

Cutting-edge Deep Learning for NLP learners.

Isabel Segura-Bedmar

1 y 2 Julio, Universidad Politécnica de Madrid (UPM)

Presentación

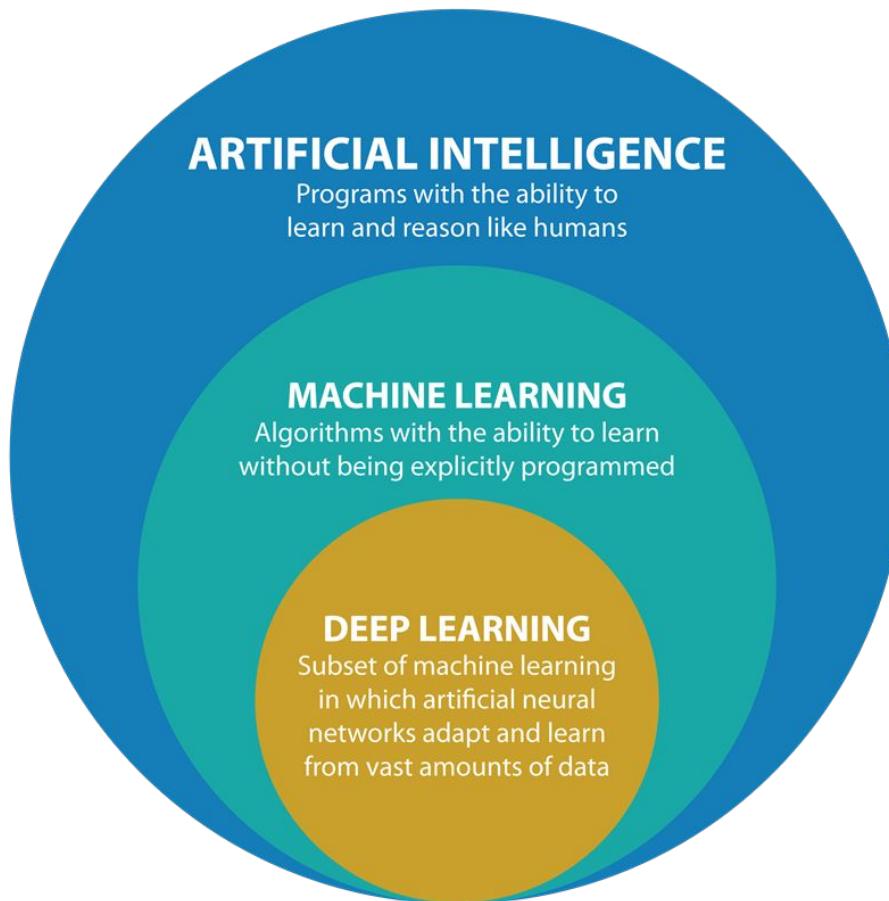
- 1998: Licenciada en Ciencias Matemáticas, Universidad Complutense de Madrid
- 1998-2004: Desarrollo Software en Telefónica I+D y Grupo Santander.
- 2004:... - Profesora Departamento Informática, Universidad Carlos III de Madrid
 - 2010: Doctora Europea en Ciencia y Tecnología Informática, Tesis: "Application of Information Extraction techniques to the pharmacological domain", Premio Extraordinario de Doctorado y Premio Investigación SEPLN.
 - 2018: Profesora Titular.
 - Investigación centrada en Procesamiento de Lenguaje Natural
 - Más de 10 proyectos de investigación competitivos. colP: DeepEMR, NLP4RARE-CM-UC3M
 - Publicaciones: <https://scholar.google.es/citations?user=ZpbtnaUAAA&hl=es>,
 - Organización de competiciones como DDIEXTRACTION 2011 Y 2013.
 - Dirección de tres tesis (dos sobre técnicas de deep learning aplicadas a tareas de PLN).
- Enamorada del PLN y fascinada por las técnicas de machine learning (y en particular, deep learning).
- Mi principal motivación: ayudar a jóvenes investigadores.

Outline (first season)

- **Introduction**
- Word Embeddings
- Deep learning architectures for NLP
 - Recurrent Neural Networks
 - Sequence to Sequence
 - Transformer
- Contextual languages models
 - BERT

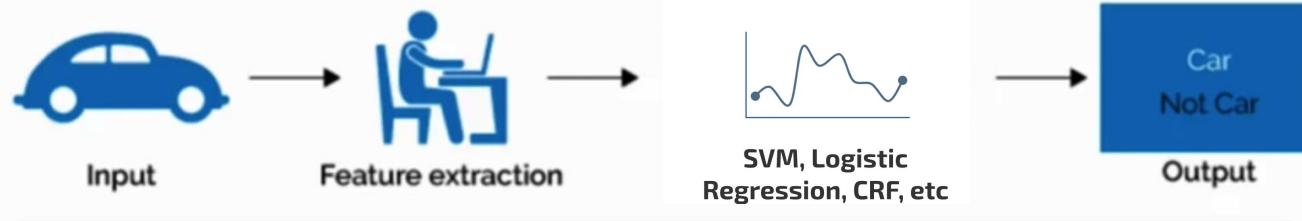
Outline (Second season)

- TASK: Sentiment Analysis of tweets about Covid-19
 - Text representation
 - basic NLP tasks
 - Word Embeddings
 - Evaluation metrics for NLP tasks
 - Machine Learning algorithms
 - SVM, Logistic Regression
 - Deep Learning models:
 - BiLSTM with random initialization
 - BiLSTM with pre-trained word embeddings.
 - BERT

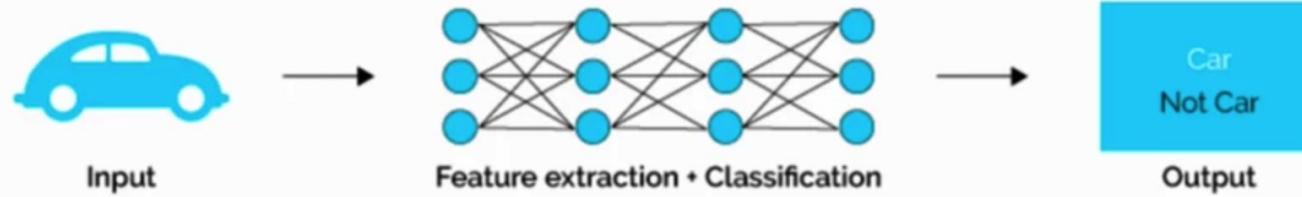


Traditional ML algorithms for NLP tasks heavily require on hand-crafted features. Feature engineering is a time-consuming process. Moreover, it is not usually robust enough (low recall)

Machine Learning

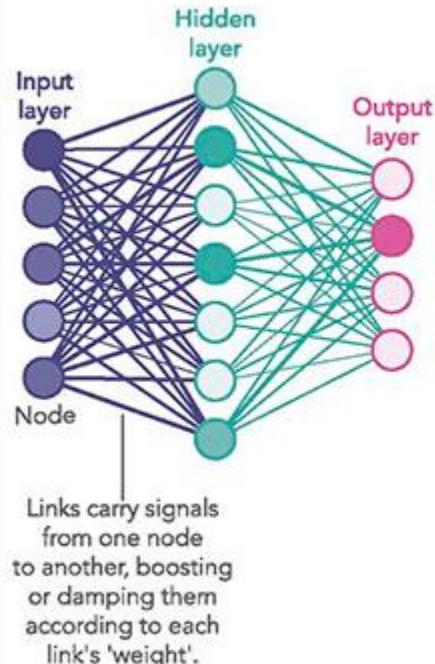


Deep Learning

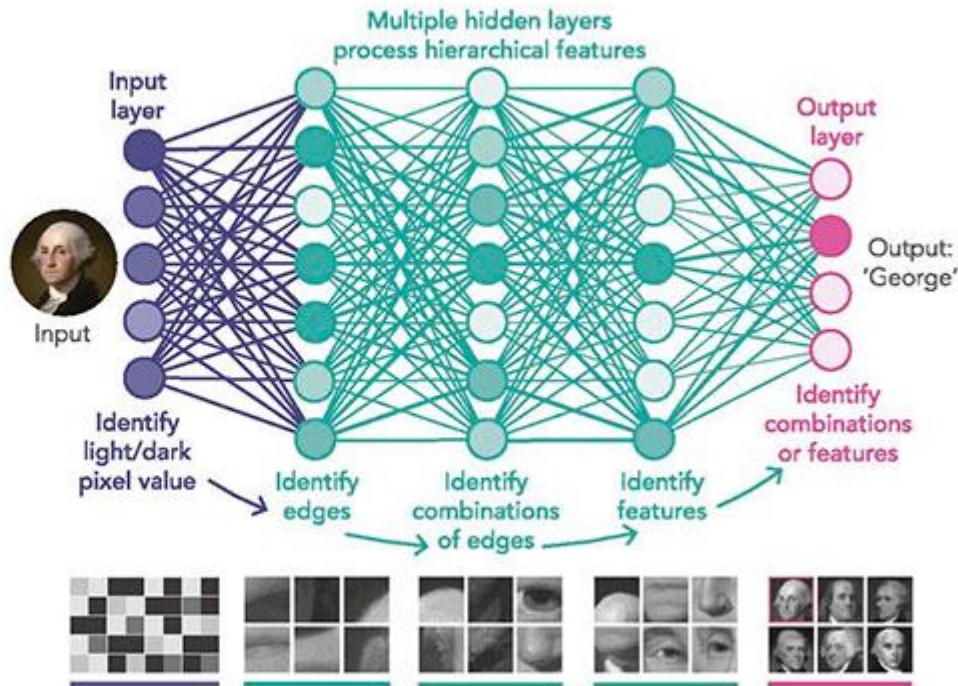


Employ multiple processing layers to learn hierarchical representations of data

1980S-ERA NEURAL NETWORK



DEEP LEARNING NEURAL NETWORK

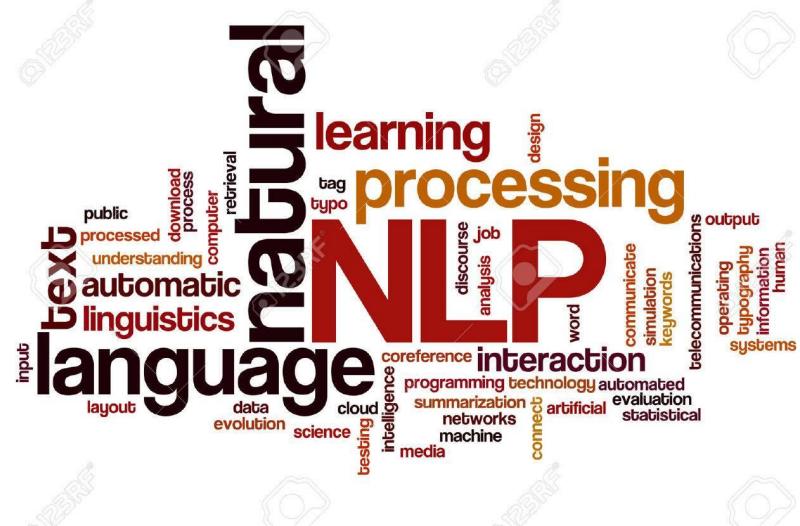


<https://www.futurespace.es/redes-neuronales-y-deep-learning-capitulo-1-preludio/>

Produce state-of-the-art results in many domains
(such as computer vision, pattern recognition)

Natural Language Processing

- Computational techniques for the automatic analysis and representation of human language
- Multidisciplinary: Linguistics, Mathematical and Computer Science, Psychology



Why NLP?

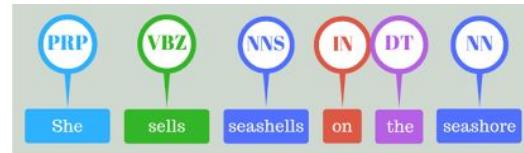


- Exponential Growth of data.
 - 2013, 3.5 ZB
 - 2022, 40 ZB
 - 2025, 180 ZB
 - > 80% in unstructured form,
primarily texts

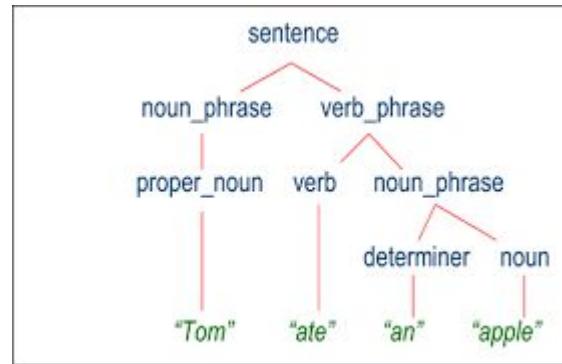
*1 ZB = 1 trillion GB

NLP applications

- NLP applications:
 - part-of-speech (POS) tagging, parsing,
 - language identification
 - Machine translation,
 - Information retrieval,
 - Information extraction,
 - Text classification,
 - Question answering,
 - Text summarization
 - Text simplification,
 - Conversational AI,
 - ...



<https://github.com/dhirajhr/POS-Tagging>



<https://forum.huawei.com>



<https://aylien.com/text-api/language-detection/>

Machine Translation (MT)

The screenshot shows the Google Translate interface. At the top, there's a navigation bar with the Google logo, a grid icon, and a "Sign in" button. Below it, the word "Translate" is displayed in red. There are two language selection bars: one for the source language (German) and one for the target language (English). The source text is in German: "Wir danken all unseren Gästen für die freundlichen Worte und kritischen Anregungen. Jede Anmerkung nehmen wir zum Anlass, unseren eigenen Qualitätsanspruch ständig neu zu prüfen und uns zu verbessern. Wir freuen uns, Sie in unserem Haus begrüßen zu dürfen!" The target text is in English: "We thank all our guests for the kind words and critical suggestions. We take every note as an occasion to constantly revise our own quality requirements and improve ourselves. We are looking forward to welcome you in our house!" Below the text are several small icons for sharing and editing.

Google Sign in

Translate Turn off instant translation

English Spanish French German - detected

English Spanish Arabic

Translate

Wir danken all unseren Gästen für die freundlichen Worte und kritischen Anregungen. Jede Anmerkung nehmen wir zum Anlass, unseren eigenen Qualitätsanspruch ständig neu zu prüfen und uns zu verbessern. Wir freuen uns, Sie in unserem Haus begrüßen zu dürfen!

We thank all our guests for the kind words and critical suggestions. We take every note as an occasion to constantly revise our own quality requirements and improve ourselves. We are looking forward to welcome you in our house!

Information Retrieval

The screenshot shows a Google search results page with the query "BERT". The search bar has "BERT" typed into it. Below the search bar, there are tabs for "Todo", "Imágenes", "Videos", "Noticias", "Maps", "Más", "Configuración", and "Herramientas". The main search results area displays approximately 424,000,000 results. The first result is a link to <https://www.inboundcycle.com> titled "Guía avanzada de Google BERT: qué es, cómo funciona y en ...". The second result is a link to <https://www.bbc.com> titled "Google: cómo funciona BERT, la mayor actualización del ...". The third result is a link to <https://www.arsys.es> titled "Así funciona Bert, el nuevo algoritmo de Google - Blog de ...". The fourth result is a link to <https://conectasoftware.com> titled "Nuevo algoritmo Google BERT: ¿Qué es? ¿Cómo ayuda a ...". On the right side of the search results, there is a sidebar with the title "BERT" and a sub-section "Modelo de lenguaje". It includes a snippet of text about BERT being a Bidirectional Encoder Representations from Transformers model created by Google in 2018. Below this, there is a section titled "También se buscó" with links to "OpenAI GPT", "Word2vec", "PyTorch", and "fastText".

Information Retrieval

The screenshot shows the PubMed search interface with the query "covid". The results page displays 147,246 articles. A prominent callout box encourages the use of COVID-19 filters from PubMed Clinical Queries. The results are sorted by Best match. The interface includes filters for year (1993 to 2022), text availability (Abstract, Free full text, Full text), article attributes (Associated data), and article types (Books and Documents, Clinical Trial, Meta-Analysis, Randomized Controlled Trial, Review). Specific search results are listed, such as "Laboratory testing for the diagnosis of COVID-19" by Lai CKC, Lam W., and "Approaches towards fighting the COVID-19 pandemic (Review)" by Tsai SC, Lu CC, Bau DT, et al.

NIH National Library of Medicine
National Center for Biotechnology Information

Log in

PubMed.gov

covid

Advanced Create alert Create RSS User Guide

Save Email Send to Sorted by: Best match Display options

MY NCBI FILTERS

RESULTS BY YEAR

1993 2022

TEXT AVAILABILITY

Abstract Free full text Full text

ARTICLE ATTRIBUTE

Associated data

ARTICLE TYPE

Books and Documents Clinical Trial Meta-Analysis Randomized Controlled Trial Review

147,246 results

Use COVID-19 filters from PubMed Clinical Queries to refine your search

Treatment Mechanism Transmission More filters

See more SARS-CoV-2 literature, sequence, and clinical content from NCBI

Laboratory testing for the diagnosis of COVID-19.

1 Lai CKC, Lam W. Biochem Biophys Res Commun. 2021 Jan 29;538:226-230. doi: 10.1016/j.bbrc.2020.10.069. Epub 2020 Oct 28. PMID: 33139015 Free PMC article. Review.

Rapid and accurate laboratory diagnosis of active COVID-19 infection is one of the cornerstones of pandemic control. ...

Approaches towards fighting the COVID-19 pandemic (Review).

2 Tsai SC, Lu CC, Bau DT, Chiu YJ, Yen YT, Hsu YM, Fu CW, Kuo SC, Lo YS, Chiu HY, Juan YN, Tsai FJ, Yang JS. Int J Mol Med. 2021 Jan;47(1):3-22. doi: 10.3892/ijmm.2020.4794. Epub 2020 Nov 20. PMID: 33236131 Free PMC article. Review.

The causative agent of COVID-19, SARS-CoV-2, is a novel coronavirus strain. To date, remdesivir has been granted emergency use authorization for use in the management of infection. Additionally, several efficient diagnostic tools are being activ ...

COVID-19: Therapeutics and interventions currently under consideration.

3 McFee RB. Dis Mon. 2020 Sep;66(9):101058. doi: 10.1016/j.dismonth.2020.101058. Epub 2020 Jul 28. PMID: 32833222 Free PMC article. Review.

Information Extraction

Text in



Data out

THE COUNTRIES WITH THE LARGEST POPULATION

China	1	1,388,232,693
India	2	1,342,512,706
United States	3	326,474,013
Indonesia	4	263,510,146
Brasil	5	174,315,386

THE COUNTRY'S FIRST LADIES

- Brigitte Macron
- Spouse: Emmanuel Macron, President of France (2017 -)
- Melania Trump
- Spouse: Donald J. Trump, U.S. President (2017 -)
- Iriana Widodo
- Spouse: Joko Widodo, President of Indonesia (2014 -)
- Also known as: "Ibu Negara" (Lady/Mother of the State)

IMDB TOP RATED TV SHOWS

- 1 Planet Earth II (2016) 9.6.
- 2 Band of Brothers (2001) 9.5.
- 3 Planet Earth (2006) 9.5.
- 4 Game of Thrones (2011) 9.4.
- 5 Breaking Bad (2008) 9.4.

IE tasks

1) Named Entity Recognition (NER)

John lives in London . He works there for Polar Bear Design .

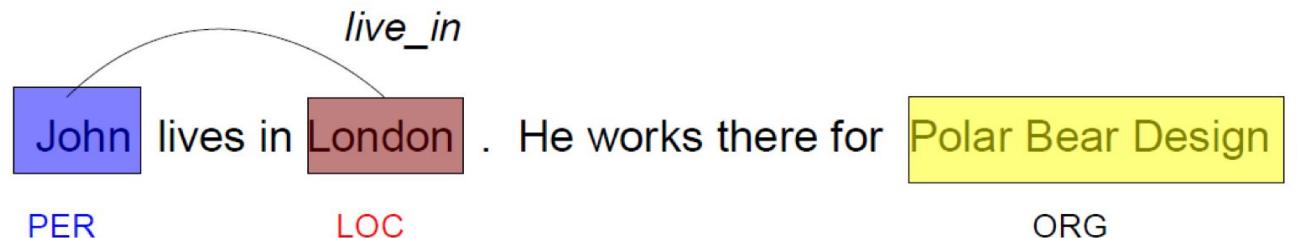
PER LOC ORG

2) Relation Extraction

John lives in London . He works there for Polar Bear Design .

