

Lecture 10:

Attention Mechanisms and Transformers

André Martins, Francisco Melo, Mário Figueiredo



Deep Learning Course, Winter 2023-2024

Today's Roadmap

Previous lecture: sequence-to-sequence models using RNNs and attention.

Today we look at **self-attention and transformers**:

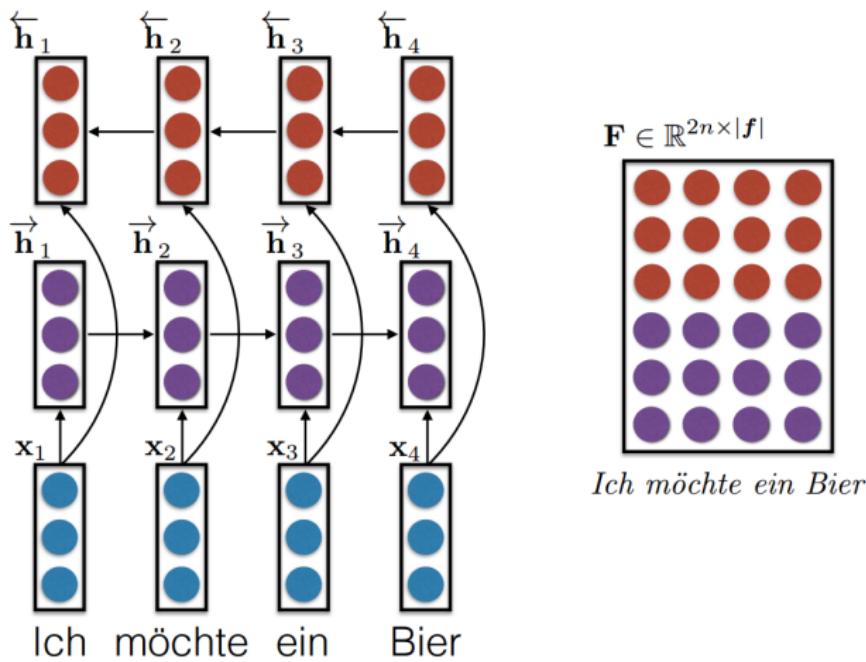
- Convolutional sequence-to-sequence models
- Self-attention
- Transformer networks
- Pre-trained models and transfer learning (next class)

Outline

- ① Convolutional Encoder-Decoder
- ② Self-Attention and Transformer Networks
- ③ Conclusions

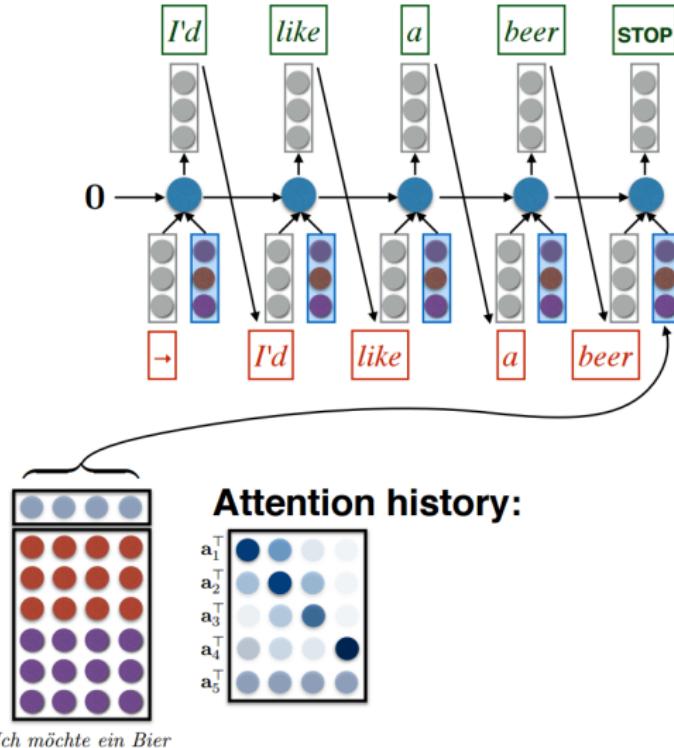
Recap: RNN with Attention (Encoder)

$$\mathbf{f}_i = [\overleftarrow{\mathbf{h}}_i; \overrightarrow{\mathbf{h}}_i]$$



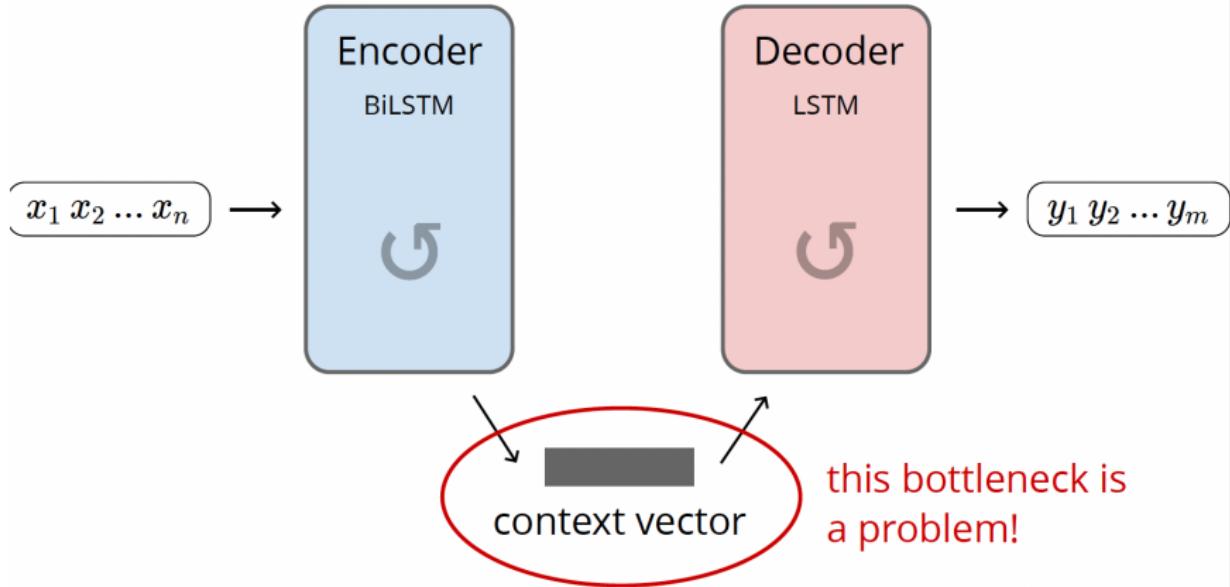
(Slide credit: Chris Dyer)

Recap: RNN with Attention (Decoder)

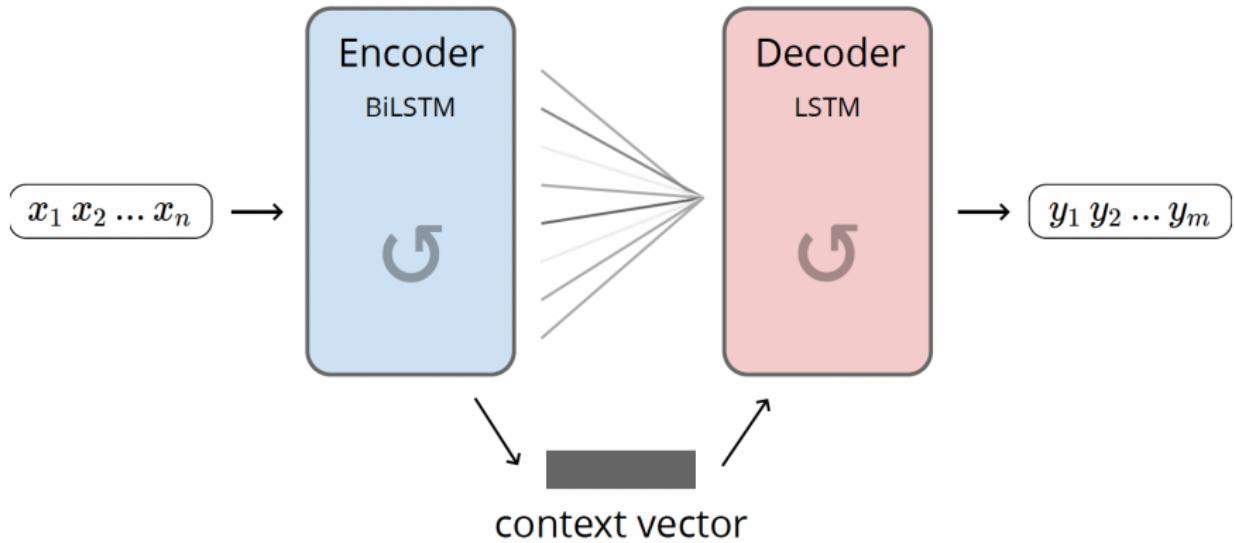


(Slide credit: Chris Dyer)

