



# LSTM Based Ensemble Network to Enhance the Learning of Long-term Dependencies in Chatbot

Shruti Patil<sup>1\*</sup>, Venkatesh Mudaliar<sup>2</sup>, Pooja Kamat<sup>3</sup>

<sup>1,3</sup> Assistant Professor, Symbiosis Institute of Technology, Pune

<sup>2</sup> Mtech Research Scholar, Symbiosis Institute of Technology

Received: 9<sup>th</sup> February 2020

Accepted: 17<sup>st</sup> May 2020

OPEN ACCESS 

**Abstract:** A chatbot is a software that can reproduce a discussion portraying a specific dimension of articulation among people and machines utilizing Natural Human Language. With the advent of AI, chatbots have developed from being minor guideline based models to progressively modern models. A striking highlight of the current chatbot frameworks is their capacity to maintain and support explicit highlights and settings of the discussions empowering them to have a human contact through the span of involvement. The paper expects to build up a detailed database with respect to the models utilized to deal with the learning of long haul conditions in a chatbot. The paper proposes a crossbreed Long Short Term Memory based Ensemble Network arrangement model to save the continuation of the specific situation. The proposed model uses a characterized number of Long Short Term Memory Networks as a major aspect of the amassed model working as one to create the aggregate forecast class for the info inquiry handled.

**Keywords:** Chatbot, AI, LSTM, Ensemble Method, GRU

## Introduction

A Conversational Agent is otherwise called 'Chatbot' is a software program which leads a discussion by means of sound-related or literary strategies in a characteristic language, for example, English. Chatbots are being coordinated universally into our lives in a type of Virtual collaborators and messaging applications. In 1950, Alan Turing proposed 'Turing Test' as a benchmark of a chatbot program to imitate a human in a discussion [3]. ELIZA, Jaberwacky, A.L.I.C.E. were not many of the underlying chatbots created dependent on principle based methodology [1]. The 'Measurable Revolution' contingent upon Machine Learning blossomed in the late 1980s and mid-1990s. There has been a significant powerful movement in the zone of chatbot inferable from the development of man-made reasoning. The presentation of the influx of Artificial Intelligence-based chatbots has introduced another time of conversational interfaces. The other factor adding to advancement is the noticeable

change of the elements of human discussion leaning toward short informing over different types of correspondence. Most chatbots are gotten to through remote helpers, informing applications or association's sites. Right now, the market of cutting edge conversational specialists is shared by IBM's Watson, Apple's Siri, Google Assistant, Amazon Alexa, Microsoft's Cortana to give some examples. Endeavours have been made to typify the usefulness of chatbot consistently into administrations alongside contracting the uniqueness contrasted with human discussions. The incorporation of an inductive memory practically equivalent to the human cerebrum into the engineering of a chatbot encourages the chatbot to keep up the edge of setting for longer durations. Protecting the highlights identified with the relationship for longer lengths is named as adapting Long term dependencies. This model acquires a factor of commonality and lucidness over the span of the discussion between the end client and chatbot. Wide and unambiguous data about the advancement of memory incorporated models utilized in chatbot could give an

exhaustive comprehension and bits of knowledge into the eventual fate of chatbot inquire about. The design and advancement strategy of a run of the mill chatbot relies upon fundamental ideas as determined in Figure 1 [2].

A brief look at the concepts to comprehend the variations possible at the formative stage is presented below.

1. Text Processing: Word embeddings are the vector representations of words within the specific vocabulary enabling better implementation and utilization of statistical machine learning models.
2. Machine Learning Model: The concept of Artificial Neural Network is extensively employed in dealing with input processing, classification and generating the most appropriate response for the input query. Chatbots like Deepprobe, Superagent utilize the Long short-term memory (LSTM) model with Seq2Seq, while Rubystar uses Seq2Seq with Gated recurrent unit (GRU) [2].
3. Knowledge Base: The dataset used for training the model can be either Open or Close in its bounds. Open domain chatbot was found to be compromised on relevance and accuracy of the responses and Closed domain chatbot performs well owing to limited yet definite confines of the dataset [2].
4. Response Generation: The response returned for input is either retrieved or generated. The former selects the appropriate response from a collection while the latter generates the response depending on the features of input vectors, dictionary a trained classifier. The hybrid RNN-Seq2Seq model has progressed to become a popular choice in chatbot architecture [3].

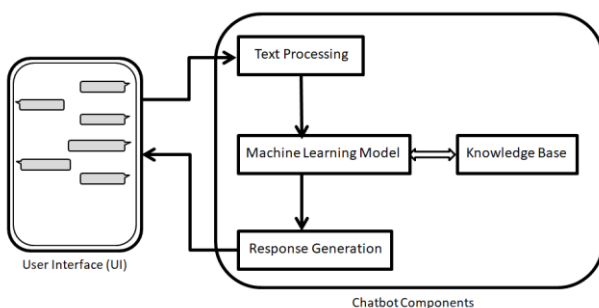


Figure 1. Chatbot Operation

## Literature Survey

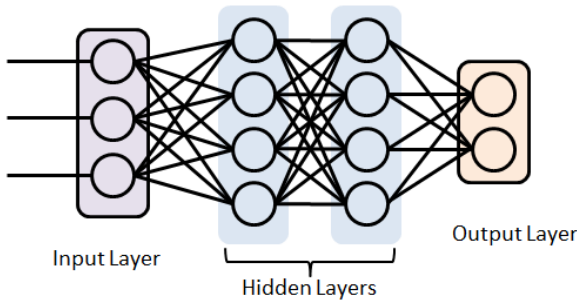
This section aims to provide an overview on different concepts employed in handling long-term dependencies and discuss their corresponding nuances. The timeline of some fundamental technologies is listed in Table 1.

Table 1. Course of development

Year	Author	Contribution
1943	W. McCulloch, W. Pitts	Artificial Neural Network (ANN)
1990	Elman	Simple Recurrent Neural Network (RNN)
1990	L. K. Hansen, P. Salamon	Ensemble Learning
1994	Y. Bengio	Issue with long term dependencies
1997	S. Hochreiter	Long Short Term Memory (LSTM)
2000	F.A. Gers	LSTM with forget gates
2014	K. Cho	Gated Recurrent Unit (GRU)
2014	A.Graves	Neural Turing Machine (NTM)

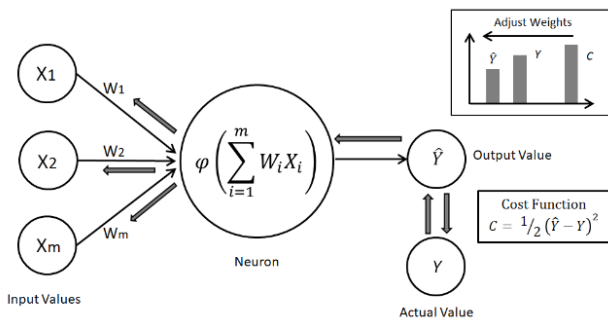
## Artificial Neural Networks

Artificial Neural Networks (ANNs) are information processing models inspired by the biological neural system and the capability of the brain to process information. The work on Neuron Circuit and Perceptron by Warren Pitts, Warren McCulloch and F. Rosenblatt served groundwork for ANN to evolve and induct over the traditional computer frameworks in the 1970s [1]. ANN is composed of a large number of densely interconnected mathematical function units called 'Neurons' clustered into three types of layers as shown in Figure 2. The input layer is responsible for the initial processing of input data whereas the output layer deals with aggregating the final outputs and presenting the result. The weighted connections between neurons in hidden layer form the basis of learning process providing variable strength to the input data traversing forward towards output neurons. An activation function like Sigmoid, ReLU or tanh is applied on the summation of weighted inputs in a neuron [4].



**Figure 2.** Artificial Neural Network structure

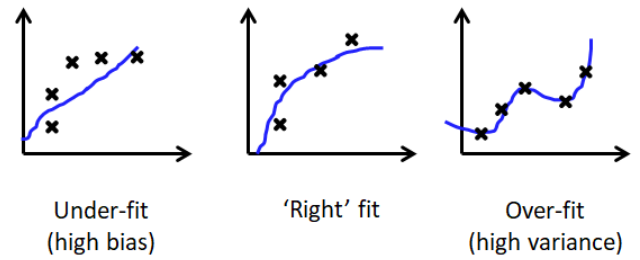
The model is trained using 'Back Propagation' where the error calculated leads to optimal updation of weights. Gradient ( $\Delta w$ ) can be calculated as change in error with respect to change in weights ( $de/dw$ ). Values for new weights is determined by adding weight ( $w$ ) and the gradient ( $\Delta w$ ). The entire process is depicted in Figure 3.



