

# Transformer

Sergey Ivanov (617)

*qbrick@mail.ru*

October 28, 2019

- 1 Recap
- 2 Attention
- 3 Multi-Head Attention
- 4 Attention is All You Need
- 5 BERT, GPT-2, ...

# Section 1

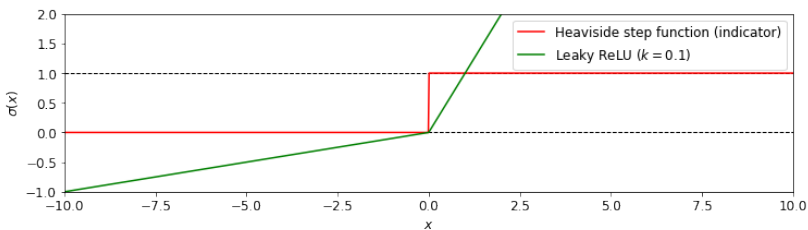
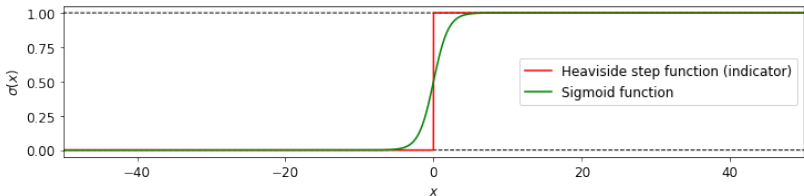
## Recap

# Neural Networks

$$y(x) = \sigma(Wx + b)$$

# Neural Networks

$$y(x) = \sigma(Wx + b)$$



# Softmax

Classification problem head:

$$\text{softmax}(x)_i \propto e^{x_i}$$

# Softmax

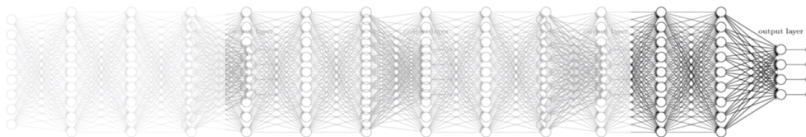
## Classification problem head:

$$\text{softmax}(x)_i \propto e^{x_i}$$

	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$
$x$	<b>-0.3</b>	<b>0</b>	<b>0.01</b>	<b>1</b>	<b>-0.1</b>
$\text{softmax}(x)$	0.116	0.157	0.158	0.426	0.142
$\text{softmax}(3x)$	0.017	0.043	0.044	0.863	0.032
$\text{softmax}(10x)$	2e-06	5e-5	5e-5	0.99999	2e-05

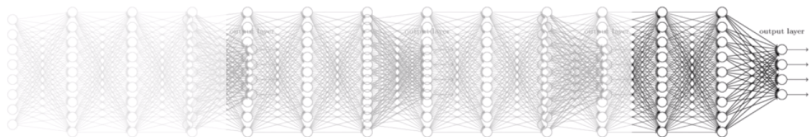
! scaling factor often influences the smoothness of approximations!

# Vanishing gradients problem

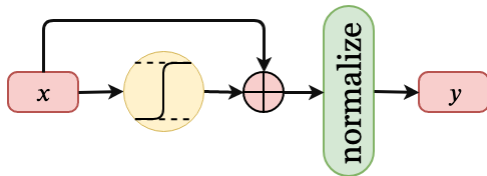




# Vanishing gradients problem

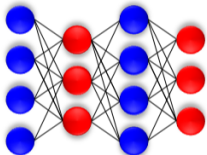


Trying to avoid the problem:

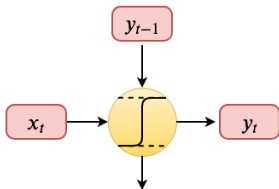
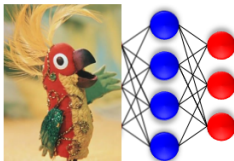