RNN-Based Encoder-Decoder



RNN-Based Encoder-Decoder

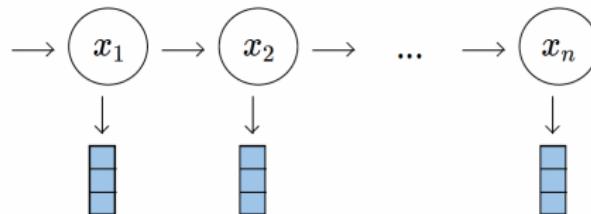


RNN-Based Encoder-Decoder

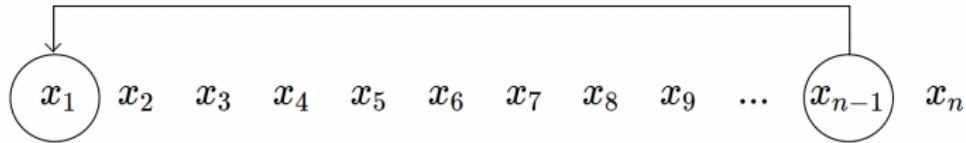


Drawbacks of RNNs

- Sequential mechanism prohibits parallelization



- Long-range dependencies are tricky, despite gating



- Possible solution: replace RNN encoder by hierarchical 1-D CNN

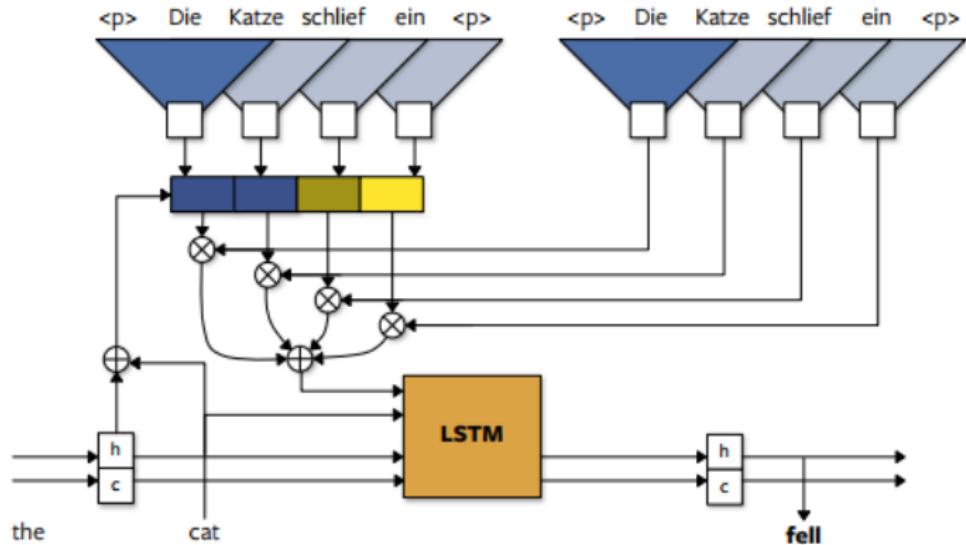
Convolutional Encoder

Convolutional
Encoder Networks

Attention Weights

Conditional
Input Computation

LSTM Decoder

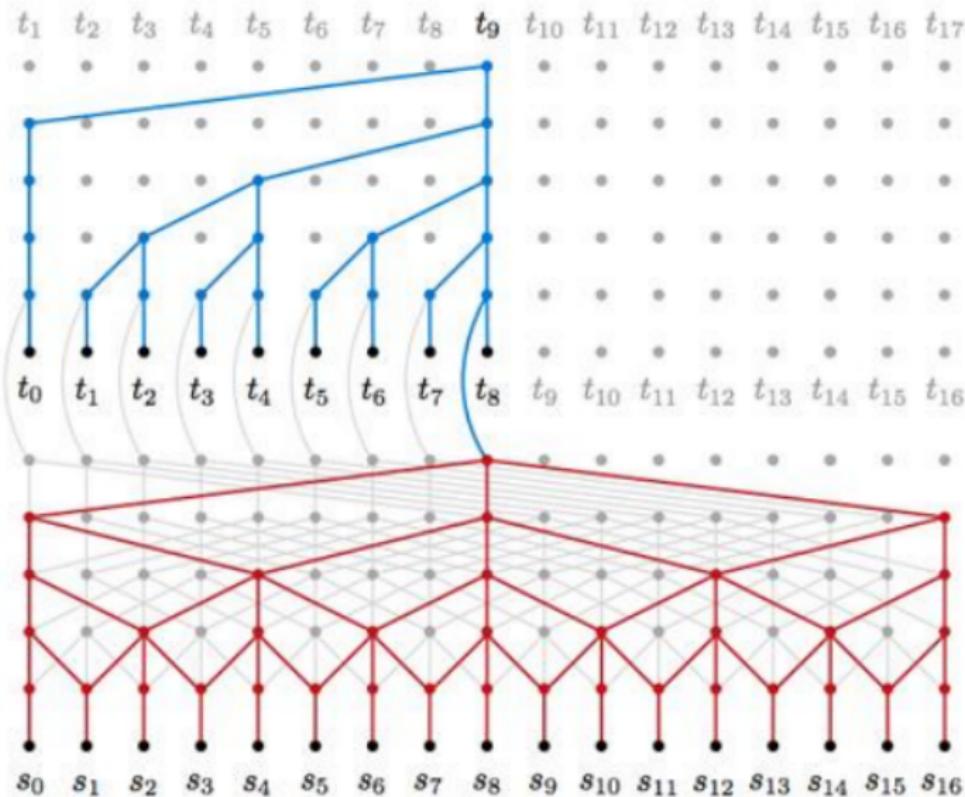


(Gehring et al., 2017)

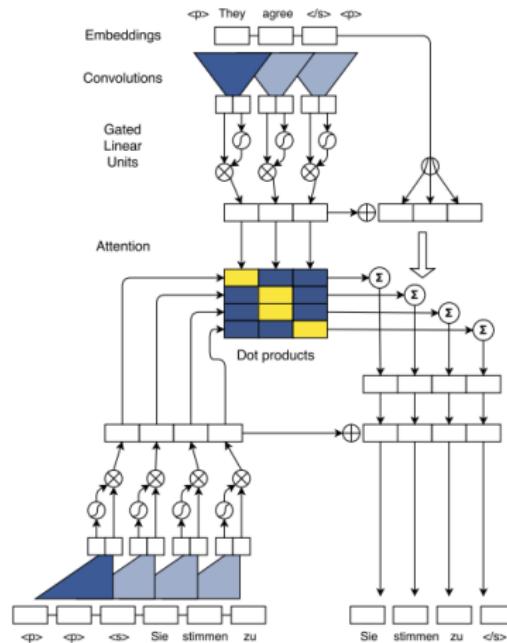
Fully Convolutional

- Can use CNN decoder too!
- Convolutions will be over output **prefixes** only
- Encoder is parallelizable, but decoder still requires sequential computation (the model is still auto-regressive)

Convolutional Sequence-to-Sequence



Convolutional Sequence-to-Sequence



(Gehring et al., 2017)

Next: Self-Attention

- Both RNN and CNN decoders require an attention mechanism
- Attention allows focusing on an arbitrary position in the source sentence, shortcircuiting the computation graph
- But if attention gives us access to any state...
...maybe we don't need the RNN?

Outline

- ① Convolutional Encoder-Decoder
- ② Self-Attention and Transformer Networks
- ③ Conclusions

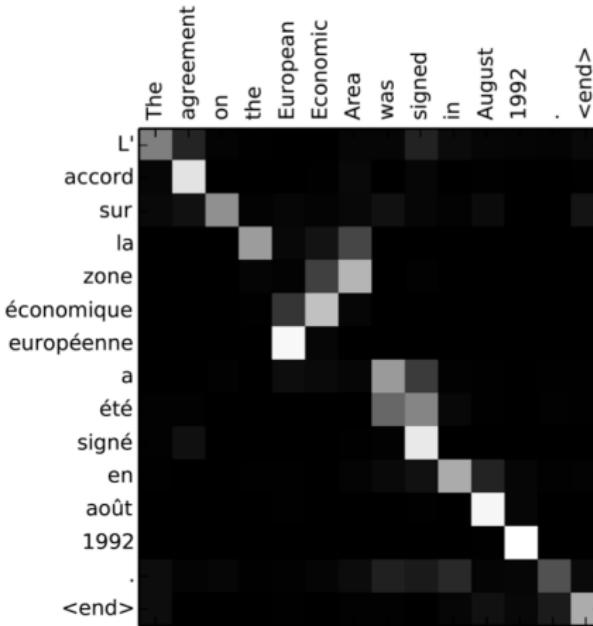
Why Attention?

We want NNs that **automatically weigh** input relevance

Main advantages:

- performance gain
- none or few parameters
- fast (easy to parallelize)
- tool for “interpreting” predictions

Example: Machine Translation



Dzmitry Bahdanau, KyungHuyn Cho, and Yoshua Bengio. Neural Machine Translation by Jointly Learning to Translate and Align. ICLR'15.

Example: Caption Generation

Attention over images:



A woman is throwing a frisbee in a park.

A dog is standing on a hardwood floor.

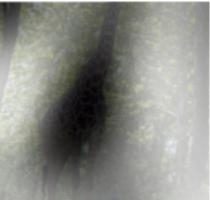
A stop sign is on a road with a mountain in the background.



A little girl sitting on a bed with a teddy bear.



A group of people sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.

(Slide credit to Yoshua Bengio)

Example: Document Classification

Task: Hotel location

you get what you pay for . not the cleanest rooms but bed was clean and so was bathroom . bring your own towels though as very thin . service was excellent , let us book in at 8:30am ! for location and price , this ca n't be beaten , but it is cheap for a reason . if you come expecting the hilton , then book the hilton ! for uk travellers , think of a blackpool b&b.

