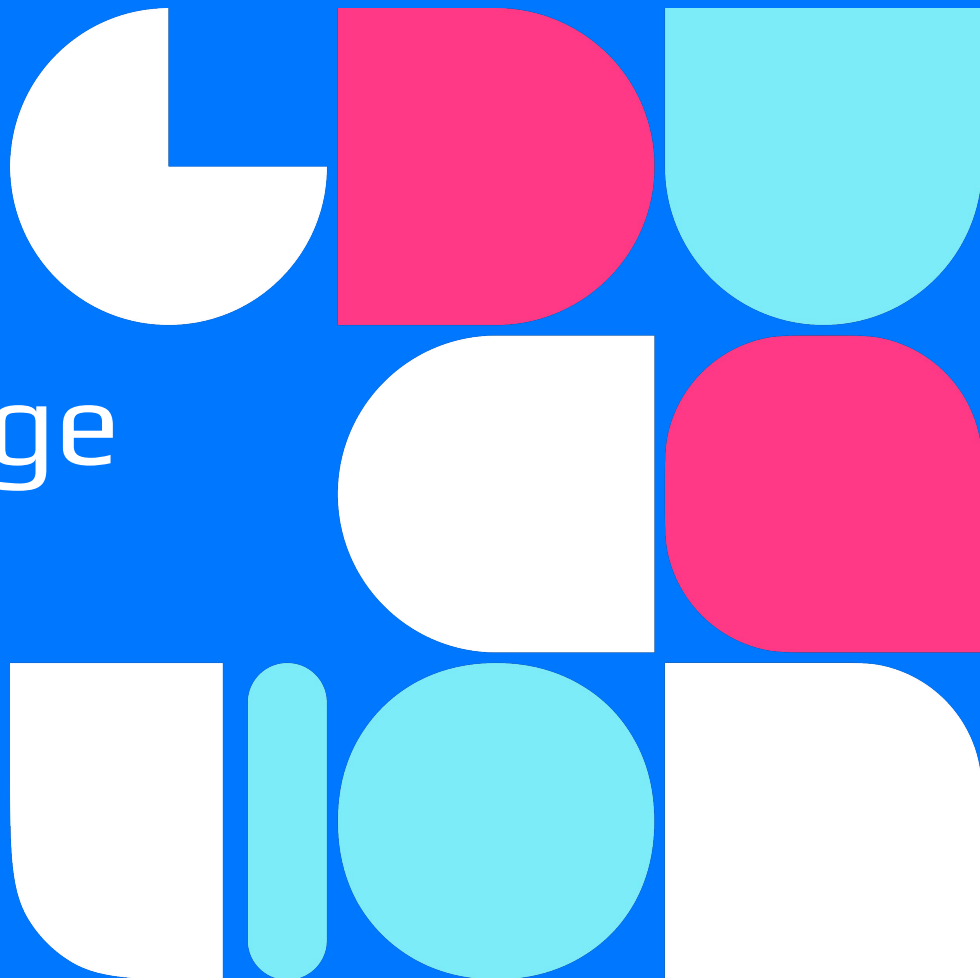




Natural Language Processing

Нейронные сети в ML

Козулин Илья



Цель занятия:

- Понять как может нейросеть оперировать текстовыми данными
- Разобрать как исторически развивались методы получения векторных представлений слов
- Увидеть что из себя представляют и как обучаются трансформеры
- Посмотреть какие задачи могут решаться при помощи текстовых данных и в чем их особенности



Регламент на занятия 29.10 и 05.11

- Ставим в чат
 - “++” если материал предельно понятен/очевиден/уже известен
 - “+” если материал более менее понятен, либо нужно еще немного над ним подумать, но можно идти дальше
 - “-” если что-то совсем не понятно и нужно остановиться и разобрать подробнее
- На лекции скидывается опрос с вопросами по материалу. Оценка за опрос идет к оценке за домашнее задание
- На семинаре будут практические задания, решаемые студентами с демонстрацией экрана. За каждое задание +1 балл к оценке за дз. Максимум +2 балла на одного человека.

О чём поговорим?

- Какие задачи охватывает NLP
- Векторные представления слов/предложений
- Моделирование языка
- Масштабирование моделей
- Parameter-Efficient Fine-Tuning

О чём поговорим?

- Какие задачи охватывает NLP
- Векторные представления слов/предложений
- Моделирование языка
- Масштабирование моделей
- Parameter-Efficient Fine-Tuning

Часть 1

Часть 2

О чём поговорим?

- Какие задачи охватывает NLP
- Векторные представления слов/предложений
- Моделирование языка
- Масштабирование моделей
- Parameter-Efficient Fine-Tuning

Часть 1

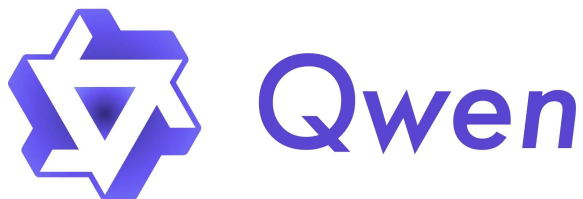
→ Как все начиналось

→ Где мы сейчас

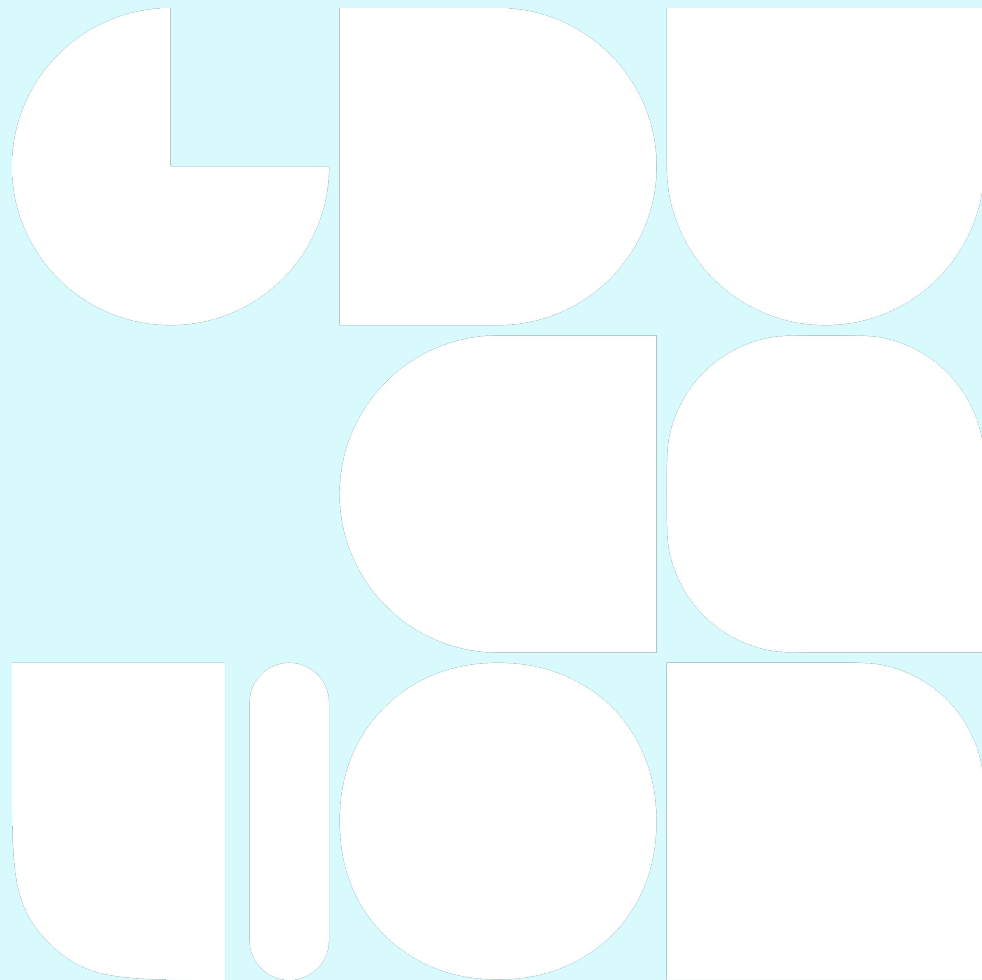
Часть 2

→ Куда мы движемся

Введение



Какие задачи
охватывает NLP



Классификация

1. Текстов/предложений

Примеры:

- I guess your parents dropped you a lot during childhood → "Hate Speech"
- It seems like the weather is going to be very unpleasant today → "Safe"

Классификация

1. Текстов/предложений

Примеры:

- I guess your parents dropped you a lot during childhood → "Hate Speech"
- It seems like the weather is going to be very unpleasant today → "Safe"

2. Токенов

Примеры:

- Named Entity Recognition (NER)

Jürgen Schmidhuber published adversarial neural networks

Jürgen | PERSON (First Name) Schmidhuber | PERSON (Last Name) published | Common Word
adversarial | Common Word neural | Common Word networks | Common Word

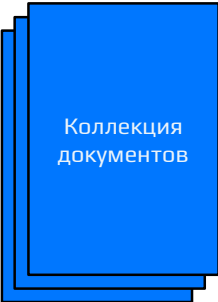
- Part-of-Speech (POS) Tagging

Jürgen Schmidhuber published adversarial neural networks

Jürgen | Proper Noun Schmidhuber | Proper Noun published | Verb adversarial | Adjective neural |
Adjective networks | Common Noun

