

Lecture 10: Attention Mechanisms and Transformers

André Martins, Francisco Melo, Mário Figueiredo



Deep Learning Course, Winter 2022-2023

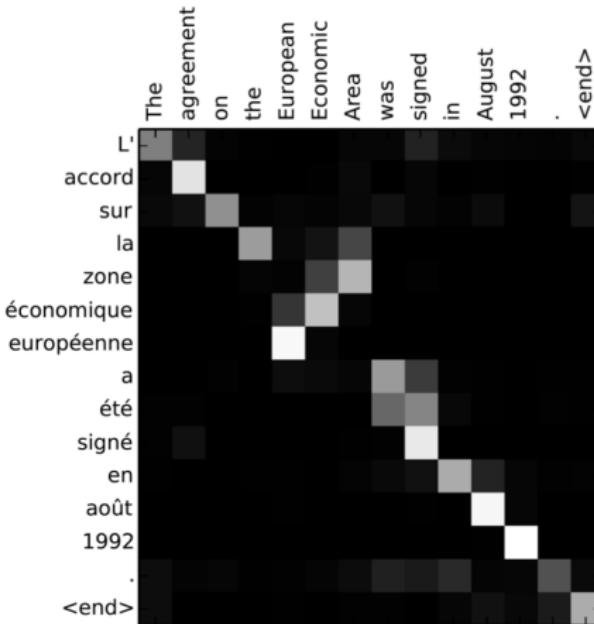
Why Attention?

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:



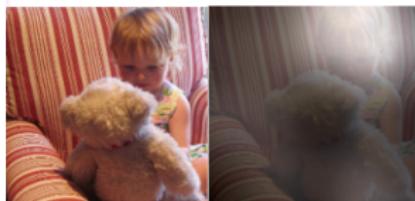
A woman is throwing a frisbee in a park.



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: 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!

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** $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_L]^\top \in \mathbb{R}^{L \times d_V}$ (e.g. 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 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.,

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

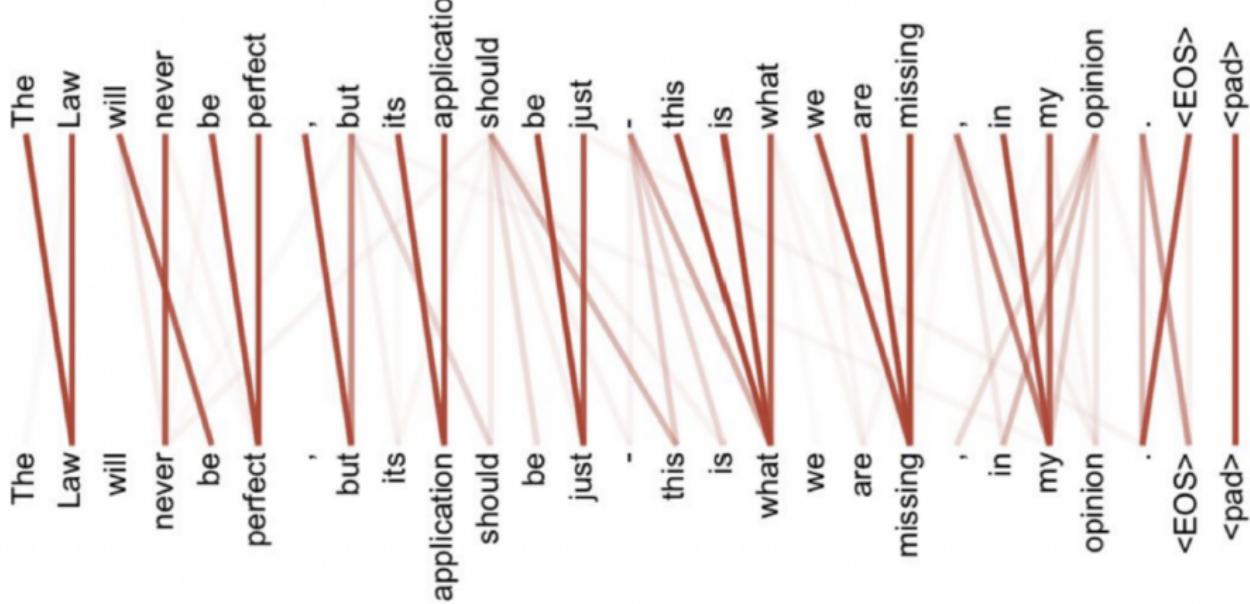
- ⑤ Convert these scores to **probabilities** (row-wise):

