

BOON : An Intelligent Chatbot using Deep Learning

Rahul K Chauhan

Saint Louis University

rahul.chauhan@slu.edu

Kevin Scannell

Saint Louis University

kevin.scannell@slu.edu

Abstract

Dialogue Generation or Intelligent Conversational Agent development using Artificial Intelligence or Machine Learning technique is an interesting problem in the field of Natural Language Processing. In many research and development projects, they are using Artificial Intelligence, Machine Learning algorithms and Natural Language Processing techniques for developing conversation/dialogue agent. Their research and development is still under progress and under experimentation. Dialogue/conversation agents are predominately used by businesses, government organizations and non-profit organizations. They are frequently deployed by financial organizations like bank, credit card companies, businesses like online retail stores and start-ups. These virtual agents are adopted by businesses ranging from very small start-ups to large corporations. There are many chatbot development frameworks available in market both code based and interface based. But they lack the flexibility and usefulness in developing real dialogues. In this project, I have developed intelligent conversational agent BOON using state of the art techniques proposed in recently published research papers. The name BOON simply means talkative. For developing intelligent chatbot, I have used Google's Neural machine Translation(NMT) Model which is based on Sequence to Sequence(Seq2Seq) modeling with encoder-decoder architecture. This encoder-decoder is using Recurrent Neural Network with bi-directional LSTM (Long-Short-Term-Memory) cells. For performance optimization, I applied Neural Attention Mechanism and Beam Search during training.

1 Introduction

Conversational agent or Chatbot is a program that generates response based on given input to emulate human conversations in text or voice

mode. These applications are designed to simulate human-human interactions. Chatbots are predominantly used in business and corporate organizations including government, non-profit and private ones. Their functioning can range from customer service, product suggestion, product inquiry to personal assistant. Many of these chat agents are built using rule based techniques, retrieval techniques or simple machine learning algorithms. In retrieval based techniques, chat agents scan for keywords within the input phrase and retrieves relevant answers based on the query string. They rely on keyword similarity and retrieved text is pulled from internal or external data sources including world wide web or organizational database. Some other advanced chatbots are developed with natural language processing(NLP) techniques and machine learning algorithms. Also, there are many commercial chat engines available, which help build chatbots based on client data input.

2 Related Works

There have been many recent development and experimentation in conversational agent system. Apart from traditional chatbot development techniques that use rule based techniques, or simple machine learning algorithms, many advanced chatbots are using advanced Natural Language Processing (NLP) techniques and Deep Learning Techniques like Deep Neural Network (DNN) and Deep Reinforcement Learning (DRL).

2.1 Sequence to Sequence (Seq2Seq)

Some of the state of the art techniques(Sojasingarayar, 2020) involve using Deep Neural Network and its architectural variations. Sequence to Sequence (Seq2Seq) model based on encoder-decoder architecture is such an architecture which is very popular for dialogue generation, language modeling and

machine translation. Seq2Seq uses Recurrent Neural Network(RNN) which is a popular Deep Neural Network architecture specially for Natural Language Processing tasks. In Sequence to Sequence (Seq2Seq) model, many to many RNN architecture is used for decoder. In this, encoder-decoder architecture, input sequence is fed as a vector representation of text to encoder. Then, encoder produces some intermediate representation of information or thought vectors. Consequently, the thought vector generated by encoder is fed into decoder as input. Finally, decoder processes the thought vector and converts the sequence

one by one word and produces multiple output from the decoder in form of target sequence. Though, vanilla RNN is default in Seq2Seq and works well for many NLP problems yet, due to higher complexity of language modeling problem, vanilla recurrent neural network cells often fails, specially, where long sequence of information needs to be remembered, as this information frequently becomes large for bigger datasets and turns to information bottleneck for the RNN network. Therefore, researchers uses variations of recurrent neural network to handle such problem.

