



Washington University in St. Louis

COLLEGE OF ARTS & SCIENCES

3 Day Seminar

LLM & Transfer Learning in Python
TRIADS Training Series

Ishita Gopal

March 25th, 26th, 27th, 2024

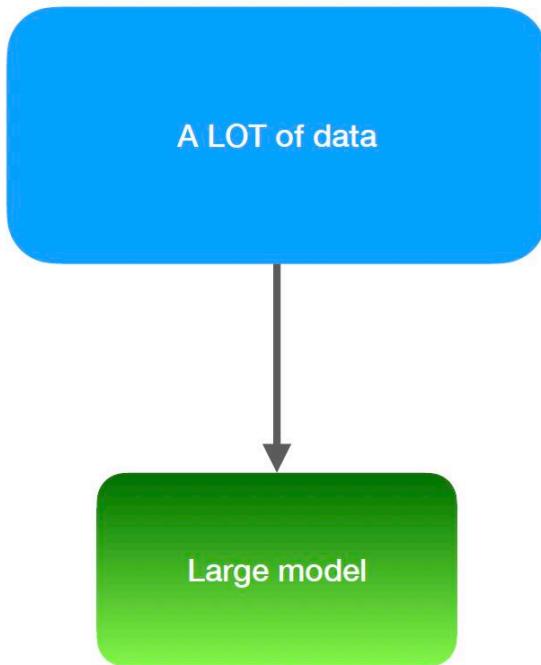
Plan for Day 1

- Overview of Transformers
 - Model Architecture
 - Attention
 - Pretraining tasks
 - Fine tuning tasks
- Colab walkthrough:
 - Introduction to the HuggingFace and the transformers library
 - Off the shelf prediction with different transformer models
 - First step for fine tuning – learning about tokenizers

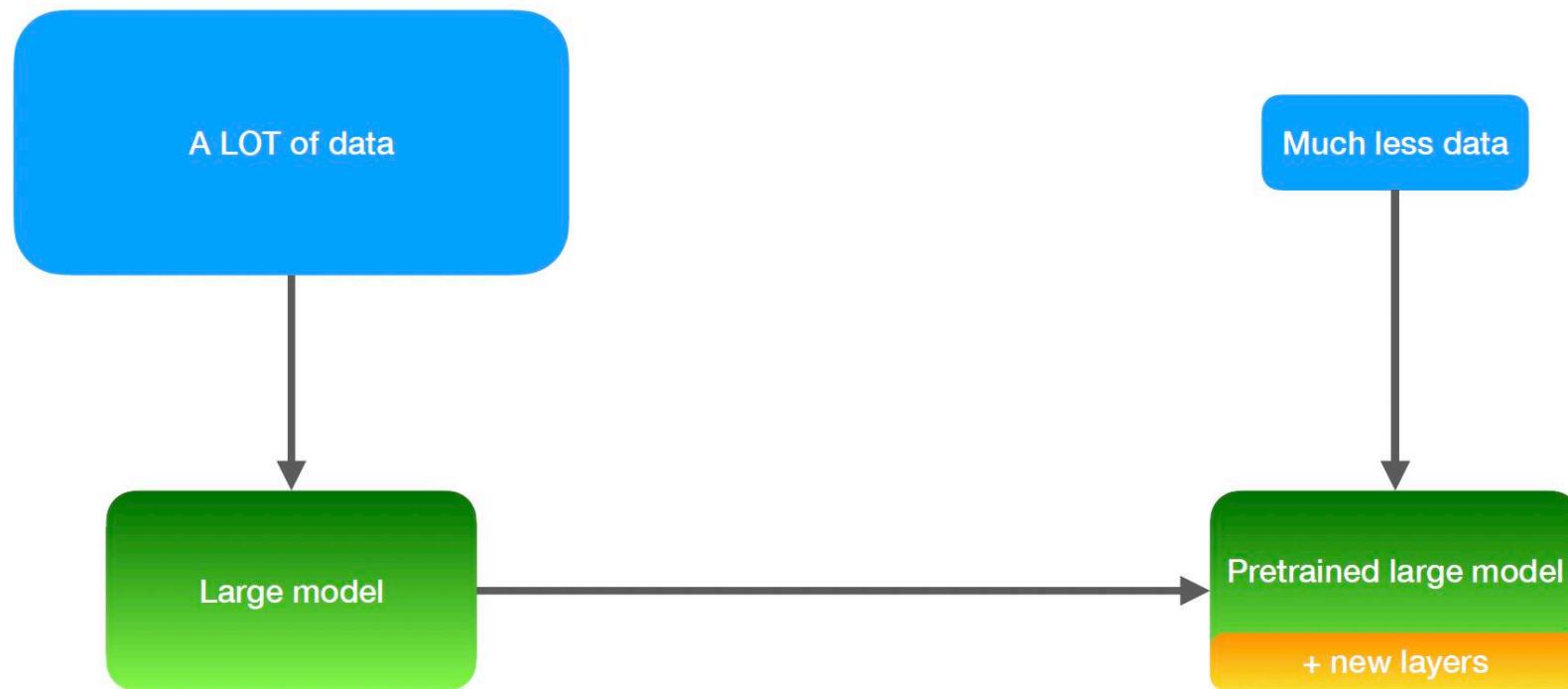
Day 2 and 3

- Day 2
 - Colab Walkthrough of how to implement fine tuning
- Day 3
 - Overview of hyperparameter tuning
 - Colab Walkthrough of how to implement hyperparameter tuning
 - Colab Walkthrough on Text Translation Using a Pretrained Transformer

Transfer Learning



**Traditional Machine Learning:
slow training on a lot of data**



Traditional Machine Learning:
slow training on a lot of data

Transfer learning:
fast training on a little data

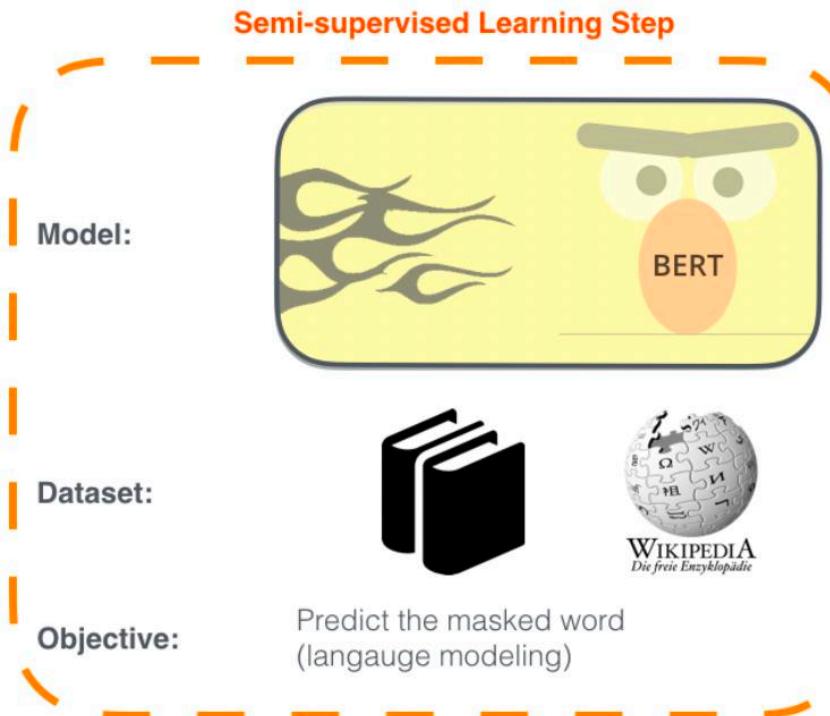
Some use cases

Transfer Learning

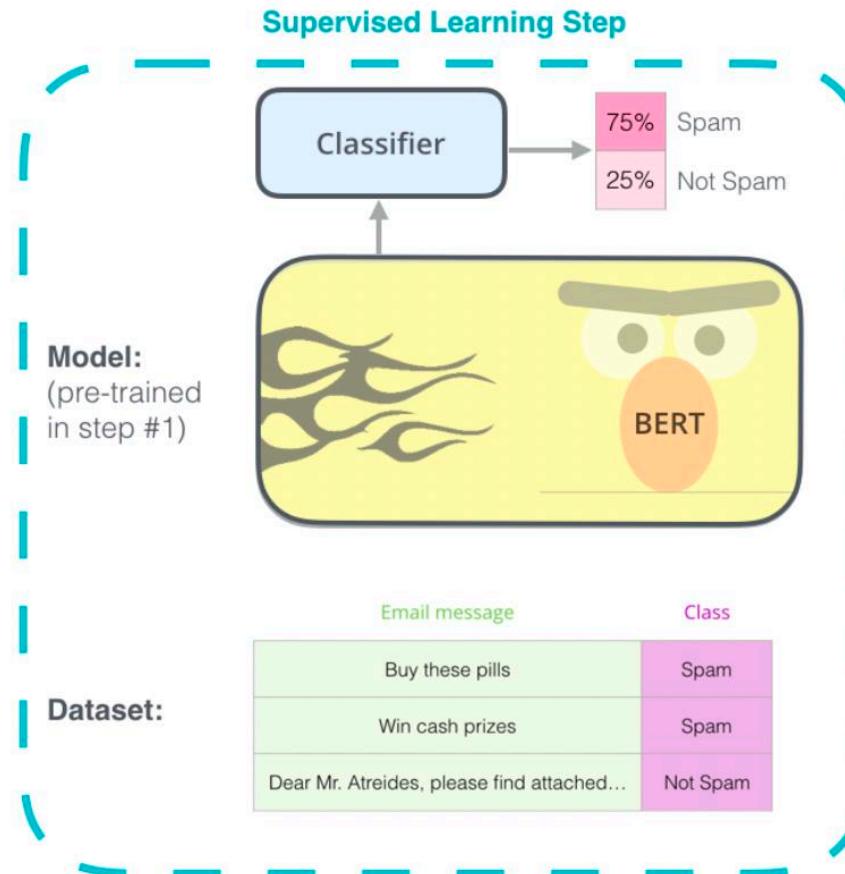
A typical use case

1 - **Semi-supervised** training on large amounts of text (books, wikipedia..etc).

The model is trained on a certain task that enables it to grasp patterns in language. By the end of the training process, BERT has language-processing abilities capable of empowering many models we later need to build and train in a supervised way.



2 - **Supervised** training on a specific task with a labeled dataset.



Source: Alammar (2018)

Online Incivility in the 2020 Congressional Elections

Political Research Quarterly
2022, Vol. 75(2) 512–526
© The Author(s) 2022
Article reuse guidelines:
sagepub.com/journals-permissions
DOI: 10.1177/10659129221078863
journals.sagepub.com/home/prq


Michael Heseltine¹  and Spencer Dorsey²

Abstract

This paper explores the prevalence and correlates of political incivility among Congressional candidates in the 2020 election cycle, focusing specifically on which types of candidates were most likely to use uncivil language in their online communications and the self-reinforcing nature of incivility between candidates. Based on a comprehensive analysis of more than two million tweets sent by major party candidates in the 2020 House and Senate races, we conclude that several individual and electoral factors were influential in driving candidate incivility. Specifically, Republicans, challengers, and candidates in less competitive races were more likely to use uncivil rhetoric. Women, racial minorities, and candidates running in open seat races were less prone to incivility. We also find that incivility begets incivility, with candidates whose opponents used higher rates of incivility also being more likely to use incivility themselves. Uncivil tweets were also found to generate significantly more likes and retweets, suggesting that incivility is a viable means of driving engagement for candidates. These results shed light on the factors behind incivility among political elites, as well as highlight the feedback effects which contribute to a self-reinforcing rise in political incivility.

