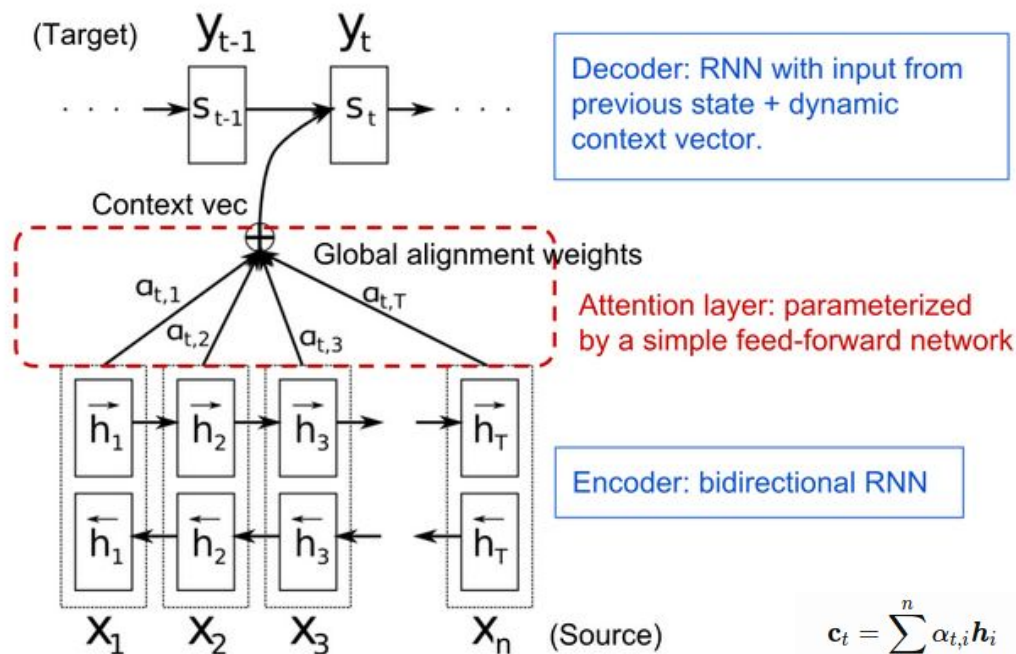


Трансформер

(в схемах и мемах)

Мотивация

- Обучение RNN - трудоёмкая задача
- Хотим избавиться от RNN, не теряя возможность учитывать дальнoдействующие зависимости между словами
- Оказывается, что можно обойтись механизмом Внимания



Additive Attention

$$\mathbf{h}_i = [\vec{h}_i^T; \overleftarrow{h}_i^T]^T, i = 1, \dots, n$$

$$\mathbf{s}_t = f(\mathbf{s}_{t-1}, \mathbf{y}_{t-1}, \mathbf{c}_t)$$

$$\mathbf{c}_t = \sum_{i=1}^n \alpha_{t,i} \mathbf{h}_i$$

; Context vector for output y_t

$$\alpha_{t,i} = \text{align}(y_t, x_i)$$

; How well two words y_t and x_i are aligned.

$$= \frac{\exp(\text{score}(\mathbf{s}_{t-1}, \mathbf{h}_i))}{\sum_{i'=1}^n \exp(\text{score}(\mathbf{s}_{t-1}, \mathbf{h}_{i'}))}$$

; Softmax of some predefined alignment score..

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

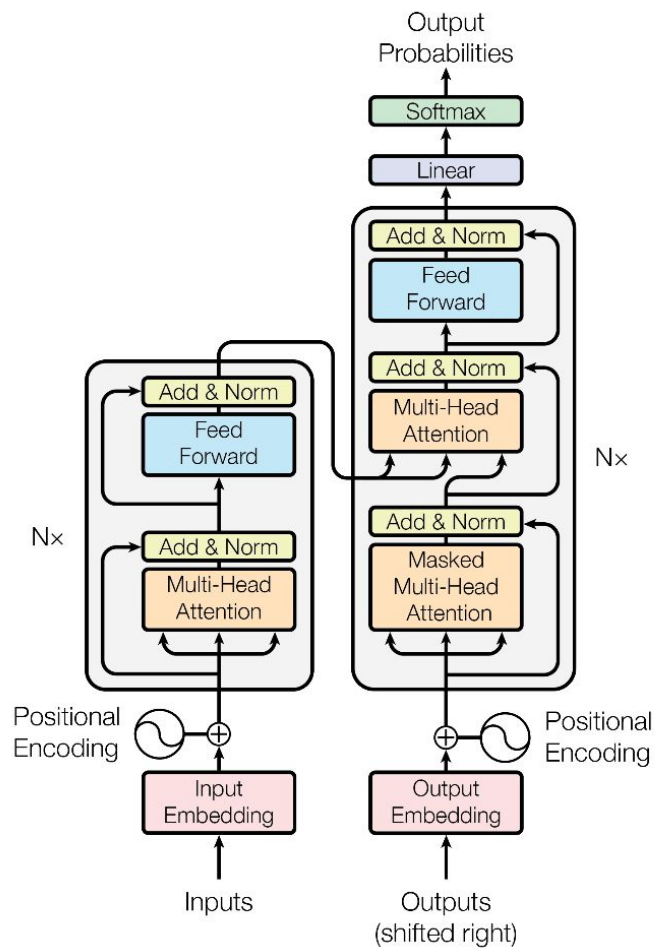
The FBI is chasing a criminal on the run .

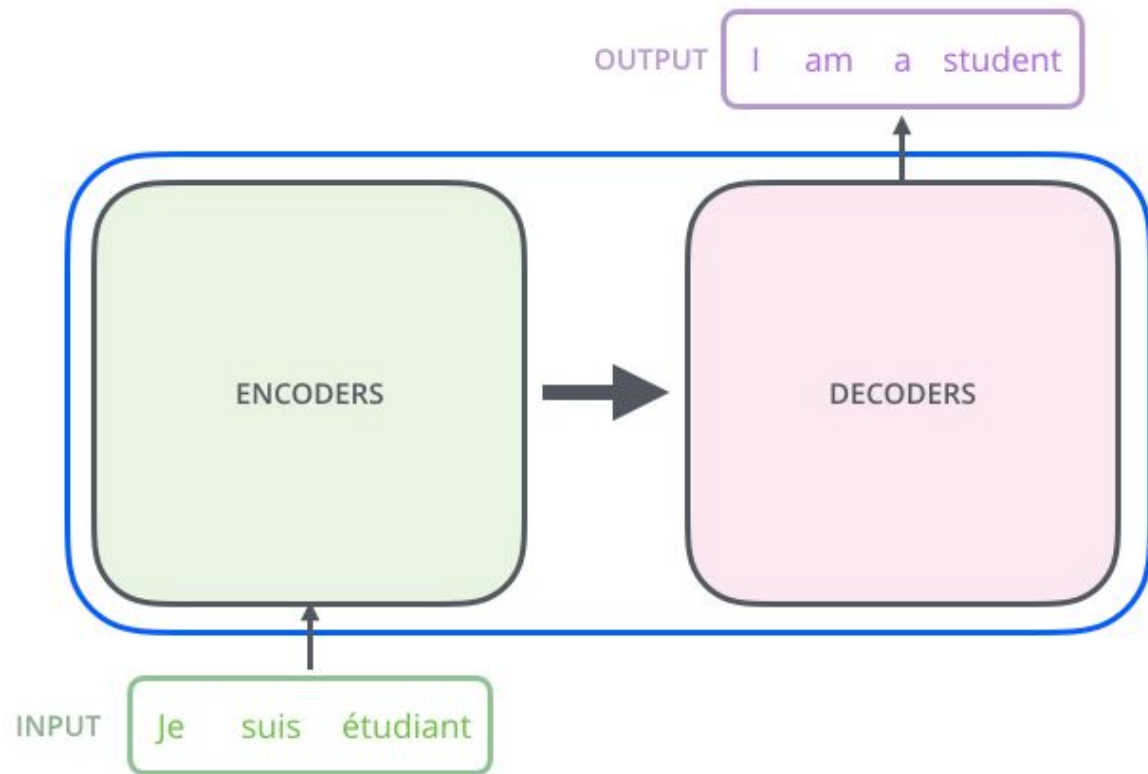
The FBI is chasing a criminal on the run .

