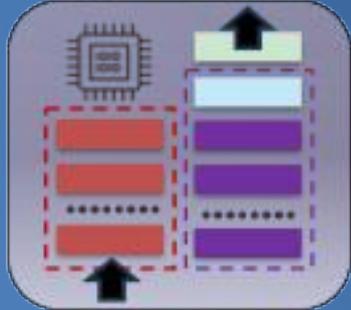


Image: Transformers: "...and one 'architecture' to rule them all" – Juxtaposition by Raaid Saqur (2024).

Transformers

CSC401/2511 – Natural Language Computing – Spring 2026

Ken Shi and Gerald Penn



Transformer networks

- Breakout paper: Vaswani et al. (2017) Attention is all you need.
- **Core idea:** replace recurrent connections with attention

Table 2: The Transformer achieves better BLEU scores than previous state-of-the-art models on the English-to-German and English-to-French newstest2014 tests at a fraction of the training cost.

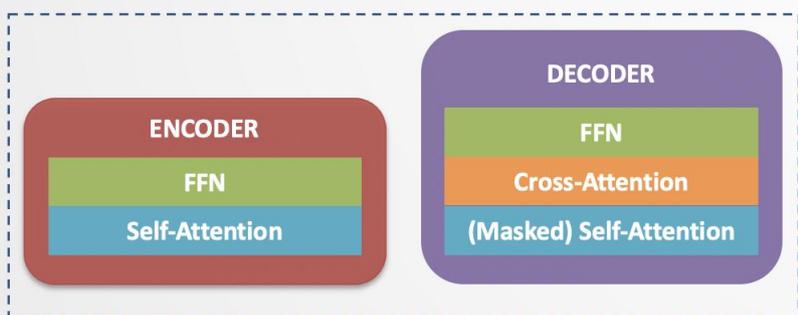
Model	BLEU		Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet [15]	23.75			
Deep-Att + PosUnk [32]		39.2		$1.0 \cdot 10^{20}$
GNMT + RL [31]	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4 \cdot 10^{20}$
ConvS2S [8]	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$
MoE [26]	26.03	40.56	$2.0 \cdot 10^{19}$	$1.2 \cdot 10^{20}$
Deep-Att + PosUnk Ensemble [32]		40.4		$8.0 \cdot 10^{20}$
GNMT + RL Ensemble [31]	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$
ConvS2S Ensemble [8]	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	41.0	$2.3 \cdot 10^{19}$	

¹ Vaswani, Ashish, et al. "Attention is all you need." *NeurIPS* (2017).

Transformer networks (abstract)

Core Idea

Now **two** kinds of attention



- Encoder uses **self-attention**

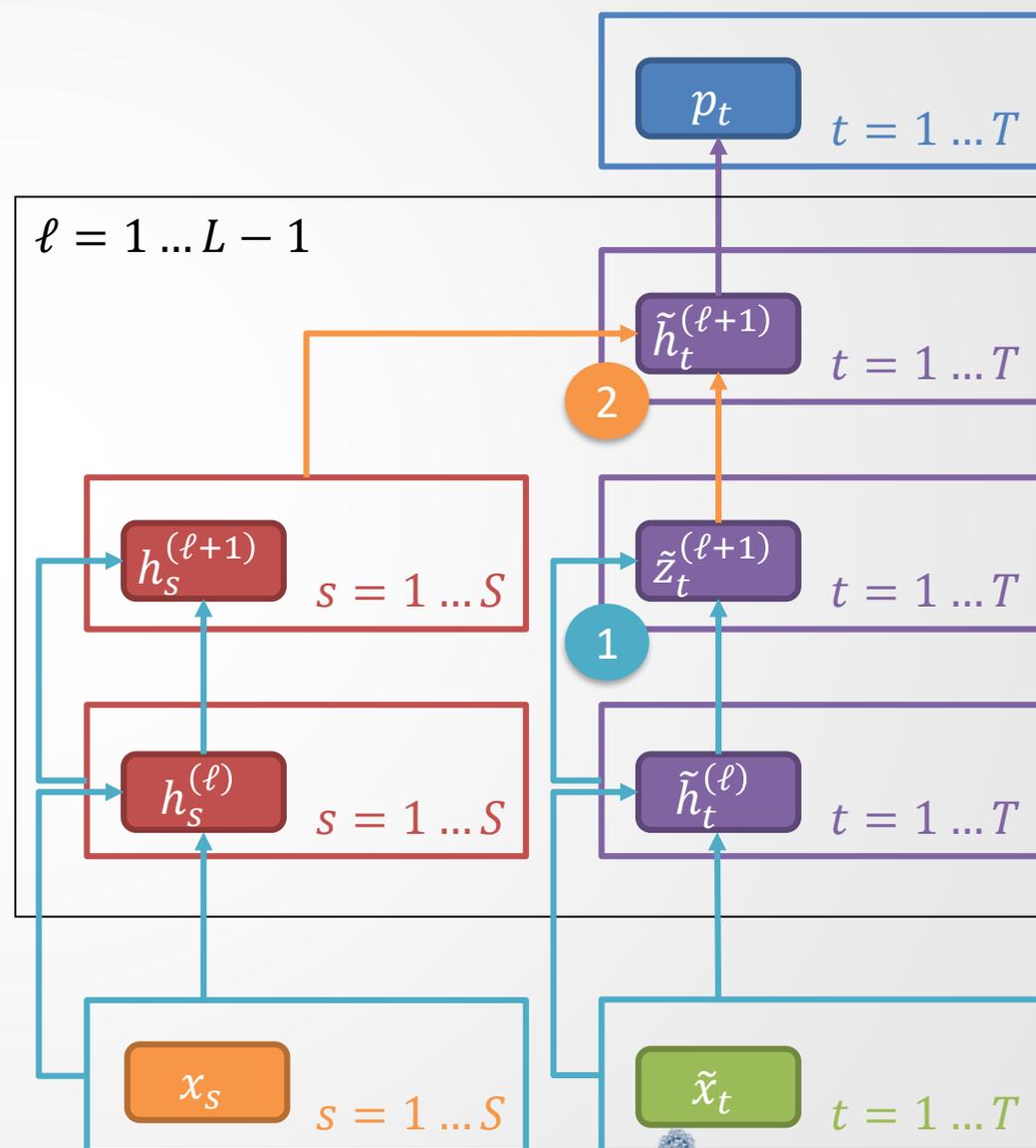
$$h_s^{(\ell+1)} \leftarrow Att_{Enc} \left(h_s^{(\ell)}, h_{1:S}^{(\ell)} \right)$$

Decoder uses **1. self-attention***

$$\tilde{z}_t^{(\ell+1)} \leftarrow Att_{Dec1} \left(\tilde{h}_t^{(\ell)}, \tilde{h}_{1:t}^{(\ell)} \right)$$

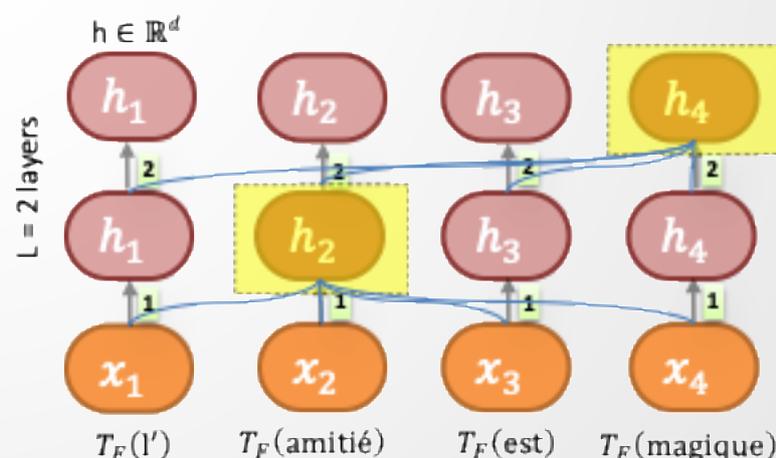
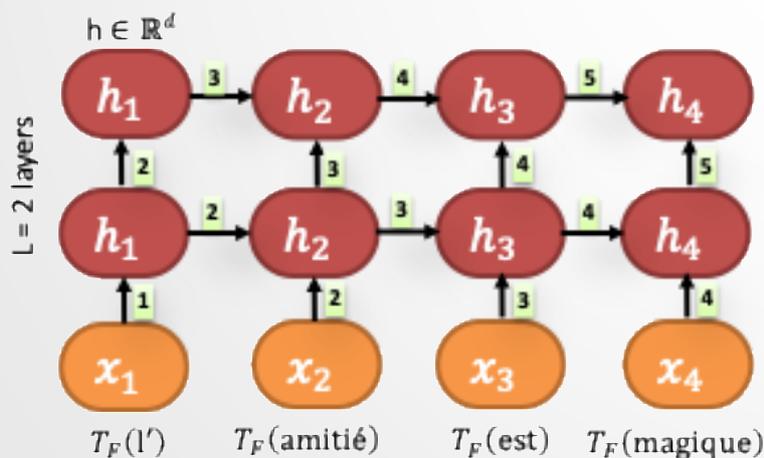
then **2. attention with encoder**

$$\tilde{h}_t^{(\ell+1)} \leftarrow Att_{Dec2} \left(\tilde{z}_t^{(\ell+1)}, h_{1:S}^{(\ell+1)} \right)$$

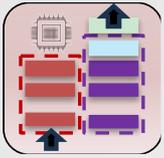


Transformer motivations

- **Limitations** of recurrent connections: long-term dependencies, lack of parallelizability, interaction distance (steps to distant tokens).
- **Attention** allows access to entire sequence
- Lots of computation can be shared, parallelized across sequence indices. Identical layers: [self, cross]-attention, feed-forward w/ tricks
- There are more *avant-garde* applications to other domains



