

## COPY RIGHT



ELSEVIER  
SSRN

**2023 IJEMR.** Personal use of this material is permitted. Permission from IJEMR must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works. No Reprint should be done to this paper, all copy right is authenticated to Paper Authors

IJEMR Transactions, online available on 10<sup>th</sup> Apr 2023. Link

[:http://www.ijiemr.org/downloads.php?vol=Volume-12&issue=Issue 04](http://www.ijiemr.org/downloads.php?vol=Volume-12&issue=Issue 04)

**10.48047/IJEMR/V12/ISSUE 04/84**

Title **VAGUS - AI LANGUAGE MODEL ASSISTANT USING NATURAL LANGUAGE PROCESSING TECHNIQUES**

Volume 12, ISSUE 04, Pages: 690-698

Paper Authors

**Mr.B. Kalyan Chakravarty, Katragadda Mukhesh Raghava, Kondeddula Karthikeya, Paruchuri Saketh**



USE THIS BARCODE TO ACCESS YOUR ONLINE PAPER

To Secure Your Paper As Per **UGC Guidelines** We Are Providing A Electronic Bar Code

## VAGUS - AI LANGUAGE MODEL ASSISTANT USING NATURAL LANGUAGE PROCESSING TECHNIQUES

**Mr.B. Kalyan Chakravarty (M.Tech)<sup>1</sup>**, assistant professor, Department of IT, Vasireddy Venkatadri Institute of Technology, Nambur, Guntur Dt., Andhra Pradesh.

**Katragadda Mukhesh Raghava<sup>2</sup>, Kondeddu Karthikeya<sup>3</sup>, Paruchuri Saketh<sup>4</sup>**  
<sup>2,3,4</sup> UG Students, Department of IT, Vasireddy Venkatadri Institute of Technology, Nambur, Guntur Dt., Andhra Pradesh.

<sup>1</sup> kalyan.battula@vvit.net, <sup>2</sup> coolmukhesh@gmail.com,

<sup>3</sup> kondeddulakarhikeya@gmail.com,

<sup>4</sup> Sak.paruchuri64@gmail.com

### Abstract

An AI language model assistant that uses natural language processing techniques is a breath-taking technological advancement that has changed the way we communicate. Natural language processing and artificial intelligence have combined to create a revolutionary tool that not only understands but also responds to human language. Natural language processing is a branch of computer science that focuses on teaching computers to understand and respond to human language. An AI language model assistant can help people write better emails, articles, gather information, generate code in various programming languages, and even make suggestions based on user queries for the most appropriate words or phrases to use. It can also help with grammar and spelling, ensuring that the text is easy to read and understand. To accomplish this, the AI language model assistant analyses the context of the text and employs statistical models to predict the best next word or phrase. This is accomplished by training the AI language model assistant on large datasets, which allows it to learn patterns and relationships among data. AI language models assist users in conveying their message and clearly expressing their thoughts. People who are not native speakers of a language can also benefit from AI language model assistants. The assistant can provide translations, correct grammar, and recommend more commonly used words in the language. Non-native speakers will find it easier to communicate and express themselves as a result of this. To summarize, AI language model assistants are powerful tools that can help people communicate more effectively and efficiently gather knowledge and information. Overall, artificial intelligence language model assistants are bridging the gap between humans and machines, making communication easier and more effective for everyone.

**Keywords:** Artificial intelligence, Deep Learning, Neural Networks, Natural Language Processing, Transformer, T5, GPT, BERT, RoBERTa, and ELECTRA

## Introduction

The importance of language in human communication cannot be exaggerated, as it allows us to express our thoughts and ideas in ways that transcend time and distance. With the advent of Artificial Intelligence (AI), we are now able to fully utilize the power of language by developing intelligent systems that can understand and respond to natural language. AI Language Model Assistants have emerged as a game-changing technology that is transforming the way we interact with language by utilizing advanced Natural Language Processing (NLP) techniques. These remarkable systems are capable of comprehending the subtleties of human speech, processing massive amounts of data, and generating intelligent responses that rival human communication capabilities. More such AI assistants are heralding a new era of intelligent communication that is anticipated to influence the way we live, work, and communicate by seamlessly integrating with language processing technologies and with advanced machine learning and deep learning algorithms.

