

# Approximation speed of quantized vs. unquantized ReLU neural networks and beyond

Antoine Gonon

ENS Lyon

Joint work with: Nicolas Brisebarre, Rémi Gribonval, Elisa Riccietti

Curves & Surfaces, Arcachon, June 22, 2022

# ReLU neural networks

$$\begin{array}{c} \text{parameters} \\ (A_1, \dots, A_L, b_1, \dots, b_L) \\ \underbrace{\hspace{1.5cm}} \quad \underbrace{\hspace{1.5cm}} \\ \text{matrices} \quad \quad \text{vectors} \end{array}$$

function represented

# ReLU neural networks

$$\underbrace{(A_1, \dots, A_L, b_1, \dots, b_L)}_{\theta}$$

function represented

# ReLU neural networks

$$\underbrace{(A_1, \dots, A_L, b_1, \dots, b_L)}_{\theta} \quad x \mapsto$$

$$\text{function represented} \\ A_1 x + b_1$$

# ReLU neural networks

$$\underbrace{(A_1, \dots, A_L, b_1, \dots, b_L)}_{\theta} \quad x \mapsto \quad \text{function represented} \quad \rho(A_1 x + b_1)$$

with  $\rho(x) = \max(0, x) = x_+$  the ReLU function

# ReLU neural networks

$$\underbrace{(A_1, \dots, A_L, b_1, \dots, b_L)}_{\theta} \quad x \mapsto \quad \text{function represented} \quad A_2 \rho(A_1 x + b_1) + b_2$$

with  $\rho(x) = \max(0, x) = x_+$  the ReLU function

# ReLU neural networks

$$\underbrace{(A_1, \dots, A_L, b_1, \dots, b_L)}_{\theta} \quad x \mapsto \quad \text{function represented} \quad \rho(A_2 \rho(A_1 x + b_1) + b_2)$$

with  $\rho(x) = \max(0, x) = x_+$  the ReLU function

## ReLU neural networks

$$\underbrace{(A_1, \dots, A_L, b_1, \dots, b_L)}_{\theta} \quad x \mapsto A_L \cdots \rho(A_2 \rho(A_1 x + b_1) + b_2) \cdots + b_L$$

with  $\rho(x) = \max(0, x) = x_+$  the ReLU function



# ReLU neural networks

$$\underbrace{(A_1, \dots, A_L, b_1, \dots, b_L)}_{\theta} \quad \underbrace{x \mapsto A_L \cdots \rho(A_2 \rho(A_1 x + b_1) + b_2) \cdots + b_L}_{R_\theta(x)}$$

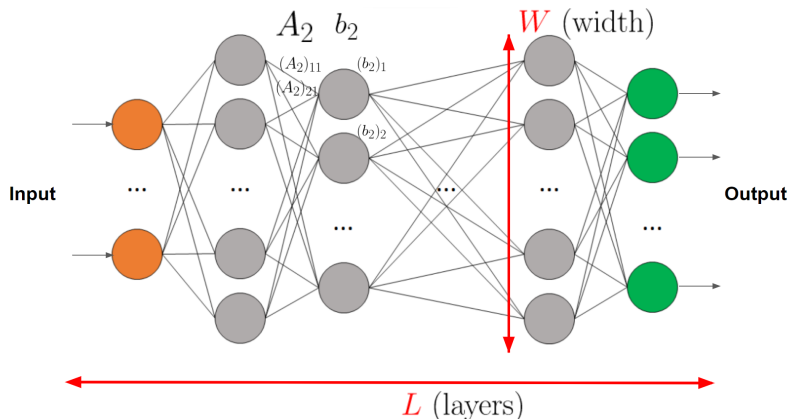
The equation shows the mapping from parameters  $\theta$  to the function  $R_\theta(x)$  represented by the neural network. The parameters are grouped under a brace labeled  $\theta$ , and the function is grouped under a brace labeled  $R_\theta(x)$ .

with  $\rho(x) = \max(0, x) = x_+$  the ReLU function

# ReLU neural networks

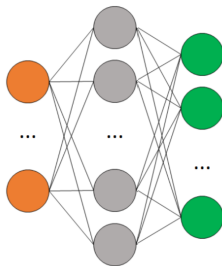
$$\underbrace{(A_1, \dots, A_L, b_1, \dots, b_L)}_{\theta} \quad \underbrace{x \mapsto A_L \cdots \rho(A_2 \rho(A_1 x + b_1) + b_2) \cdots + b_L}_{R_\theta(x)}$$

with  $\rho(x) = \max(0, x) = x_+$  the ReLU function



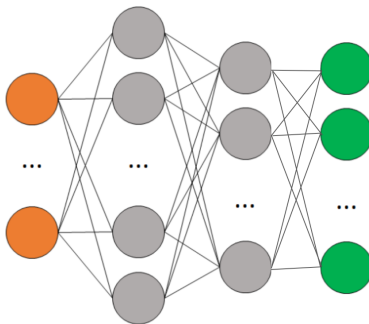
## Context: approximation with increasing complexity

**Consider:**  $\Sigma_M$  set of networks with complexity increasing with  $M$



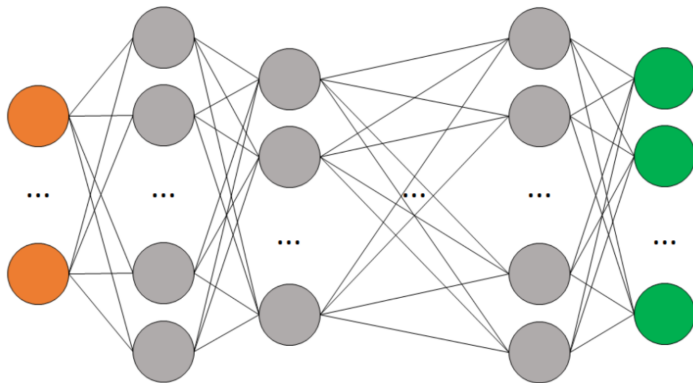
## Context: approximation with increasing complexity

**Consider:**  $\Sigma_M$  set of networks with complexity increasing with  $M$



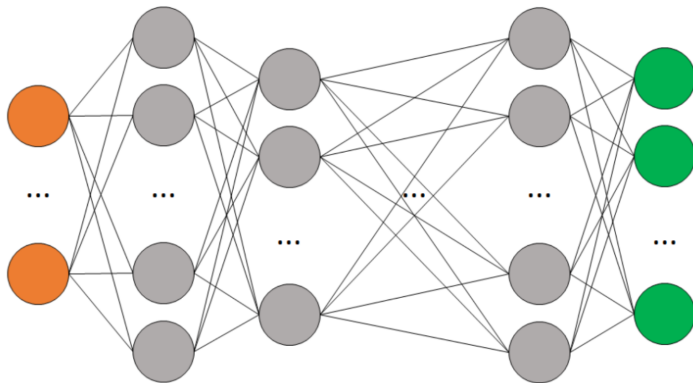
## Context: approximation with increasing complexity

