



KESHAV MEMORIAL INSTITUTE OF TECHNOLOGY
AN AUTONOMOUS INSTITUTION - ACCREDITED BY NAAC WITH 'A' GRADE
Narayananaguda, Hyderabad.

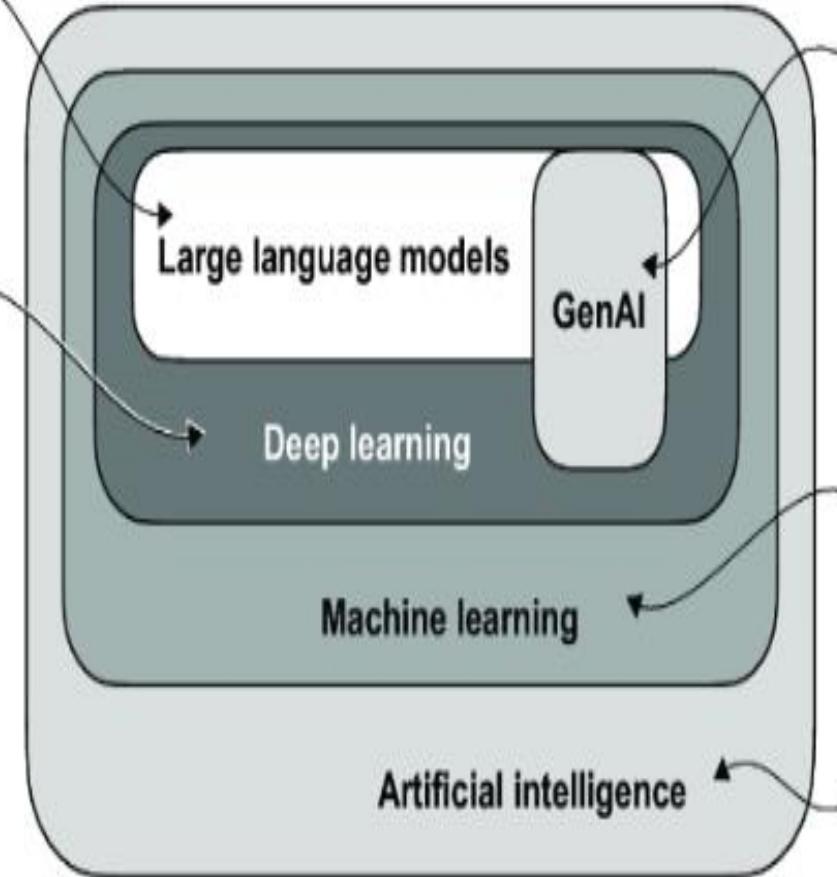
Embedded Learning TRANSFORMERS

03-02-2025

By
Asha.M
Assistant professor
CSE(AI&ML)
KMIT

Deep neural network for parsing and generating human-like text

Machine learning with neural networks consisting of many layers



GenAI involves the use of deep neural networks to create new content, such as text, images, or various forms of media

Algorithms that learn rules automatically from data

Systems with human-like intelligence

RNN



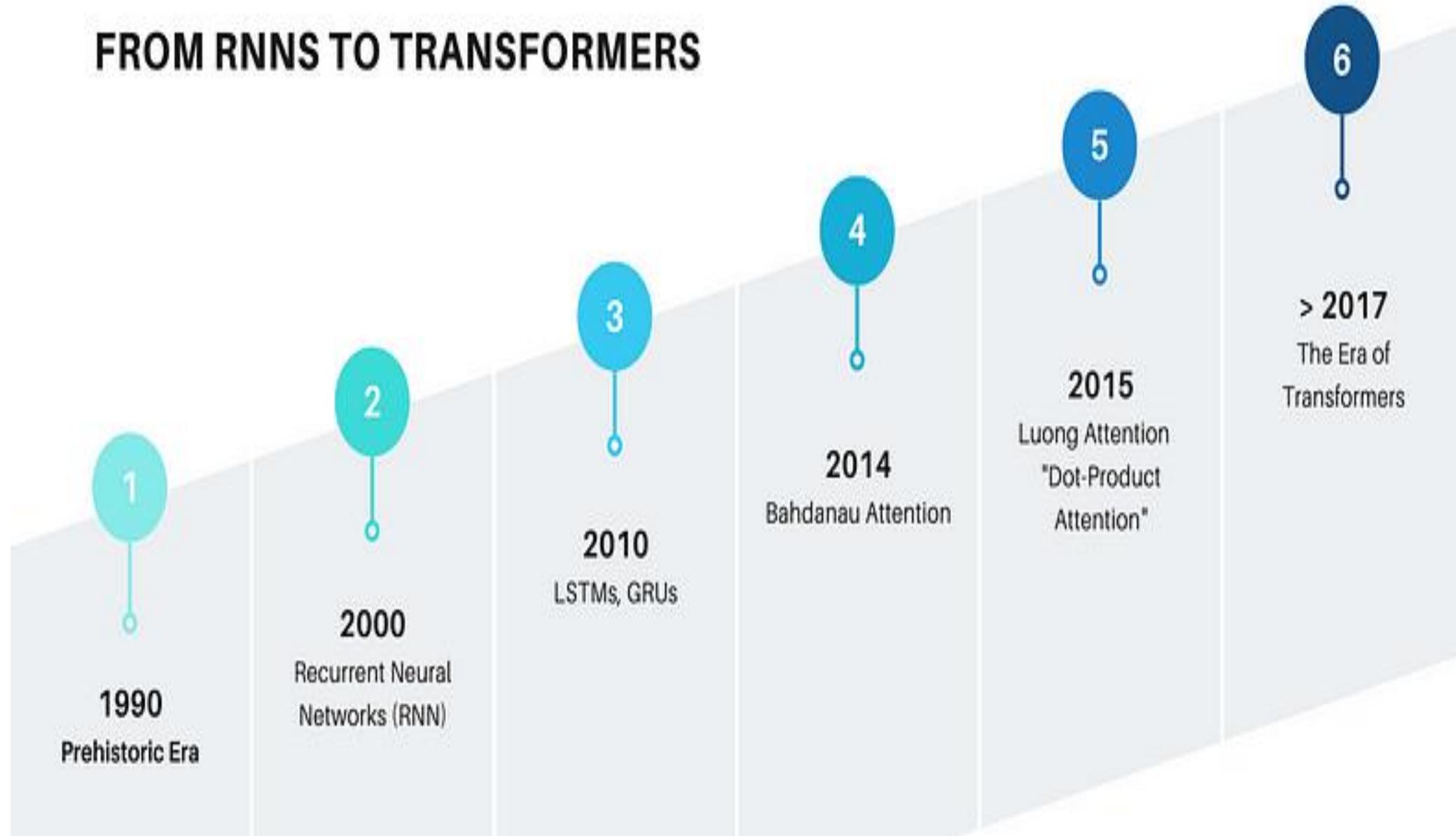
LSTM



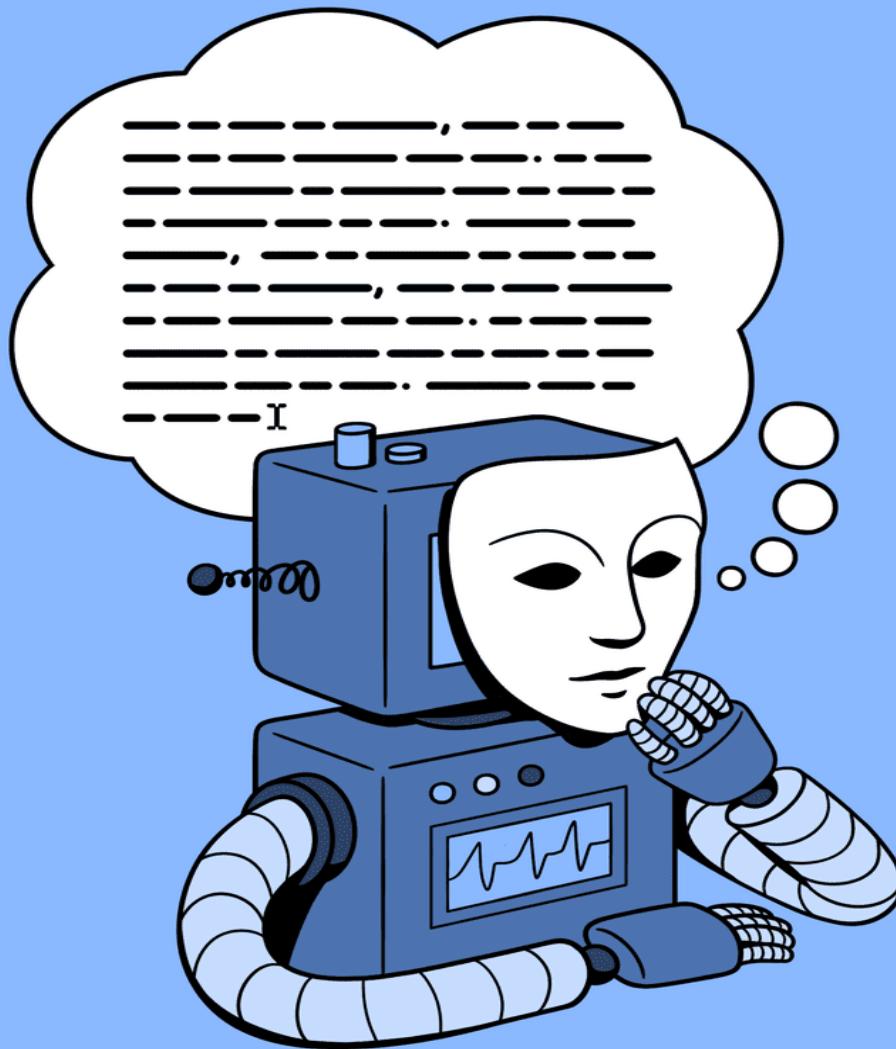
Transformer



FROM RNNs TO TRANSFORMERS







Large Language Model (LLM)

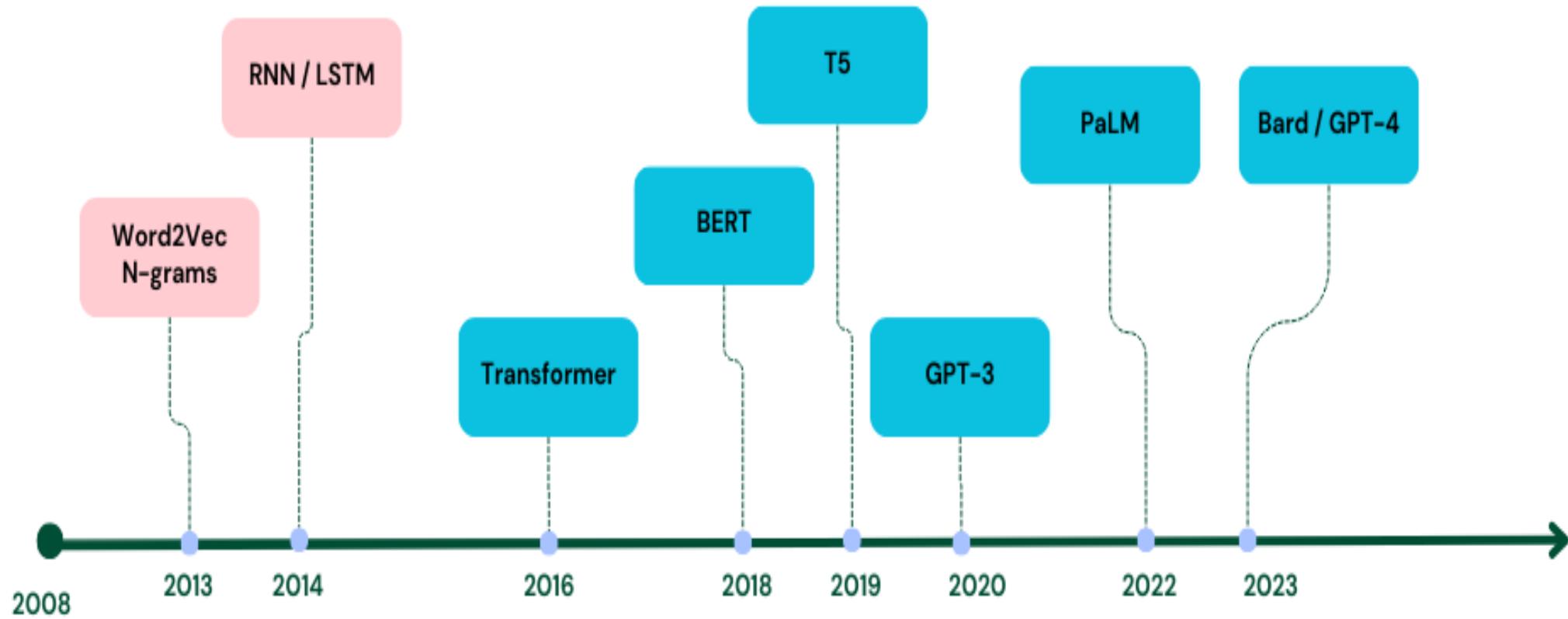
[*'lärj 'laŋ-gwij 'mä-dəl*]

A deep learning algorithm that's equipped to summarize, translate, predict, and generate human-sounding text to convey ideas and concepts.

Language Model History

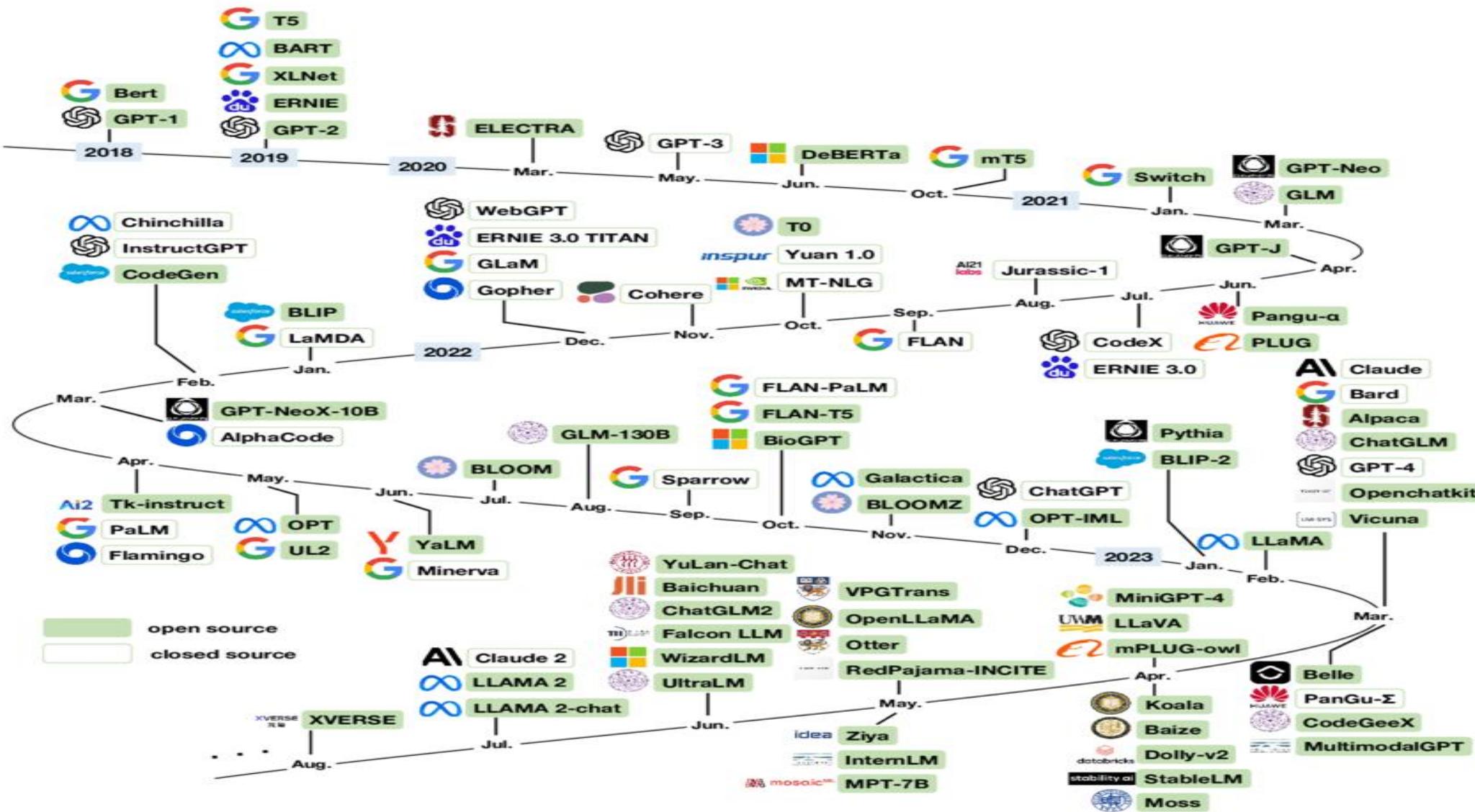
Legend:

- Before Transformer (Pink)
- After Transformer (Blue)

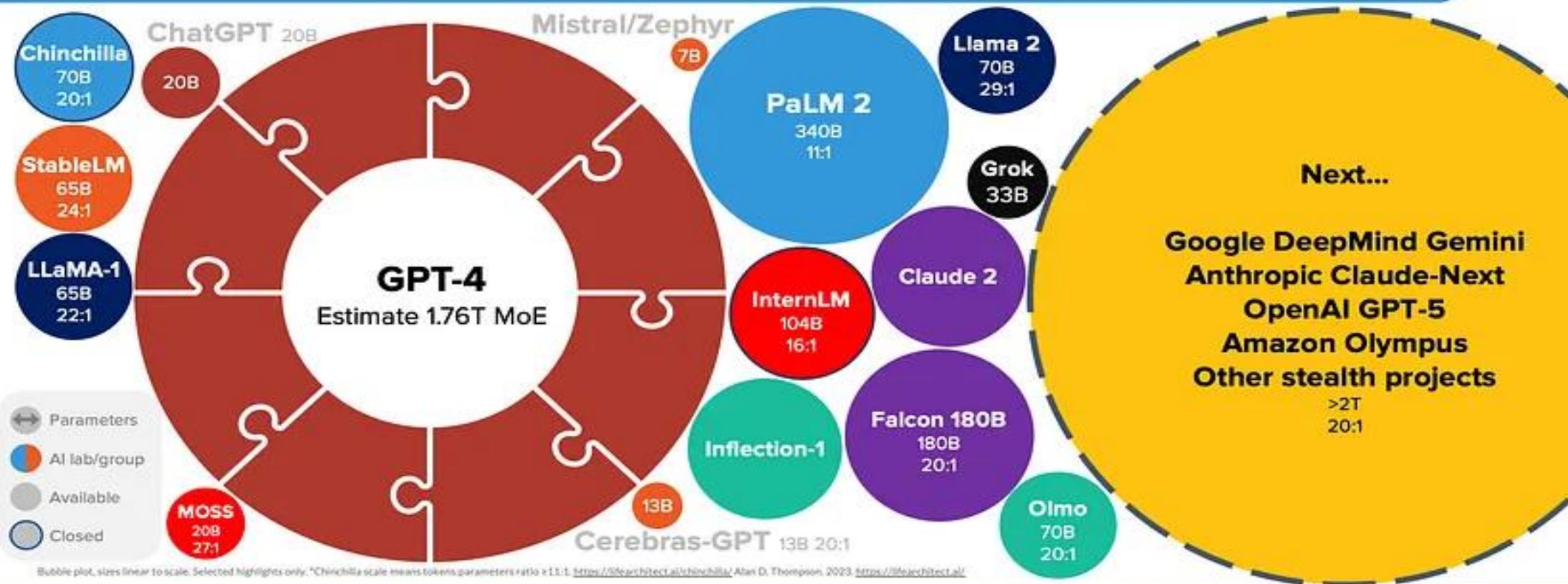


Evolution of Large Language Models



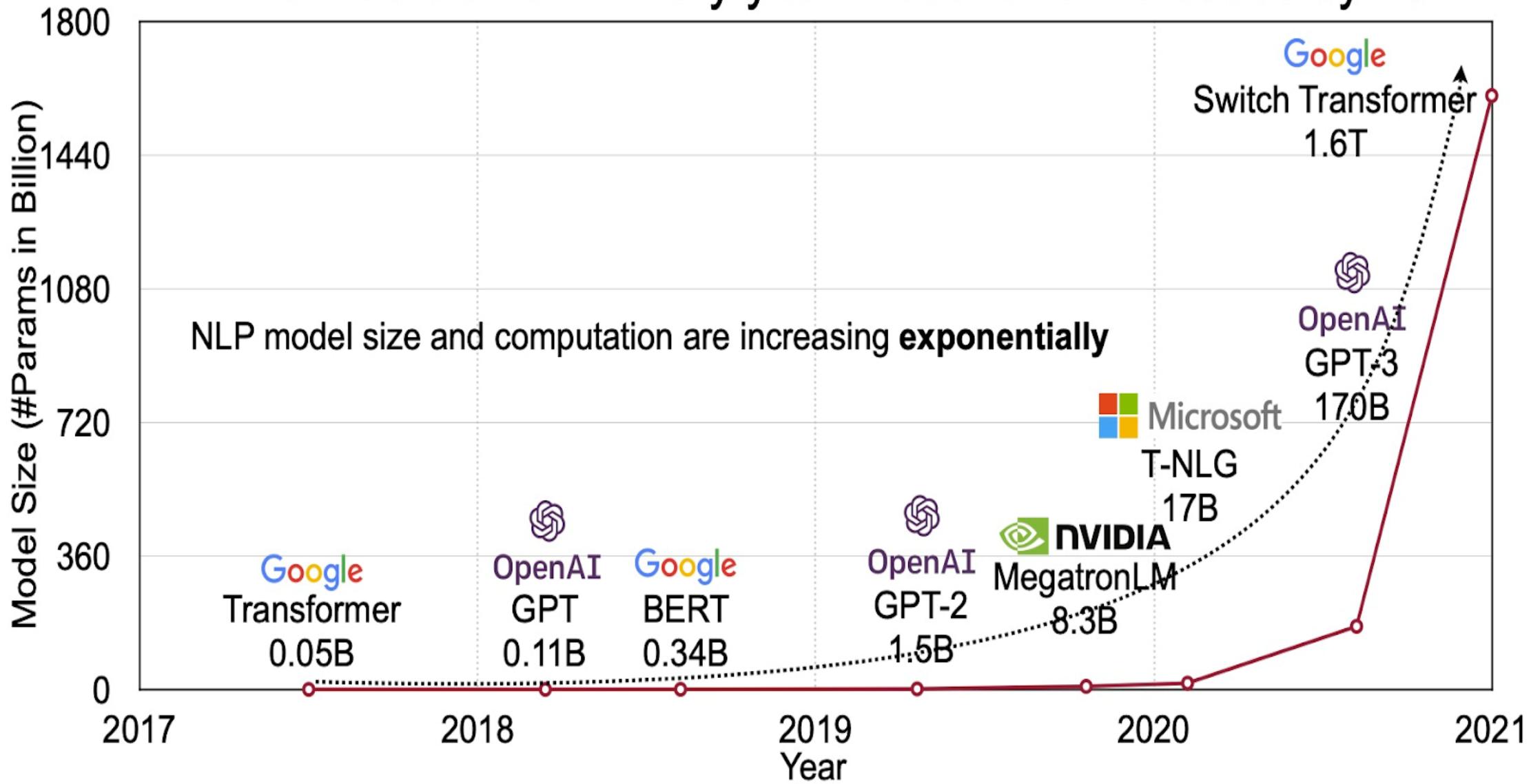


2023-2024 OPTIMAL LANGUAGE MODELS

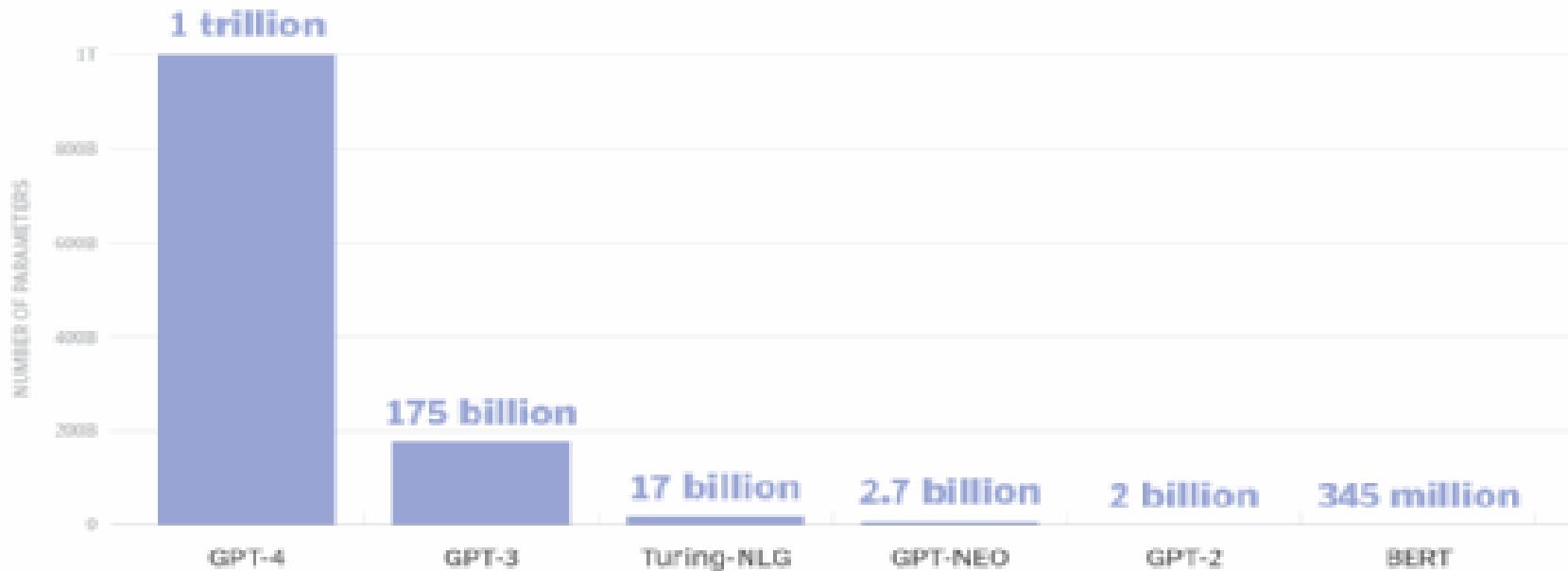
NOV/
2023

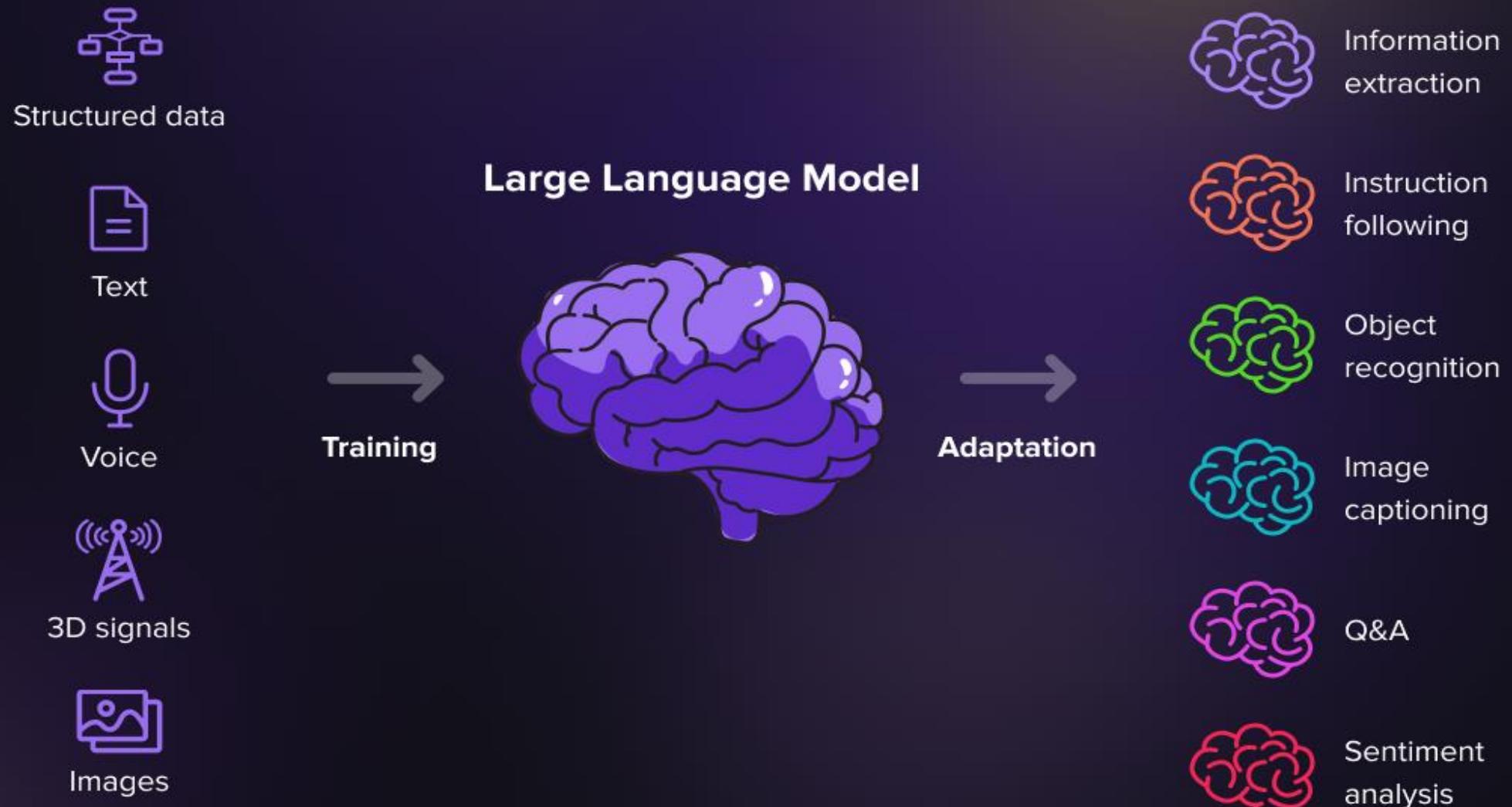
LifeArchitect.ai/models

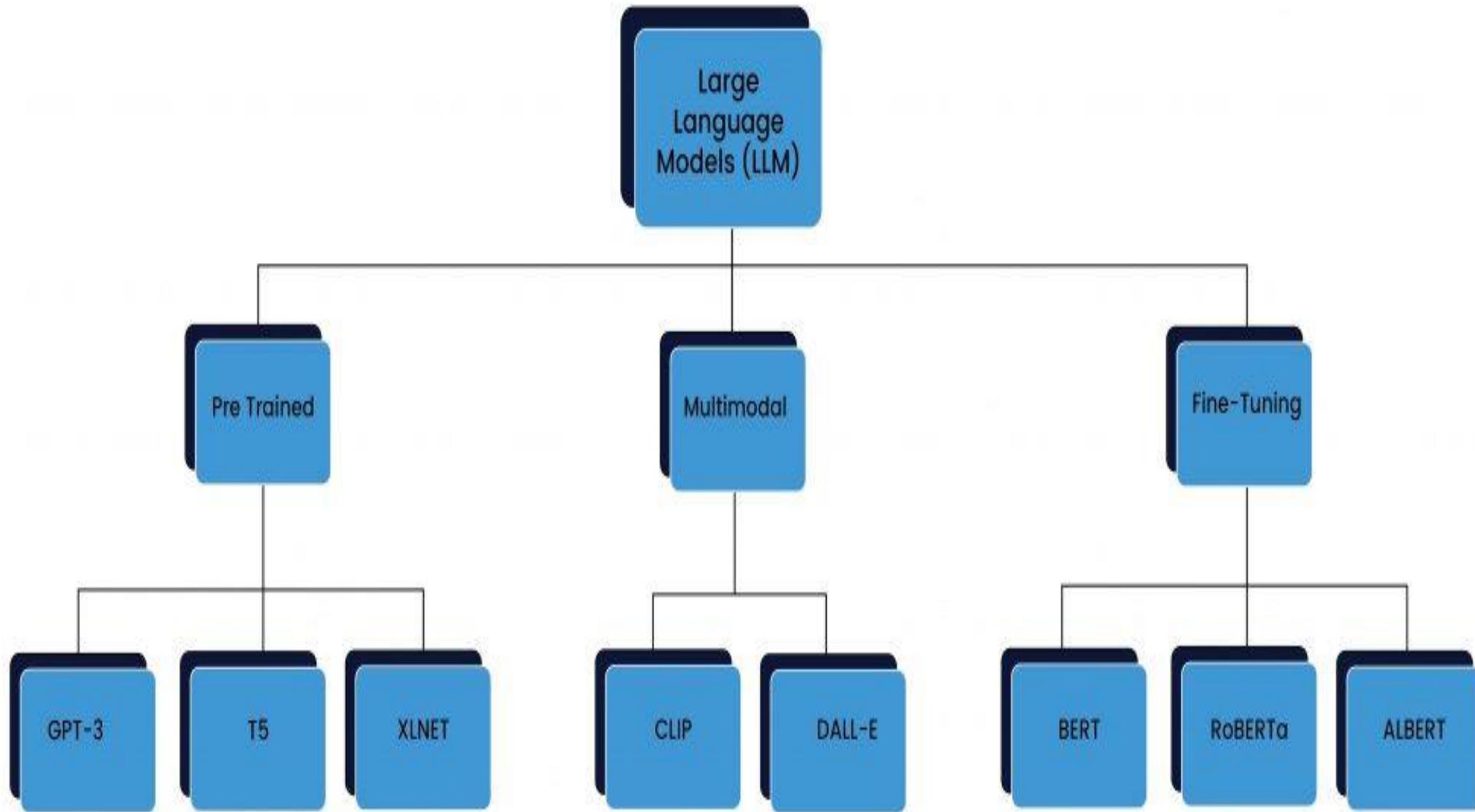
NLP's Moore's Law: Every year model size increases by 10x

