# Experimental Research on Encoder-Decoder Architectures with Attention for Chatbots

## Nitu Kumari¹, Shailesh Kumar Singh²

*¹Student M.Sc(Physics),Monad University, Hapur*
*²Associate Professor, Department of Physics, HimalayanGarhwal University Uttarakhand*

---------------------------------------------------------------------***---------------------------------------------------------------------

**Abstract -** *Chatbots goal at routinely providing a communication among a human and a computer. While there may be a protracted music of studies in rule-primarily based totally and retrieval-primarily based totally strategies, the generation-primarily based totally strategies are promisingly rising fixing problems like responding to queries in inference that had been now no longer formerly visible in improvement or schooling time. In this paper, we provide an experimental view of the way latest advances in near regions as gadget translation may be followed for chatbots. In particular, we examine how opportunity encoder-decoder deep getting to know architectures carry out withinside the context of chatbots. Our studies concludes that a totally attention-primarily based totally structure is capable of outperform the recurrent neural community baseline system.*

*Key Words***:** Chatbot, encoder-decoder, attention mechanisms.

## 1.INTRODUCTION

A chatbot stands for the fast model of chat plus robotic and it's far a pc software that conducts a human-device verbal exchange in any topic.One of the first actual chatbots become rule-primarily based totally. It become proposed in 1966 via way of means of Joseph Weizenbaum's software ELIZA [13]. Input sentences have been analyzed the usage of numerous predefined decomposition rules, and after that key phrases have been used to generate responses to them. The Artificial Intelligence Markup Language (AIML) is an evolution of those first rule-primarily based totally chatbots. This AIML follows the concept of defining written styles and the corresponding templates which might be responses to the styles. Then, in inference, if the robotic identifies a sample in a sentence from a user, the robotic is capable of respond taking the corresponding template [11].To lessen the quantity of labor that growing those styles and templates calls for, opportunity chatbots, now not rule-primarily based totally, however retrieval-primarily based totally have been proposed. These structures use exceptional talk databases to teach an facts retrieval system [2]. The massive benefit of those retrieval-primarily based totally structures is that their schooling calls for little human dedication. However, those structures nevertheless depend upon giving the maximum suitable reaction from a

hard and fast of sentences, which limits their overall performance withinside the case of unseen events.

Thanks to the emergent deep gaining knowledge of techniques, the unconventional generative-primarily based totally methods have arisen imparting chatbots which can be capable, for the primary time, to reply to non-predefined sentences. The first a success technique is primarily based totally at the famous encoder-decoder structure, which has been correctly utilized in pretty some herbal language applications, and, moreover, it's been prolonged to photo and speech processing [10, 12]. One a success implementation of this encoder-decoder structure in herbal language processing has been the latest concatenation of recurrent neural networks [7, 3]. In fact, this structure builds on pinnacle of recurrent neural language models [6] via way of means of including an encoder step and a decoder step. In the encoder step, a recurrent neural community converts an enter series into a hard and fast illustration (known as idea vector). This illustration is fed withinside the recurrent neural community from the decoder step which permits the decoder version to output greater shrewd predictions given the context from the encoding. While this implementation has proven a few effects in chatbots [10], the principle downside is that lengthy sequences aren't nicely codified right into a unmarried vector. This project is confronted thru the latest interest-primarily based totally mechanisms [1, 9] these days proposed for device translation.

The most important contribution of this paper is the software of the experimentation of those interest-primarily based totally mechanisms [1, 9] to chatbots. Taking [10] as beginning point, we evaluate the encoder-decoder structure with interest [1] and the transformer [9]. A manually carried out assessment suggests that the latter is capable of outperform the encoder-decoder with interest that is already higher than the encoder-decoder baseline structure.

