

E9: 309 Advanced Deep Learning

4-11-2020

<http://leap.ee.iisc.ac.in/sriram/teaching/ADL2020/>

Housekeeping

* 1st mini-project

✓ Deadlines

- ★ Abstract submission deadline (Nov 2nd, Monday)
 - ★ Using the google form given in the webpage
- ★ Solo projects or 2-member projects
 - ★ Indicate roles of each member in 2-member project
 - ★ 200 page abstract of the work. If modifications are needed, we will review and let you know in 2-3 days.



Recap of previous class



State of affairs

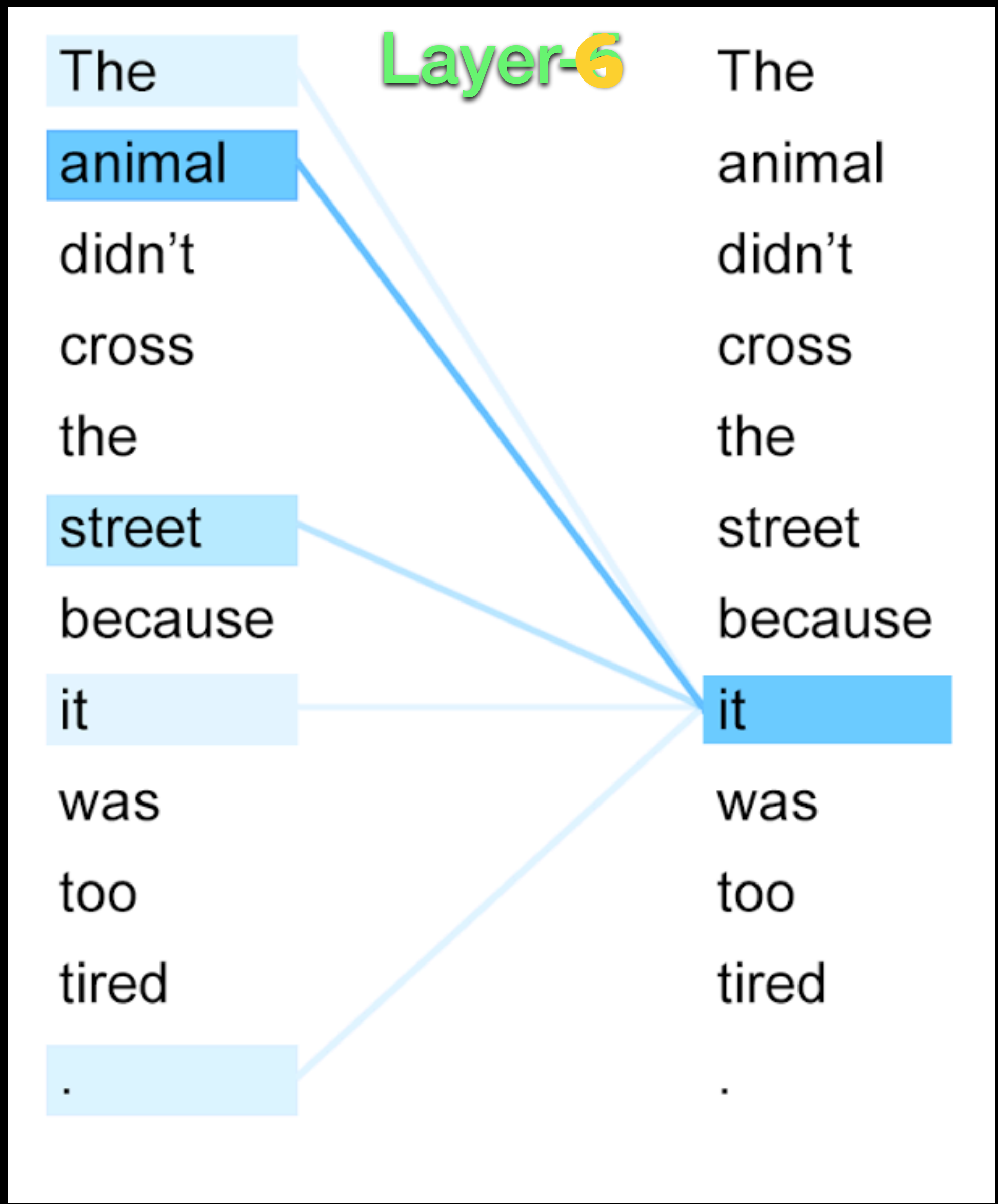
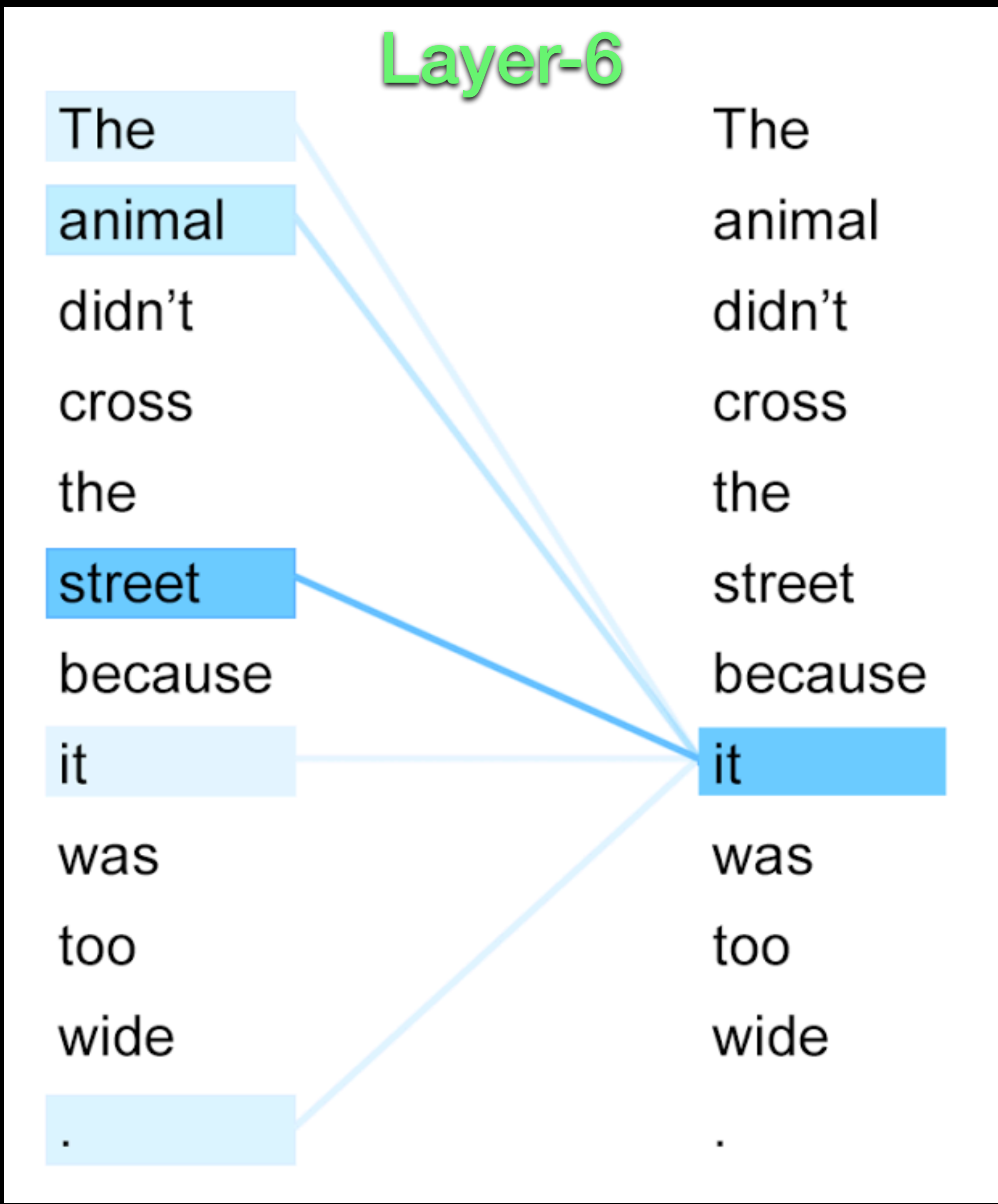
✦ Transformer models