**Neural Networks** is the most Astonishing Invention that Replicates the Human Brain. Neural networks are the AI technology that most closely resembles the human brain. Similar to how our brains absorb information, these potent computing systems are built to identify patterns, learn from data, and make predictions. Picture a network of neurons, or connected brain cells, interpreting the

environment around us. Similar methods are employed by neural networks to process enormous volumes of data and come to wise conclusions. Let's examine neural networks' operation in more detail. A neural network is first trained using a sizable dataset of samples for input and output. These examples might be any number of datasets, from a million to a billion. To increase the accuracy of its predictions and query responses, the network analyses the data and modifies its parameters. The actual magic of neural networks occurs during this procedure, which is referred to as inference. Many applications, such as speech recognition, image recognition, and natural language processing, have made use of neural networks. In order to show their adaptability and capacity for creativity and other things, they have even been used to generate music and art. The capacity of neural networks to learn and adapt over time is one of its most important advantages.

**Natural Language Processing** is a stunning technology that gives machines language power. Moreover, have you ever communicated with a chatbot or virtual assistant and felt as if you were conversing with a human? That is the power of NLP which stands for natural language processing, an enthralling technology that allows machines to understand, interpret, and generate human language. Natural language processing (NLP) is a branch of neural networks that teaches computers to read,

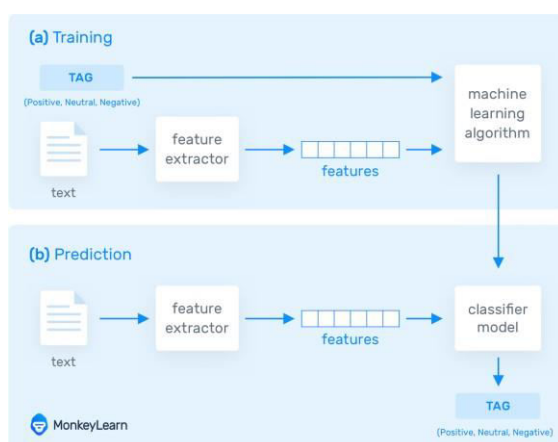
understand, and interpret human language. This technology has advanced over time, allowing machines to communicate with humans in a more natural and intuitive manner. A set of algorithms and techniques that enable machines to process and analyse large amounts of textual data is at the heart of NLP. Among these techniques are machine learning, pattern recognition, and semantic analysis. Chatbots, virtual assistants, and AI language model assistants are some of the most exciting applications of NLP. These intelligent systems are capable of understanding natural language queries and responding with pertinent information, making them extremely useful for customer service and support. NLP has numerous other applications, including sentiment analysis, which can assist businesses in gauging public opinion about their products or services. It is also useful for language translation, speech recognition, and writing assistance, making it a versatile and potent technology.

## Transformer NLP model

Google researchers unveiled the Transformer model in 2017. It was created to increase machine translation performance, which is the task of translating text from one language to another. The Transformer model, on the other hand, has since been utilized for a variety of NLP applications, including language modelling, sentiment analysis, and question answering. The Transformer model is built on a neural network known as the self-attention mechanism. While producing predictions, this method allows the model to focus on different areas of the input text. When translating a statement from English to French, for example, the model can focus on the words that are most significant for comprehending the text's meaning. One of the Transformer model's primary advantages is its ability to generate high-quality text. This is due to the model's ability to grasp human linguistic characteristics such as sarcasm, irony, and humor. This makes it suitable for activities like chatbots, where the goal is to engage consumers in a natural and engaging discussion.

## T5 NLP model

The T5 NLP model is a cutting-edge natural language processing technology developed by Google researchers. It is a transformer model, which means it is built on a neural network architecture meant to analyze massive volumes of text data. The T5 model is distinct in that it is a "text-to-text" model, which means it may accept any sort of text input and



**Fig.1.** NLP model working

produce any type of text output. As a result, it is extremely versatile and suitable for a wide range of NLP tasks, including language translation, summarization, and question answering. The T5 model is likewise extremely precise and efficient. It was trained on vast volumes of text data, allowing it to learn patterns and relationships between words and sentences. This training procedure is known as pre-training, and it is what makes the T5 model so effective.

### **GPT NLP model**

The GPT NLP model is a sort of language model created by OpenAI, a research organization which is dedicated to the advancement of artificial intelligence. The model's neural network design is referred to as GPT (Generative Pre-trained Transformer). The GPT model is intended for the processing and generation of human language. It is trained using vast volumes of text data to learn the patterns and relationships between words and phrases. This training procedure is known as pre-training, and it is what makes the GPT model so effective. After pre-trained, the GPT model can be fine-tuned for specific NLP tasks. This method of fine-tuning entails training the model on a smaller dataset that is specific to the job at hand. If the goal is language translation, for example, the model would be fine-tuned using a dataset of translated text. One of the primary features of the GPT model is its capacity to generate high-quality text. This is due to the model's ability to grasp human linguistic characteristics such as

sarcasm, irony, and humor. This makes it perfect for jobs like chatbots and conversational agents, where the goal is to engage consumers in a natural and engaging conversation.