The rest of the paper is organized as follows. Section 2 in short introduces the deep gaining knowledge of architectures used on this work, which essentially are encoder-decoder primarily based totally on recurrent neural networks (without or with interest mechanism) and the transformer which makes use of a completely interest-primarily based totally encoder-decoder with out recurrent neural networks. Section three information the

experimental framework, particularly, facts records and parameters from structures. Section four describes the guide assessment. Section five discusses insights of effects and contributions of this study

## 2. Encoder-Decoder Architectures and Attention-based Mechanisms

An autoencoder is a sort of neural community that targets at getting to know a illustration of the enter at the same time as taking into account a deciphering of this illustration via way of means of minimizing the improving error. A generalization of this structure is the encoder-decoder which lets in for inputs and outputs to be different. This structure, see a diagram of it in Figure 1 (left), has emerged as an powerful paradigm for managing variable-period inputs and outputs. Although this effectiveness, its simplicity limits their performance. For example, it appears unrealistic that the enter statistics has to suit in a set period vector of the inner illustration. It appears extra affordable to take enter statistics progressively at the same time as we're producing the output. That is why, interest-primarily based totally mechanisms have arisen as a technique to this matter.

In this section, we in brief offer a high-stage description of this encoder-decoder with recurrent neural networks plus a success encoder-decoder implementations that use interest-primarily based totally mechanisms. Among those architectures we first describe the encoder-decoder primarily based totally on recurrent neural networks with interest primarily based totally on multi-layer perceptron. And, we 2d describe the transformer encoder-decoder structure that makes use of best a mixture of feed-ahead neural networks with extra state-of-the-art interest primarily based totally on a couple of heads.

## 3. Encoder-Decoder with Recurrent Neural Networks

Given an enter sentence, the encoder iteratively computes for every phrase a hidden country vector the use of the phrase and former hidden country of the recurrent neural network (RNN). Once the complete sentence has been analyzed, the applicable data of the enter sentence is contained withinside the final hidden country of the RNN, called context or notion vector. The decoder computes, phrase with the aid of using phrase, an output withinside the unique illustration area the use of the data contained withinside the context vector and former decoded words.The structure implementation can range de-pending at the form of RNN mobileular used (authentic RNN mobileular, a LSTM mobileular [4] or a GRU mobileular [3]), range of cells according to layer or the range of hidden layers amongst different parameters. Figure 1 (left) suggests a diagram of this structure.

One of the primary drawbacks of this structure is living withinside the truth that as the scale of the enter sentence

increases, the encoder wishes to compress a massive amount of data right into a fixed-duration vector.This is awful compressing system that could yield to a negative overall performance of the chatbot.

## 4. Encoder-Decoder with Recurrent Neural Networks with Attention

To triumph over the aforementioned downside of the fundamental RNN-primarily based totally encoder-decoder approach, an interest mechanism is usually used with inside the decoder [1]. In this case, for every generated word, the decoder computes a context vector composed of the weighted sum of all hidden kingdom vectors of the encoder rather than counting on the cappotential of the encoder to compress the complete enter series into the notion vector.
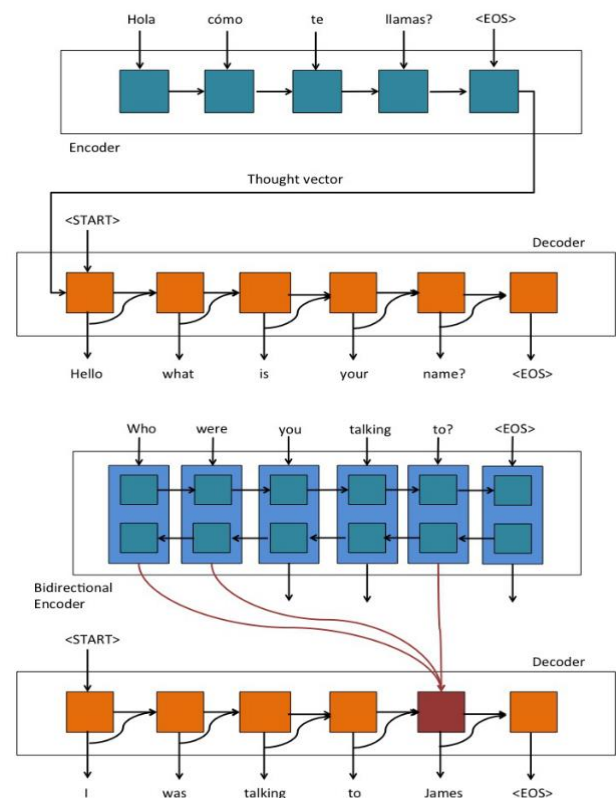