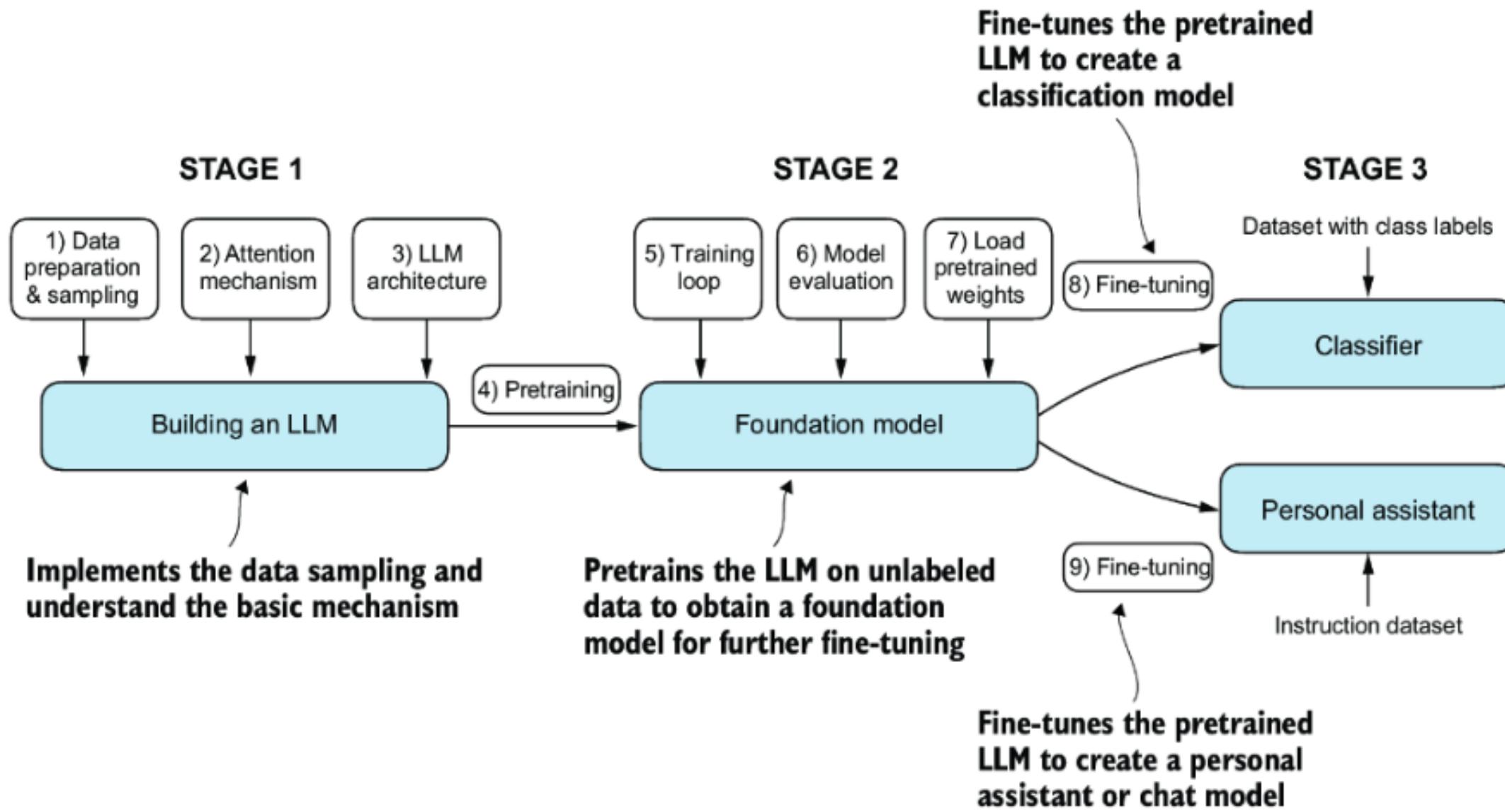


Parameters of transformer-based language models

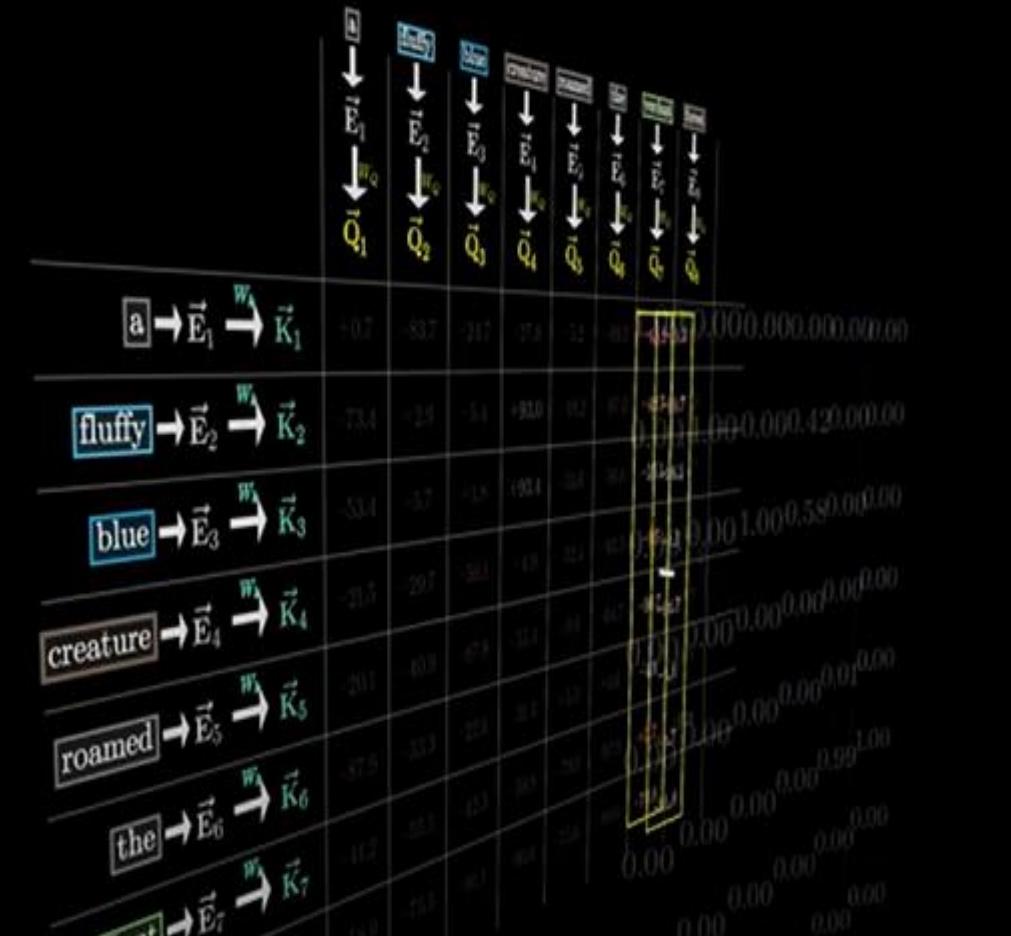
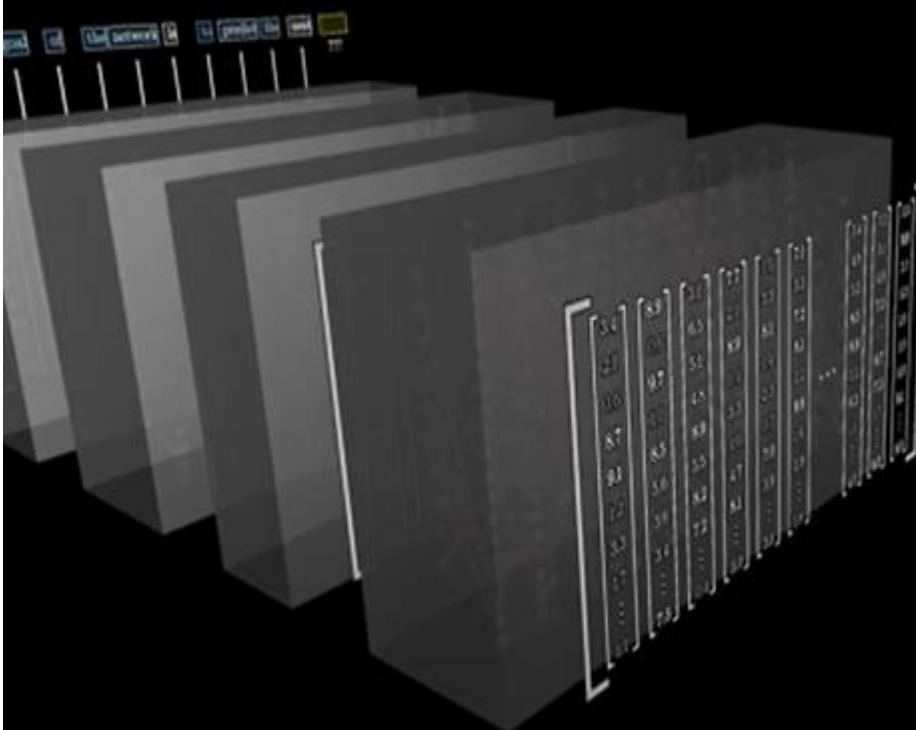








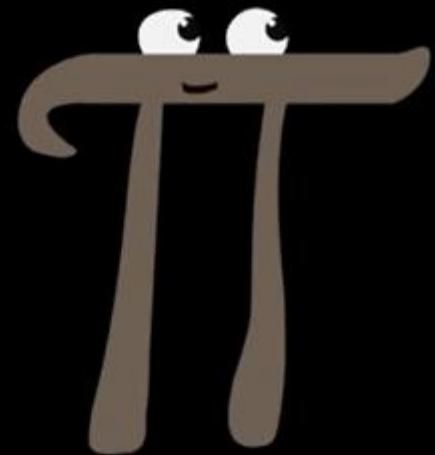
Transformer



What is a Transformer?

A **Transformer** is a deep learning model introduced in "**Attention Is All You Need**" (**Vaswani et al., 2017**). It **replaces recurrent layers (RNNs, LSTMs)** with **self-attention mechanisms**, allowing for **parallel processing** and **long-range dependencies**.

Generative Pre-trained Transformer



Generative

The most effective way to learn computer science is to actively engage with the material, practice regularly, and seek help when needed. Here are some specific steps you can take to improve your

Pre-trained

to the jostlement of all weaker people, might have seen how safe and strong he was. His way taking him past Tellson's, and he both banking at Tellson's and knowing Mr. Lorry as the intimate friend of the Manettes, it entered Mr. Stryver's mind to enter the bank, and reveal

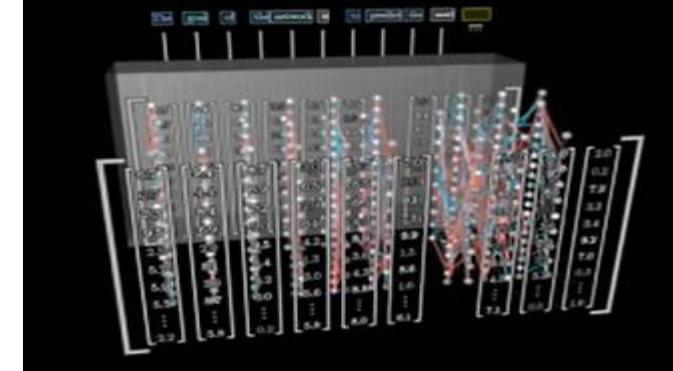


Generative → Can generate human-like text.

Pre-trained → Learns from huge datasets before fine-tuning.

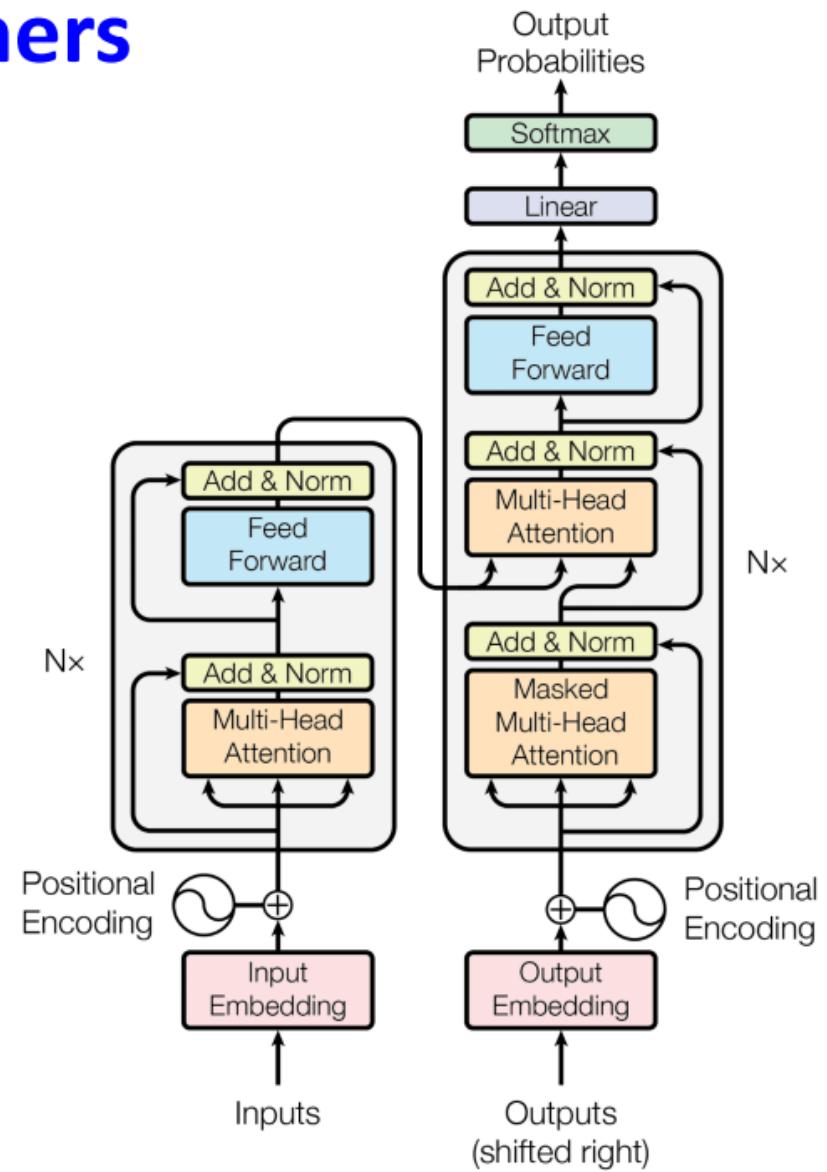
Transformer → Uses self-attention for efficient language understanding.

Transformer



Transformers

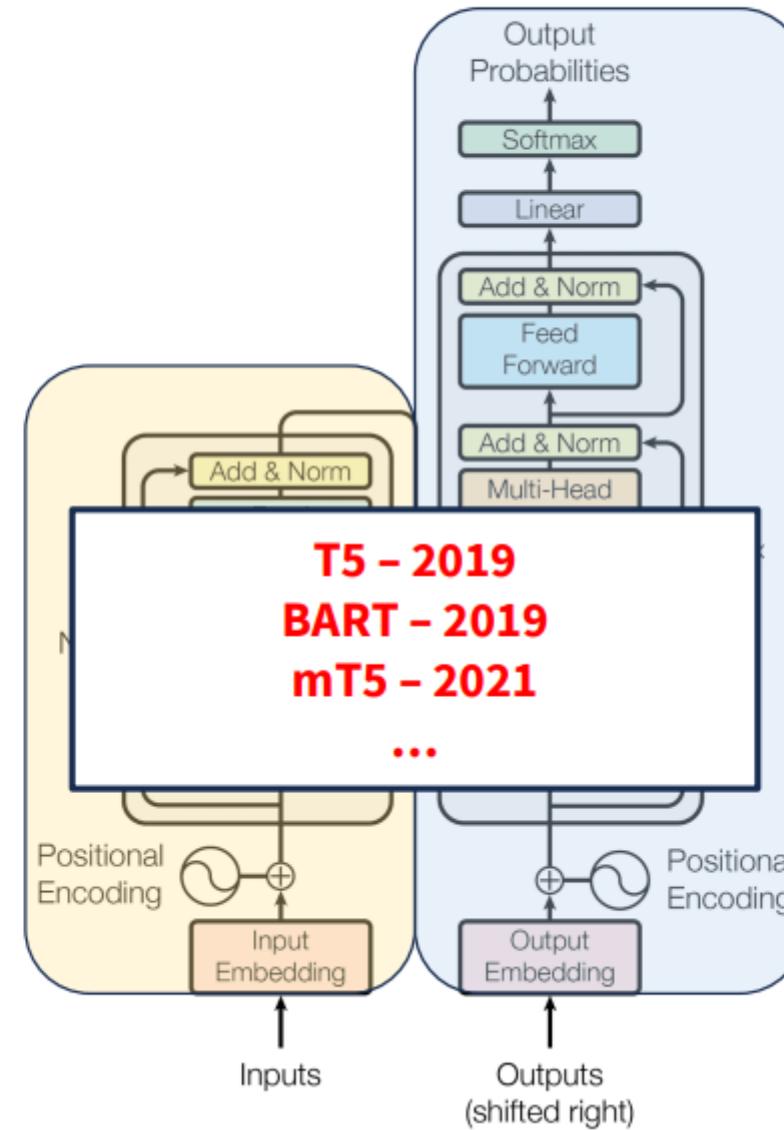
- Tokenization
- Input Embeddings
- Position Encodings
- Residuals
- Query
- Key
- Value
- Add & Norm
- Encoder
- Decoder
- Attention
- Self Attention
- Multi Head Attention
- Masked Attention
- Encoder Decoder Attention
- Output Probabilities / Logits
- Softmax
- Encoder-Decoder models
- Decoder only models



The LLM Era – Paradigm Shift in Machine Learning

BERT – 2018
DistilBERT – 2019
RoBERTa – 2019
ALBERT – 2019
ELECTRA – 2020
DeBERTa – 2020
...

Representation



GPT – 2018
GPT-2 – 2019
GPT-3 – 2020
GPT-Neo – 2021
GPT-3.5 (ChatGPT) – 2022
LLaMA – 2023
GPT-4 – 2023
...

Generation

Transformers

ENCODER ONLY

aka

auto-encoding models

TASKS

- Sentence classification
- Named entity recognition
- Extractive question-answering
- Masked language modeling

EXAMPLES

BERT, RoBERTa, distilBERT

DECODER ONLY

aka

auto-regressive models

TASKS

- Text generation
- Causal language modeling

EXAMPLES

GPT-2, GPT Neo, GPT-3

ENCODER- DECODER

aka

sequence-to- sequence models

TASKS

- Translation
- Summarization
- Generative question-answering

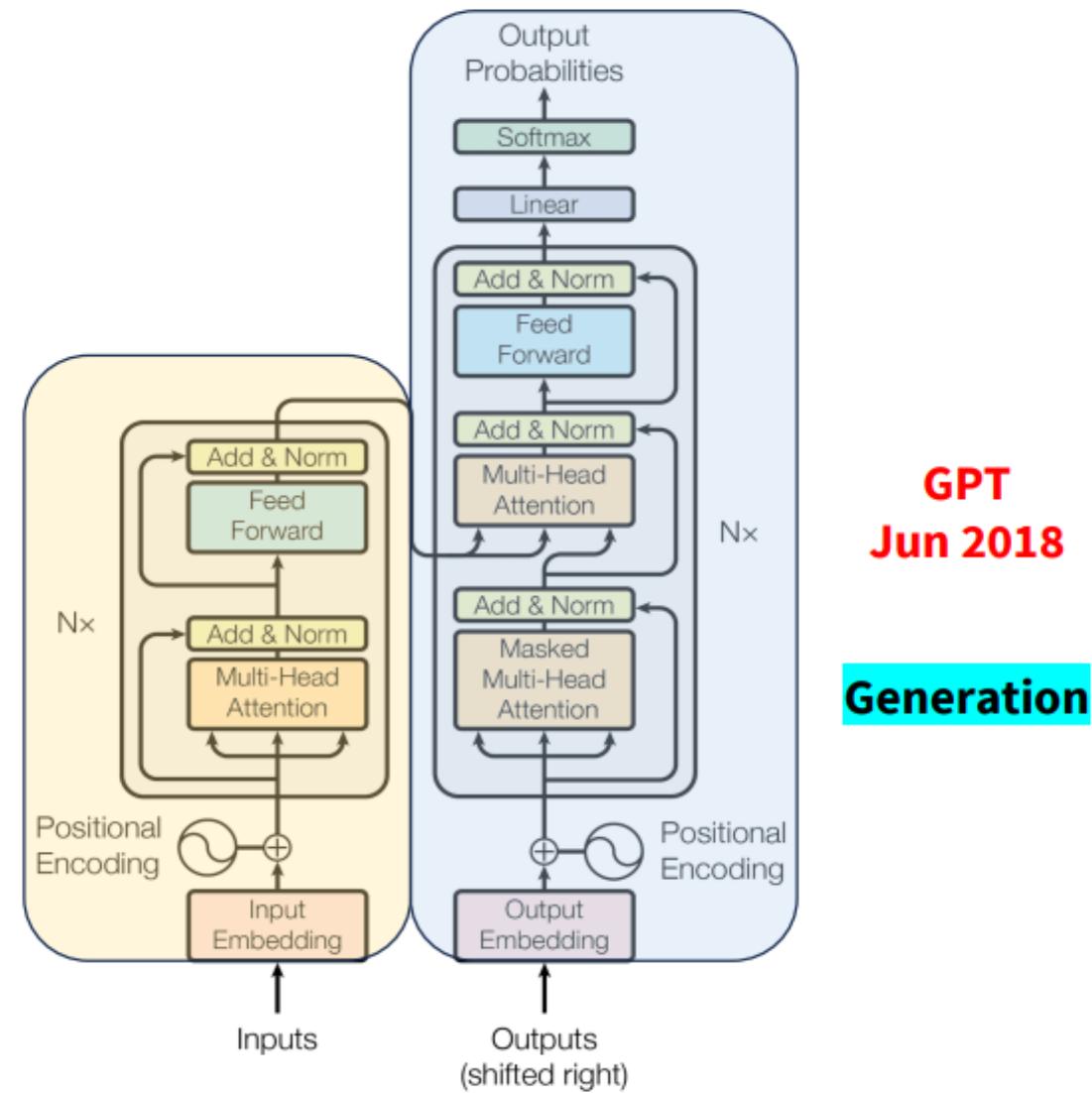
EXAMPLES

BART, T5, Marian

2018 – The Inception of the LLM Era

BERT
Oct 2018

Representation

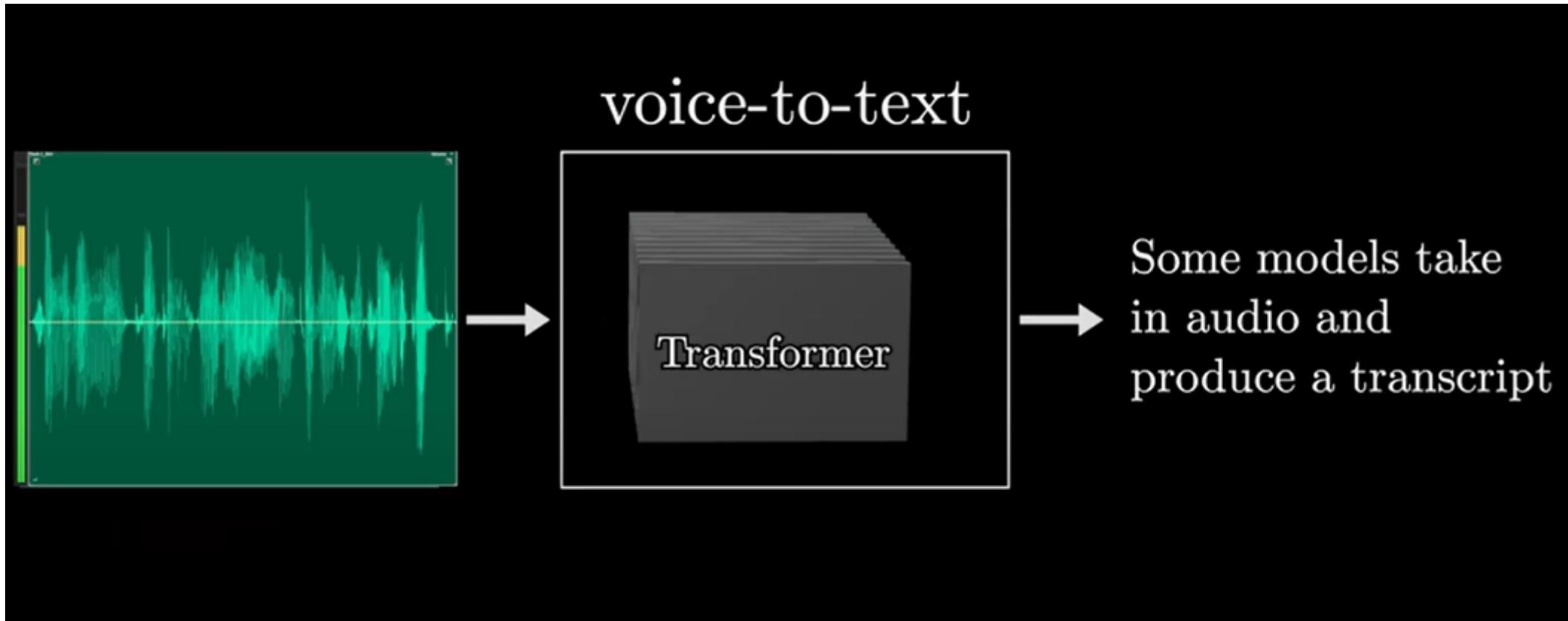


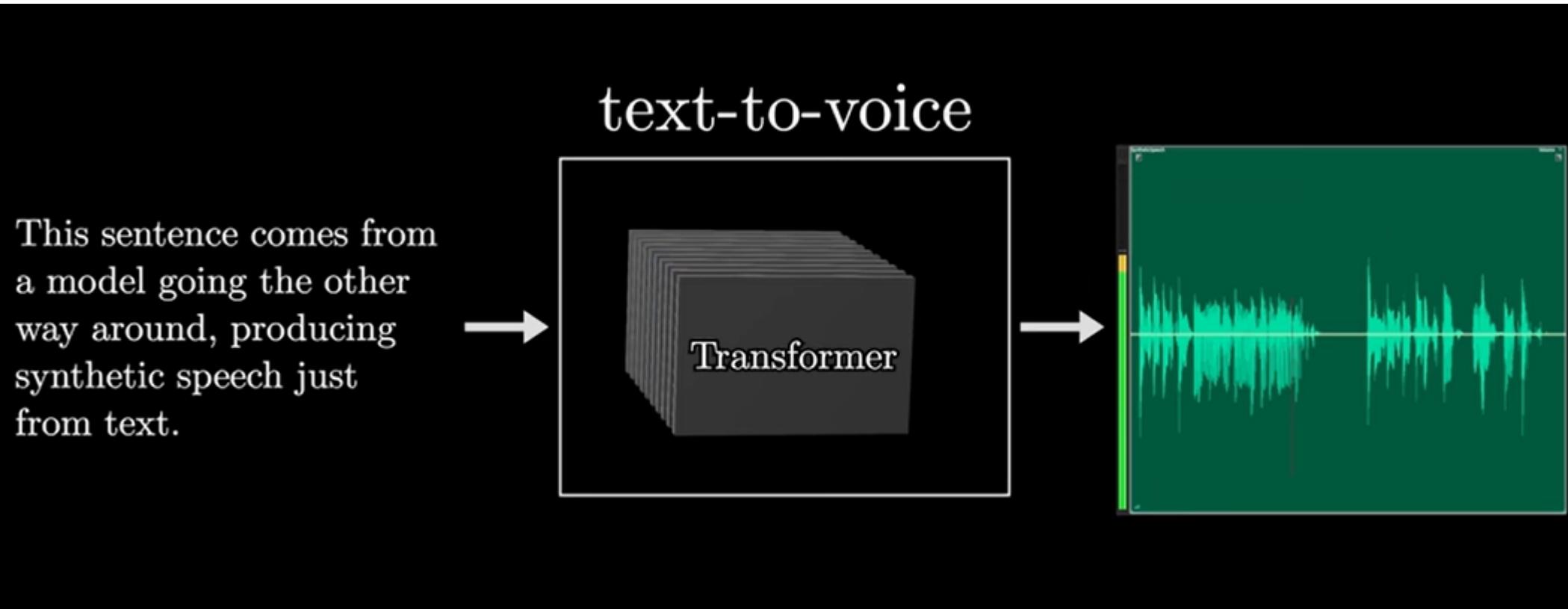
The LLM Era – Paradigm Shift in Machine Learning

- What has caused this paradigm shift?
 - Problem in recurrent networks
 - Information is effectively lost during encoding of long sequences
 - Sequential nature disables parallel training and favors late timestep inputs

The LLM Era – Paradigm Shift in Machine Learning

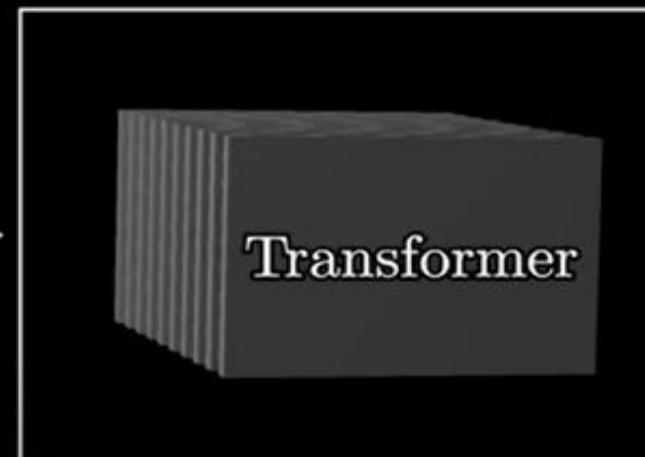
- What has caused this paradigm shift?
 - **Problem in recurrent networks**
 - Information is effectively lost during encoding of long sequences
 - Sequential nature disables parallel training and favors late timestep inputs
 - **Solution: Attention mechanism**
 - Handling long-range dependencies
 - Parallel training
 - Dynamic attention weights based on inputs





text-to-image

1960s photograph of a cute fluffy blue wild pi creature, a creature whose body is shaped like the symbol π , who is foraging in its native territory, staring back at the camera with an exotic scene in the background.

