

Основы глубинного обучения

Ульянкин Филипп

30 ноября 2023 г.

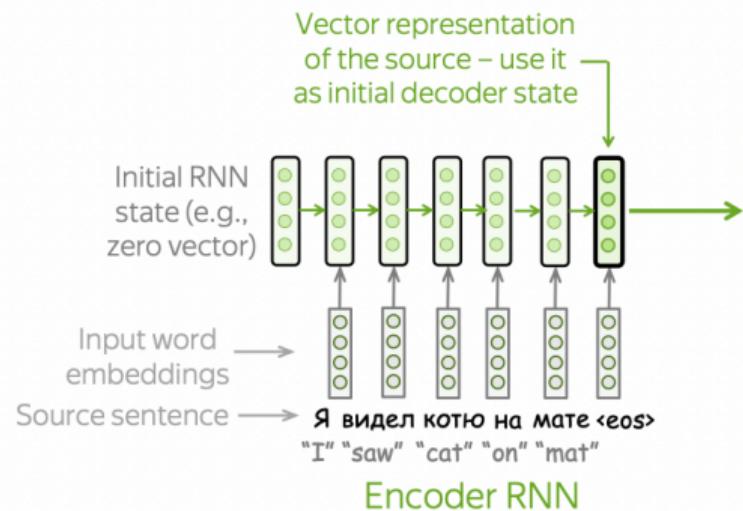
Лекция 10: seq2seq модели и механизм внимания

Agenda

- Seq2seq модели
- Механизмы внимания
- Attention is all you need
- Трансформеры
- BERT и GPT

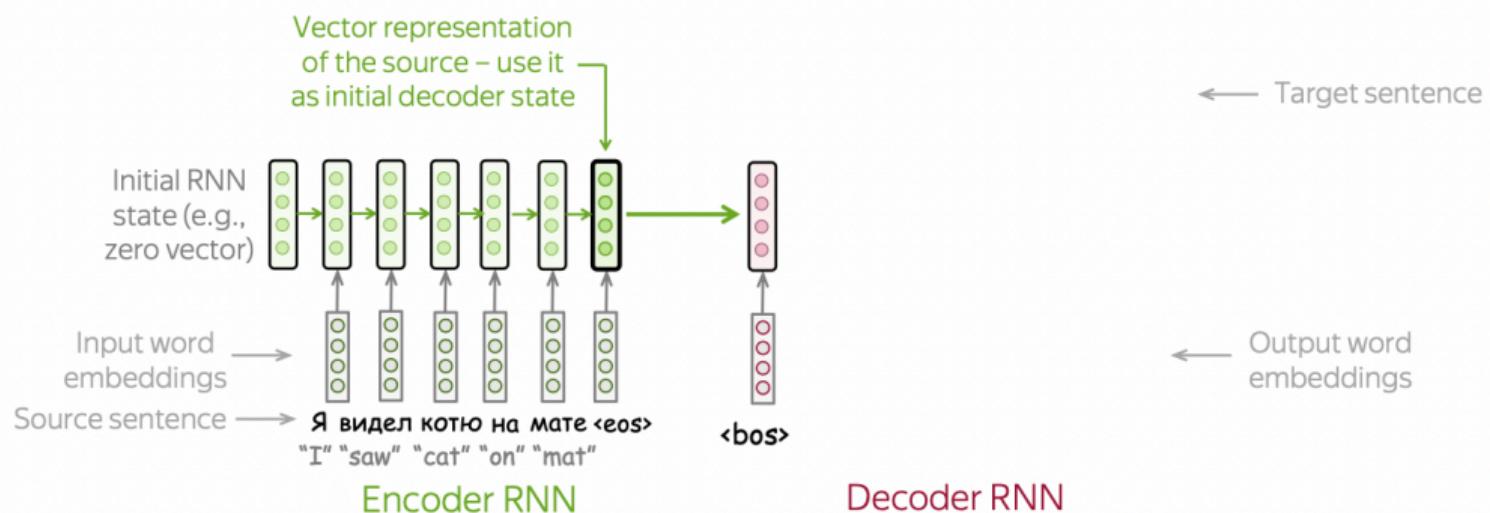
Seq2Seq перевод

Простейшая модель: RNN энкодер-декодер



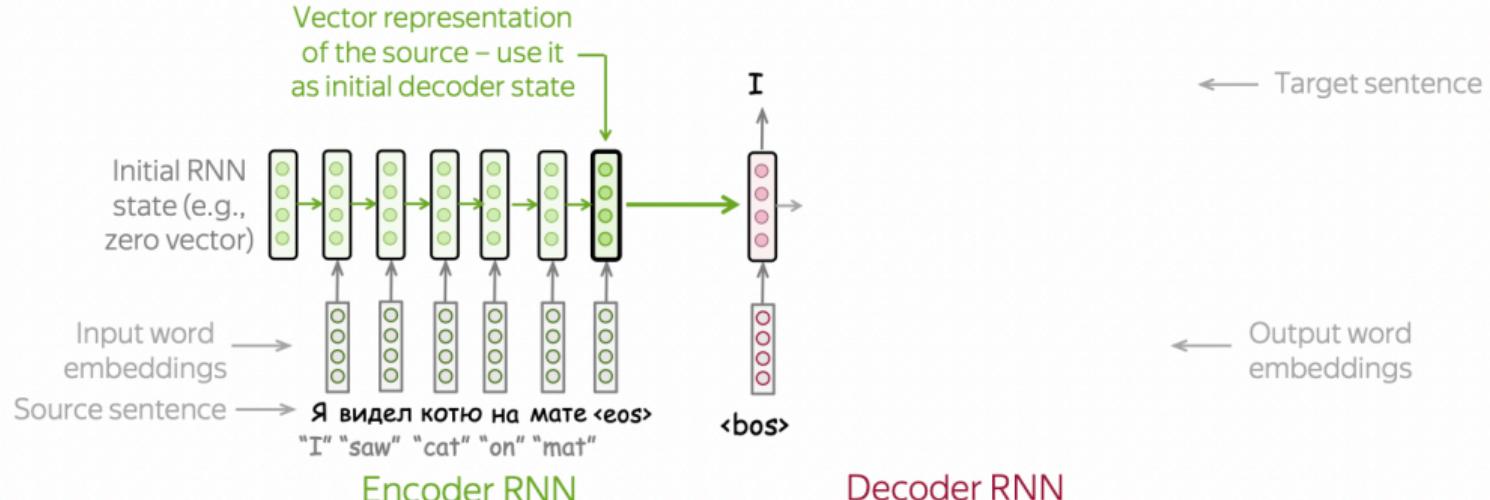
https://github.com/yandexdataschool/nlp_course/tree/2021/week04_seq2seq

Простейшая модель: RNN энкодер-декодер



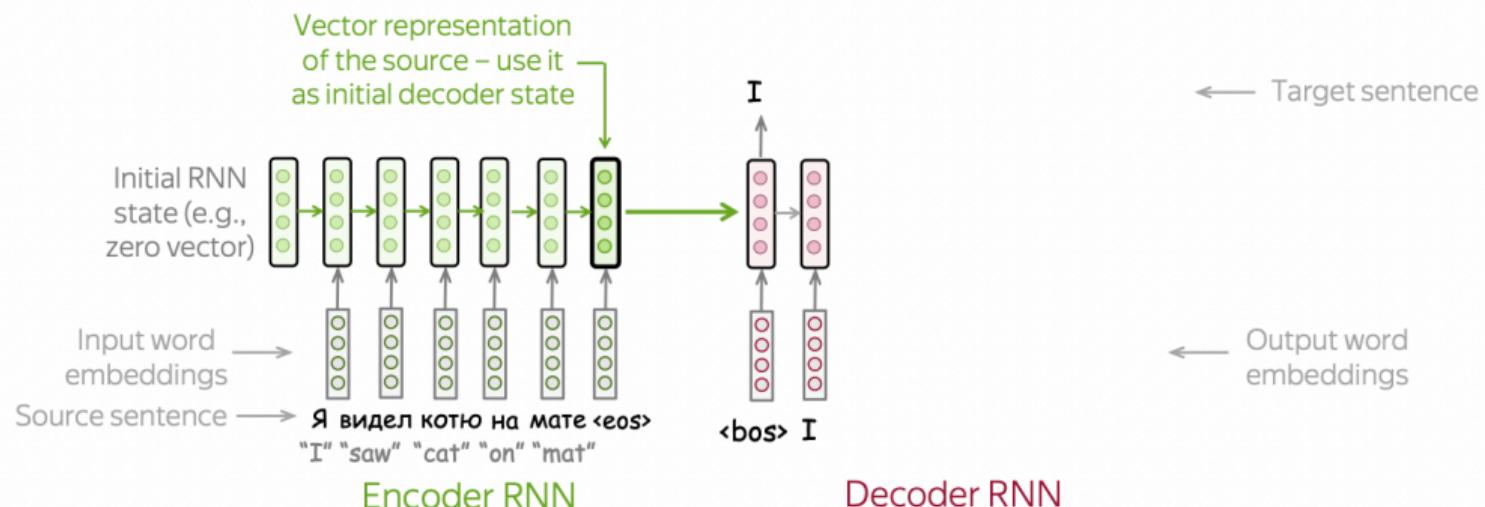
https://github.com/yandexdataschool/nlp_course/tree/2021/week04_seq2seq

Простейшая модель: RNN энкодер-декодер



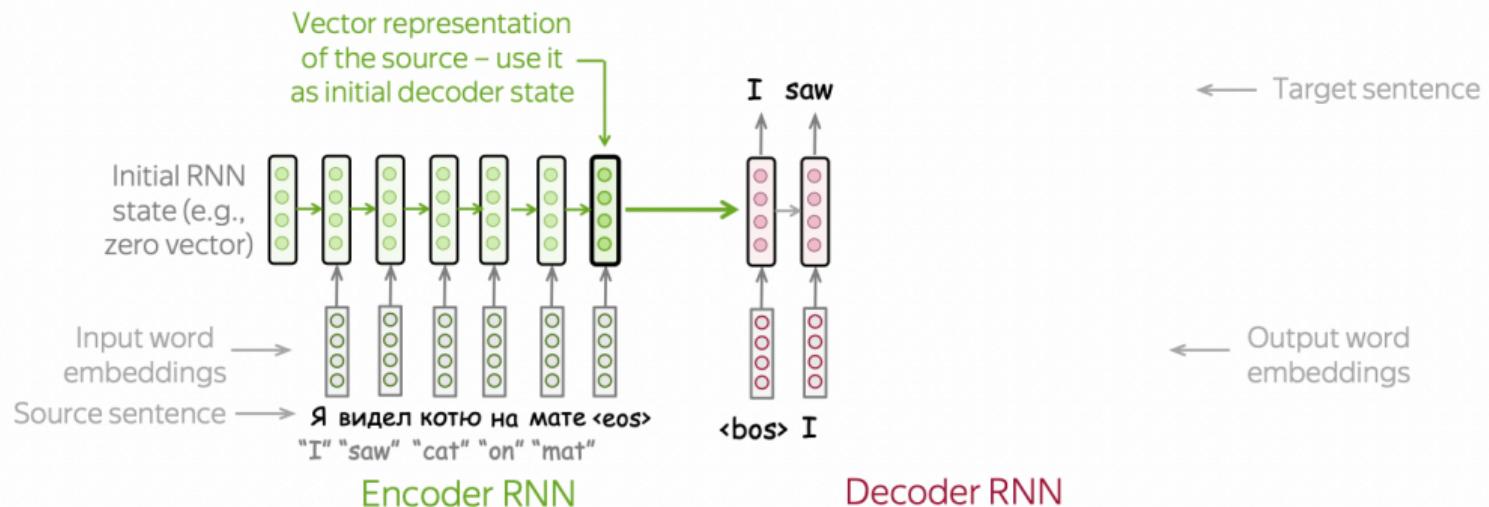
https://github.com/yandexdataschool/nlp_course/tree/2021/week04_seq2seq

Простейшая модель: RNN энкодер-декодер



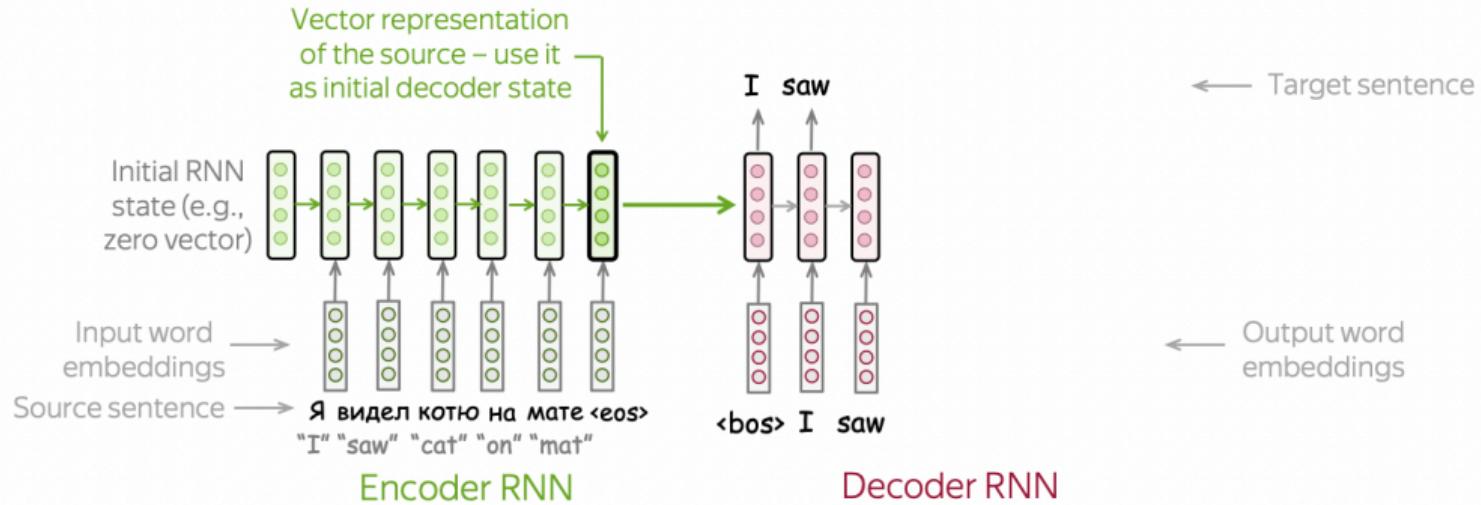
https://github.com/yandexdataschool/nlp_course/tree/2021/week04_seq2seq

Простейшая модель: RNN энкодер-декодер



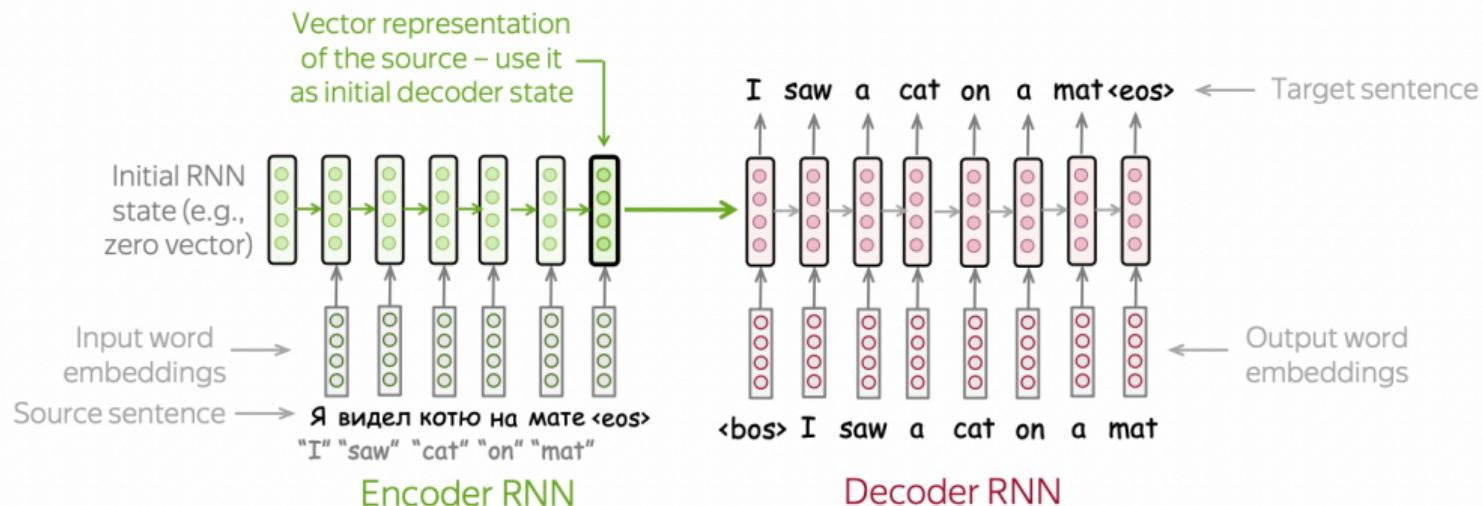
https://github.com/yandexdataschool/nlp_course/tree/2021/week04_seq2seq

Простейшая модель: RNN энкодер-декодер



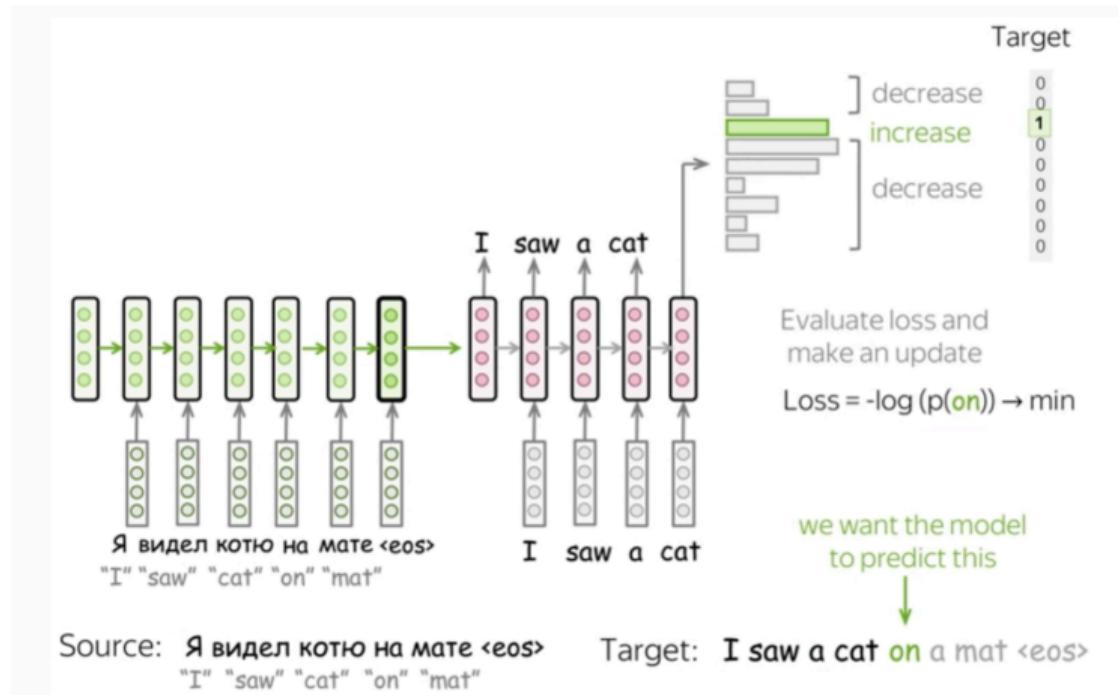
https://github.com/yandexdataschool/nlp_course/tree/2021/week04_seq2seq

Простейшая модель: RNN энкодер-декодер



https://github.com/yandexdataschool/nlp_course/tree/2021/week04_seq2seq

Как обучить модель?



https://lena-voita.github.io/resources/lectures/seq2seq/general/seq2seq_training_with_target.mp4

Как сделать прогноз?

