

# Deep Learning for Natural Language Processing (NLP)

Frédéric Armetta

MCF, LIRIS laboratory, University of Lyon1  
frederic.armetta@univ-lyon1.fr

October 8, 2024

# Overview

Introduction

Word  
encoding

RNN

RNN & NLP

Seq2Seq

Transformer

My research

Personal work  
introduction

- 1 Introduction
- 2 Word encoding
- 3 RNN
- 4 RNN & NLP
- 5 Seq2Seq
- 6 Transformer
- 7 My research
- 8 Personal work introduction

# What is NLP ?

*"Natural language processing (NLP) is a subfield of linguistics, computer science, and artificial intelligence concerned with the interactions between computers and human language, in particular how to program computers to process and analyze large amounts of natural language data." (wikipedia'21)*

# Applications of NLP

- Classification  
Sentiment Analysis, Text Classifications, Spam Detection, etc.

# Applications of NLP

- **Classification**  
Sentiment Analysis, Text Classifications, Spam Detection, etc.
- **Generation**  
Predictive Typing, Question Answering, Text Summarization, Chatbot, etc.

# Applications of NLP

- **Classification**  
Sentiment Analysis, Text Classifications, Spam Detection, etc.
- **Generation**  
Predictive Typing, Question Answering, Text Summarization, Chatbot, etc.
- **Tagging**  
Part-of-speech tagging, Spell Checking, etc.

# Applications of NLP

- **Classification**  
Sentiment Analysis, Text Classifications, Spam Detection, etc.
- **Generation**  
Predictive Typing, Question Answering, Text Summarization, Chatbot, etc.
- **Tagging**  
Part-of-speech tagging, Spell Checking, etc.
- **... but also, extended applications**  
Speech Recognition (sound processing), Character Recognition (image processing), etc.

# How to address NLP problems ?

- Symbolic AI

rule-based-approaches : grammar, vocabulary, etc.

- + : explicit, determinist
- - : inability to render complexity (semantics is complex !), specific knowledge to provide for each application field



# How to address NLP problems ?

## Introduction

### Word encoding

### RNN

### RNN & NLP

### Seq2Seq

### Transformer

### My research

### Personal work introduction

- Symbolic AI

rule-based-approaches : grammar, vocabulary, etc.

- + : explicit, determinist
- - : inability to render complexity (semantics is complex !), specific knowledge to provide for each application field

- Deep Learning

(the topic for this presentation!)

- + : generalization, self-acquired knowledge
- - : availability of datasets, black boxes, explainability

# How to address NLP problems ?

## Introduction

Word  
encoding

RNN

RNN & NLP

Seq2Seq

Transformer

My research

Personal work  
introduction

- Symbolic AI

rule-based-approaches : grammar, vocabulary, etc.

- + : explicit, determinist
- - : inability to render complexity (semantics is complex !),  
specific knowledge to provide for each application field

- Deep Learning

(the topic for this presentation!)

- + : generalization, self-acquired knowledge
- - : availability of datasets, black boxes, explainability

- ... but also Hybrid Approaches

to make benefit from the two ways as far as possible or not ...

# Datasets and transfert learning

## Introduction

### Word encoding

### RNN

### RNN & NLP

### Seq2Seq

### Transformer

### My research

### Personal work introduction

How to benefit from a general purpose dataset for a specific task ?

- (first) Self-supervised learning
  - Huge datasets to educate the network
    - wikipedia, large corpus
    - big networks, heavy time and energy consumption
- (second) Fine-tuning
  - Specialize to a specific task, mainly educate the last layers of the network
    - load a pre-trained general purpose network for the first layers
    - adapt the network (last layers) to the dedicated task
    - modest size dataset
    - restricted learning time
    - greatly enhanced accuracy !

# How to feed neural networks with texts ?

*"If the network can't efficiently differentiate letters or words,  
it's a bad start ..."*

First idea : letters or words in  $[0, 1]$

input      network      output

0.04  $\leftarrow$   $\begin{smallmatrix} \cdots \\ \cdots \\ \cdots \end{smallmatrix}$   $\begin{smallmatrix} \diagup & \diagdown \\ \diagdown & \diagup \end{smallmatrix}$   $\frac{1}{0}$  (consonne)  
("c")

First idea : letters or words in  $[0, 1]$

input      network      output

0.04    $\leftarrow$     $\begin{smallmatrix} \cdots \\ \cdots \\ \cdots \end{smallmatrix}$     $\begin{smallmatrix} \diagup & \diagdown \\ \diagdown & \diagup \end{smallmatrix}$     $\frac{1}{0}$  (consonne)

("c")

$\Rightarrow$  Difficult to separate 26 letters, very difficult to separate 5000 words !

