



Тверской  
государственный  
технический  
университет

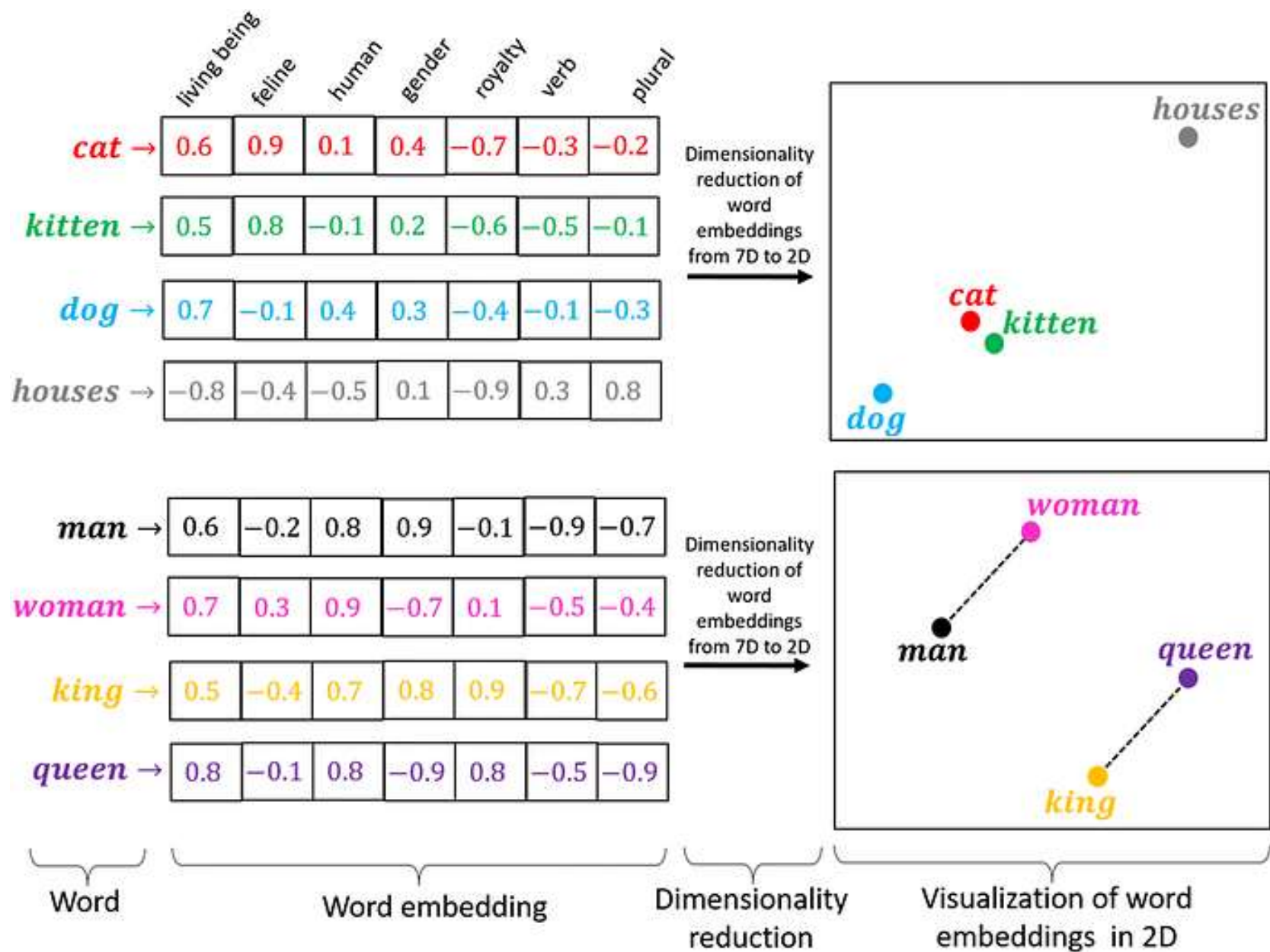
# Интеллектуальные информационные системы

## Модели внимания и трансформеры

Материалы курса доступны по ссылке:

<https://github.com/AndreyShpigar/ML-course>

2024 г.



- Текст (документ) – последовательность токенов
- Токенизация – представление текста в виде последовательности
- Токен – последовательность символов (слово или часть слова)

text → features → ML model →  $P(y|x)$



Как представить текст  
в виде, который могла  
бы понять модель?

# Представление текста

Необходимо кодировать текстовые данные – преобразовать их в вектор, прежде чем подавать на вход нейронной сети.

Существует два подхода:

1. Векторизовать текст целиком, превращая его в один вектор
2. Векторизовать отдельные структурные единицы, превращая текст в последовательность векторов (embedding)

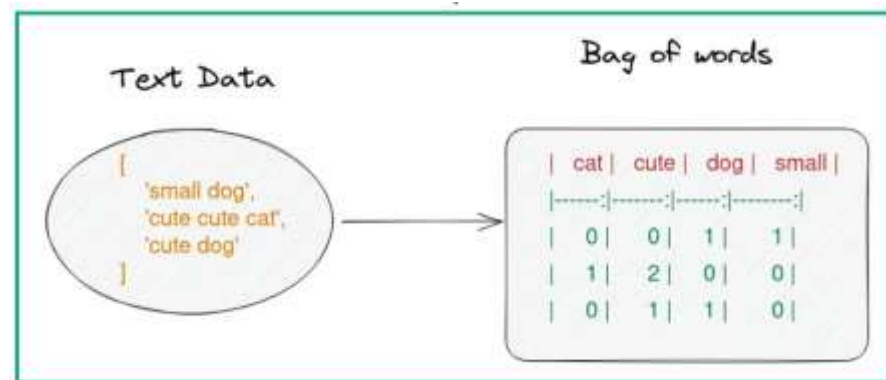
- Bag-of-Words (мешок слов) – текст представляется в виде вектора частот встречаемости каждого токена, кроме элементов заранее заданного списка «стоп-слов», в которые обычно включают самые вездесущие токены: личные местоимения, артикли и так далее
- Представление мешка слов – таблица с числами, в которой столбцы – уникальные слова, а в ячейках находится число вхождений слова в документ. В каждой строке получается набор чисел (вектор), характеризующий состав документа
- Не учитывается порядок слов – два предложения могут называться одинаковыми, если содержат один и тот же набор слов

## The Bag of Words Representation

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!



it 6  
I 5  
the 4  
to 3  
and 3  
seen 2  
yet 1  
would 1  
whimsical 1  
times 1  
sweet 1  
satirical 1  
adventure 1  
genre 1  
fairy 1  
humor 1  
have 1  
great 1  
...



- TF-IDF (Term Frequency - Inverse Document Frequency) – статистическая мера, используемая для оценки важности слова в контексте документа, являющегося частью коллекции документов. Вес некоторого слова пропорционален частоте употребления этого слова в документе и обратно пропорционален частоте употребления слова во всех документа коллекции
- TF (Term Frequency, частота термина) – обозначает, насколько часто определенное слово появляется в данном документе (частота вхождения токена в документ).
- TF измеряет важность слова в контексте отдельного документа

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

где  $n_t$  – число вхождений токена  $t$  в документ,  
а в знаменателе стоит общее число слов в данном документе

- IDF (Inverse document frequency) – измеряет, насколько уникально слово является по всей коллекции документов. Слова, которые в большинстве документов, имеют низкое IDF, так как не вносят большой информационной ценности

$$IDF(t, D) = \log \frac{|D|}{|\{d_i \in D | t \in d_i\}|}$$

где  $\{d_i \in D | t \in d_i\}$  – число документов в текстовой коллекции  $D$ , в которых встречается слово  $t$

- Этот множитель штрафует компоненты, отвечающие слишком распространенным токенам и повышает вес специфических для отдельных текстов (и, вероятно, информативных) слов

Мера TF-IDF часто используется в задача анализа текстов и информационного поиска

Преимущества TD-IDF:

- Учет важности слов – учитывает как частоту слова в документе, так и его общую редкость по всей коллекции. Выделяет ключевые слова, которые часто встречаются в данном документе, но не слишком распространены в остальных
- Устранение шума – слова, которые встречаются в большинстве документов - «стоп-слова» - имеют низкий общий вес TF-IDF

Недостатки TF-IDF:

- Отсутствие семантической информации – не учитывает семантические связи между словами
- Чувствительность к длине документа – длинные документы могут иметь более высокие значения TF, даже если ключевые слова встречаются реже



# Words embeddings

- Мотивация использовать контекст:

Допустим, что одно и то же слово будет представлено одним и тем же вектором во всех текстах и в любых позициях. Как заключить в векторе его смысл (информацию)?