Detect incivility in political tweets using BERTweet

Table A1: Custom Model (Tuned Transformer) Performance and Comparison to VADER Sentiment Model

	Precision		Recall		F1-Score		<i>n</i>
	VADER	Tuned Transformer	VADER	Tuned Transformer	VADER	Tuned Transformer	
Neutral/Civil	0.85	0.92	0.85	0.95	0.85	0.94	626
Uncivil	0.46	0.8	0.46	0.71	0.46	0.75	170
Macro Average	0.66	0.86	0.66	0.79	0.66	0.84	796
Weighted Average	0.77	0.9	0.77	0.9	0.77	0.9	796

Attention to the COVID-19 Pandemic on Twitter: Partisan Differences Among U.S. State Legislators

Taegyoon Kim✉, Nitheesha Nakka✉, Ishita Gopal✉, Bruce A. Desmarais✉, Abigail Mancinelli✉,
Jeffrey J. Harden✉, Hyein Ko✉, Frederick J. Boehmke✉

First published: 15 December 2021 | <https://doi.org/10.1111/lsq.12367> | Citations: 1

[Read the full text >](#)

 PDF  TOOLS  SHARE

Abstract

Subnational governments in the United States have taken the lead on many aspects of the response to the COVID-19 pandemic. Variation in government activity across states offers the opportunity to analyze responses in comparable settings. We study a common and informative activity among state officials—state legislators' attention to the pandemic on Twitter. We find that legislators' attention to the pandemic strongly correlates with the number of cases in the legislator's state, the national count of new deaths, and the number of pandemic-related public policies passed within the legislator's state. Furthermore, we find that the degree of responsiveness to pandemic indicators differs significantly across political parties, with Republicans exhibiting weaker responses, on average. Lastly, we find significant differences in the content of tweets about the pandemic by Democratic and Republican legislators, with Democrats focused on health indicators and impacts, and Republicans focused on business impacts and opening the economy.

Identify discussion of Covid-19 using BERT-base

	Precision (S.D)	Recall (S.D)	F-1 (S.D)
RF + Count	0.9713 (0.0187)	0.3701 (0.0343)	0.5354 (0.0384)
RF + TFIDF	0.9746 (0.0240)	0.3382 (0.0268)	0.5017 (0.0324)
RF + GloVe	0.8302 (0.0284)	0.2548 (0.0567)	0.3866 (0.065)
XGB + Count	0.9392 (0.0072)	0.6604 (0.045)	0.7744 (0.0284)
XGB + TFIDF	0.9365 (0.009)	0.6325 (0.0344)	0.7544 (0.0231)
XGB + GloVe	0.7391 (0.0324)	0.4698 (0.0538)	0.5718 (0.0366)
BERT Fine Tuned	0.8504 (0.0471)	0.8428 (0.0474)	0.8446 (0.0139)

Table S2. Performance of Classifiers (5-fold Cross Validation)

Sentiment is Not Stance: Target-Aware Opinion Classification for Political Text Analysis

Samuel E. Bestvater^{id} and Burt L. Monroe

Department of Political Science, The Pennsylvania State University, University Park, PA, USA. E-mail: seb654@psu.edu, burtmonroe@psu.edu

Abstract

Sentiment analysis techniques have a long history in natural language processing and have become a standard tool in the analysis of political texts, promising a conceptually straightforward automated method of extracting meaning from textual data by scoring documents on a scale from positive to negative. However, while these kinds of sentiment scores can capture the overall tone of a document, the underlying concept of interest for political analysis is often actually the document's stance with respect to a given target—how positively or negatively it frames a specific idea, individual, or group—as this reflects the author's underlying political attitudes. In this paper, we question the validity of approximating author stance through sentiment scoring in the analysis of political texts, and advocate for greater attention to be paid to the conceptual distinction between a document's sentiment and its stance. Using examples from open-ended survey responses and from political discussions on social media, we demonstrate that in many political text analysis applications, sentiment and stance do not necessarily align, and therefore sentiment analysis methods fail to reliably capture ground-truth document stance, amplifying noise in the data and leading to faulty conclusions.

Keywords: text-as-data, sentiment analysis, political stance, machine learning

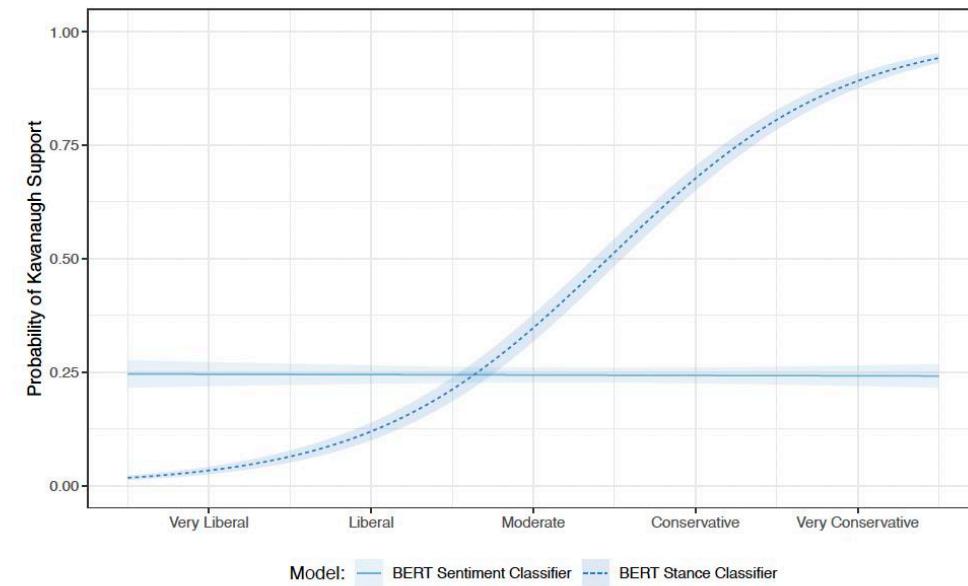
Sentiment and Stance in Tweets About the Kavanaugh Confirmation using BERT-base

Table 6. Classifier Performance: Kavanaugh Tweets

Classifier	F1 Score (Predicting Sentiment)	F1 Score (Predicting Stance)
Lexicoder	0.788 (0.005)	0.572 (0.014)
VADER	0.754 (0.005)	0.514 (0.011)
SVM (sentiment-trained)	0.943 (0.003)	0.514 (0.012)
BERT (sentiment-trained)	0.954 (0.002)	0.582 (0.005)
SVM (stance-trained)		0.935 (0.006)
BERT (stance-trained)		0.938 (0.002)

Reported figures are the average F1 score over 5-fold cross validation
Standard Errors in parentheses

Figure 5. Predicting Kavanaugh Support from Ideology



Source: Bestvater and Monroe, 2020.

How do Transformers work?

Transformers were introduced:

Attention Is All You Need

Ashish Vaswani*

Google Brain

avaswani@google.com

Noam Shazeer*

Google Brain

noam@google.com

Niki Parmar*

Google Research

nikip@google.com

Jakob Uszkoreit*

Google Research

usz@google.com

Llion Jones*

Google Research

llion@google.com

Aidan N. Gomez* †

University of Toronto

aidan@cs.toronto.edu

Lukasz Kaiser*

Google Brain

lukaszkaiser@google.com

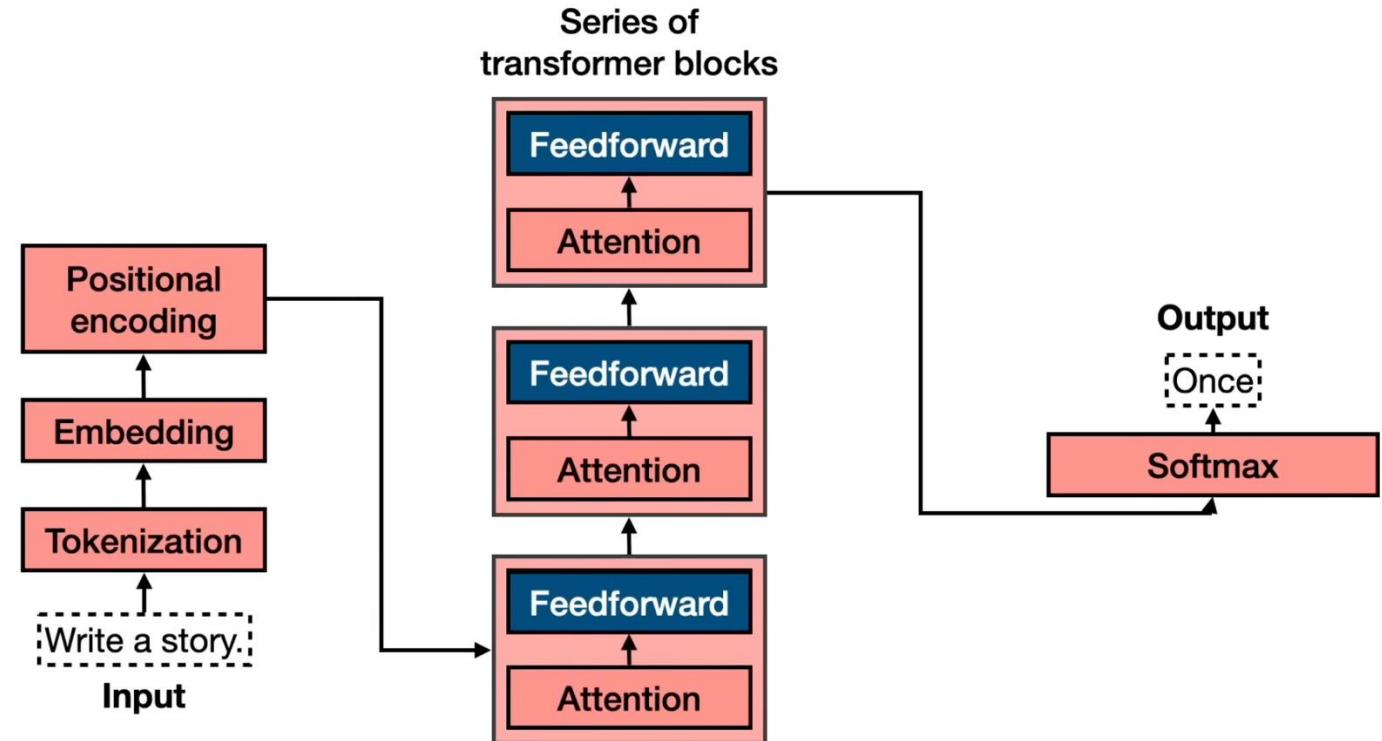
Illia Polosukhin* ‡

illia.polosukhin@gmail.com

Architecture:

Includes several steps:

- Tokenization
- Embeddings
- Attention Mechanisms
- Feedforward Neural Networks



Tokenization

Write a story.

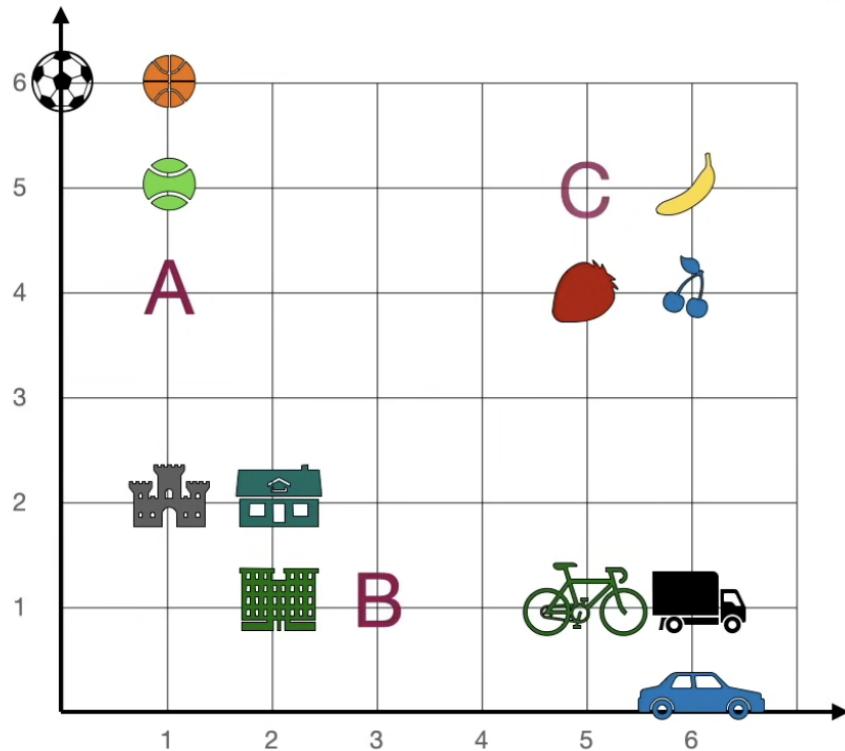


doesn't

The diagram shows the contraction "doesn't" followed by a horizontal arrow pointing to the right. To the right of the arrow are two colored boxes containing the words "does" and "n't" respectively. The word "does" is in a blue box and "n't" is in a green box.

