



Natural Language Processing

AIS

Romain Benassi

Course 4 - Part 1

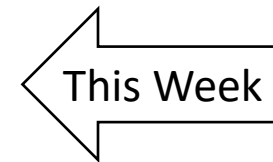
Spring 2022

Course Schedule

- **Course 1:** NLP introduction
- **Course 2:** Word embedding
- **Course 3:** Long Short-Term Memory (LSTM) architecture
- **Course 4:** “Attention” mechanism and Transformer architectures
- **Course 5:** Chatbot introduction

Course Schedule

- **Course 1:** NLP introduction
- **Course 2:** Word embedding
- **Course 3:** Long Short Term Memory (LSTM) architecture
- **Course 4: “Attention” mechanism and Transformer architectures**
 - Encoder-decoder
 - Attention Mechanism
 - Transformer architectures
 - BERT model
 - Current state-of-the-art
 - Fine tuning use case on GPT-2
- **Course 5:** Chatbot introduction

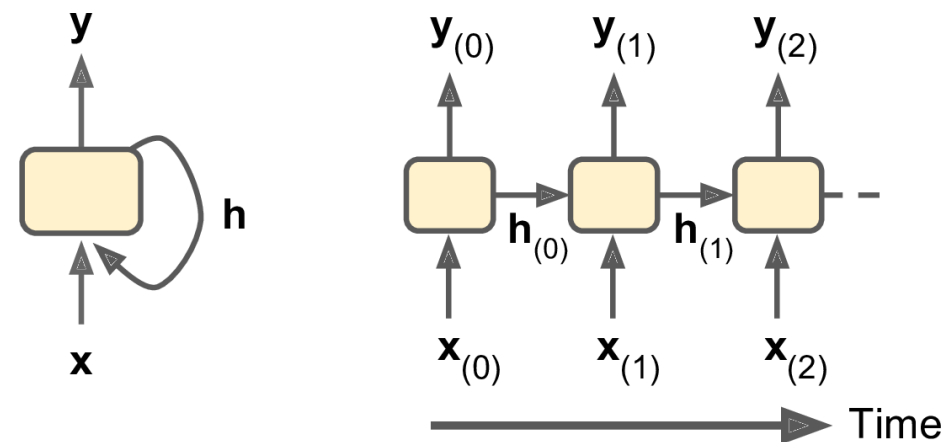


Previously in Course 3

Recurrent Neural Network (RNN): Definition

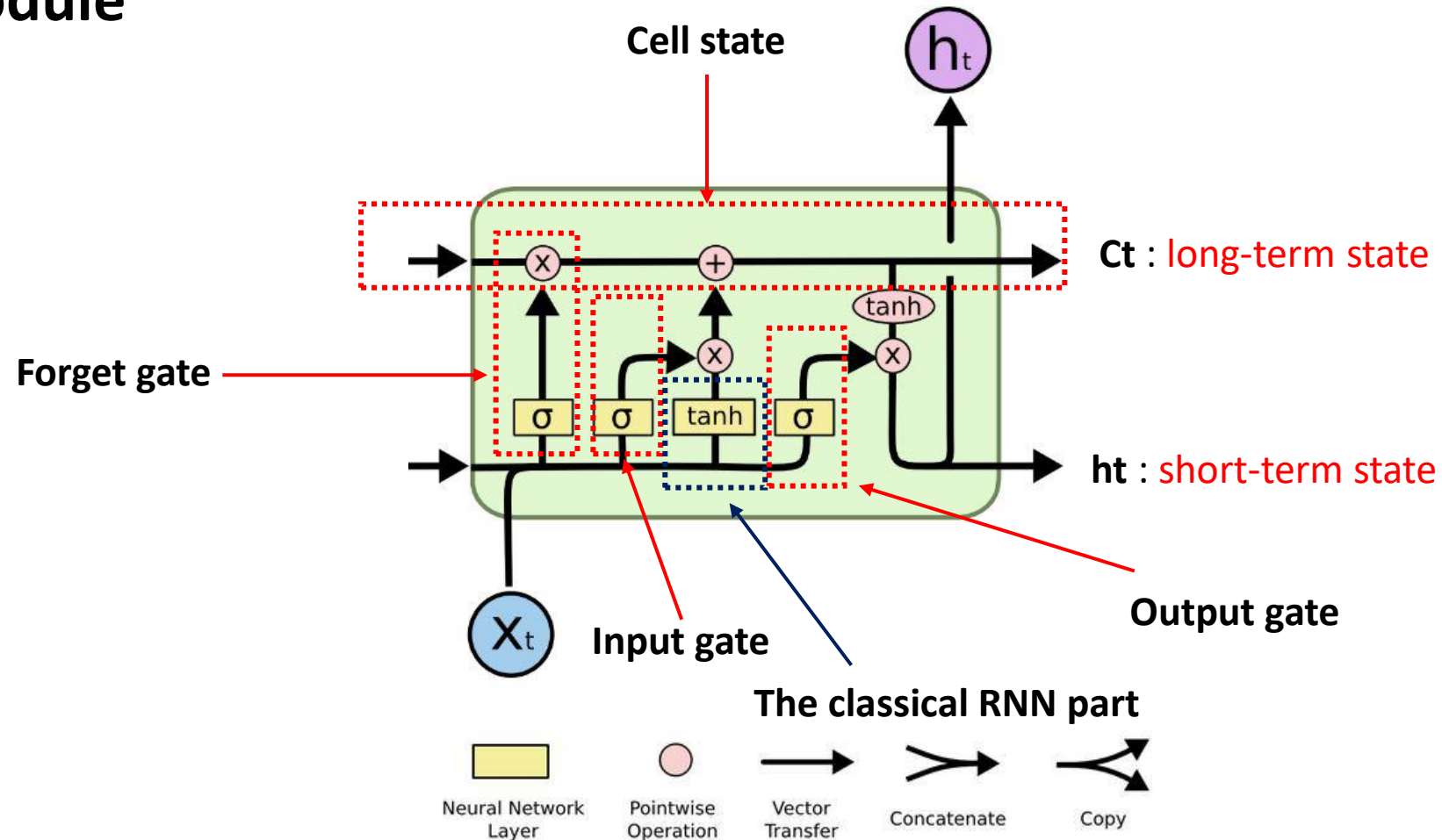
- A **RNN** is a kind of neural network dealing with **recurrent** connection
- Allows to deal with **temporal sequences**
- E.g., on the figure below, a sequence **X** is given as **input**, and we get a sequence **Y** as **output**
- A cell hidden state **h** and the output **y** may be different

A hidden state RNN



LSTM (Long Short Term Memory)

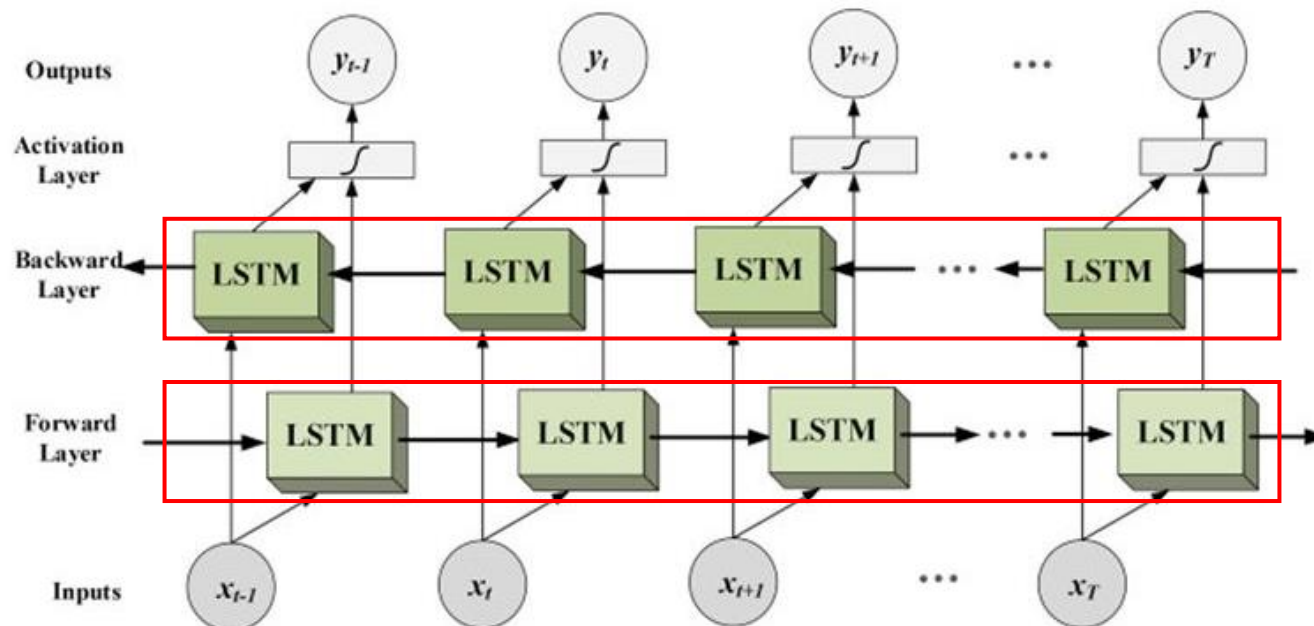
LSTM module



Bidirectional LSTM

In a **classic LSTM** architecture, the sentence/sequence is processed only in **one direction** (generally from the past to the present or future)

The **bidirectional** structure process the data in **both directions** => can be useful for some applications