Task: Hotel cleanliness

you get what you pay for . not the cleanest rooms but bed was clean and so was bathroom . bring your own towels though as very thin . service was excellent , let us book in at 8:30am ! for location and price , this ca n't be beaten , but it is cheap for a reason . if you come expecting the hilton , then book the hilton ! for uk travellers , think of a blackpool b&b.

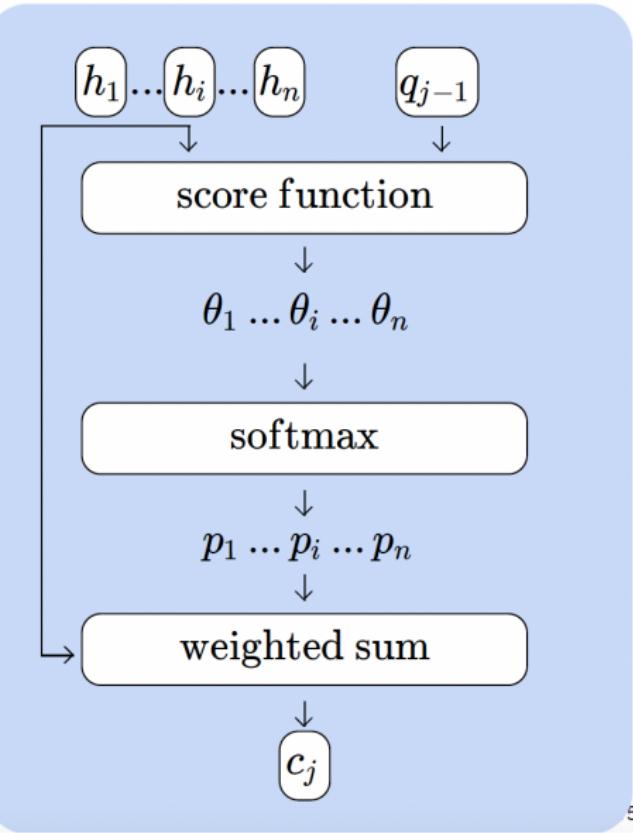
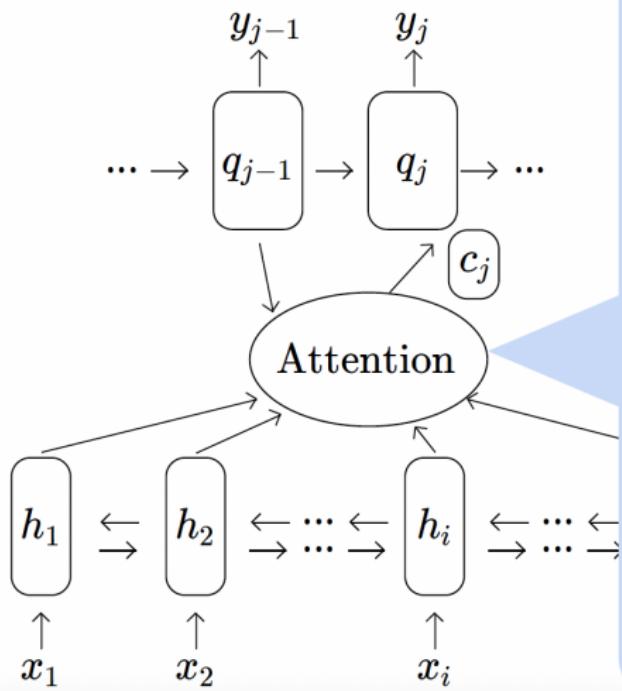
Task: Hotel service

you get what you pay for . not the cleanest rooms but bed was clean and so was bathroom . bring your own towels though as very thin . service was excellent , let us book in at 8:30am ! for location and price , this ca n't be beaten , but it is cheap for a reason . if you come expecting the hilton , then book the hilton ! for uk travellers , think of a blackpool b&b.

(Bao et al., 2018)

Attention Mechanism

- Bahdanau et al. (2015)



Attention Mechanism: Recap

Recall how attention works:

- ① We have a **query vector** \mathbf{q} (e.g. the decoder state)
- ② We have **input vectors** $\mathbf{H} = [\mathbf{h}_1, \dots, \mathbf{h}_L]^\top$ (e.g. one per source word)
- ③ We compute **affinity scores** s_1, \dots, s_L by “comparing” \mathbf{q} and \mathbf{H}
- ④ We convert these scores to **probabilities**:

$$\mathbf{p} = \text{softmax}(\mathbf{s})$$



- ⑤ We use this to output a representation as a **weighted average**:

$$\mathbf{c} = \mathbf{H}^\top \mathbf{p} = \sum_{i=1}^L p_i \mathbf{h}_i$$

Let's see these steps in detail!

Affinity Scores

Several ways of “comparing” a query \mathbf{q} and an input (“key”) vector \mathbf{h}_i :

- **Additive attention** (Bahdanau et al., 2015), what we covered in previous class:

$$s_i = \mathbf{u}^\top \tanh(\mathbf{A}\mathbf{h}_i + \mathbf{B}\mathbf{q})$$

- **Bilinear attention** (Luong et al., 2015):

$$s_i = \mathbf{q}^\top \mathbf{U}\mathbf{h}_i$$

- **Dot product attention** (Luong et al., 2015) (particular case; queries and keys must have the same size):

$$s_i = \mathbf{q}^\top \mathbf{h}_i$$

The last two are easier to batch when we have multiple queries and multiple keys.

Keys and Values

The input vectors $\mathbf{H} = [\mathbf{h}_1, \dots, \mathbf{h}_L]^\top$ appear in two places:

- They are used as **keys** to “compare” them with the query vector \mathbf{q} to obtain the affinity scores
- They are used as **values** to form the weighted average $\mathbf{c} = \mathbf{H}^\top \mathbf{p}$

To be fully general, **they don't need to be the same** – we can have:

- A **key matrix** $\mathbf{K} = [\mathbf{k}_1, \dots, \mathbf{k}_L]^\top \in \mathbb{R}^{L \times d_K}$
- A **value matrix** $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_L]^\top \in \mathbb{R}^{L \times d_V}$

Attention Mechanism: More General Version

- ① We have a **query vector** \mathbf{q} (e.g. the decoder state)
- ② We have **key vectors** $\mathbf{K} = [\mathbf{k}_1, \dots, \mathbf{k}_L]^\top \in \mathbb{R}^{L \times d_K}$ and **value vectors** $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_L]^\top \in \mathbb{R}^{L \times d_V}$ (one of each per source word)
- ③ We compute **query-key affinity scores** s_1, \dots, s_L “comparing” \mathbf{q} and \mathbf{K}
- ④ We convert these scores to **probabilities**:

$$\mathbf{p} = \text{softmax}(\mathbf{s})$$

- ⑤ We output a weighted average of the **values**:

$$\mathbf{c} = \mathbf{V}^\top \mathbf{p} = \sum_{i=1}^L p_i \mathbf{v}_i \in \mathbb{R}^{d_V}$$

Self-Attention

- So far we talked about **contextual** attention – the decoder attends to encoder states (this is called “input context”)
- The encoder and the decoder states were propagated sequentially with a RNN, or hierarchically with a CNN
- Alternative: **self-attention** – at each position, the encoder attends to the other positions in the encoder itself
- Same for the decoder.

Self-Attention Layer

Self-attention for a sequence of length L :

- ① **Query vectors** $\mathbf{Q} = [\mathbf{q}_1, \dots, \mathbf{q}_L]^\top \in \mathbb{R}^{L \times d_Q}$
- ② **Key vectors** $\mathbf{K} = [\mathbf{k}_1, \dots, \mathbf{k}_L]^\top \in \mathbb{R}^{L \times d_K}$
- ③ **value vectors** $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_L]^\top \in \mathbb{R}^{L \times d_V}$
- ④ Compute **query-key** affinity scores “comparing” \mathbf{Q} and \mathbf{K} , e.g.,

$$\mathbf{S} = \mathbf{Q}\mathbf{K}^\top \in \mathbb{R}^{L \times L} \quad (\text{dot-product affinity})$$