- Жадно (greedy decoding): на каждом шаге берём токен с самой большой вероятностью

$$\prod_{t=1}^T \arg \max_{y_t} p(y_t | y_{<t}, x) \neq \arg \max_y \prod_{t=1}^T p(y_t | y_{<t}, x)$$

- Наш прогноз это последовательность, жадный способ на каждом шаге выбирает локальный оптимум
- Перебирать все траектории, чтобы найти глобальный оптимум очень дорого

Неоптимальность Greedy Decoding

Time step	1	2	3	4
A	0.5	0.1	0.2	0.0
B	0.2	0.4	0.2	0.2
C	0.2	0.3	0.4	0.2
<eos>	0.1	0.2	0.2	0.6

$$P = 0.5 \times 0.4 \times 0.4 \times 0.6 = 0.048$$

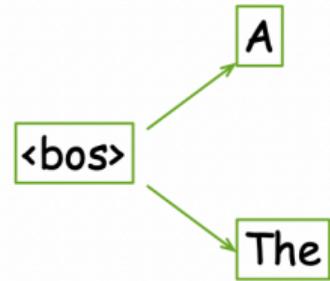
Time step	1	2	3	4
A	0.5	0.1	0.1	0.1
B	0.2	0.4	0.6	0.2
C	0.2	0.3	0.2	0.1
<eos>	0.1	0.2	0.1	0.6

$$P = 0.5 \times 0.3 \times 0.6 \times 0.6 = 0.054$$

Beam Search

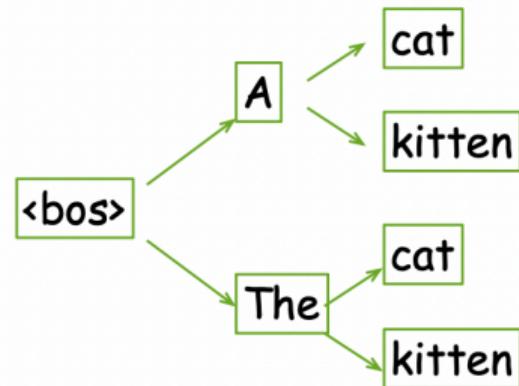
- Давайте поддерживать несколько самых вероятных траекторий
- Такая стратегия называется **Beam Search**
- Число траекторий, которое мы помним будет гиперпараметром (больше 10 брать не стоит)

Beam Search



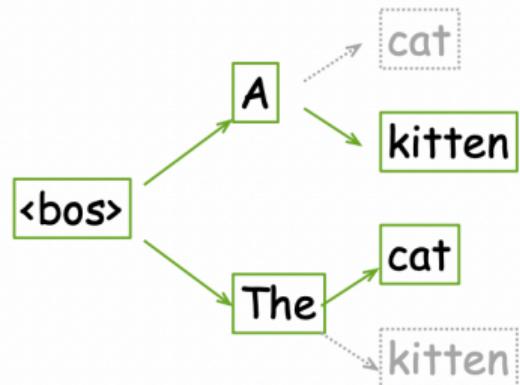
https://github.com/yandexdataschool/nlp_course/tree/2021/week04_seq2seq

Beam Search



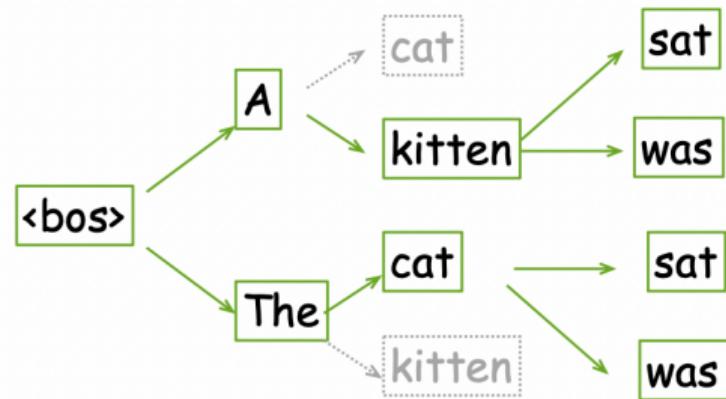
https://github.com/yandexdataschool/nlp_course/tree/2021/week04_seq2seq

Beam Search



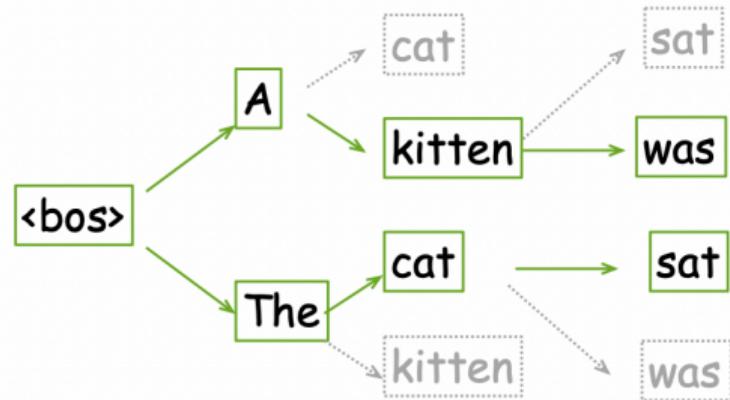
https://github.com/yandexdataschool/nlp_course/tree/2021/week04_seq2seq

Beam Search



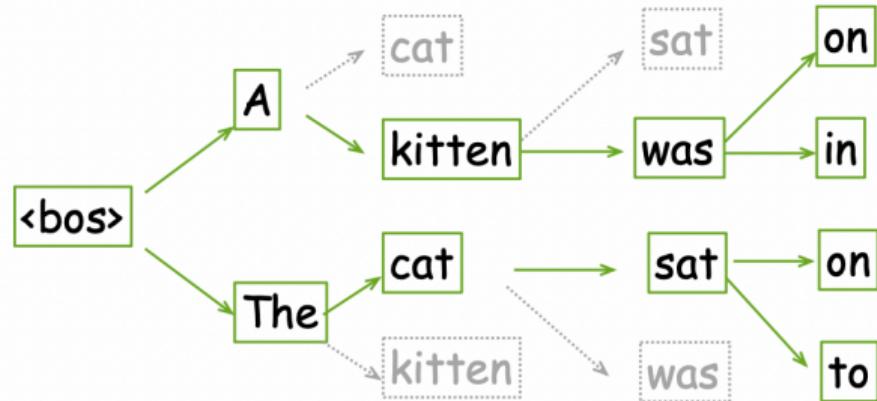
https://github.com/yandexdataschool/nlp_course/tree/2021/week04_seq2seq

Beam Search



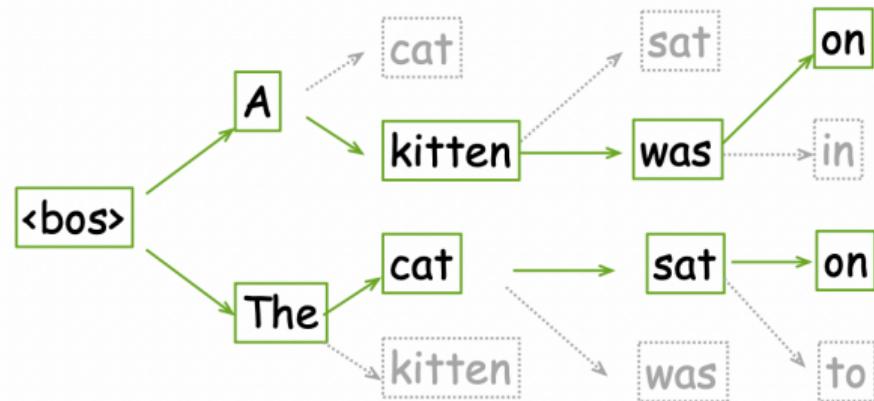
https://github.com/yandexdataschool/nlp_course/tree/2021/week04_seq2seq

Beam Search



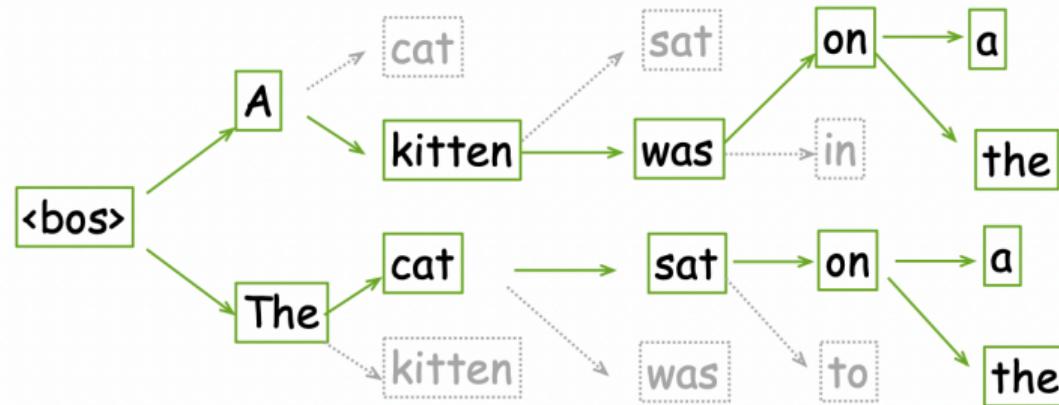
https://github.com/yandexdataschool/nlp_course/tree/2021/week04_seq2seq

Beam Search



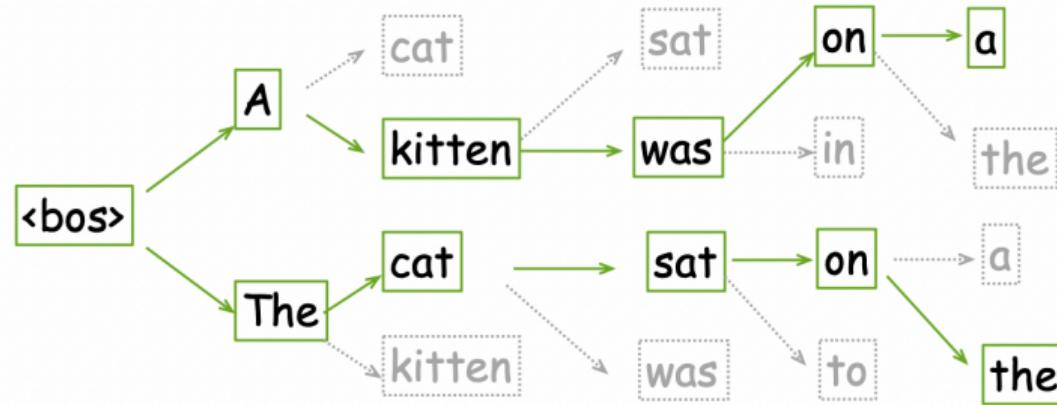
https://github.com/yandexdataschool/nlp_course/tree/2021/week04_seq2seq

Beam Search



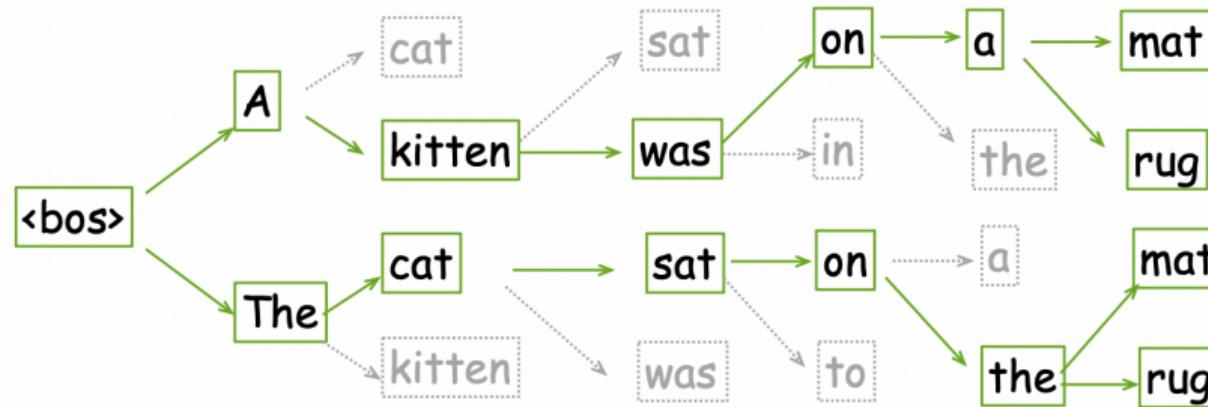
https://github.com/yandexdataschool/nlp_course/tree/2021/week04_seq2seq

Beam Search



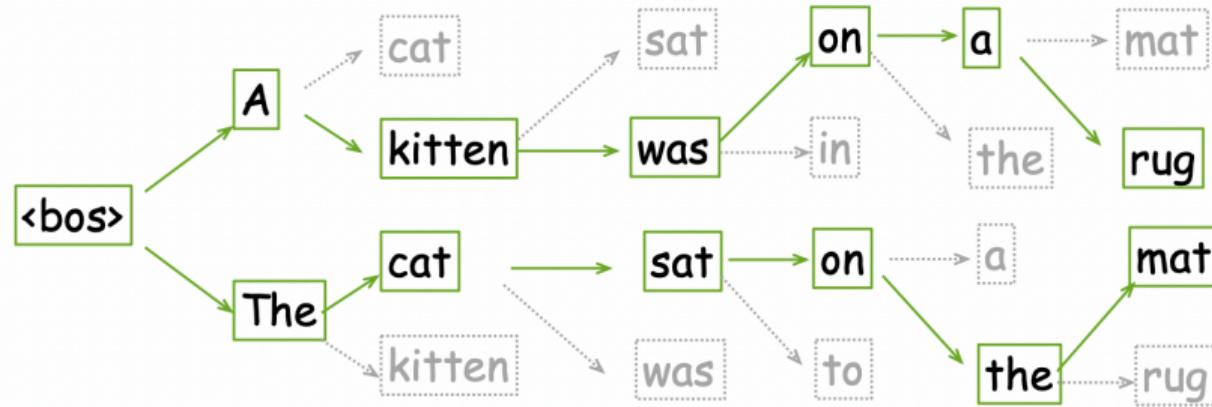
https://github.com/yandexdataschool/nlp_course/tree/2021/week04_seq2seq

Beam Search



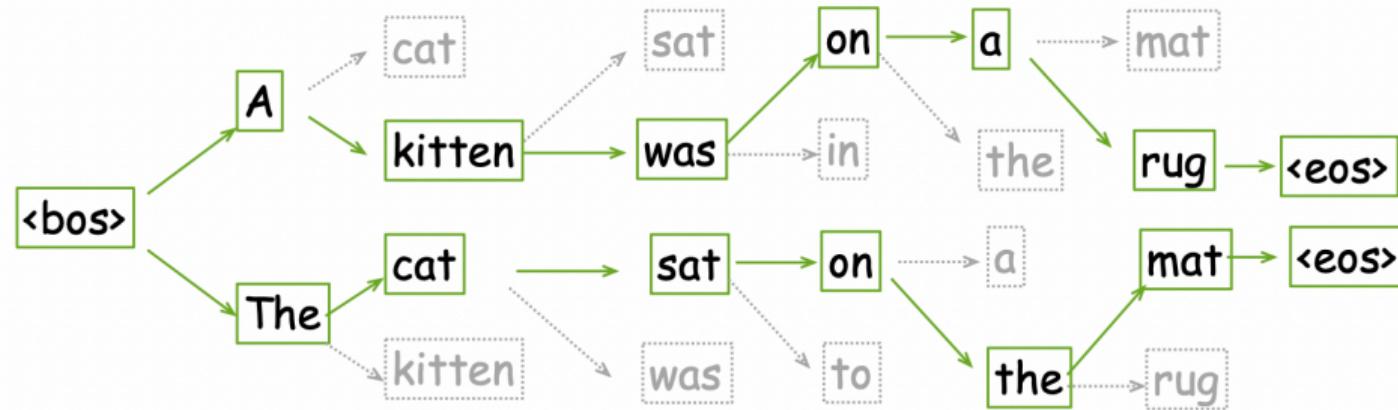
https://github.com/yandexdataschool/nlp_course/tree/2021/week04_seq2seq

Beam Search



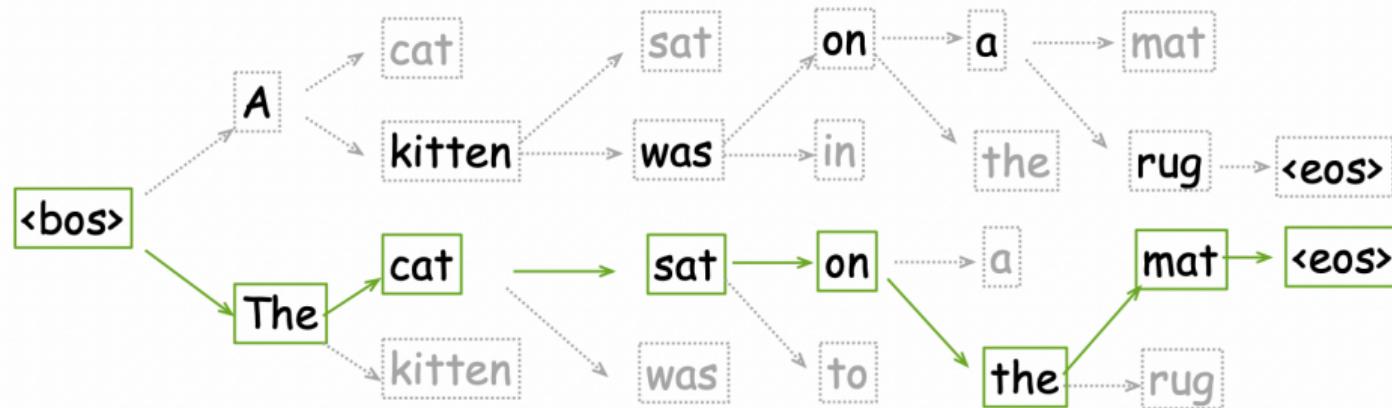
https://github.com/yandexdataschool/nlp_course/tree/2021/week04_seq2seq

Beam Search



https://github.com/yandexdataschool/nlp_course/tree/2021/week04_seq2seq

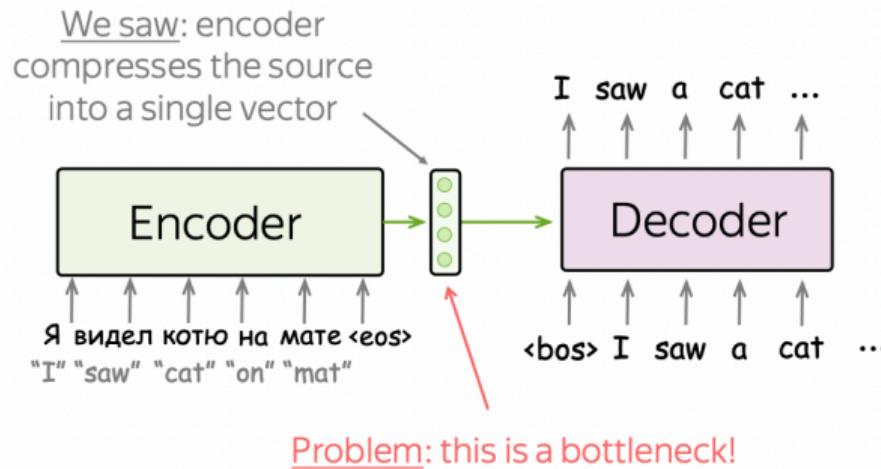
Beam Search



https://github.com/yandexdataschool/nlp_course/tree/2021/week04_seq2seq

Механизмы внимания

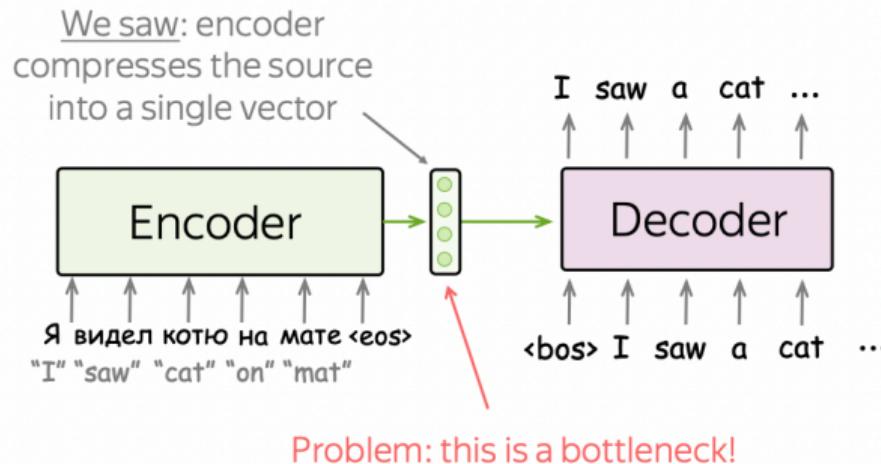
В чём проблема простой RNN модели?



- Сжать предложение в вектор сложно
- Декодеру может понадобиться информация о разных частях предложения для его расшифровки

https://github.com/yandexdataschool/nlp_course/tree/2021/week04_seq2seq

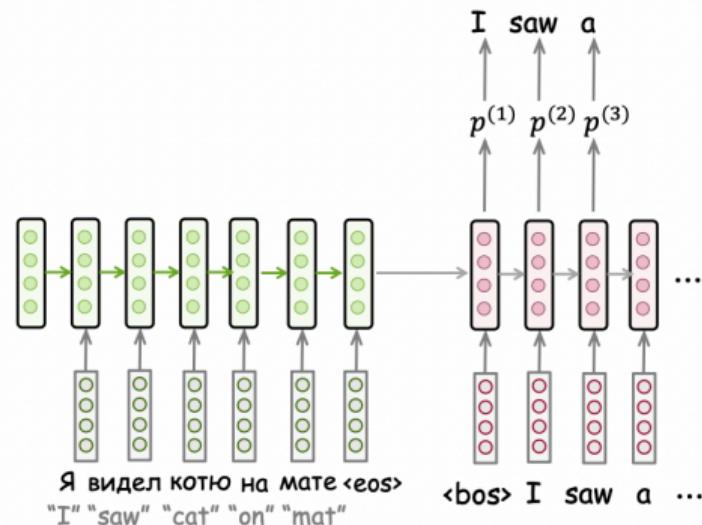
В чём проблема простой RNN модели?