Использование контекста позволяет понять смысл слова по словам, встречающимся рядом с ним

*«Я не успел к началу лекции, потому что опоздал на...»*

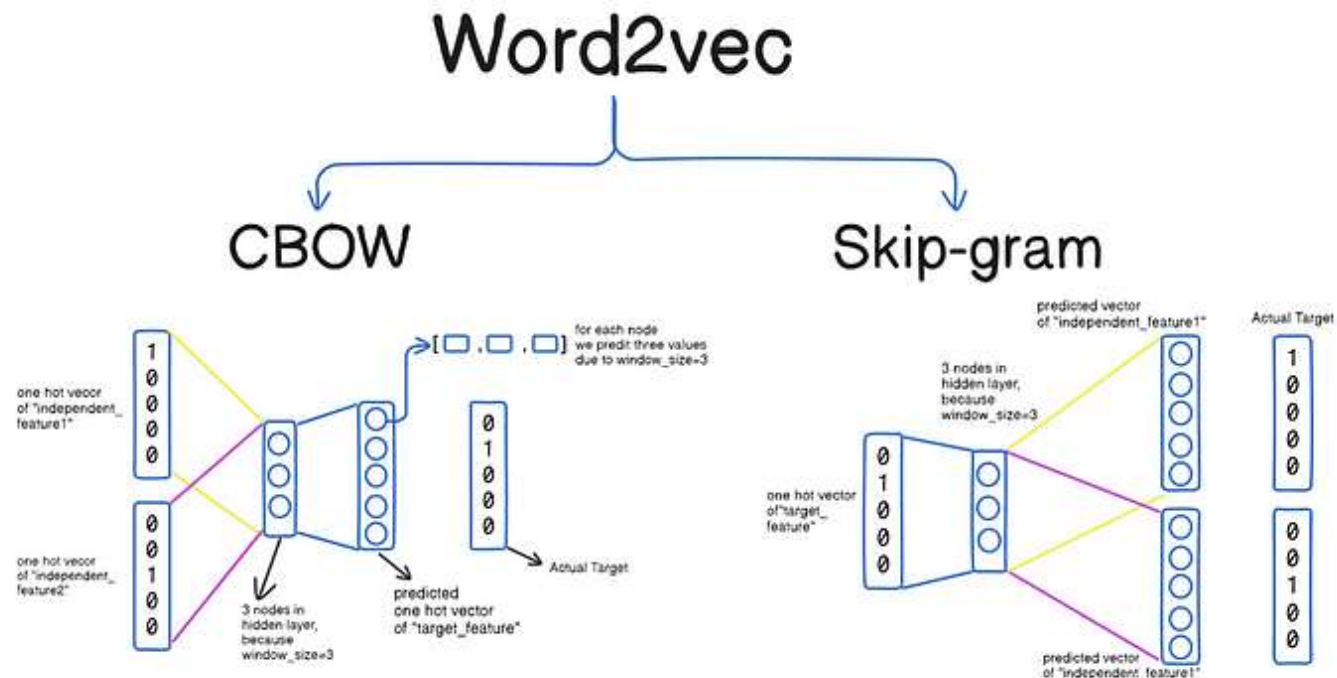
1. *Автобус – 0.6*
2. *Метро – 0.3*
3. *Такси – 0.05*
4. *Велосипед – 0.04*
5. *Самолет – 0.01*
6. *...*

# Word2vec

Реализация подхода с использованием контекста – word2vec

Предложен Томашем Миколовым в 2013 году в статье  
«Efficient Estimation of Word Representations in Vector Space»

Основная идея – упаковать информацию о контексте слов в вектора. Будем обучать вектора, пытаюсь предсказать контекст, в котором встречается слово

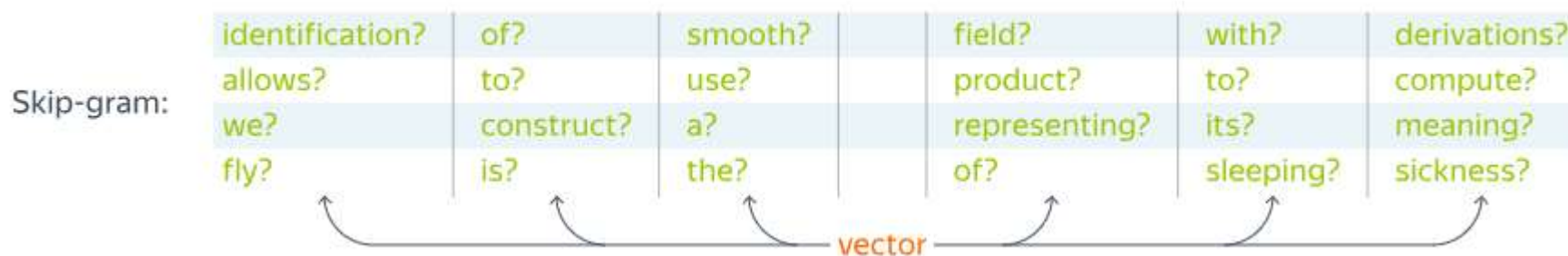


Предложено две стратегии – CBOW и Skip-gram

1. CBOW (continuous bag-of-words) – модель учится предсказывать данное (центральное) слово по контексту. Например, по двум словам перед данным и по двум словам после него.



2. Skip-gram – модель учится предсказывать контекст.  
Например, каждого из двух соседей слева и справа.



- Для каждого слова w обучается два эмбединга:  $v_u$  и  $v_w$ . Первый используется, когда w является центральным, второе – когда оно рассматривается как часть контекста.

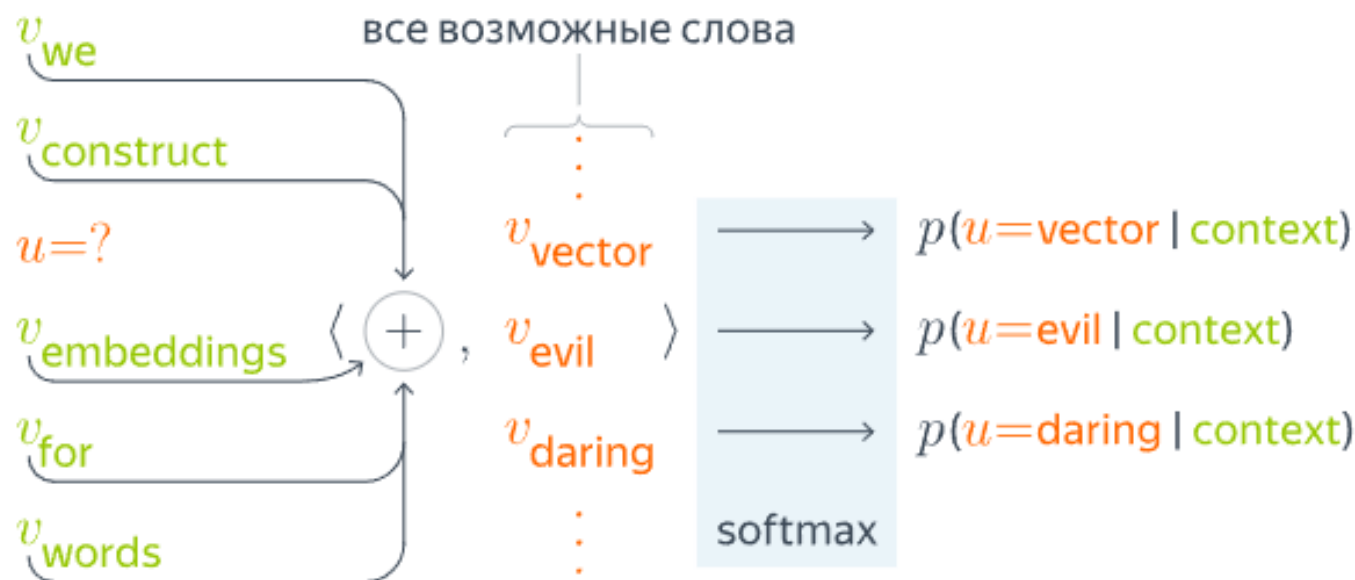
В модели CBOW при фиксированном контексте «*context*» вычисляются логиты

$$\text{logits}_u = \left\langle \sum_{w \in \text{context}} v_w, v_u \right\rangle$$