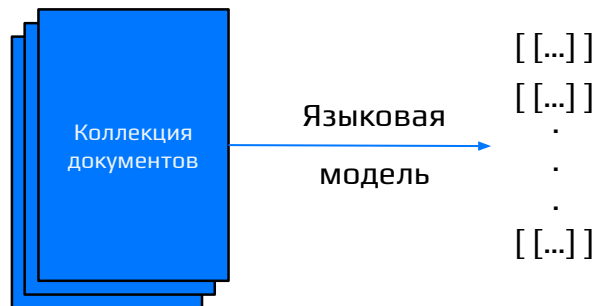
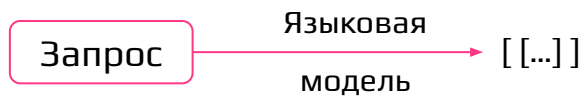
Семантический поиск (Semantic Search)

Запрос

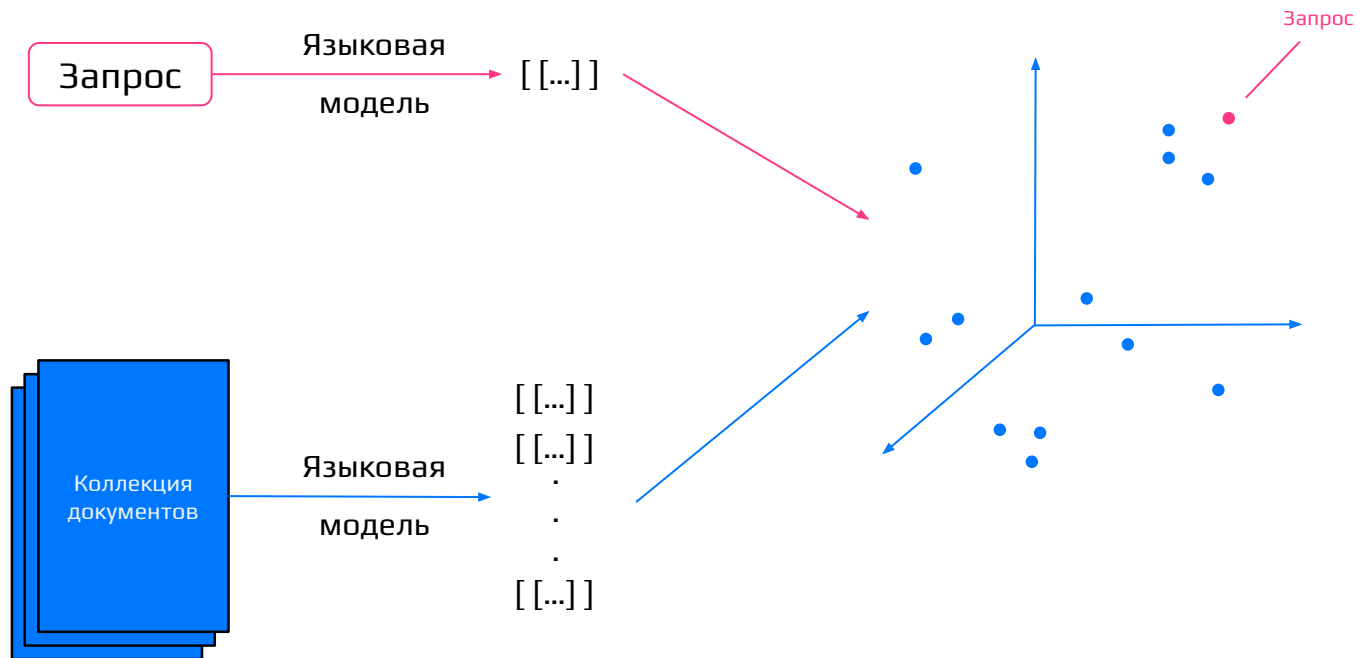


Коллекция
документов

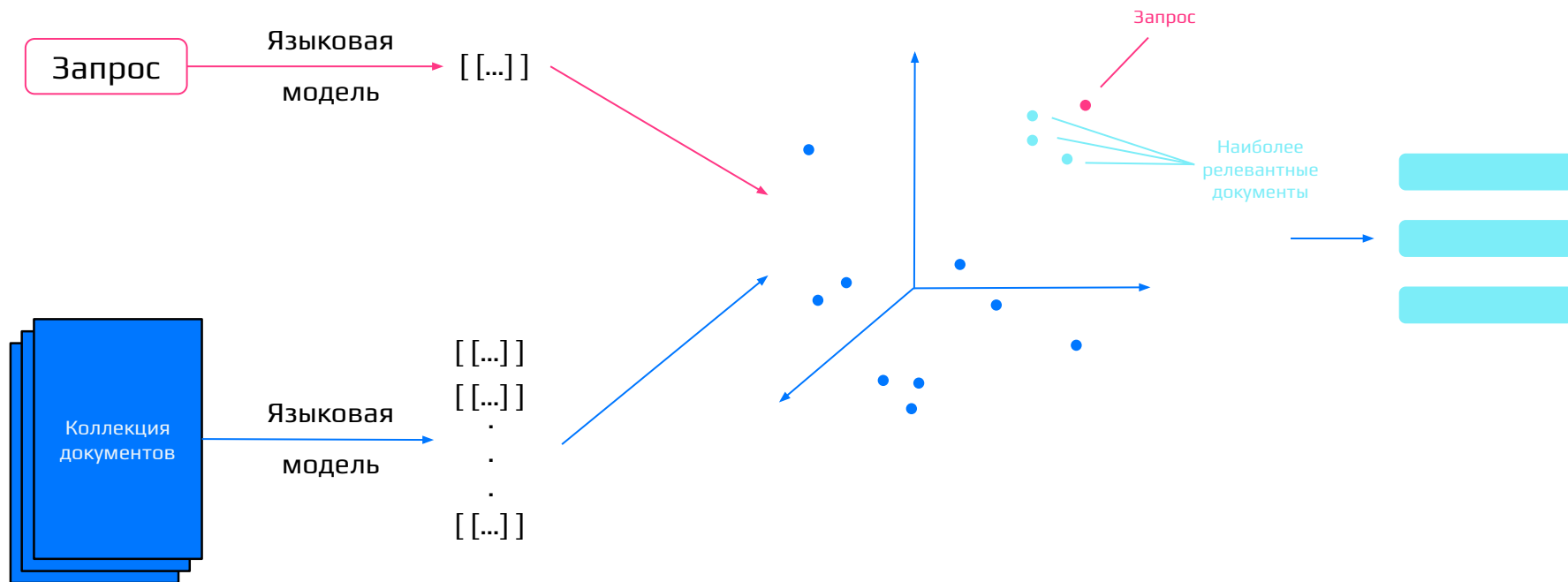
Семантический поиск (Semantic Search)



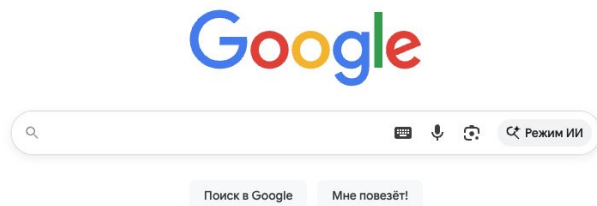
Семантический поиск (Semantic Search)



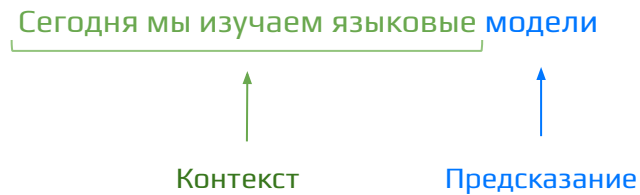
Семантический поиск (Semantic Search)



Семантический поиск (Semantic Search)



Моделирование языка (Language Modeling)



Моделирование языка (Language Modeling)

Я к вам **MASK** — чего же боле? Что я могу еще сказать



Masked Language Modeling (MLM)



пишу

Модель видит все токены, включая те, что после маски

Я к вам пишу — чего же боле? Что я могу еще **MASK**



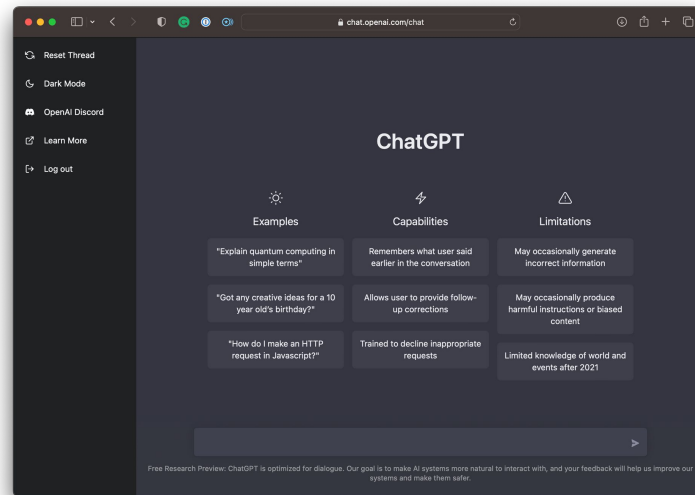
Causal Language Modeling (CLM)



