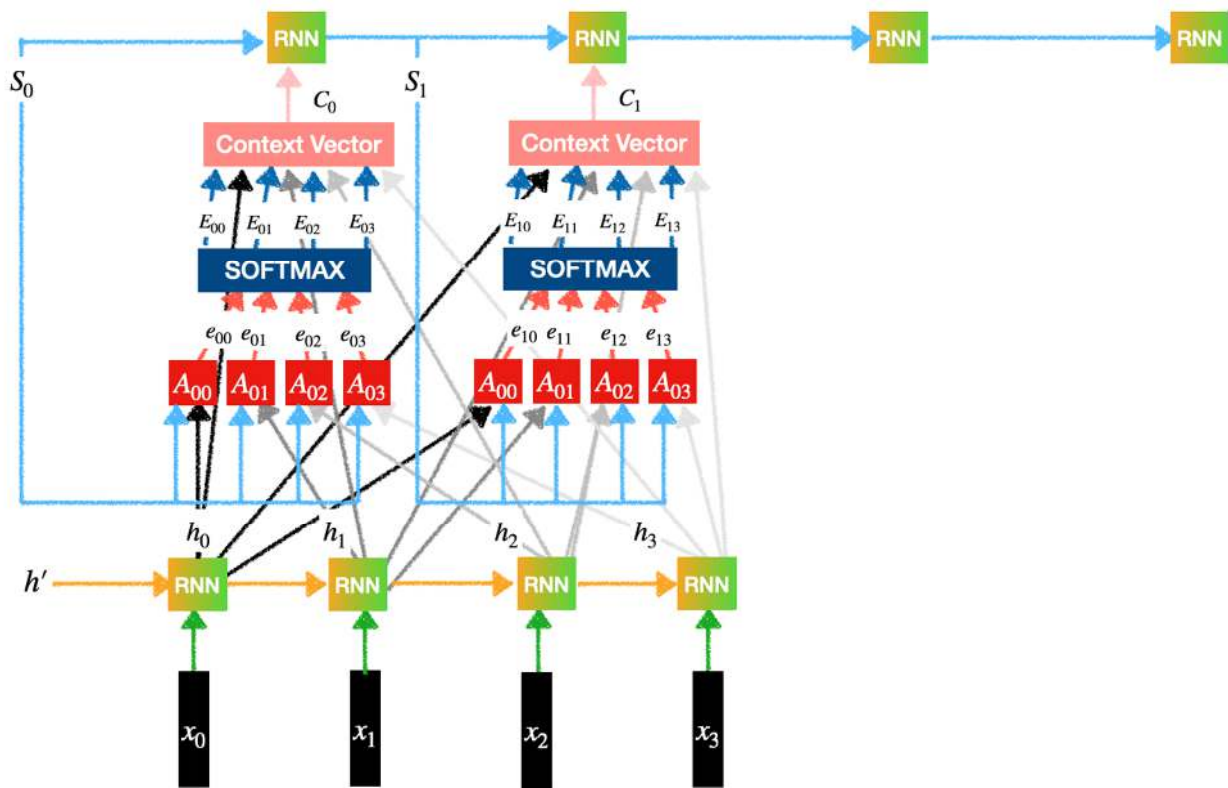


## ML - Attention mechanism

### Attention Mechanism Architecture for Machine Translation

Attention mechanism in machine translation has three main components: **Encoder**, **Attention** and **Decoder**. Here's how each of these components works:

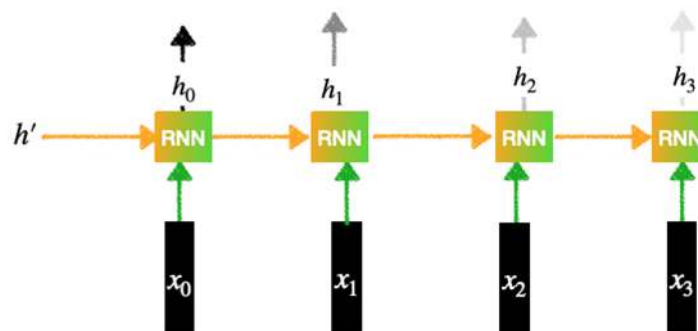


Encoder-Decoder with Attention

### 1. Encoder:

Encoder processes the input sequence which is usually a sentence and generates hidden states. It uses Recurrent Neural Networks (RNNs), LSTMs, GRUs or transformer-based models to process the input step by step. At each step encoder produces a hidden state that contains both the data from the previous hidden state and the current input token. This allows the encoder to capture the relationships between the tokens in the input sequence. As the input sequence is processed encoder generates a series of hidden states which together represent the entire input sequence. These hidden states are passed on to the attention mechanism for further processing.

It is important because it creates a representation of the input sequence that the attention mechanism and decoder will later use to generate an accurate output sequence.



Encoder

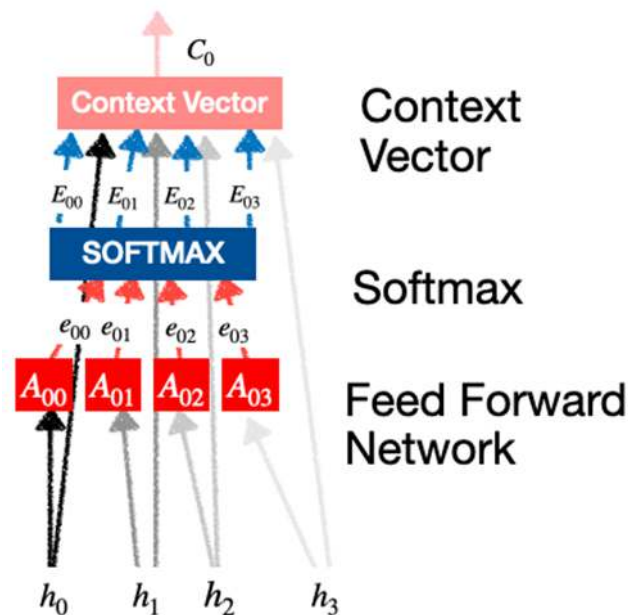
### Contains an RNN layer (Can be LSTMs or GRU):

- Let's say, there are 4 words sentence then inputs will be:  $x_0, x_1, x_2, x_3$
- Each input goes through an Embedding Layer, It can be RNN, LSTM, GRU or transformers
- Each of the inputs generates a hidden representation.
- This generates the outputs for the Encoder:  $h_0, h_1, h_2, h_3$

## 2. Attention:

Attention component find importance of each encoder's hidden state with respect to the current target hidden state. It generates a context vector that captures the relevant information from the encoder's hidden states. Its mechanism can be represented mathematically as follows:

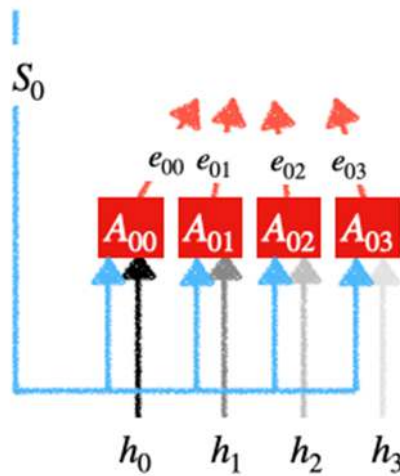
- Our goal is to generate the context vectors.
- For example, context vector  $C_1$  tells us how much importance, attention should be given the inputs:  $x_0, x_1, x_2, x_3$
- This layer in turn contains 3 subparts:
  1. Feed Forward Network
  2. Softmax Calculation
  3. Context vector generation



Attention

### Feed Forward Network:

It transforms target hidden state into a form that is compatible with the attention mechanism. This transformation is done using a linear transformation followed by a non-linear activation function like ReLU or sigmoid. This step helps in finding alignment between target and the encoder's hidden states.



Feed-Forward-Network

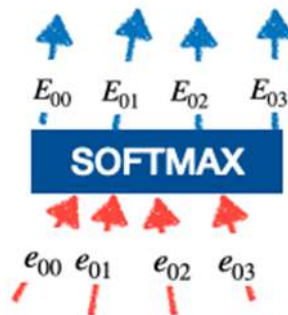
Each  $A_{00}, A_{01}, A_{02}, A_{03}$  is a simple feed-forward neural network with one hidden layer. The input for this feed-forward network is:

- Previous Decoder state
- The output of Encoder states.

Each unit generates outputs:  $e_{00}, e_{01}, e_{02}, e_{03}$  i.e  $e_{0i} = g(S_0, h_i)$ . Here  $g$  can be any activation function such as sigmoid, tanh or ReLu.

#### Softmax Calculation:

Once the feed-forward network generates the relevant scores these scores are passed through a softmax function. It converts scores into probability-like attention weights. These weights indicate the relative importance of each encoder's hidden state in generating current target token. For example higher attention weights shows that a particular encoder's hidden state is more important for the current step in the decoding process.



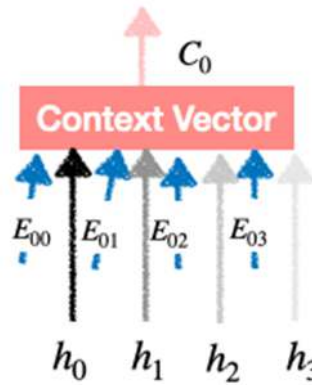
softmax calculation

$$E_{0i} = \frac{\exp(e_{0i})}{\sum_{i=0}^3 \exp(e_{0i})}$$

These  $E_{00}, E_{01}, E_{02}, E_{03}$  are called the attention weights. It decides how much importance should be given to the inputs  $x_0, x_1, x_2, x_3$ .

#### Context Vector Generation:

When the attention weights are calculated they are applied to encoder's hidden states which helps in generating a weighted sum. This weighted sum forms the **context vector** which finds most relevant information from the encoder's hidden states based on their attention weights.



context vector generation

$$C_0 = E_{00} * h_0 + E_{01} * h_1 + E_{02} * h_2 + E_{03} * h_3$$

We find  $C_1, C_2, C_3$  in the same way and feed it to different RNN units of the Decoder layer. So this is the final vector which is the product of Probability Distribution and Encoder's output which is nothing but the attention paid to the input words.

### 3. Decoder:

It receives the context vector from the attention layer which contains relevant information from the encoder's hidden states along with its own current hidden state. Using this combined information decoder predicts next token in the output sequence. Decoder's hidden state is then updated based on the predicted token and this process repeats again and again for each token in the output sequence. This helps model to generate entire output sequence like translation word by word with each step focusing on the most important parts of the input.

#### Applications of attention mechanisms

1. **Machine Translation:** It allows model to focus on different parts of the source sentence when generating each word in the target sentence which improves translation accuracy.
2. **Sentiment Analysis & Named Entity Recognition:** In tasks like sentiment analysis, question answering and named entity recognition it helps models to focus on key words that helps in entity identification.
3. **Text Summarization:** It helps in selecting important information by helping model to create short and accurate summaries of large text.
4. **Image Captioning:** It help image captioning models to focus on specific image part while generating captions which improves description's detail.

Attention mechanism will keep getting better and making models more efficient and useful. This will help them perform well in more tasks and handle even more complex problems.

### Self Attention Mechanism

This mechanism captures long range dependencies by calculating attention between all words in the sequence and helping the model to look at the entire sequence at once. Unlike traditional models that process words one by one it helps the model to find which words are most relevant to each other helpful for tasks like translation or text generation.

Here's how the self attention mechanism works:

1. **Input Vectors and Weight Matrices:** Each encoder input vector is multiplied by three trained weight matrices ( $W(Q)$ ,  $W(K)$ ,  $W(V)$ ) to generate the key, query and value vectors.
2. **Query Key Interaction:** Multiply the query vector of the current input by the key vectors from all other inputs to calculate the attention scores.

3. **Scaling Scores:** Attention scores are divided by the square root of the key vector's dimension ( $d_k$ ) usually 64 to prevent the values from becoming too large and making calculations unstable.
4. **Softmax Function:** Apply the softmax function to the calculated attention scores to normalize them into probabilities.
5. **Weighted Value Vectors:** Multiply the softmax scores by the corresponding value vectors.
6. **Summing Weighted Vectors:** Sum the weighted value vectors to produce the self attention output for the input.