После чего «вероятности» всевозможных слов и быть центральным словом для «*context*» вычисляются как  $\text{softmax}(\text{logits})$

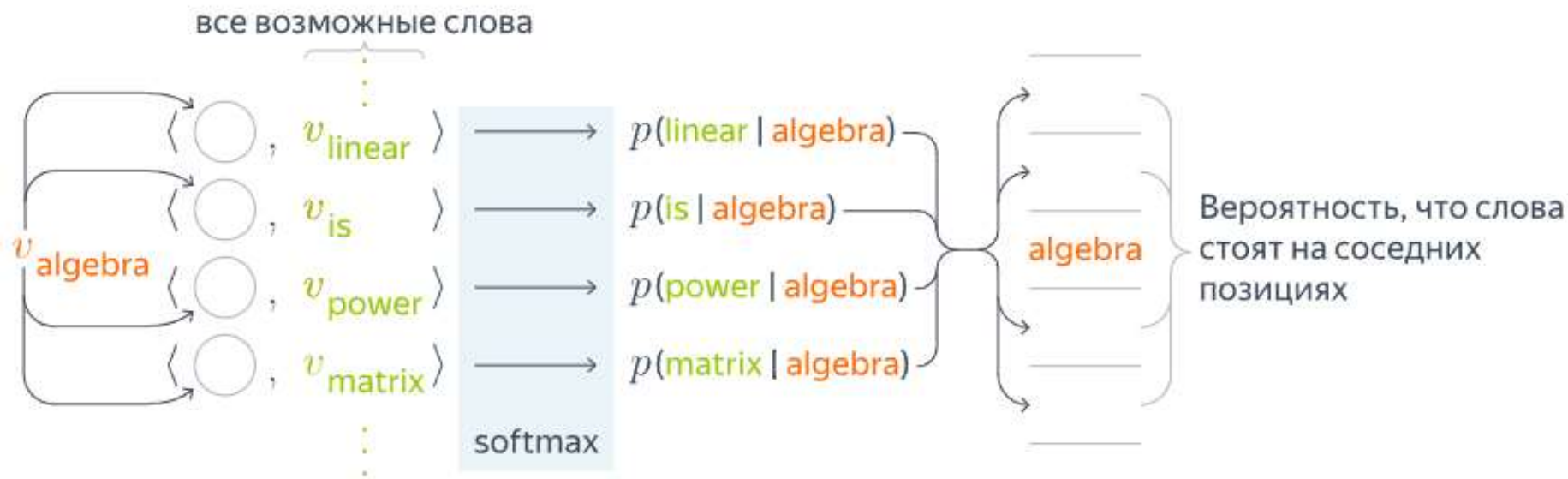
Модель учится с помощью SGD на кросс-энтропию полученного распределения с истинным распределением центральных слов

CBOW:



В модели Skip-gram по данному центральному слову  $u$  для каждой позиции контекста «context» независимо предсказывается распределение вероятностей. В качестве функции потерь выступает сумма кросс-энтропий распределений слов контекста с их истинными распределениями

Skip-gram:



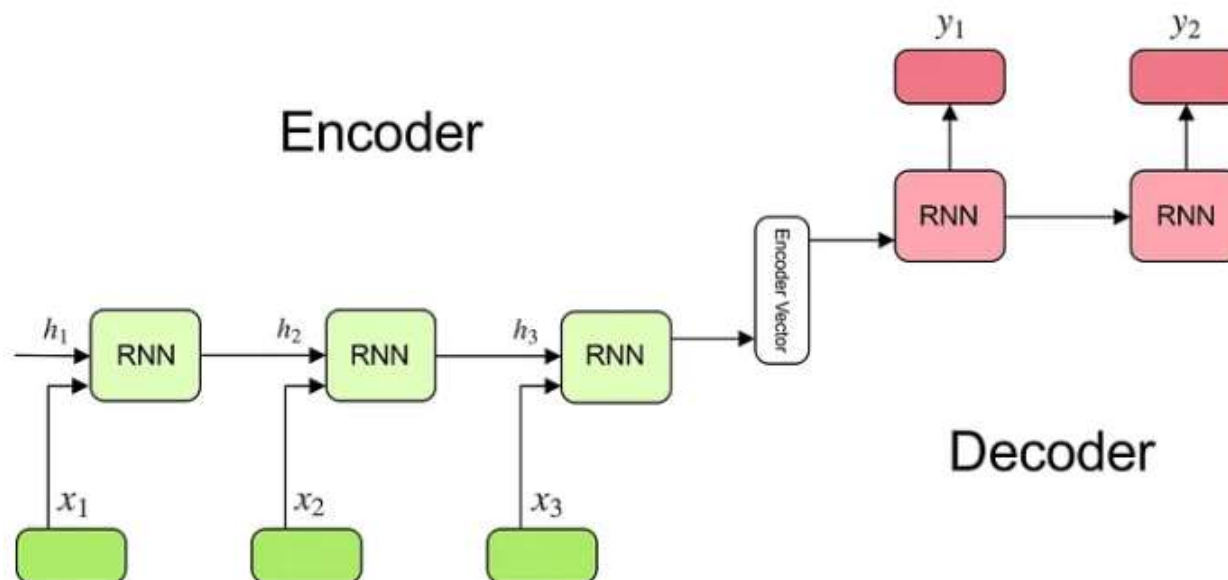
- Примеры. Возьмём несколько слов и посмотрим, как выглядят топ-10 слов, ближайших к ним в пространстве эмбедингов (обученных на одном из датасетов Quora Questions с помощью word2vec):
  1. **quantum**: electrodynamics, computation, relativity, theory, equations, theoretical, particle, mathematical, mechanics, physics;
  2. **personality**: personalities, traits, character, persona, temperament, demeanor, narcissistic, trait, antisocial, charisma;
  3. **triangle**: triangles, equilateral, isosceles, rectangle, circle, choke, quadrilateral, hypotenuse, bordered, polygon;
  4. **art**: arts, museum, paintings, painting, gallery, sculpture, photography, contemporary, exhibition, artist.
- Размерность эмбединга в каждой из архитектур — это гиперпараметр и подбирается эмпирически. В оригинальной статье предлагается взять размерность эмбединга 300. Полученные представления центральных слов могут дальше использоваться в качестве эмбедингов слов, которые сохраняют семантическую связь слов друг с другом.

# Sequence-to-sequence

Рекуррентная сеть для синтеза последовательностей (seq2sec) – архитектура, предназначенная для задач, в которых требуется преобразование одной последовательности в другую (машинный перевод, суммаризация текста, генерация описаний и т.д.). Состоит из энкодера и декодера.