Descriptions of words: Embeddings

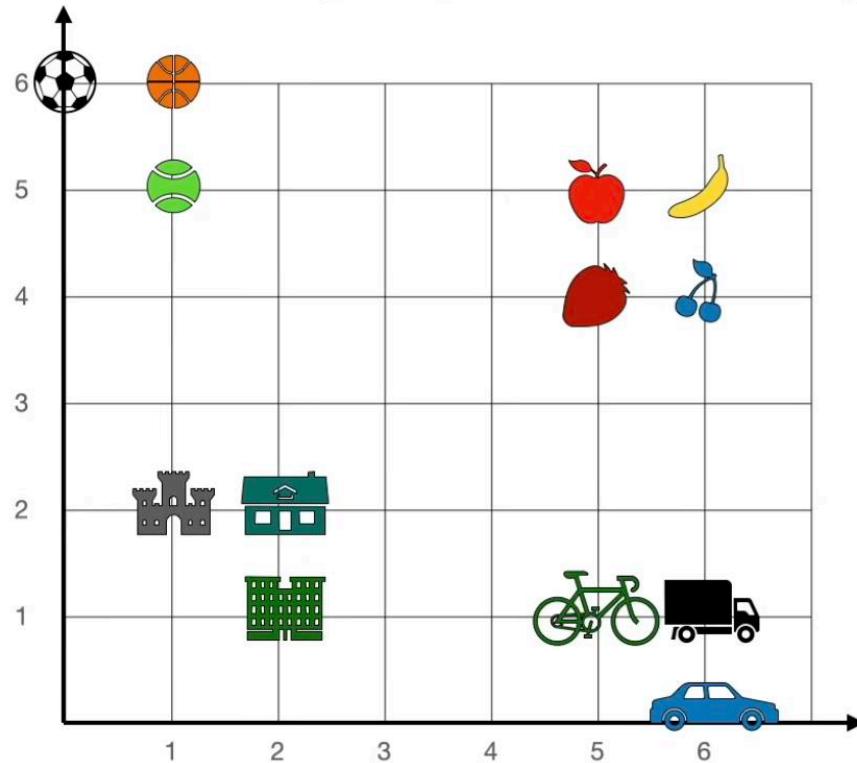
Where would you put the word apple?



Word	Numbers	
Apple	?	?
Banana	6	5
Strawberry	5	4
Cherry	6	4
Soccer	0	6
Basketball	1	6
Tennis	1	5
Castle	1	2
House	2	2
Building	2	1
Bicycle	5	1
Truck	6	1
Car	6	0

Similar words get similar numbers!

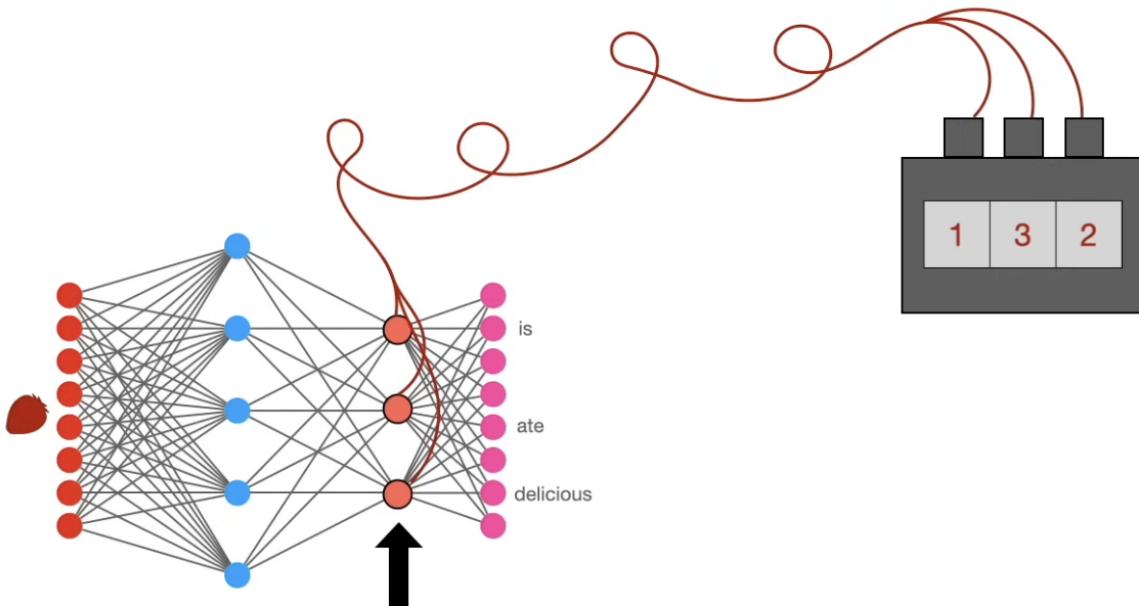
Where would you put the word apple?



Word	Numbers	
Apple	5	5
Banana	6	5
Strawberry	5	4
Cherry	6	4
Soccer	0	6
Basketball	1	6
Tennis	1	5
Castle	1	2
House	2	2
Building	2	1
Bicycle	5	1
Truck	6	1
Car	6	0

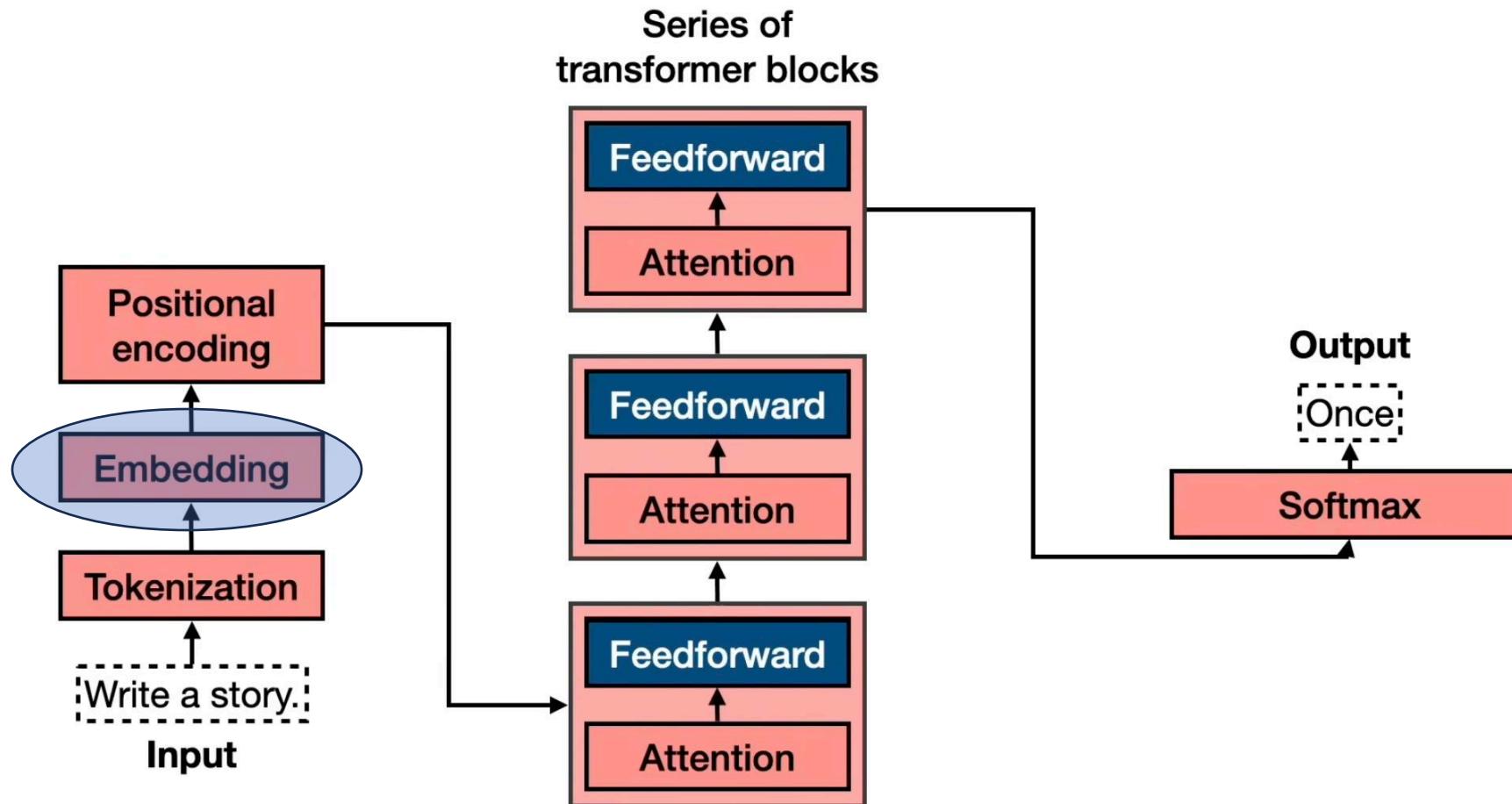
How can we get these numbers?

- Train Neural Networks
- Each layer understands and captures different properties of the word.
- The last layer understands the word pretty well!
- Each word can be associated to hundreds/thousands of numbers

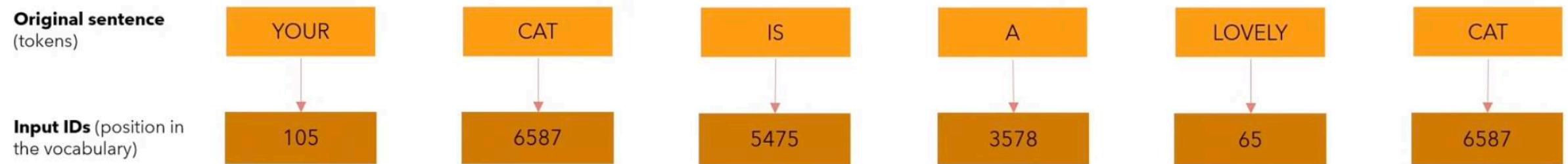


Word	Numbers		
Strawberry	1	3	2
Apple	1.1	2.9	2.2
Castle	7	-5.4	0.4

Going back to the Embeddings step in the Transformer...



Step 1: Map tokens onto numbers called “input ids” (these are fixed)



Step 2: Map Input Ids to a vector of size 512 → Embeddings!.

Original sentence (tokens)	YOUR	CAT	IS	A	LOVELY	CAT
Input IDs (position in the vocabulary)	105	6587	5475	3578	65	6587
Embedding (vector of size 512)	952.207 5450.840 1853.448 ... 1.658 2671.529	171.411 3276.350 9192.819 ... 3633.421 8390.473	621.659 1304.051 0.565 ... 7679.805 4506.025	776.562 5567.288 58.942 ... 2716.194 5119.949	6422.693 6315.080 9358.778 ... 2141.081 735.147	171.411 3276.350 9192.819 ... 3633.421 8390.473

Step 2: Map Input Ids to a vector of size 512 → Embeddings!.
- Same word → Same embedding



Step 2: Map Input Ids to a vector of size 512 → Embeddings!

