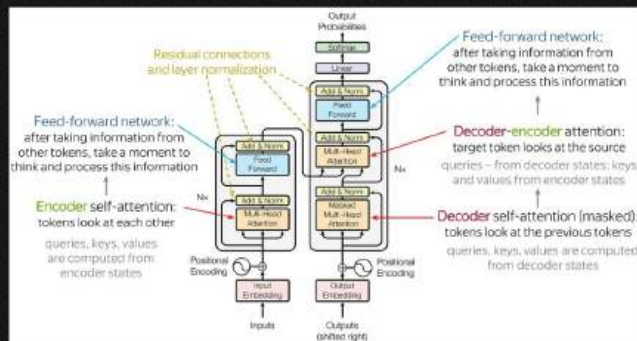# Introduction To Transformers

1) RNN/LSTM/GRU RNN
2) Encoder Decoder Architecture ⎫
3) ATTENTION MECHANISM       ⎬
4) TRANSFORMERS               ⎭

   ① Why Transformers?
   ② Architecture of Transformers?
   ③ SELF ATTENTION → $Q, K, V$
   ④ Positional Encoding
   ⑤ Multi Head ATTENTION
   ⑥ Combining the working of Transformers

## Architecture



Generative AI → LLM, Multimodel

BERT, GPT ←
      ↓
Open AI → Chat GPT
      ↓
     GPT-4o

## ① What And Why → Transformers

**Transformers in natural language processing (NLP) are a type of deep learning model that use self-attention mechanisms to analyze and process natural language data. They are encoder-decoder models that can be used for many applications, including machine translation.** ⇒ Seq2Seq Task

Eg: Language Translation ⟶ Google Translation

   English ⟶ French
   i/p ⇒ Many ⟶ O/p: Many . { Length of the sentence }

Sentence Length ↑↑          Bleau Score ↓↓

                              Length Sentence ↑↑.

**Encoder — Decoder**

Encoder                                    Decoder



Sent 1  $x_1$        $x_2$        $x_3$
     $t=1$        $t=2$        $t=3$

Additional Context ⟶ Decoder

Long Sentence. Accuracy ⇑⇑.

Hello, What's Up

## Attention Mechanism ↓

① Parallely We Cannot Send all the words in a Sentence ⟶ **Scalable**

DATASET ⟶ Huge ⟶ Scalable With Respect to Training.

TRANSFORMERS ✳ LSTM RNN
↳ Self Attention Module ← All the words will be parallely sent to encoder.

↓
Positional Encoding

Transformer ⇑⇑ DATASET ⟶ Amazing SOTA ← NLP

Transfer Learning ⟶ MultiModal Task ⟶ NLP + Image ←

Transformers ÷ AI Space ⟶ SOTA Model ⟶ } LLM's

BERT GPT ⟶ Transfer Learning ⟶ SOTA Models ⟶ DALLE } Generative AI
⇑
Train Huge Data

## ② Contextual Embedding → Self Attention

Contextual Vector

Eg: My name is KRISH And I play CRICKET

Word2vec

| Embedding Layer |
|---|

↓ ↓ ↓ ↓ ↓ ↓

□ □ □ □ □ □

| Contextual vector Embedding |
|---|

↓ ↓ ↓ ↓ ↓ ↓↓

□

## ② Basic Transformer Architecture    {Seq2Seq Task} → Language Translation {Eng → French}

I/p →

| The TRANSFORMER |
|---|

→ O/p.



↑ 6 Encoder

Encoder

Comment vas-tu?

| Encoder | → | Decoder |
|---|---|---|

HOW ARE YOU?

Encoder

O/p Encoder Decoder Architecture

← 6 Decoders

Encoder

Decoder

| Feed Forward Neural N/W |
|---|
| Self - Attention |

{ LSTM RNN

| Feed Forward |
|---|
| Encoder - Decoder Attention |
| Self - Attention |

**Encoder**

FEED FORWARD

$Z_1$ [ | ]  $Z_2$ [ | ]  $Z_3$ [ | ]  } ⇒ Contextual vector ←

SELF ATTENTION

Vectors | } 
⇩
Embedding Layer

→ [ | | ] → [ | | ] → [ | | ]
How ⤴   ARE ⤴   You ⤴
=         =

{ All the words will be passed parallely }.

**Encoder 2**

[ | | ]  [ | | ]  [ | | ]

**Encoder**

( N/w )  ( N/w )  ( N/w . )

[ | | ]  [ | | ]  [ | | ]

SELF ATTENTION

[ | | ]  [ | | ]  [ | | ]
How      ARE      You

# Self Attention At a Higher Level

Eg: The cat sat on the mat, the cat lay on the rug.

Word Embedding

↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓

```
┌─────────────────────────────────────────┐
│            SELF ATTENTION                 │
└─────────────────────────────────────────┘
```

▭ ▭ ▭ ▭ ▭ ▭

⇓

▭

The

The..

→ [Cat] - - - - - [Cat] → [ | | ]

Sat

1 → [Sat]                      { Contextual Embedding }

on

Rank 2 → On

the

3 → [the]

mat          mat → [ — | — | — ]

cat ───────────

## Self Attention In Detail

① To create 3 vectors from each of the encoder i/p: Query vector,

Key vector, value vector. ⇒ Contextual Embedding

Topis

Eg: YT ⇒ Search keywords ⇒ { Query }.

Query ⟶ { Key } ⟨ Tags
                    Description    ⇒ Value ⇒ O/P Video
                    Title

② Second step in calculating self attention is to calculate the score.

The Score determines how much focus to place on the other part of the sentence.

I/P

| How | (1×3) | GOOD | (1×3) |

(3×3) { Dependencies of each And every word }

Embedding   $X_1$ [    ]   $X_2$ [    ]

$W^Q$

Queries   $q_1$ [    ]   $q_2$ [    ]

Keys   $K_1$ [    ]   $k_2$ [    ]   $W^K$ .   Context of words

Values   $V_1$ [    ]   $V_2$ [    ]   $W^V$

Score   $q_1 \cdot K_1 =$   $q_1 \cdot k_2 =$

⊛ How much Focus to place on other part of the i/p Sentence as we encode a word at certain position

Divide $\sqrt{dk}$   —   —

⇓

{ More Stable gradients }

Dimension=3   Dimension ↑↑ =12

↓   O/p ↑↑↑

$\dfrac{Variance}{\sqrt{dk}}$ ≈ $\dfrac{Variance ↑↑↑ =}{\sqrt{dk}}$

SoftMax   [0-1]   [0-1] .

0.88   +   0.12   = 1

Softmax ⇒ The Softmax score determines how much each word will be expressed at this position

✗   ✗   ✗

Value Vectors   $V_1$ [ 1 1 ]   $V_2$ . [ 1 0 1 ]

0.88 [ 1 1 1 ]

[ 0.88 | 0.88 | 0.88 ]   +   [ 0.12 | 0.6 | 0.12 ]

↓

Contextual Vector   $Z_1$ → [  |  |  ]   $Z_2$ → [    ]

# Self Attention At Higher And Detailed Level

**Self-attention, also known as scaled dot-product attention, is a crucial mechanism in the transformer architecture that allows the model to weigh the importance of different tokens in the input sequence relative to each other**

**Scaled Dot-Product Attention**



Idea :

→ | 08 | 0.3 | 0.4 |   | 0.4 | 0.4 | 0.2 |   | 0.3 | 0.4 | 0.5 |   { Contextual Embedding }

Language Translation

Text Summarization

[ Self    Attention ]

| 1 | 0 | 0 |   | 0 | 1 | 0 |   | 0 | 0 | 1 |

The       CAT       SAT

{ Sentence, Dataset }.

