

# Attention Models in Natural Language Processing

Tomasz Rembiasz, Juliusz Straszynski & Michał Maras

HUAWEI Warsaw Research Center

Phd Open

21 April 2022, University of Warsaw



# Outline

- 1 Attention Models in Machine Learning
- 2 Recurrent Neural Networks (RNNs)
- 3 Sequence-to-Sequence Models (seq2seq)
- 4 Attention in Natural Language Processing

# Outline

- 1 Attention Models in Machine Learning
- 2 Recurrent Neural Networks (RNNs)
- 3 Sequence-to-Sequence Models (seq2seq)
- 4 Attention in Natural Language Processing

# Attention Models in Machine Learning

- **Natural Language Processing**
- Computer Vision
- Reinforcement Learning

# Applications of attention in NLP

- Neural Machine Translation
- Question Answering (Q&A)
- Text Summarization
- Autocomplete
- Chatbot\*
- Named Entity Recognition (NER)\*
- Natural Language Understanding (NLU)\*

\* done at HUAWEI WRC



# Chatbot

## More task oriented

```

>>> I'm looking for a room in the city centre.
<<< There are many options. Do you have a price range in mind?
>>> I'd prefer something cheap.
<<< I have two options for you. One is in the cheap price range and the other is in the expensive price range. Which would you prefer?
>>> Cheap price range.
<<< I have the Acorn Guest House. Would you like me to book it for you?
>>> Where is it located?
<<< It is located at 156 chesterston road.
>>> Is it near the city centre?
<<< Yes, it is in the city centre. Would you like me to book it for you?
>>> Oh, sure.
<<< How many people will be staying?
>>> Three.
<<< I have booked you a room at the Acorn Guest House. Your reference number is: 1Z5VZQV. Is there anything else I can help you with?
  
```

## More chit-chat style

```

>>> Do you have any hobby?
<<< I like to play the piano.
>>> What type of music do you play?
<<< I play the piano in a classical band.
>>> Sounds great! Will you play me something?
<<< Sure. What do you want to hear?
>>> Maybe some piece by Chopin, how about that?
<<< Sure.
  
```

Conversation(s) with one of chatbots developed by our AI Team (Internship Project Summer 2021).

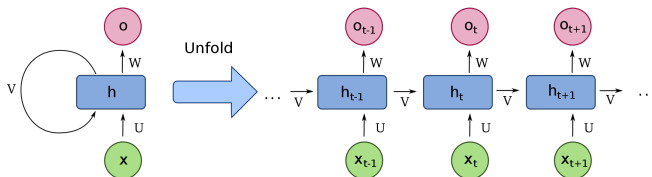
- Architecture: GPT-2
- Pretrained model: DialoGPT
- Dataset: custom mixture



# Outline

- 1 Attention Models in Machine Learning
- 2 Recurrent Neural Networks (RNNs)**
- 3 Sequence-to-Sequence Models (seq2seq)
- 4 Attention in Natural Language Processing

# Recurrent Neural Networks (RNNs)



From Wikipedia

## Selected applications:

- Machine translation
- Speech recognition
- Speech synthesis
- Grammar learning
- Music composition
- Handwriting recognition
- Protein homology detection
- Brain-computer interfaces



# Recurrent Neural Networks (RNNs)

words: one-hot vectors (of dim.  $v$ )

$$\mathbf{x}^{(t)} \in \mathbb{R}^v$$

word embeddings (of dimension  $n$ )

$$\mathbf{e}^{(t)} = \mathbf{E}\mathbf{x}^{(t)} \in \mathbb{R}^n$$

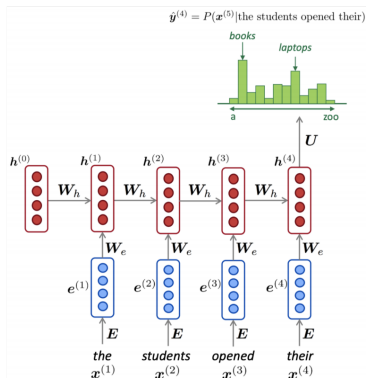
hidden states ( $\mathbf{h}^{(0)}$  being initial)

$$\mathbf{h}^{(t)} = \sigma \left( \mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_e \mathbf{e}^{(t)} + \mathbf{b}_1 \right)$$

