

SELF-ATTENTION :

Self-Attention is a Mechanism that lets a model look at other words in the same sentence (or input sequence) to decide which ones are important when processing each word.

Why it Matters ?

Helps capture Context : Eg ; in "The bank near the river", Self-attention helps understand the "bank" refers to a location, not finance.

Real-Time Example:

"The cat sat on the mat because it was tied"

Here the Tricky part:

⇒ The word "it" refers to "cat", not "mat".

Self-attention is exactly designed to figure out these relationships.

Q, K, V in Self-Attention:-

Query (Q):-

⇒ Represents the current word's question: "What am I looking for in other words?"

⇒ Example: "play" might query "Who is playing?" or "What is being played?"

Key (K):-

⇒ Represents each word's identity/features: "What do I have that others may want to know?"

⇒ Example: "I" has the feature "Subject / dear", "Football" has "Object / game".

Value (V):-

⇒ The actual Information content carried by the word.

⇒ Once Attention is calculated (Q-K → importance score), we use it to pick up values (V).

⇒ Example: if "play" attends to "Football", it grabs "Football's Meaning" through V.

Simple Analogy:

Query = What I want (my Question).

Key = What I can offer (my features).

Value = The actual info I give when someone attends to me.

Self-Attention Formula:-

$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Explanation:-

$QK^T \rightarrow$ Similarity Scores

\Rightarrow Multiply each Query (Q) with all keys (K).

\Rightarrow This gives how much one word should attend to another.

\Rightarrow Bigger Score = Stronger relationship.

\Rightarrow Example: "play". "foot ball" \rightarrow high score,

"play". "the" \rightarrow low score.

Softmax - Attention Weights

\Rightarrow Converts similarity scores into prob

\Rightarrow Each word's attention to others
0 and 1.

Why Softmax is chosen?

\Rightarrow Simple

\Rightarrow Smooth & differentiable

⇒ Probabilistic

⇒ Balances focus vs. Distribution

That's why almost Every attention-based model (BERT, GPT, Transformers in general) uses Softmax.

Multiply with V - Weighted Values

⇒ Each weight is applied to the Value (V) of the word.

⇒ So the final Embedding of a word is the weighted Sum of Values it gathered from others.

Why Divide by $\sqrt{d_k}$ in Self-Attention?

Context

In Self-attention, we compute

$$\text{Score} = Q \cdot K^T$$

This dot product can get Very large when the Vector dimension d_k is high.

Problem:

⇒ Large dot products → large Values in Softmax.

⇒ Softmax becomes too sharp → one token dominates.

Ex: $[1.0, 0, 0]$

⇒ Model becomes unstable during training.

Solution:

We scale down the dot product by dividing by

$\sqrt{d_k}$:

$$\text{Scaled Score} = \frac{Q \cdot K^T}{\sqrt{d_K}}$$

This keeps the values in a reasonable range,
so softmax stays smooth and gradients don't explode

Intuition:

Think of $\sqrt{d_K}$ as a normalizer.

prevents attention from becoming overly confident too early.

Example Calculation:

$$I = [0.2, 0.4, 0.6]$$

$$\text{play} = [0.8, 0.3, 0.3]$$

$$\text{football} = [0.1, 0.2, 0.5]$$

Stacking them gives $X \in \mathbb{R}^{3 \times 3}$

$$X = \begin{bmatrix} [0.2, 0.4, 0.6], \\ [0.8, 0.3, 0.3], \\ [0.1, 0.2, 0.5] \end{bmatrix}$$

Weights of Q, K, V :

Each map $3 \rightarrow 2$ so $Q, K, V \in \mathbb{R}^{3 \times 2}$

$$W_Q = \begin{bmatrix} [1.0, 0.0], \\ [0.0, 1.0], \\ [1.0, -1.0] \end{bmatrix}$$

$$W_K = \begin{bmatrix} [0.5, -0.5], \\ [1.0, 0.0], \end{bmatrix}$$

$$[0.0, 1.0]$$

$$WLV = \begin{bmatrix} 1.0, 1.0 \\ 0.5, -0.5 \\ 1.0, 0.0 \end{bmatrix}$$

compute Q, K, V ($Q = X \cdot W_Q = X \cdot W_L \cdot K$, $V = X \cdot W_V$)