Above procedure is applied to all the input sequences. Mathematically self attention matrix for input matrices (Q,K,V) is calculated as:

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

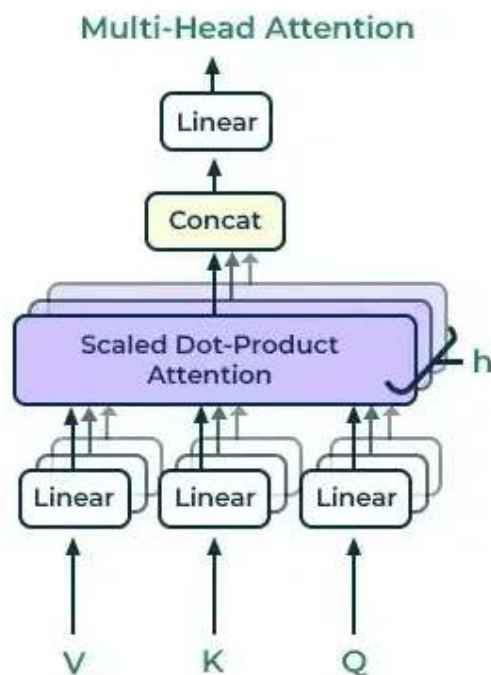
where Q,K,V are the concatenation of query, key and value vectors

### Multi Head Attention

In multi headed attention mechanism, multiple attention heads are used in parallel which allows the model to focus on different parts of the input sequence simultaneously. This approach increases model's ability to capture various relationships between words in the sequence.

$$MultiHead(Q,K,V) = concat(head_1 head_2 \dots head_n) W_O$$

Here's a step by step breakdown of how multi headed attention works:



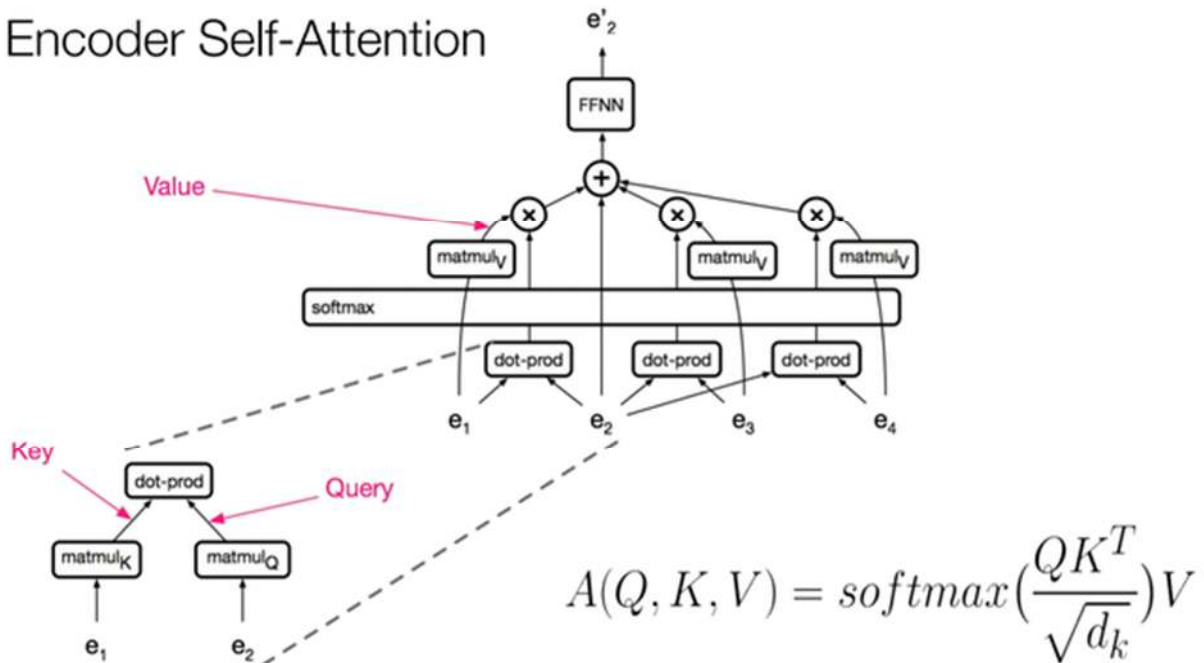
Multi-headed-attention

1. **Generate Embeddings:** For each word in the input sentence it generate its embedding representation.
2. **Create Multiple Attention Heads:** Create  $h$  (e.g.  $h=8$ ) attention heads and each with its own weight matrices  $W(Q), W(K), W(V)$ .
3. **Matrix Multiplication:** Multiply the input matrix by each of the weight matrices  $W(Q), W(K), W(V)$  for each attention head to produce key, query and value matrices.
4. **Apply Attention:** Apply attention mechanism to the key, query and value matrices for each attention head which helps in generating an output matrix from each head.
5. **Concatenate and Transform:** Concatenate the output matrices from all attention heads and apply a dot product with weight  $W_O$  to generate the final output of the multi-headed attention layer.

## Use in Transformer Architecture

- **Encoder Decoder Attention:** In this layer queries come from the previous decoder layer while the keys and values come from the encoder's output. This allows each position in the decoder to focus on all positions in the input sequence.
- **Encoder Self Attention:** This layer receives queries, keys and values from the output of the previous encoder layer. Each position in the encoder looks at all positions from the previous layer to calculate attention scores.

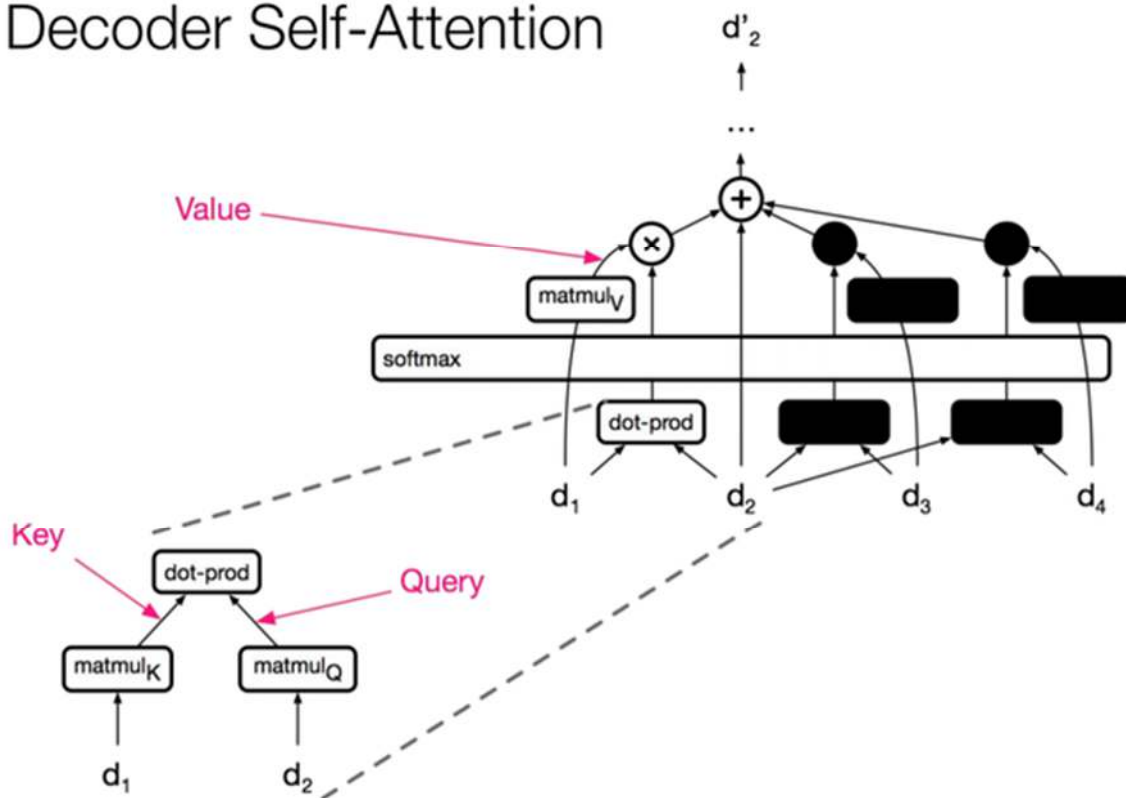
### Encoder Self-Attention



Encoder Self-Attention

- **Decoder Self Attention:** Similar to the encoder's self attention but here the queries, keys and values come from the previous decoder layer. Each position can attend to the current and previous positions but future positions are masked to prevent the model from looking ahead when generating the output and this is called masked self attention.

# Decoder Self-Attention



Decoder Self Attention

## Implementation

### Step 1: Install Necessary Libraries

This line imports numpy for matrix operations and softmax from scipy.special to convert attention scores into probability distributions.

```
import numpy as np
from scipy.special import softmax
```

### Step 2: Extract Dimensions

This function starts by extracting the input shape: batch size, sequence length and model dimension. It sets  $d_k$ , the dimension of keys and queries equal to the model dimension for simplicity.

```
def self_attention(X):
    batch_size, seq_len, d_model = X.shape
    d_k = d_model
```

### Step 3: Initialize Weight Matrix

These lines initialize random weight matrices for queries ( $W_q$ ), keys ( $W_k$ ) and values ( $W_v$ ). In real models these are learnable parameters used to project the input into Q, K, and V representations.

```
W_q = np.random.randn(d_model, d_k)
W_k = np.random.randn(d_model, d_k)
W_v = np.random.randn(d_model, d_k)
```

### Step 4: Compute Q, K, V matrices

These lines project the input X into query (Q), key (K) and value (V) matrices by multiplying with their respective weights. This transforms the input into different views used for computing attention.

```
Q = X @ W_q
```



$$K = X @ W_k$$

$$V = X @ W_v$$

### Step 5: Compute Attention scores and weights

This computes the final output by weighting the values (V) with the attention scores, aggregating relevant information from the sequence. It then returns both the attention output and the attention weights for further use or analysis.

```
output = attention_weights @ V
```

```
return output, attention_weights
```

### Step 6: Example Usage

This sets a random seed for reproducibility creates a sample input tensor with shape (1, 3, 4) runs the self-attention function on it and then prints the resulting output and attention weights.

```
np.random.seed(42)
```

```
X = np.random.rand(1, 3, 4)
```

```
output, weights = self_attention(X)
```

```
print("Output:\n", output)
```

```
print("\nAttention Weights:\n", weights)
```

**Output:**

Output:

```
[[[-1.00006959 -0.08691166  0.66098659  1.15282625]
 [-0.94070605 -0.15072265  0.6716722  1.07836487]
 [-0.94622694 -0.07915902  0.67950479  1.09794086]]]
```

Attention Weights:

```
[[[0.1910979  0.45244565 0.35645645]
 [0.24080494 0.55910874 0.20008632]
 [0.19722492 0.54711997 0.25565511]]]
```

### Advantages

1. **Parallelization:** Unlike sequential models it allows for full parallel processing which speeds up training.
2. **Long Range Dependencies:** It provides direct access to distant elements making it easier to model complex structures and relationships across long sequences.
3. **Contextual Understanding:** Each token's representation is influenced by the entire sequence which integrates global context and improves accuracy.
4. **Interpretable Weights:** Attention maps can show which parts of the input were most influential in making decisions.

### Challenges

1. **Computational Cost:** Self attention requires computing pairwise interactions between all input tokens which causes a time and memory complexity of  $O(n^2)$  where  $n$  is the sequence length. This becomes inefficient for long sequences.
2. **Memory Usage:** Large number of pairwise calculations in self attention uses high memory while working with very long sequences or large batch sizes.
3. **Lack of Local Context:** It focuses on global dependencies across all tokens but it may not effectively capture local patterns. This can cause inefficiencies when local context is more important than global context.



4. **Overfitting:** Due to its ability to model complex relationships it can overfit when it is trained on small datasets.

## Transformer Architecture

Transformers were first developed to solve the problem of sequence transduction, or neural machine translation, which means they are meant to solve any task that transforms an input sequence to an output sequence. This is why they are called “Transformers”.

- A transformer model is a neural network that learns the context of sequential data and generates new data out of it.
- *A transformer is a type of artificial intelligence model that learns to understand and generate human-like text by analyzing patterns in large amounts of text data.*

Transformers are a current state-of-the-art NLP model and are considered the evolution of the encoder-decoder architecture. However, while the encoder-decoder architecture relies mainly on Recurrent Neural Networks (RNNs) to extract sequential information, Transformers completely lack this recurrency.

They are specifically designed to comprehend context and meaning by analyzing the relationship between different elements, and they rely almost entirely on a mathematical technique called attention to do so.



At the time of the Transformer model introduction, **RNNs** were the preferred approach for dealing with sequential data, which is characterized by a specific order in its input.

RNNs function similarly to a feed-forward neural network but process the input sequentially, one element at a time.

Transformers were inspired by the encoder-decoder architecture found in RNNs. However, instead of using recurrence, the Transformer model is completely based on the Attention mechanism.

Besides improving RNN performance, Transformers have provided a new architecture to solve many other tasks, such as text summarization, image captioning, and speech recognition.

So, what are RNNs' main problems? They are quite ineffective for NLP tasks for two main reasons:

- They process the input data sequentially, one after the other. Such a recurrent process does not make use of modern graphics processing units (GPUs), which were designed for parallel computation and, thus, makes the training of such models quite slow.
- They become quite ineffective when elements are distant from one another. This is due to the fact that information is passed at each step and the longer the chain is, the more probable the information is lost along the chain.

The shift from Recurrent Neural Networks (RNNs) like LSTM to Transformers in NLP is driven by these two main problems and Transformers' ability to assess both of them by taking advantage of the Attention mechanism improvements:

- Pay attention to specific words, no matter how distant they are.
- Boost the performance speed.

Next, let's take a look at how transformers work.