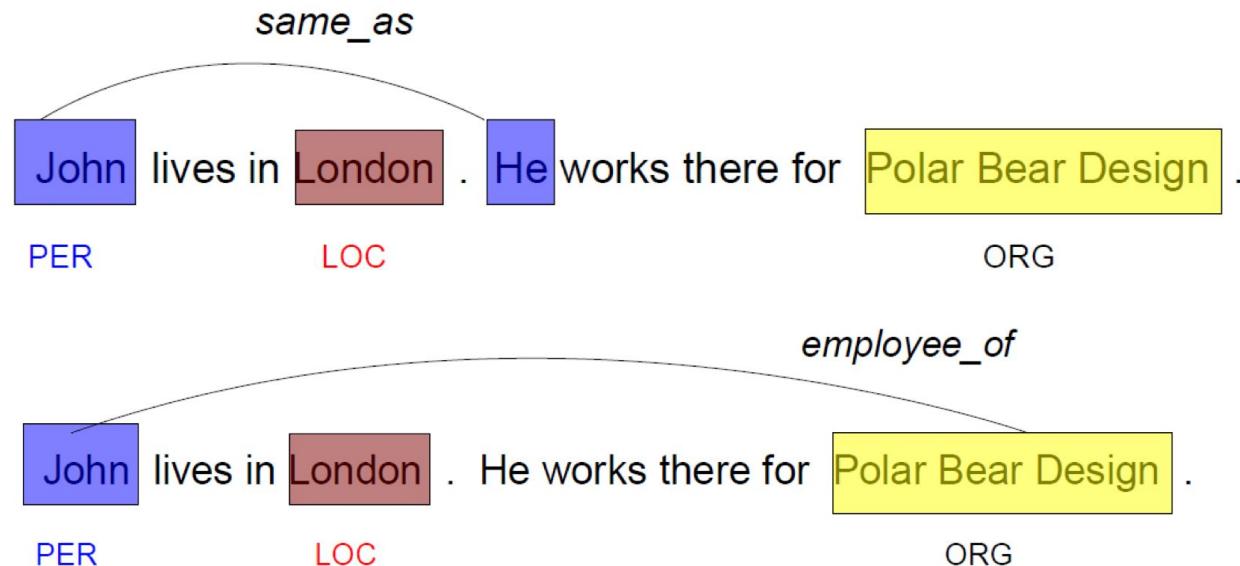
PER LOC ORG

live_in



IE tasks

3) Co-reference resolution



Question Answering

Google

cuál es el edificio más alto del mundo

Todo Imágenes Noticias Vídeos Maps Más Configuración Herramientas

Mundo / Edificios (más altos)

Burj Khalifa 828 m	Torre de Shangái 632 m	Makkah Clock Tower 601 m	Ping An Finance Cen... 555 m	Lotte World Tower 555 m	One World Trade Center 541,3 m
-----------------------	---------------------------	-----------------------------	---------------------------------	----------------------------	-----------------------------------

[¿Cuál es el edificio más alto del mundo? | Plataforma ...](#)

<https://www.plataformaarquitectura.cl> › ArchDaily › Artículos

23 ene. 2019 - En la actualidad existen instituciones especializadas que establecen los parámetros para definir objetivamente cuánto mide un edificio.

Question Answering

Quora

Follow Tech Lounge if you like this content.

Answer · Follow · Request · View all >

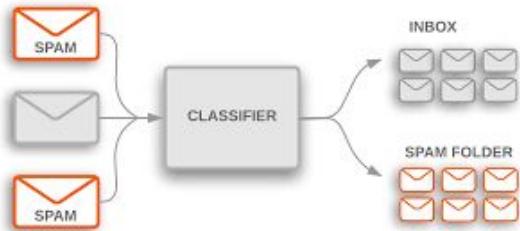
Bhushan · Follow
I speak two languages - English and Sarcasm.. · 12h ago

Pradip Bedre · Answered June 27, 2020
What is neno technology?
How it Started The ideas and concepts behind nanoscience and nanotechnology started with a talk entitled "There's Plenty of Room at the Bottom" by physicist Richard Feynman at an American Physical Society meeting at the California Inst (more)



1 | [Upvote](#) [Downvote](#) [Comment](#) [Share](#)

Text Classification



Text summarization

Source Text: Peter and Elizabeth took a taxi to attend the night party in the city.

While in the party, Elizabeth collapsed and was rushed to the hospital.

Summary: Elizabeth was hospitalized after attending a party with Peter.

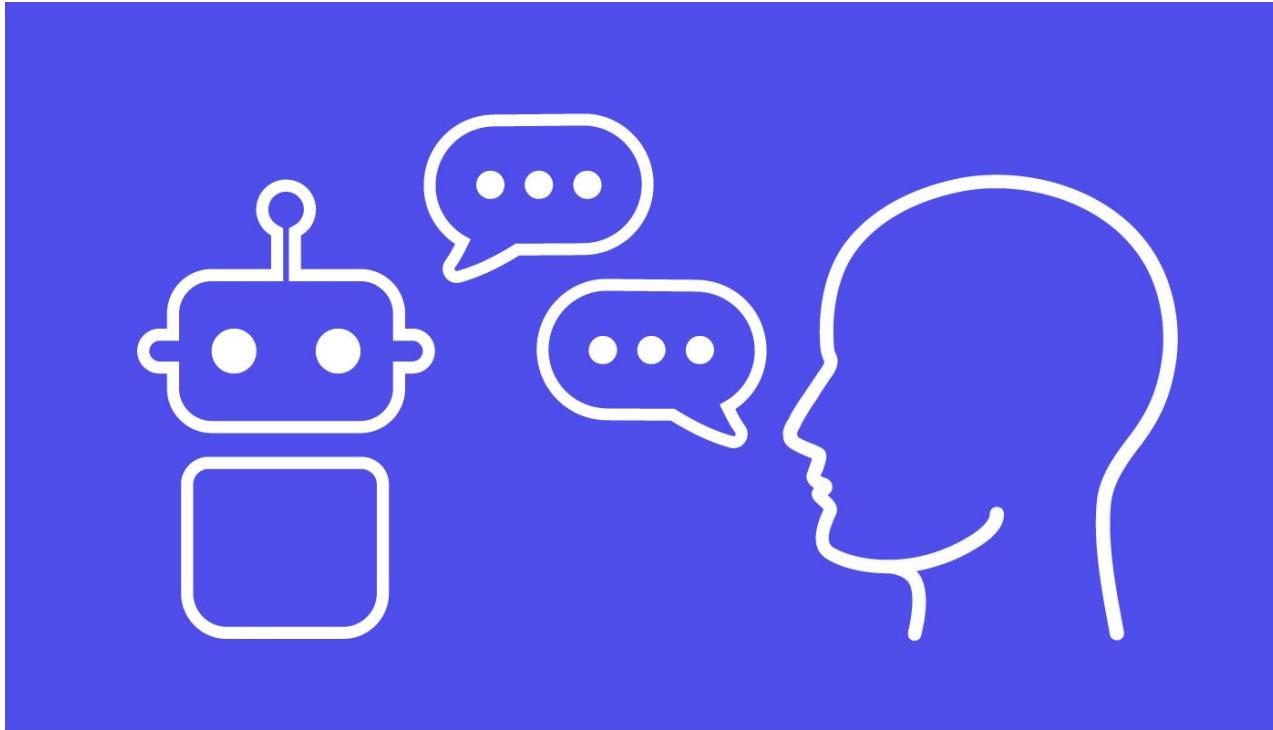


Text simplification

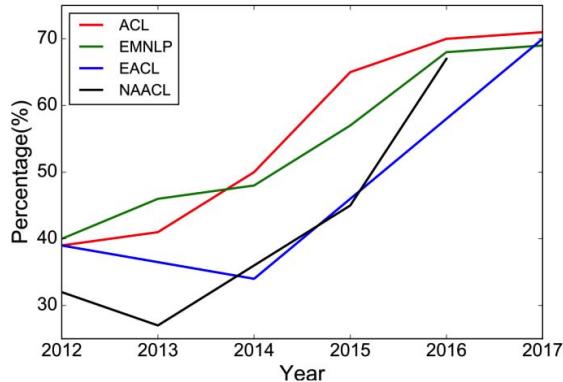
Normal: *Alfonso Perez Munoz, usually referred to as Alfonso, is a former Spanish footballer, in the striker position.*

Simple: *Alfonso Perez is a former Spanish football player.*

Conversational AI (Chatbots)



DEEP LEARNING for NLP



Deep Learning in NLP

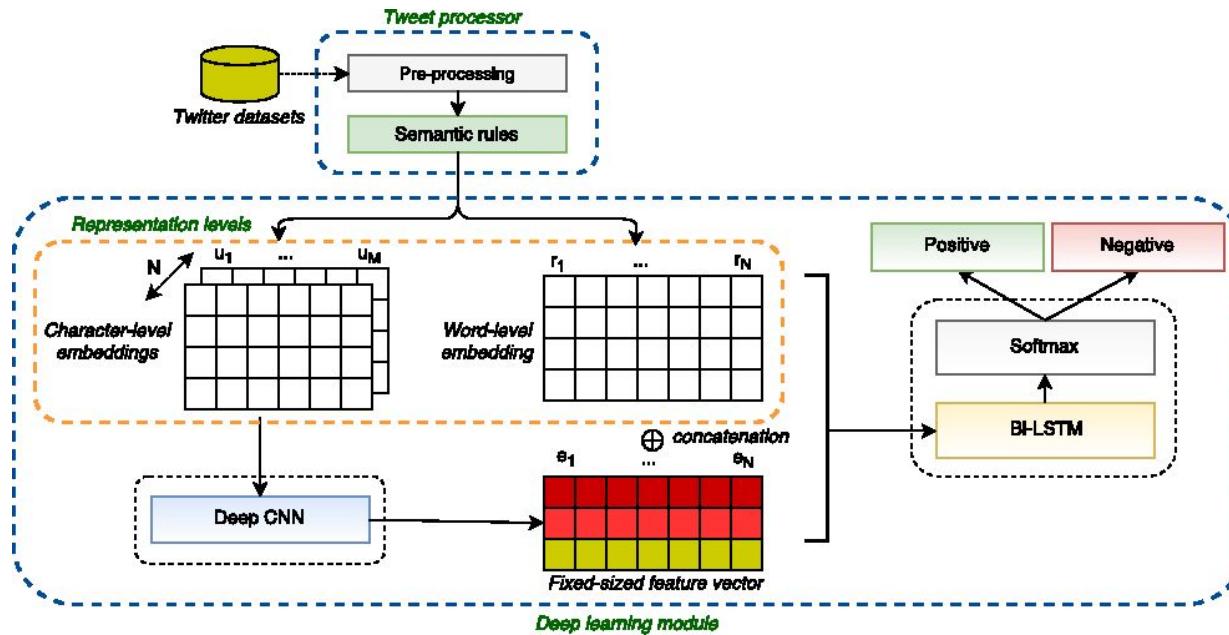
In the last years, NLP research has also exploited deep learning techniques

Produce state-of-the-art results in many NLP applications (machine translation, IE, text summarization, etc)

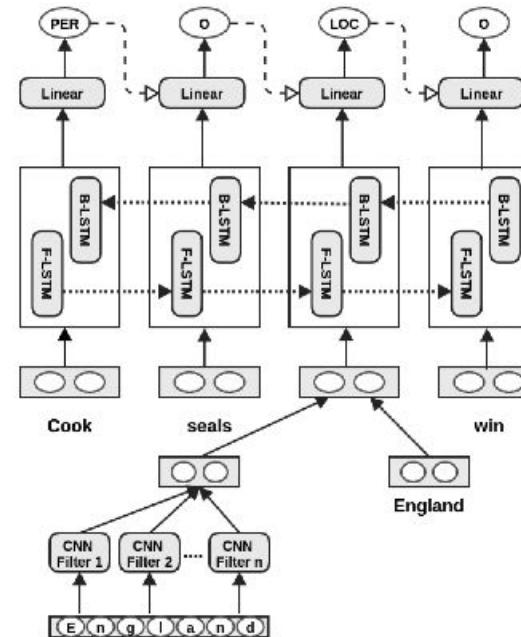
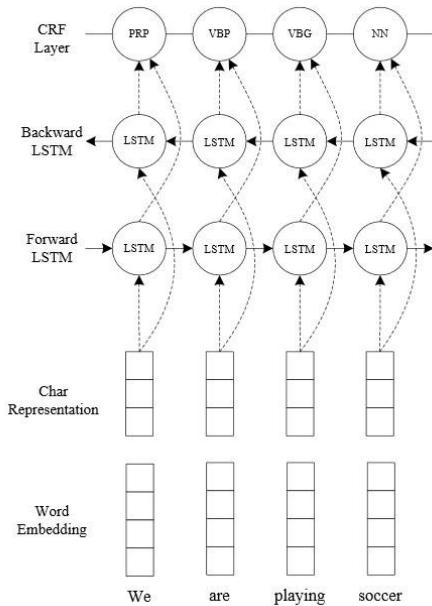
Results



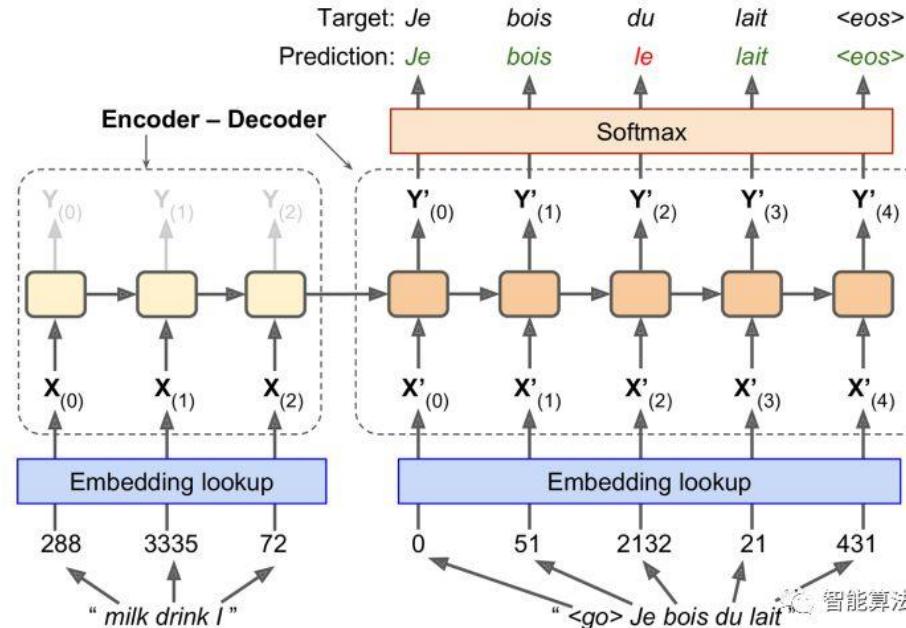
BiLSTM for Text Classification (Sentiment Analysis of Tweets)



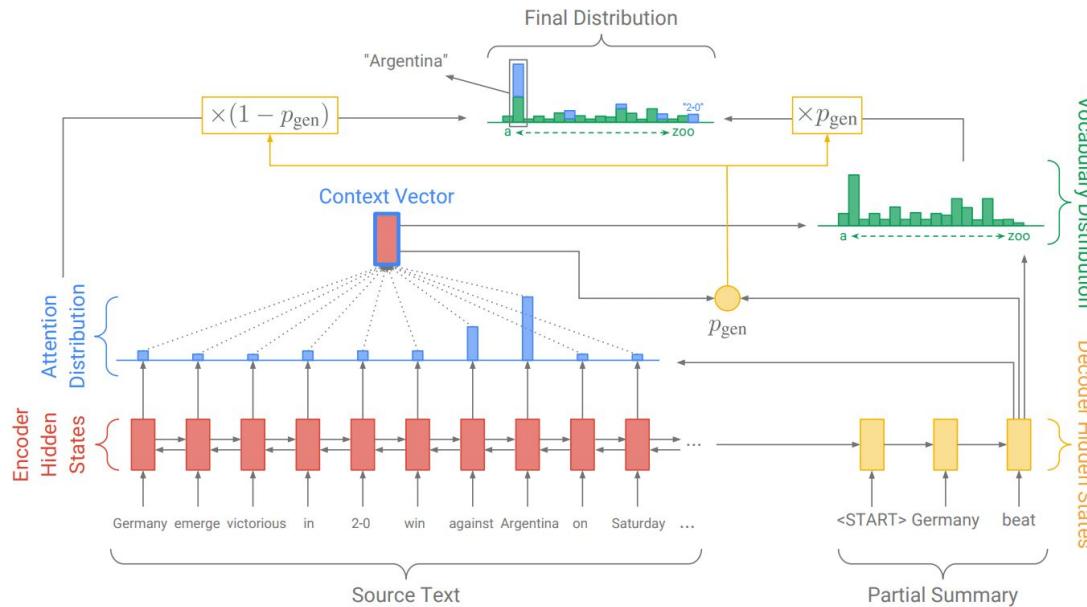
BiLSTM + CRF for Sequence Labeling tasks (PoS tagging and NER)



Seq2Seq for Machine translation



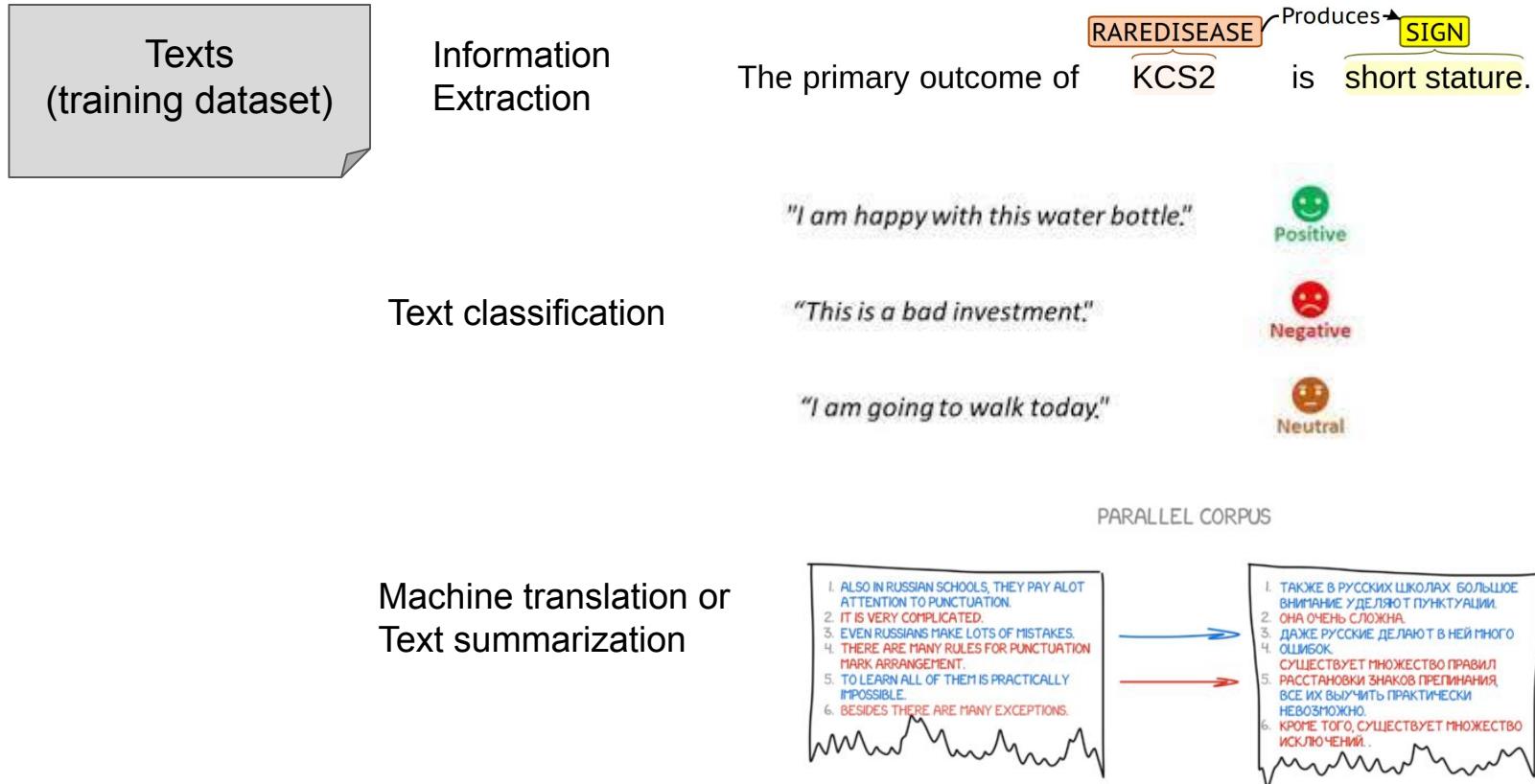
Transformers for Text summarization



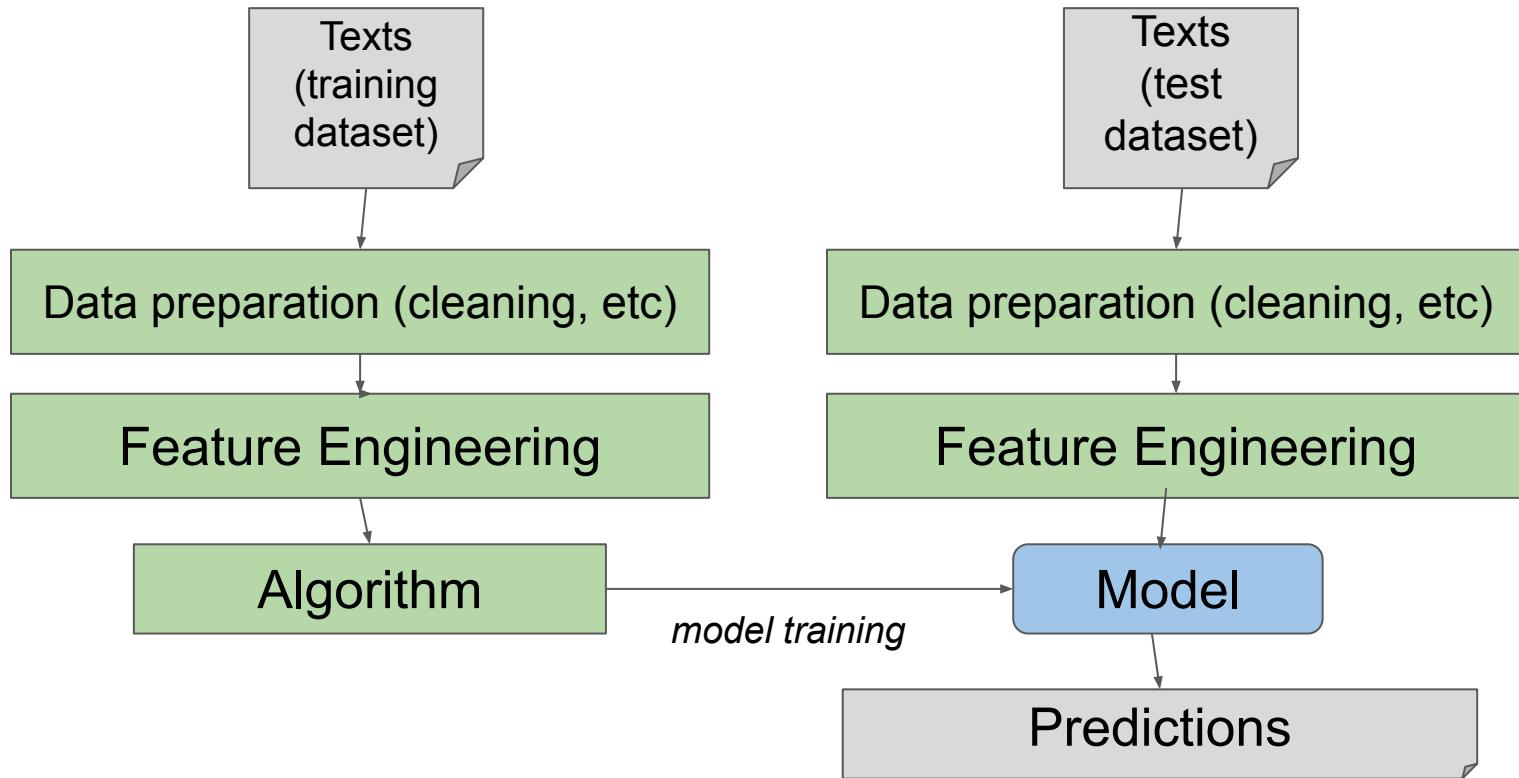
Outline (first season)

