

Anatomia degli LLMs

2. Introduzione ai Transformers



Alessio Miaschi

ItaliaNLP Lab, Istituto di Linguistica Computazionale (CNR-ILC), Pisa

alessio.miaschi@ilc.cnr.it

<https://alemmaschi.github.io/>

<http://www.italianlp.it/alessio-miaschi/>

Introduzione ai Transformers

Transformer

- L'architettura ad oggi più utilizzata per lo sviluppo di NLM contestuali è quella del **Transformer**, introdotto per la prima volta qui: [Attention is All you Need \(Vaswani et al., 2017\)](#)

Transformer

- L'architettura ad oggi più utilizzata per lo sviluppo di NLM contestuali è quella del **Transformer**, introdotto per la prima volta qui: [Attention is All you Need \(Vaswani et al., 2017\)](#)
- Il Transformer è una rete neurale (Encoder-Decoder) che sfrutta uno specifico meccanismo, l'**Attention**, per concentrarsi solo su alcune porzioni fondamentali di una frase e creare rappresentazioni contestuali delle parole

Transformer

- L'architettura ad oggi più utilizzata per lo sviluppo di NLM contestuali è quella del **Transformer**, introdotto per la prima volta qui: [Attention is All you Need \(Vaswani et al., 2017\)](#)
- Il Transformer è una rete neurale (Encoder-Decoder) che sfrutta uno specifico meccanismo, l'**Attention**, per concentrarsi solo su alcune porzioni fondamentali di una frase e creare rappresentazioni contestuali delle parole

I arrived at the **bank** after crossing thestreet? ...river?
What does **bank** mean in this sentence?



I've no idea: let's wait until I read the end

RNNs

$O(N)$ steps to process a sentence with length N



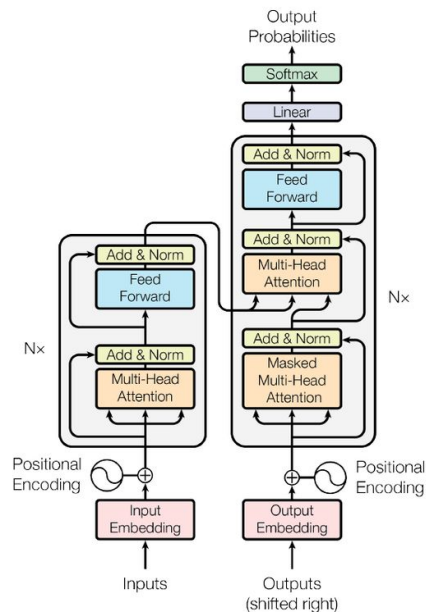
I don't need to wait - I see all words at once!

Transformer

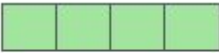
Constant number of steps to process any sentence

Transformer

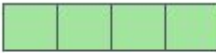
- L'architettura ad oggi più utilizzata per lo sviluppo di NLM contestuali è quella del **Transformer**, introdotto per la prima volta qui: [Attention is All you Need \(Vaswani et al., 2017\)](#)
- Il Transformer è una rete neurale (Encoder-Decoder) che sfrutta uno specifico meccanismo, l'**Attention**, per concentrarsi solo su alcune porzioni fondamentali di una frase e creare rappresentazioni contestuali delle parole



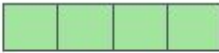
Transformer - Attention

x_1 

Je

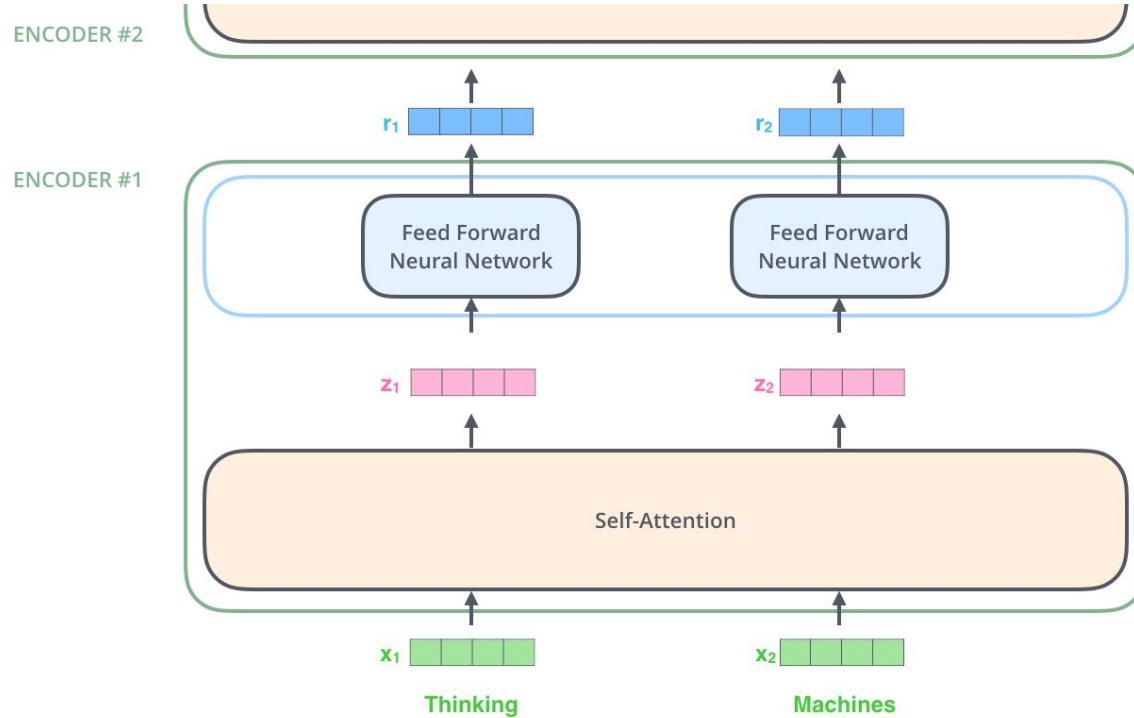
x_2 

suis

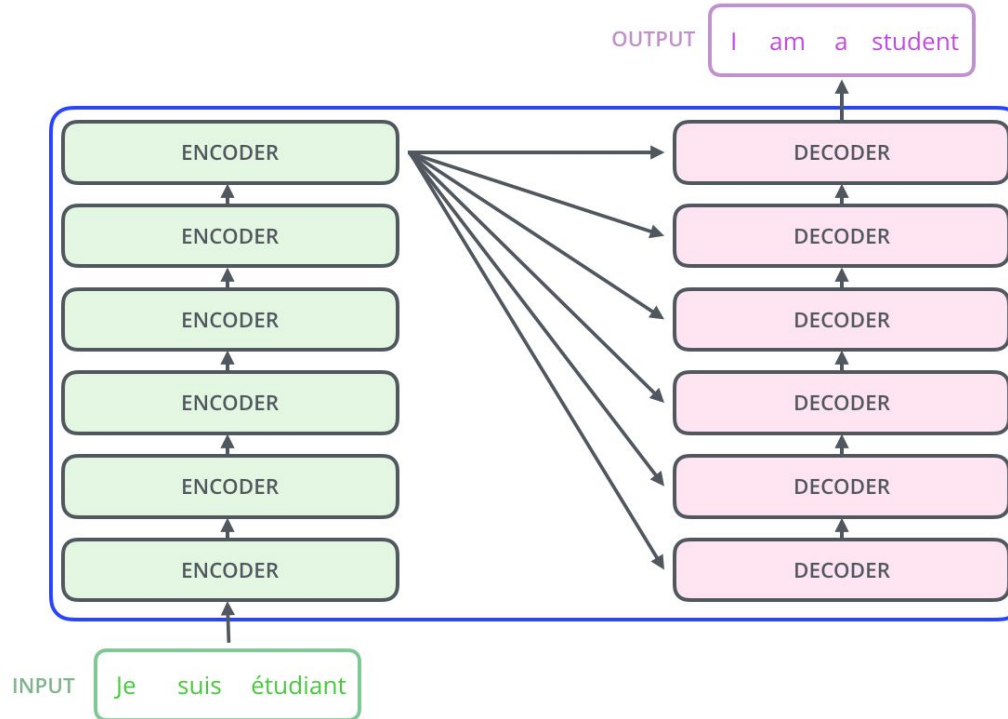
x_3 

étudiant

Transformer - Attention



Transformer - Attention



Transformer - Attention

$$\text{Attention}(q, k, v) = \overbrace{\text{softmax}\left(\frac{qk^T}{\sqrt{d_k}}\right)}^{\text{Attention weights}} v$$

from to

vector dimensionality of K, V

The diagram illustrates the Attention mechanism formula. The function $\text{Attention}(q, k, v)$ takes three inputs: q (blue), k (yellow), and v (red). Below the inputs, a blue arrow points from the word 'from' to q , and a yellow arrow points from the word 'to' to k . The formula shows q and k being multiplied (qk^T) and then divided by the square root of the vector dimensionality of k ($\sqrt{d_k}$). A grey arrow points from the text 'vector dimensionality of K, V' to d_k . The result of the softmax operation is labeled 'Attention weights' with a bracket above it. Finally, this result is multiplied by v .