## One-hot encoding (letters)

## Introduction

Word  
encoding

RNN

## RNN & NLP

Seq2Seq

## Transformer

My research

## Personal work introduction

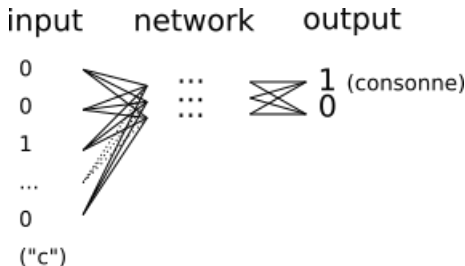
$$V = \{a, b \cdots, z\}, |V| = 26$$

$$a = (1, 0, \dots, 0)$$

$$b = (0, 1, \dots, 0)$$

• • •

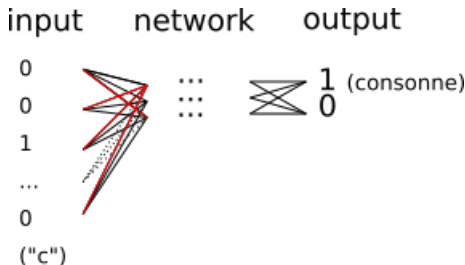
$$z = (0, 0, \dots, 1)$$



## One-hot encoding (letters)

$$V = \{a, b \dots, z\}, \quad |V| = 26$$
$$a = (1, 0, \dots, 0)$$
$$b = (0, 1, \dots, 0)$$

• • •

$$z = (0, 0, \dots, 1)$$


⇒ Easier task for the network !



# One-hot encoding (words)

Introduction

Word  
encoding

RNN

RNN &amp; NLP

Seq2Seq

Transformer

My research

Personal work  
introduction

$$V = \{a, cat \dots, cats\}, |V| = 6$$

$$"a" = (1, 0, \dots, 0)$$

$$"cat" = (0, 1, \dots, 0)$$

$$\dots$$

$$"cats" = (0, 0, \dots, 1)$$

# Word Embedding

- One hot encoding is good for network manipulation but ...
- doesn't hold any semantics !

# Word Embedding

Introduction

Word  
encoding

RNN

RNN &amp; NLP

Seq2Seq

Transformer

My research

Personal work  
introduction

- One hot encoding is good for the network manipulation but ...
- doesn't hold any semantics !

Word embedding ("Plongement lexical")

Transform words into k-dimensional vectors

$$h_{(k,1)} = W_{(k,n)} \cdot X_{(n,1)}$$

$$h_i = \sum_{j=1}^n W_{ij} \cdot x_j$$

# How to learn $W_{(k,n)}$ ?

*"A cat catches a mouse"*

Word2Vect :

- CBOW : predict a words thanks to words in close proximity  
"a", "cat", "?", "a", "mouse" → "catches"
- Skip-gram : predict proximity words thanks to an input word  
"catches" → "a", "cat", "a", "mouse"