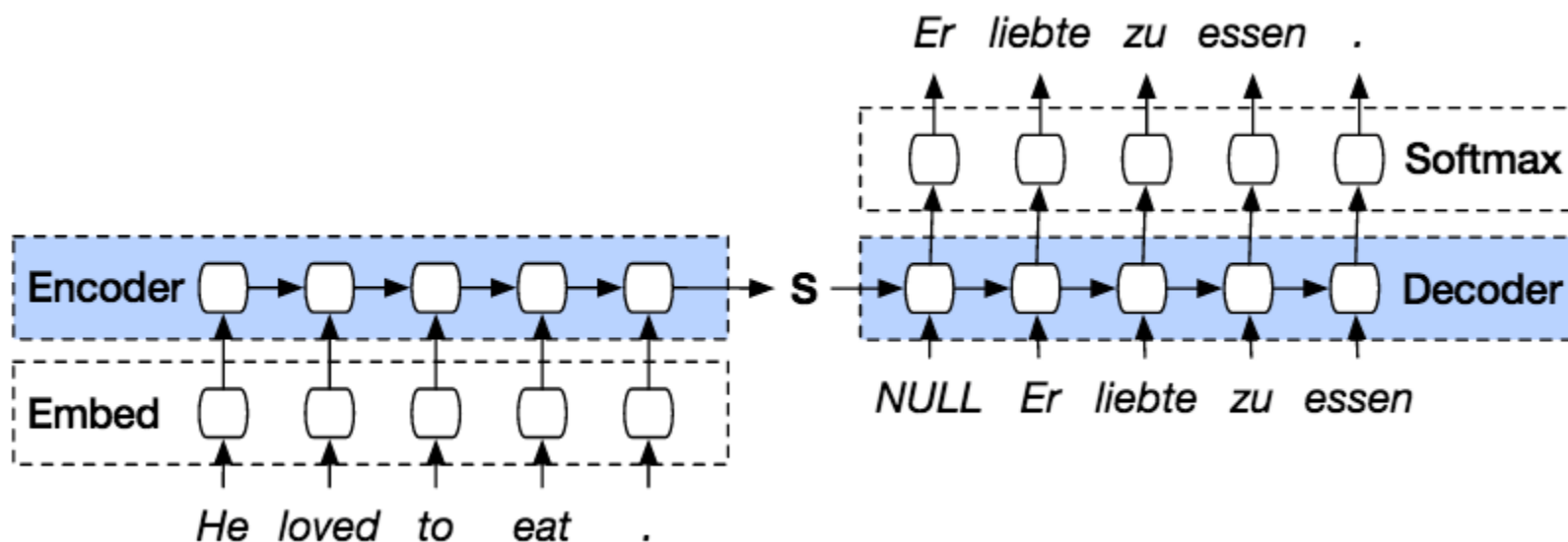
Энкодер (encoder) – кодирование информации об исходной последовательности в контекстный вектор (context vector, encoded vector)

Декодер (decoder) – преобразование закодированной энкодером информации в новую последовательность



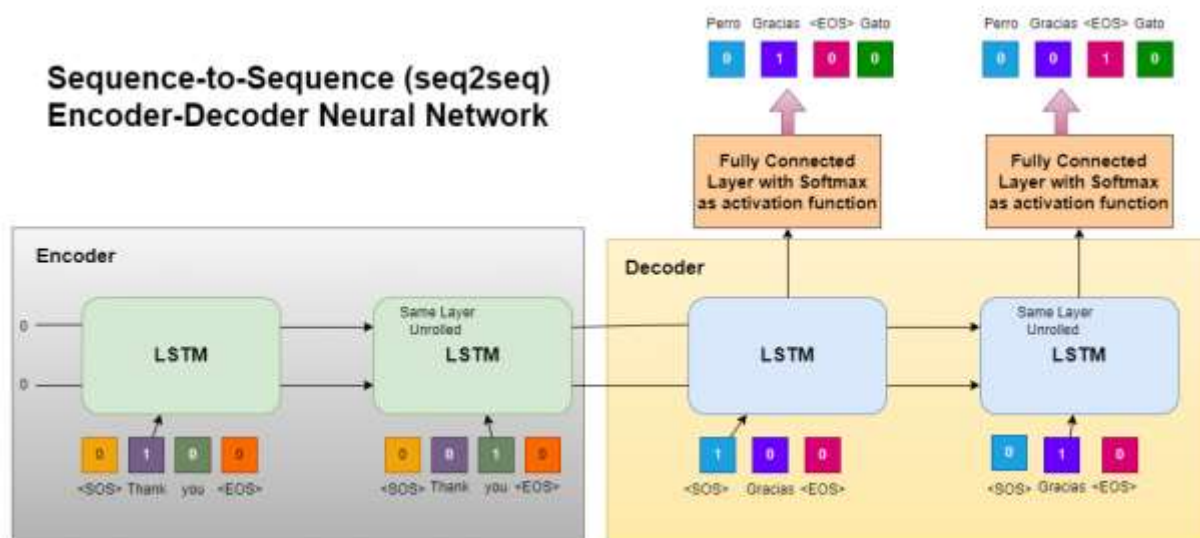


- Входные текстовые данные преобразуются в числовые представления с помощью эмбеддинга, который преобразует каждый токен во входной последовательности в вектор фиксированной размерности



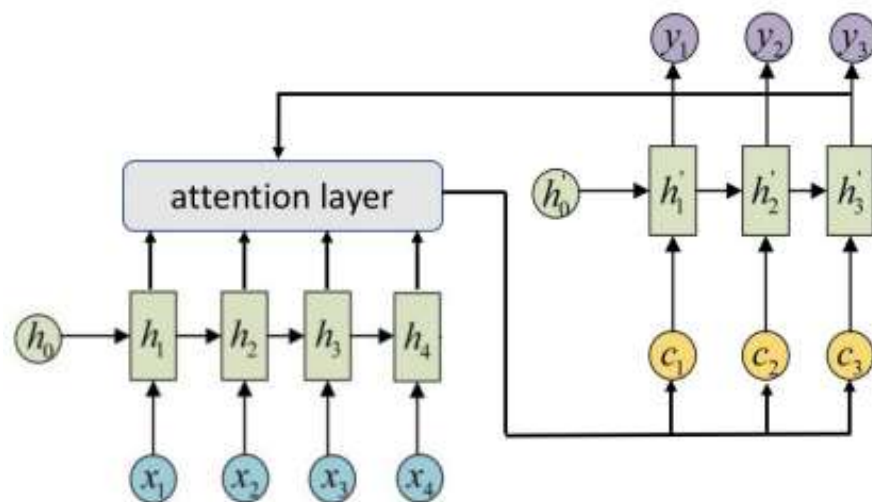
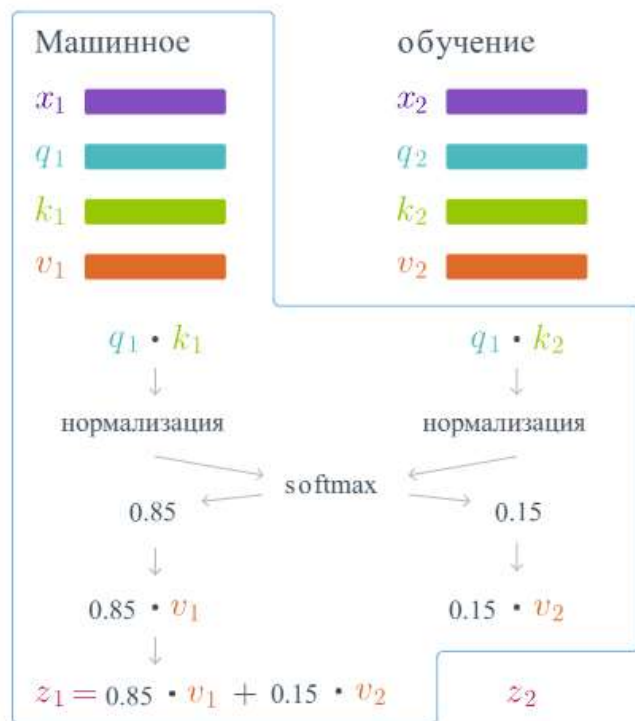
- Encoder – преобразует входную последовательность в контекстный вектор, содержащий обобщённое представление всей входной информации
- Основу энкодера RNN, обычно реализованные с использованием LSTM или GRU. Эти сети обрабатывают входную последовательность пошагово:
  1. На каждом шаге RNN принимает эмбединговое представление текущего токена и скрытое состояние от предыдущего шага
  2. Выход каждого шага включает новое скрытое состояние, которое передаётся на следующий шаг вместе со следующим токеном
- В конце последовательности RNN генерирует контекстный вектор, который является финальным скрытым состоянием. Этот вектор обобщает всю информацию из входной последовательности и передаётся в декодер для дальнейшей генерации выходной последовательности. Контекстный вектор — это своего рода сжатая версия входной последовательности, включающая в себя её смысл
- Для улучшения качества представления входной последовательности часто используются двунаправленные RNN. В этом случае два RNN работают параллельно: один — слева направо, другой — справа налево. Их состояния объединяются на каждом шаге, что позволяет учитывать как предшествующие, так и последующие слова для каждого токена в последовательности