# Architectures

## FULLY-CONNECTED



## CONVOLUTIONAL

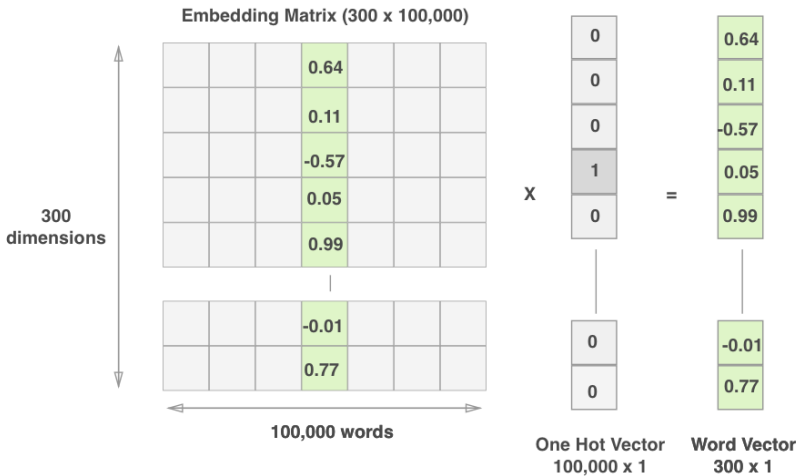


## RECURRENT

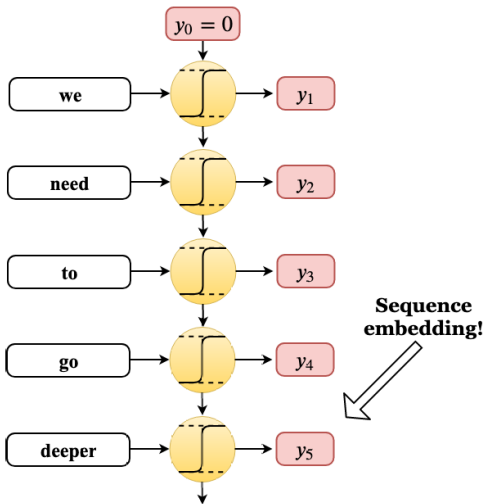


## TRANSFORMER

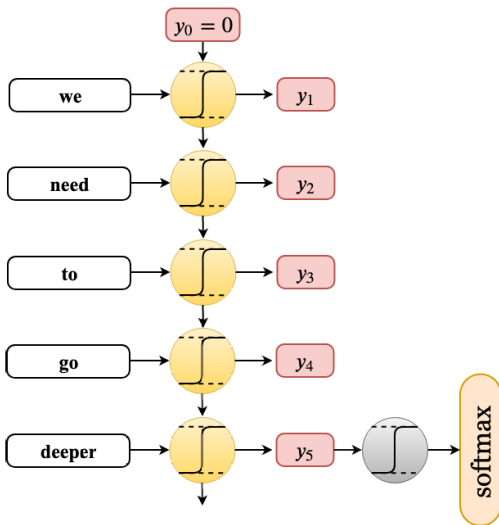
# Word embeddings



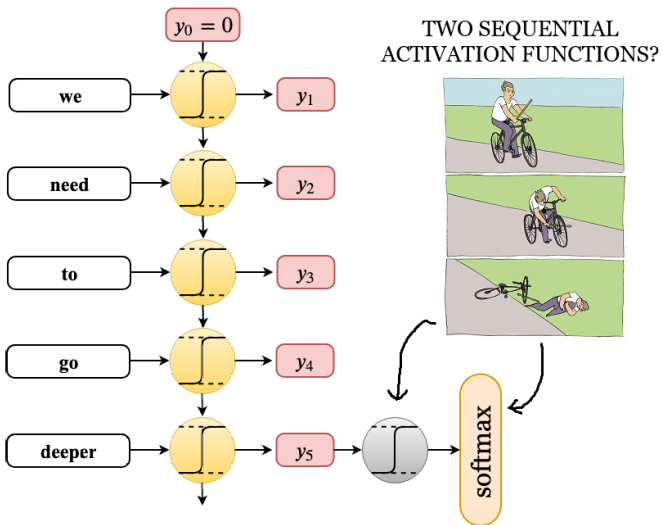
# Sequence Embedding



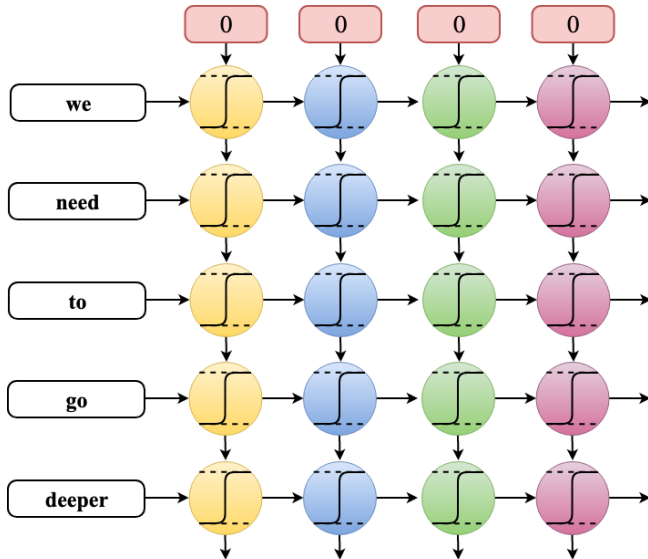
# Sequence Embedding



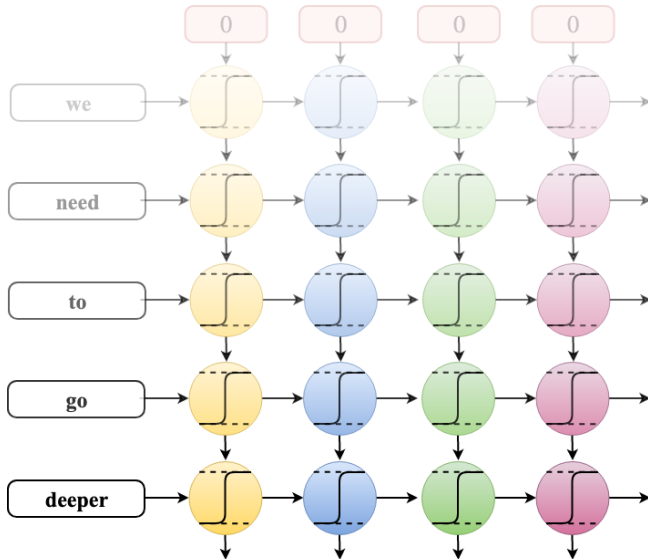
# Sequence Embedding



# Sequence Embedding

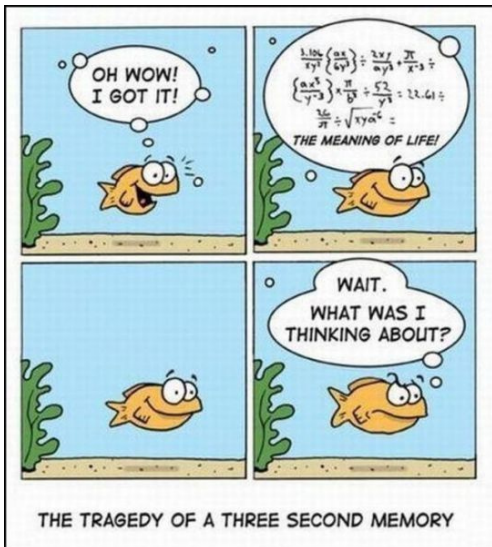


# Sequence Embedding

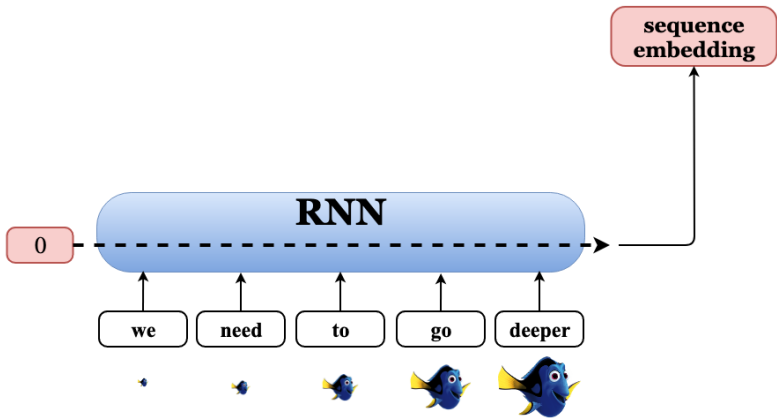




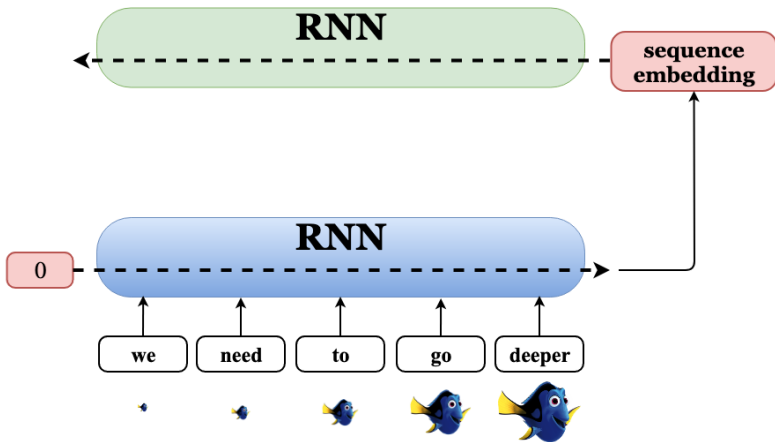
# No memory?