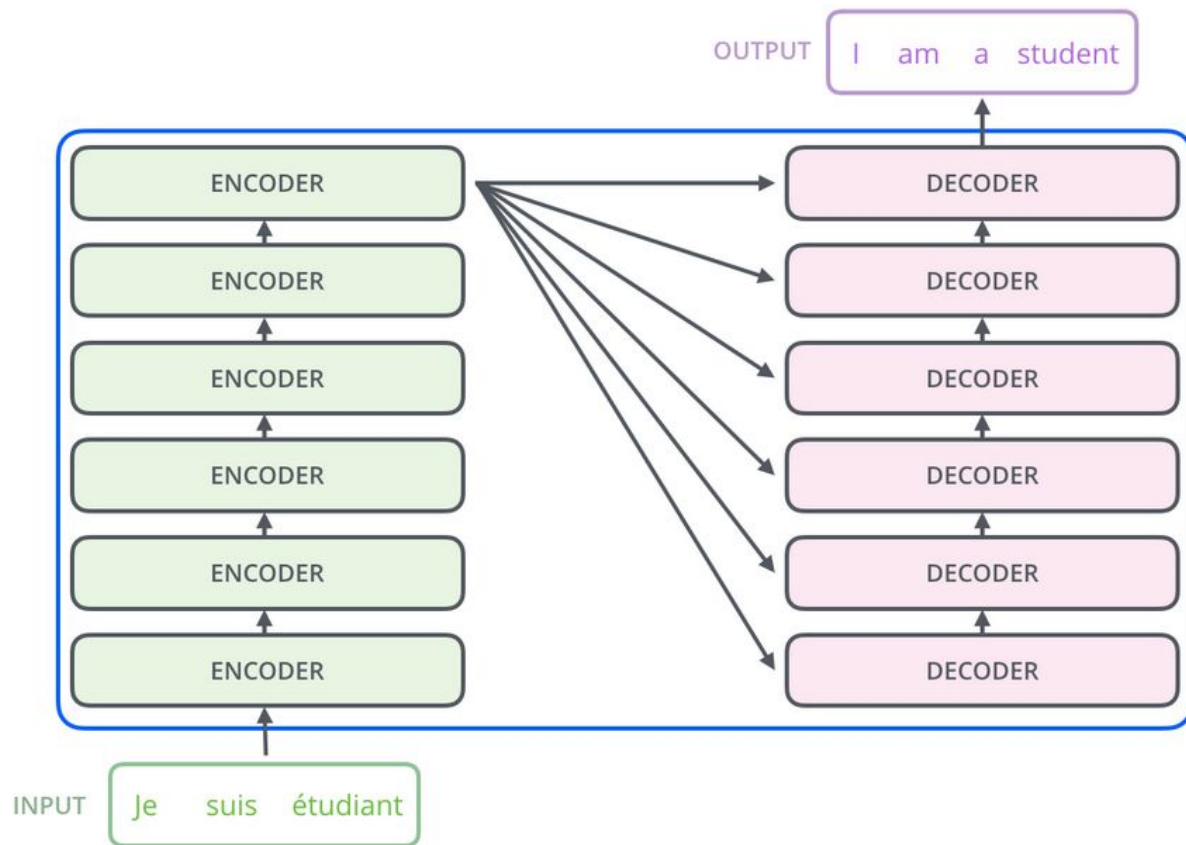
The FBI is chasing a criminal on the run .

The FBI is chasing a criminal on the run .

| Name | Alignment score function | Citation |
|------------------------|---|------------------------------|
| Content-base attention | $\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \text{cosine}[\mathbf{s}_t, \mathbf{h}_i]$ | Graves2014 |
| Additive(*) | $\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \mathbf{v}_a^\top \tanh(\mathbf{W}_a[\mathbf{s}_t; \mathbf{h}_i])$ | Bahdanau2015 |
| Location-Base | $\alpha_{t,i} = \text{softmax}(\mathbf{W}_a \mathbf{s}_t)$ Note: This simplifies the softmax alignment to only depend on the target position. | Luong2015 |
| General | $\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \mathbf{s}_t^\top \mathbf{W}_a \mathbf{h}_i$ where \mathbf{W}_a is a trainable weight matrix in the attention layer. | Luong2015 |
| Dot-Product | $\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \mathbf{s}_t^\top \mathbf{h}_i$ | Luong2015 |
| Scaled Dot-Product(^) | $\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \frac{\mathbf{s}_t^\top \mathbf{h}_i}{\sqrt{n}}$ Note: very similar to the dot-product attention except for a scaling factor; where n is the dimension of the source hidden state. | Vaswani2017 |