- **Attention:** на разных шагах надо позволить модели фокусироваться на разных входных токенах

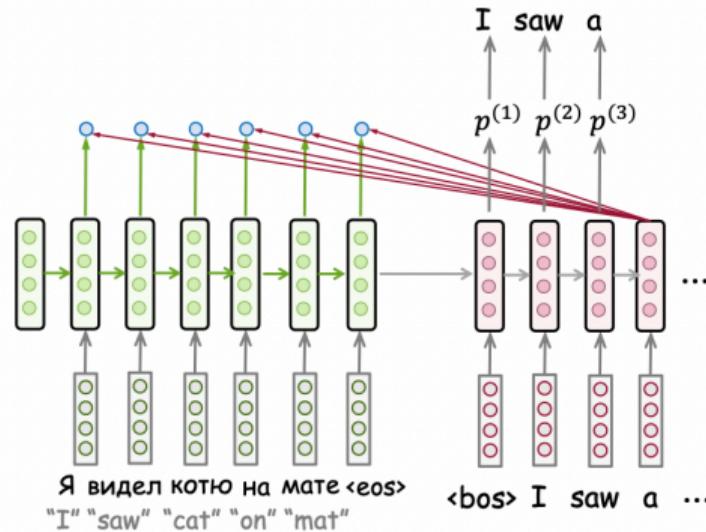
https://github.com/yandexdataschool/nlp_course/tree/2021/week04_seq2seq

Attention



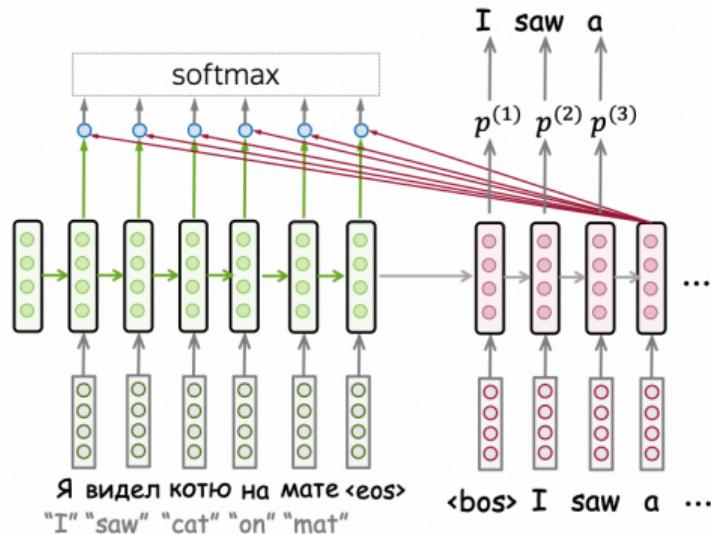
https://github.com/yandexdataschool/nlp_course/tree/2021/week04_seq2seq

Attention



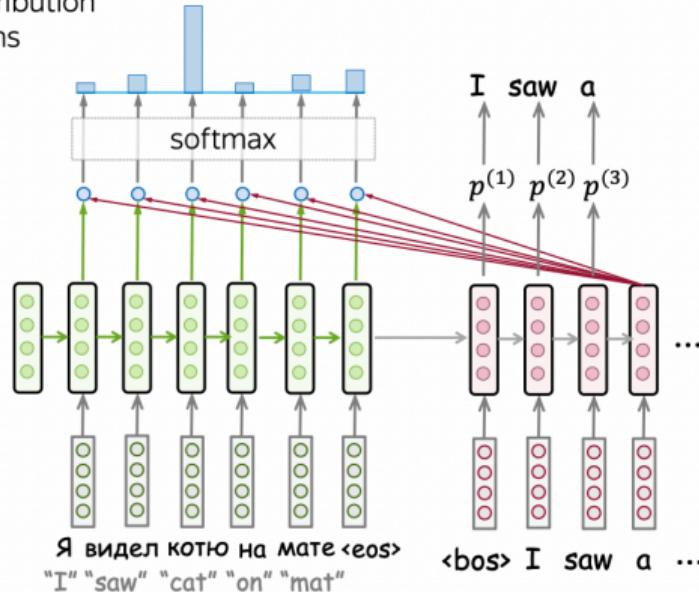
https://github.com/yandexdataschool/nlp_course/tree/2021/week04_seq2seq

Attention



Attention

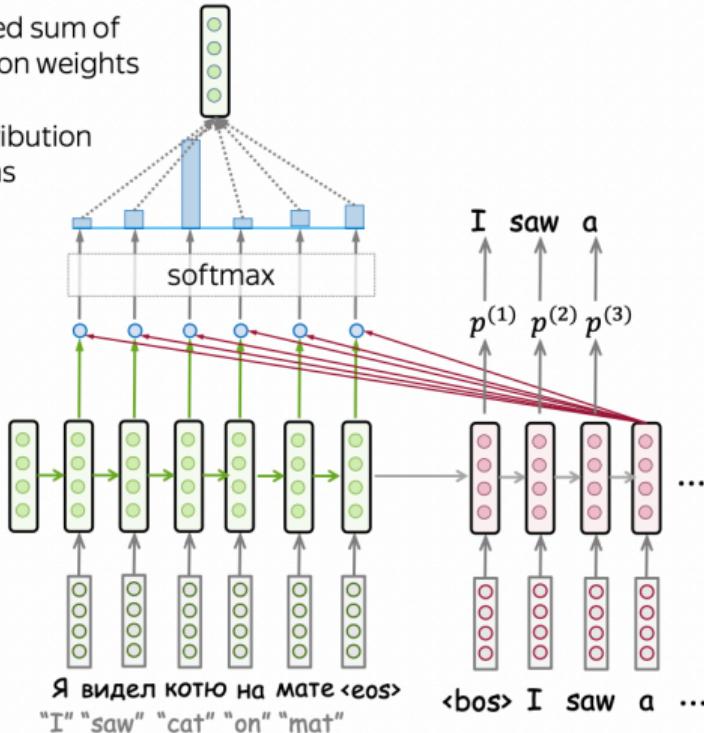
Attention weights: distribution over source tokens



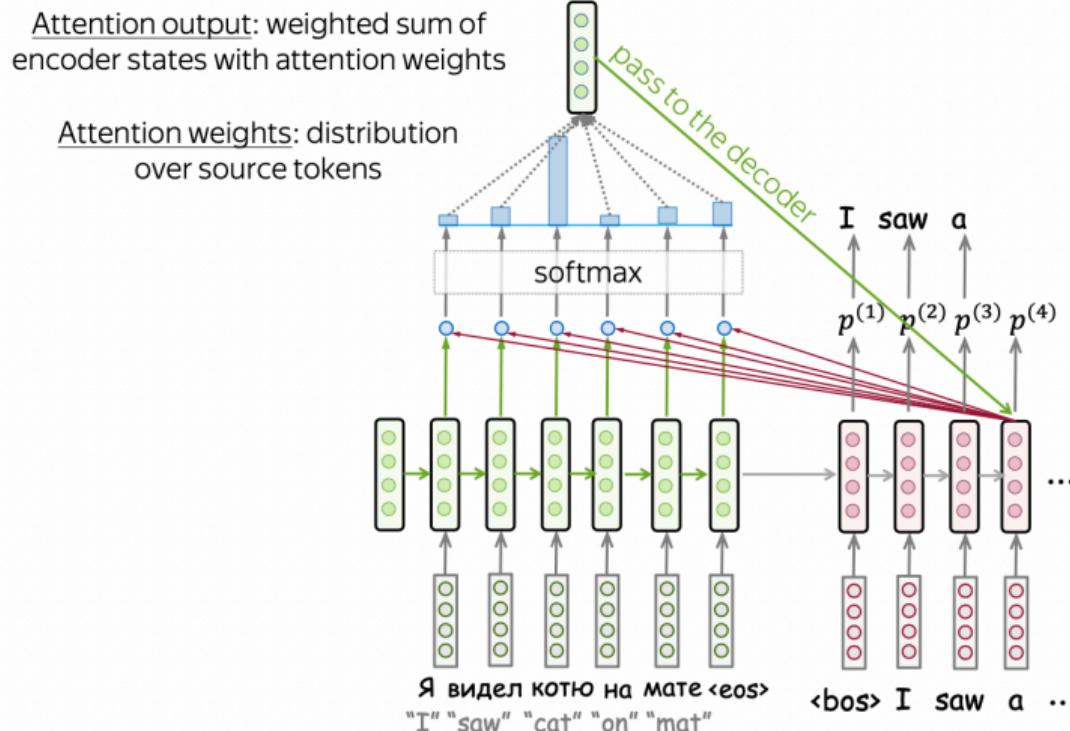
Attention

Attention output: weighted sum of encoder states with attention weights

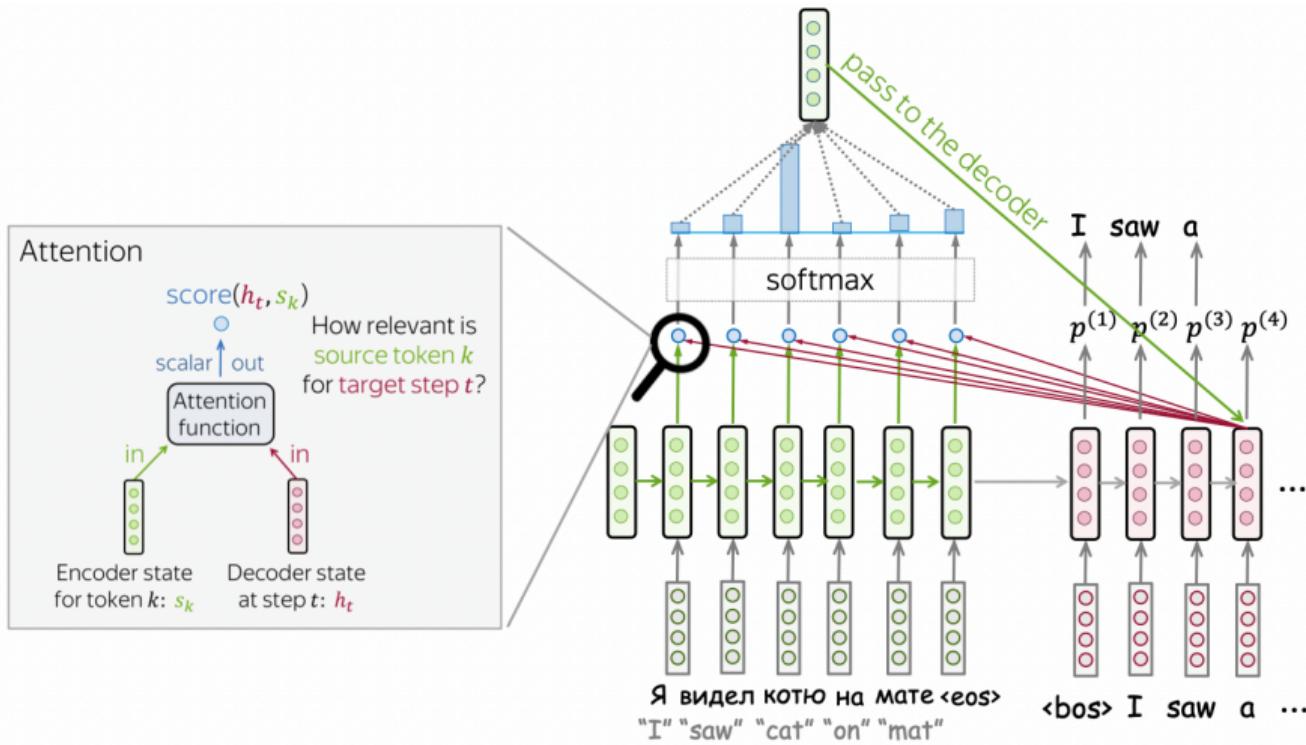
Attention weights: distribution over source tokens