- Introduction
- **Word Embeddings**
- Deep learning architectures for NLP
 - Recurrent Neural Networks
 - Sequence to Sequence
 - Transformer
- Contextual languages models
 - BERT

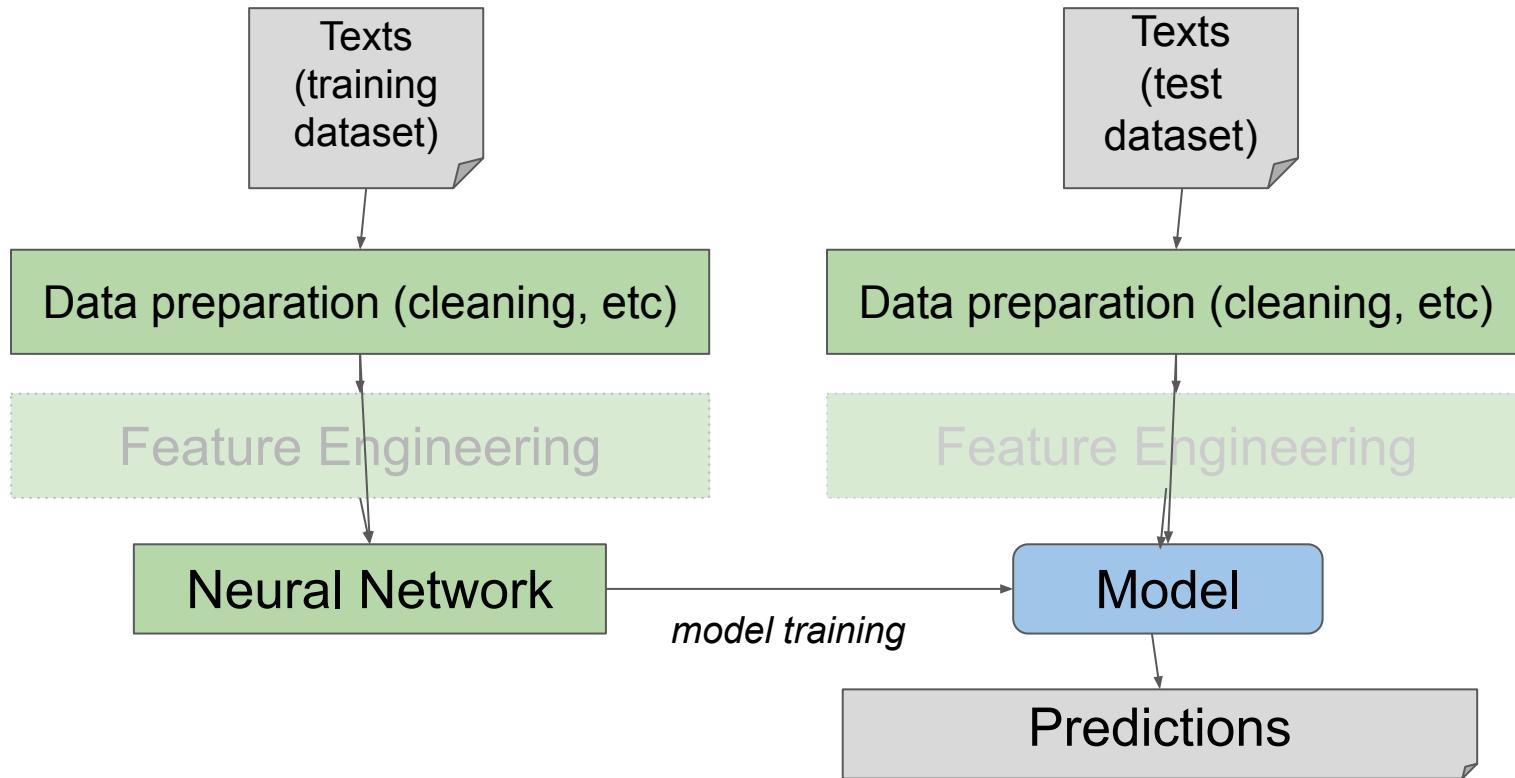
Machine Learning Pipeline



Machine Learning Pipeline



Deep Learning Pipeline



How to initialize a neural network?

- Two approaches:
 - random initialization
 - pre-trained word embeddings (vectors)

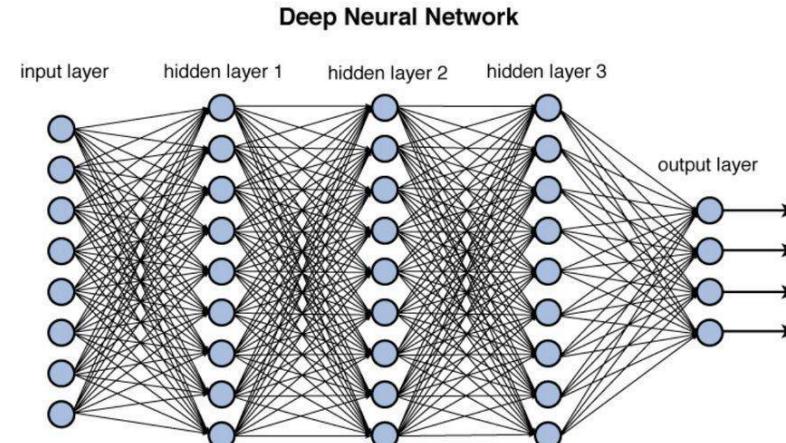
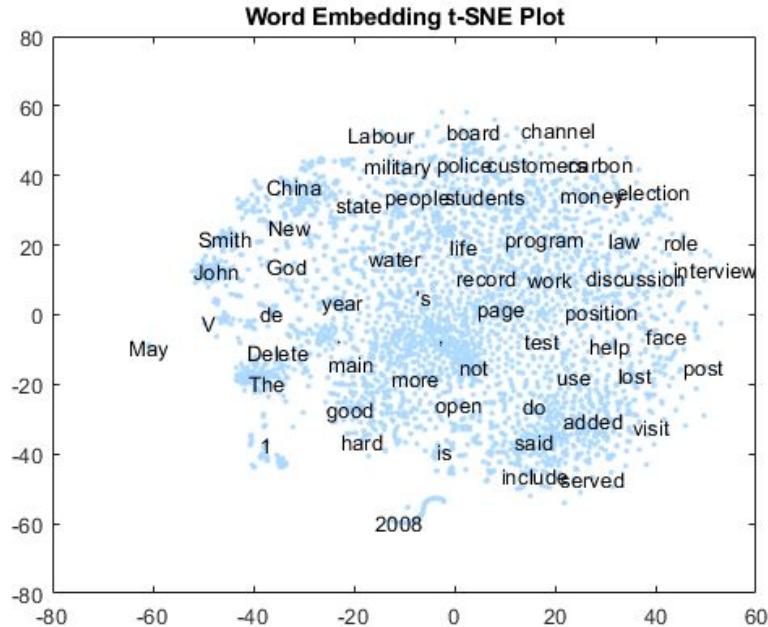


Figure 12.2 Deep network architecture with multiple layers.

<https://towardsdatascience.com/training-deep-neural-networks-9fdb1964b964>

Word embeddings

- NLP techniques where words from the vocabulary are mapped to vectors of real number
- Traditional approaches: BoW and TF-IDF.
- Neural Networks: word2vec, glove, Fast2text



<https://es.mathworks.com/help/textanalytics/ug/visualize-word-embedding-using-text-scatter-plot.html>

Bag of Words

- based on counting words in the document
- Steps:
 - Cleaning:
 - Remove stopwords, punctuation and special symbols.
 - Normalize texts (lemmatization or stemming).

Lemmatization

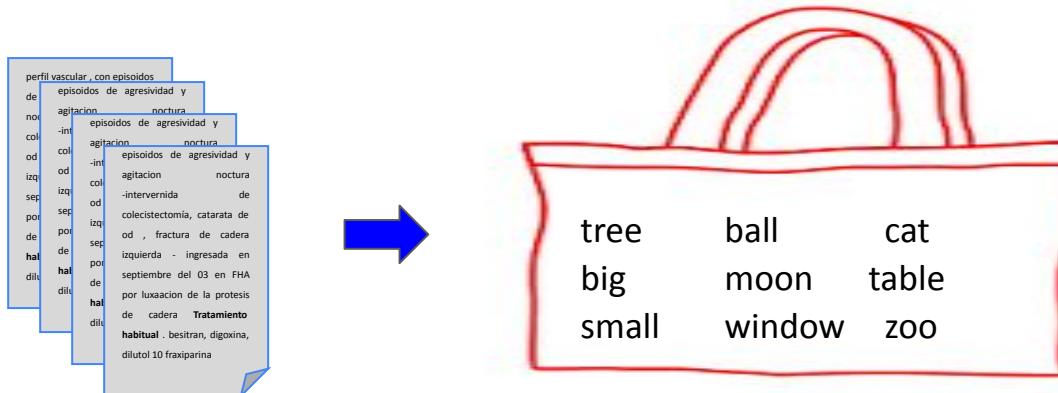
- Obtain the lemma of a word.
 - house -> house
 - *smallest* -> *small*
 - *lawyer, lawsuit* -> *law*
 - *are* -> *be*
 - *women* -> *woman*
- *Decrease the vocabulary size and improve information retrieval*
smaller, smallest -> *small*
- online lemmatizer:
 - <https://cst.dk/tools/index.php#output>
 - <http://www.gedlc.ulpgc.es/investigacion/scogeme02/lematiza.htm>

Stemming

- Obtain the stem of each word
 - *singing* -> *sing*
 - *singer* -> *sing*
 - *house* -> *hous*
 - *lawyer* -> *lawyer*
 - *lawsuit* -> *lawsuit*
 - *are* -> *are*
 - *women* -> *women*
- *Decrease the vocabulary size and improve information retrieval*
 - *poliposis nasal* = *pólipos nasales*
- Stemmer online:
 - <https://snowballstem.org/demo.html>

Bag of Words

- Cleaning
- Obtain vocabulary (unique words) from all texts.



ball	big	cat	moon	small	table	tree	window	zoo

Bag of Words

- Each text is represented as a vector with the frequencies of their words

D: The big cat is on the table and the small cat is in the window.

after cleaning:

D: The big cat is on the table and the small cat is in the window-

Vector (features):

ball	big	cat	moon	small	table	tree	window	zoo
0	1	2	0	1	1	0	1	0

Bag of Words

D1: ~~The big cat is on the table and the small cat in the window~~

D2: ~~The table and the window are small~~

D2: ~~The moon and the small tree are big~~

Vector (features):

	ball	big	cat	moon	small	table	tree	window	zoo
D1	0	1	2	0	1	1	0	1	0
D2	0	0	0	0	1	1	0	1	0
D3	0	1	0	1	1	0	1	0	0

TF-IDF

- Extended version of BoW.
- Every text is represented using tf-idf of its words
- TF-IDF decrease the weight of the very common words in the collection of texts
- Esta métrica, TF-IDF, consigue disminuir el peso de las palabras que son muy comunes en toda la colección de documentos.

TF-IDF

- Term frequency - inverse document frequency.

$$\text{TF-IDF}(W) = \text{TF}(W,d) * \text{IDF}(W)$$

- $\text{TF}(W,d)$ = term frequency of the word W in the document d .
- $\text{IDF}(W)$ = inverse document frequency. The logarithm of the quotient of the total number of documents and the number of documents that contains the word W .

$$\text{idf}(W) = \log \frac{\#\text{(documents)}}{\#\text{(documents containing word } W\text{)}}.$$

TF-IDF

D1: ~~The big cat is on the table and the small cat in the window~~

D2: ~~The table and the window are small~~

D2: ~~The moon and the small tree are big~~

Bag of Words

	ball	big	cat	moon	small	table	tree	window	zoo
D1	0	1	2	0	1	1	0	1	0
D2	0	0	0	0	1	1	0	1	0
D3	0	1	0	1	1	0	0	0	0

TF-IDF (W) = TF(W,d) * IDF(W)

$$\text{idf}(W) = \log \frac{\#\text{(documents)}}{\#\text{(documents containing word W)}}.$$

	ball	big	cat	moon	small	table	tree	window	zoo
D1	0	0.17	0.95	0	0	0.17	0	0.17	0
D2	0	0	0	0	0	0.17	0	0.17	0
D3	0	0.17	0	0.47	0	0	0.47	0	0

Drawbacks of traditional approachess

- Have high dimensionality and are very sparse.
- Don't capture semantics
 - *Edema de glotis != hinchazón de la laringe*
- Don't position of occurrence of words
 - *The hotel was very good and not expensive != The hotel was very expensive and not good*

Figure taken from <https://laptrinhx.com/word-embeddings-part-2-3902951508/>

Neural Networks for Word embeddings

- Shallow neural networks trained on a large unlabeled corpus to predict a word based on its context (or given a context, to predict the most appropriate word for it).
- Word2Vec, Glove, FastText.
- Efficient in capturing context similarity.
- Cosine distance (fast and efficient).

Distributional Hypothesis
(Harris, 1954)
Words with similar meanings tend
to occur in similar context

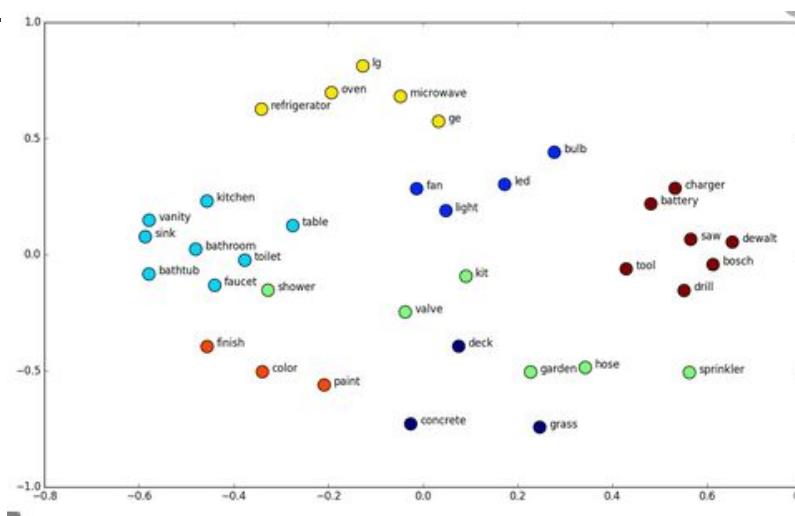


Figure taken from <https://laptrinhx.com/word-embeddings-part-2-3902951508/>

Neural Networks for Word embeddings

- Capture syntactic and semantic information.

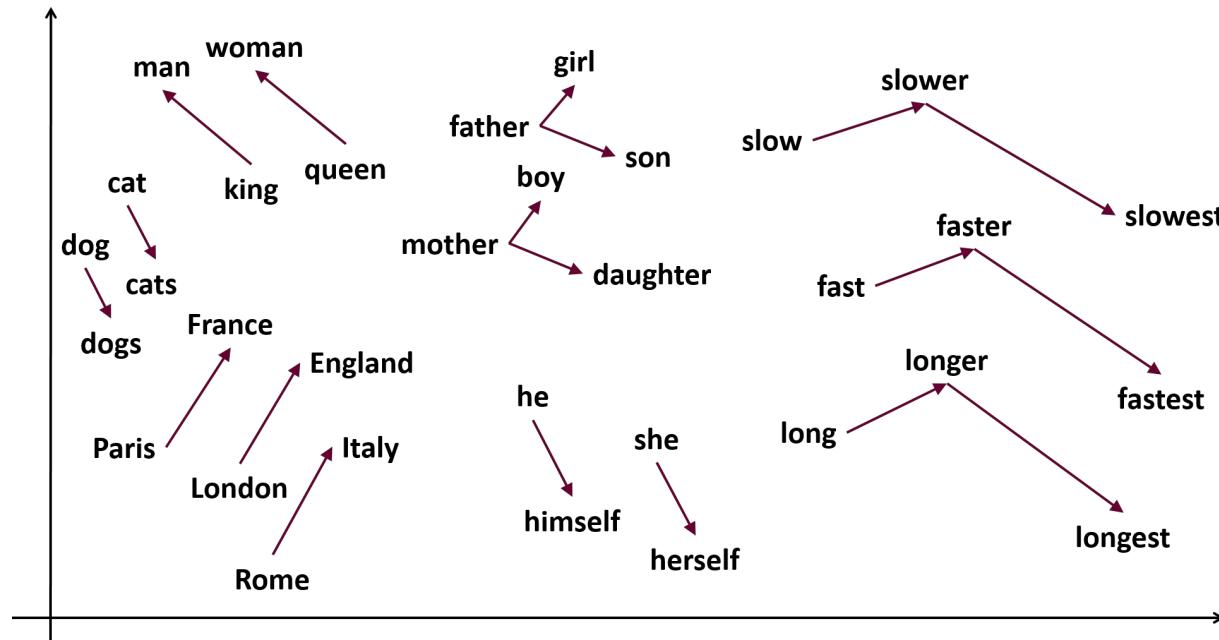
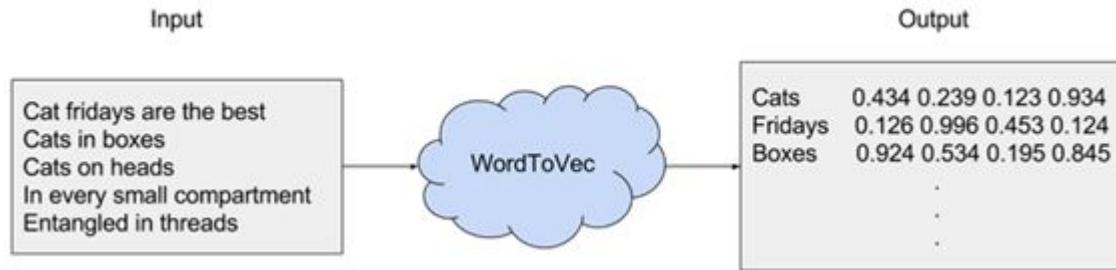


Figure taken from <https://medium.com/analytics-vidhya/implementing-word2vec-in-tensorflow-44f93cf2665f>

Word2Vec



<https://medium.com/@zafaralibagh6/simple-tutorial-on-word-embedding-and-word2vec-43d477624b6d>

Word2Vec

- Proposed two architectures to efficiently create word embeddings: continuous bag-of-words (CBOW) and skip-grams models.
- CBOW computes the conditional probability of a target word given the context words surrounding it across a window of size k.
- Skip-gram model predicts the surrounding context words given the central target word

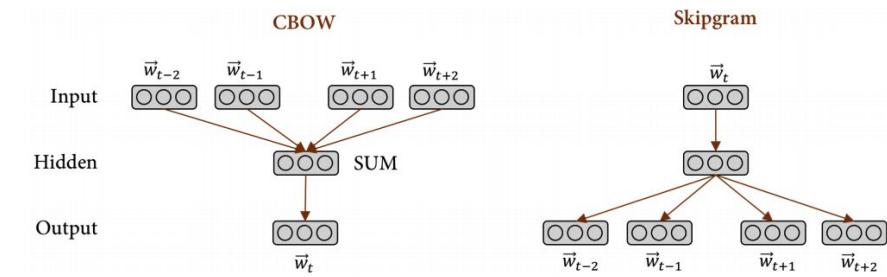


Figure taken from Thomás Mikolov, Kai Chen, Greg Corrado, Jeffrey Dean: Efficient Estimation of Word Representations in Vector Space. ICLR 2013.