Transformer - Attention

Each vector receives three representations (“roles”)

$$\begin{bmatrix} W_Q \end{bmatrix} \times \begin{bmatrix} \text{green} \\ \text{green} \\ \text{green} \end{bmatrix} = \begin{bmatrix} \text{blue} \\ \text{blue} \\ \text{blue} \end{bmatrix}$$

Query: vector **from** which the attention is looking

“Hey there, do you have this information?”

$$\begin{bmatrix} W_K \end{bmatrix} \times \begin{bmatrix} \text{green} \\ \text{green} \\ \text{green} \end{bmatrix} = \begin{bmatrix} \text{yellow} \\ \text{yellow} \\ \text{yellow} \end{bmatrix}$$

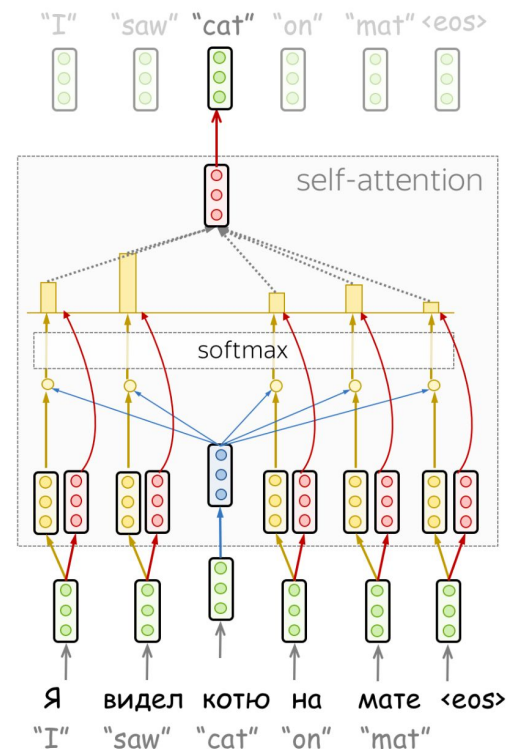
Key: vector **at** which the query looks to compute weights

“Hi, I have this information – give me a large weight!”

$$\begin{bmatrix} W_V \end{bmatrix} \times \begin{bmatrix} \text{green} \\ \text{green} \\ \text{green} \end{bmatrix} = \begin{bmatrix} \text{red} \\ \text{red} \\ \text{red} \end{bmatrix}$$

Value: their weighted sum is attention output

“Here’s the information I have!”



Transformer - Attention

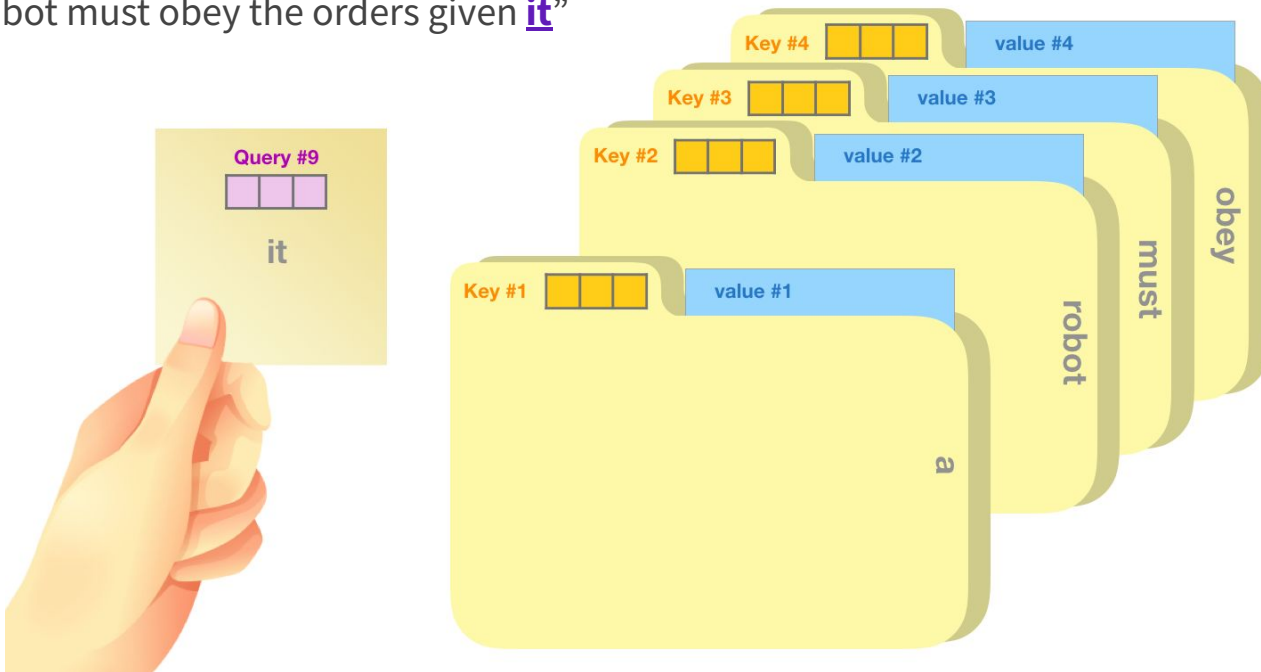
Input: “A robot must obey the orders given it”

Transformer - Attention

Input: “A robot must obey the orders given it”

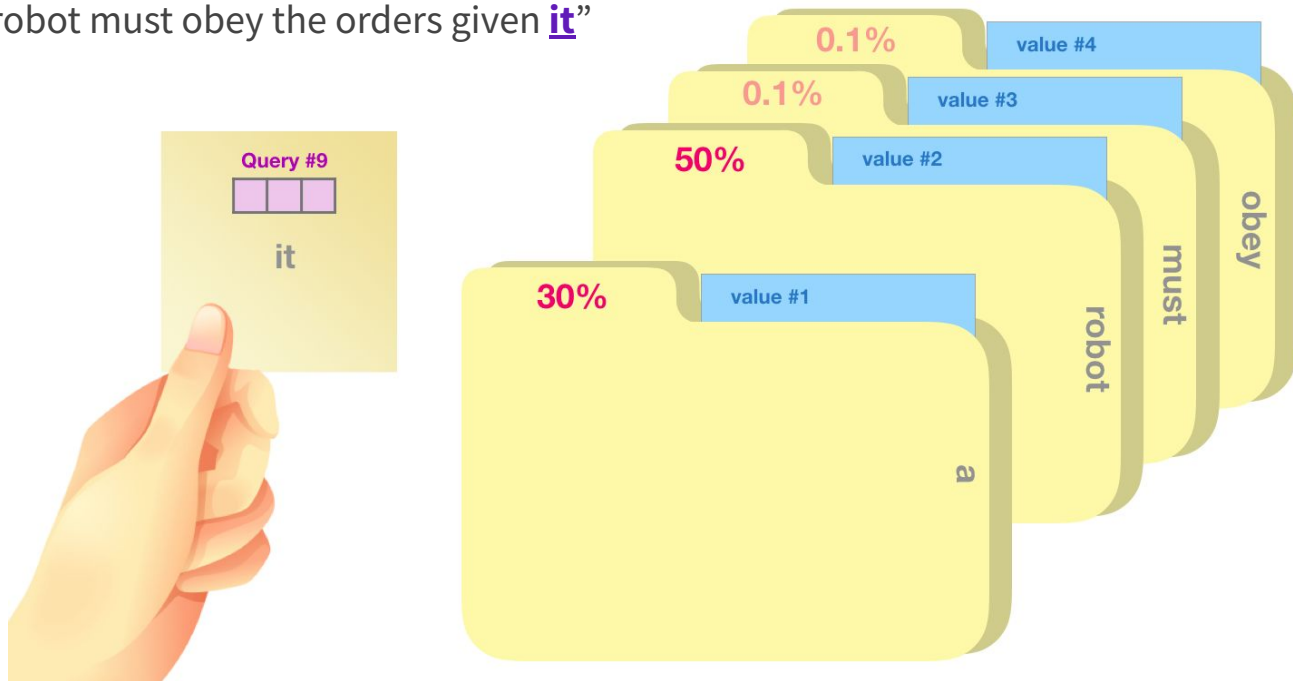
Transformer - Attention

Input: “A robot must obey the orders given **it**”

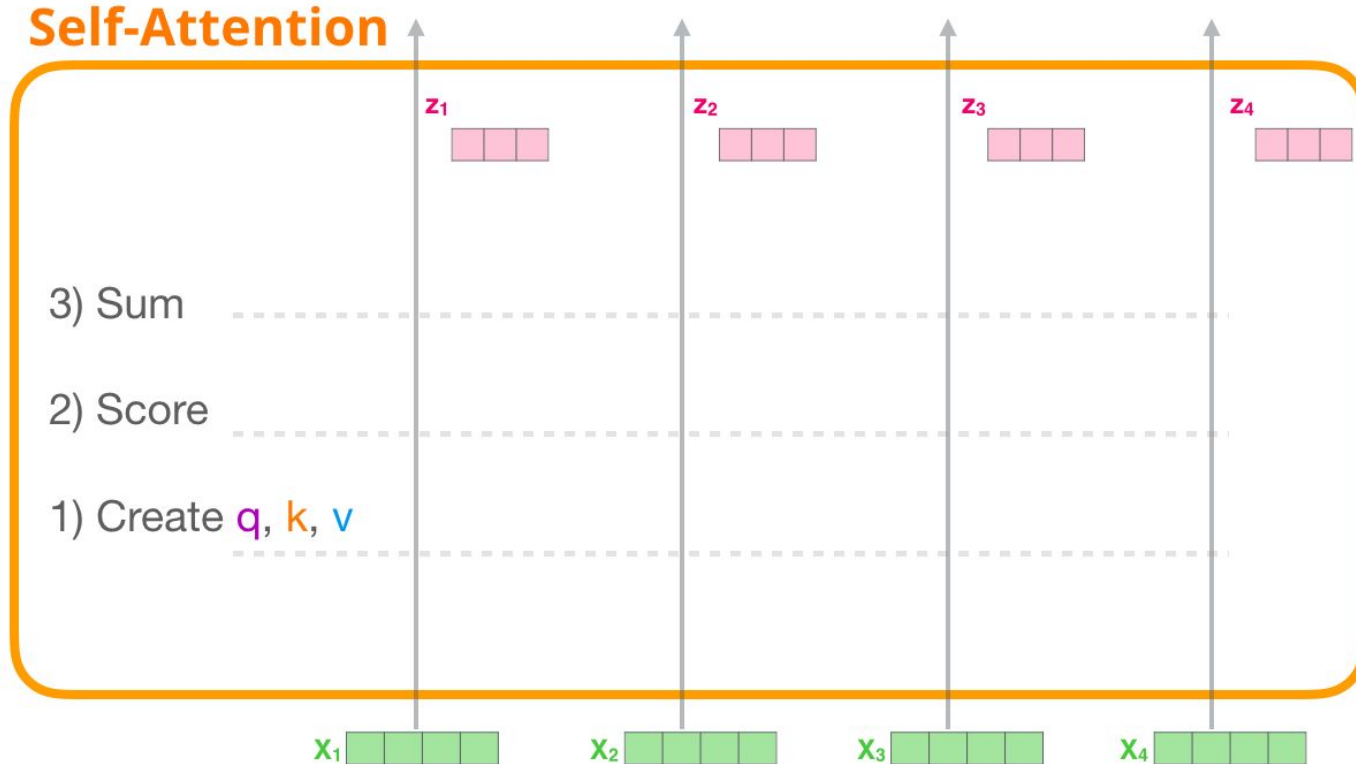


Transformer - Attention

Input: “A robot must obey the orders given **it**”