**Consider:**  $\Sigma_M$  set of networks with complexity increasing with  $M$



## Context: approximation with increasing complexity

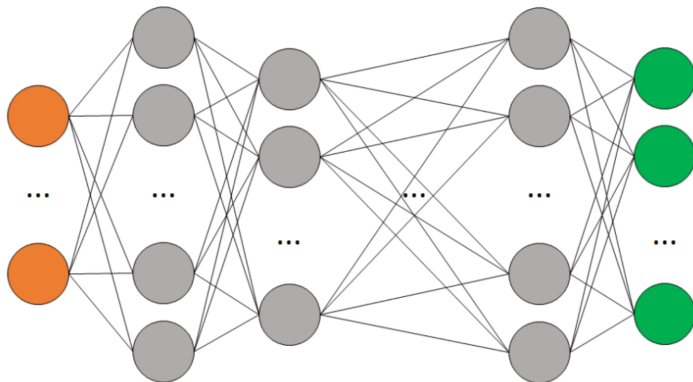
**Consider:**  $\Sigma_M$  set of networks with complexity increasing with  $M$



complexity  $\nearrow \Rightarrow$  approximation error  $\searrow$

## Context: approximation with increasing complexity

**Consider:**  $\Sigma_M$  set of networks with complexity increasing with  $M$



complexity  $\nearrow \Rightarrow$  approximation error  $\searrow$

**Typical known result:**

$$f \text{ "}\gamma\text{-smooth" } \Rightarrow d(f, \Sigma_M) \lesssim M^{-\gamma}$$

# Challenges

## Definition:

$$\gamma^{\text{*approx}}(\mathcal{C}|\Sigma) := \text{largest } \gamma > 0 \text{ s.t. } \sup_{f \in \mathcal{C}} d(f, \Sigma_M) \underset{M \rightarrow \infty}{=} \mathcal{O}(M^{-\gamma})$$



# Challenges

## Definition:

$$\gamma^{\text{*approx}}(\mathcal{C}|\Sigma) := \text{largest } \gamma > 0 \text{ s.t. } \sup_{f \in \mathcal{C}} d(f, \Sigma_M) \underset{M \rightarrow \infty}{=} \mathcal{O}(M^{-\gamma})$$

### 1 Quantization vs. approximation

# Challenges

## Definition:

$$\gamma^{\text{approx}}(\mathcal{C}|\Sigma) := \text{largest } \gamma > 0 \text{ s.t. } \sup_{f \in \mathcal{C}} d(f, \Sigma_M) \underset{M \rightarrow \infty}{=} \mathcal{O}(M^{-\gamma})$$

### 1 Quantization vs. approximation

**Context:**  $\gamma^{\text{approx}}(\mathcal{C}|\Sigma)$  known  
for ReLU neural networks with  
**weights in  $\mathbb{R}$**  and certain  $\mathcal{C}$ 's

# Challenges

## Definition:

$$\gamma^{\text{*approx}}(\mathcal{C}|\Sigma) := \text{largest } \gamma > 0 \text{ s.t. } \sup_{f \in \mathcal{C}} d(f, \Sigma_M) \underset{M \rightarrow \infty}{=} \mathcal{O}(M^{-\gamma})$$

### 1 Quantization vs. approximation

**Context:**  $\gamma^{\text{*approx}}(\mathcal{C}|\Sigma)$  known  
for ReLU neural networks with  
**weights in  $\mathbb{R}$**  and certain  $\mathcal{C}$ 's

**Question:** **Quantized weights?**

# Challenges

## Definition:

$$\gamma^{\text{approx}}(\mathcal{C}|\Sigma) := \text{largest } \gamma > 0 \text{ s.t. } \sup_{f \in \mathcal{C}} d(f, \Sigma_M) \underset{M \rightarrow \infty}{=} \mathcal{O}(M^{-\gamma})$$

① Quantization vs. approximation

② Relation between approximation- and information-theoretic quantities

**Context:**  $\gamma^{\text{approx}}(\mathcal{C}|\Sigma)$  known for ReLU neural networks with **weights in  $\mathbb{R}$**  and certain  $\mathcal{C}$ 's

**Question:** **Quantized weights?**

# Challenges

## Definition:

$$\gamma^{\text{approx}}(\mathcal{C}|\Sigma) := \text{largest } \gamma > 0 \text{ s.t. } \sup_{f \in \mathcal{C}} d(f, \Sigma_M) \underset{M \rightarrow \infty}{=} \mathcal{O}(M^{-\gamma})$$

### ① Quantization vs. approximation

**Context:**  $\gamma^{\text{approx}}(\mathcal{C}|\Sigma)$  known for ReLU neural networks with **weights in  $\mathbb{R}$**  and certain  $\mathcal{C}$ 's

**Question: Quantized weights?**

### ② Relation between approximation- and information-theoretic quantities

**Context:** known inequality in **certain situations**

$$\underbrace{\gamma^{\text{approx}}(\mathcal{C}|\Sigma)}_{\text{approximation theory}} \leq \underbrace{\text{complexity}(\mathcal{C})}_{\text{information theory}}$$

# Challenges

## Definition:

$$\gamma^{\text{approx}}(\mathcal{C}|\Sigma) := \text{largest } \gamma > 0 \text{ s.t. } \sup_{f \in \mathcal{C}} d(f, \Sigma_M) \underset{M \rightarrow \infty}{=} \mathcal{O}(M^{-\gamma})$$

### ① Quantization vs. approximation

**Context:**  $\gamma^{\text{approx}}(\mathcal{C}|\Sigma)$  known for ReLU neural networks with **weights in  $\mathbb{R}$**  and certain  $\mathcal{C}$ 's

**Question:** **Quantized weights?**

### ② Relation between approximation- and information-theoretic quantities

**Context:** known inequality in **certain situations**

$$\underbrace{\gamma^{\text{approx}}(\mathcal{C}|\Sigma)}_{\text{approximation theory}} \leq \underbrace{\text{complexity}(\mathcal{C})}_{\text{information theory}}$$

**Question:** **Unified framework?**

## Part 1: quantization vs. approximation

A key tool: bound on the Lipschitz constant of the parameterization

$$\|f - R_{Q(\theta)}\| \leq \|f - R_\theta\| + \underbrace{\|R_\theta - R_{Q(\theta)}\|}_{\text{quantization error}}$$

**Question:** Tradeoff number of bits/quantization error



# A key tool: bound on the Lipschitz constant of the parameterization

$$\|f - R_{Q(\theta)}\| \leq \|f - R_\theta\| + \underbrace{\|R_\theta - R_{Q(\theta)}\|}_{\text{quantization error}}$$

**Question:** Tradeoff number of bits/quantization error

**Known result**<sup>1</sup>: On every bounded set of parameters  $\Theta$ , there exists  $K_\Theta > 0$  s.t. for every  $\theta, \theta' \in \Theta$ :

$$\|R_\theta - R_{\theta'}\|_{L^p} \leq K_\Theta \|\theta - \theta'\|_\infty$$

*Explicit bounds on  $K_\Theta$ ?*

---

<sup>1</sup>Neural network approximation. R. DeVore et al. 2021.