- Same word → Same embedding

But!!

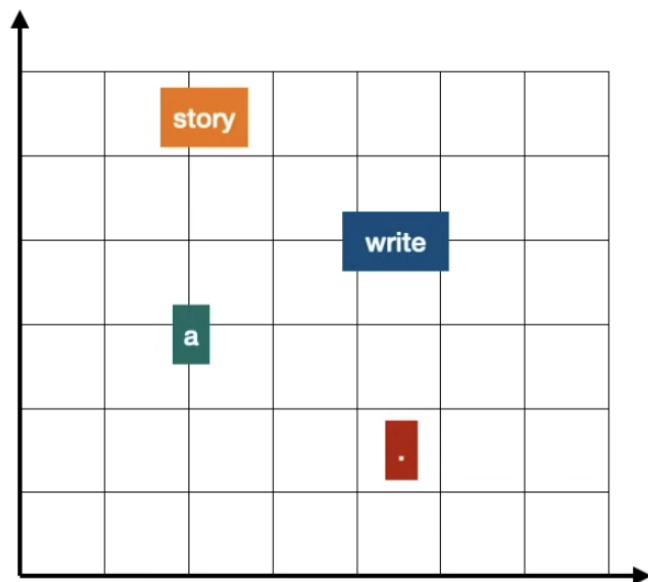
- Embeddings are not fixed. They are *parameters* for the model.

- Learned & change during training → to best represent *meaning of the word*.

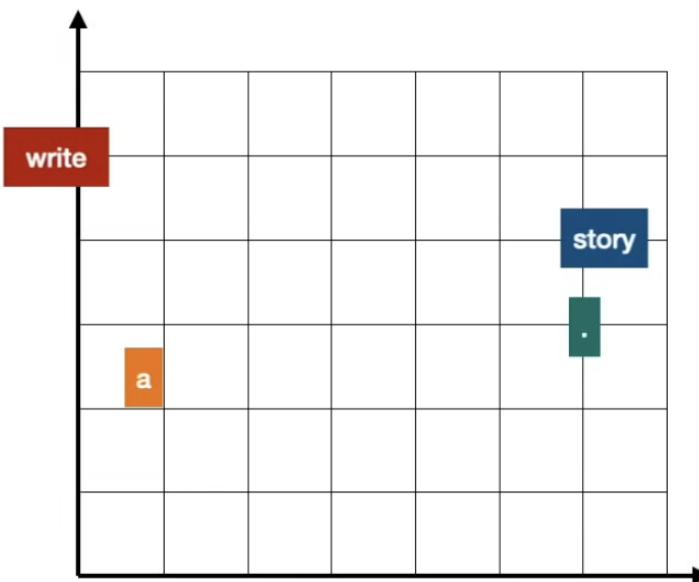


Account for word position:

write a story .



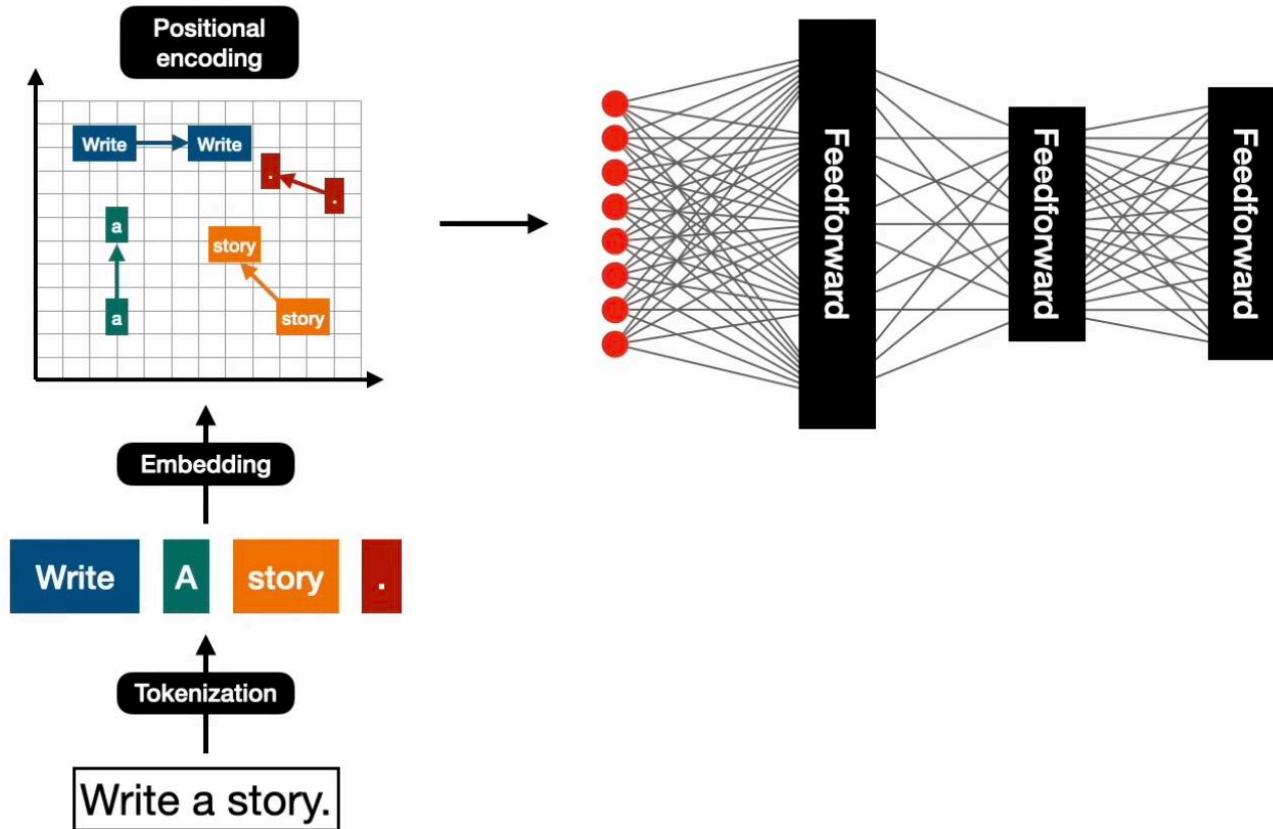
story . a write

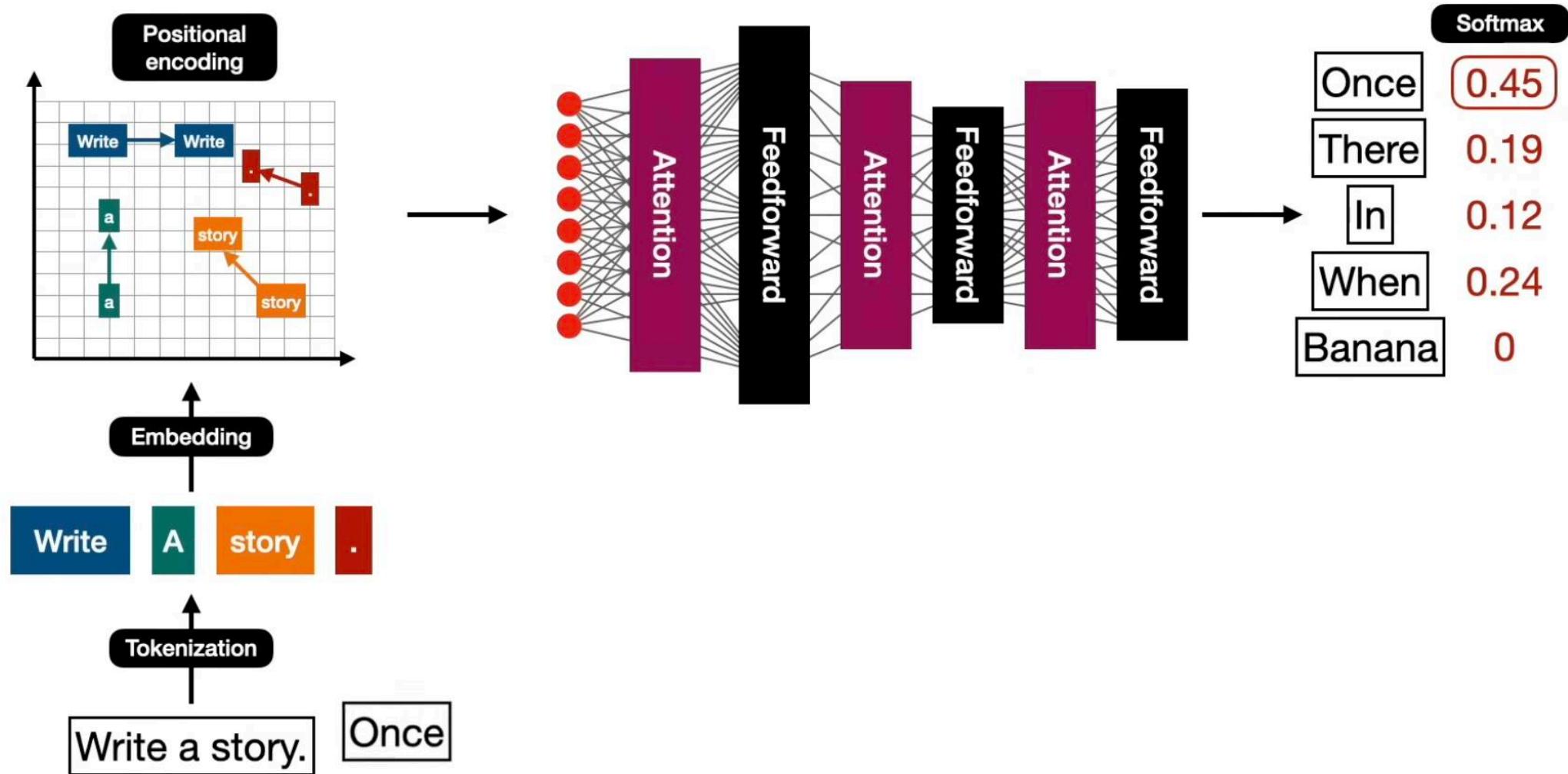


Add that to the embeddings:

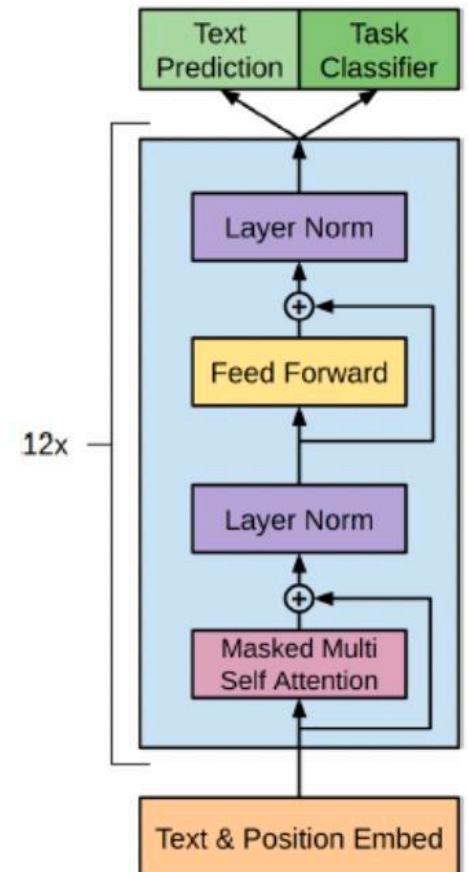
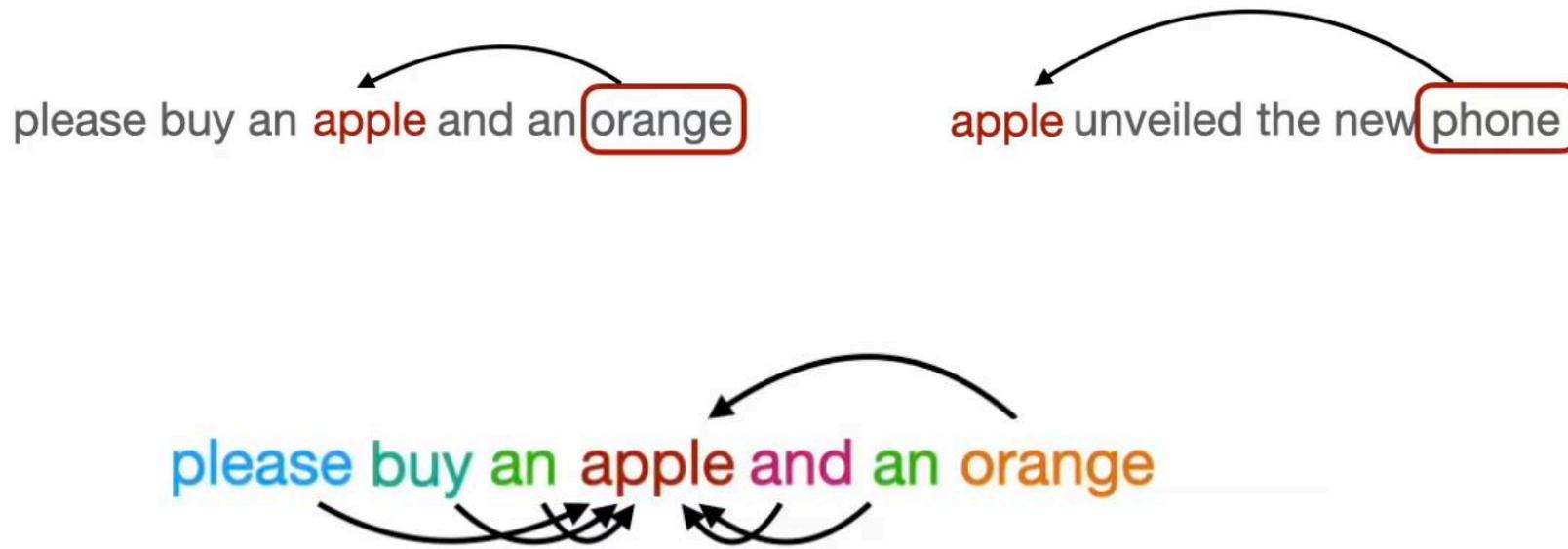
Original sentence	YOUR	CAT	IS	A	LOVELY	CAT
Embedding (vector of size 512)	952.207 5450.840 1853.448 ... 1.658 2671.529	171.411 3276.350 9192.819 ... 3633.421 8390.473	621.659 1304.051 0.565 ... 7679.805 4506.025	776.562 5567.288 58.942 ... 2716.194 5119.949	6422.693 6315.080 9358.778 ... 2141.081 735.147	171.411 3276.350 9192.819 ... 3633.421 8390.473
Position Embedding (vector of size 512). Only computed once and reused for every sentence during training and inference.	1664.068 8080.133 2620.399 ... 9386.405 3120.159	1281.458 7902.890 912.970 3821.102 1659.217 7018.620
Encoder Input (vector of size 512)	=	= 1835.479 11356.483 11813.218 ... 13019.826 11510.632	=	=	=	= 1452.869 11179.24 10105.789 ... 5292.638 15409.093

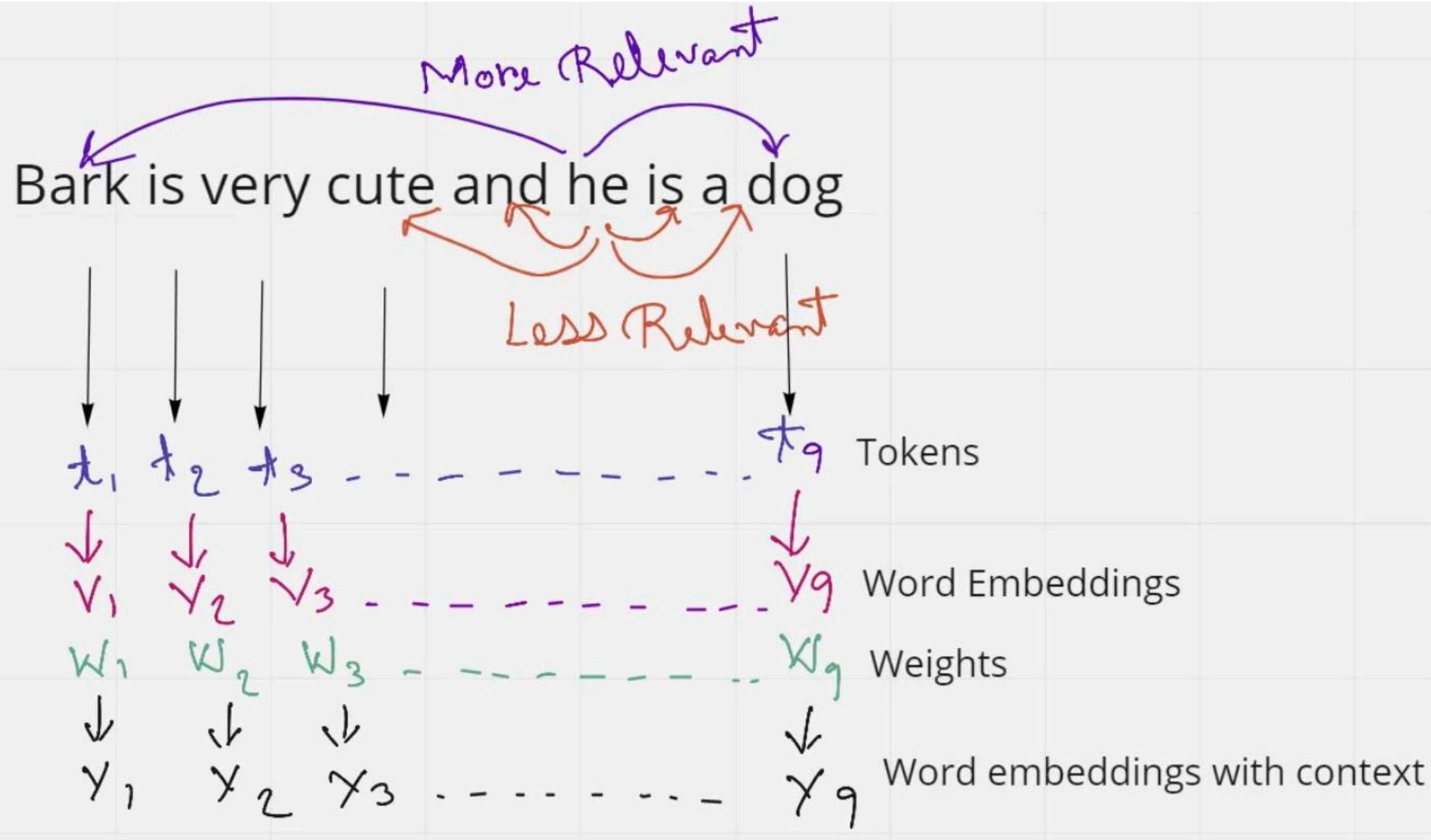
But *still* does not capture “context”





Attention for contextualized representations





1. Finding the Weights

$$V_1 V_1 = W_{11}$$

$$V_1 V_2 = W_{12}$$

$$V_1 V_3 = W_{13}$$

:

$$V_1 V_q = W_{1q}$$

Normalizing \rightarrow

$$\begin{aligned} &W_{11} \\ &W_{12} \\ &W_{13} \\ &\vdots \\ &W_{1q} \end{aligned}$$

Weights to re-weigh the first vector

2 Obtaining Embedding with context

$$W_{11} V_1 + W_{12} V_2 + W_{13} V_3 \dots + W_{1q} V_q = \underline{Y_1}$$

$$W_{21} V_1 + W_{22} V_2 + W_{23} V_3 \dots + W_{2q} V_q = Y_2$$

$$W_{q1} V_1 + W_{q2} V_2 + W_{q3} V_3 \dots + W_{qq} V_q = Y_q$$

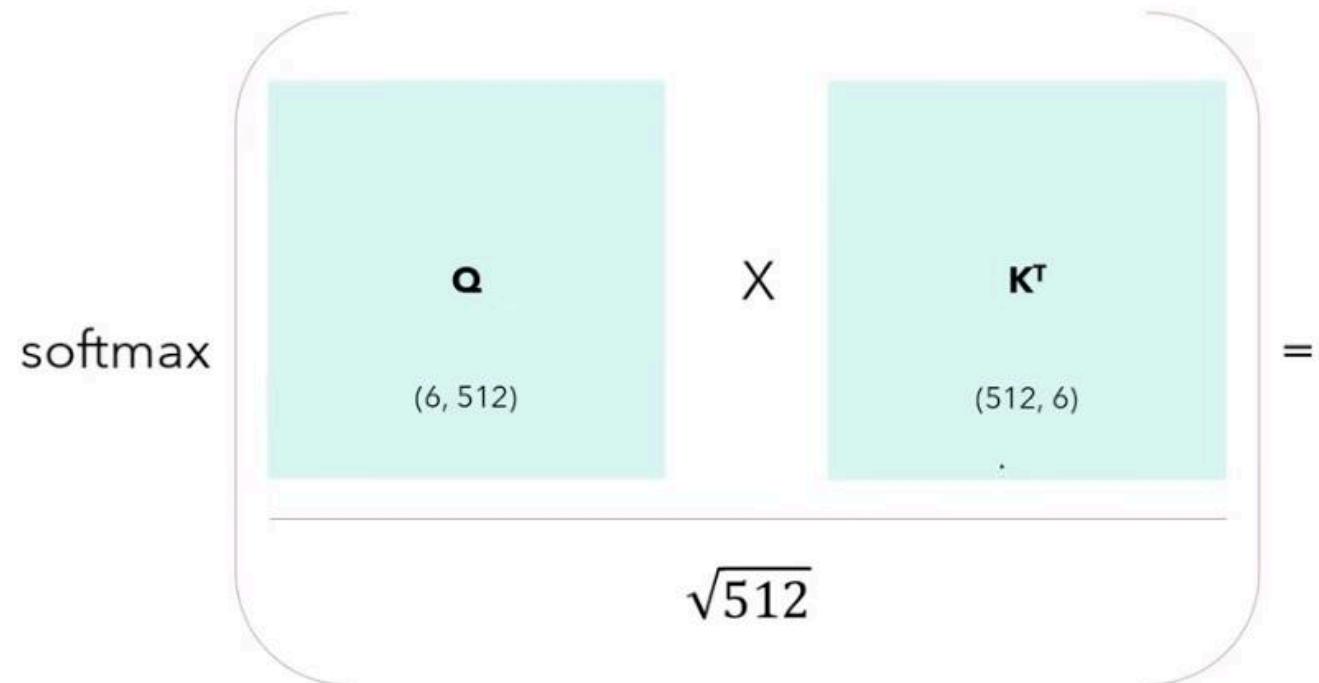
What is Self-Attention?

Self-Attention allows the model to relate words to each other.

In this simple case we consider the sequence length $\text{seq} = 6$ and $\mathbf{d}_{\text{model}} = \mathbf{d}_k = 512$.

The matrices **Q**, **K** and **V** are just the input sentence.

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



	YOUR	CAT	IS	A	LOVELY	CAT	Σ
YOUR	0.268	0.119	0.134	0.148	0.179	0.152	1
CAT	0.124	0.278	0.201	0.128	0.154	0.115	1
IS	0.147	0.132	0.262	0.097	0.218	0.145	1
A	0.210	0.128	0.206	0.212	0.119	0.125	1
LOVELY	0.146	0.158	0.152	0.143	0.227	0.174	1
CAT	0.195	0.114	0.203	0.103	0.157	0.229	1

* a

* for simplicity I considered only one head, which makes $\mathbf{d}_{\text{model}} = \mathbf{d}_k$.