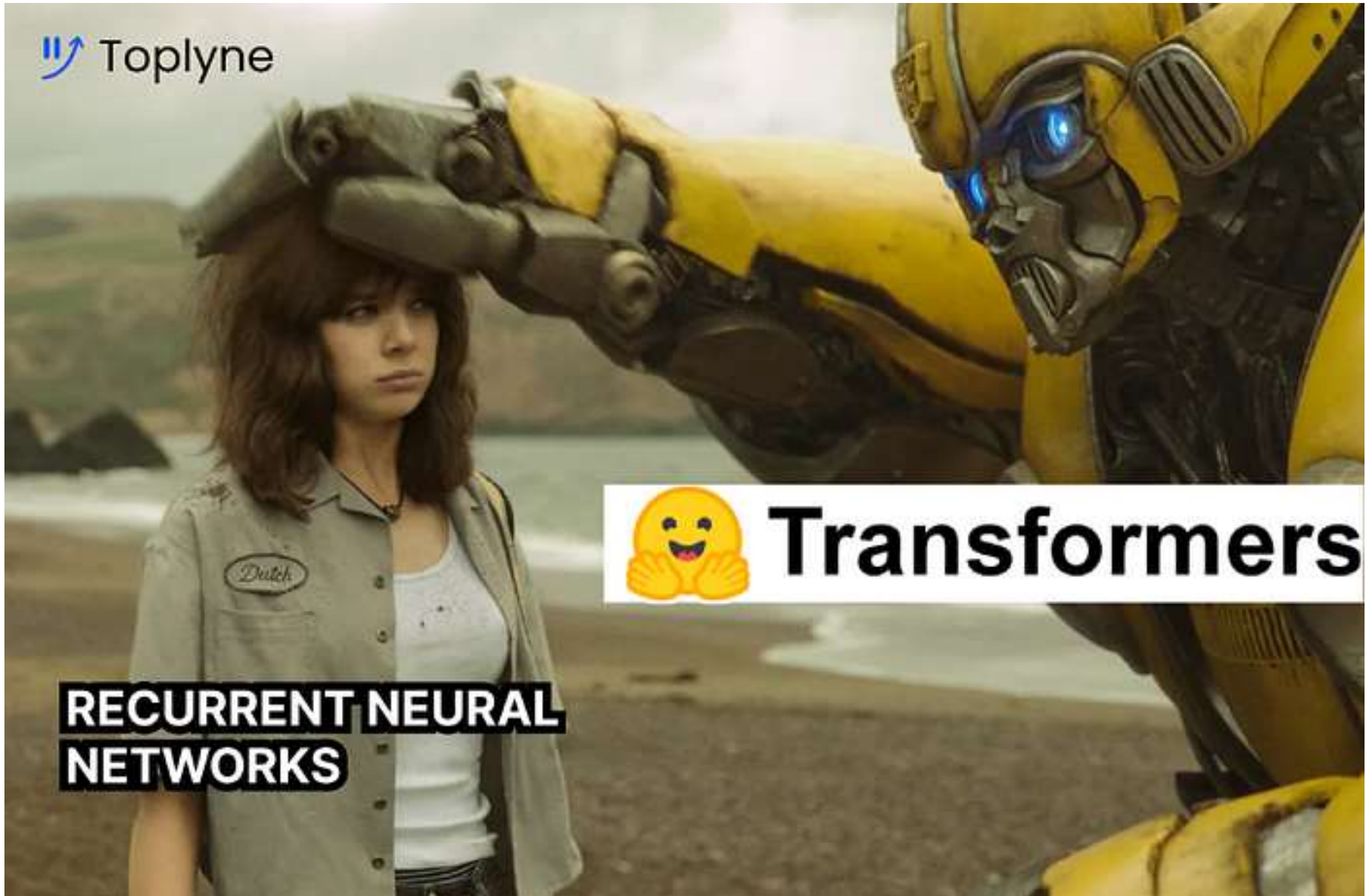
- Decoder – декодер генерирует данные на основе предыдущих предсказаний и контекстного вектора, предоставленного энкодером
- Подобно энкодеру, декодер принимает токены, которые сначала преобразуются в числовые представления с помощью эмбединга. На вход подаются не только реальные данные, но и предсказанные токены на предыдущих шагах. Реализован с использованием LSTM или GRU, на каждом шаге принимает:
  1. Контекстный вектор от энкодера
  2. Предыдущий предсказанный токен (начальный токен для первого шага)
  3. Скрытое состояние от предыдущего шага декодера. Этот выходной вектор затем преобразуется в вероятности через слой Softmax, который указывает на вероятность каждого возможного токена в выходной последовательности



# Механизм внимания (attention)

- Предоставление декодеру информации обо всех токенах исходного предложения на каждом шаге генерации для указания, на какое слово обратить внимание. На каждом шаге декодера считаются attention scores, получаем N значений, указывающих, насколько каждый из токенов с номерами (0 ... n) из исходного состояния важен для генерации токена i из выходной последовательности.





# Transformers

**RECURRENT NEURAL  
NETWORKS**

# Attention Is All You Need

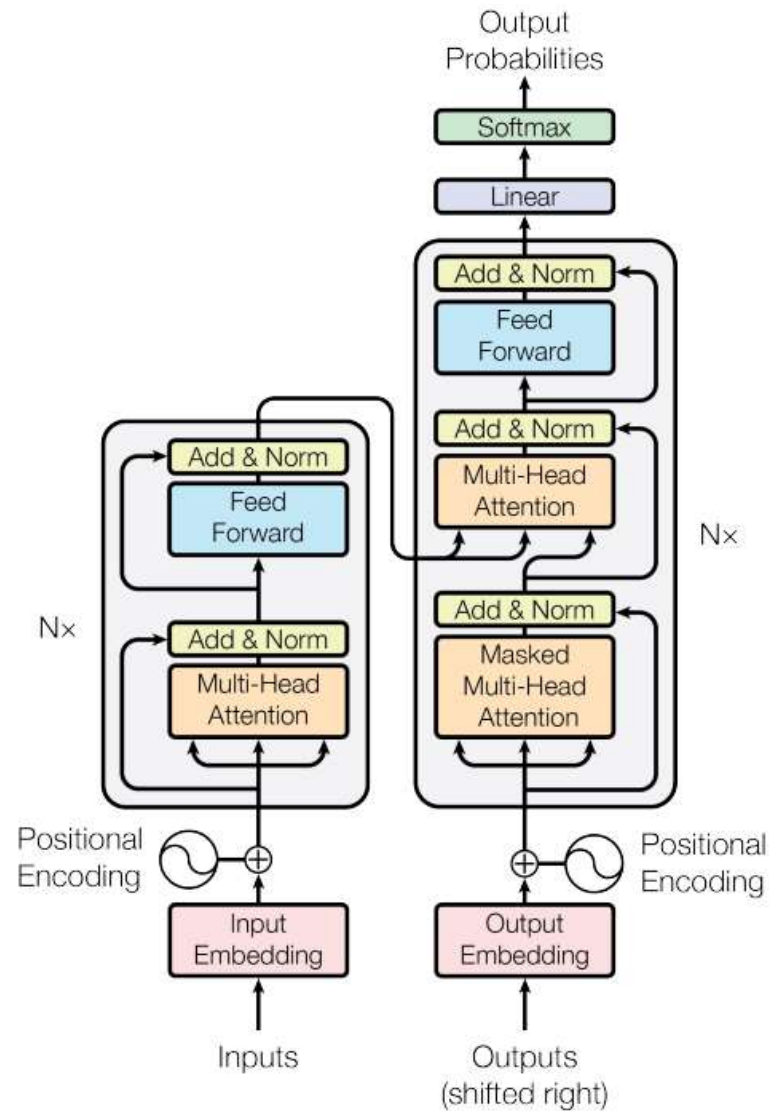


Figure 1: The Transformer - model architecture.

