

# Глубокое обучение и вообще

Бекезин Никита

8 декабря 2021 г.

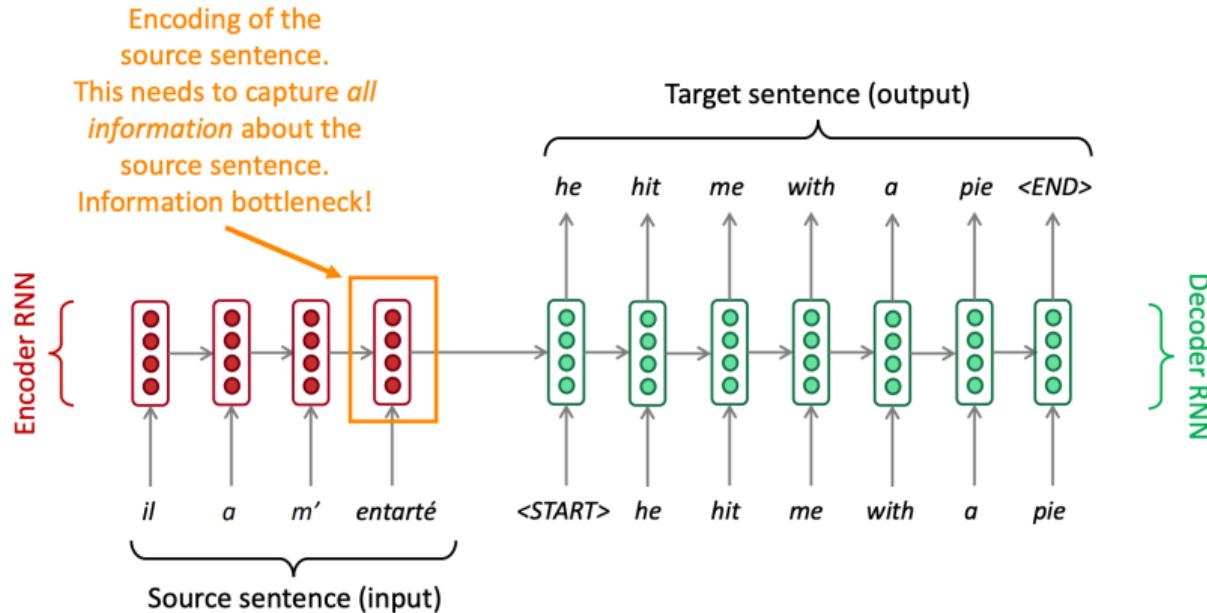
**Лекция 14:** Быстрое введение в SOTA 2

# Agenda

- Seq2seq recap
- Attention recap
- Transformer
- Self-Attention
- BERT

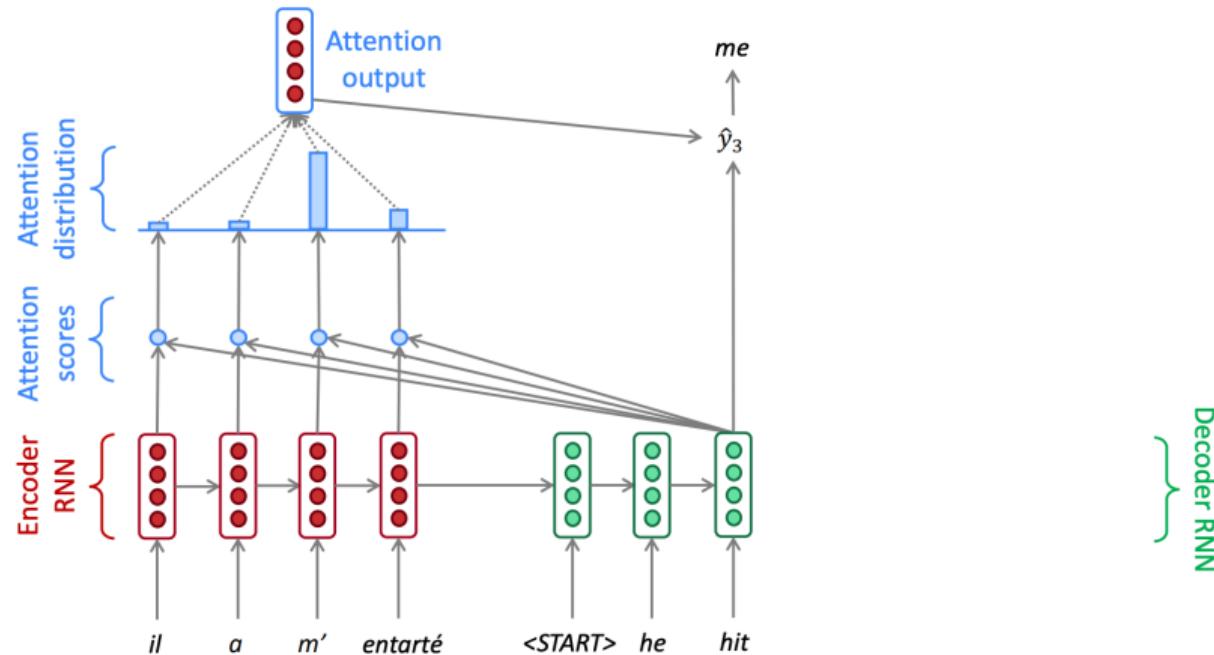
# Seq2seq recap

## Sequence-to-sequence: the bottleneck problem



# Attention recap

## Sequence-to-sequence with attention



# Attention recap

There are **several ways** you can compute  $e \in \mathbb{R}^N$  from  $\mathbf{h}_1, \dots, \mathbf{h}_N \in \mathbb{R}^{d_1}$  and  $\mathbf{s} \in \mathbb{R}^{d_2}$ :

- **Basic dot-product attention:**  $e_i = \mathbf{s}^T \mathbf{h}_i \in \mathbb{R}$ 
  - Note: this assumes  $d_1 = d_2$
  - This is the version we saw earlier
- **Multiplicative attention:**  $e_i = \mathbf{s}^T \mathbf{W} \mathbf{h}_i \in \mathbb{R}$ 
  - Where  $\mathbf{W} \in \mathbb{R}^{d_2 \times d_1}$  is a weight matrix
- **Additive attention:**  $e_i = \mathbf{v}^T \tanh(\mathbf{W}_1 \mathbf{h}_i + \mathbf{W}_2 \mathbf{s}) \in \mathbb{R}$ 
  - Where  $\mathbf{W}_1 \in \mathbb{R}^{d_3 \times d_1}$ ,  $\mathbf{W}_2 \in \mathbb{R}^{d_3 \times d_2}$  are weight matrices and  $\mathbf{v} \in \mathbb{R}^{d_3}$  is a weight vector.
  - $d_3$  (the attention dimensionality) is a hyperparameter

# Attention recap

- “Free” word alignment
- Better results on long sequences with attention

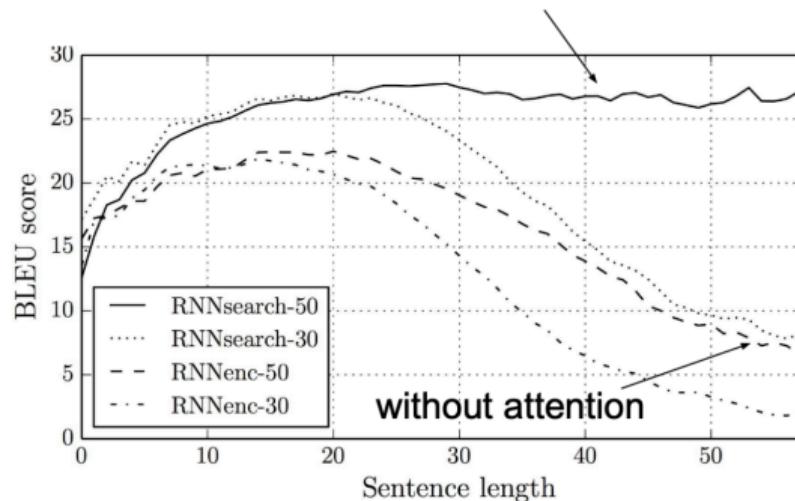
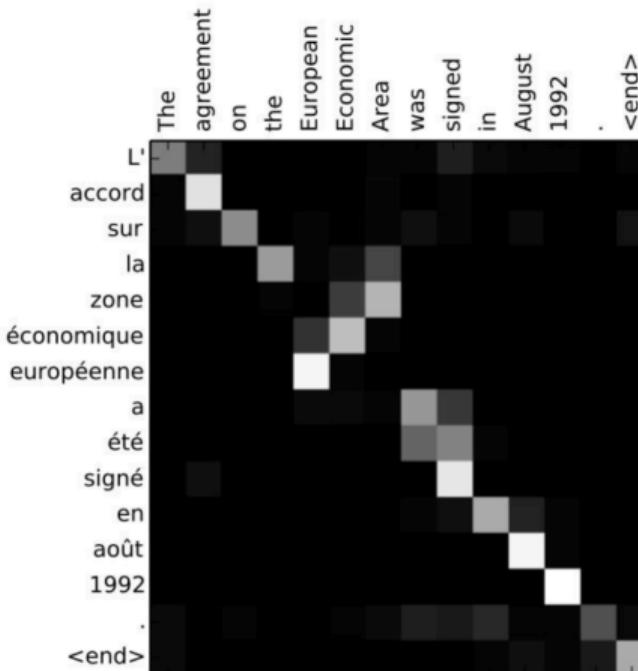