Attention



Attention



Computation Pipeline

Attention input

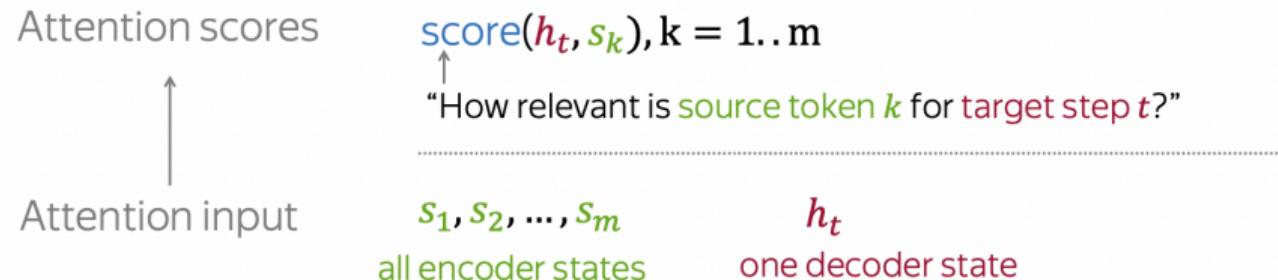
s_1, s_2, \dots, s_m

all encoder states

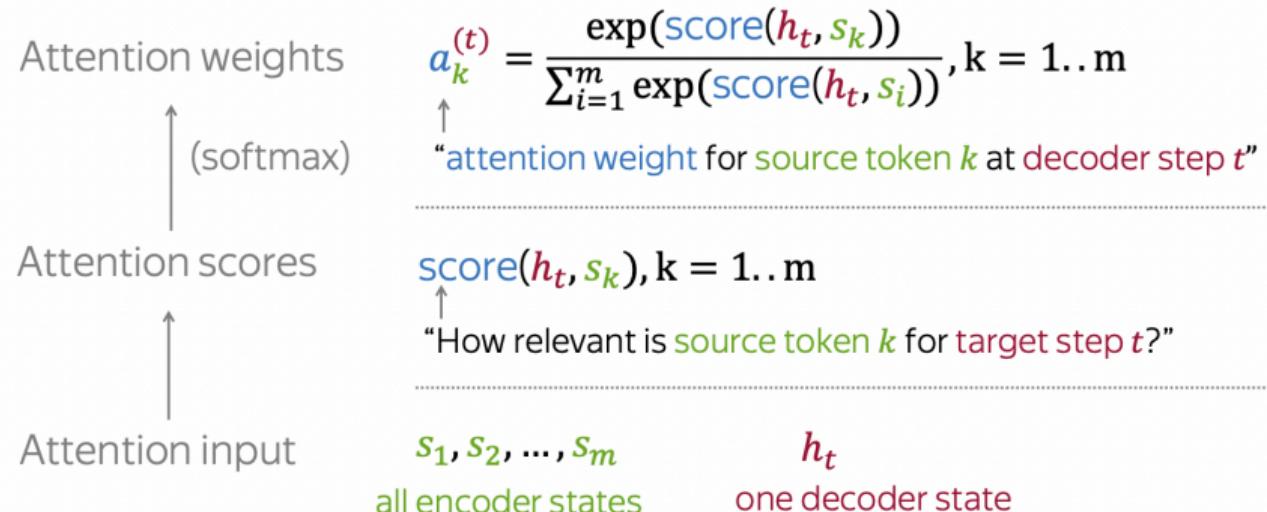
h_t

one decoder state

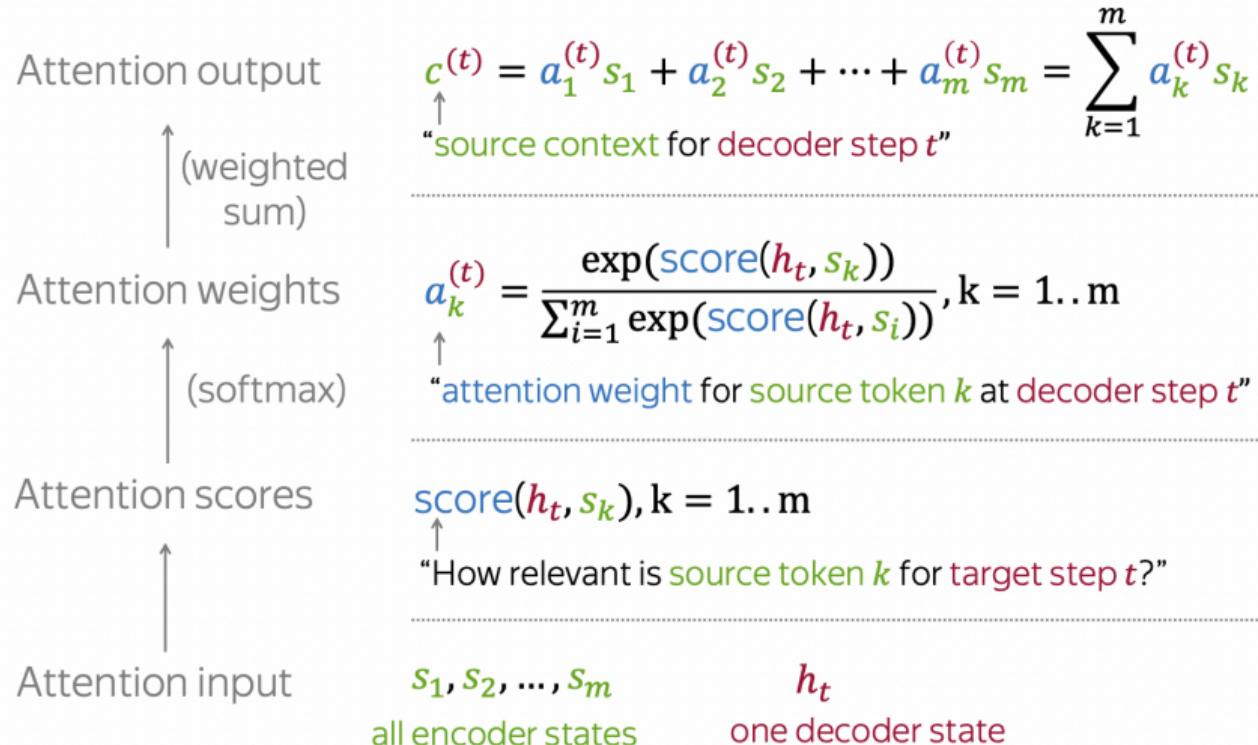
Computation Pipeline



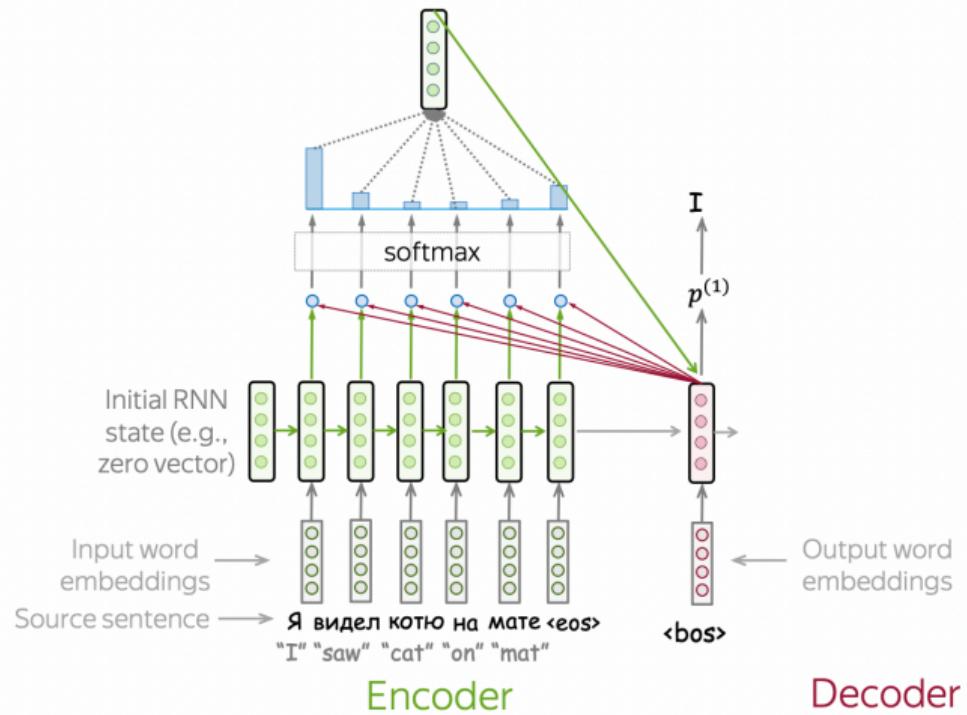
Computation Pipeline



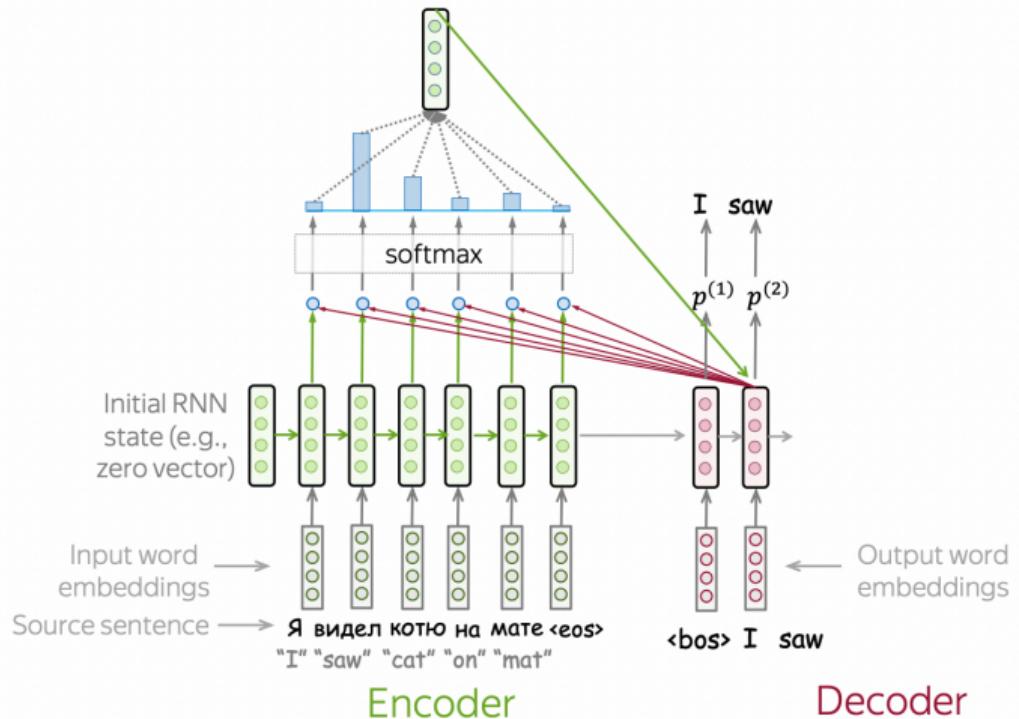
Computation Pipeline



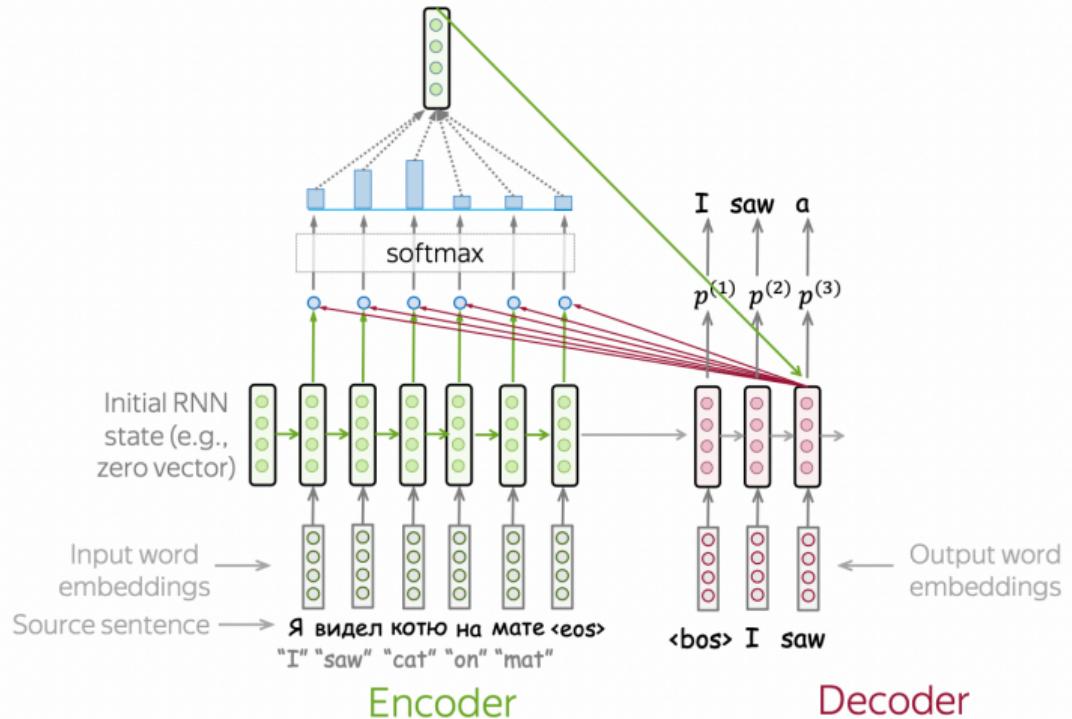
Внимание в декодере



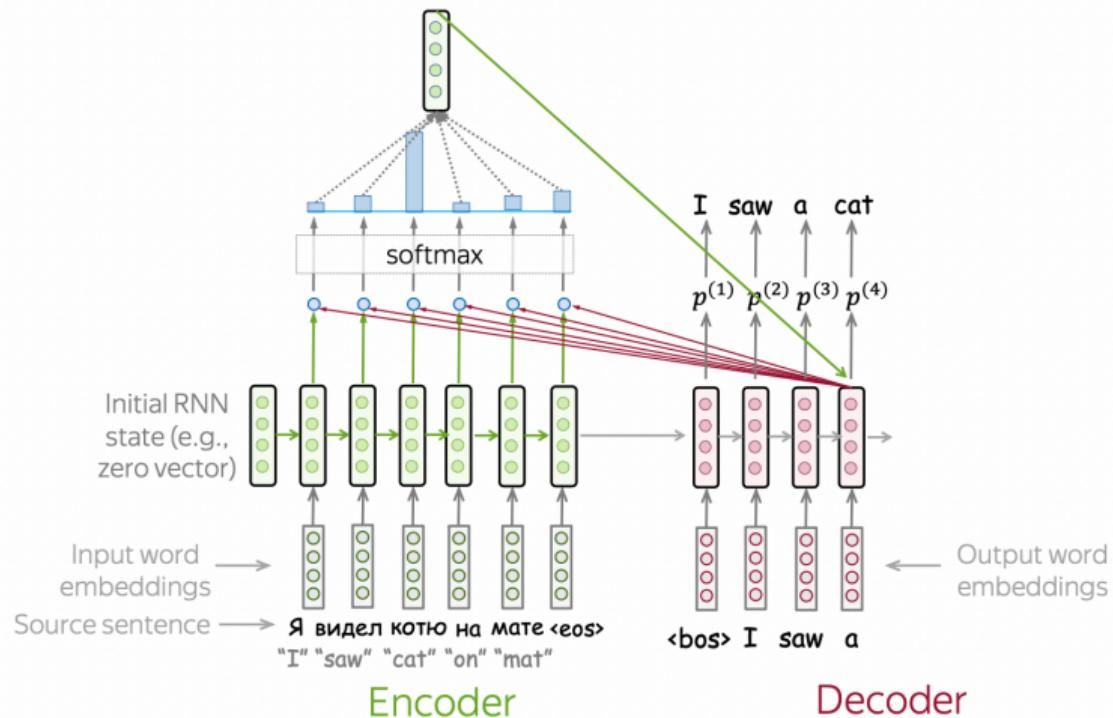
Внимание в декодере



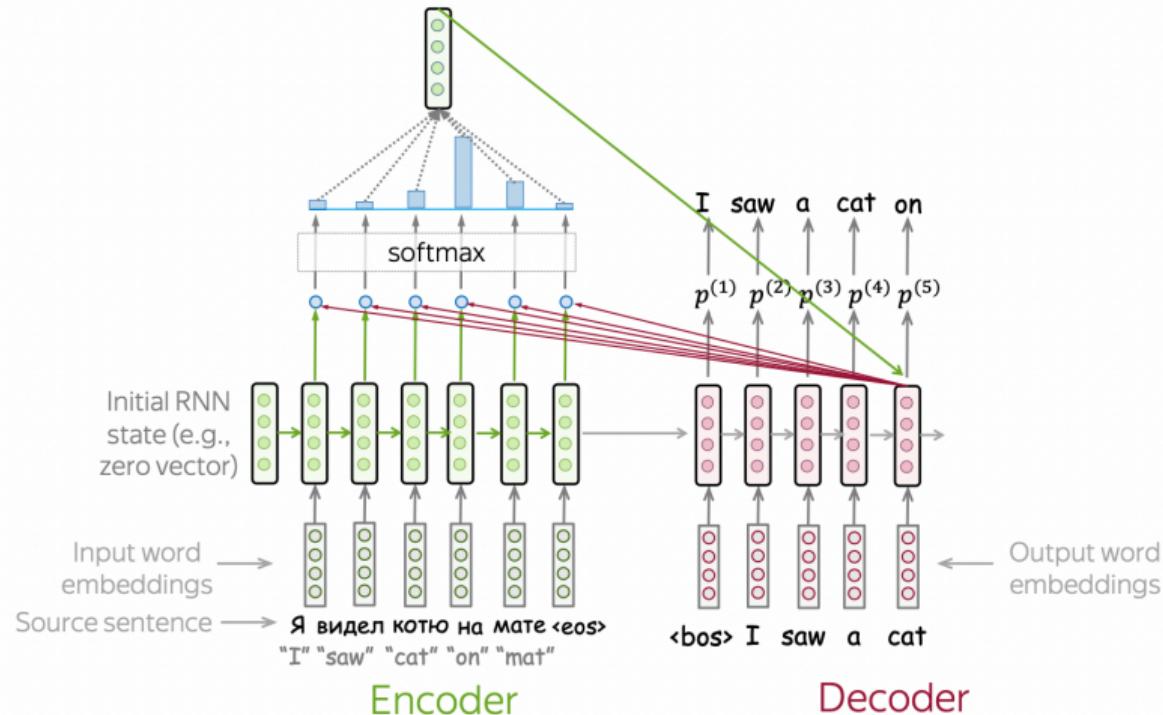
Внимание в декодере



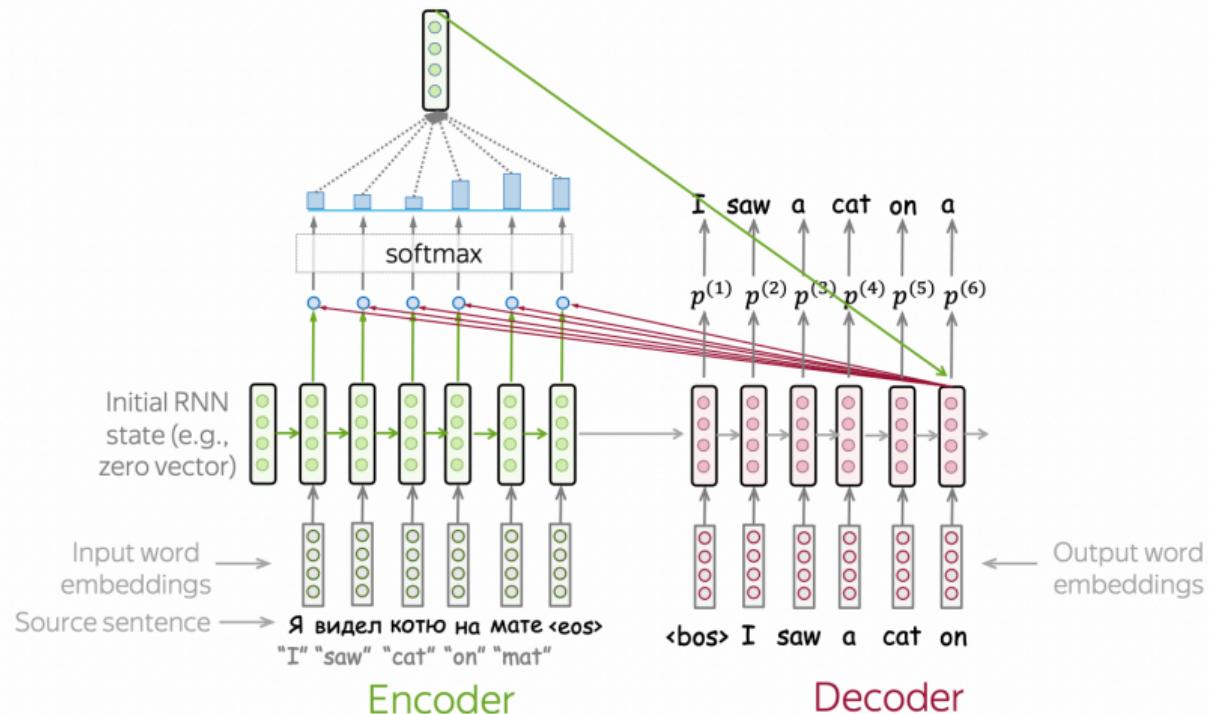
Внимание в декодере



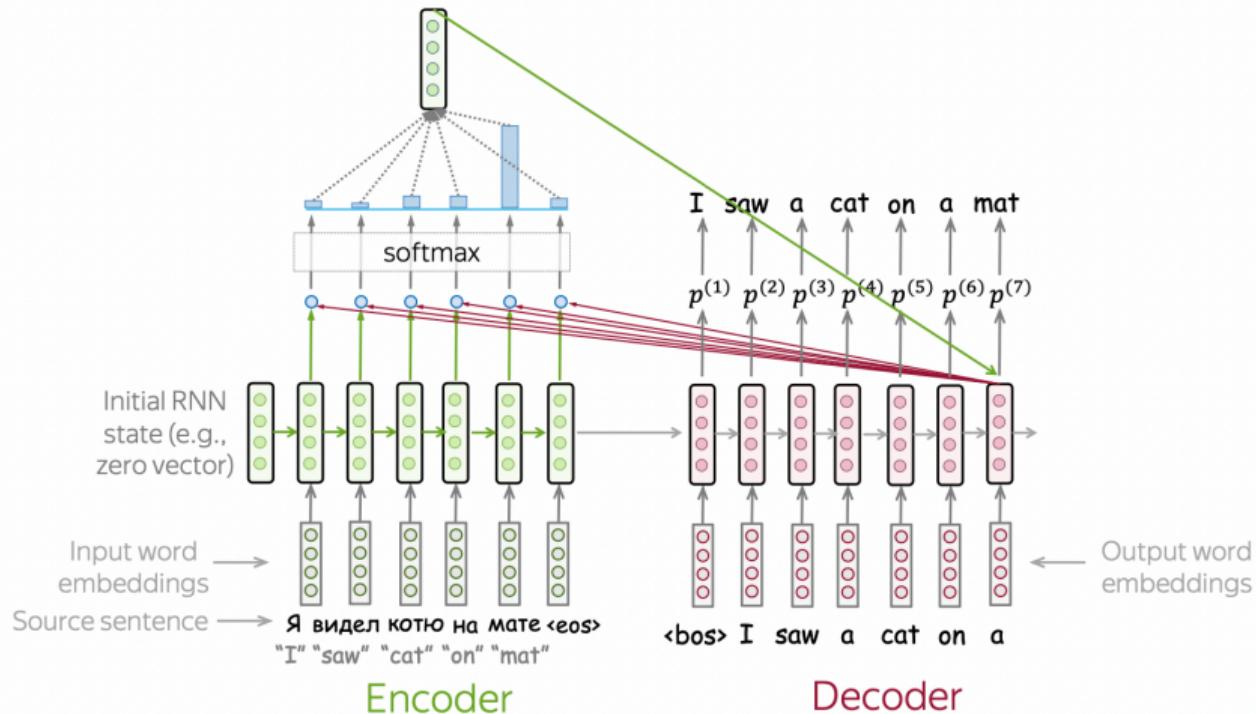
Внимание в декодере



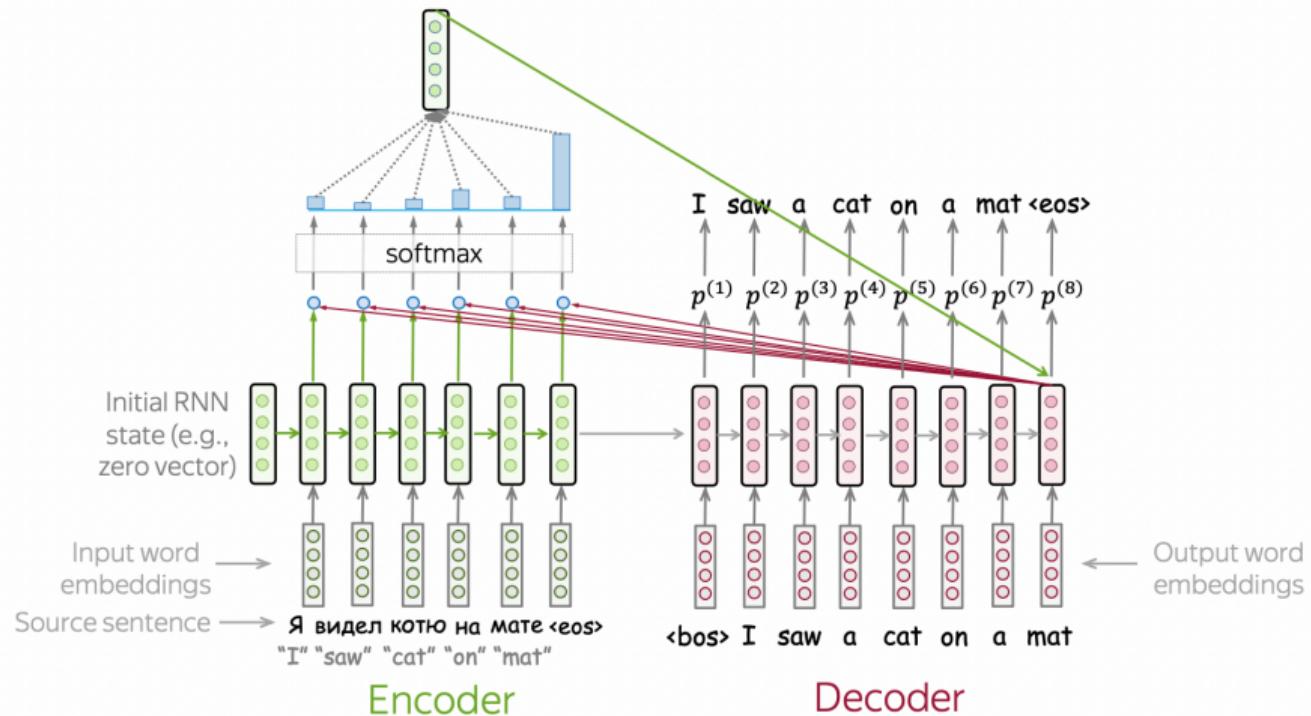
Внимание в декодере



Внимание в декодере



Внимание в декодере



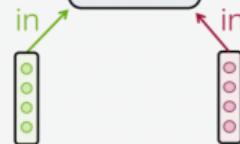
Attention Score Functions

Attention

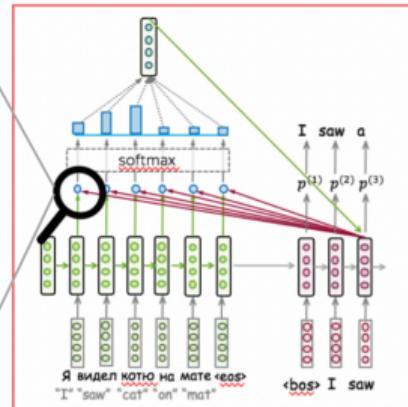
$\text{score}(h_t, s_k)$

scalar
out

Attention
function



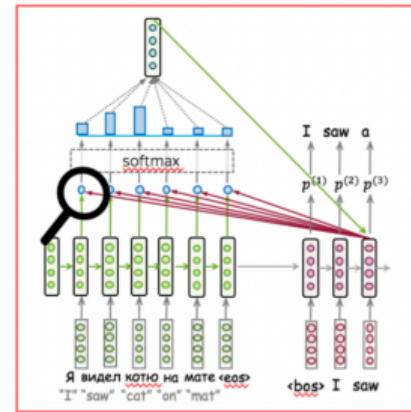
How relevant is
source token k
for target step t ?



Attention Score Functions

- Dot-product: $\text{score}(h_t, s_k) = h_t^T s_k$

$$\begin{matrix} h_t^T \\ \text{---} \end{matrix} \times \begin{matrix} s_k \\ \text{---} \end{matrix}$$



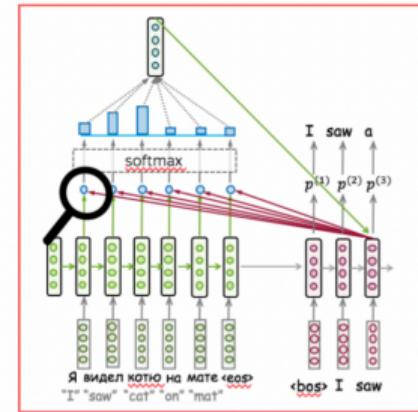
Attention Score Functions

- Dot-product: $\text{score}(h_t, s_k) = h_t^T s_k$

$$\begin{matrix} h_t^T \\ \text{---} \end{matrix} \times \begin{matrix} s_k \\ | \\ \text{---} \end{matrix}$$

- Bilinear: $\text{score}(h_t, s_k) = h_t^T W s_k$

$$\begin{matrix} h_t^T \\ \text{---} \end{matrix} \times \begin{matrix} W \\ | \\ \text{---} \end{matrix} \times \begin{matrix} s_k \\ | \\ \text{---} \end{matrix}$$



Attention Score Functions

- Dot-product: $\text{score}(h_t, s_k) = h_t^T s_k$

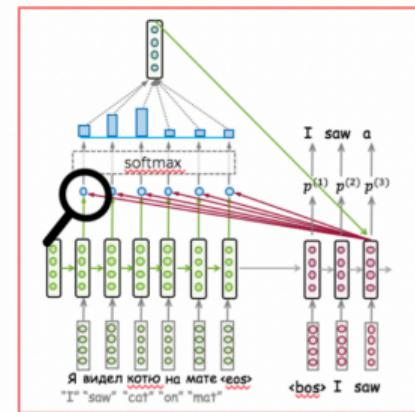
$$h_t^T \times \begin{array}{c} \text{---} \\ | \\ \text{---} \end{array} \times \begin{array}{c} \text{---} \\ | \\ \text{---} \end{array} s_k$$

- Bilinear: $\text{score}(h_t, s_k) = h_t^T W s_k$

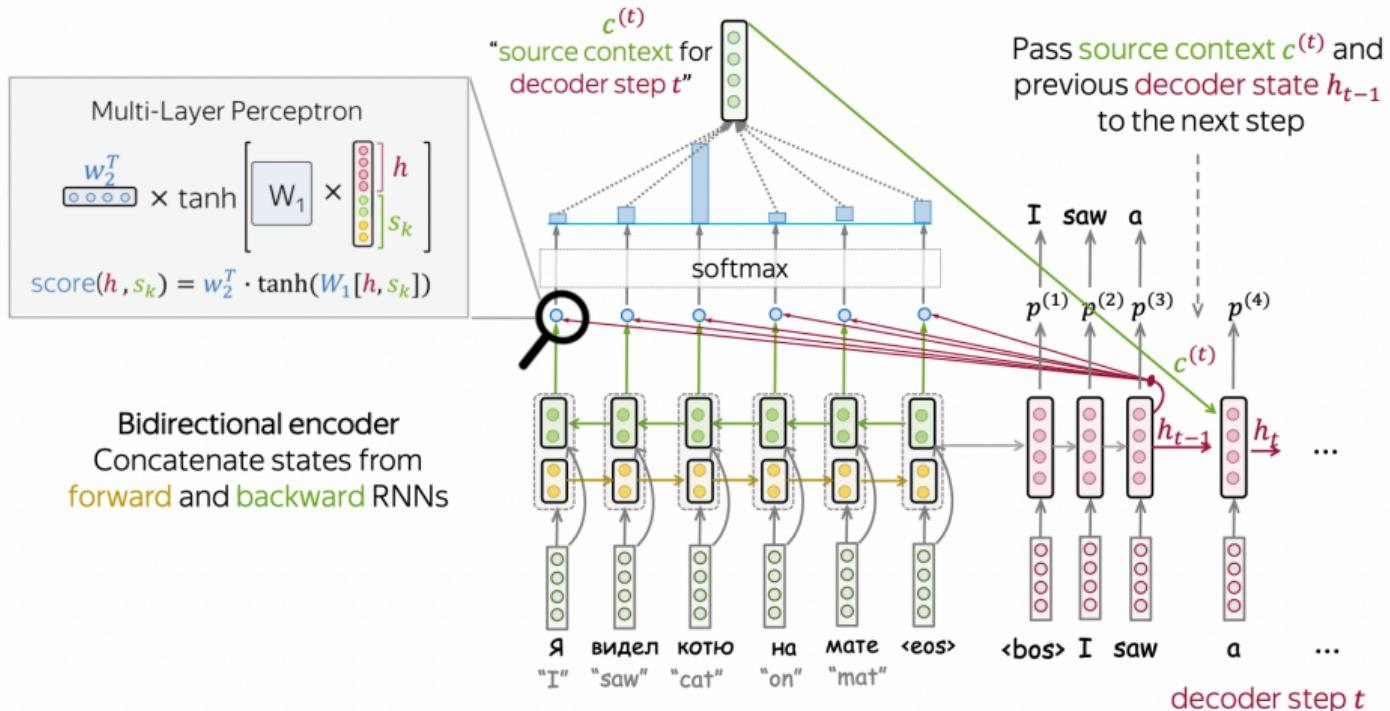
$$h_t^T \times \begin{array}{c} \text{---} \\ | \\ \text{---} \end{array} \times \begin{array}{c} \text{---} \\ | \\ \text{---} \end{array} W \times \begin{array}{c} \text{---} \\ | \\ \text{---} \end{array} \times \begin{array}{c} \text{---} \\ | \\ \text{---} \end{array} s_k$$