сказать

Модель видит только те токены, что находятся до маски

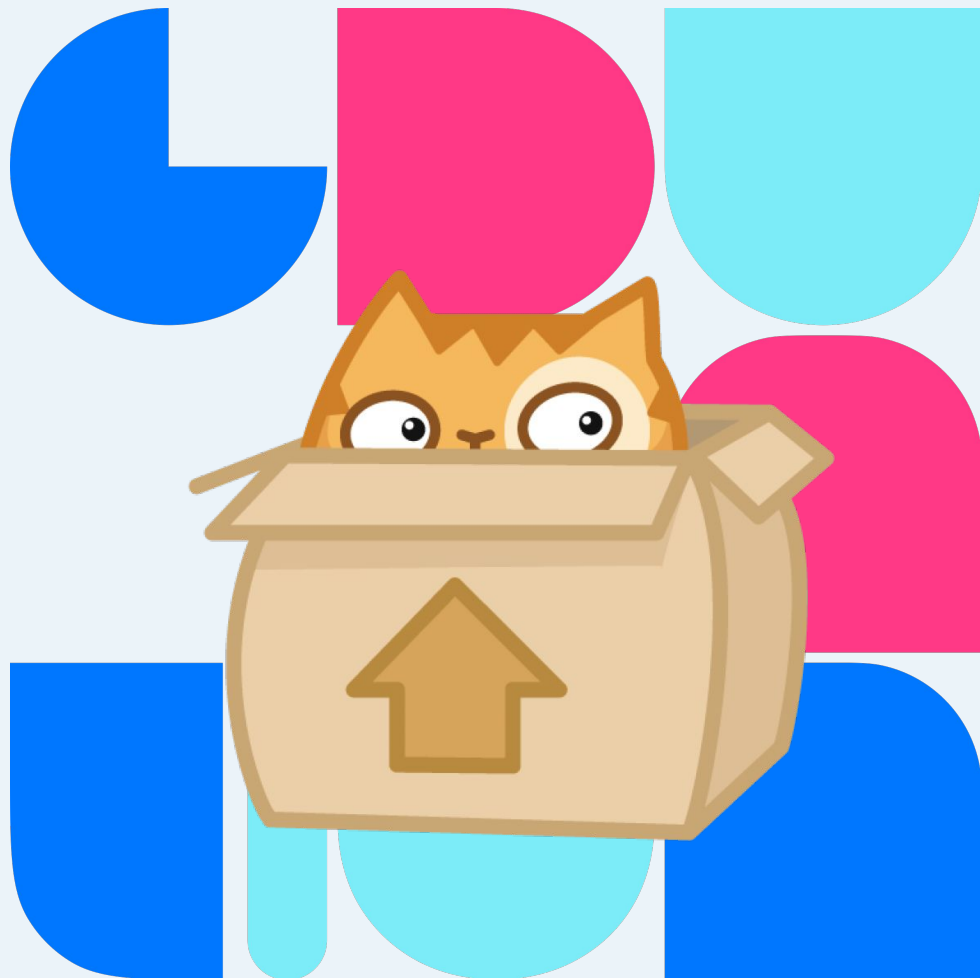
Моделирование языка (Language Modeling)



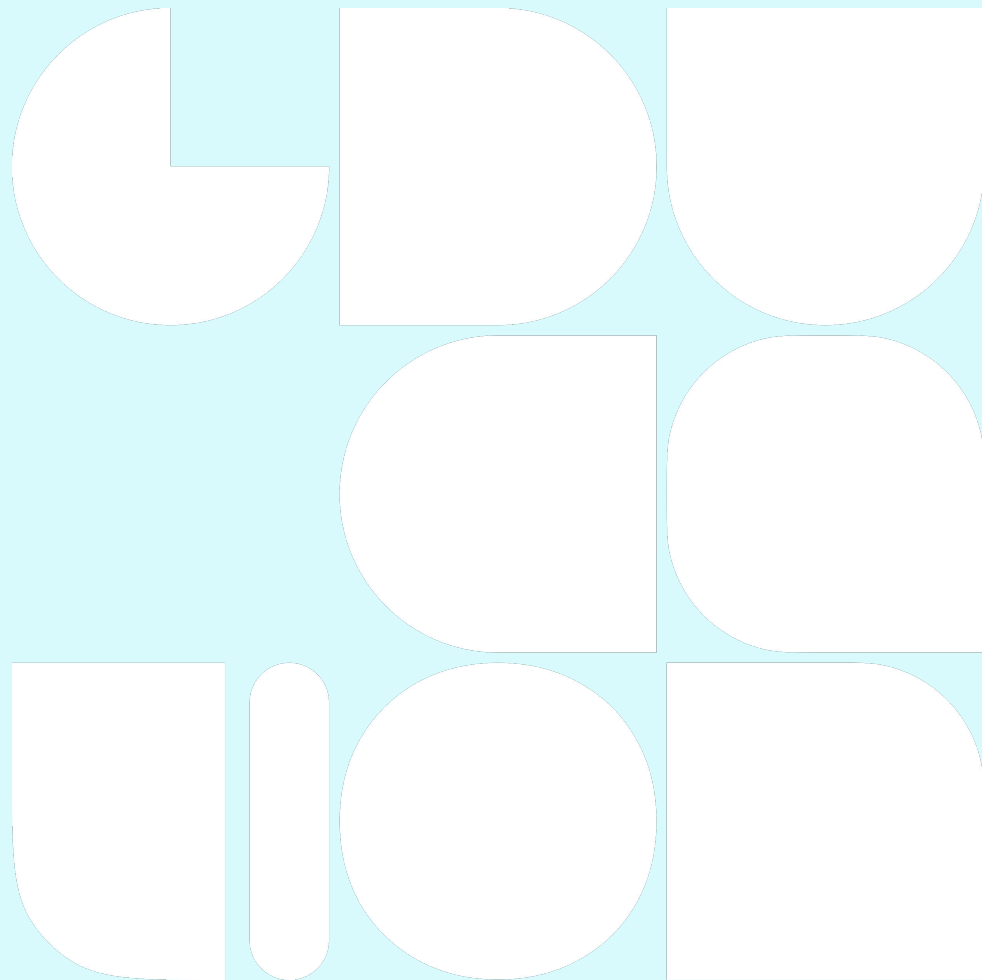
Итоги раздела:

- Задачи NLP очень разнообразны и тесно связаны с соседними доменами
- Обработка естественного языка сейчас крайне востребована и используется во многих известных продуктах
- Тексты можно обрабатывать с разной гранулярностью (полный текст/предложения/слова/токены/etc).

Вопросы?



Векторные представления слов/предложений



“Неглубокие” методы получения эмбеддингов слов

One-Hot Vectors

- 1) В прошлом году кот сбегал три раза
- 2) Быть или не быть, вот в чем вопрос

[illegible]

“Неглубокие” методы получения эмбедингов слов

One-Hot Vectors

Плюсы:

- Очень просто реализовать

Минусы:

- Всегда приходится хранить столько ключей, сколько всего уникальных слов во всех документах коллекции → огромная размерность
- Не учитывается контекст

“Неглубокие” методы получения эмбедингов слов

Bag-of-Words

- 1) Мой кот не любит купаться. В прошлом году кот сбежал три раза
- 2) Быть или не быть, вот в чем вопрос



1)

мой	кот	не	любит	купаться	в	прошлом	году	сбежал	три	раза
1	2	1	1	1	1	1	1	1	1	1

2)

быть	или	не	вот	в	чем	вопрос
2	1	1	1	1	1	1

“Неглубокие” методы получения эмбеддингов слов

Bag-of-Words

Плюсы:

- Очень просто реализовать

Минусы:

- Всегда сильный перекося в сторону частиц и союзов (и, в, бы, но и т.д.)
- Не учитывается порядок слов
- В худшем случае приходится только на один документ хранить столько ключей, сколько слов в нем

“Неглубокие” методы получения эмбеддингов слов

TF-IDF

Идея: Если слово часто встречается во всех документах, то вряд ли эти слова имеют большое значение. И наоборот, если слово встречается редко, то скорее всего оно в большей степени определяет смысл.

$$\text{tf}(t, d) = \frac{n_t}{\sum_k n_k}$$

Term Frequency

$$\text{idf}(t, D) = \log \frac{|D|}{|\{d_i \in D \mid t \in d_i\}|}$$

Inverse Document Frequency

$$\text{tf-idf}(t, d, D) = \text{tf}(t, d) \times \text{idf}(t, D)$$

n_t - число вхождений слова t в документ

D - коллекция документов

“Неглубокие” методы получения эмбеддингов слов

TF-IDF

Плюсы:

- Учитывает частоту слов как внутри документа, так и относительно всей коллекции
- Просто реализуется
- Можно интерпретировать как “веса” слов/токенов

Минусы:

- Все еще плохо хранит информацию о контексте

Word2Vec

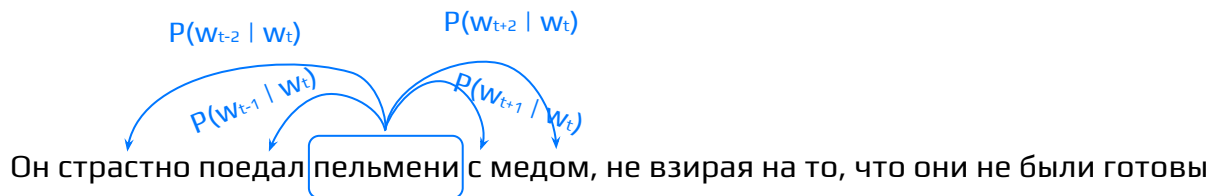
Идея

Будем обучать векторные представления слов так, чтобы по ним можно было определить контекст