## CBOW

Introduction

Word  
encoding

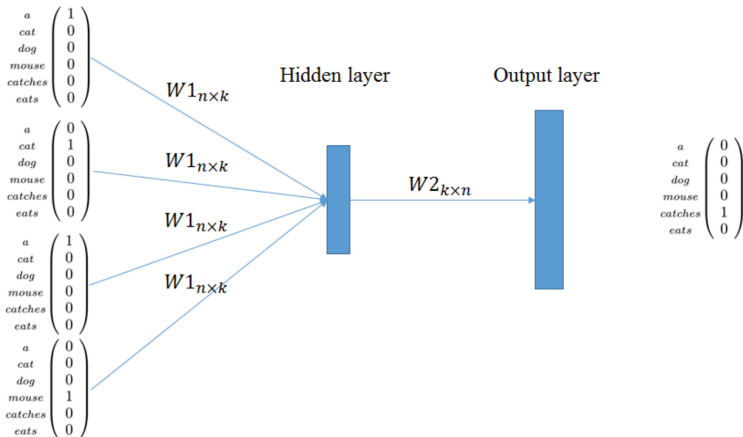
RNN

RNN &amp; NLP

Seq2Seq

Transformer

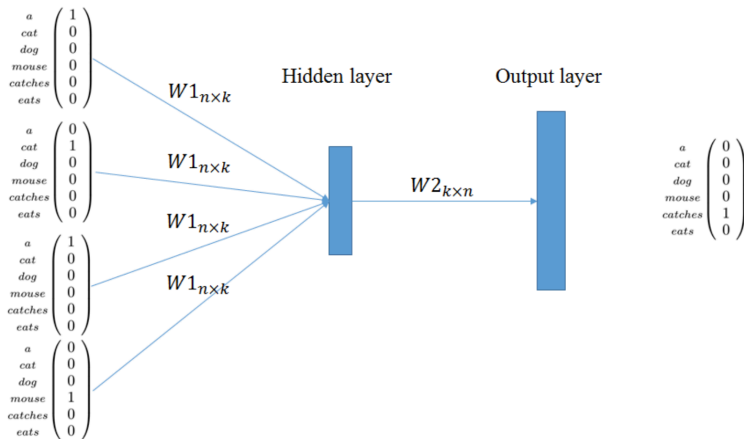
My research

Personal work  
introductionword2vec: CBOW model<sup>1</sup>

<sup>1</sup>Symbolic, Distributed and Distributional Representations for Natural Language Processing in the Era of Deep Learning: a Survey, L. Ferrone and F. M. Zanzotto, CoRR, 2017

# Skip-gram

general idea : reverse order



$outputs \leftarrow input$   
(predict hidden layers for each word of the context)

- Similar to word2vect but ...
- additional global statistics inclusion to obtain word vectors (word2vect incorporates only local statistics)
- open-source project at Stanford<sup>2</sup>

---

<sup>2</sup>GloVe: Global Vectors for Word Representation, J. Pennington and al., Proceedings of EMNLP, 2014, <https://aclanthology.org/D14-1162.pdf>

co-occurrence matrix for the sentence “a cat catches a mouse”  
with a window size of 1:

	a	cat	catches	mouse
a	0	1	1	1
cat	1	0	1	1
catches	1	1	0	0
mouse	1	0	0	0



- Semantic distance between two words thanks to a third one
- homomorphism hypothesis
- and more ...

Table 1: Co-occurrence probabilities for target words *ice* and *steam* with selected context words from a 6 billion token corpus. Only in the ratio does noise from non-discriminative words like *water* and *fashion* cancel out, so that large values (much greater than 1) correlate well with properties specific to ice, and small values (much less than 1) correlate well with properties specific of steam.

Probability and Ratio	$k = \text{solid}$	$k = \text{gas}$	$k = \text{water}$	$k = \text{fashion}$
$P(k \text{ice})$	$1.9 \times 10^{-4}$	$6.6 \times 10^{-5}$	$3.0 \times 10^{-3}$	$1.7 \times 10^{-5}$
$P(k \text{steam})$	$2.2 \times 10^{-5}$	$7.8 \times 10^{-4}$	$2.2 \times 10^{-3}$	$1.8 \times 10^{-5}$
$P(k \text{ice})/P(k \text{steam})$	8.9	$8.5 \times 10^{-2}$	1.36	0.96

context of word  $i$ .

We begin with a simple example that showcases how certain aspects of meaning can be extracted directly from co-occurrence probabilities. Consider two words  $i$  and  $j$  that exhibit a particular aspect of interest; for concreteness, suppose we are interested in the concept of thermodynamic phase, for which we might take  $i = \text{ice}$  and  $j = \text{steam}$ . The relationship of these words can be examined by studying the ratio of their co-occurrence prob-

the information present the ratio  $P_{ik}/P_{jk}$  in the word vector space. Since vector spaces are inherently linear structures, the most natural way to do this is with vector differences. With this aim, we can restrict our consideration to those functions  $F$  that depend only on the difference of the two target words, modifying Eqn. (1) to,

$$F(w_i - w_j, \tilde{w}_k) = \frac{P_{ik}}{P_{jk}}. \quad (2)$$

GloVe: Global Vectors for Word Representation, J. Pennington and al., Proceedings of EMNLP, 2014, <https://aclanthology.org/D14-1162.pdf>

# Word Embedding properties

Introduction

Word  
encoding

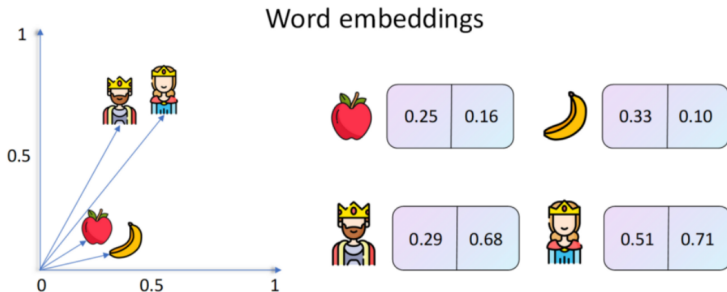
RNN

RNN &amp; NLP

Seq2Seq

Transformer

My research

Personal work  
introduction

source

<https://towardsdatascience.com/deep-learning-for-nlp-word-embeddings-4f5c90bcadab5>


