

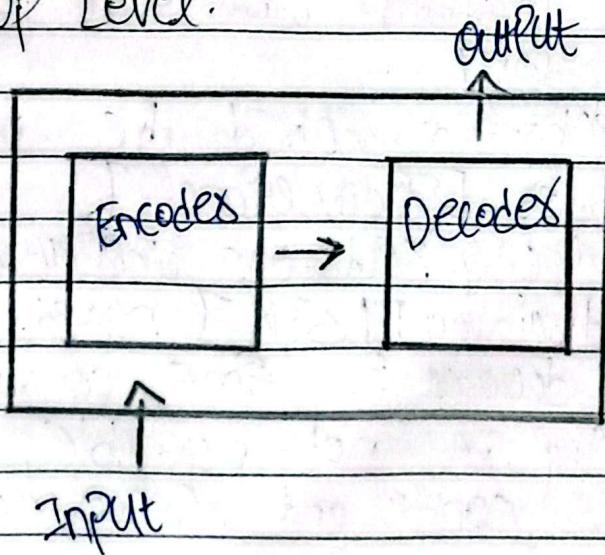
(*)

Recent improvements sector for Transformers.

- Padding
- Quantization
- Knowledge Distillation.
- Improve the multimodal capabilities
- Multi-lingual capabilities.

Transformer Architecture

(*) TOP Level.



• Encodes

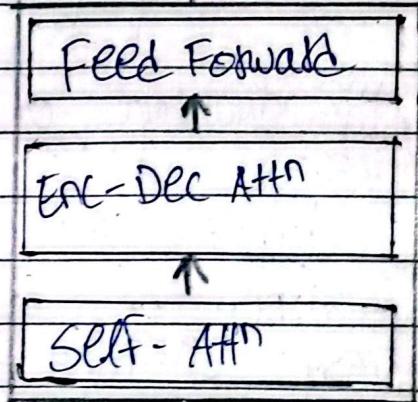
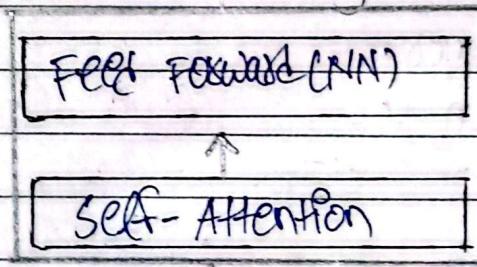
→ contains multiple Encodes.

• Decodes

→ contains multiple Decodes.

Decodes - Single

(*) Bit Detail (Fig)



Encodes - single

In the each [Encodes], takes the sequence, in parallel as input. Self Attr layer processes those sequence, this self Attr layer converts the [Raw Vector] (embedding) to a different vector known as [Contextual Vector].

(*) Contextual Vector [Contextual Embedding]:
Vector that keeps the context of a word or token relative to other token, or word.

- the generated [Contextual vector] will be passed as an input to the another Encodes in the stack.

(*) Self Attention layers

Basically it's just a function that takes the Raw Embedding [Static] and generate Dynamic Contextual Embedding.

↳ If we just use static word embedding we don't have the power to change the embeddings based on the context [Refer to apple example].

- Self Attention provides the context of each word relative to other words in a sentence. This is what we call Attention, now we can know which word has the maximum impact for understanding the [meaning] of the sentence.

Bank here, has the same [Word Embedding]
but, Self Attn creates dynamic [Contextual Embedding]

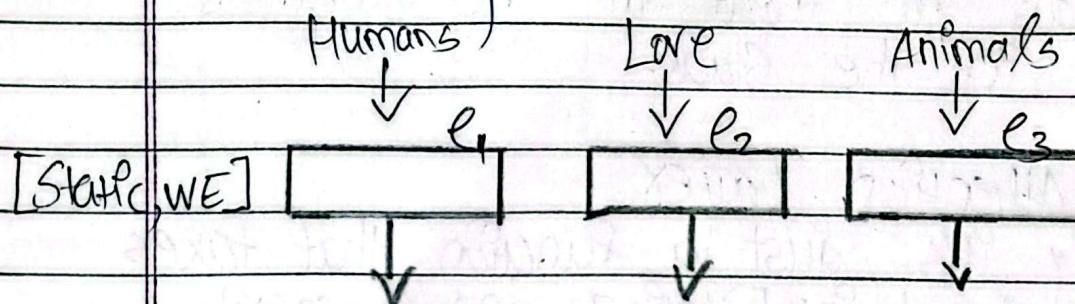
e.g [Rivers Bank], [Money Bank]

Here, the word bank can be confusing
so it is very important we need to
build a model that captures or takes
in consideration to its surrounding context
to make or understand the actual meaning
of the sentence.

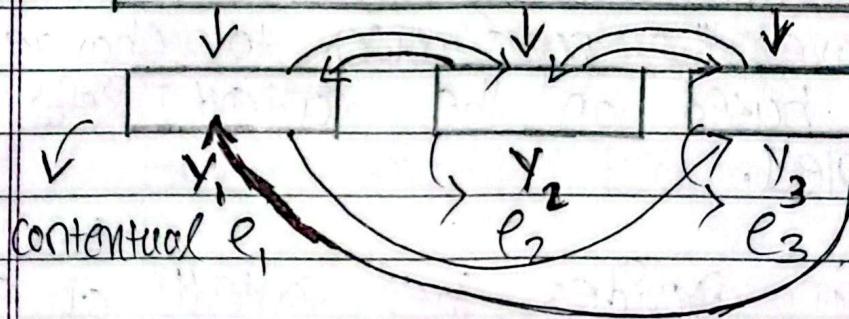
↳ Self Attention is what we're looking
for.