Word2Vec: CBOW

- A simple fully connected neural network with one hidden layer.
- The input layer takes a one-hot vector for each context word. The output layer is softmax probability over all words in the vocabulary
- Finally, each word from the vocabulary is represented as two learned vectors, corresponding to the context and target word representations.

$$\mathbf{v}_c = \mathbf{W}_{(k, \cdot)} \quad \text{and} \quad \mathbf{v}_w = \mathbf{W}'_{(\cdot, k)}$$

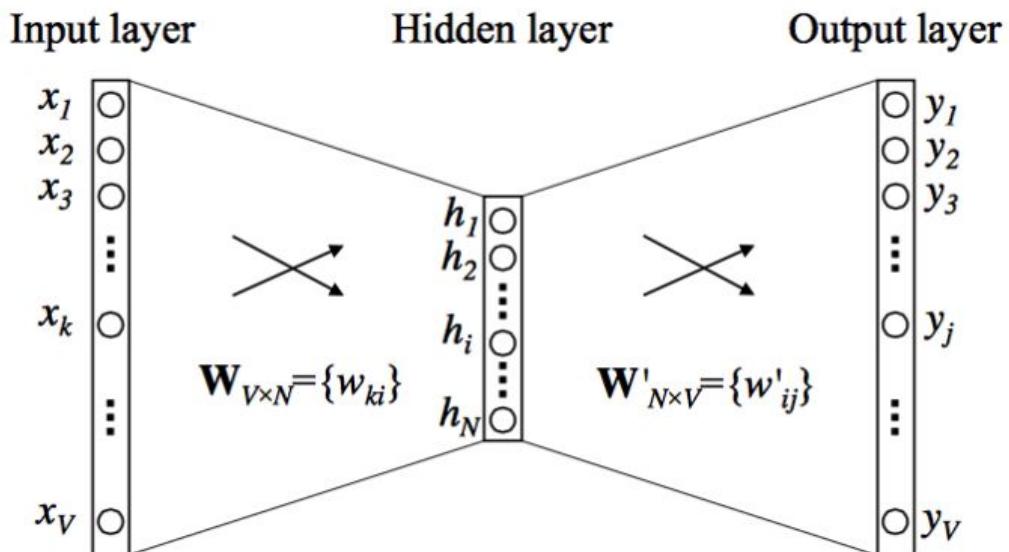
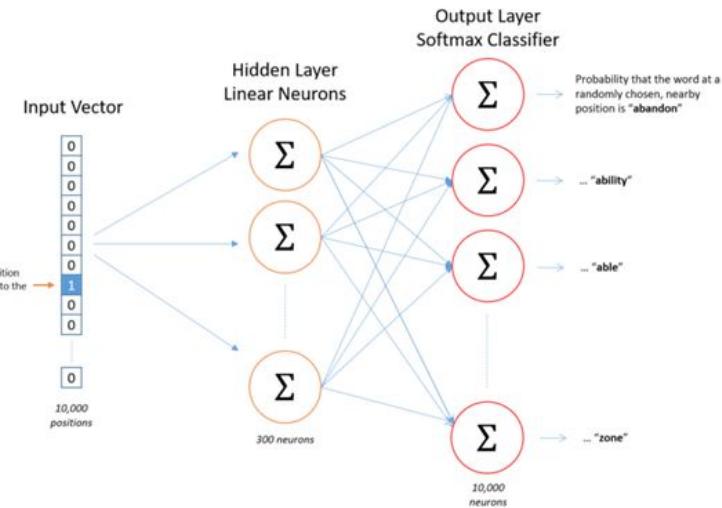


Figure taken from Thomás Mikolov, Kai Chen, Greg Corrado, Jeffrey Dean: Efficient Estimation of Word Representations in Vector Space. ICLR 2013.

Word2Vec: Skip-gram

- input word is represented with a hot-vector, with dimension the size of the vocabulary (total unique words)
- the output is a vector containing for each word in the vocabulary, its probability as neighbour of the input word.
- Produces more accurate results on large datasets.



<https://towardsdatascience.com/skip-gram-nlp-context-words-prediction-algorithm-5bbf34f84e0c>

Drawbacks of neural word embeddings

- Training word embeddings models require large amount of texts and time.
- Inability to represent phrases (“Joe Biden”, “American Airlines”, “Ford Motor Company”).
- Do not handle polysemy and homonyms correctly.
- Some words such as good and bad have very similar word embeddings (very small context window).

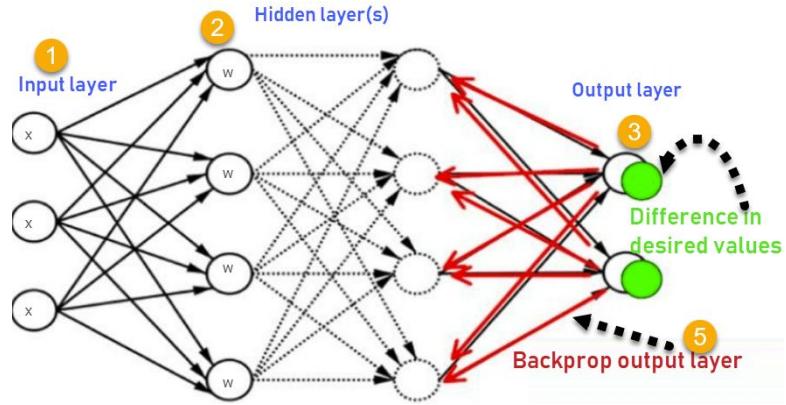
Pre-trained word embedding models

- Repository: <http://vectors.nlpl.eu/repository/>

Outline (first season)

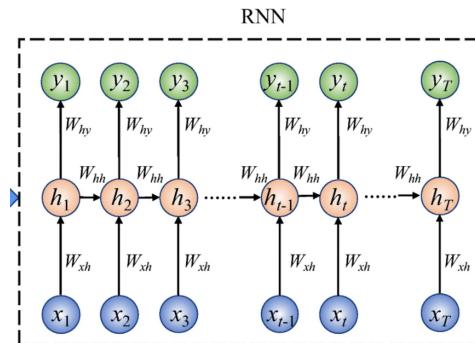
- Introduction
- Word Embeddings
- **Deep learning architectures for NLP**
 - Recurrent Neural Networks
 - Sequence to Sequence
 - Transformer
- Contextual languages models
 - BERT

Artificial Neural Network



Recurrent Neural Network

- Processes **sequences** (words, symbols, images, sensor measurements, etc) of different length.
- Remembers the past and uses it to take the next decision.
- Influenced not just by weights applied on inputs like a regular NN, but also by a “hidden” state vector representing the context based on prior input(s)/output(s)



Applications of RNNs

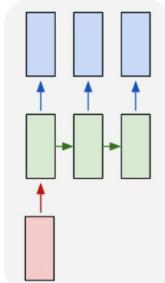
- Computer vision: image and video caption generation, image detection, etc
- Prediction problems (power demand, airline traffic volume, etc)
- Speech recognition
- NLP: *PoS tagging, NER, language modelling, text summarization, machine translation, text classification, etc*

RNN for NLP

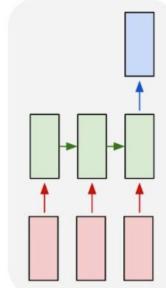
*text summarization
machine translation*

*multilabel text classification
Sequence labeling task: NER,
PoS tags*

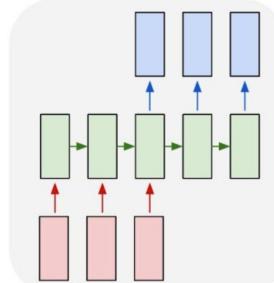
one to many



many to one



many to many



many to many

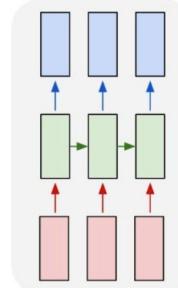
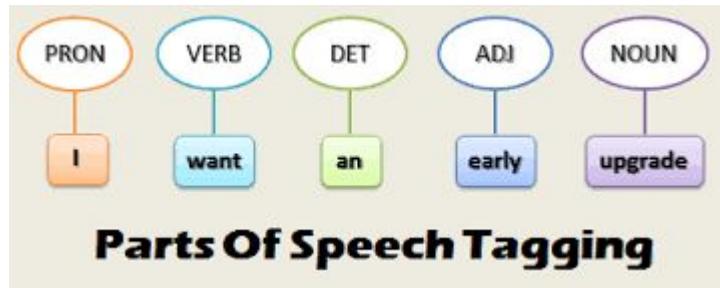


image captioning

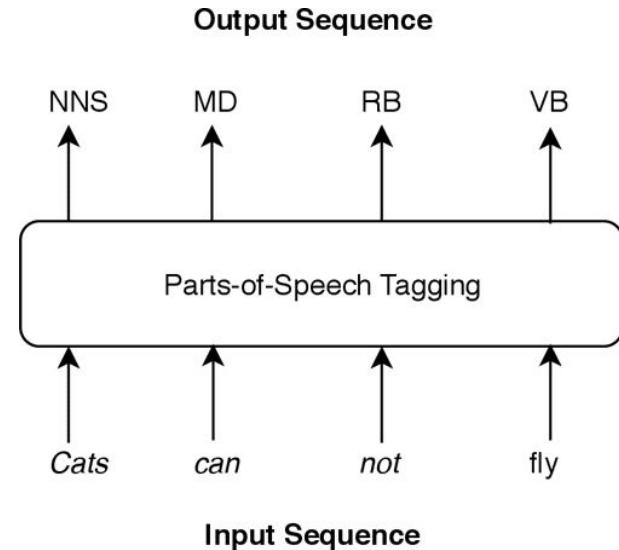
text classification

next word prediction (language modelling)

Why RNN for NLP?

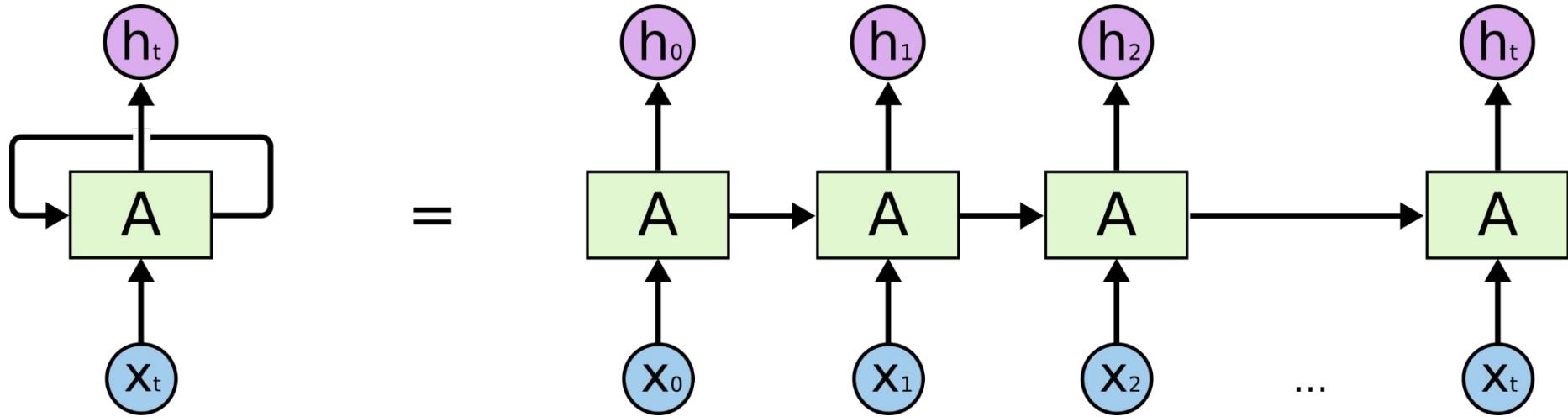


Source: <https://thinkinfi.com/extract-custom-keywords-using-nltk-pos-tagger-in-python/>

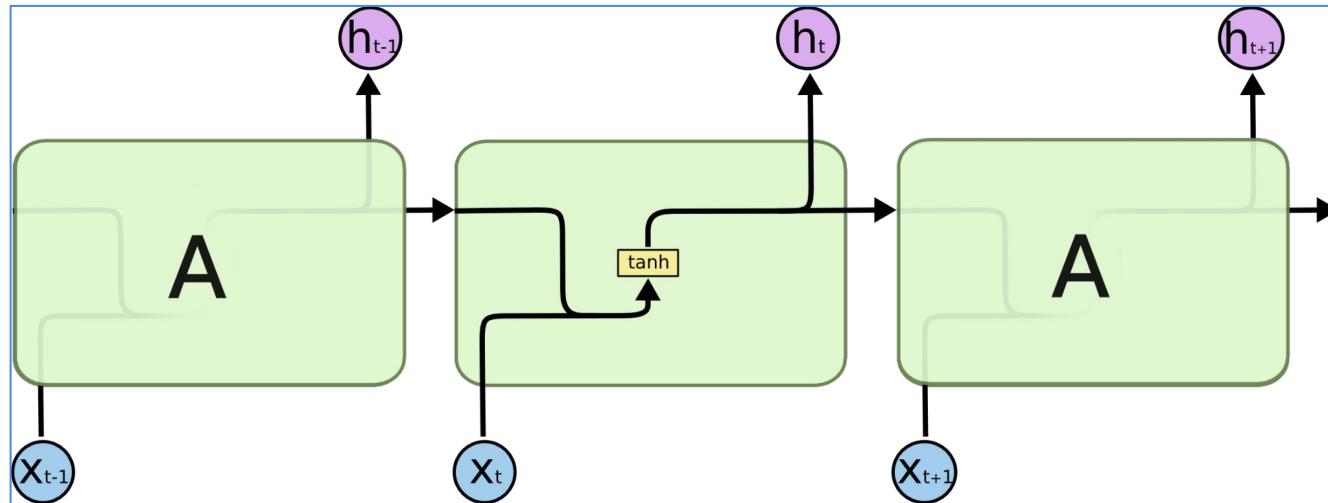


Source: Ahmed, M., Samee, M. R., & Mercer, R. (2018, December). Improving Neural Sequence Labelling Using Additional Linguistic Information. In *2018 17th IEEE International Conference on Machine Learning and Applications (ICMLA)* (pp. 650-657). IEEE.

RNN cells

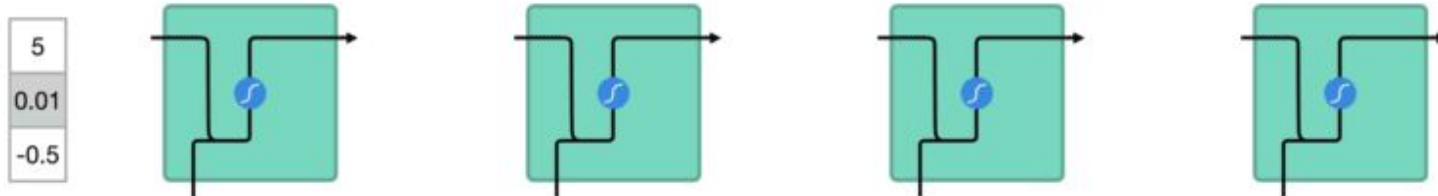


RNN cell



RNN cell

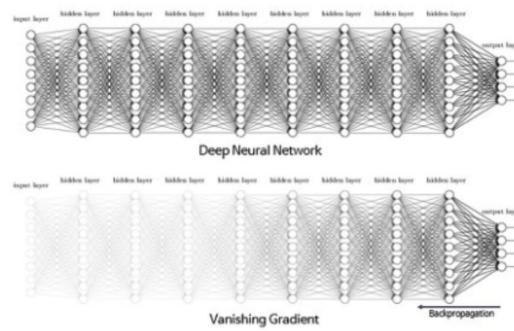
vector transformations without tanh



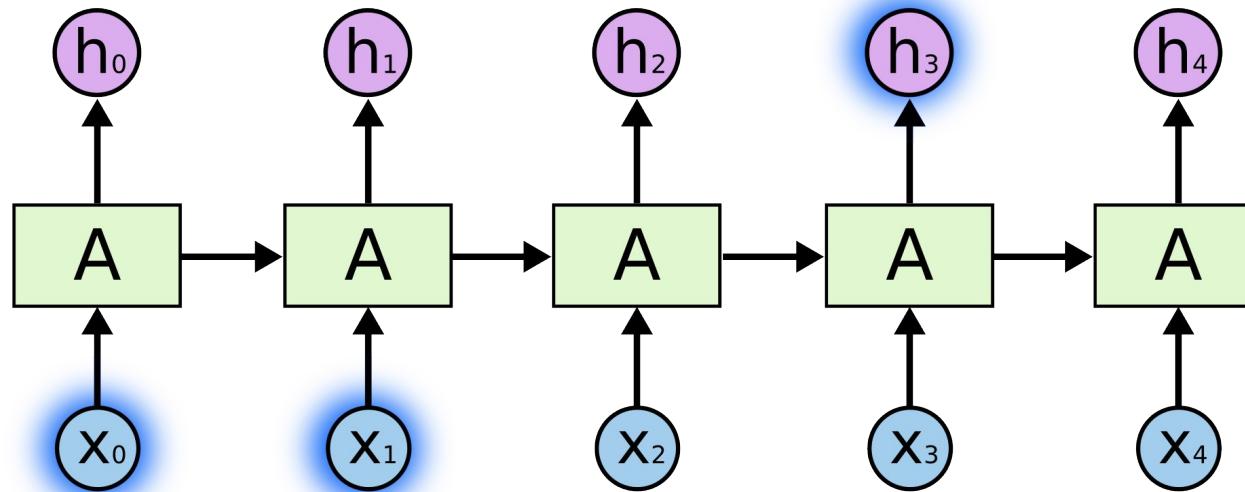
vector transformations with tanh

Vanishing gradient problem

- RNNs suffer from the vanishing gradient.
- **Vanishing gradient** leads to slow convergence

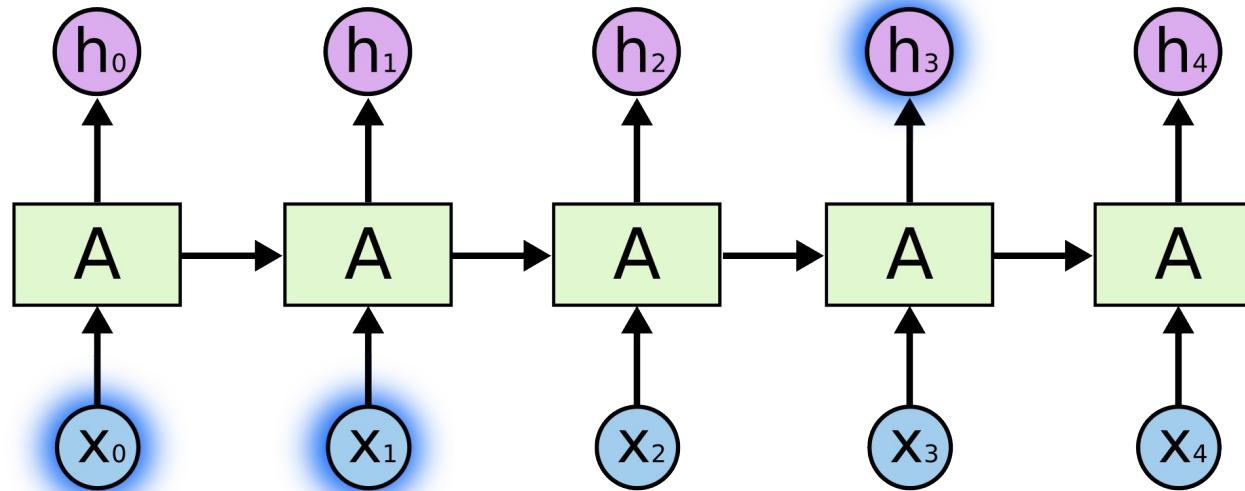


The problem of Long-term dependencies



The clouds are in

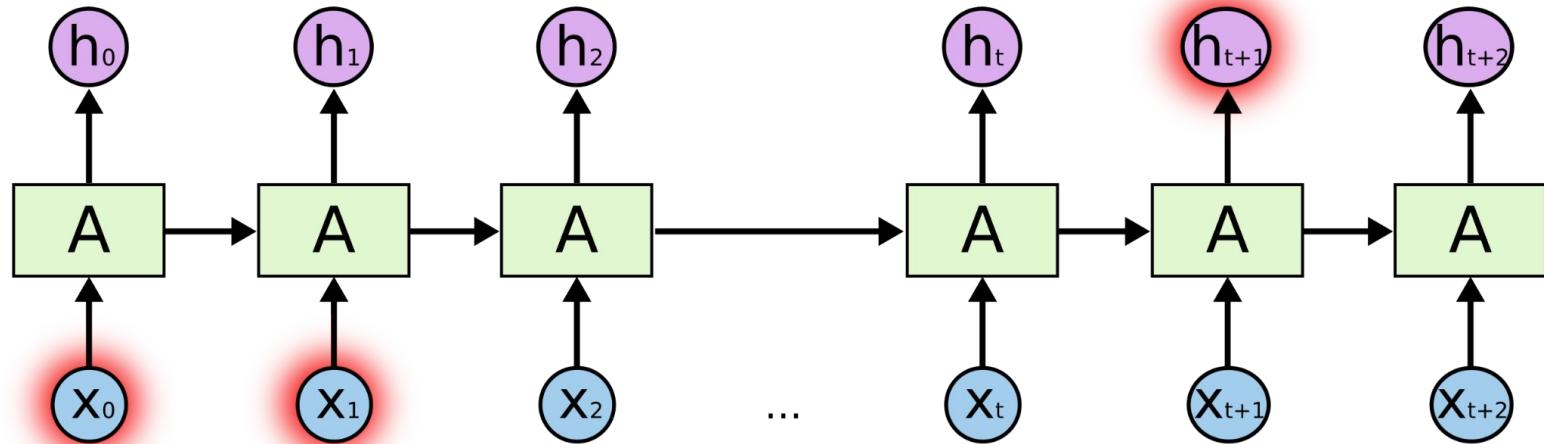
The problem of Long-term dependencies



The clouds are in **sky**

The problem of Long-term dependencies

I grew up in France, in Lyon, surrounded by vineyards, markets and my parents' bistro. I speak fluent French.