$$Q = \begin{bmatrix} 0.8, -0.2 \\ 1.1, 0.80 \\ 0.6, -0.3 \end{bmatrix}$$

$$K = \begin{bmatrix} 0.50, 0.50 \\ 0.70, -0.10 \\ 0.25, 0.45 \end{bmatrix}$$

$$V = \begin{bmatrix} 1.00, 0.00 \\ 1.25, 0.65 \\ 0.70, 0.00 \end{bmatrix}$$

i) Attention Scores = QK^T

$$Q = \begin{bmatrix} 0.8 & -0.2 \\ 1.1 & 0.0 \\ 0.6 & -0.3 \end{bmatrix}, K = \begin{bmatrix} 0.50 & 0.50 \\ 0.70 & -0.10 \\ 0.25 & 0.45 \end{bmatrix}$$

$$\text{Scores} = QK^T = \begin{bmatrix} 0.300 & 0.580 & 0.110 \\ 0.550 & 0.770 & 0.275 \\ 0.150 & 0.450 & 0.015 \end{bmatrix}$$

ii) Scale by $\sqrt{d_k}$

Here $d_k = 2 \Rightarrow \sqrt{2} \approx 1.4142$

$$\frac{\text{Scores}}{\sqrt{2}} = \begin{bmatrix} 0.2121 & 0.4101 & 0.0778 \\ 0.3889 & 0.5445 & 0.1945 \\ 0.1061 & 0.3182 & 0.0106 \end{bmatrix}$$

iii) Softmax (row-wise) \rightarrow attention weights

$$\text{weights} = \begin{bmatrix} 0.3233 & 0.3941 & 0.2826 \\ 0.3343 & 0.3905 & 0.2752 \\ 0.3179 & 0.3931 & 0.2890 \end{bmatrix}$$

Interpretation (each row sums to 1):

For I : attends $\sim 32.3\%$ to I, 39.4% to play, 28.3% to football.

For Play : $\sim 33.4\%$ I, 39.1% play, 28.9% football

For football : $\sim 31.8\%$ I, 39.3% play, 28.9% football

Weighted Sum with V:

$$V = \begin{bmatrix} 1.00 & 0.00 \\ 1.25 & 0.65 \\ 0.70 & 0.00 \end{bmatrix}$$

Final Context Vectors (one per token):

$$\text{Output} = \text{weight} \cdot V = \begin{bmatrix} 1.0137 & 0.2561 \\ 1.0151 & 0.2538 \\ 1.0116 & 0.2555 \end{bmatrix}$$

Each now is the new Embedding for [I, play, football] after this attention head is, a blend of the V Vectors weighted by the attention distributions.

Multi-Head Attention :-

Multi-head attention allows the model to attend to different parts of the input simultaneously, using multiple sets of learned projections. Each head learns a different representation of the input sequence.

Why Multi-Head Attention?

In Normal Self-Attention, each token computes attention over all others \rightarrow gives context.

Problem: a single attention head may only learn one kind of relationship (eg: word order, syntactic role).

Solution: Use multiple attention heads in parallel, each head learns different relationships (eg: grammar, semantics, long-distance dependencies).

Single-Head (vs) Multi-Head :-

Single-Head Attention (Limitation) :-

\Rightarrow In one head, queries (Q) look at keys (K) to decide "where to pay attention."

\Rightarrow This produces one set of weights per word \rightarrow one "view" of dependencies.

\Rightarrow Example :

In "The cat sat on the mat because it was tied"; a single head might learn that "it" attends strongly to "cat".

But what if we also want:

"it" \leftrightarrow "tied" (semantic relation)

"sat" \leftrightarrow "on the mat." (syntactic relation)?

\rightarrow one head is too limited.

Multi-Head Attention (solution)

MHA solves this by projecting inputs into multiple subspaces with different W^Q, W^K, W^V .

Each head:

\Rightarrow Sees the same sentence but from a different angle.

\Rightarrow Learns different attention patterns.

How multiple heads capture different relations:

Each head has its own learned weight matrices (W^Q, W^K, W^V).

So Head 1 might align dimensions to capture Syntax (who relates to who structurally).

Head 2 might align dimensions to capture Semantics (meaning similarity)

Head 3 might look for long-range dependencies (eg: start \leftrightarrow end).

Think of them like different spotlights on the sentence, each lighting up a different pattern.

Steps Involved in Multi-Head Attention:-

1. Linear Transformation: The input x is projected into multiple smaller-dimensional subspaces using different weight matrices.

$$Q_i = x W_i^Q, \quad K_i = x W_i^K, \quad V_i = x W_i^V$$

where i denotes the head index.