Understanding Self Attention.



calculations [Self Attn]



Let's take an example:

⇒ Two sentences

$S_1 \Rightarrow$ [Money Bank Grows] $S_2 \Rightarrow$ [Rivers Bank Flows]

Here, for static Word Embedding for
Bank is same.

If we use static word embedding we'll not get our desired output because both the embeddings are same.

But, what if we could represent each word relative to it's importance of other word.

So, for S_1

$$\text{Money} \Rightarrow 0.7 \text{ Money} + 0.2 \text{ Bank} + 0.1 \text{ grows}$$

$$\text{Bank} \Rightarrow 0.25 \text{ Money} + 0.7 \text{ Bank} + 0.05 \text{ grows}$$

$$\text{Grows} \Rightarrow 0.1 \text{ Money} + 0.2 \text{ Bank} + 0.7 \text{ grows}$$

and, S_2

$$\text{Rives} \Rightarrow 0.8 \text{ Rives} + 0.15 \text{ Bank} + 0.05 \text{ Flows}$$

$$\text{Bank} \Rightarrow 0.2 \text{ Rives} + 0.78 \text{ Bank} + 0.02 \text{ Flows}$$

$$\text{Flows} \Rightarrow 0.4 \text{ Rives} + 0.01 \text{ Bank} + 0.59 \text{ Flows}$$

This way we're representing each word in the form (relative) to the importance of other word.

And, if we compare two [Bank] words.
We've different Embeddings ~~for~~ each sentence.

Now, in both the sentence the word [Bank] will have new embedding (that covers the context or takes in consideration of other words as well).

- The numbers used here, such as 0.8, 0.1 etc represents the similarity between words.

We can write,

New Embedding For [Bank] as:

$$\text{e}_{\text{bank}}^{\text{new}} = [\text{e}_{\text{bank}} \cdot \text{e}^{\text{Tmoney}}] \text{money} + [\text{e}_{\text{bank}} \cdot \text{e}^{\text{Tgrows}}] \text{grows}$$

- where the score or number we used before is Dot Product i.e Represents Similarity.

- The [Scalar] represents the Similarity for two words.

↳ All the Dot Product are scaled i.e [Normalized] i.e [0 to 1] which also represents the probability.

- Softmax is used to calculate the probability for each [Dot Product].

$$f(x_i) = \left[\frac{e^{x_i}}{\sum_{i=0}^n e^{x_i}} \right] \Rightarrow \text{Softmax}$$

↳ Once we get this probability we'll multiply with the Raw of original vector.

↳ Then sum the each multiplication to represent the vector (New Embedding)

(*) Problem with Simple Self Att

- As our self-Attention layer remains static i.e there's no any trainable parameter it will have no idea about how to adapt the Attention for each word or sequence.
- ↳ So it fails to adapt to the specific patterns or nuances present in the [Training] data required for a specific task i.e "sentiment", "Translation".

To fix this, we introduce [Query, key, and value]

- For each word or embedding [3] new vectors i.e [Query vector, Key vector] and [Value vector].

So, for embedding [Bank] we'll have

$Q_{\text{bank}} = [E_{\text{bank}} \times W_Q]$
 $V_{\text{bank}} = [E_{\text{bank}} \times W_V]$

$$Q_{\text{bank}} = E_{\text{bank}} \times W_Q$$

$$Q_{\text{money}} = E_{\text{money}} \times W_Q \quad \text{so, } Q = \begin{bmatrix} Q_{\text{bank}} \\ Q_{\text{money}} \\ Q_{\text{cow}} \end{bmatrix}$$

$$Q_{\text{cow}} = E_{\text{cow}} \times W_Q$$

- If our input seq. is of 3 we'll have 3x3 vectors, i.e. 3 Query, 3 key, 3 value vectors

This [Q] matrix is collection of all our query in the sequence,

Then similarly, for [key]

$$K_{bank} = E_{bank} \times W_k$$

$$K_{money} = E_{money} \times W_k \text{ so, } K = \begin{bmatrix} K_{bank} \\ K_{money} \\ K_{flows} \end{bmatrix}$$

$$K_{flows} = E_{flows} \times W_k$$

Also, for [value]

$$V_{bank} = E_{bank} \times W_k$$

$$V_{money} = E_{money} \times W_k \text{ so, } V = \begin{bmatrix} V_{bank} \\ V_{money} \\ V_{flow} \end{bmatrix}$$

Now, once we've all the 3 vectors i.e. query, key, value.

We'll need to follow these steps to generate contextual Embedding:

- Calculate Dot Product b/w Query and key matrix for similarity check.

$$\text{i.e. } S = Q \cdot K^T \quad (3 \times 3, 3 \times 3) \Rightarrow 3 \times 3$$

- Apply Softmax to the Dot Product output matrix (Normalized, same unit)
(Scaled Dot Product Attention)

$$A = \begin{bmatrix} \text{Softmax}(a_1, a_2, a_3) \\ \text{Softmax}(b_1, b_2, b_3) \\ \text{Softmax}(c_1, c_2, c_3) \end{bmatrix}$$

(iii) Compute the weighted sum i.e Final Attention scores

$$C = A \cdot V \quad (\text{softmax output times the value matrix})$$

Now, this C matrix is the [Contentual Embedding matrix] for all the words in a sequence.

This C will be the output of a single Self-Attention layer.

Here, W_Q , W_K , and W_V are the weights for the query, key and value.

- The W_Q , W_K , and W_V weights will be changing throughout the training with Backpropagation.
- These weights are randomly initialized before training.