



# BERT

Табишева Анастасия  
ФКН ПМИ 171  
НИС 2020



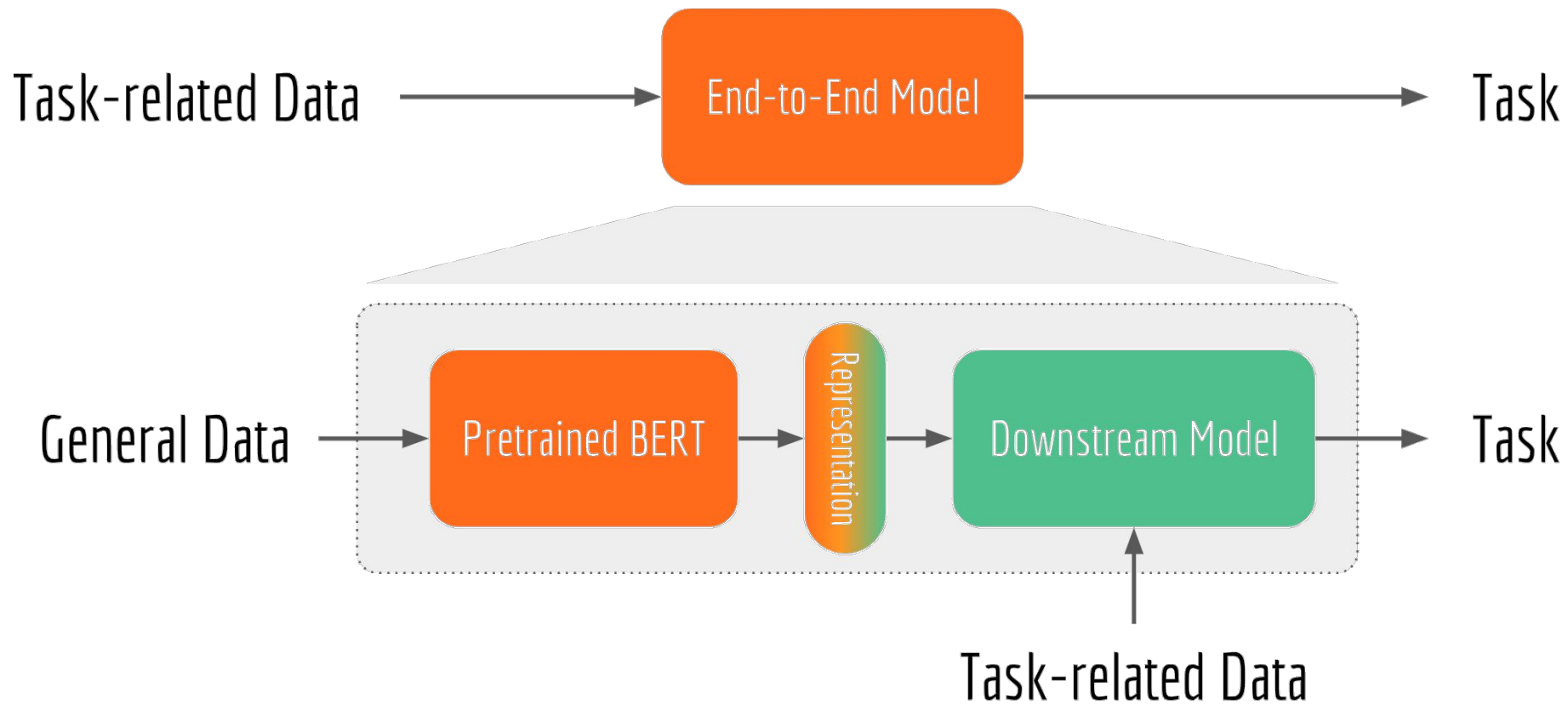
# BERT: зачем?

## Проблема

Обучать “с нуля” всю модель и собирать данные для конкретной задачи  
- долго и трудоёмко

## Идея

Будем использовать предобученную на огромном корпусе текстов модель и дообучать её под узкую задачу



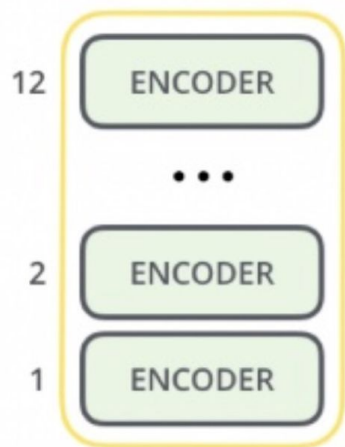
# Особенности BERT:

- Двухнаправленный трансформер
- Дообучение до разных задач
- Универсальный словарный запас
- Лучшие результаты для некоторых бенчмарков

# Bert: Model Architecture



BERT<sub>BASE</sub>



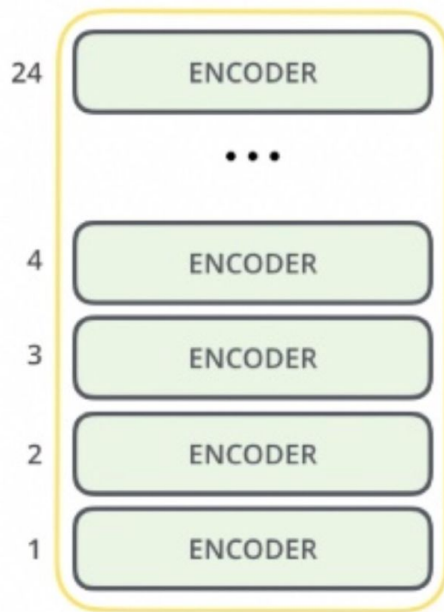
BERT<sub>BASE</sub>

## Equal to Open AI Transformer

- 12 Blocks (Heads)
- 768 Hidden Unit
- 12 Attention Heads
- 110M parameters



BERT<sub>LARGE</sub>



BERT<sub>LARGE</sub>

## State of the Art Model

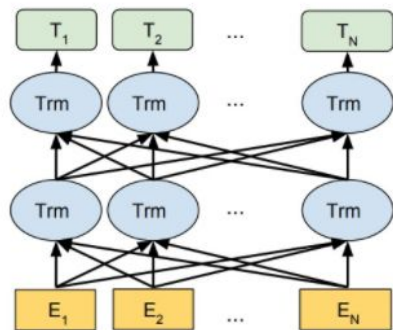
- 24 Blocks
- 1024 Hidden Unit
- 16 Attention Heads
- 340M parameters

