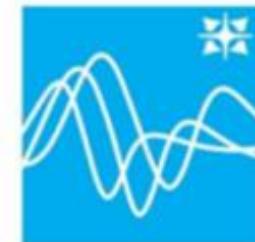


Context-based models

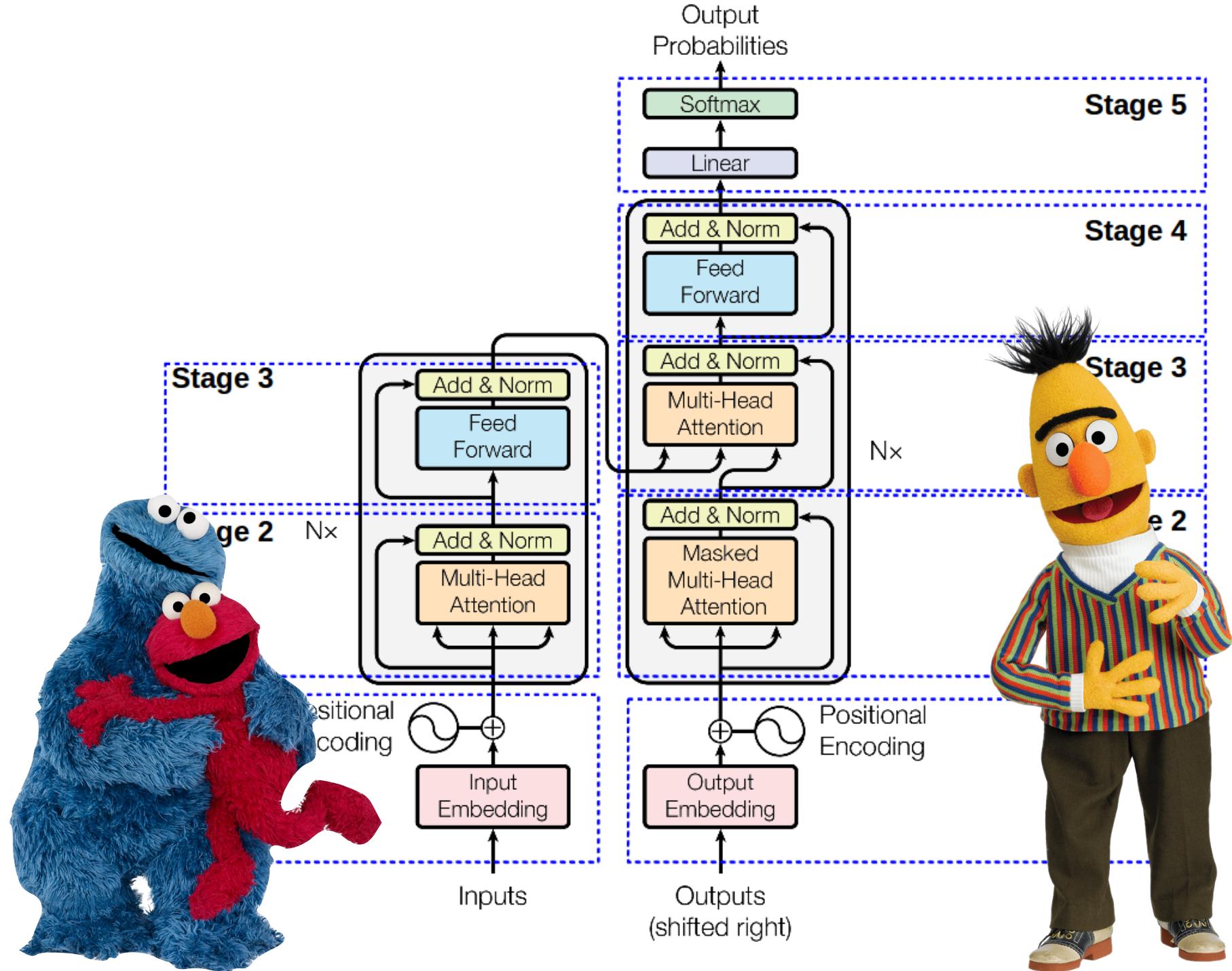
ELMo, BERT, GPT



2025



Кафедра
технологий
проектирования
сложных
технических
систем

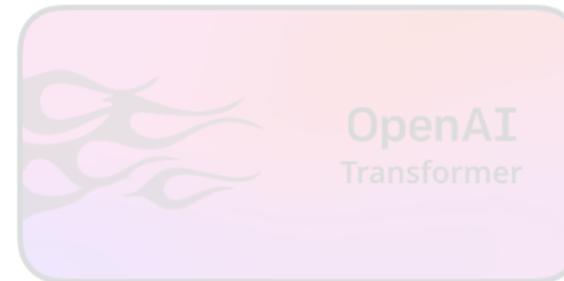
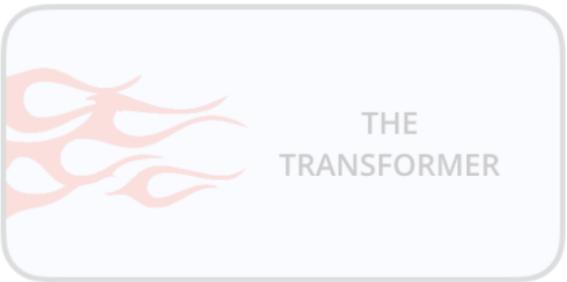


Другие модели



[The Illustrated BERT, ELMo, and co. \(How NLP Cracked Transfer Learning\)](#)

План лекции



[The Illustrated BERT, ELMo, and co. \(How NLP Cracked Transfer Learning\)](#)

ELMo



1. Expedited Labour Market Opinion
2. Electric Light Machine Organization
3. Enough Let's Move On

Special thx for @Anastasia Yanina:

https://github.com/ml-mipt/ml-mipt/blob/advanced/week05_BERT_and_LDA/Lecture_BERT_DIHT.pdf

ELMo



1. Expedited Labour Market Opinion
2. Electric Light Machine Organization
3. Enough Let's Move On

4. Embeddings from Language Models

Special thx for @Anastasia Yanina:

https://github.com/ml-mipt/ml-mipt/blob/advanced/week05_BERT_and_LDA/Lecture_BERT_DIHT.pdf



Hey ELMo, what's the embedding
of the word "stick"?

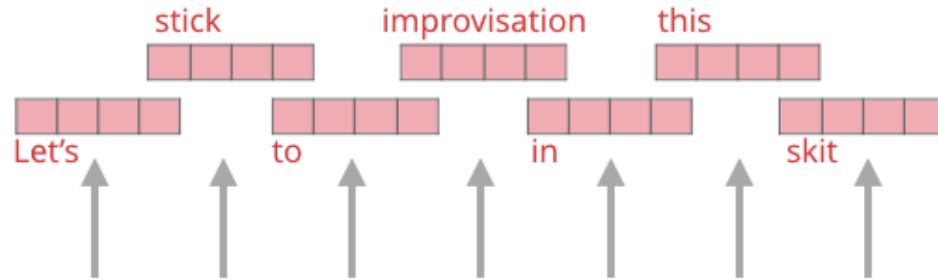
There are multiple possible
embeddings! Use it in a sentence.

Oh, okay. Here:
"Let's stick to improvisation in this
skit"

