

# **Deep Learning Methods For NLP**

**Pierre Colombo**

[pierre.colombo@centralesupelec.fr](mailto:pierre.colombo@centralesupelec.fr)

**MICS - CentraleSupélec**

**Advanced Natural Language Processing**



# Final Project & Presentations

# Lectures Outline

1. The Basics of Natural Language Processing

**2. Representing Text with Vectors**

3. Deep Learning Methods for NLP

4. Language Modeling

.....

# Today Lecture Outline

- Deep Learning Framework
- The Multi-Layer Perceptron
- Recurrent Neural Network
- Attention Mechanism
- Self-Attention Mechanism and the Transformer Architecture

# Motivations

So far, we have seen, **techniques to represent tokens with vectors**

Given a certain representations of tokens:

→ **How can we model a sequence of tokens to perform a specific task?**

In the past 10 years, a “new” class of machine learning techniques has become very popular and successful in NLP: **Deep Learning**

*In this session, we introduce Deep Learning with a focus on the methods used in NLP*

# Framework

We want to model  $(X_1, \dots, X_T)$  i.e. find the correct label  $Y$

$$dnn : \mathbb{R}^{d,T} \rightarrow \mathbb{R}^p \quad or \quad [0, |K|]^p$$

$$(X_1, \dots, X_T) \rightarrow \hat{Y}$$

- Output space is  $\mathbb{R}^p$  for **Regression** tasks
- Output space is  $[0, |K|]^p$  for **Classification** tasks

# Framework

We want to model  $(X_1, \dots, X_T)$  i.e. find the correct label

$$dnn : \mathbb{R}^{d,T} \rightarrow \mathbb{R}^p \quad or \quad [0, |K|]^p$$

$$(X_1, \dots, X_T) \rightarrow \hat{Y}$$

**Questions:** when we do Deep Learning...

- How do we **define**  $f_\theta$
- How do we **train.**  $f_\theta$  with data ?

# Framework

We want to model  $(X_1, \dots, X_T)$  i.e. find the correct label

$$\begin{aligned} dnn : \quad \mathbb{R}^{d,T} &\rightarrow \quad \mathbb{R}^p \\ (X_1, \dots, X_T) &\rightarrow \quad \hat{Y} \end{aligned}$$

Most Deep Learning Models (all the ones we will use in this course):

- are **parametric** (i.e.  $\theta \in \Theta = \mathbb{R}^D$ )
- defined as a **composition of “simple” functions (linear & non-linear)**
- are trained in an **end-to-end** fashion with **backpropagation**

NB: In Deep Learning, **the parametrization of  $f_\theta$**  is called **the Architecture**

# Different Types of Architecture

How can we define our predictive function  $dnn$

- Multi-Layer Perceptron
- Recurrent Layers
- Attention Layers
- Self-Attention Layers (in a Transformer Architecture)

# Different Types of Architecture

How can we define our predictive function *dnn*

- Multi-Layer Perceptron
- Recurrent Layers
- Attention Layers
- Self-Attention Layers (in a Transformer Architecture)

How do we **train our model** ? (i.e. estimate the parameters of the model)

- **Stochastic Gradient Descent** also called **backpropagation** in this context

# The MultiLayer Perceptron (MLP)

*aka “the Most simple Deep Learning Architecture”*

The **MLP** works **on unidimensional data** (e.g. dimension  $d$ )

We present the **MLP in the regression case** (e.g. output space is.  $\mathbb{R}^2$  ))

$$\begin{array}{ccc} dnn : & \mathbb{R}^d & \rightarrow & \mathbb{R}^2 \\ & X & \rightarrow & \hat{Y} \end{array}$$

# The MultiLayer Perceptron (MLP)

*aka “the Most simple Deep Learning Architecture”*

The **MLP** works **on unidimensional data** (e.g. dimension  $d$ )

We present the **MLP in the regression case** (e.g. output space is.  $\mathbb{R}^2$  ))

$$\begin{array}{ccc} f_{\theta} : & \mathbb{R}^d & \rightarrow \mathbb{R}^2 \\ & X & \rightarrow \hat{Y} \end{array}$$

# The MultiLayer Perceptron (MLP)

*aka “the Most simple Deep Learning Architecture”*

The **MLP** works **on unidimensional data** (e.g. dimension  $d$ )

We present the **MLP in the regression case** (e.g. output space is.  $\mathbb{R}^2$  ))

$$dnn_{(W_1, b_1, W_2, b_2)}(X) = W_2 \varphi_1(W_1 X + b_1) + b_2$$

# The MultiLayer Perceptron (MLP)

*aka “the Most simple Deep Learning Architecture”*

The **MLP** works **on unidimensional data** (e.g. dimension  $d$ )

We present the **MLP in the regression case** (e.g. output space is  $\mathbb{R}^2$  ))

$$dnn_{(W_1, b_1, W_2, b_2)}(X) = W_2 \varphi_1(W_1 X + b_1) + b_2$$

$W_1, b_1, W_2$  and  $b_2$  are trainable parameters.  $W_1 \in \mathbb{R}^{\delta \times d}$ ,  $b_1 \in \mathbb{R}^\delta$ ,  $W_2 \in \mathbb{R}^{2 \times \delta}$  and  $b_2 \in \mathbb{R}$

$\varphi_1$  is a fixed non-linear function,  $\varphi_1 : \mathbb{R}^d \rightarrow \mathbb{R}^\delta$

# The MultiLayer Perceptron (MLP)

The **MLP** works **on unidimensional data** (e.g. dimension  $d$ )

We present the **MLP in the regression case** (e.g. output space is  $\mathbb{R}^2$ )

$$dnn_{(W_1, b_1, W_2, b_2)}(X) = W_2 \varphi_1(W_1 X + b_1) + b_2$$

$W_1, b_1, W_2$  and  $b_2$  are trainable parameters.  $W_1 \in \mathbb{R}^{\delta \times d}$ ,  $b_1 \in \mathbb{R}^\delta$ ,  $W_2 \in \mathbb{R}^{2 \times \delta}$  and  $b_2 \in \mathbb{R}$

$\varphi_1$  is a fixed non-linear function,  $\varphi_1 : \mathbb{R}^d \rightarrow \mathbb{R}^\delta$

- This model is a **2-layer MLP** model

# The MultiLayer Perceptron (MLP)

The **MLP** works **on unidimensional data** (e.g. dimension  $d$ )

We present the **MLP in the regression case** (e.g. output space is  $\mathbb{R}^2$ )

$$dnn_{(W_1, b_1, W_2, b_2)}(X) = W_2 \varphi_1(W_1 X + b_1) + b_2$$

$W_1, b_1, W_2$  and  $b_2$  are trainable parameters.  $W_1 \in \mathbb{R}^{\delta \times d}$ ,  $b_1 \in \mathbb{R}^\delta$ ,  $W_2 \in \mathbb{R}^{2 \times \delta}$  and  $b_2 \in \mathbb{R}$

$\varphi_1$  is a fixed non-linear function,  $\varphi_1 : \mathbb{R}^d \rightarrow \mathbb{R}^\delta$

- This model is a ***2-layer MLP*** model
- With **1 *hidden layer*** of dimension  $\delta$

# The MultiLayer Perceptron (MLP)

The **MLP** works **on unidimensional data** (e.g. dimension  $d$ )

We present the **MLP in the regression case** (e.g. output space is  $\mathbb{R}^2$ )

$$dnn_{(W_1, b_1, W_2, b_2)}(X) = W_2 \varphi_1(W_1 X + b_1) + b_2$$

$W_1, b_1, W_2$  and  $b_2$  are trainable parameters.  $W_1 \in \mathbb{R}^{\delta \times d}$ ,  $b_1 \in \mathbb{R}^\delta$ ,  $W_2 \in \mathbb{R}^{2 \times \delta}$  and  $b_2 \in \mathbb{R}$

$\varphi_1$  is a fixed non-linear function,  $\varphi_1 : \mathbb{R}^d \rightarrow \mathbb{R}^\delta$

- This model is a **2-layer MLP** model
- With **1 hidden layer** of dimension  $\delta$
- Taking as input a vector of **dimension  $d$**  to output a vector of **dimension 2**

# The MultiLayer Perceptron (MLP)

The **MLP** works **on unidimensional data** (e.g. dimension  $d$ )

We present the **MLP in the regression case** (e.g. output space is  $\mathbb{R}^2$ )

$$dnn_{(W_1, b_1, W_2, b_2)}(X) = W_2 \varphi_1(W_1 X + b_1) + b_2$$

$W_1, b_1, W_2$  and  $b_2$  are trainable parameters.  $W_1 \in \mathbb{R}^{\delta \times d}$ ,  $b_1 \in \mathbb{R}^\delta$ ,  $W_2 \in \mathbb{R}^{2 \times \delta}$  and  $b_2 \in \mathbb{R}$

$\varphi_1$  is a fixed non-linear function,  $\varphi_1 : \mathbb{R}^d \rightarrow \mathbb{R}^\delta$

- This model is a **2-layer MLP** model
- With **1 hidden layer** of dimension  $\delta$
- Taking as input a vector of **dimension  $d$**  to output a vector of **dimension 2**
- Such a model is also referred to as a **Feed-Forward Neural Network (FNN)**

