**344.063 KV Special Topic:**
# Natural Language Processing with Deep Learning
## Transformers

Navid Rekab-saz

navid.rekabsaz@jku.at

JΣU
**JOHANNES KEPLER
UNIVERSITY LINZ**

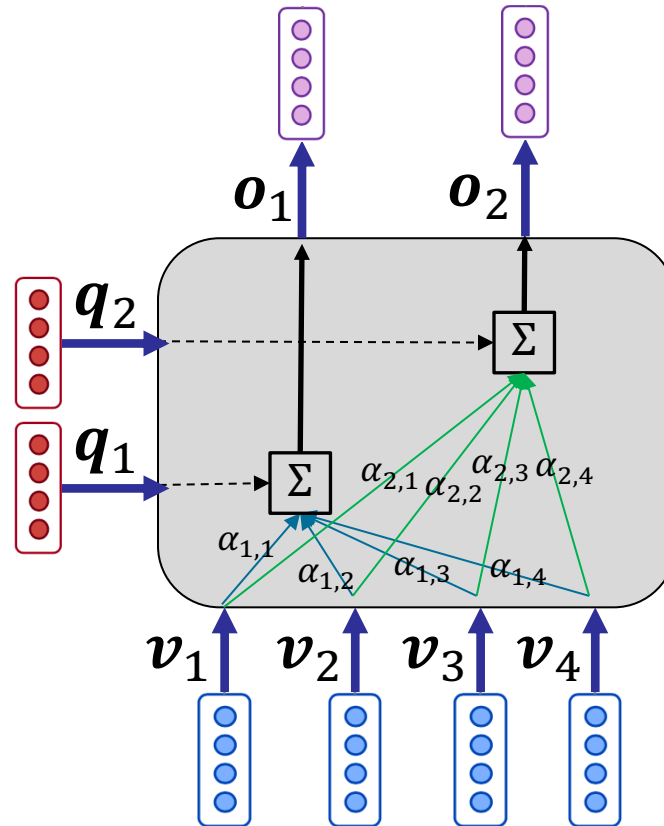Institute of
Computational
Perception

# Agenda

- Transformer encoder
- Transformer decoder
- seq2seq with Transformers

# Agenda

- **Transformer encoder**
- Transformer decoder
- seq2seq with Transformers

# Attentions! – recap



$\alpha_{i,j}$ is the attention score of query $q_i$ on value $v_j$

$\alpha_i$ is the vector of attentions of query $q_i$ over value vectors $V$ which forms a probability distribution

# Attention Networks – recap

- Given query vector $\boldsymbol{q}_i$, an attention network uses the attention similarity function $f$ to assign a non-normalized attention score $\tilde{\alpha}_{i,j}$ to value vector $\boldsymbol{v}_j$:
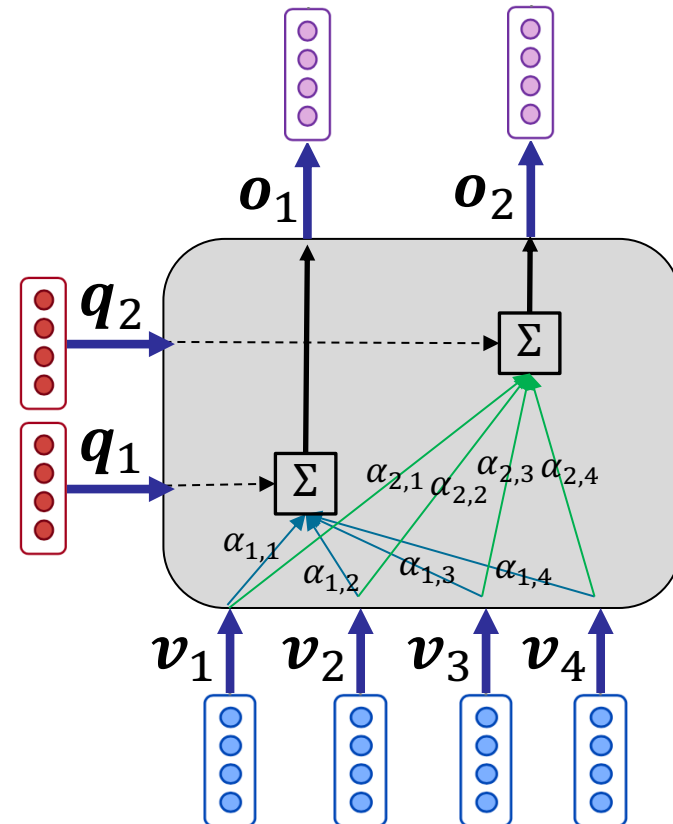
$$\tilde{\alpha}_{i,j} = f(\boldsymbol{q}_i, \boldsymbol{v}_j)$$

- Then, the attention scores over values are turned to a probability distribution using softmax:

$$\boldsymbol{\alpha}_i = \text{softmax}(\widetilde{\boldsymbol{\alpha}}_i), \qquad \sum_{j=1}^{|V|} \alpha_{i,j} = 1$$

- Finally, output vector $\boldsymbol{o}_i$ regarding query $\boldsymbol{q}_i$ is defined as the sum of the value vectors weighted by their corresponding attentions:

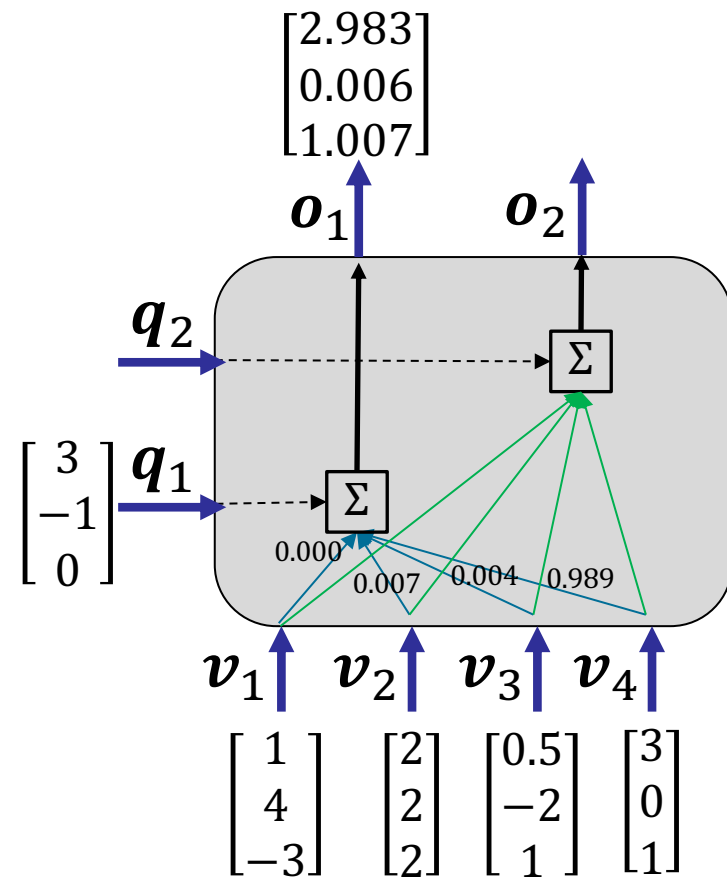$$\boldsymbol{o}_i = \sum_{j=1}^{|V|} \alpha_{i,j} \boldsymbol{v}_j$$



5

# Example – recap

$$\widetilde{\boldsymbol{\alpha}}_1 = \begin{bmatrix} \boldsymbol{q}_1 \boldsymbol{v}_1^{\mathrm{T}} = -1 \\ \boldsymbol{q}_1 \boldsymbol{v}_2^{\mathrm{T}} = 4 \\ \boldsymbol{q}_1 \boldsymbol{v}_3^{\mathrm{T}} = 3.5 \\ \boldsymbol{q}_1 \boldsymbol{v}_4^{\mathrm{T}} = 9 \end{bmatrix} \rightarrow \boldsymbol{\alpha}_1 = \begin{bmatrix} 0.000 \\ 0.007 \\ 0.004 \\ 0.989 \end{bmatrix}$$

$$\boldsymbol{o}_1 = 0.000 \begin{bmatrix} 1 \\ 4 \\ -3 \end{bmatrix} + 0.007 \begin{bmatrix} 2 \\ 2 \\ 2 \end{bmatrix} + 0.004 \begin{bmatrix} 0.5 \\ -2 \\ 1 \end{bmatrix} + 0.989 \begin{bmatrix} 3 \\ 0 \\ 1 \end{bmatrix}$$

$$\boldsymbol{o}_1 = \begin{bmatrix} 2.983 \\ 0.006 \\ 1.007 \end{bmatrix}$$
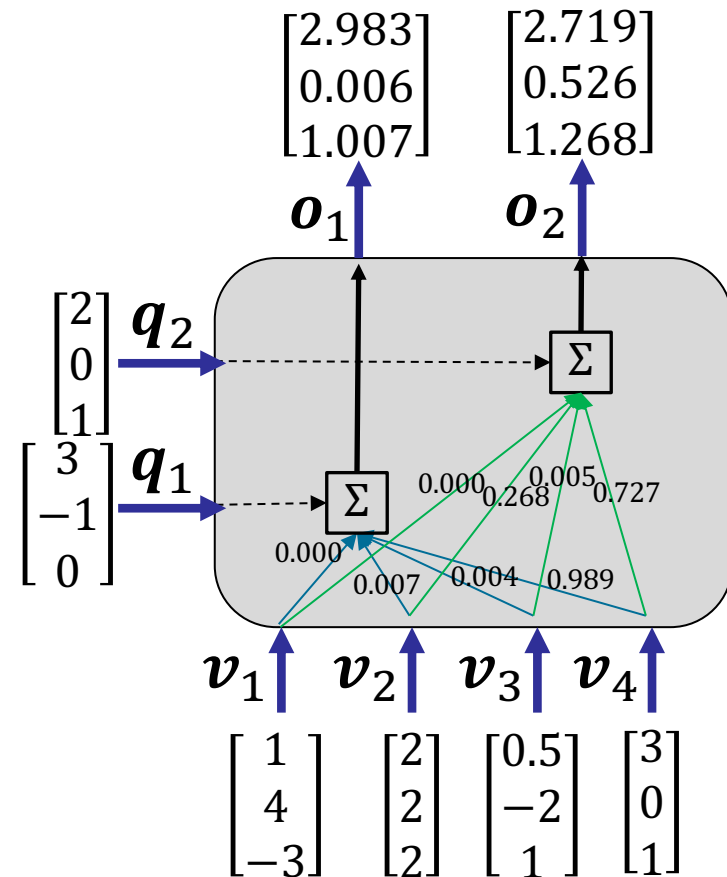
# Example – recap

$$\widetilde{\boldsymbol{\alpha}}_2 = \begin{bmatrix} \boldsymbol{q}_2\boldsymbol{v}_1^{\mathrm{T}} = -1 \\ \boldsymbol{q}_2\boldsymbol{v}_2^{\mathrm{T}} = 6 \\ \boldsymbol{q}_2\boldsymbol{v}_3^{\mathrm{T}} = 2 \\ \boldsymbol{q}_2\boldsymbol{v}_4^{\mathrm{T}} = 7 \end{bmatrix} \rightarrow \boldsymbol{\alpha}_2 = \begin{bmatrix} 0.000 \\ 0.268 \\ 0.005 \\ 0.727 \end{bmatrix}$$