## A key tool: bound on the Lipschitz constant of the parameterization

$$\|f - R_{Q(\theta)}\| \leq \|f - R_\theta\| + \underbrace{\|R_\theta - R_{Q(\theta)}\|}_{\text{quantization error}}$$

**Question:** Tradeoff number of bits/quantization error

**Known result**<sup>1</sup>: On every bounded set of parameters  $\Theta$ , there exists  $K_\Theta > 0$  s.t. for every  $\theta, \theta' \in \Theta$ :

$$\|R_\theta - R_{\theta'}\|_{L^p} \leq K_\Theta \|\theta - \theta'\|_\infty$$

*Explicit bounds on  $K_\Theta$ ?*

**To understand the tradeoff: new explicit bounds in terms of the depth, the width and a bound on the weights of the network**

---

<sup>1</sup>Neural network approximation. R. DeVore et al. 2021.

# Contribution 1: explicit bounds on the Lipschitz parameterization of ReLU networks and its consequences

Under mild assumptions, there exists  $c > 0$  s.t.

$$\frac{1}{c}LB^{L-1} \leq K_{\Theta_{L,W}(B)} \leq cWL \times LB^{L-1}$$

**Definition of  $\Theta_{L,W}(B)$ :**

- depth =  $L \in \mathbb{N}^*$
- width =  $W \in \mathbb{N}^*$
- bound  $B \geq 1$  on  $\theta = (A_1, \dots, A_L, b_1, \dots, b_L) : \|A_\ell\|_2, \|b_\ell\|_2 \leq B$

# Contribution 1: explicit bounds on the Lipschitz parameterization of ReLU networks and its consequences

Under mild assumptions, there exists  $c > 0$  s.t.

$$\frac{1}{c}LB^{L-1} \leq K_{\Theta_{L,W}(B)} \leq c \quad \underbrace{WL}_{\text{can be improved?}} \quad LB^{L-1}$$

**Definition of  $\Theta_{L,W}(B)$ :**

- depth =  $L \in \mathbb{N}^*$
- width =  $W \in \mathbb{N}^*$
- bound  $B \geq 1$  on  $\theta = (A_1, \dots, A_L, b_1, \dots, b_L) : \|A_\ell\|_2, \|b_\ell\|_2 \leq B$

# Setup to get guarantees on quantization with our results

**Recipe for quantization guarantees of  $\theta$**

# Setup to get guarantees on quantization with our results

## Recipe for quantization guarantees of $\theta$

### Fixed:

- depth =  $L$
  - width =  $W$
  - bound  $B$  on the parameters
  - desired quantization error  $\varepsilon > 0$
- 
-

# Setup to get guarantees on quantization with our results

## Recipe for quantization guarantees of $\theta$

### Fixed:

- depth =  $L$
- width =  $W$
- bound  $B$  on the parameters
- desired quantization error  $\varepsilon > 0$

---

**Result:** necessary and sufficient number of bits/weight to get  $\varepsilon$ -quantization error on  $\Theta_{L,W}(B)$

# Setup to get guarantees on quantization with our results

## Recipe for quantization guarantees of $\theta$

### Fixed:

- depth =  $L$
- width =  $W$
- bound  $B$  on the parameters
- desired quantization error  $\varepsilon > 0$

---

**Ingredient provided by our work:** bounds on  $K_{\Theta_{L,W}(B)}$

---

**Result:** necessary and sufficient number of bits/weight to get  $\varepsilon$ -quantization error on  $\Theta_{L,W}(B)$



## Consequence 1: $\varepsilon$ -accuracy with naive uniform quantization

**Fixed:** depth  $L$ , width  $W$ , bound  $B$  on the parameters, desired quantization error  $\varepsilon > 0$

**Look for:** smallest number of bits/weight to provide  $\varepsilon$ -error with  $Q_\eta(x) = \lfloor x/\eta \rfloor \eta$  applied coordinatewise:

$$\|R_\theta(x) - R_{Q_\eta(\theta)}(x)\|_p \leq \varepsilon.$$

## Consequence 1: $\varepsilon$ -accuracy with naive uniform quantization

**Fixed:** depth  $L$ , width  $W$ , bound  $B$  on the parameters, desired quantization error  $\varepsilon > 0$

**Look for:** smallest number of bits/weight to provide  $\varepsilon$ -error with  $Q_\eta(x) = \lfloor x/\eta \rfloor \eta$  applied coordinatewise:

$$\max_{x \in [0,1]^d} \|R_\theta(x) - R_{Q_\eta(\theta)}(x)\|_p \leq \varepsilon.$$

## Consequence 1: $\varepsilon$ -accuracy with naive uniform quantization

**Fixed:** depth  $L$ , width  $W$ , bound  $B$  on the parameters, desired quantization error  $\varepsilon > 0$

**Look for:** smallest number of bits/weight to provide  $\varepsilon$ -error with  $Q_\eta(x) = \lfloor x/\eta \rfloor \eta$  applied coordinatewise:

$$\max_{\theta \in \Theta_{L,W}(B)} \max_{x \in [0,1]^d} \|R_\theta(x) - R_{Q_\eta(\theta)}(x)\|_p \leq \varepsilon.$$

## Consequence 1: $\varepsilon$ -accuracy with naive uniform quantization

**Fixed:** depth  $L$ , width  $W$ , bound  $B$  on the parameters, desired quantization error  $\varepsilon > 0$

**Look for:** smallest number of bits/weight to provide  $\varepsilon$ -error with  $Q_\eta(x) = \lfloor x/\eta \rfloor \eta$  applied coordinatewise:

$$\max_{\theta \in \Theta_{L,W}(B)} \max_{x \in [0,1]^d} \|R_\theta(x) - R_{Q_\eta(\theta)}(x)\|_p \leq \varepsilon.$$