✓ Encoder-decoder model

- ★ Repeated architecture in encoder and decoder stages
- ★ Multi-head self-attention
- ★ Position wise feed-forward.
- ★ Skip connection with layer norm.

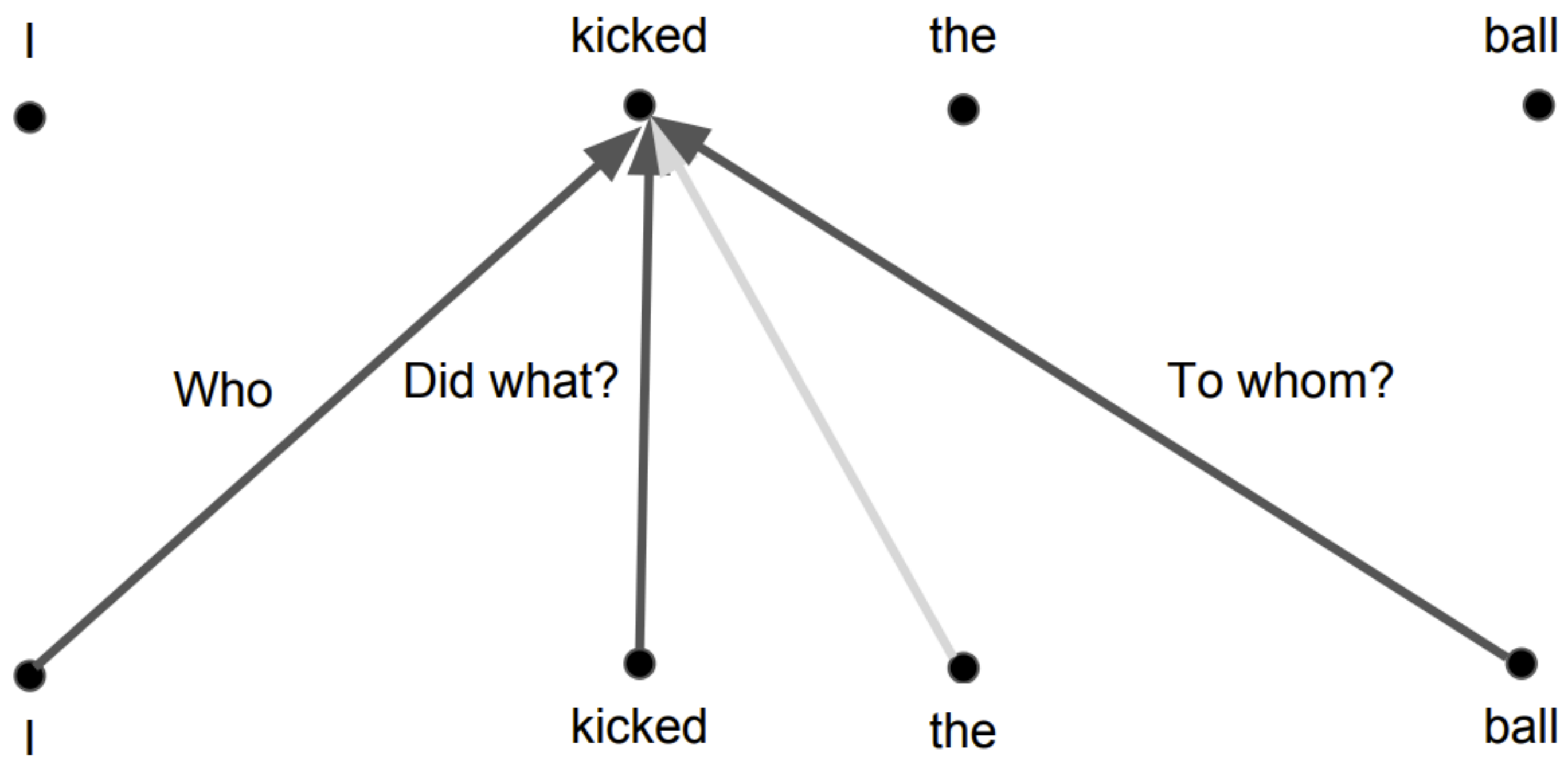


Self-attention - need for depth



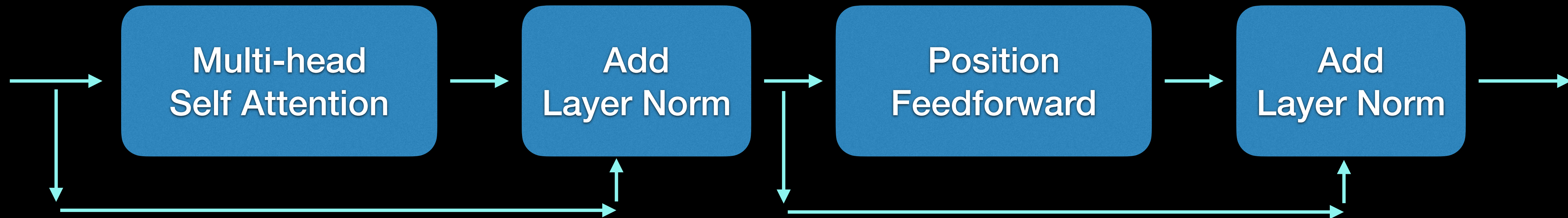
Need for multi-head attention

Self-Attention



Single layer of encoder (typical implementation)