Oh in that case, the embedding is:
-0.02, -0.16, 0.12, -0.1etc

ELMo

ELMo
Embeddings



Words to embed



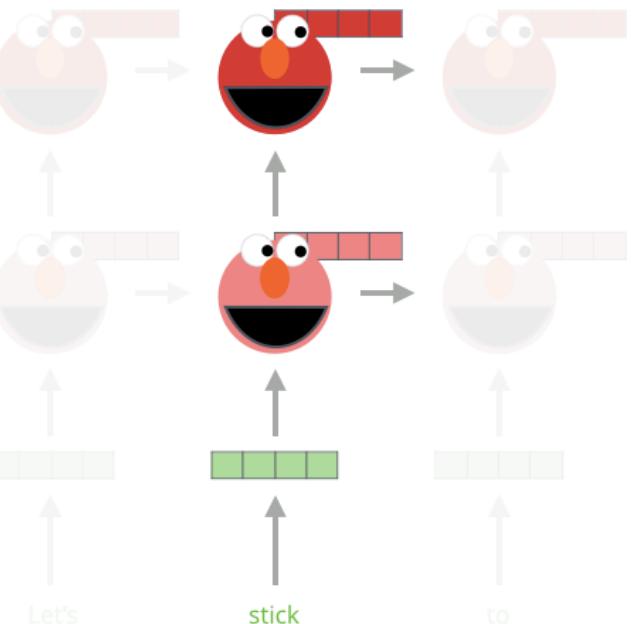
ELMo

Embedding of “stick” in “Let’s stick to” - Step #2

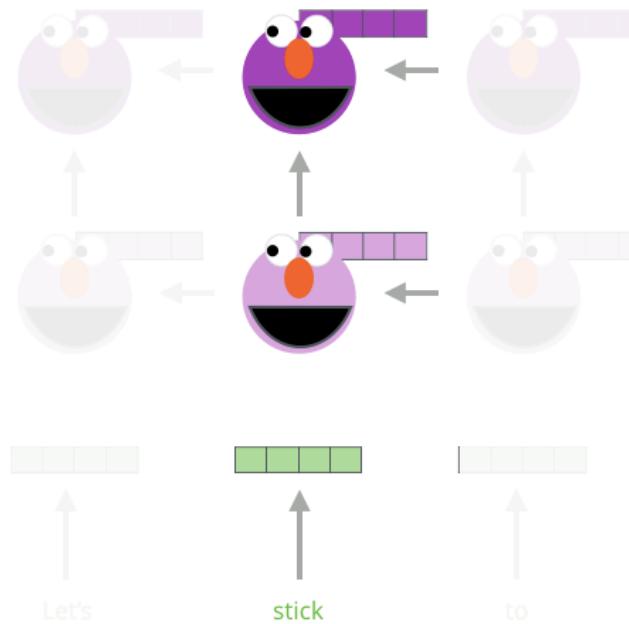
1- Concatenate hidden layers



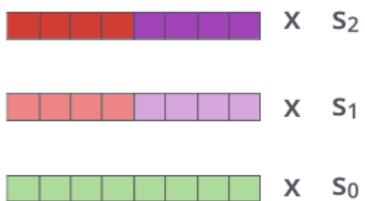
Forward Language Model



Backward Language Model



2- Multiply each vector by a weight based on the task

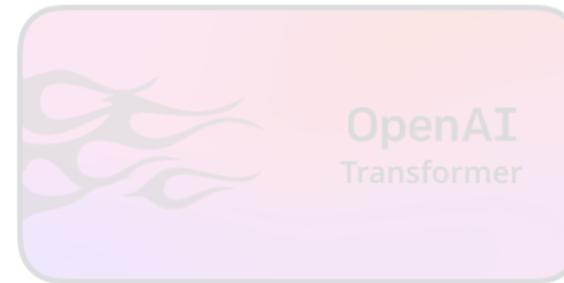
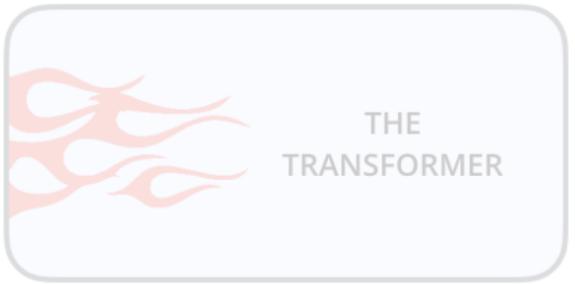


3- Sum the (now weighted) vectors



ELMo embedding of “stick” for this task in this context

Другие модели



[The Illustrated BERT, ELMo, and co. \(How NLP Cracked Transfer Learning\)](#)