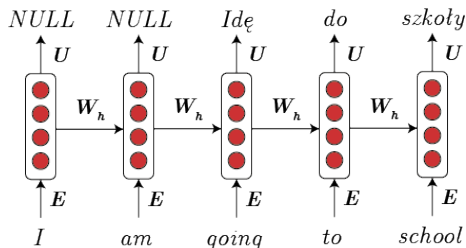
output distribution

$$\hat{\mathbf{y}}^{(t)} = \text{softmax} \left( \mathbf{U}\mathbf{h}^{(t)} + \mathbf{b}_2 \right) \in \mathbb{R}^v$$

$$\mathbf{x}^{(t)} \parallel \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \\ 0 \end{pmatrix}$$



# Machine Translation with RNNs

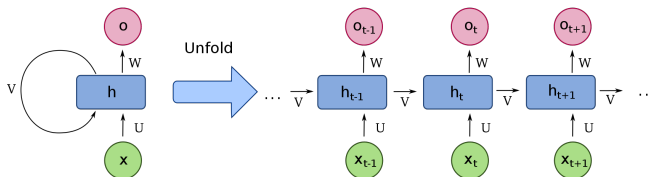


Sentences usually do not translate one-to-one.

*I am going to school.* → *Idę do szkoły.*

**Alignment** is the correspondence between particular words in the translated sentence pair.

# Drawbacks of RNNs

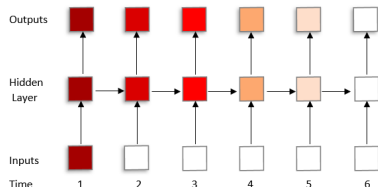
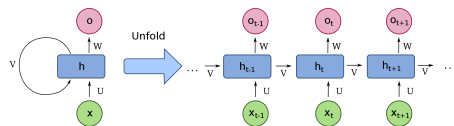


From Wikipedia

- Parallel computing difficult/impossible to implement
- Loss of information for long sequences
- Vanishing/exploding gradient problem

# Vanishing gradient problem

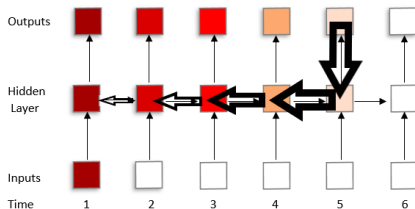
*Ich bin zu Hause, weil ich ... [VERY LONG GERMAN SENTENCE] ... bin.*



In RNNs, the hidden state gets constantly rewritten

$$\mathbf{h}^{(t)} = \sigma \left( \mathbf{V}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)} + \mathbf{b} \right).$$

# Vanishing gradient problem



$J^{(t)}(\theta)$  - loss function on step  $t$ .

$$\theta^{\text{new}} = \theta^{\text{old}} - \alpha \nabla_{\theta} J(\theta)$$

$$\frac{\partial J^{(t)}}{\partial \mathbf{h}^{(1)}} = \frac{\partial \mathbf{h}^{(2)}}{\partial \mathbf{h}^{(1)}} \frac{\partial \mathbf{h}^{(3)}}{\partial \mathbf{h}^{(2)}} \cdots \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}} \frac{\partial J^{(t)}}{\partial \mathbf{h}^{(t)}}$$

# Vanishing/exploding gradient problem

$$\mathbf{h}^{(t)} = \sigma \left( \mathbf{V}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)} + \mathbf{b} \right)$$

assume  $\sigma = \text{Id}$ , then

$$\frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}} = \mathbf{V}$$

$$\frac{\partial J^{(t)}}{\partial \mathbf{h}^{(1)}} = \frac{\partial \mathbf{h}^{(2)}}{\partial \mathbf{h}^{(1)}} \frac{\partial \mathbf{h}^{(3)}}{\partial \mathbf{h}^{(2)}} \cdots \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}} \frac{\partial J^{(t)}}{\partial \mathbf{h}^{(t)}} = \frac{\partial J^{(t)}}{\partial \mathbf{h}^{(t)}} \mathbf{V}^{t-1}$$

# Vanishing/exploding gradient problem

In basis of eigenvectors of  $\mathbf{V}$ , i.e.

$$\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_n$$

with eigenvalues

$$\lambda_1, \lambda_2, \dots, \lambda_n,$$

we can write

$$\frac{\partial J^{(t)}}{\partial \mathbf{h}^{(t)}} \mathbf{V}^{t-1} = \sum_{i=1}^n c_i \lambda_i^{t-1} \mathbf{q}_i.$$

For  $t \gg 1$ , if

$$\lambda_i < 1 \rightarrow$$

vanishing gradient (memory needed!),

$$\lambda_i > 1 \rightarrow$$

exploding gradient (gradient clipping).