- * Single encoder layer has typically self-attention skip connection, layer norm and feedforward layer



Positional encoding

* No recurrence or position awareness yet in the model

Binary format -

position can encode the rate of change of bits across time

In floating format - one can use sines and cosines

0 :	0	0	0	0	8 :	1	0	0	0
1 :	0	0	0	1	9 :	1	0	0	1
2 :	0	0	1	0	10 :	1	0	1	0
3 :	0	0	1	1	11 :	1	0	1	1
4 :	0	1	0	0	12 :	1	1	0	0
5 :	0	1	0	1	13 :	1	1	0	1
6 :	0	1	1	0	14 :	1	1	1	0
7 :	0	1	1	1	15 :	1	1	1	1



Positional encoding

* An example used in the first paper

$$\mathbf{p}(t) \in \mathcal{R}^D$$

$$p_i(t) = \begin{cases} \sin(\omega_k t), & \text{if } i = 2k \\ \cos(\omega_k t), & \text{if } i = 2k + 1 \end{cases}$$

$i = \dots$

$$k \in \{1 \dots \frac{D}{2}\}$$

$$\omega_k = \frac{1}{10000^{\frac{2k}{D}}}$$

model

$$\mathbf{x}(t) = \mathbf{x}(t) + \mathbf{p}(t)$$

$t = 1 \dots T$