2. Independent Attention Computation: Each head independently computes its own self-attention using the scaled dot-product formula.

3. Concatenation: The outputs from all heads are concatenated.

4. Final Linear Transformation: A final weight matrix is applied to transform the concatenated output into the desired dimension. (W^O)

Purpose of Final Linear Transformation W^O :-

1. Dimensionality Restoration

\Rightarrow After multi-head attention, the output shape is $h \cdot d_k$

\Rightarrow But the model expects output in shape d_{model}

\Rightarrow So $W^O \in \mathbb{R}^{(h \cdot d_k) \times d_{model}}$ projects it back.

2. Learned Fusion of Head outputs

- ⇒ Each head captures different relationships (Syntax, Semantics, Position)
- ⇒ W^O learns how to combine these diverse perspectives.
- ⇒ It's not just resizing - it's a trainable mixer.

3. Information Integration

- ⇒ Without W^O , head outputs remain isolated.
- ⇒ W^O blends them into a unified representation for the next layer.
- ⇒ Enables cross-head synergy and richer contextual understanding.

Analogy :-

- # The Final Layer is like a team lead who reads all their reports and writes a cohesive summary. It decides which insights are most relevant for the next layer.

Example Calculation:

Assume,

Input Sequence: "the", "cat", "sat"

Embedding Dimension (d_{model}): $d_{\text{model}} = 4$

Number of heads (h): $h = 2$

Dimension of each head (d_k, d_v): $d_k = d_v =$

$$d_{\text{model}} / h = 4 / 2 = 2$$

Step 1: Define Input Embeddings and Weight Matrices

$$X_{\text{the}} = [1, 0, 1, 0]$$

$$X_{\text{cat}} = [0, 1, 0, 1]$$

$$X_{\text{sat}} = [1, 1, 0, 0]$$

Next, we define the weight matrices for the two heads. These matrices are of size $d_{\text{model}} \times d_k$ (ie, 4×2) for query, key and value, and $d_{\text{model}} \times d_{\text{model}}$ for the final projection.

Head 1 weights:

$$W_1^Q = \begin{pmatrix} 0.1 & 0.2 \\ 0.3 & 0.4 \\ 0.5 & 0.6 \\ 0.7 & 0.8 \end{pmatrix}$$

$$W_1^{QK} = \begin{pmatrix} 0.8 & 0.7 \\ 0.6 & 0.5 \\ 0.4 & 0.3 \\ 0.2 & 0.1 \end{pmatrix}$$

$$W_1^K = \begin{pmatrix} 0.1 & 0.1 \\ 0.2 & 0.2 \\ 0.3 & 0.3 \\ 0.4 & 0.4 \end{pmatrix}$$

Head 2 weights:

$$W_2^Q = \begin{pmatrix} 0.9 & 0.8 \\ 0.7 & 0.6 \\ 0.5 & 0.4 \\ 0.3 & 0.2 \end{pmatrix}$$

$$W_{K2}^K = \begin{pmatrix} 0.2 & 0.3 \\ 0.4 & 0.5 \\ 0.6 & 0.7 \\ 0.8 & 0.9 \end{pmatrix}$$

$$W_2^V = \begin{pmatrix} 0.5 & 0.5 \\ 0.6 & 0.6 \\ 0.7 & 0.7 \\ 0.8 & 0.8 \end{pmatrix}$$

Final Projection Matrix:

Formula for W^O :

$$W^O \in \mathbb{R}^{(h \cdot d_K) \times d_{\text{model}}}$$

For Example:

$$d_{\text{model}} = 512$$

$$h = 8$$

$$\text{So each head as } d_K = 64 \text{ (since } 512/8 = 64)$$

Then,

$$W^O \in \mathbb{R}^{512 \times 512}$$

In our Example calculation Example,

$$d_{\text{model}} = 4$$

$$h = 2$$

$$d_K = 2 \text{ (since } 4/2 = 2)$$

Then,

$$W^O \in \mathbb{R}^{4 \times 4}$$

So our Example projection Matrix is,

$$K^0 = \begin{pmatrix} 0.1 & 0.2 & 0.3 & 0.4 \\ 0.5 & 0.6 & 0.7 & 0.8 \\ 0.9 & 1.0 & 1.1 & 1.2 \\ 1.3 & 1.4 & 1.5 & 1.6 \end{pmatrix}$$

Step 2: Calculate Q, K and V Vectors for Each Head.

