

# DIVE INTO DEEP LEARNING: A COMPREHENSIVE INTRODUCTION

## FROM AI FUNDAMENTALS TO CUTTING-EDGE DEEP LEARNING TECHNIQUES

**Alexandre Vérine,  
Research Fellow, ENS-PSL  
Université PSL**

Executive Master IASD  
Université Paris-Dauphine, PSL

February 10, 2025

# AI 101: FROM FUNDAMENTALS TO DEEP LEARNING

<b>1</b>	<b>Introduction to Artificial Intelligence . . . . .</b>	<b>7</b>
1.1	Deep Learning in the AI family . . . . .	7
1.2	Representation Learning . . . . .	12
<b>2</b>	<b>Neural Networks Fundamentals . . . . .</b>	<b>17</b>
2.1	Neurons . . . . .	18
2.2	Layers . . . . .	20
2.3	Activation Functions . . . . .	22
<b>3</b>	<b>The Multi-layer Perceptron (MLP) . . . . .</b>	<b>31</b>
3.1	The first Deep Learning Model . . . . .	32
3.2	Stochastic Gradient Descent . . . . .	33
3.3	Back-propagation . . . . .	36
3.4	Example : Image classification of handwritten digits from A to Z . . . . .	58

# DEEP LEARNING IN ACTION: FROM NEURAL NETWORKS TO TRANSFORMER MODELS

<b>1 Convolutional Neural Networks . . . . .</b>	<b>67</b>
1.1 The Two dimensional Convolution . . . . .	68
1.2 CNN : Convolutional in a network Networks . . . . .	76
1.3 CNN in practice: CIFAR 10 . . . . .	83
<b>2 Recurrent Neural Networks . . . . .</b>	<b>108</b>
2.1 Recurrent Block . . . . .	109
2.2 LSTM and GRU . . . . .	111
<b>3 Transformer and Attention Mechanism . . . . .</b>	<b>121</b>
3.1 Self-Attention Mechanism . . . . .	122
3.2 Transformers Model . . . . .	126

# TECHNIQUES TO IMPROVE DEEP LEARNING TRAINING

<b>1 Techniques to Improve Deep Learning Training . . . . .</b>	<b>128</b>
1.1 Data Augmentation . . . . .	129
1.2 Learning Rate Scheduling . . . . .	130
1.3 Early Stopping . . . . .	131
1.4 Gradient Clipping . . . . .	132
1.5 Weight Initialization . . . . .	133
1.6 Regularization . . . . .	135
1.7 GPU Acceleration . . . . .	136

# DEEP LEARNING AND APPLICATIONS

<b>1 Learning to act with Deep Reinforcement Learning . . . . .</b>	<b>138</b>
1.1 Deep Q-Learning . . . . .	138
1.2 The Mouse Game - Reinforcement Learning . . . . .	138
1.3 Q-Learning . . . . .	141
<b>2 Synthetic Data Generation with Generative Adversarial Networks . . . . .</b>	<b>146</b>
2.1 GANS Models . . . . .	146
2.2 MNIST Generation . . . . .	148
<b>3 Sentiment Analysis with Transformers and GRU . . . . .</b>	<b>151</b>
3.1 Bert . . . . .	151
3.2 Sentiment Analysis . . . . .	154
<b>4 Density Estimation with Normalizing Flows . . . . .</b>	<b>159</b>
4.1 Estimating Density . . . . .	159
4.2 Normalizing Flows . . . . .	160
<b>5 Image Segmentation with U-Net . . . . .</b>	<b>172</b>
5.1 Image Segmentation . . . . .	172
5.2 U-Net Architecture . . . . .	175

# YOUR TURN: HANDS-ON DEEP LEARNING

<b>1 Build and use an autoencoder . . . . .</b>	<b>178</b>
1.1 Formal introduction of an autoencoder . . . . .	178
1.2 Applications of U-Net . . . . .	182
<b>2 Image Segmentation with U-Net: Correctly evaluate the model . . . . .</b>	<b>184</b>
2.1 Image Segmentation with U-Net . . . . .	184
2.2 Evaluating the models . . . . .	187

## Part I

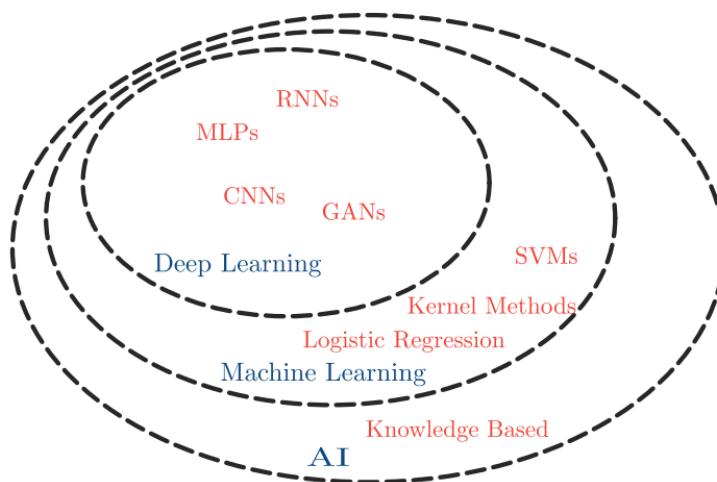
# AI 101: FROM FUNDAMENTALS TO DEEP LEARNING

# INTRODUCTION TO ARTIFICIAL INTELLIGENCE

## DEEP LEARNING IN THE AI FAMILY

In general, among all the class of AI algorithms, we make the difference between 3 sub-categories :

- ▶ **Artificial Intelligence** : human designed program and...
- ▶ **Machine Learning** : human designed features with learned mapping such as Support Vector Machine, Kernels methods, Logistic Regression and ...
- ▶ **Deep Learning**: Learned features with learned mapping such as Multilayer Perceptron, Convolutional Networks, ...



**Figure.** Subsets of Artificial Intelligence

# INTRODUCTION TO ARTIFICIAL INTELLIGENCE

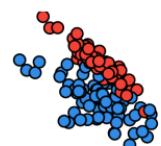
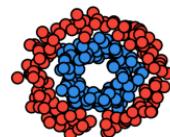
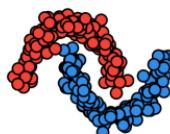
## DEEP LEARNING IN THE AI FAMILY

In the field of Artificial Intelligence, the fundamental objective is to find a function  $f$  that can perform a desired task. This function can either be set by a human or can be learned through training.

For example, in the context of a binary classification task, the goal is to determine  $f(x)$  such that  $f(x) = 0$  when the label of  $x$  is 0 and  $f(x) = 1$  when its label is 1. The choice of AI model impacts the expressivity of the function  $f$ .

For example, a logistic regression model uses a linear function to make decisions, where  $f(x) = \text{sgn}(Ax + b)$ . The expressivity of the model can be increased by using more complex functions, such as polynomials or radial basis functions.

Input data



# INTRODUCTION TO ARTIFICIAL INTELLIGENCE

## CLASSIFICATION TASK

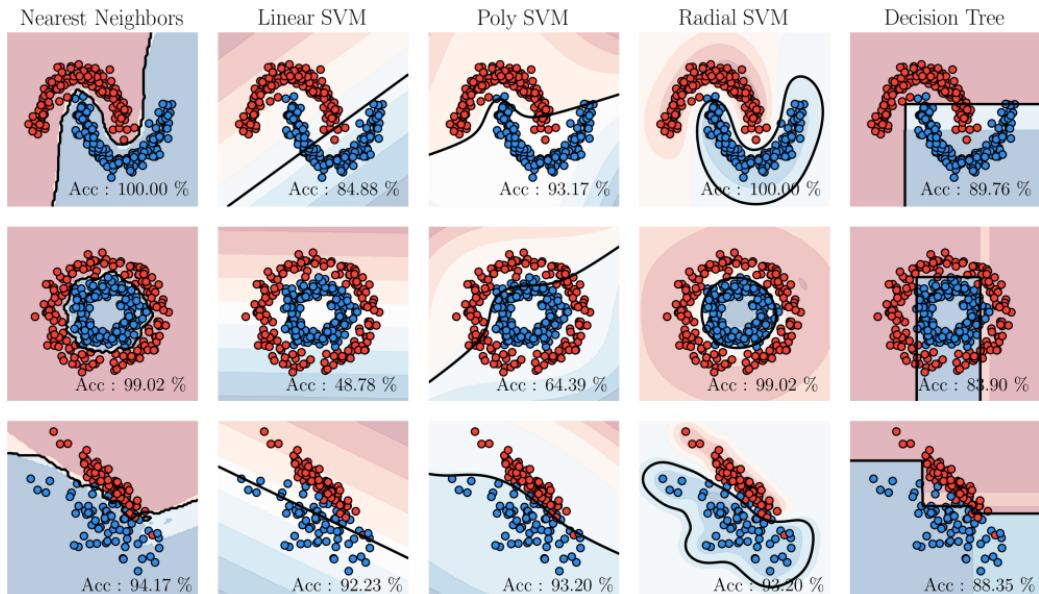


Figure. 2D classification for different AI models.

# INTRODUCTION TO ARTIFICIAL INTELLIGENCE

## THE UNIVERSAL APPROXIMATION THEOREM

The Universal Approximation Theorem is a fundamental result in the field of artificial neural networks. It states that a deep learning model can approximate any function.

### Theorem 1 (Universal Approximation Theorem)

Let  $\mathcal{X} \subset \mathbb{R}^d$  be compact,  $\mathcal{Y} \subset \mathbb{R}^m$ ,  $f : \mathcal{X} \rightarrow \mathcal{Y}$  be a continuous function and  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  be a continuous real function.

Then  $\sigma$  is not polynomial if and only if for every  $\epsilon > 0$ , there exist  $k \in \mathbb{N}$ ,  $A \in \mathbb{R}^{k \times d}$ ,  $b \in \mathbb{R}^k$  and  $C \in \mathbb{R}^{m \times k}$  such that

$$\sup_{x \in \mathcal{X}} \|f(x) - g(x)\| \leq \epsilon$$

where  $g(x) = C \times \sigma(Ax + b)$ .

# INTRODUCTION TO ARTIFICIAL INTELLIGENCE

## CLASSIFICATION TASK

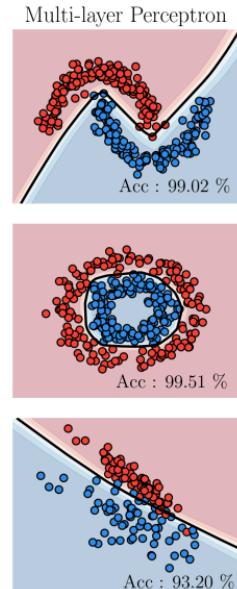


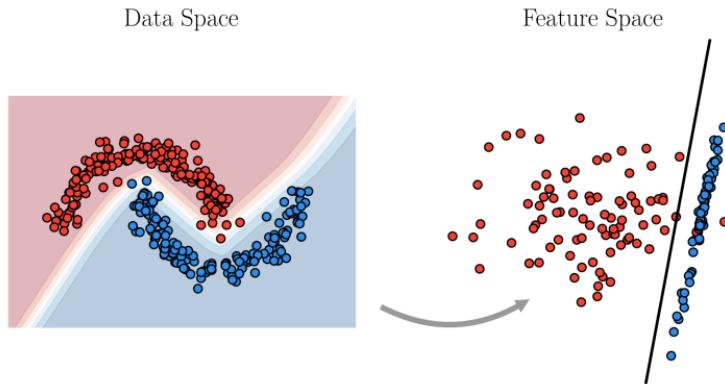
Figure. 2D classification for small Neural Network.

# INTRODUCTION TO ARTIFICIAL INTELLIGENCE

## REPRESENTATION LEARNING

How does deep learning work in practice ?

Deep learning is a subset of representation learning that uses deep neural networks to learn meaningful representations of data. In deep learning, representations are learned through a hierarchy of nonlinear transformations, where each layer of the network builds upon the previous one to extract increasingly abstract and higher-level features from the input data.



# INTRODUCTION TO ARTIFICIAL INTELLIGENCE

## EXAMPLE OF REPRESENTATION LEARNING

Consider the task of recognizing objects in images. A traditional approach would be to hand-engineer features such as edge detectors and color histograms that can be fed into a classifier.

However, with deep learning representation learning, the model learns to automatically discover these features from the data. The network might start by learning simple features such as edges and color blobs in the first layer, then build upon these to learn more complex features such as parts of objects in subsequent layers, until finally, the final layer outputs a probability distribution over classes of objects.

In this way, deep learning of representation enables the model to automatically learn a rich and meaningful representation of the data, without the need for manual feature engineering.

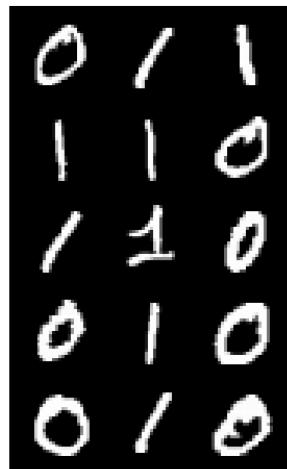


Figure. MNIST

# INTRODUCTION TO ARTIFICIAL INTELLIGENCE

## EXAMPLE OF REPRESENTATION LEARNING

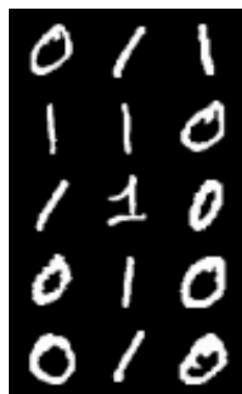


Figure. MNIST : Layer 0

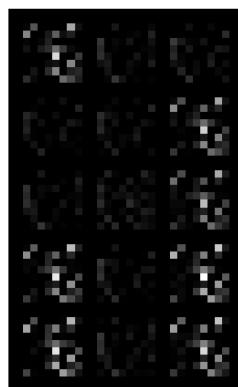


Figure. MNIST : Layer 1

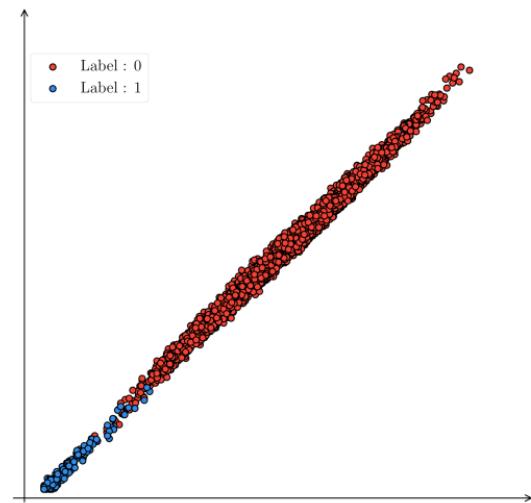


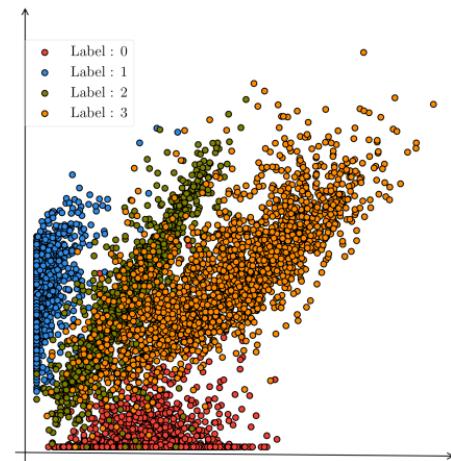
Figure. MNIST : Layer 2

# INTRODUCTION TO ARTIFICIAL INTELLIGENCE

## EXAMPLE OF REPRESENTATION LEARNING



**Figure.** MNIST : Layer 0



**Figure.** MNIST : Layer 2

# INTRODUCTION TO ARTIFICIAL INTELLIGENCE

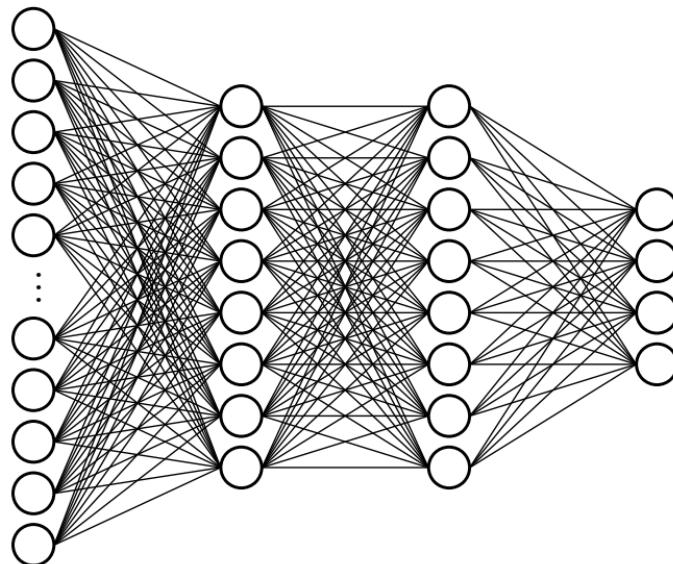
## DEEP LEARNING AND NEURAL NETWORKS

Ok, Deep Learning is a model that learns a good representation of the feature. But how?

- ▶ How does it work ?
- ▶ How can we build a model ?
- ▶ How does it learn ?

## NEURAL NETWORKS FUNDAMENTALS

Typically, a neural network is defined as a computational model composed of interconnected nodes, organised into layers, that perform transformations on input data.

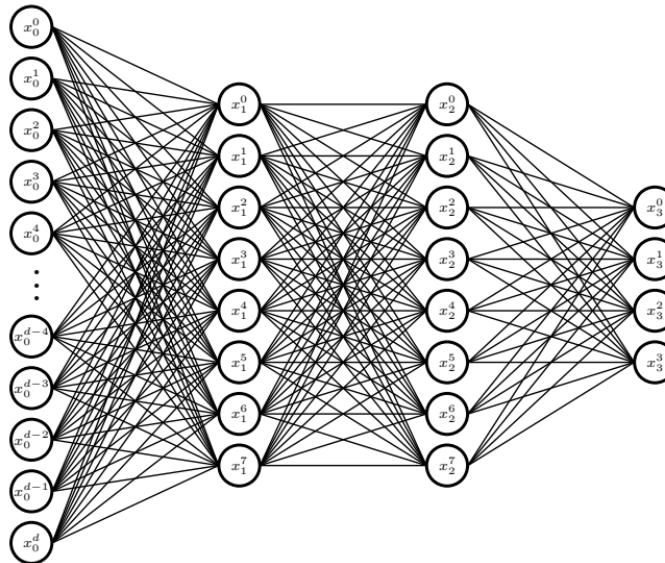


Let's see what the interconnected nodes, the layers and the transformations are.

# NEURAL NETWORKS FUNDAMENTALS

## NEURONS

If we consider that the Neural Network is a function  $f : \mathbb{R}^d \rightarrow \mathbb{R}^m$ :



A **Neuron** is a processing unit that receives input, performs a computation, and produces an output. Here, the inputs are  $x_{i-1}$  and the output is  $x_i^k$ .

# NEURAL NETWORKS FUNDAMENTALS

## NEURONS

For example, with an image dataset, the image can be flattened:

0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.99	0.91	0.02	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.99	0.45	0.18	0.66	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.99	0.07	0.00	0.00	0.99	0.00	0.00	0.00	0.00
0.00	0.00	0.30	0.44	0.00	0.00	0.00	0.99	0.00	0.00	0.00	0.00
0.00	0.00	0.33	0.00	0.00	0.00	0.99	0.00	0.00	0.00	0.00	0.00
0.00	0.00	0.33	0.99	0.99	0.77	0.00	0.00	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

$$\in [0, 1]^{d/2 \times d/2}$$

$$x_0 = [0.00, 0.00, \dots, 0.00, 0.99, 0.07 \dots, 0.00, 0.00] \in [0, 1]^d$$

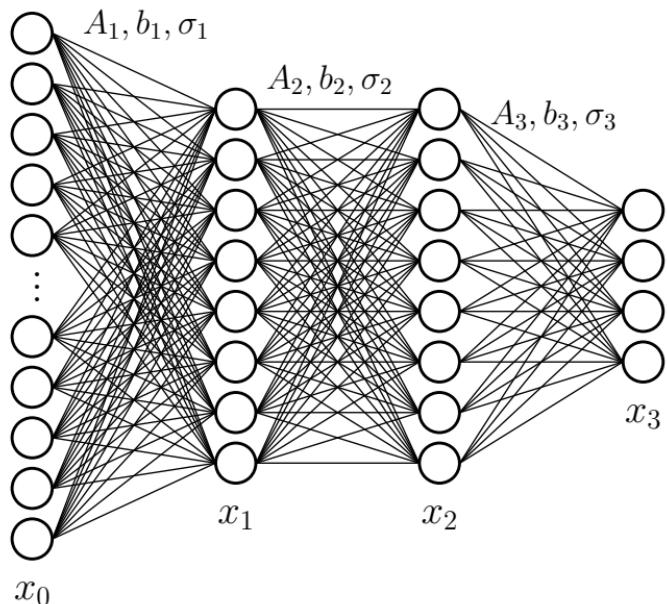
# NEURAL NETWORKS FUNDAMENTALS

## LAYERS

A layer  $i$  is defined by a matrix  $A_i \in \mathbb{R}^{k_{i-1} \times k_i}$ , a vector  $b_i \in \mathbb{R}^{k_i}$  and a nonlinear function  $\sigma_i : \mathbb{R} \mapsto \mathbb{R}$ . The transformation made by a layer is:

$$x_i = \sigma_i(A_i x_{i-1} + b_i).$$

The non-linear function  $\sigma_i$  the **activation function**.



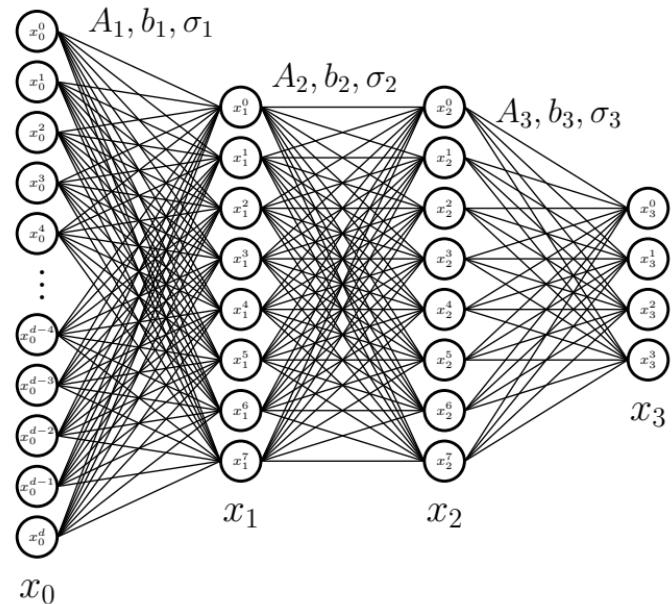
# NEURAL NETWORKS FUNDAMENTALS

## LAYERS

A layer  $i$  is defined as a matrix  $A_i \in \mathbb{R}^{k_{i-1} \times k_i}$ , a vector  $b_i \in \mathbb{R}^{k_i}$  and a nonlinear function  $\sigma_i : \mathbb{R} \mapsto \mathbb{R}$ . The transformation made by a layer is:

$$x_i^k = \sigma_i \left( \sum_{l=1}^{k_i} [A_i]_{l,k} x_{i-1} + [b_i]_k \right).$$

The non-linear function  $\sigma_i$  the activation function.



# NEURAL NETWORKS FUNDAMENTALS

## ACTIVATION FUNCTIONS

The activation functions play a crucial role in the implementation of deep neural networks, as they allow them to approximate any continuous function, as stated by the Universal Approximation Theorem. We can list some activation function that are commonly used :

- ▶ Linear
- ▶ Sigmoid
- ▶ Hyperbolic Tangent
- ▶ Rectified Linear Unit (ReLU)
- ▶ Leaky Rectified Linear Unit (Leaky ReLU)
- ▶ Exponential Linear Unit (ELU)
- ▶ Sigmoid-Weighted Linear Unit (Swish)
- ▶ Softmax

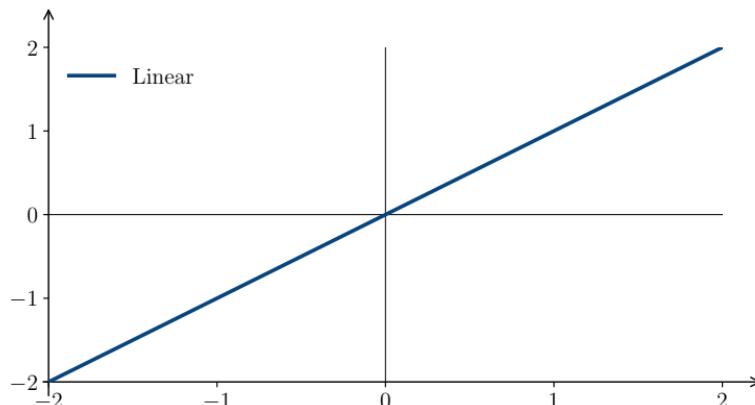
# NEURAL NETWORKS FUNDAMENTALS

## LINEAR

- ▶ Linear activation Function:

$$\sigma(x) = x$$

- ▶ Final activation
- ▶ Use case : Regression



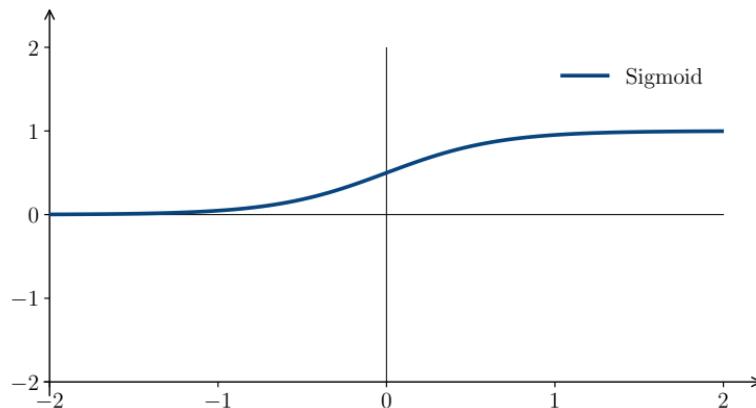
# NEURAL NETWORKS FUNDAMENTALS

## SIGMOID

- ▶ Sigmoid Function:

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

- ▶ Final activation
- ▶ Use case : Classification



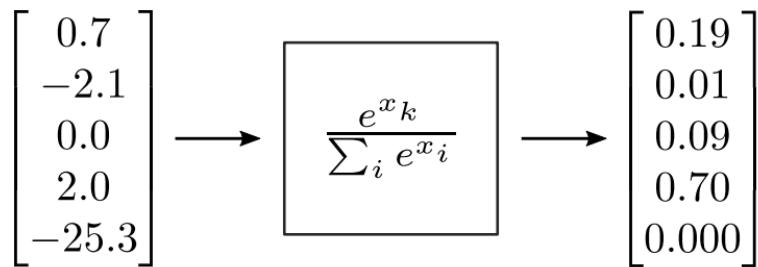
# NEURAL NETWORKS FUNDAMENTALS

## SOFTMAX

- ▶ Softmax Function:

$$\sigma(x_k) = \frac{e^{x_k}}{\sum_{i=1}^k e^{x_i}}$$

- ▶ Final activation
- ▶ Use case : Multi-class Classification



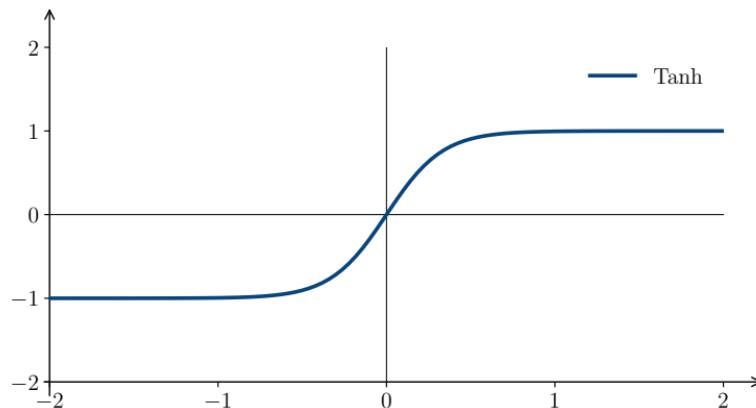
# NEURAL NETWORKS FUNDAMENTALS

## HYPERBOLIC TANGENT

- ▶ Hyperbolic Tangent

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

- ▶ Final activation
- ▶ Use case : Generative task



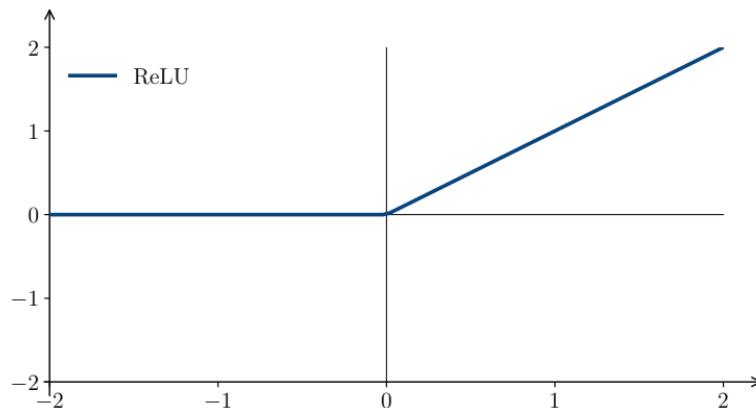
# NEURAL NETWORKS FUNDAMENTALS

## ReLU

- ▶ Rectified Linear Unit (ReLU):

$$\sigma(x) = \max\{0, x\}$$

- ▶ Intermediate activation



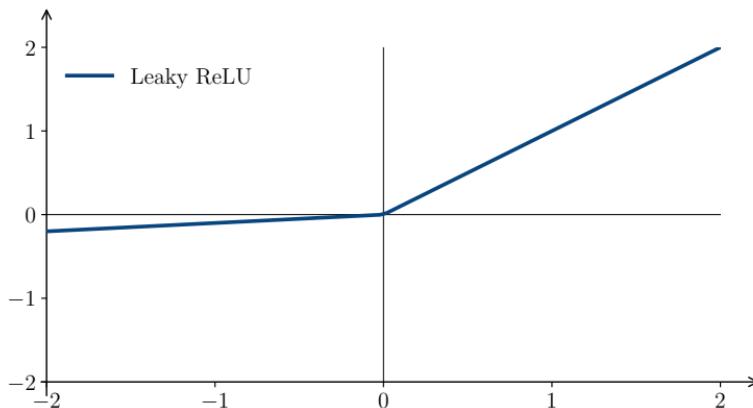
# NEURAL NETWORKS FUNDAMENTALS

## LEAKY RELU

- ▶ Leaky Rectified Linear Unit (Leaky ReLU):

$$\sigma(x) = \max\{\alpha x, x\}$$

- ▶ Intermediate activation



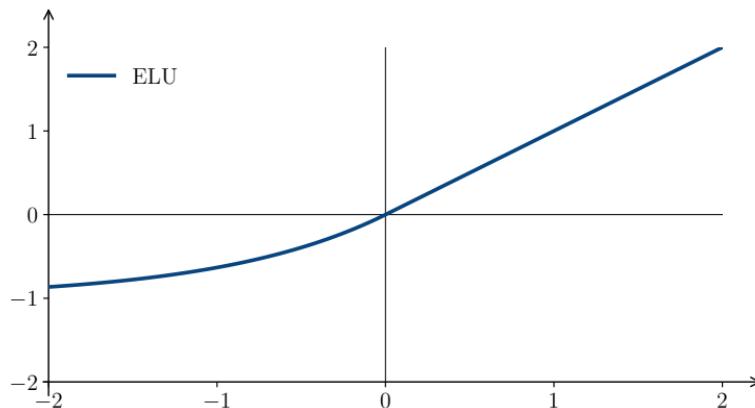
# NEURAL NETWORKS FUNDAMENTALS

## ELU

- ▶ Exponential Linear Unit (ELU):

$$\sigma(x) = \begin{cases} \alpha(e^x - 1) & \text{if } x < 0, \\ x & \text{if } x \geq 0. \end{cases}$$

- ▶ Intermediate activation



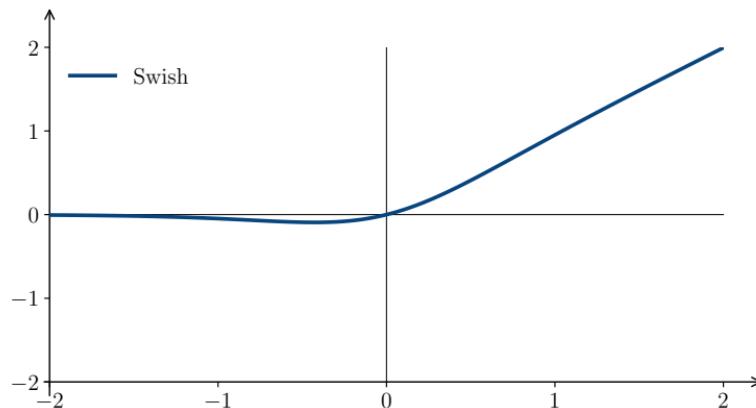
# NEURAL NETWORKS FUNDAMENTALS

## SWISH

- ▶ Sigmoid-Weighted Linear Unit (Swish):

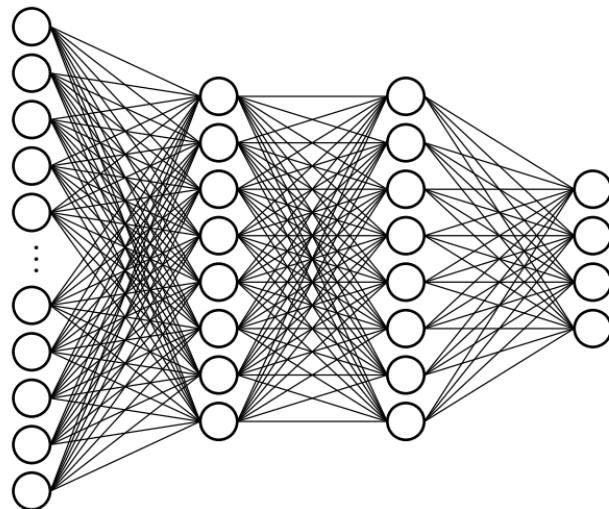
$$\sigma(x) = \frac{x}{1 + e^{-x}}$$

- ▶ Intermediate activation



## THE MULTI-LAYER PERCEPTRON (MLP)

Having discussed the structure of a neural network, we will proceed to examine the process of training a model for a specific task. As an illustration, we will consider the example of a Multilayer Perceptron. The two intermediate activation functions are ReLUs and the final activation is a softmax to perform multi-class classification on MNIST. We will consider only 4 classes.

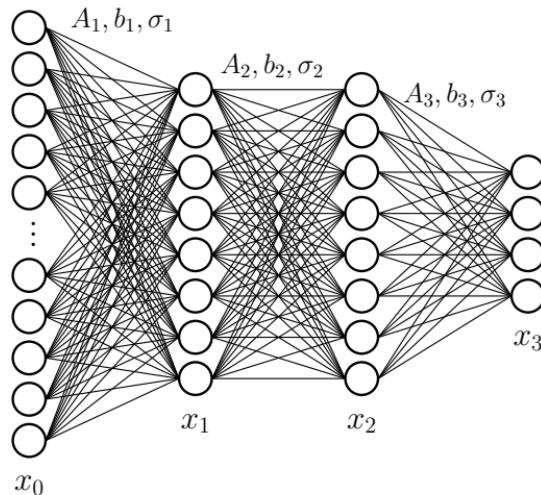


# THE MULTI-LAYER PERCEPTRON (MLP)

## THE FIRST DEEP LEARNING MODEL

To introduce the training process, we will consider a 3 layers MLP trained to minimise a loss  $\mathcal{L}$  over a given dataset  $\mathcal{D}$ . The model  $f_\theta$  is parameterised by a vector  $\theta = \{A_1, A_2, A_3, b_1, b_2, b_3\}$ :

$$\theta^* = \arg \min_{\theta} \mathcal{L}(\theta, \mathcal{D})$$



# THE MULTI-LAYER PERCEPTRON (MLP)

## STOCHASTIC GRADIENT DESCENT

Stochastic gradient descent (SGD) is widely used in deep learning instead of traditional gradient descent due to its efficiency and faster convergence rate. SGD updates the model parameters after computing the gradient of the loss function with respect to each parameter using only a single randomly selected sample. This leads to a faster convergence rate and improved optimization compared to traditional gradient descent, which uses the entire training dataset to compute the gradient at each iteration.

$$\theta^* = \arg \min_{\theta} \mathcal{L}(\theta, \mathcal{D}) = \arg \min_{\theta} \mathbb{E}_{x \sim \mathcal{D}} [l(x, f_{\theta}(x))]$$

# THE MULTI-LAYER PERCEPTRON (MLP)

## STOCHASTIC GRADIENT DESCENT

Theoretically the algorithm is the following:

**Require:** Given a loss function  $l$ , a dataset  $\mathcal{D} = \{x_1, x_2, \dots, x_N\}$  and a learning rate  $\lambda$

- 1: Initialize parameters  $\theta$
- 2: **while**  $\theta$  has not converged **do**
- 3:   **for**  $i = 1$  to  $N$  **do**
- 4:     Randomly select  $x_i$  from the dataset
- 5:     Compute gradient of the loss with respect to  $\theta$ :  $\nabla_{\theta} l(x_i, f_{\theta}(x_i))$
- 6:     Update parameters  $\theta = \theta - \lambda \nabla_{\theta} l(x_i, f_{\theta}(x_i))$
- 7:   **end for**
- 8: **end while**
- 9: **return**  $\theta$

# THE MULTI-LAYER PERCEPTRON (MLP)

## SGD IN MINI-BATCH

In practice the algorithm is modified to use mini-batches of data instead of single samples. This is done to improve the stability of the optimization process and reduce the variance of the gradient estimates. The algorithm is as follows:

**Require:** Given a loss function  $l$ , a dataset  $\mathcal{D} = \{x_1, x_2, \dots, x_N\}$ , a learning rate  $\lambda$  and a batch size  $b$

- 1: Initialize parameters  $\theta$
- 2: Initialize the number of batches  $B = \lfloor \frac{N}{b} \rfloor$
- 3: **while**  $\theta$  has not converged **do**
- 4:   **for**  $i = 1$  to  $B$  **do**
- 5:     Randomly select a mini-batch of  $b$  samples from the dataset
- 6:     Compute gradient of the loss with respect to  $\theta$ :  $\frac{1}{B} \sum_{i=1}^B \nabla_\theta l(x_i, f_\theta(x_i))$
- 7:     Update parameters  $\theta = \theta - \lambda \frac{1}{B} \sum_{i=1}^B \nabla_\theta l(x_i, f_\theta(x_i))$
- 8:   **end for**
- 9: **end while**
- 10: **return**  $\theta$

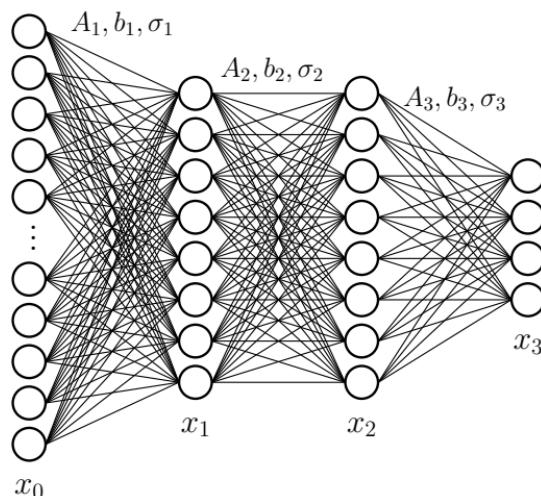
# THE MULTI-LAYER PERCEPTRON (MLP)

## BACK-PROPAGATION

At every step  $t$  of the gradient descent, setting a learning rate  $\lambda$ , the parameter  $\theta$  is updated as:

$$\theta_{t+1} = \theta_t - \lambda \nabla_{\theta} l(f(x_i), y_i)$$

But  $\theta = \{A_1, A_2, A_3, b_1, b_2, b_3\}$  and the gradient is computed with respect to each parameter.



## THE MULTI-LAYER PERCEPTRON (MLP)

### BACK-PROPAGATION

First we will consider a single data point  $x$ , the loss will depend on the output only:  $l(f(x))$ .

$f$  is a layered composed function. Let us focus on the last layer:

$$f(x) = x_3 = \sigma_3(A_3x_2 + b_3)$$

Therefore:

$$l(f(x)) = l(\sigma_3(A_3x_2 + b_3))$$

To minimise the loss, we have to act on  $A_3$ ,  $b_3$  and  $x_2$ .

# THE MULTI-LAYER PERCEPTRON (MLP)

## BACK-PROPAGATION

Let us look at the gradients with respect to  $A_3$ :

$$\begin{aligned}\frac{\partial l}{\partial A_3} &= \frac{\partial l}{\partial x_3} \frac{\partial x_3}{\partial A_3} = l'(x_3) \frac{\partial \sigma_3(A_3 x_2 + b_3)}{\partial A_3} = l'(x_3) \sigma'_3(A_3 x_2 + b_3) \frac{\partial [A_3 x_2 + b_3]}{\partial A_3} \\ &= \underbrace{l'(x_3)}_{\in \mathbb{R}} \underbrace{\sigma'_3(A_3 x_2 + b_3)}_{\in \mathbb{R}^{k_i \times 1}} \underbrace{x_2^T}_{\in \mathbb{R}^{1 \times k_{i-1}}}\end{aligned}$$

and therefore:

$$A_3 \leftarrow A_3 - \lambda l'(x_3) \sigma'_3(A_3 x_2 + b_3) x_2^T.$$

We need to keep in memory the latent values of  $x$ , i.e.  $x_2$ .

# THE MULTI-LAYER PERCEPTRON (MLP)

## BACK-PROPAGATION

Let us look at the gradients with respect to  $A_2$ :

$$\begin{aligned}\frac{\partial l}{\partial A_2} &= \frac{\partial l}{\partial x_2} \frac{\partial x_2}{\partial A_2} \\ &= \frac{\partial l}{\partial x_2} \frac{\partial \sigma_2(A_2x_1 + b_2)}{\partial A_2} \\ &= \frac{\partial l}{\partial x_2} \sigma'_2(A_2x_1 + b_2) \frac{\partial [A_2x_1 + b_2]}{\partial A_2} \\ &= \frac{\partial l}{\partial x_2} \sigma'_2(A_2x_1 + b_2) x_1^T\end{aligned}$$

which depends on  $\frac{\partial l}{\partial x_2}$ , we need to compute it.

# THE MULTI-LAYER PERCEPTRON (MLP)

## BACK-PROPAGATION

We have to compute the gradient with respect to  $x_2$ :

$$\begin{aligned}\frac{\partial l}{\partial x_2} &= \frac{\partial l}{\partial x_3} \frac{\partial x_3}{\partial x_2} = l'(x_3) \frac{\partial \sigma_3(A_3x_2 + b_3)}{\partial x_2} = l'(x_3) \frac{\partial [A_3x_2 + b_3]}{\partial x_2} \sigma'_3(A_3x_2 + b_3) \\ &= l'(x_3) A_3^T \sigma'_3(A_3x_2 + b_3)\end{aligned}$$

Therefore:

$$A_2 \leftarrow A_2 - \lambda \left[ l'(x_3) A_3^T \sigma'_3(A_3x_2 + b_3) \times \sigma'_2(A_2x_1 + b_2) x_1^T \right]$$

The update of  $A_2$  depends on  $l'(x_3)$ ,

## BACK-PROPAGATION

We have to compute the gradient with respect to  $A_1$ :

$$\begin{aligned}\frac{\partial l}{\partial A_1} &= \frac{\partial l}{\partial x_1} \frac{\partial x_1}{\partial A_1} \\ &= \frac{\partial l}{\partial x_1} \frac{\partial \sigma_1(A_1 x_0 + b_1)}{\partial A_1} \\ &= \frac{\partial l}{\partial x_1} \sigma'_1(A_1 x_0 + b_1) x_0^T,\end{aligned}$$

which depends on  $\frac{\partial l}{\partial x_1}$ , we need to compute it.

## BACK-PROPAGATION

Let us compute the gradient with respect to  $x_1$ :

$$\begin{aligned}\frac{\partial l}{\partial x_1} &= \frac{\partial l}{\partial x_2} \frac{\partial x_2}{\partial x_1} = \frac{\partial l}{\partial x_2} \frac{\partial \sigma_2 (A_2 x_1 + b_2)}{\partial x_1} = \frac{\partial l}{\partial x_2} \frac{\partial [A_2 x_1 + b_2]}{\partial x_1} \sigma'_2 (A_2 x_1 + b_2) \\ &= \frac{\partial l}{\partial x_2} A_2^T \sigma'_2 (A_2 x_1 + b_2)\end{aligned}$$

Therefore:

$$A_1 \leftarrow A_1 - \lambda \left[ l'(x_3) A_3^T \sigma'_3 (A_3 x_2 + b_3) A_2^T \sigma'_2 (A_2 x_1 + b_2) \times \sigma'_1 (A_1 x_0 + b_1) x_0^T \right]$$

## BACK-PROPAGATION

In other words, the update on the weights is:

$$A_3 \leftarrow A_3 - \lambda l'(x_3) \sigma'_3 (A_3 x_2 + b_3) x_2^T$$

$$A_2 \leftarrow A_2 - \lambda \left[ l'(x_3) A_3^T \sigma'_3 (A_3 x_2 + b_3) \times \sigma'_2 (A_2 x_1 + b_2) x_1^T \right]$$

$$A_1 \leftarrow A_1 - \lambda \left[ l'(x_3) A_3^T \sigma'_3 (A_3 x_2 + b_3) A_2^T \sigma'_2 (A_2 x_1 + b_2) \times \sigma'_1 (A_1 x_0 + b_1) x_0^T \right]$$

# THE MULTI-LAYER PERCEPTRON (MLP)

## BACK-PROPAGATION

If we look at the update of the different biases, we can easily compute the different gradient and see the updates. First, let us compute the gradient with respect to  $b_3$ :

$$\begin{aligned}\frac{\partial l}{\partial b_3} &= \frac{\partial l}{\partial x_3} \frac{\partial x_3}{\partial b_3} \\&= l'(x_3) \frac{\partial \sigma_3(A_3x_2 + b_3)}{\partial b_3} \\&= l'(x_3) \sigma'_3(A_3x_2 + b_3) \frac{\partial [A_3x_2 + b_3]}{\partial b_3} \\&= \underbrace{l'(x_3)}_{\in \mathbb{R}} \underbrace{\sigma'_3(A_3x_2 + b_3)}_{\in \mathbb{R}_i^{k \times 1}}\end{aligned}$$

And thus :

$$b_3 \leftarrow b_3 - \lambda l'(x_3) \sigma'(A_3x_2 + b_3)$$

# THE MULTI-LAYER PERCEPTRON (MLP)

## BACK-PROPAGATION

Let's move on the second layer:

$$\begin{aligned}\frac{\partial l}{\partial b_2} &= \frac{\partial l}{\partial x_2} \frac{\partial x_2}{\partial b_2} \\ &= \frac{\partial l}{\partial x_2} \frac{\partial \sigma_2 (A_2 x_1 + b_2)}{\partial b_2} \\ &= \frac{\partial l}{\partial x_2} \sigma'_2 (A_2 x_1 + b_2)\end{aligned}$$

And thus :

$$b_2 \leftarrow b_2 - \lambda \frac{\partial l}{\partial x_2} \sigma' (A_2 x_1 + b_2)$$

We need to back-propagate the term  $\frac{\partial l}{\partial x_2}$  computed for the first layer.

# THE MULTI-LAYER PERCEPTRON (MLP)

## BACK-PROPAGATION

For the first layer:

$$\begin{aligned}\frac{\partial l}{\partial b_1} &= \frac{\partial l}{\partial x_1} \frac{\partial x_1}{\partial b_1} \\ &= \frac{\partial l}{\partial x_1} \frac{\partial \sigma_1(A_1x_0 + b_1)}{\partial b_1} \\ &= \frac{\partial l}{\partial x_1} \sigma'_1(A_1x_0 + b_1)\end{aligned}$$

And thus :

$$b_1 \leftarrow b_1 - \lambda \frac{\partial l}{\partial x_1} \sigma'(A_1x_0 + b_1)$$

We need to back-propagate the term  $\frac{\partial l}{\partial x_1}$  computed for the second layer which has been computed with  $\frac{\partial l}{\partial x_2}$  back-propagated from the first layer.

# THE MULTI-LAYER PERCEPTRON (MLP)

## BACK-PROPAGATION

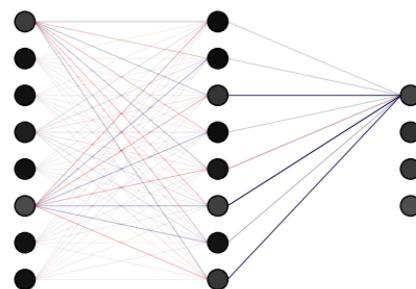
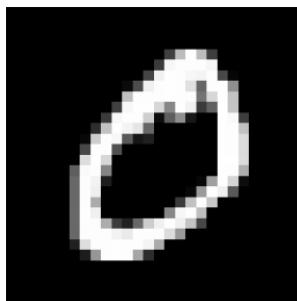
To update the weights, we need to compute the gradient of the loss with respect to the output of the network, and then **back-propagate** the gradient of the loss with respect to each activation, the  $\frac{\partial l}{\partial x_i}$ , through the network to compute the gradients with respect to the weights and biases of each layer.

# THE MULTI-LAYER PERCEPTRON (MLP)

## LAST LAYER

We can plot the current state of the network for a given input.

The red lines show positive values for  $A_i$ , the blue lines represent negative values for  $A_i$ . The level of transparency is proportional to the previous neurons.

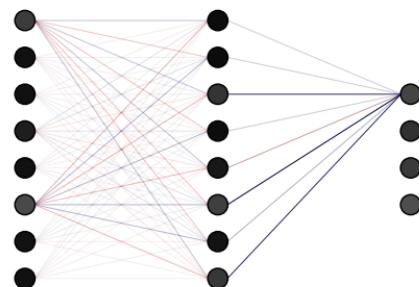
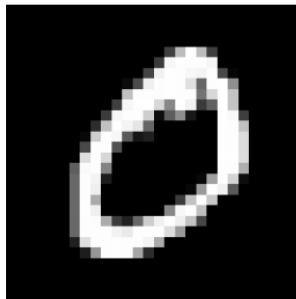


$$x_3^1 = \sigma_3 \left( A_3^{1,1} x_2^1 + A_3^{1,2} x_2^2 + \dots + A_3^{1,8} x_2^8 \right)$$

# THE MULTI-LAYER PERCEPTRON (MLP)

## BACK-PROPAGATION

Iteratively, the neural networks improves its performance.

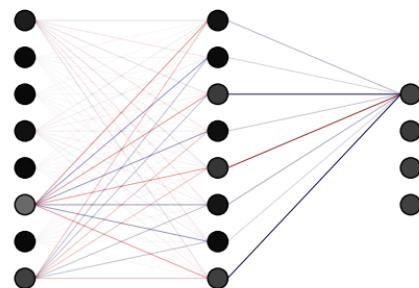
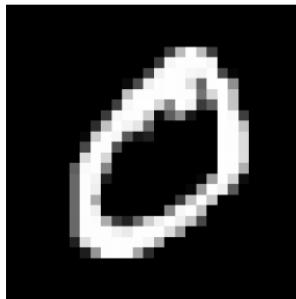


$$x_3^1 = \sigma_3 \left( A_3^{1,1} x_2^1 + A_3^{1,2} x_2^2 + \cdots + A_3^{1,8} x_2^8 \right)$$

# THE MULTI-LAYER PERCEPTRON (MLP)

## BACK-PROPAGATION

Iteratively, the neural networks improves its performance.

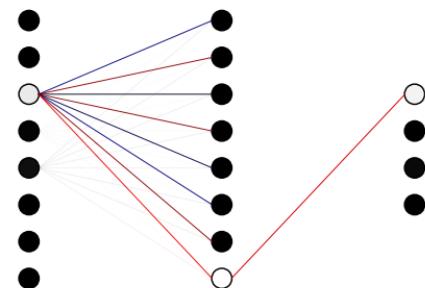
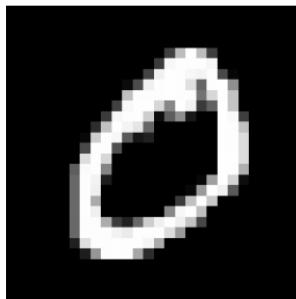


$$x_3^1 = \sigma_3 \left( A_3^{1,1} x_2^1 + A_3^{1,2} x_2^2 + \cdots + A_3^{1,8} x_2^8 \right)$$

# THE MULTI-LAYER PERCEPTRON (MLP)

## BACK-PROPAGATION

Iteratively, the neural networks improves its performance.

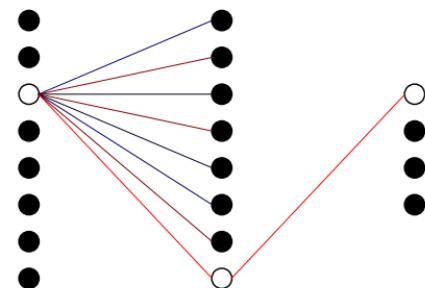
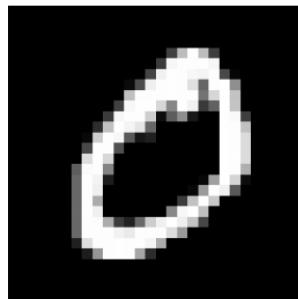


$$x_3^1 = \sigma_3 \left( A_3^{1,1} x_2^1 + A_3^{1,2} x_2^2 + \cdots + A_3^{1,8} x_2^8 \right)$$

# THE MULTI-LAYER PERCEPTRON (MLP)

## BACK-PROPAGATION

Iteratively, the neural networks improves its performance.



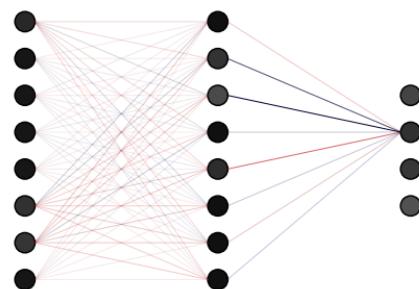
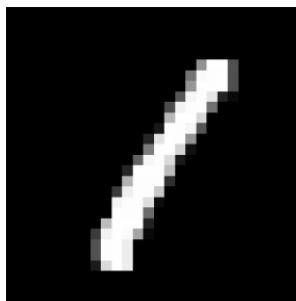
$$x_3^1 = \sigma_3 \left( A_3^{1,1} x_2^1 + A_3^{1,2} x_2^2 + \cdots + A_3^{1,8} x_2^8 \right)$$

# THE MULTI-LAYER PERCEPTRON (MLP)

## LAST LAYER

We can plot the current state of the network for a given input.

Red lines show positive values of  $A_i$ , Blue lines represent negative values of  $A_i$ . The level of transparency is proportional to the previous neurons.

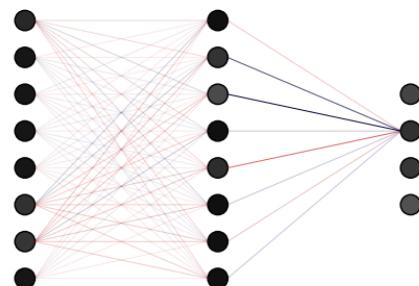
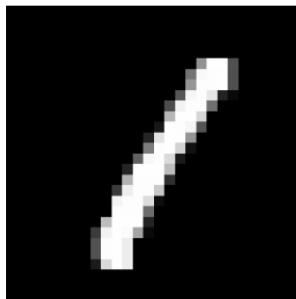


$$x_3^2 = \sigma_3 \left( A_3^{2,1}x_2^1 + A_3^{2,2}x_2^2 + \dots + A_3^{2,8}x_2^8 \right)$$

# THE MULTI-LAYER PERCEPTRON (MLP)

## BACK-PROPAGATION

Iteratively, the neural networks improves its performance.