# Word Embedding properties

Introduction

Word  
encoding

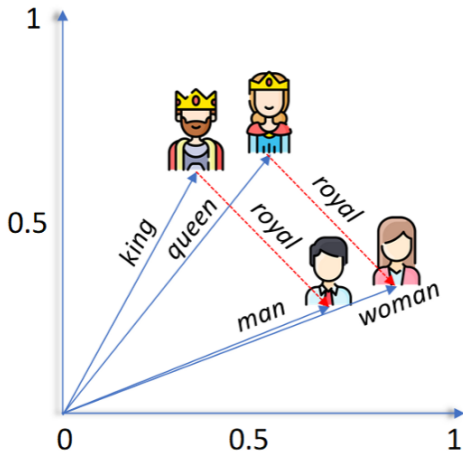
RNN

RNN &amp; NLP

Seq2Seq

Transformer

My research

Personal work  
introduction

source

<https://towardsdatascience.com/deep-learning-for-nlp-word-embeddings-4f5c90bcdab5>

# Why RNN ?

- RNN = Recurrent Neural Network
  - as an encoder: integrate in an internal state (memory) sequences of inputs
  - as a decoder: generate from an internal state (memory) sequences of outputs

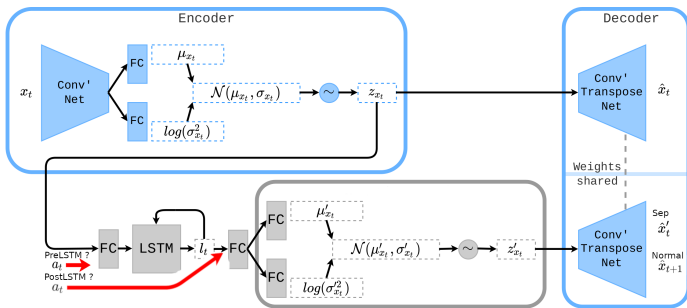
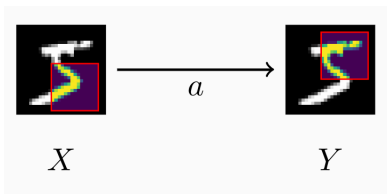


Figure: Accumulate perceptions for learning

# Why RNN and NLP ?

- RNN for NLP
  - as an encoder : build a representation for a sentence/document
  - as an encoder : build an embedding for words depending on the context (see ELMO)
  - as a decoder : text generation (text summarising, machine translation, etc.)

# Why to capture the context ?

Introduction

Word  
encoding

RNN

RNN &amp; NLP

Seq2Seq

Transformer

My research

Personal work  
introduction

- the **driver** is outdated
- the **driver** exceeds the speed limit

different contexts  $\Rightarrow$  different meanings

# Recurrent Neural Network (RNN)

Introduction

Word  
encoding

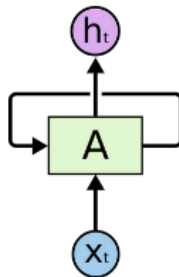
RNN

RNN &amp; NLP

Seq2Seq

Transformer

My research

Personal work  
introduction



# Recurrent Neural Network (RNN)

Introduction

Word  
encoding

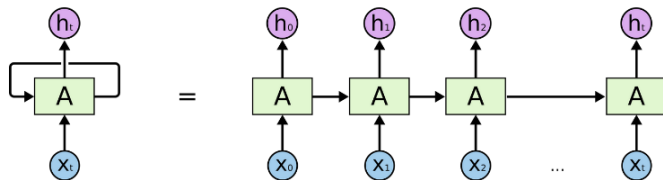
RNN

RNN &amp; NLP

Seq2Seq

Transformer

My research

Personal work  
introduction

An unrolled recurrent neural network.

---

source

<https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

# Recurrent Neural Network (RNN)

Introduction

Word  
encoding

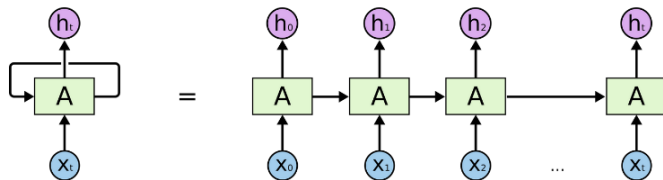
RNN

RNN &amp; NLP

Seq2Seq

Transformer

My research

Personal work  
introduction

An unrolled recurrent neural network.

exemple : text completion

---

source

<https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

# Recurrent Neural Network (RNN)

Introduction

Word  
encoding

RNN

RNN &amp; NLP

Seq2Seq

Transformer

My research

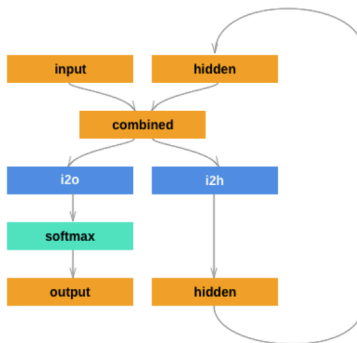
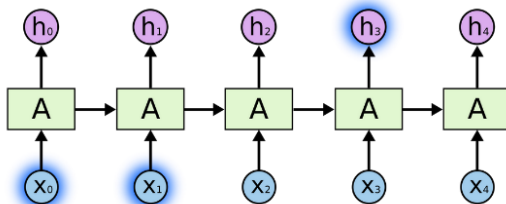
Personal work  
introduction