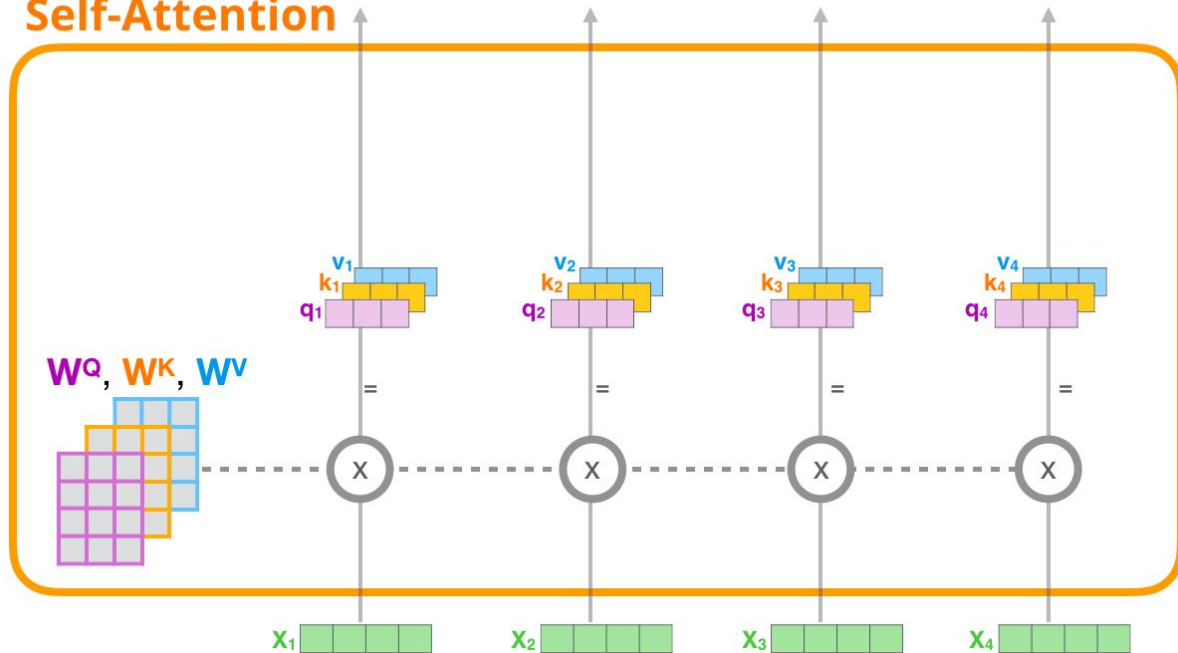
Transformer - Attention



Transformer - Attention

1) For each input token, create a **query vector**, a **key vector**, and a **value vector** by multiplying by weight Matrices W^Q , W^K , W^V

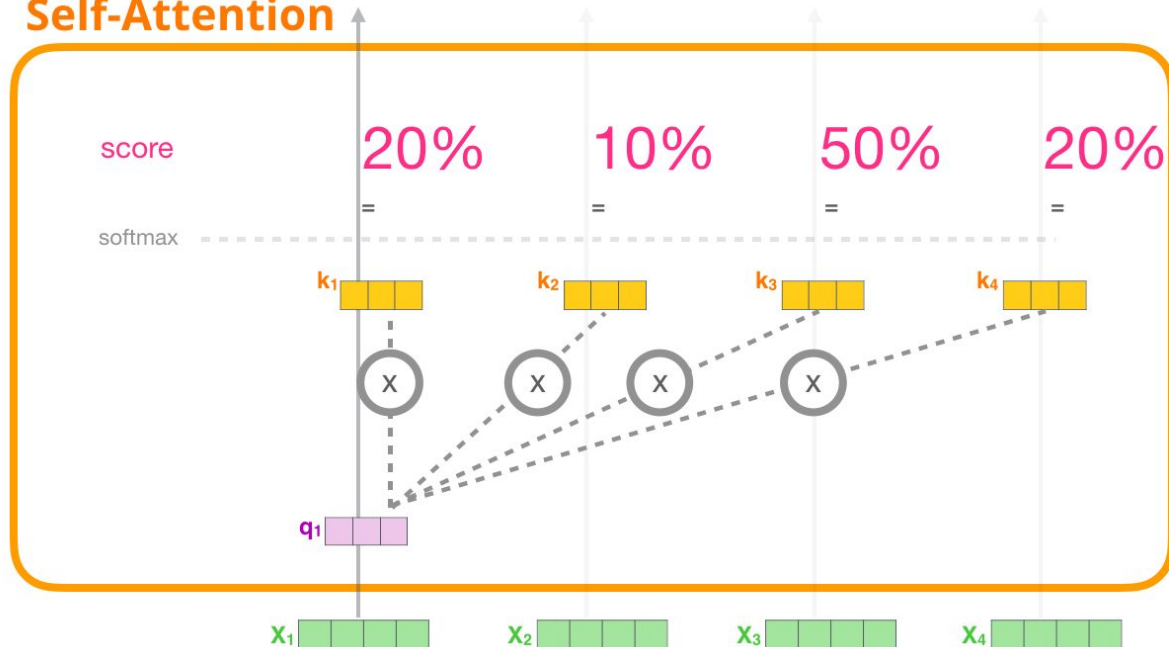
Self-Attention



Transformer - Attention

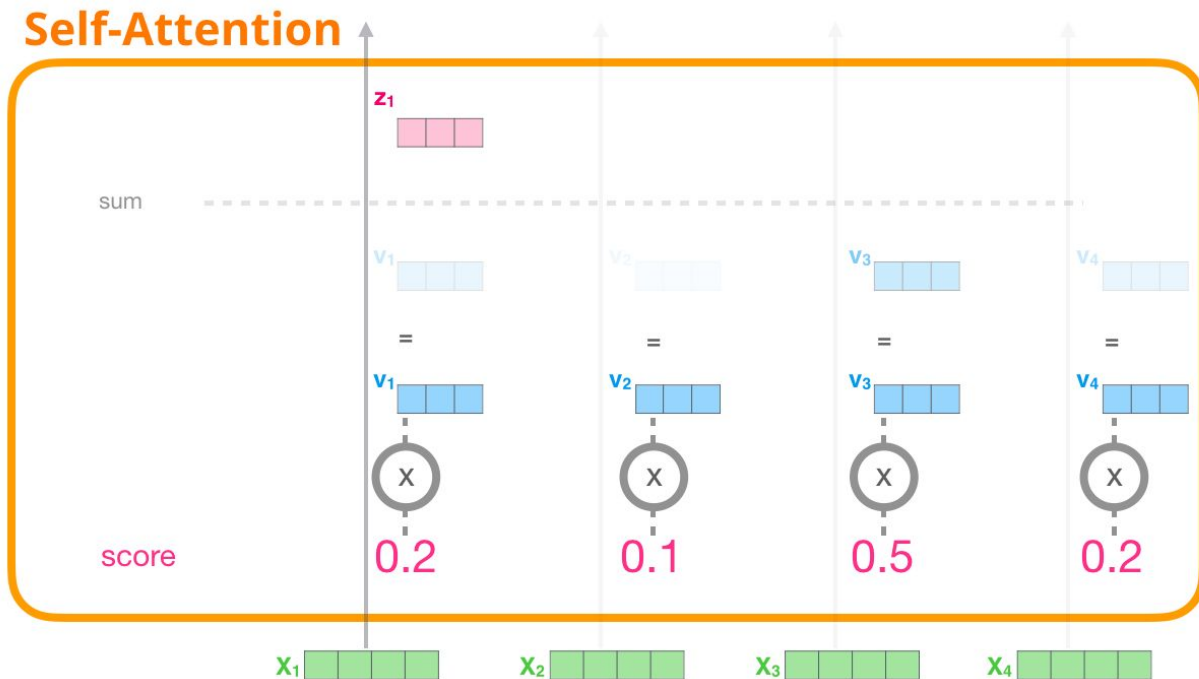
2) Multiply (dot product) the current **query vector**, by all the **key vectors**, to get a score of how well they match