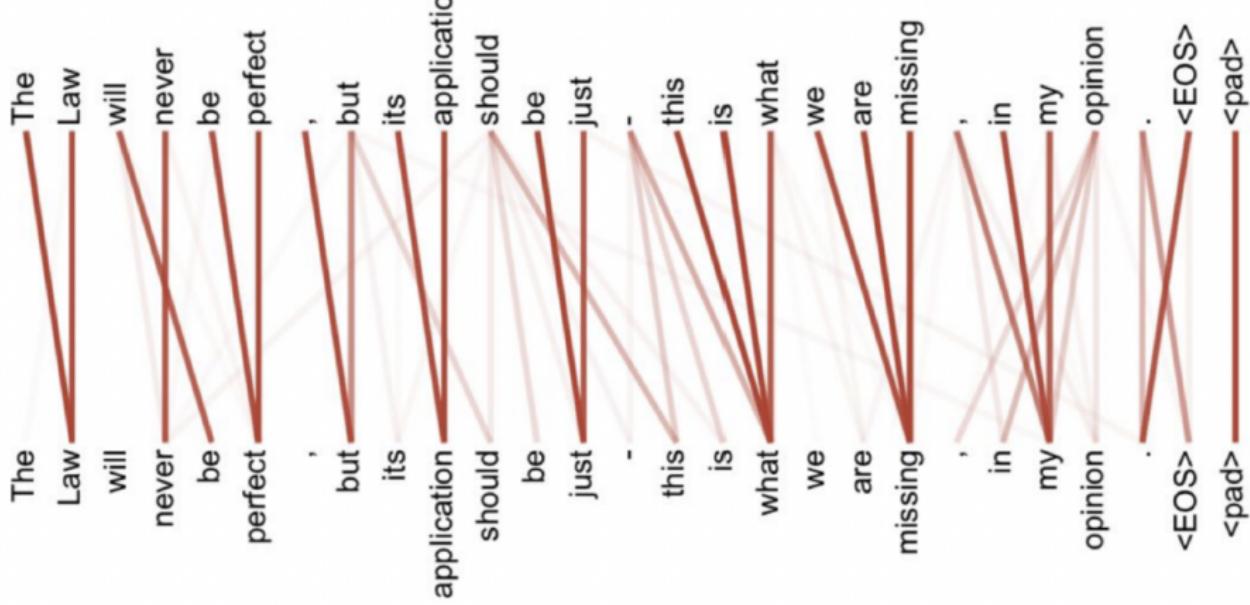
- ⑤ Convert these scores to **probabilities** (row-wise):

$$\mathbf{P} = \text{softmax}(\mathbf{S}) \in \mathbb{R}^{L \times L}$$

- ⑥ Output the weighted average of the **values**:

$$\mathbf{Z} = \mathbf{PV} = \underbrace{\text{softmax}(\mathbf{Q}\mathbf{K}^\top)}_P \mathbf{V} \in \mathbb{R}^{L \times d_V}.$$

Self-Attention



(Vaswani et al., 2017)

Transformer (Vaswani et al., 2017)

- **Key idea:** instead of RNN/CNNs, use **self-attention** in the encoder
- Each word state attends to all the other words
- Each self-attention is followed by a feed-forward transformation
- Do several layers of this
- Do the same for the decoder, attending only to already generated words.

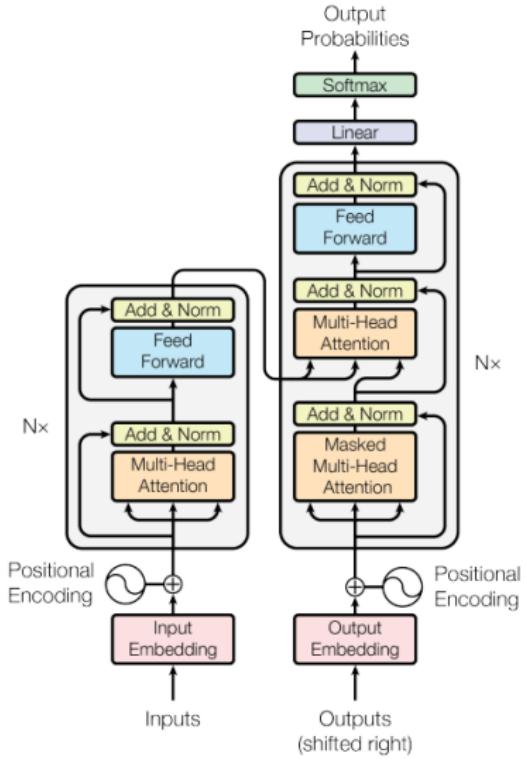
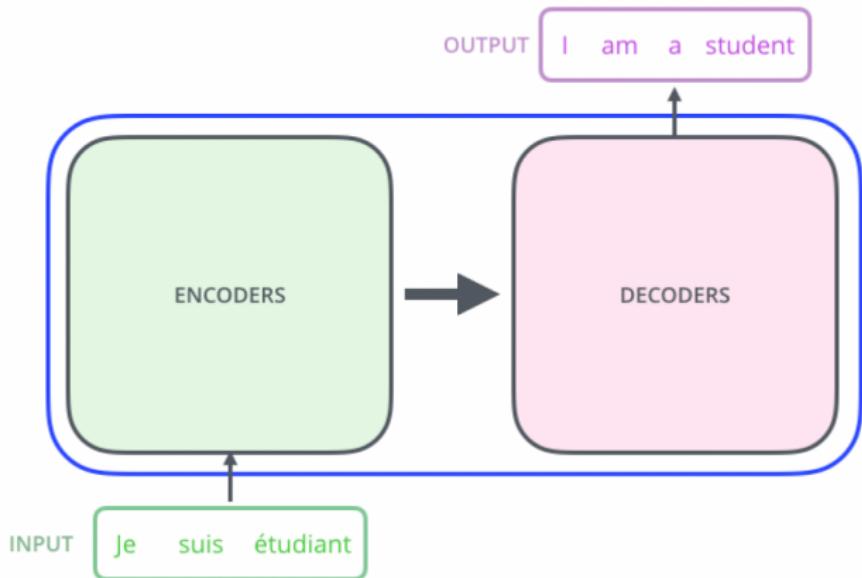


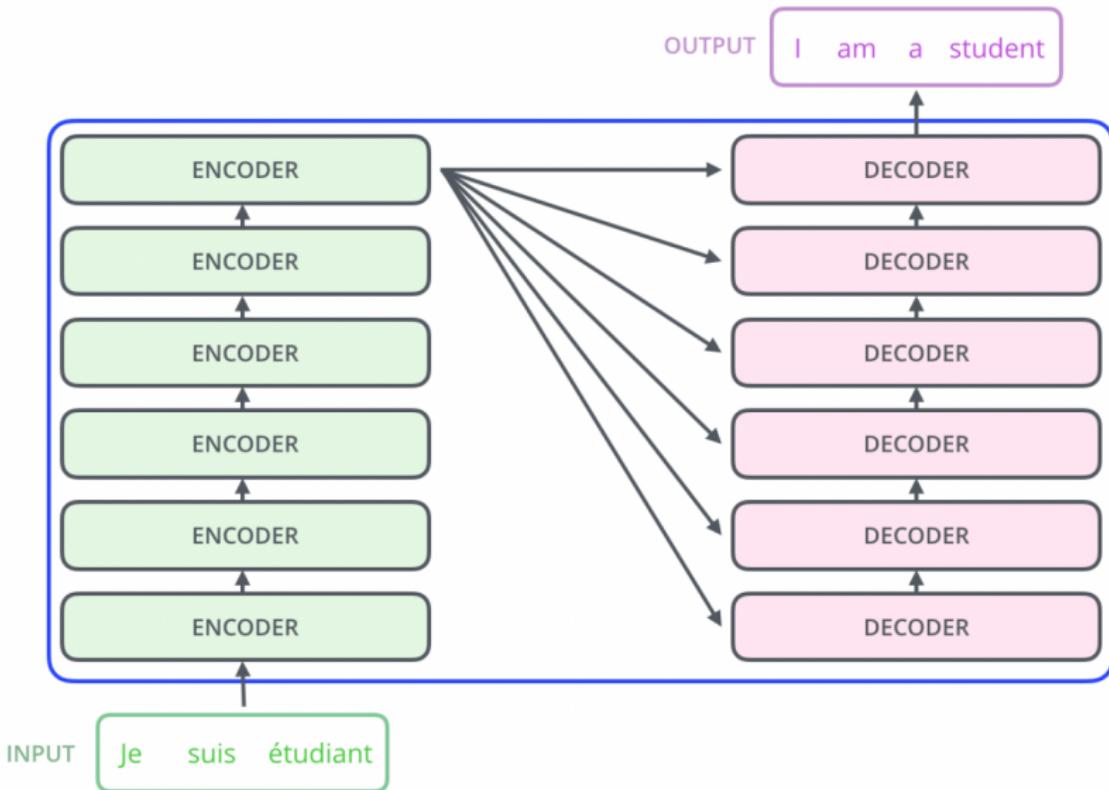
Figure 1: The Transformer - model architecture.