# Механизм внутреннего внимания (self-attention)

- Используется, чтобы посмотреть на другие слова во входной последовательности во время кодирования конкретного слова
- Из исходного вектора (эмбединга каждого токена) формируется три вектора – Query (запрос), Key (ключ), Value (значение). Они получаются с помощью умножения входного вектора на матрицы  $W_Q$ ,  $W_K$ ,  $W_V$ , веса которых учатся вместе со всеми остальными параметрами модели с помощью обратного распространения ошибки
- Выделение трех абстракций нужно, чтобы разграничить эмбединги, задающие направление внимания (query, key) и смысловую часть токена (value). Вектор query задает модальность «начальной точки» механизма внутреннего внимания (от какого токена направлено внимание), вектор key – модальность «конечной точки» (к какому токenu направлено внимание). Один и тот же токен может выступать как «начальной», так и «конечной» точкой направления внимания – self-attention вычисляется между всеми токенами в выбранном фрагменте текста

# Механизм внутреннего внимания (self-attention)

- Каждый токен фиксируется по очереди – он становится query и просчитывается его степень связанности со всеми оставшимися токенами. Для этого поочередно key-вектора всех токенов скалярно умножаются на query-вектор текущего токена. Полученные числа показывают, насколько важны остальные токены при кодировании query токена в конкретной позиции
- Полученные числа нормализуются и пропускаются через Softmax для получения распределения вероятностей. Затем подсчитывается взвешенная сумма value векторов, где в качестве весов используются полученные на предыдущем шаге вероятности. Полученный вектор и будет выходом внутреннего слоя внимания для одного токена
- На практике используются матричные вычисления вместо вычисления значений  $Q, K, V$  для отдельных токенов
- Обычно параллельно используется несколько self-attention блоков – «Multi-head self-attention»



# Механизм внутреннего внимания (self-attention)

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

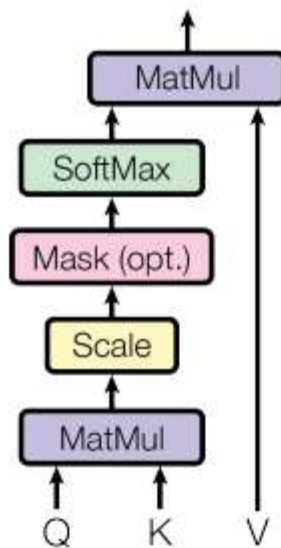
Q, K, V – матрицы запросов, ключей и значений (query, key, values)

$\sqrt{d_k}$  - нормировочная константа – корень из размерности ключей и значений

The diagram illustrates the self-attention calculation in matrix form. It shows the formula:  $\text{softmax}\left(\frac{Q \times K^T}{\sqrt{d_k}}\right) V = Z$ . The matrices are represented as grids: Q (purple, 2x3), K<sup>T</sup> (orange, 3x2), V (blue, 2x3), and Z (pink, 2x3). The multiplication of Q and K<sup>T</sup> is shown with a 'x' symbol between them, and the result is divided by  $\sqrt{d_k}$  before the softmax function is applied. The final result Z is shown as a pink 2x3 matrix.

The self-attention calculation in matrix form

### Scaled Dot-Product Attention



### Multi-Head Attention

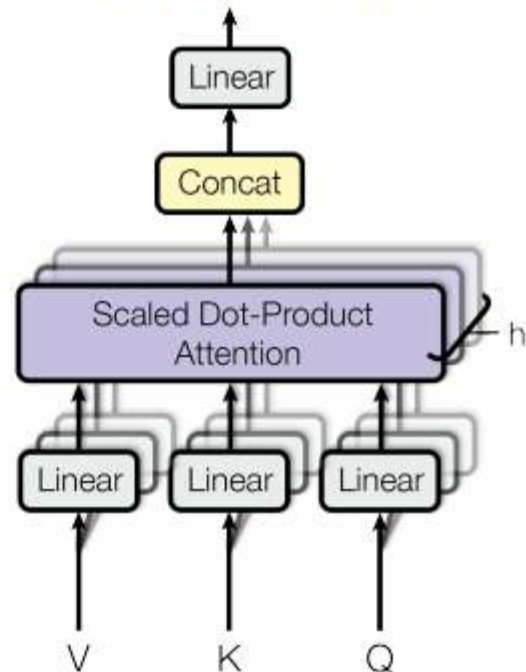
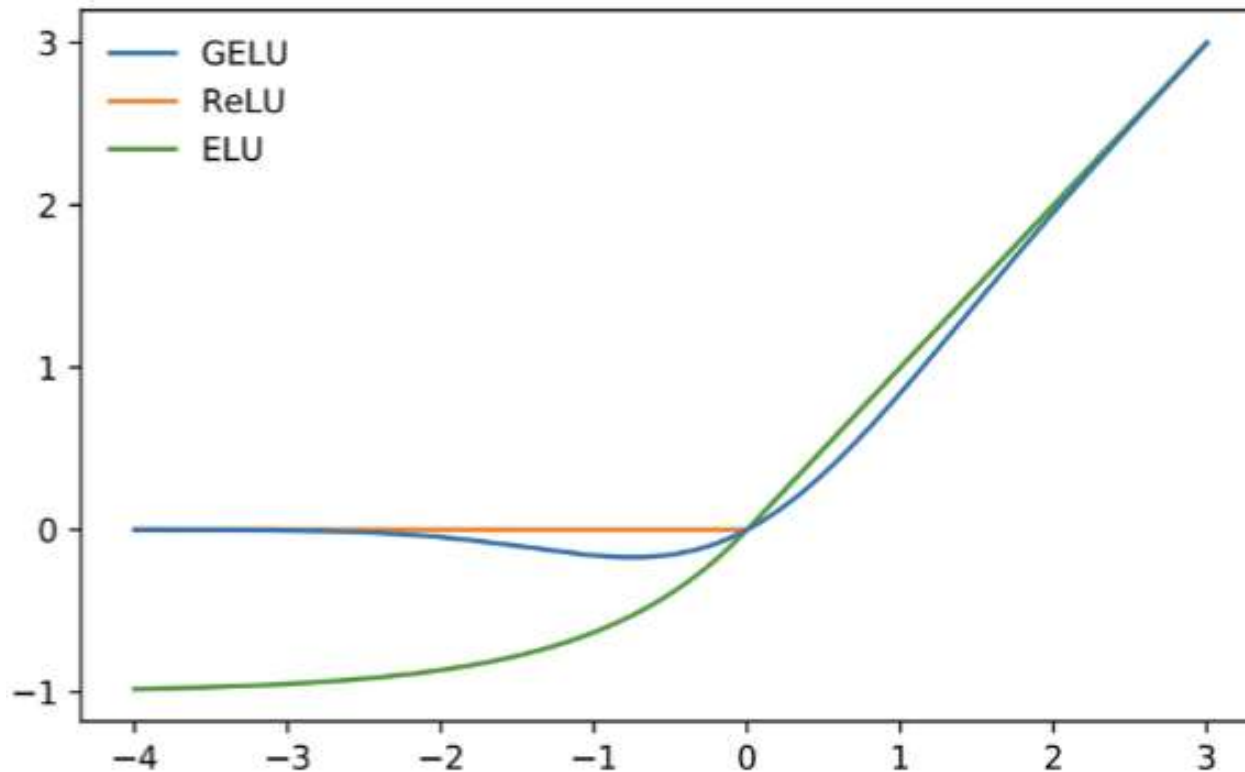


Figure 2: (left) Scaled Dot-Product Attention. (right) Multi-Head Attention consists of several attention layers running in parallel.