# The MultiLayer Perceptron: Diagram View

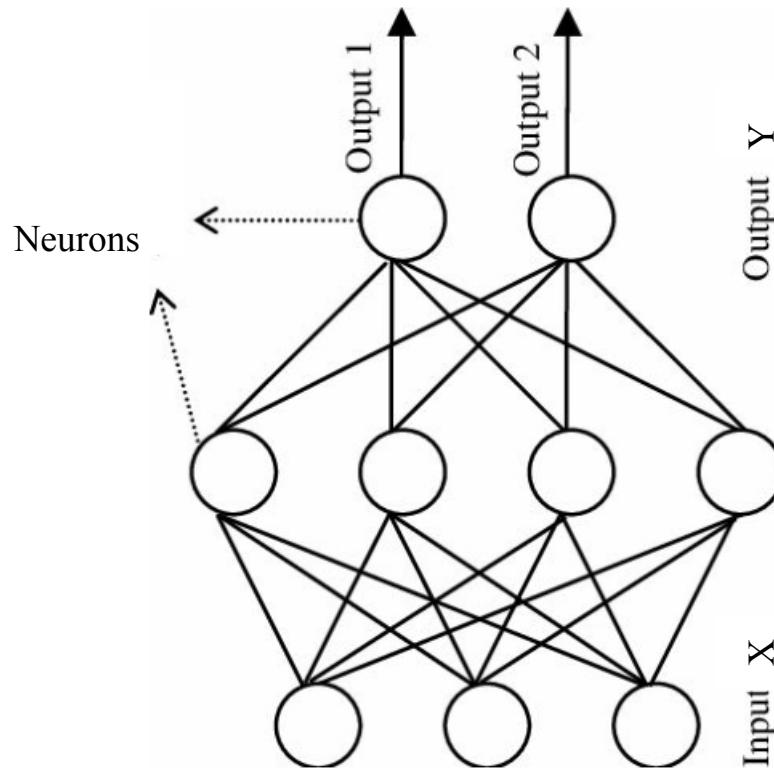
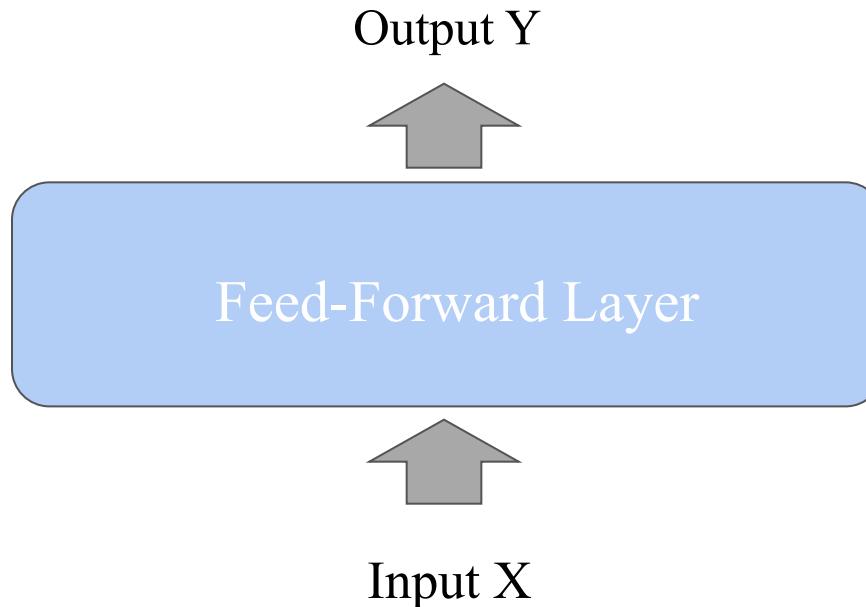


Figure from [\(R. Rezvani et. al. 2012\)](#)

In Deep Learning, it is usual to represent equations **with diagrams**

# The MultiLayer Perceptron: Diagram View



In Deep Learning, it is usual to represent equations **with diagrams**

# The MultiLayer Perceptron:

We have defined a 2-layers MLP model

We can define in the same way a **3-layers , 4-layers , L-layers** MLP

$$dnn_{(W_i \ b_i, i \in [|1, L|])}(X) = W_L \varphi_{L-1}(\dots \varphi_2 \circ W_2 \varphi_1(W_1 X + b_1) + b_2) \dots) + b_L$$

$W_l$  and  $b_l$  are trainable parameters.  $W_l \in \mathbb{R}^{\delta_{l-1} \times \delta_l}$ ,  $b_l \in \mathbb{R}^{\delta_l}$ , with  $\delta_l \in \mathbb{N}^*$ ,  $\forall l \in [|1, L|]$

$\varphi_l$  fixed non-linear functions,  $\varphi_l : \mathbb{R}^{\delta_{l-1}} \rightarrow \mathbb{R}^{\delta_l}$ ,  $\forall l \in [|1, L-1|]$

# The MultiLayer Perceptron

The same equation with a loop...

$$h_{i+1} = \varphi_i(W_i h_i + b_i), \forall i \in [|1, L - 1|]$$

with  $h_1 = X$  and  $\hat{Y} = dnn(X) = h_L$

$W_l$  and  $b_l$  are trainable parameters.  $W_l \in \mathbb{R}^{\delta_{l-1} \times \delta_l}$ ,  $b_l \in \mathbb{R}^{\delta_l}$ , with  $\delta_l \in \mathbb{N}^*$ ,  $\forall l \in [|1, L|]$

$\varphi_l$  fixed non-linear functions,  $\varphi_l : \mathbb{R}^{\delta_{l-1}} \rightarrow \mathbb{R}^{\delta_l}$ ,  $\forall l \in [|1, L - 1|]$

# The MultiLayer Perceptron

The same equation with a loop...

$$h_{i+1} = \varphi_i(W_i h_i + b_i), \forall i \in [|1, L - 1|]$$

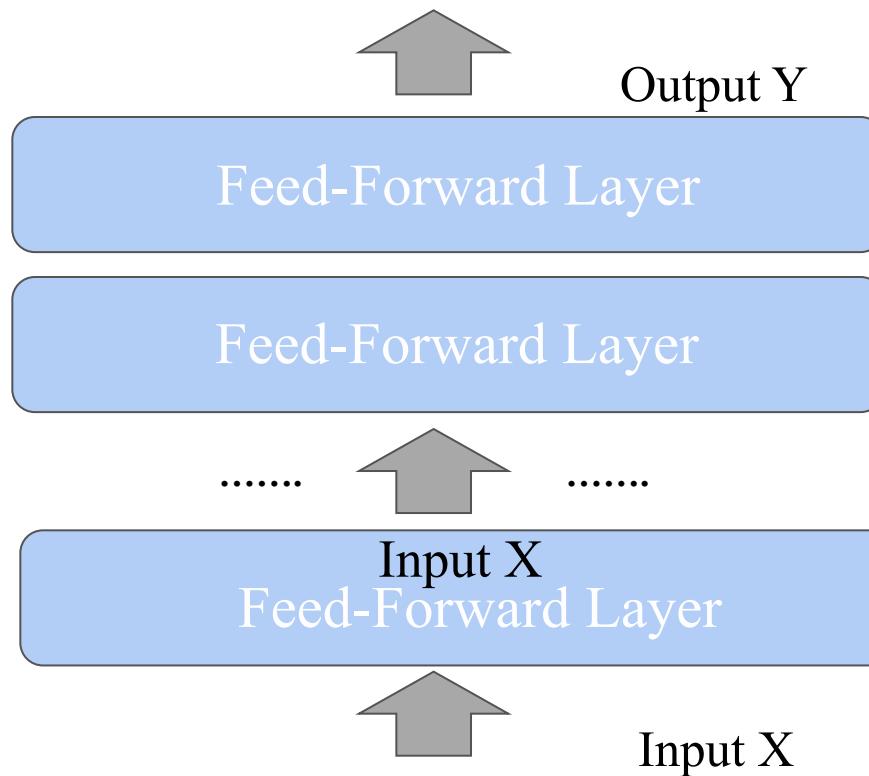
with  $h_1 = X$  and  $\hat{Y} = dnn(X) = h_L$

$W_l$  and  $b_l$  are trainable parameters.  $W_l \in \mathbb{R}^{\delta_{l-1} \times \delta_l}$ ,  $b_l \in \mathbb{R}^{\delta_l}$ , with  $\delta_l \in \mathbb{N}^*$ ,  $\forall l \in [|1, L|]$

$\varphi_l$  fixed non-linear functions,  $\varphi_l : \mathbb{R}^{\delta_{l-1}} \rightarrow \mathbb{R}^{\delta_l}$ ,  $\forall l \in [|1, L - 1|]$

$h_i$  are called hidden states ( $h_i \in \mathbb{R}^{\delta_i}$ ).

# The MultiLayer Perceptron: Diagram View



# Output Activation Function for Classification

When we do a classification task the goal is to learn a distribution of probability on the output label space

To do so, **we usually use the softmax function** as the last activation function

$$\text{softmax}(s) = \left( \frac{e^{s_i}}{\sum_k e^{s_k}} \right)_{i \in [|1, K|]}, \text{ for } s \in \mathbb{R}^K$$

# Loss Functions

Based on the task we aim at modeling, we can use:

## For Regression: Mean-Square Error

$$l(y, \hat{y}) = \|y - \hat{y}\|_2^2 = \sum_i (y_i - \hat{y}_i)^2 \text{ assuming } y_i, \hat{y}_i \in \mathbb{R}$$

## For Classification: Cross-Entropy Loss

$$l(y, \hat{y}) = CE(y, \hat{y}) = \sum_i y_i \log(\hat{y}_i) \text{ assuming } y_i, \hat{y}_i \in [0, 1]$$

Most NLP tasks will be based on the **Cross-Entropy loss**

# The MultiLayer Perceptron: Hyperparameters

- Number of **hidden layers**
- **Hidden layers dimensions**
- Initialization of the trainable parameters/weights

# The MultiLayer Perceptron: Hyperparameters

- Number of hidden layers
- Hidden layers dimensions
- **Initialization** of the **trainable parameters/weights**

# The MultiLayer Perceptron: Hyperparameters

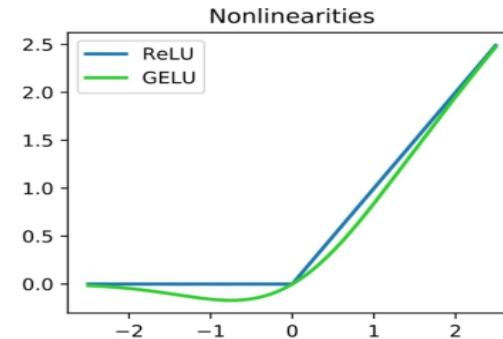
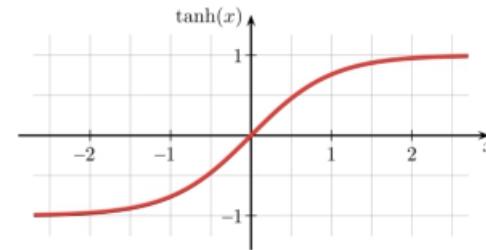
- Number of hidden layers
- Hidden layers dimensions
- Initialization
- **Activation Functions**

