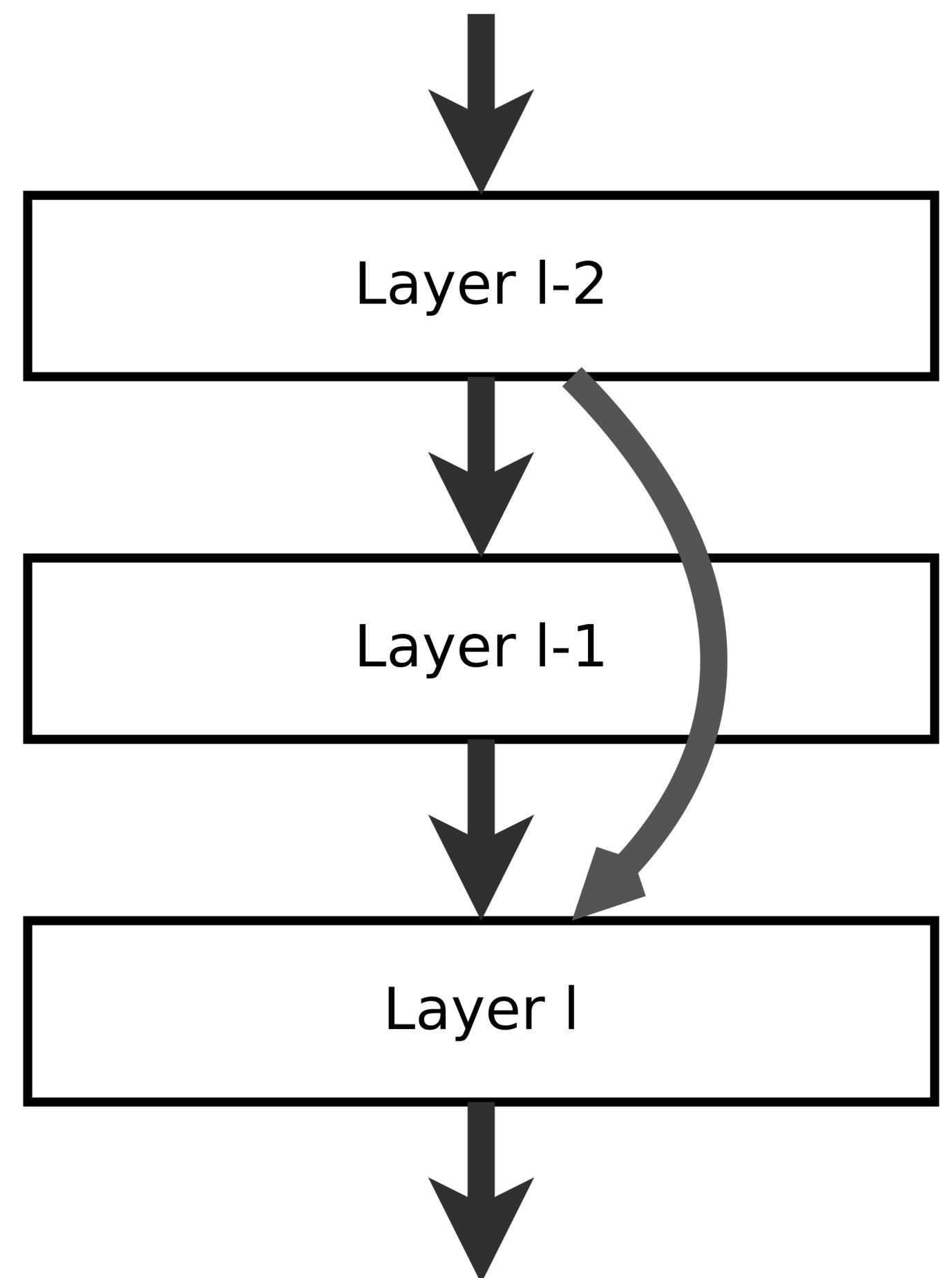
Outline

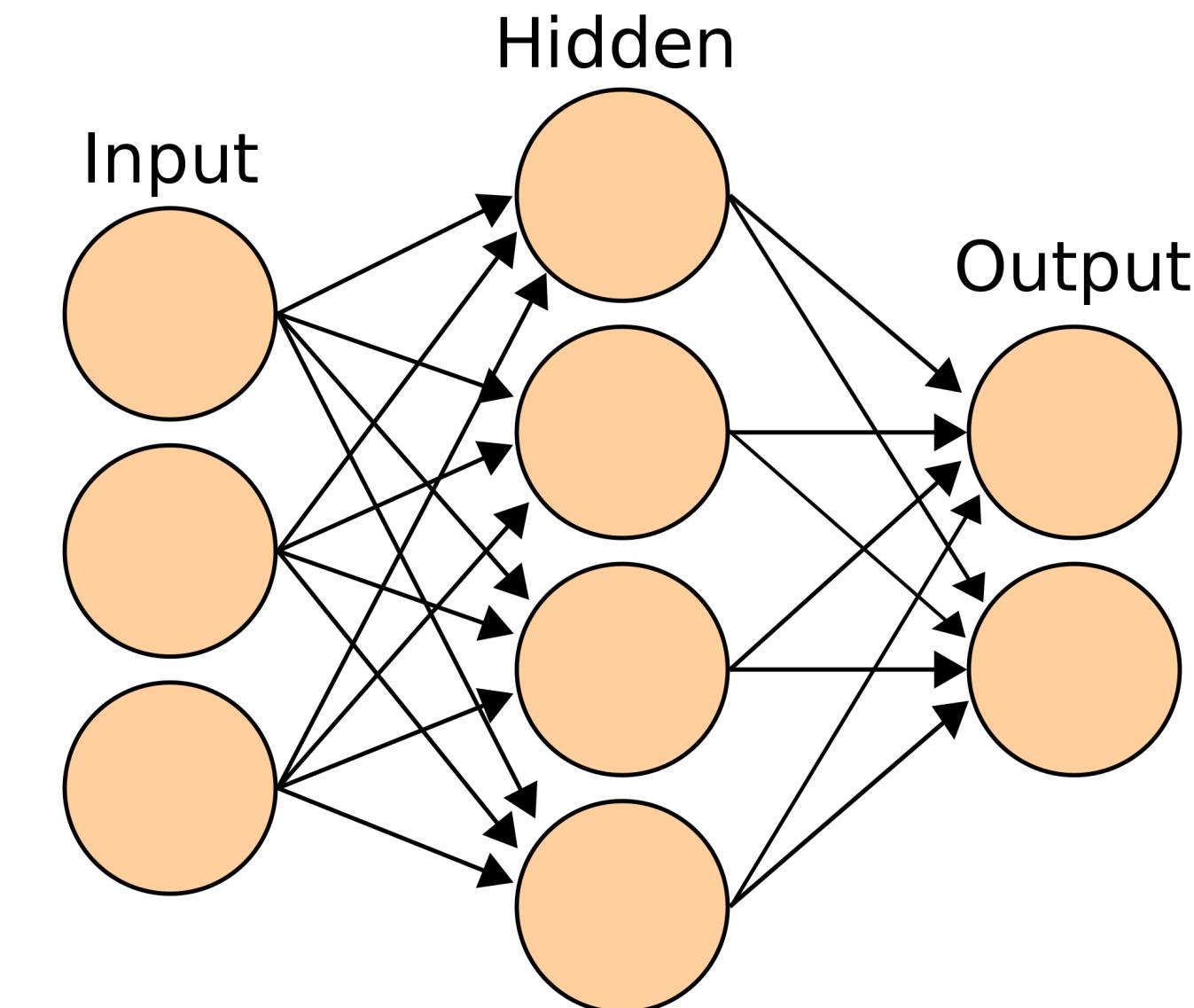
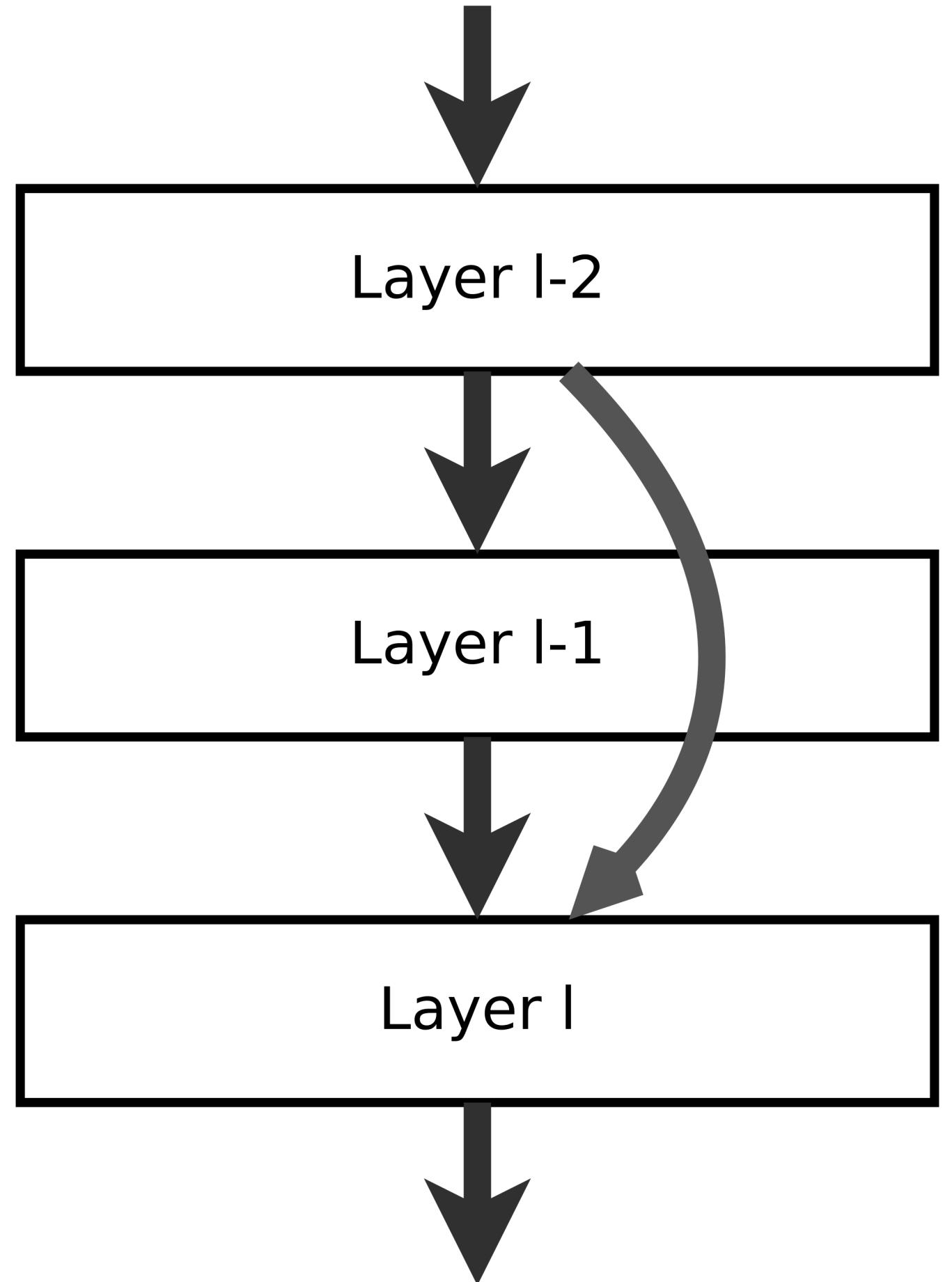
1. Single layer and putting layers together
2. How to train models
3. Deeper is better
4. Practical tips

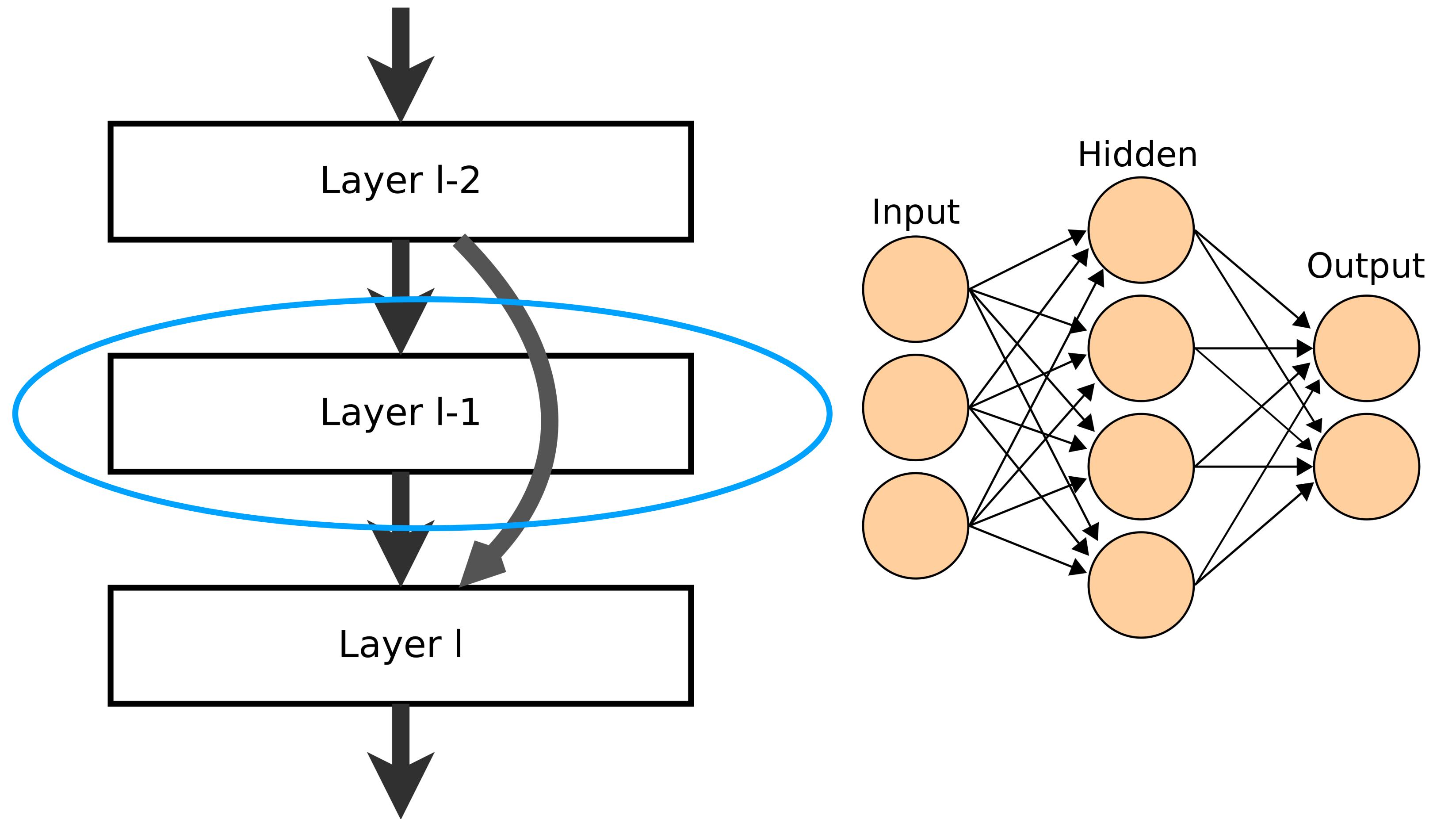
**Understanding is not
something you can
Google.**

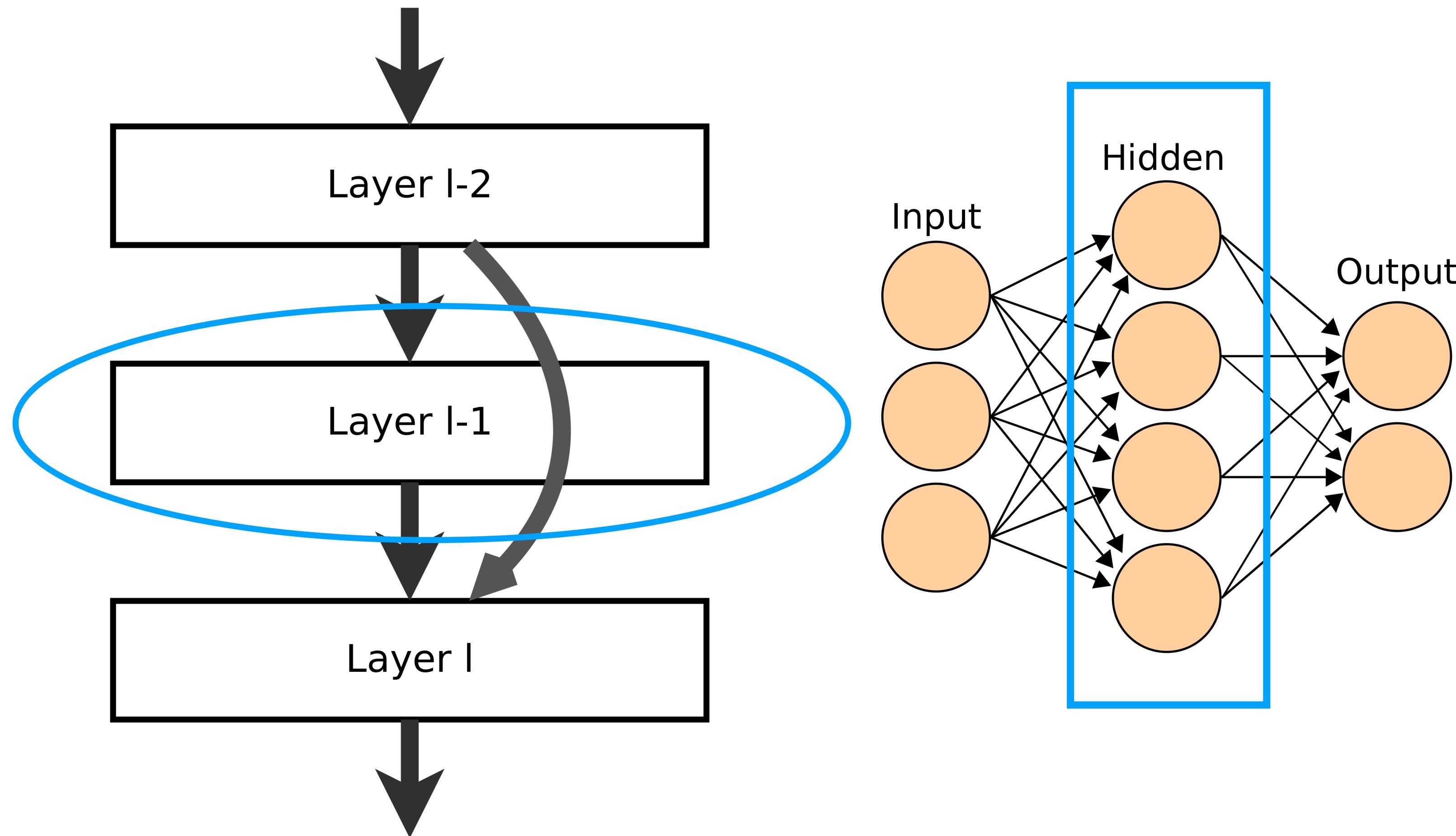
**Machine learning is
engineering-driven.**