**Figure 3.** Backpropagation process

However, the brute force approach for updating weights suffers from 'Curse of dimensionality' [5]. Gradient Descent (GD) and Stochastic Gradient (SGD) descent offer a faster way to find optimum weights. Both these methods determine the global minima by finding the point where the slope of the cost function is zero hence resulting the error to be minimum. GD and SGD are compared in Table 2.

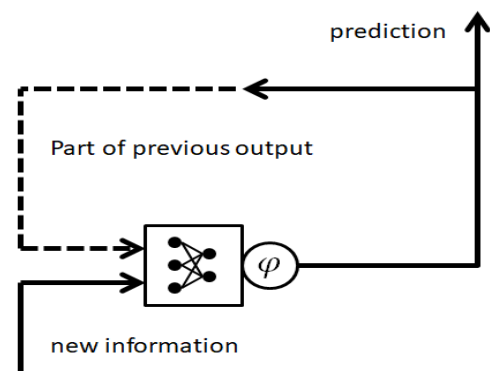
ANN suffers from both overfitting and underfitting as described in Figure 4. Overfitting is an outcome of an overly accurate or complicated model showing low bias but high variance. Underfitting is a result of a too simple model showing low variance but high bias [6]. ANNs deal with fixed sized vectors only and they do not possess a dedicated memory element to handle sequential data hence making them an inappropriate choice for a chatbot handling dependent vectors.



**Figure 4.** Overfitting and Underfitting of a model

**RNN**

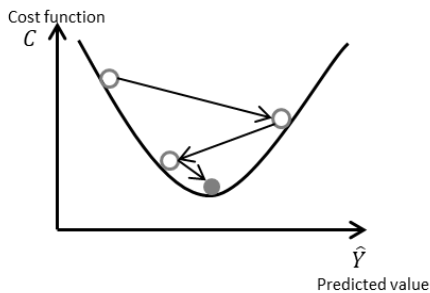
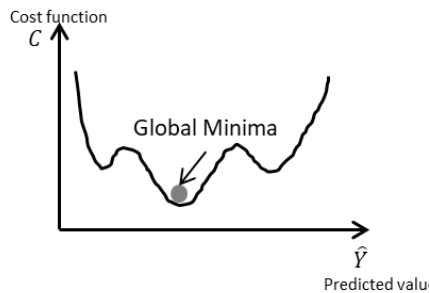
Recurrent Neural Network (RNN) is the class of Artificial Neural Network supplemented by the integration of edges spanning adjacent timestamps. Psychologist David Rumelhart's work on symbolic artificial intelligence from 1986 formed base for the development of RNN. RNN has two inputs, the present values and values from recent past enabling it to capture the dynamics of a sequence of inputs in scenarios like handwriting recognition, stock price prediction, etc. Owing to the variable size of input and output vectors RNN has shown significant improvement over traditional feed forward networks in Chatbots as RNNs are capable of exploiting a dynamically changing contextual window over input sequences. Overall architecture of RNN is specified in Figure 5.



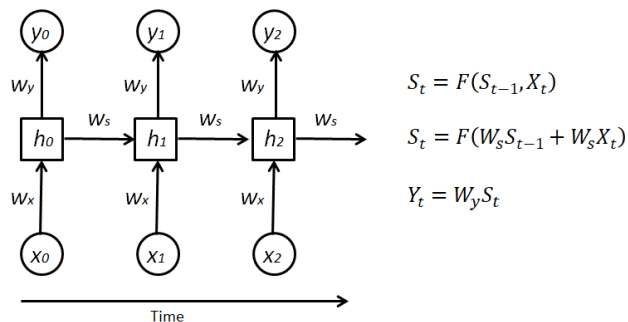
**Figure 5.** RNN overall architecture

At given time  $t$ , output for state  $S_t$  is calculated applying function on portion of output from previous state  $S_{t-1}$  and current input  $X_t$ . It can be termed mathematically as  $S_t = F(S_{t-1}, X_t)$  where  $F$  is activation function like tanh or ReLU. This process continues forming an information loop for a given state with respect to time. The unrolled structure of RNN is shown in Figure 6 along with equations.

**Table 2.** Comparison between GD and SGD

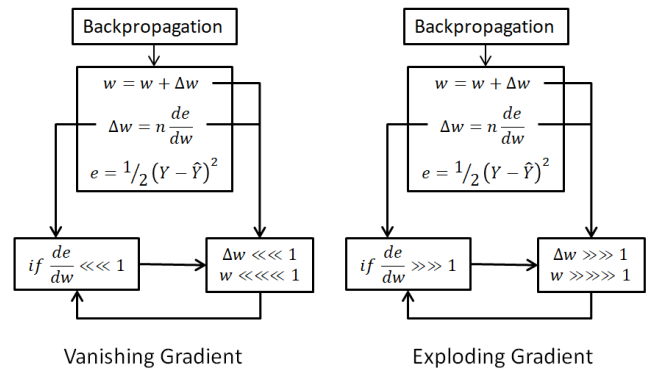
Gradient Descent	Stochastic Gradient Descent
	
GD computes gradient using a batch from dataset	SGD computes gradients using single rows of training examples
It follows a deterministic approach	It follows a random approach
It converges slower on large training samples	It converges faster on large training samples
<b>Steps:</b> For every iteration 1. Traverse entire dataset 2. Evaluate gradient 3. Return	<b>Steps:</b> For every iteration 1. Iterate over each value in dataset 2. Evaluate Gradient 3. Return

Like Feed-forward networks, RNNs use back propagation for training the difference being the additional parameter 'time', hence it is termed as 'Back propagation through time (BTT)' as shown in Figure 6[7].



**Figure 6.** RNN unrolled structure

The range of context to be used practically is limited as each prediction looks at one step prior state value. While back propagating the recurrent connections, the influence of given input vector on the corresponding hidden layer and hence overall network output either decays or blows up exponentially giving rise to Vanishing Gradient and Exploding Gradient problem respectively as shown in Figure 7. Both these problems cause the model to train poorly and performance degradation.



**Figure 7.** Vanishing Gradient and Exploding Gradient

A prediction of a state at the time 't' depends on the input presented at earlier time T where T << t. When the gap between T and t grows large, it becomes extremely difficult for the model to attain convergence causing the failure of RNN to handle 'Long Term Dependencies' which makes it unfitting model for chatbots dealing with time series conversations [8].

*LSTM*

Long Short Term Memory networks are an extension for Recurrent Neural Networks with explicitly extended memory capability well suited to handle long term dependencies [9]. LSTM networks were proposed by German researchers Sepp Hochreiter and Juergen Schmidhuber in 1997 as a solution to the vanishing

gradient problem [10]. In comparison, LSTM can learn to bridge the features in excess about 1000 definite time steps by imposing constant error flow through the units termed as 'cells' effectively dealing with Long Term dependencies [11].

LSTM contain information from a context in a gated cell. The cells control the data to be written, stored, read and erased using Forget, Input and Output gates which are implemented with element-wise multiplications by sigmoids as shown in Figure 8 [7]. The forget gate learns the weights controlling the decay rate of values stored in memory cells. For the instance when the input and output gates are off and the forget gate is not causing decay, the memory cell maintains its value over time causing the gradient of error to stay constant during backpropagation. This enables the model to remember information for longer periods. The overall architecture of LSTM is shown in Figure 8.