# Полносвязный слой

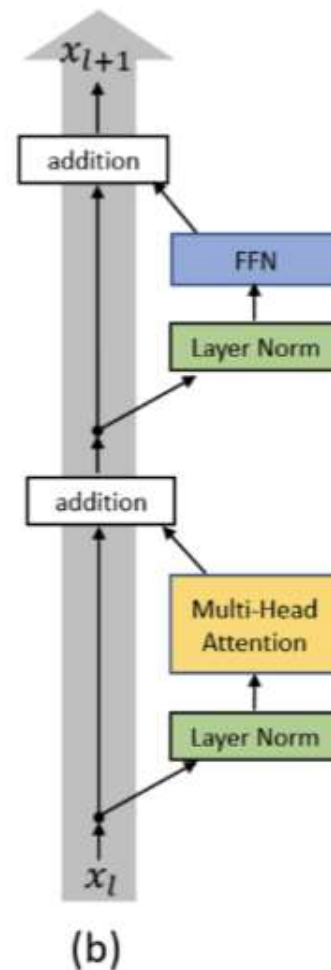
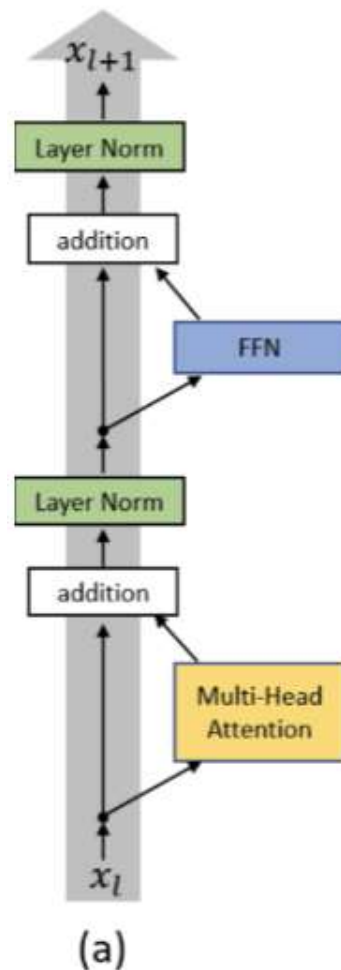
- Feed-forward network (FFN) – два полносвязных слоя, применяемых независимо к каждому элементу входной последовательности

$$FFN(x) = GELU(xW_1 + b_1)W_2 + b_2$$



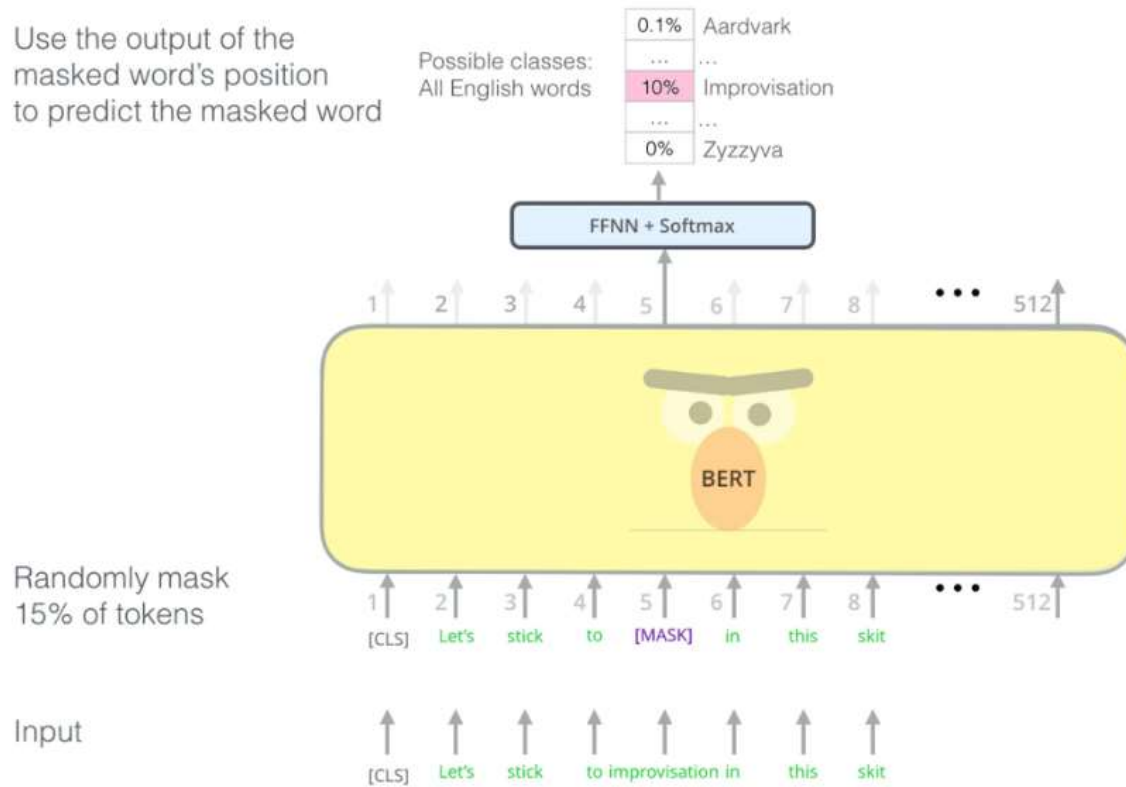
# Нормализация

- Layer normalization – нормализация применяется после остаточной связи – PostLN (a) либо ко входу residual-ветки – PreLN (b)



# BERT

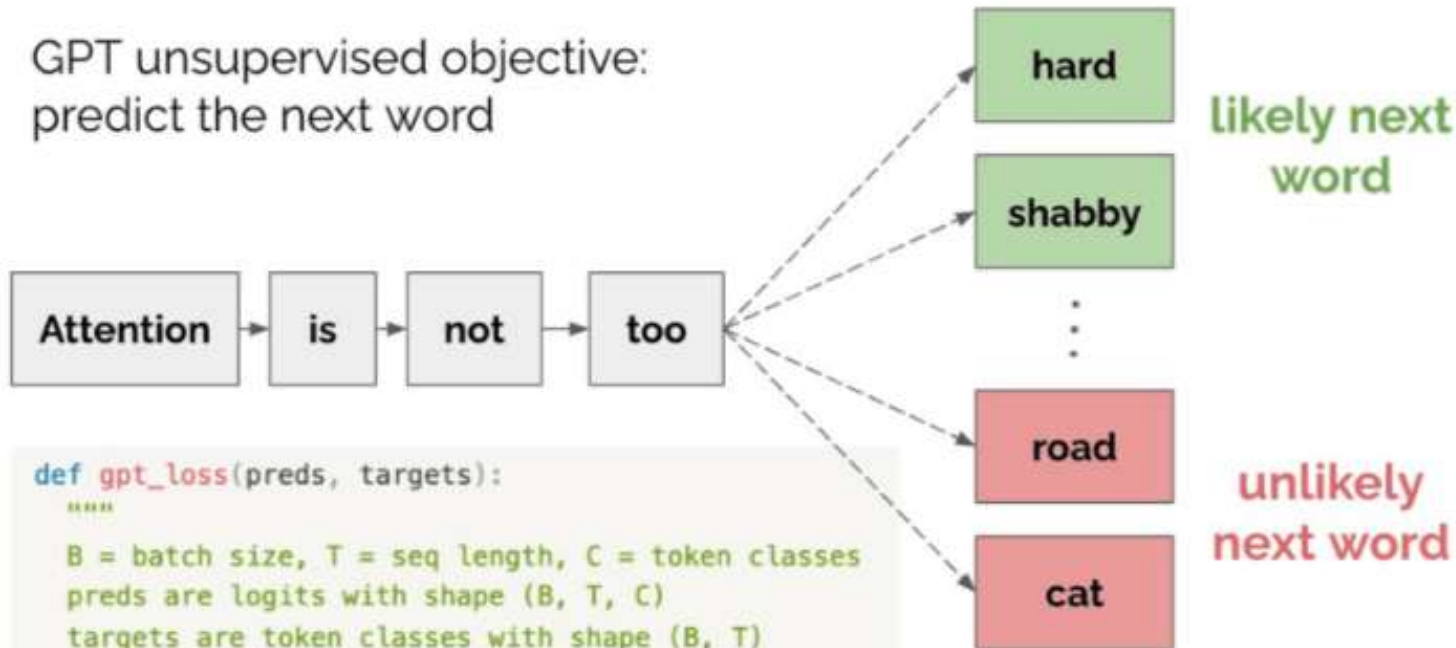
- Bidirectional Encoder Representations from Transformer – двунаправленный механизм внимания, при обработке входной последовательности все токены могут использовать информацию друг о друге
- BERT не учится генерировать тексты с нуля, первая его задача это предсказание случайно замаскированных слов по оставшимся (masked language modeling), вторая – предсказание по паре текстовых фрагментов, следуют они друг за другом или нет (next sentence prediction)



# GPT

- Generative Pretrained Transformer – языковая модель, реализованная в виде последовательности слоев декодера трансформера
- В качестве задачи при обучении выступает предсказание следующего токена (многоклассовая классификация по словарю)

