**Attention Mechanism in Deep Learning**

One of the most groundbreaking advances in deep learning is the **Attention Mechanism**, a revolutionary concept that has transformed natural language processing (NLP). By allowing models to focus on the most relevant parts of input data, it paved the way for breakthroughs like the **Transformer architecture** and **Google's BERT**. Thanks to increased computing power, attention mechanisms have become essential in modern NLP, enabling smarter, more efficient AI systems. If you're working (or planning to work) in NLP, understanding how attention works is a must—it's one of the most important innovations in deep learning over the past decade. In this article you will get understanding about the attention mechanism in deep learning, how it works also attention mechanism in computer vision.

What is Attention Mechanism?

Attention mechanisms enhance deep learning models by selectively focusing on important input elements, improving prediction accuracy and computational efficiency. They prioritize and emphasize relevant information, acting as a spotlight to enhance overall model performance.

In psychology, attention is the cognitive process of selectively concentrating on one or a few things while ignoring others.

Let me explain what this means. Let's say you are seeing a group photo of your first school. Typically, there will be a group of children sitting across several rows, and the teacher will sit somewhere in between. Now, if anyone asks the question, "How many people are there?", how will you answer it?

Simply by counting heads, right? You don't need to consider any other things in the photo. Now, if anyone asks a different question, "Who is the teacher in the photo?", your brain knows exactly what to do. It will simply start looking for the features of an adult in the photo. The rest of the features will simply be ignored. **This is the 'Attention' which our brain is very adept at implementing.**

How Attention Mechanism Works?

Here's how they work:

1. **Breaking Down the Input**: Let's say you have a bunch of words (or any kind of data) that you want the computer to understand. First, it breaks down this input into smaller pieces, like individual words.

2. **Picking Out Important Bits**: Then, it looks at these pieces and decides which ones are the most important. It does this by comparing each piece to a question or 'query' it has in mind.

3. **Assigning Importance**: Each piece gets a score based on how well it matches the question. The higher the score, the more important that piece is.

4. **Focusing Attention**: After scoring each piece, it figures out how much attention to give to each one. Pieces with higher scores get more attention, while less important ones get less attention.

5. **Putting It All Together**: Finally, it adds up all the pieces, but gives more weight to the important ones. This way, the computer gets a clearer picture of what's most important in the input.

How Attention Mechanism was Introduced in Deep Learning

The attention mechanism emerged as an improvement over the encoder decoder-based **neural machine translation system** in natural language processing (NLP). Later, this mechanism, or its variants, was used in other applications, including **computer vision**, speech processing, etc.

Before **Bahdanau et al** proposed the first Attention model in 2015, neural machine translation was based on encoder-decoder **RNNs/LSTMs**. Both encoder and decoder are stacks of LSTM/RNN units. It works in the two following steps:

1. **The encoder LSTM is used to process the entire input sentence and encode it into a context vector**, which is the last hidden state of the LSTM/RNN. This is expected to be a good summary of the input sentence. All the intermediate states of the encoder are ignored, and the final state id supposed to be the initial hidden state of the decoder

2. **The decoder LSTM or RNN units produce the words in a sentence one after another**

In short, there are two RNNs/LSTMs. One we call the encoder – this reads the input sentence and tries to make sense of it, before summarizing it. It passes the summary (context vector) to the decoder which translates the input sentence by just seeing it.

Drawbacks of the Encoder-Decoder Approach (RNN/LSTM):

1. **Dependence on Encoder Summary:**

   o If the encoder produces a poor summary, the translation will also be poor.

   o This issue worsens with longer sentences.

2. **Long-Range Dependency Problem:**

   o RNNs/LSTMs struggle to understand and remember long sentences.

   o This happens due to the **vanishing/exploding gradient problem**, making it hard to retain distant information.

   o They perform better on recent inputs but forget earlier parts.

3. **Performance Degradation with Longer Inputs:**

   o Even the original creators (Cho et al., 2014) showed that translation quality drops as sentence length increases.

4. **LSTM Limitations:**

   o While LSTMs handle long-range dependencies better than RNNs, they still **fail in certain cases**.

o They can become "forgetful" over long sequences.
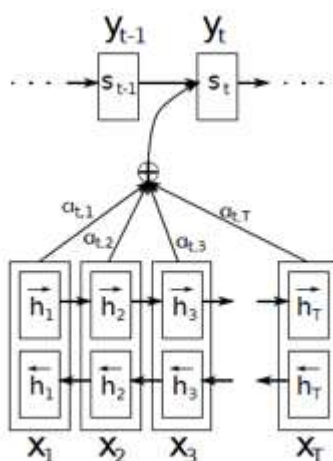
5. **Lack of Selective Attention:**

   o The model cannot **focus more on important words** in the input while translating.

   o All words are treated equally, even if some are more critical for translation.

Now, let's say, we want to predict the next word in a sentence, and its context is located a few words back. Here's an example – **"Despite originally being from Uttar Pradesh, as he was brought up in Bengal, he is more comfortable in Bengali".** In these groups of sentences, if we want to predict the word **"Bengali"**, the phrase **"brought up"** and **"Bengal"-** these two should be given more weight while predicting it. And although **Uttar Pradesh** is another state's name, it should be "ignored".

So is there any way we can keep all the relevant information in the input sentences intact while creating the context vector?

Bahdanau et al (2015) came up with a simple but elegant idea where they suggested that not only can all the input words be taken into account in the context vector, but relative importance should also be given to each one of them.