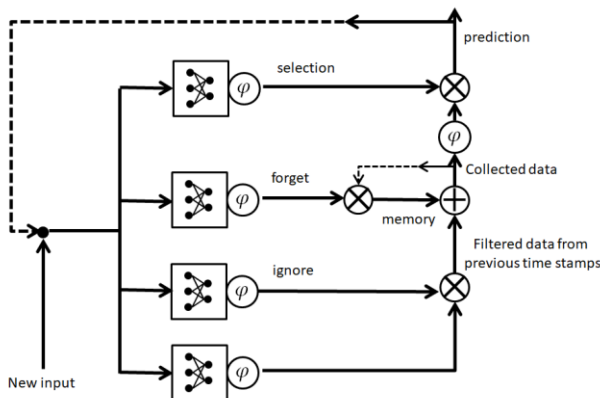


Figure 8. LSTM network

Mathematically each step can be explained as follows:

1. In the first step Forget Gate layer decides the features to be flushed out from cell state looking at  $h_{t-1}$  and new input  $x_t$ .

$$f_t = \sigma(w_f[h_{t-1}, x_t] + b_f)$$

2. In the second step, deciding the information to be stored in the cell state is done in two steps. Input Gate layer  $i_t$  which is a sigmoid layer establishes the values to be updated. Then a  $\tanh$  layer generates the vector of new candidate values  $\tilde{C}_t$ .

$$i_t = \sigma(w_i[h_{t-1}, x_t] + b_i)$$

$$c_t = \tanh(w_c[h_{t-1}, x_t] + b_c)$$

3. The old cell state  $C_{t-1}$  is updated to new cell  $C_t$

summing the output from Forget gate layer function  $f_t$  and  $i_t * \tilde{C}_t$ .

$$c_t = f_t * c_{t-1} + i_t * c_t$$

4. The output is determined in two steps - First, the sigmoid layer decides the parts of cells to output  $o_t$ . The product of new cell state  $C_t$  through  $\tanh$  and the output of sigmoid gate outputs  $h_t$  the selectively decided parts.

$$o_t = \sigma(w_o[h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(c_t)$$

Hyperparameters tuning and optimization is an arduous and experimental task [4]. The training of the LSTM model is expensive in terms of memory and computational power. In the domain of chatbots for time series conversations, LSTM is shown to perform well and maintain the context for longer durations.

### GRU

A Gated recurrent unit (GRU) is a specific model of Recurrent Neural Network introduced by Kyunghyun Cho in 2014 as a variation of an intermediate unit like LSTM enabling the recurrent unit to capture dependencies of different time steps.

Unlike LSTMs, GRU has 2 gates as Reset and Update to control the flow of information and refine the outputs. When compared to LSTM, the update gate can be considered a combination of Forget and Input gate from LSTM. Update gate determines the portion of information from previous time steps needs to be passed to the next states. This gives GRU an edge over LSTM as the model can decide to maintain all features from previous timestamps. Reset gate is used to decide the irrelevant part of the information which needs to be discarded [12]. GRU works in the following steps:

1. Update gate  $[z_t]$  at time 't' is calculated.

$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

2. Reset gate  $[r_t]$  calculates the information to be forgotten using

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

3. New memory content is introduced which uses the reset gate to store the relevant information

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

4.  $[h_t]$  is calculated which holds information for the

current unit using update gate output and memory content from previous steps [ht-1].

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

GRU exposes complete memory content without control gate when compared to controlled exposure of LSTM using Output gate. GRU explicitly controls the influx of information while calculating new memory content using the Update gate. Owing to less complex nature and few tensor operations, GRU is computationally more effective and faster to train.

### NTM

Neural Turing Machine (NTM) explores the concept of evidently extending the context accumulator of RNN with an addressable external memory. They are an example of Memory Augmented Neural Networks which decouple the computation from memory [13]. NTM have been shown to outperform LSTMs on sequence learning tasks demanding large memory for handling memorization of longer contexts.

Controller and Memory matrix are primary components in NTM as shown in Figure 9. The controller is a recurrent or feed forward neural network which takes input and returns the output. External memory unit constitutes of N\*W memory matrix where N is the number of memory locations and W is the dimension of each memory cell. The interaction between the Controller and Memory matrix is carried out by reading and write heads. The memory matrix is initialized using schemes like Constant initialization or Truncated Normal distribution [43]. The NTM model can be trained by variants of Stochastic gradients using back propagation through time in case of an RNN based controller.

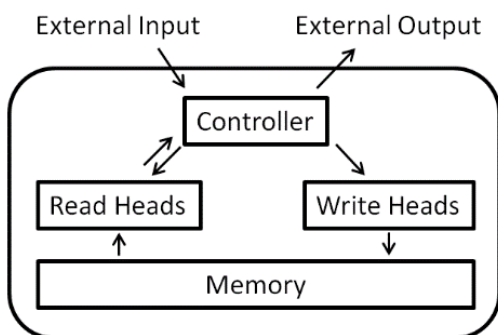


Figure 9. NTM architecture

Algorithmic tasks like priority sort, Associative Recall, Copy, Repeat Copy, etc can be performed to test if the

NTM could be trained via supervised learning for efficient performance. NTM models generalize reasonably well to longer inputs

### Ensemble Learning

The concept of ensemble learning was popularized in 1990 by Lars Kai Hansen and Peter Salamon [15] over the idea that performance of a set of classifiers outweighs that of a single classifier. The individual models work in unison where the outputs are combined with a certain decision fusion strategy to output a single answer [14]. Owing to the combination of various learning models, the generalization ability turns to be stronger. The basic architecture of Ensemble model is depicted in Figure 10. The variation in the member models is a critical factor for classification performance [16], hence strategies as follows were proposed for boosting the diversity scale among the member learners:

1. Employing different learning algorithms for different learners or using the same algorithm with variation in parameters
2. Training the members with varied datasets by subsampling or changing the attributes.
3. Combination of the above two methods is used simultaneously.

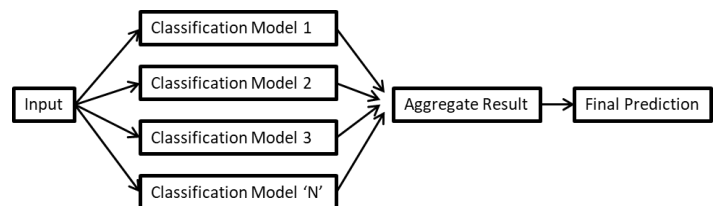


Figure 10. Ensemble Learning Model

An overall comparison between the concepts discussed along with the problem statements each individual methodology is well suited for is stated in Table 3.

### Dataset

#### Data Source

The dataset used in the entire research is Cornell Movie Dialog Corpus. Its distributed by Cornell Edu. The dataset consists of difference metadata-rich files. The conversations in the dataset are extracted from movie scripts. The dataset in whole has 220579 exchanges between 10292 characters collected from 617 movies.

Two files are used for establishing the conversation data. 'Movie\_lines.txt' contains texts from the dialogues and it has attributes like lineID, CharacterId, movieID, character name and the actual text. 'movie\_conversations.txt' forms the structure of the conversation. It maps the

conversation between two characterIDs together along with the movieID of the movie. The " +++\$+++ " acts as the field separator between the attributes mentioned for each file utilized.

**Table 3.** Overall comparison of Concepts from Literature Survey

Model	Advantages	Suited problem statements
ANN	<ul style="list-style-type: none"> <li>- Self-organizing to changes in information</li> <li>- Fault tolerant to the corruption of cells and missing input values</li> </ul>	ANN is well suited for classification and regression dealing with a large number of variables. Character recognition, Image processing are some of the applications for ANN
RNN	<ul style="list-style-type: none"> <li>- Adapts wells to quick changes in the input nodes</li> <li>- Variable size of input and output vectors</li> <li>- Works well with contextual input sequences</li> <li>-Excel at modeling temporal structure</li> </ul>	Appropriate for sequence prediction, classification prediction, Natural language processing and generative model. Hence they can be used in text generation, prediction of the values of an attribute in a problem statement.
Ensemble Learning	<ul style="list-style-type: none"> <li>- Better generalization ability</li> <li>- Weak models can be boosted to efficient learners</li> <li>- By the reason of growing computations power, the Ensemble model can be well utilized</li> </ul>	It can be used to enhance the performance of existing models like RNN, LSTM and GRU.
LSTM	<ul style="list-style-type: none"> <li>- Extended memory capabilities than RNN</li> <li>- Handles Long-term dependencies well.</li> <li>- More robust to vanishing gradients than RNN</li> </ul>	Good choice for problem statements like Time series forecasting. LSTM can be applied in Conversation agent, handwriting generation, Language translation, Image captioning.
GRU	<ul style="list-style-type: none"> <li>- Handles Long-term dependencies effectively</li> <li>- Robust to Vanishing gradient problem</li> <li>- Computationally effective than LSTM</li> </ul>	GRU can be used in applications related to time series prediction like text generations, classification, etc
NTM	<ul style="list-style-type: none"> <li>- Generalize well to longer inputs as compared to LSTM</li> <li>- Presence of external memory complements the RNNs existing memory</li> </ul>	NTM is well suited for models with heavy and longer sequences of data. NTM has demonstrated the solutions to be generalizing well for basic algorithms like copying and sorting.