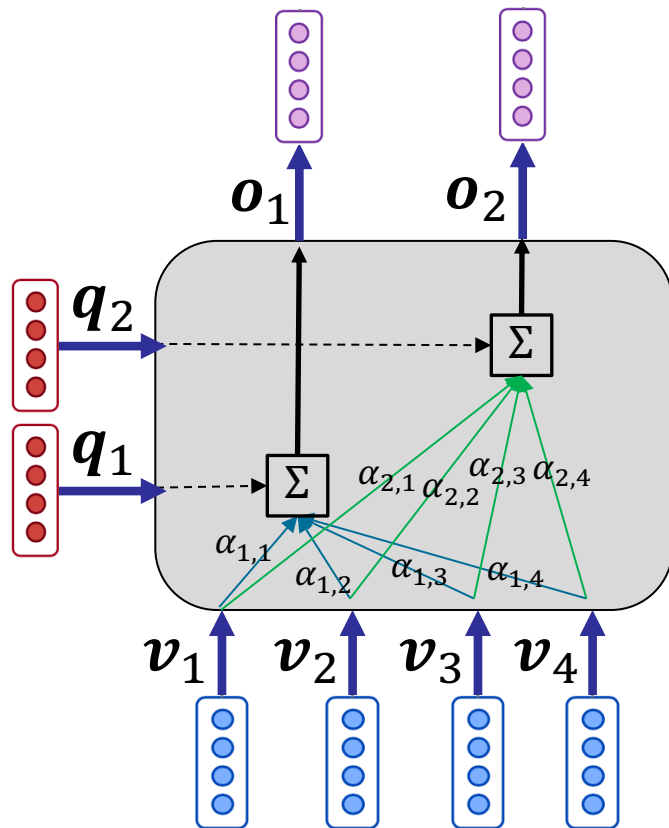$$\boldsymbol{o}_2 = 0.000 \begin{bmatrix} 1 \\ 4 \\ -3 \end{bmatrix} + 0.268 \begin{bmatrix} 2 \\ 2 \\ 2 \end{bmatrix} + 0.005 \begin{bmatrix} 0.5 \\ -2 \\ 1 \end{bmatrix} + 0.727 \begin{bmatrix} 3 \\ 0 \\ 1 \end{bmatrix}$$

$$\boldsymbol{o}_2 = \begin{bmatrix} 2.719 \\ 0.526 \\ 1.268 \end{bmatrix}$$
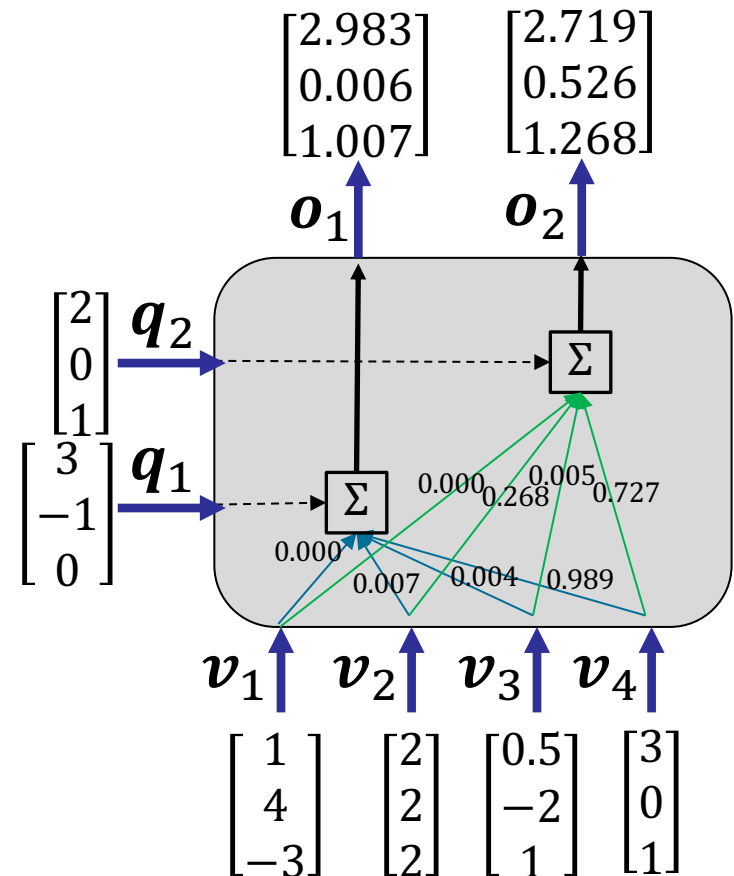
# Attention table

In the example:

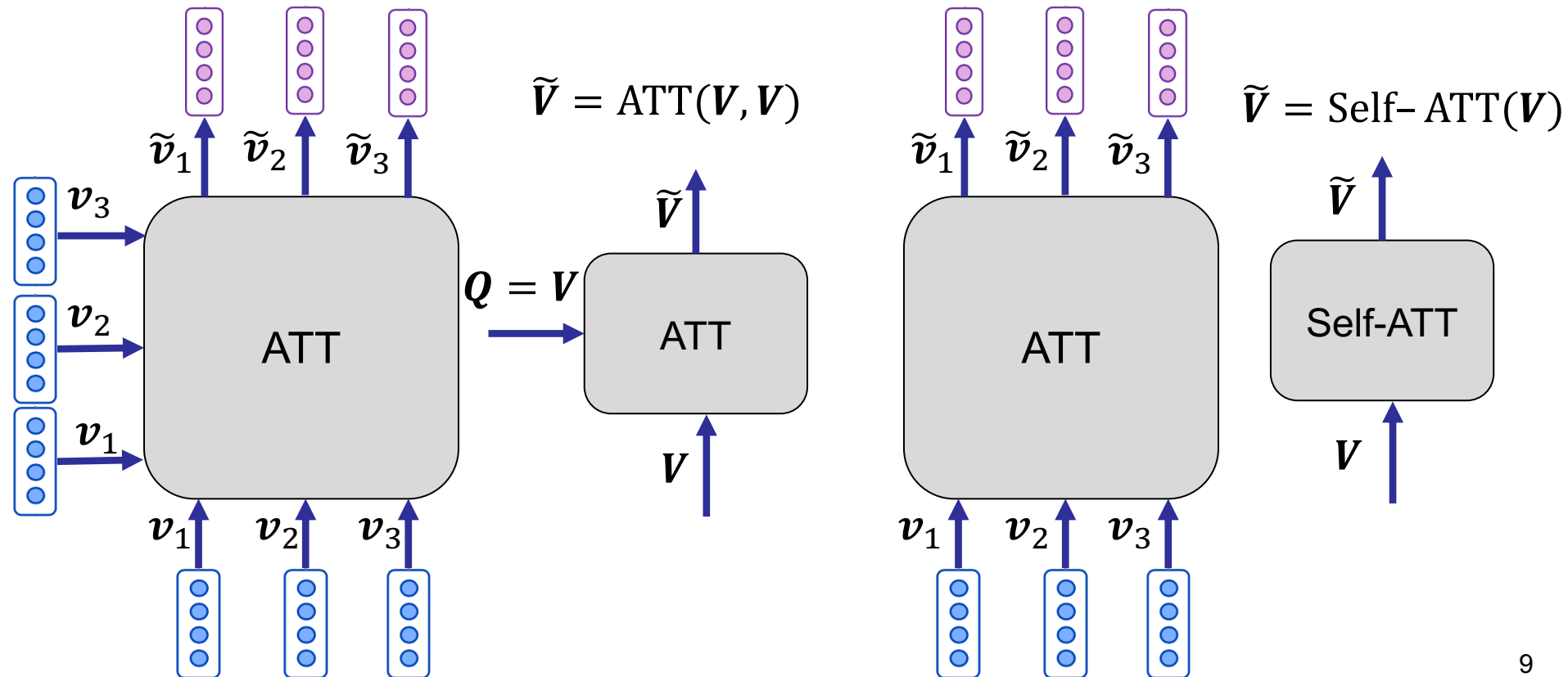| | $v_1$ | $v_2$ | $v_3$ | $v_4$ |
|---|---|---|---|---|
| $q_1$ | $\alpha_{1,1}$ | $\alpha_{1,2}$ | $\alpha_{1,3}$ | $\alpha_{1,4}$ |
| $q_2$ | $\alpha_{2,1}$ | $\alpha_{2,2}$ | $\alpha_{2,3}$ | $\alpha_{2,4}$ |

| | $v_1$ | $v_2$ | $v_3$ | $v_4$ |
|---|---|---|---|---|
| $q_1$ | 0.000 | 0.007 | 0.004 | 0.989 |
| $q_2$ | 0.000 | 0.268 | 0.005 | 0.727 |

# Self-attention

- Self-attention is when the values are also given as the queries: $Q = V$
- Self-attention encodes a sequence $V$ to a contextualized sequence $\widetilde{V}$
    - In self-attention, each input vector $v_i$ attends to all other input vectors $V$, and outputs $\widetilde{v}_i$ as a composition of input vectors
    - Output vector $\widetilde{v}_i$ is the contextual embedding of the input vector $v_i$

$$\widetilde{V} = \text{ATT}(V, V)$$

$$\widetilde{V} = \text{Self--ATT}(V)$$

# Transformers

- Attention network with DL best practices!
    - Originally introduced in the context of machine translation and is now widely adopted for sequence encoding and decoding

**Transformer Encoder**

**Transformer Decoder**

Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. "Attention is all you need." In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 2017.

# Transformer Encoder

- Transformer Encoder consists of two sub-layers:
  - 1st : Multi-head scaled dot-product self-attention
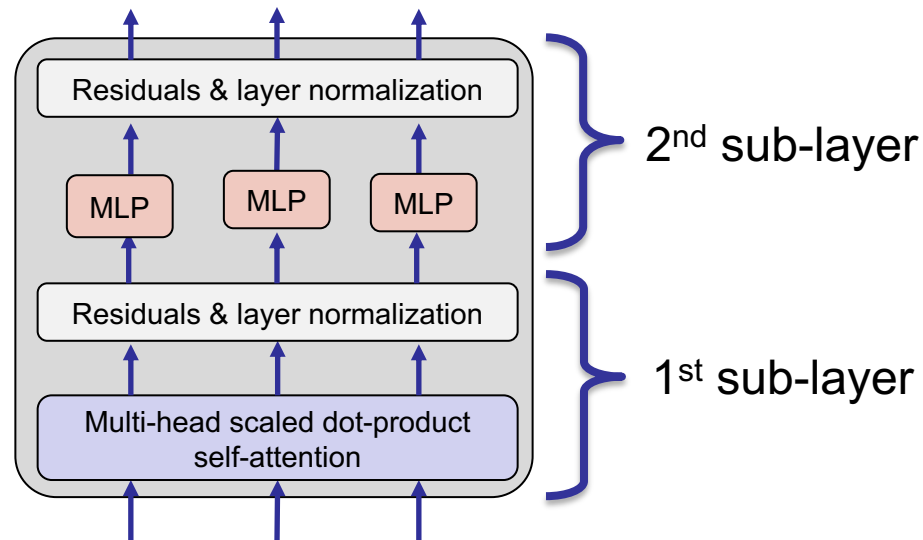  - 2nd : Position-wise multi-layer perceptron (feed forward)
- Each sub-layer is followed by residual networks and layer normalization
  - Drop-outs are applied after each computation

# Transformer Encoder

Let's start from <u>multi-head</u> <u>scaled dot-product</u> <u>self-attention</u>:
1. Scaled dot-product attention
2. Multi-head attention
3. self-attention



**Transformer Encoder**

# Transformer Encoder

Let's start from <u>multi-head</u> <u>scaled dot-product</u> <u>self-attention</u>:

1. **Scaled dot-product attention**
2. Multi-head attention
3. self-attention



**Transformer Encoder**

# Basic dot-product attention – recap

- Non-normalized attention scores:

$$\tilde{\alpha}_{i,j} = f(\boldsymbol{q}_i, \boldsymbol{v}_j)$$

$$\tilde{\alpha}_{i,j} = \boldsymbol{q}_i \boldsymbol{v}_i^{\mathrm{T}}$$

  - In this case, $d_q = d_v$
  - Attention network has no parameter to learn!

