

Trnasformers

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Lukasz Kaiser*
Google Brain
lukaszkaiser@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

Contents

1. Introduction

2. Model Architecture

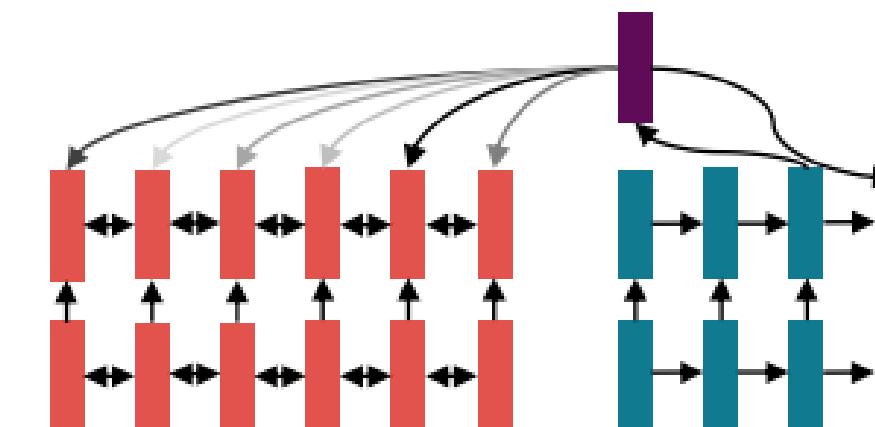
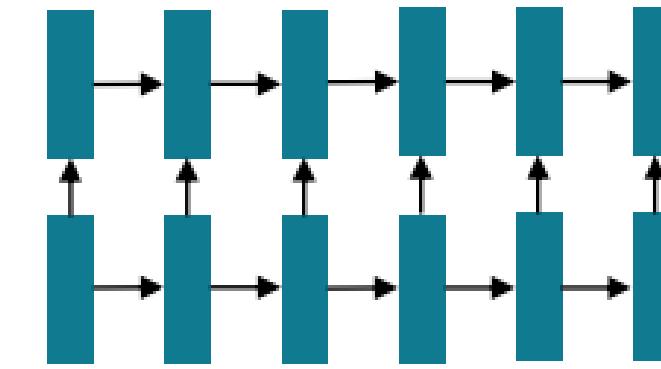
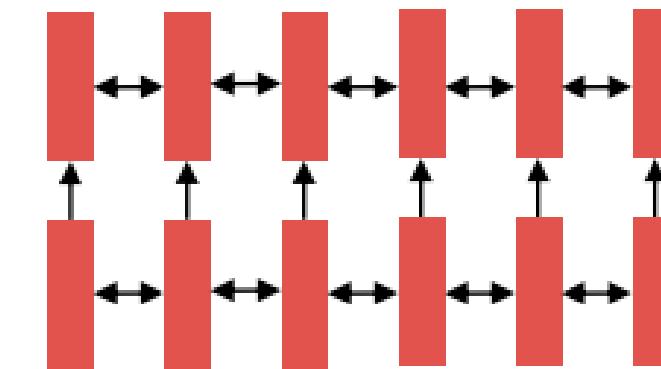
- Encoder and Decoder Stacks
- Attention
- Position-wise Feed-Forward Networks
- Embeddings and Softmax
- Positional Encoding

3. Training

4. Results

5. Conclusion

- Circa 2016, the de facto strategy in NLP is to **encode** sentences with a bidirectional LSTM:
(for example, the source sentence in a translation)
- Define your output (**parse**, sentence, summary) as a sequence, and use an LSTM to generate it.
- Use attention to allow flexible access to memory

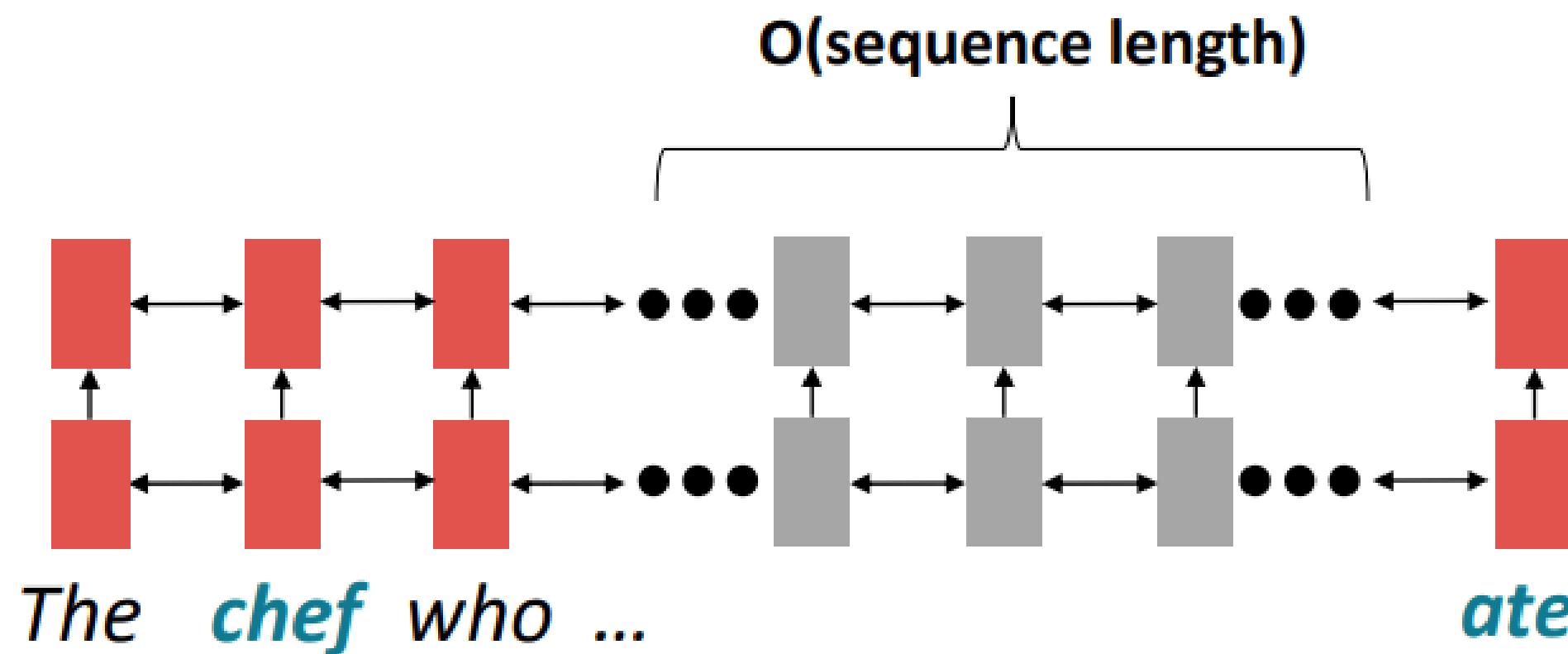
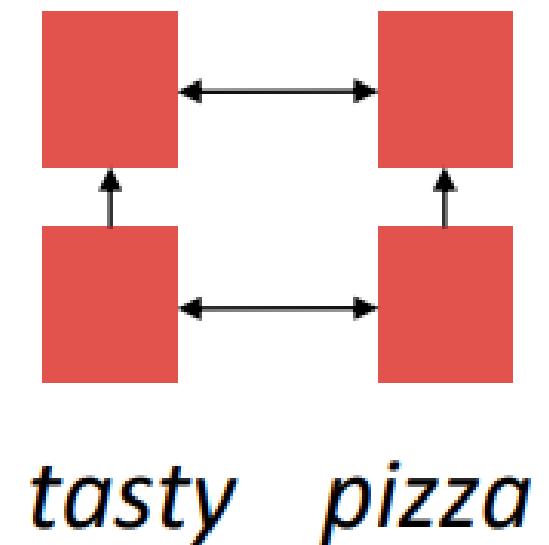


01

Introduction

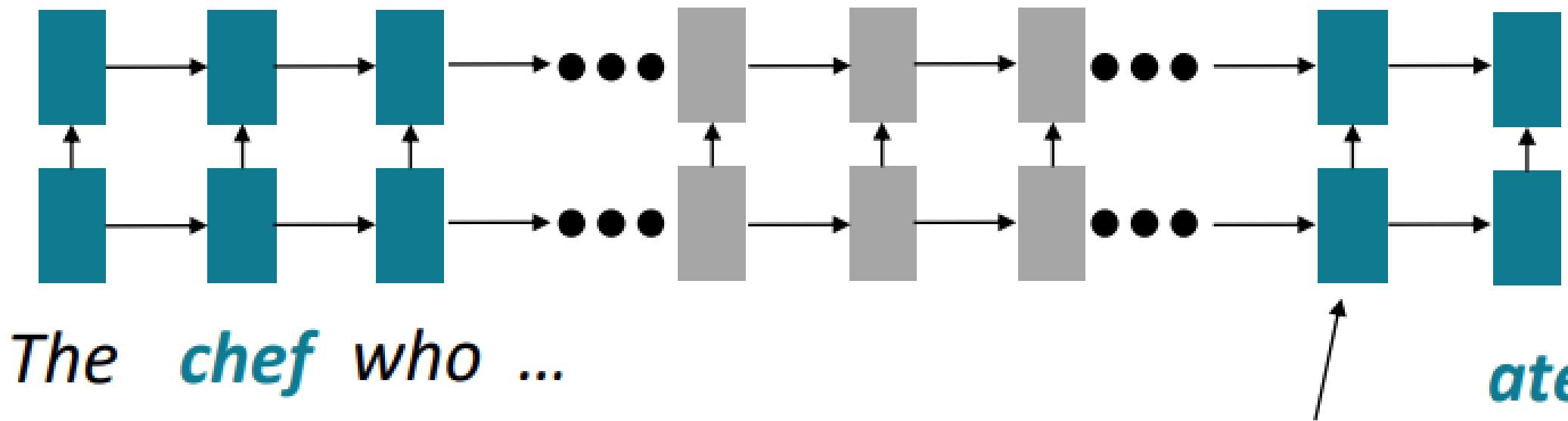
Issues with recurrent models: Linear interaction distance

- RNNs are unrolled “left-to-right”.
- It encodes linear locality: a useful heuristic!
 - Nearby words often affect each other’s meanings
- Problem: RNNs take $O(\text{sequence length})$ steps for distant word pairs to interact.



Issues with recurrent models: Linear interaction distance

- **O(sequence length)** steps for distant word pairs to interact means:
 - Hard to learn long-distance dependencies (because gradient problems!)
 - Linear order of words is “baked in”; we already know sequential structure doesn't tell the whole story...



*The **chef** who ...*

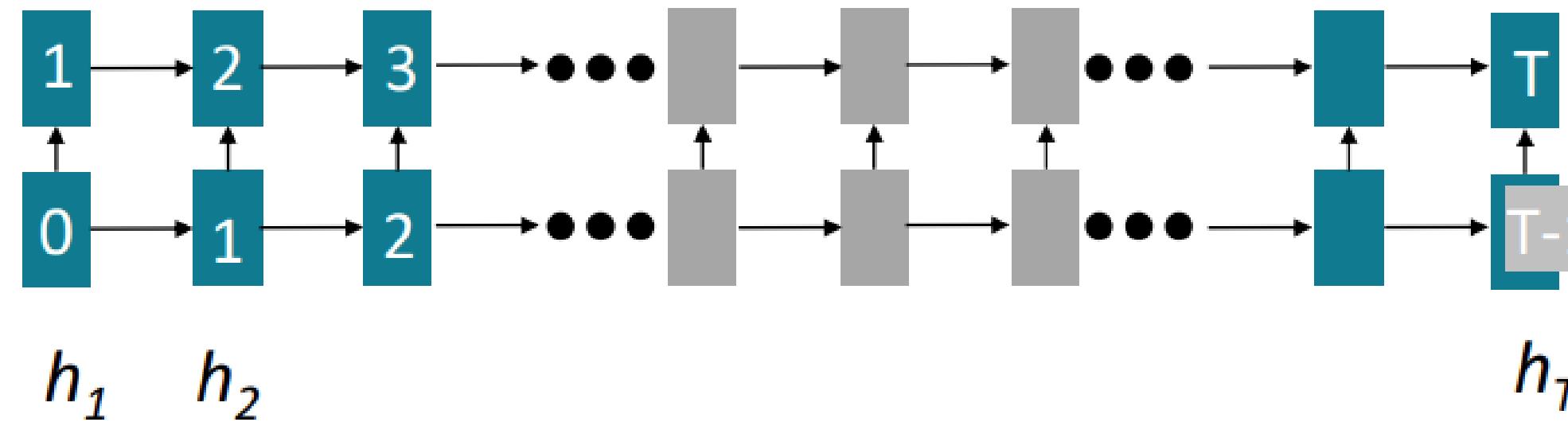
Info of **chef** has gone through
O(sequence length) many layers!

01

Introduction

Issues with recurrent models: Lack of parallelizability

- Forward and backward passes have **O(seq length)** unparallelizable operations
 - GPUs (and TPUs) can perform many independent computations at once!
 - But future RNN hidden states can't be computed in full before past RNN hidden states have been computed
 - Inhibits training on very large datasets!
 - Particularly problematic as sequence length increases, as we can no longer batch many examples together due to memory limitations



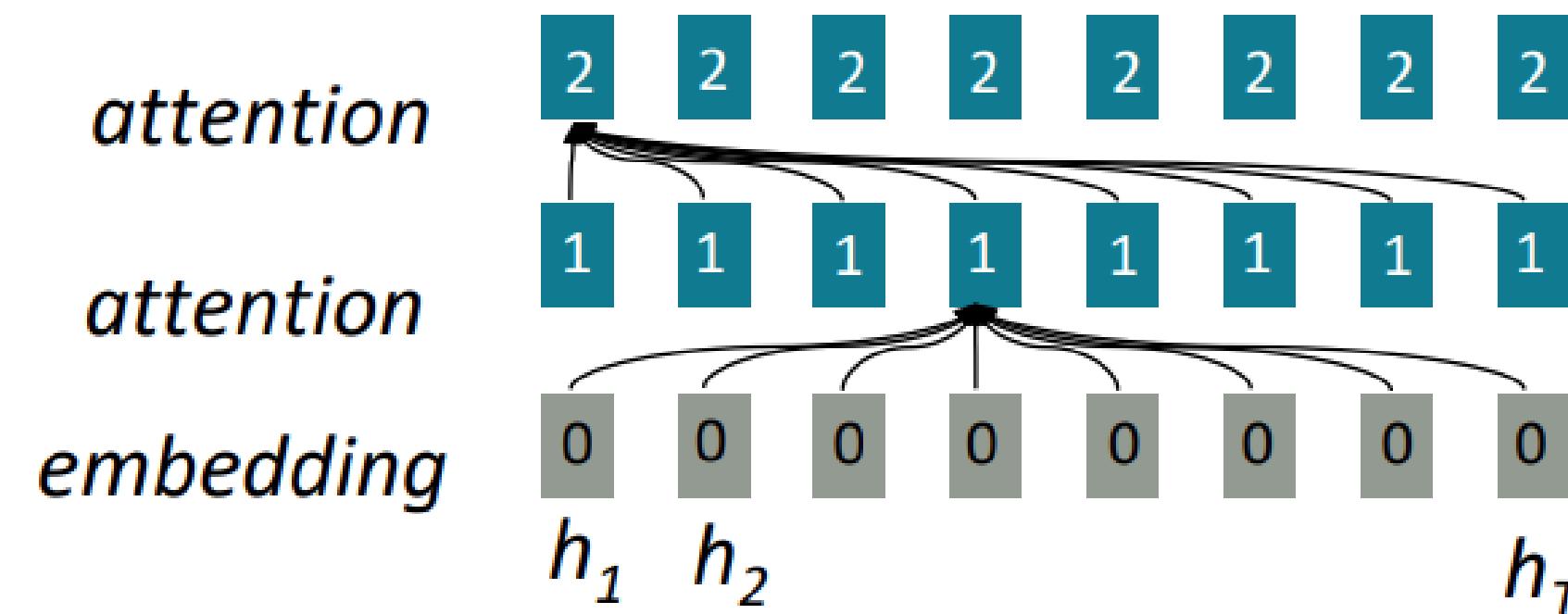
Numbers indicate min # of steps before a state can be computed

01

Introduction

If not recurrence, then what? How about (self) attention?

- To recap, **attention** treats each word's representation as a **query** to access and incorporate information from **a set of values**.
 - Last week, we saw attention from the **decoder** to the **encoder**;
 - **Self-attention** is **encoder-encoder** (or **decoder-decoder**) attention where each word attends to each other word **within the input (or output)**.



All words attend to all words in previous layer; most arrows here are omitted

Transformer Advantages:

- Number of unparallelizable operations does not increase with sequence length.
- Each "word" interacts with each other, so maximum interaction distance: $O(1)$.

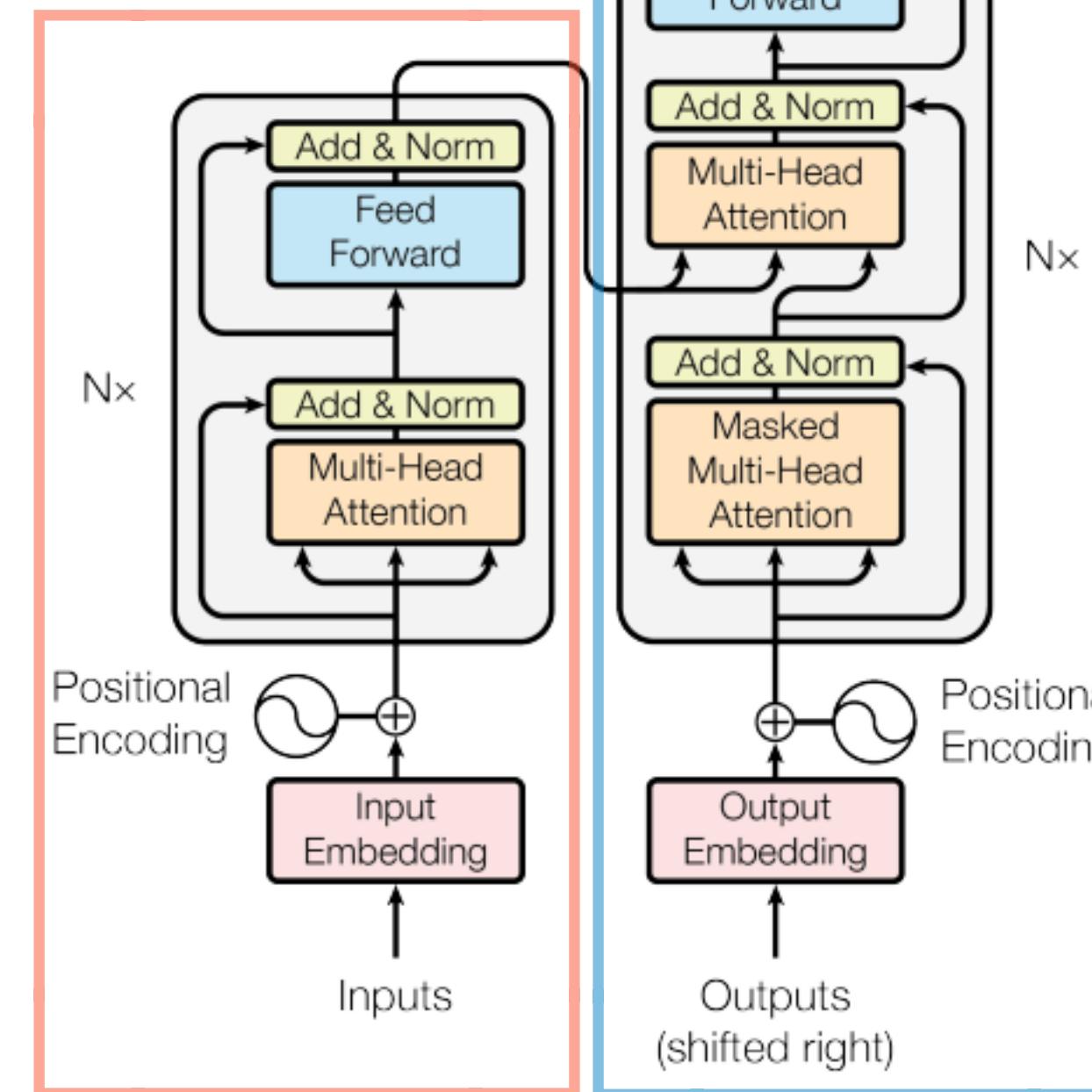
Table 1: Maximum path lengths, per-layer complexity and minimum number of sequential operations for different layer types. n is the sequence length, d is the representation dimension, k is the kernel size of convolutions and r the size of the neighborhood in restricted self-attention.