The authors now present a comprehensive review of some the recent works carried out in this domain:

**Table 4:** Comprehensive Review of recent works involving AI for chatbot implementation

Sr No.	Research Paper	Year	Algorithm Used	Research Findings
1.	M. Nuruzzaman and O. K. Hussain, "A Survey on Chatbot Implementation in Customer Service Industry through Deep Neural Networks," [22]	2018	ANN	Artificial Neural Network (ANN) owing to its capability to handle the complicated combination of features provides the most appropriate base to work upon for a problem statement such as Conversational Agents or Chatbots.
2.	Lee, M. C., Chiang, S. Y., Yeh, S. C., & Wen, T. F. "Study on emotion recognition and companion Chatbot using deep neural network" [23]	2020	RNN	RNN provides a better response to problem statements about Seq2Seq framework of RNN built over Domain-Specific Knowledgebase.
3.	Bali, M., Mohanty, S., Chatterjee, S., Sarma, M., & Puravankara, R. Diabot: "A Predictive Medical	2019	Ensemble Learning	Ensemble Learning as a meta-algorithm has the potential to provide better generalization. The increased

	Chatbot using Ensemble Learning” [24]			performance can be mapped with strong correlations with a humane sense of conversation
4.	Pathak, K., & Arya, A. “A Metaphorical Study Of Variants Of Recurrent Neural Network Models For A Context Learning Chatbot” [25]	2019	LSTM	LSTM is the most appropriate choice when the states of dialogues & responses in a conversation need to be tracked and predicted.
5.	G. Dzakwan and A. Purwarianti, "Comparative Study of Topology and Feature Variants for Non-Task-Oriented Chatbot using Sequence to Sequence Learning," [26]	2018	GRU	The Encoder-Decoder framework over GRU comes in as an alternative to LSTM for the Chatbot. GRU offers similar if not better performance over LSTM when compared over the parameters like the size of the dataset, the resources being consumed.

The " +++\$+++ " acts as the field separator between the attributes mentioned for each file utilized.

### Data Pre-processing

Python with numpy is used to preprocess the dataset to institute the conversation dictionaries. Dictionaries are created to map each line and the corresponding 'id', creating a list of all conversations, separating questions and answers. Individual conjoint words are cleaned and replaced with simple words. A dictionary is also created to map each word with its number of occurrences and for mapping the questions words and answer words to unique integer values.

### Proposed Method

Ensemble learning forms the basis of the proposed methodology. Classifiers like Support vector machines, Linear regression were used in the Ensemble model initially. With the onset of Deep learning, a more elaborate approach can be followed to improve the overall performance of the Ensemble model. The idea is to define a number of LSTM networks with variation in hyperparameters as part of the ensemble model. The member models work together in parallel and their individual outputs are aggregated to generate the output of the overall model. As a fine-tuning measure, the concept of Pruning is also employed. An architectural overview is presented in Figure 10 followed by detailed Ensemble Network algorithm. Segmentation, Vector Space Model (VSM), Classification algorithm & Response generation forms the primary components of the chatbot. The flow of operations is shown in Figure 11. After the output class is predicted, the output of Chabot is returned to the user.

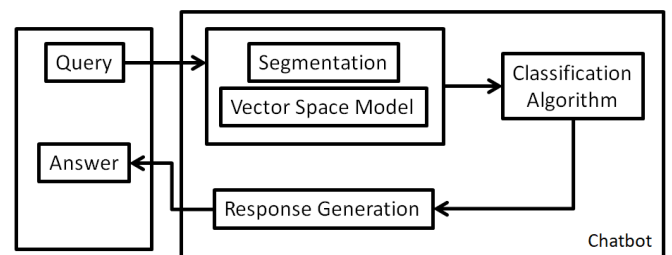


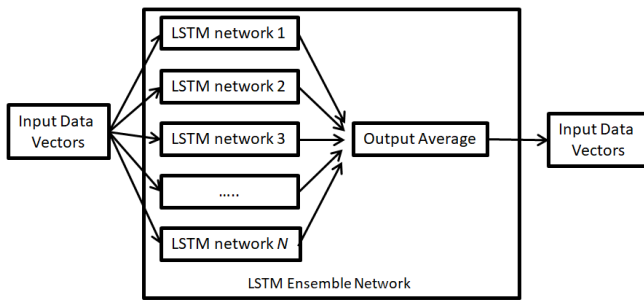
Figure 11. Proposed chatbot architecture

Components in the LSTM based Ensemble network as described in Figure 12.

- 1. Input Data:** Input sentences are segmented into terms. These terms are transformed into vectors by VSM corresponding to Vector space. The vectors generated are fed into the classifier model.
- 2. STMs based Ensemble Network:** The ensemble network consists of a defined number of LSTM networks working concurrently to generate the overall output prediction. The variation in hyperparameters like number of hidden layers, number of neurons forms the basis of the distinction between each LSTM model. This leads to the training of models with different generalization features and accuracy metrics.



- Prediction:** The LSTM models in the ensemble generate the predicted output which is converted to predicted class.



**Figure 12.** LSTM based Ensemble Network Classifier

### The Proposed Methodology

The Encoder Decoder LSTM acts the base for defining the even single LSTM as well as the combination of LSTMs acting in unison as part of the ensemble. The entire process of implementation can be broken down in certain steps as discussed below. The proposed model includes training phase and the testing phase which are shown in Figure 12 and Figure 13 respectively. In the training phase, a definite number of LSTM networks are generated and trained using variations in training data. The models with lower accuracy are filtered out. In the testing phase, the models with higher accuracy work in conjunction to predict the output class from the calculated output weights. The detail of the two phases is presented in the following sections. The notations used in the algorithm specifications are specified in Table 4.

**Table 5.** Notations used in Algorithm

Notation	Description
$N$	The number of LSTM networks
$h_{max}$	Maximum number of hidden layers in a LSTM model
$p_{max}$	Maximum number of neurons in a hidden layer
$w_{threshold}$	Threshold value for the accuracy
$h_m$	Number of hidden layers in LSTM model
$p_m$	Number of neurons in a hidden layer
$T$	Total number of partitions of training dataset
$A_m$	Average Accuracy
$k$	LSTM models remaining after pruning
$C_k$	Output of $k^{th}$ LSTM model

### Data Pre-processing

The dataset is preprocessed as defined in the section NUMBER. Different functions are created for the generation of dictionaries and cleaning the text.

### Building the Model

Encoder LSTM is responsible for reading the input sequence and encoding the same into a vector essentially to map the corresponding vector from the vector space defined. Decoder LSTM deals with decoding the vector generated and outputting the predicted sequence. Encoder Decoder LSTM generates a continual representation of data from a considerable number of data attributes from previous time stamps. This architecture of Encoder Decoder LSTM was found effective on long and continuous data influx. We split the dataset into training and validation dataset as an attempt to carry out cross-validation. Three difference decoder LSTMs are created in order to decode the training data, decode the validation data and the actual decoder for the encoder created.

### Training Phase

The hyperparameters like the number of epochs, batch size, LSTM size, number of layers in Encoder and Decoder LSTM, Learning rate are initialized for single LSTM as well as the Ensemble LSTM model. A session for training is initialized and the models are training for both portions of the dataset that ie the Training dataset and validation dataset as well. As the training progresses the model generates the weights. The model generalizes the data for patterns and features and stores it in the model to be utilized while testing. The training phase is depicted in Figure 13.