Question

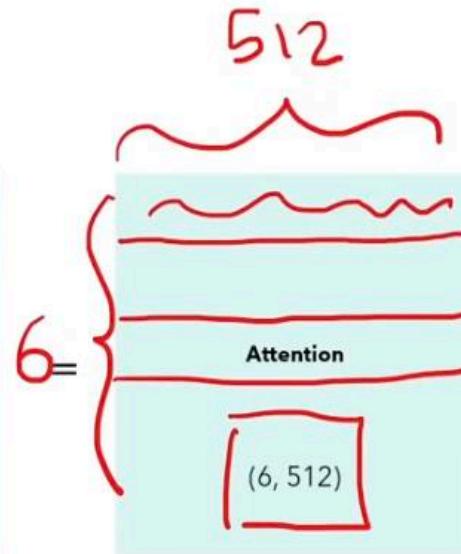
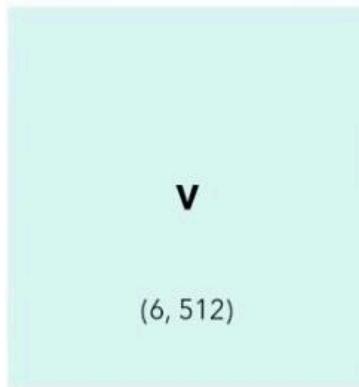
- Why are the scores not symmetric?
- After computing the raw attention scores → SoftMax normalization

- Row tokens (queries) attend to column tokens (keys)
- The matrix weights represent a way to probabilistically measure where attention is directed to when querying over keys
- i.e. to which key - and so to which token of the sentence - each query (token) mainly focuses to)

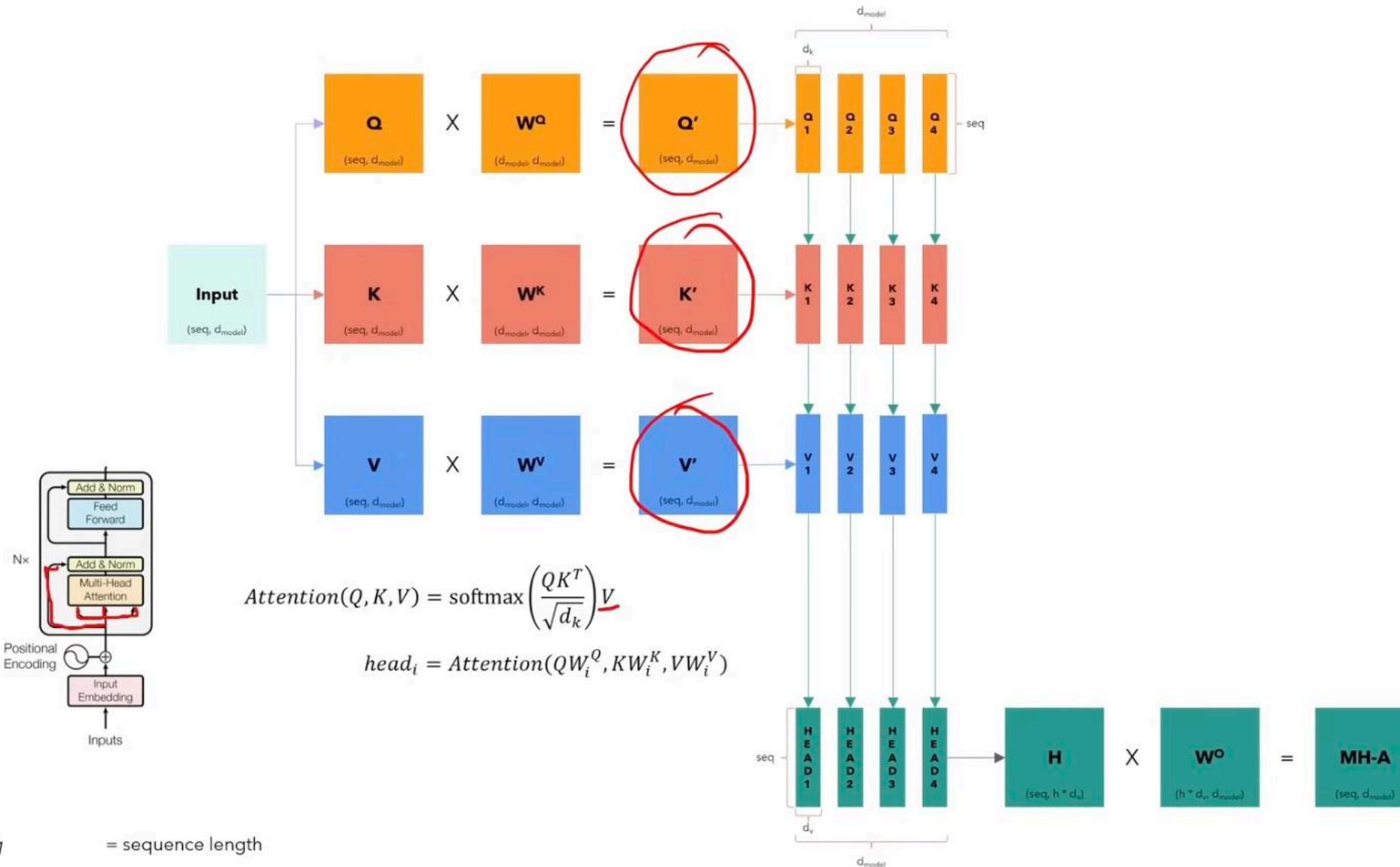
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

	YOUR	CAT	IS	A	LOVELY	CAT
YOUR	0.268	0.119	0.134	0.148	0.179	0.152
CAT	0.124	0.278	0.201	0.128	0.154	0.115
IS	0.147	0.132	0.262	0.097	0.218	0.145
A	0.210	0.128	0.206	0.212	0.119	0.125
LOVELY	0.146	0.158	0.152	0.143	0.227	0.174
CAT	0.195	0.114	0.203	0.103	0.157	0.229

X



Each row in this matrix captures not only the meaning (given by the embedding) or the position in the sentence (represented by the positional encodings) but also each word's interaction with other words.



seq = sequence length

d_{model} = size of the embedding vector

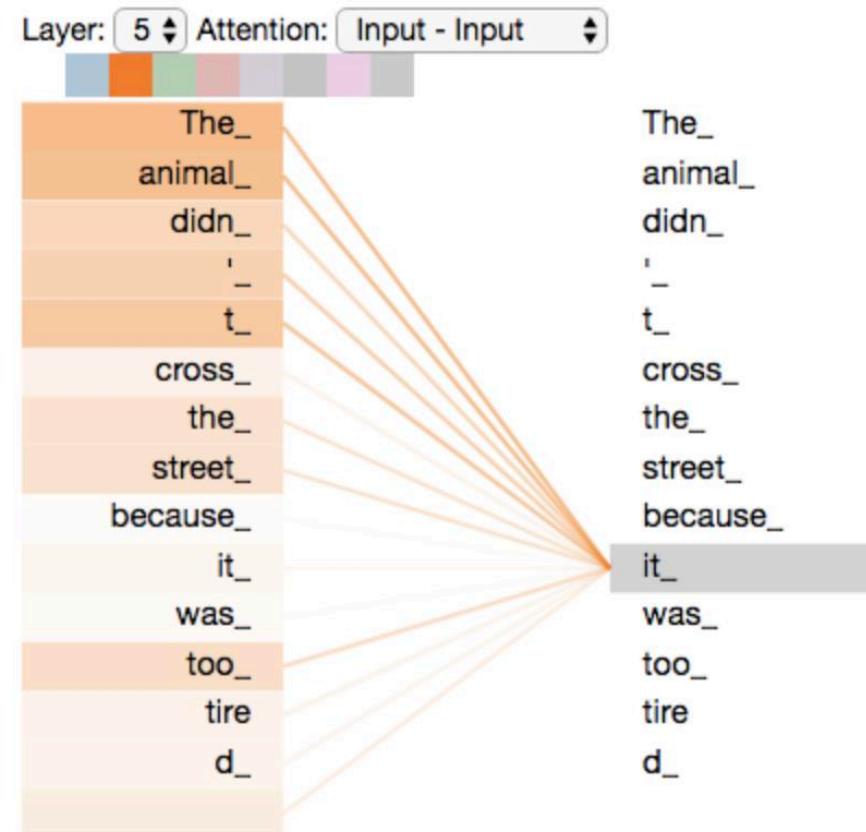
h = number of heads

$d_k = d_v = d_{model} / h$

Visualize Attention

“The animal didn’t cross the street because it was too tired”

- What does “it” in this sentence refer to?
- the street or to the animal?
- When the model is processing the word “it”, self-attention allows it to associate “it” with “animal”



Causal Masks

- In models like GPT, attention calculation is slightly different
 - GPT is an autoregressive language model based on transformers.
 - They generate text one token at a time, conditioning on the previously generated tokens.
- **Causal masks** are employed in the self-attention layers
 - to ensure that each token attends only to previous tokens during generation.
 - prevents the model from "cheating" by peeking at future tokens during training or generation

A language model is a probabilistic model that assign probabilities to sequence of words.
In practice, a language model allows us to compute the following:

$$P[\text{ "China"} | \text{ "Shanghai is a city in" }]$$

Next Token **Prompt**

Shanghai is a city in **China**, it is also a financial center.

Left context **Right context**

To model the probability distribution above, each word should only depend on words that come **before it** (*left context*).

	[SOS]	Before	my	bed	lies	a	pool	of	moon	bright
[SOS]	5.45	-∞	-∞	-∞	-∞	-∞	-∞	-∞	-∞	-∞
Before	4.28	2.46	-∞	-∞	-∞	-∞	-∞	-∞	-∞	-∞
my	8.17	3.56	5.54	-∞	-∞	-∞	-∞	-∞	-∞	-∞
bed	6.71	4.13	6.76	0.79	-∞	-∞	-∞	-∞	-∞	-∞
lies	5.43	7.59	3.91	6.14	9.03	-∞	-∞	-∞	-∞	-∞
a	4.42	4.35	7.55	3.14	1.35	7.57	-∞	-∞	-∞	-∞
pool	8.36	6.00	4.56	0.52	3.13	6.78	9.00	-∞	-∞	-∞
of	2.21	3.72	4.16	6.30	0.66	6.14	7.46	6.77	-∞	-∞
moon	4.08	6.22	5.00	4.20	5.72	5.35	7.46	3.55	4.70	-∞
bright	6.43	8.88	6.17	3.65	4.54	5.22	5.51	5.55	0.64	1.38

$$\sigma(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

	[SOS]	Before	my	bed	lies	a	pool	of	moon	bright
[SOS]	5.45	$-\infty$								
Before	4.28	2.46	$-\infty$							
my	8.17	3.56	5.54	$-\infty$						
bed	6.71	4.13	6.76	0.79	$-\infty$	$-\infty$	$-\infty$	$-\infty$	$-\infty$	$-\infty$
lies	5.43	7.59	3.91	6.14	9.03	$-\infty$	$-\infty$	$-\infty$	$-\infty$	$-\infty$
a	4.42	4.35	7.55	3.14	1.35	7.57	$-\infty$	$-\infty$	$-\infty$	$-\infty$
pool	8.36	6.00	4.56	0.52	3.13	6.78	9.00	$-\infty$	$-\infty$	$-\infty$
of	2.21	3.72	4.16	6.30	0.66	6.14	7.46	6.77	$-\infty$	$-\infty$
moon	4.08	6.22	5.00	4.20	5.72	5.35	7.46	3.55	4.70	$-\infty$
bright	6.43	8.88	6.17	3.65	4.54	5.22	5.51	5.55	0.64	1.38

$$\frac{QK^T}{\sqrt{d_k}}$$

	[SOS]	Before	my	bed	lies	a	pool	of	moon	bright
[SOS]	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Before	0.86	0.14	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
my	0.92	0.01	0.07	0.00	0.00	0.00	0.00	0.00	0.00	0.00
bed	0.47	0.04	0.49	0.00	0.00	0.00	0.00	0.00	0.00	0.00
lies	0.02	0.18	0.00	0.04	0.75	0.00	0.00	0.00	0.00	0.00
a	0.02	0.02	0.47	0.01	0.00	0.48	0.00	0.00	0.00	0.00
pool	0.31	0.03	0.01	0.00	0.00	0.06	0.59	0.00	0.00	0.00
of	0.00	0.01	0.02	0.15	0.00	0.12	0.47	0.23	0.00	0.00
moon	0.02	0.16	0.05	0.02	0.10	0.07	0.55	0.01	0.03	0.00
bright	0.07	0.71	0.05	0.03	0.01	0.02	0.03	0.03	0.02	0.03