I. Single Layer

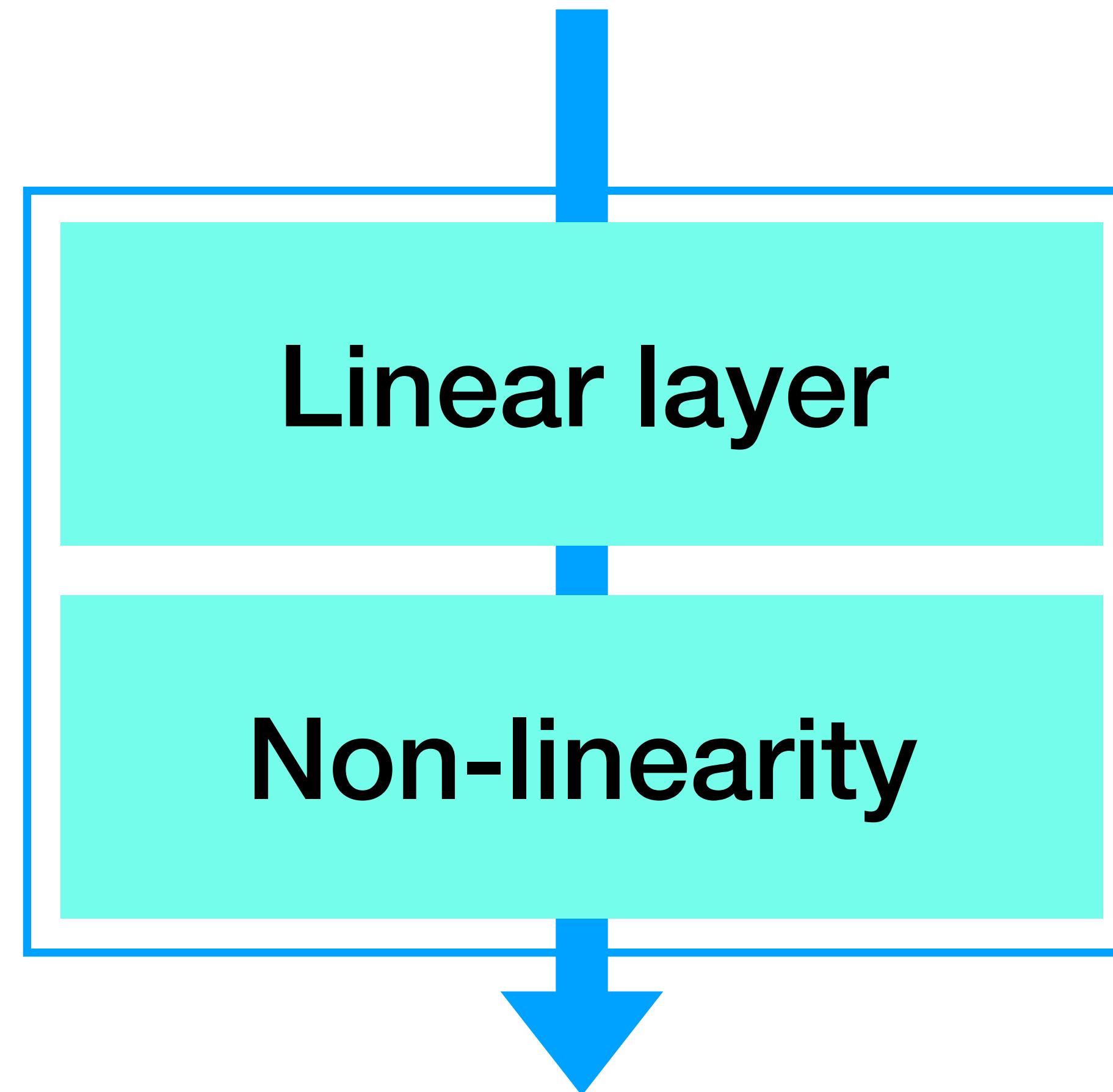






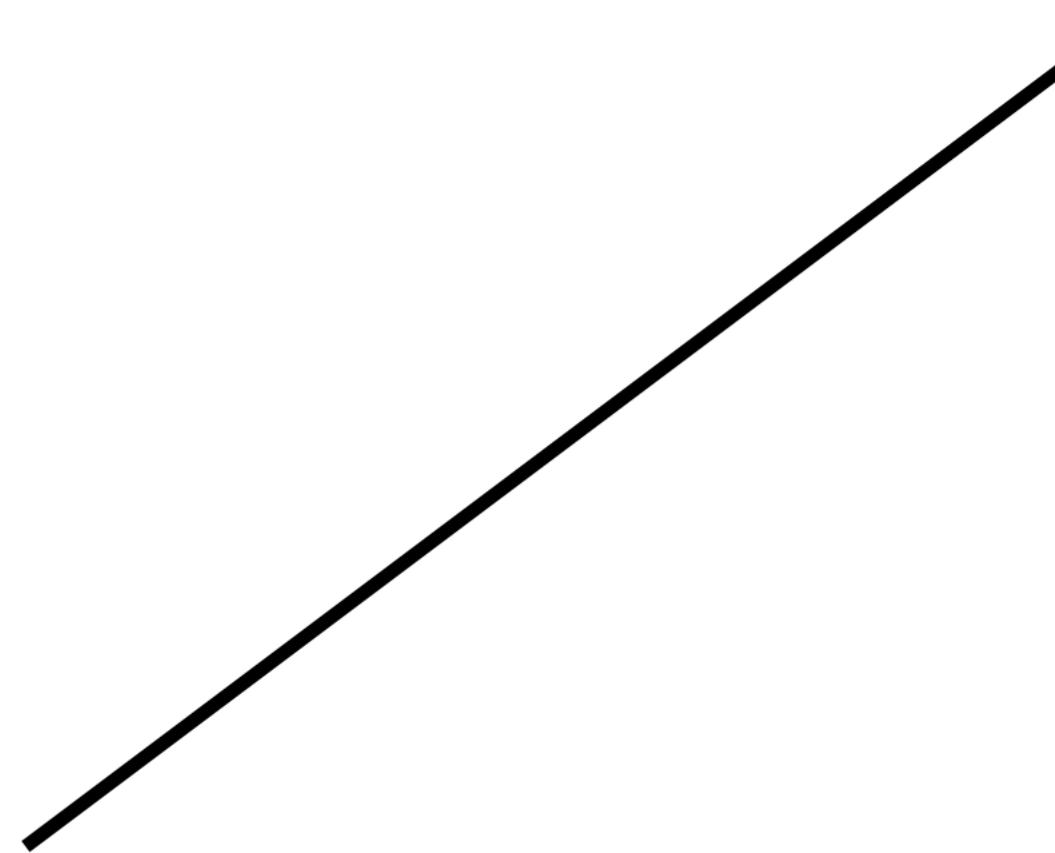


Fully connected layer

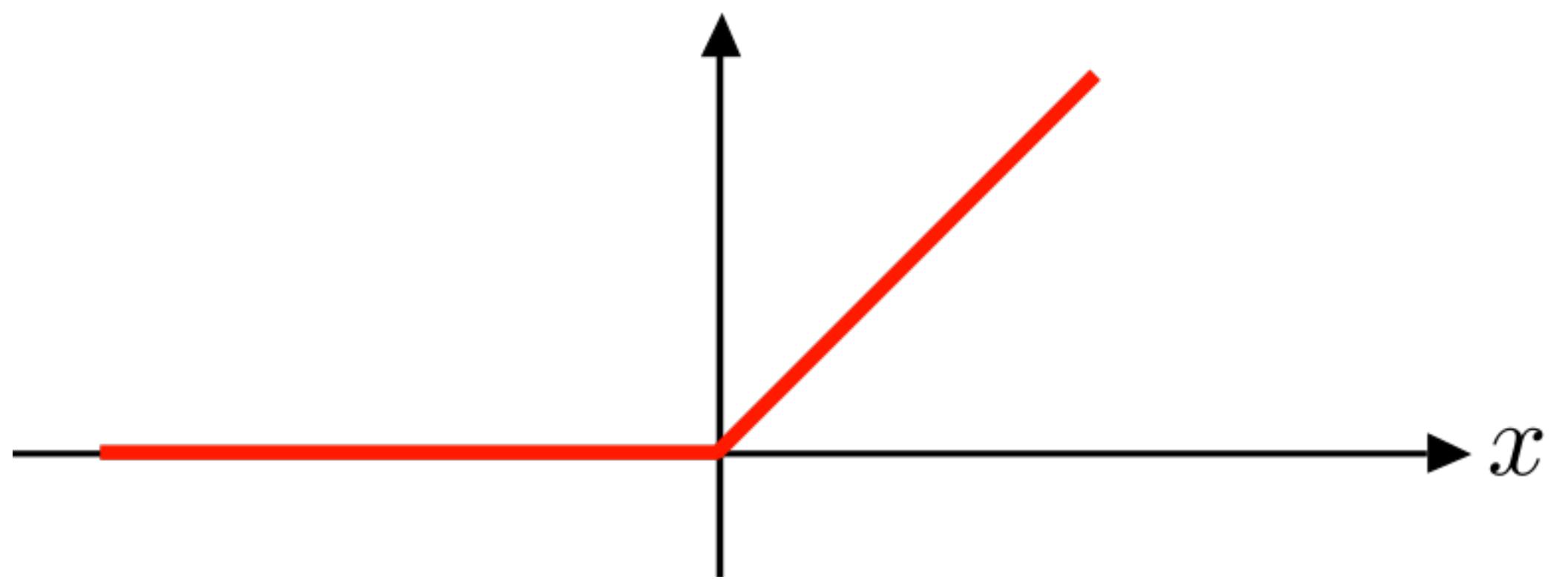


A single layer

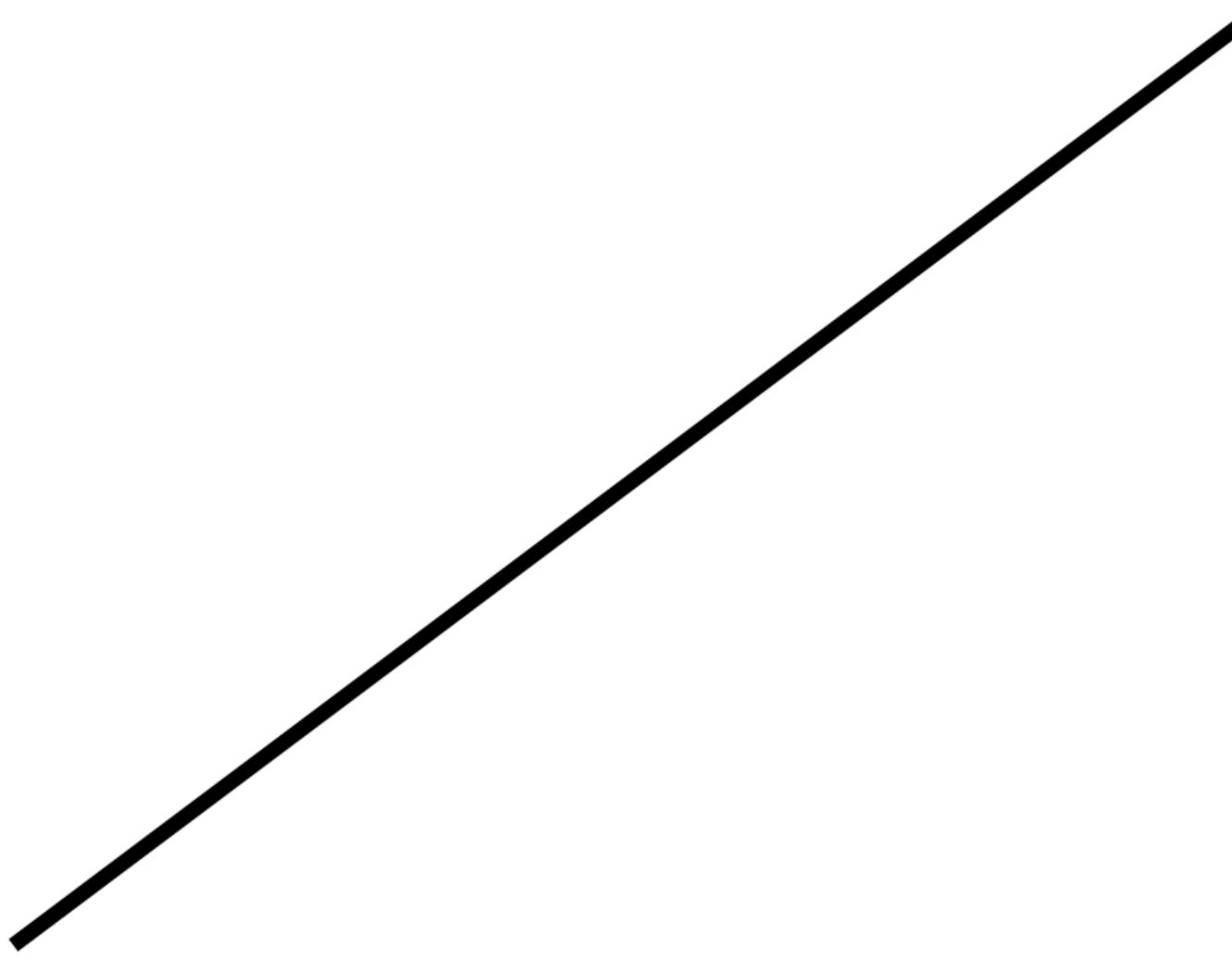
Linear layer



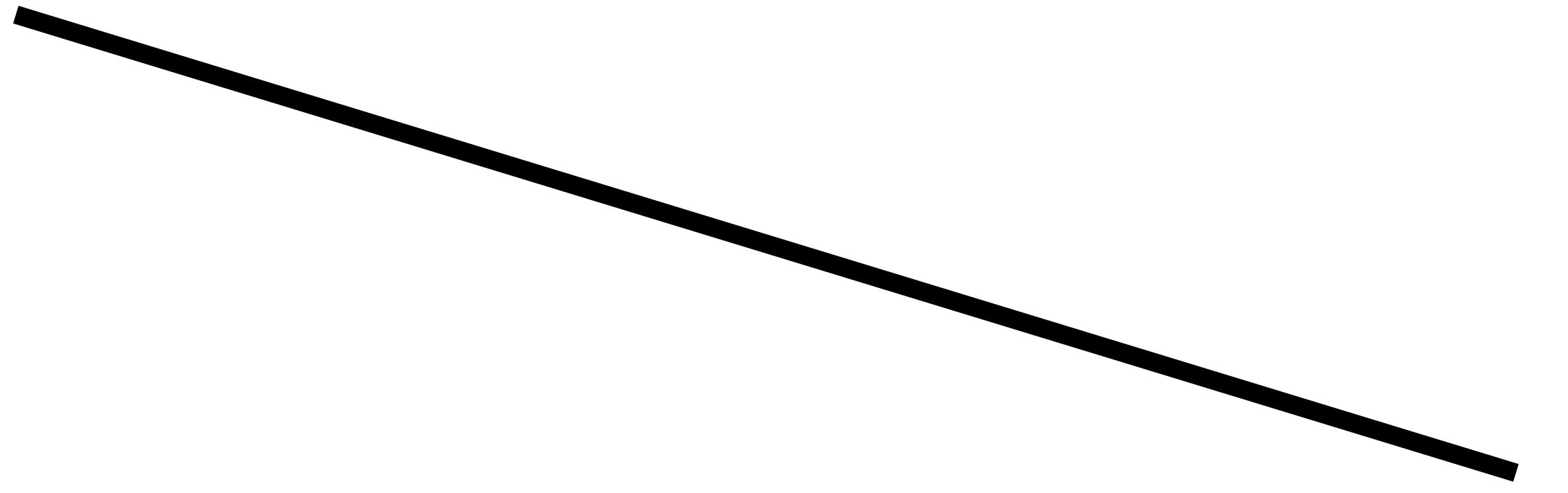
Non-linearity

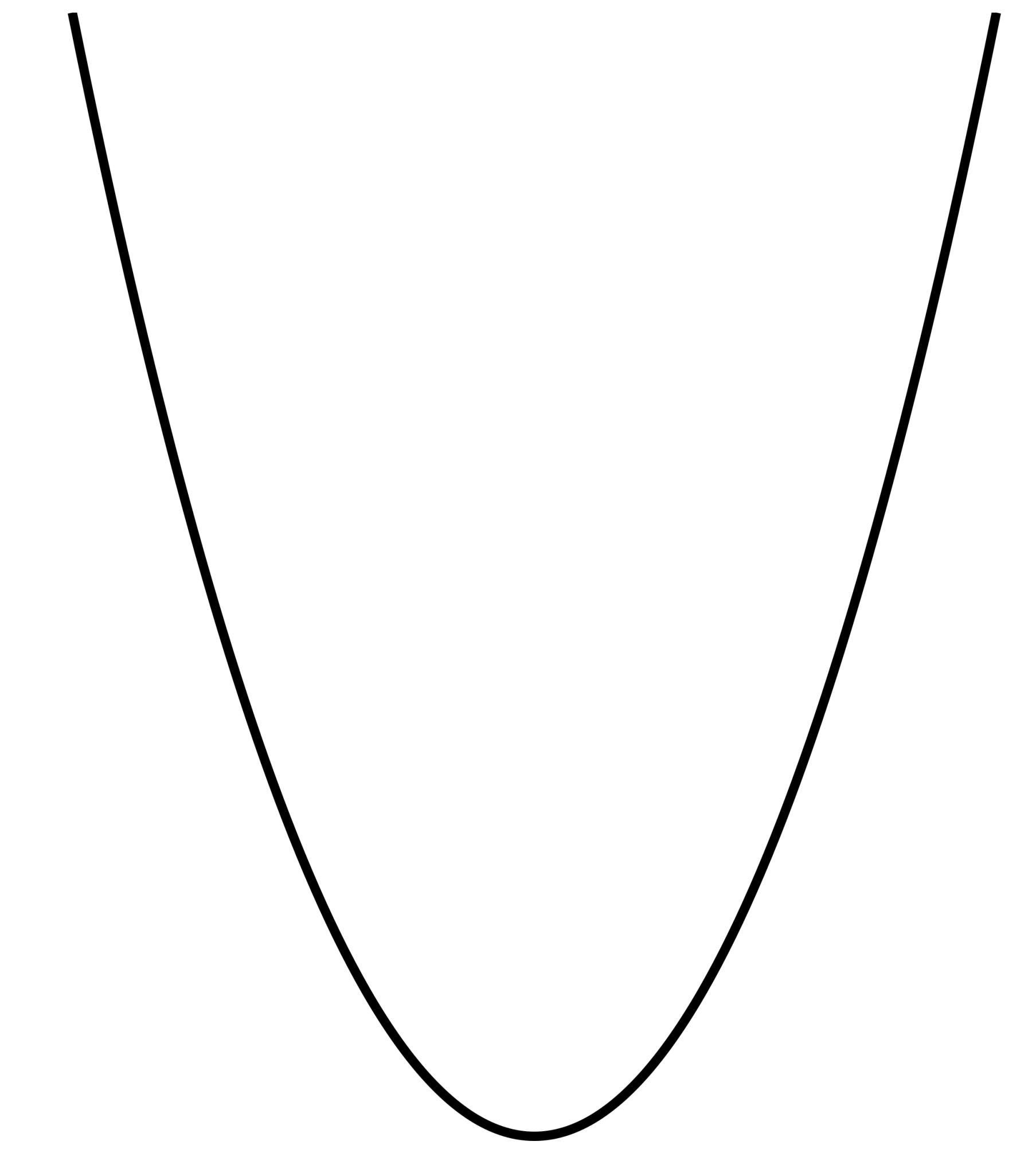


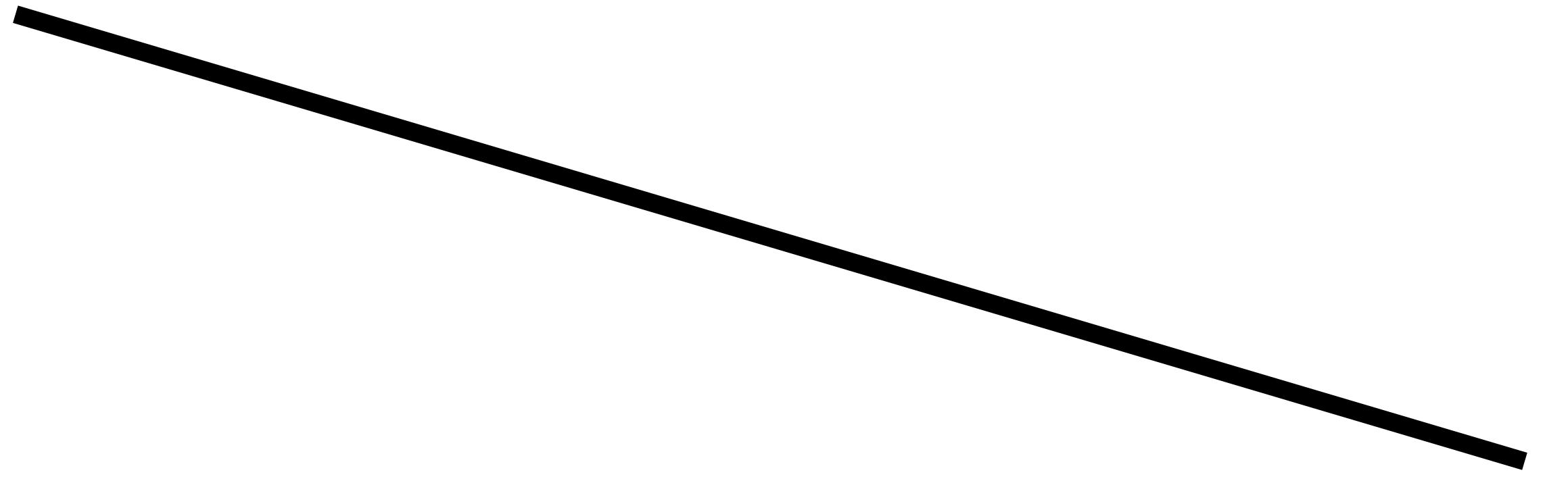
Linear layer

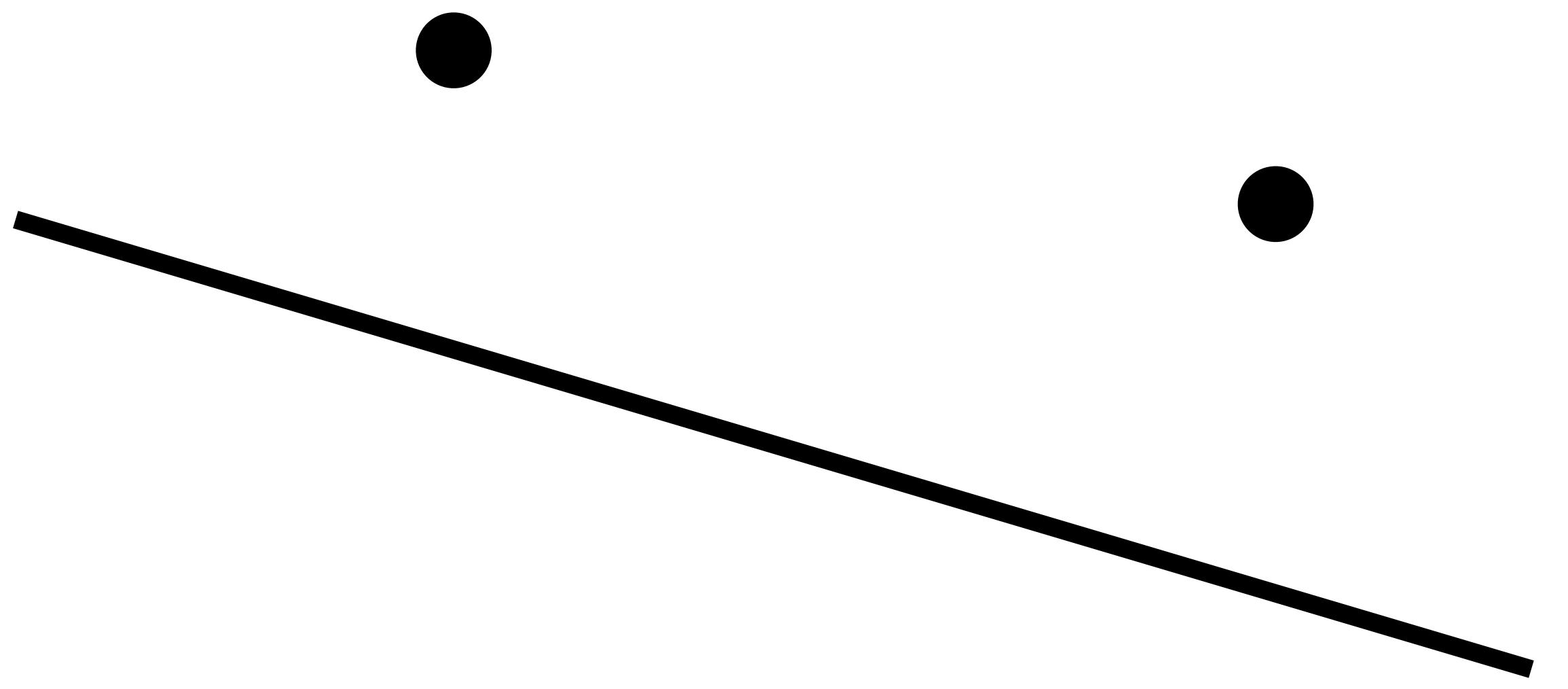


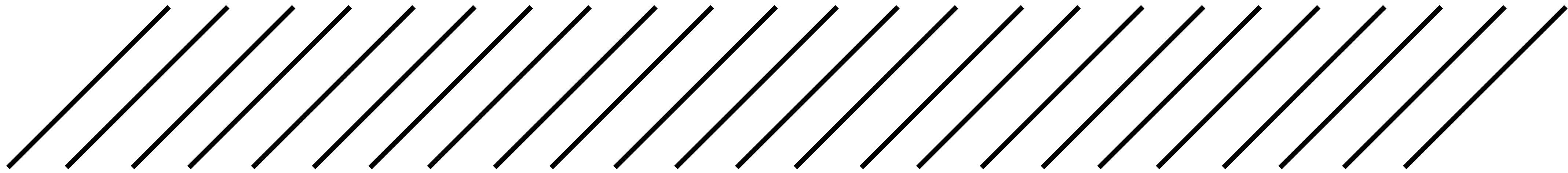
$$w_x + b$$

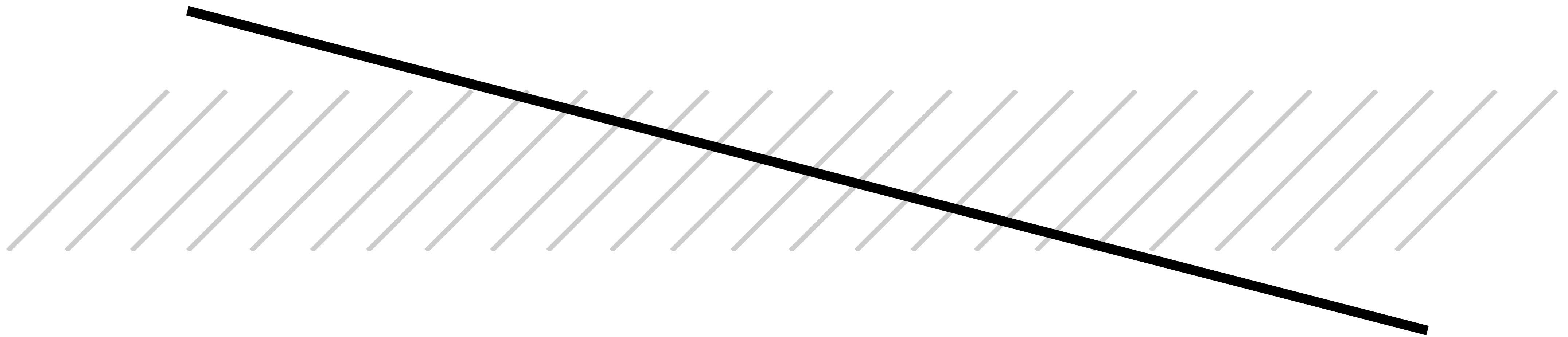


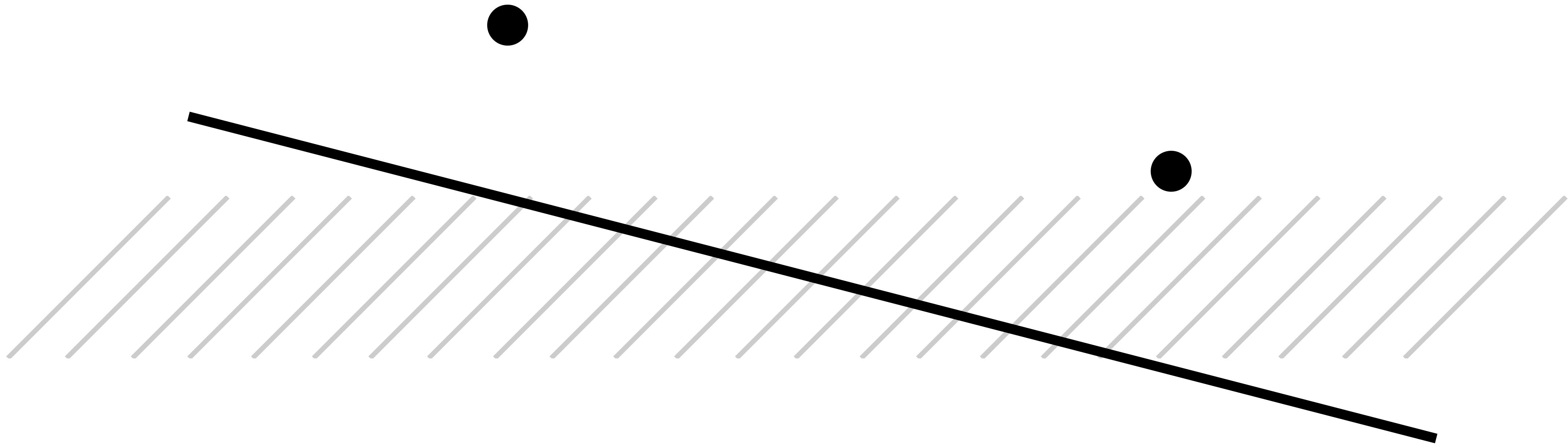


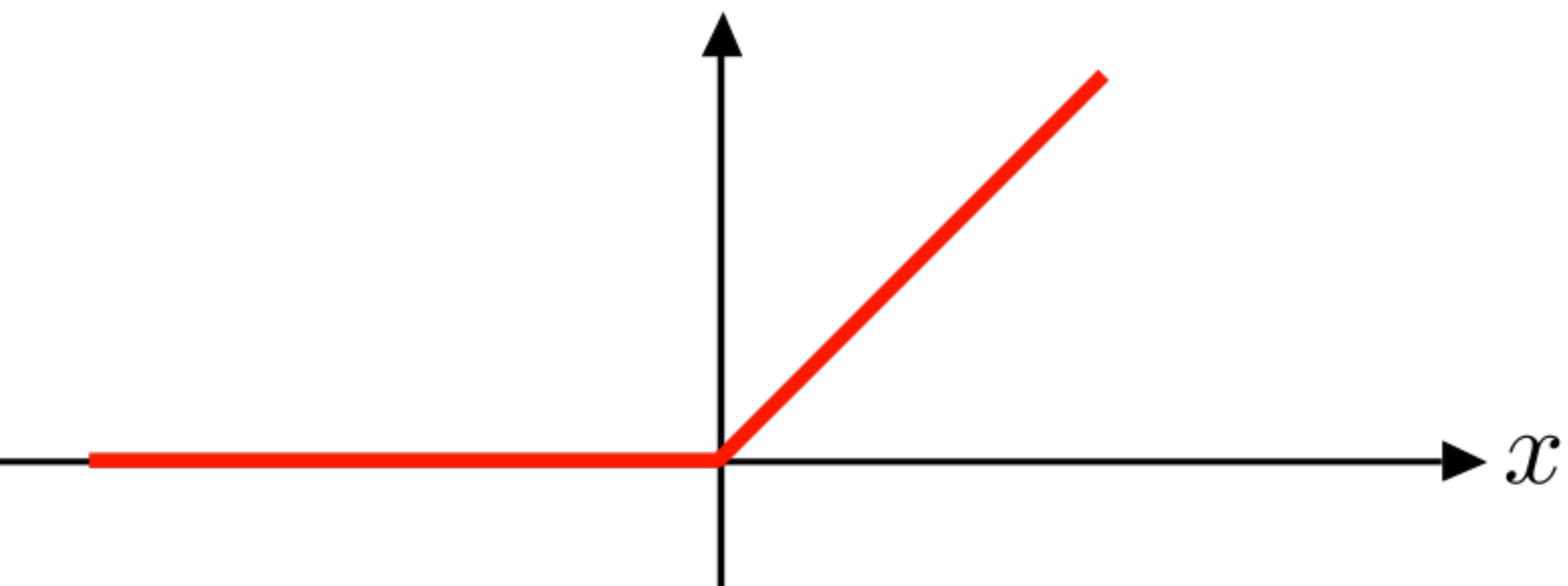


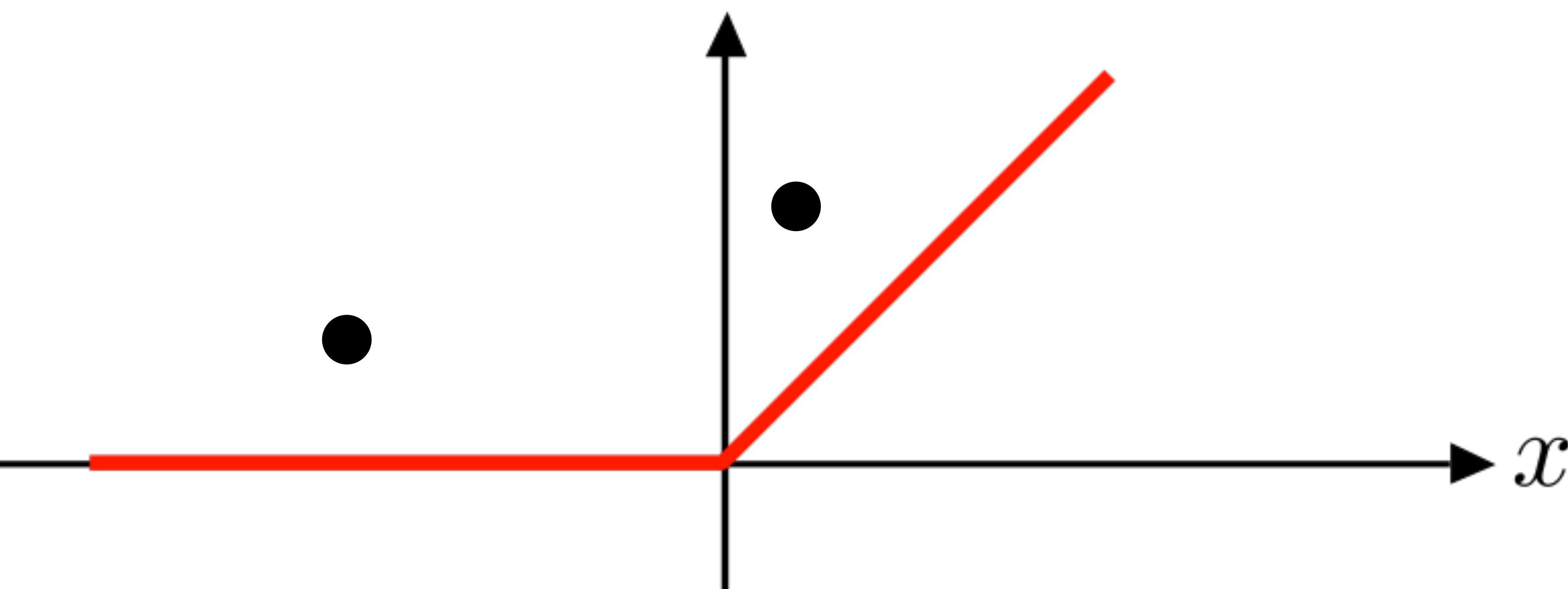






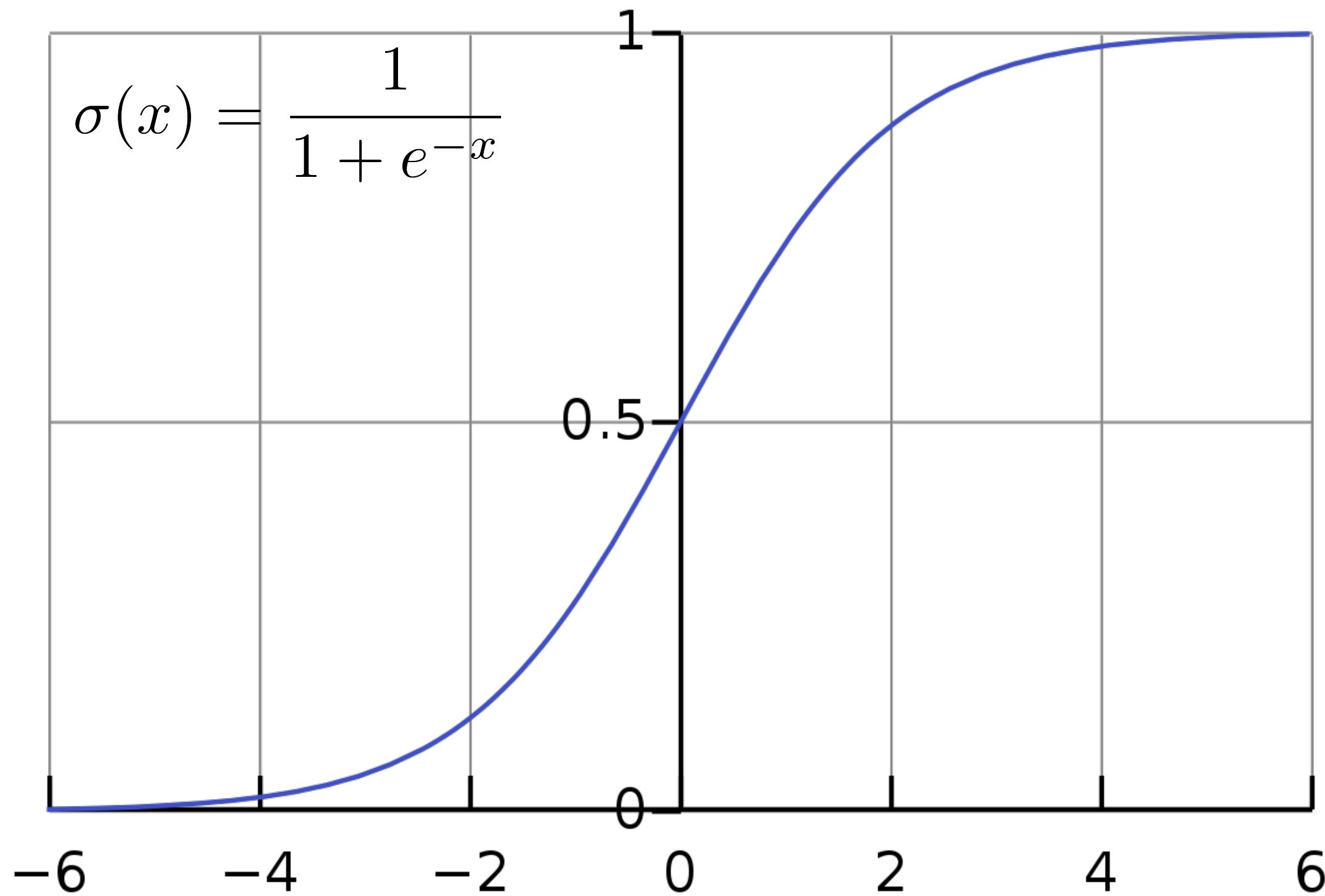




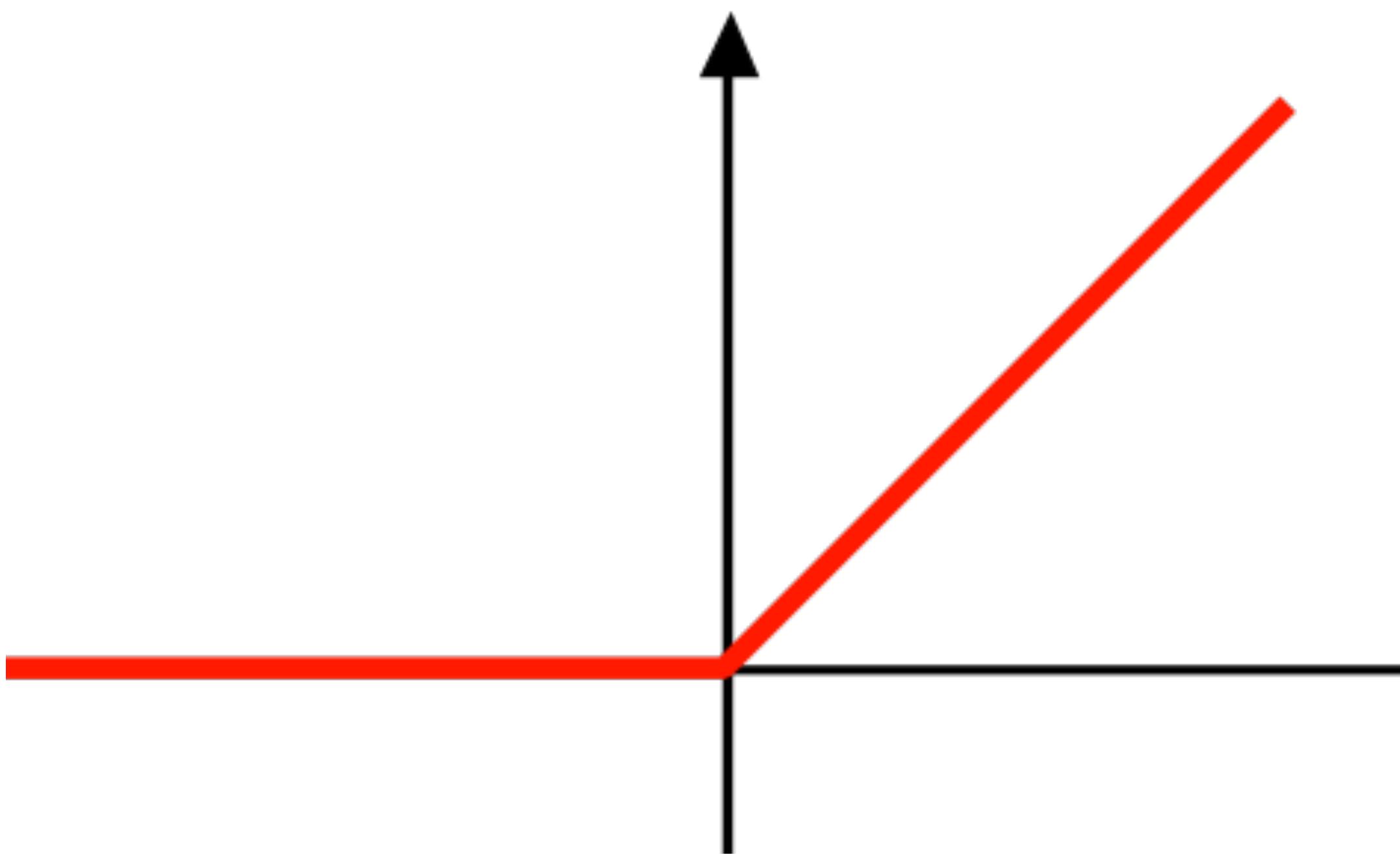


Non-linearities

Sigmoid

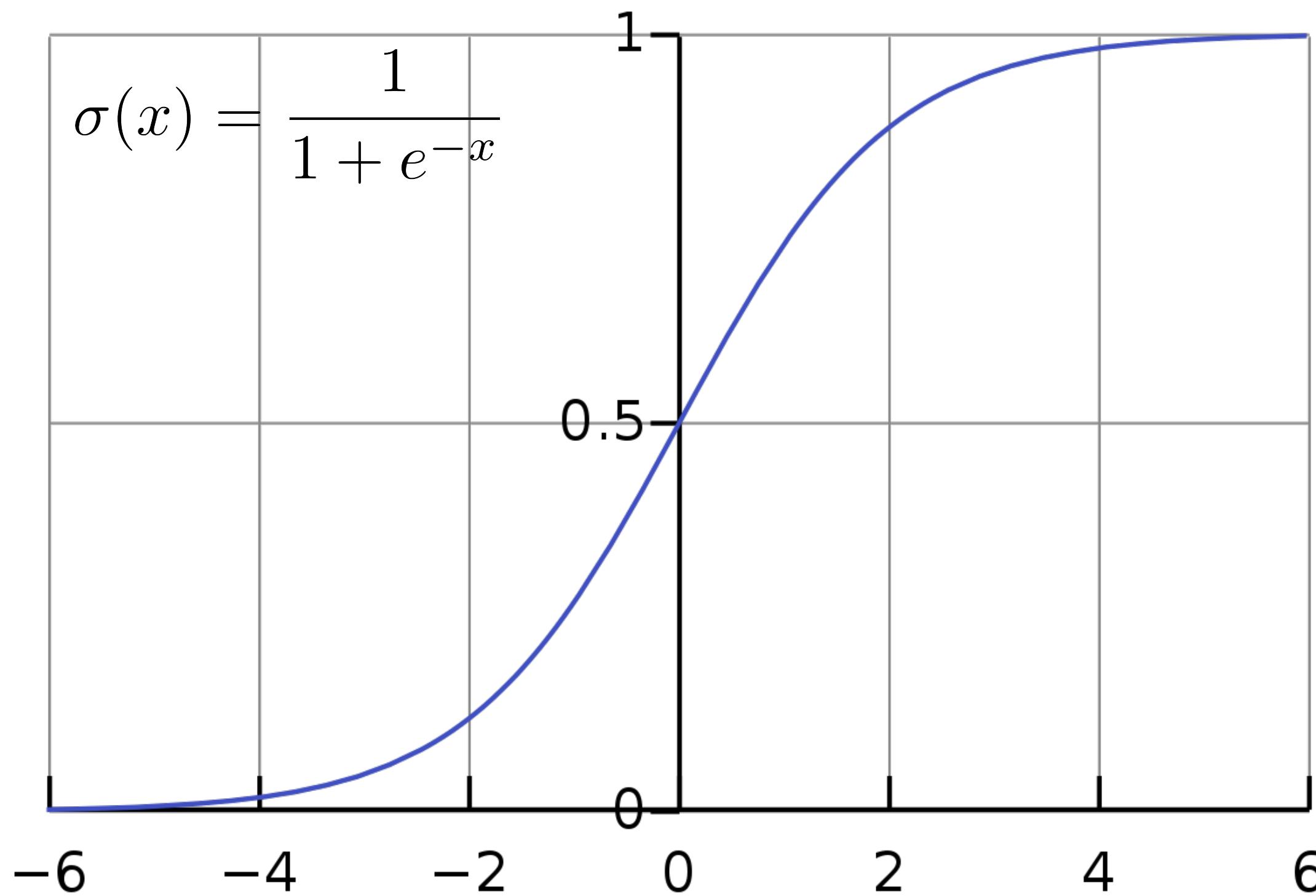


ReLU

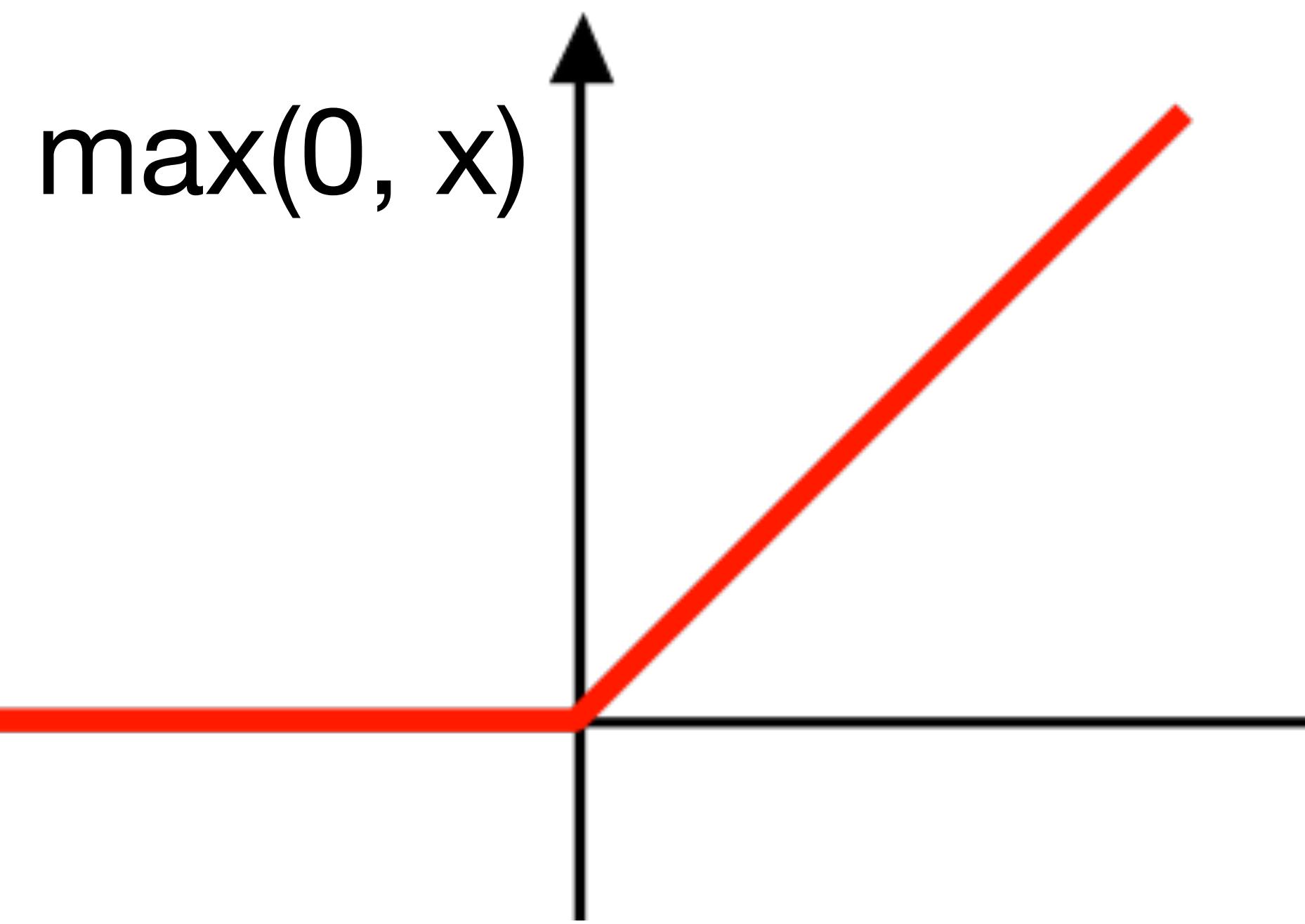


Non-linearities

Sigmoid



ReLU

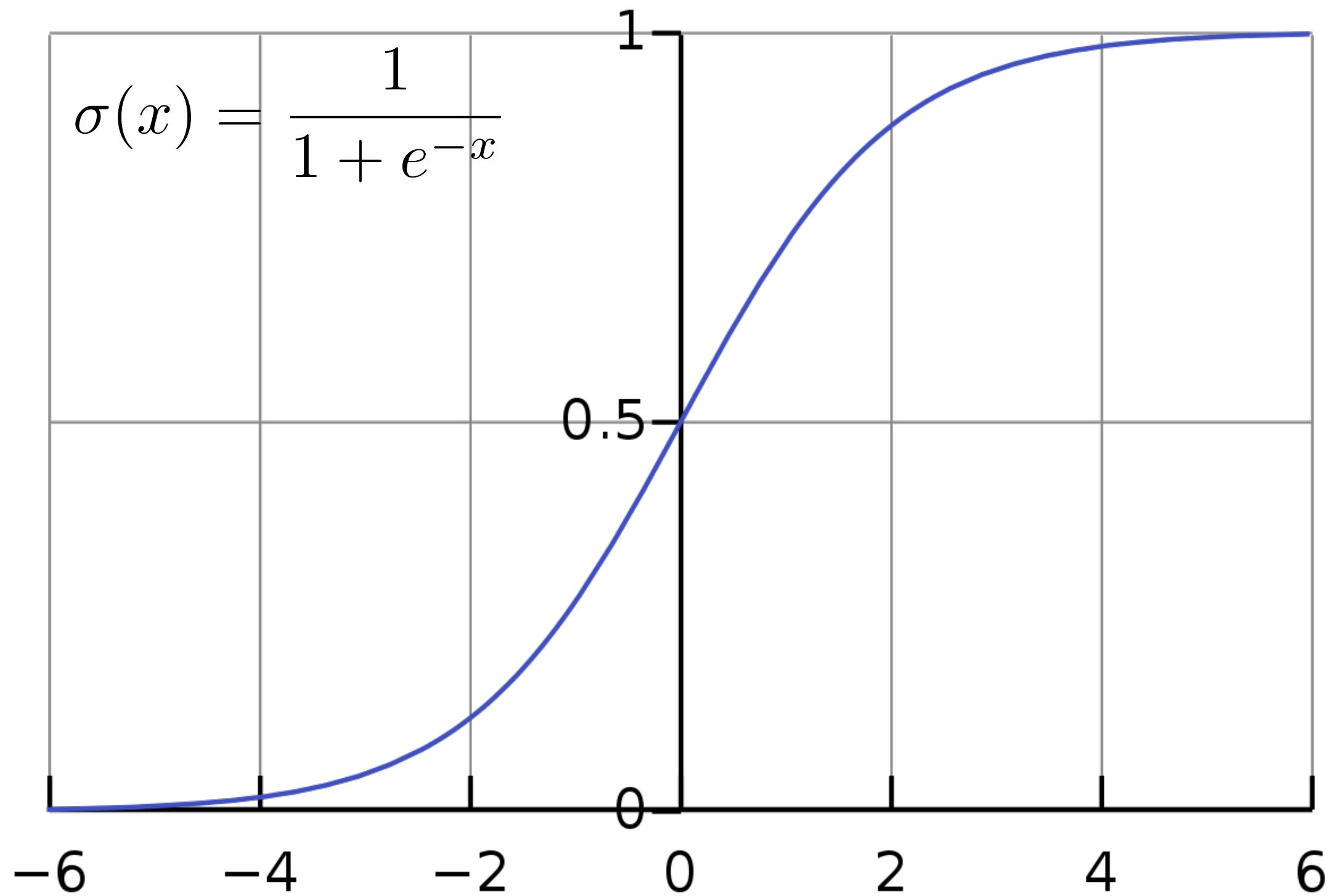


Output layer

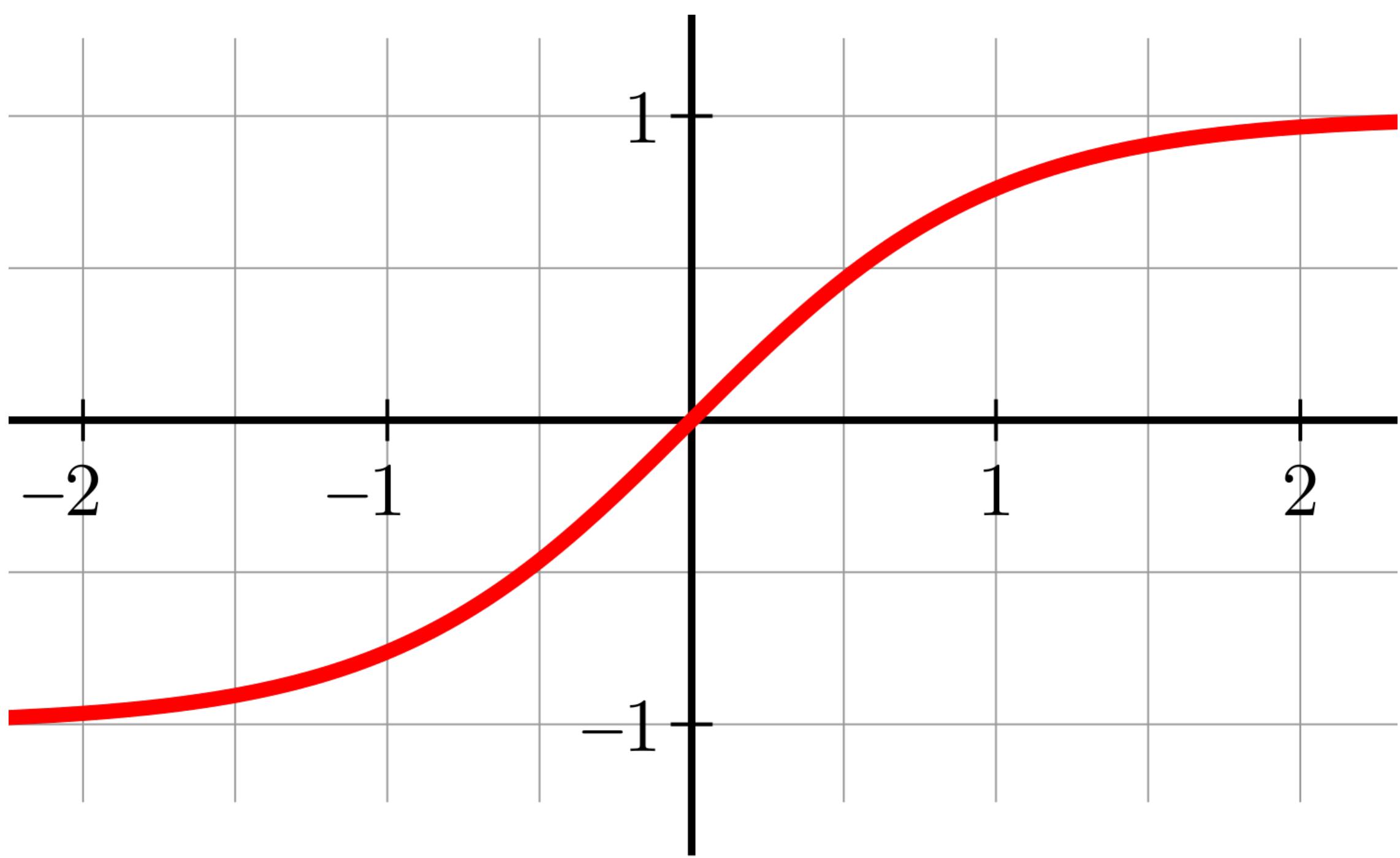




Sigmoid

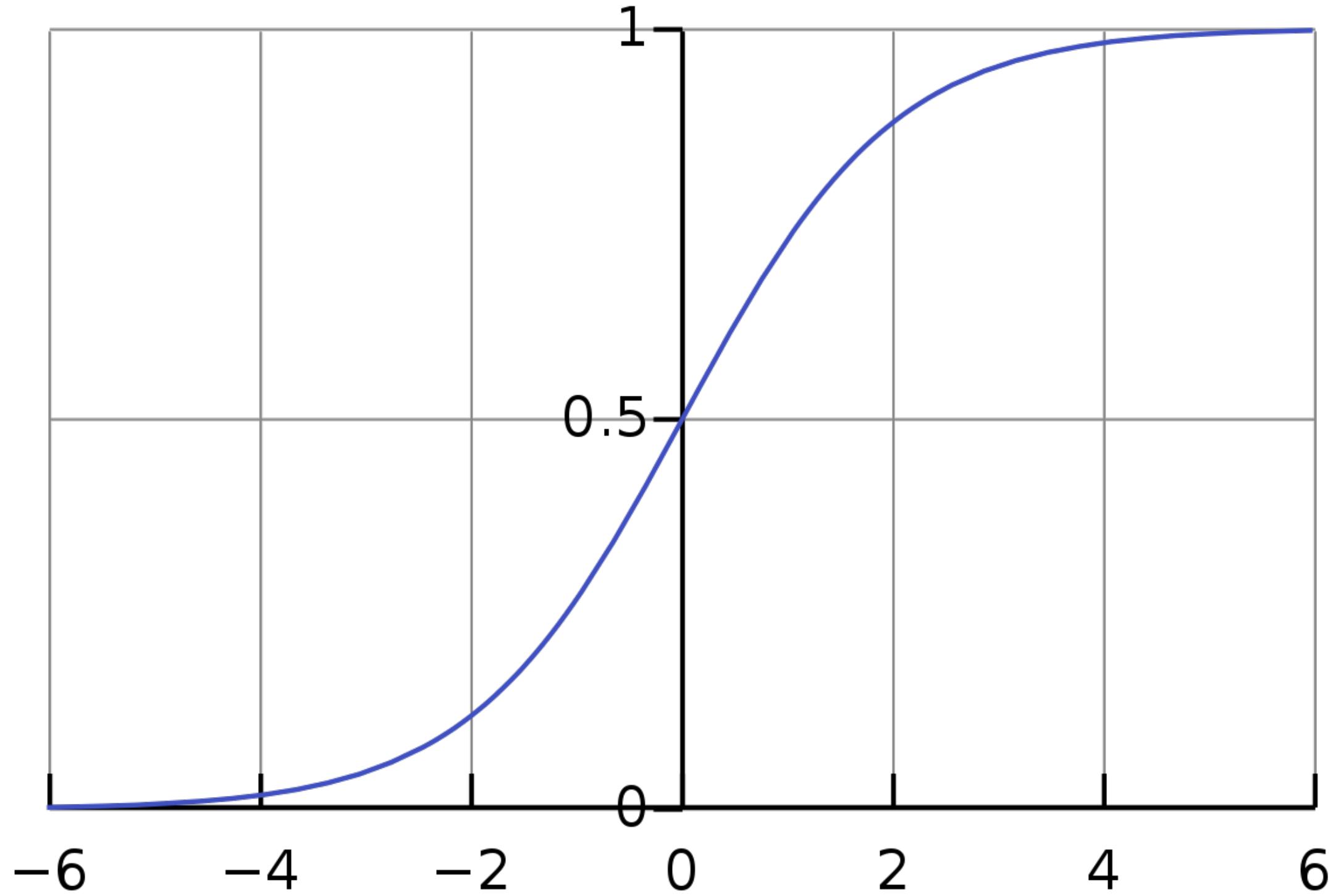


tanh

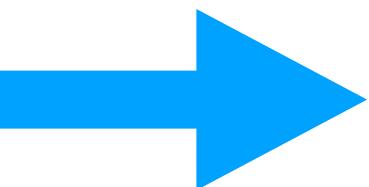




Sigmoid



Linear layer

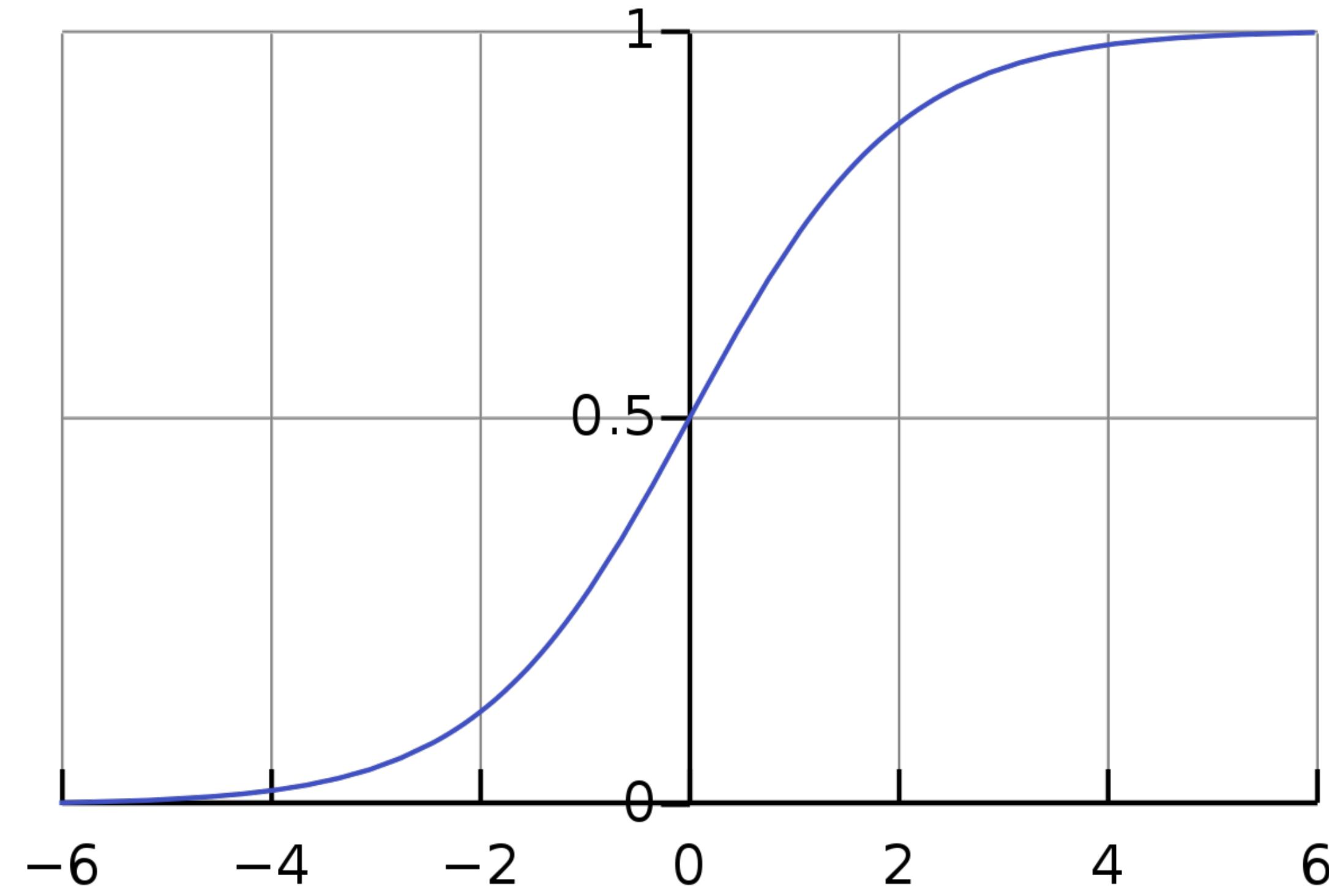


$$\mathbf{Wx} + \mathbf{b}$$

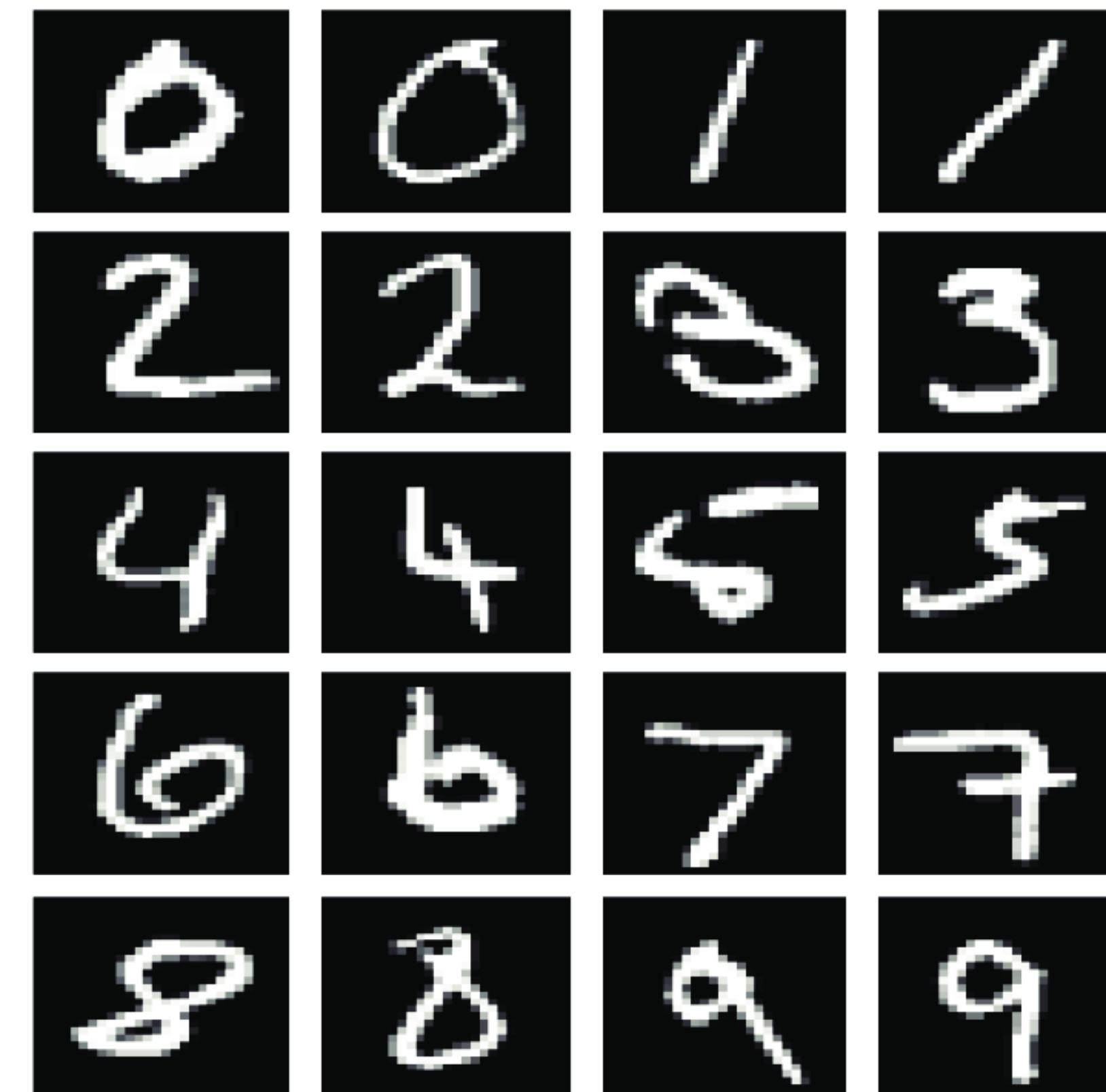


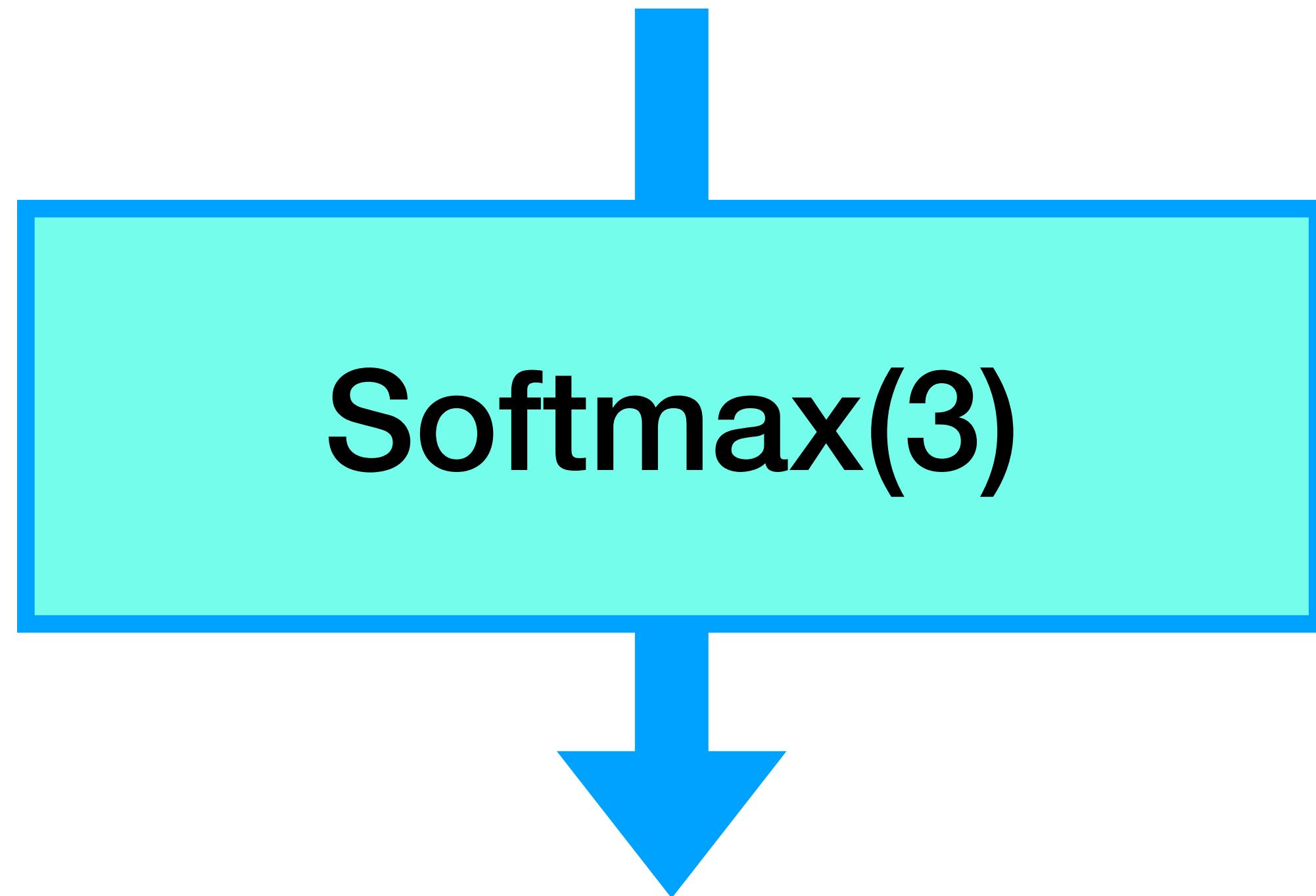


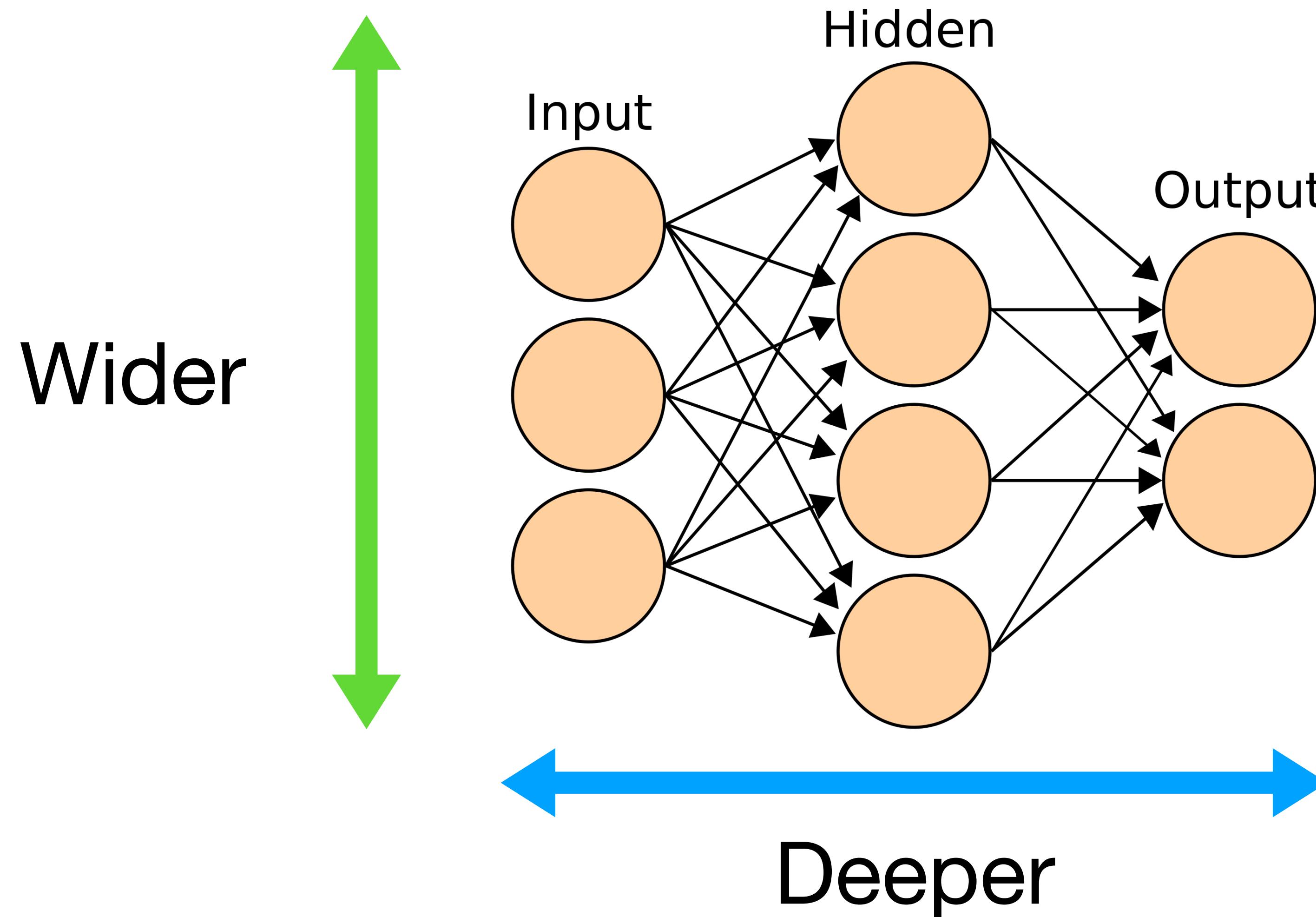
Sigmoid

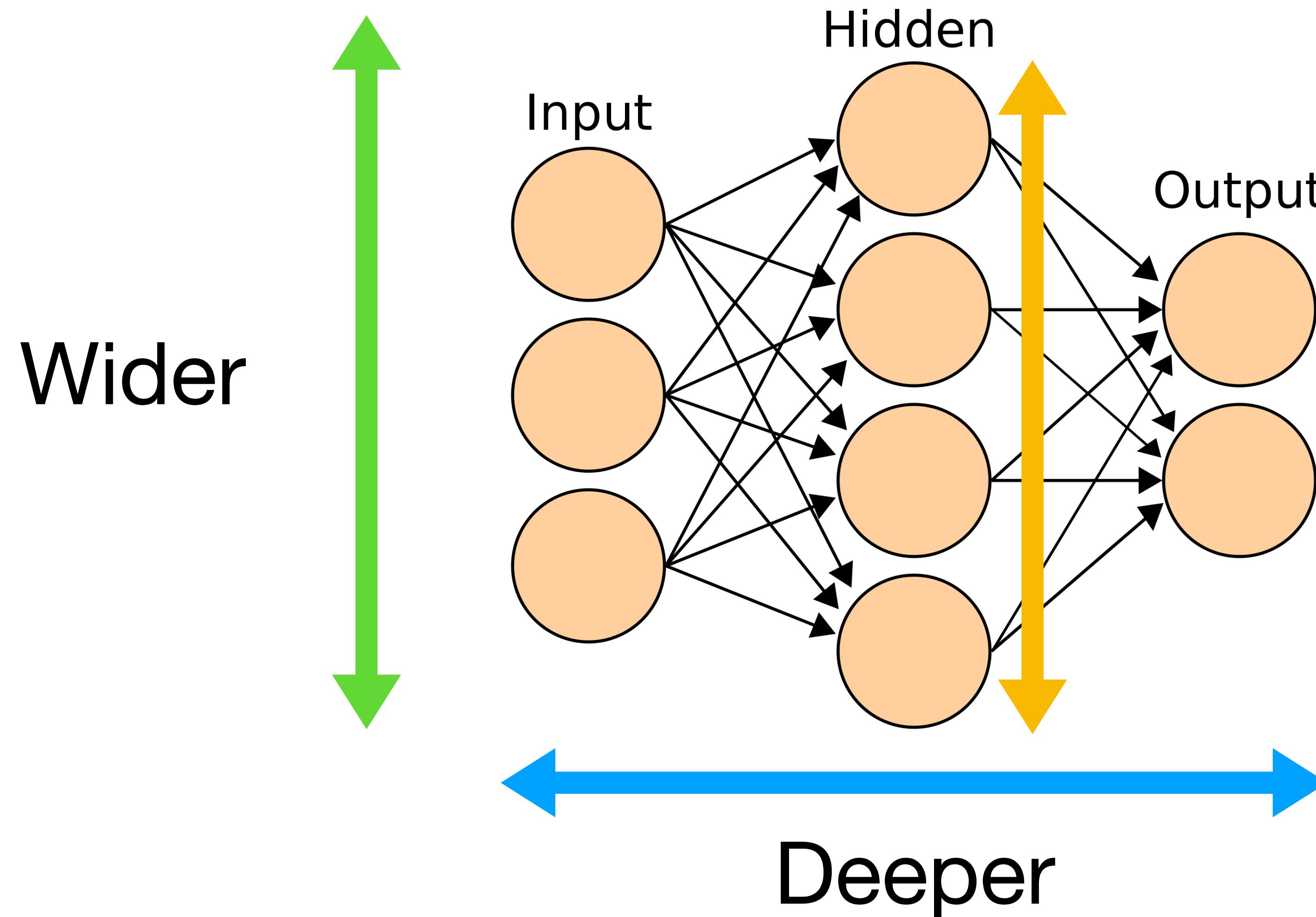


MNIST dataset



$$[-1, 4, 7]$$

$$[0.0003 \ 0.047 \ 0.95]$$

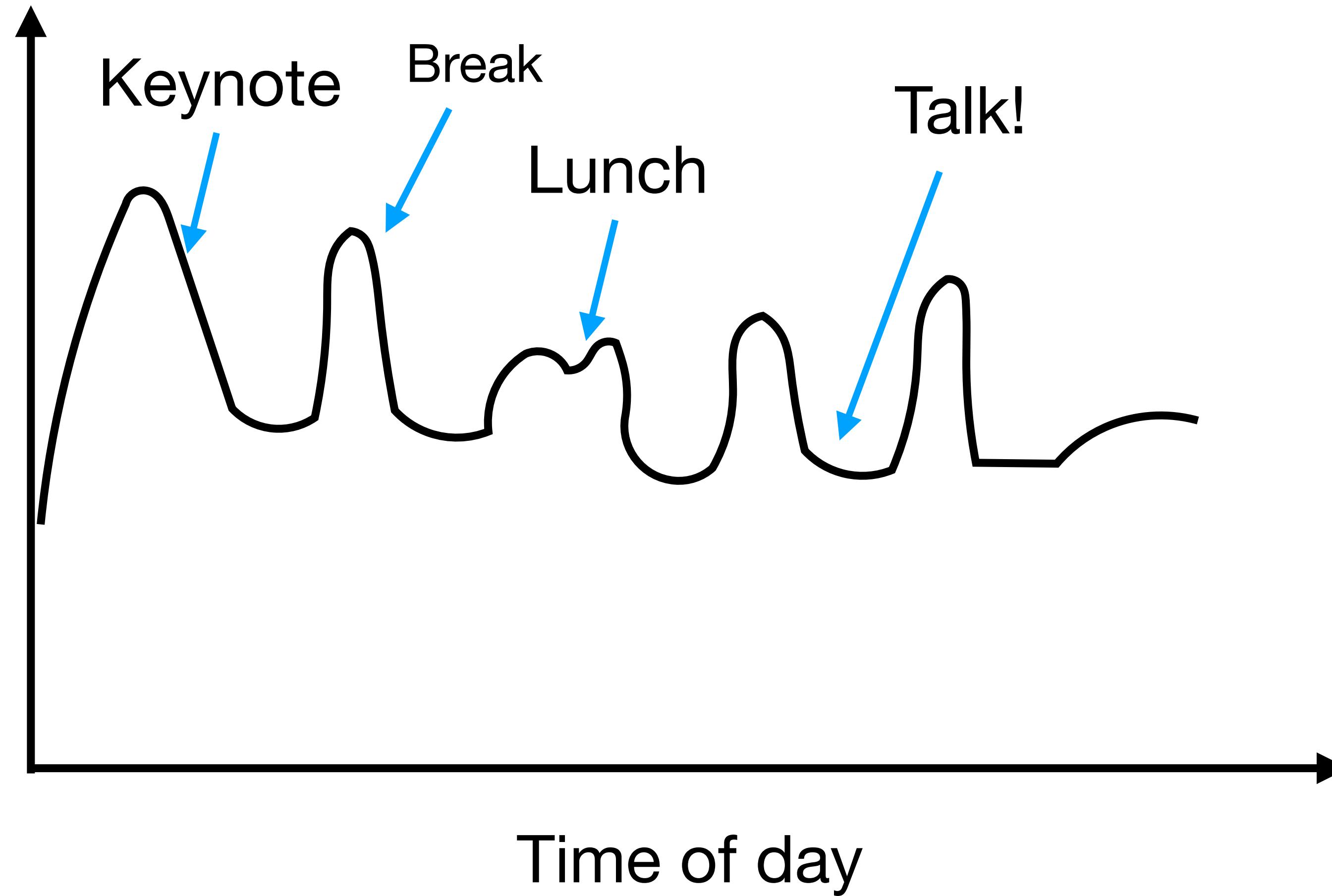


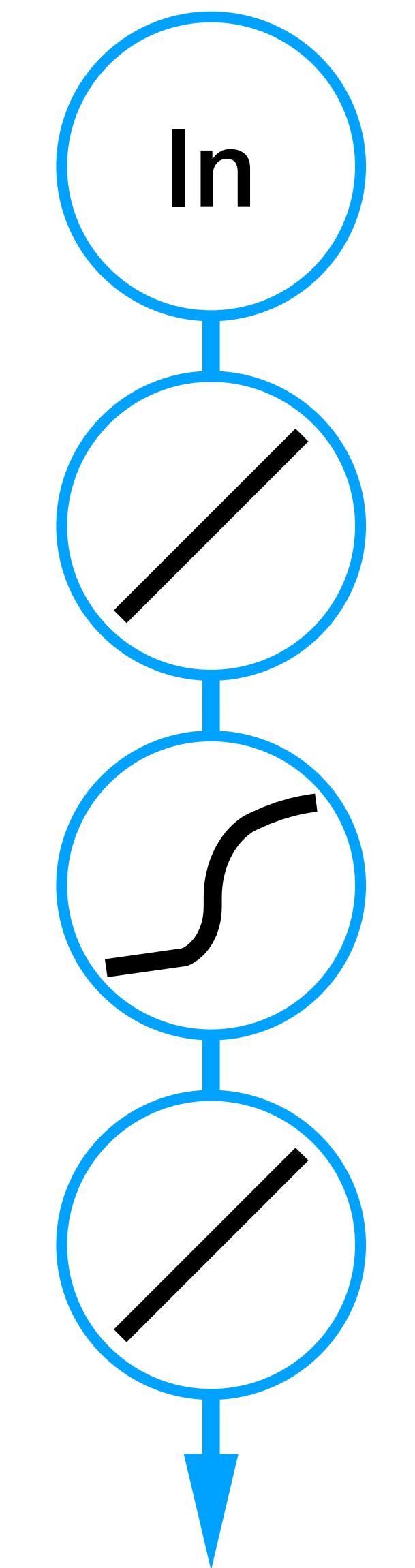


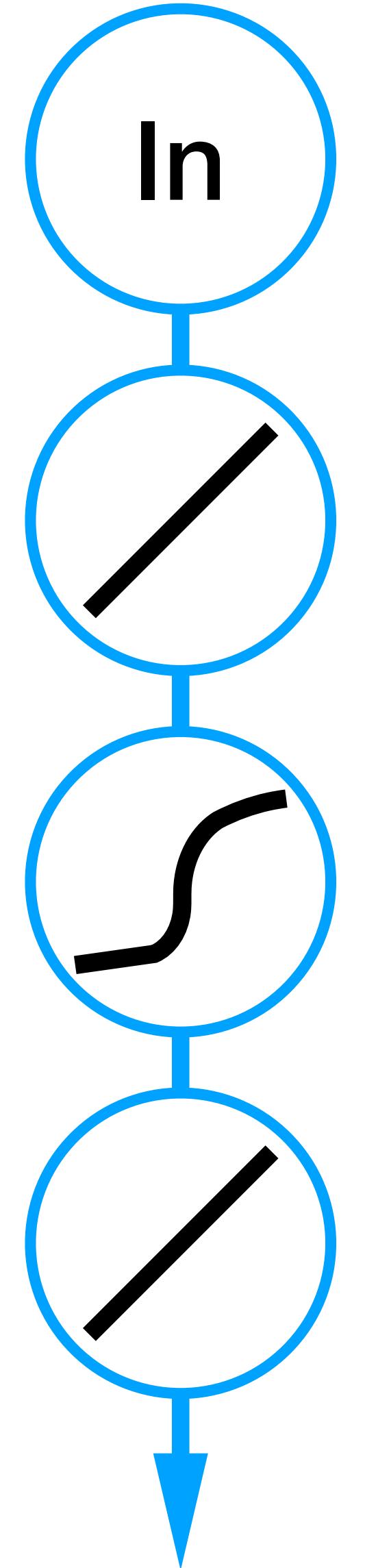




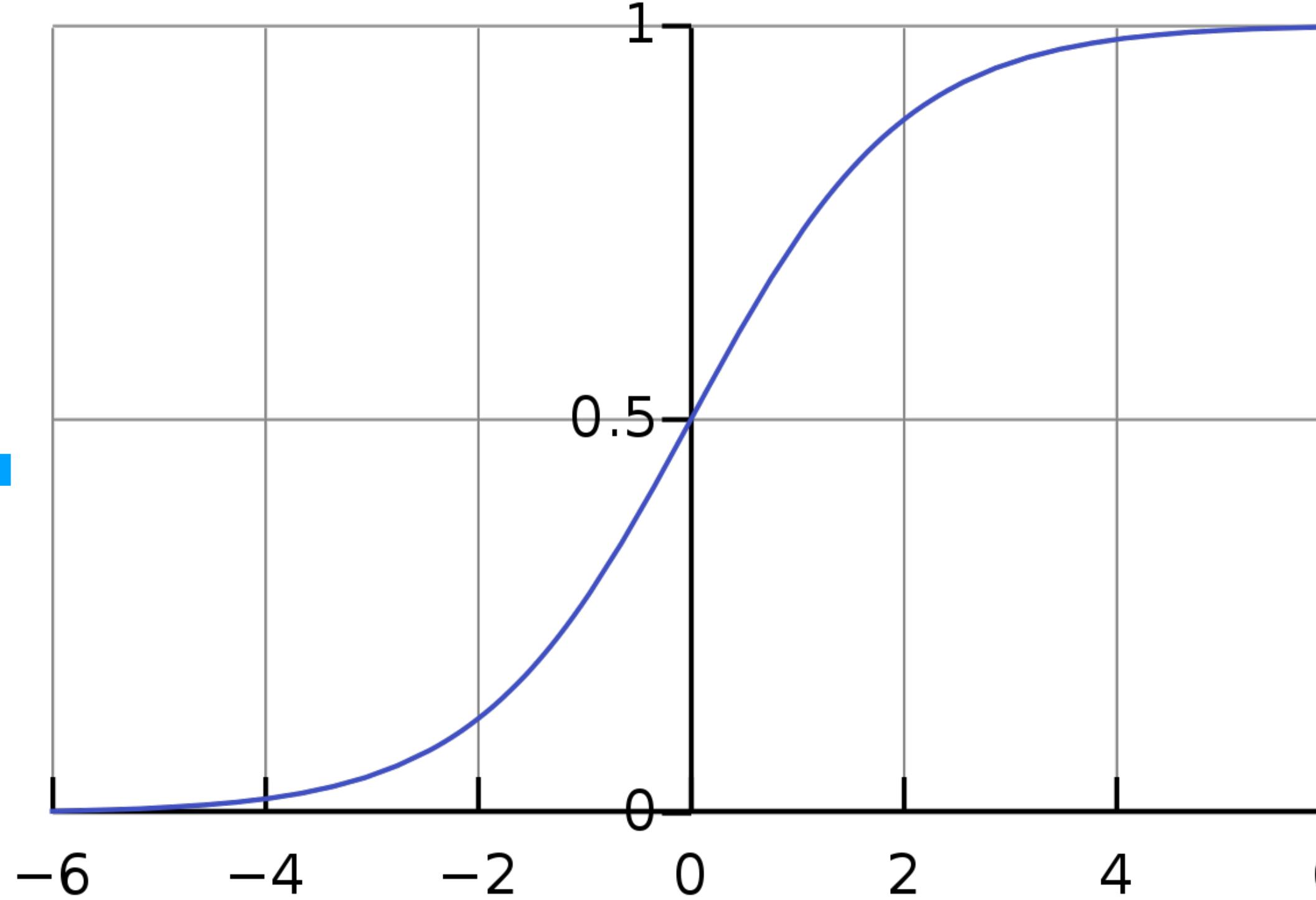
Coffee
needed

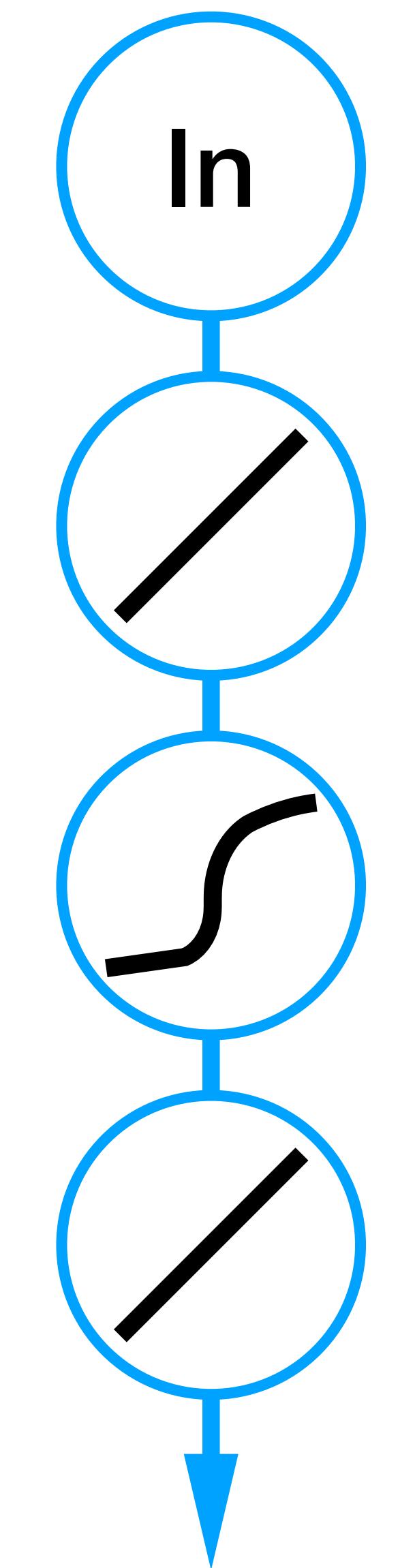


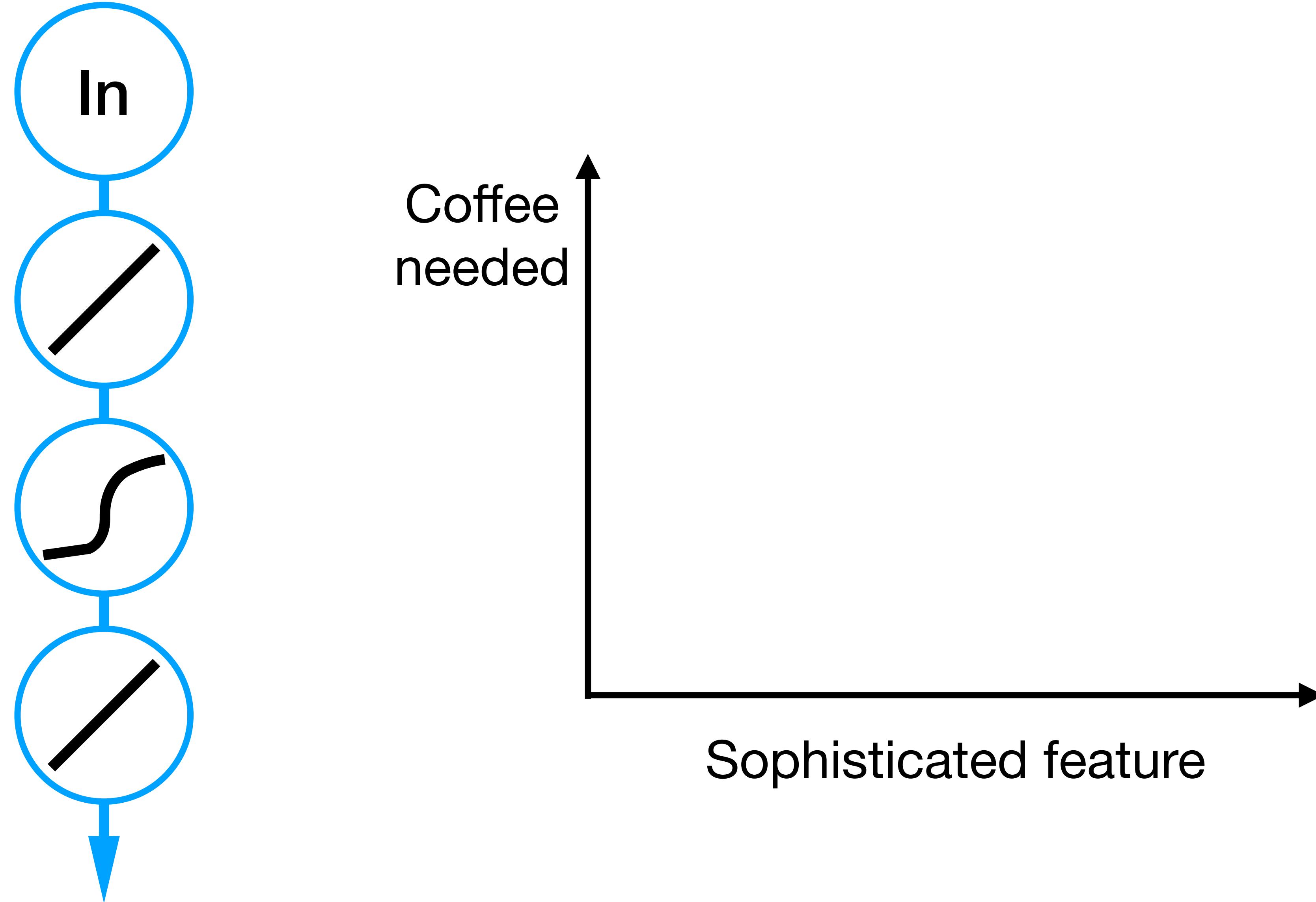


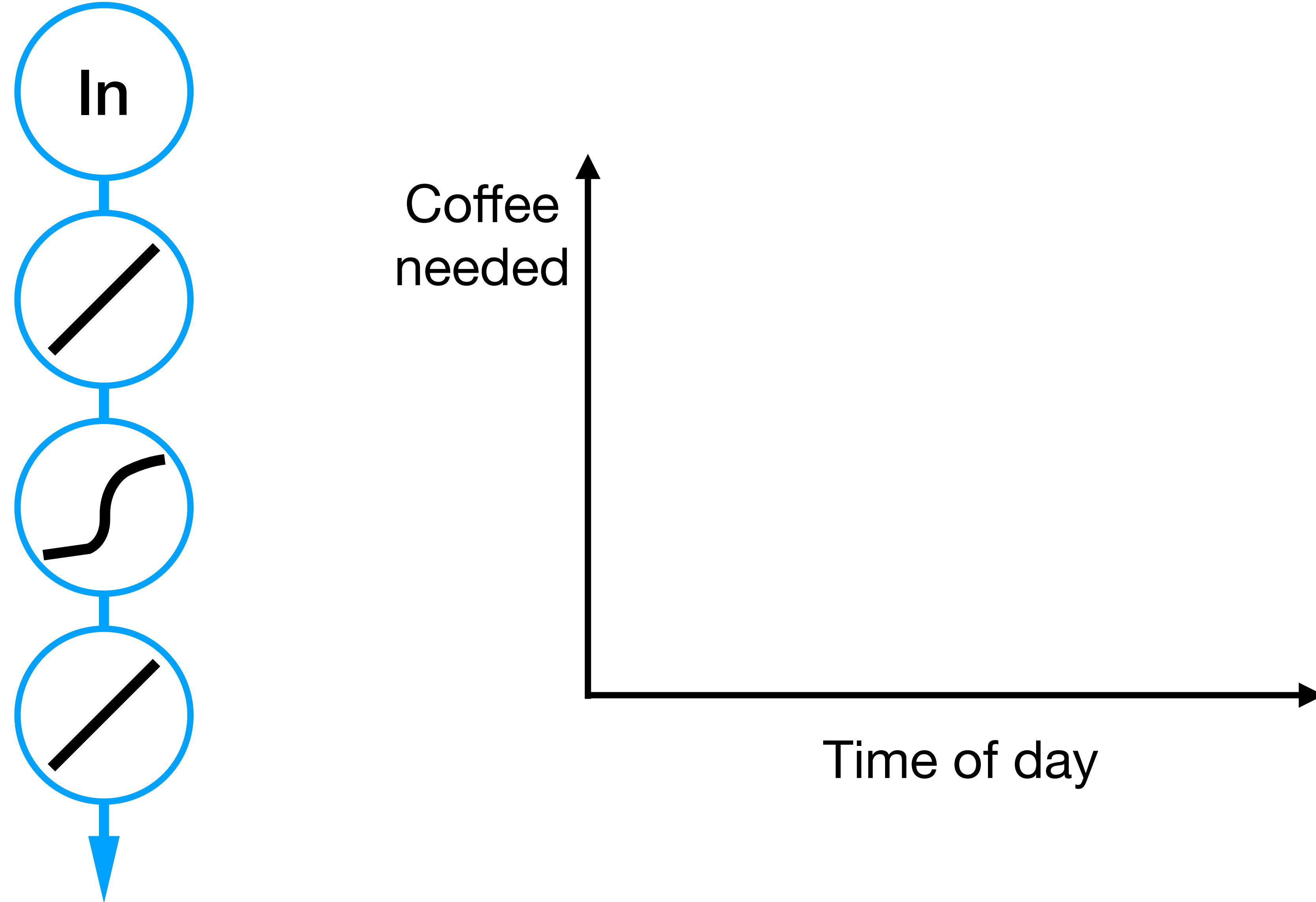


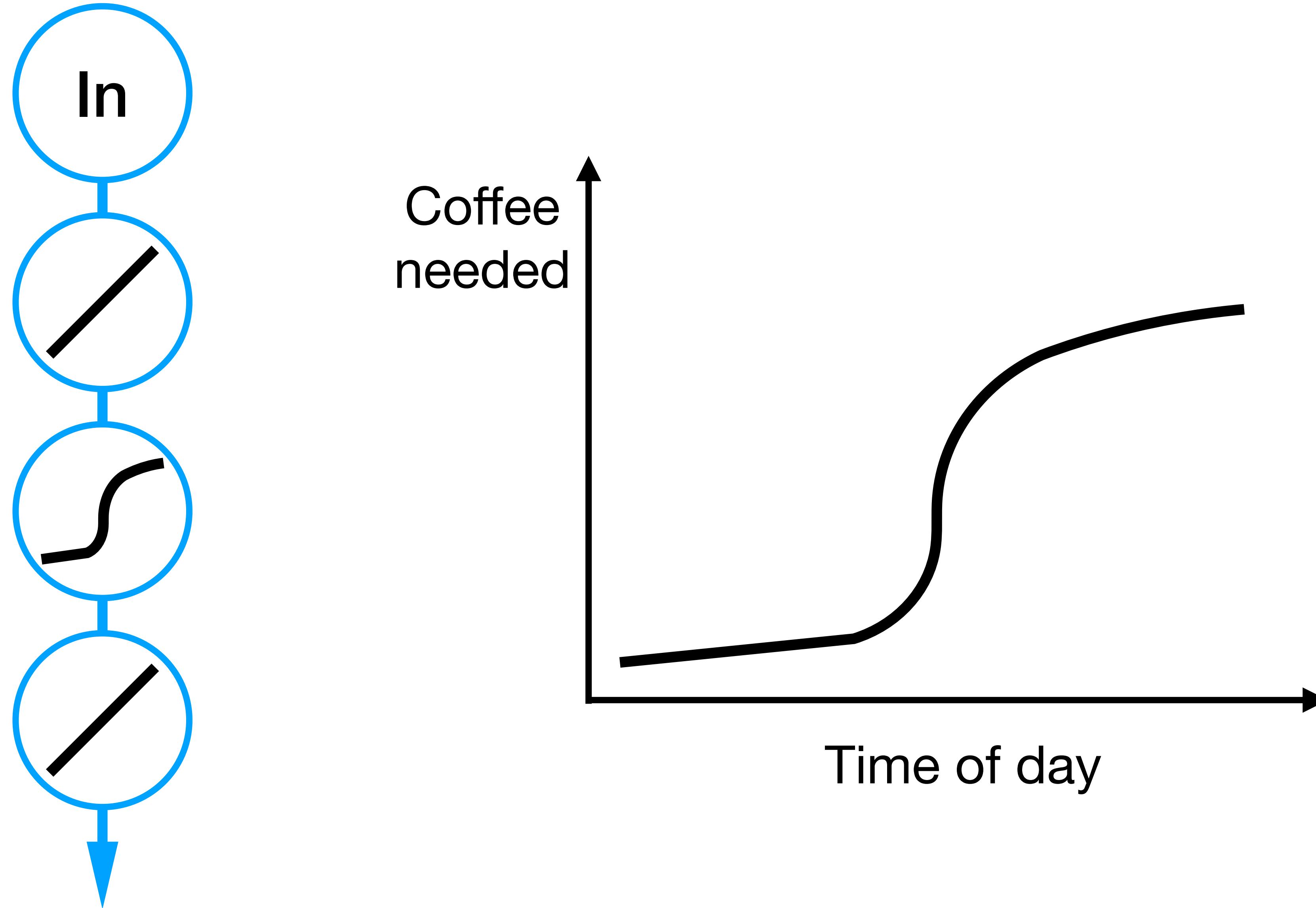
Sigmoid

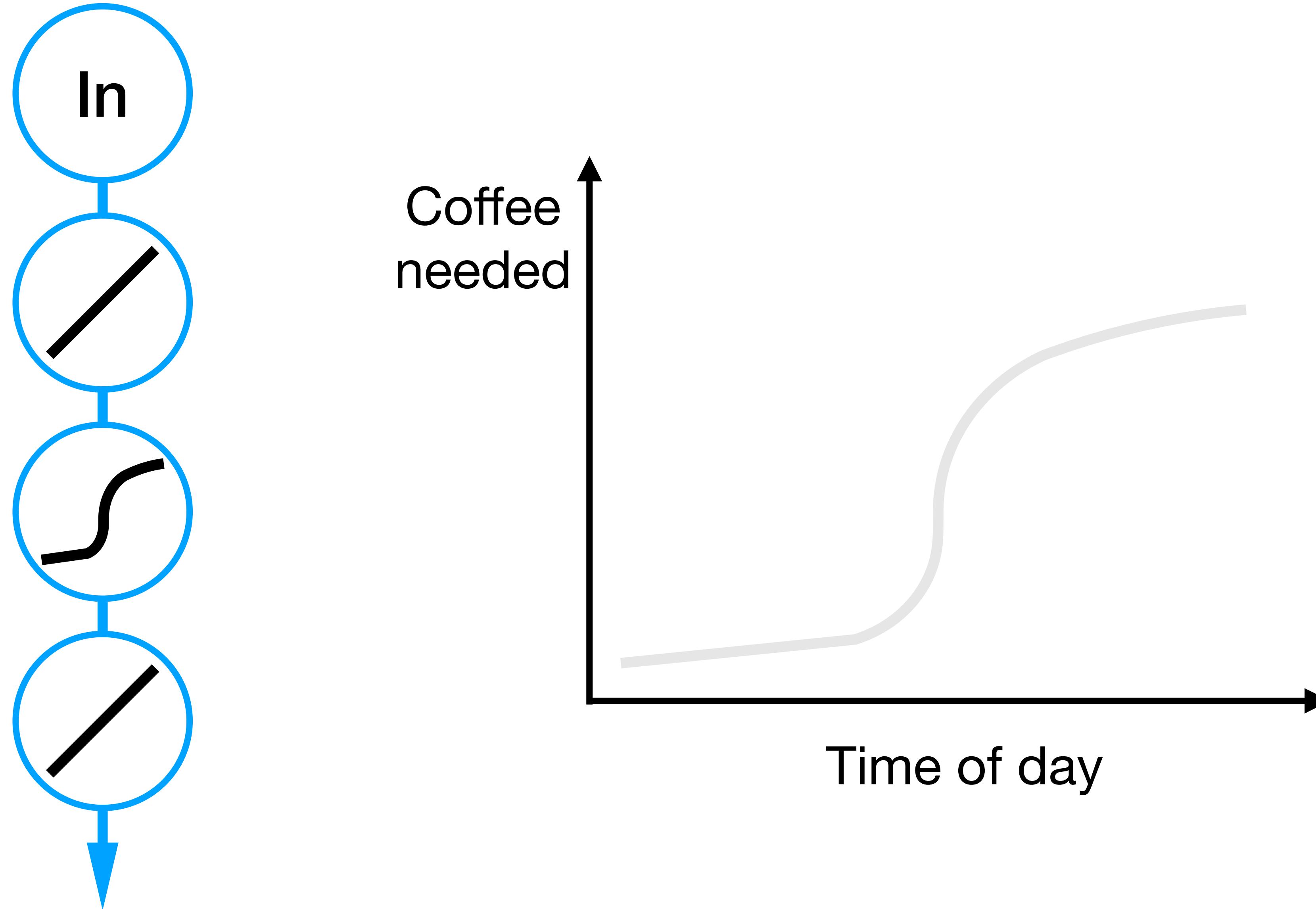


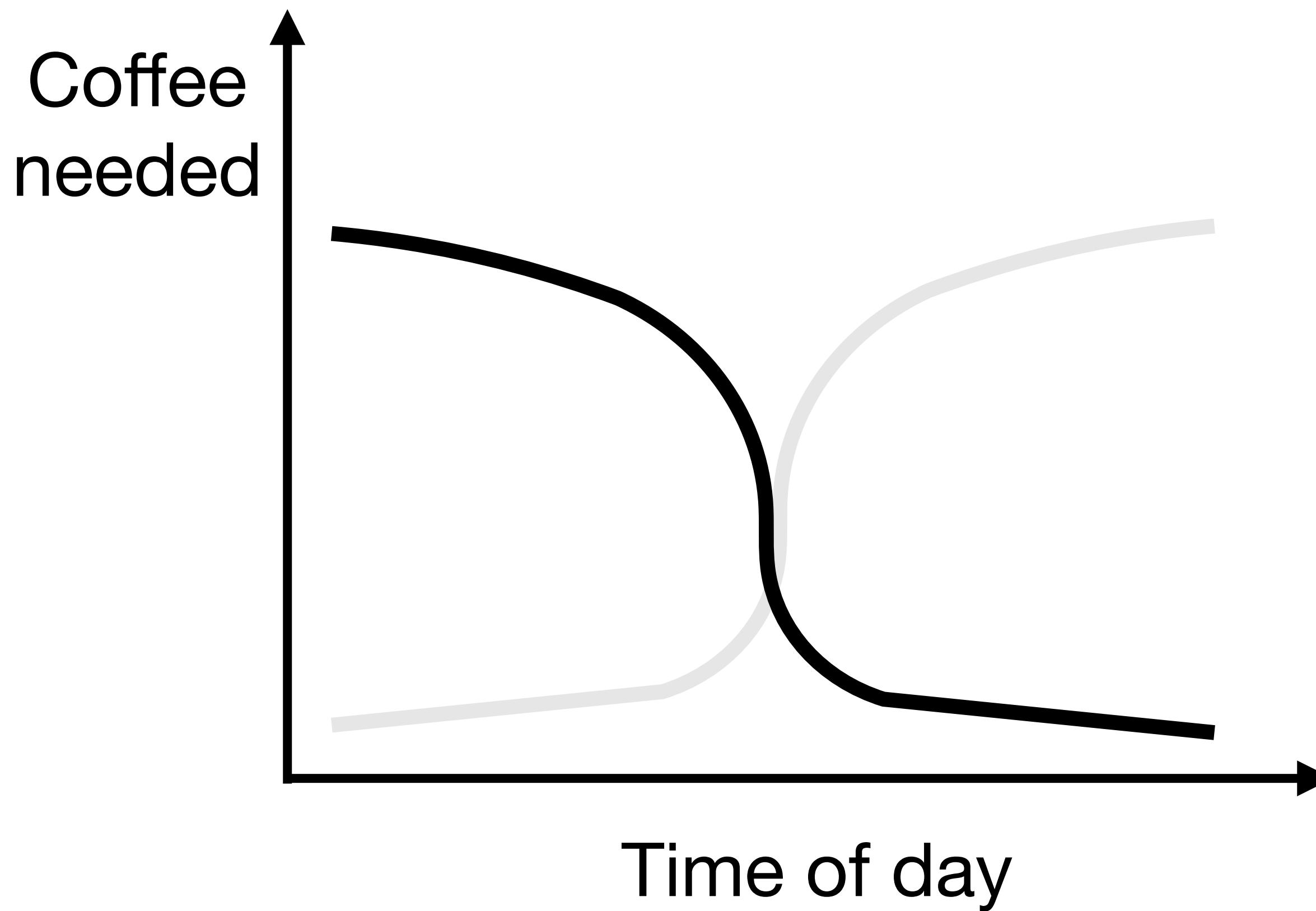
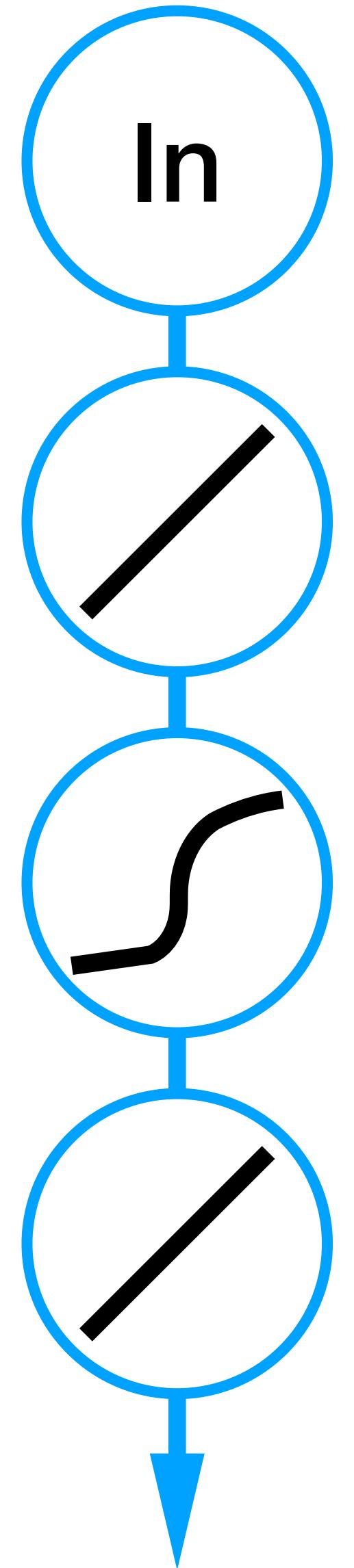


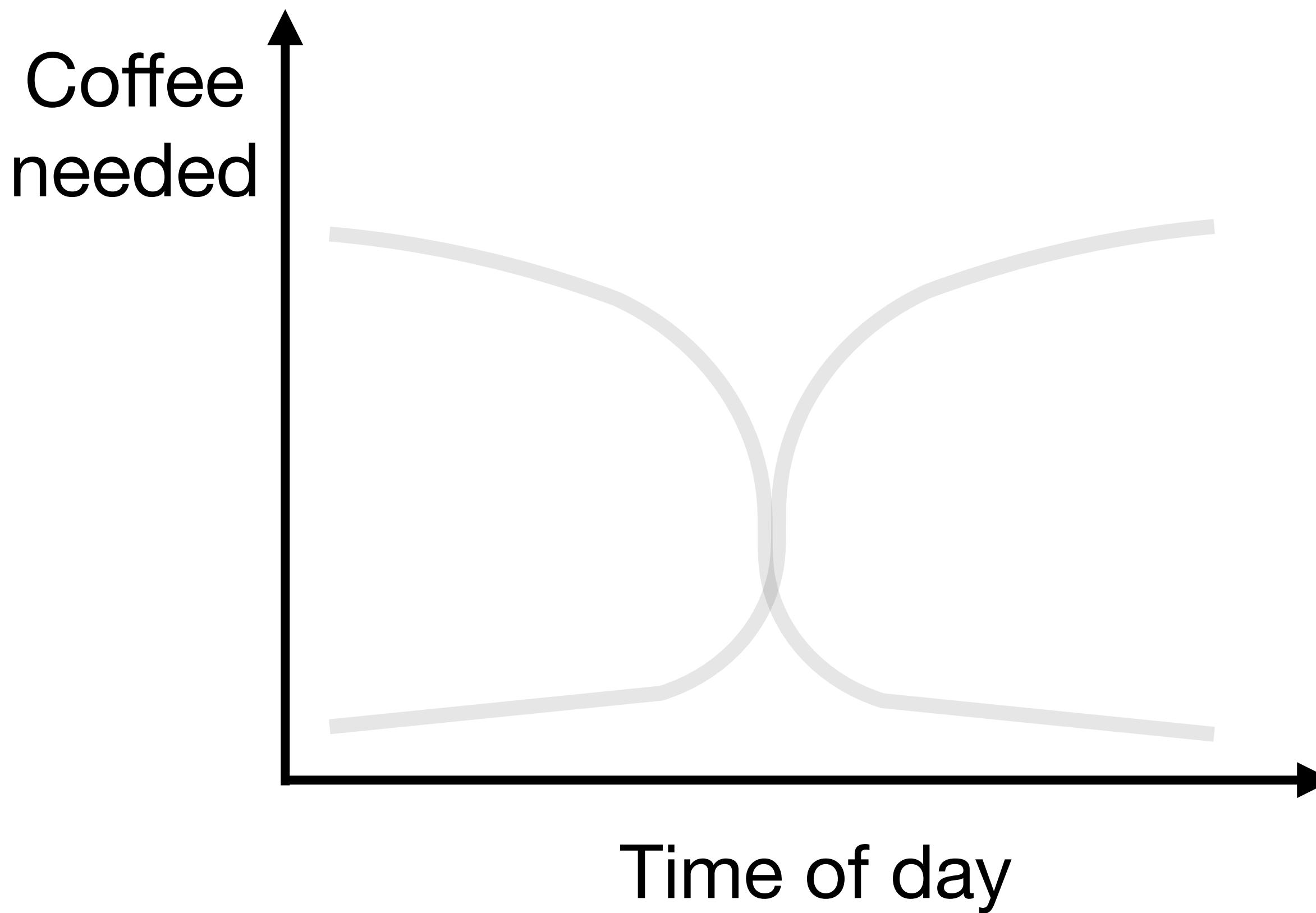
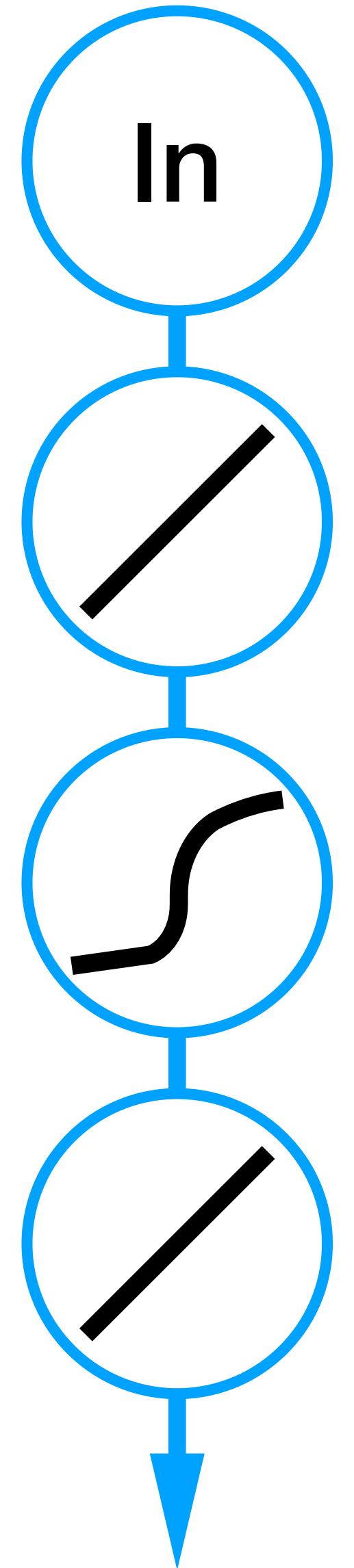


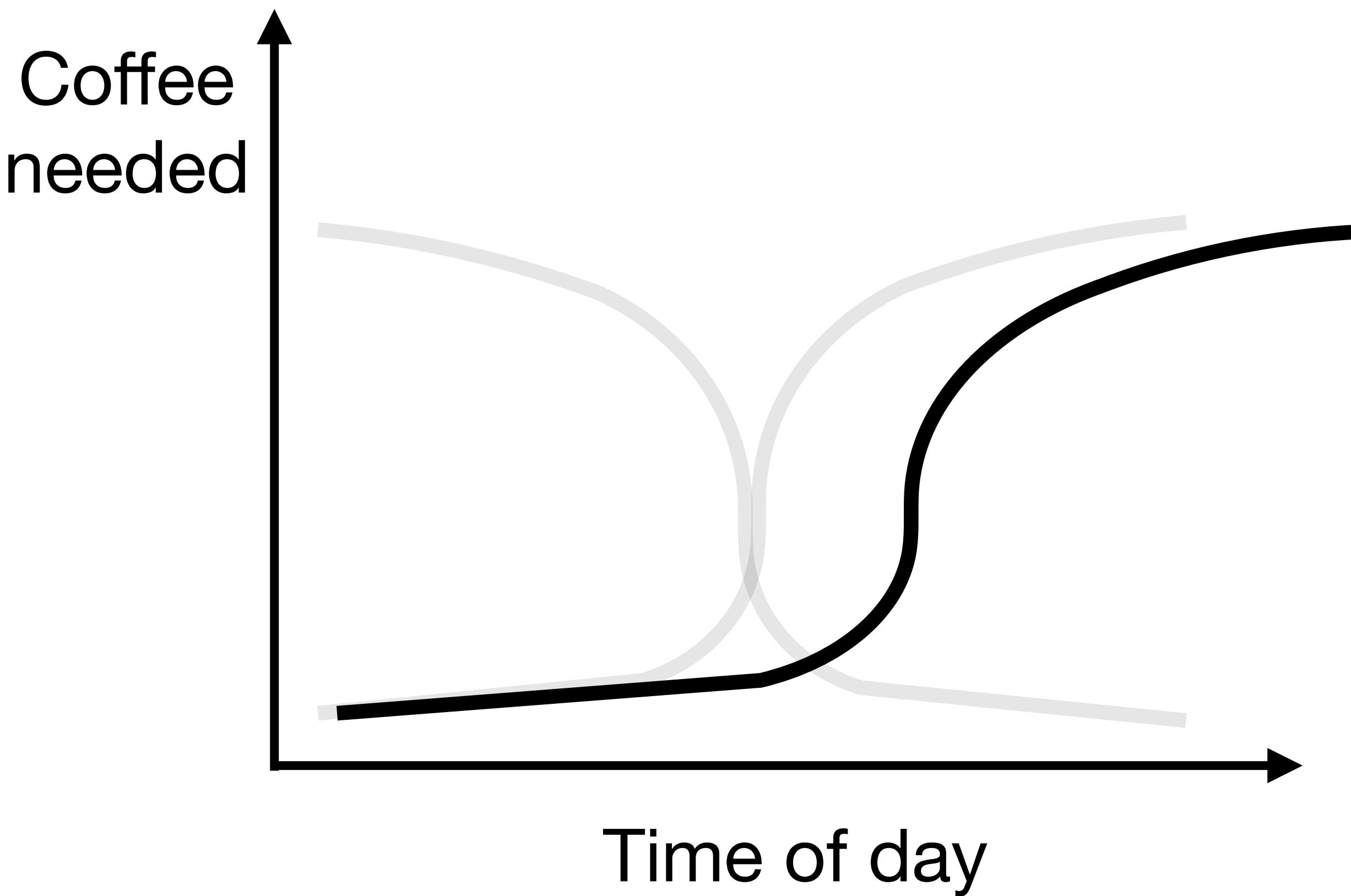
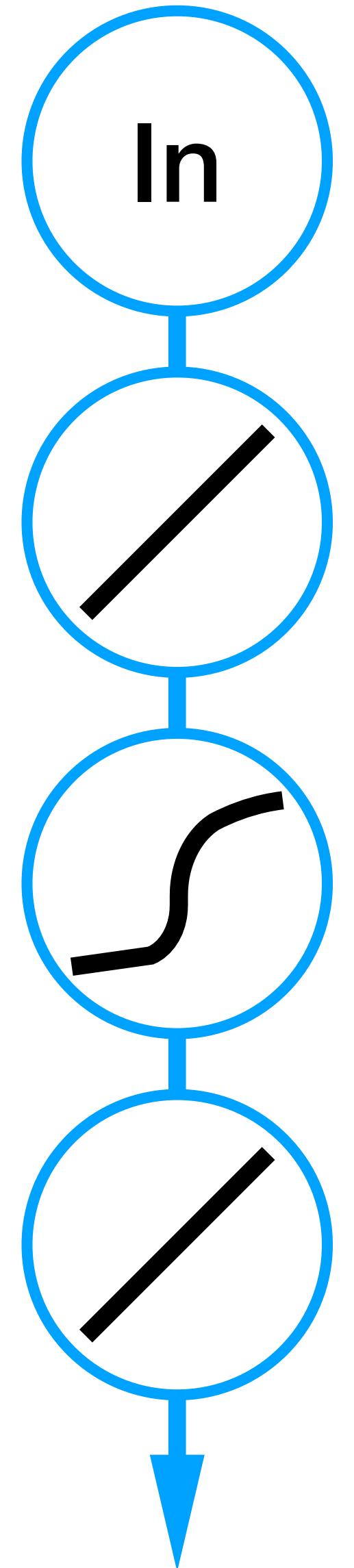


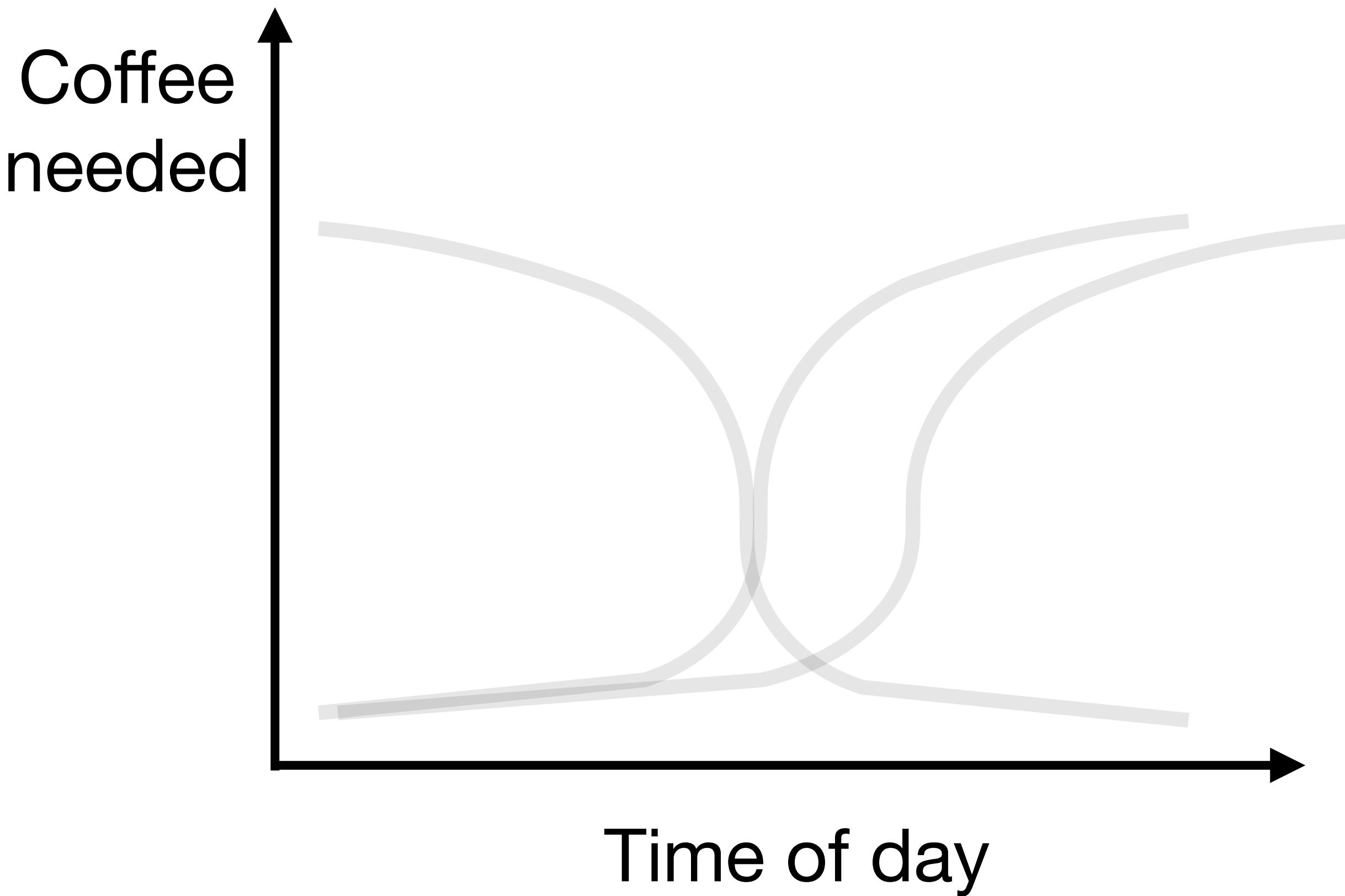
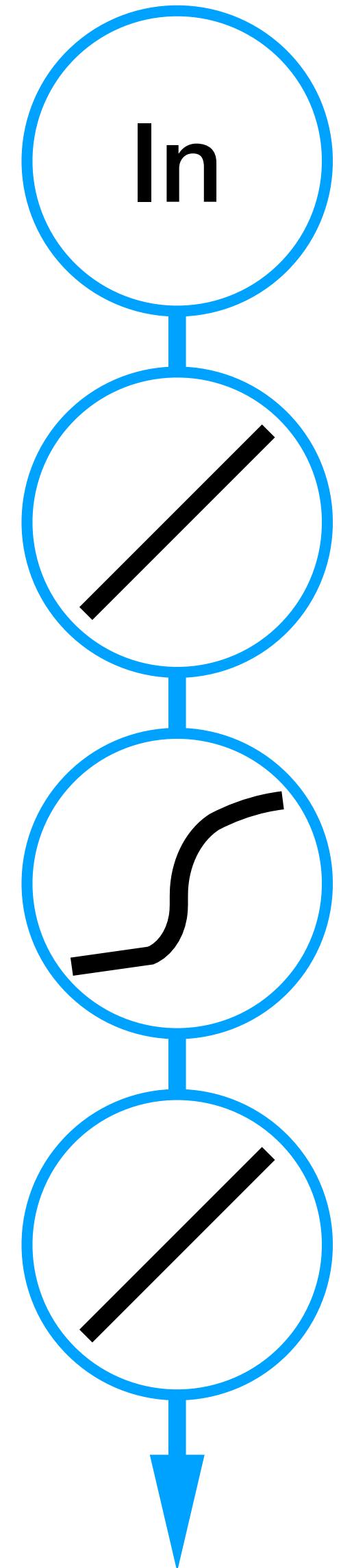


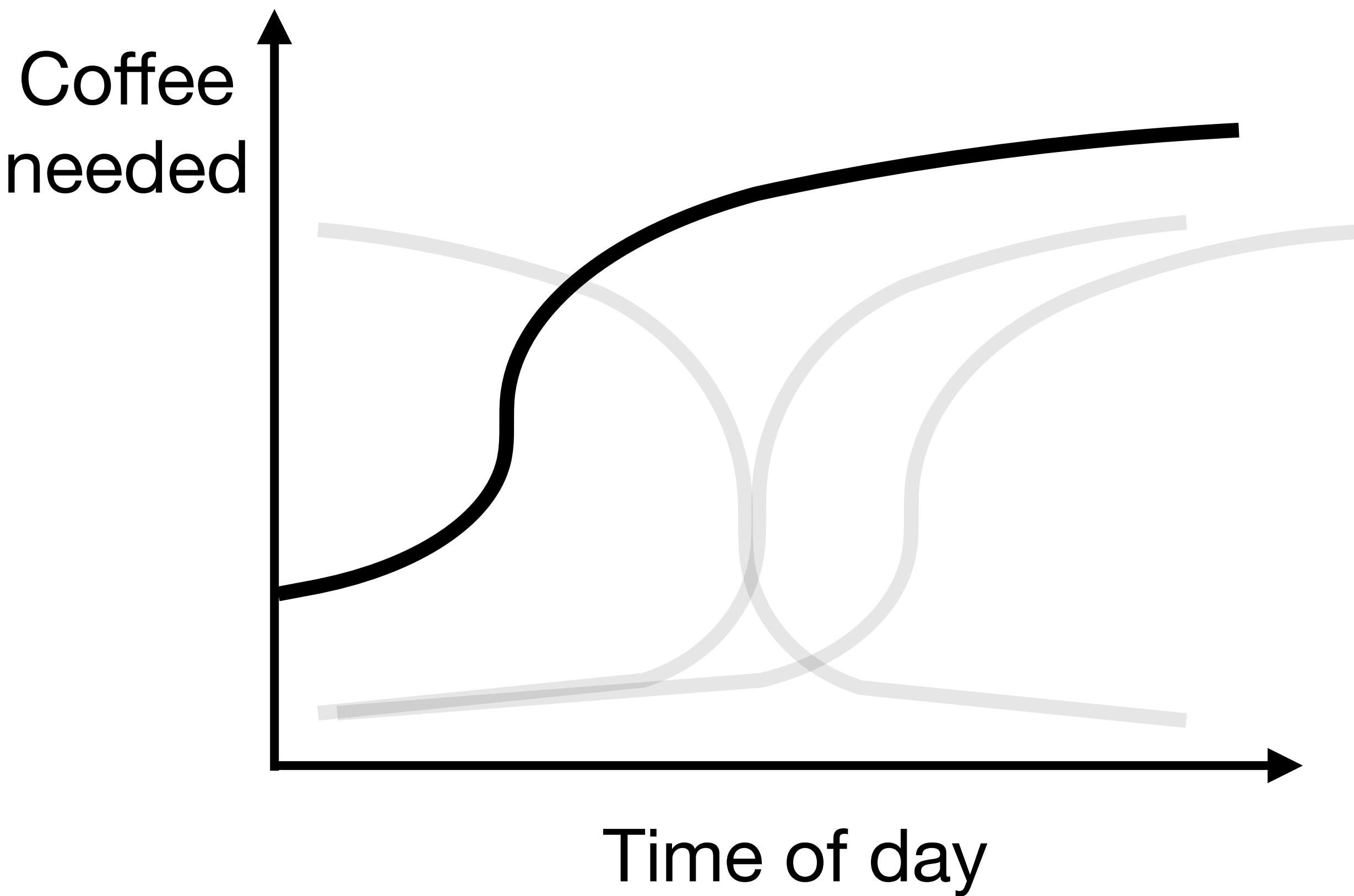
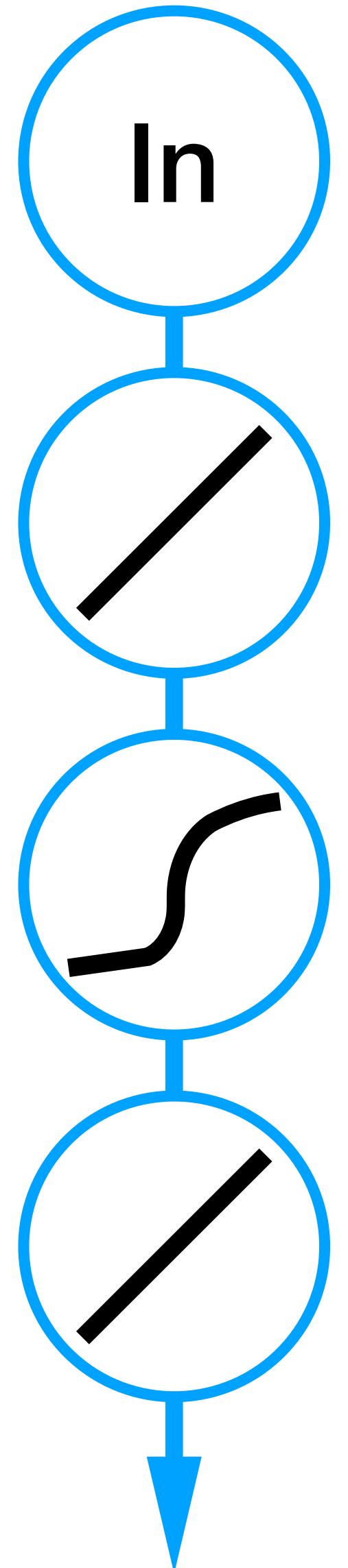