So, whenever the proposed model generates a sentence, it searches for a set of positions in the encoder hidden states where the most relevant information is available. This idea is called 'Attention'.



Understanding the Attention Mechanism

This is the diagram of the Attention model shown in Bahdanau's paper. The Bidirectional LSTM used here generates a sequence of annotations (h1, h2,….., hTx) for each input sentence. All the vectors h1,h2.., etc., used in their work are basically the concatenation of forward and backward hidden states in the encoder.

$$ h_j = \left[ \overrightarrow{h}_j^{\top} ; \overleftarrow{h}_j^{\top} \right]^{\top} $$

To put it in simple terms, all the vectors h1,h2,h3…., hTx are representations of Tx number of words in the input sentence. In the simple encoder and decoder model, only the last state of the encoder LSTM was used (hTx in this case) as the context vector.

But Bahdanau et al put emphasis on **embeddings** of all the words in the input (represented by hidden states) while creating the context vector. They did this by simply taking a weighted sum of the hidden states.

Now, the question is how should the weights be calculated? Well, the weights are also learned by a **feed-forward neural network** and I've mentioned their mathematical equation below.

The context vector ci for the output word yi is generated using the weighted sum of the annotations:

$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j.$$

The weights αij are computed by a softmax function given by the following equation:

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})},$$

$$e_{ij} = a(s_{i-1}, h_j)$$

**eij** is the output score of a feedforward neural network described by the function **a** that attempts to capture the alignment between input at **j** and output at **i**.

Basically, if the encoder produces **Tx** number of "annotations" (the hidden state vectors) each having dimension **d,** then the input dimension of the feedforward network is **(Tx , 2d)** (assuming the previous state of the decoder also has **d** dimensions and these two vectors are concatenated)**.** This input is multiplied with a matrix **Wa** of **(2d, 1)** dimensions (of course followed by addition of the bias term) to get scores **eij (having a dimension (Tx , 1))**.

On the top of these **eij** scores, a tan hyperbolic function is applied followed by a softmax to get the normalized alignment scores for output j:

**E = I [Tx*2d] * Wa [2d * 1] + B[Tx*1]**

**α = softmax(tanh(E))**

**C= IT * α**

So, **α** is a (Tx, 1) dimensional vector and its elements are the weights corresponding to each word in the input sentence.

Let **α** is **[0.2, 0.3, 0.3, 0.2]** and the input sentence is **"I am doing it".** Here, the context vector corresponding to it will be:

**C=0.2\*I"I" + 0.3\*I"am" + 0.3\*I"doing" + + 0.3\*I"it"** [**Ix** is the hidden state corresponding to the word **x**]

Global vs. Local Attention

So far, we have discussed the most basic Attention mechanism where all the inputs have been given some importance. Let's take things a bit deeper now.

The term "global" Attention is appropriate because all the inputs are given importance. Originally, the Global Attention (defined by Luong et al 2015) had a few subtle differences with the Attention concept we discussed previously.

The differentiation is that it considers all the hidden states of both the encoder LSTM and decoder LSTM to calculate a "variable-length context vector **ct**, whereas Bahdanau et al. used the previous hidden state of the unidirectional decoder **LSTM** and all the hidden states of the encoder LSTM to calculate the context vector.

In **encoder-decoder** architectures, the score generally is a function of the encoder and the decoder hidden states. Any function is valid as long as it captures the relative importance of the input words with respect to the output word.

**When a "global" Attention layer is applied, a lot of computation is incurred. This is because all the hidden states must be taken into consideration, concatenated into a matrix, and multiplied with a weight matrix of correct dimensions to get the final layer of the feedforward connection.**

So, as the input size increases, the matrix size also increases. In simple terms, the number of nodes in the feedforward connection increases and in effect it increases computation.

**Can we reduce this in any way? Yes! Local Attention is the answer**.

Intuitively, when we try to infer something from any given information, our mind tends to intelligently reduce the search space further and further by taking only the most relevant inputs.

The idea of Global and Local Attention was inspired by the concepts of Soft and Hard Attention used mainly in computer vision tasks.

Soft Attention is the global Attention where all image patches are given some weight; but in hard Attention, only one image patch is considered at a time.

But local Attention is not the same as the hard Attention used in the image captioning task. On the contrary, it is a blend of both the concepts, where instead of considering all the encoded inputs, only a part is considered for the context vector generation. This not only avoids expensive computation incurred in soft Attention but is also easier to train than hard Attention.

How can this be achieved in the first place? Here, the model tries to predict a position **pt** in the sequence of the embeddings of the input words. Around the position **pt**, it considers a window of size,