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

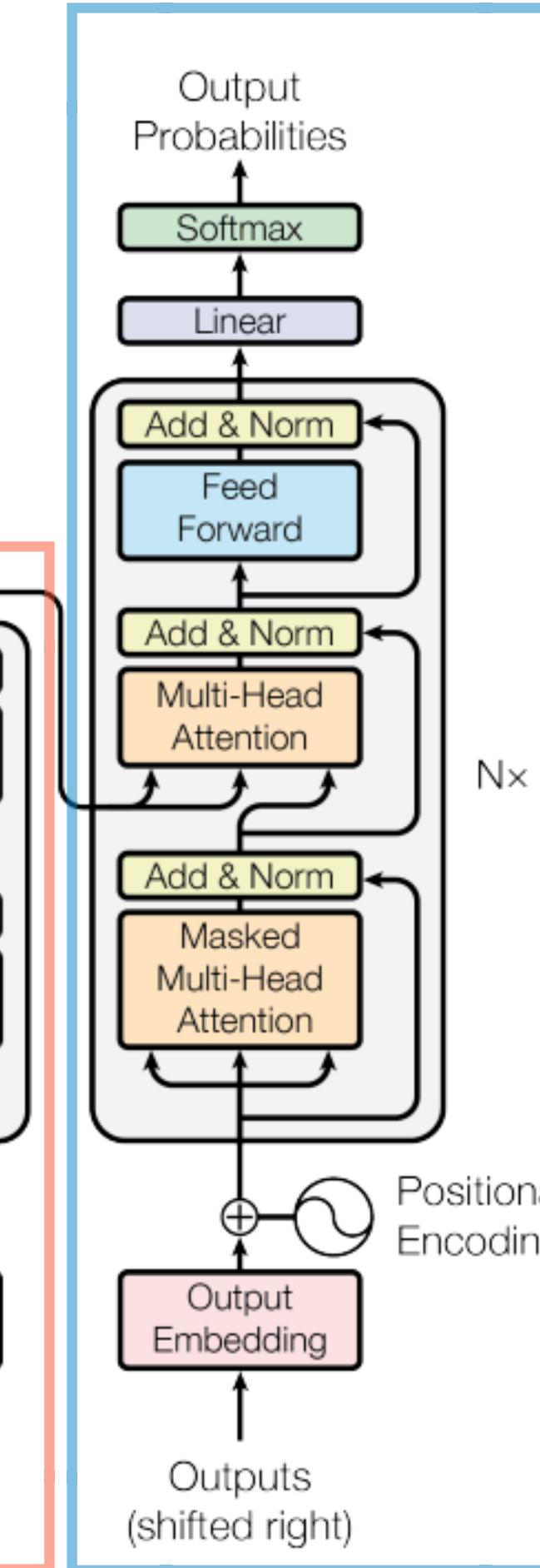
02 | Model Architecture

Encoder and Decoder Stacks

Encoder



Decoder

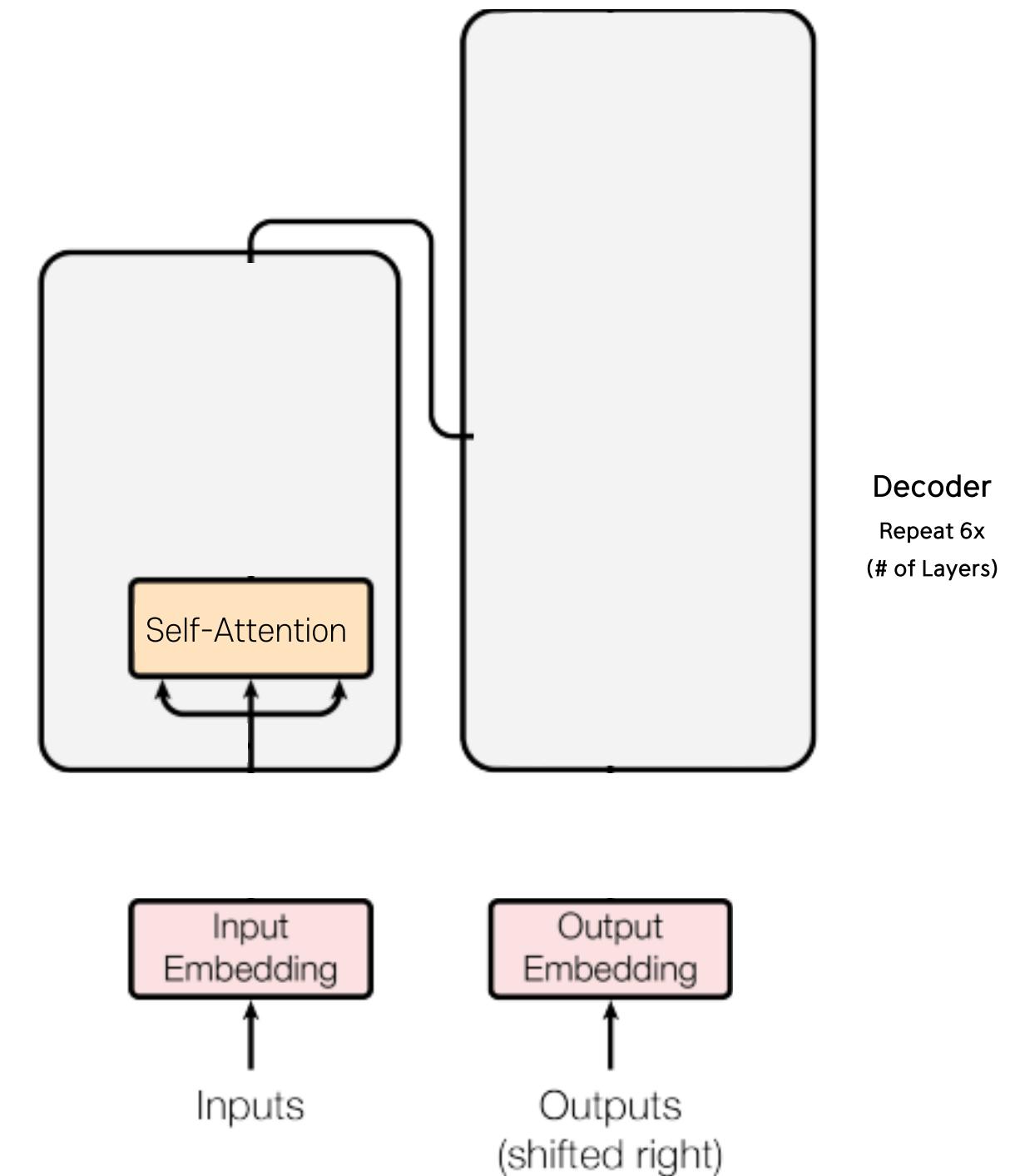


02 | Model Architecture

Encoder: Self-Attention

Output
Probabilities

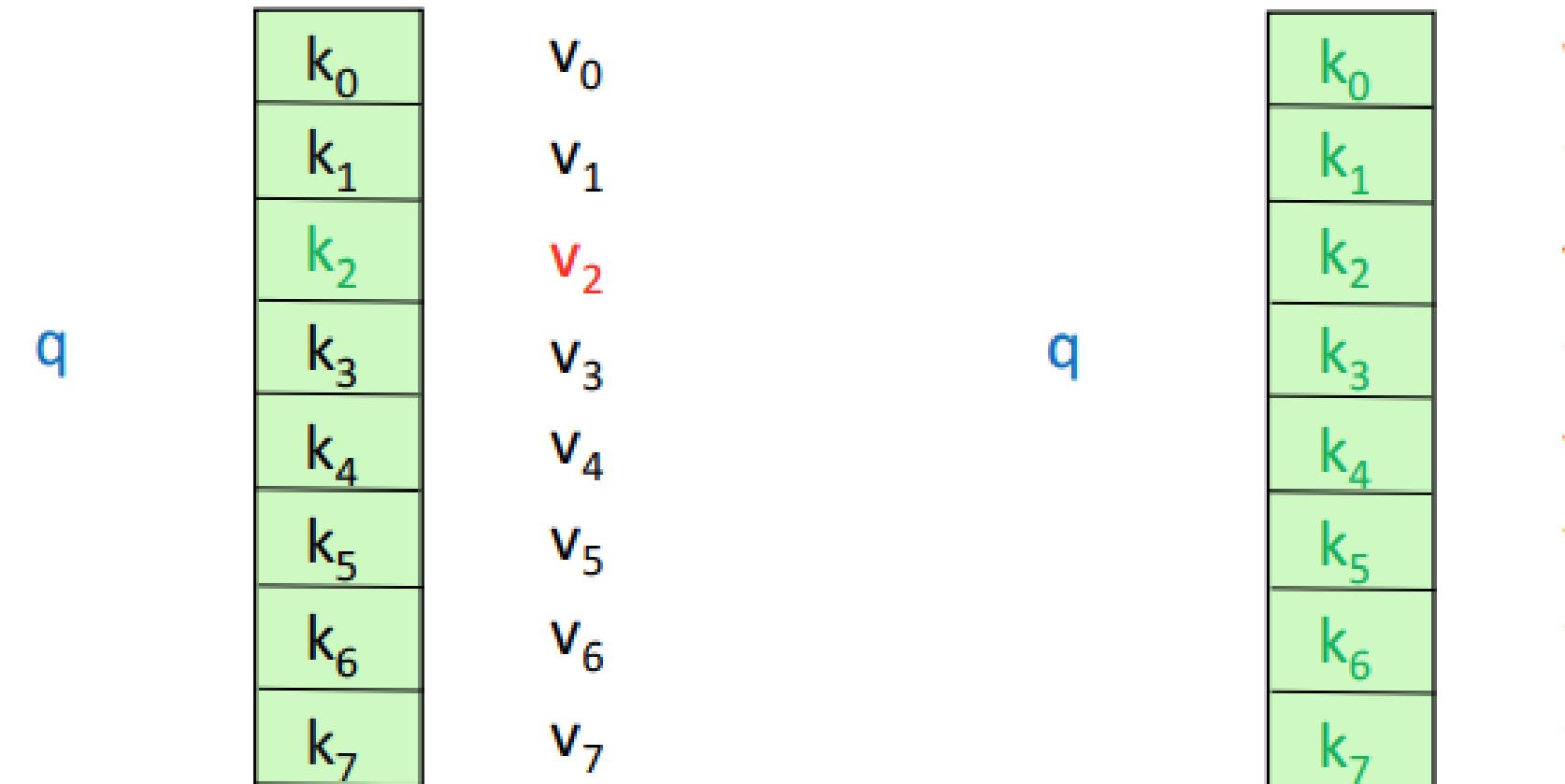
- **Self-Attention is the core building block of Transformer,
so let's first focus on that!**



02 | Model Architecture

Intuition for Attention Mechanism

- Let's think of attention as a "fuzzy" or approximate hashtable:
 - To look up a **value**, we compare a **query** against **keys** in a table.
 - In a hashtable (shown on the bottom left):
 - Each **query** (hash) maps to exactly one **key-value** pair.
 - In (self-)attention (shown on the bottom right):
 - Each **query** matches each **key** to varying degrees.
 - We return a sum of **values** weighted by the **query-key** match.



02 | Model Architecture

Recipe for Self-Attention in the Transformer Encoder

- Step 1: For each word x_i , calculate its **query**, **key**, and **value**.

$$q_i = W^Q x_i \quad k_i = W^K x_i \quad v_i = W^V x_i$$

- Step 2: Calculate attention score between **query** and **keys**.

$$e_{ij} = q_i \cdot k_j$$

- Step 3: Take the softmax to normalize attention scores.

$$\alpha_{ij} = \text{softmax}(e_{ij}) = \frac{\exp(e_{ij})}{\sum_k \exp(e_{ik})}$$

- Step 4: Take a weighted sum of **values**.

$$\text{Output}_i = \sum_j \alpha_{ij} v_j$$

k_0	v_0
k_1	v_1
k_2	v_2
k_3	v_3
k_4	v_4
k_5	v_5
k_6	v_6
k_7	v_7

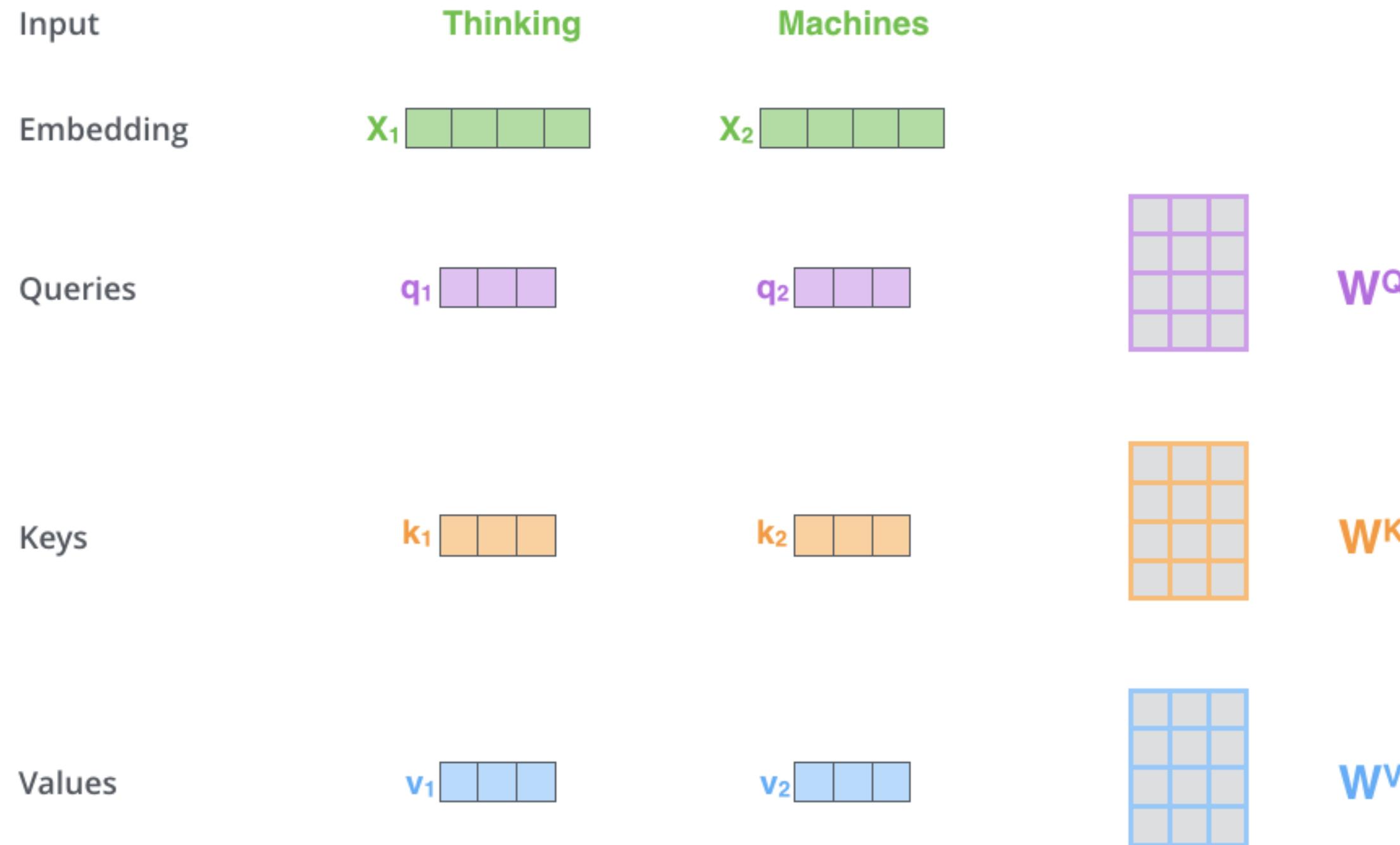
02 | Model Architecture

Self-Attention in Detail

$$\text{Output} = \text{softmax}(QK^T)V$$

- Step 1: For each word x_i , calculate its **query**, **key**, and **value**.

$$q_i = W^Q x_i \quad k_i = W^K x_i \quad v_i = W^V x_i$$



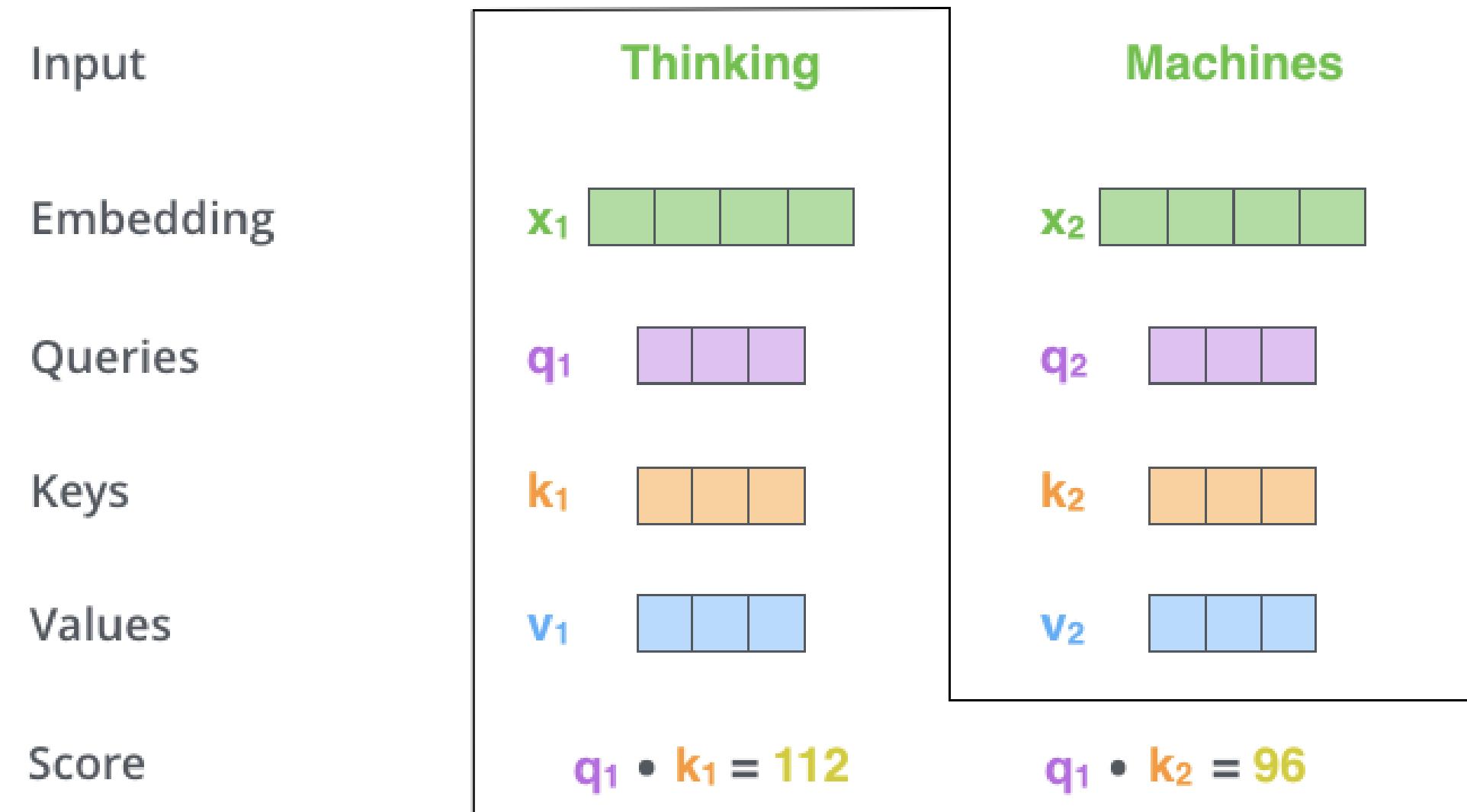
02 | Model Architecture

Self-Attention in Detail

$$\text{Output} = \text{softmax}(QK^T)V$$

- Step 2: Calculate attention score between **query** and **keys**.

$$e_{ij} = q_i \cdot k_j$$



02 | Model Architecture

Self-Attention in Detail

$$\text{Output} = \text{softmax}(QK^T)V$$

- Step 3: Take the softmax to normalize attention scores.

$$\alpha_{ij} = \text{softmax}(e_{ij}) = \frac{\exp(e_{ij})}{\sum_k \exp(e_{ik})}$$

Input	Thinking		Machines	
Embedding	x_1	[5 green squares]	x_2	[5 green squares]
Queries	q_1	[3 purple squares]	q_2	[3 purple squares]
Keys	k_1	[3 orange squares]	k_2	[3 orange squares]
Values	v_1	[3 blue squares]	v_2	[3 blue squares]
Score	$q_1 \cdot k_1 = 112$		$q_1 \cdot k_2 = 96$	
Softmax	9.99999887e-01		1.12535162e-07	