Self-Attention

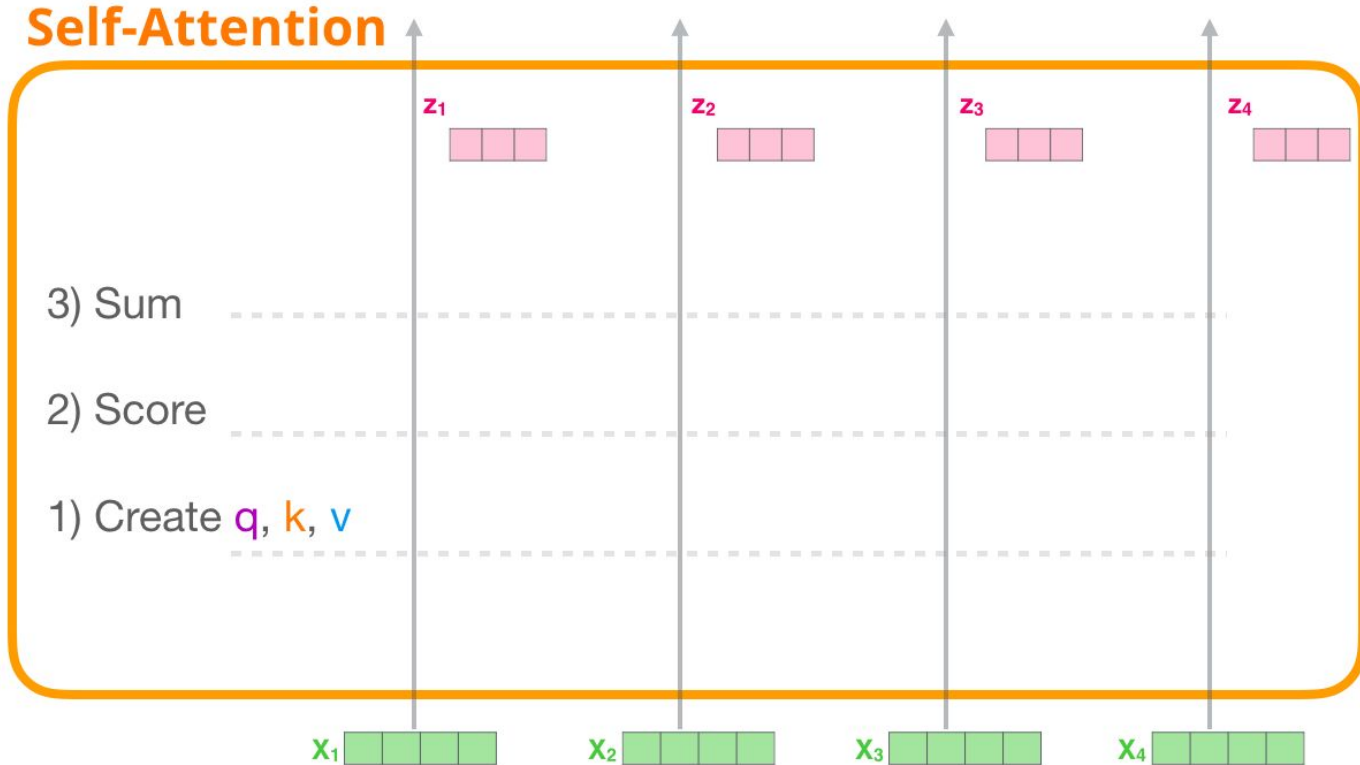


Transformer - Attention

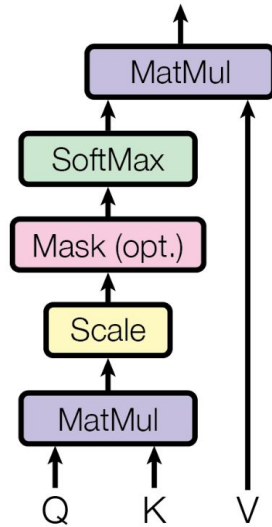
3) Multiply the **value vectors** by the **scores**, then sum up



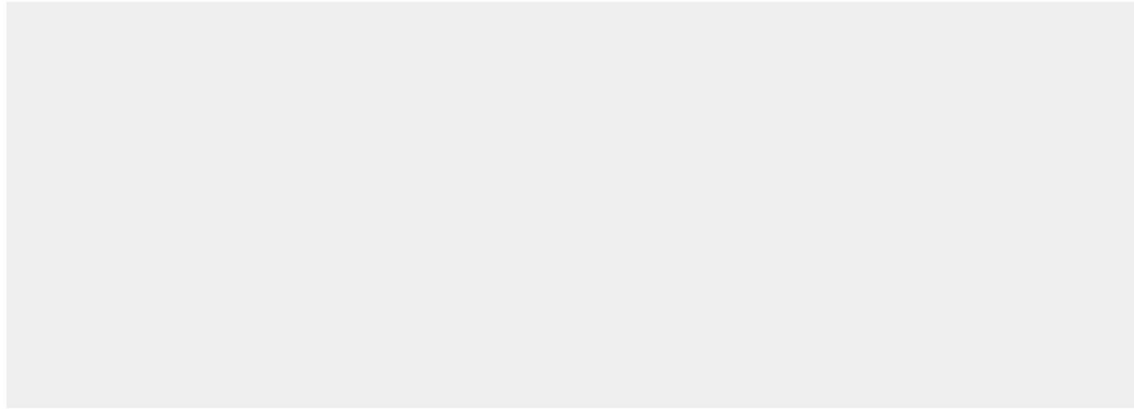
Transformer - Attention



Transformer - Attention



Self-attention



input #1

1	0	1	0
---	---	---	---

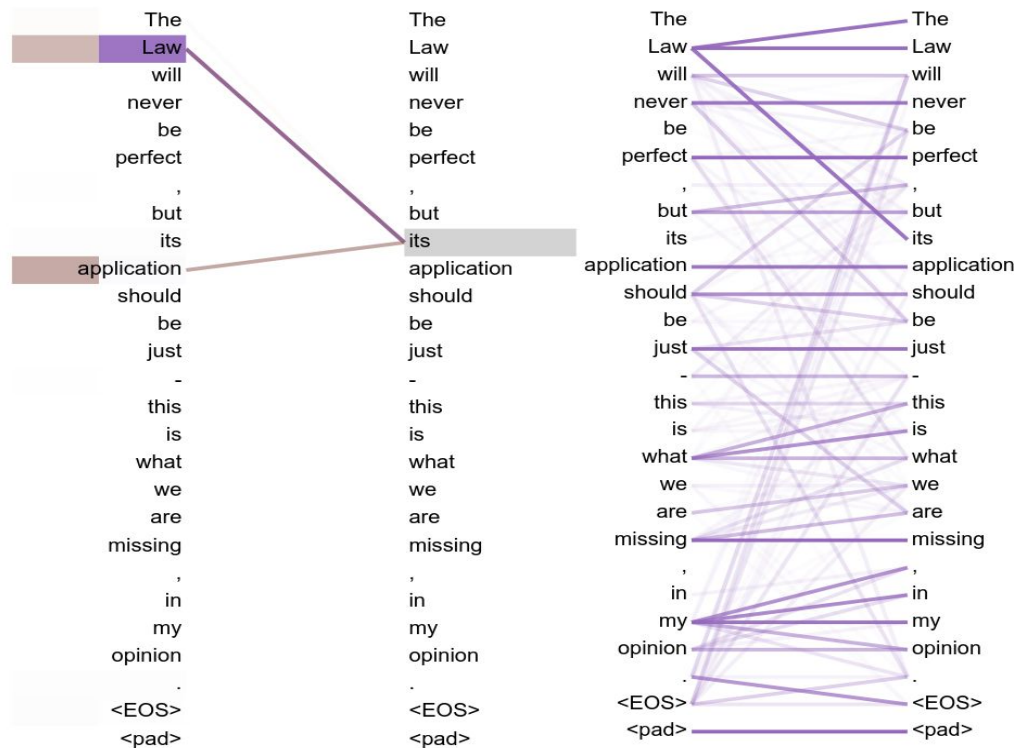
input #2

0	2	0	2
---	---	---	---

input #3

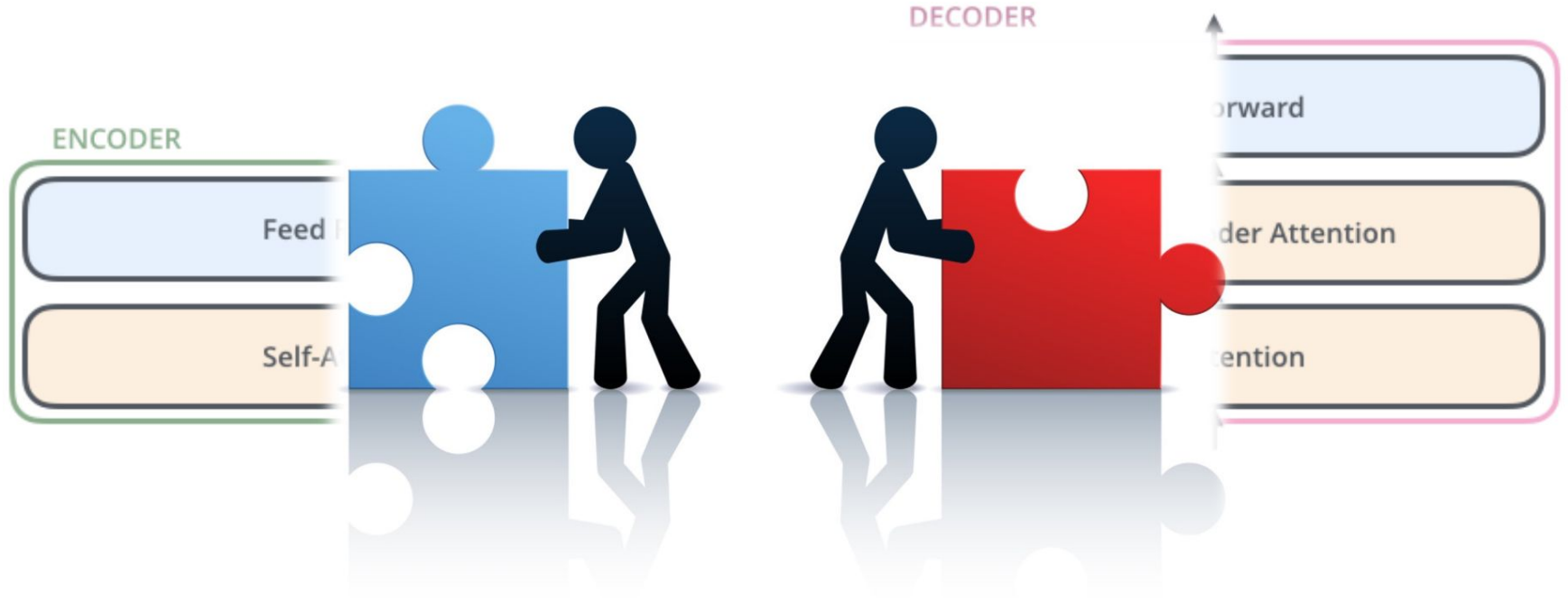
1	1	1	1
---	---	---	---

Transformer - Attention



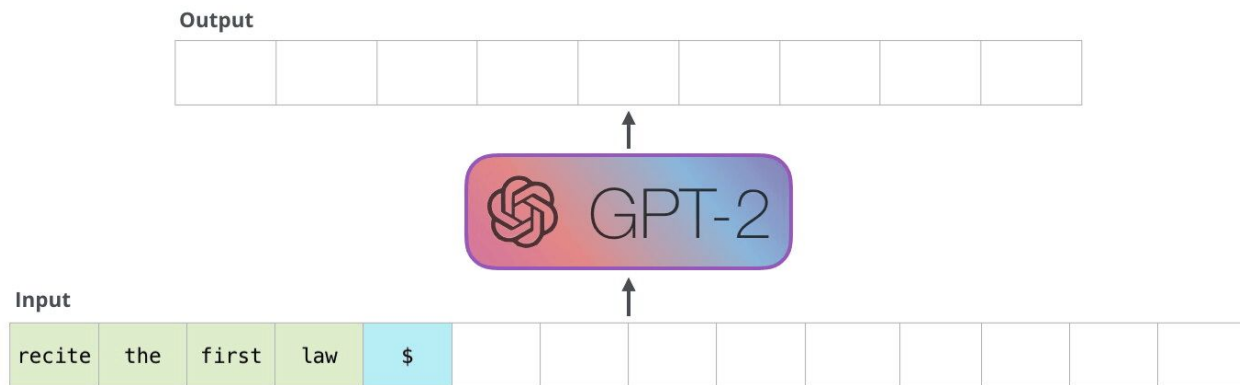
Transformer-based NLMs

Transformer-based NLMs



GPT (Radford et al, 2018), GPT-2 (Radford et al, 2019), etc

- Decoder Transformer model
- Addestrato per approssimare la funzione di **Language Modeling (LM)**
- Modello generativo → pensato per generare testo