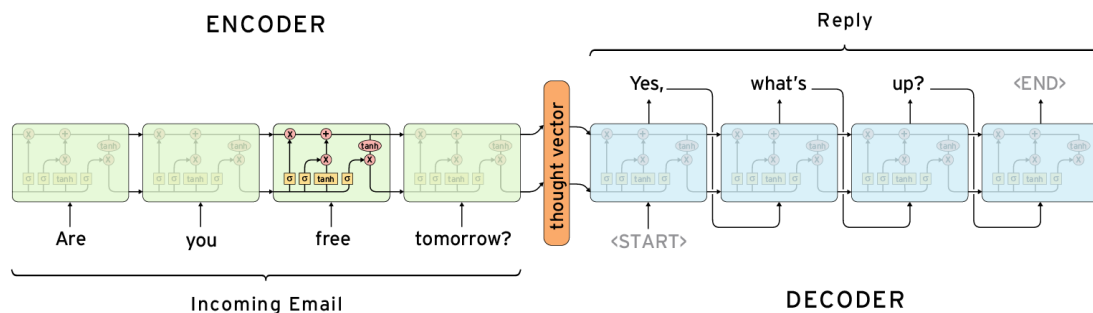


Figure 1: Sequence to Sequence Model

Long-Short-Term-Memory(LSTM) is a special variant of cell type of Recurrent Neural Network which has empirically shown to work well for language modeling. LSTM has forget gates along with input gates and output gates. This helps remember more relevant and contextual information and discards the rest of the sequence which is desirable in language modeling where dependency within sequence is sparse. Also, instead of using unidirectional cells, bidirectional LSTM cells can perform much better.

lation, text summarization and question-answering and image captioning.

Long short Term Memory

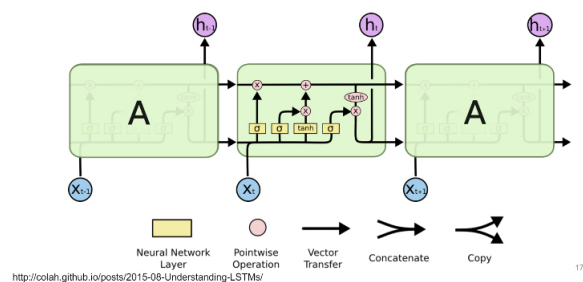


Figure 2: Ref 7

Another technique, Neural Attention Mechanism embedded in Seq2Seq module has significantly improved performance in dialogue generation system and other NLP tasks and thus become industry standard practice. In Neural attention mechanism, each hidden target compares with source hidden state, generates attention vector by calculating score and preserves the attention vector in memory to choose over other candidate. Also, other techniques like, Beam Search can help improve decoding performance further by choosing top candidates. Seq2Seq have also been applied for other NLP tasks including machine trans-

2.2 Google's Neural Machine Translation

Google's Neural Machine Translation(GNMT) model is a module for neural machine translation from and to other language and English. GNMT has also been used for dialogue generation experimentally. It is based on Seq2Seq model which is popular in dialogue generation. Also, GNMT has many techniques embedded in the module which are crucial for intelligent chatbot develop-

ment. The GNMT model includes, Sequence to Sequence modeling with encoder-decoder architecture built using uni or bi directional LSTM cells. They also have option for Neural Attention Mechanism, Beam Search, and vocabulary generation using Google's sub-word module. Also, they have option for adjusting the hyperparameters for better model training.

2.3 Deep Reinforcement Learning

"Deep Reinforcement Learning for Dialogue Generation"(Li et al., 2016) of Dan Jurafsky, Deep Reinforcement Learning (DRL) has been used for developing long conversation chatbots. Seq2Seq model can generate coherent dialogues but may produce repeated generic responses regardless of input and can get stuck in a loop in longer conversations. This occurs as Seq2Seq predicts utterances one at a time while ignoring their influence on future outcomes. Seq2Seq models tend to generate highly frequent repeated responses like "I don't know". This is due to high frequency of generic responses in the training set, also this replies are more compatible with a wide range of input text.

In Dufarsky's paper[(Li et al., 2016), they have generated intermediate response using Seq2Seq model with attention where input was raw text. Then, the intermediate generated responses were fed into Reinforcement Model and was rewarded based on Ease of answering, Information Flow and Semantic Coherence. This is forward centric model, where if generated response is easy to answer, contribute to more information compared to previous dialogue history and grammatically and semantically correct, they are rewarded.