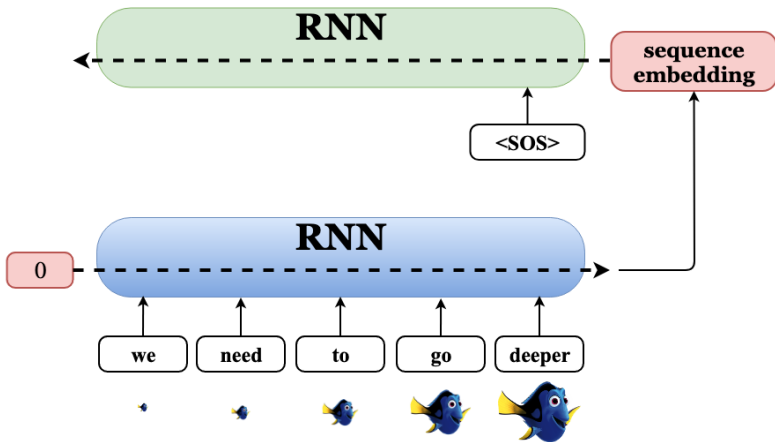
# Machine Translation



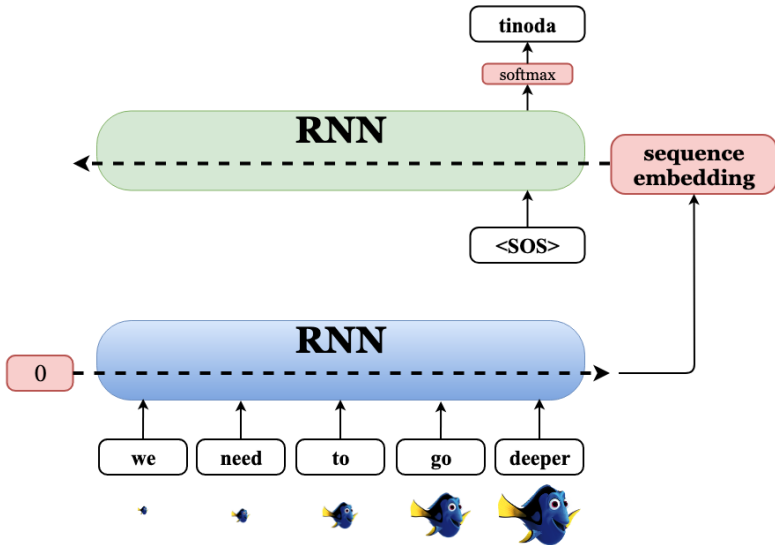
# Machine Translation



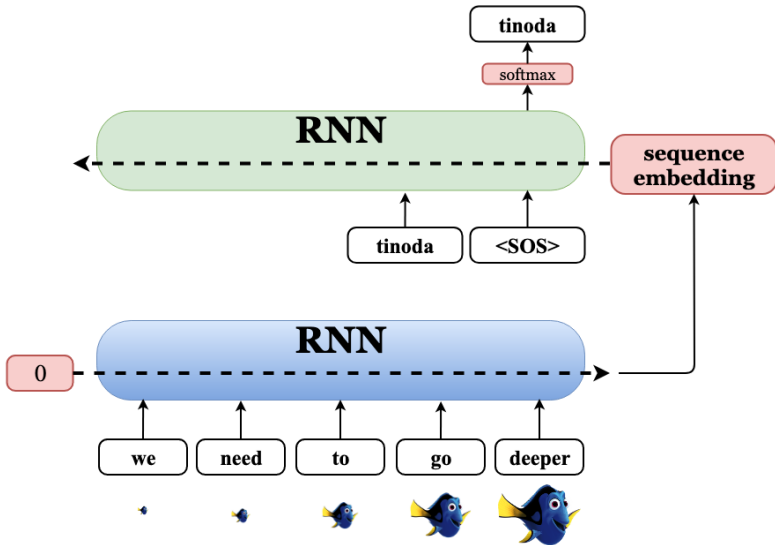
# Machine Translation



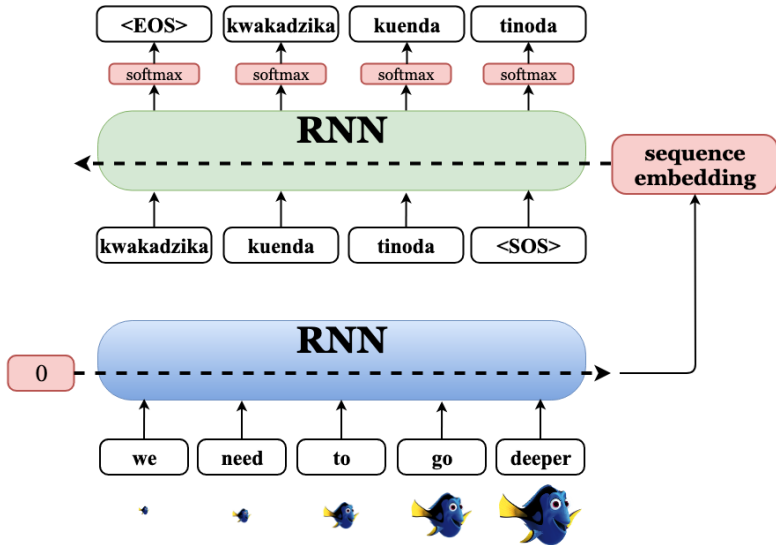
# Machine Translation



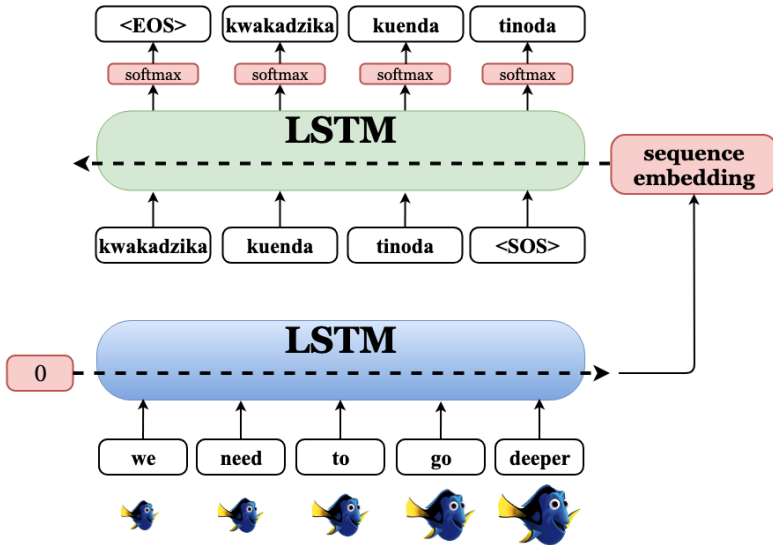
# Machine Translation



# Machine Translation

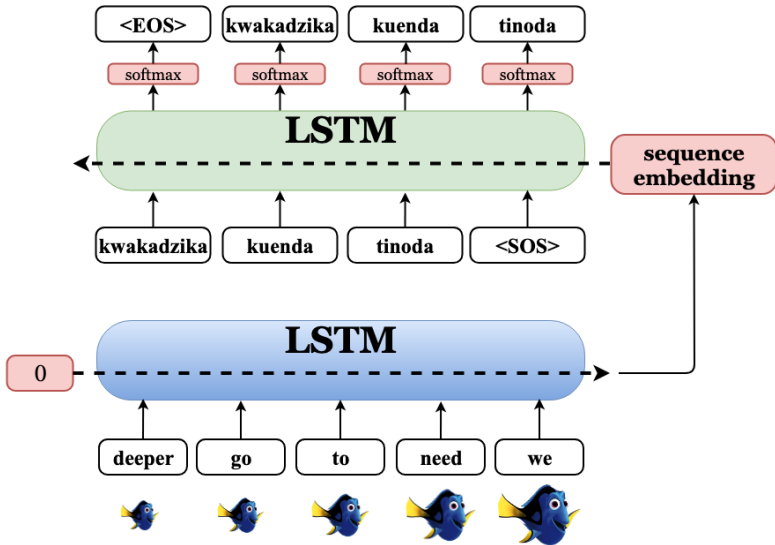


# Machine Translation

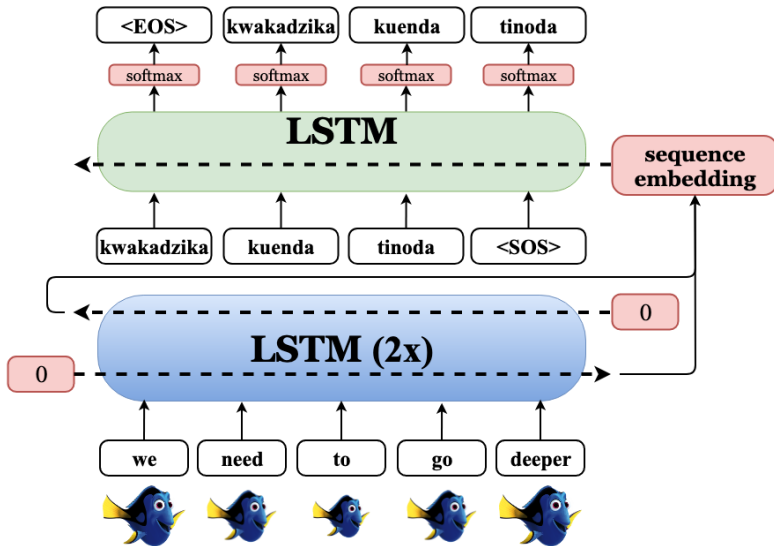




# Machine Translation



# Machine Translation



## Section 2

# Attention

# Remembering dict()

Suppose you have a Python dictionary:

```
d = dict(zip(K, V))
```

where

- $K \in \mathbb{R}^{n \times k}$  —  $n$  **keys**,  $K_i \in \mathbb{R}^k$
- $V \in \mathbb{R}^{n \times v}$  —  $n$  **values**,  $V_i \in \mathbb{R}^v$

# Remembering dict()

Suppose you have a Python dictionary:

```
d = dict(zip(K, V))
```

where

- $K \in \mathbb{R}^{n \times k}$  —  $n$  **keys**,  $K_i \in \mathbb{R}^k$
- $V \in \mathbb{R}^{n \times v}$  —  $n$  **values**,  $V_i \in \mathbb{R}^v$

Suppose you have **query**  $q \in \mathbb{R}^k$ .  $d[q]$  — ?

# Remembering dict()

Suppose you have a Python dictionary:

```
d = dict(zip(K, V))
```

where

- $K \in \mathbb{R}^{n \times k}$  —  $n$  **keys**,  $K_i \in \mathbb{R}^k$
- $V \in \mathbb{R}^{n \times v}$  —  $n$  **values**,  $V_i \in \mathbb{R}^v$

Suppose you have **query**  $q \in \mathbb{R}^k$ .  $d[q]$  — ?

Weeeell, find  $i: K_i = q$ , then the answer is  $V_i$ .

# Approximating dict()

$$K \in \mathbb{R}^{n \times k} \quad V \in \mathbb{R}^{n \times v} \quad q \in \mathbb{R}^k$$

## Solution

- 1  $w_i = \mathbb{I}[K_i = q]$
- 2  $d[q] = \sum_i^n w_i V_i$

# Approximating dict()

$$K \in \mathbb{R}^{n \times k} \quad V \in \mathbb{R}^{n \times v} \quad q \in \mathbb{R}^k$$

## Solution

- 1  $w_i = \mathbb{I}[K_i = q]$  — discrete!
- 2  $d[q] = \sum_i^n w_i V_i$



# Approximating dict()

$$K \in \mathbb{R}^{n \times k} \quad V \in \mathbb{R}^{n \times v} \quad q \in \mathbb{R}^k$$

## Solution

- 1  $w_i = \mathbb{I}[K_i = q]$  — discrete!
- 2  $d[q] = \sum_i^n w_i V_i$

Let  $\rho(K_i, q) \in \mathbb{R}$  be some measure of similarity (*compatibility function*) between  $K_i, q$ .