### **BERT and RoBERTa models**

BERT and RoBERTa are two highly effective natural language processing (NLP) models developed by Google and Facebook, respectively. These models are based on the transformer neural network architecture, which is meant to process and create human language. BERT and RoBERTa are two highly effective natural language processing (NLP) models developed by Google and Facebook, respectively. These models are based on the transformer neural network architecture, which is meant to process and create human language. In contrast, RoBERTa stands for "Robustly Optimized BERT Pretraining Approach." It is a BERT variant that has been improved for performance on a variety of NLP tasks. This optimisation entails modifying the training method and hyperparameters to increase the accuracy and efficiency of the model. Finally, BERT and RoBERTa are two cutting-edge NLP models that have transformed the field of natural language processing. Its capacity to comprehend word context and provide very accurate predictions makes them useful for a variety of applications.

### **ELECTRA NLP model**

ELECTRA is a natural language processing (NLP) model created by Google researchers in 2020. It is a form of neural network architecture intended to increase

the efficiency and accuracy of NLP activities. The acronym ELECTRA stands for "Efficiently Learning an Encoder that Correctly Classifies Token Replacements." This refers to the model's capacity to learn the links between words and phrases in a sentence efficiently. ELECTRA's capacity to undertake "discriminative pre-training" is one of its primary advantages. This entails teaching the model to differentiate between actual and phoney words in a sentence.

The T5, BERT, RoBERTa, GPT and ELECTRA models all are based on Transformer NLP model and Transformer uses Self-Attention Algorithm

### **System Implementation**

The process of incorporating natural language processing (NLP) models consists of numerous steps, which are as follows. The first step is to prepare the data. This entails gathering and cleansing the data that will be used to train the model. The next step is to choose a model. There are various NLP models available, such as GPT and BERT, each with its own set of strengths and shortcomings. The next stage is to train the model after it has been chosen. This entails putting the prepared data into the model and modifying its parameters to maximize its performance. After training, the model must be tested to ensure that it is accurate and dependable. The model can be put into the software application or system after it has been tested and validated.

**Prerequisites:** To follow along, the user will require the following:

1. NodeJS installed
2. OpenAI package installed
3. ViteJS installed
4. Rust compiler installed
5. AI language models Datamodels, pre-trained GPT, BERT, Transformer and AffectNet models available in the form of Datasets and API's

### **Rust implementation of Language**

#### **Model**

Here we are going to show BERT model as example, which is a common artificial intelligence (AI) paradigm for natural language processing (NLP) that may be used for tasks such as text classification, sentiment analysis, and question answering. Rust, on the other hand, is a fast and safe systems programming language. When the two are combined, they can produce a powerful and efficient NLP solution.

**Step1:-**Install the Required Libraries

**Step2:-**Load the BERT Model

**Step3:-**Tokenization: To prepare the input text for processing by the BERT model tokenize it next.

**Step4:-**Input Encoding: After tokenization the input text must be encoded for processing by the BERT model.

**Step5:-** Inference: Lastly, the BERT model can be used to do inference on the input text. This is a pre-trained BERT model for sequence classification tasks.

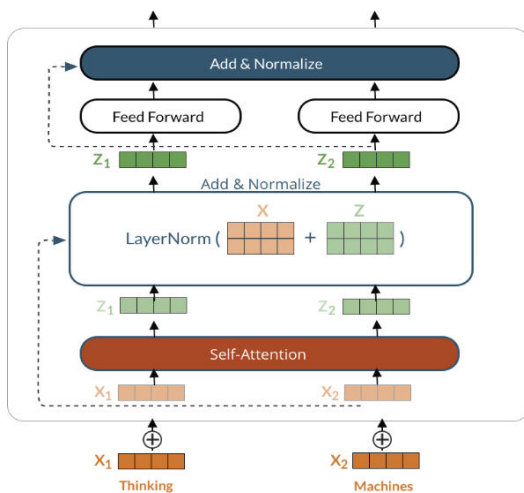
**Step6:-**Deployment: Once the BERT model has been implemented in Rust, it can be used in a variety of applications.