BERT

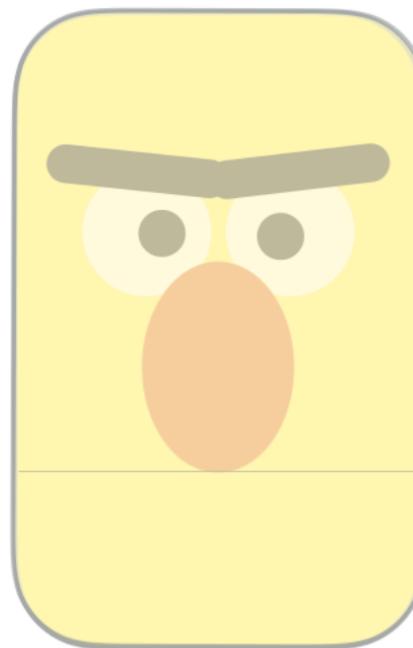


[1810.04805] BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

BERT

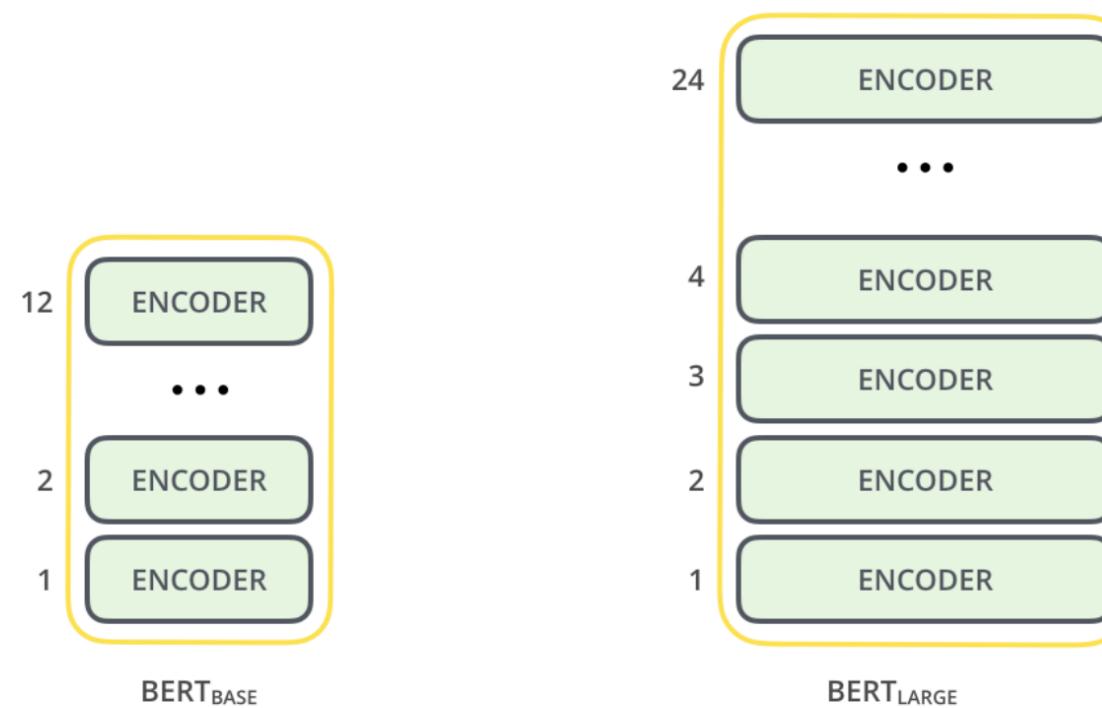


BERT_{BASE}



BERT_{LARGE}

BERT

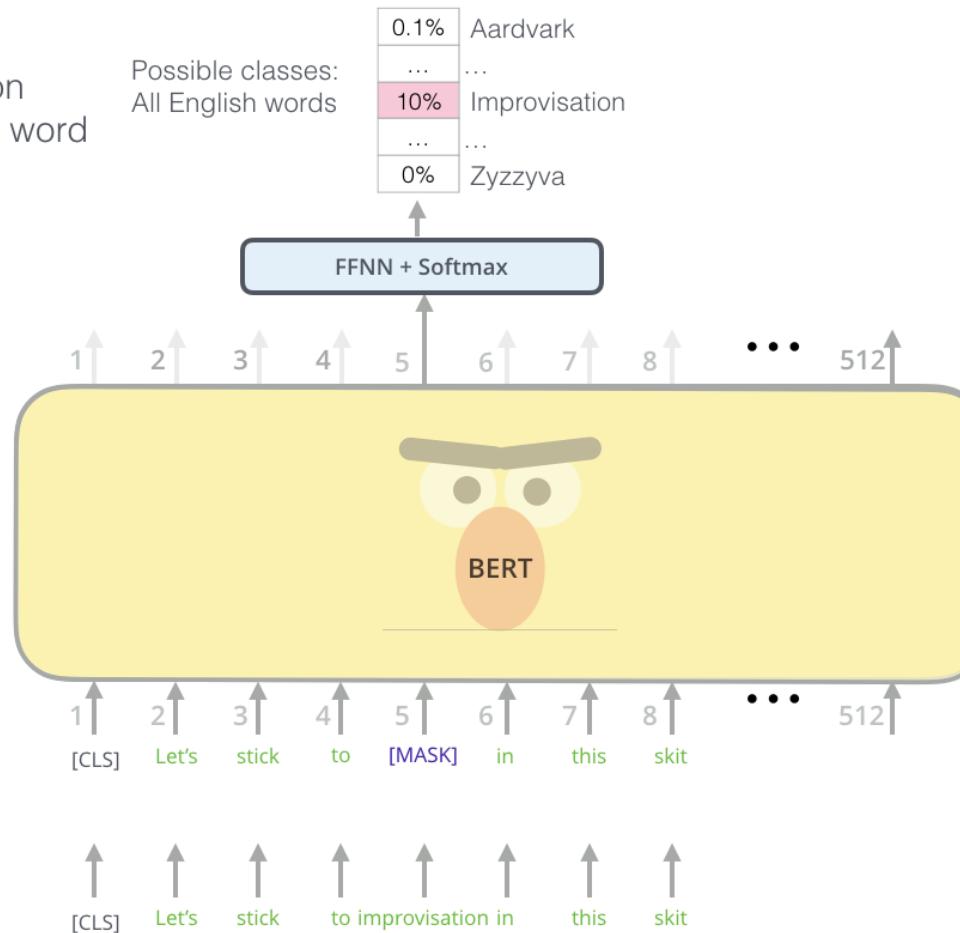


BERT

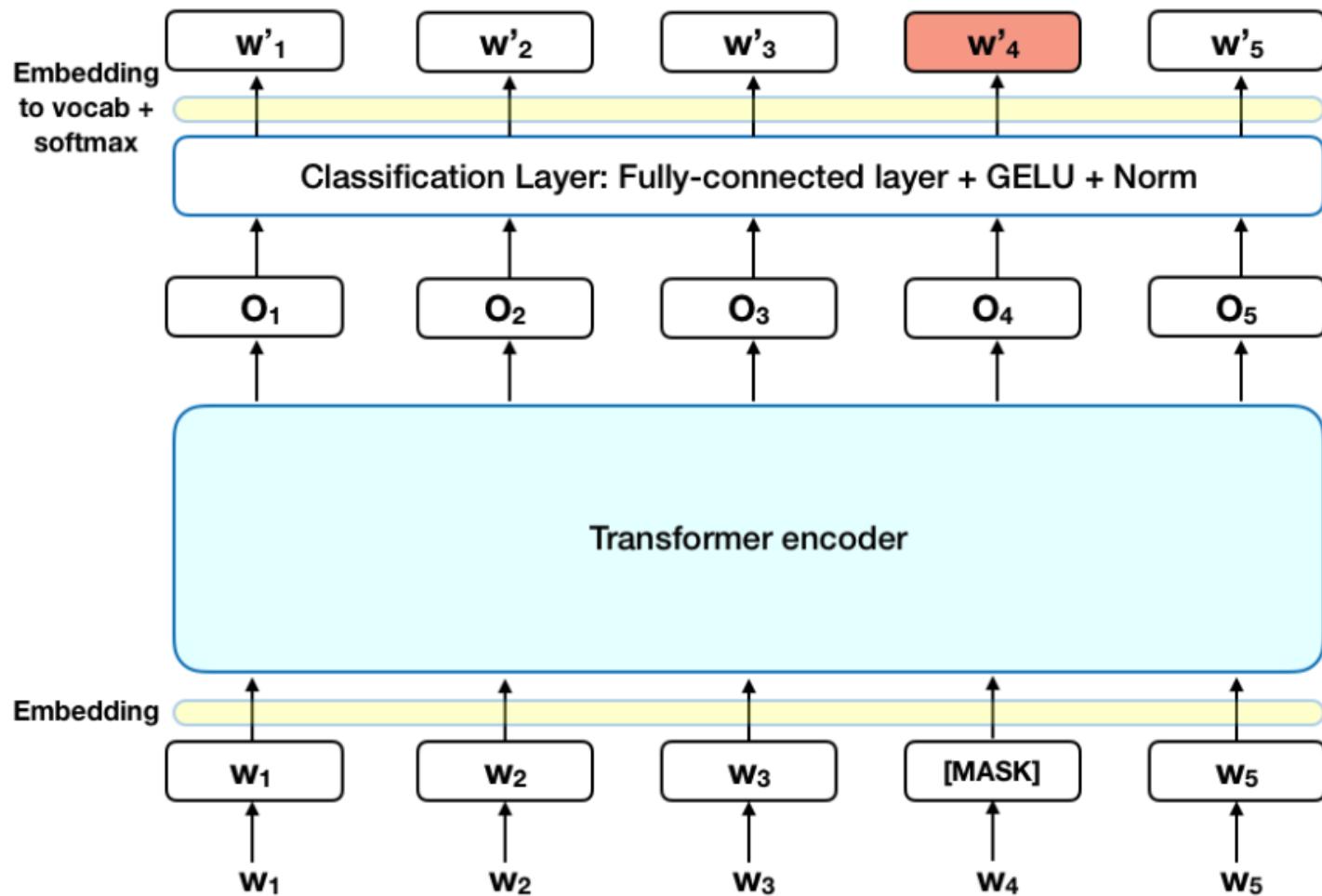
Use the output of the masked word's position to predict the masked word

Randomly mask 15% of tokens

Input



BERT



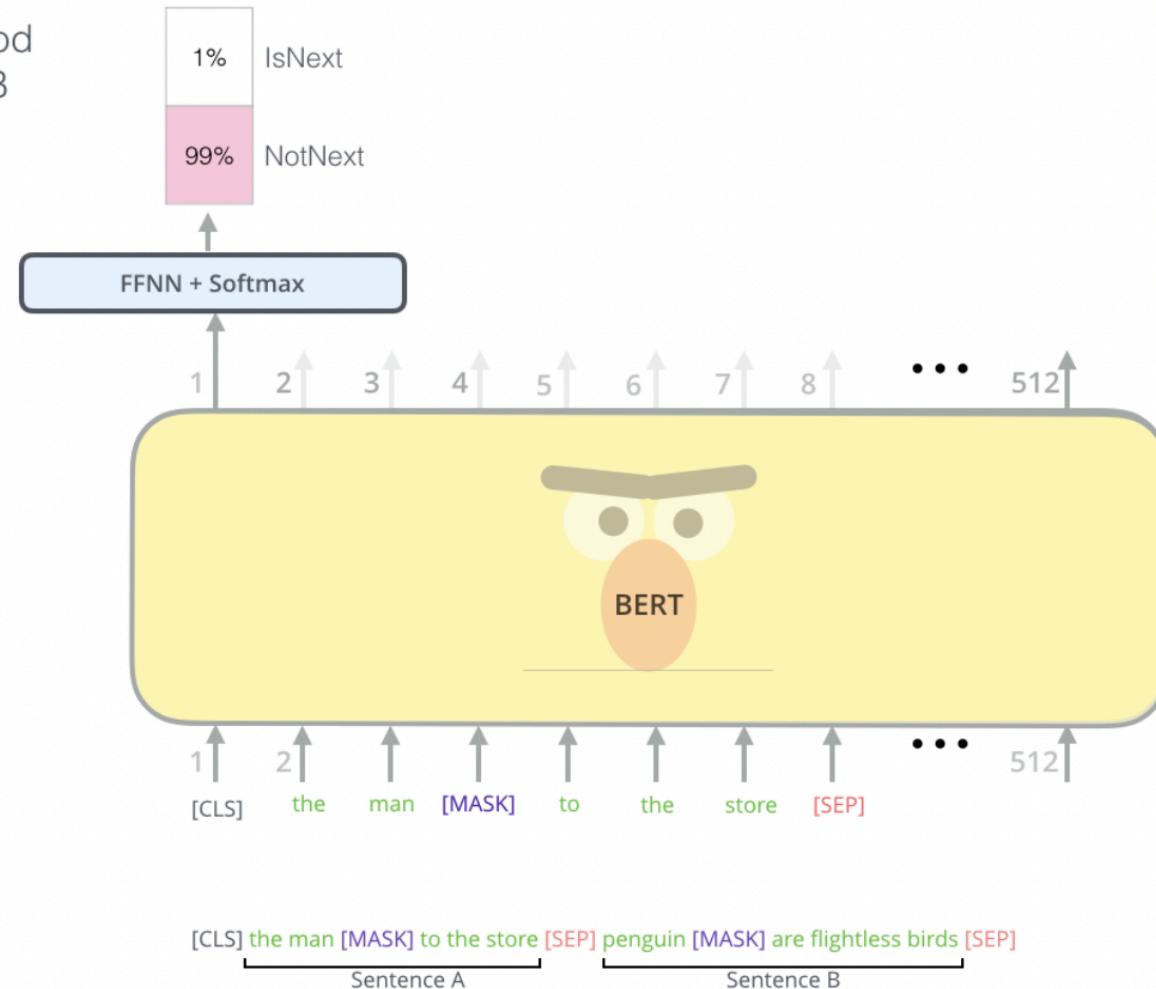
[1810.04805] BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

BERT

Predict likelihood
that sentence B
belongs after
sentence A

Tokenized
Input

Input



The second task BERT is pre-trained on is a two-sentence classification task. The tokenization is oversimplified in this graphic as BERT actually uses WordPieces as tokens rather than words --- so some words are broken down into smaller chunks.



Byte Pair Encoding

this is the hugging face course. this chapter is about tokenization. this section shows several tokenizer algorithms.

this is the hugging face course . this chapter is about tokenization . this section shows several tokenizer algorithms .

3x this 2x is 1x the
1x hugging 1x face
1x course 3x . 1x chapter
1x about 1x tokenization
1x section 1x shows
1x several 1x tokenizer
1x algorithms

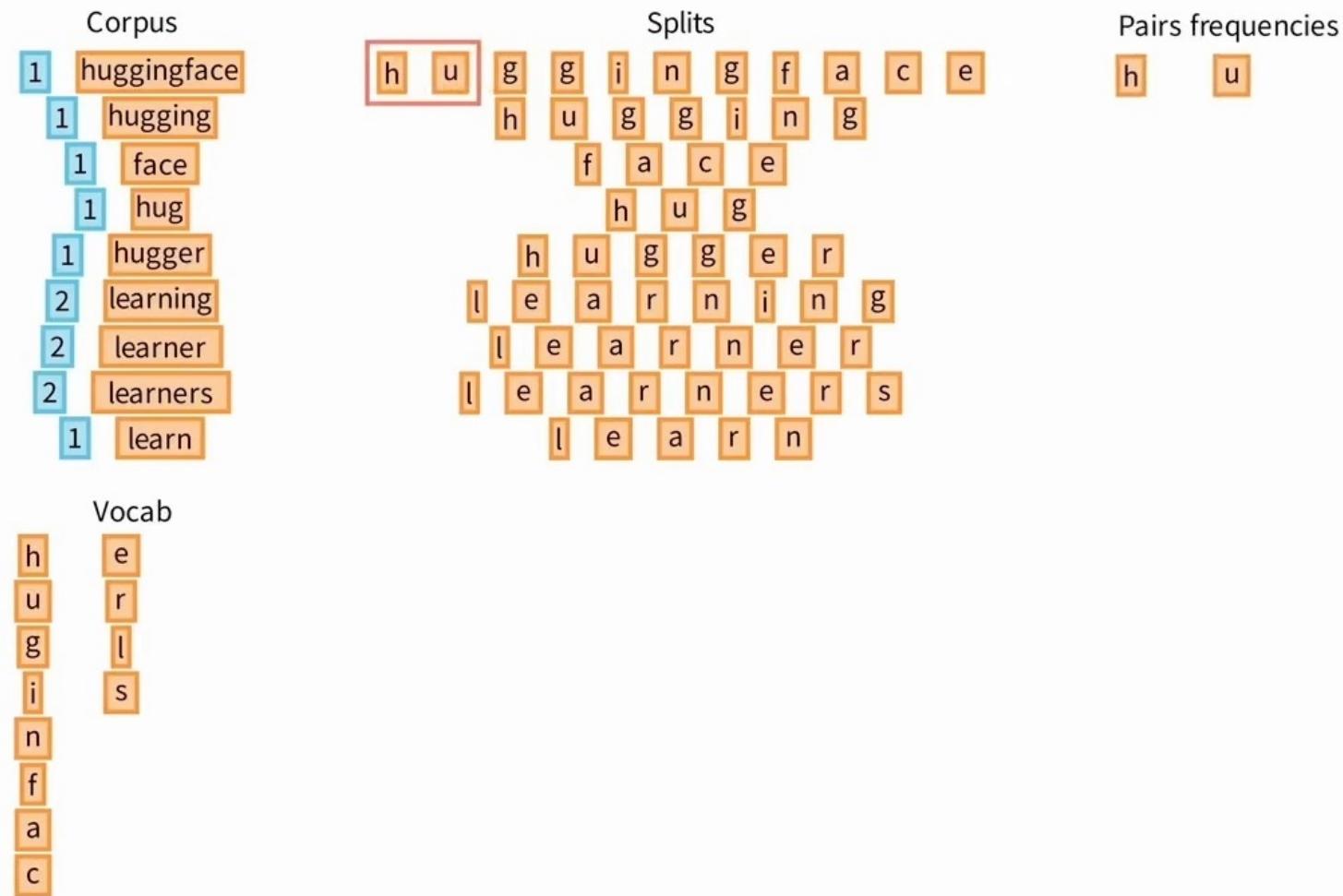
Corpus
1 huggingface
1 hugging
1 face
1 hug
1 hugger
2 learning
2 learner
2 learners
1 learn