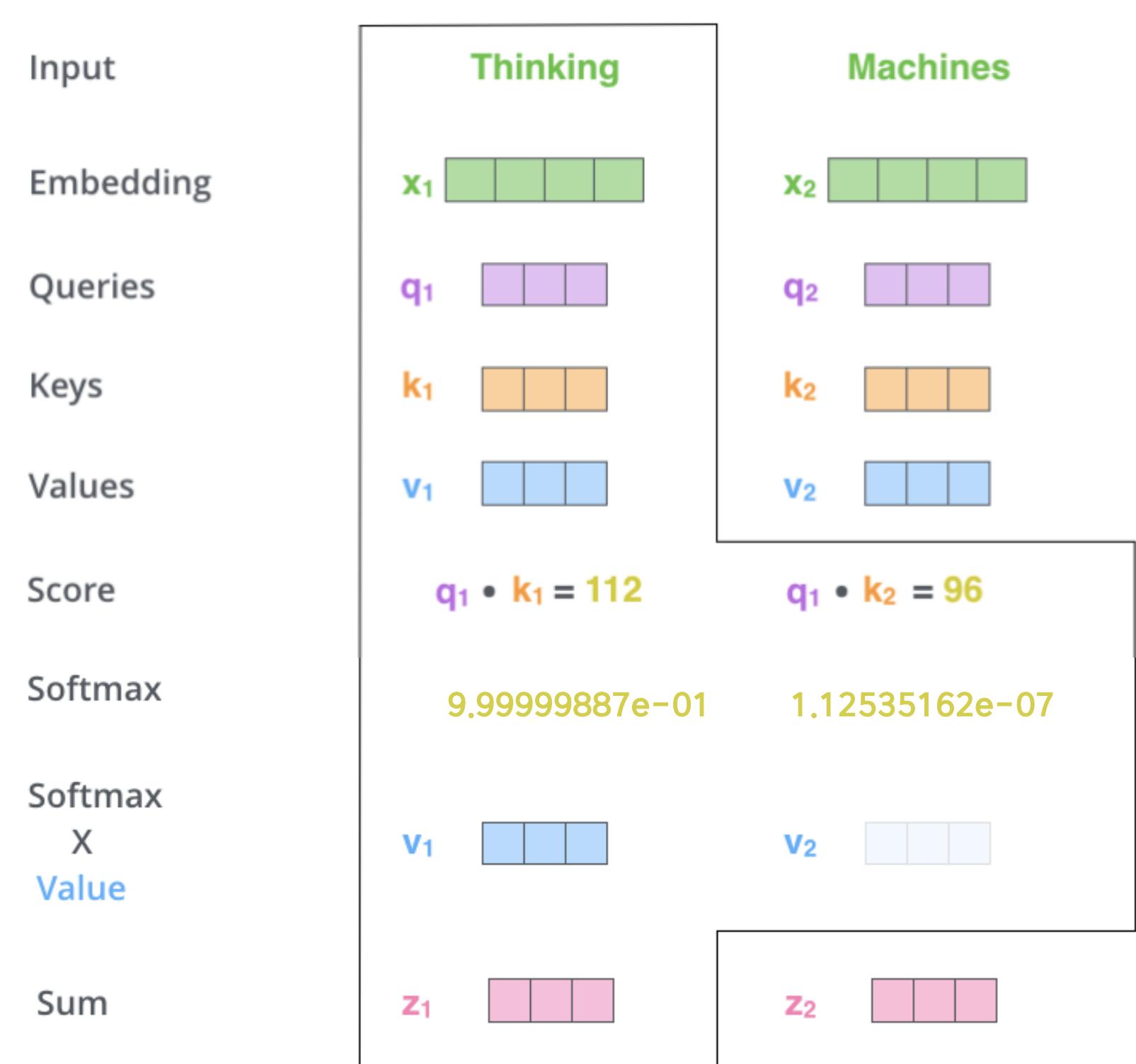
02 | Model Architecture

Self-Attention in Detail

$$Output = \text{softmax}(QK^T)V$$

- Step 4: Take a weighted sum of **values**.

$$Output_i = \sum_j \alpha_{ij} v_j$$



02 | Model Architecture

Recipe for (Vectorized) Self-Attention in the Transformer Encoder

- Step 1: With embeddings stacked in X , calculate **queries**, **keys**, and **values**.

$$Q = XW^Q \quad K = XW^K \quad V = XW^V$$

- Step 2: Calculate attention scores between **query** and **keys**.

$$E = QK^T$$

- Step 3: Take the softmax to normalize attention scores.

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

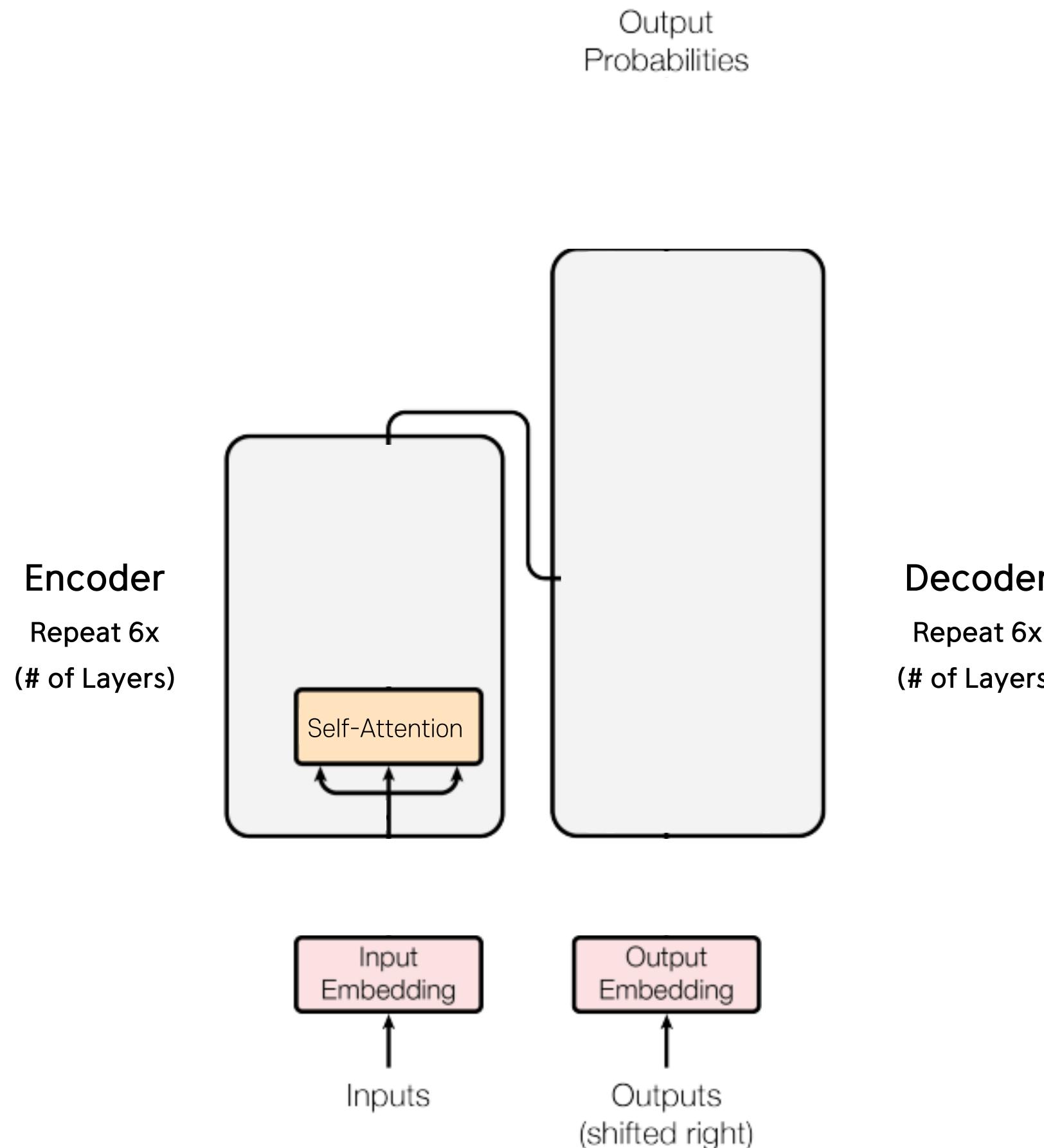
- Step 4: Take a weighted sum of **values**.

$$\text{Output} = AV$$

$$\text{Output} = \text{softmax}(QK^T)V$$

02 | Model Architecture

What We Have So Far: (Encoder) Self-Attention!



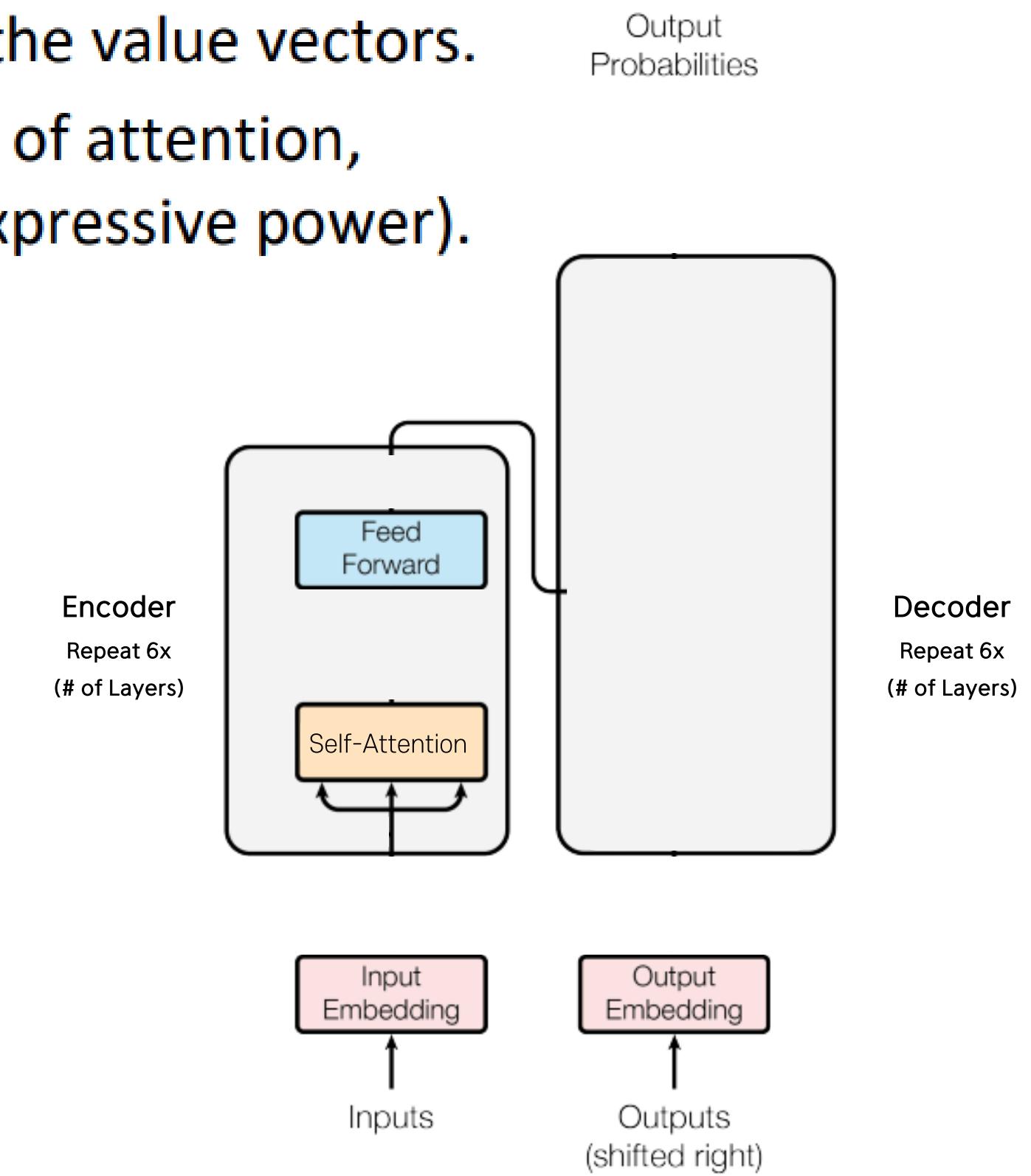
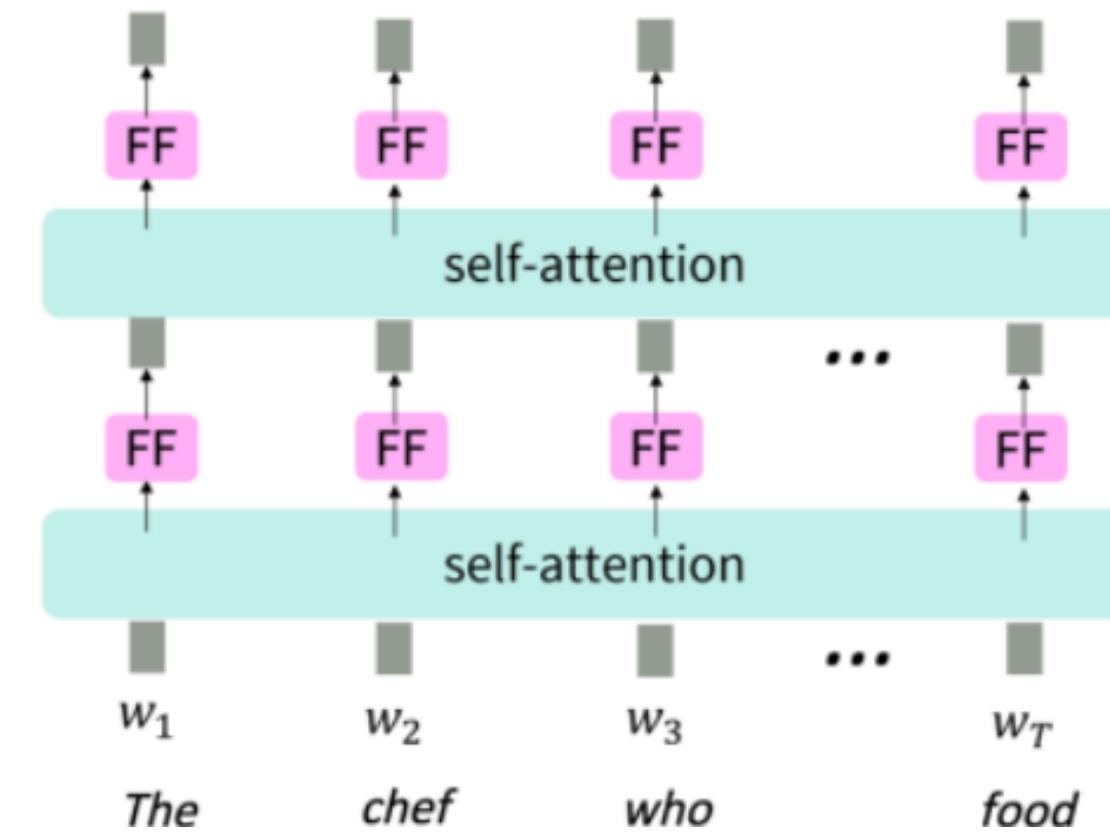
02 | Model Architecture

But attention isn't quite all you need!

- **Problem:** Since there are no element-wise non-linearities, self-attention is simply performing a re-averaging of the value vectors.
- **Easy fix:** Apply a feedforward layer to the output of attention, providing non-linear activation (and additional expressive power).

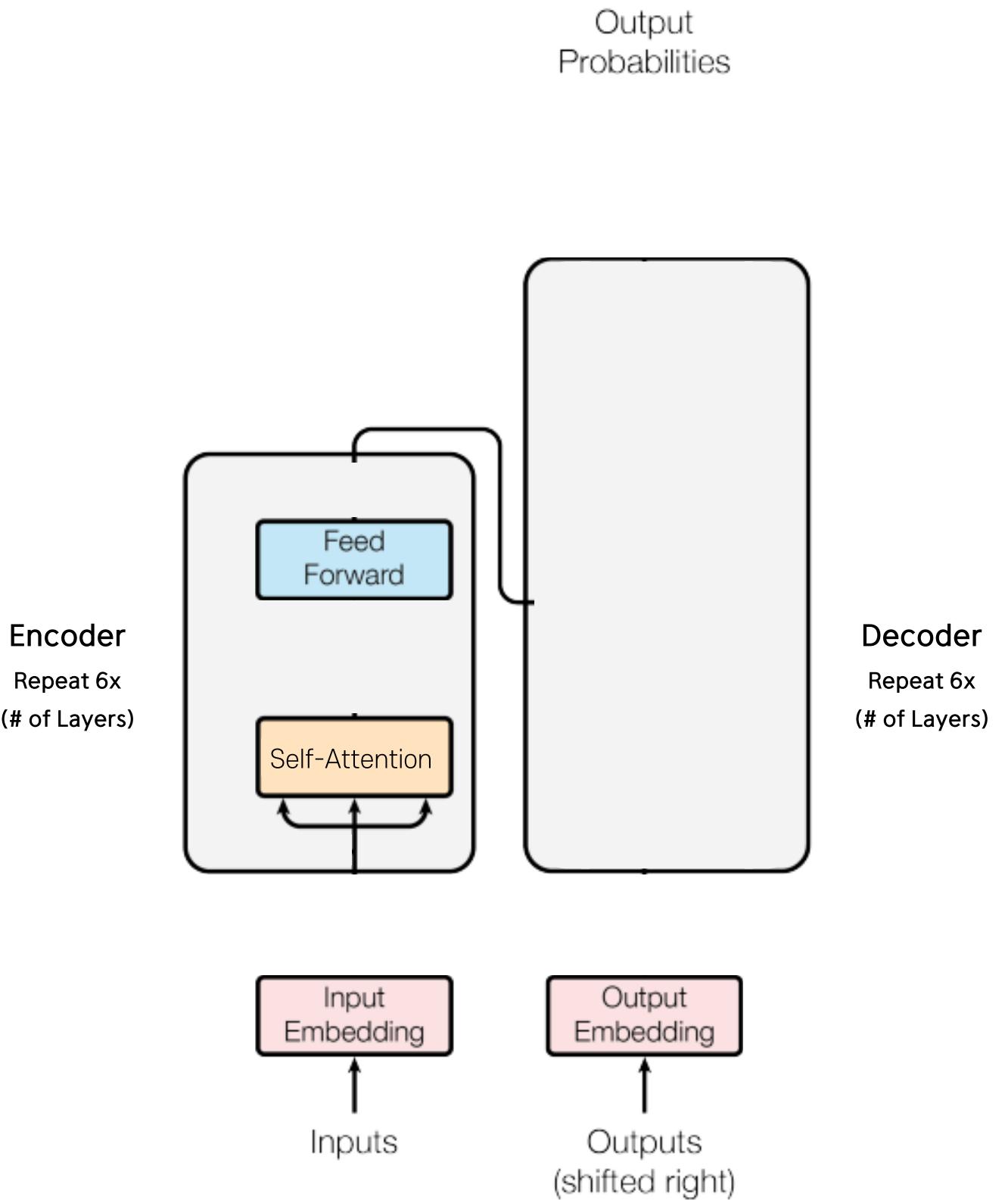
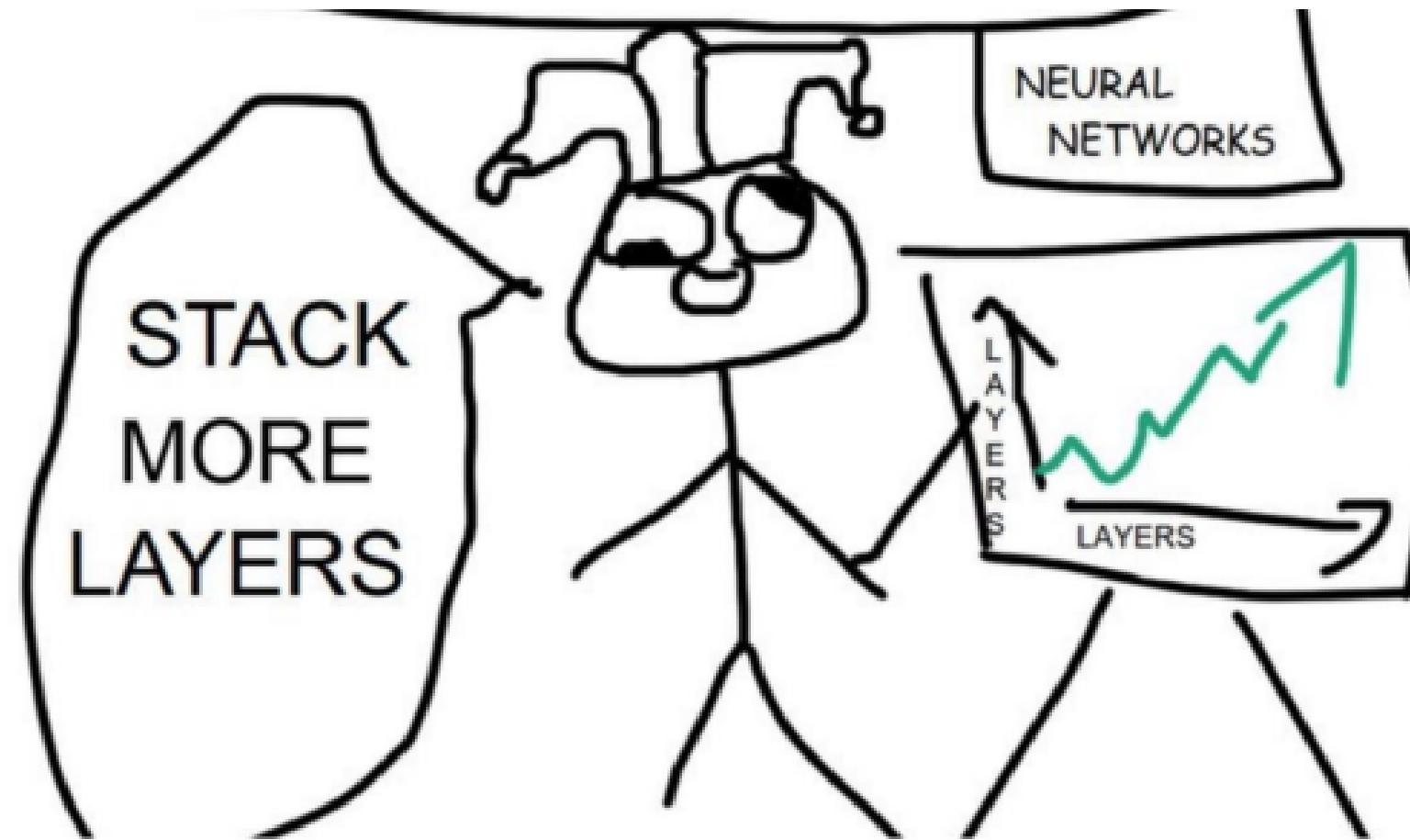
Equation for Feed Forward Layer

$$\begin{aligned} m_i &= \text{MLP}(\text{output}_i) \\ &= W_2 * \text{ReLU}(W_1 \times \text{output}_i + b_1) + b_2 \end{aligned}$$



02 | Model Architecture

But how do we make this work for deep networks?