# Approximating dict()

$$K \in \mathbb{R}^{n \times k} \quad V \in \mathbb{R}^{n \times v} \quad q \in \mathbb{R}^k$$

## Solution

- 1  $w = \operatorname{argmax}(a), \quad a_i = \rho(K_i, q)$
- 2  $d[q] = \sum_i^n w_i V_i$

Let  $\rho(K_i, q) \in \mathbb{R}$  be some measure of similarity (*compatibility function*) between  $K_i, q$ .

# Approximating dict()

$$K \in \mathbb{R}^{n \times k} \quad V \in \mathbb{R}^{n \times v} \quad q \in \mathbb{R}^k$$

## Solution

- 1  $w = \text{softmax}(a), \quad a_i = \rho(K_i, q)$
- 2  $d[q] = \sum_i^n w_i V_i$

Let  $\rho(K_i, q) \in \mathbb{R}$  be some measure of similarity (*compatibility function*) between  $K_i, q$ .

# Approximating dict()

$$K \in \mathbb{R}^{n \times k} \quad V \in \mathbb{R}^{n \times v} \quad q \in \mathbb{R}^k$$

## Solution

- 1  $w = \text{softmax}(a), \quad a_i = \rho(K_i, q)$
- 2  $d[q] = \sum_i^n w_i V_i$

Let  $\rho(K_i, q) \in \mathbb{R}$  be some measure of similarity (*compatibility function*) between  $K_i, q$ .

**Common choice:**  $\rho(K_i, q) = \langle K_i, q \rangle$

# Approximating dict()

$$K \in \mathbb{R}^{n \times k} \quad V \in \mathbb{R}^{n \times v} \quad q \in \mathbb{R}^k$$

## Solution

- 1  $w = \text{softmax}(a), \quad a = Kq$
- 2  $d[q] = \sum_i^n w_i V_i$

Let  $\rho(K_i, q) \in \mathbb{R}$  be some measure of similarity (*compatibility function*) between  $K_i, q$ .

**Common choice:**  $\rho(K_i, q) = \langle K_i, q \rangle$

# Approximating dict()

$$K \in \mathbb{R}^{n \times k} \quad V \in \mathbb{R}^{n \times v} \quad q \in \mathbb{R}^k$$

## Solution

- 1  $w = \text{softmax}(Kq)$
- 2  $d[q] = \sum_i^n w_i V_i$

Let  $\rho(K_i, q) \in \mathbb{R}$  be some measure of similarity (*compatibility function*) between  $K_i, q$ .

**Common choice:**  $\rho(K_i, q) = \langle K_i, q \rangle$

# Approximating dict()

$$K \in \mathbb{R}^{n \times k} \quad V \in \mathbb{R}^{n \times v} \quad q \in \mathbb{R}^k$$

## Solution

- 1  $w = \text{softmax}(Kq)$
- 2  $d[q] = w^T V$

Let  $\rho(K_i, q) \in \mathbb{R}$  be some measure of similarity (*compatibility function*) between  $K_i, q$ .

**Common choice:**  $\rho(K_i, q) = \langle K_i, q \rangle$

# Scalar Product Normalization

Assume that all  $q, K_i \sim \mathcal{N}(0, I_{k \times k})$ .



# Scalar Product Normalization

Assume that all  $q, K_i \sim \mathcal{N}(0, I_{k \times k})$ .

What will be distribution of  $\langle q, K_i \rangle = \sum_j q_j K_{ij} \sim ?!$

# Scalar Product Normalization

Assume that all  $q, K_i \sim \mathcal{N}(0, I_{k \times k})$ .

What will be distribution of  $\langle q, K_i \rangle = \sum_j q_j K_{ij} \sim ?!$

$$\mathbb{E} \sum_j^k q_j K_{ij} = \sum_j^k \mathbb{E} q_j K_{ij} = \{\text{independence}\} = \sum_j^k \mathbb{E} q_j \mathbb{E} K_{ij} = \sum_j^k 0 = 0$$

# Scalar Product Normalization

Assume that all  $q, K_i \sim \mathcal{N}(0, I_{k \times k})$ .

What will be distribution of  $\langle q, K_i \rangle = \sum_j q_j K_{ij} \sim ?!$

$$\mathbb{E} \sum_j^k q_j K_{ij} = \sum_j^k \mathbb{E} q_j K_{ij} = \{\text{independence}\} = \sum_j^k \mathbb{E} q_j \mathbb{E} K_{ij} = \sum_j^k 0 = 0$$

$$\begin{aligned} \mathbb{D} \sum_j^k q_j K_{ij} &= \{\text{independence}\} = \sum_j^k \mathbb{D} q_j K_{ij} = \\ &= \{\text{independence, zero expectation}\} = \sum_j^k \mathbb{D} q_j \mathbb{D} K_{ij} = \sum_j^k 1 = k \end{aligned}$$

# Scalar Product Normalization

Assume that all  $q, K_i \sim \mathcal{N}(0, I_{k \times k})$ .

What will be distribution of  $\langle q, K_i \rangle = \sum_j q_j K_{ij} \sim ?!$

$$\mathbb{E} \sum^k q_j K_{ij} = \sum^k \mathbb{E} q_j K_{ij} = \{\text{independence}\} = \sum^k \mathbb{E} q_j \mathbb{E} K_{ij} = \sum^k 0 = 0$$

$$\begin{aligned} \mathbb{D} \sum^k q_j K_{ij} &= \{\text{independence}\} = \sum^k \mathbb{D} q_j K_{ij} = \\ &= \{\text{independence, zero expectation}\} = \sum^k \mathbb{D} q_j \mathbb{D} K_{ij} = \sum^k 1 = k \end{aligned}$$

**Similarity metric normalisation:**

$$\rho(K_i, q) = \frac{\langle K_i, q \rangle}{\sqrt{k}}$$

# Scalar Product Normalization

Assume that all  $q, K_i \sim \mathcal{N}(0, I_{k \times k})$ .

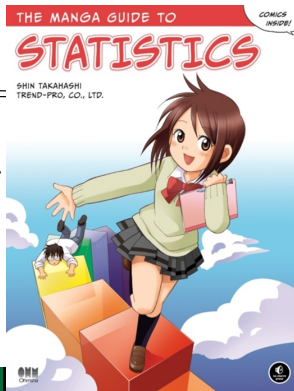
What will be distribution of  $\langle q, K_i \rangle = \sum_j q_j K_{ij} \sim ?!$

$$\mathbb{E} \sum_j^k q_j K_{ij} = \sum_j^k \mathbb{E} q_j K_{ij} = \{\text{independence}\} = \sum_j^k \mathbb{E} q_j \mathbb{E} K_{ij} = \sum_j^k 0 = 0$$

$$\begin{aligned} \mathbb{D} \sum_j^k q_j K_{ij} &= \{\text{independence}\} = \sum_j^k \mathbb{D} q_j K_{ij} = \\ &= \{\text{independence, zero expectation}\} = \sum_j^k \mathbb{D} q_j \mathbb{D} K_{ij} \end{aligned}$$

Similarity metric normalisation:

$$\rho(K_i, q) = \frac{\langle K_i, q \rangle}{\sqrt{k}}$$



# Attention

input:  $K \in \mathbb{R}^{n \times k}$ ,  $V \in \mathbb{R}^{n \times v}$ ,  $Q \in \mathbb{R}^{b \times k}$

# Attention

input:  $K \in \mathbb{R}^{n \times k}$ ,  $V \in \mathbb{R}^{n \times v}$ ,  $Q \in \mathbb{R}^{b \times k}$

Get attention weights matrix:

$$W = \text{softmax} \left( \frac{QK^T}{\sqrt{k}}, \text{dim}=1 \right) \in \mathbb{R}^{b \times n}$$