Image source: [Neural Machine Translation by Jointly Learning to Align and Translate](#)

## Attention advantages

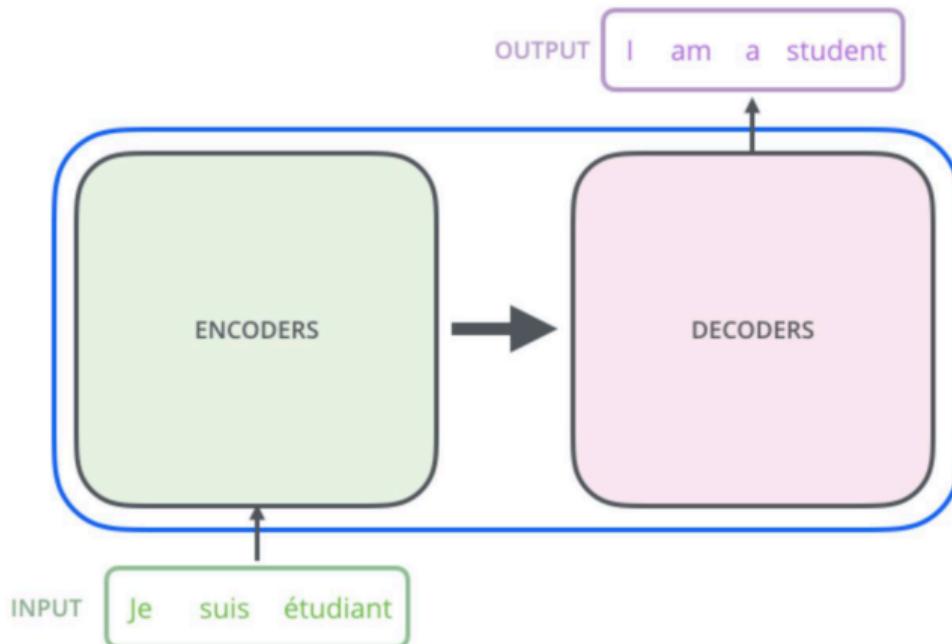


Attention is all you need!

# attention is all you need

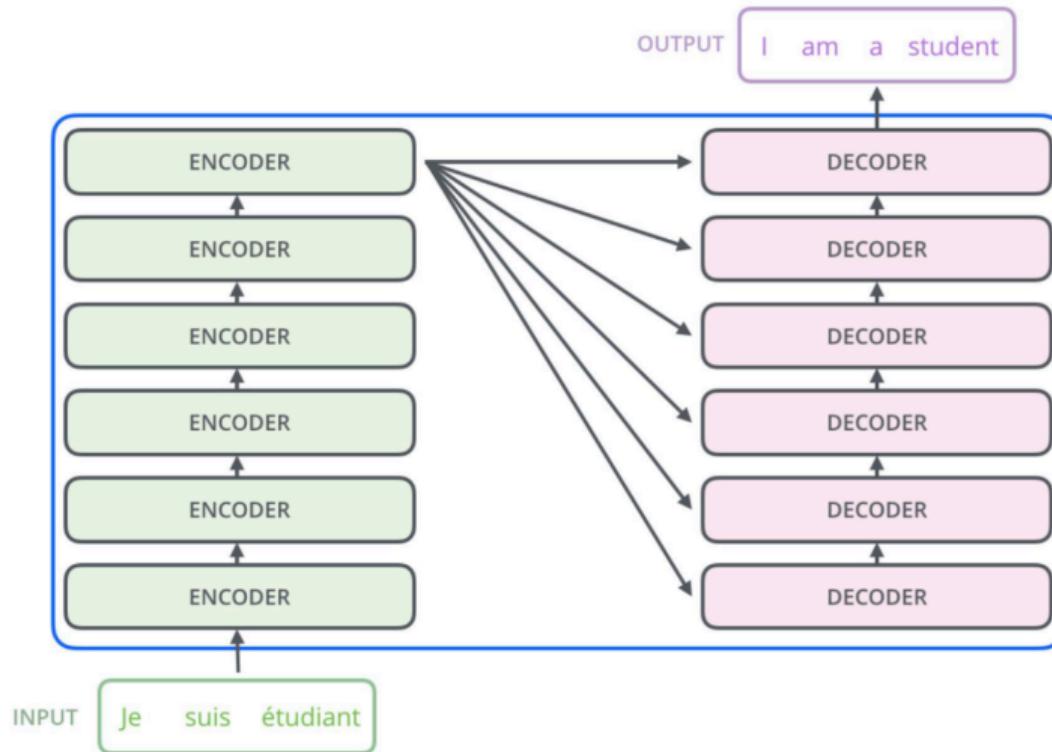
Развитие идеи внимания. Статья вышла в 2017 году и стала мамой всех текущих SOTA моделей. А зачем нам вообще что-то, кроме внимания? Давайте напихаем в энкодер и декодер как можно больше внимания и будем такой штукой его учить.

