

# DESIGNING A SEQ2SEQ CHATBOT WITH ATTENTION MECHANISM FOR REAL WORLD APPLICATIONS

<sup>1</sup>Dr.N.Sai Sindhuri, <sup>2</sup>Nelaturu Radha Krishna

<sup>1</sup>Associate Professor, <sup>2</sup>UG Student, <sup>1,2</sup>Department of Computer Science & Engineering, Geethanjali Institute of Science and Technology, Gangavaram, Andhra Pradesh, India

## Abstract

Chatbots have become widespread in various industries, serving as virtual assistants, customer support agents, and information providers. However, designing a chatbot capable of engaging in meaningful and contextually relevant conversations remains a significant challenge. In this paper, we propose a Seq2Seq (Sequence-to-Sequence) chatbot architecture enhanced with an attention mechanism for real-world applications. The Seq2Seq model, consisting of encoder and decoder components, is well-suited for generating responses based on input sequences.

We outline the architectural design of the Seq2Seq chatbot, detailing the encoder-decoder framework and the integration of attention mechanisms. We discuss the training process, including data preprocessing, model training, and optimization techniques. Additionally, we highlight the importance of incorporating large, diverse datasets to ensure the chatbot's ability to handle a wide range of user inputs and contexts.

## INTRODUCTION

The design of a Seq2Seq chatbot with attention mechanism for real-world applications involves creating an advanced conversational AI system capable of generating contextually relevant responses to user queries across diverse domains. Leveraging the sequence-to-sequence (Seq2Seq) architecture, the chatbot processes user inputs through an encoder to capture semantic meaning and contextual information, then decodes this representation using attention mechanisms to produce coherent and engaging responses. This approach enhances the chatbot's understanding of context and enables it to adapt to various conversational scenarios, ensuring a seamless and personalized user experience. By integrating with existing communication platforms and systems, the designed chatbot aims to address practical considerations such as scalability and integration while continuously improving through iterative refinement based on user feedback and performance evaluation.

## CONVERSATIONAL AI

Conversational AI is a set of technologies that work together to automate human-like communications via both speech and text between a person and machines. It combines Artificial Intelligence, Natural Language Processing, Machine Learning to understand, interpret and respond to user inputs in a way that simulates human conversation.

Aiming to connect humans and computers, it comprises a bunch of cutting-edge technologies to construct synthetic brainpower that further makes machines or chatbots capable of understanding, reading, & responding to the human language. And we call it conversational AI in technical terms.

You can find the essence of conversational AI in IVR systems, messaging platforms, voice-based communication channels, Chatbot, mobile apps, & many other channels. When employees usually spend 16% of their time in in-house communication & collaboration, conversational AI can reduce the time spent on such activities with automated & immediate referrals to customers' issues, according to Bloom fire.

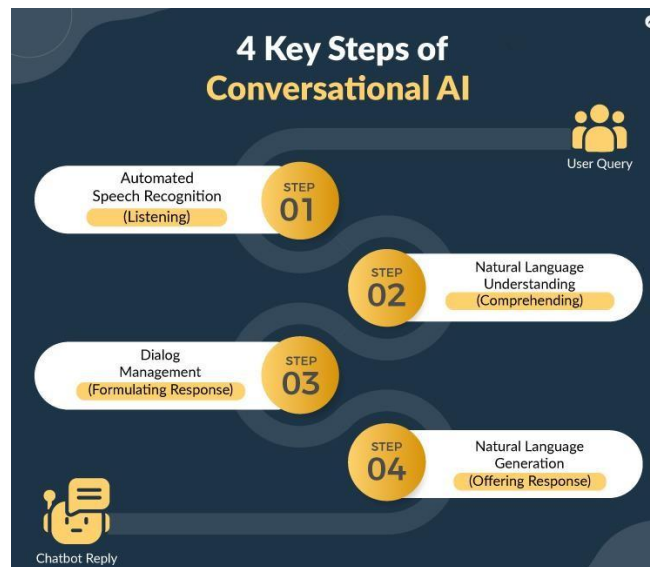


Fig.1. Steps in Conversational AI

## OBJECTIVE

The objective of designing a Seq2Seq chatbot with attention mechanism for real-world applications is to develop an advanced conversational AI system capable of generating contextually relevant and coherent responses to user queries across diverse domains. By leveraging the Seq2Seq architecture augmented with attention mechanisms, the goal is to enhance the chatbot's understanding of context, improve the quality of generated responses, and foster a more natural and engaging user experience.