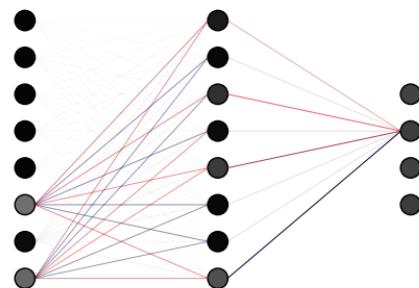
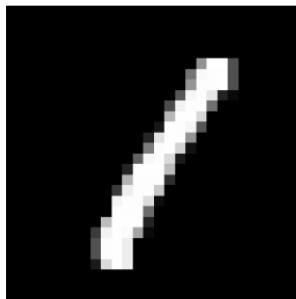


$$x_3^2 = \sigma_3 \left( A_3^{2,1} x_2^1 + A_3^{2,2} x_2^2 + \cdots + A_3^{2,8} x_2^8 \right)$$

# THE MULTI-LAYER PERCEPTRON (MLP)

## BACK-PROPAGATION

Iteratively, the neural networks improves its performance.

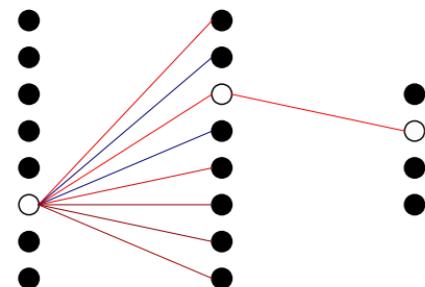
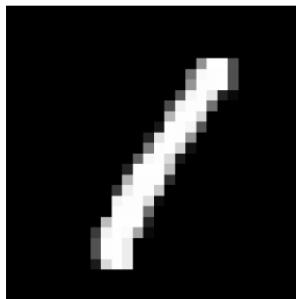


$$x_3^2 = \sigma_3 \left( A_3^{2,1} x_2^1 + A_3^{2,2} x_2^2 + \cdots + A_3^{2,8} x_2^8 \right)$$

# THE MULTI-LAYER PERCEPTRON (MLP)

## BACK-PROPAGATION

Iteratively, the neural networks improves its performance.

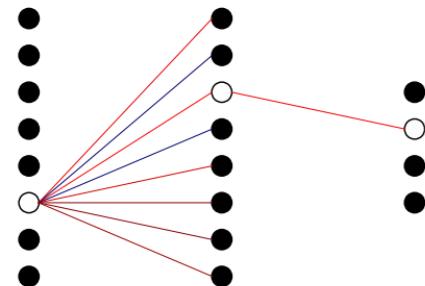
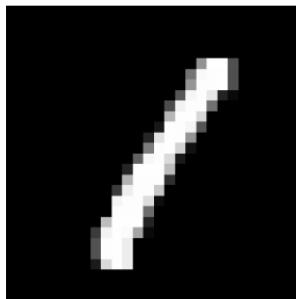


$$x_3^2 = \sigma_3 \left( A_3^{2,1} x_2^1 + A_3^{2,2} x_2^2 + \cdots + A_3^{2,8} x_2^8 \right)$$

# THE MULTI-LAYER PERCEPTRON (MLP)

## BACK-PROPAGATION

Iteratively, the neural networks improves its performance.

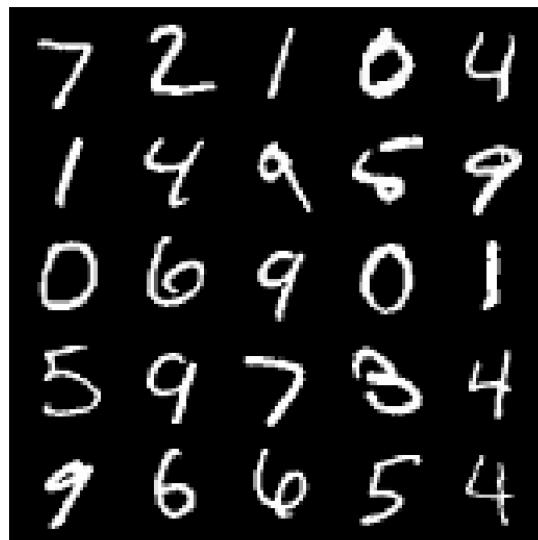


$$x_3^2 = \sigma_3 \left( A_3^{2,1} x_2^1 + A_3^{2,2} x_2^2 + \cdots + A_3^{2,8} x_2^8 \right)$$

## THE MULTI-LAYER PERCEPTRON (MLP)

EXAMPLE : IMAGE CLASSIFICATION OF HANDWRITTEN DIGITS FROM A TO Z

Having discussed the theory behind Artificial Neural Networks and the training process, we will now proceed to demonstrate a comprehensive end-to-end example of image classification on MNIST.



# THE MULTI-LAYER PERCEPTRON (MLP)

EXAMPLE : IMAGE CLASSIFICATION OF HANDWRITTEN DIGITS FROM A TO Z

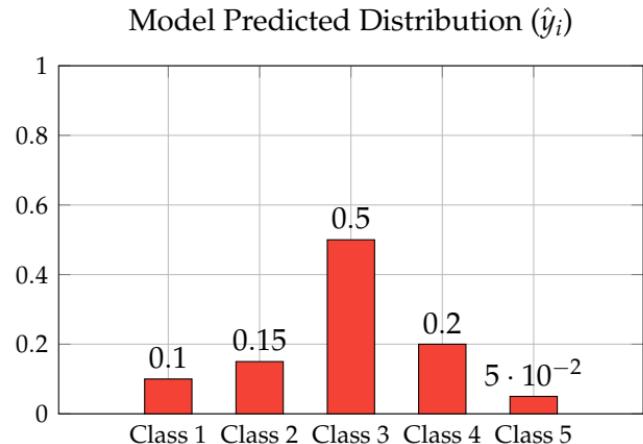
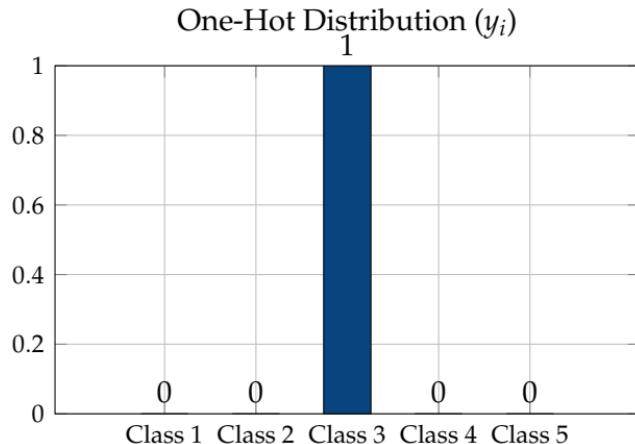
- ▶ Input shape :  $1 \times 28 \times 28$ .
- ▶ Number of Classes : 10.
- ▶ Number of training samples  $(x, y)$ : 60000.
- ▶ Number of evaluating samples: 10000.
- ▶ Loss : cross-entropy

$$L(\hat{y}, y) = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^K y_{ij} \log(\hat{y}_{ij})$$

where :

- $\hat{y} \in \mathbb{R}^{N \times K}$  is the predicted probability distribution over  $K$  classes for  $N$  samples,
- $y \in \{0, 1\}^{N \times K}$  is the ground-truth one-hot encoded label matrix,

## RECAP ON THE CROSS-ENTROPY LOSS



The cross-entropy loss for one sample is:

$$l(\hat{y}_i, y_i) = - \sum_{j=1}^K y_{ij} \log(\hat{y}_{ij}).$$

# THE MULTI-LAYER PERCEPTRON (MLP)

EXAMPLE : IMAGE CLASSIFICATION OF HANDWRITTEN DIGITS FROM A TO Z

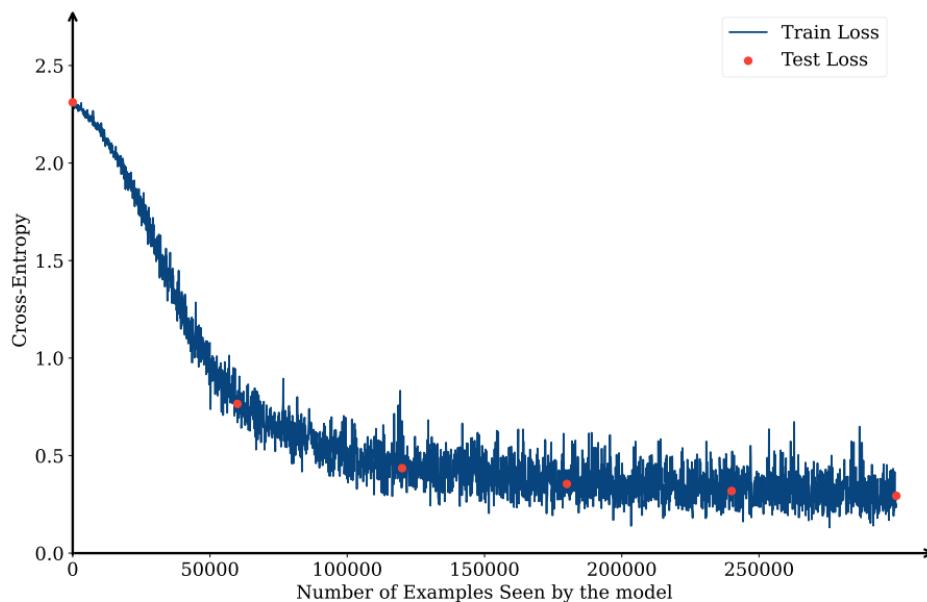
We build a 3 layers network.

- ▶ Batch size : 64
- ▶ Learning rate : 0.01
- ▶ Intermediate activation : ReLU
- ▶ Final activation : Softmax
- ▶ Number of epochs : 12
- ▶ Number of trained parameters: 52.6k



# THE MULTI-LAYER PERCEPTRON (MLP)

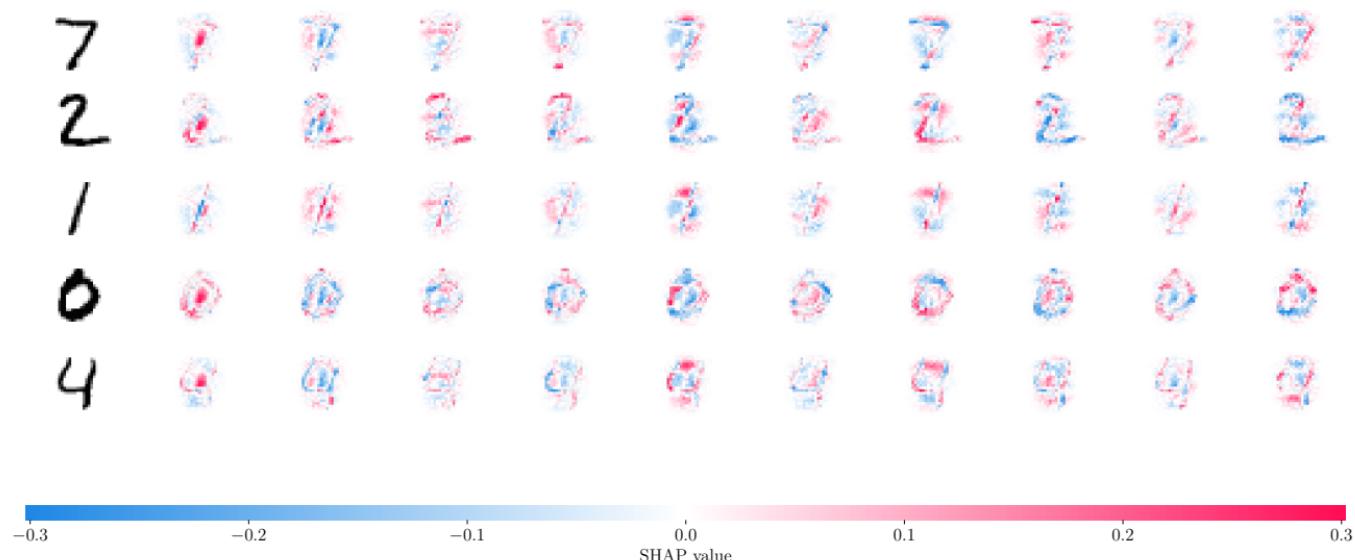
EXAMPLE : IMAGE CLASSIFICATION OF HANDWRITTEN DIGITS FROM A TO Z



# THE MULTI-LAYER PERCEPTRON (MLP)

EXAMPLE : IMAGE CLASSIFICATION OF HANDWRITTEN DIGITS FROM A TO Z

With an interpretation tool such as SHAP:



# TP1: THE MULTI-LAYER PERCEPTRON (MLP)

THE FIRST DEEP LEARNING MODEL

Link to the notebook (ipynb): [TP1.ipynb](#)

Link to the notebook (html): [TP1.html](#)

## Part II

# DEEP LEARNING IN ACTION: FROM NEURAL NETWORKS TO TRANSFORMER MODELS

Now that we have an understanding of the training procedure for Artificial Neural Networks, we shall examine several widely-utilized structures within the literature of Neural Networks, including Convolutional Neural Networks (CNN), Residual Networks (ResNet), Recurrent Neural Networks (RNN), and Transformers.

## CONVOLUTIONAL NEURAL NETWORKS

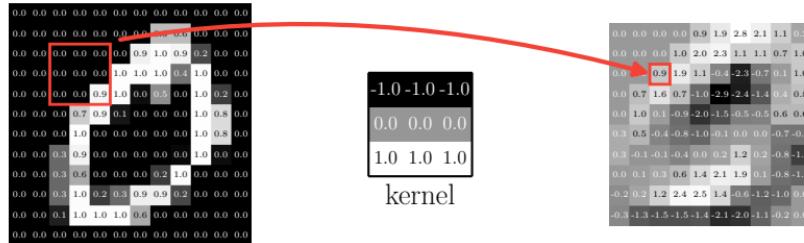
In the field of image processing, the Convolution Operators are widely considered as the most favoured approach. While it has been demonstrated that Dense blocks, or Linear blocks, are capable of accurately classifying images in the case of the MNIST dataset, the need for convolutional transformations arises when addressing wider and more intricate datasets.

# CONVOLUTIONAL NEURAL NETWORKS

## THE TWO DIMENSIONAL CONVOLUTION

A 2D convolution in a neural network context can be mathematically represented as a sliding window operation where a filter (also called kernel)  $w$  of size  $k \times k$  is applied to each  $k \times k$  sub-matrix of the input matrix  $x$ . The operation can be defined as the element-wise multiplication of the filter  $w$  and the sub-matrix followed by summing the results, i.e.

$$y_{i,j} = \sum_{m=1}^k \sum_{n=1}^k w_{m,n} \cdot x_{i+m,j+n}$$

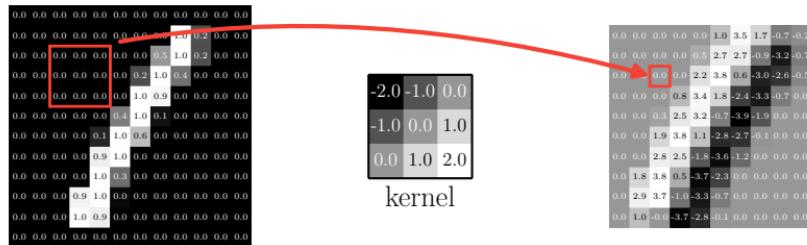


# CONVOLUTIONAL NEURAL NETWORKS

## THE TWO DIMENSIONAL CONVOLUTION

A 2D convolution in a neural network context can be mathematically represented as a sliding window operation where a filter (also called kernel)  $w$  of size  $k \times k$  is applied to each  $k \times k$  submatrix of the input matrix  $x$ . The operation can be defined as the element-wise multiplication of the filter  $w$  and the submatrix followed by summing the results, i.e.

$$y_{i,j} = \sum_{m=1}^k \sum_{n=1}^k w_{m,n} \cdot x_{i+m,j+n}$$



# CONVOLUTIONAL NEURAL NETWORKS

## KERNEL SIZE, PADDING AND STRIDE

In every Deep Learning library, the Conv2D block takes three parameters in argument:

- ▶ the Kernel's size,
- ▶ the Stride,
- ▶ the Padding.

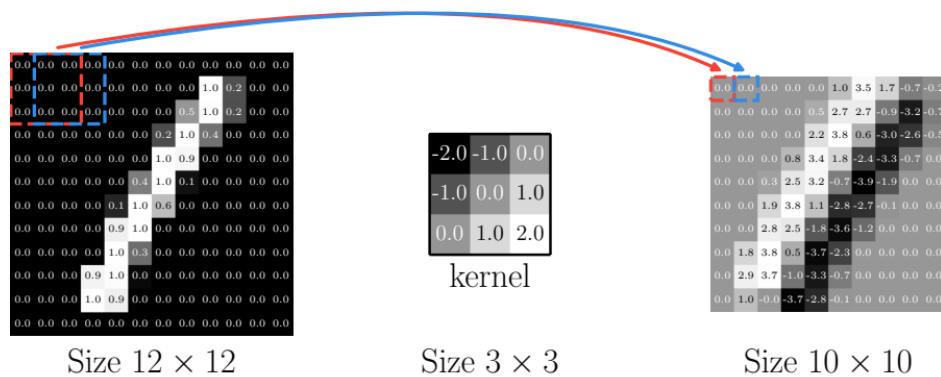
The size out the output is :

$$d_{\text{out}} = \frac{d_{\text{in}} + 2 \times \text{Padding} - \text{KernelSize}}{\text{Stride}} + 1$$

# CONVOLUTIONAL NEURAL NETWORKS

## KERNEL SIZE, PADDING AND STRIDE

Kernel size: 3, Padding: 0, Stride: 1

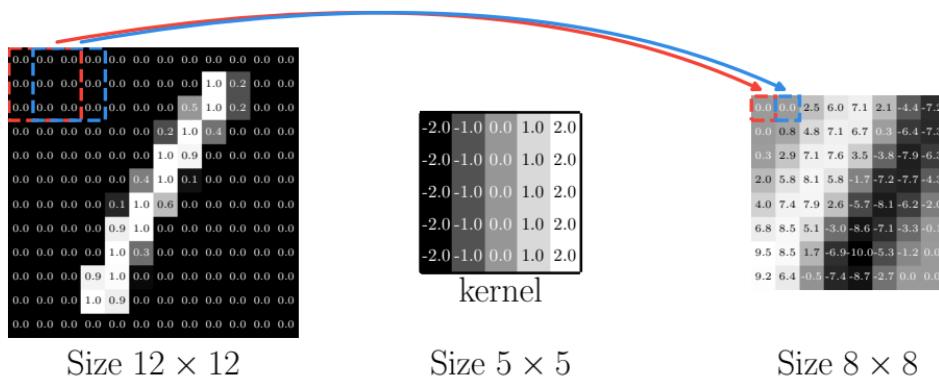


$$d_{\text{out}} = \frac{d_{\text{in}} + 2 \times \text{Padding} - \text{KernelSize}}{\text{Stride}} + 1$$

# CONVOLUTIONAL NEURAL NETWORKS

## KERNEL SIZE, PADDING AND STRIDE

Kernel size: 5, Padding: 0, Stride: 1

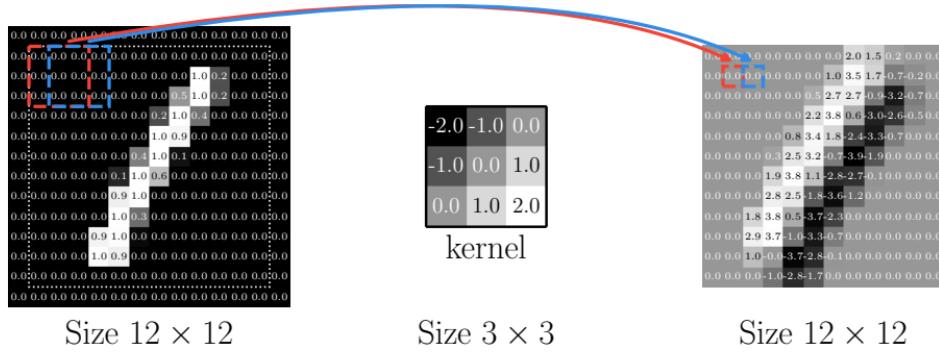


$$d_{\text{out}} = \frac{d_{\text{in}} + 2 \times \text{Padding} - \text{KernelSize}}{\text{Stride}} + 1$$

# CONVOLUTIONAL NEURAL NETWORKS

## KERNEL SIZE, PADDING AND STRIDE

Kernel size: 3, Padding: 1, Stride: 1. Padding mode can be 'zeros', 'reflect', 'replicate' or 'circular'.

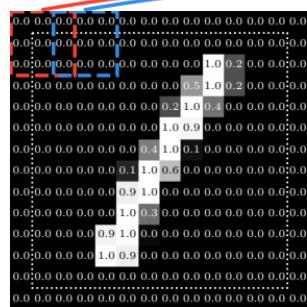


$$d_{\text{out}} = \frac{d_{\text{in}} + 2 \times \text{Padding} - \text{KernelSize}}{\text{Stride}} + 1$$

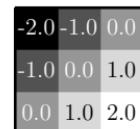
# CONVOLUTIONAL NEURAL NETWORKS

## KERNEL SIZE, PADDING AND STRIDE

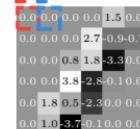
Kernel size: 3, Padding: 1, Stride: 2



Size  $12 \times 12$



kernel



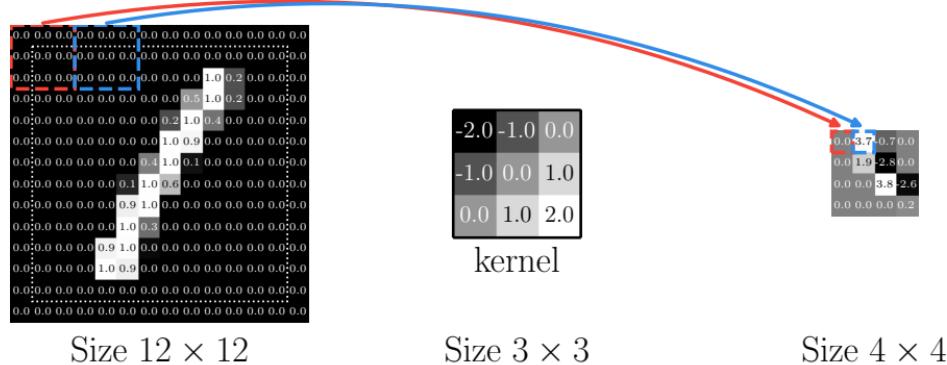
Size  $6 \times 6$

$$d_{\text{out}} = \frac{d_{\text{in}} + 2 \times \text{Padding} - \text{KernelSize}}{\text{Stride}} + 1$$

# CONVOLUTIONAL NEURAL NETWORKS

## KERNEL SIZE, PADDING AND STRIDE

Kernel size: 3, Padding: 1, Stride: 3

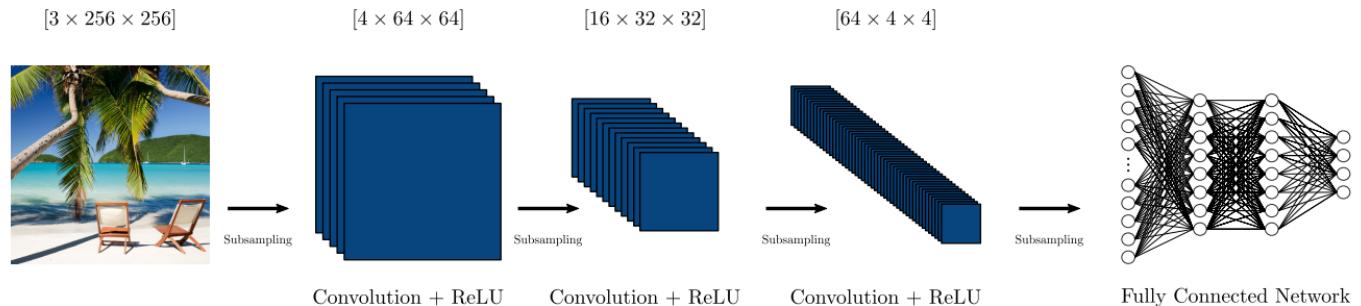


$$d_{\text{out}} = \frac{d_{\text{in}} + 2 \times \text{Padding} - \text{KernelSize}}{\text{Stride}} + 1$$

# CONVOLUTIONAL NEURAL NETWORKS

## CNN : CONVOLUTIONAL IN A NETWORK NETWORKS

We can represent a CNN as under this form:

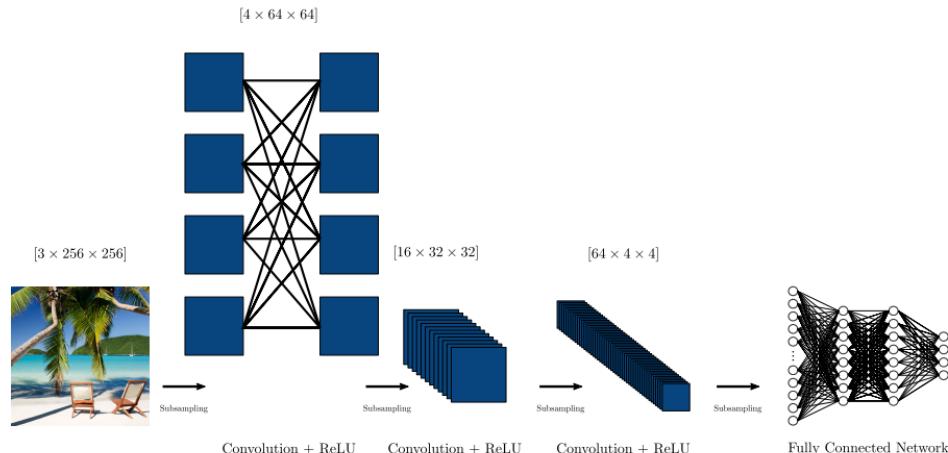


# CONVOLUTIONAL NEURAL NETWORKS

## CNN : CONVOLUTIONAL IN A NETWORK NETWORKS

Usually, the output of a convolutional block is linear combination of the Convolutional output of every previous channels and a bias:

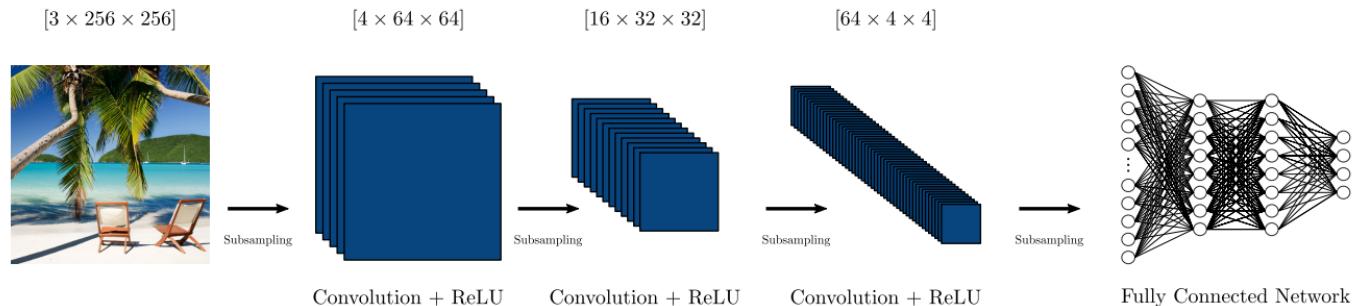
$$\text{out}_{i,j}(c_{\text{out}}) = \text{bias}(c_{\text{out}}) + \sum_{k=0}^{|c_{\text{in}}|-1} \text{Conv}(\text{input}(k), \text{kernel}_k)_{i,j}$$



# CONVOLUTIONAL NEURAL NETWORKS

## CNN : CONVOLUTIONAL IN A NETWORK NETWORKS

In practice, we split the image into multiple channels : the three channels RGB to begin with. Then we apply convolutional operation on different scales and then we use a fully connected tail. To change the scale we can use different sub-sampling : Max pooling, Average pooling or Invertible pooling.



# CONVOLUTIONAL NEURAL NETWORKS

## MAX POOLING

Max Pooling takes the maximum within a given sized sub-matrix. In practice, the matrix is size  $2 \times 2$  in order to reduce the dimension by 4 and doubling the scale.