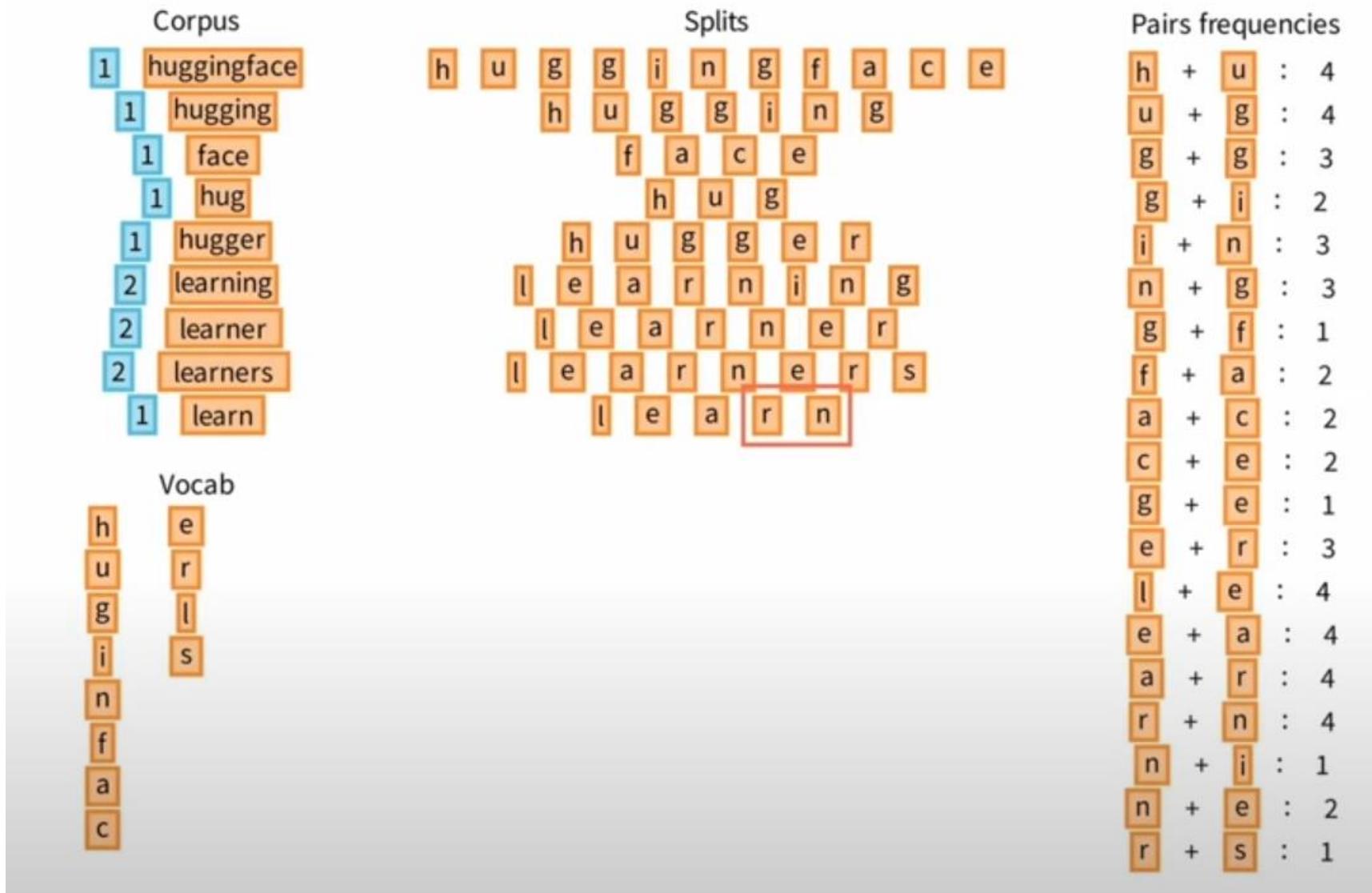
Vocab

h	e
u	r
g	l
i	s
n	
f	
a	
c	

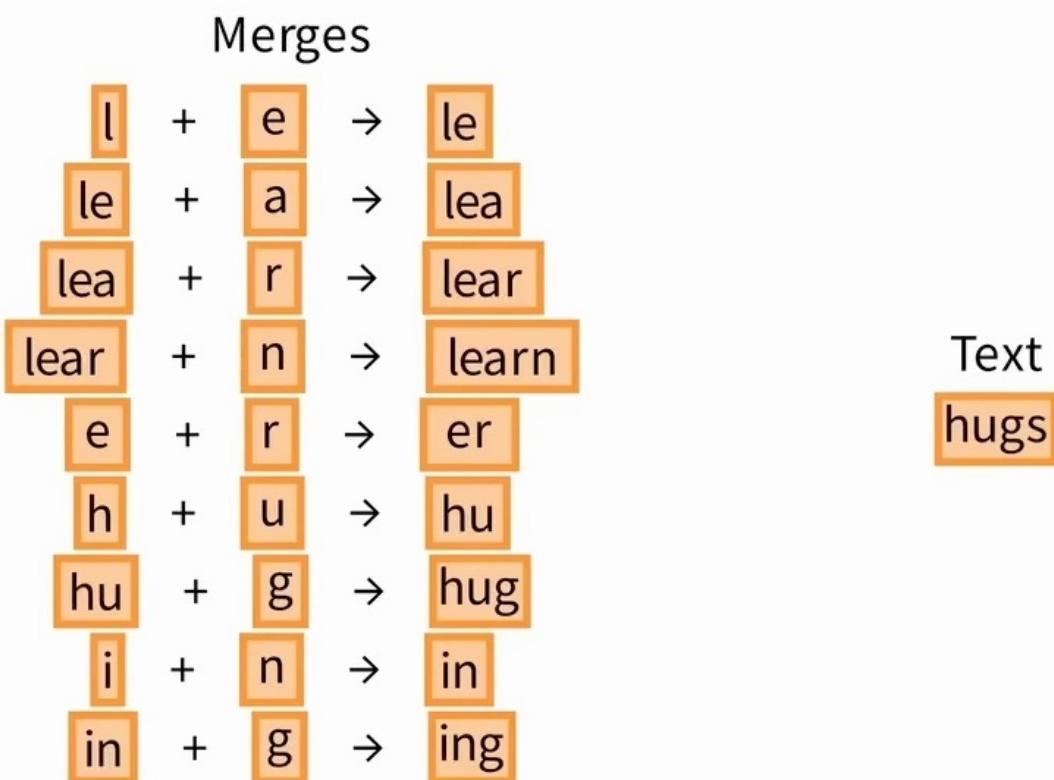
Byte Pair Encoding



Byte Pair Encoding



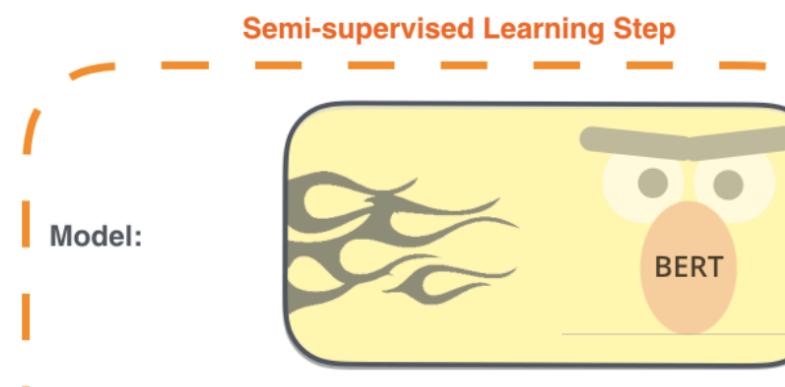
Byte Pair Encoding



BERT

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.



