# 5 – Deep Neural Networks

**Mathematics of Data Science**

**Lecturer: Dominik Dold**

University of Vienna, WiSe 2025
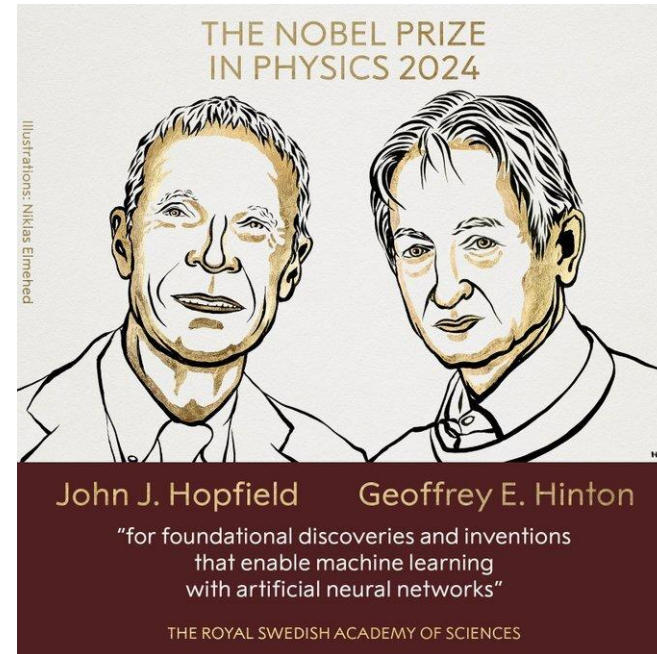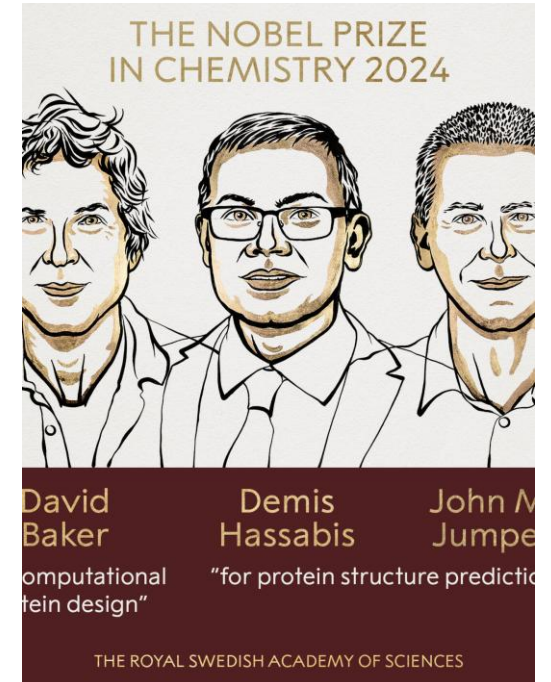
Master's programme in Data Science

# Motivation

**There is no way around deep artificial neural networks.**
They are one of the greatest achievements of our time!



*Neural networks beating up DotA 2 world champions (2019, OG vs OpenAI 5)*



*Neural network-related Nobel prizes*

**BUT:** In practice (and with sparse data), simpler models might often be the better choice!
E.g., trees, random forests, linear regression, hand-crafted features / models, etc.
So please **do not treat neural networks like a magic silver bullet**!

# Content

- Neural network definition

- Approximating Lipschitz-continuous functions

- Universal function approximation

- Affine / linear pieces

- Bounding the interpolation error

- Depth separation

- Error backpropagation and automatic differentiation

# What is a neural network?

We will now look at a very famous hypothesis set: **artificial neural networks (ANN)**

We define a single **layer** of a neural network as:

non-linearity      linear regression!

$$\Psi_{W,b}(\boldsymbol{x}) = \phi(\boldsymbol{Wx} + \boldsymbol{b})$$

where $\boldsymbol{W} \in \mathbb{R}^{n_2 \times n_1}, \boldsymbol{x} \in \mathbb{R}^{n_1}, \boldsymbol{b} \in \mathbb{R}^{n_1}$, and $\phi$ is a non-linear function (**activation function**).

We then define the hypothesis set of $L$-**layer neural networks** to consist of functions:

$$h(\boldsymbol{x}) = \Psi_{\boldsymbol{W}_L, \boldsymbol{b}_L} \circ \cdots \circ \Psi_{\boldsymbol{W}_1, \boldsymbol{b}_1}(\boldsymbol{x})$$

*Note: $\boldsymbol{W_i}$ are commonly called the weights, and $\boldsymbol{b_i}$ the biases.*

# Sketch of an artificial neural network



$x_1 \in \mathbb{R}^2$     $x_2 = \phi(W_1 x_1 + b_1)$     $x_3 = \phi(W_2 x_2 + b_2)$     $h = W_3 x_3 + b_3$

Input

Output

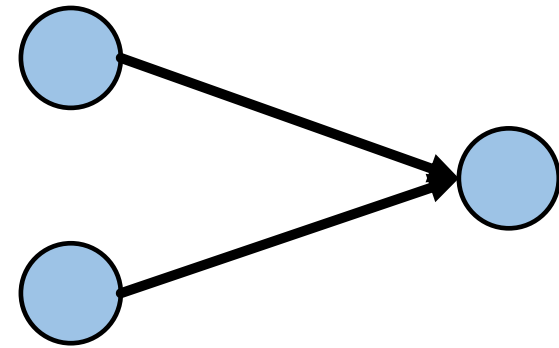$W_1 \in \mathbb{R}^{3 \times 2}$     $W_2 \in \mathbb{R}^{3 \times 3}$     $W_3 \in \mathbb{R}^{1 \times 3}$

# Question!

> **?** Why the non-linearity in each layer?

Assume we replace $\phi$ by the identity, i.e., $\phi(x) = x$. W.l.o.g., we set all biases $b_i = 0$.
Then we have:

$$h(x) = \Psi_{W_L,0} \circ \cdots \circ \Psi_{W_1,0}(x)$$

$$= W_L \cdot W_{L-1} \cdot \ldots \cdot W_1 x$$

$$= \widetilde{W} x$$

with $\widetilde{W} = W_L \cdot W_{L-1} \cdot \ldots \cdot W_1$

**That's just linear regression!**

# Warm-up: approximation properties of neural networks

So, what can we represent with neural networks?

Let's start with a specific case: **fitting Lipschitz continuous functions.**

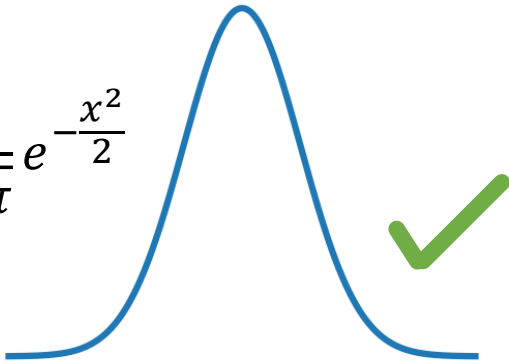A function $g: \mathbb{R} \to \mathbb{R}$ is $L$-Lipschitz if for all $x, y \in \mathbb{R}$, we have

$$|f(x) - f(y)| \leq L |x - y|$$

*Note: if the activation function is Lipschitz, then the neural network is too!*

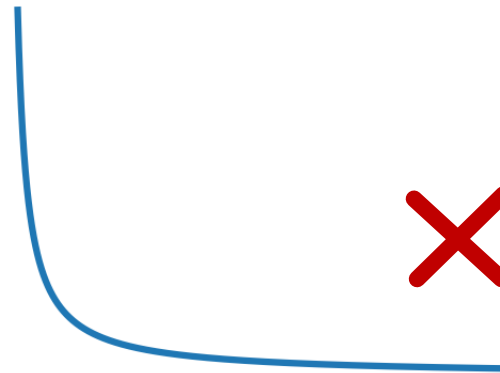**Which ones are Lipschitz continuous?**

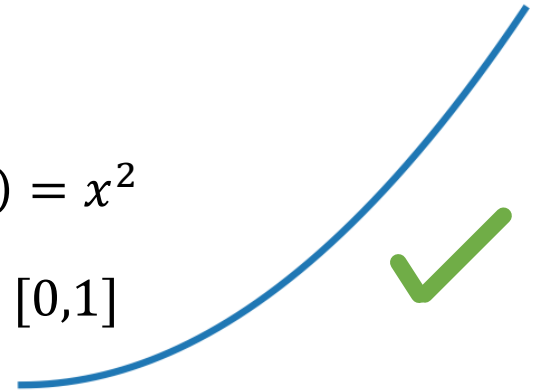$f(x) = \dfrac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$

$x \in \mathbb{R}$

$f(x) = \dfrac{1}{x}$

$x > 0$

$f(x) = x^2$

$x \in [0,1]$

# Representing Lipschitz-continuous functions

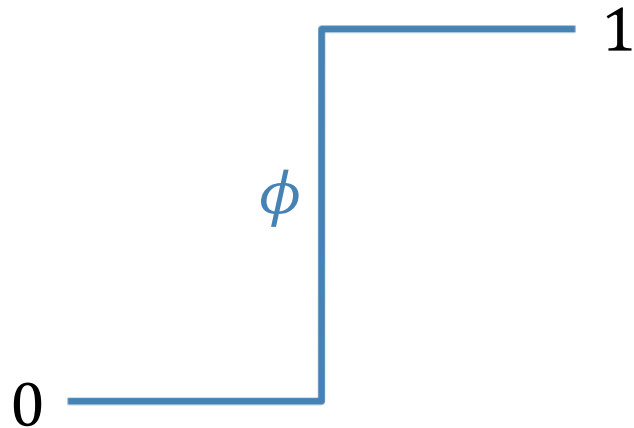Let $g: [0,1] \to \mathbb{R}$ be a $L$-Lipschitz function and $\epsilon > 0$. Then a 2-layer neural network $f$ with threshold activation function $\phi$ and $\mathcal{O}\left(\frac{L}{\epsilon}\right)$ neurons exists that satisfies:

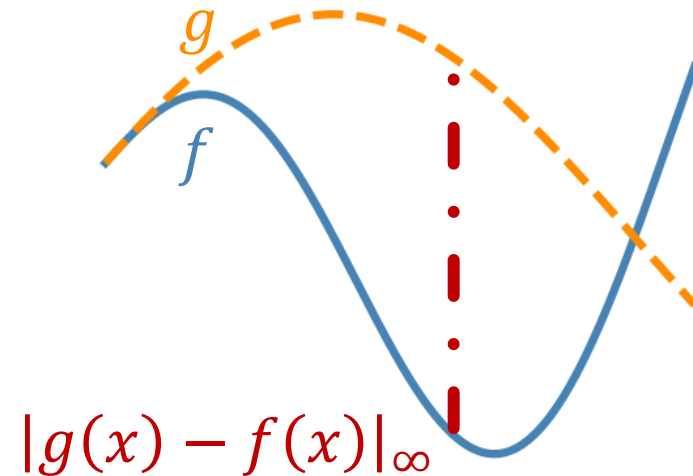$$|g(x) - f(x)|_\infty < \epsilon$$

*Proof: see handwritten notes*

**Threshold function:**

$$\phi(x) = \begin{cases} 1, & x > 0 \\ 0, & x \leq 0 \end{cases}$$



**Supremum norm:**

$$|g(x) - f(x)|_\infty = \sup_x |g(x) - f(x)|$$

# The same holds for ReLU activation functions!

Let $g: [0,1] \to \mathbb{R}$ be a $L$-Lipschitz function and $\epsilon > 0$. Then a 2-layer neural network $f$ with ReLU activation function $\phi$ and $\mathcal{O}\left(\frac{L}{\epsilon}\right)$ neurons exists that satisfies:

$$|g(x) - f(x)|_\infty < 2\epsilon$$

*Proof: see handwritten notes*

**ReLU (rectified linear unit) function:**

$$\phi(x) = \max(0, x)$$



ReLU is one of the **most** used activation functions in deep learning!
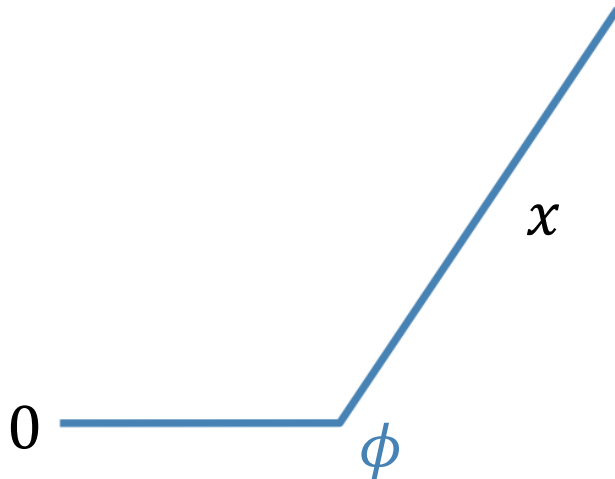
# Final warm-up: going multivariate!

Let $g: [0,1]^d \to \mathbb{R}$ be a $L$-Lipschitz function and $\epsilon > 0$. Then a 3-layer neural network $f$ with threshold activation function $\phi$ and $\mathcal{O}\left(\frac{L}{\epsilon^d}\right)$ neurons exists that satisfies:

$$|g(x) - f(x)|_\infty < \epsilon$$

Or in case of ReLU activation function:

$$|g(x) - f(x)|_\infty < 2\epsilon$$

*Proof sketch: see handwritten notes*

**Note:** We need one more layer compared to the univariate case. The first layer basically does the same as in the univariate case. The second layer takes the output of the first to identify in which hypercube we are! The third layer constructs the function from these hypercubes (as in 1D).

# Are neural networks universal function approximators?

Our previous results are limited to Lipschitz continuous functions and special neural network architectures. Can we do better?

**Yes!** Using the Stone-Weierstrass theorem!

Let $\mathcal{H}$ be the hypothesis set consisting of **shallow** neural networks, i.e., neural networks with one **hidden** layer, given by:

$$\mathcal{H} = \text{span}_{\alpha,w,b}\{f_{\alpha,w,b}\}$$

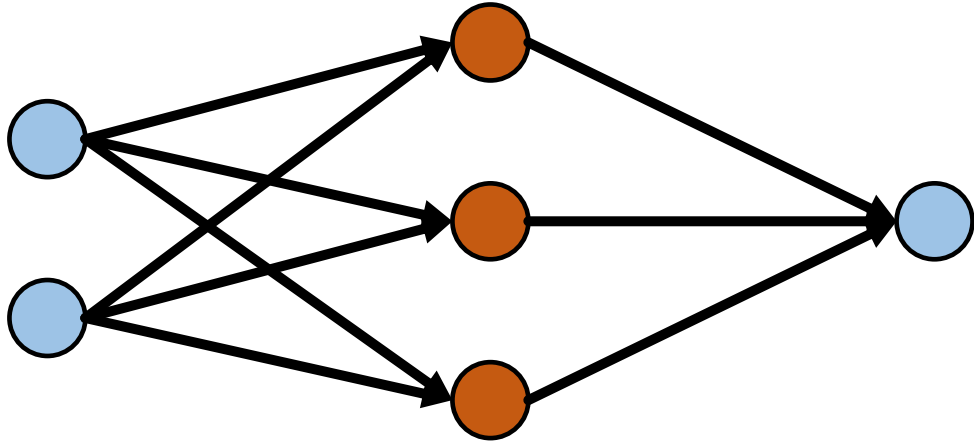with single neurons $f_{\alpha,w,b}: x \mapsto \alpha\, \phi(\boldsymbol{w}^T \boldsymbol{x} + b)$.

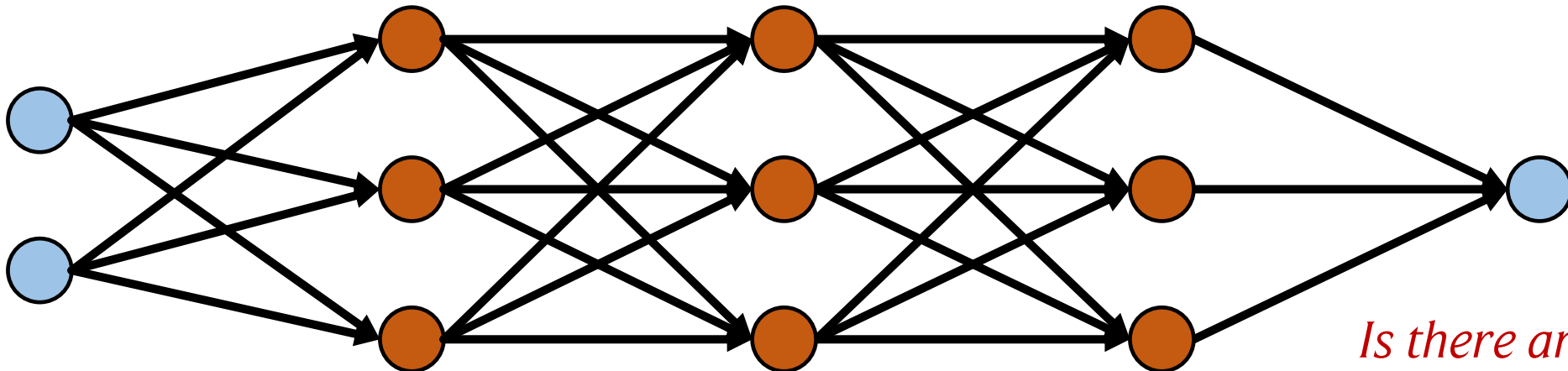e.g., $h(x) = \sum_i^3 \alpha_i\, \phi(\boldsymbol{w}_i^T \boldsymbol{x} + b_i) \in \mathcal{H}$

**Exercise:** Using Stone-Weierstrass, show that for $\phi(x) = \cos(x)$ and $\phi(x) = e^x$, $\mathcal{H}$ is a universal approximator. If $\phi$ is a polynomial, is $\mathcal{H}$ still a universal approximator? What about sigmoidal functions ($\phi$ non-decreasing, $\lim_{x \to \infty} \phi(x) = 1$, $\lim_{x \to -\infty} \phi(x) = 0$)?

K. Hornik, M. Stinchcombe, H. White, *Multilayer feedforward networks are universal approximators*, 1989

# If shallow ANNs are universal, what about depth?

**Reminder:** shallow = 1 hidden layer (what we dealt with so far)



But we can stack more hidden layers to go deep (usually done in practice).
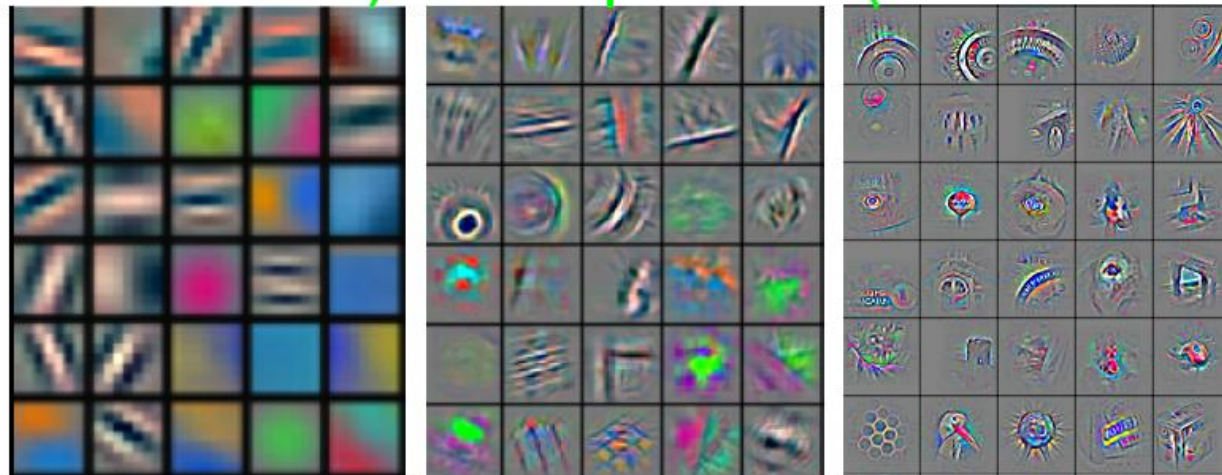


*Is there any benefit to depth?*

# Intuitive picture: hierarchy of features



Deep Learning = Learning Hierarchical Representations

Y LeCun

It's deep if it has more than one stage of non-linear feature transformation

Low-Level Feature → Mid-Level Feature → High-Level Feature → Trainable Classifier

Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]
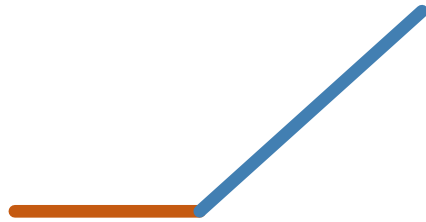
early layers

later layers

We utilized this in our constructions earlier! E.g., the 3-layer network: **simple features from the first layer were used to construct more complex ones in the second one**.

But that's just an intuitive picture... can we show clear benefits of depth?

# A new tool: affine pieces

An affine / linear piece is a linear segment of a piecewise continuous function.