Head 1:

$$\begin{aligned} \Rightarrow Q_{\text{the}, 1} &= X_{\text{the}} \cdot K^Q_1 = [1, 0, 1, 0] \begin{pmatrix} 0.1 & 0.2 \\ 0.3 & 0.4 \\ 0.5 & 0.6 \\ 0.7 & 0.8 \end{pmatrix} \\ &= [0.1 + 0.5, 0.2 + 0.6] \\ &= [0.6, 0.8] \end{aligned}$$

$$\begin{aligned} \Rightarrow K_{\text{the}, 1} &= X_{\text{the}} \cdot K^K_1 = [1, 0, 1, 0] \begin{pmatrix} 0.8 & 0.7 \\ 0.6 & 0.5 \\ 0.4 & 0.3 \\ 0.2 & 0.1 \end{pmatrix} \\ &= [0.8 + 0.4, 0.7 + 0.3] \\ &= [1.2, 1.0] \end{aligned}$$

$$\begin{aligned} \Rightarrow V_{\text{the}, 1} &= X_{\text{the}} \cdot K^V_1 = [1, 0, 1, 0] \begin{pmatrix} 0.1 & 0.1 \\ 0.2 & 0.2 \\ 0.3 & 0.3 \\ 0.4 & 0.4 \end{pmatrix} \\ &= [0.1 + 0.3, 0.1 + 0.3] \\ &= [0.4, 0.4] \end{aligned}$$

$$\Rightarrow Q_{\text{cat}, 1} = X_{\text{cat}} \cdot K^Q_1 = [0, 1, 0, 1] (\dots)$$

$$= [0.3 + 0.7, 0.4 + 0.8] = [1.0, 1.2]$$

$$\Rightarrow K_{cat, 1} = X_{cat} \cdot k_1^K = [0, 1, 0, 1] (\dots)$$

$$= [0.6 + 0.2, 0.5 + 0.1] = [0.8, 0.6]$$

$$\Rightarrow V_{cat, 1} = X_{cat} \cdot k_1^V = [0, 1, 0, 1] (\dots)$$

$$= [0.2 + 0.4, 0.2 + 0.4] = [0.6, 0.6]$$

$$\Rightarrow Q_{sat, 1} = X_{sat} \cdot k_1^Q = [1, 0, 1, 0, 0] (\dots)$$

$$= [0.1 + 0.3, 0.2 + 0.4] = [0.4, 0.6]$$

$$\Rightarrow QK_{sat, 1} = X_{sat} \cdot k_1^K = [1, 1, 0, 0] (\dots)$$

$$= [0.8 + 0.6, 0.7 + 0.5] = [1.4, 1.2]$$

$$\Rightarrow V_{sat, 1} = X_{sat} \cdot k_1^V = [1, 1, 0, 0] = [0.1 + 0.2,$$

$$0.1 + 0.2] = [0.3, 0.3]$$

Head 2: (calculations are similar, just with Head 2 weights)

$$\Rightarrow Q_{the, 2} = [1.4, 1.2],$$

$$K_{the, 2} = [0.8, 1.0]$$

$$V_{the, 2} = [1.2, 1.2]$$

$$\Rightarrow Q_{cat, 2} = [1.0, 0.8]$$

$$KQ_{cat, 2} = [1.2, 1.4]$$

$$V_{cat, 2} = [1.4, 1.4]$$

$$\Rightarrow Q_{\text{sat}, 2} = [1.6, 1.4]$$

$$K_{\text{sat}, 2} = [0.6, 0.8]$$

$$V_{\text{sat}, 2} = [1.1, 1.1]$$

Step 3: Calculate Attention Scores:

The attention scores are calculated as $\text{score} = \frac{Q \cdot K^T}{\sqrt{d_K}}$

$$\text{For } d_K = 2, \sqrt{d_K} = \sqrt{2} \approx 1.414$$

Head 1:

$$\Rightarrow S_{\text{the}} = Q_{\text{the}, 1} \cdot K_1^T = [0.6, 0.8] \cdot \begin{pmatrix} 1.2 & 0.8 & 1.4 \\ 1.0 & 0.6 & 1.2 \end{pmatrix}$$