**Necessary condition?**

## Consequence 1: $\varepsilon$ -accuracy with naive uniform quantization

**Fixed:** depth  $L$ , width  $W$ , bound  $B$  on the parameters, desired quantization error  $\varepsilon > 0$

**Look for:** smallest number of bits/weight to provide  $\varepsilon$ -error with  $Q_\eta(x) = \lfloor x/\eta \rfloor \eta$  applied coordinatewise:

$$\max_{\theta \in \Theta_{L,W}(B)} \max_{x \in [0,1]^d} \|R_\theta(x) - R_{Q_\eta(\theta)}(x)\|_p \leq \varepsilon.$$

**Necessary condition?**

**Number of bits/weight  $\propto \log(1/\varepsilon) + L \log(B)$**

## Consequence 1: $\varepsilon$ -accuracy with naive uniform quantization

**Fixed:** depth  $L$ , width  $W$ , bound  $B$  on the parameters, desired quantization error  $\varepsilon > 0$

**Look for:** smallest number of bits/weight to provide  $\varepsilon$ -error with  $Q_\eta(x) = \lfloor x/\eta \rfloor \eta$  applied coordinatewise:

$$\max_{\theta \in \Theta_{L,W}(B)} \max_{x \in [0,1]^d} \|R_\theta(x) - R_{Q_\eta(\theta)}(x)\|_p \leq \varepsilon.$$

**Necessary condition?**

**Number of bits/weight**  $\propto \log(1/\varepsilon) + L \log(B)$

**Sufficient condition?**

## Consequence 1: $\varepsilon$ -accuracy with naive uniform quantization

**Fixed:** depth  $L$ , width  $W$ , bound  $B$  on the parameters, desired quantization error  $\varepsilon > 0$

**Look for:** smallest number of bits/weight to provide  $\varepsilon$ -error with  $Q_\eta(x) = \lfloor x/\eta \rfloor \eta$  applied coordinatewise:

$$\max_{\theta \in \Theta_{L,W}(B)} \max_{x \in [0,1]^d} \|R_\theta(x) - R_{Q_\eta(\theta)}(x)\|_p \leq \varepsilon.$$

**Necessary condition?**

$\text{Number of bits/weight} \propto \log(1/\varepsilon) + L \log(B)$

**Sufficient condition?**

$\text{Number of bits/weight} \propto \log(1/\varepsilon) + L \log(B) + \underbrace{\log(WL)}_{\text{can be improved?}}$

## Consequence 1: $\varepsilon$ -accuracy with naive uniform quantization

**Fixed:** depth  $L$ , width  $W$ , bound  $B$  on the parameters, desired quantization error  $\varepsilon > 0$

**Look for:** smallest number of bits/weight to provide  $\varepsilon$ -error with  $Q_\eta(x) = \lfloor x/\eta \rfloor \eta$  applied coordinatewise:

$$\max_{\theta \in \Theta_{L,W}(B)} \max_{x \in [0,1]^d} \|R_\theta(x) - R_{Q_\eta(\theta)}(x)\|_p \leq \varepsilon.$$


**Necessary condition?**

**Number of bits/weight**  $\propto \log(1/\varepsilon) + L \log(B)$

**Sufficient condition?**

**Number of bits/weight**  $\propto \log(1/\varepsilon) + L \log(B) + \underbrace{\log(WL)}_{\text{can be improved?}}$

Improvement of the multiplicative constant compared to a known result<sup>2</sup>

<sup>2</sup>Deep Neural Network Approximation Theory. D. Elbrächter et al. 2021. 



## Consequence 2: approximation with quantized networks

**Recipe for quantized approximation of  $f$**

---

---

## Consequence 2: approximation with quantized networks

### Recipe for quantized approximation of $f$

**Fixed:** quantization error  $\varepsilon > 0$ , function  $f$

---

---

## Consequence 2: approximation with quantized networks

### Recipe for quantized approximation of $f$

**Fixed:** quantization error  $\varepsilon > 0$ , function  $f$

---

**Result:** depth  $L$ , width  $W$ , bound  $B$  on the parameters, and number of bits/weight to get error  $\varepsilon$

## Consequence 2: approximation with quantized networks

### Recipe for quantized approximation of $f$

**Fixed:** quantization error  $\varepsilon > 0$ , function  $f$

---

**Ingredient independent of our work:** depth  $L$ , width  $W$ , bound  $B$  and  $\theta \in \Theta_{L,W}(B)$  s.t.  $\|f - R_\theta\|_p \leq \varepsilon/2$

---

**Result:** depth  $L$ , width  $W$ , bound  $B$  on the parameters, and number of bits/weight to get error  $\varepsilon$

## Consequence 2: approximation with quantized networks

### Recipe for quantized approximation of $f$

**Fixed:** quantization error  $\varepsilon > 0$ , function  $f$

---

**Ingredient independent of our work:** depth  $L$ , width  $W$ , bound  $B$  and  $\theta \in \Theta_{L,W}(B)$  s.t.  $\|f - R_\theta\|_p \leq \varepsilon/2$

**Ingredient provided by our work:** Lipschitz bound to quantize  $\theta$  within error  $\varepsilon/2$

---

**Result:** depth  $L$ , width  $W$ , bound  $B$  on the parameters, and number of bits/weight to get error  $\varepsilon$

## Consequence 2: approximation with quantized networks

### Recipe for quantized approximation of $f$

**Fixed:** quantization error  $\varepsilon > 0$ , function  $f$

---

**Ingredient independent of our work:** depth  $L$ , width  $W$ , bound  $B$  and  $\theta \in \Theta_{L,W}(B)$  s.t.  $\|f - R_\theta\|_p \leq \varepsilon/2$

**Ingredient provided by our work:** Lipschitz bound to quantize  $\theta$  within error  $\varepsilon/2$

---

**Result:** depth  $L$ , width  $W$ , bound  $B$  on the parameters, and number of bits/weight to get error  $\varepsilon$

**Applying this recipe:** **Recovery of a known result**<sup>3</sup> in Sobolev spaces using an external ingredient<sup>4</sup>

---

<sup>3</sup>On the Universal Approximability and Complexity Bounds of Quantized ReLU Neural Networks. Y. Ding et al. 2019.