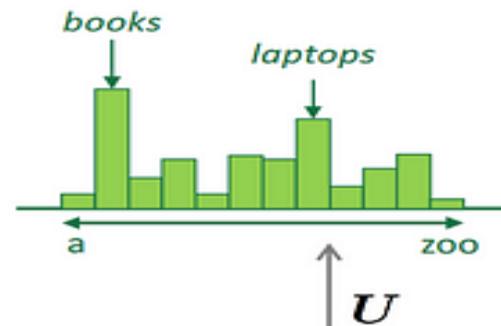


Generated by DALL-E 3

Recurrent Neural Networks

To increase the context available to our language model, we can use a recurrent neural network, a type of stateful NN. The network for each word maintains a state influenced by both the current word and previous hidden state. This architecture removes the fixed size of the context alleviating the problems of simpler neural language models.

$$\hat{y}^{(4)} = P(\mathbf{x}^{(5)} | \text{the students opened their})$$



output distribution

$$\hat{y}^{(t)} = \text{softmax} (\mathbf{U} \mathbf{h}^{(t)} + \mathbf{b}_2) \in \mathbb{R}^{|V|}$$

hidden states

$$\mathbf{h}^{(t)} = \sigma (\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_e \mathbf{e}^{(t)} + \mathbf{b}_1)$$

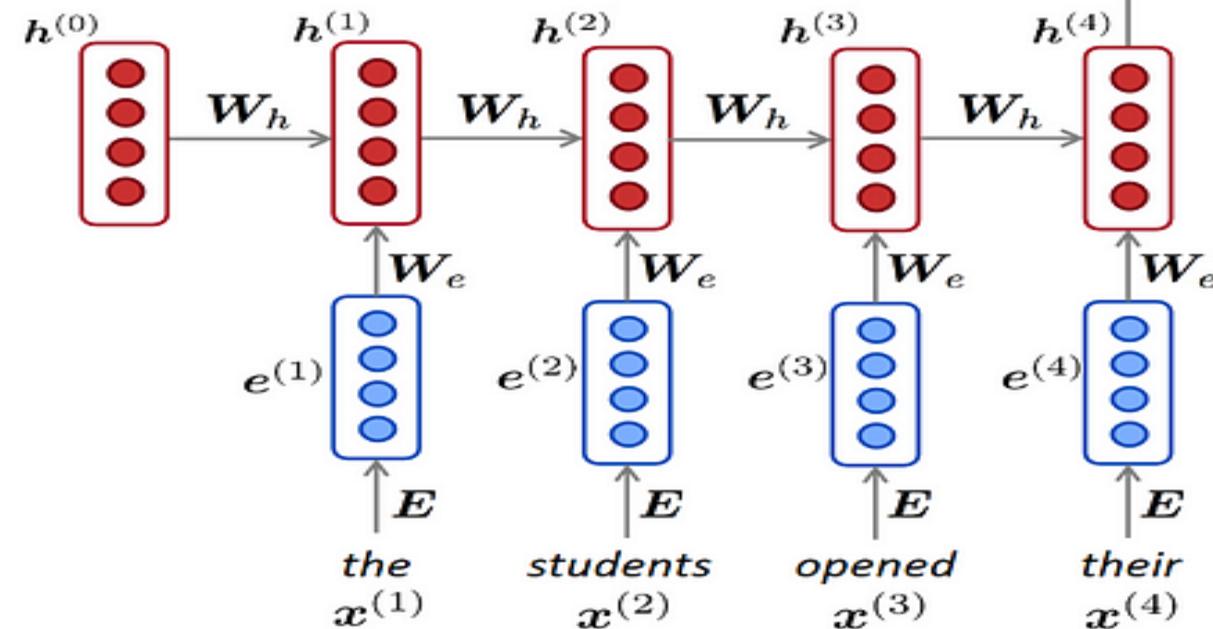
$\mathbf{h}^{(0)}$ is the initial hidden state

word embeddings

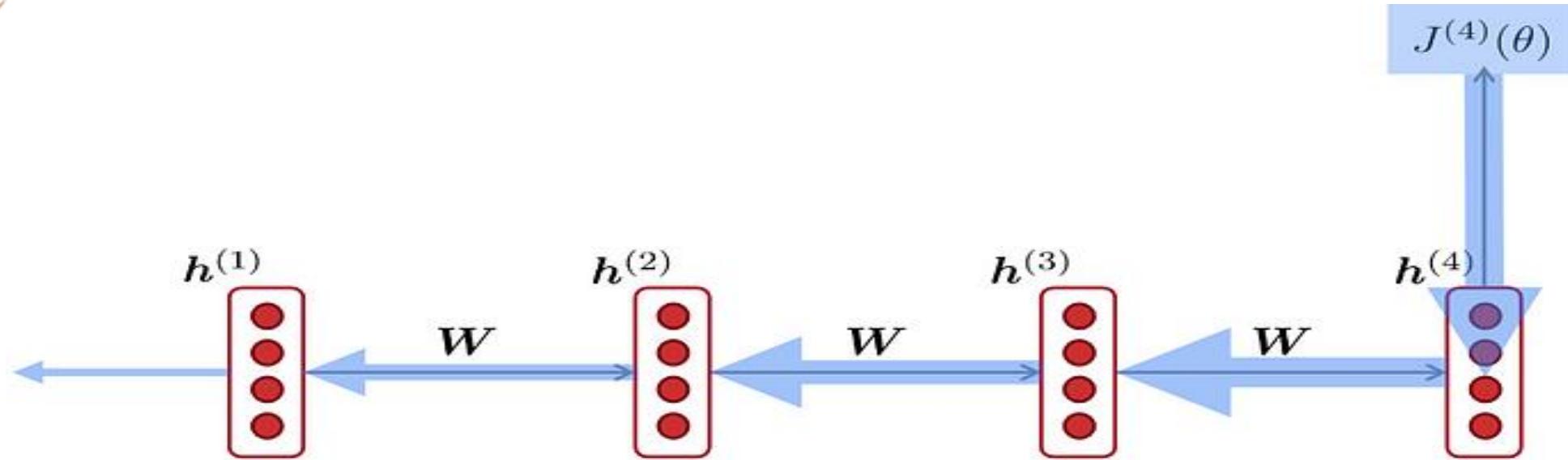
$$\mathbf{e}^{(t)} = \mathbf{E} \mathbf{x}^{(t)}$$

words / one-hot vectors

$$\mathbf{x}^{(t)} \in \mathbb{R}^{|V|}$$



Note: this input sequence could be much longer, but this slide doesn't have space!



$$\frac{\partial J^{(4)}}{\partial \mathbf{h}^{(1)}} = \boxed{\frac{\partial \mathbf{h}^{(2)}}{\partial \mathbf{h}^{(1)}} \times} \quad \boxed{\frac{\partial \mathbf{h}^{(3)}}{\partial \mathbf{h}^{(2)}} \times} \quad \boxed{\frac{\partial \mathbf{h}^{(4)}}{\partial \mathbf{h}^{(3)}} \times} \quad \frac{\partial J^{(4)}}{\partial \mathbf{h}^{(4)}}$$

What happens if these are small?

Vanishing gradient problem:
When these are small, the gradient signal gets smaller and smaller as it backpropagates further

LSTM

One of the architecture proposed to solve the vanishing gradient problem is called Long Short-Term Memory, and it works by having both a hidden state and a memory cell and three gates that form read (in the hidden state), write and delete of the cell. They are called the output, input, and forget gate, respectively. This architecture implements a dedicated way to maintain long-term dependencies, which is never forgetting the memory cell.

We have a sequence of inputs $\mathbf{x}^{(t)}$, and we will compute a sequence of hidden states $\mathbf{h}^{(t)}$ and cell states $\mathbf{c}^{(t)}$. On timestep t :

Forget gate: controls what is kept vs forgotten, from previous cell state

Input gate: controls what parts of the new cell content are written to cell

Output gate: controls what parts of cell are output to hidden state

New cell content: this is the new content to be written to the cell

Cell state: erase (“forget”) some content from last cell state, and write (“input”) some new cell content

Hidden state: read (“output”) some content from the cell

Sigmoid function: all gate values are between 0 and 1

$$\mathbf{f}^{(t)} = \sigma(\mathbf{W}_f \mathbf{h}^{(t-1)} + \mathbf{U}_f \mathbf{x}^{(t)} + \mathbf{b}_f)$$

$$\mathbf{i}^{(t)} = \sigma(\mathbf{W}_i \mathbf{h}^{(t-1)} + \mathbf{U}_i \mathbf{x}^{(t)} + \mathbf{b}_i)$$

$$\mathbf{o}^{(t)} = \sigma(\mathbf{W}_o \mathbf{h}^{(t-1)} + \mathbf{U}_o \mathbf{x}^{(t)} + \mathbf{b}_o)$$

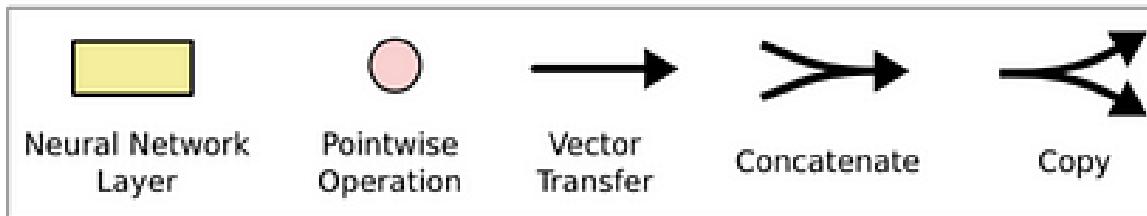
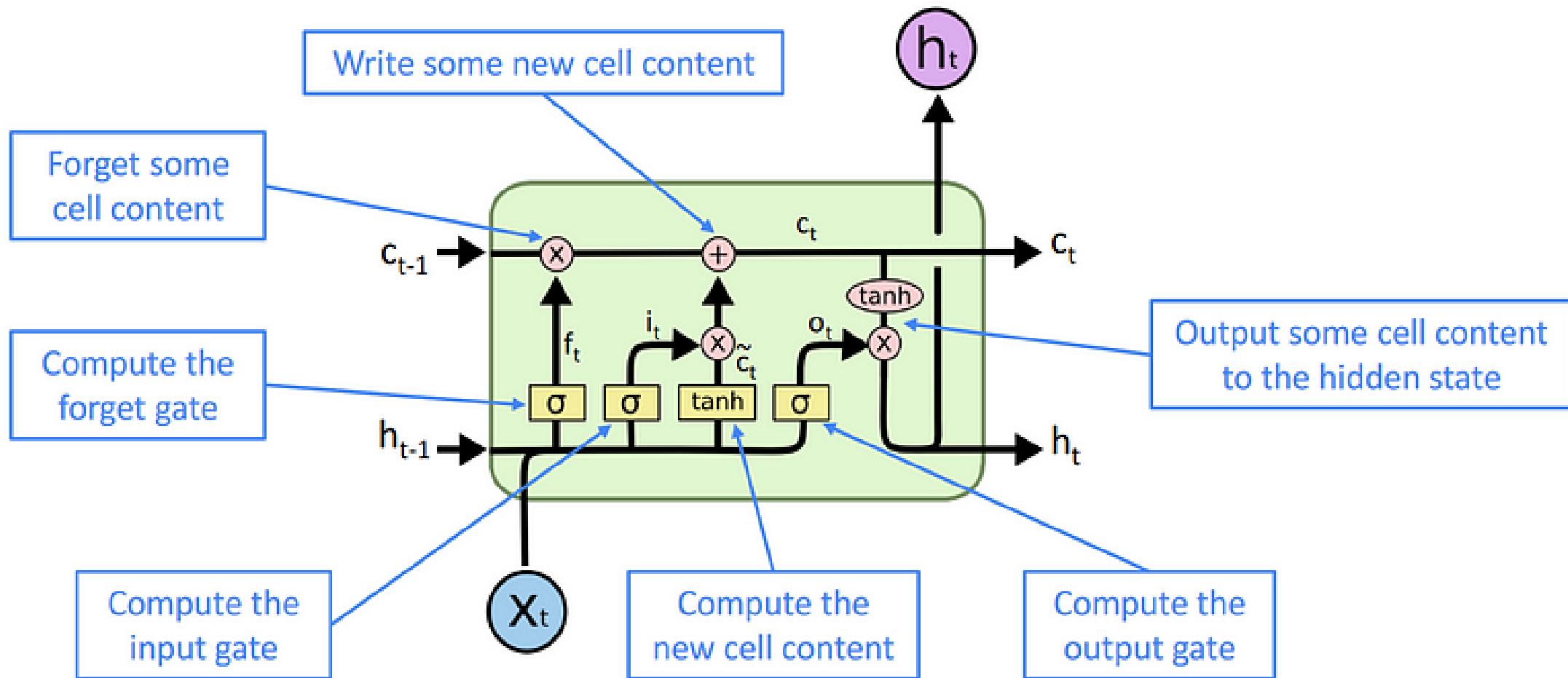
$$\tilde{\mathbf{c}}^{(t)} = \tanh(\mathbf{W}_c \mathbf{h}^{(t-1)} + \mathbf{U}_c \mathbf{x}^{(t)} + \mathbf{b}_c)$$

$$\mathbf{c}^{(t)} = \mathbf{f}^{(t)} \circ \mathbf{c}^{(t-1)} + \mathbf{i}^{(t)} \circ \tilde{\mathbf{c}}^{(t)}$$

$$\mathbf{h}^{(t)} = \mathbf{o}^{(t)} \circ \tanh \mathbf{c}^{(t)}$$

Gates are applied using element-wise product

All these are vectors of same length n



- Seq2Seq (Sequence-to-Sequence) models are a type of neural network, an exceptional Recurrent Neural Network architecture, designed to transform one data sequence into another.
- They are handy for tasks such as solving complex language problems like machine translation, question answering, creating chatbots, text summarization, etc.

Machine Language Translation

Les modèles de séquence sont super puissants → Sequence Model → *Sequence models are super powerful*

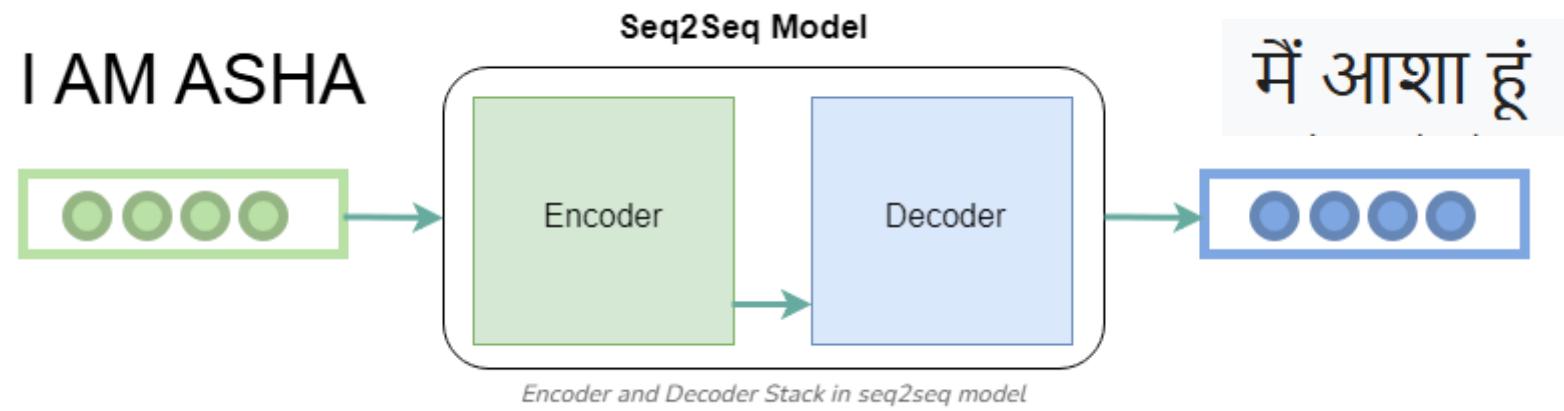
Text Summarization

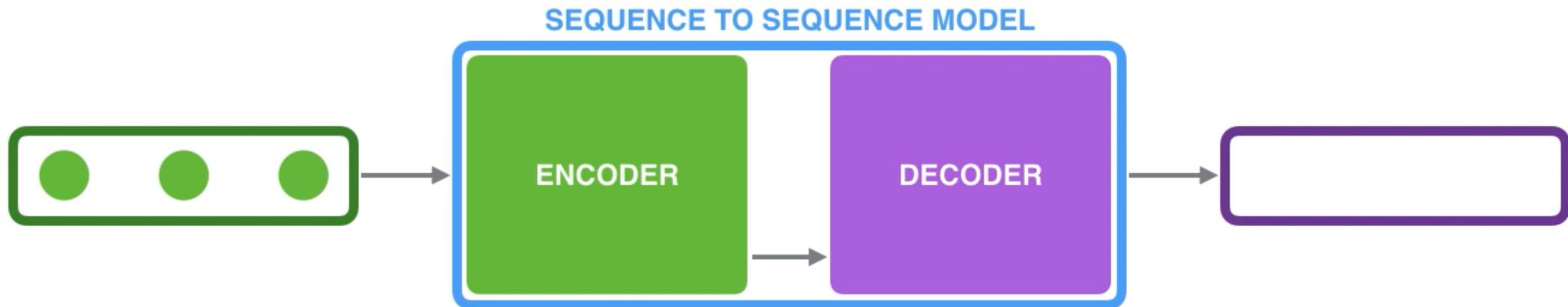
A strong analyst have 6 main characteristics. One should master all 6 to be successful in the industry :
1.
2.

→ Sequence Model → *6 characteristics of successful analyst*

Chatbot

How are you doing today? → Sequence Model → *I am doing well. Thank you.
How are you doing today?*

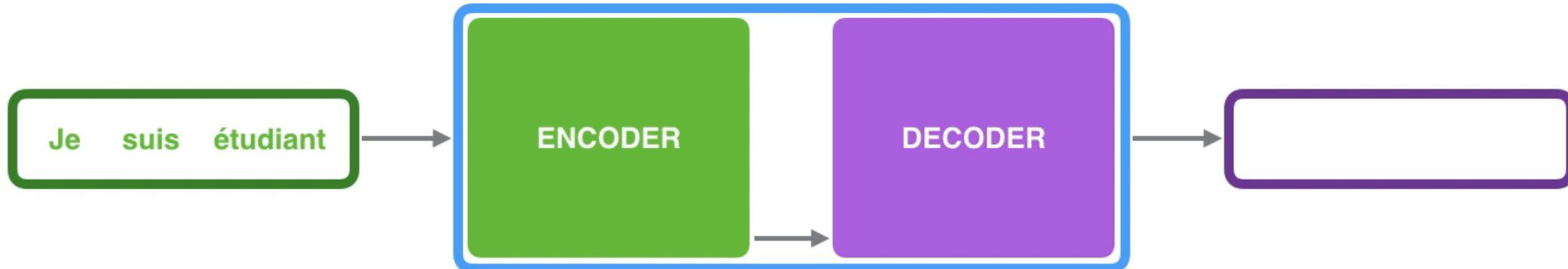




The model is composed of an **encoder** and a **decoder**.

The **encoder** processes each item in the input sequence, it compiles the information it captures into a vector (called the **context**). After processing the entire input sequence, the **encoder** sends the **context** over to the **decoder**, which begins producing the output sequence item by item.

Neural Machine Translation SEQUENCE TO SEQUENCE MODEL



The **context** is a vector (an array of numbers, basically) in the case of machine translation. The **encoder** and **decoder** tend to both be recurrent neural networks

CONTEXT

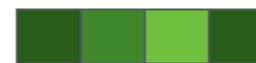
0.11
0.03
0.81
-0.62

0.11
0.03
0.81
-0.62

Input

Je
suis
étudiant

0.901	-0.651	-0.194	-0.822
-0.351	0.123	0.435	-0.200
0.081	0.458	-0.400	0.480

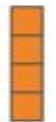


By design, a RNN takes two inputs at each time step: an input (in the case of the encoder, one word from the input sentence), and a hidden state. The word, however, needs to be represented by a vector. To transform a word into a vector, we turn to the class of methods called “[word embedding](#)” algorithms. These turn words into vector spaces that capture a lot of the meaning/semantic information of the words

Recurrent Neural Network

Time step #1:

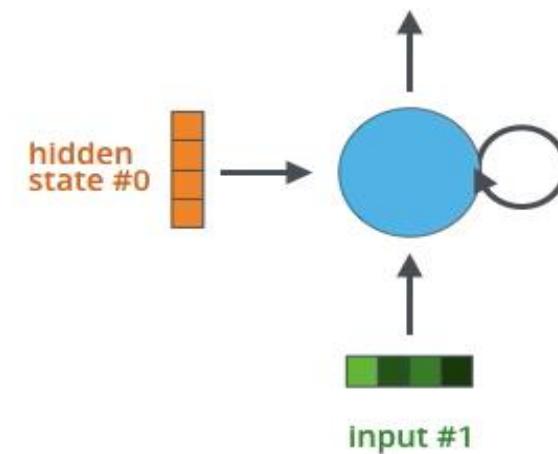
An RNN takes two input vectors:



hidden
state #0



input vector #1

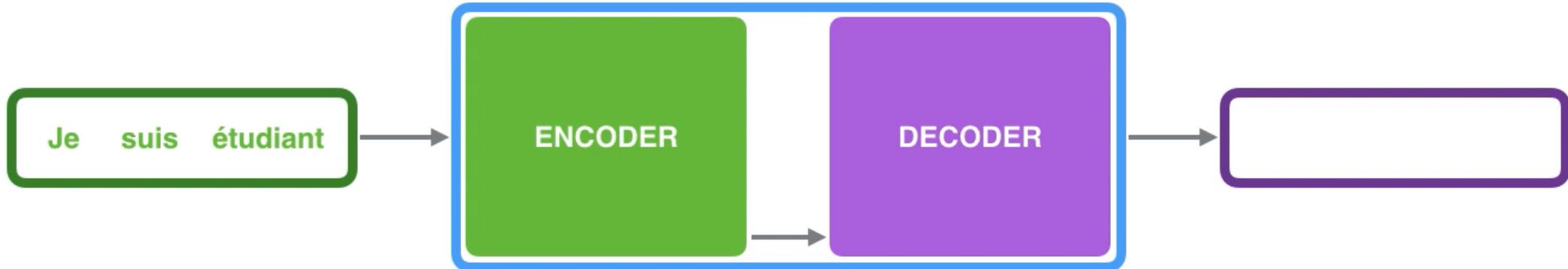


Each pulse for the **encoder** or **decoder** is that RNN processing its inputs and generating an output for that time step.

Since the **encoder** and **decoder** are both RNNs, each time step one of the RNNs does some processing, it updates its **hidden state** based on its inputs and previous inputs it has seen.

Time step:

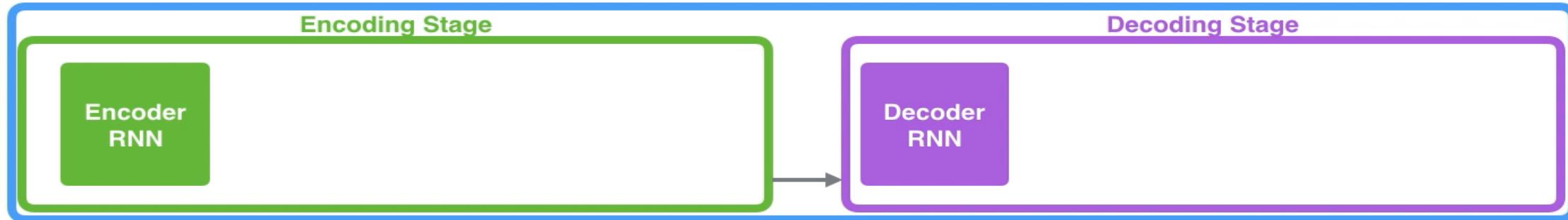
Neural Machine Translation SEQUENCE TO SEQUENCE MODEL



When we look at the **hidden states** for the **encoder**. Notice how the last **hidden state** is actually the **context** we pass along to the **decoder**.

Neural Machine Translation

SEQUENCE TO SEQUENCE MODEL



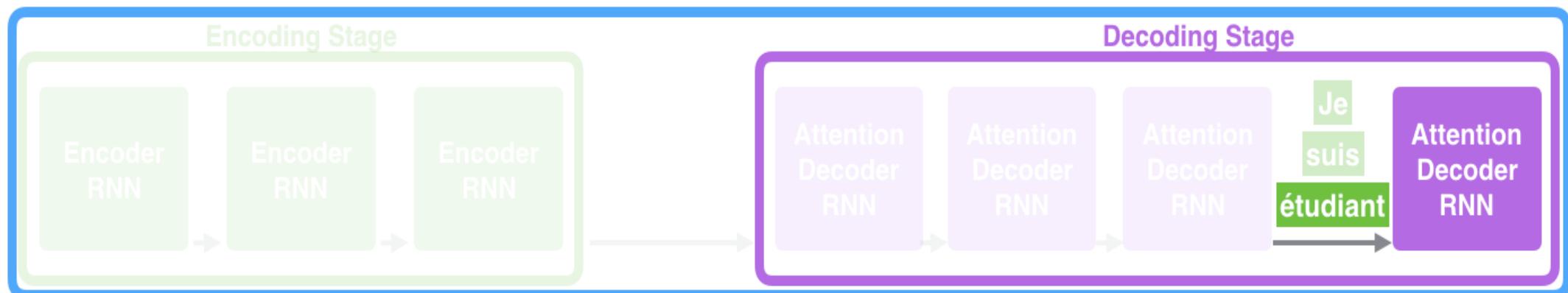
Je suis étudiant

This is called an “unrolled” view where instead of showing the one **decoder**, we show a copy of it for each time step as we can look at the inputs and outputs of each time step.

The **context** vector turned out to be a bottleneck for these types of models. It made it challenging for the models to deal with long sentences. A solution was proposed in [Bahdanau et al., 2014](#) and [Luong et al., 2015](#). These papers introduced and refined a technique called “Attention”, which highly improved the quality of machine translation systems. Attention allows the model to focus on the relevant parts of the input sequence as needed.

Time step: 7

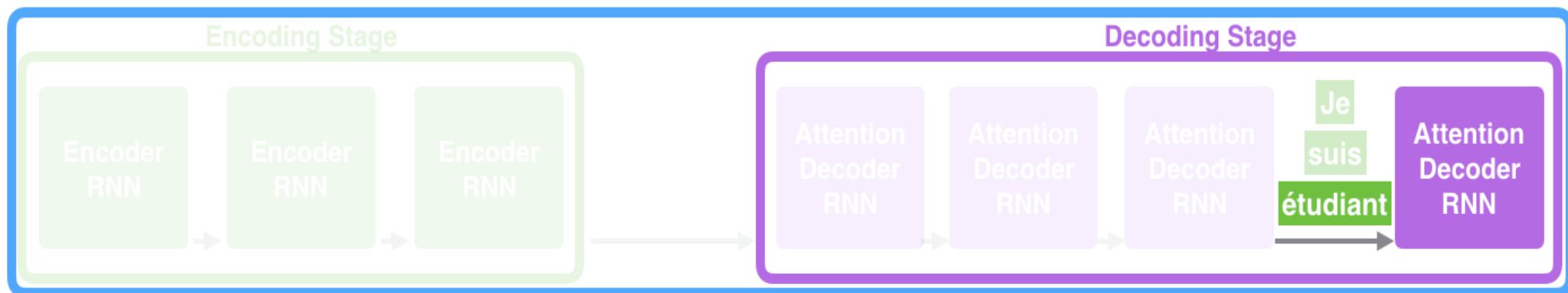
Neural Machine Translation SEQUENCE TO SEQUENCE MODEL WITH ATTENTION



At time step 7, the attention mechanism enables the decoder to focus on the word "étudiant" ("student" in french) before it generates the English translation. This ability to amplify the signal from the relevant part of the input sequence makes attention models produce better results than models without attention.

Time step: 7

Neural Machine Translation SEQUENCE TO SEQUENCE MODEL WITH ATTENTION



Neural Machine Translation

SEQUENCE TO SEQUENCE MODEL WITH ATTENTION



Je

suis

étudiant

An attention model differs from a classic sequence-to-sequence model in two main ways:

- 1) First, the **encoder** passes a lot more data to the **decoder**. Instead of passing the last hidden state of the encoding stage, the **encoder** passes *all* the **hidden states** to the **decoder**
- 2) Second, an attention **decoder** does an extra step before producing its output. In order to focus on the parts of the input that are relevant to this decoding time step, the **decoder** does the following:
 - Look at the set of encoder **hidden states** it received – each **encoder hidden state** is most associated with a certain word in the input sentence
 - Give each **hidden state** a score (let's ignore how the scoring is done for now)
 - Multiply each **hidden state** by its softmaxed score, thus amplifying **hidden states** with high scores, and drowning out **hidden states** with low scores

Attention at time step 4



- 1.The attention decoder RNN takes in the embedding of the <END> token, and an initial decoder hidden state.
- 2.The RNN processes its inputs, producing an output and a new hidden state vector (h_4). The output is discarded.
- 3.Attention Step: We use the encoder hidden states and the h_4 vector to calculate a context vector (C_4) for this time step.
- 4.We concatenate h_4 and C_4 into one vector.
- 5.We pass this vector through a feedforward neural network (one trained jointly with the model).
- 6.The output of the feedforward neural networks indicates the output word of this time step.
- 7.Repeat for the next time steps

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

Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.0 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature.