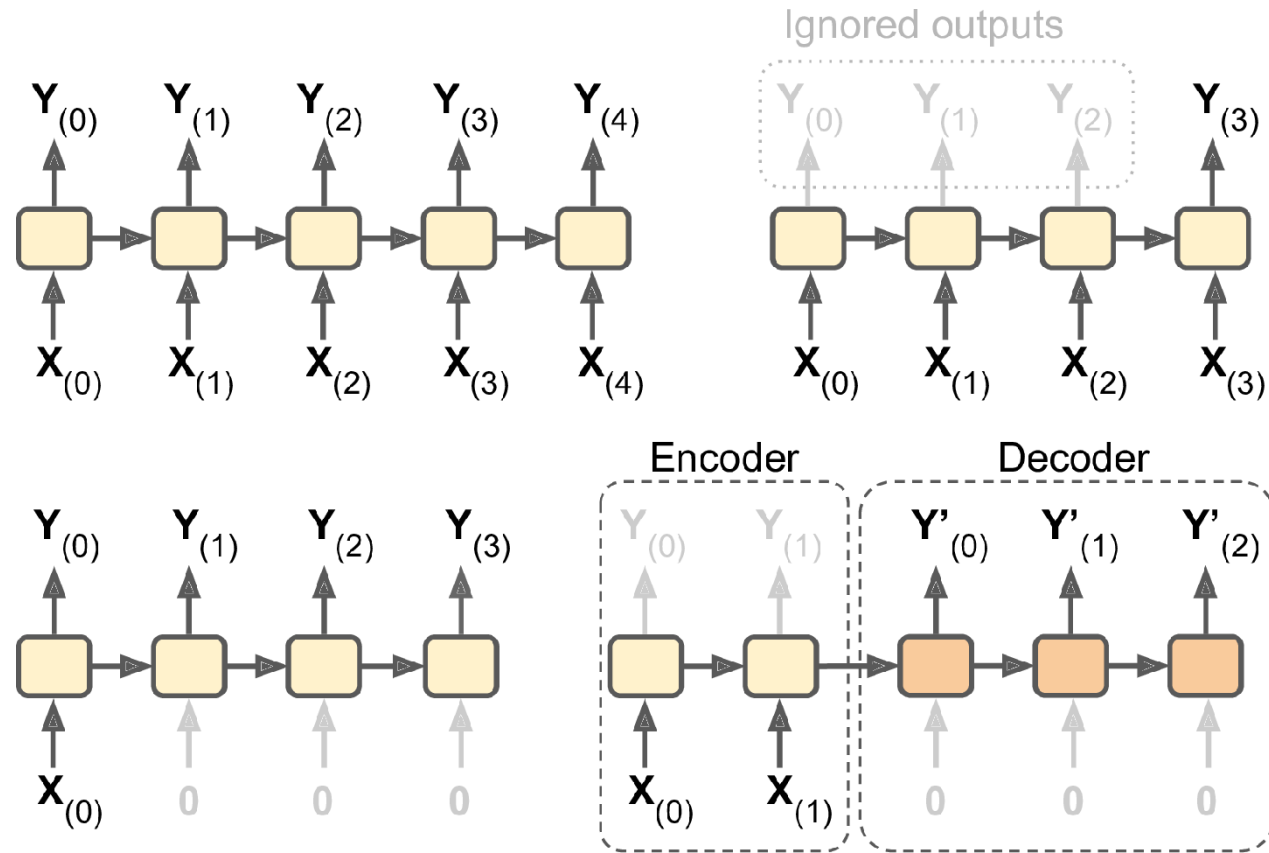
The **Queen** of the United Kingdom

The **queen** of hearts

The **queen** of bees

Here, the word queen will be encode differently if read in reverse order

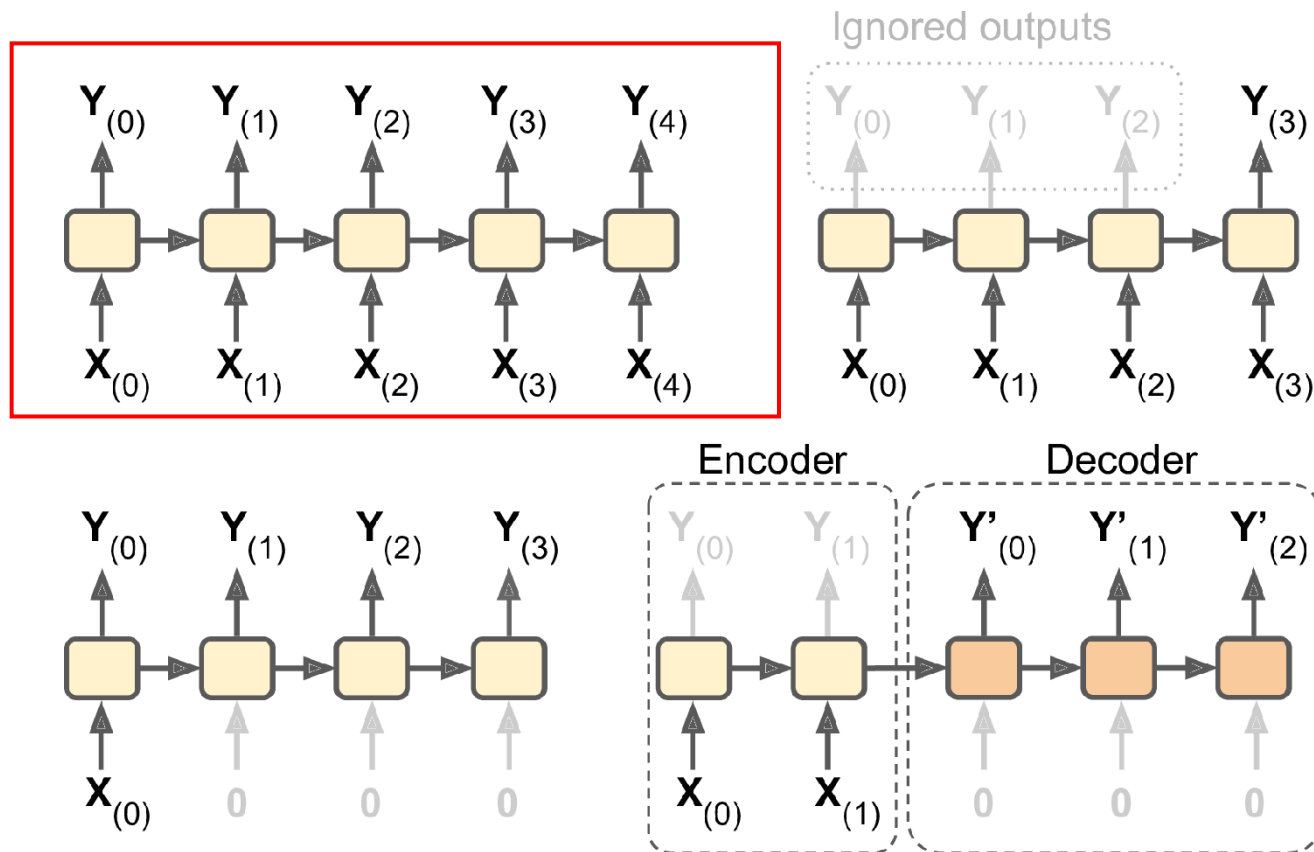
Recurrent Neural Network (RNN): Use cases



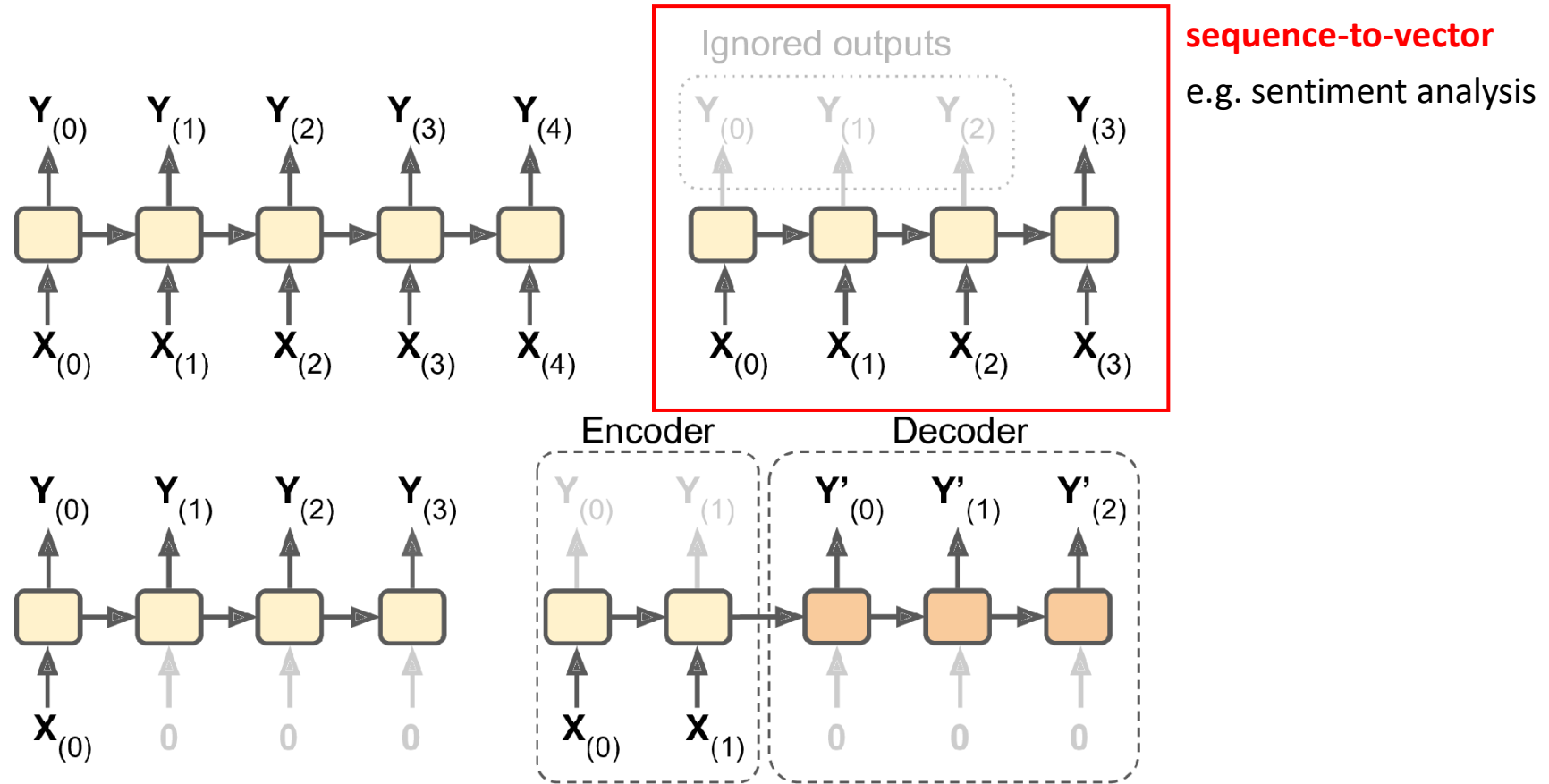
Recurrent Neural Network (RNN): Use cases

sequence-to-sequence

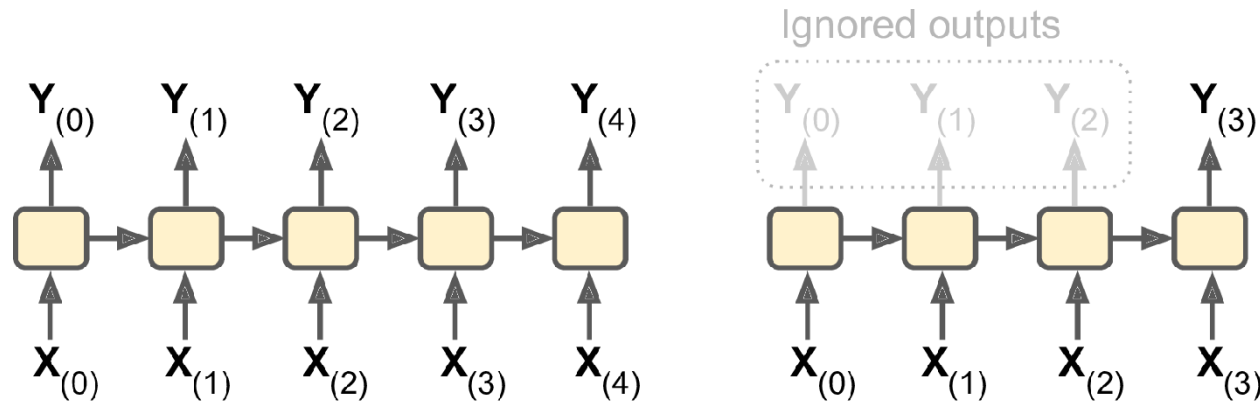
e.g. predicting time series



Recurrent Neural Network (RNN): Use cases

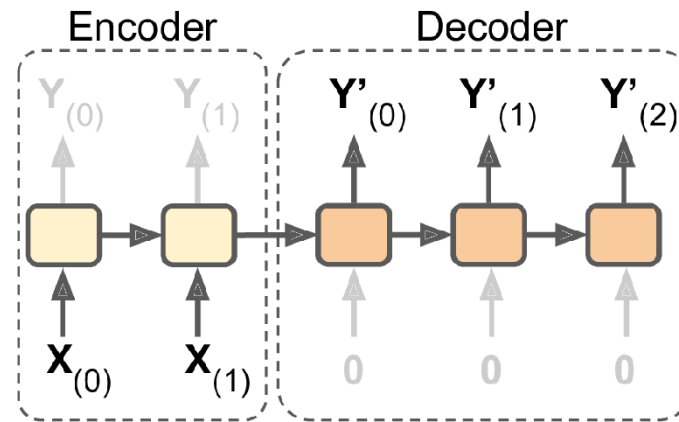
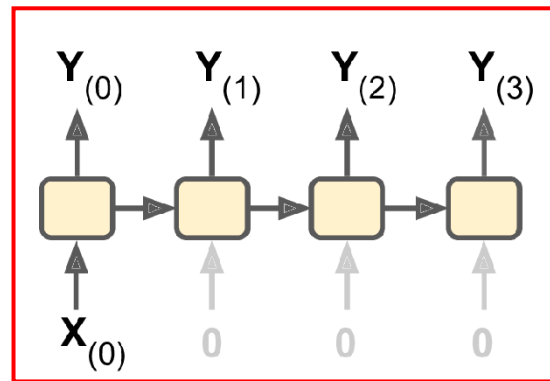


Recurrent Neural Network (RNN): Use cases

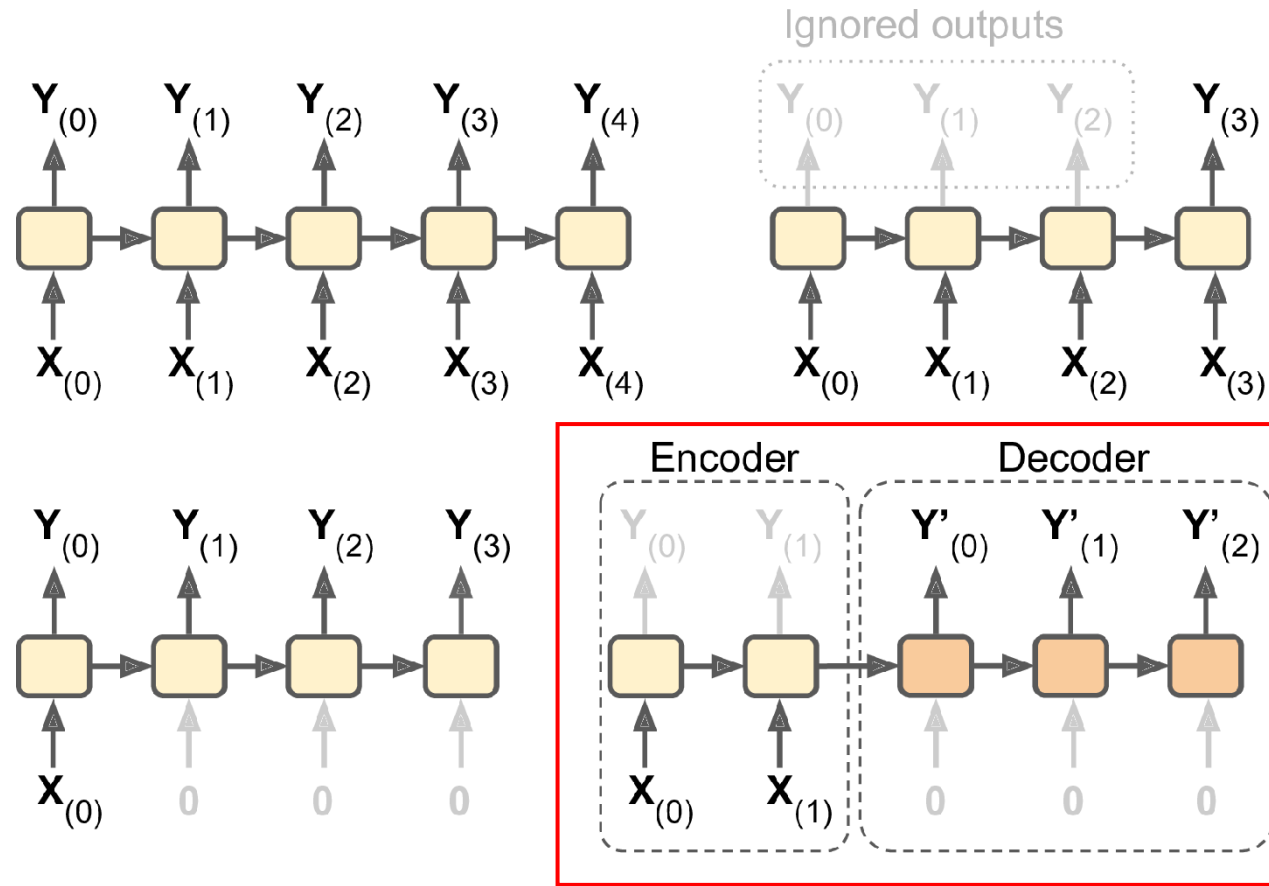


vector-to-sequence

e.g. image caption generation



Recurrent Neural Network (RNN): Use cases



sequence-to-vector
e.g. translation

Course 4: Attention Mechanism and Transformer Architectures

Encoder-decoder

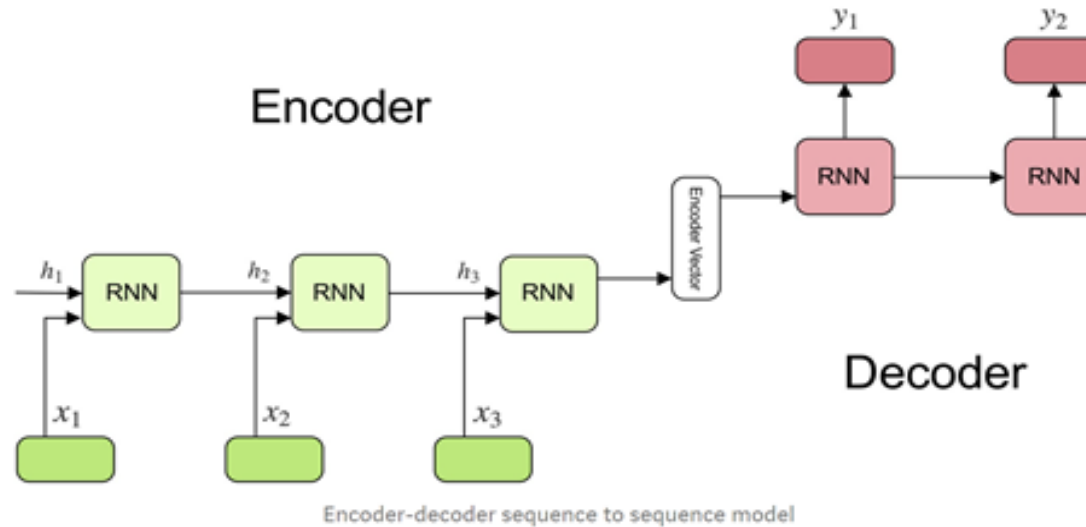
Encoder-decoder: A core concept

Principle

The goal is to project the input values into a smaller vector space before returning into the larger vector space expected