<sup>4</sup>Error bounds for approximations with deep ReLU networks. D. Yarotsky. 2017

## Consequence 3: approximation speed of quantized ReLU networks

**Consider:** Increasingly complex  $\Sigma = (\Sigma_M)_{M \in \mathbb{N}}$  with weights in  $\mathbb{R}$

**Can we guarantee:**

$$\gamma^{*\text{approx}}(\mathcal{C} | \mathbf{Q}(\Sigma)) = \gamma^{*\text{approx}}(\mathcal{C} | \Sigma)?$$

Recall the definition:

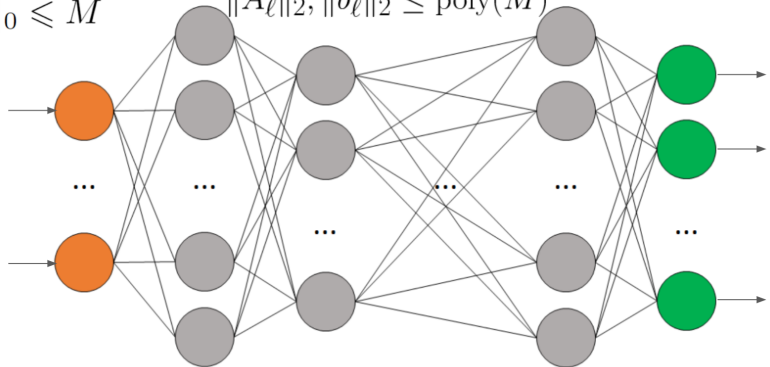
$$\gamma^{*\text{approx}}(\mathcal{C} | \Sigma) := \text{largest } \gamma > 0 \text{ s.t. } \sup_{f \in \mathcal{C}} d(f, \Sigma_M) \underset{M \rightarrow \infty}{=} \mathcal{O}(M^{-\gamma})$$

## Consequence 3: same approximation speed with quantized networks

$(\log M)^2$  **bits per weight is enough** for  $\Sigma_M :=$  functions represented by a ReLU network:

$$\|\theta\|_0 \leq M$$

$$\|A_\ell\|_2, \|b_\ell\|_2 \leq \text{poly}(M)$$





## Part 2: on a relation between approximation- and information-theoretic quantities

# Challenge: unified framework for a known inequality

**Known:** in **certain situations** (arbitrary  $\mathcal{C}$ , sequence  $\Sigma$ : dictionaries<sup>5</sup>, ReLU neural networks<sup>6</sup>)

$$\underbrace{\gamma^{\text{approx}}(\mathcal{C}|\Sigma)}_{\text{approximation theory}} \leq \underbrace{\gamma^{\text{encod}}(\mathcal{C})}_{\text{information theory}}$$

$\gamma^{\text{encod}}(\mathcal{C})$  measures how well  $\mathcal{C}$  can be encoded into bit sequences

**Challenge: Unified framework?**

---

<sup>5</sup>Optimally sparse data representations. P. Grohs. 2015

<sup>6</sup>Deep Neural Network Approximation Theory. D. Elbrächter et al. 2021

## We need to control the increasing complexity

**Can we:** identify a property of  $\Sigma$  to guarantee  $\gamma^{\text{approx}}(\mathcal{C}|\Sigma) \leq \gamma^{\text{encod}}(\mathcal{C})$

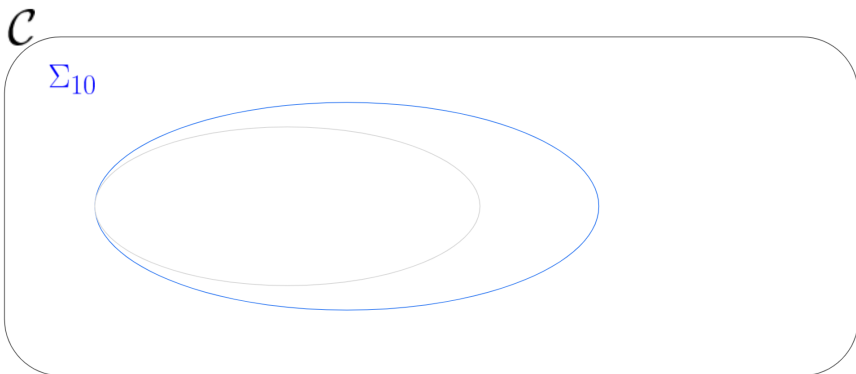
# We need to control the increasing complexity

**Can we:** identify a property of  $\Sigma$  to guarantee  $\gamma^{*\text{approx}}(\mathcal{C}|\Sigma) \leq \gamma^{*\text{encod}}(\mathcal{C})$



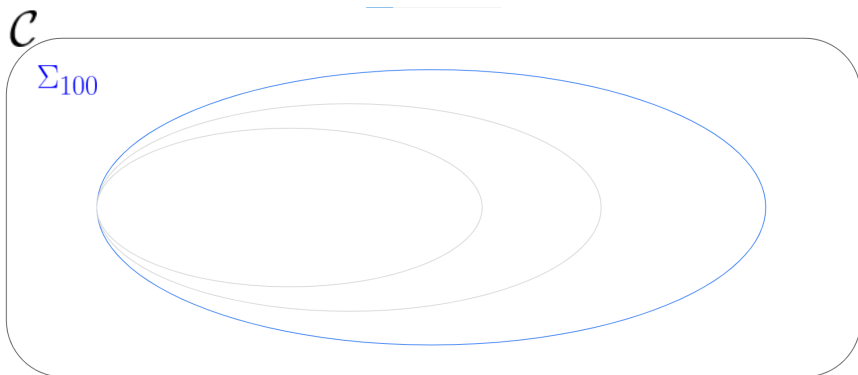
# We need to control the increasing complexity

**Can we:** identify a property of  $\Sigma$  to guarantee  $\gamma^{*\text{approx}}(\mathcal{C}|\Sigma) \leq \gamma^{*\text{encod}}(\mathcal{C})$



# We need to control the increasing complexity

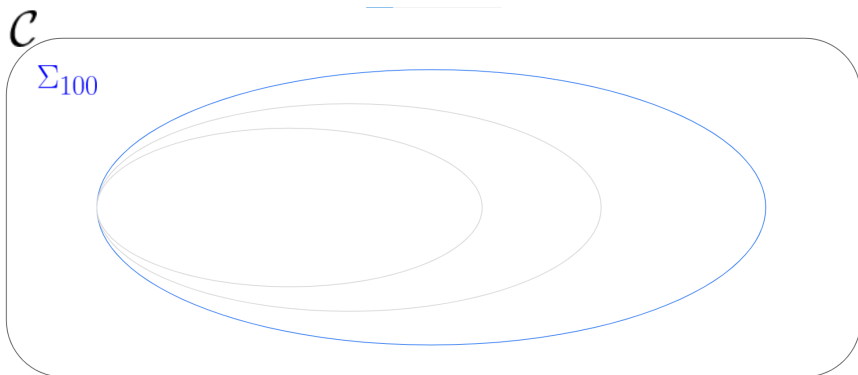
**Can we:** identify a property of  $\Sigma$  to guarantee  $\gamma^{\text{approx}}(\mathcal{C}|\Sigma) \leq \gamma^{\text{encod}}(\mathcal{C})$