Figure: A simple RNN : [https://pytorch.org/tutorials/intermediate/char\\_rnn\\_classification\\_tutorial](https://pytorch.org/tutorials/intermediate/char_rnn_classification_tutorial)

# Long-Term Dependencies

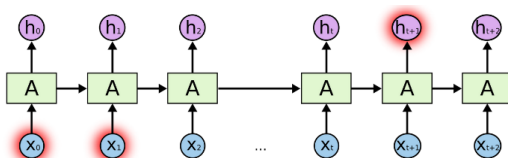


---

source

<https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

# Long-Term Dependencies



---

source

<https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

# Long short-term memory (LSTM)

Introduction

Word  
encoding

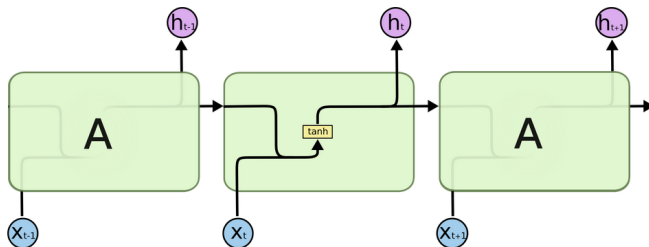
RNN

RNN &amp; NLP

Seq2Seq

Transformer

My research

Personal work  
introduction

# Long short-term memory (LSTM)

Introduction

Word  
encoding

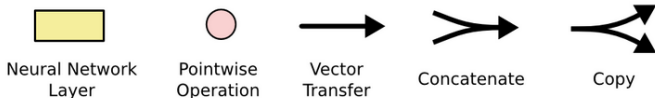
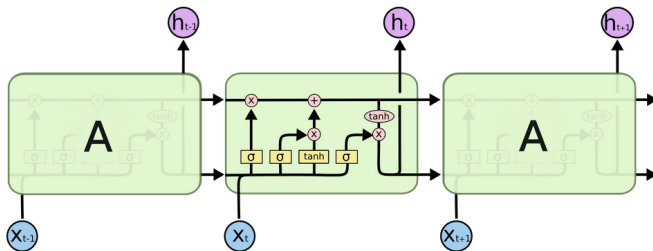
RNN