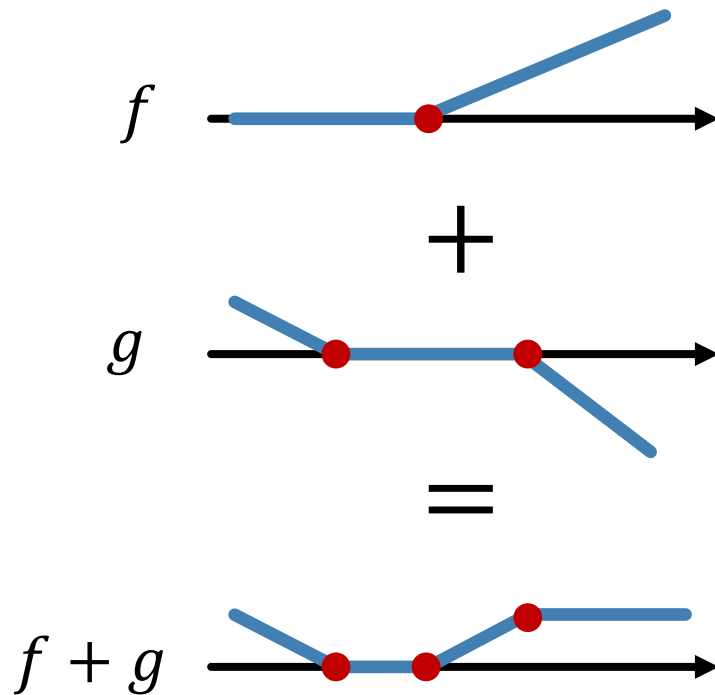**Examples:**

ReLU (2 pieces)

A function with 5 pieces.

Generally, we are interested in the **number of pieces** a function has. We will show later: the more pieces, the better we can fit functions!

# Addition and composition of piecewise functions

Adding and composing piecewise continuous functions results again in piecewise continuous functions with an **increased** (or unchanged) **number of pieces**.

**Addition:**



If $f$ has $p$ and $g$ has $q$ pieces, then $f + g$ **has at most** $p + q$ pieces.
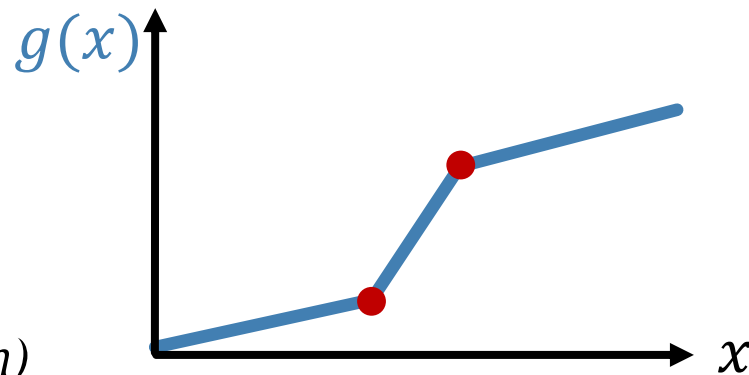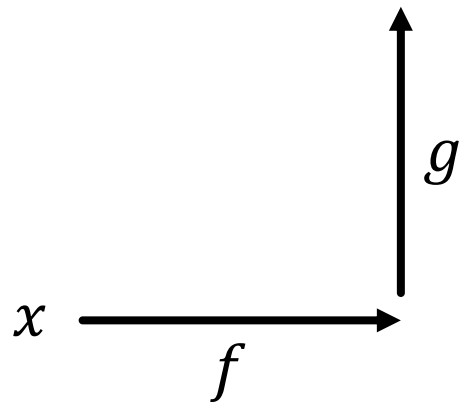
# Addition and composition of piecewise functions

Adding and composing piecewise continuous functions results again in piecewise continuous functions with an **increased** (or unchanged) **number of pieces**.
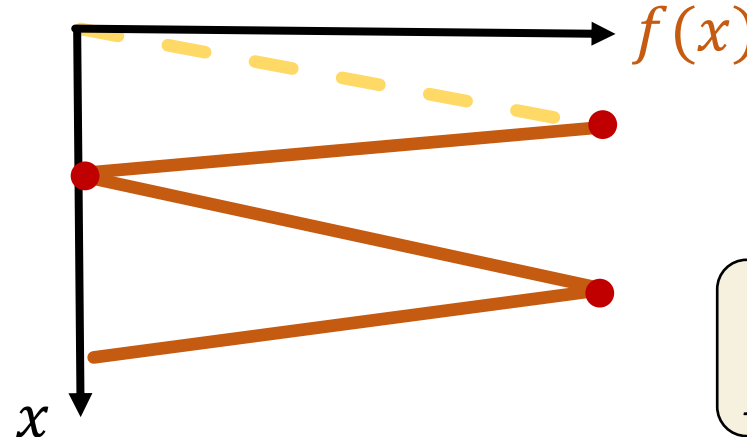
**Composition:** $g \circ f(x)$

$g(x)$

*How to read this:*
*(the output of $g \circ f$ is not shown)*

$g$

$x$ $\xrightarrow{\quad f \quad}$

$x$

$f(x)$

$x$

Consider a single piece of $f$ (yellow dashed line).

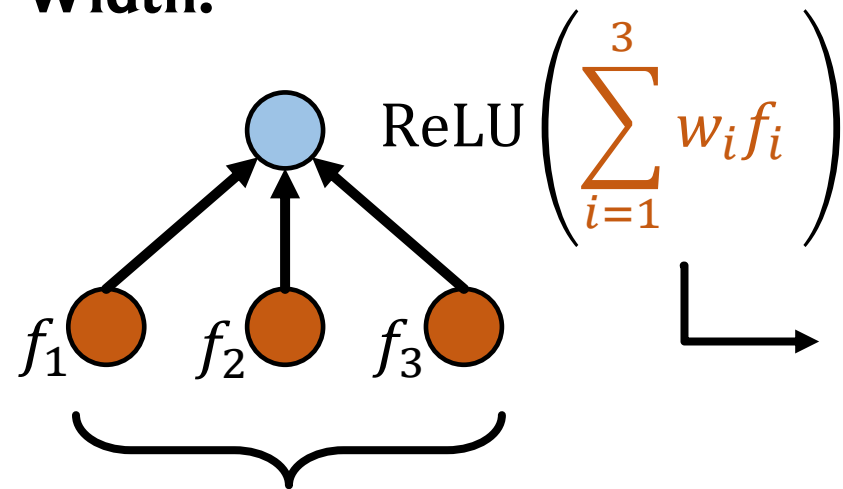On this piece, whenever $f$ crosses a kink of $g$, the piece of $f$ gets "split" (i.e., a new piece is added).

If $f$ has $p$ and $g$ has $q$ pieces, then $f \circ g$ **has at most** $p \cdot q$ pieces.

# ReLU neural networks are piecewise continuous functions

ReLU neural networks are composed of ReLU functions, which are piecewise continuous.

Hence, **ReLU neural networks are piecewise continuous functions**!

How many pieces does such a network have?

**Width:**



$$\text{ReLU}\left(\sum_{i=1}^{3} w_i f_i\right)$$

$f_1 \quad f_2 \quad f_3$

*width = # neurons in a layer*

Addition!
$\rightarrow 2 \cdot \text{width}$ pieces
*(2 pieces per neuron)*

**Depth:**

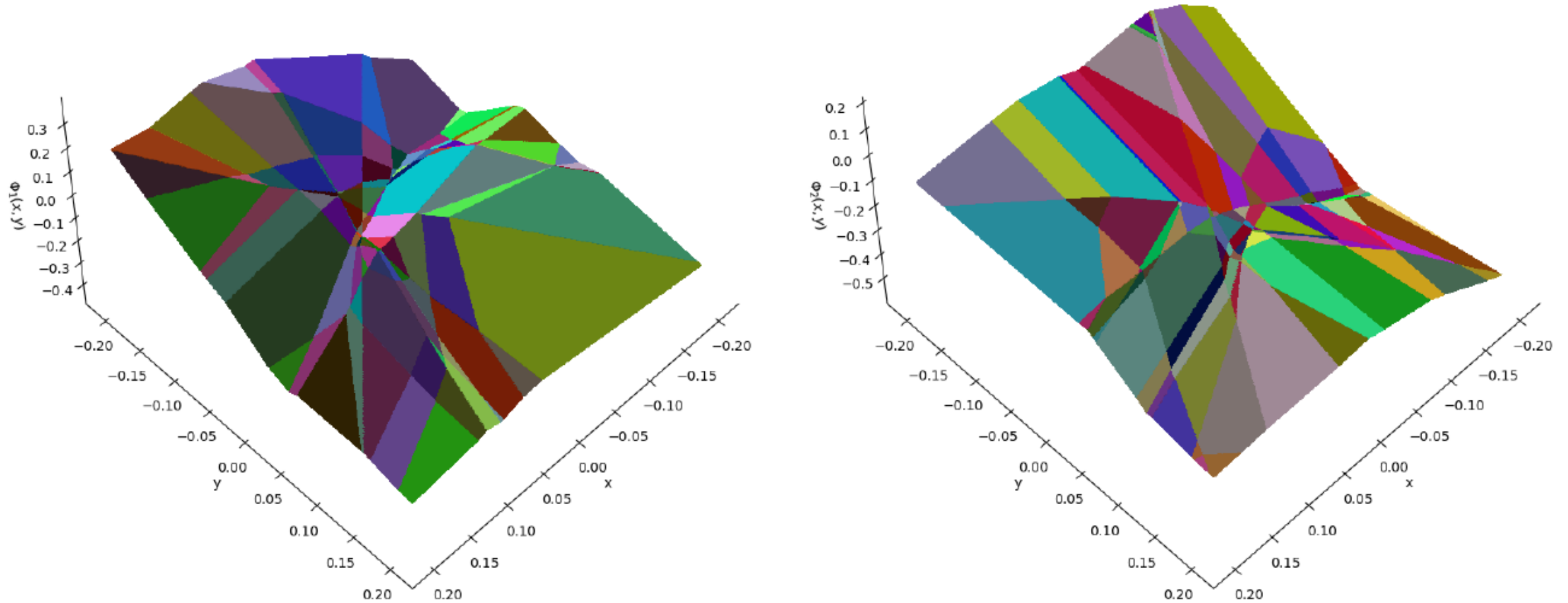$$\text{ReLU}\left(\sum_{i=1}^{3} w_i f_i\right)$$

Composition!
$\rightarrow 2 \cdot (2 \cdot \text{width})$ pieces
*(ReLU has 2 pieces)*

A **ReLU neural network** with **_L_ layers** and given **width** has at most $(\mathbf{2 \cdot width})^{\mathbf{L}}$ affine pieces.

# Example



Affine pieces of two different ReLU neural networks with two inputs and 1 output.