- Multi-Layer Perceptron: $\text{score}(h_t, s_k) = w_2^T \cdot \tanh(W_1[h_t, s_k])$

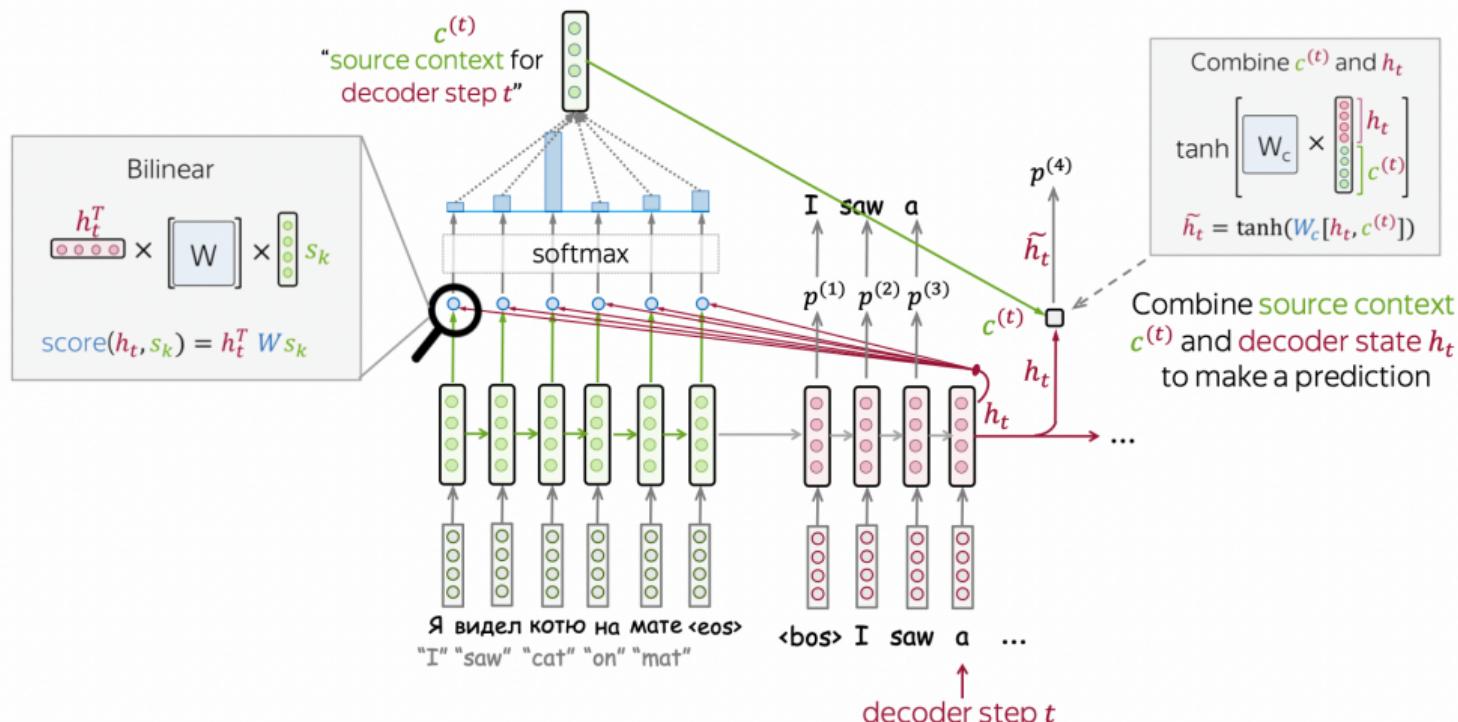
$$w_2^T \times \tanh \left[\begin{array}{c} \text{---} \\ | \\ \text{---} \end{array} W_1 \times \begin{array}{c} \text{---} \\ | \\ \text{---} \end{array} h_t \right] \begin{array}{c} \text{---} \\ | \\ \text{---} \end{array} s_k$$



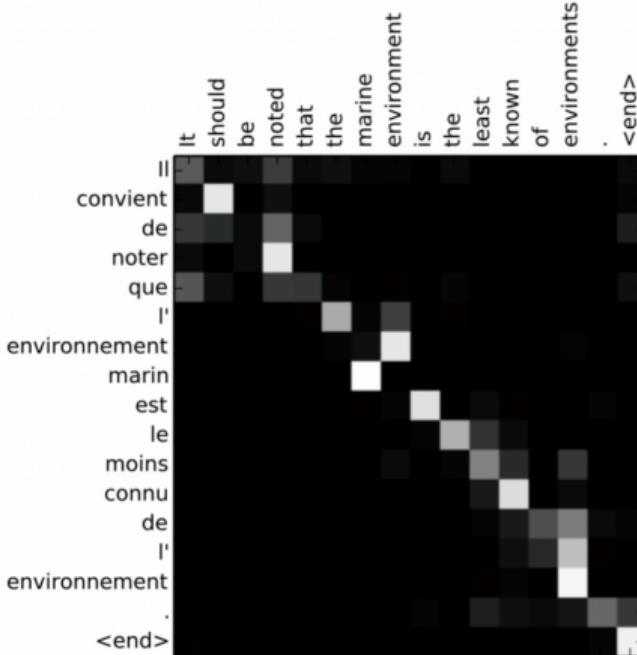
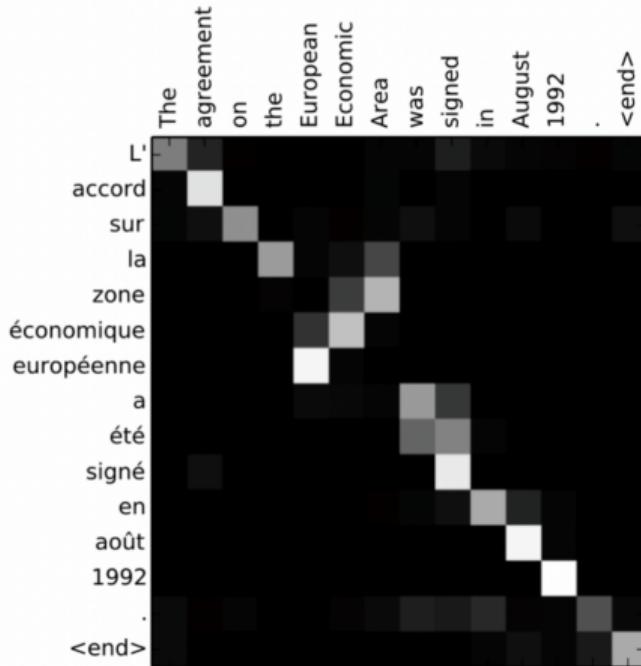
Bahdanau Model (the original attention model)



Luong Model



Attention



Attention is All You Need!

Attention is All You Need! (2017)

- RNN это очень долго! Всегда, чтобы найти следующий токен, надо знать предыдущий
- Backward pass идёт ещё и через время :(
- Transformer — нейросетевая архитектура для задач seq2seq, основанная исключительно на полно связных слоях
- Превзошла существовавшие seq2seq архитектуры как по качеству, так и по скорости работы
- Основной элемент — multi-head self-attention

Attention is All You Need! (2017)

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Llion Jones*
Google Research
llion@google.com

Noam Shazeer*
Google Brain
noam@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Niki Parmar*
Google Research
nikip@google.com

Lukasz Kaiser*
Google Brain
lukaszkaiser@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

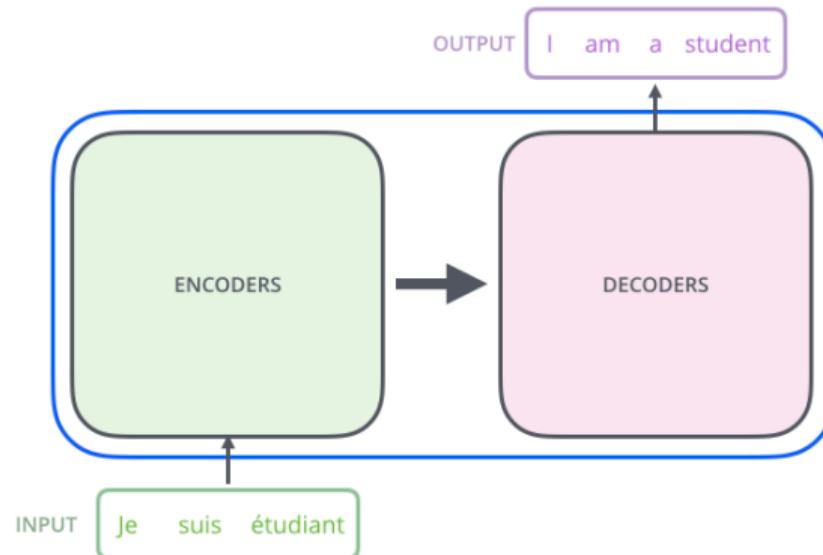
Illia Polosukhin* §
illia.polosukhin@gmail.com

Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

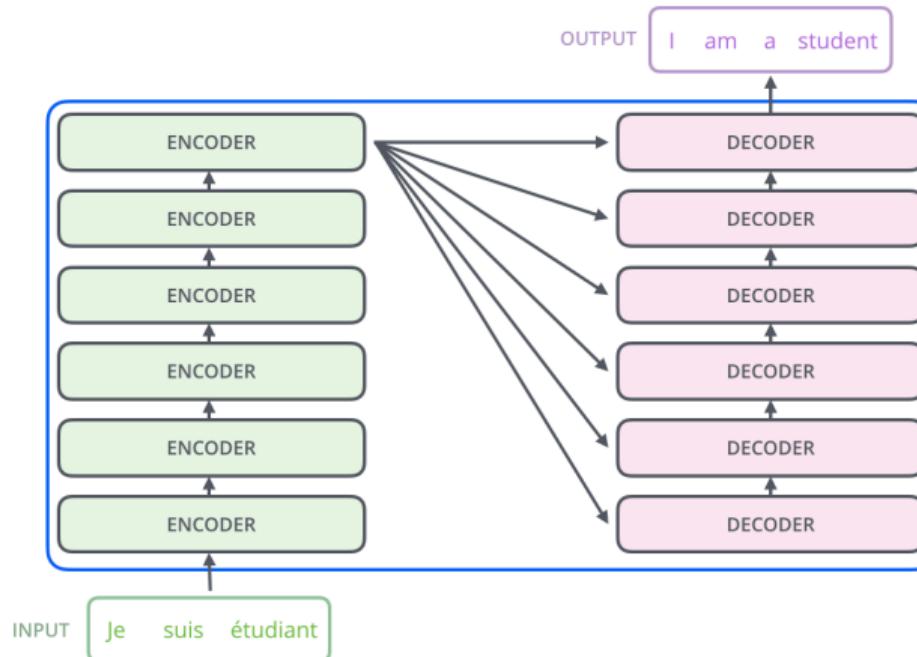
Transformer

Верхнеуровнного - это просто энкодер и декодер



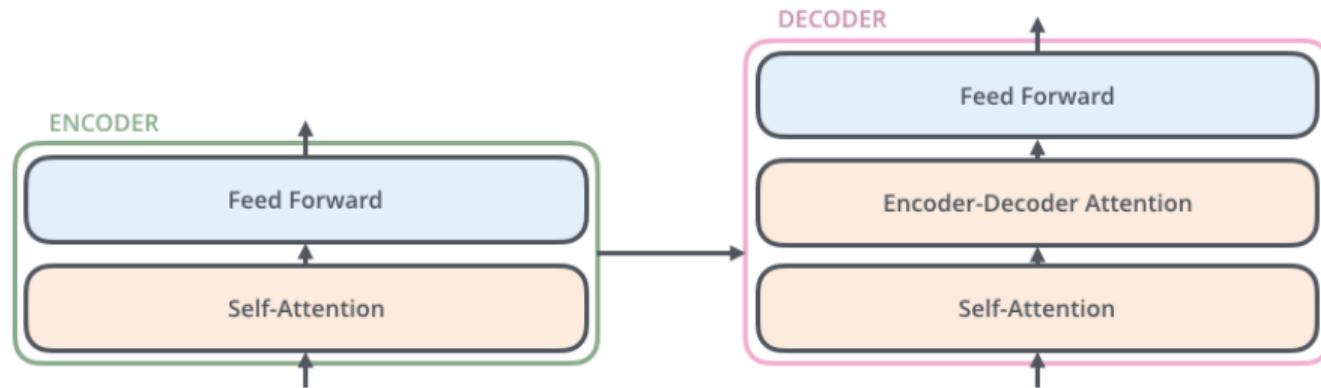
Transformer

Энкодер и декодер состоят из одинаковых блоков; веса во всех блоках разные

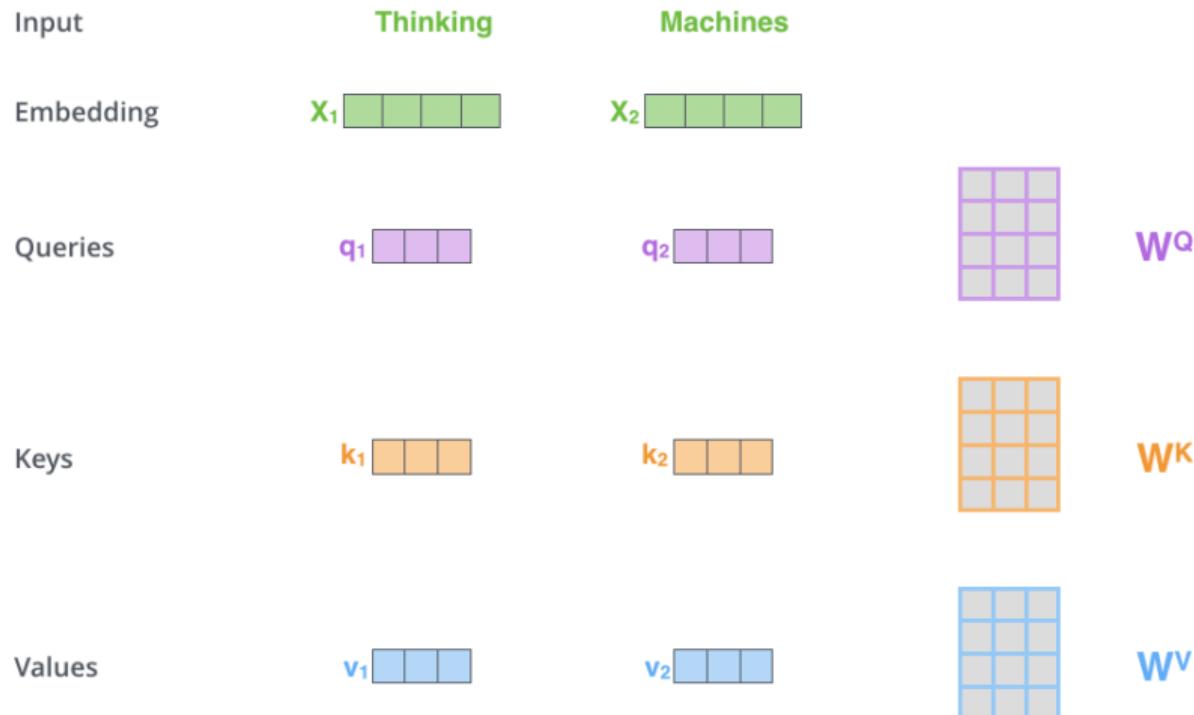


Transformer

В энкодере происходят две вещи: сначала вход прогоняется через self-attention, а затем — через полно связанный слой. В декодере помимо обычного self-attention есть ещё и attention из энкодера.



Self-attention



Абстракции

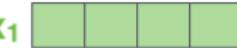
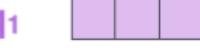
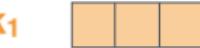
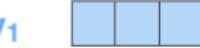
- Для каждого входного слова считаются три вектора: Query, Key и Value
- Матрицы W^Q, W^K, W^V обучаются вместе с моделью
- Value - то, что мы знаем об этом слове
- Query, Key помогают искать связи между словами, мы ходим по всем словам и пытаемся понять насколько они связаны между собой
- Query - мое текущее слово, Key - мое слово с которым я сравниваю себя

Self-attention

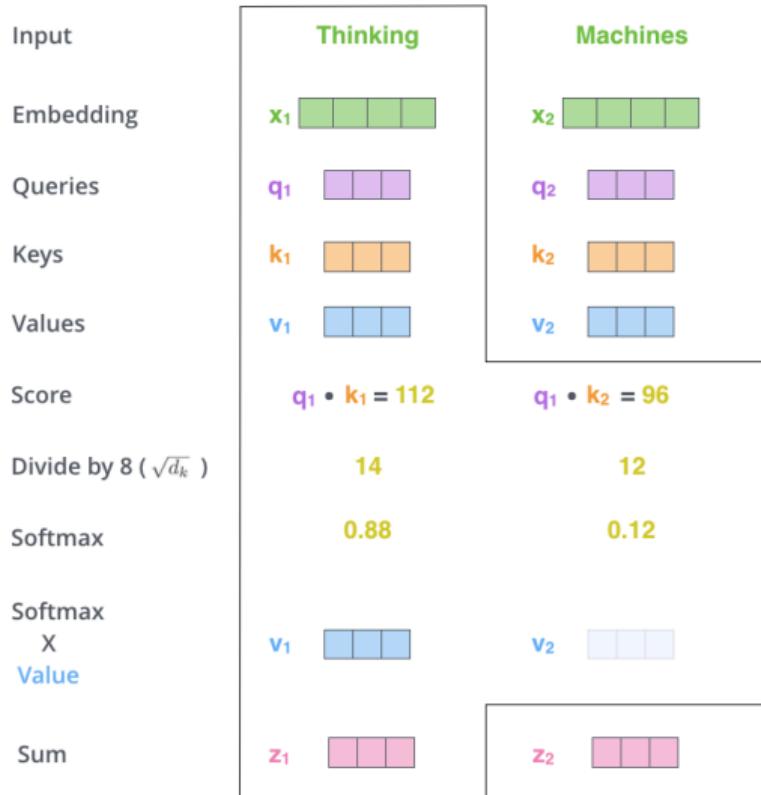
Цель этого слоя — сложить Value с некоторыми весами

$$\text{softmax} \left(\frac{\begin{array}{c} \text{Q} \quad \text{K}^T \\ \text{---} \times \text{---} \\ \begin{matrix} \text{purple} & \text{purple} & \text{purple} \\ \text{purple} & \text{purple} & \text{purple} \end{matrix} \end{array}}{\sqrt{d_k}} \right) \text{V}$$
$$= \begin{matrix} \text{Z} \\ \text{---} \\ \begin{matrix} \text{pink} & \text{pink} & \text{pink} \\ \text{pink} & \text{pink} & \text{pink} \end{matrix} \end{matrix}$$

Более детально

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

Более детально



Query, Key, Value

Each vector receives three representations (“roles”)

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

Query: vector from which the attention is looking

“Hey there, do you have this information?”

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

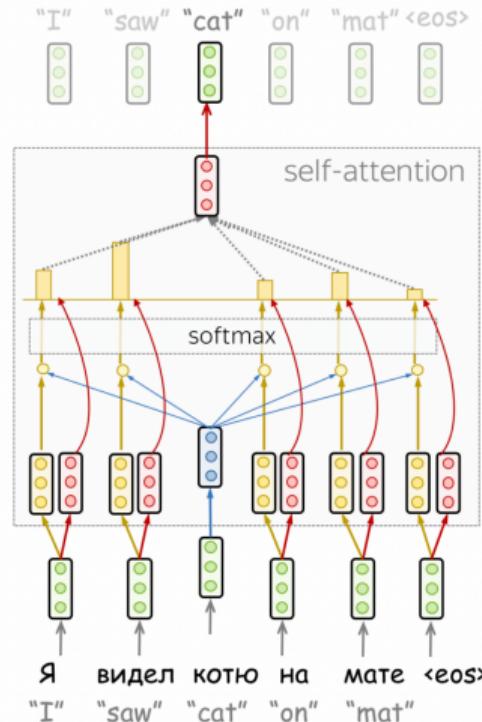
Key: vector at which the query looks to compute weights

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

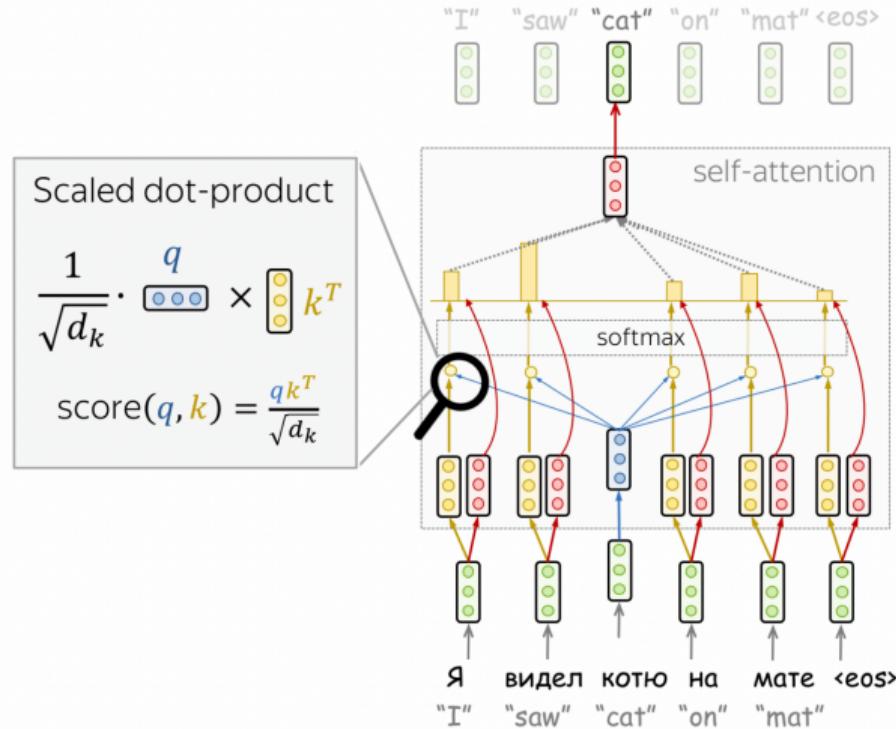
$$[W_V] \times \begin{array}{c} \text{green} \\ \text{green} \\ \text{green} \end{array} = \begin{array}{c} \text{pink} \\ \text{pink} \\ \text{pink} \end{array}$$

Value: their weighted sum is attention output

“Here’s the information I have!”