# The MultiLayer Perceptron: Hyperparameters

- Number of hidden layers
- Hidden layers dimensions
- Initialization
- **Activation Functions**
  - They should be **non-linear**
  - **Differentiable**
  - **Standard ones are:**  
*Relu, tanh, sigmoid*

# The MultiLayer Perceptron: Hyperparameters

- Number of hidden layers
- Hidden layers dimensions
- Initialization
- **Activation Functions**
  - They should be **non-linear**
  - **Differentiable**
  - **Standard ones are:**  
*Relu, tanh, sigmoid*



# The MultiLayer Perceptron: Hyperparameters

- Number of hidden layers
- Hidden layers dimensions
- Initialization
- Activation Functions

## How to define them?

- Look for **best practices** to choose which are the best
- In most DL libraries, the “**good**” hyperparameters are usually the default
- If no best practices/default: **you have to find the best ones empirically**

# Training Deep Learning Models

- Nearly all Deep Learning models are trained with (some version of)  
**Stochastic Gradient Descent (SGD)**

## Stochastic Gradient Descent

- The goal is find the set of **parameters/weights** that **minimizes the loss function**
- To do so, SGD estimates the true gradient of a function with **one sample at time**
- **Repeat** this process multiple times

**NB:** in deep learning, we usually train all the parameters together  
**“end-to-end”**

# Stochastic Gradient Descent

---

**Algorithm 2** Stochastic Gradient Descend

---

Given observations  $((x_i), (y_i))$  of two variables  $(X, Y)$

Given a loss function  $l$ . An architecture  $dnn_\theta$

**The goal is to find the best  $\theta$  s.t.  $E(l(Y, dnn_\theta(X)))$  is small.** Given a learning rate  $\alpha$

**for**  $step < max$  **do**

    Sample  $(x, y)$

    # Forward pass:

$\hat{y} = dnn_\theta(x)$  and  $l(y, \hat{y})$

    # Backward pass:

$\nabla_\theta l(y, \hat{y})$  # compute loss

$\theta := \theta - \alpha \nabla_\theta l(y, \hat{y})$  # parameter update

**end**

---

# Stochastic Gradient Descent

---

**Algorithm 2** Stochastic Gradient Descend

---

Given observations  $((x_i), (y_i))$  of two variables  $(X, Y)$

Given a loss function  $l$ . An architecture  $dnn_\theta$

**The goal is to find the best  $\theta$  s.t.  $E(l(Y, dnn_\theta(X)))$  is small.** Given a learning rate  $\alpha$

**for**  $step < max$  **do**

    Sample  $(x, y)$

    # Forward pass:

$\hat{y} = dnn_\theta(x)$  and  $l(y, \hat{y})$

    # Backward pass:

$\nabla_\theta l(y, \hat{y})$  # compute loss

$\theta := \theta - \alpha \nabla_\theta l(y, \hat{y})$  # parameter update

**end**

---

# Stochastic Gradient Descent

---

**Algorithm 2** Stochastic Gradient Descend

---

Given observations  $((x_i), (y_i))$  of two variables  $(X, Y)$

Given a loss function  $l$ . An architecture  $dnn_\theta$

**The goal is to find the best  $\theta$  s.t.  $E(l(Y, dnn_\theta(X)))$  is small.** Given a learning rate  $\alpha$

**for**  $step < max$  **do**

Sample  $(x, y)$

# Forward pass:

$\hat{y} = dnn_\theta(x)$  and  $l(y, \hat{y})$

# Backward pass:

$\nabla_\theta l(y, \hat{y})$  # compute loss

$\theta := \theta - \alpha \nabla_\theta l(y, \hat{y})$  # parameter update

**end**

---

# Stochastic Gradient Descent

---

**Algorithm 2** Stochastic Gradient Descend

---

Given observations  $((x_i), (y_i))$  of two variables  $(X, Y)$

Given a loss function  $l$ . An architecture  $dnn_\theta$

**The goal is to find the best  $\theta$  s.t.  $E(l(Y, dnn_\theta(X)))$  is small.** Given a learning rate  $\alpha$

**for**  $step < max$  **do**

    Sample  $(x, y)$

        # Forward pass:

$\hat{y} = dnn_\theta(x)$  and  $l(y, \hat{y})$

        # Backward pass:

$\nabla_\theta l(y, \hat{y})$  # compute loss

$\theta := \theta - \alpha \nabla_\theta l(y, \hat{y})$  # parameter update

**end**

---

# Stochastic Gradient Descent

---

**Algorithm 2** Stochastic Gradient Descend

---

Given observations  $((x_i), (y_i))$  of two variables  $(X, Y)$

Given a loss function  $l$ . An architecture  $dnn_\theta$

**The goal is to find the best  $\theta$  s.t.  $E(l(Y, dnn_\theta(X)))$  is small.** Given a learning rate  $\alpha$

**for**  $step < max$  **do**

    Sample  $(x, y)$

    # Forward pass:

$\hat{y} = dnn_\theta(x)$  and  $l(y, \hat{y})$

    # Backward pass:

$\nabla_\theta l(y, \hat{y})$  # compute gradients

$\theta := \theta - \alpha \nabla_\theta l(y, \hat{y})$  # parameter update

**end**

---

# Stochastic Gradient Descent

## Optimization Hyperparameters

### Learning Rate

- Can be refined with **variable learning rate**

*E.g. increasing during the first steps ( **warmup** ) then decreasing*

### Number of steps

- Usually defined with the validation loss

*When it stops decreasing we can stop training (= **early stopping** )*

# Stochastic Gradient Descent

Optimizing large Deep Learning Models **is challenging**

- **Unstable training**
- **Overfitting**
- **Take a lot of steps/epochs**

To make training better, many refinement of the SGD have been proposed

- In practice, we (nowadays) **use the ADAM optimizer**  
(cf. Kingma et. al 2015)

# Stochastic Gradient Descent for MLP

Let  $(X, Y) \in \mathbb{R}^d \times \mathbb{R}$ , the MSE loss  $l(y, \hat{y}) = (y - \hat{y})^2$ .

We define a 1-hidden-layer MLP with a RELU activation function of dimension  $\delta$ .

$$\hat{y} = dnn_{W_1, W_2}(x) = W_2 \max(W_1 x, 0) \text{ and } W_1 \in \mathbb{R}^{d \times \delta} \text{ and } W_2 \in \mathbb{R}^{1 \times \delta}$$

# Stochastic Gradient Descent for MLP

Let  $(X, Y) \in \mathbb{R}^d \times \mathbb{R}$ , the MSE loss  $l(y, \hat{y}) = (y - \hat{y})^2$ .

We define a 1-hidden-layer MLP with a RELU activation function of dimension  $\delta$ .

$$\hat{y} = dnn_{W_1, W_2}(x) = W_2 \max(W_1 x, 0) \text{ and } W_1 \in \mathbb{R}^{d \times \delta} \text{ and } W_2 \in \mathbb{R}^{1 \times \delta}$$

→ Goal: Apply SGD to *dnn*

# Stochastic Gradient Descent for MLP

Let  $(X, Y) \in \mathbb{R}^d \times \mathbb{R}$ , the MSE loss  $l(y, \hat{y}) = (y - \hat{y})^2$ .

We define a 1-hidden-layer MLP with a RELU activation function of dimension  $\delta$ .

$$\hat{y} = dnn_{W_1, W_2}(x) = W_2 \max(W_1 x, 0) \text{ and } W_1 \in \mathbb{R}^{d \times \delta} \text{ and } W_2 \in \mathbb{R}^{1 \times \delta}$$

1. Forward pass: Compute  $\hat{y}$

# Stochastic Gradient Descent for MLP

Let  $(X, Y) \in \mathbb{R}^d \times \mathbb{R}$ , the MSE loss  $l(y, \hat{y}) = (y - \hat{y})^2$ .

We define a 1-hidden-layer MLP with a RELU activation function of dimension  $\delta$ .

$$\hat{y} = dnn_{W_1, W_2}(x) = W_2 \max(W_1 x, 0) \text{ and } W_1 \in \mathbb{R}^{d \times \delta} \text{ and } W_2 \in \mathbb{R}^{1 \times \delta}$$

1. Forward pass
2. Compute Gradients

$$\nabla_{W_1} l(y, \hat{y}) \quad \nabla_{W_2} l(y, \hat{y})$$

# Stochastic Gradient Descent for MLP

Let  $(X, Y) \in \mathbb{R}^d \times \mathbb{R}$ , the MSE loss  $l(y, \hat{y}) = (y - \hat{y})^2$ .

We define a 1-hidden-layer MLP with a RELU activation function of dimension  $\delta$ .

$$\hat{y} = dnn_{W_1, W_2}(x) = W_2 \max(W_1 x, 0) \text{ and } W_1 \in \mathbb{R}^{d \times \delta} \text{ and } W_2 \in \mathbb{R}^{1 \times \delta}$$

1. Forward pass
2. Compute Gradients
3. Backward pass (parameter update)

# Stochastic Gradient Descent for MLP

Let  $(X, Y) \in \mathbb{R}^d \times \mathbb{R}$ , the MSE loss  $l(y, \hat{y}) = (y - \hat{y})^2$ .

$$\hat{y} = dnn_{W_1, W_2}(x) = W_2 \max(W_1 x, 0) \text{ and } W_1 \in \mathbb{R}^{d \times \delta} \text{ and } W_2 \in \mathbb{R}^{1 \times \delta}$$

Idea: we use **the chain rule** to decompose **the gradient starting from the top layers**

# Stochastic Gradient Descent for MLP

Let  $(X, Y) \in \mathbb{R}^d \times \mathbb{R}$ , the MSE loss  $l(y, \hat{y}) = (y - \hat{y})^2$ .

$$\hat{y} = dnn_{W_1, W_2}(x) = W_2 \underbrace{\max(W_1 x, 0)}_{h_1} \text{ and } W_1 \in \mathbb{R}^{d \times \delta} \text{ and } W_2 \in \mathbb{R}^{1 \times \delta}$$

## Compute Gradient

$$\nabla_{W_2} l(y, \hat{y}) = \frac{\partial l}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial W_2}$$