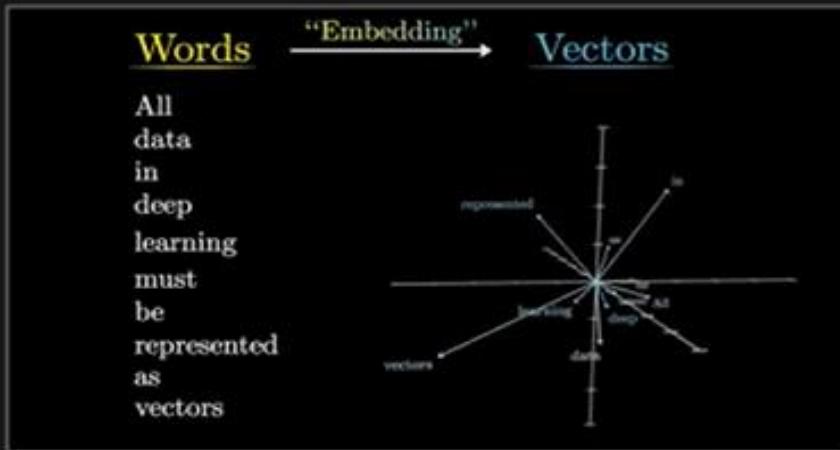
Attention is all
you need



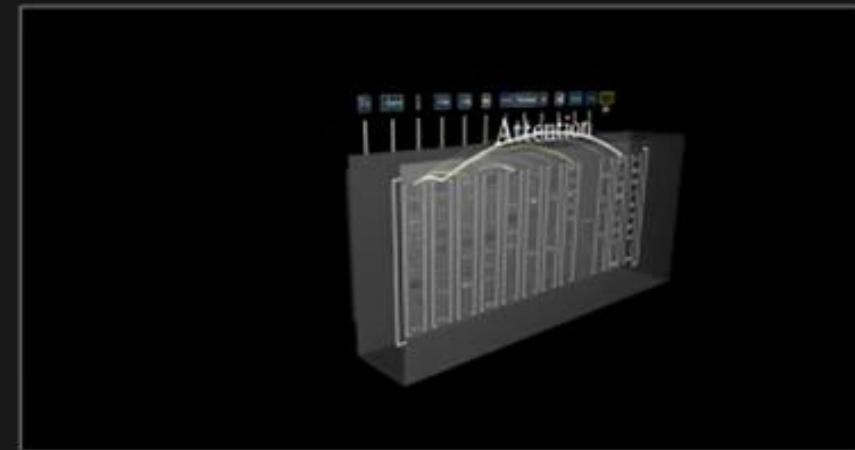
machine translation

注意力就是你所需要的一切

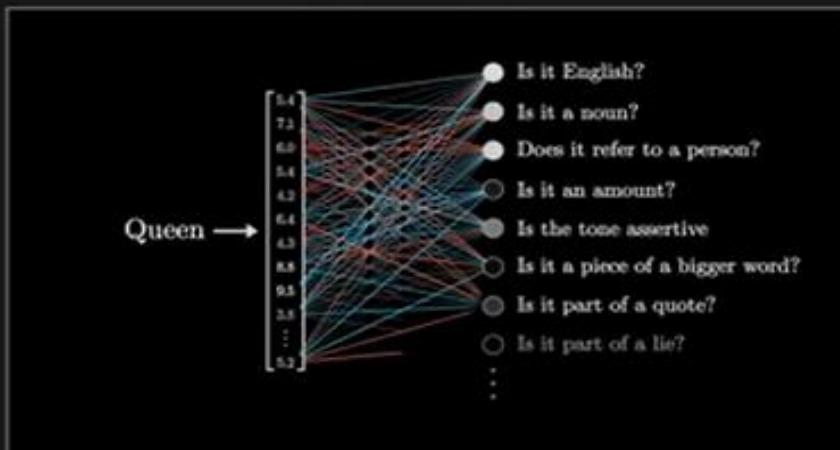
Embedding



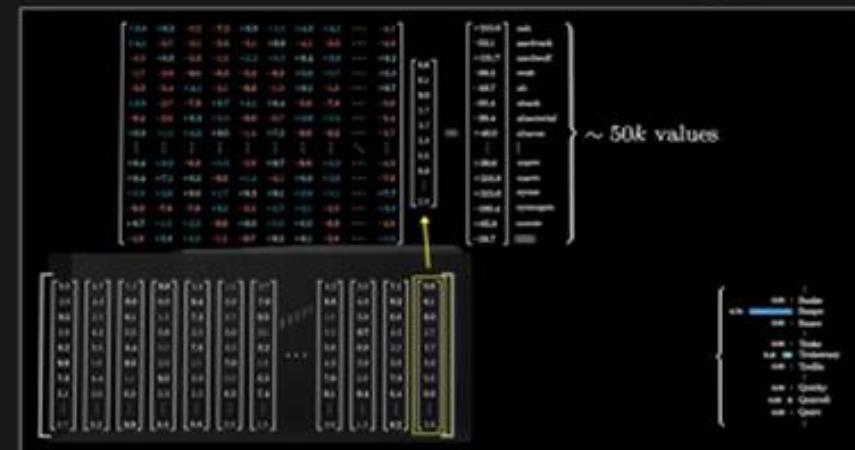
Attention



MLPs



Unembedding

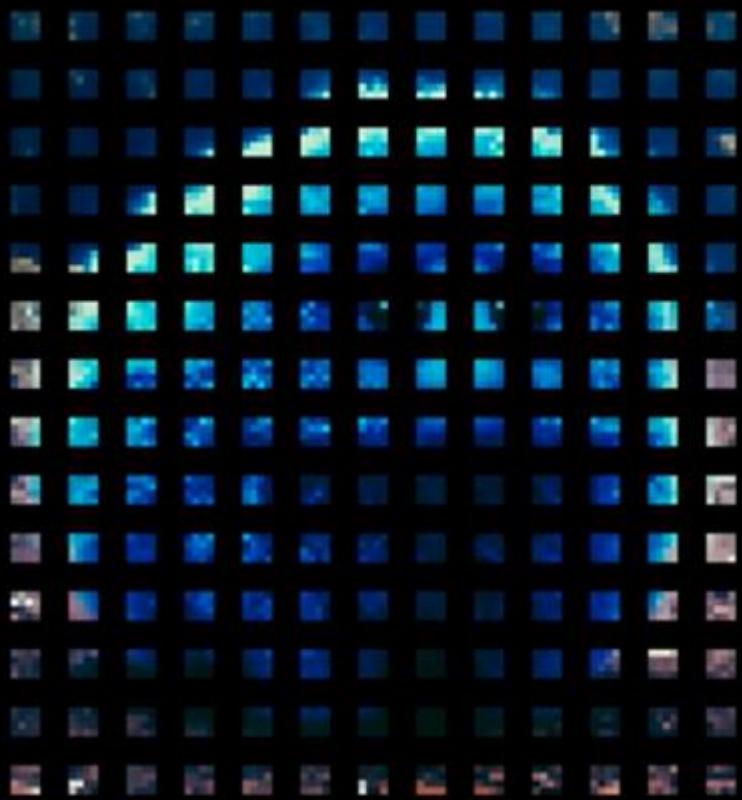


Tokens



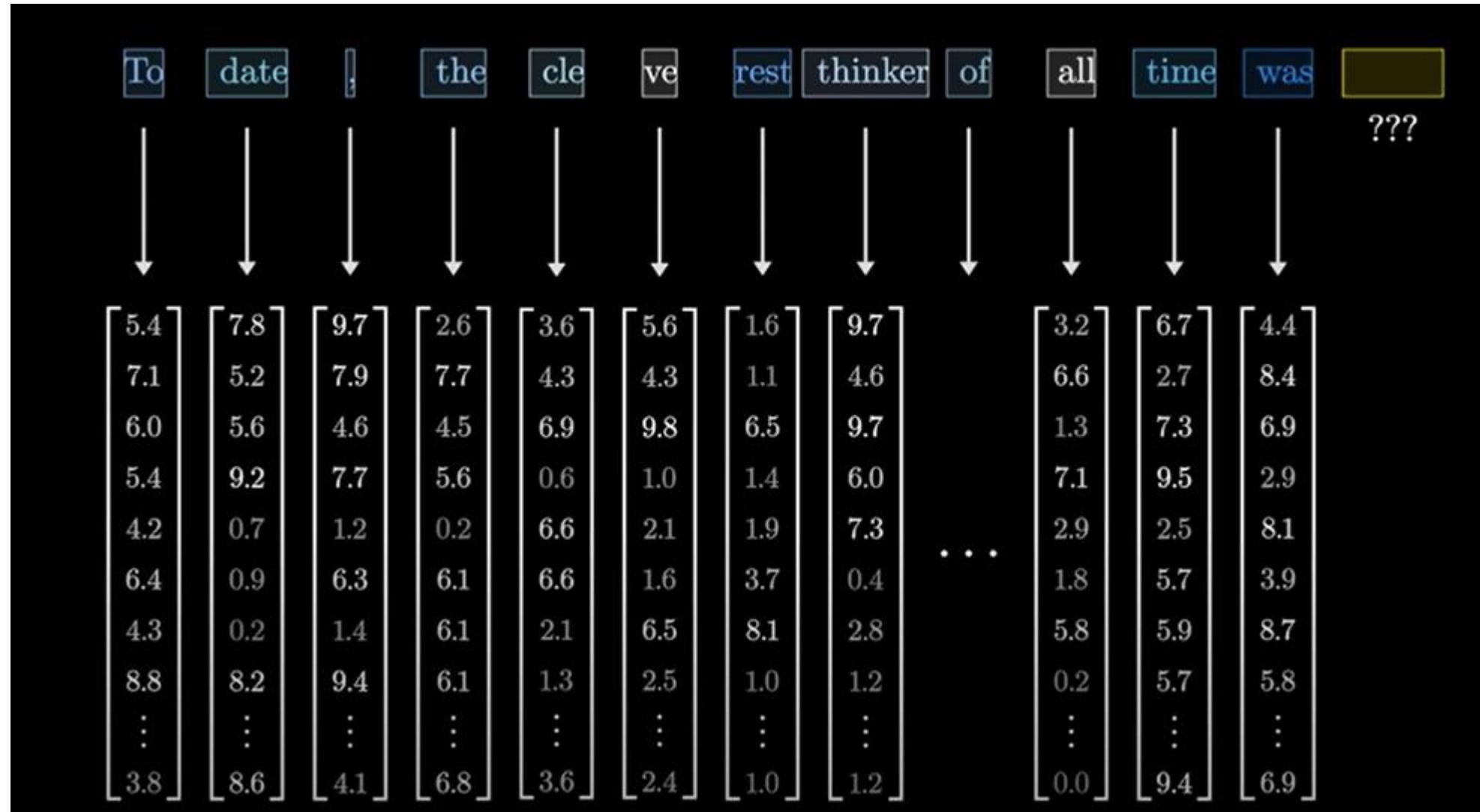
To date, the cleverest thinker of all time was ???

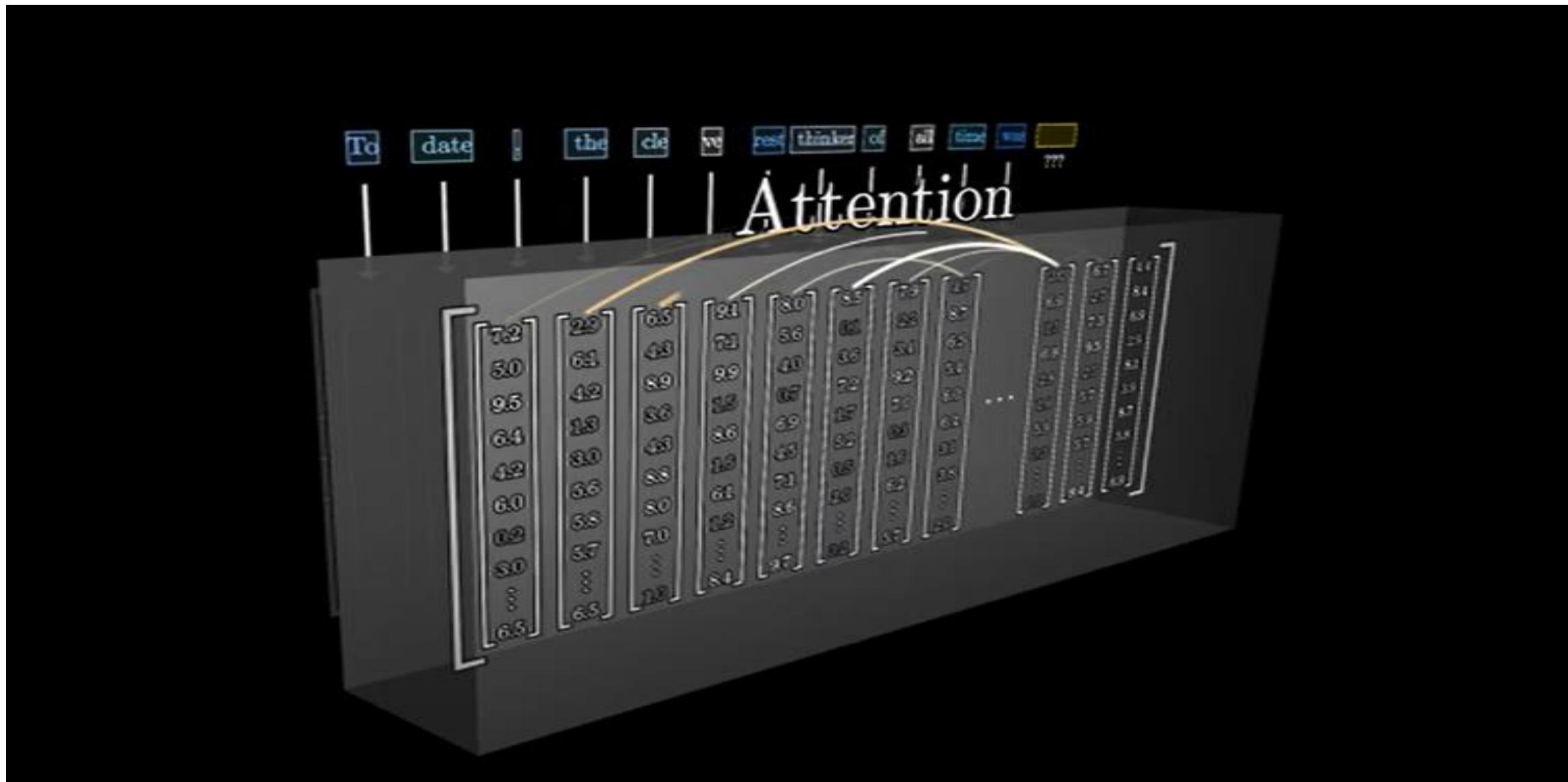
The diagram illustrates tokens as colored rectangles underneath the text "To date, the cleverest thinker of all time was ???". The word "cleverest" is highlighted in red, and the word "was" is highlighted in yellow. Three blue arrows point from the top towards the word "cleverest", indicating it is the target token for processing.

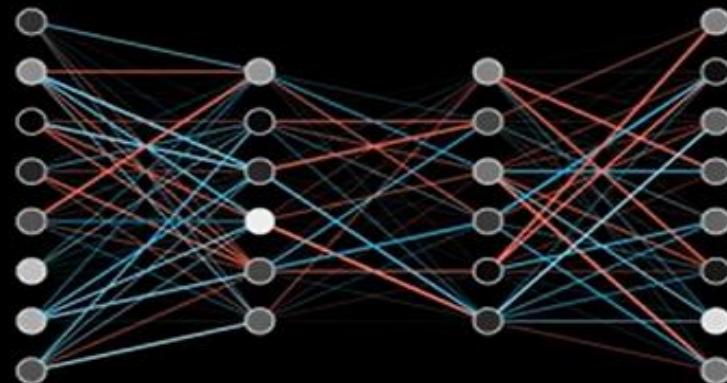


To date, the cleverest thinker of all time was

???







A machine learning model ...

$$\begin{bmatrix} 3.6 \\ 5.6 \\ 4.3 \\ 9.8 \\ 1.0 \\ \vdots \\ 2.1 \end{bmatrix} \quad \begin{bmatrix} 1.6 \\ 6.5 \\ 2.5 \\ 4.6 \\ 2.4 \\ \vdots \\ 1.6 \end{bmatrix} \quad \begin{bmatrix} 1.1 \\ 6.5 \\ 1.4 \\ 1.9 \\ 3.7 \\ \vdots \\ 8.1 \end{bmatrix} \quad \begin{bmatrix} 1.0 \\ 8.3 \\ 1.0 \\ 9.7 \\ 4.6 \\ \vdots \\ 9.7 \end{bmatrix}$$



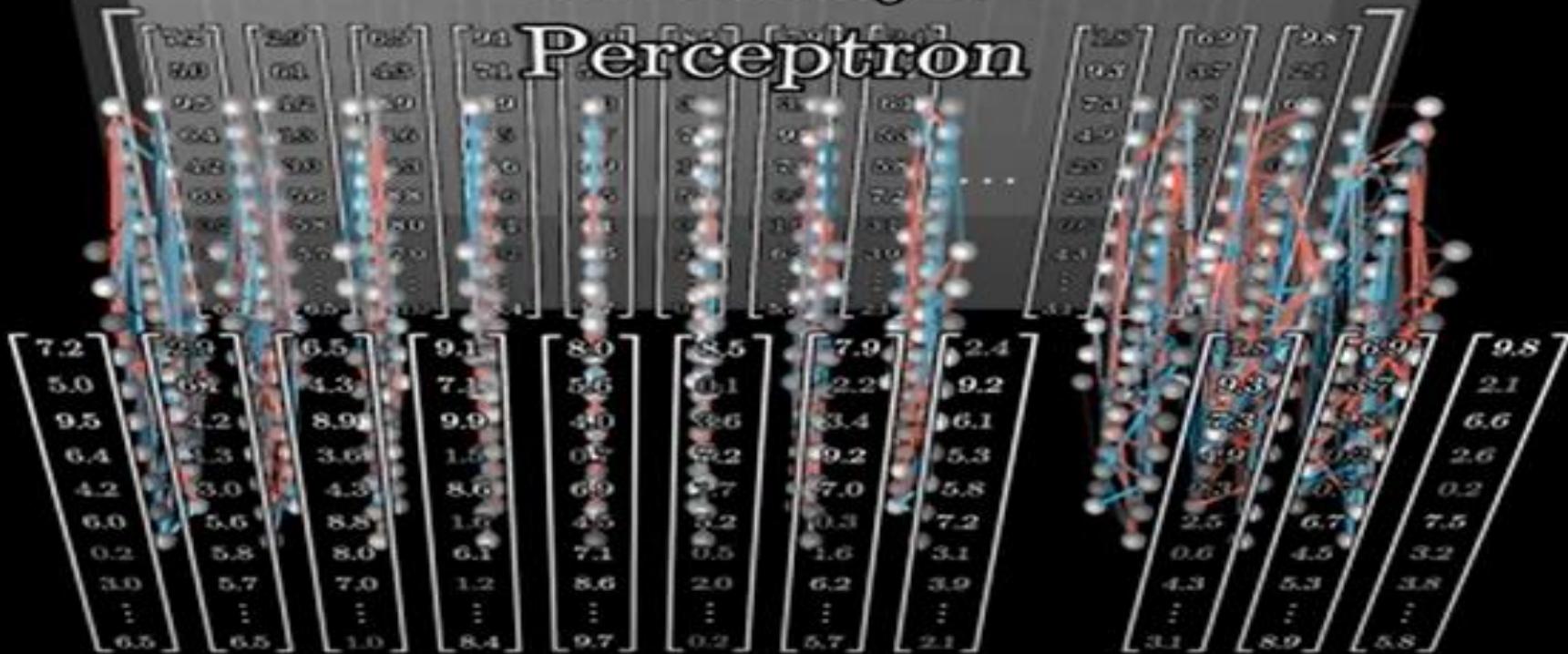
A fashion model ...

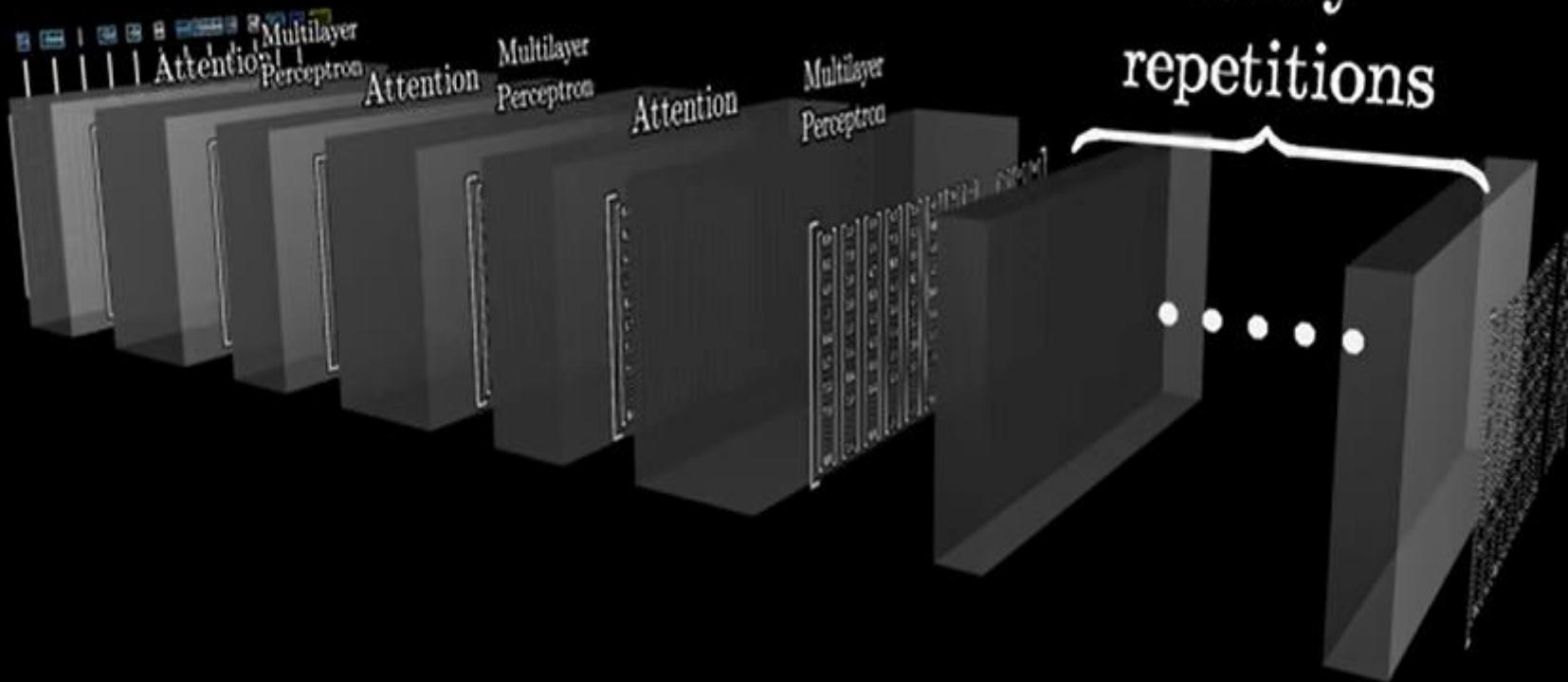
$$\begin{bmatrix} 6.0 \\ 7.3 \\ 0.4 \\ 2.8 \\ 1.2 \\ \vdots \\ 2.9 \end{bmatrix} \quad \begin{bmatrix} 1.2 \\ 3.1 \\ 4.1 \\ 0.6 \\ 6.9 \\ \vdots \\ 5.6 \end{bmatrix} \quad \begin{bmatrix} 2.6 \\ 5.2 \\ 0.9 \\ 5.7 \\ 9.2 \\ \vdots \\ 3.2 \end{bmatrix}$$

To date is the cle ve rest thinker of all time who ???
↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓

Attention

Multilayer
Perceptron





Many
repetitions

Attention Mechanism

Assign attention weight to each word, to know how much "attention" the model should pay to each word (i.e., for each word, the network learns a "context")

Originally developed for language translation:

Bahdanau, D., Cho, K., & Bengio, Y. (2014). Neural machine translation by jointly learning to align and translate. <https://arxiv.org/abs/1409.0473>

**Hidden state in a regular RNN
(RNN #2)**

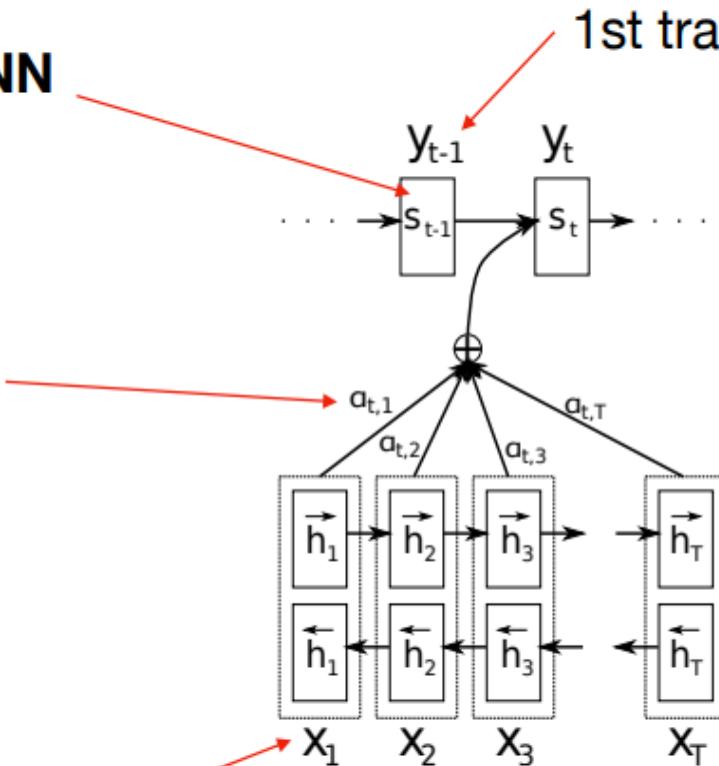
1st translated word

Attention weight

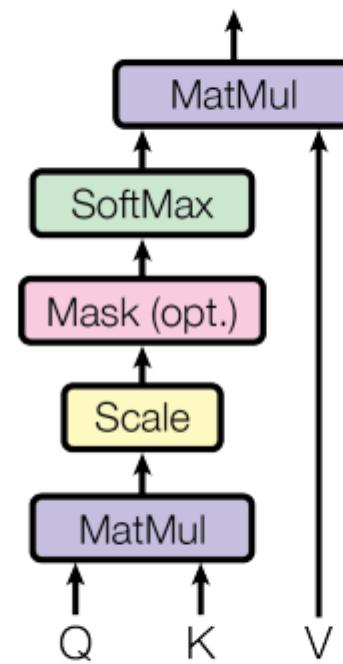
**Bidirectional RNN
(RNN #1)**

1st input word

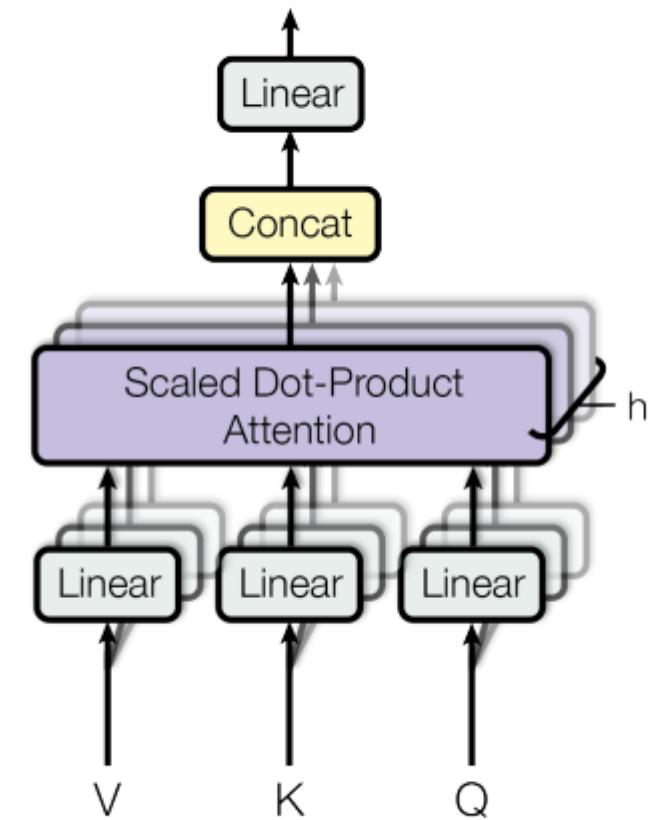
Figure 1: The graphical illustration of the proposed model trying to generate the t -th target word y_t given a source sentence (x_1, x_2, \dots, x_T) .

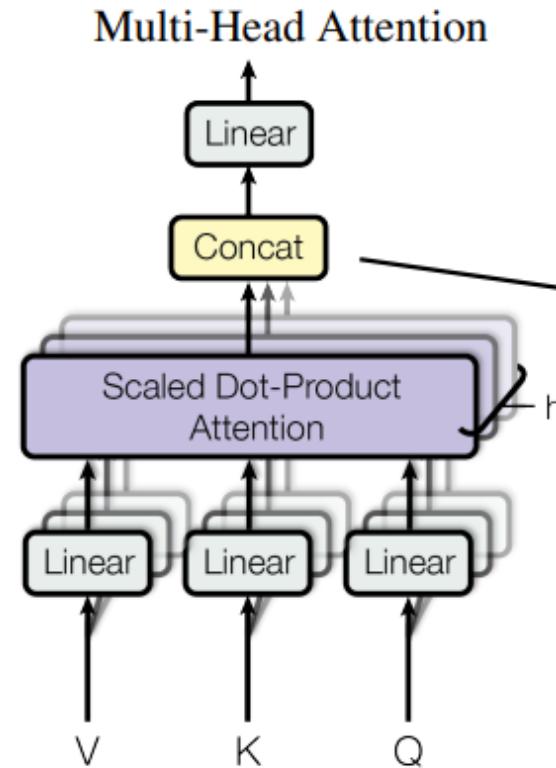


Scaled Dot-Product Attention



Multi-Head Attention





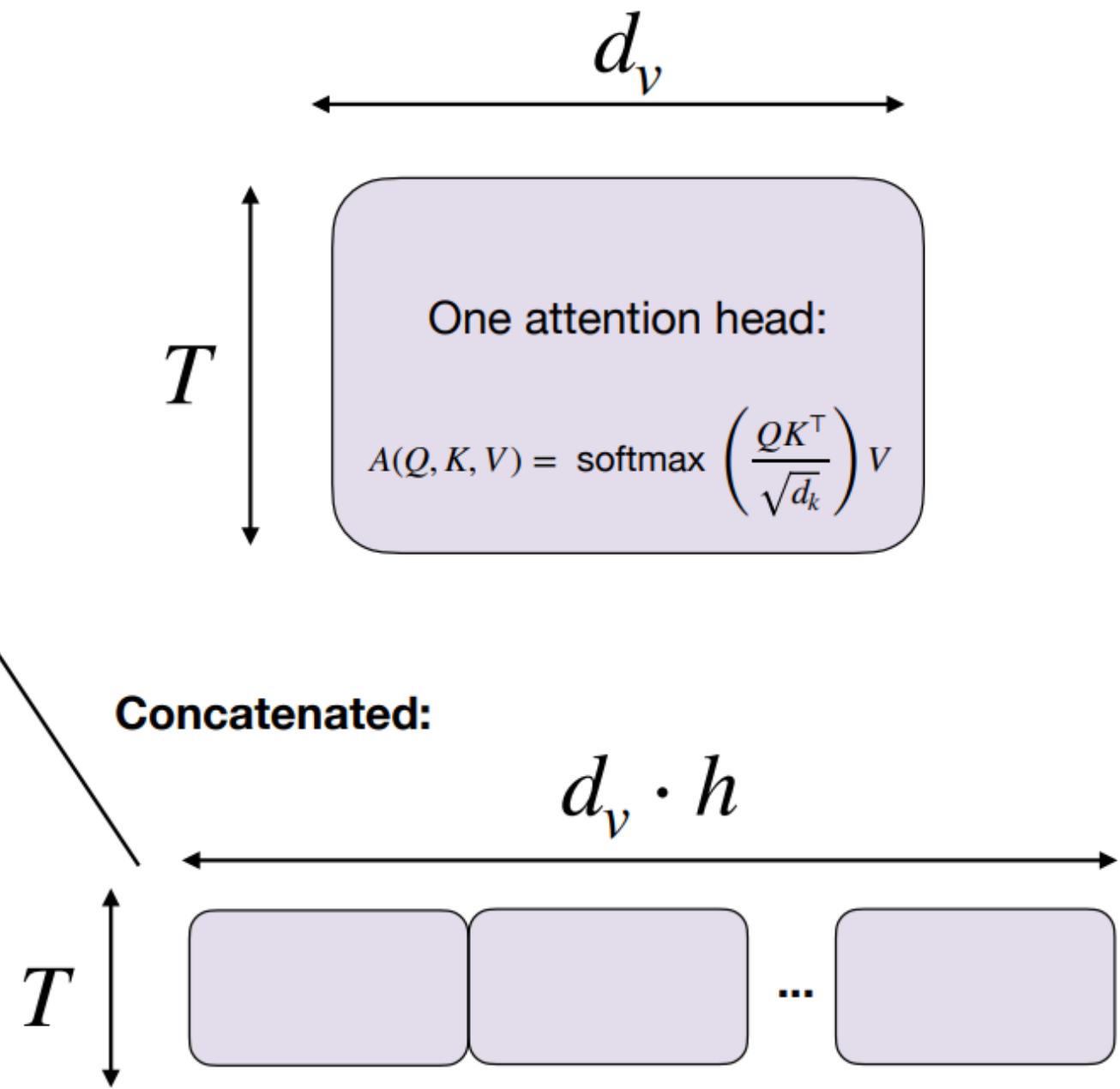
Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L. and Polosukhin, I., 2017. Attention Is All You Need.