# Why does this matter? Bounding the interpolation error!

Through the number of affine pieces, we get a lower bound on the interpolation error:

Let $f: [0,1] \to \mathbb{R}$ be a continuous and three-times differentiable function. Then there exists a constant $C$ such that for all ReLU neural networks $\Psi: \mathbb{R} \to \mathbb{R}$ with $p$ affine pieces, we have:

$$|f - \Psi|_\infty \geq \frac{C}{p^2}$$

*Note: this is also true for multivariate input along lines. The statement here is also slightly simplified.*
*Proof: see handwritten notes.*

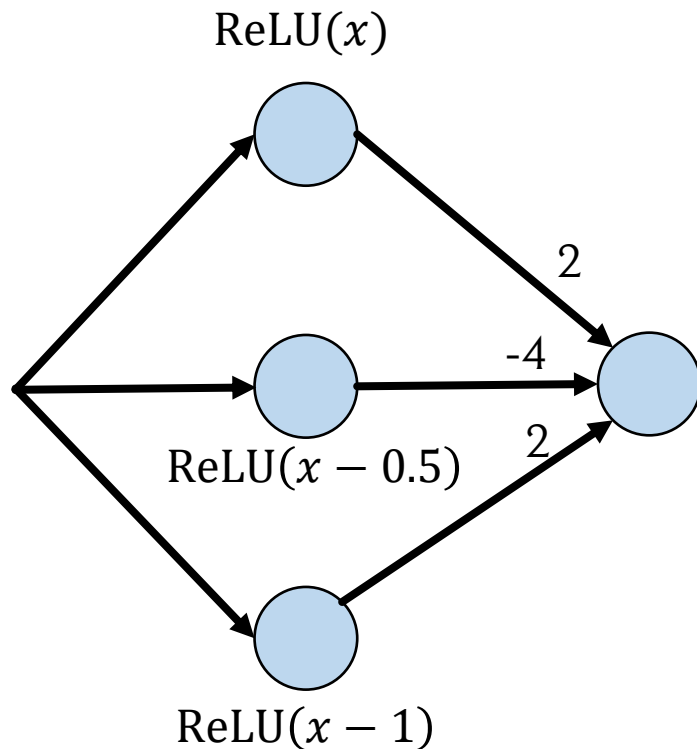Specifically, with our previous result, we get:   $|f - \Psi|_\infty \geq C \cdot (2 \cdot \text{width})^{-2L}$

This clearly indicates a strong benefit of depth, as only with $L$, the bound decays exponentially!

**BUT:** the multiplicative behavior was only a best-case scenario, does this hold in practice?!

C. L. Frenzen, T. Sasao, and J. T. Butler. "On the number of segments needed in a piecewise linear approximation" (2010)

# The hat function

Consider the following 2-layer ReLU neural network:

$$\Delta(x) = 2\left[\text{ReLU}(x) - 2 \cdot \text{ReLU}\left(x - \frac{1}{2}\right) + \text{ReLU}(x - 1)\right]$$

ReLU($x$)

ReLU($x - 0.5$)

ReLU($x - 1$)

2

-4

2

# The hat function

Consider the following 2-layer ReLU neural network:

$$\Delta(x) = \mathbf{2}\left[\mathbf{ReLU}(\boldsymbol{x}) - 2 \cdot \text{ReLU}\left(x - \frac{1}{2}\right) + \text{ReLU}(x - 1)\right]$$

# The hat function
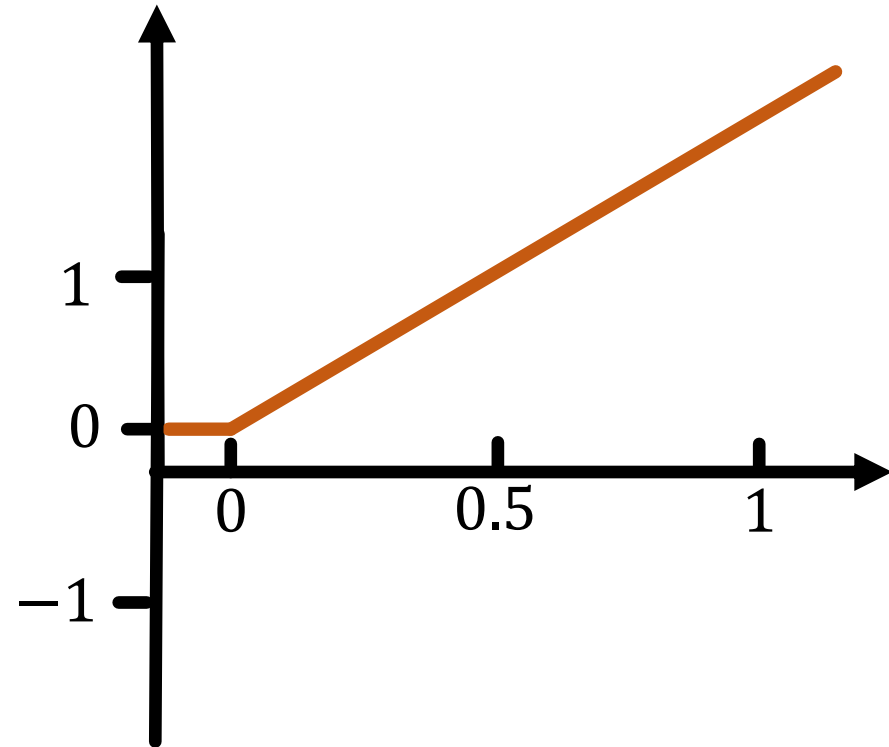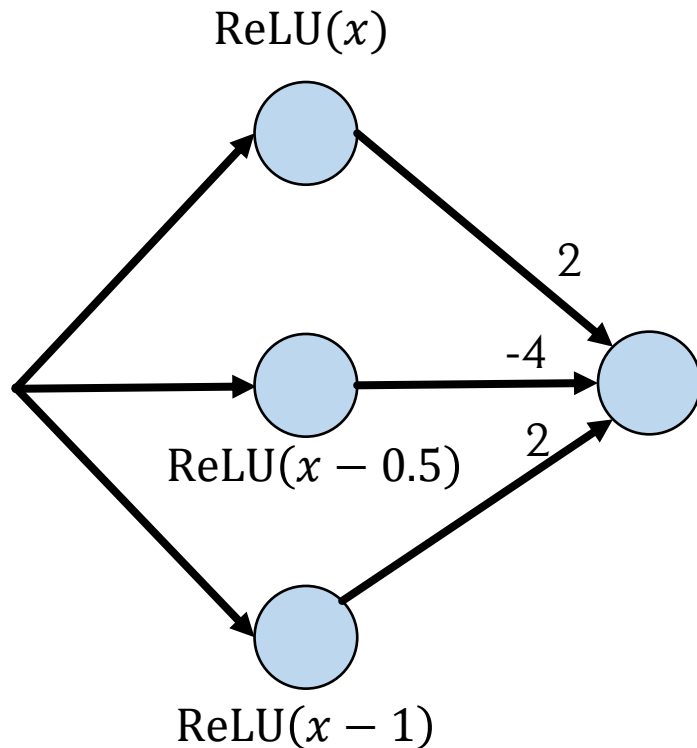
Consider the following 2-layer ReLU neural network:

$$\Delta(x) = 2\left[\text{ReLU}(x) - 2 \cdot \text{ReLU}\left(x - \frac{1}{2}\right) + \text{ReLU}(x - 1)\right]$$

# The hat function

Consider the following 2-layer ReLU neural network:
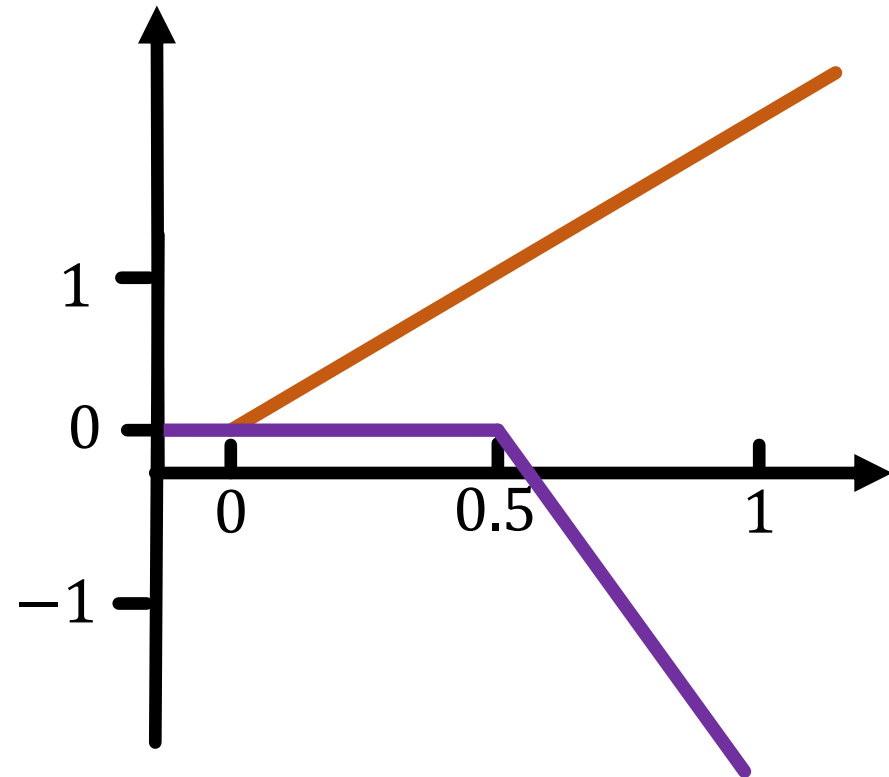
$$\Delta(x) = 2\left[\text{ReLU}(x) - 2 \cdot \text{ReLU}\left(x - \frac{1}{2}\right) + \text{ReLU}(x-1)\right]$$

# The hat function

Consider the following 2-layer ReLU neural network:

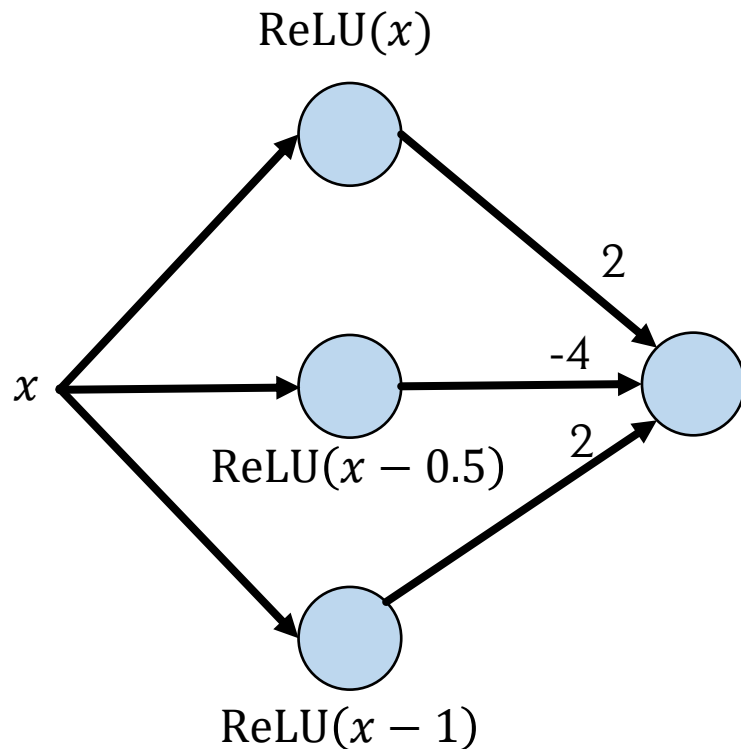$$\Delta(x) = 2 \left[ \text{ReLU}(x) - 2 \cdot \text{ReLU}\left(x - \frac{1}{2}\right) + \textbf{ReLU}(\textbf{x} - \textbf{1}) \right]$$