Embedding
Layer → Fixed Vector

## 1) Inputs: Queries, Keys, And Values

Model ——→ Queries, Keys And Values

## 1. Query Vectors (Q):

Role: Query vectors represent the token for which we are calculating the attention. They help determine the importance of other tokens in the context of the current token.

**Importance**:

**Focus Determination**: **Queries help the model decide which parts of the sequence to focus on for each specific token. By calculating the dot product between a query vector and all key vectors, the model assesses how much attention to give to each token relative to the current token.**

**Contextual Understanding**: **Queries contribute to understanding the relationship between the current token and the rest of the sequence, which is essential for capturing dependencies and context.**

## 2. Key Vectors (K):

Role: Key vectors represent all the tokens in the sequence and are used to compare with the query vectors to calculate attention scores.

Importance:

**Relevance Measurement**: **Keys are compared with queries to measure the relevance or compatibility of each token with the current token. This comparison helps in determining how much attention each token should receive.**
**Information Retrieval**: **Keys play a critical role in retrieving the most relevant information from the sequence by providing a basis for the attention mechanism to compute similarity scores.**

## 3. Value Vectors (V):

**Role: Value vectors hold the actual information that will be aggregated to form the output of the attention mechanism.**

Importance:

**Information Aggregation**: Values contain the data that will be weighted by the attention scores. The weighted sum of values forms the output of the self-attention mechanism, which is then passed on to the next layers in the network.
**Context Preservation**: By weighting the values according to the attention scores, the model preserves and aggregates relevant context from the entire sequence, which is crucial for tasks like translation, summarization, and more.

$$\text{Input Sequence} = \left[\text{"The"}, \text{"CAT"}, \text{"SAT"}\right]$$

Embedding Size $= 4$

$Q, K, V \Rightarrow 4$

$\hookrightarrow \boxed{\text{Self Att}} \rightarrow \boxed{0.4 \mid 0.2 \mid 0.3} \leftarrow$

① Token Embedding

$$E_{The} = \left[1\ 0\ 1\ 0\right]$$

$$E_{SAT} = \left[1, 1, 1, 1\right]$$

$\boxed{1 \mid 0 \mid 0}$ ☐ ☐

$\Downarrow$

Sentence, Dataset

$$E_{CAT} = \left[0\ 1\ 0\ 1\right]$$

## ② Linear Transformation

We create $Q, K, V$ by multiplying the embeddings by <u>learned</u>
weights matrices $W_Q$, $W_K$ and $W_V$.

CAT

$$\boxed{1 \quad 0 \quad 0}$$

$Q$ ☐☐☐

$K$ ☐☐☐☐

$V:$ ☐☐☐

$W_Q$

$W_K$

$W_V$

TRAINED

{ learned
  Weights. }

Lets consider

$$W_Q = W_K = W_V = I \qquad \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \leftarrow$$

$$Q_{The} = \begin{bmatrix} 1 & 0 & 10 \end{bmatrix} \cdot \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 10 \end{bmatrix}$$

$$K_{The} = \begin{bmatrix} 1 & 0 & 10 \end{bmatrix}$$

$$V_{The} = \begin{bmatrix} 1 & 0 & 10 \end{bmatrix}$$

① $Q_{The} = K_{The} = V_{The} = \begin{bmatrix} 1 & 0 & 10 \end{bmatrix}$

② $Q_{CAT} = K_{CAT} = V_{CAT} = \begin{bmatrix} 0 & 1 & 01 \end{bmatrix}$

③ $Q_{SAT} = K_{SAT} = V_{SAT} = \begin{bmatrix} 1,1,1,1 \end{bmatrix}$

③ Compute Attention Scores                    $\begin{bmatrix} 1 & 0 & 10 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$

The

$Score\ (Q_{The}, K_{The}) = \begin{bmatrix} 1 & 0 & 1 & 0 \end{bmatrix} \cdot \begin{bmatrix} 1 & 0 & 1 & 0 \end{bmatrix}^T = 2$

$Score\ (Q_{The}, K_{CAT}) = \begin{bmatrix} 1 & 0 & 1 & 0 \end{bmatrix} \cdot \begin{bmatrix} 0 & 1 & 0 & 1 \end{bmatrix}^T = 0$

$Score\ (Q_{The}, K_{SAT}) = \begin{bmatrix} 1 & 0 & 1 & 0 \end{bmatrix} \cdot \begin{bmatrix} 1 & 1 & 1 & 1 \end{bmatrix}^T = 2$

For the token CAT

$Score\ (Q_{CAT}, K_{The}) = \begin{bmatrix} 0 & 1 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1 & 0 & 1 & 0 \end{bmatrix}^T = 0$

$Score\ (Q_{CAT}, K_{CAT}) = \begin{bmatrix} 0 & 1 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 0 & 1 & 0 & 1 \end{bmatrix}^T = 2$

$Score\ (Q_{CAT}, K_{SAT}) = \begin{bmatrix} 0 & 1 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1 & 1 & 1 & 1 \end{bmatrix}^T = 2$

For the Token SAT

$\Bigl\{ \begin{array}{l} \end{array}$  $Score\ (Q_{SAT}, K_{The}) = \begin{bmatrix} 1 & 1 & 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1 & 0 & 10 \end{bmatrix}^T = 2$

$Score\ (Q_{SAT}, K_{CAT}) = 2$

$Score\ (Q_{SAT}, K_{SAT}) = 4$

④ Scaling $\div$

   We take up the scores and scale down by dividing the scores by the

   $\sqrt{d_k} \Rightarrow d_k = 4 \qquad \sqrt{d_k} = 2.$

Scaling in the attention mechanism is crucial to prevent the dot
product from growing too large. $\Rightarrow$ Ensure stable gradients during Training

$d_k$ is large $\longrightarrow$
① Gradient Exploding
② Softmax Saturation $\{ \int \int \} \rightarrow$ Vanishing Gradient Problem.

$Q = \begin{bmatrix} 2 & 3 & 4 & 1 \end{bmatrix}$ $\quad K_1 = \begin{bmatrix} 1 & 0 & 1 & 0 \end{bmatrix}$ $\quad K_2 = \begin{bmatrix} 0 & 1 & 0 & 1 \end{bmatrix}$

## Without Scaling

$Q \cdot K_1^T = 2 \times 1 + 3 \times 0 + 4 \times 1 + 1 \times 0 = 2 + 4 = 6.$

$Q \cdot K_2^T = 2 \times 0 + 3 \times 1 + 4 \times 0 + 1 \times 1 = 0 + 3 + 0 + 1 = 4$

*) Score $[6, 4]$ $\Rightarrow$ Scaling Not Applied

$\text{Softmax}([6,4]) = \begin{bmatrix} \dfrac{e^6}{e^6 + e^4} & , & \dfrac{e^4}{e^6 + e^4} \end{bmatrix} = \begin{bmatrix} \dfrac{e^6}{e^6(1 + e^{-2})} & , & \dfrac{e^4}{e^4(e^2+1)} \end{bmatrix}$

① Property of Softmax $\qquad W_Q, W_K, W_V$

$= \begin{bmatrix} \dfrac{1}{(1 + e^{-2})} & , & \dfrac{1}{(e^2 + 1)} \end{bmatrix}$

$([10, 1]) = [0.99, \boxed{0.01}]$

Dot product = Large values

$\approx [0.88, 0.12].$

Most of the attention weight is assigned to the first key vector, very little to the second vector.