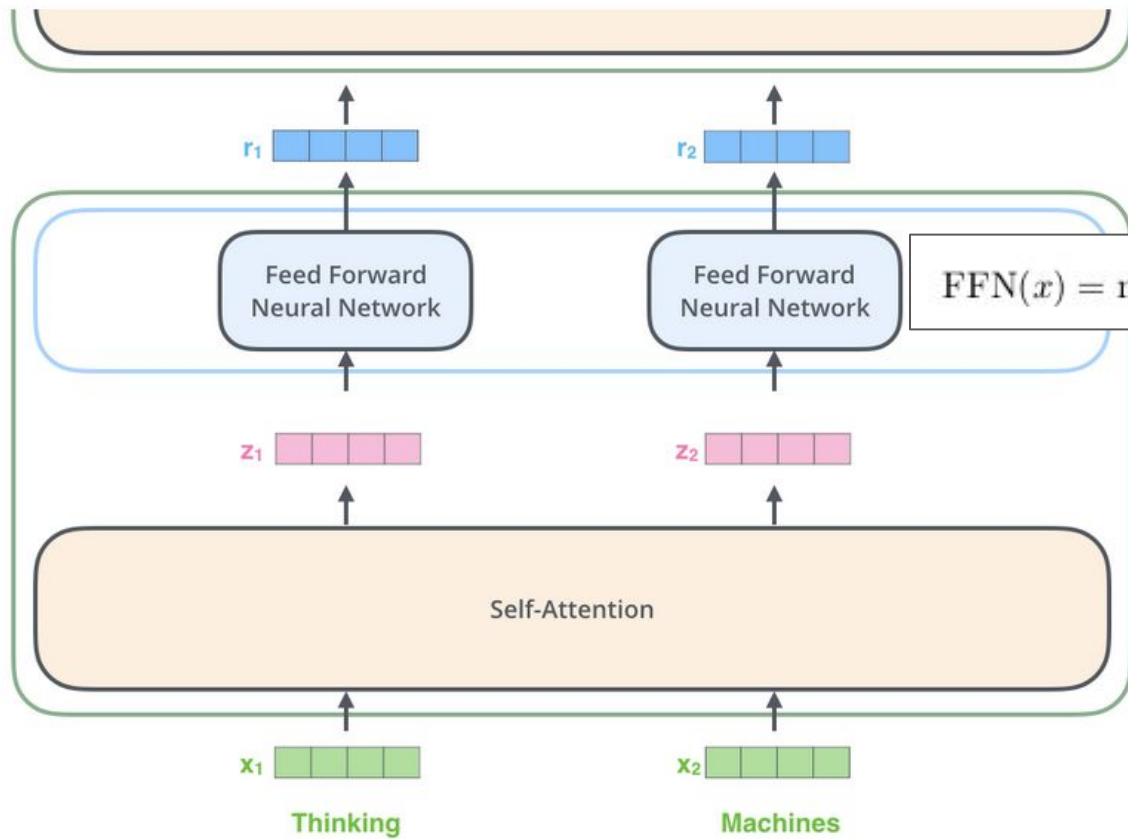






ENCODER #2

ENCODER #1



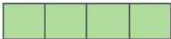
$$\text{FFN}(x) = \max(0, \mathbf{z}_k^T W_1 + b_1) W_2 + b_2$$

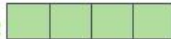
Input

Thinking


Machines

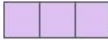
Embedding

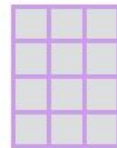
x_1 

x_2 

Queries


q_1 


q_2 

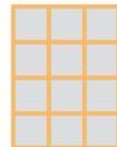


W^Q

Keys


k_1 

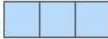
k_2 



W^K

Values

v_1 

v_2 



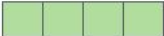
W^V

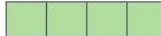
Input

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

Softmax

X

Value

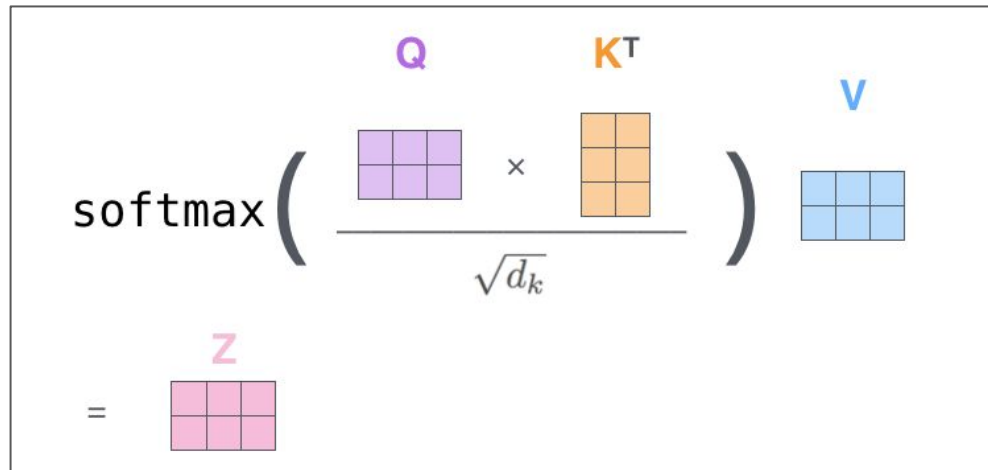
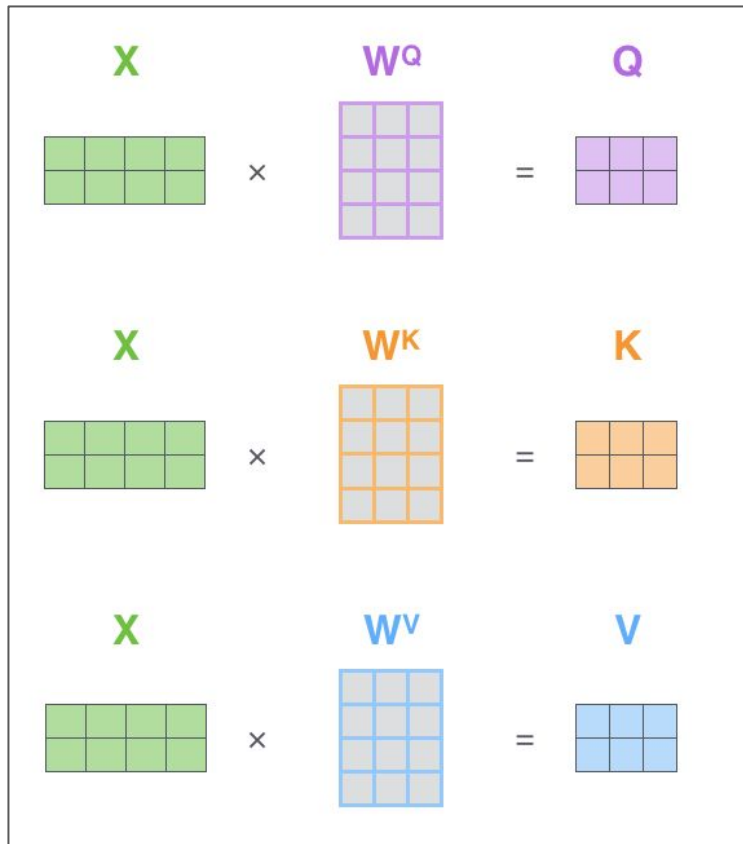
v_1 

v_2 

Sum

z_1 

z_2 



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

1) This is our input sentence*

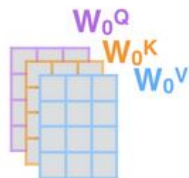
2) We embed each word*

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

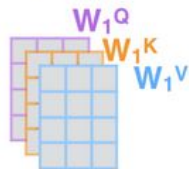
4) Calculate attention using the resulting $Q/K/V$ matrices

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

Thinking
Machines



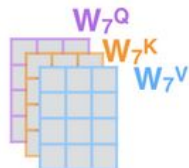
* In all encoders other than #0,
we don't need embedding.
We start directly with the output
of the encoder right below this one



...