# Transformer



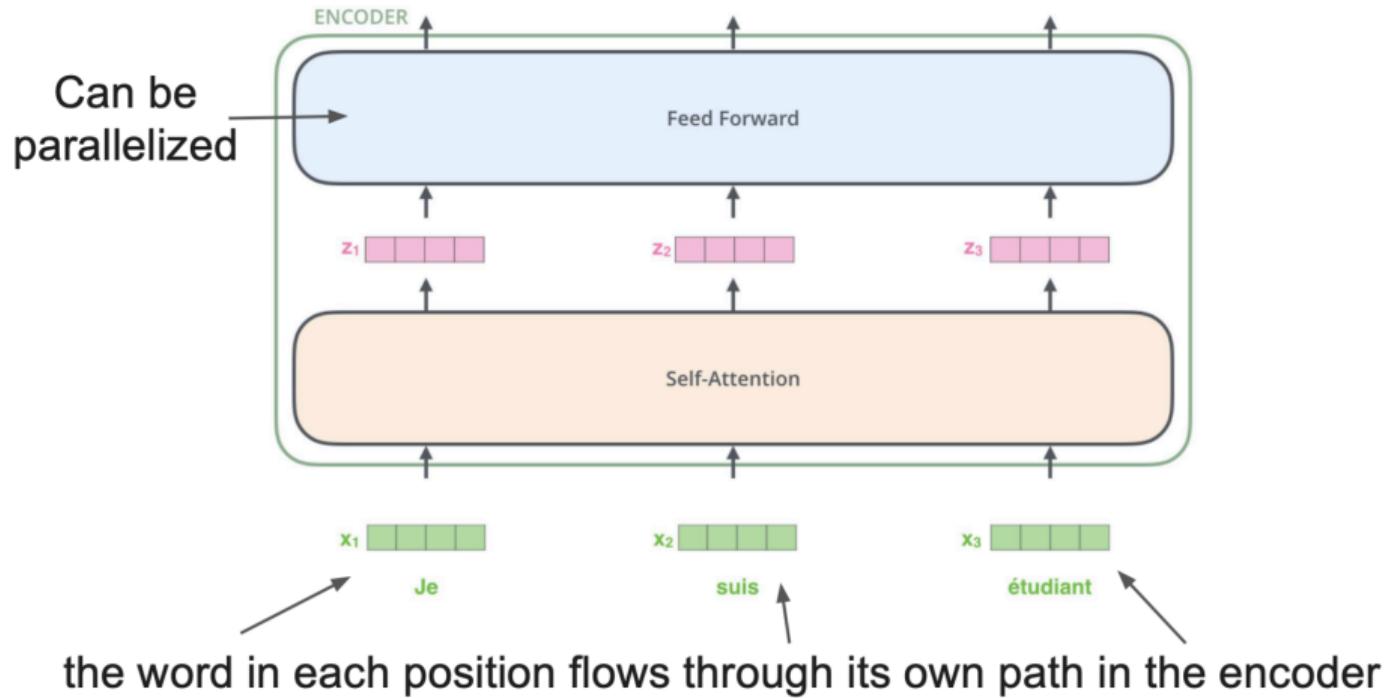
Source

# Transformer



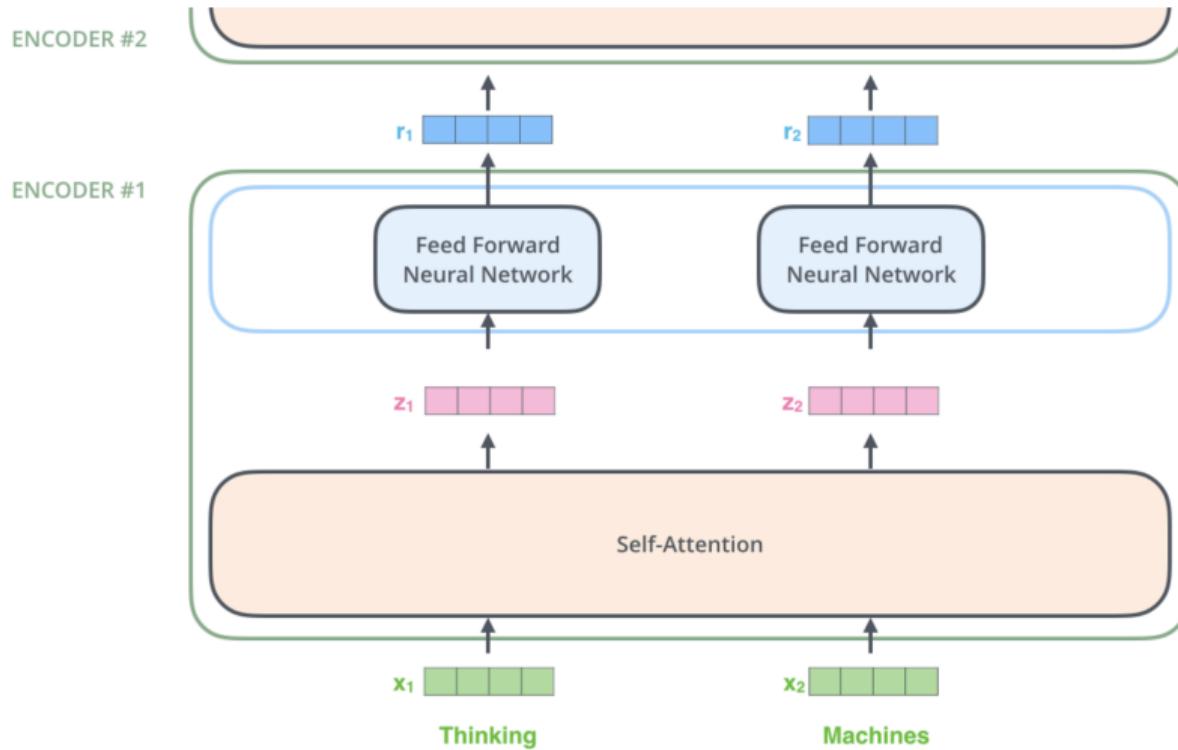
Source

# Transformer



Source

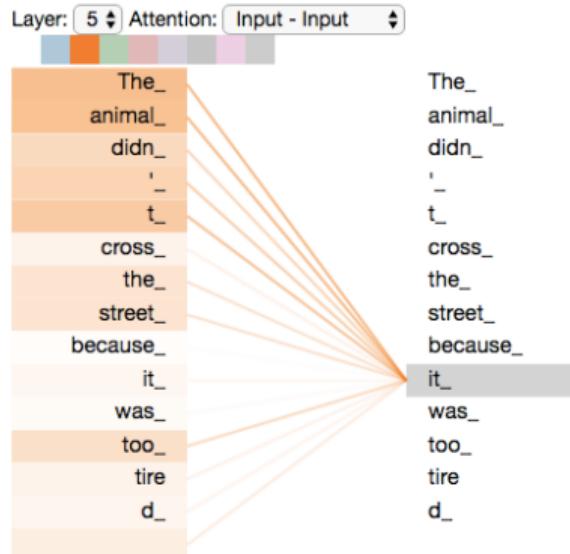
# Transformer



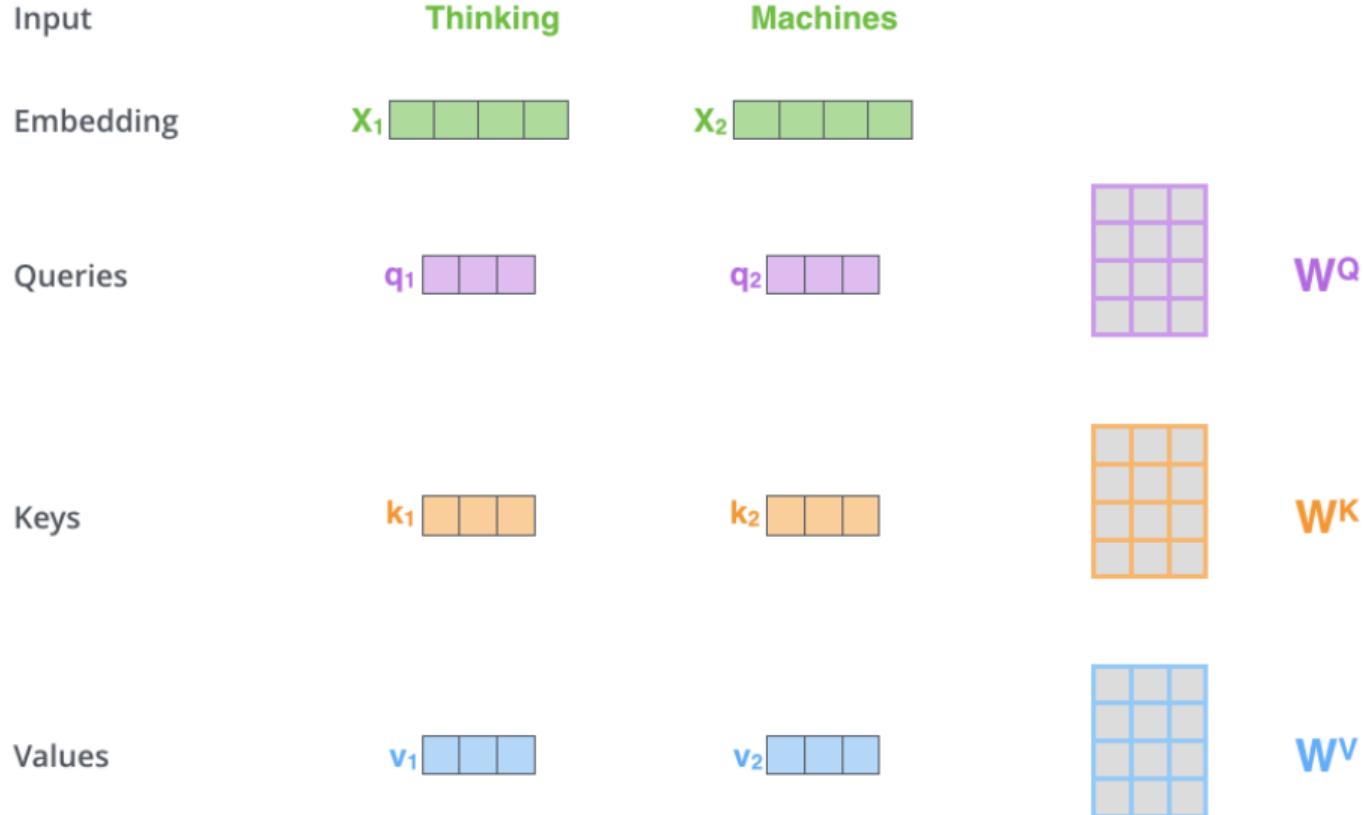
Source

# Self-attention или механизм внимания

Есть предложение: "The animal didn't cross the street because it was too tired"



# Шаг 1. Абстракции

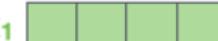
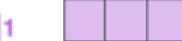
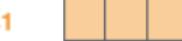
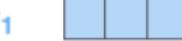


# Self-attention

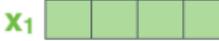
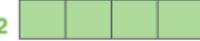
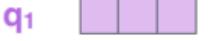
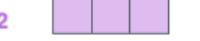
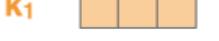
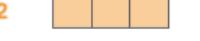
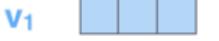
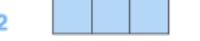
А теперь тоже самое, но словами:

1. Query, key - ищем связи между словами. Ходим по всем со всеми смотрим насколько они связаны. Query - мое текущее слово, key - мое слово с которым я сравниваю себя.
2. Value - то, что мы знаем об этом слове

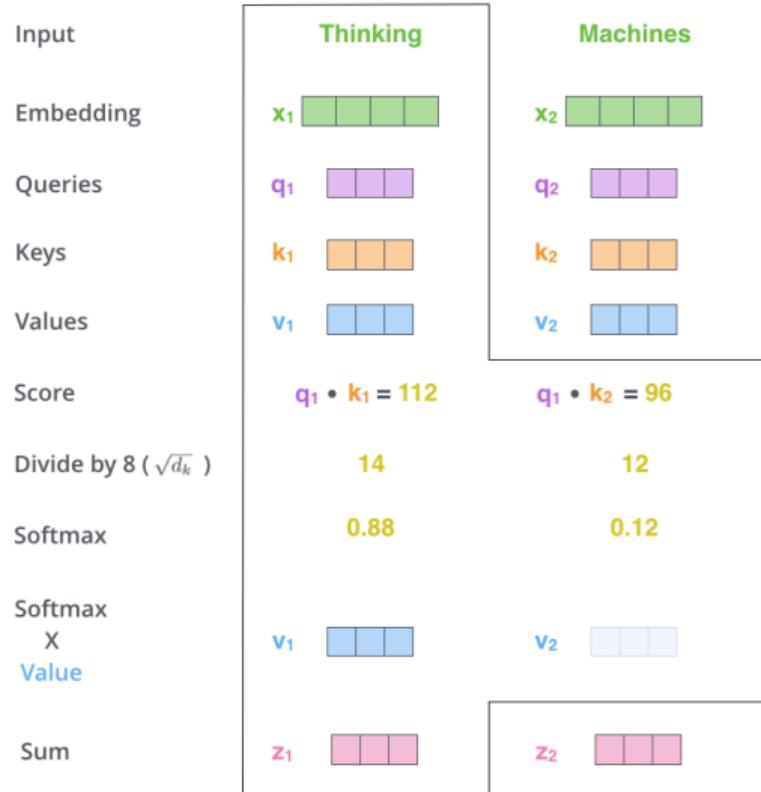
## Шаг 2. Расчет attention score

Input		
Embedding	$x_1$	
Queries	$q_1$	
Keys	$k_1$	
Values	$v_1$	
Score	$q_1 \cdot k_1 = 112$	$q_1 \cdot k_2 = 96$

## Шаг 3. Шкалирование и Softmax

Input	Thinking		Machines	
Embedding	$x_1$		$x_2$	
Queries	$q_1$		$q_2$	
Keys	$k_1$		$k_2$	
Values	$v_1$		$v_2$	
Score	$q_1 \cdot k_1 = 112$		$q_1 \cdot k_2 = 96$	
Divide by 8 ( $\sqrt{d_k}$ )	14		12	
Softmax	0.88		0.12	

## Шаг 4. Сумма взвешенных векторов Value



# Матричные перемножения

$$X \times W^Q = Q$$

A diagram illustrating matrix multiplication. On the left, a green 3x4 matrix labeled 'X' is multiplied by a purple 4x4 matrix labeled 'W<sup>Q</sup>'. The result is a purple 3x3 matrix labeled 'Q'. The matrices are represented as grids of colored squares.

$$X \times W^K = K$$

A diagram illustrating matrix multiplication. On the left, a green 3x4 matrix labeled 'X' is multiplied by an orange 4x4 matrix labeled 'W<sup>K</sup>'. The result is an orange 3x3 matrix labeled 'K'. The matrices are represented as grids of colored squares.