RNN &amp; NLP

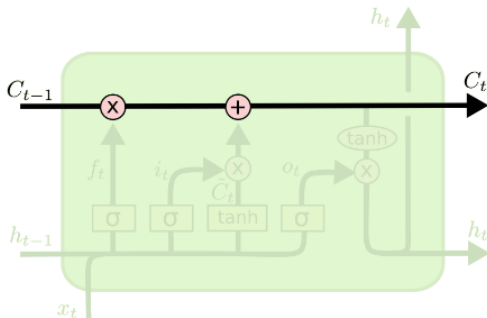
Seq2Seq

Transformer

My research

Personal work  
introduction

# Long short-term memory (LSTM)





# Long short-term memory (LSTM)

Introduction

Word  
encoding

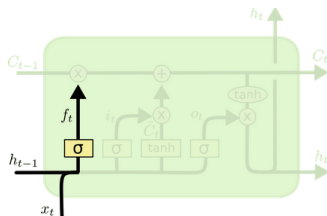
RNN

RNN &amp; NLP

Seq2Seq

Transformer

My research

Personal work  
introduction

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

# Long short-term memory (LSTM)

Introduction

Word  
encoding

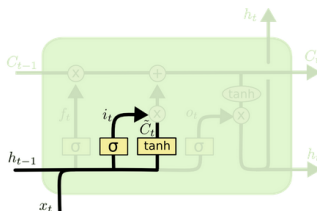
RNN

RNN &amp; NLP

Seq2Seq

Transformer

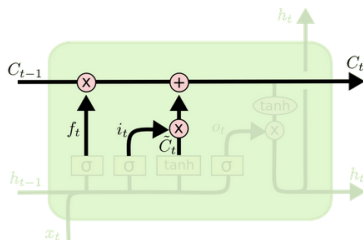
My research

Personal work  
introduction

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

# Long short-term memory (LSTM)



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

# Long short-term memory (LSTM)

Introduction

Word  
encoding

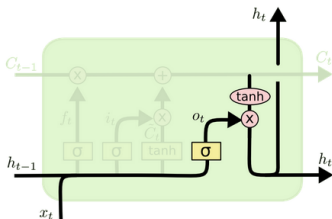
RNN

RNN &amp; NLP

Seq2Seq

Transformer

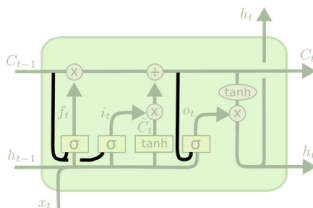
My research

Personal work  
introduction

$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

# Variant : the gate layers look at the cell state

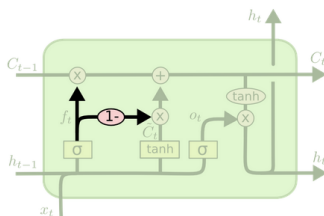


$$f_t = \sigma(W_f \cdot [C_{t-1}, h_{t-1}, x_t] + b_f)$$

$$i_t = \sigma(W_i \cdot [C_{t-1}, h_{t-1}, x_t] + b_i)$$

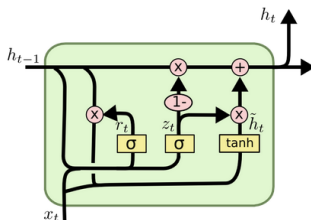
$$o_t = \sigma(W_o \cdot [C_t, h_{t-1}, x_t] + b_o)$$

# Variant : coupled forget and input gates



$$C_t = f_t * C_{t-1} + (1 - f_t) * \tilde{C}_t$$

# Gated Recurrent Unit (GRU)



$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

Figure: cell state and hidden state merged

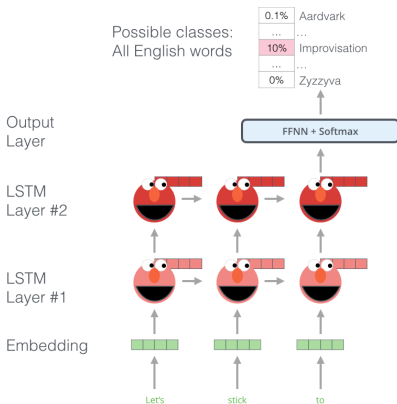
- Embeddings from Language Model
- Embeddings are context-sensitive

---

<sup>2</sup>ME Peters and al., Deep contextualized word representations, 2018,  
<https://arxiv.org/pdf/1802.05365>

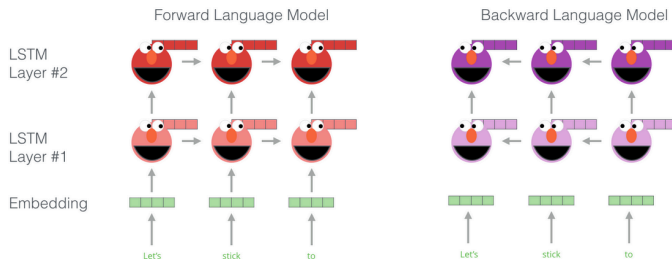


# Stacked LSTMs



A step in the pre-training process of ELMo: Given "Let's stick to" as input, predict the next most likely word – a language modeling task. When trained on a large dataset, the model starts to pick up on language patterns. It's unlikely it'll accurately guess the next word in this example. More realistically, after a word such as "hang", it will assign a higher probability to a word like "out" (to spell "hang out") than to

## BI-LSTM



# Contextualized embedding

Introduction

Word  
encoding

RNN

RNN &amp; NLP

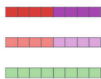
Seq2Seq

Transformer

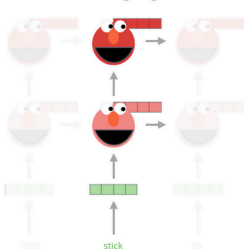
My research

Personal work  
introduction

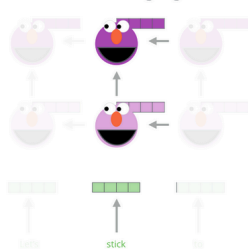
1- Concatenate hidden layers

2- Multiply each vector by  
a weight based on the task3- Sum the (now weighted)  
vectors

