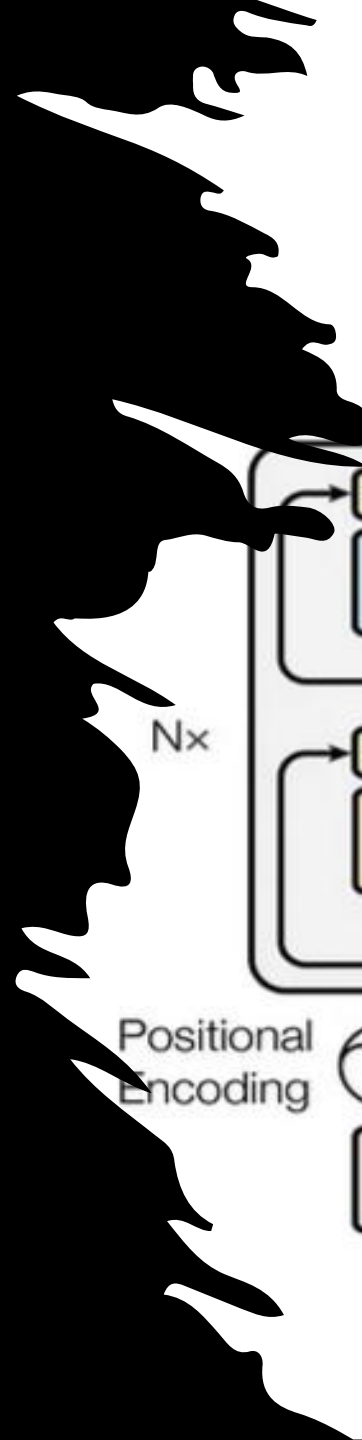


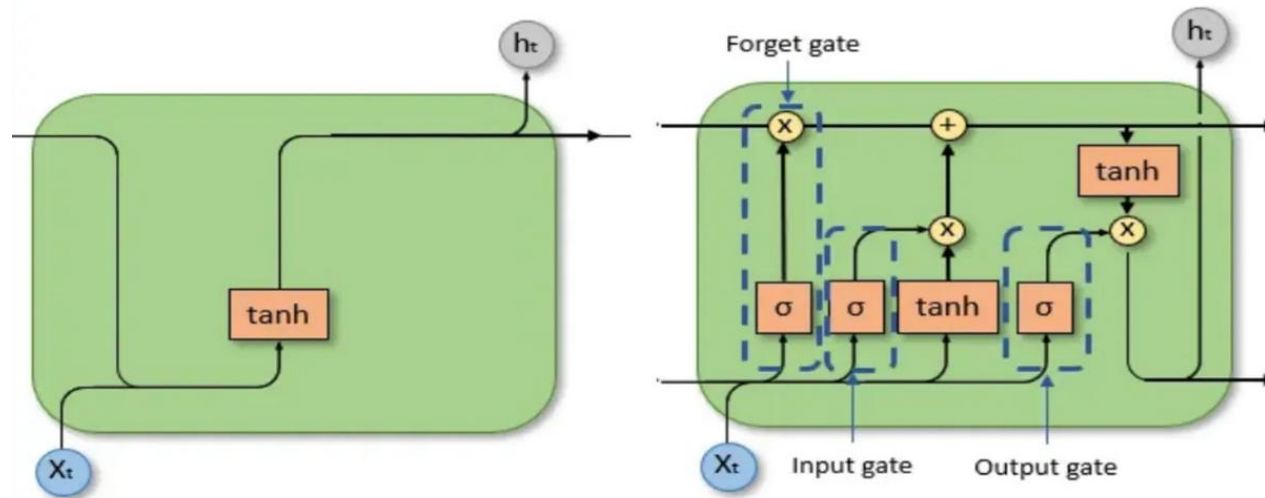
**“ATTENTION IS ALL
YOU NEED”**

TRANSFORMERS AND SELF-ATTENTION

**PRESENTER: SHAKKYA
RANASINGHE**



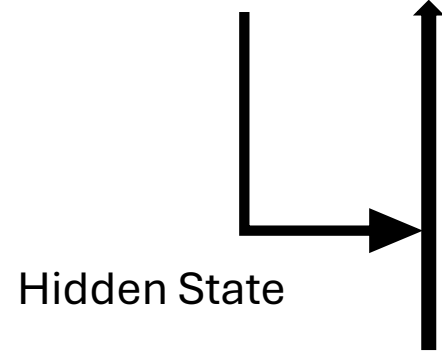
Once there was RNN's....



- Sequential Data
- Retained Information through Hidden States
- Improved with the use of cell states and gates into GRU and LSTM
 - Scalability
 - “Context”

x_{t-1} x_t

The quick brown fox jumps over the lazy dog.



“Attention is all you need”

Ashish Aswani et al. 2017
Machine Translation Paper

Natural Language Processing



Computer Vision



Generative Models



Translation



Reinforcement Learning

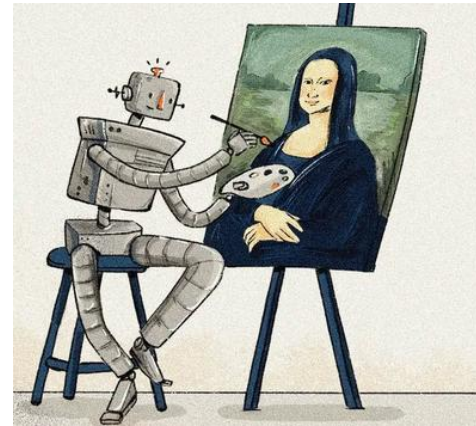
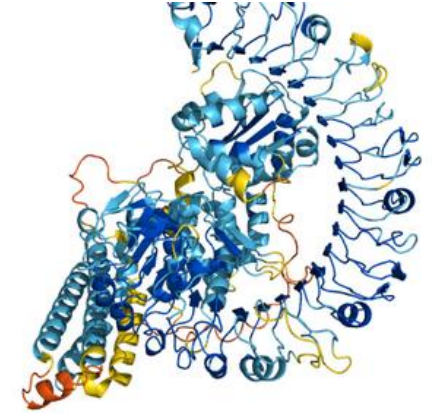