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

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

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

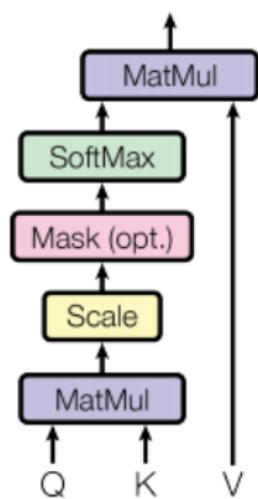
Self-Attention



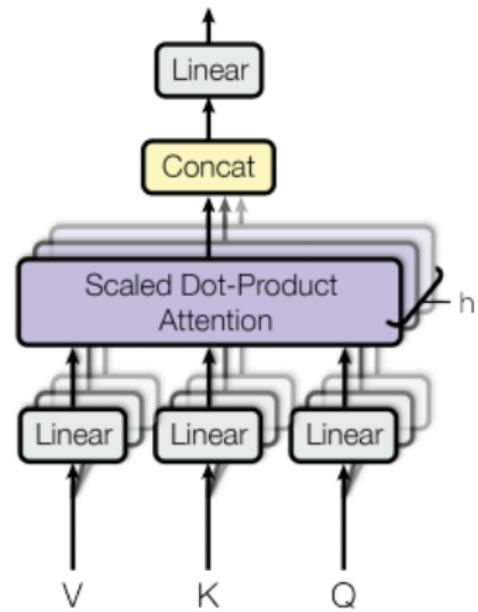
(Vaswani et al., 2017)

Scaled Dot-Product and Multi-Head Attention

Scaled Dot-Product Attention



Multi-Head Attention



(Vaswani et al., 2017)

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.

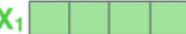
Transformer Self-Attention: Queries, Keys, Vectors

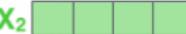
Input

Thinking

Machines

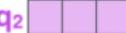
Embedding

X_1 

X_2 

Queries

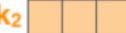
q_1 

q_2 



Keys

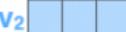
k_1 

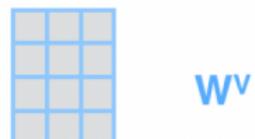
k_2 



Values

v_1 

v_2 



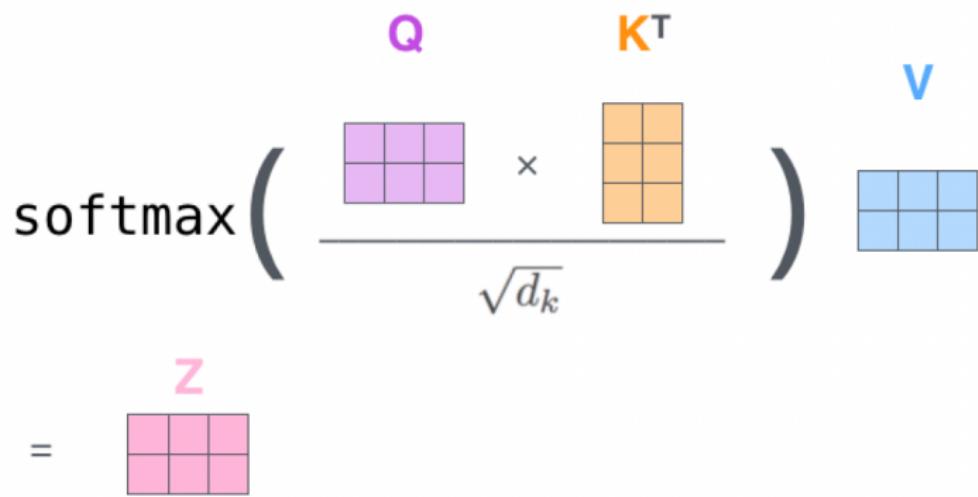
Scaled Dot-Product Attention

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:

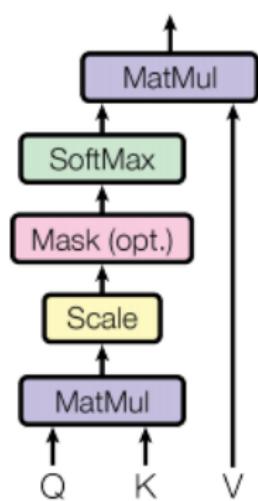
$$\mathbf{Z} = \text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_K}} \right) \mathbf{V}.$$

Scaled Dot-Product Attention

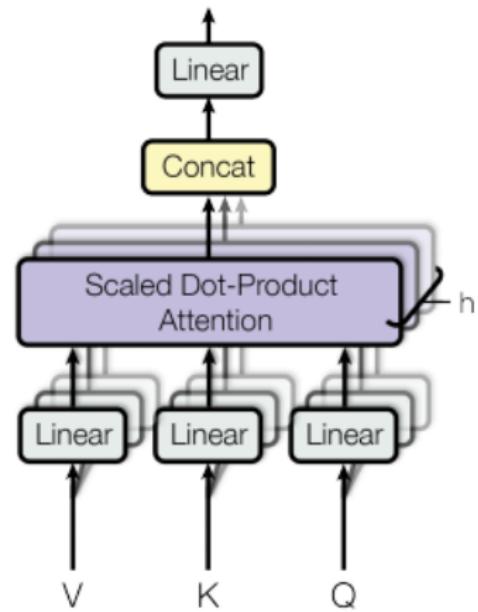
$$\text{softmax} \left(\frac{\begin{matrix} \mathbf{Q} & \mathbf{K}^T \\ \begin{matrix} \times \end{matrix} & \end{matrix}}{\sqrt{d_k}} \right) \mathbf{V}$$
$$= \mathbf{Z}$$


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**

This is vaguely similar to a **1D convolution**, but where the filter weights are “dynamic” is 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 $\mathbf{Z}_i = \text{Attention}(\underbrace{\mathbf{X}\mathbf{W}_i^Q}_{\mathbf{Q}_i}, \underbrace{\mathbf{X}\mathbf{W}_i^K}_{\mathbf{K}_i}, \underbrace{\mathbf{X}\mathbf{W}_i^V}_{\mathbf{V}_i}).$

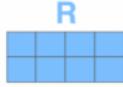
Multi-Head Attention

1) This is our input sentence*
2) We embed each word*

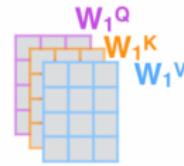
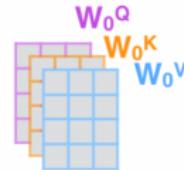
Thinking
Machines



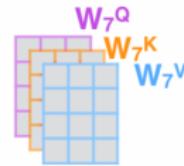
* In all encoders other than #0, we don't need embedding.
We start directly with the output of the encoder right below this one



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)