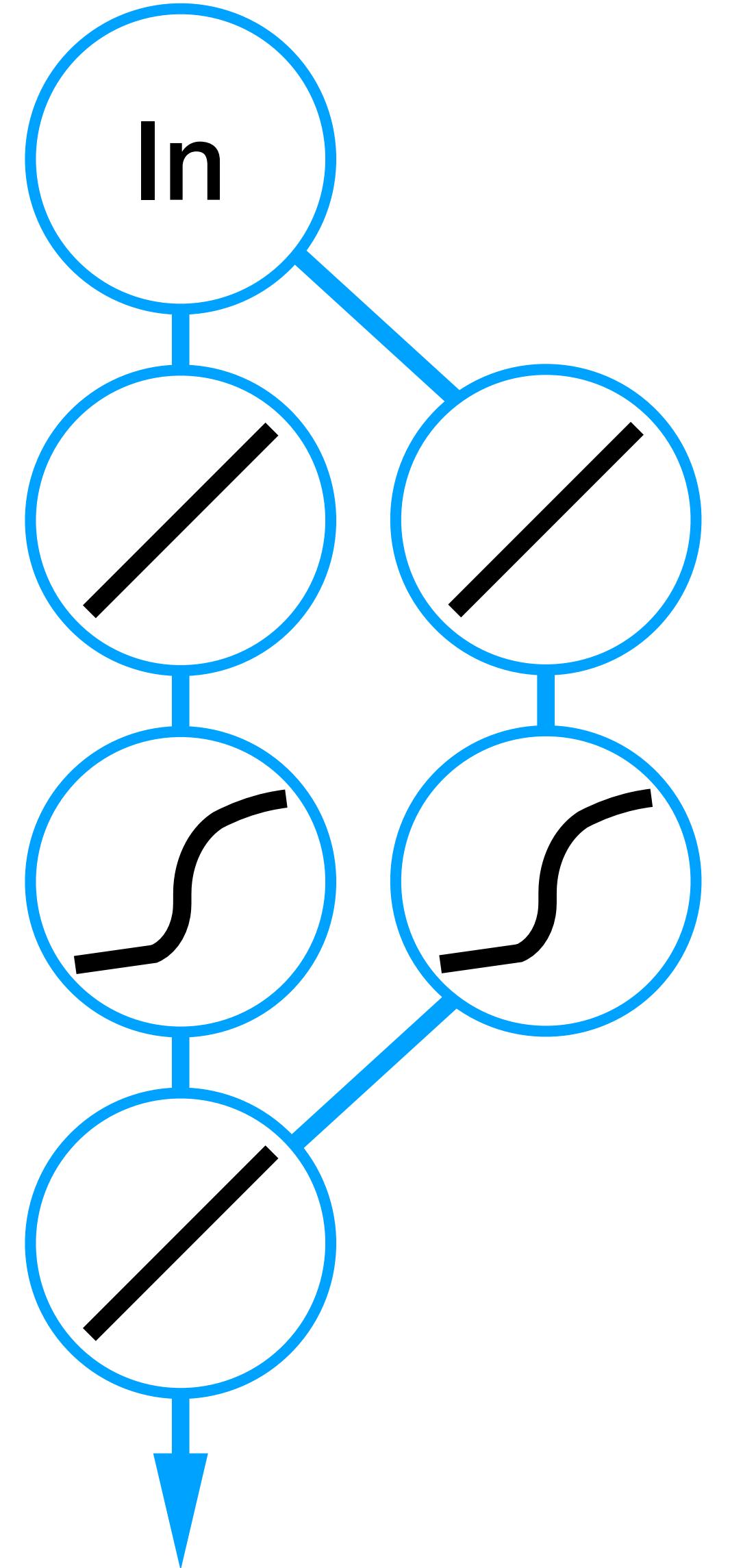


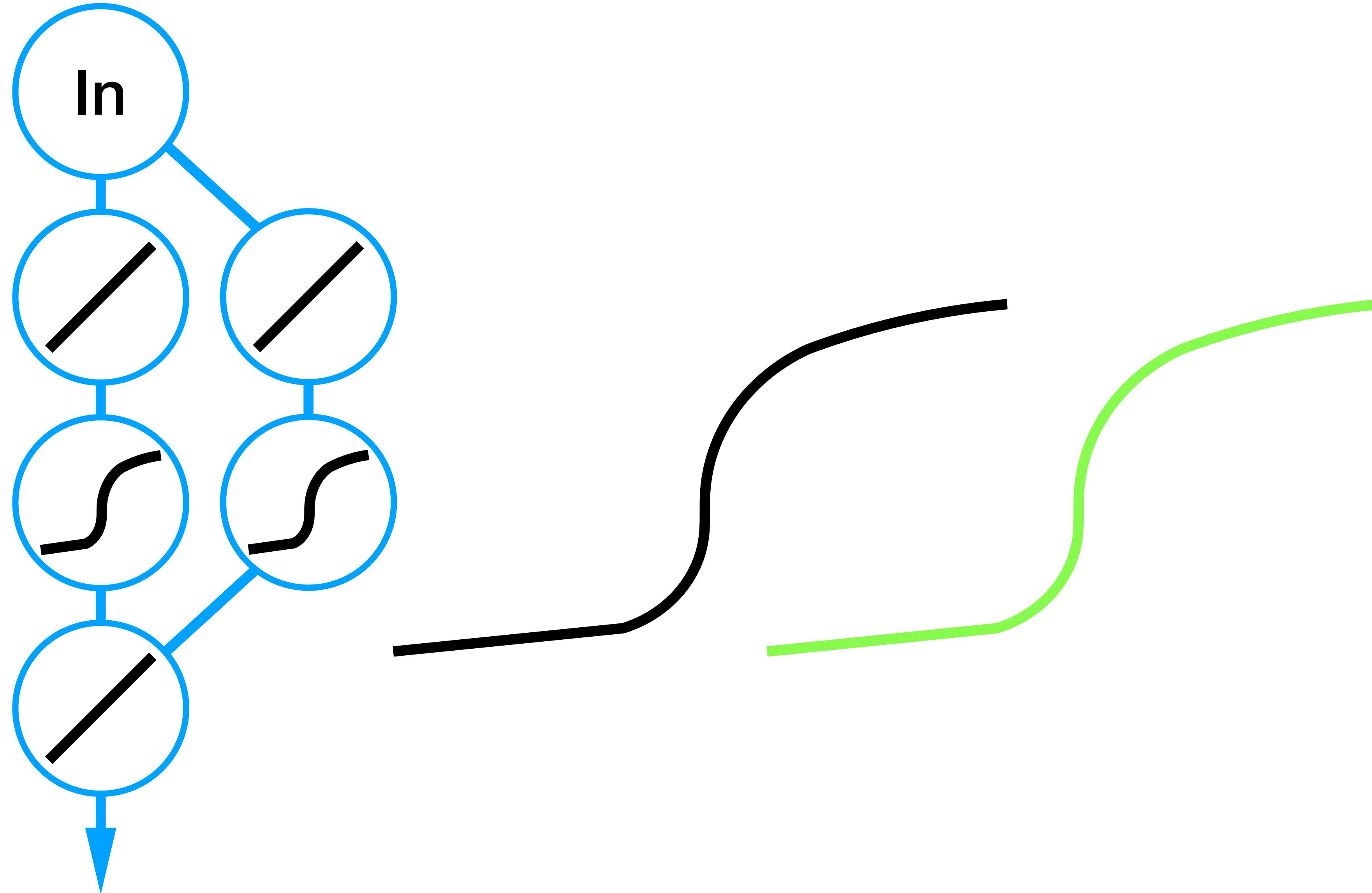


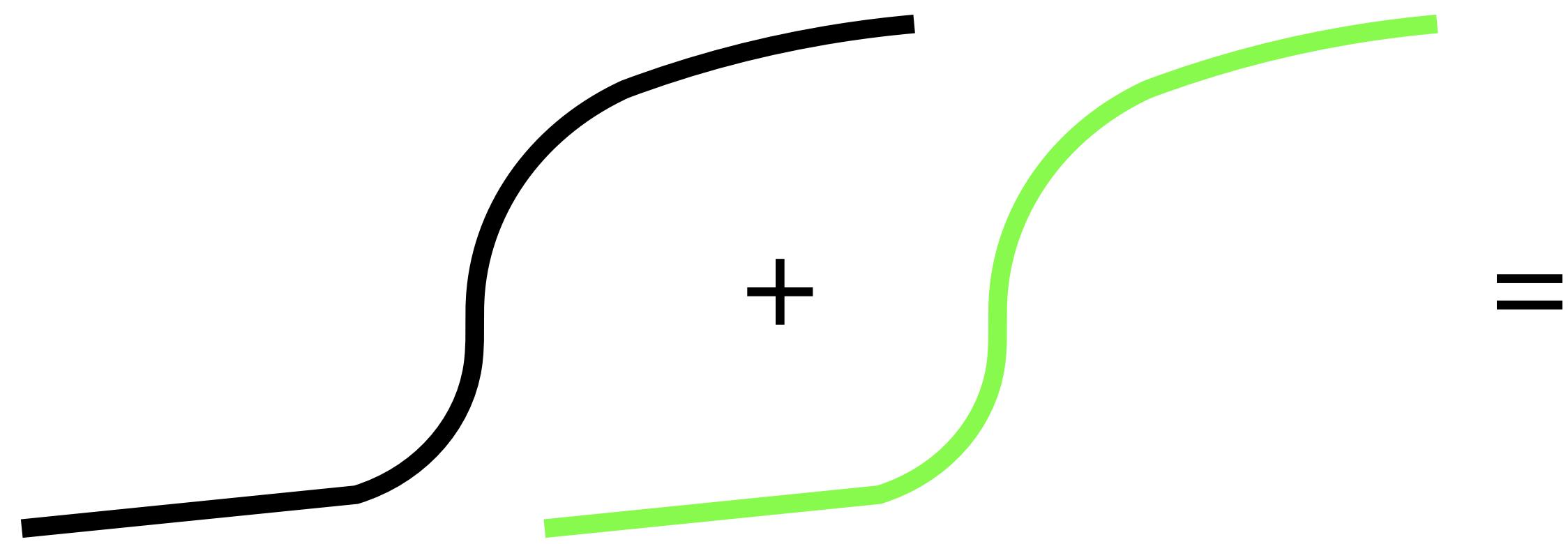


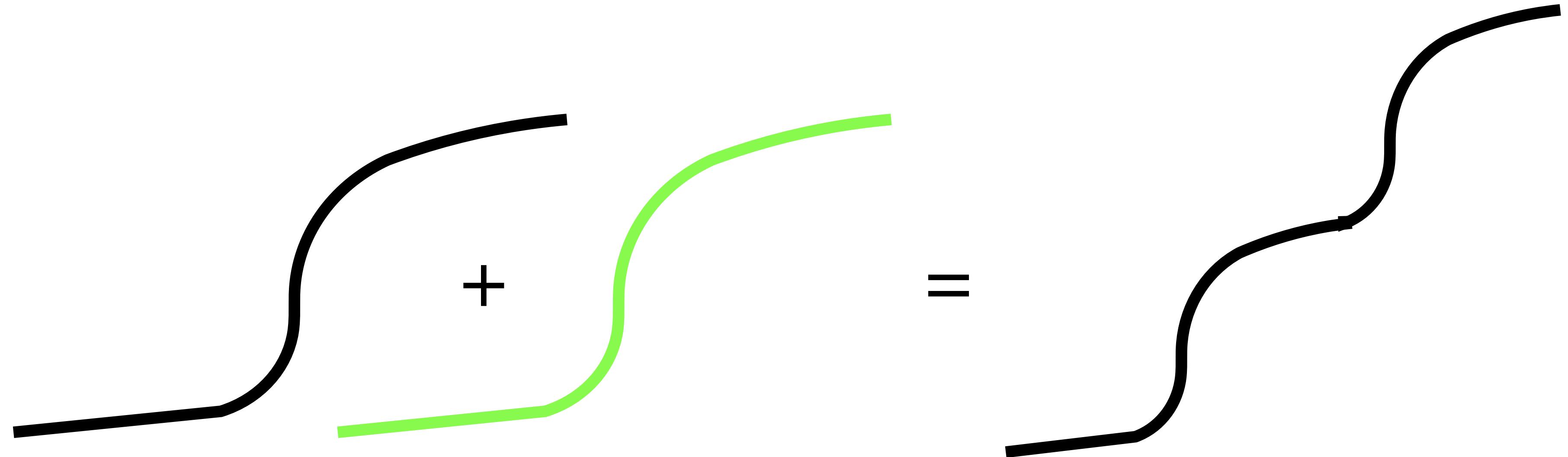


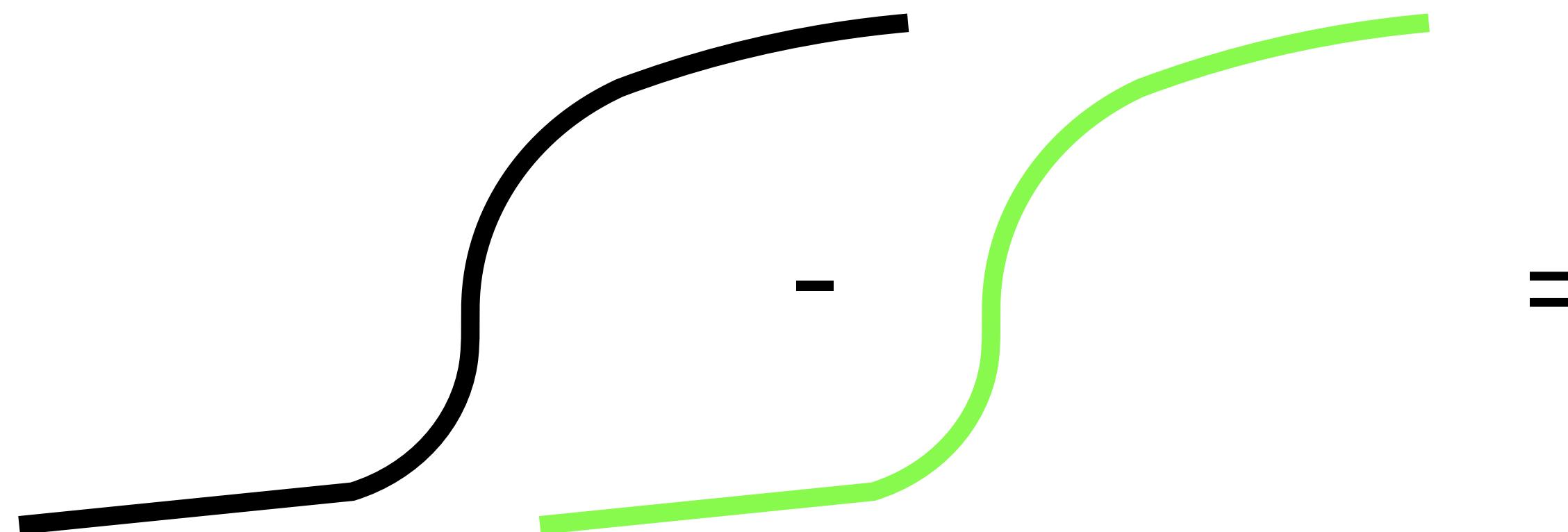
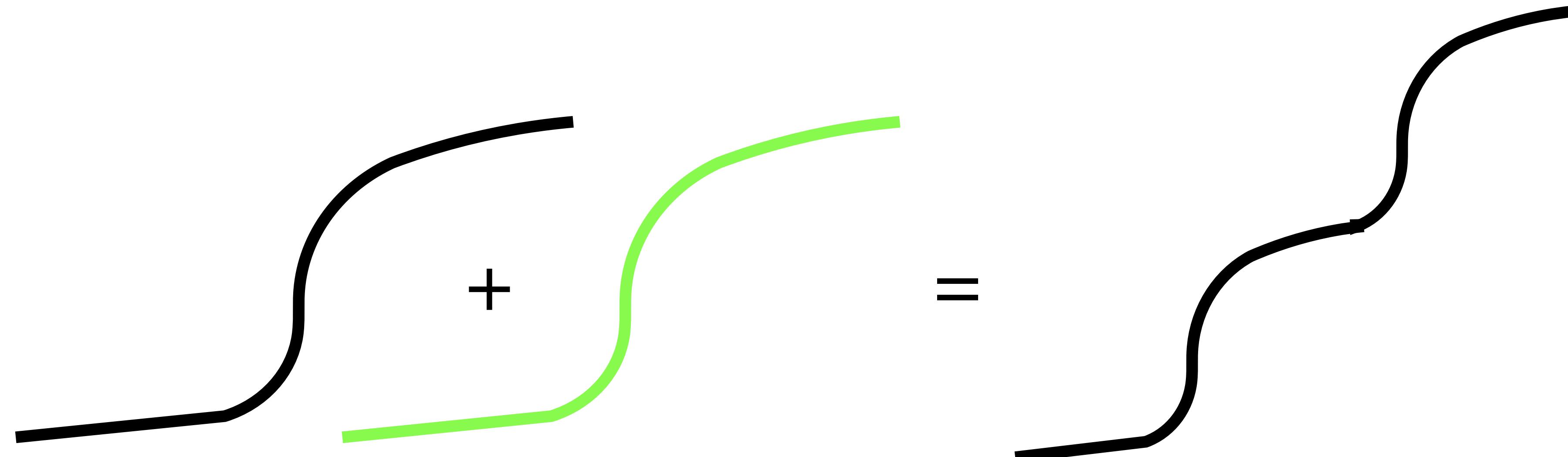