Training Trick #1: Residual Connections

Training Trick #2: LayerNorm

Training Trick #3: Scaled Dot Product Attention

02

Model Architecture

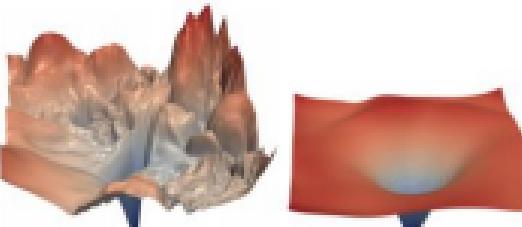
Training Trick #1: Residual Connections [He et al., 2016]

- Residual connections are a simple but powerful technique from computer vision.
- Deep networks are surprisingly bad at learning the identity function!
- Therefore, directly passing "raw" embeddings to the next layer can actually be very helpful!

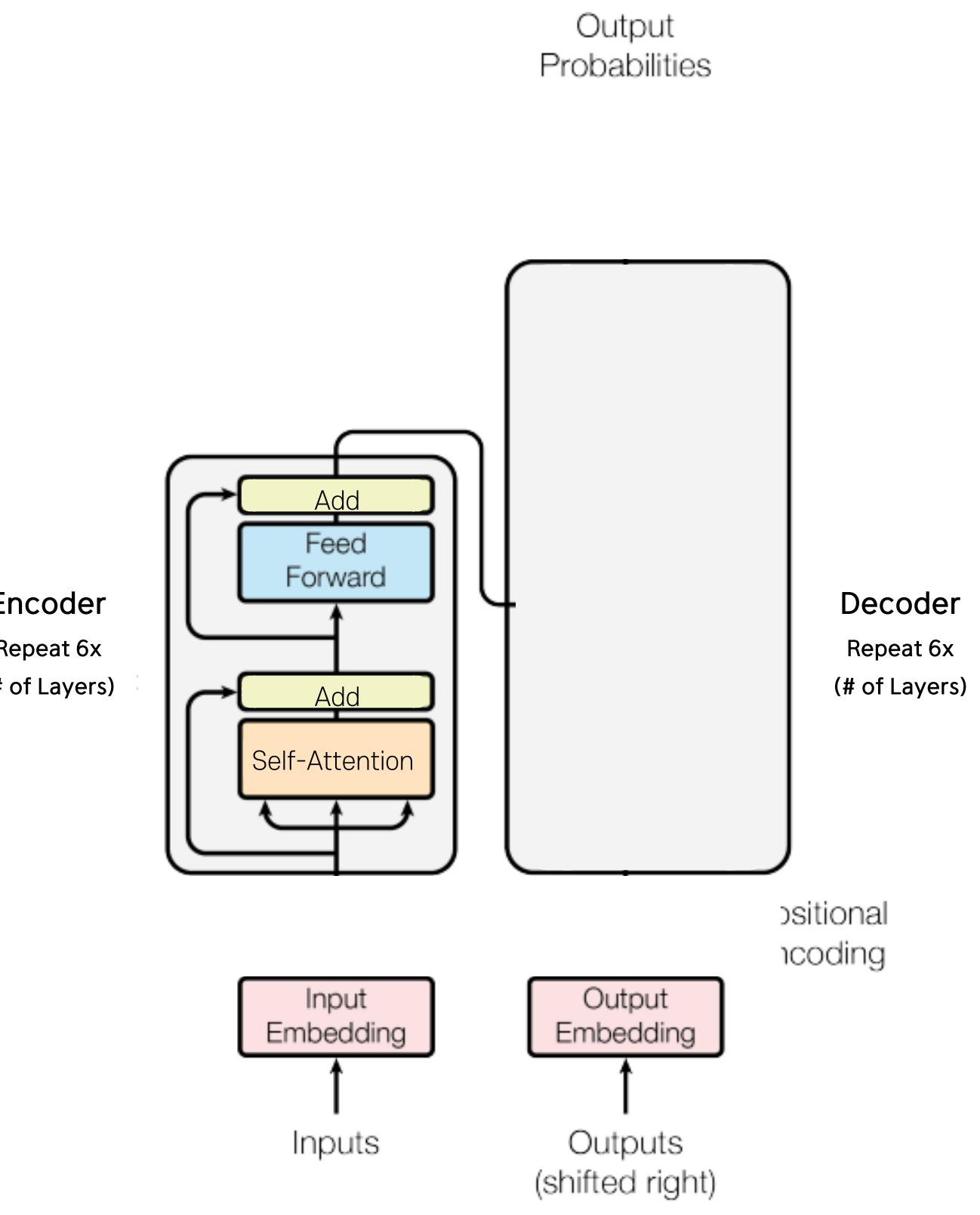
$$x_\ell = F(x_{\ell-1}) + x_{\ell-1}$$

- This prevents the network from "forgetting" or distorting important information as it is processed by many layers.

Residual connections are also thought to smooth the loss landscape and make training easier!



[Loss landscape visualization,
Li et al., 2018, on a ResNet]



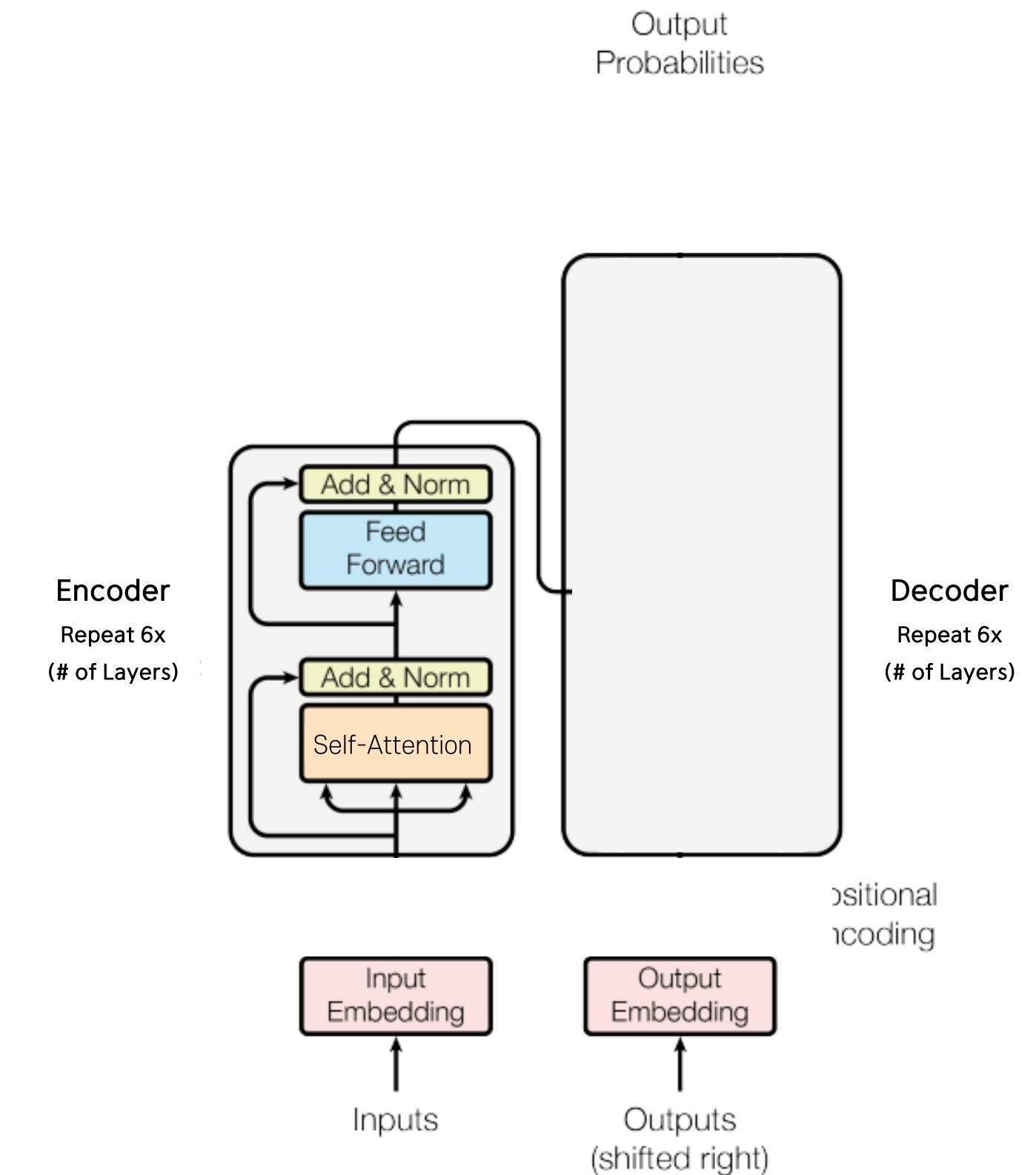
02 | Model Architecture

Training Trick #2: Layer Normalization [Ba et al., 2016]

- **Problem:** Difficult to train the parameters of a given layer because its input from the layer beneath keeps shifting.
- **Solution:** Reduce uninformative variation by **normalizing** to zero mean and standard deviation of one within each **layer**.

$$\text{Mean: } \mu^l = \frac{1}{H} \sum_{i=1}^H a_i^l \quad \text{Standard Deviation: } \sigma^l = \sqrt{\frac{1}{H} \sum_{i=1}^H (a_i^l - \mu^l)^2}$$

$$x'^l = \frac{x^l - \mu^l}{\sigma^l + \epsilon}$$



02

Model Architecture

Training Trick #3: Scaled Dot Product Attention

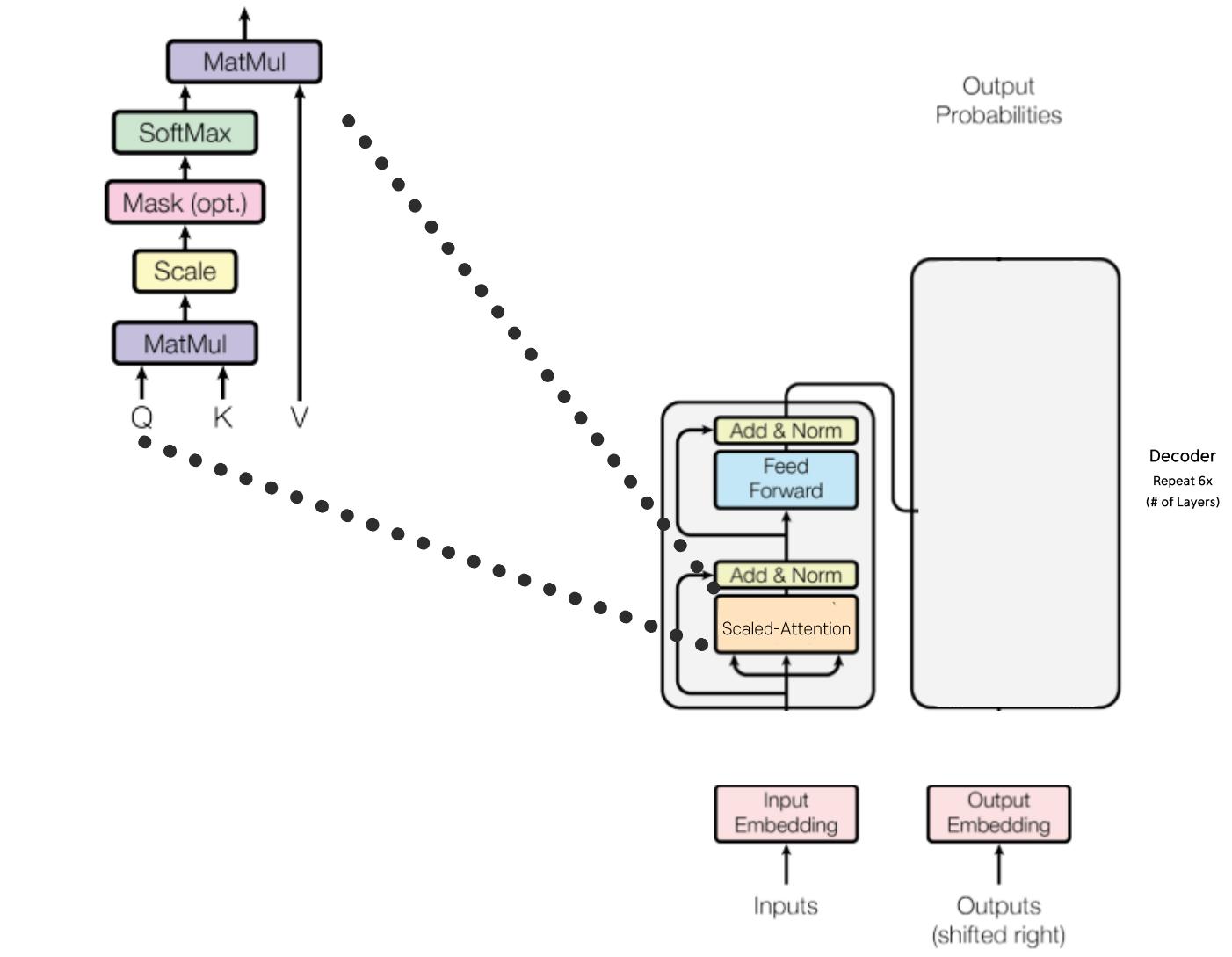
- After LayerNorm, the mean and variance of vector elements is 0 and 1, respectively.
- However, the dot product still tends to take on extreme values, as its variance scales with dimensionality d_k

Quick Statistics Review:

- Mean of sum = sum of means = $d_k * 0 = 0$
- Variance of sum = sum of variances = $d_k * 1 = d_k$
- To set the variance to 1, simply divide by $\sqrt{d_k}$!

Updated Self-Attention Equation:

$$\text{Output} = \text{softmax}\left(QK^T / \sqrt{d_k}\right)V$$



02 | Model Architecture

Major issue

- We're almost done with the Encoder, but we have a major problem!

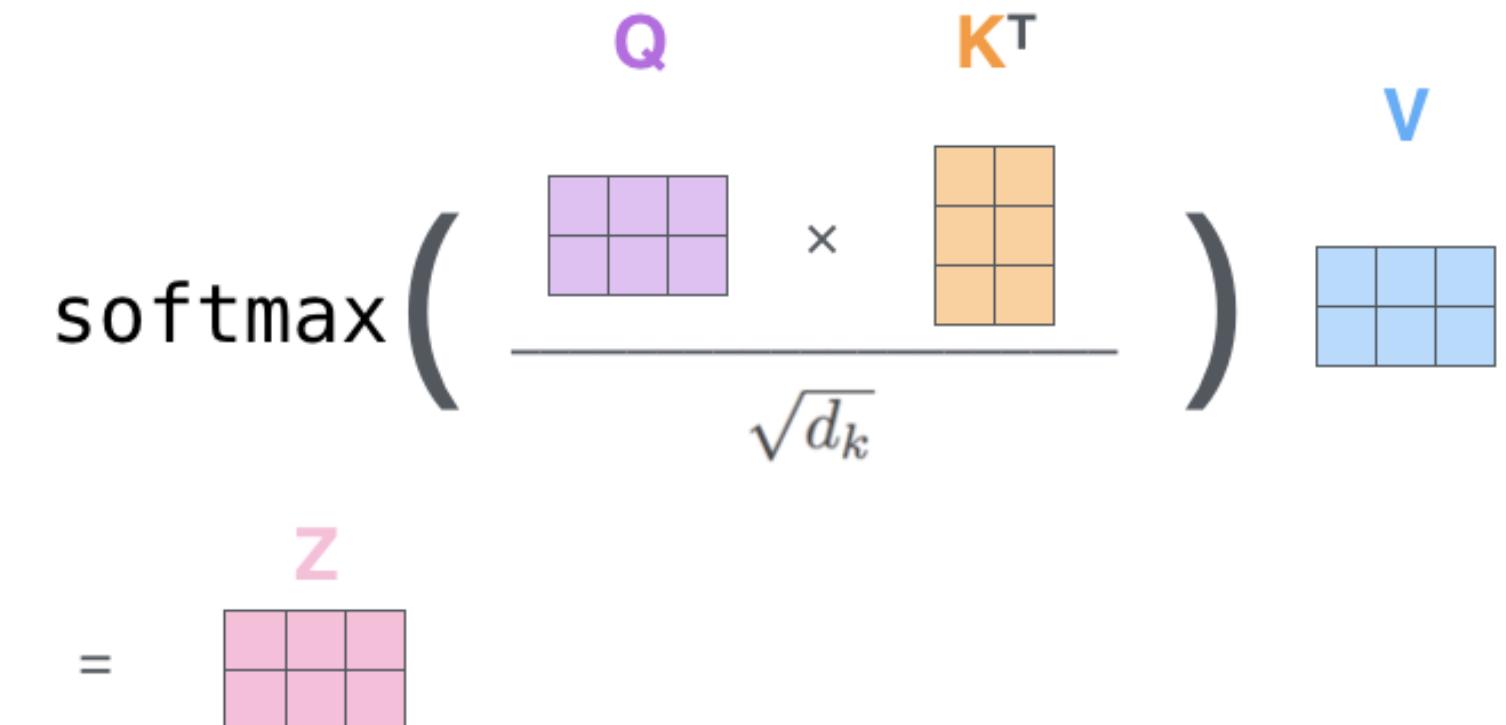
Has anyone spotted it?

- Consider this sentence:

"Man eats small dinosaur."