Input sequence dim.
in original transformer:

$$T \times d_e = T \times 512$$

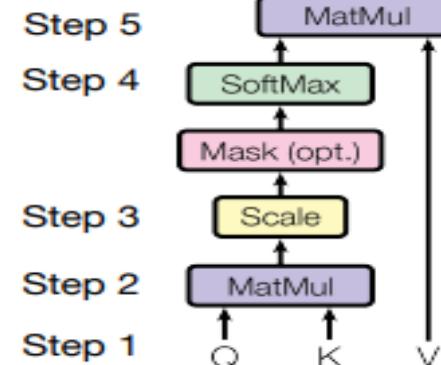
and

$$d_v = 512/h = 64$$



Scaled Dot-Product Attention Recap

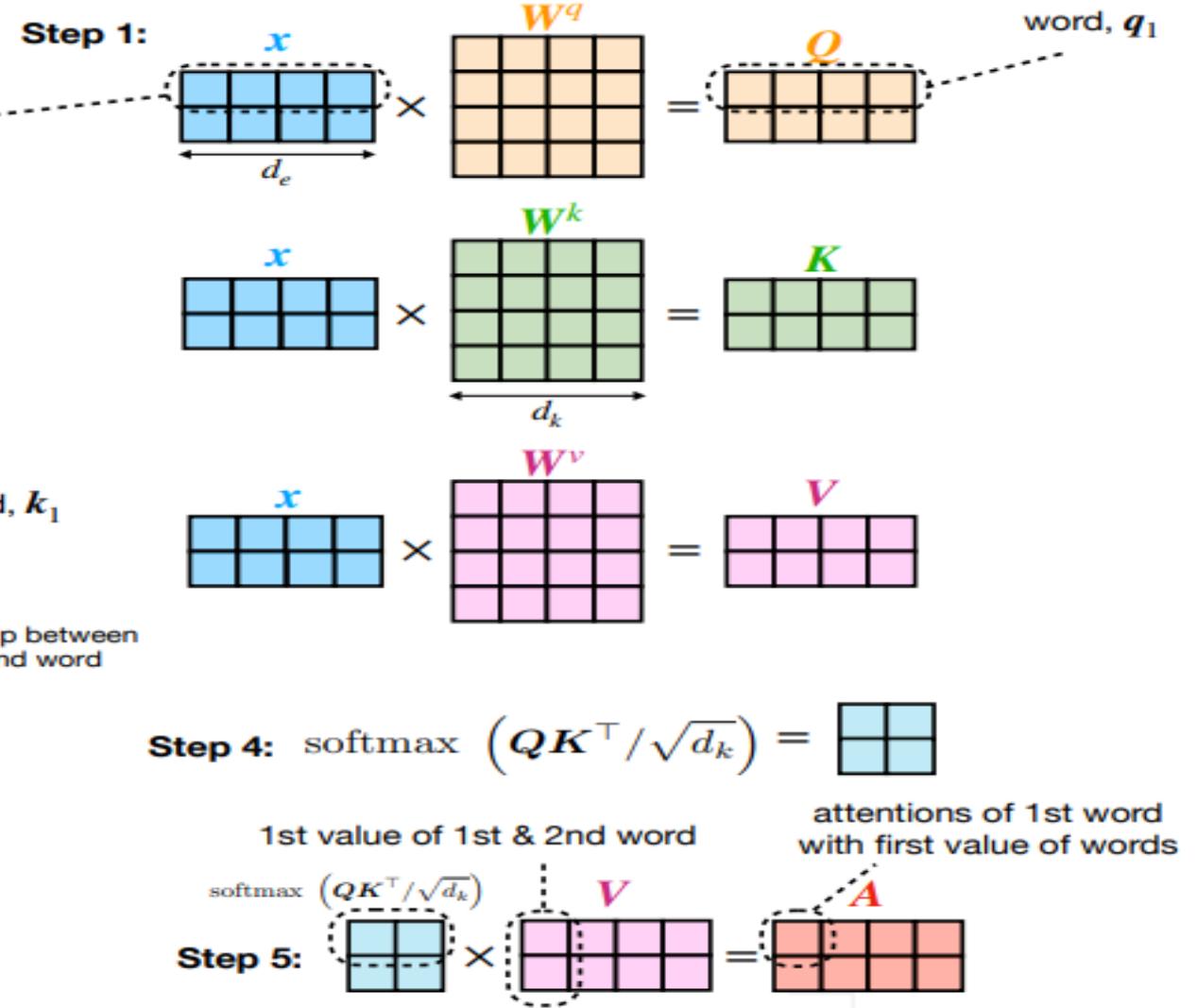
Scaled Dot-Product Attention

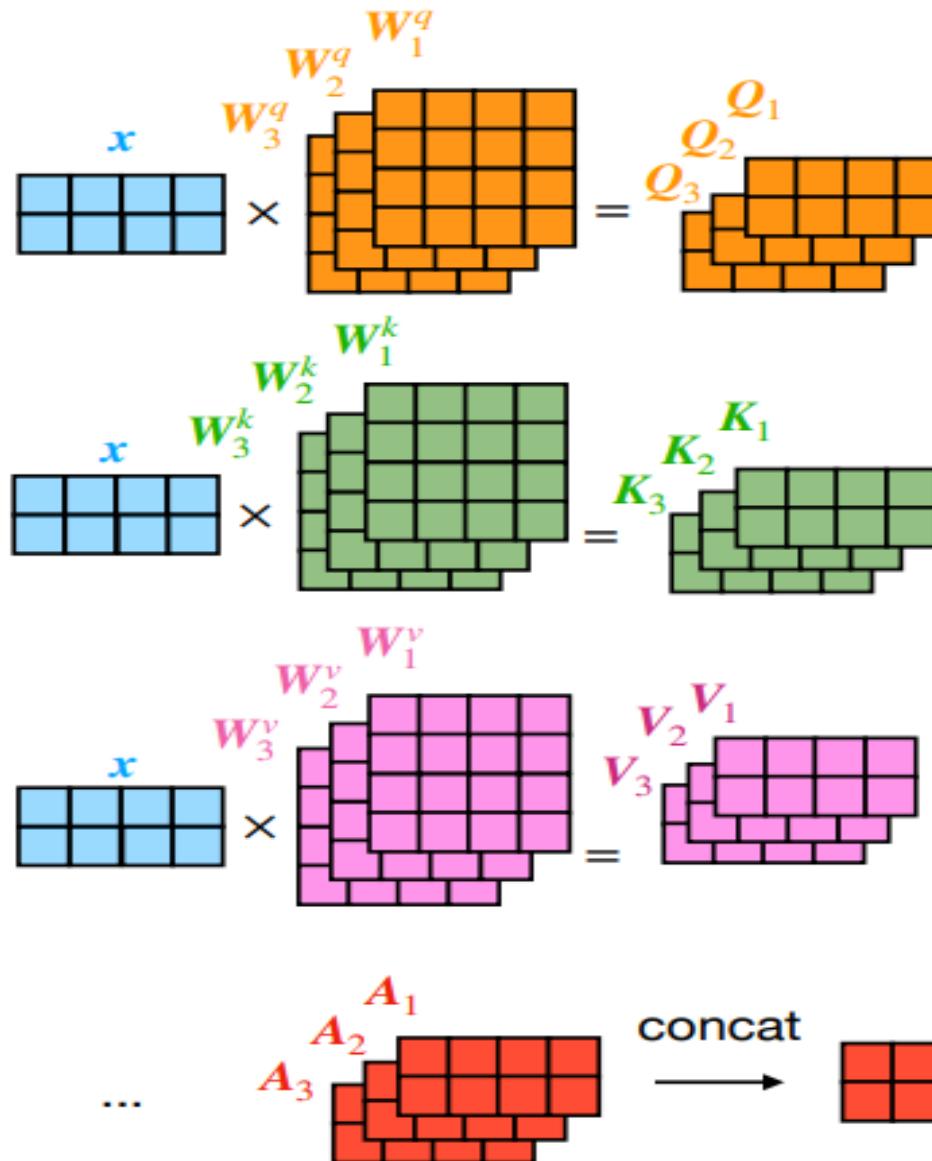


Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L. and Polosukhin, I., 2017. Attention Is All You Need.

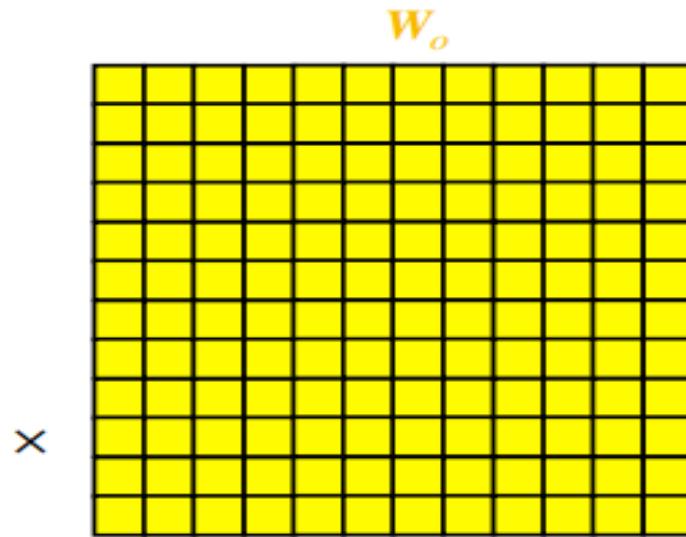
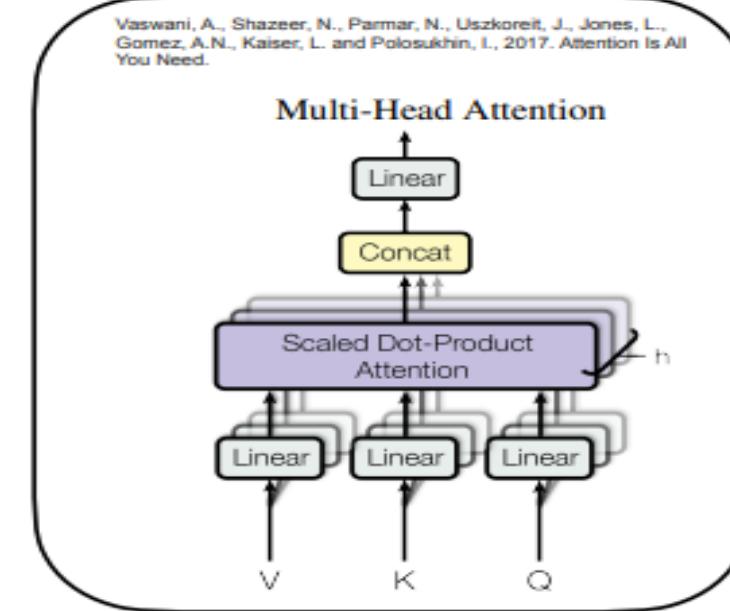
Step 2: $Q \times K^\top = QK^\top$ relationship between 1st & 2nd word

Step 3: $QK^\top / \sqrt{d_k} = \frac{QK^\top}{\sqrt{d_k}}$





Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L. and Polosukhin, I., 2017. Attention Is All You Need.



Multilayer Perceptrons

Generates output words one at a time

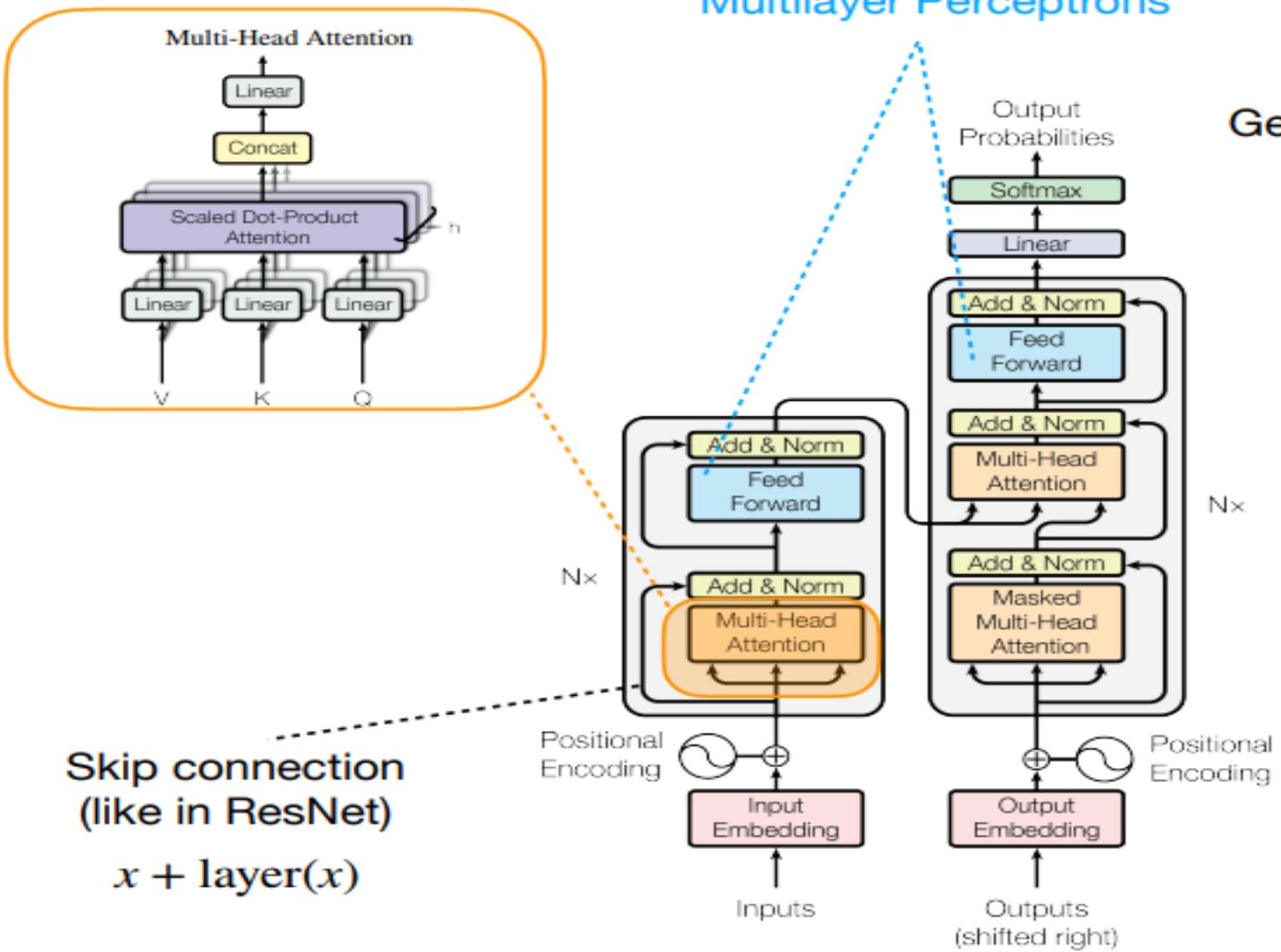
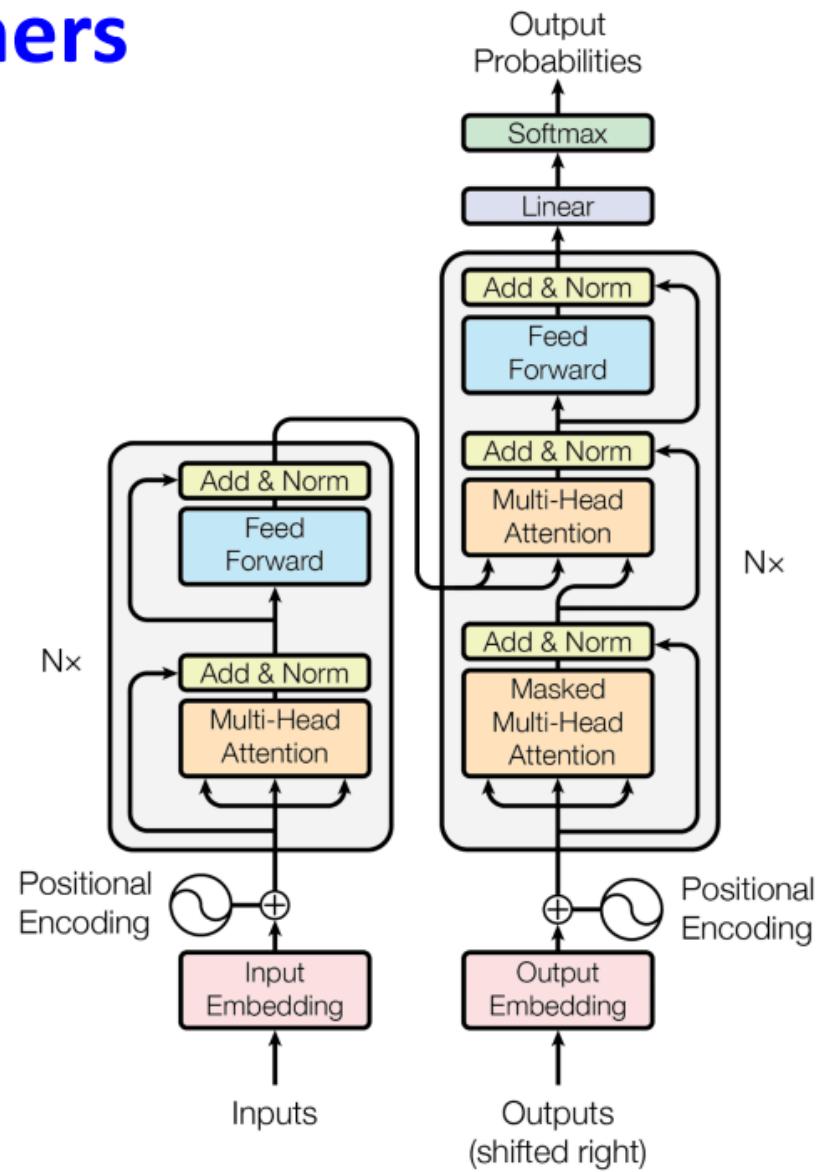


Figure 1: The Transformer - model architecture.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L. and Polosukhin, I., 2017. Attention Is All You Need.

Transformers

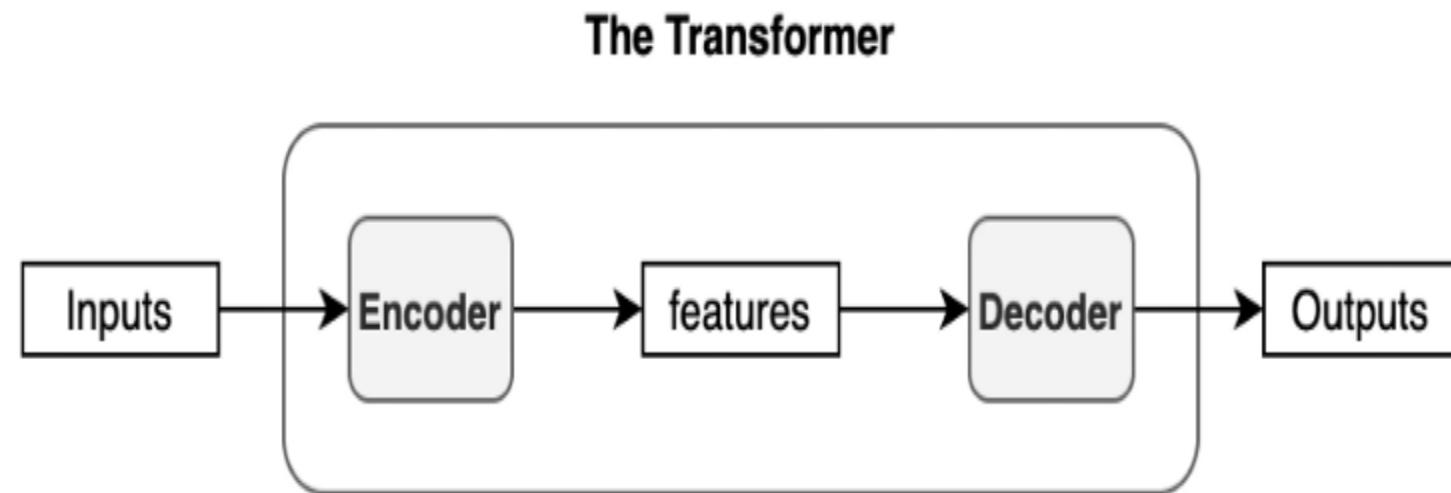
- Tokenization
- Input Embeddings
- Position Encodings
- Residuals
- Query
- Key
- Value
- Add & Norm
- Encoder
- Decoder
- Attention
- Self Attention
- Multi Head Attention
- Masked Attention
- Encoder Decoder Attention
- Output Probabilities / Logits
- Softmax
- Encoder-Decoder models
- Decoder only models



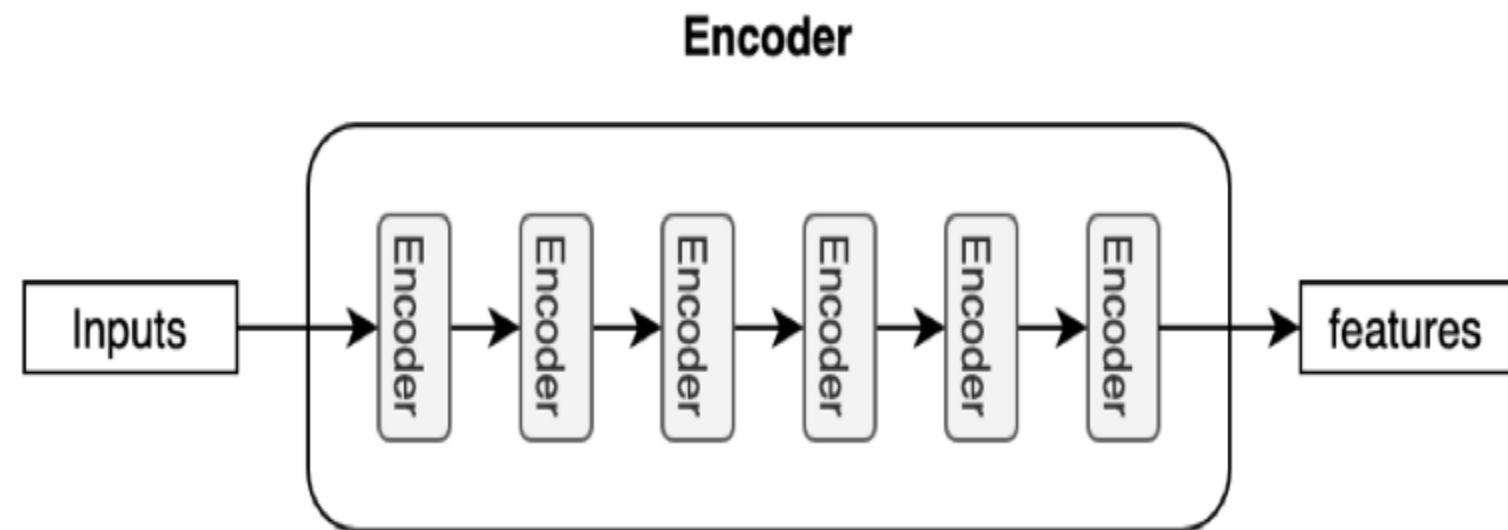
The original transformer published in the paper is a neural machine translation model. For example, we can train it to translate English into French sentences.



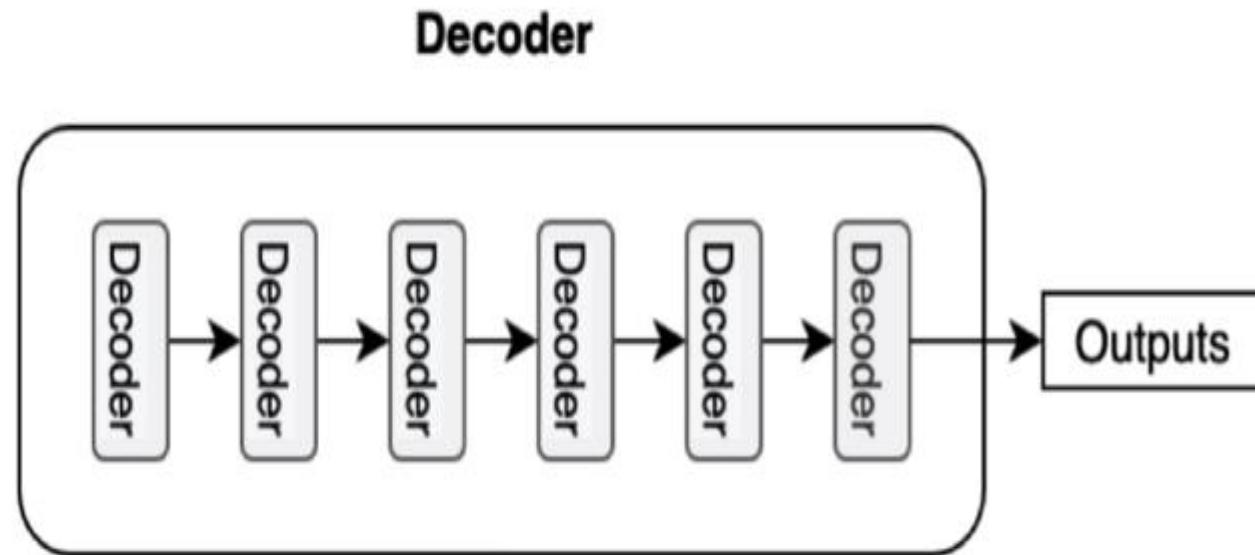
The transformer uses an encoder-decoder architecture. The encoder extracts features from an input sentence, and the decoder uses the features to produce an output sentence (translation).



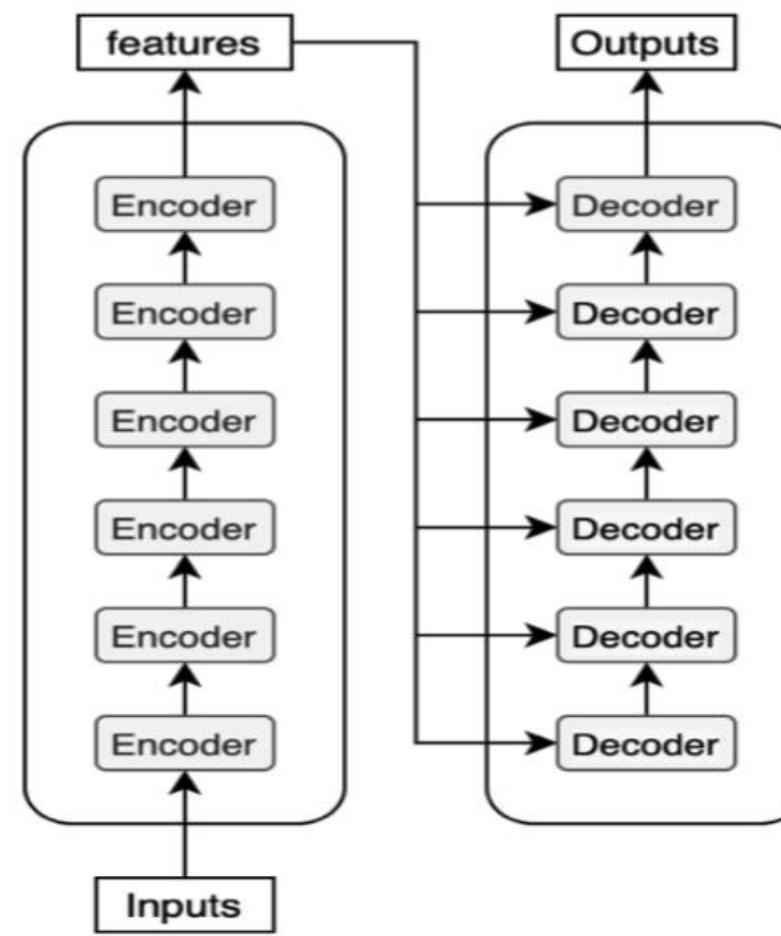
The encoder in the transformer consists of multiple encoder blocks. An input sentence goes through the encoder blocks, and the output of the last encoder block becomes the input features to the decoder.



The decoder also consists of multiple decoder blocks.

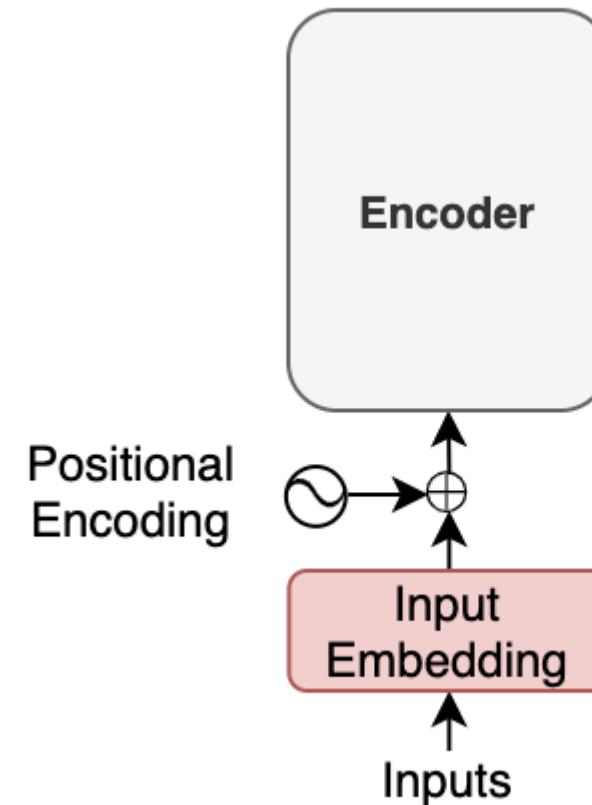


Each decoder block receives the features from the encoder. If we draw the encoder and the decoder vertically, the whole picture looks like the diagram from the paper.