**Fig. 1.** (Left) Encoder-decoder with RNNs; (Right) Encoder-de

weights are computed with the aid of using an alignment version and normalized over all values to get a percent of the way applicable the phrase from the enter sentence is, on the subject of the phrase to be decoded, see discern 1 (right) displaying the diagram of this architecture. For in addition technical rationalization of the way weights are computed see

## 5.  Transformer

While preceding structure has been efficiently implemented to system translation, there are nevertheless a few troubles to solve. The structure in exercise may be genuinely sluggish to educate and given the manner RNNs address sequences, it isn't always clean to parallelize the set of rules and take gain of new computational assets which includes Tensor Processing Units (TPUs). Motivated with the aid of using this issue, the Transformer version [9] has been proposed and it's been demonstrated to be aggressive withinside the mission of system translation. The Transformer version is capable of enhance modern outcomes in more than one instructional benchmarks at the same time as dashing up schooling with the aid of using an order of significance in evaluation to RNN-primarily based totally encoder-decoder with interest proven in preceding section.

The Transformer structure is largely an encoder-decoder which concatenates interest-primarily based totally mechanisms permitting to version relationships among phrases with out requiring recurrence. More specifically, there are 3 principal degrees withinside the encoder (see Figure 2). The first one is in which enter phrases are projected right into a vector illustration area with the aid of using an embedding matrix and then, for the reason that there's no records of the order and role of phrases withinside the enter sentence, a positional encoding is delivered to the embedded enter vectors. Note that during preceding RNN encoder/decoder models, because of their sequential nature, no positional records is required.

The 2nd degree is a multi-head interest block (of Self-Attention on this first case) that linearly initiatives the enter records into exclusive area representations and plays interest over all of them. This approach lets in the version to become aware of exclusive semantic, morphological and lexical traits of the enter series and attend them one at a time on the deciphering process. coder with RNNS and interest.
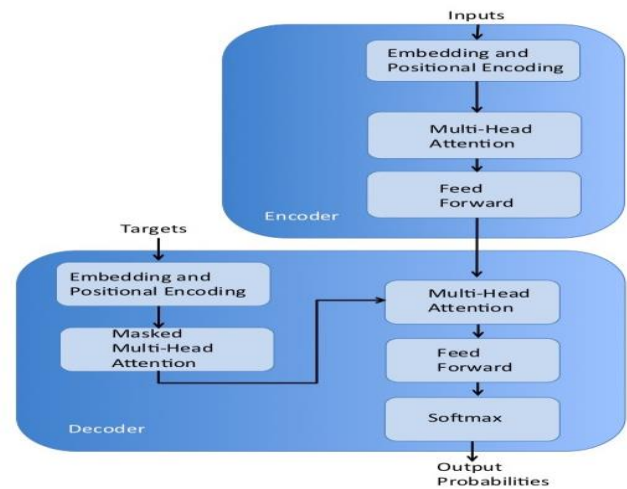


**Fig. 2.** Transformer

Finally, a function-clever feed-ahead community is used, which applies linear alterations to every function separately.

The decoder has 5 stages, the primary best used on the education phase: an output embedding and positional encoding (much like the only used withinside the encoder however for goal sentences withinside the education phase), a masked multi-head interest (additionally Self-Attention), a multi-head interest, a feed ahead community and sooner or later a softmax layer to compute the output probabilities.Given that on the interpreting method we can't realize the destiny words, the eye can best be carried out to preceding ones. This is what the masked multi-head interest does, that's a multi-head interest block with a masks that restricts the eye best to beyond words. For a deeper technical clarification of the structure see [9].

## 6.  Experimental Framework

This section reports the data, preprocessing and parameters that we used to build our chatbot systems.

### 6.1 Data and Preprocessing

Models had been examined at the OpenSubtitles dataset The Open Subtitles Corpus consists via way of means of a huge variety of film and TV collection scripts translated to a couple of languages. It is typically utilized by video structures to expose subtitles in their movies/TV collection.The subtitles do now no longer comprise identification nor flip information. Therefore, in addition to [10], we assumed that consecutive sentences had been uttered via way of means of one-of-a-kind characters. We built a dataset which include pairs of consecutive utterances, the use of each sentence two times as context and as target. Due to computing and reminiscence constrains, we extracted a subset of the primary 10 million sentences for schooling.