Architecture

- An encoder
- An encoded vector
- A decoder



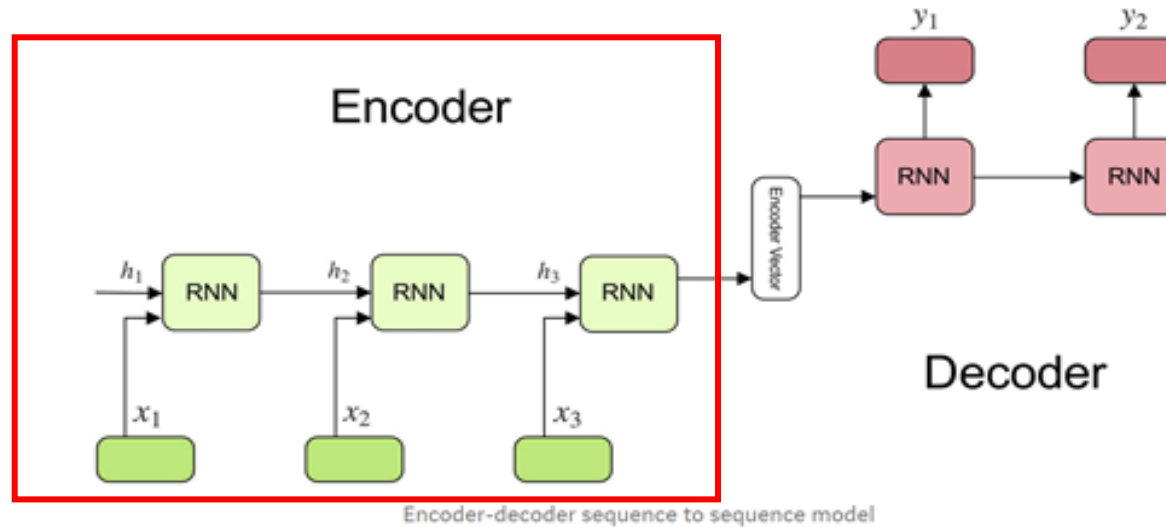
Encoder-decoder: A core concept

Principle

The goal is to project the input values into a smaller vector space before returning into the larger vector space expected

Architecture

- **An encoder**
- An encoded vector
- A decoder



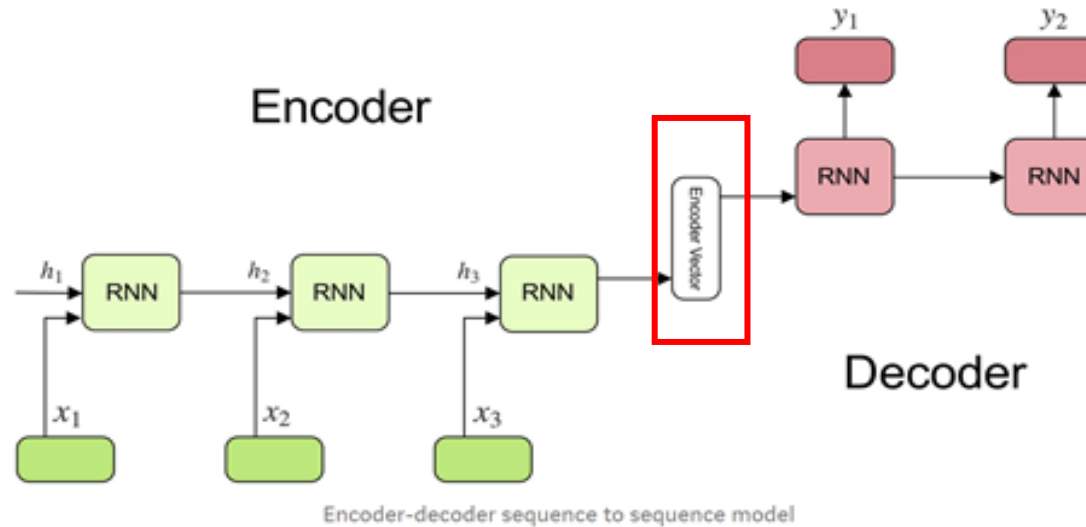
Encoder-decoder: A core concept

Principle

The goal is to project the input values into a smaller vector space before returning into the larger vector space expected

Architecture

- An encoder
- **An encoded vector**
- A decoder



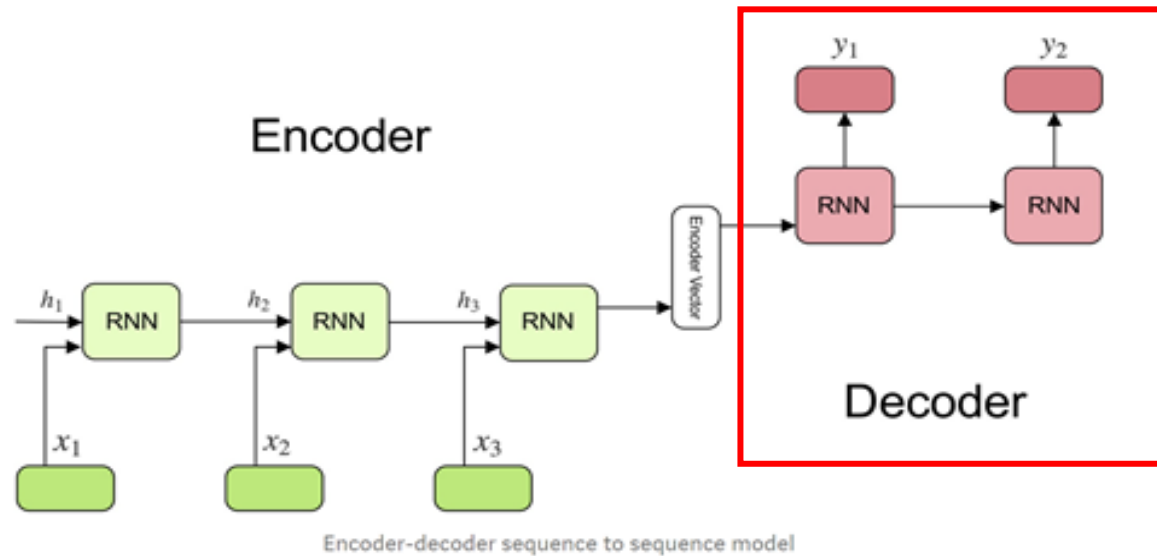
Encoder-decoder: A core concept

Principle

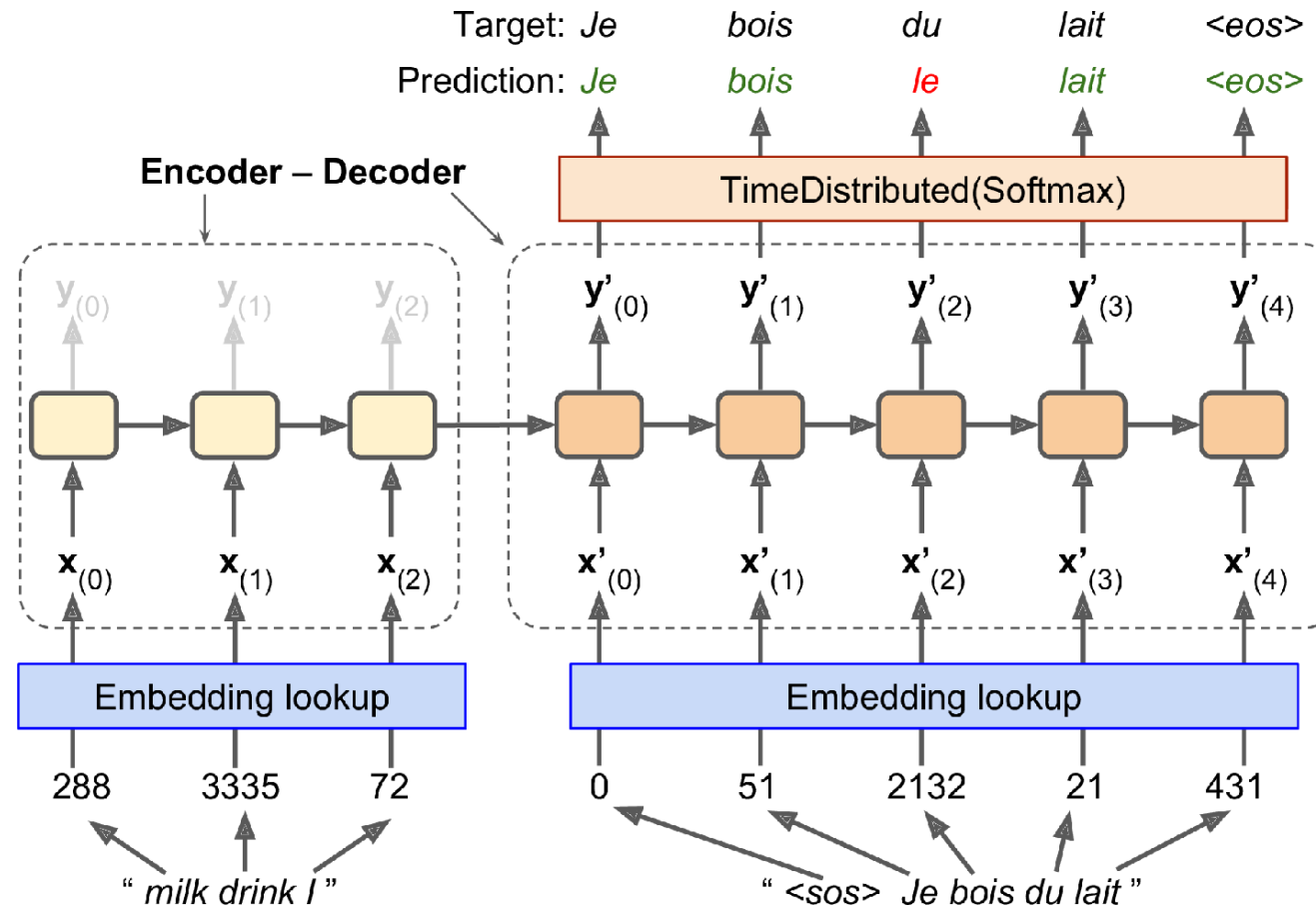
The goal is to project the input values into a smaller vector space before returning into the larger vector space expected

Architecture

- An encoder
- An encoded vector
- **A decoder**



Encoder-decoder: Translation example



Encoder-decoder: Keras

```
import tensorflow_addons as tfa

encoder_inputs = keras.layers.Input(shape=[None], dtype=np.int32)
decoder_inputs = keras.layers.Input(shape=[None], dtype=np.int32)
sequence_lengths = keras.layers.Input(shape=[], dtype=np.int32)

embeddings = keras.layers.Embedding(vocab_size, embed_size)
encoder_embeddings = embeddings(encoder_inputs)
decoder_embeddings = embeddings(decoder_inputs)

encoder = keras.layers.LSTM(512, return_state=True)
encoder_outputs, state_h, state_c = encoder(encoder_embeddings)
encoder_state = [state_h, state_c]

sampler = tfa.seq2seq.sampler.TrainingSampler()

decoder_cell = keras.layers.LSTMCell(512)
output_layer = keras.layers.Dense(vocab_size)
decoder = tfa.seq2seq.basic_decoder.BasicDecoder(decoder_cell, sampler,
                                                  output_layer=output_layer)

final_outputs, final_state, final_sequence_lengths = decoder(
    decoder_embeddings, initial_state=encoder_state,
    sequence_length=sequence_lengths)
Y_proba = tf.nn.softmax(final_outputs.rnn_output)

model = keras.Model(inputs=[encoder_inputs, decoder_inputs, sequence_lengths],
                    outputs=[Y_proba])
```

Encoder-decoder: Keras

encoder

```
import tensorflow_addons as tfa

encoder_inputs = keras.layers.Input(shape=[None], dtype=np.int32)
decoder_inputs = keras.layers.Input(shape=[None], dtype=np.int32)
sequence_lengths = keras.layers.Input(shape=[], dtype=np.int32)

embeddings = keras.layers.Embedding(vocab_size, embed_size)
encoder_embeddings = embeddings(encoder_inputs)
decoder_embeddings = embeddings(decoder_inputs)

encoder = keras.layers.LSTM(512, return_state=True)
encoder_outputs, state_h, state_c = encoder(encoder_embeddings)
encoder_state = [state_h, state_c]
```

allow to get both ht and Ct

decoder

```
decoder_cell = keras.layers.LSTMCell(512)
output_layer = keras.layers.Dense(vocab_size)
decoder = tfa.seq2seq.basic_decoder.BasicDecoder(decoder_cell, sampler,
                                                  output_layer=output_layer)

final_outputs, final_state, final_sequence_lengths = decoder(
    decoder_embeddings, initial_state=encoder_state,
    sequence_length=sequence_lengths)
Y_proba = tf.nn.softmax(final_outputs.rnn_output)
```

to tell the decoder at each step
what was the previous output

model construction

```
model = keras.Model(inputs=[encoder_inputs, decoder_inputs, sequence_lengths],
                    outputs=[Y_proba])
```

Attention Mechanism

Attention Mechanism: Introduction

Problem to solve

- Even with a LSTM architecture, it can be difficult to characterize long sentences in an efficient way