# The hat function

Consider the following 2-layer ReLU neural network:

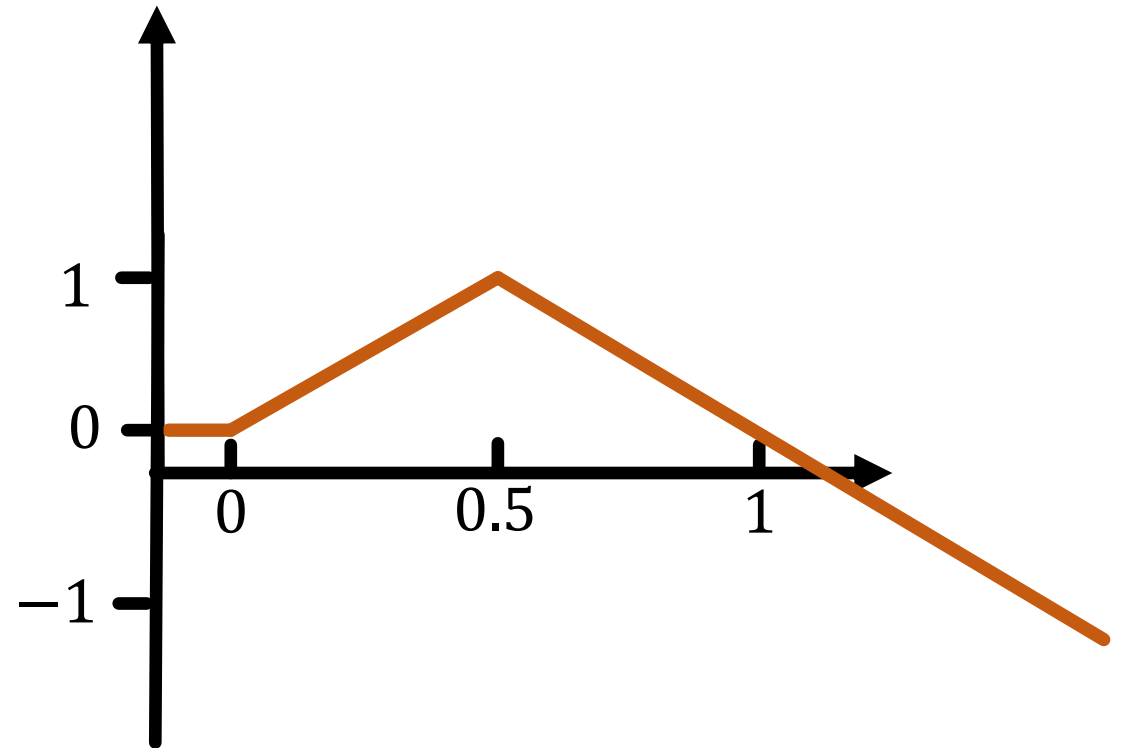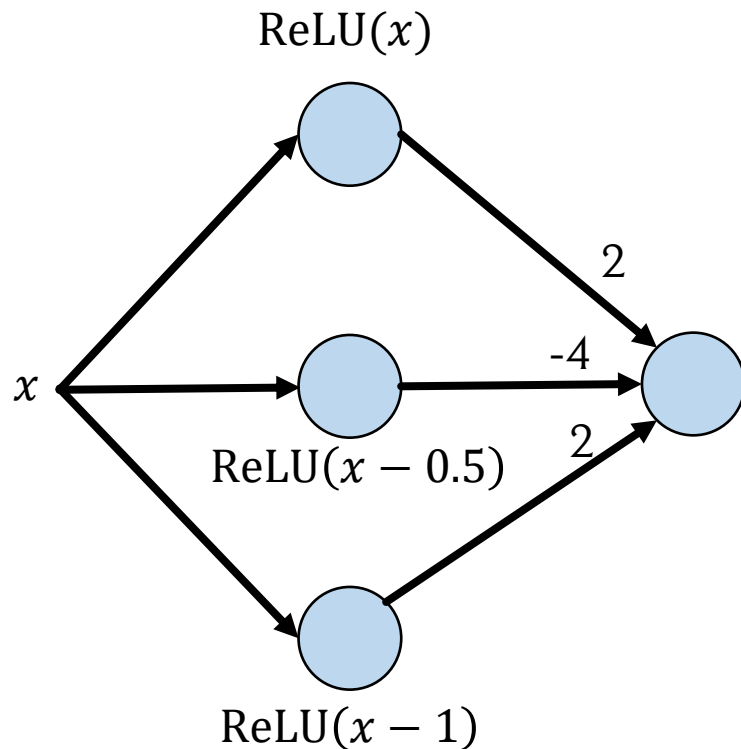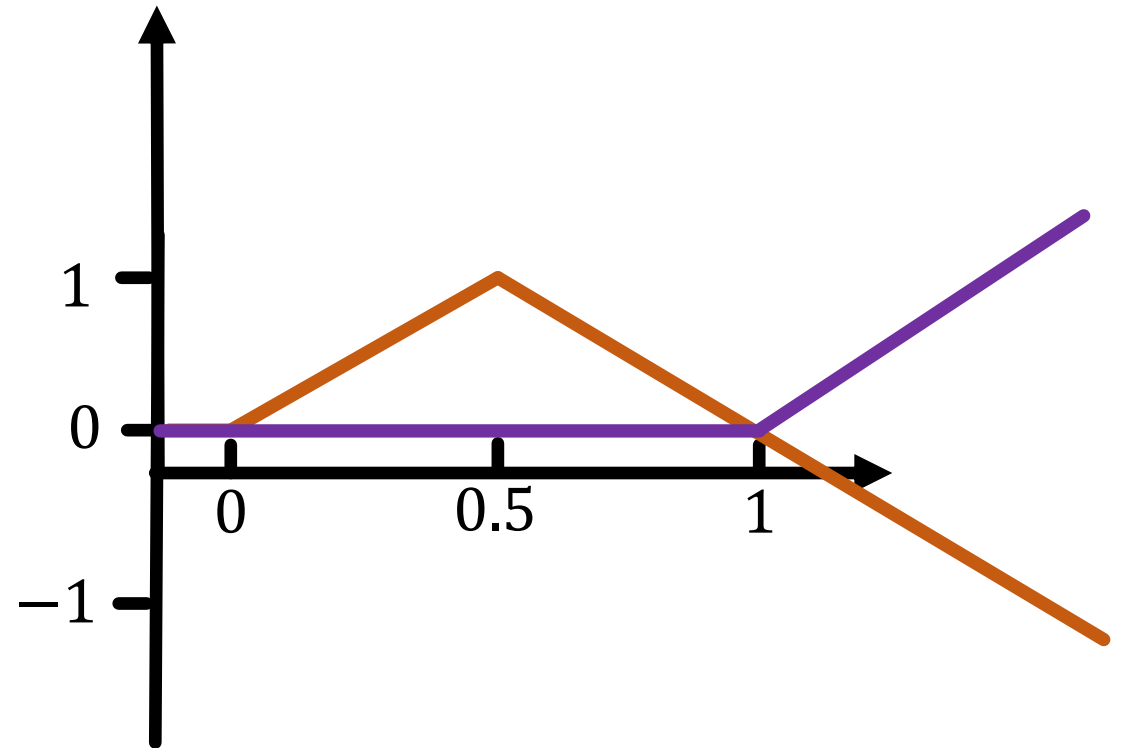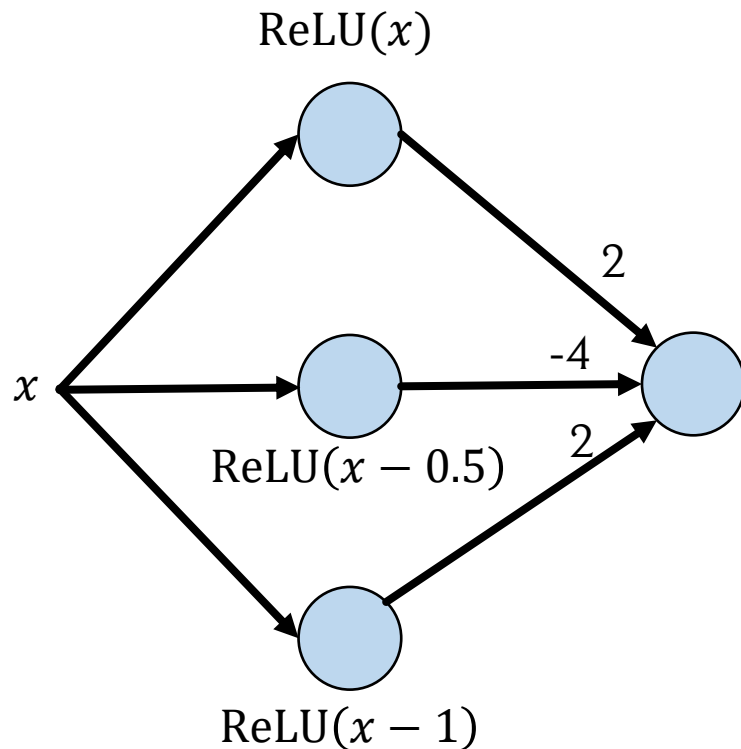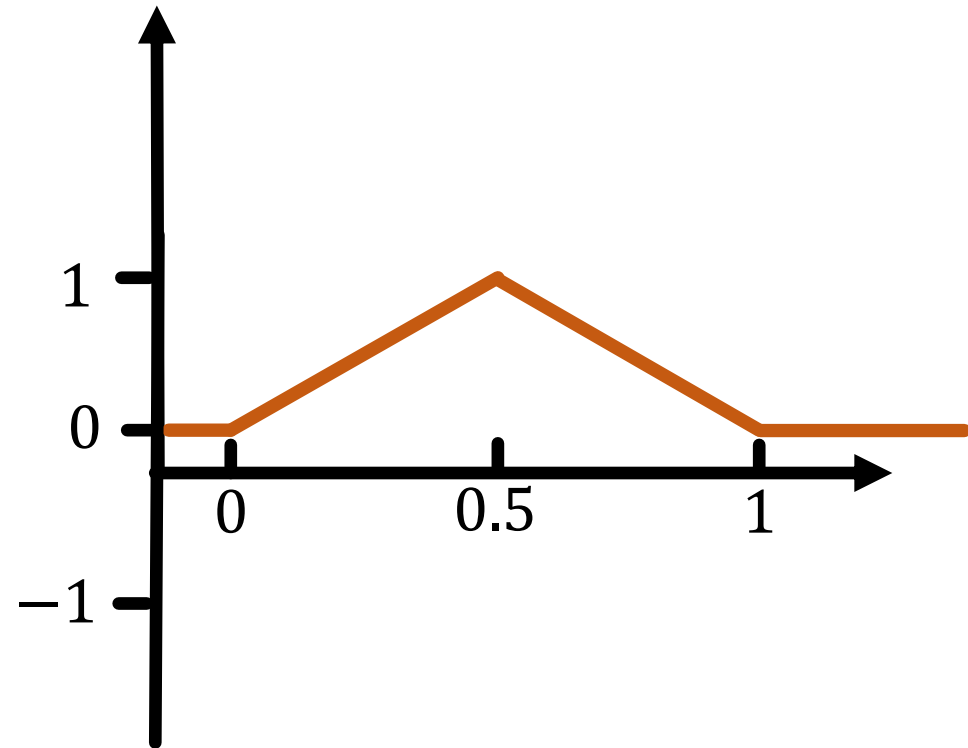$$\Delta(x) = 2\left[\text{ReLU}(x) - 2 \cdot \text{ReLU}\left(x - \frac{1}{2}\right) + \text{ReLU}(x - 1)\right]$$