# Эволюция эмбеддингов

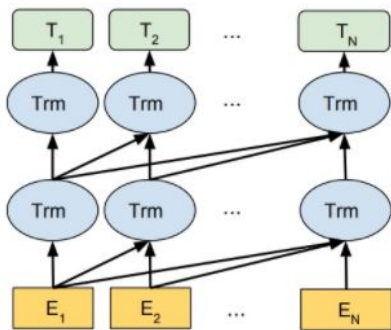
- Эмбеддинги **word2vec** не учитывают контекст
- **OpenAI GPT** - только левый (правый) контекст
- **ELMO** - независимо обучается на левом и правом контексте
- **BERT** - одновременно учитывает весь контекст

- ELMO - две независимые LSTM
- GPT - однонаправленный трансформер
- BERT - двунаправленный трансформер

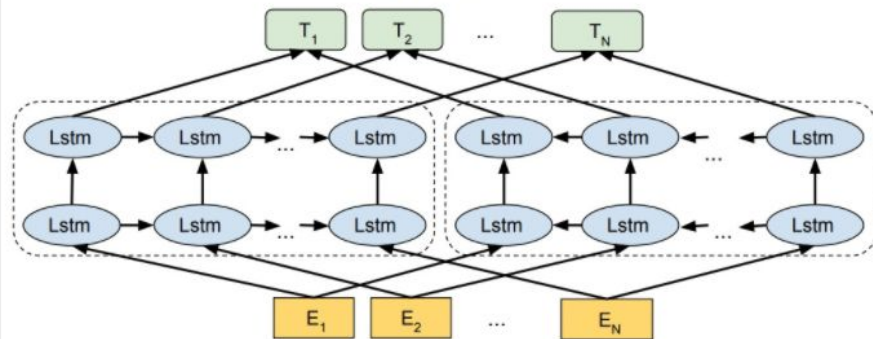
BERT (Ours)