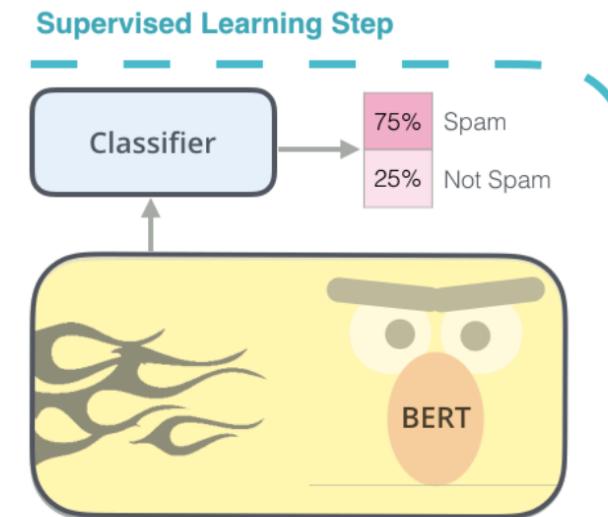
Dataset:



Predict the masked word
(language modeling)

Objective:

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

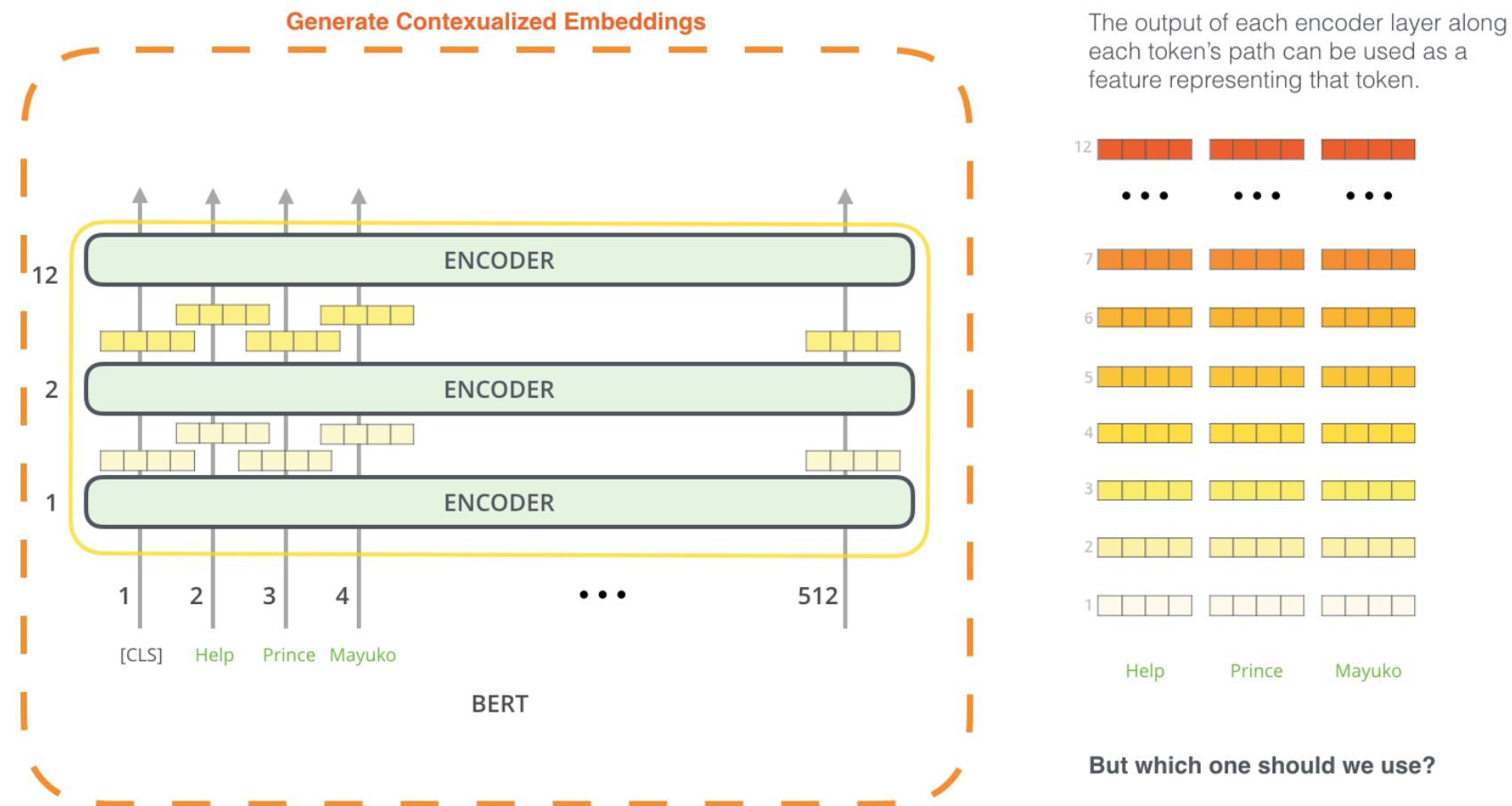


Model:
(pre-trained
in step #1)

Email message	Class
Buy these pills	Spam
Win cash prizes	Spam
Dear Mr. Atreides, please find attached...	Not Spam

Dataset:

BERT for feature extraction



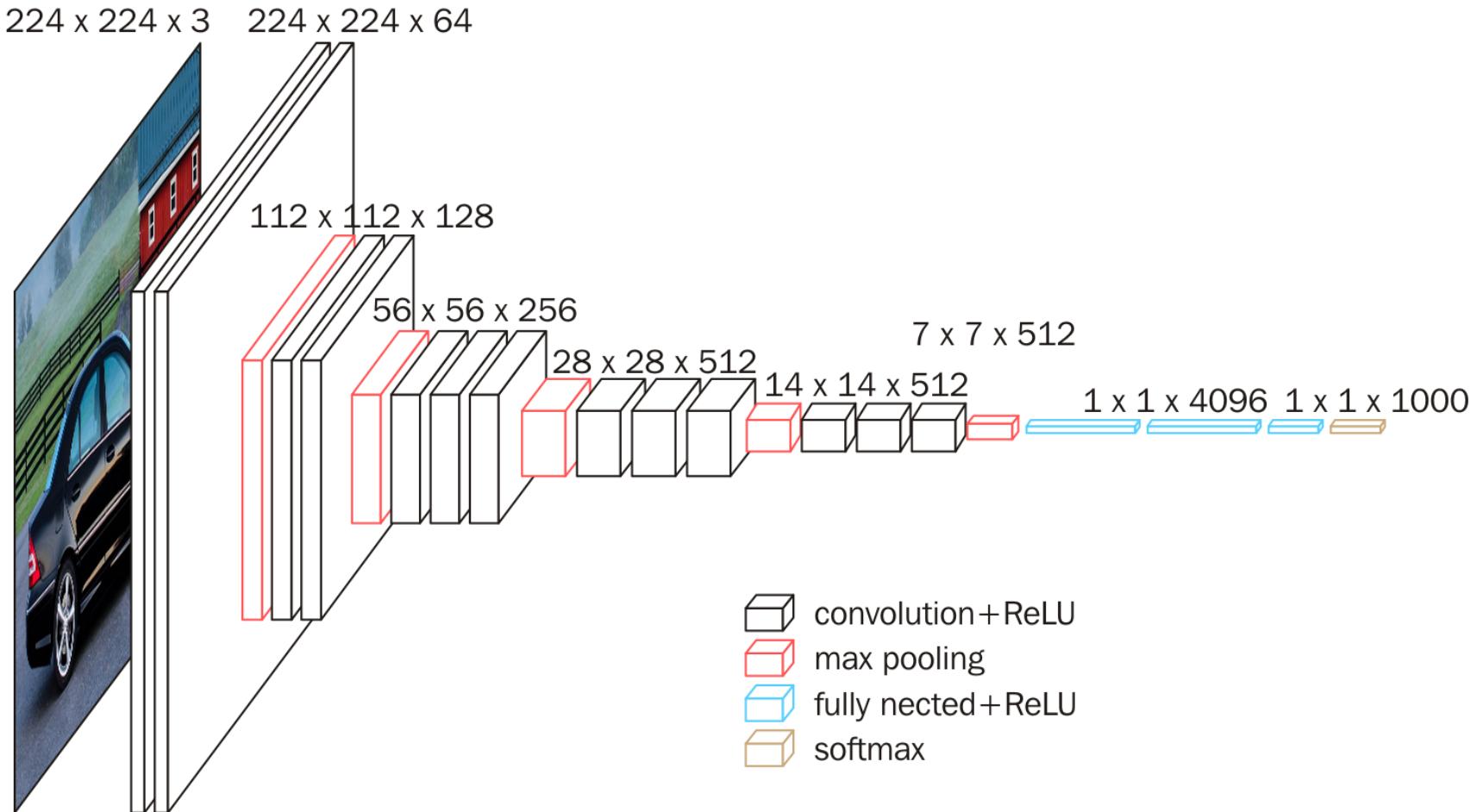
BERT for feature extraction

What is the best contextualized embedding for “Help” in that context?
For named-entity recognition task CoNLL-2003 NER

		Dev F1 Score
12	First Layer	91.0
• • •	Last Hidden Layer	94.9
7		
6	Sum All 12 Layers	95.5
5		
4	Second-to-Last Hidden Layer	95.6
3		
2		
1	Sum Last Four Hidden	95.9
Help	Concat Last Four Hidden	96.1