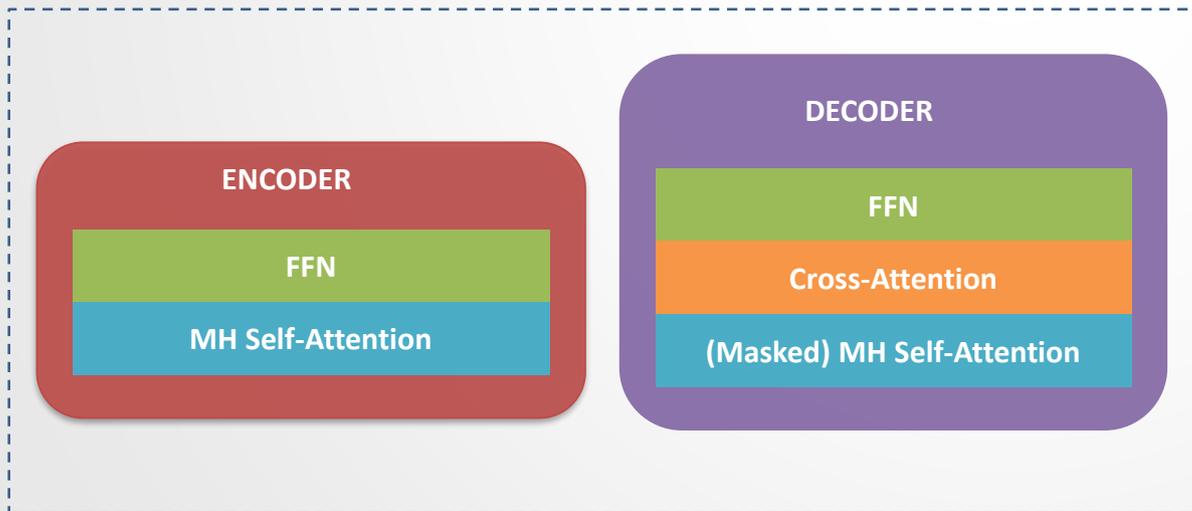
Source sentence (French): *L' amitié est magique*
Target sentence (English): *Friendship is magic*



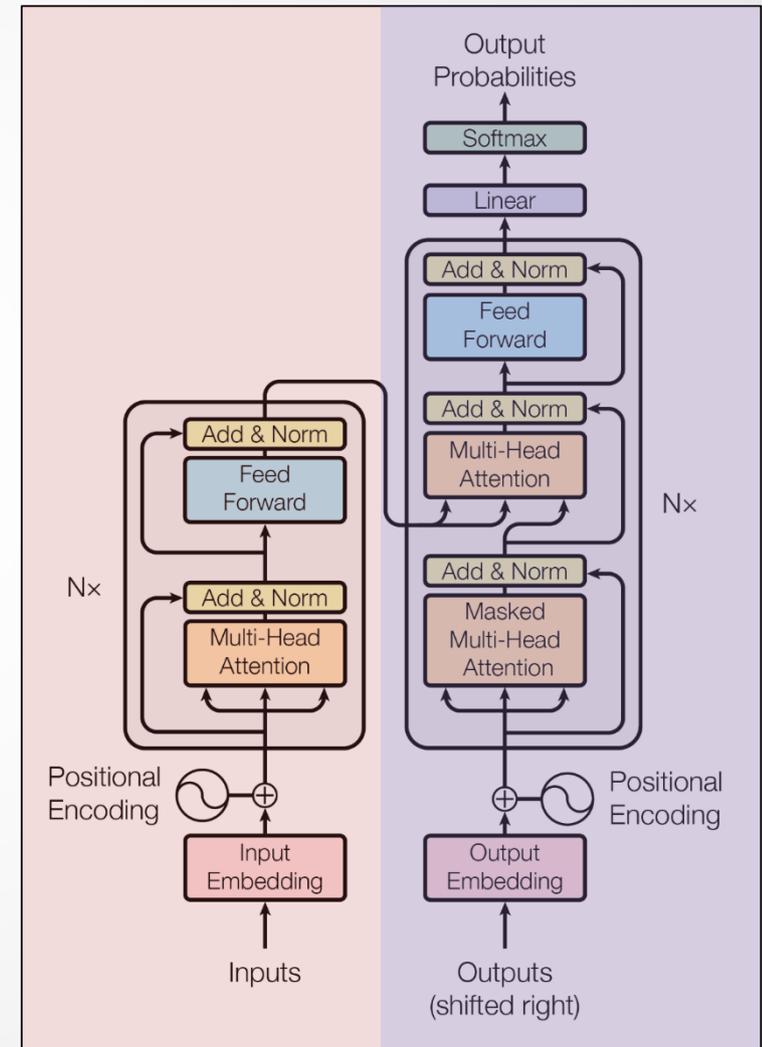
Transformer Architecture

- Architecture diagram from Vaswani^[1]
- Building blocks:
 1. Encoder
 2. Decoder

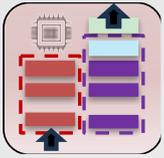
Encoder Decoder



Encoder Decoder



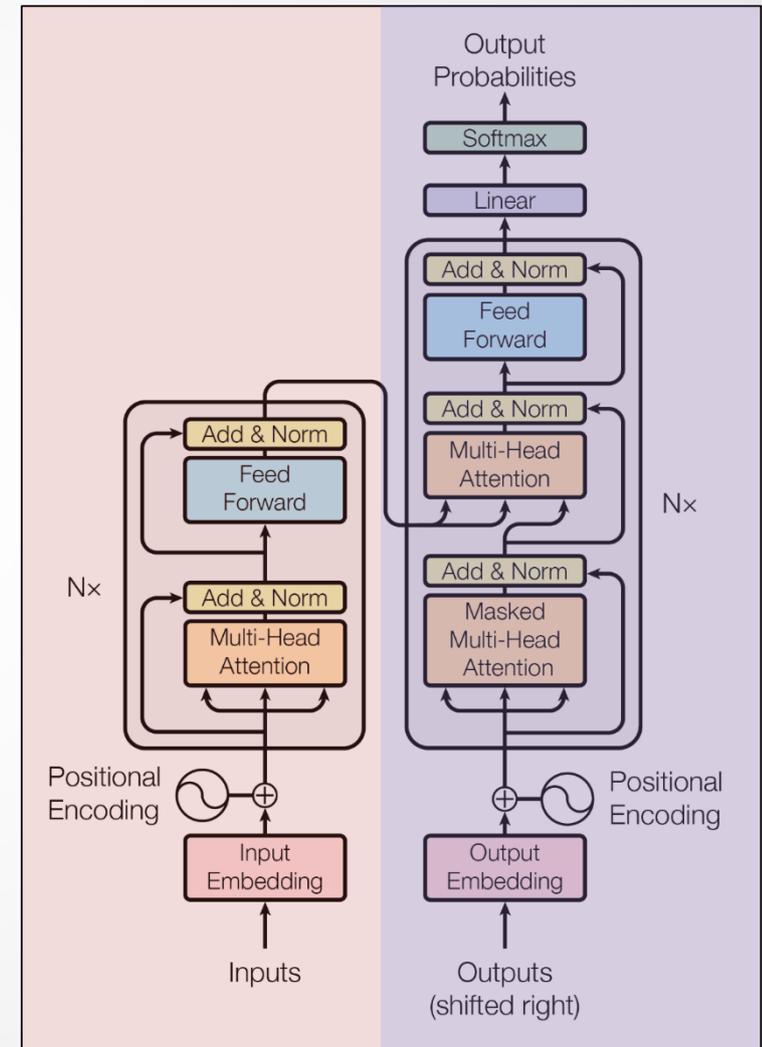
¹Vaswani, Ashish, et al. "Attention is all you need." *NeurIPS* (2017).



Transformer Architecture

- Architecture diagram from Vaswani^[1]
- Building blocks:
 1. Encoder
 2. Decoder
- Main components within building blocks:
 - Attention mechanisms:
 - single and multi-head attention
 - self, cross, and masked attention
 - Feed-forward MLPs (FFN)
 - Layer normalization (LN)
 - Positional encodings (PE)
 - Residual connections

Encoder Decoder

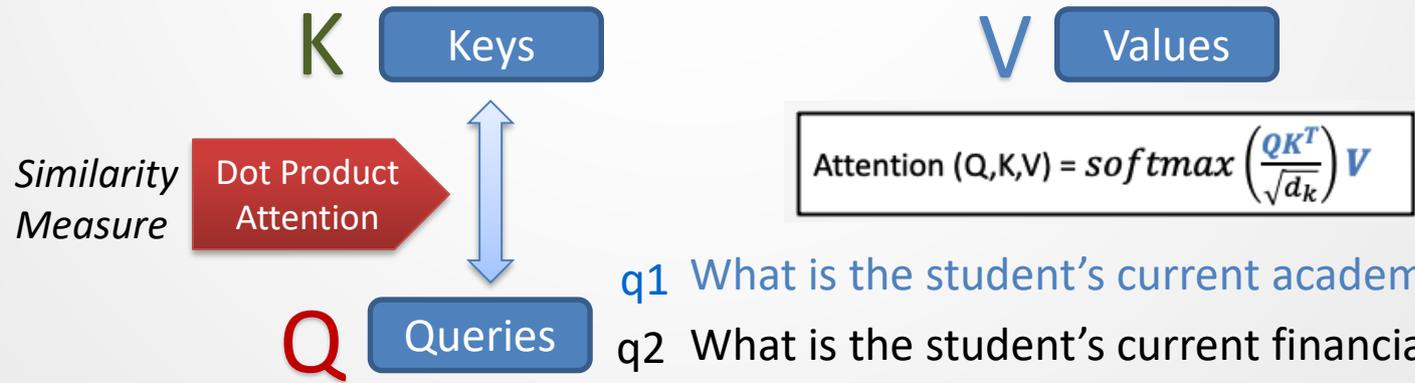
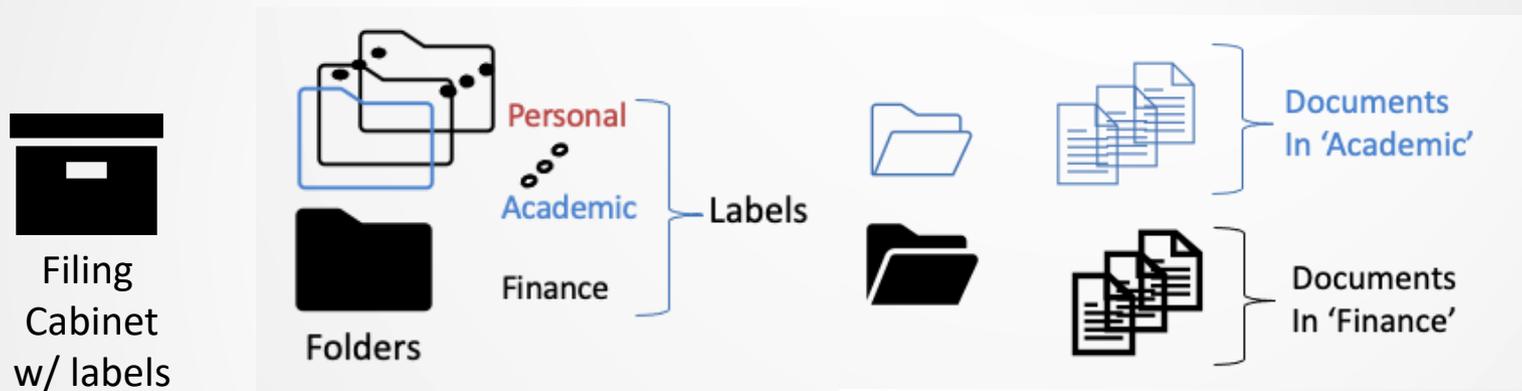
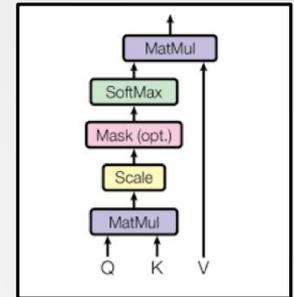


¹ Vaswani, Ashish, et al. "Attention is all you need." *NeurIPS* (2017).