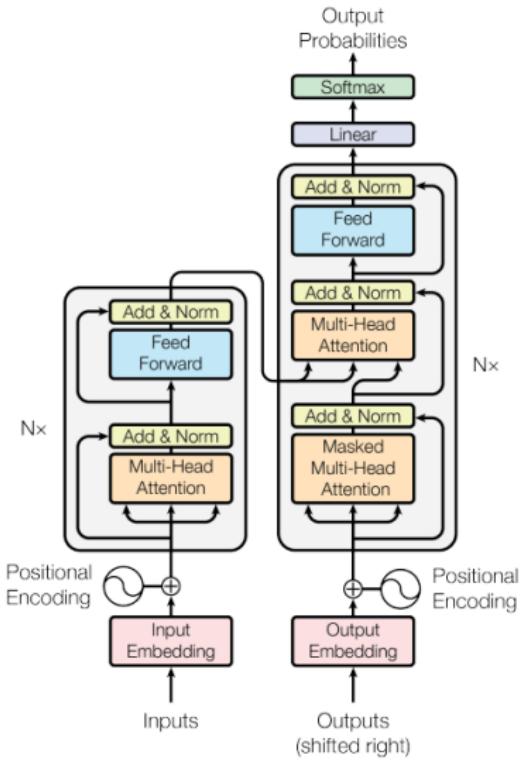
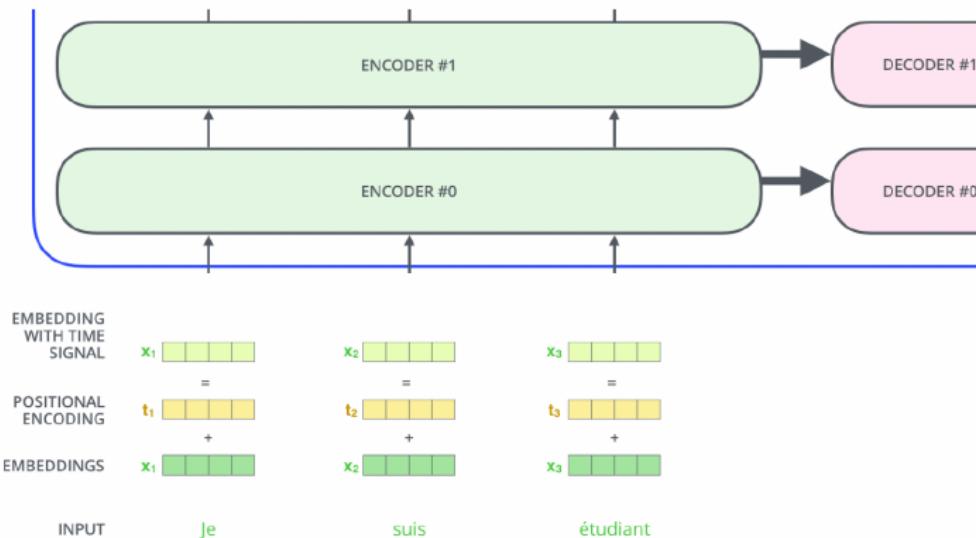


Figure 1: The Transformer - model architecture.