## LITERATURE SURVEY

In this chapter we review some papers to get knowledge and understanding on the techniques had been proposed. All those techniques have the same aim which is track the vehicle and estimate the velocity of the moving vehicle. As Archimedes once said, “Man has always learned from the past. After all, you can't learn history in reverse!” it is essential for man to learn from history. Thus, considering all past researches, the most relevant research glimpses have been picked to be explained in detail. The overview shall discuss relevant aspects contributing to our research.

**Deep Reinforcement Learning for Dialogue Generation" by Jiwei Li, Will Monroe, Tianlin Shi, et al. (2016):**

The paper builds upon existing research in natural language processing and dialogue systems, leveraging advances in deep learning architectures such as recurrent neural networks (RNNs) and sequence-to-sequence (Seq2Seq) models. The authors draw inspiration from reinforcement learning frameworks, particularly policy gradient methods, to train a dialogue generation model capable of producing contextually coherent and engaging responses. This work builds on prior research exploring various methods for dialogue generation, including rule-based systems, retrieval-based approaches, and generative models.

**J. Vanian, "Google Adds More Brainpower to Artificial Intelligence Research Unit in Canada," Fortune, 21 November 2016. [Online]:**

This article contributes to the literature on AI research and industry developments by highlighting Google's commitment to advancing AI technologies and leveraging the expertise available in Canada's research community. The expansion discussed in the article reflects the growing importance of AI research and innovation on a global scale, as major tech companies like Google invest in talent and resources to drive progress in this field. By augmenting its AI research capabilities in Canada, Google aims to strengthen its position as a leader in AI research and development, fostering collaboration with academic institutions and contributing to the growth of Canada's AI ecosystem.

**Survey on Neural Network-Based Approaches for Dialogue Systems" by Hongshen Chen, Xiaorui Liu, Dawei Yin, et al. (2017):**

By reviewing a wide range of literature, the paper synthesizes the advancements, challenges, and trends in this rapidly evolving field. It surveys various neural network architectures, including recurrent neural networks (RNNs), convolutional neural networks (CNNs), and their combinations, such as hierarchical RNNs and memory networks, for dialogue modeling and generation. The paper also explores key components of dialogue systems, such as intent detection, slot filling, and response generation, and discusses how neural network-based approaches have been applied to improve the performance of these components. Additionally, the survey provides insights into the datasets, evaluation metrics, and benchmarking efforts in the domain of dialogue systems. Overall, this paper serves as a valuable resource for researchers and practitioners interested in understanding the state-of-the-art neural network-based approaches for dialogue systems and identifying directions for future research.

**"Attention is All You Need" by Ashish Vaswani et al. (2017):**

The literature survey surrounding this paper encompasses a wide range of research areas, including neural machine translation, language understanding, and sequence modeling. Building upon the foundation laid by previous works in recurrent neural networks (RNNs) and convolutional neural networks (CNNs), Vaswani et al. introduce the Transformer architecture, which relies solely on attention mechanisms without using recurrent or convolutional layers. The paper surveys the limitations of traditional Seq2Seq models, such as long-range dependencies and sequential computation inefficiency, and demonstrates how the Transformer overcomes these challenges by enabling parallelization and capturing global dependencies through self-attention mechanisms. Furthermore, the paper discusses the broader impact of the Transformer architecture on various natural language processing tasks, including machine translation, text summarization, and language modeling. By introducing a novel and efficient architecture for sequence modeling, "Attention is All You Need" has sparked a wave of research and innovation in the deep learning community, leading to further advancements in neural network-based approaches for natural language understanding and generation.

**"A Survey of Available Corpora for Building Data-Driven Dialogue Systems" by Sashank J. Reddi, Nikola Mrkšić, Milica Gašić, et al. (2019):**

The literature survey surrounding this paper delves into the growing importance of corpora in advancing the field of dialogue systems, particularly in the context of machine learning and natural language processing. The survey discusses a wide range of dialogue corpora, encompassing different domains, languages, and conversational styles, sourced from various sources including social media, customer service interactions, and open-domain conversations. The paper highlights the diversity and richness of these corpora, exploring their suitability for different tasks such as dialogue modeling, intent detection, and response generation. Additionally, the survey examines the challenges and limitations associated with existing dialogue corpora, such as data sparsity, domain specificity, and ethical considerations. By providing a comprehensive overview of available dialogue corpora and their characteristics, the paper serves as a valuable resource for researchers and practitioners seeking to develop and evaluate data-driven dialogue systems.