Speech Recognition

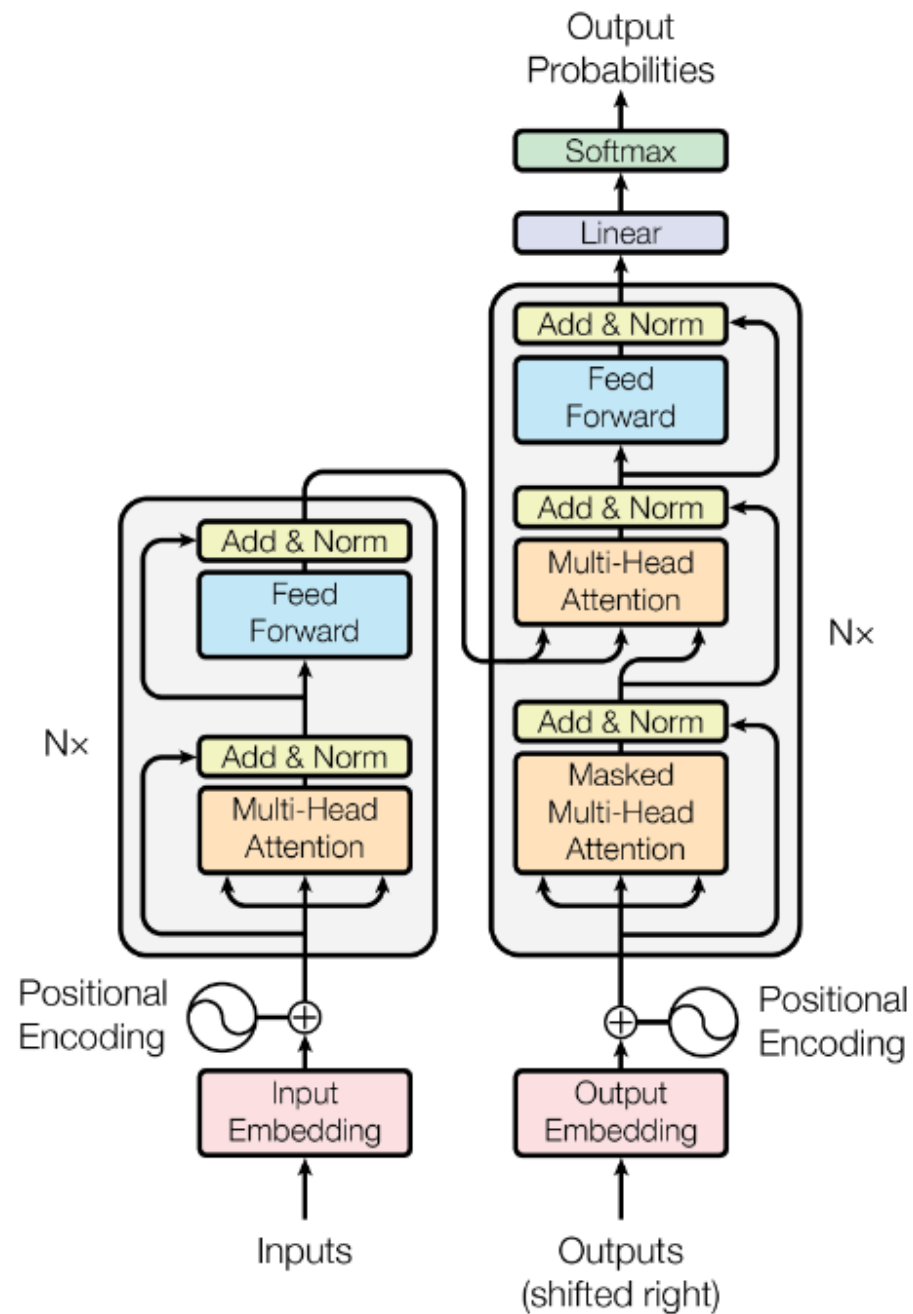


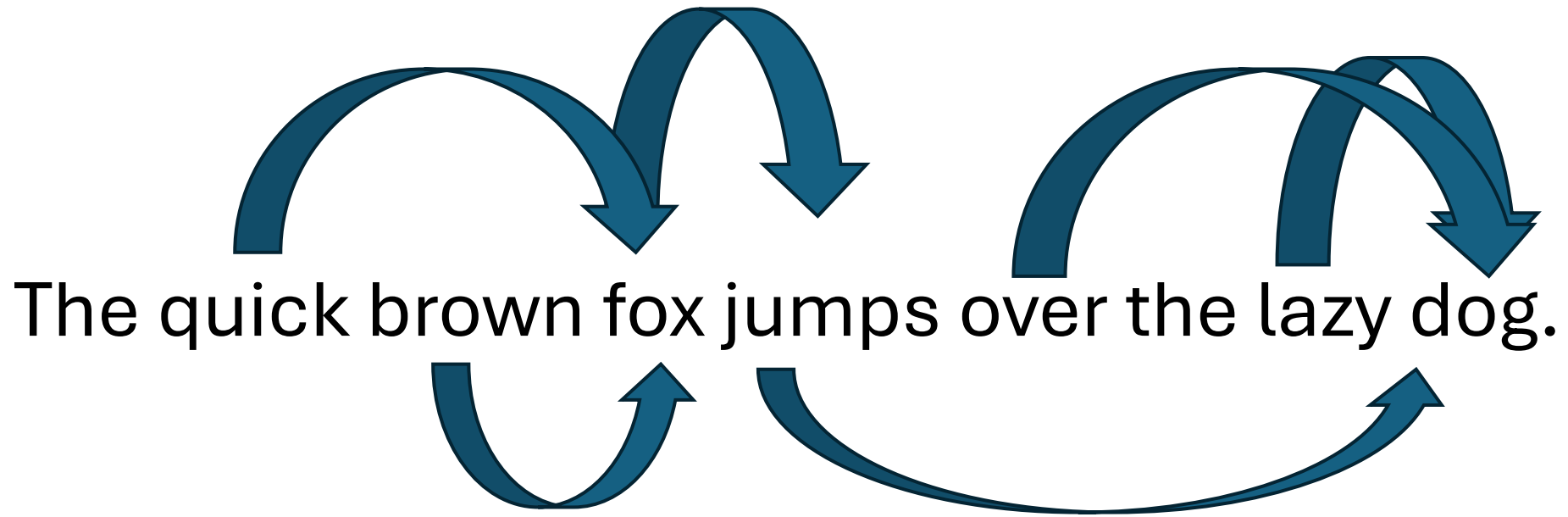
Evolution of Transformer Architecture

- Conversational AI and Agents
 - ChatGPT, Claude,
- Image Generation Models
 - DALL-E, SORA
- Healthcare
 - Alpha Fold
- Coding
 - Github Copilot
- Computer vision
 - Tesla Optimus