Preprocessing of the database consisted on doing away with XML tags, proscribing the sentence length and doing away with peculiar symbols (i.e., #, ", *, -, musical notes, etc.). For assessment we used the identical 2 hundred sentences that had been used in [10], which covers one-of-a-kind varieties of conversation (i.e. basic, philosophical, persona and popular knowledge). Details on schooling and assessment cut up are said on Table 1.

## 6.2 Parameters

At the experiments in [10], the architectures had 4096 unit cells for the complete OpenSubtitles database. Due to computational limitations, our version needed to be less complicated each via way of means of proscribing the database (as we explained) and additionally via way of means of the use of a layered LSTM version with 512 unit cells consistent with layer. Additionally, the version with interest makes use of a sampled softmax loss feature with 512 samples. All 3 fashions have a sixty four dense length for the embedding matrix. To make sure that we cowl the 99% of the dataset, we've got constrained the vocabulary length to 72,827 phrases and the duration of sentences to 24 phrases. All phrases which can be used best as soon as are discarded.

| Set | Role | Segments | Words | Vocab |
|---|---|---|---|---|
| Training | Context | 20,000,000 | 64,192,197 | 180,368 |
| | Target | | 48,174,044 | 182,404 |
| Evaluation | Context | 200 | 1,446 | 399 |

For training, we used ADAM [5] with a studying charge of 0.002, an exponential decay charge for the primary second estimates ( 1) identical to 0.nine, an exponential decay charge for the second one second estimates ( 2) identical to 0.999 and = 10 eight (offset to save you any department through zero); a batch length of 256 and a dropout charge of 0.1.The transformer structure has eight interest heads and six hidden layers with 512 units. For training, we've got extensively utilized ADAM with a studying charge of 0.2, exponential decay fees for the primary and 2d second estimates 1 = 0.nine and 2 = 0.98. and = 10 nine; a batch length of 4096 and a drop out and interest dropout fees of 0.1. For higher responses, we've got brought to the 3 fashions a beam seek set of rules with a beam length of four and a penalty period component of 0.6.

## 7. Evaluation

There isn't anyt any automated degree hooked up for chatbots. Therefore, we done the evaluationmanually. In order to evaluate our 3 conversational fashions, we requested 3 evaluators to rank chatbot responses among: proper, horrific or impartial. Additionally, we annotated

whilst the 3 evaluators agreed in rankings.Results are proven in Table 2 displaying rating from Evaluator 1, Evaluator 2, Evaluator three and settlement among them. Results display clean development of each interest mechanisms over the baseline machine. In fact, all evaluators agree that the transformer structure is ideal or impartial in 34.5% of the questions while the baseline structure is simplest proper or impartial in 23.5% of the questions and the baseline plus interest is ideal or impartial in 24.5% of the questions.

In addition to the results, we document numerous examples on one of a kind sorts of conversations gift withinside the check set: basic, philosophical, character and fashionable knowledge. These examples display the sort of development that interest is succesful to convey to the machine which in fashionable consist of greater good enough responses withinside the context of the person input.

Conversation 1: Basic
Context       what is your name ?
RNN Anouk.
+Attention    Adriano faustino sotolongo.
Transformer i don't know .
Context       what time is it ?
RNN I do n't know.
+Attention    6:30.
Transformer i don't know.
Conversation 2: Philosophical
Context       I think you are a machine
RNN I'm sorry
+Attention    Vlanne:
Transformer I'm not a machine .
Context       Life is sad .
RNN I'm not a monster.
+Attention    I do n't want to go to jail.
Transformer It's not sad .
Context       Is there a god ?
RNN No.
+Attention    Surprisingly, yes.
Transformer no .
Conversation 3: Personality
Context       How old were you when you were able to speak ?
RNN Eighteen.
+Attention    Seventeen.
Transformer six months.
Conversation 4: General knowledge
Context       ronald reagan is ...
RNN I'm not a good person.
+Attention    And i'm not a spy..
Transformer he's in the hospital.
Context       Is five plus five equal to ten ?
RNN No, no.
+Attention    That's exactly what i'm gon na do.
Transformer yes .