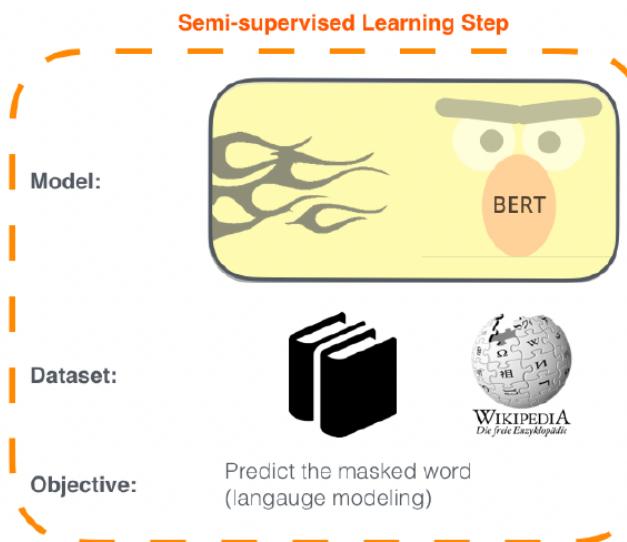
$$\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)$$

BERT - so what?

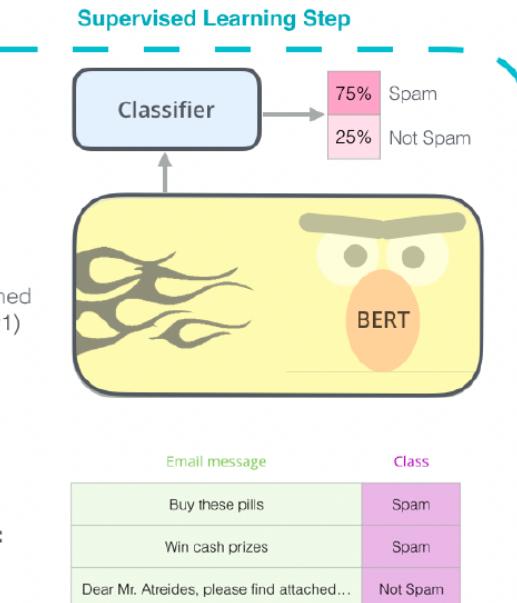
- BERT is a pretrained language model that can be fine-tuned to a specific task.

- 1 - **Semi-supervised** training on large amounts of text (books, wikipedia..etc).

The model is trained on a certain task that enables it to grasp patterns in language. By the end of the training process, BERT has language-processing abilities capable of empowering many models we later need to build and train in a supervised way.



- 2 - **Supervised** training on a specific task with a labeled dataset.



Source: Alammar (2018)