say, **2D**. Therefore, the context vector is generated as a weighted average of the inputs in a position **[pt – D,pt + D]** where **D** is empirically chosen.
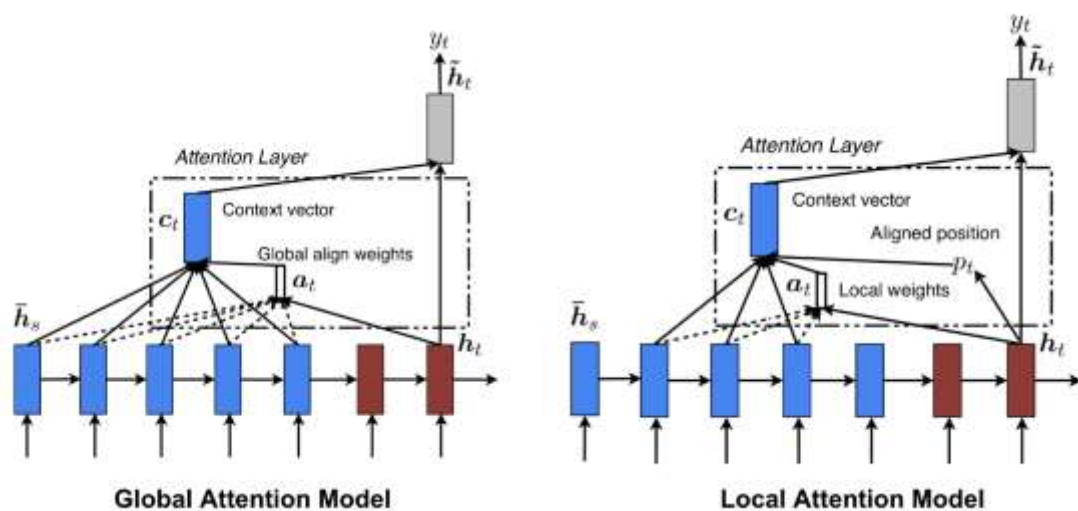
Furthermore, there can be two types of alignments:

1. **Monotonic alignment**, where pt is set to t, assuming that at time t, only the information in the neighborhood of t matters

2. **Predictive alignment** where the model itself predicts the alignment position as follows:

$$p_t = S \cdot \text{sigmoid}(\boldsymbol{v}_p^\top \tanh(\boldsymbol{W_p h_t}))$$

where 'Vp' and 'Wp' are the model parameters that are learned during training and 'S' is the source sentence length. Clearly, pt ε [0,S].

The figures below demonstrate the difference between the Global and Local Attention mechanism. Global Attention considers all hidden states (blue) whereas local Attention considers only a subset:
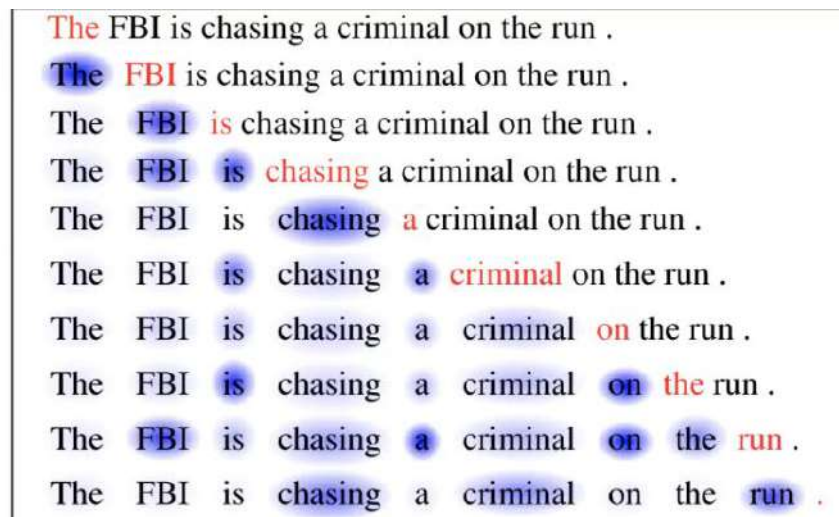


## Transformers – Attention is All You Need

The paper named **"Attention is All You Need"** by **Vaswani et al** is one of the most important contributions to Attention so far. They have redefined Attention by providing a very generic and broad definition of Attention based on **key**, **query,** and **values**. They have referenced another concept called **multi-headed Attention.** Let's discuss this briefly.

First, let's define what **"self-Attention"** is. Cheng et al, in their paper named **"Long Short-Term Memory-Networks for Machine Reading"**, defined self-Attention as the mechanism of relating different positions of a single sequence or sentence in order to gain a more vivid representation.

Machine reader is an algorithm that can automatically understand the text given to it. We have taken the below picture from the paper. The red words are read or processed at the current instant, and the blue words are the memories. The different shades represent the degree of memory activation.

When we are reading or processing the sentence word by word, where previously seen words are also emphasized on, is inferred from the shades, and this is exactly what self-Attention in a machine reader does.



Previously, to calculate the Attention for a word in the sentence, the mechanism of score calculation was to either use a dot product or some other function of the word with the hidden state representations of the previously seen words. In this paper, a fundamentally same but a more generic concept altogether has been proposed.

Let's say we want to calculate the Attention for the word "chasing". The mechanism would be to take a dot product of the embedding of "chasing" with the embedding of each of the previously seen words like "The", "FBI", and "is".

Key, Query, and Value Vectors in Attention Mechanism

1. **Three Vectors per Embedding**:

   o   Each word embedding has three ssociated vectors:

      ▪   **Key (K)**

      ▪   **Query (Q)**

      ▪   **Value (V)**

   o   These are derived using matrix multiplication.

2. **Calculating Attention**:

   o   To find the attention of a **target word** relative to input embeddings:

      ▪   Use the **Query (Q)** of the target word.

      ▪   Use the **Key (K)** of the input word.

   o   Compute a **matching score** (attention weight) from **Q** and **K**.

o   Use these scores as weights to sum the **Value (V)** vectors.

3. **What Are Key, Query, and Value?**

   o   They represent the same embedding in **different subspaces** (abstract perspectives).

   o   **Analogy**:

   ▪   **Query**: A question you ask.

   ▪   **Key**: The memory address being searched.