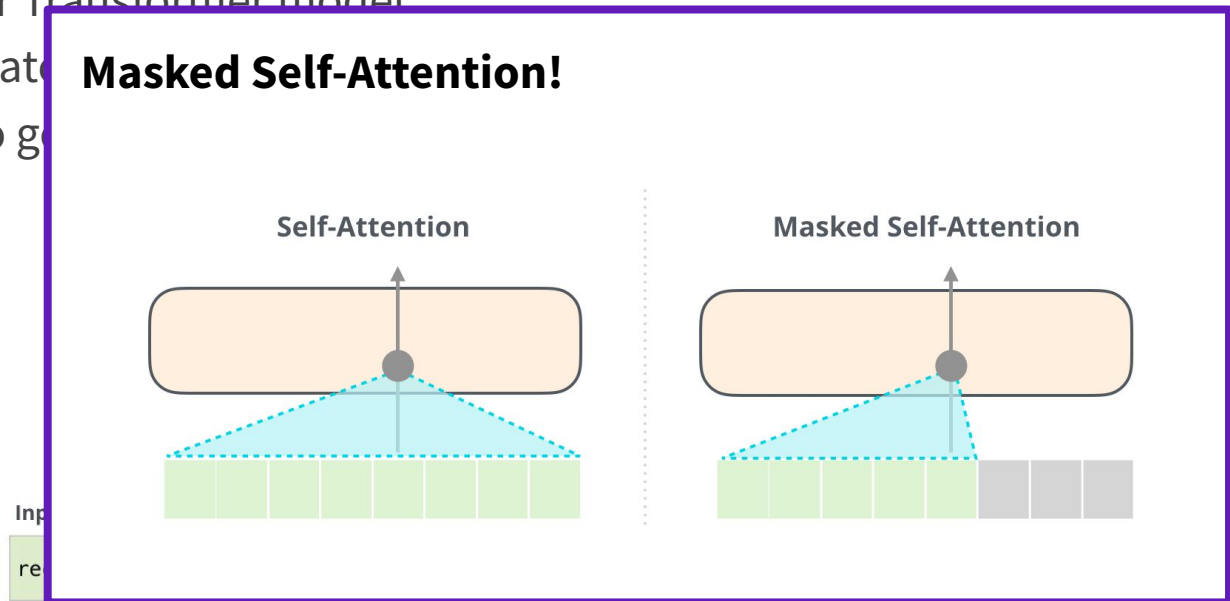


[Improving Language Understanding by Generative Pre-Training \(Radford et al., 2018\), https://openai.com/research/language-unsupervised](https://openai.com/research/language-unsupervised)

[Language Models are Unsupervised Multitask Learners \(Radford et al., 2019\), https://openai.com/research/better-language-models](https://openai.com/research/better-language-models)

GPT (Radford et al, 2018), GPT-2 (Radford et al, 2019), etc

- Decoder Transformer model
- Addestrato
- Modello g



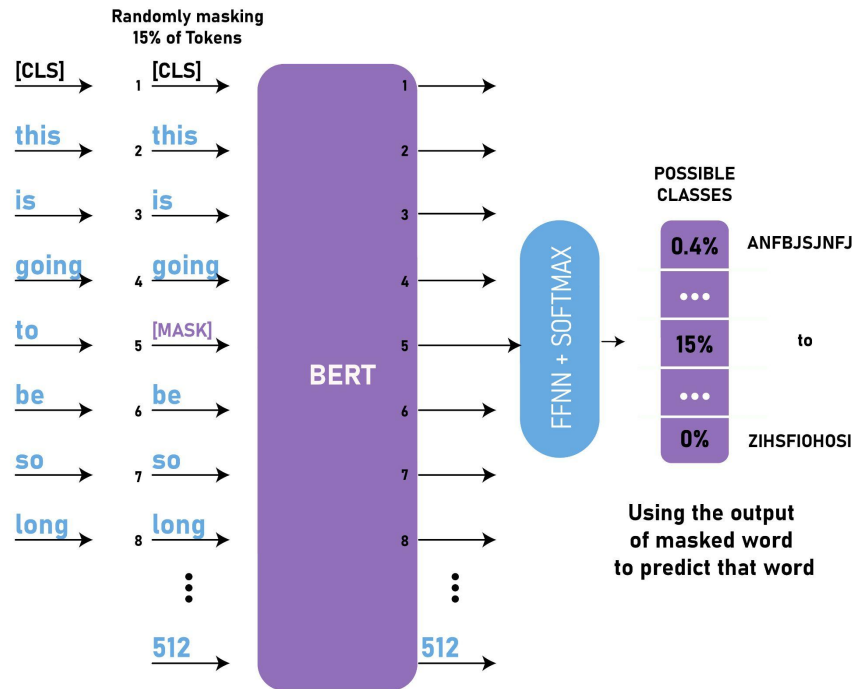
[Improving Language Understanding by Generative Pre-Training \(Radford et al., 2018\), https://openai.com/research/language-unsupervised](https://openai.com/research/language-unsupervised)

[Language Models are Unsupervised Multitask Learners \(Radford et al., 2019\), https://openai.com/research/better-language-models](https://openai.com/research/better-language-models)

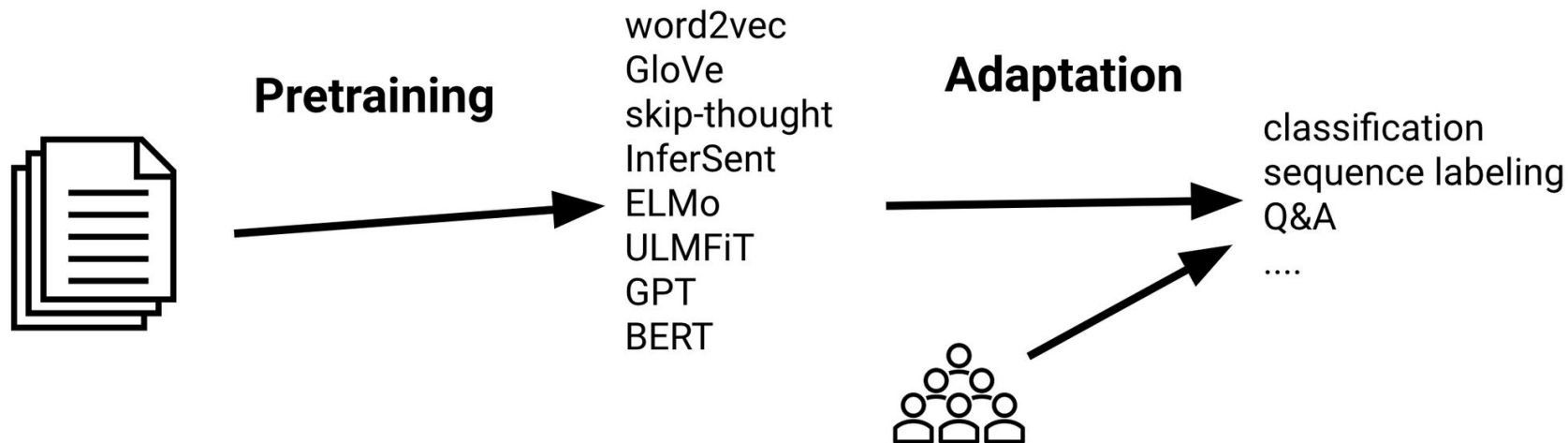
BERT (Devlin et al., 2019)



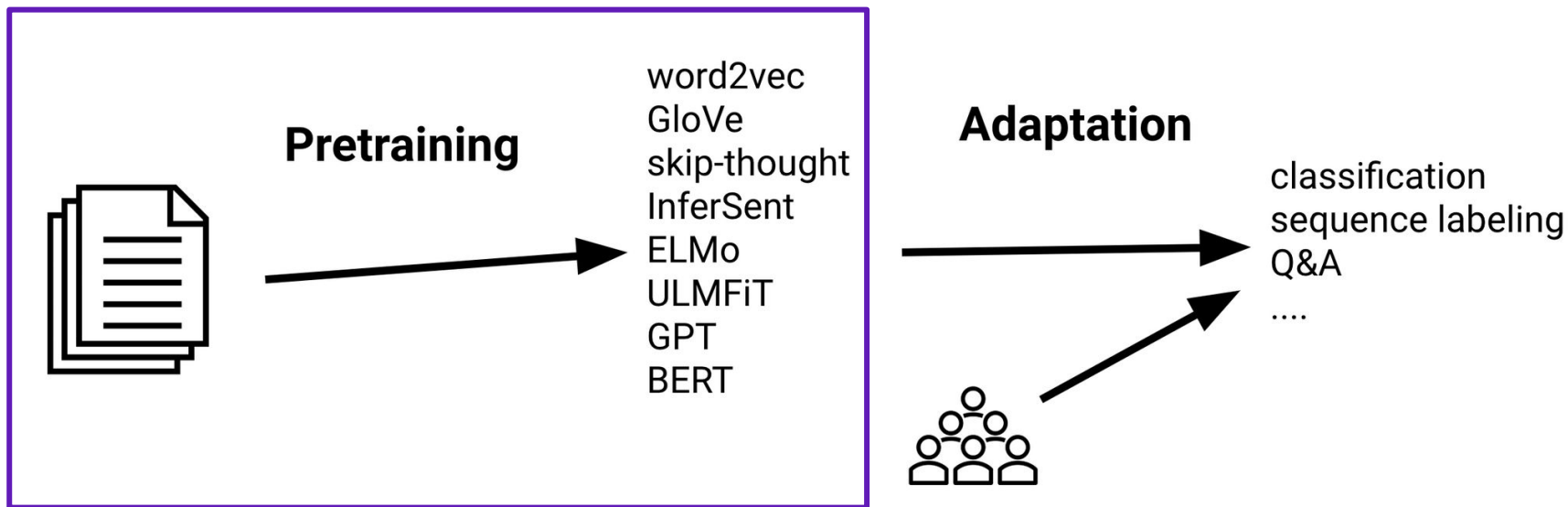
- Encoder Transformer model (12/24 layers)
- Addestrato per approssimare la funzione di **Masked Language Modeling (MLM)**
- Il modello può poi essere ri-addestrato (fine-tuning) per risolvere svariati task di NLP:
 - Sentiment analysis;
 - Question answering;
 - Textual entailment;
 - etc.