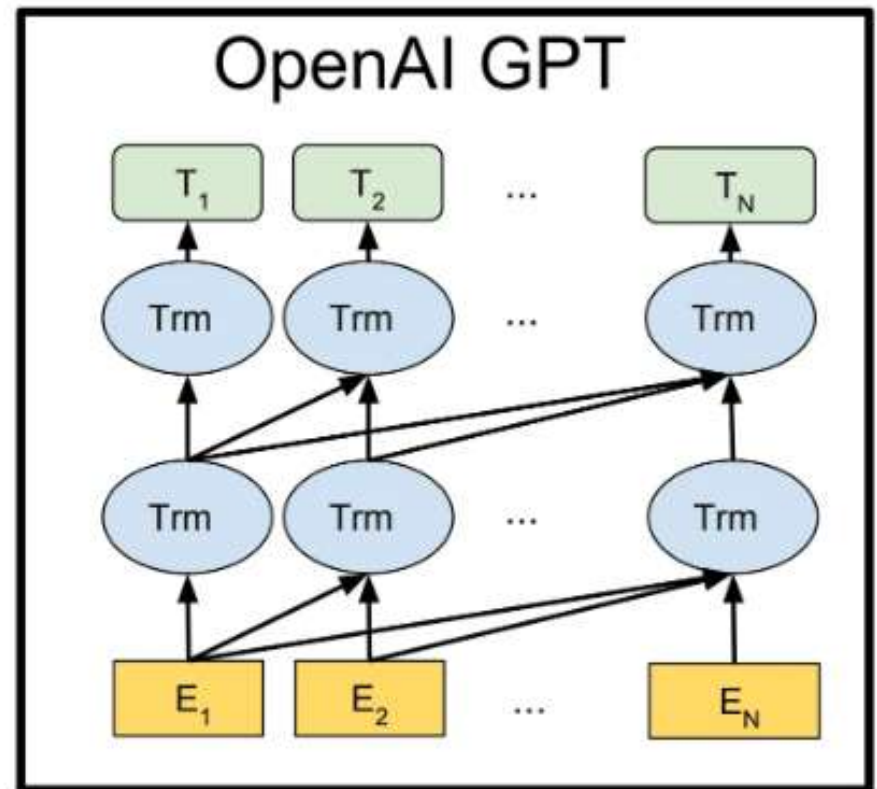
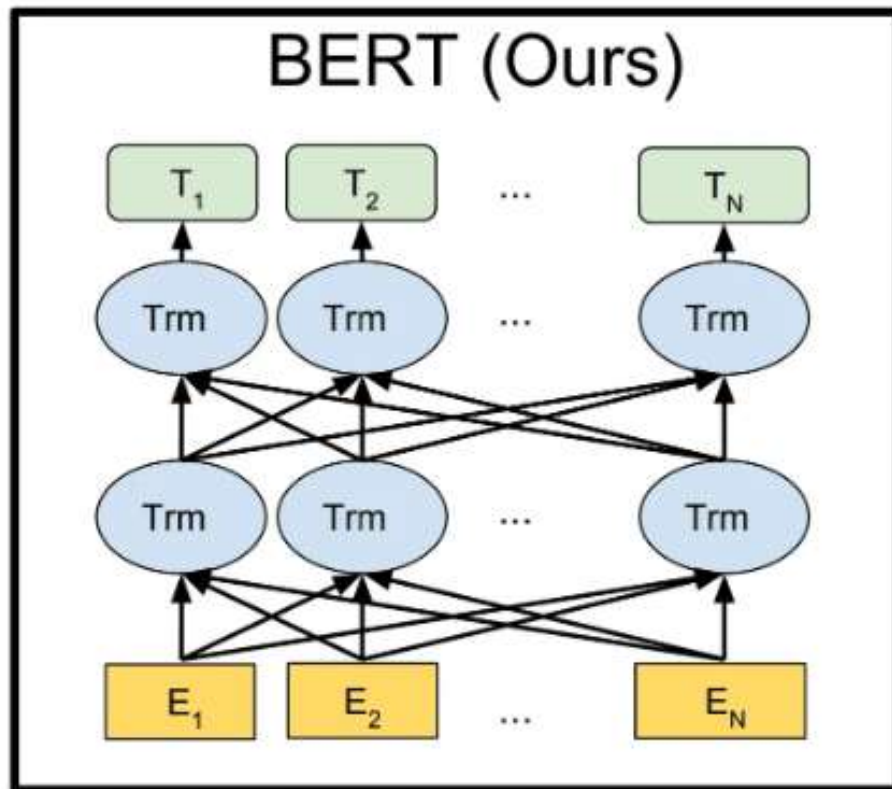
GPT unsupervised objective:  
predict the next word



```
def gpt_loss(preds, targets):  
    """  
    B = batch size, T = seq length, C = token classes  
    preds are logits with shape (B, T, C)  
    targets are token classes with shape (B, T)  
    """  
    preds = preds.view(B*T,C)  
    targets = targets.view(B*T)  
    return F.cross_entropy(preds, targets)
```

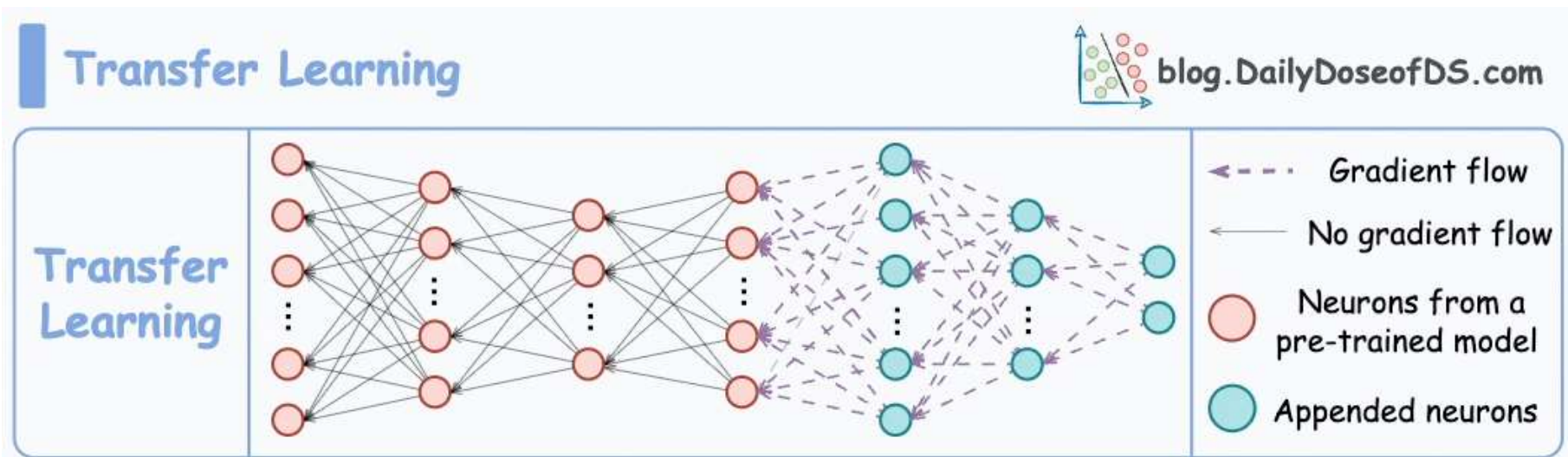
# BERT vs GPT

- Ключевое отличие – использование разных видов внимания



# Transfer learning

- Перенос обучения (transfer learning) – метод обучения, при котором модель, обученная для одной задачи (связанная задача), переиспользуется для схожей задачи (целевая задача)
- Целевая задача содержит мало данных. Например, классификация специфических изображений. Связанная содержит много данных, есть обученная модель которая хорошо решает связанную задачу.
- Заменяем от 1 до N последних слоев, обучаем только эти слои (веса в остальных слоях не изменяются)





# Fine-tuning

- Fine tuning – метод дополнительного обучения сети, при котором обученная на схожей задаче модель обучается на новом наборе данных. При этом веса не инициализируются случайно, а используются уже оптимизированные в процессе предыдущего обучения веса нейронной сети
- Можно добавить несколько слоев на выходе нейронной сети
- Обучаются все слои