...

...



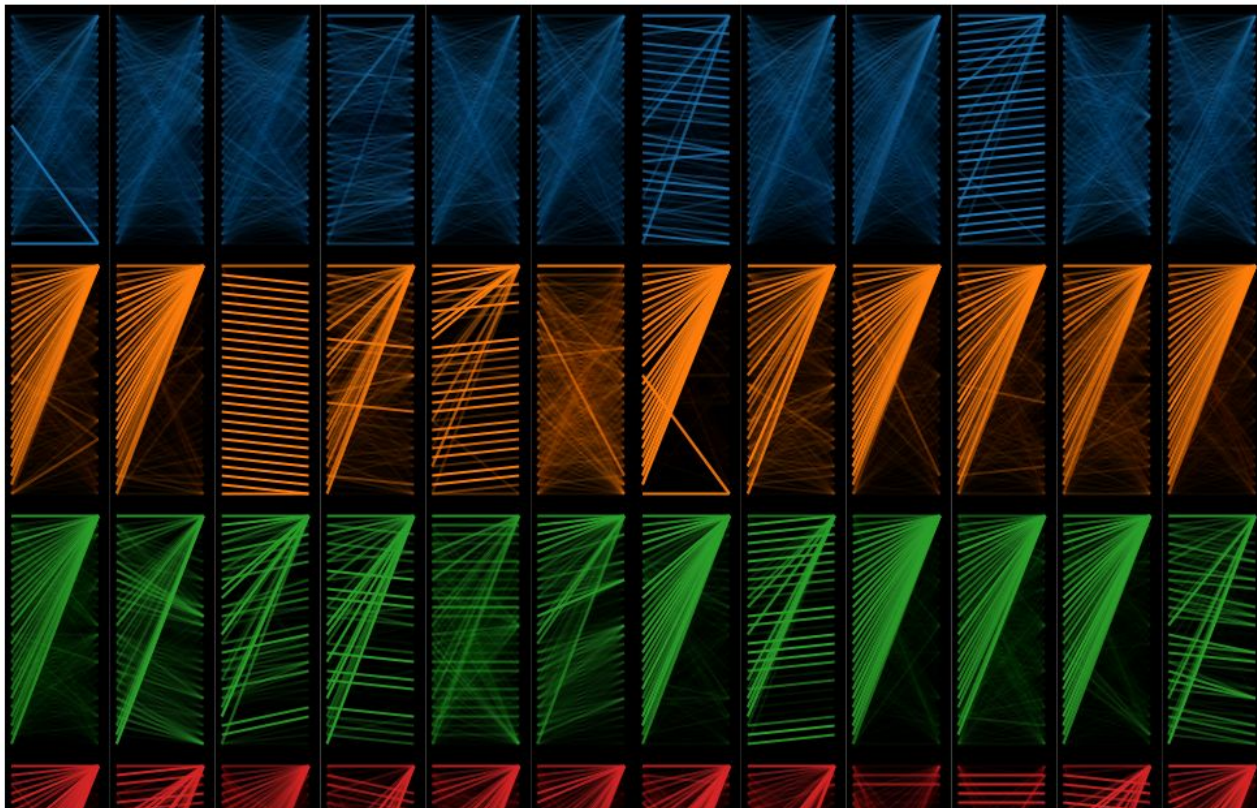
$$\text{MultiHead}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = [\text{head}_1; \dots; \text{head}_h] \mathbf{W}^O$$

$$\text{where head}_i = \text{Attention}(\mathbf{QW}_i^Q, \mathbf{KW}_i^K, \mathbf{VW}_i^V)$$

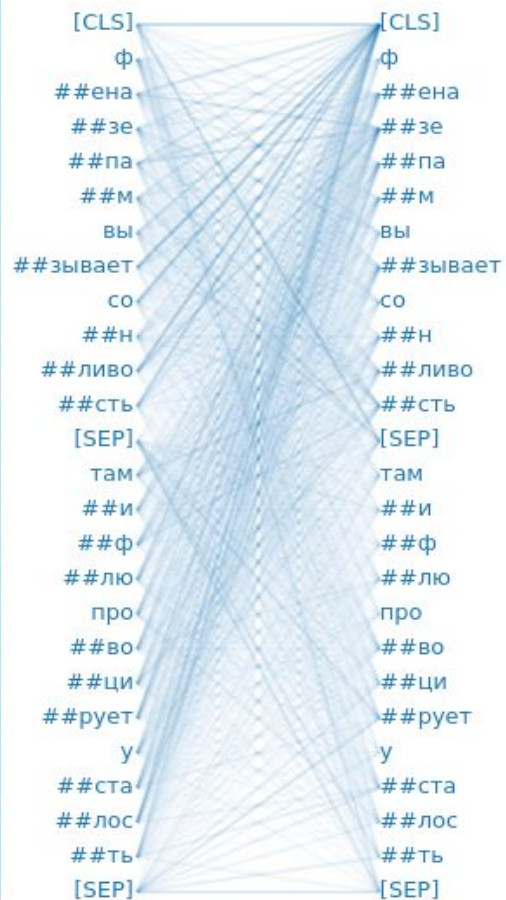
where \mathbf{W}_i^Q , \mathbf{W}_i^K , \mathbf{W}_i^V , and \mathbf{W}^O are parameter matrices to be learned.

```
from bertviz import model_view
```

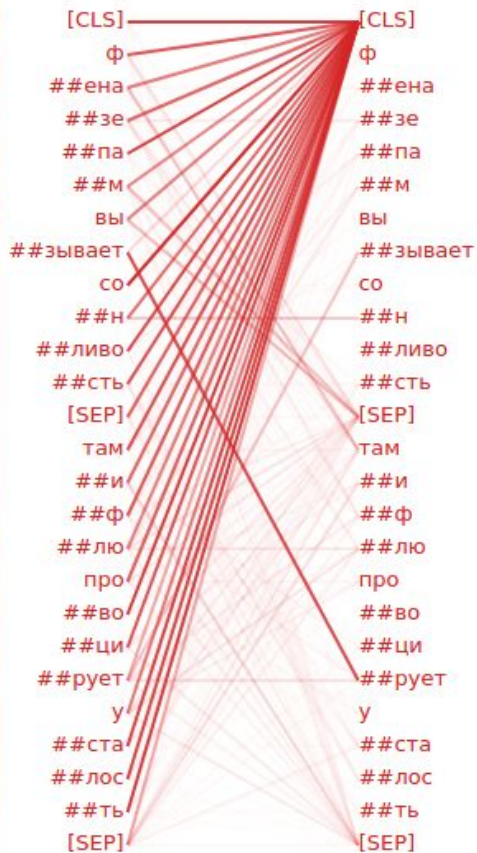
```
call_html()  
model_view(attention, tokens)
```