Word2Vec

Идея

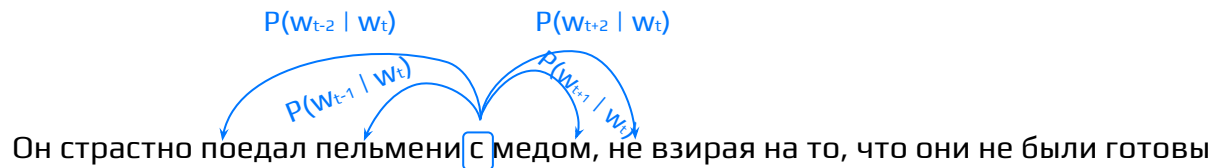
Будем обучать векторные представления слов так, чтобы по ним можно было определить контекст



Word2Vec

Идея

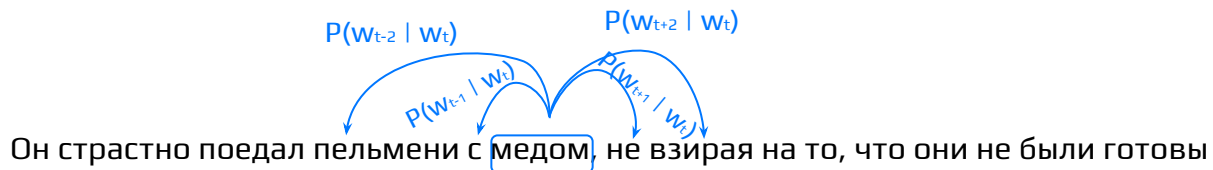
Будем обучать векторные представления слов так, чтобы по ним можно было определить контекст



Word2Vec

Идея

Будем обучать векторные представления слов так, чтобы по ним можно было определить контекст



Word2Vec

Objective

$$\text{Likelihood} = L(\theta) = \prod_{t=1}^T \prod_{\substack{-m \leq j \leq m, \\ j \neq 0}} P(w_{t+j} | w_t, \theta)$$

$$\text{Loss} = J(\theta) = -\frac{1}{T} \log L(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m, \\ j \neq 0}} \log P(w_{t+j} | w_t, \theta)$$

Word2Vec

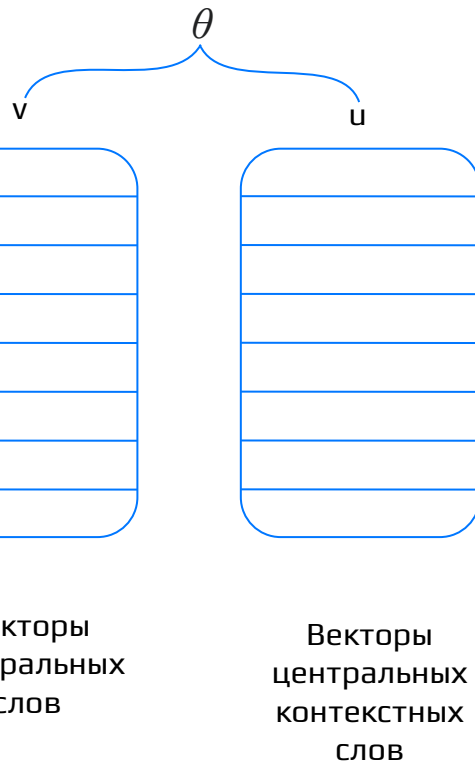
Objective

$$\text{Likelihood} = L(\theta) = \prod_{t=1}^T \prod_{\substack{-m \leq j \leq m, \\ j \neq 0}} P(w_{t+j} | w_t, \theta)$$

$$\text{Loss} = J(\theta) = -\frac{1}{T} \log L(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m, \\ j \neq 0}} \log P(w_{t+j} | w_t, \theta)$$

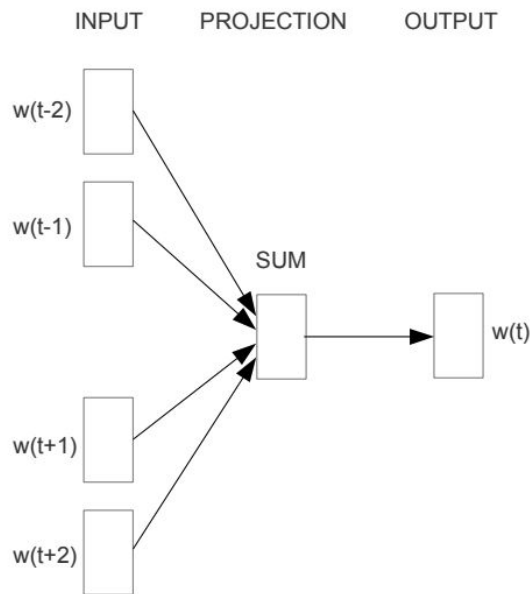
vocab_size

$$P(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$$

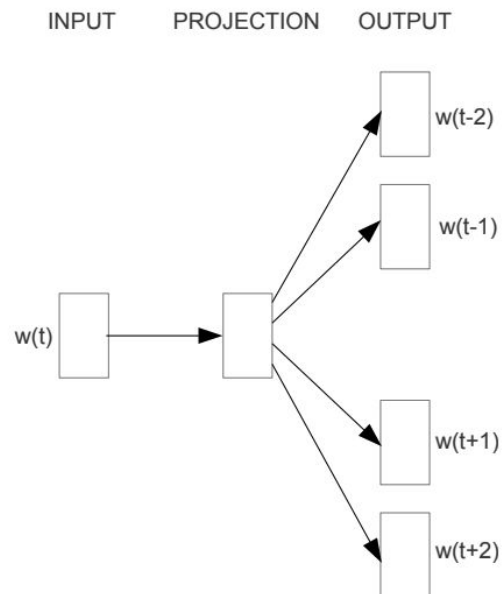


Word2Vec

Continuous Bag-of-Words (CBOW) vs Skip-Gram



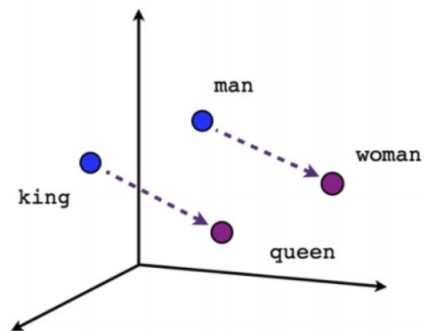
CBOW



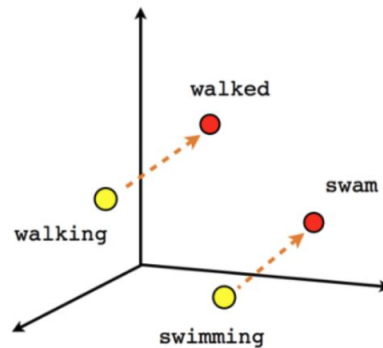
Skip-gram

Word2Vec

Арифметика эмбеддингов



Male-Female



Verb tense

$$\begin{aligned} e(\langle \text{king} \rangle) - e(\langle \text{man} \rangle) + e(\langle \text{woman} \rangle) &\simeq e(\langle \text{queen} \rangle) \\ e(\langle \text{swimming} \rangle) + e(\langle \text{walked} \rangle) - e(\langle \text{walking} \rangle) &\simeq e(\langle \text{swam} \rangle) \end{aligned}$$

Word2Vec

Плюсы:

- Умеет хранить семантическую информацию о словах

Минусы:

- Обычно использует целые слова, что сильно подгоняет модель под заданный корпус текстов. Чтобы использовать модель на новом наборе данных, нужно ее снова обучать на нем
- Все еще не использует всю информацию о контексте. Только в пределах заданного скользящего окна.