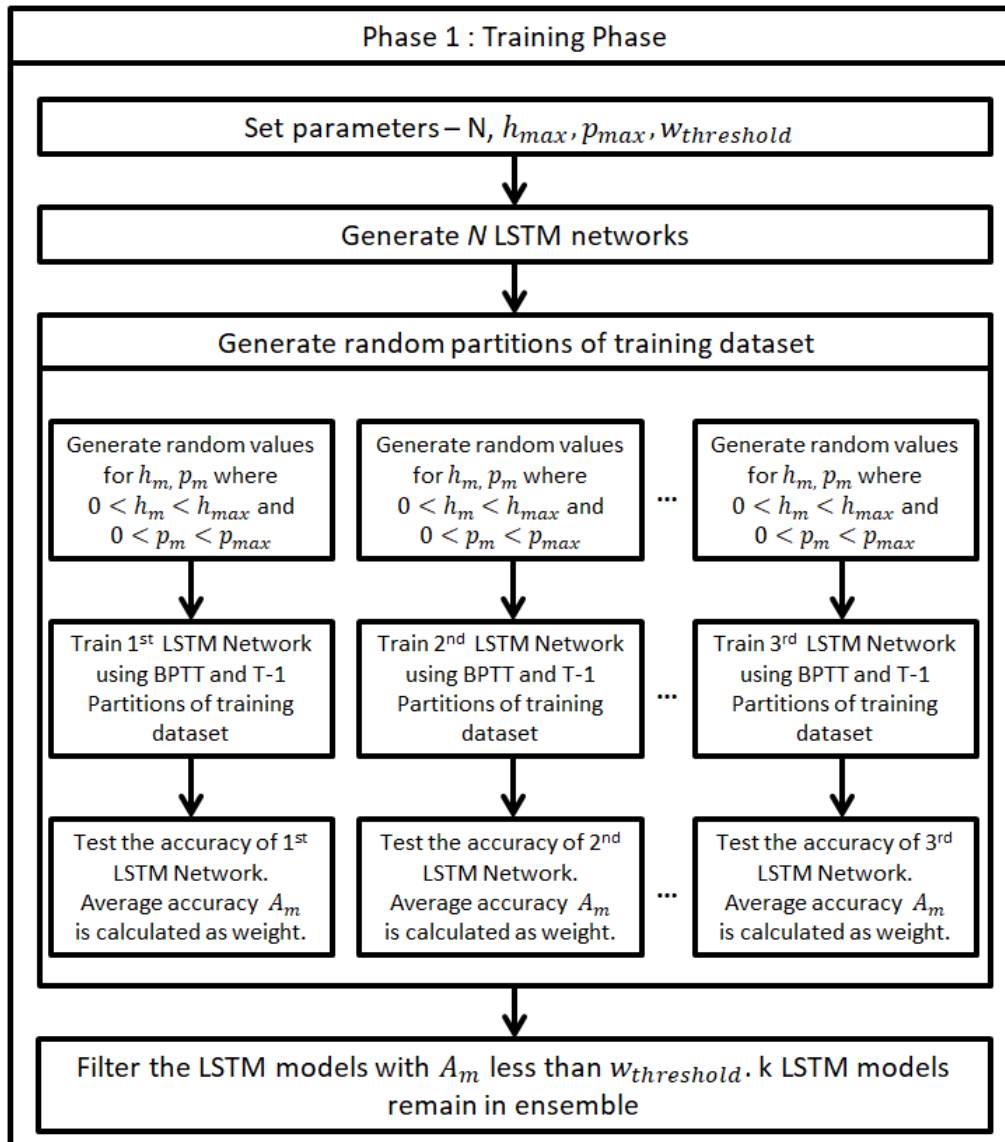
**Step 1:** The parameters like the number of LSTM networks  $N$ , the maximum number of hidden layers in the model  $h_{max}$ , a maximum number of neurons in a hidden layer  $p_{max}$  and accuracy threshold for each LSTM model  $w_{threshold}$  are initialized.

**Step 2:** Generate  $N$  LSTM models with variations in hyperparameters. The number of hidden layers  $h_m$  and the number of neurons in a hidden layer  $p_m$  are assigned random values between 0 and corresponding maximum bounds.

**Step 3:** The training dataset is split into  $T$  partitions.  $T - 1$  partitions are used to train each LSTM model. Every distinct LSTM model trains and learns using BPTT.

**Step 4:** The performance of each LSTM model along with the average accuracy  $A_m$  is evaluated using the remaining single partition from testing dataset.

**Step 5:** The LSTM models with accuracy lower than the accuracy threshold  $w_{threshold}$  are dropped leaving  $k$  models in the ensemble network.



**Figure 13.** Training Phase algorithm

### Testing Phase

The weights generated in training are loaded and a session is initiated running the Encoder Decoder LSTM model as a part of both the single model and Ensemble LSTM. The incoming queries are cleaned and preprocessed using the functions defined in Data Preprocessing. The predicted answer is returned to the user. This phase deals with applying the patterns and features learned on the testing dataset. The testing dataset is fed into the Ensemble model to predict the

output classes. The testing phase is depicted in Figure 14.

**Step 1:** The testing dataset is retrieved and fed into the LSTM models of the ensemble.

**Step 2:** The output  $C_k$  for the individual model is calculated.  $C_k$  is multiplied with corresponding weights of each model. A weighted average of the sum of total weighted outputs is calculated.

**Step 3:** The weighted average is used to predict the response class.

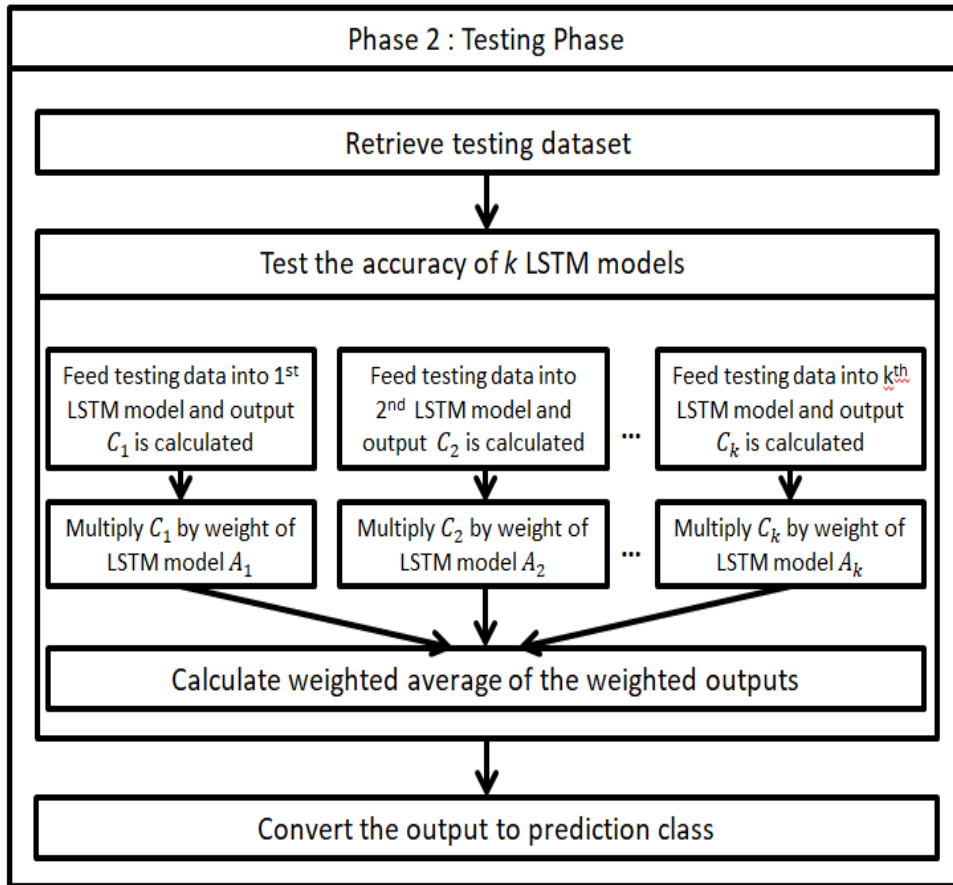


Figure 14. Testing Phase algorithm

An example of the variations in the hypermeters for multiple LSTM multiples as a part of the Ensemble Network is presented in Figure 15.