# Long Short-Term Memory (LSTM)

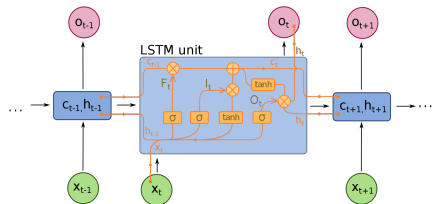
Hochreiter & Schmidhuber 1997

forget gate	$\mathbf{f}^{(t)} = \sigma \left( \mathbf{W}_f \mathbf{h}^{(t-1)} + \mathbf{U}_f \mathbf{x}^{(t)} + \mathbf{b}_f \right)$
input gate	$\mathbf{i}^{(t)} = \sigma \left( \mathbf{W}_i \mathbf{h}^{(t-1)} + \mathbf{U}_i \mathbf{x}^{(t)} + \mathbf{b}_i \right)$
output gate	$\mathbf{o}^{(t)} = \sigma \left( \mathbf{W}_o \mathbf{h}^{(t-1)} + \mathbf{U}_o \mathbf{x}^{(t)} + \mathbf{b}_o \right)$
new cell content	$\tilde{\mathbf{c}}^{(t)} = \tanh \left( \mathbf{W}_c \mathbf{h}^{(t-1)} + \mathbf{U}_c \mathbf{x}^{(t)} + \mathbf{b}_c \right)$
cell state	$\mathbf{c}^{(t)} = \mathbf{f}^{(t)} \odot \mathbf{c}^{(t-1)} + \mathbf{i}^{(t)} \odot \tilde{\mathbf{c}}^{(t)}$
hidden state	$\mathbf{h}^{(t)} = \mathbf{o}^{(t)} \odot \tanh \mathbf{c}^{(t)}$

$$\mathbf{f}^{(t)}, \mathbf{i}^{(t)}, \mathbf{o}^{(t)}, \tilde{\mathbf{c}}^{(t)}, \mathbf{c}^{(t)}, \mathbf{h}^{(t)}, \mathbf{x}^{(t)} \in \mathbb{R}^n$$



# Long Short-Term Memory (LSTM)



Idea: Hochreiter & Schmidhuber 1997. Figure: Wikipedia

- cells store long-term information
- read, write and erase options

$$\mathbf{f}^{(t)} = \sigma \left( \mathbf{W}_f \mathbf{h}^{(t-1)} + \mathbf{U}_f \mathbf{x}^{(t)} + \mathbf{b}_f \right)$$

$$\mathbf{i}^{(t)} = \sigma \left( \mathbf{W}_i \mathbf{h}^{(t-1)} + \mathbf{U}_i \mathbf{x}^{(t)} + \mathbf{b}_i \right)$$

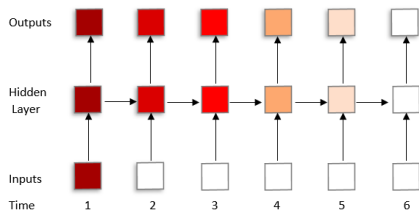
$$\mathbf{o}^{(t)} = \sigma \left( \mathbf{W}_o \mathbf{h}^{(t-1)} + \mathbf{U}_o \mathbf{x}^{(t)} + \mathbf{b}_o \right)$$

$$\tilde{\mathbf{c}}^{(t)} = \tanh \left( \mathbf{W}_c \mathbf{h}^{(t-1)} + \mathbf{U}_c \mathbf{x}^{(t)} + \mathbf{b}_c \right)$$