Attention mechanism principle

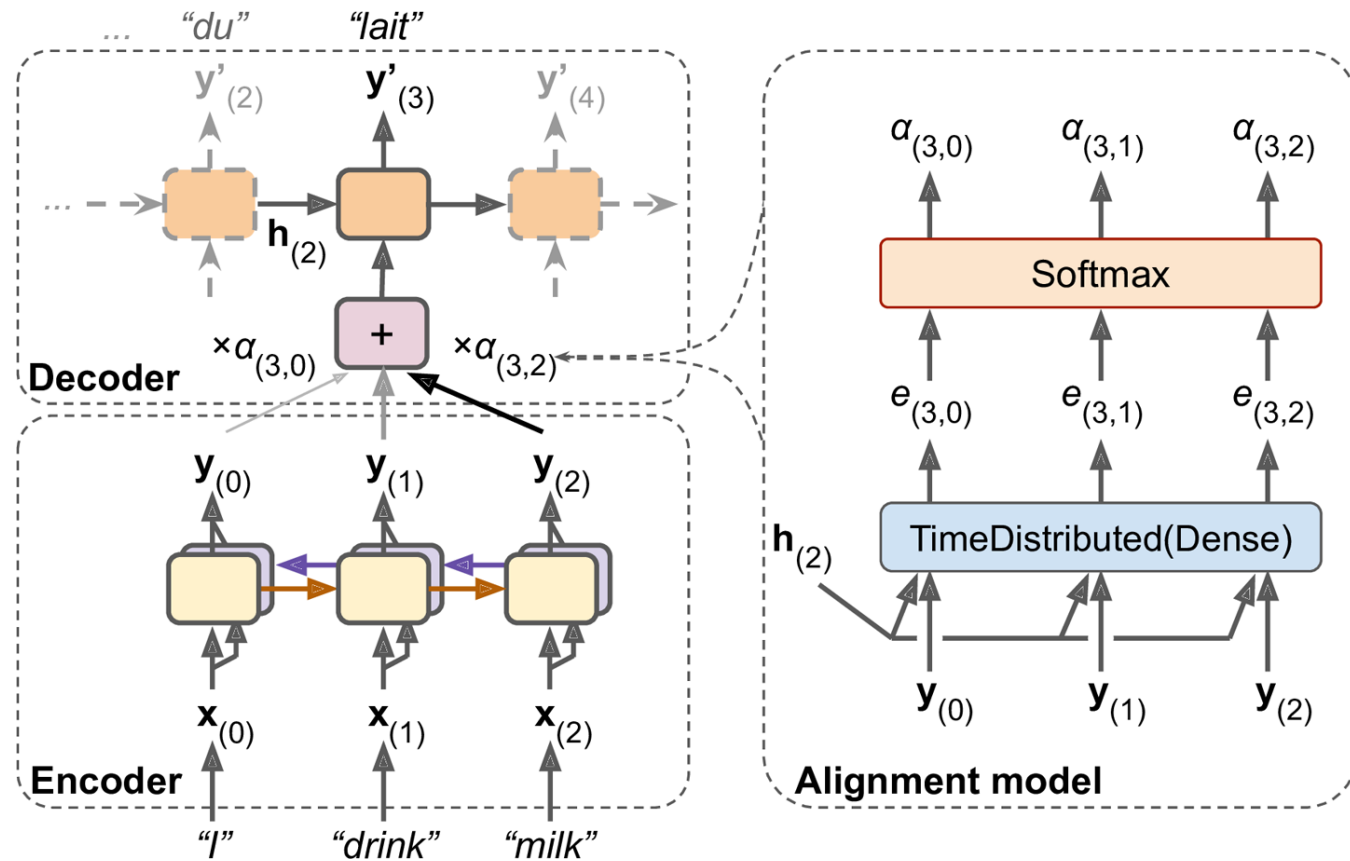
- To put more weight on specific word positions of the sentence depending **both** on the **context** and the **relative positions of words**

Example

Despite being from Uttar Pradesh, as he was brought up in Bengal, he is more comfortable in Bengali.

To predict the word *Bengali*, the expressions *brought up* and *Bengal* must be associated to higher weights.

Attention Mechanism: Translation illustration

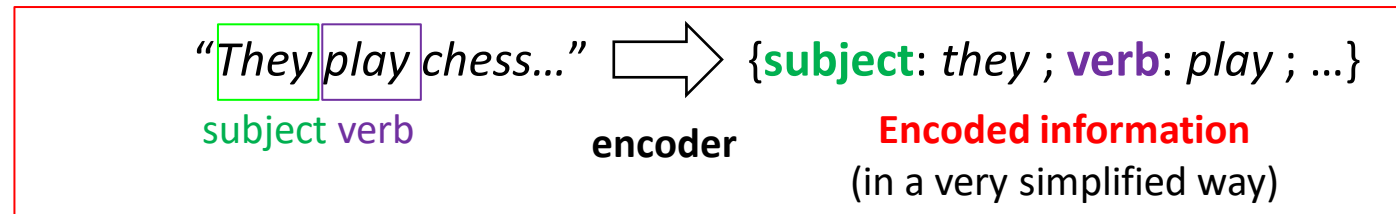


Attention Mechanism: Get the intuition

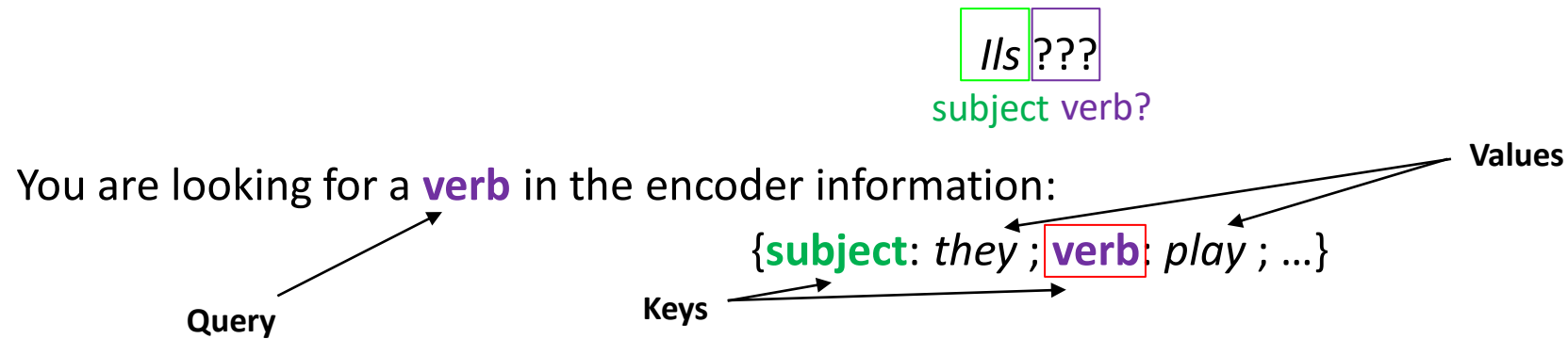
Imagine you want to translate the sentence using an **encoder-decoder**:

"They play chess..."

Actually, you expect the encoder to deal with it that way:



Just imagine you already translated the **first word** and are looking for the **next one**:

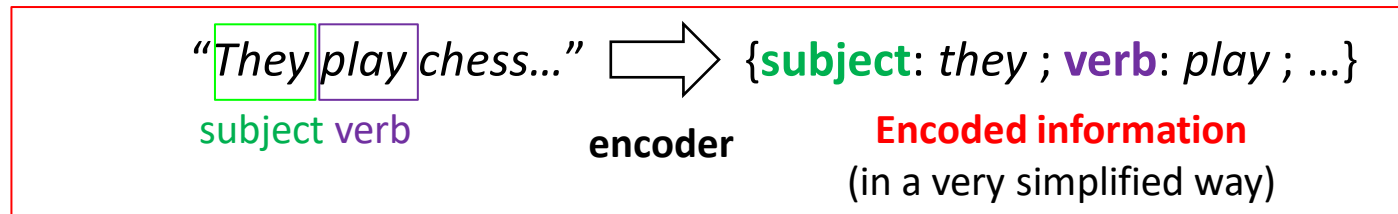


Attention Mechanism: Get the intuition

Imagine you want to translate the sentence using an **encoder-decoder**:

"They play chess..."

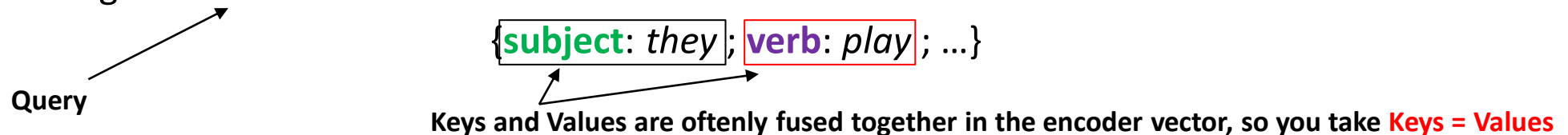
Actually, you expect the encoder to deal with it that way:



Just imagine you already translated the **first word** and are looking for the **next one**:

//s ???
subject verb?

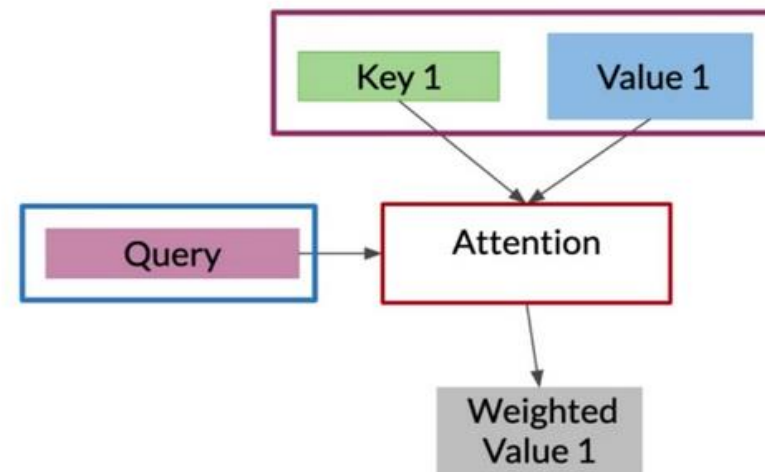
You are looking for a **verb** in the encoder information:



Attention Mechanism: Principle

Formalism

- To associate a *query* and a set of (*key-value*) tuples with an *output*
- The *query*, the (*key-value*) set and the *output* are all vectors
- The output corresponds to the **weighted sum** of values, with weights determined from a **compatibility function** between the query and the corresponding key



Attention Mechanism: Formula

Formula

The mathematical expression generally used to compute the attention function is

$$\text{softmax}(QK^T)V$$

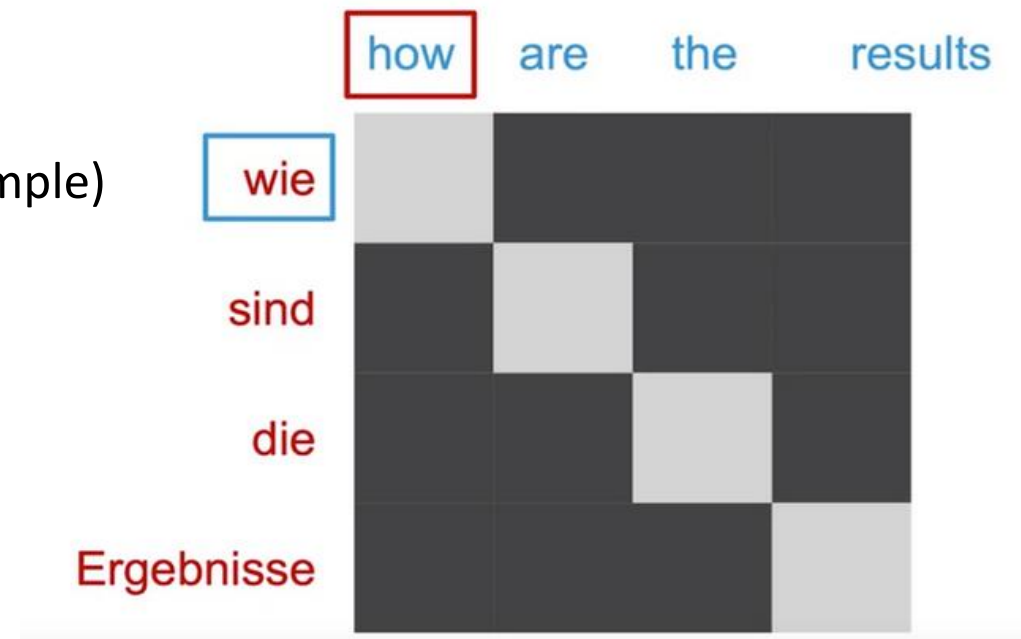
With Q , K and V respectively the *query*, *key* and *value*, matrices.