Transfer Learning



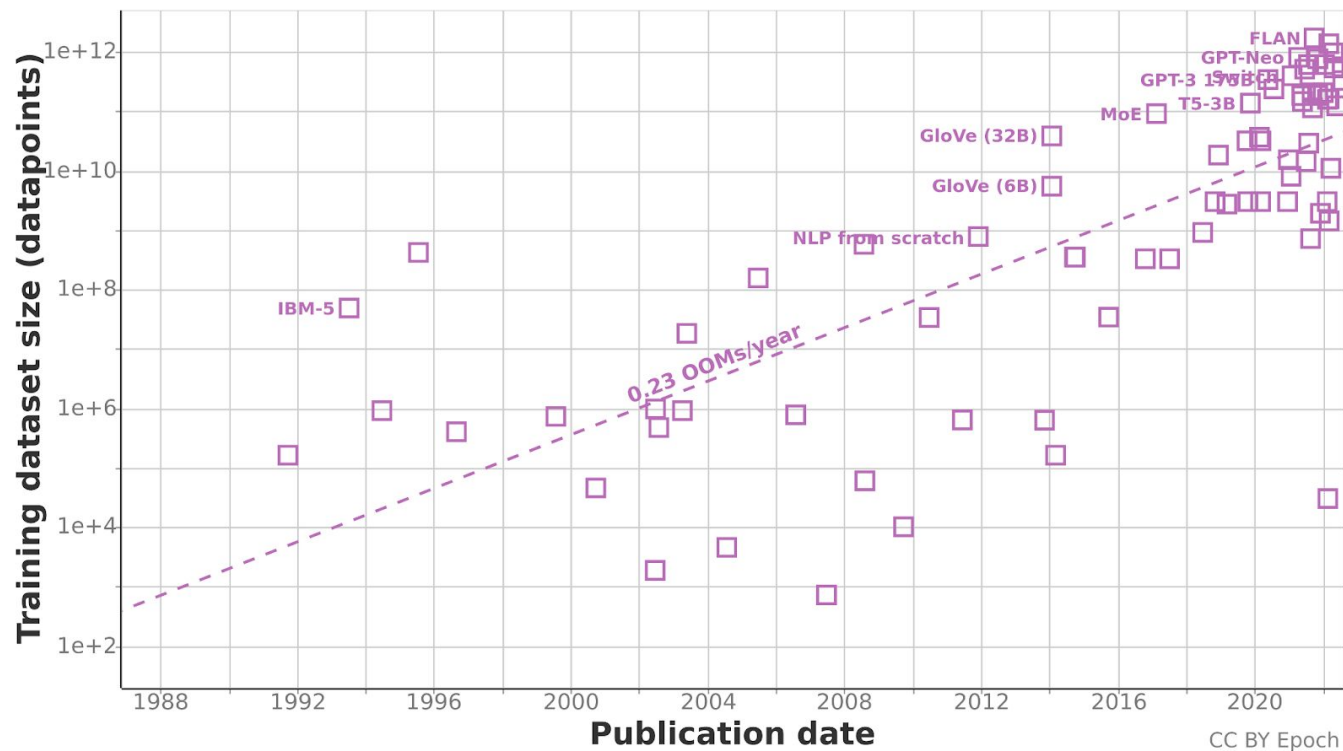
Transfer Learning



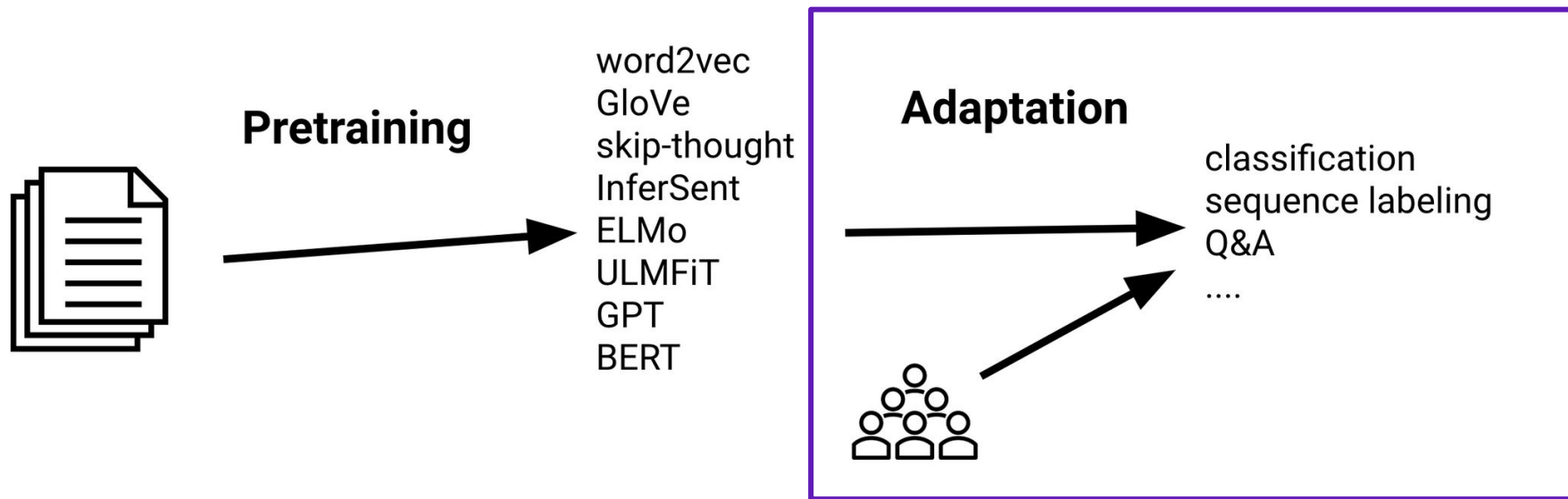
Pre-training

- Durante la fase di “*Pre-training*”, il modello viene addestrato in maniera unsupervised (e.g. LM, MLM) su una grande quantità di dati grezzi
- Alcuni esempi:
 - **Training di BERT**: BookCorpus (800M di parole) e Wikipedia Inglese (2500M di parole)
 - **Training di GPT-3**: CommonCrawl + WebText2 + Books1 + Books2 + Wikipedia (circa 500B di parole)

Pre-training

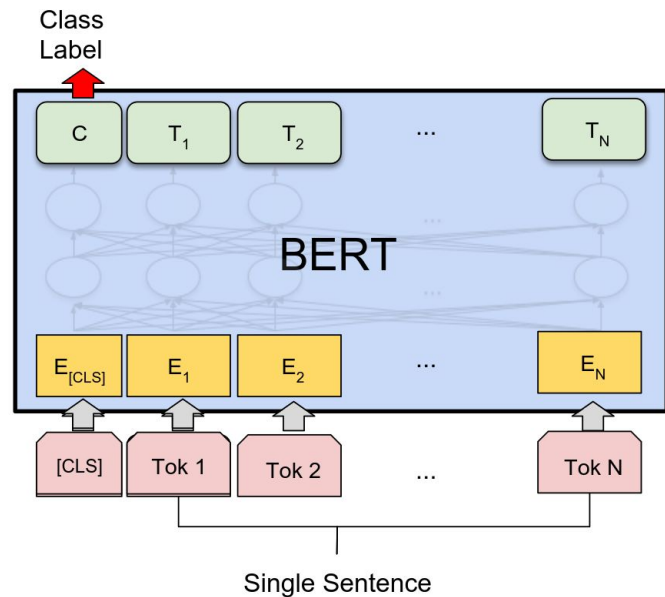


Transfer Learning



Fine-tuning (Adaptation)

- Durante la fase di “*Fine-tuning*”, si va a specializzare il modello su un determinato task
 - In altri termini, si prende il modello pre-trainato e si continua l’addestramento su un nuovo dataset e modificando la sua funzione obiettivo (e.g. sentiment analysis, textual entailment, sentence complexity, etc.)



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

Prompting → Large Language Models (LLMs)

- Negli ultimi anni, lo sviluppo dei NLMs si è spostato verso la creazione di modelli generativi:
 - Scopo principale: considerarsi qualsiasi task (e.g. classificazione, translation, question answering, etc) come task di **generazione**

Prompting → Large Language Models (LLMs)

- Negli ultimi anni, lo sviluppo dei NLMs si è spostato verso la creazione di modelli generativi:
 - Scopo principale: considerarsi qualsiasi task (e.g. classificazione, translation, question answering, etc) come task di **generazione**

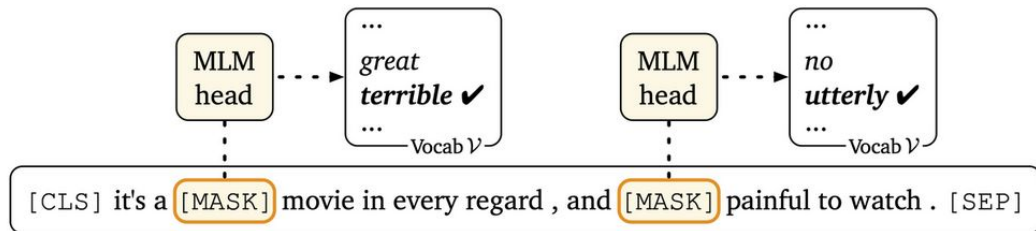
Prompting

“A prompt is a piece of text inserted in the input examples, so that the original task can be formulated as a (masked) language modeling problem.”