0.2	1.0	0.3	0.8	0.1	0.8	0.6	0.9
0.7	0.3	0.3	0.5	0.4	0.3	0.3	0.4
0.2	0.6	0.7	0.9	0.9	0.1	0.3	0.5
0.8	0.7	0.1	0.3	0.3	0.6	0.9	0.5
0.7	0.1	0.1	0.2	0.8	0.4	0.9	0.7
0.9	0.1	0.8	0.9	0.6	0.3	0.1	0.9
0.8	0.7	0.2	0.7	0.0	0.6	0.9	0.5
0.9	0.5	0.1	0.2	0.0	0.1	0.1	0.4

Max Pooling  
Subsampling

1.0	0.8	0.8	0.9
0.8	0.9	0.9	0.9
0.9	0.9	0.8	0.9
0.9	0.7	0.6	0.9

# CONVOLUTIONAL NEURAL NETWORKS

## AVERAGE POOLING

The Average pooling takes the average value within the sub-matrix.

0.2	1.0	0.3	0.8	0.1	0.8	0.6	0.9
0.7	0.3	0.3	0.5	0.4	0.3	0.3	0.4
0.2	0.6	0.7	0.9	0.9	0.1	0.3	0.5
0.8	0.7	0.1	0.3	0.3	0.6	0.9	0.5
0.7	0.1	0.1	0.2	0.8	0.4	0.9	0.7
0.9	0.1	0.8	0.9	0.6	0.3	0.1	0.9
0.8	0.7	0.2	0.7	0.0	0.6	0.9	0.5
0.9	0.5	0.1	0.2	0.0	0.1	0.1	0.4

Average Pooling  
→  
Subsampling

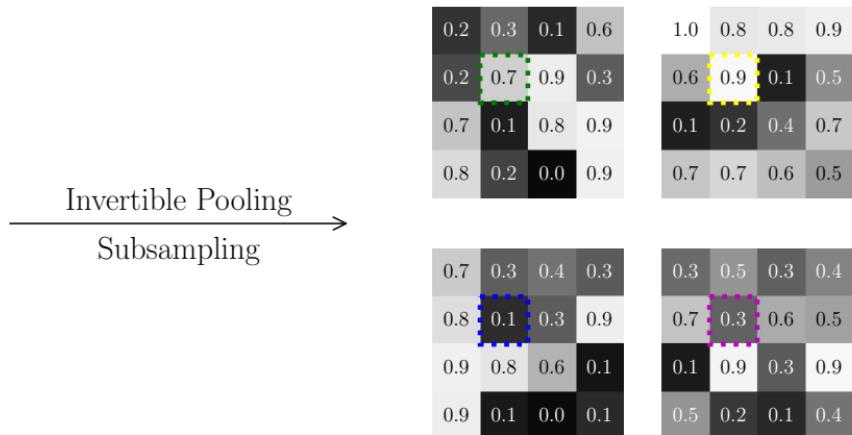
0.5	0.5	0.4	0.6
0.6	0.5	0.5	0.5
0.4	0.5	0.5	0.6
0.7	0.3	0.2	0.5

# CONVOLUTIONAL NEURAL NETWORKS

## INVERTIBLE POOLING

For Invertible Networks, we can use Invertible Pooling, aka Squeeze. It preserves the information contained in the channels and keeps the dimension constant.

0.2	1.0	0.3	0.8	0.1	0.8	0.6	0.9
0.7	0.3	0.3	0.5	0.4	0.3	0.3	0.4
0.2	0.6	0.7	0.9	0.9	0.1	0.3	0.5
0.8	0.7	0.1	0.3	0.3	0.6	0.9	0.5
0.7	0.1	0.1	0.2	0.8	0.4	0.9	0.7
0.9	0.1	0.8	0.9	0.6	0.3	0.1	0.9
0.8	0.7	0.2	0.7	0.0	0.6	0.9	0.5
0.9	0.5	0.1	0.2	0.0	0.1	0.1	0.4



# CONVOLUTIONAL NEURAL NETWORKS

## CNN : CONVOLUTIONAL IN A NETWORK NETWORKS

Convolutional Neural Networks are more suitable for image processing compared to fully connected networks due to their ability to efficiently handle the spatial relationships between pixels in an image. This is achieved through the use of convolutional layers that apply filters to small portions of an image, rather than fully connected layers that process the entire image as a single vector. Additionally, the shared weights in convolutional layers allow for learning of hierarchical features, reducing the number of parameters in the network and increasing its ability to generalize to new images.

# CONVOLUTIONAL NEURAL NETWORKS

## CNN IN PRACTICE: CIFAR 10

- ▶ Input shape :  $3 \times 32 \times 32$ .
- ▶ Number of Classes : 10.
- ▶ Number of training samples ( $x, y$ ): 50000.
- ▶ Number of evaluating samples: 10000.



# CONVOLUTIONAL NEURAL NETWORKS

## CNN IN PRACTICE: CIFAR 10

We will compare three different models:

- ▶ Model 1 : Fully Connected Neural Network with 3.4 million parameters.
- ▶ Model 2 : CNN with 62 thousand parameters.
- ▶ Model 3 : Wider and longer CNN with 5.8 million parameters.

# CONVOLUTIONAL NEURAL NETWORKS

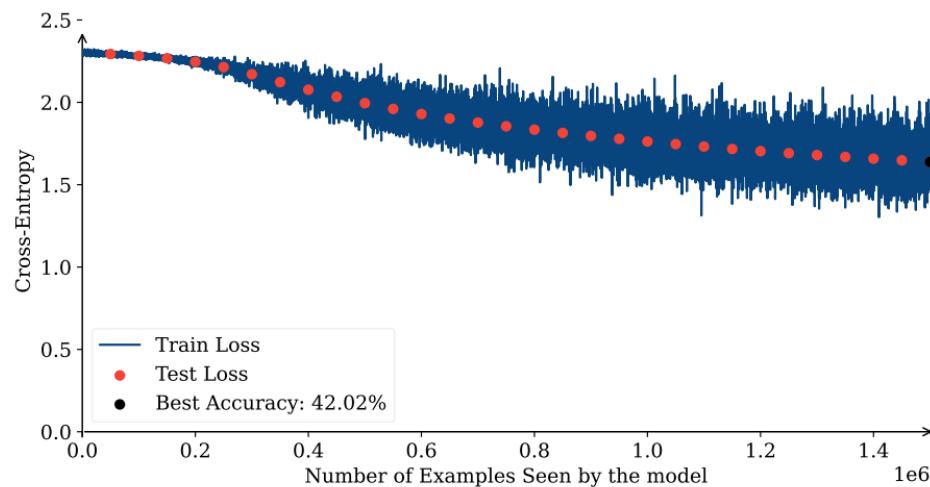
## MODEL 1 : FULLY CONNECTED NEURAL NETWORK

The Net is composed of 4 linear layers with ReLU activations:

- ▶ Linear  $3072 \mapsto 1024 + \text{ReLU}$
- ▶ Linear  $1024 \mapsto 256 + \text{ReLU}$
- ▶ Linear  $256 \mapsto 64 + \text{ReLU}$
- ▶ Linear  $64 \mapsto 10 + \text{SoftMax}$

# CONVOLUTIONAL NEURAL NETWORKS

## MODEL 1



# CONVOLUTIONAL NEURAL NETWORKS

## USING A BETTER OPTIMIZATION ALGORITHM ?

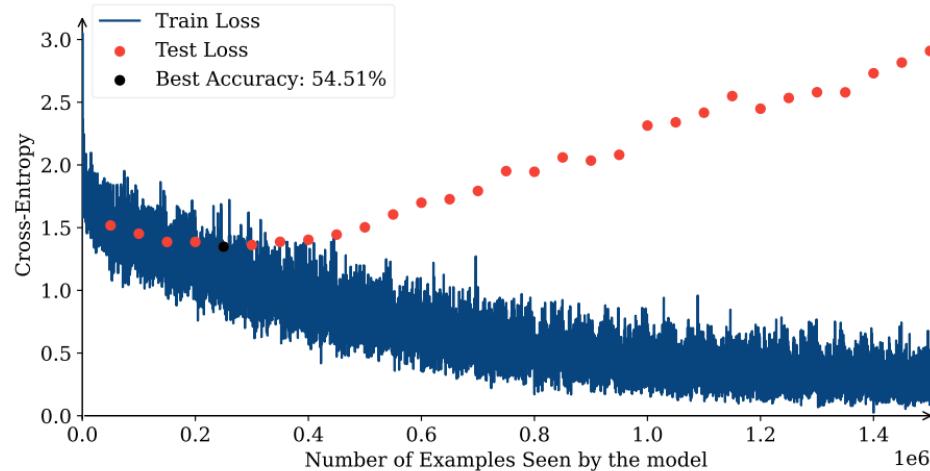
**Adam:** Adam is an optimization algorithm that can be used instead of the classical stochastic gradient descent procedure to update network weights iteratively based on training data. It combines the advantages of two other extensions of stochastic gradient descent, AdaGrad and RMSProp, and is designed to work well for a wide range of applications:

$$\begin{aligned}\theta_{t+1} &= \theta_t - \frac{\lambda}{\sqrt{\hat{v}_t + \varepsilon}} \hat{m}_t \\ \hat{m}_t &= \beta_1 m_{t-1} + (1 - \beta_1) \nabla_\theta l(\theta_t) \\ \hat{v}_t &= \beta_2 v_{t-1} + (1 - \beta_2) \nabla_\theta l(\theta_t)^2\end{aligned}$$

- ▶  $\lambda$  is the learning rate.
- ▶  $\hat{m}_t$  is the moving average of the gradient. It allows to smooth the gradient step.
- ▶ By dividing the gradient by the square root of the moving average of the squared gradient, we can adapt the learning rate for each parameter to improve the convergence speed.

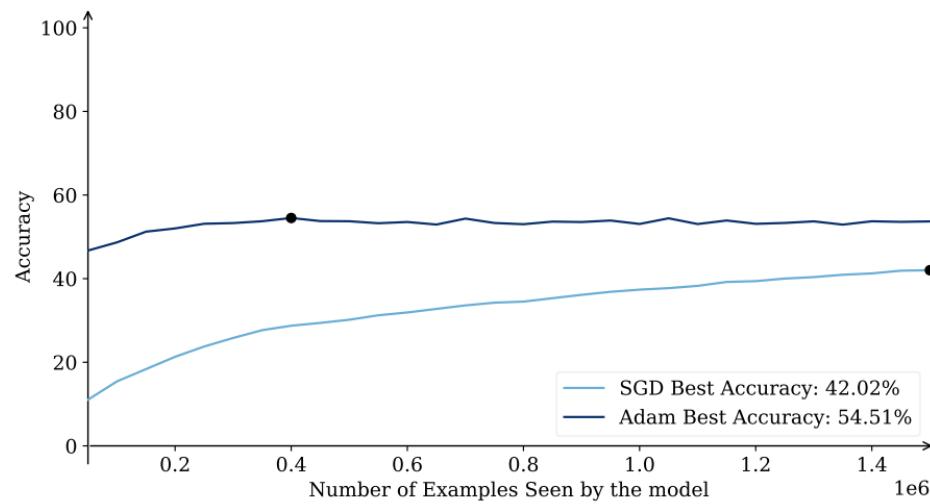
# CONVOLUTIONAL NEURAL NETWORKS

## MODEL 1



# CONVOLUTIONAL NEURAL NETWORKS

## MODEL 1



# CONVOLUTIONAL NEURAL NETWORKS

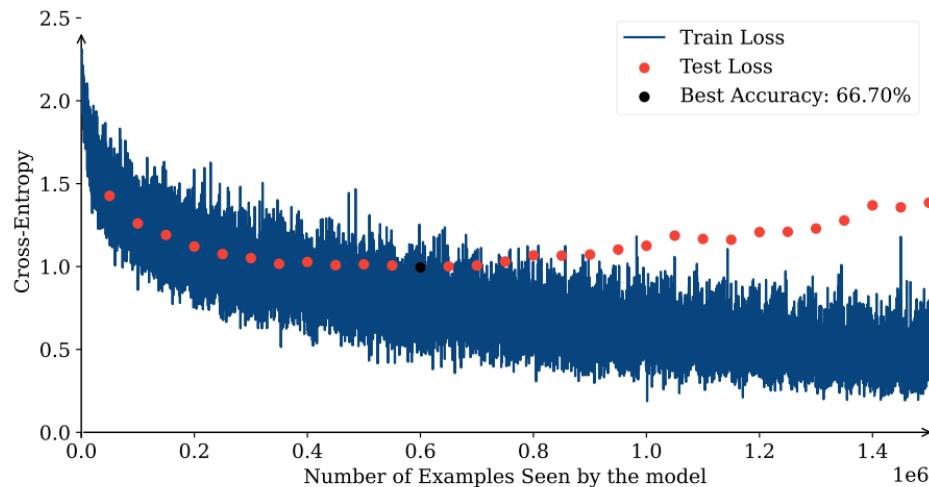
## MODEL 2

The Net is composed 2 convolutional layers and 2 linear layers:

- ▶ Conv  $3 \times 32 \times 32 \mapsto 6 \times 28 \times 28$  + ReLU
- ▶ Max Pooling  $6 \times 28 \times 28 \mapsto 6 \times 14 \times 14$
- ▶ Conv  $6 \times 14 \times 14 \mapsto 16 \times 10 \times 10$  + ReLU
- ▶ Max Pooling  $16 \times 10 \times 10 \mapsto 16 \times 5 \times 5$
- ▶ Linear  $400 \mapsto 120$  + ReLU
- ▶ Linear  $120 \mapsto 84$  + ReLU
- ▶ Linear  $84 \mapsto 10$  + SoftMax

# CONVOLUTIONAL NEURAL NETWORKS

## MODEL 2



# CONVOLUTIONAL NEURAL NETWORKS

## MODEL 3

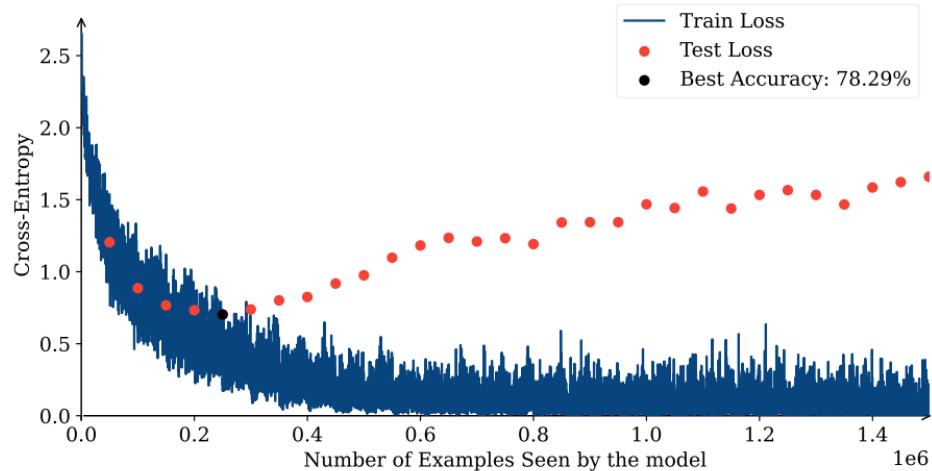
The Net is composed 6 convolutional layers and 3 linear layers:

- ▶ Conv  $3 \times 32 \times 32 \mapsto 32 \times 32 \times 32$  + BatchNorm2d + ReLU
- ▶ Conv  $32 \times 32 \times 32 \mapsto 64 \times 32 \times 32$  + ReLU
- ▶ Max Pooling  $64 \times 32 \times 32 \mapsto 64 \times 16 \times 16$
- ▶ Conv  $64 \times 16 \times 16 \mapsto 128 \times 16 \times 16$  + BatchNorm2d + ReLU
- ▶ Conv  $128 \times 16 \times 16 \mapsto 128 \times 16 \times 16$  + ReLU
- ▶ Max Pooling  $128 \times 16 \times 16 \mapsto 128 \times 8 \times 8$
- ▶ Conv  $128 \times 8 \times 8 \mapsto 256 \times 8 \times 8$  + BatchNorm2d + ReLU
- ▶ Conv  $256 \times 8 \times 8 \mapsto 256 \times 8 \times 8$  + ReLU
- ▶ Max Pooling  $256 \times 8 \times 8 \mapsto 256 \times 4 \times 4$  + DropOut  $p = 0.05$
- ▶ Linear  $4096 \mapsto 1024$  + ReLU
- ▶ Linear  $1024 \mapsto 512$  + ReLU + DropOut  $p = 0.05$
- ▶ Linear  $512 \mapsto 10$  + SoftMax

We have added Batch Normalization to improve the training stability and Drop Out to reduce overfitting.

# CONVOLUTIONAL NEURAL NETWORKS

## MODEL 3 WITHOUT BATCH NORMALIZATION AND DROP OUT



# CONVOLUTIONAL NEURAL NETWORKS

## BATCH NORMALIZATION

Batch normalization is used to normalize the activations of a layer within a batch of data.

- ▶ This helps to prevent the problem of vanishing or exploding gradients and also speeds up the training process.
- ▶ By normalizing the activations, batch normalization helps to stabilize the distribution of the inputs to each layer, reducing the covariate shift and allowing the network to learn more effectively.

# CONVOLUTIONAL NEURAL NETWORKS

## BATCH NORMALIZATION

- 1: **for** each  $x_i$  in a mini-batch  $B$  of size  $b$  **do**
- 2:   Compute the mean  $\mu_B$  and variance  $\sigma_B^2$  of the features in the mini-batch  $B$ .

$$\mu_B = \frac{1}{b} \sum_i x_i \quad \text{and} \quad \sigma_B^2 = \frac{1}{m} \sum_i (x_i - \mu_B)^2$$

- 3:   Normalize each feature  $x_i$  in the mini-batch  $B$  using  $\mu_B$  and  $\sigma_B^2$ .

$$\bar{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \varepsilon}}$$

- 4:   Scale and shift each normalized feature  $x_i$  using two learnable parameters  $\gamma$  and  $\beta$  respectively.

$$y_i = \gamma \bar{x}_i + \beta$$

- 5:   Update the moving average of the mean and variance for inference:

$$\mu = (1 - \gamma)\mu + \gamma\mu_B \quad \text{and} \quad \sigma^2 = (1 - \gamma)\sigma^2 + \gamma\sigma_B^2$$

- 6: **end for**

**Algorithm 1:** Batch Normalization during training

# CONVOLUTIONAL NEURAL NETWORKS

## BATCH NORMALIZATION

At inference:

1: **for** each  $x_i$  **do**

2:   Normalize each feature  $x_i$  using the moving average of the mean and variance:

$$\bar{x}_i = \frac{x_i - \mu}{\sqrt{\sigma^2 + \varepsilon}}$$

3:   Scale and shift each normalized feature  $x_i$  using two learnable parameters  $\gamma$  and  $\beta$  respectively.

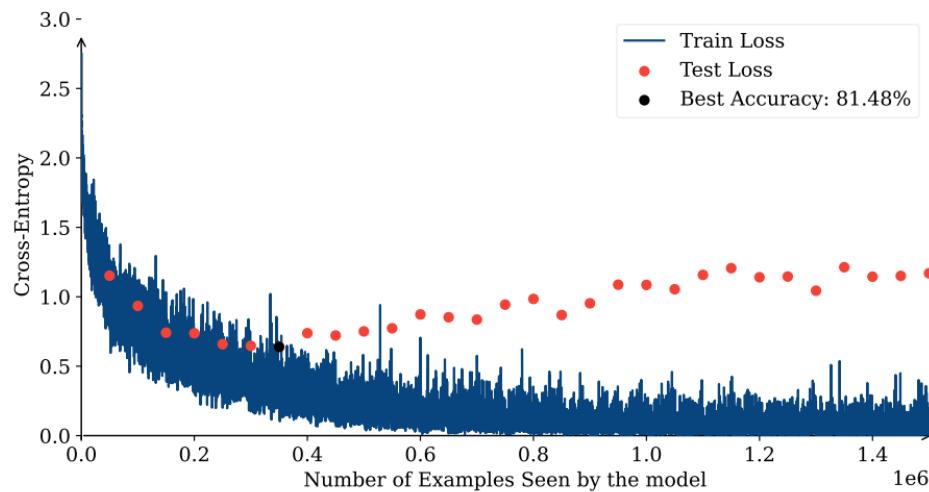
$$y_i = \gamma \bar{x}_i + \beta$$

4: **end for**

**Algorithm 2:** Batch Normalization at inference

# CONVOLUTIONAL NEURAL NETWORKS

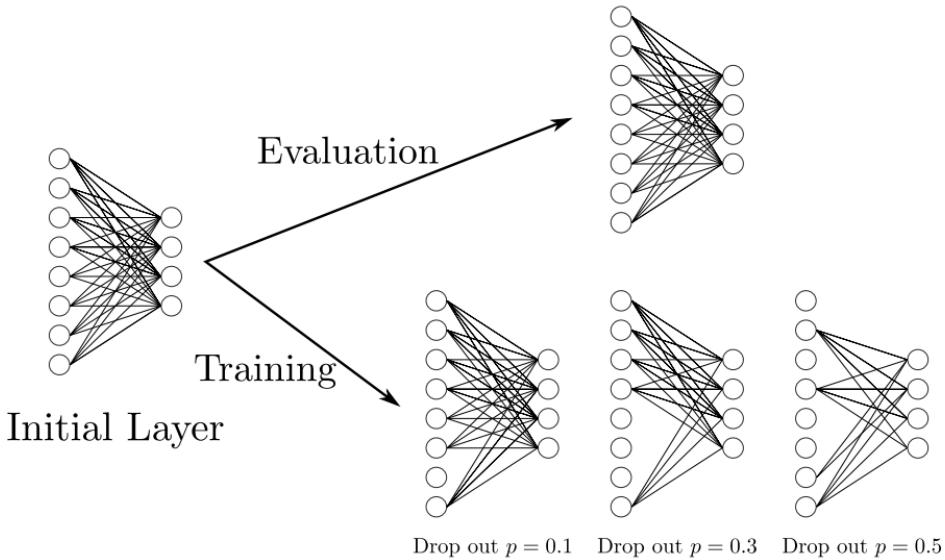
## MODEL 3 WITHOUT DROP OUT



# CONVOLUTIONAL NEURAL NETWORKS

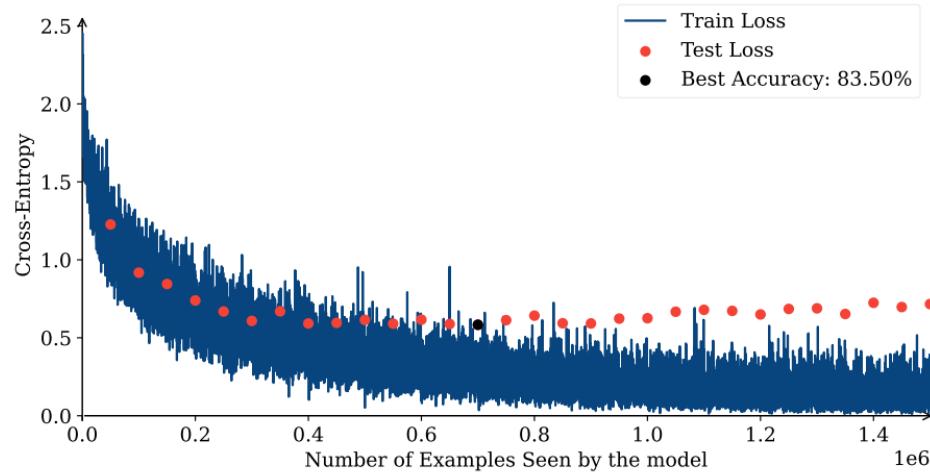
## DROP OUT

Dropout is a regularization technique in neural networks where *during training*, a portion of the nodes are randomly "dropped out" or ignored during each iteration. This helps prevent over-fitting by preventing the model from relying too heavily on any one node. The result is a more robust and generalizable model that can better handle unseen data.



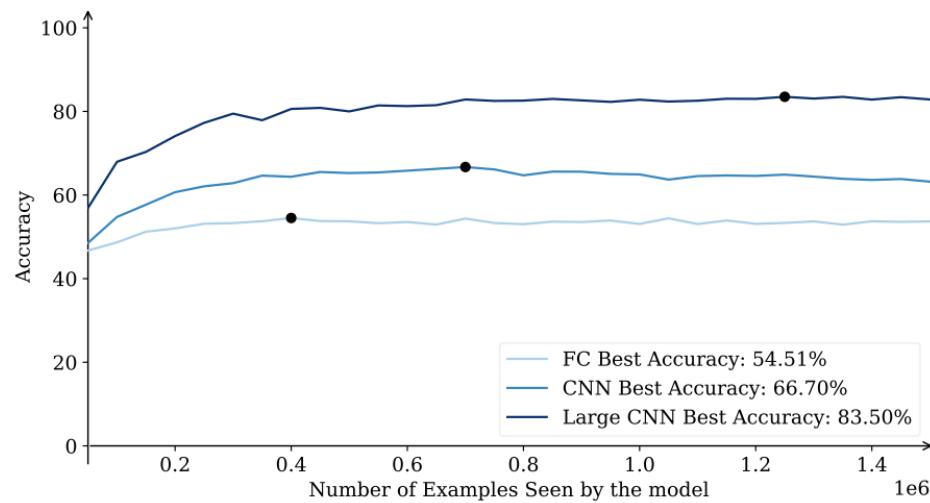
# CONVOLUTIONAL NEURAL NETWORKS

## MODEL 3



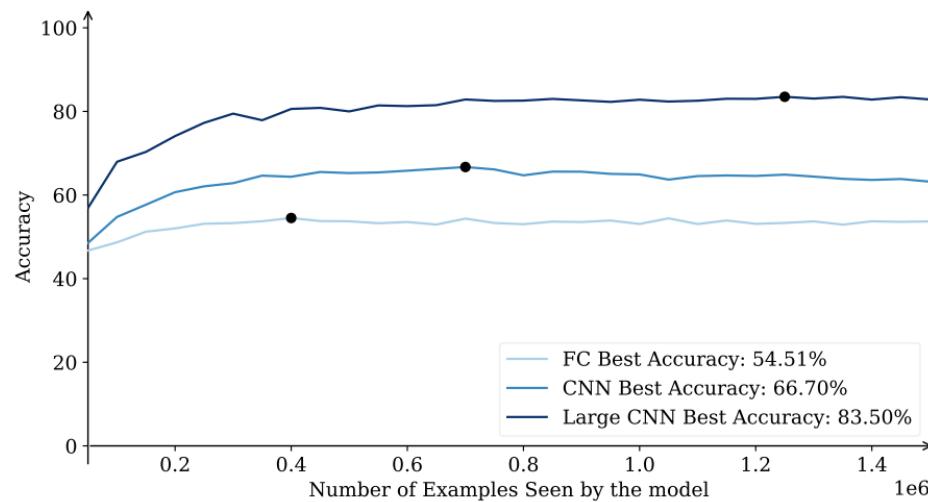
# CONVOLUTIONAL NEURAL NETWORKS

## MODEL 3



# CONVOLUTIONAL NEURAL NETWORKS

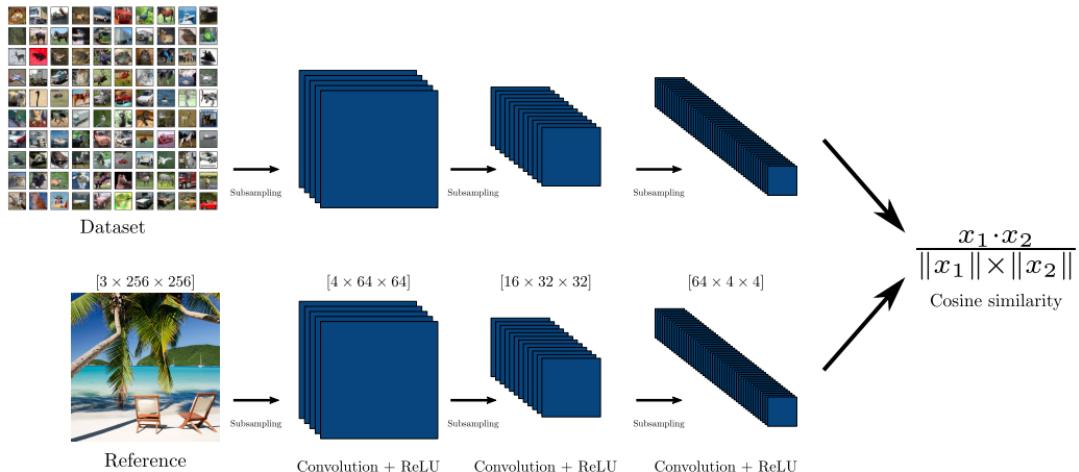
## CNN IN PRACTICE: CIFAR 10



# CONVOLUTIONAL NEURAL NETWORKS

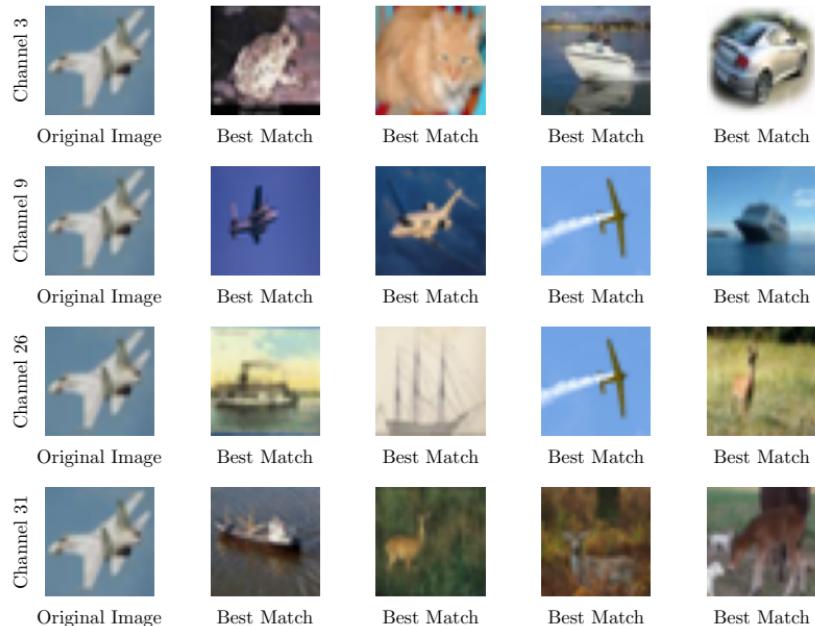
## INTUITION BEHIND CHANNELS

To examine the information captured by different channels in a Neural Network, we can compare their output on a dataset. For a given input  $x$ , we can compute the similarity between the output of a specific channel and the same channel for other images in the dataset.