Transformer Architecture





The quick brown **fox**

The quick brown fox **jumps**

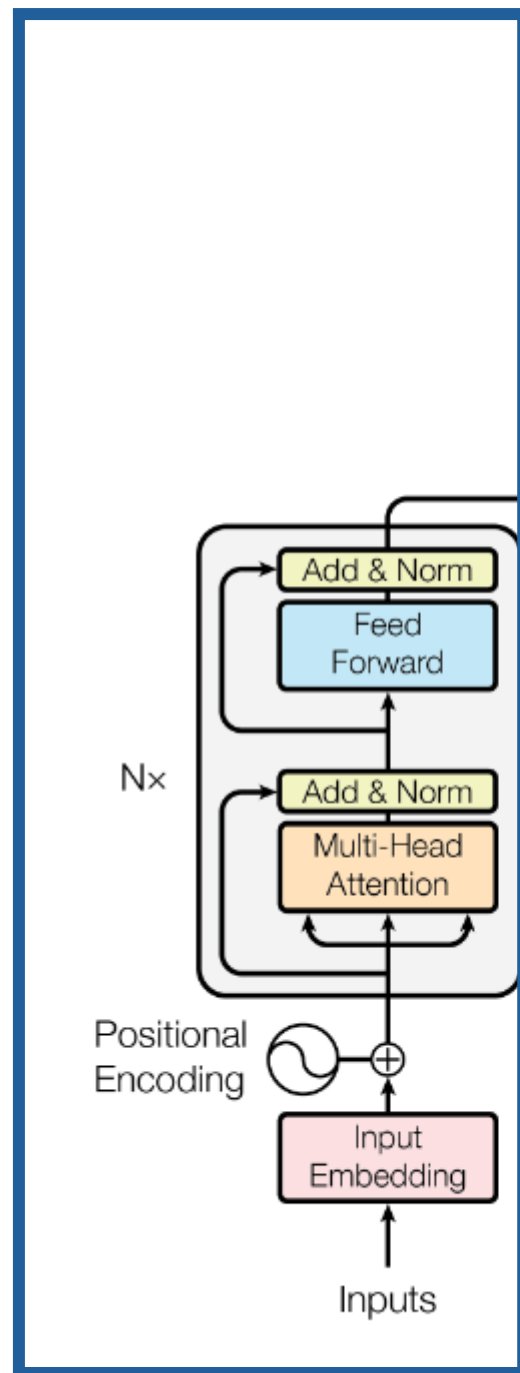
The quick brown fox jumps **over**

The quick brown fox jumps over **the**

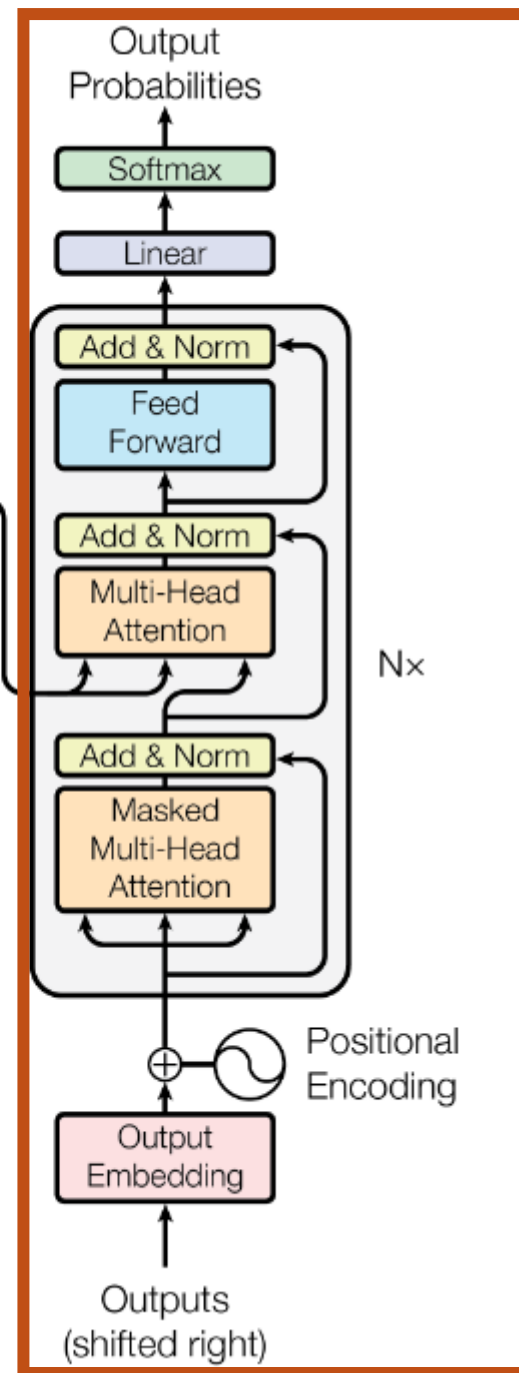
The quick brown fox jumps over the **lazy**

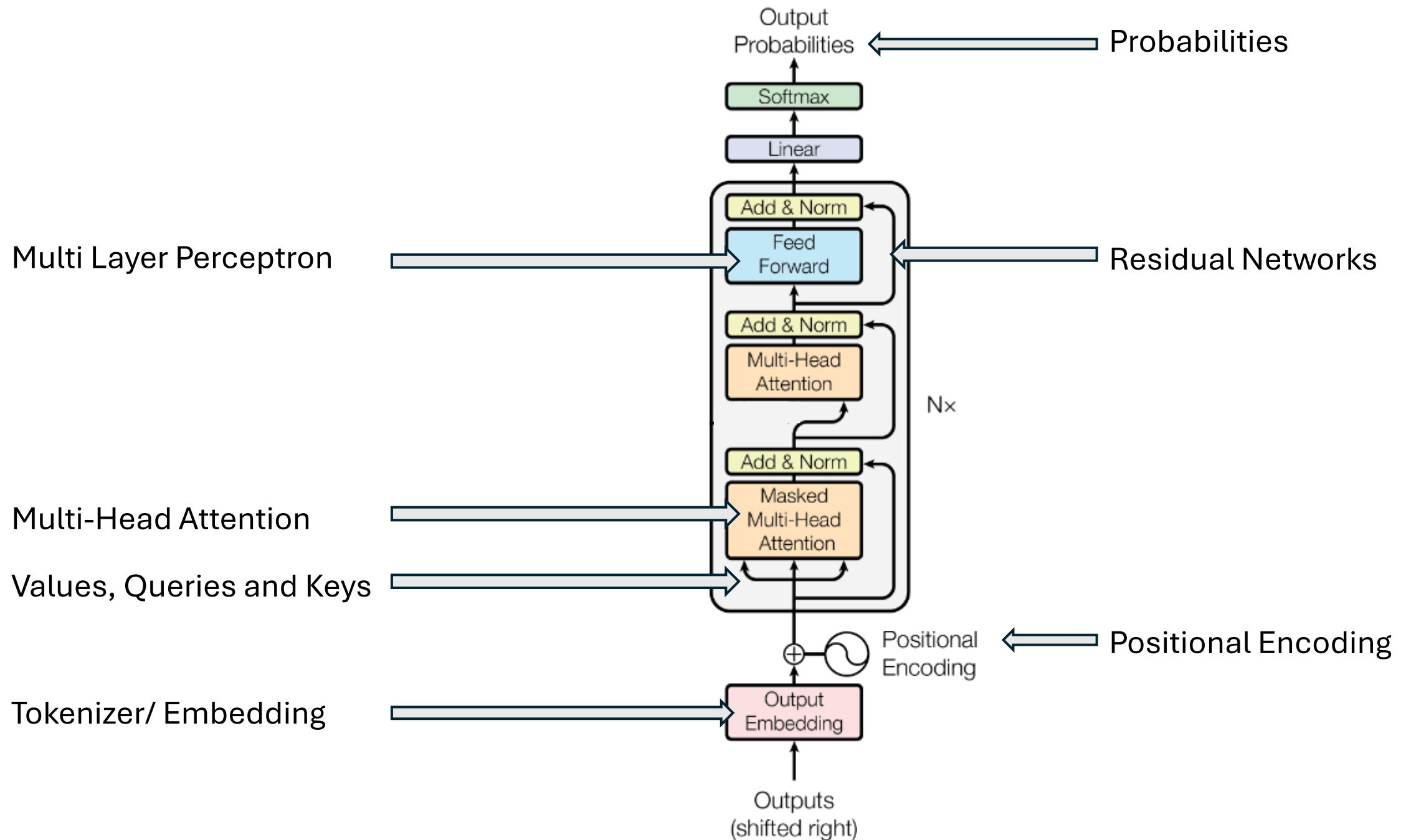
The quick brown fox jumps over the lazy **dog**