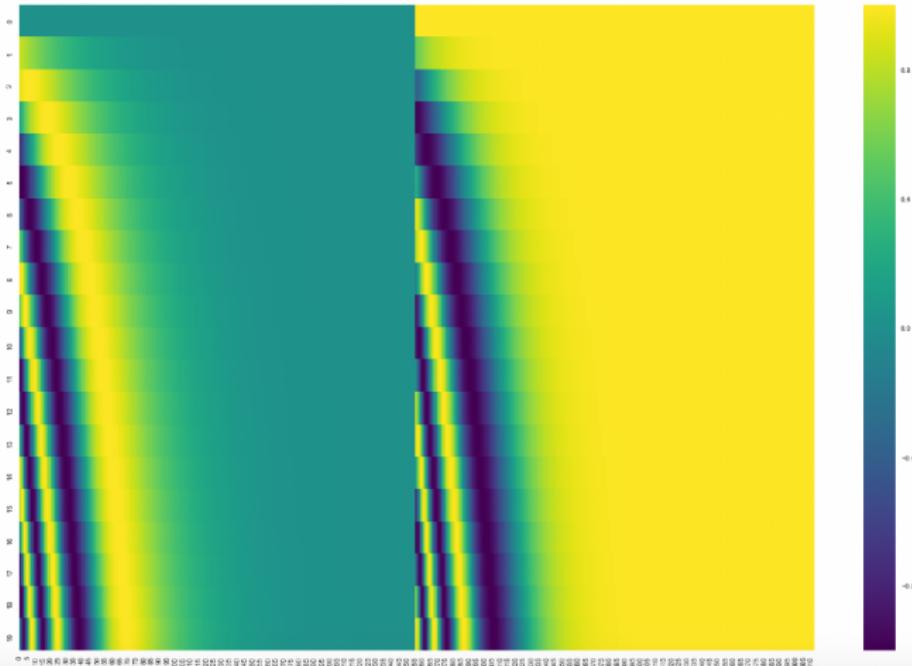
Positional Encodings

- 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)