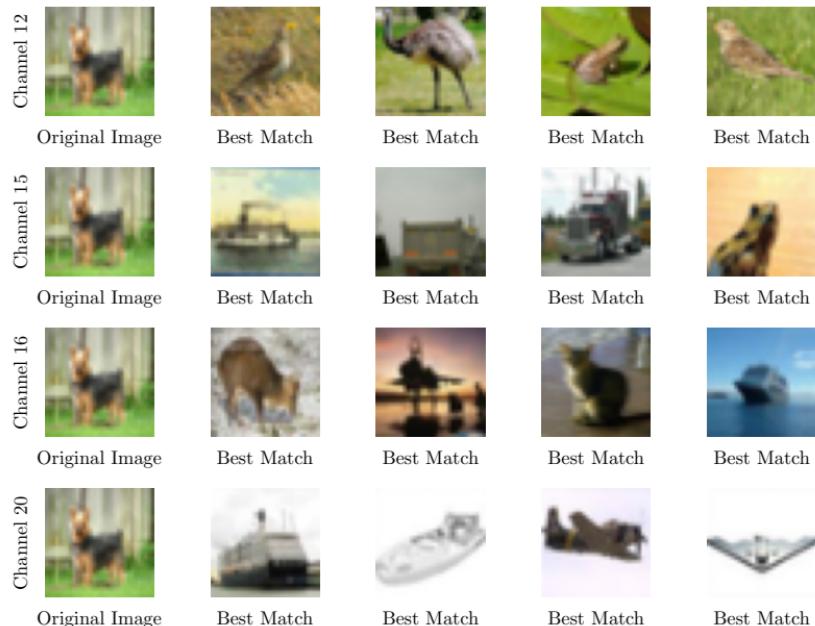
# CONVOLUTIONAL NEURAL NETWORKS

## INTUITION BEHIND CHANNELS



# CONVOLUTIONAL NEURAL NETWORKS

## INTUITION BEHIND CHANNELS



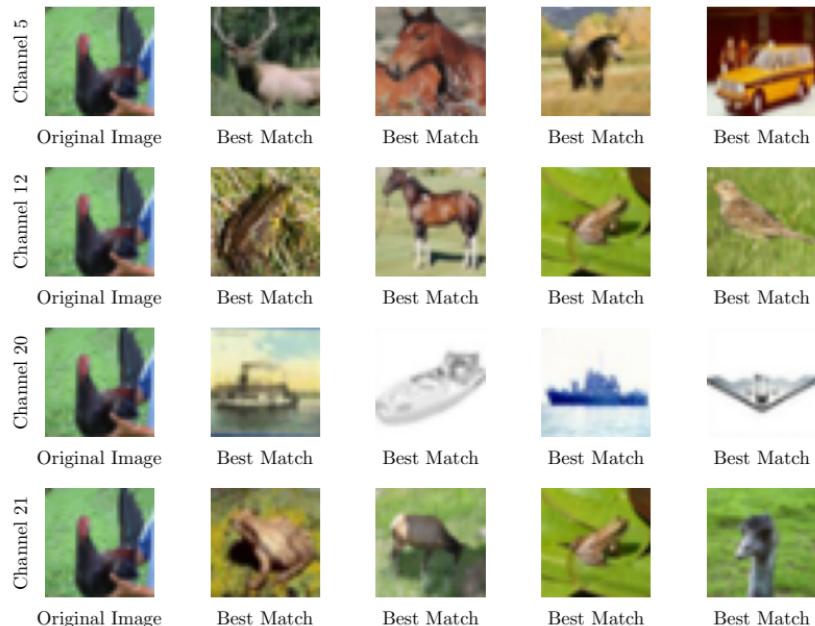
# CONVOLUTIONAL NEURAL NETWORKS

## INTUITION BEHIND CHANNELS

	Channel 2	Channel 12	Channel 13	Channel 24
Original Image				
Best Match				
Best Match				
Best Match				
Best Match				

# CONVOLUTIONAL NEURAL NETWORKS

## INTUITION BEHIND CHANNELS



## TP2: THE CONVOLUTIONAL NEURAL NETWORK

### INTRODUCTION TO CNN AND CIFAR-10 DATASET

Link to the notebook (ipynb): [TP2.ipynb](#)

Link to the notebook (html): [TP2.html](#)

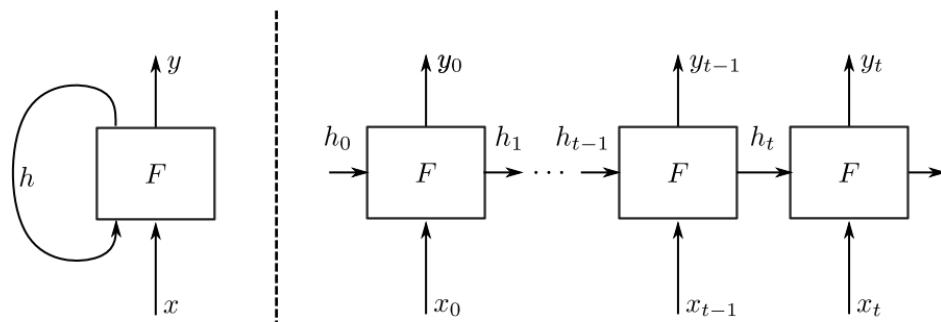
## RECURRENT NEURAL NETWORKS

Recurrent Networks (RNNs) are a type of neural network that are specifically designed to handle sequential data, whereas CNNs are more suited for image and grid-like data. The main difference between RNNs and CNNs lies in the way they process data, with RNNs considering the sequence of elements and their interdependencies, while CNNs focus on capturing local patterns within the input.

# RECURRENT NEURAL NETWORKS

## RECURRENT BLOCK

A Recurrent Network is a type of neural network that contains a loop mechanism, allowing previous outputs to be used as inputs for future computations. This creates a form of memory that allows the network to process sequential data with variable-length sequences.



Rolled

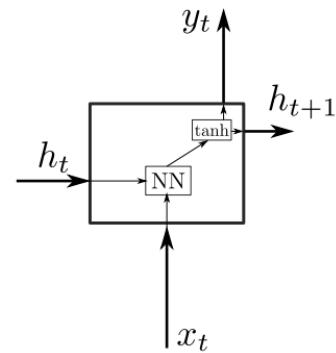
Unrolled

# RECURRENT NEURAL NETWORKS

## RECURRENT BLOCK

Some of the limitations of Vanilla RNNs:

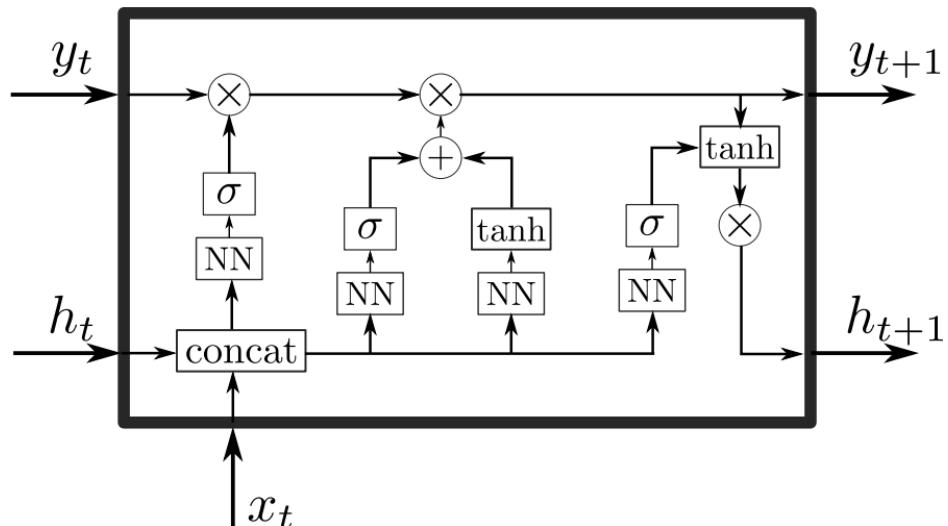
- ▶ Vanishing gradient problem: The gradient signals used to update the weights during training can become very small, making it difficult to train RNNs effectively.
- ▶ Exploding gradient problem: On the other hand, gradients can become too large and cause numeric instability, making it difficult to train RNNs effectively.
- ▶ Short-term memory: Vanilla RNNs have difficulty retaining information over long periods of time, making them unsuitable for tasks that require remembering information from previous time steps.
- ▶ Computational limitations: RNNs can be computationally intensive, making it difficult to apply them to large sequences of data.
- ▶ Difficulty with parallelization: The sequential nature of RNNs can make it difficult to take advantage of parallel processing to speed up training and inference.



# RECURRENT NEURAL NETWORKS

## LSTM

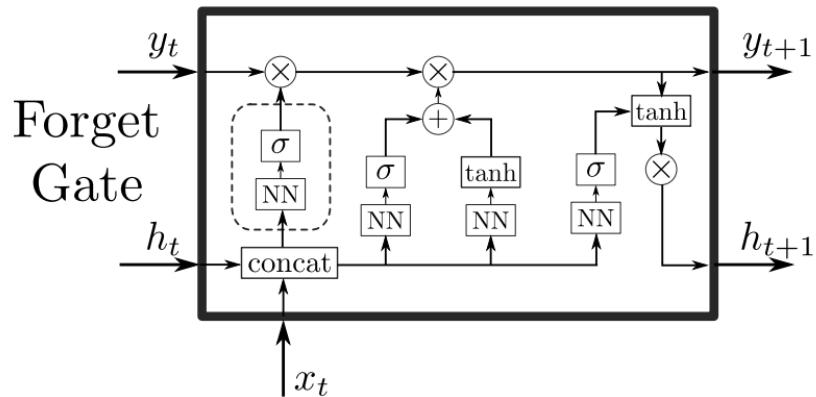
Long Short-Term Memory (LSTM) networks are a variant of recurrent neural networks (RNNs) that overcome some of the limitations of traditional RNNs, such as the vanishing gradient problem and difficulty in learning long-term dependencies. LSTM networks introduce memory cells, gates, and a process for updating cells, which allows them to selectively preserve information from previous time steps.



# RECURRENT NEURAL NETWORKS

## LSTM

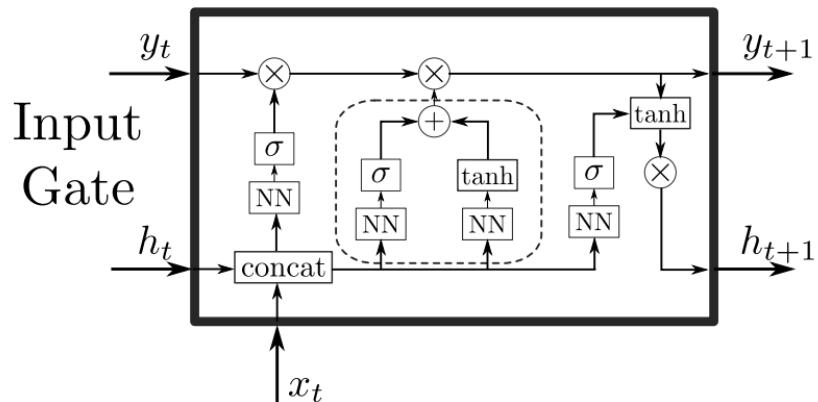
Long Short-Term Memory (LSTM) networks are a variant of recurrent neural networks (RNNs) that overcome some of the limitations of traditional RNNs, such as the problem of vanishing gradients and the difficulty of learning long-term dependencies. LSTM networks introduce memory cells, gates, and a process for updating cells, which allows them to selectively preserve information from previous time steps.



# RECURRENT NEURAL NETWORKS

## LSTM

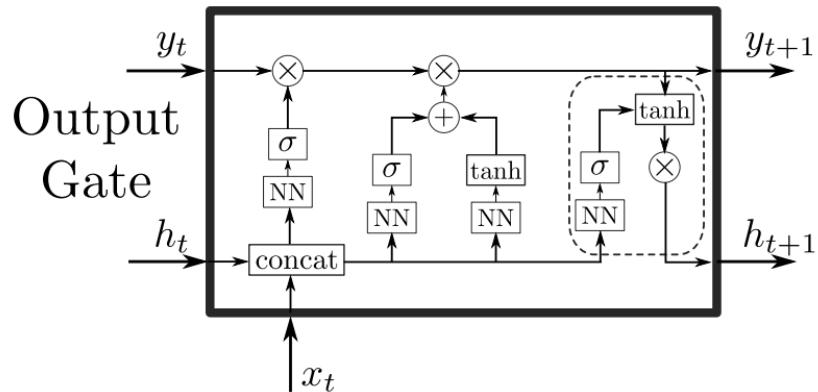
Long Short-Term Memory (LSTM) networks are a variant of recurrent neural networks (RNNs) that overcome some of the limitations of traditional RNNs, such as the vanishing gradient problem and difficulty in learning long-term dependencies. LSTM networks introduce memory cells, gates, and a process for updating the cells, which allows them to selectively preserve information from previous time steps.



# RECURRENT NEURAL NETWORKS

## LSTM

Long Short-Term Memory (LSTM) networks are a variant of recurrent neural networks (RNNs) that overcome some of the limitations of traditional RNNs, such as the vanishing gradient problem and difficulty in learning long-term dependencies. LSTM networks introduce memory cells, gates, and a process for updating the cells, which allows them to selectively preserve information from previous time steps.



# RECURRENT NEURAL NETWORKS

## LIMITS OF LSTM

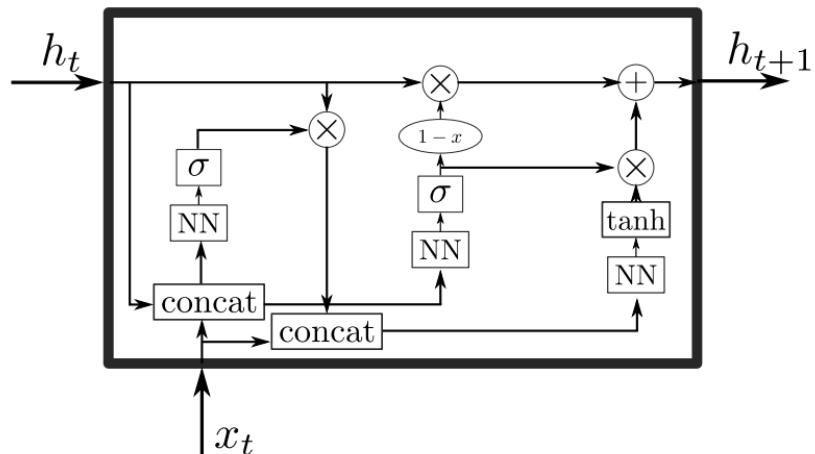
Limitations of LSTM RNNs:

- ▶ High computational cost: LSTMs are computationally more expensive compared to other traditional neural network models due to the presence of multiple gates and their sequential processing nature.
- ▶ Vanishing Gradient Problem: LSTMs, like any other RNNs, are prone to the vanishing gradient problem when the sequences are too long, making it difficult for the model to learn long-term dependencies.
- ▶ Overfitting: LSTMs are complex models and are more susceptible to overfitting compared to simple feedforward networks.
- ▶ Difficult to parallelize: Due to the sequential nature of LSTMs, they are difficult to parallelize and can take longer to train.
- ▶ Gradient Explosion: LSTMs can also suffer from the gradient explosion problem, where the gradients can become too large and cause numerical instability during training.

# RECURRENT NEURAL NETWORKS

## GRU

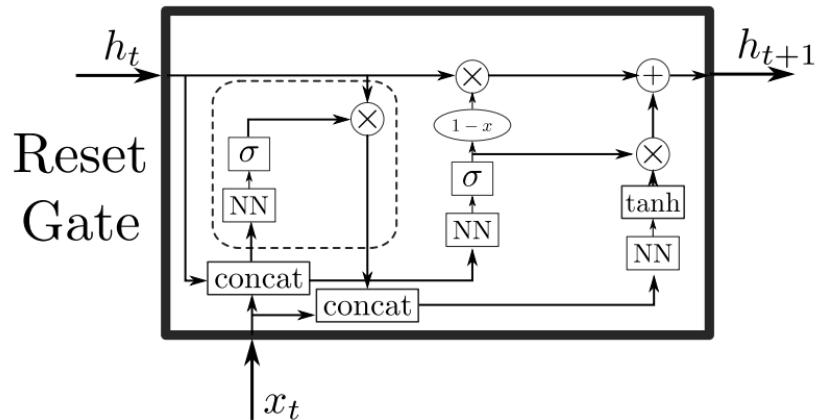
GRU blocks, or Gated Recurrent Units, are a type of recurrent neural network architecture that are similar to LSTMs in their function and ability to process sequential data. GRUs were introduced as a simplification of LSTMs, with the aim of reducing the number of parameters in the network and improving computational efficiency. GRUs achieve this by merging the forget and input gates in LSTMs into a single update gate, effectively combining the two operations in a single step.



# RECURRENT NEURAL NETWORKS

## GRU

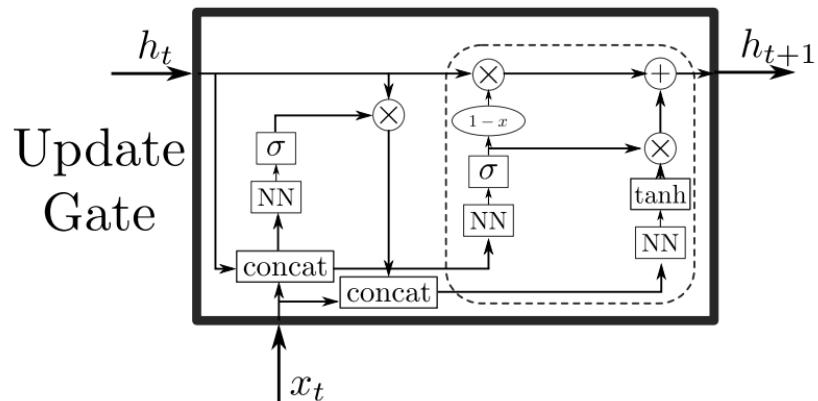
GRU blocks, or Gated Recurrent Units, are a type of recurrent neural network architecture that are similar to LSTMs in their function and ability to process sequential data. GRUs were introduced as a simplification of LSTMs, with the aim of reducing the number of parameters in the network and improving computational efficiency. GRUs achieve this by merging the forget and input gates in LSTMs into a single update gate, effectively combining the two operations in a single step.



# RECURRENT NEURAL NETWORKS

## GRU

GRU blocks, or Gated Recurrent Units, are a type of recurrent neural network architecture that are similar to LSTMs in their function and ability to process sequential data. GRUs were introduced as a simplification of LSTMs, with the aim of reducing the number of parameters in the network and improving computational efficiency. GRUs achieve this by merging the forget and input gates in LSTMs into a single update gate, effectively combining the two operations in a single step.



# RECURRENT NEURAL NETWORKS

## LSTM AND GRU

Limitations of GRU RNNs:

- ▶ Computational complexity: GRUs are more computationally efficient than LSTMs but still more complex than feedforward neural networks.
- ▶ Long-term dependencies: GRUs may struggle with capturing long-term dependencies in sequences, although they perform better in this regard than vanilla RNNs.
- ▶ Vanishing gradient problem: GRUs can still be affected by the vanishing gradient problem that plagues all RNN models. This problem makes it difficult for the model to learn from long sequences.
- ▶ Non-stationary data: GRUs may struggle with nonstationary data, where the statistical properties of the data change over time.

# RECURRENT NEURAL NETWORKS

## APPLICATION OF RNNs

Applications of RNNs:

- ▶ Natural language processing (NLP): Using RNNs for text classification, language translation, and text generation.
- ▶ Time-series prediction: Using RNNs to make predictions based on sequential data, such as stock prices and weather patterns.
- ▶ Speech recognition: Using RNNs for speech-to-text conversion.

## TRANSFORMER AND ATTENTION MECHANISM

Transformers and Attention Mechanisms are recent advancements in deep learning, widely used for sequential data processing, especially in NLP. Unlike RNNs, which apply the same weights iteratively, Transformers leverage self-attention to dynamically assign importance to different sequence elements. This enhances their ability to capture long-range dependencies, improving NLP performance.

# TRANSFORMER AND ATTENTION MECHANISM

## SELF-ATTENTION MECHANISM

The self-attention mechanism in Transformers computes the importance of each token relative to others in a sequence. It generates a weighted representation where the most relevant tokens contribute more. Mathematically, self-attention is computed as:

$$\text{Attention}(Q, K, V) = \text{Softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right) V$$

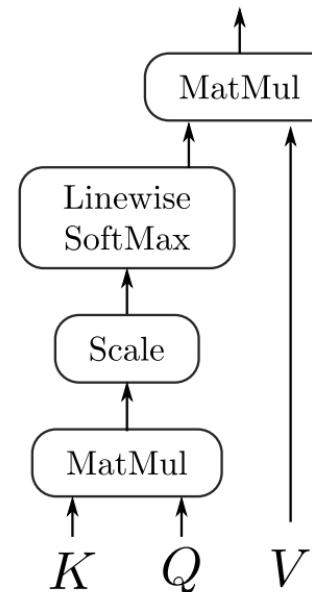
where  $Q \in \mathbb{R}^{m \times d_k}$ ,  $K \in \mathbb{R}^{n \times d_k}$ , and  $V \in \mathbb{R}^{n \times d_v}$  are the query, key, and value matrices derived from the input sequence. The softmax function ensures proper weighting, producing the final representation.

# TRANSFORMER AND ATTENTION MECHANISM

## SELF-ATTENTION MECHANISM

### Intuition behind the self-attention mechanism:

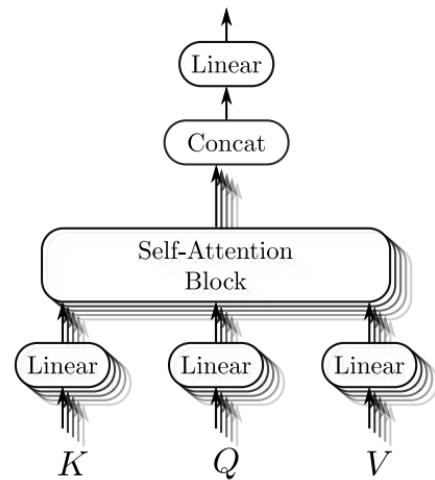
- ▶ The **query** and **key** matrices measure the similarity between different tokens in the sequence.
- ▶ These similarity scores determine how much each **value** contributes to the final output.
- ▶ The resulting **weighted sum** of the values represents the output of the self-attention mechanism.
- ▶ This process captures the relationships between different parts of the input sequence, allowing the model to focus on the most relevant elements.



# TRANSFORMER AND ATTENTION MECHANISM

## MULTI-HEAD ATTENTION

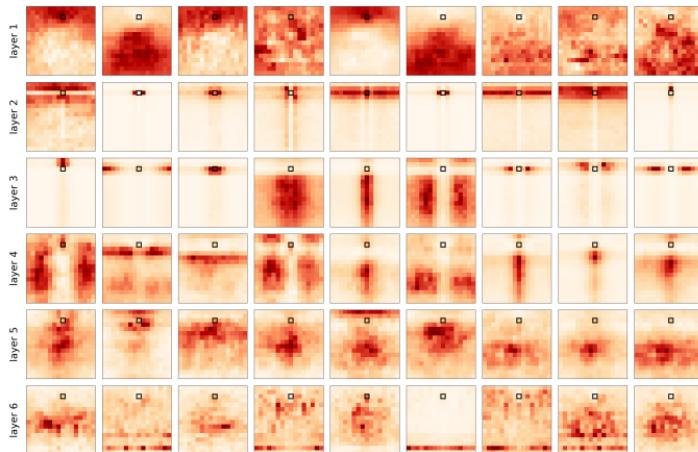
In Multi-head Attention, the self-attention mechanism is performed multiple times in parallel with different weight matrices, before being concatenated and once again projected, leading to a more robust representation of the input sequence. The intuition behind the three matrices ( $Q$ ,  $K$ ,  $V$ ) remains the same as in self-attention, with  $Q$  representing the query,  $K$  the key and  $V$  the value. Each head performs an attention mechanism on the input sequence, capturing different aspects and dependencies of the data, before being combined to form a more comprehensive representation of the input.



# TRANSFORMER AND ATTENTION MECHANISM

## VISUALIZING MULTI-HEAD ATTENTION

Visualizing Self-Attention for  
Image:  
Link



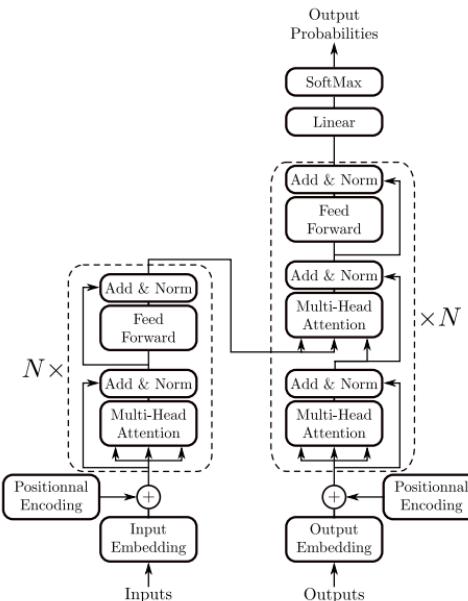
# TRANSFORMER AND ATTENTION MECHANISM

## TRANSFORMERS MODEL

Transformers are neural network models that use an encoder-decoder architecture. The encoder takes the input sequence and converts it into a continuous hidden representation, which is then passed to the decoder to generate the output sequence. The architecture of the transformer model is designed to allow the model to process the entire sequence in parallel, rather than processing one element at a time like in traditional RNNs.

Training of transformers involves optimizing a loss function that measures the difference between the model predictions and the true outputs. This loss function is usually based on the cross entropy between the predicted and true sequences.

The encoder-decoder mechanism is commonly referred to as the seq2seq mechanism.