Input Embedding and Positional Encoding

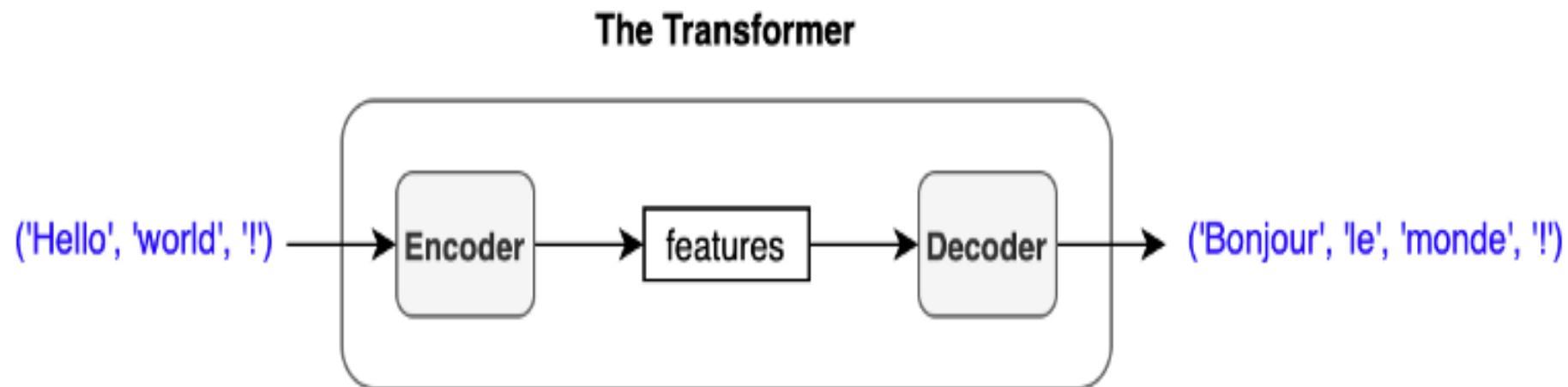
Like any neural translation model, we often tokenize an input sentence into distinct elements (tokens). A tokenized sentence is a fixed-length sequence.



- with the transformer, we inject **positional encoding** into each embedding so that the model can know word positions without recurrence.
- The input to the transformer is not the characters of the input text but a sequence of embedding vectors.
- Each vector represents the semantics and position of a token. The encoder performs linear algebra operations on those vectors to extract the contexts for each token from the entire sentence and enrich the embedding vectors with helpful information for the target task through the multiple encoder blocks.

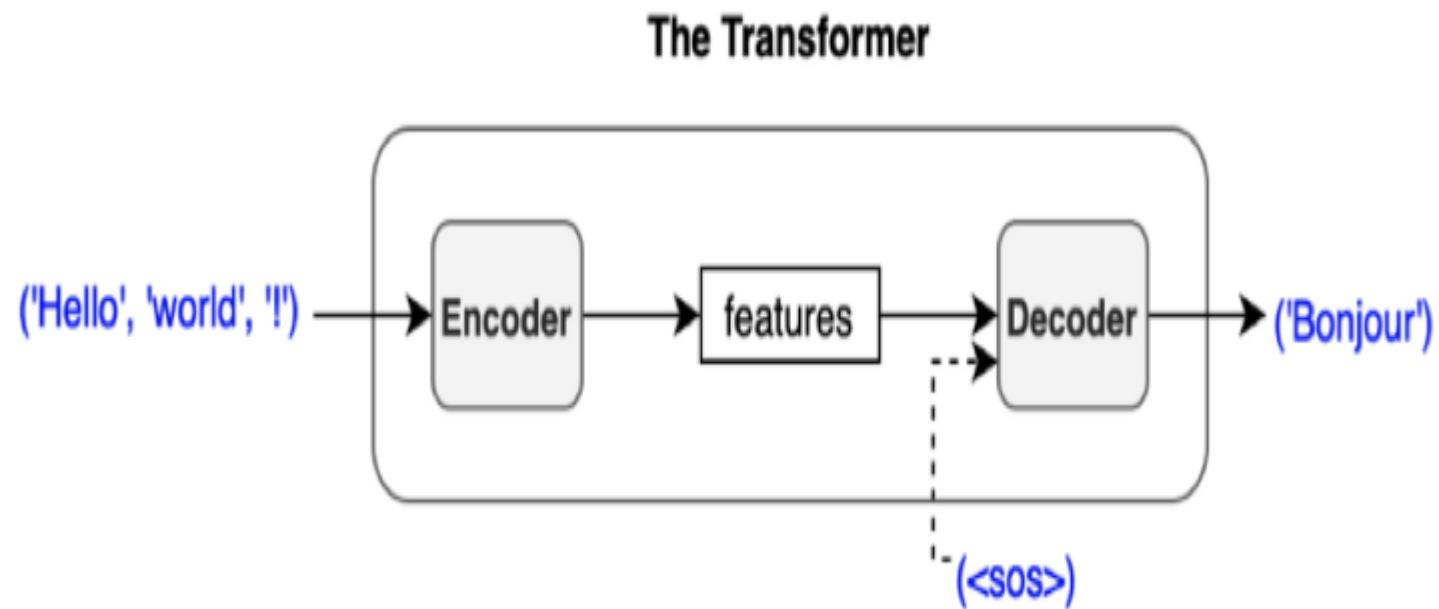
Softmax and Output Probabilities

The decoder uses input features from the encoder to generate an output sentence. The input features are nothing but enriched embedding vectors.



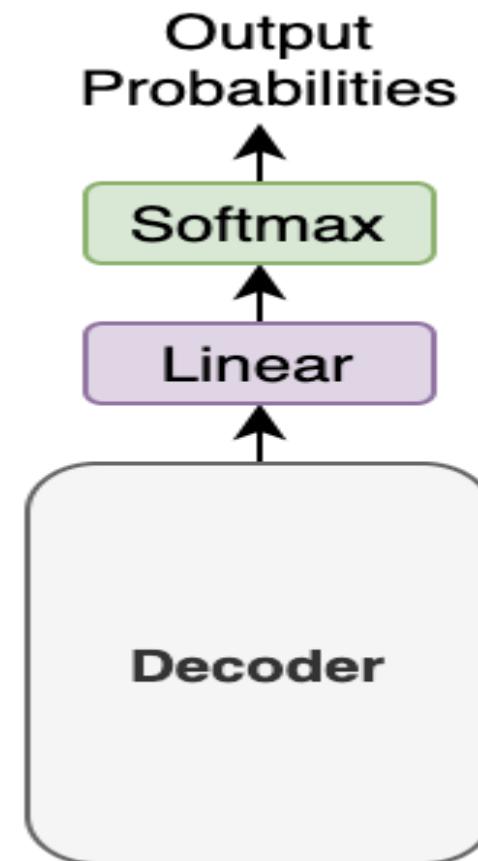
The decoder outputs one token at a time. An output token becomes the subsequent input to the decoder. In other words, a previous output from the decoder becomes the last part of the next input to the decoder. This kind of processing is called “**auto-regressive**”—a typical pattern for generating sequential outputs and not specific to the transformer. It allows a model to generate an output sentence of different lengths than the input.

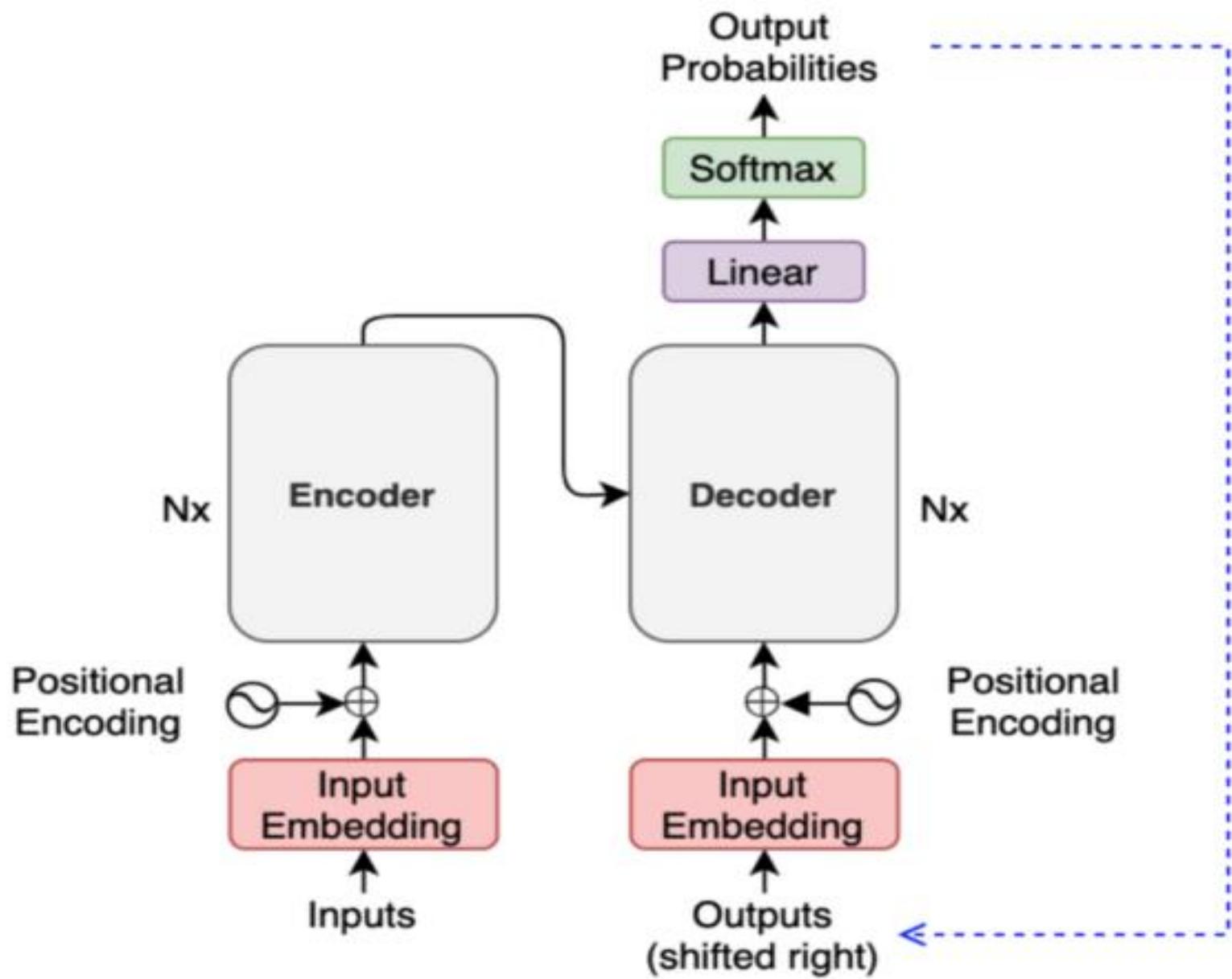
However, we have no previous output at the beginning of a translation. So, we pass the start-of-sentence marker <SOS> to the decoder to initiate the translation.



The decoder uses multiple decoder blocks to enrich <SOS> with the contextual information from the input features. In other words, the decoder transforms the embedding vector <SOS> into a vector containing information helpful to generating the first translated word (token).

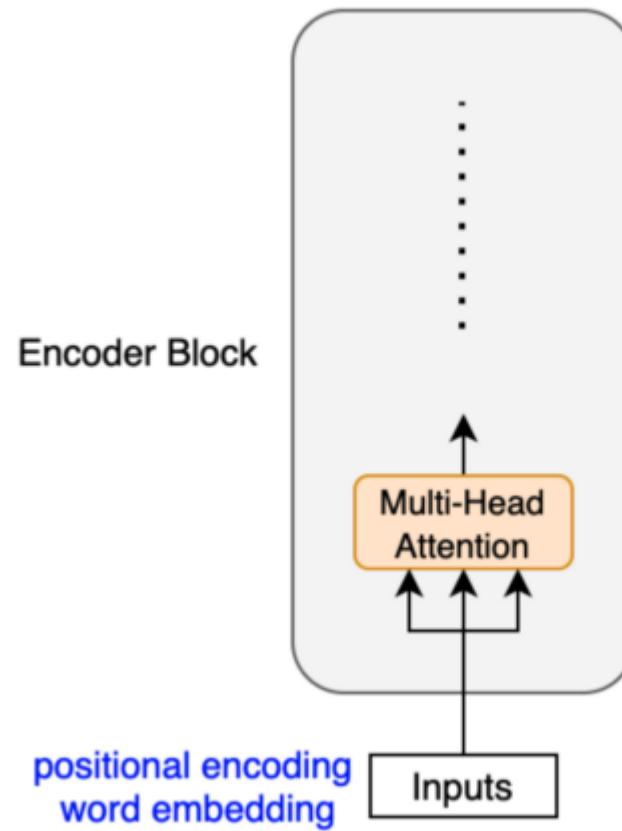
Then, the output vector from the decoder goes through a linear transformation that changes the dimension of the vector from the embedding vector size (512) into the size of vocabulary (say, 10,000). The softmax layer further converts the vector into 10,000 probabilities.



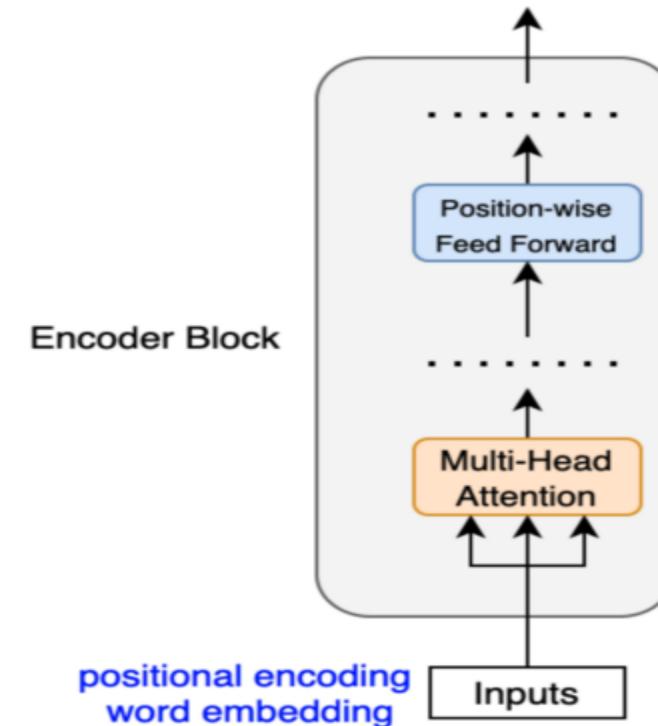


Encoder Block Internals

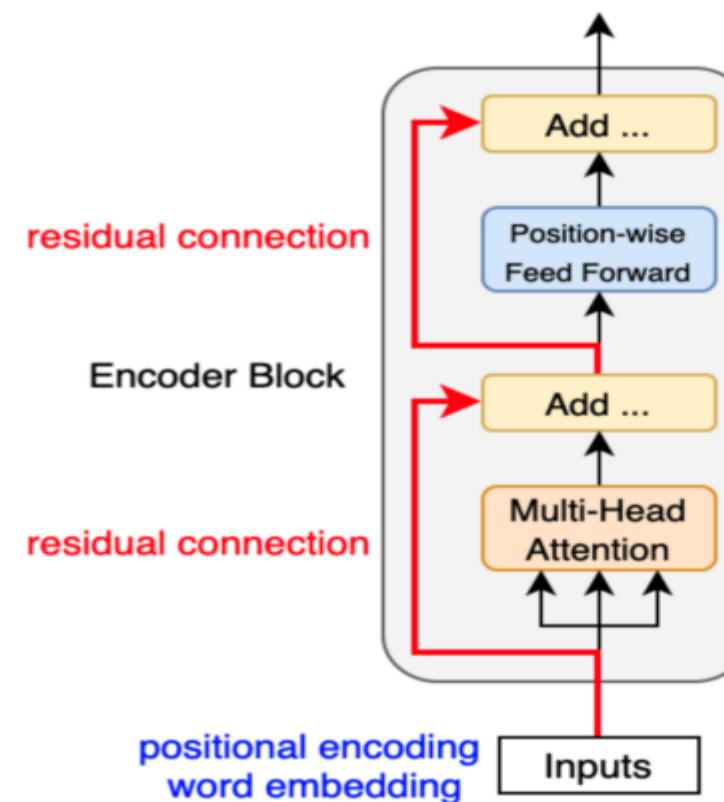
The encoder block uses the **self-attention mechanism** to enrich each token (embedding vector) with contextual information from the whole sentence.



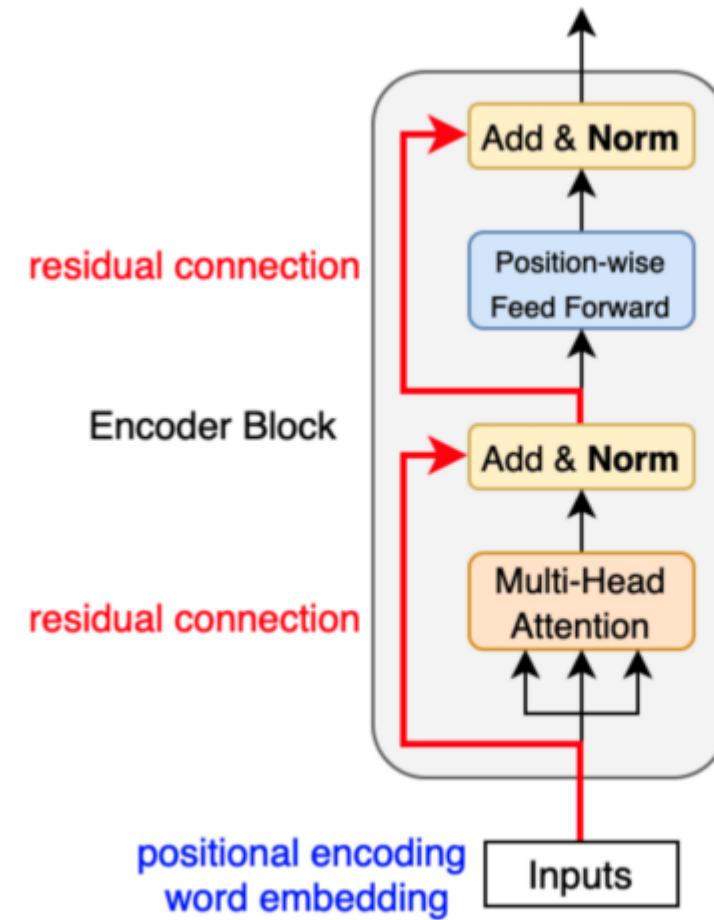
The **position-wise feed-forward network (FFN)** has a linear layer, ReLU, and another linear layer, which processes each embedding vector independently with identical weights. So, each embedding vector (with contextual information from the multi-head attention) goes through the position-wise feed-forward layer for further transformation.



Residual connections carry over the previous embeddings to the subsequent layers. As such, the encoder blocks enrich the embedding vectors with additional information obtained from the multi-head self-attention calculations and position-wise feed-forward networks.



After each residual connection, there is a **layer normalization**:

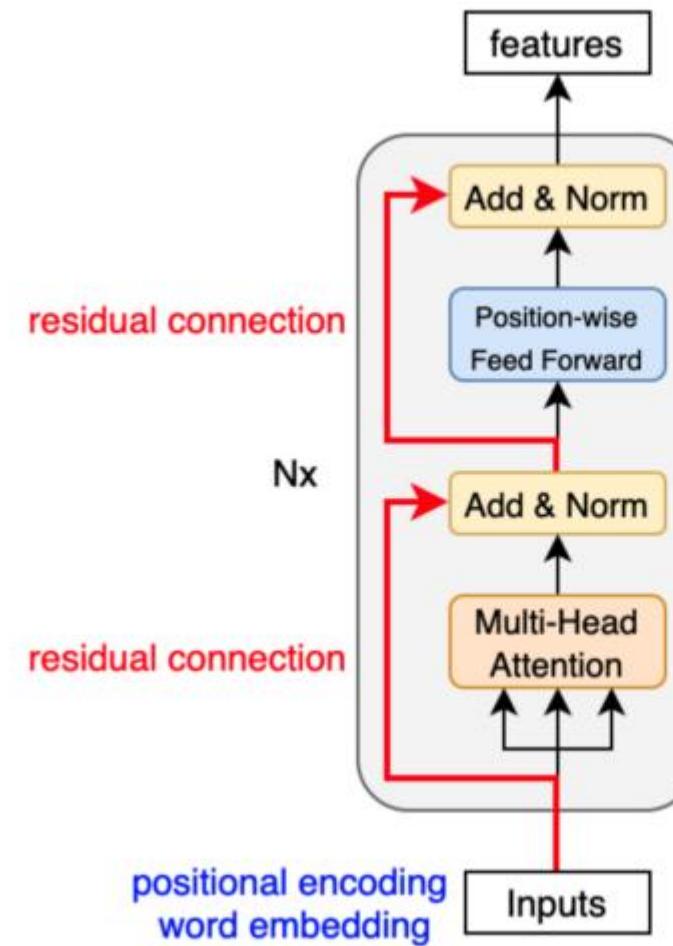


Like batch normalization, layer normalization aims to reduce the effect of covariant shift.

In other words, it prevents the mean and standard deviation of embedding vector elements from moving around, which makes training unstable and slow (i.e., we can't make the learning rate big enough).

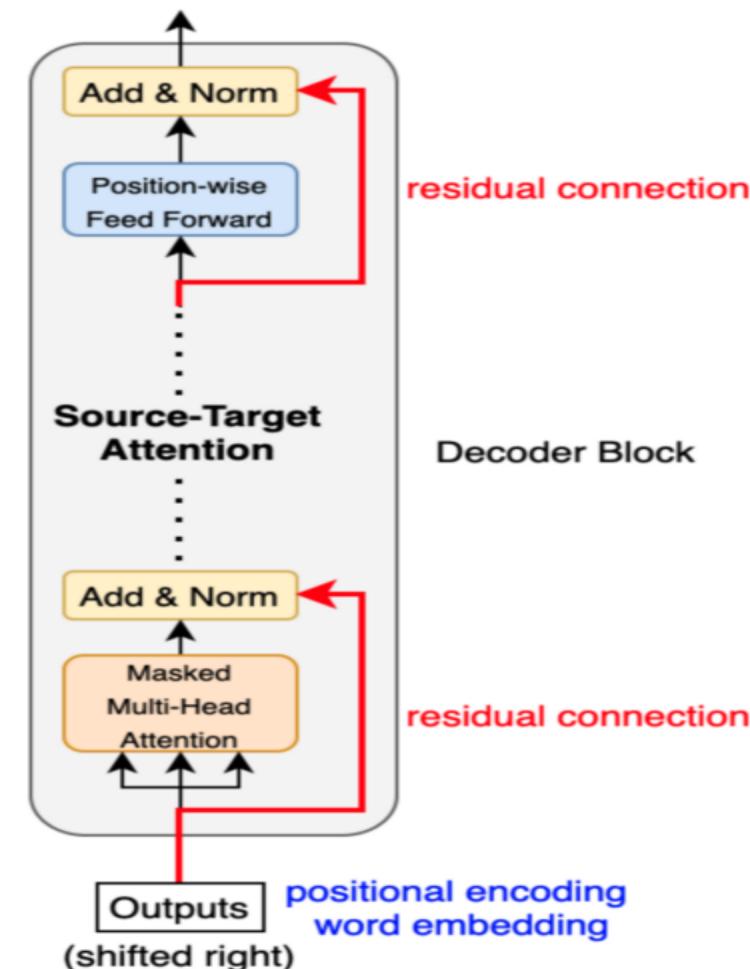
Unlike batch normalization, layer normalization works at each embedding vector (not at the batch level).

The transformer uses six stacked encoder blocks. The outputs from the last encoder block become the input features for the decoder.



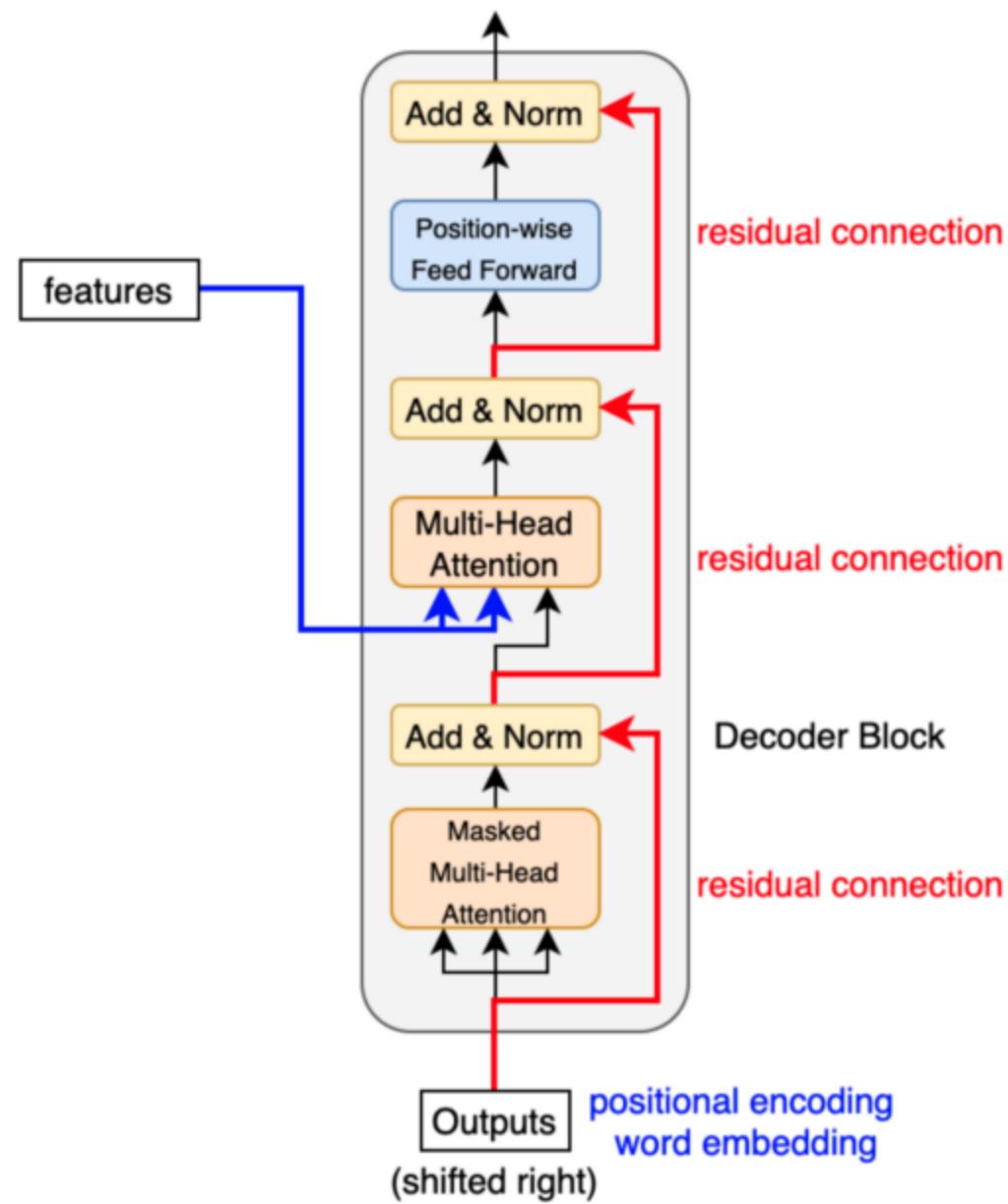
Decoder Block Internals

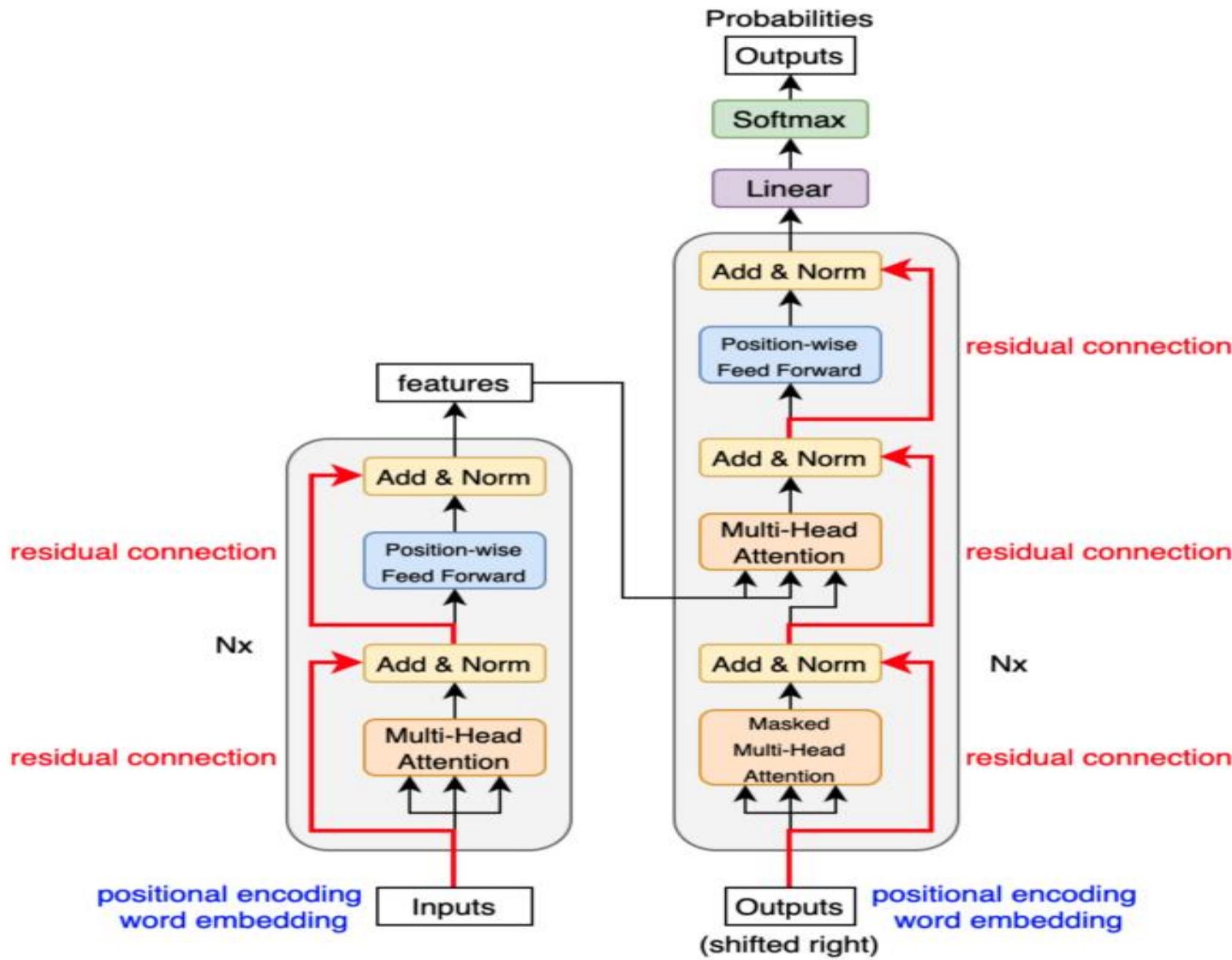
The decoder block is similar to the encoder block, except it calculates the source-target attention.