- Softmax over value vectors:

$$\boldsymbol{\alpha}_i = \mathrm{softmax}(\widetilde{\boldsymbol{\alpha}}_i)$$

- Output (weighted sum): $\boldsymbol{o}_i = \sum_{j=1}^{|V|} \alpha_{i,j} \boldsymbol{v}_j$

# Scaled dot-product attention

- Problem with basic dot-product attention:
    - As $d$ gets large, the variance of $\tilde{\alpha}_{i,j}$ increases …
    - … this makes softmax very peaked for some values of $\widetilde{\boldsymbol{\alpha}}_i$ …
    - … and hence its gradient gets smaller
- One approach: normalize/scale $\tilde{\alpha}_{i,j}$ by vector size $d$

**<u>Scaled dot-product attention</u>**

- Non-normalized attention scores:

$$\tilde{\alpha}_{i,j} = \frac{\boldsymbol{q}_i \boldsymbol{v}_j^{\mathrm{T}}}{\sqrt{d}}$$

- Softmax over values: $\boldsymbol{\alpha}_i = \mathrm{softmax}(\widetilde{\boldsymbol{\alpha}}_i)$

- Output: $\boldsymbol{o}_i = \sum_{j=1}^{|V|} \alpha_{i,j} \boldsymbol{v}_j$

# Transformer Encoder

Let's start from <u>multi-head</u> <u>scaled dot-product</u> <u>self-attention</u>:

1. Scaled dot-product attention
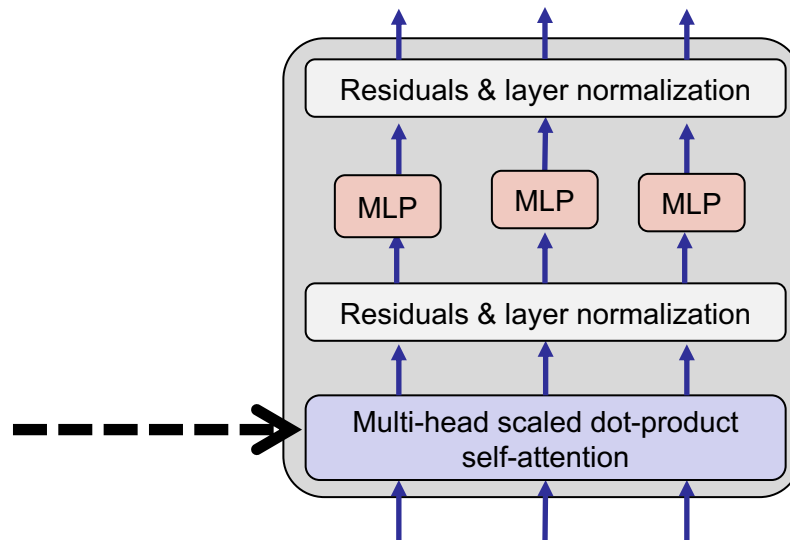2. **Multi-head attention**
3. self-attention



**Transformer Encoder**

# Softmax bottleneck!

- <u>Softmax</u> is applied to non-normalized attention vectors
  - Recall: softmax makes the maximum value much higher than the other

$$z = \begin{bmatrix} 1 & 2 & 5 & 6 \end{bmatrix} \rightarrow \text{softmax}(z) = \begin{bmatrix} 0.004 & 0.013 & 0.264 & 0.717 \end{bmatrix}$$

- Common in language, a word may be related to <u>several</u> other words in a sequence, each through a specific concept
  - Like the relations of a verb to its subject and object
- However, normal (single-head) attention network aggregates all concepts in one set
- In this case, due to softmax, value vectors must compete for the attention of query vector → softmax bottleneck

# Multi-head attention

- Multi-head attention approaches *softmax bottleneck* by calculating multiple sets of attentions between a query and values

**Multi-head attention:**

1. Transfer each query/value vector to $h$ query/value subspaces, each called a head

2. In each subspace, apply a normal (single-head) attention network using the queries and values transferred to the subspace to achieve the output vectors of that head

3. Concatenate the output vectors of all heads in respect to a query to achieve the final output of the query

- In multi-head attention, each head (and each subspace) can specialize on capturing a specific kind of relation

# Multi-head attention

# Multi-head attention – formulation

- Transfer every query $q_i$ to $h$ vectors, each with size ${d}/{h}$:

size: ${d}/{h}$

$$q_i^1 = q_i W_Q^1 \quad \ldots \quad q_i^h = q_i W_Q^h$$

Matrix size: $d \times {d}/{h}$

- Transfer every value $v_j$ to $h$ vectors, each with size ${d}/{h}$:

size: ${d}/{h}$

$$v_j^1 = v_j W_V^1 \quad \ldots \quad v_j^h = v_j W_V^h$$

Matrix size: $d \times {d}/{h}$

- Calculate outputs of subspaces corresponding to $q_i$:

size: ${d}/{h}$

$$o_i^1 = \mathrm{ATT}(q_i^1, V^1) \quad \ldots \quad o_i^h = \mathrm{ATT}(q_i^h, V^h)$$

- Concatenate outputs of subspaces for $q_i$ as its final output:

size: $d$

$$o_i = W_O \, [o_i^1; \ldots ; o_i^h]$$

Size: $d \times d$
This matrix linearly combines the dimensions of the concatenated vectors

Parameters are shown in red

20

# Multi-head attention – graphic in original paper



- Default number of heads in Transformers: $h = 8$
- Recall: Attentions (and Transformers) in fact have <u>three inputs</u> (not two), namely queries, <u>keys</u>, and values.
  - <u>Keys</u> are used to calculate attentions
  - <u>Values</u> are used to produce outputs

Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. "Attention is all you need." In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 2017.