The problem of Long-term dependencies

Customers Review 2,491

Thanos

September 2018

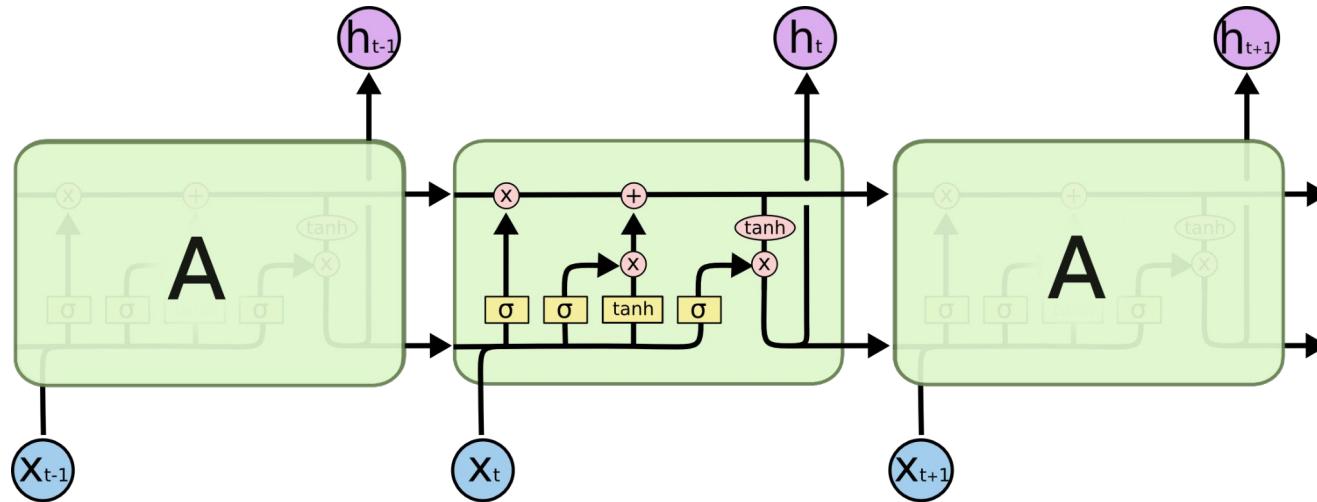
Verified Purchase

Amazing! This box of cereal gave me a perfectly balanced breakfast, as all things should be. I only ate half of it but will definitely be buying again!

A Box of Cereal
\$3.99

Long-short Term Memory Networks (LSTM)

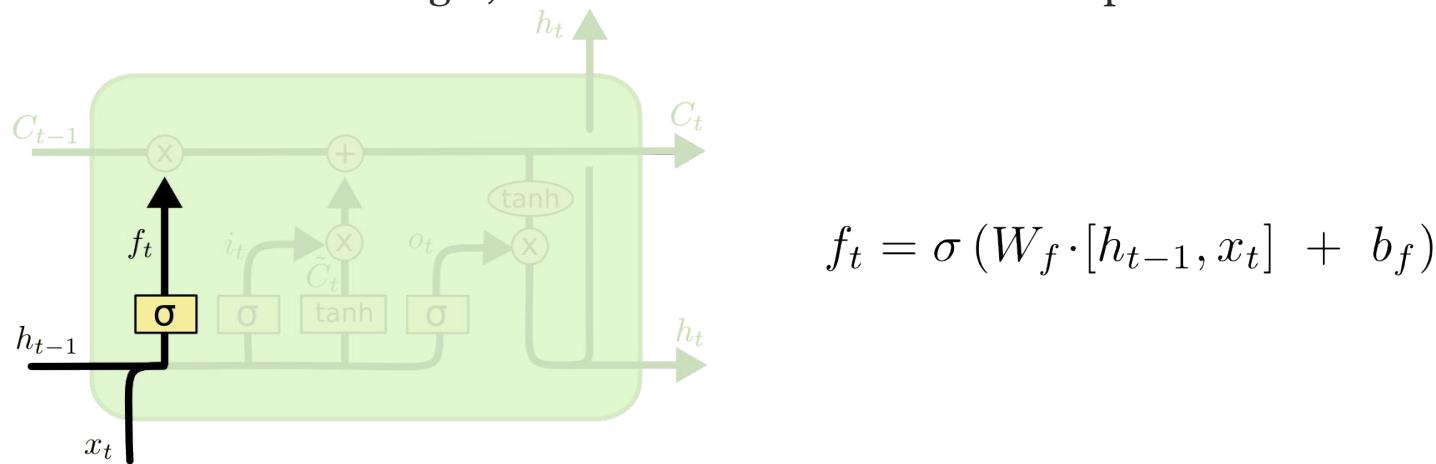
Gates allow LSTM cells to control what information to remove, store, and output to the next cell



Forget gate layer

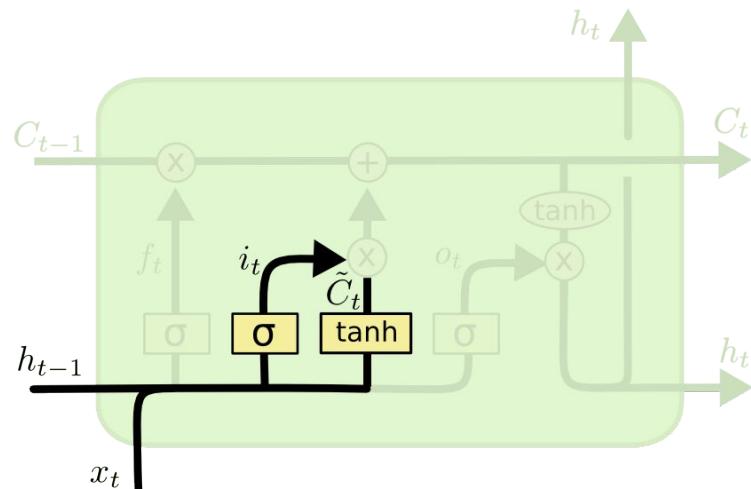
which information to forget

To do this, it uses a sigmoid layer that outputs a number between 0 and 1: the closer to 0 means to forget, and the closer to 1 means to keep.



Input gate layer

what information is relevant to update in the current cell state of the **LSTM** unit.

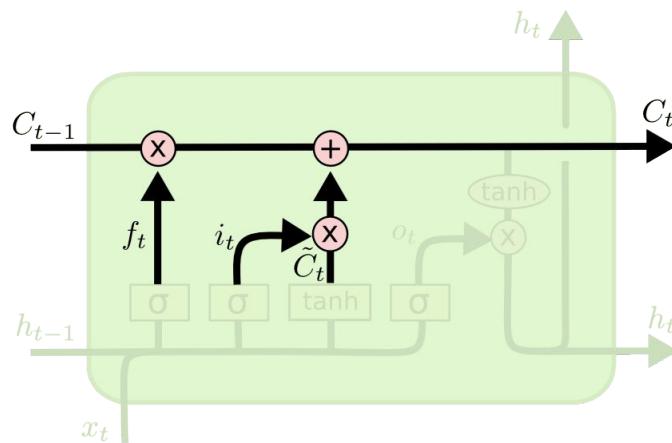


$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

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

The new cell state

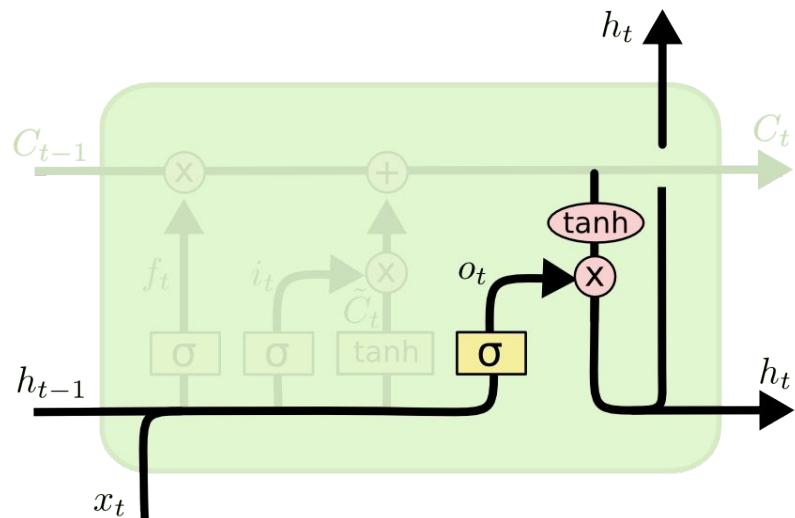
We update the old cell state, C_{t-1} into the new cell state C_t (multiplying the old state by f_t , forgetting the things we decided to forget earlier). Then we add $i_t * \tilde{C}_t$, which is the new candidate values.



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

Output gate layer

what information will be passed to the next

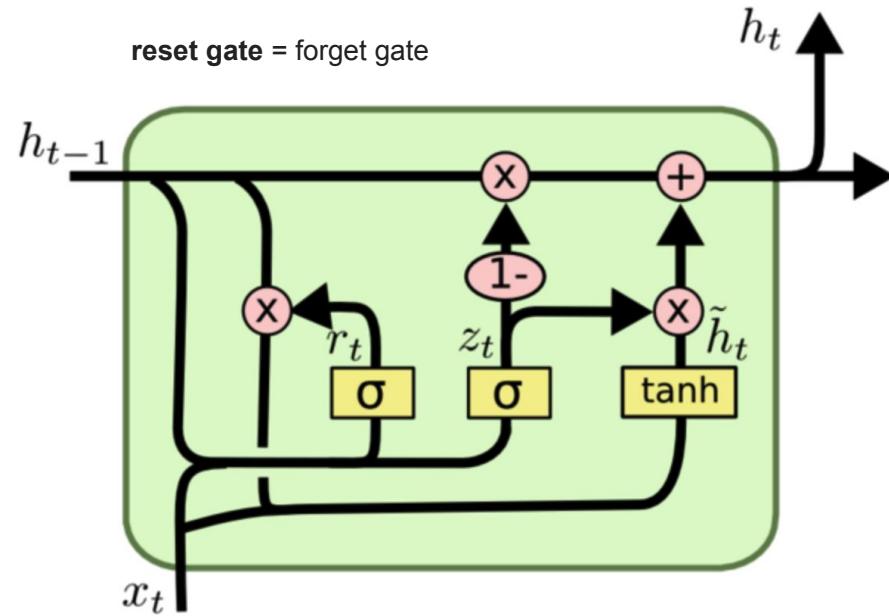


$$o_t = \sigma (W_o [h_{t-1}, x_t] + b_o)$$

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

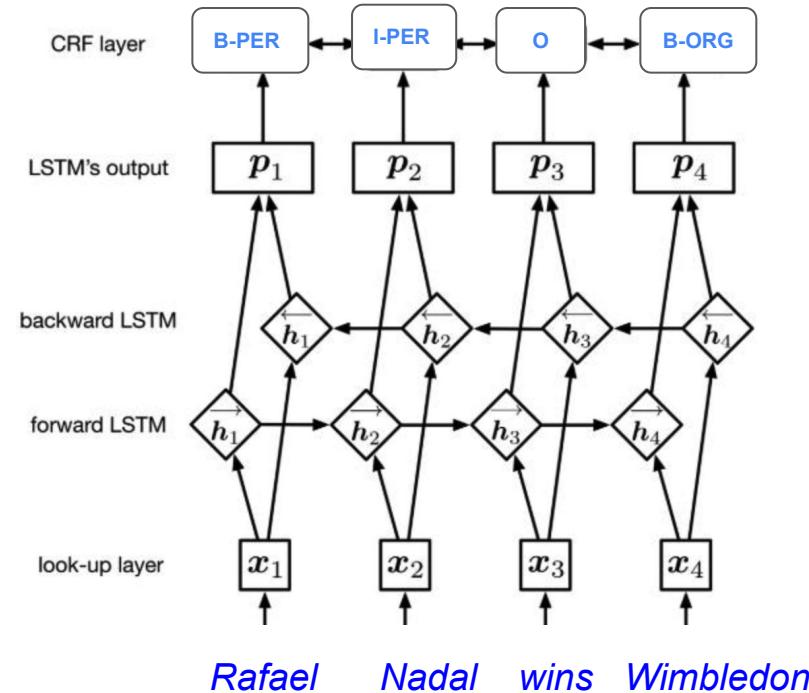
GRU unit

GRUs are slightly faster and easier to run than LSTM

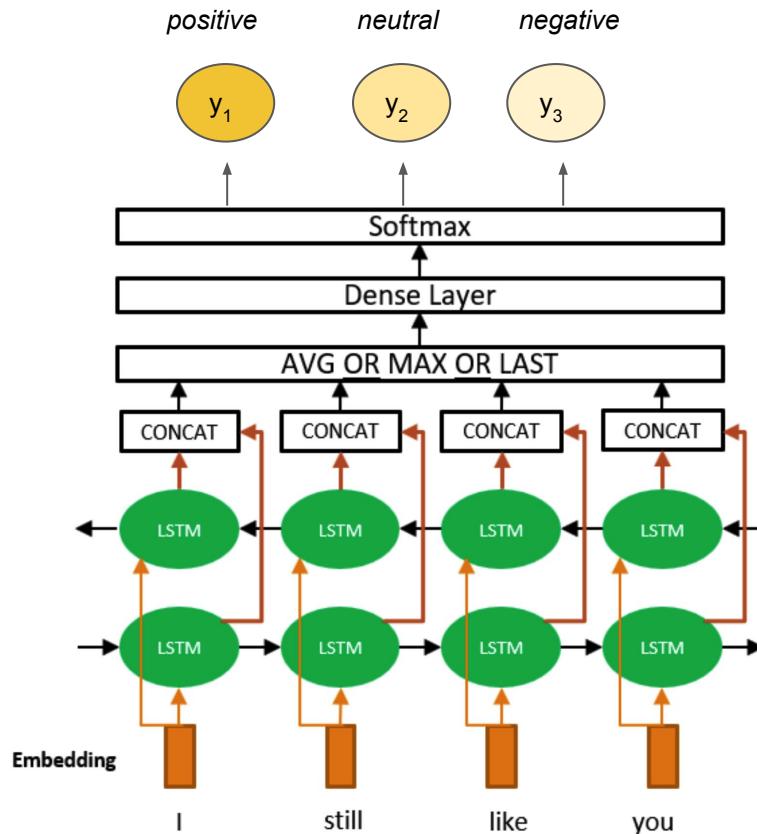


update gate determines both how much information to keep from the last state and how much information to let in from the previous layer.

Recurrent Neural Network (bidirectional)



Recurrent Neural Network (bidirectional)



Saroufim, C., Almatarfy, A., & Hady, M. A. (2018, October). Language independent sentiment analysis with sentiment-specific word embeddings. In *Proceedings of the 9th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis* (pp. 14-23).

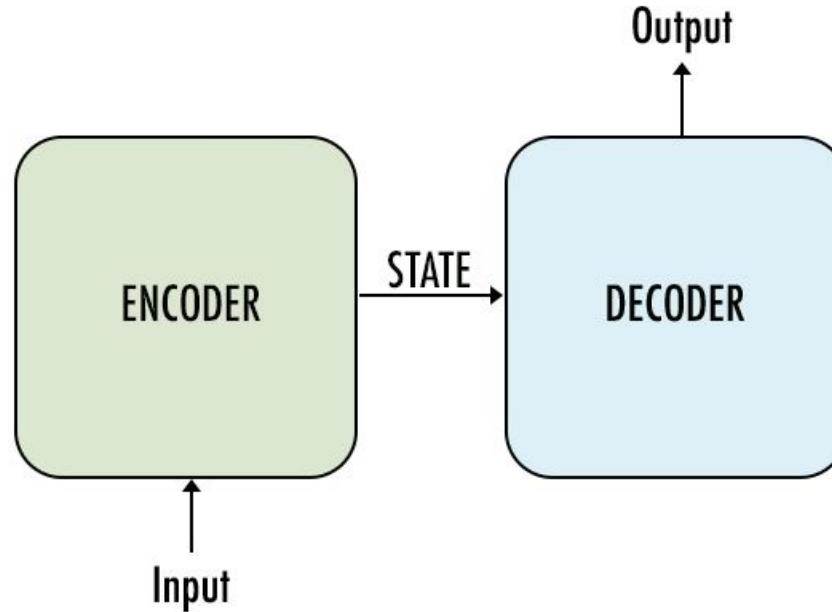
Keys ideas in RNN

- RNN can process sequences.
- RNN takes its decisions based on the learnt from the past.
- Vanilla RNN (\tanh) suffers from vanishing gradient problem.
- LSTM capable to learn long term dependencies
- GRU is simpler and more efficient than LSTM
- The main drawback of RNNs is that they **cannot be parallelized**. Then, **computational cost** becomes **critical** dealing with long sequences.

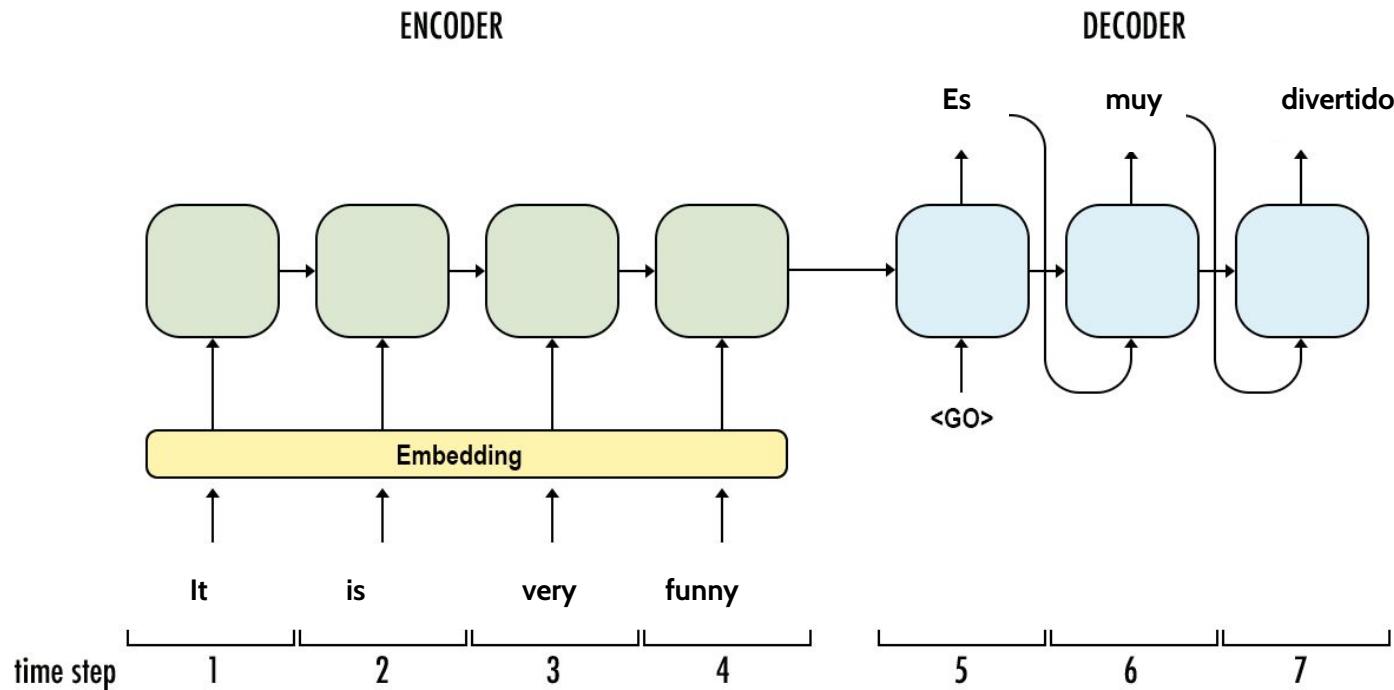
Outline (first season)

- Introduction
- Word Embeddings
- Deep learning architectures for NLP
 - Recurrent Neural Networks
 - **Sequence to Sequence**
 - Attention mechanism
 - Transformer
- Contextual languages models
 - BERT