Transformer

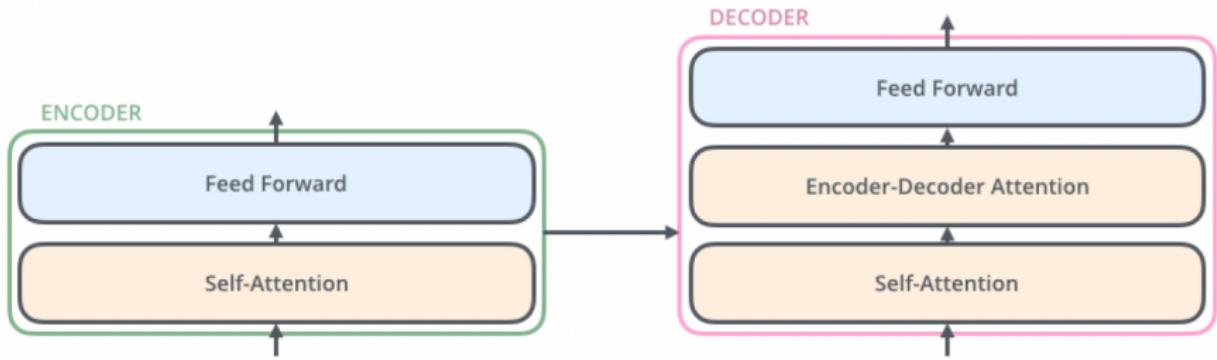


$$x_1 \ x_2 \dots x_n \xrightarrow{\text{encode}} \mathbf{r}_1 \ \mathbf{r}_2 \dots \mathbf{r}_n \xrightarrow{\text{decode}} y_1 \ y_2 \dots y_m$$

Transformer



Transformer Blocks

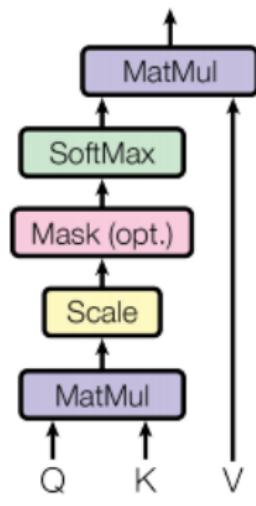


(Illustrated transformer: <http://jalammar.github.io/illustrated-transformer/>)

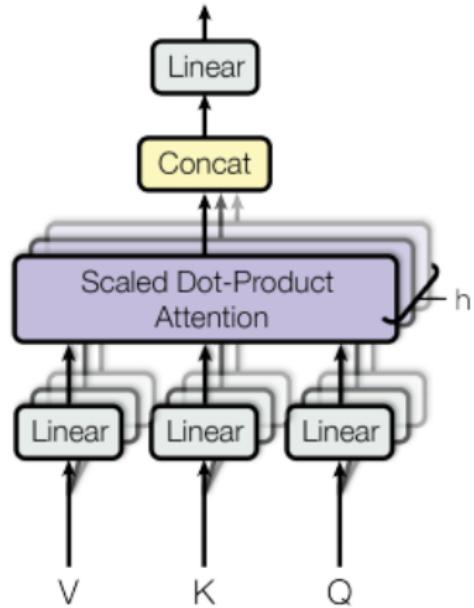
Scaled Dot-Product and Multi-Head Attention

Two innovations: scaled dot-product attention; multi-head attention.

Scaled Dot-Product Attention



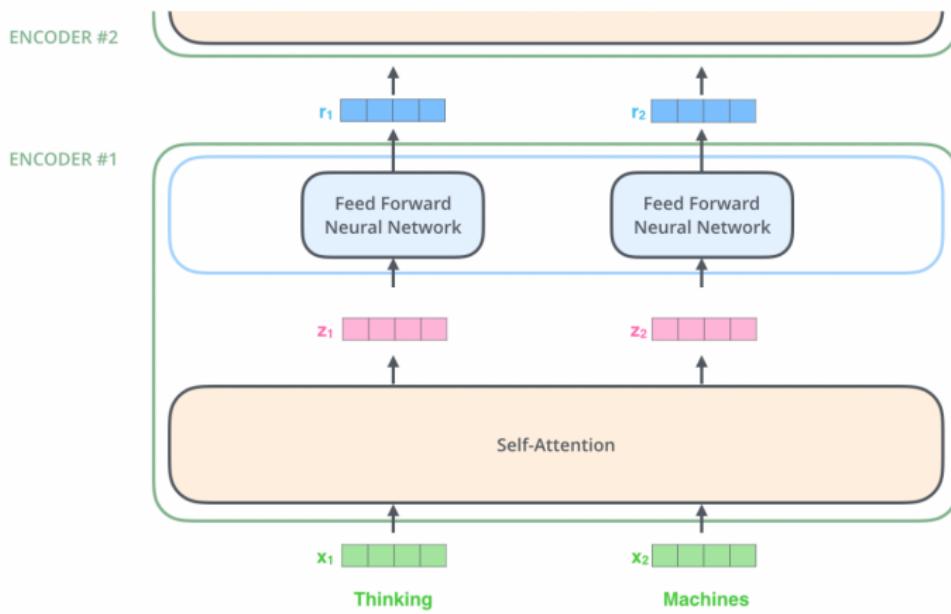
Multi-Head Attention



(Vaswani et al., 2017)

The Encoder

Example for a sentence with 2 words:



Transformer Self-Attention: Queries, Keys, Vectors

- Obtained by projecting the embedding matrix $\mathbf{X} \in \mathbb{R}^{L \times e}$ to a lower dimension:

$$\mathbf{Q} = \mathbf{XW}^Q$$

$$\mathbf{K} = \mathbf{XW}^K$$

$$\mathbf{V} = \mathbf{XW}^V.$$

$$\mathbf{X} \times \mathbf{W}^Q = \mathbf{Q}$$

$$\mathbf{X} \times \mathbf{W}^K = \mathbf{K}$$

$$\mathbf{X} \times \mathbf{W}^V = \mathbf{V}$$

- The projection matrices \mathbf{W}^Q , \mathbf{W}^K , \mathbf{W}^V are model parameters.

Transformer Self-Attention: Queries, Keys, Vectors

Input

Thinking

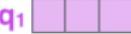
Machines

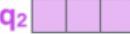
Embedding

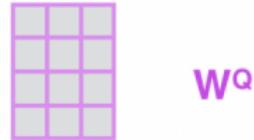
X_1 

X_2 

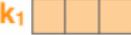
Queries

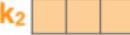
q_1 

q_2 



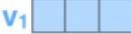
Keys

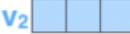
k_1 

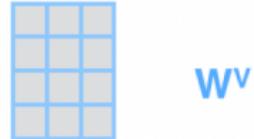
k_2 



Values

v_1 

v_2 



Scaled Dot-Product Attention

Problem: As d_K gets large, the variance of $\mathbf{q}^\top \mathbf{k}$ increases, the softmax gets very peaked, hence its gradient gets smaller.

Solution: scale by length of query/key vectors:

$$\mathbf{Z} = \text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_K}} \right) \mathbf{V}.$$

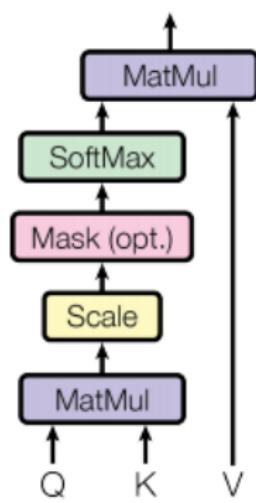
Scaled Dot-Product Attention

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

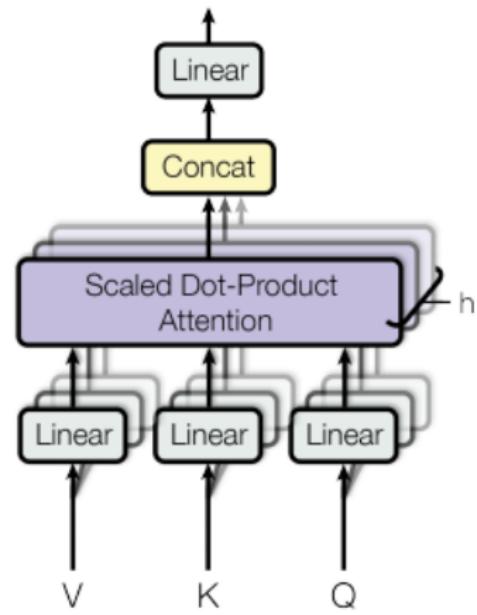
The diagram illustrates the computation of scaled dot-product attention. It shows three input matrices: \mathbf{Q} (purple, 3x3), \mathbf{K}^T (orange, 3x3), and \mathbf{V} (blue, 3x3). The \mathbf{Q} and \mathbf{K}^T matrices are multiplied together, and the result is divided by the square root of the dimension d_k . This result is then passed through a softmax function to produce the output matrix \mathbf{Z} (pink, 3x3).

Scaled Dot-Product and Multi-Head Attention

Scaled Dot-Product Attention



Multi-Head Attention



(Vaswani et al., 2017)

Multi-Head Attention

Self-attention: each word forms a **query vector** and attends to the other words' **key vectors**

Somewhat similar to a **1D convolution**, with “adaptive” weights and the window size spans the entire sentence

Problem: only one channel for words to interact with one-another

Solution: **multi-head attention!**

- define h attention heads, each with their own projection matrices (e.g. $h = 8$)
- apply attention in multiple channels, concatenate the outputs and pipe through linear layer:

$$\text{MultiHead}(\mathbf{X}) = \text{Concat}(\mathbf{Z}_1, \dots, \mathbf{Z}_h) \mathbf{W}^O,$$

where $\mathbf{Z}_i = \text{Attention}(\underbrace{\mathbf{X}\mathbf{W}_i^Q}_{\mathbf{Q}_i}, \underbrace{\mathbf{X}\mathbf{W}_i^K}_{\mathbf{K}_i}, \underbrace{\mathbf{X}\mathbf{W}_i^V}_{\mathbf{V}_i}).$

Multi-Head Attention

1) This is our input sentence* each word*

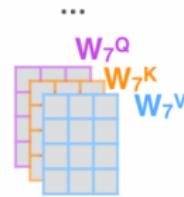
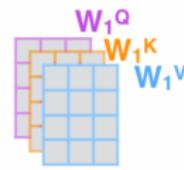
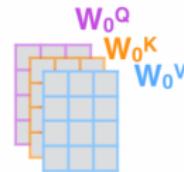


* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one

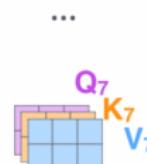


2) We embed each word*

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



W^o



Other Tricks

- Self-attention blocks are repeated several times (e.g. 6 or 12)
- Residual connections on each attention block
- Layer normalization
- Positional encodings (to distinguish word positions)

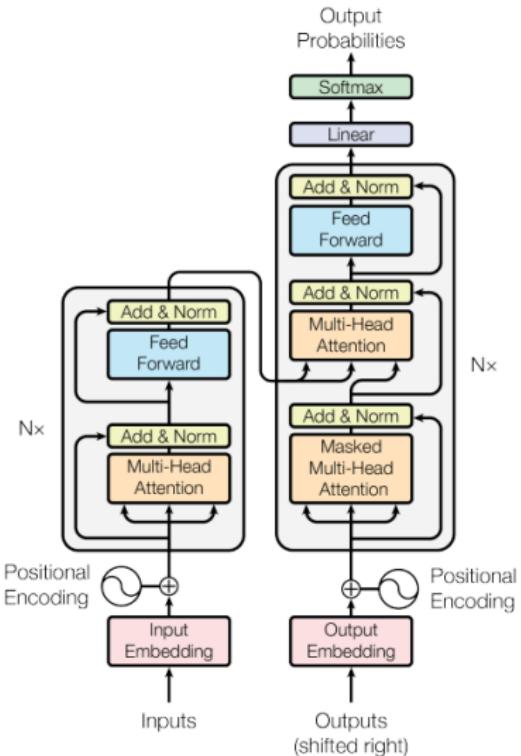
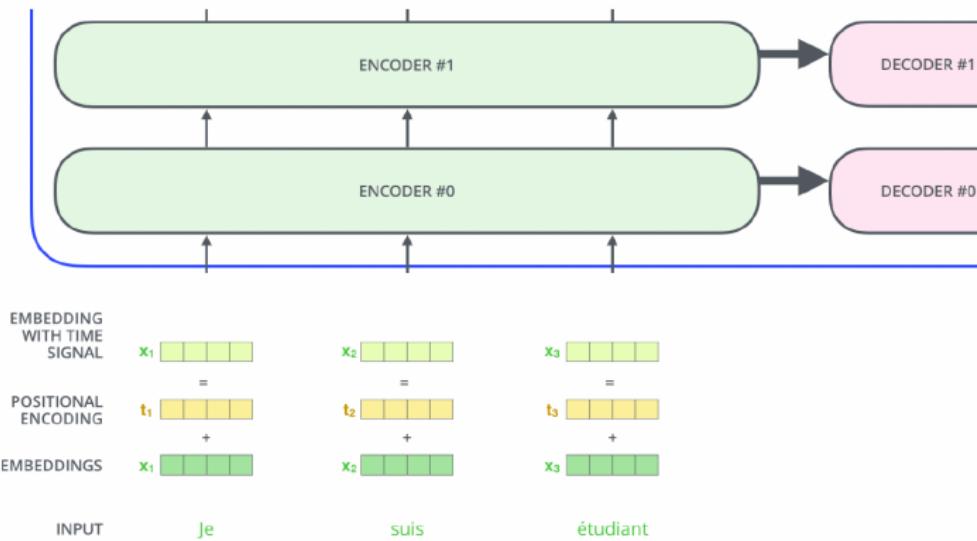


Figure 1: The Transformer - model architecture.

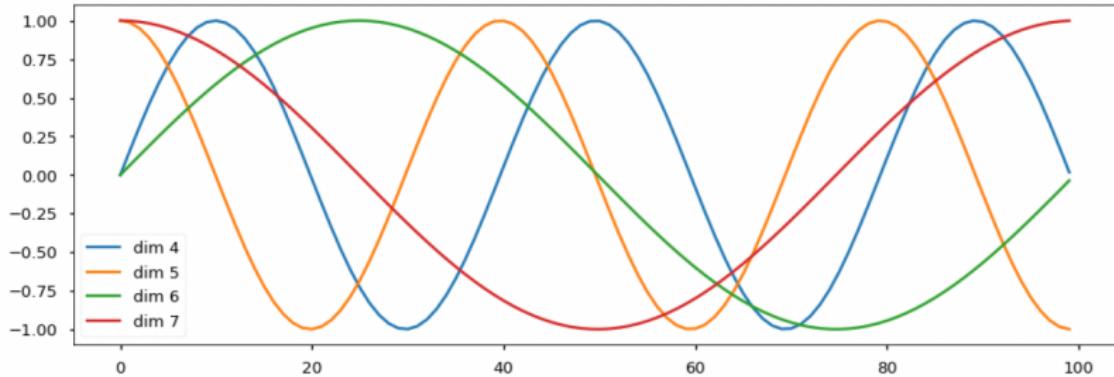
Positional Encodings

- As just described, the transformer is insensitive to word order!
 - queries attend to keys regardless of their position in the sequence
- To make it sensitive to order, we add **positional encodings**
- Two strategies: learn one embedding for each position (up to a maximum length) or use sinusoidal positional encodings (next)

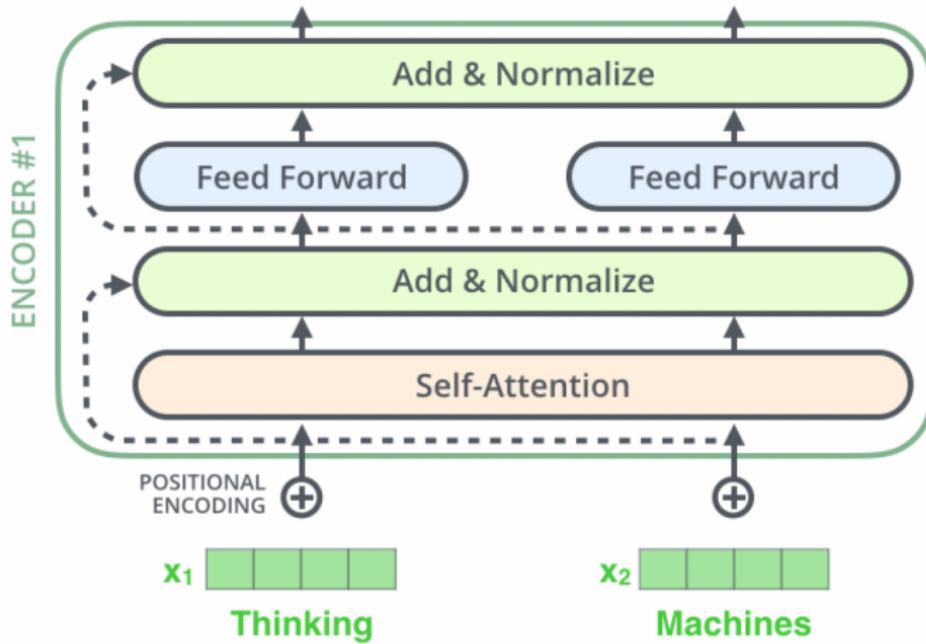


Sinusoidal Positional Encodings

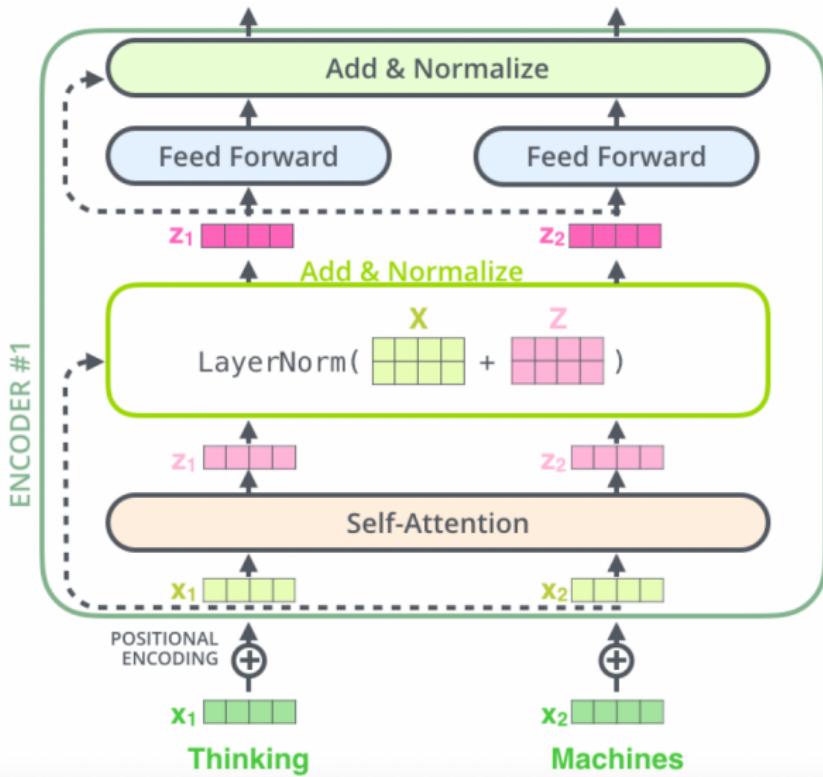
$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{2i/d}}\right) \quad PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000^{2i/d}}\right)$$