3 Architecture

In this project, I have developed Intelligent conversational agent following state of the art techniques proposed in recently published research papers(Wu et al., 2016). I have used Google's Neural Machine Translation(GNMT) module for building dialogue generator. Although, Google's Neural Machine Translation(GNMT) module is primarily used for Machine Translation tasks, they have empirically shown to be successful in other NLP tasks including dialogue generation and text summerization. GNMT has rich Seq2Seq module with many additional features for dialouge genera-

tion. Seq2Seq is a industry standard choice for dialogue generation and many NLP tasks. Although, there exists a separate Seq2Seq module(seq2seq) which was the the early robust Seq2Seq module, but is not currently compatible with latest version of Google's machine learning framework Tensorflow and is only compatible with Tensorflow v1.0 whereas current Tensorflow is v1.6. Also, earlier version of Tensorflow had a minimalistic Seq2Seq module, it is not efficient for building robust model. Currently, Google does not offer any robust Seq2Seq module in it's official Machine Learning Framework Tensorflow and has moved many functionality of Seq2Seq to NMT model. Moreover, GNMT has many techniques embedded for dialogue generation agent development. GNMT model includes, Sequence to Sequence modeling (Seq2Seq) with encoder-decoder architecture based on uni or bi directional LSTM cells. GNMT also have option for Neural Attention Mechanism, Beam Search, and vocabulary generation using Google's sub-word module.

3.1 Google's Neural Machine Translation

3.1.1 Sequence to Sequence (Seq2Seq) Architecture

Google's Neural Machine Translation (GNMT) is primarily used for machine translation. But, GNMT contains Sequence to sequence module with many enhancement techniques which can help build good dialogue generator. For translation, GNMT, does not apply traditional phrase-based translation systems where translation is performed by breaking up source sentences into multiple chunks and then translate phrase-by-phrase. It rather uses more human like translation approach.

GNMT is based on Seq2Seq architecture composed of encoder-decoder unit. GNMT can be used with with different variation of architectures. The Seq2Seq module is composed of encoder and decoder. Encoder takes source text as input and processes the text to generate intermediate representation of input text called thought vector. The thought vector is then fed into the decoder input unit. The decoder now processes the thought vector and generates outputs. In case of dialogue generation problem the output is a response and for machine translation problem output is the target text. This architecture have shown to be capable of processing long-range dependencies and produce

more fluent translations or responses.

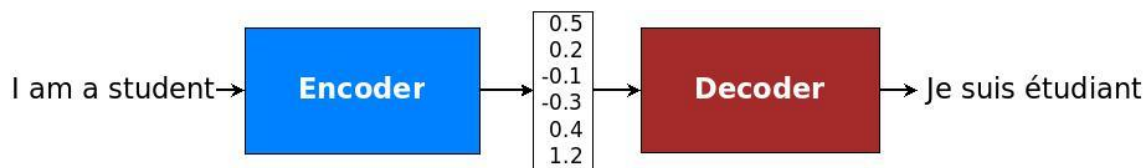


Figure 3: Encoder-decoder architecture – example of a general approach for NMT. An encoder converts a source sentence into a "meaning" vector which is passed through a decoder to produce a translation.