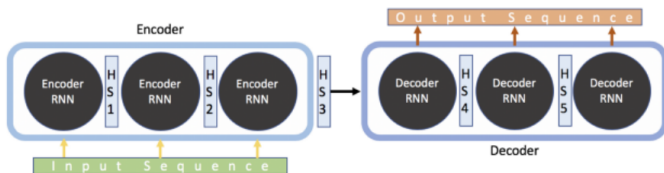
Forward Language Model



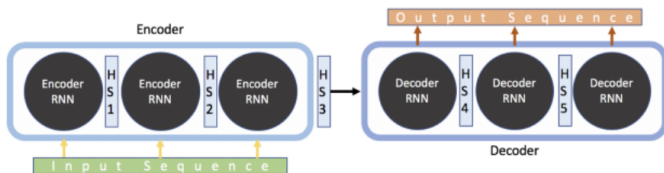
Backward Language Model



# Sequence-2-Sequence Model (Seq2Seq)

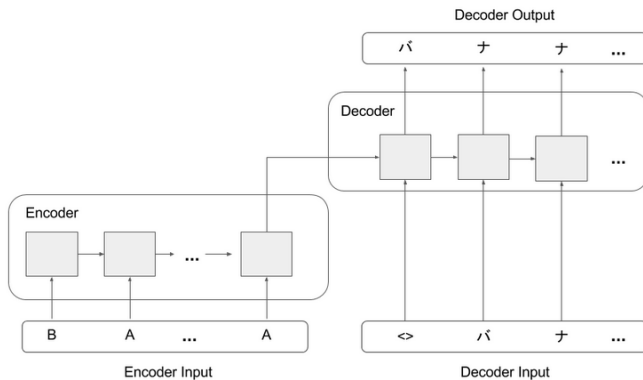


# Sequence-2-Sequence Model (Seq2Seq)



exemple : question answering, text traduction, text summarization

# Teacher forcing



---

source

<https://www.kaggle.com/residentmario/seq-to-seq-rnn-models-attention-teacher-forcing>

# Attention mechanism

Introduction

Word  
encoding

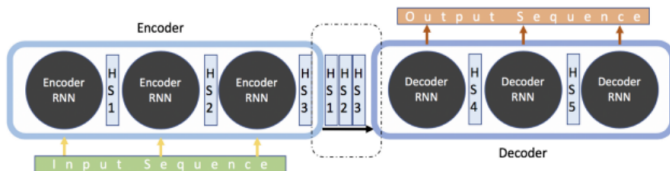
RNN

RNN &amp; NLP

Seq2Seq

Transformer

My research

Personal work  
introduction

# Attention mechanism

Introduction

Word  
encoding

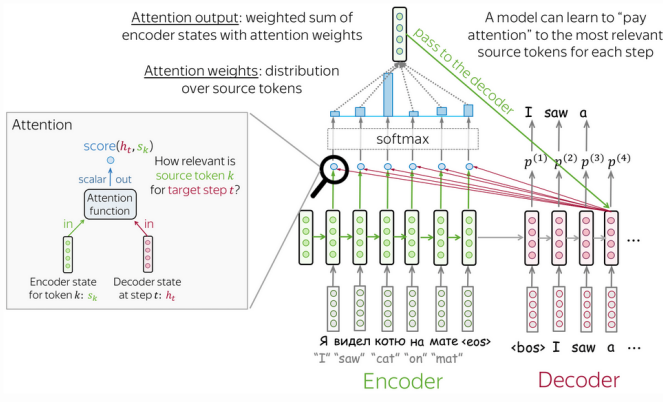
RNN

RNN &amp; NLP

Seq2Seq

Transformer

My research

Personal work  
introduction



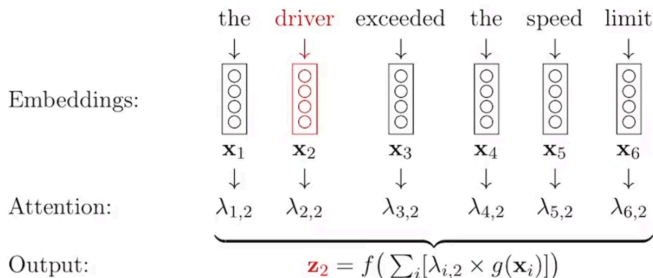
# Attention is all you need<sup>3</sup>

Recurrence no more necessary

<sup>3</sup> Attention Is All You Need, A. Vaswani et al. , NeurIPS proceedings, 2017

# Self-attention

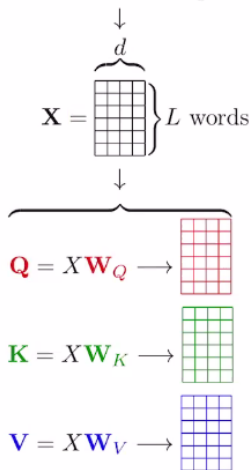
Consider the word **driver**:



- $(\lambda_{i,j})$  are the attention coefficients,  $\sum_i \lambda_{i,j} = 1$ , and
- Reflects the influence of  $\mathbf{x}_i$  on  $\mathbf{x}_j$  (transformed version)