OpenAI GPT



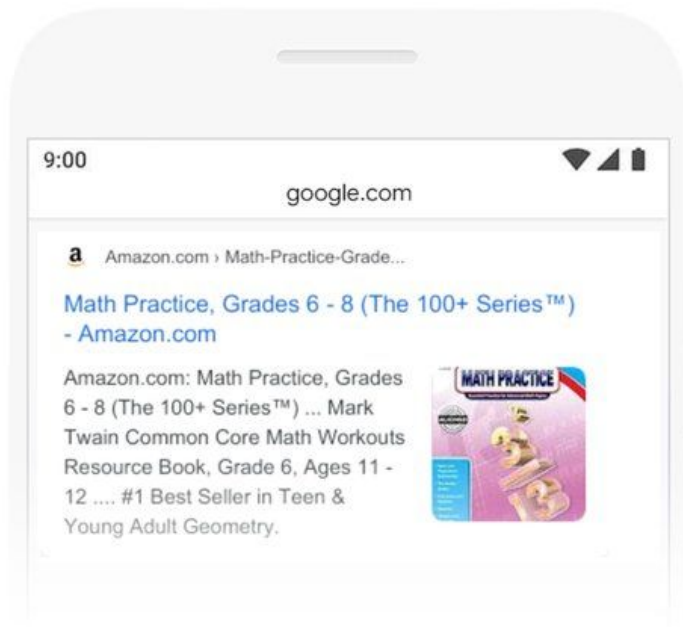
ELMo



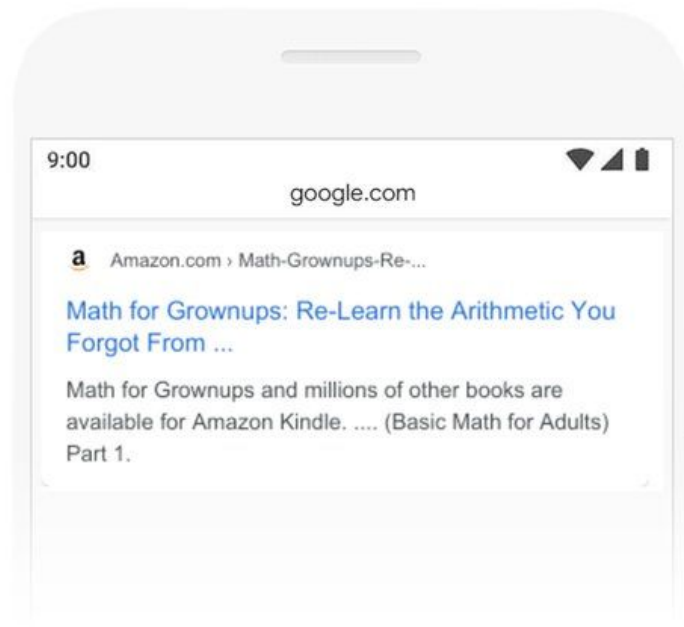


math practice books for adults

BEFORE



AFTER



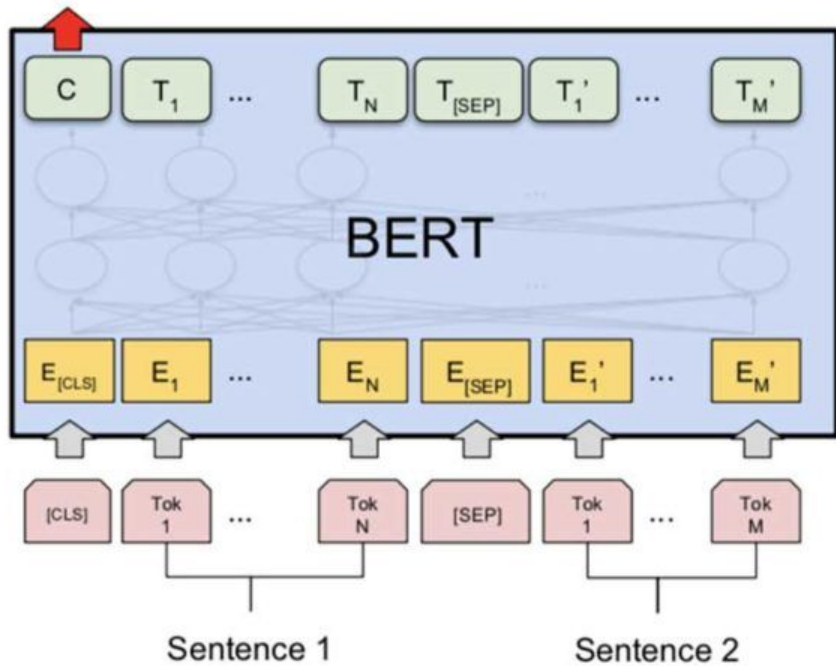


# Представление токенов для BERT

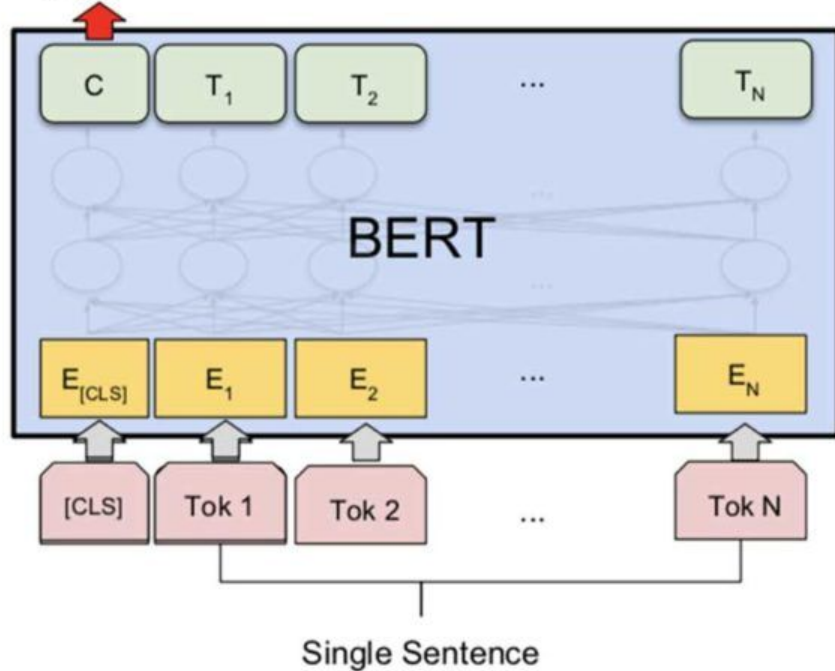
- WordPiece токенизация слова
- Индикатор одного из двух предложений
- Позиция слова в предложении

Input	[CLS]	my	dog	is	cute	[SEP]	he	likes	play	##ing	[SEP]
Token Embeddings	$E_{[\text{CLS}]}$	$E_{\text{my}}$	$E_{\text{dog}}$	$E_{\text{is}}$	$E_{\text{cute}}$	$E_{[\text{SEP}]}$	$E_{\text{he}}$	$E_{\text{likes}}$	$E_{\text{play}}$	$E_{\text{##ing}}$	$E_{[\text{SEP}]}$
	+	+	+	+	+	+	+	+	+	+	+
Segment Embeddings	$E_A$	$E_A$	$E_A$	$E_A$	$E_A$	$E_A$	$E_B$	$E_B$	$E_B$	$E_B$	$E_B$
	+	+	+	+	+	+	+	+	+	+	+
Position Embeddings	$E_0$	$E_1$	$E_2$	$E_3$	$E_4$	$E_5$	$E_6$	$E_7$	$E_8$	$E_9$	$E_{10}$

Class  
Label



Class  
Label



# Стадии обучения

- Pre-training
  - ★ Маскированная языковая модель
  - ★ Предсказание следующего предложения
- Fine-tuning

# Маскированная языковая модель (MLM)

Хотим научить нашу модель понимать контекст вокруг слов

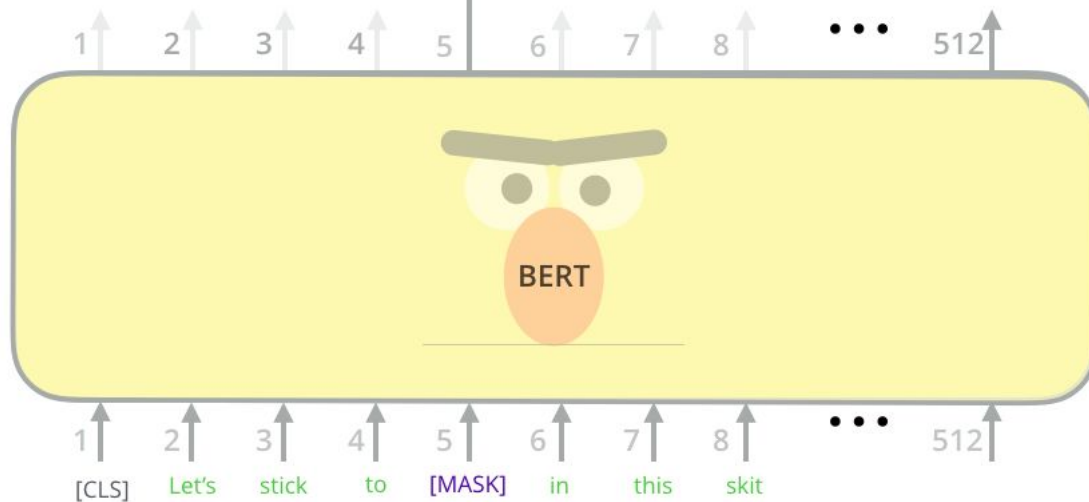
1. Случайно заменяем 15% токенов на [MASK]
  - 80% времени заменяем токен на [MASK]
  - 10% времени заменяем на случайный токен
  - 10% времени токен остается исходным
2. Модель предсказывает только замаскированные токены, но у нее нет информации, какие именно это токены