The ***softmax*** function gives normalized values between 0 and 1

Attention Mechanism: Illustration 1

Example for the case of English-German translation

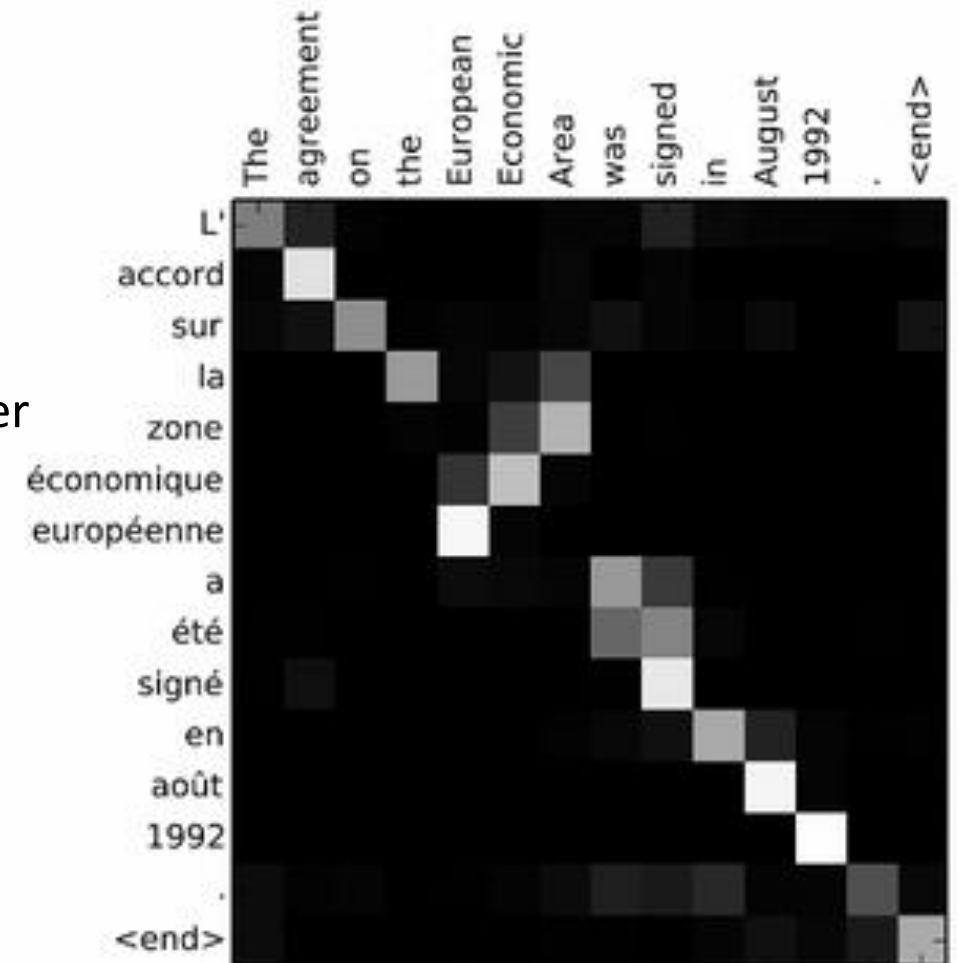
- The bright squares show where the algorithm « **look** » to translate a word
- Here, words in English and German **share the same order**
- So, the match is only on the **diagonal**
- No need to look other word positions (on this specific example)



Attention Mechanism: Illustration 2

Example for the case of English-French translation

- Here, the match is **not on a specific word only**
- The figure shows the contribution of each word
- **Several words** can contribute to the translation of another



Attention Mechanism: Examples

Intuition behind **query** and (**key**, **value**) concepts

Depends on the application

In NLP, usually, the **key** and **value** correspond to the **same vectors**:

- **Translation**: the **query** can be the encoded sentence vector in **one language** and both **key** and **value** can be the **encoded** sentence vector in the **other language**
- **Text similarity**: the **query** can be the sequence embedding of the **first piece of text** and both **key** and **value** can be the sequence embedding of the **second piece of text**

However, in some cases, the concepts key and value **can be different**:

- if you are looking for a video in *Youtube*, the **query** can be derived from the **text you entered**, the **keys** can correspond to **the video descriptions** and the **values** can be **the videos themselves**

Attention Mechanism: Keras Layer

An Attention Layer is available in Keras

Here we can see an illustration for translation

[illegible]

Attention Mechanism: Keras Layer

An Attention Layer is available in Keras

Here we can see an illustration for translation

```
# Attention layer
attention_layer = AttentionLayer(name='attention_layer')
attention_out, attention_states = attention_layer({"values": encoded_en,
                                                  "query": encoded_fr})
```

Here, the **values** correspond to the encoded **English** sentences

As usual in that kind of context, the **keys** are implicitly the same as the **values** (the encoded English sentences)

Here, the **queries** correspond to the encoded **French** sentences

Visual Attention: Illustration

Goal: Generate automatically the caption of a picture with an attention mechanism



“A woman is throwing a frisbee in a park”

Visual Attention: Illustration

Goal: Generate automatically the caption of a picture with an attention mechanism



*“A woman is throwing a **frisbee** in a park”*

Transformer architecture

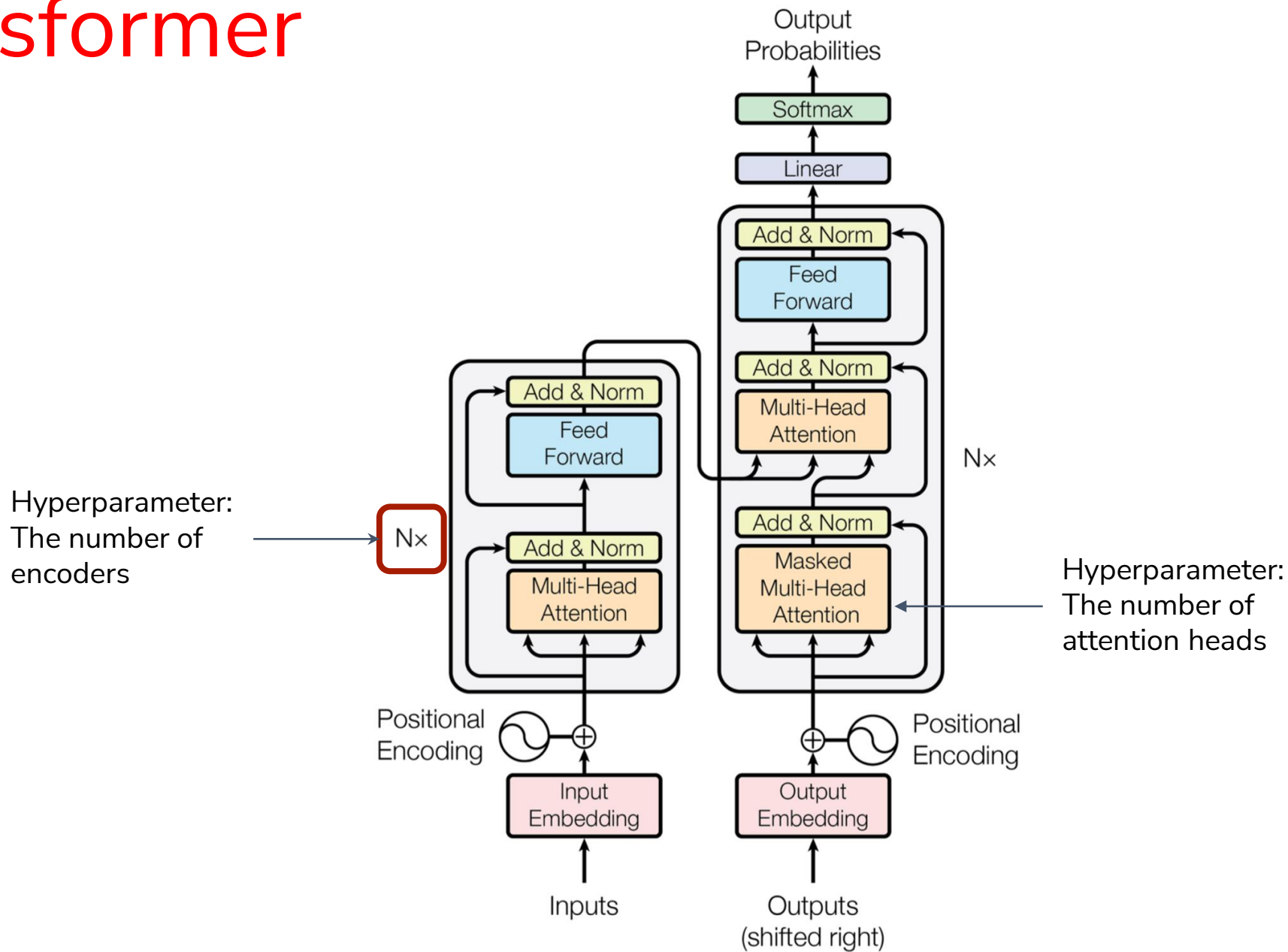
Transformer: Introduction

Principle

- A quite new deep learning model (2017), introduced in the **game changing** article *Attention is all you need*
- Relies heavily on the **Attention mechanism**
- **Does not need to process a sequence in a specific order**
- Solves the issue of keeping into memory the information related to distant words (and that even **without LSTM**)
- Makes possible a significant use of **parallelization** computing
- At every moment, the algorithm can access the complete set of the successive states visited during the procedure

The main idea is that the **Attention mechanism alone**, without any recurrent sequential procedure, **is powerful enough to reach the state of the art**

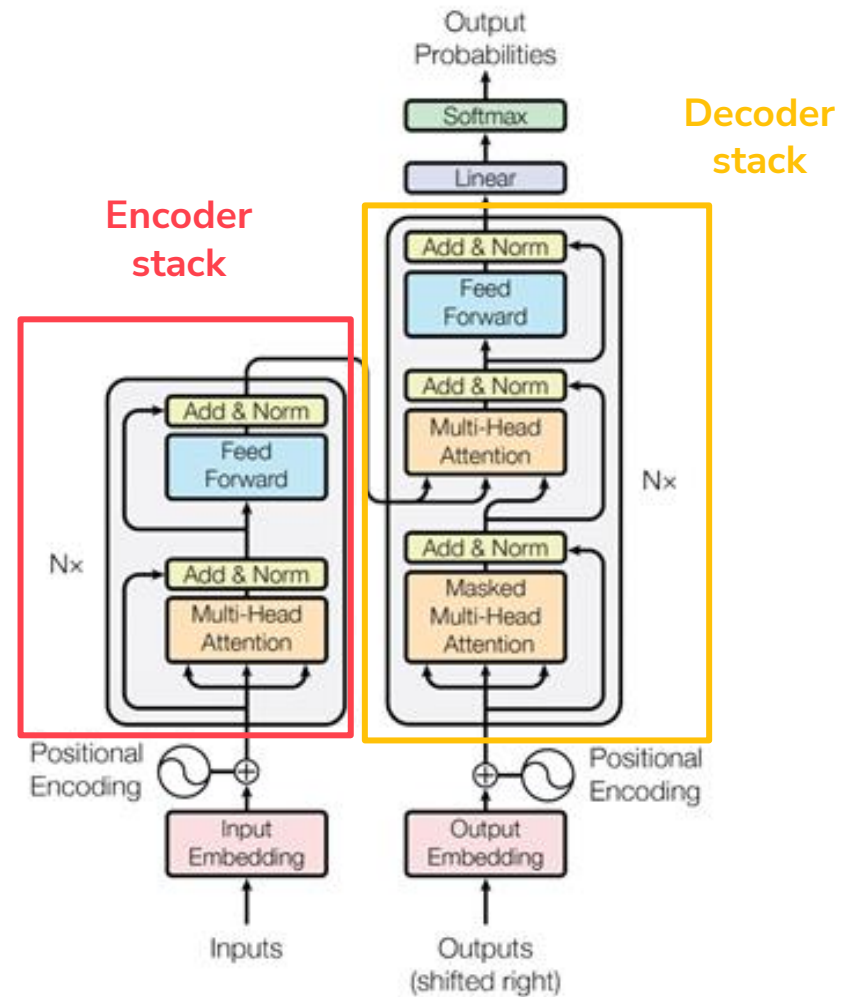
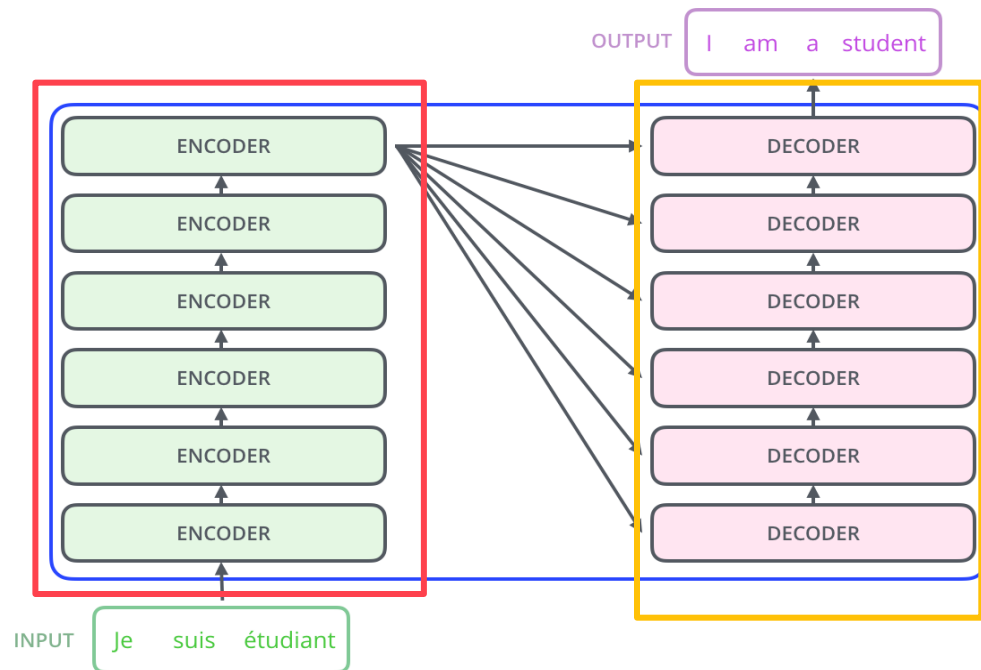
Transformer



Transformer: Architecture

The Transformer is made of two main components:

- A stack of **identical encoders** (independents from each other)
- Followed by a stack of **identical decoders** (independents from each other)

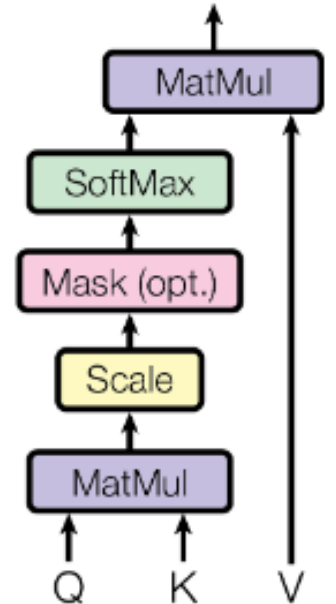


Transformer: Scaled Dot-Product Attention

Formula

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

d_k corresponds to the query and key dimension



- The **normalization** prevents the result of the dot-product to reach high values and, doing so, **prevents the gradient to vanish** into small values
- An **additive** variant of the attention function exists, but the dot-product version is preferred here for the efficiency of the matrix calculations it allows