Transformer Attention(Q,K,V) : Intuition

In the classical roboquity era (2100-2250 AD), humans in designated zones/zoos are only allowed *filing cabinets* and *paper documents* to store information.

ACORN has been terminated, and UofT students' info (financial, *academic*, *personal*) retrieval works as follows:

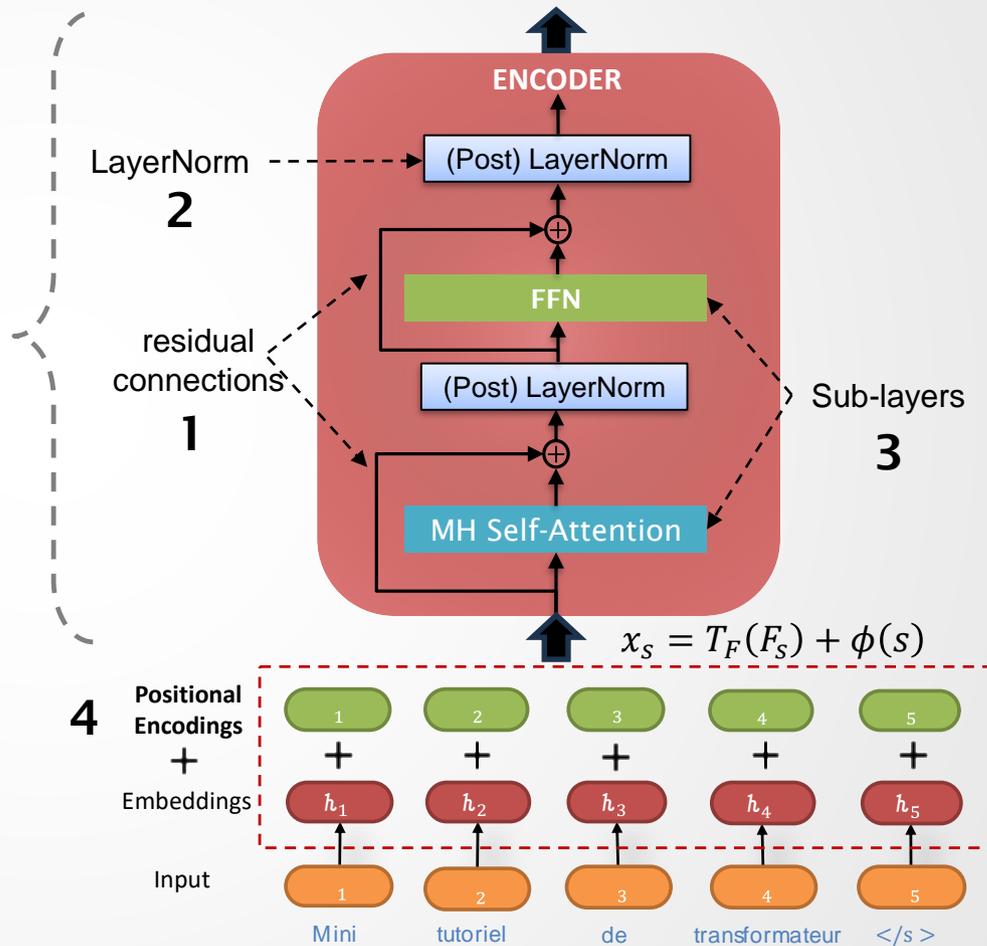
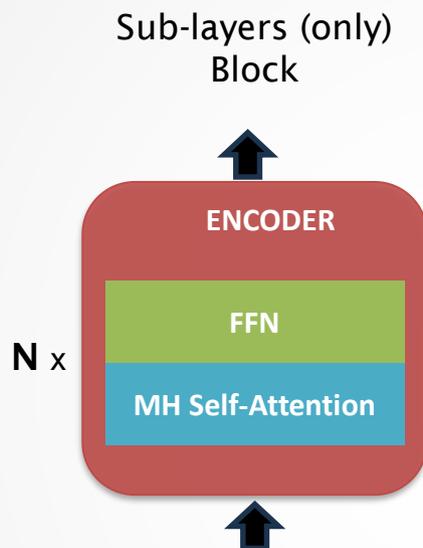
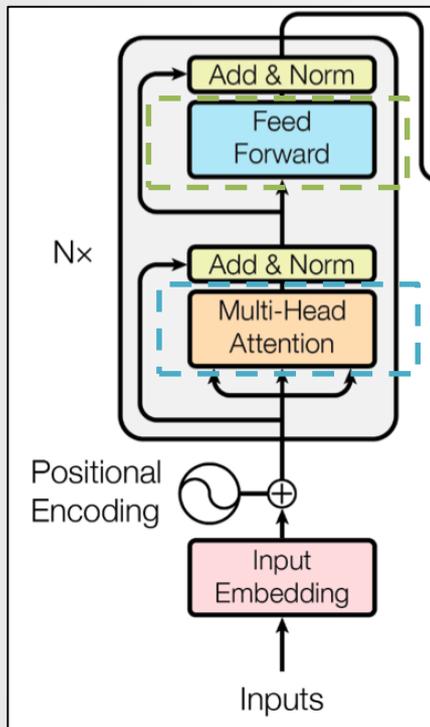


- q1 What is the student's current academic standing?
- q2 What is the student's current financial status?
- q3 What is the student's residency status in Canada?

Transformer Encoder

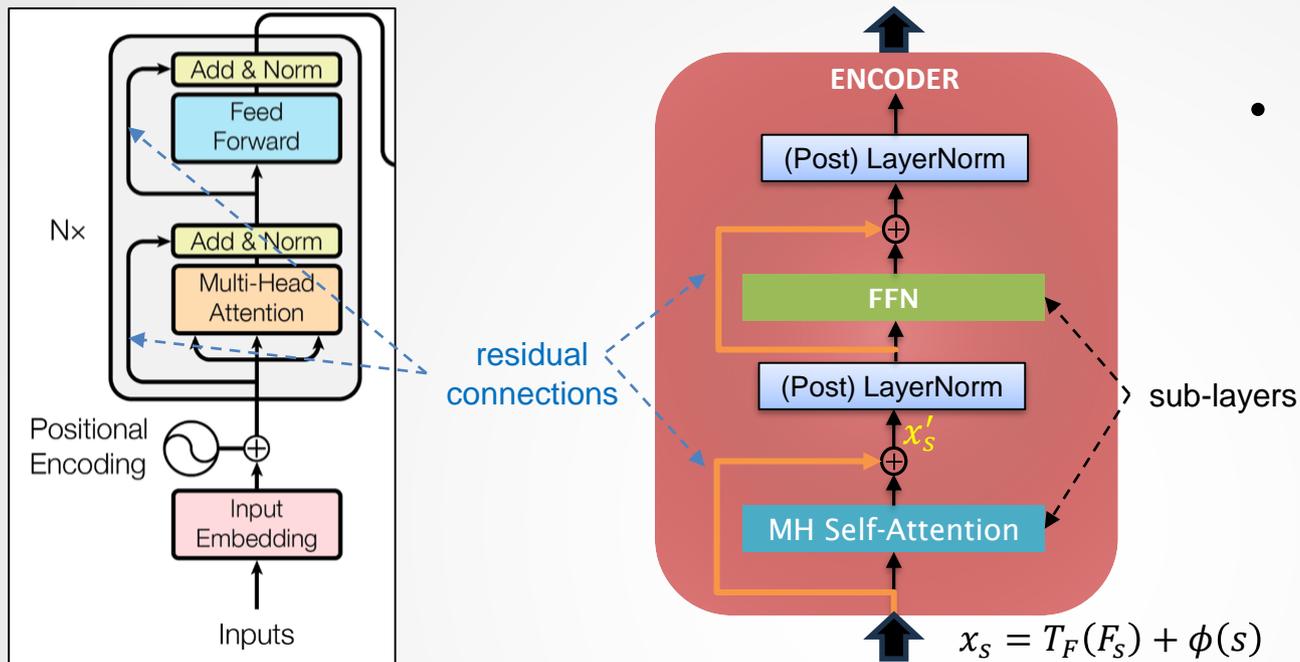
F_S : Mini tutoriel de transformateur

Encoder  Decoder E_T : Tiny transformer tutorial



- **Plan:** discuss building blocks:
 1. Residual connections
 2. LayerNorm
 3. Attention and FFN sub-layers
 4. Positional encodings

Residual Connections



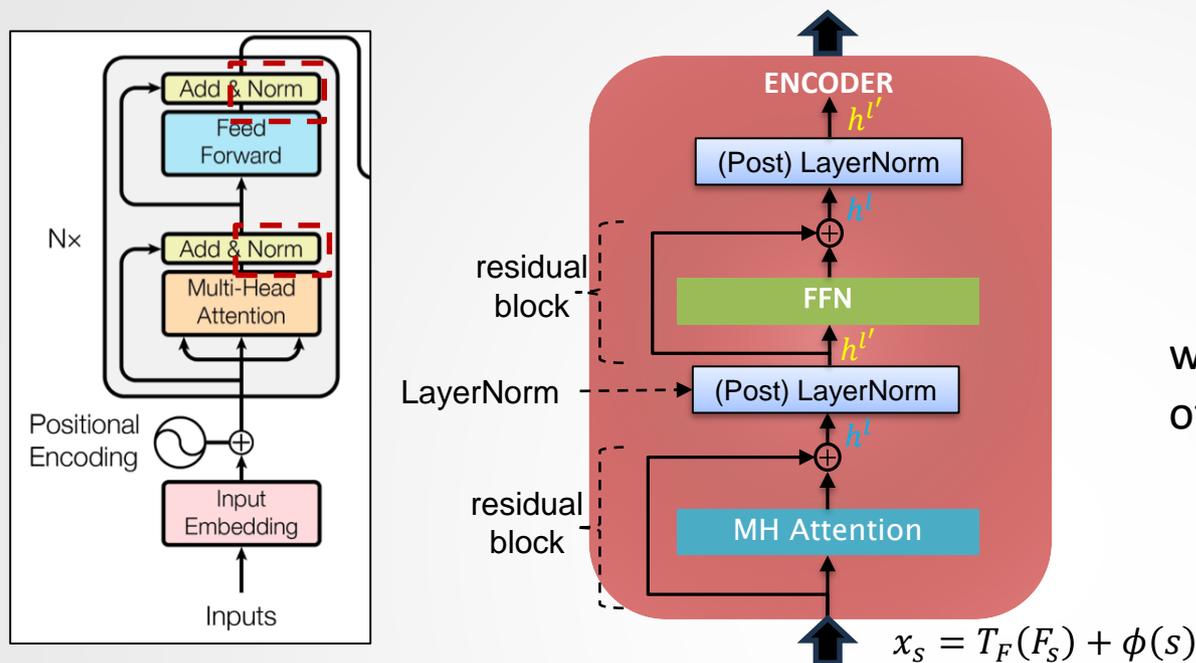
- Old idea used in speech:

$$x_{l+1} = \mathcal{F}(x_l) + x_l$$
 + Helps smoothen loss curvature allowing better backprop.