## The Transformer Architecture

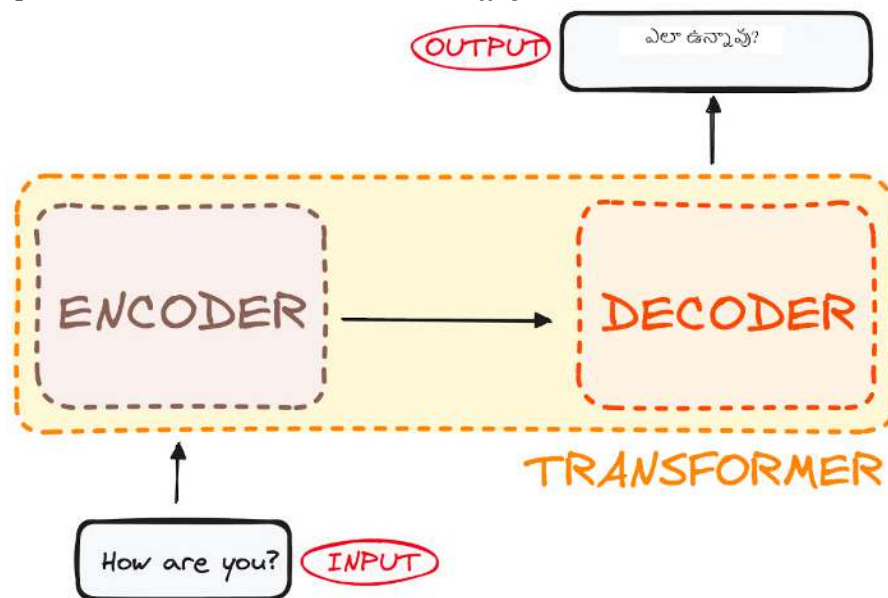
Originally devised for sequence transduction or neural machine translation, transformers excel in converting input sequences into output sequences. It is the first transduction model relying entirely on self-attention to compute representations of its input and output without using sequence-aligned RNNs or convolution. The main core characteristic of the Transformers architecture is that they maintain the encoder-decoder model.

If we start considering a Transformer for language translation as a simple black box, it would take a sentence in one language, English for instance, as an input and output its translation in Telugu.



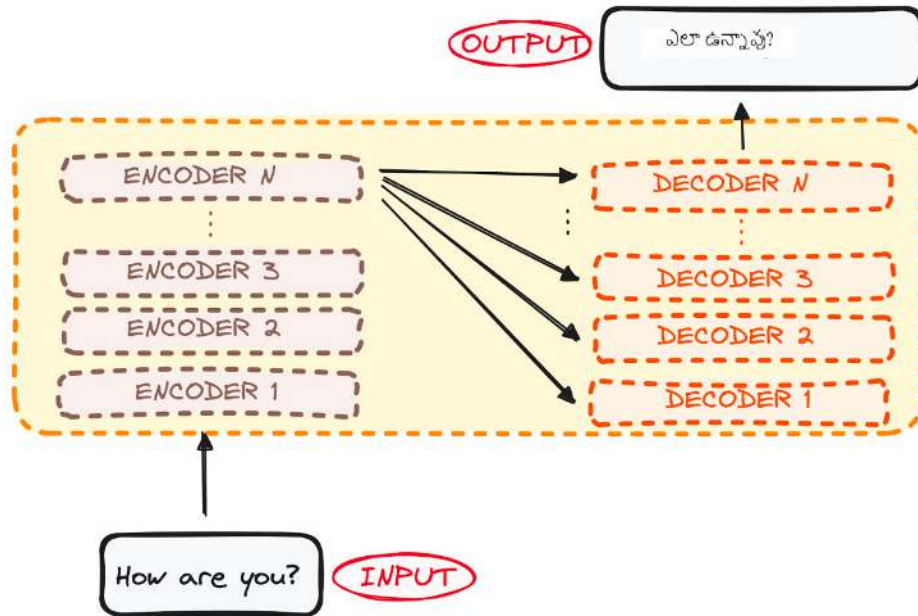
If we dive a little bit, we observe that this black box is composed of two main parts:

- The encoder takes in our input and outputs a matrix representation of that input. For instance, the English sentence "How are you?"
- The decoder takes in that encoded representation and iteratively generates an output. In our example, the translated sentence ఎలా ఉన్నావు?



However, both the encoder and the decoder are actually a stack with multiple layers (same number for each). All encoders present the same structure, and the input gets into each of them and is passed to the next one. All decoders present the same structure as well and get the input from the last encoder and the previous decoder.

The original architecture consisted of 6 encoders and 6 decoders, but we can replicate as many layers as we want. So let's assume N layers of each.



So now that we have a generic idea of the overall Transformer architecture, let's focus on both Encoders and Decoders to understand better their working flow:

### The Encoder WorkFlow

The encoder is a fundamental component of the Transformer architecture. The primary function of the encoder is to transform the input tokens into contextualized representations. Unlike earlier models that processed tokens independently, the Transformer encoder captures the context of each token with respect to the entire sequence.

Its structure composition consists as follows:

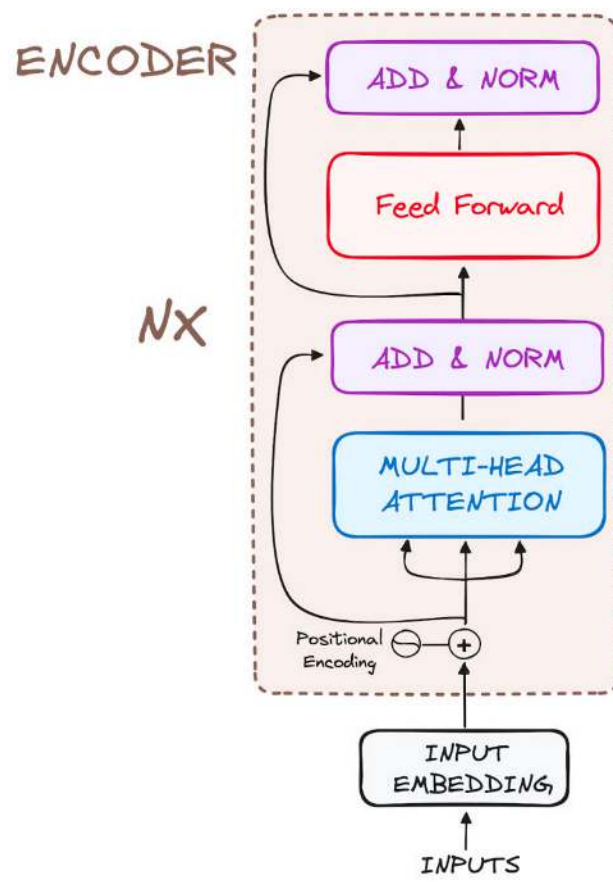


Fig: Global structure of Encoders.

So let's break its workflow into its most basic steps:

### STEP 1 - Input Embeddings

The embedding only happens in the bottom-most encoder. The encoder begins by converting input tokens - words or subwords - into vectors using embedding layers. These embeddings capture the semantic meaning of the tokens and convert them into numerical vectors.

All the encoders receive a list of vectors, each of size 512 (fixed-sized). In the bottom encoder, that would be the word embeddings, but in other encoders, it would be the output of the encoder that's directly below them.

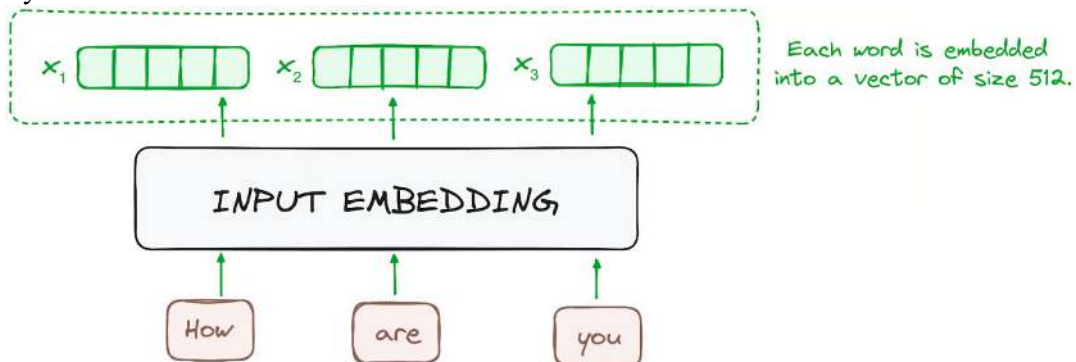
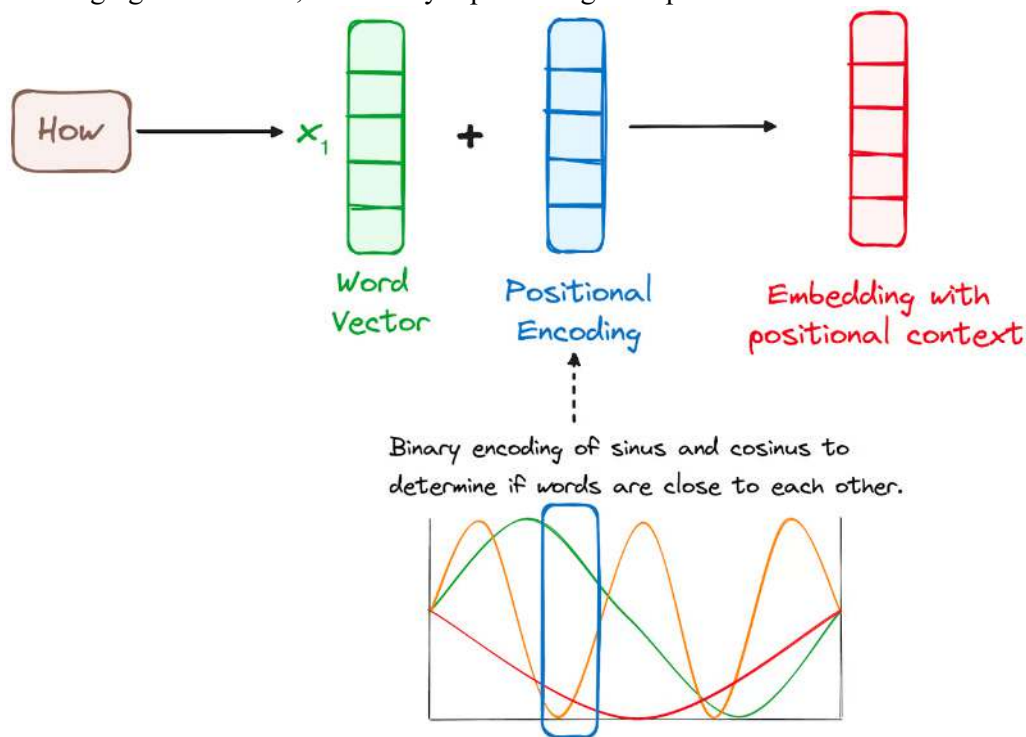


Fig:Encoder's workflow. Input embedding.

### STEP 2 - Positional Encoding

Since Transformers do not have a recurrence mechanism like RNNs, they use positional encodings added to the input embeddings to provide information about the position of each token in the sequence. This allows them to understand the position of each word within the sentence. To do so, the researchers suggested employing a combination of various sine and cosine functions to create positional vectors, enabling the use of this positional encoder for sentences of any length. In this approach, each dimension is represented by unique frequencies and offsets of the wave, with the values ranging from -1 to 1, effectively representing each position.



*Fig: Encoder's workflow. Positional encoding.*

### STEP 3 - Stack of Encoder Layers

The Transformer encoder consists of a stack of identical layers (6 in the original Transformer model).

The encoder layer serves to transform all input sequences into a continuous, abstract representation that encapsulates the learned information from the entire sequence. This layer comprises two sub-modules:

- A multi-headed attention mechanism.
- A fully connected network.

Additionally, it incorporates residual connections around each sublayer, which are then followed by layer normalization.