**Observation:** intuitively, if  $\Sigma_M$  "grows" too fast then  $\gamma^{\text{approx}}(\mathcal{C}|\Sigma)$  may be unreasonably large

# We need to control the increasing complexity

**Can we:** identify a property of  $\Sigma$  to guarantee  $\gamma^{\text{approx}}(\mathcal{C}|\Sigma) \leq \gamma^{\text{encod}}(\mathcal{C})$



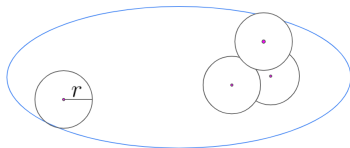
**Observation:** intuitively, if  $\Sigma_M$  "grows" too fast then  $\gamma^{\text{approx}}(\mathcal{C}|\Sigma)$  may be unreasonably large

Extreme case :  $\Sigma_1 = \dots = \Sigma_M = \dots = \mathcal{C}$  where  $\gamma^{\text{approx}}(\mathcal{C}|\Sigma) = \infty$

## Contribution 2: encodability property that allows unification

**A way to measure the "size" of  $\Sigma_M$ :** covering numbers

$\Sigma_1$

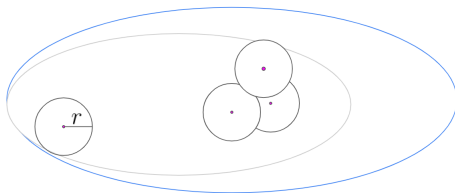




## Contribution 2: encodability property that allows unification

**A way to measure the "size" of  $\Sigma_M$ :** covering numbers

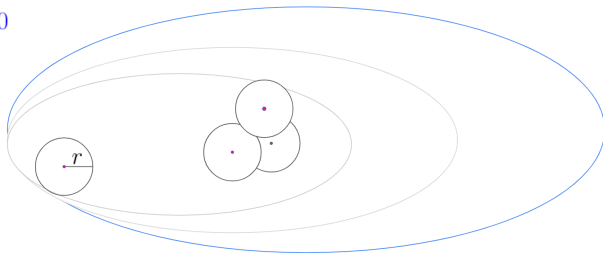
$\Sigma_{10}$



## Contribution 2: encodability property that allows unification

**A way to measure the "size" of  $\Sigma_M$ :** covering numbers

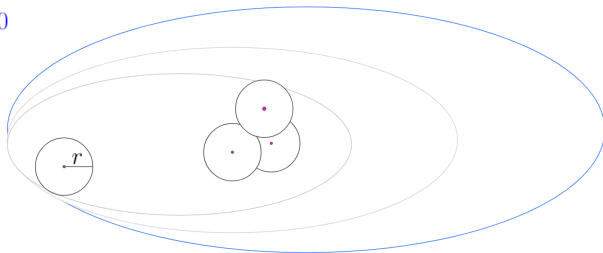
$\Sigma_{100}$



## Contribution 2: encodability property that allows unification

**A way to measure the "size" of  $\Sigma_M$ :** covering numbers

$\Sigma_{100}$

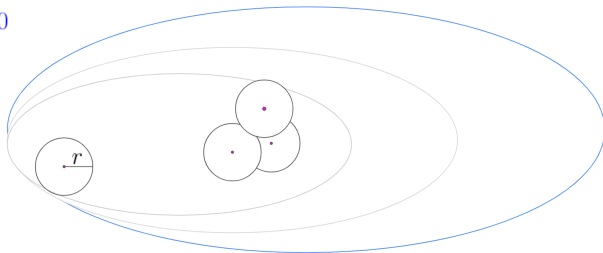


**Introduced property:** **controlled growth (with respect to  $M$ ) of the covering numbers** of  $\Sigma_M$

## Contribution 2: encodability property that allows unification

**A way to measure the "size" of  $\Sigma_M$ :** covering numbers

$\Sigma_{100}$



**Introduced property:** **controlled growth (with respect to  $M$ ) of the covering numbers** of  $\Sigma_M$

**Unification with this encodability property**

## Contribution 2: recovery and generalization of

$$\gamma^{*\text{approx}}(\mathcal{C}|\Sigma) \leq \gamma^{*\text{encod}}(\mathcal{C})$$

① encodability property  $\implies \gamma^{*\text{approx}}(\mathcal{C}|\Sigma) \leq \gamma^{*\text{encod}}(\mathcal{C})$  for every  $\mathcal{C}$

## Contribution 2: recovery and generalization of

$$\gamma^{*\text{approx}}(\mathcal{C}|\Sigma) \leq \gamma^{*\text{encod}}(\mathcal{C})$$

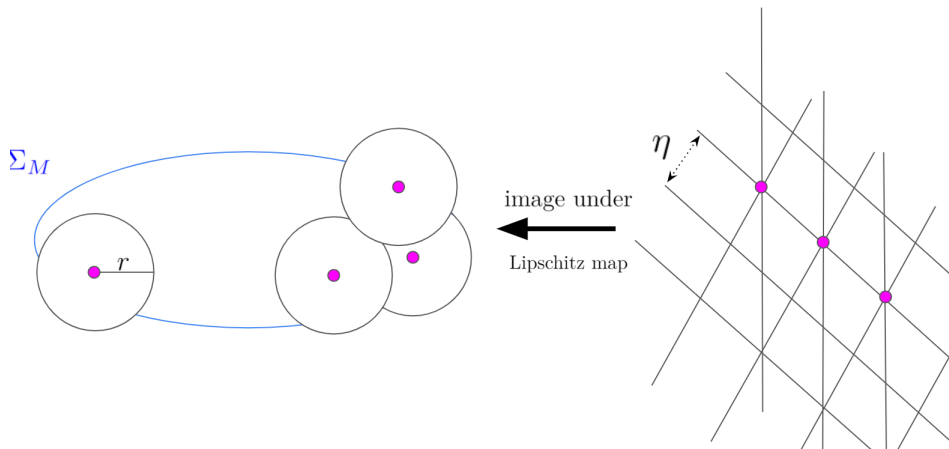
- ① encodability property  $\implies \gamma^{*\text{approx}}(\mathcal{C}|\Sigma) \leq \gamma^{*\text{encod}}(\mathcal{C})$  for every  $\mathcal{C}$
- ② satisfied for  $\Sigma$ : dictionaries, ReLU neural networks