# Attention

**input:**  $K \in \mathbb{R}^{n \times k}$ ,  $V \in \mathbb{R}^{n \times v}$ ,  $Q \in \mathbb{R}^{b \times k}$

Get attention weights matrix:

$$W = \text{softmax} \left( \frac{QK^T}{\sqrt{k}}, \text{dim}=1 \right) \in \mathbb{R}^{b \times n}$$

**output:**  $WV \in \mathbb{R}^{b \times v}$



# Attention

**input:**  $K \in \mathbb{R}^{n \times k}$ ,  $V \in \mathbb{R}^{n \times v}$ ,  $Q \in \mathbb{R}^{b \times k}$

Get attention weights matrix:

$$W = \text{softmax} \left( \frac{QK^T}{\sqrt{k}}, \text{dim}=1 \right) \in \mathbb{R}^{b \times n}$$

**output:**  $WV \in \mathbb{R}^{b \times v}$

**parameters:** :(

# Attention

**input:**  $K \in \mathbb{R}^{n \times k}$ ,  $V \in \mathbb{R}^{n \times v}$ ,  $Q \in \mathbb{R}^{b \times k}$

Get attention weights matrix:

$$W = \text{softmax} \left( \frac{QK^T}{\sqrt{k}}, \text{dim}=1 \right) \in \mathbb{R}^{b \times n}$$

**output:**  $WV \in \mathbb{R}^{b \times v}$

**parameters:** :(

Today...

$$K \equiv V$$

# Attention

input:  $K \in \mathbb{R}^{n \times k}$ ,  $V \in \mathbb{R}^{n \times v}$ ,  $Q \in \mathbb{R}^{b \times k}$

Get attention weights matrix:

$$W = \text{softmax} \left( \frac{QK^T}{\sqrt{k}}, \text{dim}=1 \right) \in \mathbb{R}^{b \times n}$$

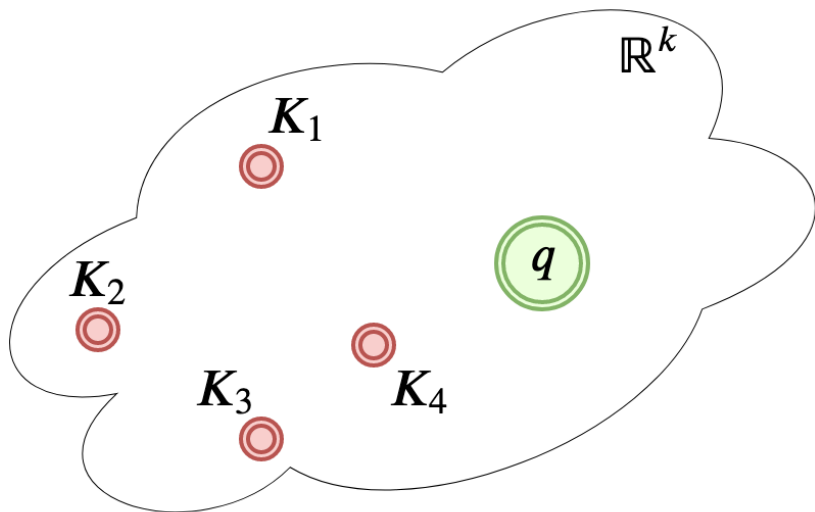
output:  $WV \in \mathbb{R}^{b \times v}$

parameters: :(

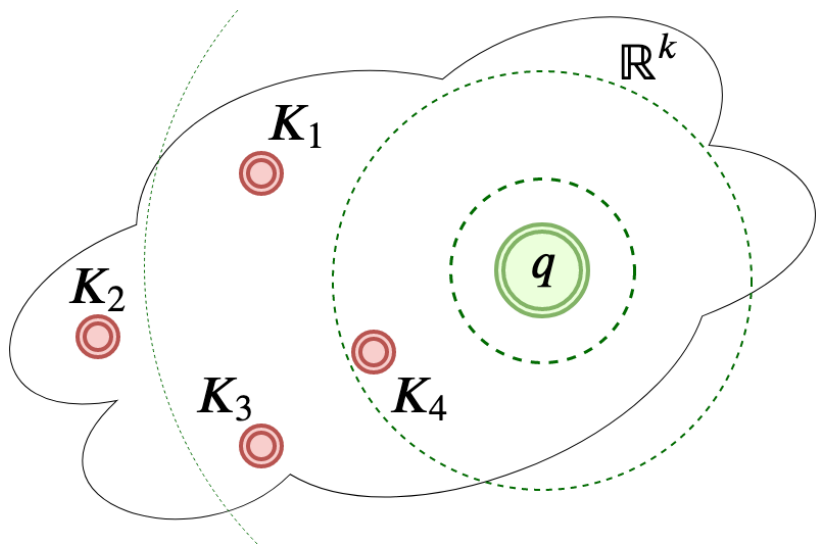
Today...