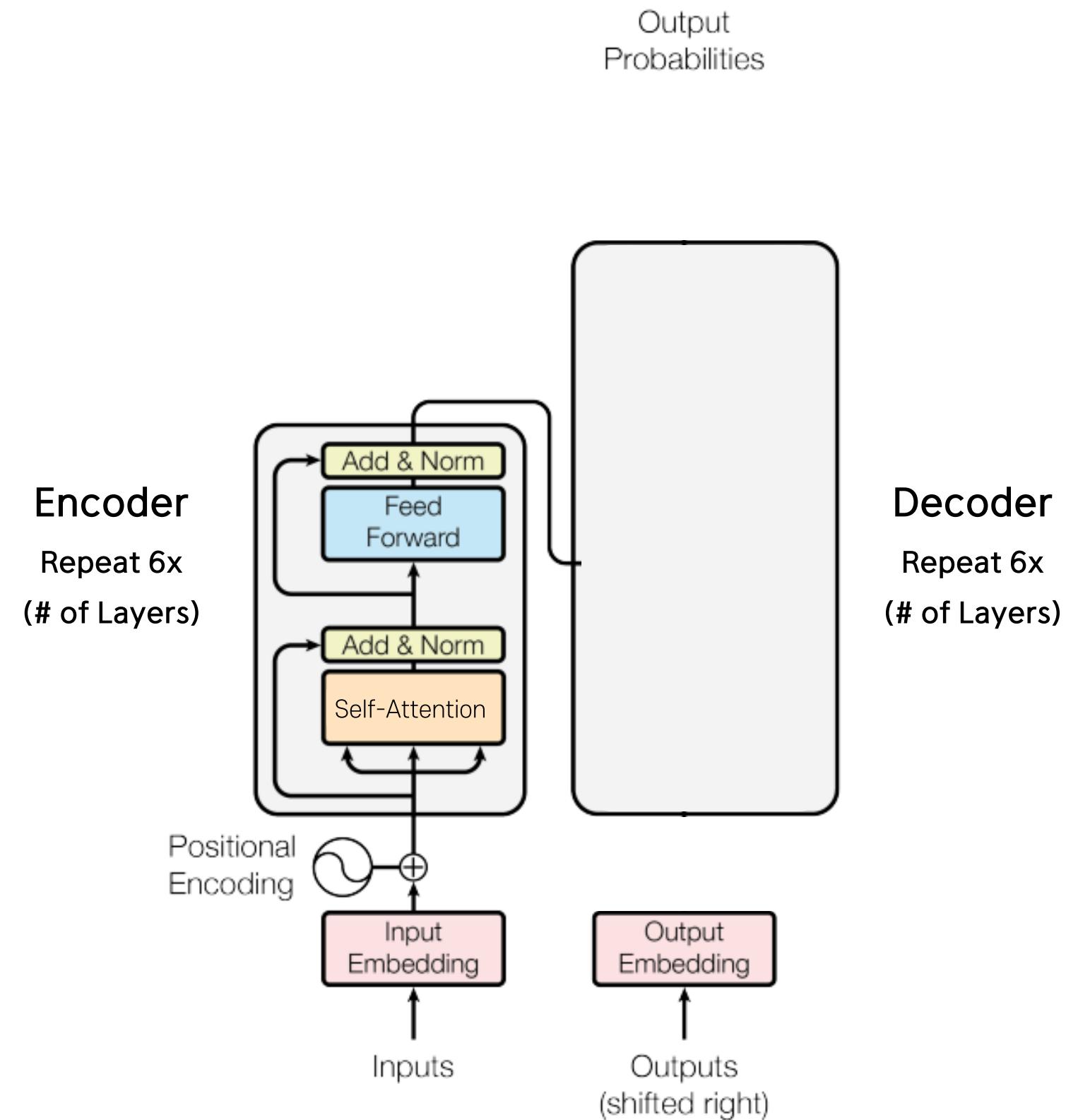
- Wait a minute, order doesn't impact the network at all!
- This seems wrong given that word order does have meaning in many languages, including English!

$$\boxed{Output = \text{softmax}\left(QK^T / \sqrt{d_k}\right)V}$$

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


02 | Model Architecture

Solution: Inject Order Information through Positional Encodings



Fixing the first self-attention problem: sequence order

- Since self-attention doesn't build in order information, we need to encode the order of the sentence in our keys, queries, and values.
- Consider representing each **sequence index** as a **vector**

$p_i \in \mathbb{R}^d$, for $i \in \{1, 2, \dots, T\}$ are position vectors

- Don't worry about what the p_i are made of yet!
- Easy to incorporate this info into our self-attention block: just add the p_i to our inputs!
- Let $\tilde{v}_i, \tilde{k}_i, \tilde{q}_i$ be our old values, keys, and queries.

$$\begin{aligned} v_i &= \tilde{v}_i + p_i \\ q_i &= \tilde{q}_i + p_i \\ k_i &= \tilde{k}_i + p_i \end{aligned}$$

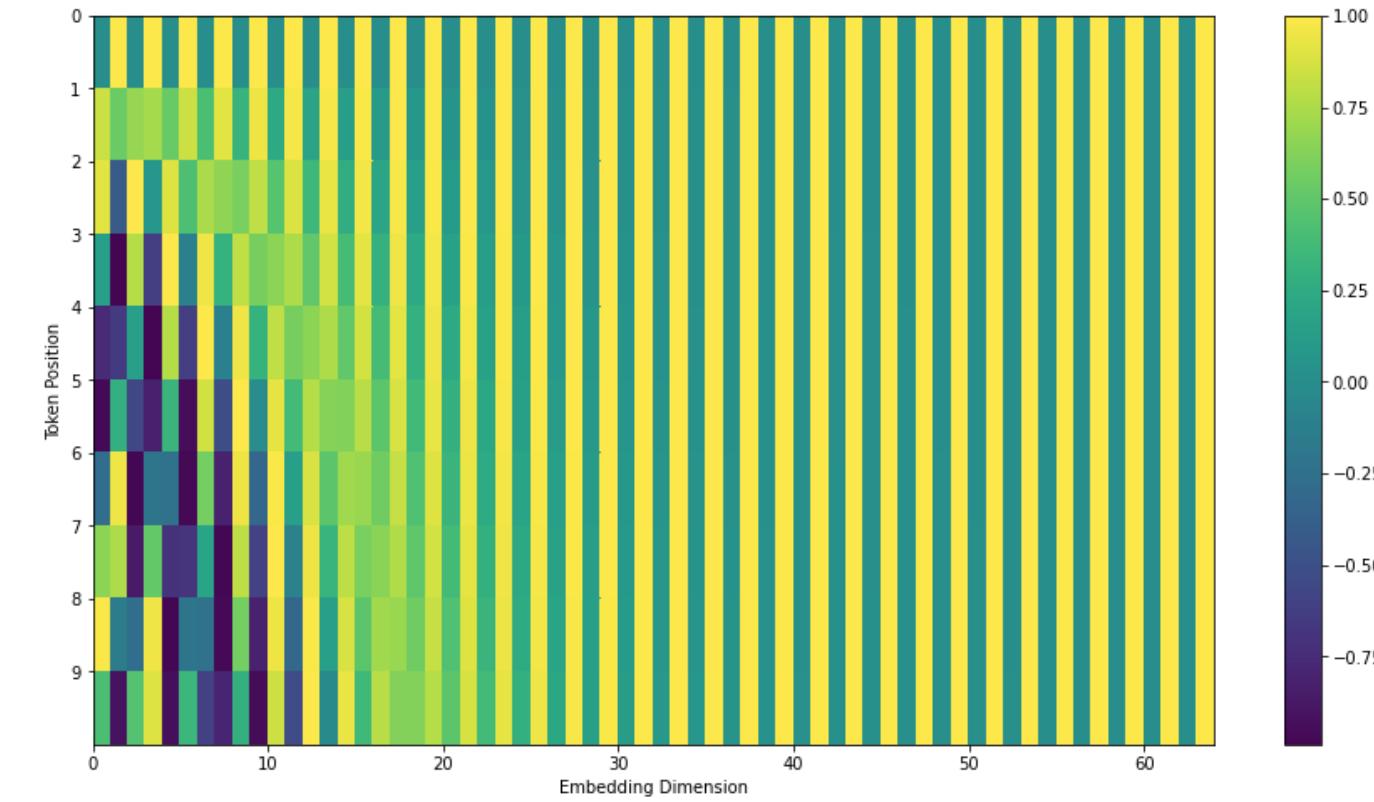
In deep self-attention networks, we do this at the first layer! You could concatenate them as well, but people mostly just add...

02 | Model Architecture

Position representation vectors through sinusoids

- **Sinusoidal position representations:** concatenate sinusoidal functions of varying periods:

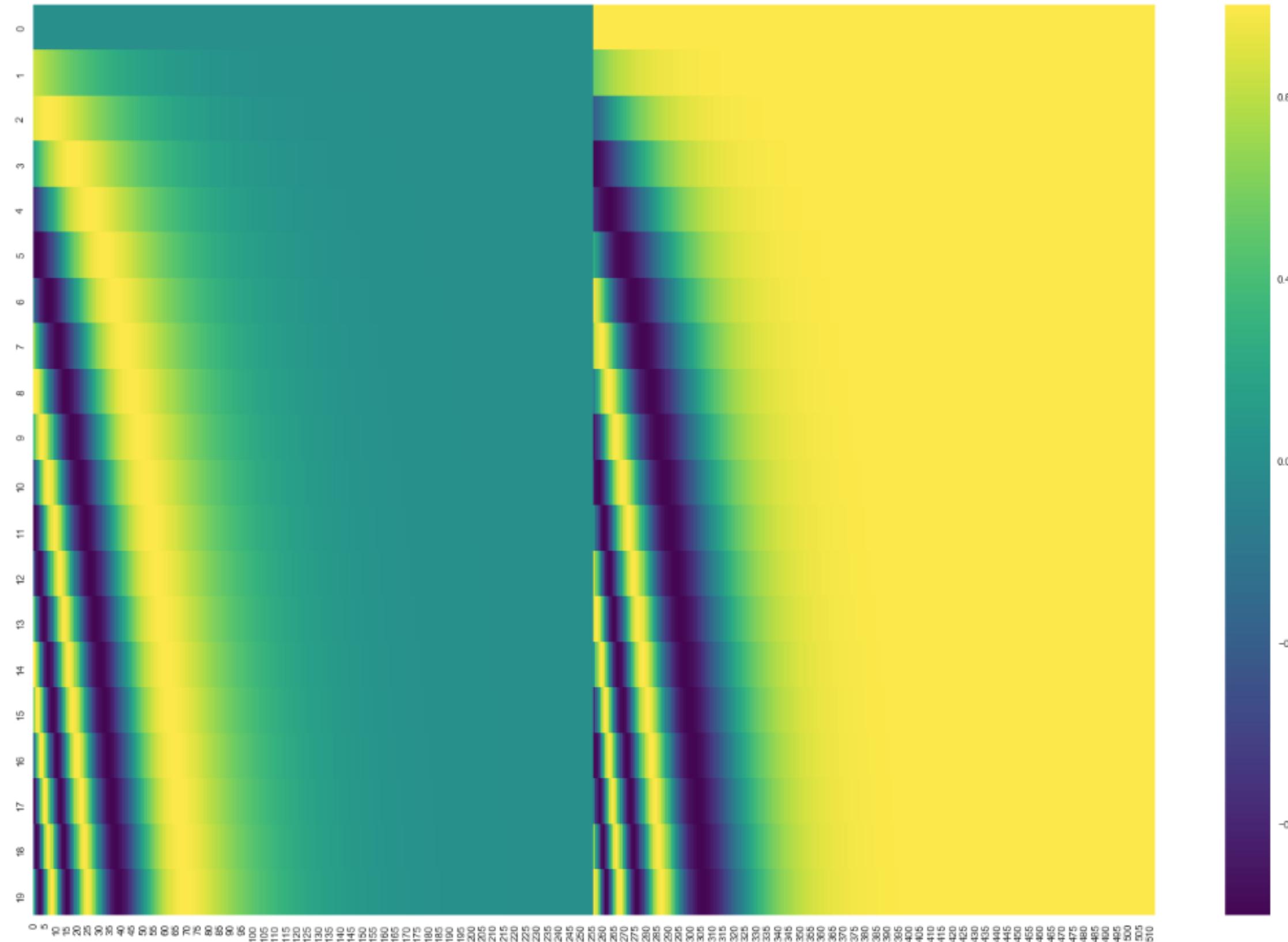
$$p_i = \begin{pmatrix} \sin(i/10000^{2*1/d}) \\ \cos(i/10000^{2*1/d}) \\ \vdots \\ \sin(i/10000^{2*\frac{d}{2}/d}) \\ \cos(i/10000^{2*\frac{d}{2}/d}) \end{pmatrix}$$



- Pros:
 - Periodicity indicates that maybe “absolute position” isn’t as important
 - Maybe can extrapolate to longer sequences as periods restart
- Cons:
 - Not learnable; also the extrapolation doesn’t really work

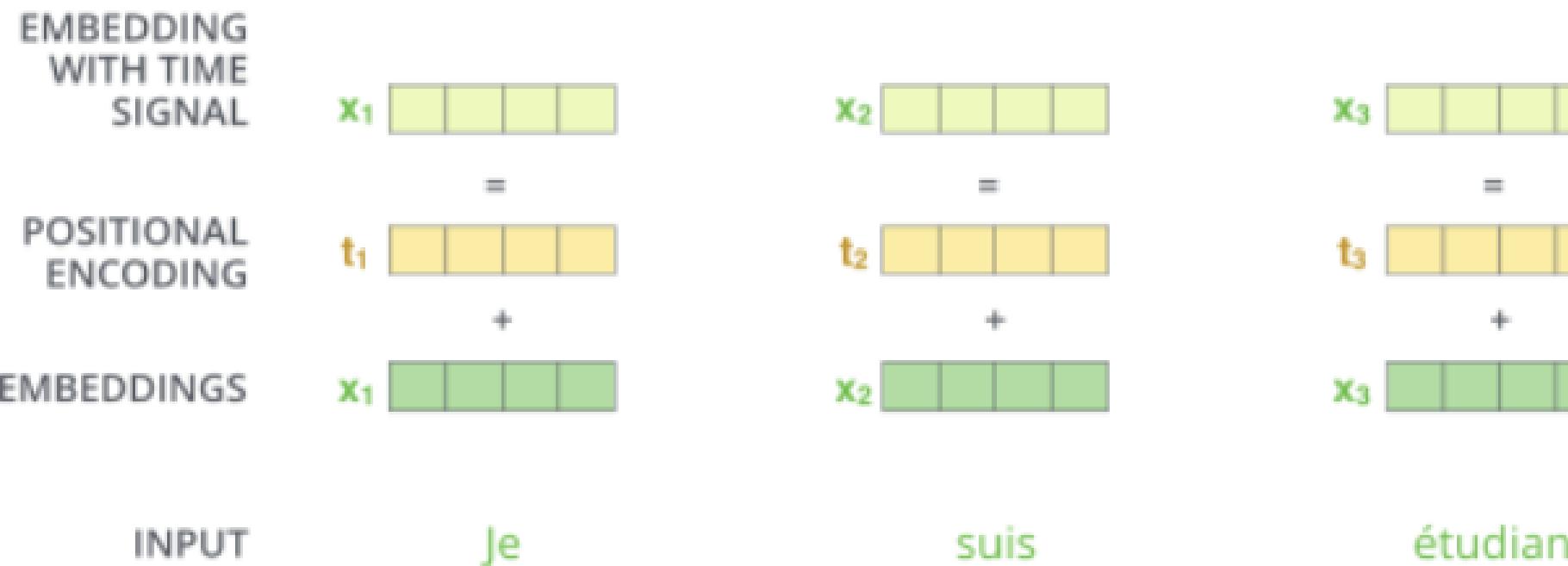
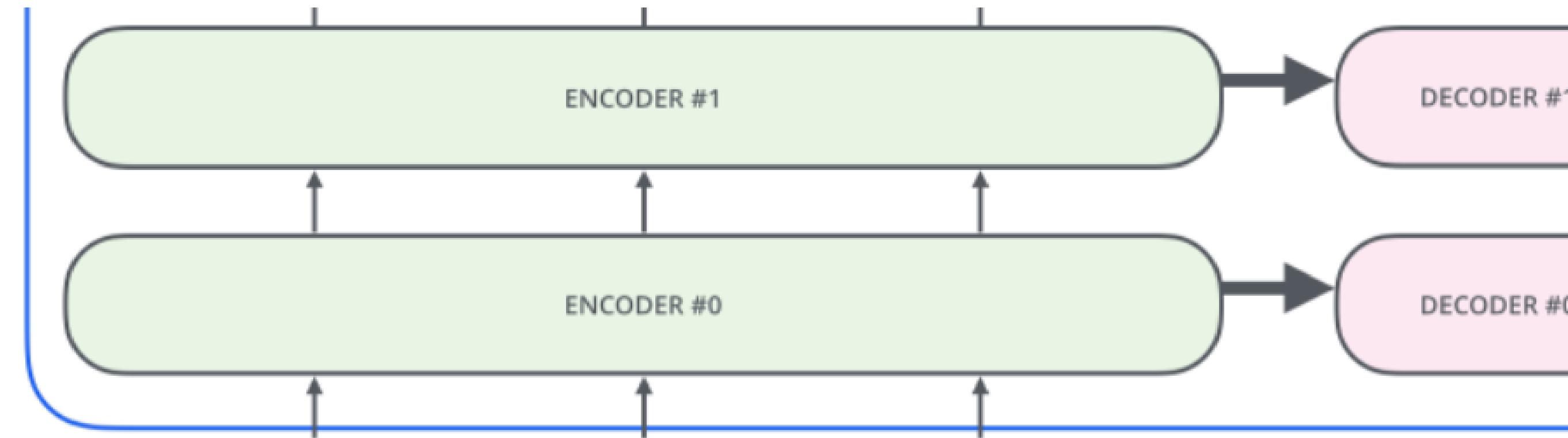
02 | Model Architecture

Position representation vectors through sinusoids



02 | Model Architecture

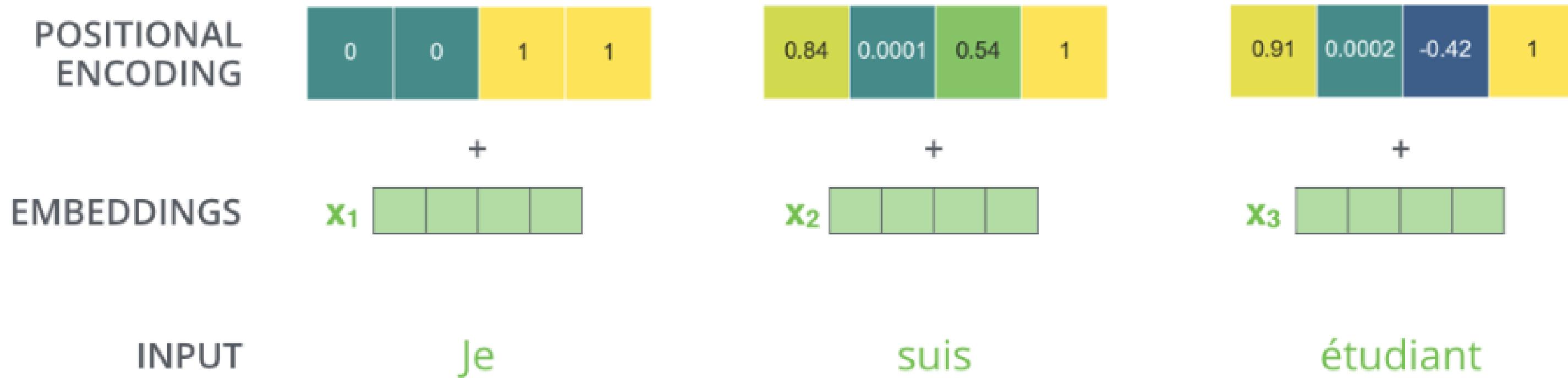
Position representation vectors through sinusoids



02 | Model Architecture

Position representation vectors through sinusoids

If we assumed the embedding has a dimensionality of 4, the actual positional encodings would look like this:



02 | Model Architecture

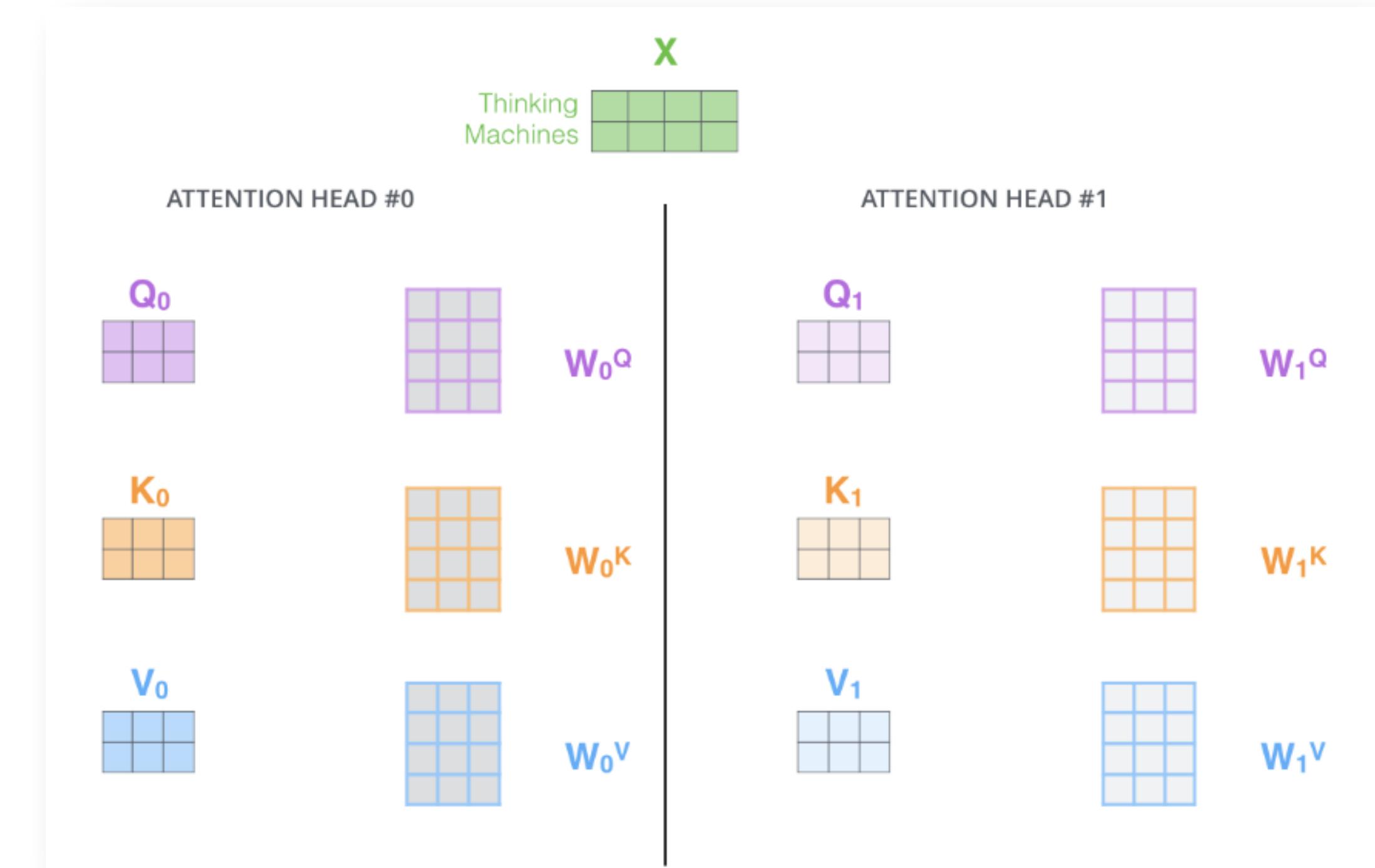
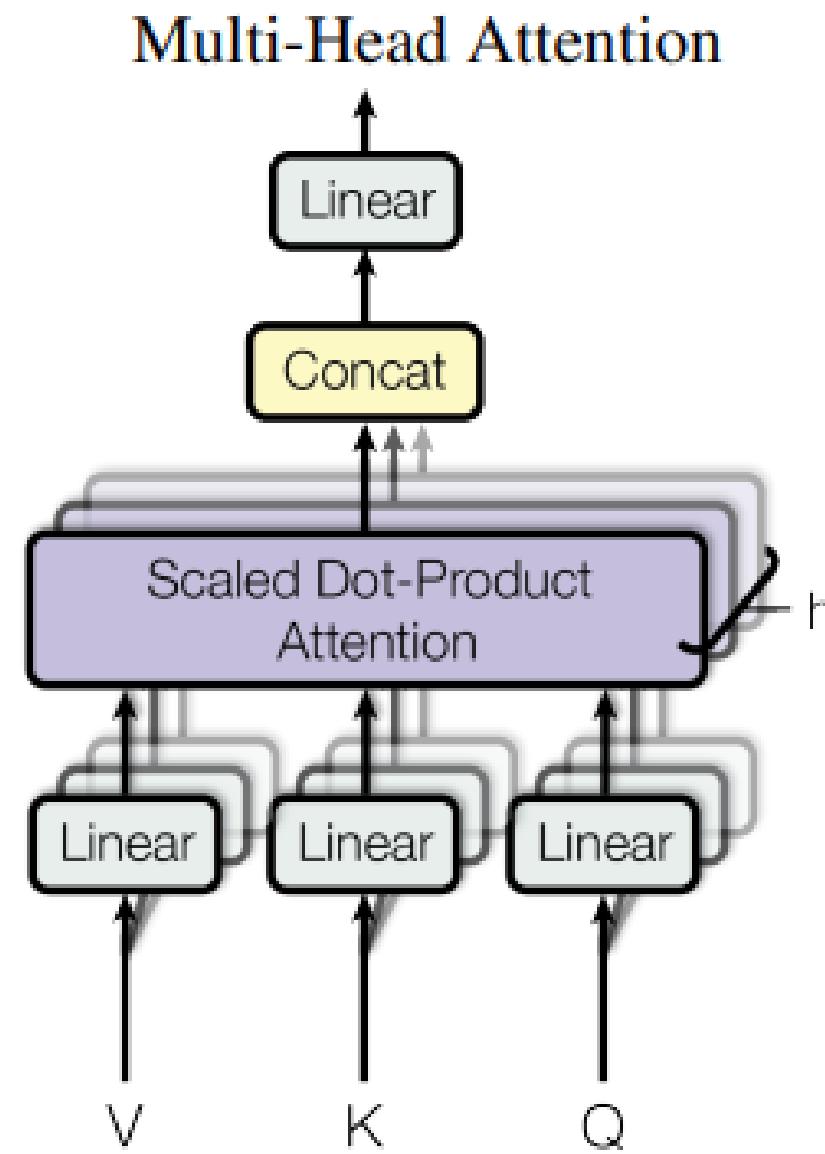
Position representation vectors learned from scratch