## EXISTING SYSTEM

The framework used in existing system relies on a basic structure and with limited natural language capabilities. Before the arrival of Seq2Seq models, the machine translation systems relied on statistical methods and phrase-based approaches. The most popular approach was the use of phrase-based statistical machine translation (SMT) systems. That was not able to handle long-distance dependencies and capture global context. However, this chatbot assistance is a communication channel that is improved in the understanding of natural language processing and understanding capabilities. The chatbot is created using

## PROPOSED SYSTEM

Our proposed system aims to develop an advanced conversational AI chatbot leveraging the Seq2Seq architecture with an attention mechanism for real-world applications. Seq2Seq models have demonstrated effectiveness in generating coherent and contextually relevant responses in conversational contexts. By integrating attention mechanisms into the model, we enhance its ability to focus on salient parts of the input sequence, improving the quality and relevance of generated responses. The proposed system will consist of an

encoder-decoder framework, where the encoder processes input sequences and the decoder generates corresponding responses. The encoder will employ recurrent neural networks (RNNs) or transformer architectures to encode input sequences into a fixed-length vector representation. The decoder, augmented with attention mechanisms, will use this representation to generate output sequences one token at a time.

**SYSTEM ARCHITECTURE**

The System consists of the following steps :- 1. Encoder  
2. Decoder

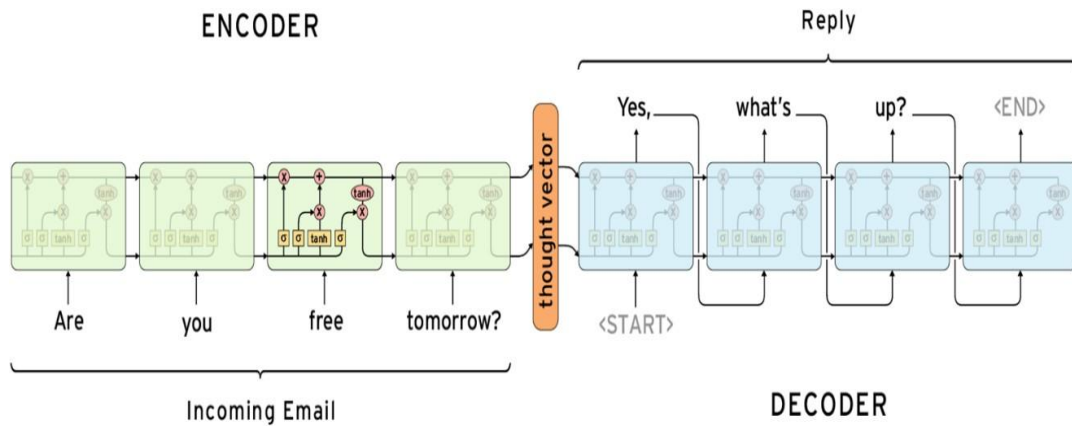


Fig.2. System Architecture

Designing a Seq2Seq chatbot with an attention mechanism for real-world applications requires a meticulously crafted architecture to ensure robust performance. The architecture typically consists of an encoder-decoder framework augmented with attention mechanisms. The encoder component processes input sequences, such as user queries or messages, into fixed-length vector representations, capturing the semantic information effectively. Meanwhile, the decoder generates responses based on the encoded information, iteratively producing tokens while attending to relevant parts of the input sequence using attention mechanisms. This attention mechanism enables the model to focus on relevant context during both encoding and decoding stages, enhancing the chatbot's ability to generate contextually appropriate responses.

Moreover, integrating the chatbot architecture with techniques like bidirectional recurrent neural networks (RNNs) or transformer models can further improve performance. Bidirectional RNNs enable the model to capture contextual information from both past and future tokens in the input sequence, enhancing its understanding of the conversation context. On the other hand, transformer models leverage self-attention mechanisms to capture global dependencies efficiently, enabling the model to attend to all input tokens simultaneously.

**ENCODER**

The Encoder plays a crucial role in the architecture of a Seq2Seq chatbot with attention mechanism, serving as the initial processing unit for input sequences. Typically implemented as recurrent neural networks (RNNs), Long Short-Term Memory (LSTM), or Transformer models, the Encoder processes the input text, such as user queries or messages, into fixed-length vector representations, capturing the semantic information effectively. It reads the input sequence token by token, updating its internal state at each time step to encode the information. The Encoder's primary objective is to transform variable-length input sequences into fixed-length continuous representations called context vectors, which encapsulate the semantic meaning of the input sequence. This transformation enables the subsequent decoding phase to generate meaningful responses based on the encoded context. Additionally, modern architectures often incorporate bidirectional RNNs or self-attention mechanisms within the Encoder to capture contextual information from both past and future tokens simultaneously, enhancing the model's ability to understand and represent the input sequence comprehensively.