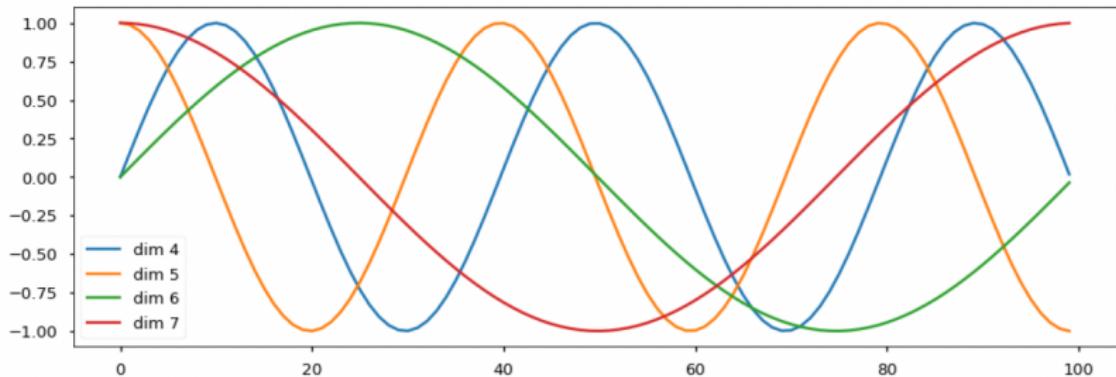
Sinusoidal Positional Encodings

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

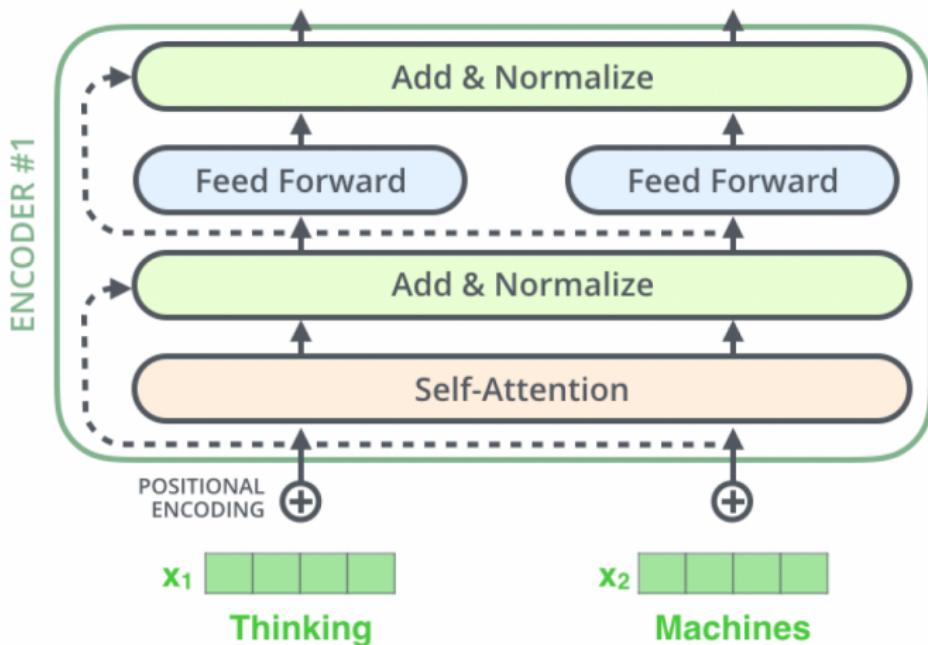


Sinusoidal Positional Encodings

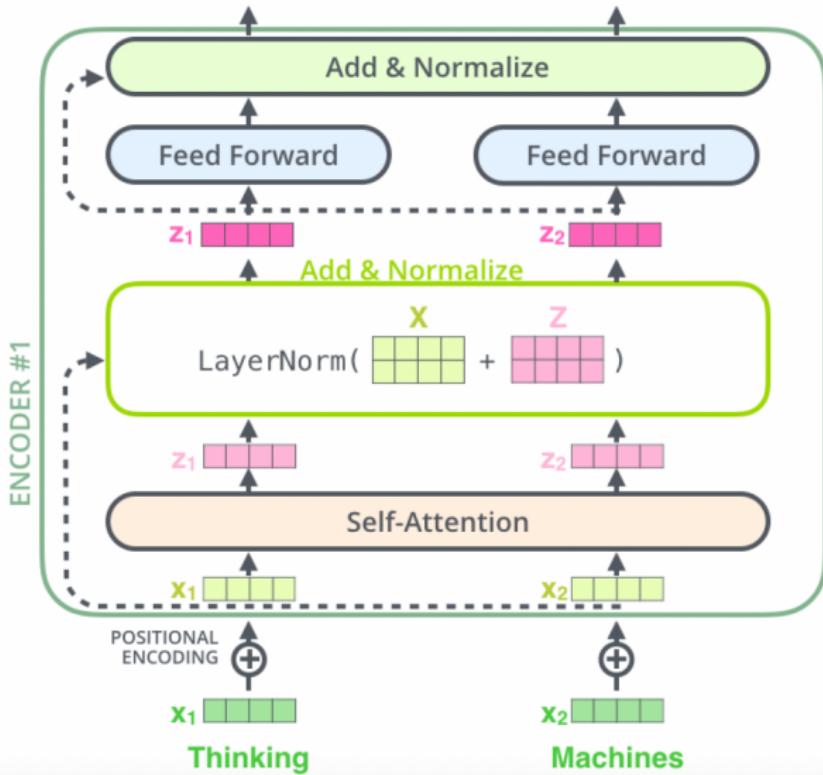
$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{2i/d}}\right) \quad PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000^{2i/d}}\right)$$



Residuals and Layer Normalization



Residuals and Layer Normalization



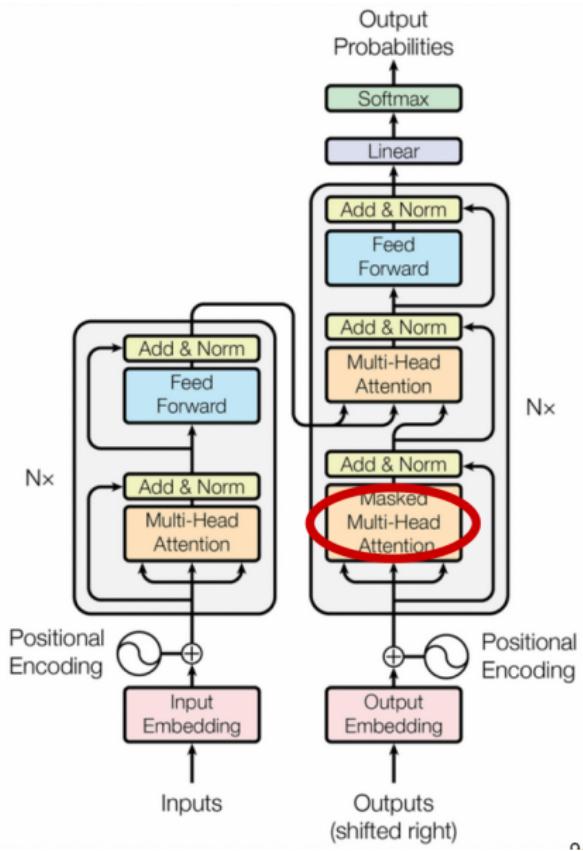
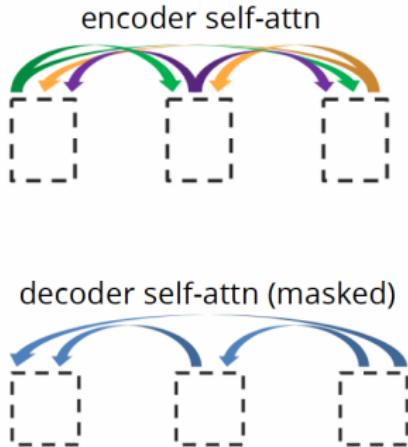
The Decoder