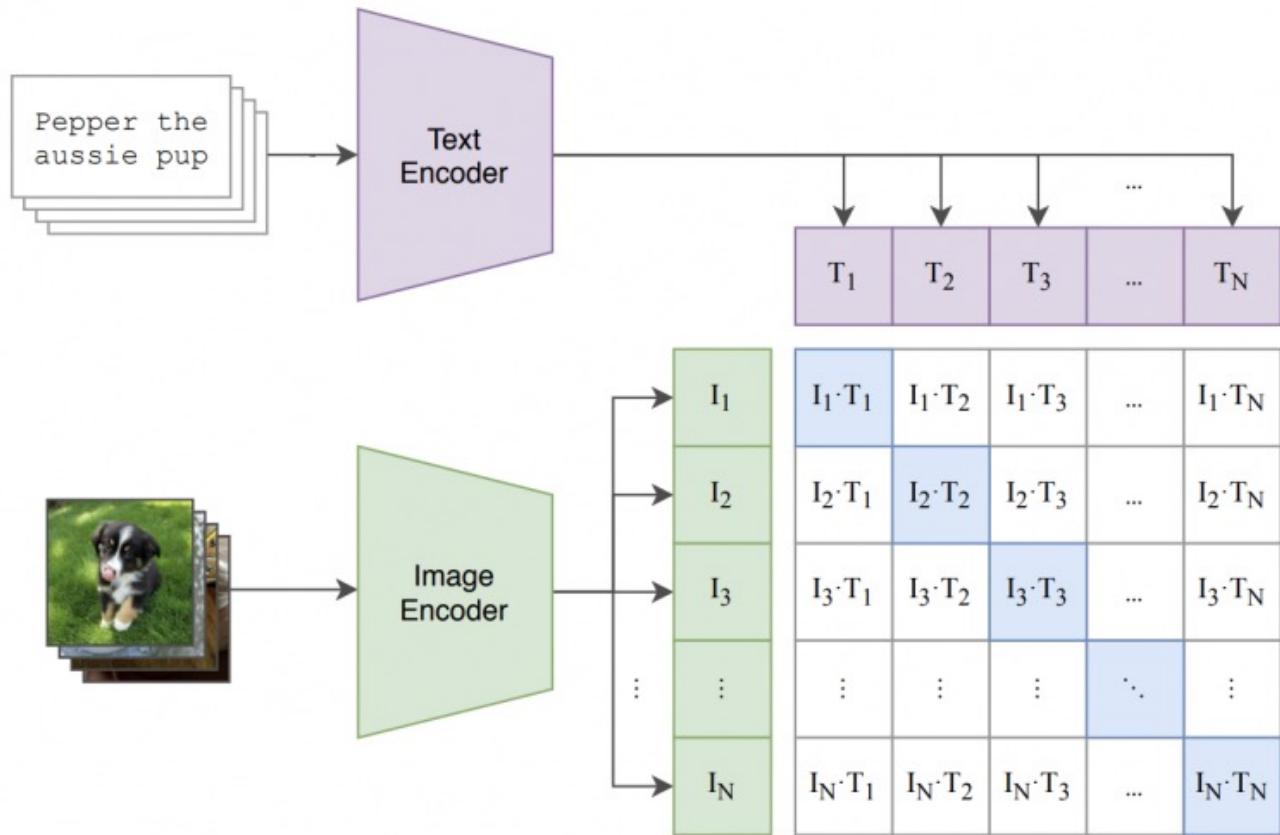
The figure shows the Dev F1 Scores for different BERT feature extraction methods. The methods are listed from top to bottom: First Layer, Last Hidden Layer, Sum All 12 Layers, Second-to-Last Hidden Layer, Sum Last Four Hidden, and Concat Last Four Hidden. Each method is accompanied by a diagram showing the dimensionality of the embeddings at each layer. The 'First Layer' has 4 red squares. The 'Last Hidden Layer' has 4 red squares. The 'Sum All 12 Layers' method shows a stack of 12 layers: layer 12 (4 red), layer 11 (4 red), layer 10 (4 orange), and layer 9 (4 orange). Below these are three intermediate layers (2, 3, 4) each with 4 yellow squares, followed by a layer 1 with 4 white squares. The 'Second-to-Last Hidden Layer' has 4 red squares. The 'Sum Last Four Hidden' method shows a stack of 4 layers: layer 12 (4 red), layer 11 (4 red), layer 10 (4 orange), and layer 9 (4 orange). Below these are three intermediate layers (2, 3, 4) each with 4 yellow squares, followed by a layer 1 with 4 white squares. The 'Concat Last Four Hidden' method shows a stack of 4 layers: layer 9 (4 red), layer 10 (4 orange), layer 11 (4 orange), and layer 12 (4 orange).

Recap: Applying transfer learning in CV

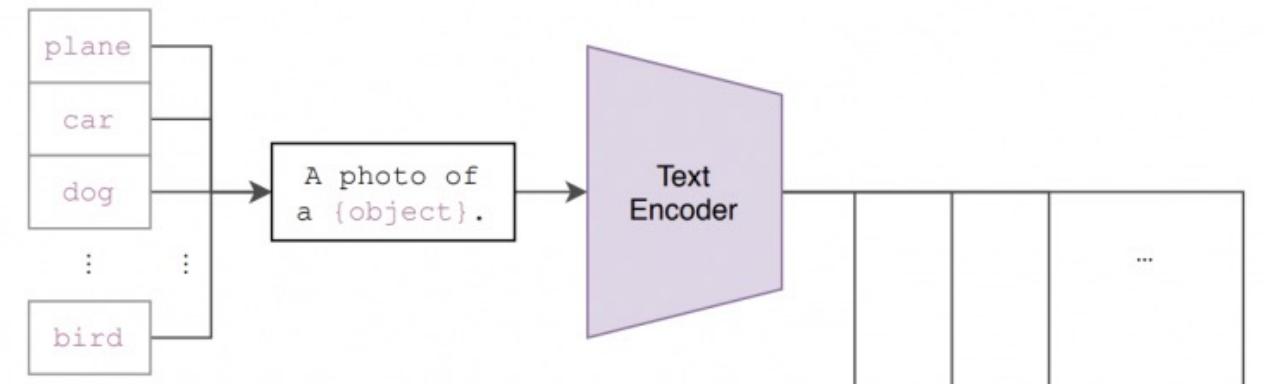


CLIP

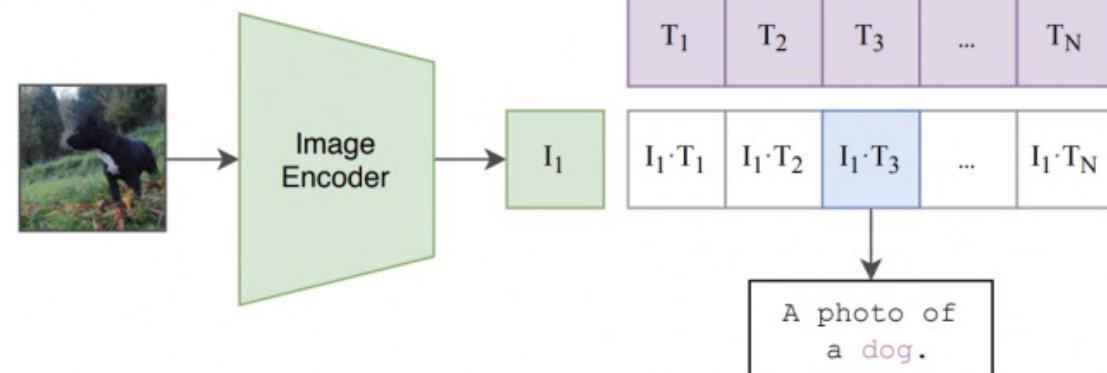
(1) Contrastive pre-training



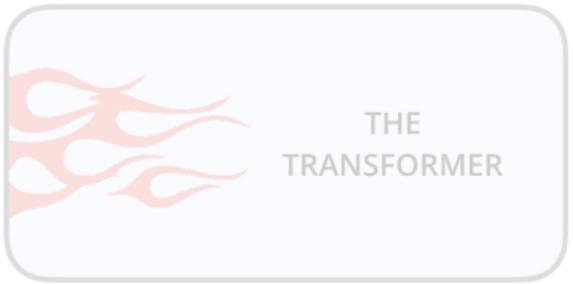
(2) Create dataset classifier from label text



(3) Use for zero-shot prediction



Другие модели

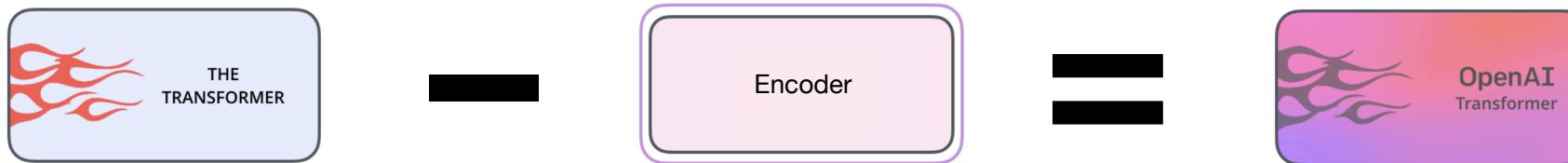


[The Illustrated BERT, ELMo, and co. \(How NLP Cracked Transfer Learning\)](#)