$x(1), x(2), \dots, x(T)$

$$X = \{x(1) \dots x(T)\}$$

$\in \mathbb{R}^{T \times D}$

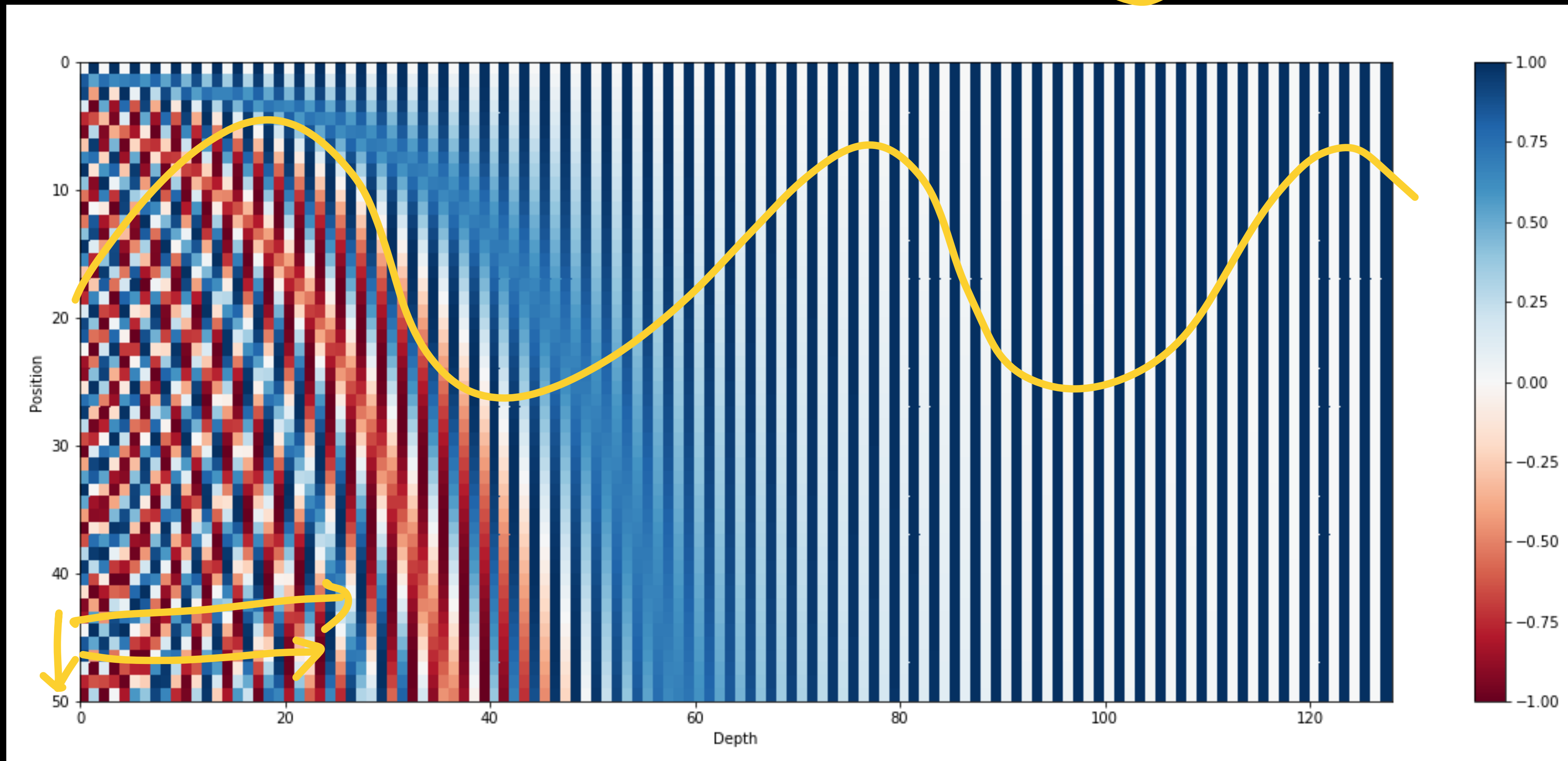
$x(T)$



Positional encoding

$T=20$
 $T=100$

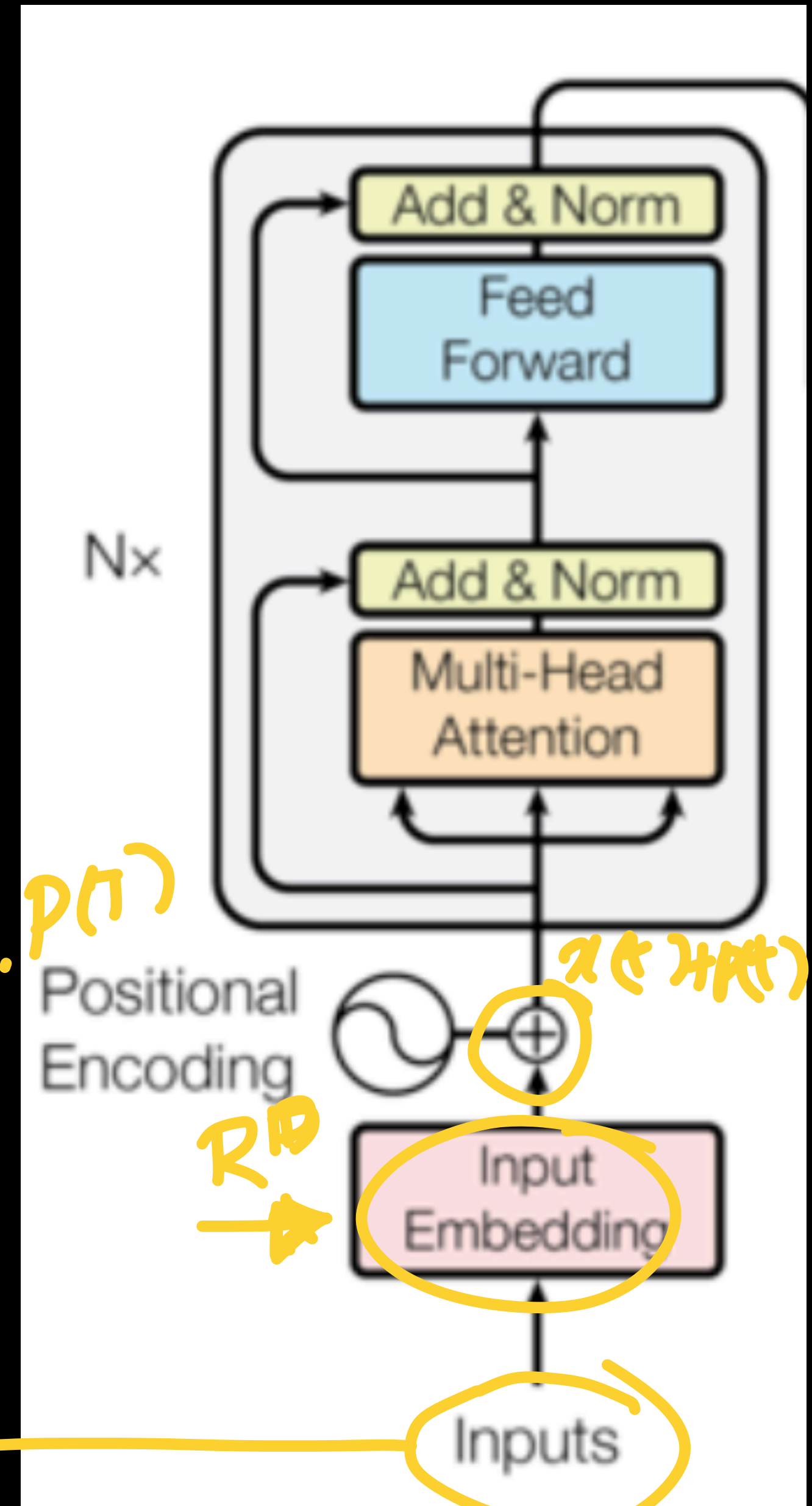
* An example used in the first paper $\mathbf{p}(t) \in \mathcal{R}^D$ [$T=50$, $D=128$]



$(1), (2), (3) \dots (T)$
 $p(1), p(2), p(3) \dots p(T)$

Transformer encoder - overview

Reading Assignment - "Attention is All You Need"
<https://arxiv.org/pdf/1706.03762.pdf>



$p(1) \dots p(T)$
 T, D

$$X = \left\{ \underset{\substack{\uparrow \\ 128D}}{x^{(1)}} \dots \underset{\substack{\uparrow \\ 128D}}{x^{(T)}} \right\}$$



Transformer decoder

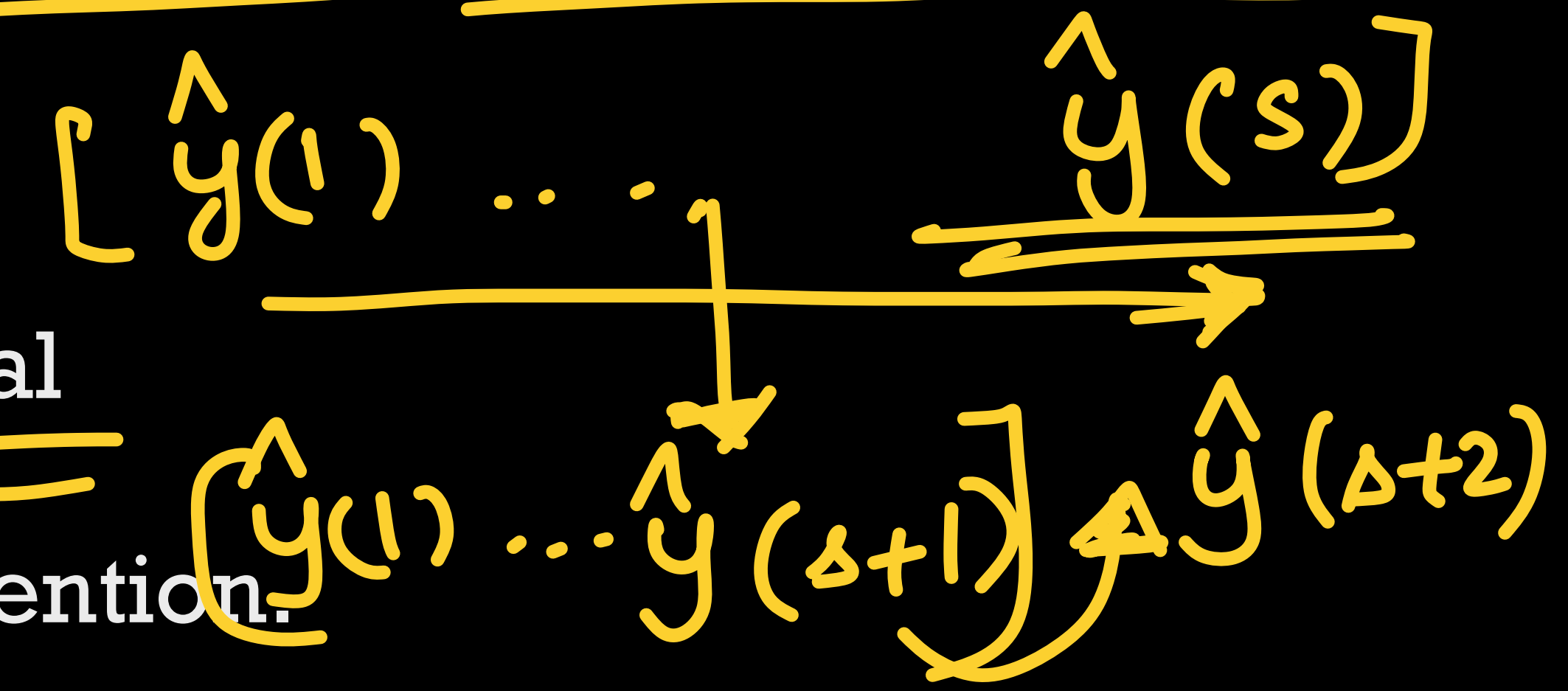
$\hat{y}(s) \dots \hat{y}(1)$

$x(1) \dots x(T)$

* Masked self-attention layer -

✓ Mask makes the output dependencies causal

★ Only the past is used to encode the attention.