# Stochastic Gradient Descent for MLP

Let  $(X, Y) \in \mathbb{R}^d \times \mathbb{R}$ , the MSE loss  $l(y, \hat{y}) = (y - \hat{y})^2$ .

$$\hat{y} = dnn_{W_1, W_2}(x) = W_2 \underbrace{\max(W_1 x, 0)}_{h_1} \text{ and } W_1 \in \mathbb{R}^{d \times \delta} \text{ and } W_2 \in \mathbb{R}^{1 \times \delta}$$

## Compute Gradient

$$\nabla_{W_2} l(y, \hat{y}) = \frac{\partial l}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial W_2} = 2(y - \hat{y}) h_1$$

# Stochastic Gradient Descent for MLP

Let  $(X, Y) \in \mathbb{R}^d \times \mathbb{R}$ , the MSE loss  $l(y, \hat{y}) = (y - \hat{y})^2$ .

$$\hat{y} = dnn_{W_1, W_2}(x) = W_2 \underbrace{\max(W_1 x, 0)}_{h_1} \text{ and } W_1 \in \mathbb{R}^{d \times \delta} \text{ and } W_2 \in \mathbb{R}^{1 \times \delta}$$

Compute Gradient:

$$\nabla_{W_2} l(y, \hat{y}) = \frac{\partial l}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial W_2} = 2(y - \hat{y}) h_1$$

$$\nabla_{W_1} l(y, \hat{y}) = \frac{\partial l}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial h_1} \frac{\partial h_1}{\partial W_2}$$

# Stochastic Gradient Descent for MLP

Let  $(X, Y) \in \mathbb{R}^d \times \mathbb{R}$ , the MSE loss  $l(y, \hat{y}) = (y - \hat{y})^2$ .

$$\hat{y} = dnn_{W_1, W_2}(x) = W_2 \underbrace{\max(W_1 x, 0)}_{h_1} \text{ and } W_1 \in \mathbb{R}^{d \times \delta} \text{ and } W_2 \in \mathbb{R}^{1 \times \delta}$$

**Compute Gradient:**

$$\nabla_{W_2} l(y, \hat{y}) = \frac{\partial l}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial W_2} = 2(y - \hat{y}) h_1$$

$$\nabla_{W_1} l(y, \hat{y}) = \frac{\partial l}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial h_1} \frac{\partial h_1}{\partial W_2} = 2(y - \hat{y}) W_2 1_{W_1 X > 0}$$

# Backpropagation and Deep Learning in practice

In practice, we use Deep Learning Libraries

- Define **the Architecture with *tensor* operators**
- Backpropagation is done **seamlessly using automatic differentiation**

# Deep Learning & Backpropagation in practice

In practice, we use Deep Learning Libraries (e.g. pytorch, tensorflow, jax)

- Define the Architecture with *tensor* operators
- Backpropagation is done **seamlessly using automatic differentiation**
- Standard layers **are pre-implemented** (Feed-Forward Layers, LSTM, Attention, Self-Attention...)

[See code example with pytorch](#)

# Recurrent Neural Network

# Vanilla Recurrent Neural Network

We would like to model sequences (e.g. words)  $(X_1, \dots, X_T)$  in  $\mathbb{R}^{d,T}$

We can introduce **a recurrence relation** into our MLP to model it:

$$h_{i+1,t+1} = \varphi_i(W_i h_{i,t} + U_i h_{i+1,t} + b_i), \forall i \in [|1, L-1|]$$

with  $h_{1,t} = X_t$  and  $\hat{Y}_t = dnn(X_t) = h_{L,t} \forall t \in [|1, T-1|]$

# Recurrent Neural Network

## Illustration of a 1-layer Recurrent Neural Network

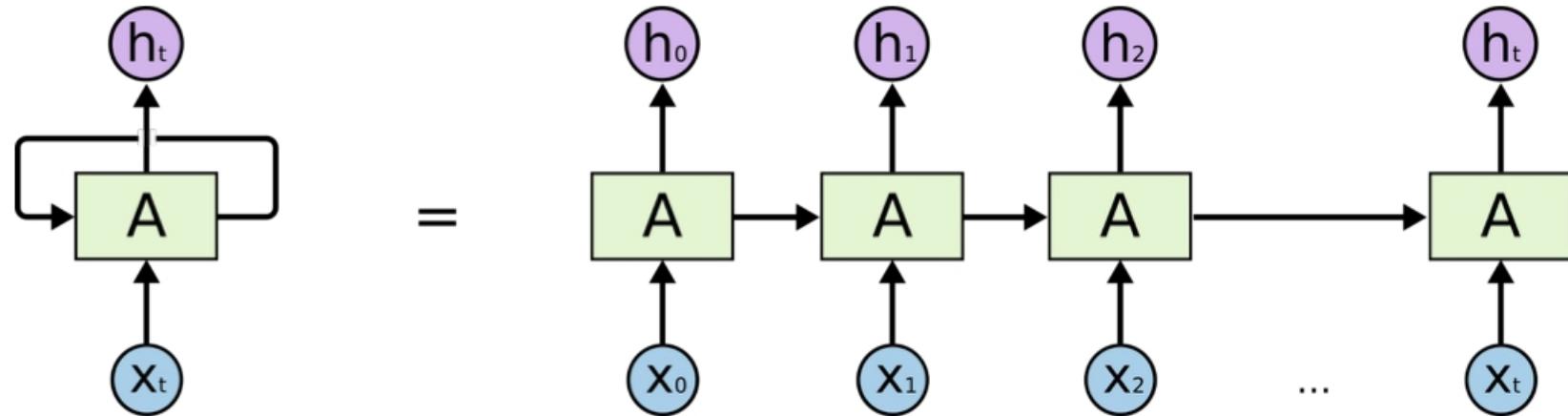


Figure from [colah](#)

# Recurrent Neural Network

## Illustration of a 1-layer Recurrent Neural Network

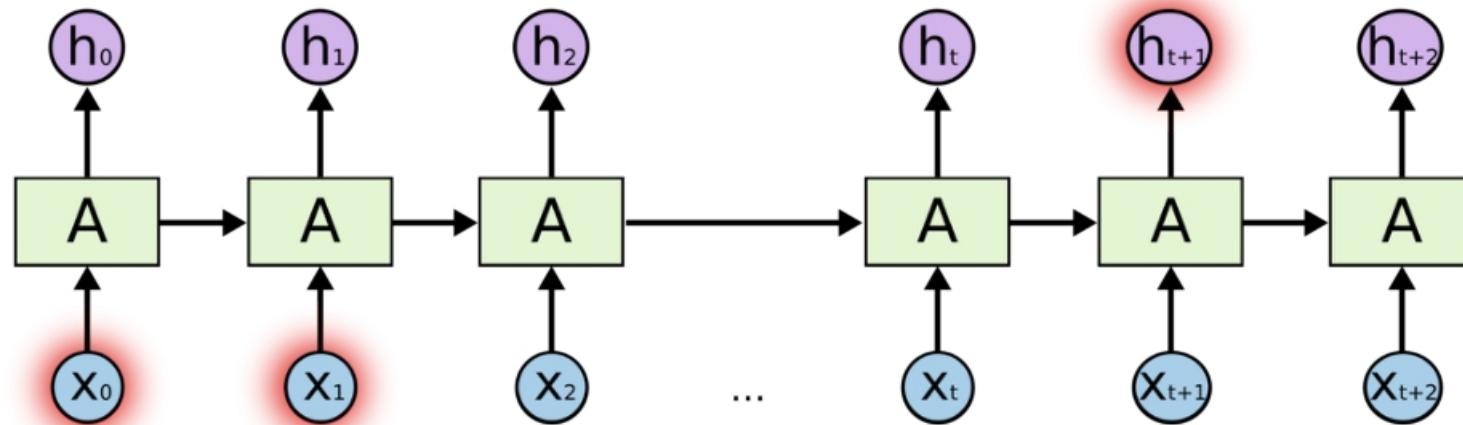


Figure from [colah](#)

# Training Recurrent Neural Network

Recurrent Neural Network are trained with an extension of the Backpropagation algorithm

→ Backpropagation Through Time (BPTT)

BPTT follows exactly the same ideas as backpropagation

- SGD
- Chain Rule starting from the last layer and the last hidden state
- **With extra derivative dependencies between state  $t$  and  $t+1$**

# Limits of Recurrent Neural Networks

Vanilla Recurrent Neural Network have trouble to capture long-term dependencies

Idea:

- Encode **explicitly in a vector a “memory” in the recurrent architecture**
- Control what is memorized and forgotten
- Train all those parameters **end-to-end**

# Attention Mechanism

## Motivation for Attention Mechanisms

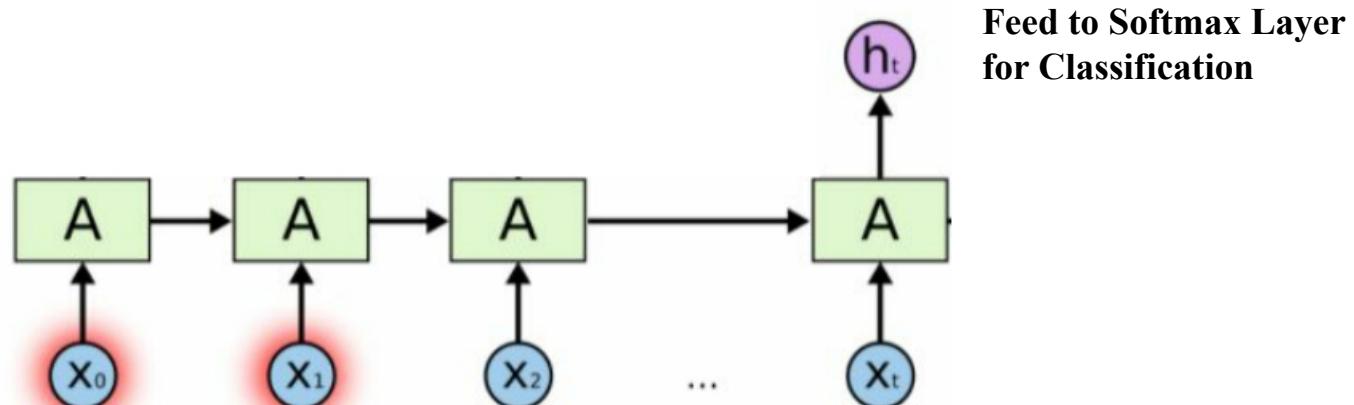
- The Deep Learning Architecture that we have seen so far are **hard to interpret (black-box)**
- Recurrent Network provide a fixed vector encoding of a sequence at each step