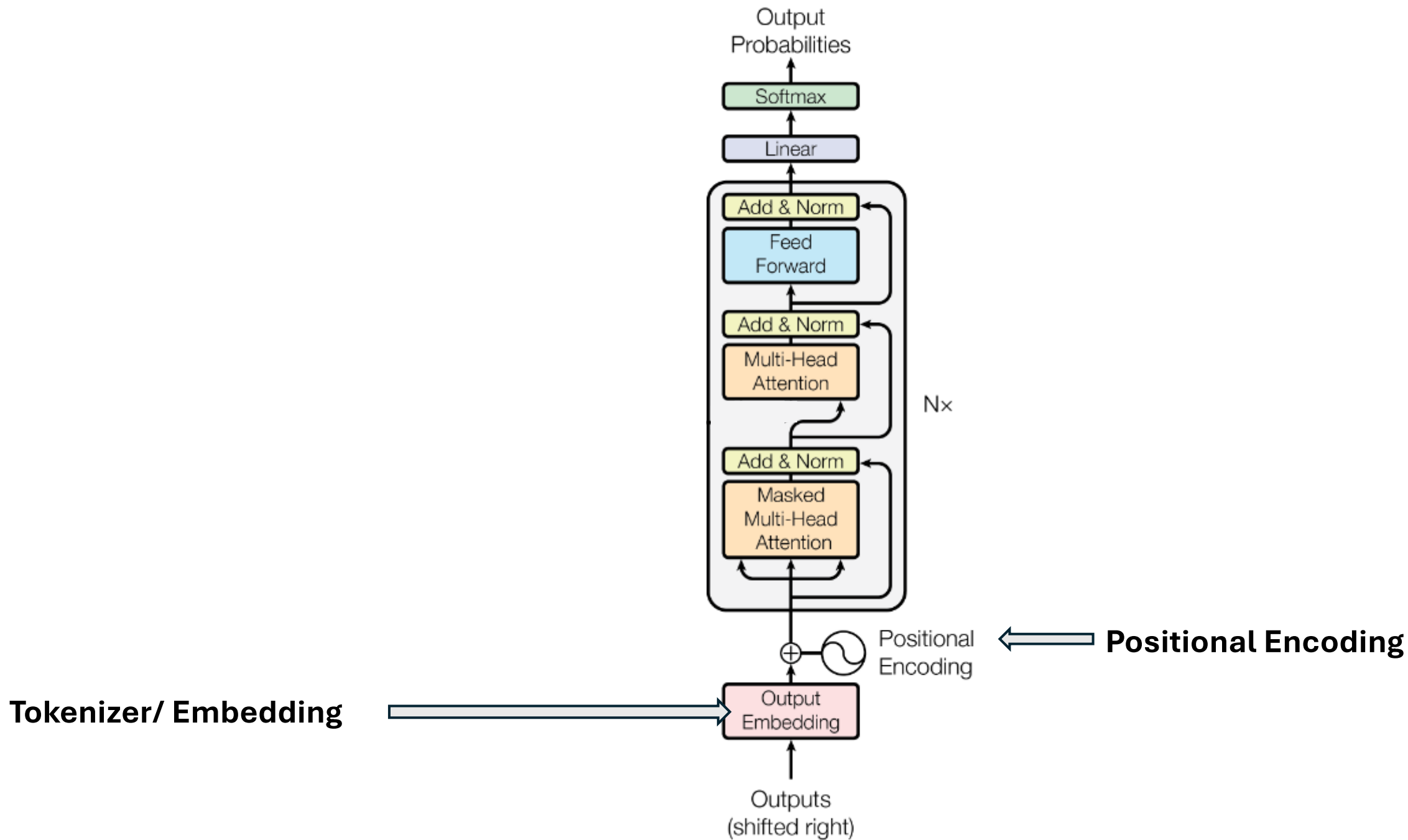
Encoder



Decoder

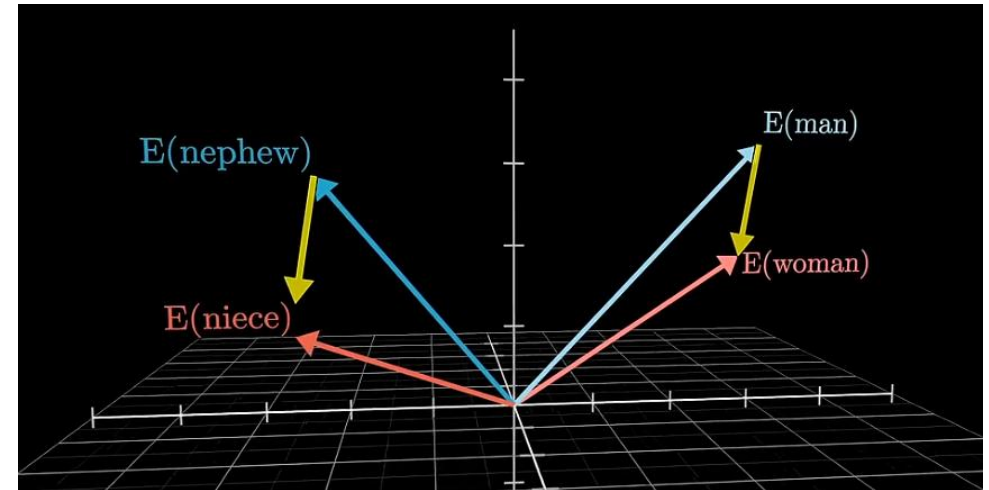




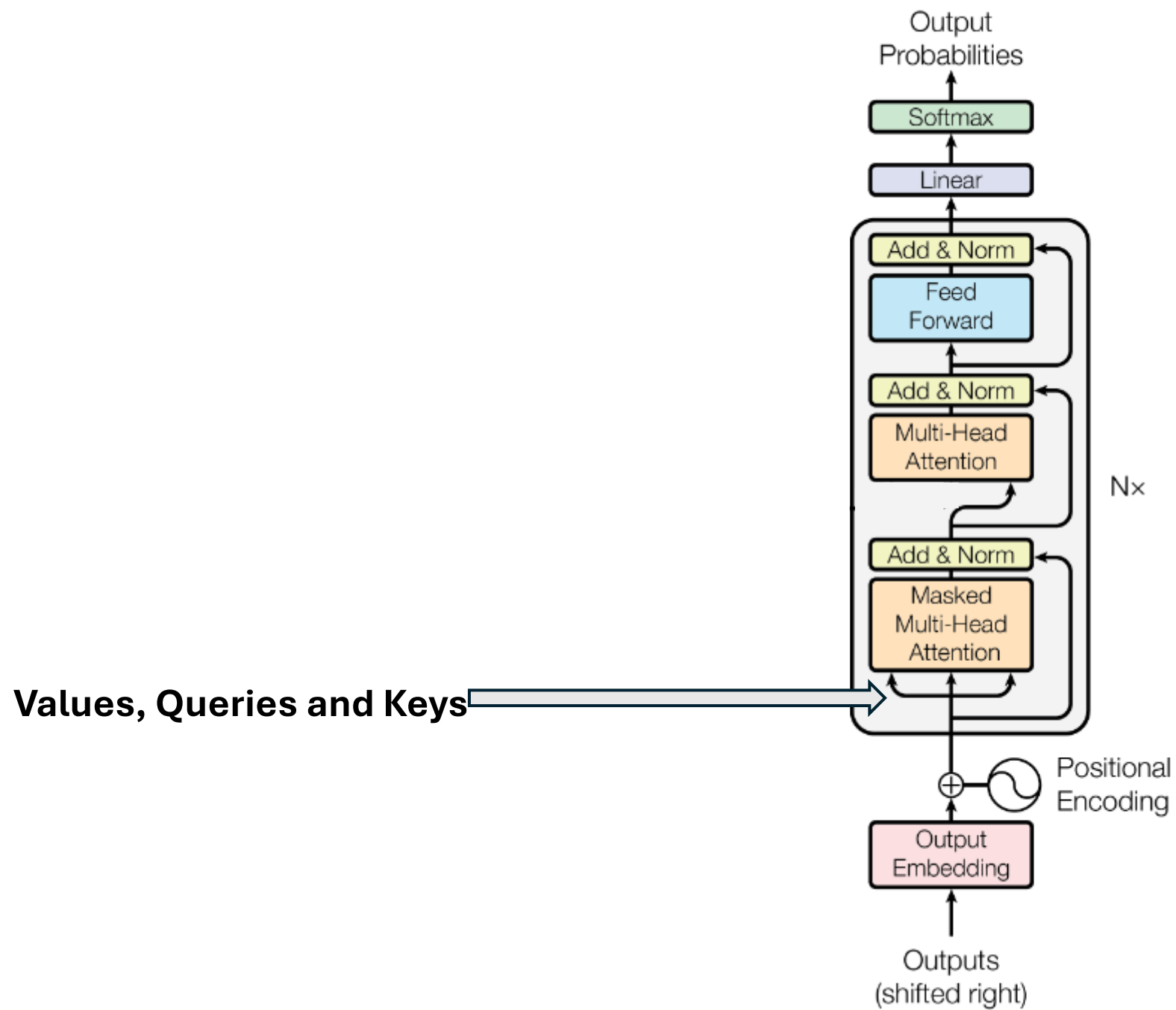


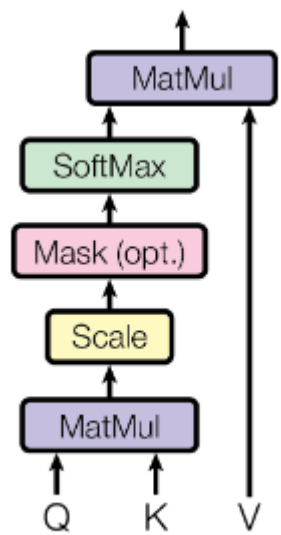
The quick brown fox jumps over the lazy dog

0.4571
0.9855
1.3480
-80.230
-0.2392