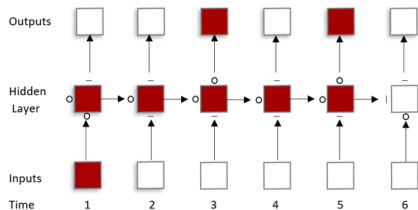
$$\mathbf{c}^{(t)} = \mathbf{f}^{(t)} \odot \mathbf{c}^{(t-1)} + \mathbf{i}^{(t)} \odot \tilde{\mathbf{c}}^{(t)}$$

$$\mathbf{h}^{(t)} = \mathbf{o}^{(t)} \odot \tanh \mathbf{c}^{(t)}$$

# Vanishing gradient problem and LSTM



RNN without LSTM



RNN with LSTM

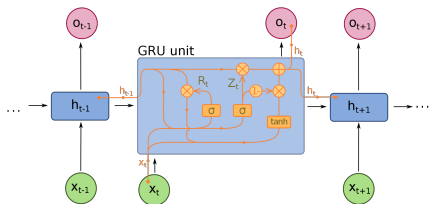
# Gated Recurrent Units (GRUs)

Cho et al. 2014

update gate	$\mathbf{u}^{(t)} = \sigma \left( \mathbf{W}_u \mathbf{h}^{(t-1)} + \mathbf{U}_u \mathbf{x}^{(t)} + \mathbf{b}_u \right)$
reset gate	$\mathbf{r}^{(t)} = \sigma \left( \mathbf{W}_r \mathbf{h}^{(t-1)} + \mathbf{U}_r \mathbf{x}^{(t)} + \mathbf{b}_r \right)$
new hidden state content	$\tilde{\mathbf{h}}^{(t)} = \tanh \left( \mathbf{W}_h (\mathbf{r}^{(t)} \mathbf{h}^{(t-1)} + \mathbf{U}_h \mathbf{x}^{(t)} + \mathbf{b}_h) \right)$
hidden state	$\mathbf{h}^{(t)} = (1 - \mathbf{u}^{(t)}) \odot \mathbf{h}^{(t-1)} + \mathbf{u}^{(t)} \odot \tilde{\mathbf{h}}^{(t)}$

$$\mathbf{u}^{(t)}, \mathbf{r}^{(t)}, \tilde{\mathbf{h}}^{(t)}, \mathbf{h}^{(t)} \in \mathbb{R}^n$$

# Gated Recurrent Units (GRUs)



Idea: Cho et al. 2014. Figure: Wikipedia

$$\mathbf{u}^{(t)} = \sigma \left( \mathbf{W}_u \mathbf{h}^{(t-1)} + \mathbf{U}_u \mathbf{x}^{(t)} + \mathbf{b}_u \right)$$

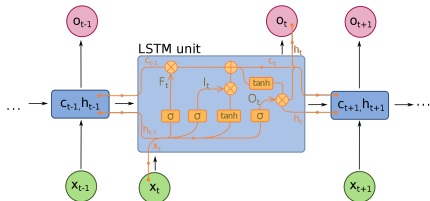
$$\mathbf{r}^{(t)} = \sigma \left( \mathbf{W}_r \mathbf{h}^{(t-1)} + \mathbf{U}_r \mathbf{x}^{(t)} + \mathbf{b}_r \right)$$

$$\tilde{\mathbf{h}}^{(t)} = \tanh \left( \mathbf{W}_h (\mathbf{r}^{(t)} \mathbf{h}^{(t-1)} + \mathbf{U}_h \mathbf{x}^{(t)} + \mathbf{b}_h) \right)$$

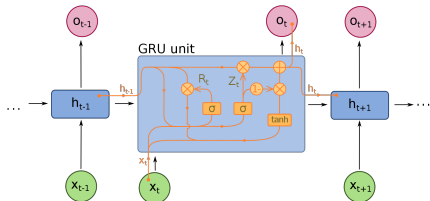
$$\mathbf{h}^{(t)} = (1 - \mathbf{u}^{(t)}) \odot \mathbf{h}^{(t-1)} + \mathbf{u}^{(t)} \odot \tilde{\mathbf{h}}^{(t)}$$