→ **Attention Mechanisms**

# Attention Mechanism for Sequence Classification

We want to classify  $(X_0, X_t)$  sequences (e.g. sentiment analysis)

**Solution 1:** Use a LSTM model → Problem (not interpretable)

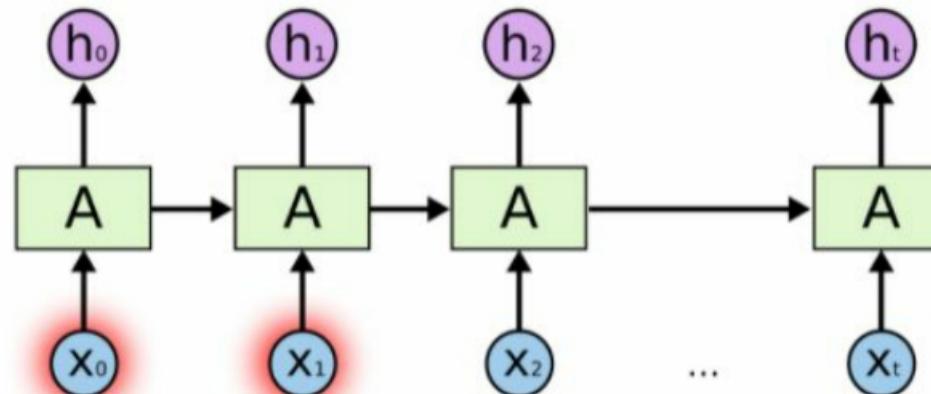


# Attention Mechanism for Sequence Classification

We want to classify  $(X_0, X_t)$  sequences (e.g. sentiment)

**Solution 2:** Integrate an Attention Mechanism to interpret what input impacts the prediction

→ Learn a ponderation/weighting of the hidden states  $h_t$



# Attention Mechanism for Sequence Classification

We want to classify ( $X_0, X_t$ ) sequences (e.g. sentiment)

How to learn this weighting?

1. Define a specific type of layer to learn the ponderation
2. Train this layer end-to-end with all the other parameters of the model

# Attention Mechanism for Sequence Classification

We want to classify ( $X_0, X_t$ ) sequences (e.g. sentiment)

How to learn this weighting?

Given  $(h_1, \dots, h_T)$  hidden representations of  $(x_1, \dots, x_T)$  (e.g. output of a LSTM Layer).

$$q_i = \tanh(W_a h_i + b_a), \text{ with } W_a \in \mathbb{R}^{\delta \times \delta_a}$$

$$s_t = \frac{e^{q_t q_T}}{\sum_j e^{q_j q_T}}, \text{ i.e. } \sum_{t \in [1, T]} s_t = 1$$

$$\tilde{h}_T = \sum_{t \in [1, T]} s_t \cdot h_t$$

# Attention Mechanism for Sequence Classification

We want to classify ( $X_0, X_t$ ) sequences (e.g. sentiment)

How to learn this weighting?

Given  $(h_1, \dots, h_T)$  hidden representations of  $(x_1, \dots, x_T)$  (e.g. output of a LSTM Layer).

$$q_i = \tanh(W_a h_i + b_a), \text{ with } W_a \in \mathbb{R}^{\delta \times \delta_a}$$

$$s_t = \frac{e^{q_t q_T}}{\sum_j e^{q_j q_T}}, \text{ i.e. } \sum_{t \in [1, T]} s_t = 1$$

$$\tilde{h}_T = \sum_{t \in [1, T]} s_t \cdot h_t$$

# Attention Mechanism for Sequence Classification

We want to classify  $(X_0, X_t)$  sequences (e.g. sentiment)

How to learn this weighting?

Given  $(h_1, \dots, h_T)$  hidden representations of  $(x_1, \dots, x_T)$  (e.g. output of a LSTM Layer).

$$q_i = \tanh(W_a h_i + b_a), \text{ with } W_a \in \mathbb{R}^{\delta \times \delta_a}$$

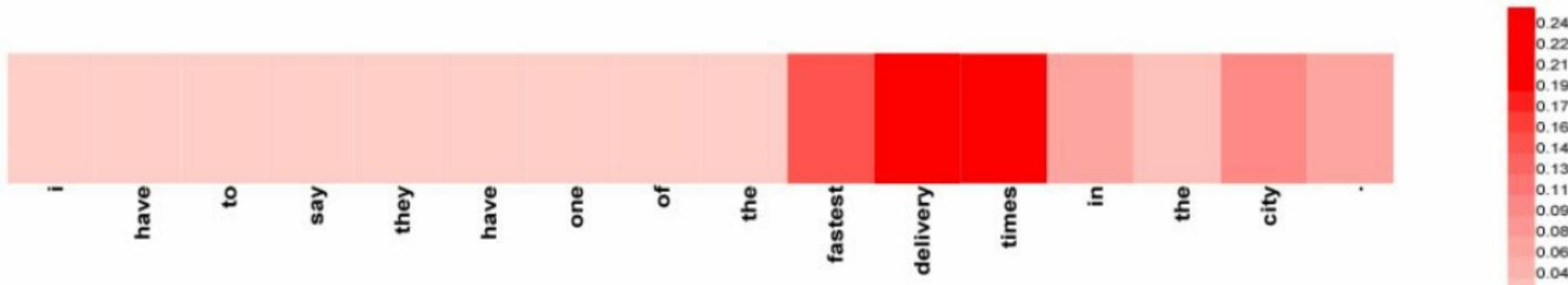
$$s_t = \frac{e^{q_t q_T}}{\sum_j e^{q_j q_T}}, \text{ i.e. } \sum_{t \in [1, T]} s_t = 1$$

$$\tilde{h}_T = \sum_{t \in [1, T]} s_t \cdot h_t$$

# Attention Mechanism for Sequence Classification

We want to classify  $(X_0, X_t)$  sequences (e.g. sentiment classification)

After we trained the model, **Attention scores** can be used to interpret the model behavior and what input vector impacted the decision



(Wang et. al 2016)

# Attention Mechanism for Sequence Classification

Many variant of Attention Mechanisms (in combination with LSTM layers) have been designed

## Design Choices

- How to define the *query vectors* ?
- How to define the *scoring function* ?

Many variants exists but the principles are the same.

# The Transformer Architecture

# *Attention might be all we need*

**Do we really need recurrent layers?**

RNN models (such as vanilla RNN, LSTM...) were designed to model sequential data

Still, for most tasks, we **need both left and right context** (e.g. **sequence classification, sequence labelling..**)

Why not modelling sequences in a bi-directional way directly  
→ **Using Self-Attention Mechanism**

# Self-Attention Layers

Given a sequence of input vectors  $(x_1, \dots, x_T) \in \mathbb{R}^\delta$  (noted  $(h_{0,1}, \dots, h_{0,T})$ ).

## Objective:

- Build a representation of the input vectors based on the **surrounding vectors** (both right-and left-context)

## Idea:

- **No need of recurrent cells**
- **Self-Attention**

# Self-Attention Layers: Intuition

Given a sequence of input vectors  $X = (x_1, \dots, x_T) \in \mathbb{R}^\delta$  (noted  $H = (h_{0,1}, \dots, h_{0,T})$ ).

We build 3 new vectorial representation of our sequence  $H = (h_1, \dots, h_T)$ .

The *query*  $Q = (q_1, \dots, q_T)$ , the *key*  $K = (k_1, \dots, k_T)$  and the *value*  $V = (v_1, \dots, v_T)$  vectors.

- For a given vector  $h_t$  and its query vector  $q_t$  we want to build the new representation vector  $\tilde{h}_t$
- Using the best ponderation of the information encoded in  $(v_1, \dots, v_T)$
- This ponderation being computed by finding the key vectors in  $(k_1, \dots, k_T)$  that are more similar to the query vector  $q_t$  (that encodes relevant information from  $h_t$ ).

# Self-Attention Layers: Intuition

Given a sequence of input vectors  $X = (x_1, \dots, x_T) \in \mathbb{R}^\delta$  (noted  $H = (h_{0,1}, \dots, h_{0,T})$ ).

We build 3 new vectorial representation of our sequence  $H = (h_1, \dots, h_T)$ .

The *query*  $Q = (q_1, \dots, q_T)$ , the *key*  $K = (k_1, \dots, k_T)$  and the *value*  $V = (v_1, \dots, v_T)$  vectors.

- For a given vector  $h_t$  and its query vector  $q_t$  we want to build the new representation vector  $\tilde{h}_t$
- Using the best ponderation of the information encoded in  $(v_1, \dots, v_T)$
- This ponderation being computed by finding the key vectors in  $(k_1, \dots, k_T)$  that are more similar to the query vector  $q_t$  (that encodes relevant information from  $h_t$ ).

# Self-Attention Layers: Intuition

Given a sequence of input vectors  $X = (x_1, \dots, x_T) \in \mathbb{R}^\delta$  (noted  $H = (h_{0,1}, \dots, h_{0,T})$ ).

We build 3 new vectorial representation of our sequence  $H = (h_1, \dots, h_T)$ .

The *query*  $Q = (q_1, \dots, q_T)$ , the *key*  $K = (k_1, \dots, k_T)$  and the *value*  $V = (v_1, \dots, v_T)$  vectors.

- For a given vector  $h_t$  and its query vector  $q_t$  we want to build the new representation vector  $\tilde{h}_t$
- Using the best ponderation of the information encoded in  $(v_1, \dots, v_T)$
- This ponderation being computed by finding the key vectors in  $(k_1, \dots, k_T)$  that are more similar to the query vector  $q_t$  (that encodes relevant information from  $h_t$ ).

# Self-Attention Layers: Intuition

Given a sequence of input vectors  $X = (x_1, \dots, x_T) \in \mathbb{R}^\delta$  (noted  $H = (h_{0,1}, \dots, h_{0,T})$ ).