- **Learned absolute position representations:** Let all p_i be learnable parameters!
Learn a matrix $p \in \mathbb{R}^{d \times T}$, and let each p_i be a column of that matrix!
- Pros:
 - Flexibility: each position gets to be learned to fit the data
- Cons:
 - Definitely can't extrapolate to indices outside $1, \dots, T$.
- Most systems use this!
- Sometimes people try more flexible representations of position:
 - Relative linear position attention [\[Shaw et al., 2018\]](#)
 - Dependency syntax-based position [\[Wang et al., 2019\]](#)

02 | Model Architecture

Multi-Headed Self-Attention: k heads are better than 1

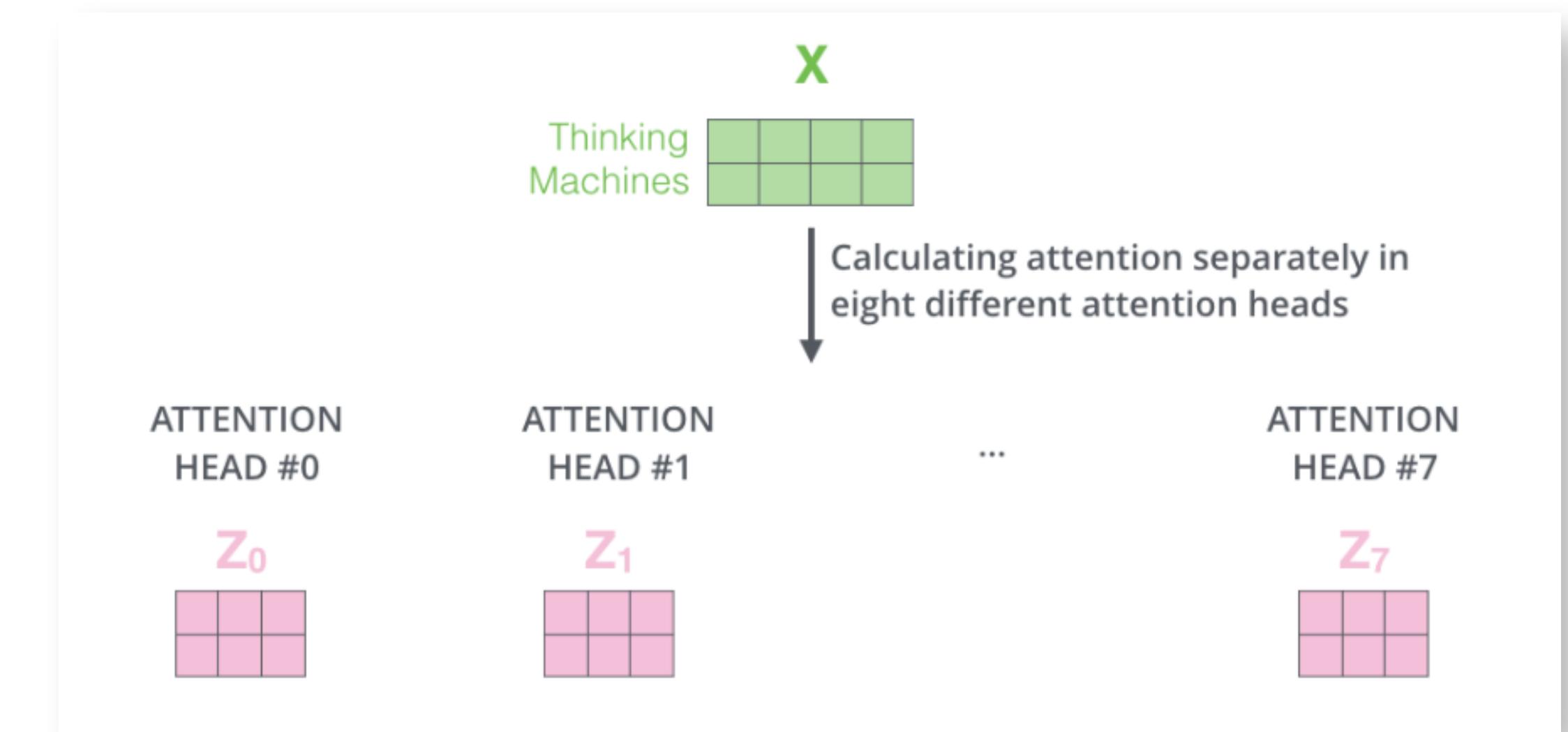
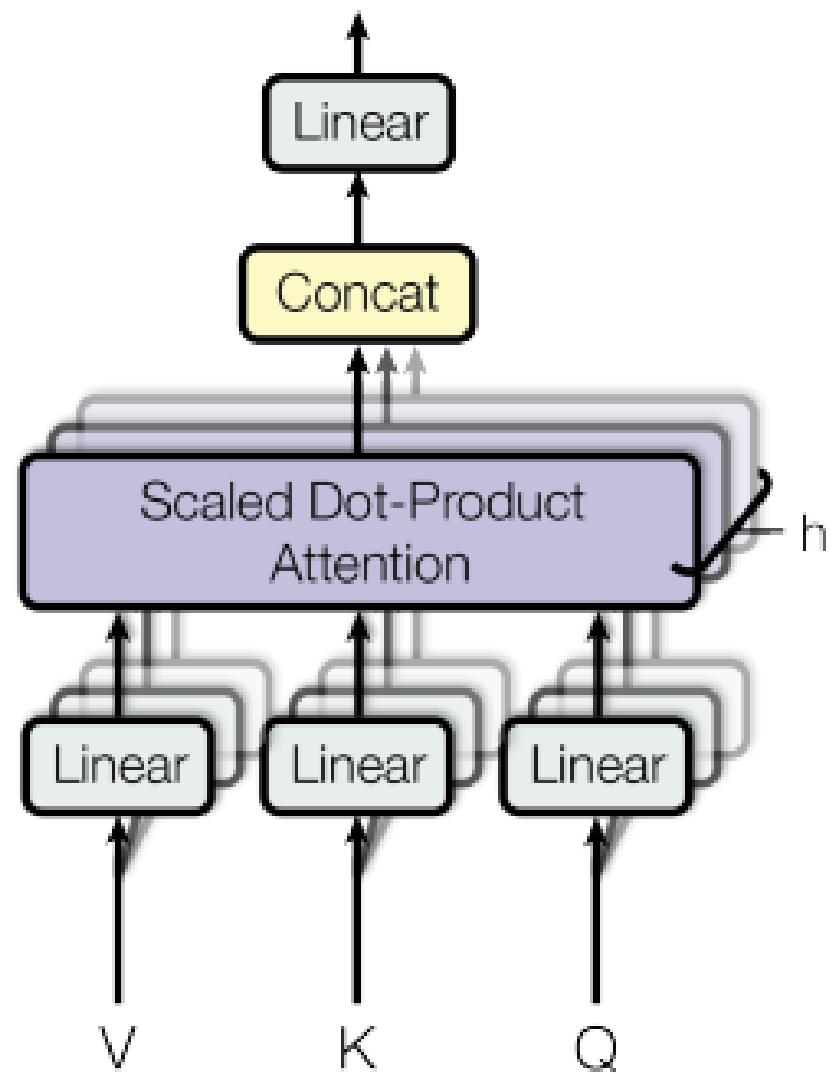
- High-Level Idea: Let's perform self-attention multiple times in parallel and combine the results.



02 | Model Architecture

Multi-Headed Self-Attention: k heads are better than 1

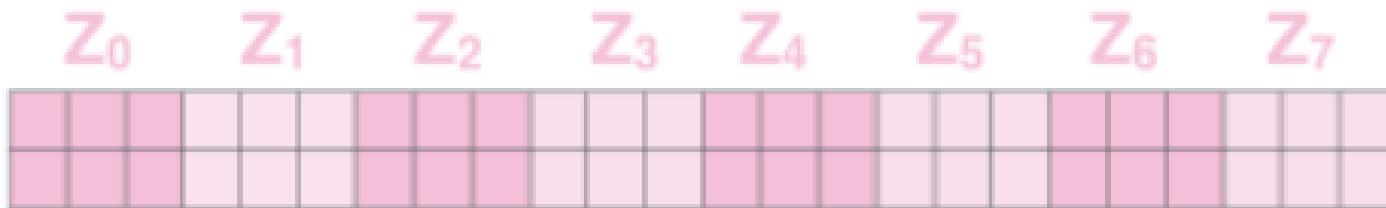
Multi-Head Attention



02 | Model Architecture

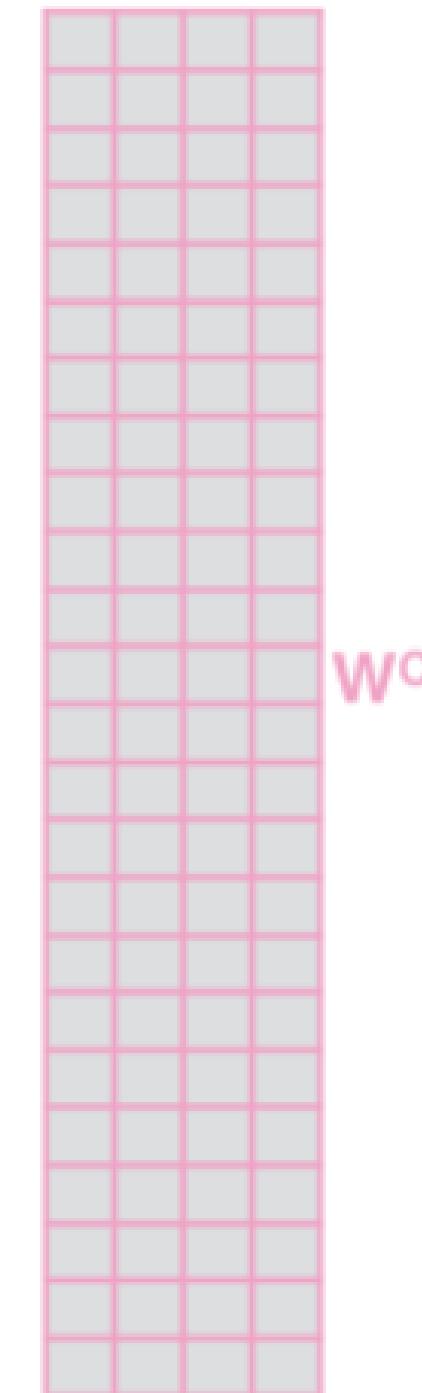
Multi-Headed Self-Attention: k heads are better than 1

1) Concatenate all the attention heads



2) Multiply with a weight matrix W^o that was trained jointly with the model

\times



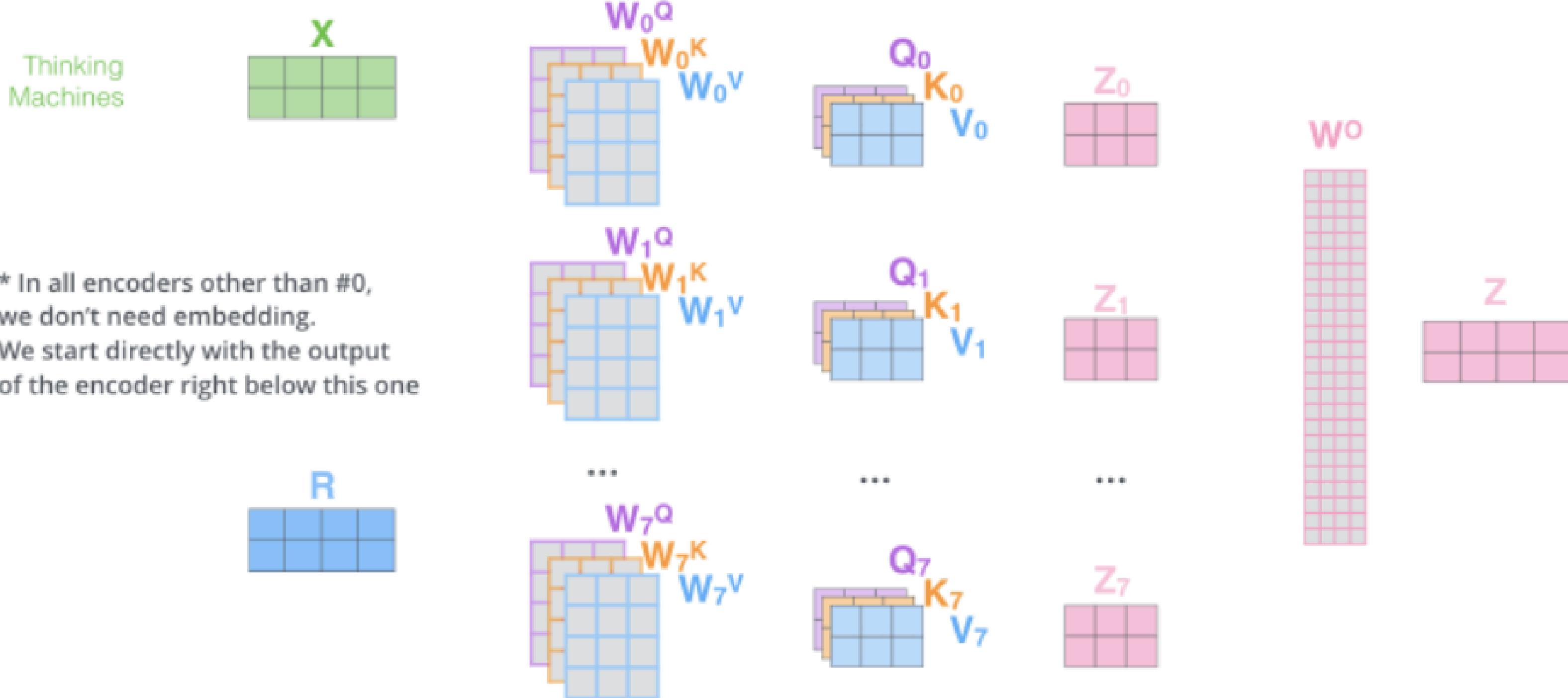
3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN

$$= \begin{matrix} Z \\ \hline \text{---} \\ \begin{matrix} \text{---} & \text{---} & \text{---} & \text{---} \end{matrix} \end{matrix}$$

02 | Model Architecture

Multi-Headed Self-Attention: k heads are better than 1

- 1) This is our input sentence* X
- 2) We embed each word* R
- 3) Split into 8 heads. We multiply X or R with weight matrices W_0^Q, W_0^K, W_0^V
- 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

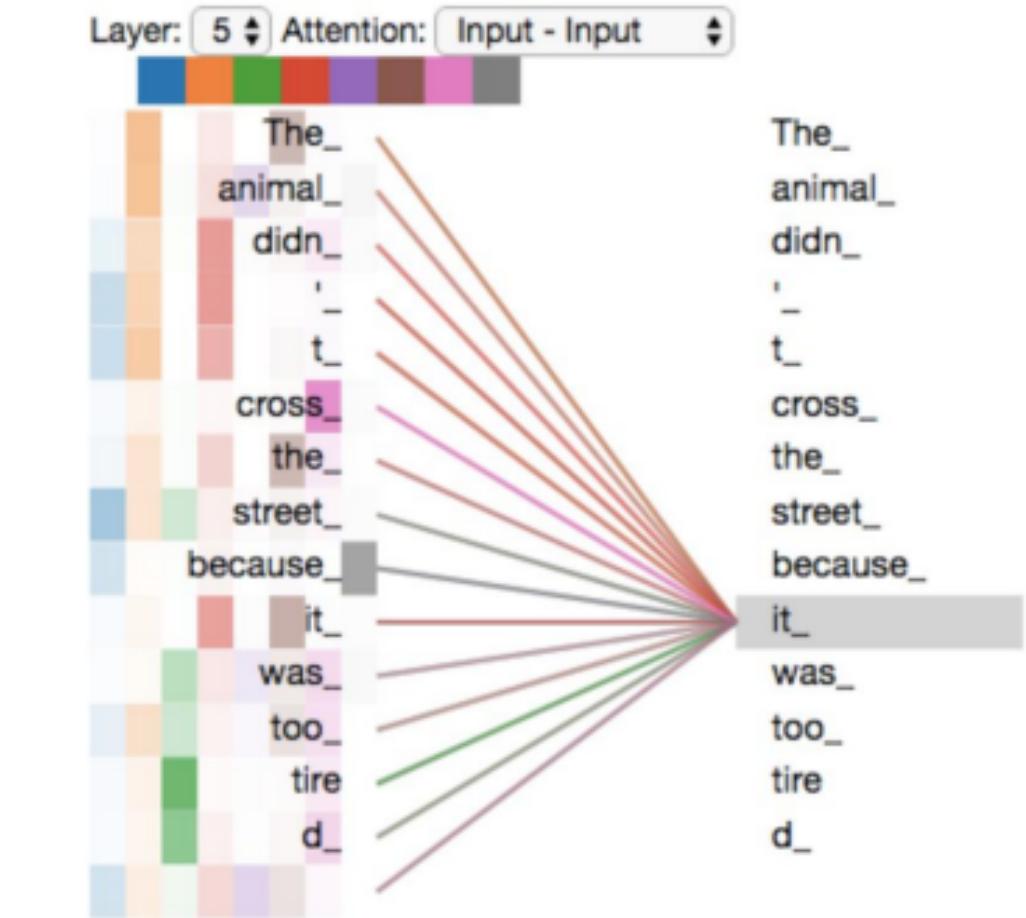


02 | Model Architecture

The Transformer Encoder: Multi-headed Self-Attention

- What if we want to look in multiple places in the sentence at once?
 - For word i , self-attention “looks” where $x_i^T Q^T K x_j$ is high, but maybe we want to focus on different j for different reasons?
- We’ll define **multiple attention “heads”** through multiple Q, K, V matrices
- Let, $Q_\ell, K_\ell, V_\ell \in \mathbb{R}^{d \times \frac{d}{h}}$, where h is the number of attention heads, and ℓ ranges from 1 to h .
- Each attention head performs attention independently:
 - $\text{output}_\ell = \text{softmax}(X Q_\ell K_\ell^T X^T) * X V_\ell$, where $\text{output}_\ell \in \mathbb{R}^{d/h}$
- Then the outputs of all the heads are combined!
 - $\text{output} = Y[\text{output}_1; \dots; \text{output}_h]$, where $Y \in \mathbb{R}^{d \times d}$

Visualizing 8 attentions



Credit to <https://jalammar.github.io/illustrated-transformer/>

02 | Model Architecture

Decoder: Masked Multi-Head Self-Attention

- Problem:

How do we keep the decoder from "cheating"?