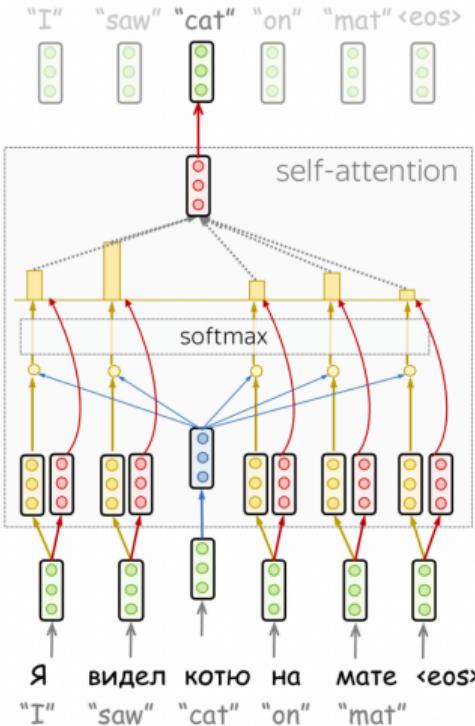
Query, Key, Value



Query, Key, Value

$$\text{Attention}(q, k, v) = \underbrace{\text{softmax}\left(\frac{qk^T}{\sqrt{d_k}}\right)}_{\text{Attention weights}} v$$

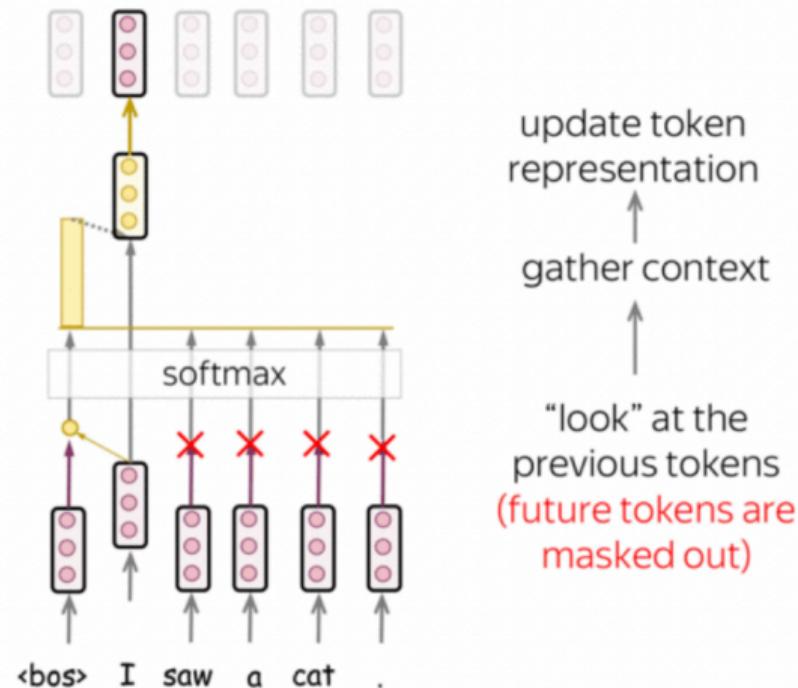
from to
vector dimensionality of K, V



Masked Self-Attention

In the decoder, we forbid looking at future tokens – we don't know them

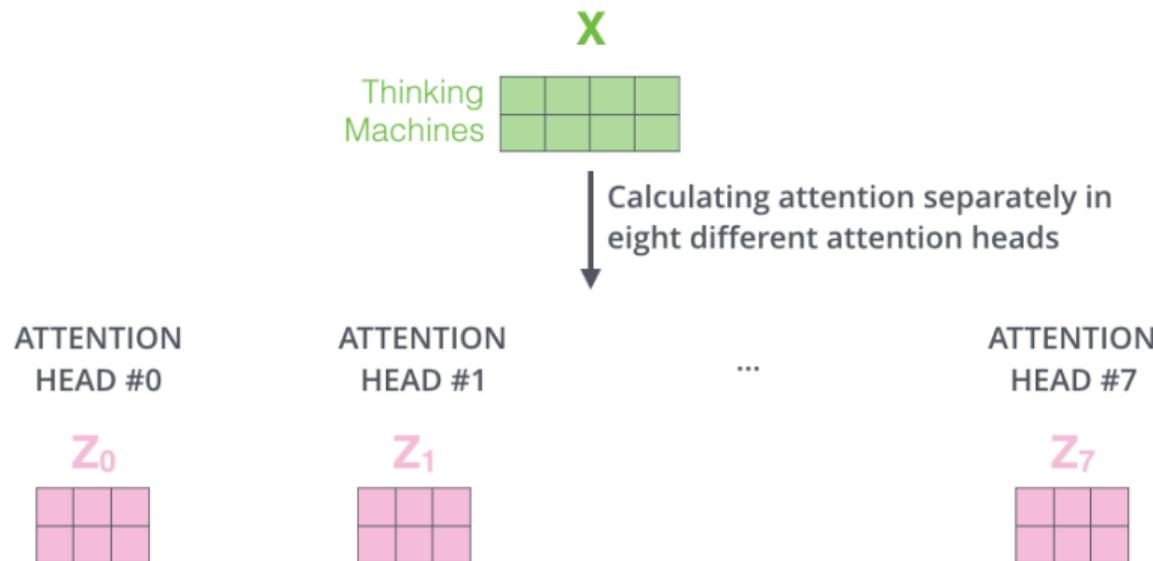
Note: in training, decoder processes all target tokens at once – without masks, it would see future



Зверь с кучей голов

Зверь с кучей голов

Несколько голов обеспечивают разное внимание



Слой целиком

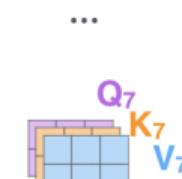
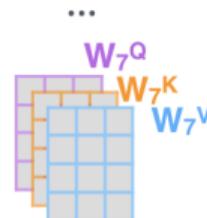
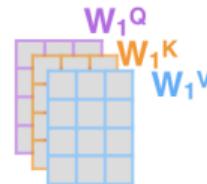
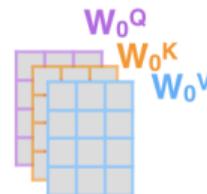
1) This is our input sentence*
2) We embed each word*

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

4) Calculate attention using the resulting $Q/K/V$ matrices

5) Concatenate the resulting Z matrices, then multiply with weight matrix W^o to produce the output of the layer

Thinking
Machines



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



...

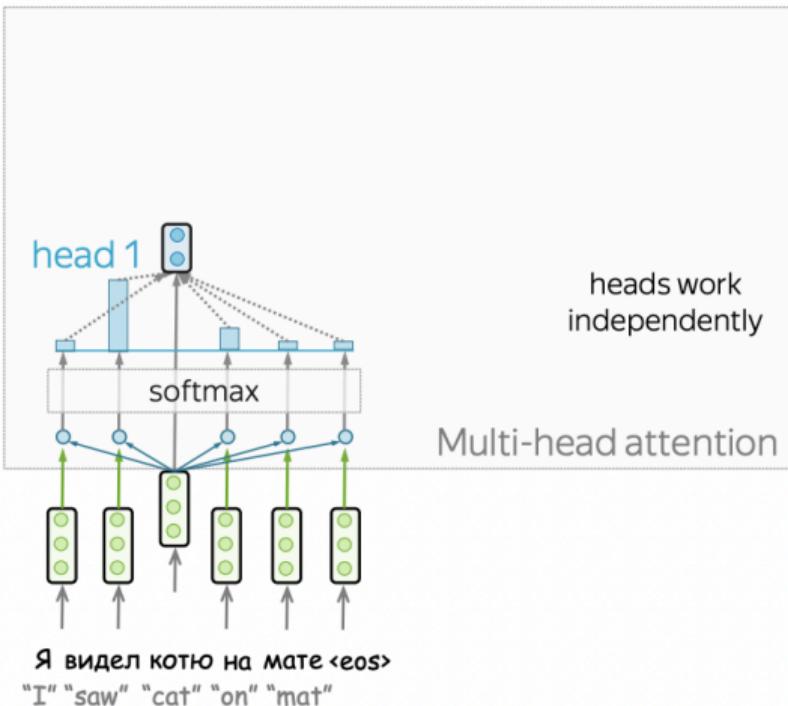


...

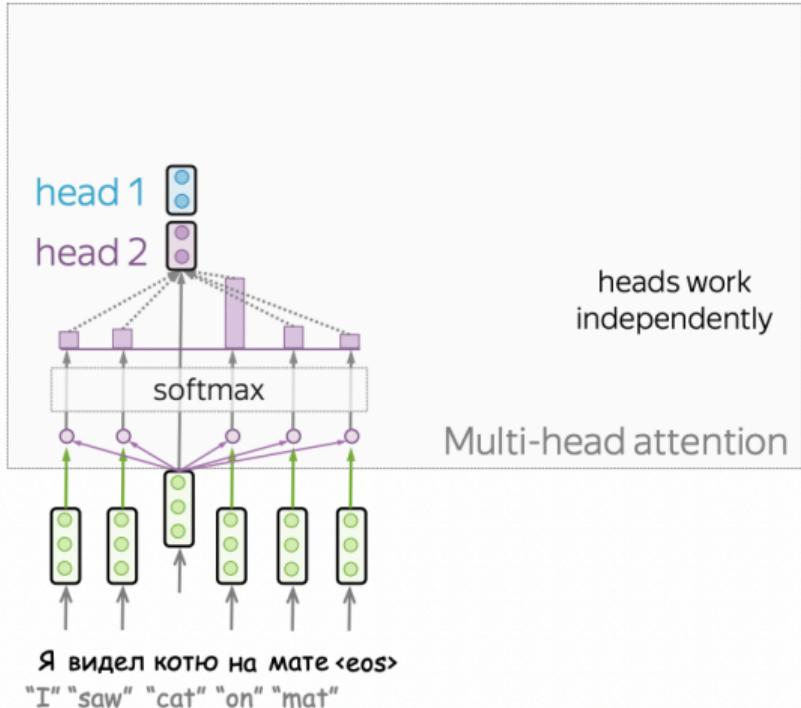
...

...

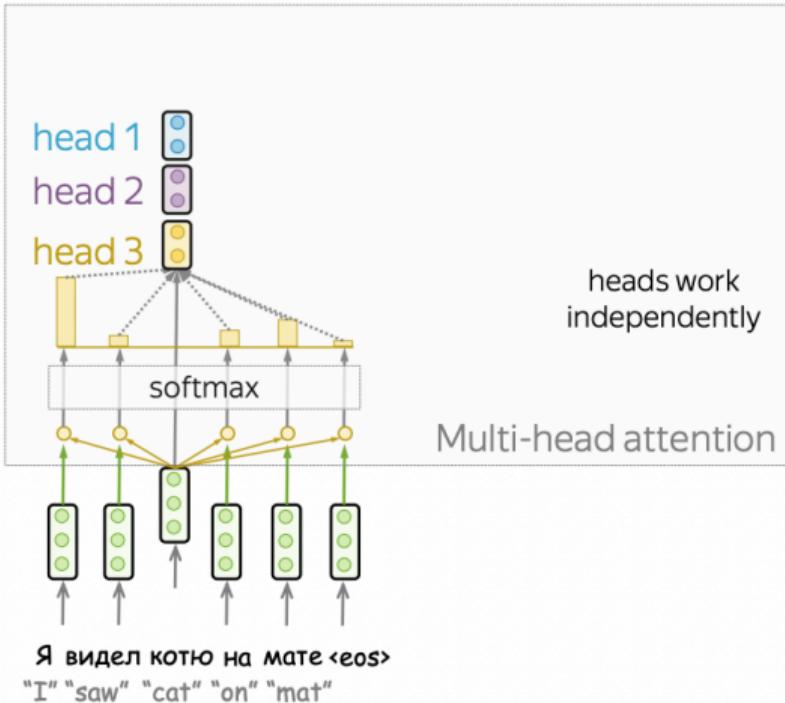
Multi-Head Attention



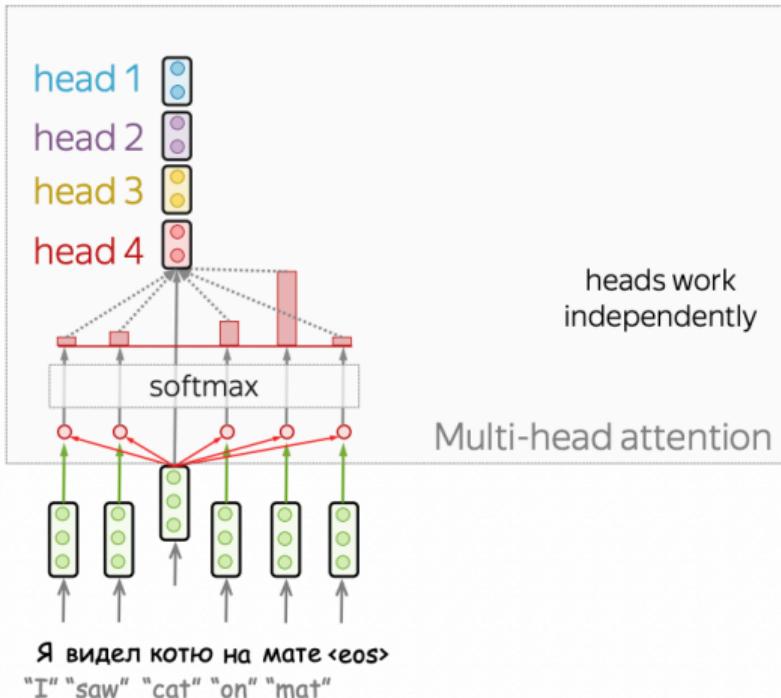
Multi-Head Attention



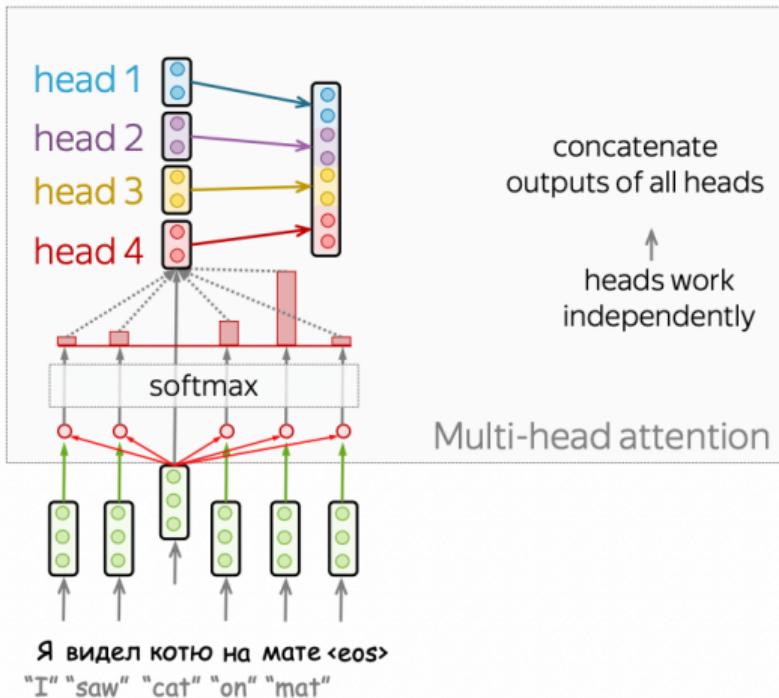
Multi-Head Attention



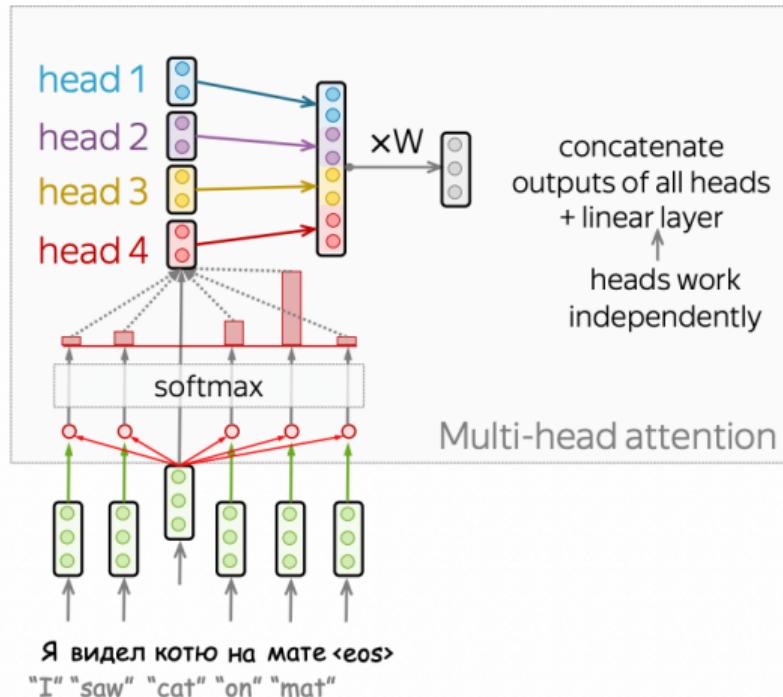
Multi-Head Attention



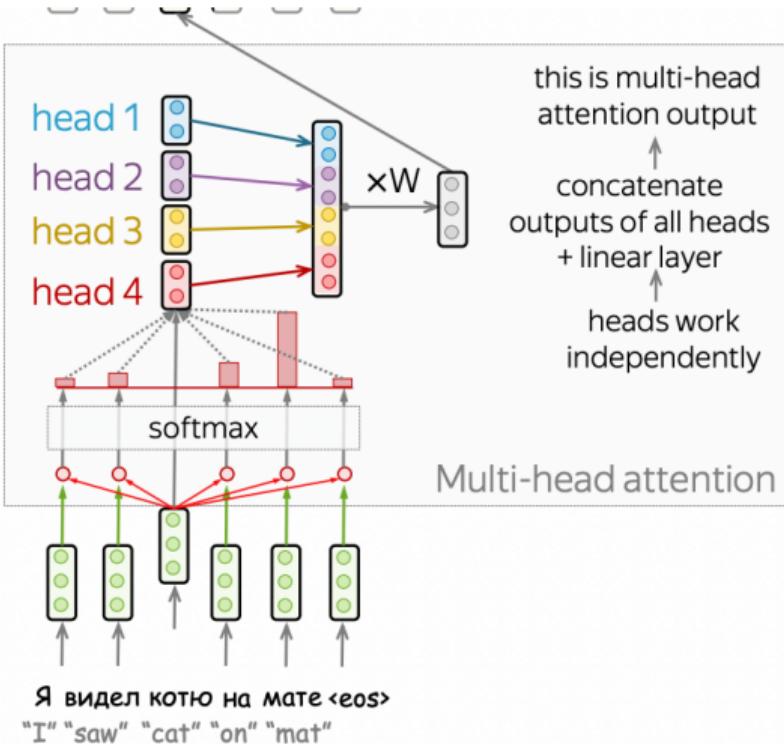
Multi-Head Attention



Multi-Head Attention



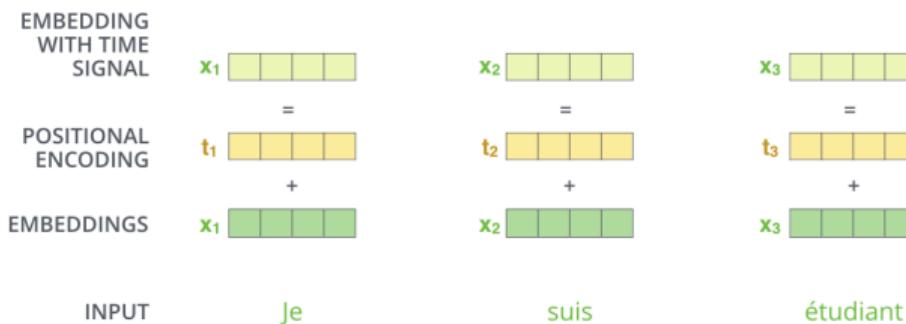
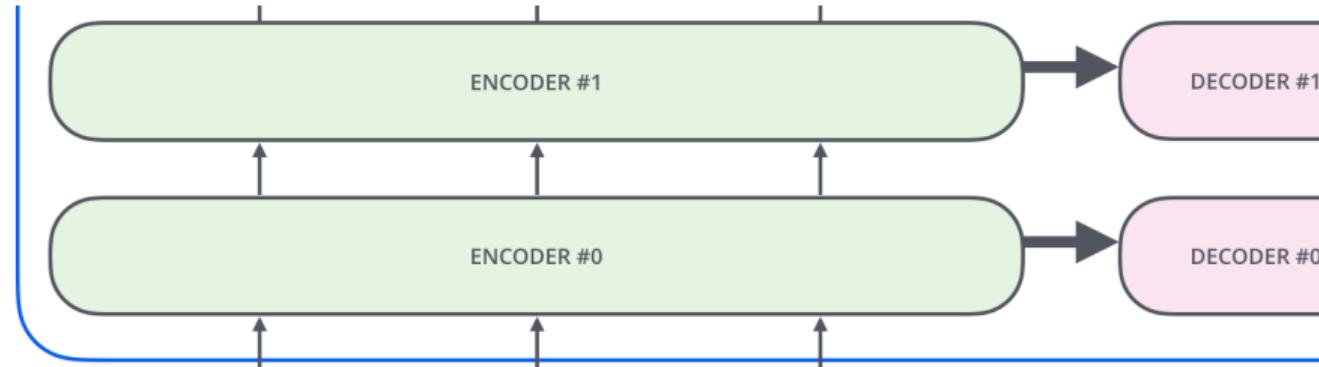
Multi-Head Attention