- **Problem:** NNs can't recover some important information from their input
- **Solution:** Add back the input embeddings to the sub-layer's output moving up

$$x'_s = \text{Sublayer}(x_s) + x_s$$
- **Analogy:** think of broadcasting through a noisy-channel. Big help to pass information through a side-channel without distortion.

Layer Normalization: default (Post-) LN



$$h^{l'} = \text{LayerNorm}(h^l)$$

$$= \gamma \left(\frac{h^l - \mu^l}{\sigma^l} \right) + \beta$$

where μ, σ are mean and std. dev. of features in h^l . γ, β are scale, bias params.

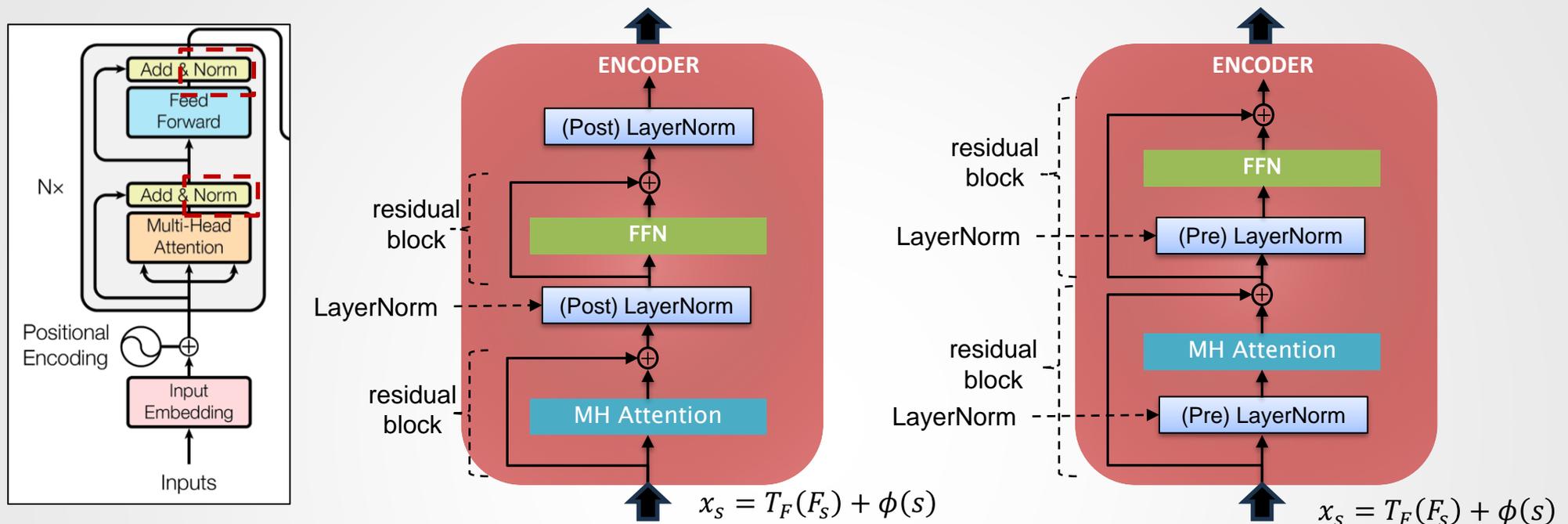
$$\mu = \frac{1}{d} \sum_{k=1}^d h_k^l$$

$$\sigma^2 = \frac{1}{d} \sum_{k=1}^d (h_k^l - \mu)^2$$

- Layer Normalization^[1]:
 - **Normalize** input layer's distribution to 0 mean and 1 standard deviation.
 - Removes uninformative variation in layer's features

¹Ba, Jimmy Lei, Jamie Ryan Kiros, and Geoffrey E. Hinton. "Layer normalization. 2016" [\[link\]](#)

Layer Normalization Variant: Pre-LN



- Layer Normalization: two popular variants

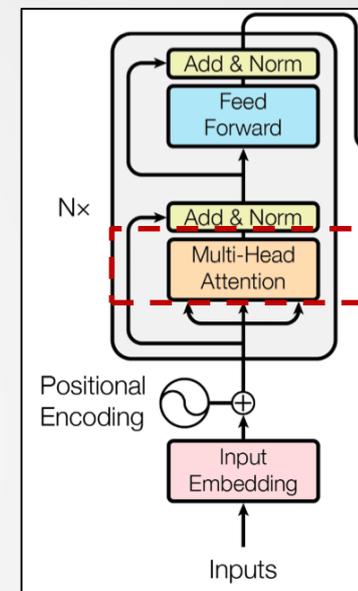
- *Post layer normalization* (Post-LN): original Transformer model: requires learning rate warm-up due to initial instability of large output gradients.
- *Pre layer normalization* (**Pre-LN**): puts layer-norm within the residual block. Allows removing warm-up stage. More stable training initialization.^[2]

- Secret sauce: Fisher information matrix (based on KL-divergence!)

² Xiong, Ruibin, et al. "On layer normalization in the transformer architecture." *ICML*. PMLR, 2020

Transformer Encoder - Self Attention

- Recall our discussion of attention in transformer LMs
- Steps:
 1. Calculate the query, **key**, and **value** for each token
 - Attention of each query (q_i) against all the keys ($k_{1:j}$)
 2. Calculate the **attention score** between query and keys
 3. **Normalize** the attention scores by applying softmax
 4. Calculate values by taking a **weighted sum**



$$\begin{aligned}q_i &= W^Q x_i \\k_i &= W^K x_i \\v_i &= W^V x_i\end{aligned}$$

$$\begin{aligned}a_{i,j} &= \text{score}(q_i, k_j) \\a_{i,j} &= q_i \cdot k_j \\a_{i,j} &= \frac{q_i \cdot k_j}{\sqrt{d_k}}\end{aligned}$$

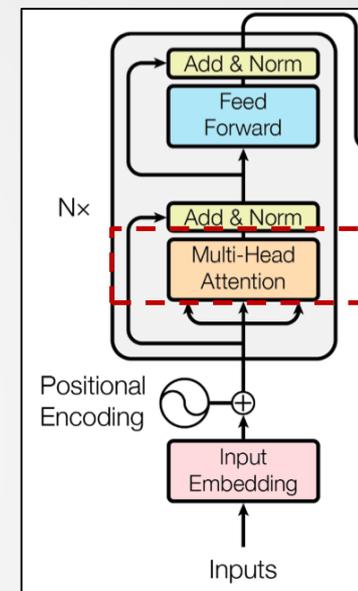
$$\begin{aligned}a_{i,j} &= \text{softmax}(a_{i,1:K}) \\ \alpha_{i,j} &= \frac{\exp(a_{i,j})}{\sum_{k=1}^K \exp(a_{i,k})}\end{aligned}$$

$$c_i = \sum_j \alpha_{i,j} v_j$$

¹ Vaswani, Ashish, et al. "Attention is all you need." *NeurIPS* (2017).

Transformer Encoder - Self Attention

- Recall our discussion of attention in transformer LMs
- Steps:
 1. Calculate the **query**, **key**, and **value** for each token
 - Attention of each query (q_i) against all the keys ($k_{1:j}$)
 2. Calculate the **attention score** between query and keys
 3. **Normalize** the attention scores by applying softmax
 4. Calculate values by taking a **weighted sum**



Vectorized notation:

$$\begin{aligned} Q &= XW^Q \\ K &= XW^K \\ V &= XW^V \end{aligned}$$

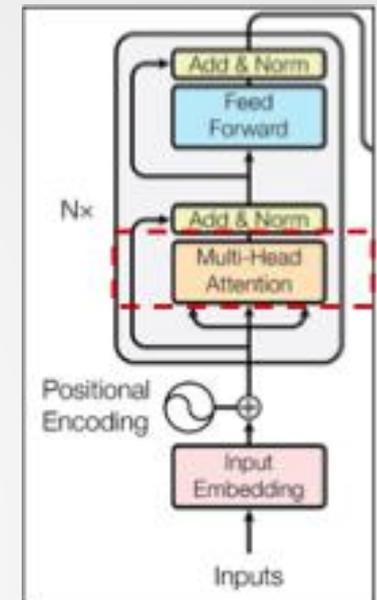
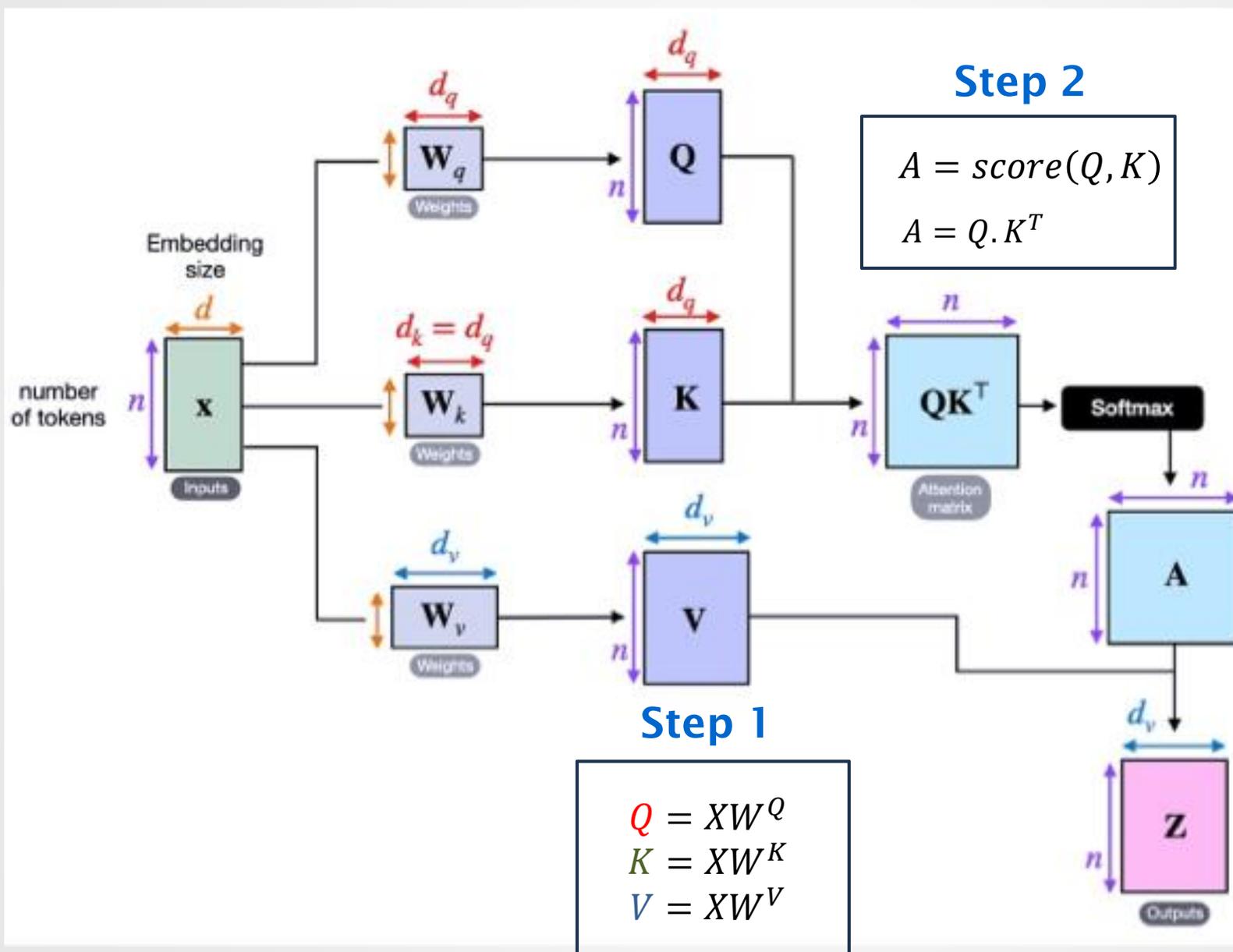
$$\begin{aligned} A &= \text{score}(Q, K) \\ A &= Q \cdot K^T \\ A &= \frac{Q \cdot K^T}{\sqrt{d_k}} \end{aligned}$$

$$\begin{aligned} A &= \text{softmax}(A) \\ A &= \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) \end{aligned}$$

$$Z = A \cdot V$$

¹Vaswani, Ashish, et al. "Attention is all you need." *NeurIPS* (2017).

Transformer Encoder - Self Attention



Step 3

$$A = \text{softmax}(A)$$

Step 4

$$Z = A \cdot V$$

¹Vaswani, Ashish, et al. "Attention is all you need." *NeurIPS* (2017).

Multi-head Self Attention (MHA)

- Multi-head attention (MHA) jointly attends to information from different representation subspaces – hopefully the extra heads will get something right that the earlier heads missed.

$$MHA(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

where $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$

And projections are parameter matrices:

$$\begin{aligned} W_i^Q &\in \mathbb{R}^{d_{\text{model}} \times d_k} & W^O &\in \mathbb{R}^{hd_v \times d_{\text{model}}} \\ W_i^K &\in \mathbb{R}^{d_{\text{model}} \times d_k} & d_k = d_v &= \frac{d_{\text{model}}}{h} \\ W_i^V &\in \mathbb{R}^{d_{\text{model}} \times d_v} \end{aligned}$$

¹Vaswani, Ashish, et al. "Attention is all you need." *NeurIPS* (2017).

Feed-forward (FFN) layers

- Attention only **re-weights** the **value vectors**
- We still need more degrees of freedom (and *non-linearity*) to learn
- The feed-forward layer(s) (**FFN**) provide these non-linearities to attention layer outputs
- Specifically, each output x undergoes two (layer) linear transformations with a ReLU activation in between. Pointwise:

$$FFN(x_i) = \max(0, x_i W_1 + b_1) W_2 + b_2$$

$$z_i = \sum_j \alpha_{i,j} v_j$$

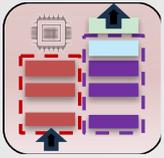
- FFN sub-layer is applied to each token pos. separately and identically
- Given x , a sequence of tokens (x_1, \dots, x_S) :

$$FFN(x) = \text{ReLU}(xW_1 + b_1) W_2 + b_2$$

$$Z = A \cdot V$$

where W_1, W_2, b_1 and b_2 are parameters

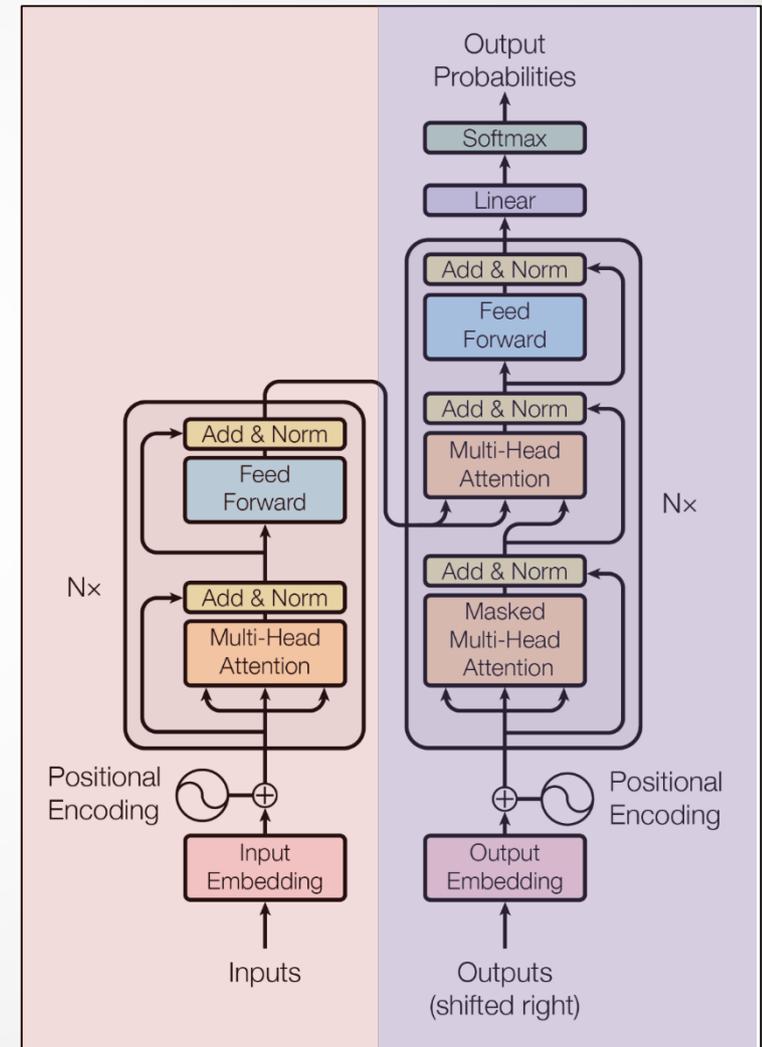
¹ Vaswani, Ashish, et al. "Attention is all you need." *NeurIPS* (2017).



Recap: Transformer Architecture

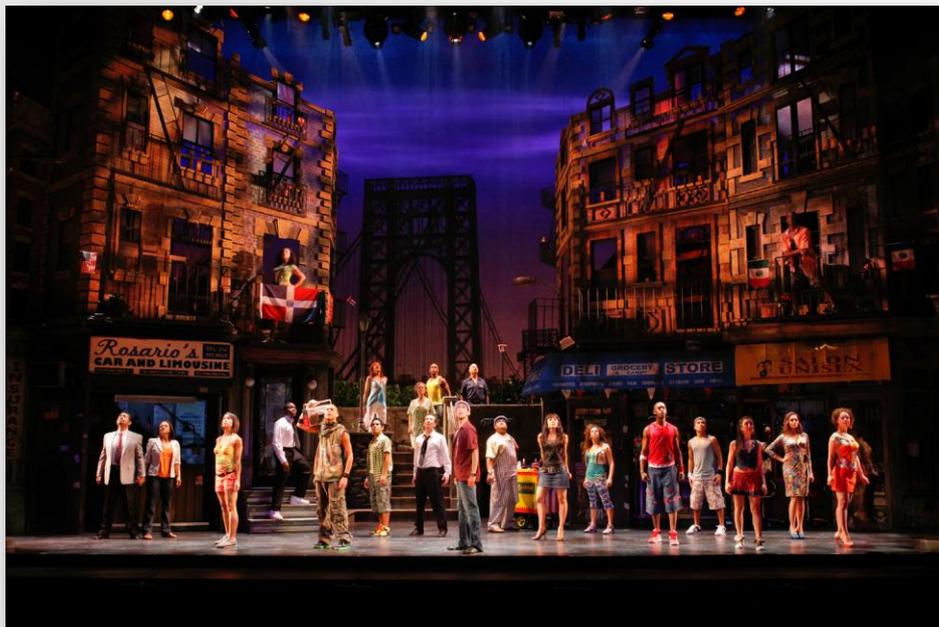
- Architecture diagram from Vaswani^[1]
- Building blocks:
 - ✓ 1. Encoder
 - 2. Decoder
- Main components within building blocks:
 - Attention mechanisms:
 - ✓ single and multi-head attention
 - self, **cross, and masked attention**
 - ✓
 - Feed-forward MLPs (FFN)
 - Layer normalization (LN)
 - Positional encodings (PE)
 - Residual connections

Encoder Decoder



Aside: Contextualized embeddings

- What does the word *play* mean?