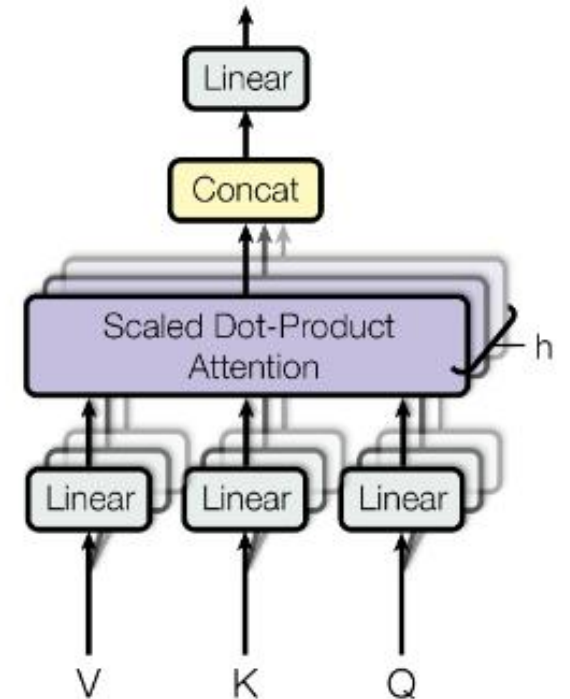
Transformer: Multi-Head Attention

Principle

- **Several Attention** layers computed in parallel
 - the *queries*, *keys* and *values* are projected h times
 - the h results are concatenated
 - then projected one last time
- Enable the use of different sub-space information
- The mathematical expression is the following

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W^O$$

where $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$

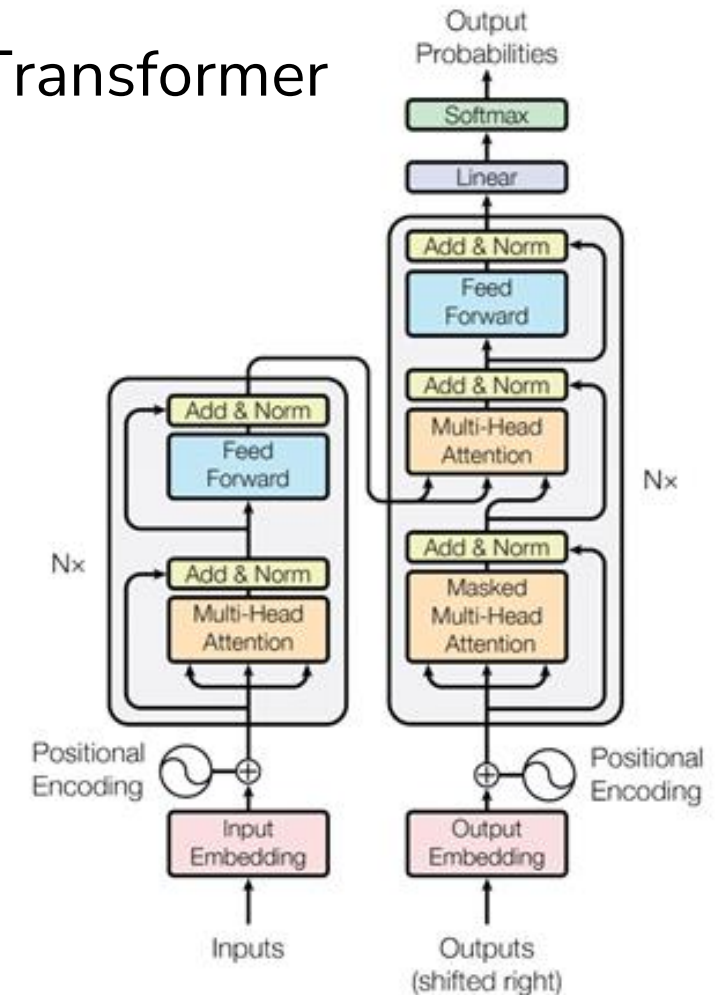


Where the W_i^* corresponds to the respective projection matrices

Transformer: Multi-Head Attention

The **Multi-Head Attention** is used three times in the Transformer

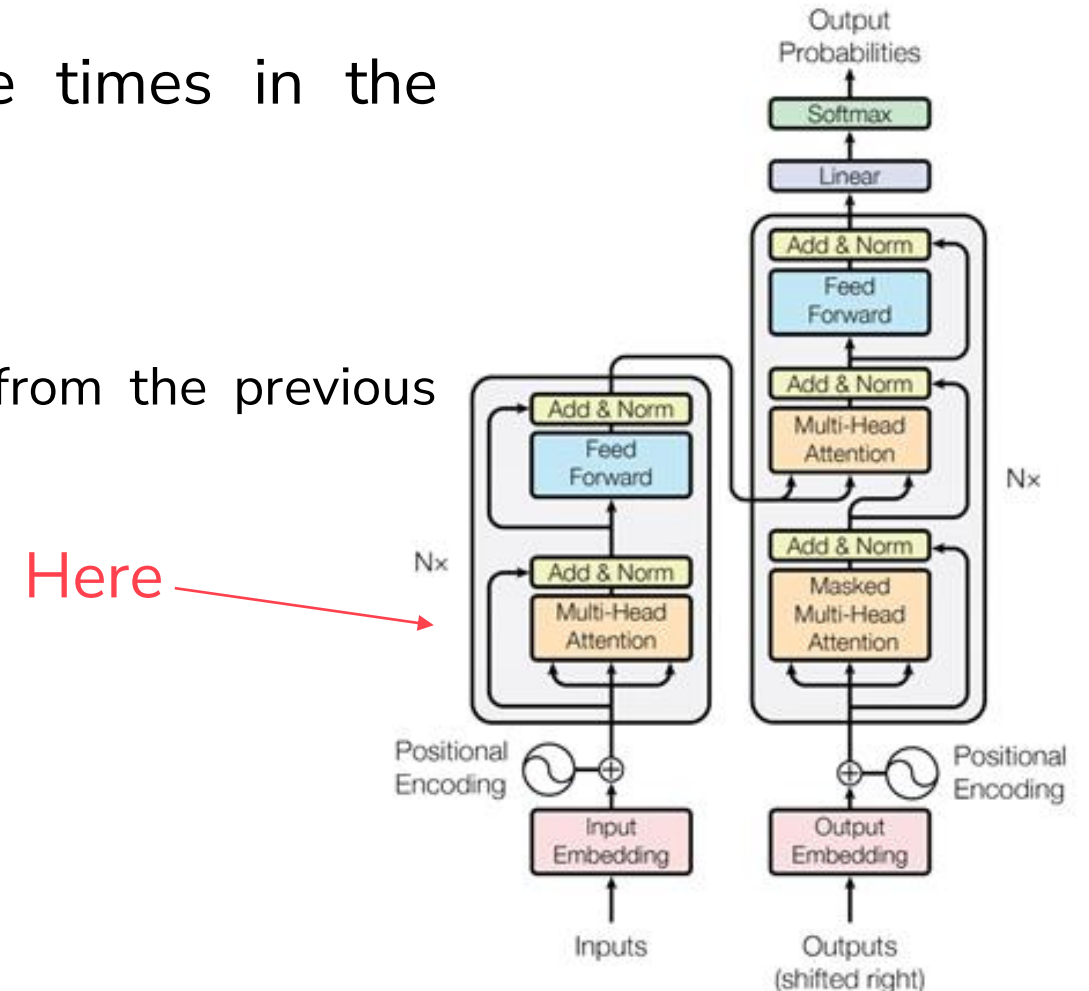
- Inside the encoder
- Inside the decoder
- At the encoder-decoder interface



Transformer: Multi-Head Attention

The **Multi-Head Attention** is used three times in the Transformer

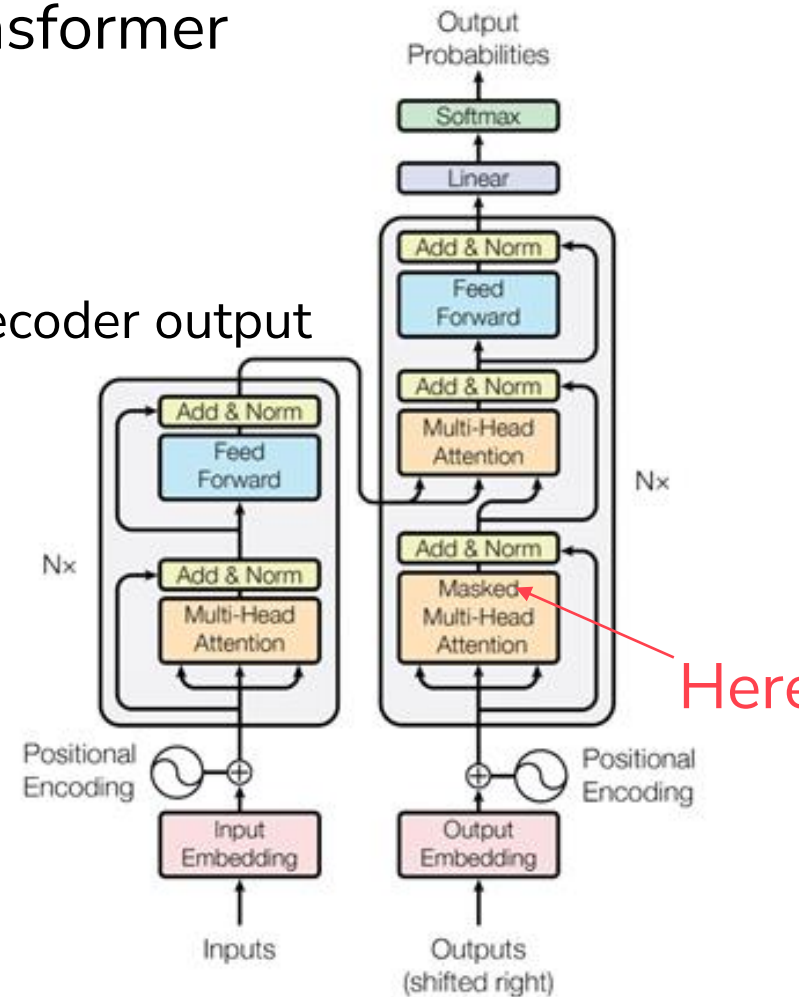
- **Inside the encoder**
 - the queries, keys and values all come from the previous encoder output
- Inside the decoder
- At the encoder-decoder interface



Transformer: Multi-Head Attention

The **Multi-Head Attention** is used three times in the Transformer

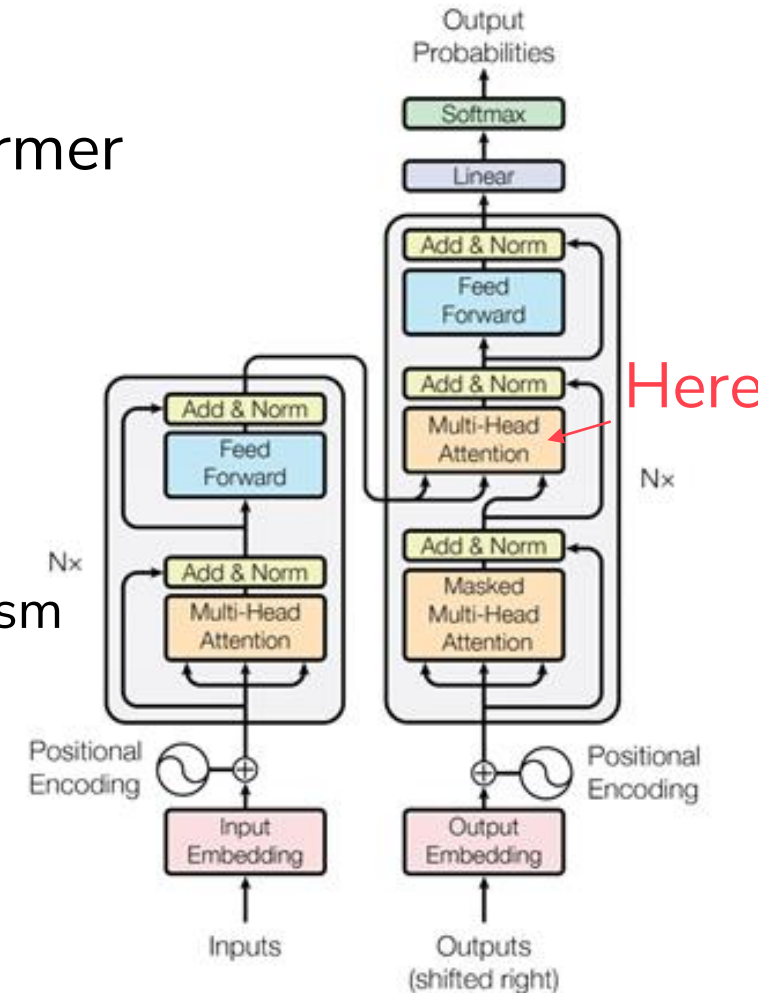
- Inside the encoder
- **Inside the decoder**
 - the queries, keys and values all come from the previous decoder output
 - use of a « mask » to decode only from known elements
- At the encoder-decoder interface



Transformer: Multi-Head Attention

The **Multi-Head Attention** is used three times in the Transformer

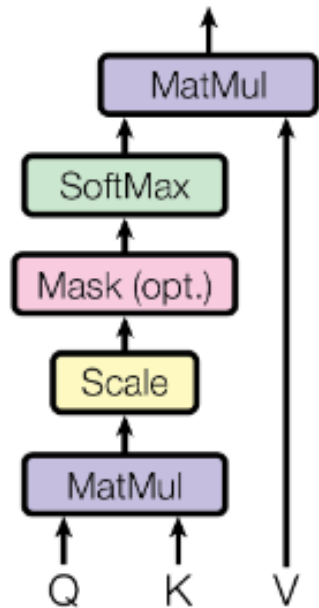
- Inside the encoder
- Inside the decoder
- **At the encoder-decoder interface**
 - the queries come from the previous decoder layer
 - the keys and values come from the encoder outputs
 - identical to the classical “encoder-decoder” Attention mechanism