Pretraining #1: Masked Language Modeling

Use the output of the masked word's position to predict the masked word

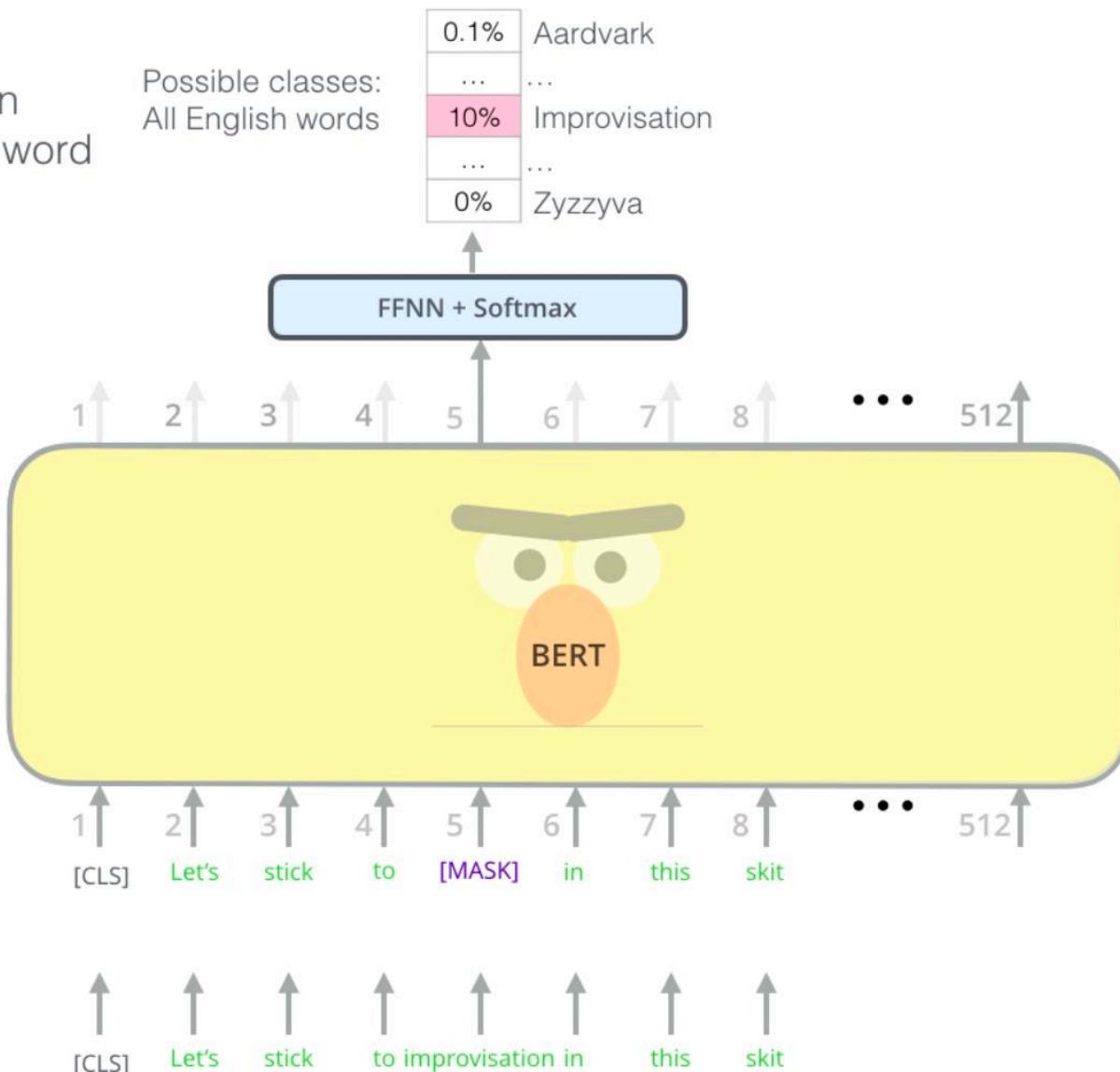
Possible classes:
All English words

0.1%	Aardvark
...	...
10%	Improvisation
...	...
0%	Zyzyva

FFNN + Softmax

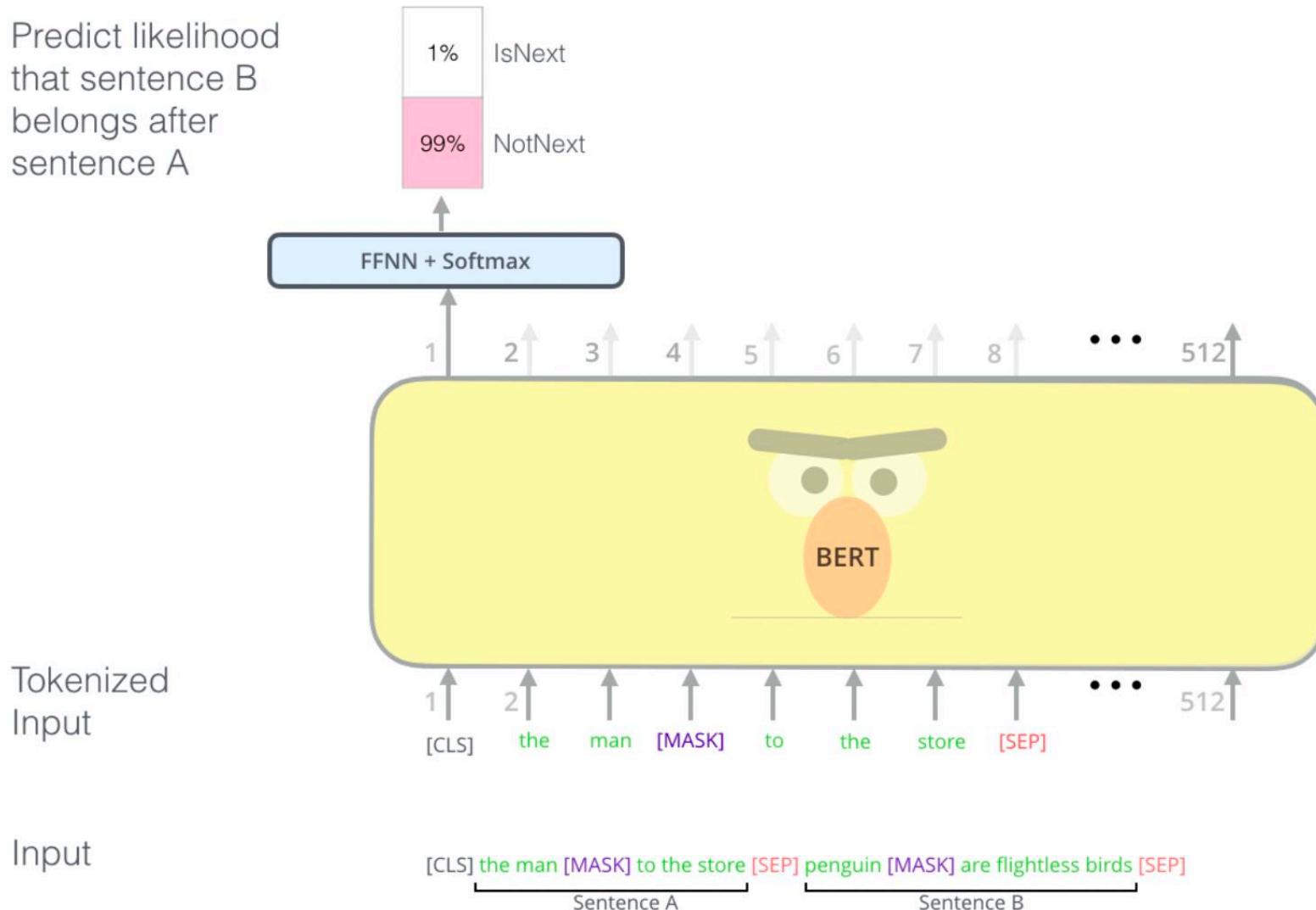
Randomly mask
15% of tokens

Input



Pretraining #2: Next Sentence Prediction

Predict likelihood
that sentence B
belongs after
sentence A



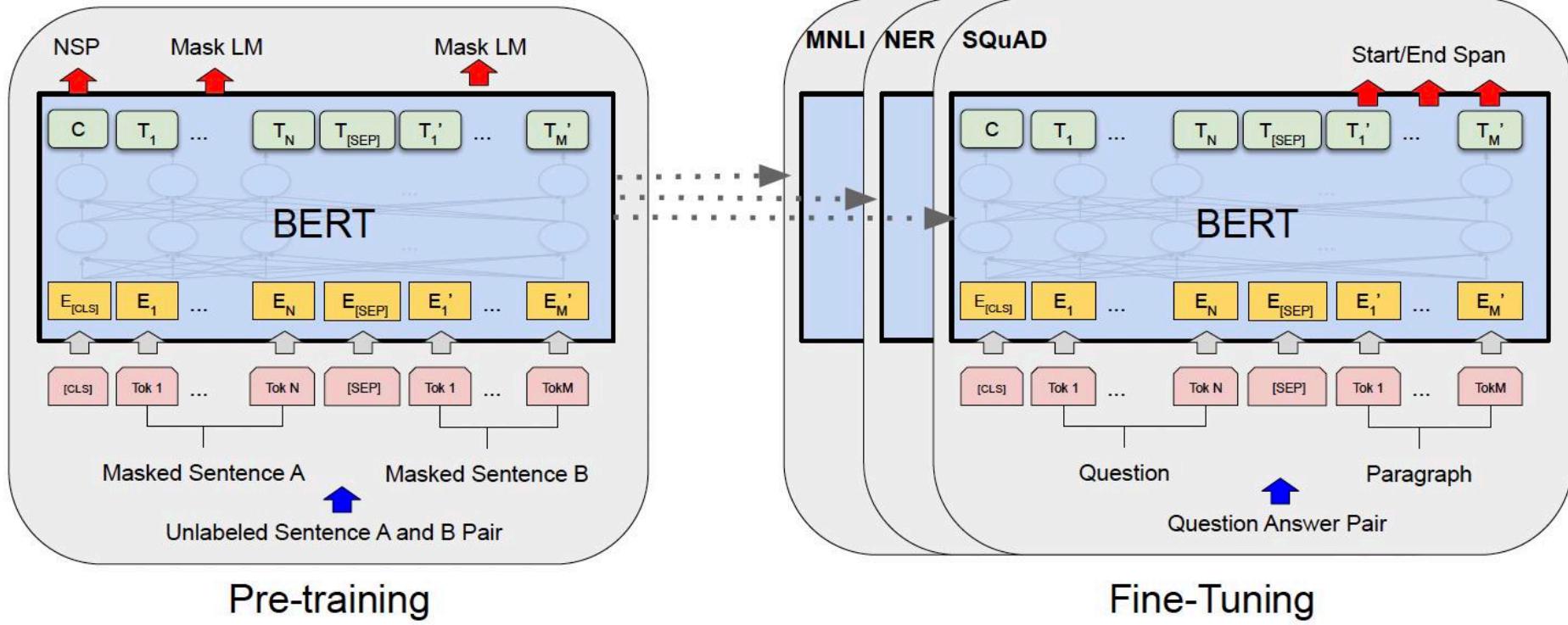


Figure 1: Overall pre-training and fine-tuning procedures for BERT. Apart from output layers, the same architectures are used in both pre-training and fine-tuning. The same pre-trained model parameters are used to initialize models for different down-stream tasks. During fine-tuning, all parameters are fine-tuned. [CLS] is a special symbol added in front of every input example, and [SEP] is a special separator token (e.g. separating questions/answers).

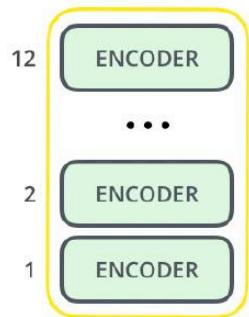
BERT (Google, 2018)



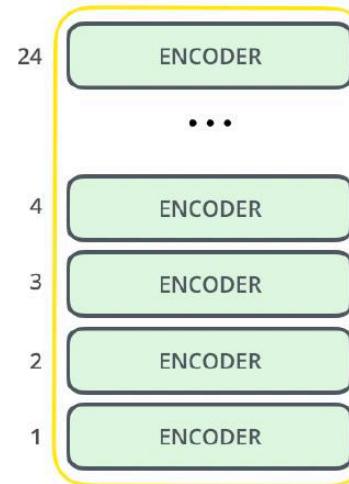
BERT_{BASE}



BERT_{LARGE}



BERT_{BASE}



BERT_{LARGE}

With BERT

- Work with fine-tuning
 - Text classification, Question Answering
- Does not work with prompts (GPT, Lamma)
 - Zero shot, Few shot, Chain of thought

Text Classification

- We want to automatically assign labels to a piece of text

My router's led is not working, I tried changing the power socket but still nothing.

Hardware

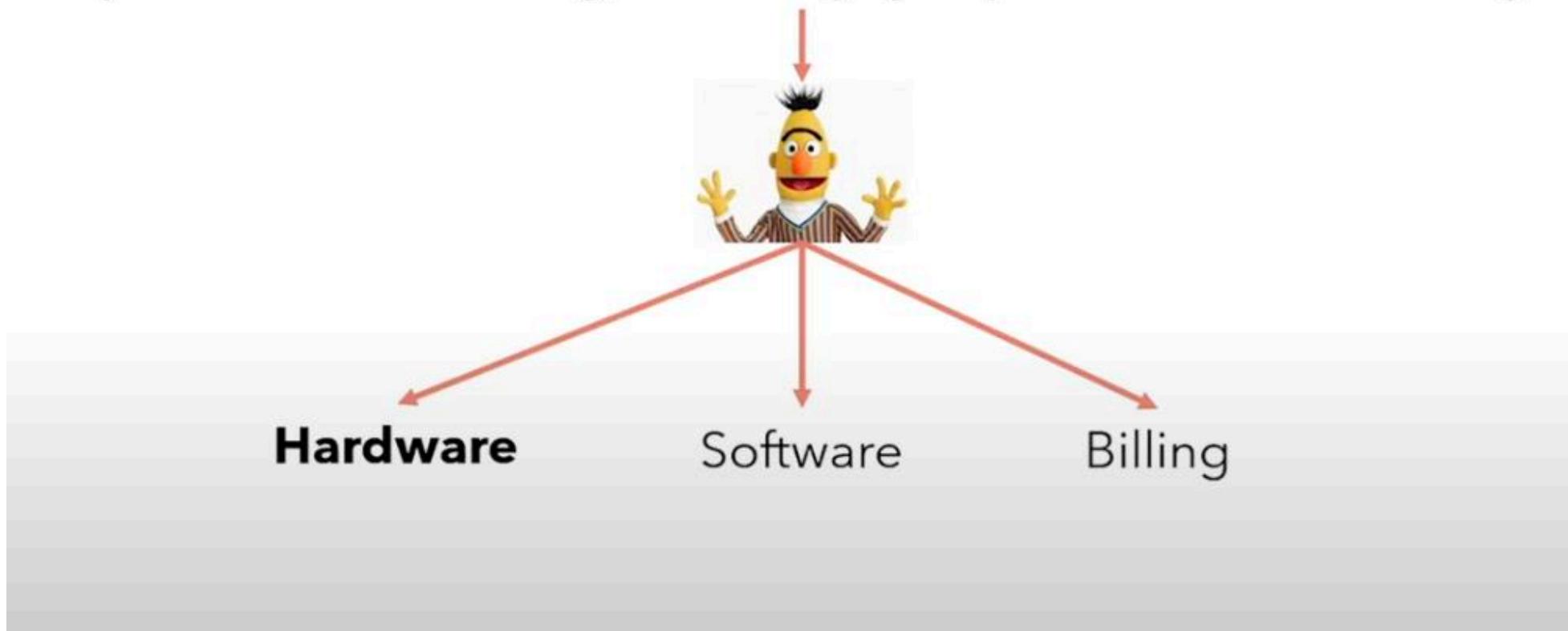
My router's web page doesn't allow me to change password anymore... I tried restarting it but nothing.

Software

In this month's bill I have been charged 100\$ instead of the usual 60\$, why is that?

Billing

My router's led is not working, I tried changing the power socket but still nothing.



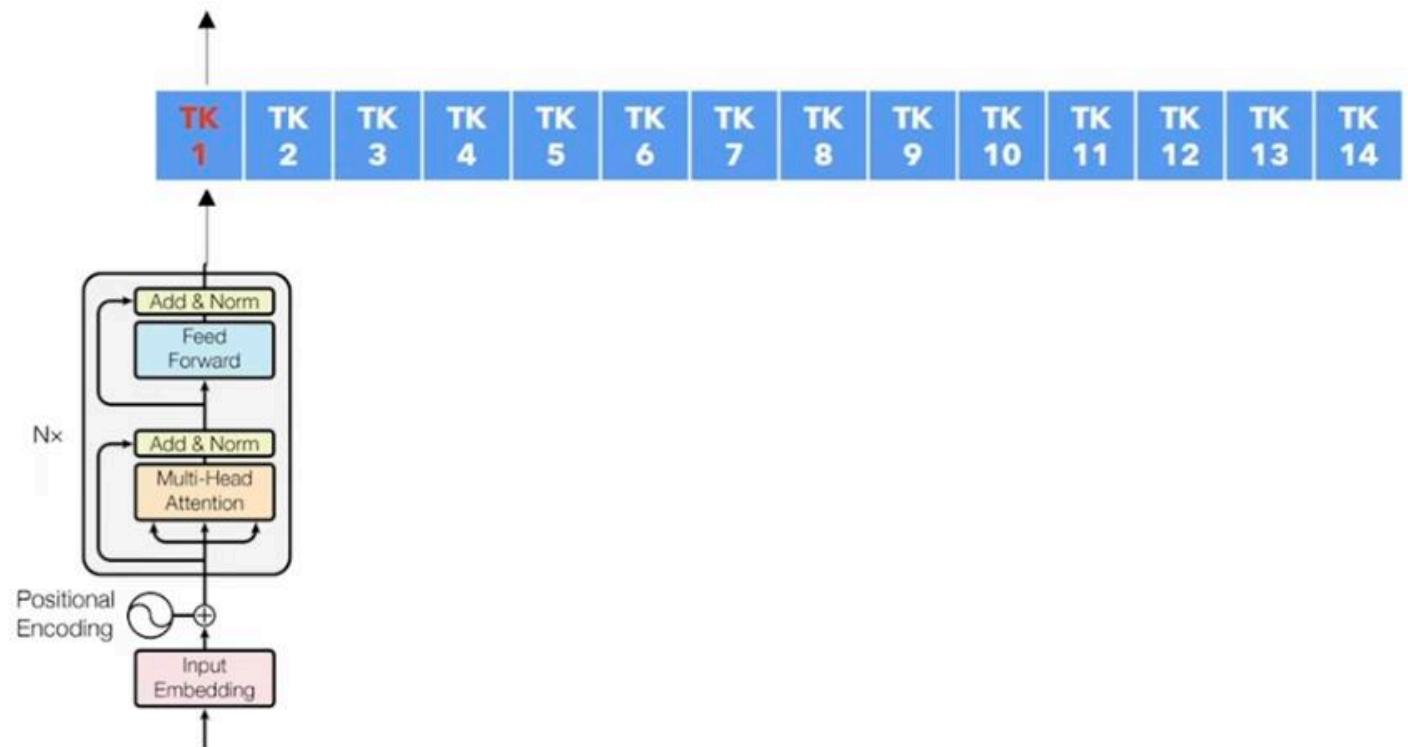
Target (1 token):

I dislike this



Output (16 tokens):

TK 1 | TK 2 | TK 3 | TK 4 | TK 5 | TK 6 | TK 7 | TK 8 | TK 9 | TK 10 | TK 11 | TK 12 | TK 13 | TK 14



Input (16 tokens):

[CLS] My router's led is not working, I tried changing the power socket but still nothing.

Arms Race

GPT-3 - 175000!



Let's look at the Colab Notebook here:

<https://tinyurl.com/triads-llm-workshop>

<https://tinyurl.com/triads-llm-translation>