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%	Zyzzzva

FFNN + Softmax



Randomly mask  
15% of tokens

Input

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

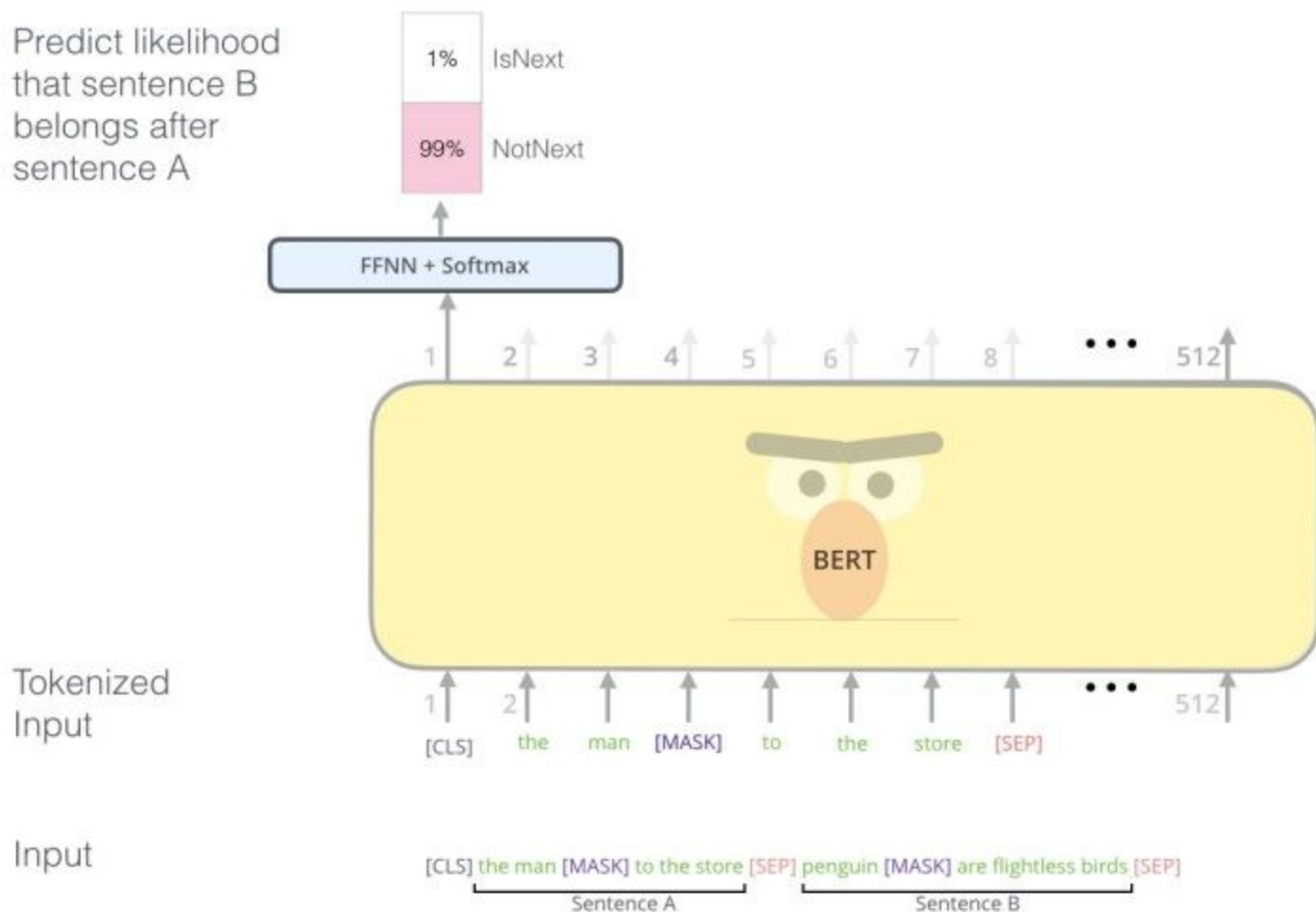
# Предсказание следующего предложения (NSP)

Хотим понимать связь между двумя предложениями

- 50% времени даем следующее предложение
- 50% времени даем несвязанные предложения

Модель должна предсказать, следует ли второе предложение из первого

Predict likelihood  
that sentence B  
belongs after  
sentence A





# Пример ввода для NSP

**Input** = [CLS] the man went to [MASK] store [SEP]

he bought a gallon [MASK] milk [SEP]

**Label** = IsNext

**Input** = [CLS] the man [MASK] to the store [SEP]

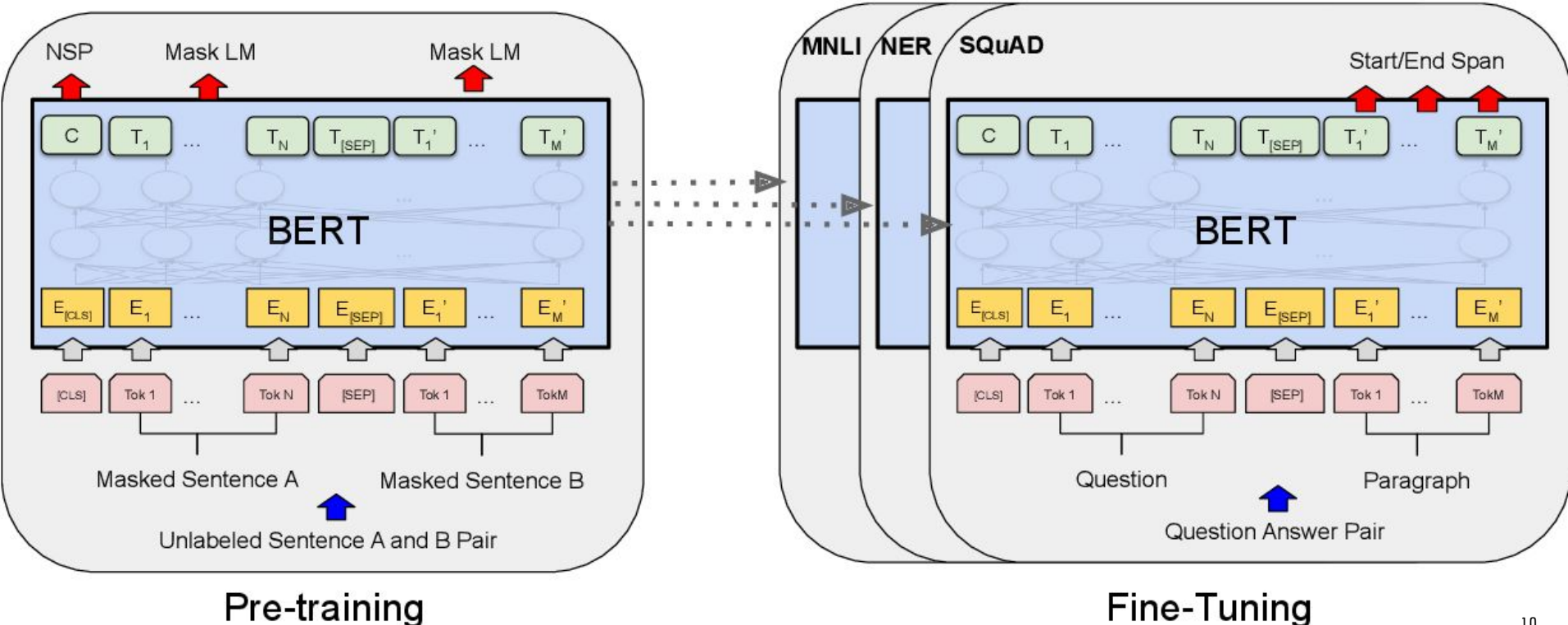
penguin [MASK] are flight ##less birds [SEP]