**DECODER**

The Decoder is a pivotal component in the Sequence-to-Sequence (Seq2Seq) architecture, responsible

for generating output sequences based on the encoded information provided by the Encoder. Typically implemented as recurrent neural networks (RNNs), Long Short-Term Memory (LSTM), or Transformer models, the Decoder receives the context vector, which encapsulates the semantic meaning of the input sequence, from the Encoder. During the decoding process, the Decoder uses this context vector along with its internal state to generate the output sequence token by token. It operates in an autoregressive manner, meaning that it predicts one token at a time, considering the previously generated tokens. Attention mechanisms are often integrated into the Decoder to allow the model to focus on relevant parts of the input sequence while generating each token of the output sequence. This attention mechanism enables the model to capture dependencies between elements in the input and output sequences effectively, improving the model's ability to generate accurate and contextually relevant responses.

### SEQUENCE TO SEQUENCE ALGORITHM

The sequence-to-sequence (Seq2Seq) algorithm is a popular approach in natural language processing (NLP) for tasks like machine translation and chatbot development.

**Encode:** First, you have an "encoder" that takes in your input sequence, like a sentence in one language. Imagine it as a conveyor belt where each word in the sentence is placed one after another. The encoder processes these words step by step, keeping track of their meanings and relationships. At the end of the conveyor belt, the encoder summarizes everything it has seen into a single vector, capturing the essence of the input sequence.

**Decode:** Next, you have a "decoder" that takes this summarized information (the single vector) and uses it to generate an output sequence, like a translation in another language or a response in a conversation. It's like having a blank canvas, and the decoder paints one word at a time based on the information it received from the encoder. It consults the single vector from the encoder and starts producing the output sequence word by word, trying to capture the meaning of the input sequence.

**Attention:** It helps the decoder focus on the most relevant parts of the input sequence at each step of generating the output sequence. This makes the translation or response more accurate and contextually relevant.

**Training:** To teach this system, you give it many examples of input-output pairs. During training, the system learns to adjust its encoder and decoder so that the generated output sequences closely match the desired ones.

**Inference:** Once trained, you can use this system to translate new sentences or generate responses in conversations. You feed a new input sequence into the encoder, get the summarized vector, and then let the decoder generate the output sequence based on that vector. With the attention mechanism, it can focus on the right parts of the input sequence as needed.

### Results & Analysis

#### Evaluation Criteria

In this study, we construct a chatbot for real-time product filtration and we use it to display the details of required product of user that are available in various e-commerce platforms such as flipkart, amazon, etc., In this process, we: a) Use chatbot for user interface. b) RASA framework to effectively provide product details. c) Use Sequence to Sequence algorithm to process the conversation.

The execution of the process will be explained clearly with the help of the continuous screenshots. The whole process in the execution is giving input to the chatbot and it will automatically process the input and provide the respective output of product based on user requirements. This whole process is done in four simple steps. Each figure mentioned below are the simultaneous process of outputs.

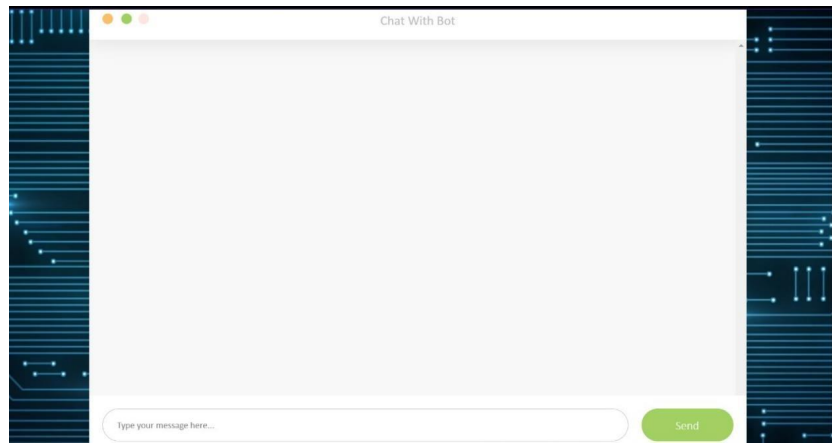


Fig : 3 The Web Application/Chat Interface

**Description :** Fig 3, describes how the website of project looks like.

In this section a detailed analysis of the trained models is given. They are also compared with a baseline seq2seq model.