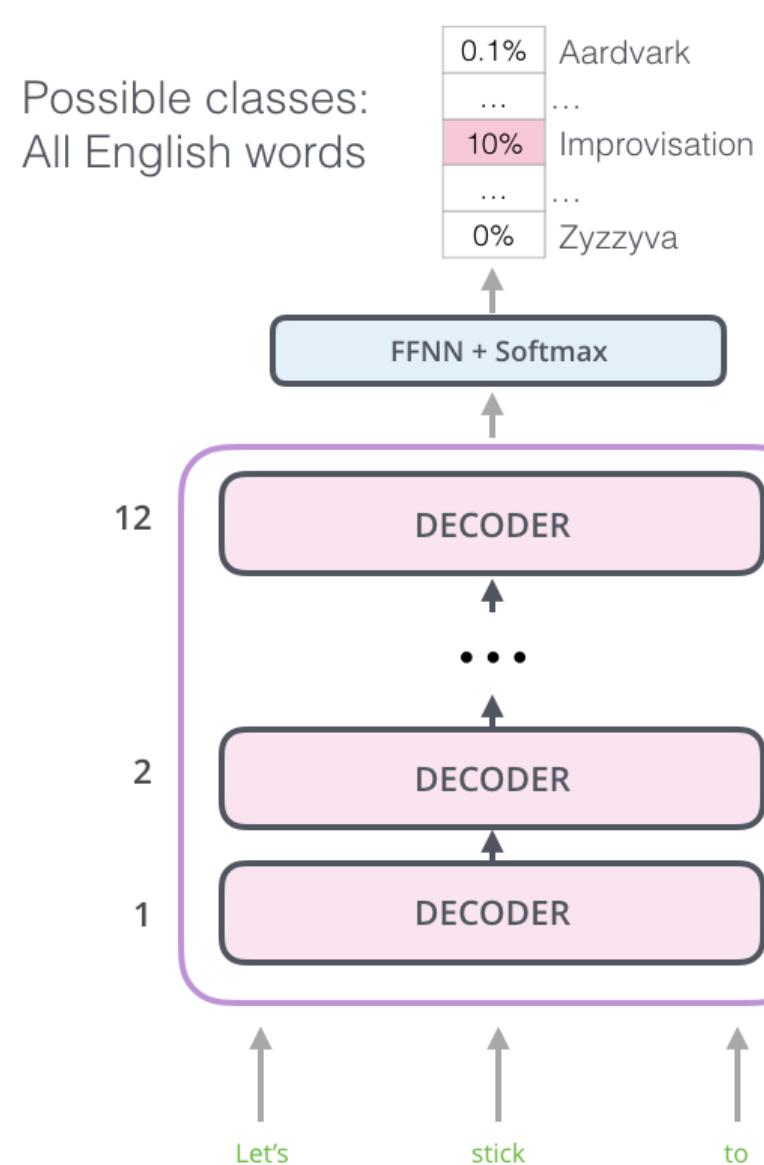
OpenAI Transformer



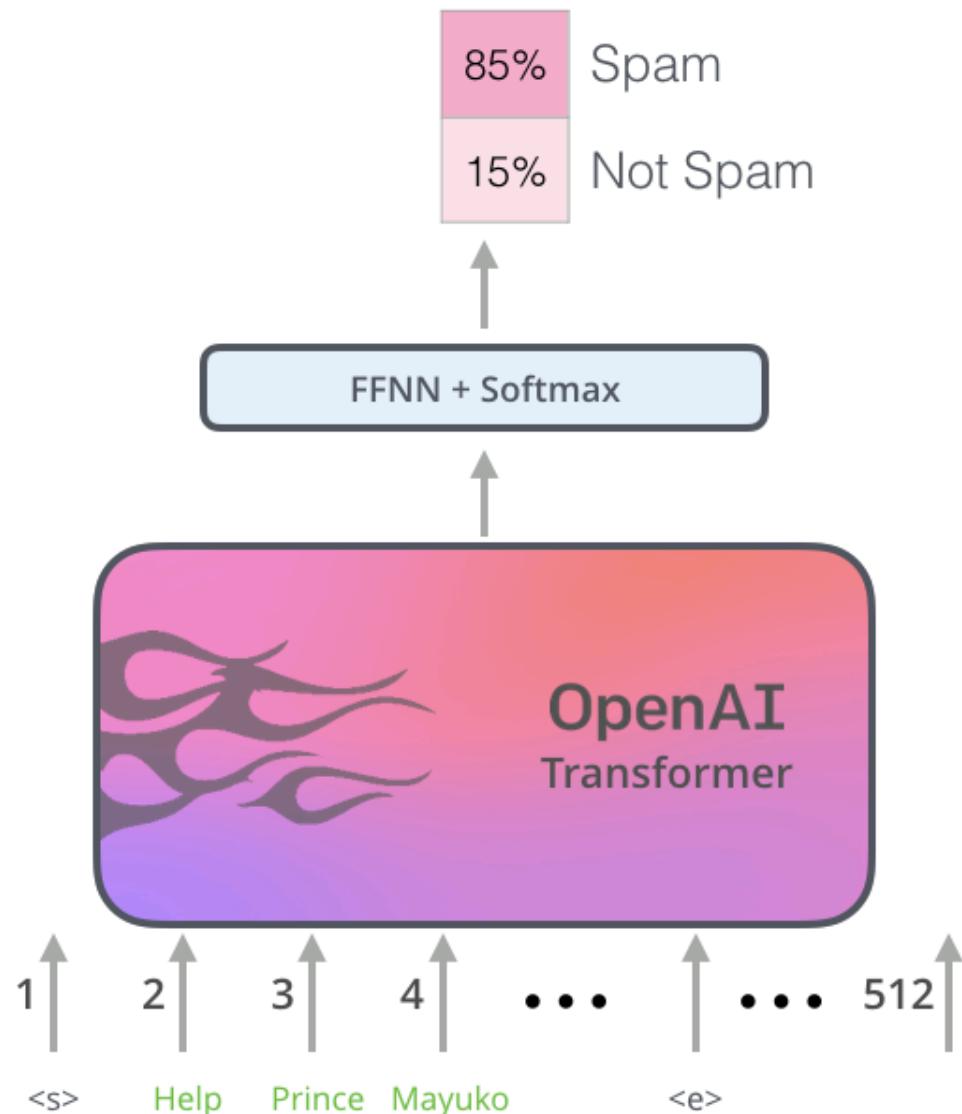
Pre-training a Transformer Decoder for
Language Modeling