**Label** = NotNext

# Fine-tuning

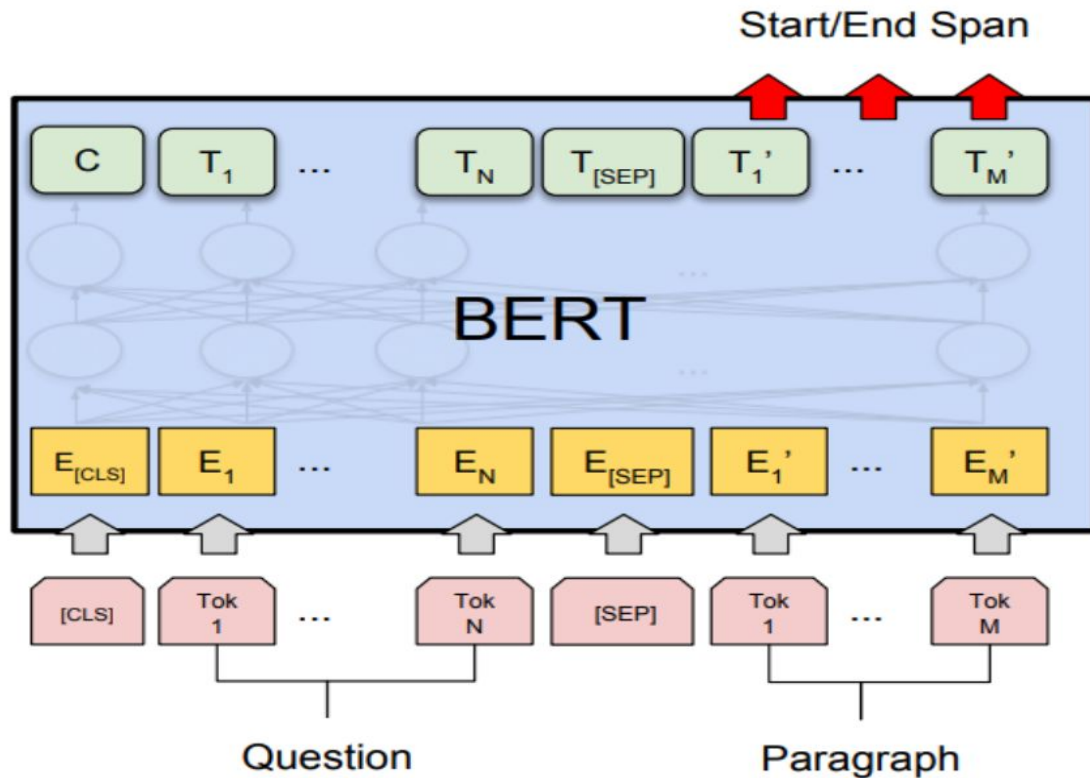
- Классификация предложений (анализ тональности)
- Классификация пар предложений (эквивалентность)
- Создание ответов на вопросы (ответ содержится в тексте)
- Таггинг предложения (предсказание именованных сущностей)