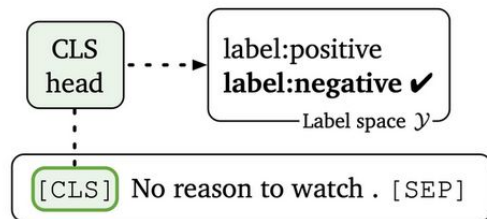
([Prompting: Better Ways of Using Language Models for NLP Tasks, The Gradient](#))

Prompting → Large Language Models (LLMs)

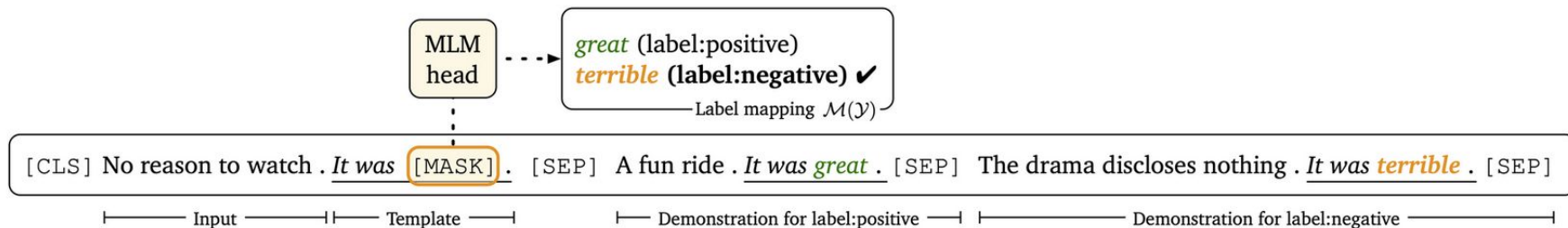
Why Prompts?



(a) MLM pre-training



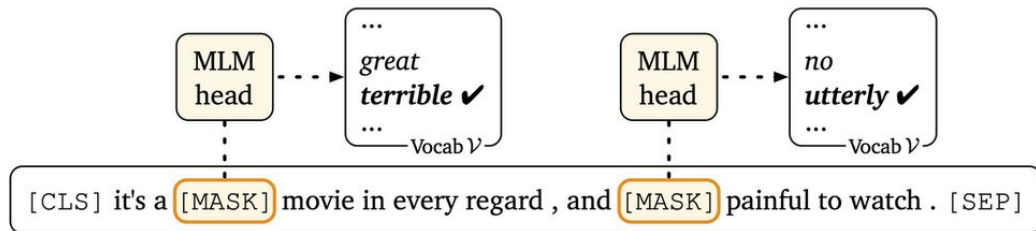
(b) Fine-tuning



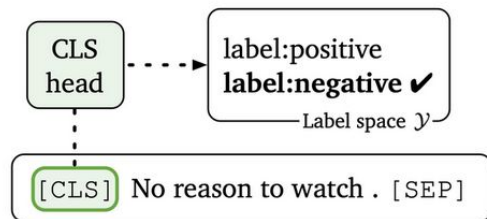
(c) Prompt-based fine-tuning with demonstrations (our approach)

Prompting → Large Language Models (LLMs)

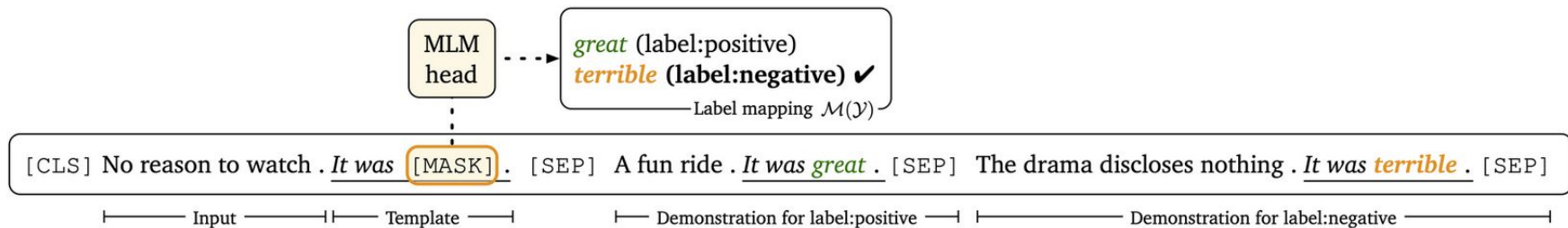
Why Prompts?



(a) MLM pre-training



(b) Fine-tuning



(c) Prompt-based fine-tuning with demonstrations (our approach)

T5 (Raffel et al., 2020)

- Encoder-Decoder Transformer model
- Pensato con l'intenzione di riformulare tutti i task di NLP in un formato di tipo “text-to-text”, dove input e output sono quindi sempre stringhe di testo

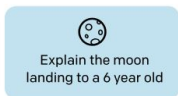


Instruction Tuning e RLHF: da GPT-3 a InstructGPT

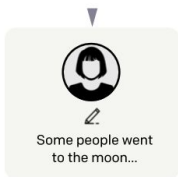
Step 1

Collect demonstration data, and train a supervised policy.

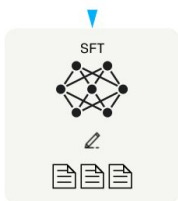
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



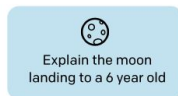
This data is used to fine-tune GPT-3 with supervised learning.



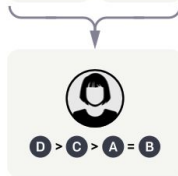
Step 2

Collect comparison data, and train a reward model.

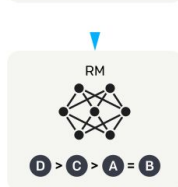
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



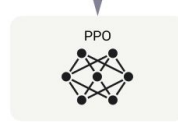
Step 3

Optimize a policy against the reward model using reinforcement learning.

A new prompt is sampled from the dataset.



The policy generates an output.

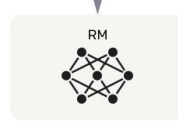


Once upon a time...

The reward model calculates a reward for the output.



The reward is used to update the policy using PPO.

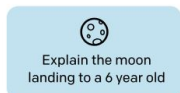


Instruction Tuning e RLHF: da GPT-3 a InstructGPT

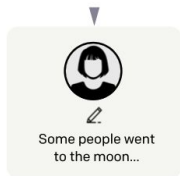
Step 1

Collect demonstration data, and train a supervised policy.

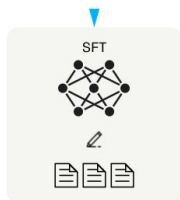
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



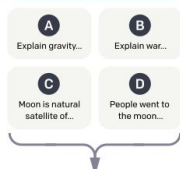
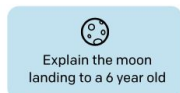
This data is used to fine-tune GPT-3 with supervised learning.



Step 2

Collect comparison data, and train a reward model.

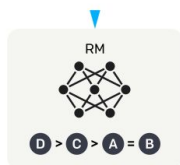
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



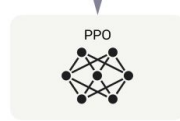
Step 3

Optimize a policy against the reward model using reinforcement learning.

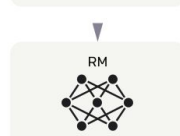
A new prompt is sampled from the dataset.



The policy generates an output.



The reward model calculates a reward for the output.



The reward is used to update the policy using PPO.

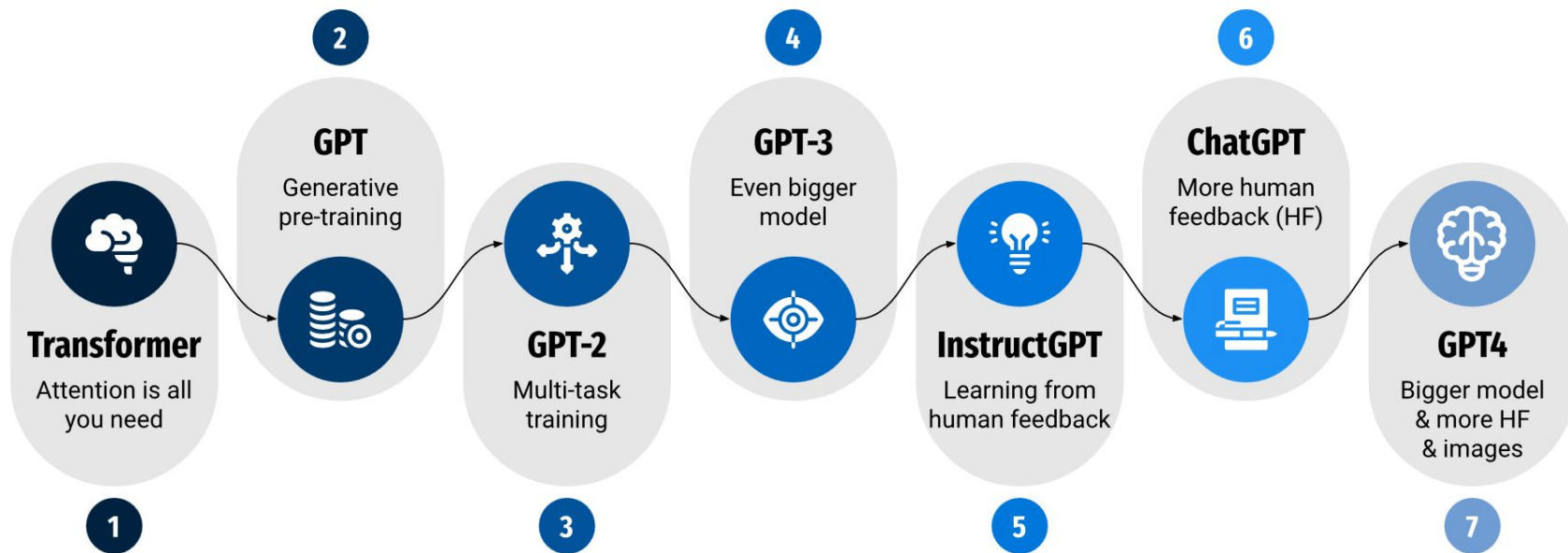


Reinforcement Learning from Human Feedback (RLHF)