<s>	der	schnelle	braune	fuchs	</s>
<s>					
der					
schnelle		100	20	80	
braune			100	80	
fuchs					
</s>					

$e^{-inf} = 0$

0	-inf	-inf	-inf	-inf	-inf
0	0	-inf	-inf	-inf	-inf
0	0	0	-inf	-inf	-inf
0	0	0	0	-inf	-inf
0	0	0	0	0	-inf
0	0	0	0	0	0

Mask Shape = TxT = 6x6

<s>	der	schnelle	braune	fuchs	</s>
<s>	-inf	-inf	-inf	-inf	-inf
der		inf	-inf	-inf	-inf
schnelle		100	-inf	-inf	-inf
braune			100	-inf	-inf
fuchs					-inf
</s>					

QxK^T / \sqrt{d}

QxK^T / \sqrt{d}



Transformer decoder

* Masked self-attention layer -

✓ Mask makes the output dependencies causal

★ Only the past is used to encode the attention.

$$\text{Softmax} \left\{ \frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}} \right\} \mathbf{V} \xrightarrow{A_n} \text{Softmax} \left\{ \text{Mask} + \frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}} \right\} \mathbf{V}$$

★ Make the attention matrix to be lower triangular



Transformer decoder

* Masked self-attention layer -

✓ Mask makes the output dependencies causal

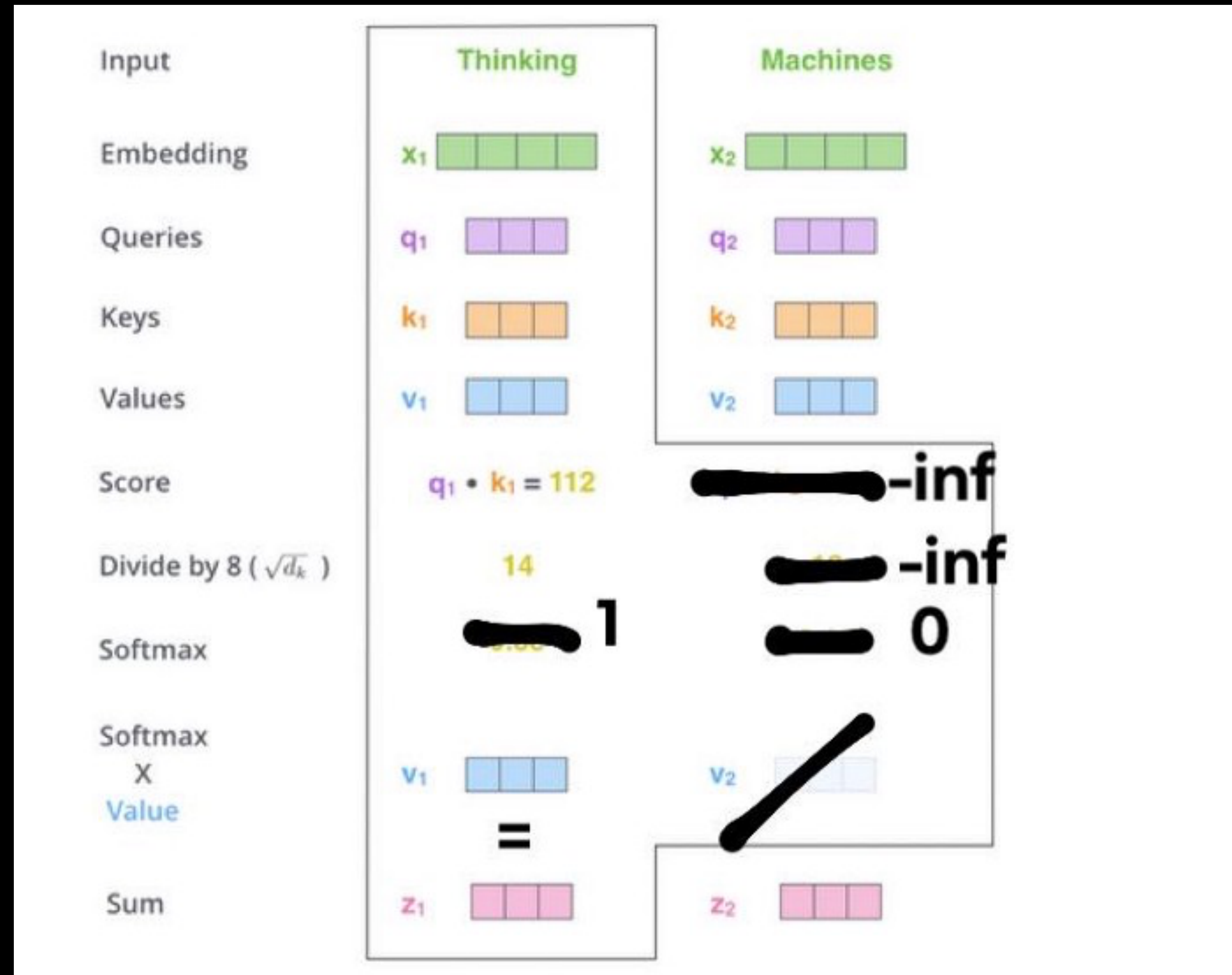
★ Only the past is used to encode the attention.