- No cell state
- Faster than LSTMs
- Fewer parameters

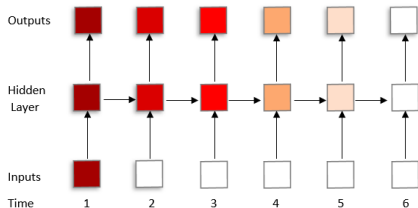
# Vanishing gradient problem and LSTM & GRU gates



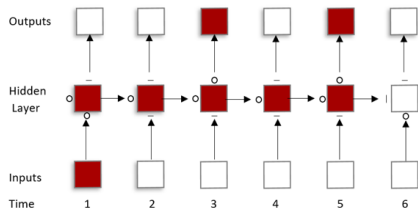
Idea: Hochreiter & Schmidhuber 1997. Figure: Wikipedia



Idea: Cho et al. 2014. Figure: Wikipedia



Vanishing gradient problem



RNN with LSTM

# Outline

- 1 Attention Models in Machine Learning
- 2 Recurrent Neural Networks (RNNs)
- 3 Sequence-to-Sequence Models (seq2seq)**
- 4 Attention in Natural Language Processing

# RNN encoder-decoder architecture

- Encoder:

input sentence  $\mathbf{x} = (x_1, \dots, x_{T_x})$  via

$$h_t = f(x_t, h_{t-1})$$

to context vector  $c$ :

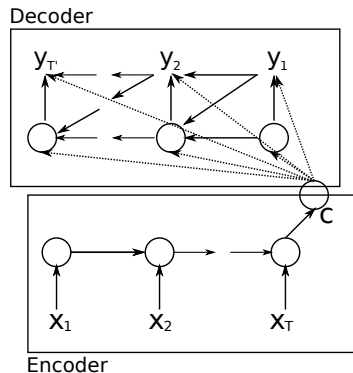
$$c = q(\{h_1, \dots, h_{T_x}\}).$$

- Decoder:

next word  $y_t$ :

$$p(y_t \mid \{y_1, \dots, y_{t-1}\}, c) = g(y_{t-1}, s_t, c)$$

from previous ones and (fixed)  $c$ .



From Cho et al. 2014

# RNN encoder-decoder architecture

**Encoder** converts the input sentence, a sequence of vectors  $\mathbf{x} = (x_1, \dots, x_{T_x})$ , into a *context vector*  $c$ .

First,  $h_t \in \mathbb{R}^n$ , encoder hidden states at time  $t$  are calculated:

$$h_t = f(x_t, h_{t-1}),$$

where  $f$  is a nonlinear function (eg. LSTM or GRU).

Next

$$c = q(\{h_1, \dots, h_{T_x}\}),$$

where  $g$  is a nonlinear function, eg.  $q(\{h_1, \dots, h_{T_x}\}) = h_{T_x}$ .



# RNN encoder-decoder architecture

**Decoder** predicts the next word,  $y_t$ , given the (fixed) context vector,  $c$ , and all the previously predicted words  $\{y_1, \dots, y_{t-1}\}$ :

$$p(\mathbf{y}) = \prod_{t=1}^{T_y} p(y_t \mid \{y_1, \dots, y_{t-1}\}, c).$$

Each conditional probability is modeled as:

$$p(y_t \mid \{y_1, \dots, y_{t-1}\}, c) = g(y_{t-1}, s_t, c),$$

where  $s_t$  is a decoder hidden state:

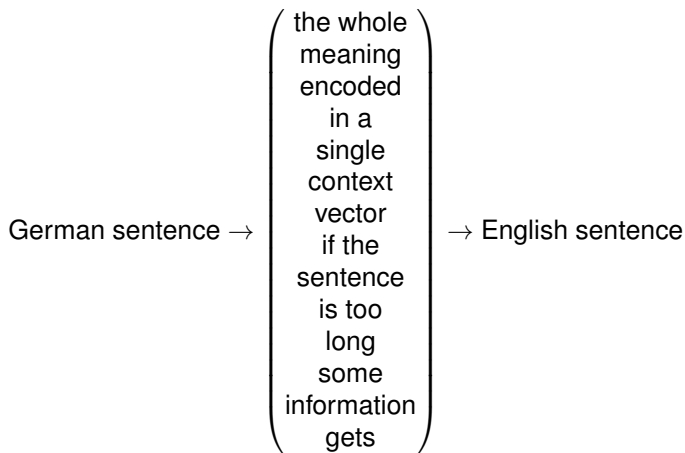
$$s_t = f(s_{t-1}, y_{t-1}, c),$$

$f$  and  $g$  are nonlinear functions.