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

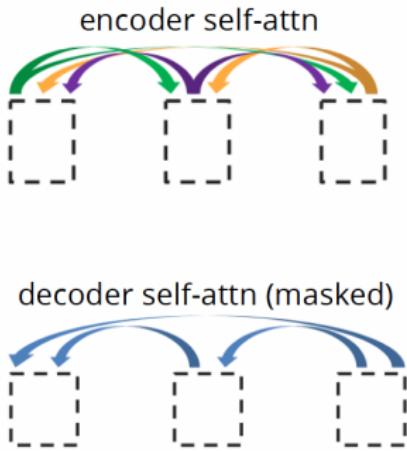
Everything is pretty much the same as in the encoder, with two twists:

- 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**).

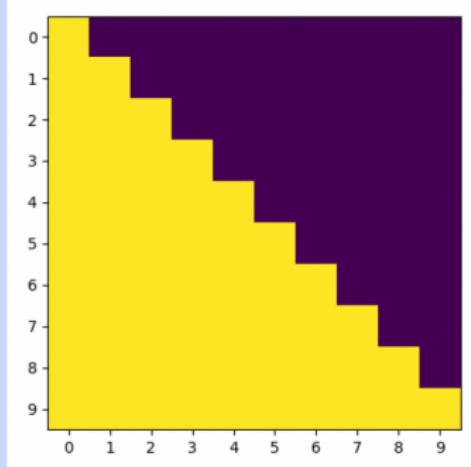
The Decoder



The Decoder



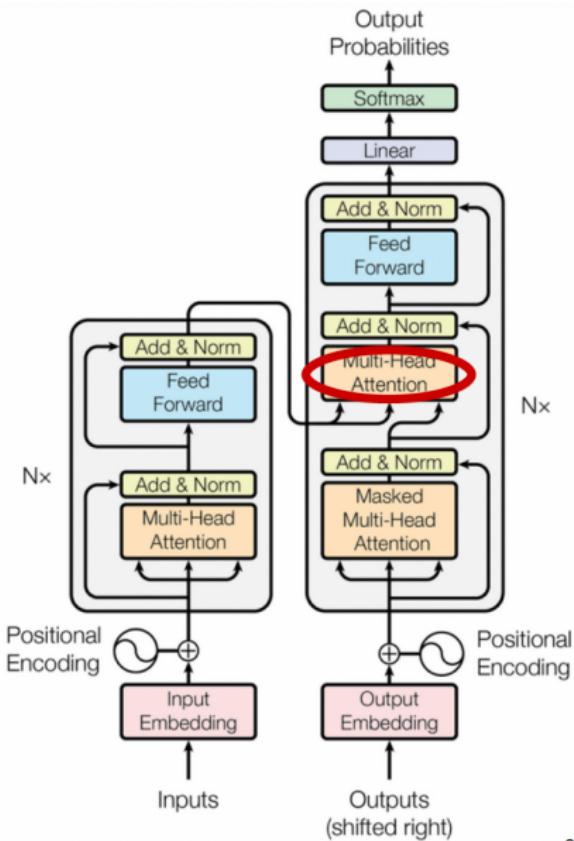
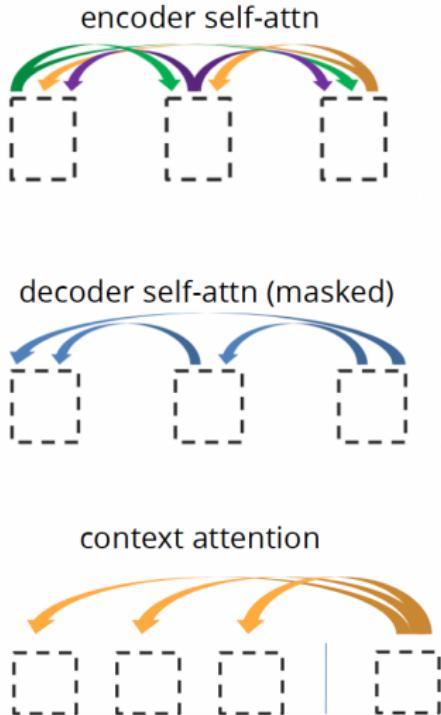
- Mask subsequent positions (before softmax)



