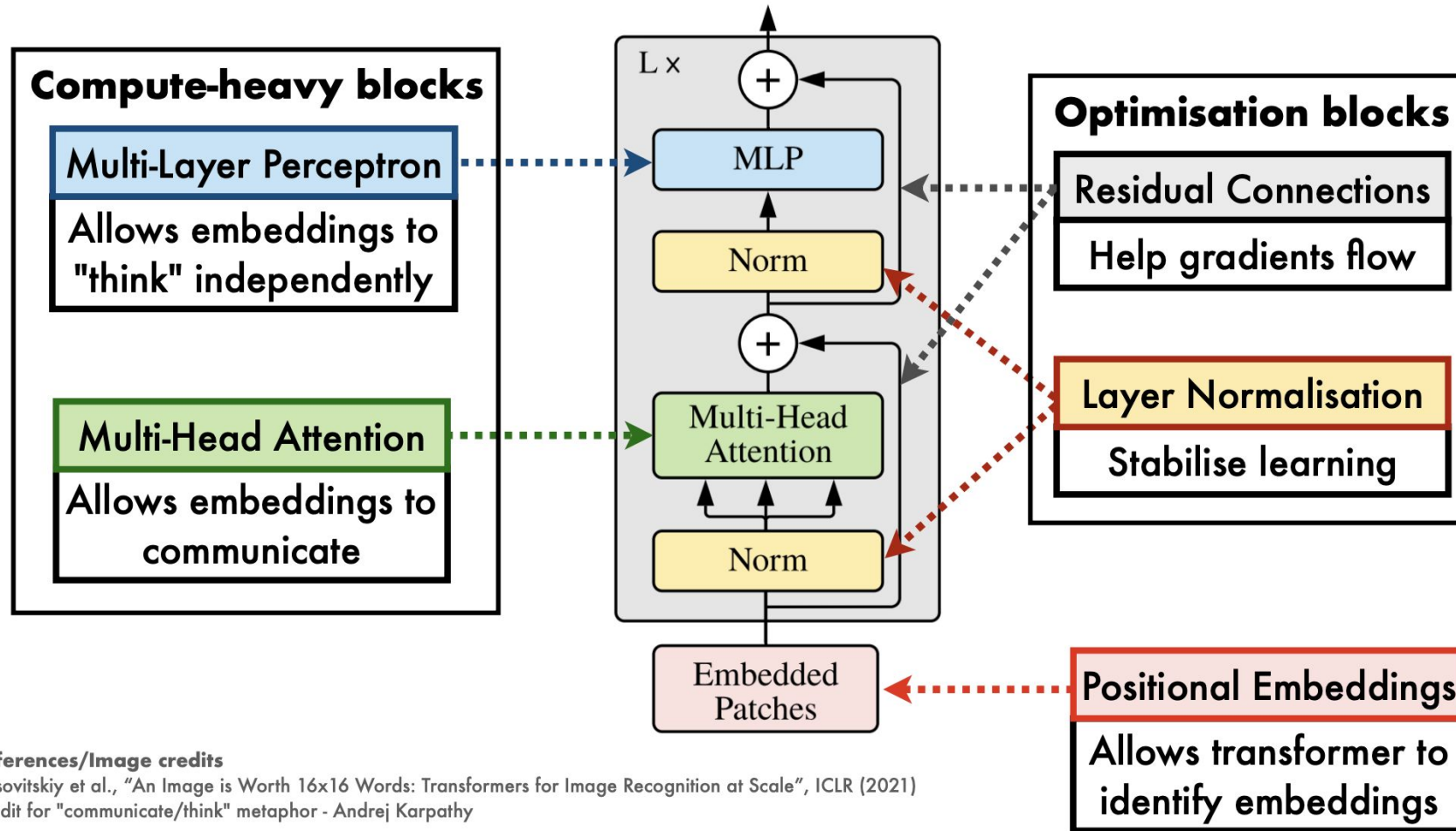


Transformers

Transformers

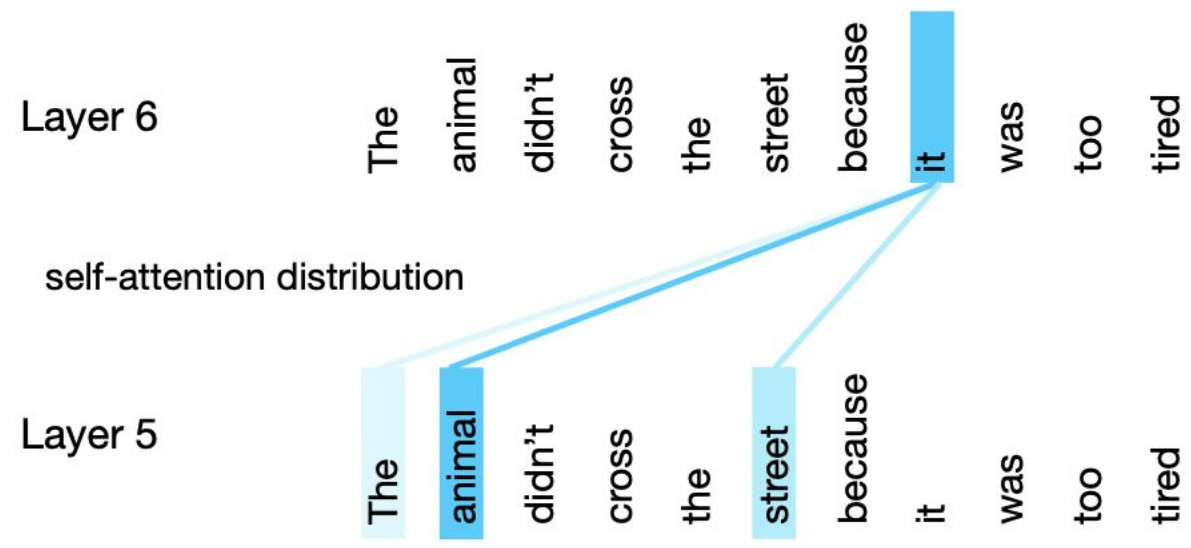
Transformer Encoder

Five key ideas



Blank Slide

Self Attention

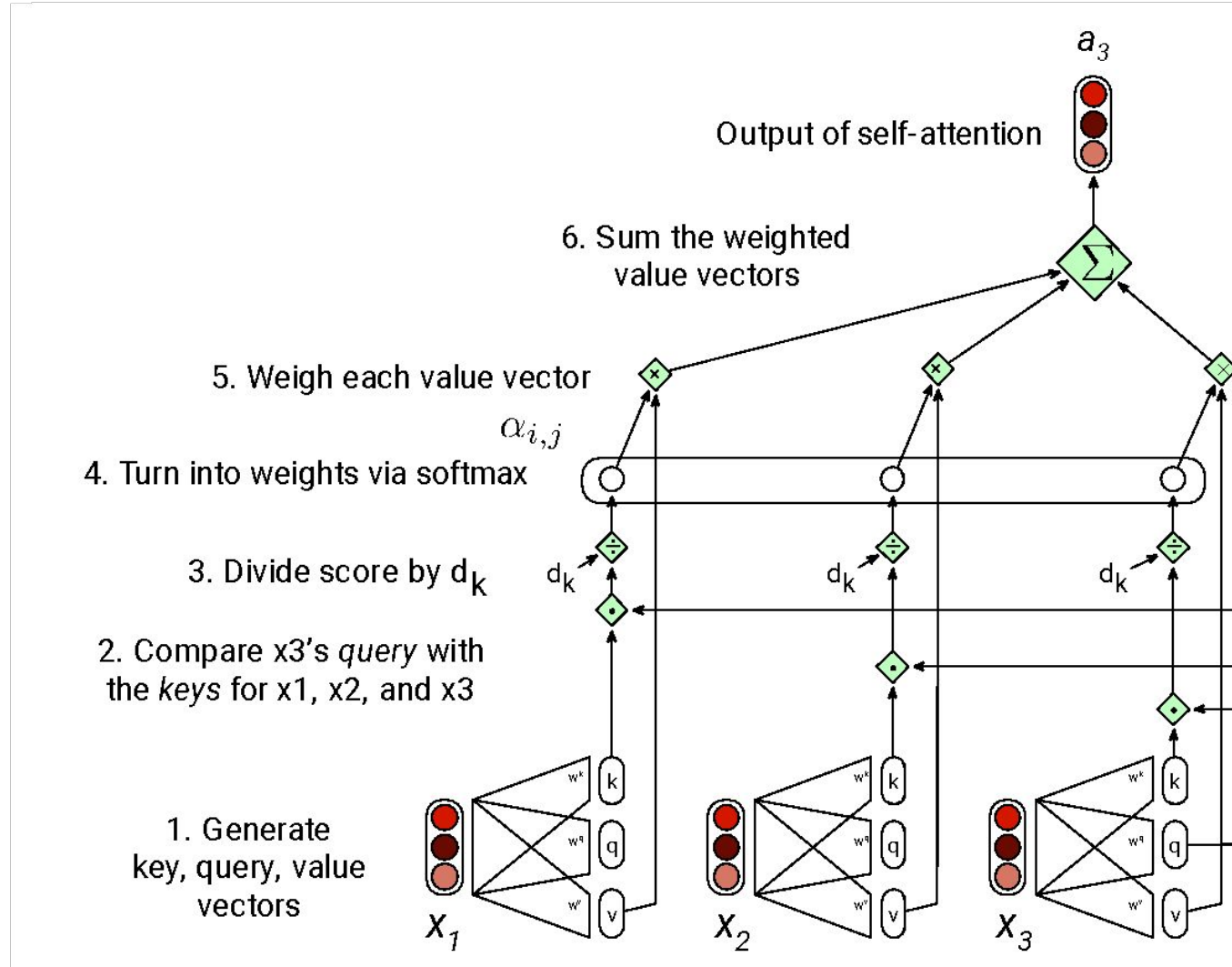


Query, Key and Value

- As *the current focus of attention* when being compared to all of the other preceding inputs. We'll refer to this role as a **query**.
- In its role as *a preceding input* being compared to the current focus of attention. We'll refer to this role as a **key**.
- And finally, as a **value** used to compute the output for the current focus of attention.

$$\mathbf{q}_i = \mathbf{x}_i \mathbf{W}^Q; \quad \mathbf{k}_i = \mathbf{x}_i \mathbf{W}^K; \quad \mathbf{v}_i = \mathbf{x}_i \mathbf{W}^V$$

Example Computation



Computation

$$q_i = x_i W^Q; k_i = x_i W^K; v_i = x_i W^V$$

$$\text{score}(x_i, x_j) = \frac{q_i \cdot k_j}{d_k}$$

$$a_{ij} = \frac{\text{softmax}(\text{score}(x_i, x_j))}{X} \quad \forall j \leq i$$

$$a_i = \sum_{j \leq i} a_{ij} v_j$$