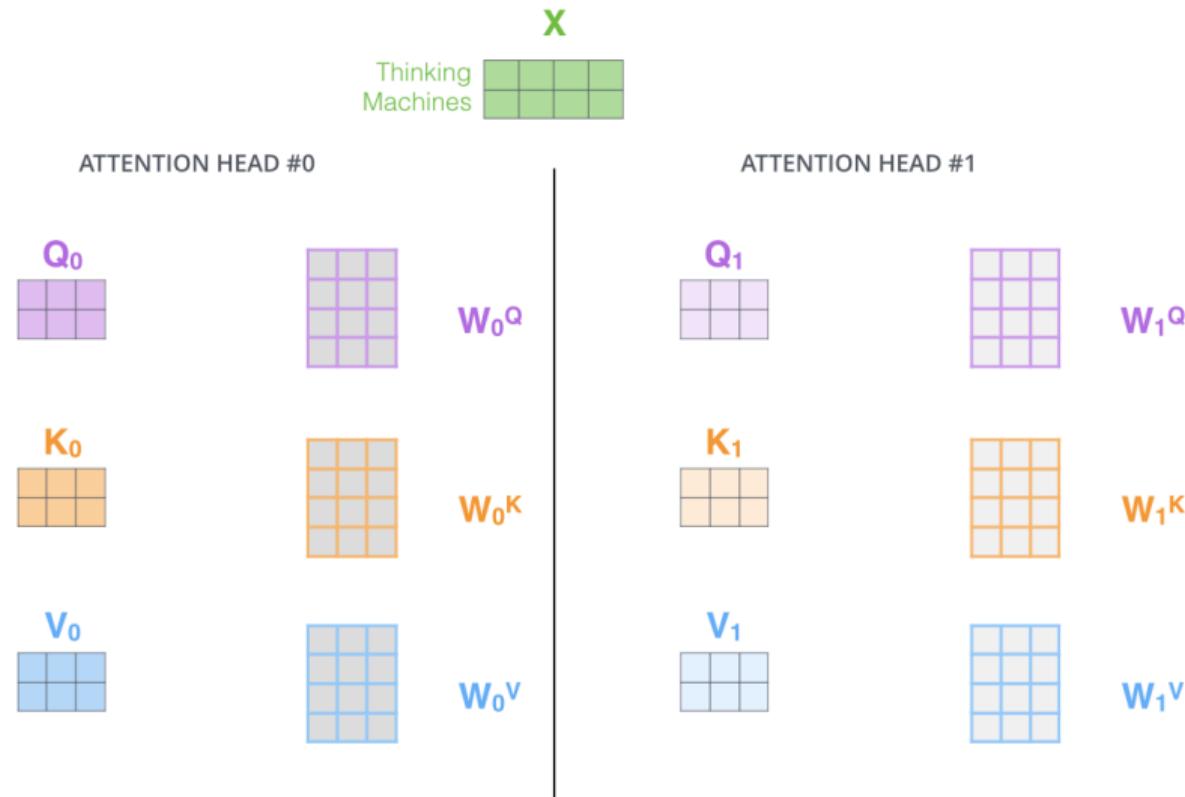
$$X \times W^V = V$$

A diagram illustrating matrix multiplication. On the left, a green 3x4 matrix labeled 'X' is multiplied by a blue 4x4 matrix labeled 'W<sup>V</sup>'. The result is a blue 3x3 matrix labeled 'V'. The matrices are represented as grids of colored squares.

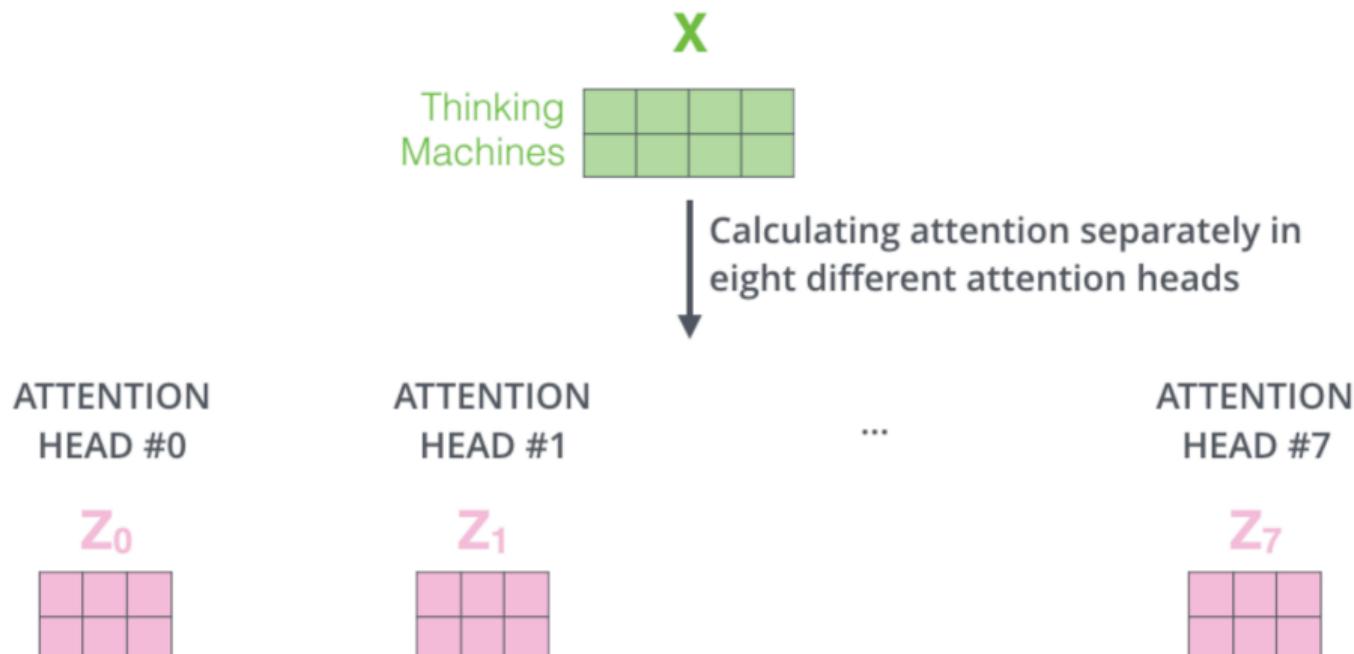
# All together

$$\text{softmax} \left( \frac{\begin{matrix} \mathbf{Q} & \times & \mathbf{K}^T \\ \begin{matrix} \text{---} \end{matrix} & \times & \begin{matrix} \text{---} \end{matrix} \end{matrix}}{\sqrt{d_k}} \right) \mathbf{V}$$
$$= \mathbf{z}$$

# Multi-Head Attention



# Multi-Head Attention



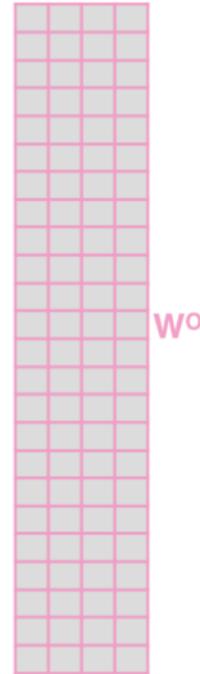
# Соединяем!