### **Python implementation of GPT-Neo model**

**Step 1:-** Importing Library Files

**Step 2:-** Load the GPT-Neo model

**Step 3:-** Set seed for reproducibility: To ensure that the results are reproducible, you can set a seed using the set\_seed function.



**Fig.2.** GPT-Neo model

**Step 4:-** Generate text: Once the model is loaded, you can generate text using the generator object and specifying the input prompt, maximum length of the generated text, and other parameters such as do\_sample and temperature.

### Self-Attention NLP algorithm

Self-Attention is a neural network layer used in natural language processing (NLP) models, notably those based on Transformers. By allowing it to weigh the relevance of different words while computing the output representation of each word, the Self-Attention mechanism assists the model in understanding the context and links between words in a phrase or sequence.

The following is how the Self-Attention algorithm works:

- 1.) For each word in the input sequence, compute the Query, Key, and Value vectors.
- 2.) Using the dot product, compute the similarity between the Query vector and all of the Key vectors.
- 3.) Scale the similarity scores by dividing them by the square root of the Key vector's dimensionality.
- 4.) Use a soft max function to generate a probability distribution over the input sequence's values.
- 5.) Multiply the probability distribution by the Value vectors to get the weighted sum, which is the attention layer's output.
- 6.) On top of the attention output, apply feedforward neural networks to further process the input.
- 7.) Repeat steps 1-6 for multiple layers of self-attention and feedforward neural networks to capture different levels of

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

abstraction in the input data.

If it is Multi-head then,

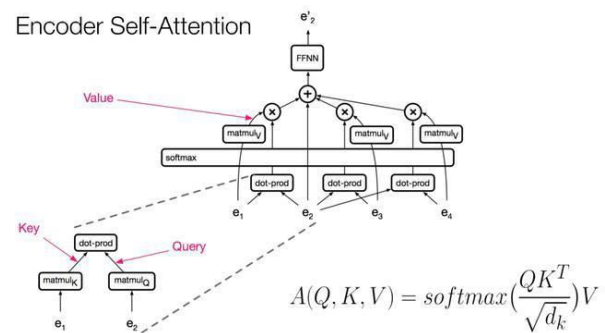
$$Multi\text{-}head(Q, K, V) = \text{concat}(\text{head}1, \text{head}2, \dots, \text{head}n, )Wo$$

Where, head i = Attention (QW base (i) power Q, KW base (i) power K and VW base (i) power V)

### Case Study

The **encoder** is a critical component of the Self-Attention NLP algorithm that processes input sequences of tokens such as words or subwords. To construct a context-aware representation of each token in the input sequence, the encoder

employs a multi-layered neural network design with self-attention mechanisms. The encoder is made up of several layers, each of which performs two basic functions: self-attention and feedforward neural network. The encoder computes the relevance of each token in the input sequence relative to all other tokens during the self-attention operation by comparing their Query, Key, and Value vectors using dot products. This enables the encoder to record long-term token dependencies and build context-aware representations. The output of the encoder's final layer is a context-aware representation of each token in the input sequence. These encoded representations can then be transmitted through the NLP model's subsequent layers to accomplish tasks like classification, translation, and creation. The encoder's capacity to capture long-range connections between tokens in the input sequence, which was difficult for previous NLP models, is one of the encoder's advantages in Self-Attention NLP algorithms. The Self-Attention method enables the model to learn contextual associations between words without being restricted by fixed position-based encoding techniques by allowing the encoder to consider the value of distinct words in a phrase or sequence when generating or predicting the next word. This enables the model to make more accurate predictions on a variety of NLP tasks, such as machine translation, question answering, and sentiment analysis.