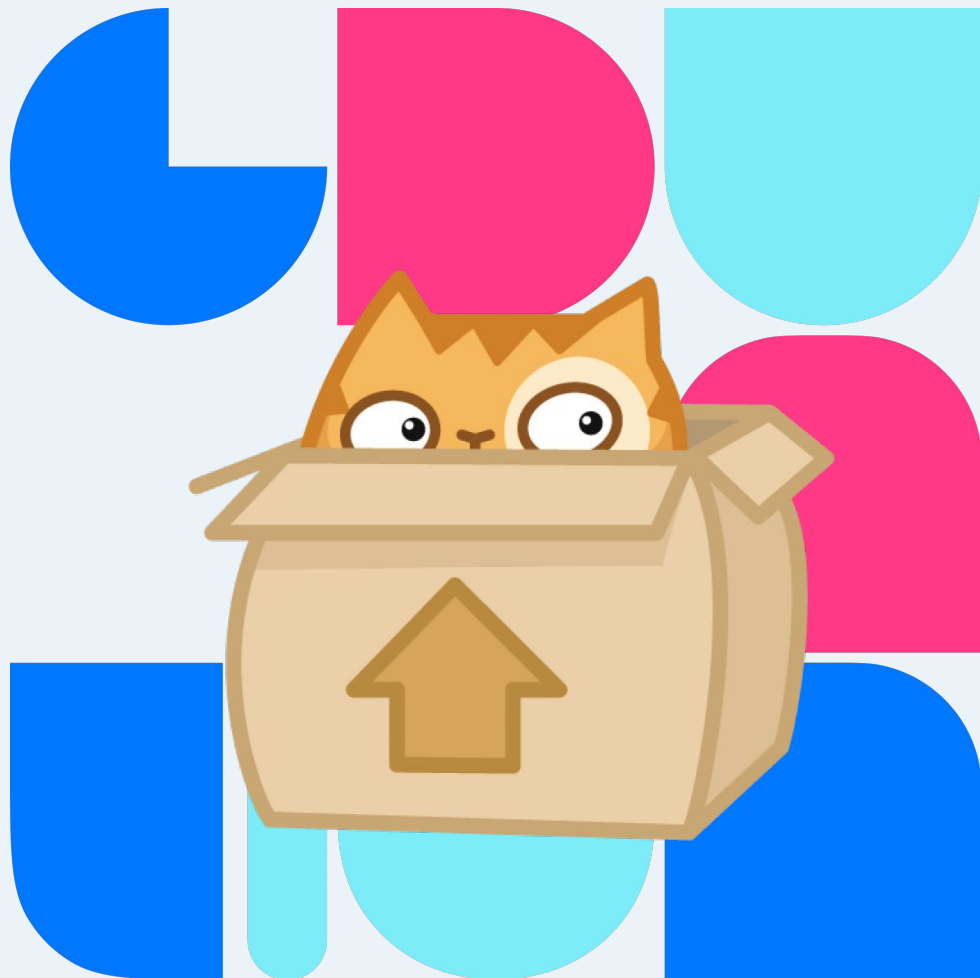
fasttext

Идея: используем в Word2Vec вместо слов n-граммы

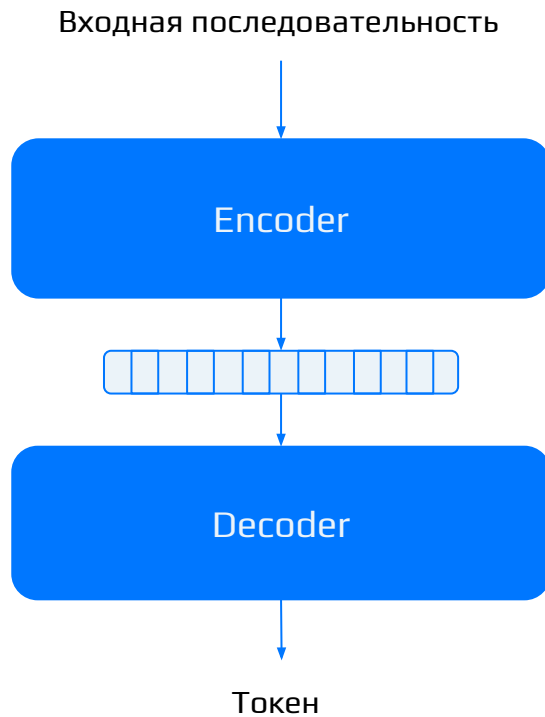
Пример: 3-граммы

“Нейросеть” → [“Ней”, “ейр”, “йро”, “рос”, “осе”, “сет”, “еть”]

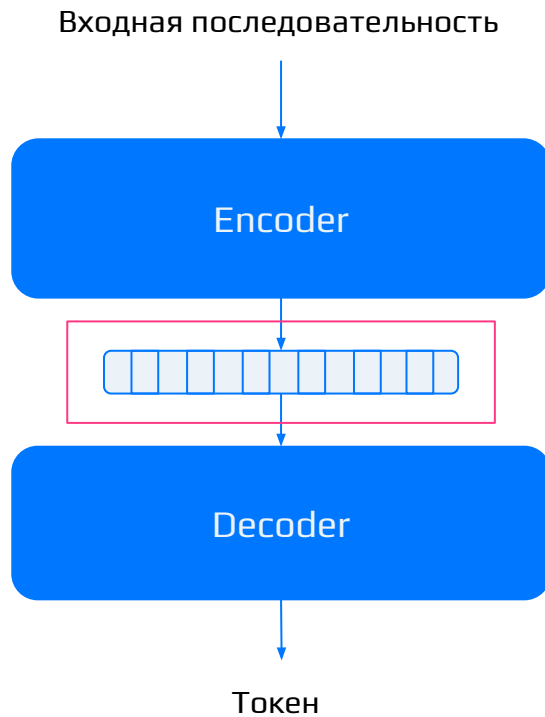
Вопросы?



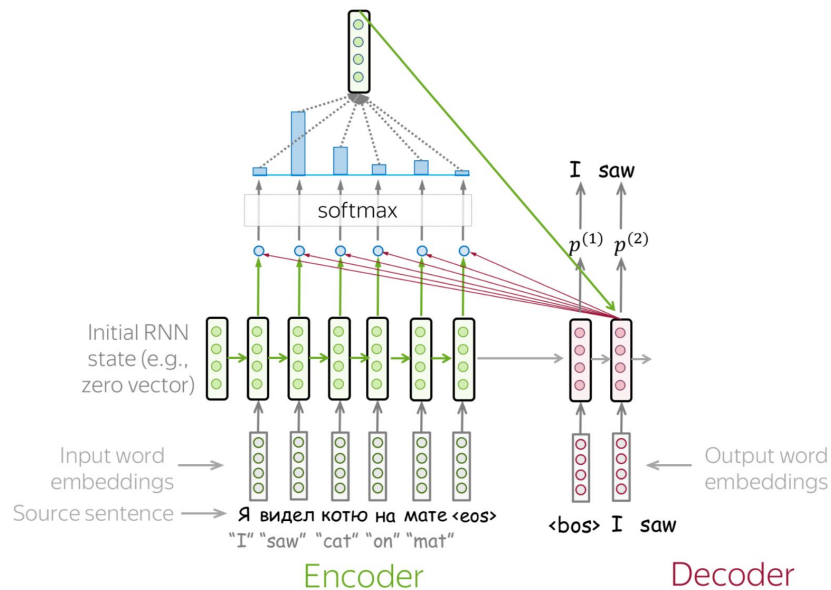
Attention



Attention



Attention



Self-Attention

Each vector receives three representations ("roles")

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

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

"Hey there, do you have this information?"

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

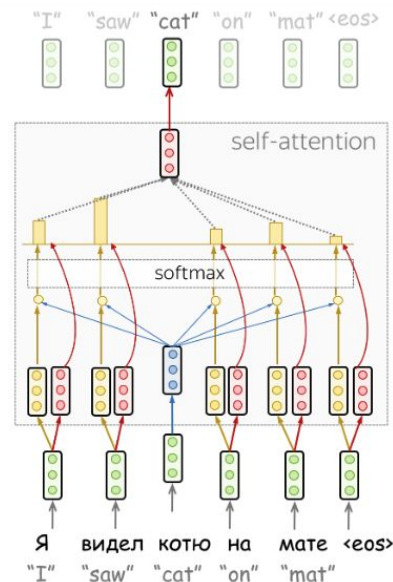
Key: vector **at** which the query looks to compute weights

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

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

Value: their weighted sum is attention output

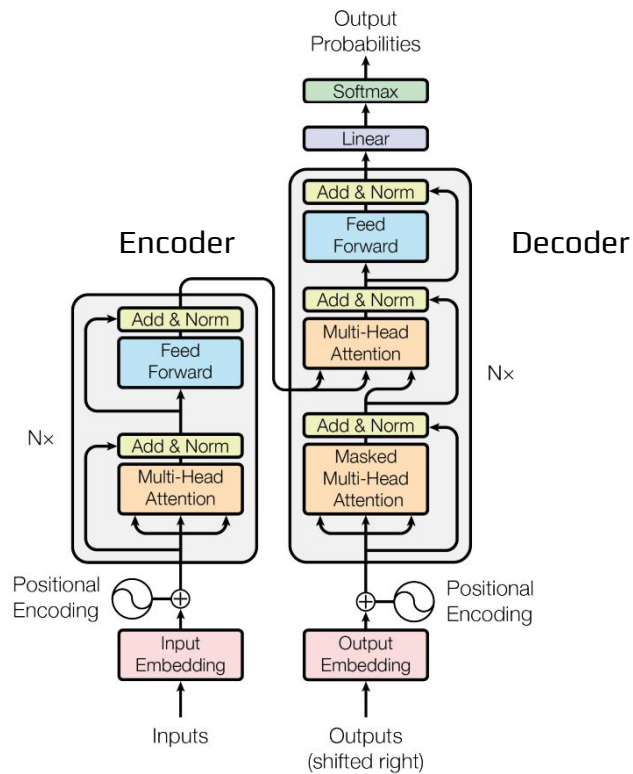
"Here's the information I have!"