1) Concatenate all the attention heads



2) Multiply with a weight matrix  $W^o$  that was trained jointly with the model

$\times$



3) The result would be the  $Z$  matrix that captures information from all the attention heads. We can send this forward to the FFNN

$$= \begin{matrix} Z \\ \hline \end{matrix}$$

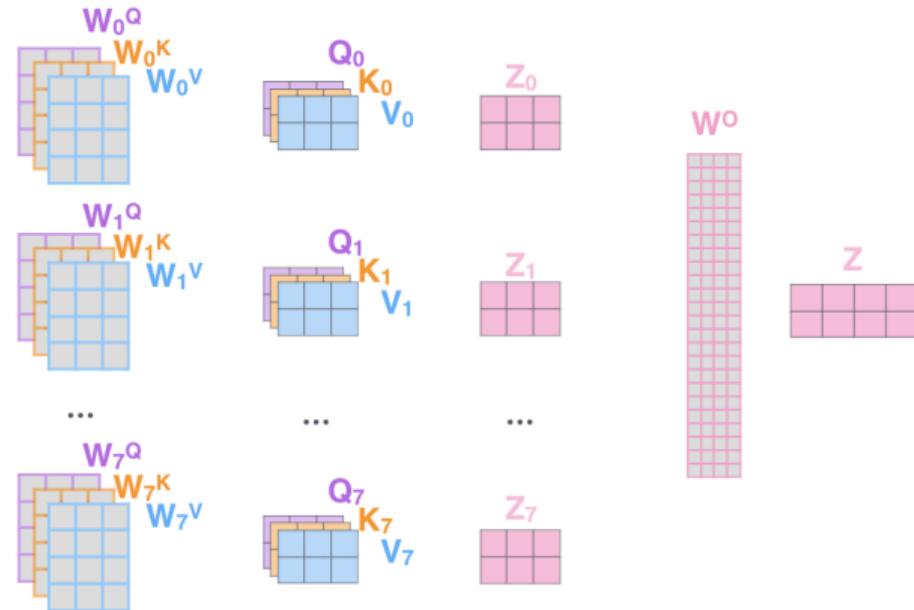
# ИТОГО

1) This is our input sentence\*  
each word\*



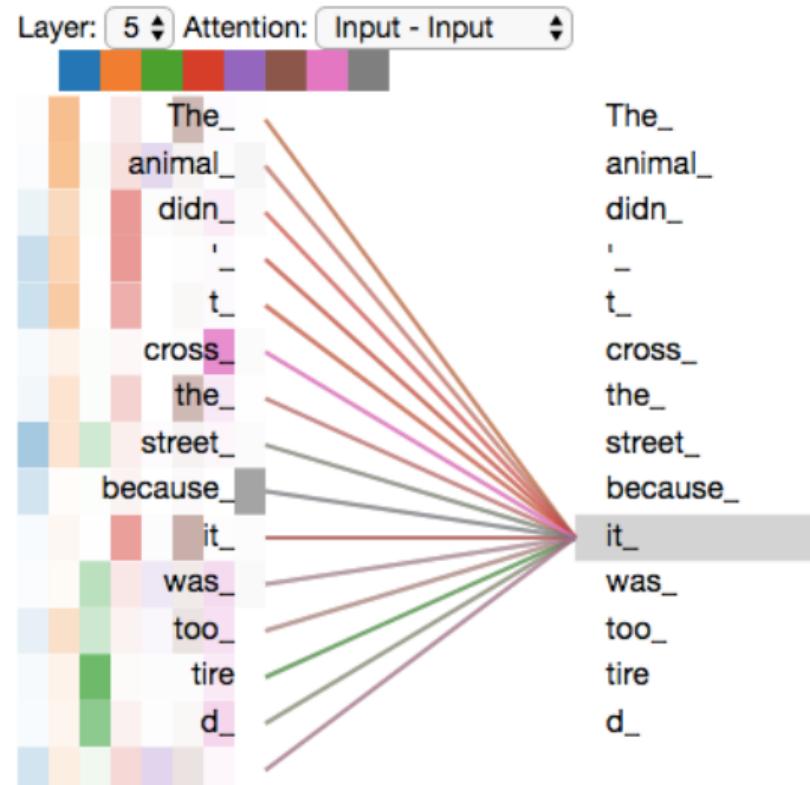
2) We embed each word\*

3) Split into 8 heads.  
We multiply  $X$  or  $R$  with weight matrices

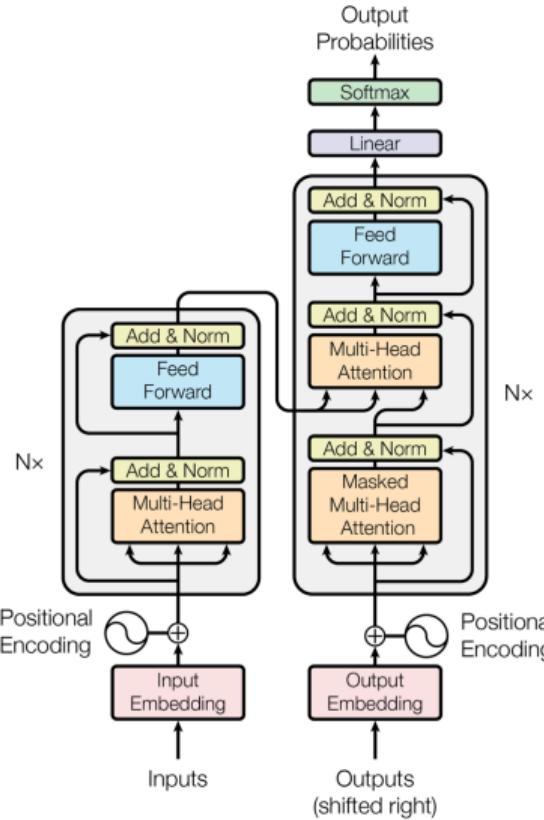


\* In all encoders other than #0, we don't need embedding.  
We start directly with the output of the encoder right below this one