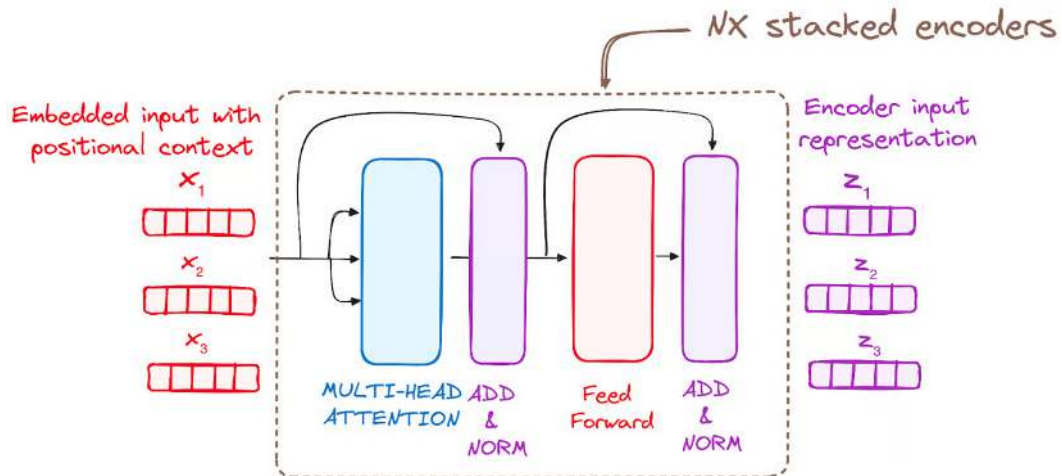


Fig: Encoder's workflow. Stack of Encoder Layers

### STEP 3.1 Multi-Headed Self-Attention Mechanism

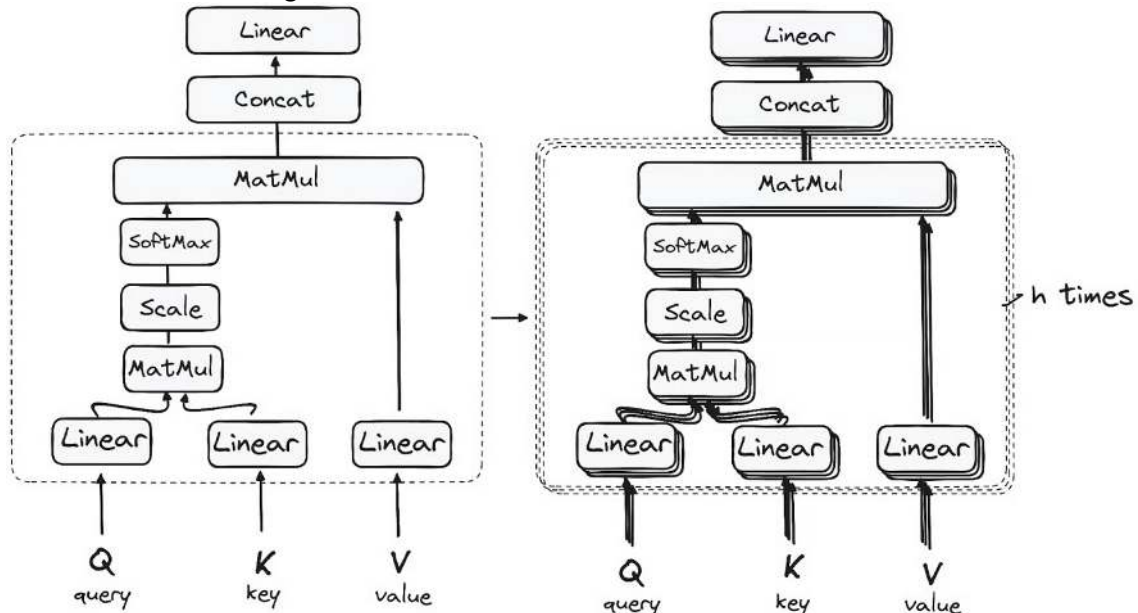
In the encoder, the multi-headed attention utilizes a specialized attention mechanism known as self-attention. This approach enables the models to relate each word in the input with other words. For instance, in a given example, the model might learn to connect the word “are” with “you”.

This mechanism allows the encoder to focus on different parts of the input sequence as it processes each token. It computes attention scores based on:

- A query is a vector that represents a specific word or token from the input sequence in the attention mechanism.
- A key is also a vector in the attention mechanism, corresponding to each word or token in the input sequence.
- Each value is associated with a key and is used to construct the output of the attention layer. When a query and a key match well, which basically means that they have a high attention score, the corresponding value is emphasized in the output.

This first Self-Attention module enables the model to capture contextual information from the entire sequence. Instead of performing a single attention function, queries, keys and values are linearly projected  $h$  times. On each of these projected versions of queries, keys and values the attention mechanism is performed in parallel, yielding  $h$ -dimensional output values.

The detailed architecture goes as follows:





### Matrix Multiplication (MatMul) - Dot Product of Query and Key

Once the query, key, and value vectors are passed through a linear layer, a dot product matrix multiplication is performed between the queries and keys, resulting in the creation of a score matrix. The score matrix establishes the degree of emphasis each word should place on other words. Therefore, each word is assigned a score in relation to other words within the same time step. A higher score indicates greater focus.

This process effectively maps the queries to their corresponding keys.

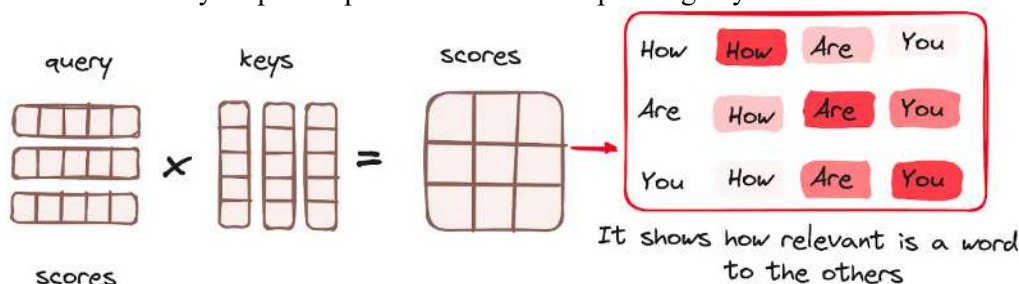


Fig: Encoder's workflow. Attention mechanism - Matrix Multiplication.

### Reducing the Magnitude of attention scores

The scores are then scaled down by dividing them by the square root of the dimension of the query and key vectors. This step is implemented to ensure more stable gradients, as the multiplication of values can lead to excessively large effects.

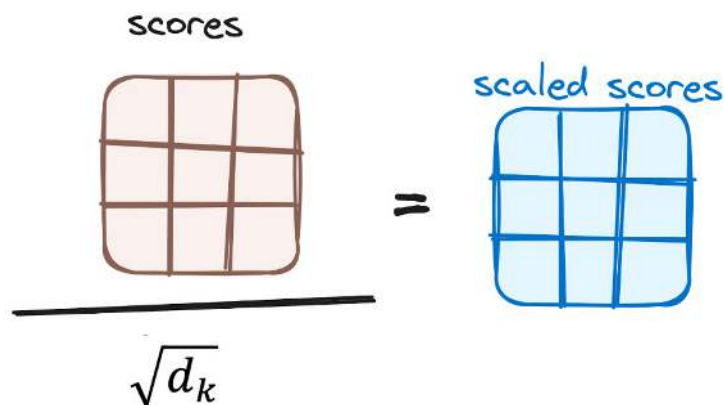


Fig: Encoder's workflow. Reducing the attention scores.

### Applying Softmax to the Adjusted Scores

Subsequently, a softmax function is applied to the adjusted scores to obtain the attention weights. This results in probability values ranging from 0 to 1. The softmax function emphasizes higher scores while diminishing lower scores, thereby enhancing the model's ability to effectively determine which words should receive more attention.

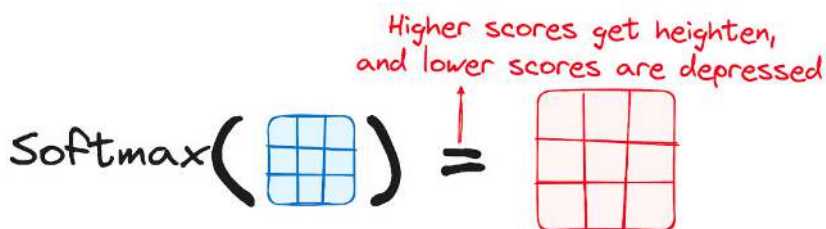


Fig: Encoder's workflow. Softmax adjusted scores.

### Combining Softmax Results with the Value Vector

The following step of the attention mechanism is that weights derived from the softmax function are multiplied by the value vector, resulting in an output vector.

In this process, only the words that present high softmax scores are preserved. Finally, this output vector is fed into a linear layer for further processing.



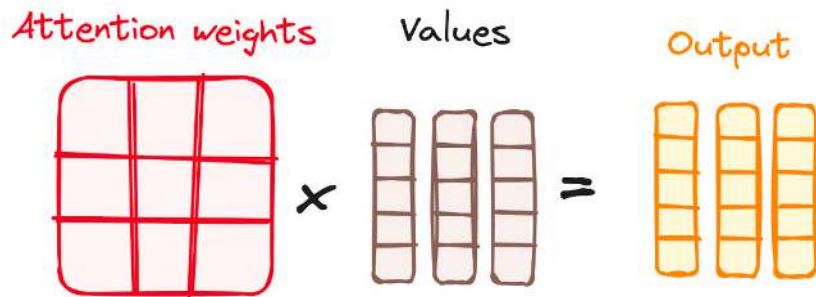


Fig: Encoder's workflow. Combining Softmax results with the value vector.

And we finally get the output of the Attention mechanism!

So, you might be wondering why it's called Multi-Head Attention?

Remember that before all the process starts, we break our queries, keys and values  $h$  times. This process, known as self-attention, happens separately in each of these smaller stages or 'heads'. Each head works its magic independently, conjuring up an output vector.

This ensemble passes through a final linear layer, much like a filter that fine-tunes their collective performance. The beauty here lies in the diversity of learning across each head, enriching the encoder model with a robust and multifaceted understanding.

### STEP 3.2 Normalization and Residual Connections

Each sub-layer in an encoder layer is followed by a normalization step. Also, each sub-layer output is added to its input (residual connection) to help mitigate the vanishing gradient problem, allowing deeper models. This process will be repeated after the Feed-Forward Neural Network too.

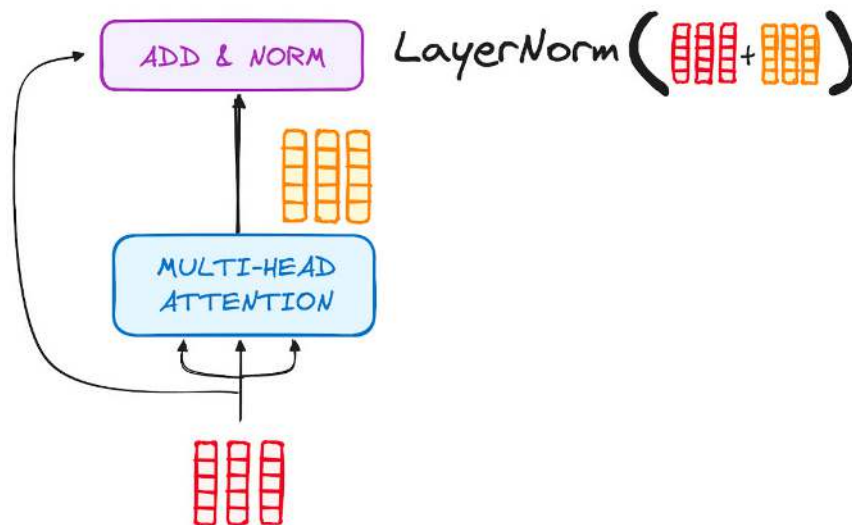


Fig: Encoder's workflow. Normalization and residual connection after Multi-Head Attention.

### STEP 3.3 Feed-Forward Neural Network

The journey of the normalized residual output continues as it navigates through a pointwise feed-forward network, a crucial phase for additional refinement.

Picture this network as a duo of linear layers, with a ReLU activation nestled in between them, acting as a bridge. Once processed, the output embarks on a familiar path: it loops back and merges with the input of the pointwise feed-forward network.

This reunion is followed by another round of normalization, ensuring everything is well-adjusted and in sync for the next steps.

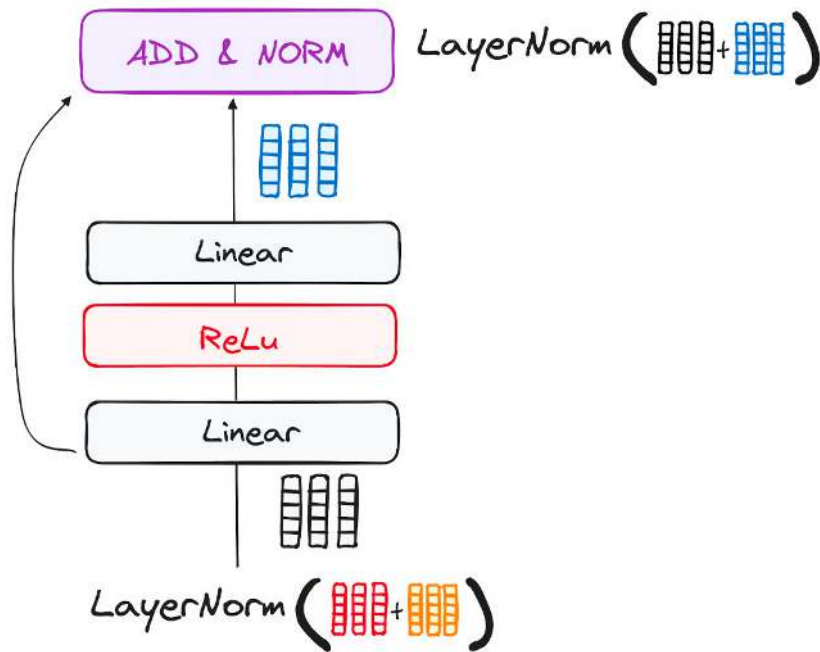


Fig: Encoder's workflow. Feed-Forward Neural Network sub-layer.

#### STEP 4 - Output of the Encoder

The output of the final encoder layer is a set of vectors, each representing the input sequence with a rich contextual understanding. This output is then used as the input for the decoder in a Transformer model.

This careful encoding paves the way for the decoder, guiding it to pay attention to the right words in the input when it's time to decode.