We build 3 new vectorial representation of our sequence  $H = (h_1, \dots, h_T)$ .

The *query*  $Q = (q_1, \dots, q_T)$ , the *key*  $K = (k_1, \dots, k_T)$  and the *value*  $V = (v_1, \dots, v_T)$  vectors.

- For a given vector  $h_t$  and its query vector  $q_t$  we want to build the new representation vector  $\tilde{h}_t$
- Using the best ponderation of the information encoded in  $(v_1, \dots, v_T)$
- This ponderation being computed by finding the key vectors in  $(k_1, \dots, k_T)$  that are more similar to the query vector  $q_t$  (that encodes relevant information from  $h_t$ ).

# Self-Attention Layers

Given a sequence of input vectors  $X = (x_1, \dots, x_T) \in \mathbb{R}^\delta$  (noted  $H = (h_{0,1}, \dots, h_{0,T})$ ).

We build 3 new vectorial representation of our sequence  $H = (h_1, \dots, h_T)$ .

The *query*  $Q = (q_1, \dots, q_T)$ , the *key*  $K = (k_1, \dots, k_T)$  and the *value*  $V = (v_1, \dots, v_T)$  vectors.

$$q_t = W_Q h_t, \forall t \in [|1, T|] \text{ with } W_Q \in \mathbb{R}^{\delta_q \times \delta}$$

$$k_t = W_K h_t, \forall t \in [|1, T|] \text{ with } W_K \in \mathbb{R}^{\delta_k \times \delta}$$

$$v_t = W_V h_t, \forall t \in [|1, T|] \text{ with } W_V \in \mathbb{R}^{\delta_v \times \delta}$$

# Self-Attention Layers

Given a sequence of input vectors  $X = (x_1, \dots, x_T) \in \mathbb{R}^\delta$  (noted  $H = (h_{0,1}, \dots, h_{0,T})$ ).

We build 3 new vectorial representation of our sequence  $H = (h_1, \dots, h_T)$ .

The *query*  $Q = (q_1, \dots, q_T)$ , the *key*  $K = (k_1, \dots, k_T)$  and the *value*  $V = (v_1, \dots, v_T)$  vectors.

$$q_t = W_Q h_t, \forall t \in [|1, T|] \text{ with } W_Q \in \mathbb{R}^{\delta_q \times \delta}$$

$$k_t = W_K h_t, \forall t \in [|1, T|] \text{ with } W_K \in \mathbb{R}^{\delta_k \times \delta}$$

$$v_t = W_V h_t, \forall t \in [|1, T|] \text{ with } W_V \in \mathbb{R}^{\delta_v \times \delta}$$

# Self-Attention Layers

Given a sequence of input vectors  $X = (x_1, \dots, x_T) \in \mathbb{R}^\delta$  (noted  $H = (h_{0,1}, \dots, h_{0,T})$ ).

We build 3 new vectorial representation of our sequence  $H = (h_1, \dots, h_T)$ .

The *query*  $Q = (q_1, \dots, q_T)$ , the *key*  $K = (k_1, \dots, k_T)$  and the *value*  $V = (v_1, \dots, v_T)$  vectors.

$$q_t = W_Q h_t, \forall t \in [|1, T|] \text{ with } W_Q \in \mathbb{R}^{\delta_q \times \delta}$$

$$k_t = W_K h_t, \forall t \in [|1, T|] \text{ with } W_K \in \mathbb{R}^{\delta_k \times \delta}$$

$$v_t = W_V h_t, \forall t \in [|1, T|] \text{ with } W_V \in \mathbb{R}^{\delta_v \times \delta}$$

# Self-Attention Layers

Given a sequence of input vectors  $X = (x_1, \dots, x_T) \in \mathbb{R}^\delta$  (noted  $H = (h_{0,1}, \dots, h_{0,T})$ ).

We build 3 new vectorial representation of our sequence  $H = (h_1, \dots, h_T)$ .

The *query*  $Q = (q_1, \dots, q_T)$ , the *key*  $K = (k_1, \dots, k_T)$  and the *value*  $V = (v_1, \dots, v_T)$  vectors.

$$\tilde{H} = \text{softmax}\left(\frac{Q K^T}{\sqrt{\delta_K}}\right)V$$

i.e.  $\tilde{h}_t = \text{softmax}\left(\frac{q_t K^T}{\sqrt{\delta_K}}\right).V = \sum_{t'} s_{t'} v_{t'}$  with  $s_{t'} = \frac{e^{q_t k_{t'}}}{\sum_t e^{q_t k_t}}$

# Self-Attention Layers

Given a sequence of input vectors  $X = (x_1, \dots, x_T) \in \mathbb{R}^\delta$  (noted  $H = (h_{0,1}, \dots, h_{0,T})$ ).

We build 3 new vectorial representation of our sequence  $H = (h_1, \dots, h_T)$ .

The *query*  $Q = (q_1, \dots, q_T)$ , the *key*  $K = (k_1, \dots, k_T)$  and the *value*  $V = (v_1, \dots, v_T)$  vectors.

$$\tilde{H} = \text{softmax}\left(\frac{Q K^T}{\sqrt{\delta_K}}\right)V$$

i.e.  $\tilde{h}_t = \text{softmax}\left(\frac{q_t K^T}{\sqrt{\delta_K}}\right).V = \sum_{t'} s_{t'} v_{t'} \text{ with } s_{t'} = \frac{e^{q_t' k_t}}{\sum_t e^{q_t' k_t}}$

# Self-Attention Layers

Given a sequence of input vectors  $X = (x_1, \dots, x_T) \in \mathbb{R}^\delta$  (noted  $H = (h_{0,1}, \dots, h_{0,T})$ ).

We build 3 new vectorial representation of our sequence  $H = (h_1, \dots, h_T)$ .

The *query*  $Q = (q_1, \dots, q_T)$ , the *key*  $K = (k_1, \dots, k_T)$  and the *value*  $V = (v_1, \dots, v_T)$  vectors.

$$\tilde{H} = \text{softmax}\left(\frac{Q K^T}{\sqrt{\delta_K}}\right)V$$

i.e.  $\tilde{h}_t = \text{softmax}\left(\frac{q_t K^T}{\sqrt{\delta_K}}\right).V = \sum_{t'} s_{t'} v_{t'}$  with  $s_{t'} = \frac{e^{q_t k_{t'}}}{\sum_t e^{q_t k_t}}$

# Self-Attention Layers

Attention is a building block.

Think of it as a "soft" **kv** dictionary lookup:

1. Attention weights  $a_{1:N}$  are query-key similarities:

$$\hat{a}_i = \mathbf{q} \cdot \mathbf{k}_i$$

Normalized via softmax:  $a_i = e^{\hat{a}_i} / \sum_j e^{\hat{a}_j}$

2. Output **z** is attention-weighted average of values  $v_{1:N}$ :

$$\mathbf{z} = \sum_i \hat{a}_i \mathbf{v}_i = \hat{\mathbf{a}} \cdot \mathbf{v}$$

3. Usually, **k** and **v** are derived from the same input **x**:

$$\mathbf{k} = \mathbf{W}_k \cdot \mathbf{x}$$

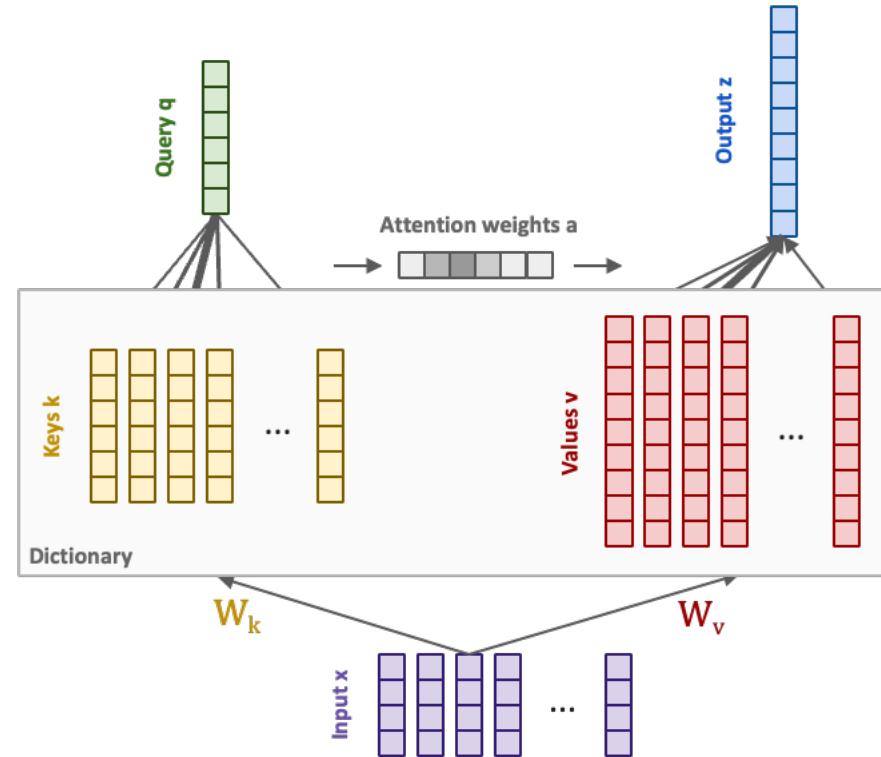
$$\mathbf{v} = \mathbf{W}_v \cdot \mathbf{x}$$

The query **q** can come from a separate input **y**:

$$\mathbf{q} = \mathbf{W}_q \cdot \mathbf{y}$$

Or from the same input **x**! Then we call it "self attention":

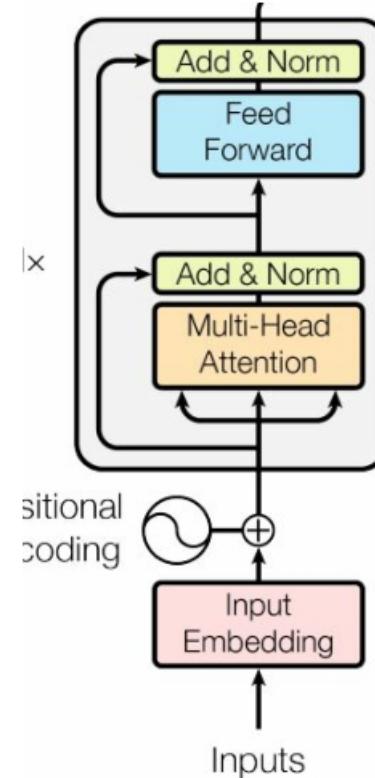
$$\mathbf{q} = \mathbf{W}_q \cdot \mathbf{x}$$