## With Scaling

① Compute Scaled Dot Product

$\sqrt{d_k} = \sqrt{4} = \sqrt{2}$ $\Rightarrow$ Variance $_{2,3}$

dimension ↑ Variance↑

$[6, 4] \Rightarrow \text{Scale} \Rightarrow \begin{bmatrix} \dfrac{6}{2} & , & \dfrac{4}{2} \end{bmatrix} = \begin{bmatrix} 3, 2 \end{bmatrix}.$

$\sqrt{d_k}$↑ same

$\text{Softmax}([3,2]) = \begin{bmatrix} \dfrac{e^3}{e^3 + e^2} & , & \dfrac{e^2}{e^3 + e^2} \end{bmatrix} = \begin{bmatrix} \dfrac{e^{3}}{e^{3}(1 + e^{-1})} & , & \dfrac{e^{2}}{e^{2}(e^{1} + 1)} \end{bmatrix}$

$= [0.73, 0.27] \Rightarrow$ Attention Weights

(A) Here, the attention weights are more balanced compared to the unscaled case

Summary of Importance
**Stabilizing Training**: Scaling prevents extremely large dot products, which helps in stabilizing the gradients during backpropagation, making the training process more stable and efficient.

**Preventing Saturation**: By scaling the dot products, the softmax function produces more balanced attention weights, preventing the model from focusing too heavily on a single token and ignoring others.
Improved Learning: Balanced attention weights enable the model to learn better representations by considering multiple relevant tokens in the sequence, leading to better performance on tasks that require context understanding.

**Scaling ensures that the dot products are kept within a range that allows the softmax function to operate effectively, providing a more balanced distribution of attention weights and improving the overall learning process of the model.**

(4) $\underline{\text{Scaling}}$ $= \sqrt{d_k} = \sqrt{4} \Rightarrow 2$

Similarly Scaling will be done for all other Tokens.

$$\text{Scaled-Score}\left(Q_{The}, K_{The}\right) = \frac{2}{2} = \underline{1}$$

$$\text{Scaled-Score}\left(Q_{The}, K_{CAT}\right) = 0/2 = 0$$

$$\text{Scaled-Score}\left(Q_{The}, K_{SAT}\right) = \frac{2}{2} = \underline{1}$$

(5) $\underline{\text{Apply Softmax}}$

$$\text{ATTENTION WEIGHTS}_{\text{"The"}} = \text{softmax}\left([1,0,1]\right) = [0.4223, 0.1554, 0.4223]$$

$$\text{ATTENTION WEIGHTS}_{\text{"CAT"}} = \text{softmax}\left([0,2,2]\right) = [0.1554, 0.4223, 0.4223]$$

$$\text{ATTENTION WEIGHTS}_{\text{"SAT"}} = \text{softmax}\left([2,2,4]\right) = [0.2119, 0.2119, 0.5762]$$

(6) $\underline{\text{Weight Sum of Values}}$

We multiply the attention weights by corresponding value vectors

For the token The =

$$Output_{(The)} = 0.4223 * V_{The} + 0.1554 * V_{CAT} + 0.4223 \cdot V_{sat.}$$

$$= 0.4223 \begin{bmatrix} 1 & 0 & 1 & 0 \end{bmatrix} + 0.1554 * \begin{bmatrix} 0 & 1 & 0 & 1 \end{bmatrix} + 0.4223 * \begin{bmatrix} 1 & 1 & 1 & 1 \end{bmatrix}$$

$$= \begin{bmatrix} 0.4223, & 0, & 0.4223, & 0 \end{bmatrix} + \begin{bmatrix} 0, & 0.1554, & 0, & 0.1554 \end{bmatrix} + \begin{bmatrix} 0.4223, & 0.4223, \\ 0.4223, & 0.4223 \end{bmatrix}$$

$$= \begin{bmatrix} 1.2669, & 0.9999, & 1.2669, & 0.9999 \end{bmatrix}.$$

Contextual
↓ Vector

The $\boxed{1 \mid 0 \mid 1 \mid 0}$ ⟹ Self Attention ⟹ $\begin{bmatrix} 1.2669, 0.9999, 1.2669, \\ 0.9999 \end{bmatrix}$.

① ↳ $Q, K, V \; [W^Q, W^K, W^V]$
② ↳ Attention Score
③ ↳ Scaled
④ ↳ Softmax
⑤ ↳ Weighted Sum of Value $(Softmax \times V)$

④ Multi Head Attention

$$\text{softmax}\left(\frac{Q \times K^T}{\sqrt{d_k}}\right) \cdot V$$

↳ scaling

$= \boxed{Z}$ =) Vectors.

↓

Attention Head.

→ Self Attention with Multi Heads



ATTENTION HEAD #0 — Car SAT MAT — Thinking Machines X
ATTENTION HEAD #1

$Q_0$ ✓   $W_0^Q$    $Q_1$    $W_1^Q$
$K_0$ ✓.   $W_0^K$    $K_1$ ✓   $W_1^K$ ✓
$V_0$ ✓   $W_0^V$    $V_1$ ✓.   $W_1^V$ ✓

It expands model's ability to focus on different position of terms.

Score, Softmax X V

$[ \; \_ \; \_ \; \_ \; ]$    ↓    $[ \; \_ \; \_ \; \_ \; ]$
$Z_0$    Vectors    $Z_1$

Multi Head Attention



Thinking Machines X

Calculating attention separately in eight different attention heads

$Q, K, V$    $Q, K, V$ (WK Wq Wv)

ATTENTION HEAD #0    ATTENTION HEAD #1    ...    ATTENTION HEAD #7

$Z_0$    $Z_1$    $Z_7$

$\boxed{Z}$    $\boxed{Z_k}$    $\boxed{Z}$

(f) Feed Forward Neural N/W

FEED FORWARD

$Z_1$  $Z_2$  $Z_3$  $Z_4$ $Z = Z_1$      $Z_1$ □|□|□      $Z_2$ □|□|□      $Z_3$ □|□|□   } ⟹ Coni
□ □ □ □□ □

SELF  ATTENTION

→ □|□|□    → □|□|□    → □|□|□    { All
→ How ⟲      ARE ⟲       You ⟲       Will
             —            —

X

Thinking
Machines  ▦

$Q, K, V$        $Q, K, V$          Calculating attention separately in   $Q, K, V$.
□                □                  eight different attention heads        □

ATTENTION        ATTENTION          ...                ATTENTION
HEAD #0          HEAD #1                               HEAD #7

$Z_0$            $Z_1$                                 $Z_7$
▦                ▦                                     ▦

1) Concatenate all the attention heads    7 Head Attention    2) Multiply with a weight
                                                               matrix $W^0$ that was trained
$Z_0$  $Z_1$  $Z_2$  $Z_3$  $Z_4$  $Z_5$  $Z_6$  $Z_7$        jointly with the model
CAT  →  ▦▦▦▦▦▦▦▦▦▦▦▦▦▦▦▦                        X

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

                                                          $W^0$
Z
= ▦                                         Feed
                                            Forward
                                            NN.

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 → X

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

R

$W_0^Q$ $W_0^K$ $W_0^V$ → $Q_0$ $K_0$ $V_0$ → $Z_0$

$W_1^Q$ $W_1^K$ $W_1^V$ → $Q_1$ $K_1$ $V_1$ → $Z_1$

...