Think of it like building a tower, where you can stack up N encoder layers. Each layer in this stack gets a chance to explore and learn different facets of attention, much like layers of knowledge. This not only diversifies the understanding but could significantly amplify the predictive capabilities of the transformer network.

#### The Decoder WorkFlow

The decoder's role centers on crafting text sequences. Mirroring the encoder, the decoder is equipped with a similar set of sub-layers. It boasts two multi-headed attention layers, a pointwise feed-forward layer, and incorporates both residual connections and layer normalization after each sub-layer.

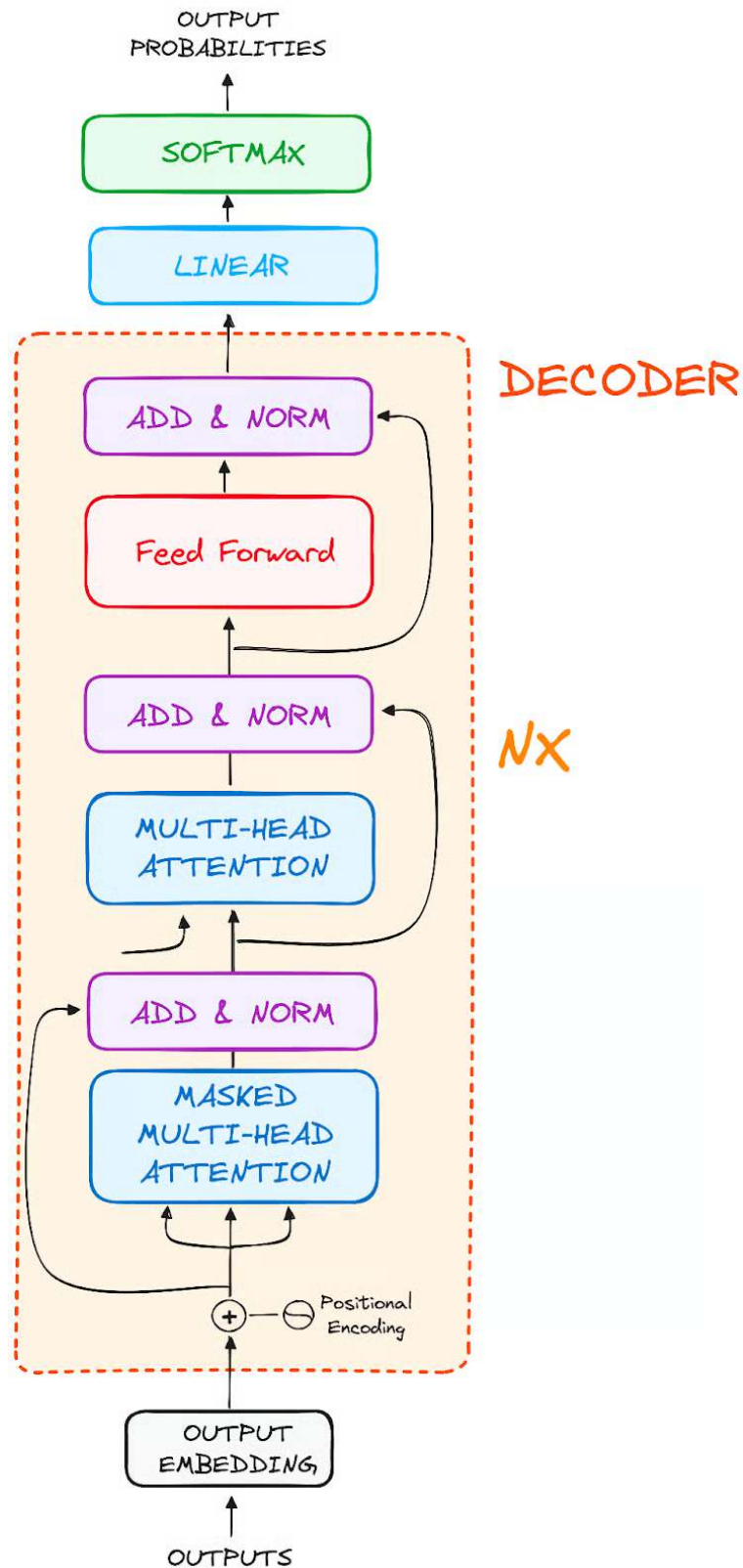


Fig: Global structure of Encoders.

These components function in a way akin to the encoder's layers, yet with a twist: each multi-headed attention layer in the decoder has its unique mission.

The final of the decoder's process involves a linear layer, serving as a classifier, topped off with a softmax function to calculate the probabilities of different words.

The Transformer decoder has a structure specifically designed to generate this output by decoding the encoded information step by step.

It is important to notice that the decoder operates in an autoregressive manner, kickstarting its process with a start token. It cleverly uses a list of previously generated outputs as its inputs, in tandem with the outputs from the encoder that are rich with attention information from the initial input.

This sequential dance of decoding continues until the decoder reaches a pivotal moment: the generation of a token that signals the end of its output creation.

### STEP 1 - Output Embeddings

At the decoder's starting line, the process mirrors that of the encoder. Here, the input first passes through an embedding layer

### STEP 2 - Positional Encoding

Following the embedding, again just like the encoder, the input passes by the positional encoding layer. This sequence is designed to produce positional embeddings.

These positional embeddings are then channeled into the first multi-head attention layer of the decoder, where the attention scores specific to the decoder's input are meticulously computed.

### STEP 3 - Stack of Decoder Layers

The decoder consists of a stack of identical layers (6 in the original Transformer model). Each layer has three main sub-components:

#### STEP 3.1 Masked Self-Attention Mechanism

This is similar to the self-attention mechanism in the encoder but with a crucial difference: it prevents positions from attending to subsequent positions, which means that each word in the sequence isn't influenced by future tokens.

For instance, when the attention scores for the word "are" are being computed, it's important that "are" doesn't get a peek at "you", which is a subsequent word in the sequence.

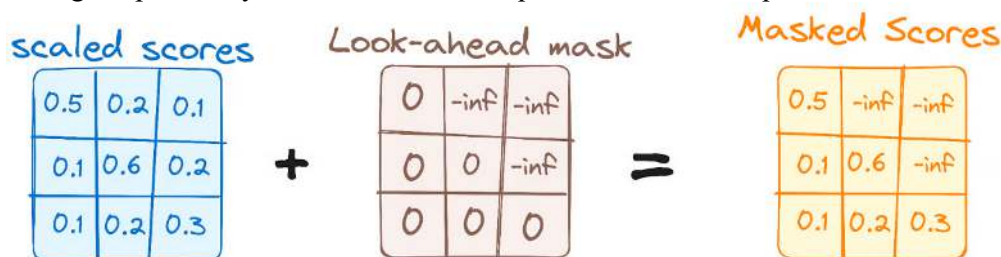


Fig: Decoder's workflow. First Multi-Head Attention Mask.

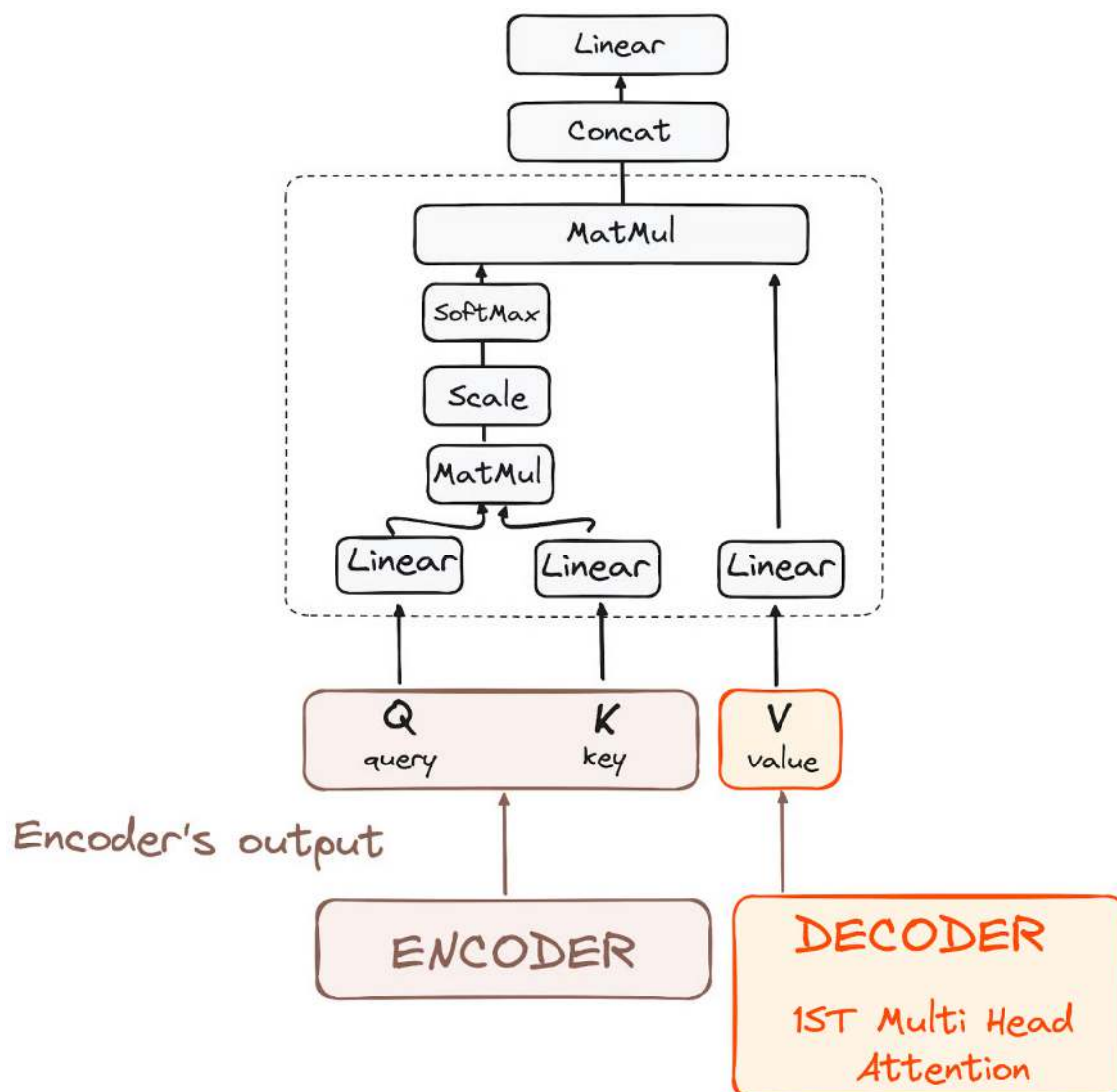
This masking ensures that the predictions for a particular position can only depend on known outputs at positions before it.

#### STEP 3.2 - Encoder-Decoder Multi-Head Attention or Cross Attention

In the second multi-headed attention layer of the decoder, we see a unique interplay between the encoder and decoder's components. Here, the outputs from the encoder take on the roles of both queries and keys, while the outputs from the first multi-headed attention layer of the decoder serve as values.

This setup effectively aligns the encoder's input with the decoder's, empowering the decoder to identify and emphasize the most relevant parts of the encoder's input.

Following this, the output from this second layer of multi-headed attention is then refined through a pointwise feedforward layer, enhancing the processing further.



*Fig: Decoder's workflow. Encoder-Decoder Attention.*

In this sub-layer, the queries come from the previous decoder layer, and the keys and values come from the output of the encoder. This allows every position in the decoder to attend over all positions in the input sequence, effectively integrating information from the encoder with the information in the decoder.