$$= [0.6(1.2) + 0.8(1.0), 0.6(0.8) + 0.8(0.6), 0.6(1.4) + 0.8(1.2)]$$

$$\Rightarrow S_{\text{the}} = [1.52, 0.96, 1.8]$$

$$\Rightarrow S_{\text{cat}} = Q_{\text{cat}, 1} \cdot K_1^T = [1.0, 1.2] \cdot \begin{pmatrix} 1.2 & 0.8 & 1.4 \\ 1.0 & 0.6 & 1.2 \end{pmatrix}$$

$$= [2.4, 1.52, 2.84]$$

$$\Rightarrow S_{\text{sat}} = Q_{\text{sat}, 1} \cdot K_1^T = [0.4, 0.6] \cdot \begin{pmatrix} 1.2 & 0.8 & 1.4 \\ 1.0 & 0.6 & 1.2 \end{pmatrix}$$

$$= [1.08, 0.68, 1.28]$$

Now we divide each score by $\sqrt{2}$:

$$\Rightarrow \text{Scaled-}S_{\text{the}, 1} = [1.07, 0.68, 1.27]$$

$$\Rightarrow \text{Scaled-}S_{\text{cat}, 1} = [1.70, 1.07, 2.01]$$

$$\Rightarrow \text{Scaled-}S_{\text{sat}, 1} = [0.76, 0.48, 0.91]$$

Head 2 :

$$\Rightarrow \text{Scaled-}\delta_{\text{the}}, 2 = [1.06, 2.11, 0.88]$$

$$\Rightarrow \text{Scaled-}\delta_{\text{cat}}, 2 = [2.11, 2.76, 1.63]$$

$$\Rightarrow \text{Scaled-}\delta_{\text{sat}}, 2 = [0.88, 1.63, 0.77]$$

Step 4: Apply Softmax to Get Attention weights

We apply the Softmax function to each row of the Scaled Score matrix. Softmax is $\frac{e^{x_i}}{\sum e^{x_i}}$

Head 1:

$$\Rightarrow w_{\text{the}}, 1 = \text{Softmax}([1.07, 0.68, 1.27])$$

$$= [0.43, 0.29, 0.28]$$

$$\Rightarrow w_{\text{cat}}, 1 = \text{Softmax}([1.70, 1.07, 2.01])$$

$$= [0.38, 0.20, 0.42]$$

$$\Rightarrow w_{\text{sat}}, 1 = \text{Softmax}([0.76, 0.48, 0.91])$$

$$= [0.39, 0.29, 0.32]$$

Head 2:

$$\Rightarrow w_{\text{the}}, 2 = \text{Softmax}([1.06, 2.11, 0.88])$$

$$= [0.24, 0.66, 0.10]$$

$$w_{\text{cat}}, 2 = \text{Softmax}([2.11, 2.76, 1.63])$$

$$= [0.30, 0.57, 0.13]$$

$$w_{\text{sat}}, 2 = \text{Softmax}([0.88, 1.63, 0.77])$$

$$= [0.27, 0.56, 0.17]$$

Step 5: Calculate the Weighted Sum of Values

Now we multiply the attention weights by the Value Vectors of each head.

Head 1 output (I_1):

$$\Rightarrow I_{the,1} = 0.43 \cdot V_{the,1} + 0.29 \cdot V_{cat,1} + 0.28 \cdot$$

$$V_{sat,1} = 0.43 [0.4, 0.4] + 0.29 [0.6, 0.6] + 0.28 [0.3, 0.3] = [0.172 + 0.174 + 0.084,$$

$$0.172 + 0.174 + 0.084] =$$

$$[0.43, 0.43]$$

$$\Rightarrow I_{cat,1} = 0.38 \cdot V_{the,1} + 0.20 \cdot V_{cat,1} + 0.42 \cdot V_{sat,1}$$

$$= [0.38(0.4) + 0.2(0.6) + 0.42(0.3), \dots]$$

$$= [0.34, 0.34]$$

$$\Rightarrow I_{sat,1} = 0.39 \cdot V_{the,1} + 0.29 \cdot V_{cat,1} + 0.32 \cdot V_{sat,1}$$

$$= [0.39(0.4) + 0.29(0.6) + 0.32(0.3), \dots]$$

$$= [0.41, 0.41]$$

Head 2 output (I_2):

$$\Rightarrow I_{the,2} = 0.24 \cdot V_{the,2} + 0.66 \cdot V_{cat,2} + 0.10 \cdot V_{sat,2} =$$