# Transformer Encoder

Let's start from <u>multi-head</u> <u>scaled dot-product</u> <u>self-attention</u>:
1.  Scaled dot-product attention
2.  Multi-head attention
3.  **Self-attention**



**Transformer Encoder**

# Self-attention (recap)

- Values are the same as queries
- Each output vector is the contextual embedding of the corresponding input vector
  - $\widetilde{v}_i$ is the contextual embedding of $v_i$



$$\widetilde{V} = \text{Self-ATT}(V)$$

# Residuals

- Residual (short-cut) connection:
$$\text{output} = f(x) + x$$

- Learn in detail:
  - He, Kaiming; Zhang, Xiangyu; Ren, Shaoqing; Sun, Jian (2016). "Deep Residual Learning for Image Recognition" . In proc. of CVPR
  - Srivastava, Rupesh Kumar; Greff, Klaus; Schmidhuber, Jürgen (2015). "Highway Networks". https://arxiv.org/pdf/1505.00387.pdf



**Transformer Encoder**

# Layer normalization

- Layer normalization changes the activations of each vector to have mean 0 and variance 1 …
  - … and learns two parameters per layer to shift the mean and variance



**Transformer Encoder**

# Multi-layer perceptron on embedding

- A two-layer multi-layer perceptron (with ReLU) is applied to each output embedding
    - This layer provides the capacity for a non-linear transformation over each (contextualized) embedding



The same feed forward network is applied to every embedding

**Transformer Encoder**

# Transformer Encoder – all together

■ Transformer Encoder receive input embeddings and outputs the corresponding contextualized embeddings

- Processing all inputs happen at the same time → non auto-regressive

# Transformer Encoder – summary

- A self-attention model using
  - multi-head scaled dot-product attention
  - followed by the same feed-forward layer applied to each embedding
  - all packed with residuals, layer norms, and dropouts

**Transformers as in attentions …**

- do not have locality (position) bias
  - A long-distance context has "equal opportunity"
- process all the input together with a single computation per each layer
  - Friendly with parallel computations in GPU

Learn more and study the PyTorch implementation: http://nlp.seas.harvard.edu/2018/04/03/attention.html

# Position embeddings

- Transformers are agnostic to the position of tokens
  - A context token in long-distance has the same effect as the one in short-distance (no *locality bias*)
- However, the positions of tokens in a sequence might be informative and important in some tasks

**Position embeddings – a common approach in Transformers:**

- Create embeddings representing positions in a sequence, and add the values of the position embeddings to the token embeddings at corresponding positions
  - Position embedding is usually created using a sine/cosine function
    - It can also be learned end-to-end with the model parameters
  - Using position embeddings, the same token at different positions of a sequence will have different final representations

# Position embeddings – examples

**An example of embeddings with four dimensions:**



Position embedding for location 0

Position embeddings

Position embedding for location 20

Values from -1 (dark) to +1 (light)

Dimensions (512)

# Transformer Encoder with position embedding

$\tilde{\boldsymbol{e}}_1$ $\tilde{\boldsymbol{e}}_2$ $\tilde{\boldsymbol{e}}_3$

Residuals & layer normalization

MLP  MLP  MLP

Residuals & layer normalization

Multi-head scaled dot-product self-attention

| 0.00 | | 0.84 | | 0.91 |
| 0.00 | | 0.00 | | 0.00 |
| 1.00 | | 0.54 | | −0.4 |
| 1.00 | | 1 | | 1 |

+       +       +

$\boldsymbol{e}_1$   $\boldsymbol{e}_2$   $\boldsymbol{e}_3$

brown   lazy    dog

# Transformer Encoder with position embedding

Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. "Attention is all you need." In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 2017.

# Agenda

- Transformer encoder
- **Transformer decoder**
- seq2seq with Transformers

# Transformer Decoder

Transformer Decoder consists of three sub-layers:

- 1st : Masked multi-head self-attention
  - Exactly like Transformer Encoder but also with a *masking* functionality

- 2nd : Multi-head cross attention
  - Values are given from outside
    - Like from the outputs of a Transformer Encoder
  - Queries are the outputs of the 1st sub-layer

- 3rd : Position-wise multi-layer perceptron
  - Exactly like Transformer Encoder

3rd sub-layer

Residuals & layer normalization

MLP MLP MLP MLP

Residuals & layer normalization

Multi-head scaled dot-product cross-attention

2nd sub-layer

Residuals & layer normalization

Multi-head scaled dot-product self-attention

1st sub-layer

# Transformer Decoder with position embedding

Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. "Attention is all you need." In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 2017.

# Agenda

- Transformer encoder
- Transformer decoder
- **seq2seq with Transformers**

# Sequence-to-sequence modeling – recap

- Given the source sequence $X = \{x^{(1)}, x^{(2)}, ..., x^{(L)}\}$, …

- generate the target sequence $Y = \{y^{(1)}, y^{(2)}, ..., y^{(T)}\}$

- A seq2seq model estimates the conditional probability:

$$P(Y|X)$$

- and at inference time, it generates a new sequence $Y^*$ such that:

$$Y^* = \underset{Y}{\mathrm{argmax}}\, P(Y|X)$$

# Seq2seq with Transformers – training

$\hat{\boldsymbol{z}}^{(i)}$: predicted probability distribution of the next target word

$\hat{\boldsymbol{z}}^{(1)}$  $\hat{\boldsymbol{z}}^{(2)}$  $\hat{\boldsymbol{z}}^{(3)}$  $\hat{\boldsymbol{z}}^{(4)}$

$\boldsymbol{W}$  $\boldsymbol{W}$  $\boldsymbol{W}$  $\boldsymbol{W}$

$\boldsymbol{s}^{(1)}$  $\boldsymbol{s}^{(2)}$  $\boldsymbol{s}^{(3)}$  $\boldsymbol{s}^{(4)}$

**Transformer Decoder**

Residuals & layer normalization

MLP  MLP  MLP  MLP

Residuals & layer normalization

Multi-head scaled dot-product cross-attention

Contextualized <u>target</u> embeddings as **queries**

Contextualized source embeddings

$\tilde{\boldsymbol{e}}^{(1)}$  $\tilde{\boldsymbol{e}}^{(2)}$  $\tilde{\boldsymbol{e}}^{(3)}$