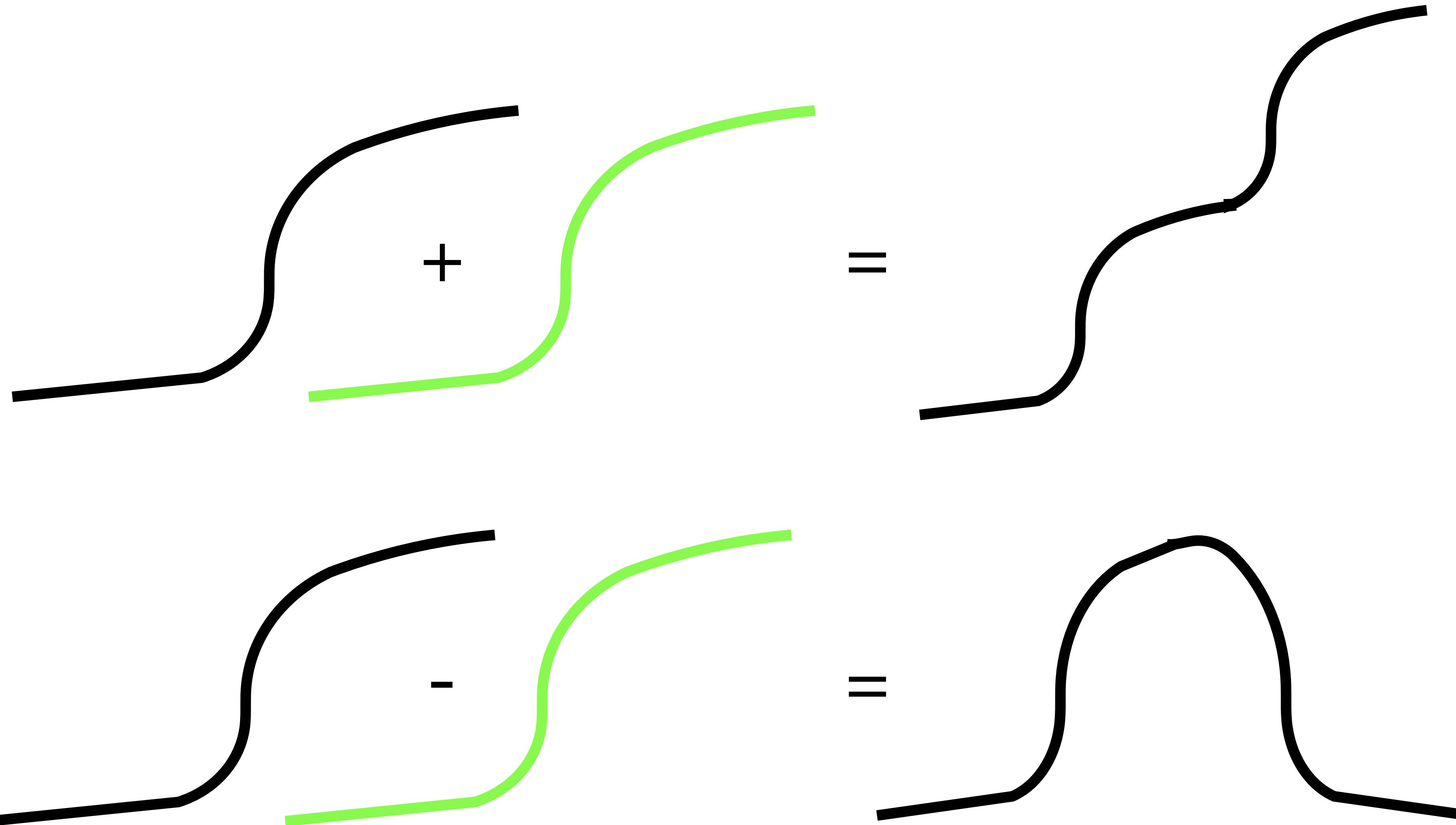


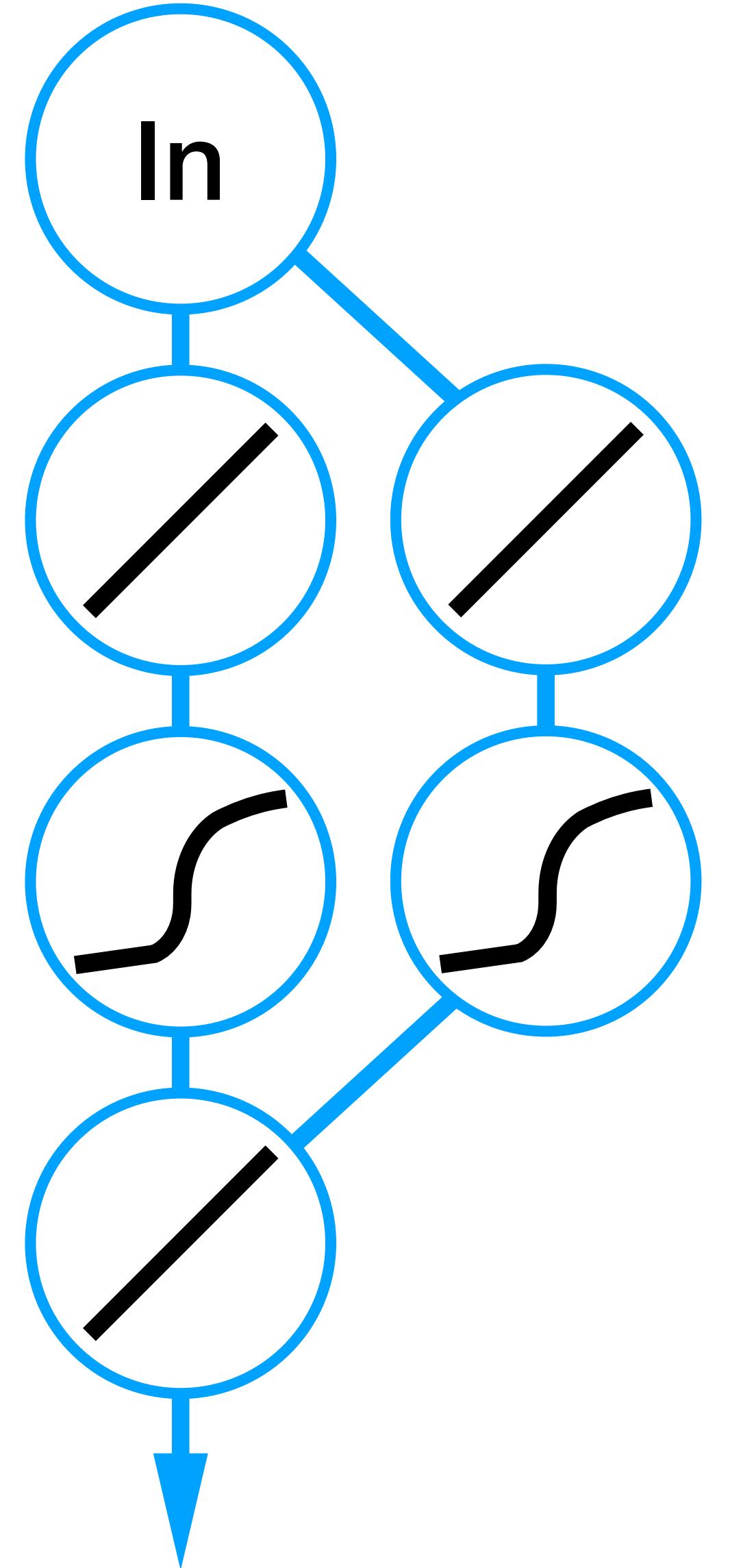


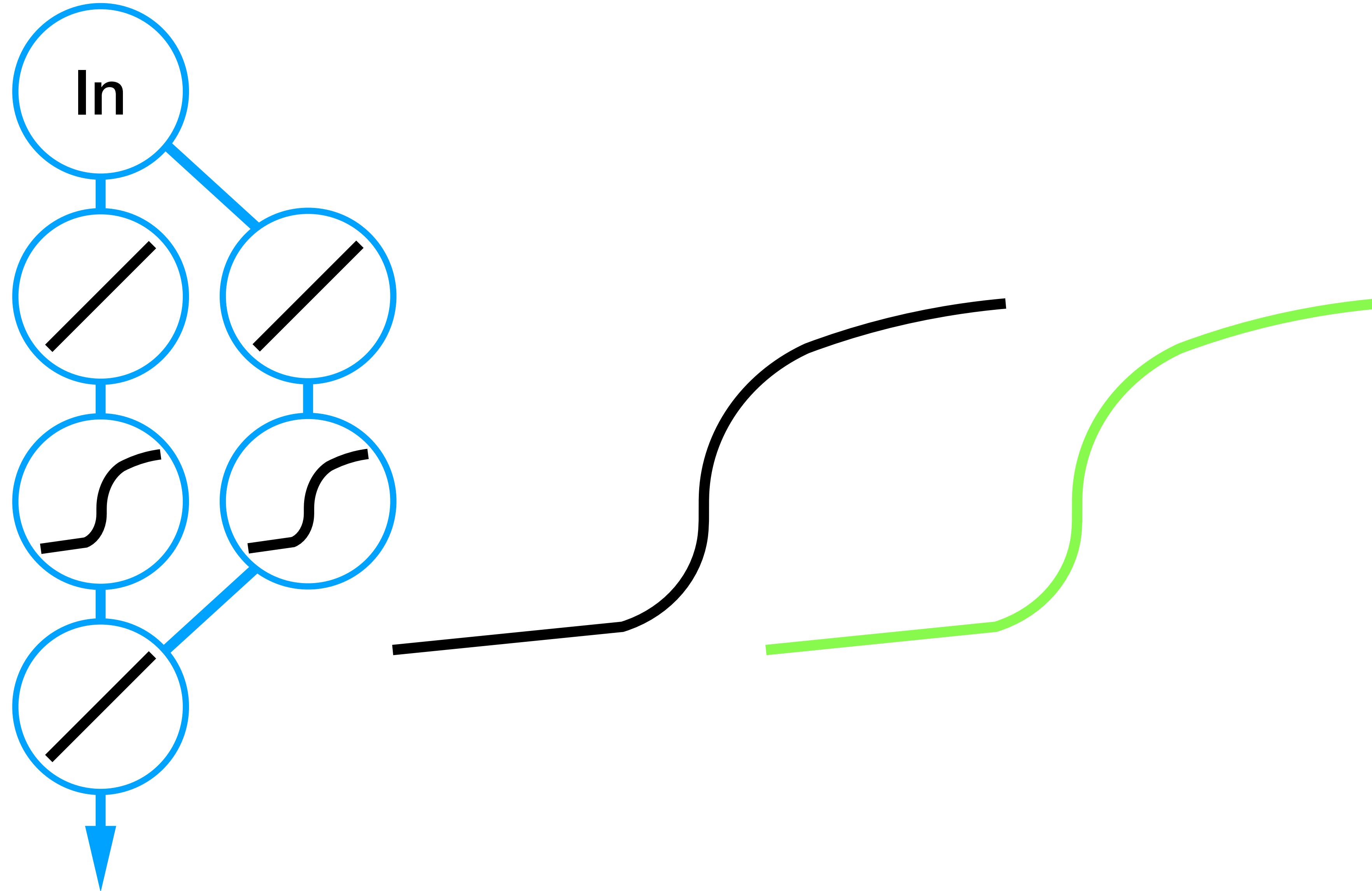


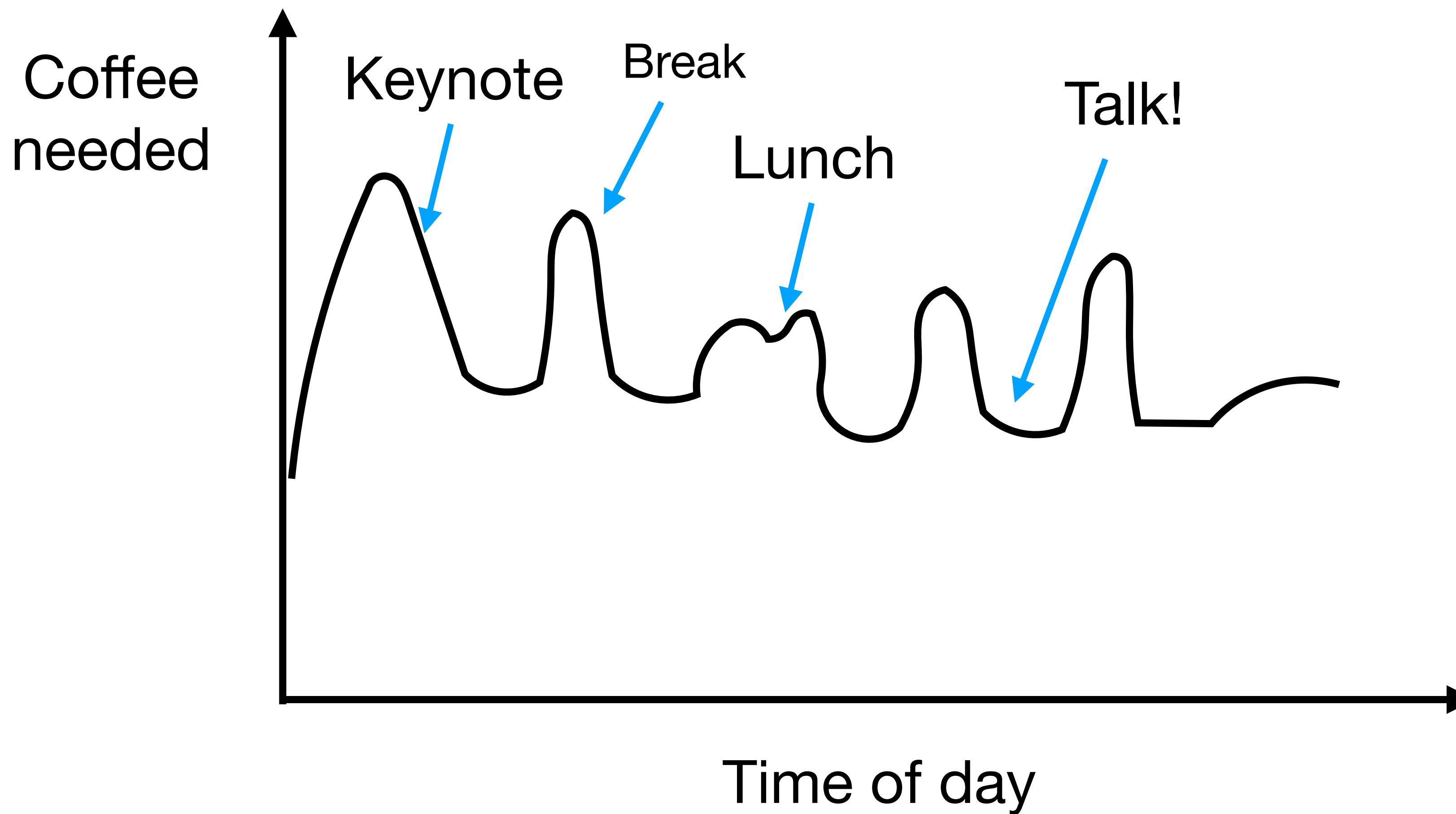


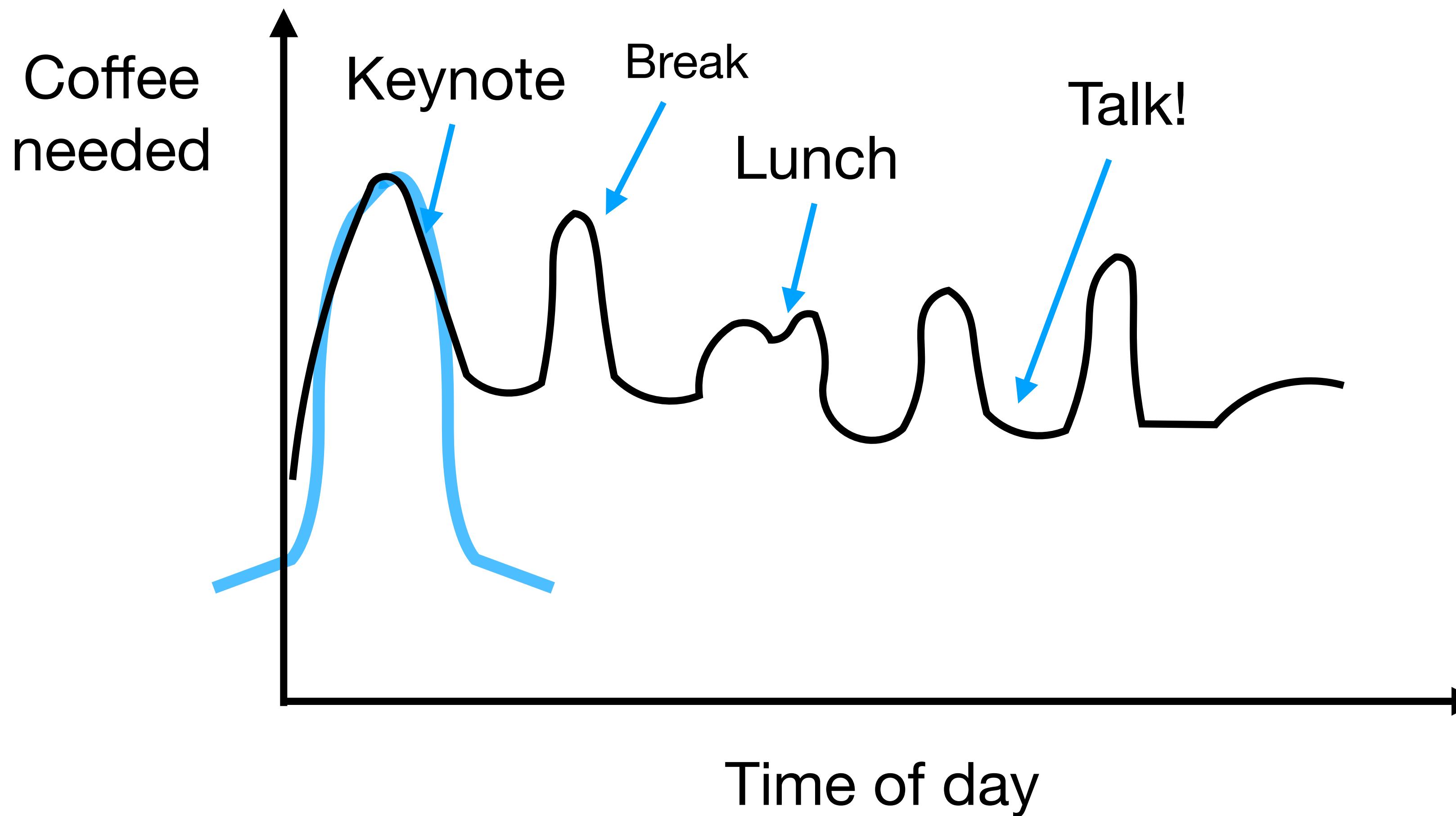


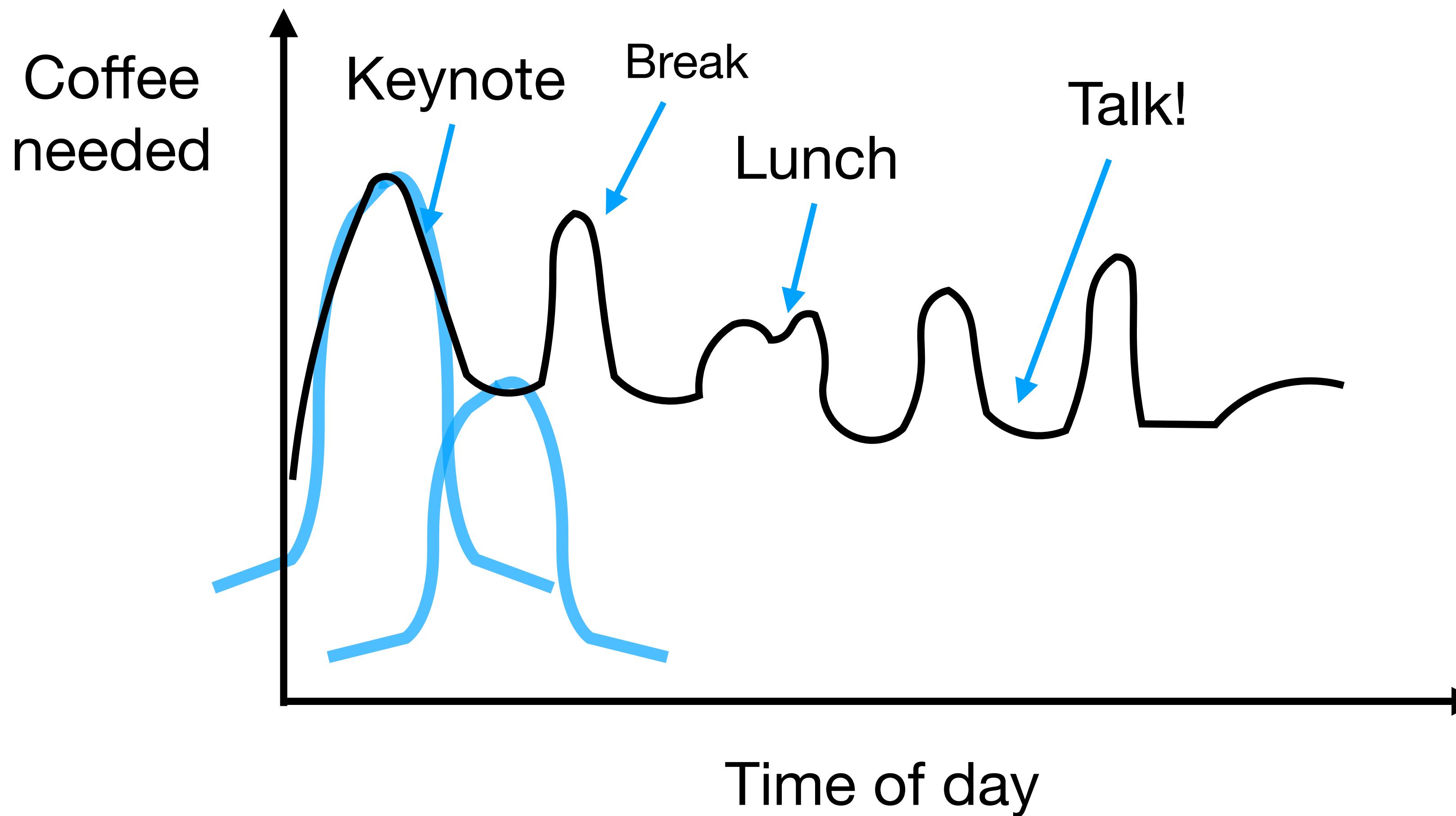


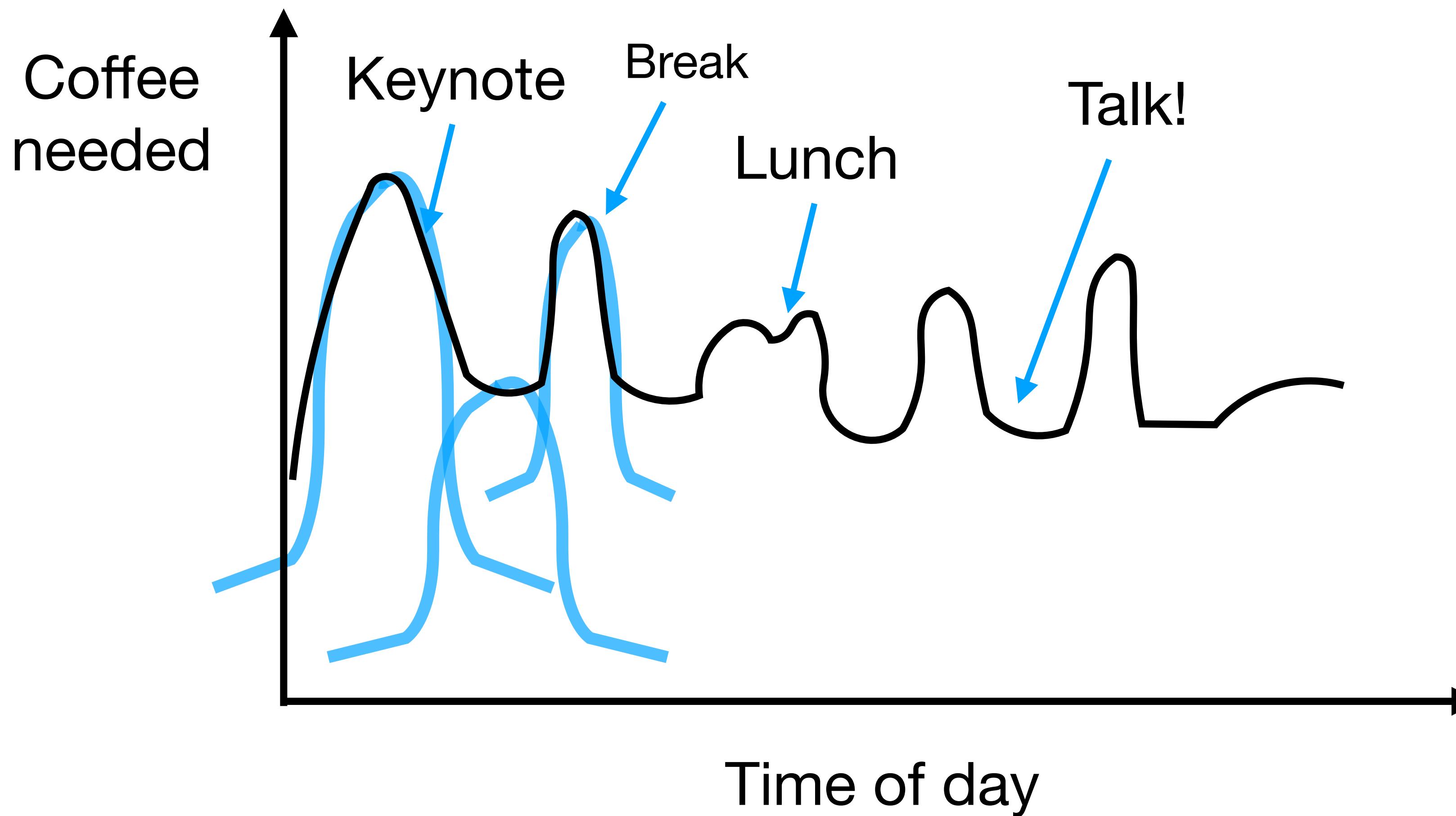


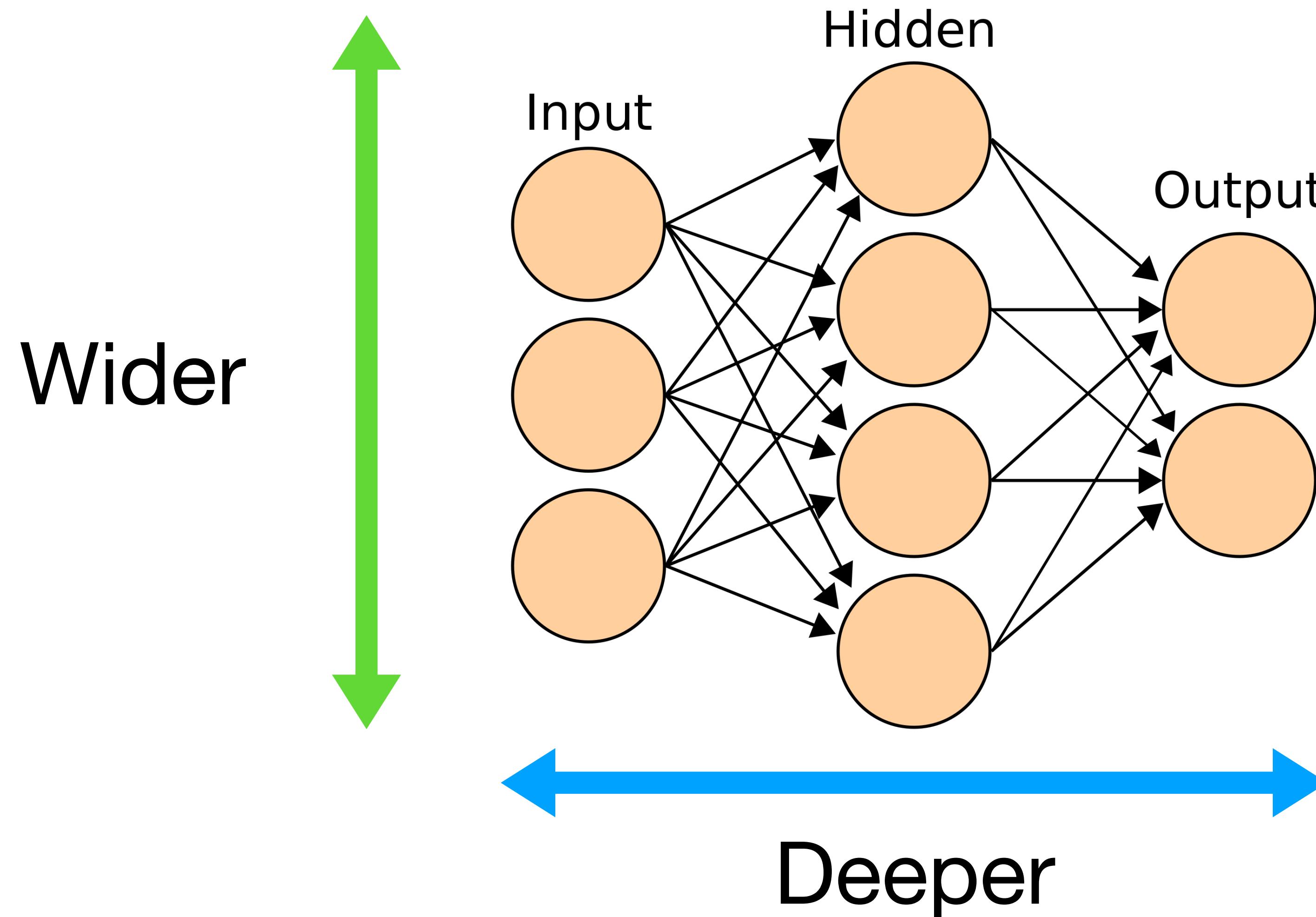








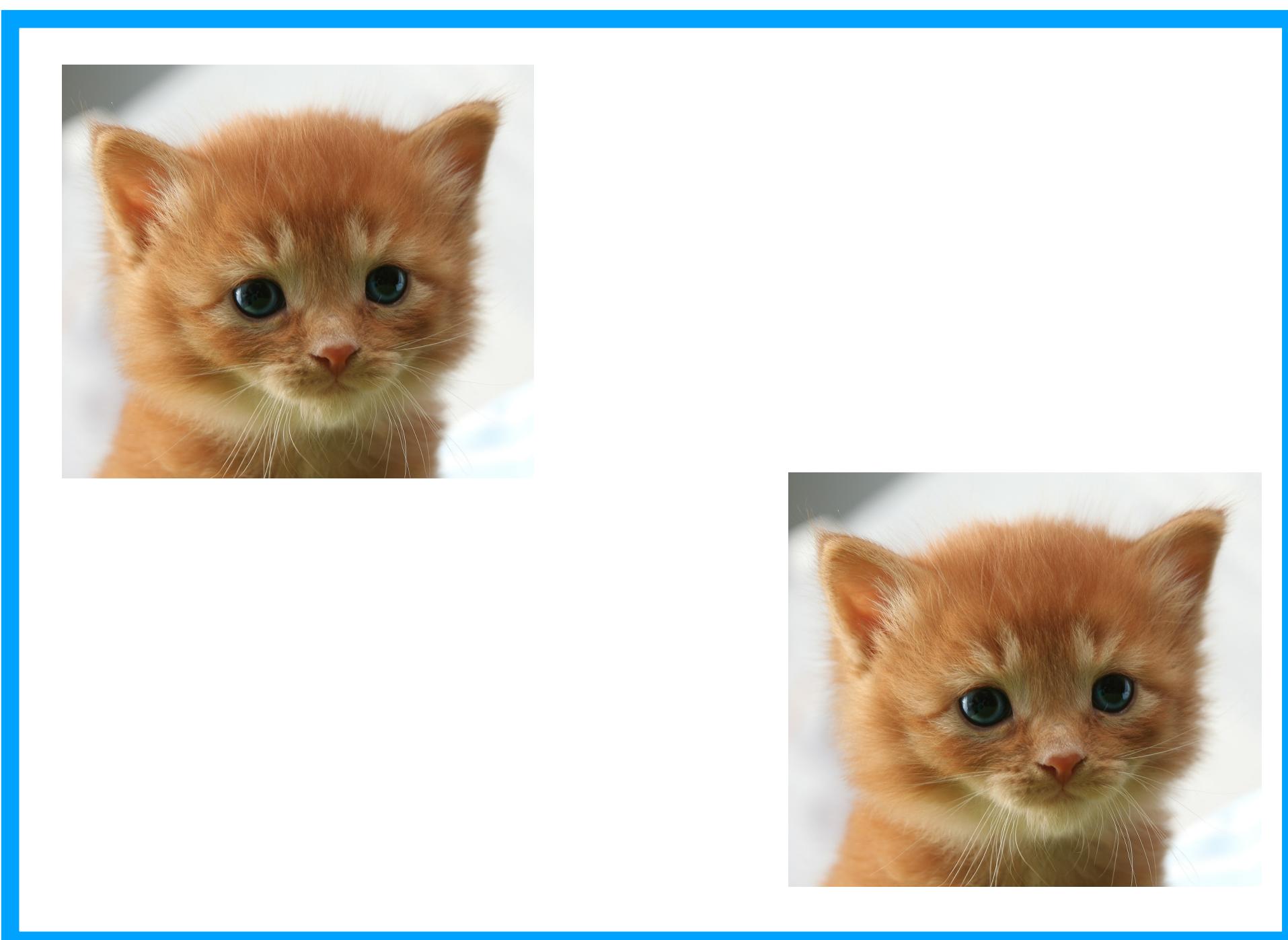




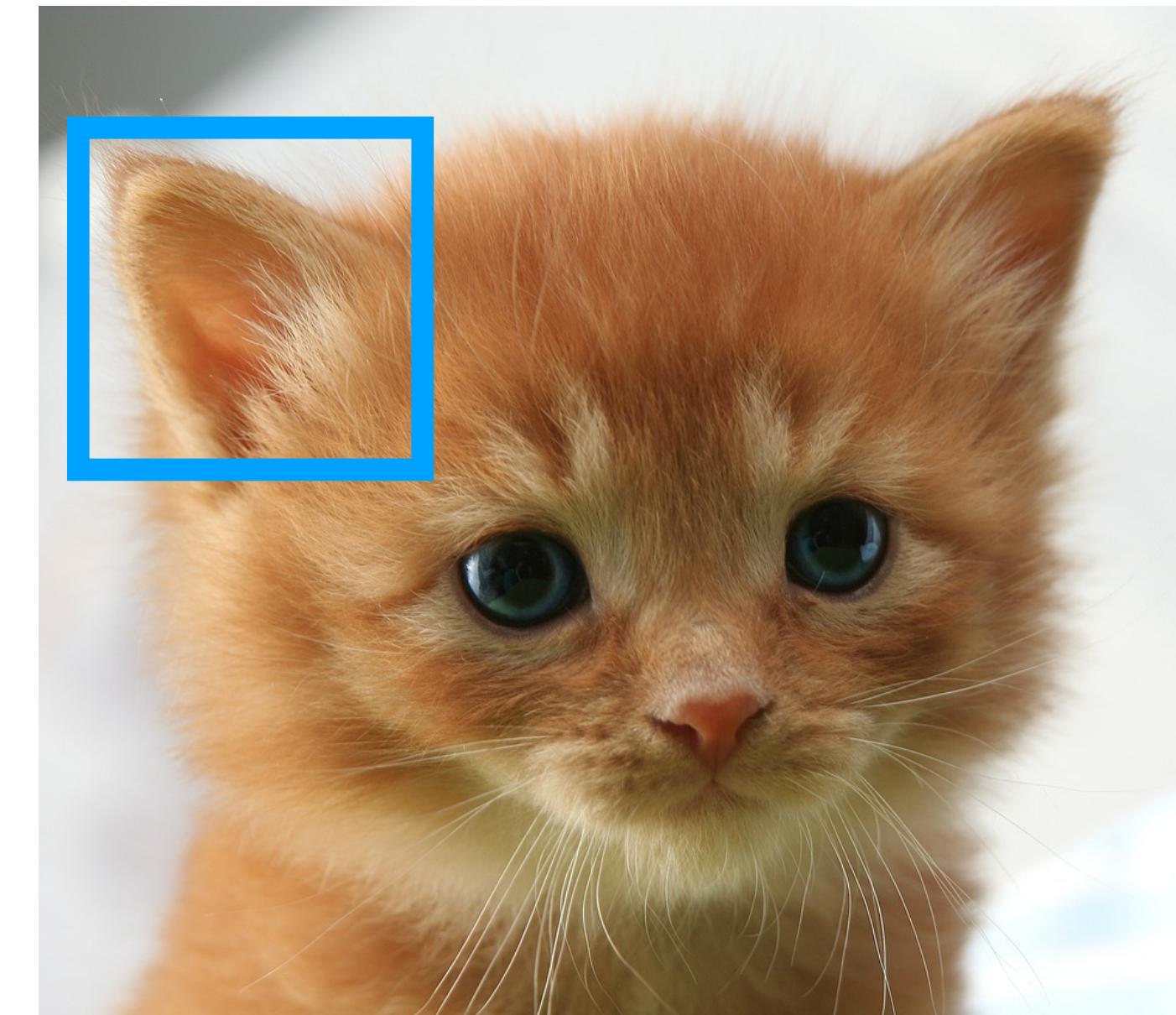
Foundational for CNNs and RNNs

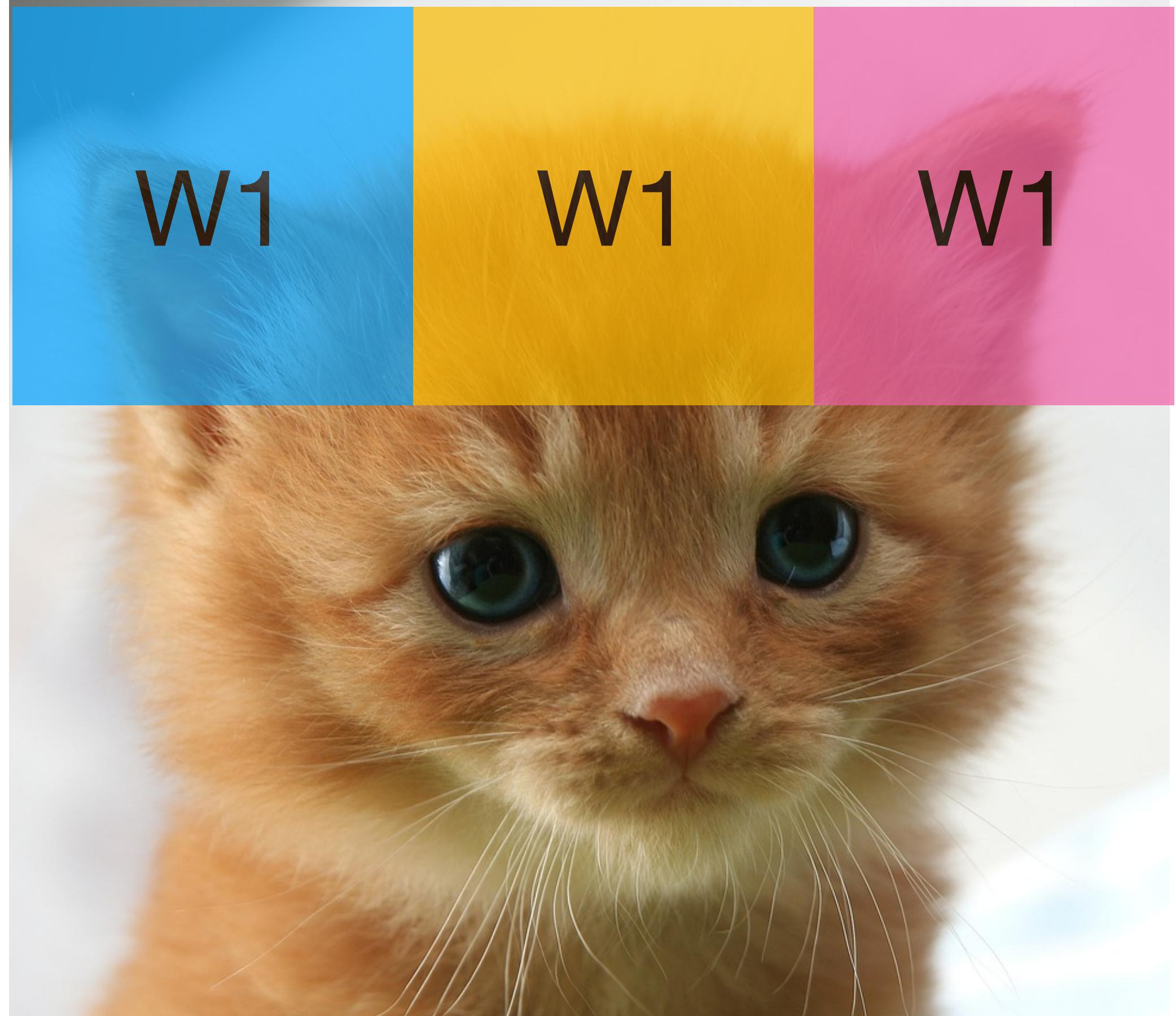
Convolutional layers

1. Translation invariance

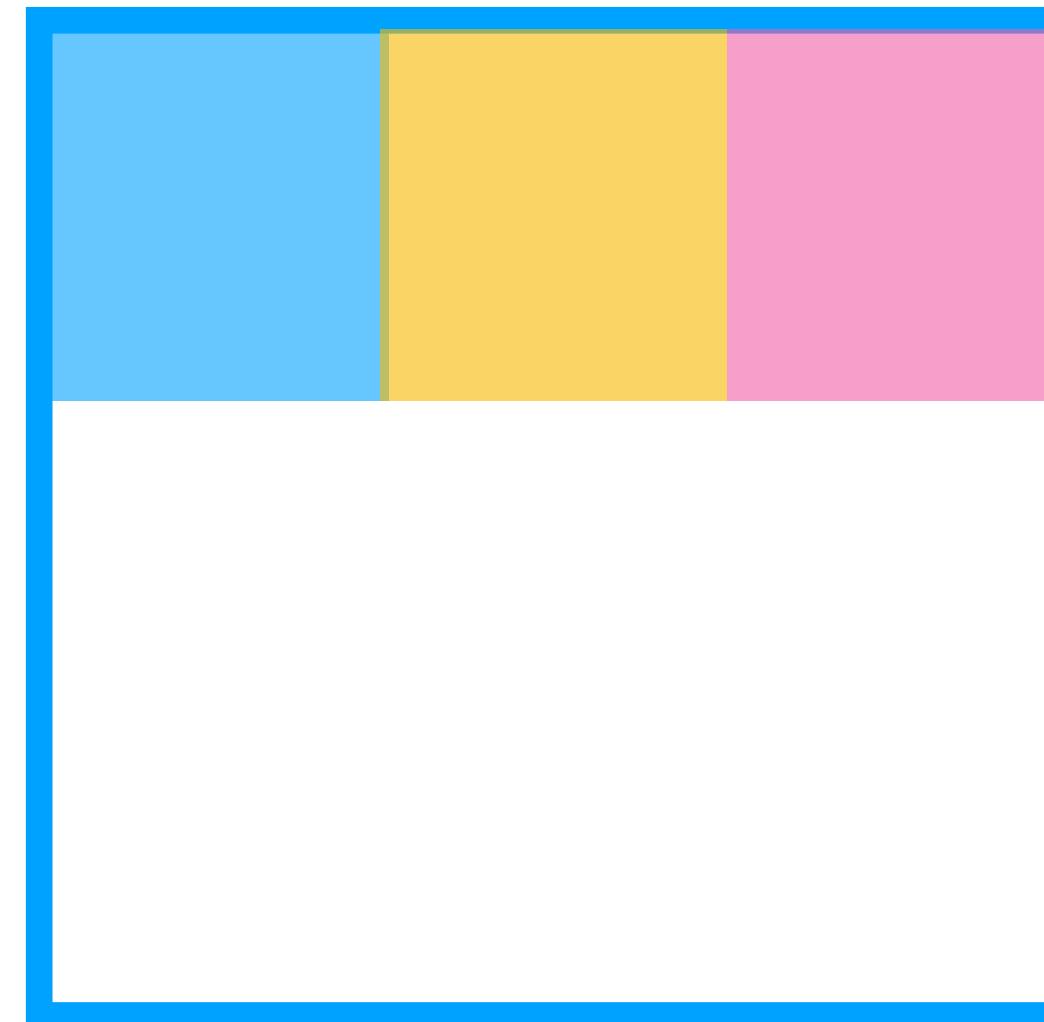


2. Locality

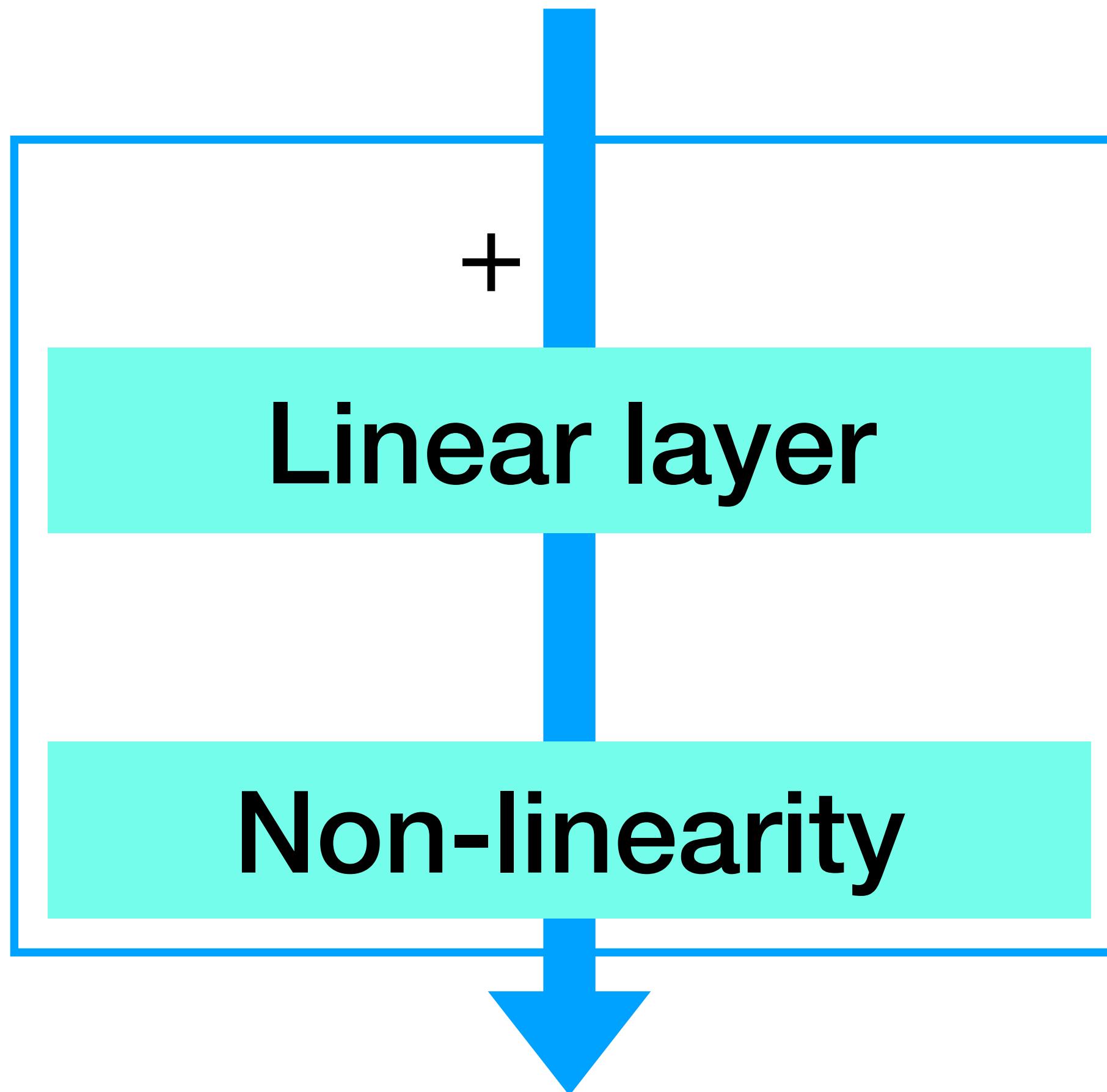




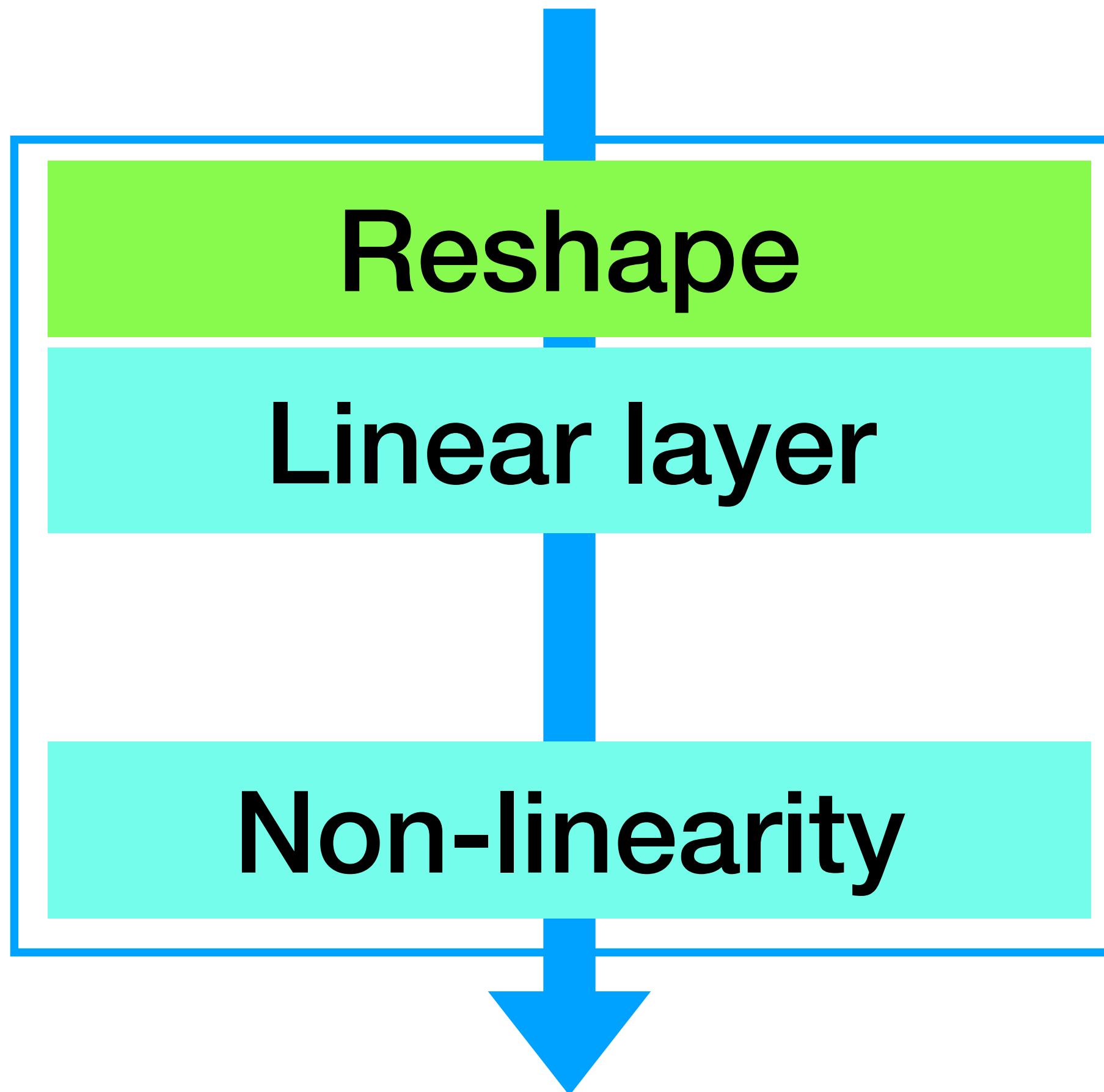
Output



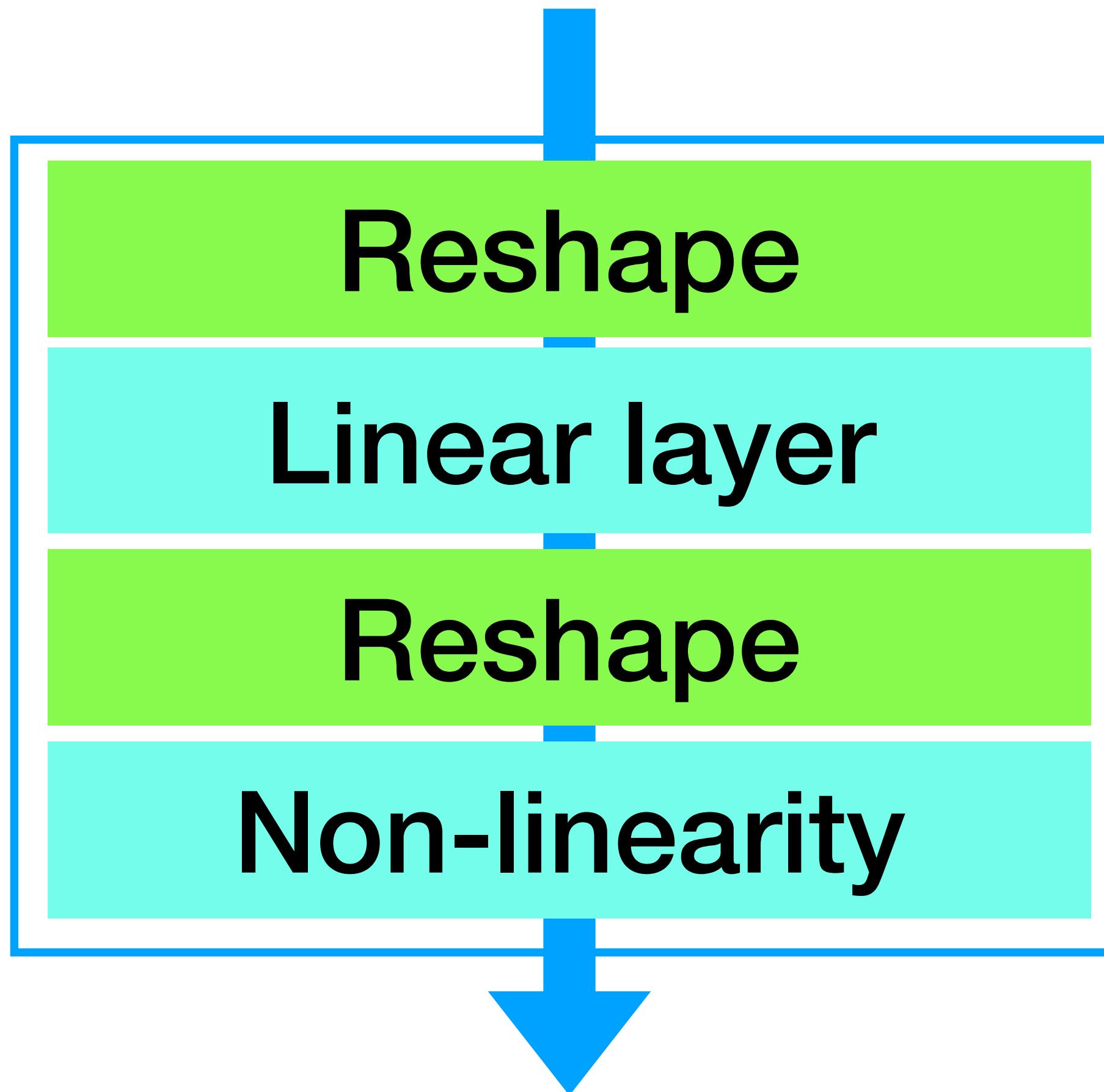
Convolutional layers



Convolutional layers

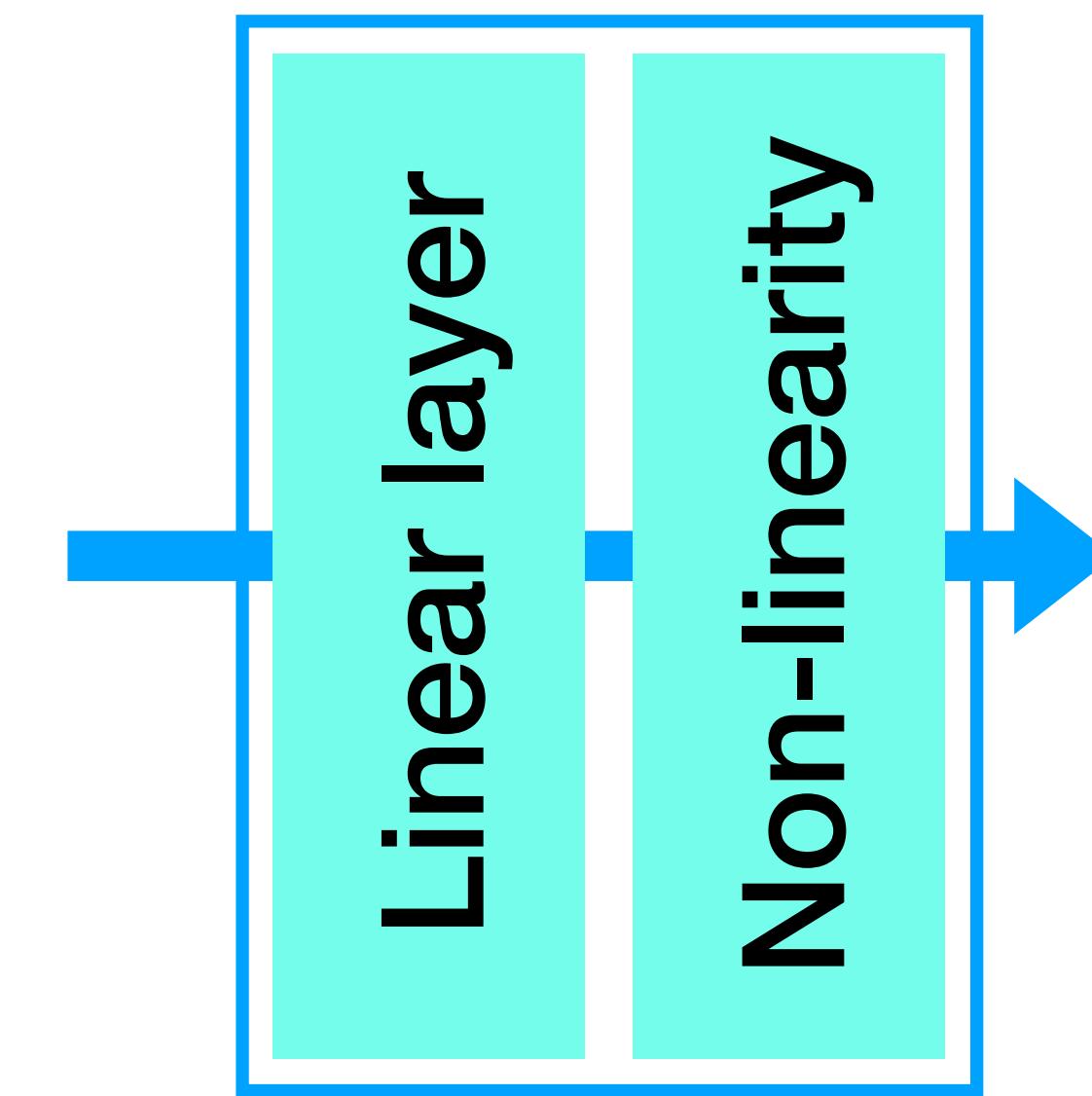


Convolutional layers

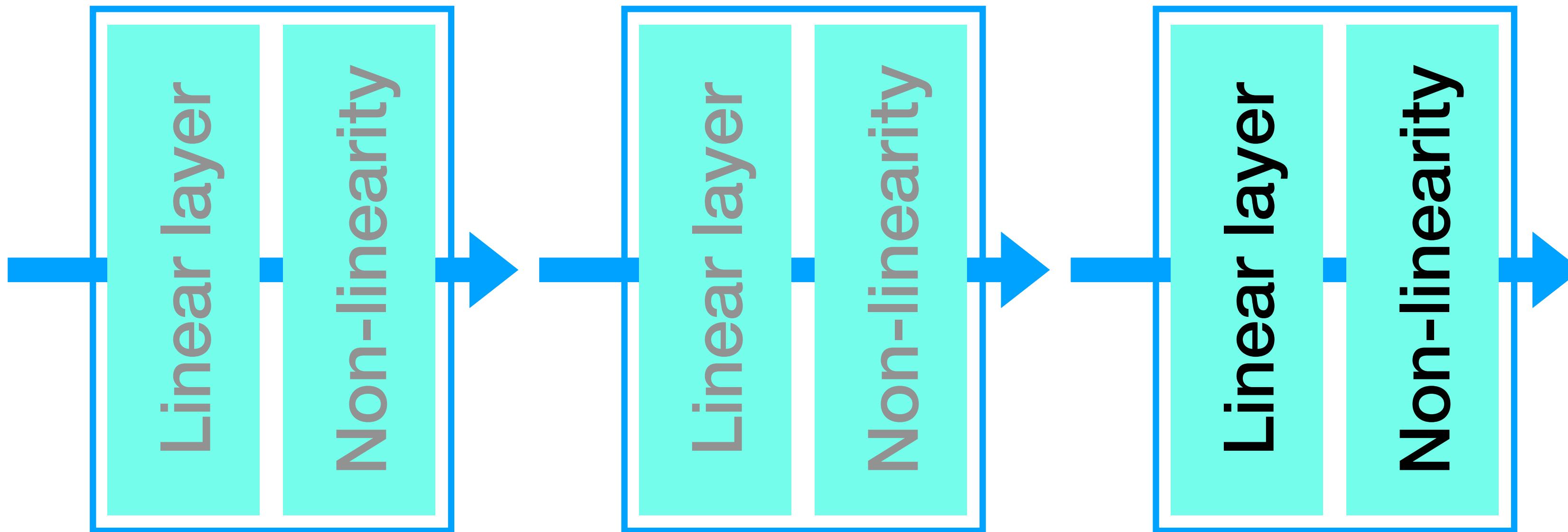


Recurrent layers

Recurrent layers

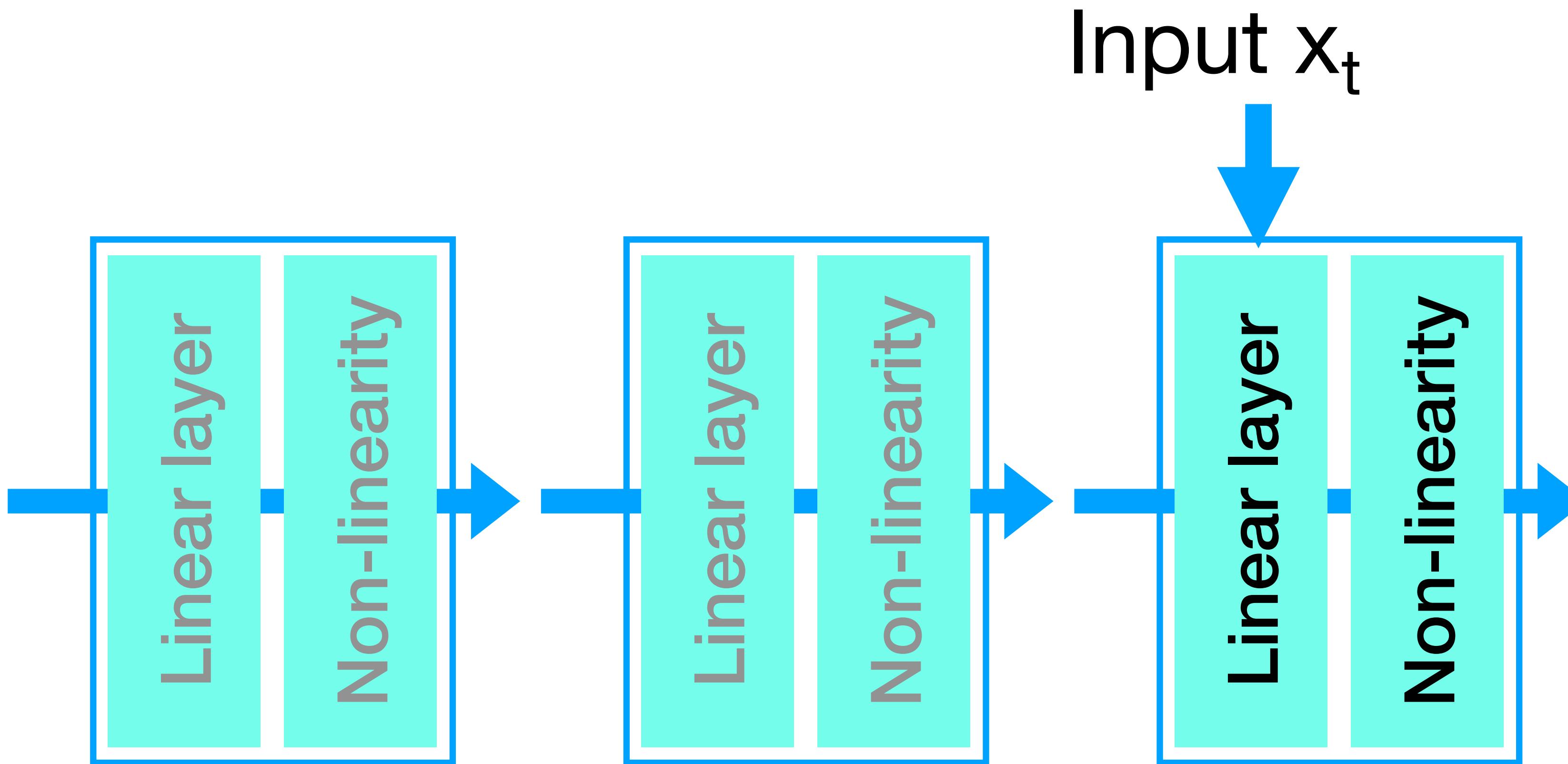


Recurrent layers



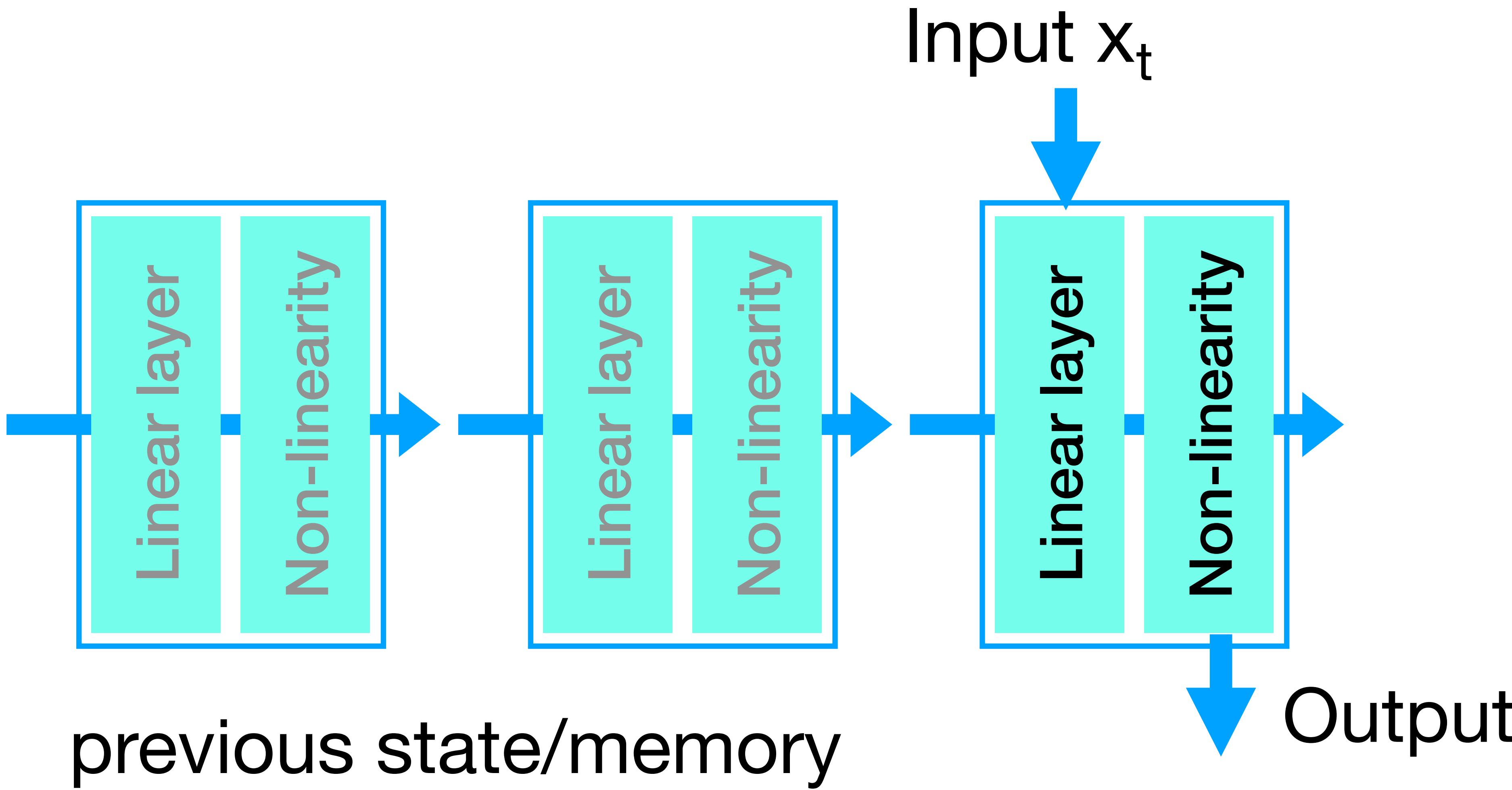
previous state/memory

Recurrent layers

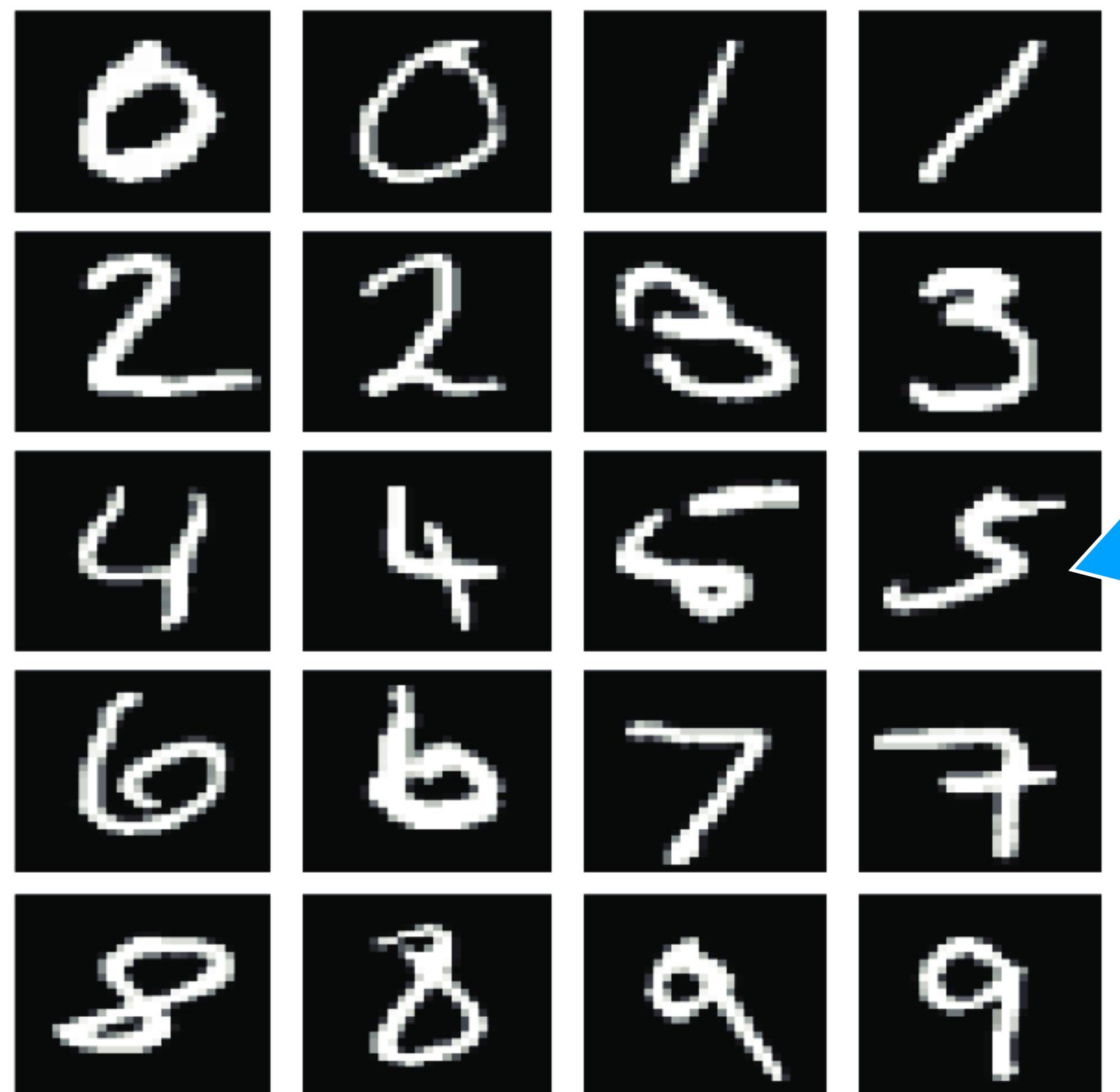


previous state/memory

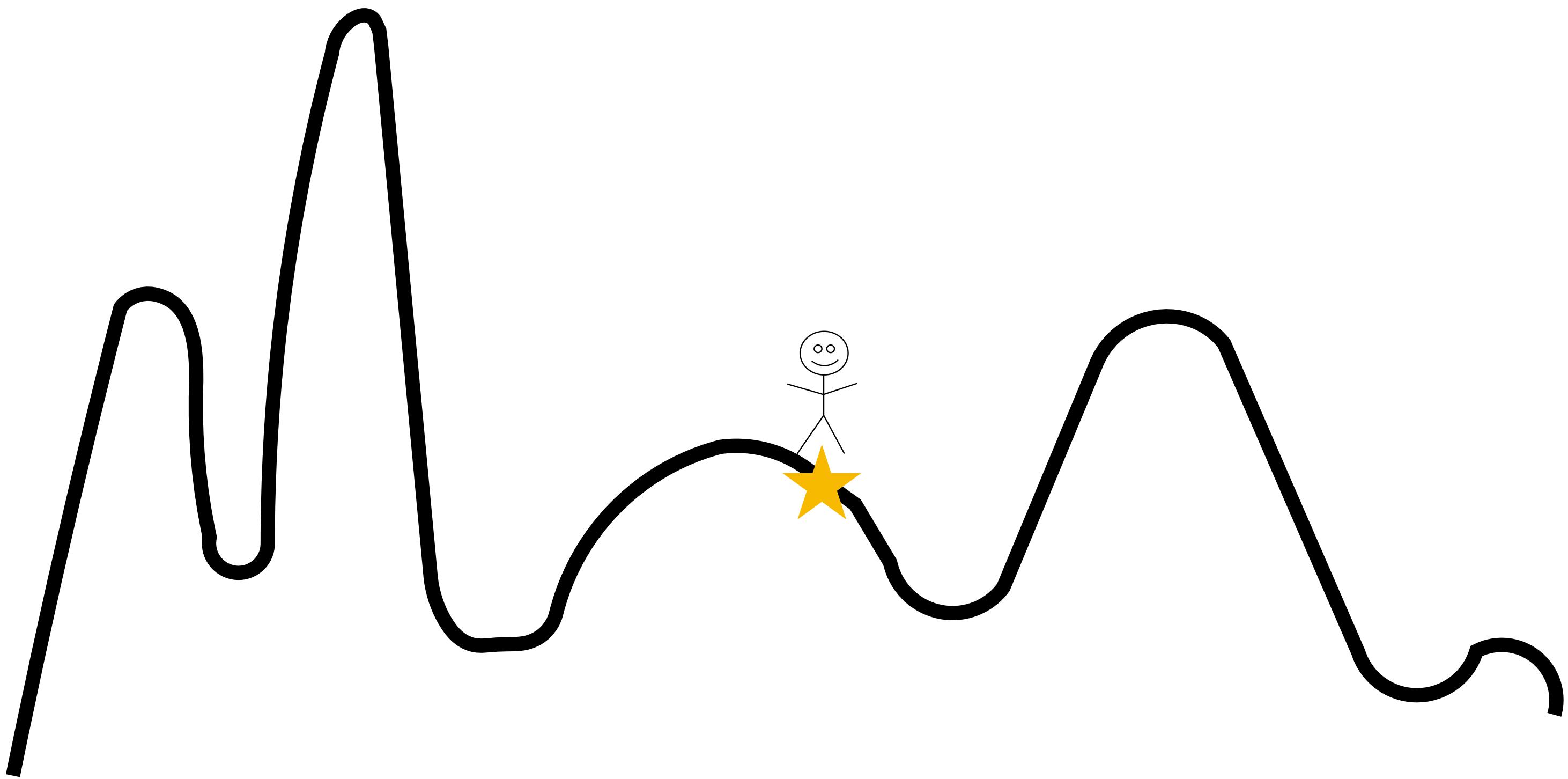
Recurrent layers

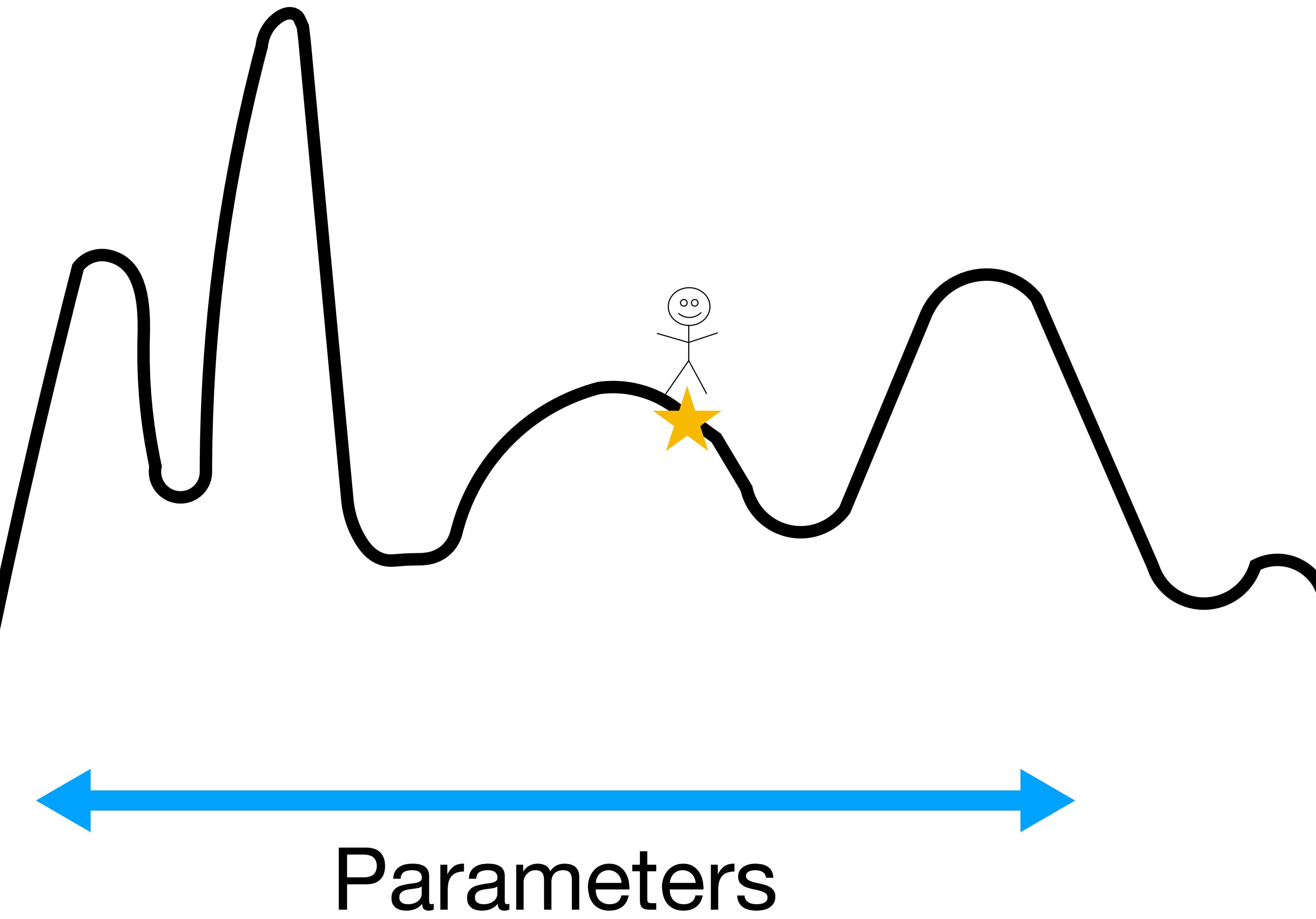


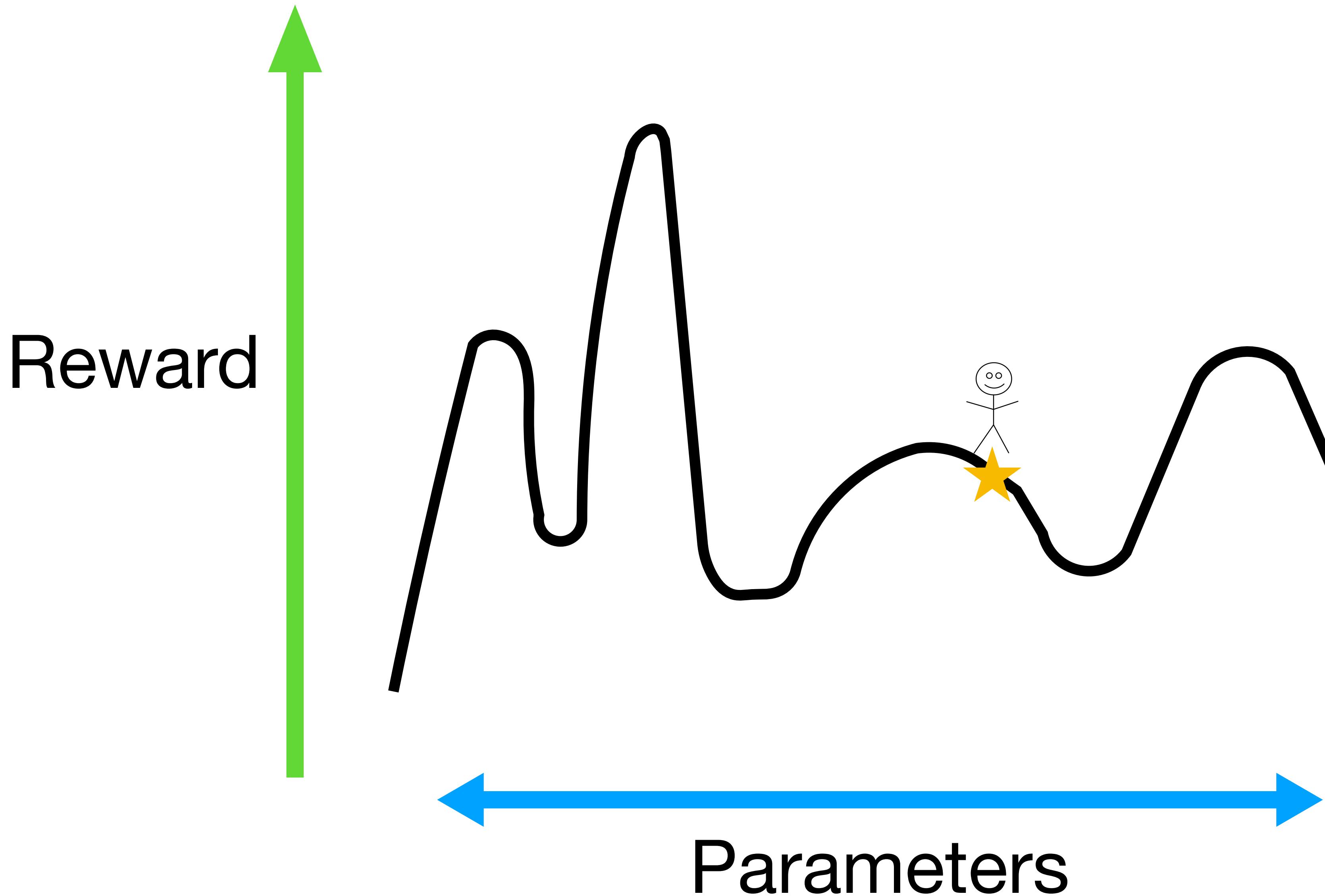
II. How train?