*GPT3 Model embeddings has 12,288 dimensions





Query: Question to be asked on the Data

Key: Answer Position of the Data

Dot Product

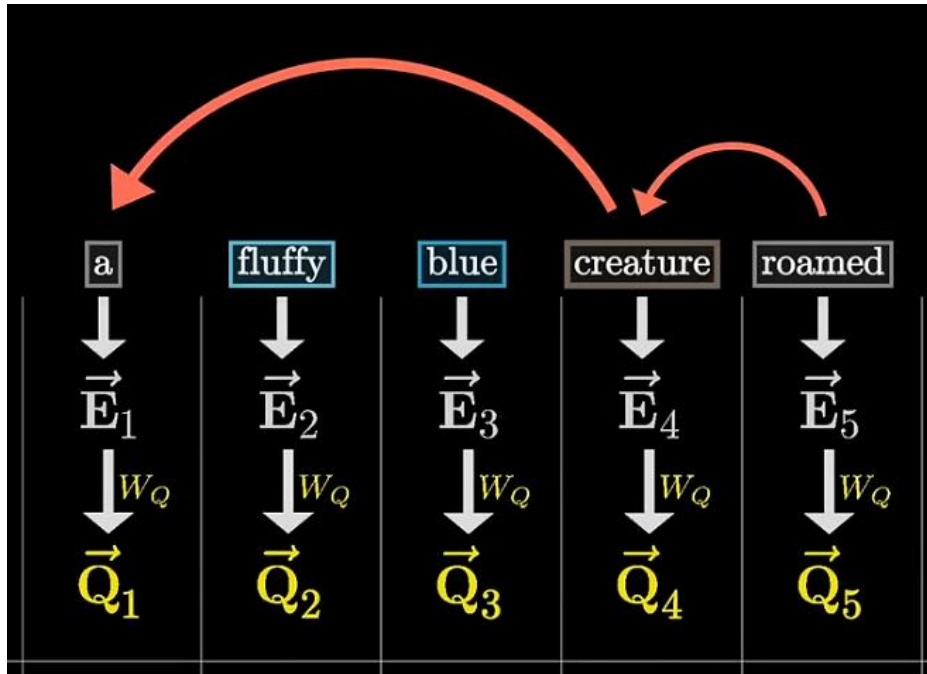
Value: Importance/relevance between Key & Query for each word: “Attention”

Softmax(Q, K, V)