## Contribution 2: recovery and generalization of

$$\gamma^{*\text{approx}}(\mathcal{C}|\Sigma) \leq \gamma^{*\text{encod}}(\mathcal{C})$$

- ① encodability property  $\implies \gamma^{*\text{approx}}(\mathcal{C}|\Sigma) \leq \gamma^{*\text{encod}}(\mathcal{C})$  for every  $\mathcal{C}$
- ② satisfied for  $\Sigma$ : dictionaries, ReLU neural networks

**Proof of the encodability property for ReLU networks:**



## Summary and perspectives



# Conclusion

- 1 **Quantization vs. approximation**

# Conclusion

## ① Quantization vs. approximation

New key tool: **Numbers of bits  
vs. quantization error controlled**  
via new explicit Lipschitz bounds

# Conclusion

## ① Quantization vs. approximation

New key tool: **Numbers of bits vs. quantization error controlled**  
via new explicit Lipschitz bounds

New consequences:

- **characterize  $\varepsilon$ -quantization error** for  $Q_\eta(x) = \lfloor x/\eta \rfloor \eta$   
applied coordinatewise

# Conclusion

## 1 Quantization vs. approximation

New key tool: **Numbers of bits vs. quantization error controlled**  
via new explicit Lipschitz bounds

New consequences:

- **characterize  $\varepsilon$ -quantization error** for  $Q_\eta(x) = \lfloor x/\eta \rfloor \eta$  applied coordinatewise
- sufficient number of bits/weight to keep **same approximation speed with quantized networks** for every function class  $\mathcal{C}$

# Conclusion

## ① Quantization vs. approximation

New key tool: **Numbers of bits vs. quantization error controlled** via new explicit Lipschitz bounds

New consequences:

- **characterize  $\varepsilon$ -quantization error** for  $Q_\eta(x) = \lfloor x/\eta \rfloor \eta$  applied coordinatewise
- sufficient number of bits/weight to keep **same approximation speed with quantized networks** for every function class  $\mathcal{C}$

## ② Relation between approximation- and information-theory

# Conclusion

## ① Quantization vs. approximation

New key tool: **Numbers of bits vs. quantization error controlled** via new explicit Lipschitz bounds

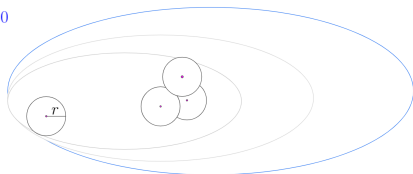
New consequences:

- **characterize  $\varepsilon$ -quantization error** for  $Q_\eta(x) = \lfloor x/\eta \rfloor \eta$  applied coordinatewise
- sufficient number of bits/weight to keep **same approximation speed with quantized networks** for every function class  $\mathcal{C}$

## ② Relation between approximation- and information-theory

New key tool: **Encodability property**

$\Sigma_{100}$



# Conclusion

## ① Quantization vs. approximation

New key tool: **Numbers of bits vs. quantization error controlled** via new explicit Lipschitz bounds

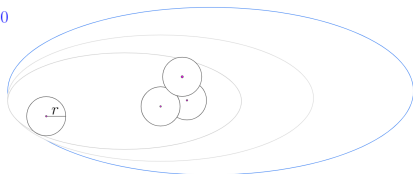
New consequences:

- **characterize  $\varepsilon$ -quantization error** for  $Q_\eta(x) = \lfloor x/\eta \rfloor \eta$  applied coordinatewise
- sufficient number of bits/weight to keep **same approximation speed with quantized networks** for every function class  $\mathcal{C}$

## ② Relation between approximation- and information-theory

New key tool: **Encodability property**

$\Sigma_{100}$



New consequence: **Unifies** situations where

$$\gamma^{\text{approx}}(\mathcal{C}|\Sigma) \leq \text{complexity}(\mathcal{C})$$

## Possible future directions on quantization

**Current theory = number of bits/weight must be linear in the depth:** if  $Q_\eta(x) = \lfloor x/\eta \rfloor \eta$  gives  $\varepsilon$ -quantization error

$$\max_{\theta \in \Theta_{L,W}(B)} \max_{x \in [0,1]^d} \|R_\theta(x) - R_{Q_\eta(\theta)}(x)\|_p \leq \varepsilon$$

**Practice = 1 bit/weight is enough<sup>7</sup>:** quantization-aware training, same performance on MNIST for a 3 hidden-layers with 1024 neurons per layer

---

<sup>7</sup>BinaryConnect: Training Deep Neural Networks with binary weights during propagations. Courbariaux et al. 2015.



## Possible future directions on quantization

**Current theory = number of bits/weight must be linear in the depth:** if  $Q_\eta(x) = \lfloor x/\eta \rfloor \eta$  gives  $\varepsilon$ -quantization error

$$\max_{\theta \in \Theta_{L,W}(B)} \max_{x \in [0,1]^d} \|R_\theta(x) - R_{Q_\eta(\theta)}(x)\|_p \leq \varepsilon$$

**Practice = 1 bit/weight is enough<sup>7</sup>:** quantization-aware training, same performance on MNIST for a 3 hidden-layers with 1024 neurons per layer

**Fill the gap theory/practice?**

---

<sup>7</sup>BinaryConnect: Training Deep Neural Networks with binary weights during propagations. Courbariaux et al. 2015.

# Possible future directions on quantization

**Current theory = number of bits/weight must be linear in the depth:** if  $Q_\eta(x) = \lfloor x/\eta \rfloor \eta$  gives  $\varepsilon$ -quantization error

$$\max_{\theta \in \Theta_{L,W}(B)} \max_{x \in [0,1]^d} \|R_\theta(x) - R_{Q_\eta(\theta)}(x)\|_p \leq \varepsilon$$

**Practice = 1 bit/weight is enough<sup>7</sup>:** quantization-aware training, same performance on MNIST for a 3 hidden-layers with 1024 neurons per layer

## Fill the gap theory/practice?

- Less naive quantization?
- $\varepsilon$ -quantization error on a smaller set  $\Theta \subset \Theta_{L,W}(B)$ ? e.g., *parameters that can be learned in practice*
- not interested in every  $\varepsilon > 0$ ?

---

<sup>7</sup>BinaryConnect: Training Deep Neural Networks with binary weights during propagations. Courbariaux et al. 2015.