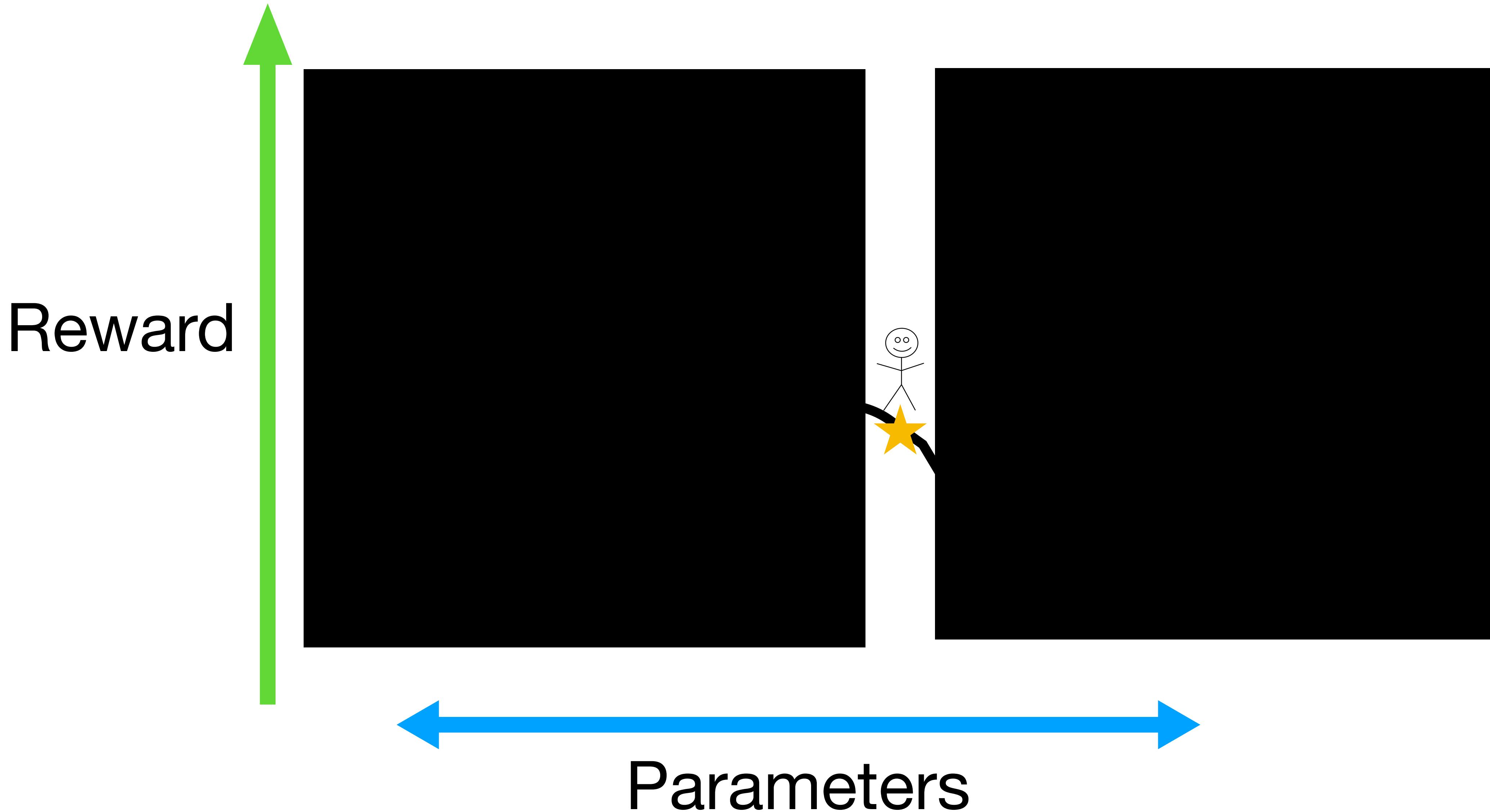


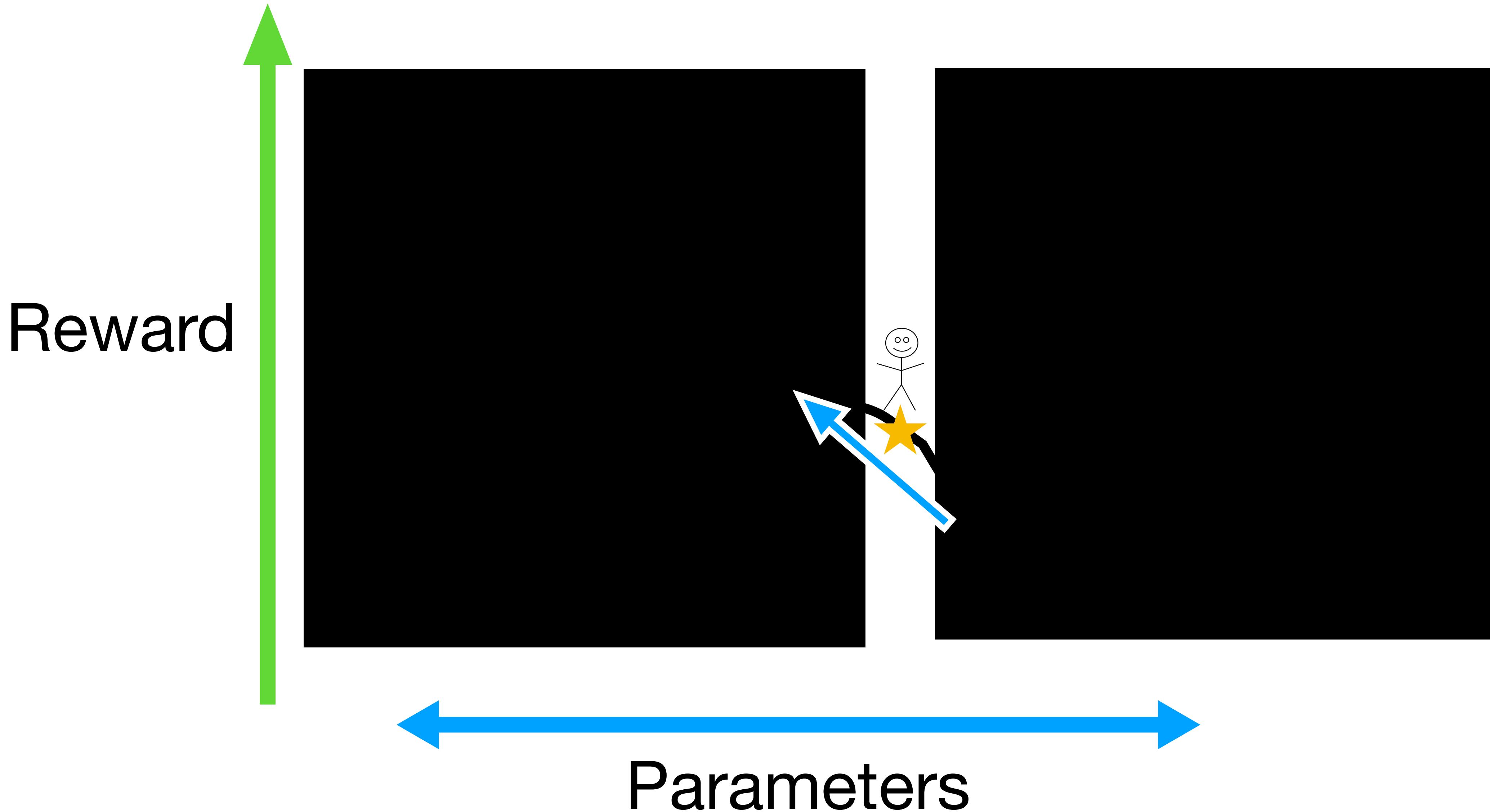
Metric:
~Prediction error

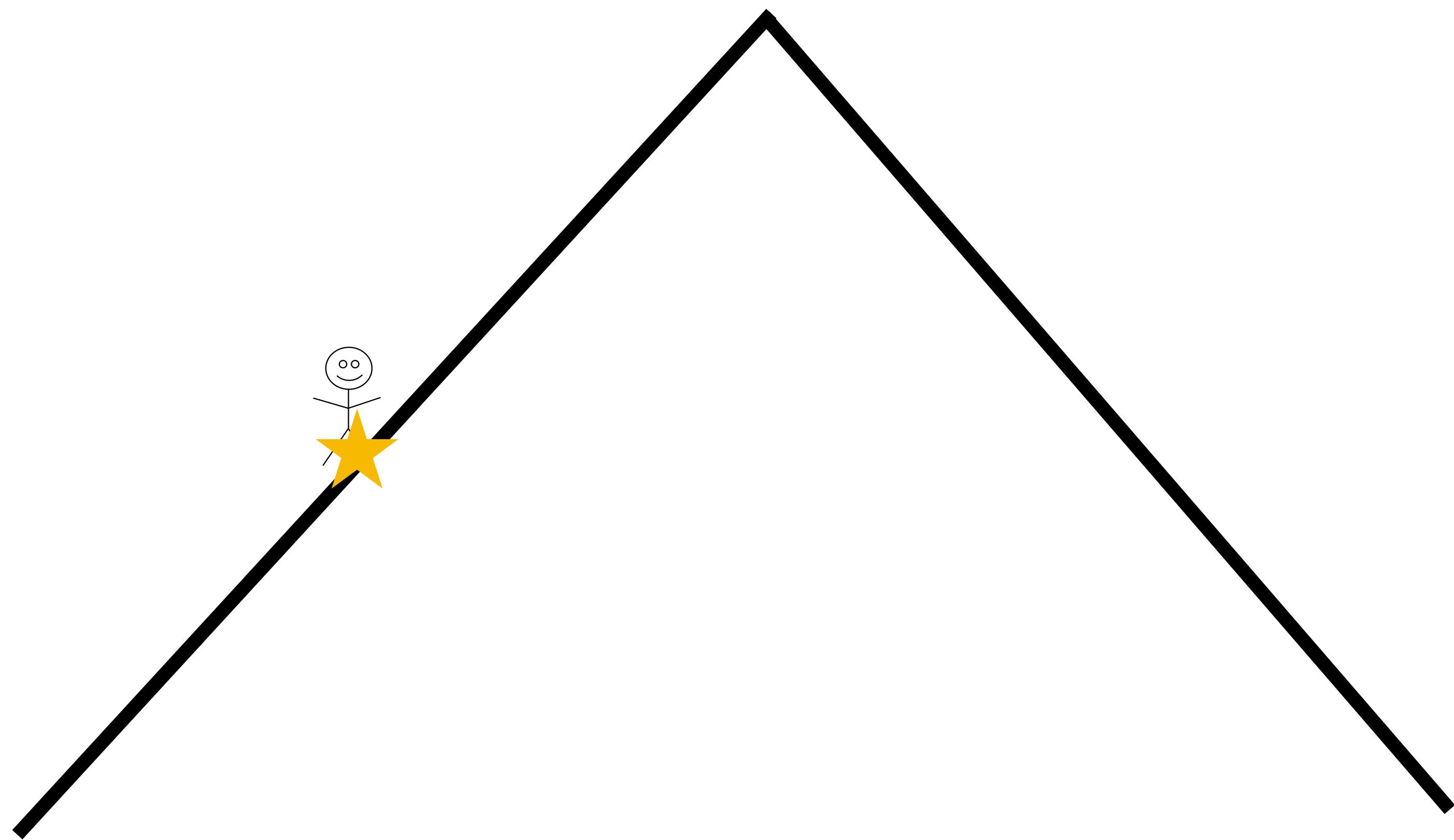


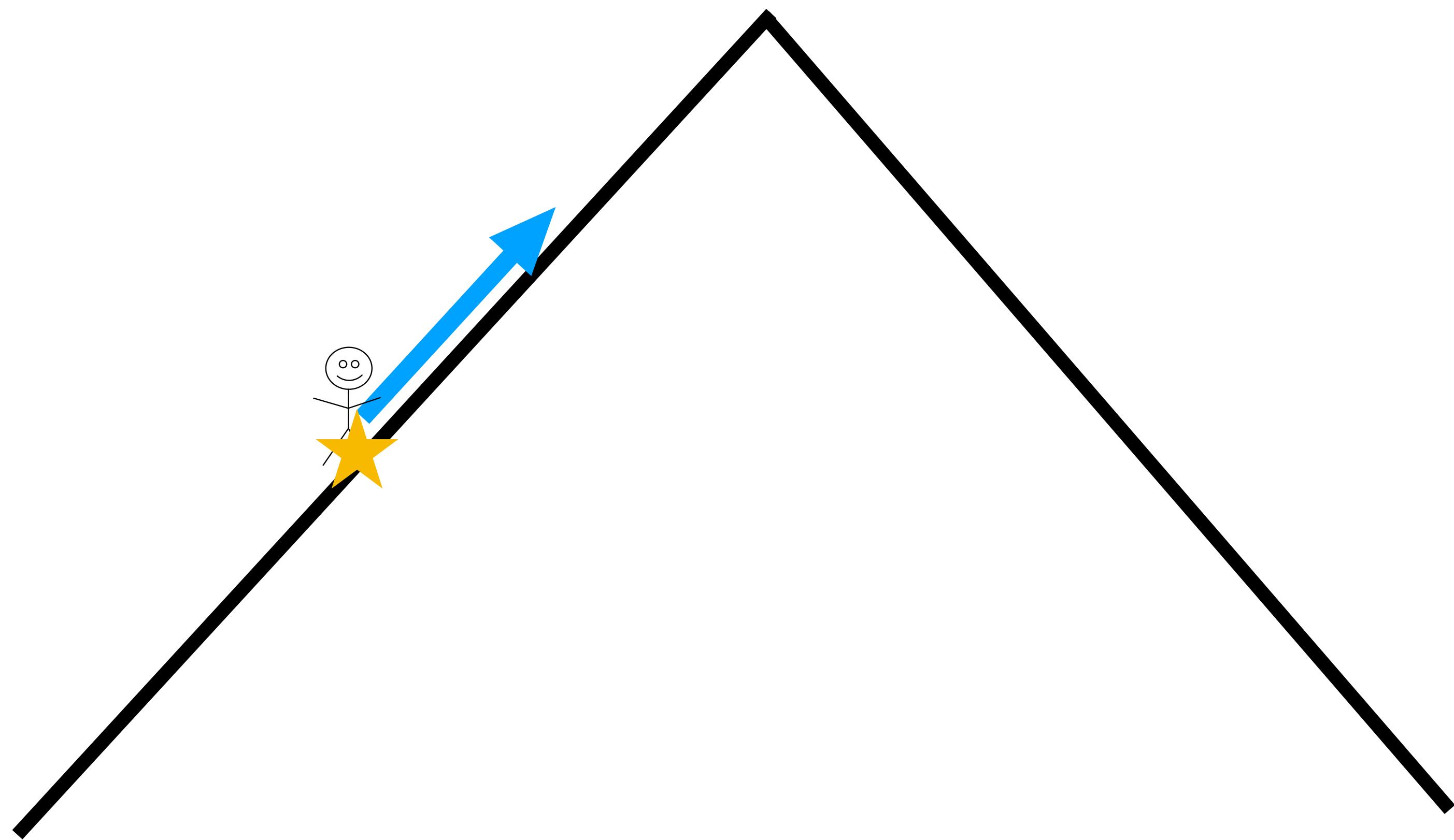


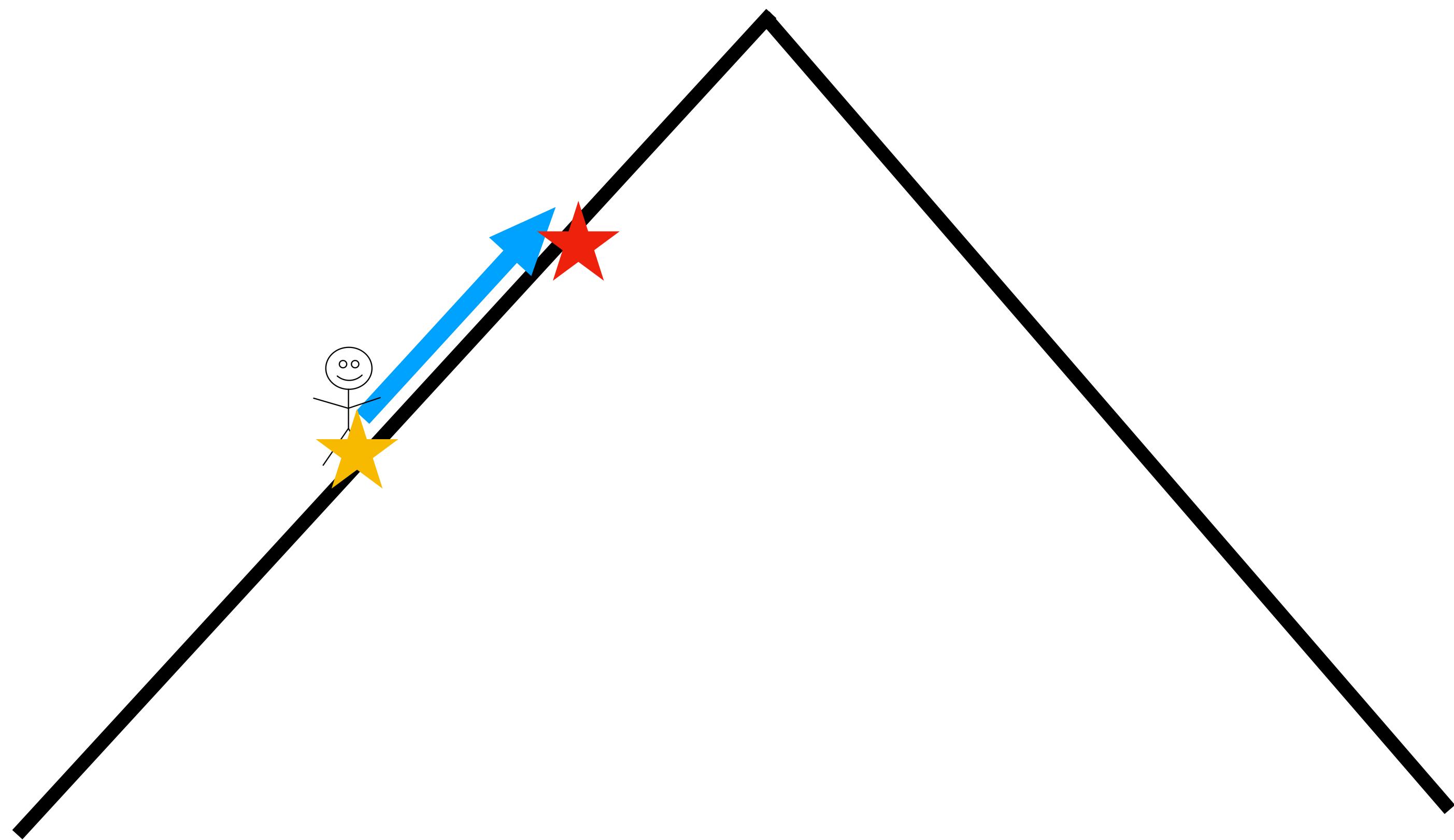


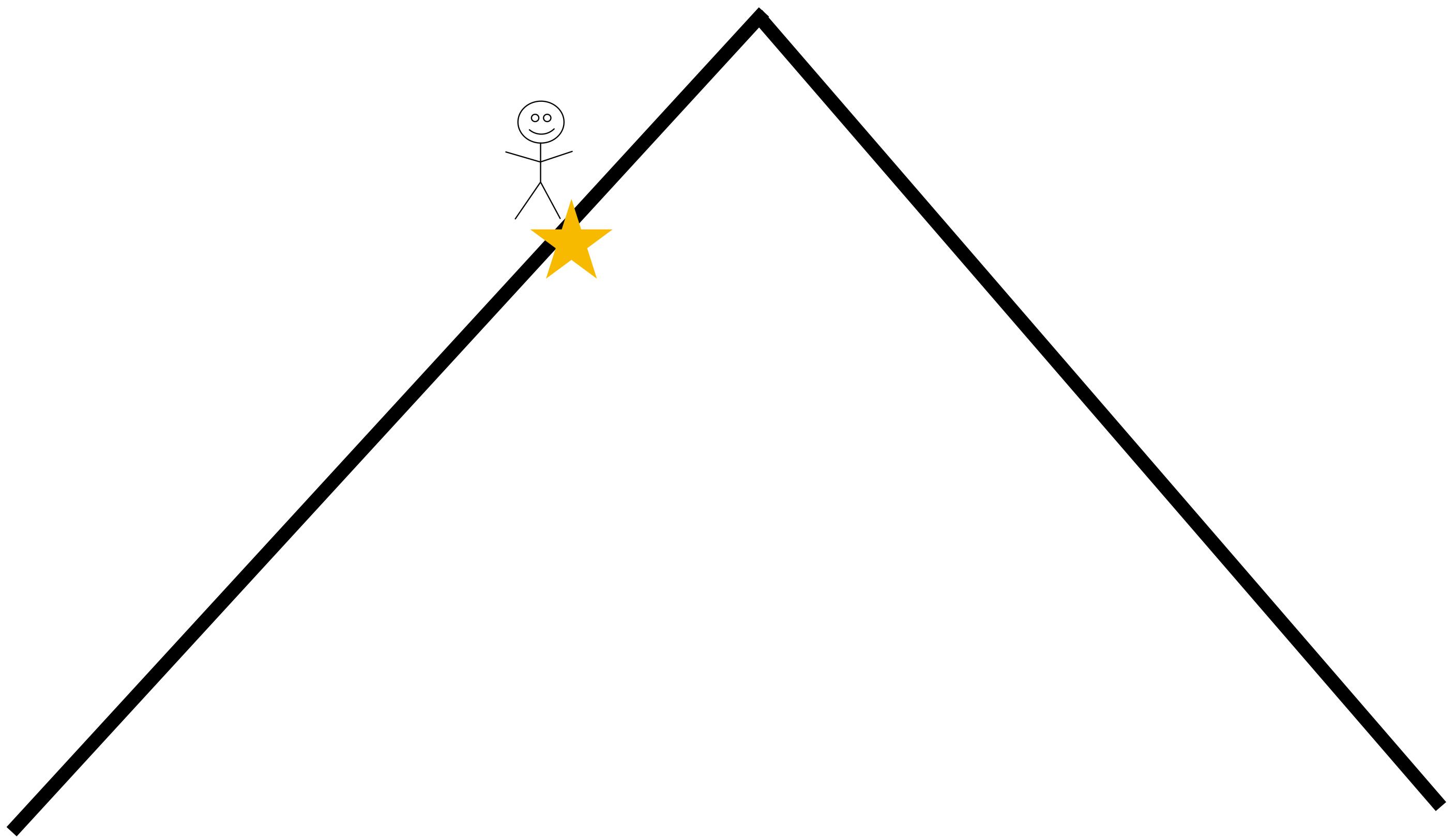


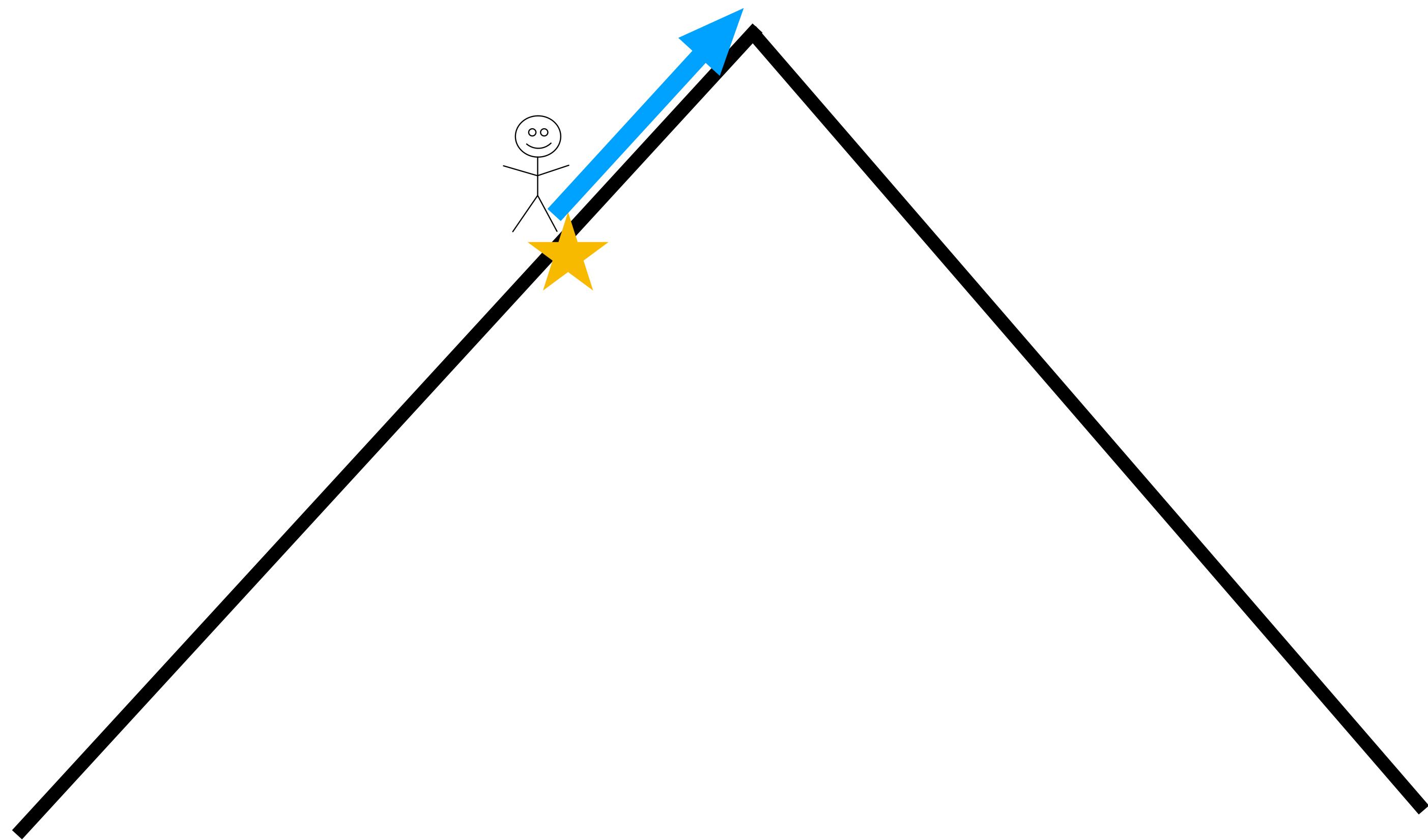


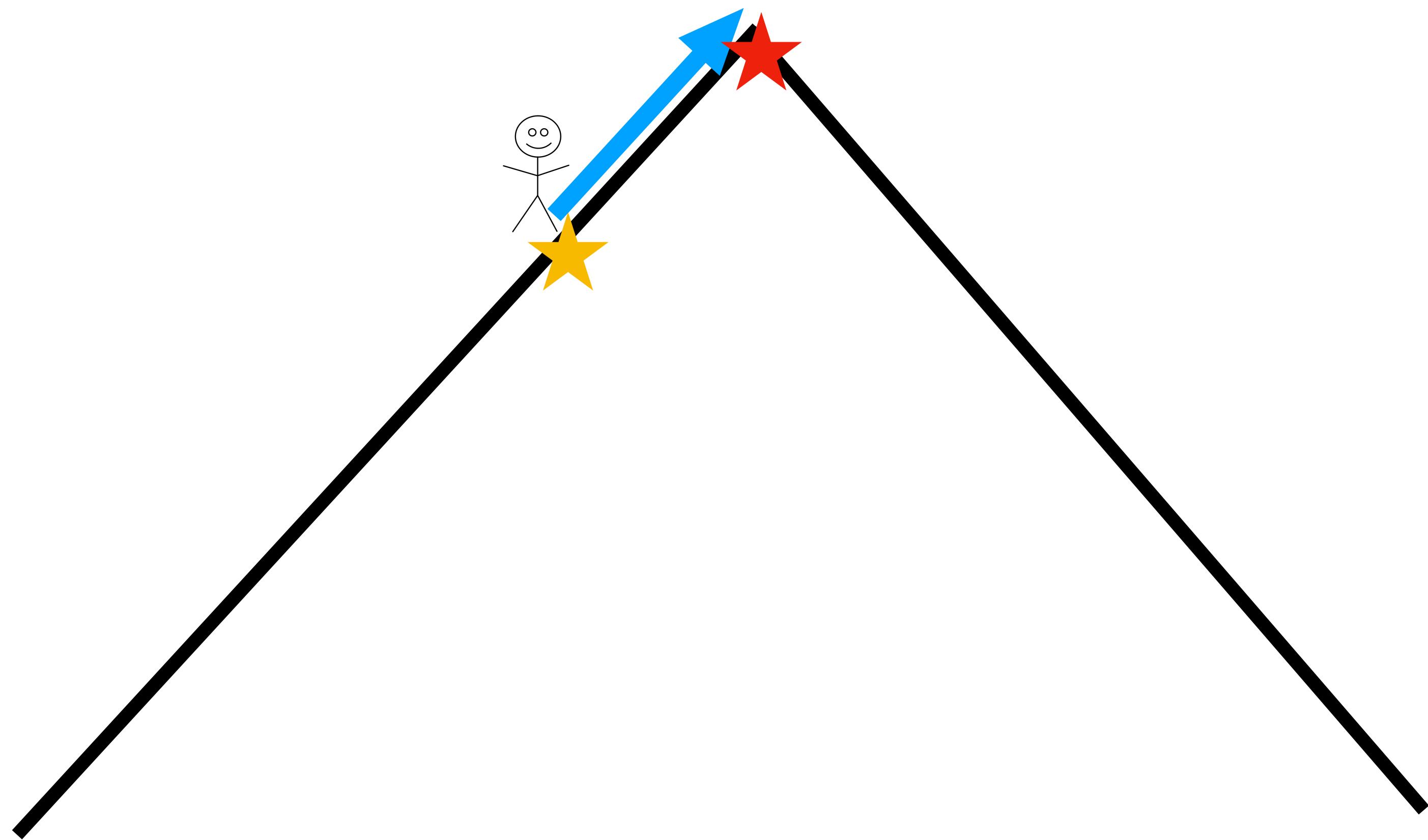


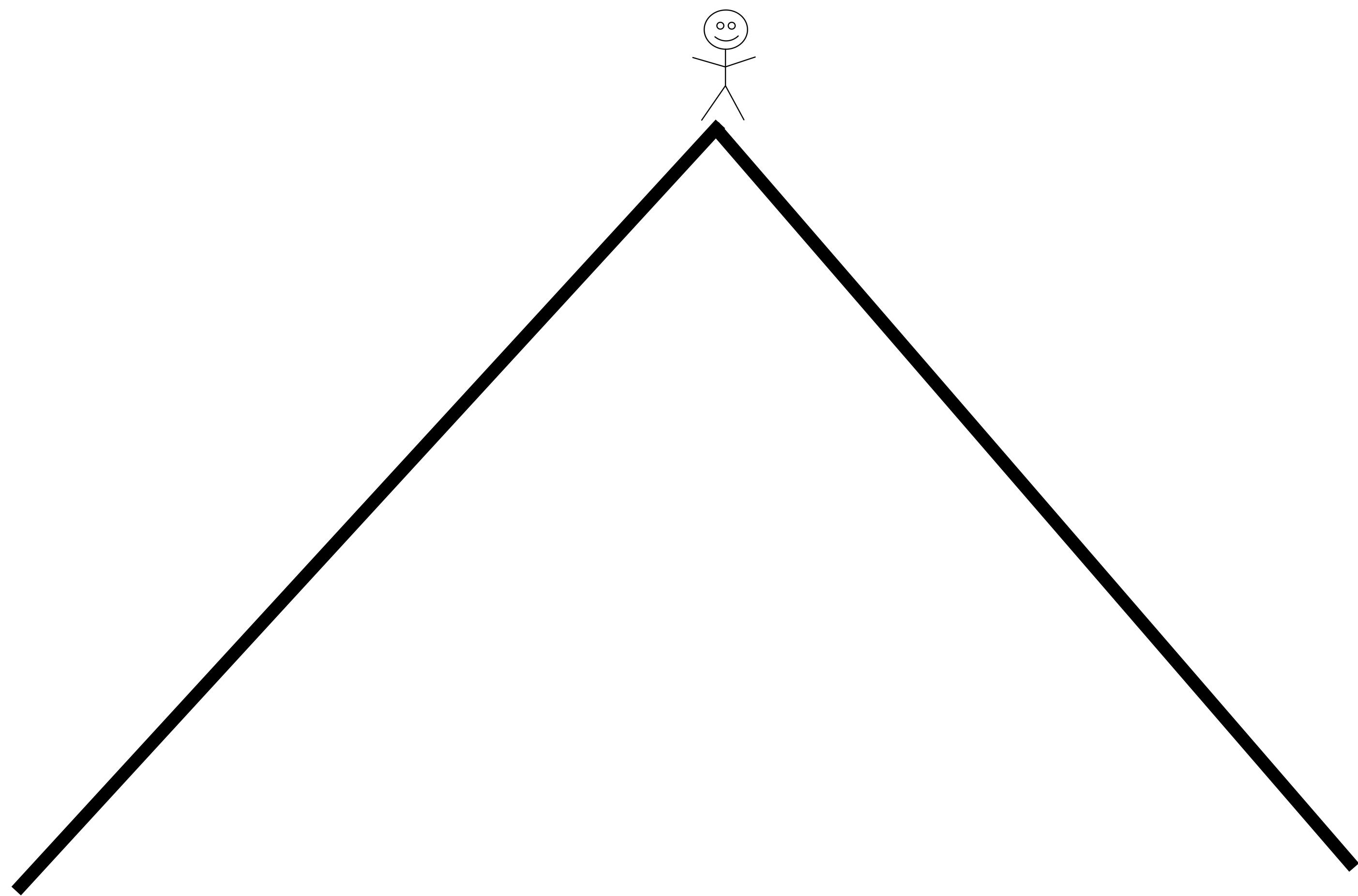


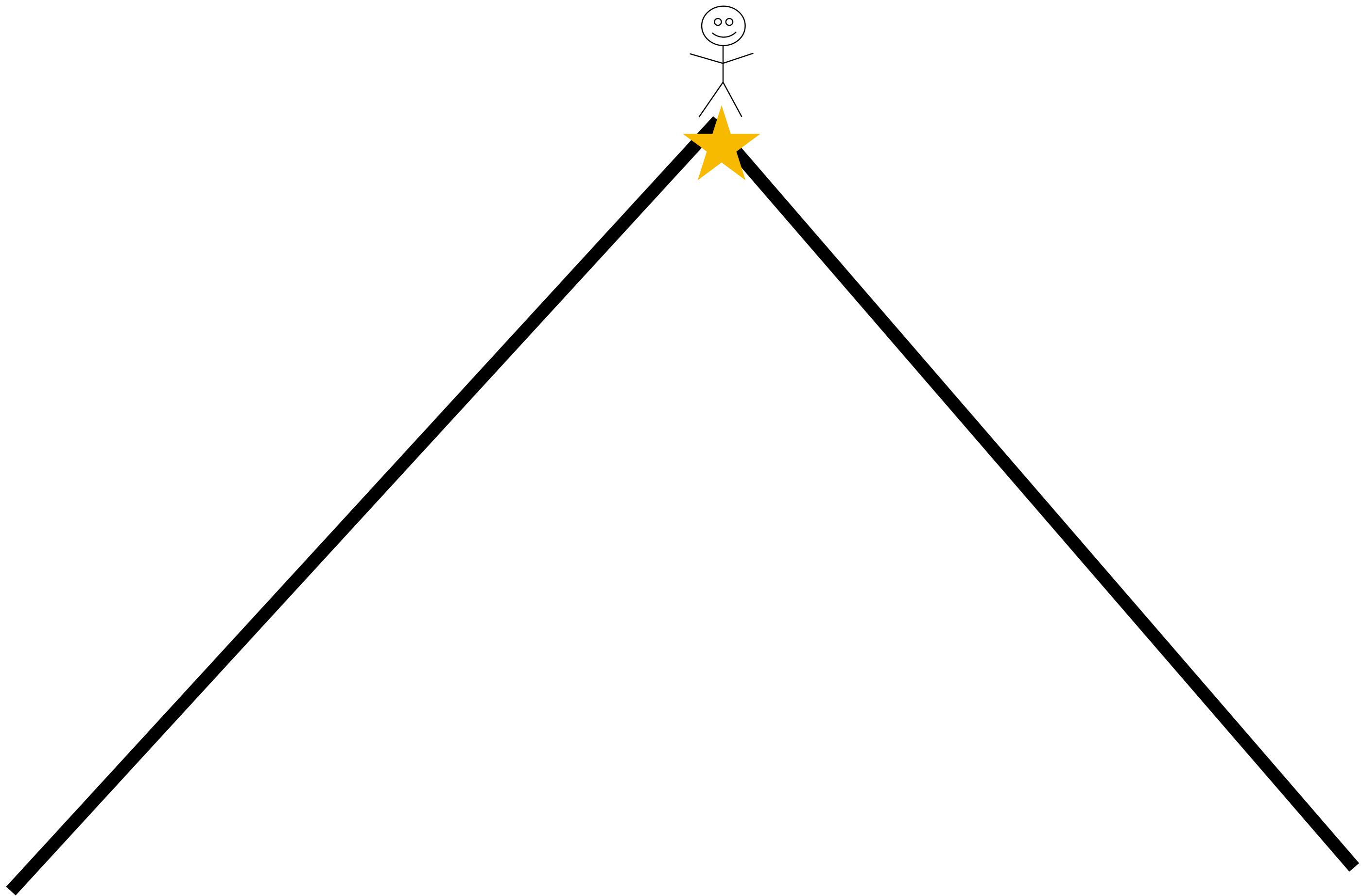


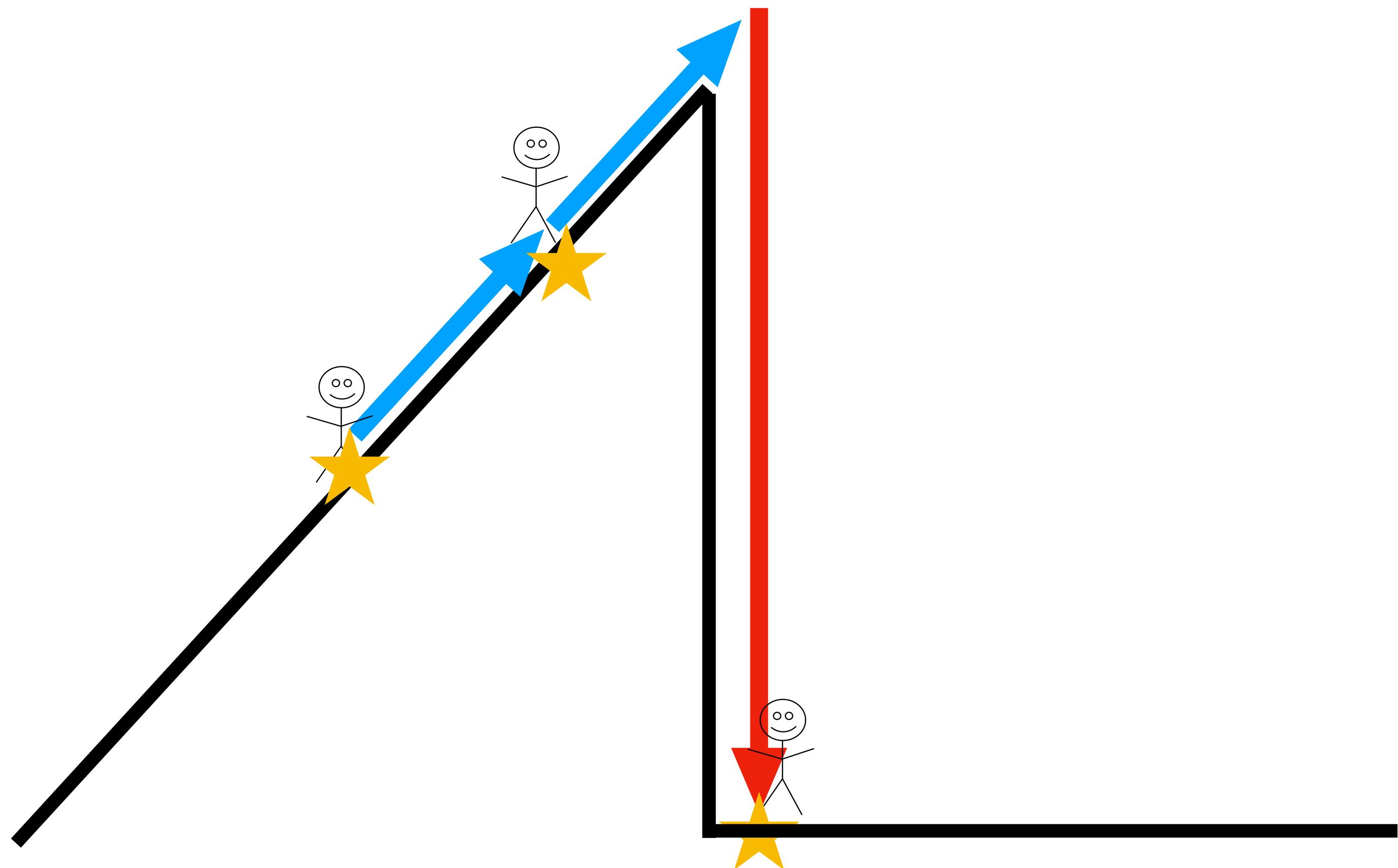


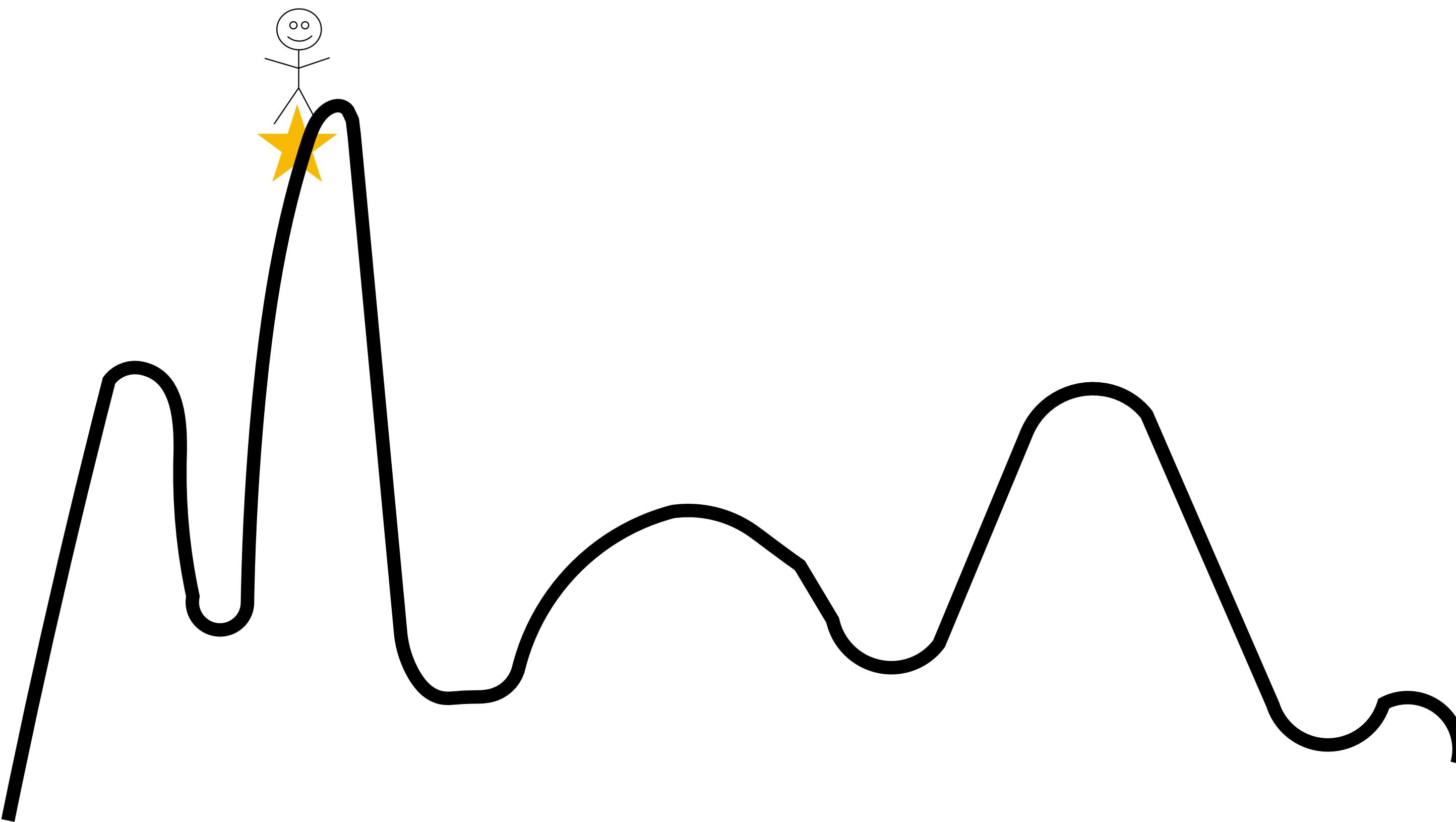


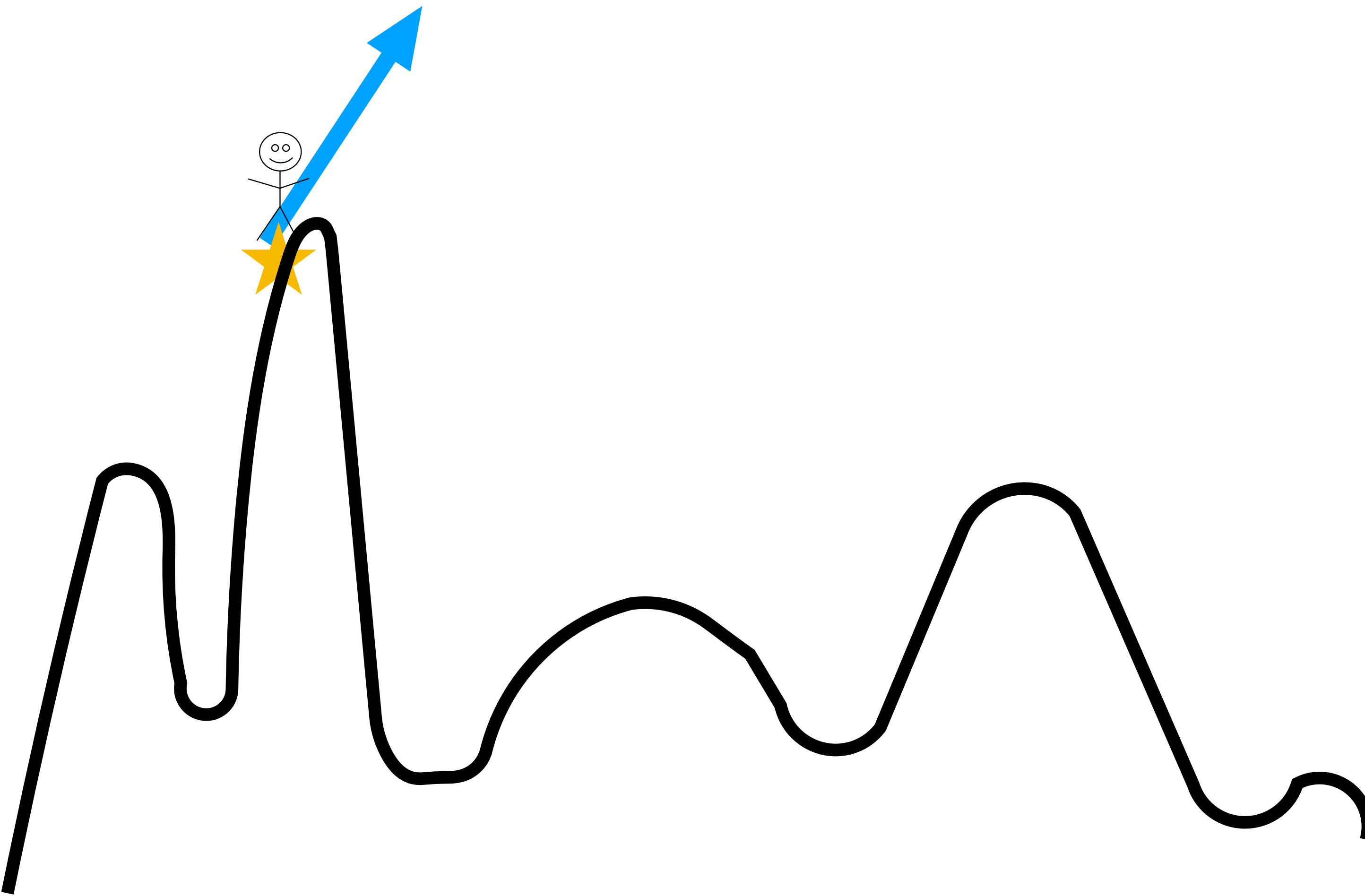


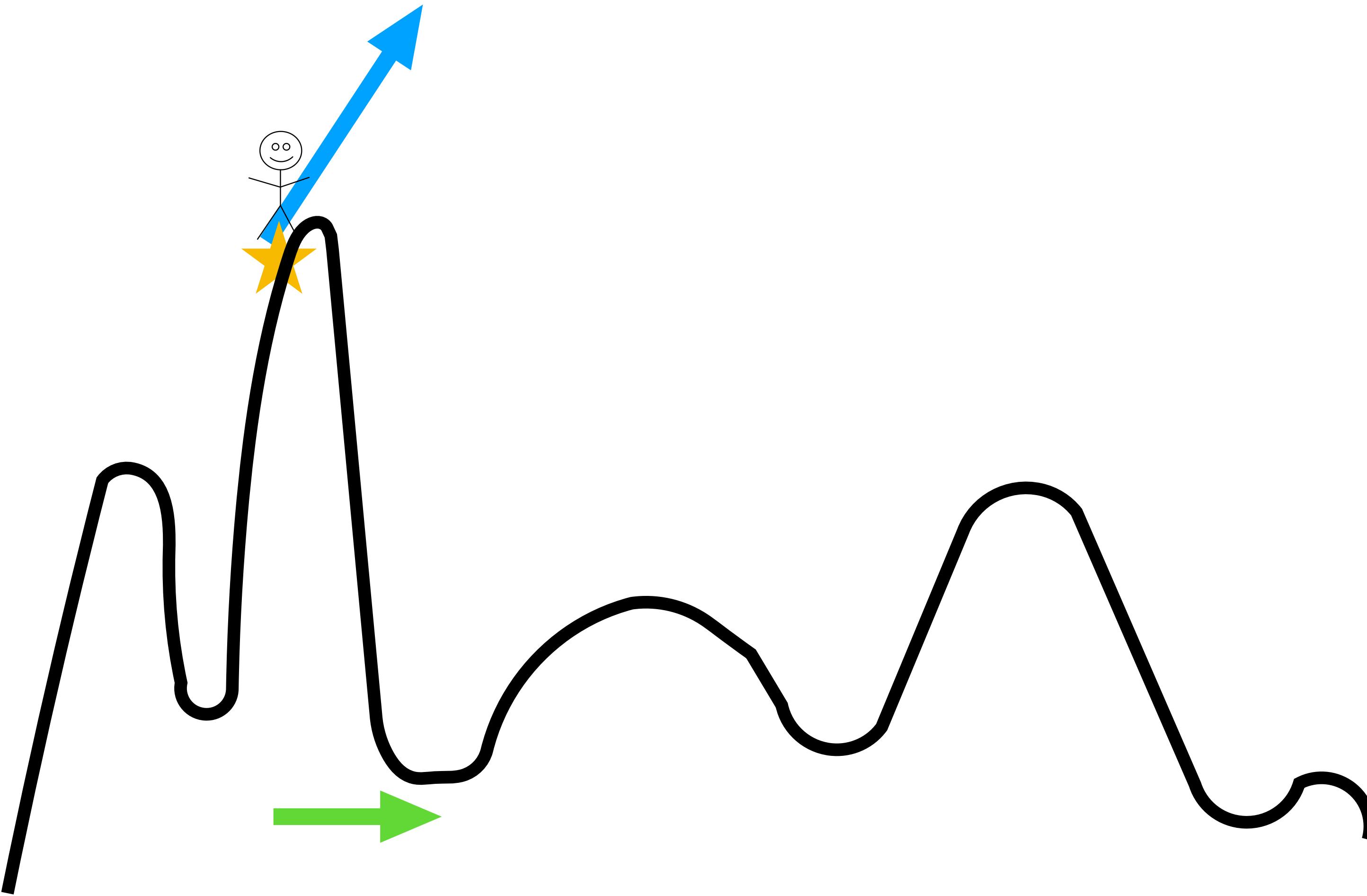


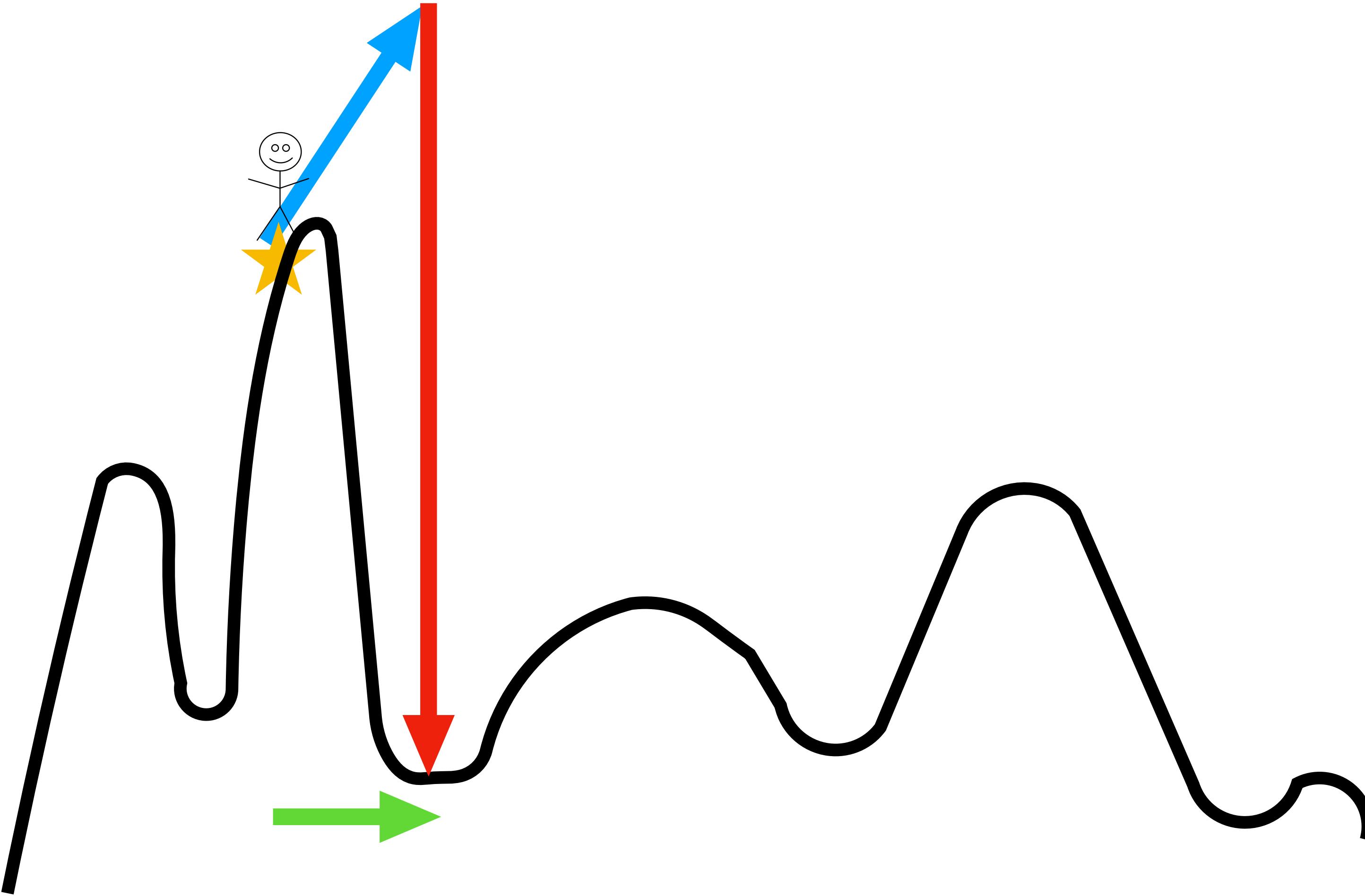


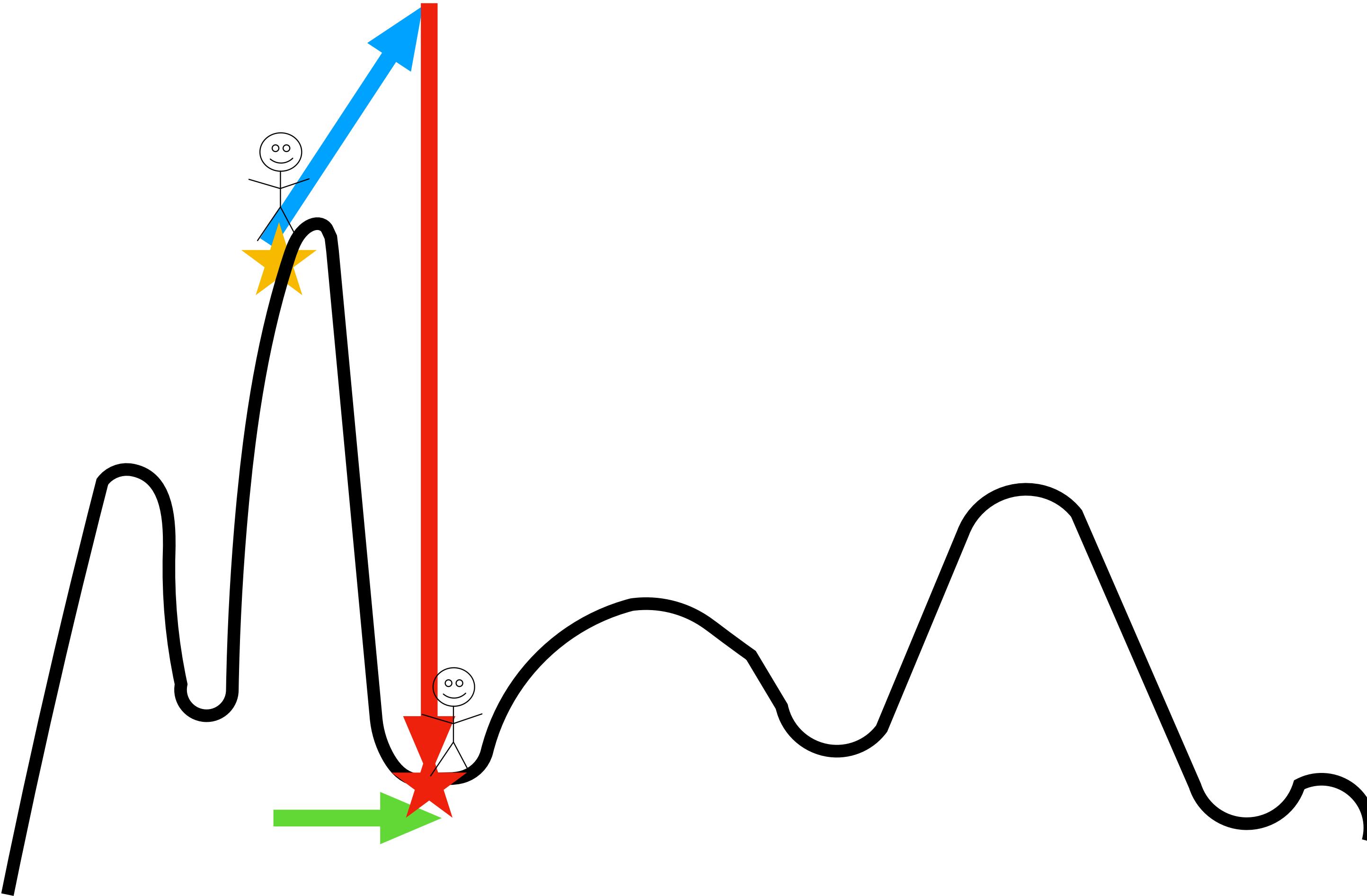


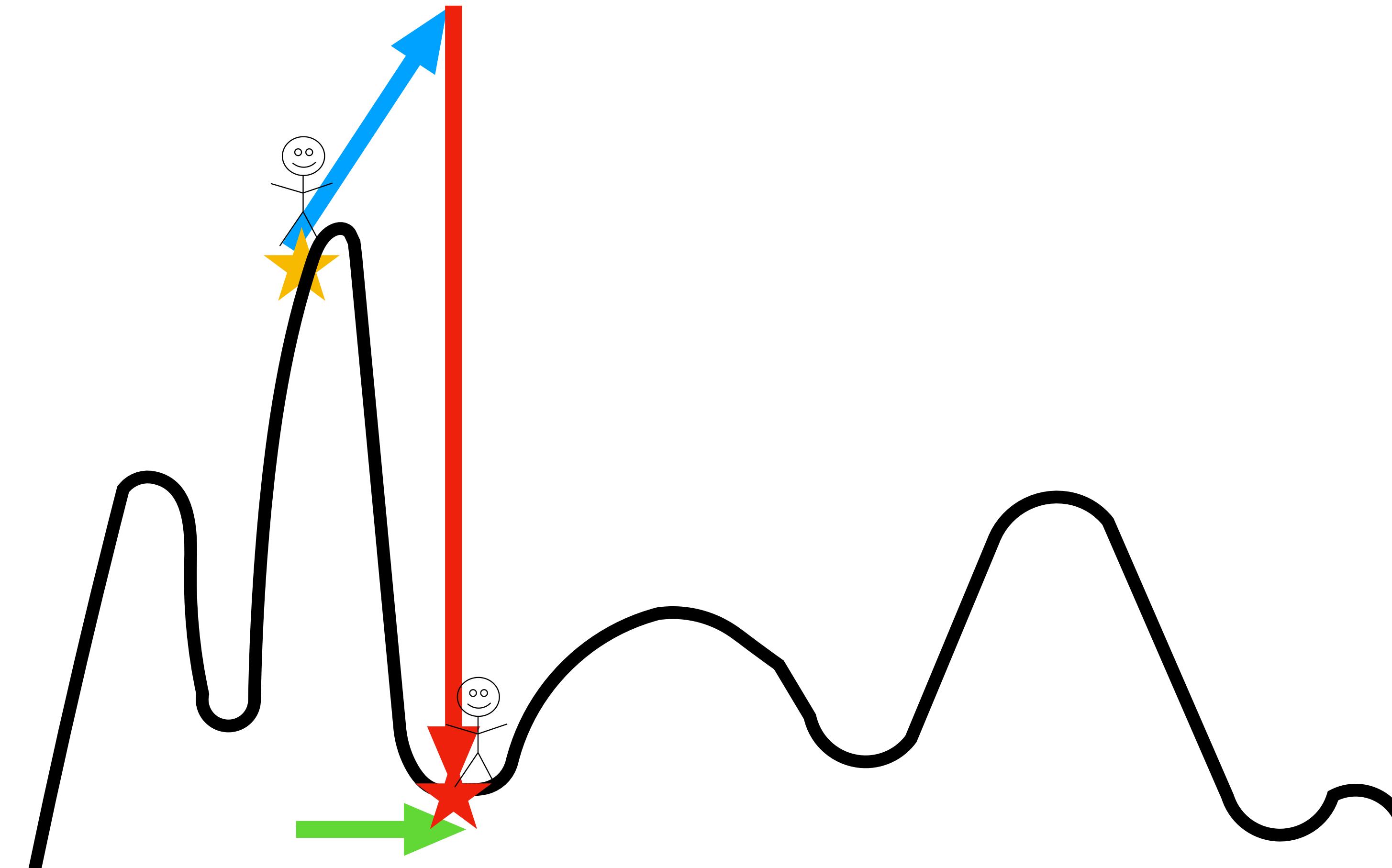












step size x gradient

Gradient descent

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

Gradient descent

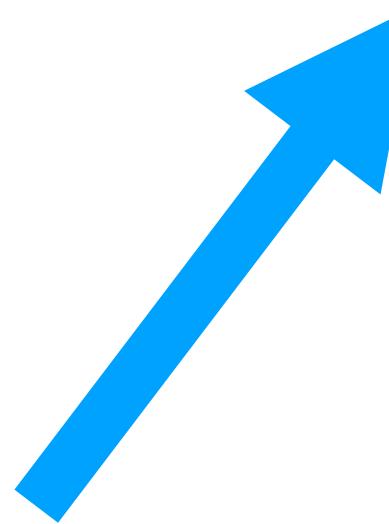
$$W = W - \alpha \nabla_W f$$



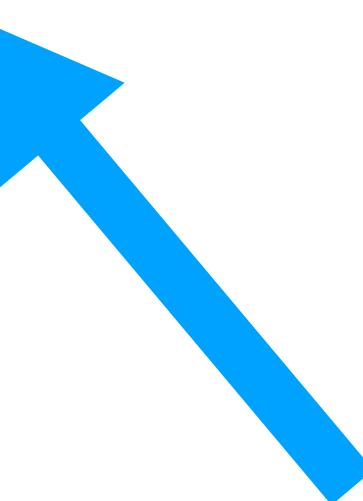
Parameters

Gradient descent

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



Parameters



Gradient

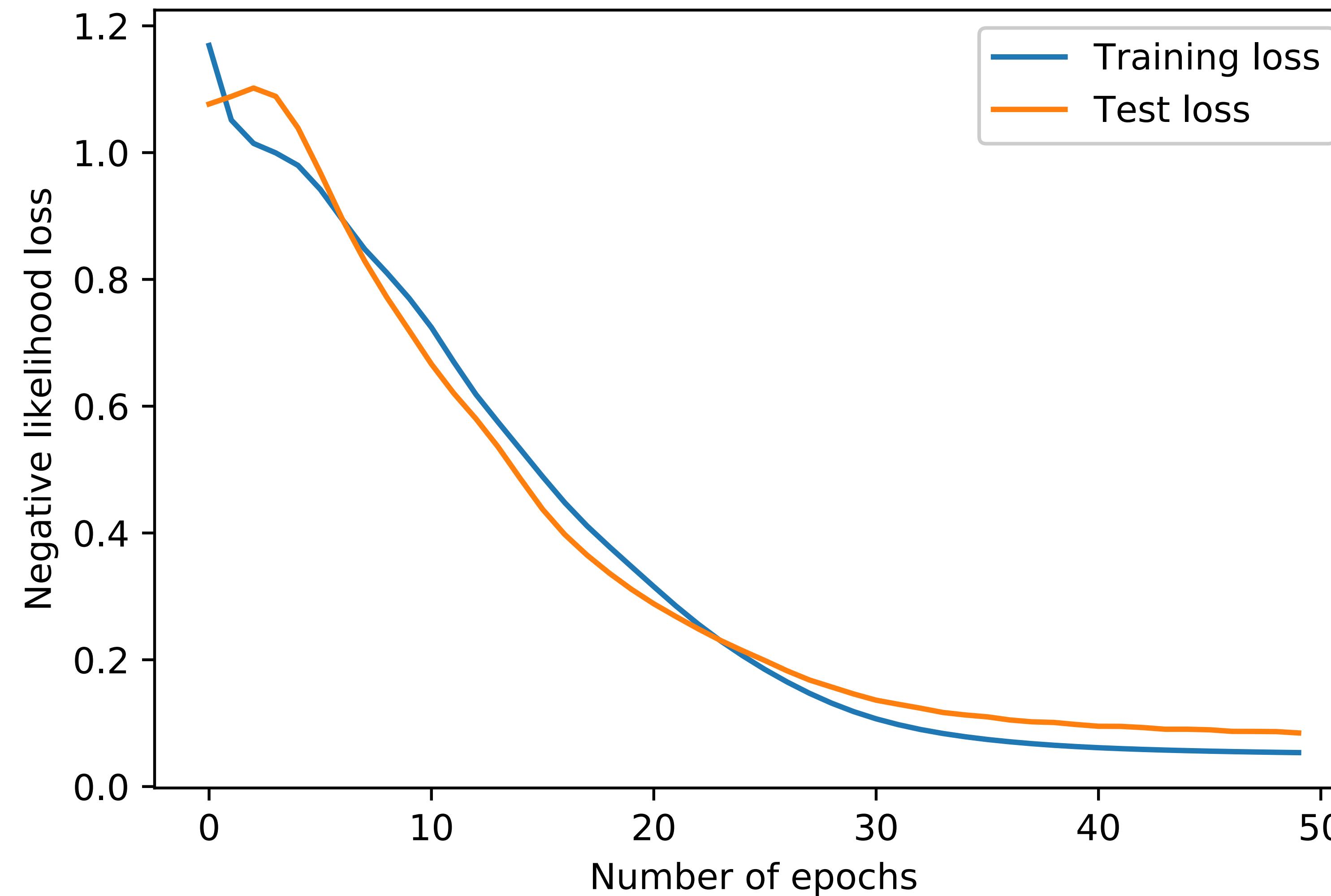
Gradient descent

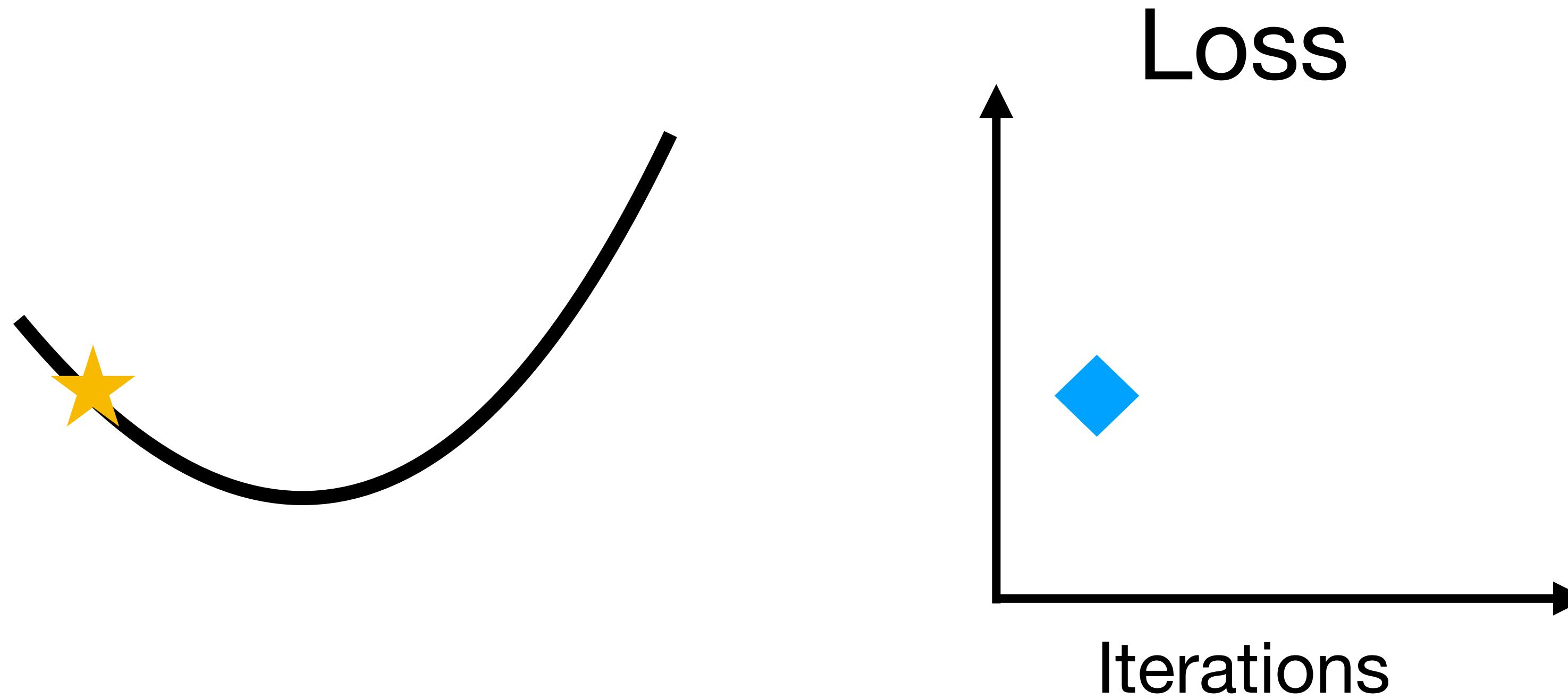
$$W = W - \alpha \nabla_W f$$

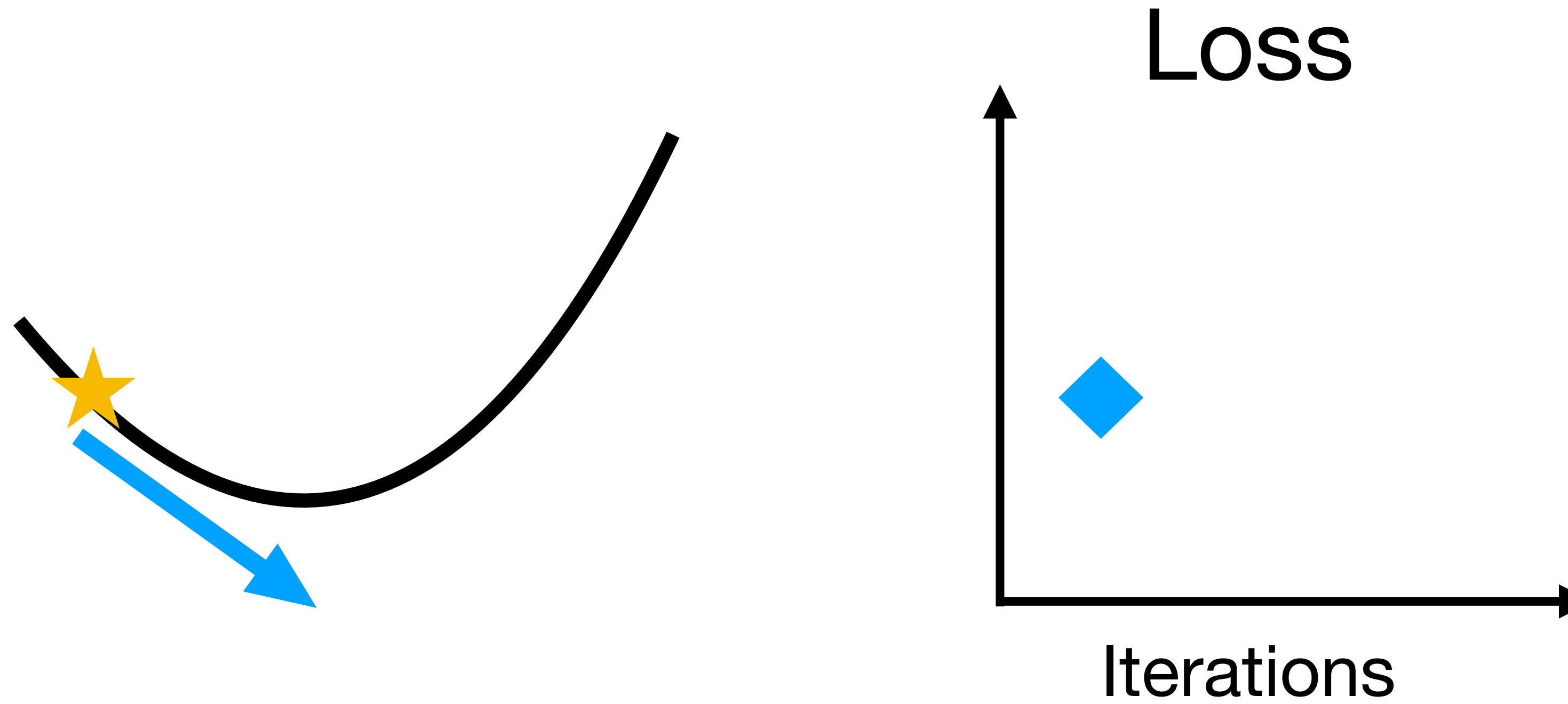
Parameters Learning rate (step size) Gradient

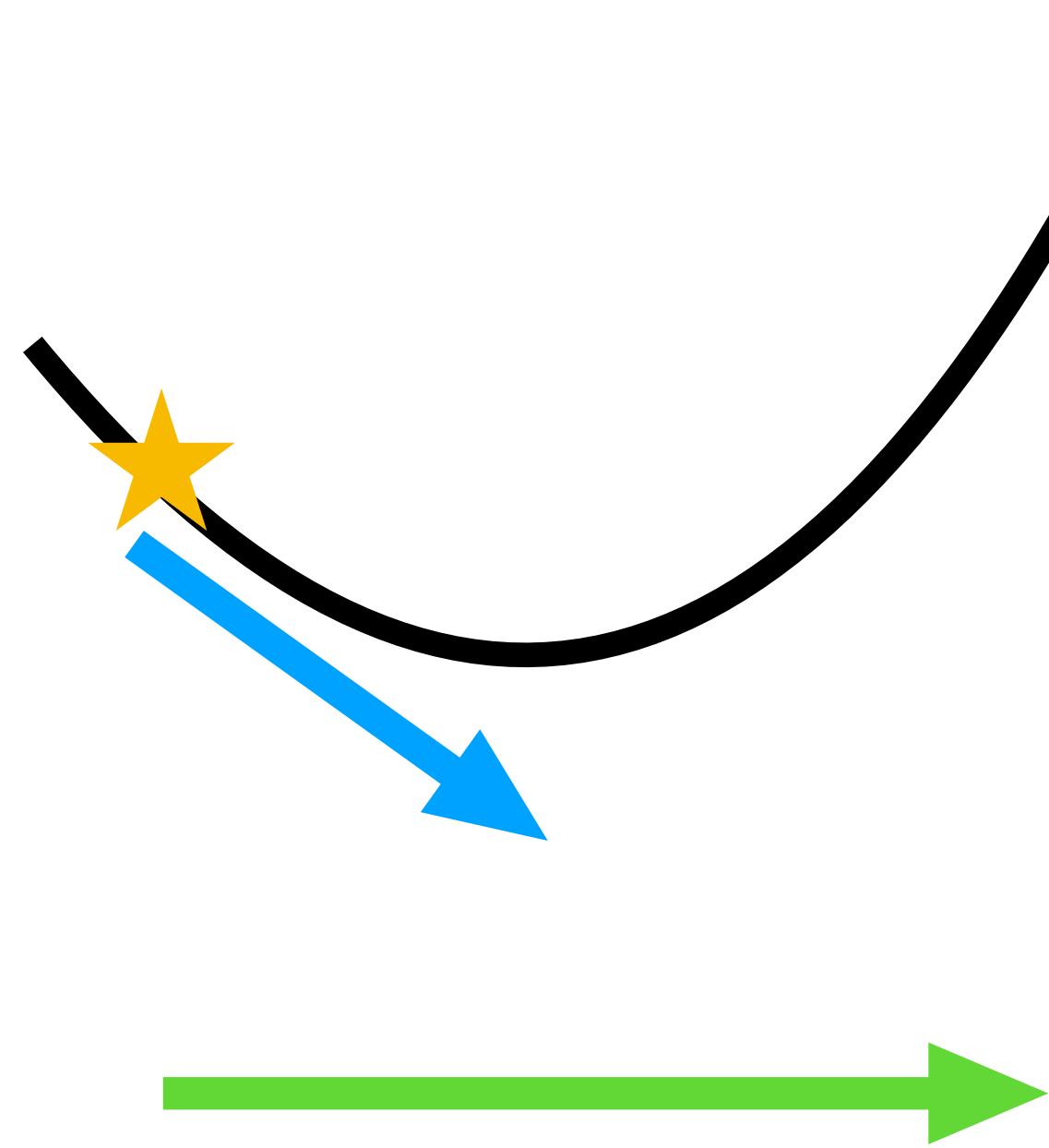
The diagram illustrates the gradient descent update rule. It shows the formula $W = W - \alpha \nabla_W f$. Three blue arrows point to the components of the formula: one arrow points to the variable W , another to the learning rate α , and a third to the gradient term $\nabla_W f$. The gradient term is highlighted with a blue rectangular box.

Good loss curves

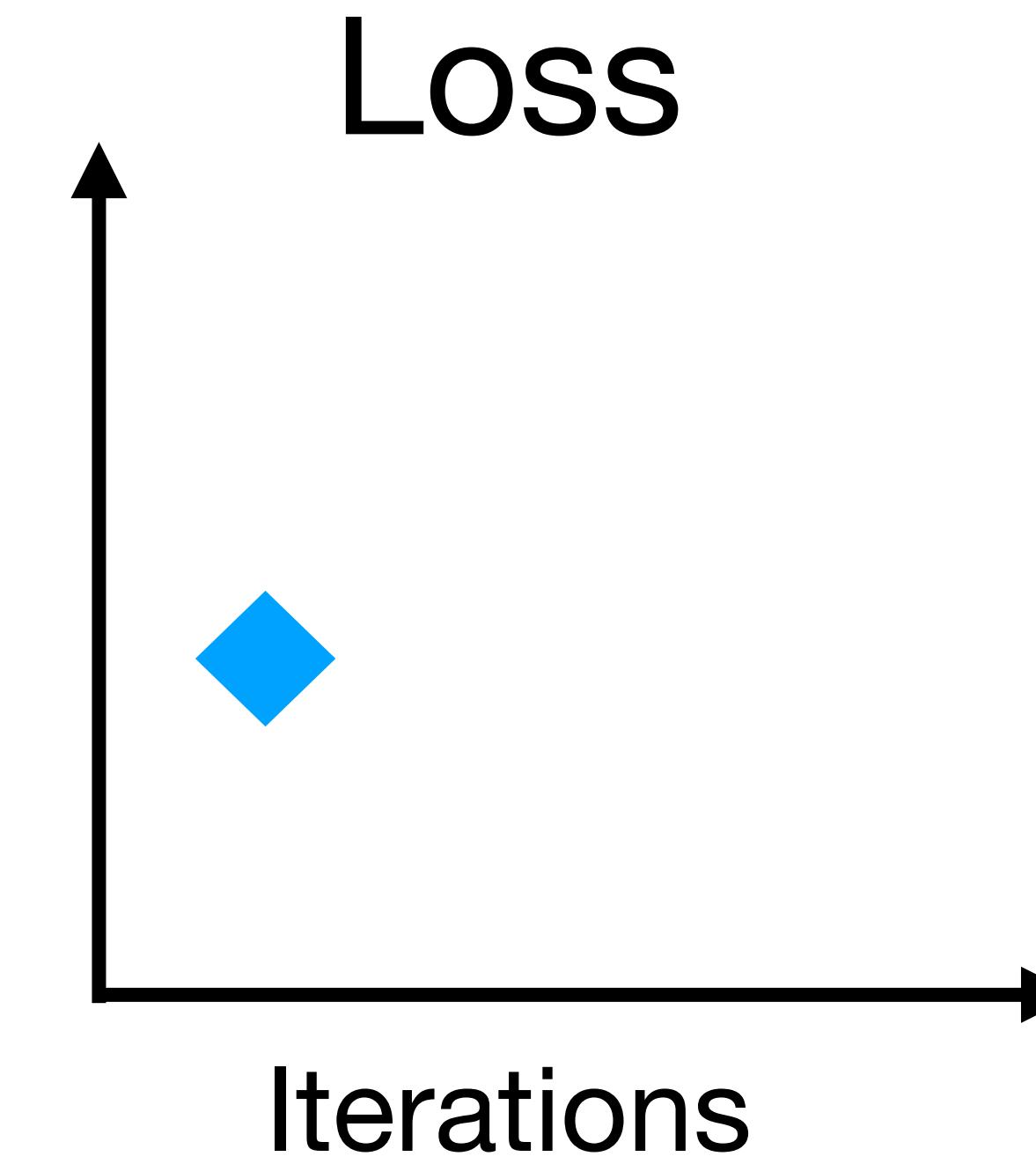


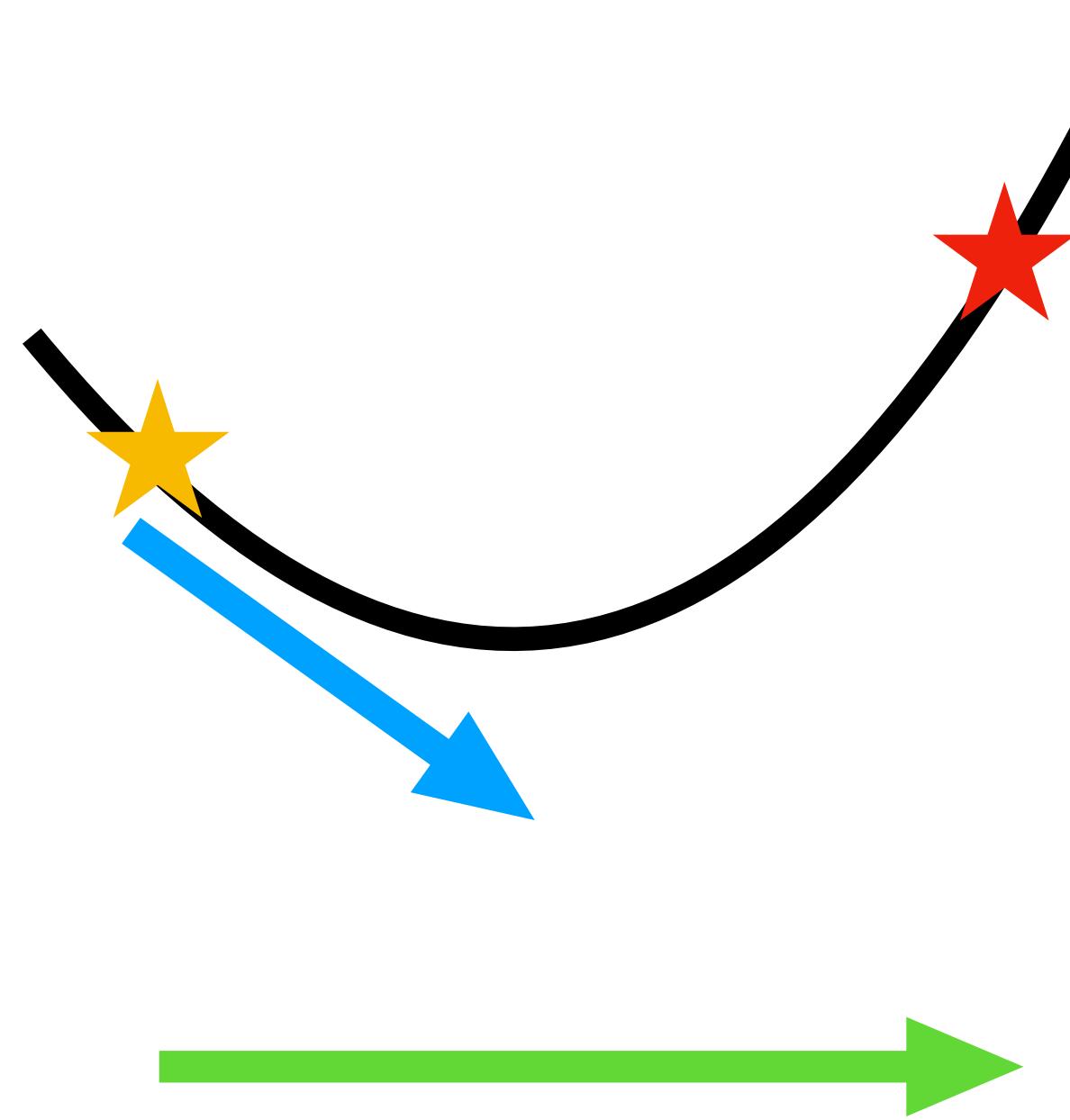




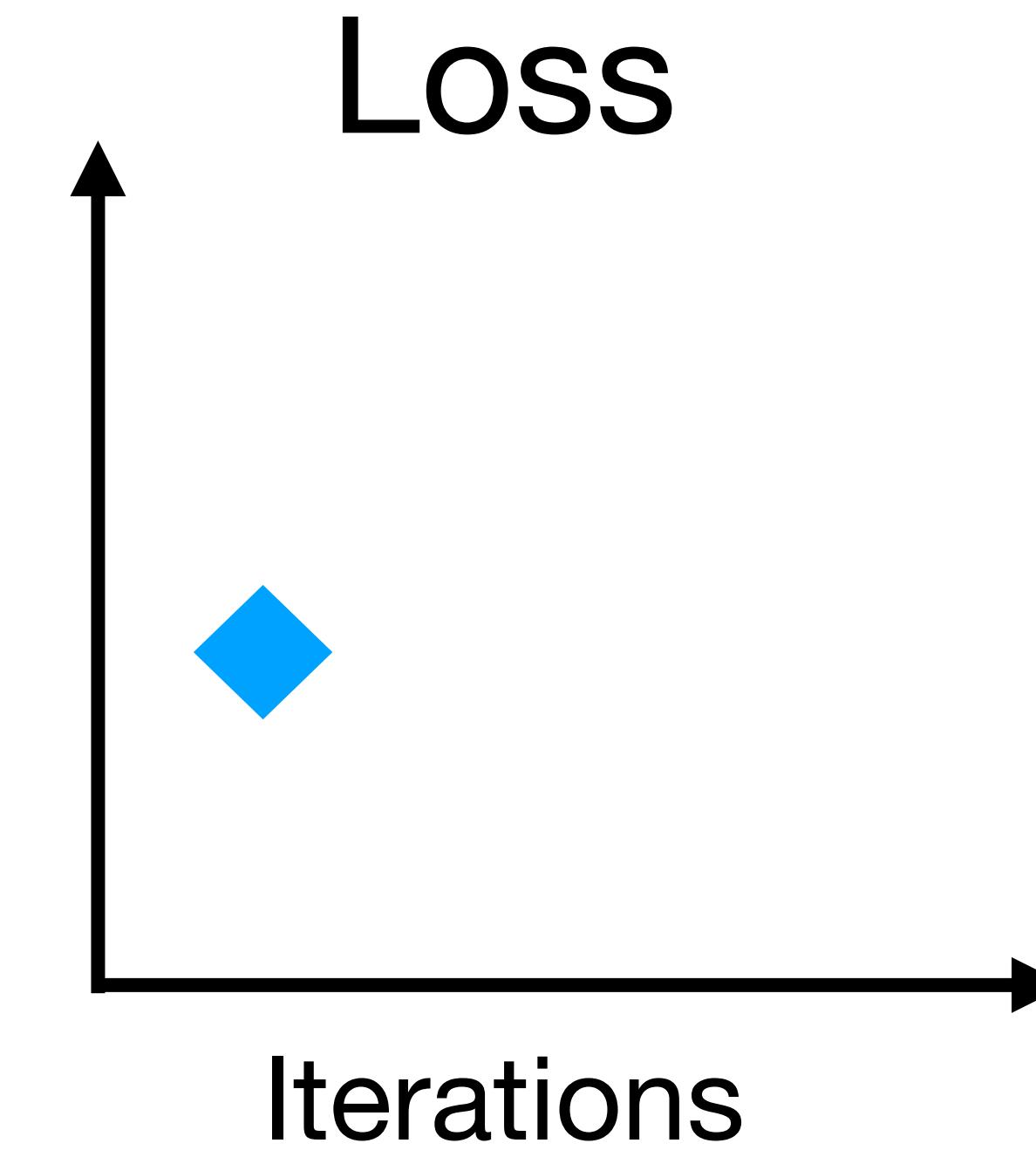


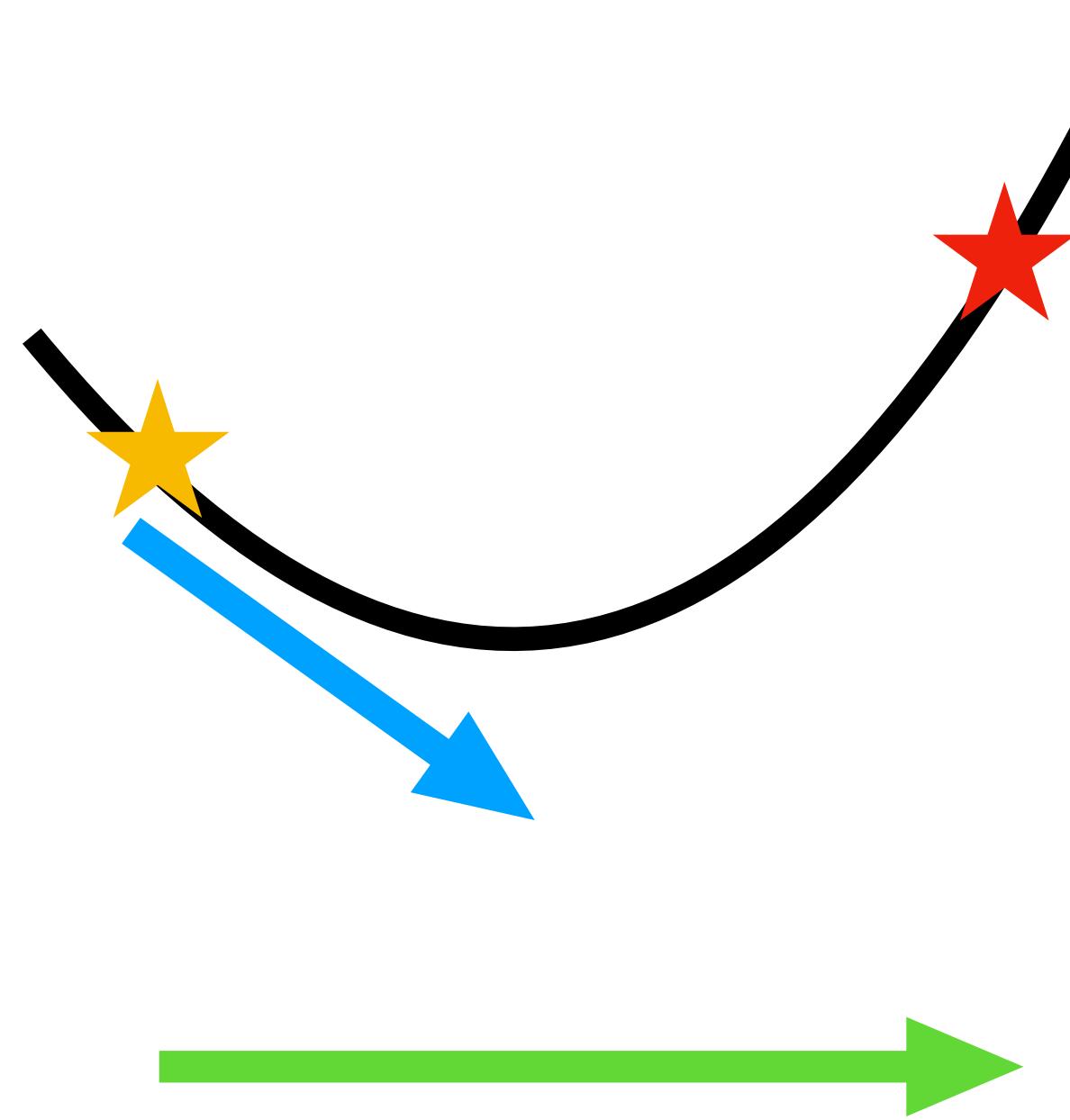
gradient x step size



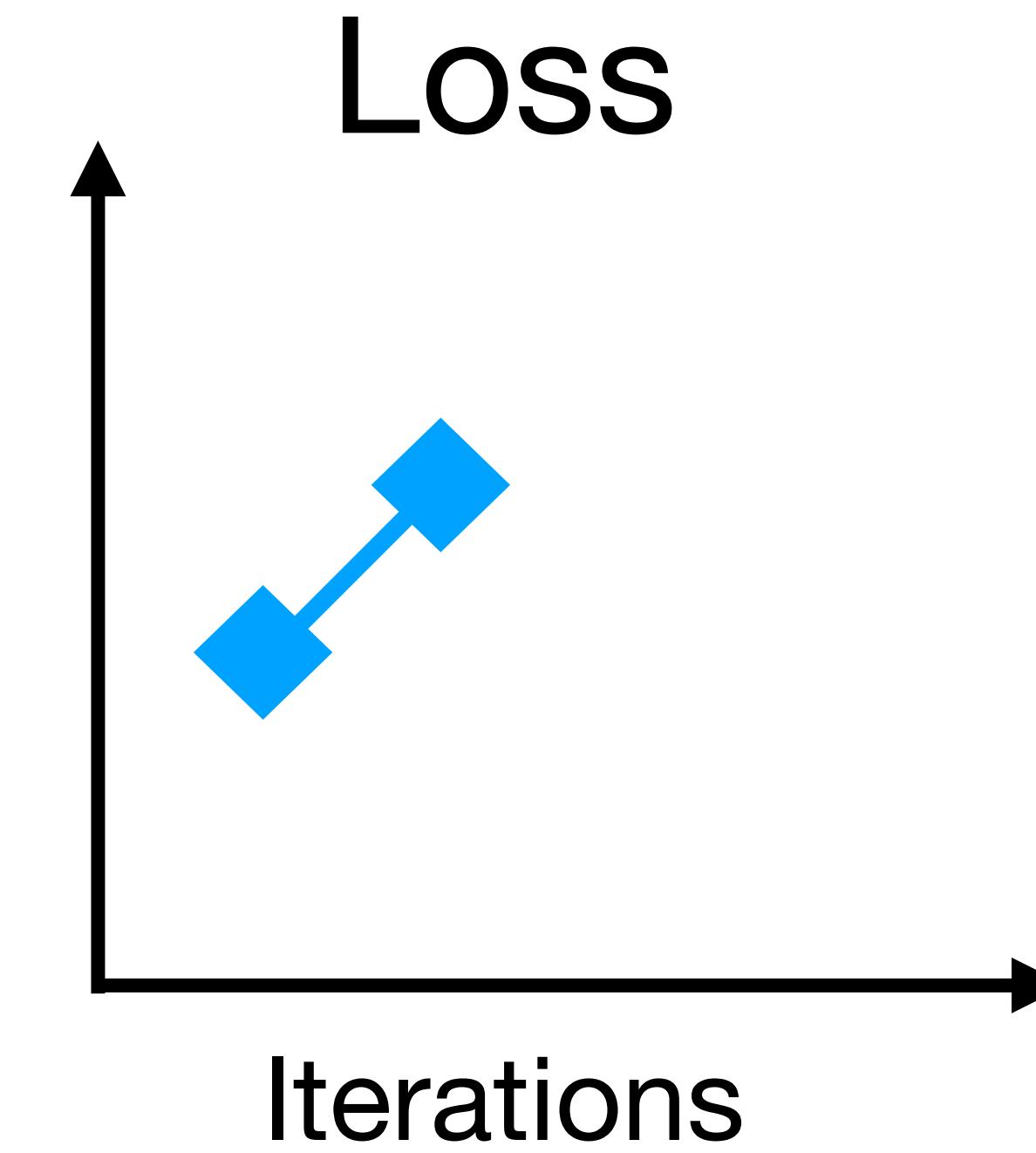


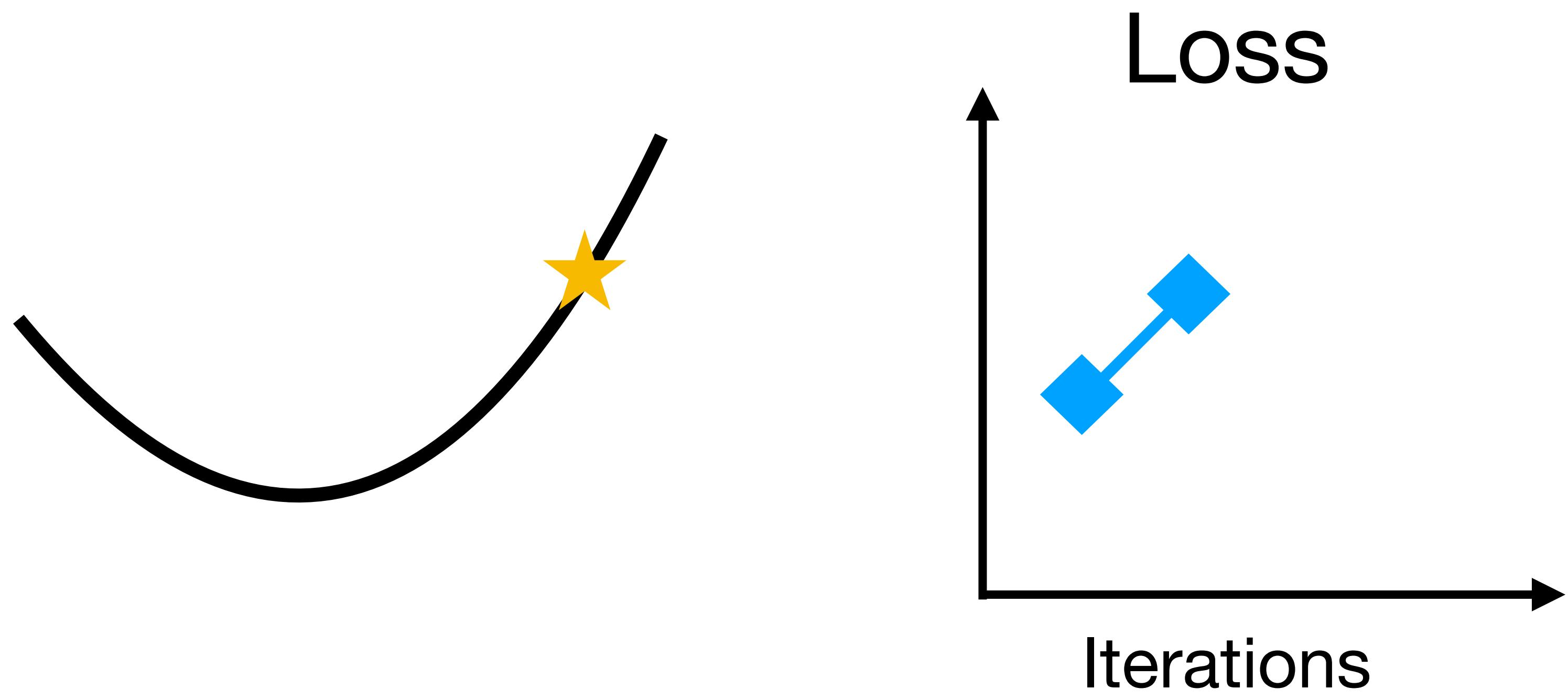
gradient x step size

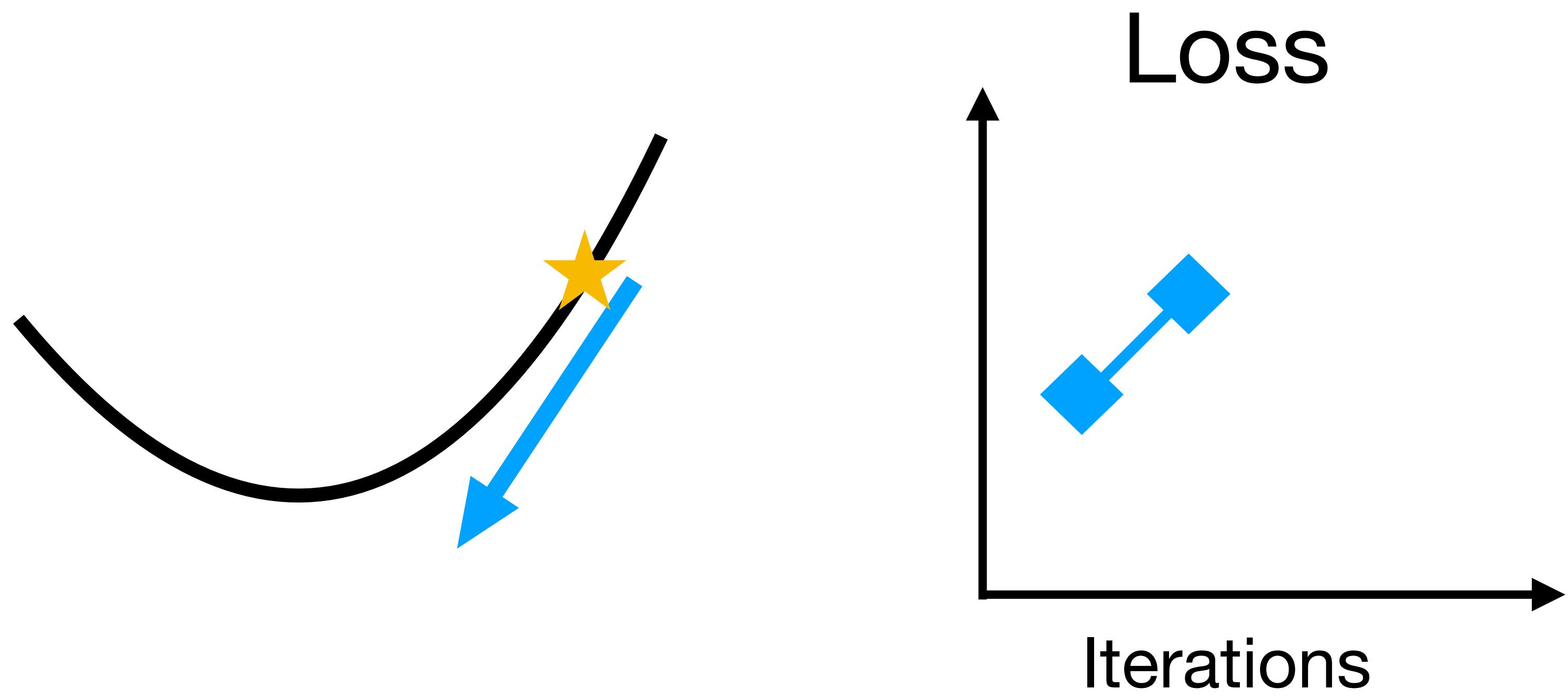


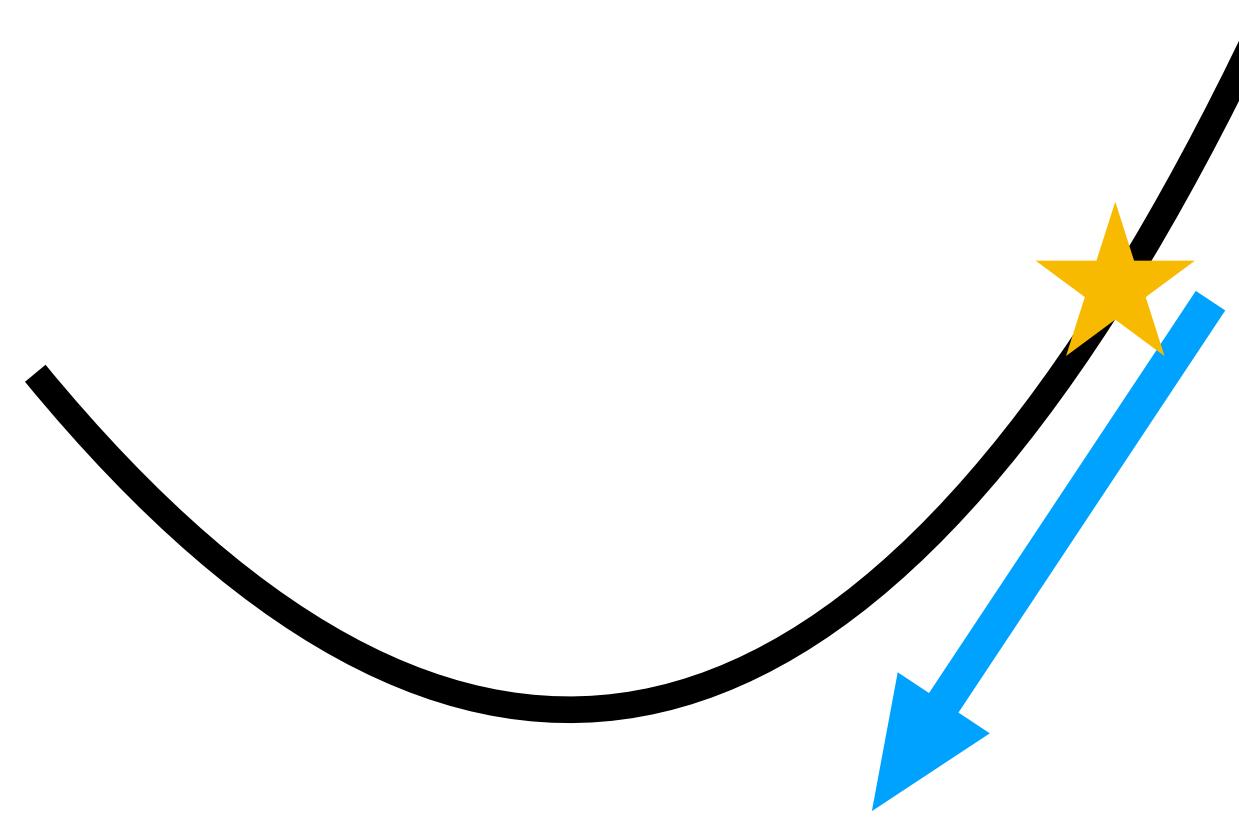


gradient x step size

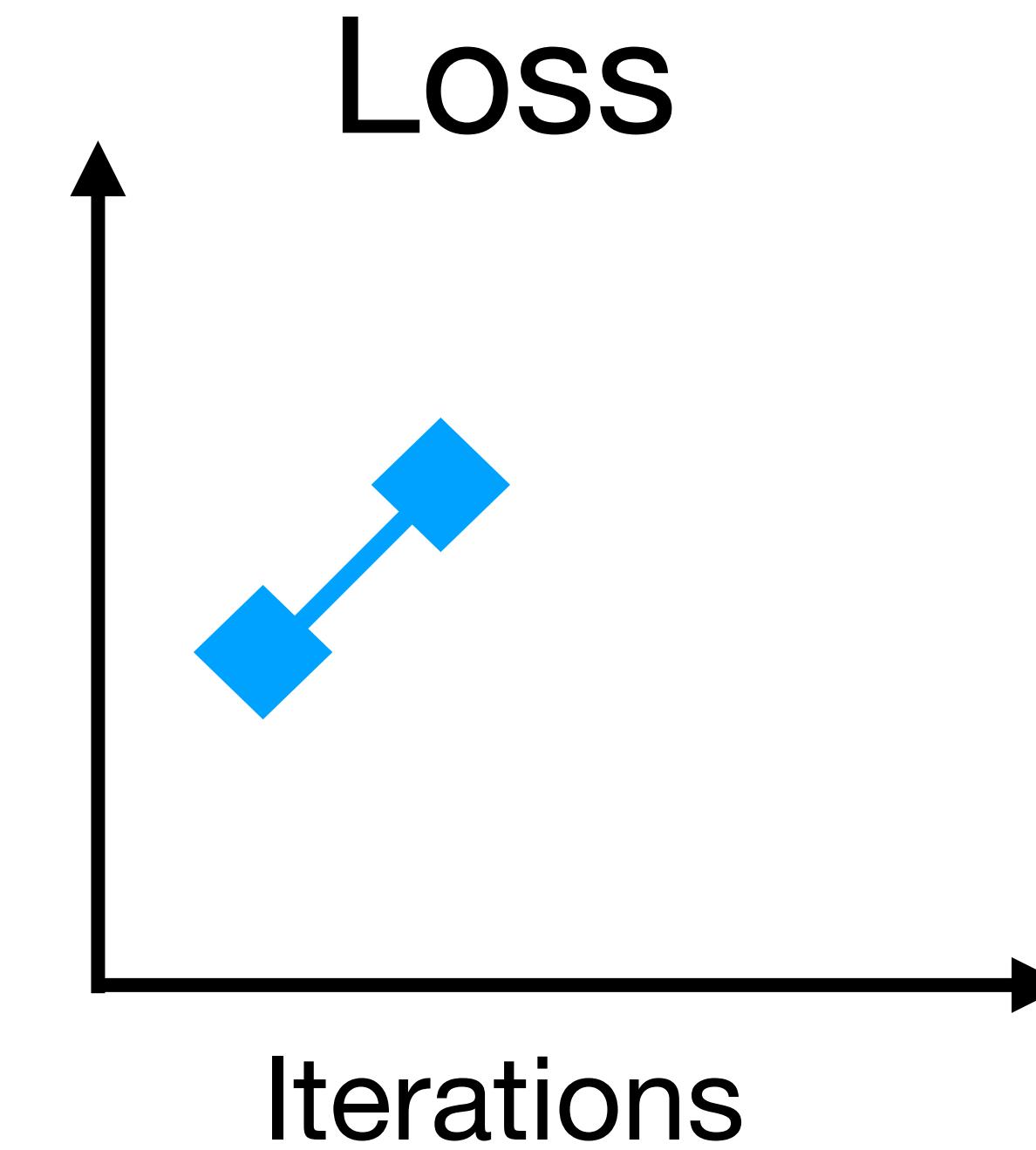


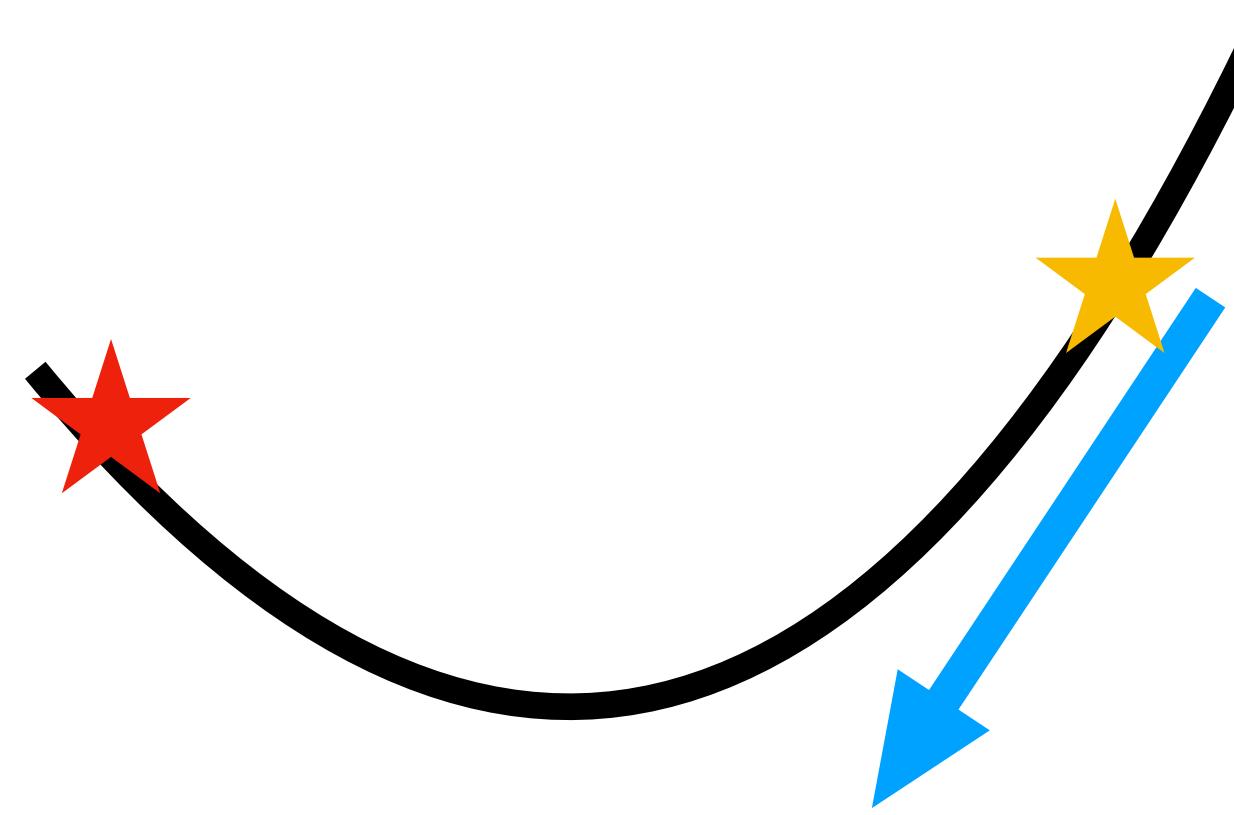




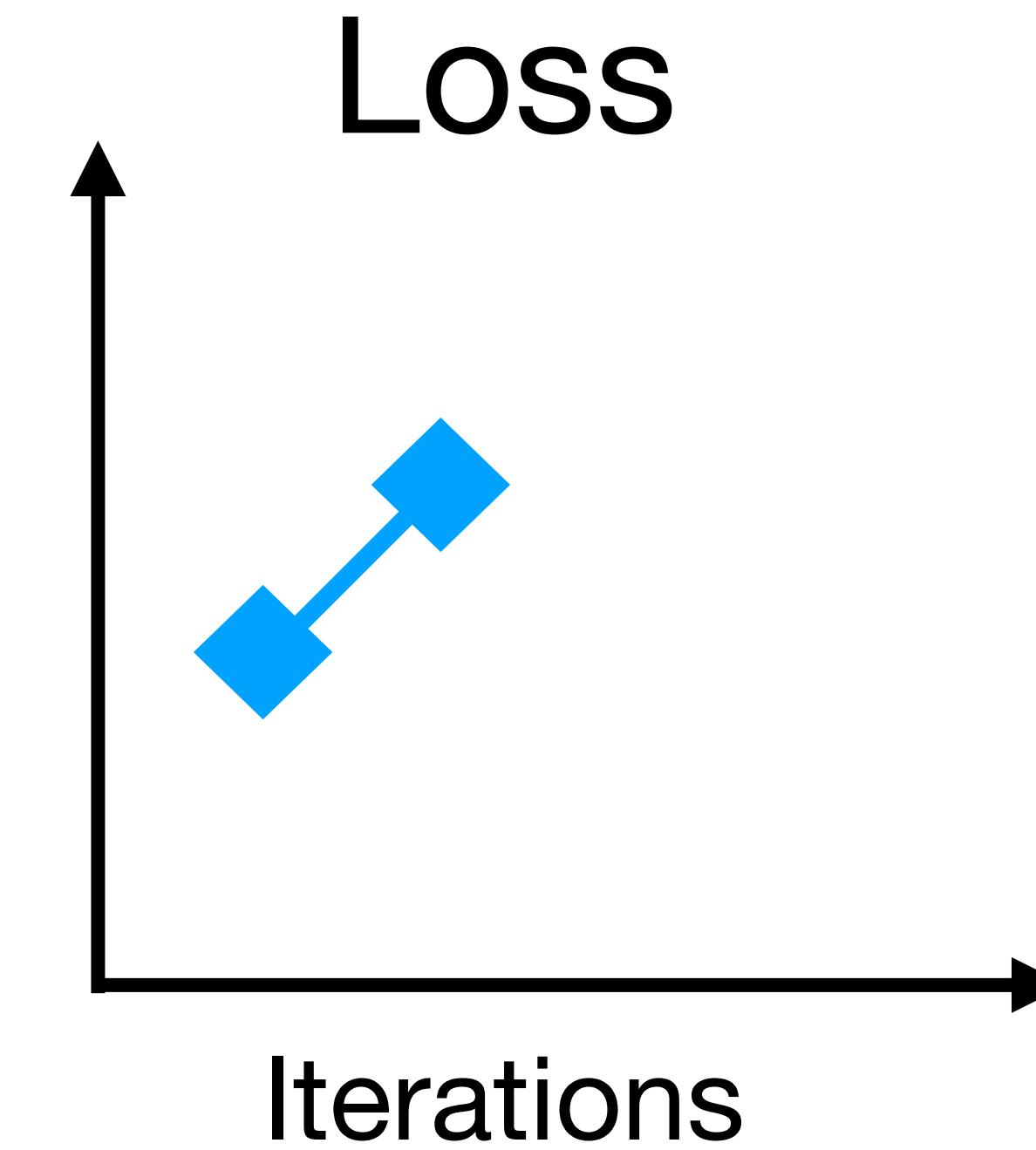


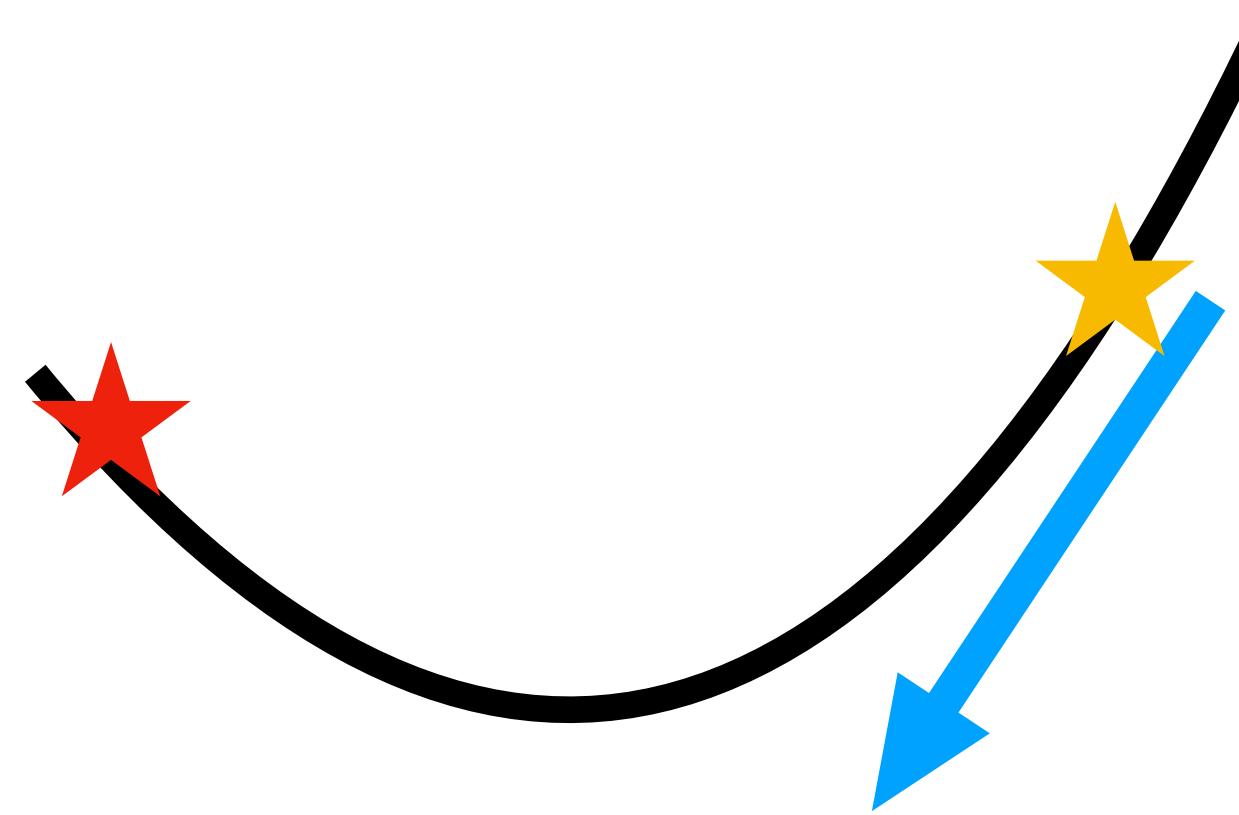
gradient x step size





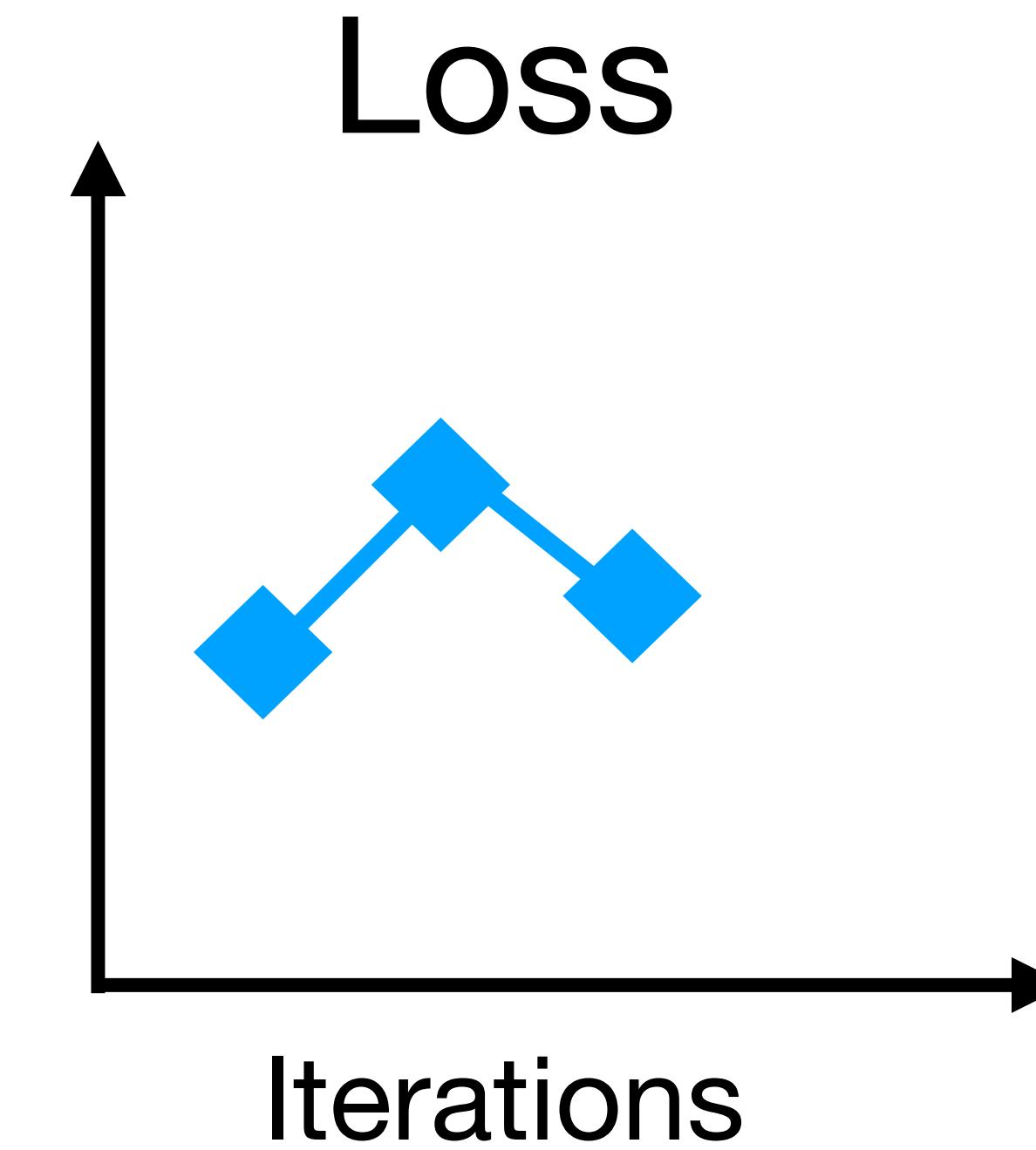
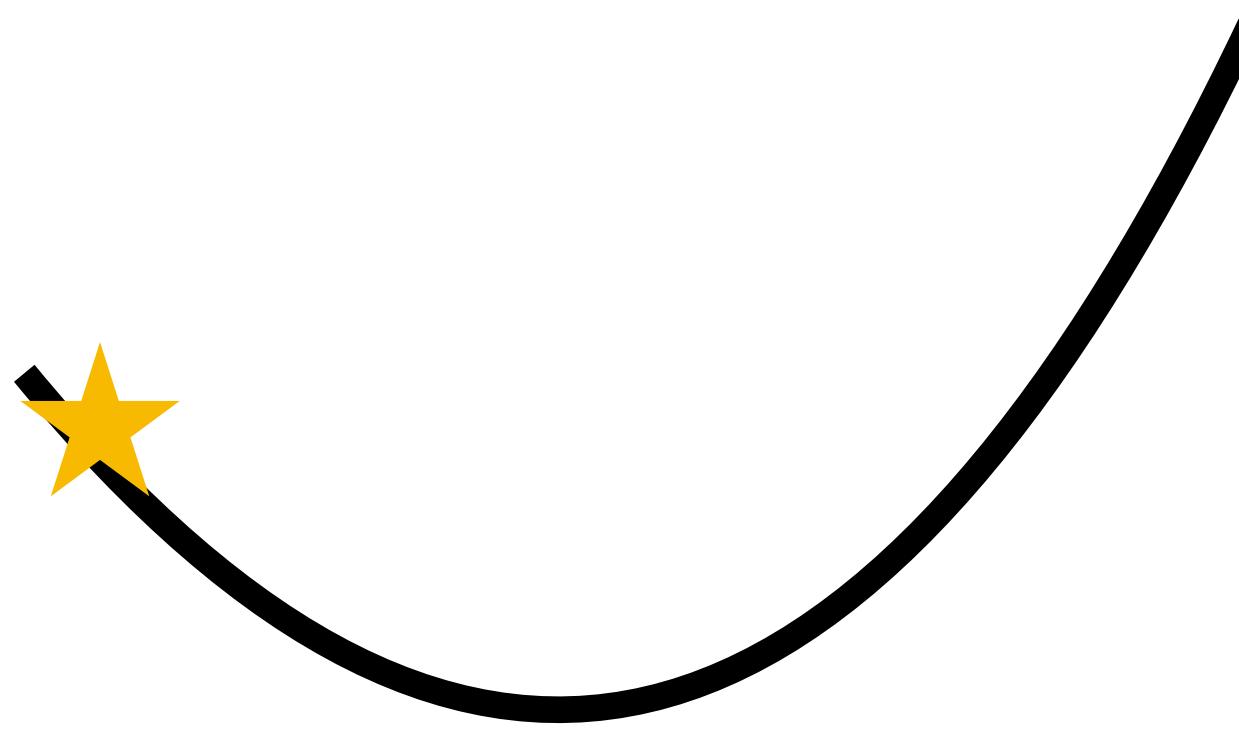
gradient x step size

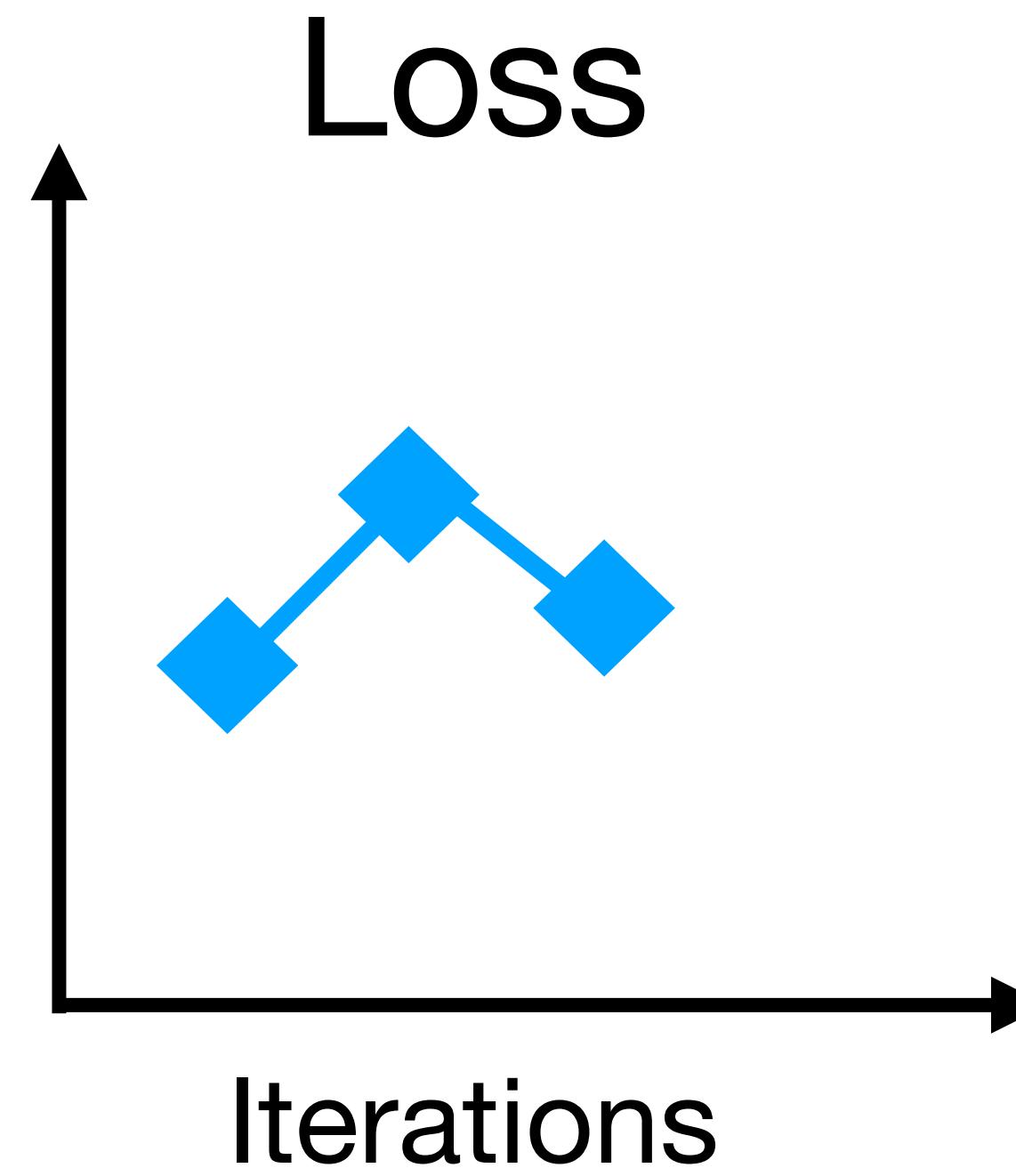
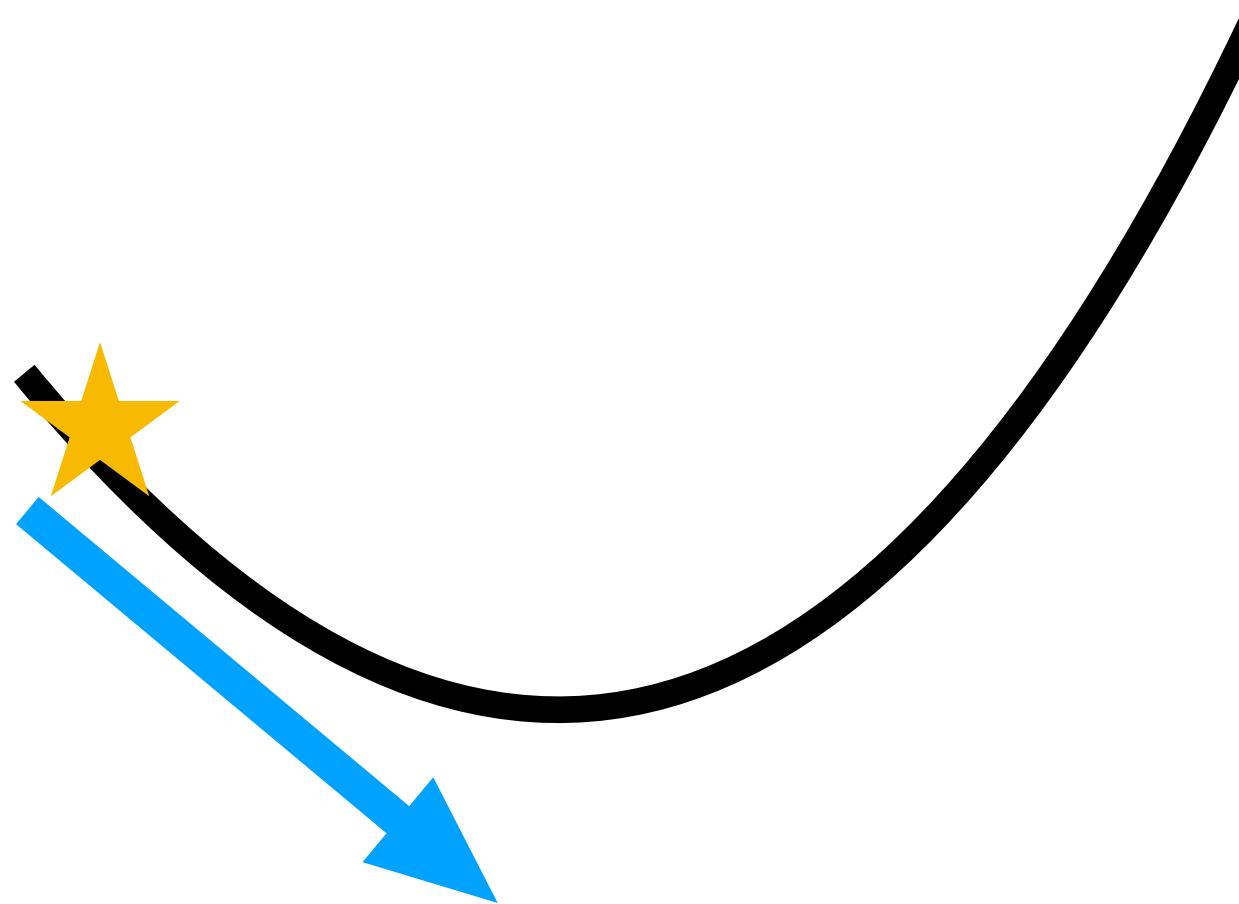