Layer 0, Head 5

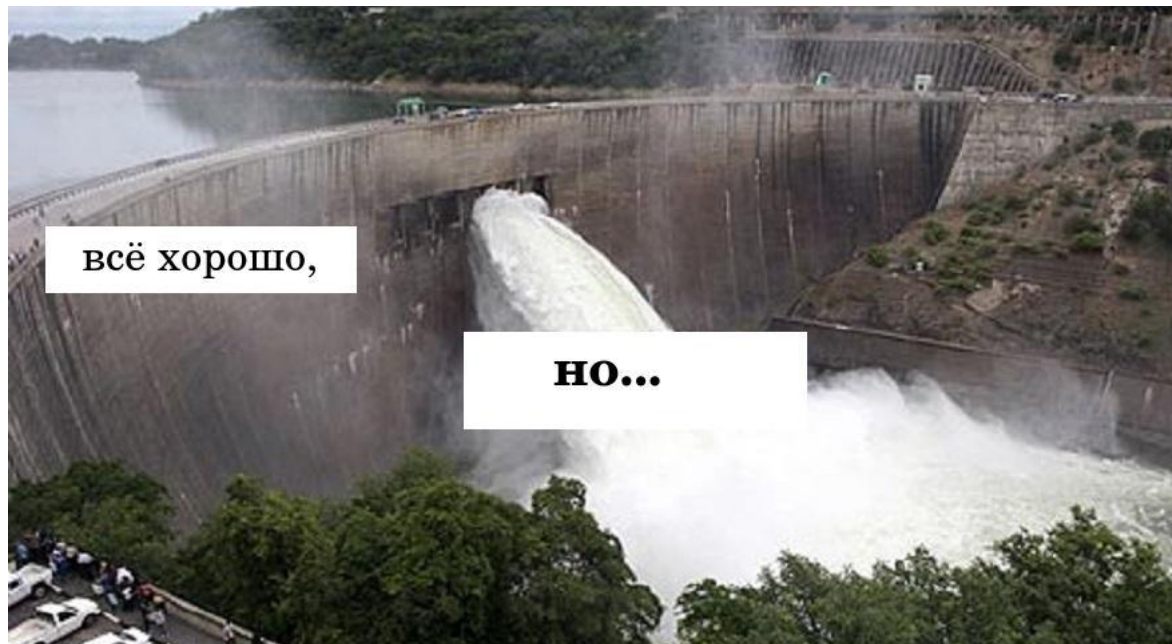


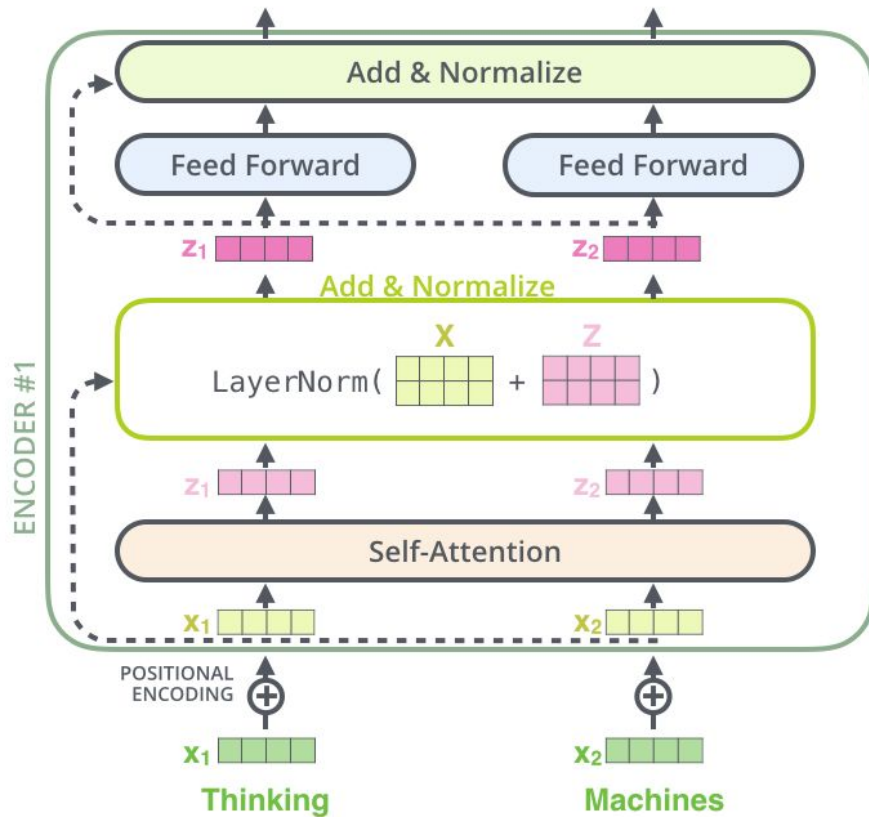
Layer 3, Head 0

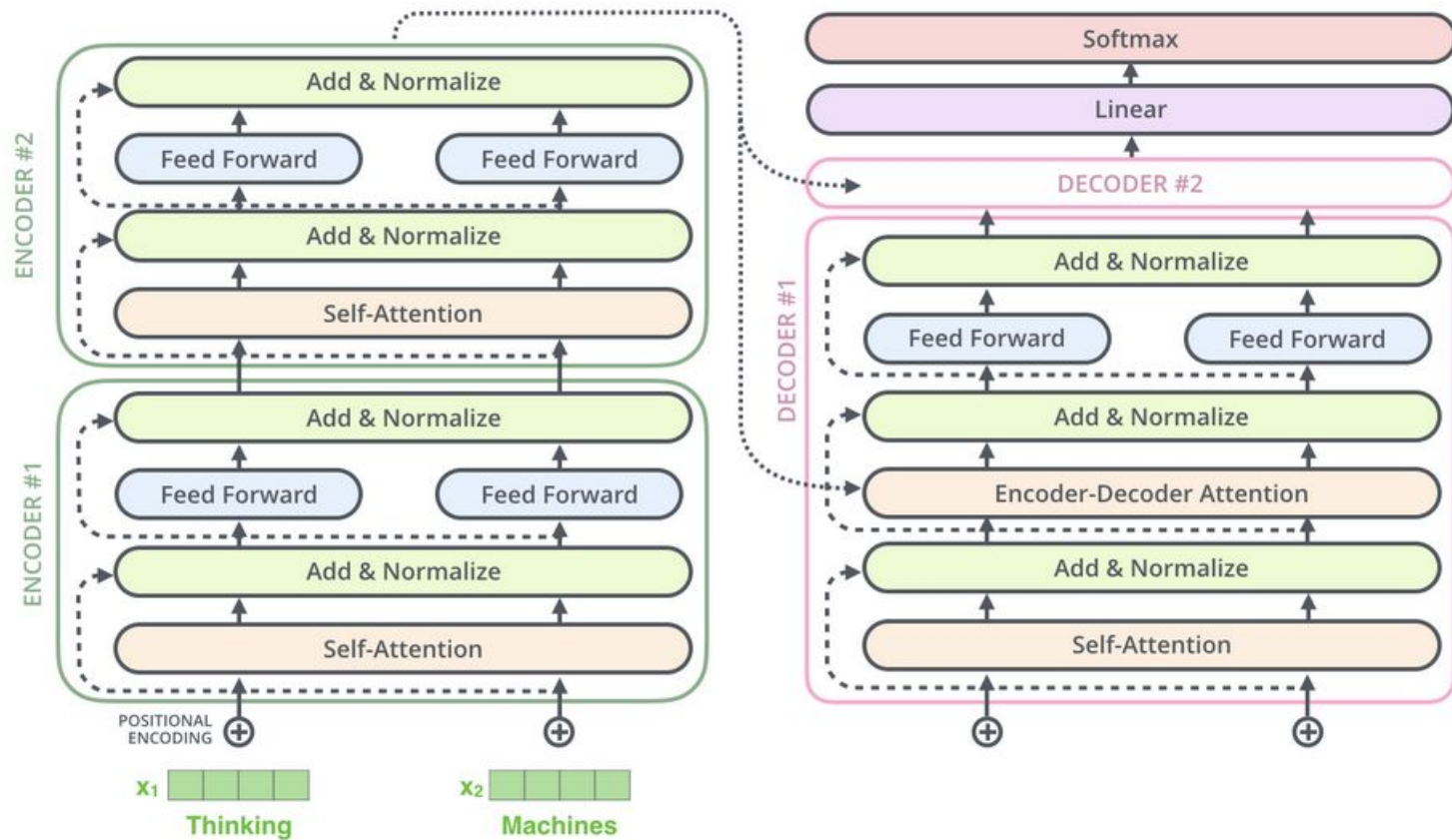


Ну всё, можно уже стакать слои?

Ну всё, можно уже стакать слои?

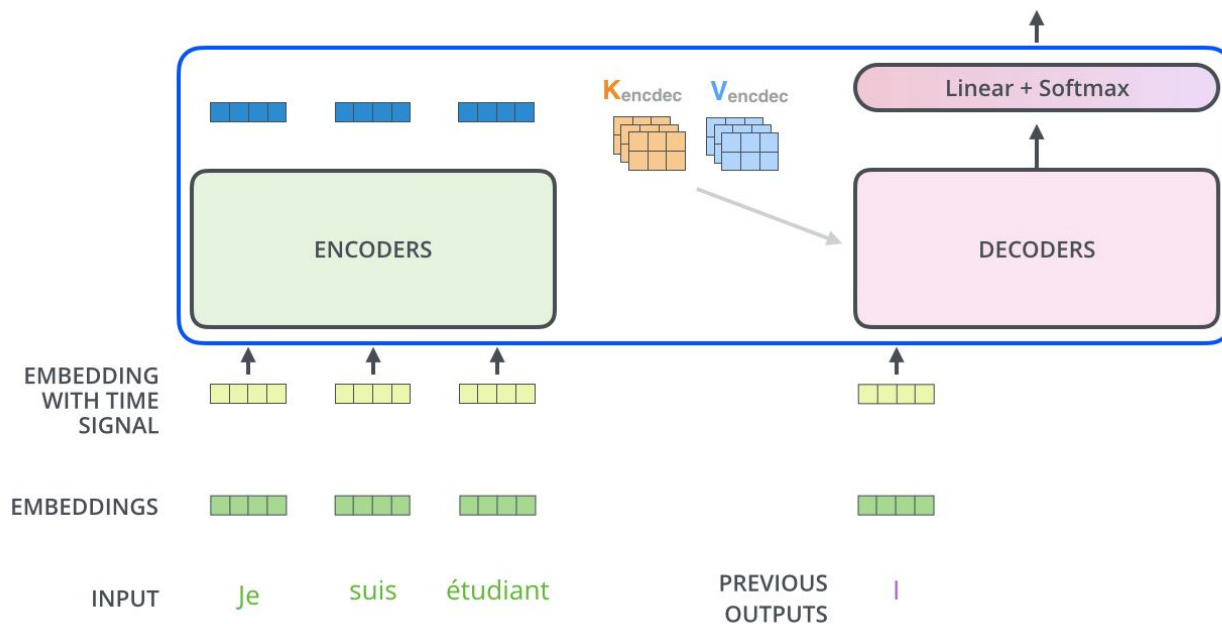






Decoding time step: 1 2 3 4 5 6

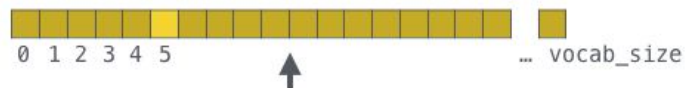
OUTPUT |



Which word in our vocabulary
is associated with this index?