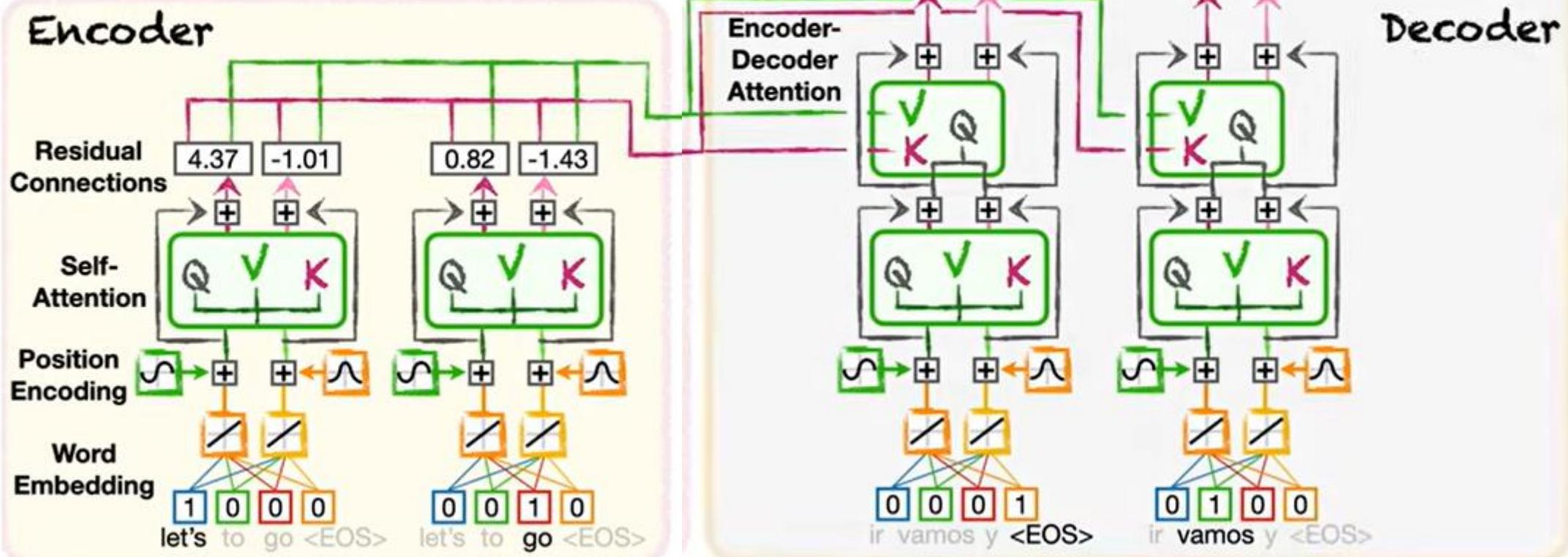


The input to the decoder is an output shifted right, which becomes a sequence of embeddings with positional encoding. So, we can think of the decoder block as another encoder generating enriched embeddings useful for translation outputs.

Masked multi-head attention means the multi-head attention receives inputs with masks so that the attention mechanism does not use information from the hidden (masked) positions. The paper mentions that they used the mask inside the attention calculation by setting attention scores to negative infinity (or a very large negative number). The softmax within the attention mechanisms effectively assigns zero probability to masked positions.







Prompt

Write an email apologizing to Sarah for the tragic gardening mishap. Explain how it happened.



Transformer LLM

Generation

Dear Sarah,
I'm writing to apologize for the incident last week. [...]

Prompt

Write an email apologizing to Sarah for the tragic gardening mishap. Explain how it happened.

**Transformer LLM****Generation**

#1

Dear

#2

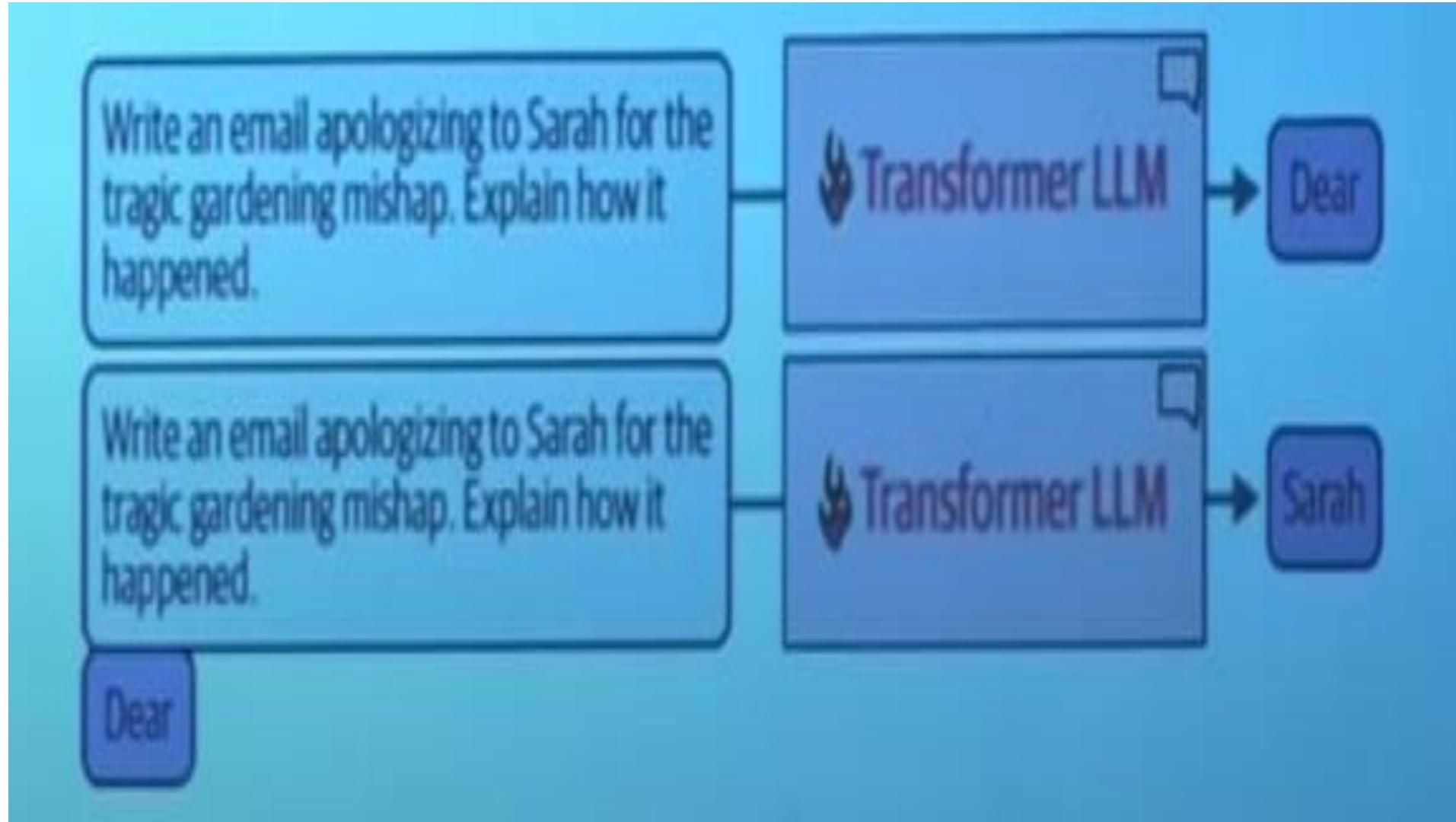
Sarah

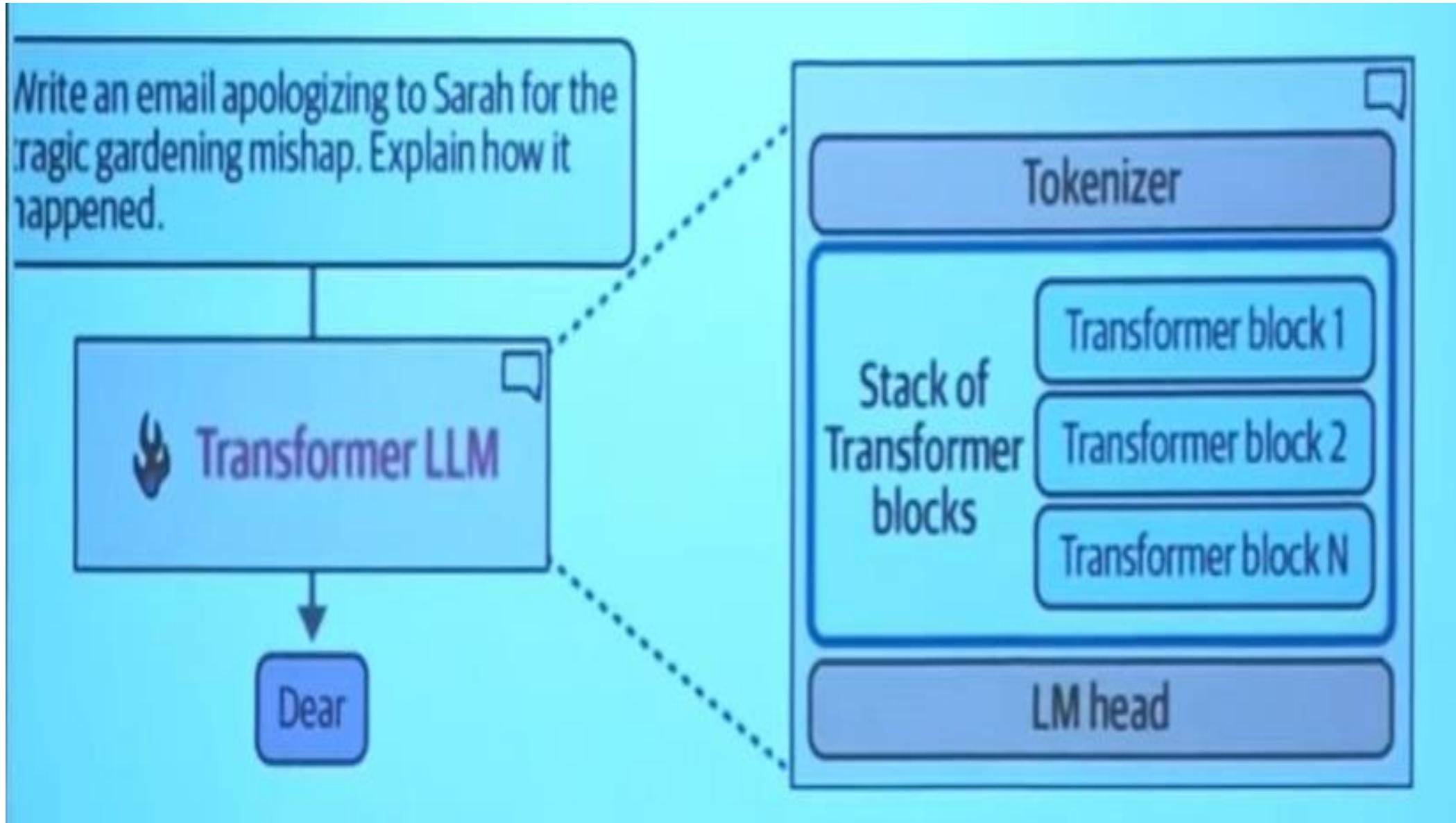
#3

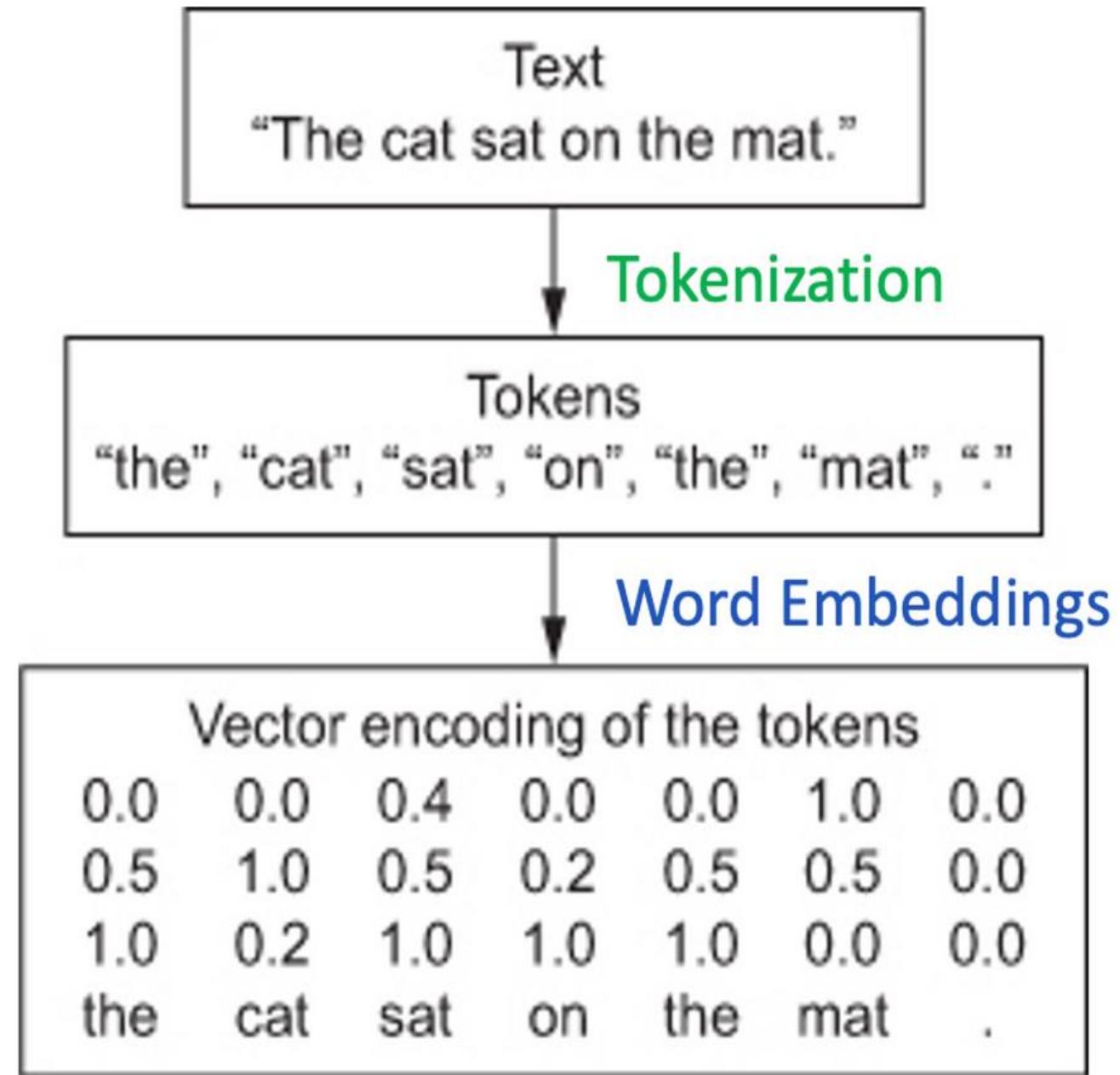
.

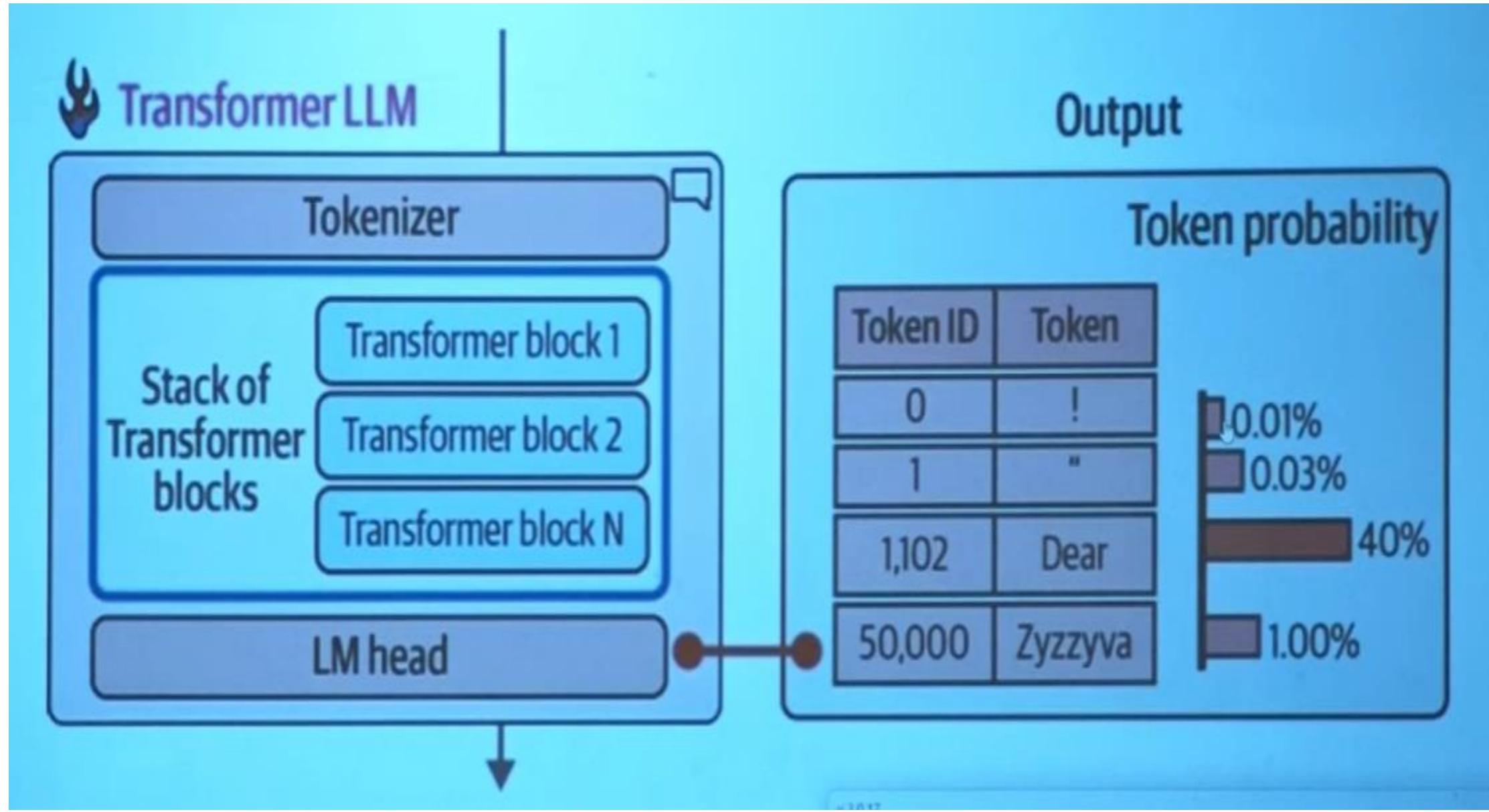
#4

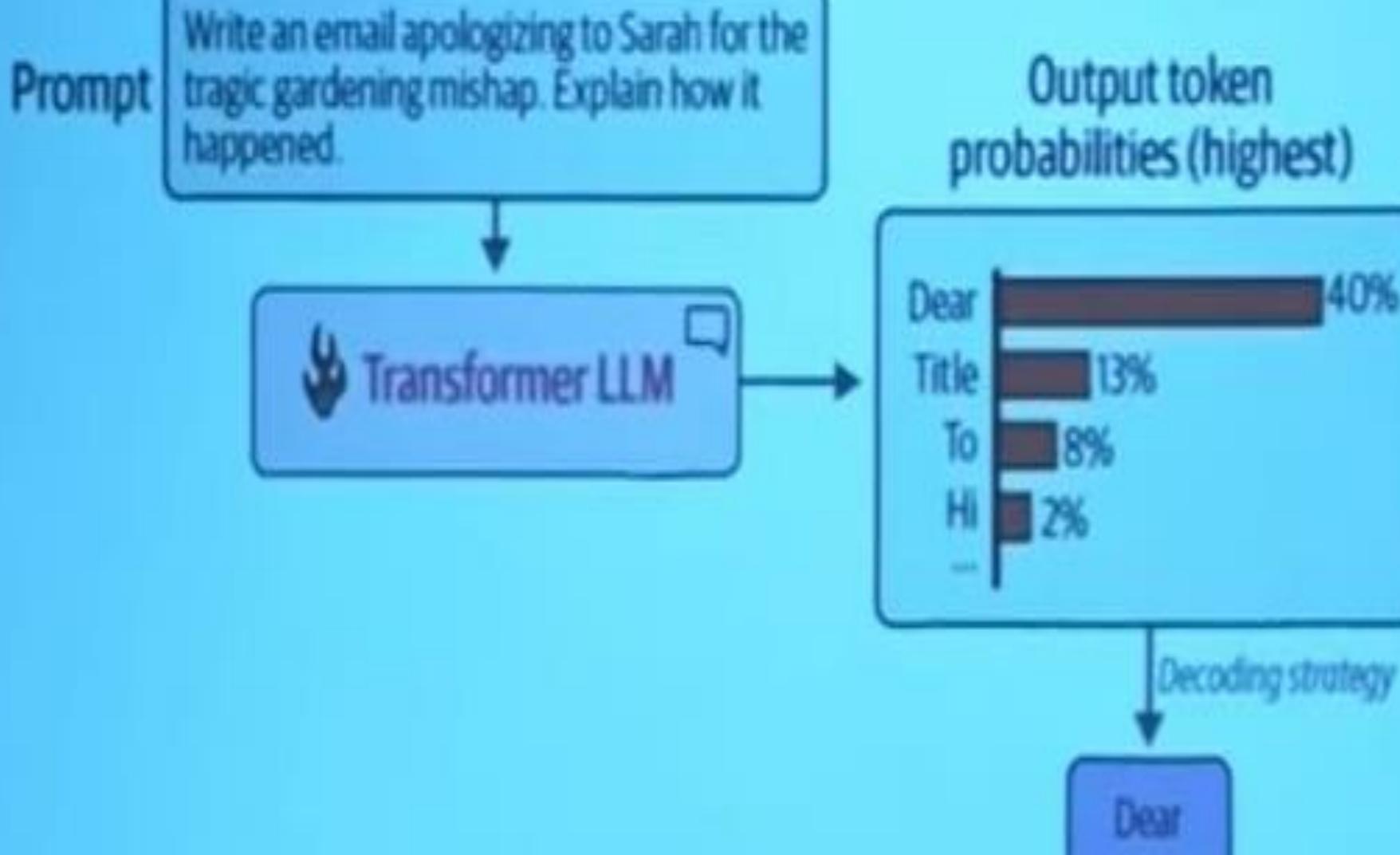
\n





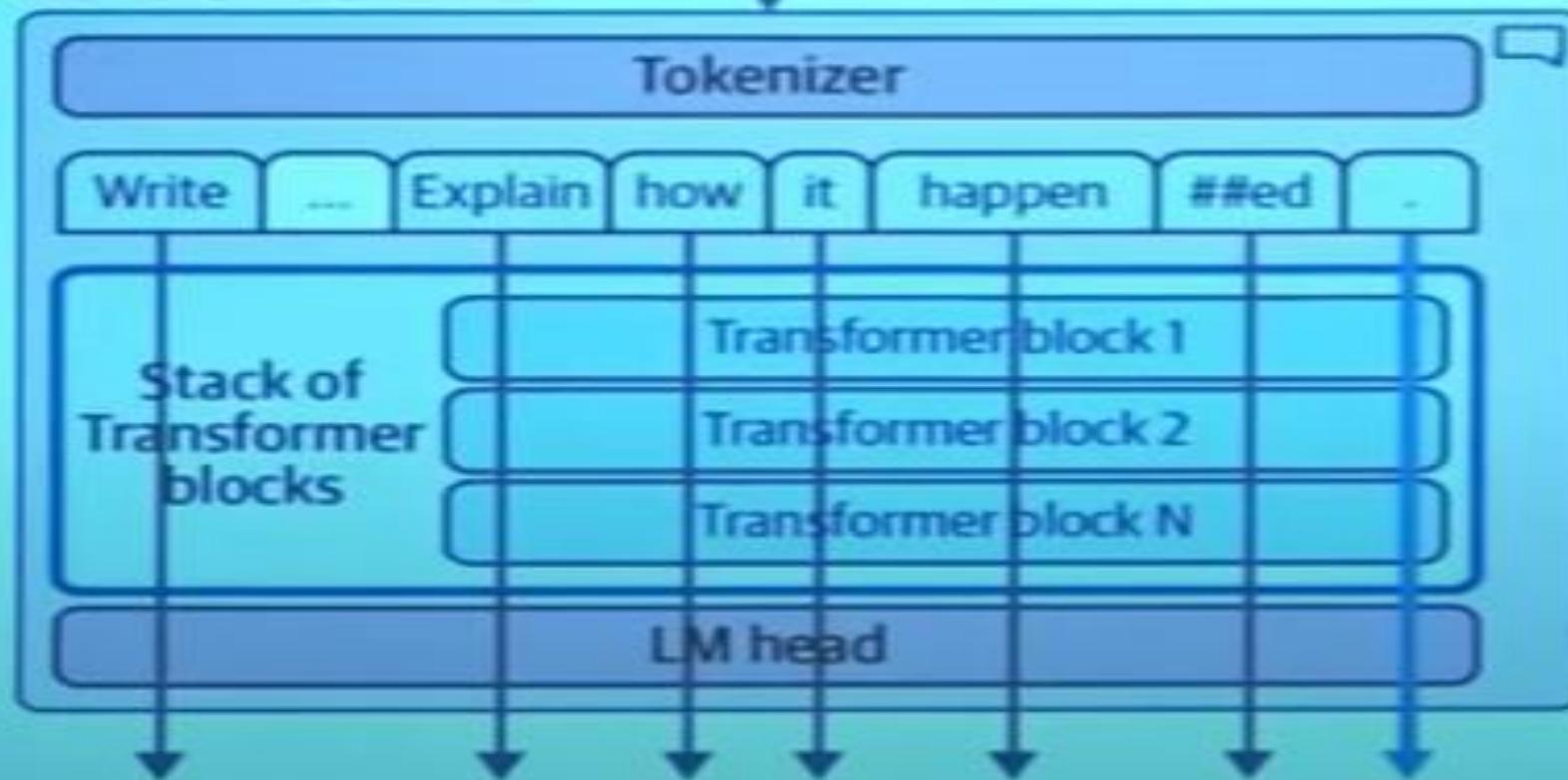






Write an email apologizing to Sarah for the tragic gardening mishap. Explain how it happened.