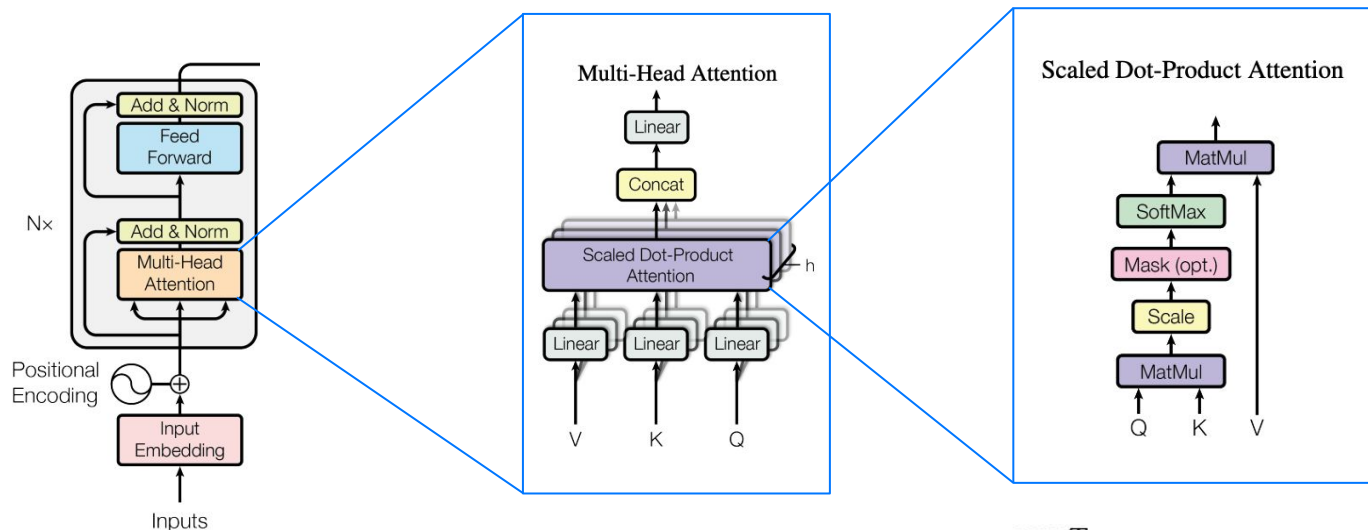
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Layer Type	Complexity per Layer
Self-Attention	$O(n^2 \cdot d)$
Recurrent	$O(n \cdot d^2)$
Convolutional	$O(k \cdot n \cdot d^2)$
Self-Attention (restricted)	$O(r \cdot n \cdot d)$

Transformer



Transformer

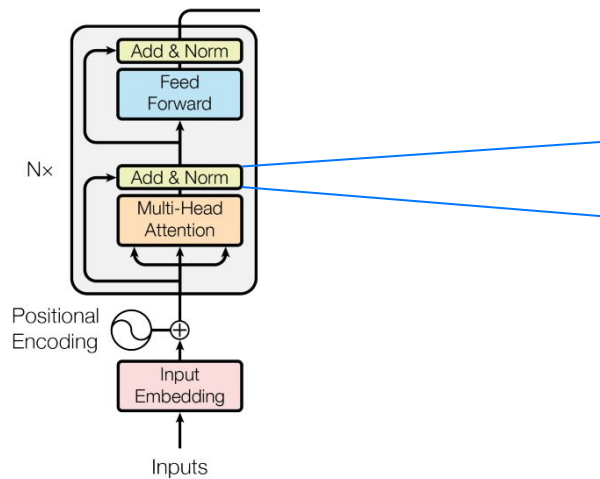


$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

where $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$

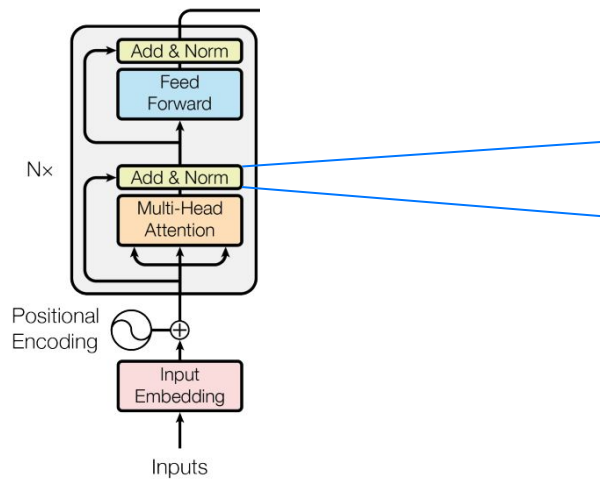
Transformer



Add == Skip-connection element-wise addition

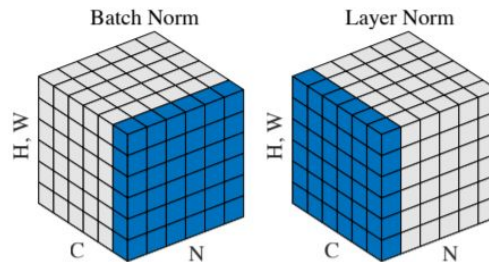
Norm == LayerNorm

Transformer



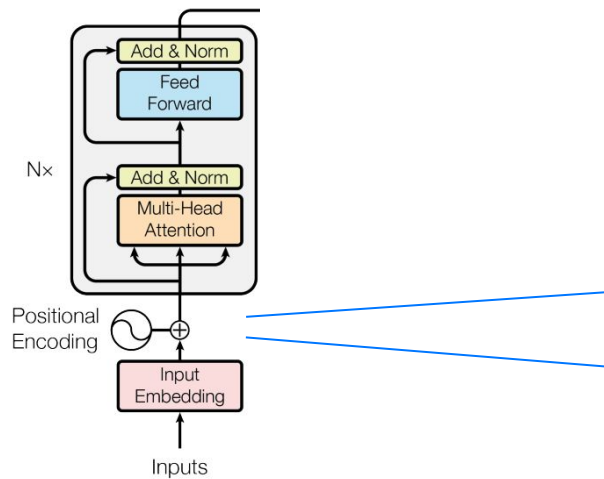
Add == Skip-connection element-wise addition

Norm == LayerNorm



$$y = \frac{x - \mathbf{E}[x]}{\sqrt{\mathbf{Var}[x] + \epsilon}} * \gamma + \beta$$