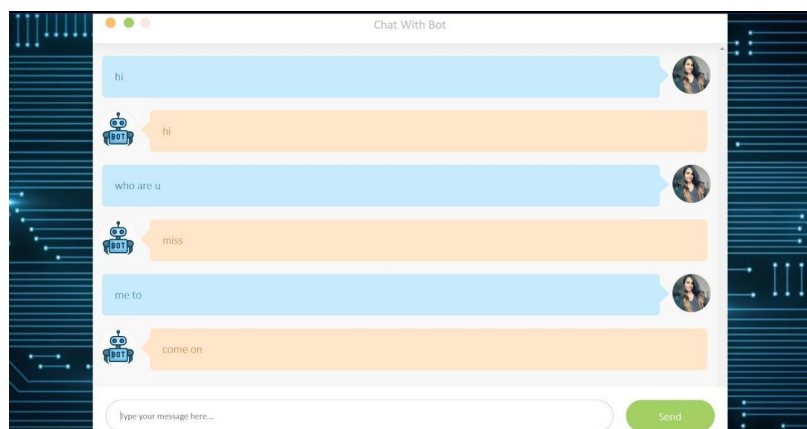


Figure : 4 Conversation Between User and Chatbot

**Description :** The above fig 4 represents the conversation between user and chatbot.

The user enters the input in user interface. Next the input is trained by Encoder and the Encoder understand the user intents and create the output response based on input. The output is generated by Decoder.

## CONCLUSION

The project “Designing seq2seq Chatbot with Attention Mechanism for Real World Applications” presents a robust solution for the interaction between human and a chatbot with a emotional bonding across language like English. By leveraging pre-trained movie conversations, Recurrent Neural Network(RNN), Long Short term Memory(LSTM) using Encoder-Decoder, the chatbot can preprocess the data into numerical form for better understanding with the help of Tensor Flow libraries. Through rigorous testing and evaluation, the system demonstrates its effectiveness in real world applications.

## FUTURE ENHANCEMENTS

After many efforts we had successfully detected and recognized the NLP chatbot. Designing seq2seq Chatbot with Attention Mechanism for Real World Applications has some of the future enhancements they are:

**Language Expansion:** Include other languages using seq2seq of Recurrent Neural Network(RNN) to enhance users usability.

**Improving UI:** By exploring the new features and improvisation we can improve the designing of Chatbot UI with the help of Flask Framework through front-end web development page.

## References

- [1] Dataset collect and information about Cornell movie dialog corpus dataset available at <https://www.cs.cornell.edu/cristian/CornellMovieDialogsCorpus.htm>.
- [2] Neural Machine Translation by Jointly Learning to Align and Translate Dzmitry Bahdanau, Kyunghyun Cho, Yoshua Bengio (Submitted on 1 Sep 2014 (v1), last revised 19 May 2016 (this version, v7)).
- [3] Sequence-to-Sequence Chatbot build using tensorflow update June 28,2016 [Online], Available:<http://complx.me/2016-06-28-easy-seq2seq/>.
- [4] Richard Krisztian Csaky, Department Of Automation And Applied Informatics” Deep Learning Based Chatbot Models Scientific Students’associations Report” 2017. [Paper], Available:<https://tdk.bme.hu/VIK/DownloadPaper/asdad>.
- [5] Oriol Vinyals and Quoc V. Le. A neural conversational model. CoRR, abs/1506.05869, 2015 and Oriol Vinyals, et al., Show and Tell: A Neural Image Caption Generator, 2014.
- [6] J. Vanian, ”Google Adds More Brainpower to Artificial Intelligence Research Unit in Canada,” Fortune, 21 November 2016. [Online]
- [7] Anjana Tiha(April 26,2019), Intelligent Chatbot using Deep Learning,[Paper] Available:[https://www.researchgate.net/publication/328582617\\_Intelligent\\_Chatbot\\_using\\_Deep\\_Learning](https://www.researchgate.net/publication/328582617_Intelligent_Chatbot_using_Deep_Learning).
- [8] Barak Turovsky (November 15, 2016), ”Found in translation: More accurate, fluent sentences in Google Translate”, Google Blog, Retrieved January 11, 2017.
- [9] Mike Schuster, Melvin Johnson, and Nikhil Thorat (November 22, 2016), ”Zero-Shot Translation with Google’s Multilingual Neural Machine Translation System”, Google Research Blog, Retrieved January 11, 2017.
- [10] Gil Fewster (January 5, 2017), ”The mind-blowing AI announcement from Google that you probably missed”, freeCodeCamp, Retrieved January 11, 2017.
- [11] Kyunghyun Cho, et al., Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation, 2014.