Transformer: Attention

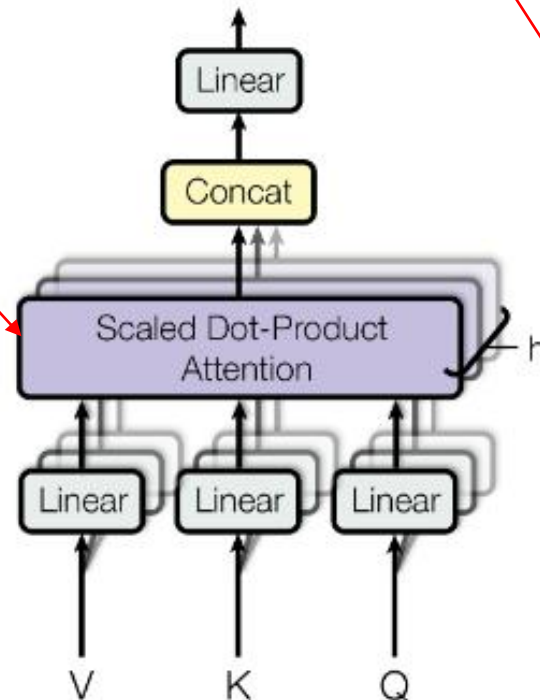
Transformer

Attention



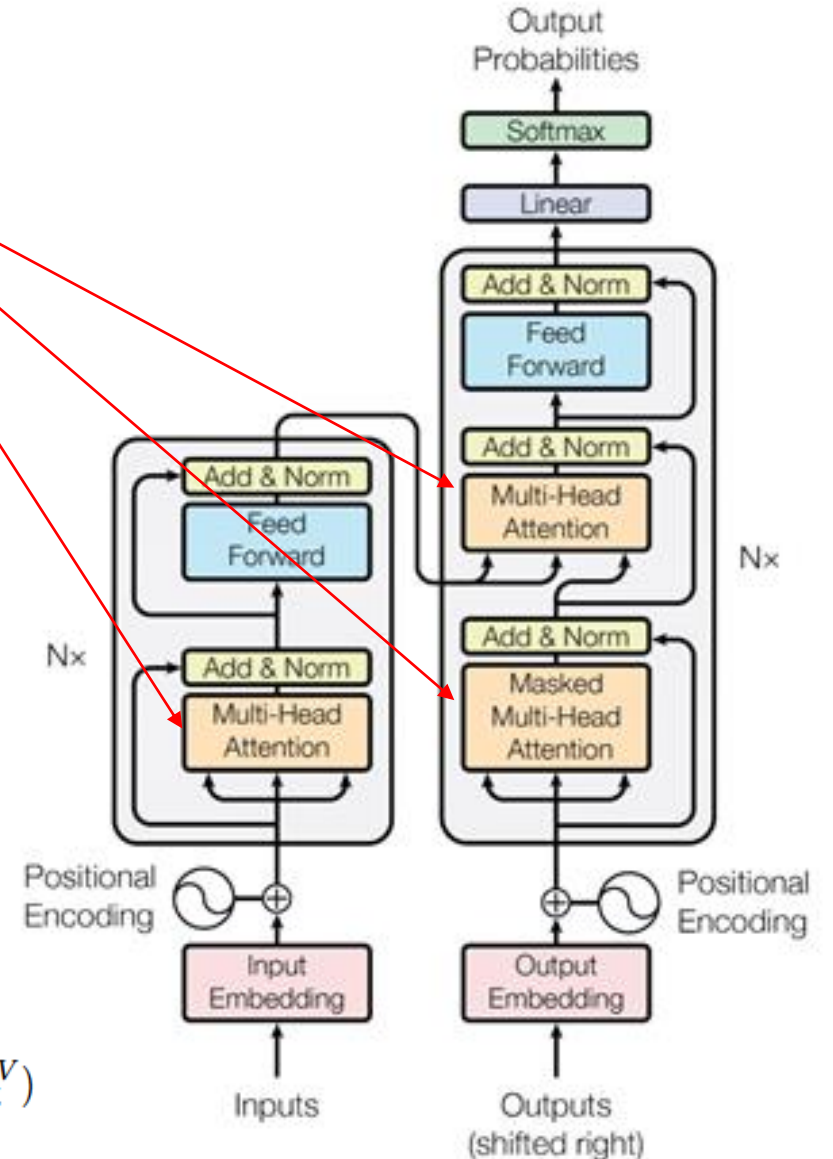
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

multi-head Attention



$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

where $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$



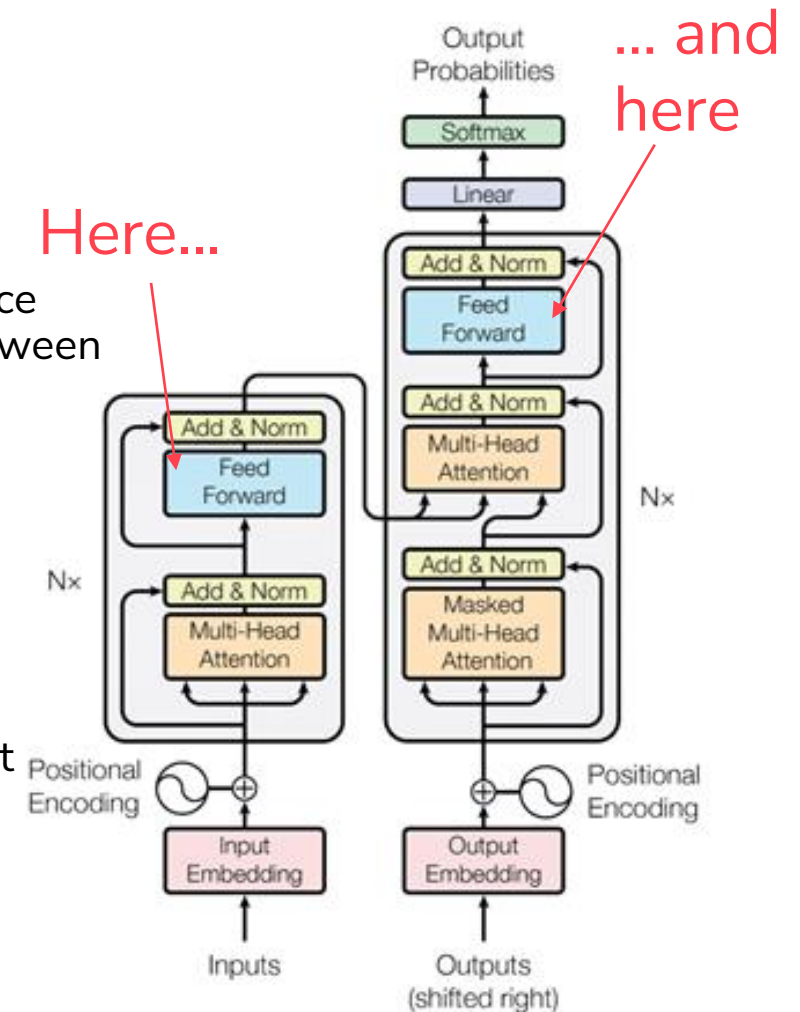
Transformer: Feed-Forward Networks

Principle

- **Feed-Forward network** in each **encoder** and each **decoder** layer
- Each of these Feed-Forward network is completely connected
- Applied in an independent and identical way to each element in the sequence
- Consist in two linear transformations, with the use of a ReLU function in between

$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

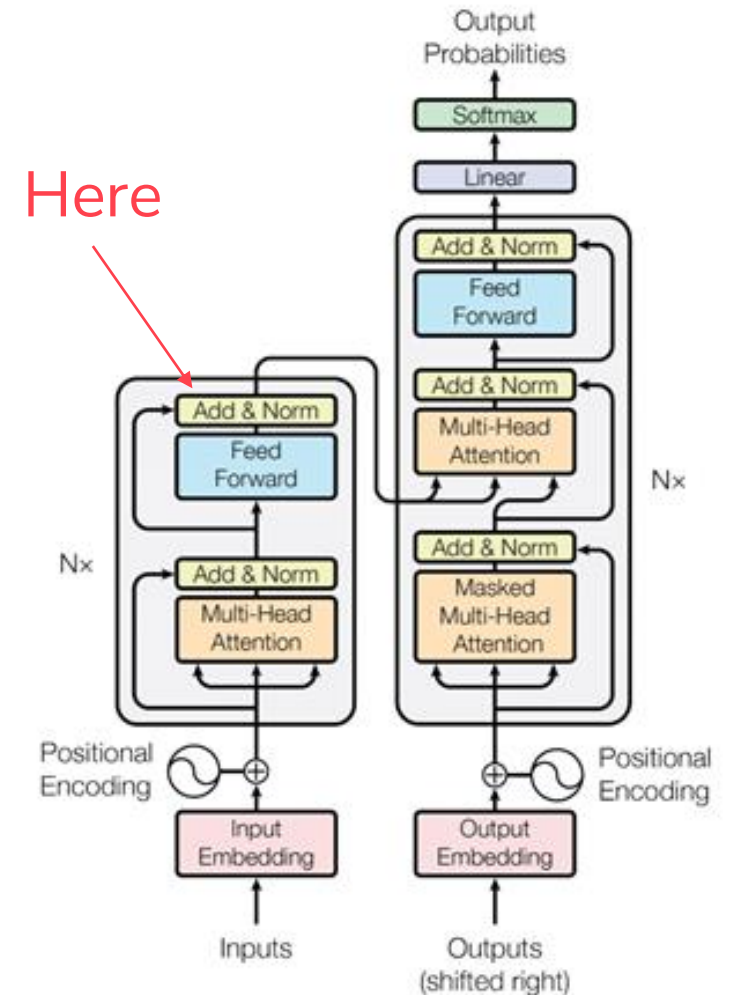
- The linear transformations are identical for each element in the sequence...
- ... but are different for each layer
- Its role is to process the output from one attention layer to better fit the next



Transformer: The residual connection

Principle

- **Add**: To strengthen the Transformer memory.
- **Normalisation**: To reduce the needed steps to optimize the network.



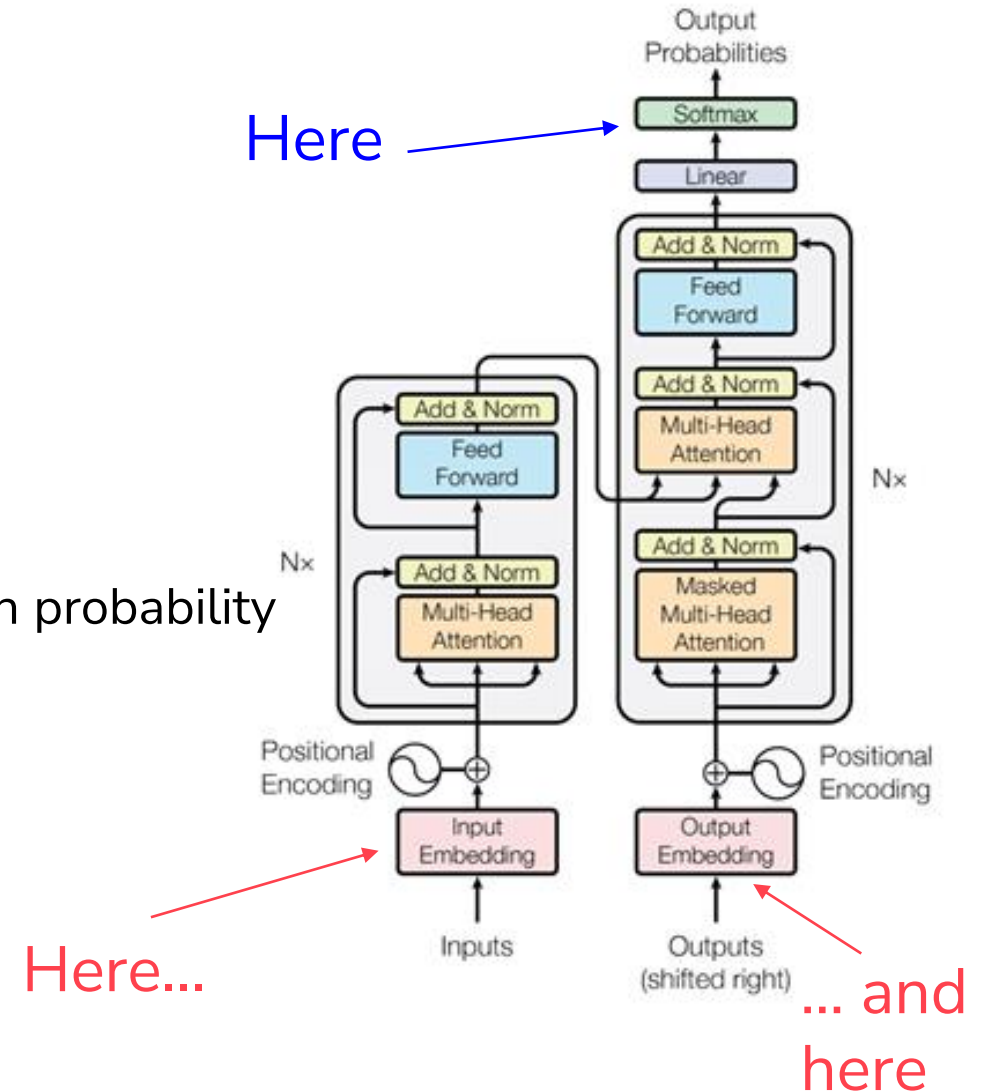
Transformer: Embeddings et Softmax

Embedding learning in order to convert

- The input *tokens*
- The output *tokens*

Use of a *Softmax* layer

- Convert the *decoder* output into the next token prediction probability



Transformer: Positional Encoding

Principle

- Keep in memory the *token* positions in the sequence
- Add of “*positional encodings*” to the input *embeddings*
- Used just before the *encoder-decoder* stacks

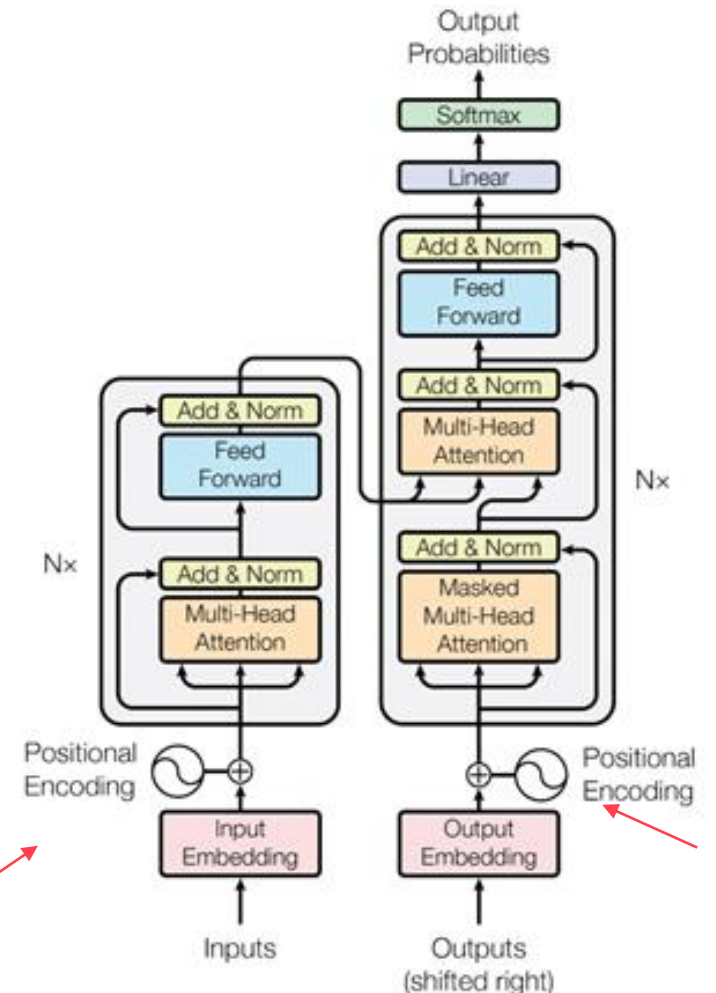
Formulas

$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

- with pos the position and i the dimension
- Each dimension corresponds to a sinusoid
- For each k , PE_{pos+k} is a linear function of PE_{pos}

Here...