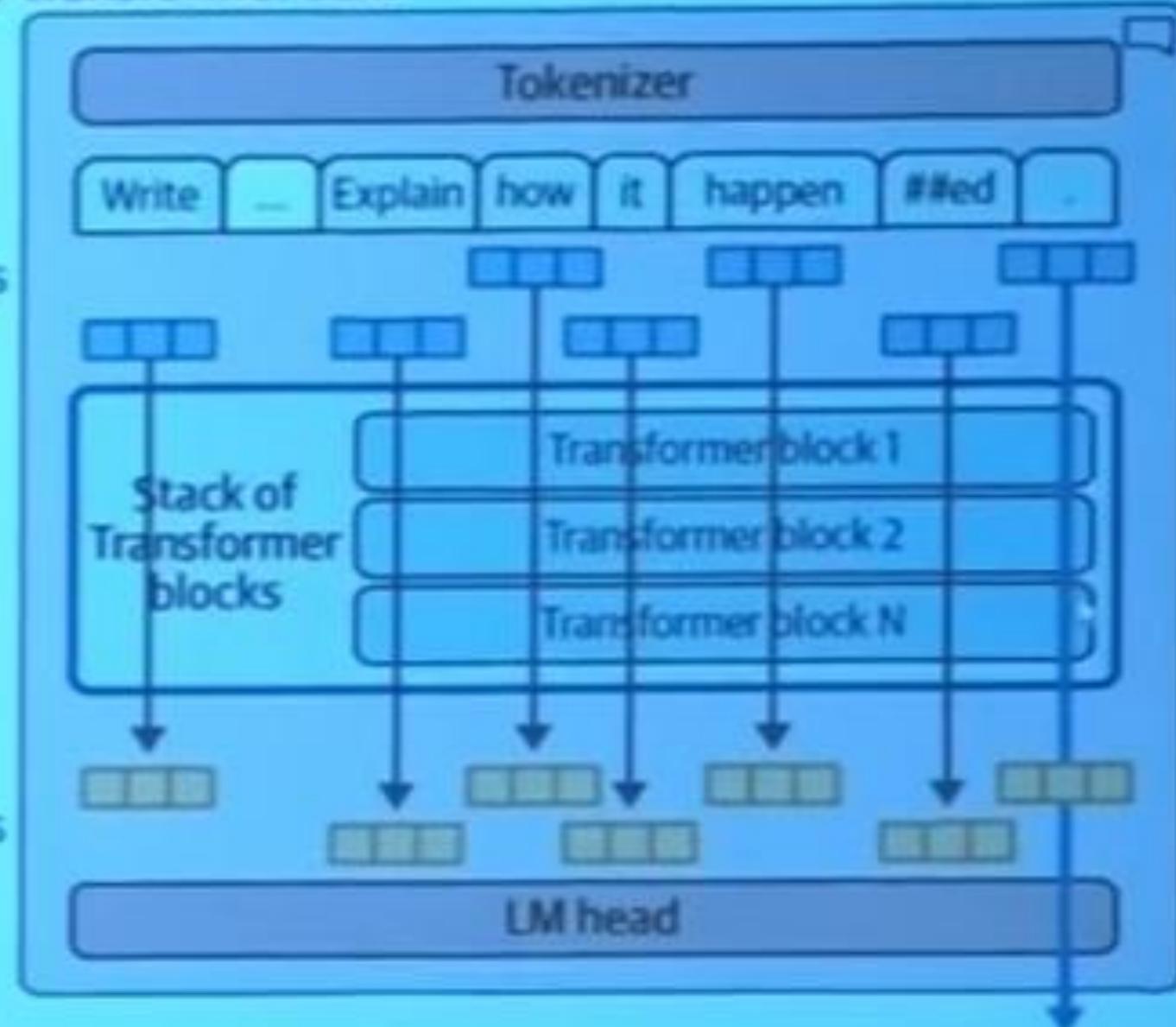
Transformer LLM





Transformer LLM

Embeddings
Output vectors





Transformer LLM

Embeddings

Output vectors

Tokenizer

Write ... Explain how it happen #filled ... Dear

Cached calculation

LM head

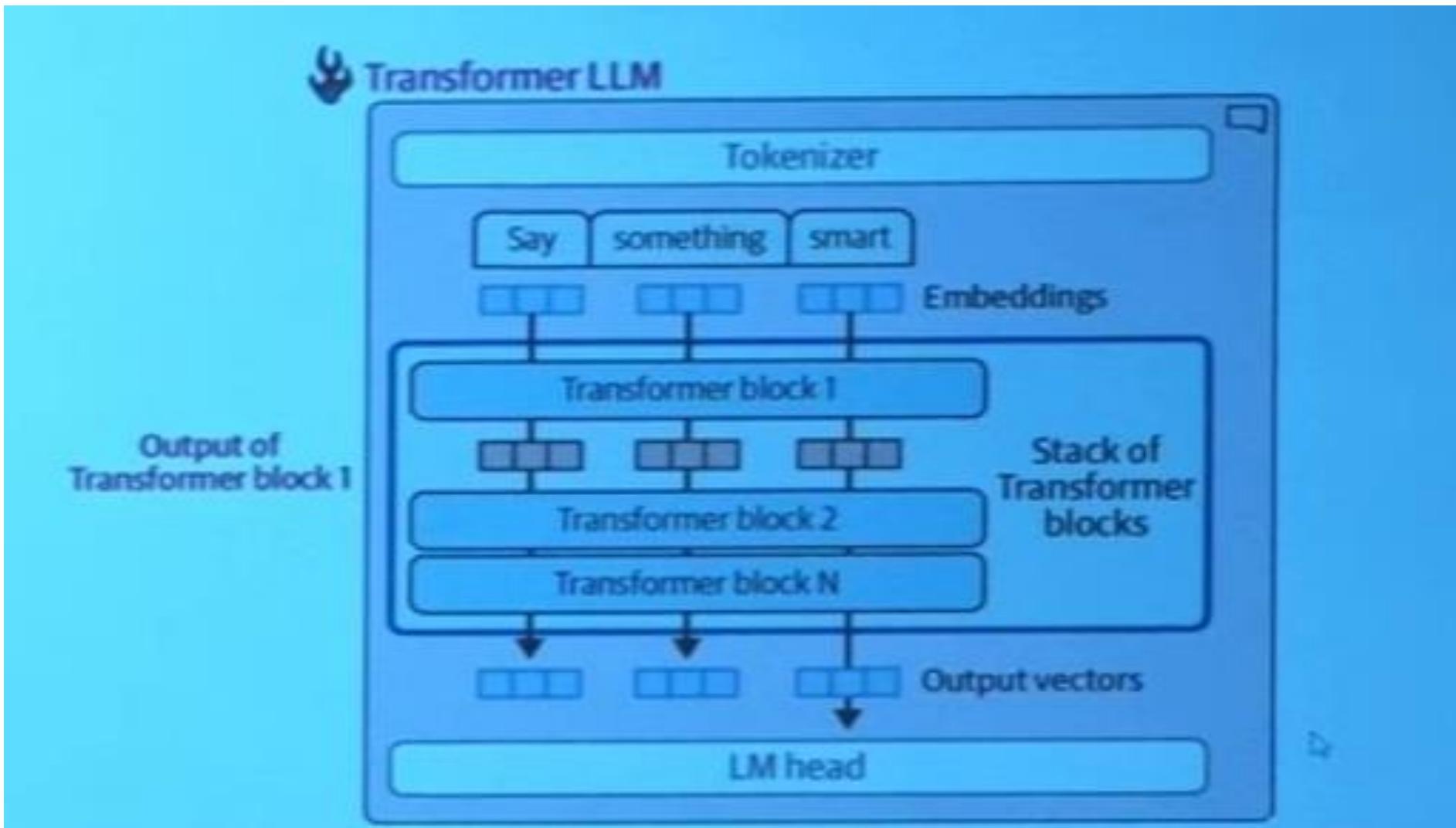
[UNK]

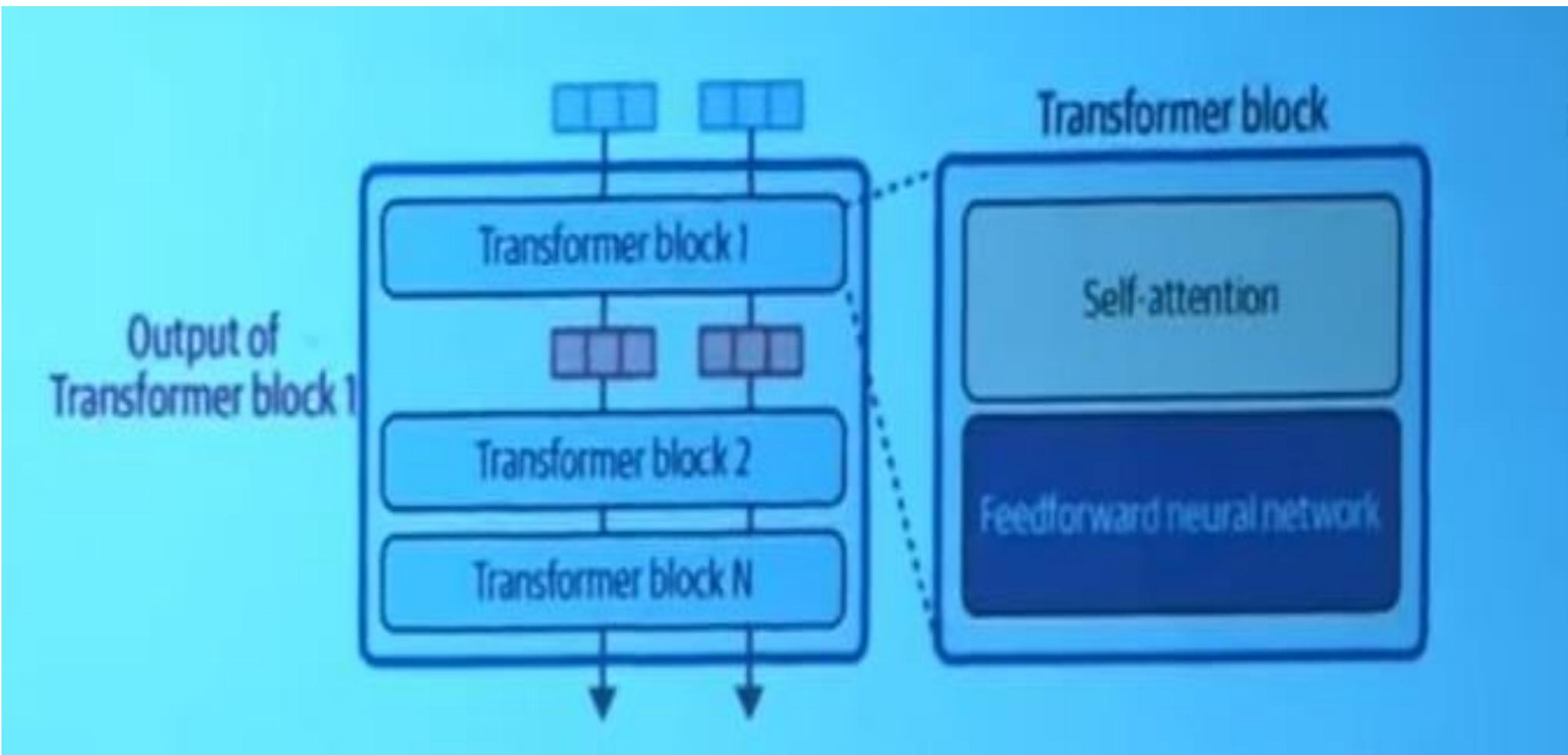
[UNK]

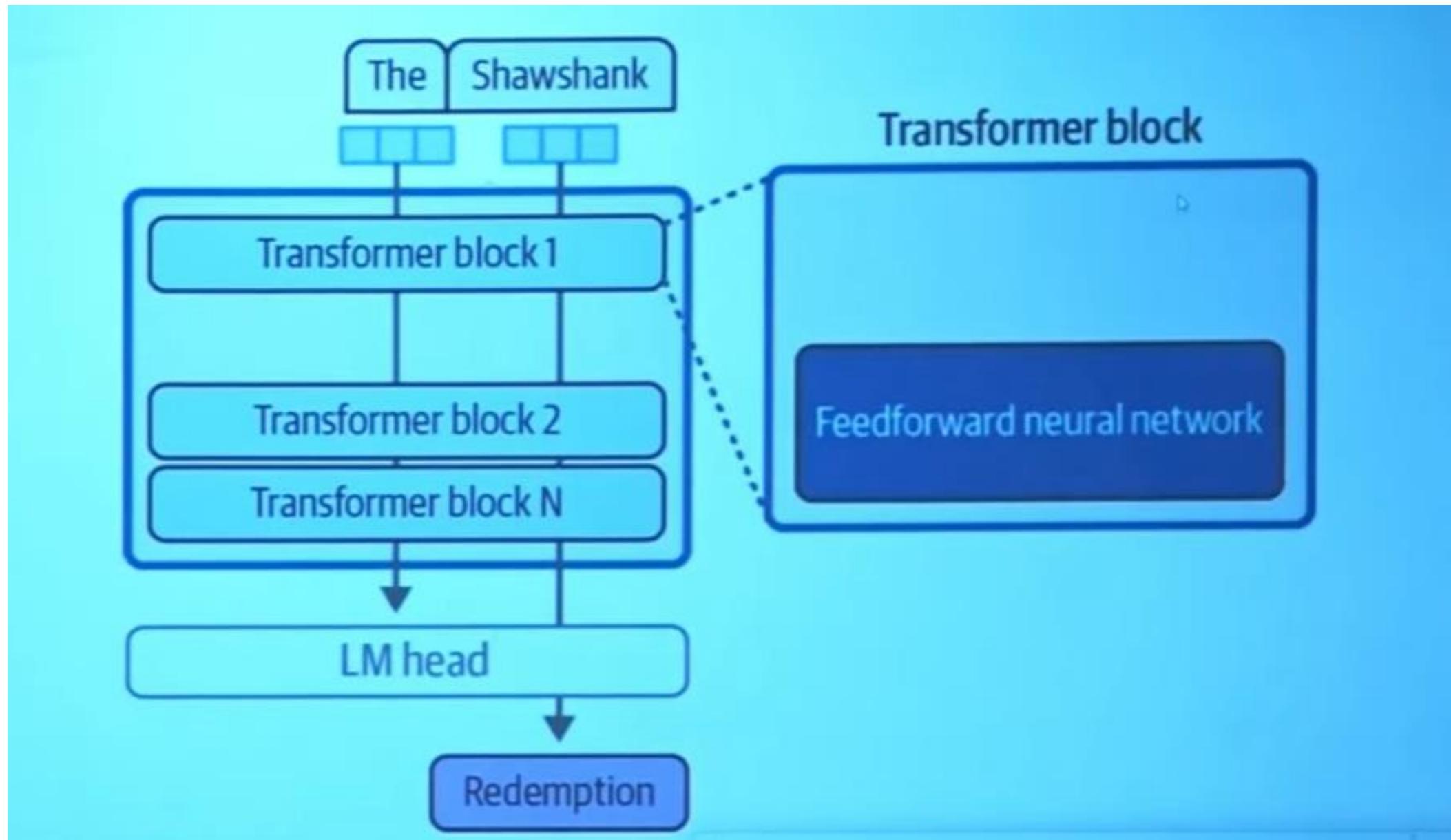
[UNK]

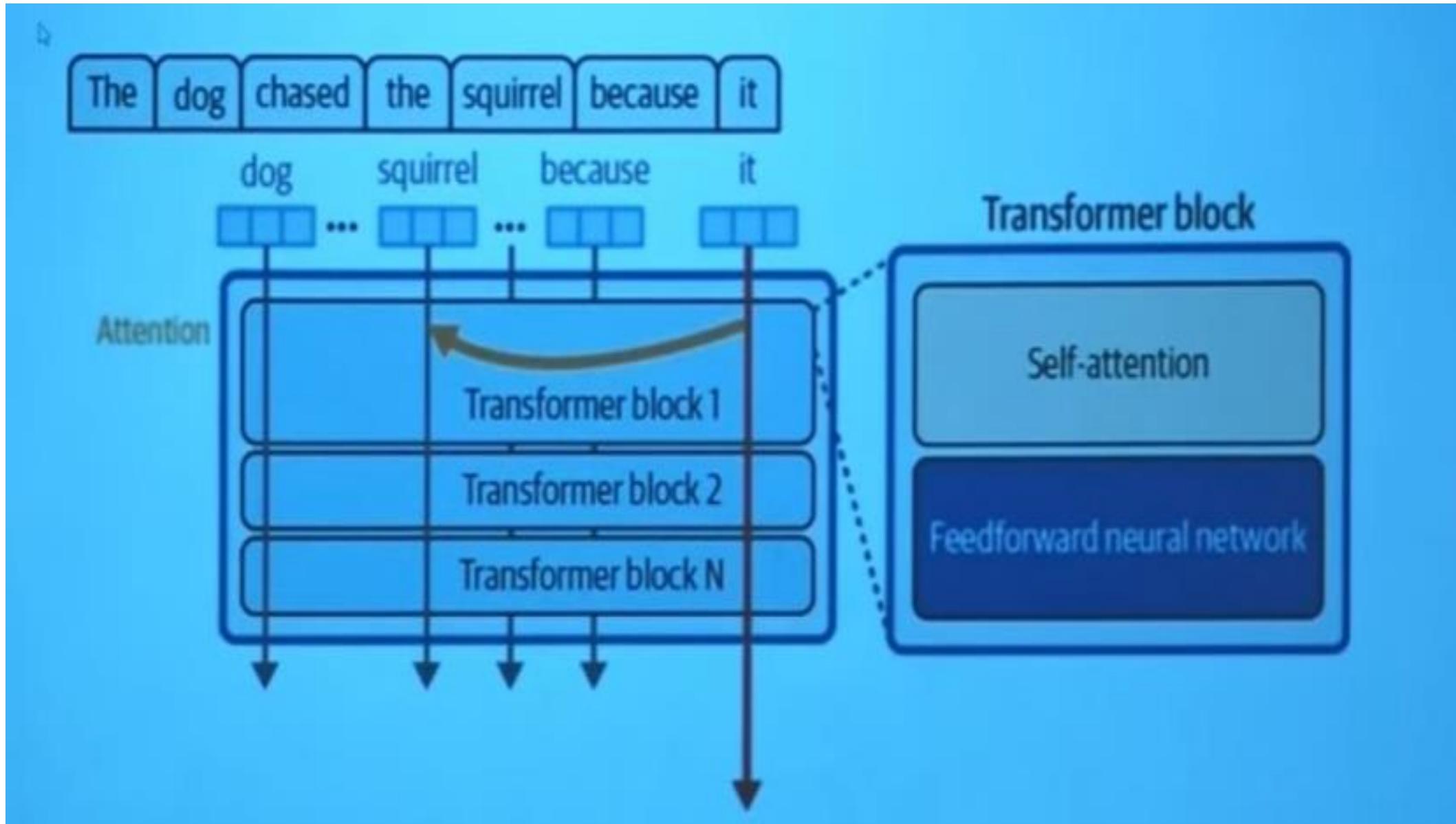
[UNK]

[UNK]









Step 1—Defining our Dataset

The dataset used for creating ChatGPT is 570 GB. On the other hand, for our purposes, we will be using a very small dataset to perform numerical calculations visually.

Dataset (corpus)

I drink and I know things.

When you play the game of thrones, you win or you die.

The true enemy won't wait out the storm, He brings the storm.

Our entire dataset containing only three sentences

Step 2— Finding Vocab Size

The vocabulary size determines the total number of **unique words** in our dataset. It can be calculated using the below formula, where N is the total number of words in our dataset.

$$\text{vocab size} = \text{count}(\text{set}(N))$$

vocab_size formula where N is total number of words

In order to find N, we need to break our dataset into individual words.

Dataset (Corpus)

I drink and I know things.

When you play the game of thrones, you win or you die.

The true enemy won't wait out the storm, He brings the storm.

$$\rightarrow N = [\text{I, drink, and, I, Know, things,} \\ \text{When, you, play, the, game, of, thrones, you, win, or, you, die,} \\ \text{The, true, enemy, won't, wait, out, the, storm, He, brings, the, storm}]$$

calculating variable N

After obtaining N, we perform a set operation to remove duplicates, and then we can count the unique words to determine the vocabulary size.

$$\text{vocab size} = \text{count}(\text{set}(N))$$

↳ set (I, drink, and, I, Know, things,
When, you, play, the, game, of, thrones, you, win, or, you, die,
The, true, enemy, won't, wait, out, the, storm, He, brings, the, storm)

↳ count (I, drink, and, Know, things, When, you, play, the, game, of,
thrones, win, or, die, true, enemy, won't, wait, out, storm, He,
brings)

↳ = 23

finding vocab size

Therefore, the vocabulary size is **23**, as there are 23 unique words in our dataset.

Step 3 — Encoding

Now, we need to assign a unique number to each unique word.

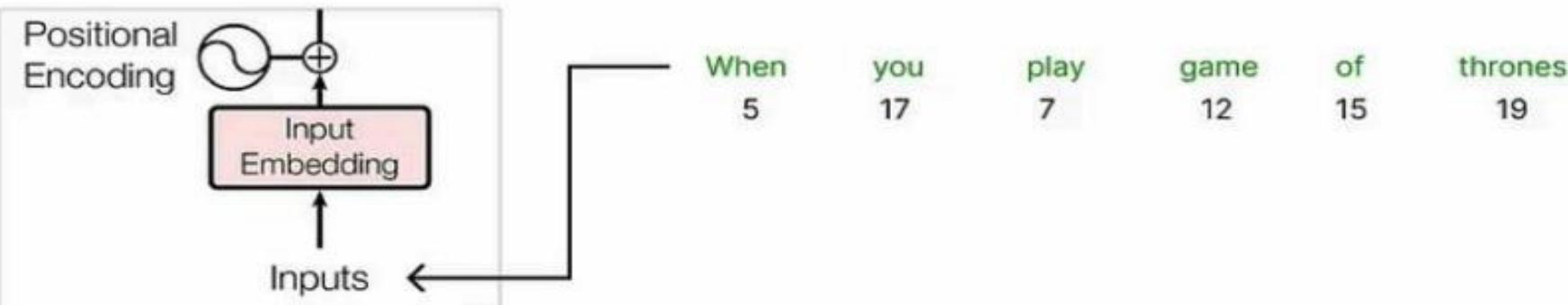
1	2	3	4	5	6	7	8	9	10	11	12
I	drink	things	Know	When	won't	play	out	true	storm	brings	game
13	14	15	16	17	18	19	20	21	22	23	
the	win	of	enemy	you	wait	thrones	and	or	die	He	

encoding our unique words

Step 4 — Calculating Embedding

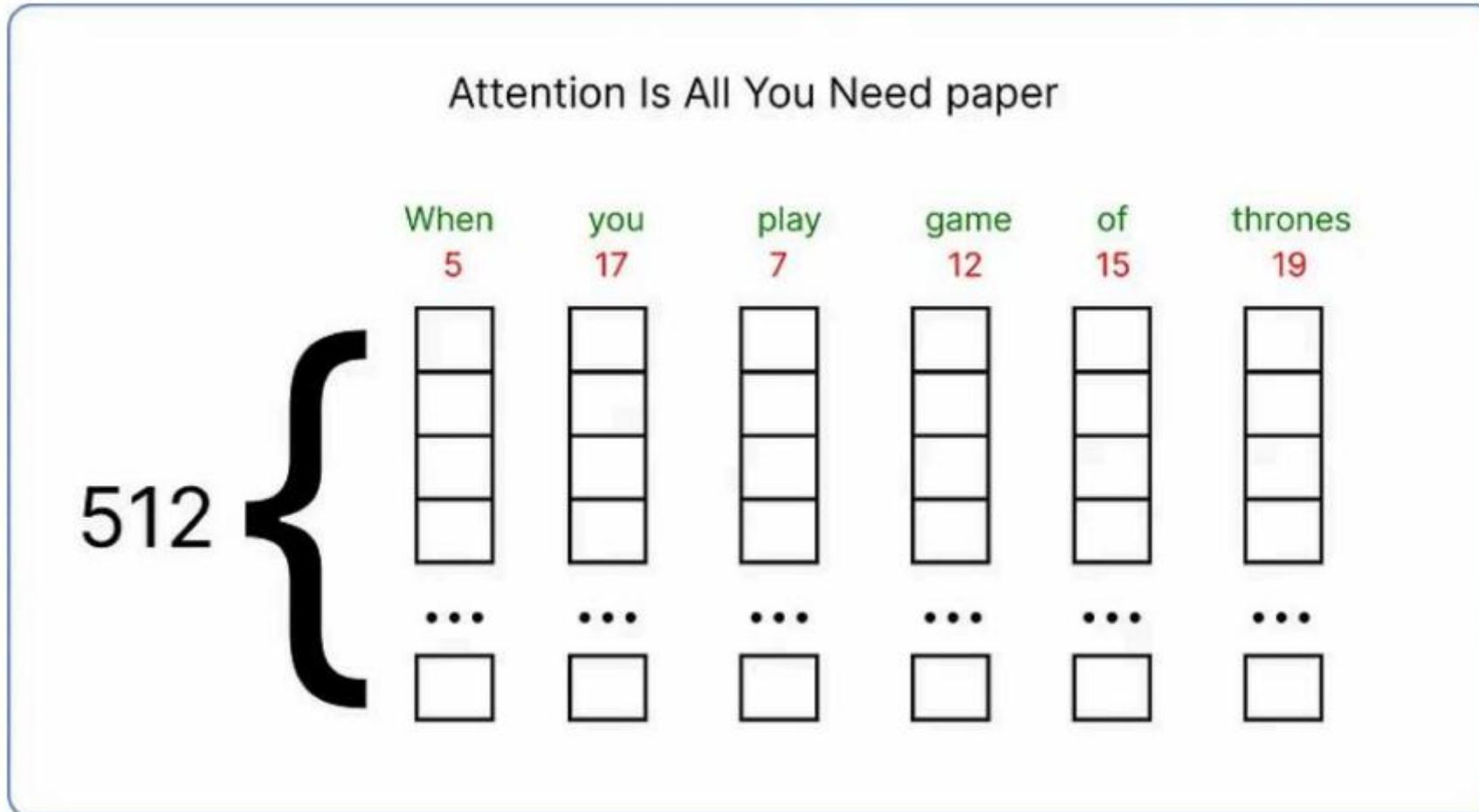
Let's select a sentence from our corpus that will be processed in our transformer architecture.

1	2	3	4	5	6	7	8	9	10	11	12
I	drink	things	Know	When	won't	play	out	true	storm	brings	game
13	14	15	16	17	18	19	20	21	22	23	
the	win	of	enemy	you	wait	thrones	and	or	die	He	



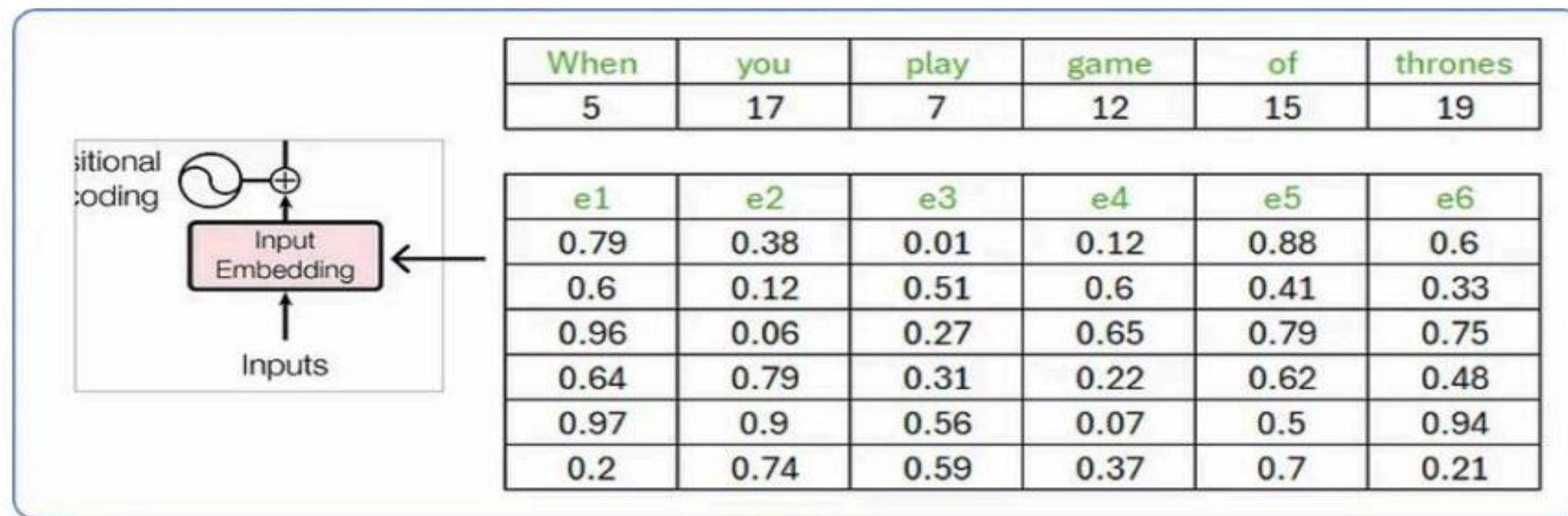
Input sentence for transformer

We have selected our input, and we need to find an embedding vector for it. The original paper uses a **512-dimensional embedding vector** for each input word.



Original Paper uses 512 dimension vector

Since, for our case, we need to work with a smaller dimension of embedding vector to visualize how the calculation is taking place. So, we will be using a dimension of 6 for the embedding vector.



Embedding vectors of our input

Step 5 — Calculating Positional Embedding

Now we need to find positional embeddings for our input. There are two formulas for positional embedding depending on the position of the i th value of that embedding vector for each word.

Embedding vector for any word

even position
odd position
even position
odd position
even position
...

For even position

$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

For odd position

$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

Positional Embedding formula

When
5

i	e1	Position	Formula	p1
0	0.79	Even	$\sin(0/10000^{(2*0/6)})$	0
1	0.6	Odd	$\cos(0/10000^{(2*1/6)})$	1
2	0.96	Even	$\sin(0/10000^{(2*2/6)})$	0
3	0.64	Odd	$\cos(0/10000^{(2*3/6)})$	1
4	0.97	Even	$\sin(0/10000^{(2*4/6)})$	0
5	0.2	Odd	$\cos(0/10000^{(2*5/6)})$	1

d (dim) 6

POS 0

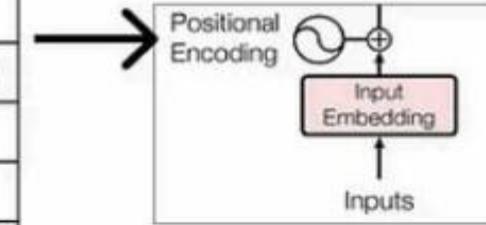
Positional Embedding for word: When

Similarly, we can calculate positional embedding for all the words in our input sentence.

When	you	play	game	of	thrones
5	17	7	12	15	19

i	p1	p2	p3	p4	p5	p6
0	0	0.8415	0.9093	0.1411	-0.7568	-0.9589
1	1	0.0464	0.9957	0.1388	0.1846	0.9732
2	0	0.0022	0.0043	0.0065	0.0086	0.0108
3	1	0.0001	1	0.0003	0.0004	1
4	0	0	0	0	0	0
5	1	0	1	0	0	1

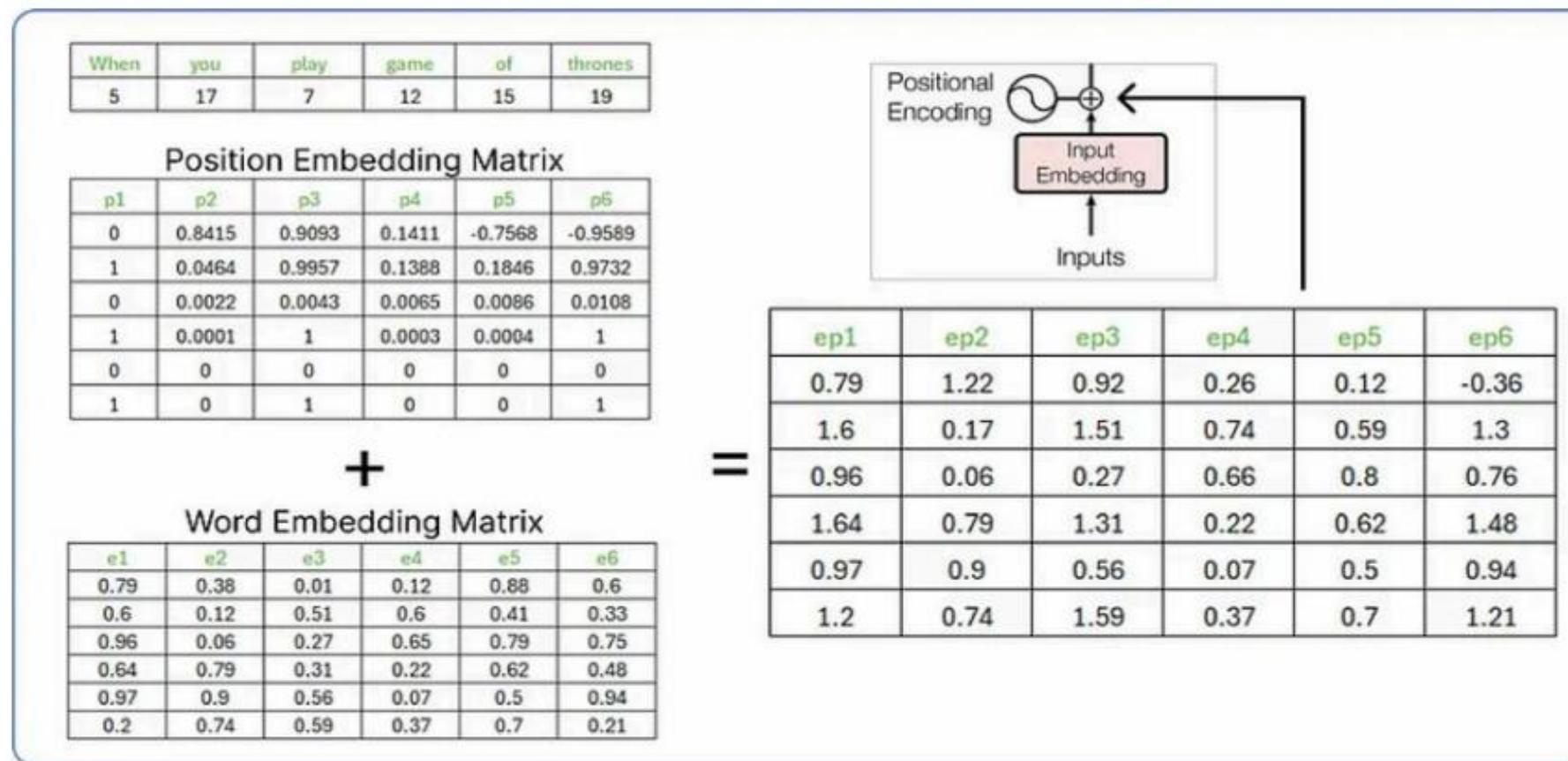
d (dim)	6	6	6	6	6	6
POS	0	1	2	3	4	5



Calculating Positional Embeddings of our input (**The calculated values are rounded**)

Step 6 — Concatenating Positional and Word Embeddings

After calculating positional embedding, we need to add word embeddings and positional embeddings.



concatenation step

This resultant matrix from combining both matrices (Word embedding matrix and positional embedding matrix) will be considered as an input to the encoder part.