If we have a language modeling objective,
can't the network just look ahead and "see" the answer?

- Solution:

Masked Multi-Head Attention.

At a high-level, we hide (mask) information about future tokens from the model.

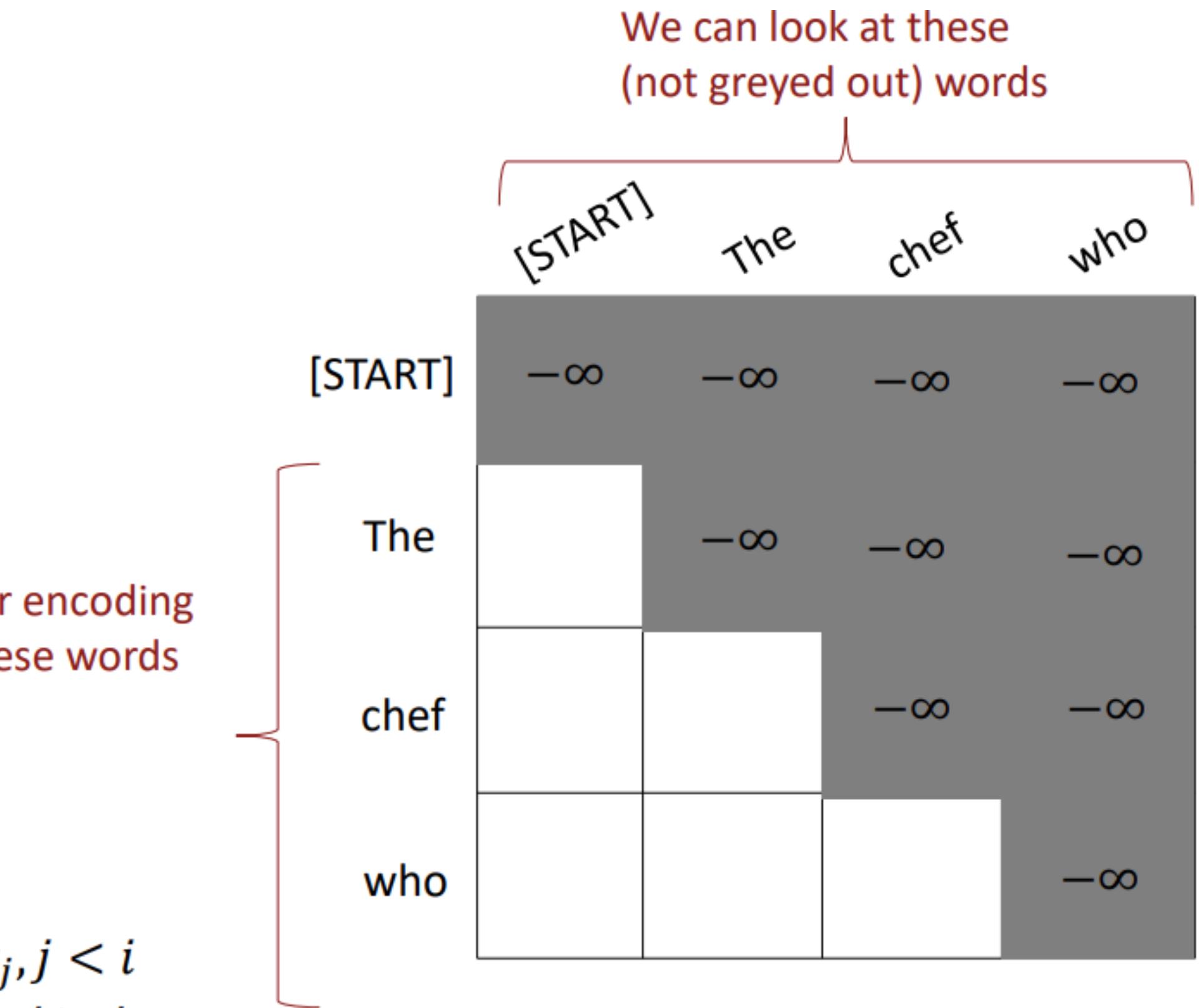
02

Model Architecture

Masking the future in self-attention

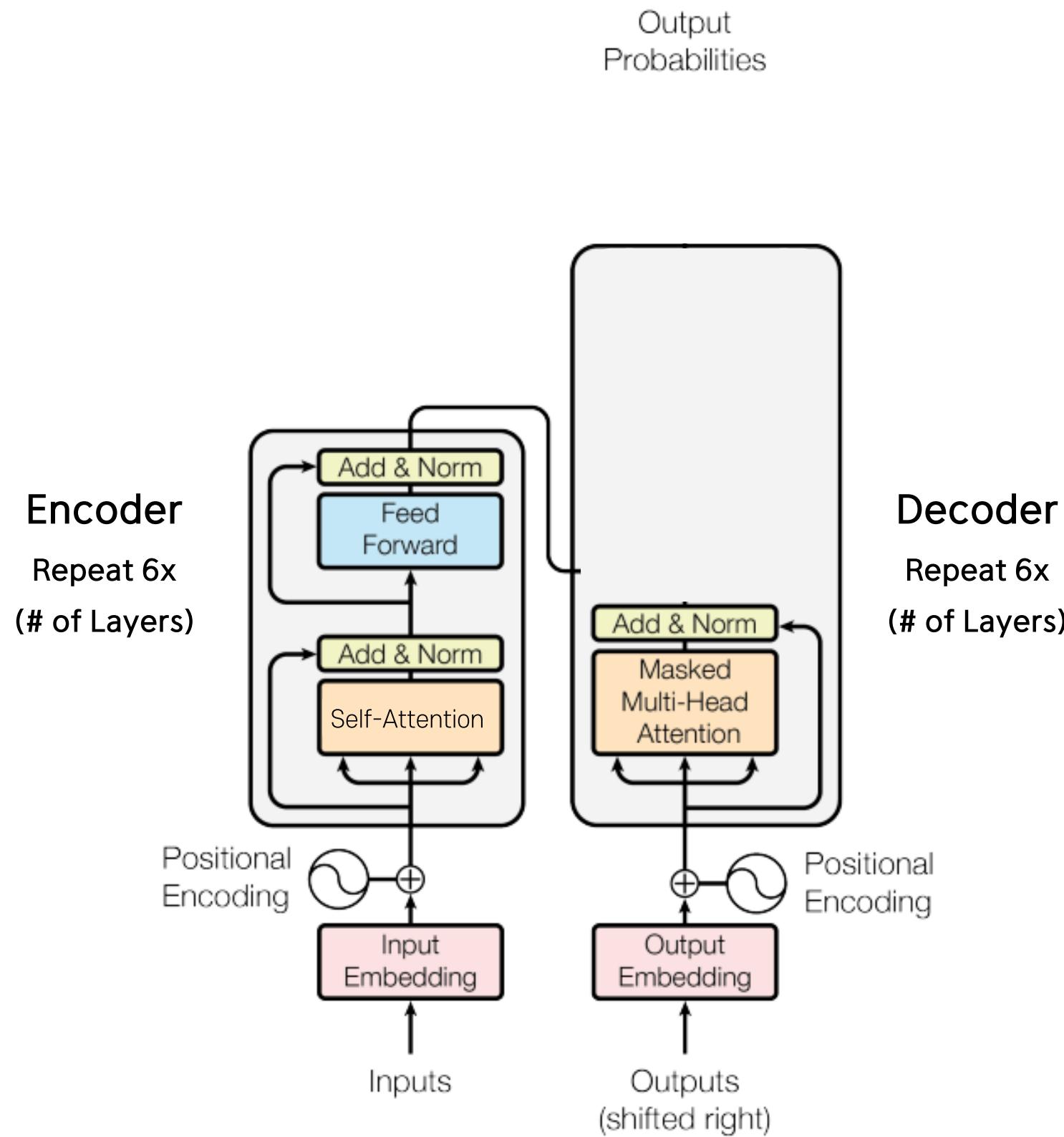
- To use self-attention in **decoders**, we need to ensure we can't peek at the future.
- At every timestep, we could change the set of **keys and queries** to include only past words. (Inefficient!)
- To enable parallelization, we **mask out attention** to future words by setting attention scores to $-\infty$.

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



02 | Model Architecture

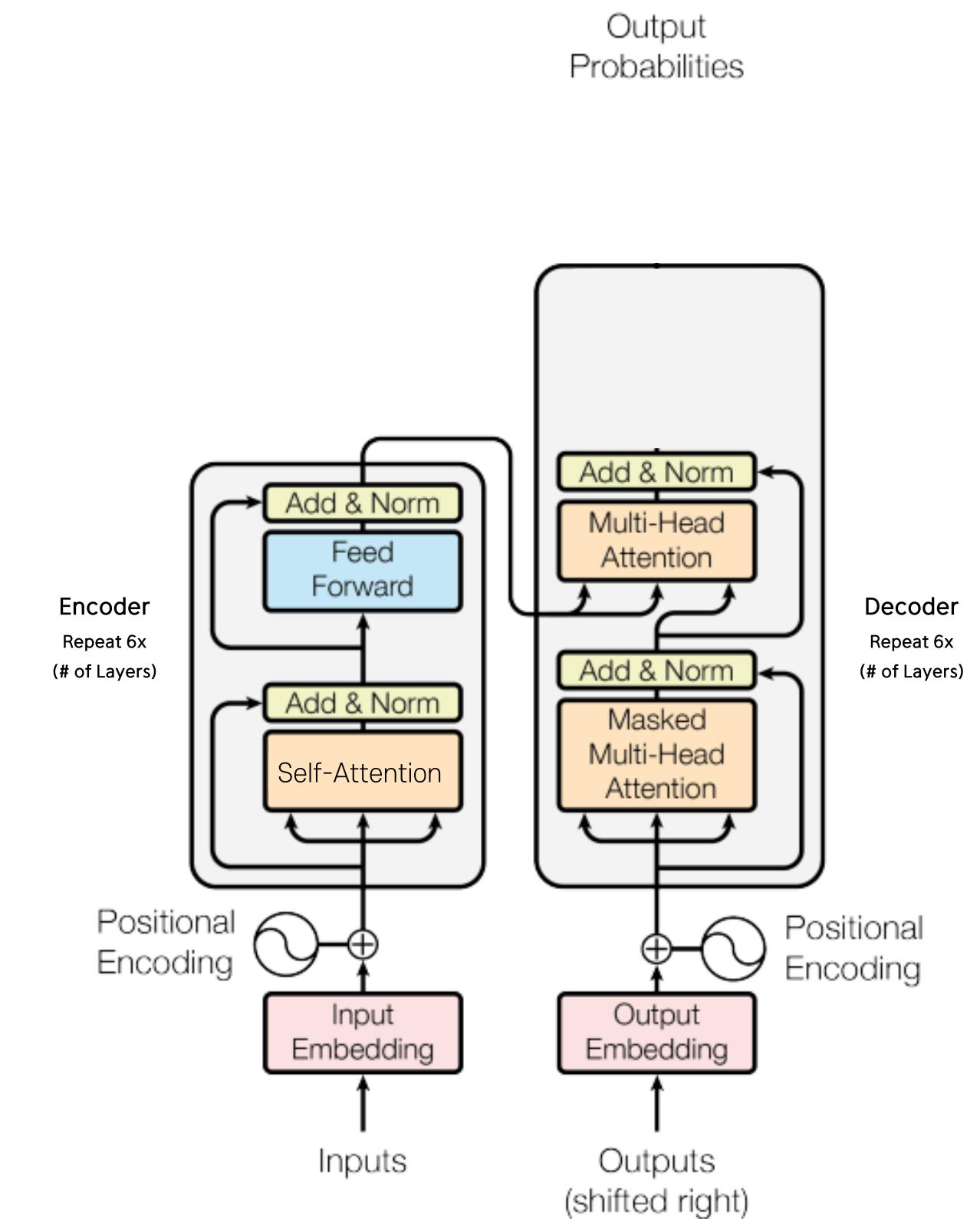
Decoder: Masked Multi-Headed Self-Attention



02 | Model Architecture

Decoder: Masked Multi-Headed Self-Attention

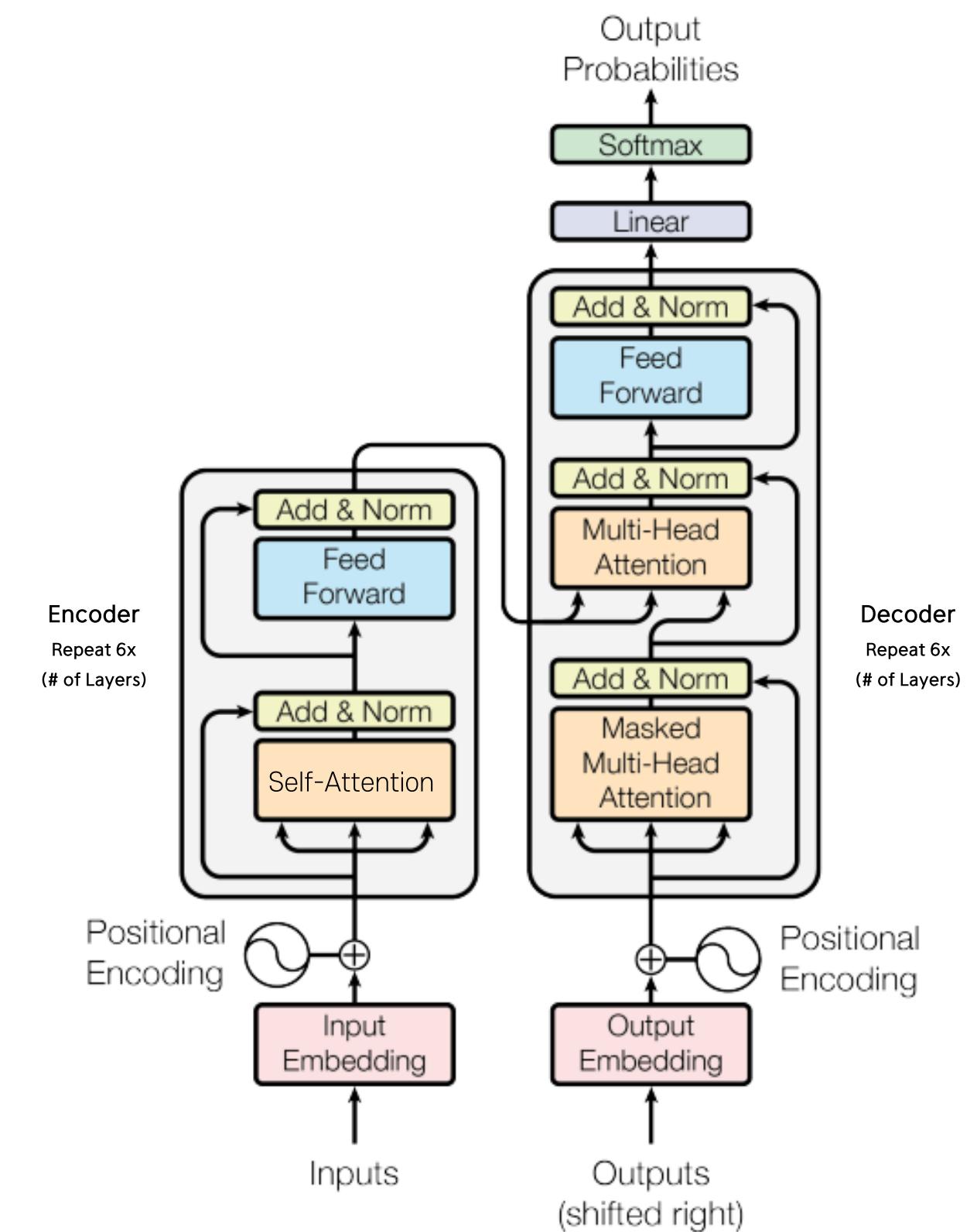
- We saw that self-attention is when keys, queries, and values come from the same source.
- In the decoder, we have attention that looks more like what we saw last week.
- Let h_1, \dots, h_T be output vectors from the Transformer **encoder**; $x_i \in \mathbb{R}^d$
- Let z_1, \dots, z_T be input vectors from the Transformer **decoder**, $z_i \in \mathbb{R}^d$
- Then keys and values are drawn from the **encoder** (like a memory):
 - $k_i = Kh_i, v_i = Vh_i$.
 - And the queries are drawn from the **decoder**, $q_i = Qz_i$.



02 | Model Architecture

Decoder: Finishing touches!

- Add a feed forward layer (with residual connections and layer norm)
- Add a final linear layer to project the embeddings into a much longer vector of length vocab size (logits)
- Add a final softmax to generate a probability distribution of possible next words!



02 | Model Architecture

What would we like to fix about the Transformer?

- **Quadratic compute in self-attention (today):**
 - Computing all pairs of interactions means our computation grows **quadratically** with the sequence length!
 - For recurrent models, it only grew linearly!
- **Position representations:**
 - Are simple absolute indices the best we can do to represent position?
 - Relative linear position attention [\[Shaw et al., 2018\]](#)
 - Dependency syntax-based position [\[Wang et al., 2019\]](#)

- WMT 2014 English-German dataset consisting of about 4.5 million sentence pairs.
 - Sentences were encoded using byte-pair encoding, which has a shared sourcetarget vocabulary of about 37000 tokens.
- WMT 2014 English-French dataset consisting of 36M sentences and split tokens into a 32000 word-piece vocabulary.
- Sentence pairs were batched together by approximate sequence length.
- Each training batch contained a set of sentence pairs containing approximately 25000 source tokens and 25000 target tokens.

- 8 NVIDIA P100 GPUs.
- For our base models using the hyperparameters described throughout the paper, each training step took about 0.4 seconds.
- We trained the base models for a total of 100,000 steps or 12 hours.
- For our big models,(described on the bottom line of table 3), step time was 1.0 seconds. The big models were trained for 300,000 steps (3.5 days).

03 | Training

Optimizer

- Adam optimizer
 - $\beta_1 = 0.9, \beta_2 = 0.98$ and $\epsilon = 10^{-9}$
- Learning rate:

$$lrate = d_{\text{model}}^{-0.5} \cdot \min(step_num^{-0.5}, step_num \cdot warmup_steps^{-1.5})$$

- Residual Dropout We apply dropout to the output of each sub-layer, before it is added to the sub-layer input and normalized.
- In addition, we apply dropout to the sums of the embeddings and the positional encodings in both the encoder and decoder stacks.
- For the base model, we use a rate of $P_{drop} = 0.1$.

03 | Training

Label Smoothing

- During training, we employed label smoothing of value $ls = 0.1$.

This hurts perplexity, as the model learns to be more unsure,
but improves accuracy and BLEU score.

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 [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		$2.3 \cdot 10^{19}$

04

Results

Model Variations

Table 3: Variations on the Transformer architecture. Unlisted values are identical to those of the base model. All metrics are on the English-to-German translation development set, newstest2013. Listed perplexities are per-wordpiece, according to our byte-pair encoding, and should not be compared to per-word perplexities.