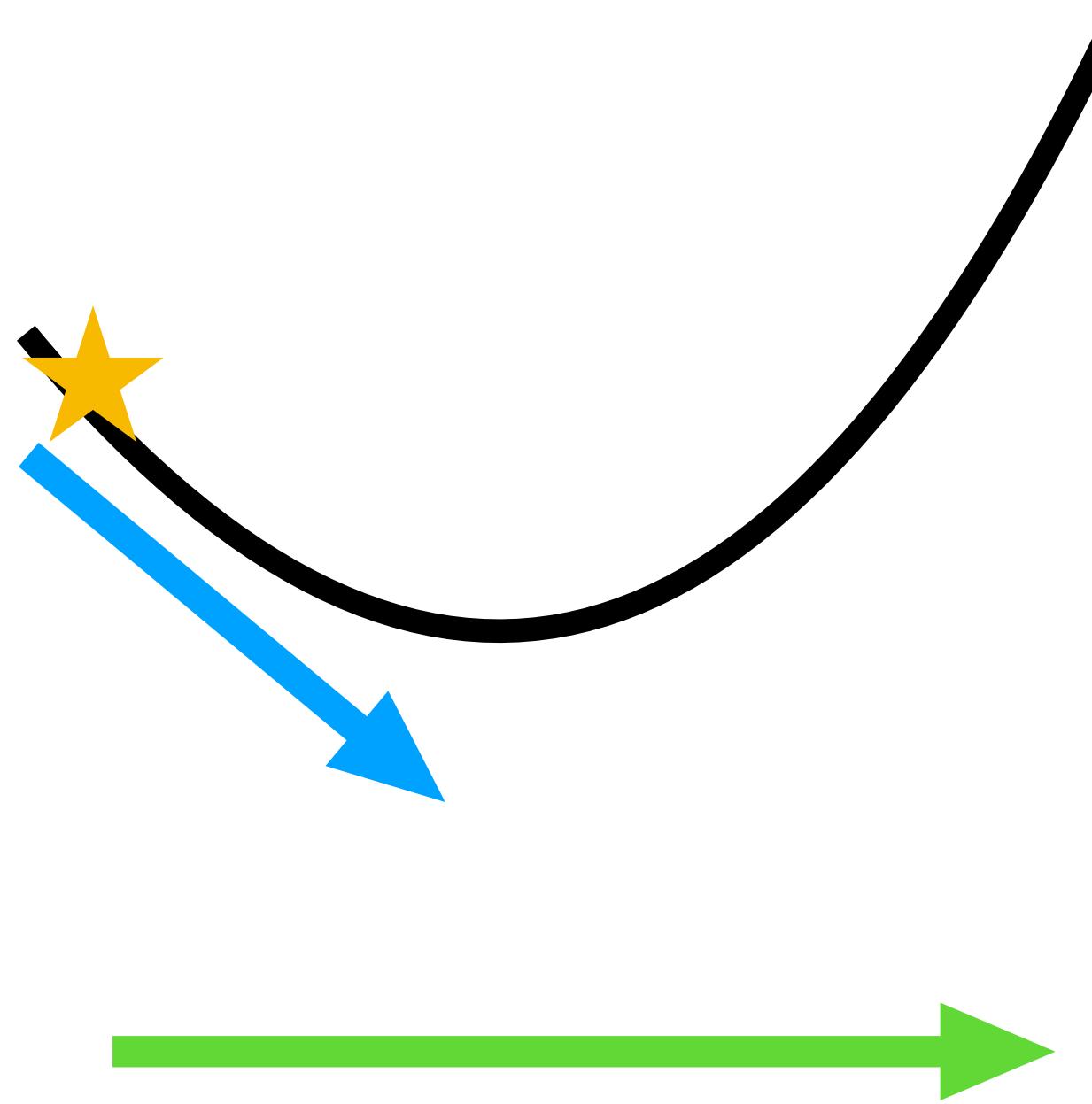


gradient x step size

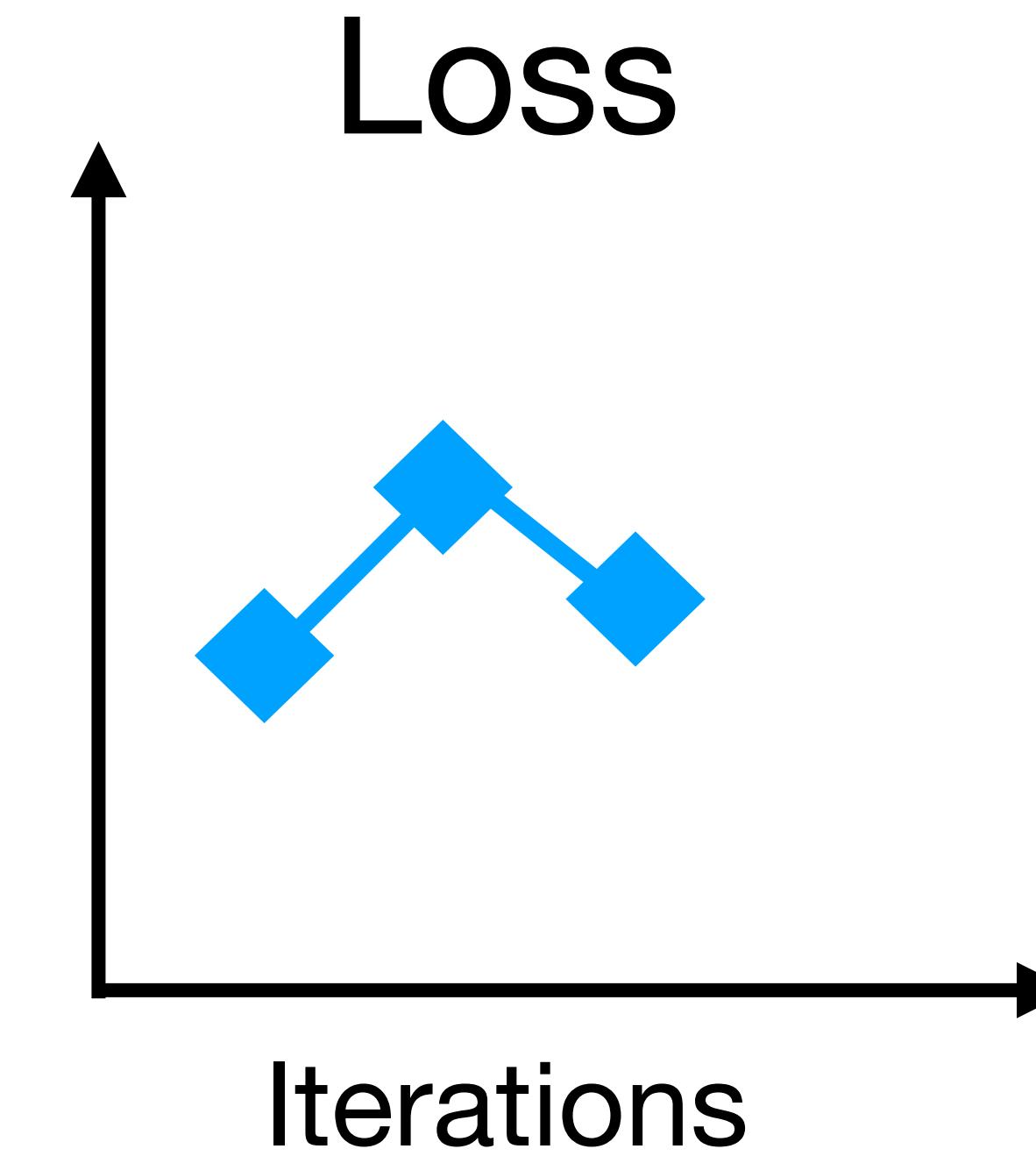


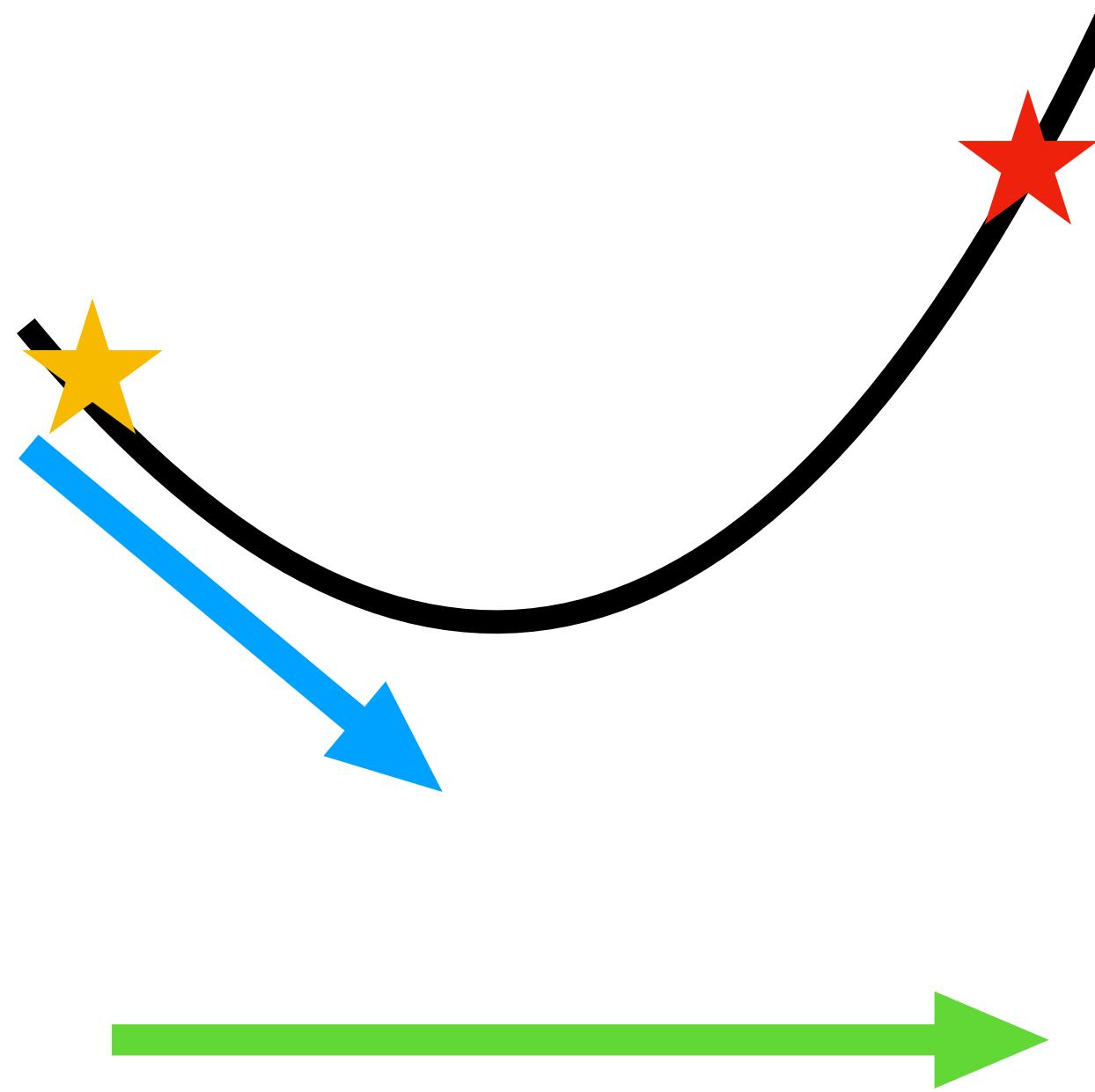




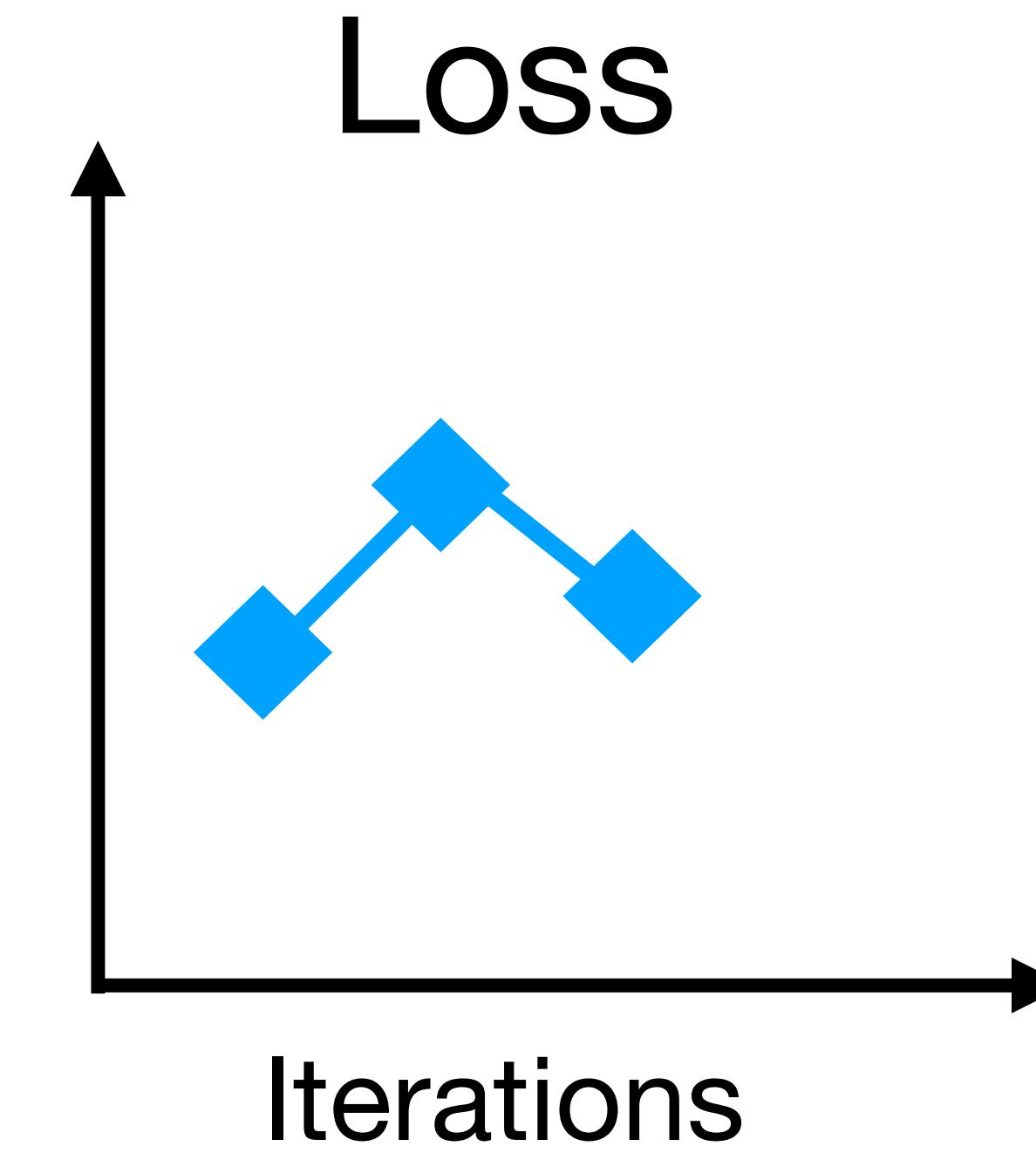


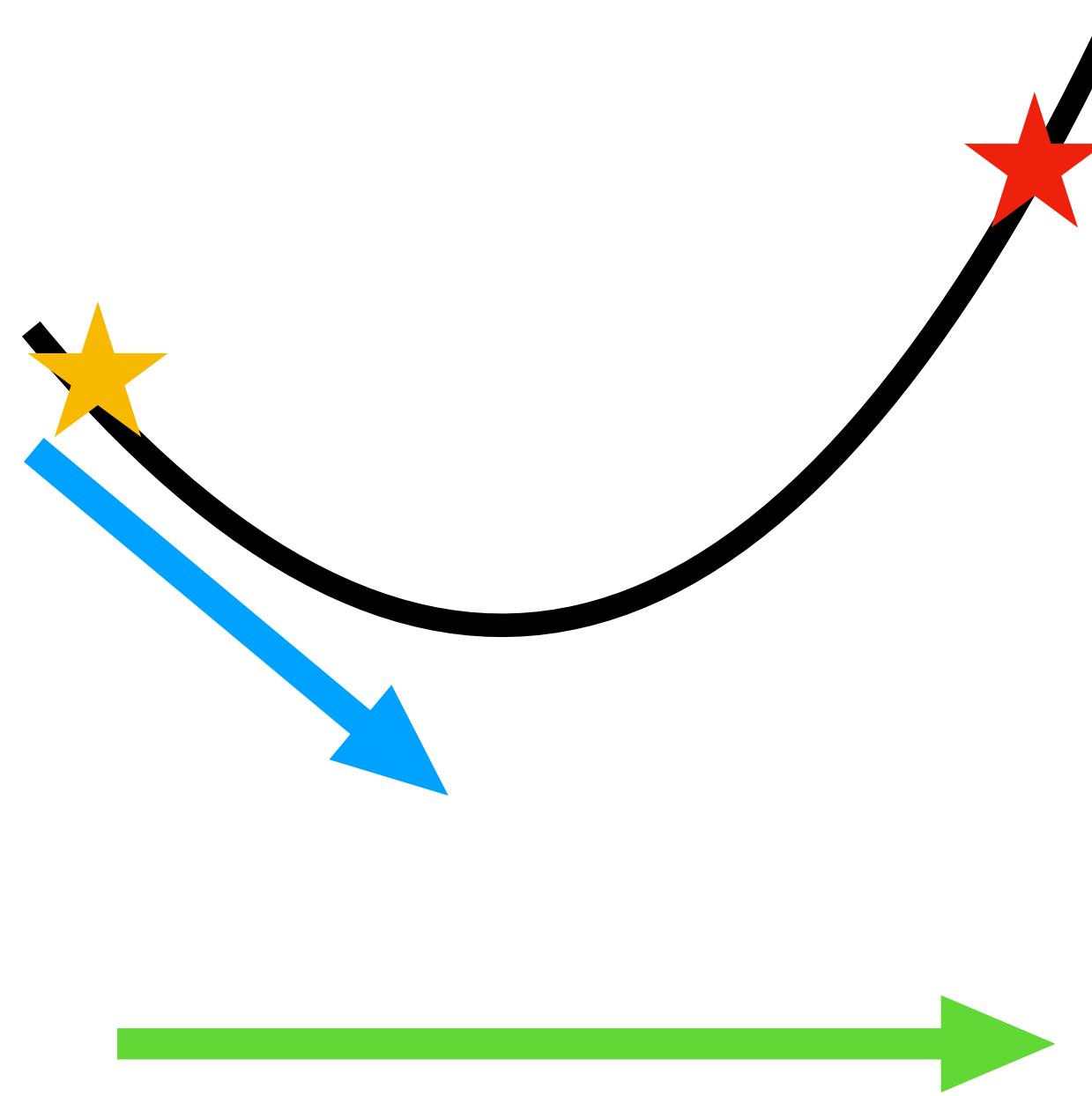
gradient x step size



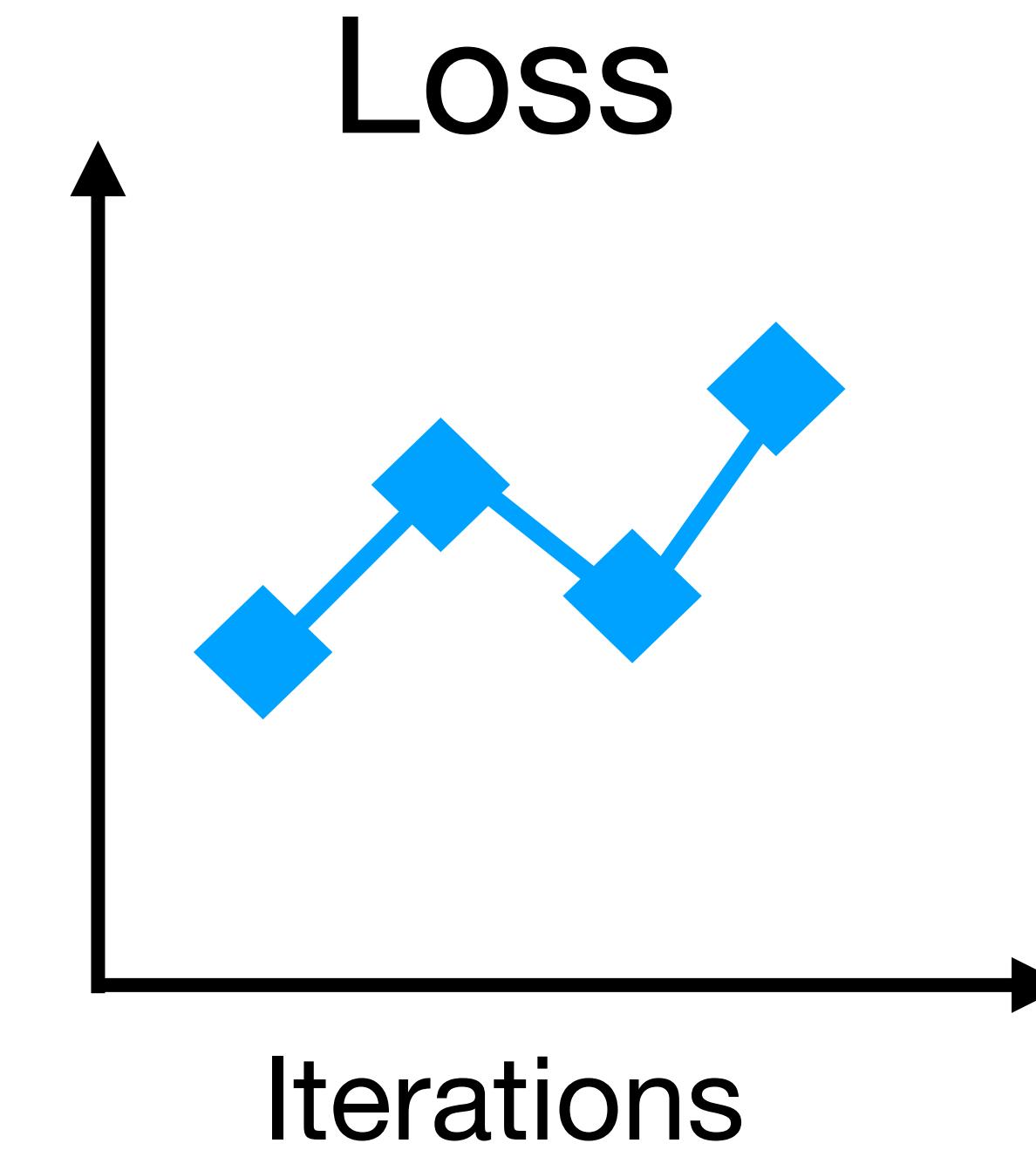


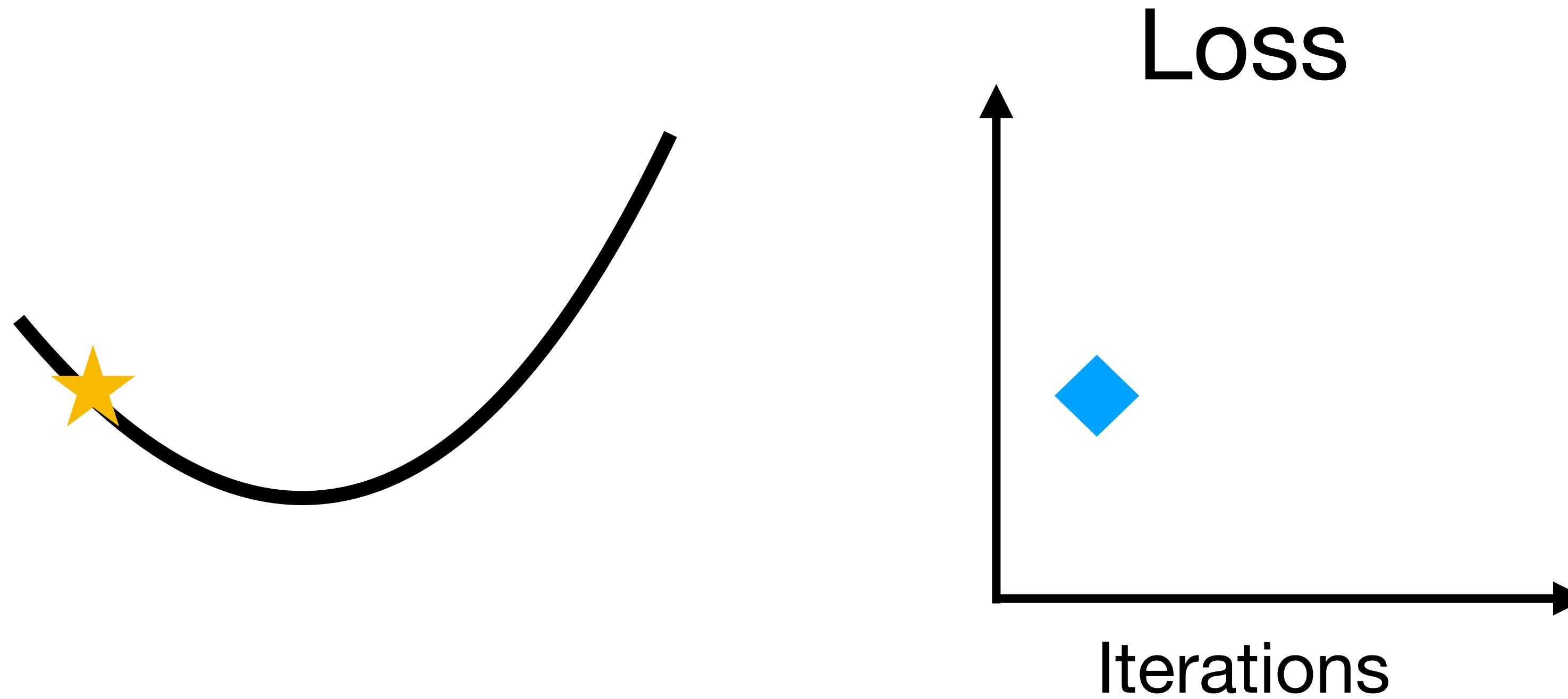
gradient x step size

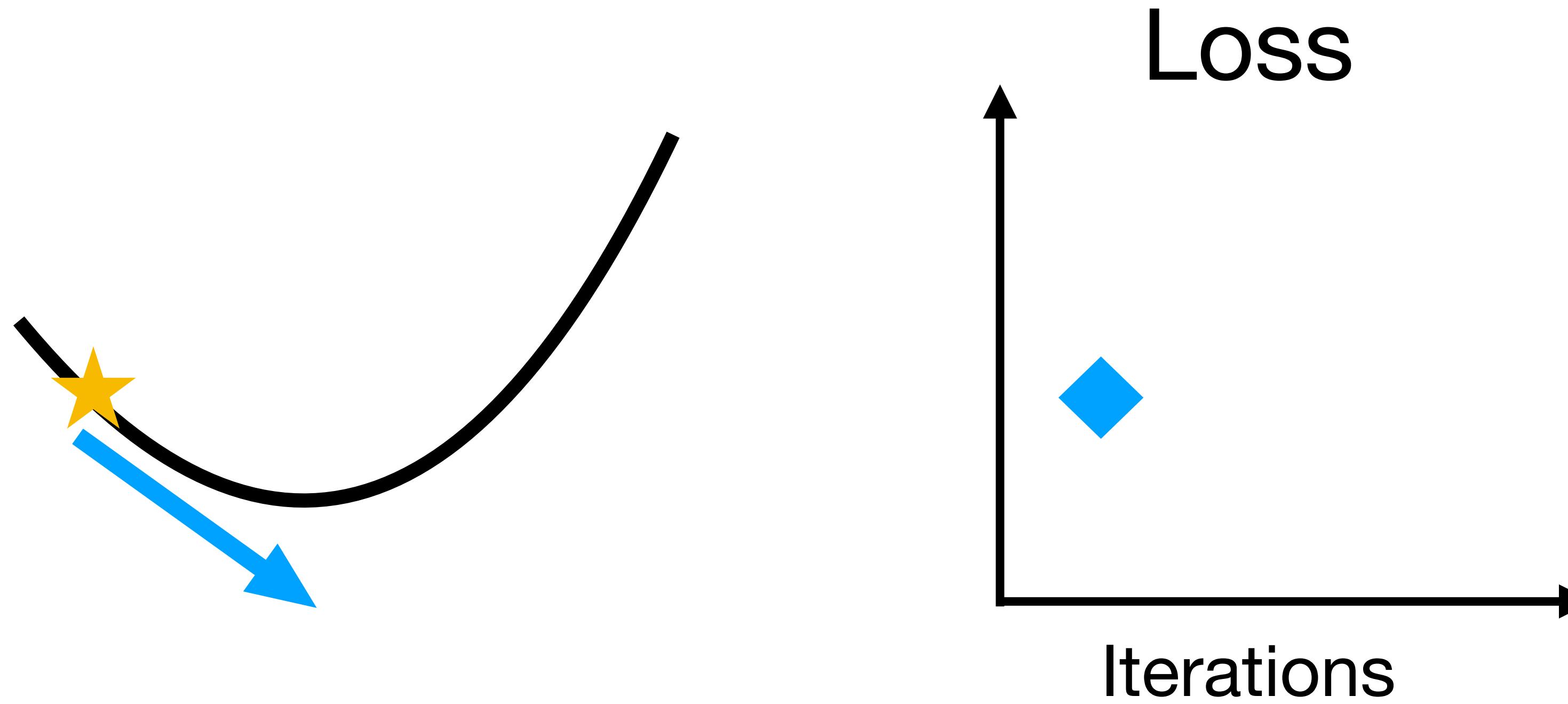


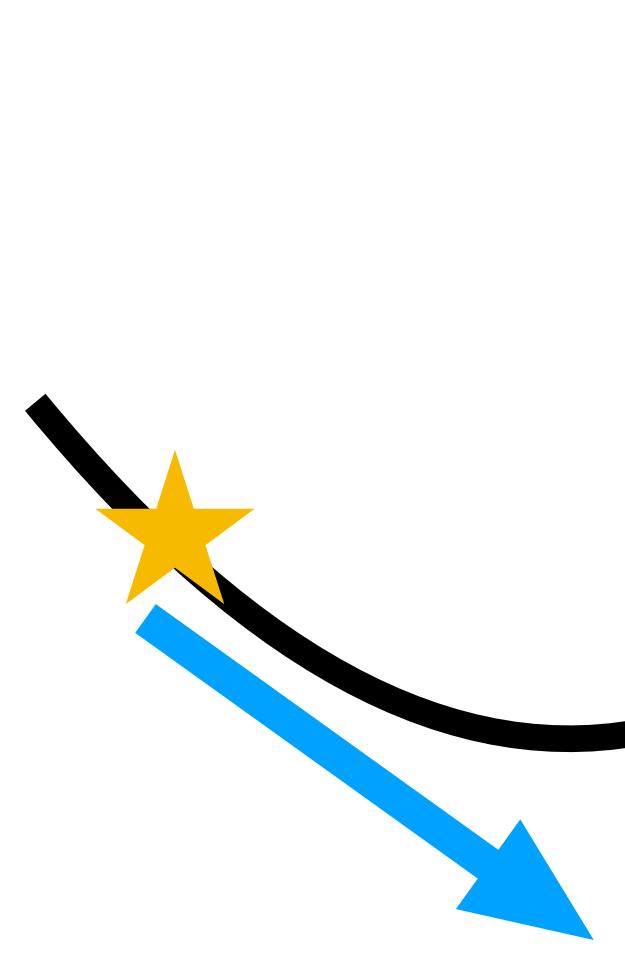


gradient x step size

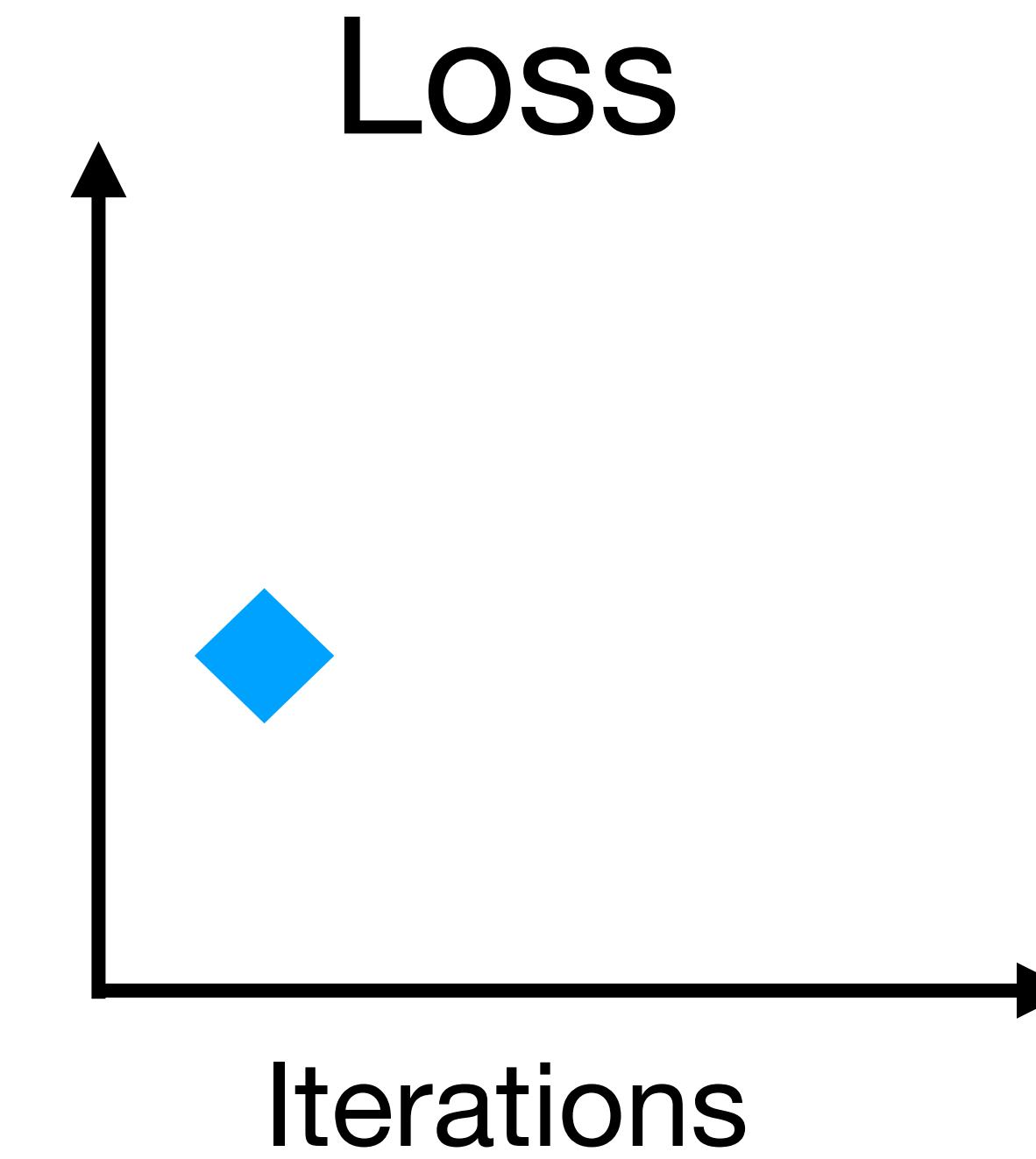


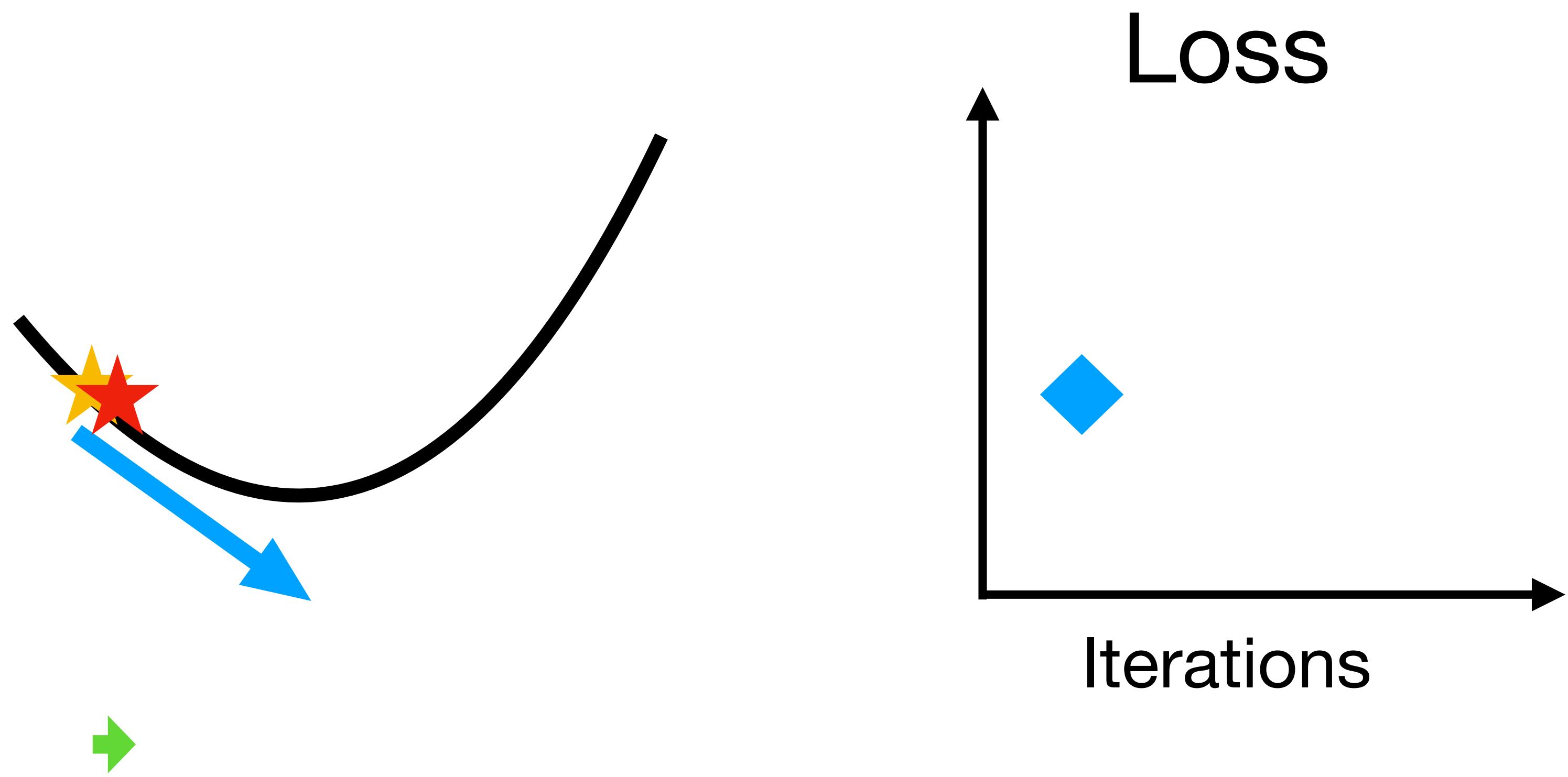




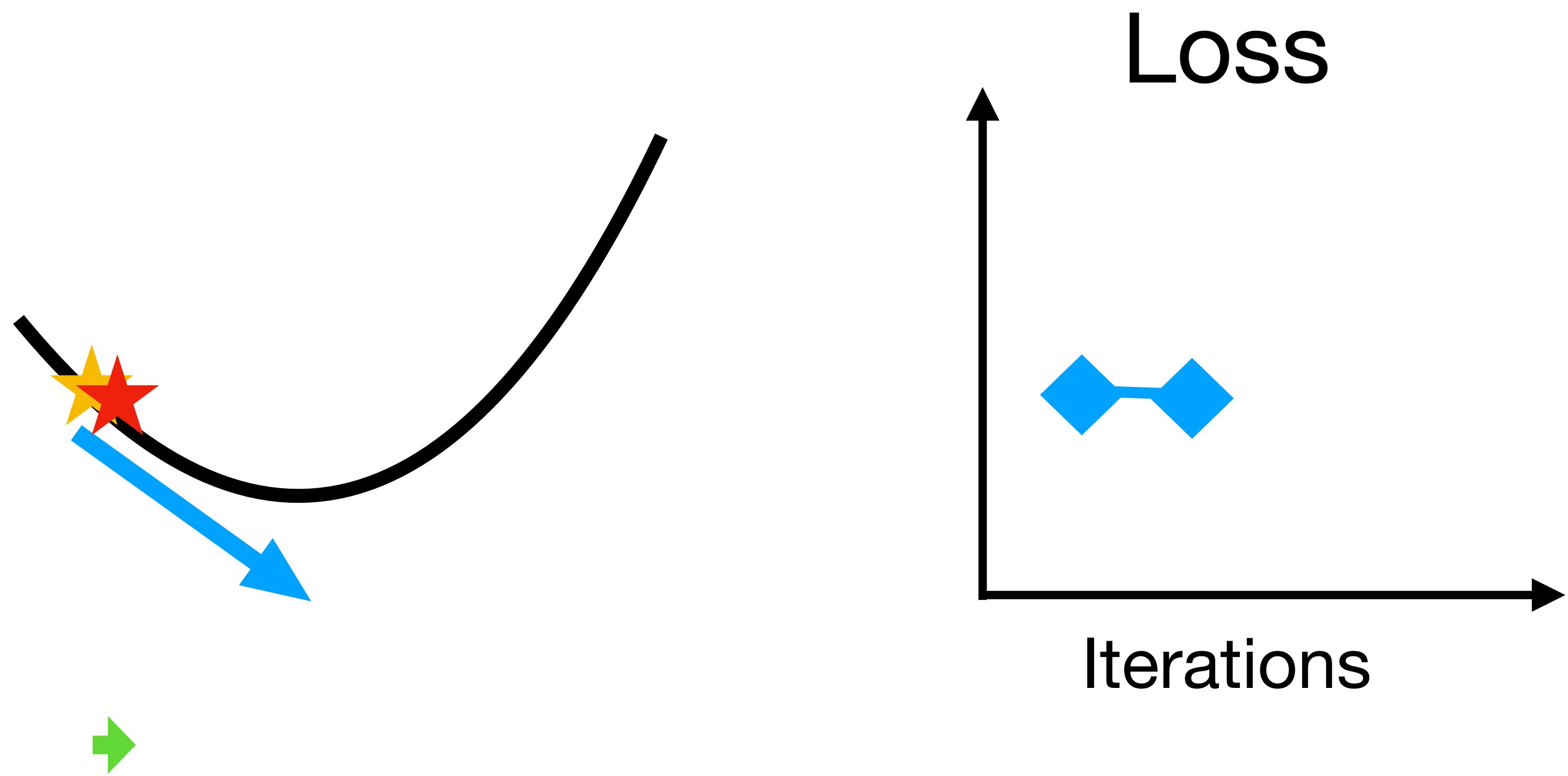


gradient x step size



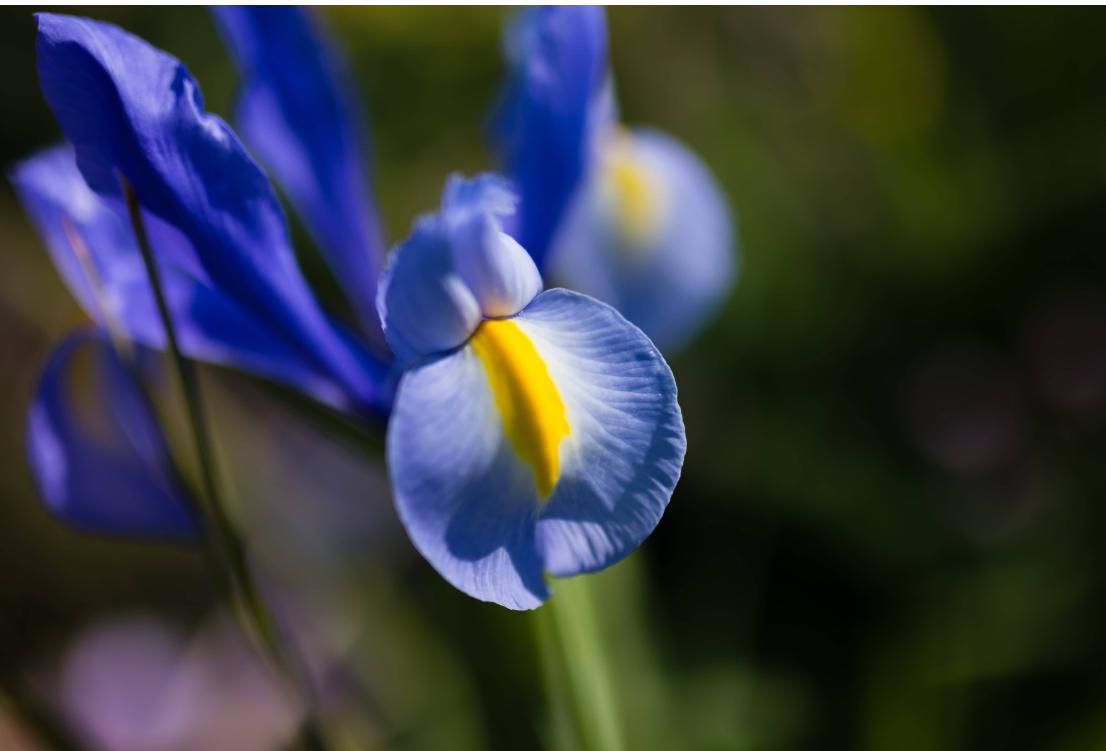


gradient x step size



gradient x step size

Iris dataset



Input

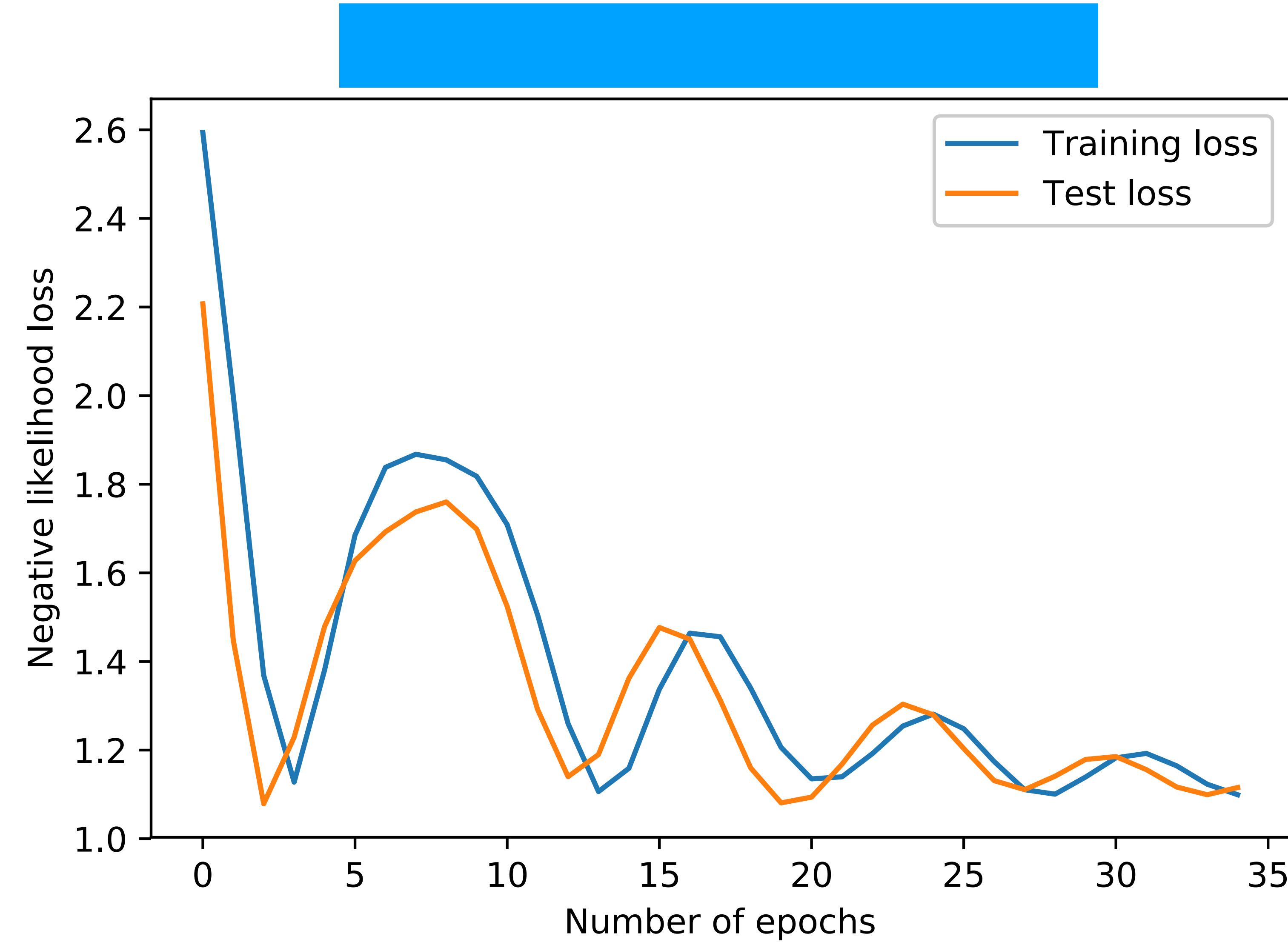
Dense 1

Dense 2

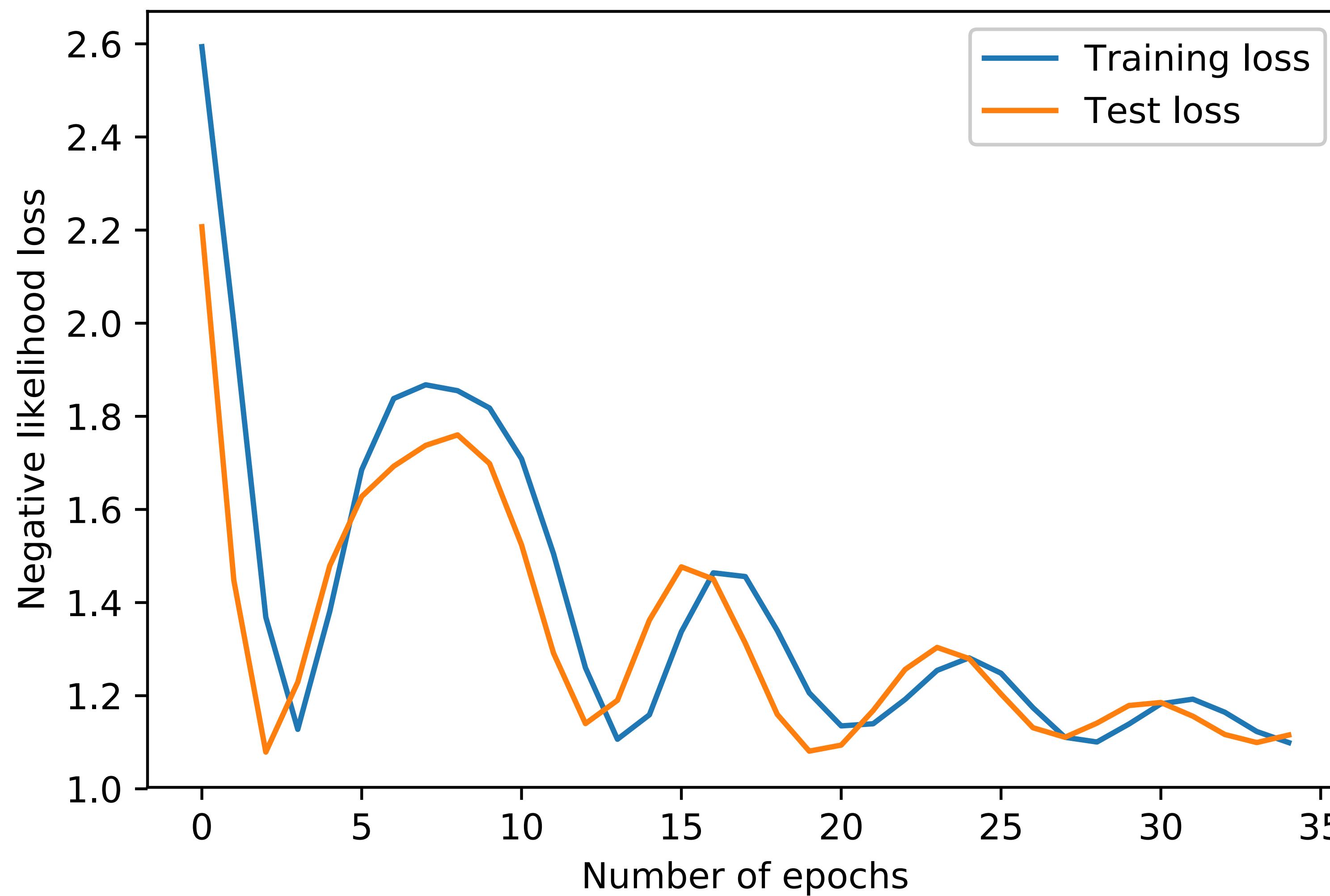
Linear

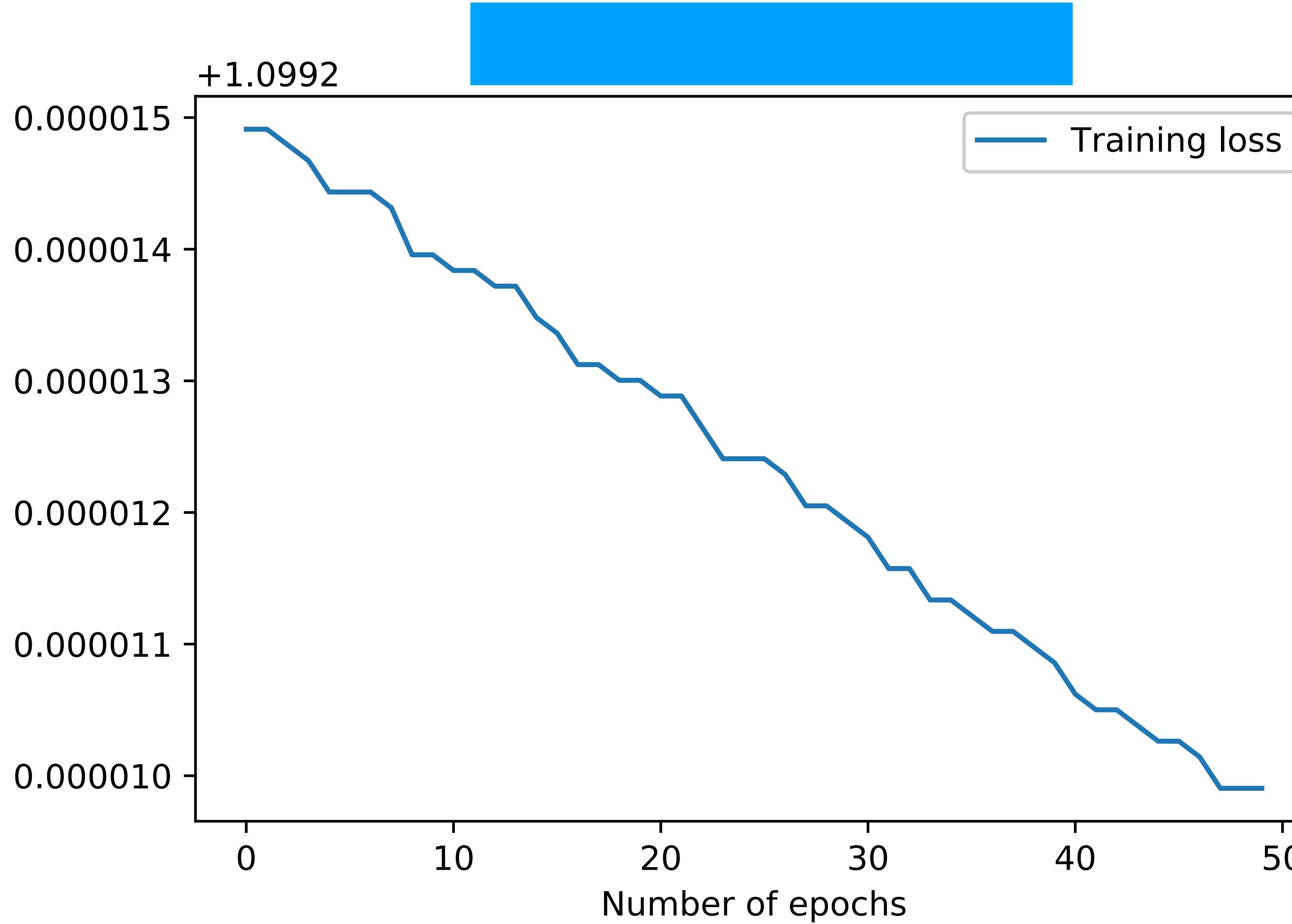
Softmax(3)

Output

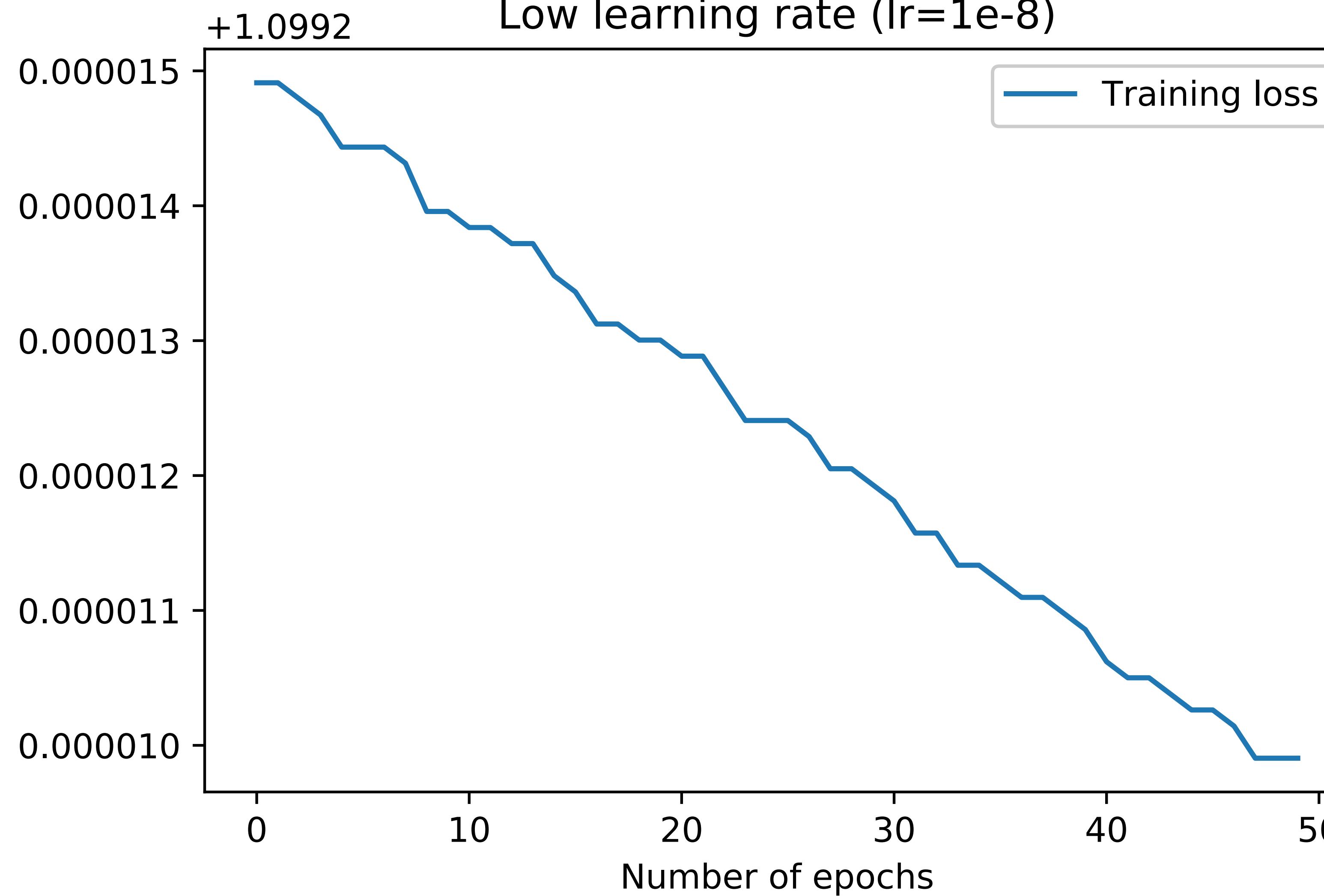


High learning rate loss curves(lr=2)

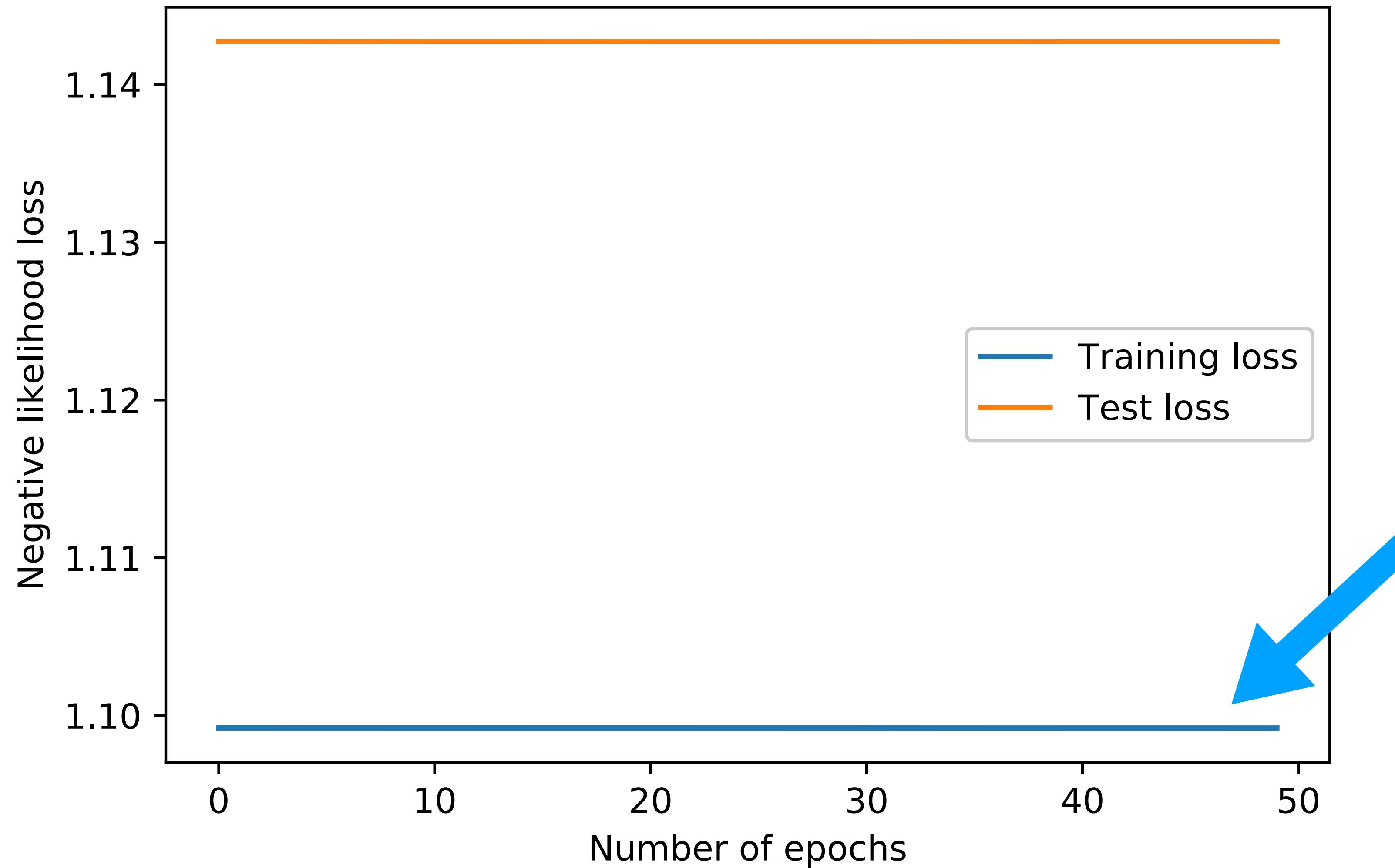


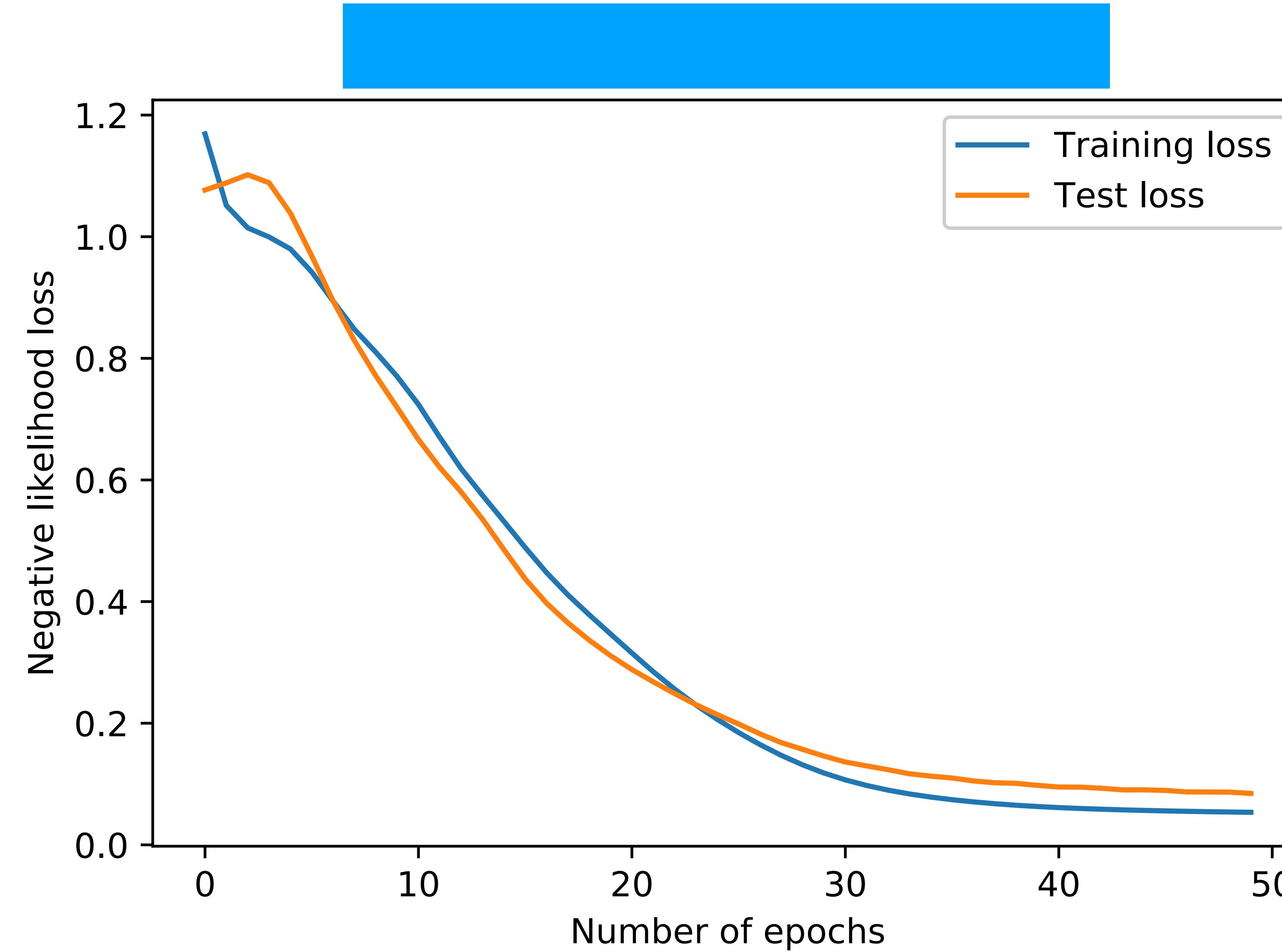


Low learning rate ($lr=1e-8$)

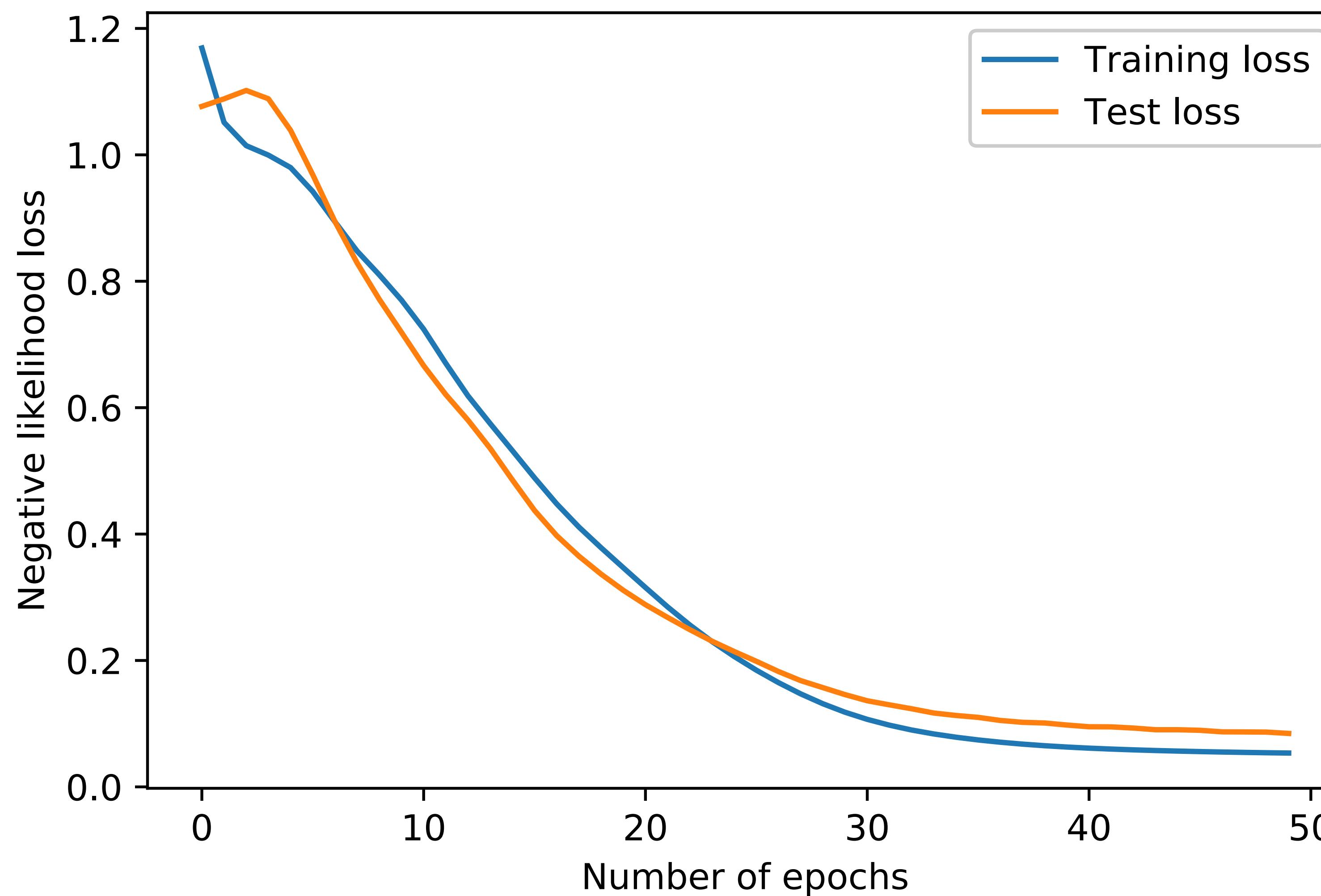


Low learning rate ($lr=1e-8$)

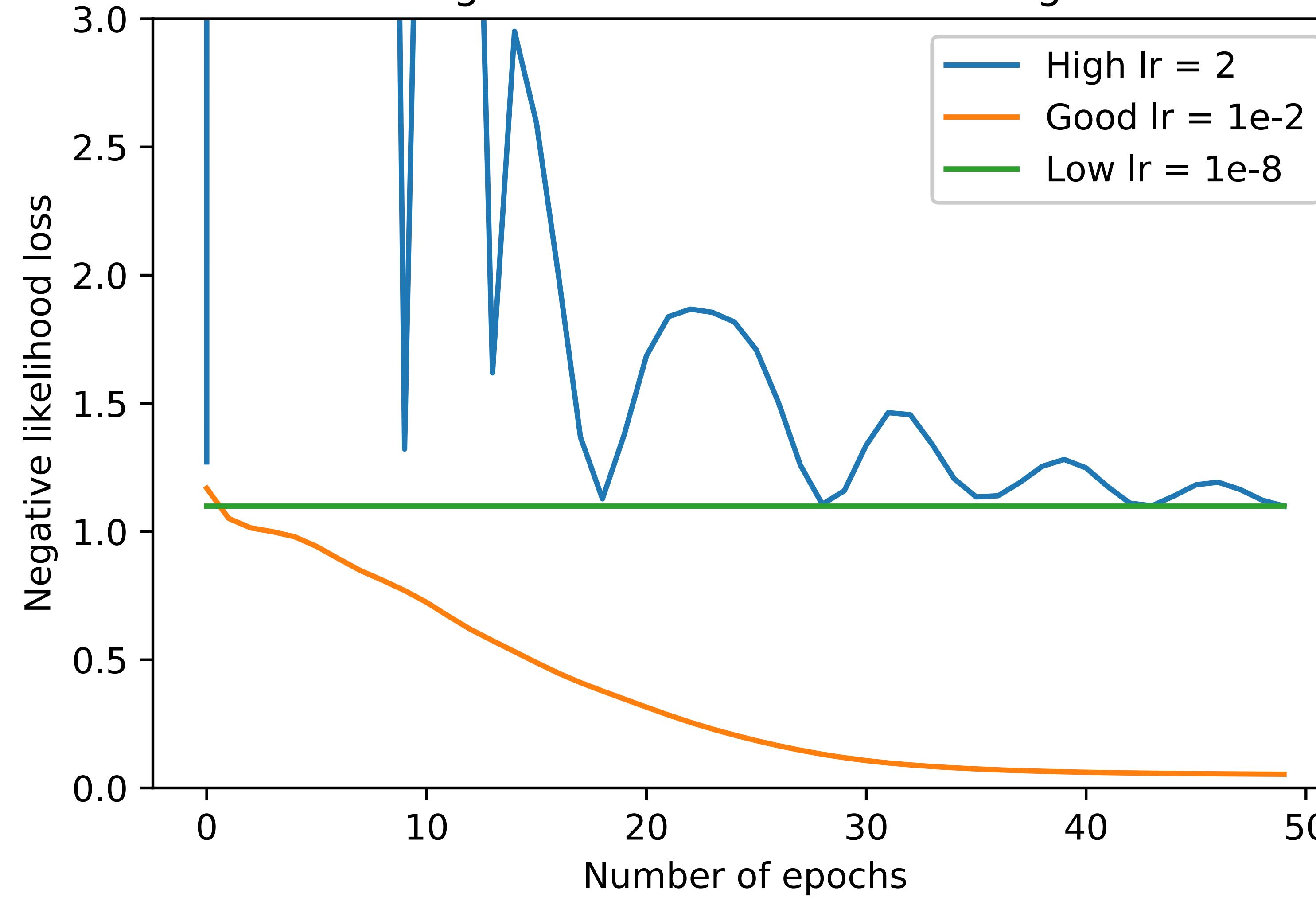




Good loss curves ($\text{lr}=1\text{e-}2$)



Training losses at different learning rates

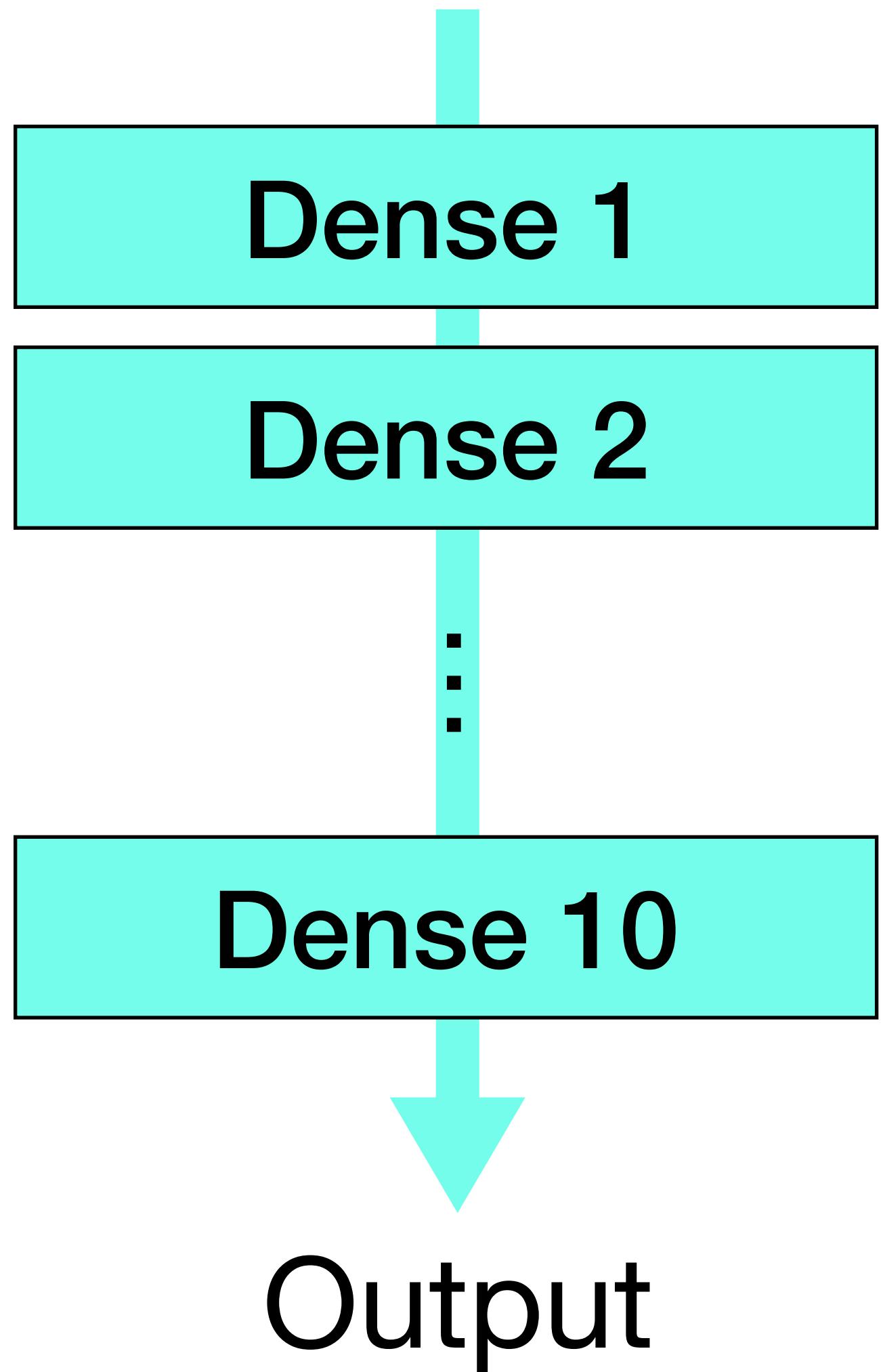


How calculate gradient?