$$K \equiv V$$

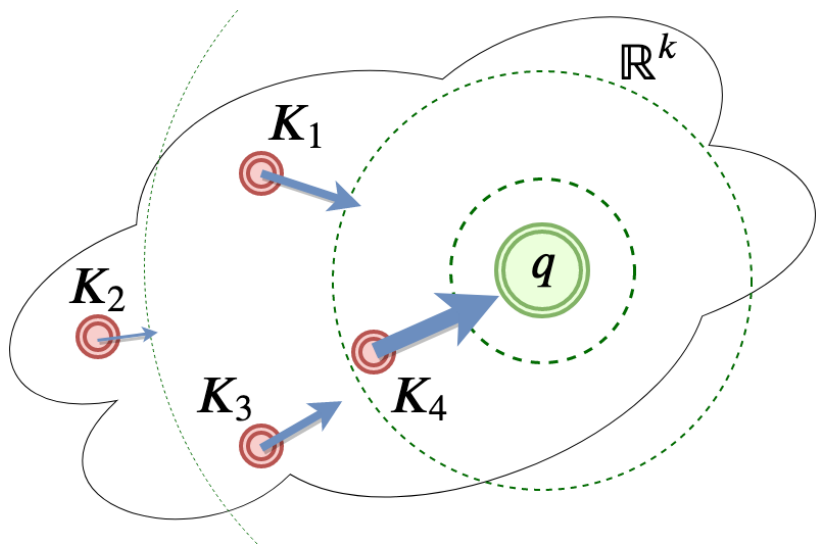
# Attention intuition



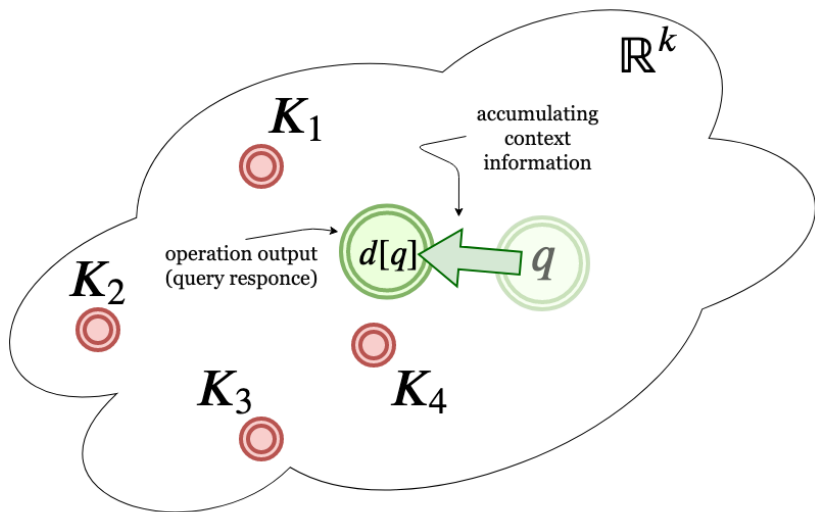
# Attention intuition



# Attention intuition



# Attention intuition

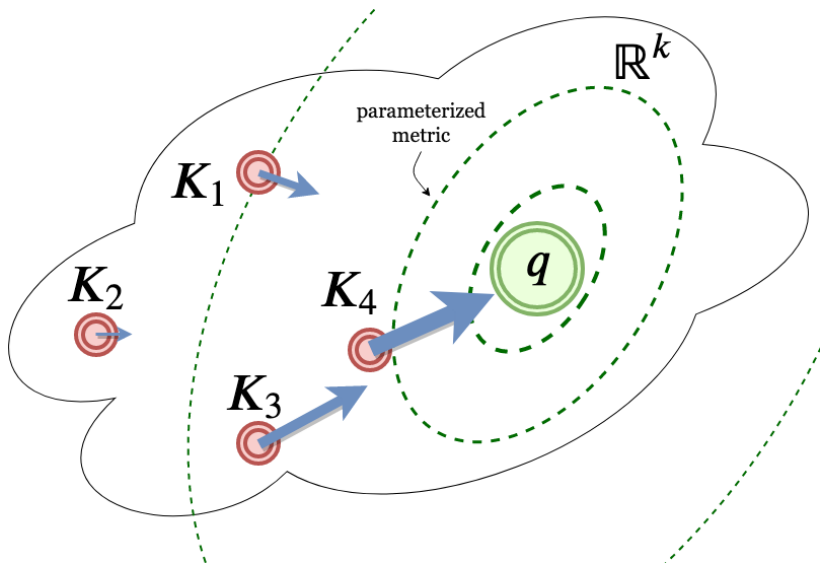


## Section 3

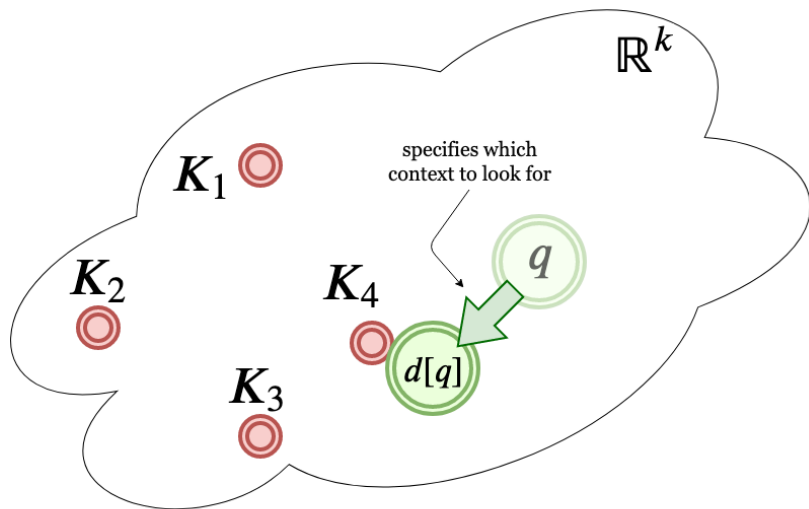
# Multi-Head Attention



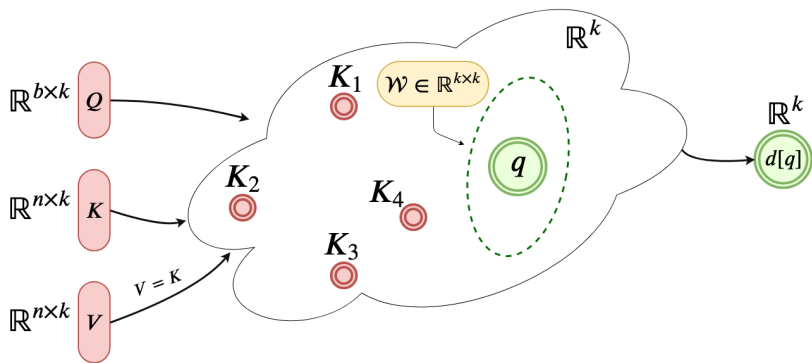
# Metric parametrization



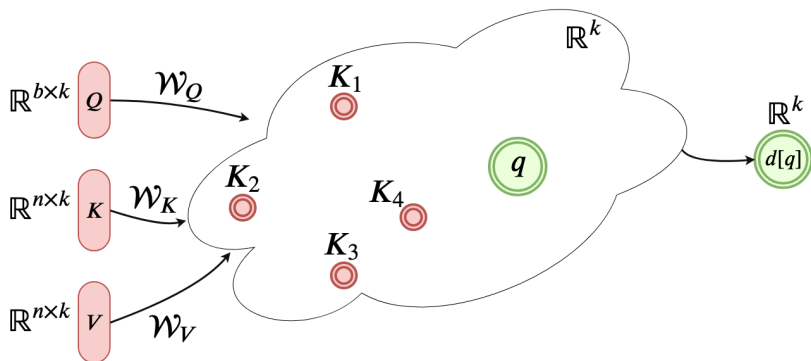
# Metric parametrization



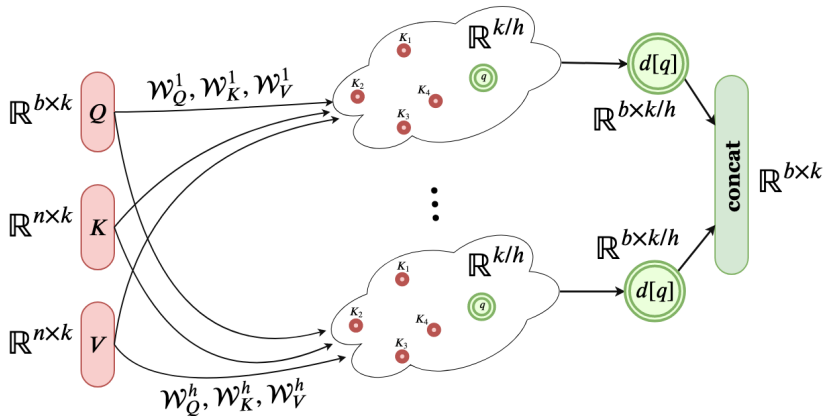
# Multi-Head Attention



# Multi-Head Attention



# Multi-Head Attention



# Self-attention

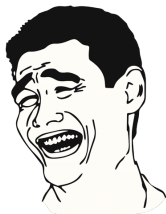
## Self-attention

$$Q \equiv K \equiv V$$

# Self-attention

## Self-attention

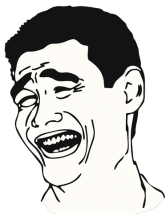
$$Q \equiv K \equiv V$$



# Self-attention

## Self-attention

$$Q \equiv K \equiv V$$



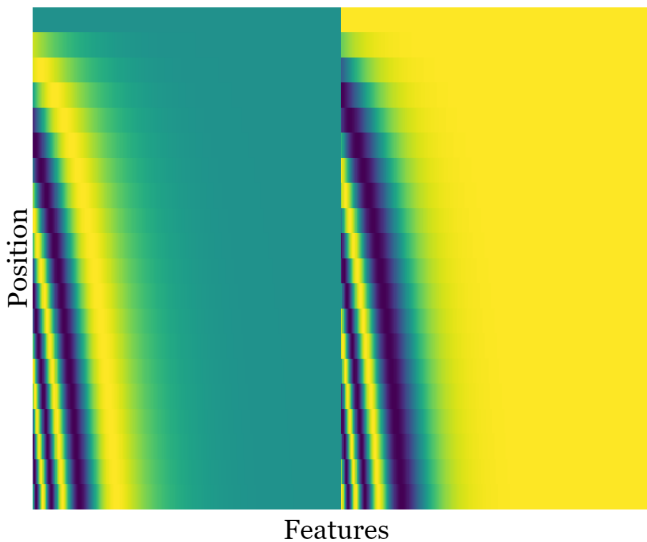
! not all multi-head attention blocks in Transformer are self-attention!