## Пример задачи для Fine-tuning - SQuAD (Stanford Question Answering Dataset)

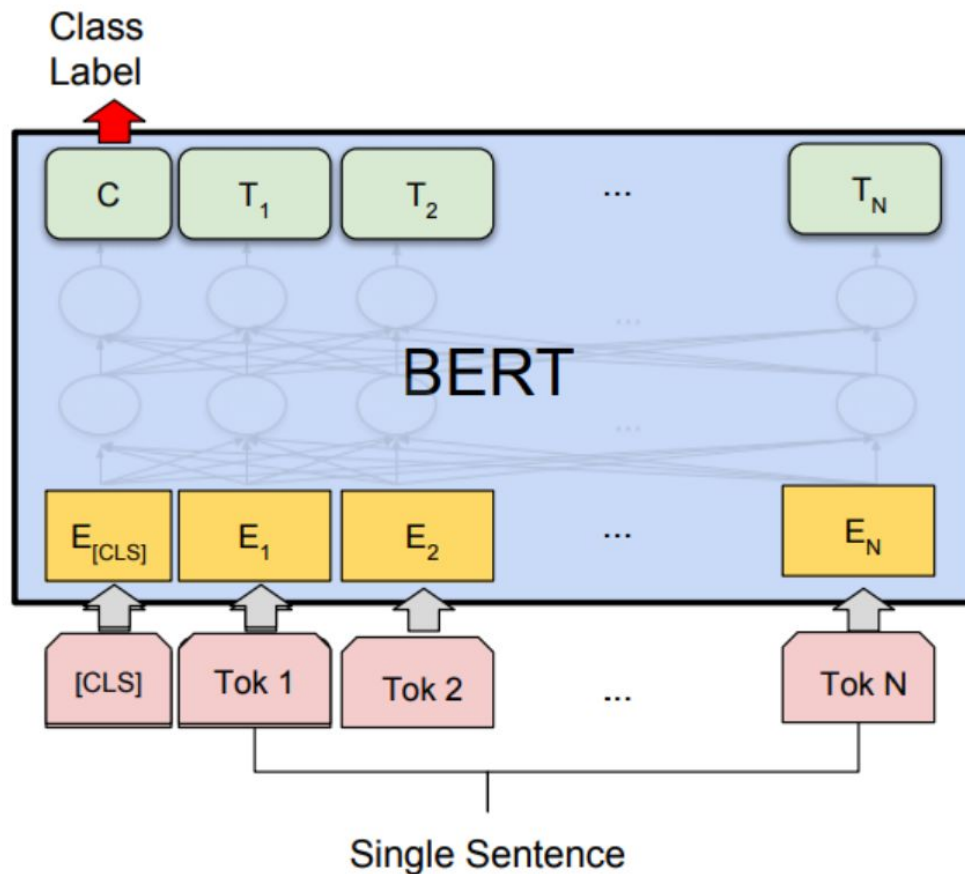


**Первая последовательность - вопрос**

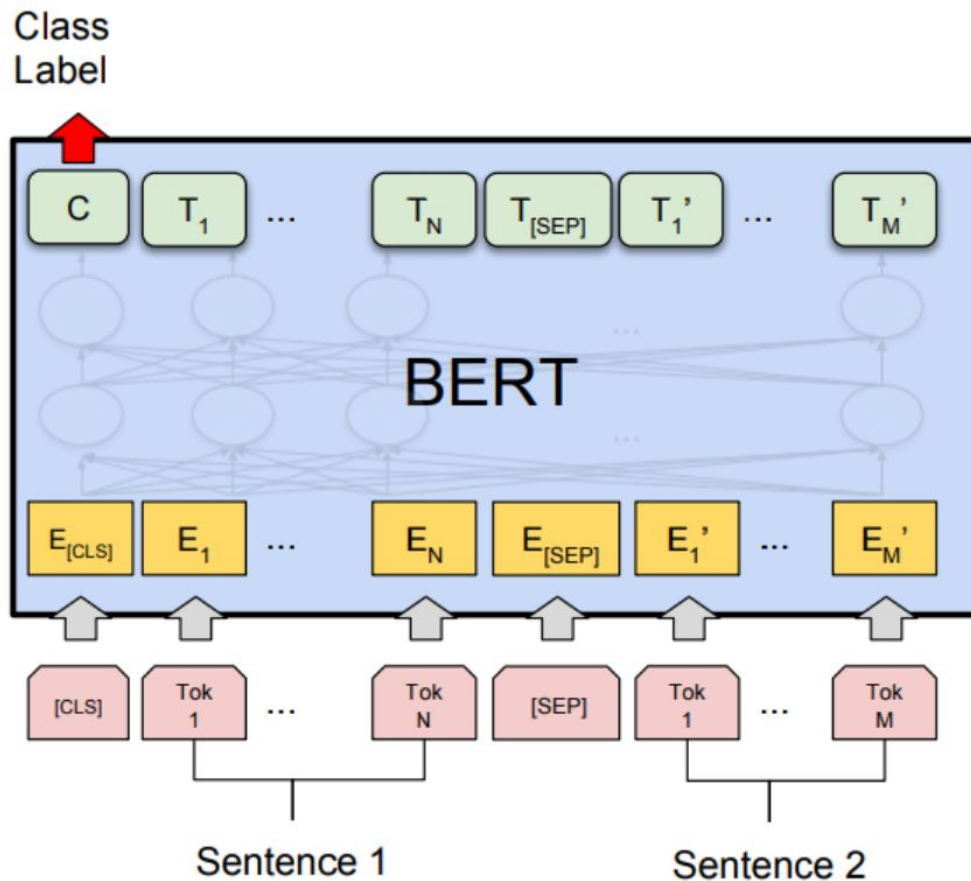
**Вторая последовательность - текст с ответом**



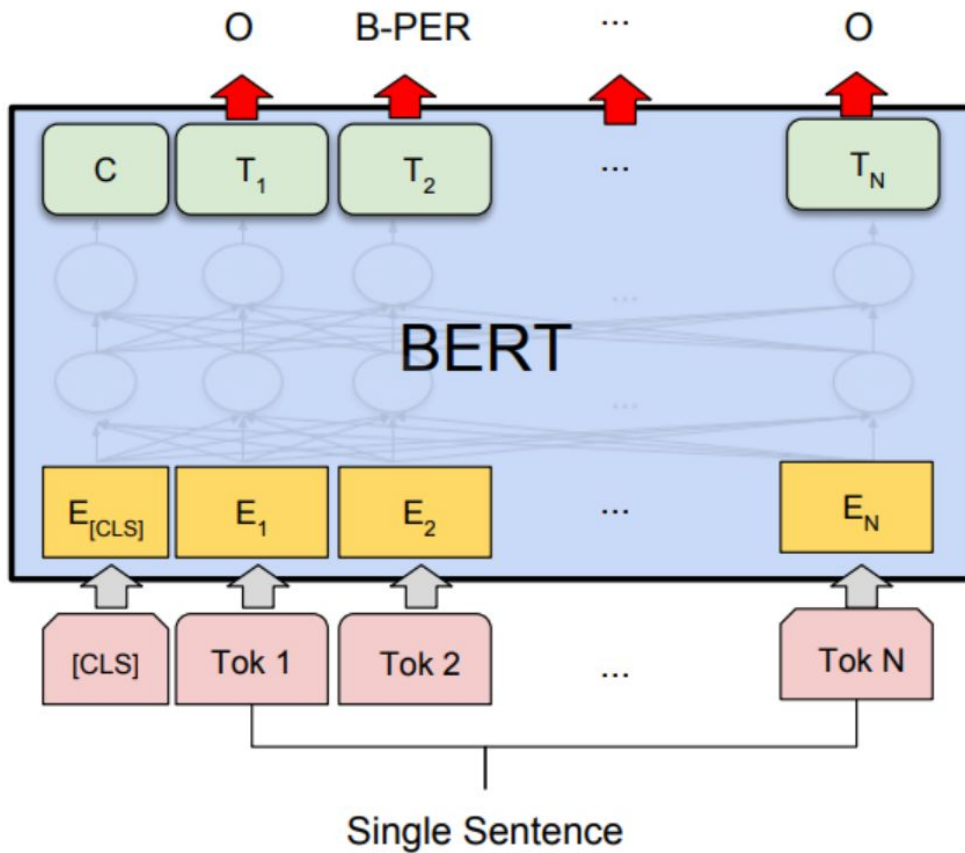
# Классификация предложений



# Классификация пар предложений



# Таггинг предложения



# Недостатки BERT

- Каждое скрытое слово предсказывается в отдельности. Мы теряем информацию о возможных связях между маскированными словами
- Несоответствие между тренировкой модели (есть [MASK] токены) и использованием предобученной модели (таких токенов нет)
- Слишком много параметров