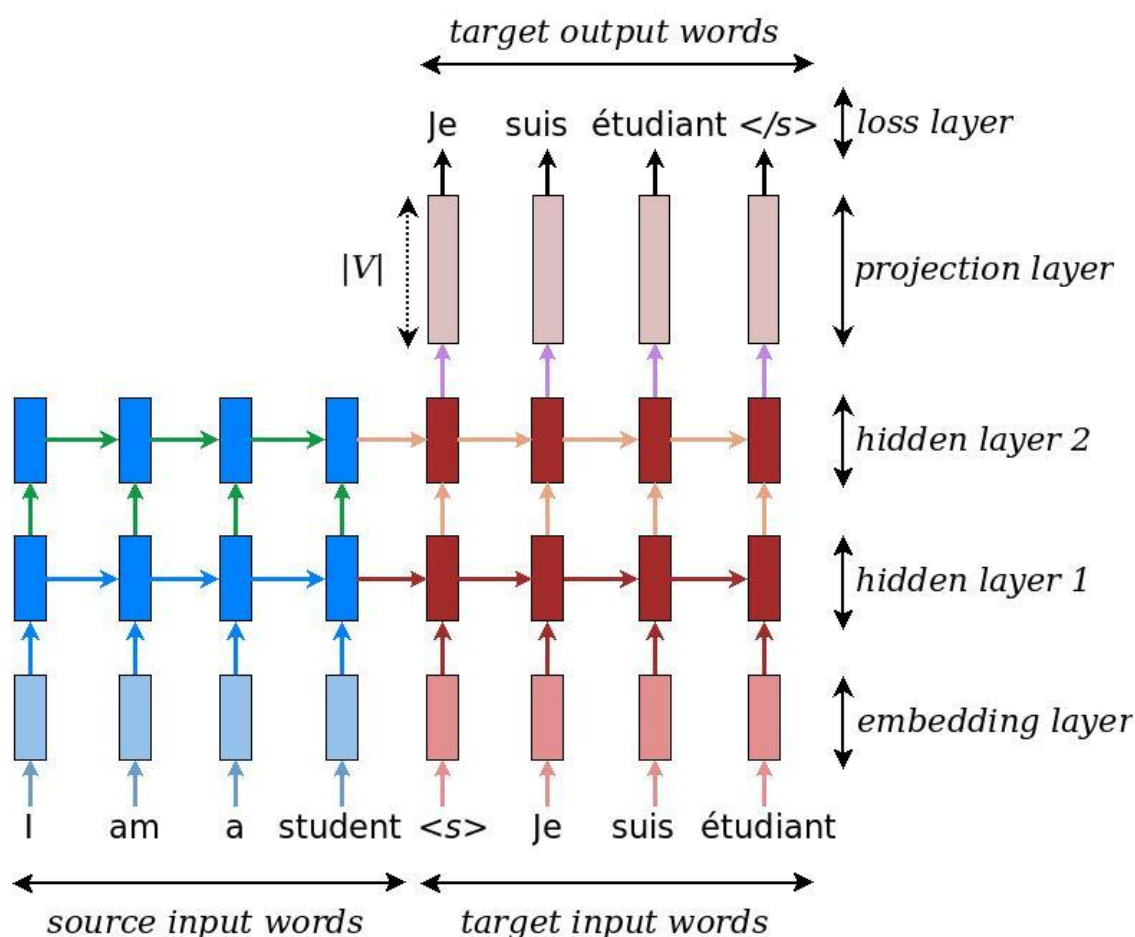


Figure 4: Neural machine translation – example of a deep recurrent architecture proposed by for translating a source sentence "I am a student" into a target sentence "Je suis étudiant". Here, "< s >" marks the start of the decoding process while "< /s >" tells the decoder to stop.

3.1.2 Neural Attention Mechanism (Google)

To build state-of-the-art neural machine translation systems, the attention mechanism works specially well, which was first introduced in 2015. The idea of the attention mechanism is to form direct short-cut connections between the target and the source by paying "attention" to relevant source content as we translate(Google).

In the vanilla seq2seq model, the last source state from the encoder to the decoder is passed

when starting the decoding process. This works well for short and medium-length text string. However, for long sentences, the single fixed-size hidden state becomes an information bottleneck. Instead of discarding all of the hidden states computed in the source RNN, the attention mechanism provides an approach that allows the decoder to peek at them (treating them as a dynamic memory of the source information). This attention mechanism improves the translation of longer sentences. Attention mechanisms are the industry standard

and applied to other NLP tasks including image caption generation, speech recognition, and text summarization.

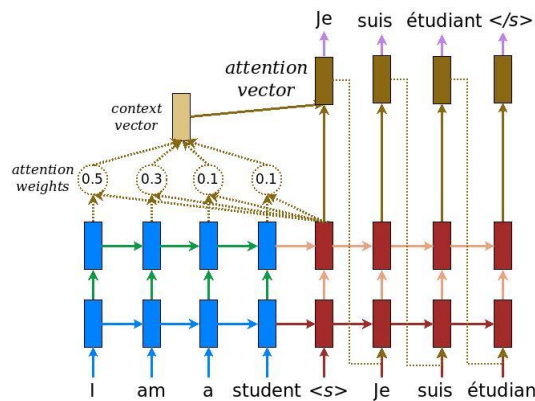


Figure 5: Attention mechanism – example of an attention-based NMT system.