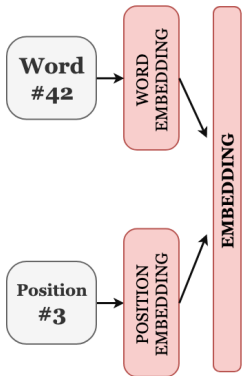
## Section 4

# Attention is All You Need

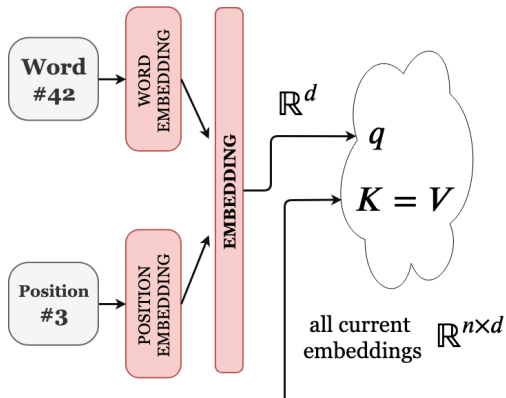
# Positional Embeddings



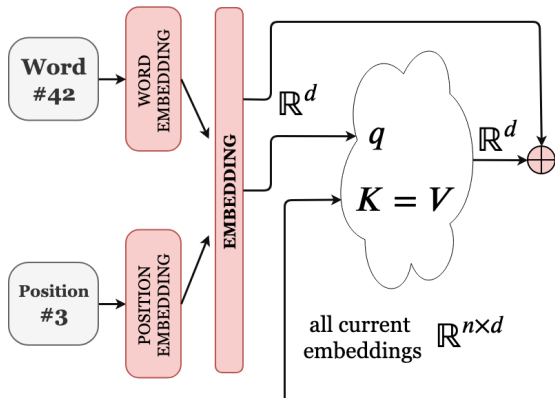
# Encoder



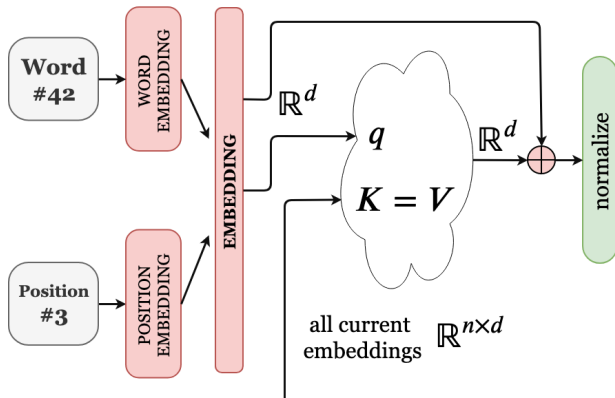
# Encoder



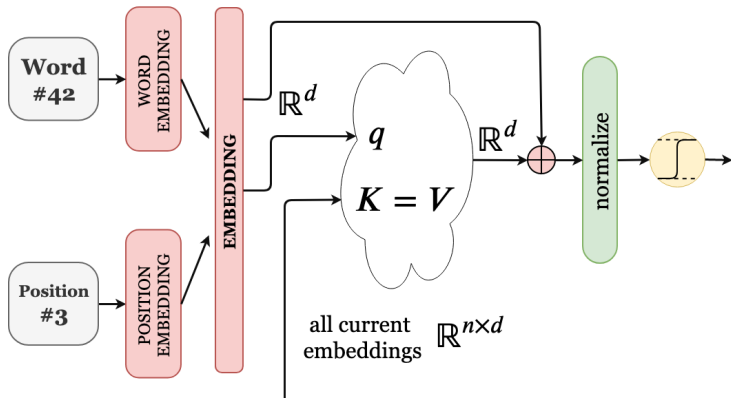
# Encoder



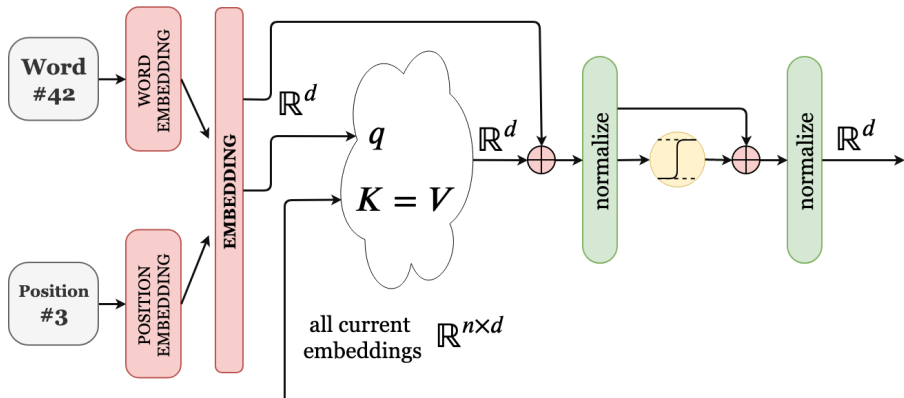
# Encoder



## Encoder

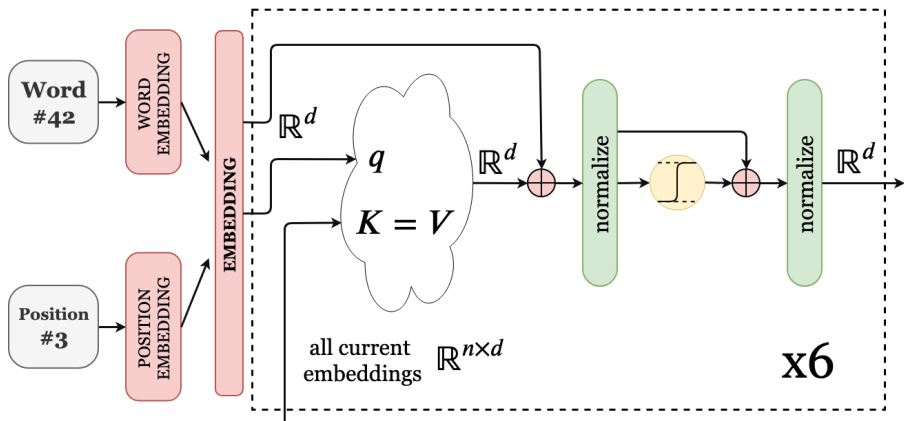


## Encoder

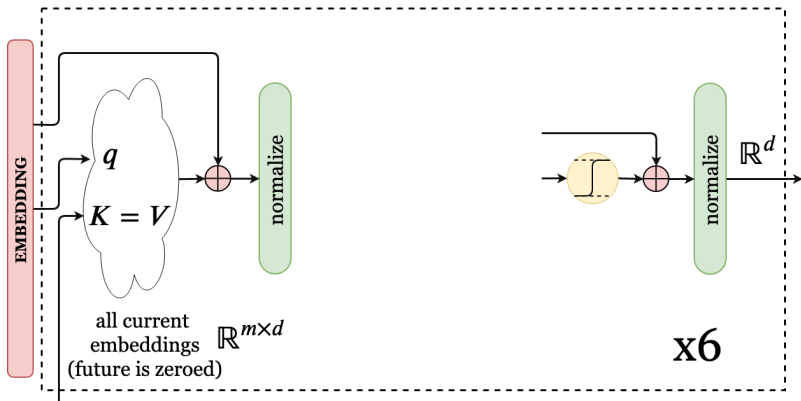




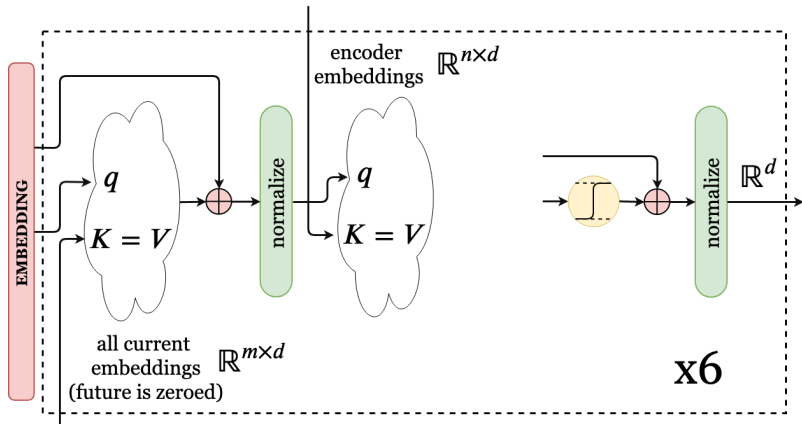
## Encoder



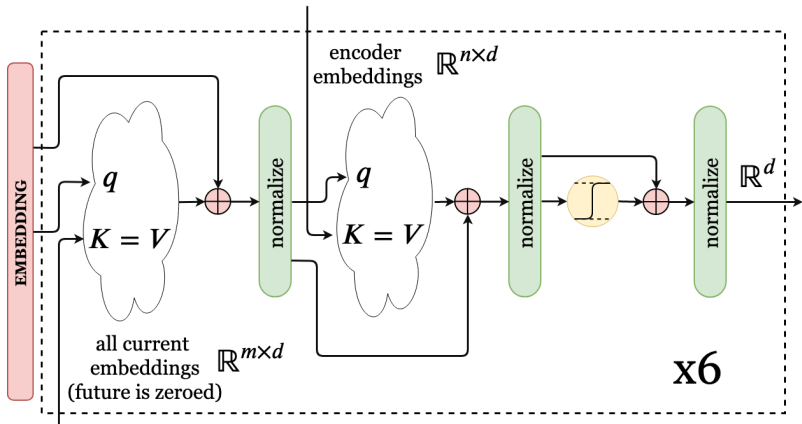
# Decoder



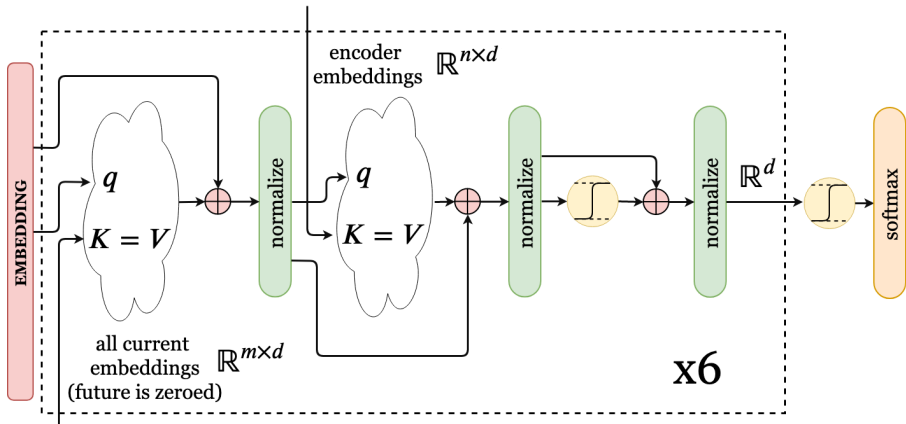
# Decoder



# Decoder



# Decoder



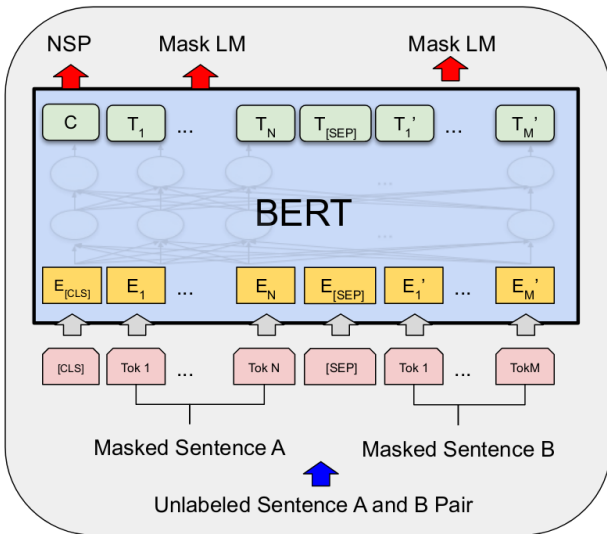
## More sources about Transformer

- **Animation of intuition:** [▶ Link](#)
- **Illustrated Transformer**  
`https://jalammr.github.io/illustrated-transformer/`
- **OpenAI Blog**  
`https://ai.googleblog.com/2017/08/transformer-novel-neural-network.html`
- **MIPT Lecture (RUS)**  
`https://www.youtube.com/watch?v=Bg8Y5q10iP0`

## Section 5

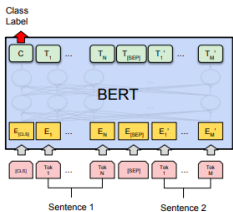
# BERT, GPT-2, ...

# BERT: pre-training

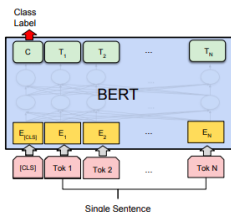




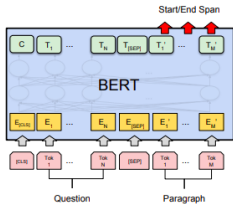
# BERT: fine-tuning



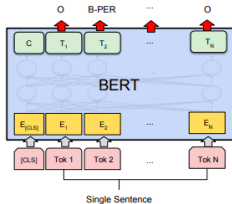
(a) Sentence Pair Classification Tasks:  
MNLI, QQP, QNLI, STS-B, MRPC,  
RTE, SWAG



(b) Single Sentence Classification Tasks:  
SST-2, CoLA



(c) Question Answering Tasks:  
SQuAD v1.1



(d) Single Sentence Tagging Tasks:  
CoNLL-2003 NER