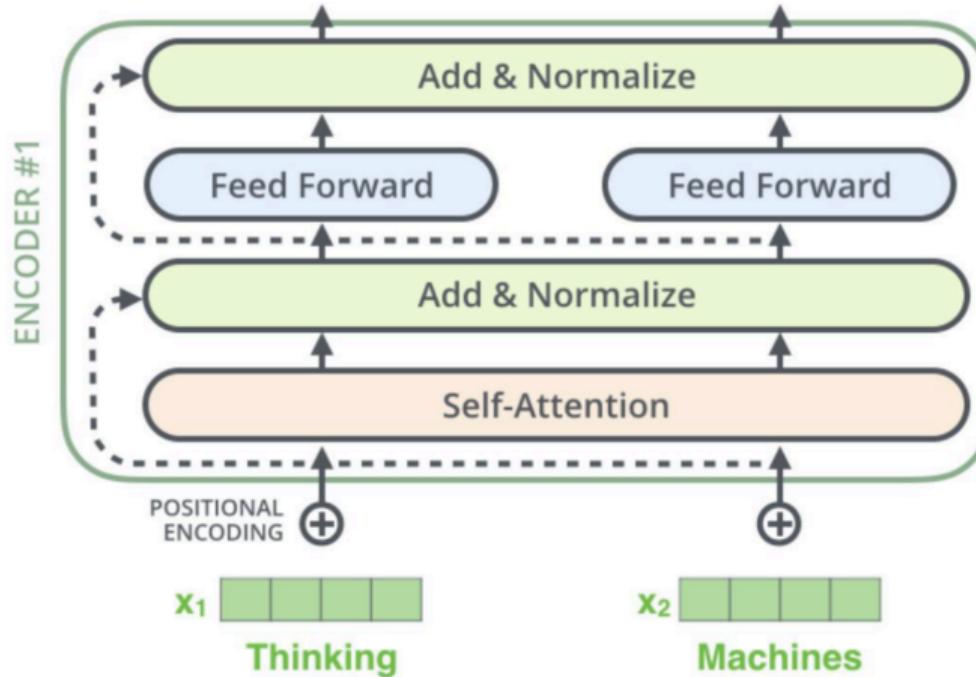
# ИТОГО



# Transformer full architecture



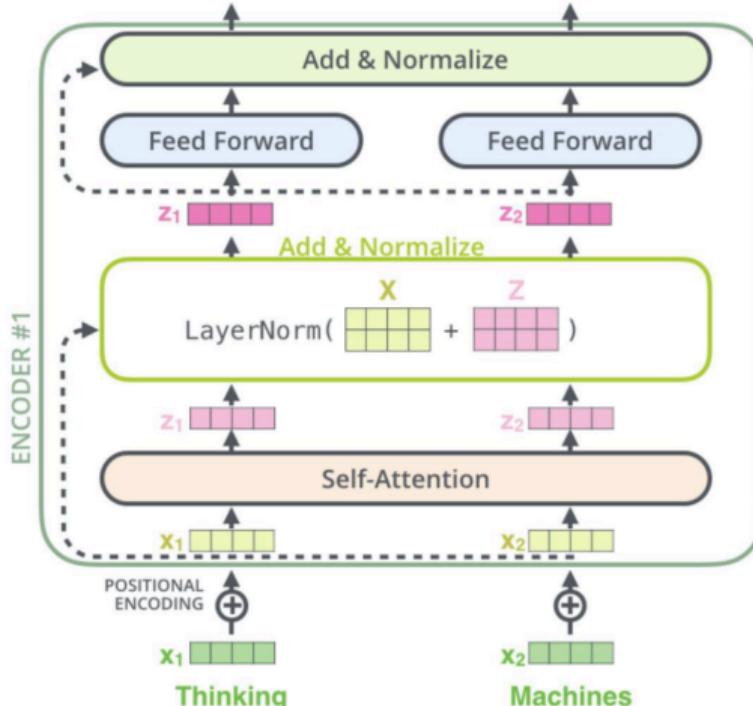
# Layer normalization



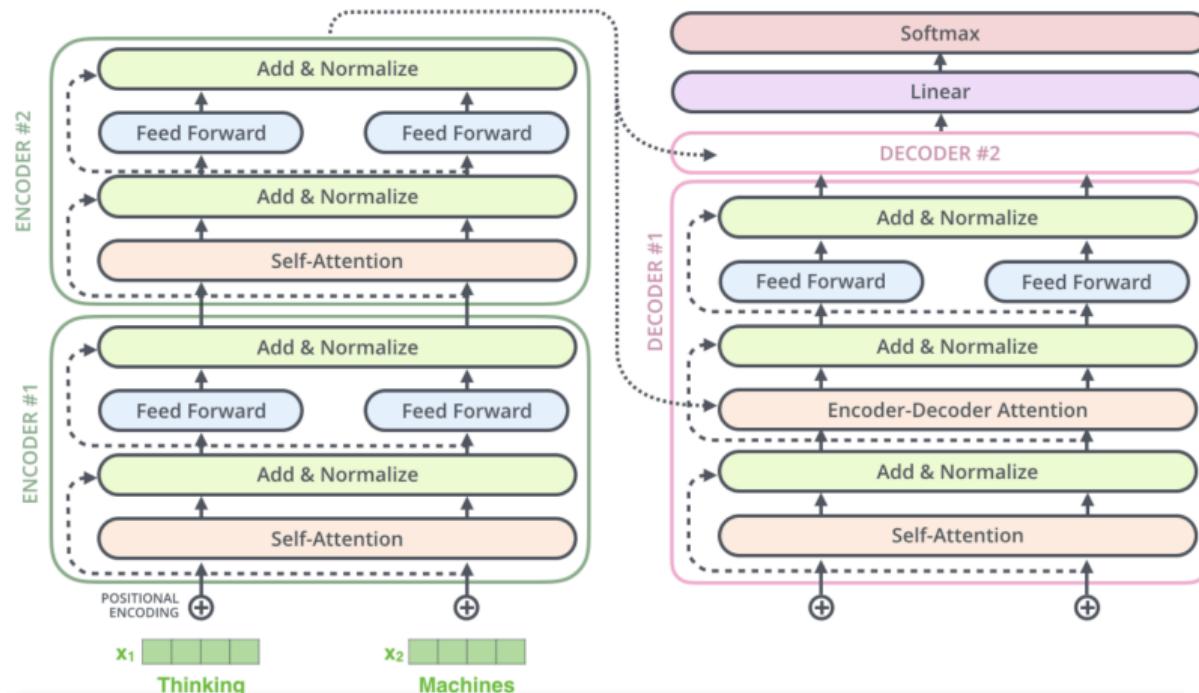
# Layer normalization

Like BatchNorm

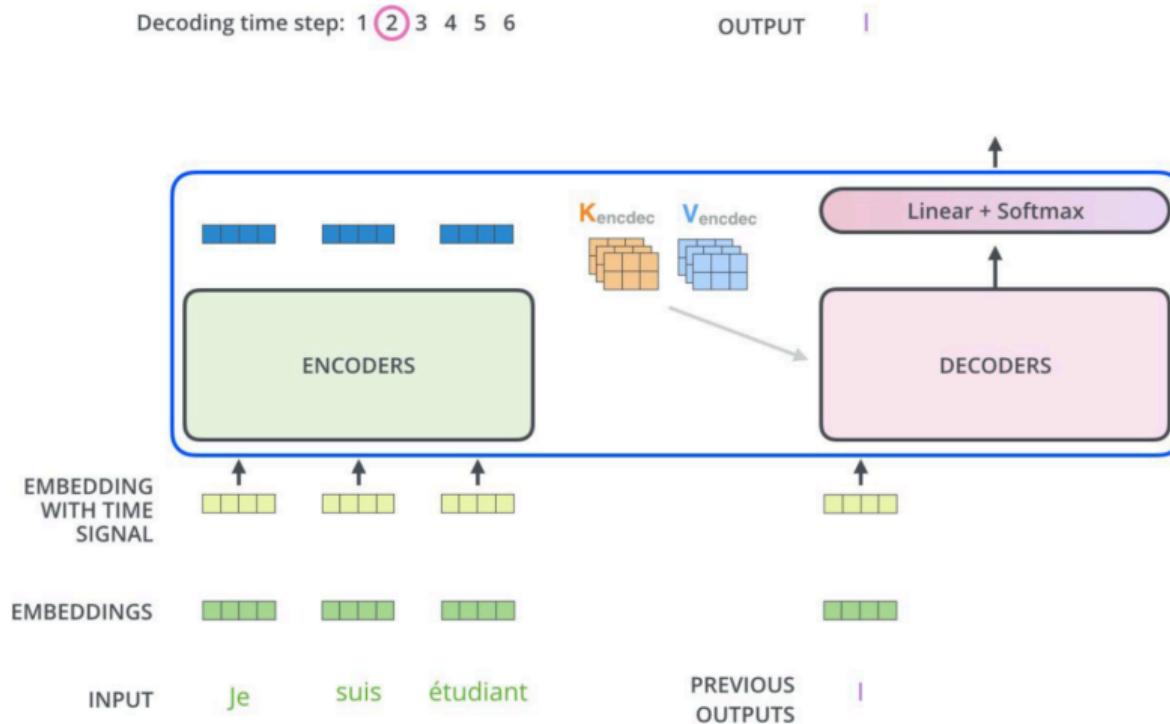
but normalize along  
all features  
representing latent  
vector



# Transformer detailed



# Decoder



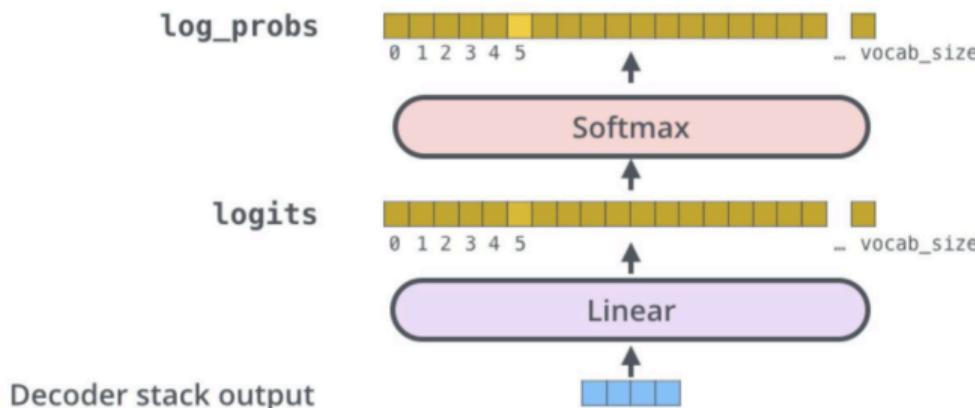
# Final steps

Which word in our vocabulary  
is associated with this index?

am

Get the index of the cell  
with the highest value  
(argmax)

5



# Итого

1. У нас нет никаких слоев, кроме dense
2. Учится очень классно, находит множество взаимосвязей
3. position encoding позволяет учитывать позицию в тексте



**И понеслась!!!** (развитие дальше - инженерные хаки и закидывание железом)

# BERT

Крутой обзорчик с техническими деталями - живет в лекциях МФТИ.  
Многие моменты этой лекции заимствованы из него.

## **Обзор**

# BERT

Шел 2018 год и гугл сказал - наши комьютеры самые мощные, а данные самые большие!

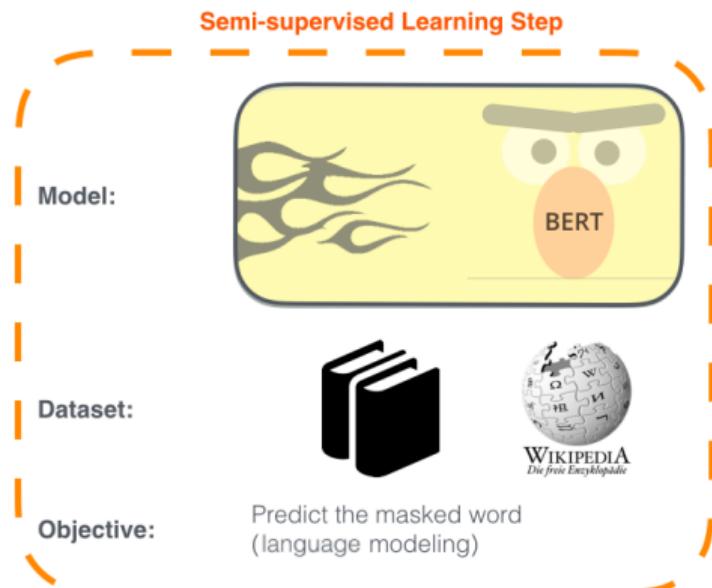
BERT - Bidirectional Encoder Representations from Transformers

Почему круто - придумали как предобучать без учителя (да, вот оно, вот он наш космос). И потом переиспользовать веса!