   ▪   **Value**: The data stored at that address.

4. **Mathematical Derivation**:

   o   If the embedding (**E**) has dimension **(D, 1)**:

   ▪   **Key (K)**: Multiply **E** by matrix **Wk** of size **(D/3, D)** → **K = Wk × E**

   ▪   **Query (Q)**: Multiply **E** by **Wq** → **Q = Wq × E**

   ▪   **Value (V)**: Multiply **E** by **Wv** → **V = Wv × E**

Now, to calculate the Attention for the word **"chasing"**, we need to take the dot product of the **query** vector of the embedding of **"chasing"** to the **key** vector of each of the previous words, i.e., the key vectors corresponding to the words **"The"**, **"FBI"** and **"is"**. Then these values are divided by D (the dimension of the embeddings) followed by a **softmax** operation. So, the operations are respectively:

- **softmax(Q"chasing"** . K"The" / D)

- **softmax(Q"chasing"** .K"FBI" / D)

- **softmax(Q"chasing"** . K"is" / D)

Basically, this is a function **f(Qtarget, Kinput)** of the query vector of the target word and the key vector of the input embeddings. It doesn't necessarily have to be a dot product of **Q** and **K.** Anyone can choose a function of his/her own choice.

Next, let's say the vector thus obtained is **[0.2, 0.5, 0.3]**. These values are the **"alignment scores"** for the calculation of Attention. These alignment scores are multiplied with the **value vector** of each of the input embeddings and these weighted value vectors are added to get the **context vector**:

**C"chasing"= 0.2 * VThe + 0.5* V"FBI" + 0.3 * V"is"**

Practically, all the embedded input vectors are combined in a single matrix **X,** which is multiplied with common weight matrices **Wk**, **Wq, Wv** to get **K, Q and V** matrices respectively. Now the compact equation becomes:

**Z=Softmax(Q*KT/**D**)V**

**Therefore, the context vector is a function of Key, Query and Value F(K, Q, V).**

Key Highlights of Attention Mechanisms

1. **Generalization of Attention Mechanisms**

   o The Bahdanau Attention and other previous attention works are special cases of the attention mechanisms discussed here.

   o **Key Feature**: A single embedded vector serves as the **Key (K)**, **Query (Q)**, and **Value (V)** vectors simultaneously.
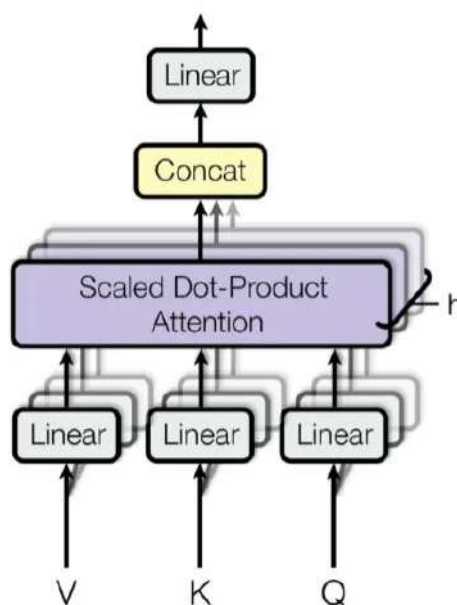
2. **Multi-Headed Attention Mechanism**

   o The input matrix **X** is multiplied by different weight matrices (**Wk, Wq, Wv**) to produce separate **K, Q, and V** matrices.

   o This results in multiple **Z matrices**, meaning each input word embedding is projected into different **"representation subspaces"**.
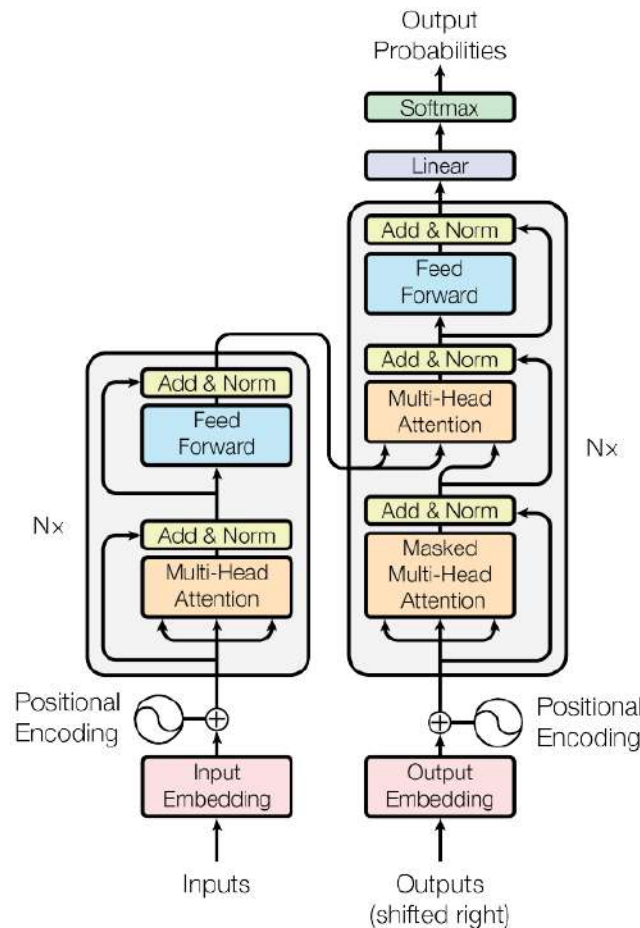
3. **Multi-Headed Self-Attention Example (3 Heads)**

   o For a word like **"chasing"**, there will be **3 different Z matrices** (called **"Attention Heads"**).

   o These heads are **concatenated** and then multiplied by a **single weight matrix**.

   o The final output is a **single Attention head** that combines information from all individual heads.