## Part III

# TECHNIQUES TO IMPROVE DEEP LEARNING TRAINING

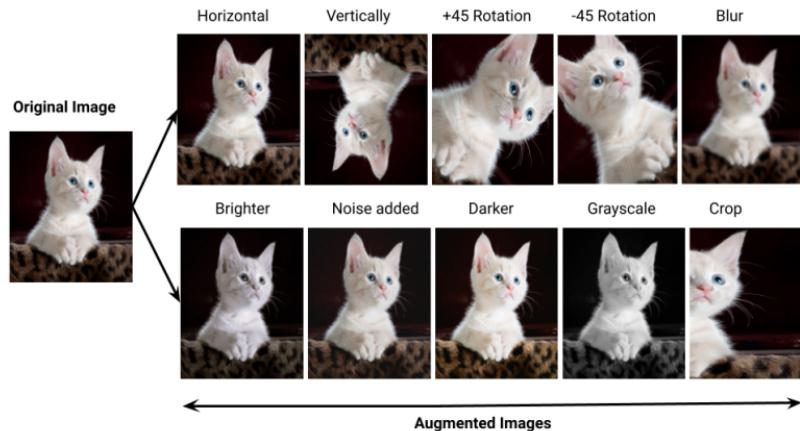
# TECHNIQUES TO IMPROVE DEEP LEARNING TRAINING

- ▶ **Batch Normalization** - Seen
- ▶ **Dropout** - Seen
- ▶ **Data Augmentation**
- ▶ **Learning Rate Scheduling**
- ▶ **Early Stopping**
- ▶ **Gradient Clipping**
- ▶ **Weight Initialization**
- ▶ **Regularization**
- ▶ **GPU Acceleration**

# TECHNIQUES TO IMPROVE DEEP LEARNING TRAINING

## DATA AUGMENTATION

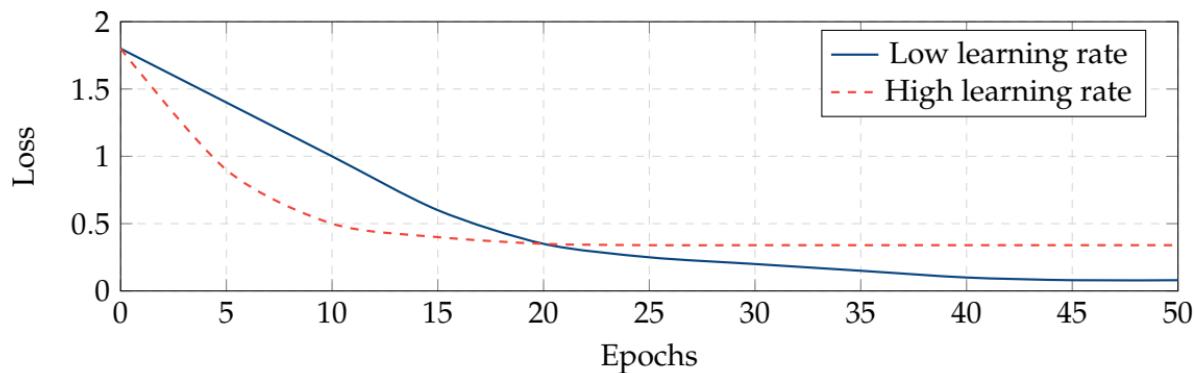
- ▶ Data Augmentation is a technique to increase the diversity of your training set by applying random (but realistic) transformations to the training images.
- ▶ The goal is to train a model that is robust to these transformations.
- ▶ For example, you can randomly rotate, scale, and flip the images in your training set.
- ▶ This helps expose the model to different aspects of the data and reduce overfitting.



## TECHNIQUES TO IMPROVE DEEP LEARNING TRAINING

### LEARNING RATE SCHEDULING

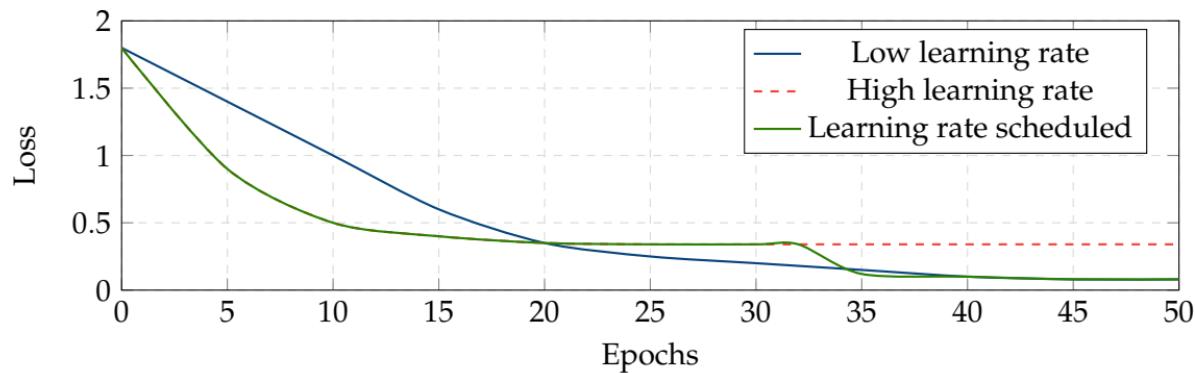
- ▶ The learning rate is one of the most important hyperparameters to tune for your deep learning model. The learning rate determines how quickly the model learns the optimal weights and how refined the gradient descent process is.
- ▶ If the learning rate is too high, the model may not converge or converge to a higher loss. If it is too low, the model may take too long to train.



# TECHNIQUES TO IMPROVE DEEP LEARNING TRAINING

## LEARNING RATE SCHEDULING

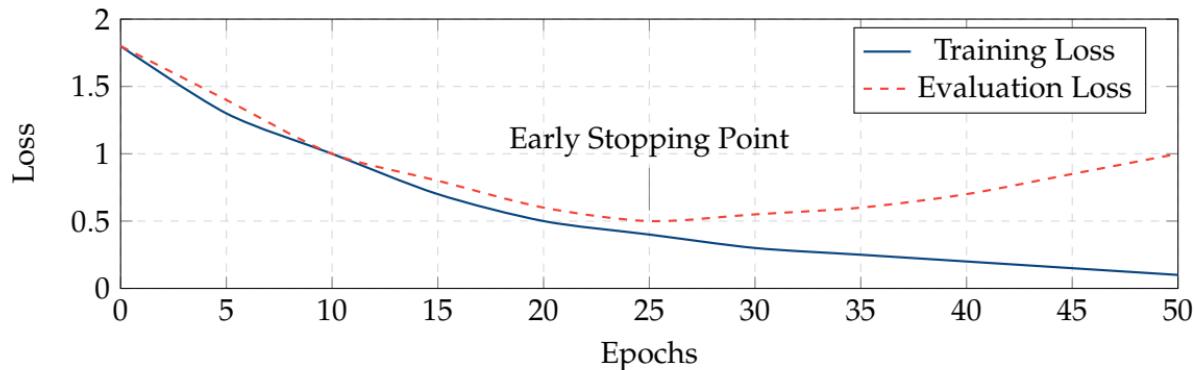
- ▶ The learning rate is one of the most important hyperparameters to tune for your deep learning model. The learning rate determines how quickly the model learns the optimal weights and how refined the gradient descent process is.
- ▶ If the learning rate is too high, the model may not converge or converge to a higher loss. If it is too low, the model may take too long to train.
- ▶ Learning rate scheduling is a technique to adjust the learning rate during training. For example, you can start with a high learning rate and then decrease it over time.



# TECHNIQUES TO IMPROVE DEEP LEARNING TRAINING

## EARLY STOPPING

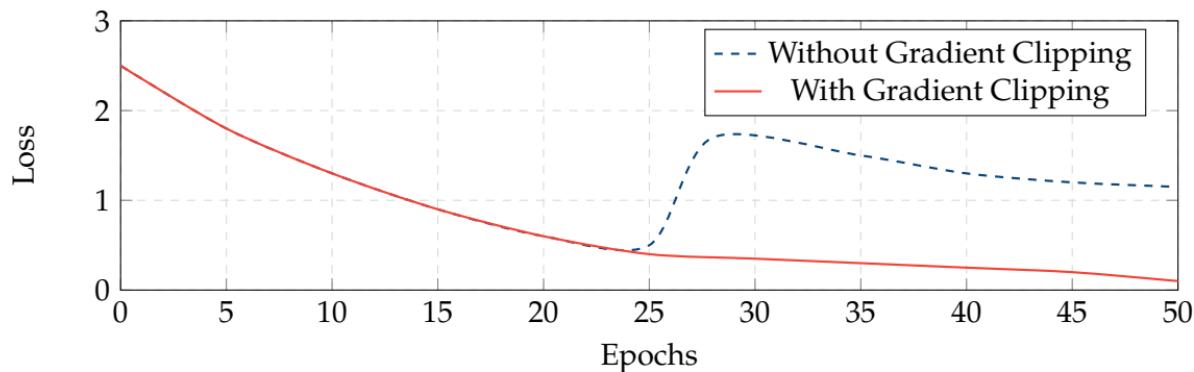
- ▶ Early stopping is a technique to prevent overfitting by stopping the training process when the model's performance on the validation set starts to degrade.
- ▶ The idea is to monitor the validation loss during training and stop training when the validation loss stops decreasing.



# TECHNIQUES TO IMPROVE DEEP LEARNING TRAINING

## GRADIENT CLIPPING

- ▶ Gradient clipping is a technique to prevent exploding gradients during training.
- ▶ Exploding gradients occur when the gradients of the loss function with respect to the model's parameters are too large.
- ▶ This can cause the model to diverge and fail to learn.
- ▶ Gradient clipping involves scaling the gradients if their norm exceeds a certain threshold.



# TECHNIQUES TO IMPROVE DEEP LEARNING TRAINING

## WEIGHT INITIALIZATION

- ▶ Weight initialization is a technique to set the initial values of the weights in the model.
- ▶ The initial values of the weights can have a significant impact on the training process and the final performance of the model.
- ▶ If the weights are initialized too small, the model may not learn effectively. If they are initialized too large, the model may not converge.
- ▶ Common weight initialization techniques include Xavier/Glorot initialization and He initialization.

# TECHNIQUES TO IMPROVE DEEP LEARNING TRAINING

## WEIGHT INITIALIZATION

- ▶ Xavier/Glorot initialization: The weights are initialized from a normal distribution with mean 0 and variance  $2/(n_{\text{in}} + n_{\text{out}})$ , where  $n_{\text{in}}$  and  $n_{\text{out}}$  are the number of input and output units, respectively. It helps prevent the gradients from vanishing or exploding during training by ensuring that the gradients have a similar scale.
- ▶ He initialization: The weights are initialized from a normal distribution with mean 0 and variance  $2/n_{\text{in}}$ , where  $n_{\text{in}}$  is the number of input units. It is commonly used for ReLU activation functions.

# TECHNIQUES TO IMPROVE DEEP LEARNING TRAINING

## REGULARIZATION

- ▶ Regularization is a technique to prevent overfitting by adding a penalty term to the loss function that discourages the model from learning complex patterns that may not generalize well.
- ▶ L1 regularization adds a penalty term to the loss function that is proportional to the absolute value of the weights. It encourages sparsity in the weights.
- ▶ L2 regularization adds a penalty term to the loss function that is proportional to the square of the weights. It encourages the weights to be small.

# TECHNIQUES TO IMPROVE DEEP LEARNING TRAINING

## GPU ACCELERATION

CPUs and GPUs are very different in terms of architecture and performance. CPUs are more suited for general-purpose computing tasks, while GPUs are optimized for parallel processing of simple operations, making them ideal for deep learning tasks.

- ▶ GPUs are much faster than CPUs for deep learning tasks because they have many more cores and can perform many more operations in parallel.
- ▶ Deep learning frameworks like PyTorch and TensorFlow are designed to take advantage of GPUs to accelerate the training process.
- ▶ Only the forward and backward passes of the model are executed on the GPU. The data loading and preprocessing are still done on the CPU.

# Part IV

## DEEP LEARNING AND APPLICATIONS

# LEARNING TO ACT WITH DEEP REINFORCEMENT LEARNING

## THE MOUSE GAME

- ▶ A Mouse has to feed on food (red) and avoid poison (blue).
- ▶ It has a vision range of 2 squares. So it can see the 25 cells around.
- ▶ The reward for a cheese cell is 0.5, while the reward for eating poison is  $-1$ .

On this example, the mouse behaves randomly.

# REINFORCEMENT LEARNING

## NOTATIONS

The environment is a finite Markov Decision Process (MDP) of size  $T$  defined by the tuple  $(S, A, P, R, \gamma)$ :

- ▶  $S$ : set of states  $s$ .
- ▶  $A$ : set of actions  $a$ .
- ▶  $P(s'|s, a)$ : transition probability. The environment can be deterministic if the mouse moves in the direction it wants to go and thus  $P(s'|s, a) = 1$  if the mouse can move in the direction  $a$  and 0 otherwise.
- ▶  $r(s, a, s')$ : reward function.
- ▶  $\gamma$ : discount factor. The discount factor is used to give more importance to the immediate reward than the future reward.
- ▶  $T$ : horizon.

# REINFORCEMENT LEARNING

## EXPECTED REWARD

In RL, the goal is to find a policy  $\pi(a|s)$ , the probability of taking action  $a$  in state  $s$ , that maximizes the expected reward:

$$R(\pi) = \mathbb{E}_{s_0, P, \pi} \left[ \sum_{t=0}^T \gamma^t r(s_t, a_t, s_{t+1}) \right],$$

and if  $\gamma$  is close to 0, the reward is the immediate reward, and if  $\gamma$  is close to 1, the reward is the cumulative reward.

One way to find the optimal policy  $\pi^*$  is the Q-learning algorithm. The Q-value is the expected reward of taking action  $a$  in state  $s$  and following the policy  $\pi$ . The Q-value is defined as:

$$Q^\pi(s, a) = \mathbb{E}_{P, \pi} \left[ \sum_{t=0}^T \gamma^t r(s_t, a_t, s_{t+1}) \middle| s_0 = s, a_0 = a \right].$$

# LEARNING TO ACT WITH DEEP REINFORCEMENT LEARNING

## DEEP Q-LEARNING

The Q-learning algorithm learns the optimal action-value function  $Q^*(s, a)$ , using the Bellman equation:

$$Q^*(s, a) = \mathbb{E}_{s' \sim \pi^*} \left[ r(s, a, s') + \gamma \max_{a'} Q^*(s', a') \right].$$

To approximate the Q-value, a deep learning model is used to predict the Q-value for each action in a given state. The model is trained to minimize the difference between the predicted Q-value and the target Q-value:

$$\mathcal{L}(\theta) = \mathbb{E}_{s,a,s'} \left[ \left( r(s, a, s') + \gamma \max_{a'} Q_{\theta^-}(s', a') - Q_\theta(s, a) \right)^2 \right].$$

# LEARNING TO ACT WITH DEEP REINFORCEMENT LEARNING

## DEEP Q-LEARNING

We can use the following algorithm to train the agent:

1. Initialize the replay buffer  $D$ .
2. Initialize the Q-network  $Q_\theta$ .
3. For each episode:
  - 3.1 Initialize the state  $s$ .
  - 3.2 For each step:
    - i. The agent chooses an action  $a$  and get the reward  $r$  and the next state  $s'$ .
    - ii. Store the transition  $(s, a, r, s')$  in the replay buffer  $D$ .
    - iii. Sample a batch of transitions from the replay buffer  $D$ . (*Batched*)
    - iv. Compute the target  $y(a)$  using the target network  $Q_{\theta-}$ . (*Batched*)
    - v. Add  $\gamma + \max_{a'} y(a')$  to the target  $y(a)$  if the episode is not done. (*Batched*)
    - vi. Update the Q-network  $Q_\theta$  using the loss function. (*Batched*)

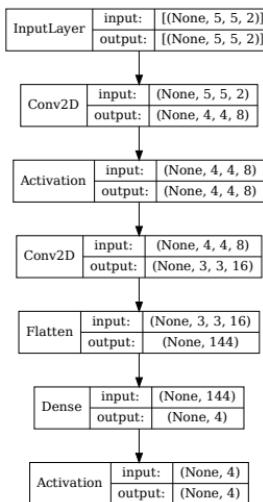
With a Q-network learned, we can use the following policy to choose the action:

$$\pi(a|s) = \begin{cases} 1 & \text{if } a \in \arg \max_{a'} Q_\theta(s, a') \\ 0 & \text{otherwise} \end{cases}.$$

# LEARNING TO ACT WITH DEEP REINFORCEMENT LEARNING

## MODEL 2

Convolutional network:



## LEARNING TO ACT WITH DEEP REINFORCEMENT LEARNING

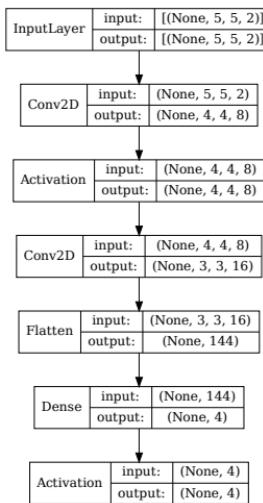
### $\epsilon$ -GREEDY ALGORITHM

The  $\epsilon$ -greedy algorithm is a common exploration strategy used in Deep Q learning. Balances exploration, where the agent tries out new actions and collects new data, and exploitation, where the agent uses the information it already has to select the action with the highest expected reward. The algorithm selects a random action with probability  $\epsilon$  and the action with the highest Q value with probability  $1 - \epsilon$ . The value of  $\epsilon$  decreases over time to gradually shift the focus from exploration to exploitation.

# LEARNING TO ACT WITH DEEP REINFORCEMENT LEARNING

## MODEL 2+ INCORPORATED $\epsilon$ -EXPLORATION

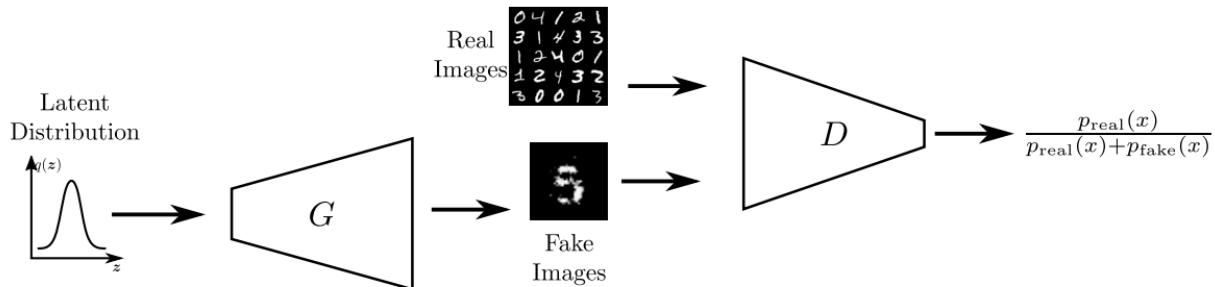
Convolutional network +  $\epsilon$ -greedy  
Algorithm during training:



# SYNTHETIC DATA GENERATION WITH GENERATIVE ADVERSARIAL NETWORKS

## GANS MODELS

Generative Adversarial Networks (GANs) are a type of deep learning architecture composed of two neural networks, the generator and the discriminator, that are trained in a adversarial manner. The generator network is trained to generate fake data that appears similar to real data, while the discriminator network is trained to distinguish between real and fake data. The loss for GANs is defined as a min-max game, where the generator minimizes the loss function, and the discriminator maximizes it.  $D$  and  $G$  represent the discriminator and generator networks, respectively, and the goal is to find the optimal configuration for  $D$  and  $G$  such that the generated samples appear indistinguishable from real data.

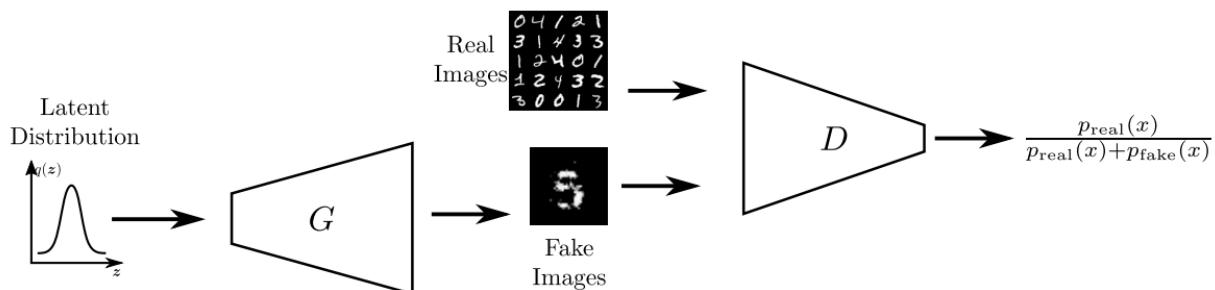


# SYNTHETIC DATA GENERATION WITH GENERATIVE ADVERSARIAL NETWORKS

## GANS MODELS

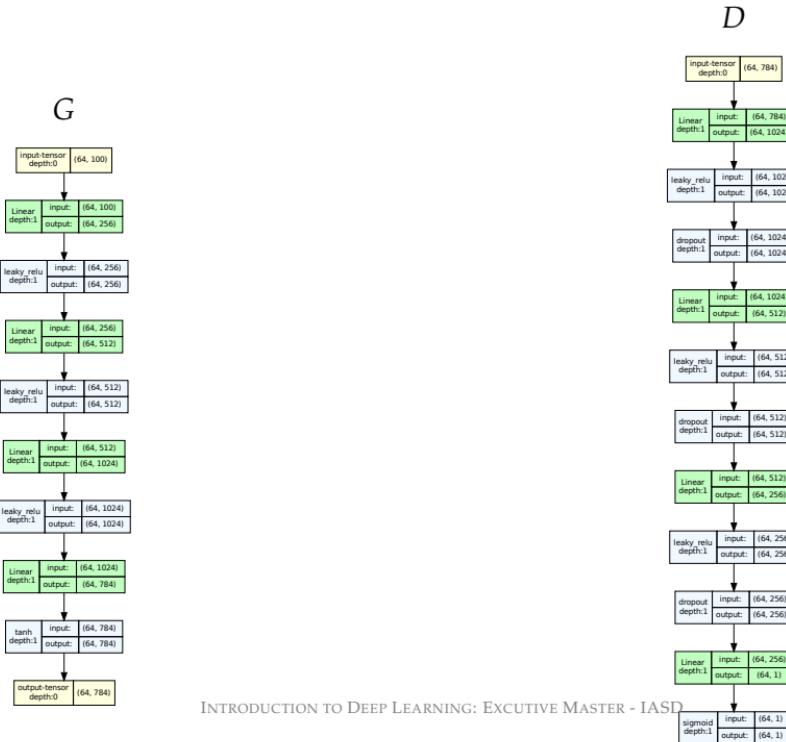
The loss of a (GAN) is defined as a minimax game between the generator and discriminator models. The generator aims to generate samples that are indistinguishable from real samples, while the discriminator aims to distinguish the generated samples from real samples. The loss function for the generator is defined as  $-\log(D(G(z)))$ , where  $D$  is the discriminator model and  $G(z)$  is the generator's output for a random noise vector  $z$ . The loss function for the discriminator is defined as  $\log(D(x)) + \log(1 - D(G(z)))$ , where  $x$  is a real sample.

$$\min_D \max_G \mathbb{E}_{x \sim p_{\text{real}}} [\log(D(x))] + \mathbb{E}_{z \sim q} [\log(1 - D(G(z)))]$$



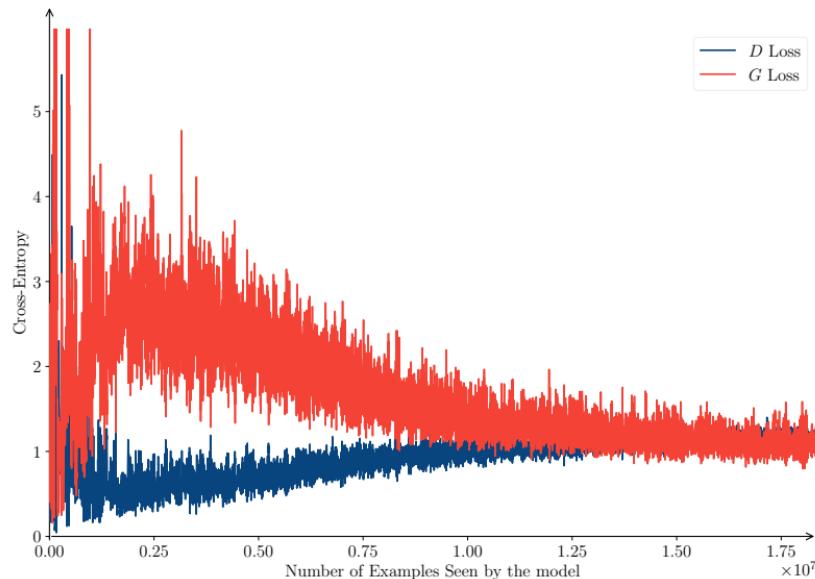
# SYNTHETIC DATA GENERATION WITH GENERATIVE ADVERSARIAL NETWORKS

## MNIST GENERATION



# SYNTHETIC DATA GENERATION WITH GENERATIVE ADVERSARIAL NETWORKS

## MNIST GENERATION



# SYNTHETIC DATA GENERATION WITH GENERATIVE ADVERSARIAL NETWORKS

## MNIST GENERATION

# SENTIMENT ANALYSIS WITH TRANSFORMERS AND GRU

## BERT

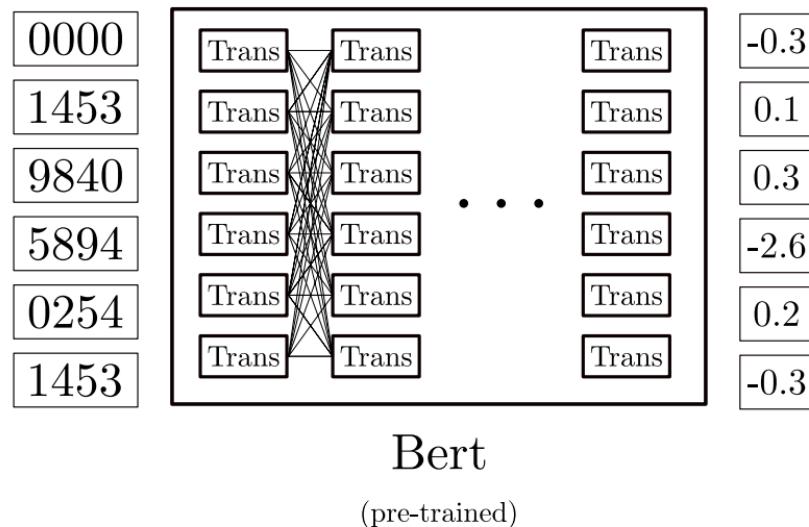
BERT (Bidirectional Encoder Representations from Transformers) is a pre-trained language model developed by Google in 2018. It is a transformer-based architecture that uses a masked language modeling task to generate a deep understanding of the contextual relationships between words in a sentence. BERT can be fine-tuned for various NLP tasks such as sentiment classification by adding a classification layer on top of its pre-trained representations. The model has achieved state-of-the-art performance in a wide range of NLP tasks, making it a popular choice for sentiment analysis.

A bidirectional transformer is a type of transformer architecture in natural language processing (NLP) where information from both past and future contexts is taken into consideration when making predictions. This is achieved by processing the input sequence in two directions, starting from the beginning and the end of the sequence, and concatenating the outputs to obtain the final representation. This allows the model to capture the context both in the forward and backward directions, providing a more comprehensive representation of the input sequence.

## SENTIMENT ANALYSIS WITH TRANSFORMERS AND GRU

### BERT

BERT model consists of multiple transformer encoder blocks, with a self-attention mechanism, a feedforward neural network and layer normalization, stacked on top of each other. It also includes a positional encoding component to capture the relative position of tokens in a sequence, and a segment encoding component to differentiate between different sequences within the same input.



# SENTIMENT ANALYSIS WITH TRANSFORMERS AND GRU

## TRAINING BERT

Bert is trained on two tasks, the masked language model and the next sentence prediction. In the masked language model task, a portion of the input sequence is masked and the model must predict the original token based on its context. In the next sentence prediction task, the model receives a pair of sentences and must predict whether the second sentence follows the first one in the context of the input text. Both of these tasks are used to train Bert to understand the context of words in a sentence and how they relate to each other, allowing it to perform well on a wide range of natural language processing tasks.