$W_7^Q$ $W_7^K$ $W_7^V$ → $Q_7$ $K_7$ $V_7$ → $Z_7$

$W^O$

Z





① Positional Encoding — Representing Order of Sequence

① Lion Kills Tiger ⎫
② Tiger Kills Lion ⎭

⟱

Missing

Self Attention

$x_1$ ↑    ↑ $x_2$

Advantage

① Word Tokens it can process parallely

⟱

DRAW BACK

Lack the sequential Structure of the words

{order}

Attention Is All
YOU NEED

Positional
Encoded
Vector

BOOK, Journal
Novel
↳ 1 lakh ⎫
    words ⎬ .
         ⎭

⇓

Back propogaten

[ - | - | - | 1 ]    [ - | - | - | 2 ]    [ - | - | - | 3 ] ←

Types of Position Encoding

✓1) Sinsusoidal Position Encoding ⟶

2) learned positional Encoding ⇒) Positional Encoding Arc learned during
                                      Traning .⇐

① Sinsusoidal Positional Encoding ÷ It uses Sine and Cosine functions

Of different frequencies to create positional encodings

Formula ÷                                [ - | - | - | - ]

$$P.E_{(pos, 2i)} = Sin\left(\frac{pos}{10000^{2i/dmodel}}\right)$$      Where pos is the position
                                                              i is the dimension

$$P.E_{(pos, 2i+1)} = Cos\left(\frac{pos}{10000^{2i/dmodel}}\right).$$   dmodel is the dimensionality
                                                                Of the embeddings.

Eg: The cat Sat

The ⟶ [ 0.1  0.2  0.3  0.4 ]              [ | | | | ]              { Miss the order of
                                                                    the elements}
CAT ⟶ [ 0.5  0.6   0.7  0.8 ]         [ - | - | - | - ]    +1
                                        0   1  2  3
SAT ⟶ [ 0.9  1.0   1.1  1.2 ].

                                                                1 —

[ | | | | ]

$$P.E_{(pos,2i)} = Sin\left(\frac{pos}{10000^{2i/dmodel}}\right) \qquad P.E_{(pos,2i+1)} = Cos\left(\frac{pos}{10000^{2i/dmodel}}\right).$$
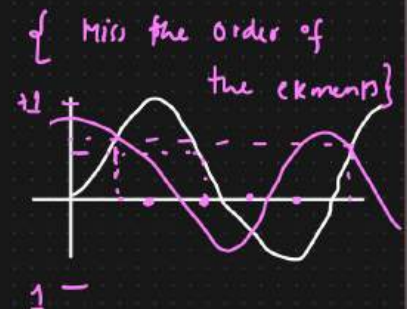
For our example $\quad dmodel = 4$

For position $\quad pos = 0$

$$P.E(0,0) = Sin\left(\frac{0}{10000^{0/4}}\right) = Sin(0) = 0$$

$$PE(0,1) = Cos\left(\frac{0}{-}\right) = Cos(0) = \underline{1}$$

$$P.E(0,2) = Sin\left(\frac{0}{10000^{4/4}}\right) = Sin(0) = 0$$

$$P.E(0,3) = Cos\left(0\right) = \underline{1}$$

$$P.E = [0,1,0,1] \qquad P.E_{(pos,2i)} = Sin\left(\frac{pos}{10000^{2i/dmodel}}\right)$$
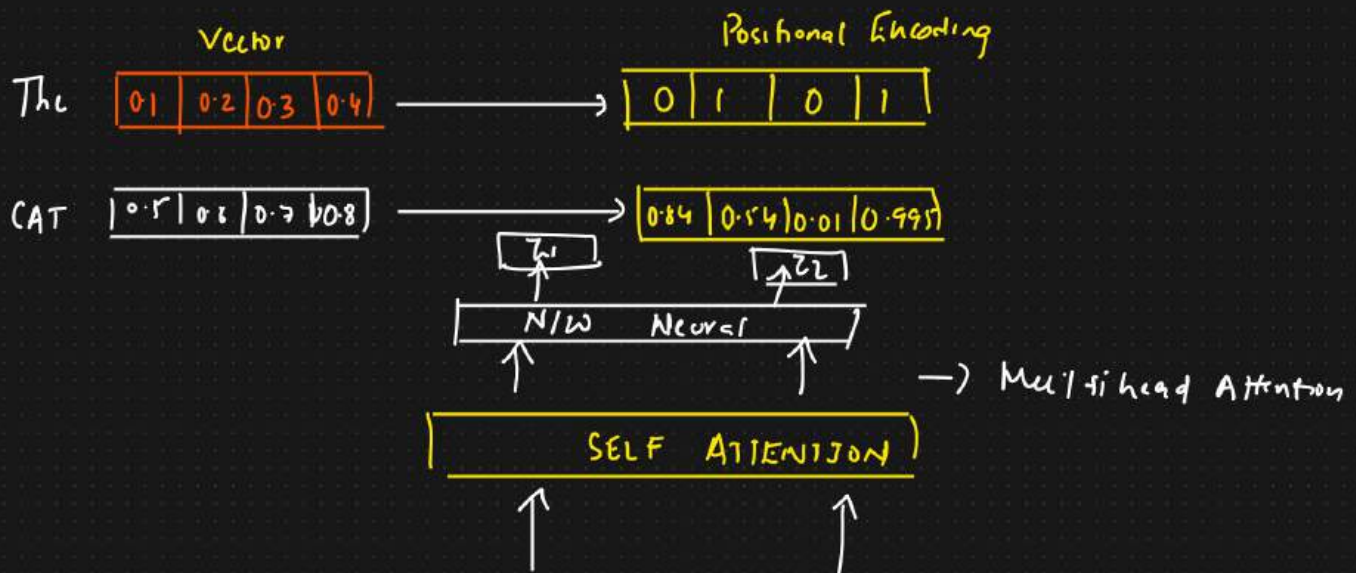
For $pos = 1$

$$P.E(1,0) = Sin\left(\frac{1}{10000^{0/4}}\right) = Sin(1) = 0.8415 \qquad P.E_{(pos,2i+1)} = Cos\left(\frac{pos}{10000^{2i/dmodel}}\right)$$

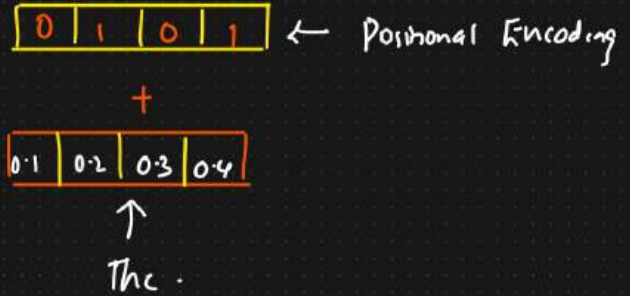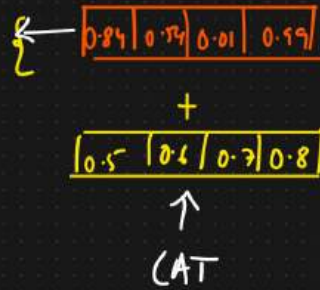$$PE(1,1) = Cos\left(\frac{1}{10000^{2/4}}\right) \approx 0.5403$$

$$PE(1,2) = Sin\left(\frac{1}{10000^{2/4}}\right) \approx 0.01$$

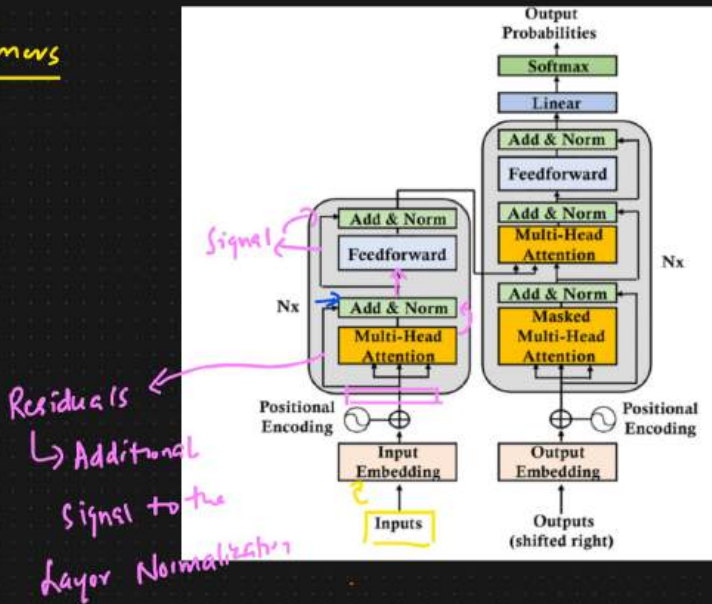$$PE(1,3) = Cos \qquad\qquad \approx 0.99995$$