The picture below depicts the multi-head Attention. You can see that there are multiple Attention heads arising from different V, K, Q vectors, and they are concatenated:

The actual transformer architecture is a bit more complicated. You can read it in much more detail here



This image above is the transformer architecture. We see that something called 'positional encoding' has been used and added with the embedding of the inputs in both the encoder and decoder.

The models that we have described so far had no way to account for the order of the input words. They have tried to capture this through positional encoding. **This mechanism adds a vector to each input embedding, and all these vectors follow a pattern that helps to determine the position of each word, or the distances between different words in the input.**

As shown in the figure, on top of this positional encoding + input embedding layer, there are two sublayers:

1.  In the first sublayer, there is a multi-head self-attention layer. There is an additive residual connection from the output of the positional encoding to the output of the multi-head self-attention, on top of which they have applied a layer normalization layer. The **layer normalization** is a technique **(Hinton, 2016)** similar to batch normalization where instead of considering the whole minibatch of data for calculating the normalization statistics, all the hidden units in the same layer of the network have been considered in the calculations. This overcomes the drawback of estimating the statistics for the summed input to any neuron over a minibatch of the training samples. Thus, it is convenient to use in RNN/LSTM.

2. In the second sublayer, instead of the multi-head self-attention, there is a feedforward layer (as shown), and all other connections are the same.

On the decoder side, apart from the two layers described above, there is another layer that applies multi-head Attention on top of the encoder stack. Then, after a sublayer followed by one linear and one softmax layer, we get the output probabilities from the decoder.

**What are the Different Types of Attention Mechanisms?**

Introduction

Imagine standing in a dimly lit library, struggling to decipher a complex document while juggling dozens of other texts. This was the world of Transformers before the "Attention is All You Need" paper unveiled its revolutionary spotlight – the **attention mechanism**.
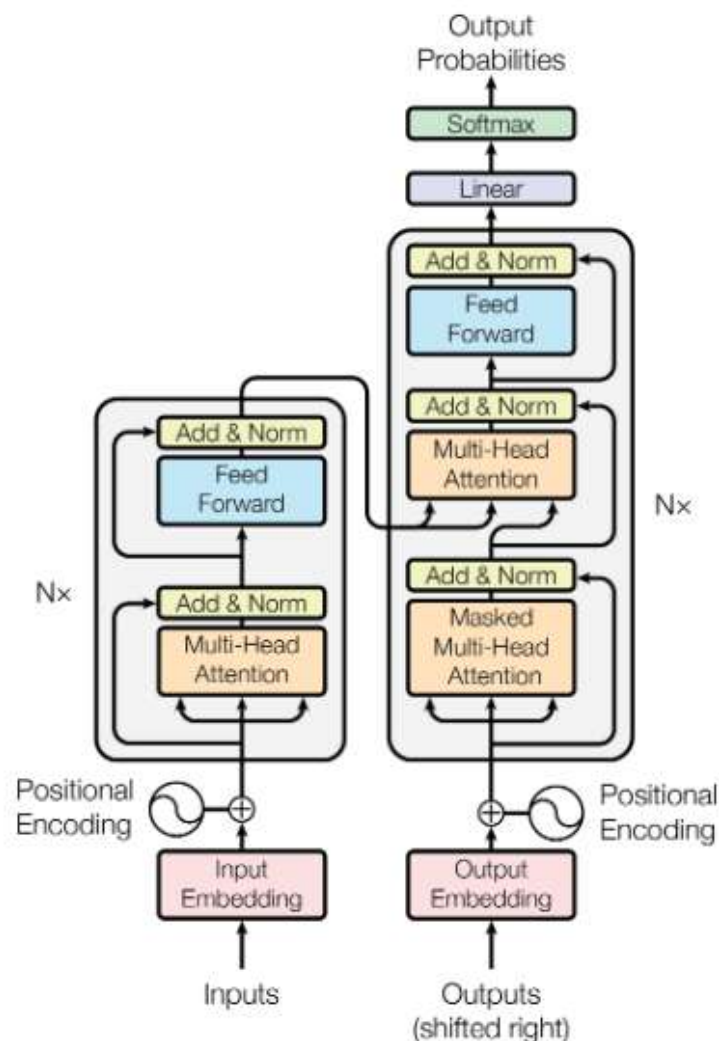


Figure 1: The Transformer - model architecture.

Limitations of RNNs

Traditional sequential models, like [Recurrent Neural Networks (RNNs)](#), processed language word by word, leading to several limitations:

- **Short-range dependence:** RNNs struggled to grasp connections between distant words, often misinterpreting the meaning of sentences like "the man who visited the zoo yesterday," where the subject and verb are far apart.

- **Limited parallelism:** Processing information sequentially is inherently slow, preventing efficient training and utilization of computational resources, especially for long sequences.

- **Focus on local context:** RNNs primarily consider immediate neighbors, potentially missing crucial information from other parts of the sentence.

These limitations hampered the ability of Transformers to perform complex tasks like machine translation and natural language understanding. Then came the **attention mechanism**, a revolutionary spotlight that illuminates the hidden connections between words, transforming our understanding of language processing. But what exactly did attention solve, and how did it change the game for Transformers?