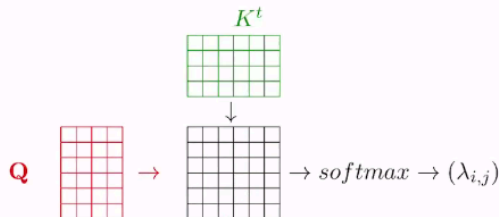
# Queries, Keys, Values

the driver exceeded the speed limit



# Transformer

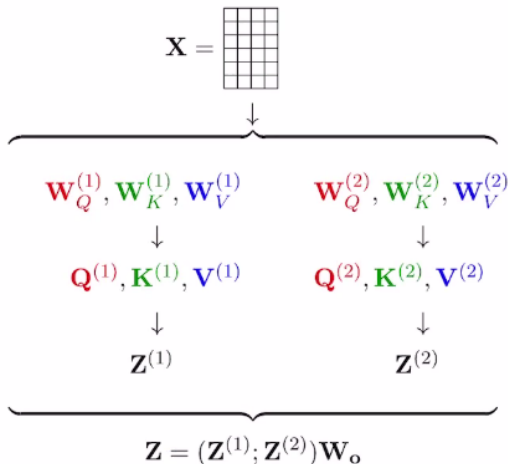
The distance matrix between  $Q$  and  $K$



Scaled Dot-Product Attention

$$\mathbf{Z} = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^t}{\sqrt{d}}\right)\mathbf{V} = \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix}$$

## Multi-head attention



# Transformer

Introduction

Word  
encoding

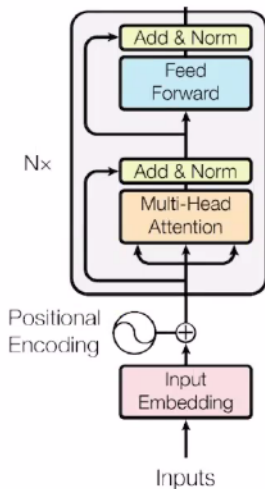
RNN

RNN &amp; NLP

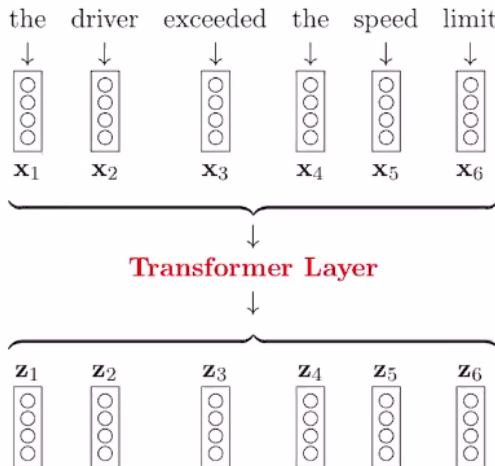
Seq2Seq

**Transformer**

My research

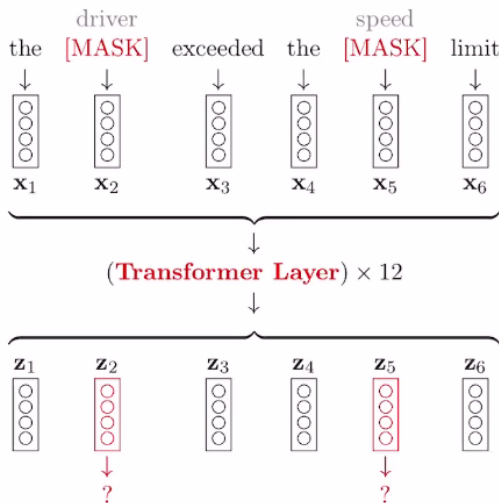
Personal work  
introduction

# Transformer



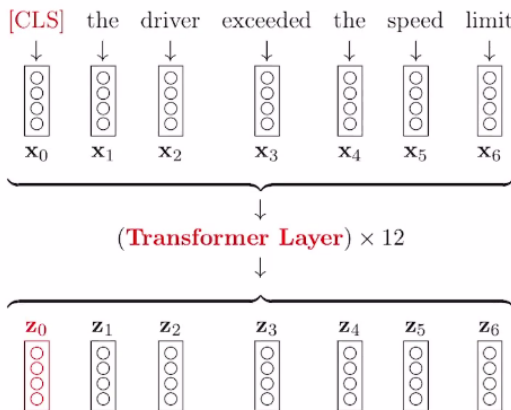
Transformer layers can be stacked !

## Learning





## Fine tuning



NLP & DL

F. Armetta

Introduction

Word  
encoding

RNN

RNN & NLP

Seq2Seq

Transformer

**My research**

Personal work  
introduction

# My research

NLP & DL

F. Armetta

Introduction

Word  
encoding

RNN

RNN & NLP

Seq2Seq

Transformer

My research

**Personal work  
introduction**

# Sentiment classification