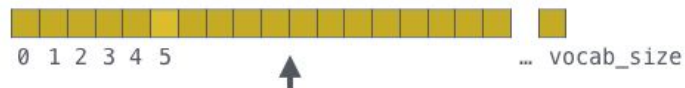
Get the index of the cell
with the highest value
(**argmax**)

log_probs



Softmax

logits



Linear

Decoder stack output



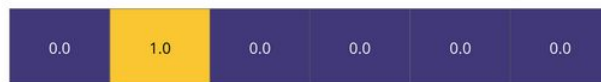
am

5

Output Vocabulary

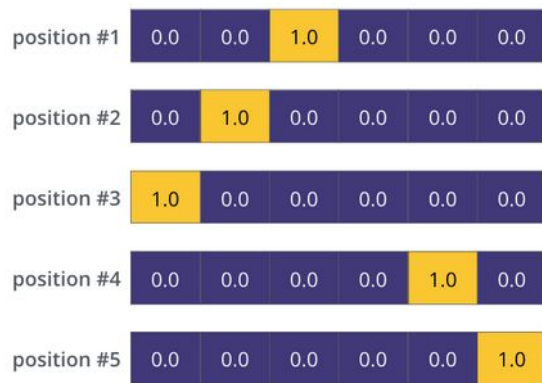
| WORD | a | am | I | thanks | student | <eos> |
|-------|---|----|---|--------|---------|-------|
| INDEX | 0 | 1 | 2 | 3 | 4 | 5 |

One-hot encoding of the word "am"