... and here

Transformer: Results on the Benchmark

BLEU

Model	BLEU		Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet	23.75			
Deep-Att + PosUnk		39.2		$1.0 \cdot 10^{20}$
GNMT + RL	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4 \cdot 10^{20}$
ConvS2S	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$
MoE	26.03	40.56	$2.0 \cdot 10^{19}$	$1.2 \cdot 10^{20}$
Deep-Att + PosUnk Ensemble		40.4		$8.0 \cdot 10^{20}$
GNMT + RL Ensemble	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$
ConvS2S Ensemble	26.36	41.29	$7.7 \cdot 10^{19}$	$1.2 \cdot 10^{21}$
Transformer (base model)	27.3	38.1	$3.3 \cdot 10^{18}$	
Transformer (big)	28.4	41.8	$2.3 \cdot 10^{19}$	

Transformer: Tutorial notebook

course4_transformer_tutorial.ipynb

Goal: Illustration of how a Transformer can be build step by step in Python

Remarks:

- This notebook comes from a **tutorial notebook** **freely accessible** (the original reference is given at the very beginning, in the first cell)

BERT model

BERT: Introduction

BERT (2018) is for ***B**idirectional **E**ncoder **R**epresentations from **T**ransformers*

- A Transformer variant with a **Bidirectional** specificity
- Thanks to it, there is no need for the end sequence token masking constraint anymore (into the decoders Attention mechanism)
- Build an efficient and generic language representation

The two main BERT steps

- A **pre-training** step on non-labelled data
- A **specialization** step (“***fine tuning***”) for a specialized task, from pre-trained parameters
 - There is only need to train **the « last » layers** of the model to get such a specialization

French versions of BERT exists, the most famous one being probably **CamemBERT**

BERT: Architecture

The initial **BERT** architecture (2018) is the following:

- Length of the processed sequences: **512** (Limitation to contain a quadratic explosion)
- Stack of **Nx = 12 blocks** of encoders and decoders (a 24-blocks variant exists)
- **h = 12** attention heads
- It means there are $h \times Nx = 12 \times 12 = 144$ distinct attention mechanisms
- Each head gives an output vector of size: **64**
- A hidden layer of $h \times 64 = 768$ dimensions
- Roughly **110 millions** of parameters for the basic version (12 encoders) and **340 millions** for the large version (24 encoders)

Reminder, the initial Transformer architecture (2017)

- Length of the processed sequences: **512** as well
- Stack of **Nx = 6 blocks** of encoders and decoders
- **h = 8** attention heads

BERT: Pre-training

BERT is pre-trained from two unsupervised tasks:

Masked Language Modeling (MLM)

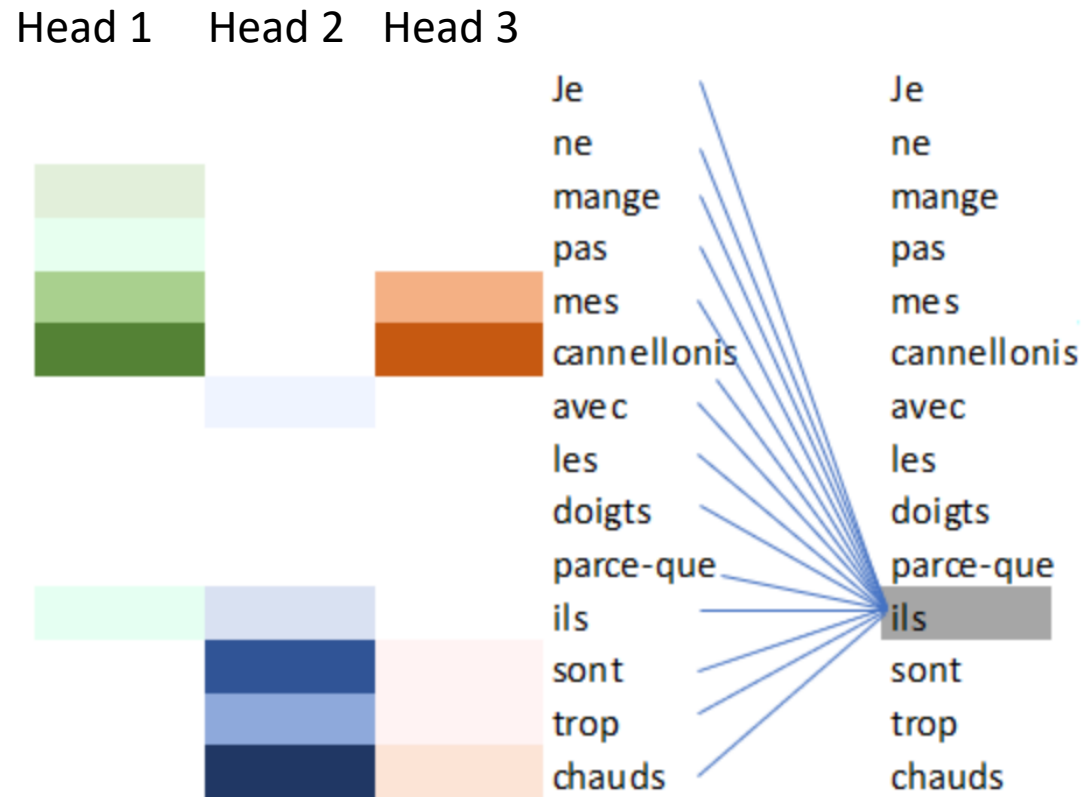
- The model is pre-trained from sequences decomposed into *tokens*
- The goal is to predict some of these tokens (15%) randomly masked in the text

Next Sentence Prediction (NSP)

- The goal is to understand the relationship between two successive sentences
- To do that, the algorithm is trained in order to predict the sentence following a given sentence
- The training set is composed of sentence tuples (A,B)
 - In 50% of the cases, B is the **sentence following** the sentence A (the tuple is then labelled *IsNext*)
 - In 50% of the cases, B is a sentence **randomly** chosen in the corpus (*NotNext* label)

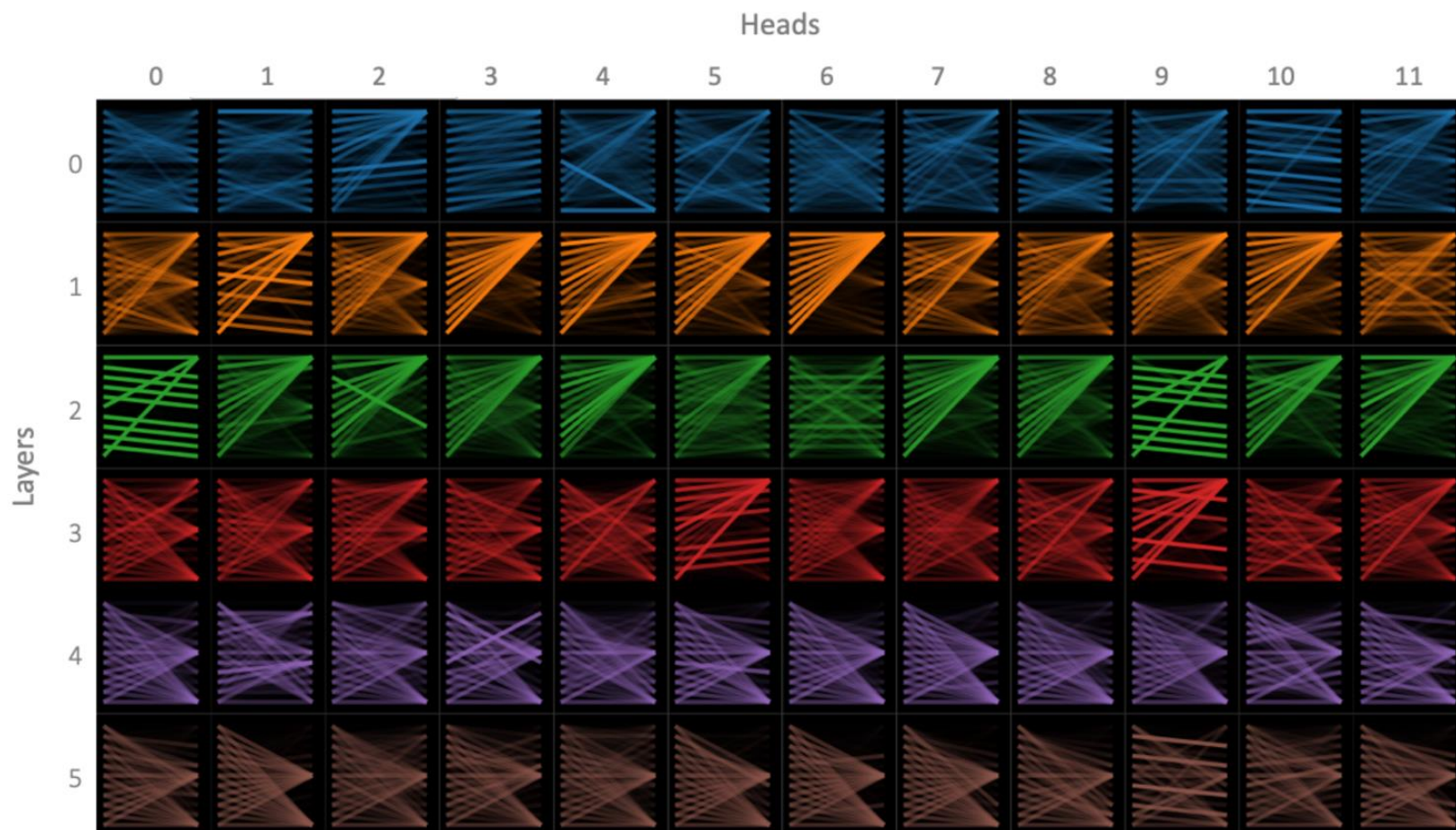
CamemBERT and Attention mechanism

Attention level, for each head, relatively to the french word “ils” (« they »). Each head “see” **different** things from others (**complementarity**)



CamemBERT and Attention mechanism

On the first 6 layers (out of 12), attention mechanism visualization example over the 12 heads



BERT: Fine-tuning

The initial pre-training step needs:

- A **huge** corpus of **generic** texts
- Very **high** computational resources

The fine-tuning step relies on the already pre-trained algorithm, so it only needs :

- A **reduce** and **specialized** corpus
- Significantly **reduces** resources

The aim is obviously to move from a **generic language** modelization to a **task specific one**

Hugging Face



- Hugging Face is an open-source provider of NLP technologies for Python
- A Hugging Face **Transformer** package is available and quite popular
- It gives access to pretrained models such as BERT but also to modules to perform fine-tuning
- It previously supported only **Pytorch**, but now **Tensorflow** is more and more supported as well

BERT implementation: Tutorial notebook

course4_transformer_model_libraries.ipynb

Goal: Illustration of how pre-trained Transformer can be loaded and used from open-source libraries such as Hugging Face

Remarks:

- This notebook comes from a **tutorial notebook** freely accessible (the original reference is given at the very beginning, in the first cell)

Take-away from Course 4

- **LSTM** architectures are **historically important** but there **are not** the state-of-the-art **any more**
- The **“Attention” mechanism** allowed a **significant performance improvement** than using LSTM only
- It has been demonstrated in the seminal paper ***Attention is all you need (2017)***, thanks to the **Transformer architecture** introduced in it, that “Attention” **alone** is enough to reach the state-of-the-art (no LSTM needed and no need to process a sequence in a specific order)
- **BERT** is an example of such **pre-trained Transformer models**, and it can be **finetuned** for specific tasks
- **Hugging Face** offers quite convenient tools **for using and fine-tuning BERT**
- The best models available (like **GPT-3**) can achieve very impressive results on a lot of tasks, but **are not open source**
- However, some **fine-tuning can be done** and sometimes quite **easily** (see **GPT-2**) for interesting results

References

Book

A. Géron, *Hands-On Machine Learning with Scikit-Learn, Keras, and Tensorflow: Concepts, Tools, and Techniques to Build Intelligent Systems* (2019)

Online formations

- <https://www.udemy.com/course/nlp-natural-language-processing-with-python>
- <https://www.coursera.org/specializations/natural-language-processing>
- <https://www.coursera.org/learn/natural-language-processing-tensorflow>

Internet sites

- <https://arxiv.org/abs/1706.03762> (*Attention is all you need*)
- https://www.youtube.com/watch?reload=9&v=OyFJWRnt_AY
- <http://jalammar.github.io/illustrated-transformer/>
- <https://medium.com/dissecting-bert/dissecting-bert-part-1-d3c3d495cdb3>
- <https://camembert-model.fr/>
- <https://huggingface.co/>
- <https://lesdieuxducode.com/blog/2019/4/bert--le-transformer-model-qui-sentraine-et-qui-represente>
- https://www.youtube.com/watch?v=S0zQoTkX_YQ (github copilot)

Online notebooks

- <https://colab.research.google.com/github/tensorflow/text/blob/master/docs/tutorials/transformer.ipynb>
- https://colab.research.google.com/github/sarthakmalik/GPT2.Training.Google.Colaboratory/blob/master/Train_a_GPT_2_Text_Generating_Model_w_GPU.ipynb
- https://colab.research.google.com/github/ziadloo/attention_keras/blob/master/examples/colab/LSTM.ipynb