### Sentence

There	is	no	curse	in	Elvish,	Entish,	or	the	tongues	of	Men	for	this	treachery.
-------	----	----	-------	----	---------	---------	----	-----	---------	----	-----	-----	------	------------

### Basic-level masking

There	is		curse	in	Elvish,		or	the	tongues	of	Men	for	this	
-------	----	--	-------	----	---------	--	----	-----	---------	----	-----	-----	------	--

### Entity-level masking

There	is			in	Elvish,	Entish,	or	the				for	this	treachery.
-------	----	--	--	----	---------	---------	----	-----	--	--	--	-----	------	------------

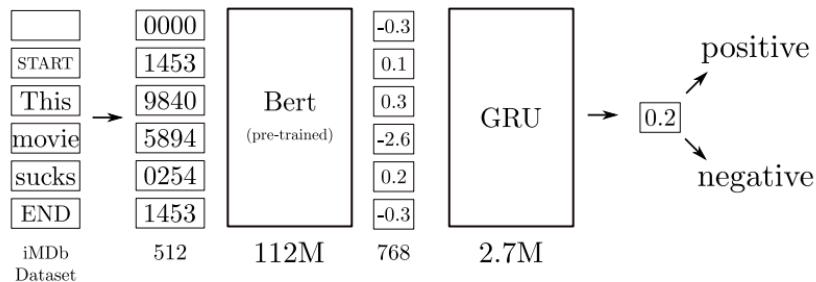
### Phrase-level masking

			curse	in	Elvish,	Entish,	or					for	this	treachery.
--	--	--	-------	----	---------	---------	----	--	--	--	--	-----	------	------------

# SENTIMENT ANALYSIS WITH TRANSFORMERS AND GRU

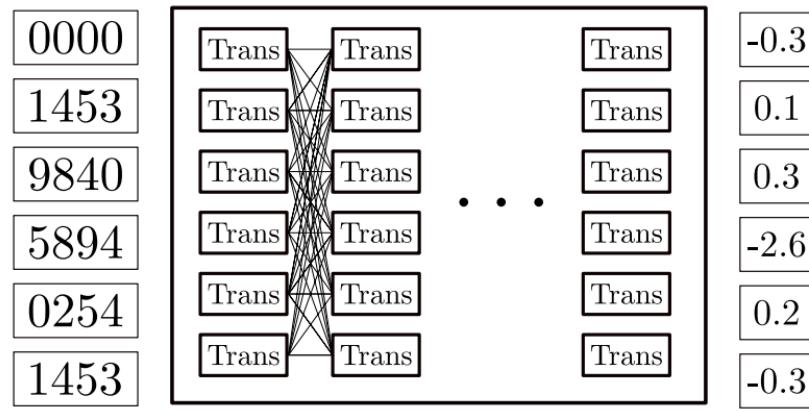
## SENTIMENT ANALYSIS

Bert can be fine-tuned for sentiment analysis by adding a classifier layer on top of the pretrained Bert model. The layer is trained on a labeled sentiment analysis dataset to predict the sentiment of a given input sequence, which can be a sentence, paragraph, or document. Fine-tuning the model allows it to learn the specific nuances of sentiment in the target task and produce improved results compared to using the pre-trained model alone.



# SENTIMENT ANALYSIS WITH TRANSFORMERS AND GRU

## SENTIMENT ANALYSIS

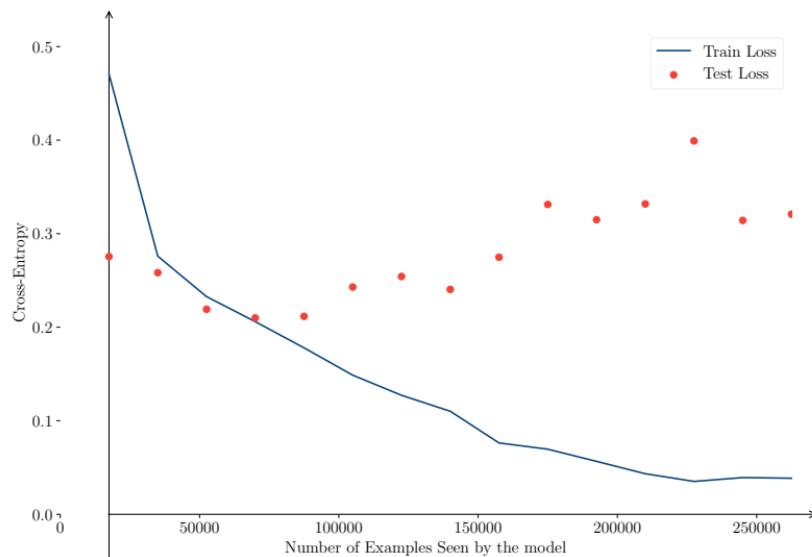


Bert

(pre-trained)

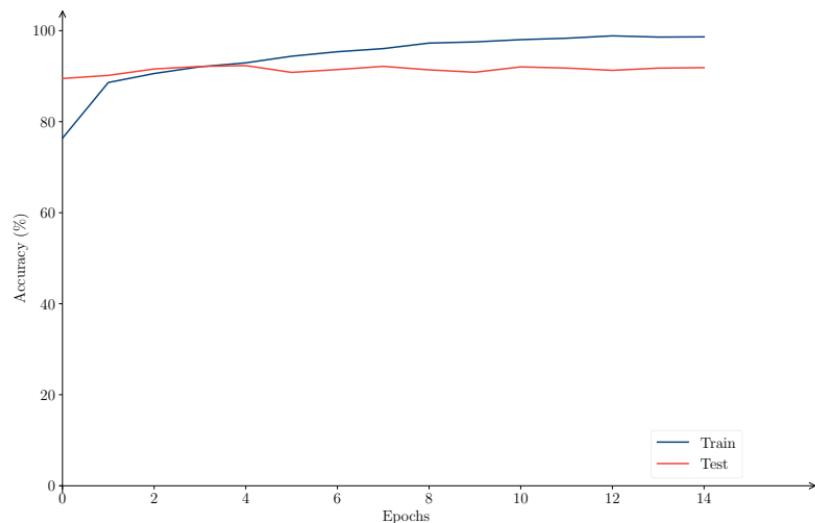
# SENTIMENT ANALYSIS WITH TRANSFORMERS AND GRU

## SENTIMENT ANALYSIS



# SENTIMENT ANALYSIS WITH TRANSFORMERS AND GRU

## SENTIMENT ANALYSIS



# SENTIMENT ANALYSIS WITH TRANSFORMERS AND GRU

## SENTIMENT ANALYSIS

```
In [117]: predict_sentiment(model, tokenizer, "This film is terrible")
Out[117]: 0.028073227033019066

In [118]: predict_sentiment(model, tokenizer, "This film is great")
Out[118]: 0.9749133586883545

In [119]: predict_sentiment(model, tokenizer, "This lecture was quite boring")
Out[119]: 0.0537508986890316

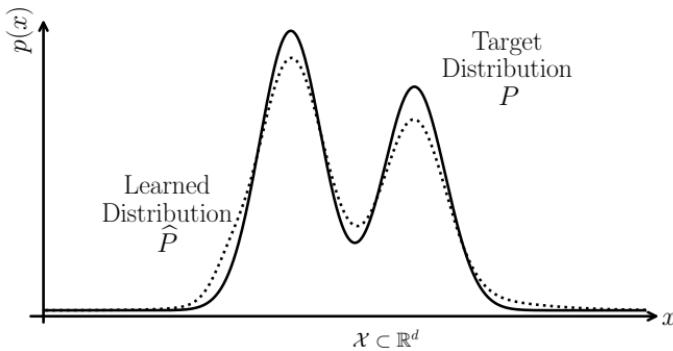
In [120]: predict_sentiment(model, tokenizer, "This lecture was challenging")
Out[120]: 0.6184227466583252

In [121]: predict_sentiment(model, tokenizer, "This lecture was dense but interesting")
Out[121]: 0.7105910181999207
```

# DENSITY ESTIMATION WITH NORMALIZING FLOWS

## ESTIMATING DENSITY

In Machine Learning, the task of *density estimation* consists in estimating the probability density function of a random variable from a set of observations. This is a fundamental problem in statistics and machine learning, with applications in a wide range of fields, including anomaly detection, clustering, and generative modeling. There are several methods for density estimation, including parametric models, non-parametric models, and deep learning models. In this section, we will focus on a specific type of deep learning model called Normalizing Flows, which is used for density estimation and generative modeling.



# DENSITY ESTIMATION WITH NORMALIZING FLOWS

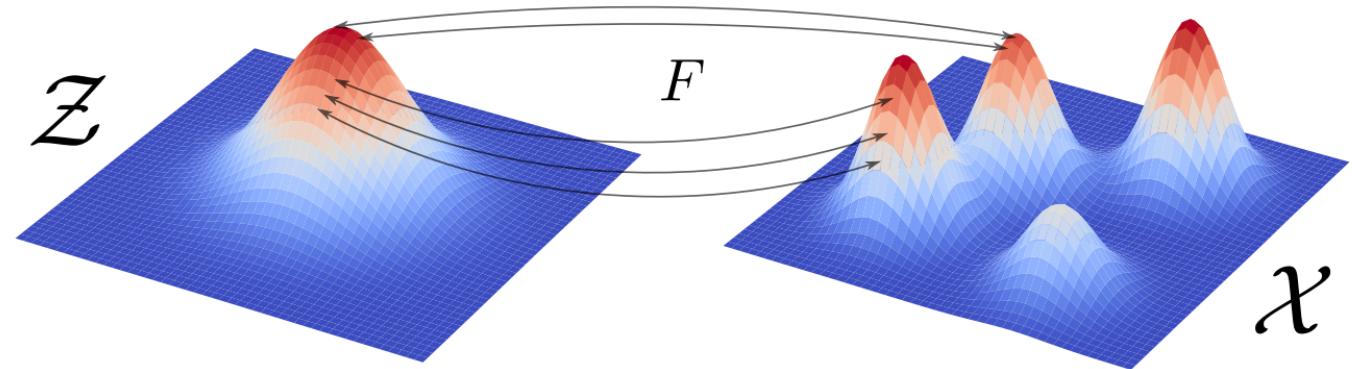
## OVERVIEW

A Normalizing Flow is usually seen as:

- ▶ a *generative model*,
- ▶ a *bijection mapping*,
- ▶ an *invertible neural network*,
- ▶ a *density estimator*.

# DENSITY ESTIMATION WITH NORMALIZING FLOWS

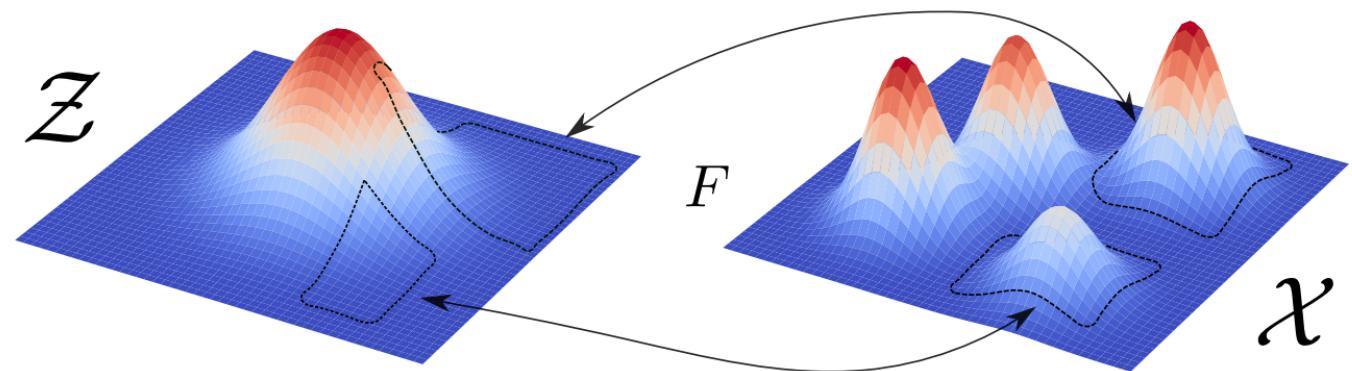
## OVERVIEW



**Figure.** A mapping between two probability distributions  
Point to point

# DENSITY ESTIMATION WITH NORMALIZING FLOWS

## OVERVIEW



**Figure.** A mapping between two probability distributions  
Subset to subset

# DENSITY ESTIMATION WITH NORMALIZING FLOWS

## MATHEMATICAL FRAMEWORK

### Normalizing Flow

A Normalizing Flow is a bijective function between a data space  $\mathcal{X}$  and a latent space  $\mathcal{Z}$ , both subset of  $\mathbb{R}^d$ .

$$\begin{aligned} F : \quad \mathcal{X} &\longmapsto \mathcal{Z} \\ x &\longmapsto z = F(x) \end{aligned}$$

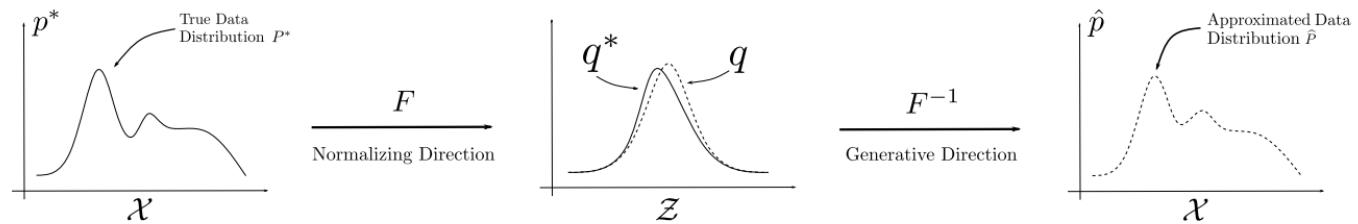
### Data and Latent Distributions

In theory, a NF maps a target distribution  $P$ , ie the data distribution to a simple latent distribution  $Q$ . Usually,  $Q$  is set to be a Normal Gaussian multivariate distribution  $\mathcal{N}(0_d, I_d)$ .  $p$  and  $q$  are respectively the probability densities of  $P$  and  $Q$ .

# DENSITY ESTIMATION WITH NORMALIZING FLOWS

## HOW DOES IT WORK ?

In practice, the mapping is *not perfect*.  $P^*$  induces a distribution  $Q$  and similarly, the latent distribution  $Q$  induces  $\hat{P}$ , which is the learned distribution. The forward pass  $F$  is called the *Normalizing* direction while the inverse pass  $F^{-1}$  is called the *Generative* direction.



**Figure.** 1D Normalizing Flow process.

# DENSITY ESTIMATION WITH NORMALIZING FLOWS

## INDUCED PROBABILITIES ?

### Change of Variable Formula

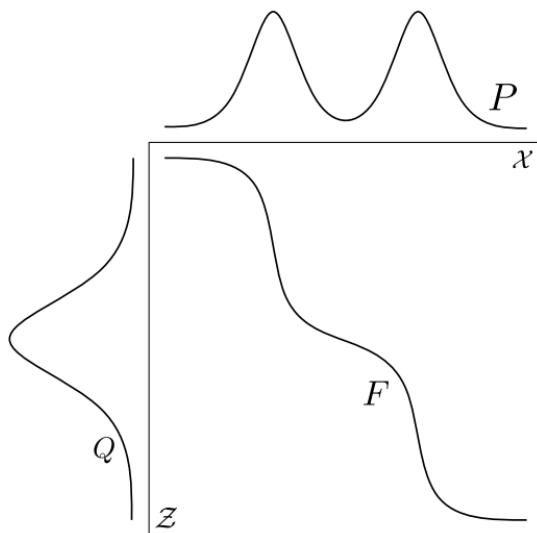
For a bijective and continuous function  $F$  and a latent distribution  $Q$ , the distribution induced by  $Q$  and  $F$  is defined through the *change of variable formula*:

$$\forall x \in \mathcal{X}, \quad \hat{p}(x) = |\det \text{Jac}_F(x)| q(F(x)). \quad (1)$$

# DENSITY ESTIMATION WITH NORMALIZING FLOWS

## INDUCED PROBABILITIES ?

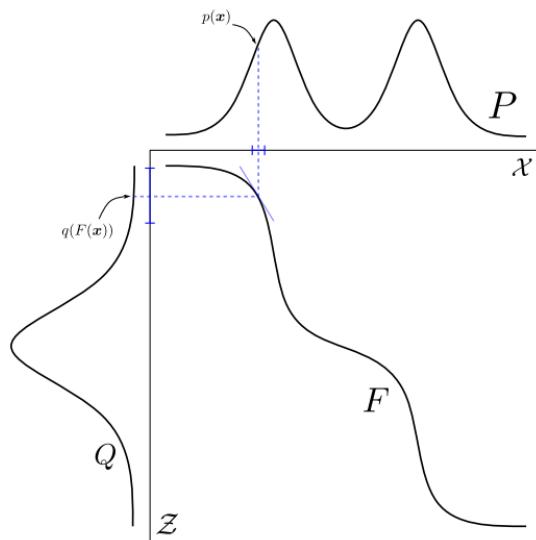
$$\forall x \in \mathcal{X}, \quad \hat{p}(x) = |\det \text{Jac}_F(x)| q(F(x)).$$



# DENSITY ESTIMATION WITH NORMALIZING FLOWS

## INDUCED PROBABILITIES ?

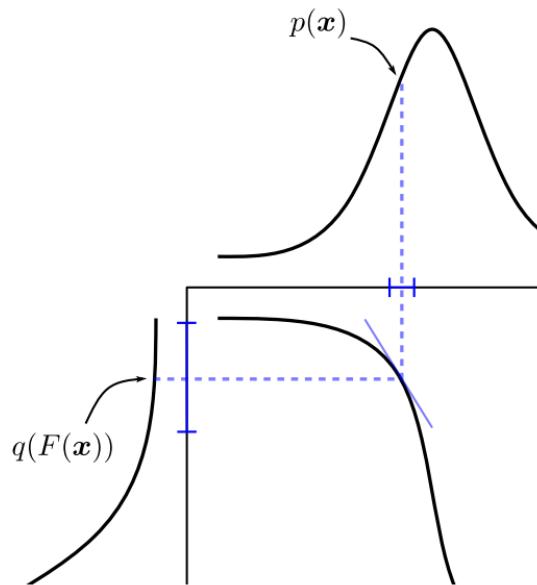
$$\forall x \in \mathcal{X}, \quad \hat{p}(x) = |\det \text{Jac}_F(x)| q(F(x)).$$



# DENSITY ESTIMATION WITH NORMALIZING FLOWS

## INDUCED PROBABILITIES ?

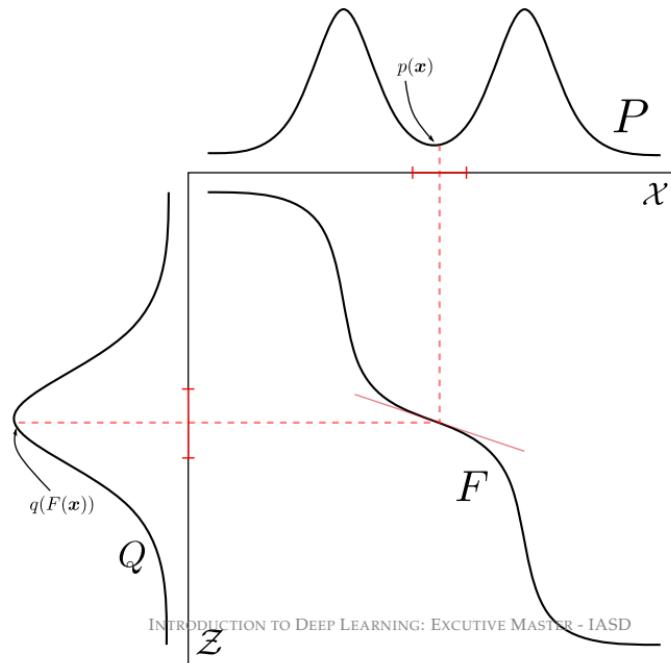
$$\forall x \in \mathcal{X}, \quad \hat{p}(x) = |\det \text{Jac}_F(x)| q(F(x)).$$



# DENSITY ESTIMATION WITH NORMALIZING FLOWS

## INDUCED PROBABILITIES ?

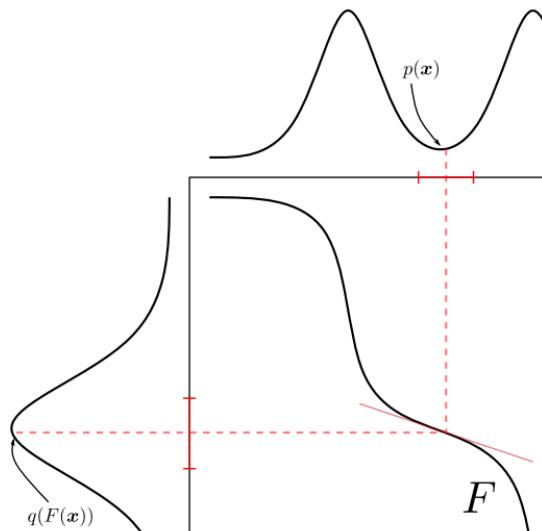
$$\forall x \in \mathcal{X}, \quad \hat{p}(x) = |\det \text{Jac}_F(x)| q(F(x)).$$



# DENSITY ESTIMATION WITH NORMALIZING FLOWS

## INDUCED PROBABILITIES ?

$$\forall x \in \mathcal{X}, \quad \hat{p}(x) = |\det \text{Jac}_F(x)| q(F(x)).$$

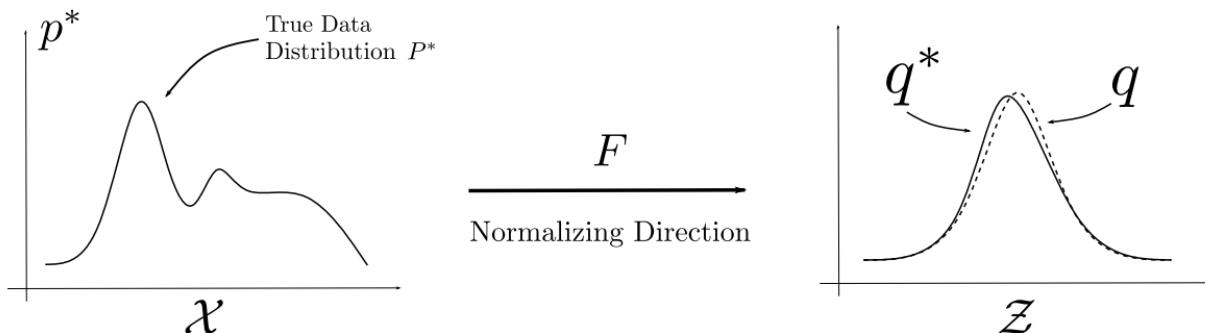


# DENSITY ESTIMATION WITH NORMALIZING FLOWS

## DENSITY ESTIMATION

To perform density estimation:

1. Draw  $x \sim P^*$ ,
2. Compute  $F(x)$  and  $|\det \text{Jac}_F(x)|$ ,
3. Compute  $\hat{p}(x) = q(F(x))|\det \text{Jac}_F(x)|$ .



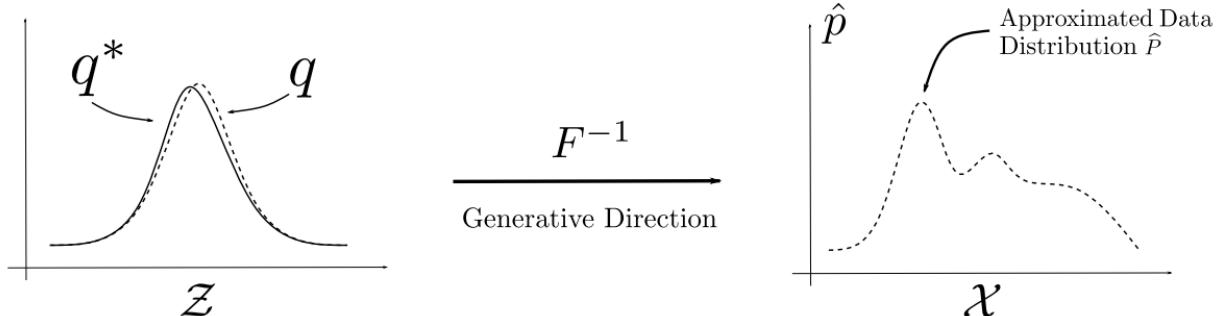
**Figure.** 1D Normalizing Process of Density Estimation.

# DENSITY ESTIMATION WITH NORMALIZING FLOWS

## DATA GENERATION

To perform data generation:

1. Draw  $z \sim Q$ ,
2. Compute  $x = F^{-1}(z)$ .



**Figure.** 1D Normalizing Flow process of Generation.

# DENSITY ESTIMATION WITH NORMALIZING FLOWS

## LOG-LIKELIHOOD

### Loss

The objective is to approximate  $P^*$  with  $\hat{P}$ . We can minimize the Kullback-Leiber Divergence :

$$\theta = \arg \min_{\theta} \mathcal{D}_{\text{KL}}(P^* \parallel \hat{P}).$$

This is equivalent to maximizing the log likelihood :

$$\theta = \arg \max_{\theta} \mathbb{E}_{x \sim \mathcal{X}} [\log \hat{p}(x)].$$

# DENSITY ESTIMATION WITH NORMALIZING FLOWS

## LOG-LIKELIHOOD

$$\mathcal{D}_{\text{KL}}(P^* \parallel \hat{P}) = \int_{\mathcal{X}} p^*(x) \log \left( \frac{p^*(x)}{\hat{p}(x)} \right) dx$$

$$\text{nll} = -\mathbb{E}_{x \sim \mathcal{X}} [\log \hat{p}(x)]$$

# DENSITY ESTIMATION WITH NORMALIZING FLOWS

## LEARNING STEPS

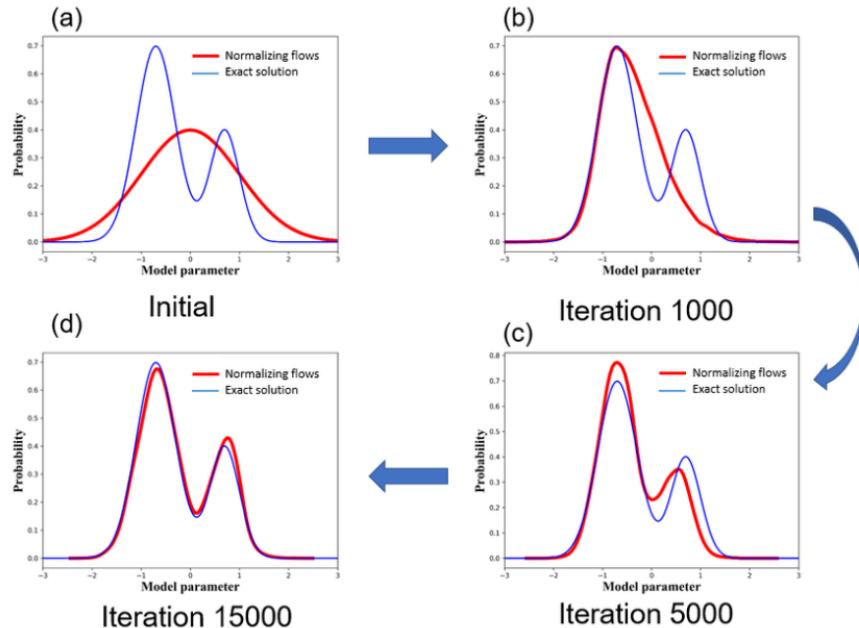


Figure. Learning Process for a 1D Normalizing Flow.

# IMAGE SEGMENTATION WITH U-NET

## IMAGE SEGMENTATION

Image segmentation is the process of partitioning an image into multiple segments or regions based on the characteristics of the pixels. It is a fundamental task in computer vision and has applications in various fields, including medical imaging, autonomous driving, and satellite image analysis. There are several methods for image segmentation, including thresholding, clustering, and deep learning-based approaches. In this section, we will focus on a deep learning model called U-Net, which is commonly used for image segmentation tasks.



# IMAGE SEGMENTATION WITH U-NET

## APPLICATIONS OF IMAGE SEGMENTATION

Image segmentation has a wide range of applications in computer vision and image processing, including:

- ▶ Medical Imaging: Segmentation of organs, tumors, and other structures in medical images.
- ▶ Autonomous Driving: Segmentation of objects such as cars, pedestrians, and road signs in images captured by autonomous vehicles.
- ▶ Satellite Image Analysis: Segmentation of land cover types, buildings, and other features in satellite images.
- ▶ Object Detection: Segmentation of objects in images to localize and classify them.
- ▶ Image Editing: Segmentation of objects for image editing tasks such as background removal and image compositing.