Blank Slide

Parallelization of Computation

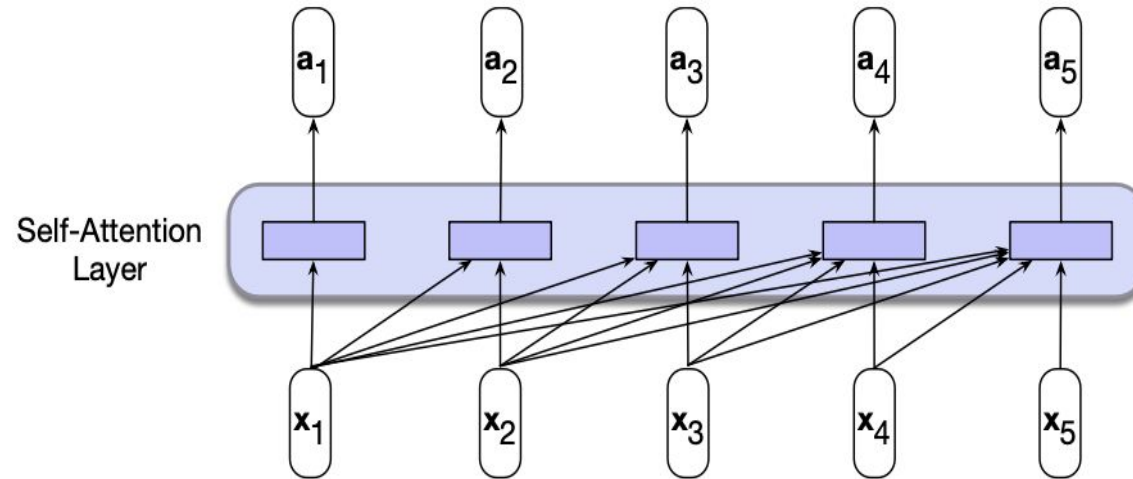
$$\mathbf{Q} = \mathbf{XW}^{\mathbf{Q}}; \quad \mathbf{K} = \mathbf{XW}^{\mathbf{K}}; \quad \mathbf{V} = \mathbf{XW}^{\mathbf{V}}$$

$$\mathbf{A} = \text{SelfAttention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{QK}^{\top}}{\sqrt{d_k}}\right) \mathbf{V}$$

Blank Slide

Self Attention: Casual/Masked

Masking out the future



Unlike RNNs, computation at each step are independent of all other time steps and can be done parallelly.

$$QK^T$$

N	q1•k1	—∞	—∞	—∞	—∞
	q2•k1	q2•k2	—∞	—∞	—∞
	q3•k1	q3•k2	q3•k3	—∞	—∞
	q4•k1	q4•k2	q4•k3	q4•k4	—∞
	q5•k1	q5•k2	q5•k3	q5•k4	q5•k5
N					

What is the typical value of N?
4096 tokens

Masked Attention

Don't let vectors "look ahead" in the sequence

Inputs:

Input vectors: X (Shape: $N_X \times D_Q$)

Key matrix: W_K (Shape: $D_X \times D_Q$)

Value matrix: W_V (Shape: $D_X \times D_V$)

Query matrix: W_Q (Shape: $D_Q \times D_Q$)

Computation:

Query Vectors $Q = XW_Q$

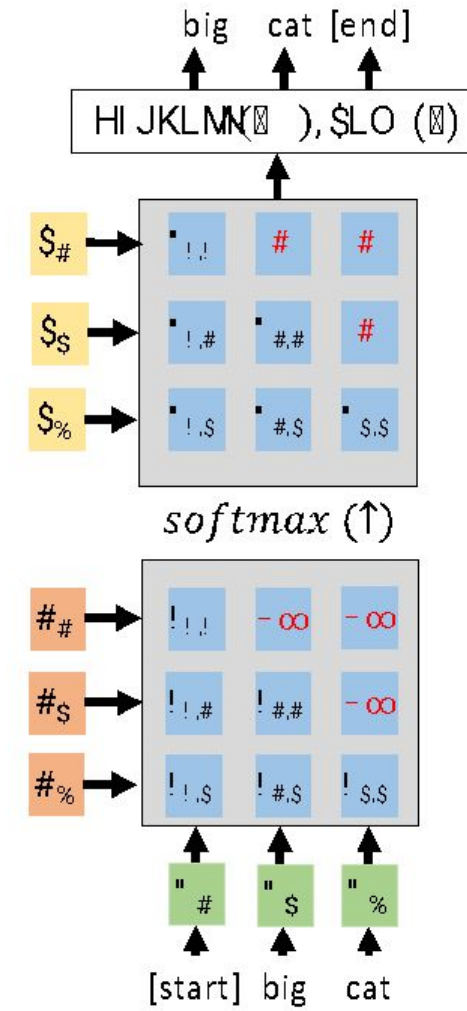
Key vectors: $K = XW_K$ (Shape: $N_X \times D_K$)

Value Vectors: $V = XW_V$ (Shape: $N_X \times D_V$)

Similarities: $E = \frac{Q \cdot K^T}{\sqrt{D_K}}$ (Shape: $N_X \times N_X$) $E_{ij} = (Q_i \cdot K_j) / \sqrt{D_K}$

Attention weights: $A = \text{softmax}(E, \text{dim} = 1)$ (Shape: $N_X \times N_X$)

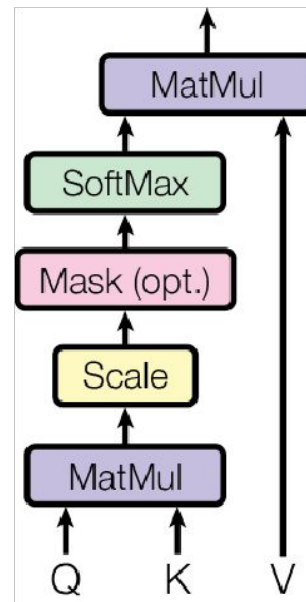
Output vectors: $Y = AV$ (Shape: $N_X \times D_V$) $Y_i = \sum_j A_{ij} V_j$



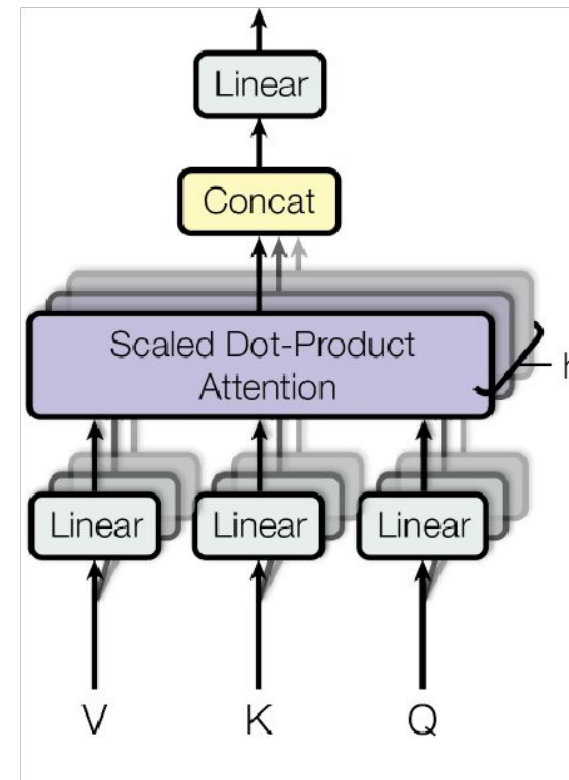
Multi Head Attention: Why?

- Why are we limited to one query, key and value?
- Multiple Feature Maps (like Convolution Layers)

Scaled Dot-Product Attention



Multi-Head Attention



Masked Multi-Head Attention

- Masked multi-head attention: multi-head where some values are masked (i.e., probabilities of masked values are nullified to prevent them from being selected).
- When decoding, an output value should only depend on previous outputs (not future outputs). Hence we mask future outputs.

$$attention(Q, K, V) = softmax\left(\frac{Q^T K}{\sqrt{d_k}}\right) V$$

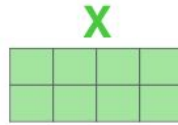
$$maskedAttention(Q, K, V) = softmax\left(\frac{Q^T K + M}{\sqrt{d_k}}\right) V$$

where M is a mask matrix of 0's and $-\infty$'s

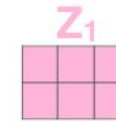
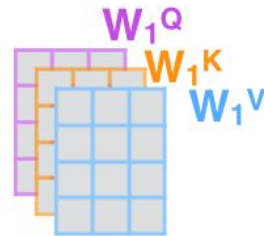
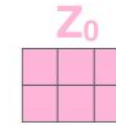
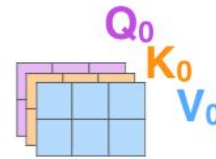
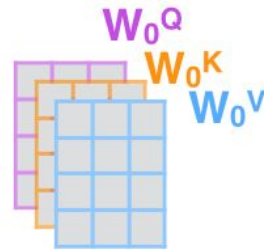
Multi-Head Attention

- 1) This is our input sentence*
- 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

Thinking
Machines



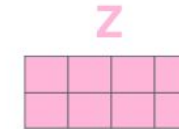
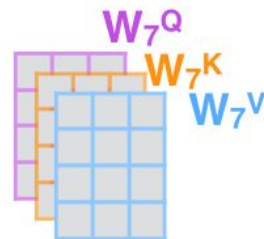
* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one



...