Table 2. Manual Evaluation. Best results in bold

| System | Ranking | Evaluator 1 | Evaluator 2 | Evaluator 3 | Agreement |
|--------|---------|-------------|-------------|-------------|-----------|
| RNN | Good | 41 | 47 | 54 | 32 |
| +Attention | | 61 | 51 | 71 | 43 |
| Transformer | | 74 | 57 | 70 | 51 |
| RNN | Bad | 75 | 123 | 53 | 46 |
| +Attention | | 90 | 116 | 57 | 46 |
| Transformer | | 25 | 110 | 25 | 20 |
| RNN | Neutral | 84 | 30 | 93 | 15 |
| +Attention | | 49 | 33 | 72 | 6 |
| Transformer | | 101 | 33 | 105 | 18 |

## 8.   Conclusions

ttention-primarily based totally mechanisms are revolutionizing herbal language, speech and photograph processing applications. In this paper, we're imposing multiple lately proposed interest mechanisms into the chatbot application.

Experiments educated on an open-area data-base display that a totally interest-primarily based totally structure plays appreciably higher in quite a few contexts together with basic, philosophical, persona and trendy knowledge. Three evaluators agreed on score the responses of the completely interest-primarily based totally mechanism 34.5% of the time both accurate or neutral, whilst the responses of the baseline structure with interest changed into rated in that phrases on a 24.5% of the time and the responses of the baseline device had been handiest 23.5% of the time both accurate or neutral. Taking gain of this typical encoder-decoder structure, amongst similarly research, we plan to introduce similarly contexts whilst schooling the device that allows you to permit the device to maintain coherence in longer dialogues and to teach our device on a couple of languages.

## REFERENCES

[1]  Banchs, R. E. & Li, H. (2012). IRIS: a chat-oriented dialogue system based on the vector space model. The 50th Annual Meeting of the Association for Computational Linguistics, Proceedings of the System Demonstrations, Jeju Island, Korea,37–42.

[2]  Cho, K., van Merrienboer, B., Gulc¨¸ehre, C¸., Bahdanau, D., Bougares, F., Schwenk, H., & Bengio, Y. (2014). Learning phrase representations using RNN encoder-decoder for statistical machine translation. Proceedings of the 2014

[3]  Conference on Empirical Methods in Natural Language Processing, EMNLP, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL, pp. 1724–1734.

[4]  Hochreiter, S. & Schmidhuber, J. (1997). Long short-term memory. Neural Comput., Vol. 9, No. 8,1735–1780.

[5]  Kingma, D. P. & Ba, J. (2014). Adam: A method for stochastic optimization. CoRR, Vol. abs/1412.6980

[6]  Experimental Research on Encoder-Decoder Architectures with Attention for Chatbots 1239

[7]  Mikolov, T., Karafiat,´ M., Burget, L., Cernocky,´ J., & Khudanpur, S. (2010). Recurrent neural network based language model. INTERSPEECH 2010, 11th Annual Conference of the International Speech Communication Association, Makuhari, Chiba, Japan, pp. 1045–1048.

[8]  Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Se-quence to sequence learning with neural networks. Advances in Neural Information Processing Sys-tems 27: Annual Conference on Neural Information Processing Systems, Montreal, Quebec, Canada,3104–3112.

[9]  Tiedemann, J. (2009). News from OPUS: A collection of multilingual parallel corpora with tools and interfaces. In Recent Advances in Natural Language Processing, volume V. John Benjamins,237–248.

[10]  Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit,J., Jones,L., Gomez, A. N., Kaiser, L.,& Polosukhin, I. (2017). Attention is all you

[11]  need. Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems, Long Beach, CA, USA, pp. 6000–6010.

[12]  Vinyals, O. & Le, Q. V. (2015). A neural conversational model. CoRR, Vol. abs/1506.05869.

[13]  Wallace, R. (2003). The elements of aiml style.Weiss, R. J., Chorowski, J., Jaitly, N., Wu, Y., & Chen, Z. (2017). Sequence-to-sequence models can directly transcribe foreign speech. CoRR, Vol. abs/1703.08581.

[14]  Weizenbaum, J. (1966). ELIZA: a computer program for the study of natural language communi-cation between man and machine. Commun. ACM, Vol. 9, No. 1, pp. 36–45.