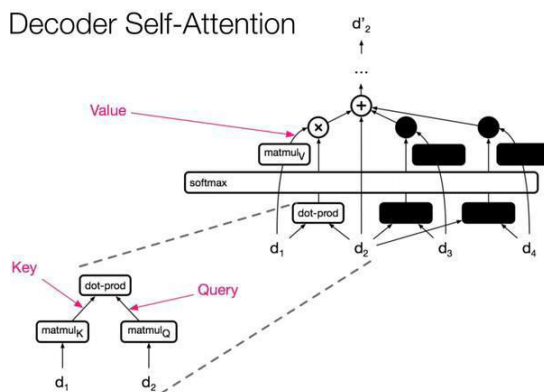


**Fig.3.** Encoder in self-attention

The **decoder** in self-attention NLP algorithms is in charge of generating the output sequence based on the encoded input sequence and a previously generated token. The self-attention mechanism is used by the decoder to attend to the relevant elements of the input sequence and generate the next token in the output sequence. Several layers of self-attention and feed-forward neural networks are commonly used in the decoder. The previous created token, the encoded input sequence, and the context vector generated by the preceding layer are all inputs to the decoder at each layer. After that, the decoder computes the attention scores between the query (prior token) and the keys and values (encoded input sequence) and generates a weighted sum of the values to obtain the context vector for the current layer. To generate the output token probabilities, the context vector is passed through a feed-forward neural network with residual connections and layer normalisation. The decoder selects the next token in the output sequence with the highest probability and sends it back as input to the next layer. Overall, the decoder in self-



attention NLP algorithms is critical in generating coherent output sequences by attending to relevant sections of the input sequence and using prior created tokens to inform the development of the future tokens.



**Fig.4.** Decoder in self-attention

## Conclusion

In this article, we created a friendly conversational language model based on NLP, and the model is trained with a massive amount of data to perform various NLP tasks such as question answering, text generation, text summarization, e-mail writing, content writing, article writing, gathering information, gaining knowledge, solving mathematical problems, generating algorithms and code in various programming languages, and more. To conclude the development of AI language model assistants utilizing NLP techniques is a wonderful and stunning breakthrough that has the potential to change the way we engage with technology. These assistants, which are powered by advanced algorithms and techniques like as transformers, GPT, and BERT, can grasp human language and

create increasingly realistic and intuitive responses. We should anticipate to see even more intriguing and inventive applications as these models continue to grow, making our interactions with technology more fluid, efficient, and productive than ever before. Natural language processing (NLP), a subfield of neural networks that is a part of Deep Learning, which is a branch of Machine Learning and artificial intelligence, is one of the most powerful approaches utilized in these language model assistants. NLP algorithms process and analyze vast volumes of data using machine learning, pattern recognition, and semantic analysis, allowing robots to communicate with people in a more natural and understandable manner.

Ultimately, we arrive to the name **“VAGUS”**. The vagus nerve is a key nerve that connects the brain to other sections of the body, including the heart and activating it will result in the whole degree of human intelligence being unlocked. This nerve is vital to our emotional intelligence. Our AI language model assistant is designed to connect humans with technology in a more natural and intuitive way, similar to how the vagus nerve connects different parts of the body and intelligence. It is also built using Neural Networks which is yet another AI concept inspired by the workings of neurons in the human brain. Moreover, it is always crucial to choose a name that accurately symbolizes the purpose and essence of our creation. Moreover the word VAGUS is easy to pronounce.

## References

- [1] aswani, A.Shazeer, N.Parmar, N.Uszkoreit, J.Jones, L.Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. In *Advances in neural information processing systems* (pp. 5998-6008).
- [2] Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- [3] Radford, A., Narasimhan, K., Salimans, T., & Sutskever, I. (2018). Improving language understanding by generative pre-training. *OpenAI Blog*, 1(8).
- [4] Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., .& Amodei, D. (2020). Language models are few-shot learners. In *Advances in Neural Information Processing Systems* (pp.1877-1901)
- [5] Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen D.& Stoyanov, V. (2019). Roberta: A robustly optimized bert pretraining approach.*arXiv preprint arXiv:1907.11692*.
- [6] "Attention Is All You Need" by Vaswani et al. (2017): This paper introduced the Transformer architecture, which revolutionized natural language processing (NLP) and led to state-of-the-art performance in many language tasks.
- [7] "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding" by Devlin et al. (2018): This paper introduced the Bidirectional Encoder Representations from Transformers (BERT) model, which achieved state-of-the-art performance on a wide range of NLP tasks.
- [8] T5: Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer" by Raffel et al. (2020): This paper introduced T5, a text-to-text transformer that achieved state-of-the-art performance on several NLP tasks and showed that it is possible to learn a single model that can be fine-tuned for a wide range of language tasks.
- [9] RoBERTa: A Robustly Optimized BERT Pretraining Approach" by Liu et al. (2019): This paper introduced RoBERTa, which achieved state-of-the-art performance on several NLP tasks by optimizing the BERT pre-training process.
- [10] "ELECTRA: Pre-training Text Encoders as Discriminators Rather Than Generators" by Clark et al. (2020): This paper introduced ELECTRA, which trains a generator model to produce corrupted input, and a discriminator model to distinguish between real and corrupted input, achieving state-of-the-art results on several NLP tasks.