4 Data

4.1 Data Collection

In this project, the dataset "Cornell Movie Subtitle Corpus" has been primarily used for final model training. Few other dialogue corpus based on movie subtitles have been preprocessed and cleaned for GNMT training including "Open Movie Subtitle Corpus" and "Movie Subtitle Corpus". But due to lack of data quality, they have been eliminated from final training.

4.1.1 Cornell Movie Subtitle Corpus

For developing final chatbot, popular movie subtitle corpus "Cornell movie subtitle corpus" has been used. This corpus contains metadata-rich large collection of conversations extracted from raw movie scripts from popular movies.

The following are found in corpus-

- 220,579 conversational exchanges between 10,292 pairs of movie characters.
- involves 9,035 characters from 617 movies
- in total 304,713 utterances

Other movie meta-data included genres, release year, IMDB rating, number of IMDB votes, IMDB rating

The data corpus can be found in http://www.cs.cornell.edu/~cristian/Cornell_Movie-Dialogs_Corpus.html.

4.2 Data Preprocessing

Data was processed to prepare for input pipeline of the Google's Neural Machine Translation(GNMT) model. In the original GNMT model, there were two input data files and two vocabulary file generated from input files. The two input files were translation from and translation to language input data file. The vocabulary files contained the processed vocabulary for the two input data file of two different language respectively. Also, there were separate test and development file for source and target.

4.2.1 Preprocessing of Cornell Movie Subtitle Corpus

Conversation data in the movie corpus contained Movie ID, Character ID, and Movie Line ID was separated by "+++++".

For preprocessing, conversation data was cleaned to remove these meta-data(eg. movie ID, character ID, Line ID). Also, data separators("+++++") were eliminated. Additionally, some of the character in the data contained unsupported encoding format by UTF-8 standard and was hence removed.

Finally, data was separated in two different files to assimilate with the format of Google's Neural Machine Translation(GNMT) model input pipeline format where first file is the dialogue 1 and second one was response to dialogue 1.

After separating the two files, data in both file was cleaned simultaneously. Everything except alphabetical character, and some punctuation (. , ?!) was removed as they hold little meaning in conversation. Also, all the text was converted to lowercase. Then, multiple consequent occurrence of these punctuation (. , ?!) was reduced to one in order to reduce punctuation overload. Next, all the punctuation except (') was separated with single space before and after for better performance in GNMT module. Finally, all the consequent multiple space was reduced to single space and each text string was trimmed to remove before and after space.

Also, data was cleaned for removing extraneous dialogues. If multiple consequent utterance from single person was present everything except

the last utterance for the person was stored. Initially, utterance with more than 100 length was discarded for both text dialogue and their reply as with increase of length the text, context relevance starts to drop due to diversity and limited data. But later full text length was embedded.

After through clenaing, the source and target text was splitted for training, testing and development/validation set with source and target format and was saved in files for final input pipeline feed.

For vocabulary generation, Google's Sub-word Neural Machine Translation(NMT) module was used as suggested by the Google Tensorflow and Google's Neural Machine Translation module documentation. The sub-word application was only applied on training files source and target files.

5 Implementation Summary

- Deep Learning Module : Google's Neural Machine Translation (NMT).
- Algorithm : Deep Neural Network (DNN), Recurrent Neural Network (RNN)
- Main Technique : Sequence to Sequence (Seq2seq) modeling
- Enhancement Techniques : Long Short Term Memory (LSTM) based RNN cell
- Main Technique : Sequence to Sequence (Seq2seq) modeling, Bidirectional LSTM, Neural Attention Model and Beam Search.

6 Training

6.1 Training Dataset

Training has been completed on 225000*2 utterance of "Cornell movie subtitle corpus" conversation and has been tested with 5000+5000 utterance and validated with 5000+5000 utterance.