Target Model Outputs

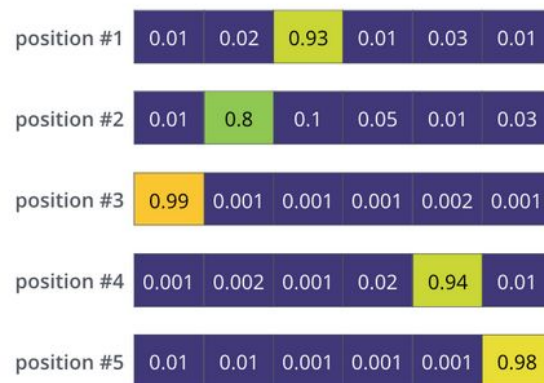
Output Vocabulary: a am I thanks student <eos>



a am I thanks student <eos>

Trained Model Outputs

Output Vocabulary: a am I thanks student <eos>

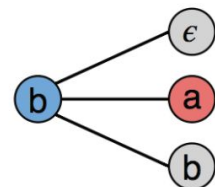
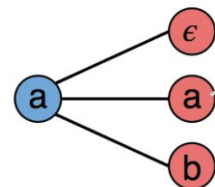
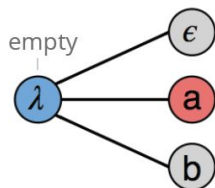


a am I thanks student <eos>

T = 2

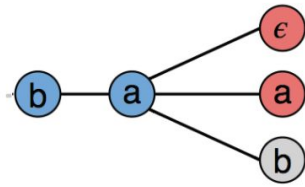
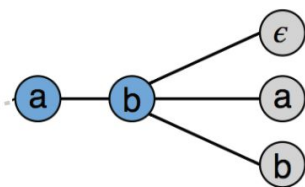
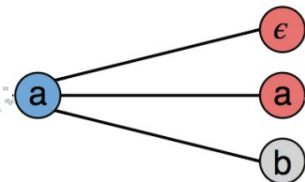
current hypotheses

proposed candidates

**T = 3**

current hypotheses

proposed candidates

**T = 4**

current hypotheses



CTC beam search algorithm with an output alphabet $\{\epsilon, a, b\}$ and a beam size of three.

Multiple extensions can merge in to the same hypothesis.

Positional Encoding

- Смотрим на весь ввод сразу - теряем информацию о порядке слов
- Нужна функция, сопоставляющая слову его относительную позицию во входной последовательности
- Наивные варианты имеют ряд недостатков

Positional Encoding

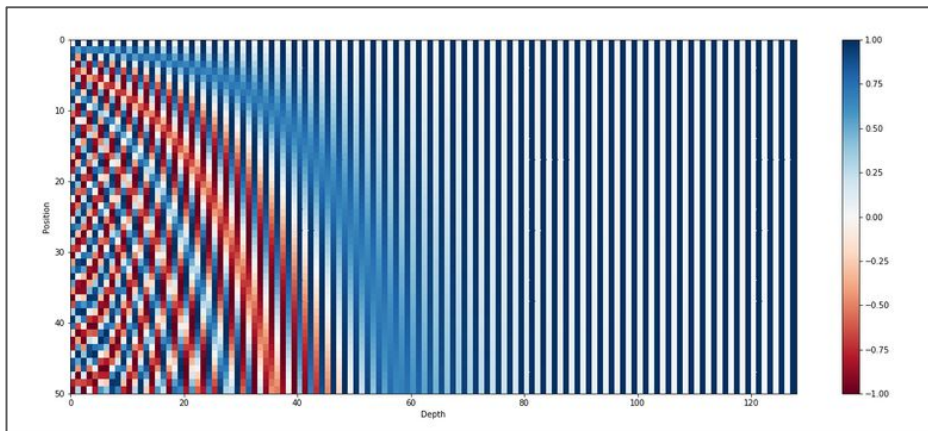
Что требуется от функции?

- Должна сопоставлять уникальную кодировку для каждого временного шага (позиции слова в предложении).
- Расстояние между любыми двумя временными шагами должно быть одинаковым в предложениях разной длины.
- Наша модель должна без каких-либо усилий обобщаться на более длинные предложения. Значения функции должны быть ограниченными.
- Отображение должно быть детерминированным.

Желаемая функция

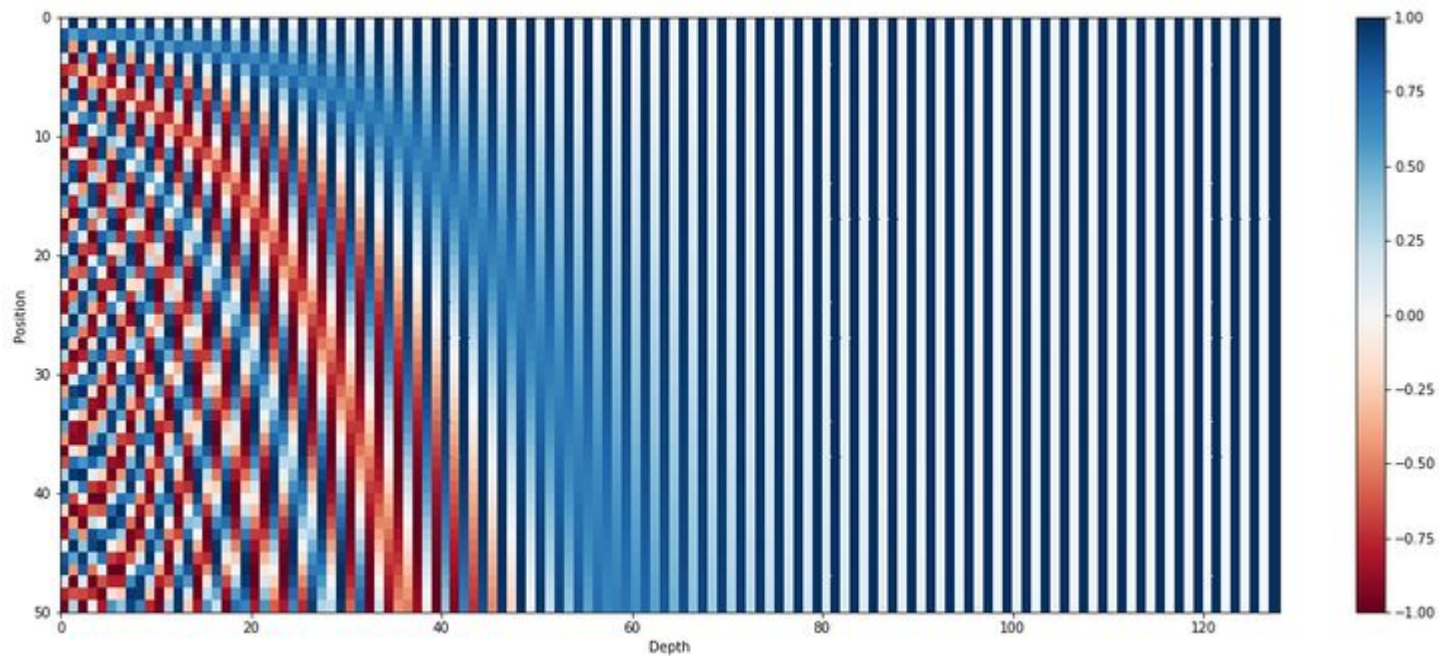
$$\vec{p}_t^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k \cdot t), & \text{if } i = 2k \\ \cos(\omega_k \cdot t), & \text{if } i = 2k + 1 \end{cases}$$