# Результаты

- Показал state-of-the-art результаты в 11 бенчмарках NLP

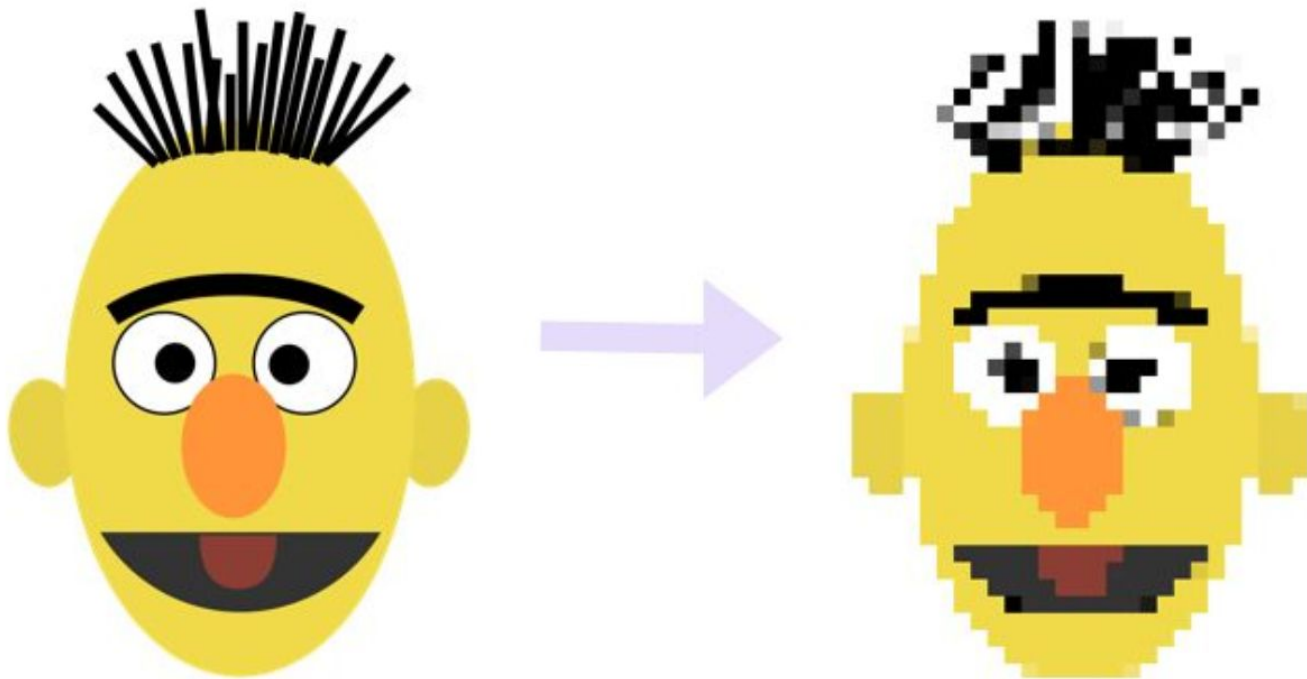
System	MNLI-(m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Average -
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
BERT <sub>BASE</sub>	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERT <sub>LARGE</sub>	<b>86.7/85.9</b>	<b>72.1</b>	<b>92.7</b>	<b>94.9</b>	<b>60.5</b>	<b>86.5</b>	<b>89.3</b>	<b>70.1</b>	<b>82.1</b>

## Не слишком ли много?

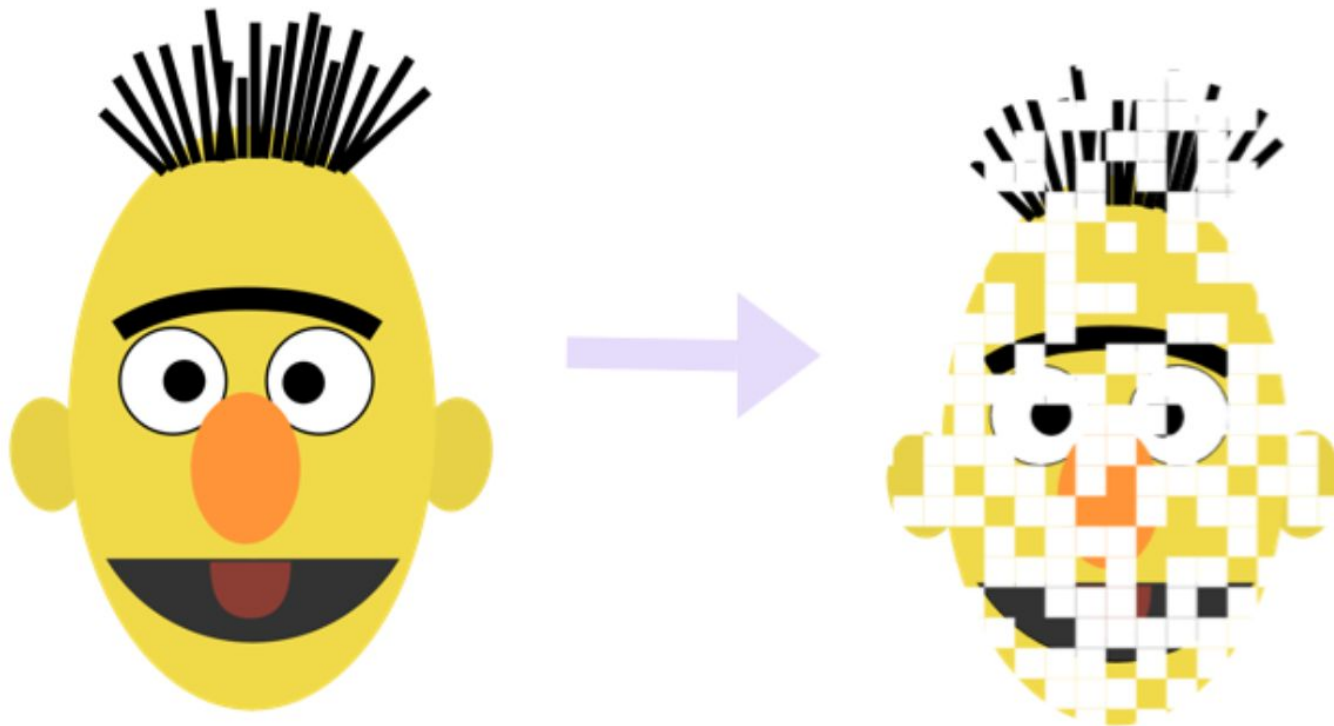
- BERT base - 12 слоев (блоков трансформера), 12 attention heads, 110 млн параметров
- BERT Large - 24 слоя, 16 attention heads, 340 млн параметров