# IMAGE SEGMENTATION WITH U-NET

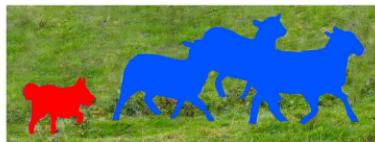
## TYPE OF SEGMENTATIONS

Image segmentation can be broadly classified into two types: semantic segmentation and instance segmentation.

- ▶ **Semantic Segmentation:** Semantic segmentation assigns a class label to each pixel in an image, such as road, car, person, etc. The goal is to partition the image into semantically meaningful regions.
- ▶ **Instance Segmentation:** Instance segmentation goes a step further than semantic segmentation by distinguishing between different instances of the same class. It assigns a unique label to each object instance in the image.



Original Image



Semantic Segmentation

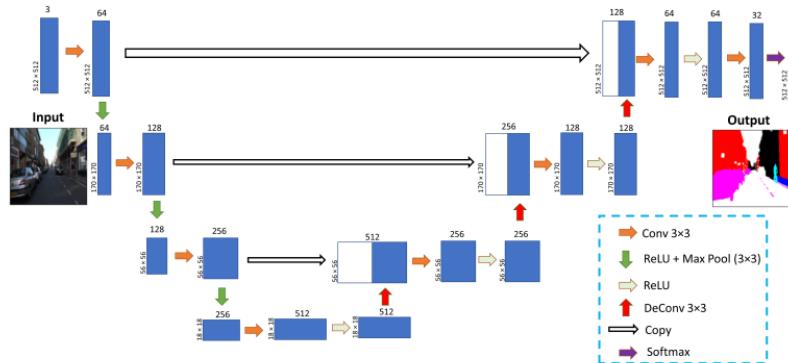


Instance Segmentation

# IMAGE SEGMENTATION WITH U-NET

## U-NET ARCHITECTURE

U-Net is a convolutional neural network architecture designed for image segmentation tasks. It consists of an encoder-decoder structure with skip connections that allow the model to capture both local and global features in the input image.



# IMAGE SEGMENTATION WITH U-NET

## LOSS FUNCTIONS

3 main loss functions are used for image segmentation tasks:

- ▶ **Categorical Cross-Entropy Loss:** Used for multi-class segmentation tasks where each pixel is classified into one of multiple classes.
- ▶ **Dice Loss:** A similarity-based loss function that measures the overlap between the predicted segmentation mask and the ground truth mask.
- ▶ **Shaped aware Loss:** A loss function that penalizes the model for making errors near the object boundaries, where the segmentation is most critical.

# IMAGE SEGMENTATION WITH U-NET

## LOSS FUNCTIONS

3 main loss functions are used for image segmentation tasks:

- ▶ **Categorical Cross-Entropy Loss:** Used for multi-class segmentation tasks where each pixel is classified into one of multiple classes.

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N \sum_{c=1}^C y_{i,c} \log(\hat{y}_{i,c})$$

- ▶ **Dice Loss:** A similarity-based loss function that measures the overlap between the predicted segmentation mask and the ground truth mask.
- ▶ **Shaped aware Loss:** A loss function that penalizes the model for making errors near the object boundaries, where the segmentation is most critical.

# IMAGE SEGMENTATION WITH U-NET

## LOSS FUNCTIONS

3 main loss functions are used for image segmentation tasks:

- ▶ **Categorical Cross-Entropy Loss:** Used for multi-class segmentation tasks where each pixel is classified into one of multiple classes.
- ▶ **Dice Loss:** A similarity-based loss function that measures the overlap between the predicted segmentation mask and the ground truth mask.

$$\mathcal{L} = 1 - \frac{2 \sum_{i=1}^N y_i \hat{y}_i}{\sum_{i=1}^N y_i + \sum_{i=1}^N \hat{y}_i}$$

- ▶ **Shaped aware Loss:** A loss function that penalizes the model for making errors near the object boundaries, where the segmentation is most critical.

# IMAGE SEGMENTATION WITH U-NET

## LOSS FUNCTIONS

3 main loss functions are used for image segmentation tasks:

- ▶ **Categorical Cross-Entropy Loss:** Used for multi-class segmentation tasks where each pixel is classified into one of multiple classes.
- ▶ **Dice Loss:** A similarity-based loss function that measures the overlap between the predicted segmentation mask and the ground truth mask.
- ▶ **Shaped aware Loss:** A loss function that penalizes the model for making errors near the object boundaries, where the segmentation is most critical.

$$\mathcal{L} = -w(x) \left[ \frac{1}{N} \sum_{i=1}^N \sum_{c=1}^C y_{i,c} \log(\hat{y}_{i,c}) \right]$$

with  $w(x)$  a weight function that assigns higher weights to pixels near the object boundaries.

## Part V

# YOUR TURN: HANDS-ON DEEP LEARNING

# BUILD AND USE AN AUTOENCODER

## FORMAL INTRODUCTION OF AN AUTOENCODER

### Definition

An autoencoder is a type of artificial neural network used to learn efficient codings of unlabeled data. It consists of two main components:

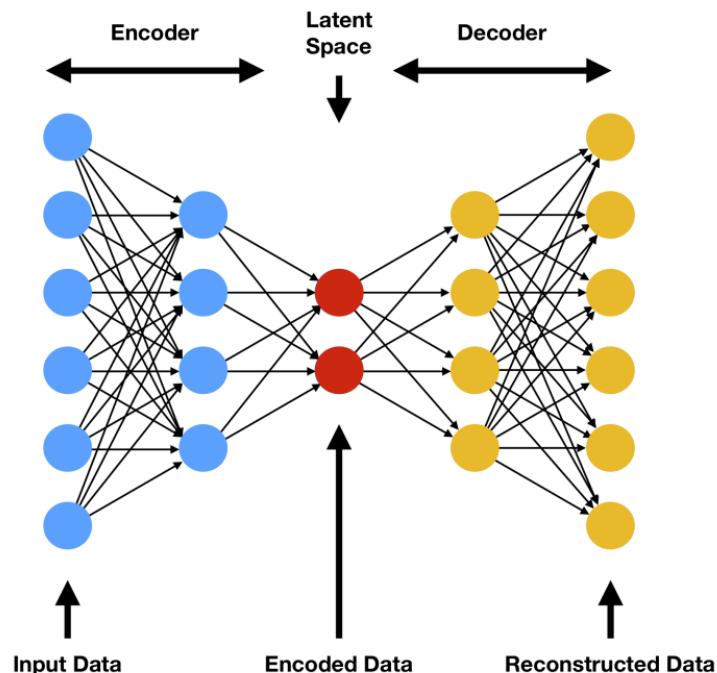
- ▶ An **encoder** function:  $\text{encoder}(x) : \mathbb{R}^d \rightarrow \mathbb{R}^m$   
Maps an input  $x$  from the input space  $\mathbb{R}^d$  to a hidden representation space  $\mathbb{R}^m$ .
- ▶ A **decoder** function:  $\text{decoder}(z) : \mathbb{R}^m \rightarrow \mathbb{R}^d$   
Maps the hidden representation  $z$  back to the original input space  $\mathbb{R}^d$ .

### Goal

The primary goal of an autoencoder is to learn a representation (encoding) for a set of data, typically for the purpose of dimensionality reduction or feature learning. Through training, the autoencoder learns to compress the data from  $\mathbb{R}^d$  to  $\mathbb{R}^m$  (where  $m < d$ ) and then reconstruct the data back to  $\mathbb{R}^d$  as accurately as possible. This process forces the autoencoder to capture the most important features of the data in the hidden representation  $z$ .

# BUILD AND USE AN AUTOENCODER

## FORMAL INTRODUCTION OF AN AUTOENCODER



# BUILD AND USE AN AUTOENCODER

## AUTOENCODERS AND UNSUPERVISED LEARNING

### Unsupervised Learning

Autoencoders are a classic example of [unsupervised learning](#). In unsupervised learning, the goal is to learn patterns from unlabelled data. Autoencoders learn to compress and decompress the input data without any explicit labels, aiming to capture the underlying structure of the data.

### Objective Function

The learning process of an autoencoder is guided by the minimization of a loss function, typically involving a norm that measures the difference between the input and the reconstructed output. Formally, the objective is to minimize:

$$\min_{\theta} \mathbb{E}_{x \sim P} [l(x, \text{decoder}_{\theta}(\text{encoder}_{\theta}(x)))]$$

where  $x$  is the input data,  $\theta$  represents the parameters of the encoder and decoder, and  $l$  is a loss function.

# BUILD AND USE AN AUTOENCODER

## AUTOENCODERS AND UNSUPERVISED LEARNING

### Unsupervised Learning

Autoencoders are a classic example of [unsupervised learning](#). In unsupervised learning, the goal is to learn patterns from unlabelled data. Autoencoders learn to compress and decompress the input data without any explicit labels, aiming to capture the underlying structure of the data.

### Objective Function

The learning process of an autoencoder is guided by the minimization of a loss function, typically involving a norm that measures the difference between the input and the reconstructed output. Formally, the objective is to minimize:

$$\min_{\theta} \mathbb{E}_{x \sim P} \left[ \|x - \hat{x}\|_2^2 \right]$$

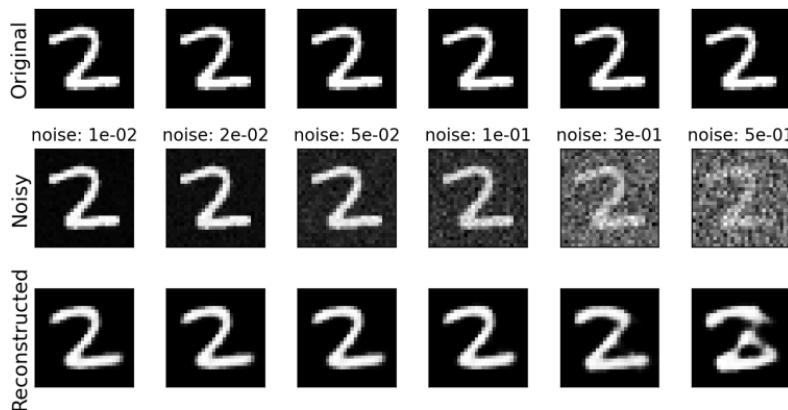
where  $x$  is the input data,  $\theta$  represents the parameters of the encoder and decoder and  $\hat{x} = \text{decoder}_\theta(\text{encoder}_\theta(x))$  in the reconstruction.

## APPLICATIONS OF AUTOENCODERS

- ▶ **Reduce the size of the data to transfer:** Autoencoders can compress data into a lower-dimensional space, facilitating faster data transfer by reducing the amount of data that needs to be transmitted.

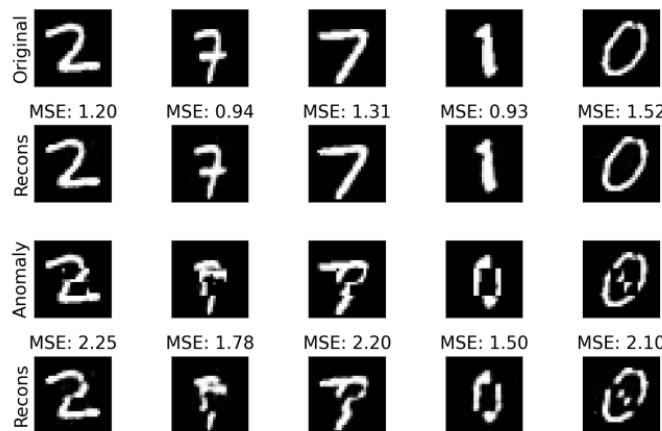
## APPLICATIONS OF AUTOENCODERS

- ▶ **Reduce the size of the data to transfer:** Autoencoders can compress data into a lower-dimensional space, facilitating faster data transfer by reducing the amount of data that needs to be transmitted.
- ▶ **Denoise image:** By learning to ignore the "noise" in the input data, autoencoders can reconstruct cleaner versions of noisy images, effectively removing the noise.



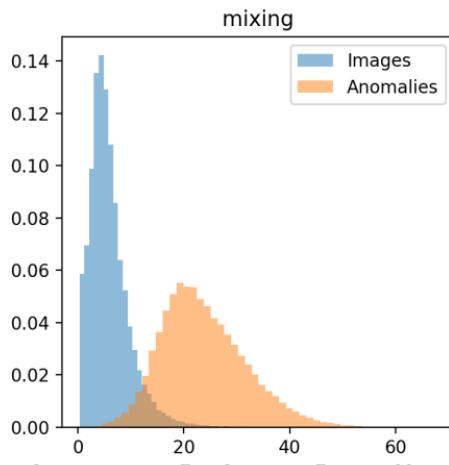
## APPLICATIONS OF AUTOENCODERS

- ▶ **Reduce the size of the data to transfer:** Autoencoders can compress data into a lower-dimensional space, facilitating faster data transfer by reducing the amount of data that needs to be transmitted.
- ▶ **Denoise image:** By learning to ignore the "noise" in the input data, autoencoders can reconstruct cleaner versions of noisy images, effectively removing the noise.
- ▶ **Anomaly detection:** Autoencoders can learn the normal patterns within data. Deviations from these patterns, when the reconstruction error is high, can indicate anomalies or outliers in the data.



## APPLICATIONS OF AUTOENCODERS

- ▶ **Reduce the size of the data to transfer:** Autoencoders can compress data into a lower-dimensional space, facilitating faster data transfer by reducing the amount of data that needs to be transmitted.
- ▶ **Denoise image:** By learning to ignore the "noise" in the input data, autoencoders can reconstruct cleaner versions of noisy images, effectively removing the noise.
- ▶ **Anomaly detection:** Autoencoders can learn the normal patterns within data. Deviations from these patterns, when the reconstruction error is high, can indicate anomalies or outliers in the data.



## BUILD AND USE AN AUTOENCODER

YOUR TURN !

Get the Notebook on the course website and start working on it.

- ▶ <https://www.alexverine.com>
- ▶ Teaching
- ▶ Introduction to Deep Learning
- ▶ Lien Notebook

# IMAGE SEGMENTATION WITH U-NET: CORRECTLY EVALUATE THE MODEL

## IMAGE SEGMENTATION

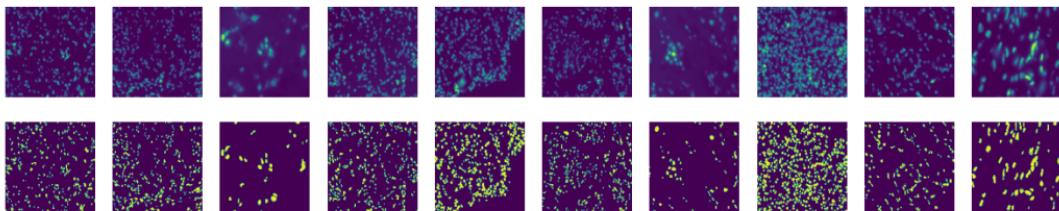
Image segmentation is the process of partitioning an image into multiple segments or regions based on certain characteristics. It is a fundamental task in computer vision and has numerous applications, such as object detection, image analysis, and medical imaging. For instance, it can be used for:

- ▶ Identifying and delineating objects in an image.
- ▶ Extracting features from images for further analysis.
- ▶ Medical image analysis, such as identifying tumors or lesions.
- ▶ Autonomous driving, where segmentation is used to detect and classify objects on the road.
- ▶ Image editing and manipulation, such as background removal or object replacement.

# IMAGE SEGMENTATION WITH U-NET: CORRECTLY EVALUATE THE MODEL

## THE DATASET

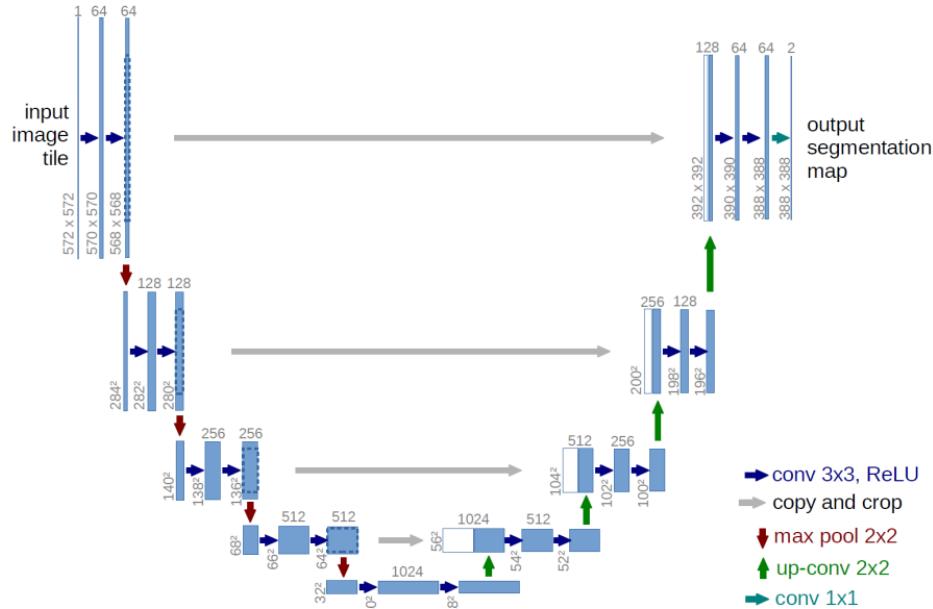
The dataset presented here is hosted in the BioStudies database under accession number S-BSST265 18. It consists of 79 fluorescence images of immuno and DAPI stained samples containing 7813 nuclei in total. The images are derived from one Schwann cell stroma-rich tissue (from a ganglioneuroblastoma) cryosection (10 images/2773 nuclei), seven neuroblastoma (NB) patients (19 images/931 nuclei), one Wilms patient (1 image/102 nuclei), two NB cell lines (CLB-Ma, STA-NB10) (8 images/1785 nuclei) and a human keratinocyte cell line (HaCaT) (41 images/2222 nuclei). In addition, the dataset is heterogenous in aspects such as type of preparation, imaging modality, magnification, signal-to-noise ratio and other technical aspects.



# IMAGE SEGMENTATION WITH U-NET: CORRECTLY EVALUATE THE MODEL

## THE U-NET ARCHITECTURE

U-Net is a convolutional neural network architecture designed for encoding/decoding images. It is widely used for image segmentation tasks:



## IMAGE SEGMENTATION WITH U-NET: CORRECTLY EVALUATE THE MODEL

### EVALUATING THE MODELS

The model is traditionally trained with the Binary Cross-Entropy loss function. Each pixel  $y_{i,j}$  is a label in  $\{0, 1\}$ , where 0 represents the background and 1 represents the object. The loss function is defined as:

$$\text{BCE}(y, \hat{y}) = -\frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N y_{i,j} \log(\hat{y}_{i,j}) + (1 - y_{i,j}) \log(1 - \hat{y}_{i,j})$$

where  $N$  is the number of pixels in the image,  $y$  is the ground truth mask, and  $\hat{y}$  is the predicted mask. The BCE loss function is used to measure the pixel-wise binary classification error between the ground truth and the predicted mask.

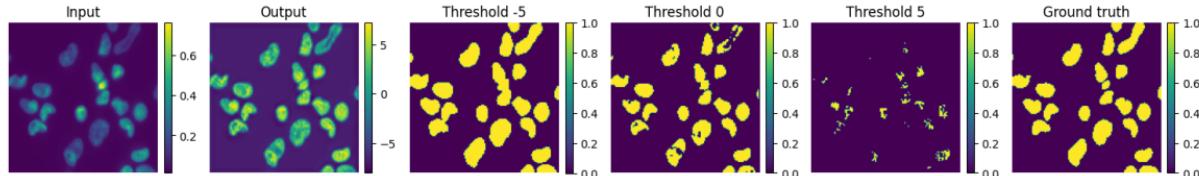
What can we use to evaluate the performance of the model?

# IMAGE SEGMENTATION WITH U-NET: CORRECTLY EVALUATE THE MODEL

## INTERPRETING THE OUTPUT

The first metrics that comes to mind is cross-entropy on the test set. However, it is not the best metric to evaluate the performance of a segmentation model. The reason is that the model predicts probabilities for each pixel, and the thresholding of these probabilities is not taken into account in the BCE loss.

The model predicts values in  $\mathbb{R}$  for each pixel. We can use a threshold to binarize the predictions:

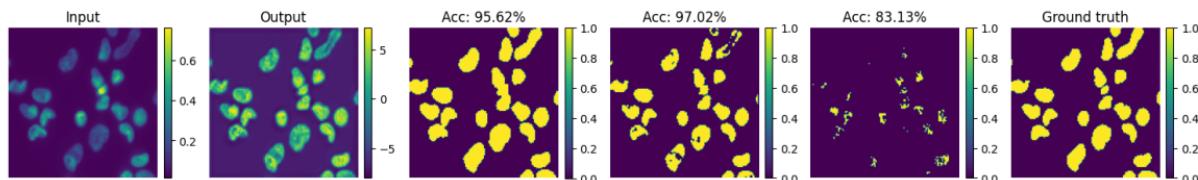


# IMAGE SEGMENTATION WITH U-NET: CORRECTLY EVALUATE THE MODEL

## THE ACCURACY

The accuracy is a common metric used to evaluate the performance of classification models. The accuracy is defined as:

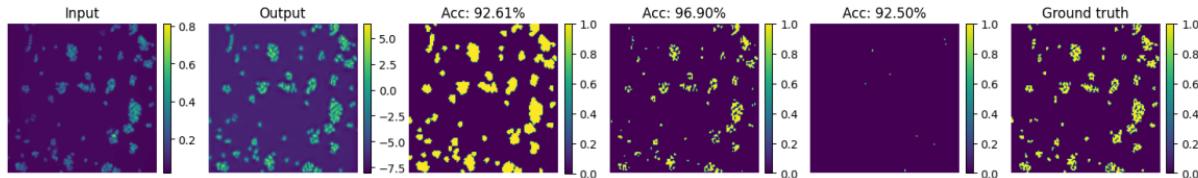
$$\text{Accuracy} = \frac{\text{Number of correct predictions}}{\text{Total number of predictions}}$$



# IMAGE SEGMENTATION WITH U-NET: CORRECTLY EVALUATE THE MODEL THE ACCURACY

The accuracy is a common metric used to evaluate the performance of classification models. The accuracy is defined as:

$$\text{Accuracy} = \frac{\text{Number of correct predictions}}{\text{Total number of predictions}}$$



However, missing a label 0 is not the same as missing a label 1 !

# IMAGE SEGMENTATION WITH U-NET: CORRECTLY EVALUATE THE MODEL

## CONFUSION MATRIX

The confusion matrix is a useful tool for evaluating classification models, particularly for segmentation. It helps to distinguish between different types of errors.

- ▶ **True Positives (TP):**  
Correctly predicted positives.
- ▶ **False Positives (FP):**  
Incorrectly predicted positives.
- ▶ **False Negatives (FN):**  
Missed positives.
- ▶ **True Negatives (TN):**  
Correctly predicted negatives.

	Positive	Negative
Predicted Positive	TP	FP
Predicted Negative	FN	TN

# IMAGE SEGMENTATION WITH U-NET: CORRECTLY EVALUATE THE MODEL

## PRECISION AND RECALL

	Positive	Negative
Predicted Positive	TP	FP
Predicted Negative	FN	TN

To evaluate the performance of a segmentation model, we can use the following metrics:

$$\text{Accuracy} = \frac{TP + TN}{P + N}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{P}$$

# IMAGE SEGMENTATION WITH U-NET: CORRECTLY EVALUATE THE MODEL

## CONFUSION MATRIX

Examples with a cancer detection model with 100 patients (20 with cancer and 80 without cancer):

		Cancer (P)	Non-Cancer (N)
Predicted Cancer (PP)	Predicted Cancer (PP)	10 (TP)	1 (FP)
	Predicted Non-Cancer (PN)	10 (FN)	79 (TN)

$$\text{Accuracy: } \frac{TP+TN}{P+N} = \frac{10+79}{20+80} = 0.89 \quad \text{Precision: } \frac{TP}{TP+FP} = \frac{10}{10+1} = 0.90 \quad \text{Recall: } \frac{TP}{P} = \frac{10}{20} = 0.5$$

# IMAGE SEGMENTATION WITH U-NET: CORRECTLY EVALUATE THE MODEL

## CONFUSION MATRIX

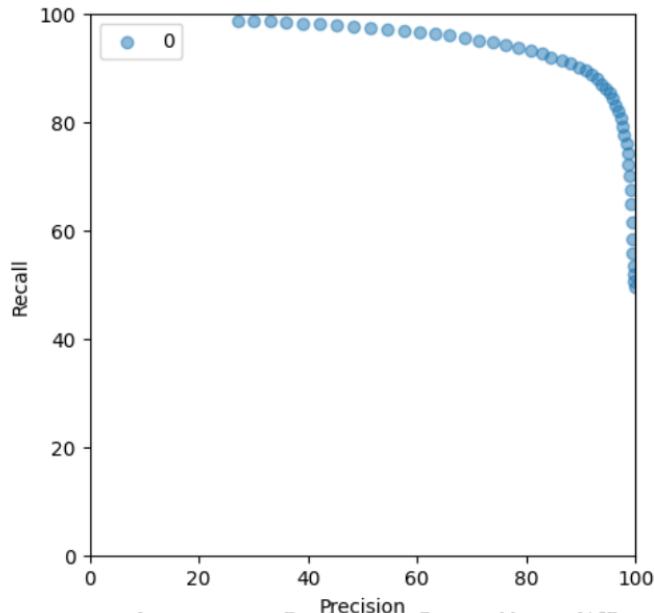
Examples with a financial operation models with 100 operations (20 high reward and 80 high loss):

	Reward (P)	Loss (N)
Predicted Reward (PP)	19 (TP)	15 (FP)
Predicted Loss (PN)	1 (FN)	65 (TN)

$$\text{Accuracy: } \frac{TP+TN}{P+N} = \frac{19+65}{20+80} = 0.84 \quad \text{Precision: } \frac{TP}{TP+FP} = \frac{19}{19+15} = 0.56 \quad \text{Recall: } \frac{TP}{P} = \frac{19}{20} = 0.95$$

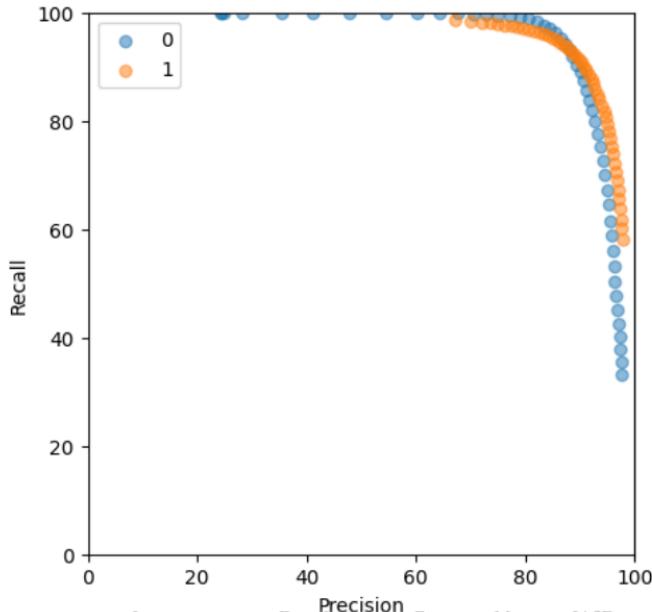
## IMAGE SEGMENTATION WITH U-NET: CORRECTLY EVALUATE THE MODEL PRECISION-RECALL CURVE

In Image segmentation tasks, by adjusting the threshold we can change the precision and recall of the model. The Precision-Recall curve is a useful tool to evaluate the performance of a model across different thresholds:



## IMAGE SEGMENTATION WITH U-NET: CORRECTLY EVALUATE THE MODEL PRECISION-RECALL CURVE

In Image segmentation tasks, by adjusting the threshold we can change the precision and recall of the model. The Precision-Recall curve is a useful tool to evaluate the performance of a model across different thresholds:



# IMAGE SEGMENTATION WITH U-NET: CORRECTLY EVALUATE THE MODEL

## YOUR TURN !

Get the Notebook on the course website and start working on it.

- ▶ <https://www.alexverine.com>
- ▶ Teaching
- ▶ Introduction to Deep Learning
- ▶ Lien Notebook

During this hands-on session we will compare the speed of the model with and without the GPU, and we will evaluate the model with the metrics we have seen. Don't hesitate to use Colab or your own computer to run the notebook.