# What happens if we compose the hat function with itself?

For $x \in [0,1]$, $\Delta(x)$ covers the range $[0,1]$ twice in its output. $\longrightarrow$ $\Delta \circ \Delta(x) = 2$ hats

$\Delta \circ \Delta(x)$ covers the range $[0,1]$ four times in its output. $\longrightarrow$ $\Delta \circ \Delta \circ \Delta(x) = 4$ hats



$n = 1$  $\qquad$  $n = 2$  $\qquad$  $n = 3$  $\qquad$  $n = 6$

$\Delta^n(x)$ **has $2^n$ pieces** *(each hat has 2)*

$\Delta^n(x)$ *is also known as the "sawtooth function"*

# The power of deep neural networks

To calculate $\Delta^n(x)$, we just have to stack our small ReLU neural network $n$ times.



Thus, there is a deep ReLU ANN with $\mathcal{O}(n)$ neurons and $n$ layers that perfectly fits $\Delta^n(x)$. In contrast, a shallow network has $2 \cdot \textbf{width}$ pieces $\rightarrow 2^{n-1}$ neurons needed!

This is a clear separation in the capabilities of shallow and deep neural networks!

# The square function

**But there is more!** We can use the hat function to build $x^2$ with exponential error rate (in $n$)!

$$s_n(x) = x - \sum_{j=1}^{n} \frac{\Delta^j(x)}{4^j} \qquad \text{for } x \in [0,1]$$



*D. Yarotsky, Error bounds for approximations with deep ReLU networks (2017)*

# A neural network approximating the square function

This can again be realized efficiently by a deep ReLU neural network:

$$s_n(x) = x - \sum_{j=1}^{n} \frac{\Delta^j(x)}{4^j}$$

$$= s_{n-1}(x) - \frac{\Delta^n(x)}{4^n}$$

$$= \text{ReLU}\big(\underbrace{s_{n-1}(x)}_{\text{since this is } > 0}\big) - \frac{\Delta \circ \Delta^{n-1}(x)}{4^n}$$



$s_0(x) = \text{ReLU}(x)$  $s_1(x) = s_0(x) - \dfrac{\Delta(x)}{4}$  $s_2(x)$  $s_{n-1}(x)$

$x$  $\cdots$  $s_n(x)$

$\Delta(x)$  $\Delta^2(x)$  $\Delta^3(x)$

Thus, there is a **deep ReLU ANN** with $\mathcal{O}(n)$ neurons and $n$ layers that approximates $x^2$ with error $\epsilon$ **that decays exponentially in $n$**.
For shallow neural networks, the decay is only polynomial : $|x^2 - \Psi|_\infty \geq C \cdot (2 \cdot \text{width})^{-2}$

# But wait… there is more!

We can use a little trick: **the polarization identity**

$$x \cdot y = \left(\frac{x+y}{2}\right)^2 - \left(\frac{x-y}{2}\right)^2$$

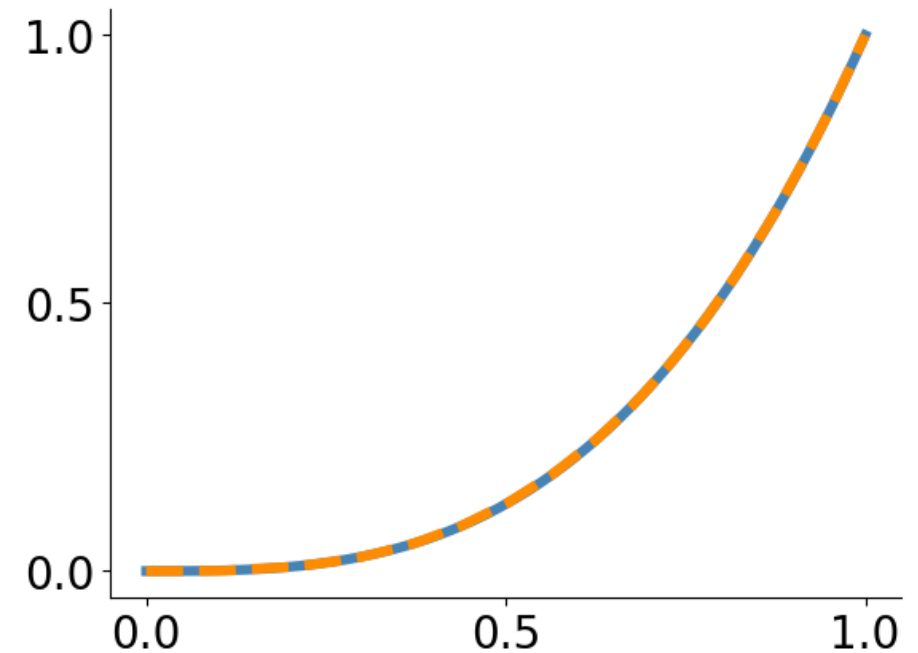This is only addition and squaring, two operations that (as we showed) can be represented using a deep ReLU neural network!

Thus, we can model any multiplication, such as $x^3 = x^2 \cdot x$!



$x^3$ modelled using our ReLU square circuits!

# And… more!

Via multiplication, **we can build polynomials**!

> Thus, **we can approximate polynomials efficiently** with deep ReLU neural networks! In fact, there exists a deep ReLU neural network with $\mathcal{O}\big(k \cdot \log(k/\epsilon)\big)$ parameters and $\mathcal{O}\big(\log(k/\epsilon)\big)$ layers that approximates a **degree $k \geq 2$ polynomial** with **error $\epsilon$**.

In fact, this extends to **Taylor series expansions**, and hence to **any analytic function**!

> To conclude, we now found a rather big and highly relevant set of functions where a shallow neural network requires far more parameters than a deep neural network!

*Note: this is another way to show universal approximation, as we showed that we can approximate polynomials, which are universal approximators, efficiently with deep ReLU neural networks.*

# The reality gap…

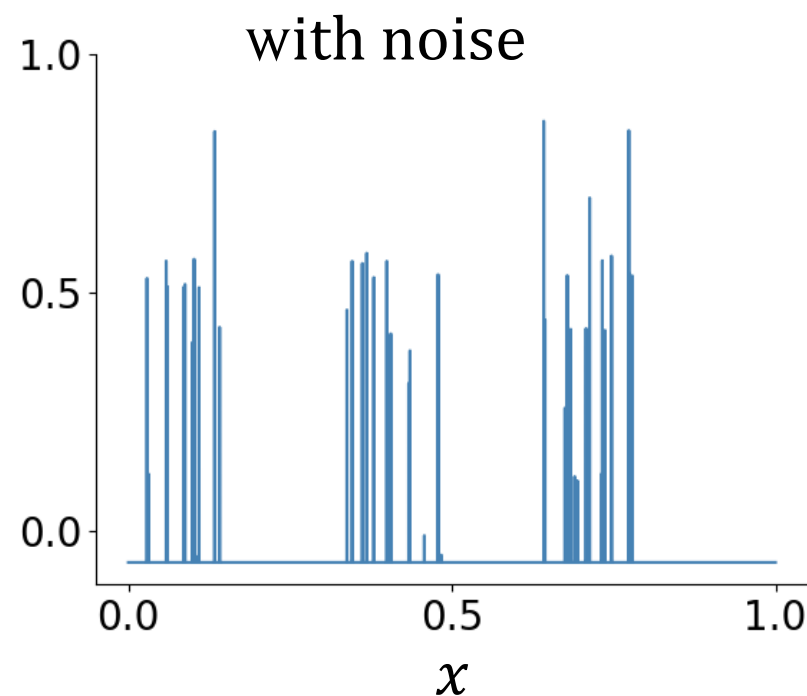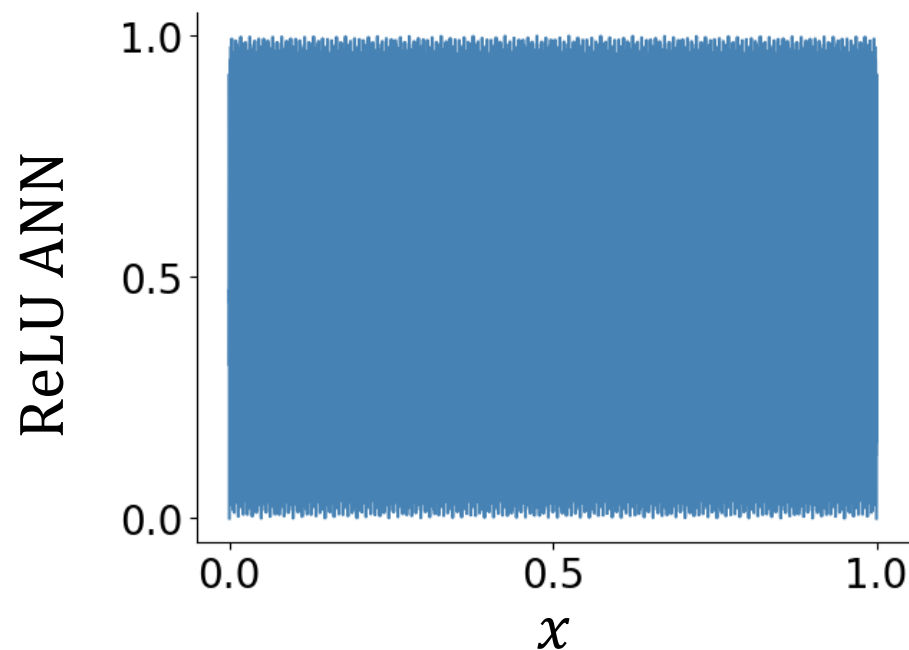**One might ask:** do we find these neural networks with exponential many pieces in practice?

This is an active area of research, but preliminary results suggest: rarely…

In fact, even our construction for $\Delta^n(x)$ collapses to only linearly many pieces if we noise the weights and biases a bit…!