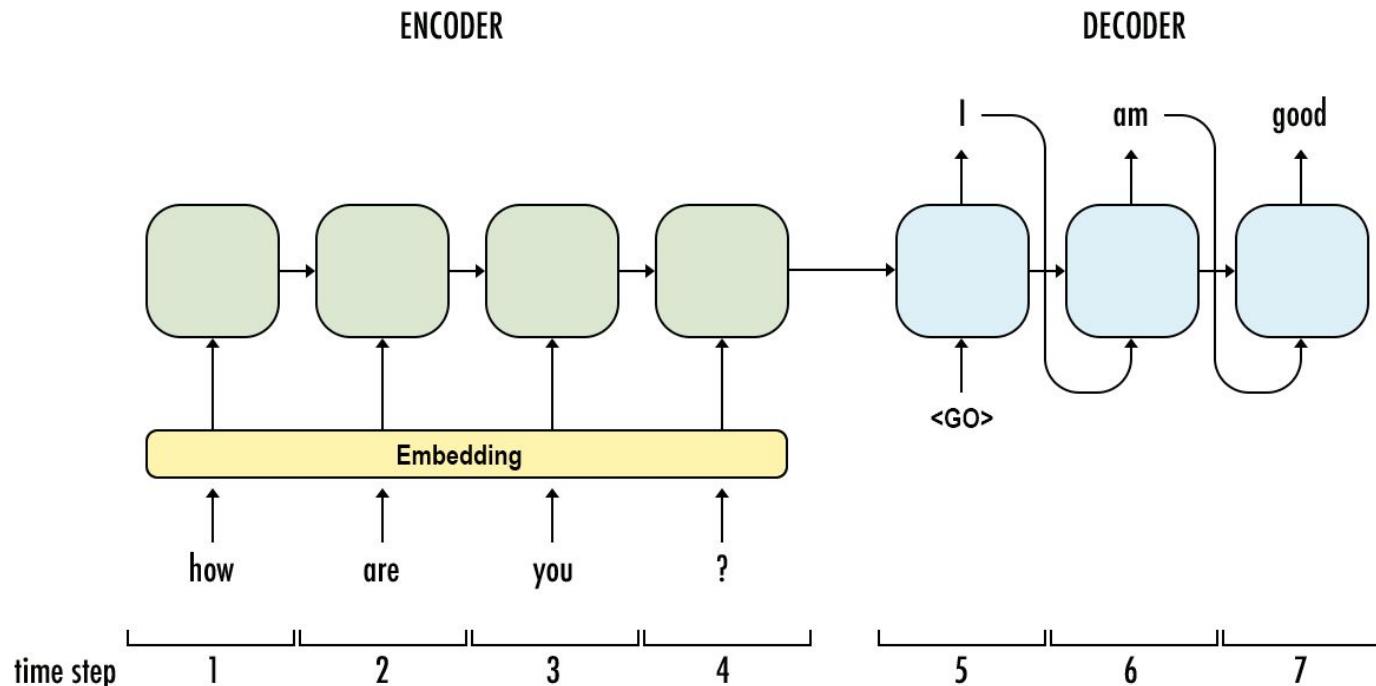
Sequence to Sequence Model (Seq2Seq)



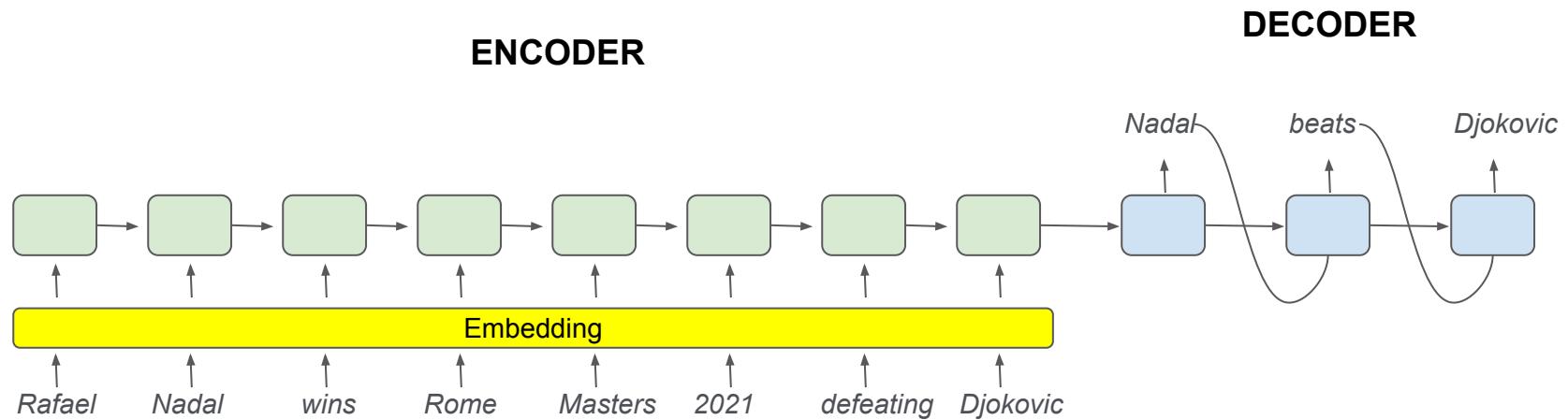
Seq2Seq for machine translation



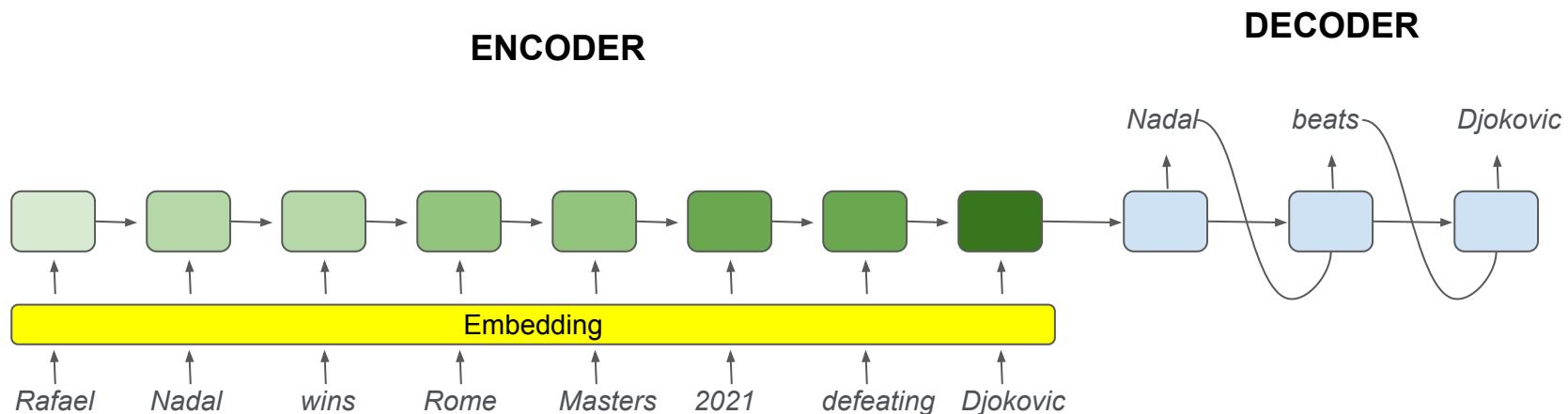
Seq2Seq for chatbot



Seq2Seq for text summarization



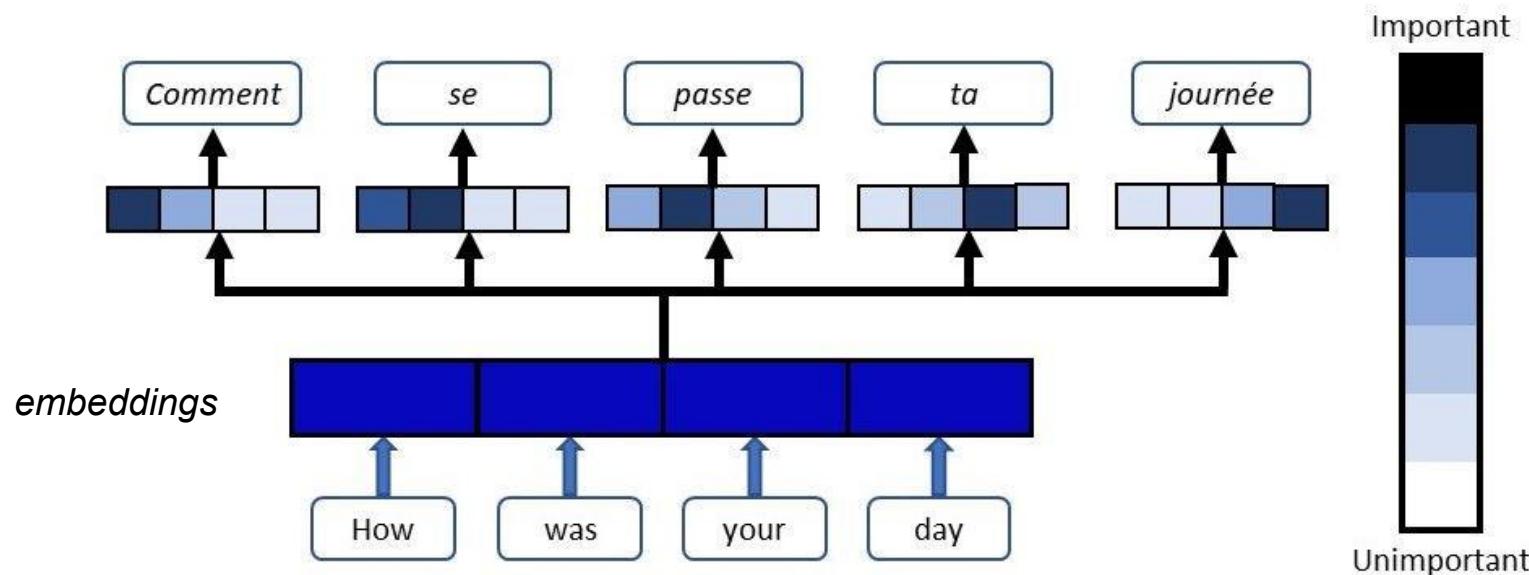
Seq2Seq



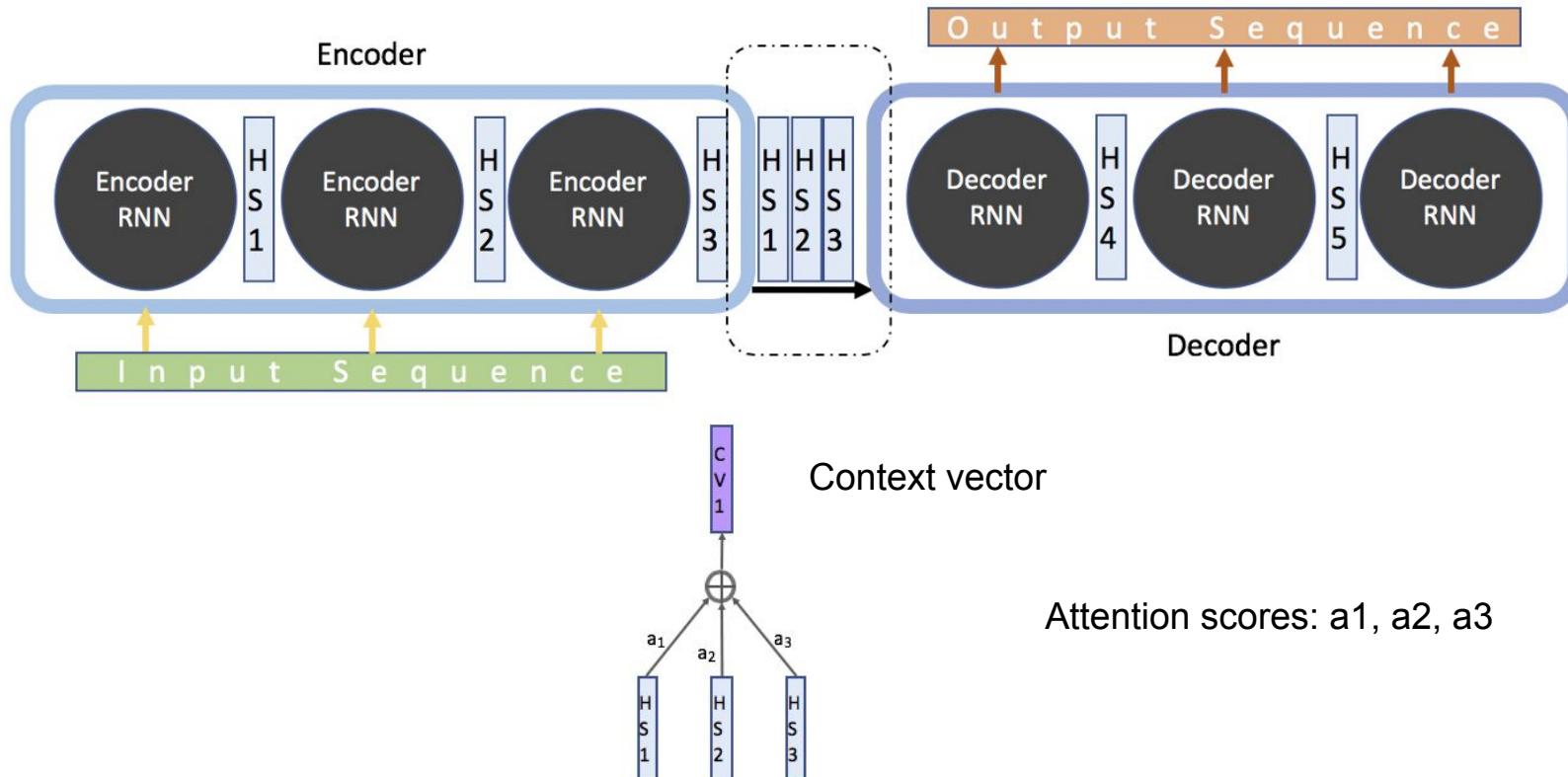
Drawbacks:

- 1) The output sequence depends heavily on the context defined by the hidden state in the final output of the encoder, making it challenging for the model to **deal with long sequences**.
- 2) **RNN steps cannot** be parallelized. Then, **computational cost** becomes **critical** for long sequences.

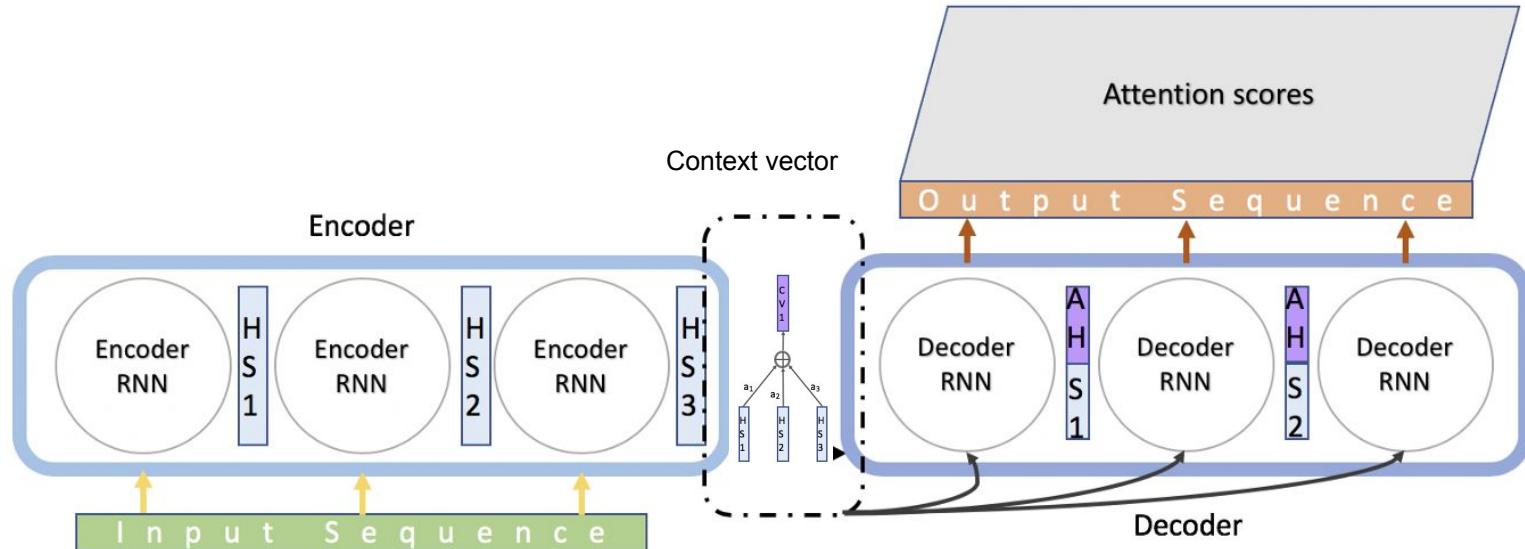
Basic idea of attention mechanism



Attention mechanism

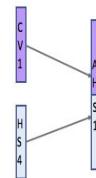


Seq2Seq with attention (attention scoring)

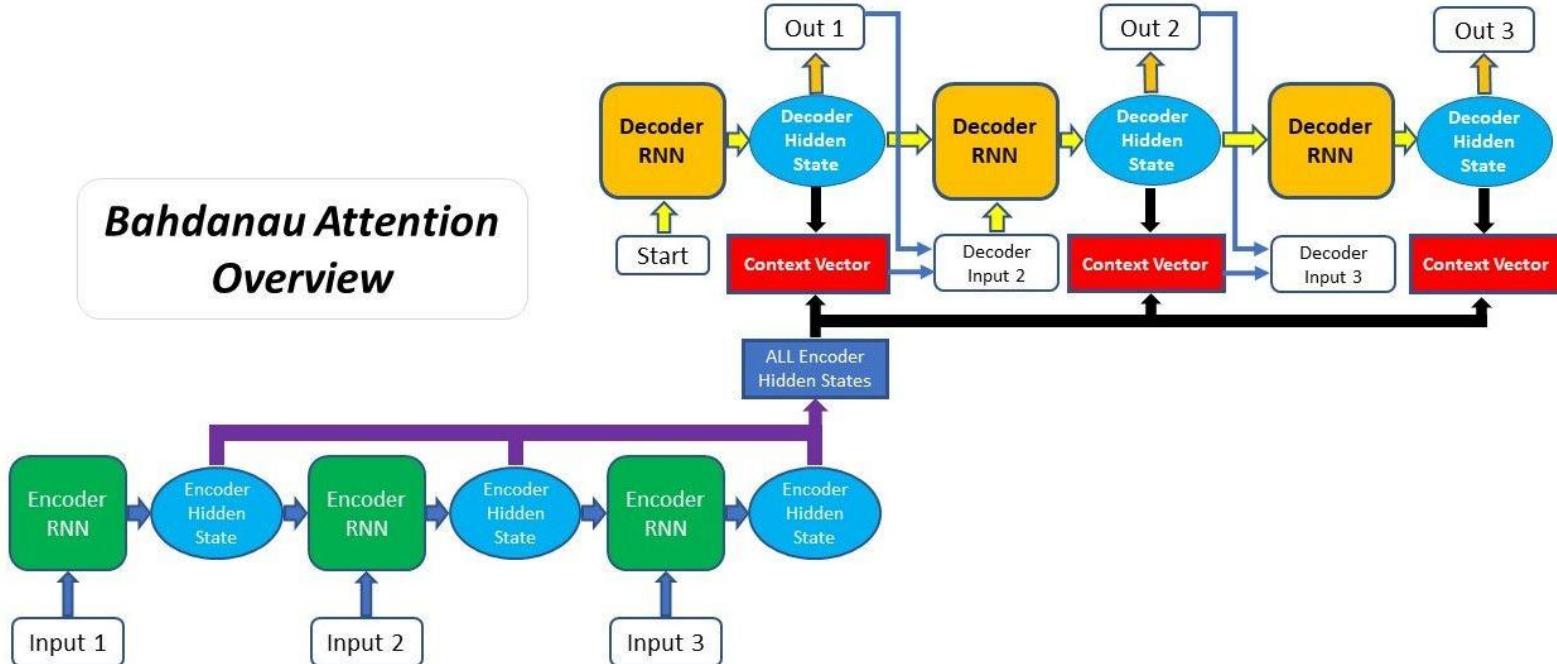


Computational cost are still **critical** for long sequences.

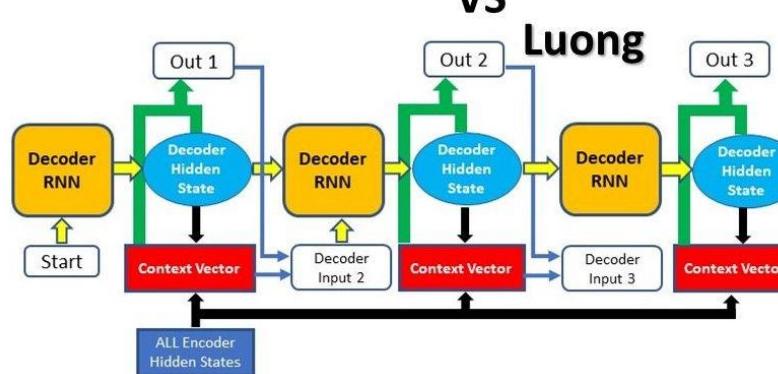
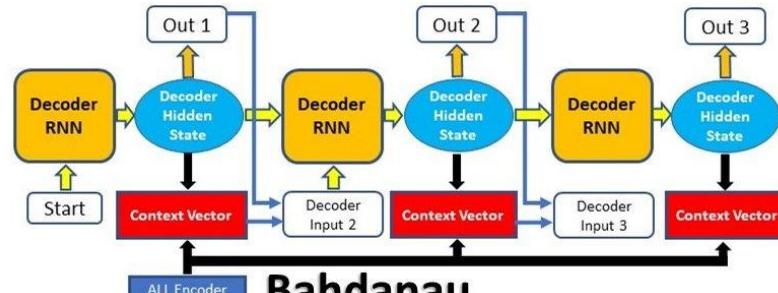
attention hidden vector



Bahdanau's attention mechanism

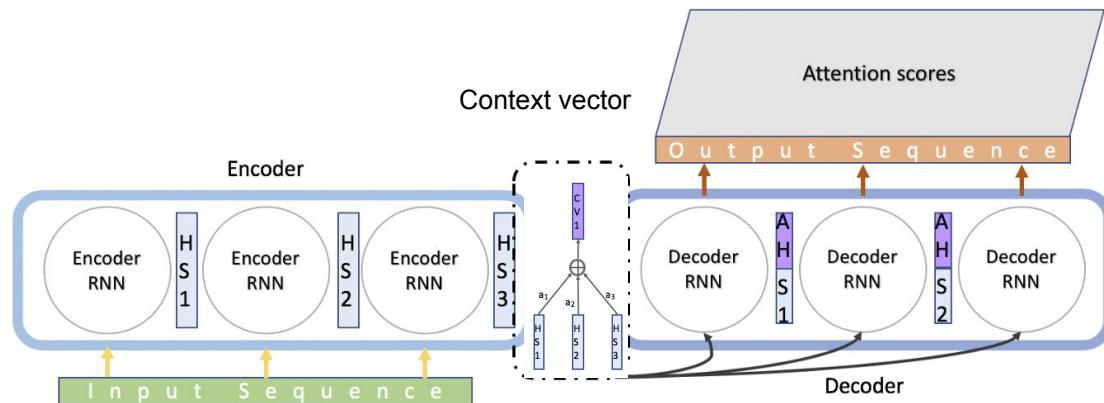


Bahdanau's attention versus Luong's attention



Seq2Seq with attention (attention scoring)

Computational cost are still **critical** for long sequences!!!.