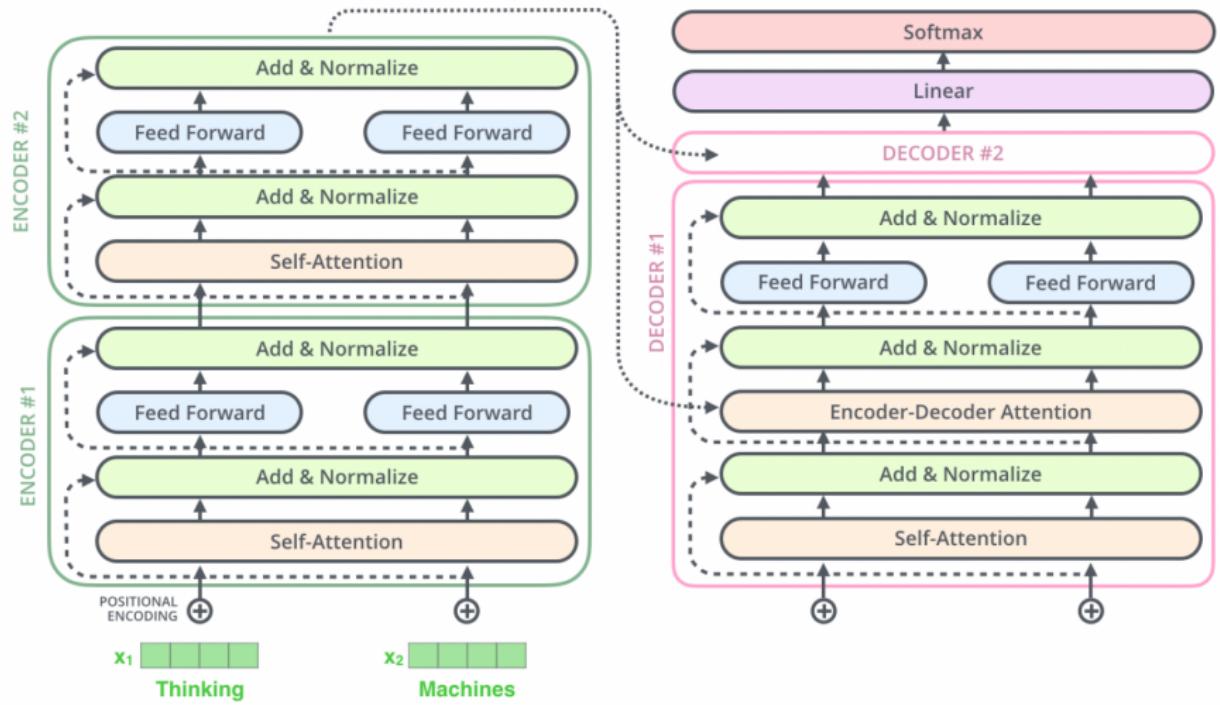
Residuals and Layer Normalization



Residuals and Layer Normalization



Residuals and Layer Normalization



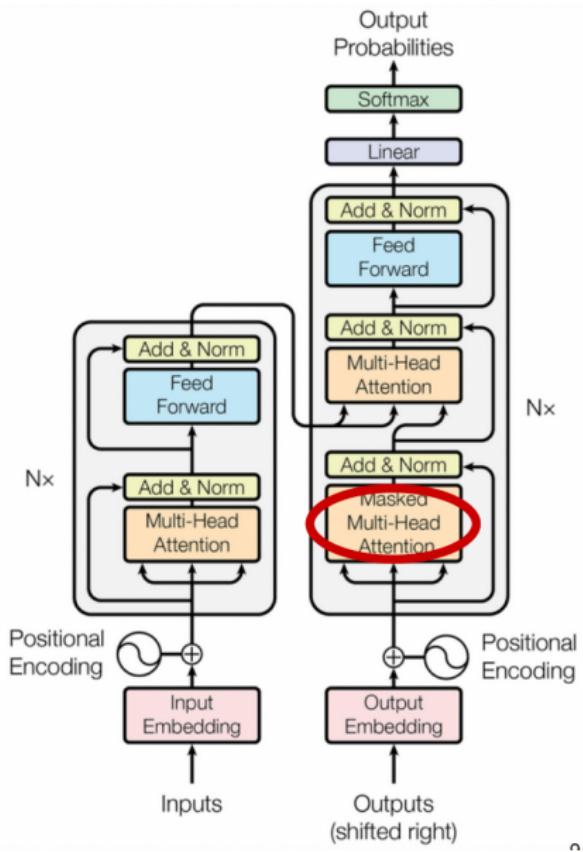
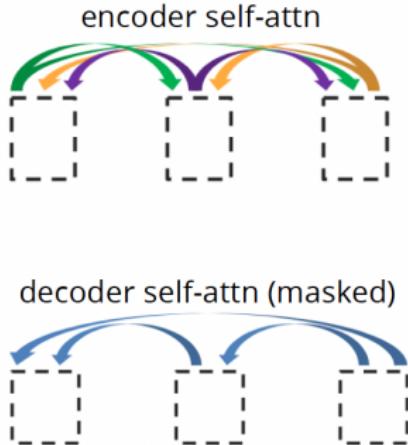
The Decoder

What about the self-attention blocks in the **decoder**?

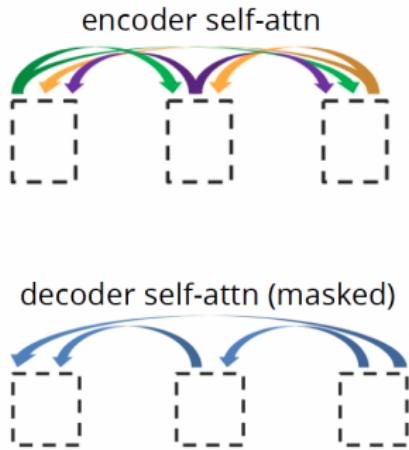
Everything is pretty much the same as in the encoder, with two twists:

- The decoder cannot see the future! Use “**causal**” masking
- The decoder should attend to itself (**self-attention**), but also to the encoder states (**contextual attention**).

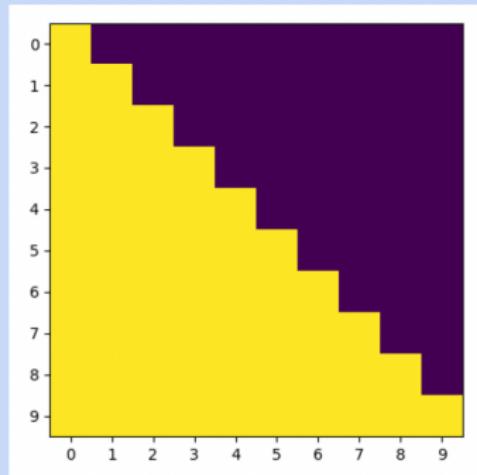
The Decoder



The Decoder



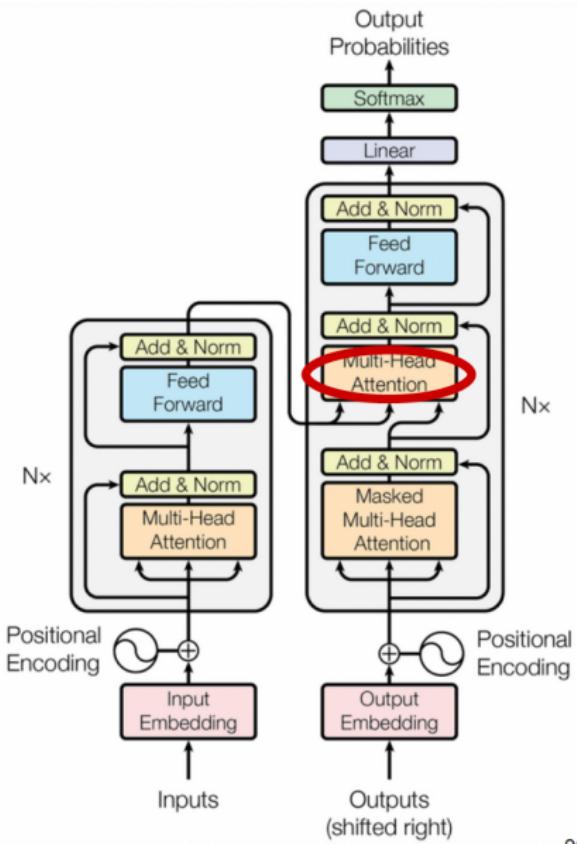
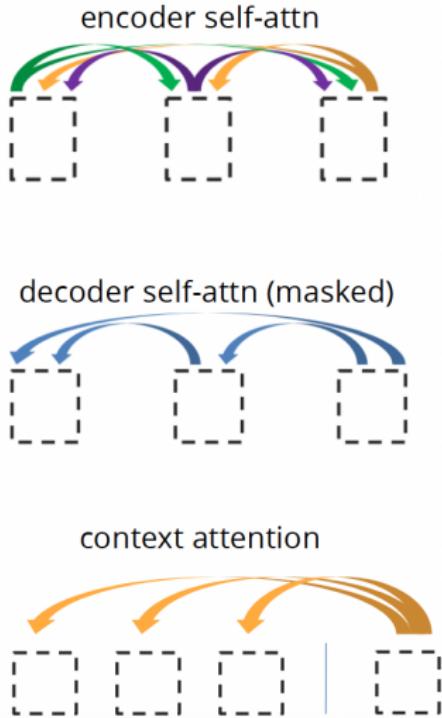
- Mask subsequent positions (before softmax)