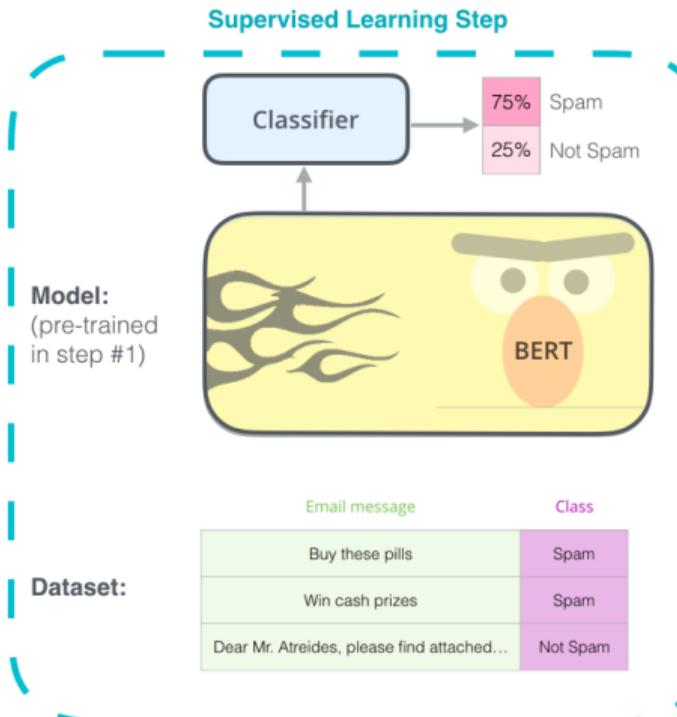
# BERT

## 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.

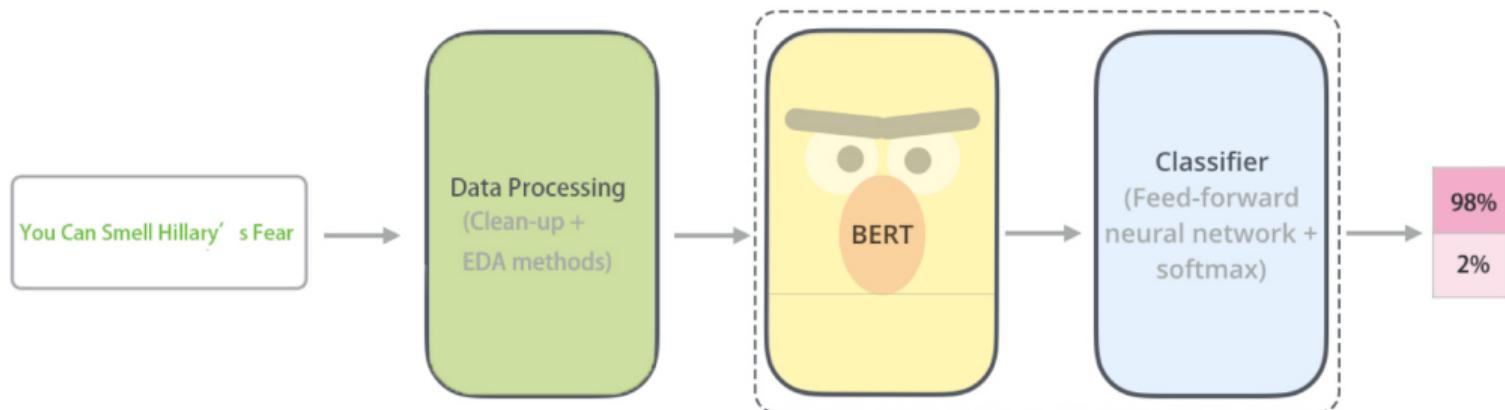


## 2 - Supervised training on a specific task with a labeled dataset.

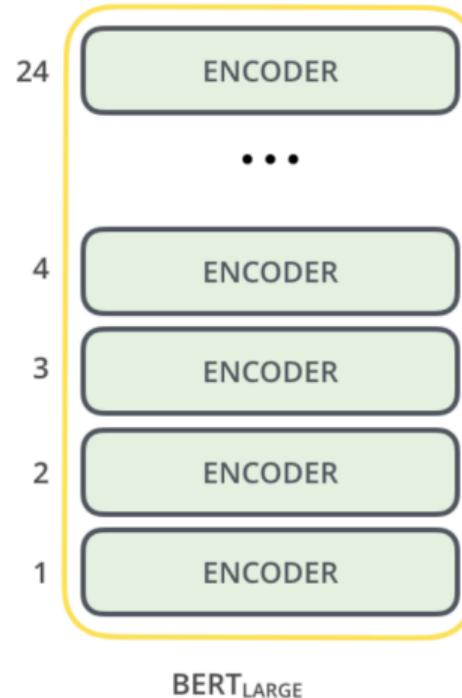
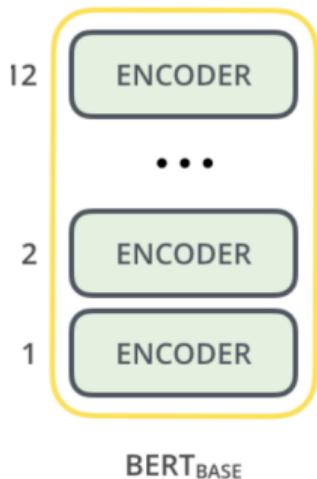


# BERT

Input Features      Output Prediction



# BERT: Base, Large



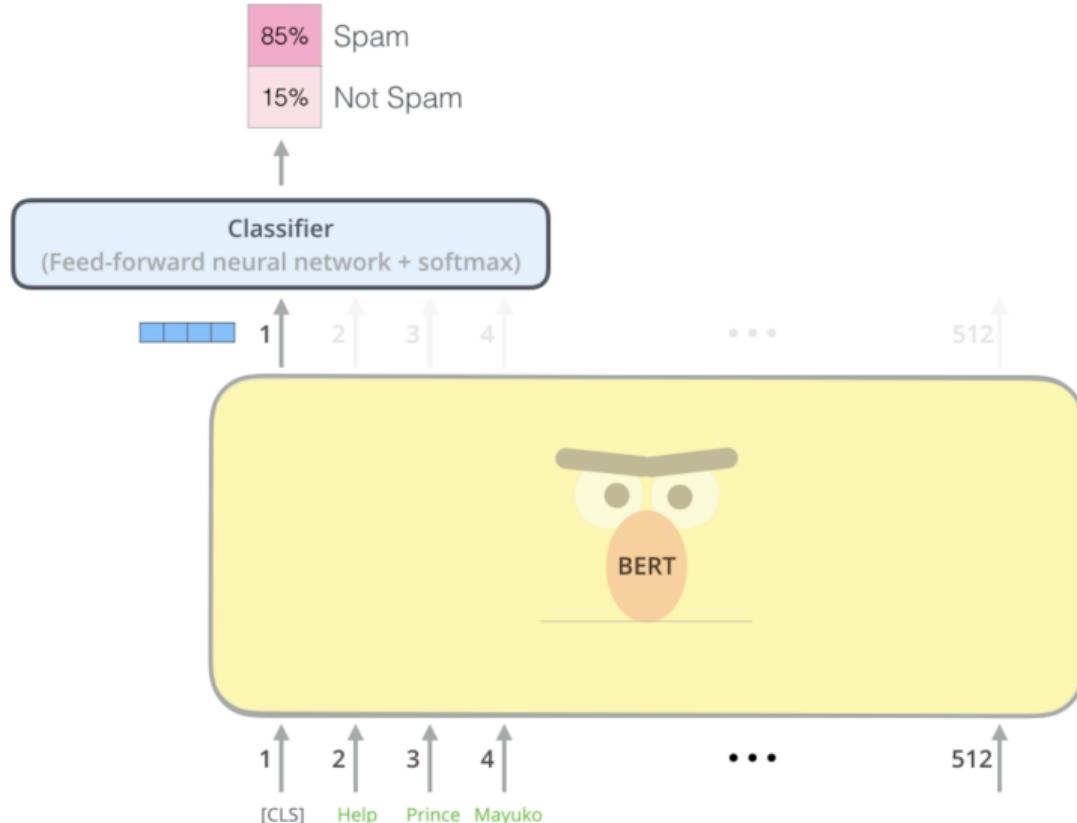
# Лежащие внутри идеи и почему он популярный

1. Предобучаем по двум задачам - берем корпус текстов и маскируем часть предложений, заставляем учить и предсказывать маску.
2. И вторая идея - предсказываем следующее слово в предложении
3. Он из коробки знает язык, ему 1-2 эпохи надо подсказать, что с этим знанием делать
4. В готовых либах лежат много готовых под задачи бортов - классификация, вопросно - ответные системы и тому подобное.
5. Опять же - подаем слова кусочками, чтобы как-то решать проблему оов.

Все мы всех победили - нам не нужно размечать данные для обучения, мы счастливы!