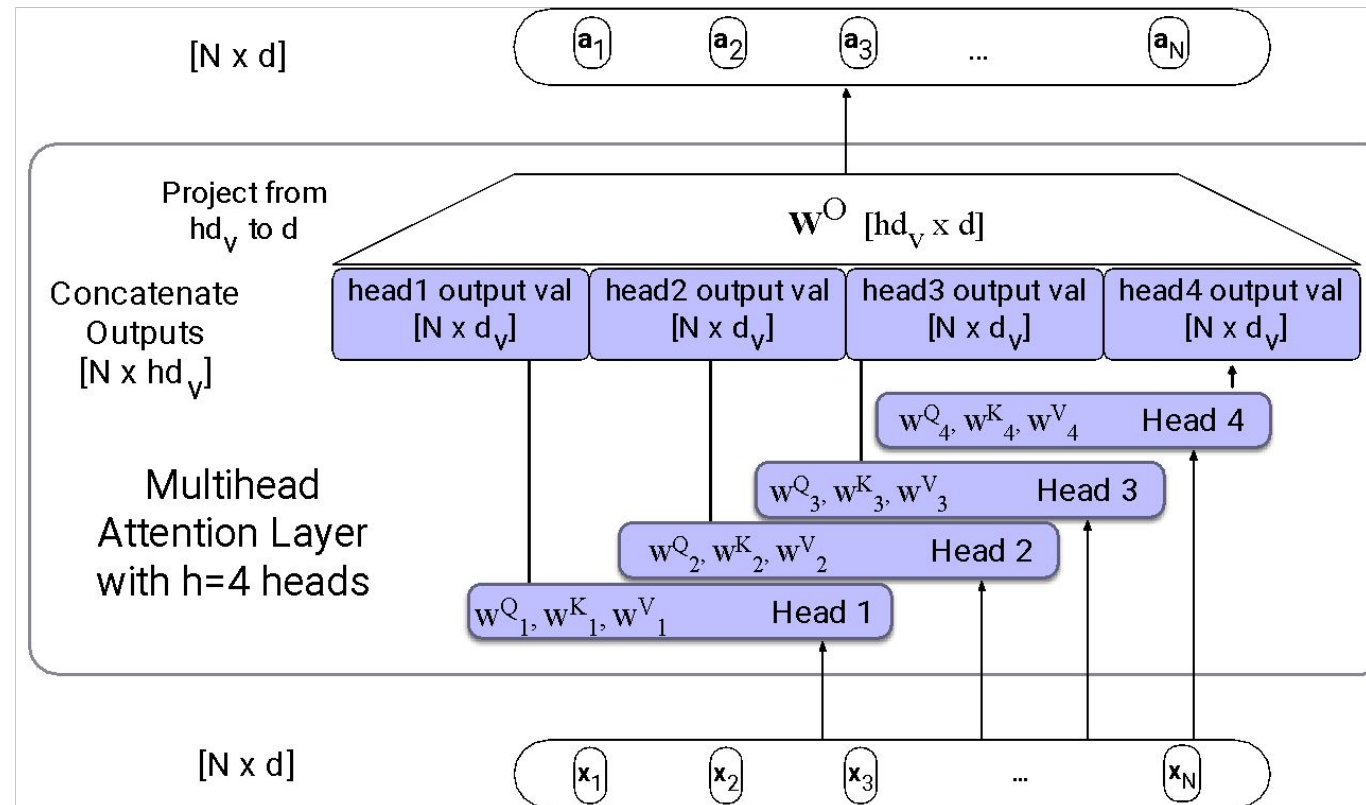
...

...



Blank Slide

Multi Head Attention

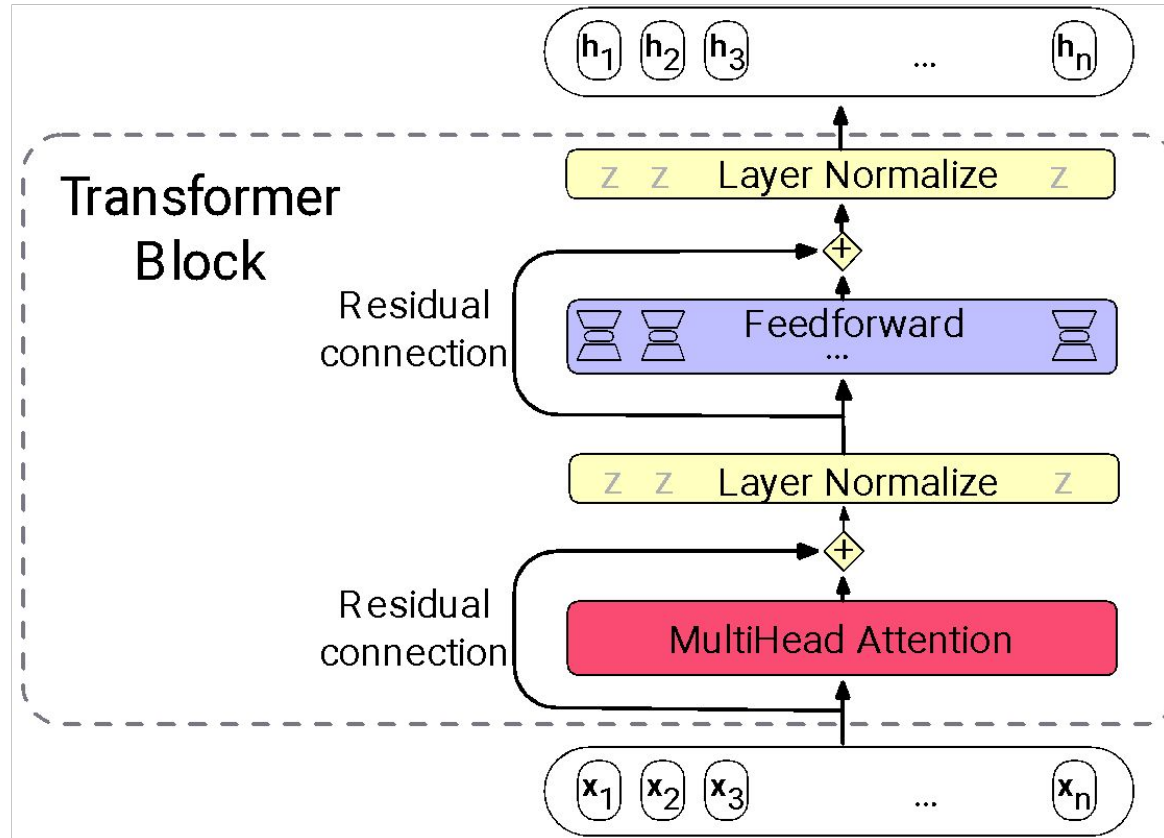


$$\mathbf{Q} = \mathbf{XW}_i^Q ; \mathbf{K} = \mathbf{XW}_i^K ; \mathbf{V} = \mathbf{XW}_i^V$$

$$\text{head}_i = \text{SelfAttention}(\mathbf{Q}, \mathbf{K}, \mathbf{V})$$

$$\mathbf{A} = \text{MultiHeadAttention}(\mathbf{X}) = (\text{head}_1 \oplus \text{head}_2 \dots \oplus \text{head}_h) \mathbf{W}^O$$