Vector           Positional Encoding

The | 0.1 | 0.2 | 0.3 | 0.4 |  $\longrightarrow$ | 0 | 1 | 0 | 1 |

CAT | 0.5 | 0.6 | 0.7 | 0.8 |  $\longrightarrow$ | 0.84 | 0.54 | 0.01 | 0.995 |

$Z_1$       $Z_2$

N/W Neural

$\longrightarrow$ Multihead Attention

SELF ATTENTION

**Positional Encoding** { ← | 0.84 | 0.74 | 0.01 | 0.99 |

| 0 | 1 | 0 | 1 | ← Positional Encoding

+

| 0.5 | 0.6 | 0.7 | 0.8 |

| 0.1 | 0.2 | 0.3 | 0.4 |

↑
CAT

↑
The.

---

**(*) Layer Normalization In Transformers**

Transformers

Residuals

Signal

Residuals
↳ Additional
   Signal to the
   Layer Normalization



① Self Attention layer
② Multi head Attention
③ Positional Encoding
④ Layer Normalization
   ↰
ADD AND Normalize



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 $W^0$ 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

Layer Normalization.
Add And Normalize
Information

---

Normalization → Batch Normalization
             → Layer Normalization

$f_1$    $f_2$

| House Size | No.of Rooms | Price |
|------------|-------------|-------|
| 1200 | 2 | 45 |
| 1500 | 3 | 70 |

$f_1$
$f_2$ → O/p.

ANN

**Normalization : Standard Scaling**

$\mu, \sigma \leftarrow | \underline{2000} | \qquad 3.5 \qquad 80$

$$Z_{score} = \frac{x_i - \mu}{\sigma} \implies \boxed{\mu=0, \sigma=1} \qquad f_1 \longrightarrow f_1' \Leftarrow \mu=0, \sigma=1$$
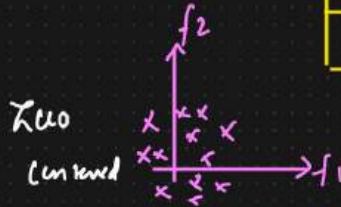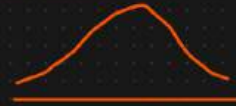
**Deep Learning :** i/p Image $\implies$ Min Max Scaler

$0-255 \leftarrow$ 

$\implies$ Min Max Scaler $\implies$



**Advantages**

① Improved Training Stability
$\downarrow$
Vanishing and exploding Gradient problem

Zero centered

$\{\mu=0\}$  $\mu=0, \sigma=1$

$\{\sigma=1\}$

② Faster Convergence



Loss / W

$\implies$ Back propogation $\implies$ Stable update.

{Batch Normalization}

| | f1 House size | f2 Rooms | Price | Normalisation $Z_1$ | $Z_2$ |
|---|---|---|---|---|---|
| | 0.45 | 0.55 | 45 | = | = |
| | 0.60 | 0.20 | — | — | — |
| | — | — | — | — | — |



$0.45 \quad f_1$
$0.55 \quad f_2$

$Z_1 \to $ value 1
$b_1$
$Z_2 \Rightarrow$ Values.
$b_2$

$$Z_1 = \sigma\left[(0.45 * W_1 + 0.55 * W_3) + b_1\right] = \text{Value 1} \}$$

Does this distribution

$$Z_2 = \sigma( \qquad\qquad + b_2) = \text{Value 2} \}. \qquad \boxed{\mu=0, \sigma=1}$$

$\mu_1, \sigma_1 \qquad \mu_2, \sigma_2$

**Batch Normalization Vs Layer Normalization**

f1   f2   $Z_1$   $Z_2$



$\to \mu_1, \sigma_1$

$\to \mu_2, \sigma_2$

$\to \mu_3, \sigma_3$

$Z_{score} = \dfrac{x_i - \mu_1}{\sigma_1}$

$Z_{score} = \dfrac{x_i - \mu_2}{\sigma_2}$

$Z_{score} = \dfrac{x_i - \mu_3/\sigma_3}{}$ ..

Layer Normalization

$\gamma, \beta \rightarrow$ learnable parameters

$\boxed{z_1}$ $\boxed{z_2}$

$\rightarrow$ layers

$\rightarrow$ layer

$\rightarrow 0$

$\rightarrow 0 \quad 0 \rightarrow \boxed{\mu=0, \sigma=0}$

$\rightarrow 0 \quad 0$

$0$

Normalizing

$\boxed{\mu, \sigma}$

Batch Normalization

**Normalization**

$\Downarrow$

$\gamma, \beta$

$z_1 = \sigma \left[ w_1^T x + b_1 \right]$

$\hookrightarrow$ learnable parameters $\searrow$

$y = \boxed{\gamma} \left[ \dfrac{z_1 - \mu_1}{\sigma_1} \right] + \boxed{\beta}$

$\Downarrow$

Scale And Shift parameters

$\nearrow$ Normalized

1) "CAT" $= [2\cdot0, 4\cdot0, 60, 8\cdot0] \leftarrow \{$ Vectors $\}$

2) Parameters $= \gamma = [1\cdot0, 1\cdot0, 1\cdot0, 1\cdot0] \rightarrow$ Learned Scale

$\beta = [0\cdot0, 0\cdot0, 0\cdot0, 0\cdot0] \rightarrow$ Shift

$\Big\} \Rightarrow$ Scale And Shift param

$z_{score} = \dfrac{x_i - \mu}{\sigma}$

i) Compute the mean

$\mu = \dfrac{1}{4} (2\cdot0 + 4\cdot0 + 6\cdot0 + 8\cdot0)$

$= \dfrac{20\cdot0}{4} = 5\cdot0$

ii) Compute the variance $(\sigma^2)$

$\sigma^2 = \dfrac{1}{4} \left[ (2\cdot0 - 5\cdot0)^2 + (4\cdot0 - 5\cdot0)^2 + (6\cdot0 - 5\cdot0)^2 + (8\cdot0 - 5\cdot0)^2 \right] = 5\cdot0$

iii) Normalize the i/p

$\hat{x}_1 = \dfrac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}} \qquad \epsilon = 1e^{-5} \Rightarrow$ Avoid division by 0

$\sqrt{\sigma^2 + \epsilon} = \sqrt{5\cdot0 + 1e-5} \approx \sqrt{5\cdot00001} = 2\cdot236$

$$\hat{x}_1 = \frac{2 \cdot 0 - 5 \cdot 0}{2 \cdot 236} \approx -1.34$$

Normalized vector

$$\hat{x}_2 = \frac{4 \cdot 0 - 5 \cdot 0}{2 \cdot 236} \approx -0.45$$

$$\hat{x} = \left[ -1.34, -0.45, 0.45, 1.34 \right]$$

$$\hat{x}_3 = \frac{6 \cdot 0 - 5 \cdot 0}{2 \cdot 236} \approx 0.45$$

$$\hat{x}_4 = \frac{8 \cdot 0 - 5 \cdot 0}{2 \cdot 236} \approx 1.34$$

4) <u>Scale And Shift</u>

$$y_i = \gamma_i \hat{x}_i + \beta_i$$

$$\gamma = \left[ 1.0, 1.0, 1.0, 1.0 \right] \quad \beta = \left[ 0.0, 0.0, 0.0, 0.0 \right]$$

$$y = \left[ -1.34, -0.45, 0.45, 1.34 \right]$$

⑤ <u>Encoder Architecture</u> [Research Paper]



Figure 1: The Transformer - model architecture.