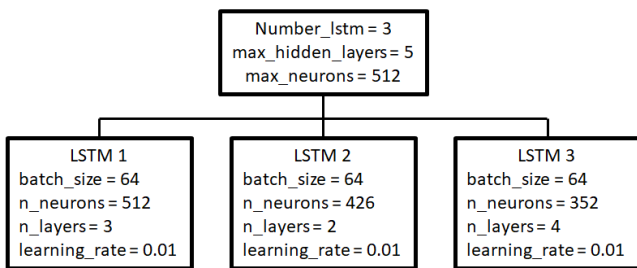


Figure 15. Hyper parameter variations in LSTM

### Performance Analysis

RNN, LSTM and GRU serve as the best choice for the classifier in Ensemble Network. LSTM and GRU provide an edge over RNN owing to the presence of a dedicated memory control unit enabling the learning of long term dependencies. The selection of an appropriate model between them depends on the key differences and the dataset. GRU exposes complete memory content without

control gate when compared to controlled exposure of LSTM using Output gate. LSTM doesn't control the amount of information flowing in from previous time steps while computing new memory content. On the other hand, GRU explicitly controls the influx of information while calculating new memory content using the Update gate. An experiment was performed to compare LSTM and GRU for their performance in the time series prediction.

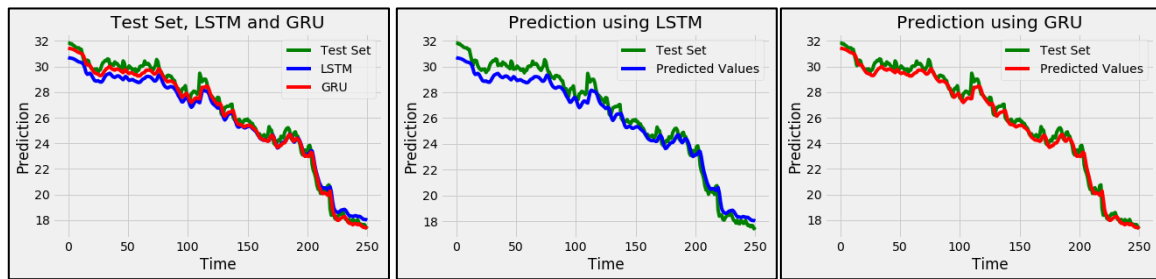
### Prediction Comparison

LSTM and GRU are closely related mechanisms for handling long term dependencies. A comparison between both for their performance provided important insight for the selection of LSTM over GRU as seen in Figure 16. Two models constituting single LSTM and a single GRU were created for the comparison. To make the comparison more just, both LSTM and GRU models had 4 hidden layers with 50 neurons each with 0.2 dropout rate. Both the models were trained with 50 epochs and a batch size of 32. The dataset used to select the appropriate model for time series data analysis is Stock Price Dataset. The Dataset is split into portions randomly to generate a more

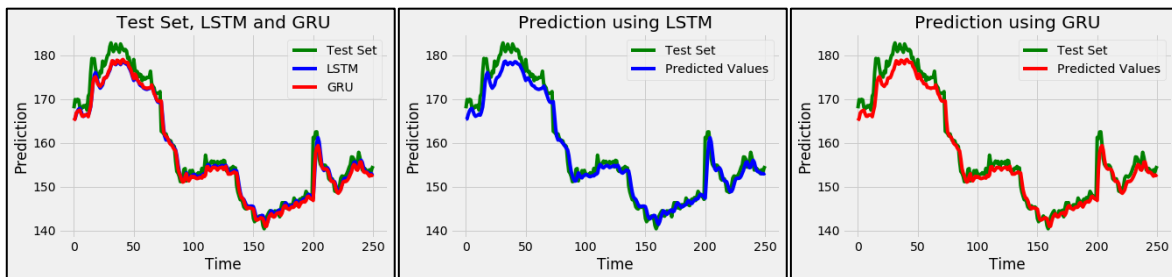
stable and evenly spread out output values. In the first split of the dataset, there are 700 entries in total. In the second split, the dataset contains 1400 entries of data points. The final dataset consists of 3000 data points. The models developed were executed on the three variations of the dataset to analyze the parameters like Mean

Squared Error, Accuracy, Loss and Time taken for training and ultimately deciding the most applicable model for handling time series data as well as long term dependencies.

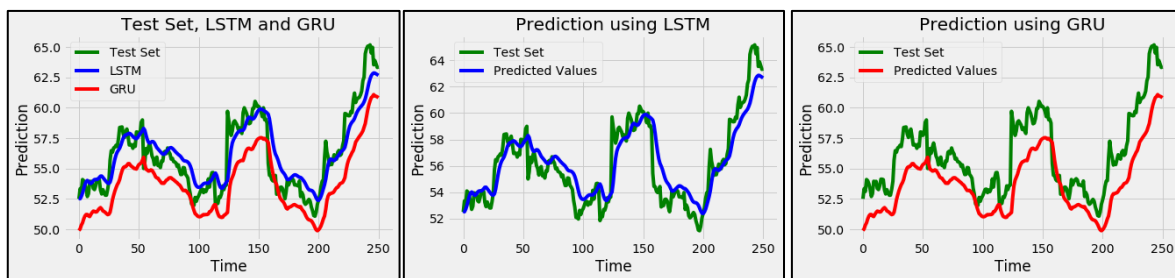
**(a) Dataset 1**



**(b) Dataset 2**



**(c) Dataset 3**



**Figure 16.** Outputs for dataset variations

For Dataset 1 with 700 data points, the performance of GRU seems more better than LSTM as the predicted values by GRU and the Testing values map well than that of LSTM. In iteration for Dataset 2 with 1400 data points, the performance of LSTM and GRU both seem almost equal as the values predicted by both the models correspond with the testing values. In the iteration for Dataset 3, GRU shows more error and deflect away from the Testing Data points. On the other hand LSTM shows better performance over 3000 data points from Dataset 3.

*Mean Squared Error Analysis*

The mean squared error is calculated for each dataset variation for both LSTM and GRU model. Mean squared error is one way to calculate the error during Back propagation which is the basis for or training for both LSTM and GRU. Higher error value indicates performance degradation and improper training. The mean squared error for both models on three datasets is presented in Table 5.

**Table 6.** Mean Squared Error for LSTM and GRU on Datasets.

	LSTM	GRU
<b>Dataset 1</b>	0.7516197811968257	0.4601261472286205
<b>Dataset 2</b>	2.374904252912631	2.2920254537142664
<b>Dataset 3</b>	1.4342299451204816	2.6693156947773167

From the values of Mean Squared Error we can conclude that GRU performs well and generates low error for datasets with small size with lesser data points. LSTM was found performing well for datasets with large data points.

*Average Loss Calculation*

The loss values are calculated for each epoch for both LSTM and GRU. For a given model, the lesser the loss the better is the training of the model. Loss is summation of errors made for each batch of training dataset over an epoch. The average of loss values over the complete training procedure for 50 epochs is specified in the Table 6 for the variations in dataset.

**Table 7.** Loss values for each dataset

	Dataset 1	Dataset 2	Dataset 3
<b>LSTM</b>	0.002012	0.013321	0.001093
<b>GRU</b>	0.00151	0.00129	0.02667

Owing to less computational steps, GRU generates less error than LSTM. LSTM executes well over the last dataset variation having the highest number of data points.

*Training Time Comparison*

The training time for the complete training process in minutes is specified for both LSTM and GRU in Table 7.

**Table 8.** Training Time in minutes for LSTM and GRU

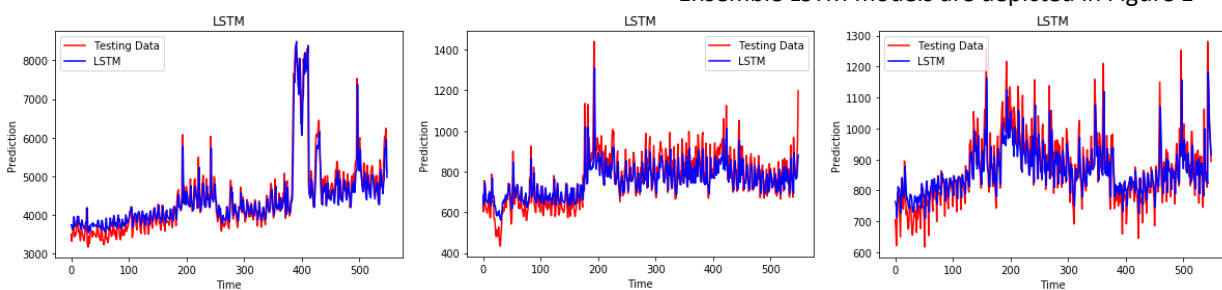
	LSTM	GRU
<b>Dataset 1</b>	9.51	8.28
<b>Dataset 2</b>	12.91	11.59
<b>Dataset 3</b>	16.8	14.97

For all the disparity in the dataset, LSTM model takes more time to train over GRU. When LSTM and GRU were used to train the chatbot models, the model with LSTM showed overall better accuracy of 71.69% over 70.12% accuracy of GRU during Training. During Testing the LSTM model showed better accuracy of 71.59% over 70.75% of GRU model.