<https://huggingface.co/blog/rlhf>

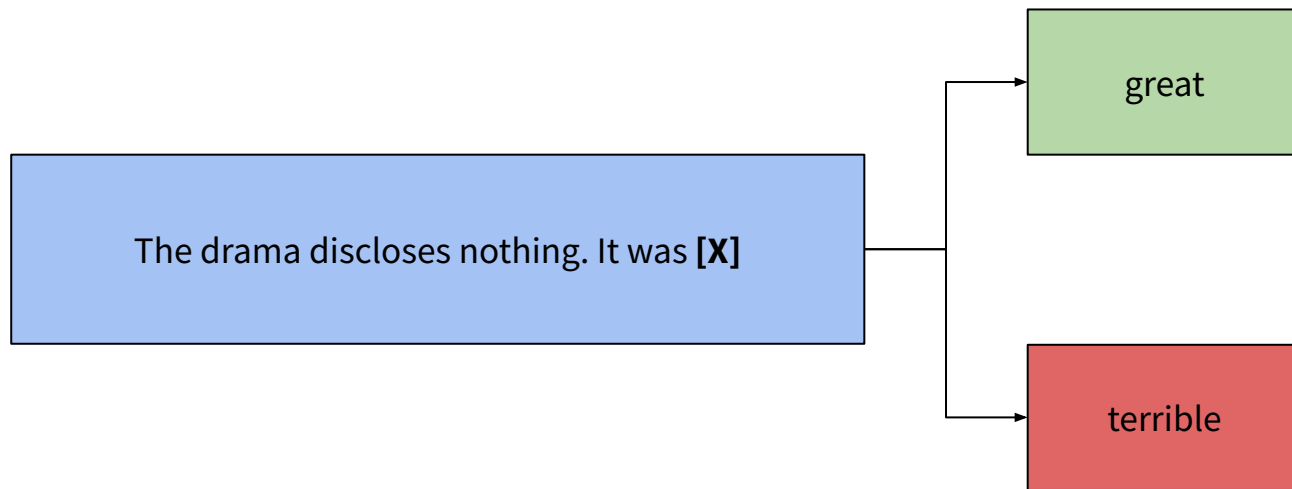
Dal Transformer a GPT4

Evolution from Transformer architecture to ChatGPT



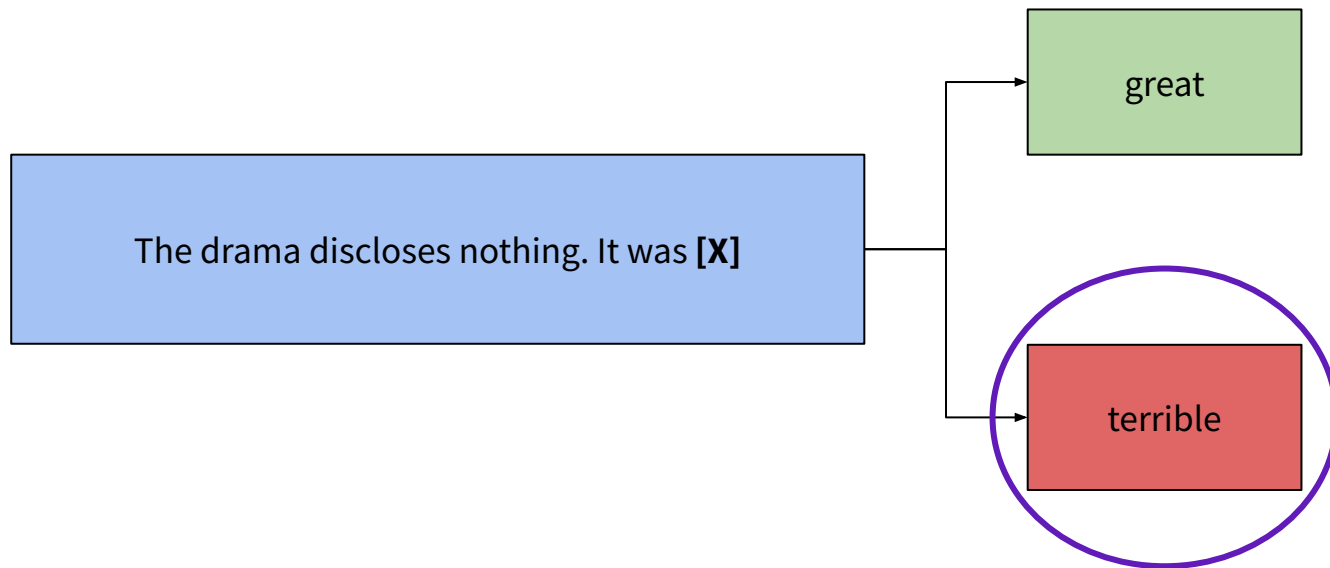
Large Language Models (LLMs)

Zero-Shot Text Classification



Large Language Models (LLMs)

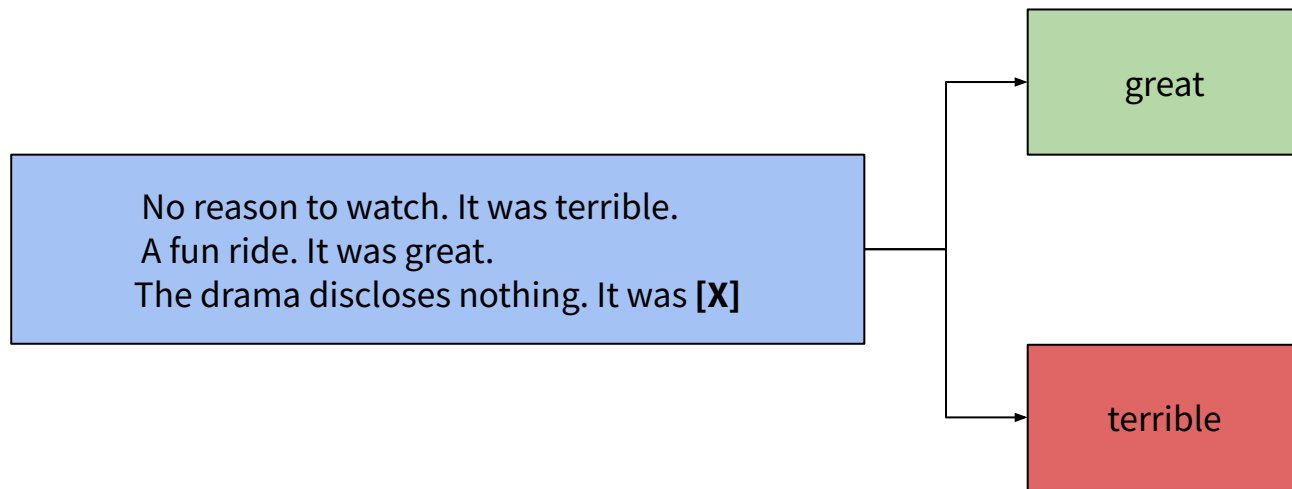
Zero-Shot Text Classification



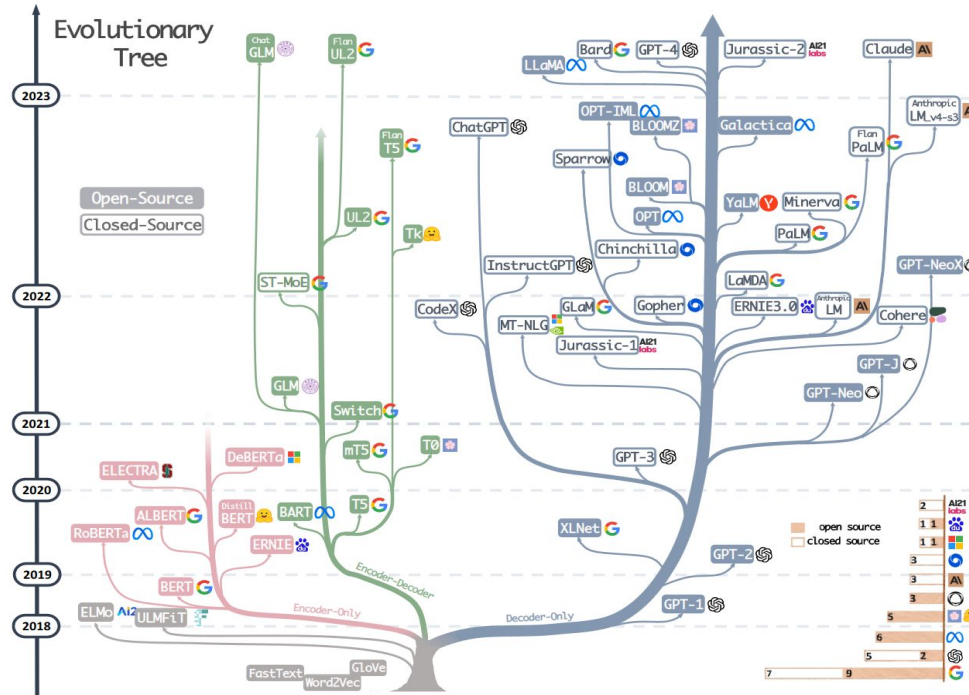
Si confronta la probabilità che il modello generi il token "great" rispetto al token "terrible".

Large Language Models (LLMs)

Few-Shot Text Classification



“Evolutionary Tree” dei recenti NLMs



Harnessing the Power of LLMs in Practice: A Survey on ChatGPT and Beyond (Yang et al., 2024), <https://dl.acm.org/doi/10.1145/3649506>

Un po' di pratica

<https://github.com/michelepapucci/llms-anatomy-course>