### STEP 3.3 Feed-Forward Neural Network

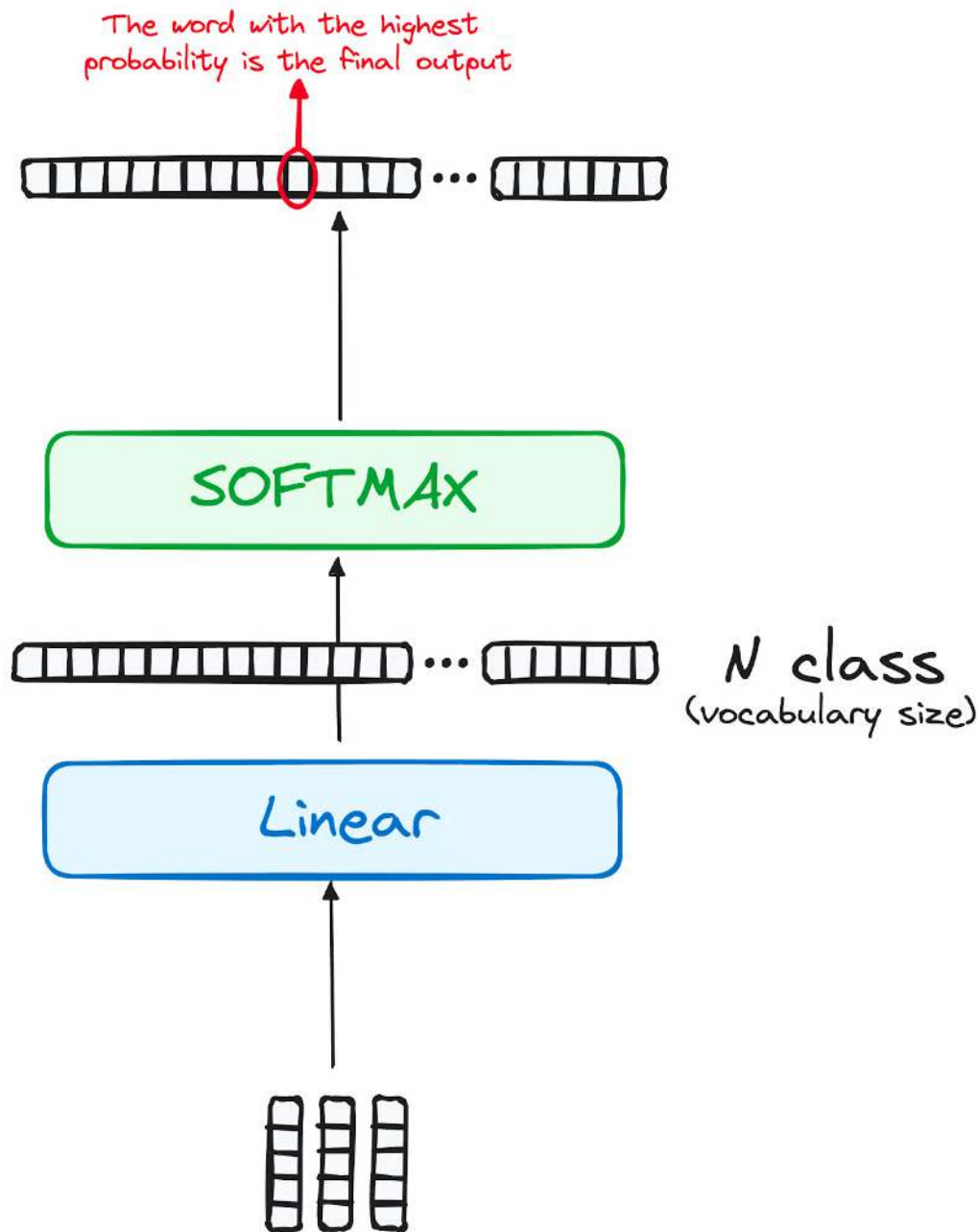
Similar to the encoder, each decoder layer includes a fully connected feed-forward network, applied to each position separately and identically.

### STEP 4 Linear Classifier and Softmax for Generating Output Probabilities

The journey of data through the transformer model culminates in its passage through a final linear layer, which functions as a classifier.

The size of this classifier corresponds to the total number of classes involved (number of words contained in the vocabulary). For instance, in a scenario with 1000 distinct classes representing 1000 different words, the classifier's output will be an array with 1000 elements.

This output is then introduced to a softmax layer, which transforms it into a range of probability scores, each lying between 0 and 1. The highest of these probability scores is key, its corresponding index directly points to the word that the model predicts as the next in the sequence.



*Fig: Decoder's workflow. Transformer's final output.*

### Normalization and Residual Connections

Each sub-layer (masked self-attention, encoder-decoder attention, feed-forward network) is followed by a normalization step, and each also includes a residual connection around it.

### Output of the Decoder

The final layer's output is transformed into a predicted sequence, typically through a linear layer followed by a softmax to generate probabilities over the vocabulary.

The decoder, in its operational flow, incorporates the freshly generated output into its growing list of inputs, and then proceeds with the decoding process. This cycle repeats until the model predicts a specific token, signaling completion.

The token predicted with the highest probability is assigned as the concluding class, often represented by the end token.

Again remember that the decoder isn't limited to a single layer. It can be structured with  $N$  layers, each one building upon the input received from the encoder and its preceding layers. This layered



architecture allows the model to diversify its focus and extract varying attention patterns across its attention heads.

Such a multi-layered approach can significantly enhance the model's ability to predict, as it develops a more nuanced understanding of different attention combinations.

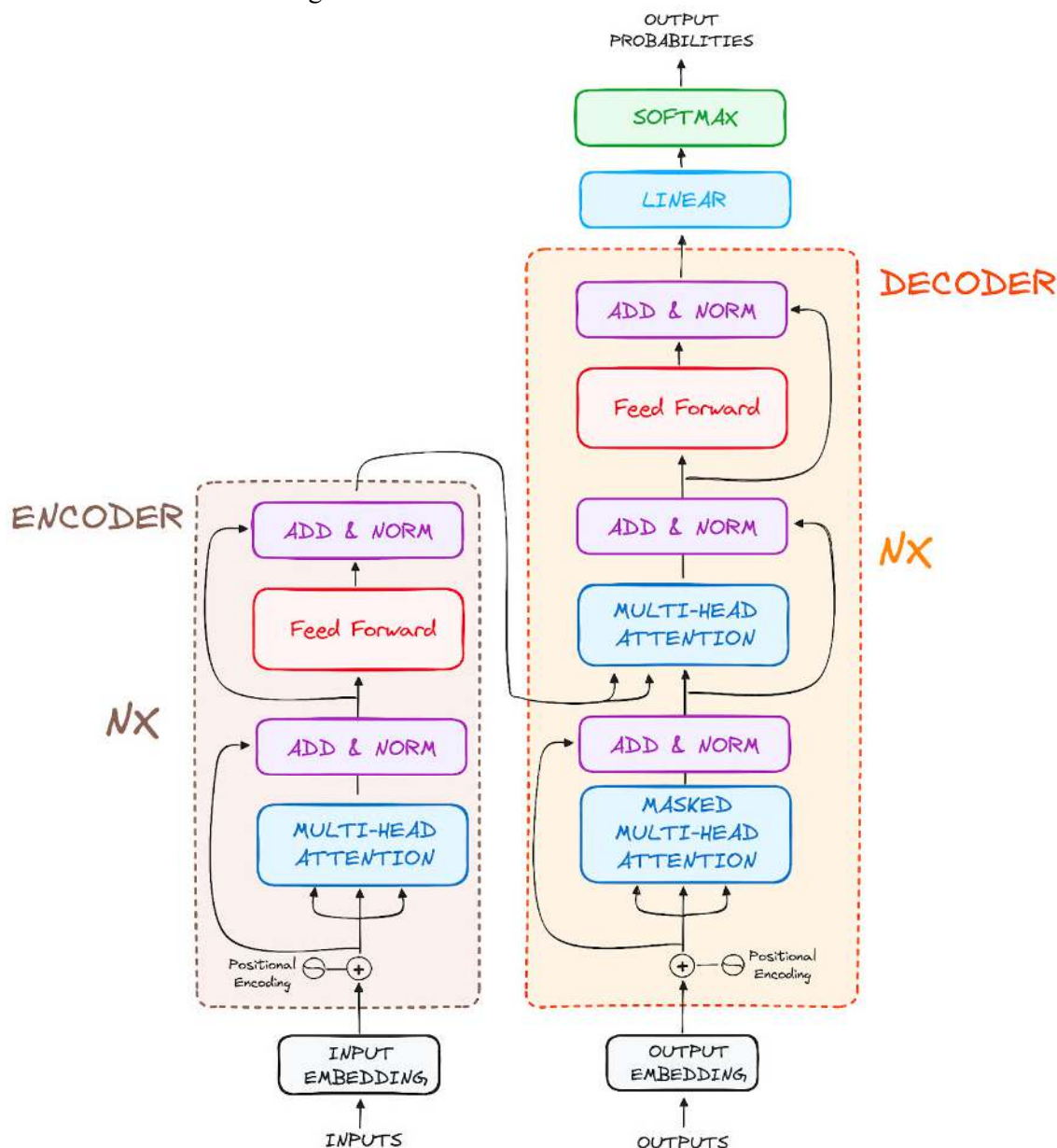


Fig: Original structure of Transformers.

## Real-Life Transformer Models

### BERT

Google's 2018 release of BERT, an open-source natural language processing framework, revolutionized NLP with its unique bidirectional training, which enables the model to have more context-informed predictions about what the next word should be.

By understanding context from all sides of a word, BERT outperformed previous models in tasks like question-answering and understanding ambiguous language. Its core uses Transformers, connecting each output and input element dynamically.

BERT, pre-trained on Wikipedia, excelled in various NLP tasks, prompting Google to integrate it into its search engine for more natural queries. This innovation sparked a race to develop advanced



language models and significantly advanced the field's ability to handle complex language understanding.

## LaMDA

LaMDA (Language Model for Dialogue Applications) is a Transformer-based model developed by Google, designed specifically for conversational tasks, and launched during the 2021 Google I/O keynote. They are designed to generate more natural and contextually relevant responses, enhancing user interactions in various applications.

LaMDA's design enables it to understand and respond to a wide range of topics and user intents, making it ideal for applications in chatbots, virtual assistants, and other interactive AI systems where a dynamic conversation is key.

This focus on conversational understanding and response marks LaMDA as a significant advancement in the field of natural language processing and AI-driven communication.

## GPT and ChatGPT

GPT and ChatGPT, developed by OpenAI, are advanced generative models known for their ability to produce coherent and contextually relevant text. GPT-1 was its first model launched in June 2018 and GPT-3, one of the most impactful models, was launched two years later in 2020.

These models are adept at a wide range of tasks, including content creation, conversation, language translation, and more. GPT's architecture enables it to generate text that closely resembles human writing, making it useful in applications like creative writing, customer support, and even coding assistance. ChatGPT, a variant optimized for conversational contexts, excels in generating human-like dialogue, enhancing its application in chatbots and virtual assistants.

## Other Variations

The landscape of foundation models, particularly transformer models, is rapidly expanding. A study identified over 50 significant transformer models, while **the Stanford group evaluated 30 of them**, acknowledging the field's fast-paced growth. NLP Cloud, an innovative startup part of NVIDIA's Inception program, utilizes around 25 large language models commercially for various sectors like airlines and pharmacies.

There is an increasing trend towards making these models open-source, with platforms like Hugging Face's model hub leading the way. Additionally, numerous Transformer-based models have been developed, each specialized for different NLP tasks, showcasing the model's versatility and efficiency in diverse applications.

You can learn more from all existing **Foundation Models** in a separate article which talks about what they are and which ones are the most used.

## Pre-trained models:

A pre-trained model, having been trained on extensive data, serves as a foundational model for various tasks, leveraging its learned patterns and features. In natural language processing (NLP), these models are commonly employed as a starting point for tasks like language translation, sentiment analysis, and text summarization.

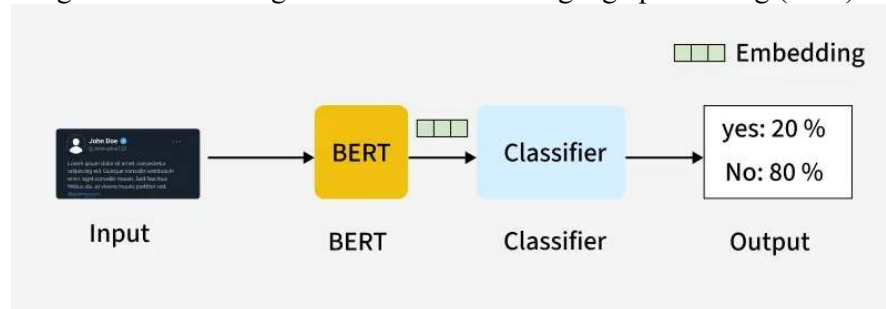
### Why do we use Pretrained Models?

Pretrained models are effective and efficient solutions for developers, researchers and businesses as they eliminate the need to write and train the code from scratch, saving time. Other advantages of using pre-trained models for projects are:

- Reduces the computational burden required for initial model training hence, making development more accessible.
- The learned knowledge can be used for various applications.
- Models can be fine-tuned according to the task and can result in superior performance to training from the initial point.
- Less labelled data is required for fine-tuning specific tasks.

## BERT:

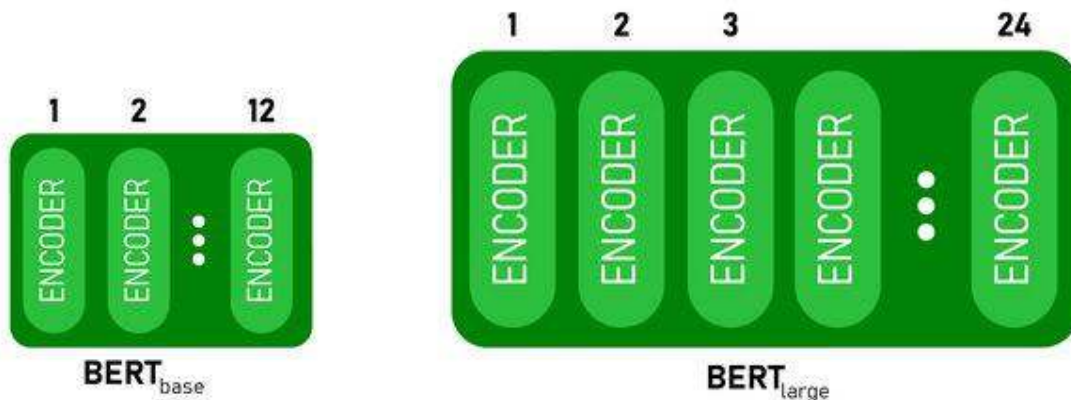
BERT (Bidirectional Encoder Representations from Transformers) stands as an open-source machine learning framework designed for the natural language processing (NLP).