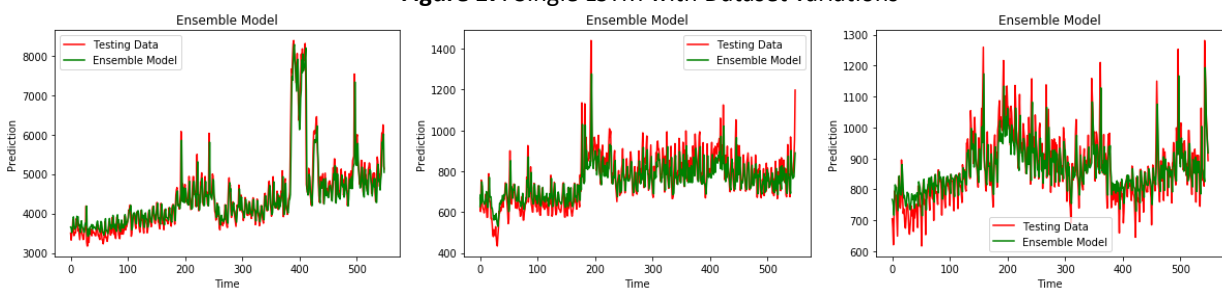
Considering four parameters like Mean Squared Error, Loss, Training Time and Accuracy we can deduce LSTM is a better choice over GRU as whole.

*Comparison between LSTM and Ensemble LSTM*

Following graphs present a comparative look at the performance of LSTM and Ensemble Model for the same time series data analysis. The Ensemble Model consists of three LSTMs with variation in hyperparameters. The graphs depict the mapping between the values from the testing dataset and the predicted values from the corresponding model. The graphs related to LSTM models are specified in Figure 17 and the graphs related to Ensemble LSTM models are depicted in Figure 1



**Figure 17.** Single LSTM with Dataset variations



**Figure 18.** Ensemble LSTM Model with Dataset Variations

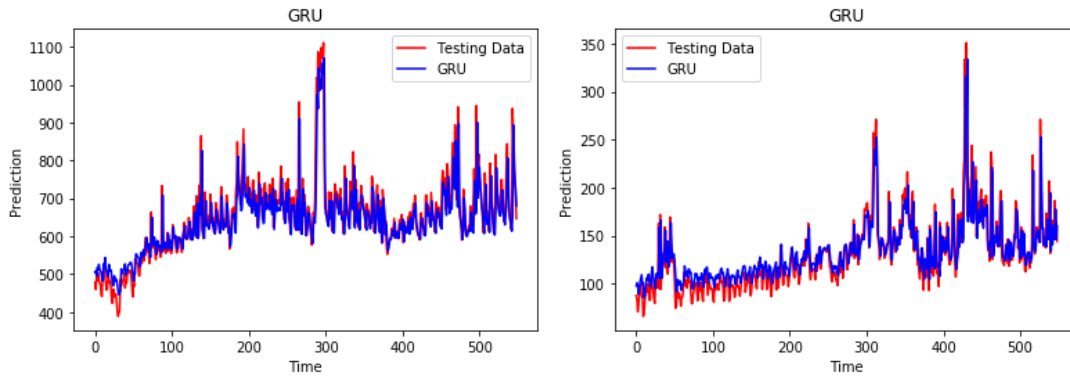


Figure 19. Single GRU model with Dataset variations

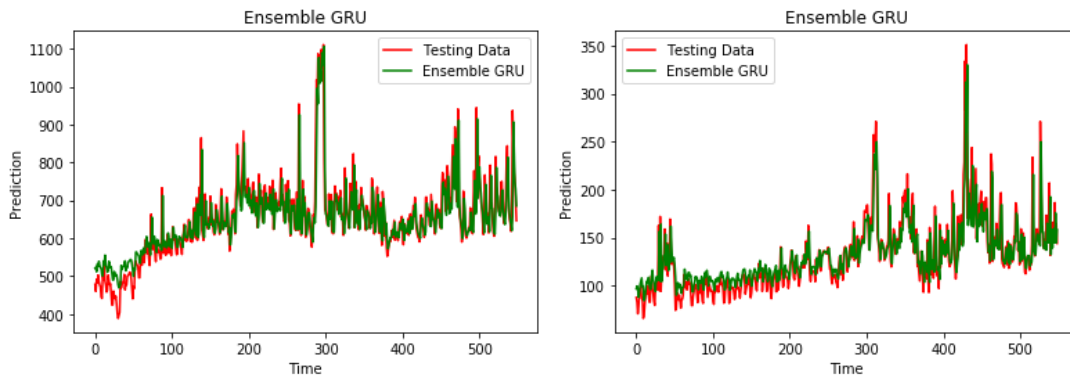


Figure 20. Ensemble GRU with Dataset variations

From the graphs we can conclude the performance of Ensemble Model is similar in the aspect of the performance of LSTM model and in some cases better. With fine tuning the Ensemble models, the performance can be improved over standalone LSTM model.

*Comparison between GRU and Ensemble GRU*

The following graphs in Figure 19 and 20 provide overall performance analysis for single GRU model and Ensemble GRU model in respective sections. The Ensemble GRU model consists of three individual GRU models with variations in hyperparameters.

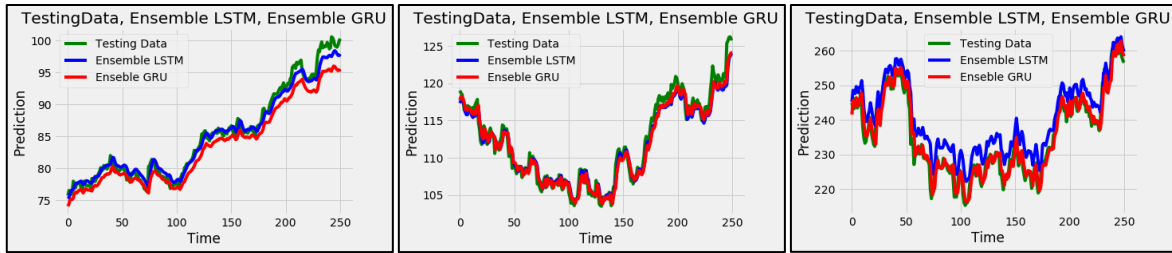
We derive from the graphs, the performance of single GRU and the Ensemble GRU is comparatively similar. The performance could be better enhanced with the diversity in the hyperparameters of singular models constituting the Ensemble model.

*Comparison between Ensemble LSTM and Ensemble GRU*

This section provides the differentiation in the performance between the Ensemble LSTM model and Ensemble GRU model against the prediction values from the testing dataset. The training was carried out on variations of the same dataset hence employing the strategy of cross-validation to make it most even comparison.

*Prediction Comparison – Ensemble LSTM, Ensemble GRU*

From Figure 21, it can be deduced from the graph i and ii that the performance of both ensemble models are almost equal in terms of close mapping with the prediction values. However, in graph iii, it is observed that Ensemble GRU performance better than Ensemble LSTM. The values predicted by Ensemble GRU are closely mapped with the actual values from the testing dataset.



**Figure 21.** Results over dataset variations for Ensemble LSTM and Ensemble GRU

### Mean Squared Error Analysis

The final Mean Squared Error for both Ensemble models is calculated. Lower the values, better the performance of the model. The values are defined in the Table 8.

**Table 9.** Mean Squared Error Analysis for Ensemble LSTM and Ensemble GRU

	Ensemble LSTM	Ensemble GRU
<b>Dataset 1</b>	1.060990337370966	2.0204261732384396
<b>Dataset 2</b>	1.249621120822286	1.1770008219011074
<b>Dataset 3</b>	5.661777383126808	2.6198470374584173

From the values calculated we can deduce that no specific model performs better than the other in every variation of the dataset. In two cases, the performance of both models was comparatively equal while in one case the Ensemble GRU performs better than Ensemble LSTM.