Let's focus on three key areas:

Long-range Dependency

- **Problem:** Traditional models often stumbled on sentences like "the woman who lived on the hill saw a shooting star last night." They struggled to connect "woman" and "shooting star" due to their distance, leading to misinterpretations.

- **Attention Mechanism:** Imagine the model shining a bright beam across the sentence, connecting "woman" directly to "shooting star" and understanding the sentence as a whole. This ability to capture relationships regardless of distance is crucial for tasks like machine translation and summarization.

Parallel Processing Power

- **Problem:** Traditional models processed information sequentially, like reading a book page by page. This was slow and inefficient, especially for long texts.

- **Attention Mechanism:** Imagine multiple spotlights scanning the library simultaneously, analyzing different parts of the text in parallel. This dramatically speeds up the model's work, allowing it to handle vast amounts of data efficiently. This parallel processing power is essential for training complex models and making real-time predictions.

Global Context Awareness

- **Problem:** Traditional models often focused on individual words, missing the broader context of the sentence. This led to misunderstandings in cases like sarcasm or double meanings.

- **Attention Mechanism:** Imagine the spotlight sweeping across the entire library, taking in every book and understanding how they relate to each other. This global context awareness allows the model to consider the entirety of the text when interpreting each word, leading to a richer and more nuanced understanding.

Disambiguating Polysemous Words

- **Problem:** Words like "bank" or "apple" can be nouns, verbs, or even companies, creating ambiguity that traditional models struggled to resolve.

- **Attention Mechanism:** Imagine the model shining spotlights on all occurrences of the word "bank" in a sentence, then analyzing the surrounding context and relationships with other words. By considering grammatical structure, nearby nouns, and even past sentences, the attention mechanism can deduce the intended meaning. This ability to disambiguate polysemous words is crucial for tasks like machine translation, text summarization, and dialogue systems.

These four aspects – long-range dependency, parallel processing power, global context awareness, and disambiguation – showcase the transformative power of attention mechanisms. They have propelled Transformers to the forefront of natural language processing, enabling them to tackle complex tasks with remarkable accuracy and efficiency.
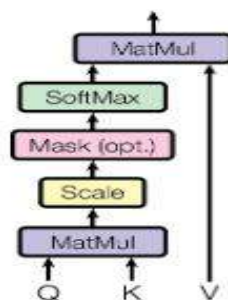
As NLP and specifically LLMs continue to evolve, attention mechanisms will undoubtedly play an even more critical role. They are the bridge between the linear sequence of words and the rich tapestry of human language, and ultimately, the key to unlocking the true potential of these linguistic marvels. This article delves into the various types of attention mechanisms and their functionalities.

1. Self-Attention: The Transformer's Guiding Star

Imagine juggling multiple books and needing to reference specific passages in each while writing a summary. Self-attention or Scaled Dot-Product attention acts like an intelligent assistant, helping models do the same with sequential data like sentences or time series. It allows each element in the sequence to attend to every other element, effectively capturing long-range dependencies and complex relationships.

Here's a closer look at its core technical aspects:



Scaled Dot-Product Attention

Vector Representation

Each element (word, data point) is transformed into a high-dimensional vector, encoding its information content. This vector space serves as the foundation for the interaction between elements.

QKV Transformation

Three key matrices are defined:

- **Query (Q):** Represents the "question" each element poses to the others. Q captures the current element's information needs and guides its search for relevant information within the sequence.

- **Key (K):** Holds the "key" to each element's information. K encodes the essence of each element's content, enabling other elements to identify potential relevance based on their own needs.

- **Value (V):** Stores the actual content each element wants to share. V contains the detailed information other elements can access and leverage based on their attention scores.

Attention Score Calculation

The compatibility between each element pair is measured through a dot product between their respective Q and K vectors. Higher scores indicate a stronger potential relevance between the elements.

Scaled Attention Weights

To ensure relative importance, these compatibility scores are normalized using a softmax function. This results in attention weights, ranging from 0 to 1, representing the weighted importance of each element for the current element's context.

Weighted Context Aggregation

Attention weights are applied to the V matrix, essentially highlighting the important information from each element based on its relevance to the current element. This weighted sum creates a contextualized representation for the current element, incorporating insights gleaned from all other elements in the sequence.

Enhanced Element Representation

With its enriched representation, the element now possesses a deeper understanding of its own content as well as its relationships with other elements in the sequence. This transformed representation forms the basis for subsequent processing within the model.

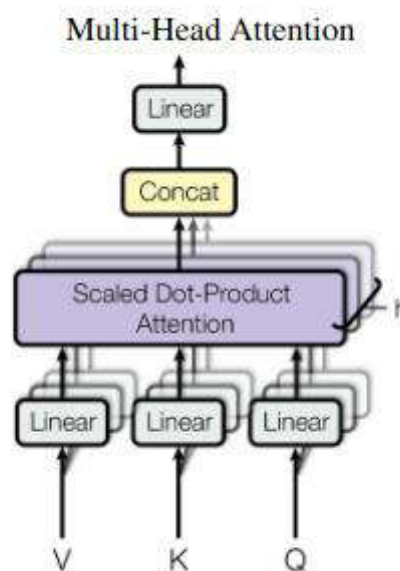This multi-step process enables self-attention to:

- **Capture long-range dependencies:** Relationships between distant elements become readily apparent, even if separated by multiple intervening elements.