6.2 Training Model and Parameters

For training, GNMT module based on Bidirectional LSTM cells have been used. Neural Attention Mechanism have been applied to improve performance. Beam search was also applied on training but was later discarded. Total 3200 iteration with 512 hidden unit was trained with 2 layers.

Parameter Name	Value
Num of Train Steps	3200
Steps Per Stats	100
Num of Units	512
Num of Layers	2
Cell Type	LSTM
Encoder type	Bidirectional
Neural Attention	scaled luong

Parameter Name	Value
Optimizer	adam
Learning Rate	0.001
Decay Steps	1
Start Decay Step	1
Beam Width	10
Dropout	0.2
Metrics	bleu

7 Result

Following are some response derived after training on full dataset with 34MB of training text with 225000+225000 utterance. The initial test result produced moderately coherent sentences. The following responses were generated after inference from trained model. In inference, trained model produced 10-30 candidate response for each input. Hence, the following examples are most suited from 10-30 candidate responses during inference.

Input (Person 1)	GNMT Model Output (Person 2 Response)
how are you	how am i?
whats your plan for the weekend ?	i dont know
would you like some tea	would you like some help ?
You got something on your mind	yeah
would you like to go to watch movie ?	no .
where would you like to go for lunch today	i know where to go for lunch today .
i know where to go for lunch today .	you know where to go ?
Where is the library located ?	is that where you were ?
weather seem warm would you like to get some cold drinks	would you like to get some cold coffee ?
did you like the new startrek movie sequel	i did .
band played really well i am going for their next tour	i am going for their next interview .
am i disturbing you ?	no . not at all . where are you ?
i'm heading out . how about you ?	i got to wait for a call .

Following are the Perplexity and Bleu for test and development dataset.

Table 1: Evaluation

Dataset	Perplexity	Bleu
eval dev	50.76	10.1
eval test	46.82	10.6

8 Challenges

The challenge in developing chatbot or dialogue generator lies in developing coherent dialogue generation system. As the model used in this experiment is for machine translation, the dialogue generation is treated as translation problem, where history of earlier conversations are not taken into account. Hence, the model can be limited in performance regarding long conversation.

Also, training is a long process which demands higher processing power and configured computing machine. Another problem is finding right hyper parameters to optimize the translation module for chat bot or dialogue generation system.

After training, chatbot produced results with moderate relevancy. But many of the output were repetitive and generic. Also, due to lack of real-life quality data the chatbot performed somehow below optimum for imitating human interaction. Also, many utterance was discarded due to longer length or discrepancy. And, number of training utterance was much less than required and test and development dataset was quite larger in comparison which might have caused the model to underperform. Also, as data was limited, longer period of training may not have suited the dialogue generation problem.

9 Future Work

The chatbot developed using Google's Neural Machine Translation Model(GNMT) can be further improved with more robust, high quality real-life conversational datasets which could better emulate human interaction. Also, hyper-parameters of the GNMT model can be further fine-tuned and optimized for performance enhancement. Based on available opportunity to further advance the project, Deep Reinforcement Learning(RL) can be applied that could significantly improve performance as shown empirically in Dufarsky's paper.

Reinforcement Learning algorithm can be applied after the initial training using Google's Neural Machine Translation.

10 Conclusion

The training on Cornell Movie Subtitle corpus produced result which needs further improvement and more attention and speculation on training parameters. Adding more quality data will further improve performance. Also, the training model should be trained with other hyper-parameters and different dataset for further experimentation. This was an attempt to experiment with Deep Neural Network for dialogue generation in order to develop intelligent chatbot.

References

- Google. [Google's Neural machine translation\(NMT\)](#). Google.
- Jiwei Li, Will Monroe, Alan Ritter, Michel Galley, Jianfeng Gao, and Dan Jurafsky. 2016. [Deep reinforcement learning for dialogue generation](#). *CoRR*, abs/1606.01541.
- seq2seq. [seq2seq GitHub Repository](#).
- Abonia Sojasingarayar. 2020. [Seq2seq ai chatbot with attention mechanism](#).
- Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Lukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, and Jeffrey Dean. 2016. [Google's neural machine translation system: Bridging the gap between human and machine translation](#). *CoRR*, abs/1609.08144.