Sequence to
Sequence
⇓
Complex

Feed Forward NN. → ooooooooooooo

(Add And Norm) ────→ • • • • • •

Multi Head Attention → $z_1, z_2, z_3, z_4, z_5, z_6, z_7, z_8$ {Research paper}.

Residuals

Text Embeddings + Positional Encoding

J/p Sequence

⇒ 512. → {Research paper}.

every word = 512.

$Q_1 = 64$

$K = 64$

$V = 64$

$\sqrt{64} = 8$ //.

# Self Attention to Feed Forward Neural N/W



1) This is our input sentence*

2) We embed each word*

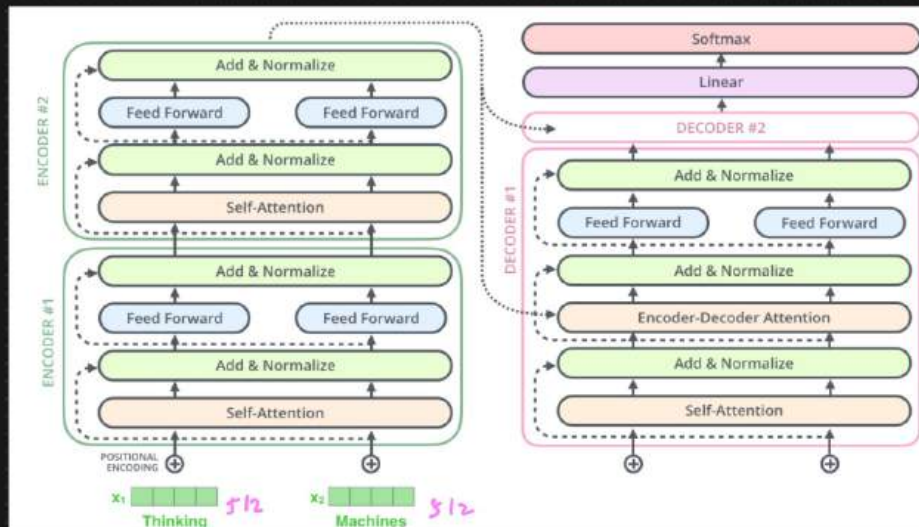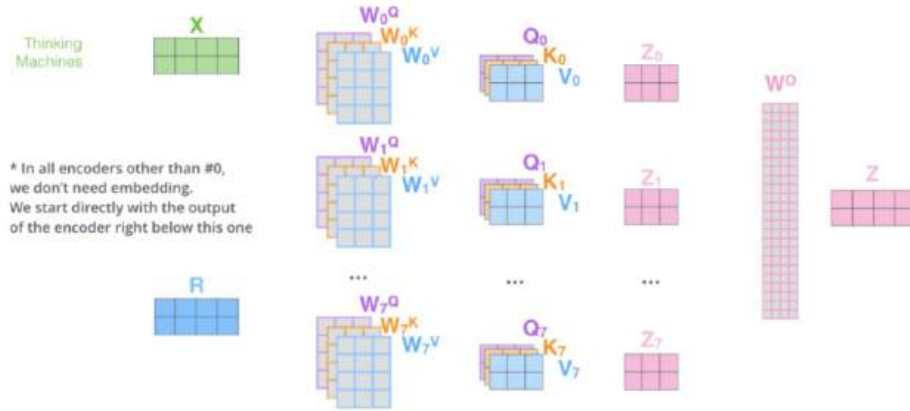3) Split into 8 heads. We multiply X or R with weight matrices

4) Calculate attention using the resulting Q/K/V matrices

5) Concatenate the resulting Z matrices, then multiply with weight matrix W⁰ to produce the output of the layer

Thinking Machines

X

$W_0^Q$
$W_0^K$
$W_0^V$

$Q_0$
$K_0$
$V_0$

$Z_0$

W⁰

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

$W_1^Q$
$W_1^K$
$W_1^V$

$Q_1$
$K_1$
$V_1$

$Z_1$

Z

R

$W_7^Q$
$W_7^K$
$W_7^V$

$Q_7$
$K_7$
$V_7$

$Z_7$



Embedding Vector: 512

$Q, K, V = 64$

Head Attention: 8

⊛ **Residual** connection : Skip connection NN.

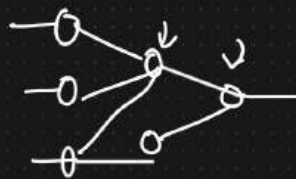1) Addressing the Vanishing Gradient Problem

**Residuals** : Residual connection create a short paths for gradients to flow directly through the n/w. Gradient remains sufficiently large.

2) Improve Gradient flow

       Convergence will be faster.

3) Enables Training of Deeper Networks.

**① Feed Forward NN**

**① Adding Non Linearity**

**② Processing Each position Independently.**

Self Attention ⟶ (Capture relationships)

FFN ⟶ Each token representation Independently :
⇓
Transforming these representation furthers

and allows the model to learn

Richer Representation.

Linear function problem

{Non Linear function}.

ANN ⟹

**③ FFN ⟶ ,Deeper ⟹ Adds Depth to the Model.**

Depth ↑↑ ⟹ More Learnings ⟶ DATA

**✓ Decoders In Transformers**

**3 main Components**

**The transformer decoder is responsible for generating the output sequence one token at a time, using the encoder's output and the previously generated tokens.**

**① Masked Multi Head Self Attention ✓**

**② Multi Head Attention (Encode Decoder Attention)**
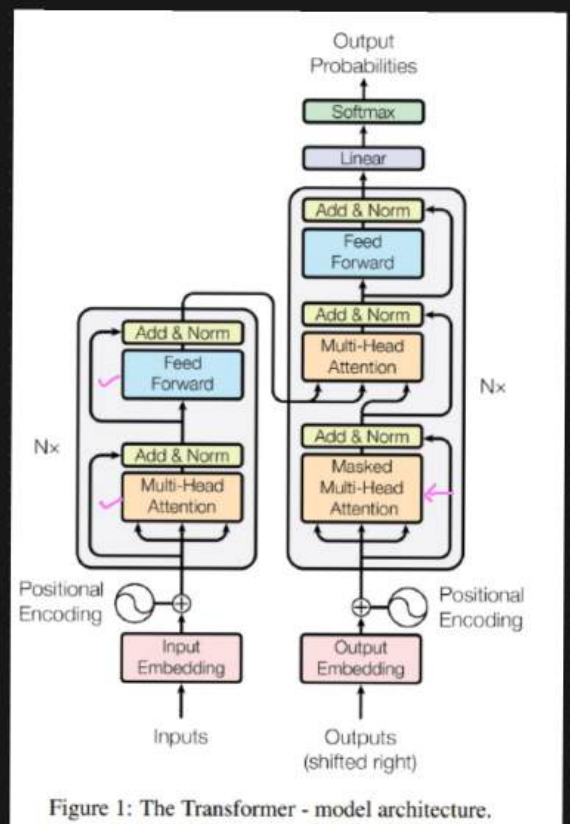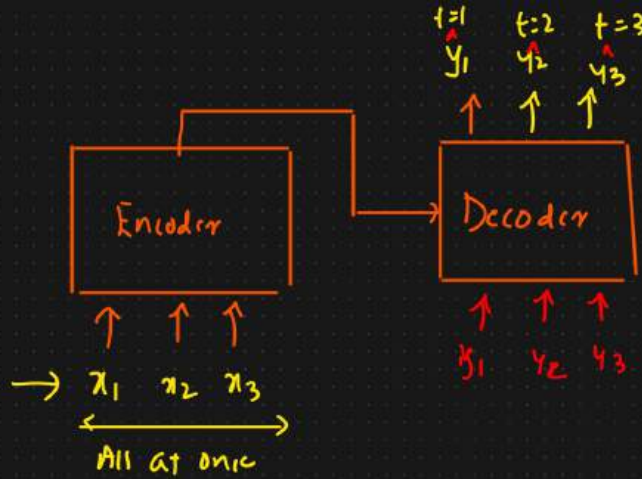
**③ Feed Forward Neural Network.**