Outline (first season)

- Introduction
- Word Embeddings
- Deep learning architectures for NLP
 - Recurrent Neural Networks
 - Sequence to Sequence
 - **Transformer**
 - self-attention and multi-head attention mechanisms
- Contextual languages models
 - BERT

Transformer

- It is a type of network that applies **attention mechanisms** to gather information about the **relevant context of a given word**, and then **encode that context** in a rich vector that smartly represents the word.
- It is also a Seq2Seq model (encoder-decoder) that only uses attention mechanisms instead of RNNs.
 - **self-attention**
 - **multi-head attention**

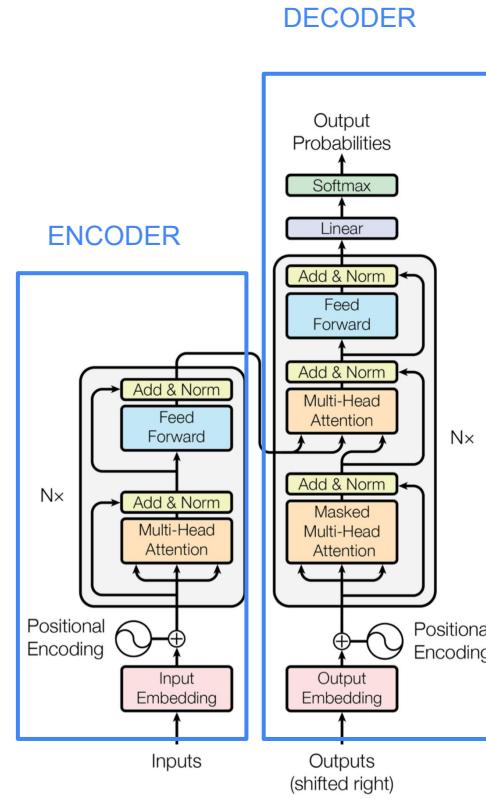
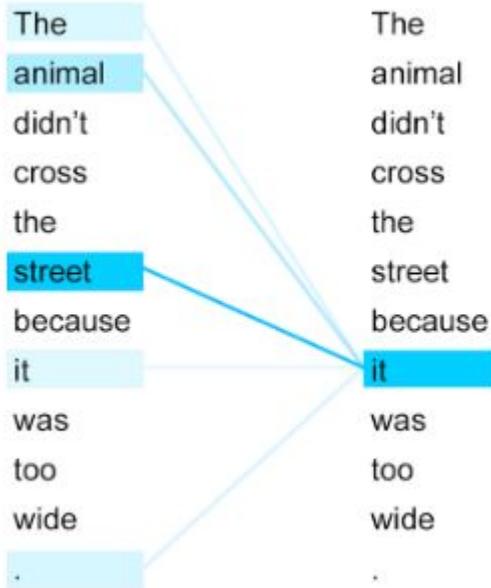


Figure 1: The Transformer - model architecture.

Basic idea of self-attention



Advantages:

- Can compute the weights in **parallel**, for all tokens at once.
- As it also sees all the other inputs, it can easily **preserve long-term dependencies** in long text sequences.
- Therefore, **self-attention** could completely replace RNN

Encoder Self-attention and Decoder self-attention

Encoder Self-Attention computes the interaction between each input word with other input words.

Decoder Self-Attention operates on each word of the target sequence. That is, it computes the interaction between each target word with other target words

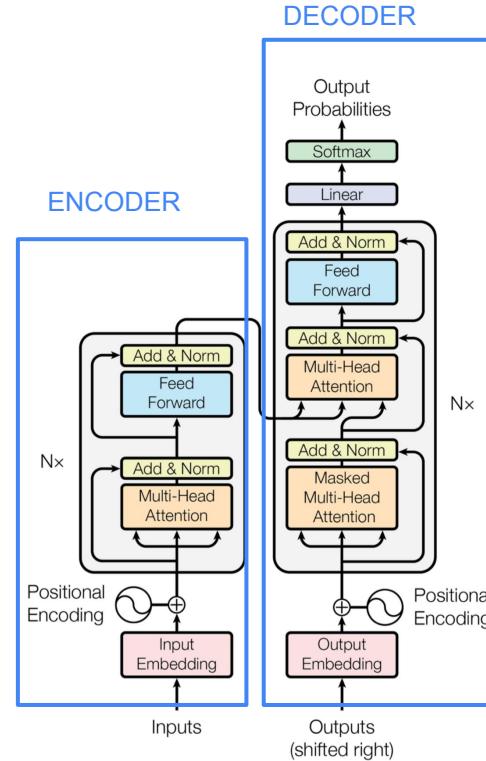
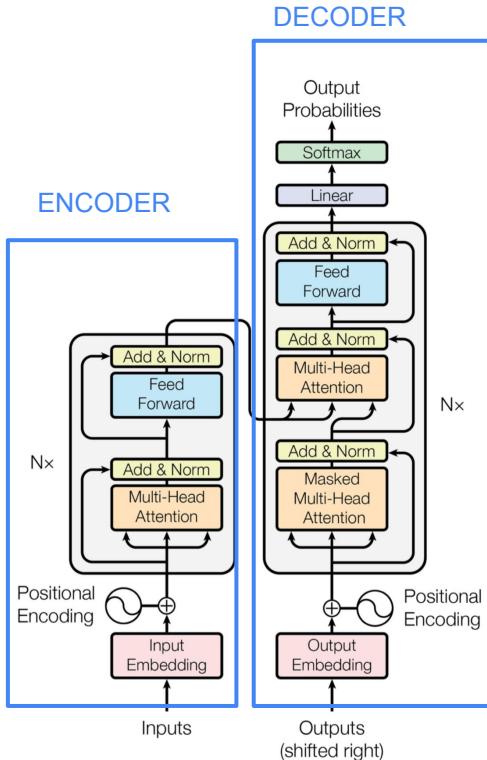


Figure 1: The Transformer - model architecture.

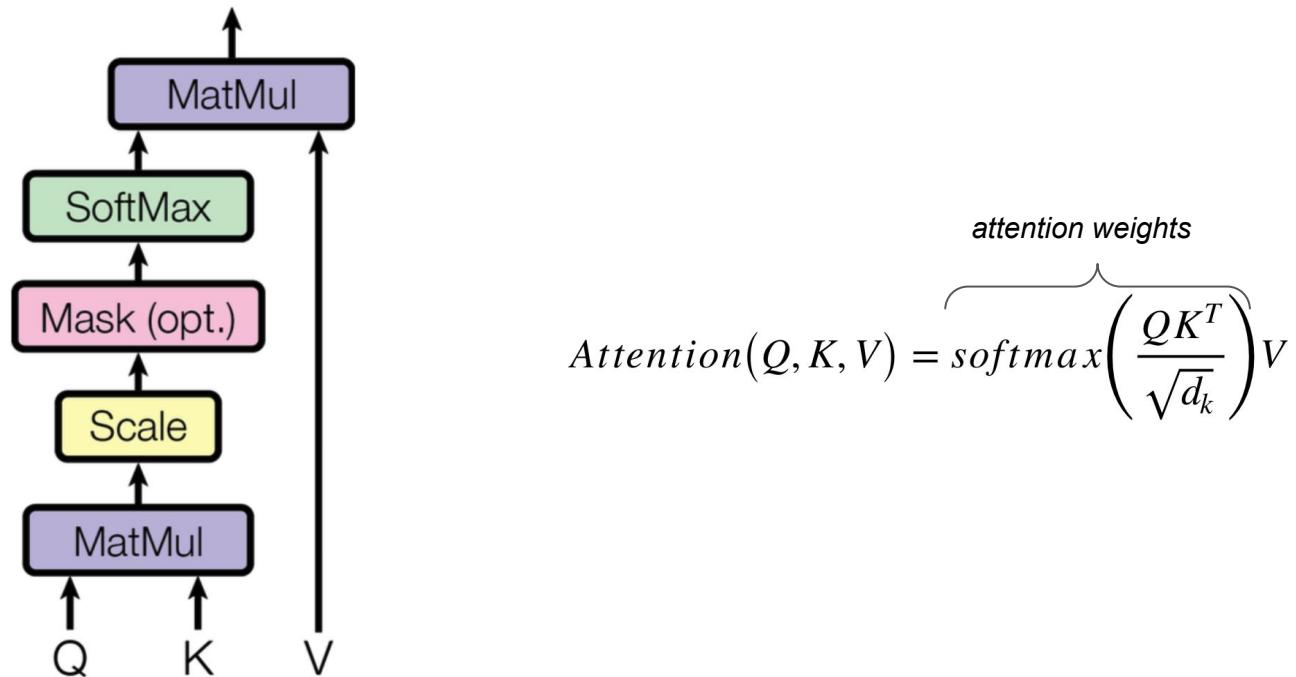
Encoder-Decoder Attention



The **Encoder-Decoder Attention** computes the interaction between each target word with each input word. Therefore, it masks out the later words in the target output.

Figure 1: The Transformer - model architecture.

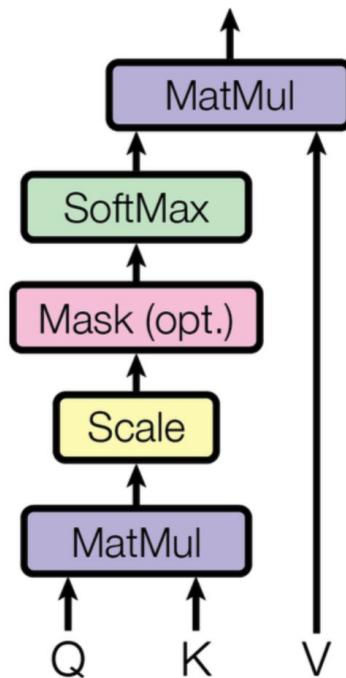
How to implement self-attention



<https://medium.com/inside-machine-learning/what-is-a-transformer-d07dd1fbec04>

<https://medium.com/analytics-vidhya/masking-in-transformers-self-attention-mechanism-bad3c9ec235c>

Why masking?



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

- all vectors must have the same dimension (sequences could have different number of tokens)
- the **Encoder-Decoder Attention** should only use the previous tokens to the current token to predict its output. Thus, we have to hide the next tokens.

<https://medium.com/inside-machine-learning/what-is-a-transformer-d07dd1fbec04>

<https://medium.com/analytics-vidhya/masking-in-transformers-self-attention-mechanism-bad3c9ec235c>

How to implement self-attention

X: embedding matrix (each row is the word embedding of a word in the input sequence)
dimension= 512 (only 4 in figure)

$$\begin{matrix} \text{X} \\ \begin{matrix} \textcolor{green}{\square} & \textcolor{green}{\square} & \textcolor{green}{\square} & \textcolor{green}{\square} \\ \textcolor{green}{\square} & \textcolor{green}{\square} & \textcolor{green}{\square} & \textcolor{green}{\square} \\ \textcolor{green}{\square} & \textcolor{green}{\square} & \textcolor{green}{\square} & \textcolor{green}{\square} \\ \textcolor{green}{\square} & \textcolor{green}{\square} & \textcolor{green}{\square} & \textcolor{green}{\square} \end{matrix} \end{matrix} \times \begin{matrix} \text{W}^Q \\ \begin{matrix} \textcolor{purple}{\square} & \textcolor{purple}{\square} & \textcolor{purple}{\square} & \textcolor{purple}{\square} \\ \textcolor{purple}{\square} & \textcolor{purple}{\square} & \textcolor{purple}{\square} & \textcolor{purple}{\square} \\ \textcolor{purple}{\square} & \textcolor{purple}{\square} & \textcolor{purple}{\square} & \textcolor{purple}{\square} \\ \textcolor{purple}{\square} & \textcolor{purple}{\square} & \textcolor{purple}{\square} & \textcolor{purple}{\square} \end{matrix} \end{matrix} = \begin{matrix} \text{Q} \\ \begin{matrix} \textcolor{purple}{\square} & \textcolor{purple}{\square} & \textcolor{purple}{\square} & \textcolor{purple}{\square} \\ \textcolor{purple}{\square} & \textcolor{purple}{\square} & \textcolor{purple}{\square} & \textcolor{purple}{\square} \end{matrix} \end{matrix}$$

$$\begin{matrix} \text{X} \\ \begin{matrix} \textcolor{green}{\square} & \textcolor{green}{\square} & \textcolor{green}{\square} & \textcolor{green}{\square} \\ \textcolor{green}{\square} & \textcolor{green}{\square} & \textcolor{green}{\square} & \textcolor{green}{\square} \\ \textcolor{green}{\square} & \textcolor{green}{\square} & \textcolor{green}{\square} & \textcolor{green}{\square} \\ \textcolor{green}{\square} & \textcolor{green}{\square} & \textcolor{green}{\square} & \textcolor{green}{\square} \end{matrix} \end{matrix} \times \begin{matrix} \text{W}^K \\ \begin{matrix} \textcolor{orange}{\square} & \textcolor{orange}{\square} & \textcolor{orange}{\square} & \textcolor{orange}{\square} \\ \textcolor{orange}{\square} & \textcolor{orange}{\square} & \textcolor{orange}{\square} & \textcolor{orange}{\square} \\ \textcolor{orange}{\square} & \textcolor{orange}{\square} & \textcolor{orange}{\square} & \textcolor{orange}{\square} \\ \textcolor{orange}{\square} & \textcolor{orange}{\square} & \textcolor{orange}{\square} & \textcolor{orange}{\square} \end{matrix} \end{matrix} = \begin{matrix} \text{K} \\ \begin{matrix} \textcolor{orange}{\square} & \textcolor{orange}{\square} & \textcolor{orange}{\square} & \textcolor{orange}{\square} \\ \textcolor{orange}{\square} & \textcolor{orange}{\square} & \textcolor{orange}{\square} & \textcolor{orange}{\square} \end{matrix} \end{matrix}$$

$$\begin{matrix} \text{X} \\ \begin{matrix} \textcolor{green}{\square} & \textcolor{green}{\square} & \textcolor{green}{\square} & \textcolor{green}{\square} \\ \textcolor{green}{\square} & \textcolor{green}{\square} & \textcolor{green}{\square} & \textcolor{green}{\square} \\ \textcolor{green}{\square} & \textcolor{green}{\square} & \textcolor{green}{\square} & \textcolor{green}{\square} \\ \textcolor{green}{\square} & \textcolor{green}{\square} & \textcolor{green}{\square} & \textcolor{green}{\square} \end{matrix} \end{matrix} \times \begin{matrix} \text{W}^V \\ \begin{matrix} \textcolor{blue}{\square} & \textcolor{blue}{\square} & \textcolor{blue}{\square} & \textcolor{blue}{\square} \\ \textcolor{blue}{\square} & \textcolor{blue}{\square} & \textcolor{blue}{\square} & \textcolor{blue}{\square} \\ \textcolor{blue}{\square} & \textcolor{blue}{\square} & \textcolor{blue}{\square} & \textcolor{blue}{\square} \\ \textcolor{blue}{\square} & \textcolor{blue}{\square} & \textcolor{blue}{\square} & \textcolor{blue}{\square} \end{matrix} \end{matrix} = \begin{matrix} \text{V} \\ \begin{matrix} \textcolor{blue}{\square} & \textcolor{blue}{\square} & \textcolor{blue}{\square} & \textcolor{blue}{\square} \\ \textcolor{blue}{\square} & \textcolor{blue}{\square} & \textcolor{blue}{\square} & \textcolor{blue}{\square} \end{matrix} \end{matrix}$$

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

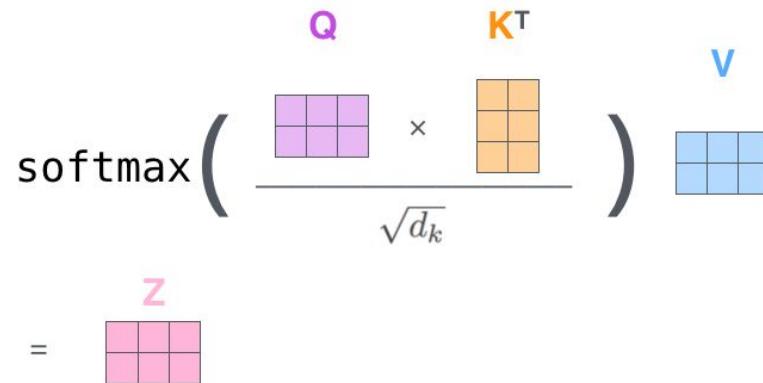
WQ, WK, WV: weight matrices we've trained

Three vectors for each of the input vector: the query vector, the key vector, and the value vector.

- (Q) represents the current word to be represented
- the key (K) are the input vectors for all words in the input sequence. It will represent the relevance to the query,
- the value (V) represent the actual contents of the input.

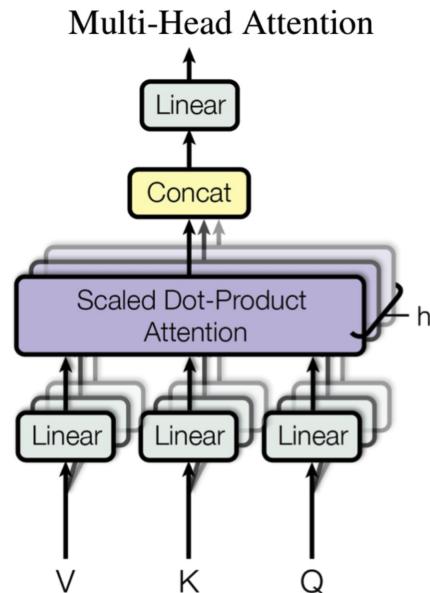
How to implement self-attention

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



The output of the attention head is a weighted sum of the values vector

Multi-head attention



Attention module repeats its computations multiple times in parallel. Each of these is called an Attention Head.

Keys ideas of transformers

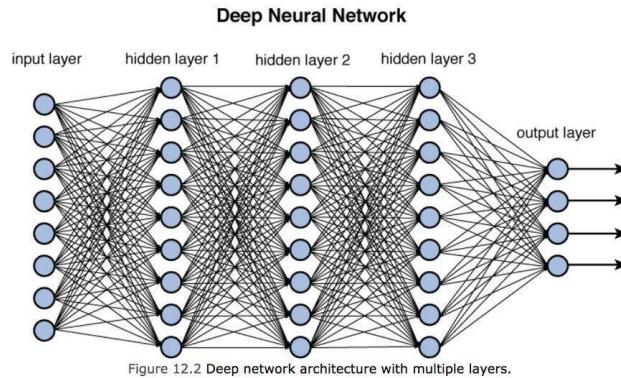
- Transformers is a Seq2Seq model (encoder-decoder) that only uses attention mechanisms instead of RNNs.
- Self-attention allows to preserve long-term dependencies.
- Multi-head attention allows to jointly attend to information from different representation **subspaces** at different positions and more parallelization.

Outline (first season)

- Introduction
- Word Embeddings
- Deep learning architectures for NLP
 - Recurrent Neural Networks
 - Sequence to Sequence
 - Transformer
- **Contextual languages models: BERT**

Why are important model languages?

- Proper language representation is key for developing general-purpose language understanding methods.



Types of language models

- Approaches:
 - Context-free language models
 - Contextual language models

Contextual-free models (such as Word2Vec or Glove)

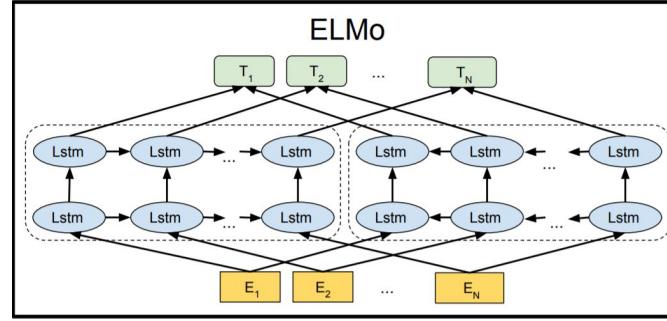
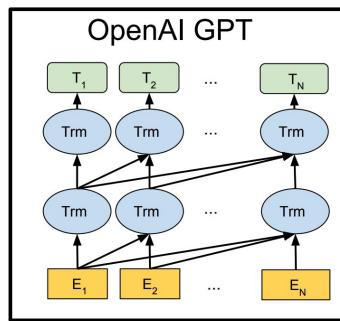
- Generate a single word embedding representation for each word.
- Cannot appropriately represent **polysemy words**.
- **Ambiguity** is one of the **biggest** challenges in **NLP**

- Work out the *solution* in your head.
- Heat the *solution* to 75° Celsius.
- The *key* broke in the lock.
- The *key* problem was not one of quality but of quantity.
- There are many non-native *pupils* in the class.
- *Pupils'* size changes according to the brightness of light.