For $n = 30$:



with noise

*B. Hanin & D. Rolnick, Complexity of Linear Regions in Deep Networks (2019)*

# How are neural networks trained in practice?

**Using empirical risk minimization**, i.e., we minimize a loss function.
In case of neural networks, this is usually done via gradient descent.



$$L = \sum_i \|\Psi(x_i) - y_i\|^2 \qquad \longrightarrow \qquad W_j \leftarrow W_j - \eta \nabla_{W_j} L$$

**Intuition:** neural network = **function composed of many small simple functions** (neurons).
Changing the parameters of all these functions a bit leads to a small change in the output.
Thus, we can do **many small steps** to **slowly walk the output to its target!**

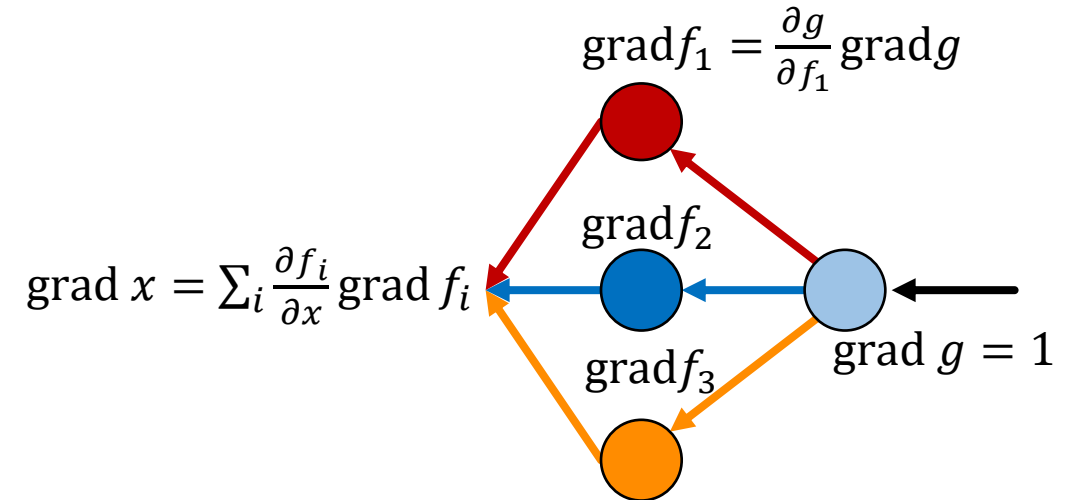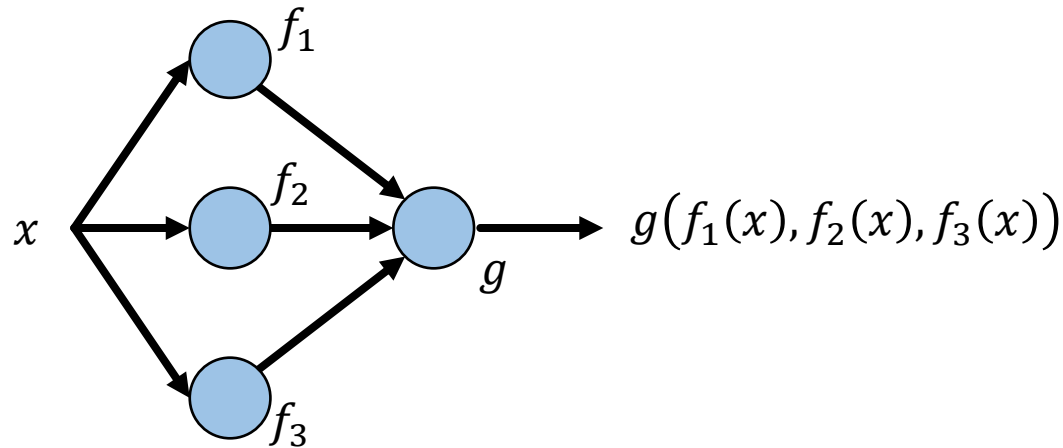# Automatic Differentiation

How can we calculate the gradients efficiently?
Since neural networks are compositions of simple operations, let's first have a look at the derivative of a function composition!



$$g\big(f_1(x), f_2(x), f_3(x)\big)$$

$$\text{grad } x = \sum_i \frac{\partial f_i}{\partial x} \text{grad } f_i$$

$$\text{grad} f_1 = \frac{\partial g}{\partial f_1} \text{grad} g$$

$$\text{grad} f_2$$

$$\text{grad} f_3$$

$$\text{grad } g = 1$$

Using the chain rule: $\dfrac{\partial g}{\partial x} = \displaystyle\sum_{i=1}^{3} \dfrac{\partial f_i}{\partial x}\dfrac{\partial g}{\partial f_i} = \color{red}{\dfrac{\partial f_1}{\partial x}\dfrac{\partial g}{\partial f_1} \cdot 1} + \color{blue}{\dfrac{\partial f_2}{\partial x}\dfrac{\partial g}{\partial f_2} \cdot 1} + \color{orange}{\dfrac{\partial f_3}{\partial x}\dfrac{\partial g}{\partial f_3} \cdot 1}$

# Automatic Differentiation

**Let's have a closer look!**

**Along edges,** we propagate products
(like message passing)!

$$\mathrm{grad}f_1 = \frac{\partial g}{\partial f_1}\,\mathrm{grad}g$$

$$\frac{\partial g}{\partial f_1}$$

grad $g$

# Automatic Differentiation

**Let's have a closer look!**

**Along edges,** we propagate products (like message passing)!

$$\mathrm{grad} f_1 = \frac{\partial g}{\partial f_1} \mathrm{grad} g$$

$$\frac{\partial f_1}{\partial x}$$

$$\frac{\partial f_1}{\partial x} \, \mathrm{grad} f_1$$

$$\frac{\partial g}{\partial f_1}$$

grad $g$

# Automatic Differentiation

**Let's have a closer look!**

**Along edges,** we propagate products
(like message passing)!



$$\frac{\partial f_3}{\partial x}\ \text{grad} f_3$$

$$\frac{\partial f_3}{\partial x}$$

$$\text{grad}\,g$$
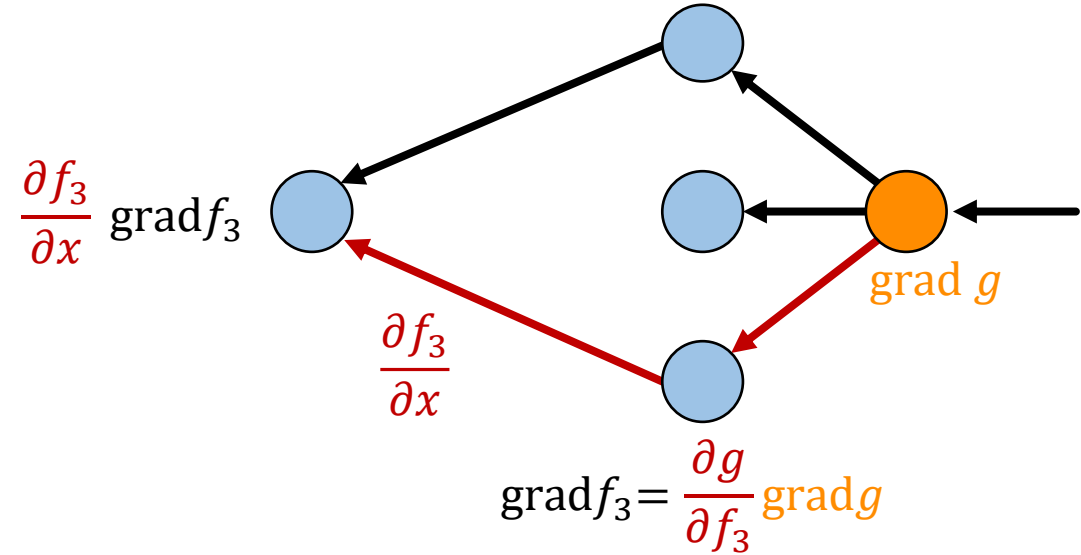
$$\text{grad} f_3 = \frac{\partial g}{\partial f_3}\,\text{grad} g$$

# Automatic Differentiation

**Let's have a closer look!**

**Along edges,** we propagate products (like message passing)!

**At nodes,** we sum up all incoming messages.

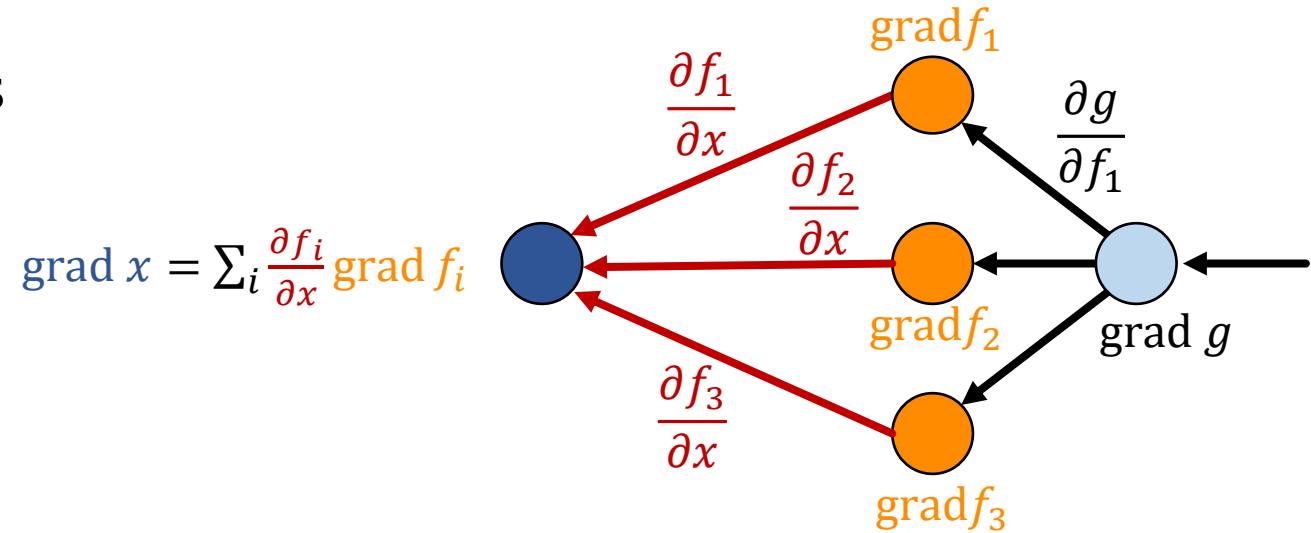$$\text{grad } x = \sum_i \frac{\partial f_i}{\partial x} \text{grad } f_i$$

# Automatic Differentiation

**Let's have a closer look!**

**Along edges,** we propagate products (like message passing)!

**At nodes,** we sum up all incoming messages.

$$\text{grad } x = \sum_i \frac{\partial f_i}{\partial x} \text{ grad } f_i$$



**Thus:** We can represent the derivative as a graph which is the "inverse" of the original computational ("forward") graph! You usually find the following structure:

**Forward function:** our original function, from which we constructed the computational graph.
**Backward function:** calculates the gradient for each node in the graph. Its structure is inverse to the forward graph and gradients are propagated as inputs!
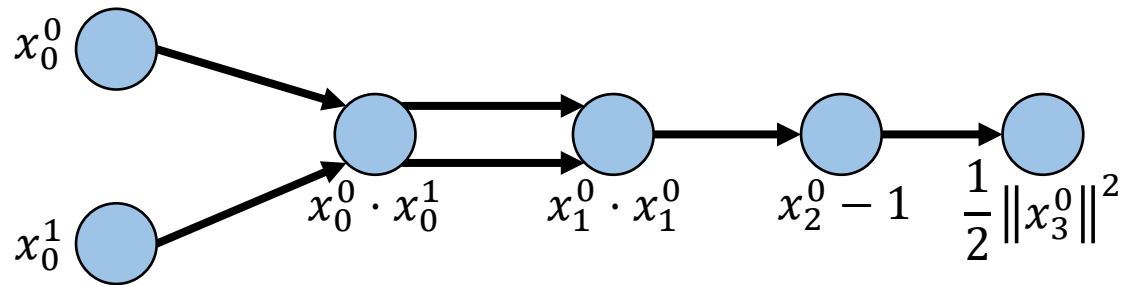
**In tools like PyTorch, every node** in the graph has a **forward** and **backward** function from which the computational graph is constructed node by node!