# Outline

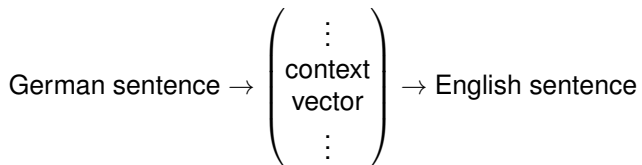
- 1 Attention Models in Machine Learning
- 2 Recurrent Neural Networks (RNNs)
- 3 Sequence-to-Sequence Models (seq2seq)
- 4 Attention in Natural Language Processing**

# Machine translation without attention



# Attention

## Without Attention



## With Attention

*Ich bin zu Hause, weil ich krank **bin**.*  $\rightarrow$  *I am home, because I **am** sick.*

- Sequence overload prevented.
- Model focuses on the likeliest word at each step.
- Uses a **query** (English word) to look in the right place to find the **key** (German word).

# Attention, Alignment

- Decoder without attention (alignment):

$$p(y_t | \{y_1, \dots, y_{t-1}\}, c) = g(y_{t-1}, s_t, c),$$

where  $s_t$  is a decoder hidden state:

$$s_t = f(s_{t-1}, y_{t-1}, c),$$

and  $c$  is a fixed context vector.

- Decoder with **attention** (alignment):

$$p(y_i | \{y_1, \dots, y_{i-1}\}, \mathbf{x}) = \tilde{g}(y_{i-1}, s_i, c_i),$$

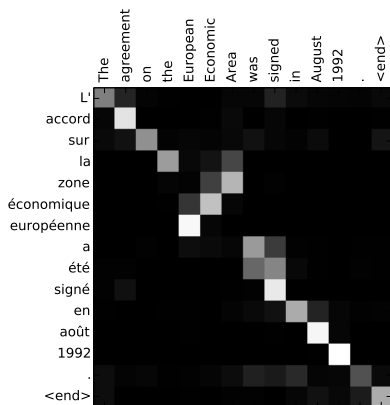
where  $s_i$  is a decoder hidden state:

$$s_i = \tilde{f}(s_{i-1}, y_{i-1}, c_i),$$

and  $c_i$  is a distinct context vector for each target word  $y_i$ .

# Neural Machine Translation

$\alpha_{ij} \rightarrow j$  source (English)  
 $\downarrow i$  target (French)



Context vector  $c_i$

$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j,$$

weights  $\alpha_{ij}$  of annotations  $h_j$

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})},$$

$$e_{ij} = a(s_{i-1}, h_j),$$

$a$  alignment model (FF NN)

$h_j$  encoder hidden state

$s_{i-1}$  decoder hidden state

From Bahdanau et al. 2015.

## Dot-Product Attention

### Dot-Product Attention

$$\text{Att}(Q, K, V) = \text{softmax}(QK^T) V,$$

where

$Q \in \mathbb{R}^{n_{\text{out}} \times d_m}$  - queries (decoder)

$K \in \mathbb{R}^{n_{\text{in}} \times d_m}$  - keys (encoder)

$V \in \mathbb{R}^{n_{\text{in}} \times d_m}$  - values (encoder)

$n_{\text{in}}, n_{\text{out}}$  - sequence lengths,

$d_m$  - embedding dimension

## Additive Attention

Context vector  $c_i$

$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j,$$

weights  $\alpha_{ij}$  of annotations  $h_j$

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})},$$

$$e_{ij} = a(s_{i-1}, h_j),$$

$a$  alignment model (FF NN)

$h_j$  encoder hidden state

$s_{i-1}$  decoder hidden state

# Dot-Product Attention

$$\text{Att}(Q, K, V) = \text{softmax}(QK^T) V,$$

where the dot-product of the  $i$ -th query  $Q_{i*} \in \mathbb{R}^{1 \times d_m}$  ( $i$ -th output token) with  $j$ -th key  $K_{j*} \in \mathbb{R}^{1 \times d_m}$  ( $j$ -th input token) is

$$(QK^T)_{ij} = \sum_{l=1}^{d_m} Q_{il}K_{jl}.$$

Function

$$\text{softmax}(QK^T)_{ij} = \frac{\exp\left[(QK^T)_{ij}\right]}{\sum_{l=1}^{n_{\text{in}}} \exp\left[(QK^T)_{il}\right]}$$

measures similarity between  $Q_{i*}$  and  $K_{j*}$ .