- **Model complex interactions:** Subtle dependencies and correlations within the sequence are brought to light, leading to a richer understanding of the data structure and dynamics.

- **Contextualize each element:** The model analyzes each element not in isolation but within the broader framework of the sequence, leading to more accurate and nuanced predictions or representations.

Self-attention has revolutionized how models process sequential data, unlocking new possibilities across diverse fields like machine translation, natural language generation, time series forecasting, and beyond. Its ability to unveil the hidden relationships within sequences provides a powerful tool for uncovering insights and achieving superior performance in a wide range of tasks.

2. Multi-Head Attention: Seeing Through Different Lenses

Self-attention provides a holistic view, but sometimes focusing on specific aspects of the data is crucial. That's where multi-head attention comes in. Imagine having multiple assistants, each equipped with a different lens:
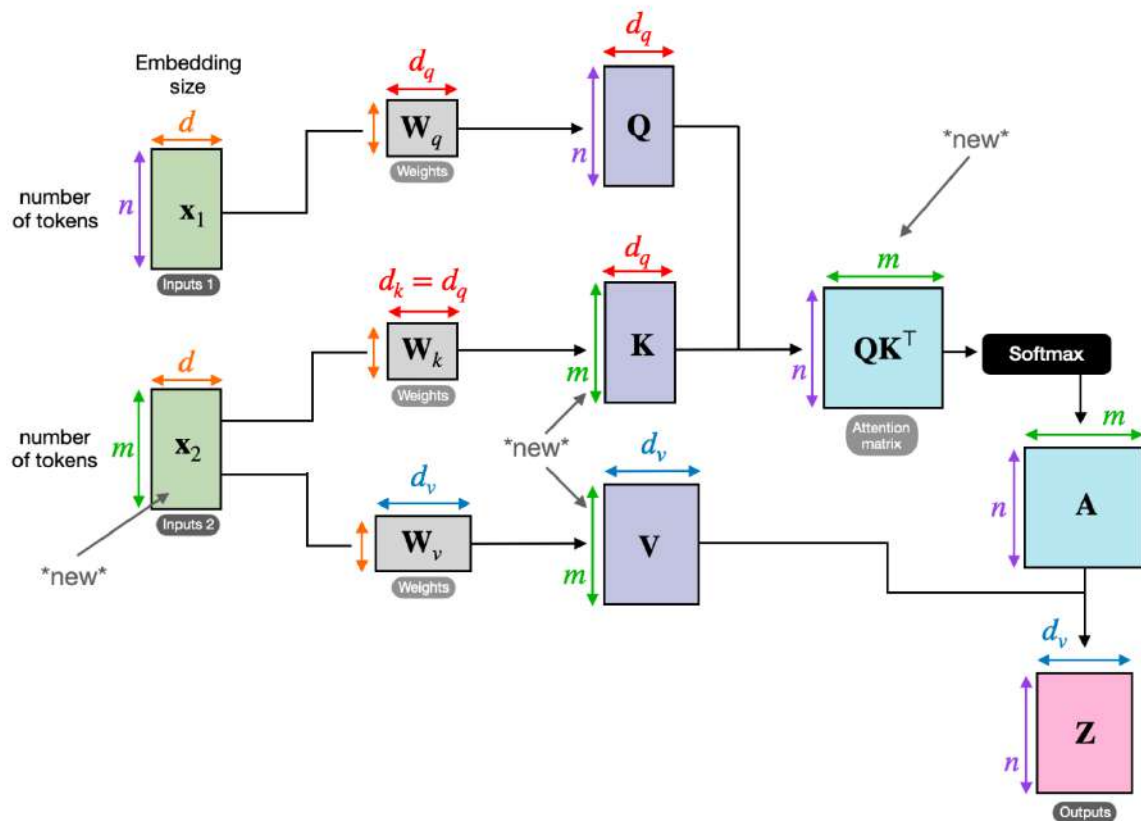


Multi-Head Attention

**Multiple "heads"** are created, each attending to the input sequence through its own Q, K, and V matrices.

- Each head learns to focus on different aspects of the data, like long-range dependencies, syntactic relationships, or local word interactions.

- The outputs from each head are then concatenated and projected to a final representation, capturing the multifaceted nature of the input.

This allows the model to simultaneously consider various perspectives, leading to a richer and more nuanced understanding of the data.

3. Cross-Attention: Building Bridges Between Sequences

The ability to understand connections between different pieces of information is crucial for many NLP tasks. Imagine writing a book review – you wouldn't just summarize the text word for word, but rather draw insights and connections across chapters. Enter **cross-attention**, a potent mechanism that builds bridges between sequences, empowering models to leverage information from two distinct sources.



In encoder-decoder architectures like Transformers, the **encoder** processes the input sequence (the book) and generates a hidden representation.

- The **decoder** uses cross-attention to attend to the encoder's hidden representation at each step while generating the output sequence (the review).