# Example

Let's try this for the following function: $L(x) = \frac{1}{2}\|(w \cdot x)^2 - 1\|^2$

**Build the forward graph.** We denote by $x_j^i$ the output of the $i$'th node in layer $j$.



$x_0^0$

$x_0^1$

$x_0^0 \cdot x_0^1$

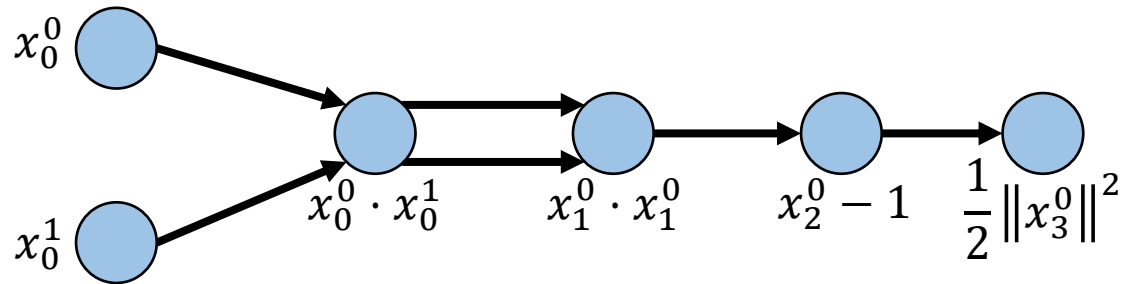$x_1^0 \cdot x_1^0$

$x_2^0 - 1$

$\frac{1}{2}\|x_3^0\|^2$

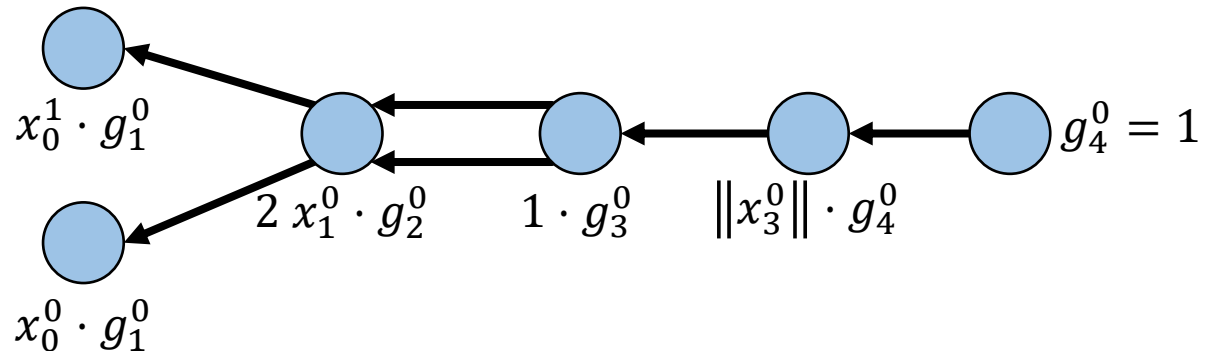*$L(x)$ is recovered when setting $x_0^0 = x$ and $x_0^1 = w$.*

# Example

Let's try this for the following function: $L(x) = \frac{1}{2}\|(w \cdot x)^2 - 1\|^2$

**Build the forward graph.** We denote by $x_j^i$ the output of the $i$'th node in layer $j$.



$x_0^0 \cdot x_0^1 \qquad x_1^0 \cdot x_1^0 \qquad x_2^0 - 1 \qquad \frac{1}{2}\|x_3^0\|^2$

*$L(x)$ is recovered when setting $x_0^0 = x$ and $x_0^1 = w$.*

**Build the backward graph.** We denote layer inputs by $g_j^i$.



$x_0^1 \cdot g_1^0$

$2\, x_1^0 \cdot g_2^0 \qquad 1 \cdot g_3^0 \qquad \|x_3^0\| \cdot g_4^0 \qquad g_4^0 = 1$

$x_0^0 \cdot g_1^0$

# Example

Let's try this for the following function: $L(x) = \frac{1}{2}\|(w \cdot x)^2 - 1\|^2$

**Build the forward graph.** We denote by $x_j^i$ the output of the $i$'th node in layer $j$.
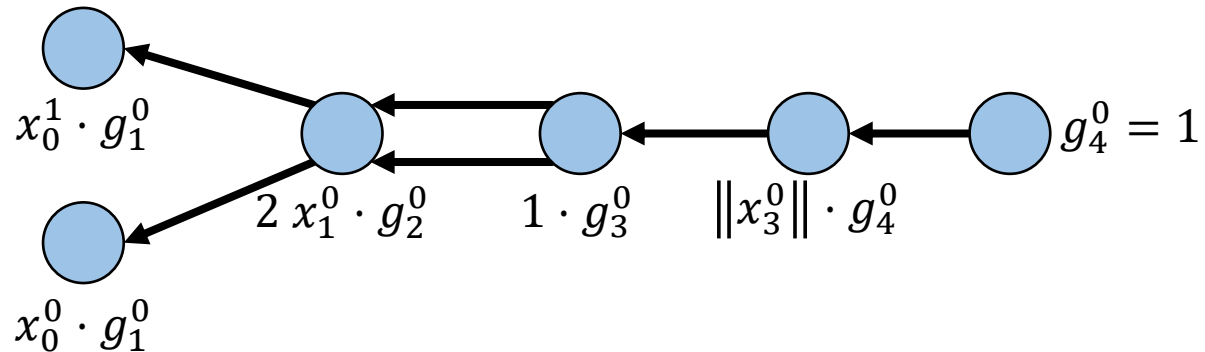
$x_0^0$ ② 

② → ④ → ③ → 4.5

$x_0^1$ ①

$x_0^0 \cdot x_0^1 \qquad x_1^0 \cdot x_1^0 \qquad x_2^0 - 1 \qquad \frac{1}{2}\|x_3^0\|^2$

**Build the backward graph.** We denote layer inputs by $g_j^i$.

$x_0^1 \cdot g_1^0$

$2\,x_1^0 \cdot g_2^0 \qquad 1 \cdot g_3^0 \qquad \|x_3^0\| \cdot g_4^0 \qquad g_4^0 = 1$

$x_0^0 \cdot g_1^0$

# Example

Let's try this for the following function: $L(x) = \frac{1}{2}\|(w \cdot x)^2 - 1\|^2$

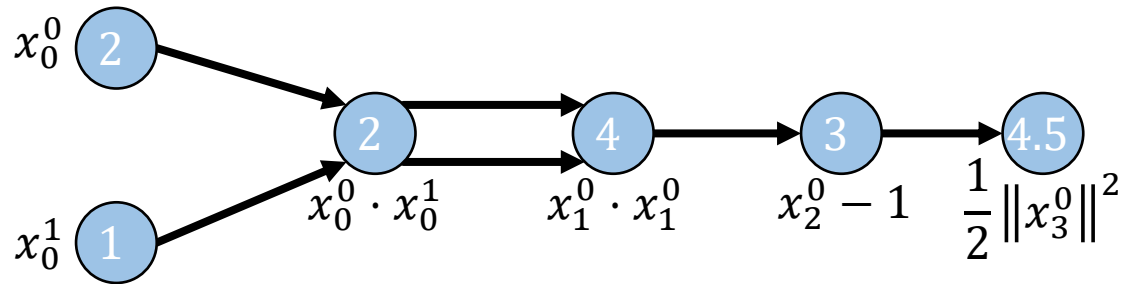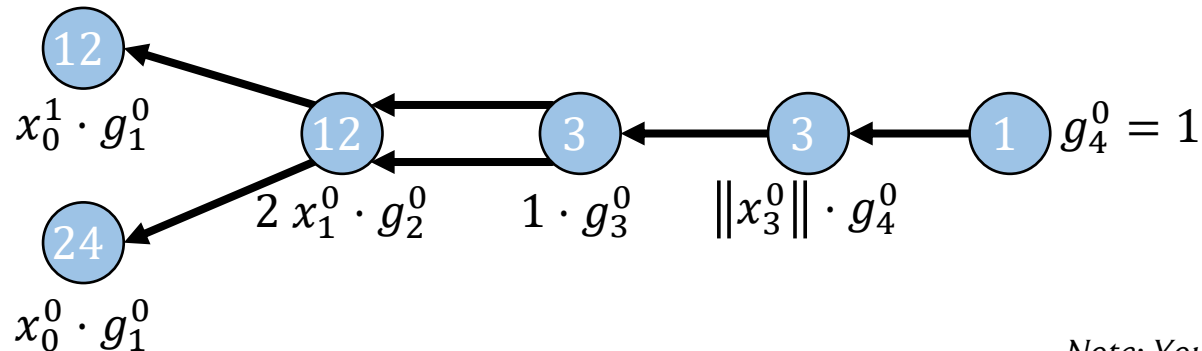**Build the forward graph.** We denote by $x_j^i$ the output of the $i$'th node in layer $j$.



When our function is a deep neural network, then this method is called **error backpropagation**.

Why? Because the gradient of the last node is the **error**! In the example here, that's $\|x_3^0\|$

**Build the backward graph.** We denote layer inputs by $g_j^i$.

*See also: coding examples.*



*Note: You can also construct a forward graph to calculate the gradients. This is called forward-mode automatic differentiation. But commonly, backward-mode is used in practice.*