<http://esl.fis.edu/teachers/support/vocabPoly.htm>

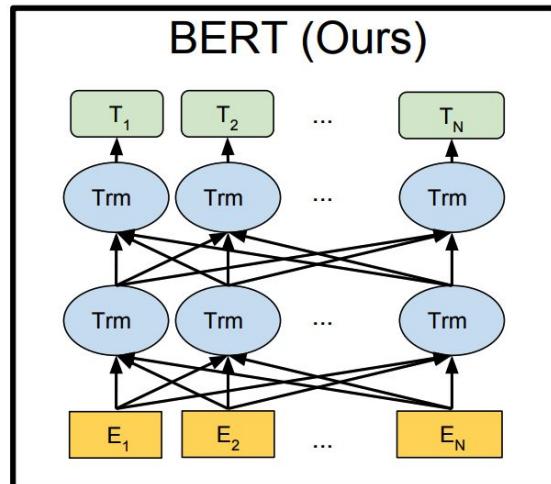
Contextual language models

- The **word vector** is based on the **context** of the word
- Previous contextual language models were **unidirectional** (ULMFit, GPT, OpenAI GPT) or shallow bidirectional model (ELMo).



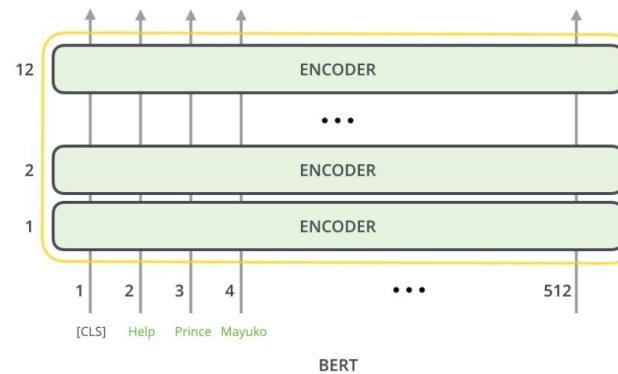
Contextual language models: BERT

- BERT is fully bidirectional
- Based on **attention mechanisms** to gather information about the **relevant context of a given word**
- BERT has become a state-of-the-art in many NLP tasks



BERT architecture

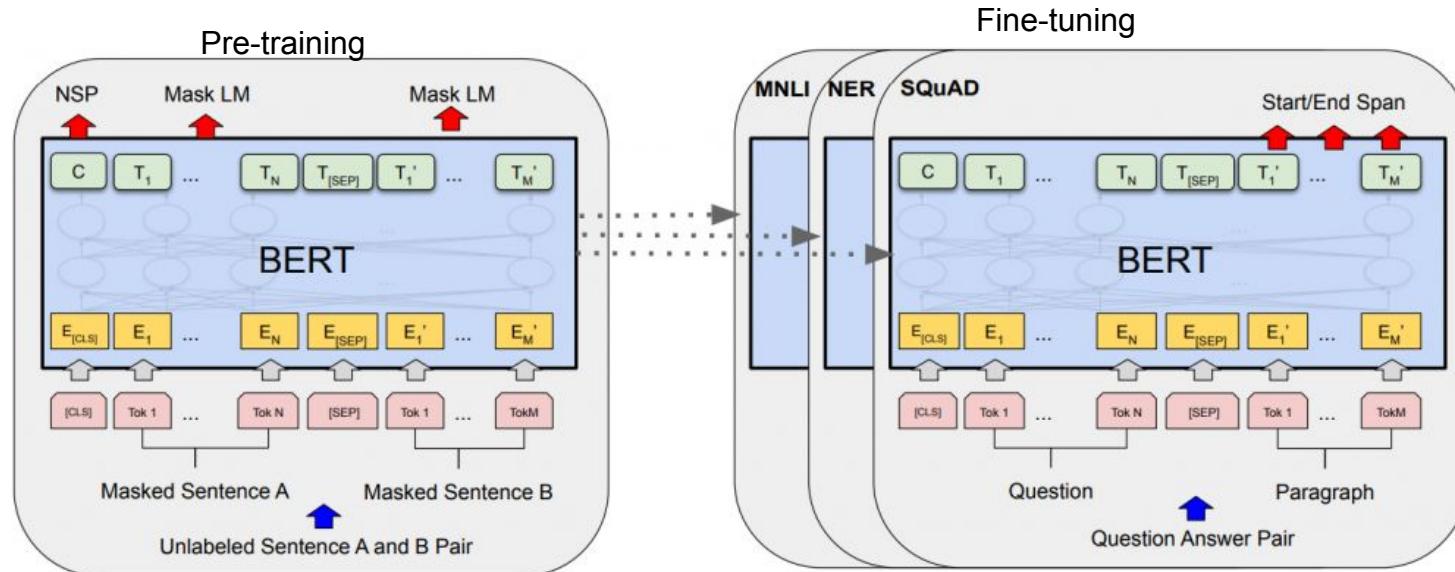
- A stack of encoders
- A vector is generated for each input token.



source: <http://alammar.github.io/illustrated-bert/>

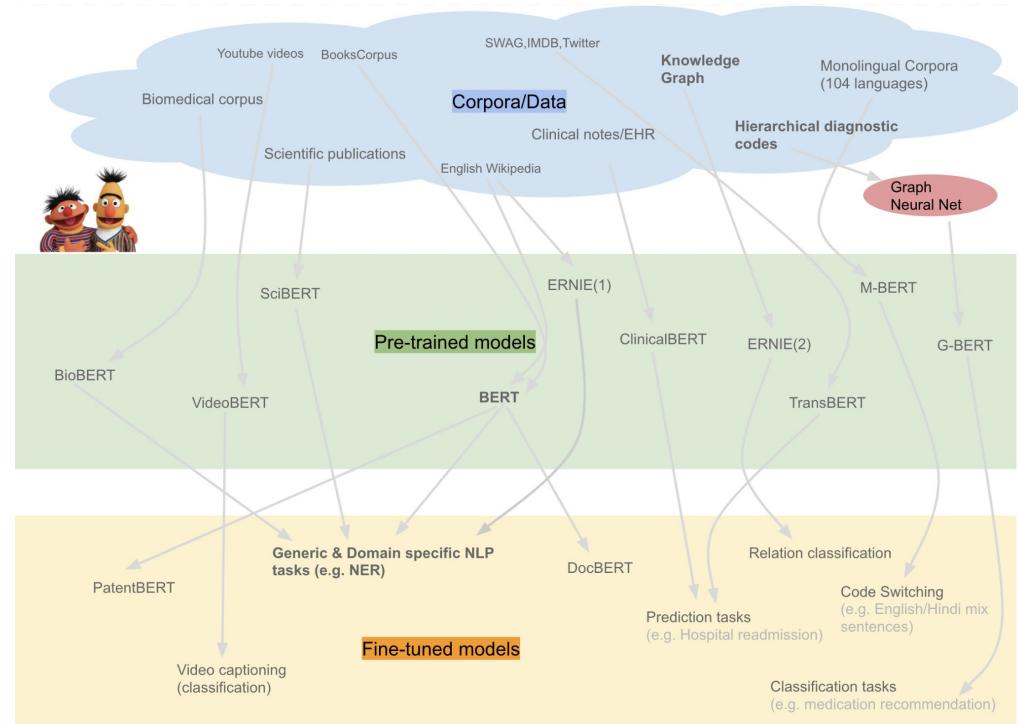
BERT

- These vectors can be used in many different NLP tasks: text classification, text summarization, NER, etc.



Training BERT

BERT was trained on Wikipedia (2,500 million words) and BookCorpus (800 million words)

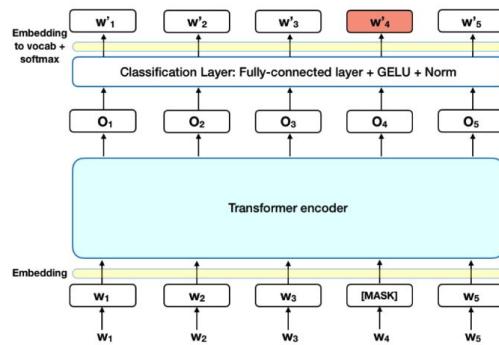


Pre-Training BERT

- To train BERT, two approaches are used simultaneously:
 - Masked Language Model (MLM)
 - Next Sentence Prediction (NSP)

Masked Language Model (MLM)

- The model is fed with a sentence such that 15% of the words in the sentence are masked.
- Then, BERT has to predict the masked words correctly given the context of unmasked words.



Input: The $[MASK]_1$ is not working. It's unable to $[MASK]_2$

Labels: $[MASK]_1$ = computer; $[MASK]_2$ = start.

Next Sentence Prediction (NSP)

- The model is fed with 2 sentences.
- The model has to predict the order of the 2 sentences.

Sentence A = The computer is not working.

Sentence B = It's unable to start.

Label = IsNextSentence

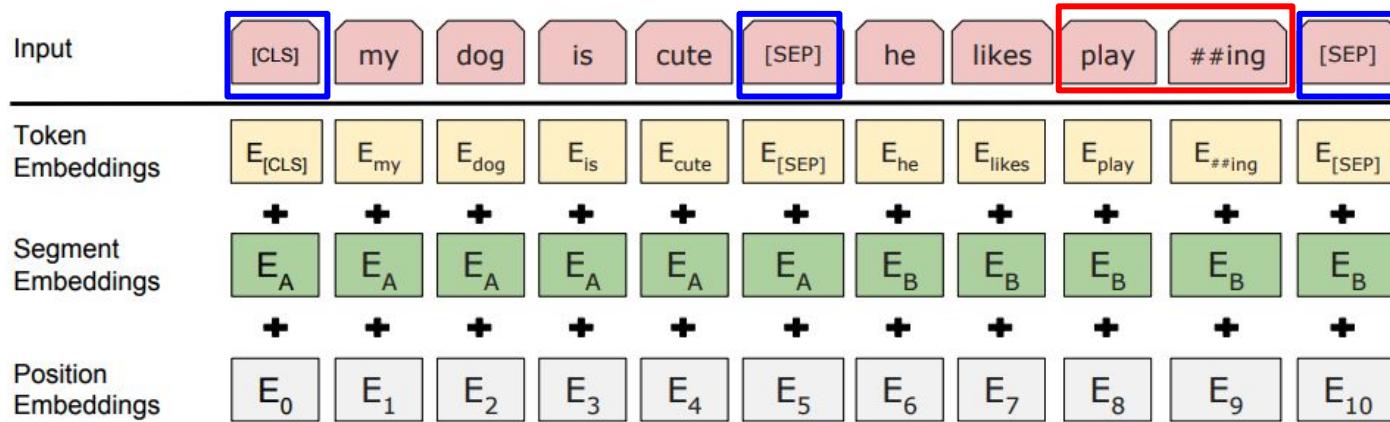
Sentence A = The computer is not working.

Sentence B = Coffee is very tasty.

Label = NotNextSentence

BERT input

[CLS] marks the beginning of an input sentence, [SEP] marks the separation/end of sentences.



Input: "My dog is cute. He likes playing" => ['[CLS]', 'My', 'dog', 'is', 'cute' [SEP], 'He', 'likes', 'play', '# #ing', '[SEP]']

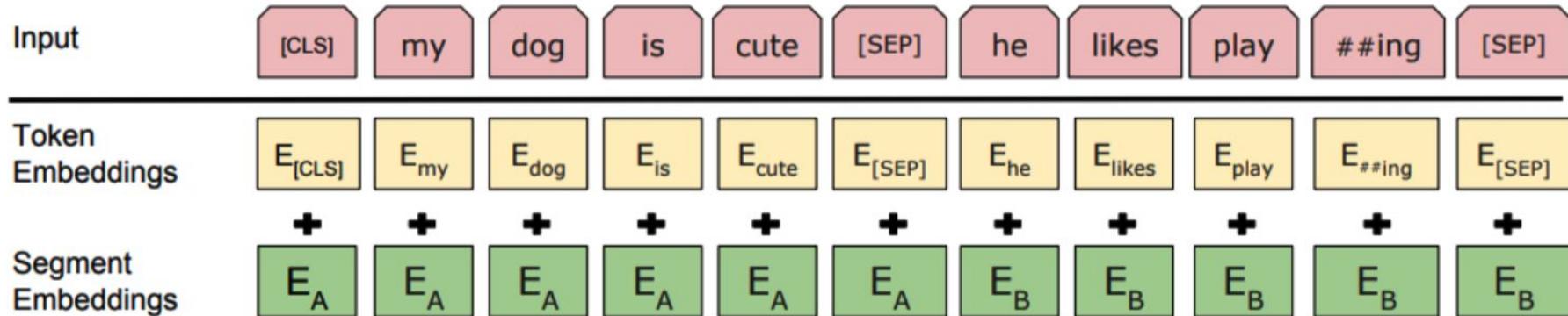
BERT input - token embeddings

- generates the token embedding (a numerical vector for each word in the input sentence).

Input	[CLS]	my	dog	is	cute	[SEP]	he	likes	play	##ing	[SEP]
Token Embeddings	$E_{[CLS]}$	E_{my}	E_{dog}	E_{is}	E_{cute}	$E_{[SEP]}$	E_{he}	E_{likes}	E_{play}	$E_{##ing}$	$E_{[SEP]}$

BERT input - segment embeddings

- help to distinguish between the different sentences in a single input.
- For the input:
[‘[CLS]’, ‘my’, dog, ‘is’, ‘cute’, ‘[SEP]’, ‘he’, ‘likes’, ‘play’, ‘##ing’, ‘[SEP]’],
the segment embeddings will be:
[0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1]

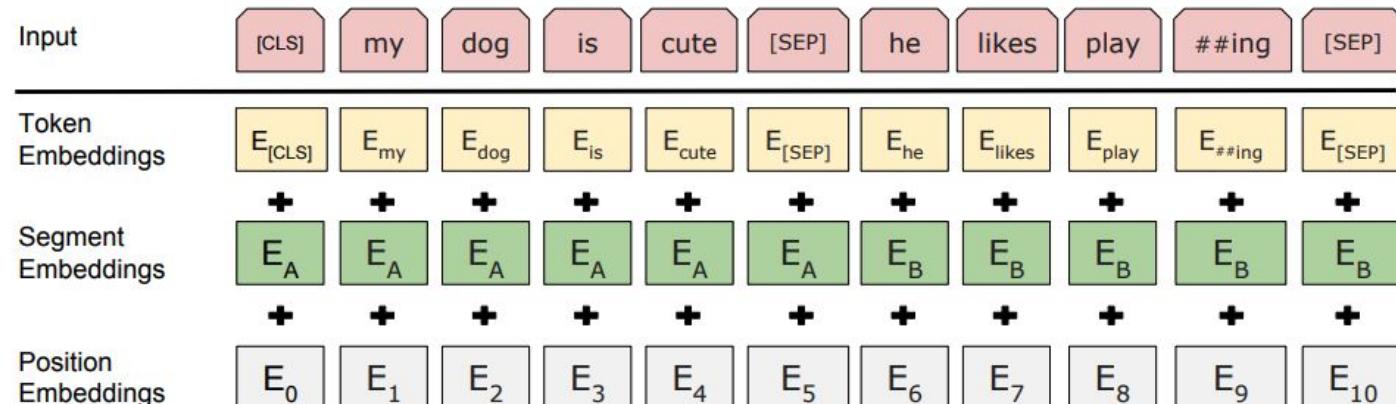


BERT input - mask embeddings

- The input size for BERT is **512**.
- We must add padding of size $512 - \text{len}(\text{our input})$ at the end.
- We also generate a mask token of size **512** in which the index corresponding to the relevant words will have **1s** and the index corresponding to padding will have **0s**

BERT, position embeddings

- Position Embeddings generated internally in BERT and that provide the input data a sense of order.



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

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

Different versions of BERT

- BERT base: 12 encoders, 12 head attention heads. It has an output size of 768 dimensions.
- BERT large: 24 encoders, 16 head attentions. It has an output size of 1024 dimensions.

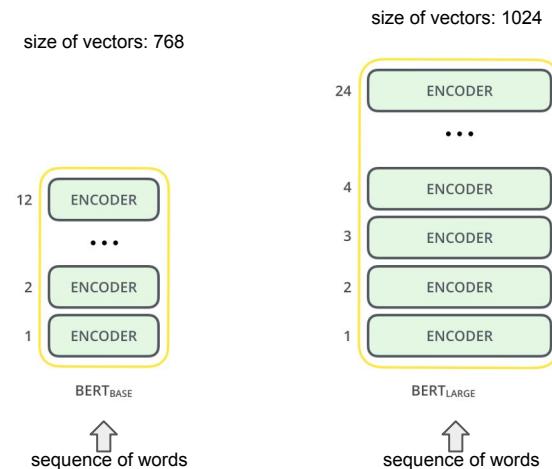
Download from:

<https://huggingface.co/>

<https://github.com/google-research/bert>

Spanish version of BERT:

<https://github.com/dccuchile/beto>

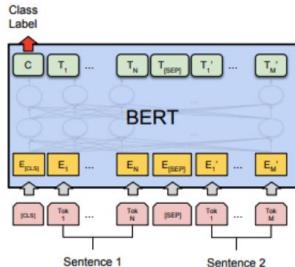


Fine-tuning for a specific task

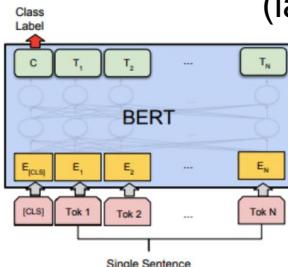
- As we feed input data, the entire pre-trained BERT model and **the additional untrained classification layer** is trained on our specific task.
- For example, if the task is text classification, we can add a linear softmax layer on top of the BERT for predicting the class label of the input text.
- If our task is NER, we can add a CRF layer on the top of BERT.
- Training the model is relatively inexpensive: the bottom layers (BERT) have already great words representation, and we only really need to train the last layer .

Fine-tuning for a specific task

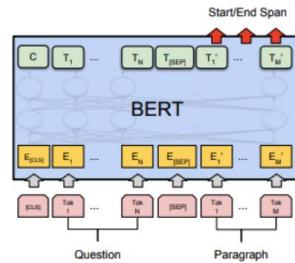
Text similarity



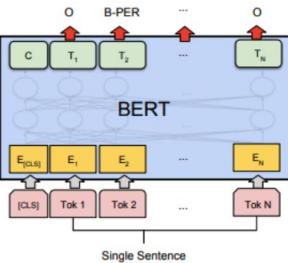
(a) Sentence Pair Classification Tasks:
MNLI, QQP, QNLI, STS-B, MRPC,
RTE, SWAG



(b) Single Sentence Classification Tasks:
SST-2, CoLA



(c) Question Answering Tasks:
SQuAD v1.1



(d) Single Sentence Tagging Tasks:
CoNLL-2003 NER

QA
Text summarization
Machine Translation

output dimension=768 (base), 1024 (large)

Text classification: spam detection, fake news detection, hate speech detection, sentiment analysis,...

Sequence labelling tasks: NER, PoS tagging

Later models based on BERT

- ROBERTa
- Distilbert
- XLM/mBERT
- ALBERT

ROBERTa

- Developed by Facebook,
- Based on BERT.
- Do not use Next Sentence Prediction (NSP)
- Use dynamic amskin during the training epochs.

DistilBERT

- Developed by HuggingFace,
- Learns a distilled (approximate) version of BERT, retaining 95% performance but using only half the number of parameters
- The concept is that once a large neural network has been trained, its full output distributions can be approximated using a smaller network.

XLM/mBERT

- Developed by Facebook,
- XLM a dual-language training mechanism with BERT in order to learn relations between words in different languages.
- The model outperforms other models in a multi-lingual classification task
- Significantly improves machine translation if a pre-trained model is used for the initialization of the translation model.

ALBERT

- Developed by Google Research and Toyota Technological Institute,,
- Based on BERT.
- It is much smaller and lighter and smarter than BERT.
- Improves parameter efficiency by sharing all parameters (feed forward and attention parameters), across all layers.
- Replace NSP with **Sentence-Order Prediction (SOP)**
- ALBERT represents a new state of the art for NLP on several benchmarks and a new state of the art for parameter efficiency.

Conclusions

- Deep learning can learn unsupervised features effectively and provides state-of-the-art techniques for many NLP.
- Proper language models are key for developing NLP systems.
- Recurrent Neural Network requires high computation cost to process long sequences.
- Transforms based on attention mechanisms can replace RNN
- BERT is a contextual language model capable to correctly represent polysemy words.
- BERT is fully bidirectional and obtains state-of-the-art results in many NLP tasks.
- Fine-tuning a BERT model for a specific task is relatively inexpensive

Thank you
Question time!!!

isegura@inf.uc3m.es

<https://hulat.inf.uc3m.es/nosotros/miembros/isegura>

<https://github.com/isegura>