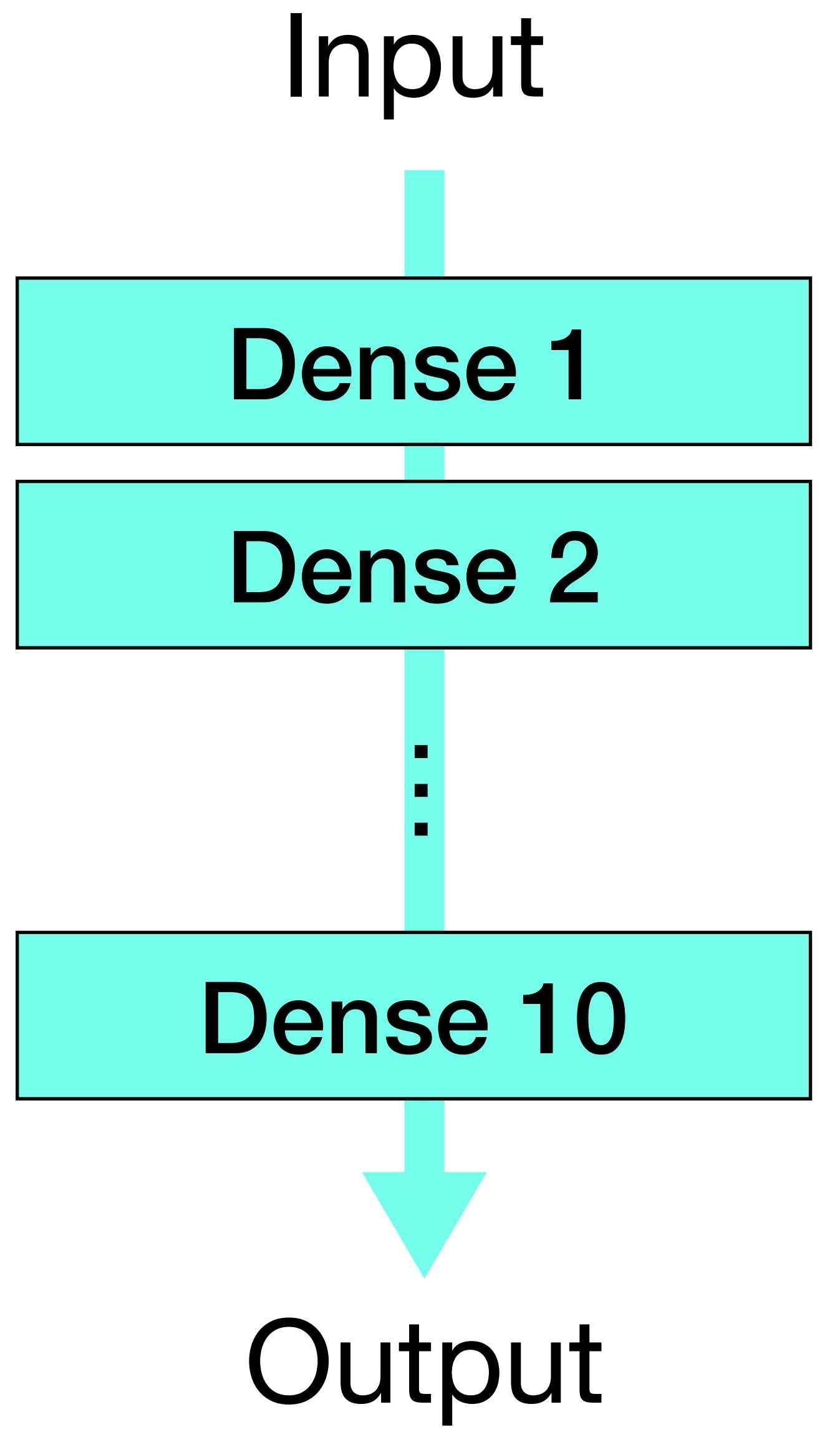


How calculate gradient?

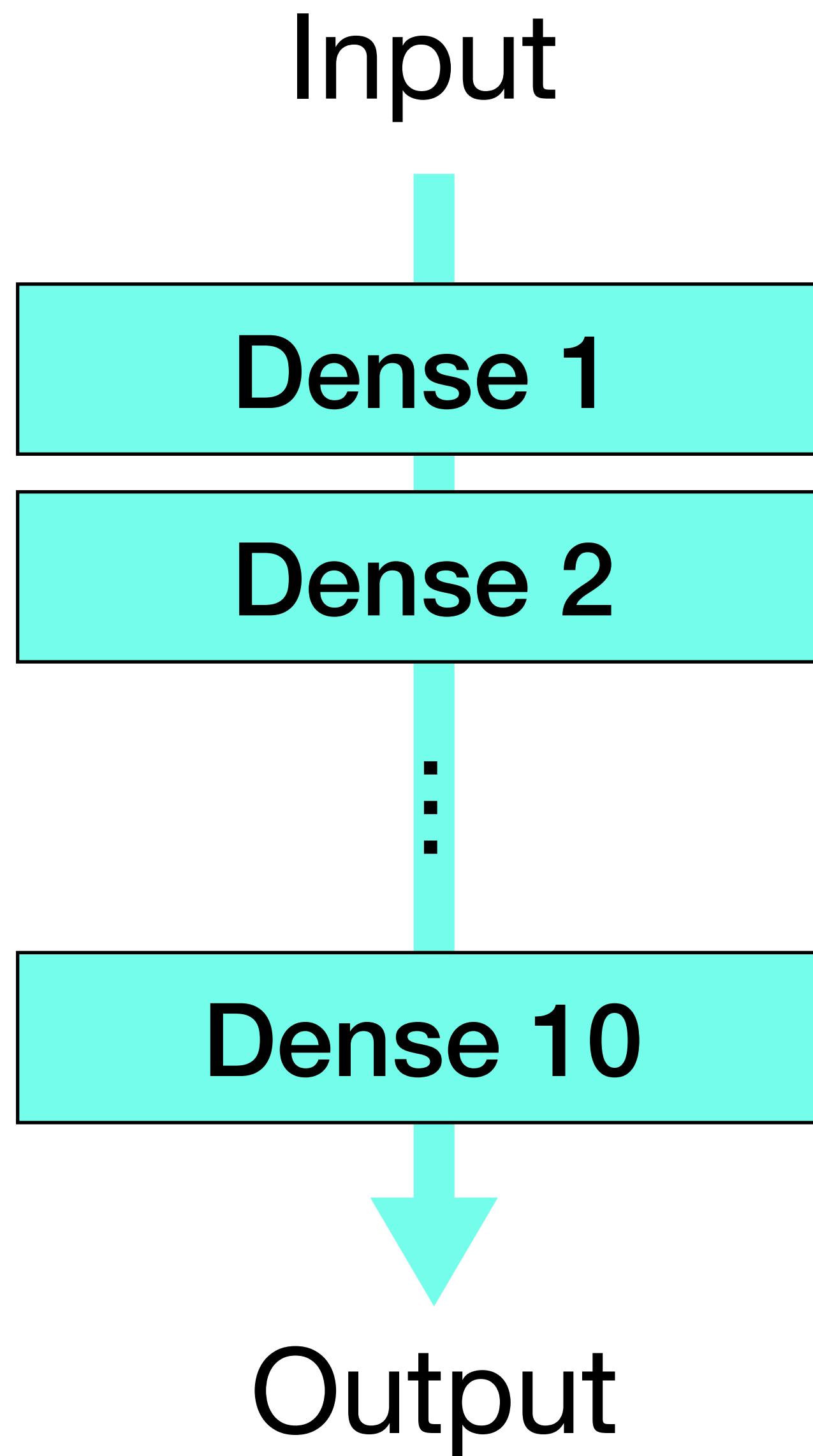
Input



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



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

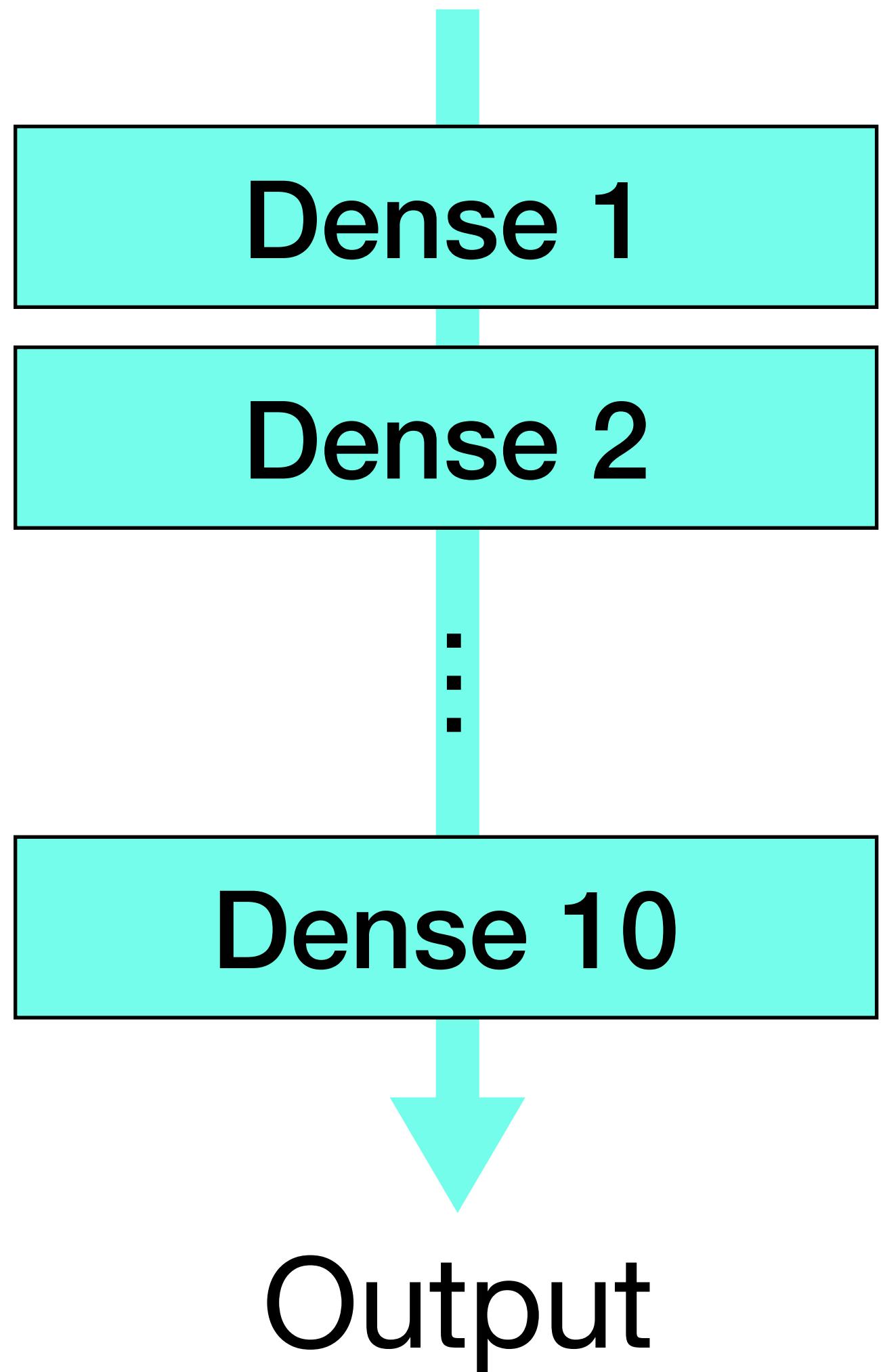


$$\frac{\partial \text{LOSS}}{\partial W_1}$$

$$= \frac{\partial \text{Loss}}{\partial O_{10}} \frac{\partial}{\partial W_1} (\sigma(W_{10} \cdot O_9))$$

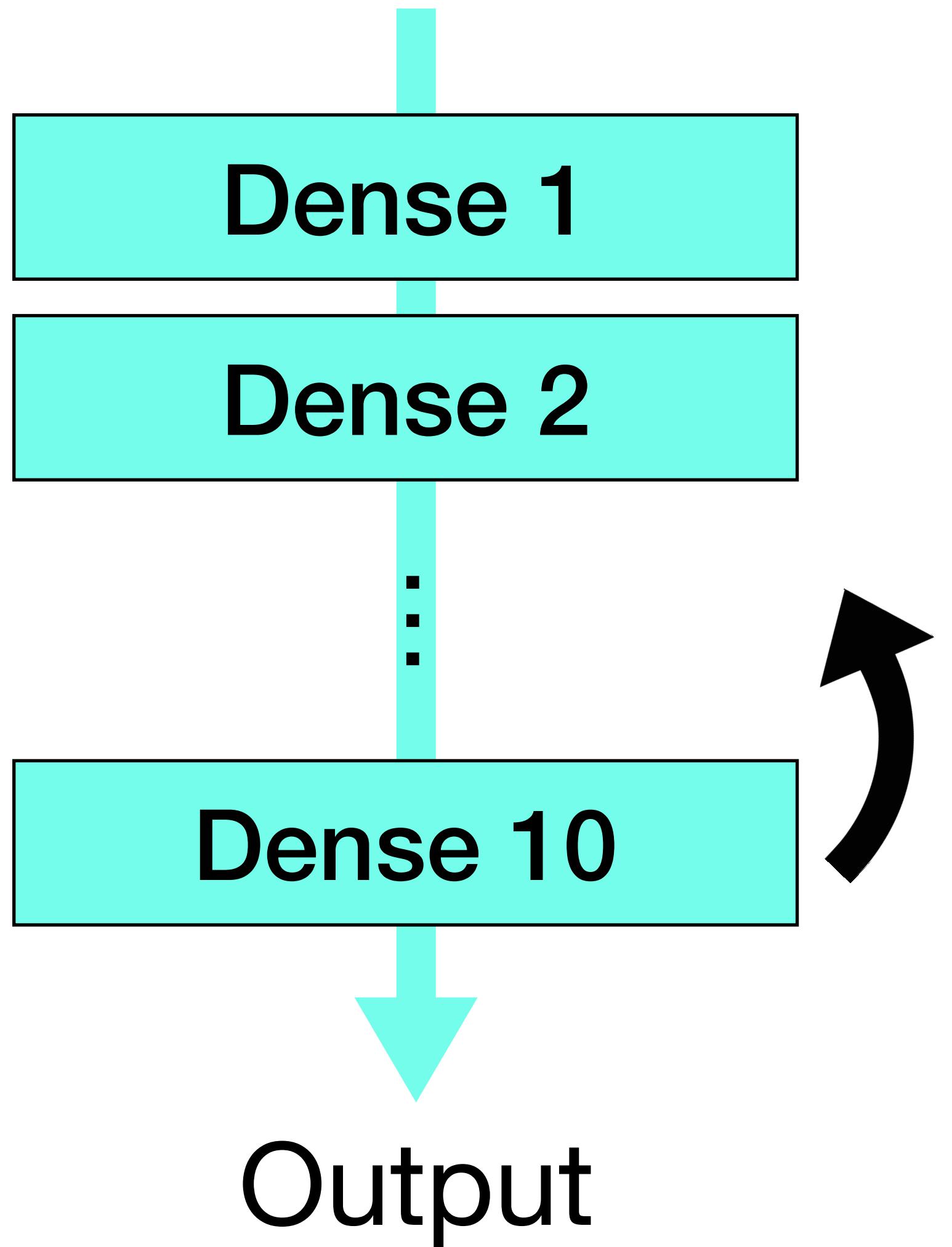
CENSORED

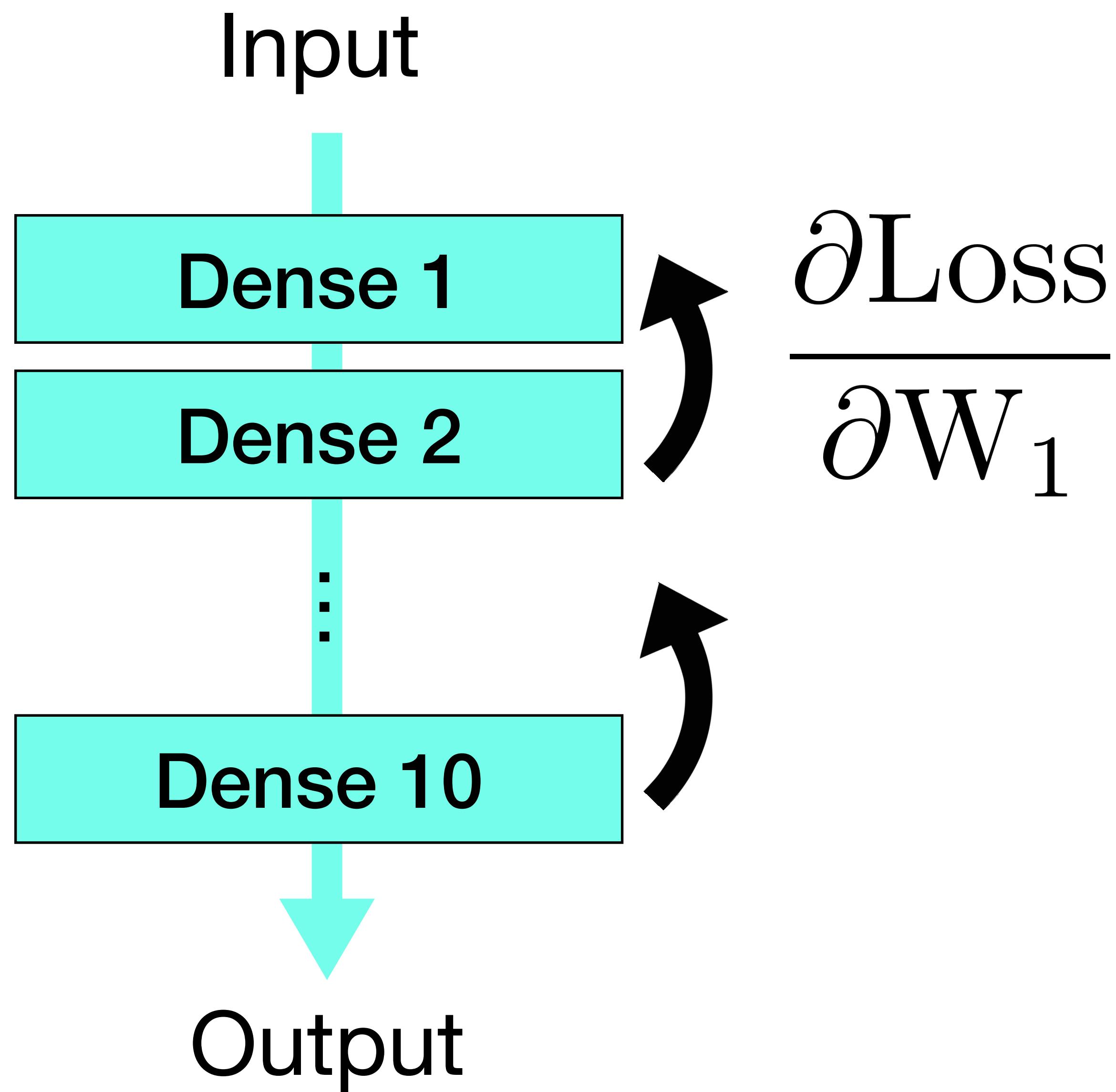
Input

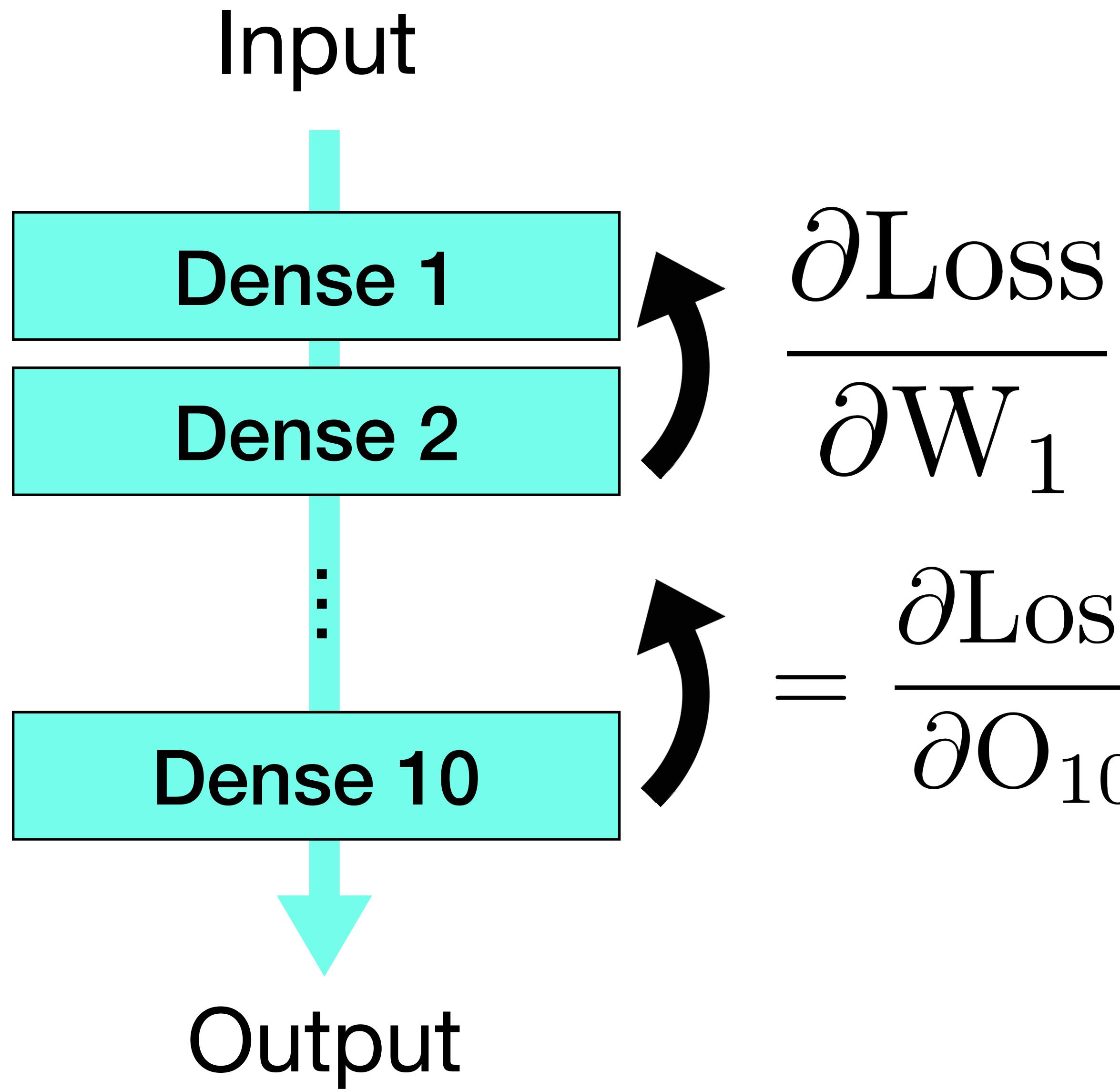


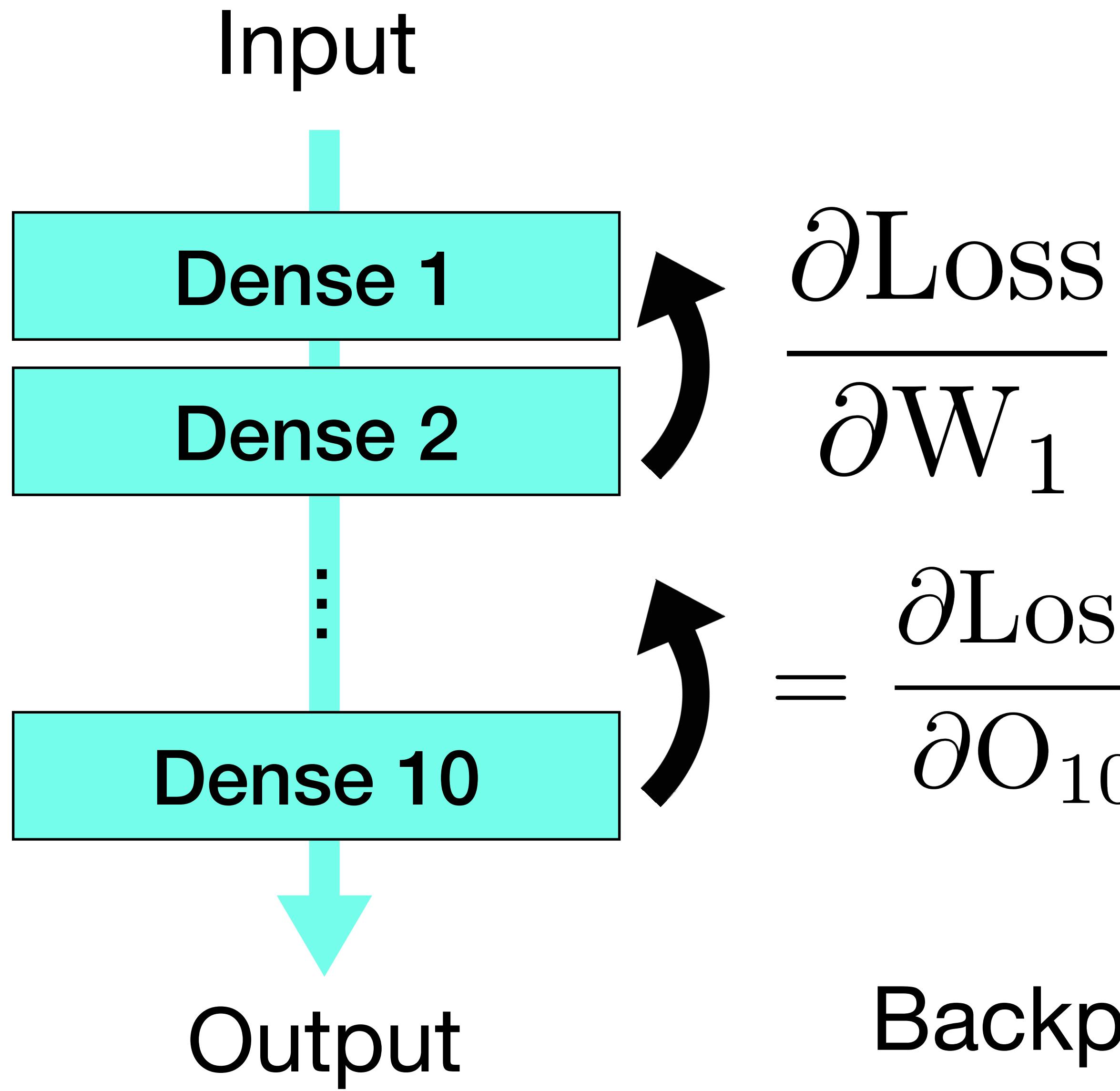
$$\frac{\partial \text{LOSS}}{\partial W_1}$$

Input

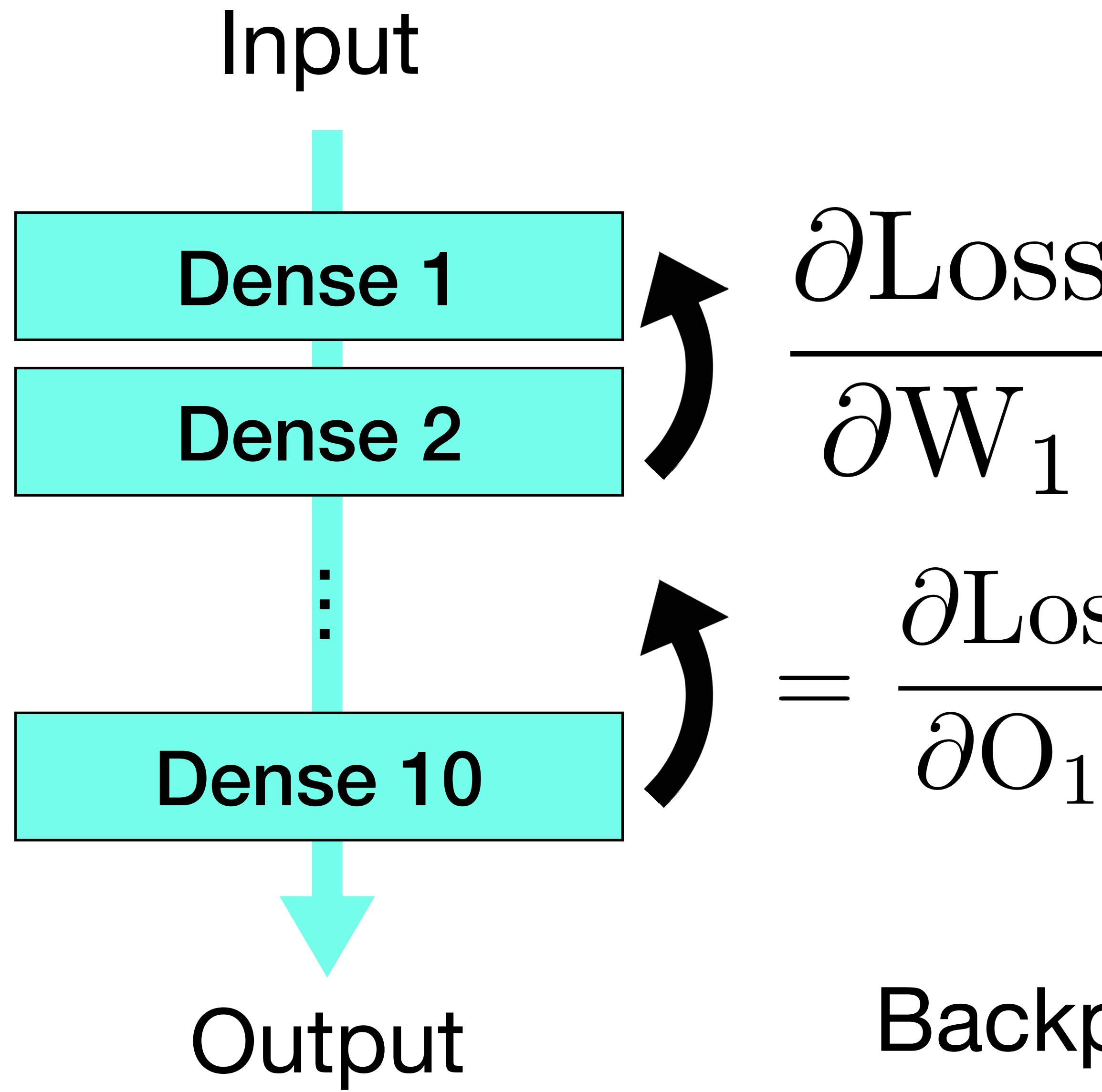






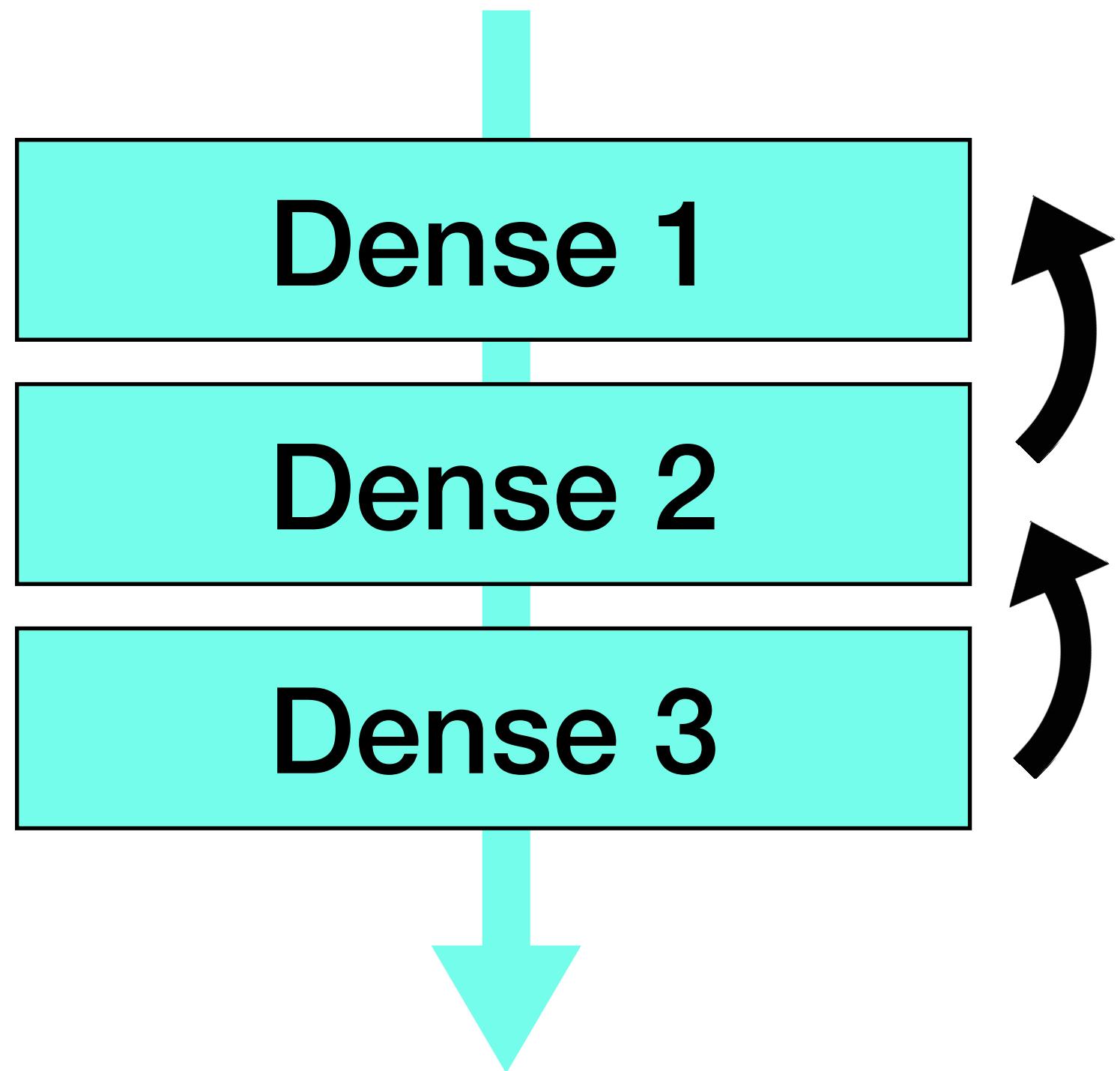


Backpropagation



Backpropagation

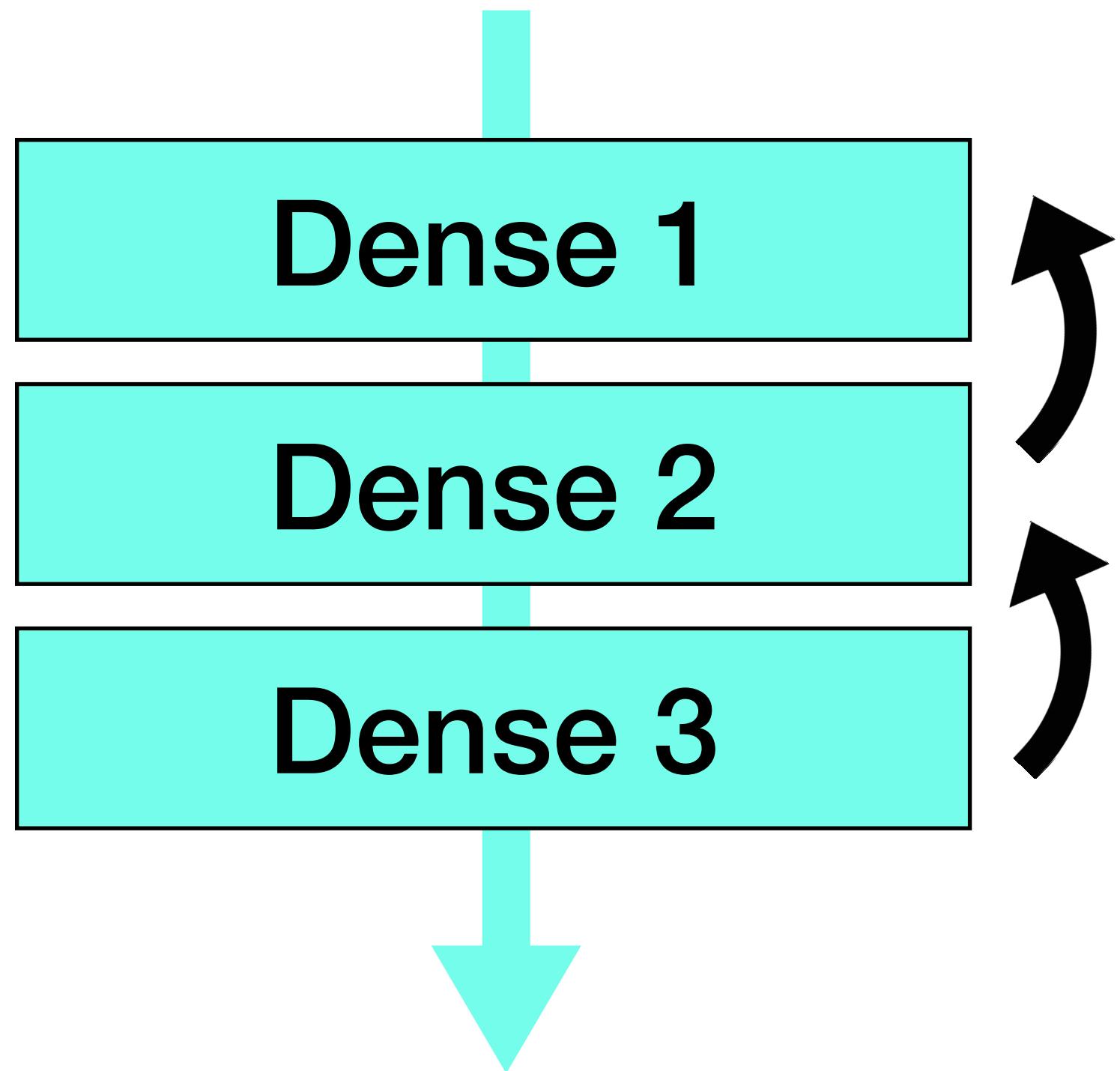
Input



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



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

Output

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

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

A single layer

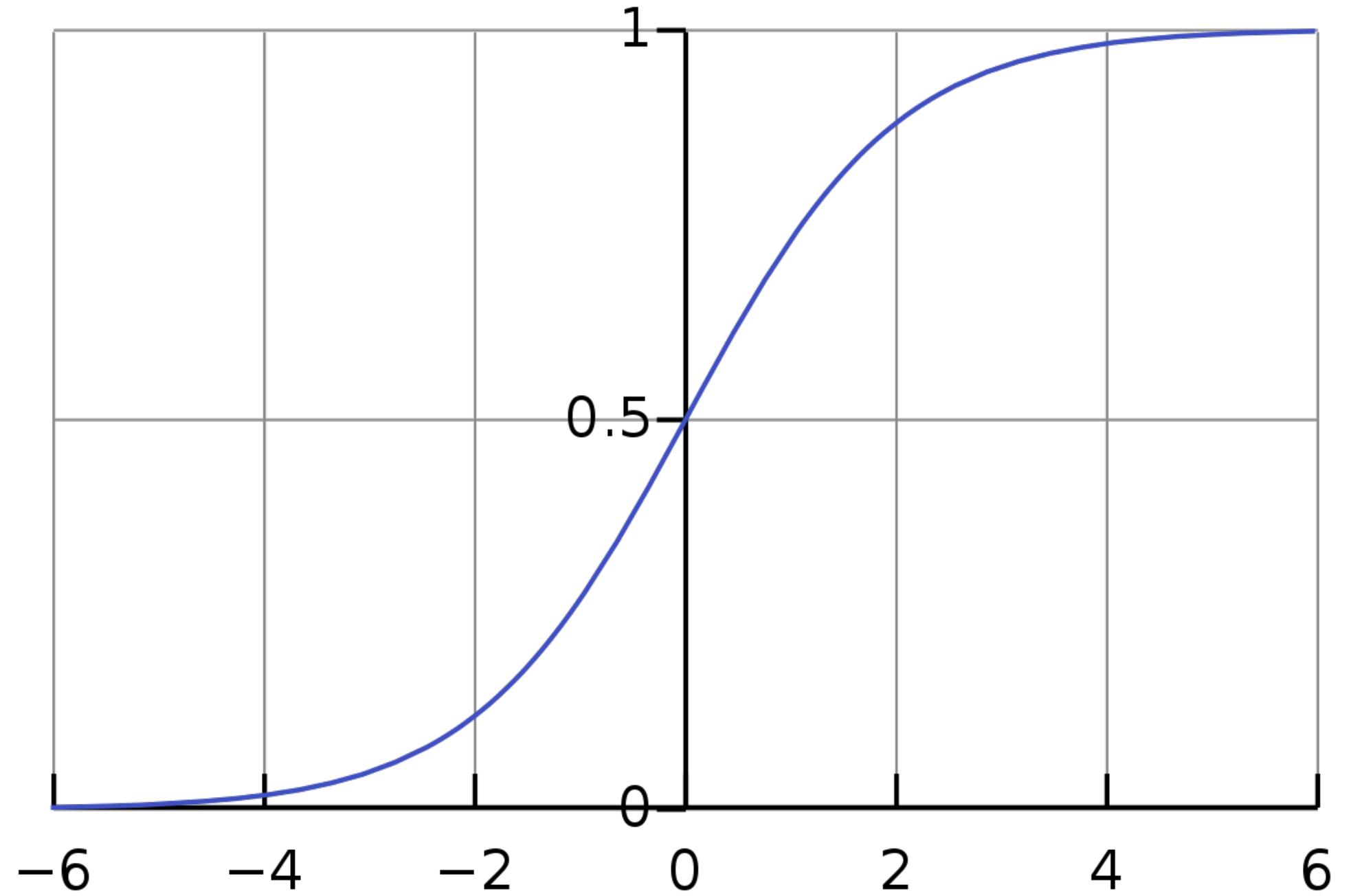
$$\frac{\partial O_3}{\partial O_2}$$

A single layer

Linear layer

$$Wx + b$$

Non-linearity



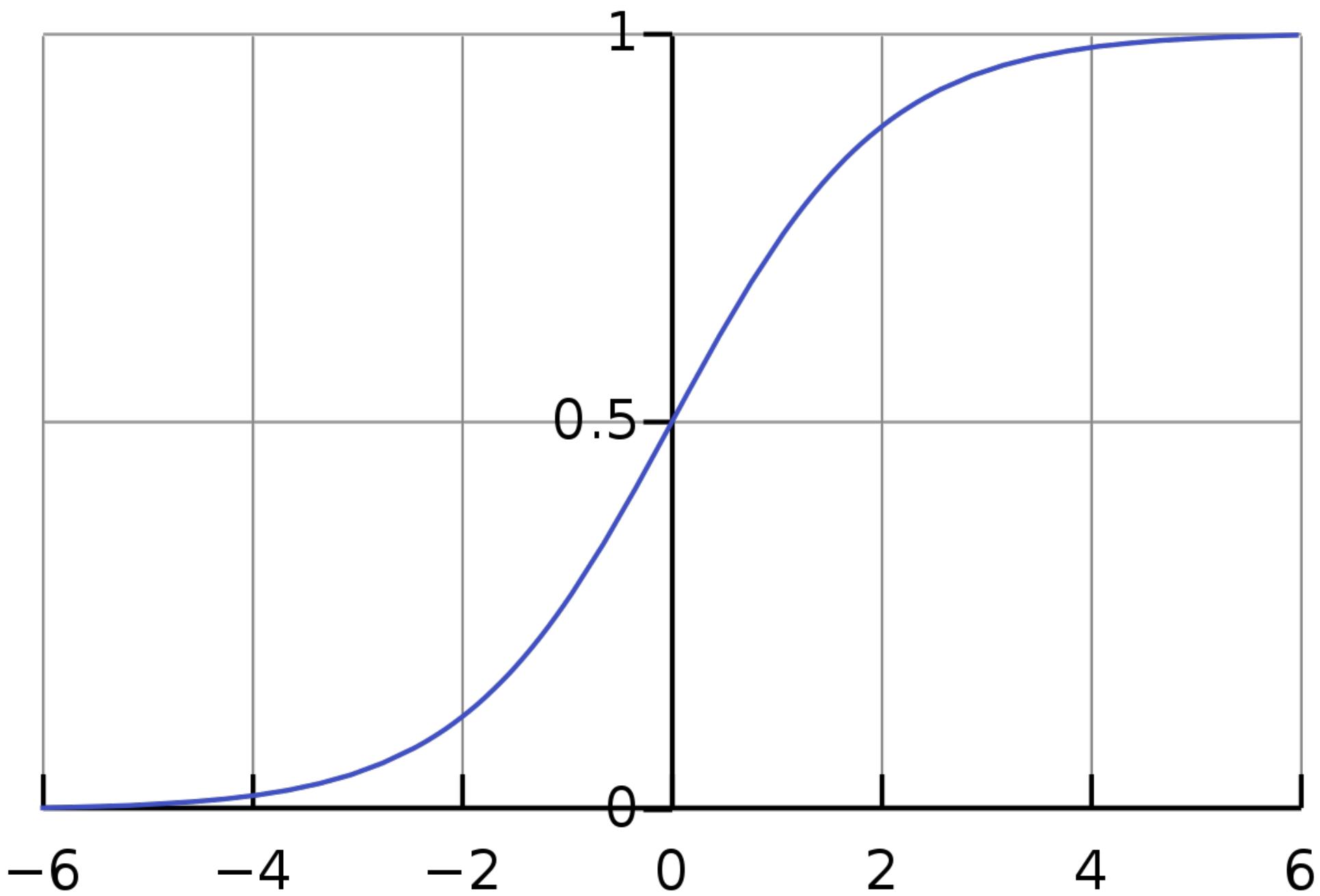
$$\frac{\partial O_3}{\partial O_2}$$

A single layer

Linear layer

$$Wx + b$$

x

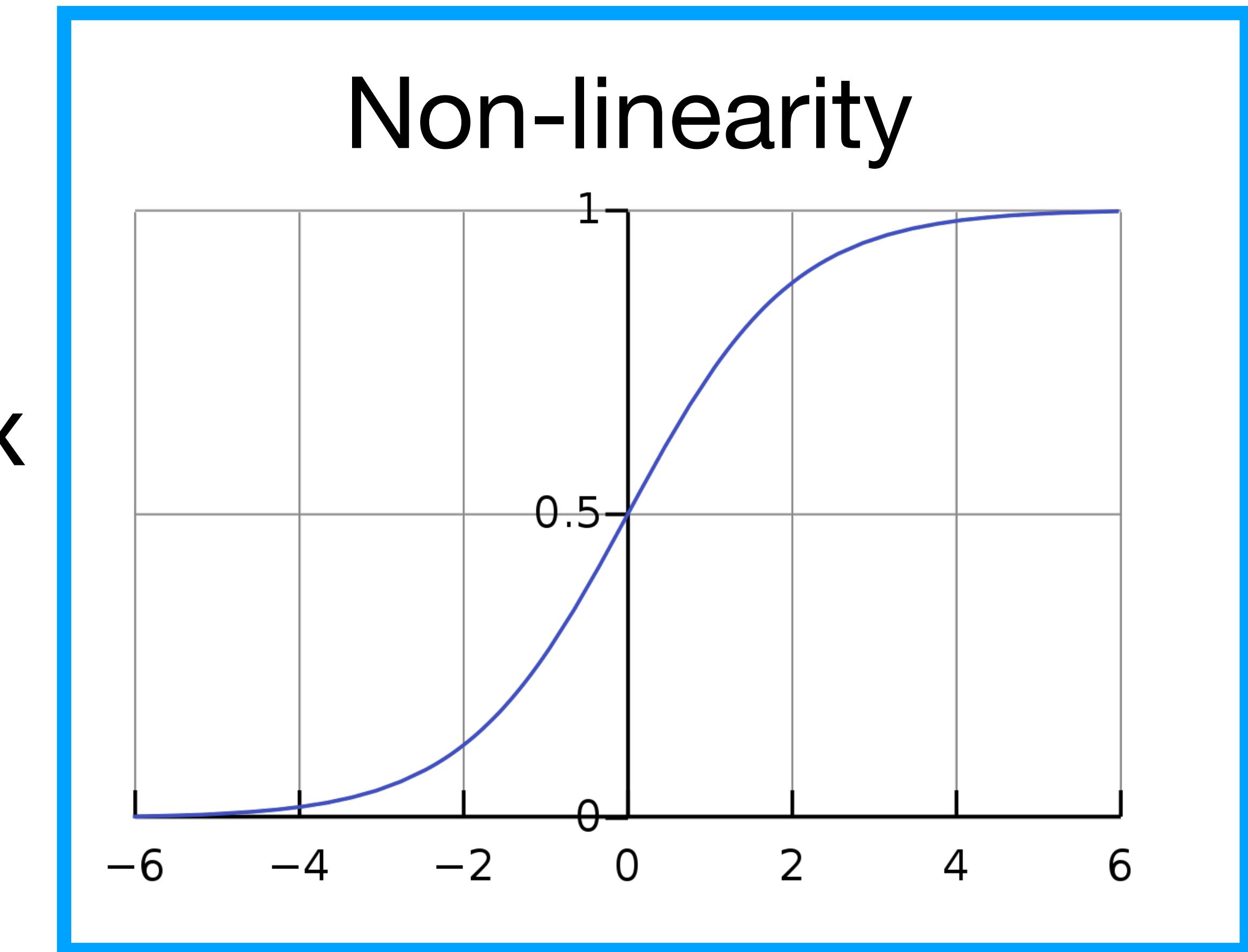


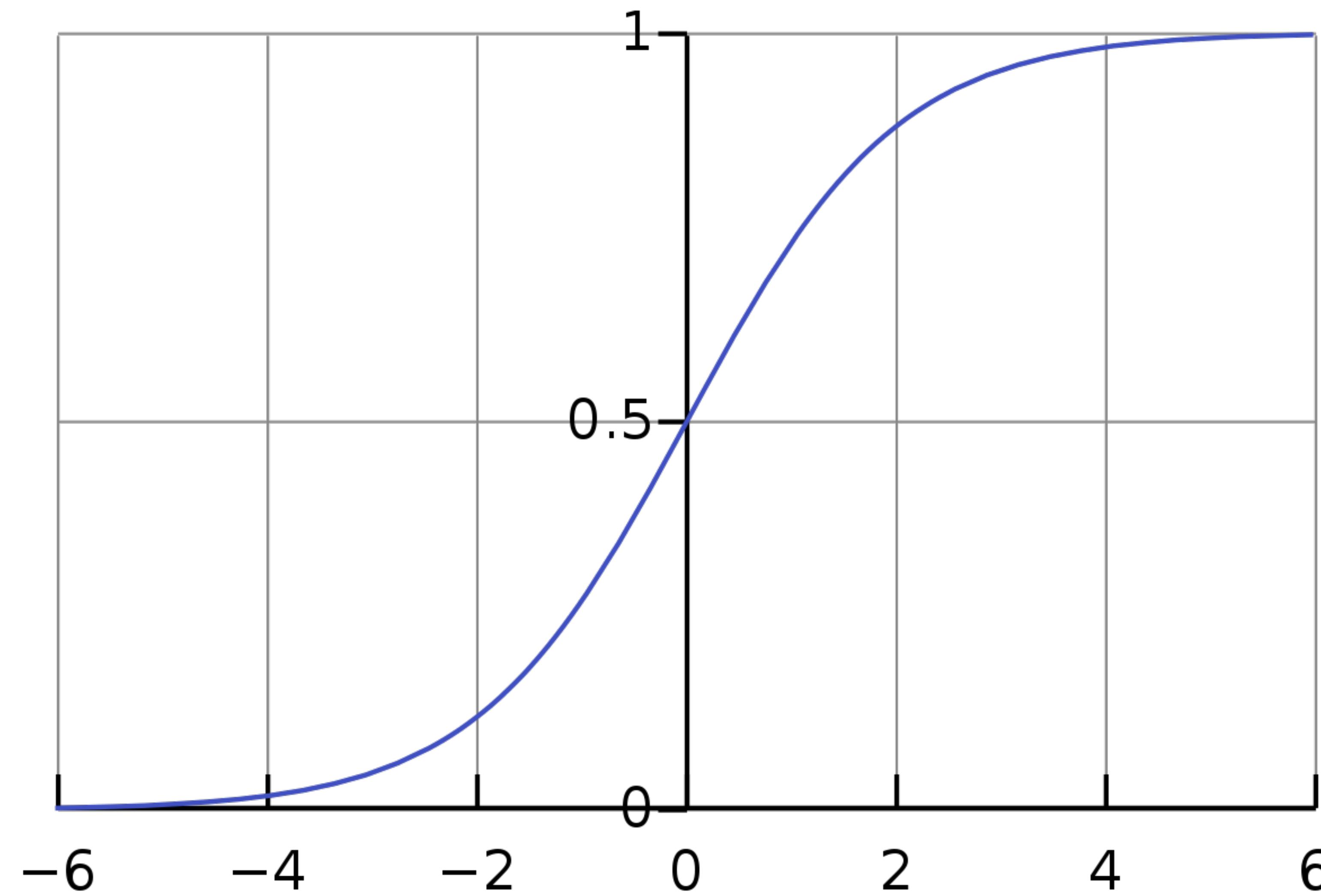
$$\frac{\partial O_3}{\partial O_2}$$

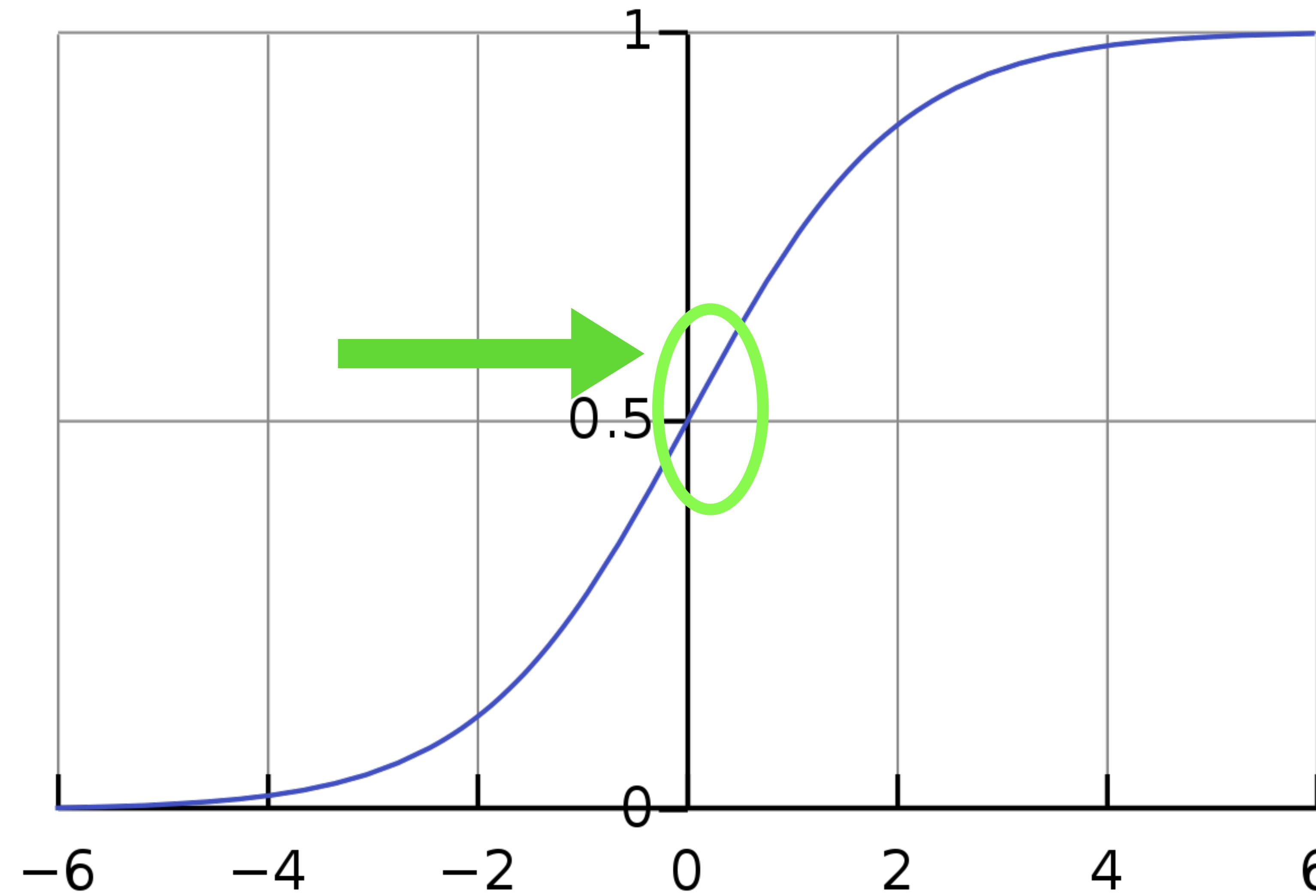
A single layer

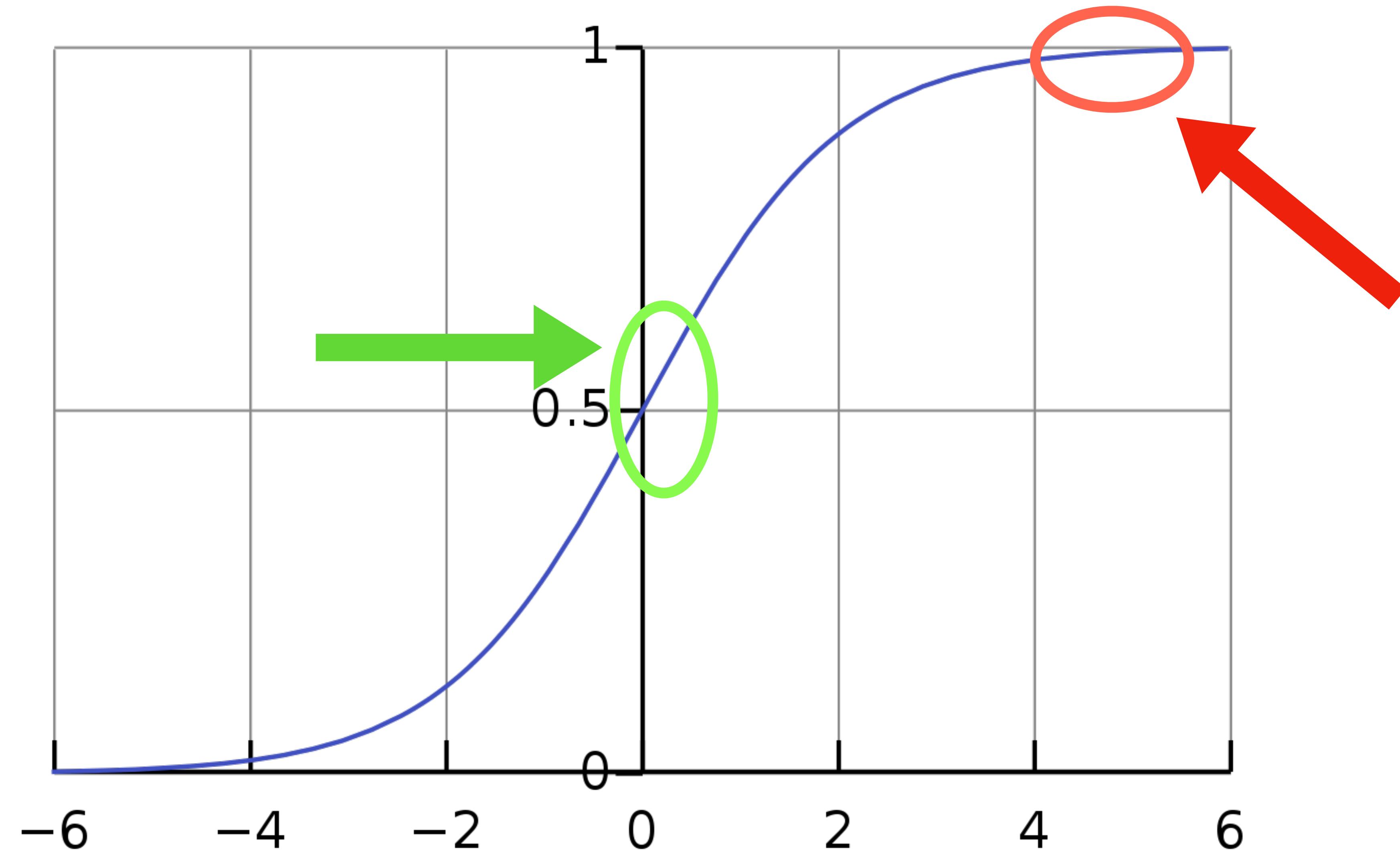
Linear layer

$$Wx + b$$









**Initialise weights with
SMALL variance!**

**Initialise weights with
SMALL variance!**

e.g. Glorot initialisation

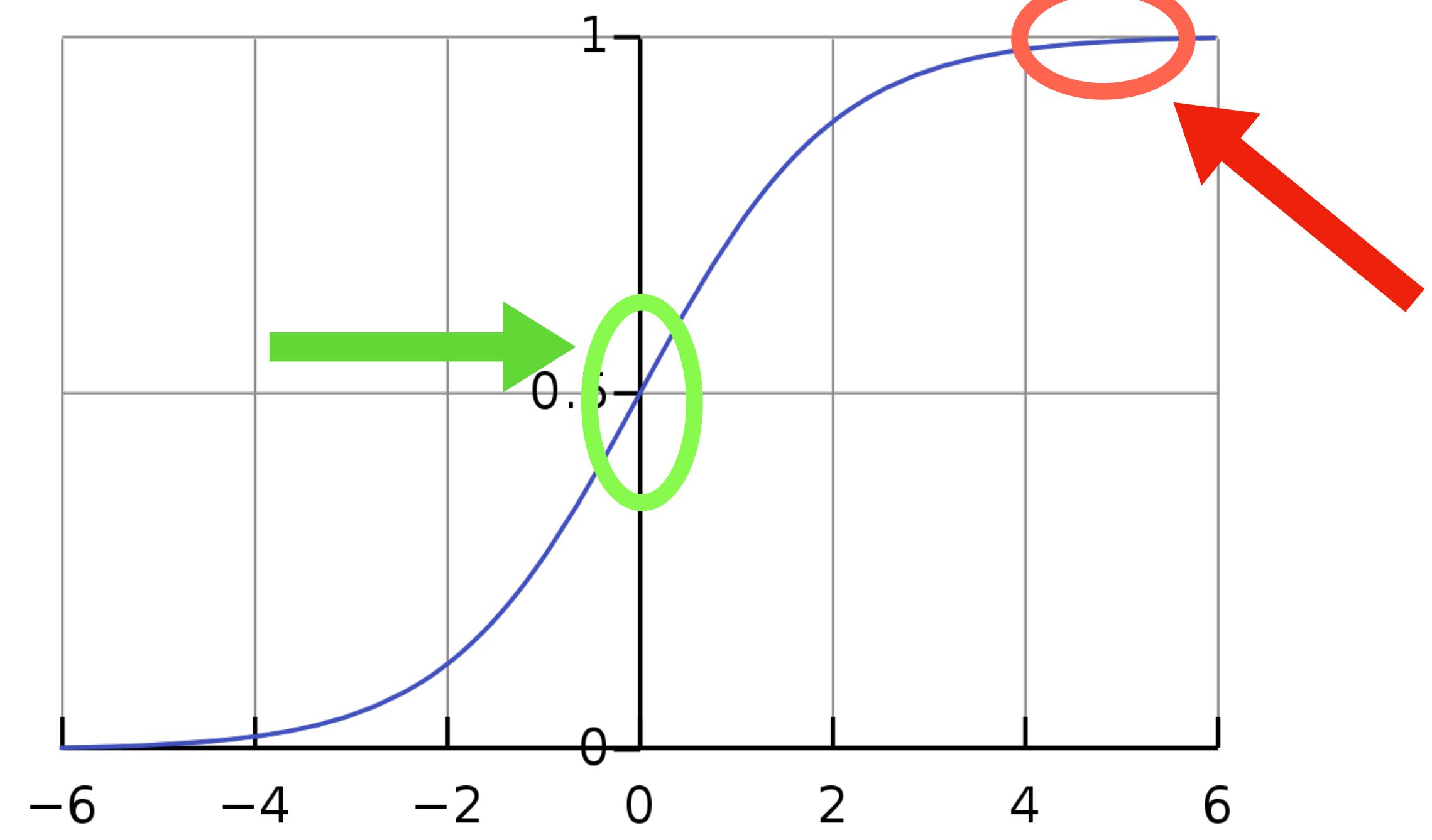
$$\frac{\partial O_3}{\partial O_2}$$

A single layer

Linear layer

$$Wx + b$$

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

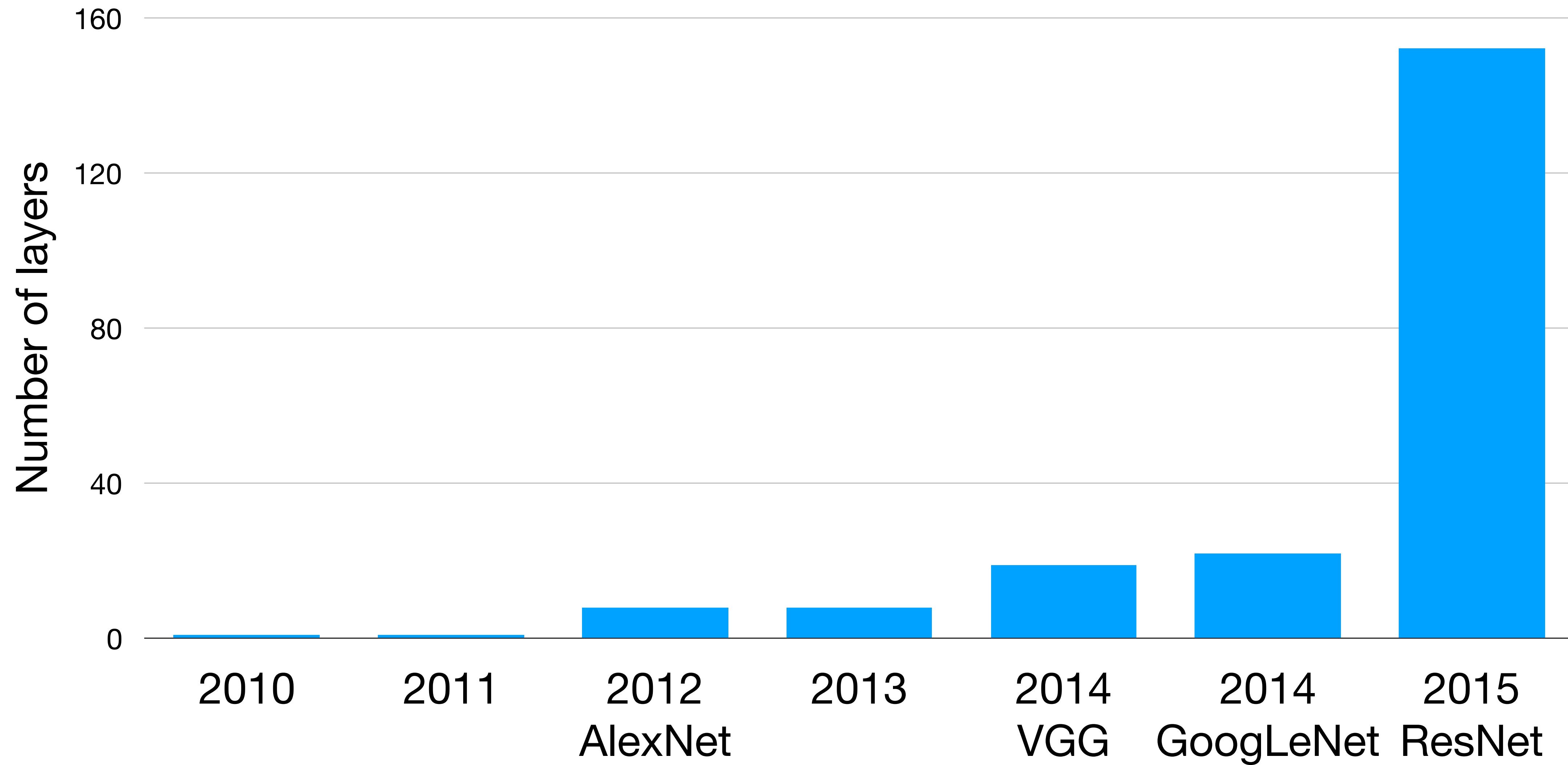
leopard
cat snow leopard

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.

$$W = W - \alpha \boxed{\nabla_W f}$$

Parameters

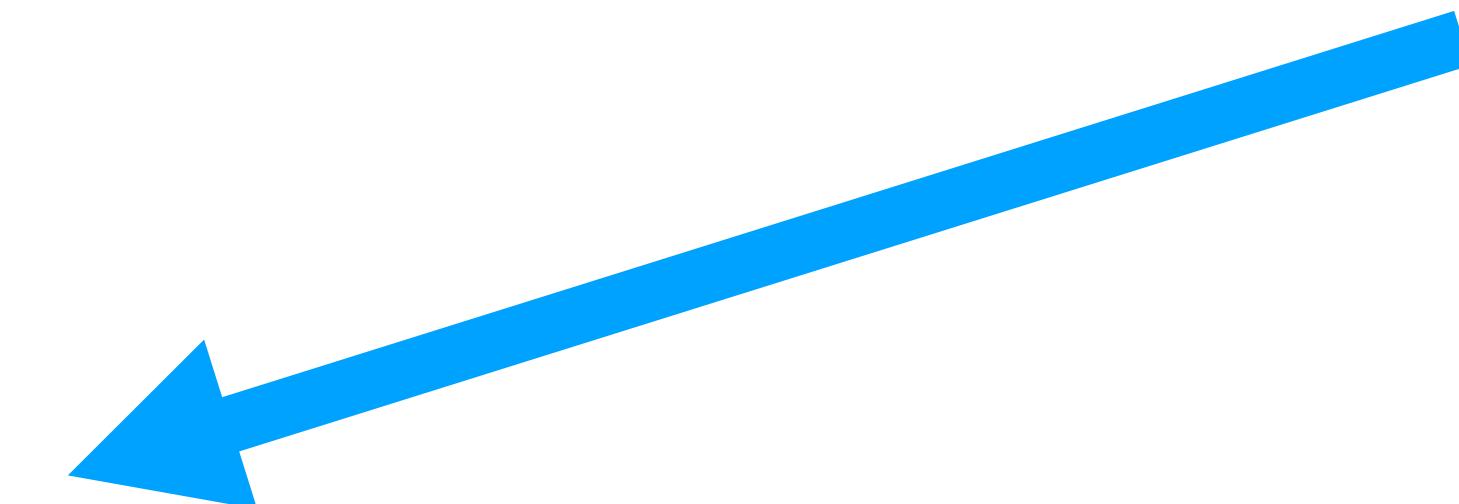
Learning rate
(step size)

Gradient

The diagram illustrates the gradient descent update rule. It shows the formula $W = W - \alpha \nabla_W f$. Three blue arrows point to the components: one to the vector W , one to the gradient term $\nabla_W f$, and one to the learning rate α .

$$W = W - \alpha \boxed{\nabla_W f}$$

$$W = W - \alpha \boxed{\nabla_W f}$$



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

$$W = W - \alpha \boxed{\nabla_W f}$$

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

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

$$0.1 \boxed{\frac{\partial O_{100}}{\partial O_{99}}}$$

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

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

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

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

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

$$= 10^{-100}$$

Vanishing gradients

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

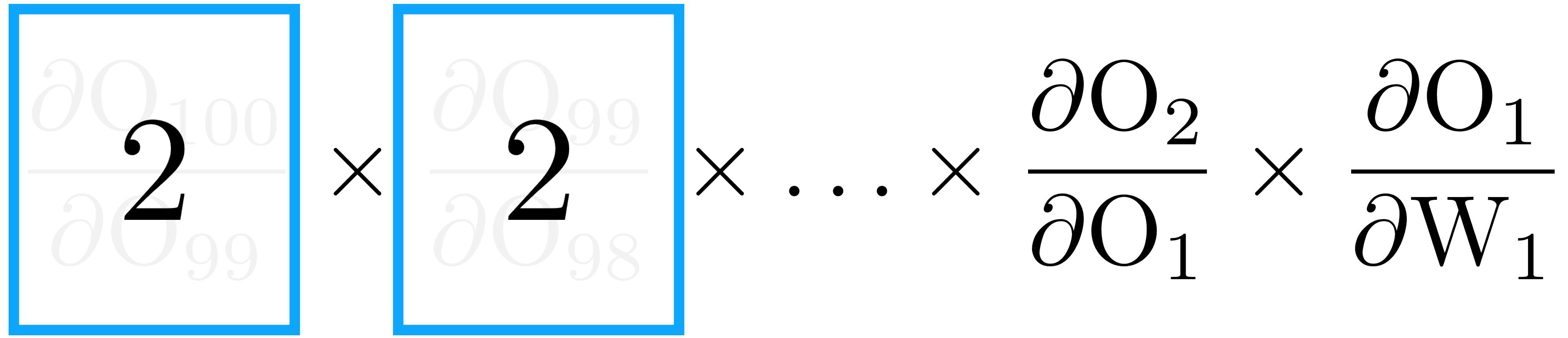
$$= 10^{-100}$$

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

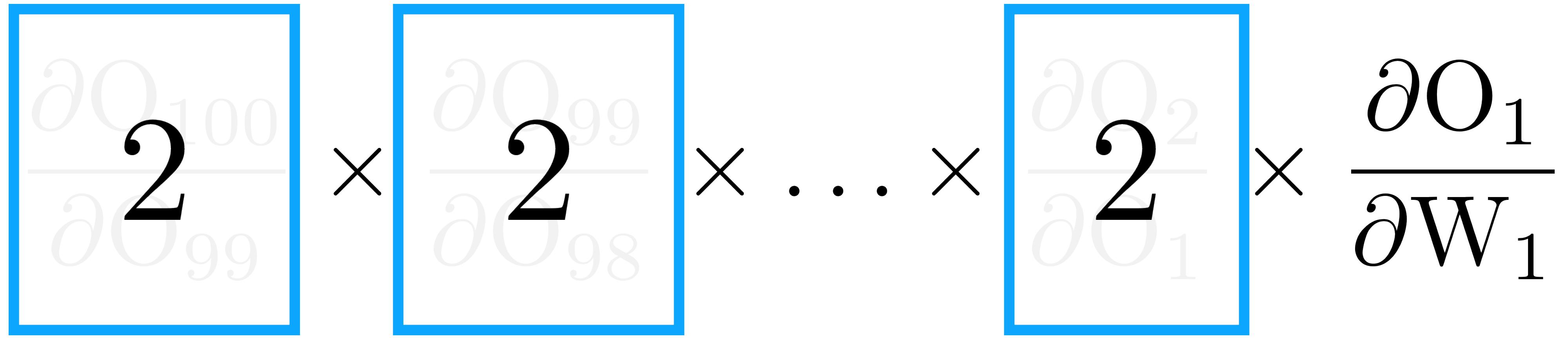
$$\frac{\partial O_{100}}{\partial O_{99}}$$

2

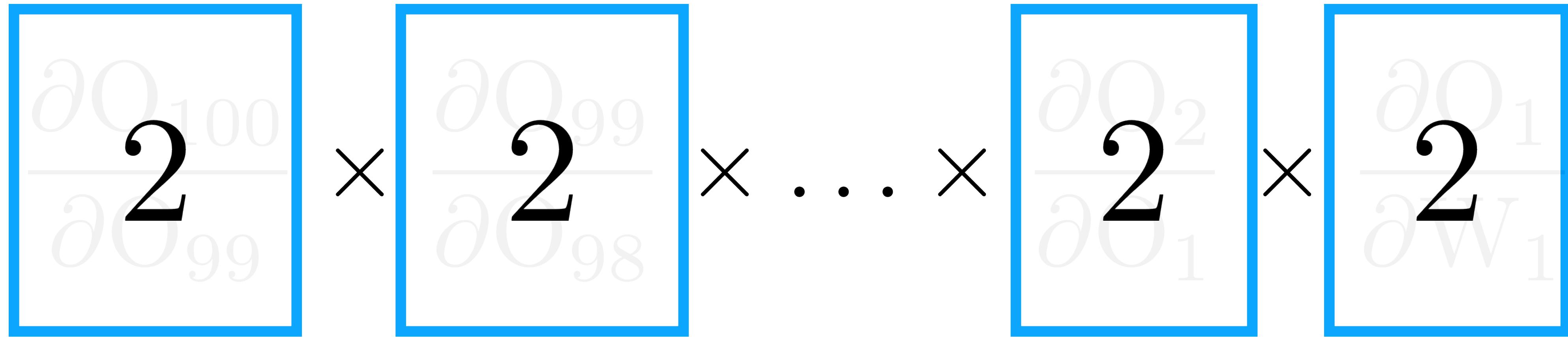
$$\times \frac{\partial O_{99}}{\partial O_{98}} \times \dots \times \frac{\partial O_2}{\partial O_1} \times \frac{\partial O_1}{\partial W_1}$$

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


The diagram illustrates a mathematical product of partial derivatives. It consists of a sequence of terms separated by multiplication signs. Each term is a fraction where the numerator is a variable (O_i) and the denominator is the previous variable in the sequence (O_{i-1}). There are two prominent visual elements: two white rectangular boxes with a blue border, each containing the number '2'. The first box is located to the left of the first fraction, and the second box is located to the left of the second fraction. These boxes likely represent specific values or coefficients in the derivative chain.