- In PyTorch

```
scores.masked_fill_(~mask, float('-inf'))
```

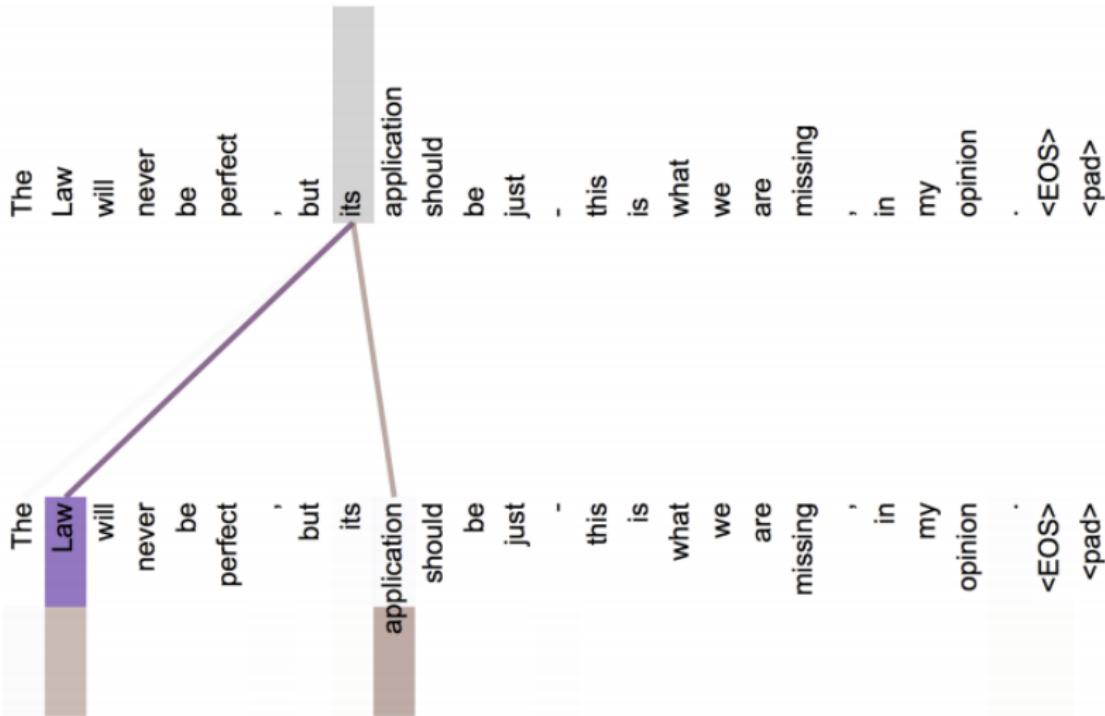
The Decoder



Attention Visualization Layer 5

The figure displays a neural network architecture for processing the sentence "It is in this spirit that a majority of American governments have passed new laws since 2009 making the registration or voting process more difficult". The words are represented as vectors in a 2D space, with arrows indicating their relationships. The word "making" is highlighted with a large purple box and has multiple outgoing arrows pointing to other words: "the" (dark brown), "registration" (light blue), "or" (light green), "voting" (pink), "process" (orange), "more" (red), and "difficult" (green). Other words like "spirit", "that", "a", "majority", "of", "American", "governments", "have", "passed", "new", "laws", "since", "2009", "the", "or", "voting", "process", "more", "difficult", and punctuation like "" and "<pad>" are also shown as vectors in the same 2D space.

Implicit Anaphora Resolution



Computational Cost

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)	$O(r \cdot n \cdot d)$	$O(1)$	$O(n/r)$

n = seq. length

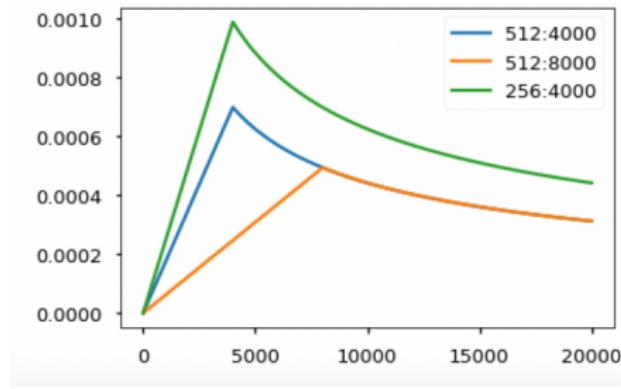
d = hidden dim

k = kernel size

- Faster to train (due to self-attention parallelization)
- More expensive to decode
- Scale quadratically with respect to sequence length (problematic for long sequences).

Other Tricks

- Label smoothing (e.g., replace one-hot encoding by smoothed version)
- Dropout at every layer before residuals
- Beam search with length penalty
- Adam optimizer with learning-rate decay



Overall, transformers are harder to optimize than RNN seq2seq models
They don't work out of the box: hyperparameter tuning is very important.

Transformer Results

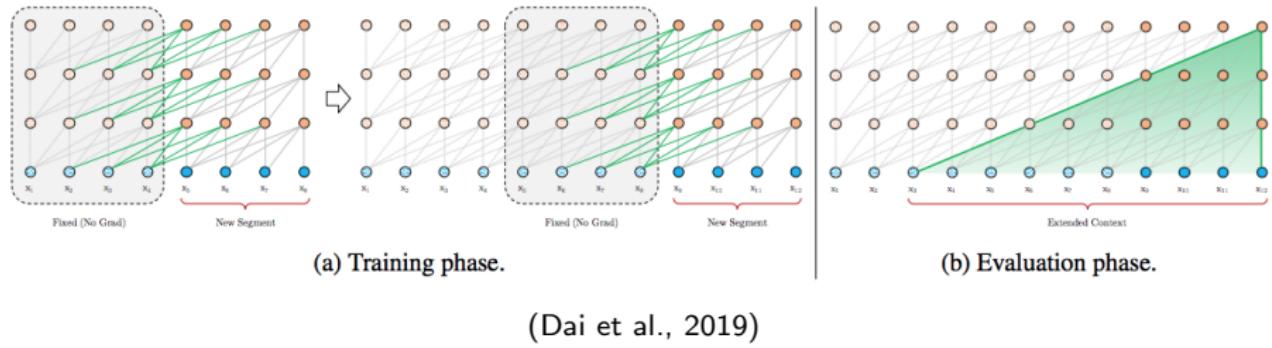
Model	BLEU		Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet [18]	23.75			
Deep-Att + PosUnk [39]		39.2		$1.0 \cdot 10^{20}$
GNMT + RL [38]	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4 \cdot 10^{20}$
ConvS2S [9]	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$
MoE [32]	26.03	40.56	$2.0 \cdot 10^{19}$	$1.2 \cdot 10^{20}$
Deep-Att + PosUnk Ensemble [39]		40.4		$8.0 \cdot 10^{20}$
GNMT + RL Ensemble [38]	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$
ConvS2S Ensemble [9]	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	41.8		$2.3 \cdot 10^{19}$

(Vaswani et al., 2017)'s "Attention Is All You Need"

TransformerXL

Big transformers can look at larger contexts.

TransformerXL: enables going beyond a fixed length without disrupting temporal coherence:



Outline

- ① Convolutional Encoder-Decoder
- ② Self-Attention and Transformer Networks
- ③ Conclusions

Conclusions

- RNN-based seq2seq models require sequential computation and have difficulties with long range dependencies
- Attention mechanisms allow focusing on different parts of the input
- Encoders/decoders can be RNNs, CNNs, or self-attention layers
- Transformers are the current state of the art for many tasks in NLP and vision
- Other applications: speech recognition, image captioning, etc.
- Next lecture: pretrained models and transfer learning (BERT, GPT-2, GPT-3, etc.)

Pointers for Today's Class

- Lena Voita's seq2seq with attention: https://lena-voita.github.io/nlp_course/seq2seq_and_attention.html
- Marcos Treviso lecture on attention mechanisms:
<https://andre-martins.github.io/docs/dsl2020/attention-mechanisms.pdf>
- John Hewitt's lecture on self-attention and transformers:
<http://web.stanford.edu/class/cs224n/slides/cs224n-2021-lecture09-transformers.pdf>
- Illustrated transformer:
<http://jalammar.github.io/illustrated-transformer/>
- Annotated transformer:
<https://nlp.seas.harvard.edu/2018/04/03/attention.html>

Thank you!

Questions?



References I

- Bahdanau, D., Cho, K., and Bengio, Y. (2015). Neural Machine Translation by Jointly Learning to Align and Translate. In *International Conference on Learning Representations*.
- Bao, Y., Chang, S., Yu, M., and Barzilay, R. (2018). Deriving machine attention from human rationales. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1903–1913.
- Dai, Z., Yang, Z., Yang, Y., Cohen, W. W., Carbonell, J., Le, Q. V., and Salakhutdinov, R. (2019). Transformer-xl: Attentive language models beyond a fixed-length context. *arXiv preprint arXiv:1901.02860*.
- Gehring, J., Auli, M., Grangier, D., Yarats, D., and Dauphin, Y. N. (2017). Convolutional sequence to sequence learning. *arXiv preprint arXiv:1705.03122*.
- Luong, M.-T., Pham, H., and Manning, C. D. (2015). Effective approaches to attention-based neural machine translation. *arXiv preprint arXiv:1508.04025*.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L. u., and Polosukhin, I. (2017). Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 5998–6008.