Figure 1: The Transformer - model architecture.

t=1   t=2   t=3
$\hat{y}_1$   $\hat{y}_2$   $\hat{y}_3$

$<y_1, y_2, y_3>$

```
┌──────────┐        ┌──────────┐
│          │        │          │
│ Encoder  │───────▶│ Decoder  │
│          │        │          │
└──────────┘        └──────────┘
```

① Training Mechanism
② Inference Mechanism

$\uparrow$  $\uparrow$  $\uparrow$
$y_1$   $y_2$   $y_3$

$\rightarrow$ $x_1$  $x_2$  $x_3$

$\longleftrightarrow$
All at once

---

⊛ Masked Multi Head Self Attention

① I/P Embedding And Positional Embedding ✓ → Zero padding → Sequence Length Equal

② Linear Projection for Q,K,V

③ Scaled Dot Product Attention

✓④ Mask Application } ⟹ Try to understand the imp.

⑤ Multi Head Attention

⑥ Concatenation And Final Linear Projection

⑦ Residual Connection And Layer Normalization



Figure 1: The Transformer - model architecture.

Dataset

Eng                    Hindi

$<x_1, x_2, x_3>$      $<y_1, y_2>$  ←
                              ⟱
                       $<y_1, y_2 \ 0>$
                              ↓
                            Zero
                           padding

① Linear Projections  Q,K,V
② Scaled Dot Product Attention
③ Mask Application → Look Ahead Mask
                   → Padding Mask

I/P
$[4 \quad 5 \quad \overline{6 \quad 7}]$

O/P
$[1 \quad 2 \quad 3]$

$\begin{bmatrix} 1 & 2 & 3 & 0 \end{bmatrix}$
↓
4 dimension vector

Masked
Multi Head
Attention.

i) <u>Input Embedding and Positional Encoding</u>

<u>Output Embedding</u>  [STEP 1]

$$\begin{bmatrix} [0.1, 0.2, 0.3, 0.4], \\ [0.5, 0.6, 0.7, 0.8], \quad + \quad P.E \Rightarrow 0 \\ [0.9, 1.0, 1.1, 1.2], \\ [0.0 \quad 0.0 \quad 0.0 \quad 0.0] \end{bmatrix}$$

Step 2 : Linear Projection for Q, K, and V.

$W_Q = W_K = W_V = I$

Create query (Q), Key(K) and value(V) vectors

Q = Output Embedding $*$ W_Q = Output Embedding.

K =            "            $*$ W_K =        "            "

V =            "            $*$ W_V =        "            "

$Q = K = V = \begin{bmatrix} [0.1, 0.2, 0.3, 0.4], \\ \rightarrow [0.5, 0.6, 0.7, 0.8], \\ [0.9, 1.0, 1.1, 1.2], \\ [0.0 \quad 0.0 \quad 0.0 \quad 0.0] \end{bmatrix}$ $\begin{bmatrix} 0.1 \\ 0.2 \\ 0.3 \\ 0.4 \end{bmatrix}^{K^T} \begin{bmatrix} 0.5 \\ 0.6 \\ 0.7 \\ 0.8 \end{bmatrix} \cdot \begin{bmatrix} 0.9 \\ 1.0 \\ 1.1 \\ 1.2 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$

③ Scaled Dot Product Attention Calculation

Scores $= Q * K^T / \sqrt{d_k}$

$= Q * K^T / 2$.

$\begin{bmatrix} [0.1 * 0.1 + 0.2 * 0.2 + 0.3 * 0.3 + 0.4 * 0.4, \\ 0.1 * 0.5 + 0.2 * 0.6 + 0.3 * 0.7 + 0.4 * 0.8, \\ 0.1 * 0.9 + 0.2 * 1.0 + 0.3 * 1.1 + 0.4 * 1.2 \\ 0.1 * 0 + 0.2 * 0 + 0.3 * 0 + 0.4 * 0 \end{bmatrix}$

Scores :
$$\begin{bmatrix} [0.3, & 0.7, & 1.1, 0.0] \\ [0.7, & 1.9, & 3.1, 0.0] \\ [1.1 & 3.1 & 5.1 & 0.0] \\ [0.0 & 0.0 & 0.0 & 0.0] \end{bmatrix}$$

$$\begin{bmatrix} [ & & & ] \\ [ & & & ] \\ [ & & & ] \end{bmatrix}$$

## ⊕ Masked Application

It helps manage the structure of the sequences being processed and ensures the models behaves correctly during training And Inferencing

### Reasons

① Handling Variable length Sequences with Padding MASK

#### Purpose

① To handle sequences of different length in batch

② To ensure that padding tokens, which are added to make sequences of uniform length, do not affect the model prediction

O/P
$\rightarrow Y_1, Y_2 \ 0 \ 0 \ 0 \ 0 \ 0$

Eg: i/p ← Sequence 1  $[1, 2, 3]$

O/p ← Sequence 2  $[4, 5, \boxed{0}]$   0 is the padding token
$\rightarrow 100$       ↑
              └→ Influence the Attention Mechanism
                         ⇓
              lead to Incorrect or biased predictions.

A padding mask ⇒ The tokens Are ignored.

Masking ⬚ → Padding Mask           Padding Mask  $[1 \quad 1 \quad 1]$
        → Look Ahead Mask .  }
                                              $[1 \quad 1 \quad 0]$

② **Look Ahead Mask** → Maintain Auto Regressive Property.

① To ensure that each position in the decoder
output sequence can only attend to
the previous position, but no future position

Decoder

How Are Yo

② Sequence → Language Modelling, Translation

→ 1 D Mask.

Eg: $[4, 5, 0]$ → $[1, 1, 0]$

Convert 1D to 2D Mask

$$\rightarrow \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

Token 1 attends

Attention ← to token 1,2

Mechanism      1,2

For each token in the sequence,
the mask should indicate
which tokens it can attend
to.