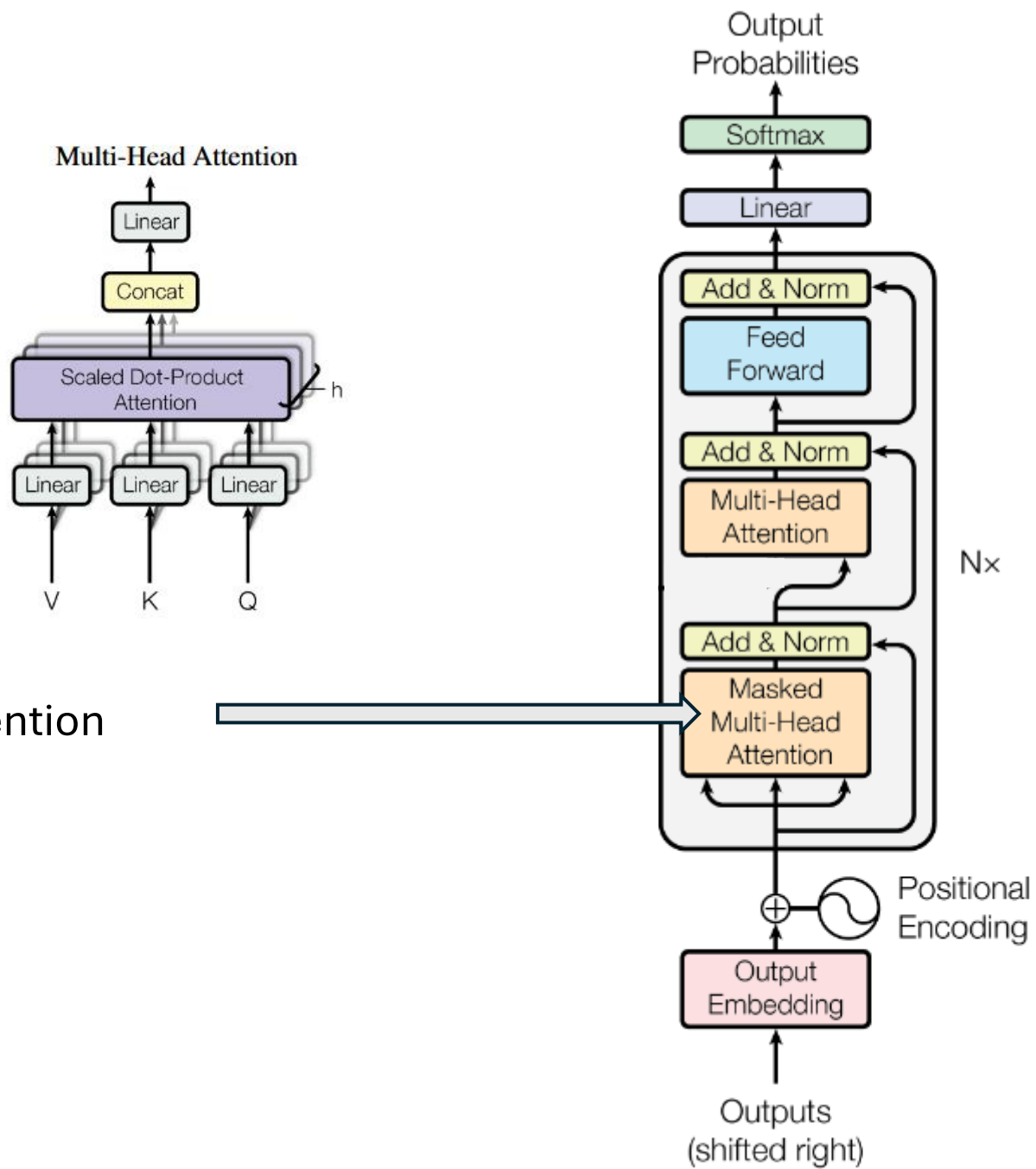
Probability Distribution of possible words

	a	fluffy	blue	creature	roamed	the	verdant
	\vec{E}_1	\vec{E}_2	\vec{E}_3	\vec{E}_4	\vec{E}_5	\vec{E}_6	\vec{E}_7
	\vec{Q}_1	\vec{Q}_2	\vec{Q}_3	\vec{Q}_4	\vec{Q}_5	\vec{Q}_6	\vec{Q}_7
a $\rightarrow \vec{E}_1 \xrightarrow{w_k} \vec{K}_1$	$\vec{K}_1 \cdot \vec{Q}_1$	$\vec{K}_1 \cdot \vec{Q}_2$	$\vec{K}_1 \cdot \vec{Q}_3$	$\vec{K}_1 \cdot \vec{Q}_4$	$\vec{K}_1 \cdot \vec{Q}_5$	$\vec{K}_1 \cdot \vec{Q}_6$	$\vec{K}_1 \cdot \vec{Q}_7$
fluffy $\rightarrow \vec{E}_2 \xrightarrow{w_k} \vec{K}_2$	$\vec{K}_2 \cdot \vec{Q}_1$	$\vec{K}_2 \cdot \vec{Q}_2$	$\vec{K}_2 \cdot \vec{Q}_3$	$\vec{K}_2 \cdot \vec{Q}_4$	$\vec{K}_2 \cdot \vec{Q}_5$	$\vec{K}_2 \cdot \vec{Q}_6$	$\vec{K}_2 \cdot \vec{Q}_7$
blue $\rightarrow \vec{E}_3 \xrightarrow{w_k} \vec{K}_3$	$\vec{K}_3 \cdot \vec{Q}_1$	$\vec{K}_3 \cdot \vec{Q}_2$	$\vec{K}_3 \cdot \vec{Q}_3$	$\vec{K}_3 \cdot \vec{Q}_4$	$\vec{K}_3 \cdot \vec{Q}_5$	$\vec{K}_3 \cdot \vec{Q}_6$	$\vec{K}_3 \cdot \vec{Q}_7$
creature $\rightarrow \vec{E}_4 \xrightarrow{w_k} \vec{K}_4$	$\vec{K}_4 \cdot \vec{Q}_1$	$\vec{K}_4 \cdot \vec{Q}_2$	$\vec{K}_4 \cdot \vec{Q}_3$	$\vec{K}_4 \cdot \vec{Q}_4$	$\vec{K}_4 \cdot \vec{Q}_5$	$\vec{K}_4 \cdot \vec{Q}_6$	$\vec{K}_4 \cdot \vec{Q}_7$
roamed $\rightarrow \vec{E}_5 \xrightarrow{w_k} \vec{K}_5$	$\vec{K}_5 \cdot \vec{Q}_1$	$\vec{K}_5 \cdot \vec{Q}_2$	$\vec{K}_5 \cdot \vec{Q}_3$	$\vec{K}_5 \cdot \vec{Q}_4$	$\vec{K}_5 \cdot \vec{Q}_5$	$\vec{K}_5 \cdot \vec{Q}_6$	$\vec{K}_5 \cdot \vec{Q}_7$
the $\rightarrow \vec{E}_6 \xrightarrow{w_k} \vec{K}_6$	$\vec{K}_6 \cdot \vec{Q}_1$	$\vec{K}_6 \cdot \vec{Q}_2$	$\vec{K}_6 \cdot \vec{Q}_3$	$\vec{K}_6 \cdot \vec{Q}_4$	$\vec{K}_6 \cdot \vec{Q}_5$	$\vec{K}_6 \cdot \vec{Q}_6$	$\vec{K}_6 \cdot \vec{Q}_7$
verdant $\rightarrow \vec{E}_7 \xrightarrow{w_k} \vec{K}_7$	$\vec{K}_7 \cdot \vec{Q}_1$	$\vec{K}_7 \cdot \vec{Q}_2$	$\vec{K}_7 \cdot \vec{Q}_3$	$\vec{K}_7 \cdot \vec{Q}_4$	$\vec{K}_7 \cdot \vec{Q}_5$	$\vec{K}_7 \cdot \vec{Q}_6$	$\vec{K}_7 \cdot \vec{Q}_7$