Residuals & layer normalization

**Transformer Encoder**

MLP  MLP  MLP

Residuals & layer normalization

Residuals & layer normalization

$\tilde{\boldsymbol{u}}^{(1)}$  $\tilde{\boldsymbol{u}}^{(2)}$  $\tilde{\boldsymbol{u}}^{(3)}$  $\tilde{\boldsymbol{u}}^{(4)}$

Multi-head scaled dot-product self-attention

Contextualized <u>source</u> embeddings as **values**

Multi-head scaled dot-product self-attention

Contextualized target embeddings

Position embedding

$\boldsymbol{e}^{(1)}$  $\boldsymbol{e}^{(2)}$  $\boldsymbol{e}^{(3)}$

$\boldsymbol{u}^{(1)}$  $\boldsymbol{u}^{(2)}$  $\boldsymbol{u}^{(3)}$  $\boldsymbol{u}^{(4)}$

38

# Seq2seq with Transformers – training

- Two sets of vocabularies
  - $\mathbb{V}_e$ is the set of vocabularies for source sequences
  - $\mathbb{V}_d$ is the set of vocabularies for target sequences
- Source sequence $X$ and target sequence $Y$
  - Both are typically started/ended with $< \text{bos} >/< \text{eos} >$

## **Encoder**

- Transformer encoder
  - passes source embeddings $[e^{(1)}, \dots, e^{(L)}]$ and creates contextualized source embeddings: $[\tilde{e}^{(1)}, \dots, \tilde{e}^{(L)}]$



39

# Seq2seq with Transformers – training

## Decoder

- Transformer Decoder self-attention layer

  - passes target embeddings $[u^{(1)}, ..., u^{(T)}]$ and creates contextualized target embeddings: $[\widetilde{u}^{(1)}, ..., \widetilde{u}^{(T)}]$

- Transformer Decoder cross-attention layer

  - applies attention with $[\widetilde{u}^{(1)}, ..., \widetilde{u}^{(T)}]$ as queries. and $[\tilde{e}^{(1)}, ..., \tilde{e}^{(L)}]$ as values (and keys)

- Transformer Decoder output

  - A set of vectors $[s^{(1)}, ..., s^{(T)}]$

**Incomplete version!**

40

# Seq2seq with Transformers – training

## Decoder (cont.)

- Decoder output prediction
  - uses $\left[ \boldsymbol{s}^{(1)}, ..., \boldsymbol{s}^{(T)} \right]$ to calculate $\left[ \hat{\boldsymbol{z}}^{(1)}, ..., \hat{\boldsymbol{z}}^{(T)} \right]$, the vectors of the predicted probability distribution at the next position:

$$\hat{\boldsymbol{z}}^{(t)} = \text{softmax}\left( \boldsymbol{W} \boldsymbol{s}^{(t)} + \boldsymbol{b} \right) \in \mathbb{R}^{|\mathbb{V}_d|}$$

- Training loss for each position $t$
  - NLL of the predicted probability of the next target word $y^{(t+1)}$

$$\mathcal{L}^{(t)} = -\log \hat{z}^{(t)}_{y^{(t+1)}}$$

- Overall loss is the average of loss values over the target sequence:

$$\mathcal{L} = \frac{1}{T} \sum_{t=1}^{T} \mathcal{L}^{(t)}$$

# Let's revisit the decoder!

**Decoder**

- Transformer Decoder self-attention layer

  - passes target embeddings $[\boldsymbol{u}^{(1)}, ..., \boldsymbol{u}^{(T)}]$ and creates contextualized target embeddings: $[\widetilde{\boldsymbol{u}}^{(1)}, ..., \widetilde{\boldsymbol{u}}^{(T)}]$

- Transformer Decoder cross-attention layer

  - applies attention with $[\widetilde{\boldsymbol{u}}^{(1)}, ..., \widetilde{\boldsymbol{u}}^{(T)}]$ as queries, and $[\tilde{\boldsymbol{e}}^{(1)}, ..., \tilde{\boldsymbol{e}}^{(L)}]$ as values (and keys)

- Transformer Decoder output

  - A set of vectors $[\boldsymbol{s}^{(1)}, ..., \boldsymbol{s}^{(T)}]$

**Problem:** in self-attention part, every token looks at all other tokens, namely the previous ones <u>but also the next tokens</u>!

- Every token has access to what it suppose to predict!

**Incomplete version!**



42

# Masking attentions

- In seq2seq with Transformers, we mask the attentions to every future token according to the <u>self-attentions</u> table of the <u>Transformer Decoder</u>

**Example**

- Non-normalized self-attention scores of Transformer Decoder:

attends to …
other target embeddings

|  | $u^{(1)}$ | $u^{(2)}$ | $u^{(3)}$ | $u^{(4)}$ |
|---|---|---|---|---|
| $u^{(1)}$ | 5 | 3 | 1 | -4 |
| $u^{(2)}$ | 1 | 4 | -2 | 3 |
| $u^{(3)}$ | 0 | 2 | 2 | -3 |
| $u^{(4)}$ | 3 | -1 | 1 | 4 |

Each target embedding

Non-normalized self-attention scores

| | $u^{(1)}$ | $u^{(2)}$ | $u^{(3)}$ | $u^{(4)}$ |
|---|---|---|---|---|
| $u^{(1)}$ | 5 | 3 | 1 | -4 |
| $u^{(2)}$ | 1 | 4 | -2 | 3 |
| $u^{(3)}$ | 0 | -2 | 2 | -3 |
| $u^{(4)}$ | 3 | -1 | 1 | 4 |

attentions masks

| | $u^{(1)}$ | $u^{(2)}$ | $u^{(3)}$ | $u^{(4)}$ |
|---|---|---|---|---|
| $u^{(1)}$ | 1 | 0 | 0 | 0 |
| $u^{(2)}$ | 1 | 1 | 0 | 0 |
| $u^{(3)}$ | 1 | 1 | 1 | 0 |
| $u^{(4)}$ | 1 | 1 | 1 | 1 |