- The decoder's Q matrix interacts with the encoder's K and V matrices, allowing it to focus on relevant parts of the book while writing each sentence of the review.
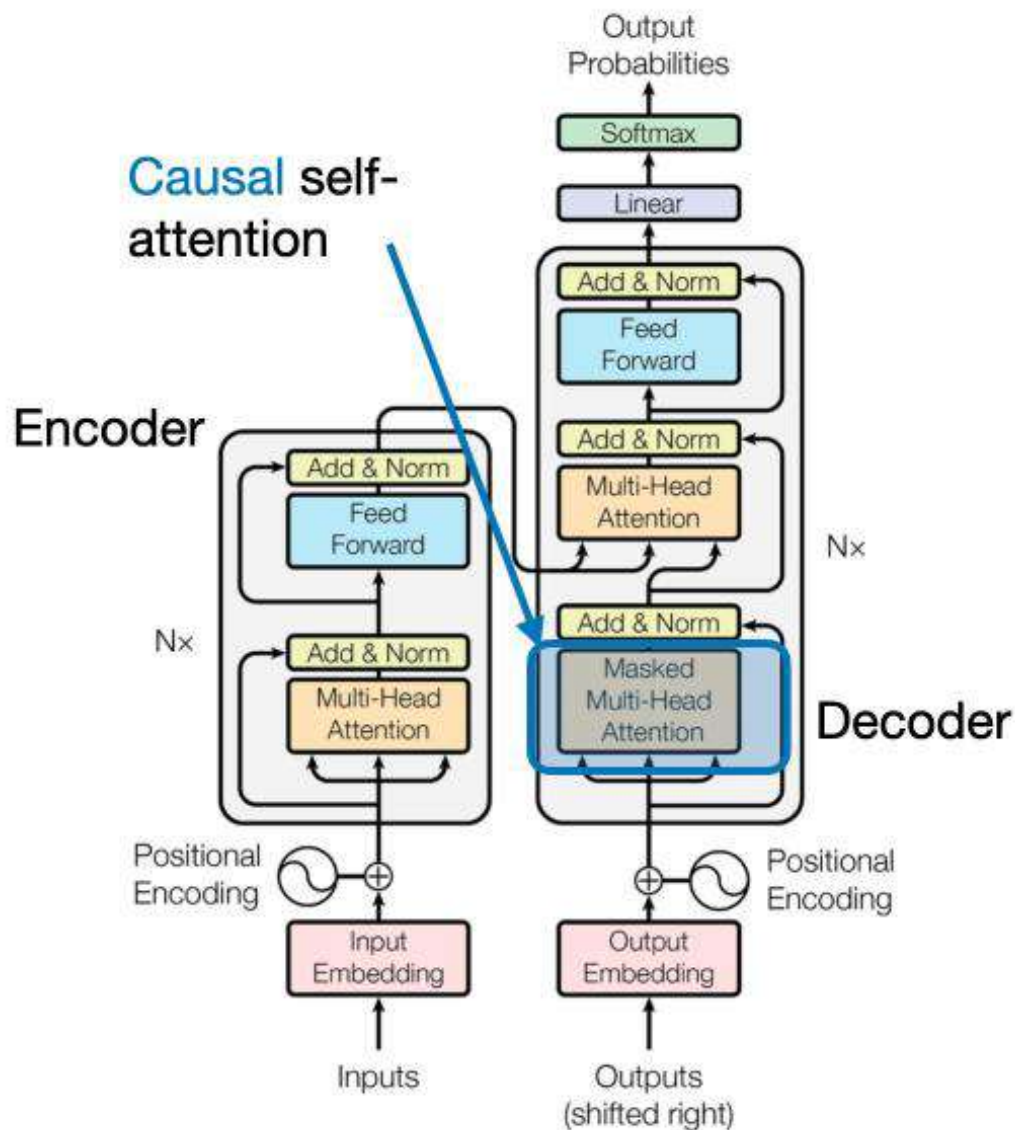
This mechanism is invaluable for tasks like machine translation, summarization, and question answering, where understanding the relationships between input and output sequences is essential.

4. Causal Attention: Preserving the Flow of Time

Imagine predicting the next word in a sentence without peeking ahead. Traditional attention mechanisms struggle with tasks that require preserving the temporal order of information, such as text generation and time-series forecasting. They readily "peek ahead" in the sequence, leading to inaccurate predictions. Causal attention addresses this limitation by ensuring predictions solely depend on previously processed information.

Here's How it Works

- **Masking Mechanism:** A specific mask is applied to the attention weights, effectively blocking the model's access to future elements in the sequence. For instance, when predicting the second word in "the woman who…", the model can only consider "the" and not "who" or subsequent words.

- **Autoregressive Processing:** Information flows linearly, with each element's representation built solely from elements appearing before it. The model processes the sequence word by word, generating predictions based on the context established up to that point.



- Causal attention is crucial for tasks like text generation and time-series forecasting, where maintaining the temporal order of the data is vital for accurate predictions.

5. Global vs. Local Attention: Striking the Balance

Attention mechanisms face a key trade-off: capturing long-range dependencies versus maintaining efficient computation. This manifests in two primary approaches: **global attention** and **local**

**attention**. Imagine reading an entire book versus focusing on a specific chapter. Global attention processes the whole sequence at once, while local attention focuses on a smaller window:

- **Global attention** captures long-range dependencies and overall context but can be computationally expensive for long sequences.

- **Local attention** is more efficient but might miss out on distant relationships.

The choice between global and local attention depends on several factors:

- **Task requirements**: Tasks like machine translation require capturing distant relationships, favoring global attention, while sentiment analysis might favor local attention's focus.

- **Sequence length**: Longer sequences make global attention computationally expensive, necessitating local or hybrid approaches.

- **Model capacity**: Resource constraints might necessitate local attention even for tasks requiring global context.

To achieve the optimal balance, models can employ:

- **Dynamic switching**: use global attention for key elements and local attention for others, adapting based on importance and distance.

- **Hybrid approaches**: combine both mechanisms within the same layer, leveraging their respective strengths.

.