# The Transformer Architecture

## The Transformer Architecture is

- Stack of [Self-Attention + FF Layer]
- With Skip-Layer and Normalization Layers in between
- Encoding the position with positional vector



# Multi-head Self-Attention Layers

1. We usually use **many queries**  $q_{1:M}$ , not just one.

Stacking them leads to the Attention matrix  $A_{1:N,1:M}$

and subsequently to many outputs:

$$z_{1:M} = \text{Attn}(q_{1:M}, x) = [\text{Attn}(q_1, x) | \text{Attn}(q_2, x) | \dots |$$

2.  $\text{Attn}(q_M, x)$  We usually use "**multi-head**" attention. This means

the

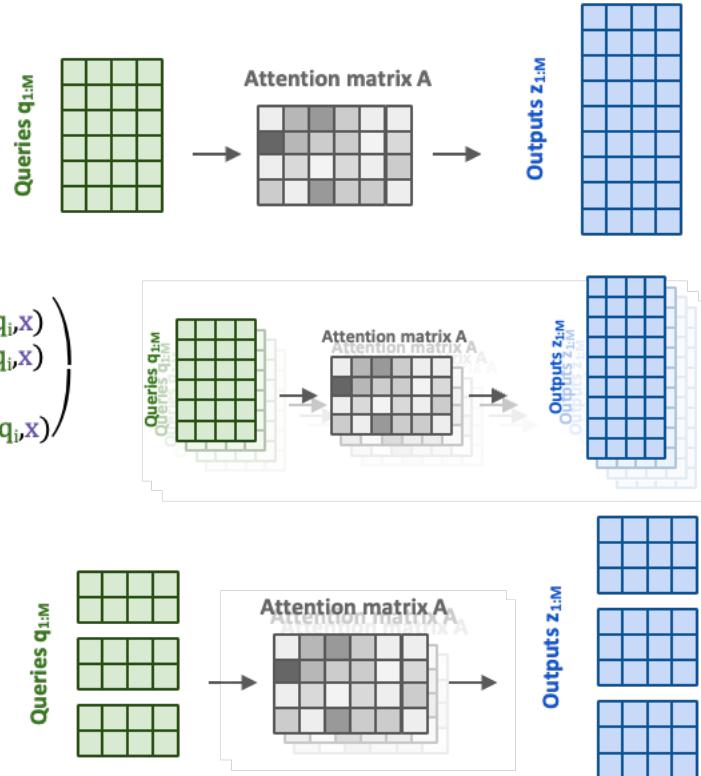
operation is repeated K times and the results are

concatenated along the feature dimension. Ws differ.

3. The most commonly seen formulation:

$$z = \text{softmax}(QK'/\sqrt{d_{\text{key}}})V$$

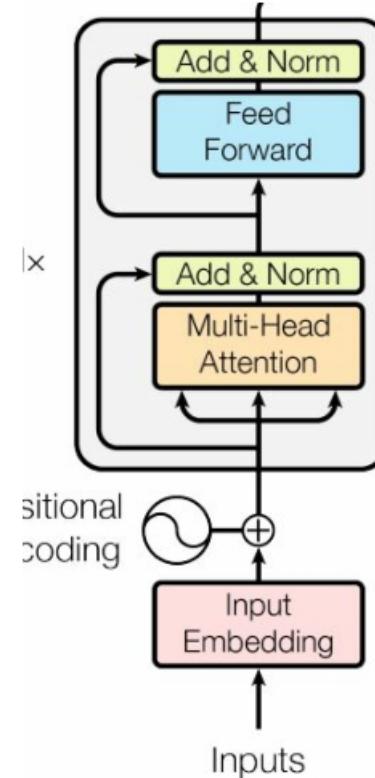
Note that the complexity is  $O(N^2)$



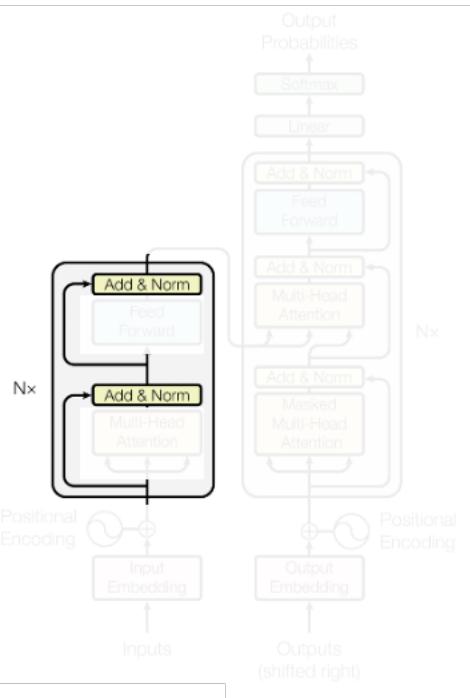
# The Transformer Architecture

## The Transformer Architecture is

- Stack of [Self-Attention + FF Layer]
- With Skip-Layer and Normalization Layers in between
- Encoding the position with positional vector



# Residual Connections



## Residual/skip connections

Each module's output has the exact same shape as its input.

Following ResNets, the module computes a "residual" instead of a new value:

$$z_i = \text{Module}(x_i) + x_i$$

This was shown to dramatically improve trainability.

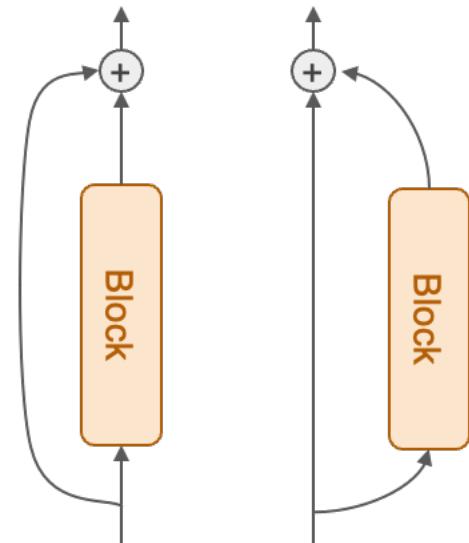
## LayerNorm

Normalization also dramatically improves trainability.

There's **post-norm** (original) and **pre-norm** (modern)

$$\text{Module}(x_i) + x_i$$

"Skip connection" == "Residual block"

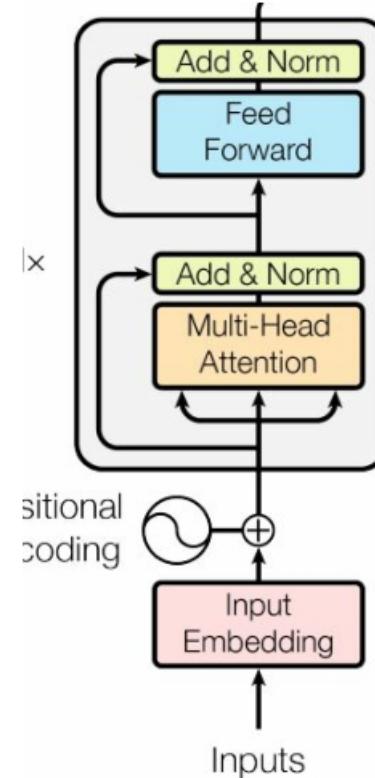


$$z_i = \text{Module}(\text{LN}(x_i)) + x_i$$

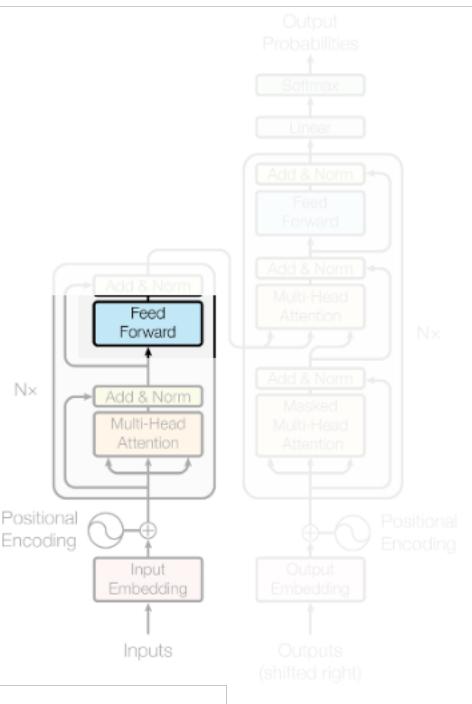
# The Transformer Architecture

## The Transformer Architecture is

- Stack of [Self-Attention + FF Layer]
- With Skip-Layer and Normalization Layers in between
- Encoding the position with positional vector



# MLP



## Point-wise MLP

A simple MLP applied to each token individually:

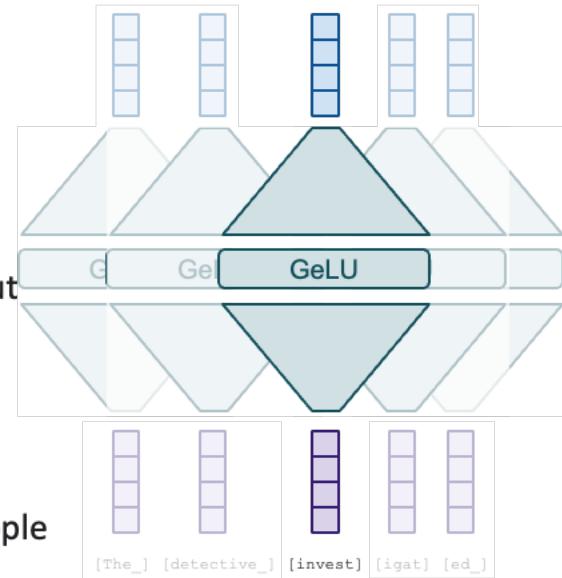
$$z_i = W_2 \text{GeLU}(W_1 x + b_1) + b_2$$

Think of it as each token pondering for itself about what it has observed previously.