Transformer Block



$$\mathbf{T}^1 = \text{SelfAttention}(\mathbf{X})$$

$$\mathbf{T}^2 = \mathbf{X} + \mathbf{T}^1$$

$$\mathbf{T}_3 = \text{LayerNorm}(\mathbf{T}^2)$$

$$\mathbf{T}^4 = \text{FFN}(\mathbf{T}^3)$$

$$\mathbf{T}^5 = \mathbf{T}^4 + \mathbf{T}^3$$

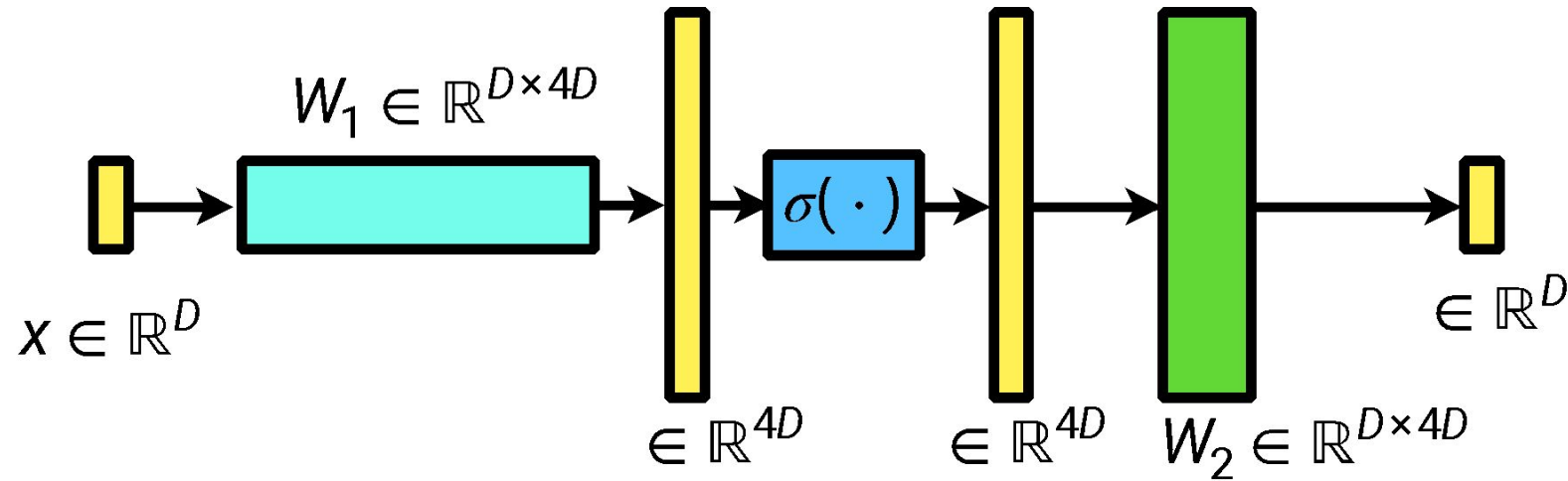
$$\mathbf{H} = \text{LayerNorm}(\mathbf{T}^5)$$

Both X and H are of N X d

Blank Slide

Blank Slide

MLP/Feed Forward Layer



$$\text{MLP}(x) = W_2 \sigma(W_1 x + b_1) + b_2$$

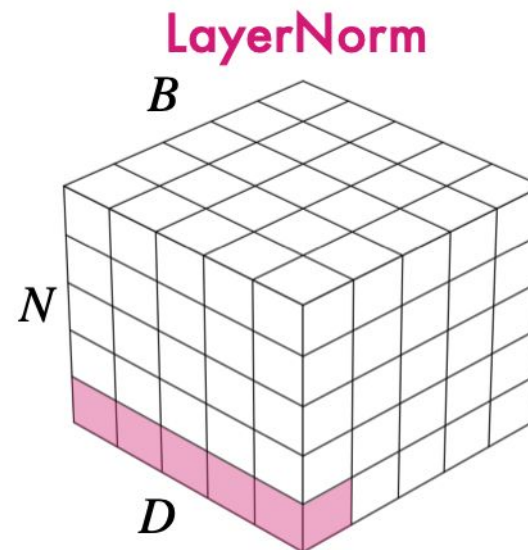
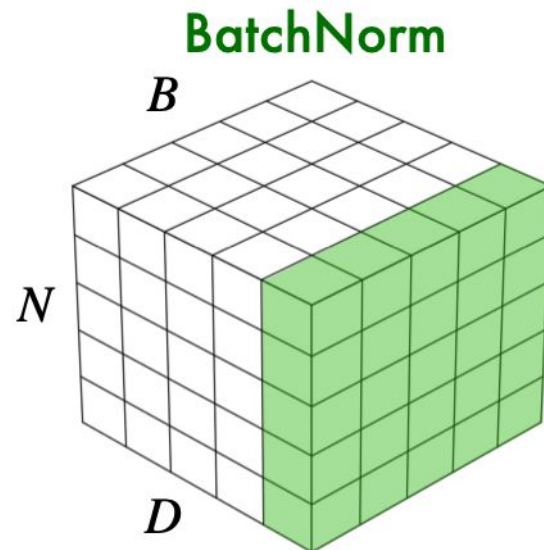
Layer Norm Enhances the Stability

$$\mu = \frac{1}{d_h} \sum_{i=1}^{d_h} x_i$$

$$\sigma = \sqrt{\frac{1}{d_h} \sum_{i=1}^{d_h} (x_i - \mu)^2}$$

$$\hat{\mathbf{x}} = \frac{(\mathbf{x} - \mu)}{\sigma}$$

$$\text{LayerNorm} = \gamma \hat{\mathbf{x}} + \beta$$



LayerNorm has

- No dependence on batch dim.
- Same procedure at train/test time

Positional Encoding

- Position embedding soon after input embedding
- “Bag of words” to “ordered words”
- Integer (position) to a vector of say 100s