## BERT Architecture

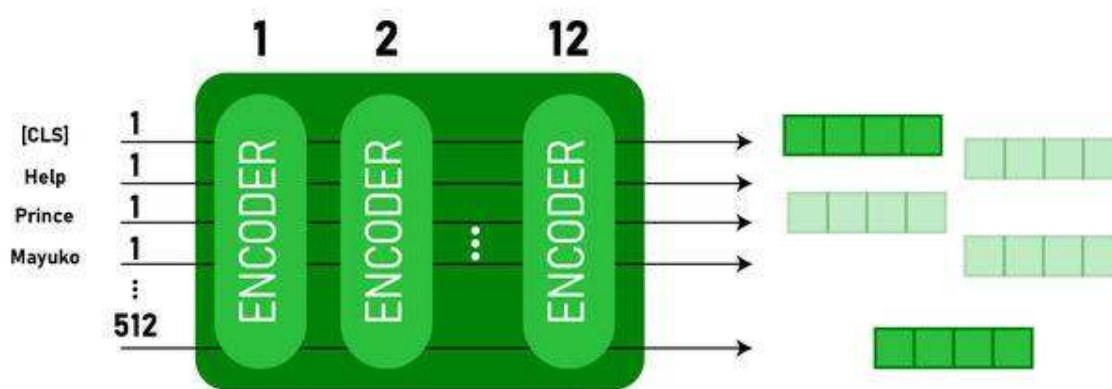
The architecture of BERT is a multilayer bidirectional transformer encoder which is quite similar to the transformer model. A transformer architecture is an encoder-decoder network that uses self-attention on the encoder side and attention on the decoder side.

1. BERTBASE has 12 layers in the Encoder stack while BERTLARGE has 24 layers in the Encoder stack. These are more than the Transformer architecture described in the original paper (6 encoder layers).
2. BERT architectures (BASE and LARGE) also have larger feedforward networks (768 and 1024 hidden units respectively), and more attention heads (12 and 16 respectively) than the Transformer architecture suggested in the original paper. It contains 512 hidden units and 8 attention heads.
3. BERTBASE contains 110M parameters while BERTLARGE has 340M parameters.



## BERT BASE and BERT LARGE architecture

This model takes the CLS token as input first, then it is followed by a sequence of words as input. Here CLS is a classification token. It then passes the input to the above layers. Each layer applies self-attention and passes the result through a feedforward network after then it hands off to the next encoder. The model outputs a vector of hidden size (768 for BERT BASE). If we want to output a classifier from this model we can take the output corresponding to the CLS token.



### BERT output as Embeddings

Now, this trained vector can be used to perform a number of tasks such as classification, translation, etc. For Example, the paper achieves great results just by using a single layer Neural Network on the BERT model in the classification task.

### How to use BERT model in NLP?

BERT can be used for various natural language processing (NLP) tasks such as:

#### 1. Classification Task

- BERT can be used for classification task like sentiment analysis, the goal is to classify the text into different categories (positive/ negative/ neutral), BERT can be employed by adding a classification layer on the top of the Transformer output for the [CLS] token.
- The [CLS] token represents the aggregated information from the entire input sequence. This pooled representation can then be used as input for a classification layer to make predictions for the specific task.

#### 2. Question Answering

- In question answering tasks, where the model is required to locate and mark the answer within a given text sequence, BERT can be trained for this purpose.
- BERT is trained for question answering by learning two additional vectors that mark the beginning and end of the answer. During training, the model is provided with questions and corresponding passages, and it learns to predict the start and end positions of the answer within the passage.

#### 3. Named Entity Recognition (NER)

- BERT can be utilized for NER, where the goal is to identify and classify entities (e.g., Person, Organization, Date) in a text sequence.
- A BERT-based NER model is trained by taking the output vector of each token from the Transformer and feeding it into a classification layer. The layer predicts the named entity label for each token, indicating the type of entity it represents.

### Fine-Tuning on Labeled Data

We perform Fine-tuning on labeled data for specific NLP tasks.

- After the pre-training phase, the BERT model, armed with its contextual embeddings, is fine-tuned for specific natural language processing (NLP) tasks. This step tailors the model to more targeted applications by adapting its general language understanding to the nuances of the particular task.
- BERT is fine-tuned using labeled data specific to the downstream tasks of interest. These tasks could include sentiment analysis, question-answering, named entity recognition, or any other NLP application. The model's parameters are adjusted to optimize its performance for the particular requirements of the task at hand.

BERT's unified architecture allows it to adapt to various downstream tasks with minimal modifications, making it a versatile and highly effective tool in natural language understanding and processing.

## How to Tokenize and Encode Text using BERT?

To tokenize and encode text using BERT, we will be using the 'transformer' library in Python.

### Command to install transformers:

*pip install transformers*

- We will load the pretrained BERT tokenizer with a cased vocabulary using `BertTokenizer.from_pretrained("bert-base-cased")`.
  - `tokenizer.encode(text)` tokenizes the input text and converts it into a sequence of token IDs.
  - `print("Token IDs:", encoding)` prints the token IDs obtained after encoding.
  - `tokenizer.convert_ids_to_tokens(encoding)` converts the token IDs back to their corresponding tokens.
  - `print("Tokens:", tokens)` prints the tokens obtained after converting the token IDs
- from transformers import BertTokenizer

```
tokenizer = BertTokenizer.from_pretrained("bert-base-cased")
```

```
text = 'ChatGPT is a language model developed by OpenAI, based on the GPT (Generative Pre-trained Transformer) architecture. '
```

```
# Tokenize and encode the text
```

```
encoding = tokenizer.encode(text)
```

```
print("Token IDs:", encoding)
```

```
# Convert token IDs back to tokens
```

```
tokens = tokenizer.convert_ids_to_tokens(encoding)
```

```
print("Tokens:", tokens)
```

### Output

```
Token IDs: [101, 24705, 1204, 17095, 1942, 1110, 170, 1846, 2235, 1872, 1118, 3353, 1592, 2240, 117, 1359, 1113, 1103, 15175, 1942, 113, 9066, 15306, 11689, 118, 3972, 13809, 23763, 114, 4220, 119, 102]
```

```
Tokens: ['[CLS]', 'Cha', '##t', '##GP', '##T', 'is', 'a', 'language', 'model', 'developed', 'by', 'Open', '##A', '##I', ',', 'based', 'on', 'the', 'GP', '##T', '(', 'Gene', '##rative', 'Pre', '-', 'trained', 'Trans', '##former', ')', 'architecture', '.', '[SEP]']
```

The `tokenizer.encode` method adds the special [CLS] - classification and [SEP] - separator tokens at the beginning and end of the encoded sequence. In the token IDs section, token id: 101 refers to the start of the sentence and token id: 102 represents the end of the sentence.

### Application of BERT

BERT is used for various applications. Some of these are:

1. **Text Representation:** BERT is used to generate word embeddings or representation for words in a sentence.
2. **Named Entity Recognition (NER):** BERT can be fine-tuned for named entity recognition tasks, where the goal is to identify entities such as names of people, organizations, locations, etc., in a given text.
3. **Text Classification:** BERT is widely used for text classification tasks, including sentiment analysis, spam detection, and topic categorization. It has demonstrated excellent performance in understanding and classifying the context of textual data.
4. **Question-Answering Systems:** BERT has been applied to question-answering systems, where the model is trained to understand the context of a question and provide relevant answers. This is particularly useful for tasks like reading comprehension.
5. **Machine Translation:** BERT's contextual embeddings can be leveraged for improving machine translation systems. The model captures the nuances of language that are crucial for accurate translation.

6. **Text Summarization:** BERT can be used for abstractive text summarization, where the model generates concise and meaningful summaries of longer texts by understanding the context and semantics.
7. **Conversational AI:** BERT is employed in building conversational AI systems, such as chatbots, virtual assistants, and dialogue systems. Its ability to grasp context makes it effective for understanding and generating natural language responses.
8. **Semantic Similarity:** BERT embeddings can be used to measure semantic similarity between sentences or documents. This is valuable in tasks like duplicate detection, paraphrase identification, and information retrieval.

### BERT vs GPT

The difference between BERT and GPT are as follows:

	BERT	GPT
<b>Architecture</b>	Bidirectional; predicts masked words based on left, right context.	Unidirectional; predicts next word given preceding context.
<b>Pre-training Objectives</b>	BERT is pre-trained using a masked language model objective and next sentence prediction.	GPT is pre-trained using Next word prediction only.
<b>Context Understanding</b>	Strong at understanding and analyzing text.	Strong in generating coherent and contextually relevant text.
<b>Tasks and Use Cases</b>	Commonly used in tasks like text classification, NER, sentiment analysis, and QA	Applied to tasks like text generation, chat, summarization, etc.
<b>Fine-tuning vs Few-Shot Learning</b>	Fine-tuning with labeled data to adapt its pre-trained representations to the task at hand.	GPT is designed to perform few-shot or zero-shot learning, where it can generalize with minimal task-specific data.

### Question Answering using Transformers:

Question Answering (QA) in Natural Language Processing (NLP) is a task that involves using computational methods to automatically answer questions posed in natural language. This task involves identifying the relevant information within a large corpus of text and extracting the answer that best addresses the question. It involves a combination of various NLP techniques, such as information retrieval, information extraction, and text understanding, to accurately answer questions based on the context and meaning of the input text. QA systems are widely used in applications such as customer service, knowledge management, and education, providing quick and accurate answers to a wide range of questions.

#### How Question Answering works in NLP:

Question Answering (QA) in NLP involves identifying the answer to a question from a large corpus of text, such as a document or a collection of documents. The process typically involves the following steps:

1. **Question Analysis:** The question is analyzed to identify its type (e.g., factual, definition, comparison, etc.) and the information required to find the answer (e.g., named entities, keywords, etc.).
2. **Information Retrieval:** The corpus is searched to retrieve the relevant documents or text snippets that contain the information needed to answer the question.
3. **Candidate Answer Generation:** Possible answers are extracted from the retrieved text snippets and candidate answers are generated.
4. **Answer Selection:** The candidate answers are evaluated against the question and the best answer is selected based on various criteria, such as relevance, accuracy, and confidence.

The process can be aided by techniques such as natural language processing (NLP), information retrieval (IR), and machine learning (ML), among others. The performance of QA systems can be improved by using large pre-trained language models, such as BERT or GPT, to encode the context of the question and candidate answers.

### Applications of the Question Answering

Question Answering (QA) systems have a wide range of applications, including:

1. **Customer Service:** QA systems can be used to automate customer service by providing instant answers to common questions.
2. **Knowledge Management:** QA systems can be used to manage and retrieve information from large knowledge bases, such as company FAQs and product manuals.
3. **Education:** QA systems can be used to provide students with instant answers to their questions and help them learn more effectively.
4. **Healthcare:** QA systems can be used to provide healthcare professionals with instant answers to clinical questions, improving patient outcomes and reducing diagnostic errors.
5. **News and Media:** QA systems can be used to provide quick answers to factual questions about current events, politics, sports, and more.
6. **Legal:** QA systems can be used to provide lawyers and legal researchers with instant answers to legal questions, improving the speed and efficiency of their work.

QA systems aim to improve decision-making, increase productivity, and provide fast and accurate answers to questions, making it a valuable tool in a wide range of industries and domains.

### Advantages of using Transformers for Question Answering

Transformers are widely used for Question Answering (QA) for several reasons:

1. **Contextual Embeddings:** Transformers can capture the context of a question and candidate answer, allowing them to better understand the relationship between the two and make more accurate predictions.
2. **Pre-training:** Transformers can be pre-trained on large amounts of text data, giving them a strong understanding of language and enabling them to perform well on a wide range of tasks, including QA.
3. **Attention Mechanisms:** Transformers use attention mechanisms, which allow them to focus on specific parts of the input when making predictions, leading to improved performance on QA tasks.
4. **Transfer Learning:** Transformers can be fine-tuned for specific QA tasks, allowing for transfer learning from a pre-trained model and reducing the amount of labeled data needed for training.
5. **Performance:** Transformers have demonstrated state-of-the-art performance on a wide range of QA tasks, outperforming traditional methods in terms of accuracy and speed.

The use of transformers for QA allows for more accurate, fast, and scalable solutions compared to traditional methods, making it a popular choice in many applications.

### Example:

```
from transformers import pipeline
```

```
# create pipeline for QA  
qa = pipeline('question-answering')
```

```
ctx = "My name is Ganesh and I am studying Data Science"
que = "What is Ganesh studying?"
qa(context = ctx, question = que)
```

```
[5] ctx = "My name is Ganesh and I am studying Data Science"
    que = "What is Ganesh studying?"
    qa(context = ctx, question = que)
```

```
{'score': 0.9988707900047302, 'start': 36, 'end': 48, 'answer': 'Data Science'}
```

```
question = "What is the capital city of Japan?"
qa(context = context, question = question)
```

```
{'score': 0.3223932087421417, 'start': 725, 'end': 747, 'answer': 'Pretoria, South Africa'}
```

## Text Summarization:

Text summarization is the process of condensing a large text document into a shorter version while preserving its key information and meaning. The goal of text summarization is to extract the most important information from a text document and present it in a concise and comprehensible form. This can be done through techniques such as keyword extraction, sentence extraction, or abstractive summarization. Text summarization has various applications including news aggregation, content analysis, and information retrieval.