# Dot-Product Attention

Hence, the Dot-Product Attention

$$\text{Att}(Q, K, V) = \text{softmax}(QK^T) V,$$

for the  $i$ -th query  $Q_{i*}$ , based on the most similar to it key  $K_{j*}$ , predominantly picks the  $j$ -th value  $V_{j*}$ .

In other words, for each query  $Q_{i*} \in \mathbb{R}^{1 \times d_m}$ , the Dot-Product Attention  $\text{Att}(Q, K, V)_{i*} \rightarrow \mathbb{R}^{1 \times d_m}$  is a weighted average of all values  $V_{j*}$

$$\text{Att}(Q, K, V)_{i*} = \sum_{j=1}^{n_{in}} W_{ij} V_{j*},$$

with weights  $W_{ij} = \text{softmax}(QK^T)_{ij}$ .

# Additive Attention vs Dot-Product Attention

- Similar in theoretical complexity.
- Additive attention [AA] uses a feed-forward network with a single hidden layer.
- Dot-product (multiplicative) attention [DPA] can be faster and more space-efficient in practice when implemented with a highly optimized matrix multiplication code.
- For small  $d_m$ , both attentions perform similarly.
- For larger  $d_m$ , DPA gets slower than AA.
- To counteract it, Vaswani et al. 2017 introduced Scaled DPA .

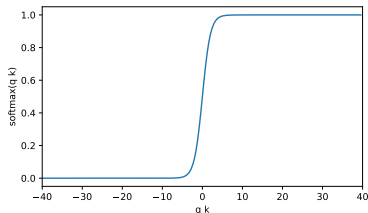


# Scaled Dot-Product Attention

## Dot-Product Attention

$$\text{Att}(Q, K, V) = \text{softmax}(QK^T) V.$$

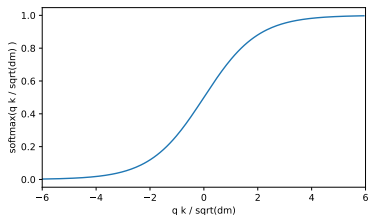
Each element  $(QK^T)_{ij}$  has variance  $\sim d_m$ .



## Scaled Dot-Product Attention

$$\widetilde{\text{Att}}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_m}}\right) V.$$

Each  $(QK^T / \sqrt{d_m})_{ij}$  has variance  $\sim 1$ .



# Variance in Dot-Product Attention

Matrix elements of the dot-product  $QK^T$  can get large for large  $d_m$ . For arbitrarily fixed  $i$  and  $j$ , define  $q = (q_1, \dots, q_{d_m}) \equiv Q_{i*}$  and  $k \equiv K_{j*}$ . Assume all  $q_\alpha$  and  $k_\beta$  are independent random variables with mean  $E(q_\alpha) = E(k_\beta) = 0$  and variance  $\text{Var}[q_\alpha] = \text{Var}[k_\beta] = 1$ . Then

$$E[q \cdot k] = E \left[ \sum_{\alpha=1}^{d_m} q_\alpha k_\alpha \right] = 0,$$

because  $\text{Cov}[q_\alpha, k_\alpha] = 0$  and

$$\text{Var}[q \cdot k] = E[(q \cdot k)^2] = E \left[ \sum_{\alpha=1}^{d_m} \sum_{\beta=1}^{d_m} q_\alpha k_\alpha q_\beta k_\beta \right] = \sum_{\alpha=1}^{d_m} E[q_\alpha^2] E[k_\alpha^2] = d_m,$$

because  $\text{Cov}[q_\alpha k_\alpha, q_\beta k_\beta] = \delta_{\alpha\beta}$  and  $\text{Cov}[q_\alpha^2, k_\alpha^2] = 0$ .



Thank you for your attention