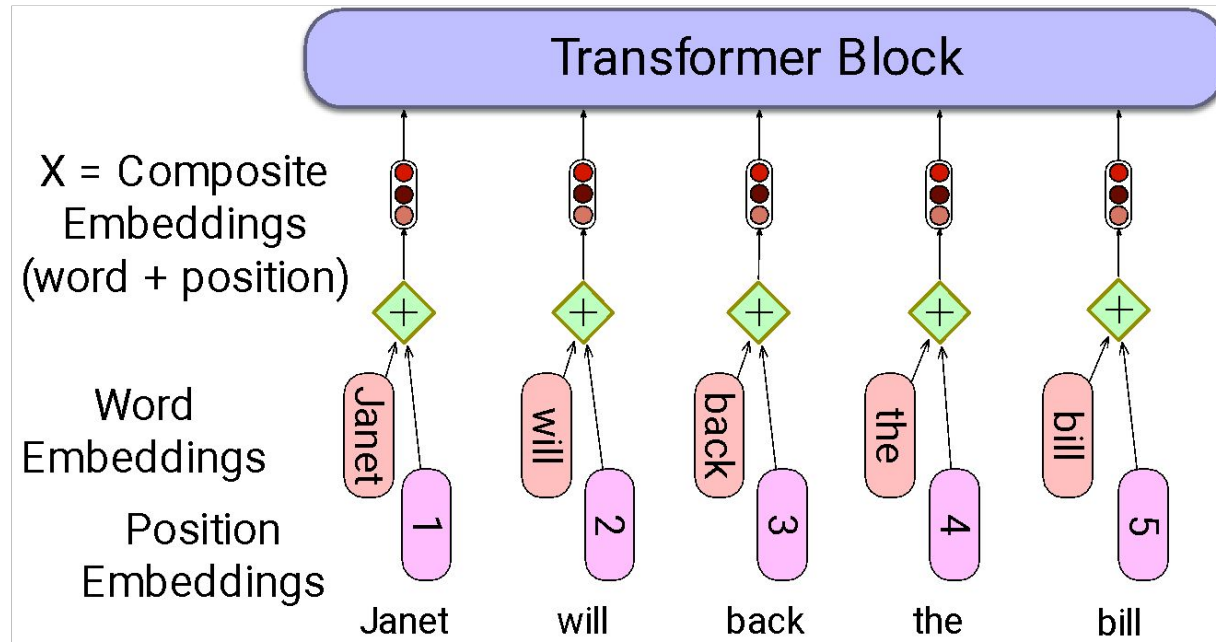
Positional embedding

- Embedding to distinguish each position

$$PE_{position,2i} = \sin(position/10000^{2i/d})$$

$$PE_{position,2i+1} = \cos(position/10000^{2i/d})$$

Position Embedding

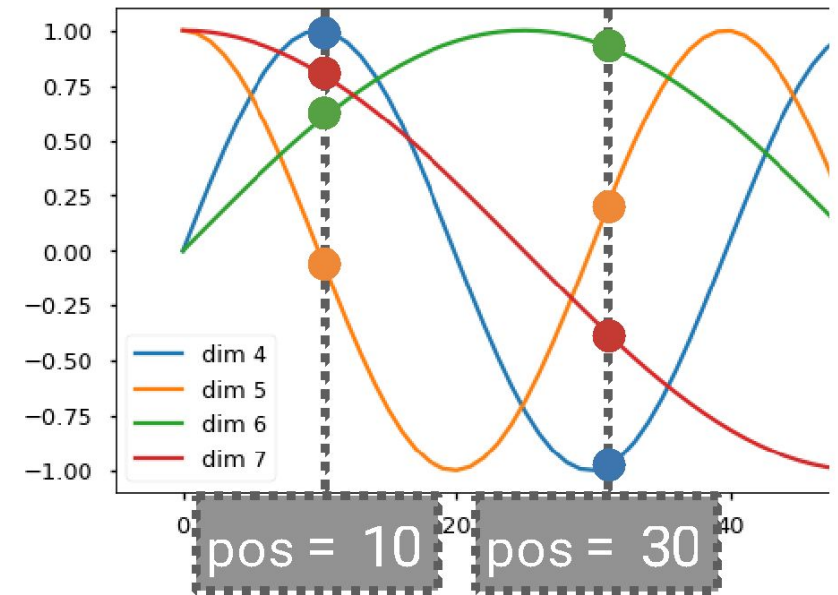


Absolute Positions

Choose static functions (say sinusoids) that map integer inputs to real valued vectors that captures inherent relationship between vectors

$$PE(pos, 2i) = \sin\left(\frac{pos}{10,000^{2i/D}}\right)$$

$$PE(pos, 2i + 1) = \cos\left(\frac{pos}{10,000^{2i/D}}\right)$$



Blank Slide

puts:

key matrix: $W_K(Shape: D_X \times D_Q)$

value matrix: W_V (*Shape:* $D_X \times D_V$)

Query matrix: \mathbf{W}_Q (Shape: $D_Q \times D_Q$)

Query Vectors $Q = XW_1$

low Vectors: $V = XW_*$ (Shape: $N_* \times D_*$)

$$\text{similarities: } E = \frac{1}{\sqrt{S_n}} (\text{Shape: } N_n \times N_n) E_{\phi} = (Q_{\phi} \cdot K) / \sqrt{D_{\phi}}$$

tention weights: $A = \text{softmax}(E, \text{dim} = 1)$ (*Shape: $N \times N$*)

Input vectors: $Y = AV$ (Shape: $N \times D_*$) $Y_{\%} = \Sigma A_{\%} V$

In order to make processing position aware, concatenate or add positional encoding to the input

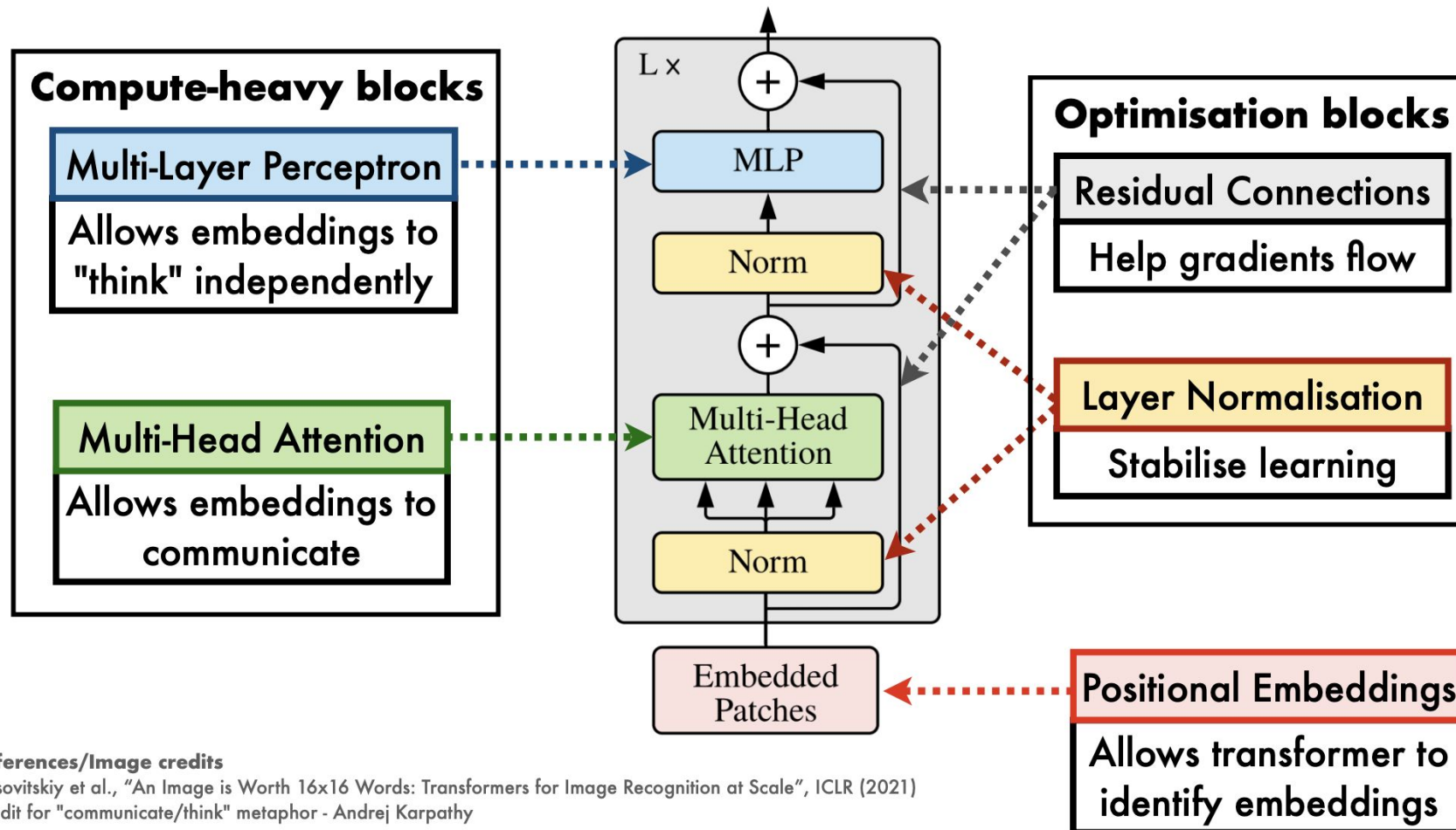
- can be learned lookup table, or fixed function