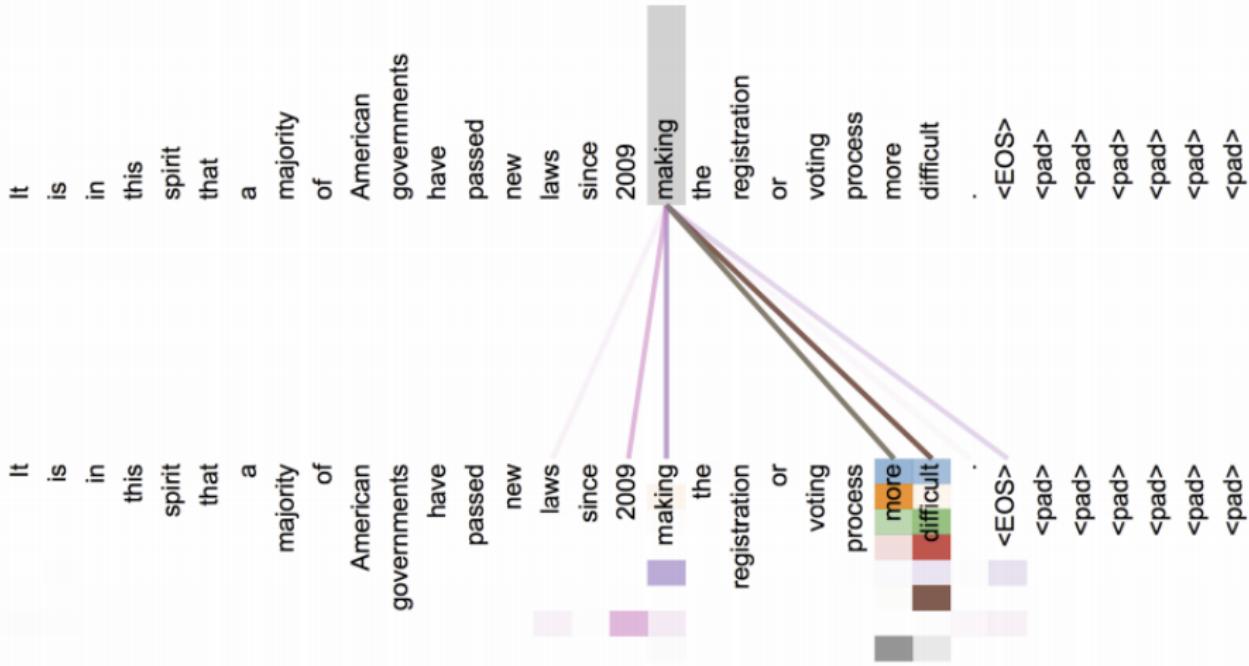
- In PyTorch

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

The Decoder



Attention Visualization Layer 5



Implicit Anaphora Resolution