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


The diagram illustrates a product of three terms, each enclosed in a blue-bordered box. The first term is $\frac{\partial O_{100}}{\partial O_{99}}$, the second is $\frac{\partial O_{99}}{\partial O_{98}}$, and the third is $\frac{\partial O_2}{\partial O_1}$. Each term is multiplied by a factor of 2, represented by a large black '2' inside the box. The boxes are connected by multiplication signs (\times).

$$\frac{\partial \theta_{100}}{\partial \theta_{99}} \times \frac{\partial \theta_{99}}{\partial \theta_{98}} \times \dots \times \frac{\partial \theta_2}{\partial \theta_1} \times \frac{\partial \theta_1}{\partial w_1}$$


$$\frac{2}{\partial \theta_{100}} \times \frac{2}{\partial \theta_{99}} \times \dots \times \frac{2}{\partial \theta_1} \times \frac{2}{\partial w_1}$$

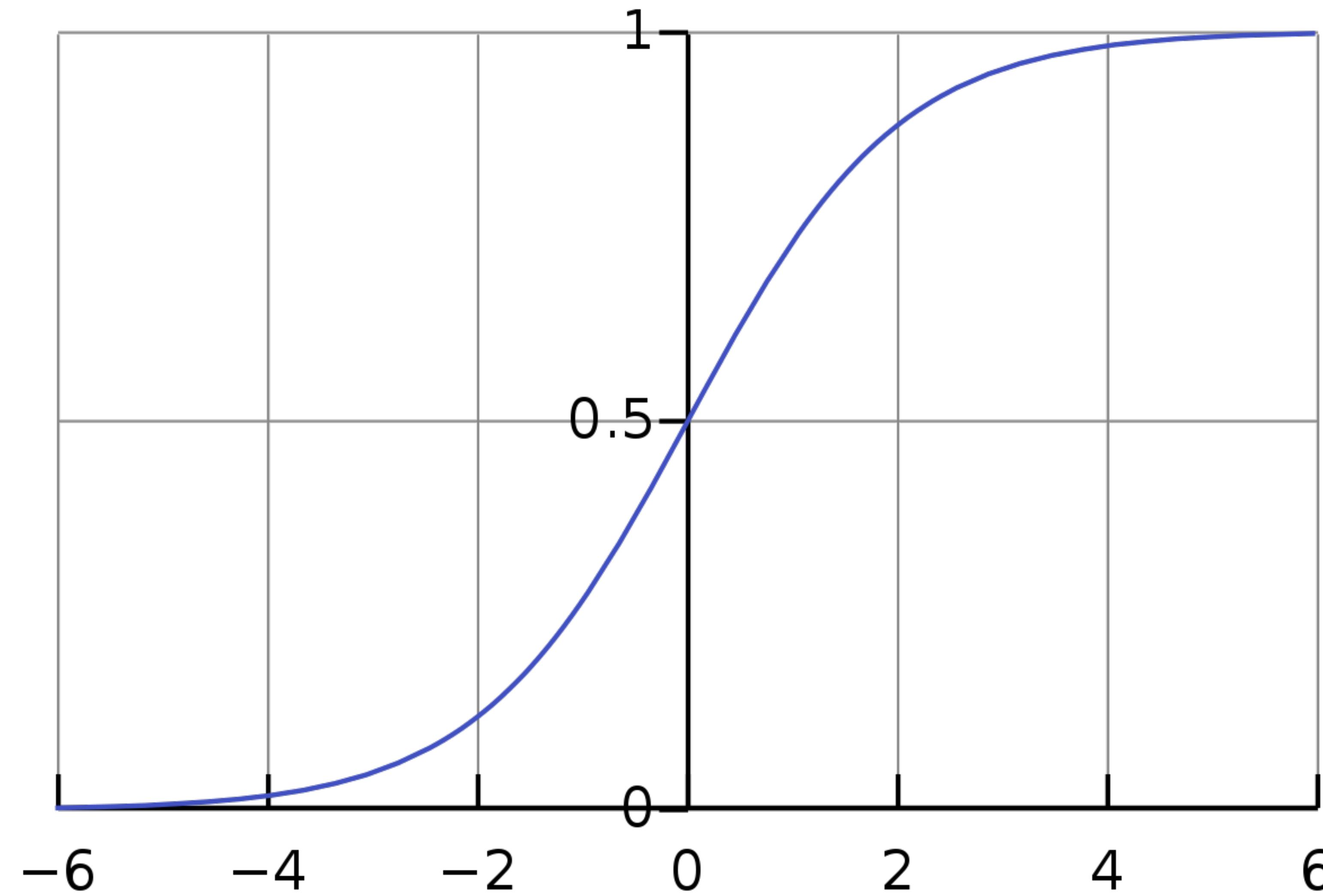
$> 10^{30}$

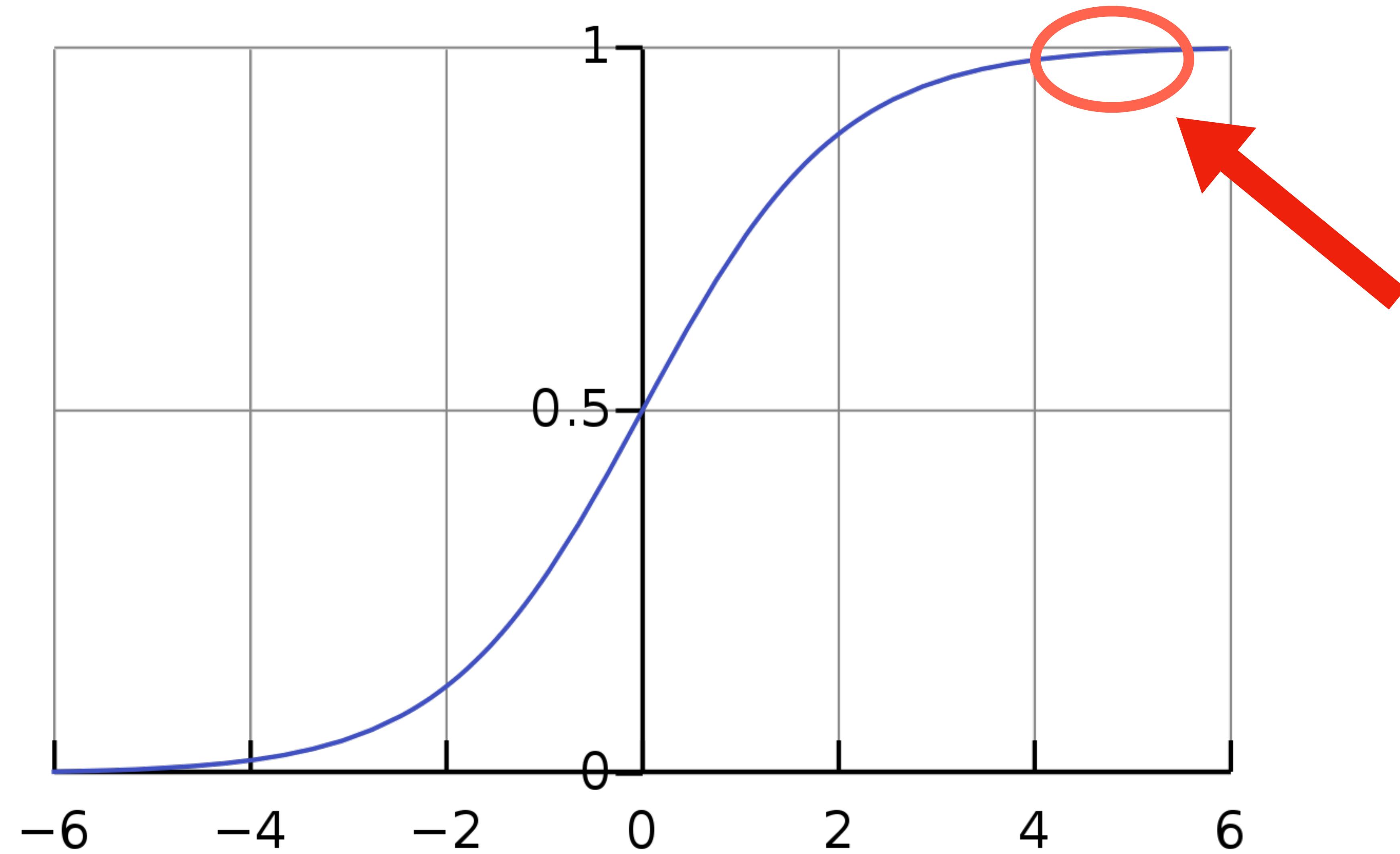
Exploding gradients

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

The diagram illustrates the concept of exploding gradients in a neural network. It shows a sequence of four derivative terms, each represented by a blue-bordered box containing a large black '2'. The first term is $\frac{\partial O_{100}}{\partial O_{99}}$, the second is $\frac{\partial O_{99}}{\partial O_{98}}$, the third is \dots , and the fourth is $\frac{\partial O_1}{\partial W_1}$. The multiplication symbol '×' is placed between each pair of adjacent terms.

$> 10^{30}$



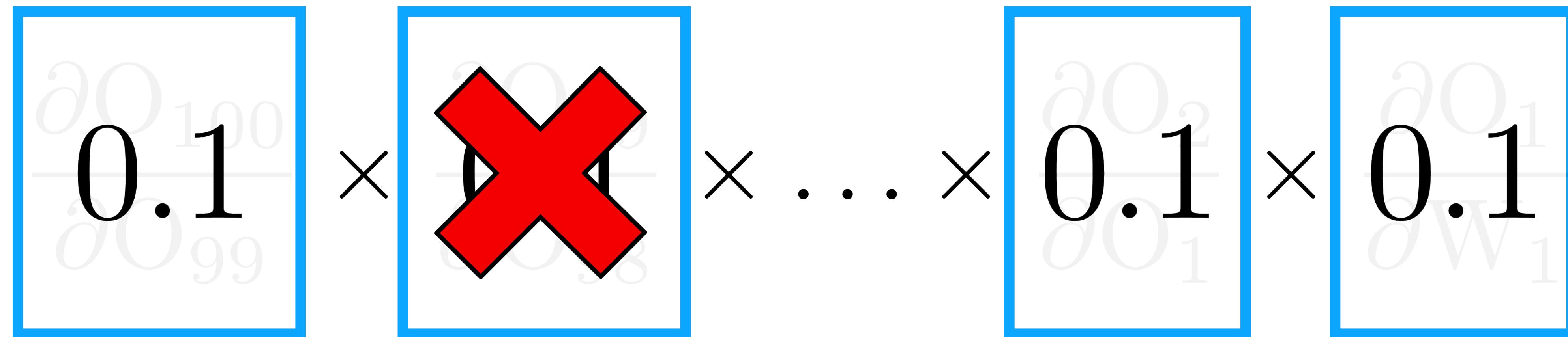


Vanishing gradients

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

$$= 10^{-100}$$

Vanishing gradients



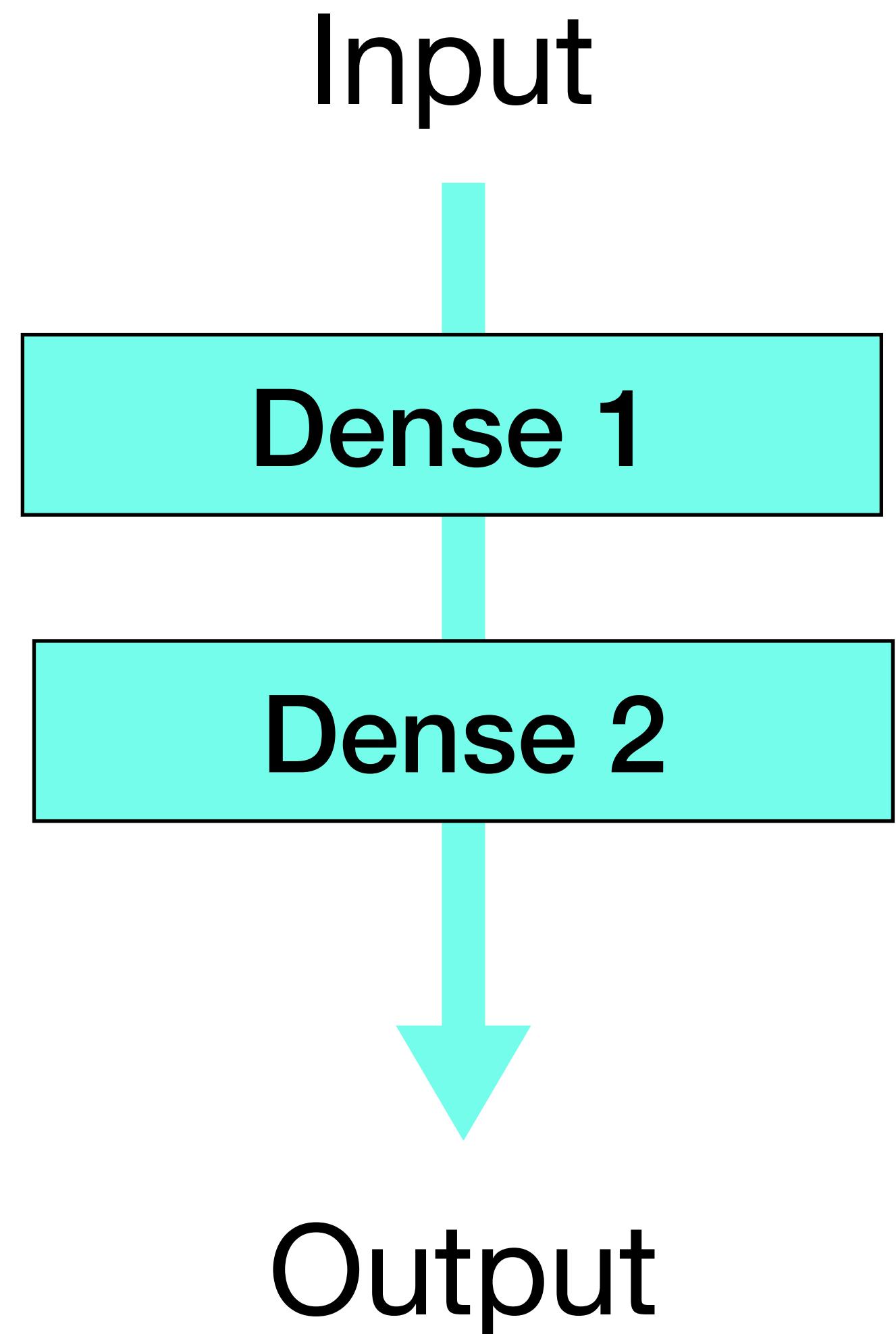
$$= 10^{-100}$$

Vanishing gradients

$$0.1 \times \begin{matrix} 1 \\ \vdots \end{matrix} \times \dots \times 0.1 \times 0.1$$

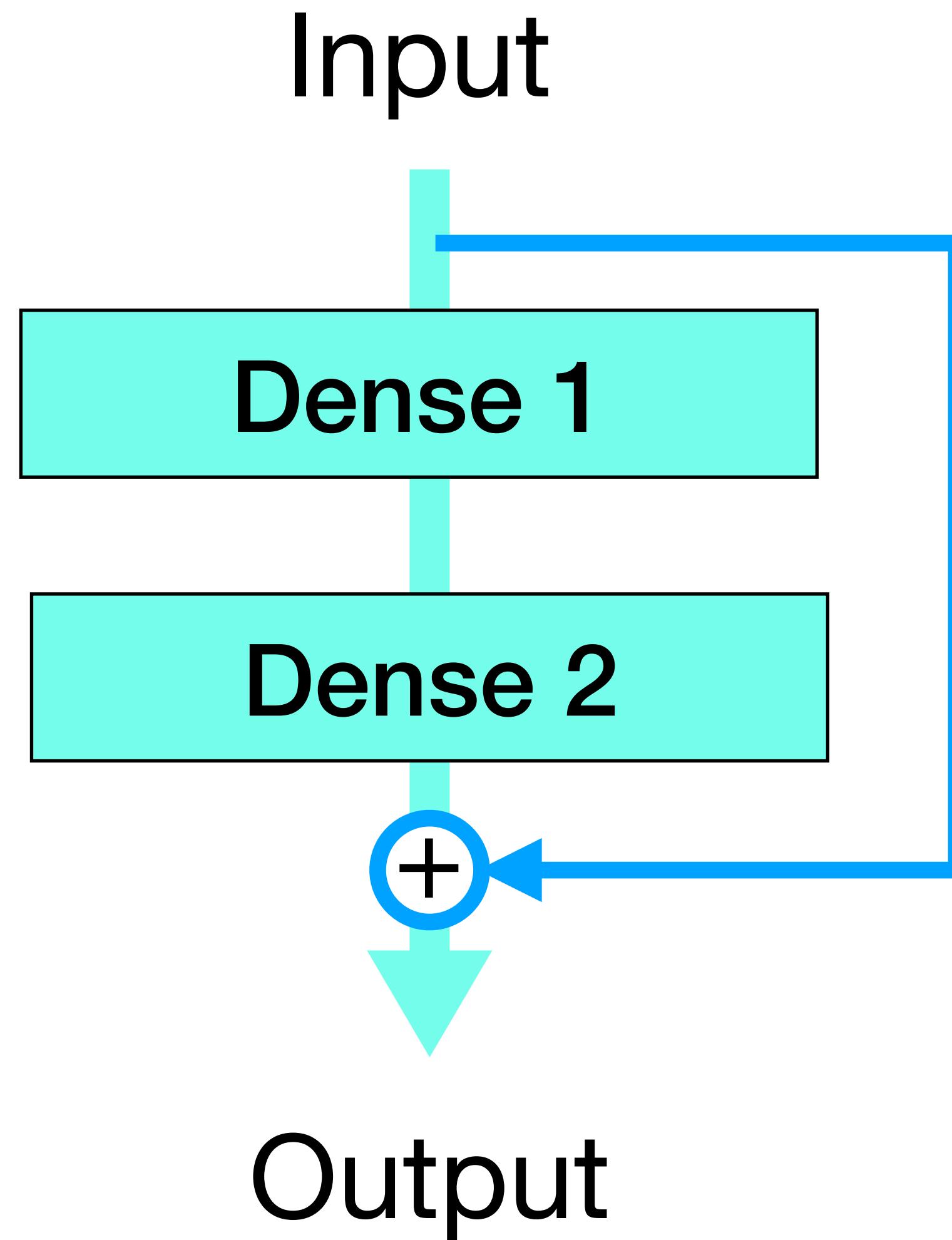
$$= 10^{-100}$$

Skip connections



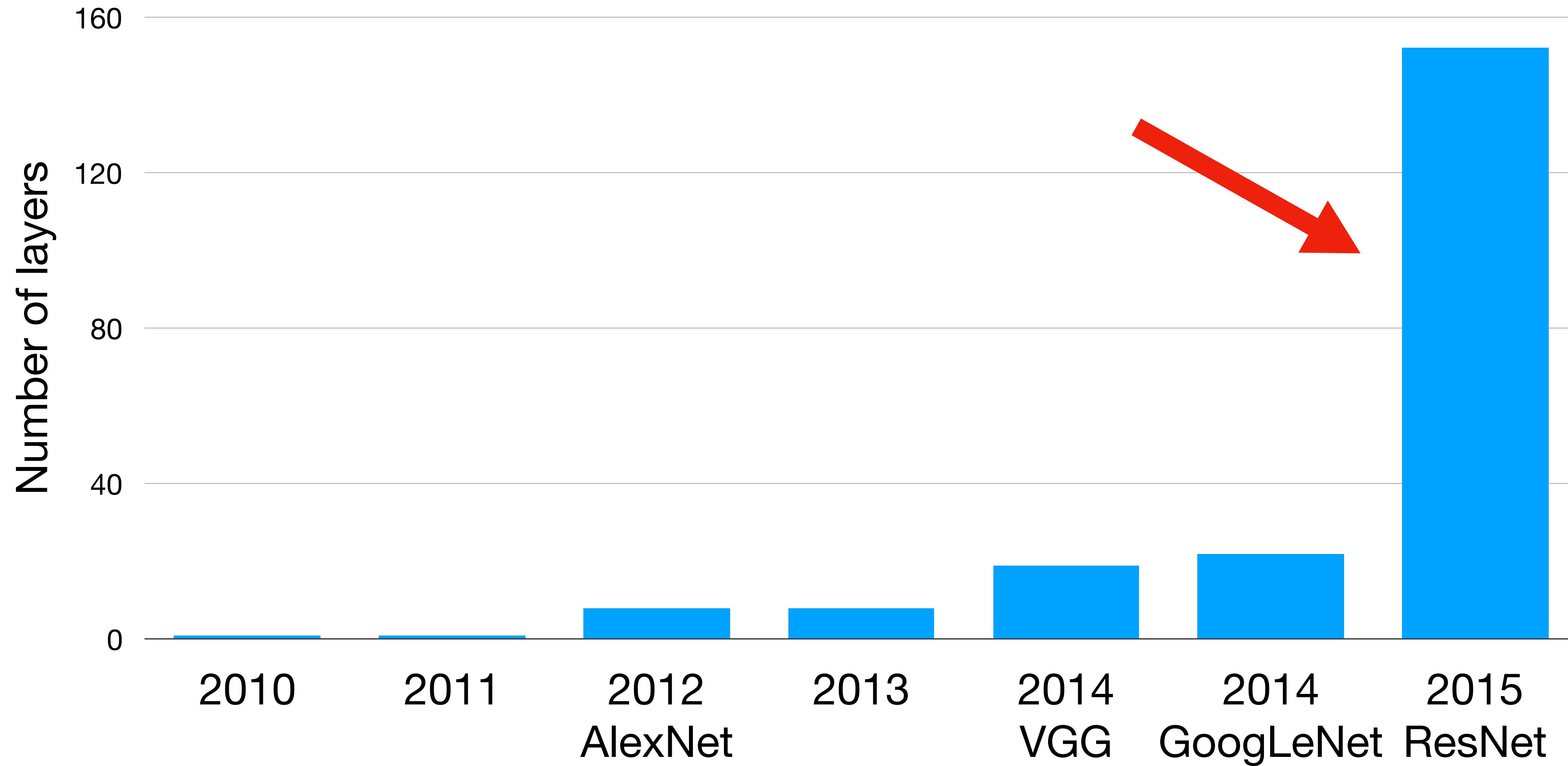
ResNets (He et. al., 2015)

Skip connections



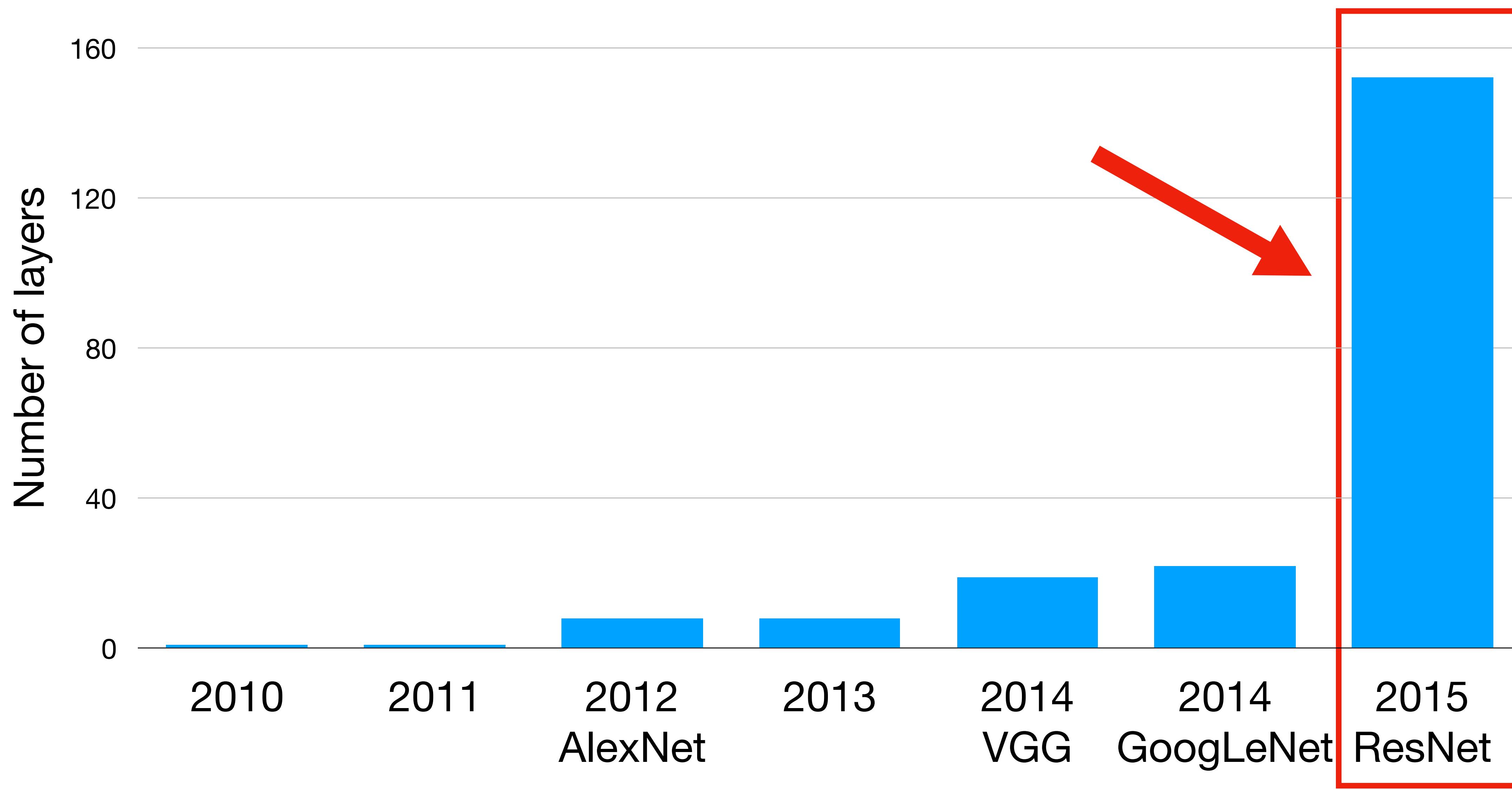
ResNets (He et. al., 2015)

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.

Revolution of Depth: Top models on ImageNet

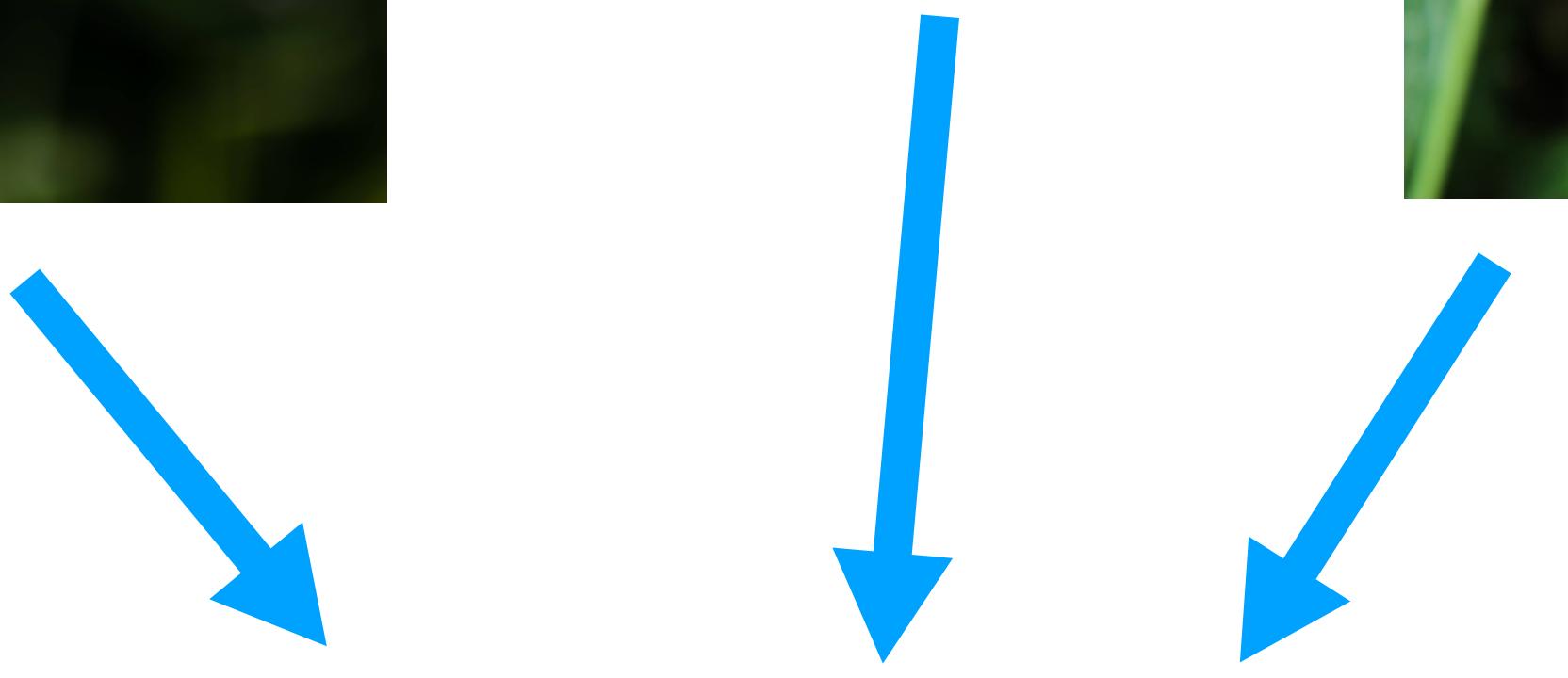
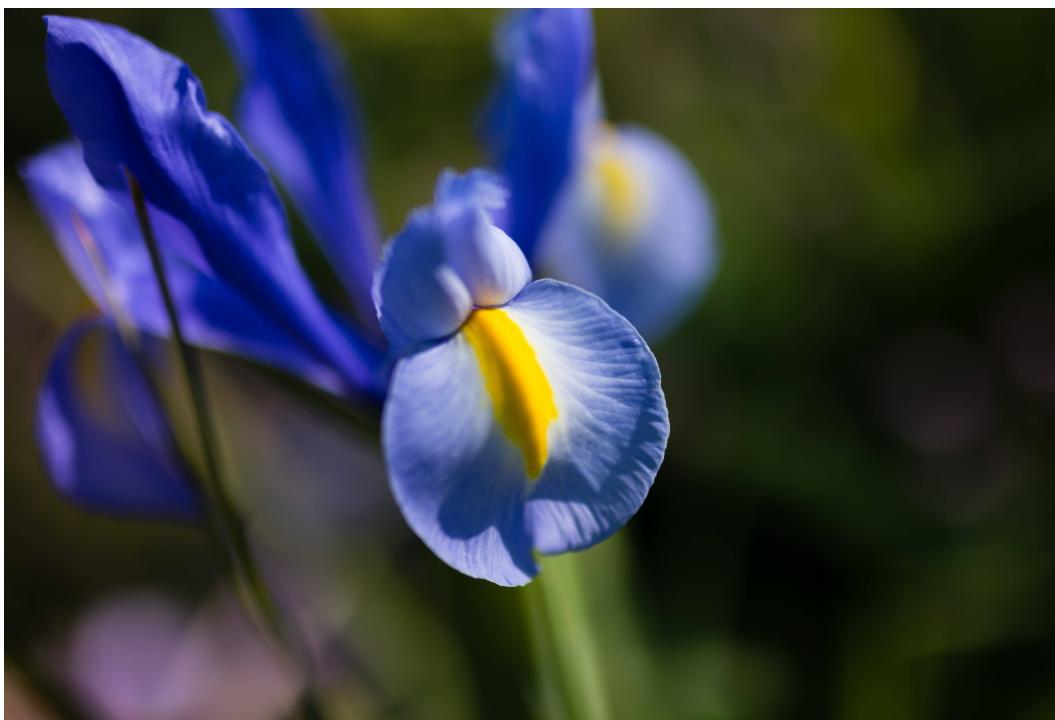


Data from Kaiming He, Xiangyu Zhang, Shaoqing Ren and Jian Sun. 'Deep Residual Learning for Image Recognition', CVPR 2016.

IV. Practical Tips

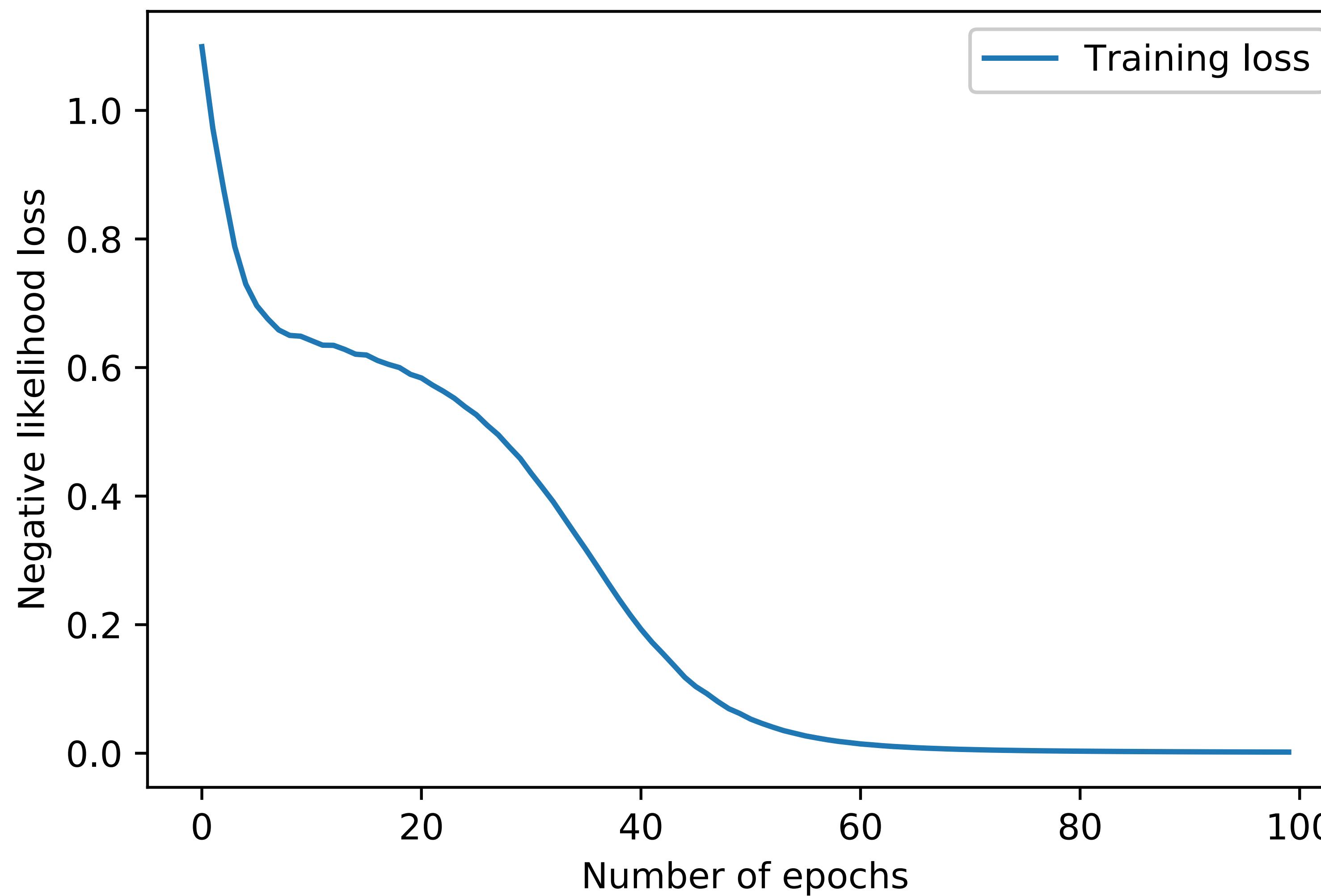
Overfit First

a few examples

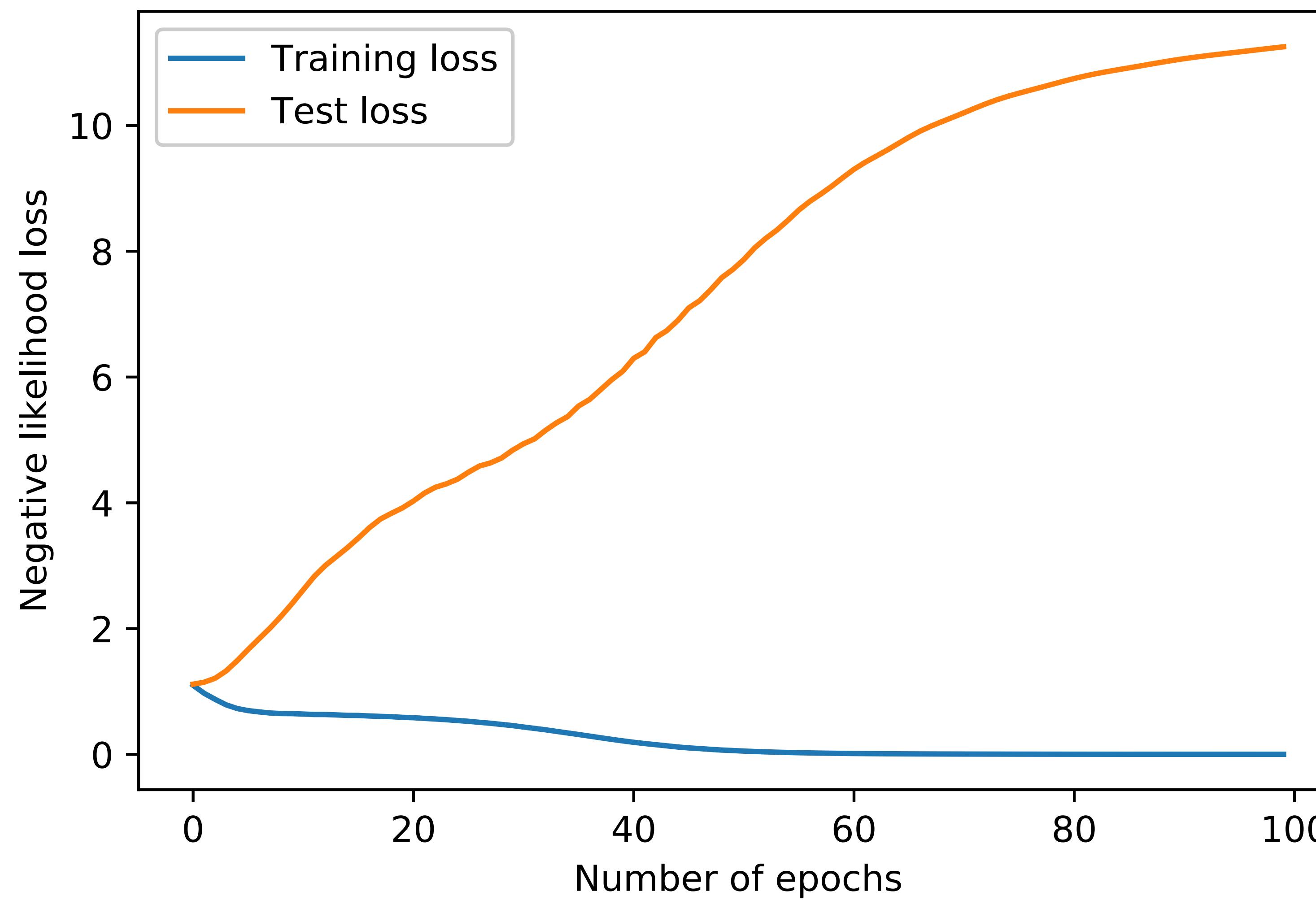


Model

Train with 3 examples

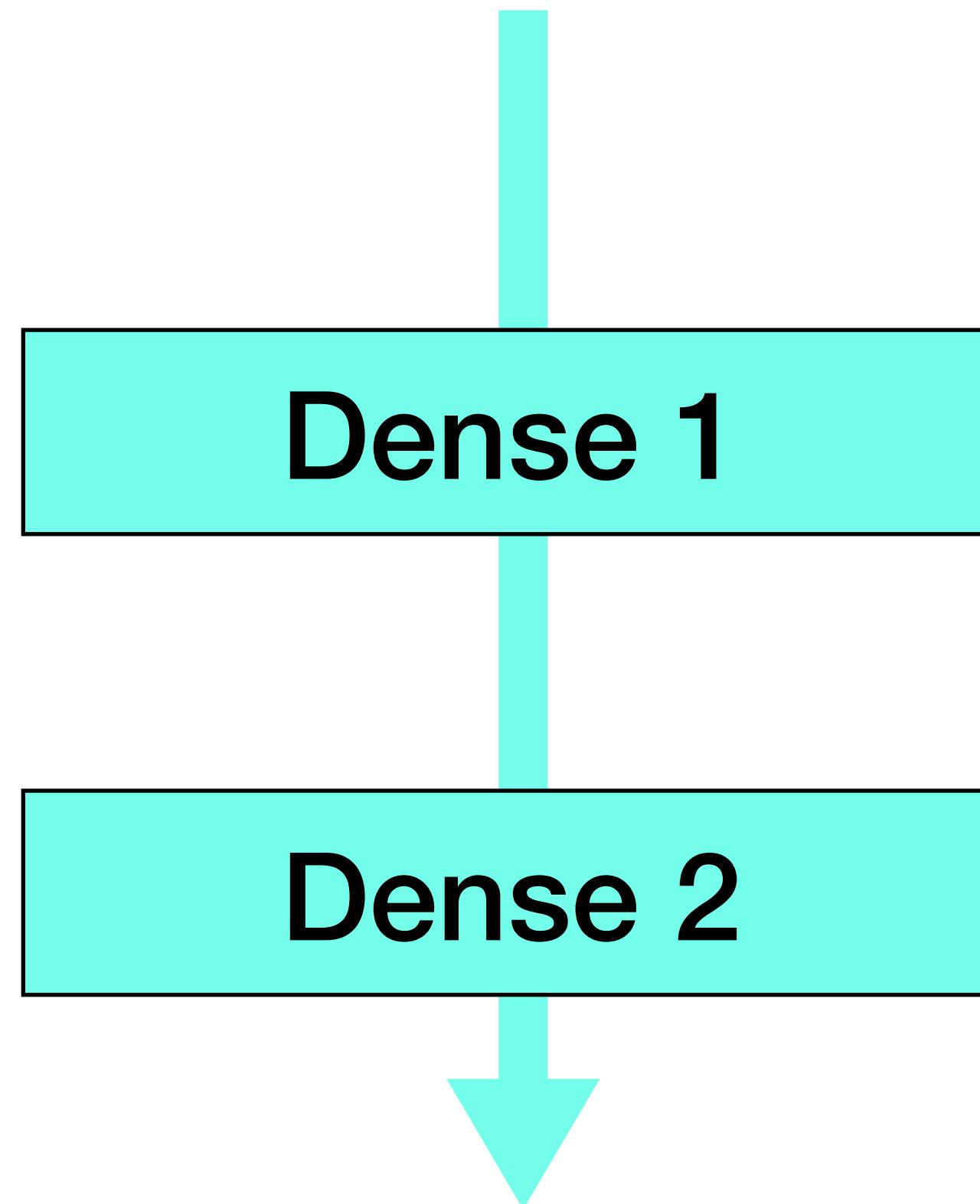


Train with 3 examples



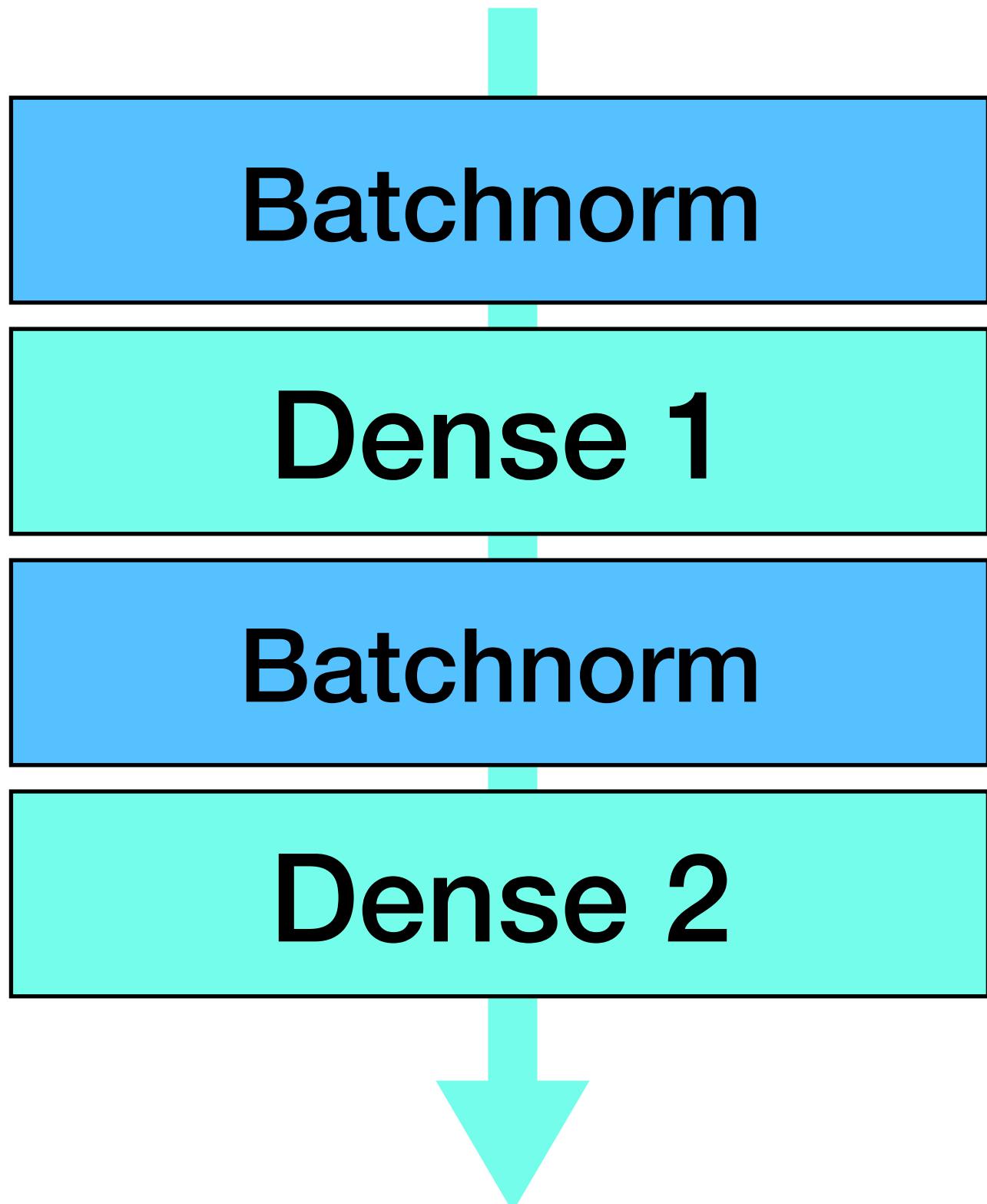
Batchnorm

Input



Output

Input

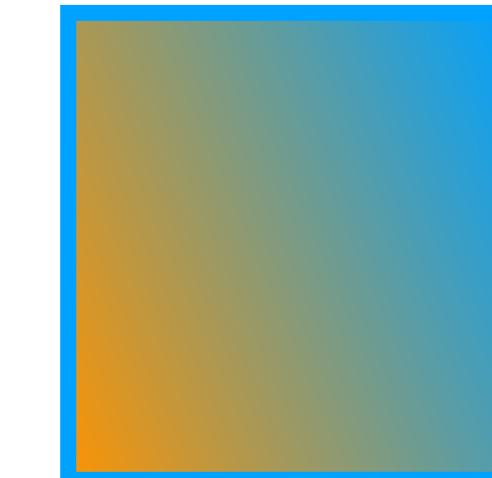
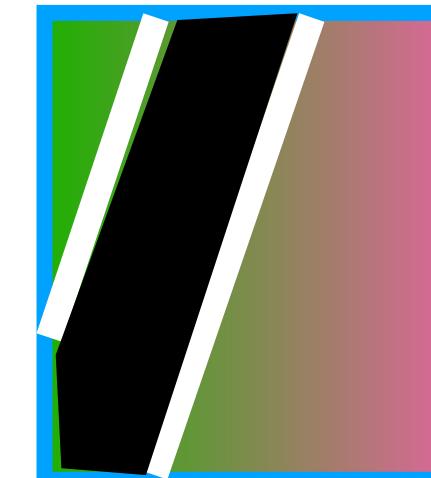
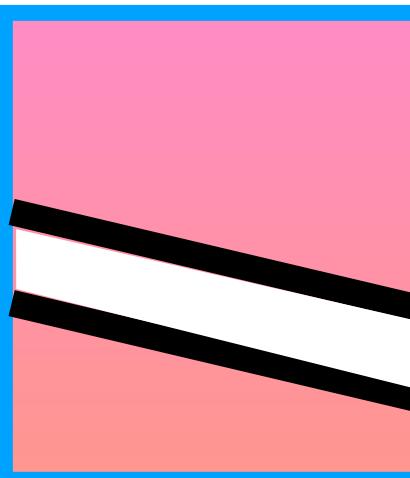


Output

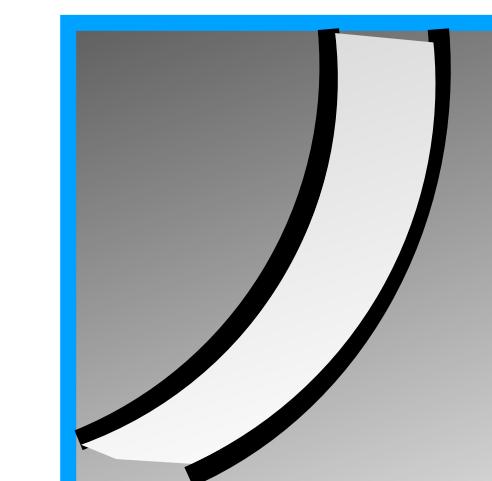
Transfer Learning



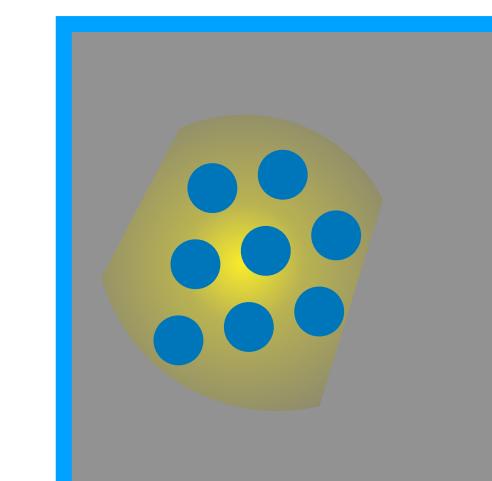
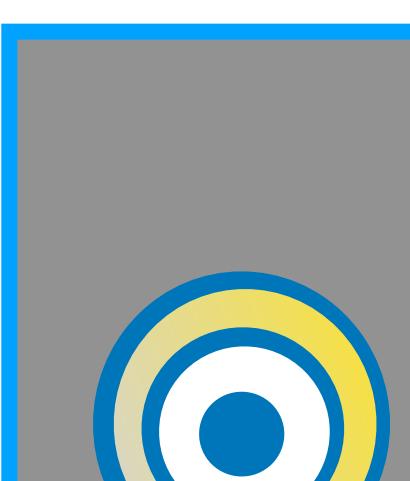
Low



Mid



High



Based on Zeiler and Fergus, 2013

**Trained
network**

Last layer(s)

e.g. ResNet

Mid/high level
features

Trained
network

Last layer(s)

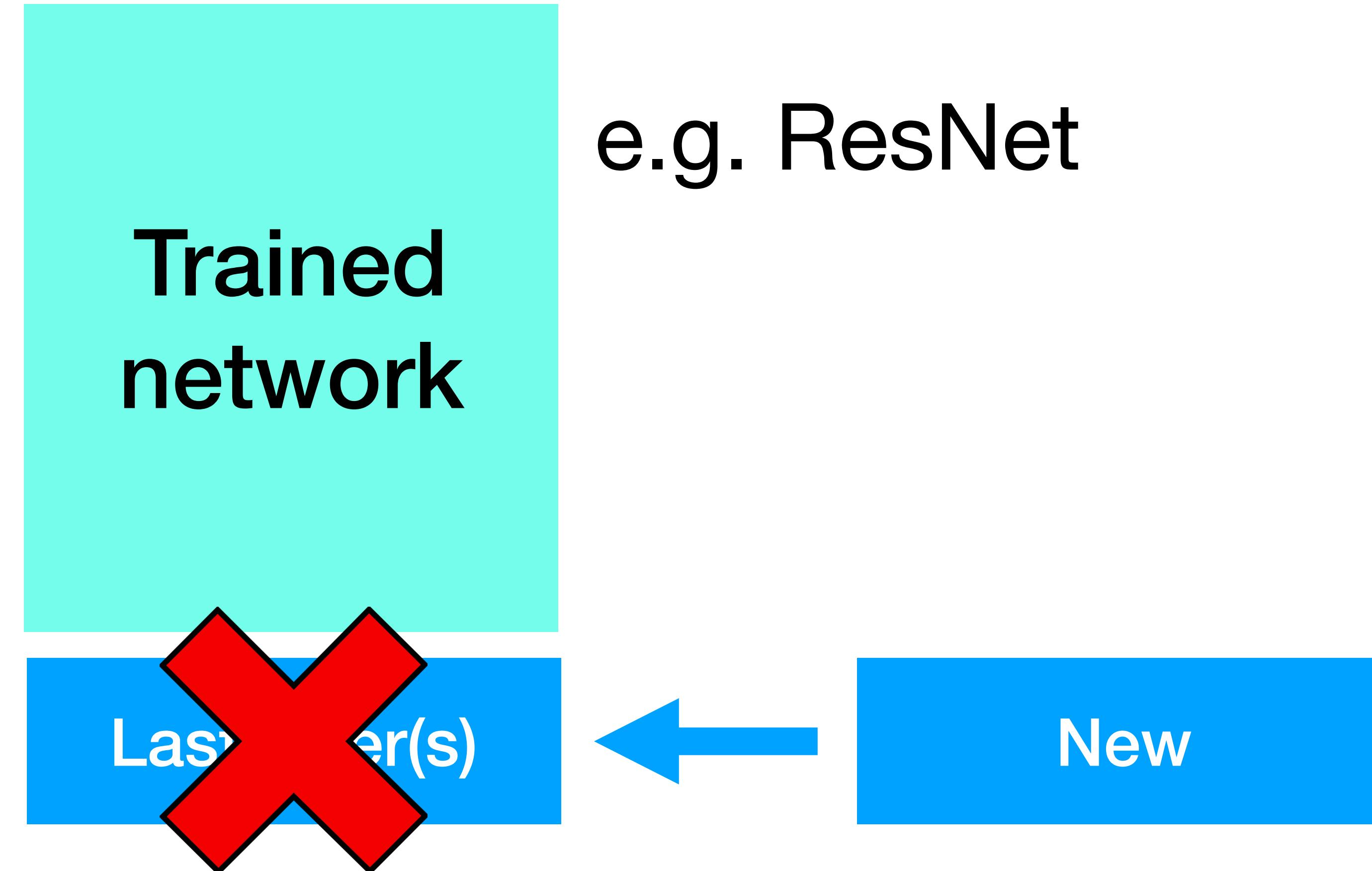
e.g. ResNet

Mid/high level
features

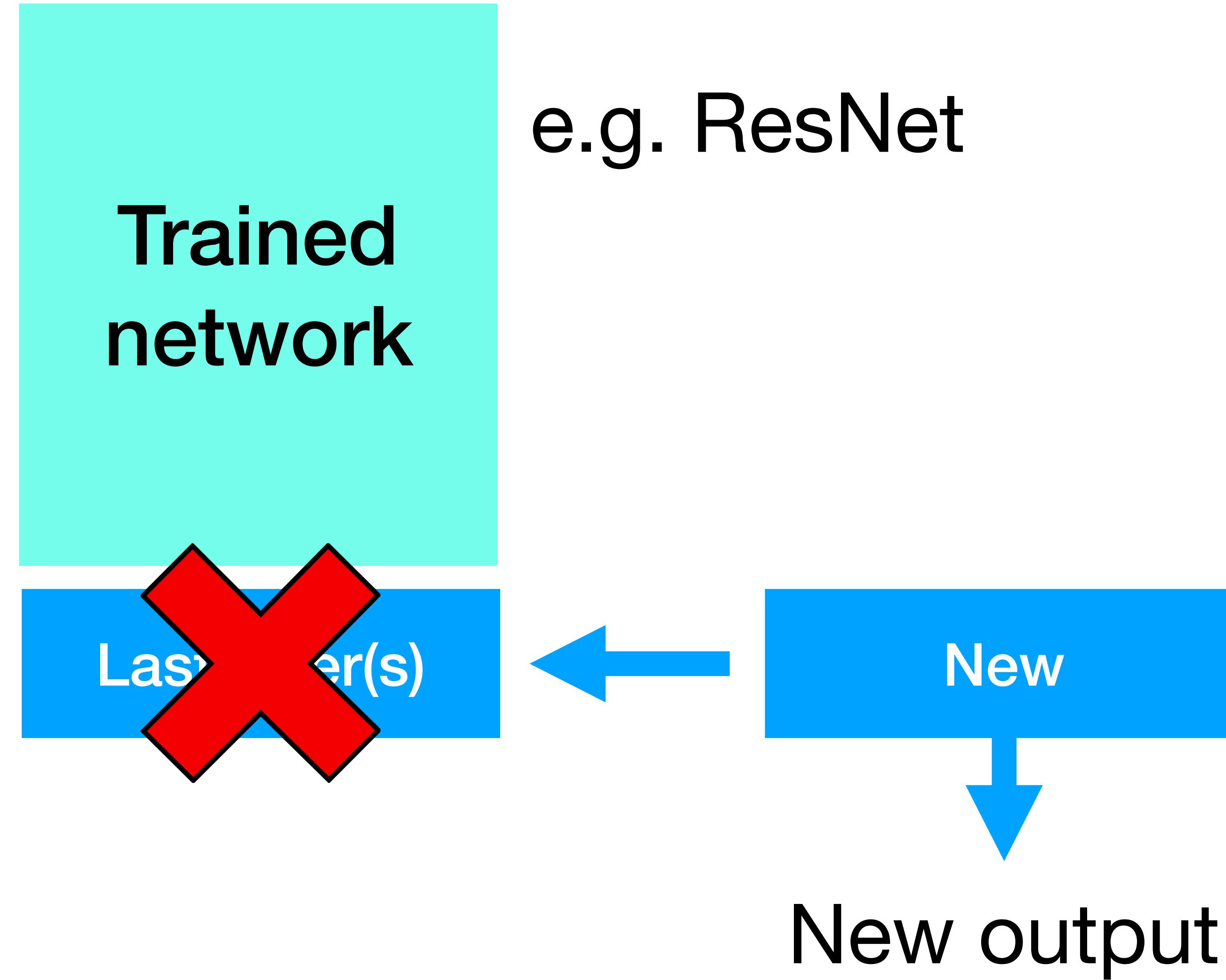


e.g. ResNet

Mid/high level
features

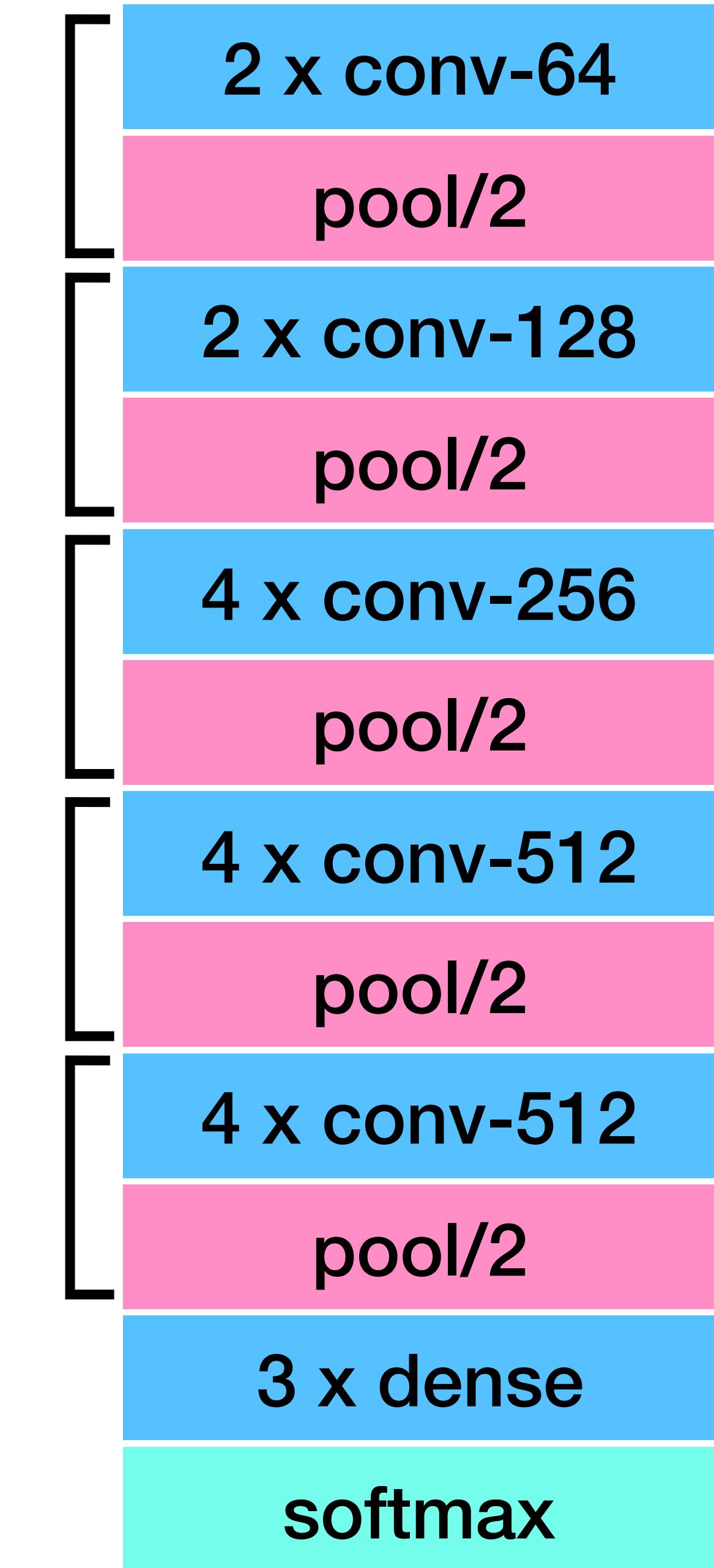


Mid/high level
features

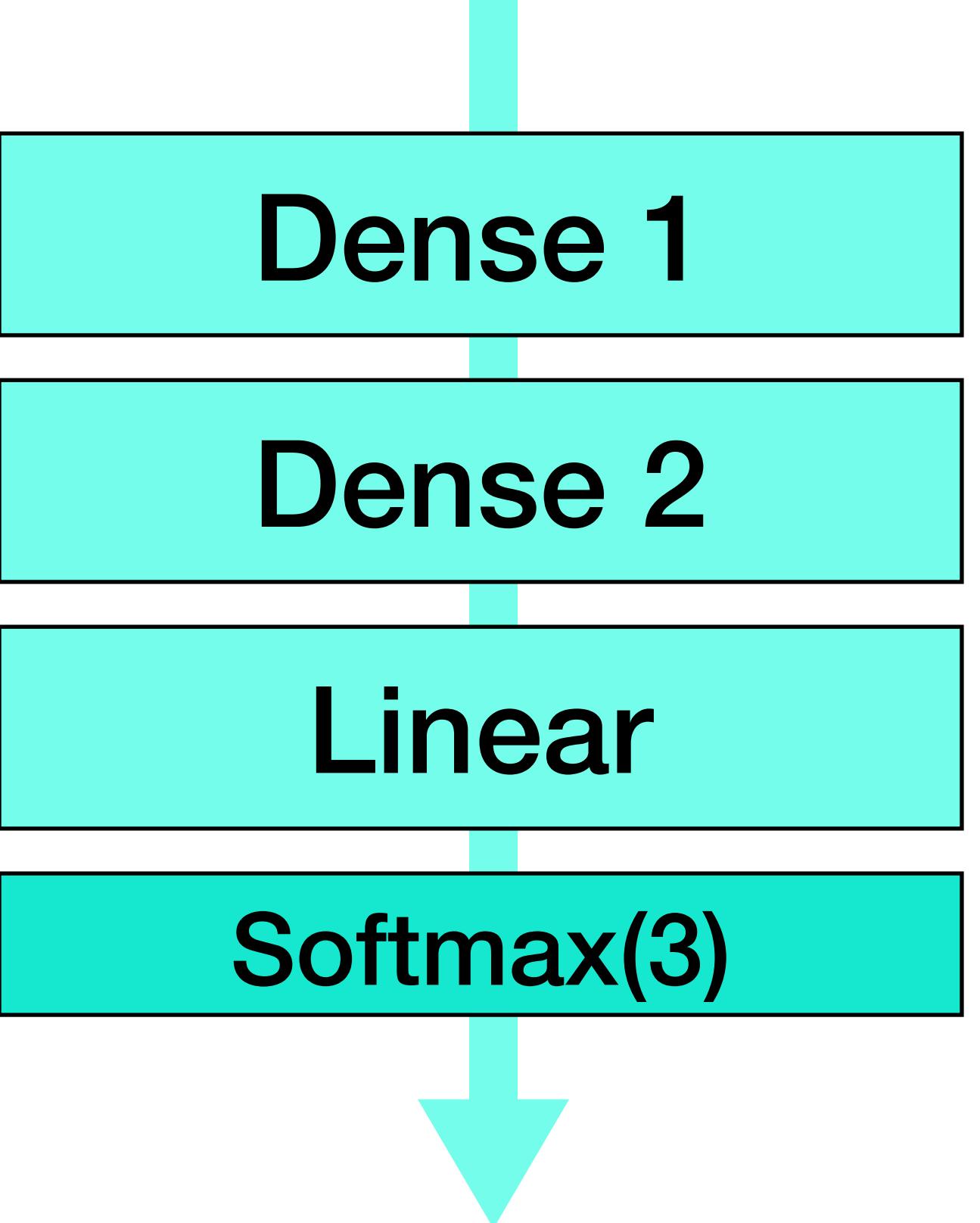


Simplify

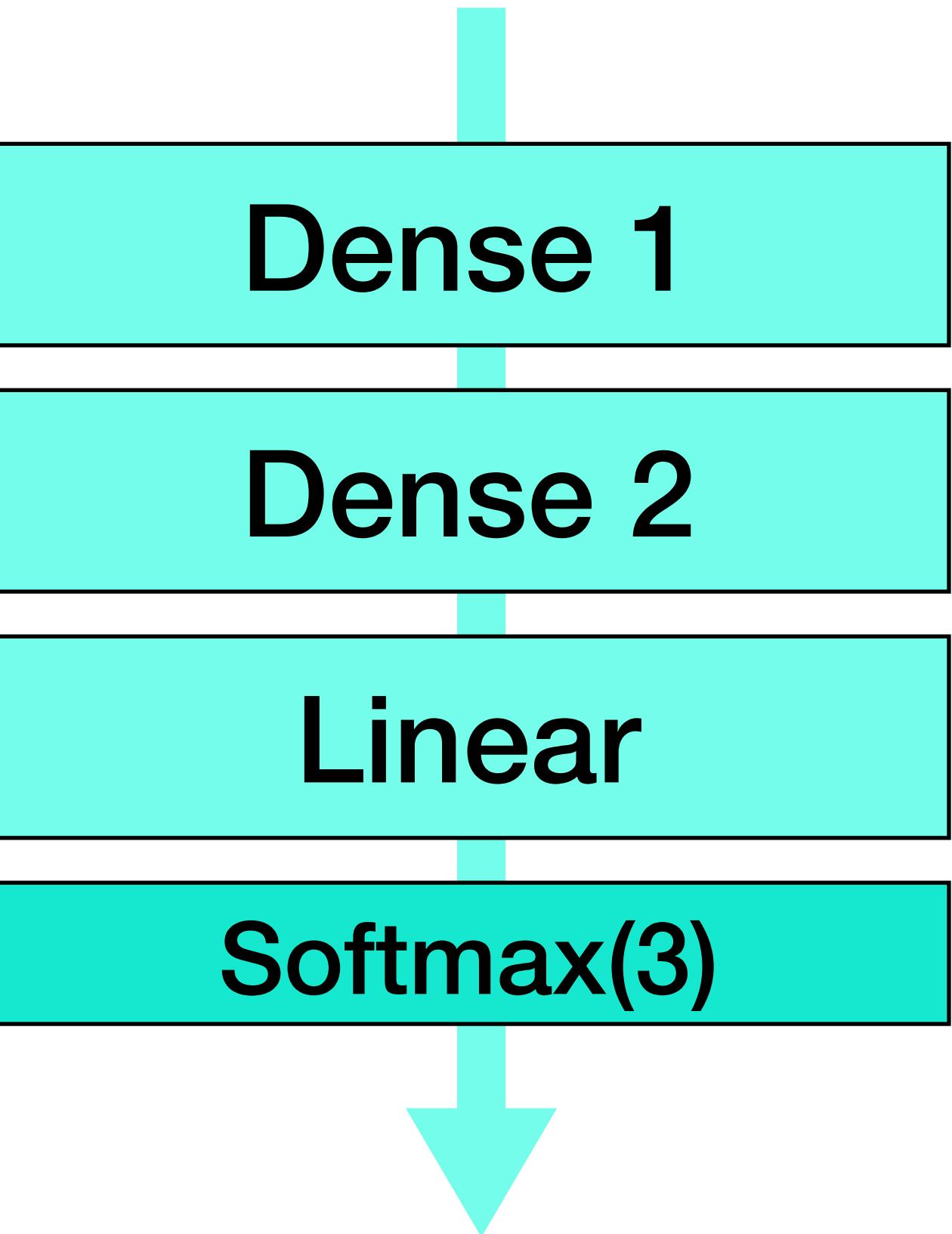
VGG-19



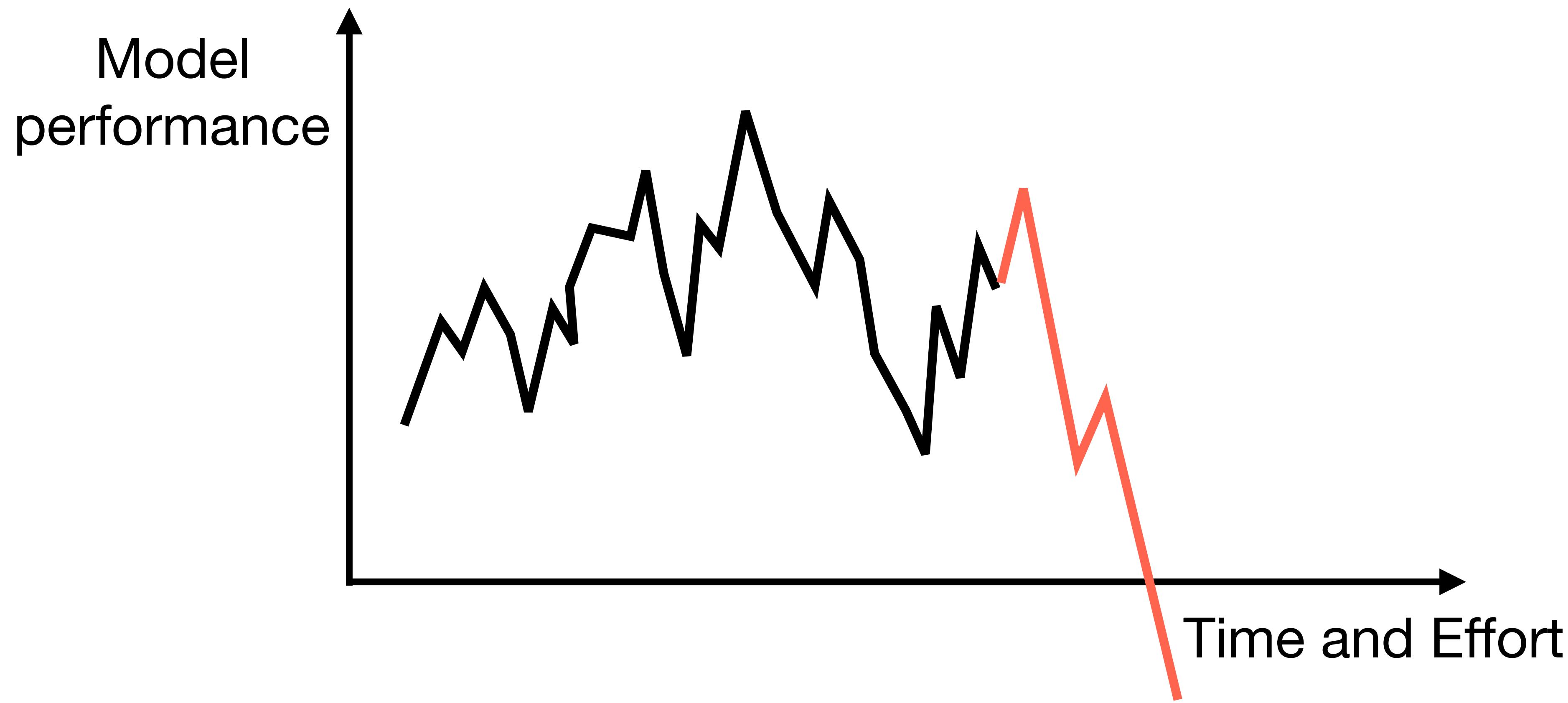
Dense 1

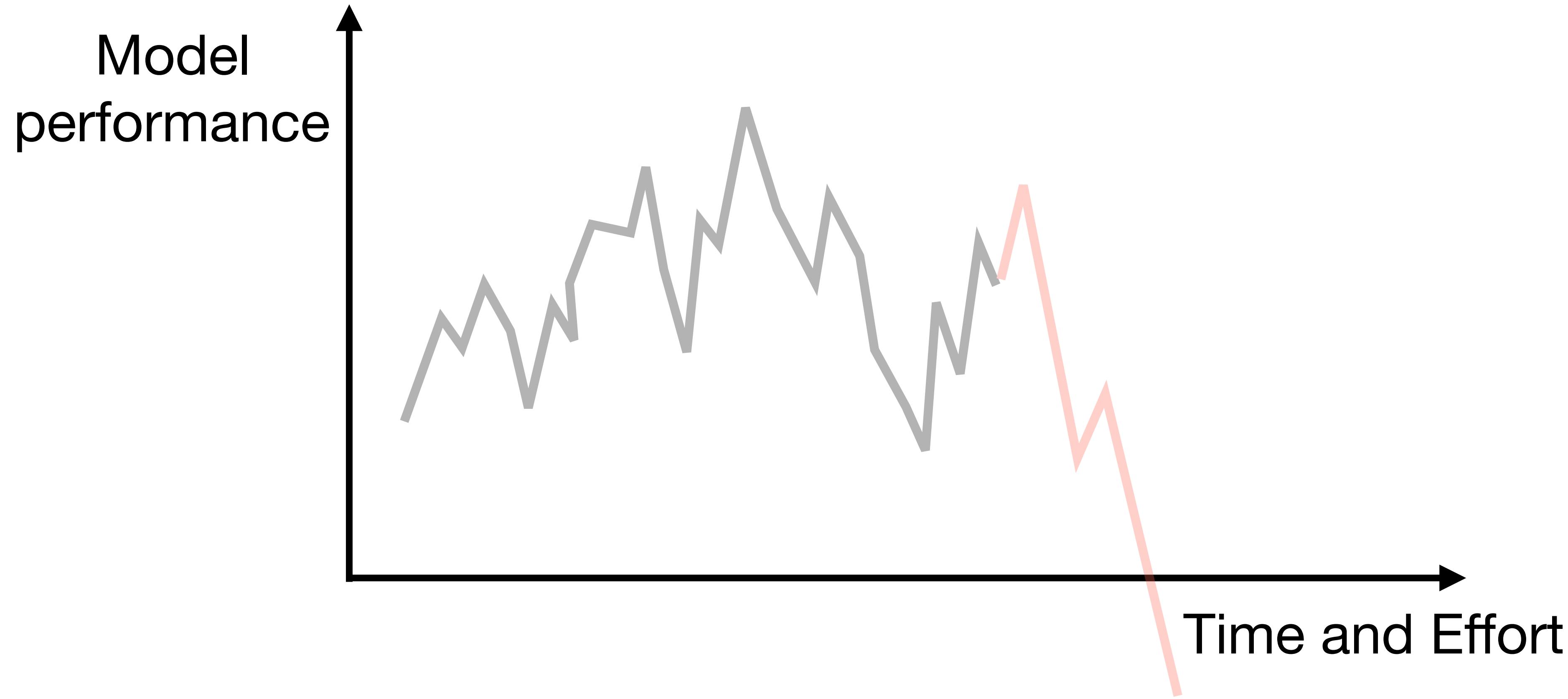


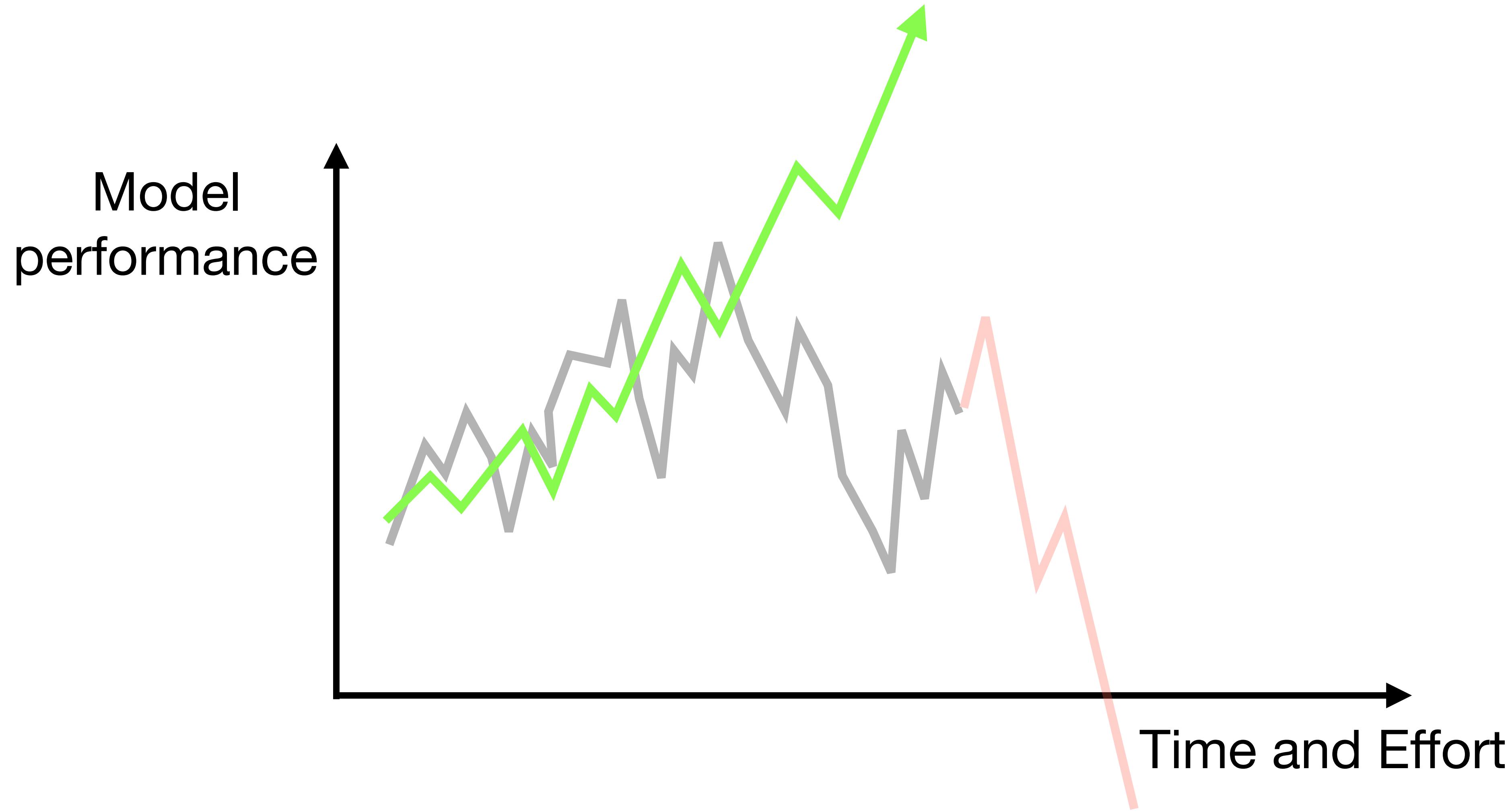
Input



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)