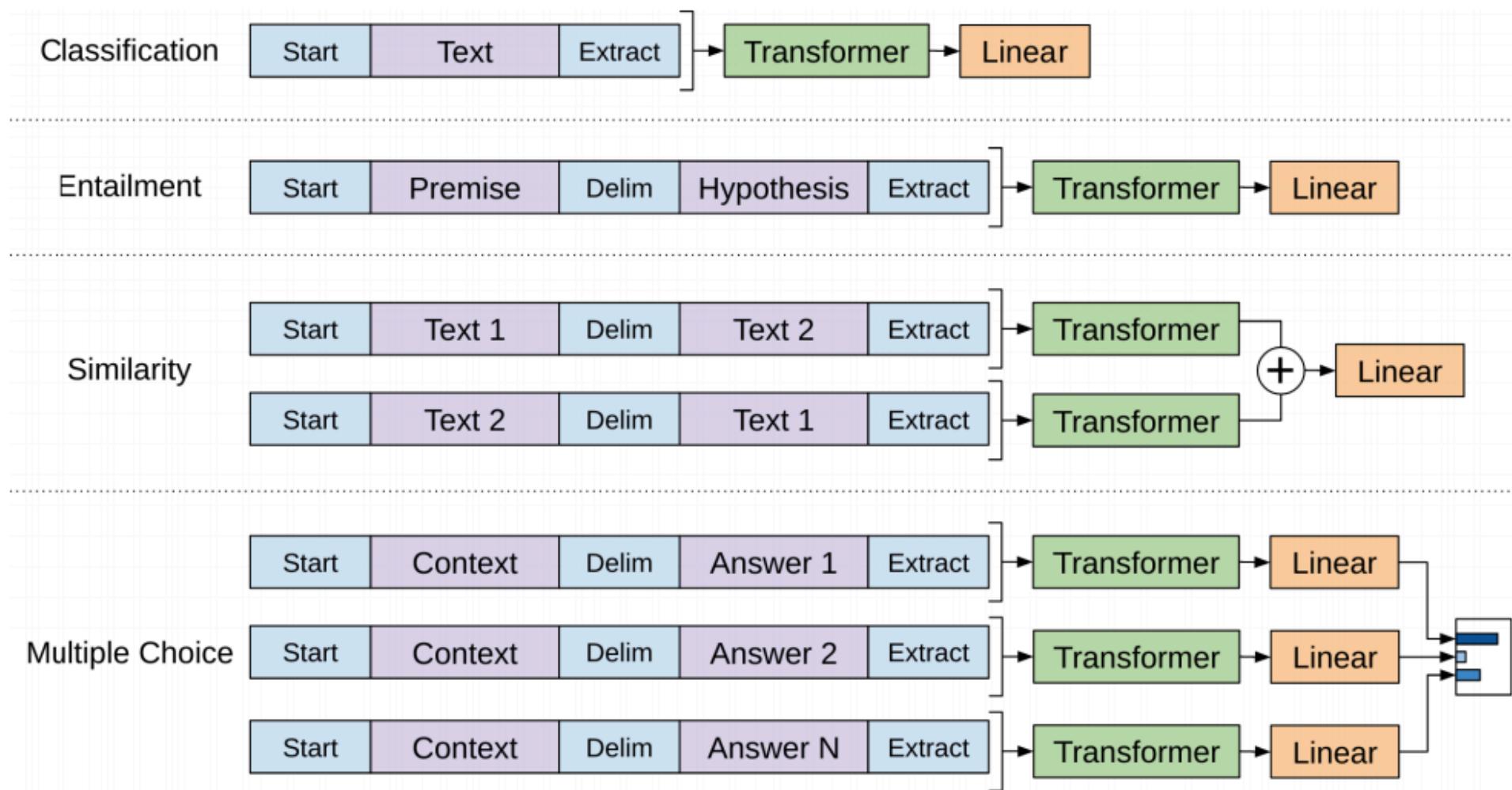
OpenAI Transformer



OpenAI Transformer



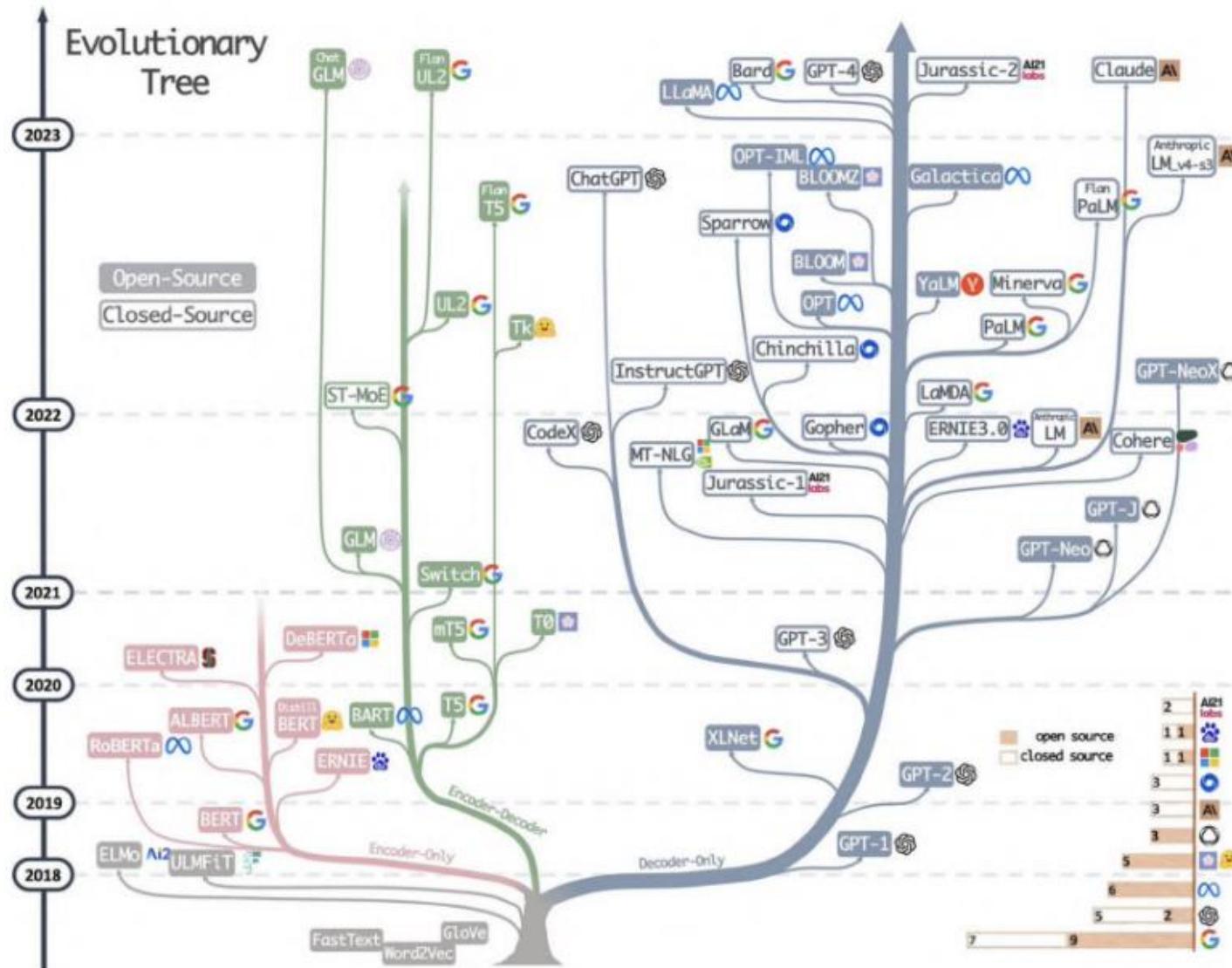
OpenAI Transformer



Сейчас «LLM» \approx «AI»

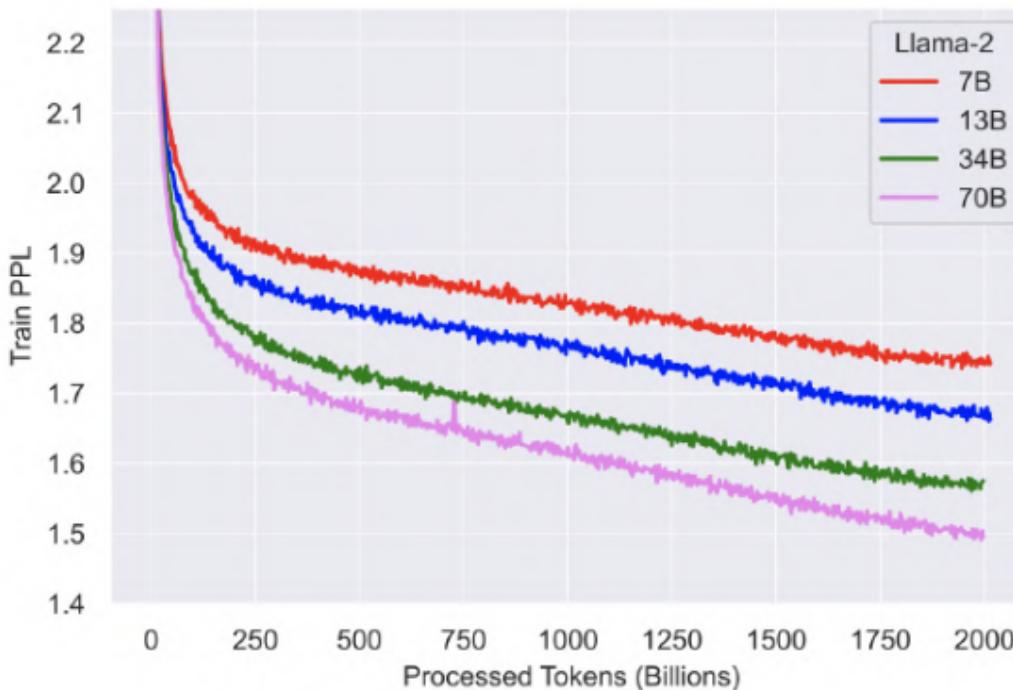
Современные AI-сервисы:

- ▶ **общаются с людьми на любые темы на естественном языке**
- ▶ **решают школьные и университетские задачи по разным дисциплинам**
- ▶ **понимают и генерируют тексты, изображения, аудио**



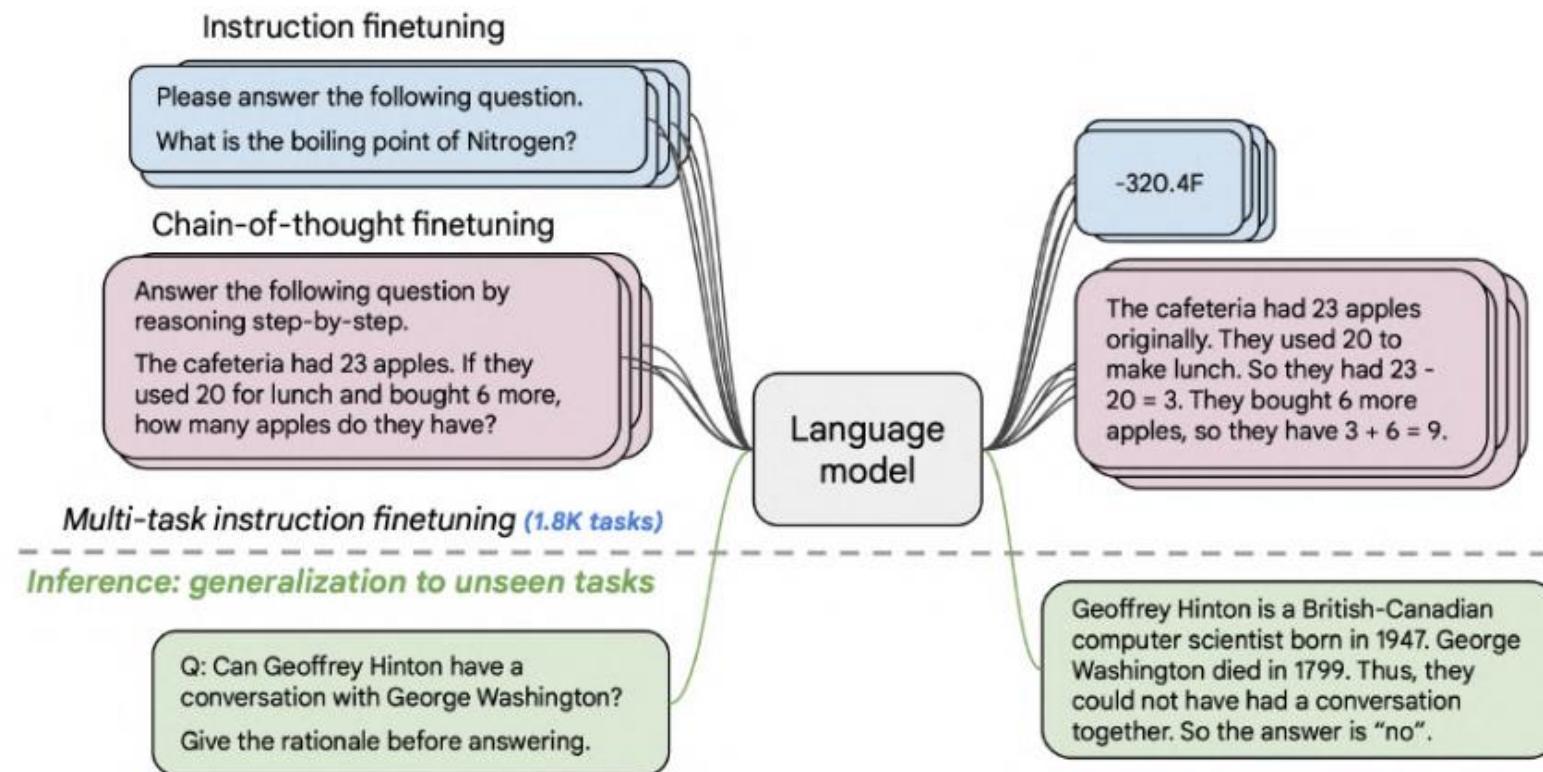
Этапы обучения

- ▶ Первая стадия обучения LLM — предобучение (pre-train)
- ▶ Модель учится предсказывать следующий токен по контексту слева
- ▶ Если учить с Teacher Forcing — контекст берётся из обучения, если без — из того, что сгенерировала в процессе сама модель (комбинируют)
- ▶ На этом этапе приобретает основные знания о языке и мире



Этапы обучения

- ▶ Вторая стадия — Instruction Tuning (SFT)
- ▶ Модель учится понимать и исполнять запросы людей на естественном языке и вести диалоги
- ▶ Например: pre-train — LLaMA, instruct-tuned — Alpaca или Vicuna



Этапы обучения

- ▶ Третий, optionalный, шаг — выравнивание (alignment)
- ▶ Диалоговая модель дообучается для генерации более корректных, полезных и безопасных ответов
- ▶ Популярная техника, использованная в Instruct GPT — RLHF:
 - ▶ обученная LLM генерирует на тестовом наборе инструкций ответы
 - ▶ ответы размечаются ассессорами, на их ответах учится сильная reward-модель, она оценивает по тексту его качество
 - ▶ заводятся две копии модели (A) и (B), учится (A)
 - ▶ обе модели генерируют ответы на каждый промпт, ответ (A) оценивается reward-моделью
 - ▶ веса (A) обновляются так, чтобы максимизировать reward и не давать ответы, очень далёкие от исходной (B)
 - ▶ расстояние определяется по KL-дивергенции между выходными распределениями моделей
 - ▶ обновление весов идёт по заданному алгоритму (PPO или A2C)

Projections of the stock of public text and data usage



Effective stock (number of tokens)