$$0.24 [1.2, 1.2] + 0.66 [1.4, 1.4] + 0.10 [1.1, 1.1]$$

$$= [0.288 + 0.924 + 0.11, 0.288 + 0.924 + 0.11]$$

$$= [1.322, 1.322]$$

$$\Rightarrow I_{cat, 2} = 0.30 \cdot \sqrt{the, 2} + 0.57 \cdot \sqrt{cat, 2} + 0.13 \cdot \sqrt{sat, 2}$$

$$\sqrt{sat, 2} = [0.3(1.2) + 0.57(1.4) + 0.13(1.1) \dots] = [1.36, 1.36]$$

$$\Rightarrow I_{sat, 2} = 0.27 \cdot \sqrt{the, 2} + 0.56 \cdot \sqrt{cat, 2} + 0.17 \cdot \sqrt{sat, 2}$$

$$\sqrt{sat, 2} = [0.27(1.2) + 0.56(1.4) + 0.17(1.1) \dots] = [1.30, 1.30]$$

Step 6: Concatenate and project

We concatenate the outputs of the two heads for each word.

$$\Rightarrow I_{the} = [I_{the, 1}, I_{the, 2}] = [0.43, 0.43, 1.322, 1.322]$$

$$\Rightarrow I_{cat} = [I_{cat, 1}, I_{cat, 2}] = [0.34, 0.34, 1.36, 1.36]$$

$$\Rightarrow I_{sat} = [I_{sat, 1}, I_{sat, 2}] = [0.41, 0.41, 1.30, 1.30]$$

Finally we multiply the concatenated vectors by the final projection matrix W^o

$$\Rightarrow \text{Output}_{the} = I_{the} \cdot W^o = [0.43, 0.43, 1.322, 1.322]$$

$$\begin{pmatrix} 0.1 & 0.2 & 0.3 & 0.4 \\ 0.5 & 0.6 & 0.7 & 0.8 \\ 0.9 & 1.0 & 1.1 & 1.2 \\ 1.3 & 1.4 & 1.5 & 1.6 \end{pmatrix}$$

$$= [2.63, 2.92, 3.2, 3.5]$$

$$\Rightarrow \text{Output}_{\text{cat}} = I_{\text{cat}} \cdot k^0 = [0.34, 0.34, 1.36, 1.36]$$

$$(\dots)$$

$$= [2.67, 2.97, 3.27, 3.57]$$

$$\text{Output}_{\text{sat}} = I_{\text{sat}} \cdot k^0 = [0.41, 0.41, 1.80, 1.80]$$

$$(\dots)$$

$$= [2.58, 2.87, 3.16, 3.45]$$

The Final output is a new 4-dimensional Vector for each word in the sequence, which now contains a richer representation that has been informed by all other words.

Note :-

In Original Transformer paper, they Experimented with different numbers of heads.

In the Base Model :

$$\text{Number of heads} = 8$$

$$d_{\text{model}} = 512 \text{ (Embedding Size)}$$

$$\text{Each head dimension} = 64 \text{ (Since } 512/8 = 64 \text{)}$$

In the Large Model :

$$\text{Number of heads} = 16$$

$$d_{\text{model}} = 1024$$

$$\text{Each head dimension} = 64$$

$$= [2.63, 2.92, 3.2, 3.5]$$

$$\Rightarrow \text{Output}_{\text{cat}} = \mathbf{I}_{\text{cat}} \cdot \mathbf{K}^0 = [0.34, 0.34, 1.36, 1.36]$$

$$(\dots)$$

$$= [2.67, 2.97, 3.27, 3.57]$$

$$\text{Output}_{\text{sat}} = \mathbf{I}_{\text{sat}} \cdot \mathbf{K}^0 = [0.41, 0.41, 1.30, 1.36]$$

$$(\dots)$$

$$= [2.58, 2.87, 3.16, 3.45]$$

The Final output is a new 4-dimensional vector for each word in the sequence, which now contains a richer representation that has been informed by all other words.

Note :-

- # In Original Transformer paper, they Experimented with different numbers of heads.
- # In the Base Model:

$$\text{Number of heads} = 8$$

$$d_{\text{model}} = 512 \text{ (Embedding Size)}$$

$$\text{Each head dimension} = 64 \text{ (Since } 512/8 = 64 \text{)}$$

- # In the Large Model:

$$\text{Number of heads} = 16$$

$$d_{\text{model}} = 1024$$

$$\text{Each head dimension} = 64$$