	N	d_{model}	d_{ff}	h	d_k	d_v	P_{drop}	ϵ_{ls}	train steps	PPL (dev)	BLEU (dev)	params $\times 10^6$	
base	6	512	2048	8	64	64	0.1	0.1	100K	4.92	25.8	65	
(A)				1	512	512				5.29	24.9		
				4	128	128				5.00	25.5		
				16	32	32				4.91	25.8		
				32	16	16				5.01	25.4		
(B)					16					5.16	25.1	58	
					32					5.01	25.4	60	
(C)				2						6.11	23.7	36	
				4						5.19	25.3	50	
				8						4.88	25.5	80	
				256		32	32			5.75	24.5	28	
				1024		128	128			4.66	26.0	168	
				1024				5.12	25.4	53			
				4096				4.75	26.2	90			
(D)							0.0			5.77	24.6		
							0.2			4.95	25.5		
							0.0			4.67	25.3		
							0.2			5.47	25.7		
(E)	positional embedding instead of sinusoids									4.92	25.7		
big	6	1024	4096	16				0.3	300K	4.33	26.4	213	

04

Results

English Constituency Parsing

Table 4: The Transformer generalizes well to English constituency parsing (Results are on Section 23 of WSJ)

Parser	Training	WSJ 23 F1
Vinyals & Kaiser el al. (2014) [37]	WSJ only, discriminative	88.3
Petrov et al. (2006) [29]	WSJ only, discriminative	90.4
Zhu et al. (2013) [40]	WSJ only, discriminative	90.4
Dyer et al. (2016) [8]	WSJ only, discriminative	91.7
Transformer (4 layers)	WSJ only, discriminative	91.3
Zhu et al. (2013) [40]	semi-supervised	91.3
Huang & Harper (2009) [14]	semi-supervised	91.3
McClosky et al. (2006) [26]	semi-supervised	92.1
Vinyals & Kaiser el al. (2014) [37]	semi-supervised	92.1
Transformer (4 layers)	semi-supervised	92.7
Luong et al. (2015) [23]	multi-task	93.0
Dyer et al. (2016) [8]	generative	93.3

05 | Conclusion

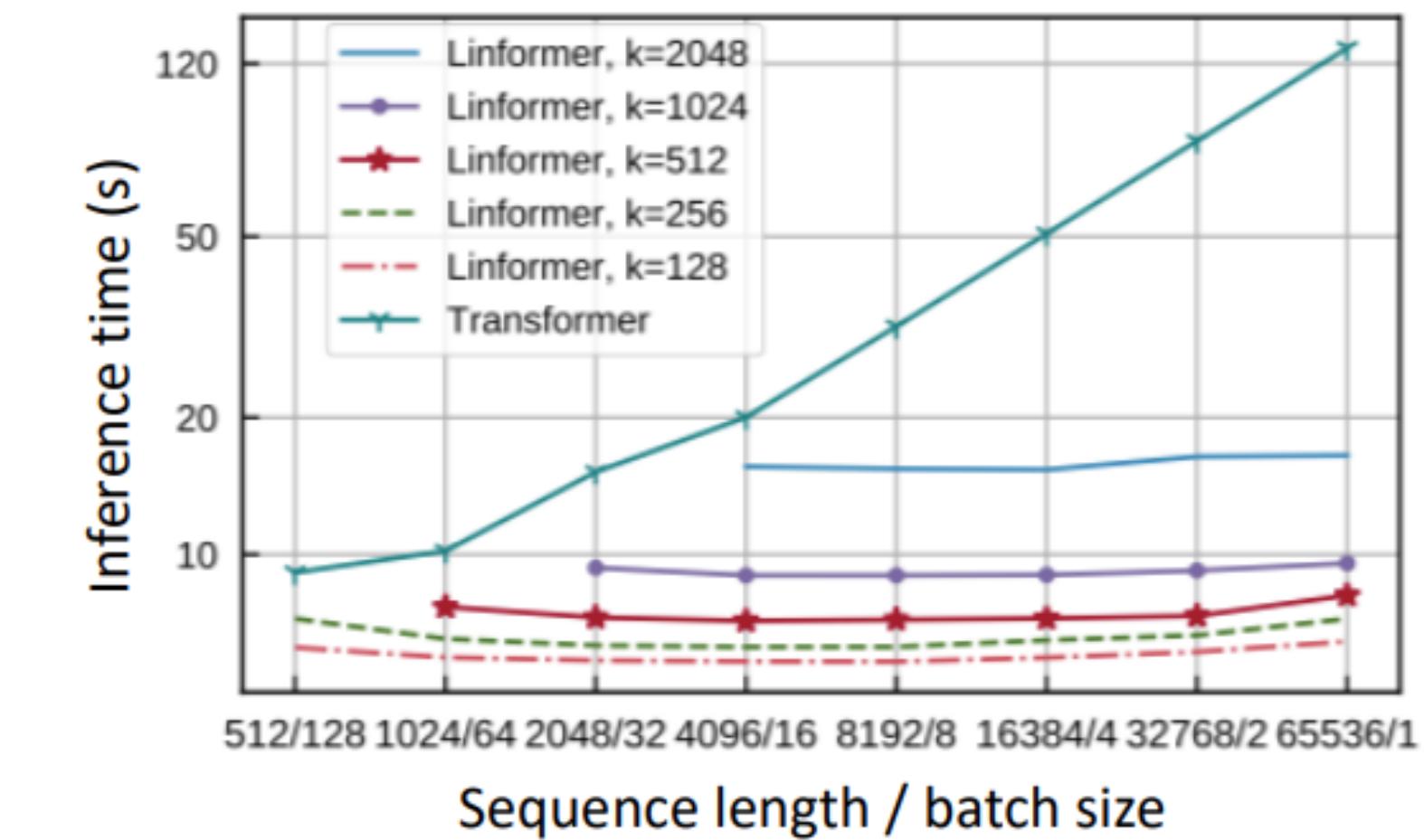
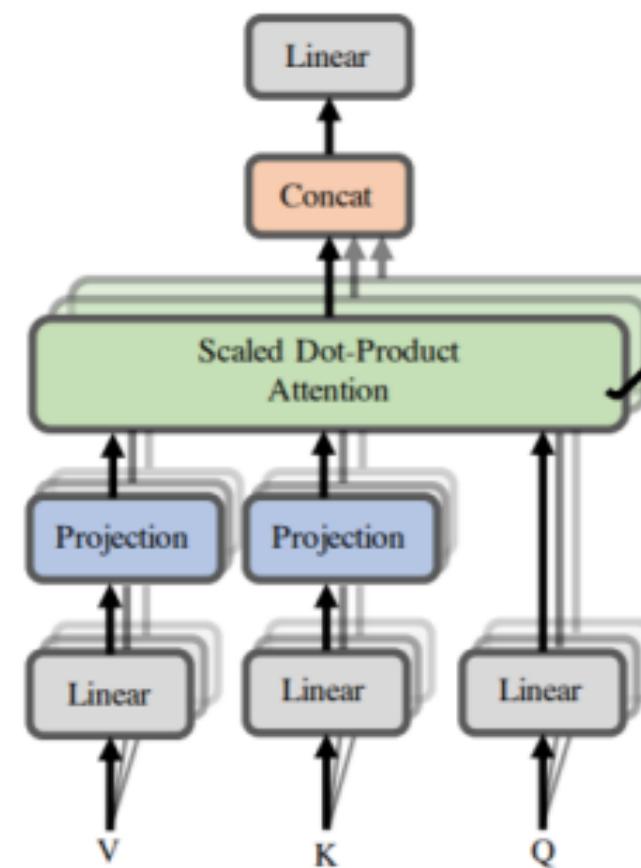
- The first sequence transduction model based entirely on attention, replacing the recurrent layers most commonly used in encoder-decoder architectures with multi-headed self-attention.
- For translation tasks, the Transformer can be trained significantly faster than architectures based on recurrent or convolutional layers.

05 | Conclusion

Recent work on improving on quadratic self-attention cost

- Considerable recent work has gone into the question, *Can we build models like Transformers without paying the $O(T^2)$ all-pairs self-attention cost?*
- For example, **Linformer** [[Wang et al., 2020](#)]

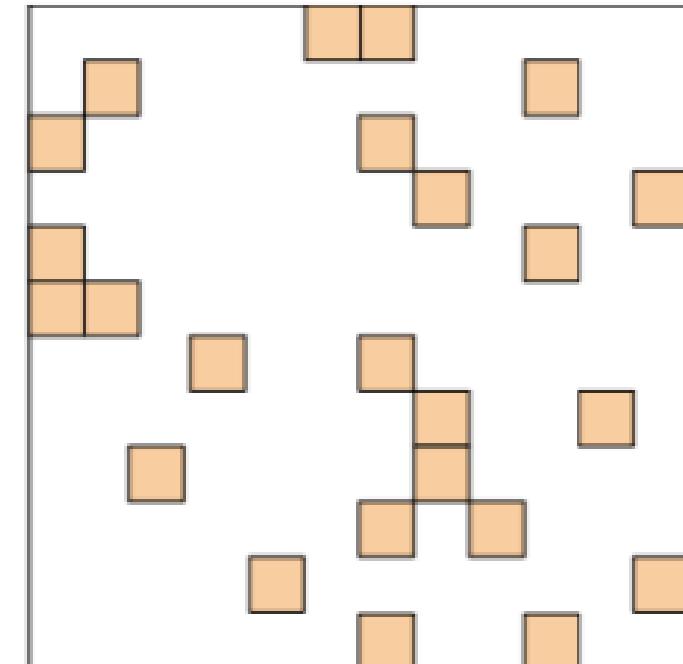
Key idea: map the sequence length dimension to a lower-dimensional space for values, keys



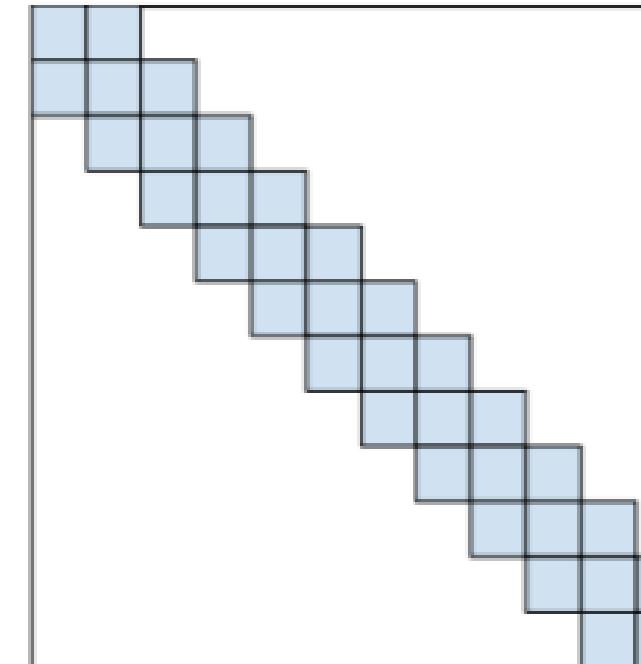
Recent work on improving on quadratic self-attention cost

- Considerable recent work has gone into the question, *Can we build models like Transformers without paying the $O(T^2)$ all-pairs self-attention cost?*
- For example, **BigBird** [[Zaheer et al., 2021](#)]

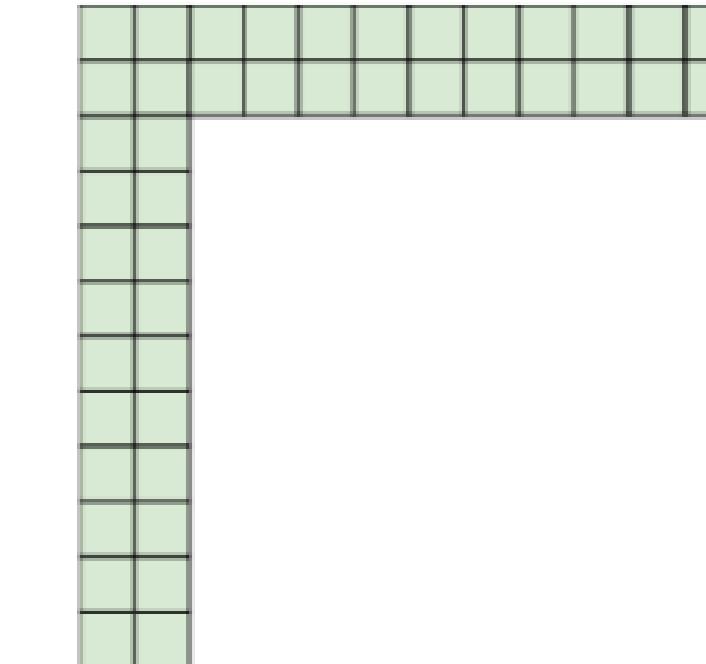
Key idea: replace all-pairs interactions with a family of other interactions, **like local windows, looking at everything, and random interactions.**



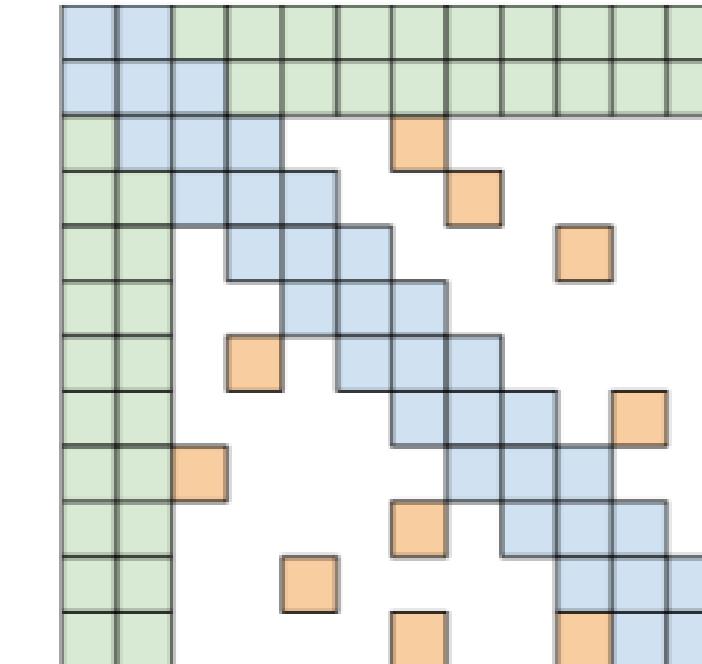
(a) Random attention



(b) Window attention



(c) Global Attention



(d) BIGBIRD

05 | Conclusion

Do Transformer Modifications Transfer?

- "Surprisingly, we find that most modifications do not meaningfully improve performance."

Model	Params	Ops	Step/s	Early loss	Final loss	RoGLOUE	RoScore	WikiQ	WMT EnDe
Vanilla Transformer	220M	11.17 [*]	0.30	2.182 ± 0.005	1.818	71.00	17.78	23.02	26.62
Gelu	220M	11.17 [*]	0.38	2.179 ± 0.002	1.818	71.78	17.88	23.12	26.47
Sinh	220M	11.17 [*]	0.62	2.180 ± 0.003	1.817	71.77	17.74	23.04	26.75
ELU	220M	11.17 [*]	0.56	2.270 ± 0.007	1.812	87.63	16.73	23.02	26.38
GLU	220M	11.17 [*]	0.39	2.174 ± 0.003	1.814	74.29	17.42	24.54	27.12
GeGLU	220M	11.17 [*]	0.55	2.180 ± 0.006	1.799	71.44	18.37	24.87	26.47
ReGLU	220M	11.17 [*]	0.57	2.145 ± 0.004	1.808	71.17	18.36	24.87	27.02
SeLU	220M	11.17 [*]	0.55	2.235 ± 0.004	1.818	87.70	22.75	23.09	26.39
SwiGLU	220M	11.17 [*]	0.53	2.127 ± 0.003	1.799	76.09	18.29	24.38	27.02
LoGLU	220M	11.17 [*]	0.39	2.149 ± 0.005	1.798	71.34	17.97	24.34	26.53
Sigmoid	220M	11.17 [*]	0.63	2.281 ± 0.019	1.817	74.31	17.51	23.02	26.30
Softplus	220M	11.17 [*]	0.47	2.297 ± 0.011	1.805	73.45	17.65	24.54	26.89
RMS Norm	220M	11.17 [*]	0.68	2.187 ± 0.006	1.821	75.45	17.84	24.97	27.14
ResNet	220M	11.17 [*]	0.51	2.262 ± 0.003	1.819	81.09	15.94	26.00	26.37
ResNet + LayerNorm	220M	11.17 [*]	0.26	2.221 ± 0.006	1.814	70.42	17.58	23.02	26.29
ResNet + RMS Norm	220M	11.17 [*]	0.54	2.221 ± 0.009	1.815	70.33	17.32	23.02	26.19
Pump	220M	11.17 [*]	0.95	2.380 ± 0.012	1.807	88.56	14.47	23.02	26.31
24 layers, $d_h = 1536$, $H = 6$	224M	11.17 [*]	0.39	2.280 ± 0.007	1.813	74.89	17.75	25.18	26.89
16 layers, $d_h = 3072$, $H = 8$	224M	11.17 [*]	0.38	2.185 ± 0.005	1.821	76.45	18.83	24.54	27.10
8 layers, $d_h = 6144$, $H = 16$	224M	11.17 [*]	0.69	2.180 ± 0.005	1.817	74.58	17.09	23.08	26.86
6 layers, $d_h = 6144$, $H = 24$	224M	11.17 [*]	0.79	2.201 ± 0.010	1.817	73.35	17.59	24.80	26.89
Block sharing	61.6M	11.17 [*]	0.91	2.497 ± 0.027	1.814	84.50	14.52	25.94	27.45
+ Factorized embeddings	61.6M	9.47 [*]	4.21	2.621 ± 0.005	1.803	80.41	14.09	19.64	25.27
+ Factorized & shared embeddings	61.6M	9.17 [*]	4.27	2.397 ± 0.013	1.805	83.80	11.37	19.64	25.19
Encoder only block sharing	176M	11.17 [*]	0.68	2.289 ± 0.023	1.829	85.60	16.23	23.02	26.23
Decoder only block sharing	144M	11.17 [*]	0.79	2.352 ± 0.029	1.802	87.93	16.13	23.02	26.08
Factorized Embedding	227M	9.47 [*]	0.89	2.289 ± 0.006	1.813	70.41	15.82	22.75	26.50
Factorized & shared embeddings	227M	9.17 [*]	0.92	2.220 ± 0.010	1.802	86.49	16.33	22.75	26.44
Text encoder/decoder input embeddings	248M	11.17 [*]	0.55	2.182 ± 0.002	1.840	71.79	17.72	24.34	26.49
Text decoder input and output embeddings	248M	11.17 [*]	0.37	2.187 ± 0.007	1.827	74.46	17.74	24.87	26.47
Unified embeddings	272M	11.17 [*]	0.53	2.180 ± 0.005	1.824	72.99	17.58	23.28	26.49
Adaptive input embeddings	291M	9.27 [*]	0.55	2.250 ± 0.002	1.809	86.37	16.21	24.97	26.88
Adaptive softmax	291M	9.27 [*]	0.69	2.384 ± 0.005	1.802	73.81	16.87	25.18	26.56
Adaptive softmax without projection	224M	9.47 [*]	0.63	2.289 ± 0.009	1.811	71.43	17.18	23.02	26.73
Mixture of softmaxes	224M	10.37 [*]	2.36	2.297 ± 0.017	1.821	76.77	17.62	22.75	26.82
Transparent attention	223M	11.17 [*]	0.33	2.181 ± 0.014	1.811	84.31	18.49	25.18	26.89
Dynamic convolution	297M	11.17 [*]	0.65	2.603 ± 0.009	1.807	88.30	12.67	25.18	27.03
Lightweight convolution	294M	10.47 [*]	4.07	2.187 ± 0.010	1.809	81.07	14.96	23.02	24.73
Evolved Transformer	227M	9.37 [*]	0.69	2.220 ± 0.003	1.803	73.67	18.76	24.97	26.56
Synthesizer (dense)	224M	11.17 [*]	0.47	2.354 ± 0.025	1.802	81.03	14.27	18.14	26.63
Synthesizer (dense plus)	243M	11.17 [*]	0.37	2.181 ± 0.019	1.813	78.59	18.98	24.84	26.71
Synthesizer (dense plus alpha)	243M	11.17 [*]	0.61	2.180 ± 0.007	1.828	74.33	17.03	23.08	26.41
Synthesizer (factorized)	297M	10.37 [*]	0.66	2.381 ± 0.017	1.808	83.79	15.39	23.05	26.42
Synthesizer (random)	254M	10.17 [*]	0.66	2.326 ± 0.012	1.809	84.27	18.35	19.54	26.44
Synthesizer (random plus)	292M	11.17 [*]	0.63	2.189 ± 0.004	1.812	78.32	17.04	24.87	26.43
Synthesizer (random plus alpha)	293M	11.17 [*]	0.47	2.186 ± 0.007	1.828	75.34	17.08	24.06	26.39
Universal Transformer	64M	40.37 [*]	0.88	2.480 ± 0.026	1.810	70.13	14.09	19.05	25.91
Mixture of experts	648M	11.17 [*]	0.30	2.148 ± 0.006	1.798	74.15	18.13	24.08	26.84
Switch Transformer	1100M	11.17 [*]	0.18	2.133 ± 0.007	1.758	75.38	18.02	24.19	26.83
Panel Transformer	223M	1.37 [*]	0.30	2.289 ± 0.008	1.818	87.30	16.28	23.75	26.20
Weighted Transformer	290M	71.17 [*]	0.39	2.278 ± 0.023	1.803	89.34	16.98	23.02	26.30
Product key memory	621M	386.67 [*]	0.25	2.153 ± 0.003	1.798	75.16	17.01	23.05	26.73

Do Transformer Modifications Transfer Across Implementations and Applications?

Sharan Narang^{*} Hyung Won Chung Yi Tay William Fedus

Thibault Fevry[†] Michael Matena[†] Karishma Malkan[†] Noah Fiedel

Noam Shazeer Zhenzhong Lan[†] Yanqi Zhou Wei Li

Nan Ding Jake Marcus Adam Roberts Colin Raffel[†]

References

1. <http://web.stanford.edu/class/cs224n/slides/cs224n-2022-lecture07-nmt.pdf>
2. <http://web.stanford.edu/class/cs224n/slides/cs224n-2022-lecture09-transformers.pdf>
3. <https://jalammar.github.io/illustrated-transformer/>
4. <https://ai.googleblog.com/2017/08/transformer-novel-neural-network.html>
5. <https://arxiv.org/pdf/1706.03762.pdf>

THANK YOU

Q & A