Transformers

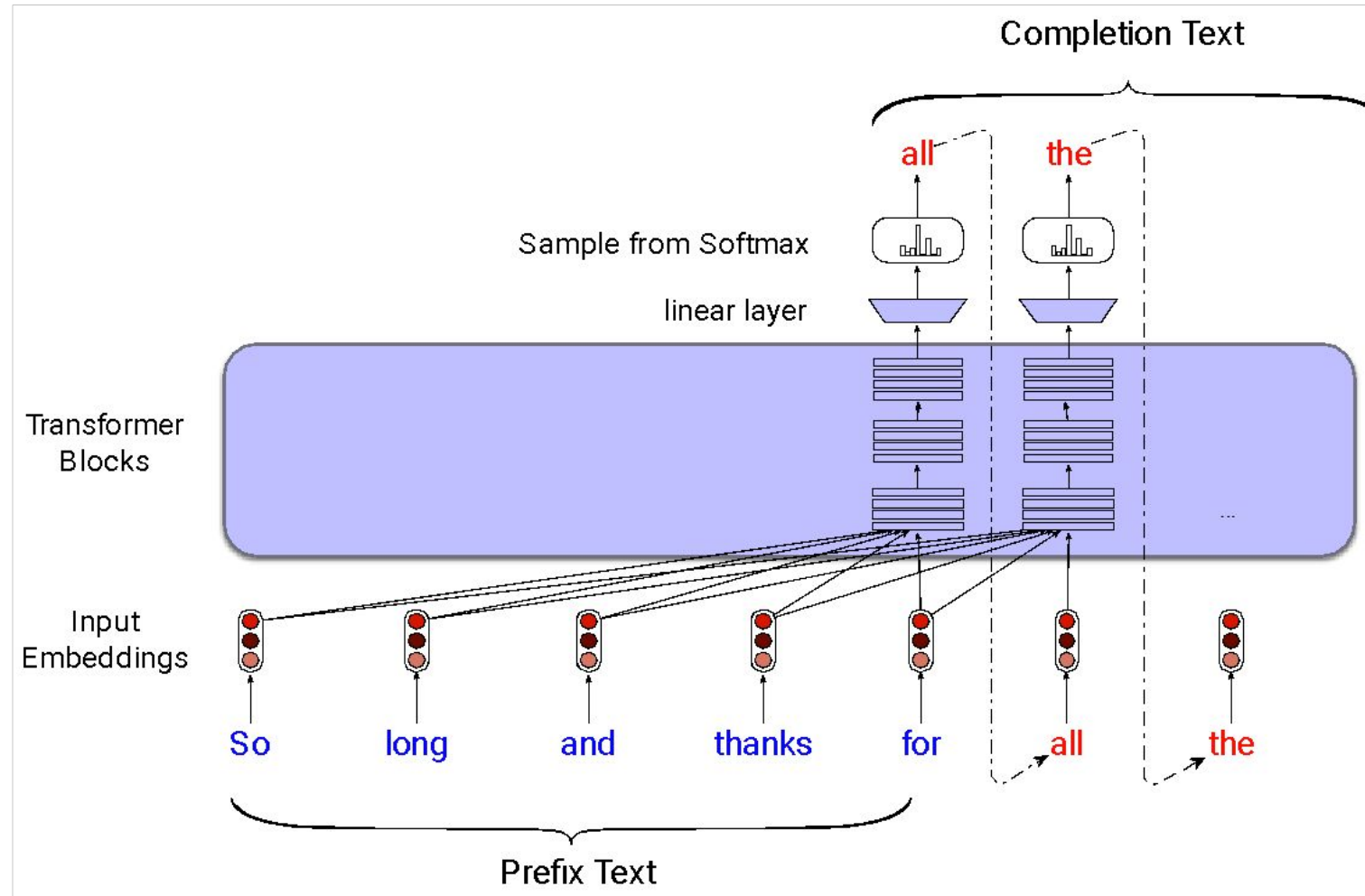
Transformer Encoder

Five key ideas

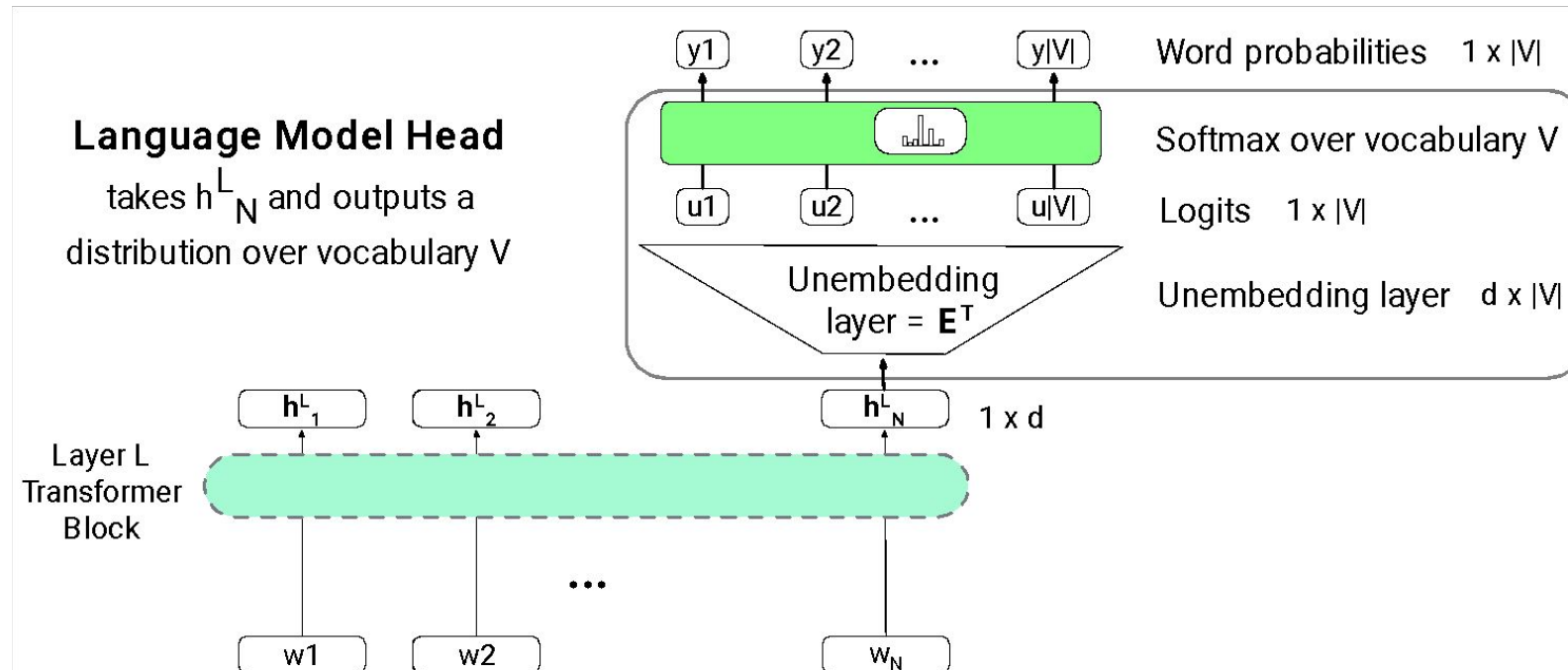


Blank Slide

Sentence Completion



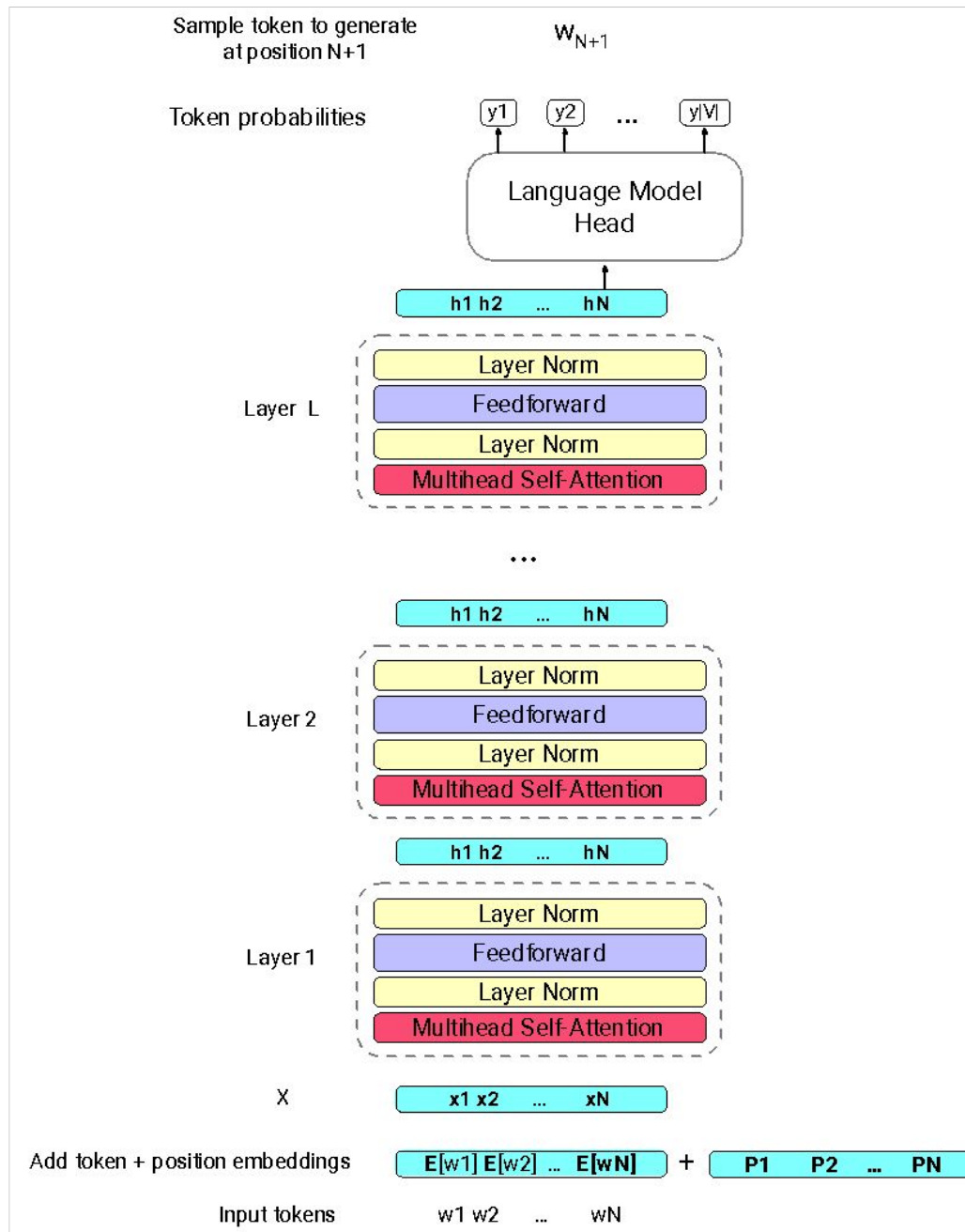
Language Modelling Head