**<http://jalammar.github.io/illustrated-bert/>**

# Pretreaining. Task 1: Masked LM



# Pretaining. Task 2: Next Sentence Prediction

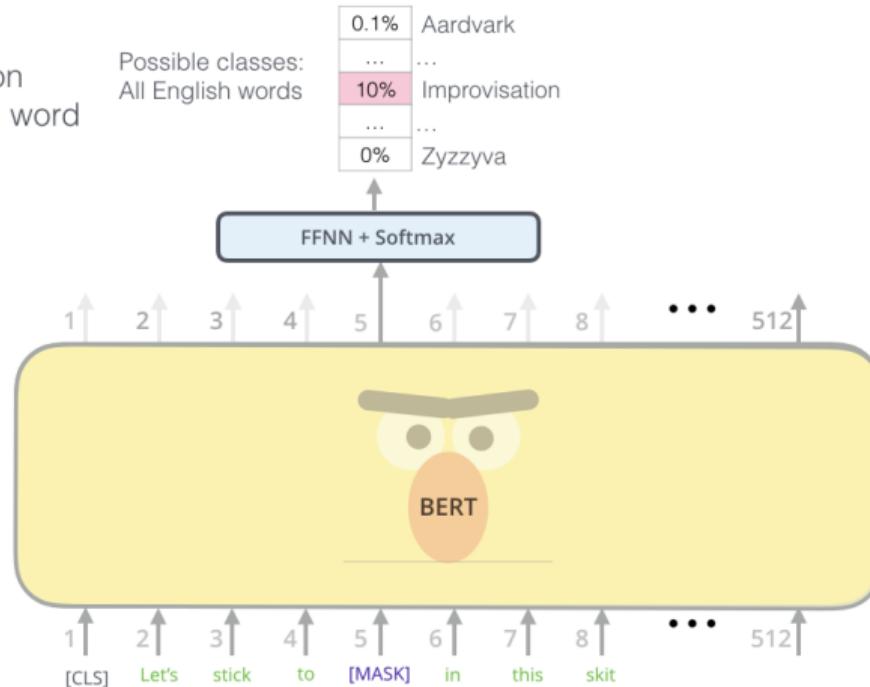
Use the output of the masked word's position to predict the masked word

Possible classes:  
All English words

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

FFNN + Softmax

Randomly mask  
15% of tokens



Input

[CLS] Let's stick to improvisation in this skit

# Где использовать?

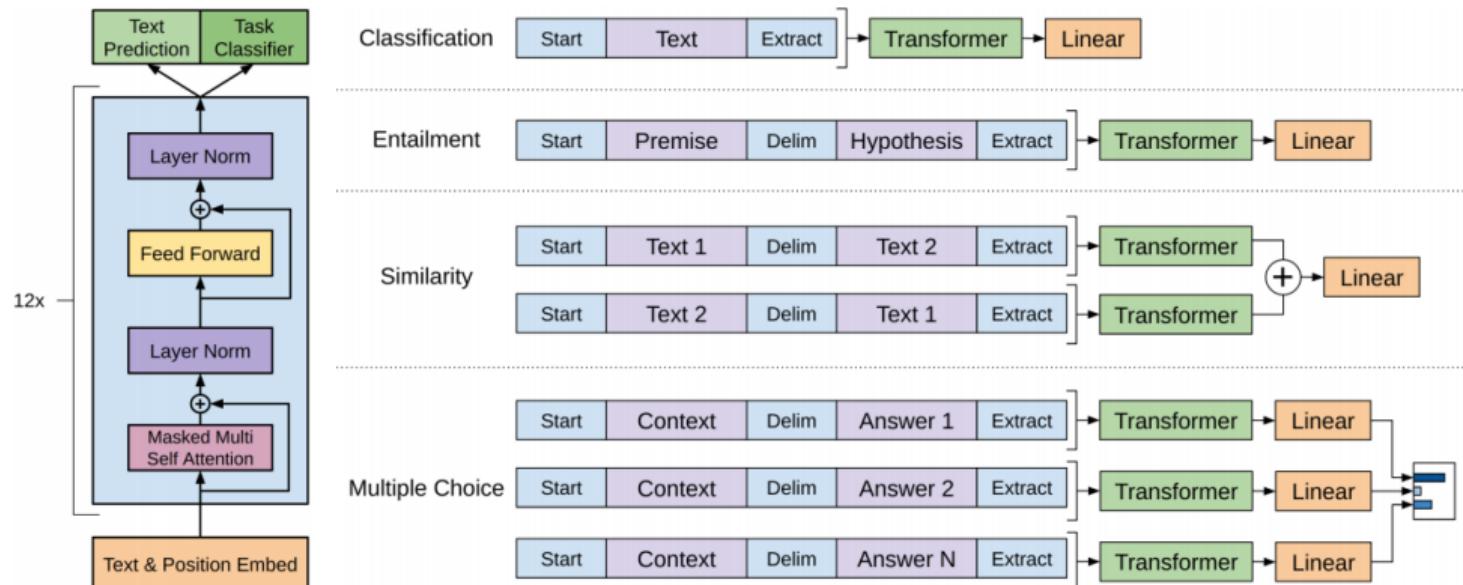
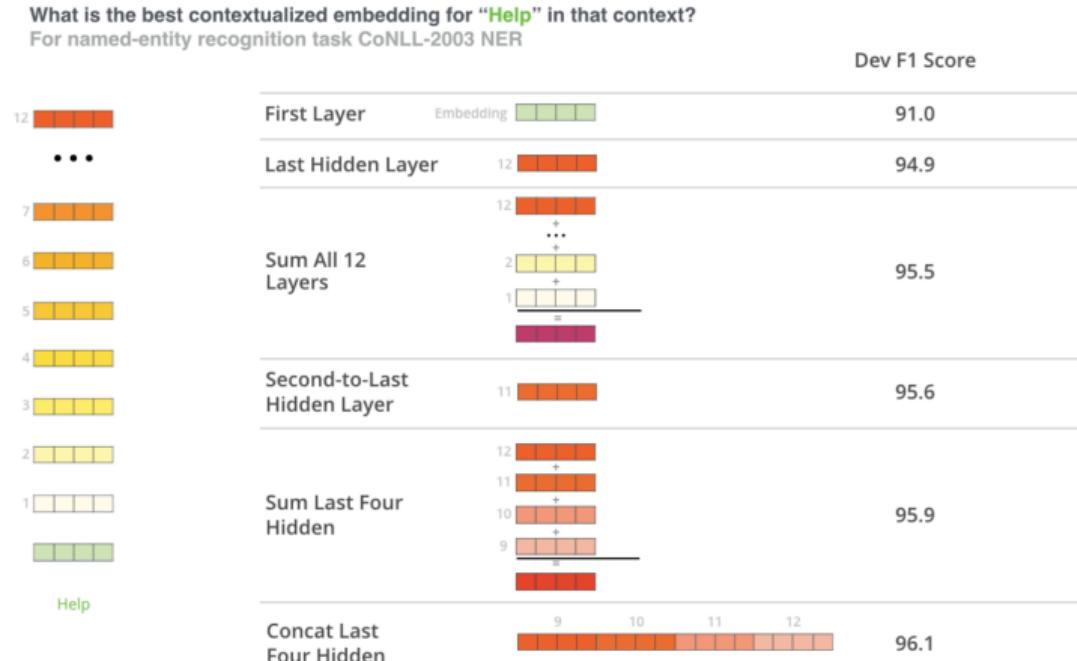


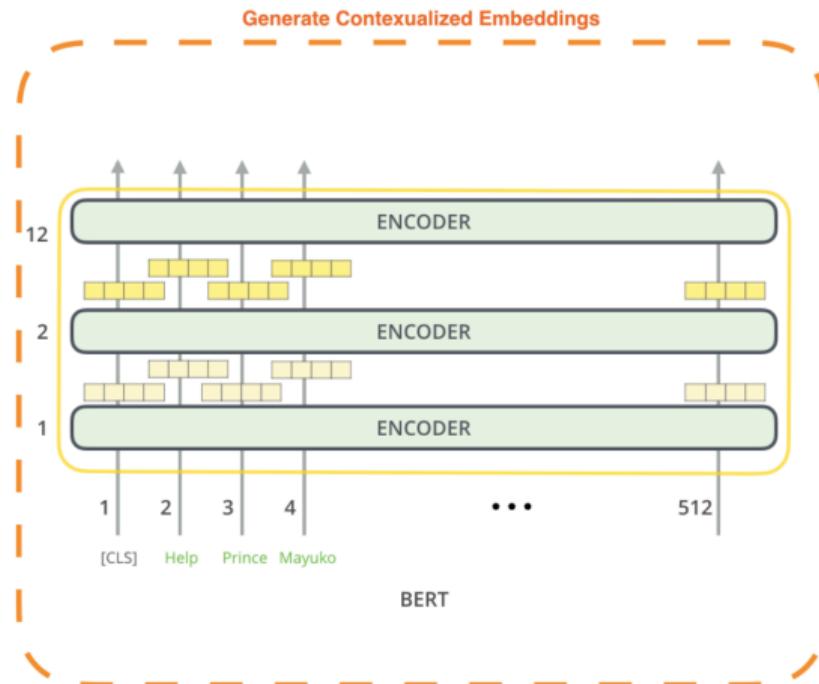
Figure 1: **(left)** Transformer architecture and training objectives used in this work. **(right)** Input transformations for fine-tuning on different tasks. We convert all structured inputs into token sequences to be processed by our pre-trained model, followed by a linear+softmax layer.

# BERT for feature extraction

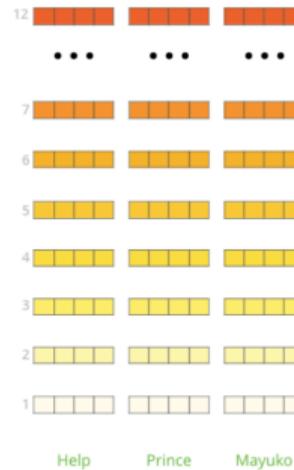
The fine-tuning approach isn't the only way to use BERT. Just like ELMo, you can use the pre-trained BERT to create contextualized word embeddings.



# BERT for feature extraction



The output of each encoder layer along each token's path can be used as a feature representing that token.



But which one should we use?

# Немного статистики

<b>Model Type</b>	<b>Vocab Size</b>	<b>Hidden Dim</b>	<b># Params</b>	<b>Model Size (MB)</b>	<b>FLOPS ratio</b>
BERT <sub>DISTILLED</sub>	4928	48	1,775,910	6.8	1.3%
		96	5,665,926	22	1.32%
		192	19,169,094	73	4.49%
BERT <sub>BASE</sub>	30522	768	110,106,428	420	100%