⊛ **Look Ahead Mask** → Decoder Output

$$\begin{bmatrix} [1 & 0 & 0] \\ [1 & 1 & 0] \\ [1 & 1 & 1] \end{bmatrix}$$

④ **Combine Padding And Looking Ahead Mask**

Element Wise multiplication of 2 mask

$$Combine\ Mask = \begin{bmatrix} [1,0,0] \\ [1,1,0] \\ [0,0,0] \end{bmatrix}$$

# ① MASK

Scores : $\begin{bmatrix} [0.3, & 0.7, & 1.1, & 0.0] \\ [0.7, & 1.9, & 3.1, & 0.0] \\ [1.1 & 3.1 & 5.1 & 0.0] \\ [0.0 & 0.0 & 0.0 & 0.0] \end{bmatrix}$

## Look Ahead Mask

$\begin{bmatrix} [1 & 0 & 0 & 0] \\ [1 & 1 & 0 & 0] \\ [1 & 1 & 1 & 0] \\ [1 & 1 & 1 & 1] \end{bmatrix}$

## Padding Masking [extended to 2 D Format]

$\begin{bmatrix} [1,1,1,0] \\ [1,1,1,0] \\ [1,1,1,0] \\ [0,0,0,0] \end{bmatrix}$

## Combined Mask : Look Ahead Mask + Padding Mask

$\begin{bmatrix} [1*1, 0*1, 0*1, 0*0] \\ [1*1, 1*1, 0*1, 0*0] \\ [1*1, 1*1, 1*1, 0*0] \\ [1*1, 1*1, 1*1, 1*0] \end{bmatrix} = \begin{bmatrix} [1,0,0,0], \\ [1,1,0,0], \\ [1,1,1,0], \\ [1,1,1,0] \end{bmatrix} \Rightarrow \begin{bmatrix} [1,-\infty,-\infty,-\infty]. \\ [1,1,-\infty,-\infty], \\ [1,1,1,-\infty] \\ [1,1,1,-\infty] \end{bmatrix}$

## Masked Score

$\pi^0$

$\begin{bmatrix} [0.3, -\infty, -\infty, -\infty] \\ [0.7, 1.9, -\infty, -\infty] \\ [1.1, 3.1, 5.1, -\infty] \\ [0.0, 0.0, 0.0, -\infty] \end{bmatrix}$

Zero out the influence when the Softmax is applied.

Attention weight.

(*) **Softmax**

$$\text{Softmax Score} = \text{softmax}(\text{Masked Scores})$$

$$= \begin{bmatrix} [1.0, 0.0, 0.0, 0.0], \\ [0.3, 0.7, 0.0, 0.0], \\ [0.1, 0.3, 0.6, 0.0], \\ [1.0, 0.0, 0.0, 0.0] \end{bmatrix}$$

(*) **Weigth Sum of Value**

$$\text{Attention O/p} = \text{Softmax Scores} * V.$$

## Masking

Masking in the transformer architecture is essential for several reasons. It helps manage the structure of the sequences being processed and ensures the model behaves correctly during training and inference. Here are the key reasons for using masking:

**1. Handling Variable-Length Sequences with Padding Mask**

Purpose
To handle sequences of different lengths in a batch.
To ensure that padding tokens, which are added to make sequences of uniform length, do not affect the model's predictions.

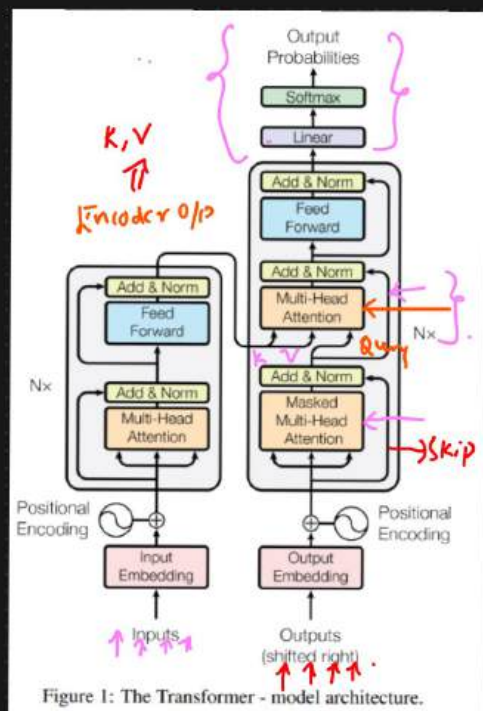**2. Maintaining Autoregressive Property with Look-Ahead Mask**
Purpose
To ensure that each position in the decoder output sequence can only attend to previous positions and itself, but not future positions.
This is crucial for sequence generation tasks like language modeling and translation, where the model should not have access to future tokens when predicting the current token.

# A Encoder Decoder Multi Head Attention



Figure 1: The Transformer - model architecture.

K,V

↑

Encoder O/p

① Encoder O/p → Set of Attention vector K & V

② Masked Multihead → Attention vector Q {Query Vector}

These are to be used by each decoder in its "Encoder–decoder" Attention layer

⇓

Helps the Decoder to focus on appropriate places in the i/p Sequence

→Skip Connection

Dataset

i/P                     O/P

$\langle x_1, x_2, x_3 \rangle$        $\langle y_1, y_2, y_3 \rangle$



Decoding time step: 1 ② 3 4 5 6        OUTPUT |

$K_{encdec}$  $V_{encdec}$    Linear + Softmax

ENCODERS        DECODERS

EMBEDDING WITH TIME SIGNAL

EMBEDDINGS

INPUT   Je   suis   étudiant      PREVIOUS OUTPUTS   I

# ⓧ The Final Linear And Softmax Layer  { Vectors ⟶ o/p word}
{ BLOG } ⟹ Transformers



Figure 1: The Transformer - model architecture.

Vectors ⇓ Word



Which word in our vocabulary is associated with this index?  **am**

Get the index of the cell with the highest value (argmax)  Word  5

log_probs — No-of words vocab

Softmax → Multi class Classification

logits — Score of a unique word.

Linear

Decoder stack output

**linear** ⟹ The linear layer is a simple fully connected neural n/w that projects the vector produced by the stack of Decoder ⟹ logits vector ⟹

**Model** ⟹ 10,000 ⟹ Vocabulary ⟹ logits vector = 10000 cells wide

② Softmax layer turns those scores into probabilities (all add upto 1·0).
The cell with the highest probability is chosen, and the word associated with it is produced as the o/p. ⟹ time stamp.

## Recap of Training

| Output Vocabulary | | | | | | |
|---|---|---|---|---|---|---|
| **WORD** | a | am | I | thanks | student | \<eos\> |
| **INDEX** | 0 | 1 | 2 | 3 | 4 | 5 |

## Output Vocabulary

| WORD | a | am | I | thanks | student | <eos> |
|------|---|----|----|--------|---------|-------|
| INDEX | 0 | 1 | 2 | 3 | 4 | 5 |

OHE ←

One-hot encoding of the word "am"

| 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 |
|-----|-----|-----|-----|-----|-----|

Merci → Thanks

I am a Student <eos>.

---

**Untrained Model Output**

| 0.2 | 0.2 | 0.1 | 0.2 | 0.2 | 0.1 |
|-----|-----|-----|-----|-----|-----|

**Correct and desired output**

| 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 |
|-----|-----|-----|-----|-----|-----|

| a | am | I | thanks | student | <eos> |

=) Loss function ↓↓.

Back Propagation

---

**i/p**   **O/p**

a, am, i, Student    I am a Student
↓   ↓   ↓   ↓
OHE  OHE  OHE  OHE

---

**Target Model Outputs**

Output Vocabulary:

| | a | am | I | thanks | student | <eos> |
|---|---|----|----|--------|---------|-------|
| position #1 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 |
| position #2 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| position #3 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| position #4 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 |
| position #5 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 |
| | a | am | I | thanks | student | <eos> |

Loss

**Trained Model Outputs**

Output Vocabulary:

| | a | am | I | thanks | student | <eos> |
|---|---|----|----|--------|---------|-------|
| position #1 | 0.01 | 0.02 | 0.93 | 0.01 | 0.03 | 0.01 |
| position #2 | 0.01 | 0.8 | 0.1 | 0.05 | 0.01 | 0.03 |
| position #3 | 0.99 | 0.001 | 0.001 | 0.001 | 0.002 | 0.001 |
| position #4 | 0.001 | 0.002 | 0.001 | 0.02 | 0.94 | 0.01 |
| position #5 | 0.01 | 0.01 | 0.001 | 0.001 | 0.001 | 0.98 |
| | a | am | I | thanks | student | <eos> |