There's some weak evidence this is where "world knowledge" is stored, too.

It contains the bulk of the parameters. When people make giant models and sparse/moe, this is what becomes giant.

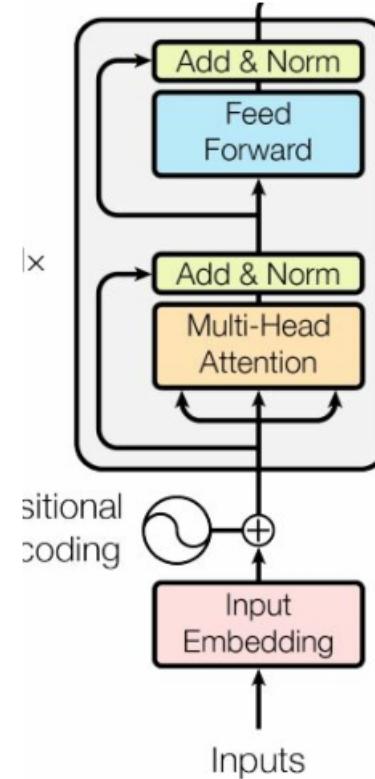
Some people like to call it  $1 \times 1$  convolution.



# The Transformer Architecture

## The Transformer Architecture is

- Stack of [Self-Attention + FF Layer]
- With Skip-Layer and Normalization Layers in between
- Encoding the position with positional vector



# Positional Embedding Vector

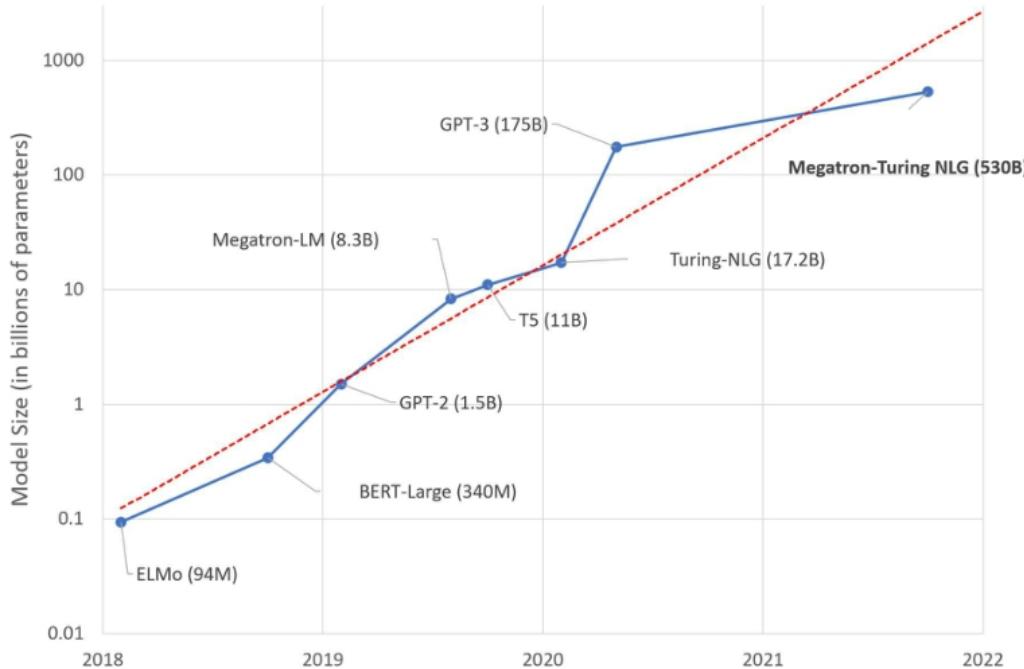
- **Limitation:** self attention does not take position into account!
- Indeed, shuffling the input gives the same results
- **Solution:** add position encodings.
- Replace the matrix  $\mathbf{W}$  by  $\mathbf{W} + \mathbf{E}$ , where  $\mathbf{E} \in \mathbb{R}^{d \times T}$
- $\mathbf{E}$  can be learned, or defined using sin and cos:

$$e_{2i,j} = \sin\left(\frac{j}{10000^{2i/d}}\right)$$
$$e_{2i+1,j} = \cos\left(\frac{j}{10000^{2i/d}}\right)$$

# Scaling Laws Intuition

- The larger the dimension of the weight matrices
- The larger the number of parameters in the model
- The more “expressive” is the model
- The better it will generalize

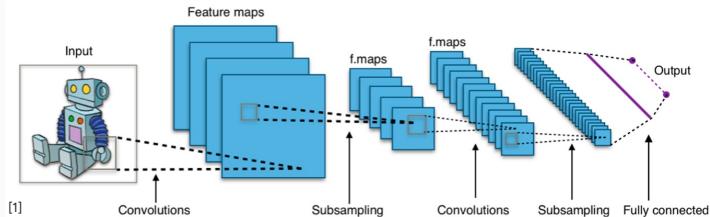
# Typical Architecture Sizes



**The classic landscape:  
One architecture  
per "community"**

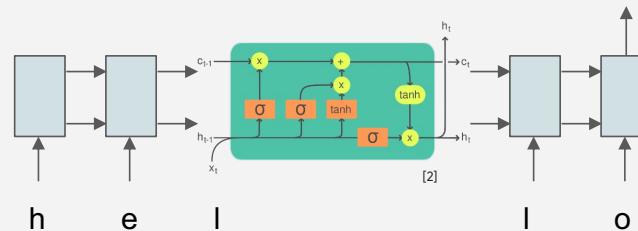
# Computer Vision

## Convolutional NNs (+ResNets)



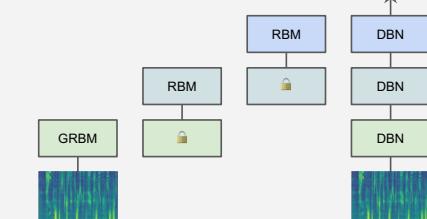
# Natural Lang. Proc.

## Recurrent NNs (+LSTMs)



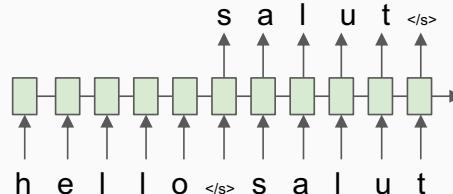
# Speech

## Deep Belief Nets (+non-DL)



# Translation

## Seq2Seq



[1] CNN image CC-BY-SA by Aphex34 for Wikipedia [https://commons.wikimedia.org/wiki/File:Typical\\_cnn.png](https://commons.wikimedia.org/wiki/File:Typical_cnn.png)  
[2] RNN image CC-BY-SA by GChe for Wikipedia [https://commons.wikimedia.org/wiki/File:The\\_LSTM\\_Cell.svg](https://commons.wikimedia.org/wiki/File:The_LSTM_Cell.svg)

# RL

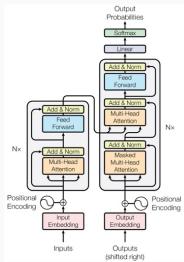
## BC/GAIL

### Algorithm 1 Generative adversarial imitation learning

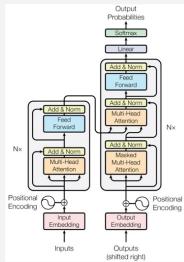
- 1: **Input:** Expert trajectories  $\tau_E \sim \pi_E$ , initial policy and discriminator parameters  $\theta_0, w_0$
- 2: **for**  $i = 0, 1, 2, \dots$  **do**
- 3:   Sample trajectories  $\tau_i \sim \pi_{\theta_i}$
- 4:   Update the discriminator parameters from  $w_i$  to  $w_{i+1}$  with the gradient
$$\hat{\mathbb{E}}_{\tau_i}[\nabla_w \log(D_w(s, a))] + \hat{\mathbb{E}}_{\tau_E}[\nabla_w \log(1 - D_w(s, a))] \quad (17)$$
- 5:   Take a policy step from  $\theta_i$  to  $\theta_{i+1}$ , using the TRPO rule with cost function  $\log(D_{w_{i+1}}(s, a))$ . Specifically, take a KL-constrained natural gradient step with
$$\hat{\mathbb{E}}_{\tau_i}[\nabla_\theta \log \pi_\theta(a|s)Q(s, a)] - \lambda \nabla_\theta H(\pi_\theta), \quad (18)$$
where  $Q(\bar{s}, \bar{a}) = \hat{\mathbb{E}}_{\tau_i}[\log(D_{w_{i+1}}(s, a)) | s_0 = \bar{s}, a_0 = \bar{a}]$
- 6: **end for**

# The Transformer's takeover: One community at a time

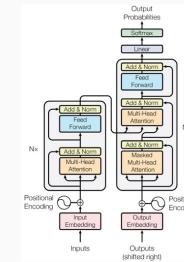
## Computer Vision



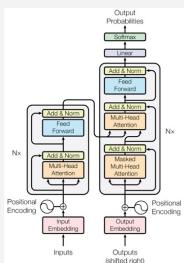
## Natural Lang. Proc.



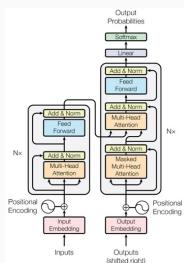
## Reinf. Learning



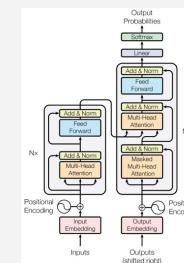
## Speech



## Translation



## Graphs/Science



Transformer image source: "Attention Is All You Need" paper

# Lecture Summary

**Deep Learning is a powerful and general modelling approach**

- **Designing Architectures** , i.e. composition of linear transformation and non-linear transformation (possibly including recurrences)
- All those transformations **should be differentiable**
- All the parameters of the model **are trained with backpropagation**
- **Toward a specific task** s.t. regression or classification
- All the hyperparameters are chosen based on **best-practices** or empirical research

# Bibliography and Acknowledgment

All these class have been taken from <https://nlp-ensae.github.io/materials/> and is taken from Benjamin Muller