# Possible future directions on quantization

**Current theory = number of bits/weight must be linear in the depth:** if  $Q_\eta(x) = \lfloor x/\eta \rfloor \eta$  gives  $\varepsilon$ -quantization error

$$\max_{\theta \in \Theta_{L,W}(B)} \max_{x \in [0,1]^d} \|R_\theta(x) - R_{Q_\eta(\theta)}(x)\|_p \leq \varepsilon$$

**Practice = 1 bit/weight is enough<sup>7</sup>:** quantization-aware training, same performance on MNIST for a 3 hidden-layers with 1024 neurons per layer

## Fill the gap theory/practice?

- Less naive quantization?
- $\varepsilon$ -quantization error on a smaller set  $\Theta \subset \Theta_{L,W}(B)$ ? e.g., *parameters that can be learned in practice*
- not interested in every  $\varepsilon > 0$ ?

More on [agonon.github.io](https://agonon.github.io)

---

<sup>7</sup>BinaryConnect: Training Deep Neural Networks with binary weights during propagations. Courbariaux et al. 2015.

# Possible future directions on quantization

**Current theory = number of bits/weight must be linear in the depth:** if  $Q_\eta(x) = \lfloor x/\eta \rfloor \eta$  gives  $\varepsilon$ -quantization error

$$\max_{\theta \in \Theta_{L,W}(B)} \max_{x \in [0,1]^d} \|R_\theta(x) - R_{Q_\eta(\theta)}(x)\|_p \leq \varepsilon$$

**Practice = 1 bit/weight is enough<sup>7</sup>:** quantization-aware training, same performance on MNIST for a 3 hidden-layers with 1024 neurons per layer

## Fill the gap theory/practice?

- Less naive quantization?
- $\varepsilon$ -quantization error on a smaller set  $\Theta \subset \Theta_{L,W}(B)$ ? e.g., *parameters that can be learned in practice*
- not interested in every  $\varepsilon > 0$ ?

More on [agonon.github.io](https://agonon.github.io)

**Thank you!**

<sup>7</sup>BinaryConnect: Training Deep Neural Networks with binary weights during propagations. Courbariaux et al. 2015.

# Definition of the encoding speed

$$L(\varepsilon, \mathcal{C}) := \inf \left\{ \ell \in \mathbb{N}, \exists E : \mathcal{C} \mapsto \{0, 1\}^\ell, \right. \\ \left. \exists D : \{0, 1\}^\ell \mapsto \mathcal{F}, \sup_{f \in \mathcal{C}} d(f, D(E(f))) \leq \varepsilon \right\}$$


$$\gamma^{\text{*encod}}(\mathcal{C}) := \sup \left\{ \gamma > 0, L(\varepsilon, \mathcal{C}) = \mathcal{O}_{\varepsilon \rightarrow 0} \left( \varepsilon^{-1/\gamma} \right) \right\}.$$

# Known approximation speeds<sup>8</sup>

| $\mathcal{C} :=$ unit ball of |                       | $\Sigma$      | $\gamma^{*\text{approx}}(\mathcal{C} \Sigma) = \gamma^{*\text{encod}}(\mathcal{C})$ |
|-------------------------------|-----------------------|---------------|---|
| $\alpha$ -Hölder              | $C^\alpha([0, 1])$    | Wavelet basis | $\alpha$  |
| $L^p$ -Sobolev <sup>a</sup>   | $W_p^m([0, 1]^d)$     | Wavelet frame | $\frac{m}{d}$   |
| Besov <sup>b</sup>            | $B_{p,q}^m([0, 1]^d)$ | Wavelet frame | $\frac{m}{d}$   |

<sup>a</sup>where  $p \in [1, \infty]$ ,  $m > d \max(1/p - 1/2, 0)$

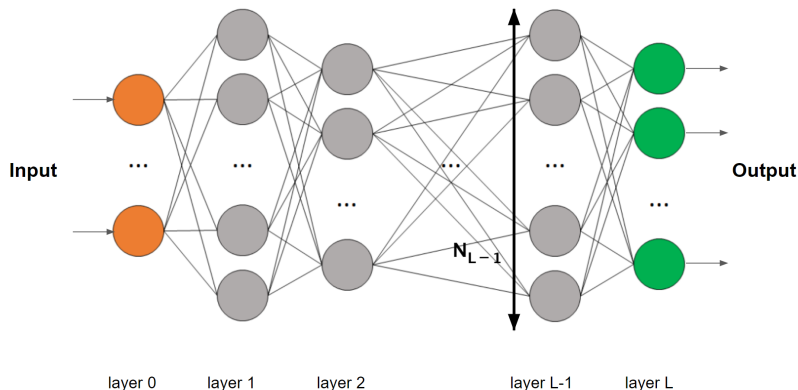
<sup>b</sup>where  $p, q \in (0, \infty]$ ,  $m > d \max(1/p - 1/2, 0)$

<sup>8</sup>Deep Neural Network Approximation Theory. D. Elbrächter et al. 2021. 

# Notations : ReLU neural networks

**Architecture:**  $(L, \mathbf{N})$  where

- $L \in \mathbb{N}$  is the number of layers : the **depth**
- $\mathbf{N} = (N_0, \dots, N_L) \in \mathbb{N}^{L+1}$  with  $N_\ell$  the **width** (number of neurons) of layer  $\ell$



# Notations : ReLU neural networks

**Architecture:**  $(L, \mathbf{N})$  where

- $L \in \mathbb{N}$  is the number of layers : the **depth**
- $\mathbf{N} = (N_0, \dots, N_L) \in \mathbb{N}^{L+1}$  with  $N_\ell$  the **width** (number of neurons) of layer  $\ell$

**Parameters:**  $\theta = (A_1, \dots, A_L, b_1, \dots, b_L)$  with  $A_\ell \in \mathbb{R}^{N_\ell \times N_{\ell-1}}$  and  $b_\ell \in \mathbb{R}^{N_\ell}$



# Notations : ReLU neural networks

**Architecture:**  $(L, \mathbf{N})$  where

- $L \in \mathbb{N}$  is the number of layers : the **depth**
- $\mathbf{N} = (N_0, \dots, N_L) \in \mathbb{N}^{L+1}$  with  $N_\ell$  the **width** (number of neurons) of layer  $\ell$

**Parameters:**  $\theta = (A_1, \dots, A_L, b_1, \dots, b_L)$  with  $A_\ell \in \mathbb{R}^{N_\ell \times N_{\ell-1}}$  and  $b_\ell \in \mathbb{R}^{N_\ell}$

**Realization of the network:**  $R_\theta : \mathbb{R}^{N_0} \rightarrow \mathbb{R}^{N_L}$  the realization of  $\theta$ :

$$R_\theta(x) = A_L \rho(\cdots (A_2 \rho(A_1 x + b_1) + b_2) \cdots) + b_L \text{ avec } \rho(x) = \max(0, x).$$