## Можно ли меньше?

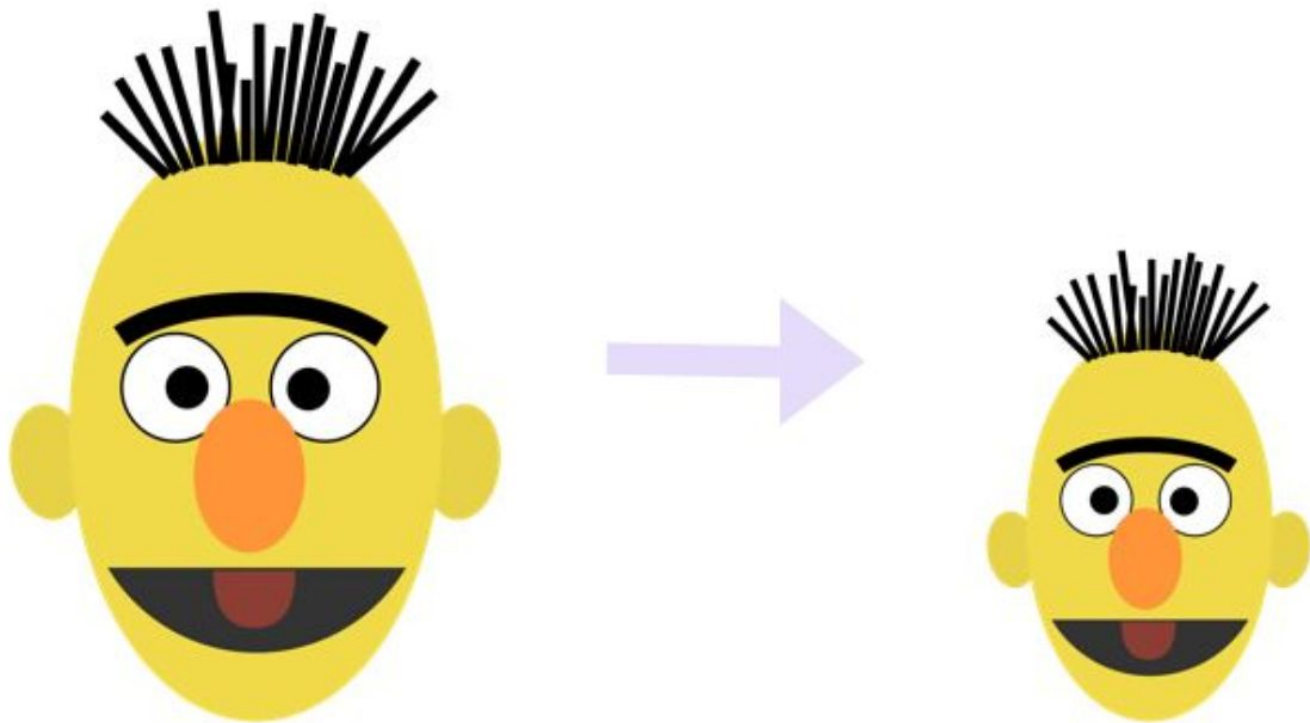
# Квантизация



# Прунинг



# Дистилляция знаний



	Compression	Performance	Speedup	Model	Evaluation	
Distillation	DistilBERT (Sanh et al., 2019)	$\times 2.5$	90%	$\times 1.6$	BERT <sub>6</sub>	All GLUE tasks
	BERT <sub>6</sub> -PKD (Sun et al., 2019a)	$\times 1.6$	97%	$\times 1.9$	BERT <sub>6</sub>	No WNLI, CoLA and STS-B
	BERT <sub>3</sub> -PKD (Sun et al., 2019a)	$\times 2.4$	92%	$\times 3.7$	BERT <sub>3</sub>	No WNLI, CoLA and STS-B
	(Aguilar et al., 2019)	$\times 2$	94%	-	BERT <sub>6</sub>	CoLA, MRPC, QQP, RTE
	BERT-48 (Zhao et al., 2019)	$\times 62$	87%	$\times 77$	BERT <sub>12</sub> <sup>*†</sup>	MNLI, MRPC, SST-2
	BERT-192 (Zhao et al., 2019)	$\times 5.7$	94%	$\times 22$	BERT <sub>12</sub> <sup>*†</sup>	MNLI, MRPC, SST-2
	TinyBERT (Jiao et al., 2019)	$\times 7.5$	96%	$\times 9.4$	BERT <sub>4</sub> <sup>*†</sup>	All GLUE tasks
	MobileBERT (Sun et al.)	$\times 4.3$	100%	$\times 4$	BERT <sub>24</sub> <sup>†</sup>	No WNLI
	PD (Turc et al., 2019)	$\times 1.6$	98%	$\times 2.5^3$	BERT <sub>6</sub> <sup>†</sup>	No WNLI, CoLA and STS-B
	MiniBERT(Tsai et al., 2019)	$\times 6^{\S}$	98%	$\times 27^{\S}$	mBERT <sub>3</sub> <sup>†</sup>	CoNLL-2018 POS and morphology
	BiLSTM soft (Tang et al., 2019)	$\times 110$	91%	$\times 434^{\ddagger}$	BiLSTM <sub>1</sub>	MNLI, QQP, SST-2
Other Quant.	Q-BERT (Shen et al., 2019)	$\times 13$	99%	-	BERT <sub>12</sub>	MNLI, SST-2
	Q8BERT (Zafrir et al., 2019)	$\times 4$	99%	-	BERT <sub>12</sub>	All GLUE tasks
	ALBERT-base (Lan et al., 2019)	$\times 9$	97%	$\times 5.6$	BERT <sub>12</sub> <sup>**</sup>	MNLI, SST-2
	ALBERT-xxlarge (Lan et al., 2019)	$\times 0.47$	107%	$\times 0.3$	BERT <sub>12</sub> <sup>**</sup>	MNLI, SST-2
	BERT-of-Theseus (Xu et al., 2020)	$\times 1.6$	98%	-	BERT <sub>6</sub>	No WNLI



**NO, you cannot understand the meaning  
of a text without explicitly evaluating its  
linguistic constituents and defining  
grammar rules!**



**haha gpus go brrrrrrrrr**

# Вопросы

- Из каких частей состоят эмбединги для модели?
- Опишите задачи, на которых предобучается BERT
- Для каких задач может быть дообучен BERT, приведите конкретный пример, опишите, что подается на вход и какие нужно добавить выходные слои



# Всем спасибо!