$$\omega_k = \frac{1}{10000^{2k/d}}$$

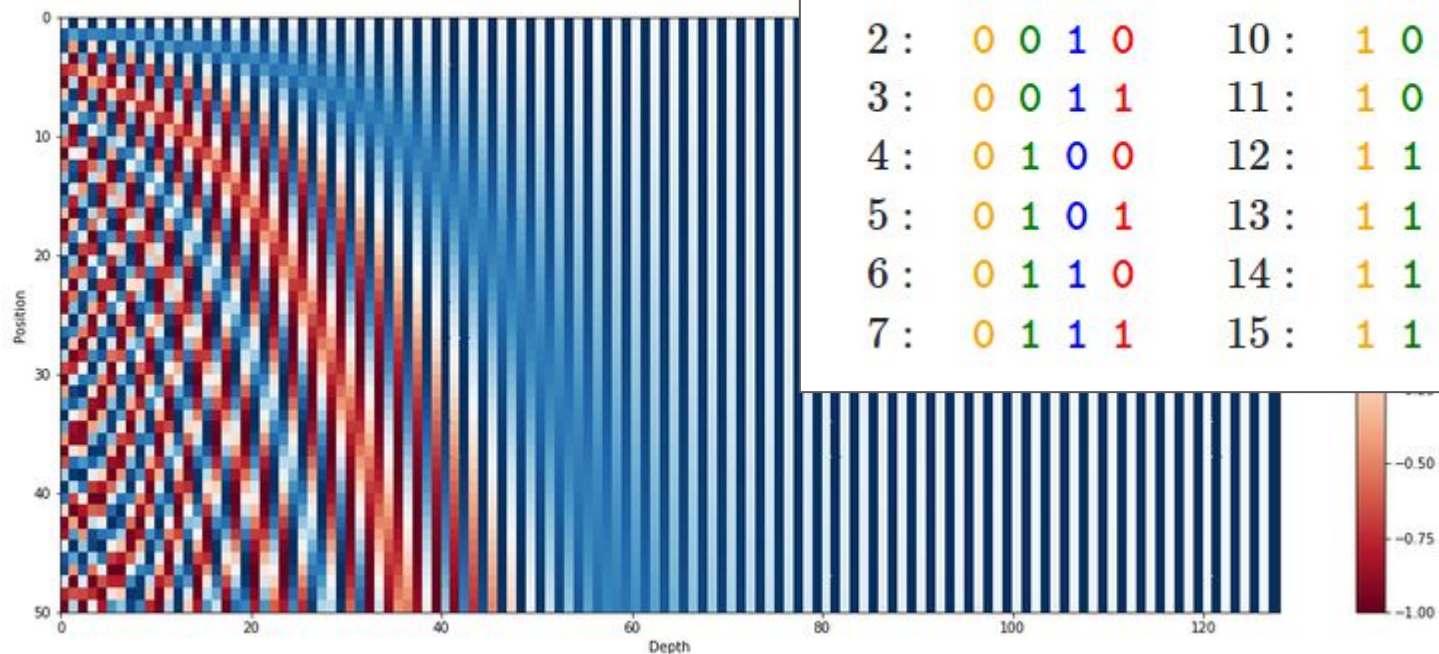


$$\vec{p}_t = \begin{bmatrix} \sin(\omega_1 \cdot t) \\ \cos(\omega_1 \cdot t) \\ \\ \sin(\omega_2 \cdot t) \\ \cos(\omega_2 \cdot t) \\ \\ \vdots \\ \\ \sin(\omega_{d/2} \cdot t) \\ \cos(\omega_{d/2} \cdot t) \end{bmatrix}_{d \times 1}$$

...И её визуализация



...И её визуализация



Label Smoothing + Residual Dropout

Formula of Label Smoothing

Label smoothing replaces one-hot encoded label vector y_{hot} with a mixture of y_{hot} and the uniform distribution:

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

where K is the number of label classes, and α is a hyperparameter that determines the amount of smoothing. If $\alpha = 0$, we obtain the original one-hot encoded y_{hot} . If $\alpha = 1$, we get the uniform distribution.

Заключение

Table 1: Maximum path lengths, per-layer complexity and minimum number of sequential operations for different layer types. n is the sequence length, d is the representation dimension, k is the kernel size of convolutions and r the size of the neighborhood in restricted self-attention.

| Layer Type | Complexity per Layer | Sequential Operations | Maximum Path Length |
|-----------------------------|--------------------------|-----------------------|---------------------|
| Self-Attention | $O(n^2 \cdot d)$ | $O(1)$ | $O(1)$ |
| Recurrent | $O(n \cdot d^2)$ | $O(n)$ | $O(n)$ |
| Convolutional | $O(k \cdot n \cdot d^2)$ | $O(1)$ | $O(\log_k(n))$ |
| Self-Attention (restricted) | $O(r \cdot n \cdot d)$ | $O(1)$ | $O(n/r)$ |

Заклучение

Table 2: The Transformer achieves better BLEU scores than previous state-of-the-art models on the English-to-German and English-to-French newstest2014 tests at a fraction of the training cost.

| Model | BLEU | | Training Cost (FLOPs) | |
|---------------------------------|-------------|--------------|---------------------------------------|---------------------|
| | EN-DE | EN-FR | EN-DE | EN-FR |
| ByteNet [15] | 23.75 | | | |
| Deep-Att + PosUnk [32] | | 39.2 | | $1.0 \cdot 10^{20}$ |
| GNMT + RL [31] | 24.6 | 39.92 | $2.3 \cdot 10^{19}$ | $1.4 \cdot 10^{20}$ |
| ConvS2S [8] | 25.16 | 40.46 | $9.6 \cdot 10^{18}$ | $1.5 \cdot 10^{20}$ |
| MoE [26] | 26.03 | 40.56 | $2.0 \cdot 10^{19}$ | $1.2 \cdot 10^{20}$ |
| Deep-Att + PosUnk Ensemble [32] | | 40.4 | | $8.0 \cdot 10^{20}$ |
| GNMT + RL Ensemble [31] | 26.30 | 41.16 | $1.8 \cdot 10^{20}$ | $1.1 \cdot 10^{21}$ |
| ConvS2S Ensemble [8] | 26.36 | 41.29 | $7.7 \cdot 10^{19}$ | $1.2 \cdot 10^{21}$ |
| Transformer (base model) | 27.3 | 38.1 | $3.3 \cdot 10^{18}$ | |
| Transformer (big) | 28.4 | 41.0 | $2.3 \cdot 10^{19}$ | |

Ссылки на материалы

Вот они, слева направо!

<https://arxiv.org/abs/1706.03762>

<https://www.youtube.com/watch?v=dQw4w9WgXcQ>

<https://lilianweng.github.io/lil-log/2018/06/24/attention-attention.html#whats-wrong-with-seq2seq-model>

<http://jalammar.github.io/illustrated-transformer/>

https://kazemnejad.com/blog/transformer_architecture_positional_encoding/#what-is-positional-encoding-and-why-do-we-need-it-in-the-first-place

Thank you for your Attention ;)