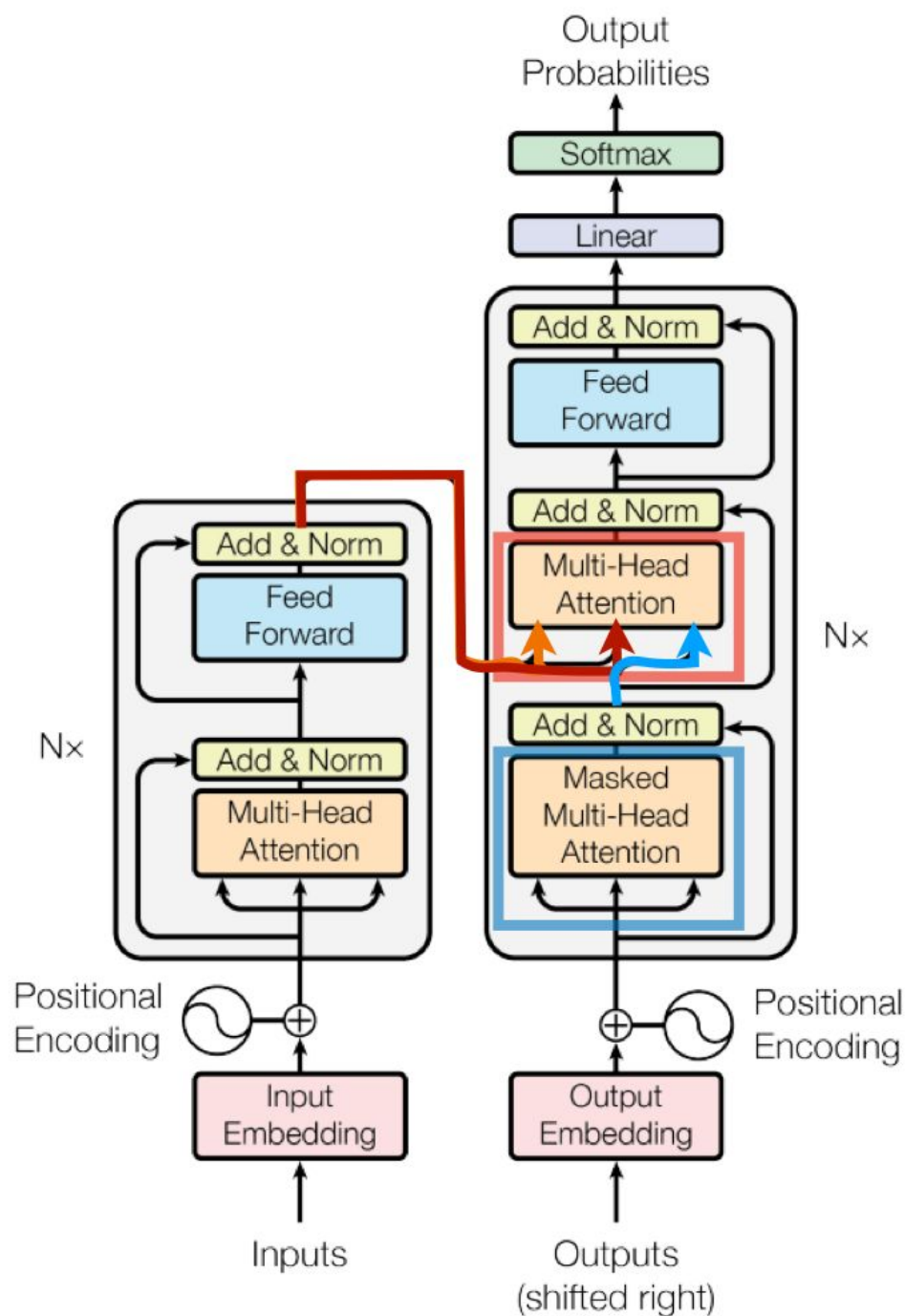
$$u = h_N^L E^T$$

$$y = \text{softmax}(u)$$

Blank Slide



Eg. A Complete Decoder Model



Cross Attention

Query from one sequence and key + value from other (in translation)

Thanks!!

Questions?