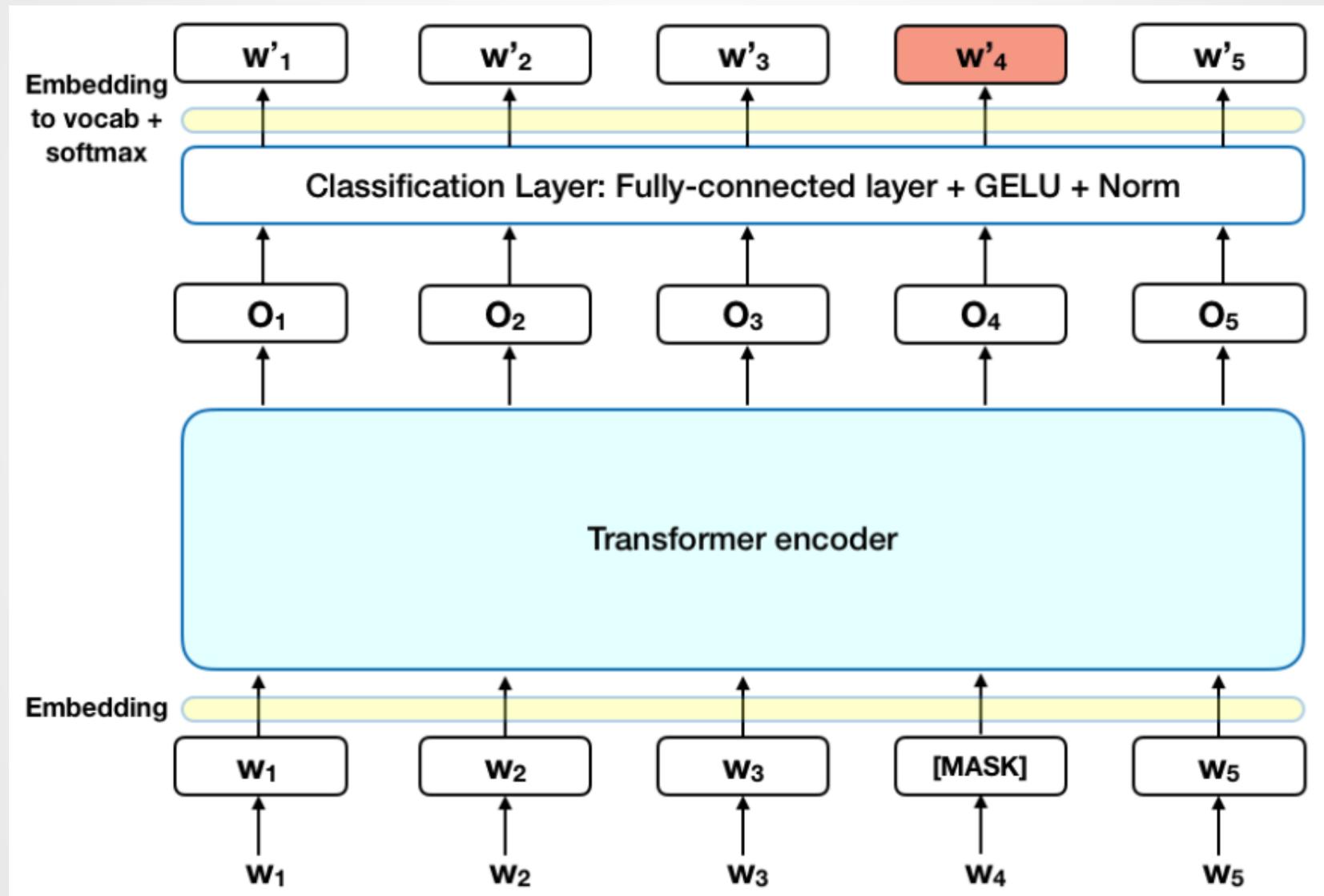


AllenNLP

Peters ME, Neumann M, Iyyer M, *et al.* (2018) Deep contextualized word representations.

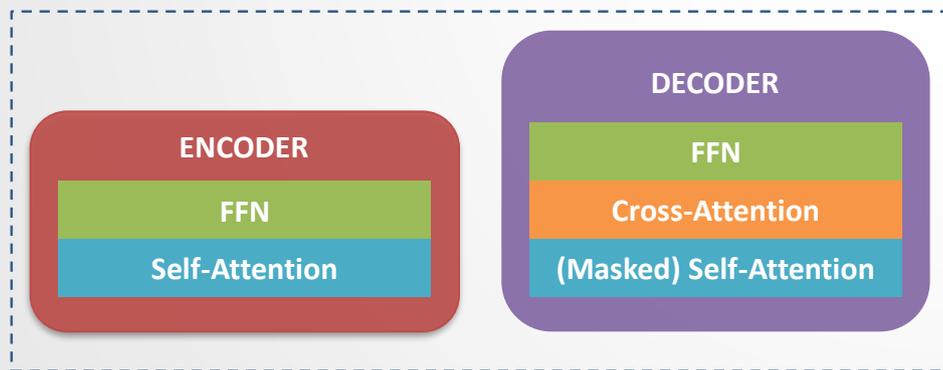
Published Online First: 2018. doi:10.18653/v1/N18-1202; <http://arxiv.org/abs/1802.05365>

BERT: Masked Language Modelling

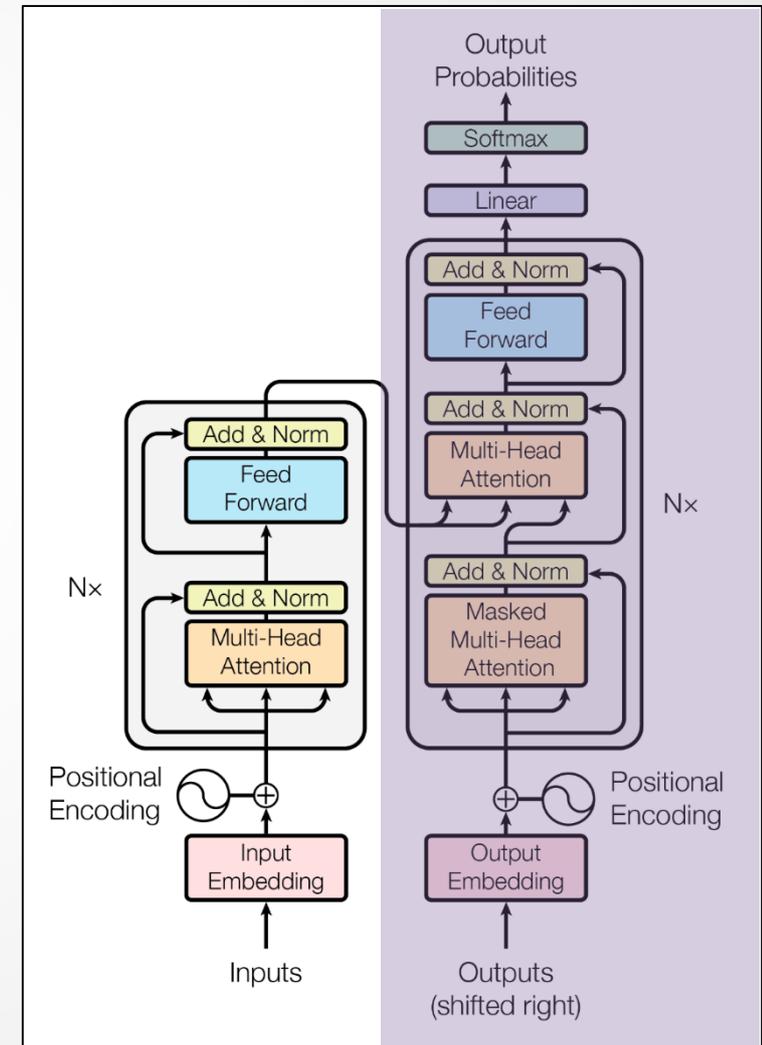


Next block: Transformer Decoder

- Layer normalization, residual connections, FFNs are identical to the encoder block
- Thus, we focus on remaining:
 - Masked/Causal self-attention sub-layer
 - Cross-attention



Encoder  Decoder



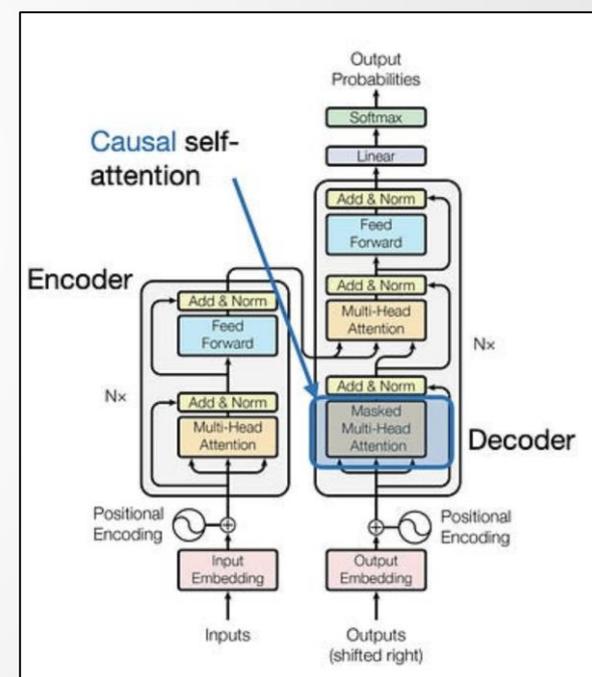
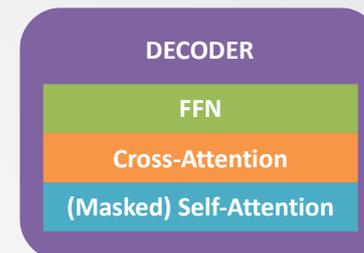
¹ Vaswani, Ashish, et al. "Attention is all you need." *NeurIPS* (2017).

Decoder – Masked Self-Attention

- Masked (Multi-head) self-attention:
 - Enforce auto-regressive language modeling objective. The decoder cannot peek and pay attention to the (unknown) future words
 - Solution:** use a *look-ahead mask* M , by setting attention scores of future tokens to $-\infty$.