Applying masks to attention scores
- adds $-\infty$ for every mask value 0
- adds 0 for every mask value 1

| | $u^{(1)}$ | $u^{(2)}$ | $u^{(3)}$ | $u^{(4)}$ |
|---|---|---|---|---|
| $u^{(1)}$ | 5 | $-\infty$ | $-\infty$ | $-\infty$ |
| $u^{(2)}$ | 1 | 4 | $-\infty$ | $-\infty$ |
| $u^{(3)}$ | 0 | -2 | 2 | $-\infty$ |
| $u^{(4)}$ | 3 | -1 | 1 | 4 |

softmax

Final self-attention scores

| | $u^{(1)}$ | $u^{(2)}$ | $u^{(3)}$ | $u^{(4)}$ |
|---|---|---|---|---|
| $u^{(1)}$ | 1.00 | 0.00 | 0.00 | 0.00 |
| $u^{(2)}$ | 0.04 | 0.96 | 0.00 | 0.00 |
| $u^{(3)}$ | 0.11 | 0.01 | 0.86 | 0.00 |
| $u^{(4)}$ | 0.25 | 0.01 | 0.34 | 0.70 |

☞ In Transformers, there are $h$ times of such attention matrices. The same masking is applied to each of them.

44

# Seq2seq with Transformers – training

**<u>Decoder</u>**

- Transformer Decoder self-attention layer
  - passes target embeddings $\left[\boldsymbol{u}^{(1)}, ..., \boldsymbol{u}^{(T)}\right]$ and creates contextualized target embeddings: $\left[\widetilde{\boldsymbol{u}}^{(1)}, ..., \widetilde{\boldsymbol{u}}^{(T)}\right]$ **while masking future tokens**

- Transformer Decoder cross-attention layer
  - applies attention with $\left[\widetilde{\boldsymbol{u}}^{(1)}, ..., \widetilde{\boldsymbol{u}}^{(T)}\right]$ as queries and $\left[\widetilde{\boldsymbol{e}}^{(1)}, ..., \widetilde{\boldsymbol{e}}^{(L)}\right]$ as values (and keys)

- Transformer Decoder output
  - A set of vectors $\left[\boldsymbol{s}^{(1)}, ..., \boldsymbol{s}^{(T)}\right]$

**Complete version!**

# Inference (decoding)

- During inference, as in training, the encoding of input sequence is done with a single computation (non-autoregressive)

- However, as in seq2seq with RNNs, decoding of seq2seq with Transformers is done in autoregressive fashion (one token after each other):
  - Pass the 1st target token ($< bos >$), generate the 2nd token
  - Pass the 1st token + the 2nd generated target tokens, generate the 3rd token
  - Pass the 1st token + the 2nd and 3rd generated target tokens, generate the 4th token
  - …

# Seq2seq with Transformers – inference (decoding)

$\hat{y}^{(2)}$ a sampled word from $\hat{z}^{(1)}$

$\hat{z}^{(1)}$

$W$

$s^{(1)}$

Residuals & layer normalization

MLP

Residuals & layer normalization

Multi-head scaled dot-product cross-attention

Residuals & layer normalization

Multi-head scaled dot-product self-attention

$\tilde{e}^{(1)}$    $\tilde{e}^{(2)}$    $\tilde{e}^{(3)}$

Residuals & layer normalization

MLP    MLP    MLP

Residuals & layer normalization

Multi-head scaled dot-product self-attention

$e^{(1)}$ $\oplus$    $e^{(2)}$ $\oplus$    $e^{(3)}$ $\oplus$

$u^{(1)}$ $\oplus$

$\hat{y}^{(2)}$

< bos >

47

# Seq2seq with Transformers – inference (decoding)



$\hat{y}^{(3)}$ a sampled word from $\boldsymbol{z}^{(2)}$

$\hat{\boldsymbol{z}}^{(2)}$

$W$

$\boldsymbol{s}^{(1)}$ $\boldsymbol{s}^{(2)}$

Residuals & layer normalization

MLP    MLP

Residuals & layer normalization

Multi-head scaled dot-product cross-attention

Residuals & layer normalization

Multi-head scaled dot-product self-attention

$\tilde{\boldsymbol{e}}^{(1)}$ $\tilde{\boldsymbol{e}}^{(2)}$ $\tilde{\boldsymbol{e}}^{(3)}$

Residuals & layer normalization

MLP    MLP    MLP

Residuals & layer normalization

Multi-head scaled dot-product self-attention

$\boldsymbol{e}^{(1)}$ $\boldsymbol{e}^{(2)}$ $\boldsymbol{e}^{(3)}$

$\boldsymbol{u}^{(1)}$ $\boldsymbol{u}^{(2)}$

$\hat{y}^{(3)}$

< bos >

48

**Seq2seq with Transformers – inference (decoding)**

49

# Seq2seq with Transformers – code

- Each Transformer encoder/decoder is a block. You can stack them several times and make the network deep!

| | | |
|---|---|---|
| **CLASS** `torch.nn.TransformerEncoder(encoder_layer, num_layers, norm=None)` | | [SOURCE] |
| **CLASS** `torch.nn.TransformerEncoderLayer(d_model, nhead, dim_feedforward=2048, dropout=0.1, activation='relu')` | | [SOURCE] |

| | | |
|---|---|---|
| **CLASS** `torch.nn.TransformerDecoder(decoder_layer, num_layers, norm=None)` | | [SOURCE] |
| **CLASS** `torch.nn.TransformerDecoderLayer(d_model, nhead, dim_feedforward=2048, dropout=0.1, activation='relu')` | | [SOURCE] |
| `forward(tgt, memory, tgt_mask=None, memory_mask=None, tgt_key_padding_mask=None, memory_key_padding_mask=None)` | | [SOURCE] |

Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. "Attention is all you need." In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 2017.