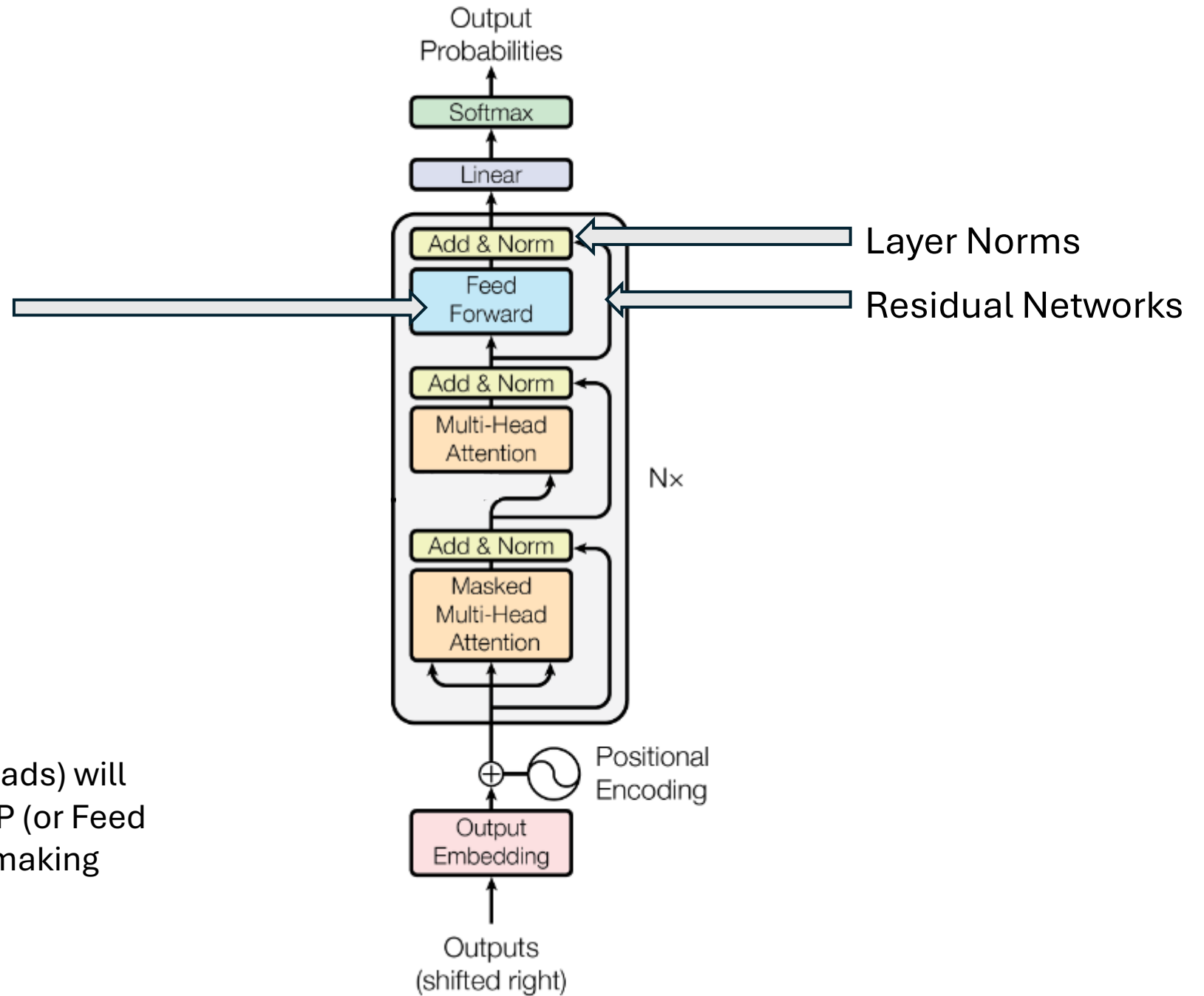
Masking: Following Words Cannot influence previous words



Unnormalized Attention Pattern						Normalized Attention Pattern						
+3.53	+0.80	+1.96	+4.48	+3.74	-1.95	softmax →	1.00	0.75	0.69	0.92	0.46	0.00
−∞	−0.30	−0.21	+0.82	+0.29	+2.91		0.00	0.25	0.08	0.02	0.01	0.46
−∞	−∞	+0.89	+0.67	+2.99	−0.41		0.00	0.00	0.24	0.02	0.22	0.02
−∞	−∞	−∞	+1.31	+1.73	−1.48		0.00	0.00	0.00	0.04	0.06	0.01
−∞	−∞	−∞	−∞	+3.07	+2.94		0.00	0.00	0.00	0.00	0.24	0.48
−∞	−∞	−∞	−∞	−∞	+0.31		0.00	0.00	0.00	0.00	0.00	0.03