	<s>	der	schnelle	braune	fuchs	</s>
<s>		0	0	0	0	0
der			0	0	0	0
schnelle				.75	0	0
braune					.85	0
fuchs						0
</s>						

$\text{softmax}(QxK^T / \sqrt{d})$



Transformer decoder



Encoder-decoder attention

$$x = \{x^{(1)} \dots x^{(T)}\}$$

$$\hat{y} = \{\hat{y}^{(1)} \dots \hat{y}^{(S)}\}$$

$T \times D$
 $S \times D$

* Use the Key and Value matrices from the last layer of the encoder

$$Q_h^p = D^{p-1} W_h^{p,Q} + \mathbf{1} (b_h^{p,Q})^T \in \mathcal{R}^{S \times d}$$

$$K_h^p = E^L W_h^{p,K} + \mathbf{1} (b_h^{p,K})^T \in \mathcal{R}^{T \times d}$$

$$V_h^p = E^L W_h^{p,V} + \mathbf{1} (b_h^{p,V})^T \in \mathcal{R}^{T \times d}$$

$p = 1 \dots P$
↑
decoder layer

$$D_h^p = \text{softmax} \left(\frac{Q_h^p (K_h^p)^T}{\sqrt{d}} \right) V_h^p \in \mathcal{R}^{S \times d}$$