Transformer

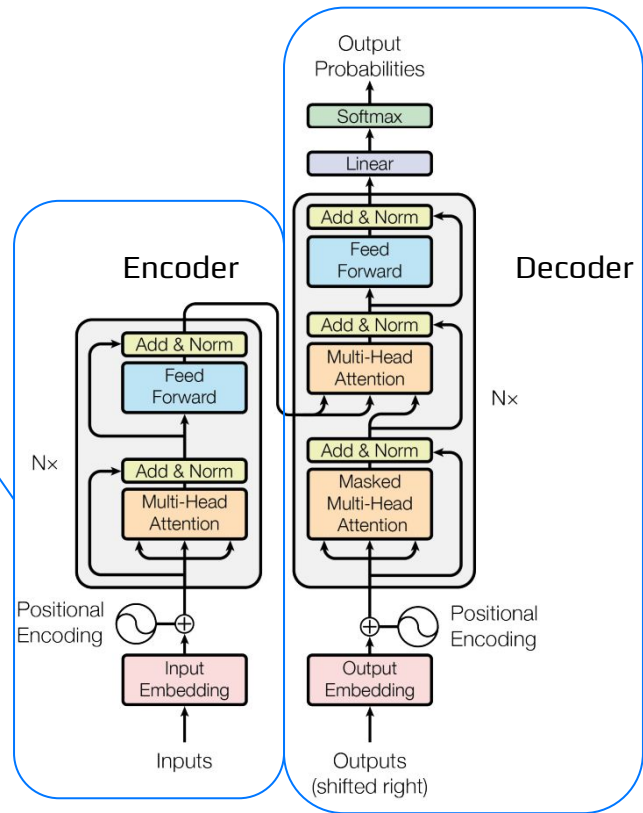


$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

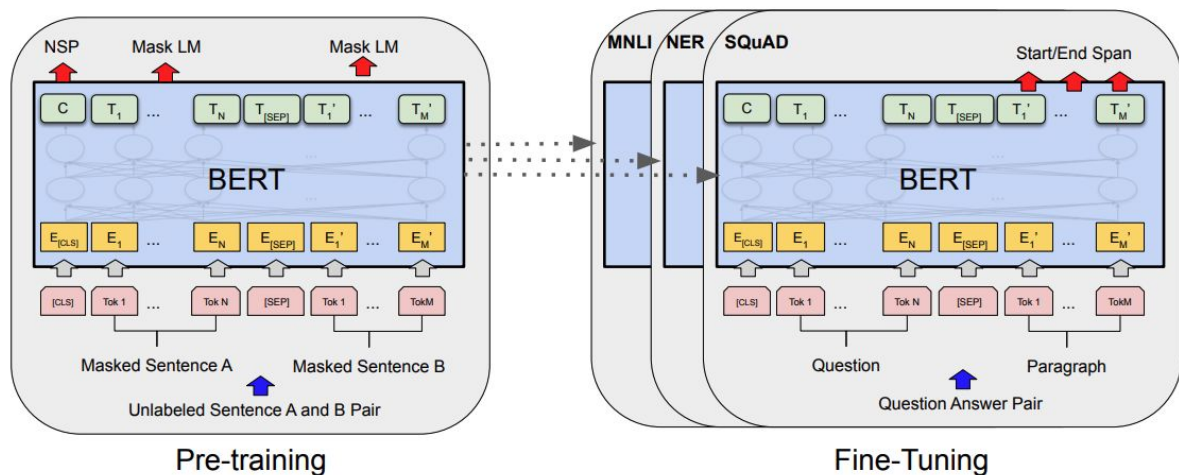
Transformer

BERT-like
encoders



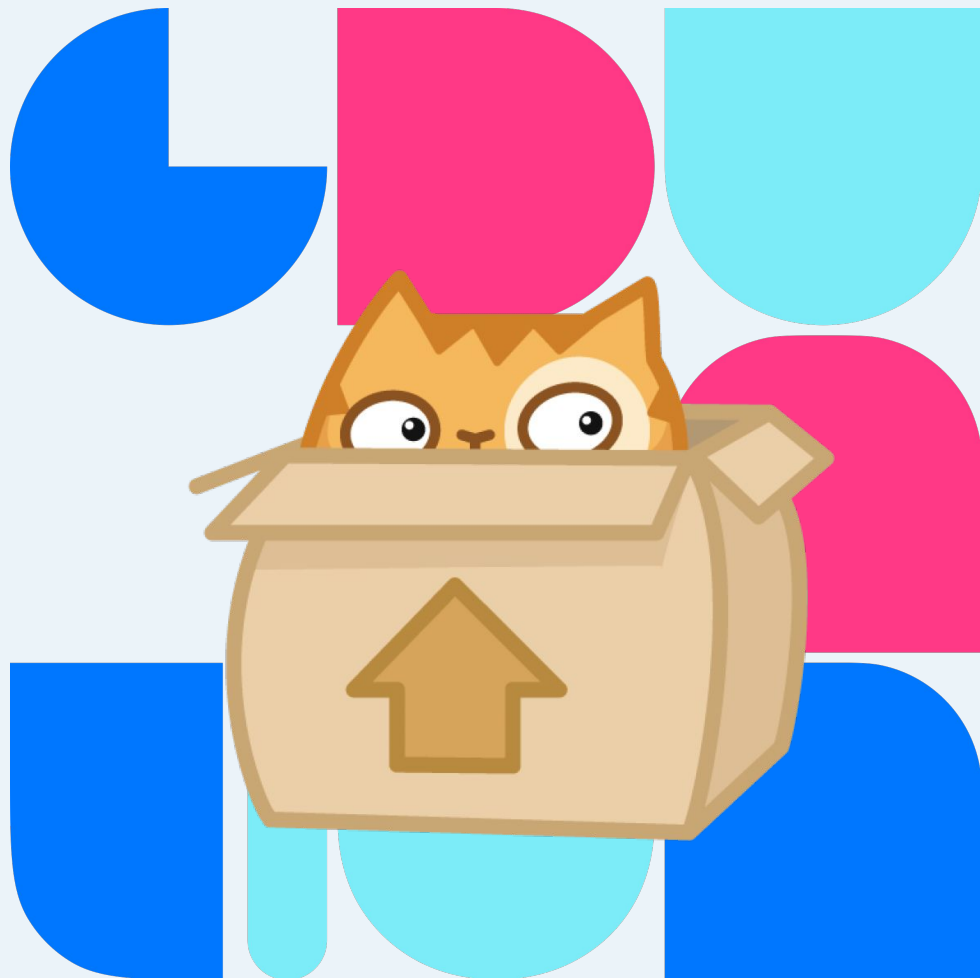
GPT-like
decoders

Bidirectional Encoder Representations from Transformers (BERT)



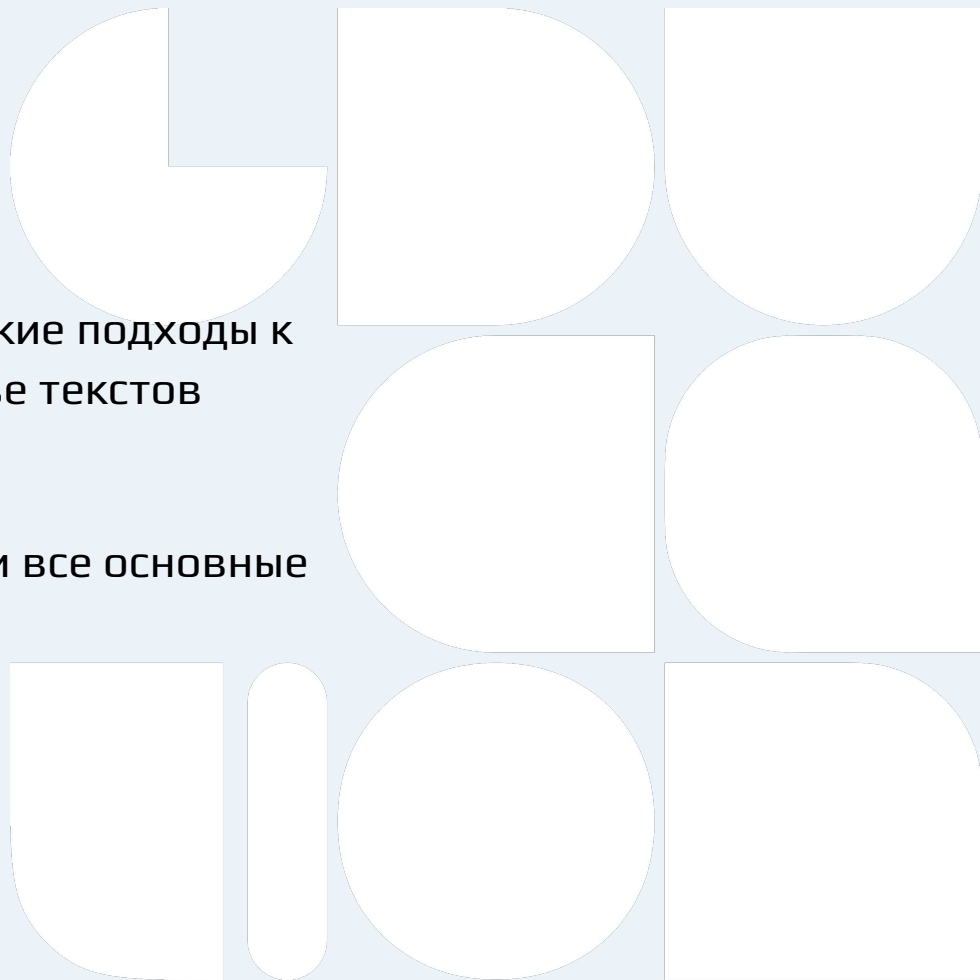
- MASK - некоторые слова заменяем на токен неизвестного слова и пытаемся их восстановить.
- NSP - Для пары предложений пытаемся предсказать, правда ли, что B следует за A. (берем B случайно в 50% случаев)
- Нужно для улучшения модели языка и вопросно-ответных задач. Обучающее множество включает всю английскую википедию и книги не защищенные авторским правом. Для большой модели надо 4 дня на 16-и cloud TPU.

Вопросы?

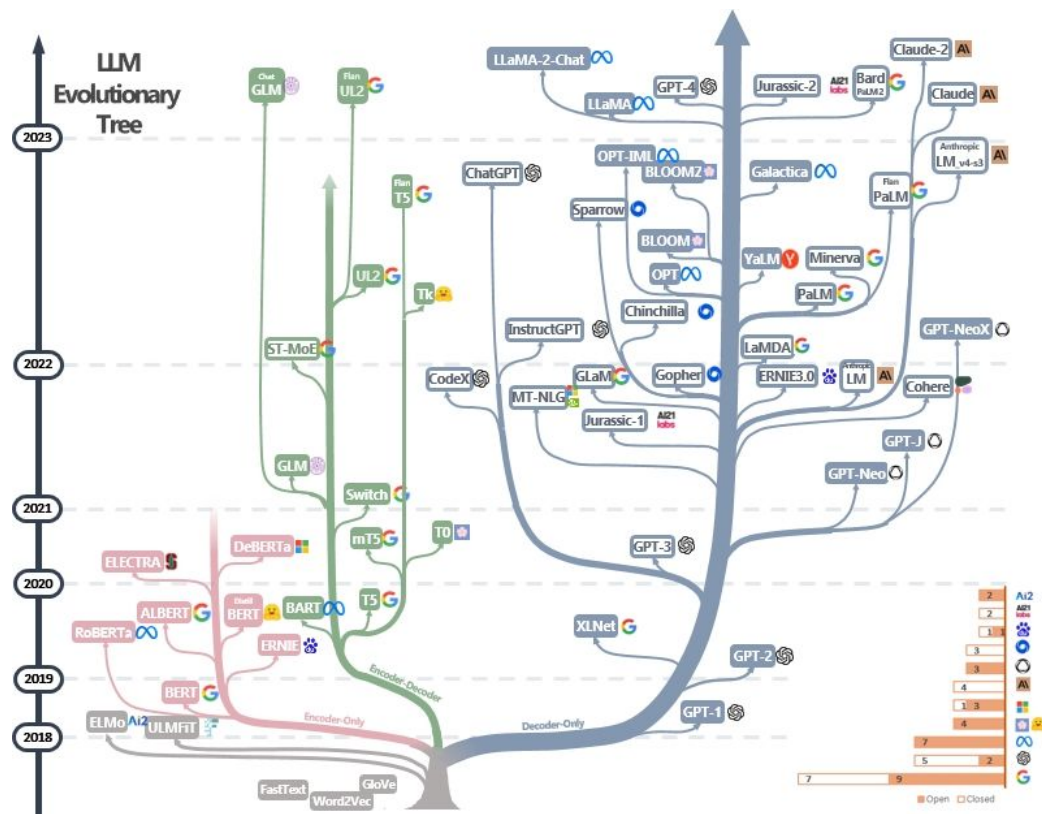


Итоги раздела:

- Рассмотрены все основные неглубокие подходы к формированию признаков на основе текстов
- Пройден путь от Word2Vec до BERT
- Разобраны attention, self-attention и все основные компоненты трансформеров



Что дальше?



Что дальше?