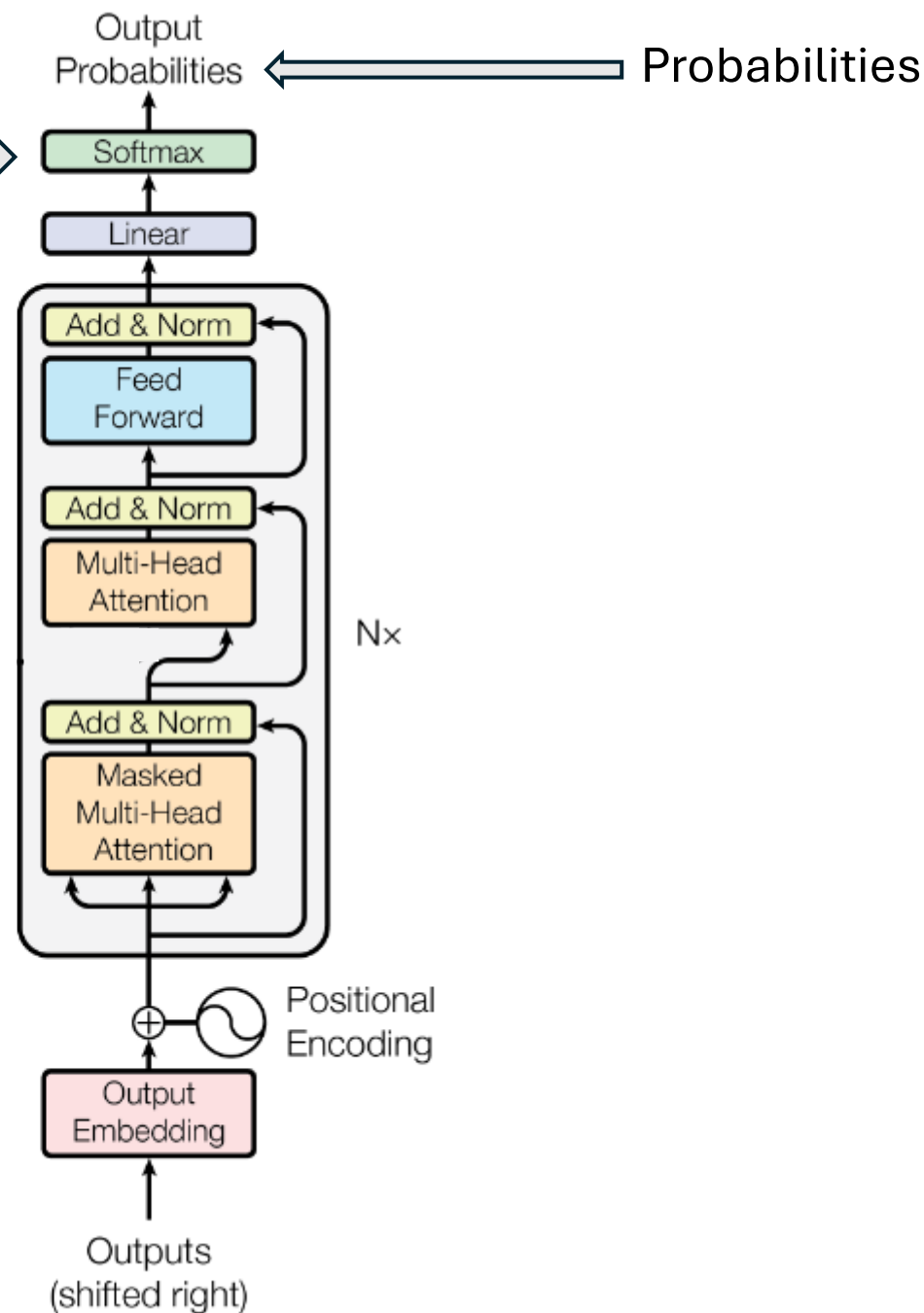
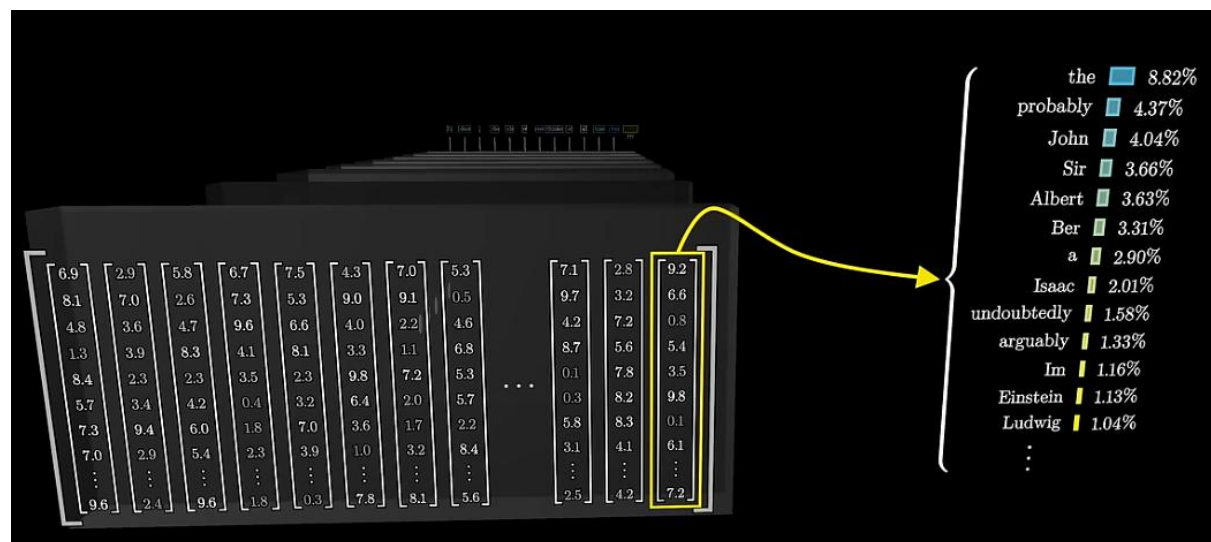


Multi Layer Perceptron

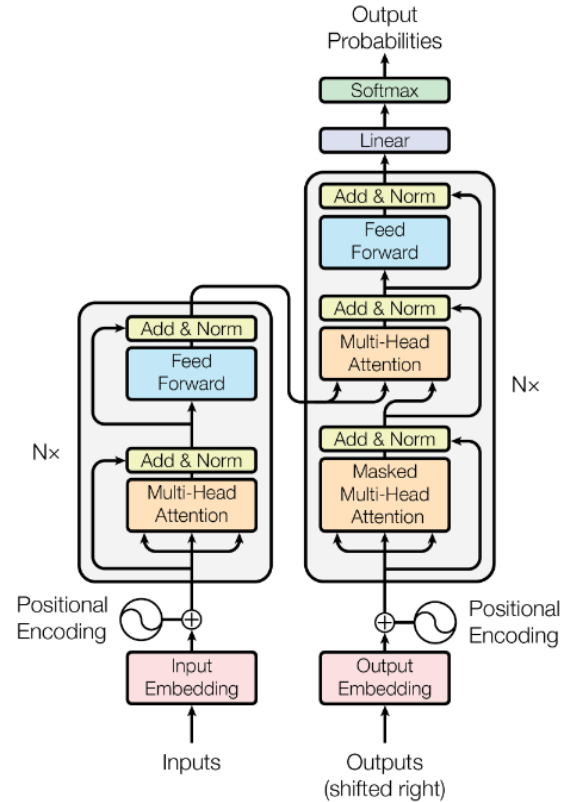


Self Attention blocks (or multi-heads) will gather the data, whereas the MLP (or Feed Forward networks), will work on making sense of the data

- Softmax
 - Rigid
- Temperature Softmax
 - Creativity !!



Thank you !!



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
lukaszkaizer@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

Reference: Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.