$S \times T$ $S \times d$

$$h = \{1..H\}$$

heads

$$d = \frac{D}{H}$$

$$(S \times d)(d \times d) \rightarrow S \times d$$

$S \times D$

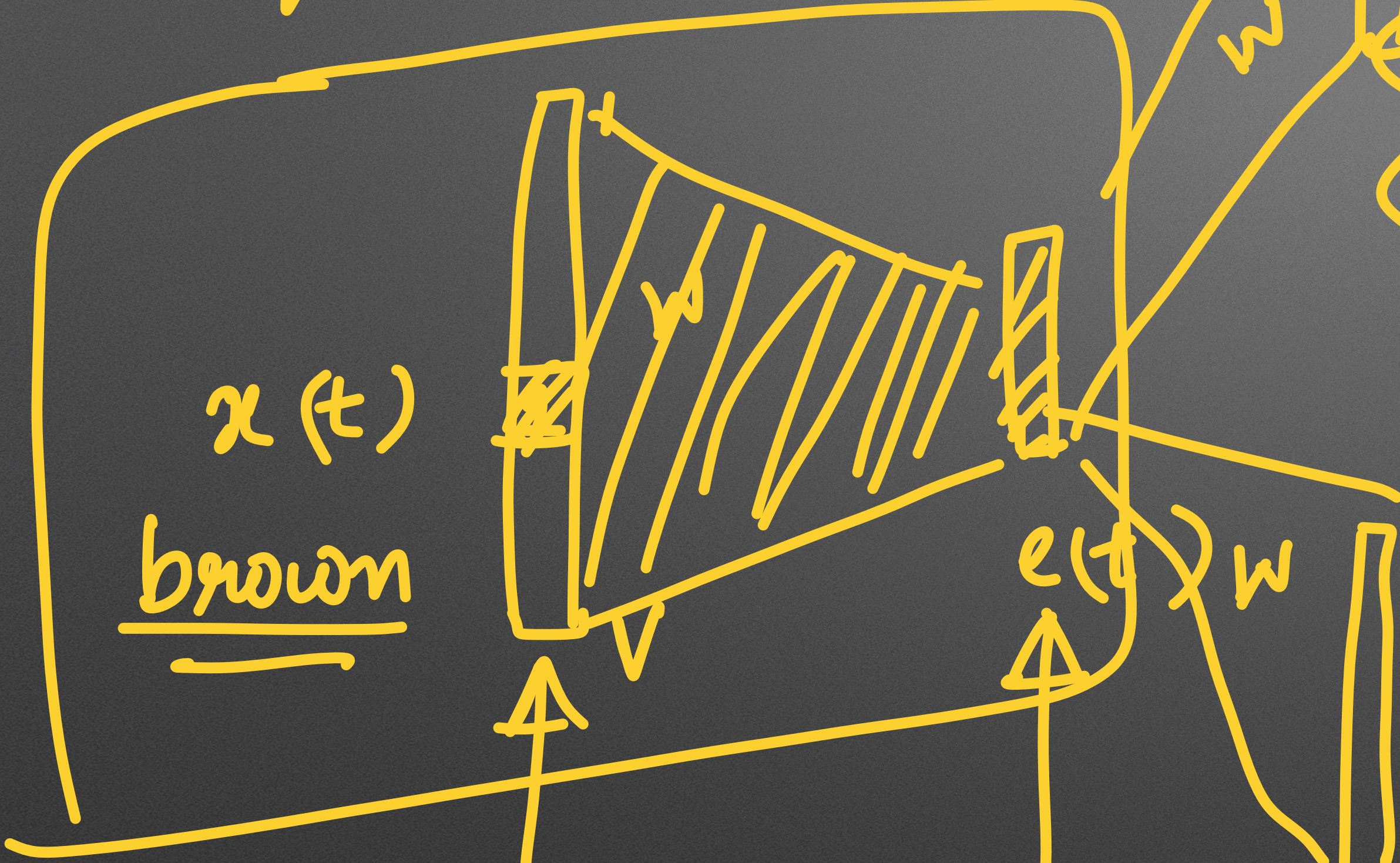


The quick brown fox

Eng V = 10K

x(t-1)
softmax
quick {2000}

V_1
 V_2
 $V = \{V_1, V_2\}$



$x(t+1)$
softmax v

one hot

embeddings for x {8862}
D

D << V

Transformer - decoder

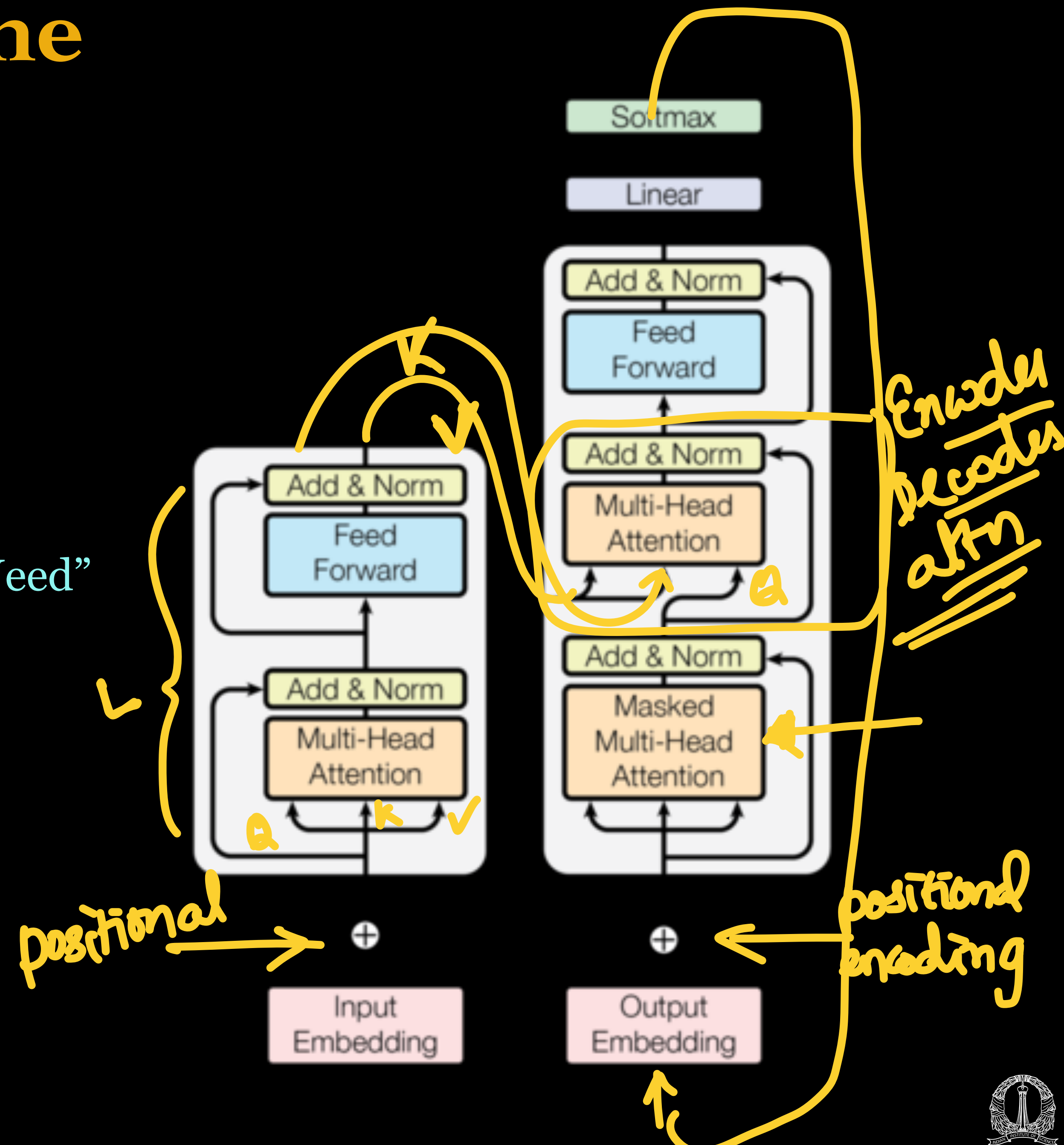
* Decoder Layer Output

$$* \quad [\mathbf{d}^p(1) \dots \mathbf{d}^p(S)] = \text{ReLU} \left(\mathbf{D}_{ff}^p \mathbf{W}_{of}^p + \mathbf{1} (\mathbf{b}_{of}^p)^T \right) \in \mathcal{R}^{S \times D}$$



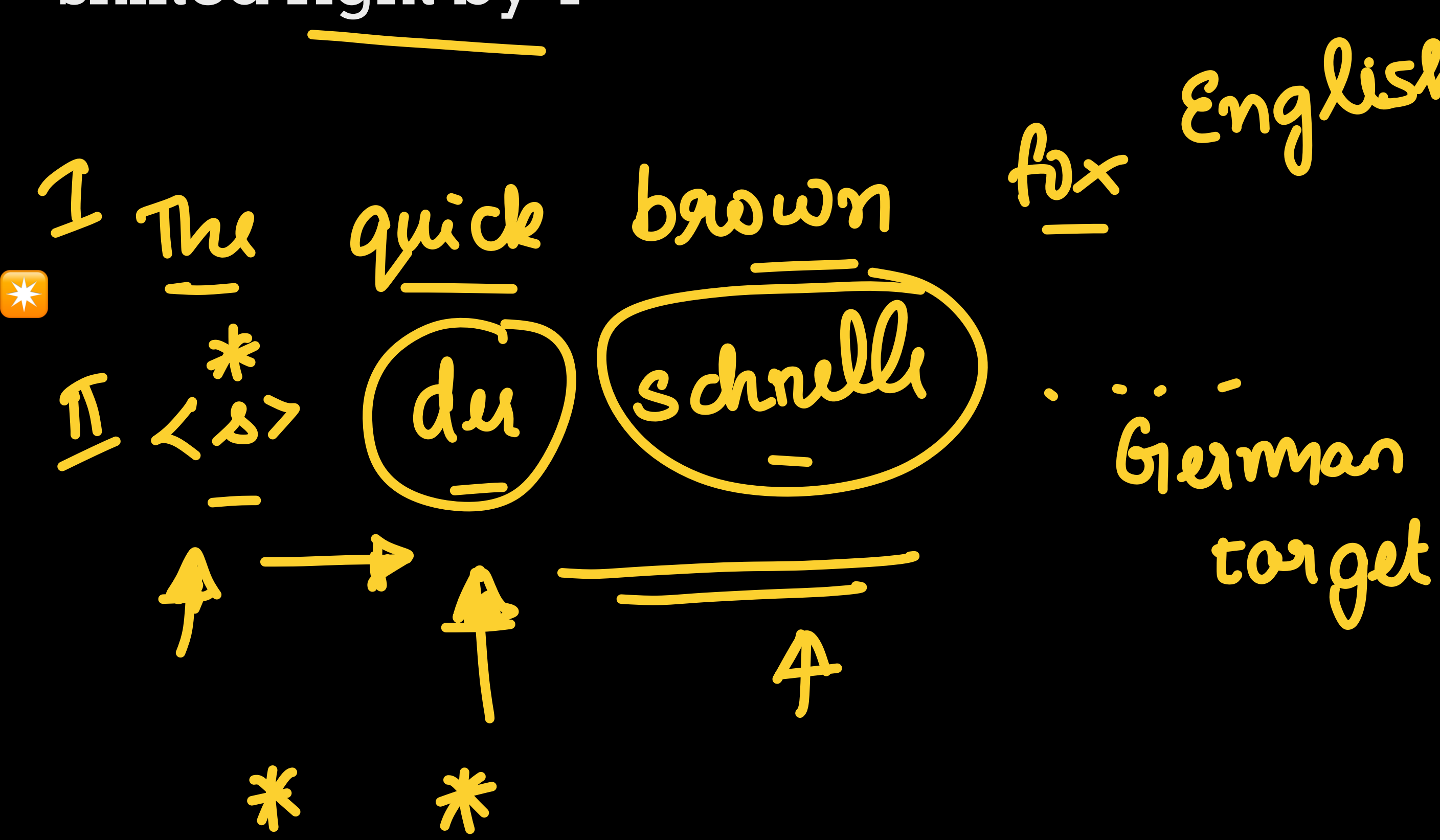
Transformer - full pipeline

Reading Assignment - "Attention is All You Need"
<https://arxiv.org/pdf/1706.03762.pdf>