Positional Encoding

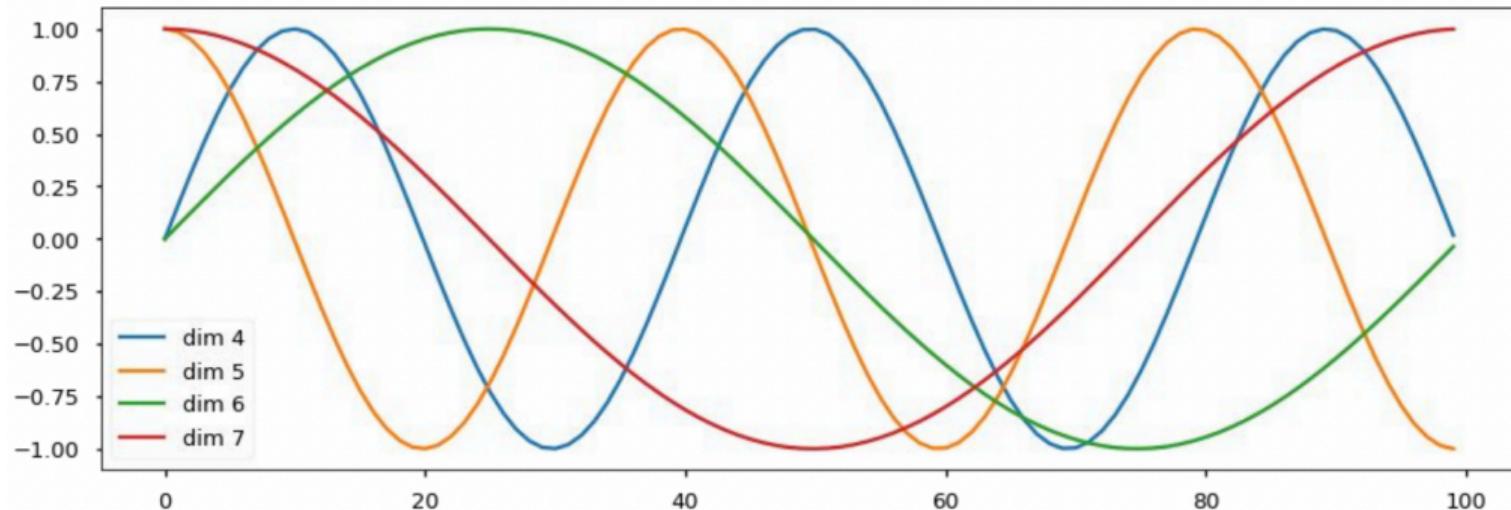
Positional encoding

Для учёта позиции можно приплюсовать дополнительный вектор



Positional encoding

Можно закодировать позицию с помощью синусоиды или какого-нибудь ОНЕ-вектора



Positional encoding

Fixed encodings:

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

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

pos – position, i - dimension

“token x on position k ” →

Input is sum of two embeddings: for token and position

tokens →

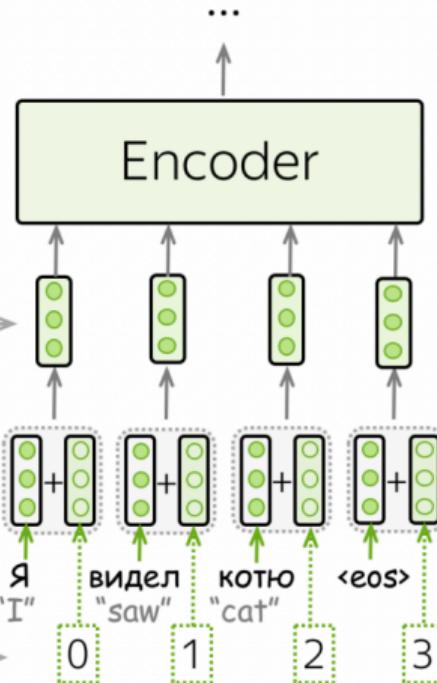
Я
“I”
0

positions →

видел
“saw”
1

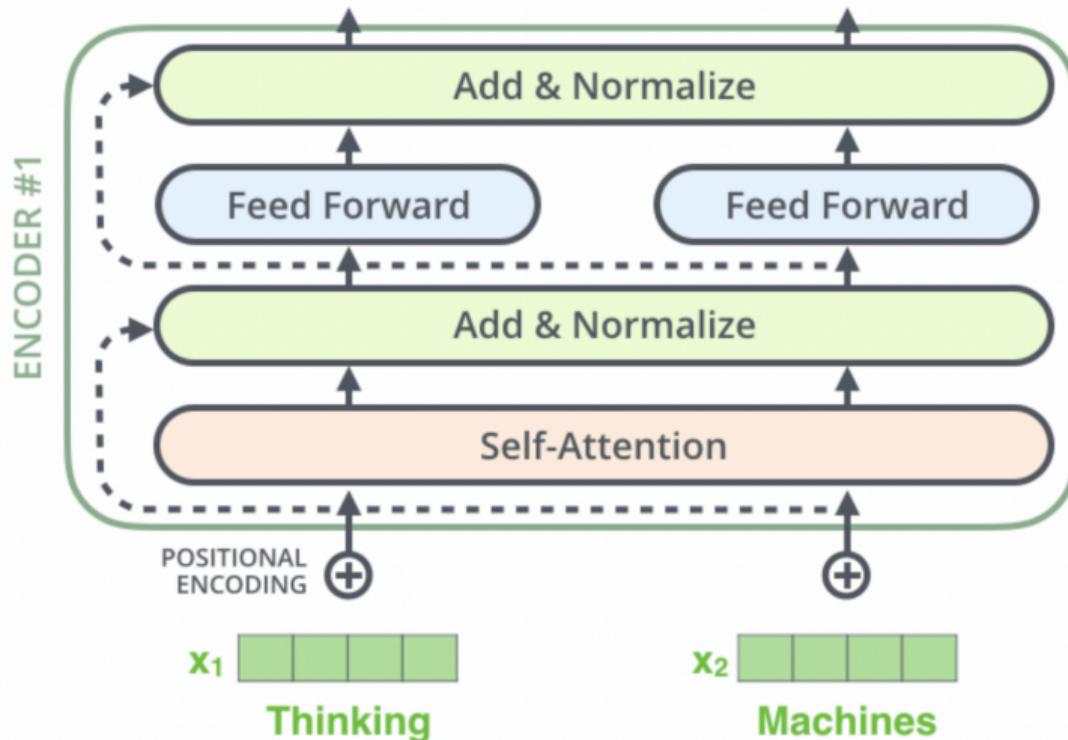
котю
“cat”
2

<eos>
3

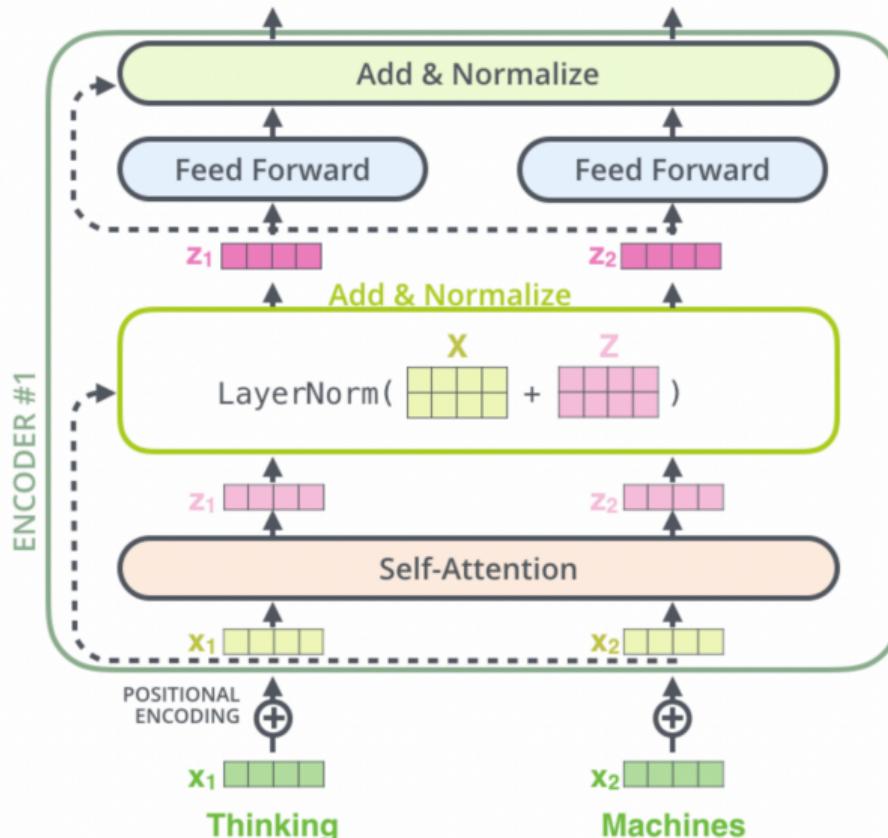


Энкодер и декодер

Что происходит в энкодере

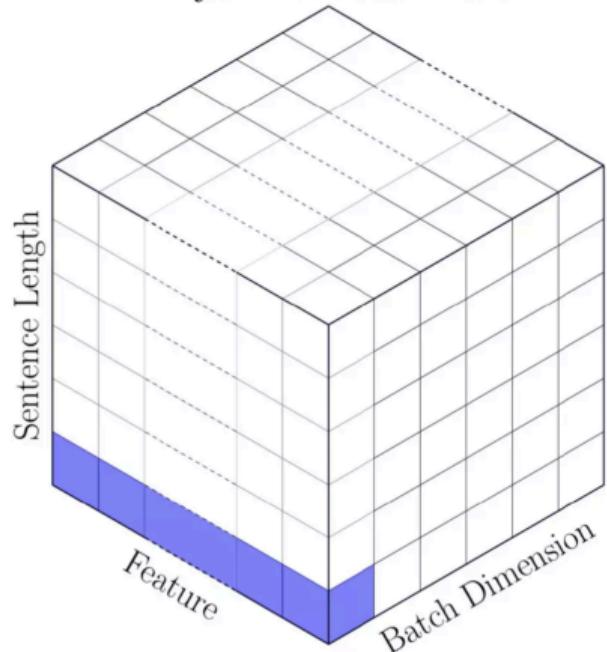


Что происходит в энкодере

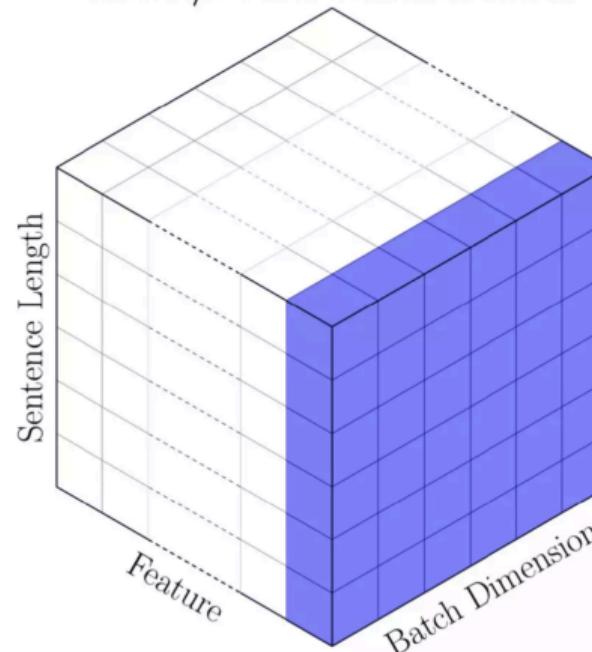


Layer Normalization

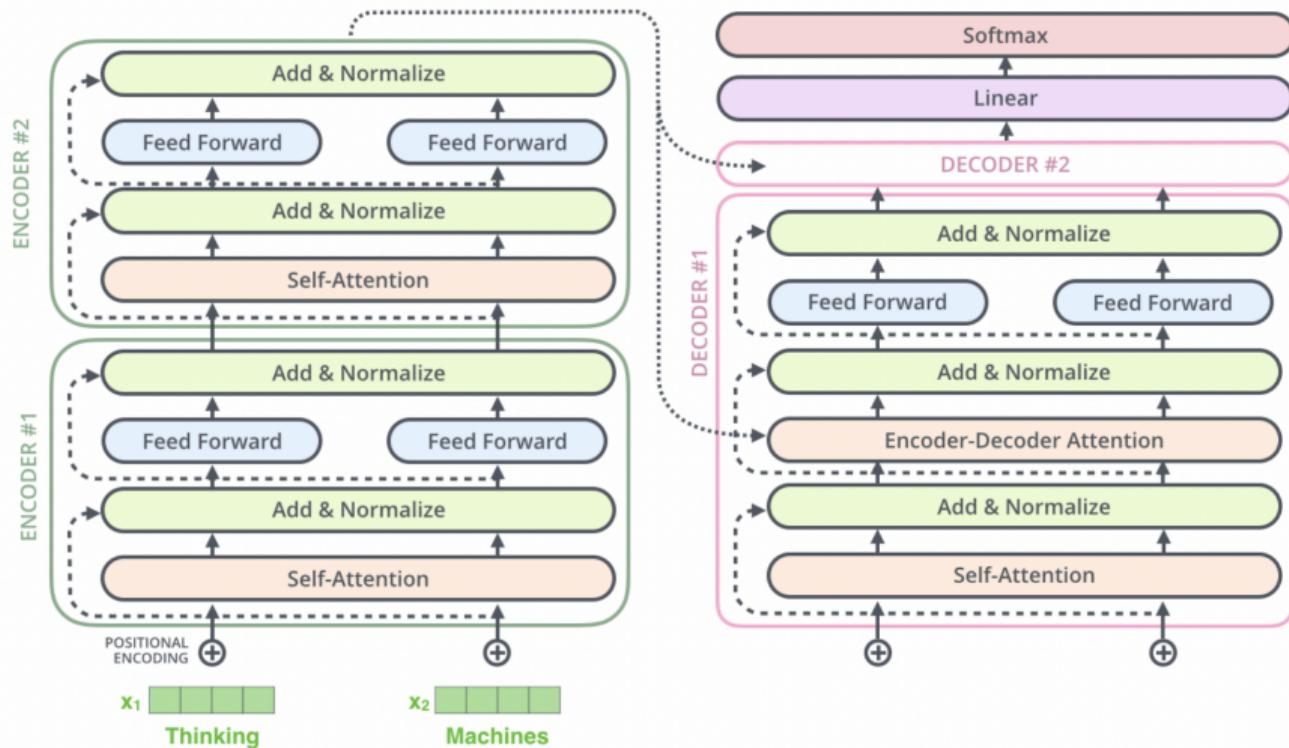
Layer Normalization



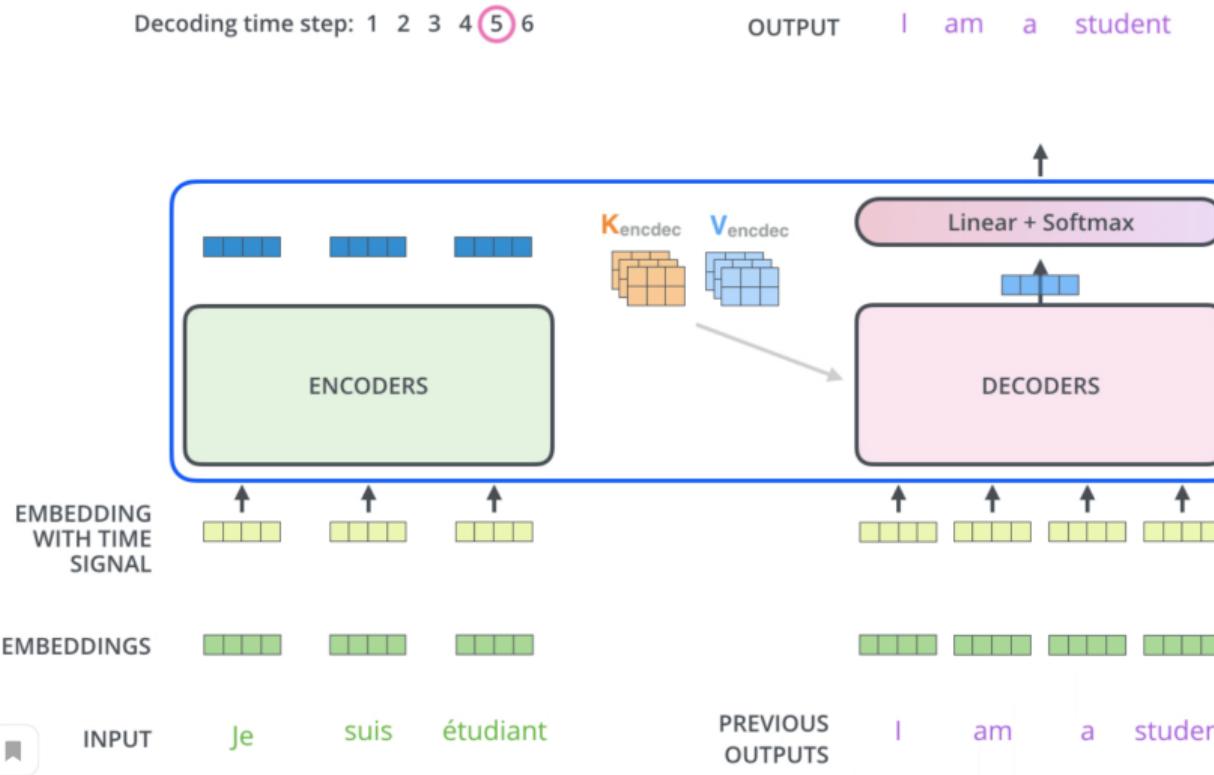
Batch/Power Normalization



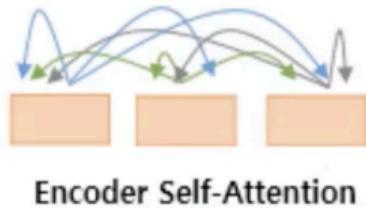
Две модели рядом



Что происходит в декодере?



Encoder-Encoder attention



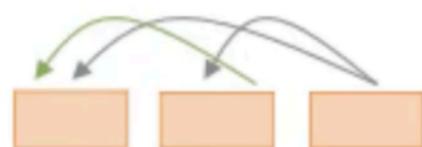
The one we have discussed earlier:

- Multi-head
- Scaled dot-product
- Self-attention: Q, K, V are computed from the same input matrix X

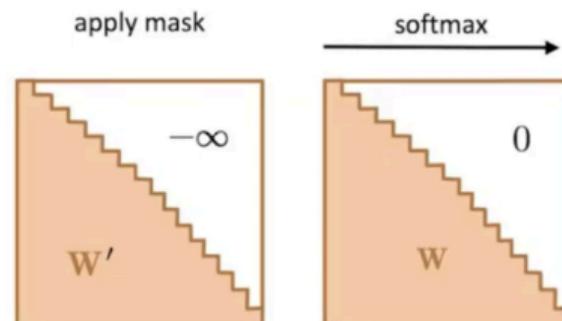
Update embeddings of tokens from the input sequence

Decoder-Decoder attention

We need mask attention matrix to not look at future tokens during training:



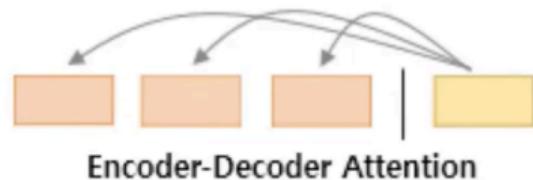
Masked Decoder Self-Attention



$$\mathbf{W}'_{ij} \leftarrow -\infty \text{ if } j > i \quad \mathbf{W} = \text{softmax}(\mathbf{W}')$$

Attend to the previous words in the generated output sequence

Encoder-Decoder attention



This is not self-attention:

- Q from decoder
- K, V from encoder

Attend to tokens from the input sequence relevant
for the generation of the next output token

Обучение

Trainig details

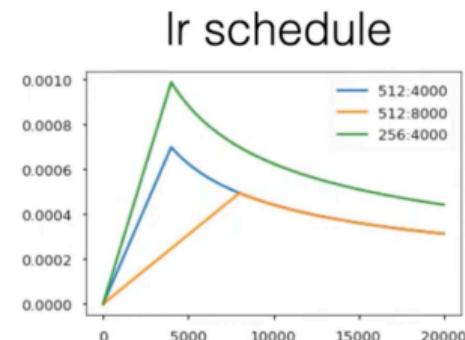
- Loss - standard NLL (cross-entropy) - lower is better:

$$NLL(y_{1:M}) = - \sum_{t=1}^M \log p(y_t | t_{t-1})$$

Instead perplexity is usually reported - higher is better:
 $Perplexity(y_{1:M}) = 2^{\frac{1}{M} NLL(y_{1:M})}$

- Teacher forcing for decoder
- Adam optimizer
- Learning rate schedule with warm-up:

$$lr = d_{model}^{-0.5} \cdot \min(step_num^{-0.5}, step_num \cdot warmup^{-1.5})$$



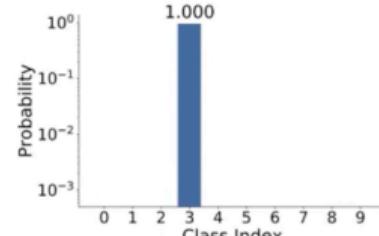
Trainig details

- Label smoothing

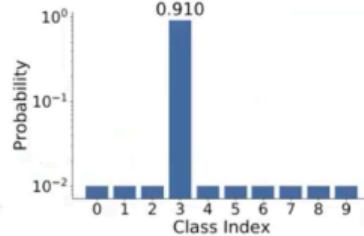
$$y_{ls} = (1 - \alpha) \cdot y_{hot} + \alpha/K$$

- Residual dropout - to the output of each sub-layer (before add+norm) + to the sums of the embeddings and the positional encodings
- BPE/Word-piece (shared embeddings for input/output)
- Model averaging (average last k checkpoints - SWA)

label smoothing



(a) Hard Label



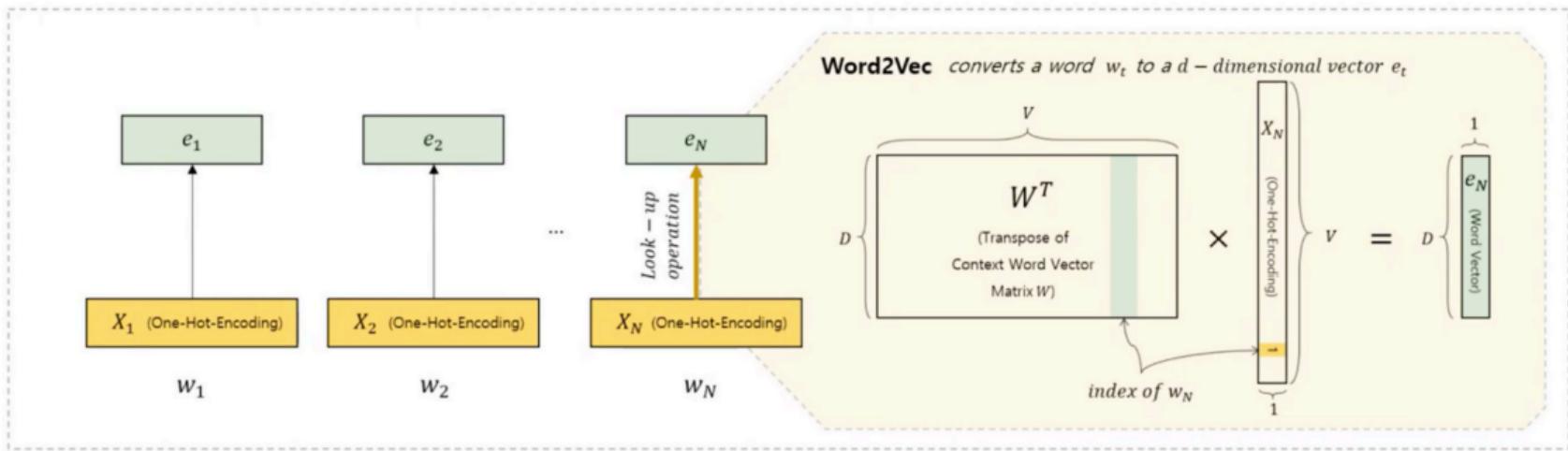
(b) LS

Contextualized Word Embeddings

Embeddings in NLP

- Большие объёмы неразмеченных данных в интернете в разных доменах (книги, новости, википедия, иные тексты из интернет-страниц)
- Размеченных данных мало. Качественная разметка дорогая и долгая
- Много вычислительных ресурсов, GPU, TPU, фреймворки распределённых вычислений
- Может ли мы как-то заиспользовать имеющиеся ресурсы?