# GPT-2 (Generative Pre-Training)

- language model based on masked multi-head self-attention
- with 1.5 billions of parameters (!)
- (rumors) 2048 TPU days to train
- which is able to generate pretty realistic texts

# GPT-2 (Generative Pre-Training)

- language model based on masked multi-head self-attention
- with 1.5 billions of parameters (!)
- (rumors) 2048 TPU days to train
- which is able to generate pretty realistic texts



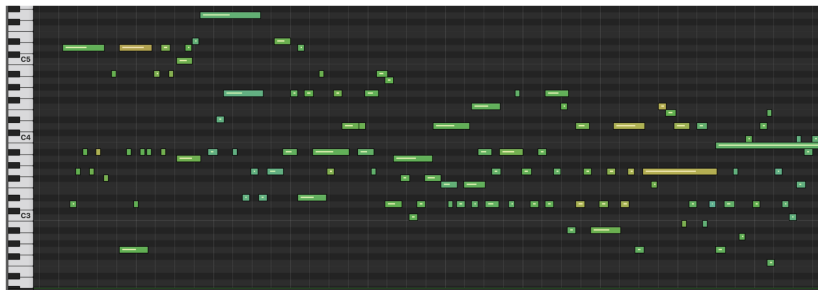
# How to apply this to music?



# How to apply this to music?



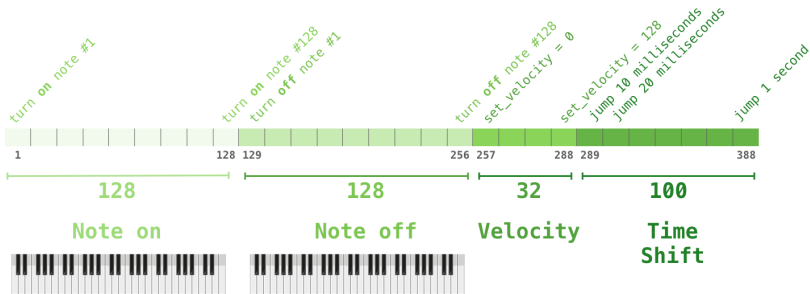
# MIDI Music Representation



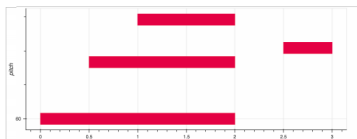
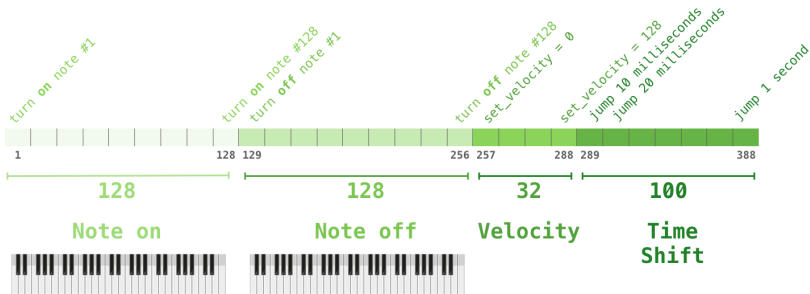
## Each note has:

- beginning (milliseconds)
- end (milliseconds)
- pitch (key): 128 possible options
- velocity (128 possible options; but we can take smaller grid)

# Notes is Language



# Notes is Language



```

SET_VELOCITY<80>, NOTE_ON<60>
TIME_SHIFT<500>, NOTE_ON<64>
TIME_SHIFT<500>, NOTE_ON<67>
TIME_SHIFT<1000>, NOTE_OFF<60>, NOTE_OFF<64>,
NOTE_OFF<67>
TIME_SHIFT<500>, SET_VELOCITY<100>, NOTE_ON<65>
TIME_SHIFT<500>, NOTE_OFF<65>
  
```



# Looking for more details...

- **Generalized Language Modeling (BERT section)**  
<https://lilianweng.github.io/lil-log/2019/01/31/generalized-language-models.html#bert>
- **Illustrated GPT**  
<http://jalamar.github.io/illustrated-gpt2/>
- **Transformer-XL**  
<https://arxiv.org/abs/1901.02860>
- **Music Transformer**  
<https://magenta.tensorflow.org/music-transformer>
- **Generated music visualisation:** [▶ Link](#)