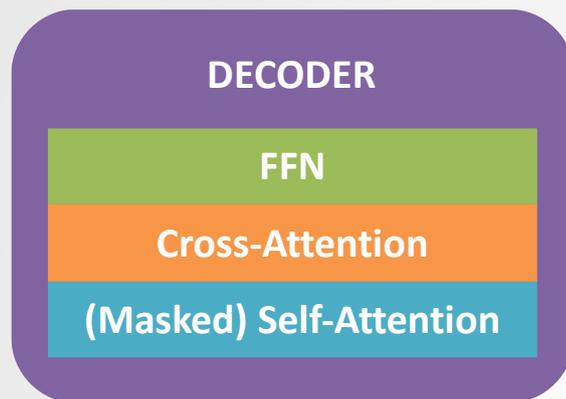
$$\begin{array}{c}
 a^K \quad b^K \quad c^K \quad D^K \\
 \begin{array}{c}
 a^Q \\
 b^Q \\
 c^Q \\
 D^Q
 \end{array}
 \begin{bmatrix}
 0 & -\infty & -\infty & -\infty \\
 0 & 0 & -\infty & -\infty \\
 0 & 0 & 0 & -\infty \\
 0 & 0 & 0 & 0
 \end{bmatrix}
 \end{array}$$

$$a_{ij} = \begin{cases} q_i^T \cdot k_j, & j < i \\ -\infty, & j \geq i \end{cases}$$



¹ Vaswani, Ashish, et al. "Attention is all you need." *NeurIPS* (2017).

Decoder – Masked Self-Attention



E_T : $\langle s \rangle$ tiny transformer tutorial

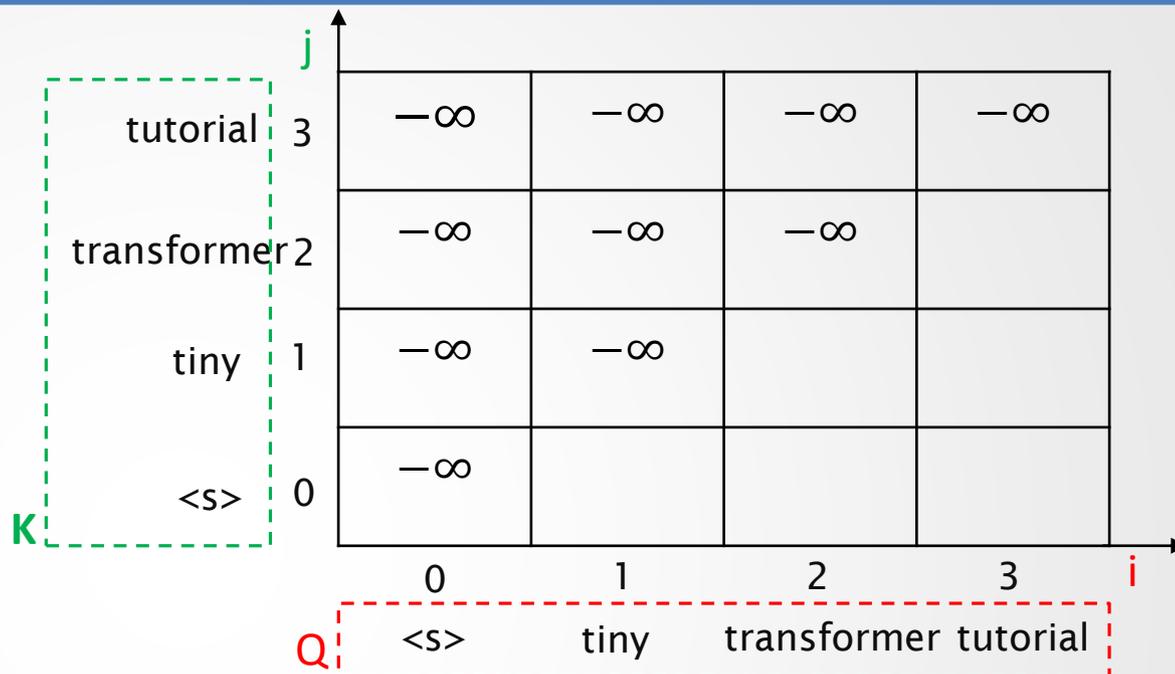


Recall

$$a_{i,j} = \text{score}(q_i, k_j)$$

or,

$$A = \text{score}(Q, K)$$

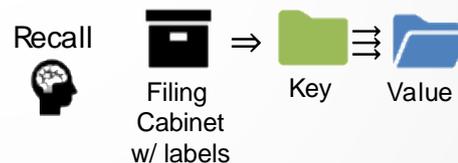


¹ Vaswani, Ashish, et al. "Attention is all you need." *NeurIPS* (2017).

Encoder-Decoder (Cross) Attention

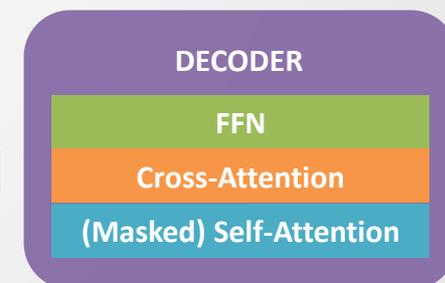
- In self-attention: **Q**, **K** and **V** have same source (tokens)
- **Cross attention** is encoder \leftrightarrow decoder attention between output vectors
- Using our running example and notation:
 - Let h_1, \dots, h_S be **encoder** output vectors, where $h_i \in \mathbb{R}^{d_k}$
 - Let $\tilde{h}_1, \dots, \tilde{h}_T$ be **decoder** output vectors, where $\tilde{h}_i \in \mathbb{R}^{d_q}$

- Then, keys and values: **K** and **V** comes from **encoder** (or, *memory*):



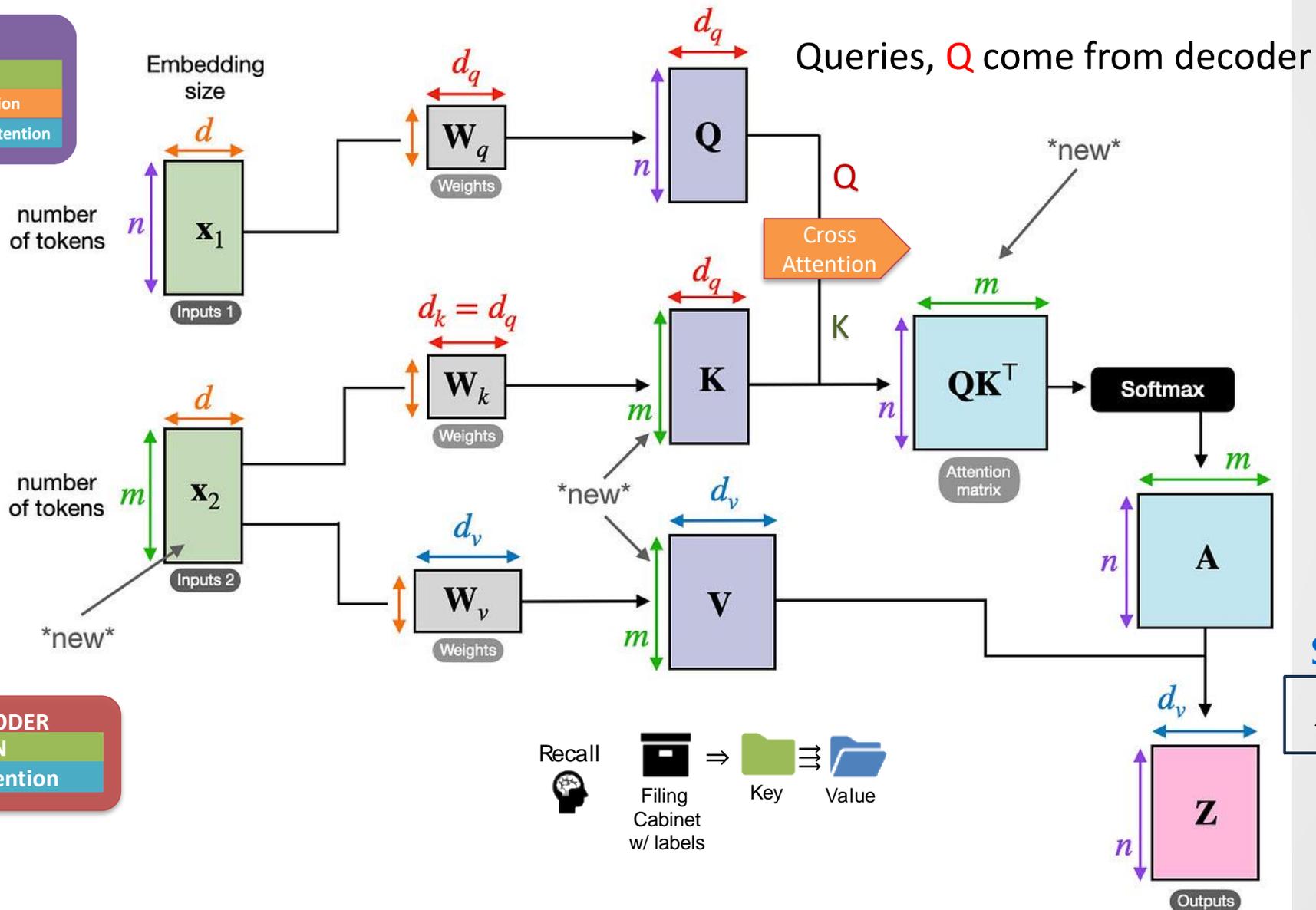
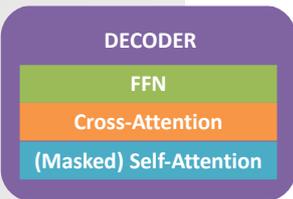
- $k_i = Kh_i$, $v_i = Vh_i$

- Queries, **Q** comes from decoder:

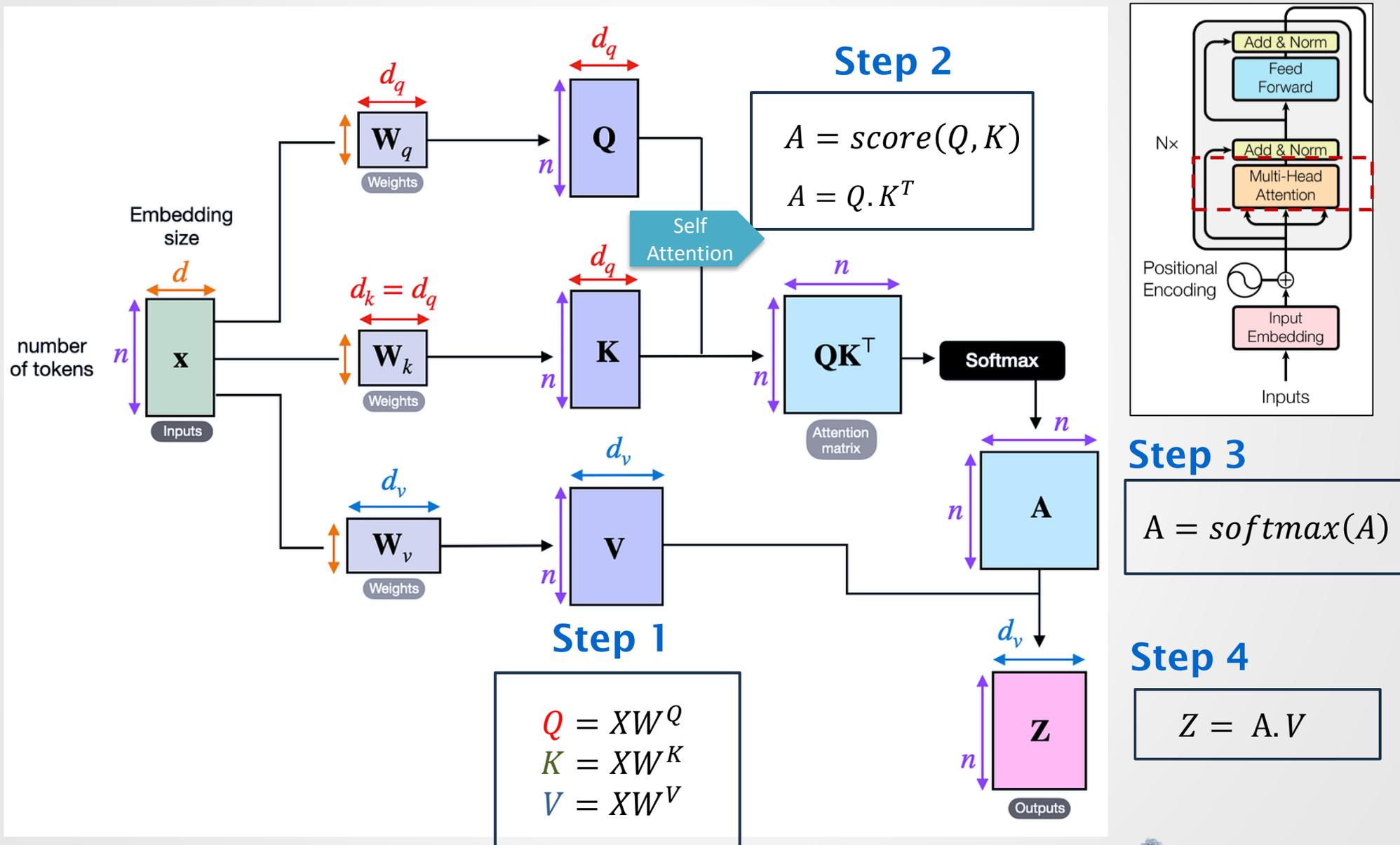


- $q_i = Q\tilde{h}_i$

Encoder-Decoder (Cross) Attention



(Compare) Encoder - Self Attention



¹Vaswani, Ashish, et al. "Attention is all you need." *NeurIPS* (2017).

Recap

- We have now covered all the primitives you need for building a transformer!
- All of these equations look abstract, but not to worry ...
- Assignment 2 was designed for you to implement all these concepts into a working MT model of your own.

Tasks	Section	Class	critierion	Max mark	Sub-Total	File
1	Building Blocks	LayerNorm	:forward	2	12	a2_transformer_model.py
2		MultiHeadAttention	:attention	4		
3			:forward	5		
4			FeedForwardLayer	:forward		
5	Architecture	TransformerEncoderLayer	:pre_layer_norm_forward	2	20	
6			:post_layer_norm_forward	1		
7			:__init__	4		
8		TransformerDecoderLayer	:pre_layer_norm_forward	2		
9			:post_layer_norm_forward	2		
10		TransformerDecoder	:forward	3		
11		TransformerEncoderDecoder	:create_pad_mask	1		
12			:create_causal_mask	2		
13			:forward	3		
15				greedy_decode		5

Transformers - Drawbacks

- Attention's **quadratic** computation **cost**
 - Function of sequence length N , and token dimension d
 - Computing all token pairs mean the function grows quadratically with N , $O(N^2d)$ unlike RNNs: $O(Nd)$

$$\begin{array}{ccc} XQ \in \mathbb{R}^{N \times d} & & \in \mathbb{R}^{N \times N} \\ \color{red}{\boxed{XQ}} & \color{purple}{\boxed{K^T X^T}} = & \color{lightpurple}{\boxed{XQK^T X^T}} \\ & \in \mathbb{R}^{d \times N} & \end{array}$$

- Can you see why this could be the biggest hurdle for increasing a transformer LM's **context length** (i.e., the size of input it can process)?

Transformers - Drawbacks

- Context (input) size limitation:
 - Dimension d in modern LLMs are $\sim > 3K$
 - If length of one sentence is $\sim 10-30$ word tokens, then computation scales with 10^2-30^2 times d
 - Thus, many encoders set a bound on N (usually 500-700 tokens)
 - But, some people want N to be much larger, e.g., processing a document ($N > 10K$) at one go (instead of chunking by N for every call)
- Active research area: improving the quadratic cost of attention, like self-attention with linear complexity^[1]
 - + Roformer, flash attention, sliding-window etc.

[1] Wang, Sinong, et al. "Linformer: Self-attention with linear complexity." (2020).

Transformers - Drawbacks

- Other drawbacks / improvement areas:
 - **Positional encoding representations:**
 - Do we need absolute indices to represent position?

$$f_{t:t \in \{q,k,v\}}(\mathbf{x}_i, i) := \mathbf{W}_{t:t \in \{q,k,v\}}(\mathbf{x}_i + \mathbf{p}_i).$$

- Slew of variants proposed to the (sinusoidal, absolute) positional encoding we saw
- General trend:
 - Move towards **relative position encoding**
 - E.g. Relative linear position attention [Shaw et al. 2018]

$$\begin{aligned} f_q(\mathbf{x}_m) &:= \mathbf{W}_q \mathbf{x}_m \\ f_k(\mathbf{x}_n, n) &:= \mathbf{W}_k(\mathbf{x}_n + \tilde{\mathbf{p}}_r^k) \\ f_v(\mathbf{x}_n, n) &:= \mathbf{W}_v(\mathbf{x}_n + \tilde{\mathbf{p}}_r^v) \end{aligned}$$

Relative distance between pos. m and n

$$r = \text{clip}(m - n, r_{\min}, r_{\max})$$

where $\tilde{\mathbf{p}}_r^k, \tilde{\mathbf{p}}_r^v \in \mathbb{R}^d$ are trainable relative position embeddings

Position (in)dependence

- To a great extent, attention is agnostic to order
 - For permutation vector v on $(1, 2, \dots, V)$:

$$\text{Att}(a, b_v) = \text{Att}(a, b_{1:V})$$

- But *word order* **matters** in language
- **Solution**: encode position in input:

$$x_s = T_F(F_s) + \phi(s)$$

- More on this soon...

Aside: Rotary Position Embeddings

- RoPE – perhaps the most common in modern LLMs
 - Encodes absolute position with a rotation matrix
 - RoPE + Transformer = **RoFormer**

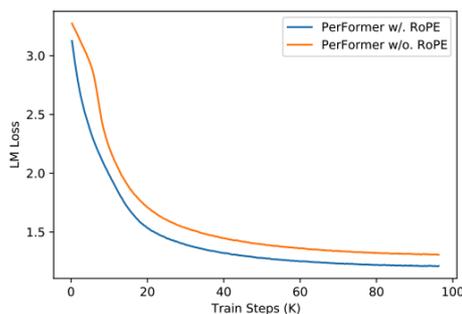
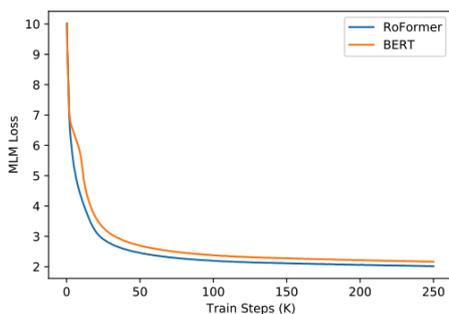
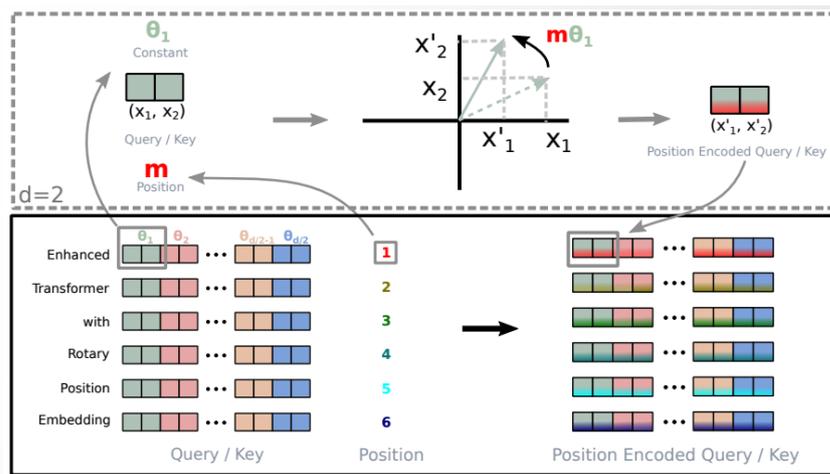
$$f_{\{q,k\}}(\mathbf{x}_m, m) = \mathbf{R}_{\Theta, m}^d \mathbf{W}_{\{q,k\}} \mathbf{x}_m$$

$$f_{\{q,k\}}(\mathbf{x}_m, m) = \begin{pmatrix} \cos m\theta & -\sin m\theta \\ \sin m\theta & \cos m\theta \end{pmatrix} \begin{pmatrix} W_{\{q,k\}}^{(11)} & W_{\{q,k\}}^{(12)} \\ W_{\{q,k\}}^{(21)} & W_{\{q,k\}}^{(22)} \end{pmatrix} \begin{pmatrix} x_m^{(1)} \\ x_m^{(2)} \end{pmatrix}$$

Model	BLEU
Transformer-base Vaswani et al. [2017]	27.3
RoFormer	27.5

Table 2: Comparing RoFormer and BERT by fine tuning on downstream GLEU tasks.

Model	MRPC	SST-2	QNLI	STS-B	QQP	MNLI(m/mm)
BERT Devlin et al. [2019]	88.9	93.5	90.5	85.8	71.2	84.6/83.4
RoFormer	89.5	90.7	88.0	87.0	86.4	80.2/79.8



[1] **RoPE**: Su, Jianlin, et al. "Roformer: Enhanced transformer with rotary position embedding." (2021)