Individual embeddings



word2vec, GLOVE, fasttext ...

Contextualised embeddings

- Обучаем большую модель трансформер на какой-нибудь unsupervised задаче на очень больших данных (очень долго, порядка нескольких недель);
- Дообучаем модель на конкретную задачку на малом корпусе размеченных данных (очень быстро, порядка 1 часа на одной ГПУ).

Contextualised embeddings

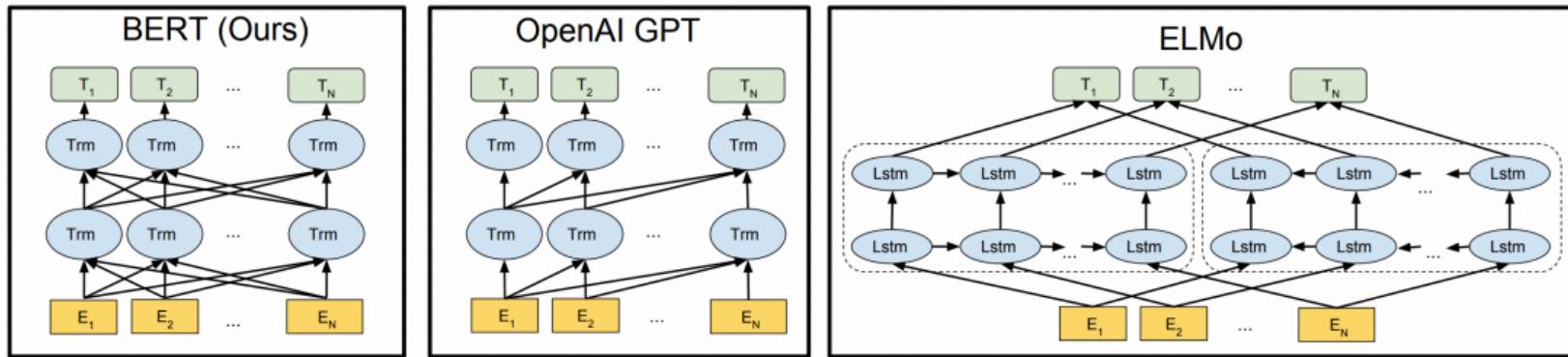


Figure 3: Differences in pre-training model architectures. BERT uses a bidirectional Transformer. OpenAI GPT uses a left-to-right Transformer. ELMo uses the concatenation of independently trained left-to-right and right-to-left LSTMs to generate features for downstream tasks. Among the three, only BERT representations are jointly conditioned on both left and right context in all layers. In addition to the architecture differences, BERT and OpenAI GPT are fine-tuning approaches, while ELMo is a feature-based approach.

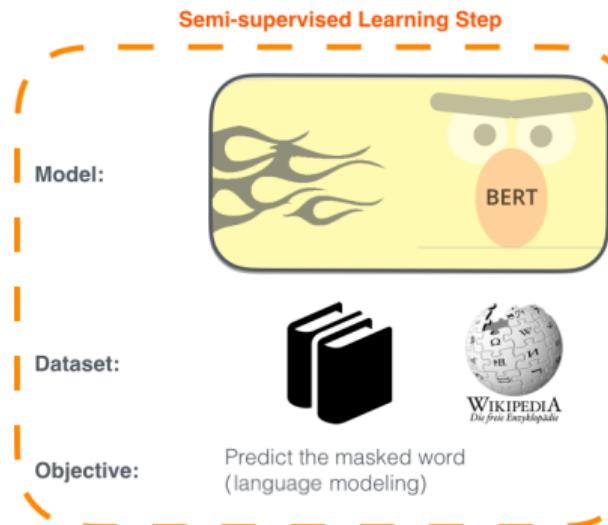
BERT (2018)

- Bidirectional Encoder Representations from Transformers
- Предобученный на искусственно придуманной задаче энкодер
- Его можно файн-тюнить для решения других задач

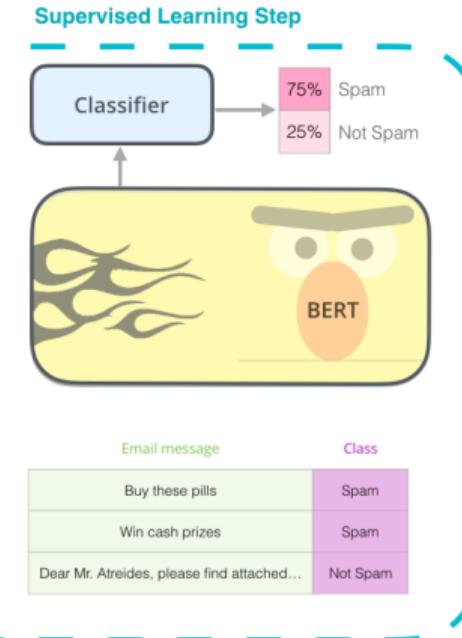
BERT (2018)

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

- **Masked Language Model:** выкидываем часть слов (обычно 15%), модель должна их восстановить:
 - 80% меняем на токен <MASK>
 - 10% меняем на случайные
 - 10% не меняем

Мы предсказываем пропуски, считаем Loss по всем 15%. Случайные и неизменные токены позволяют не переобучиться под токен <MASK>.

- **Next Sentence Prediction:** подаём на вход два предложения, модель должна понять в каком порядке они идут. Правда ли второе является продолжением первого. Эта задача помогает модели понять концепцию предложения.

CLS-токен

Пример токенезированного предложения:

[CLS] my dog is cute [SEP] he likes play ##ing [SEP]

- [CLS] — специальный токен, который мы добавляем к началу любого предложения, он несет информацию о всей входной последовательности (а вход состоит из двух предложений)
- [SEP] — специальный токен, который говорит, что слева от него первое предложение, а справа — второе
- ## — специальный токен, обозначающий, что это кусок слова

CLS-токен

Пример токенезированного предложения:

[CLS] my dog is cute [SEP] he likes play ##ing [SEP]

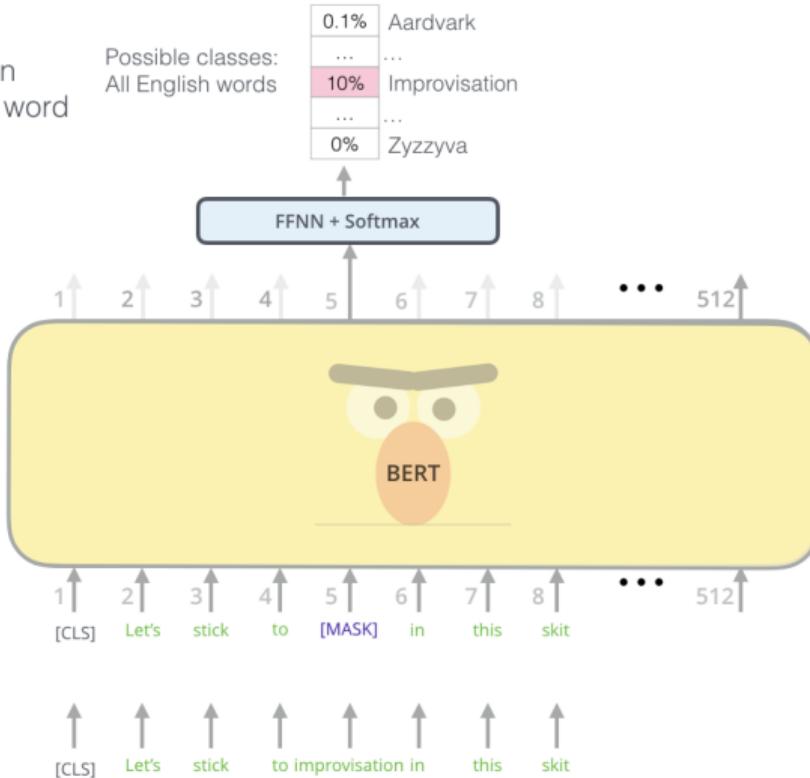
Мы токенезируем предложение, для каждого токена учим эмбеддинг. Поскольку у нас трансформер, мы добавляем к эмбеддингам слов позиционный эмбеддинг и дополнительно добавляем эмбеддинг, отвечающий за то, в каком из двух предложений мы находимся. Все три эмбеддинга суммируются.

$$E = E_{word} + E_{pos} + E_{sent}$$

BERT

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

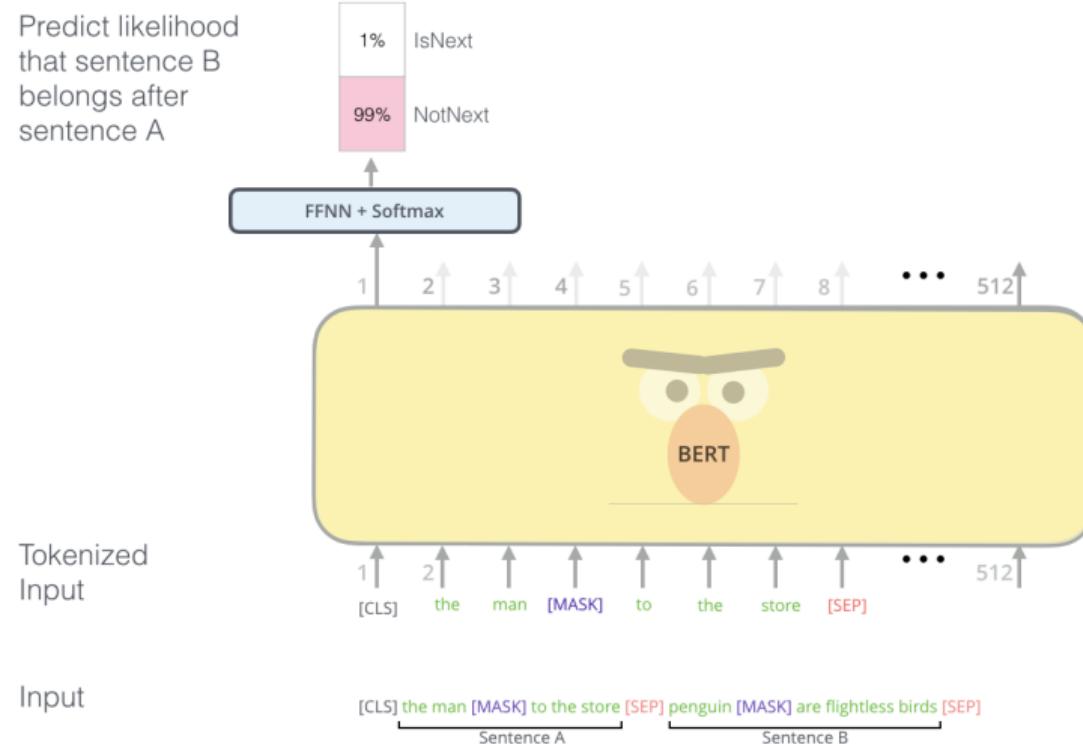
Possible classes:
All English words



Randomly mask
15% of tokens

Input

BERT



Как применять (finetuning)

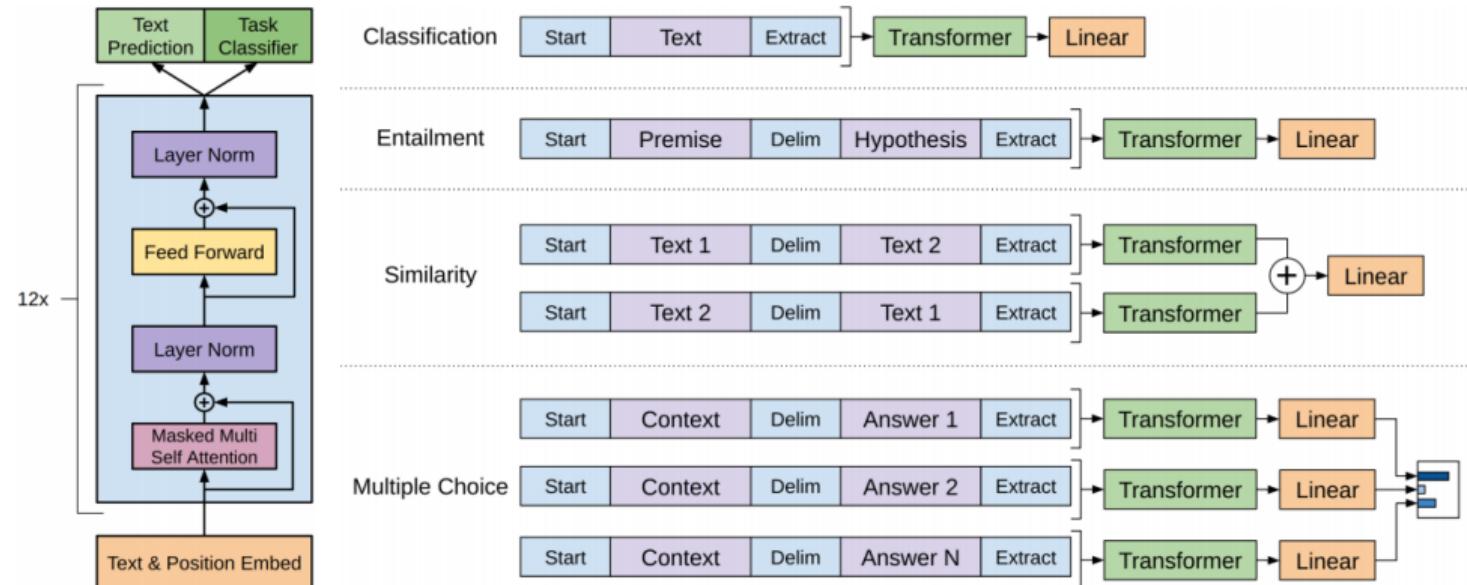


Figure 1: (**left**) Transformer architecture and training objectives used in this work. (**right**) Input transformations for fine-tuning on different tasks. We convert all structured inputs into token sequences to be processed by our pre-trained model, followed by a linear+softmax layer.

RoBERTa: A Robustly Optimized BERT (2019)

BERT был не дотюнен

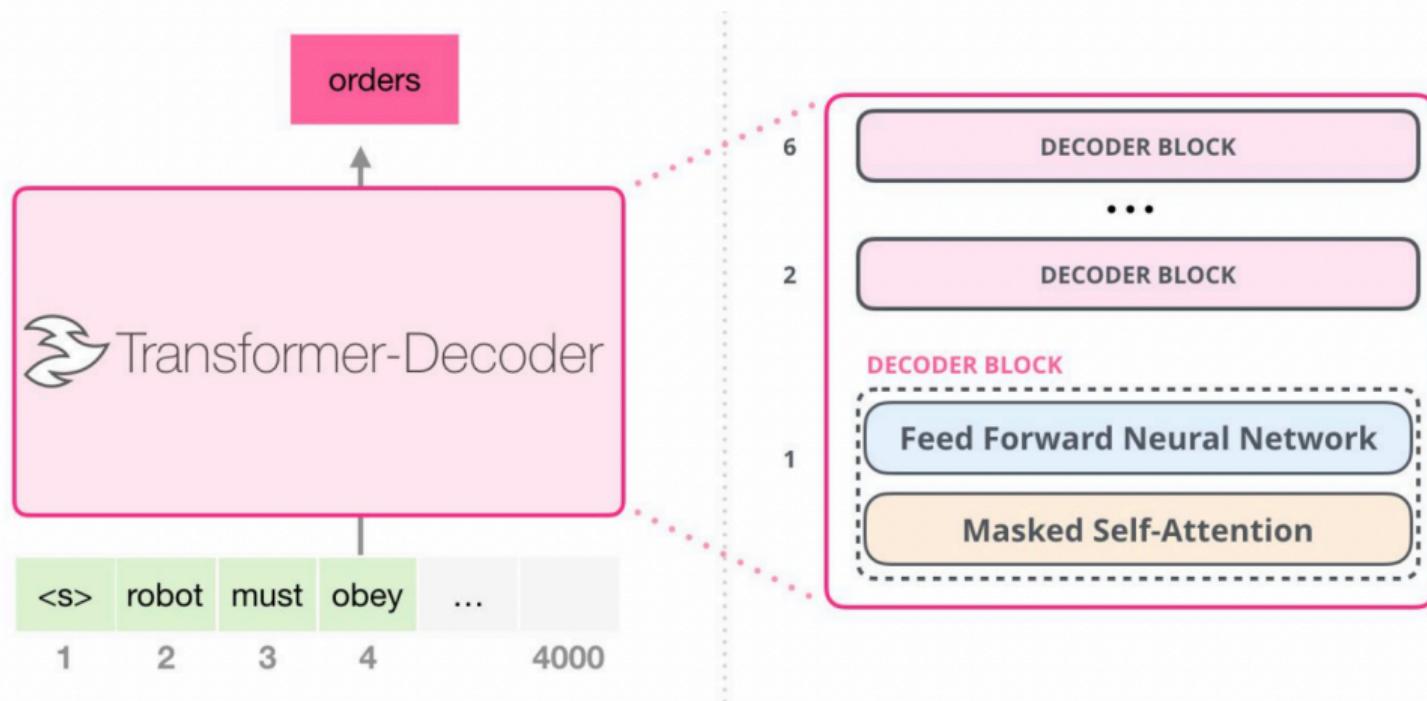
- Взять больше данных, тренировать дольше
- Next sentence prediction лишний
- Более длинные предложения
- Большие батчи
- Динамическое маскирование

<https://arxiv.org/abs/1907.11692>

GPT (Generative Pre-trained Transformer)

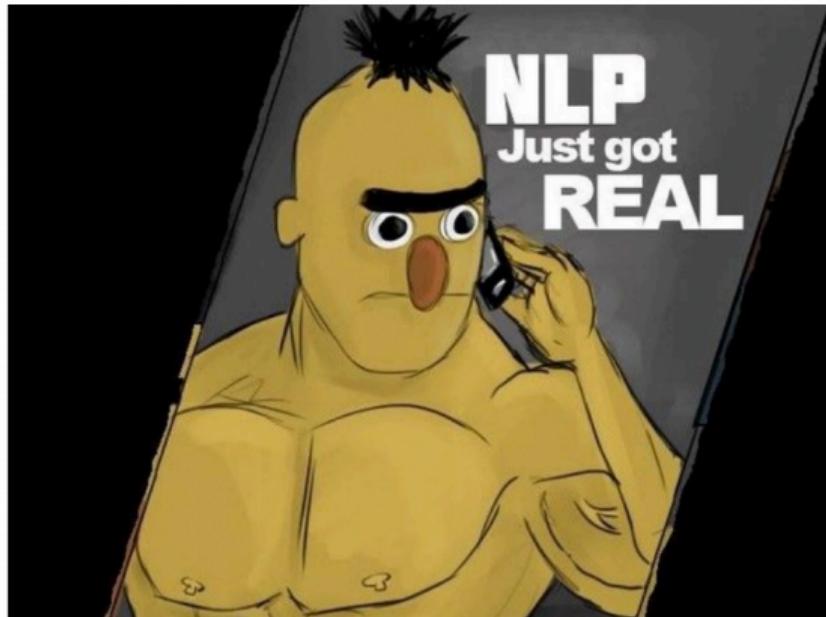
- Языковая модель на декодере трансформера
- Умеет генерить продолжение текста. Настолько хорошо, что OpenAI отказался её публиковать и устроил мощный PR
- Публикует понемногу, начиная с маленьких моделей
- Разные языковые модели на трансформерах можно попробовать здесь: <https://transformer.huggingface.co/>

GPT



Language Model Zoo

- ELMo
- ULMFiT
- GPT
- BERT (BioBERT, ClinicalBERT, ...)
- ERNIE
- XLNet
- KERMIT
- ERNIE 2.0
- GPT-2
- ALBERT
- ...



<https://arxiv.org/pdf/1810.04805.pdf>