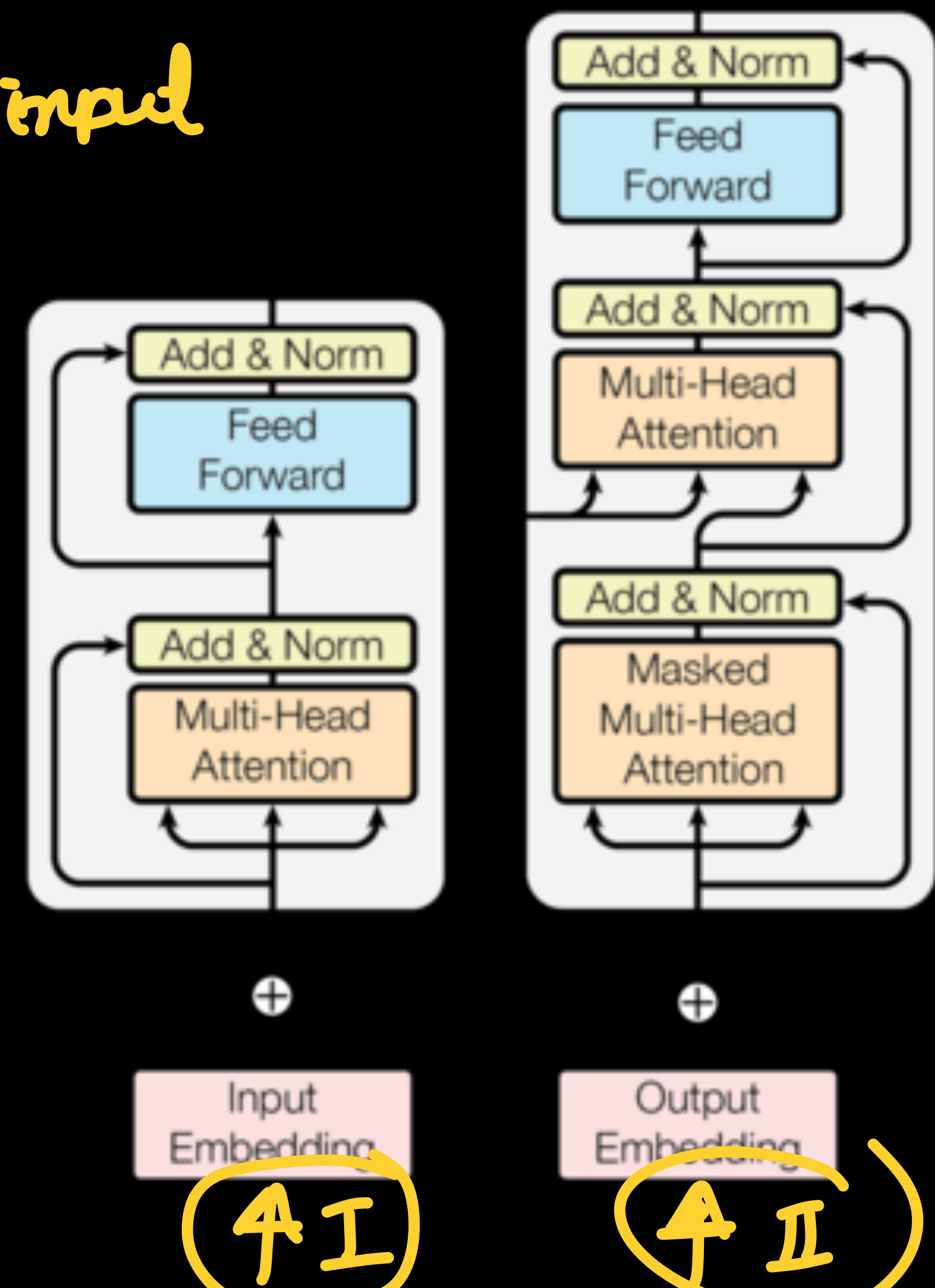


Transformer - training

* Input sequence and the output sequence shifted right by 1



$\pi <$ one hot target

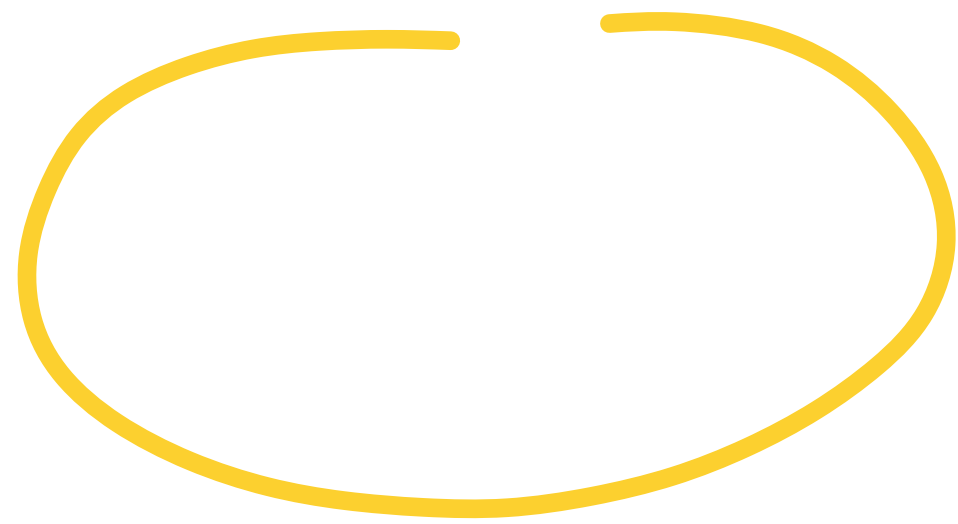


Transformer decoder

* Inference



Greedy decoding



Topics thus far ...

Visual and Time Series Modeling: Semantic Models, Recurrent neural models and LSTM models, Encoder-decoder models, Attention models.

Representation Learning, Causality And Explainability: t-SNE visualization, Hierarchical Representation, semantic embeddings, gradient and perturbation analysis, Topics in Explainable learning, Structural causal models.

Unsupervised Learning: Restricted Boltzmann Machines, Variational Autoencoders, Generative Adversarial Networks.

New Architectures: Capsule networks, End-to-end models, Transformer Networks.

Applications: Applications in NLP, Speech, Image/Video domains in all modules.