### Text summarization performed using Transformers:

Text summarization using Transformers can be performed in two ways: extractive summarization and abstractive summarization.

1. **Extractive summarization:** In this approach, the most important sentences or phrases from the original text are selected and combined to form a summary. This can be done using algorithms such as TextRank, which uses graph-based algorithms to rank sentences based on their relevance and importance. Transformers can be used to process the text, extract features, and perform sentence ranking.
2. **Abstractive summarization:** In this approach, a new summary is generated by understanding the context of the original text and generating new phrases and sentences that summarize its content. This can be done using techniques such as encoder-decoder models, where the encoder processes the input text to extract its features and the decoder generates the summary. Transformers can be used as the encoder or decoder in this architecture.

In both extractive and abstractive summarization, Transformers can be trained on large amounts of text data to learn the patterns and relationships between words, sentences, and documents, making them well-suited for text summarization tasks.

Transformers are used for text summarization because they are highly effective in processing and understanding large amounts of text data. There are several reasons why Transformers are a popular choice for text summarization:

1. **Contextual understanding:** Transformers use attention mechanisms to understand the context of words, sentences, and documents, which is crucial for text summarization tasks. This allows Transformers to accurately identify the most important information in a text document.
2. **Large language model:** Transformers have been trained on vast amounts of text data, which enables them to have a deep understanding of language patterns and relationships. This makes them well-suited for text summarization tasks where a comprehensive understanding of language is required.
3. **Scalability:** Transformers can process large amounts of text data in parallel, making them well-suited for summarizing long documents or large volumes of text data.
4. **End-to-end training:** Transformers can be trained end-to-end on text summarization tasks, which allows them to optimize their performance for the specific task at hand.



5. **State-of-the-art results:** Transformers have achieved state-of-the-art results on a wide range of natural language processing tasks, including text summarization.

By using Transformers for text summarization, organizations can benefit from their ability to understand and process large amounts of text data, leading to more accurate and reliable summarization results.

### Use of Text summarization

Text summarization has various uses, including:

1. **News aggregation:** Summarizing news articles to provide a quick overview of the most important information.
2. **Content analysis:** Reducing large volumes of text data to identify patterns, trends, and insights.
3. **Information retrieval:** Summarizing search results to provide users with a concise and easy-to-understand overview of relevant information.
4. **Document management:** Summarizing long documents to make them easier to manage and search through.
5. **Meeting minutes:** Summarizing the key points discussed in a meeting to provide a concise and organized record.
6. **Customer feedback:** Summarizing customer feedback to identify common themes and issues.
7. **Legal contracts:** Summarizing legal contracts to provide a quick overview of their key terms and conditions.
8. **Customer service:** Summarizing customer support requests to quickly identify the root cause of an issue.

By condensing large amounts of text into a more manageable form, text summarization can help users quickly understand the content of a document and make more informed decisions.

### Example:

```
# import required libraries
from transformers import pipeline
import textwrap
import numpy as np
import pandas as pd
from pprint import pprint
```

```
df = pd.read_csv('bbc_text_cls.csv?dl=0')
```

```
df.head()
```

	text	labels
0	Ad sales boost Time Warner profit\n\nQuarterly...	business
1	Dollar gains on Greenspan speech\n\nThe dollar...	business
2	Yukos unit buyer faces loan claim\n\nThe owner...	business
3	High fuel prices hit BA's profits\n\nBritish A...	business
4	Pernod takeover talk lifts Domecq\n\nShares in...	business

```
doc = df[df.labels == 'business']['text'].sample(random_state=42)
```

```
# text wrapping function
```

```
def wrap(x):
```

```
    return textwrap.fill(x, replace_whitespace = False, fix_sentence_endings = True)
```

```
print(wrap(doc.iloc[0]))
```

```

# Create a pipeline for summarization
summarizer = Pipeline([
    ('tokenizer', TfidfTokenizer(analyzer='word', token_pattern=r'\b(?:[a-z]+|'[a-z]{4,})\b')),
    ('vectorizer', TfidfVectorizer(tokenizer=tokenizer.tokenizer_, max_df=0.9, min_df=2, max_features=10000,
                                  stop_words='english')),
    ('model', Summarizer()))

# Apply the pipeline to the document
summary_text = summarizer.transform([doc]).toarray()[0]

# Print the summary text
print(summary_text)

```

```
summarizer = pipeline('summarization')
```

```
summarizer(doc.iloc[0].split('\n',1)[1])
```

```
5] summarizer(doc.iloc[0].split('\n',1)[1])
```

```

[{'summary_text': 'Retail sales dropped by 1% on the month in December, after a 0.6% rise in November
. Clothing retailers and non-specialist stores were the worst hit with only internet retailers showing
any significant growth . The last time retailers endured a tougher Christmas was 23 years ago, when
sales plunged 1.7% .'}]

```

summarized text

## Sentiment Analysis

Sentiment analysis is the process of determining the emotion or opinion expressed in a piece of text, such as a tweet or a review. It is a subfield of natural language processing (NLP) that involves using machine learning and NLP techniques to classify the sentiment of a piece of text as positive, negative, or neutral. The goal of sentiment analysis is to extract insights from large amounts of unstructured text data and understand the overall sentiment of a population towards a particular topic.

### Sentiment Analysis as a classification task

- Binary classification (positive, negative)
- Multiclass: positive, negative or neutral
- Fine-grain: very positive, positive, neutral, negative, very negative
- Dependent on the dataset you're using or how you made your dataset
- If we use prebuild classifier (e.g. Hugging Face) we are stuck with whatever the model was trained to do (in this case binary classification)

### Use of Sentiment Analysis

1. **Improve customer experience:** Sentiment analysis can be used to understand customer opinions and feedback on products and services, allowing companies to improve the customer experience and build stronger customer relationships.
2. **Market research:** Sentiment analysis can be used to monitor public sentiment towards a particular topic or brand, providing valuable insights into market trends and customer preferences.
3. **Improve product design:** Sentiment analysis can be used to understand customer preferences and opinions on product features, allowing companies to design products that better meet customer needs.
4. **Improve social media monitoring:** Sentiment analysis can be used to monitor social media conversations and understand the tone and sentiment behind them, providing valuable insights into public perception and sentiment.
5. **Improve decision making:** Sentiment analysis can provide a broad overview of the sentiment of a population towards a particular topic or brand, allowing decision-makers to make informed decisions based on data-driven insights.

Traditionally, sentiment analysis was performed using shallow learning techniques such as Support Vector Machines (SVMs) and Naive Bayes Classifiers. However, these techniques had several limitations, such as the inability to capture context and the need for a large annotated dataset for training.

### Transformers for Sentiment Analysis

Enter Transformers, a breakthrough in NLP that has revolutionized the field of sentiment analysis. They are based on the Transformer architecture and use self-attention mechanisms to process the input sequence, allowing the model to capture the context and dependencies between words.

1. **Improve context representation:** Transformers use self-attention mechanisms to process the input sequence, allowing the model to capture the context and dependencies between words, leading to improved representation of the input text.
2. **Handling long sequences:** Transformers are able to handle long sequences of text, making them well-suited for sentiment analysis where the context and dependencies between words play an important role.
3. **Large number of parameters:** Transformers have a large number of parameters, allowing them to learn complex representations of the input text, resulting in improved performance compared to shallow learning methods.
4. **Pre-training:** Many Transformer-based models, such as BERT, are pre-trained on large annotated datasets, providing a strong base for fine-tuning on specific NLP tasks, including sentiment analysis.
5. **State-of-the-art performance:** Transformer-based models, such as BERT, have achieved state-of-the-art results on a variety of NLP tasks, including sentiment analysis, making them a popular choice for the task.

#### Example:

```
from transformers import pipeline
# create pipeline for sentiment analysis
classification = pipeline('sentiment-analysis')
```

```
#let's try on simple movie reviews
classification("I thoroughly enjoyed this movie!")
classification("I did not understand anything in this movie.")
```

```
[6] classification("I thoroughly enjoyed this movie!")

[{'label': 'POSITIVE', 'score': 0.9998749494552612}]
```

```
[7] classification("I did not understand anything in this movie.")

[{'label': 'NEGATIVE', 'score': 0.9994907379150391}]
```

```
# let's try on difficult movie reviews
classification("Although the movie had its flaws and some scenes were slow-paced, I still found myself captivated by the strong performances and unique storyline. The cinematography was stunning and added another level of depth to the film. Overall, I would still recommend it for those looking for a thought-provoking movie experience.")
```

```
classification("At first, I was intrigued by the premise and was excited to see where the story would go. However, as the movie progressed, I became disappointed by the lackluster execution and underdeveloped characters. Despite a few moments of promise, the movie ultimately fell flat for me and failed to deliver on its potential. Not recommended.")
```

```
[8] classification('Although the movie had its flaws and some scenes were slow-paced,
I still found myself captivated by the strong performances and unique storyline.
The cinematography was stunning and added another level of depth to the film. Overall,
I would still recommend it for those looking for a thought-provoking movie experience.')

[{'label': 'POSITIVE', 'score': 0.999847412109375}]

[9] classification('At first, I was intrigued by the premise and was excited to see where the story would go.
However, as the movie progressed, I became disappointed by the lackluster execution and
underdeveloped characters. Despite a few moments of promise, the movie ultimately fell
flat for me and failed to deliver on its potential. Not recommended.')

[{'label': 'NEGATIVE', 'score': 0.9997879862785339}]
```

*classification(["I thoroughly enjoyed this movie!",  
"I did not understand anything in this movie."])*

```
✓ [10] classification(["I thoroughly enjoyed this movie!",  
0s "I did not understand anything in this movie."])

[{'label': 'POSITIVE', 'score': 0.9998749494552612},  
{ 'label': 'NEGATIVE', 'score': 0.9994907379150391}]
```

#### Multiple results

### Zero-Shot Classification Using Transformers

Zero-shot classification is a paradigm of machine learning where a model can classify an input into one of multiple classes, even if it has never seen any examples from those classes during training. In other words, the model can generalize to new, unseen categories based on the knowledge it has acquired during training.

### Transformers and Zero-Shot Classification

Transformers, which rely on attention mechanisms and self-attention, have become the go-to architecture for NLP tasks. Their capacity to understand context and relationships in text makes them incredibly adept at zero-shot classification.

Here's how transformers achieve zero-shot classification:

1. **Pre-training:** Transformers are initially pre-trained on vast corpora of text data. During this stage, they learn the inherent structure, grammar, and semantics of language. This pre-training equips the model with a fundamental understanding of language, which is essential for generalization.
2. **Fine-tuning:** After pre-training, transformers are fine-tuned on specific tasks, such as sentiment analysis, text summarization, or named entity recognition. During this fine-tuning process, the model becomes specialized for the task at hand. While it becomes an expert in this particular domain, it retains its general language understanding capabilities.
3. **Zero-shot classification:** The remarkable feature of transformers is their ability to perform zero-shot classification. By providing a few examples from each class and a brief description of the task, the model can classify text into the specified classes, even if it has never seen those classes during training. This capability arises from the model's general language understanding and its fine-tuning on other related tasks.

### Applications of Zero-Shot Classification

Zero-shot classification using transformers has a wide range of applications across various domains:

1. **Text categorization:** It can be used for classifying news articles, user reviews, or social media posts into predefined categories, making it useful for content recommendation systems and content moderation.
2. **Sentiment analysis:** Transformers can classify text as positive, negative, or neutral without the need for specific training data for each sentiment category.
3. **Named entity recognition:** They can recognize entities in text, such as names of people, organizations, and locations, even if the model was not explicitly trained for these entities.
4. **Content tagging:** Zero-shot classification is beneficial in automatically tagging content with relevant keywords or labels.
5. **Language translation:** It can help in automatically determining the target language of text for translation, enabling more efficient multilingual translation services.

## Challenges and Limitations

While zero-shot classification using transformers is a powerful tool, it has its challenges and limitations. Some key considerations include:

1. **Limited fine-tuning data:** The quality and quantity of the fine-tuning data can significantly impact the model's performance on zero-shot tasks.
2. **Ambiguity:** Transformers may struggle with highly ambiguous text or complex multi-class categorizations.
3. **Biases:** The models may inherit biases present in the training data, potentially affecting their zero-shot classification performance.

## Code Example

# Python code for Zero-Shot Classification using Transformers

```
# Import necessary libraries
from transformers import pipeline

# Load the zero-shot classification pipeline
classifier = pipeline("zero-shot-classification")

# Define your input text and possible labels (classes)
input_text = "Astronomy is the study of stars and planets."
possible_labels = ["Science", "History", "Sports"]

# Perform zero-shot classification
result = classifier(input_text, possible_labels)

# Print the result
print("Input Text:", input_text)
print("Predicted Class:", result["labels"][0])
print("Confidence Score:", result["scores"][0])
```