### Average Loss Calculation

For each epoch during the training, Loss values are calculated. Lesser the loss values, better is the training of the model. The average loss values calculated for each dataset variation over the complete training is presented in the Table 9.

**Table 10.** Average Loss Calculation for Ensemble LSTM and Ensemble GRU

	Dataset 1	Dataset 2	Dataset 3
<b>Ensemble LSTM</b>	0.001069	0.034829	0.000179
<b>Ensemble GRU</b>	0.002512	0.00278	0.01917

From the values in the table, we deduce the average loss during the training of Ensemble GRU is lesser than the values of Ensemble LSTM. That indicates Ensemble GRU trains well with the dataset compared to Ensemble LSTM. Based on the analysis of three parameters, we can conclude the performance of both Ensemble LSTM and Ensemble GRU is not definitively better than each other.

In some cases Ensemble GRU performed better while in some cases Ensemble LSTM. With tweaking the parameters of individual models working together as Ensemble, the best performance can be achieved.

### Conclusion

This paper presents a review of the evolution of technologies applied in Chatbots handling time series conversations in the labels of Architectural Design and Implementation. The paper also intends to contribute in developing a sturdy groundwork on the concepts utilized in learning long term dependencies hence providing a roadmap towards further enhancements being inclined towards minimalistic yet alike requisites. These primal conditions can be considered as 1. Designing the word embedding schema not constrained by the knowledge base. 2. Flexible and accurate conversational model. 3. Reaching the true peak of imitating the human conversation requiring no human intervention. The proposed LSTM based Ensemble Network architecture attempts at enhancing the user experience by providing a sense of continuance of context in a series of conversations. The algorithm does so by generalizing the features imperative to making the conversation humane.

### References

[1] Io, H. N., & Lee, C. B. (2017). Chatbots and conversational agents: A bibliometric analysis. 2017 IEEE International Conference on Industrial Engineering and Engineering Management (IEEM). <https://doi.org/10.1109/ieem.2017.8289883>

[2] Lokman, A. S., & Amedeen, M. A. (2018). Modern Chatbot Systems: A Technical Review. Proceedings of the Future Technologies Conference (FTC) 2018 Advances in Intelligent Systems and Computing, 1012-1023. [https://doi.org/10.1007/978-3-030-02683-7\\_75](https://doi.org/10.1007/978-3-030-02683-7_75)

[3] Goyal, P., Pandey, S., & Jain, K. (2018). Developing a Chatbot. Deep Learning for Natural Language

Processing,169-229.

[https://doi.org/10.1007/978-1-4842-3685-7\\_4](https://doi.org/10.1007/978-1-4842-3685-7_4)

[4] Salehinejad, H., Sankar, S., Barfett, J., Colak, E., & Valaee, S. (2017). Recent Advances in Recurrent Neural Networks. arXiv.org.

[5] Verleysen, M., & François, D. (2005). The Curse of Dimensionality in Data Mining and Time Series Prediction. Computational Intelligence and Bioinspired Systems Lecture Notes in Computer Science,758-770. [https://doi.org/10.1007/11494669\\_93](https://doi.org/10.1007/11494669_93)

[6] Lawrence, S., & Giles, C. (2000). Overfitting and neural networks: Conjugate gradient and backpropagation. Proceedings of the IEEE-INNS-ENNS International Joint Conference on Neural Networks. IJCNN 2000. Neural Computing: New Challenges and Perspectives for the New Millennium. <https://doi.org/10.1109/ijcnn.2000.857823>

[7] Lipton, Zachary. (2015). A Critical Review of Recurrent Neural Networks for Sequence Learning.

[8] Gradient Flow in Recurrent Nets: The Difficulty of Learning Long Term Dependencies. (2009). A Field Guide to Dynamical Recurrent Networks. <https://doi.org/10.1109/9780470544037.ch14>

[9] Olah, Christopher. Understanding LSTM Networks.n.d. 2019.

[10] Hochreiter, S., & Schmidhuber, J. (1997). Long Short-Term Memory. Neural Computation,9(8), 1735-1780. <https://doi.org/10.1162/neco.1997.9.8.1735>

[11] Gers, F. (1999). Learning to forget: Continual prediction with LSTM. 9th International Conference on Artificial Neural Networks: ICANN 99. <https://doi.org/10.1049/cp:19991218>

[12] Chung, J., Gulcehre, C., Cho, K. and Bengio, Y. Chung, J., Gulcehre, C., Cho, K., & Bengio, Y. (2014). Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling arXiv:1412.3555

[13] Collier, M. and Beel, J. Collier, M., & Beel, J. (2018). Implementing Neural Turing Machines. arXiv:1807.08518

[14] Opitz, D., & Maclin, R. (1999). Popular Ensemble Methods: An Empirical Study. Journal of Artificial Intelligence Research,11, 169-198.

<https://doi.org/10.1613/jair.614>

[15] Hansen, L., & Salamon, P. (1990). Neural network ensembles. IEEE Transactions on Pattern Analysis and Machine Intelligence,12(10), 993-1001. <https://doi.org/10.1109/34.58871>

[16] Chen, C., Wu, C., Lo, C., & Hwang, F. (2017). An Augmented Reality Question Answering System Based on Ensemble Neural Networks. IEEE Access,5, 17425-17435. <https://doi.org/10.1109/access.2017.2743746>

[17] M. Nuruzzaman and O. K. Hussain, "A Survey on Chatbot Implementation in Customer Service Industry through Deep Neural Networks," 2018 IEEE 15th International Conference on e-Business Engineering (ICEBE), Xi'an, 2018, pp. 54-61. <https://doi.org/10.1109/ICEBE.2018.00019>

[18] Lee, M. C., Chiang, S. Y., Yeh, S. C., & Wen, T. F. (2020). Study on emotion recognition and companion Chatbot using deep neural network. MULTIMEDIA TOOLS AND APPLICATIONS

[19] Pathak, K., & Arya, A. (2019, November). A Metaphorical Study Of Variants Of Recurrent Neural Network Models For A Context Learning Chatbot. In 2019 4th International Conference on Information Systems and Computer Networks (ISCON) (pp. 768-772). IEEE.

[20] G. Dzakwan and A. Purwarianti, "Comparative Study of Topology and Feature Variants for Non-Task-Oriented Chatbot using Sequence to Sequence Learning," 2018 5th International Conference on Advanced Informatics: Concept Theory and Applications (ICAICTA), Krabi, 2018, pp. 135-140. <https://doi.org/10.1109/ICAICTA.2018.8541285>

[21] Bali, M., Mohanty, S., Chatterjee, S., Sarma, M., & Puravankara, R. Diabot: A Predictive Medical Chatbot using Ensemble Learning.

[22] M. Nuruzzaman and O. K. Hussain, "A Survey on Chatbot Implementation in Customer Service Industry through Deep Neural Networks," 2018 IEEE 15th International Conference on e-Business Engineering (ICEBE), Xi'an, 2018, pp. 54-61. <https://doi.org/10.1109/ICEBE.2018.00019>

[23] Lee, M. C., Chiang, S. Y., Yeh, S. C., & Wen, T. F. (2020). Study on emotion recognition and companion Chatbot using deep neural network. MULTIMEDIA TOOLS AND APPLICATIONS.



[24] Bali, M., Mohanty, S., Chatterjee, S., Sarma, M., & Puravankara, R. Diabot: A Predictive Medical Chatbot using Ensemble Learning.

[25] Pathak, K., & Arya, A. (2019, November). A Metaphorical Study Of Variants Of Recurrent Neural Network Models For A Context Learning Chatbot. In 2019 4th International Conference on Information Systems and Computer Networks (ISCON) (pp. 768-772). IEEE.


[26] G. Dzakwan and A. Purwarianti, "Comparative Study of Topology and Feature Variants for Non-Task-Oriented Chatbot using Sequence to Sequence Learning," 2018 5th International Conference on Advanced Informatics: Concept Theory and Applications (ICAICTA), Krabi, 2018, pp. 135-140.

<https://doi.org/10.1109/ICAICTA.2018.8541285>

---

**Publisher:** Chinese Institute of Automation Engineers (CIAE)

**ISSN:** 2223-9766 (Online)

 **Copyright:** The Author(s). This is an open access article distributed under the terms of the [Creative Commons Attribution License \(CC BY 4.0\)](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are cited.