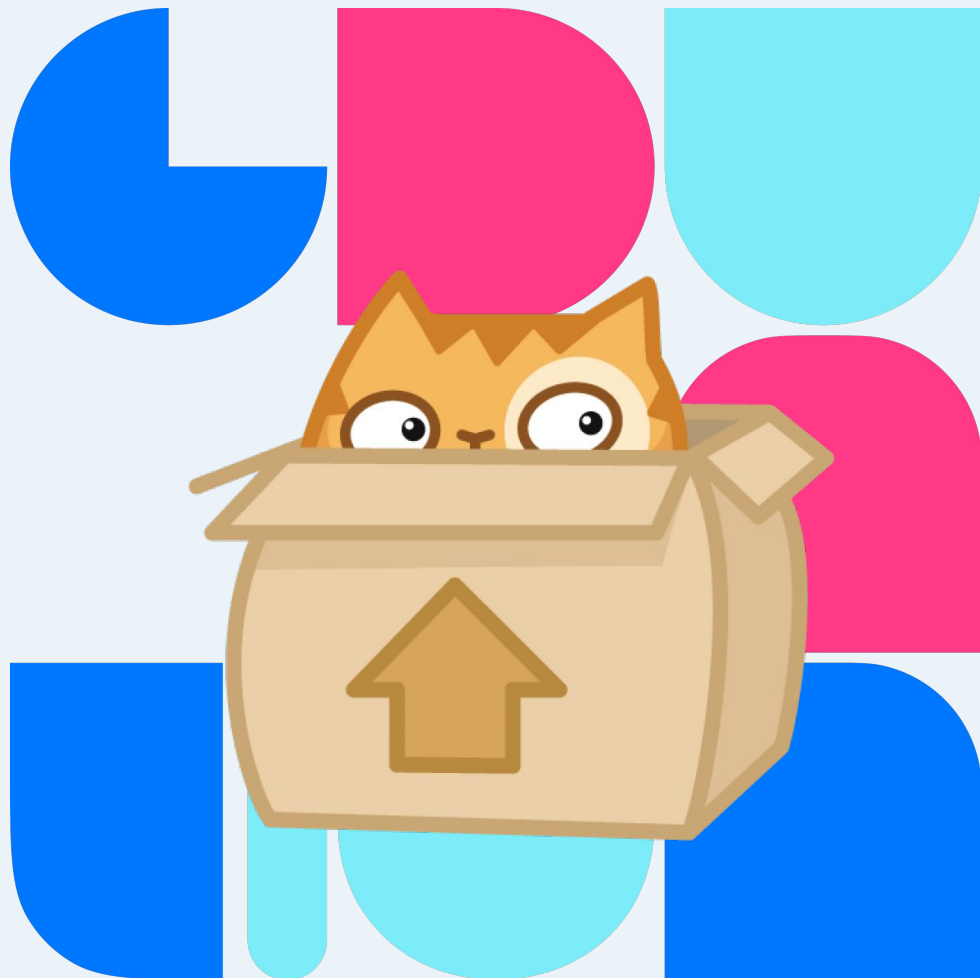
На следующей неделе

- Causal Language Modeling & GPT 1-3
- Метрики для оценки качества языковых моделей
- Становление LLM
- KV-cache, paged attention, sliding-window attention, prefix-caching
- PEFT (LoRA, p-tuning, prefix-tuning)

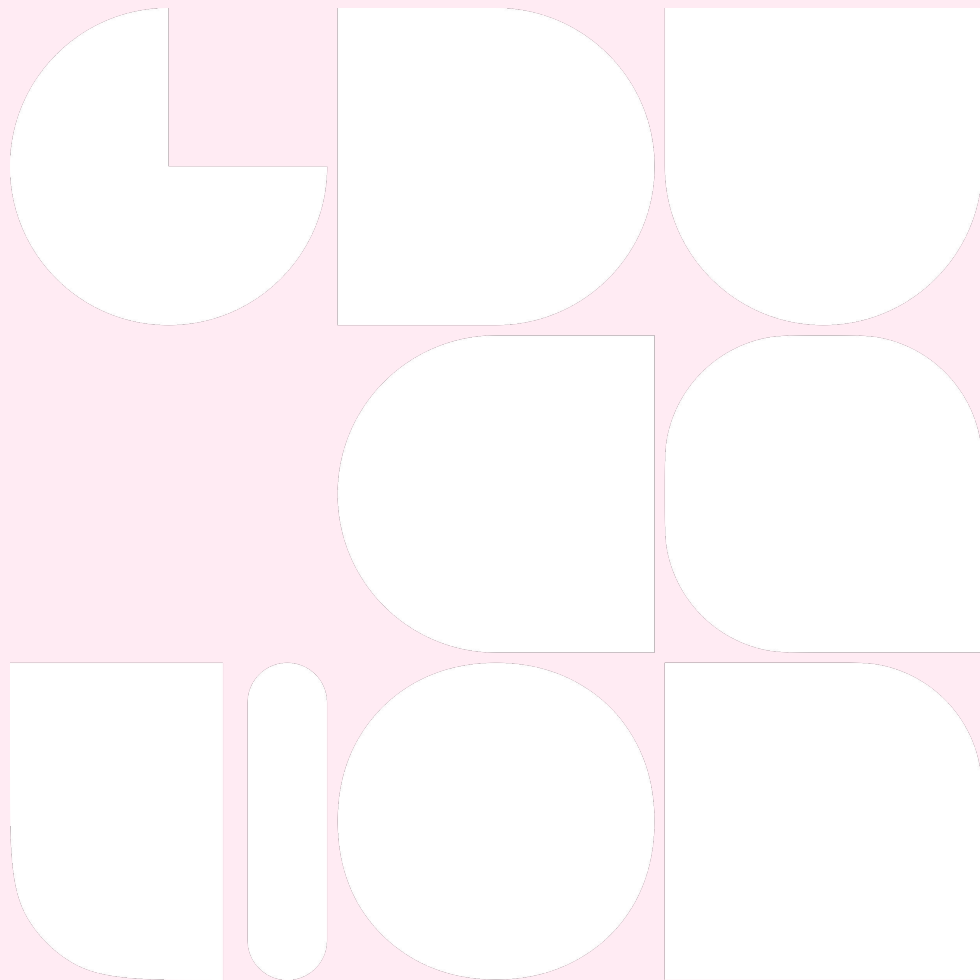
Подведём итоги:

- NLP покрывает все, что связано с векторными представлениями текстов, семантическим поиском и моделированием естественного языка
- Есть ряд необучаемых методов, которые до сих пор актуальны
- Word2Vec – первая нейросетевая модель для текстовых эмбеддингов, fasttext – его улучшенная версия на основе n-грамм
- Для учета информации из разных частей контекста используется оператор Attention
- Применяя Attention ко всем состояниям без RNN/CNN можно получить более гибкий учет контекста
- Трансформер – основная база всех современных языковых моделей, как для Masked LM, так и для Causal LM
- BERT – одна из ведущих архитектур среди моделей для получения текстовых эмбеддингов

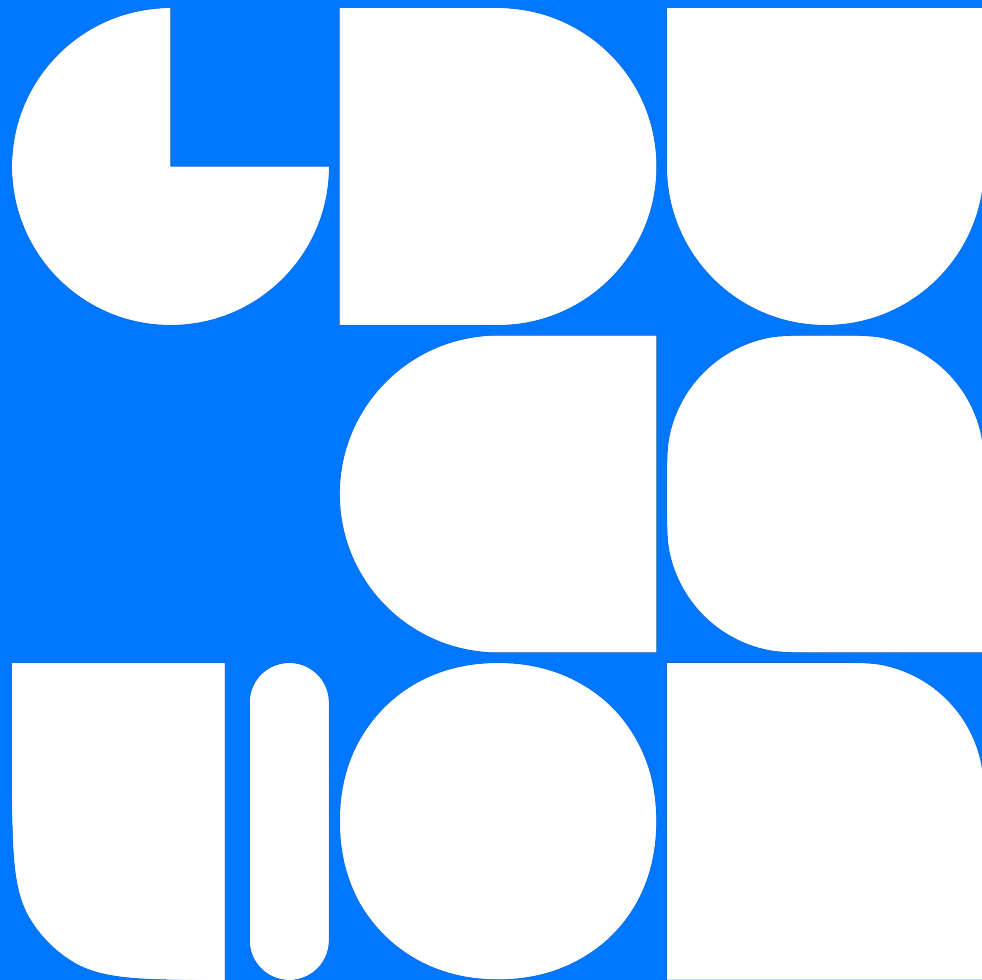
Вопросы?



Перерыв



Практика



Спасибо за
внимание!

